

# Integration of multi-omics data with artificial intelligence

Citation for published version (APA):

Ochoteco Asensio, J. (2022). Integration of multi-omics data with artificial intelligence: studying the toxic effects on the post-transcriptional regulation. [Doctoral Thesis, Maastricht University]. Maastricht University. https://doi.org/10.26481/dis.20221109ja

Document status and date: Published: 01/01/2022

DOI: 10.26481/dis.20221109ja

**Document Version:** Publisher's PDF, also known as Version of record

#### Please check the document version of this publication:

 A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.

• The final author version and the galley proof are versions of the publication after peer review.

 The final published version features the final layout of the paper including the volume, issue and page numbers.

Link to publication

#### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these riahts.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
   You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.umlib.nl/taverne-license

#### Take down policy

If you believe that this document breaches copyright please contact us at:

repository@maastrichtuniversity.nl

providing details and we will investigate your claim.

Department of Toxicogenomics, Maastricht University

Integration of multi-omics data with artificial intelligence: studying the toxic effects on the post-transcriptional regulation

Doctoral Thesis Manuscript

Juan Ochoteco Asensio 7-1-2022

© Juan Ochoteco Asensio, Maastricht, the Netherlands 2022

Cover design by Juan Ochoteco Asensio Layout by Fred Feij Printed by Proefschriftenprinten.nl

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without prior permission of the author or the copyright-owning journals for previously published papers.

#### Supervisor

Prof. dr. J.C.S. Kleinjans

#### **Co-supervisor**

Dr. F. Caiment

#### Assessment committee

- Prof. dr. H.J.M. Smeets, Full Professor of Clinical Genomics (Department of Toxicogenomics, Maastricht University)
  Prof. dr. J. Quackenbush, Professor of Computational Biology and Bioinformatics (Department of Biostatistics, Harvard University)
- Prof. dr. R. Peeters, Full professor of Mathematics in Knowledge Engineering (Department of Data Science and Knowledge Engineering, Maastricht University)
- Dr. S. Hayat, Group Leader of Translational Data Science (Uniklinik RWTH Aachen)
- Dr. L. Eijssen, Assistant Professor (Department of Bioinformatics, Maastricht University)

The research described in this thesis was conducted at GROW School for Oncology and Reproduction of Maastricht University.

# Chapter 1:

# **General Introduction**

# Introduction

In the current age of big data, several technologies are helping generate more and more data for the 2 advancement of science. The extraction of data from the cell is highly relevant, as bodily changes are in 3 the end cell changes, especially in the case of toxicology. Therefore, technologies such as transcriptomics 4 5 and proteomics are very informative on the changes that a cell suffers due to toxicity. The sheer size of 6 that data is manually impossible to analyze, but the computational improvements have not only improved 7 produce that data but eased its analysis too. This data abundance also favors machine learning algorithms, 8 which benefit from an increasing number of data observations. Those algorithms can find trends in the 9 data, build models that predict them when these are missing, or even classify what is generally too 10 complex to understand.

#### 11 Sequencing RNA

12 RNA sequencing (RNA-seq) is a high-throughput sequencing technology that allows the identification and quantification of RNA in a cell or tissue. More than a decade old, this technique has become a 13 mainstream and cheap manner to analyze cells under different conditions at a subcellular dimension. The 14 15 method consists of 3 main steps: library preparation, sequencing, and data analysis (bioinformatics). Library preparation consists of the combination of all input items (library) required for sequencing. 16 Illumina will be used as an example, being one of the most popular sequencing technologies. First, total 17 RNA is extracted. To improve its sensitivity, mRNA (messenger RNA) enrichment or rRNA (ribosomal 18 19 RNA) depletion are performed. Afterward, this RNA is fragmented, as the sequencing technology is limited in the number of bases it can sequence. It is later retrotranscribed to cDNA (complementary 20 21 DNA), a much more stable molecule that can be amplified in a polymerase chain reaction (PCR). The following step is the ligation to sequencing adaptors, which will allow the fragments to be amplified, 22 23 identified, and sequenced. Once the cDNA material has been amplified by PCR, and both library concentration and fragment lengths have been verified (quality control), the library preparation is 24 25 complete. These PCR-amplified cDNA fragments can be sequenced by detecting and signaling each of the present nucleotides. The way the cDNA is sequenced depends on the technology used. The three more 26 27 popular technologies are Illumina, Pacific Biosciences, and Oxford Nanopore.

In the **Illumina** workflow, the cDNA molecules are fastened to a flow cell floor, where sequencing is performed by synthesis of the complementary strand of each of those fragments. The PCR is performed in this case through bridge amplification, where the fragment is bound from both ends to the floor, while the reverse strand is being produced. The resulting reverse fragment will also be replicated, leading to 32 multiple copies of both forward and reverse sequences of the same fragment (clonal amplification). After 33 that, all reverse fragments will be washed out, for the sequencing signal of that cluster of fragments to be 34 equal across all copies. For the sequencing step, a reverse strand will again be synthesized, but this time using fluorescent nucleotides, which when they successfully pair to one of the fragments, a particular 35 fluorescence will be emitted. Combining all fluorescent signals (50-500 base pairs) results in a read 36 sequence. In the simplest case scenario, each amplified fragment is sequenced once from one of its ends 37 38 (single-end sequencing). That same fragment may be sequenced from both ends (paired-end sequencing), providing a higher coverage of the sequence. 39

The **Pacific Biosciences** workflow adds sequencing adaptors that circularize the cDNA fragments. These fragments are added to a sequencing chip, which contains as many nanowells as fragments will be sequenced. At the bottom of such nanowells, an immobile polymerase will synthesize a new strand using fluorescent nucleotides (as in Illumina), which will be detected and generate a read (< 50kb).

For the **Oxford Nanopore** library preparation, molecules (aside from the adapter sequence ligation) are attached to a motor protein. When added to the flow cell, the motor protein will dock to a nanopore (embedded into an electrically resistant membrane) and will move a strand of the cDNA molecule through the pore. The passing of an ionic current through the nanopore will be disrupted differently based on the nucleotide going through the pore. This disruption will be measured in a signal trace, which generates reads between 1-10kb.

50 Different technologies are characterized by different sequencing read lengths. Although the older and 51 most used transcriptomics technologies are short-read sequencing, they have some limitations. Transcripts 52 expressed from the same or homologous genes tend to present different isoform sequences. Because such 53 transcripts share a pronounced proportion of the sequence, the genome mapping of the reads generated 54 from them (identifying the genomic locus of the original transcript) becomes complex.

55 The most common objective behind RNA-seq is differential expression analysis. Such analysis aims to 56 identify transcripts or genes that are differentially expressed across different conditions in a statistically significant manner. To do so, the analysis starts with the sequencing results in some form of preliminary 57 58 data, which needs to be pre-processed. The obtained sequenced fragments need to be linked/mapped to 59 their gene of origin. After all fragments have been located in the genome/transcriptome, a quantification informs about the number of amplified transcripts per gene. These raw quantities, though, are not directly 60 61 comparable across samples due to several factors (such as the total number of fragments per sample), 62 which leads to the need for normalization methods. Having the data pre-processed and normalized, 63 statistical methods are applied to evaluate how significant the expression level differences between 64 conditions are. Frequently, the significant results are not taken directly at face value using standard 65 thresholds (p < 0.05). Due to the high number of genes (~20000) and transcripts, the expected number of 66 false positives under such a threshold is substantially high (1000). Therefore, multiple testing correction, 67 such as the Benjamini–Hochberg one (or False Discovery Rate- FDR), are applied to the probability 68 values.

69 The relevance of the RNA changes is generally not related to the intrinsic function of those RNAs, but the 70 effect that those changes have when some of those RNAs (messenger RNAs) are translated to proteins. Protein level changes are important at a cellular level because the function to which those are associated 71 72 will be affected, be it from an excessive generation of their products to a complete lack thereof. The use of transcriptomics as a proxy for the study of proteins, though, is limited<sup>1</sup>. This is partially due to the 73 74 complex nature of post-transcriptional regulation. Post-transcriptional regulation refers to systems and 75 conditions that enhance or inhibit the translation of transcripts. Examples of this regulation are diverse: 76 from the structure itself of the transcript (%GC content and poly(A) tail), protein regulatory elements (RNA-binding proteins), the efficiency of translation (number of ribosomes simultaneously bound per 77 transcript), or even other transcripts (such as microRNAs). 78

#### 79 MicroRNA

80 Even though messenger RNAs have a key function as intermediaries for protein synthesis, most transcribed RNAs do not code for proteins. Around ~95% of the total expressed RNA consists of non-81 coding RNAs (ncRNA) essential for translation (ribosomal RNA and transfer RNA)<sup>2</sup>. Other ncRNAs 82 types, instead, inhibit that process: microRNAs<sup>3</sup>. MicroRNAs (miRNAs) are a class of non-coding RNAs 83 of around 22 nucleotides in length. Their main repressive mechanism in mammals relies on their seed 84 85 region, which is located between the second and seventh nucleotides. This region presents a sequence complementarity specific to one or more protein-coding RNAs (generally in their 3' UTR). Even though 86 87 this complementarity is per se highly probable to exist by chance due to the short seed sequence, a high number of the miRNA targets are conserved: they are more frequently found than expected by chance. In 88 89 addition, miRNAs tend to have several simultaneous targets (> 400 conserved targets per miRNA family).

90 Although the repressive effect by a single miRNA seed is relatively weak, targets that contain conserved 91 sites tend to present on average 4-5 of them for the same or different miRNAs, which leads to cumulative 92 repression. Even so, conservative sites are less frequent than non-conservative ones. MiRNAs do not 93 exert their function in a stand-alone form and require specific processing steps after transcription. The 94 first maturation step occurs in the nucleus, where the pri-miRNA contains at least one hairpin loop 95 structure. Drosha, an RNA-cutting enzyme, will separate the hairpin structure from the rest of the

96 transcript, leading to the pre-miRNA stage. After its nuclear export into the cytoplasm, RNase Dicer cuts 97 away the loop connecting both strands of the pre-miRNA, getting a miRNA duplex as a product. Even 98 though both strands can be functional, generally only one will be incorporated into the RNA-induced silencing complex (RISC). In it, a protein from the Argonaute family (Ago2) will bind the miRNA, the 99 100 latter functioning as a guide to find the complementary target. The Argonaute protein will then also bind the target RNA, which will lead to translation inhibition. There are several suggested mechanisms for 101 102 how this inhibition may occur. One of them is deadenylation of the target (via a deadenylase complex), leading to the 3' polyA end shortening (mRNA destabilization<sup>4</sup>). This shortening will eventually cause a 103 104 5' Cap loss, which enables the 5 to 3' mRNA decay. Thus, miRNA differential expression might affect mRNA target levels. In addition, translational repression mechanisms (decreased rate of mRNA 105 translation into proteins within cells) have also been described, such as Cap-40S initiation inhibition, 60S 106 107 ribosomal unit joining inhibition, elongation inhibition, and ribosome drop-off. Some studies suggest that the influence on translational efficiency is rather secondary, except for early embryonic states. In that 108 109 stage, a high correlation exists between translational efficiency and poly(A) tail length, thus the shortening of the 3' end does affect their translational process<sup>5</sup>. Other mechanisms, although also leading 110 to a lower protein output, do not affect translation directly. Examples of these are co-translational nascent 111 112 protein degradation and sequestration in P-bodies.

In a less common mechanism, some miRNAs exhibit a nearly perfect complementarity of their whole sequence to their target, where the latter is cleaved (primary mechanism in plants). This mechanism presents a high resemblance with silencing RNAs (siRNAs). Some evidence suggests that this miRNA cleavage mechanism might be responsible for the cleavage performed in experimental siRNA treatments, as species with defective Argonaute function present minimal response in such experiments<sup>6</sup>.

MiRNA post-transcriptional regulation is important both biologically (development, differentiation, cancer, and disease) and analytically. At the biological level, drastic changes in the expression or function of this RNA type can lead to unforeseen consequences to the transcripts it regulates. At a sequencing analysis level, RNA-Seq interpretation does not consider this regulation, thus opening the possibility of false conclusions due to the miRNA function. An extra layer of complexity is added when taking into account factors that regulate miRNAs themselves. Alongside typical transcript turnover, another novel non-coding transcript might be able to decrease miRNA's inhibitory potential: circular RNAs.

#### 125 Circular RNA

Messenger RNAs require some maturation steps before their translation: splicing (eliminate introns and join exons together), 5' cap addition, and 3' polyadenylation. Some transcripts mature using alternative

splicing mechanisms, which lead to a different structure: circular RNAs (circRNAs). CircRNAs are a result of back-splicing, a splicing event where both 5' and 3' ends join together, thus forming a backspliced junction, which leads to a circular structure<sup>7</sup>. These transcripts present neither 5' capping nor polyadenylation, thus alternative library preparations, such as ribo-depletion, are necessary for their identification. The circRNA sequence depends strictly on the alternative splicing event, where the resulting transcript may contain one or several exons, and sometimes even retained introns as well.

134 The most general hypothesized trigger for circularization is based on the introns flanking the circRNA135 sequence. The relatively long flanking introns have enriched ALU repeats that can base-pair to each

136 other, leading to the circularization (Figure 1: Circularization event).



Figure 1: Circularization event

To regulate such circularization, all examples found so far involve RNA-binding proteins (RBP). It has 138 been hypothesized that ADAR proteins (which are RBPs) are involved, as they can mutate the flanking 139 introns<sup>8</sup>, weakening and ultimately decreasing the probability of the intron duplex to occur<sup>9,10</sup>. Another 140 mechanism found involves Quaking (QKI), which is also an RBP, but in this case, it favors 141 circularization. It does so when two proteins bind to flanking QKI-binding sites in the linear transcript, 142 and their later dimerization gets both ends of the future BSJ in close contact. In the case of the 143 Drosophila muscleblind (mbl) gene<sup>11</sup>, the linear transcripts code for the MBL proteins (also an RBP), 144 145 which in turn stimulate the circularization of its linear transcripts when binding to their specific sites, thus impeding their translation in an autoregulatory or negative feedback system. 146

147 As a relatively new discovered type of transcript, several functions have been hypothesized, which are 148 generally related to (post-)transcriptional regulation. Even so, very few examples are known so far that provide evidence of those functions, which could be classified on whether they induce or inhibit eithertranscription or translation.

In the group of inhibition, as mentioned before, the formation (in the case of MBL, via protein-binding) of circular RNAs can act as a competing regulator of the linear transcript expression. As generally all of them are transcribed from genes that contain linear isoforms, this effect seems likely. In addition, the ability to function as an RBP sponge could have alternative molecular consequences. Evidence of such a sponging effect is not so straightforward: even though some RBPs have shown a higher cluster density in circRNAs<sup>12</sup>, their computational analysis seems to show a lower RBP-binding density<sup>13</sup>.

157 In the inducement function class, a few circular RNAs have been proven to be able to compete as targets 158 for miRNA binding. Some even present copies of the same sequence complementary to a miRNA seed target, allowing a single circular transcript to bind more than one miRNA. When that occurs, circRNAs 159 are said to function as miRNA sponges, adding a new layer of complexity to the post-transcriptional 160 161 regulation of miRNAs. Again, whether this function is the primordial one is still unclear, as circRNAs are 162 also expressed in organisms that lack miRNA-like repression pathways. In addition, the computational 163 search of circRNAs with enriched seed sites (compared to linear ones) delivered limited additional 164 examples of that function<sup>12,13</sup>. Interestingly, although circRNAs have been, since their discovery, considered a subtype of non-coding RNAs, new evidence via ribosome profiling suggests that they could 165 possess cap-independent protein-coding capabilities<sup>14,15</sup>. It was found that some of them, as some RNA 166 viruses, presented IRES (internal ribosome entry site) activity, which allows eIF3 recruitment (and, thus, 167 translation) independently of the 5'-cap modification<sup>16</sup>. This characteristic could imply that these 168 transcripts could be used as a medium-term source of translation/proteins, as they are more stable (and 169 170 therefore have a longer half-life) than their linear counterparts; which in turn could help mitigate the current transcriptomics-proteomics gap. Additionally, some circRNAs<sup>17</sup> can regulate their own parental 171 genes' regulation. A subgroup of circRNAs, named exon-intron circRNAs, is located at the nucleus, 172 173 where they promote the expression of their parental genes by interacting with the U1 small nuclear 174 ribonucleoprotein<sup>18</sup>. Thus, the ability to bind with RBPs might have both transcriptional and posttranscriptional functions. CircRNAs, therefore, might have a role in pre- and post-transcriptional 175 regulation, and because these RNAs can be analyzed via RNA-Seq, their influence on the protein 176 177 expression level can also be studied.

178

#### 179 **Proteomics**

180 Proteomics refers to the large-scale study of the "proteome", i.e., the set of all proteins present in a cell, 181 including their expression level, modifications, and interactions. The study of the proteome can be considered more complex than the study of the genome, as the proteome is a consequence of several 182 183 complex regulatory networks: starting with the genome itself, followed by the intermediary 184 transcriptome, in addition to the post-translational regulation, localization, and possible modifications. Often, instead of proteomics, transcriptomics is used as a proxy to evaluate the molecular changes that 185 186 could elucidate the phenotypical effects seen at larger systems (on a cellular or tissue level). As already 187 mentioned, RNA levels do not accurately represent either the expression or activity levels of their corresponding proteins, thus the need for proteomics for a more direct approach to the analysis of possible 188 189 molecular changes<sup>19</sup>.

190 Out of all the different techniques used in proteomics, the most popular one is mass spectrometry (MS), 191 from which other techniques also are derived. Mass spectrometry measures the mass-to-charge ratio of 192 ions. In proteomics, mass spectrometry can be used based on two different methods: top-down or bottomup proteomics. Top-down proteomics starts with a step of separation and quantitation of the proteins in 193 the sample, followed by MS for identification<sup>20</sup>. In bottom-up proteomics, the sample proteins are instead 194 195 first enzymatically digested, and then the product peptides are used as input for MS identification<sup>21</sup>. Top-196 down proteomics is usually used for the identification of a small subset of proteins, while bottom-up has 197 generally a larger sensitivity. In MS, the input peptides are ionized, which can then be split based on their 198 mass-to-charge ratio (m/z). This leads to a mass spectrum, which can be used to identify and quantify 199 proteins. Even so, some technological limitations still exist.

Even though several MS methods have been improving through the years, both reproducibility and 200 repeatability are low across both the same and different MS technologies<sup>22</sup>. Reproducibility, understood 201 as the similarity of results obtained from different technologies, tends to be lower than repeatability, 202 defined as the similarity of results obtained from the same technology, whose wide range comprises from 203 30-60%<sup>22</sup>. Current proteomics does not generally identify (in a large-scale quantitative manner) post-204 205 translational modifications, such as phosphorylation and ubiquitination, which inform of their activity status or their near-future degradation. In addition to that, its sensitivity is greatly limited<sup>23</sup>, as the mass 206 spectrometer generally quantifies around 2 to 3 thousand proteins (10-15%), far from the 20,386 207 manually annotated proteins in humans<sup>24</sup>. Related to sensitivity, dynamic range, referring to the log scale 208 between the most and least abundant proteins, is also one of the metrics aimed to be improved, since the 209 210 highly expressed proteins tend to completely overwhelm the spectrum to the expense of the low expressed 211 peptides.

212 In summary, proteomics is one of the technologies that most closely can represent the cell's environment, 213 but the technology itself is limited by its sensitivity. This combination of advantages and disadvantages is 214 curiously the opposite of the ones found in transcriptomics, where the technology allows for high sensitivity while only giving a proxy of the potential protein levels. In addition, considering all the post-215 transcriptional regulation networks that exist between the two, the ability to take all these considerations 216 simultaneously is manually unpractical. Therefore, the use of computational tools is a must at the 217 218 aforementioned level of analysis. Among the wide range of bioinformatics approaches, the selection of 219 machine learning algorithms (ML) was preferred due to three main characteristics. First, certain ML 220 models portray exceptional regression capabilities, crucial for the imputation performed in Chapter 4. In addition, a selection of ML models is also capable of classification tasks as exceptionally performed in 221 Chapter 5. The third characteristic is that most of these models permit the evaluation of the most 222 223 informative variables/features for the task at hand, which in turn inform the researcher to which degree each feature is related to the desired outcome (as performed in Chapters 4 and 5). 224

#### 225 Machine learning

Machine learning involves all algorithms that learn or adapt based on the input data<sup>25</sup>. Two of the main groups in which these algorithms can be categorized are supervised and unsupervised. The difference between the two groups relies on whether a target observation is present (supervised<sup>26</sup>) or absent (unsupervised), and therefore their aim changes too: prediction (i.e., whether a compound is toxic) versus grouping (i.e., grouping samples by similarity). Thus, supervised models aim to achieve the highest accuracy or lowest residual values (difference between observed and predicted values).

Supervised machine learning algorithms can also be categorized into two groups based on the nature of the target value: classification and regression. Classification models predict discrete target values (ill or healthy, dog breeds, etc.), while regression ones predict continuous values (height, salary, etc.). Supervised algorithms "learn" from the data by trying to minimize the residual size or error. Building a model to minimize the error in a limited dataset, though, has its limitations. The resulting accuracy can only be attributed to the specific dataset the model has been trained in, usually called the training dataset.

One of the simplest ways to evaluate the accuracy of a model for unseen data is to split the available data into two differing-size datasets (holdout method). The bigger one, which normally consists of around 75 to 80% of the data, is the one used to train the models (training dataset). The complementary small subset is the one used to test the accuracy of the model for unseen data, thus named testing dataset. The model might be tweaked by changing its hyperparameters (parameters that control the learning process) to try to optimize the accuracy of the testing dataset. Again, to avoid bias with such a testing dataset, the dataset 244 might be subsetted into a third group, the validation set. The latter is used to evaluate the model once the 245 hyperparameters have been optimally set. In this splitting data strategy, there is an inherent potential bias 246 based on which subset is selected as the testing dataset, even when chosen randomly.

247 For this reason, a more sophisticated method is normally used: k-fold cross-validation<sup>27</sup>. In crossvalidation, the dataset is also split into k folds, where k is the number of folds (subsets) the dataset is 248 divided into. The model is trained using all but one folds (k - 1) and tested on the remaining one (out-of-249 250 bag fold). The difference between the previous method and cross-validation is that in the latter, such 251 training and testing are performed with all possible combinations: in each iteration of the loop, one of the 252 k folds will be used as the testing dataset, giving a specific error or metric value. This results in a list of errors/values as large as the number of folds (k errors). Finally, this list is averaged to get a final accuracy 253 254 metric. Cross-validation, therefore, avoids a possible bias of a specific testing dataset by generating several of them, which gives a more accurate representation of how the model will predict with new input 255 256 data (unless all the original input data is somehow biased).

257 Several types of models exist for the prediction based on data learning. Although the number of 258 algorithms keeps increasing at a fast rate, some of the most popular algorithm groups existing in 259 supervised learning can be described based on their prediction strategy:

- Regression analysis: statistical methods that try to find the relationship between an independent variable (target to be predicted) and one or more dependent variables (features or variables). The most common one is linear regression<sup>28</sup>, in which the optimal solution is found using ordinary least squares. Different subtypes exist based on the nature of such relationship: linear (where the relationship between the dependent and independent variables is fit with a line), polynomial (based on a polynomial relationship of any degree), or logistic regression<sup>29</sup> (used to predict categorical data).
- Support-vector machines (SVM): as a set of methods used commonly for classification (but also for regression), they are defined as a non-probabilistic binary linear classifier<sup>30</sup>. Thus, these methods search to define the optimal line that separates two different classes by maximizing the margin (maximizing the distance between the line and the closest data point from each class).
   When the separation between two classes is non-linear, SVMs use the kernel trick, in which the data points are transformed to a higher dimension (feature space), where that classification may be performed by the use of a plane (or hyperplane, if the dimensions are above 3).

• **Decision trees**: models that apply decisions on each observation, whose final decision leads to 275 the target value. Each of the decisions (based on the features) are the so-called "branches", while each of the final predictions (after all the previous decisions) are named "leaves"<sup>31</sup>. Depending on
their target, decision tree models are either classification trees or regression trees. These methods
are very popular due to their simplicity and interpretability (the user can straightforwardly
visualize how the model reaches a prediction). They are though limited by their tendency to
overfit (improve training data metrics at the cost of testing data metrics). Advanced versions such
as *random forests* can correct for that, at the cost of their original interpretability.

- **Bayesian networks**: probabilistic models graphed via a directed acyclic graph (DAG)<sup>32</sup>. These 282 models use the probability of random variables, while also taking into account the dependencies 283 and combinations between those probabilities, to predict the probability of an outcome. The 284 outcome with the highest probability is the one predicted to be true. The relationships between 285 286 the different random variables are graphed using DAGs. In such graphs, the relationship between 287 two vertices (variables) and a single edge (arrow) is understood as the vertex as the end of the 288 edge presents a probability conditioned on the existence of the other vertex, and such a relationship is unidirectional. 289
- 290 Artificial neural networks: groups of models inspired by the biological neurological system, in which neurons are connected via dendrites<sup>33</sup>. In these models, there are at least three layers of 291 neurons or vertices: input, hidden, and output layers. The input neurons contain the input features 292 293 of the dataset. Each of those observations, before reaching the hidden layer(s), is transformed by 294 the weights and biases. These values are then used in the hidden layer, where a value is extracted 295 based on an activation function chosen beforehand (such as the softplus, rectified linear unit, or sigmoid functions). Again, these values are affected by additional weights before getting to the 296 output layer. In the output layer, the values from the hidden layers are summed together (and 297 corrected by a final bias). Adding additional hidden layers allows fitting data to even more 298 299 complicated non-linear relationships between inputs and outputs.
- 300 K-nearest neighbors (k-NN): in this type of model, which was originally designed for classification, a prediction is made based on how similar the unknown data to predict is to the 301 302 known training data. Different methods exist to evaluate the neighborhood of a data point, and some do so by clustering, such as Principal Component Analysis or hierarchical clustering. K 303 304 refers to the number of nearest neighbors to the predicted data point. Those k neighbors are differently classified, and the most popular class in that k-sized population is the one predicted. 305 306 Regression is also possible, where the prediction results from the average of the k-nearest neighbors values<sup>34</sup>. 307

#### 308 Aims and outline of the thesis

The aim of this thesis is twofold. First, we want to understand the changes that are a consequence of toxic treatments in the human organism at the molecular level. Namely, the elements of the post-transcriptional regulation responsible for the differences between the transcriptome and proteome at a specific point in time and tissue. Second, and in close relationship to the first aim, we want to further advance the development of data analysis tools that will allow us, and more generally, the experts responsible for risk assessment, to evaluate the effects of any treatment at the molecular level using high-throughput methods such as transcriptomics and proteomics.

316 Circular RNAs (circRNAs) have recently been shown to be deregulated in different disease scenarios. However, in toxicological conditions, these have still not been well studied. In Chapter 2, the expression 317 profiles of several circRNAs affected after cardiotoxicant treatments were analyzed. Was proceeded with 318 an additional validation on 12 differentially expressed circRNAs, and two were selected for further 319 investigation: circCDYL and circGNAS. CircCDYL levels were diminished after anthracycline treatment, 320 321 which putatively results in de-repression of particular sponged miRNAs. Therefore, the proteins encoded 322 by these miRNA-aimed mRNAs might be lowered. CircGNAS showed a significant cumulative rise across time after being treated with Amiodarone, which was experimentally validated. For both 323 circRNAs, the expression profile of their potential sponged miRNAs and their corresponding mRNA 324 targets were also analyzed. In conclusion, the presence of a regulatory axis involving circCDYL/miR-325 326 145-5p/TJP1 was discovered in cardiac cells upon anthracycline treatment, in addition to the validation of the upregulation of circGNAS due to a toxic dose of Amiodarone in different biological replicates, thus 327 328 indicating the reproducibility of circRNA studies through the use of iPSC-derived cardiomyocytes.

329 Transcriptomics is at the moment frequently applied as an analytical apparatus to examine the extent of 330 cell expression alterations between two phenotypes or between different conditions. Nevertheless, a substantial section of the significant changes detected in transcriptomics at the gene level is typically not 331 332 reliably identified at the protein level by proteomics. This weak correlation between the measured 333 transcriptome and proteome is perhaps largely due to post-transcriptional regulation, among which miRNA and circRNA have been suggested to play an important role. Hence, since both miRNA and 334 335 circRNA are also measured by transcriptomics, a model was proposed to be built in Chapter 3. It takes factors related to the post-transcriptional regulation into account to estimate, for each transcript, the 336 337 fraction of transcripts that would be available for translation. Using a dataset of cells exposed to diverse compounds, the model was evaluated to observe how and whether it was able to improve the correlation 338 339 between the assessed transcriptome and proteome expression level. The results showed that the model

improved the correlation for a subgroup of genes, possibly due to the regulation of several miRNAsacross the genome.

342 Proteins are often deemed the main biological component in charge of the different functions and 343 structures of a cell. Nevertheless, proteomics, the global study of all proteins which is frequently executed 344 by mass spectrometry, is restricted by its stochastic sampling and can only measure a limited number of 345 proteins per sample. Transcriptomics, which permits an extensive analysis of all expressed transcripts, is 346 regularly used as a surrogate. However, the transcript quantity does not present a high degree of 347 correlation with the subsequent protein quantity, notably due to the existence of several posttranscriptional regulatory mechanisms. In Chapter 4 is hypothesized that the absent protein values in 348 proteomics could be calculated using machine learning regression techniques, trained with numerous 349 variables extracted from transcriptomics, including previously identified translational regulatory elements 350 351 such as microRNAs and circular RNAs, among others. After taking into account different machine learning algorithms applied to two different splitting approaches, the random forest algorithm was 352 reported to be able to predict proteins in new samples out of several omics data with good accuracy. 353

In next-generation transcriptomics, differential expression analysis is established as one of the main 354 approaches to assess the outcomes of two biological conditions on the gene expression of different 355 356 biological samples. Nevertheless, the current statistical roadmaps offer very mild standardized filters for 357 the selection of differentially expressed genes, leading to a substantial number of false positives. Authors 358 typically incorporate their specific arbitrary thresholds, often reliant on their criteria and prospects of the 359 quantity of differently expressed genes to be attained. This leads to the inclusion of statistically 360 significant genes with expression profiles that, if individually examined, would not be considered biologically relevant to study further. In Chapter 5, we focused on developing AutoRel, a machine 361 learning model that incorporates not only the most conventional statistical assessments but also all the 362 363 complexities that distinguish biologically relevant changes based on manual examinations. AutoRel, 364 which categorizes each evaluated gene into "relevant", "irrelevant" or "dubious", informed of the most 365 crucial variables for the selection of relevant genes. The value of the number of replicates on the performance of the model was assessed, through the use of simulated datasets and the biological 366 367 interpretation of the chosen genes.

In **Chapter 6**, a general discussion of chapters 2 to 5 is presented. First, the discussion begins with a reflection on the dataset used for all chapters, HeCaToS, including its weaknesses and potential improved designs. Following the order of the chapters, the discussion continues with Chapter 2, focusing on the new quantification method used and suggestions for future validation experiments. It continues with the difficulties of formulating the complexity of post-transcriptional regulation found in Chapter 3. It 373 connects with one of its alternatives, the use of machine learning, which led to the work and later 374 publication of Chapter 4. The latter Chapter is also discussed, including its limitations regarding the use 375 of very stable conditions. The use of machine learning in Chapter 4 goes hand in hand with Chapter 5, 376 where there is an assessment of alternative methods to generate the training datasets. The General 377 Discussion ends with a short conclusion paragraph summarizing the work performed.

- 378 Lastly, an impact paragraph describes the potential future uses of the work presented in this manuscript.
- 379

#### 380 **Bibliography**

- Gygi, S. P., Rochon, Y., Franza, B. R. & Aebersold, R. Correlation between Protein and mRNA
   Abundance in Yeast. *Mol. Cell. Biol.* 19, 1720–1730 (1999).
- Westermann, A. J., Gorski, S. A. & Vogel, J. Dual RNA-seq of pathogen and host. *Nature Reviews Microbiology* vol. 10 618–630 (2012).
- 385 3. Baek, D. et al. The impact of microRNAs on protein output. Nature 455, 64–71 (2008).
- 386 4. Eichhorn, S. W. *et al.* MRNA Destabilization Is the dominant effect of mammalian microRNAs by
  387 the time substantial repression ensues. *Mol. Cell* 56, 104–115 (2014).
- Subtelny, A. O., Eichhorn, S. W., Chen, G. R., Sive, H. & Bartel, D. P. Poly(A)-tail profiling
  reveals an embryonic switch in translational control. *Nature* 508, 66–71 (2014).
- Bartel, D. P. MicroRNAs: Target Recognition and Regulatory Functions. *Cell* vol. 136 215–233
  (2009).
- 392 7. Barrett, S. P. & Salzman, J. Circular RNAs: Analysis, expression and potential functions. *Dev.*393 143, 1838–1847 (2016).
- Nishikura, K. Functions and regulation of RNA editing by ADAR deaminases. *Annual Review of Biochemistry* vol. 79 321–349 (2010).
- Ivanov, A. *et al.* Analysis of intron sequences reveals hallmarks of circular RNA biogenesis in animals. *Cell Rep.* 10, 170–177 (2015).
- Rybak-Wolf, A. *et al.* Circular RNAs in the Mammalian Brain Are Highly Abundant, Conserved,
  and Dynamically Expressed. *Mol. Cell* 58, 870–885 (2014).

- 400 11. Houseley, J. M. *et al.* Noncanonical RNAs from transcripts of the Drosophila muscleblind gene. *J.*401 *Hered.* 97, 253–260 (2006).
- 402 12. Guo, J. U., Agarwal, V., Guo, H. & Bartel, D. P. Expanded identification and characterization of
  403 mammalian circular RNAs. *Genome Biol.* 15, 409 (2014).
- You, X. *et al.* Neural circular RNAs are derived from synaptic genes and regulated by
  development and plasticity. *Nat. Neurosci.* 18, 603–610 (2015).
- 406 14. Pamudurti, N. R. et al. Translation of CircRNAs. Mol. Cell 66, 9-21.e7 (2017).
- 407 15. Legnini, I. *et al.* Circ-ZNF609 Is a Circular RNA that Can Be Translated and Functions in
  408 Myogenesis. *Mol. Cell* 66, 22-37.e9 (2017).
- 409 16. Meyer, K. D. *et al.* 5' UTR m6A Promotes Cap-Independent Translation. *Cell* 163, 999–1010
  410 (2015).
- 411 17. Hansen, T. B. *et al.* miRNA-dependent gene silencing involving Ago2-mediated cleavage of a
  412 circular antisense RNA. *EMBO J.* 30, 4414–4422 (2011).
- 413 18. Li, Z. *et al.* Exon-intron circular RNAs regulate transcription in the nucleus. *Nat. Struct. Mol.*414 *Biol.* 22, 256–264 (2015).
- 415 19. Cox, J. & Mann, M. Is Proteomics the New Genomics? Cell vol. 130 395–398 (2007).
- 416 20. Kelleher, N. L. Peer Reviewed: Top-Down Proteomics. Anal. Chem. 76, 196 A-203 A (2004).
- 417 21. Chait, B. T. Mass spectrometry: Bottom-up or top-down? Science vol. 314 65–66 (2006).
- Tabb, D. L. *et al.* Repeatability and reproducibility in proteomic identifications by liquid
  chromatography-tandem mass spectrometry. *J. Proteome Res.* 9, 761–776 (2010).
- Westont, A. D. & Hood, L. Systems Biology, Proteomics, and the Future of Health Care: Toward
  Predictive, Preventative, and Personalized Medicine. *Journal of Proteome Research* vol. 3 179–
  196 (2004).
- 423 24. reviewed: yes AND organism: "Homo sapiens (Human) [9606]" in UniProtKB.
- 424 https://www.uniprot.org/uniprot/?query=\*&fil=organism%3A%22Homo+sapiens+%28Human
- 425 %29+%5B9606%5D%22+AND+reviewed%3Ayes.
- 426 25. Tom Michael Mitchell. Machine Learning. McGraw-Hill Education 414

427 https://books.google.nl/books?

428 id=xOGAngEACAAJ&dq=isbn:0070428077&hl=nl&sa=X&redir esc=y (1997).

429	26.	Brewka, G. Artificial intelligence—a modern approach by Stuart Russell and Peter Norvig,
430		Prentice Hall. Series in Artificial Intelligence, Englewood Cliffs, NJ. Knowl. Eng. Rev. 11, 78-79
431		(1996).

- 432 27. Xu, Y. & Goodacre, R. On Splitting Training and Validation Set: A Comparative Study of Cross433 Validation, Bootstrap and Systematic Sampling for Estimating the Generalization Performance of
  434 Supervised Learning. J. Anal. Test. 2, 249–262 (2018).
- 435 28. Schneider, A., Hommel, G. & Blettner, M. Lineare regressionsanalyse Teil 14 der serie zur
  436 bewertung wissenschaftlicher publikationen. *Deutsches Arzteblatt* vol. 107 776–782 (2010).
- 437 29. Nick, T. G. & Campbell, K. M. Logistic regression. *Methods in molecular biology (Clifton, N.J.)*438 vol. 404 273–301 (2007).
- 439 30. Noble, W. S. What is a support vector machine? *Nature Biotechnology* vol. 24 1565–1567 (2006).
- 440 31. Myles, A. J., Feudale, R. N., Liu, Y., Woody, N. A. & Brown, S. D. An introduction to decision
  441 tree modeling. *J. Chemom.* 18, 275–285 (2004).
- 442 32. Friedman, N., Geiger, D. & Goldszmidt, M. Bayesian Network Classifiers. *Mach. Learn.* 29, 131–
  443 163 (1997).
- Abbod, M. F., Catto, J. W. F., Linkens, D. A. & Hamdy, F. C. Application of Artificial
  Intelligence to the Management of Urological Cancer. *Journal of Urology* vol. 178 1150–1156
  (2007).
- 447 34. Altman, N. S. An introduction to kernel and nearest-neighbor nonparametric regression. *Am. Stat.*448 46, 175–185 (1992).

449

450

# Chapter 2:

# Effect of cardiotoxicants on circRNAs and their function in posttranscriptional regulation

# <sup>1</sup> Effect of cardiotoxicants on circRNAs and

# <sup>2</sup> their function in post-transcriptional

# <sup>3</sup> regulation

- 4 Juan Ochoteco Asensio<sup>1</sup>, Jelmer Faber<sup>1</sup>, Jos Kleinjans<sup>1</sup>, Twan v.d. Beucken<sup>1</sup>, Florian
- 5 Caiment<sup>1</sup>
- 6 <sup>1</sup>Department of Toxicogenomics, School of Oncology and Developmental Biology
- 7 (GROW), Maastricht University, Maastricht, The Netherlands

#### 8 (Manuscript under preparation)

## 9 1. Abstract

10 Circular RNAs (circRNAs) have recently been shown to be deregulated in different disease scenarios. 11 However, in toxicological conditions, these have still not been well studied. In this study, we analyzed the expression profiles of several circRNAs affected after cardiotoxicant treatments. We proceeded with 12 13 additional validation on 12 differentially expressed circRNAs, and we selected two for further investigation: circCDYL and circGNAS. CircCDYL levels were decreased after anthracycline treatment, 14 which hypothetically results in de-repression of specific sponged miRNAs. Consequently, the proteins 15 16 encoded by these miRNA targeted mRNAs might be reduced. CircGNAS presented a significant 17 cumulative increase across time after being treated with Amiodarone, which we were able to validate 18 experimentally. For both circRNAs, we also analyzed the expression profile of their potential sponged 19 miRNAs and their corresponding mRNA targets. In summary, we discovered the presence of a regulatory axis involving circCDYL/miR-145-5p/TJP1 in cardiac cells upon anthracycline treatment, in addition to 20 21 validating the upregulation of circGNAS due to a toxic dose of Amiodarone in different biological replicates, thus demonstrating the reproducibility of circRNA studies by using iPSC-derived 22 cardiomyocytes. 23

## 24 2. Introduction

25 Circular RNAs (circRNAs) are a subtype of non-coding RNA with a circular structure, created by 26 covalent binding of the 3' and 5' ends of the RNA molecule<sup>1</sup>. Although circRNAs were firstly discovered 27 in other species several decades ago, their existence and potential functions in humans have been only 28 recently explored<sup>2</sup>. Around 140,000 potential circular RNAs have been predicted in human cells and 29 deposited in databases such as circBase<sup>3</sup>, though many of them tend to be tissue-specific<sup>4</sup>. Interestingly, 30 circRNAs are hypothesized to be stable due to their closed structure, protecting them from exonuclease activity. Although circRNA functions are still poorly understood, one of the most studied functions is 31 32 their ability to act as miRNA sponges. MiRNAs are another class of non-coding RNAs characterized by a small sequence length (~22 base pairs), whose main function is to bind to messenger RNAs (mRNAs) via 33 34 base pair complementarity, decreasing the translation efficiency of the latter. In this context, miRNA 35 sponging refers to the capability of some circRNAs to present one or more miRNA binding site, thus 36 sequestering (a.k.a. "sponging") those miRNAs, alleviating the translational inhibition of the coding 37 RNAs.

Recently, a limited number of research articles have been published concerning circular RNAs in the 38 39 context of cardiac toxic reactions. Although some focus on how these are regulated by RNA-binding proteins<sup>5</sup>, most are focused on their miRNA sponging capabilities<sup>6-8</sup>. The upregulation<sup>6,8</sup> or 40 downregulation<sup>7,9</sup> of an individual or multiple circRNAs have been hypothesized to be related to the toxic 41 42 effects of cardiotoxicant compounds. To this effect, these studies tended to focus on the relationship 43 between the expression of the circRNAs, expression of the miRNAs, and potential toxic effects. Even so, 44 some additional variables are considered to be primordial to understand the regulatory functions of these 45 circRNAs, such as the expression of the gene of origin of these circRNAs and the effect at the proteomics 46 level result of the combination of the different molecular factors.

For the detection or identification of circular RNAs in the RNA-Sequencing data, bioinformatics tools 47 such as CIRI-AS<sup>10</sup> and FUCHS<sup>11</sup> identify back-spliced junctions (BSJs), allowing the closure of the 48 circular structure, which are characteristic of these transcripts. CircExplorer2<sup>12</sup>, while including changes 49 to further discover circularized transcripts in both known and *de novo* assembled annotations, also relies 50 51 on the discovery of BSJs. Two main issues arise from such methodologies: first, the number of circular 52 RNAs found relies on the number of reads that span the BSJ, and because a read spanning a BSJ is a 53 relatively rare event, the probability to find them is highly dependent on the sequencing depth of the 54 experiment. Second, the de novo identification of these BSJs and their corresponding circular RNAs leads to either random identification names or naming based on genome location, which in turn also relies on 55 56 the specific genome version used. This complicates the association of these transcripts with their circular 57 RNA databases.

58 To study the potential role of circular RNAs in eliciting cardiotoxicity, RNA-Sequencing (RNA-Seq), 59 Small RNA Sequencing, and mass spectrometry were performed on induced Pluripotent Stem Cell-60 derived cardiac cells after exposure to several known cardiotoxicants. We were able to identify and 61 quantify circular RNAs with a transcriptome annotation file derived from the Ensembl and circBase 62 databases. After the quantification, a list of circRNAs that were differentially expressed after treatment 63 was obtained for each comparison. Among these, we investigated two specific circRNAs in the lab. First, circCDYL, whose main "sponged" miRNAs were found to be differentially expressed. Following the 64 post-transcriptional network, a fraction of the genes targeted by these miRNAs were also transcriptionally 65 and significantly affected. TJP1, a protein coded by one of the latter genes, showed a time-dependent 66 67 increase for TJP1 at the protein level. Lastly, we studied circGNAS, for which we were able to experimentally validate its significant increase, and we analyzed the related miRNAs, coding transcript 68 69 targets, and proteomics.

## 70 3. Methods

#### 71 3.1. Samples

The analyzed data consists of 3D microtissues containing stem-cell-derived cardiomyocytes and 72 73 fibroblasts in a 4:1 ratio from InSphero. These microtissues were exposed to eight compounds: 74 Fluorouracil (5FU), Amiodarone (AMI), Celecoxib (CEL), Docetaxel (DOC), Mitoxantrone (MXT), Paclitaxel (PTX), Doxorubicin (DOX), and Epirubicin (EPI), plus a fluctuating DMSO control (DF2). 75 The dosing profile was established via the use of the physiologically based pharmacokinetic (PBPK) 76 modeling software PK-Sim, to simulate exposure levels under physiological conditions<sup>13</sup>. For each 77 compound, there were 2 doses: Therapeutic and Toxic. The exposures were done in triplicates, and the 78 data extraction was performed at 8-time points: 0, 2, 8, 24, 72, 168\*, 240\*, and 336\* hours; resulting in 79 80 24 data points per dose (except for Doxorubicin and Epirubicin, which did not include a 0h timepoint). \* 81 The marked time points were not available for all compounds, as some of them were so toxic that 82 insufficient biological material remained.

#### 83 3.2. RNA Sequencing (RNA-Seq)

Total RNA from the exposed microtissues was isolated using the Qiagen AllPrep Universal Kit (Cat #80224). Ribo-depletion was achieved by using the Illumina RiboZero Gold kit (Cat #MRZG12324), and the libraries were prepared using the Lexogen SENSE total RNA kit (Cat #009.96). All libraries were then sequenced on an Illumina HiSeq 2000 at 100 bp paired-end at an average sequencing depth of 21 million reads (after Salmon quantification). The adaptors were removed through Trimmomatic version 0.33<sup>2</sup>. We used the following parameters: paired-end, ILLUMINACLIP: TruSeq3-PE.fa:2:30:10, LEADING:3, TRAILING:3, SLIDINGWINDOW:4:15, MINLEN:36, HEADCROP:12.

#### 91 3.3. (non-)coding RNA quantification and analysis

#### 92 3.3.1. CircRNA prediction with CIRI2

We used CIRI2<sup>14</sup> to predict and quantify potential circular RNAs in our sequencing data. Only circular
RNAs with at least 2 reads mapping to the back-splice junction (BSJ) were quantified. The total reads for

each circular RNA were obtained by adding both BSJ reads and non-BSJ reads.

96 3.3.1.1 Differential Expression Analysis with data quantified with CIRI2

97 Samples that presented a relatively low sequencing depth (less or equal to 25% of the average, ~5 million 98 mapped reads) in comparison with their in-group levels were excluded. We used as a control group for all 99 comparisons both the fluctuating DMSO samples and the timepoint 0h samples available (on each 100 compound experiment). The differential expression analysis was performed by DESeq2 (version 101  $(1.24.0)^{15}$ . We considered circular RNAs with a p. adjusted value < 0.05 significant.

**102** 3.3.2. CircRNA quantification with Salmon & circBase

103 The genome version used for all the transcriptomics analyses was the Genome Reference Consortium
104 Human Build 38 (GRCh38.p12). For the circRNA quantification, we used as a reference transcriptome

the hg19 circRNAs putative spliced sequence from circBase (Jul 2017 update)<sup>3</sup>.

We combined the transcriptomic libraries for both coding (all cDNA) and non-coding (all ncRNA)
transcripts from Ensembl (release 92, April 2018 Ensembl Archive)<sup>16</sup> with the library for mature circRNA
sequences from circBase into a single library, which we set as the global transcriptome reference to
Salmon<sup>17</sup>. We then used Salmon to quantify the RNA-Seq data.

#### 110 3.3.2.1 Differential Expression Analysis with data quantified with Salmon

111 We also utilized DeSeq2 as mentioned in 3.3.1.1. Afterward, we filtered the results according to several 112 thresholds. First, we required a p. adjusted value < 0.01 to filter by significance. Second, the median of all 113 samples needed to be non-zero, thus requiring that at least half of all samples expressed such molecule. 114 Third, the first quantile of both control & treatment groups also needed to be non-zero, which requires 115 that at least 75% of the samples in each group present an expression. Fourth, for a group to be considered significantly over-expressed against another, its first quantile needed to be higher than the third quantile 116 of the other group, thus requiring that at least 75% of the first group expression must be higher than the 117 75% of the other. The previously mentioned filters were inspired by the R-ODAF workflow<sup>18</sup>. Finally, the 118 119 difference between the 1/7 (14%) quantile and 6/7 (86%) quantile of the control group had to be smaller 120 than its median, therefore filtering out molecules that are highly sensitive to batch effects.

121 After grouping all the filtered results for all comparisons, we ranked the circular RNAs based on their 122 median expression in decreasing order, and subset the first 100 (the top 100 expressed). The changes were 123 then examined manually, selecting those that expressed both a very significant differential expression and/or a potential time-dependent increase or decrease. That examination resulted in a sub-selection of 10 124 **RNAs** (circBase hsa circ 0010791, hsa circ 0026129, 125 circular IDs): hsa circ 0034356, hsa circ 0055922, hsa circ 0060999, hsa circ 0076194, hsa circ 0078905, hsa circ 0090448, 126 hsa circ 0090904, and hsa circ 0102325. On top of those 10, we also added two circular RNAs known 127 to be expressed in iPSC-derived cardiomyocytes for further validation: circCDYL (differentially 128 129 expressed in CIRI2 quantification) and circSMARCA5<sup>19</sup>.

130

#### 131 3.4. Small RNA Sequencing and Quantification

Starting from the same total RNA isolated for the ribo-depleted libraries, an aliquot was size selected and ligated using the TruSeq Small RNA Library Prep Kit (Illumina®). After sequencing on the HiSeq 2500 at an average of 3.6 million reads per sample (after quantification), we quantified the resulting data using miRge2 (last change: 05/06/2018)<sup>3</sup>. miRge2 used the MiRBase database as the reference library (miRBase v22), IsomiRs were not considered for this analysis.

For miRNA data, we also used DeSeq2 as mentioned in 3.3.1.1. We also evaluated as significant those
circular RNAs with a p. adjusted value < 0.05.</li>

#### 139 3.5. MiRNA-gene interactions

140 To identify which genes were targeted by miRNAs, we decided to source such information from the 141 TargetScan predicted targets (version 7.1)<sup>20,21</sup>. We evaluated the importance of the inhibition based on the 142 'Cumulative weighted context++ score', which is the value used by default to sort those targets.

#### 143 3.6. MiRNA-circRNA table

For the prediction of miRNA-circRNA interactions, we used miRanda (version 3.3a, strict condition activated)<sup>22</sup> between human miRNA sequences from miRBase (version 22)<sup>23</sup> and human circRNA sequences from circBase (Jul 2017 update)<sup>3</sup>. This led to 21 million possible interactions between all the molecules, of which around 10 million were unique between both RNA types.

148 3.7. Experimental validation (qPCR)

#### 149 3.7.1. Cell culture

Induced Pluripotent Stem Cells (iPSCs, CARIM001A) were maintained on Geltrex (Gibco, Thermo
Fisher Scientific) coated 6-well tissue culture plates containing 2 ml of E8 Flex culture (Gibco, Thermo

152 Fisher Scientific) media. Geltrex was diluted at 1:100 using DMEM:F12 (Gibco, Thermo Fisher 153 Scientific). Plates were incubated with 5% CO2 at 37°C and cells were passaged when they reached 154 approximately 70% confluency or when spontaneous differentiation was observed in the center of the colonies. For passaging, cells were washed once with magnesium/calcium-free Dulbecco's phosphate-155 buffered saline (DBPS, Gibco, Thermo Fisher Scientific) and subsequently treated with 50 mM EDTA 156 diluted in magnesium/calcium-free DPBS and incubated at 5% CO2 at 37°C until colony edges started to 157 158 lift. To remove and break up colonies into smaller clumps, iPSCs were washed off the plate using 1 ml of RT E8 Flex containing 10 µM of Y-27632 ROCK inhibitor. iPSCs were subsequently plated on GelTrex-159 160 coated 6-well tissue culture plates and maintained in 1.5 ml of E8 Flex media supplemented with 10 µM of Y-27632 for the first 24h to aid cell adhesion and survival. After 24h the media was changed to 2 ml of 161 162 E8 Flex media.

163 SV40 immortalized human cardiomyocytes (ABMGood) were cultured in RPMI-1640 (Gibco, thermo 164 Fisher Scientific) supplemented with 10% FBS (<u>supplier</u>) and grown on applied extracellular matrix 165 (ABMGood) diluted 1:10 in 20 mM acetic acid. Cells were passaged when reaching approximately 90% 166 confluency. To passage, cells were washed once with DBPS and subsequently incubated with 40  $\mu$ l/cm2 167 accutase (Corning) for 5 minutes 5% CO2 at 37°C.

#### 168 3.7.2. Cardiomyocyte differentiation

Generation of cardiomyocytes from iPSCs was performed based on a protocol adapted from the one 169 published by Lian et al<sup>24</sup>. When iPSCs reached approximately 70% confluence, cells were dissociated in 170 RT E8 Flex supplemented with 10 µM of Y-27632. Cell clumps were plated on Geltrex-coated 12-well 171 tissue culture plates and maintained in 1 ml of E8 Flex media supplemented with 10 µM of Y-27632 for 172 the first 24h at a density of 140,000 cells per well. After 24h the cell culture media was changed for 1.5 173 ml of E8 Flex media. Once cells reached approximately 70 percent confluence (usually after 3 days), the 174 175 culture media was changed for 1.5 ml of RPMI-1640 + GlutaMAX (Gibco, Thermo Fisher Scientific), 176 supplemented with B27 minus insulin (Gibco, Thermo Fisher Scientific) and 10 µM of CHIR99021 177 (Tocris Bioscience) for 24h. After 24h of exposure, the media was changed to RPMI-1640 + GlutaMAX 178 supplemented with B27 minus insulin (differentiation day 0). On the third day after initial exposure to 179 CHIR99021, a combined media was prepared by collecting 50% of media from each well of the 12-well tissue culture plate and adding 50% of fresh RPMI-1640 + GlutaMAX supplemented with B27 minus 180 181 insulin. IWP2 (Tocris Bioscience) was added to a final concentration of 5  $\mu$ M. Left-over media was 182 aspirated from the cells and the combined media was added. 5 days after initial exposure to CHIR99021, the combined media was removed and replaced with fresh RPMI-1640 + GlutaMAX supplemented with 183 184 B27 minus insulin. On day 7, the media was replaced with fresh RPMI-1640 + GlutaMAX supplemented

with B27 plus insulin (Gibco, Thermo Fisher Scientific) and replaced every two days. On Day 14 a
cardiomyocyte selection medium was prepared by adding L-(+)-Lactic acid (Santa Cruz Biotechnology)
to RPMI-1640 without glucose (Gibco, Thermo Fisher Scientific) to a final concentration of 4 mM, and 2
ml were added to the cardiomyocytes after washing twice with magnesium/calcium-free DBPS. After
three days the media was replaced with fresh RPMI-1640 + GlutaMAX supplemented with B27 plus
insulin and replaced every two days until cardiomyocytes were used in experiments.

#### **191** 3.7.3. Validation of Amiodarone Toxic dose effect on circGNAS

192 Cardiomyocytes were washed three times with RT magnesium/calcium-free DPBS before 1 ml per well of pre-warmed TrypLe (Gibco, Thermo-Fisher Scientific) supplemented with 0.5 U/ml Liberase (Roche) 193 194 and incubated at 5% CO2 at 37°C for 5 minutes. Cell clumps were dissociated by repetitive pipetting followed by another 5 minutes of incubation at 5% CO2 at 37°C. Cell clumps were again dissociated by 195 196 repetitive pipetting and pelleted at 250 x g for 2.5 minutes. The supernatant was discarded and the pellet was redissolved in fresh RPMI-1640 + GlutaMAX supplemented with B27 plus insulin supplemented 197 198 with 10 µM of Y-27632. Cells were seeded on Geltrex-coated 48-well tissue culture plates and maintained in 1 ml of E8 Flex media supplemented with 10 µM of Y-27632 for the first 24h at a density 199 of 300.000 cells per well. Cells were left over the weekend before Amiodarone treatment commenced. 200

Amiodarone hydrochloride was aseptically dissolved to a stock concentration of 10 mM in DMSO (Sigma-Aldrich) and stored at -20 °C until needed. Amiodarone dilutions were prepared using prewarmed RPMI-1640 + GlutaMAX supplemented with B27 plus insulin and cells were washed once with RT PBS before 500 µL of amiodarone-containing medium was added to the cells.

#### 205 3.7.4. RNA extraction and cDNA synthesis

After 24h of exposure, the culture medium was discarded and 500  $\mu$ L of Qiazol (Qiagen) was added to each well to lyse the cells. RNA extraction was carried out according to the manufacturer's protocol and RNA yield and quality were assessed using a Nanodrop 1000 spectrophotometer (Thermo-Fisher Scientific). cDNA was synthesized using the iScript cDNA Synthesis Kit (Bio-Rad Laboratories) according to the manufacturer's protocol.

#### 211 3.7.5. CircRNA primer design

cDNA FASTA sequences were extracted for circRNAs of interest and primers were designed such that the BSJ (back-spliced junction) would be amplified. This was done so by designing a primer at each end of the sequence, and each of them was directed to the end closest to them. In this manner, only transcripts that contain a BSJ between the start and the end of that sequence will be amplified. Primer sequences can be found in supplementary table 1. The samples used for validation by qPCR were not the ones used forthe bioinformatics analysis.

#### 218 3.7.6. SYBR Green RT-qPCR Validation

219 Real-time quantitative PCR (RT-qPCR) reactions were run in triplicates using the CFX Connect Real-

220 Time PCR Detection System (Bio-Rad laboratories) and data were collected with the Bio-Rad CFX

manager software (v3.1, Bio-Rad laboratories). cDNA amplification was measured using the IQ Sybr
 Green Supermix kit (Bio-Rad laboratories) and the abovementioned custom-designed probes. Expression

223 levels were calculated using the 2(-Delta Delta C(T)) method with  $\beta$ -Actin as endogenous control.

## 224 3.8. CIRI2 quantification results of validated circRNAs

To investigate whether the validated circRNAs were also quantified by CIRI2, we retrieved the 225 spliced sequence from circBase for each validated circRNA. Moreover, we used BLASTN 226 (Nucleotide-Nucleotide BLAST 2.12.0+)<sup>25</sup> to search the location of those sequences in our 227 genome FASTA file. To filter across possible homolog sequences in different chromosomes, we 228 229 only retrieved those coordinates that matched the chromosome where the circRNA locus was present (according to circBase). We pooled together all CIRI2 quantification results, searching 230 for predicted circRNAs located in the same loci. In that regard, only predicted circRNAs that 231 matched the chromosome number, circRNA start (locus of the first nucleotide), circRNA end 232 233 (locus of the last nucleotide), and strand with the BLASTN results were considered both experimentally validated and quantified by CIRI2. 234

## 235 3.9. Additional differential expression analyses

As performed with the miRNA data, we utilized the DESeq2 pipeline with a threshold of p.
adjusted value < 0.05 to analyze the potential differential expression of both genes and (linear)</li>
transcripts.

## 239 3.10. Proteomics analysis

Protein samples were isolated and diluted to a concentration below 0.2M. Subsequently, they were subject to digestion by trypsin and then cleaned up with the usage of Sep-Pak tC18 cartridges (Waters) according to the manufacturer's instructions. The samples were then dried using a vacuum centrifuge and quantified by the usage of an Orbitrap Fusion mass spectrometer (Thermo Fisher Scientific), coupled to a NanoLC-2D HPLC system (Eksigent). The raw MS data processing was performed with Genedata Expressionist software (v.11.0). Noise-reduction and normalization were applied to the LC-MS peaks, whose properties were later collected (m/z and RT boundaries, m/z and RT center values, intensity). The annotation of the individual MS/MS spectra was obtained using Mascot 2.6. The peak clusters were grouped using protein interference (peptide and protein annotations), and the protein intensities obtainment was performed with the Hi3 method. Proteomics data were normalized using a modified median of medians approach<sup>26</sup>.

## 252 4. Results

To assess the putative effect of toxicants on circRNAs expression, we quantified relevant elements to the circRNA post-transcriptional regulation (circRNAs, linear RNAs at both transcript and gene levels, microRNAs, and proteins) in both therapeutic and toxic doses of eight compounds (5FU, AMI, CEL, DOC, DOX, EPI, MXT & PTX) across 7 timepoints (2, 8, 24, 72, 168, 240, and 336h) using omics technologies.

To categorize circRNAs in our ribo-depleted data, we used two different methodologies. First, we predicted *in silico* the presence of circRNAs using CIRI2, which predicts and quantifies circular RNAs from the transcriptomics input data based on their back-spliced junction (BSJ). Alternatively, we mapped and quantified all sequenced reads from the ribo-depleted libraries on circRNA sequences extracted from circBase (remapping). We thus describe below the results found in both.

## 264 4.1. Differentially expressed circular RNAs

#### 265 4.1.1. In silico prediction and quantification of circRNAs using CIRI2

We performed a circRNA in silico prediction on all treatment samples and assessed their 266 differential expression using as controls all timepoint 0h samples added to control DMSO 267 samples. Depending on the treatment, this analysis detected between 11 to ~16 thousand unique 268 circRNAs (Table 1). As expected, the number of constitutive circRNAs (circular RNAs 269 quantified in all samples in the comparison) was very low. The number (and proportion) of 270 Differentially Expressed CircRNAs (DECs) was the highest in the anthracycline treatments 271 (Doxorubicin and Epirubicin). Even so, the percentage of samples where the DECs were 272 detected was on average between 22 and 10% across all treatments, which was substantially low 273 274 considering that most DECs were downregulated (higher expression in control samples). In other

words, the majority of differentially expressed circRNAs were not detected in either treatment or control groups for most of the samples.

277 Table 1: Summary of the CIRI2 and DESeq2 results. 'N quantified circRNA' refers to the number of unique circRNAs that were

278 detected by CIRI2 in any sample involved in the comparison between that treatment and the control. 'N quantified constitutive

279 circRNAs' refers to the number of unique circRNAs that were detected by CIRI2 in all samples involved in the comparison

280 between that treatment and the control. 'N DECs' refers to the number of Differentially Expressed CircRNAs (DECs) in the

**281** specific treatment when compared to the control samples (p. adjusted value < 0.05). '% DECs' refers to the percentage of DECs

taking the 'N quantified circRNA' as the total/denominator. 'Avg % samples a DEC is quantified' refers to the percentage of

283 samples that detect a circRNA on average when it is differentially expressed (i.e., if the comparison in Epi The spanned a total

284 of 100 samples, and we analyzed the quantification values for an average DEC, only 18 samples out of the 100 would detect that

285 *circular RNA*).

Treatment	Ν	Ν	N DECs	% DECs	Avg % samples a
	quantified circRNA	quantified constitutive circRNAs			DEC is quantified
Epi_The	13356	1	1135	8.50%	18.37%
Dox_The	11480	1	892	7.77%	19.71%
Dox_Tox	11201	1	764	6.82%	22.04%
Epi_Tox	11713	2	559	4.77%	21.23%
CEL_The	11002	7	391	3.55%	15.61%
MXT_Tox	11553	6	377	3.26%	13.81%
PTX_Tox	11359	6	351	3.09%	14.67%
CEL_Tox	11243	6	303	2.70%	14.10%
PTX_The	11609	8	294	2.53%	12.66%
MXT_The	12256	8	277	2.26%	13.29%
DOC_The	11916	7	248	2.08%	13.94%
DOC_Tox	13595	8	212	1.56%	12.58%
5FU_Tox	14532	10	173	1.19%	10.07%
AMI_The	13596	6	160	1.18%	13.12%
5FU_The	13950	9	150	1.08%	11.01%
AMI_Tox	15695	10	138	0.88%	10.18%

286

To further investigate this high inconsistency in DEC quantification, we pooled the quantification of all the samples and analyzed the overlap in circRNA detection across all samples. The consequent results (Figure 1) showed that 59% of all detected circRNAs (median = 1) were only quantified in a single sample, 90% were quantified in 8 samples or less, and only one circRNA (circRMRP, or hsa\_circ\_0001853) was quantified in all 309 samples.

#### **292** *4.1.1.1 circCDYL (hsa\_circ\_0008285)*

Across the differential expression results, we searched for some of the most recently studied circRNAs. Among them, we found circCDYL differentially expressed in 3 anthracycline treatments: the therapeutic (Figure 1) and toxic dose of Doxorubicin, and the therapeutic dose of Epirubicin. All three treatments showed strong down-regulation of circCDYL expression (log2 fold change < -2.5). The sequence predicted by CIRI2 for circCDYL matched perfectly to hsa circ 0008285<sup>27</sup> from circBase.

298



<sup>299</sup> 

To confirm that the change shown was not due to a down-regulation at the gene level, we quantified (using the remapping strategy in **3.3.2**) the expression of the linear CDYL mRNA transcripts. The results (Figure 2) showed no down-regulation for any of the 3 treatments. Instead, the therapeutic dose treatment of both Doxorubicin and Epirubicin presented a significant up-regulation of the linear transcripts. This suggests that the reduced circCDYL levels we observed upon doxorubicin treatment were not due to transcriptional repression of the CDYL gene.

**Figure 1: circCDYL (circRNA\_ID: 6:4891713\4892379) quantification for the Control versus the therapeutic dose of Doxorubicin (p.adj. < 0.001) comparison.** Each dark grey bar represents a control sample, while each pink bar represents a treated sample. The quantification shown was the result of the normalization performed by the DESeq2 pipeline. The 'circRNA\_ID' is the identification given by CIRI2 based on the genome location. The names in the X-axis refer to the sample names: the first term indicates the treated compound, the second term refers to the dose used, the third term refers to the timepoint (in hours), and the fourth term signals the triplicate number.





Figure 2: CDYL (ENSG00000153046) quantification by Salmon (circBase) for the Control versus the therapeutic dose of
 Doxorubicin (p.adj. < 0.01) comparison. Each dark grey bar represents a control sample, while each pink bar represents a</li>
 treated sample. The quantification shown was the result of the normalization performed by the DESeq2 pipeline.

#### **316** 4.1.2. Exhaustive quantification using Salmon & circBase

Since *in silico* prediction is solely based on the detection of BSJ sequences, which is substantially 317 dependent on coverage, we decided to re-quantify all transcripts by remapping against a common 318 319 reference made of all linear transcripts (both protein and non-coding) and all circRNAs obtained from circBase. The detection range was between 39 and 57 thousand unique quantified circRNAs across all 320 treatments (Table 2). As expected, the number of constitutive circRNAs was higher (576 on average) 321 322 compared to the *in-silico* prediction (6 on average). The number (and proportion) of DECs was the 323 highest in the anthracycline treatments (Doxorubicin and Epirubicin), similarly to what was observed in 324 Table 1. Even so, the percentage of DECs is higher (12.23% on average) than the one *in silico* prediction 325 (3.33% on average). This might be due to the higher percentage of samples where the DECs were 326 detected (29.07% on average compared to 14.77% on average in CIRI2), where the increased sensitivity 327 might help elucidate the possible changes when comparing the treatment and control groups.

**328** Table 2: Summary of the Salmon + circBase and DESeq2 results. 'N quantified circRNA' refers to the number of unique

circRNAs that were detected by Salmon in any sample involved in the comparison between that treatment and the control. 'N
 quantified constitutive circRNAs' refers to the number of unique circRNAs that were detected by Salmon in all samples involved

quantified constitutive circRNAs' refers to the number of unique circRNAs that were detected by Salmon in all samples involved
 in the comparison between that treatment and the control. 'N DECs' refers to the number of Differentially Expressed CircRNAs

**332** (DECs) in a specific treatment when compared to the control samples (p. adjusted value < 0.05). " DECs' refers to the

**333** percentage of DECs taking the 'N quantified circRNA' as the total/denominator. 'Avg % samples a DEC is quantified' refers to

334 *the percentage of samples where a DEC is detected on average.* 

Treatment	N quantified	N quantified	N DECs	% DECs	Avg % samples a DEC is quantified
	circRNA	constitutive circRNA			·
Epi_Tox	39690	94	12438	31.34%	24.62%
Dox_Tox	39323	109	10479	26.65%	36.63%
Dox_The	39645	121	9489	23.93%	35.41%
Epi_The	41529	166	9454	22.76%	35.07%
PTX_Tox	43729	575	6949	15.89%	23.66%
CEL_The	41140	632	4910	11.93%	25.53%
MXT_Tox	43278	713	4662	10.77%	21.59%
DOC_The	39091	683	3207	8.20%	21.04%
AMI_The	40121	645	2615	6.52%	26.41%
5FU_Tox	44673	819	2745	6.14%	27.83%
PTX_The	57661	746	3542	6.14%	33.22%
DOC_Tox	46112	813	2666	5.78%	30.26%
MXT_The	53610	745	2993	5.58%	28.53%
CEL_Tox	54307	725	2689	4.95%	31.02%
AMI_Tox	46755	791	2244	4.80%	24.99%
5FU_The	45670	839	1920	4.20%	39.33%

#### 335

To further investigate the increased sensitivity found in this method, we pooled the quantification of all the samples and analyzed the overlap in circRNA detection across all samples. The resulting data showed that only 8% of all detected circRNAs were quantified in a single sample, 50% (median) were quantified in 17 samples, and 9 circRNAs were quantified in all samples.

#### 340 4.1.3. Selection of circRNAs for experimental validation

We selected for experimental validation 10 DECs across all different treatments based on expression, 341 and time effect: hsa circ 0010791 (circCDC42), hsa circ 0026129 342 differential expression, (circTUBA1A), hsa circ 0034356 (circACTC1), hsa circ 0055922 (circFHL2), hsa circ 0060999 343 hsa circ 0076194 (circCDKN1A), hsa circ 0078905 344 (circGNAS), (intergenic circRNA), hsa circ 0090448 (circTIMP1), hsa circ 0090904 (circMSN), and hsa circ 0102325 (circHIF1A). 345 hsa circ 0008285 (circCDYL) and hsa circ 0001445 (circSMARCA5) have been shown to increase 346

their expression after cardiac differentiation<sup>19</sup> and were thus added for the experimental validation as
positive controls.

#### **349** *4.1.3.1 circGNAS (hsa\_circ\_0060999)*

- 350 CircGNAS was an example of a highly expressed circRNA with a significant differential expression and
- 351 time-effect accumulation when treated with a toxic dose of amiodarone. As seen in Figure 5, the circRNA
- 352 expression increased significantly to levels not found in control conditions. Moreover, the increase was
- 353 cumulative until 168 hours. At 240 hours (10 days), the expression appeared to return to control levels.



#### 354

**Figure 3: circGNAS (hsa\_circ\_0060999) quantification for the Control versus the toxic dose of Amiodarone comparison** (p.adj. < 0.001). Each dark grey bar represents a control sample, each light grey bar represents a control sample differentiated in the same conditions as the treatment, and each pink bar represents a treated sample. The quantification shown was the result of the normalization performed by the DESeq2 pipeline. The names in the X-axis refer to the sample names: the first term indicates the treated compound, the second term refers to the dose used, the third term refers to the time-point (in hours), and the fourth term signals the triplicate number.

As performed for circCDYL, we quantified and analyzed the expression of the linear transcripts originated from the GNAS (Guanine Nucleotide binding protein, Alpha Stimulating activity polypeptide) gene to investigate whether the incremental expression shown was due to an up-regulation at the gene level. The results (Figure 4) showed no significant change for the analyzed treatment (toxic dose of Amiodarone).





Figure 4: GNAS (ENSG0000087460) quantification by Salmon (circBase) for the Control versus the toxic dose of
 Amiodarone (p.adj. > 0.05) comparison. Each dark grey bar represents a control sample, each light grey bar represents a
 control sample differentiated in the same conditions as the treatment, and each pink bar represents a treated sample. The
 quantification shown was the result of the normalization performed by the DESeq2 pipeline.

#### 371 4.2. qPCR validation

#### 372 4.2.1. circRNA existence validation

To investigate whether the predicted circRNAs exist and are expressed in cardiomyocytes, two batches of cardiomyocytes were derived from iPSCs and checked for expression of our circRNAs of interest. As previously mentioned, circCDYL and circSMARCA5 have already been shown to be expressed in iPSCderived cardiomyocytes and were therefore included in these experiments as controls. Furthermore, to investigate whether our circRNAs of interest are cardiomyocyte-enriched, their expression was compared to the initial iPSCs level.

The results confirmed the expression of three (out of 12) circRNAs existed in myocardial cells after iPSC differentiation (Figure 5A and Figure 5B): circGNAS, circCDYL and circSMARCA5. Aside from the existence of circCDYL and circSMARCA5, our results also confirmed previous results that both circRNA expressions increased along with cardiac differentiation. In addition, we also observed a novel relationship between circGNAS expression and cardiac differentiation. The six transcripts (circTUBA1A, circTIMP1, circHIF1A, circFHL2, circACTC1, and hsa\_circ\_0078905) that are labelled in both Figures (Figure 5A and Figure 5B) were also quantified, but the results in gel showed several bands, thus the
quantification values were a result of several circRNA molecules, which most probably were isoforms of
the ones we were studying. In the case of circCDC42, the expression was only detected for the second
batch (Figure 5B). CircCDKN1A and circMSN were not validated through qPCR.

389

In addition, we also observed a novel relationship between circGNAS expression and cardiac differentiation. For six other circRNAs of interest (circTUBA1A, circTIMP1, circHIF1A, circFHL2, circACTC1, and hsa\_circ\_0078905) multiple PCR products were formed. this is possibly caused by other circular isoforms but we were not able to rule out aspecific amplification and were therefore excluded from further analysis. Furthermore, circCDC42 was only detected in the second batch of iPSC-CMs. CircCDKN1A and circMSN were not validated though qPCR.



Figure 5: Relative expression of circRNAs by qPCR. A: 9 circRNAs amplified in the first batch. B: 10 circRNAs (9 + circCDC42) amplified in the second batch. In the vertical axis, Relative FC refers to the relative fold change of expression. In the horizontal axis, each couple of bars is identified by their circular RNA name based on their gene of origin. For the last circRNA, we used the circBase ID due to its intergenic nature. 'CM' stands for CardioMyocytes and 'iPSC' for induced Pluripotent Stem Cells.

Taking into account the nine circRNAs (and potential isoforms) we experimentally quantified, we
evaluated which circRNAs were detected in at least one sample by each of the two methodologies used.
This required the cross-reference between gene loci (CIRI2) and circBase IDs (remapping), that is, we

- 406 transformed the circRNA names from one methodology to another. We observed that only two circRNAs
- 407 (circHIF1A and circCDYL) were identified by CIRI2, while all nine validated circRNAs were detected
- 408 by the remapping strategy.
- 409 4.2.2. Validation of Amiodarone effect on circGNAS
- 410



412 Figure 6: Relative expression of circGNAS upon treatment with three different doses of Amiodarone. Cmax concentration

413 represents the therapeutically active average plasma maximum concentration values derived from recommended therapeutic
414 doses upon a single-dose administration.

After observing differential circGNAS expression in our transcriptomics analysis upon amiodarone 415 treatment, experimental validation was carried out. To investigate a possible dynamic dose-response link 416 417 between circGNAS expression and amiodarone exposure, iPSC-derived cardiomyocytes were exposed to 418 three different Cmax-based concentrations of amiodarone for 24 hours (Figure 6). Cmax concentration 419 represents the therapeutically active average plasma maximum concentration values derived from 420 recommended therapeutic doses upon a single-dose administration. In the case of amiodarone, the 421 selected concentration (0.807  $\mu$ M) was obtained from a study carried out by Wink et al<sup>28</sup>. The three 422 different concentrations of amiodarone were selected at 1, 5, and 10 times Cmax (0.807, 4.035, and 8.07 µM respectively). After 24 hours of exposure to amiodarone, iPSC-derived cardiomyocytes were 423 harvested for RNA and circGNAS expression was quantified using RT-qPCR (Figure 6). CircGNAS 424 425 showed a dose-dependent increase expression upon Amiodarone exposure.

#### 426 4.3. Investigation of the putative function of circRNAs

427 Next, we wanted to study the ability of circCDYL and circGNAS to act as miRNA sponges. Having 428 access to different omics datasets obtained from the exact same samples, allowed us to robustly 429 investigate the effect of circRNAs on miRNA expression levels and expression of miRNA-target mRNAs 430 themselves as well as their encoded proteins.

431 First, we investigated which miRNAs could potentially bind to these circRNAs and whether their

- 432 expression was also affected by the compounds. In addition, we examined the potential mRNA targets of
- 433 these miRNAs and their expression changes. Finally, we examined potential changes at the proteomics
- 434 level as a consequence of post-transcriptional dysregulation.

#### 435 4.4. CircCDYL







#### 443 4.4.1. MiRNAs associated with circCDYL

- 444 Due to circCDYL being a previously studied circular RNA, we selected miRNAs that had been
- previously hypothesized to target this circRNA. The selected miRNAs were: hsa-miR-190a-3p<sup>29</sup>, hsa-

446 miR-185-5p<sup>30</sup>, hsa-miR-4793-5p<sup>31</sup>, hsa-miR-150-5p<sup>32</sup>, hsa-miR-892a<sup>33</sup>, hsa-miR-328-3p<sup>33</sup>, hsa-miR-92b-447  $3p^{34}$ , hsa-miR-145-5p<sup>35</sup>, and hsa-miR-1180-3p<sup>36</sup>.

First, we calculated the average expression for each of the miRNAs across the 3 treatments of interest (therapeutic and toxic dose of Doxorubicin, and therapeutic dose of Epirubicin) to filter out miRNAs with low expression, as their influence on their target's translation would be minimal. As a result, we excluded 5 miRNAs (Figure 8) with minimal (less than one normalized read) or no expression on average: hsamiR-185-5p, hsa-miR-190a-3p, hsa-miR-150-5p, hsa-miR-892a, and hsa-miR-4793-5p. Therefore, only the expression of the remaining 4 miRNAs (hsa-miR-92b-3p, hsa-miR-145-5p, hsa-miR-328-3p, and hsamiR-1180-3p) was analyzed.





456

457 *Figure 8: Bar plot of the average quantified expression of miRNAs associated with circCDYL*. The normalization and resulting
 458 values were obtained as a result of using the DESeq2 pipeline.

For each miRNA, we obtained the p. adjusted value, while also plotting the normalized expression of all samples involved in the comparison. Hsa-miR-92b-3p (Figure 9A) was the only miRNA that showed a significant decrease as a function of time after exposure to doxorubicin at a therapeutic dose. For the other three miRNAs (Figure 9B, Figure 9C, and Figure 9D), the therapeutic dose of Doxorubicin and Epirubicin always presented a significant up-regulation, while upon the toxic dose of Doxorubicin, only hsa-miR-328-3p was significantly up-regulated in a time-dependent manner (Figure 9C). Interestingly
enough, the latter miRNA also showed a time-dependent increase in expression for all three treatments
(Supplementary Results).





## 470

475 shows the miRNA ID next to the treatment name.

<sup>471</sup> Figure 9: Normalized read counts of four miRNAs by miRge2 for the Control versus the therapeutic dose of Doxorubicin

**<sup>472</sup>** comparison. A: hsa-miR-92b-3p (p.adj. < 0.01). B: hsa-miR-145-5p (p.adj. < 0.05). C: hsa-miR-328-3p (p.adj. < 0.001). D:

<sup>473</sup> hsa-miR-1180-3p (p.adj. < 0.001). Each dark grey bar represents a control sample, while each pink bar represents a treated

<sup>474</sup> sample. The quantification shown was the result of the normalization performed by the DESeq2 pipeline. The title of the plot

#### 477 4.4.2. Target transcript expression

We investigated which targets were associated in earlier studies with the previously mentioned expressed miRNAs. We excluded the targets of hsa-miR-92b-3p, as it was found to be able to degrade circCDYL, and thus the relationship between both entities could not be classified as the classical "miRNA sponging". The differential expression analysis of the target genes (Table 3) showed that all were significantly downregulated for the therapeutic dose of Epirubicin (Figure 10A and Figure 10B), while only YAP1 was significantly downregulated in all 3 treatments (Figure 10D).

- **484** *Table 3: Overall statistics of known circRNA-regulated genes. P. adjusted values in bold are significant (p. adj. value < 0.05).*
- 485 <sup>1</sup>: Paralog TJP1 gene (ENSG00000104067).<sup>2</sup>: Paralog TJP1 gene (ENSG00000277401).<sup>3</sup>: Weighted average of both TJP1 gene
- **486** paralogs (weights were derived from the sum of the three mean expressions for each paralog). L2FC: Log2 Fold Change.

Gene Name	Mean expression	Mean expression	Mean expression	P. adj. value	P. adj. value	P. adj. value
	(Dox_The	(Dox_Tox)	(Epi_The)	(L2FC)	(L2FC)	(L2FC)
	)			[Dox_The	[Dox_Tox]	[Epi_The]
				]		
$TJP1^{1}$	666	689	695	2.24x10 <sup>-1</sup> (-	4.53x10 <sup>-1</sup> (-	6.86x10 <sup>-3</sup>
				0.81)	0.55)	(-1.38)
$TJP1^2$	722	731	770	4.72x10 <sup>-2</sup> (-	2.06x10 <sup>-2</sup> (-	4.36x10 <sup>-4</sup>
				1.48)	1.39)	(-2.12)
TJP1 <sup>3</sup>	695	711	734	1.32x10 <sup>-1</sup> (-	2.28x10 <sup>-1</sup> (-	3.52x10 <sup>-3</sup>
				1.16)	0.99)	(-1.76)
HIF1AN	171	151	156	6.02x10 <sup>-1</sup> (-	9.23x10 <sup>-2</sup> (-	3.04x10 <sup>-3</sup>
				0.50)	1.44)	(-2.36)
YAP1	3331	3511	4020	1.36x10 <sup>-10</sup>	3.79x10 <sup>-3</sup> (-	<b>1.97x10</b> -6
				(-1.36)	0.90)	(-0.74)

487

By sponging miR-145-5p, circCDYL was suggested to have a regulatory effect on the expression of TJP1 in a circCDYL/miR-145-5p/TJP1 axis<sup>35</sup>. From the three treatments, only the samples exposed to a therapeutic dose of Epirubicin showed significant downregulation of TJP1 gene expression in both paralogs (Table 3, Figure 10A and 9B). The other two treatments, though significantly downregulated in the second paralog (ENSG00000277401), were not significant when taking the weighted average effect of both treatments (Table 3).







- 497 *Figure 10: Gene expression levels after exposure to anthracycline compounds. A&B*: *TJP1 gene expression for both paralogs*498 *after exposure to the therapeutic dose of Epirubicin. Taking the weighted average (Table 3, TJP1<sup>3</sup>) for both paralog genes (1:*
- 499 ENSG00000104067 and 2: ENSG00000277401), only the therapeutic dose of Epirubicin was significantly downregulated for
- 500 TJP1. C: HIF1AN gene expression after exposure to the therapeutic dose of Doxorubicin (p.adj. > 0.05). **D**: YAP1 gene
- 501 expression after exposure to the therapeutic dose of Doxorubicin (p.adj. < 0.001). Each dark grey bar represents a control
- 502 sample, while each pink bar represents a treated sample. The quantification shown was the result of the normalization performed
- 503 *by the DESeq2 pipeline.*

505 CircCDYL was suggested to interact with mRNAs encoding hypoxia-inducible factor asparagine 506 hydroxylase (HIF1AN) by acting as the sponge of miR-328-3p<sup>33</sup>. Even though HIF1AN gene expression 507 was not statistically significantly affected by the Doxorubicin treatments (Table 3), we observed a clear 508 down-regulation in all treatments (Figure 10C). The "mild" gene down-regulation, in addition to the up-509 regulation of miR-328-3p, and down-regulation of circCDYL might result in decreased translation of the 510 target.

511 CircCDYL was also suggested to be able to absorb miR-1180, thus alleviating the repression of miR-512 1180 on YAP<sup>36</sup>. The gene expression of YAP was significantly downregulated for all 3 treatments (Table 513 3, Figure 10D). Thus, the accumulated effects of the downregulation of circCDYL (Figure 1), the 514 upregulation of miR-1180 for most treatments (Figure 9D), and the downregulation of the gene 515 expression (Figure 10D) should have had a significant effect on YAP protein expression.

#### 516 4.4.3. Proteomics analysis

517 Translational inhibition is one of the primary functions of miRNAs, thus we expected to observe the 518 potential post-transcriptional regulation effects at the protein level. Consequently, we analyzed the protein 519 products of the gene targets mentioned above.

TJP1 (Q07157) 100 80 60 40 20 0 DOX\_The\_002\_3 DOX\_The\_008\_3 DOX\_The\_024\_3 DOX\_The\_024\_2 DOX\_The\_072\_2 DOX\_The\_072\_3 DOX\_The\_168\_3 DOX\_The\_168\_2 DOX\_The\_240\_2 DOX\_The\_336\_3 DOX\_The\_336\_2 DOX\_The\_002\_2 DOX\_The\_008\_2 DOX\_The\_008\_1 DOX\_The\_072\_1 DOX\_The\_240\_3 DOX\_The\_002\_1 DOX\_The\_024\_1 DOX\_The\_168\_1 DOX\_The\_240\_1 DOX\_The\_336\_1 В TJP1 (Q07157) 120 100 80 60 40 20 0 EPI\_The\_240\_2 EPI\_The\_002\_1 EPI\_The\_002\_2 EPI\_The\_002\_3 EPI\_The\_008\_1 EPI\_The\_008\_2 EPI\_The\_008\_3 EPI\_The\_024\_1 EPI\_The\_024\_2 EPI\_The\_024\_3 EPI\_The\_072\_2 EPI\_The\_072\_3 EPI\_The\_168\_2 EPI\_The\_168\_3 EPI\_The\_240\_3 EPI\_The\_336\_2 EPI\_The\_336\_3 EPI\_The\_072\_1 EPI\_The\_168\_1 EPI\_The\_240\_1 EPI\_The\_336\_1

520

А



Figure 11: Proteomics levels of TJP1 (UniProt ID: Q07157). A: After exposure to the therapeutic dose of Doxorubicin. B: After
 exposure to the therapeutic dose of Epirubicin.

TJP1 was proposed to be targeted by miR-145-5p, the latter has also been hypothesized to be regulated by
circCDYL, leading to a circCDYLßmiR-145-5pàTJP1 regulation. So far, we found that at the gene level,
only the therapeutic dose of Epirubicin showed a significant downregulation. When analyzing the

proteomics expression of TJP1-coded proteins, a clear time-dependent decrease was detected in the therapeutic dose of Doxorubicin (Figure 11A), where no significant decrease was found on average for the target gene expression. In the other 2 anthracycline treatments (Doxorubicin toxic dose and Epirubicin therapeutic dose), the protein expression was seldom detected, and only a decrease between 8 and 168h timepoints was identified (Figure 11B).

532 Due to the stochastic properties and limited sensitivity of mass spectrometry, no other target-coded 533 protein was quantified. In proteomics, this is especially the case for less abundant proteins, which was the 534 expected case for the analyzed proteins if their translation was inhibited.

#### 535 4.5. CircGNAS



#### 536

Figure 12: Graphical summary of the expression changes in the post-transcriptional regulation network of circGNAS after
 Amiodarone treatment with a toxic dose. The red arrows pointing down represent a significant expression decrease of that
 molecule. The green arrows pointing up represent a significant expression increase of that molecule. The black line arrows

540 represent the binding and/or translational inhibition of miRNAs to their targets. A green line arrow inside an XY plot represents

541 *a time-dependent increase in the abundance of a molecule.* 

#### 542 4.5.1. MiRNAs associated with circGNAS

543 No research has been published regarding circGNAS so far. Thus, we used miRanda to find miRNAs that

544 could potentially be sponged by circGNAS. Out of the ~ 140 potential miRNAs to bind circGNAS, four

- 545 of them were the most expressed and presented the highest expression consistency across all samples.
- 546 Three of them (hsa-miR-218-5p, hsa-miR-125a-3p, and hsa-miR-155-5p) showed differential expression,

- 547 one increasing and two decreasing respectively (Table 4). The remaining miRNA (hsa-miR-34a-5p) did
- 548 not present differential expression, even though it did present a time-dependent decrease.
- 549 *Table 4: Overall statistics of expressed miRNAs that were potentially regulated by circGNAS*. *P. adjusted values in bold are* 550 *significant (p. adj. value < 0.05)*. *These statistical results were obtained from DESeq2*.

miRNA Name	Mean expression (AMI Tox)	Log2 FoldChange	P. adj. value (AMI Tox)
		(AMI_Tox)	
hsa-miR-218-5p	688	1.25	1.16x10 <sup>-17</sup>
hsa-miR-125a-	44	-0.87	3.86x10 <sup>-8</sup>
3р			
hsa-miR-155-5p	187	-1.26	8.13x10 <sup>-4</sup>
hsa-miR-34a-5p	132	0.19	1.42x10 <sup>-1</sup>

552 Hsa-miR-218-5p expression showed a significant up-regulation after exposure to the toxic dose of

553 Amiodarone (Figure 13). In addition, a clear time increase was perceived.







557 Figure 13: Normalized counts through quantification by miRge2 for the Control versus the toxic dose of Amiodarone

comparison. A: hsa-miR-218-5p (p.adj. < 0.001). B: hsa-miR-125a-3p (p.adj. < 0.001). C: hsa-miR-155-5p (p.adj. < 0.001).</li>
Each dark grey bar represents a control sample, each light grey bar represents a control sample differentiated in the same conditions as the treatment, and each pink bar represents a treated sample. The quantification shown was the result of the

**561** *normalization performed by the DESeq2 pipeline. The title of the plot shows the miRNA ID next to the treatment name.* 

Hsa-miR-125a-3p also presented a significant differential expression, but as a down-regulation (Figure
13B). For this miRNA, we also observed a time-dependency, where the expression gradually decreased
over time.

- 565 For hsa-miR-155-5p, even though it was significantly down-regulated according to the p. adjusted value,
- analyzing the values we observed that this significance was mainly due to three triplicate samples with
- 567 outlier expressions (Figure 13C). In addition, we also observed (as in hsa-miR-125a-3p) a negative
- 568 cumulative effect across time.
- 569 4.5.2. Target transcript and protein expression
- 570 To assess which targets might be specifically regulated by circGNAS, we assembled all the potential
- 571 targets of the aforementioned four miRNAs based on TargetScan (7.1) data. The target overlap across the
- 572 four miRNAs can be visualized in the Venn Diagram (Figure 14). Eight transcripts were targeted by all 4
- 573 of the miRNAs selected.



575 Figure 14: Venn Diagram of the four miRNA targets.

576 Out of those eight transcripts, only three of them were sufficiently expressed: SAMD12-202 577 (ENST00000409003), PURB-201 (ENST00000395699), and PDE3A-201 (ENST00000359062). Two out 578 of the three did significantly differ from the overall control expression: PDE3A-201 (p.adj. < 0.001) and 579 PURB-201 (p. adj. < 0.05). Even so, we observed that the exposure expression of both transcripts was not 580 significantly different from the control technical replicates (AMI\_TOX\_000\_1 and AMI\_TOX\_000\_2). 581 In addition, we were able to note a tendency to slightly increase over time (Figure 15) for the three of 582 them, although the timeframe where this increase occurred varied between them.





Figure 15: PDE3A-201 (ENST00000359062) expression in control and toxic dose of Amiodarone treatment samples (p. adj.
 value < 0.001)</li>

587 We also analyzed the proteomics values of all 8 proteins coded by those 4 transcripts. None of the 588 proteins were detected in enough samples to be statistically analyzed.

## 589 5. Discussion

We wanted to evaluate whether the expression of any circular RNAs was affected by different 590 cardiotoxicants in an in vitro iPSC-derived cardiomyocyte model. We were able to identify and validate 591 the existence of several circular RNAs that showed differential expression after different compound 592 593 treatments. For circCDYL, we were able to detect a significant decrease in expression in response to 594 several anthracycline treatments. In addition, most miRNAs hypothesized to be sponged by this circRNA 595 were also found to be over-expressed, thus stimulating post-transcriptional regulation of their targets. The 596 targets' mRNA expression was either not affected or significantly decreased, and one of the proteins 597 (TJP1) coded by such target was identified to be decreased across time where no significant decrease was 598 found for the gene expression; thus, potentially linking the decrease to the studied post-transcriptional regulation. Interestingly, the regulatory axis involving circCDYL/miR-145-5p/TJP1 that we discovered in 599 600 anthracycline treated cardiac cells, has previously been identified in Wilms' Tumor (cancer of the kidneys<sup>35</sup>). A decrease in the abundance of TJP1 may be inducive to an impairment of the atrium-601

ventricle electrical conduction, as the knockout of other tight junction proteins such as CAR
 (Coxsackievirus-adenovirus receptor) are hypothesized to present cross talk effects with gap junctions<sup>37</sup>.

604 For circGNAS, a significant increase in expression was found after exposure to a toxic dose of Amiodarone. Additionally, these results were successfully reproduced using a different iPSC line and 605 cardiomyocytes derived from that iPSC-line, showing the robustness of iPSC-derived cardiomyocytes for 606 607 studying circRNAs. Several potential microRNAs were investigated as potential sponged candidates, 608 either showing differential expression or not. To understand what potential synergetic effect could be 609 derived from such different behaviors, we analyzed the targets of the microRNAs overlapping across each 610 other. Although the targets that were expressed did not show any differential expression (as miRNAs do not generally degrade mRNAs directly), we found a mild increasing tendency across time for some of 611 612 them. Unfortunately, none of the abundances of the proteins coded by those targets was quantified by 613 proteomics due to its limited sensitivity.

614 To quantify circular RNAs, we used CIRI2. CIRI2, a classical circRNA quantification tool, searches for 615 at least two reads spanning the back-spliced junctions to quantify such circular RNA. We realized that the 616 resulting quantification was very inconsistent: most of the quantified circRNAs were only detected in a 617 single sample, while only a single circular RNA was quantified in all samples. The repercussions were also seen in the statistical results, as most differentially expressed circular RNAs presented missing values 618 for several of the samples, making them less reliable. To increase the sensitivity, we developed a novel 619 620 strategy. We extracted mature circular RNA sequences from circBase and included them in a classical 621 transcriptome comprising linear mRNAs and used this as a reference for Salmon. The quantification 622 results showed a substantial sensitivity increase (compared to CIRI2), where only a minority of circular 623 RNAs (8%) were quantified in a single sample, though the number of constitutively expressed circRNAs was minimally increased (1 versus 9 circRNAs respectively). The sensitivity increase was probably due to 624 625 the sacrifice of CIRI2's specificity, as Salmon (circBase) does not require the presence of reads that span 626 the back-spliced junctions.

627 Interestingly enough, substantial experimental evidence validated the existence of circular RNAs 628 quantified by the remapping strategy. This was surprising considering the stochastic quantification of 629 circRNAs from CIRI2, and the hypothesis by some researchers that circRNAs are mostly a product of erroneous splicing. Even more interestingly, most of the circRNAs experimentally validated were not 630 631 found in our CIRI2 results, certainly for having a lack of sample coverage preventing a consistent identification of back-spliced junctions. Even so, we recognized that using exhaustive mapping for 632 633 circRNA quantification had its limitations. For this strategy, quantifying a linear RNA and a circular one 634 with the same exon sequence was equivalent, since the circRNA sequences extracted from circBase were

in the traditional 5' to 3' FASTA format and did not span the back-spliced junction as a continuoussequence.

Hence, for a consistent quantification with tools such as CIRI2, a high sequencing depth is required to
increase the chance of sequencing reads that span the required back-spliced junction. Pooling together the
sequencing reads of all samples might help to increase the chance of identifying existing circular RNAs.
Even so, the quantification at the sample level will still be limited, as the *in silico* predictions will be
derived independently from each sample.

Further research should investigate the biological impact of the validated circular RNAs. This could be accomplished either by overexpressing such circRNA (as is the case for circGNAS) or selectively repressing their expression (as we see in circCDYL). One of the main aims would be to analyze the regulatory potency of those circRNAs in relation to the translation of the targets involved. In addition, one could then further elucidate whether such an effect might be one of the main drivers of the cascade leading to the toxicity exerted by the compound.

In conclusion, we have identified, quantified, and validated several circular RNAs whose expression has significantly differed after being exposed to certain cardiotoxicants, in particular circCDYL and circGNAS. In addition, we analyzed the differential expression of potentially sponged miRNAs, their correspondent targets, and the proteins coded by those targets, finding for some of them a potential additive effect in different post-transcriptional regulatory events. These circRNAs might help understand the toxic effects/consequences those compounds induce at the molecular level.

654

655

## 656 6. Supplementary data

### 657 6.1. QPCR Primer Sequences

circBase ID	Gene of origin	5'-3'	3'-5'
		ACAGGCTTAGCTGTTAAC	GTCATAGCCTTTCCACCG
hsa_circ_0008285	circCDYL	GGGA	AACC
hsa_circ_0010791	CDC42	CCCACCTTCCCAAACCTA	CCCAACAAGCAAGAAAG

		AT	GAG
		CAGACCCAAGCTGTCCAT	CTCCAGCTTGGACTTCTT
hsa_circ_0026129	TUBA1A	TT	GC
		GCCCTGGATTTTGAGAAT	ATGGACAGGGTCAGTTG
hsa_circ_0034356	ACTC1	GA	GAG
		CGAGTAAGGCACACCCA	GACTCCTGGCTTTTCAGC
hsa_circ_0055922	FHL2	AAT	AAC
		AAAACAGCAGCAGCAAA	CAGCCATCTGTTGTTCCA
hsa_circ_0060999	circGNAS	CAA	GA
		CAGGGACCACACCCTGT	GTTCTGACATGGCGCTTA
hsa_circ_0076194	CDKN1A	ACT	CA
		CCCCTCACACACTTGGTT	TACCTGCCCCAGACTGAC
hsa_circ_0078905	None	TT	TT
		TGTTCCCACTCCCATCTT	GCTATCAGCCACAGCAAC
hsa_circ_0090448	TIMP1	TC	АА
		TTTGGAGGGGTTTATGCT	TGGATCATGTCATTGGCA
hsa_circ_0090904	MSN	CA	GT
		CCAAACAGAGCAGGAAA	GGTGAGGGGGAGCATTAC
hsa_circ_0102325	HIF1A	AGG	ATC
	circSMARCA	CCAAGATGGGCGAAAGT	AGATTCTGATCCACAAGC
hsa_circ_0001445	5	TCAC	CTCC
			ATGAGGTAGTCTGTCAGG
None	B-Actin	ATGGATGACGATATCGCT	Т

# 659 6.2. Sanger Primer Sequences

circBase ID	Gene of	Sequence	Directi
	origin		on

hsa_circ_0008		CAGGAAACAGCTATGACCGTCATAGCCTTTCCAC	
285	circCDYL	CGAACC	3'-5'
hsa_circ_0060		TGTAAAACGACGGCCAGTAAAACAGCAGCAGCA	
999	circGNAS	AACAA	5'-3'
hsa_circ_0001	circSMARC	TGTAAAACGACGGCCAGTCCAAACAGAGCAGGA	
445	A5	AAAGG	5'-3'

# 662 7. Bibliography

Panda, A. C., Grammatikakis, I., Munk, R., Gorospe, M. & Abdelmohsen, K. in *Wiley Interdisciplinary Reviews: RNA* Vol. 8 (Blackwell Publishing Ltd, 2017).

- Memczak, S. *et al.* Circular RNAs are a large class of animal RNAs with regulatory potency. *Nature* 495, 333-338, doi:10.1038/nature11928 (2013).
- Glažar, P., Papavasileiou, P. & Rajewsky, N. in *RNA* Vol. 20 1666 (Cold Spring Harbor
  Laboratory Press, 2014).
- Maass, P. G. *et al.* in *Journal of Molecular Medicine* Vol. 95 1179-1189 (Springer Verlag, 2017).
- 671 5 Gupta, S. K. *et al.* in *Circulation Research* Vol. 122 246-254 (Lippincott Williams and Wilkins,
  672 2018).
- 673 6 Wang, X. et al. in Life Sciences Vol. 265 (Elsevier Inc., 2021).
- 674 7 Yang, M. H. *et al.* in *Aging* Vol. 12 2530-2544 (Impact Journals LLC, 2020).
- 675 8 Deng, Y., Wang, J., Xie, G., Zeng, X. & Li, H. in International Journal of Biological Sciences
- 676 Vol. 15 2484-2496 (Ivyspring International Publisher, 2019).
- 677 9 Zhao, Y. et al. in Journal of Applied Toxicology (John Wiley and Sons Ltd, 2021).
- Gao, Y. *et al.* Comprehensive identification of internal structure and alternative splicing events in
  circular RNAs. *Nat Commun* 7, 12060, doi:10.1038/ncomms12060 (2016).
- Metge, F., Czaja-Hasse, L. F., Reinhardt, R. & Dieterich, C. FUCHS-towards full circular RNA
  characterization using RNAseq. *PeerJ* 5, e2934, doi:10.7717/peerj.2934 (2017).
- K. O. *et al.* Diverse alternative back-splicing and alternative splicing landscape of circular
  RNAs. *Genome Res* 26, 1277-1287, doi:10.1101/gr.202895.115 (2016).
- Kuepfer, L. *et al.* A model-based assay design to reproduce in vivo patterns of acute drug-induced
  toxicity. *Arch Toxicol* 92, 553-555, doi:10.1007/s00204-017-2041-7 (2018).

686	14	Gao, Y., Zhang, J. & Zhao, F. Circular RNA identification based on multiple seed matching.
687		Brief Bioinform 19, 803-810, doi:10.1093/bib/bbx014 (2018).
688	15	Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for
689		RNA-seq data with DESeq2. Genome Biol 15, 550, doi:10.1186/s13059-014-0550-8 (2014).
690	16	Cunningham, F. et al. Ensembl 2019. Nucleic Acids Res 47, D745-D751,
691		doi:10.1093/nar/gky1113 (2019).
692	17	Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and bias-
693		aware quantification of transcript expression. Nat Methods 14, 417-419, doi:10.1038/nmeth.4197
694		(2017).
695	18	Verheijen, M. C. et al. R-ODAF: Omics data analysis framework for regulatory application.
696		Regul Toxicol Pharmacol 131, 105143, doi:10.1016/j.yrtph.2022.105143 (2022).
697	19	Siede, D. et al. Identification of circular RNAs with host gene-independent expression in human
698		model systems for cardiac differentiation and disease. J Mol Cell Cardiol 109, 48-56,
699		doi:10.1016/j.yjmcc.2017.06.015 (2017).
700	20	
701	21	Agarwal, V., Bell, G. W., Nam, J. W. & Bartel, D. P. Predicting effective microRNA target sites
702		in mammalian mRNAs. Elife 4, doi:10.7554/eLife.05005 (2015).
703	22	John, B. et al. Human MicroRNA targets. PLoS Biol 2, e363, doi:10.1371/journal.pbio.0020363
704		(2004).
705	23	Griffiths-Jones, S., Grocock, R. J., van Dongen, S., Bateman, A. & Enright, A. J. miRBase:
706		microRNA sequences, targets and gene nomenclature. Nucleic Acids Res 34, D140-144,
707		doi:10.1093/nar/gkj112 (2006).
708	24	Lian, X. et al. Directed cardiomyocyte differentiation from human pluripotent stem cells by
709		modulating Wnt/beta-catenin signaling under fully defined conditions. Nat Protoc 8, 162-175,
710		doi:10.1038/nprot.2012.150 (2013).
711	25	Altschul, S. F. et al. Gapped BLAST and PSI-BLAST: a new generation of protein database
712		search programs. Nucleic Acids Res 25, 3389-3402, doi:10.1093/nar/25.17.3389 (1997).
713	26	Selevsek, N. et al. Network integration and modelling of dynamic drug responses at multi-omics
714		levels. Commun Biol 3, 573, doi:10.1038/s42003-020-01302-8 (2020).
715	27	
716	28	Wink, S., Hiemstra, S. W., Huppelschoten, S., Klip, J. E. & van de Water, B. Dynamic imaging
717		of adaptive stress response pathway activation for prediction of drug induced liver injury. Arch
718		Toxicol 92, 1797-1814, doi:10.1007/s00204-018-2178-z (2018).

719	29	Wang, S. et al. circCDYL Acts as a Tumor Suppressor in Triple Negative Breast Cancer by
720		Sponging miR-190a-3p and Upregulating TP53INP1. Clin Breast Cancer 20, 422-430,
721		doi:10.1016/j.clbc.2020.04.006 (2020).
722	30	Bian, W. X., Xue, F., Wang, L. Y. & Xing, X. F. Circular RNA CircCDYL Regulates
723		Proliferation and Apoptosis in Non-Small Cell Lung Cancer Cells by Sponging miR-185-5p and
724		Upregulating TNRC6A. Cancer Manag Res 13, 633-642, doi:10.2147/CMAR.S280315 (2021).
725	31	Zhang, M. et al. Circular RNA (circRNA) CDYL Induces Myocardial Regeneration by ceRNA
726		After Myocardial Infarction. Med Sci Monit 26, e923188, doi:10.12659/MSM.923188 (2020).
727	32	Cui, W., Dai, J., Ma, J. & Gu, H. circCDYL/microRNA-105-5p participates in modulating
728		growth and migration of colon cancer cells. Gen Physiol Biophys 38, 485-495,
729		doi:10.4149/gpb2019037 (2019).
730	33	Wei, Y. et al. A Noncoding Regulatory RNAs Network Driven by Circ-CDYL Acts Specifically
731		in the Early Stages Hepatocellular Carcinoma. <i>Hepatology</i> 71, 130-147, doi:10.1002/hep.30795
732		(2020).
733	34	Liang, G. et al. MiR-92b-3p Inhibits Proliferation of HER2-Positive Breast Cancer Cell by
734		Targeting circCDYL. Front Cell Dev Biol 9, 707049, doi:10.3389/fcell.2021.707049 (2021).
735	35	Zhou, R. et al. CircCDYL Acts as a Tumor Suppressor in Wilms' Tumor by Targeting miR-145-
736		5p. Front Cell Dev Biol 9, 668947, doi:10.3389/fcell.2021.668947 (2021).
737	36	Chen, F. et al. Circular RNA circ-CDYL sponges miR-1180 to elevate yes-associated protein in
738		multiple myeloma. Exp Biol Med (Maywood) 245, 925-932, doi:10.1177/1535370220918191
739		(2020).
740	37	Lisewski, U. et al. The tight junction protein CAR regulates cardiac conduction and cell-cell
741		communication. J Exp Med 205, 2369-2379, doi:10.1084/jem.20080897 (2008).
740		
/4Z		

# <sup>1</sup> Supplementary Results (Chapter 2)

# 2 CircCDYL



6:4891713|4892379:+ (Dox\_The)



6:4891713|4892379:+ (Epi\_The)



5

6 CDYL

ENSG00000153046 (Dox\_The)





ENSG00000153046 (Epi\_The)



















Epi\_The\_336\_ Epi\_The\_336\_ 072 168 168 The The he Ъе The ЩЩ 88 8 8 ž

18



hsa-miR-328-3p (Dox\_Tox)



hsa-miR-328-3p (Dox\_The)





- 23 hsa-miR-1180-3p
- 24



# hsa-miR-1180-3p (Dox\_The)



hsa-miR-1180-3p (Epi\_The)



## 28 TJP1 (ENSG00000104067)



29

ENSG00000104067 (Dox\_Tox)



# 32 TJP1 (ENSG00000277401)








ENSG00000166135 (Dox\_The)

38

ENSG00000166135 (Dox\_Tox)















#### 46 TJP1 (Q07157)

#### 47 Normalized Log2 Values



Q07157 (ZO1\_HUMAN)



48

Q07157 (ZO1\_HUMAN)



- 50
- 51 Normalized Values
- 52



Q07157 (ZO1\_HUMAN)



Q07157 (ZO1\_HUMAN)

54



Q07157 (ZO1\_HUMAN)



hsa\_circ\_0060999 (AMI\_Tox)

hsa-miR-218-5p 58

57









62 hsa-miR-155-5p









#### 68 PURB-201 (ENST00000395699)







# Chapter 3:

# Quantifying the number of translatable transcripts through the use of OMICs involved in posttranscriptional regulation

- 1 Quantifying the number of translatable transcripts
- <sup>2</sup> through the use of OMICs involved in post-

# <sup>3</sup> transcriptional regulation

4 Juan Ochoteco Asensio<sup>1</sup>, Jos Kleinjans<sup>1</sup>, Florian Caiment<sup>1</sup>

5 <sup>1</sup>Department of Toxicogenomics, School of Oncology and Developmental Biology (GROW),

6 Maastricht University, Maastricht, The Netherlands

#### 7 (Manuscript preprint available)

#### 8 Abstract

9 Transcriptomics is nowadays frequently used as an analytical tool to study the extent of cell 10 expression changes between two phenotypes or between different conditions. However, an important 11 portion of the significant changes observed in transcriptomics at the gene level is usually not 12 consistently detected at the protein level by proteomics. This poor correlation between the measured 13 transcriptome and proteome is probably mainly due to post-transcriptional regulation, among which 14 miRNA and circRNA have been proposed to play an important role. Therefore, since both miRNA 15 and circRNA are also quantified by transcriptomics, we proposed to build a model taking those 16 factors into account to estimate, for each transcript, the fraction of transcripts that would be available 17 for translation. Using a dataset of cells exposed to diverse compounds, we evaluated how our model 18 was able to improve the correlation between the assessed transcriptome and proteome expression 19 levels. The results show that the model improved the correlation for a subset of genes, probably due 20 to the regulation of different miRNAs across the genome.

## 21 Introduction

The central dogma of molecular biology states a straightforward flow of information: for a gene to transmit the information it contains, its DNA needs to be transcribed to RNA to produce the desired protein. Each of these biological pools of molecules has its field of study for both their characterization and quantification, which are commonly known as omics. Transcriptomics, for example, refers to the large-scale study of transcripts, and the same applies to proteomics (proteins) or metabolomics (metabolites).

29 Proteins, being the functional molecules of the cell, are of high interest to be analyzed, as they 30 accurately represent the phenotypic changes of the studied cell or tissue. Unfortunately, the 31 technology' sensitivity and reproducibility behind proteomics, namely mass spectrometry, is currently 32 still limited, especially when the aim is to get an exhaustive protein expression analysis. This is 33 mainly due to the mass spectrometer only being able to measure a fraction of the eluted peptides, 34 hence giving as an output a slight portion of the entire population of proteins<sup>1-3</sup>. To circumvent this 35 issue, an alternative strategy often used is the expression analysis of their precursors: the transcripts. 36 Transcriptomics is indeed usually performed as the surrogate to analyze different disease states or cell 37 conditions. Even so, the fact that transcripts are indeed the cause and origin of proteins is not reflected by a high correlation between both  $omics^{4-8}$ . This is rather unsurprising, due to the multiplicity and 38 39 complexity of the factors playing a role in post-transcriptional regulation, in addition to intrinsic 40 variable characteristics of both molecules (such as half-life<sup>9</sup>), and some of those regulatory factors are 41 even transcripts themselves.

42

43 One such factor is microRNAs (miRNAs), whose effect on the expression levels of proteins, although mild, has been characterized for more than a decade<sup>10,11</sup>. MiRNAs are short non-coding RNA 44 45 sequences playing a role in translation regulation<sup>12</sup>. They contain the so-called "seed" region<sup>13</sup>, a short 46 sequence with perfect Watson-Crick complementarity with their target (primarily to their 3' UTR region). The binding between the two molecules impairs the translation of the target<sup>14</sup>. Due to the 47 48 short length of the seed region, a single miRNA often targets several transcripts, and numerous 49 miRNAs can target a single transcript. This inhibitory relationship may differ depending on the level 50 of expression of the miRNA and its target molecules, the number of regions in a single mRNA 51 complementary to the miRNA present, and the combinatorial effect of several miRNAs targeting the 52 same transcript. Taking into account these parameters is necessary to have a better comprehension of 53 the mechanisms behind post-transcriptional regulation.

54

55 Another element to be considered in post-transcription regulation is circular RNAs (circRNAs). 56 CircRNAs are transcripts that possess a circular structure due to the covalent binding of their 5' and 3' ends<sup>15,16</sup>, hence obtaining a circular structure that bestows them resistance to exonuclease activity<sup>17</sup>. 57 One of their recently discovered functions is described as 'miRNA sponges'<sup>18</sup>. This refers to their 58 59 ability to contain several copies of short sequences that are complementary to the miRNAs' seed 60 region. Consequently, they also play a role in the regulation of translation machinery by competing against other transcripts for miRNA binding<sup>19</sup>. The number of miRNAs captured by circRNAs 61 62 influence the number of transcripts available for translation. 63

64 CircRNAs are not the only targets competing for miRNA, as other (long) non-coding RNAs 65 (ncRNAs) can also present a seed target for some particular miRNAs<sup>20</sup>. This collection of coding and

- 66 non-coding transcripts can be conceptualized as a single group of targets for a shared miRNA, which
- 67 led them to be generally called competing endogenous RNAs  $(ceRNAs)^{21}$ .
- 68

69 Thus, in the current age of sequencing, accurately associating transcriptomics to the actual phenotype 70 of the cell is a major challenge. Currently, many publications base their biological interpretation on 71 the differentially expressed genes in two groups of samples as transcriptomics is more sensitive and 72 allows the assessment of almost all possible RNA molecules in a single experiment. In this context, 73 we wondered if it would be possible to integrate all available transcriptomics information into a model 74 able to estimate the level of translation of any expressed coding transcript, namely, the fraction of 75 translatable transcripts (TrT). This value, computed for each protein-coding transcript, is an 76 estimation of the number of transcripts that are free to be translated after taking into account the 77 aforementioned post-transcriptional regulation.

78

To design and assess the model, we used an *in vitro* dataset obtained from human cardiac microtissues exposed to a range of compounds at different doses, which were analyzed with proteomics, RNA-Seq (ribo-depleted libraries), and miRNA-Seq methods. From the RNA sequencing data, we first identified and quantified different transcript biotypes: protein-coding, non-coding, miRNA, and circRNAs. The proposed model was then applied to generate the TrT score from all possible interactions between those molecules. In this manuscript, we analyze the possible benefits of such an approach compared to the state-of-the-art methods for gene expression analysis.

86

#### 87 Methods

#### 88 Samples

89 The analyzed data consists of 3D microtissues containing stem-cell-derived cardiomyocytes and 90 fibroblasts in a 4:1 ratio obtained from InSphero. These microtissues were exposed to 8 compounds: 91 Fluorouracil (5FU), Amiodarone (AMI), Celecoxib (CEL), Docetaxel (DOC), Doxorubicin (DOX), 92 Epirubicin (EPI), Mitoxantrone (MXT), and Paclitaxel (PTX), in addition to a control group 93 (Untreated/UNTR). The dosing profile was established via the use of the physiologically based 94 pharmacokinetic (PBPK) modeling software PK-Sim, to simulate exposure levels under physiological 95 conditions<sup>22</sup>. For each compound, 2 doses were applied: Therapeutic and Toxic. The exposures were 96 done in triplicates, and the data extraction was performed at 7 time-points: 2, 8, 24, 72, 168, 240, and 97 336 hours; resulting in 21 data-points per dose. Four samples were excluded due to low sequencing 98 depth: three samples for having too low miRNA sequencing depth 99 (CEL\_Tox\_002\_3/MXT\_Tox\_002\_3/PTX\_Tox\_002\_3; toxic dose, time-point 2h, 3<sup>rd</sup> triplicate) and
100 one for its RNA-Seq low read count (UNTR 002 3).

#### 101 RNA Sequencing (RNA-Seq)

Total RNA from the exposed microtissues was isolated using the Qiagen AllPrep Universal Kit (Cat #80224). Ribo-depletion was achieved by using the Illumina RiboZero Gold kit (Cat #MRZG12324),
and the libraries were prepared using the Lexogen SENSE total RNA kit (Cat #009.96). All libraries were then sequenced on an Illumina HiSeq 2000 at 100 bp paired-end at an average coverage range between 20 and 30 million reads. The adaptors were removed through Trimmomatic version 0.33<sup>23</sup>.
We used the following parameters: paired-end, ILLUMINACLIP: TruSeq3-PE.fa:2:30:10, LEADING:3, TRAILING:3, SLIDINGWINDOW:4:15, MINLEN:36, HEADCROP:12.

#### 109 Proteomics

110 Proteins were isolated and diluted to a concentration below 0.2M. The peptides were digested by 111 trypsin, and peptides were cleaned-up using Sep-Pak tC18 cartridges (Waters) according to the 112 manufacturer's instructions. A vacuum centrifuge was used to dry the peptides, before measuring 113 them on an Orbitrap Fusion mass spectrometer (Thermo Fisher Scientific), which was coupled to a 114 NanoLC-2D HPLC system (Eksigent). To process the raw MS data, Genedata Expressionist software 115 (v.11.0) was used. The noise in the LC-MS peaks was reduced and normalized, and afterward, their 116 properties were obtained (m/z and RT boundaries, m/z and RT center values, intensity). To annotate 117 the individual MS/MS spectra, Mascot 2.6 was used. To group peak clusters, protein interference was 118 used (based on peptide and protein annotations), and, using the Hi3 method, protein intensities were 119 computed.

120

#### 121 miRNA analysis

Starting from the same total RNA isolated for the ribo-depleted libraries, an aliquot was size selected and ligated using the TruSeq Small RNA Library Prep Kit (Illumina®). After sequencing on the HiSeq 2000 at 3.6 million reads per sample (after quantification), we quantified the resulting data using miRge2 (last change: 05/06/2018)<sup>24</sup>. miRge2 used the MirBase database as the reference library (miRBase v22), bowtie-1.1.1<sup>25</sup> as the mapper, miRge2/sp as the miRge library, human as the species selected, and illumina for the adapter to be removed. The output results were in gff format. Isomirs were not considered for this analysis.

#### 130 (non-)coding RNA analysis

131 The genome version used for all the transcriptomics analyses was the Genome Reference Consortium 132 Human Build 38 (GRCh38.p12). For the identification of circRNAs, we performed a de novo 133 prediction from our dataset. For this, we concatenated all RNA-Seq forward data (R1) across all 240 134 Cardiac samples into a single FASTQ file. Afterward, this file was inputted to two recent circRNA prediction software: circExplorer2<sup>26</sup> and CIRI2<sup>27</sup>. For circExplorer2, we first used BWA version 135 0.7.17<sup>28</sup> to index the genome and align (minimum score to output: 19) the reads to the indexed 136 137 genome. Afterward, we used circExplorer2 to parse the aligned reads and annotate the circRNAs. For 138 CIRI2, we also required the aligning step via BWA, and then the annotation step through its internal algorithm (-S (single-end), -U 3 -B 13, to set mapping quality thresholds of a junction read and help 139 140 control False Discovery Rate (FDR)). We then extracted the overlap of identified molecules between 141 both outputs for decreasing the amount of false-positive predictions. The circRNA IDs were 142 comprised of the chromosomal and strand locus of the predicted molecules. Based on this information, we extracted their genetic sequence using BEDTools<sup>29</sup>, obtaining a final circRNA 143 144 transcriptome library.

145

146 We downloaded the transcriptomes for both coding (all cDNA) and non-coding (all ncRNA) 147 transcripts from Ensembl (release 96)<sup>30</sup>. We combined these 2 libraries with our predicted circRNAs 148 library into a single library, which we set as the global transcriptome reference for Salmon<sup>31</sup>, with 149 which we quantified the RNA-Seq data. Salmon output contained both the number of reads and TPM 150 (Transcript per Million) values for each transcript, the latter of which was used for the analysis. The 151 benefit of using TPM values, instead of raw counts, relies on the inherent normalization performed for 152 both sequencing depth and transcript length. If raw reads would have been used instead, the variability 153 across samples and compounds would have made it impossible to optimize the model consistently. 154 Moreover, since our model computes interactions between molecules of different sizes, using read 155 counts would have greatly biased the output.

#### 156 CircRNA predictors selection

157 We selected CIRI2 and circExplorer2 out of 4 possible predictors, which also included find circ18 and circRNA finder<sup>32</sup>. The parameters for find circ were set as default, and the samtools version was 1.3. 158 For circRNA finder, it involved 2 steps: running STAR (-c --runThreadN 4 --chimSegmentMin 20 --159 160 chimScoreMin 1 --alignIntronMax 500000 --outFilterMismatchNmax 4 --161 BAM alignTranscriptsPerReadNmax 100000 --twopassMode Basic --outSAMtype 162 SortedByCoordinate --chimOutType SeparateSAMold --outFilterMultimapNmax 2) and the post-163 processing script (--minLen 100). We formatted all the output equally so that we were able to 164 compare them across tools.

#### 165 Interaction tables

To evaluate all possible bindings between the molecules, a table containing all possible interactions 166 167 between miRNAs and any putative competing ceRNAs was generated, each row representing a unique interaction. For miRNA, both identifiers (IDs) and sequences were sourced from miRBase<sup>33</sup>. 168 169 For all transcripts (both coding and non-coding), we extracted the Ensembl transcript IDs through biomaRt<sup>34</sup>. We obtained the sequences corresponding to these IDs through SAMtools<sup>35</sup>, which were 170 searched in the genome FASTA file by using the chromosomal coordinates of the transcripts. For 171 172 mRNAs, we used only the 3' UTR region. To establish the microRNA target interactions (or MTIs), we used miRanda (version 3.3a)<sup>36</sup> with the -strict parameter to force strict 5' seed pairing. We set the 173 174 score threshold to 140.0 (corresponding to a full 7 base pair seed with no mismatch). We also used –

175 noenergy, disabling the thermodynamics performance.

#### 176 Translatable Transcripts (TrT) formulation

We designed a formula (Formula 1) that aimed to estimate the fraction of translatable transcripts based on a basic principle: the number of translatable transcripts is equal to the difference between the total expression of a given protein-coding transcript and its inhibited molecules. The number of inhibited transcripts depends on the number of miRNAs and the probability that these will bind to the target transcript. The probability of a target being targeted by miRNA depends on how abundant the target is in proportion to the number of all the possible targets. For all data analyses, we used R<sup>37,38</sup> as the main programming language.

184

$$TrT_{transcript} = TPM_{transcript} - miRNA_{factor} * \sum_{i=miRNA}^{N} RPM_i * \frac{TPM_{transcript} * Seeds_{i,transcript}}{\sum_{j=target}^{M} TPM_j * Seeds_{i,j}}$$
185

**186**Formula 1: TrT formula

187

188

189 TrT = Translatable Transcripts, TPM = Transcripts Per Million, RPM = Reads Per Million, Seeds = 190 the number of seeds present in a specific transcript targeted by a specific miRNA, N = number of 191 miRNAs targeting the transcript, M = number of targets for a specific miRNA<sub>i</sub>, miRNA<sub>factor</sub> is a 192 scaling factor, set to 0.1 (TrT model design).

193 TrT example

For example, consider a protein-coding transcript expressed at 100 TPM targeted by two miRNAs,
expressed at 132.5 and 227.5 Reads Per Million (RPM) respectively. Each of them has a single seed
or sequence region presented by the target that they can bind to. The first miRNA can also interact

- 197 with a single circRNA, expressed at 2 TPM on 3 perfect seed regions, while the second miRNA can
- bind to an ncRNA with a TPM value of 30 (Equation b).
- 199

201 Those miRNAs can bind to any of the aforementioned ceRNAs (other mRNA targets, circRNAs, or 202 ncRNAs). To know how many miRNAs will bind to each of them, we need to know the probability of 203 each interaction. In our study, we hypothesized that the probability that a miRNA will bind to a 204 molecule is directly proportional to how abundant that molecule is in relation to all possible targets 205 (0.94 and 0.77 in this example, Equation c). Consequently, we estimate that 300 of these miRNA 206 molecules will bind to the target. This number is multiplied by the miRNA factor (Equation d), and 207 we subtract it from the original TPM value of our target, which will be the expected TrT value 208 (Equation e).

#### 209 Correlation between Transcriptomics and Proteomics

210 To perform correlation analysis between our transcriptomics and proteomics dataset, we used the 211 untreated samples of the first 4 time-points (2, 8, 24, and 72 hours) for transcriptomics, and all the

- normalized (Proteomics) time-points for proteomics (2, 8, 24, 72, 168, 240, and 336 hours).
- 213

To be able to correlate protein and transcript expressions, we merged protein-coding transcripts by their protein product, and we grouped all values per time-point while maintaining the same order to perform a Pearson correlation at the time-point level. For 168h, we omitted the first triplicate due to a

217 significant batch effect, correlating only the other 2 triplicates with the other 2 triplicates of all other

- 218 time-points. We did the same with the RNA-Seq sample of UNTR\_002\_3, because of the low
- 219 sequencing depth, as mentioned before. The graphical displays were generated thanks to the 'corrplot'  $\frac{39}{1000}$
- 220 package<sup>39</sup>.
- 221 Differential Expression Analysis
- 222 Proteomics

223 To obtain differentially expressed proteins (DEPs), we followed the steps performed by Selevsek et al 224 (2020). Briefly, we first log-2 transformed all the proteomic expression values. We then calculated the 225 median for each control sample. Afterward, we calculated the median of all the medians (MoM) and 226 shifted the control samples so that they all shared this median. For every treatment/dose combination, 227 we determined the set of proteins common in both control and treatment samples. This set of proteins 228 was used to determine the MoM on the normalized control samples. We then shifted the distribution 229 of the data by matching their median values to this MoM. Finally, we performed paired t-tests to 230 evaluate the significance between doses, considering as significant (or DEPs) the ones for which the 231 p-values were lower than 0.05.

#### 232 Transcriptomics

233 As we aim at evaluating the TrT fraction of the coding transcripts, Salmon mapping output was 234 filtered out for non-coding transcripts. All remaining coding transcript expression values were 235 summed based on their gene of origin. While usually, differential expression analyses are performed 236 via a dedicated sequencing analysis pipeline on raw read count using a pipeline specific for negative 237 binomial distribution (such as DESeq2 or EdgeR), this approach could not be applied here. Indeed, in our case, the values to be compared (TPM and TrT) did not allow such pipelines, either because the 238 239 input was expected to be un-normalized<sup>40</sup> and/or because other normalization steps were performed instead<sup>41</sup>. We then log-2 transformed both TPM and TrT values, which, like proteomics, also 240 241 presented originally a negative binomial distribution. We performed t-tests between therapeutic and 242 toxic doses. Due to the high percentage of significant observations and multiple testing, we also 243 applied an FDR/BH p-adjustment. We identified a differentially expressed gene (or DEG) as a gene 244 with a p-adjusted value lower than 0.05

#### 245 Biological Interpretation with GOrilla

Using as input the lists of genes ranked by increasing p-adjusted value (for each comparison and compound), we ran a gene ontology enrichment analysis with GOrilla<sup>42</sup> using the 'Single ranked list

- 248 of genes' mode. We studied the GO terms of interest and analyzed some of their genes' expression
- 249 (both in TPM and TrT). We also explored how the MTIs (using the miRTarBase<sup>43</sup> list) related to such

genes (with at least weak evidence of having a regulatory effect) were expressed in comparison towhat TPM and TrT were representing.

#### 252 Sensitivity, Specificity, and Accuracy

253 A prediction is normally evaluated by contrasting it to the reality it is trying to model. Such an 254 evaluation focuses on different aspects of the predictor: how many of the positives that have been 255 predicted are real positives (sensitivity), how many of the predicted negatives are indeed negative 256 (specificity); and how many cases are correct out of all cases, independently of whether they are 257 positive or negative (accuracy), making use of binary classification terms. In our case, proteomics was 258 considered the true condition, while transcriptomics was the predicted condition. True or false 259 referred then to the correct or incorrect representation of proteomics by transcriptomics, while 260 positive or negative referred to the presence or absence of differential expression, respectively. For 261 example, if both proteomics and transcriptomics showed differential expression in the same manner, 262 such gene was considered a true positive. If transcriptomics did not present a differential expression 263 while proteomics did, it was considered a false negative. Once these terms were defined, the 264 calculation of such evaluators was performed according to the standard statistical measures of a 265 binary classification test performance.

#### 266 Results

#### 267 Test dataset for TrT model assessment

268 To develop a TrT model, we used a dataset as input derived from a 3D cardiomyocyte culture, from 269 which both transcriptomics (ribo-depleted and small RNA libraries) and proteomics (LC-MS) data 270 were generated. The cell cultures were composed of 8 individual compound treatments, in addition to 271 an untreated control. For every compound, 3 replicates were measured in 7 time-points (2, 8, 24, 72, 272 168, 240, and 336 hours) for every dose (therapeutic and toxic), resulting in 42 samples per 273 compound. This dataset, of a total of 240 samples, presented the added value to have been generated 274 by the same technician, and all proteomics and transcriptomics samples were generated from the same 275 run, thus reducing an important source of bias. The three different technologies (RNA-Seq, miRNA-276 seq, and proteomics) were quantified based on state-of-the-art procedures (Methods). While our 277 RNA-Seq data analysis quantified the protein-coding and non-coding transcripts available in the 278 Ensembl database, circular RNA identification required additional steps such as de novo prediction.

#### 279 CircRNA predictors selection

CircRNA prediction programs are known for generating a high level of false positives<sup>44</sup>. To identify 280 the circular RNAs in this transcriptomics dataset, we decided to use an overlapping approach between 281 several tools. For this, we considered the following four tools: CIRI2<sup>27</sup>, find circ<sup>18</sup>, circRNA finder<sup>32</sup>, 282 and CIRCexplorer2<sup>26</sup>. When determining the optimal overlap, we interpreted the amount of uniquely 283 284 predicted identities by a tool as an indicator of the false-positive ratio. We could observe that most of 285 the circRNAs were only predicted by circRNA finder, which suggested this algorithm has the highest 286 rate of false positives. The other 3 tools showed similar false-positive ratios, so to compare them we 287 focused on the overlap that maximized the number of identified molecules: circExplorer2 and CIRI2. 288 The list of predicted circRNA molecules was then added to the reference transcriptome used to 289 quantify all expressed RNA molecules from the ribo-depleted libraries.



290 291

Figure 1: Venn Diagram between 4 circRNA prediction tools output

Such transcriptome was structured in the following manner: 126831 of the identifiers were noncoding RNA (57.37%), 85108 were protein-coding RNAs (38.50%), and 9138 were circRNAs
(4.13%) (Figure 2A).

296

In control conditions, for RNA-Seq, we detected 92969 transcripts expressed in at least 1 sample > 0 reads, with 4815 (5%) of them present in all 240 samples. In microRNA sequencing (miRNA-seq), out of the 1348 quantified molecules in any sample, only 55 miRNAs (4%) were detected constitutively. In proteomics, there were 1392 proteins quantified, of which 362 (26%) were detected in all samples. As expected, proteomics, having one order of magnitude less quantified molecules (both in Total and Constitutively) than RNA-Seq (Figure 2B), was the limiting factor when comparing both. Also, when comparing proteomics and miRNA-seq, although they have similar total

**<sup>292</sup>** General Expression Analysis

quantification numbers, the difference in their constitutively quantified shows that the former has a
 higher proportion of highly expressed molecules, while miRNA-seq displays a broader representation
 of all molecules (in addition to its more volatile/temporary function).

307

308



#### 309



<sup>313</sup> 

been quantified in all samples by such technology.

#### 314 TrT model design

315 With all individual RNA molecules from the transcriptome characterized and quantified, we 316 investigated the possibility to develop a model to predict the expected translatable fraction (TrT) of 317 any given coding transcript. For this, we hypothesized that taking into account as many post318 transcriptional regulation factors as possible would be fundamental. Therefore, we conceptualized a 319 model based on several key features of the various biological factors available. First, miRNAs 320 decreased a transcript's chance of being translated. An increase in miRNA expression (controlling for 321 all other conditions) would then decrease the number of mRNAs available for translation<sup>45</sup>. Second, 322 the probability that a miRNA transcript interaction (MTI) occurred is directly proportional to the 323 expression level of a target in proportion to the expression level of all possible targets. Lastly, both 324 coding and non-coding RNAs (including circRNAs) that were able to interact with a miRNA 325 inhibiting the coding transcript of interest could work as ceRNA: the higher the expression of a single target, the lower the inhibition for the rest of the targets<sup>46-48</sup>. In other words, we assumed that a given 326 327 miRNA would be equally distributed among all its possible targets (taking into account the number of individual seeds), and only the fraction available to interact with our gene of interest should be 328 329 considered to have an inhibiting effect.

330

331 The central concept behind this model was that the number of transcripts to be translated was the 332 difference between the total amount of transcripts (in TPM) and the ones that would not be translated 333 due to miRNA inhibition. The number of transcripts that would be inhibited was then based on two 334 main factors: the level of expression of each miRNA that had the mRNA as a target (RPMi), and the 335 probability that each miRNA would interact with the target. We hypothesized that the probability that 336 a miRNA would bind to a target was proportional to how prevalent that target was compared to all 337 possible targets. A target did not always equal a single molecule (a circRNA could present several 338 seed regions/targets). For this reason, we represented a target by its total amount of seed regions for a 339 specific MTI (transcript expression level times its number of seed regions).

340

341 Because miRNA libraries present a smaller density than ribo-depleted RNA libraries, the raw number 342 of reads generated for the miRNAs was in general higher than the quantification levels of other 343 transcripts. This difference did not allow us to subtract one from another directly, and miRNA 344 expression should be scaled to a comparable level of the coding RNA transcripts. To know which 345 value should be given for this scaling factor, we first analyzed which factors would be optimal for 346 TrT, that is, in which miRNA had enough power to show a difference between TrT and TPM, but not 347 so powerful that it reduced all expression values to zero. We searched such optimal value by 348 evaluating a range of values between 0 and 1 (0 leading to no miRNA effect at all, and 1 leading to no 349 scaling), and selected the values in which TrT maximized the correction benefit in comparison to 350 TPM (taking proteomics as a reference). The optimal range we observed was between 0.1 and 0.27. 351 Besides, initial investigation of a newly developed sequencing method named Combo-Seq, which 352 allows sequencing both transcripts and miRNAs altogether in a single library preparation, showed us 353 that the proportion observed of miRNA represented around 10% of all transcripts sequenced.

354 Consequently, we decided to set the miRNA factor to '0.1' (Translatable Transcripts (TrT) 355 formulation)

#### 356 Correlation between Transcriptomics (TPM/TrT) and Proteomics

Having our method described and finished, we applied our method to all the aforementioned data. Afterward, we aimed to compare it with the state-of-the-art TPM to investigate in which manner TrT could be an improvement. Initially, we wanted to confirm the initial issue (transcriptomics does not accurately represent the proteome), and compare such results with TrT. For that purpose, using the untreated samples, we performed a correlation analysis between proteomics and TPM, followed by its analogous analysis between proteomics and TrT.

363

364 The correlation analysis confirmed the low correlation between Transcriptomics (TPM) and 365 proteomics values (Figure 3A), the average correlation of which was  $0.39 (\pm 0.06 \text{ SD})$ . There was no 366 specific time interval between both omics that presented an improved correlation. For example, 367 independently of the time difference between both omics, 24h samples (be it proteomics or 368 transcriptomics) showed the highest correlations. Simultaneously, 2h proteomics and transcriptomics 369 samples presented some of the lowest correlations. Even so, the correlation values were not randomly 370 distributed, but substantially influenced by which proteomics time-point they were related to. As an 371 example, time-point 24h presented the highest values regardless of the transcriptomics time-point 372 (avg.: 0.46), while time-point 2 hours presented the lowest (avg.: 0.26). To a lower extent, a similar 373 effect was observed for the transcriptomics time-points. When evaluating TrT at the same level, we saw a similar low correlation with proteomics (Figure 3B), confirming the fact that a representative 374 375 portion of transcripts (80%) presented equal values between TPM and TrT, due to not presenting any 376 predicted inhibition by any miRNA expressed in our dataset.



377

Figure 3: Correlation analysis between Transcriptomics and proteomics. A: TPM vs Proteomics. B: TrT vs Proteomics.
 TPM: Transcripts Per Million. TrT: Translatable Transcripts. Timepoint samples are compared across each other (from 2 to 336 hours). The darker and bigger a circle is, the higher the absolute value of correlation it represents. Blue stands for

381 *positive correlation, red for negative correlation.* 

#### 382 Differential Expression Analysis

383 The high similarity between both correlation analyses was not surprising, as we predicted and 384 formulated for miRNA to have generally a small regulatory effect. Even so, we hypothesized that for 385 a portion of the transcripts regulated by miRNA, TrT would be a better proxy for proteomics. To 386 verify this, we ran our formula through all our data. Later, we investigated in which cases traditional 387 transcriptomics was falsely reporting as a proxy by taking proteomics as a reference. The analysis 388 involved the identification of differentially expressed molecules across doses. Two possible scenarios 389 unfolded: there could be a change in transcriptomics not reflected in proteomics, and vice versa: a 390 change in proteomics that was not reflected in transcriptomics. For example, in fluorouracil's (5-FU) 391 exposure, 5.48% of genes presented a change in transcriptomics and not in proteomics, while 4.02% 392 of them did not present a change while proteomics did. Therefore, we focused on the genes for which 393 TrT could correct such cases.

394

The myosin heavy chain 9 (MYH9) gene was one of them. We could observe a significant increase (p.value < 0.05) (Figure 4A) of the protein at a Toxic dose of 5-FU, while for TPM it was not (p.value

- 397 > 0.05) (Figure 4B). TrT, on the contrary, did reflect a significant increase (p. adjusted < 0.05) (Figure</li>
- 398

4C).

- 399
- 400





405 406

*TrT expression values: Gray: TPM, Black: TrT. C: Boxplots of the expression values for proteomics, TPM, and TrT. TPM: Transcripts Per Million. TrT: Translatable Transcripts.* 







416 *Figure 5: Expression values of TGFBI. A:* Proteomics expression values. *B:* TPM and TrT expression values: Gray: TPM,
417 *Black: TrT. C: Boxplots of the expression values for Proteomics, TPM, and TrT. TPM: Transcripts Per Million. TrT:*418 *Translatable Transcripts.*

#### 419 Biological Interpretation (GOrilla)

420 After observing the effects TrT could have, we investigated whether those could also affect the 421 biological interpretation. Specifically, we focused on comparing the enriched gene ontology sets (GO-422 sets) of both TPM and TrT via GOrilla. The results showed that the number of GO terms was always 423 higher in TPM than in TrT (Figure 6). Such a decrease in GO terms for TrT, though, was not related 424 to a lower number of DEGs from that quantifier. When analyzing the commonalities and differences 425 between TPM and TrT across compounds, we observed 2 different behaviors (Figure 6). For some 426 compounds (5FU, AMI, EPI & MXT), the number of TrT GO terms shared with TPM was either
427 greater or equal to the exclusive ones, while most of the TPM terms were exclusive. On the other
428 compounds, the GO terms shared between them were rather the minority for both sets.

#### 429





431 Figure 6: Number of exclusive GO terms for each quantifier (TPM & TRT) and the ones included in both (BOTH).
432 TPM: Transcripts Per Million. TrT: Translatable Transcripts. 5FU: Fluorouracil. AMI: Amiodarone. CEL: Celecoxib.
433 DOC: Docetaxel. DOX: Doxorubicin. EPI: Epirubicin. MXT: Mitoxantrone. PTX: Paclitaxel.

434 An example of how the misclassification of a DEG may affect the biological interpretation of a 435 comparison (with enough genes being misclassified) was the Cardiac-Specific Homeo Box (NKX2-5) 436 gene. According to TPM, the expression of this gene was significantly higher in the UNTR samples 437 when compared to its corresponding samples in the therapeutic dose in 5FU, AMI, and DOC, leading 438 to the enrichment of the Cardiac Muscle Tissue Morphogenesis GO set (of genes). TrT, though, 439 showed no difference between both scenarios (Figure 7A). This meant that the miRNA regulation was 440 predicted to be stronger in the control samples than the treatment ones, to the point where the 441 difference between both conditions was not significant. To contrast that prediction, we searched for 442 all 18 MTIs with either weak or strong evidence of their inhibitory regulation. We observed that the 443 expression of all those miRNAs in the Control samples was greater or equal to the therapeutic ones, 444 confirming the TrT prediction (Figure 7B).



445

446 Figure 7: NKX2-5 gene: Expression and Regulation between Untreated (UNTR) and Therapeutic (The) samples treated
447 with Fluorouracil (5FU). A. TPM (left) and TrT (right) expression represented in boxplots. B. Examples of the expression
448 in Reads Per Million (RPM) of miRNAs that regulate the NKX2-5 gene.

#### 449 Sensitivity, Specificity, and Accuracy results in TPM and TRT

450 So far, we saw that, at least in some cases, TrT could be a better proxy for proteomics than TPM, and 451 that the biological interpretation was simplified when TrT was applied. Afterward, we wanted to

- assess the global accuracy of TrT and TPM by comparing it to proteomics.
- 453

454 We calculated the accuracy of both transcriptomics values for all compounds in 3 different

455 comparisons: therapeutic versus toxic doses, untreated control versus therapeutic dose, and untreated

- 456 control versus toxic dose. To evaluate the differential expression of the last 2 comparisons, we used
- 457 the same methodology as the one used for therapeutic versus toxic doses.





459 Figure 8: Difference in Sensitivity, Specificity, and Accuracy when changing the quantifier from TPM to TrT, each
460 comparison encompassing 8 compounds: Fluorouracil (5FU), Amiodarone (AMI), Celecoxib (CEL), Docetaxel (DOC),
461 Doxorubicin (DOX), Epirubicin (EPI), Mitoxantrone (MXT), and Paclitaxel (PTX). A: Therapeutic versus Toxic doses. B:

**462** Untreated versus Therapeutic dose. C: Untreated versus Toxic dose.

We observed an ambiguous effect of TrT in relation to TPM in all 3 comparisons from a global perspective (Figure 8). Individually, though, we observed a compound effect reflected on different accuracy differences for each of the 8 compounds, which could be categorized into 3 groups when analyzing the accuracy differences for each compound. Most of them showed ambiguous results, in which both improvement and decrease in accuracy happened. EPI, instead, showed a continuous improvement across all comparisons. The last group, contrarily (CEL & MXT), showed a continuous decline across all comparisons.

#### 470 Discussion

The correlation between transcriptomics and proteomics has always tended to be low. This is true both in previous and in the current study, proving transcriptomics as a poor proxy of the phenotypical changes of the cell. Trying to predict the regulatory effect of miRNA (TrT) to improve such proxy (TPM) proved beneficial for a portion of the gene set. This is unsurprising, knowing that not all genes are regulated by miRNAs, and only a small portion of proteins are quantified through the proteomics pipeline. TrT also showed a simplified biological interpretation of the changes across conditions, although its results compared to TPM varied compound-wise.

- We did not only take miRNAs into account, though. Other ncRNAs, such as circRNAs (although recent studies have observed some of them as coding), were also quantified. Due to the novelty of circRNAs, in contrast with other OMICs, there was no golden standard database to refer to when seeking their specific IDs nor the way they should be quantified, among other reasons, because the existing databases were seldom updated<sup>18,49</sup> nor online<sup>50</sup>.
- 484

485 One of the most recent recommendations<sup>44</sup> in these circumstances was the combined use of at least 2 486 different circRNA identification software tools. The reason behind this procedure was based on the 487 flawed accuracy of these programs due to their inherent biases. Fortunately, even if this bias differed 488 for every tool, it had been shown that the convergence between the 2 tools offered a much lower 489 proportion of false positives in contrast to their original population of predictions<sup>44</sup>.

490

491 We decided to pool all samples so that we had the maximum sequencing depth available for the 492 identification process. Otherwise, selecting a random sample would lead to an under-representation of 493 some circRNA molecules, which would not be identified as such, even if they would be in the rest of 494 the samples.

495

496 The argumentation behind the selection of CIRI2 and circExplorer2 was twofold. First, their accuracy 497 outcomes in previous reviews portrayed them to be some of the best predictors currently available. 498 Secondly, our comparison of several predictors revealed that both output a substantial proportion of 499 molecules that were also outputted by other predictors (low false-positive rate), without sacrificing sensitivity. This was not the case for circRNA finder<sup>32</sup>, which found several times more molecules 500 501 than the other predictors with a very low overlap ( $\sim 7\%$ ). find circ<sup>18</sup>, although having a similar 502 proportion of overlap as circExplorer2 and CIRI2, did not present such a good accuracy in previous 503 studies.

504

505 The model we applied considered several factors in post-transcriptional regulation. We built it by 506 trying to reach an equilibrium at the complexity level: enough complexity to maximize the 507 representation of the molecular reality of the cell, but without including so many variables that would 508 have made the model either unbuildable or unfeasible to run in practical terms. Our model also 509 reflected how little is still known about post-transcriptional regulation, and recent discoveries are 510 constantly being made public around it<sup>42,51</sup>.

511 This is especially the case with miRNA regulation, where the complementarity and regulatory 512 strength seem to not only be dependent on the seed region but other loci such as the 3' end of 513 miRNAs<sup>42</sup>. These new observations are currently being incorporated into the interaction predictive 514 tools, making probably in the future the more classical tools obsolete.

- 516 The genes exemplified as good proxies for proteomics are indeed regulated by miRNAs. MYH9 is a
- 517 non-muscle myosin chain IIA gene which plays a major role in early mammalian development, while
- 518 in the adult heart is only expressed in the non-myocyte cells<sup>52</sup>. Several miRNAs have been either
- 519 positively or negatively associated with the gene's expression<sup>53</sup>. For TGFBI, miRNA regulation has
- 520 also been shown, to the point of even being related to chemoresistance<sup>54</sup>.
- 521

522 Our model was based on miRNA regulation on a 1-on-1 basis: if a miRNA can bind a transcript (after 523 taking into account the endogenous competition), the latter will not be able to be translated. Other 524 recent models prefer to focus on the targeting efficacy (as a function of the affinity between the two 525 molecules) while ignoring the ceRNAs<sup>51</sup>. We acknowledge that both models lack each other's 526 strengths, and a combination of the two would have been optimal.

527

In addition to that, there could have been other important factors that could have been taken into account. One example could be the translation ratio, i.e., the number of proteins translated per transcript, which depends on the number of ribosomes that may bind to it, and its half-life. Indeed, half-lives, both of the transcript and the protein, make it difficult to transient from one OMICs to another: the increase in proteomics might be by an increase in transcriptomics with a short half-life that cannot be seen, or a long half-life of proteins may delay the decrease in their abundance due to a decrease in translational output.

535

Proteomics, although being the technology used as a reference to represent the phenotypical changes in the cell, is quite limited in achieving so. This is especially the case due to its lack of sensitivity, where generally only around 1000-3000 proteins are quantified, making the analysis of the proteome quite challenging. In addition to that, the output scale of such quantifications can be quite dramatic when the technology or machinery involved changes.

541

542 Our next steps would include trying to breach the gap between both OMICs by, not only including 543 important parameters that can improve the proteomics prediction, but also changing the algorithms 544 involved in the model, namely, using a supervised quantitative machine learning model. That way, the 545 use of new parameters is easier to both implement and evaluate. This is especially relevant for the 546 imputation of proteomics data, prone to missing values due to the stochastic nature of mass 547 spectrometry. Our work concerning this strategy can be found in our published article<sup>57</sup>.

548 We created a formula to model post-transcriptional regulation and its effects on protein expression. 549 We assessed its behavior in both general and specific cases. As expected, due to the variability of each 550 gene's regulation, the formula benefited only a subset of genes. We portrayed the importance of 551 reaching a more accurate description of the cellular changes, and how far we still are from that objective. Therefore, we must not simply continue with the most popular technology, but try to reacha better approach to the current era of big data.

# 554 Supplementary data

555 Supplementary data can be found in the online preprint version of this manuscript<sup>58</sup>.

## 556 Bibliography

- Zhang, Z., Wu, S., Stenoien, D. L. & Paša-Tolić, L. High-throughput proteomics. *Annu. Rev. Anal. Chem. (Palo Alto. Calif).* 7, 427–54 (2014).
- Peng, J., Elias, J. E., Thoreen, C. C., Licklider, L. J. & Gygi, S. P. Evaluation of
  multidimensional chromatography coupled with tandem mass spectrometry (LC/LC-MS/MS)
  for large-scale protein analysis: The yeast proteome. J. Proteome Res. 2, 43–50 (2003).
- 3. Washburn, M. P., Wolters, D. & Yates, J. R. Large-scale analysis of the yeast proteome by
  multidimensional protein identification technology. *Nat. Biotechnol.* 19, 242–247 (2001).
- 564 4. Cagney, G. *et al.* Human tissue profiling with multidimensional protein identification
  565 technology. *J. Proteome Res.* 4, 1757–67 (2005).
- 566 5. Chen, G. *et al.* Discordant protein and mRNA expression in lung adenocarcinomas. *Mol. Cell.*567 *Proteomics* 1, 304–13 (2002).
- 568 6. Lemée, J.-M. *et al.* Integration of transcriptome and proteome profiles in glioblastoma: looking
  569 for the missing link. *BMC Mol. Biol.* 19, 13 (2018).
- 570 7. Rogers, S. *et al.* Investigating the correspondence between transcriptomic and proteomic
  571 expression profiles using coupled cluster models. *Bioinformatics* 24, 2894–2900 (2008).
- 572 8. Dhingra, V., Gupta, M., Andacht, T. & Fu, Z. F. New frontiers in proteomics research: A
  573 perspective. *International Journal of Pharmaceutics* vol. 299 1–18 (2005).
- 574 9. Belle, A., Tanay, A., Bitincka, L., Shamir, R. & O'Shea, E. K. Quantification of protein half575 lives in the budding yeast proteome. *Proc. Natl. Acad. Sci. U. S. A.* 103, 13004–13009 (2006).
- 576 10. Baek, D. et al. The impact of microRNAs on protein output. Nature 455, 64–71 (2008).
- 577 11. Selbach, M. *et al.* Widespread changes in protein synthesis induced by microRNAs. *Nature*578 455, 58–63 (2008).
- 579 12. Ambros, V. The functions of animal microRNAs. *Nature* 431, 350–355 (2004).
- 580 13. Lewis, B. P., Shih, I., Jones-Rhoades, M. W., Bartel, D. P. & Burge, C. B. Prediction of
- 581 Mammalian MicroRNA Targets. *Cell* **115**, 787–798 (2003).
- 582 14. Lim, L. P. *et al.* Microarray analysis shows that some microRNAs downregulate large
  583 numbers of-target mRNAs. *Nature* 433, 769–773 (2005).
- 584 15. Zaphiropoulos, P. G. Exon skipping and circular RNA formation in transcripts of the human

- 585 cytochrome P-450 2C18 gene in epidermis and of the rat androgen binding protein gene in
  586 testis. *Mol. Cell. Biol.* 17, 2985–2993 (1997).
- 587 16. Chen, L. L. & Yang, L. Regulation of circRNA biogenesis. RNA Biol. 12, 381–388 (2015).
- 588 17. Jeck, W. R. *et al.* Circular RNAs are abundant, conserved, and associated with ALU repeats.
  589 *RNA* 19, 141–157 (2013).
- 590 18. Memczak, S. *et al.* Circular RNAs are a large class of animal RNAs with regulatory potency.
  591 *Nature* 495, 333–338 (2013).
- 592 19. Rong, D. *et al.* An emerging function of circRNA-miRNAs-mRNA axis in human diseases.
  593 *Oncotarget* 8, (2017).
- Wang, K. *et al.* The long noncoding RNA CHRF regulates cardiac hypertrophy by targeting
  miR-489. *Circ. Res.* 114, 1377–88 (2014).
- 596 21. Salmena, L., Poliseno, L., Tay, Y., Kats, L. & Pandolfi, P. P. A ceRNA hypothesis: the Rosetta
  597 Stone of a hidden RNA language? *Cell* 146, 353–8 (2011).
- 598 22. Kuepfer, L. *et al.* A model-based assay design to reproduce in vivo patterns of acute drug599 induced toxicity. *Archives of Toxicology* vol. 92 553–555 (2018).
- Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence
  data. *Bioinformatics* 30, 2114–20 (2014).
- Baras, A. S. *et al.* miRge A Multiplexed Method of Processing Small RNA-Seq Data to
  Determine MicroRNA Entropy. *PLoS One* 10, e0143066 (2015).
- Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient
  alignment of short DNA sequences to the human genome. *Genome Biol.* 10, R25 (2009).
- 26. Zhang, X.-O. *et al.* Diverse alternative back-splicing and alternative splicing landscape of
  circular RNAs. *Genome Res.* 26, 1277–87 (2016).
- 608 27. Gao, Y., Zhang, J. & Zhao, F. Circular RNA identification based on multiple seed matching.
  609 *Brief. Bioinform.* 19, 803–810 (2018).
- 610 28. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
  611 *Bioinformatics* vol. 25 1754–1760 https://pubmed.ncbi.nlm.nih.gov/19451168/ (2009).
- 612 29. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic
  613 features. *Bioinformatics* 26, 841–2 (2010).
- 614 30. ENSEMBL. FTP Download. https://www.ensembl.org/info/data/ftp/index.html (2020).
- Batro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and
  bias-aware quantification of transcript expression. *Nat. Methods* 14, 417–419 (2017).
- 617 32. Westholm, J. O. *et al.* Genome-wide Analysis of Drosophila Circular RNAs Reveals Their
  618 Structural and Sequence Properties and Age-Dependent Neural Accumulation. *Cell Rep.* 9,
- 619 1966–1980 (2014).
  620 33. Griffiths-Jones, S., Grocock, R. J., van Dongen, S., Bateman, A. & Enright, A. J. miRBase:
  621 microRNA sequences, targets and gene nomenclature. *Nucleic Acids Res.* 34, D140–D144

622		(2006).
623	34.	Durinck, S. M. Y. et al. BioMart and Bioconductor: a powerful link between biological
624		databases and microarray data analysis. Bioinformatics 21, 3439-3440 (2005).
625	35.	Li, H. et al. The Sequence Alignment/Map format and SAMtools. Bioinformatics 25, 2078-
626		2079 (2009).
627	36.	John, B. et al. Human microRNA targets. PLoS Biol. 2, (2004).
628	37.	R Core Team. R: A Language and Environment for Statistical Computing. (2019).
629	38.	RStudio Team. RStudio: Integrated Development Environment for R. (2019).
630	39.	Wei, T. & Simko, V. R package 'corrplot': Visualization of a Correlation Matrix. (2017).
631	40.	Analyzing RNA-seq data with DESeq2.
632		https://bioconductor.org/packages/release/bioc/vignettes/DESeq2/inst/doc/DESeq2.html#why-
633		un-normalized-counts.
634	41.	Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression
635		analysis of RNA-seq data. Genome Biol. 11, 1-9 (2010).
636	42.	Chipman, L. B. & Pasquinelli, A. E. miRNA Targeting: Growing beyond the Seed. Trends in
637		Genetics vol. 35 215–222 (2019).
638	43.	Chou, C. H. et al. MiRTarBase update 2018: A resource for experimentally validated
639		microRNA-target interactions. Nucleic Acids Res. (2018) doi:10.1093/nar/gkx1067.
640	44.	Hansen, T. B. Improved circRNA identification by combining prediction algorithms. Front.
641		<i>Cell Dev. Biol.</i> <b>6</b> , 20 (2018).
642	45.	Bartel, D. P. MicroRNAs: target recognition and regulatory functions. Cell 136, 215-33
643		(2009).
644	46.	Arvey, A., Larsson, E., Sander, C., Leslie, C. S. & Marks, D. S. Target mRNA abundance
645		dilutes microRNA and siRNA activity. Mol. Syst. Biol. 6, 363 (2010).
646	47.	Franco-Zorrilla, J. M. et al. Target mimicry provides a new mechanism for regulation of
647		microRNA activity. Nat. Genet. 39, 1033-7 (2007).
648	48.	MS, E., JR, N. & PA, S. MicroRNA Sponges: Competitive Inhibitors of Small RNAs in
649		Mammalian Cells. Nat. Methods 4, (2007).
650	49.	Li, J. H., Liu, S., Zhou, H., Qu, L. H. & Yang, J. H. StarBase v2.0: Decoding miRNA-ceRNA,
651		miRNA-ncRNA and protein-RNA interaction networks from large-scale CLIP-Seq data.
652		Nucleic Acids Res. 42, D92-7 (2014).
653	50.	Liu, Y. C. et al. CircNet: A database of circular RNAs derived from transcriptome sequencing
654		data. Nucleic Acids Res. 44, D209–D215 (2016).
655	51.	McGeary, S. E. et al. The biochemical basis of microRNA targeting efficacy. Science (80 ).
656		<b>366</b> , eaav1741 (2019).
657	52.	Ma, X. & Adelstein, R. S. In vivo studies on nonmuscle myosin II expression and function in
658		heart development. Front. Biosci. 17, 545-555 (2012).

- 53. Yu, M. *et al.* Prognostic impact of MYH9 expression on patients with acute myeloid leukemia. *Oncotarget* 8, 156–163 (2017).
- 661 54. Bissey, P. A. *et al.* Dysregulation of the MiR-449b target TGFBI alters the TGFβ pathway to
  662 induce cisplatin resistance in nasopharyngeal carcinoma. *Oncogenesis* 7, (2018).
- 663 55. Nakamoto, M., Jin, P., O'Donnell, W. T. & Warren, S. T. Physiological identification of
- human transcripts translationally regulated by a specific microRNA. *Hum. Mol. Genet.* 14,
  3813–3821 (2005).
- 666 56. Moore, M. J. *et al.* MiRNA-target chimeras reveal miRNA 3'-end pairing as a major
  667 determinant of Argonaute target specificity. *Nat. Commun.* 6, 1–17 (2015).
- 668 57. Ochoteco Asensio J, Verheijen M, Caiment F. Predicting missing proteomics values using
  669 machine learning: Filling the gap using transcriptomics and other biological features. *Comput*
- 670 *Struct Biotechnol J.* 2022 Apr 22;20:2057-2069. doi: 10.1016/j.csbj.2022.04.017. PMID:
- 671 35601960; PMCID: PMC9077535.
- 672 58. Ochoteco Asensio, J., Kleinjans, J. & Caiment, F. Quantifying the number of translatable
  673 transcripts through the use of OMICs involved in post-transcriptional regulation. *bioRxiv*,
  674 doi:<u>https://doi.org/10.1101/2022.06.20.496876</u> (2022).

675

676
### Chapter 4:

## Predicting missing proteomics values using machine learning: Filling the gap using transcriptomics and other biological features

# Predicting missing proteomics values using 2 machine learning

3 Filling the gap using transcriptomics and other biological features

#### 4

- 5 Juan Ochoteco Asensio<sup>1</sup>, Marcha Verheijen<sup>1</sup>, Florian Caiment<sup>1,\*</sup>
- Department of Toxicogenomics, School of Oncology and Developmental Biology (GROW), Maastricht University,
   Maastricht, The Netherlands
- 8 \* Corresponding Author: florian.caiment@maastrichtuniversity.nl
- 9 (Manuscript <u>published</u>)

#### 10 **ABSTRACT**

Proteins are often considered the main biological element in charge of the different functions and 11 structures of a cell. However, proteomics, the global study of all expressed proteins, often performed 12 by mass spectrometry, is limited by its stochastic sampling and can only quantify a limited amount of 13 protein per sample. Transcriptomics, which allows an exhaustive analysis of all expressed transcripts, is 14 often used as a surrogate. However, the transcript level does not present a high level of correlation with 15 the corresponding protein level, notably due to the existence of several post-transcriptional regulatory 16 mechanisms. In this publication, we hypothesize that the missing protein values in proteomics could be 17 predicted using machine learning regression methods, trained with many features extracted from 18 19 transcriptomics, including known translational regulatory elements such as microRNAs and circular RNAs. After considering different machine learning algorithms applied on two different splitting 20 strategies, we report that random forest can predict proteins in new samples out of transcriptomics data 21 with good accuracy. The proposed pre-processing and model building scripts can be accessed on 22 23 GitHub: https://github.com/jochotecoa/ml proteomics

#### 24 **1 INTRODUCTION**

For a cell to react and adapt to any variation of its environment, including for instance the exposure to a 25 foreign compound, a cascade of events leading ultimately to the production of proteins occurs. For that 26 27 purpose, the cell usually initiates the transcription of its genes (such as transcription factor), and the resulting transcripts containing an open reading frame are translated into proteins. Even though such a 28 schematic view of molecular biology appears straightforward, each of those steps is controlled and 29 affected by a myriad of factors. This complexity led to the development of advanced technologies, 30 named "omics", allowing to deeply study a particular class of biological entity: transcriptomics 31 (characterization and quantification of transcripts), proteomics (proteins), metabolomics (metabolites), 32 33 etc.

Among those different classes of molecules, proteins are particularly relevant, as their expression level 34 35 and activity inform profoundly about how the cell is functioning and reacting to its environment, especially when those changes may pose a risk to the integrity and functionality of the whole system, 36 37 either due to a disease or an infection. To analyze the expression of proteins in different conditions, proteomics (mass spectrometry or MS) is usually applied. Unfortunately, its sensitivity is limited<sup>1-3</sup>, and 38 39 thus only a small subset of proteins (with the highest abundance) can be studied at a time. In addition, the stochastic sampling generates missing identifications across samples, particularly for proteins with 40 an abundance close to the detection limit; even though workflows such as DIA (Data-Independent 41 Acquisition)-MS workflow can increase reproducibility. New technologies are not exempt of these 42 43 limitations: the latest single-cell proteomics strategies (such as SCoPE2<sup>2</sup>) and newest experimental and computational workflows<sup>3</sup> only obtain ~1000 proteins per cell on average (not including their own 44 limitations<sup>4</sup>), even though their dynamic range allows for the quantification of 3000 distinct proteins. 45

Proteins are mainly translated from messenger RNAs (mRNAs), which are much easier to analyze. 46 Indeed, while having a shorter half-life than proteins, mRNA transcriptomics has become 47 overwhelming sensitive and cost-efficient over the years with the invention of next-generation 48 49 sequencing. For these reasons, RNA-Sequencing techniques are usually preferred to statistically study cell changes at the molecular level. However, a given mRNA is not an excellent proxy of its 50 corresponding protein expression level, which is reflected in a very low correlation between 51 transcriptomics and proteomics technologies<sup>5-9</sup>. While the reasons behind this gap can be multiple, the 52 main factors can be categorized into post-transcriptional regulation. By different mechanisms in such 53 regulation, the cell controls the final level of translation of each mRNA into proteins. These factors can 54

be either determined by the molecules themselves (such as the transcript's or protein's half-life<sup>10</sup>) or by
the interaction with external elements.

57 MicroRNAs (miRNAs), short non-coding transcripts of around 22 nucleotides of length, play an 58 important role in post-transcriptional regulation. They can act as inhibitors of translation<sup>11,12</sup> by base-59 pairing their seed region<sup>13</sup> (nucleotide 2 to 8) to the target mRNA, usually in their 3` UTR region. 60 While often considered mild individually, the interaction of multiple miRNAs (either the same miRNA 61 or different miRNAs) on the same 3'UTR target can have a significant effect on protein level 62 expression <sup>14,15</sup>. Considering the relatively short length of the seed region, miRNAs can target an 63 average of 200 different targets. Even so, miRNAs are not the only transcripts regulating translation.

Another newly discovered category of RNAs, named circular RNAs (circRNAs), are characterized by 64 65 their circular form, which is generated by the binding of their 5 and 3' end during splicing (backsplicing)<sup>16,17</sup>, forming the so-called back-spliced junction. Due to this particular structure, they are not 66 easily degraded due to the absence of transcript extremities, rendering them immune to exonuclease 67 activity<sup>18</sup>. Several functions have been proposed for these circRNAs, including regulating miRNA 68 activities. It has been demonstrated that circRNAs, which can contain repetitions, could present the 69 same target regions present in miRNA targets, and sometimes several times per molecule. This leads to 70 a target competition<sup>19</sup>, where circRNAs bind most miRNAs, which gave to circRNAs the function of 71 'miRNA sponges'<sup>20</sup>. The post-transcriptional regulation complexity starts to unfold once one realizes 72 73 that each transcript can be inhibited by several miRNAs, and at the same time, each of those miRNAs can be "sponged" by one or more circRNAs. 74

75 The final expression level of a protein results thus from the integration inside the cell of many factors related to transcripts: the level of expression of mRNAs, the number of possible seeds with miRNAs, 76 77 the expression level of miRNAs, and the expression level (and "sponging" capacity) of circRNAs able to capture these miRNAs. Many other features could also play a role in this final protein expression 78 level. For instance, the GC content of an mRNA has been observed to interfere with the mRNA half-79 life<sup>21</sup>, and thus the total number of proteins formed from a single mRNA. All these RNA elements or 80 characteristics just mentioned could be identified and quantified by transcriptomics with RNA-Se-81 quencing. Since the protein expression level is the most important factor for biological interpretation, 82 and considering the limited sensitivity and stochastic sampling of proteomics in addition to the very 83 low correlation of the mRNA/protein expression level, we considered the possibility of obtaining pre-84 dicted protein expression levels from the integration of as many possible features available from sev-85

eral OMICs data. Although we recognized that methods such as match-between-runs (MBR)<sup>22</sup>, DARTID<sup>23</sup>, and IceR<sup>24</sup> have already been developed (and their limitations<sup>25</sup>), including a deep learning approach to extrapolate proteomics values from transcriptomics values<sup>26</sup>, none utilized a complex multiomics strategy to approach in a novel manner the limitations of proteomics.

90 The amount and complexity of the data render impossible the task of manually integrating all these 91 parameters. Even when inputting such data digitally, it is not straightforward to visualize which is the optimal manner to predict proteomics values. This problem is characteristic of the current big data era, 92 which in turn, has led to the rise of algorithms that use straightforward optimization strategies to 93 rapidly process thousands or millions of observations. Some of those can be categorized as machine 94 learning (ML), which consists of a set of computer algorithms built to automatically improve their pre-95 diction with increasing volumes of data<sup>27</sup>. Specifically, the algorithms focused on predicting are part of 96 the supervised learning algorithms, as they require a training phase in which they are exposed to the 97 value to be predicted (target) in conjunction with other variables associated with it (features). Two ma-98 jor classes of machine learning algorithms exist: when predicting categories or labels (qualitative val-99 ues), algorithms will perform a classification; while when what is predicted are quantitative values, 100 algorithms will perform a regression. The improvement in the accuracy of these models can be evalu-101 ated based on how similar the predictions are to the actual observations. The accuracy is only relevant 102 to evaluate with new data (testing dataset), and not with the data used to train the model (training 103 104 dataset), in order to avoid the generation of a biased model due to overfitting.

In this manuscript, we hypothesized that using machine learning algorithms would allow us to estimate 105 106 the expression level of the protein not detected by proteomics out of all available data. For the omics data, we made use of an in vitro dataset obtained from primary human hepatocytes microtissues which 107 includes 3 omics datasets obtained from the exact same samples batch: RNA-Seq (ribo-depleted li-108 braries), miRNA-Seq (small RNA libraries), and proteomics (mass spectrometry). Both mRNA and 109 circRNAs quantification were extracted from the RNA-Seq data. We thus assessed the accuracy of di-110 verse machine learning predictive models based on different algorithms and data-splitting strategies 111 with the ultimate goal to predict protein expression value from transcriptomics and other mRNA fea-112 tures. 113

#### 114 2 METHODS

#### 115 2.1 DATASET & FEATURES

116 The description of the biological samples used, in addition to the proteomics and transcriptomics 117 protocols followed to obtain protein and RNA expression values, can be found in the Supplementary 118 Methods.

Proteomics expression values were set as the target to be predicted. We set as features protein 119 properties with nominal values extracted from UniProt that might affect their half-life. The features 120 were the following: protein length (Length), mass (Mass), quantity of each amino acid (Aa X), 121 organism (Organism), location on which the original gene was encoded (Gene.encoded.by), and the 122 database version of the protein sequence (Version..sequence.). From those, we also derived additional 123 124 features: linear density (mass divided by length) and proportion of each amino acid based on the protein's length (Aa X prop). Finally, we added some irrelevant features (protein sequence version) as 125 negative controls to inform us of the model reliability (based on the importance these features would be 126 given by those models). Concerning protein stability, we included all nine features extracted from the 127 supplementary table: R1-R7, PSI, and SD. 128

129 The expression values (in TPM) of protein-related transcripts were added as a feature. Furthermore, we also added diverse transcript properties: strand, transcript length, percentage gene GC content, CDS 130 length, UTR length (or non-CDS length), and proportion of UTR length (UTR length divided by the 131 transcript length). MiRNA expression was also added as a feature, linking it to the transcript targets 132 they could potentially regulate. For this, we used the miRDB's MiRNA Target Interaction (MTI) score 133 in two features in the ML algorithm: one feature with only miRNAs that presented a high probability 134 of targeting such target ('stringent', score  $\geq 80$ ), and another considering all possible regulations, 135 independently of their score ('all'). CircRNA expression as a feature ('circ') was linked to the 136 proteomics values based on the miRNA sponging effect of the former. We only utilized the expression 137 of those circRNAs that presented more than 7 targeting sites with a specific miRNA. We also added the 138 sponging effect of circRNAs as the feature 'circ score'. 139

Transcripts were named based on their Ensembl ID, while proteins were labeled with UniProt IDs. A single UniProt protein could be associated with more than one ENST transcript, potentially with very different features (expression level, transcript length, etc.). Therefore, we needed to summarize the value from all linked transcripts in a single feature. As there was no clear advantage to select a

particular summary method over another, we created a feature for each of those different methods: 144 mean, median, minimum, maximum, sum, and standard deviation. This approach was not only applied 145 to features associated with transcripts coupled to proteins (and their log2-transformed values) but also 146 to the ones associated with miRNAs and circRNAs (and their log2-transformed values as well). Indeed, 147 this problem was also applicable to those molecules (to even a greater extent) when linking them to a 148 single proteomics value: each transcript can be inhibited by several miRNAs and each of those 149 miRNAs can be sponged by several circRNAs. We also extended this strategy to those features that 150 presented a multiplicity of values for a single observation, such as the protein stability data. The 151 combination of all discussed variables led to a total of 196 features. 152

#### 153 2.2 PRE-PROCESSING

For both the pre-processing of the data and the construction of the machine learning models, we used
the R library 'caret'<sup>28</sup>.

#### 156 2.2.1 Creating dummy variables

Since categorical data (such as gender) cannot be inputted directly into a model, they needed to be transformed into dummy variables. Dummy variables are binary features that indicate the presence (1) or absence (0) of a categorical value. In our data, the dummy variables created were related to strand information (positive (+) or negative (-) strand) and protein version sequence (presence or absence of versions 1 to 7).

#### 162 2.2.2 Identifying (Near) Zero-Variance and Correlated Predictors

163 To identify variables with no variance (Zero-Variance or ZV) or insignificant variance (Near Zero-164 Variance), we used the function 'nzv' (frequency ratio > 95/5, percentage unique < 10%) described in 165 'caret'<sup>29</sup>. We then discarded those predictors from the dataset. To identify correlated variables 166 (correlation > 0.75), we used the function 'findCorrelation' also from the 'caret' package. We 167 discarded the identified correlated predictors from the dataset. The correlation plot was designed using 168 the 'corrplot' package.

#### 169 2.2.3 Centering and Scaling

170 Centering refers to the data transformation where the means of all features are set to a specific value 171 (i.e., 0) while scaling refers to the transformation where the standard deviation is also set to a constant 172 value (i.e., 1). These data transformations avoid a feature importance bias due to value size or scale. No imputation was performed, but instead, all observations with any missing value were removed from thedataset.

#### 175 2.2.4 Data splitting and algorithms used

The data split between the training dataset (80% of the whole dataset) and the testing dataset (20% of 176 the whole dataset) was performed based on 2 different strategies: sample names and protein names. For 177 each algorithm used, we performed recursive feature elimination using the 'rfe' function with (10-fold) 178 cross-validation (CV) resampling and the training dataset. After recursive feature elimination, the 179 model with the optimal subset size of variables for each algorithm was selected to predict the testing 180 dataset. As validation, we also used 'rfe' (10-fold cross-validation) for the whole dataset. To split the 181 dataset accordingly, we first generated the 10 folds using the 'groupKFold' function based on the 182 indicated categories (samples and proteins). These folds were used as input in the 'folds' parameter in 183 the 'rfeControl' function. 184

The algorithms tested were: Boosted Tree ('bstTree'), Random Forest ('rf'), Bagged Model ('bag'), Boosted Tree ('blackboost'), Lasso and Elastic-Net Regularized Generalized Linear Model ('glmnet'), k-Nearest Neighbors ('kknn'), Cubist ('cubist'), and Linear Regression ('lm'). All algorithms were used via 'caret', and thus, the default parameters used by 'caret' were utilized.

#### 189 2.2.5 Performance based on GO terms

We selected the cardiac dataset, and subselected one sample as testing dataset, while the model training 190 191 was proceeded with the rest of samples using the 10 features shown in the results. After the training, we predicted the testing dataset with the resulting random forest model, and combined the predictions with 192 the testing observations. We then extracted the GO terms associated for each protein in the testing 193 dataset, which we also combined with the observations and predictions. We discarded GO terms that 194 were categorized in less than 10 proteins. We evaluated the R<sup>2</sup> metrics for each of the groups of pro-195 teins associated to each GO term. We ranked the GO term groups from best to worst performing based 196 197 on R<sup>2</sup>. All the code used can be located in the following script on GitHub: 'script/go terms analysis/ rsquared on different go terms.R'. 198

#### 199 2.2.6 Imputation: a potential use of the random forest model

We also selected the cardiac dataset, but in this case including all proteomics missing values. We subselected all Untreated (UNTR) samples. The training dataset only contained observations with quantified proteomics values, and the 10 features mentioned in the results. We used the random forest algorithm for the training of the model. We then predicted the missing proteomics values using the newly trained model. We combined the results with the observed data, and sampled proteins with different proportions of missing data. All code run can be found in the following script on GitHub: 'script/imputation/imputing cardiac values.R'.

#### 207 **3 RESULTS**

To assess the ability of the regression ML algorithm to estimate the level of proteins, we produced a dataset that presented the added value of having transcriptomics (both ribo-depleted and small RNA libraries) and proteomics (LC/MS), all generated from the exact same sample batches to maximize the interpretability of the interactions. This dataset was composed of a total of 115 *in vitro* samples (61 cardiac and 54 hepatic). The processing of all these samples (Methods) characterized an amount of expressed biological entities summarized in Table 1. The total number of expressed biological entities was 48 266 and 48 715 for the hepatic and cardiac tissues respectively.

215 *Table 1: Summary table of all quantified biological entities. Total refers to all possible entities to be* 

216 identified. Expressed (N) refers to the number of entities that were quantified in at least 1 sample. Con-

217 stitutive (N) refers to the number of entities that were quantified in all samples. Expressed (%) and

218 *Constitutive (%) refer to the percentage of (constitutively) expressed entities based on the total number* 

219 of entities. Constitutive (% Expressed) refers to the percentage of constitutively expressed entities

220 *based on the number of expressed entities.* 

	Tissue	Total	Ex-	Constitu-	Ex-	Constitutive	Constitutive
			pressed	tive (N)	pressed	(%)	(% Expressed)
			(N)		(%)		
Proteomics	Hepatic	1806	1806	283	100.00%	15.67%	15.67%
Proteomics	Cardiac	2217	2217	247	100.00%	11.14%	11.14%
Linear transcripts	Hepatic	211939	135655	894	64.01%	0.42%	0.66%
Linear transcripts	Cardiac	211939	136860	933	64.58%	0.44%	0.68%
MicroR- NAs	Hepatic	2744	1561	280	56.89%	10.20%	17.94%

MicroR-	Cardiac	2744	1510	250	55.03%	9.11%	16.56%
NAs							
Circular-	Hepatic	140317	95106	151	67.78%	0.11%	0.16%
ized tran-							
scripts							
Circular-	Cardiac	140317	100416	156	71.56%	0.11%	0.16%
ized tran-							
scripts							

222 To assess the possibility to predict protein expression levels for all genes using ML algorithm, we needed to assemble a list of features, either parametric or categorical. From all the table 1 data, we ex-223 tracted 12 features focused on the expression level of linear transcripts, 24 features on miRNA expres-224 sion, and 12 features on circular RNA expression. We added 36 features on transcript characteristics 225 226 (strand, transcript length, etc.), 48 features on protein characteristics (Protein Mass, Protein Length, etc.), 12 features on MTI (miRNA target interaction) scores, a feature on RNA-Sequencing depth, 6 227 features on circular scores (number of miRNA binding site per circular RNA), 12 features on circular 228 RNA expression, and 45 features on protein stability. This led to a total of 196 features on the raw 229 dataset. Even so, some of those features might be deemed irrelevant due to their multiplicity and inher-230 ent structure. Those features might affect machine learning processes, depending on the algorithms' 231 inherent functionality, by decreasing their accuracy<sup>30</sup>. To avoid their inclusion, we applied several pre-232 processing filters that removed non-informative features, which are described below. 233

#### 234 **3.1** ZERO- AND NEAR ZERO-VARIANCE VARIABLES

Some predictors can have a unique value for all observations (Species: Human), which can make 235 models unstable or decrease their fitness. Those features can be named as Zero-Variance (ZV) 236 variables, and they are generally removed. Similarly, Near Zero-Variance (NZV) variables refer to 237 features that present a value in an overwhelming majority of observations (i.e., genes coded in the 238 nucleic genome vs genes coded in the mitochondrial DNA (Table 2)). These features are generally not 239 helpful in a cost/benefit ratio, as the underrepresented values might have an artificially bigger impact, 240 and these values may not even appear in the subpopulations generated by sub-sampling strategies, 241 generating a ZV variable. 242

Table 2: Examples of Zero- and Near Zero-Variance variables. The 'Organism' variable contains a
single unique value ('Human'), thus this value has no predicting value. The 'Gene encoded by' variable contains 2 possible values, of which 'Nucleus' represents more than 99% of all observations.
Even though this variable does indeed have more than a single value, the frequency of its values renders it non-informative.

Protein ID	Organism	Gene encoded by
A – Sample 1	Human	Nucleus
B – Sample 1	Human	Nucleus
C – Sample 1	Human	Nucleus
D – Sample 2	Human	Nucleus
E – Sample 2	Human	Nucleus
F – Sample 2	Human	Nucleus
G – Sample 3	Human	Mitochondrion

248

Due to both ZV and NZV filters, 44 features were removed from the dataset. Only a few were labeled as ZV, examples of which were 'circ\_min' (minimum circular expression) and its log2 transformed version 'circ\_min\_log2'. Some categories of variables were frequently labeled as NZV: almost all features related to miRNA scores; all maximum, median, and minimum miRNA expressions (nontransformed, log2-transformed, stringent, and all scores); some related to circular scores and some related to circular expression (Supplementary Table 1).

255 *Table 3: Examples of ZV and NZV features with their respective frequency ratios and unique percent-*256 *ages. The metrics for all NZ features were identical, as they only reported a single value (Inf: Infinite).*

uges. The metrics for an indication were identical, as they only reported a single value (ing. ingitute).

**257** For NZV values, they all presented a frequency ratio above 19 (95/5) and a percentage unique below

**258** *10*.

Feature name	freqRatio	percentUnique	zeroVar	nzv
circ_min	Inf	0.004	TRUE	TRUE
circ_min_log2	Inf	0.004	TRUE	TRUE
Organism	Inf	0.004	TRUE	TRUE

strand_sd	839.433	0.008	FALSE	TRUE
transcript_length_sd	140.014	3.494	FALSE	TRUE
percentage_gene_gc_content_sd	515.167	0.107	FALSE	TRUE
cds_length_sd	64.757	2.098	FALSE	TRUE
noncds_length_sd	112.344	3.427	FALSE	TRUE
proportion_noncds_length_sd	336.033	3.693	FALSE	TRUE
Gene.encoded.by	1198.048	0.008	FALSE	TRUE

#### 260 **3.2** Identification of correlated variables

261 Having correlated predictors is generally uninformative and sometimes detrimental to build models. For this reason, we removed features that presented a correlation above 0.75. For each pair of 262 correlated features, the feature labeled as 'highly correlated' was the one that presented a higher 263 264 correlation with the rest of the variables. Having the target inside the dataset would imply that the 265 features that showed a higher correlation with the target would get removed. To avoid this, we removed 266 the target from the dataset before filtering the highly correlated variables. In total, 93 features were removed due to high correlation (Supplementary Table 2). As expected, the abundances of most amino 267 acids were highly correlated to each other, and to the protein mass and length. The same results were 268 not true for the proportion of each amino acid, as they more accurately represent their presence 269 270 independently of the protein's size. More surprisingly, among all non-filtered features, we observed all possible grouping systems (minimum, mean, median, maximum, standard deviation, and sum of the 271 272 values they represented), with no clear predominance for any of them, and thus none appeared to present a tendency to be the most informative (i.e., the one with the lowest overall correlation with all 273 274 features).



Figure 1: Correlation plot between kept features (horizontal axis) and filtered features (vertical axis).
The scale unit on the right side of the figure indicates the correlation values between the features
shown based on a range of colors: from dark red (extreme negative correlation) to dark blue (extreme
positive correlation), where lighter colors represent a lower absolute correlation value.

#### 280 **3.3** Splitting strategies

In both hepatic and cardiac datasets, the observations were part of two distinct groups: proteins and samples. Having a random split of our data to form both training and testing datasets would not have enabled us to elucidate the actual accuracy of the model. In a random splitting strategy, the training dataset was highly probable to include most proteins and samples in their observations, rendering the data split futile. Instead, we split (and trained) our models separately in two manners: splitting by sample and splitting by protein (Figure 2). This strategy was applied in the hepatic dataset for both training and testing, and in the cardiac dataset for validation.





Figure 2: Splitting strategies. For all splitting strategies, 80% of the data is used to train the models 291

292 (training dataset), while the other 20% is used for testing the trained models (testing dataset). A. Ran-

dom splitting strategy, where the algorithm is trained and tested with observations from all proteins 293

and samples. B. Sample-splitting strategy: the trained models are tested with 20% of the samples. C. 294

*Protein-splitting strategy: the trained models are tested with 20% of the proteins.* 295

#### 296 **3.4** MODEL TRAINING AND TESTING USING HEPATIC SAMPLE-SPLIT DATA

When splitting by sample, 80% of hepatic samples were used as the training dataset, while the other 297 20% was used as the training dataset. For every algorithm, the training dataset was inputted through 298 RFE (cross-validated with 10-fold). Out of all models trained with different subsets of features, the one 299 with the best accuracy was used for the testing step (Figure 3). In terms of root-mean-square error 300 (RMSE, Figure 3A), both k-Nearest Neighbors ('kknn') and Random Forest ('rf') showed the highest 301 302 accuracies ( $\sim 1.25$ ), the latter having a bigger deviation between training and testing RMSE values. To evaluate these results in a more standard and informative manner, we also analyzed the R squared met-303 304 rics (Figure 3B). In this figure, we observed that rf and kknn also showed the best performance ( $R^2$ close to 0.7), showing rf better performance in this case. 305



Figure 3: Accuracy results when splitting by sample. A: The blue bar refers to the RMSE value (left
vertical axis) after training the model with 80% of the samples, and the orange bar refers to the RMSE
value after testing the model with the other 20% of the samples. The gray line refers to the percentual
change of RMSE (right vertical axis) between training and testing. B: The blue bar refers to the R<sup>2</sup>
value (left vertical axis) after training the model with 80% of the samples, and the orange bar refers to the R<sup>2</sup>

- 312 the  $R^2$  value after testing the model with the other 20% of the samples. The gray line refers to the per-
- **313** *centual change of*  $R^2$  *(right vertical axis) between training and testing.*
- After validating the aforementioned results by using RFE (10-fold cross-validation) for the whole dataset (Supplementary Figures 1 and 2), we selected random forest ('rf') as the best performing model when splitting by sample. The optimal subset size of features was 51 features, but after close examination of the RFE results (Supplementary Figure 3), we determined that subset sizes above 10 features had a minimal impact on RMSE. The 10 features were selected based on the ranking of feature importance reported by the RFE analysis (Figure 4).



Figure 4: Top 10 features based on Overall importance by RFE when using the rf algorithm. These
 values represent how important (on average) each feature is to the model, and thus which are the main
 features used by the model to predict new proteomics values.

#### 324 **3.5** MODEL TRAINING AND TESTING USING HEPATIC PROTEIN-SPLIT DATA

Similar to the splitting by sample strategy, a fifth of all proteins were split to be used as the testing dataset, while the other 4 fifths were used as the training dataset. RFE (10-fold CV) was also performed with similar optimal results as in the training dataset of the sample-splitting strategy (Figure 5). In this case, the best RMSEs in the testing dataset include 'bstTree' and 'rf' ( $\approx$ 2), which almost doubled the error shown when splitting by sample (Figure 5A). To understand how relevant this error increase was, we also evaluated the R-squared values of those values (Figure 5B). We observed that a systematic gap

- existed between the training and testing steps, leading to minimal R-squared values ( $R^2 = 0.15$  for rf).



Figure 5: Accuracy results when splitting by protein. A: The blue bar refers to the RMSE value (left
vertical axis) after training the model with 80% of the proteins, and the orange bar refers to the RMSE
value after testing the model with the other 20% of the proteins. The gray line refers to the percentual
change of RMSE (right vertical axis) between training and testing. B: The blue bar refers to the R<sup>2</sup>
value (left vertical axis) after training the model with 80% of the proteins, and the orange bar refers to

338 the  $R^2$  value after testing the model with the other 20% of the proteins. The gray line refers to the per-339 centual change of  $R^2$  (right vertical axis) between training and testing.

For all the results shown above (Figures 3 and 4), we also validated the results using RFE with the whole dataset (no training-testing split), where the folds or splits in the cross-validation step (10-fold) contained exclusively a set of proteins (Supplementary Figures 4 and 5).

#### 343 3.6 RANDOM FOREST MODEL VALIDATION WITH A CARDIAC SAMPLE-SPLIT DATA

Random forest being the best performing model, we decided to validate its accuracy to predict new samples using a cardiac dataset, which was built in the same manner as the hepatic one. The validation included using the same algorithm (rf) with the same top 10 features (Figure 4), and training and validating it with the cardiac data (27602 observations). The resampling was performed via Cross-Validation (10-fold). Using the cardiac data and the specified model, we validated that the accuracy remained robust across different cell types (RMSE = 1.04, R<sup>2</sup> = 0.75; Supplementary Figure 6).

The only remarkable difference was the feature importance ranking given by the RFE in the hepatic data (Figure 4) compared to the feature importance ranking given by the model itself with the cardiac data (Figure 6). In the latter, linear\_density is given the utmost importance, and the importance of the three RNA subtypes relate to how close they are to the protein level: mRNA level, followed by miRNA levels, and finally circRNA levels.



Figure 6: Feature Overall Importance for the rf algorithm when trained in Cardiac data. These values
represent how important (on average) each feature is to the model, and thus which are the main features used by the model to predict new proteomics values.

Therefore, the RMSE and R-squared metrics for both cardiac and hepatic models showed that building a random forest model using the aforementioned features allowed to predict with high accuracy full proteomics' samples. Comparing the testing results between sample- and protein-splitting, we observed that the high accuracy was especially due to the prediction of proteins that have already been trained on. Observing the feature importance ranking (Figure 6), we could observe that different biological entities presented a different relevance to the model's accuracy, thus missing some variables will have a minimal effect on the decided outcome.

#### 367 **3.7** Performance based on GO terms

368 Even though we obtained good substantial results for the prediction of proteomics values at a sample level, these369 results were an overall representation of all proteomics values, and thus did not inform which protein groups

370 would be better or worse predicted by our model. For this reason, we decided to stratify the predictions based on GO terms, and then evaluate their R<sup>2</sup> metrics when compared to their counterpart observations. The overall met-371 ric for the testing data/sample in this experiment was  $R^2 \sim 0.82$ . What we observed (Figure 7) is that there were 372 373 considerable differences in R<sup>2</sup> depending on the GO term the proteins were associated to. While the 6 best-performing GO terms (inflammatory response, magnesium ion binding, mitochondrial nucleoid, unfolded protein 374 binding, ATPase, and negative regulation of cell growth) had near perfect results ( $R^2 > 0.9$ ), the worst perform-375 ing ones (ligase activity, polysomal ribosome, small ribosomal unit, stress fiber, cell migration, and proteasome 376 complex) showed metrics half the performance shown in the overall results ( $R^2 \le 0.4$ ). 377



#### Best and worse predicted proteins grouped by GO terms (Rsquared)

378

Figure 7: R<sup>2</sup> results categorized in GO terms. The X axis represents the R<sup>2</sup> values, while the bar labels represent the GO terms. Only the 6
best and 6 worst GO terms are depicted.

#### 381 **3.8** IMPUTATION: A POTENTIAL USE OF THE RANDOM FOREST MODEL

As the model showed a promising accuracy for predicting whole replicate samples, we hypothesized that the model could also be used for imputation of missing values for proteins that were at least present in one of the samples of the training data. To showcase a possible example, we trained a random forest model with all the Untreated samples (UNTR) and the corresponding 10 features. The example (Table 4) showed that the proteomics values imputed fitted the range of quantification observed in the quantified values of the same protein, while differing from each other from sample to sample. We also observed that in these samples, values tend to be missing simoultaneously for samples taken at the same time.

- 389 Table 4: Imputation of Proteomics Cardiac samples. Every row is identified with a UniProt ID, and represents a protein quantified in at
- **390** *least one of the untreated samples of the cardiac dataset. Each column represents each Untreated (UNTR) sample from the cardiac*
- 391 dataset. On the column names, the first number represents the hour at which the sample was taken (2h, 8h, etc.), while the second identi-
- **392** *fies the replicate number (002\_1 was the first replicate sample taken after 2 hours). The proteins (rows) are sorted by proportion of*
- 393 missing data in a increasing order. Values with a dark green background were quantified by proteomics. Values with a light green back-

**394** ground were imputed/predicted by the random forest model.

	002_1	002_2	008_1	008_2	008_3	024_1	024_2	024_3	072_1	072_2	072_3	168_3	240_1	240_2	336_1	336_2	336_3
P22695	15.70	16.09	14.58	14.23	13.91	14.15	14.26	13.83	13.83	13.95	13.89	14.90	14.72	14.35	16.75	16.83	16.94
P51553	12.41	11.03	13.87	13.99	13.76	14.70	14.67	14.82	13.55	13.64	13.34	13.26	13.28	13.45	14.37	14.15	14.27
P62910	15.52	14.49	14.78	14.45	14.68	14.75	14.87	15.04	13.76	13.98	13.67	14.61	15.52	13.81	14.57	14.85	14.49
Q15185	14.43	13.89	14.39	14.28	14.16	14.01	13.99	14.14	14.33	14.17	14.04	11.92	13.48	13.70	13.79	13.95	13.99
P45974	14.66	14.51	11.04	11.26	11.43	13.29	13.24	12.74	14.10	13.16	13.73	13.85	13.07	12.82	13.38	13.65	14.04
P16070	12.72	12.89	13.47	13.55	13.79	13.45	13.22	13.21	13.25	13.44	13.31	12.65	12.92	13.17	12.96	13.28	12.97
P54136	10.93	10.78	11.86	11.91	11.76	12.43	12.57	12.18	13.21	13.13	13.26	12.56	13.23	12.06	11.39	9.76	11.81
043681	13.08	13.37	13.37	13.77	13.48	13.40	13.31	13.29	12.63	13.06	12.06	12.04	11.83	12.10	11.84	12.14	12.01
P50440	13.22	12.88	13.02	13.05	13.45	12.35	12.45	13.93	14.03	11.90	13.19	13.17	13.16	13.09	13.19	13.18	13.27
P47897	12.93	12.56	13.03	12.90	12.79	12.62	12.62	12.62	12.62	12.62	12.62	12.62	12.62	12.62	12.07	12.29	12.53
Q14141	14.81	14.44	13.49	14.06	13.70	13.48	14.25	13.89	13.83	13.74	14.10	13.41	14.06	14.49	14.17	13.63	13.73
P40763	12.12	11.82	12.29	11.98	11.54	11.47	11.53	11.50	11.13	11.88	12.03	12.18	11.63	11.65	12.24	11.40	11.37
P04844	12.23	14.51	13.74	13.74	13.78	14.13	13.63	13.77	13.94	13.92	13.71	13.75	13.88	13.89	14.02	13.66	14.42
Q9Y6E2	11.52	11.90	12.53	12.32	12.10	12.37	12.22	12.18	12.57	11.89	11.71	12.21	11.04	10.25	12.02	12.17	12.10
P13796	10.98	11.25	11.55	11.37	11.30	11.32	10.97	11.24	10.97	11.25	10.82	11.26	11.25	11.27	11.95	9.84	9.83
Q96KP4	12.09	10.96	11.82	11.97	11.88	11.82	12.05	12.11	11.94	12.51	11.89	12.23	11.95	12.46	12.41	12.55	12.53
Q96RQ3	13.00	12.79	12.98	12.79	12.59	12.98	12.93	12.86	12.59	12.79	12.79	12.24	12.59	12.80	12.93	12.59	12.86

395

#### 396 4 **DISCUSSION**

We wanted to build a machine learning model that tightened the gap between transcriptomics and proteomics, using the former as a predictor of the latter. The results indicate that a random forest model, by using only 10 features, can predict with good accuracy ( $R^2 = 0.74$ ) proteomics values from samples in similar circumstances to the ones where it has been trained on. However, predicting protein expression by training the model on other proteins was highly inefficient ( $R^2 = 0.15$ ).

Interestingly, 7 out of the 10 features used by the model were related to RNA expression (Fisher's 402 403 Exact Test for Count Data, p-value = 0.0027). Out of these 7, the most important (as expected) was mRNA expression, which is directly linked to translation, and thus, to protein expression. Followed in 404 feature importance came 3 features related to miRNA expression, which is known to inhibit translation 405 to a vast number of coding transcripts. The least important features related to the 3 RNA subtypes 406 referred to circular RNA expression. Circular RNAs have been hypothesized to work as miRNA 407 sponges, and so even though they are involved in post-transcriptional regulation, they have a more 408 indirect effect. It is postulated that most circular RNAs are by-products of faulty splicing<sup>31</sup>, and thus 409

their regulation might just be mainly due to the regulation of their host gene. Even so, their consistentexpression would still allow them to have an impact on post-transcriptional regulation.

Linear density (mass of a protein divided by its length) and the proportions of both Aspartic Acid and 412 Methionine were the most important features for the final random forest model. One hypothesis to 413 414 explain such model behavior was that these three features (and especially linear density) helped to categorize observations protein-wise: an observation with similar values across the three top features 415 could be likely categorized as a similar protein, and thus, also presenting a close expression value. This 416 already made the model highly accurate when trained and tested with similar samples. The other 417 418 features (related to the current transcript expression level) might have helped to succinctly tune the protein expression already observed in similar proteins during the training step. Another hypothesis, 419 only relevant to linear density, was linked to the proteomics technology itself: linear density was 420 directly linked to protein mass, which is used (along with charge) to identify and quantify protein in 421 mass spectrometry; hence, its relevance as a feature. In addition, having linear density as one of the 422 main features underlines the importance of the training data for our model. A random forest model can 423 only predict values learned beforehand, thus we hypothesize that linear density helps the model to find 424 the most similar protein when predicting. Thus, the use of this model should be to predict proteins that 425 are already quantified in some of the samples, limiting the effect of potential false positives, and 426 427 therefore also limiting the potential false biological significances created by false positives due to differences that only occur at the transcriptional level. 428

The observed divergence between the feature importance ranking in RFE and the validation model may be due to how RFE evaluates features while using cross-validation. At the beginning of the process, RFE built 10 different training-testing combinations (based on the 10 folds), and, based on the initial ranking of all the features in each of those combinations, features were removed from least to most important. Each feature was ranked based on the average of all the rankings performed during the feature elimination. In the validation model, instead, the feature importance ranking represented the concrete importance of each variable for that specific model and algorithm.

436 Considering the relatively high accuracy of the random forest model to impute protein expression from 437 a reduced subset of features, we see an application of this proposed strategy to contribute to 438 compensating for the lack of depth of proteomics. Indeed, since proteomics only allows the analysis of 439 a subset of proteins per sample, with usually only a partial overlap between samples (even at the replicate level), our model would be able to predict and fill those values, increasing the strength of thestatistical analysis of such proteins across treatments.

However, as shown in the GO-term-performance results, the metrics are not uniform for all categories 442 of proteins, and this should be taken into consideration when performing analysis with a specific focus 443 444 on a certain protein category. This difference may be the result of three different causes: 1/ the correlation of protein abundance with their coding RNA levels may differ across GO categories, 2/ as 445 different GO categories contained an unequal number of proteins, the size of a GO category was 446 inversely proportional to the R<sup>2</sup> metric (a smaller random set of values has a higher chance of obtaining 447 a high R<sup>2</sup>, and vice versa), 3/ GO categories with stable protein abundances (and mRNA levels) 448 performed better than otherwise. 449

450 An important detail to consider is that drastically different data is generated when utilizing different methods to quantify proteomics intensities: from values that correlate with absolute abundance based 451 on the MS signal of histones (also referred to as the "proteomic ruler" approach<sup>32</sup>), going through 452 intensities inferred based on the ratio of detected peptides (pertaining to each protein) between samples 453 (MaxLFQ<sup>33</sup>), to isobaric proteomics data (TMT/iTRAQ); wherein changes in the peptide intensity from 454 one sample has a ripple effect on the intensities from all the co-isolated samples<sup>34</sup>. In our study, the Hi3 455 label free method<sup>35</sup> was used to quantify protein intensities, hence values from absolute abundance 456 methodologies are expected to perform similarly. Despite that, isobaric proteomics methods should not 457 458 be entirely dismissed, as the range of values predicted by a random forest model is highly dependent on the range of the data the model is trained on. The compositional nature of isobaric proteomics 459 experiments results in signals that are highly batch-dependent. Our predictions would not take the 460 batch structure into account, and as a result, a correction would be required. Thus, the inability of 461 random forest models to extrapolate does make them an appealing option for compositional data, but 462 simoultaneously may be a limiting factor for absolute intensity values. 463

Based on the inefficient accuracy for all models tested in the protein-splitting strategy, we hypothesize that even though we tried to include as much information related to protein expression as possible (transcript expression, transcript properties, protein characteristics, and stability), predicting protein expression anew (without ever training the model with that protein's data) may have required of an even more complete (i.e. RNA binding proteins, long non-coding RNAs, transcript half-life, etc.) or different set of features. For example, a study by Barzine et al<sup>26</sup> showed improved results (R<sup>2</sup> = 0.51) extrapolating proteomics values while only using gene expression data, GO terms, and UniProt keywords. Future research should focus on either including the last two features as features to the
dataset, or improving their deep learning model by including our (or other) post-transcriptional
features.

In conclusion, after developing different machine learning models to predict proteomics values out of transcriptomics ones, we have achieved to build a random forest model that can predict with significant accuracy the protein expression of a new sample. Building a random forest model with the selected features can thus be used to predict the missing data inherent in proteomics studies, independently of the cell's nature. The code used for the pre-processing of data and the model building process is available on Github (<u>https://github.com/jochotecoa/ml\_proteomics</u>)<sup>36</sup>.

#### 480 **5 ACKNOWLEDGEMENTS**

The current research was funded by the European Union Seventh Framework Programme HeCaToS
(FP7/2007-2013) under the [Grant Agreement No. 602156].

#### 483 6 DATA AVAILABILITY

484 Data has been submitted to the BioStudies repository (<u>https://www.ebi.ac.uk/biostudies/</u>) and is avail485 able under the following accession numbers:

- Hepatic data: S-HECA33, S-HECA34, S-HECA47, S-HECA158, S-HECA457, S-HECA460,
   S-HECA463
- Cardiac data: S-HECA1, S-HECA9, S-HECA18, S-HECA139, S-HECA447, S-HECA449, S HECA453
- 490 Supplementary data can be found in the <u>online version</u> of this publication.

#### 491 **7 BIBLIOGRAPHY**

Zhang, Z., Wu, S., Stenoien, D. L. & Pasa-Tolic, L. High-throughput proteomics. *Annu Rev Anal Chem (Palo Alto Calif)* 7, 427-454, doi:10.1146/annurev-anchem-071213-020216 (2014).
 Specht, H. *et al.* Single-cell proteomic and transcriptomic analysis of macrophage heterogeneity
 using SCoPE2. *Genome Biol* 22, 50, doi:10.1186/s13059-021-02267-5 (2021).

- Schoof, E. M. *et al.* Quantitative single-cell proteomics as a tool to characterize cellular hierarchies. *Nat Commun* 12, 3341, doi:10.1038/s41467-021-23667-y (2021).
- 498 4 Cheung, T. K. *et al.* Defining the carrier proteome limit for single-cell proteomics. *Nat Methods*499 18, 76-83, doi:10.1038/s41592-020-01002-5 (2021).
- 500 5 Cagney, G. *et al.* Human tissue profiling with multidimensional protein identification technol-501 ogy. *J Proteome Res* **4**, 1757-1767, doi:10.1021/pr0500354 (2005).
- 502 6 Chen, G. *et al.* Discordant protein and mRNA expression in lung adenocarcinomas. *Molecular*503 & *cellular proteomics : MCP* 1, 304-313, doi:10.1074/mcp.m200008-mcp200 (2002).
- Lemée, J.-M. *et al.* Integration of transcriptome and proteome profiles in glioblastoma: looking
  for the missing link. *BMC molecular biology* 19, 13, doi:10.1186/s12867-018-0115-6 (2018).
- Rogers, S. *et al.* Investigating the correspondence between transcriptomic and proteomic expression profiles using coupled cluster models. *Bioinformatics* 24, 2894-2900, doi:10.1093/
  bioinformatics/btn553 (2008).
- 509 9 Dhingra, V., Gupta, M., Andacht, T. & Fu, Z. F. New frontiers in proteomics research: a per510 spective. *Int J Pharm* 299, 1-18, doi:10.1016/j.ijpharm.2005.04.010 (2005).
- Belle, A., Tanay, A., Bitincka, L., Shamir, R. & O'Shea, E. K. Quantification of protein halflives in the budding yeast proteome. *Proc Natl Acad Sci U S A* 103, 13004-13009,
- 513 doi:10.1073/pnas.0605420103 (2006).
- 514 11 Ambros, V. The functions of animal microRNAs. *Nature* **431**, 350-355,
- 515 doi:10.1038/nature02871 (2004).
- Lim, L. P. *et al.* Microarray analysis shows that some microRNAs downregulate large numbers
  of target mRNAs. *Nature* 433, 769-773, doi:10.1038/nature03315 (2005).
- Lewis, B. P., Shih, I.-h., Jones-Rhoades, M. W., Bartel, D. P. & Burge, C. B. Prediction of
  Mammalian MicroRNA Targets. *Cell* 115, 787-798, doi:10.1016/S0092-8674(03)01018-3
  (2003).
- 521 14 Baek, D. *et al.* The impact of microRNAs on protein output. *Nature* 455, 64-71, doi:10.1038/
  522 nature07242 (2008).
- 523 15 Selbach, M. *et al.* Widespread changes in protein synthesis induced by microRNAs. *Nature*524 455, 58-63, doi:10.1038/nature07228 (2008).
- 52516Zaphiropoulos, P. G. Exon skipping and circular RNA formation in transcripts of the human526cytochrome P-450 2C18 gene in epidermis and of the rat androgen binding protein gene in
- 527 testis. *Mol Cell Biol* **17**, 2985-2993, doi:10.1128/MCB.17.6.2985 (1997).

- 528 17 Chen, L. L. & Yang, L. Regulation of circRNA biogenesis. *RNA Biol* 12, 381-388,
  529 doi:10.1080/15476286.2015.1020271 (2015).
- Jeck, W. R. *et al.* Circular RNAs are abundant, conserved, and associated with ALU repeats.
   *RNA* 19, 141-157, doi:10.1261/rna.035667.112 (2013).
- Rong, D. *et al.* An emerging function of circRNA-miRNAs-mRNA axis in human diseases. *Oncotarget* 8, doi:10.18632/oncotarget.19154 (2017).
- 534 20 Memczak, S. *et al.* Circular RNAs are a large class of animal RNAs with regulatory potency.
  535 *Nature* 495, 333-338, doi:10.1038/nature11928 (2013).
- 536 21 Courel, M. *et al.* GC content shapes mRNA storage and decay in human cells. *Elife* 8, doi:10.7554/eLife.49708 (2019).
- Tyanova, S., Temu, T. & Cox, J. The MaxQuant computational platform for mass spectrometrybased shotgun proteomics. *Nat Protoc* 11, 2301-2319, doi:10.1038/nprot.2016.136 (2016).
- Chen, A. T., Franks, A. & Slavov, N. DART-ID increases single-cell proteome coverage. *PLoS Comput Biol* 15, e1007082, doi:10.1371/journal.pcbi.1007082 (2019).
- Kalxdorf, M., Muller, T., Stegle, O. & Krijgsveld, J. IceR improves proteome coverage and data
  completeness in global and single-cell proteomics. *Nat Commun* 12, 4787, doi:10.1038/s41467021-25077-6 (2021).
- Lim, M. Y., Paulo, J. A. & Gygi, S. P. Evaluating False Transfer Rates from the Match-between-Runs Algorithm with a Two-Proteome Model. *J Proteome Res* 18, 4020-4026,
  doi:10.1021/acs.jproteome.9b00492 (2019).
- 548 26 Barzine, M. P. *et al.* Using Deep Learning to Extrapolate Protein Expression Measurements.
  549 *Proteomics* 20, e2000009, doi:10.1002/pmic.202000009 (2020).
- 550 27 Mitchell, T. M. Machine Learning. (McGraw-Hill, 1997).
- 551 28 CRAN Package caret, <<u>https://cran.r-project.org/package=caret</u>>(
- 552 29 Kuhn, M. 3 Pre-Processing | The caret Package, <<u>https://topepo.github.io/caret/pre-process-</u>
   553 ing.html#zero--and-near-zero-variance-predictors> (2019).
- Butcher, B. & Smith, B. J. Feature Engineering and Selection: A Practical Approach for Predictive Models. *The American Statistician* 74, 308-309, doi:10.1080/00031305.2020.1790217
  (2020).
- Barrett, S. P. & Salzman, J. Circular RNAs: Analysis, expression and potential functions. *Development (Cambridge)* 143, 1838-1847, doi:10.1242/dev.128074 (2016).

- Wisniewski, J. R., Hein, M. Y., Cox, J. & Mann, M. A "proteomic ruler" for protein copy number and concentration estimation without spike-in standards. *Mol Cell Proteomics* 13, 34973506, doi:10.1074/mcp.M113.037309 (2014).
- 562 33 Cox, J. *et al.* Accurate proteome-wide label-free quantification by delayed normalization and
  563 maximal peptide ratio extraction, termed MaxLFQ. *Mol Cell Proteomics* 13, 2513-2526,
  564 doi:10.1074/mcp.M113.031591 (2014).
- O'Brien, J. J. *et al.* Compositional Proteomics: Effects of Spatial Constraints on Protein Quantification Utilizing Isobaric Tags. *J Proteome Res* 17, 590-599, doi:10.1021/
  acs.jproteome.7b00699 (2018).
- Silva, J. C., Gorenstein, M. V., Li, G. Z., Vissers, J. P. & Geromanos, S. J. Absolute quantification of proteins by LCMSE: a virtue of parallel MS acquisition. *Mol Cell Proteomics* 5, 144-
- 570 156, doi:10.1074/mcp.M500230-MCP200 (2006).
- 571 36 *GitHub jochotecoa/ml\_proteomics*, <<u>https://github.com/jochotecoa/ml\_proteomics</u>>(

### Chapter 5:

# AutoRel: Machine Learning model for selecting differentially expressed genes

#### 1 AutoRel: Machine Learning model for selecting differentially expressed

#### 2 genes.

3 Juan Ochoteco Asensio<sup>1</sup>, Jelmer Faber<sup>1</sup>, Twan van den Beucken<sup>1</sup>, Marcha Verheijen<sup>1</sup>, Florian Caiment<sup>1</sup>

4 <sup>1</sup>Department of Toxicogenomics, School of Oncology and Developmental Biology (GROW), Maastricht

5 University, Maastricht, The Netherlands

6 (Manuscript under submission)

#### 7 Abstract

8 In next-generation transcriptomics, differential expression analysis has established itself as one of the 9 main strategies to evaluate the effects of two biological conditions on the gene expression of different 10 biological samples. However, the current statistical workflows provide very lenient standardized filters for the selection of differentially expressed genes, leading to an important number of false positives. 11 12 Authors usually include their arbitrary thresholds, leading to statistically significant genes with expression profiles that, if individually analyzed, would not be considered biological relevant to investigate further. 13 In this study, we focused on the development of a machine learning model, called AutoRel, which 14 encompasses not only the most common statistical evaluations but also all the intricacies that characterize 15 16 biologically relevant changes based on manual observations. AutoRel classifies each assessed gene into "relevant", "irrelevant" or "dubious". We evaluated the model using both simulated datasets and 17 of 18 biological interpretation the selected genes. GitHub Repository at https://github.com/jochotecoa/AutoRel.git. 19

#### 20 Introduction

21 Next-generation sequencing (or NGS) of RNA (RNA-seq) has brought the capacity to sequence and 22 quantify virtually all RNA molecules within a given sample, being a great advancement from single or 23 few molecule identification methods such as Northern Blot or even a limited set of transcripts in microarray technologies. RNA-seq is especially relevant for contrasting different biological conditions, 24 where the effects at the molecular level can be analyzed in different conditions, such as cells during 25 tumorigenesis or upon exposure to toxic compounds. For the evaluation of these differences, multiple 26 27 tools have been developed since the dawn of transcriptomics, however, their usage can be mishandled 28 without a proper understanding of the mechanisms behind it. This problem is especially relevant for the 29 selection of differentially expressed genes (DEGs) based solely on statistical results. Using standard 30 statistical criteria, the selected DEGs between two biological conditions frequently display low read 31 abundance, inconsistent expression among the replicates, or the presence of outliers. Establishing a 32 reference and guideline to avoid such pitfalls is critical. We recently collaborated in the R-ODAF 33 framework<sup>1</sup>, which aimed to standardize omics data analyses for regulatory applications, specifically in 34 the toxicogenomics discipline<sup>2</sup>. Selecting a list of genes affected by a compound exposure for regulatory 35 decision-making requires a high level of stringency, as false interpretation could have drastic 36 consequences.

The development of R-ODAF originated primarily from the necessity to implement novel filtering criteria 37 38 based on the expression profiles of the selected DEGs. The selection process was performed by analyzing the significantly different genes and evaluating their biological relevance by considering the 39 quantification values across all samples and groups. Metrics such as fold change and p-values<sup>3-5</sup> can be 40 misleading (and lead to false positives) if not evaluated in their proper context due to common 41 42 characteristics such as the presence of outliers or a low read quantification. Systems to prevent the 43 inclusion of statistically significant genes that presented such properties (thus probably not biologically relevant) were set up in the R-ODAF pipeline. However, despite being designed to decrease the frequency 44 45 of irrelevant genes from the DEG set, the R-ODAF is still based on an adjusted p-value threshold (FDR < 46 0.01, FDR: False Discovery Rate adjustment), which by definition creates an arbitrary cut-off. Moreover, 47 applying a stringent pipeline essentially generates more false negatives, especially in the typical context 48 of experiments performed in triplicates. A filtering-based method only includes the most frequent 49 characteristics of genes that would be manually discarded, thus frequently requiring a manual selection afterward. 50

An alternative to the threshold-based selection of DEGs would be to apply a machine learning method 51 trained on gene expression profiles considered "relevant" by experts (relevant defined as "worthy of 52 53 further investigation" after an analysis of the quantification values and their distribution in each 54 condition). Advancements in artificial intelligence allow practitioners to automatize labeling processes 55 (classification) or predictions of numerical values (regression) through the use of enough data<sup>6-9</sup>. In 56 classification tasks, data is required for the nascent model to differentiate between labels via their different features or characteristics. These methodologies can learn the synchronicities of the labels across 57 58 numerous examples without the need for explicitly and manually specifying such differences. In addition, 59 the selection of the variables will rely on how informative they are for the classification of each label, and 60 thus by analyzing the predictive model one can also derive which features and to which degree each of 61 them is important.

62 In this study, we developed and tested machine learning models to automate the biologist-supervised 63 DEG selection process. The best performing model, called AutoRel, informed us of the importance of 64 each of our previously designed norms, while simultaneously providing a novel system to select biologically relevant genes without the expertise needed to evaluate them or based purely on a few 65 arbitrary thresholds. Due to the subjective nature of the notion of what is "biologically relevant", AutoRel 66 outputs not only a "relevant" and "irrelevant" list of genes, but also a class named "dubious" for the gene 67 profiles subject to debate during the training. Ultimately, the AutoRel selected list of genes (named 68 "relevant" list) is compared to conventional methods. 69

#### 70 Results

71 The transcriptomics dataset consisted of 60683 assessed genes in 42 samples: 21 APAP-treated samples 72 and 21 control samples. From this original dataset, we built three datasets based on different comparisons 73 based on the number of replicates: 21 versus 21 (all time points), 9 versus 9 (2h, 8h, 24h), and 3 versus 3 (2h). For this analysis, we considered all samples in each group (treatment or control) as replicates. For 74 75 each of these datasets, the quantification barplots of individual genes were used to manually label hundreds of genes into two categories: 'irrelevant' or 'relevant', based on whether the gene expression 76 profile was considered worthy of further investigation due to biological differences. Observations that did 77 not suit either of those classes were classified as "dubious". An example of each class is displayed in 78 Figure 1. After labeling genes for all three datasets, we assembled all labeled genes into a new dataset, 79 80 identified as 'MixR'. Next, we derived features related to the expression levels that might help train a machine learning model to predict the categories of new genes based on the manually annotated data. In 81 82 total, 160 features were obtained for all datasets (Supplementary Methods). Before the model training, each dataset was pre-processed to remove uninformative or redundant features. We trained and tested all 83 84 datasets with 11 machine learning algorithms, whose results can be found below.



86

87 *Figure 1: Quantified expression of genes from different classes. A: Example of a gene with a relevant biological difference.*88 *Example of a gene with an irrelevant biological difference.*C: *Example of a gene with a dubious biological difference.*

#### 89 Selecting the best-performing models in the 3R, 9R, 21R, and MixR datasets

To decide which was the best performing model, we did not utilize accuracy, which is easily biased towards the class with the highest number of observations. Instead, we averaged the sensitivity, specificity, PPV, and NPV for each class, which we named "equilibrated accuracy". We will use the MixR dataset as an example to show the process of the selection of the best model (in each dataset), but all datasets were pre-processed and trained in the same manner, including the same algorithms. The results for the models built in the other three datasets (3R, 9R, and 21R) can be found in the **Supplementary Results**.

97 The 'treebag' algorithm showed the best performance out of 10 tested algorithms in total (Figure 2A),98 closely followed by 'rf' (both algorithms also performed the best in the 3R and 21R datasets). Even with
99 the biggest training dataset, four algorithms ('rpart2', 'lssvmRadial', 'pam', and 'CSimca') were not able
100 to predict 'dubious'-class observations.



Figure 2: Selection and evaluation of the best-performing model using MixR as the training dataset. A: Equilibrated
 accuracies for the models built with 11 machine learning algorithms using the MixR dataset. The models were sorted by the

average of equilibrated accuracies in descending order. B: Variable importance for the 10 most important features sorted in
 decreasing order for the 'treebag' model trained with the combined dataset.

107 For the variable importance of the model (Figure 2B), 'pvalue' and 'padj' became again (as in the 108 'treebag 3R' model) the most important features. Similarly, variables featured as the most important for 109 the best-performing models in differing training datasets were also found in the top 10 features for the 110 combined dataset: 'fdrlowerthan0.01', 'threequartilediff rule', 'q3belmin', 'rule cpm 0.75 above 1', 'lfcSE', and 'baseMean'. Curiously, two additional features were also part of the top 10 that were not 111 112 observed before: the minimum expression value for the treatment group ('quantile 0. APA The'), and 113 the proportion of non-expressed samples in the treatment group 114 ('Proportion nonexpressed samples APA The').

115 Testing the best-performing models with all the testing datasets

Even though the four best-performing models were trained and tested with their characteristic datasets, we hypothesized that these models might present different predicting capabilities for datasets not similar to the ones they were trained for. To test the hypothesis, we tested each of the best-performing models with each of the corresponding test datasets and calculated the Cohen's Kappa Coefficient.

120 Table 1: Kappa score for all best-performing models. Each row represented one of the best-performing models, labeled by the

121 name of the algorithm and its training dataset. Each column, in turn, represented the testing dataset utilized to obtain the Kappa

122 score.

		Testing datasets				
		3R	9R	21R	MixR	Average score
Best model in each dataset	treebag (3R)	0.970	0.894	0.485	0.849	0.800
	sparseLDA (9R)	0.614	0.918	0.316	0.609	0.614
	rf (21R)	0.538	0.919	0.915	0.658	0.757
	treebag (MixR)	0.974	0.955	0.912	0.946	0.947
	Average score	0.774	0.922	0.657	0.766	

123

As seen in the results (Table 1), substantial differences were shown when training and testing with different datasets. The best performing model for predicting any contrast independently of the number of replicates was the treebag (MixR) model, which we referred to as 'AutoRel'. Remarkably, the latter model presented a higher Kappa Score when testing the 3R and 9R datasets than the models trained with those datasets. In turn, the testing dataset with the highest average Kappa Score (the testing dataset best predicted by all models) was '9R'. The quality of the dataset did not appear to influence directly both training and testing steps, as the average training and average testing scores presented a Pearsoncorrelation of -0.521.

## 132 Model comparison

133 To further investigate how the 'treebag\_MixR' model classified genes as 'relevant', we decided to

134 compare the predictions of this model with several methods that could be used to select genes of interest

135 in a transcriptomics study: FDR (p.adjusted value) < 0.05, FDR < 0.05 in combination with a log2 fold

136 change > 1.5, FDR < 0.01, FDR < 0.01 in combination with a log2 fold change > 1.5, and R-ODAF. The

137 statistical values were obtained using DESeq2. The input used to compare the different methods were the

- 138 full 3R, 9R, and 21R datasets.
- 139 For the first comparison, we made use of the 3R dataset for which results showed the highest difference
- 140 between 'relevant' and significant genes (Figure 3A). In addition to the 'relevant' set size being 1.8 times
- 141 bigger than the significant set size, only 39.4% of the 'relevant' genes were significant (FDR < 0.05), and
- of these, 69.7% presented an L2FC > 1.5. Although only 28.1% of the 'relevant' genes were detected by
- 143 R-ODAF, 100% of the R-ODAF genes were 'relevant'.





147 *Figure 3: UpSet plot of the intersections between the evaluated methods using the diverse datasets. A:* 3*R dataset. B:* 9*R*148 *dataset. C:* 21*R dataset. Set Size referred to the total number of genes positively labeled or predicted by any of the methods*

For the 9R dataset (Figure 3B), the 'relevant' class was no more the smallest group by Set Size. The impact of significance for the 'relevant' genes was decreased: 80.7% of them were significant (FDR < 0.05). The 9R results, when compared with the 21R results, showed that the 'relevant' class was less reliant on p-values the fewer replicates were used for the contrast. The decreased reliance on statistical values also decreased the proportion of R-ODAF genes in the 'relevant' set (60.3% of 'relevant' genes detected by R-ODAF), although simultaneously almost all genes in the R-ODAF set were also 'relevant' (94.8% of the R-ODAF set of genes were also 'relevant'). 72.6% of the significant (p.adj < 0.05)</li>
'relevant' genes presented an L2FC > 1.5. For the significant genes (FDR < 0.05), 68.1% of them were</li>
also 'relevant'.

Lastly, we performed the same analysis using the 21R dataset as input (Figure 3C). For this number of 158 replicates, the genes labeled as 'relevant' appeared to be a more restricted subset of genes compared to 159 any other significant set of genes (as seen in the Set Size). The 'relevant' genes were mostly significant 160 (99.7% of them with an FDR < 0.05), and slightly biased to have a high L2FC (64.3% of the significant 161 162 'relevant' genes presented an L2FC > 1.5). The 'relevant' genes also presented a high overlap with R-ODAF (87.0% of them detected by R-ODAF), but R-ODAF identified in total more genes (76.3% of the 163 latter were 'relevant'). For the significant genes (FDR < 0.05), 46.15% of them were also 'relevant', 164 showing that for this number of replicates, significance was a necessary condition for 'relevant' genes, 165 but not sufficient. 166

167 Analysis of exclusively 'relevant' genes

The results of the model in the 3R dataset showed that the majority of the 'relevant' genes were not statistically significant using an FDR threshold of 0.05. We thus further analyzed these gene expressions to visualize how a gene expression change could be simultaneously relevant and not significant. For this reason, we extracted the normalized reads of the first three genes from the "exclusively relevant" set (M6PR, DBNDD1, and RBM5), and generated a barplot and boxplot for each of them (Figure 4A, Figure 4B, and Figure 4C).











177 Figure 4: Barplot and boxplot of the normalized count expression of exclusively relevant genes. A: AM6PR

**178** (ENSG00000003056, FDR = 6.25E-01. L2FC = 0.38). **B**: DBNDD1 (ENSG00000003249, FDR = 4.68E-01, L2FC = -0.67). **C**:

- **179** *RBM5* (*ENSG00000003756*, *FDR* = 6.52E-02, L2FC = 1.73). In the barplot (left half), the gray columns represented the gene
- 180 expression in the control samples, while the pink columns represented the gene expression in the treated samples. The boxplot
- 181 (right half) represented the same values while summarizing the values for each condition in a boxplot.

182 M6PR gene expression (Figure 4A) showed that, in a triplicate setting, a single sample 183 (ConDMSO\_024\_3) could affect substantially the statistical values, while still being potentially relevant 184 for further investigation. For DBNDD1 (Figure 4B), a pattern could be seen for both groups triplicate-185 wise, while expressing in different scales, the low number of replicates makes these two conditions not 186 significantly different.

187 Lastly, the expression levels of RBM5 (Figure 4C) showed a clear difference between both conditions,

including a high log2 fold change (1.73). Even so, the use of a standard 5% FDR threshold identified the

189 expression levels in the treatment group as not significantly different from the control condition.

# 190 Effect of relevant genes on biological interpretation

191 The AutoRel model detected 3839, 2293, and 2686 relevant genes in the 3R, 9R, and 21R datasets 192 respectively. To evaluate the effect on the biological interpretation of the genes labeled by the model, we 193 decided to use the relevant gene set derived from the 21R dataset. As a contrast gene set, we selected the 194 significant (p. adj. < 0.05) genes from the same dataset, with a set size of 5811 genes.

195 The analysis showed that the GO terms in the three different ontologies (Biological Process, Biological Function, and Biological Component) were simplified in the relevant genes in comparison to the 196 significant ones (Figure 5). In the Biological Process ontology, GO terms exclusive to the significant gene 197 set were very variable and unrelated to hepatocytes exposed to APAP. Some examples were 198 spermatogenesis (p-value = 4.51E-5), male gamete generation (p-value = 6.04E-5), regulation of neuron 199 differentiation (p-value = 1.77E-4), and oocyte development (p-value = 2.54E-4). Out of the exclusive 200 201 GO terms in the relevant gene set, the most significant were response to stimulus (p-value = 9.51E-5) and 202 biological process (p-value = 5.02E-4). The other three GO terms were related to coagulation (negative regulation of coagulation, negative regulation of blood coagulation, and negative regulation of 203 hemostasis), which has been hypothesized to be enhanced via p62 in APAP-induced liver injury<sup>29</sup>. 204





Figure 5: UpSet plot of the different enriched GO terms between significant (FDR 0.05) and relevant genes. Set Size shows the
 number of enriched GO terms for a specific method (FDR < 0.05 or RELEVANT) in a concrete ontology (Process, Function, or</li>
 Component). The Intersection Size shows how many enriched GO terms were exclusive for that method (single dot) or shared
 with the equivalent method for the same ontology (two dots connected by a line).

For the Biological Function ontology, the only shared GO term between both methods was anion binding. The significant gene set showed a diverse range of terms related to the catalytic activity (phosphotransferase activity, kinase activity, sterol 14-demethylase activity) and molecule binding (peptide hormone receptor binding, CD4 receptor binding, drug binding), while the relevant set was limited to cytoskeletal protein binding. Lastly, the Biological Component ontology showed, despite the 215 lower number of GO terms in the relevant set, the highest overlap of the three ontology classes.216 Interestingly, the significant set presented 'neuron part' as an enriched GO term, but not the relevant set.

#### 217 Dataset simulation

To evaluate the predictive value of the model, we generated simulated datasets with different numbers of replicates: 50 (50R), 21 (21R), 9 (9R), and 3 (3R). For every number of replicates setting, 100 iterations were performed to avoid skewed results. As predictors, we selected significant genes (FDR < 0.05), relevant genes, and a combination of relevant and dubious genes. We compared the predicted results with the original population difference and evaluated it using several parameters: Accuracy, Kappa, Sensitivity (or Recall), Specificity, PPV (or Precision), NPV, and Balanced Accuracy.

In the 50R setting (Supplementary Figure 7), relevant genes presented a similar accuracy to significant 224 225 ones. Kappa and sensitivity showed the biggest decrease, the latter of which (in combination with a slightly increased specificity) led to a decrease in the balanced accuracy. PPVs were instead increased, in 226 combination with a slight decrease in NPVs. Including the dubious genes (Supplementary Figure 8) 227 228 decreased the differences with the significant genes among all metrics. In the 21R setting (Supplementary Figures 9 and 10), highly similar results to the 50R setting were found (for both relevant and 229 230 dubious+relevant groups), which suggested that similar results might be interpolated in the range between 231 21 and 50 replicates.

In the 9R setting (Supplementary Figure 11), the differences across metrics generally decreased. Accuracy continued to be similar between relevant and significant genes. Kappa was also still decreased for relevant genes, but the difference was smaller. Sensitivity was also slightly decreased, but in combination with an almost unchanged specificity, led to an almost equally balanced accuracy between both groups. The PPV was also mostly decreased, while the NPV was overwhelmingly unaffected. Including dubious genes increased sensitivity but at the cost of PPV (Supplementary Figure 12).

In the 3R setting (Supplementary Figure 13), we found the biggest differences across groups. Accuracy was slightly decreased, and kappa, although also decreased, showed a much bigger variance depending on the dataset/iteration. Sensitivity was strongly increased, while specificity was slightly decreased, thus balanced accuracy was only marginally decreased. The increase in sensitivity was at the cost of PPV, which strongly decreased. Adding dubious genes mostly decreased most metrics: accuracy, kappa, specificity, PPV, and balanced accuracy (Supplementary Figure 14).

# 244 Discussion

As an alternative to the typical selection of differentially expressed genes using a statistical threshold, we chose to build a model that would select genes with biologically relevant differences across treatments trained on the criteria selection we generally use for further validation. For this, the model was trained with our decision-making based on the characteristics observed in the expression of hundreds of genes. Training a 'treebag' algorithm with observations from diverse datasets resulted in the optimal model to classify 'relevant', 'dubious', and 'irrelevant' genes.

Among the best-performing algorithms trained on distinct training datasets (3R, 9R, 21R, and MixR), the model that showed the best overall Kappa accuracy was the one trained with the MixR dataset using the 'treebag' algorithm (which we named 'AutoRel'). This was consistent with our suppositions, as the training dataset contained the highest number of observations in addition to the most variance among them.

256 Analyzing the variable importance of the AutoRel model allowed us to evaluate what factors were critical for selecting genes as potentially relevant in a biological setting. As expected, P-values and FDR values 257 were vital predictors for the model, as extreme changes statistically tended to refute the null hypothesis. 258 Secondly, the 'threequartilerule' feature (Supplementary Methods), which refers to gene changes where 259 260 there are three quartile differences across groups (such as the minimum of one group presenting a value higher than the 3<sup>rd</sup> quartile of the other group), also appeared highly important in AutoRel. This might be 261 262 explained by the fact that these genes, while being substantially different across groups, were more 263 informative than more extreme rules (such as 'fourquartilediff'), as the former difference occurred more 264 frequently than the latter.

265 Even though statistical values were critical for the AutoRel model, the overlap between relevant and significant genes decreased drastically as the number of replicates involved decreased (, Figure 3). This 266 267 was especially the case with the dataset with the lowest number of replicates (3R, Figure 3A). This might 268 have been due to the higher impact of every sample for the whole contrast. As seen in Figure 4, the 269 expression of a single replicate (the third control triplicate) resulted in a high p-value, independently of a 270 clear difference between most samples in both groups. This was also an example of how using L2FC as a 271 statistical threshold might have discarded this gene: the average of the control group was highly affected by a single outlier, decreasing the L2FC value. In addition, M6PR, with a relatively high gene expression, 272 was modestly increased in relative terms to the control average, even though (discarding the outlier), there 273 274 was at least a 500 normalized count difference between both groups. Thus, we observed that AutoRel 275 selected relevant genes based on the differences between the majority of the replicates from both groups,

presenting resilience to the presence of potential outliers, in contrast to the use of statistical p-values orfold changes.

278 Additionally, the effects of manually selecting biologically relevant genes through the developed model were also perceivable during the biological interpretation phase. When evaluating the differences between 279 significant (FDR < 0.05) and relevant genes, we observed that in all the three main GO categories 280 281 (Process, Function, and Component) there was a reduction in the number of enriched GO terms in the relevant set. The shrinkage did not appear to be stochastic, as some GO terms completely unrelated to the 282 283 hepatic toxicity of APAP disappeared (spermatogenesis, oocyte development), while other terms, such as 284 negative regulation of coagulation, were exclusive to the relevant set. Previous research has found an inducement of coagulation in APAP-induced liver injury<sup>29</sup>, thus compensatory mechanisms might be in 285 286 effect in therapeutic doses.

A potential improvement of the model could include previously selected differentially expressed genes, especially if those were later validated (or not) *in vitro*. In addition, providing an online tool for quick assessments of random quantification barplots to increase the size of the training dataset and help diversify the selection criteria outside of our department would be beneficial. This would of course require strict criteria to recruit external experts in the training, to avoid the addition of noisy data to the dataset.

293 In conclusion, we generated a machine learning model named AutoRel that can artificially simulate the manual selection often required for further research on a select number of differentially expressed genes 294 295 based on our understanding of the biology and sequencing processes. AutoRel was more stringent than 296 standard FDR thresholds in experiments with a high number of replicates, decreasing the number of false positives. For a low number of replicates, the effects were inverted, where the number of relevant genes 297 298 was superior to the number of significant ones, decreasing the number of false negatives in a low informative experiment. The model (and the code used to generate it) is publicly available on GitHub: 299 https://github.com/jochotecoa/AutoRel.git. 300

# 301 Bibliography

- Verheijen, M. C. *et al.* R-ODAF: Omics data analysis framework for regulatory application.
   *Regul Toxicol Pharmacol* 131, 105143, doi:10.1016/j.yrtph.2022.105143 (2022).
- Lovett, R. A. Toxicogenomics. Toxicologists brace for genomics revolution. *Science* 289, 536537, doi:10.1126/science.289.5479.536 (2000).
- Wasserstein, R. L., Schirm, A. L. & Lazar, N. A. Moving to a World Beyond "p < 0.05". *The American Statistician* 73, 1-19, doi:10.1080/00031305.2019.1583913 (2019).

308 4 Hurlbert, S. H., Levine, R. A. & Utts, J. Coup de Grâce for a Tough Old Bull: "Statistically 309 Significant" Expires. The American Statistician 73, 352-357, 310 doi:10.1080/00031305.2018.1543616 (2019). 5 McShane, B. B., Gal, D., Gelman, A., Robert, C. & Tackett, J. L. Abandon Statistical 311 Significance. The American Statistician 73, 235-245, doi:10.1080/00031305.2018.1527253 312 (2019). 313 Ochoteco Asensio, J., Verheijen, M. & Caiment, F. Predicting missing proteomics values using 314 6 machine learning: Filling the gap using transcriptomics and other biological features. Comput 315 Struct Biotechnol J 20, 2057-2069, doi:10.1016/j.csbj.2022.04.017 (2022). 316 7 Weis, C. et al. Direct antimicrobial resistance prediction from clinical MALDI-TOF mass spectra 317 using machine learning. Nat Med 28, 164-174, doi:10.1038/s41591-021-01619-9 (2022). 318 319 8 Zhang, J. et al. Rapid Antibiotic Resistance Serial Prediction in Staphylococcus aureus Based on Large-Scale MALDI-TOF Data by Applying XGBoost in Multi-Label Learning. Front Microbiol 320 13, 853775, doi:10.3389/fmicb.2022.853775 (2022). 321 322 9 Jabal, M. S. et al. Interpretable Machine Learning Modeling for Ischemic Stroke Outcome Prediction. Front Neurol 13, 884693, doi:10.3389/fneur.2022.884693 (2022). 323 324 10 Kuepfer, L. et al. Applied Concepts in PBPK Modeling: How to Build a PBPK/PD Model. CPT 325 Pharmacometrics Syst Pharmacol 5, 516-531, doi:10.1002/psp4.12134 (2016). 326 11 Team, R. C. R: A Language and Environment for Statistical Computing, <a href="https://www.R-team.action.org">https://www.R-team.action.org</a> 327 project.org/> (2021). 328 12 Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for 329 RNA-seq data with DESeq2. Genome Biol 15, 550, doi:10.1186/s13059-014-0550-8 (2014). 330 13 Max Kuhn [aut, c., Jed Wing [ctb], Steve Weston [ctb], Andre Williams [ctb], Chris Keefer [ctb], Allan Engelhardt [ctb], Tony Cooper [ctb], Zachary Mayer [ctb], Brenton Kenkel [ctb], R Core 331 Team [ctb], Michael Benesty [ctb], Reynald Lescarbeau [ctb], Andrew Ziem [ctb], Luca Scrucca 332 [ctb], Yuan Tang [ctb], Can Candan [ctb], Tyler Hunt [ctb]. caret: Classification and Regression 333 334 *Training*, <<u>https://github.com/topepo/caret/</u>>(2019). 14 Todorov, V. rrcov: Scalable Robust Estimators with High Breakdown Point, 335 336 <<u>https://github.com/valentint/rrcov</u>>(2021). 15 Klaus Schliep [aut, c., Klaus Hechenbichler [aut], Antoine Lizee [ctb]. kknn: Weighted k-Nearest 337 *Neighbors*, <<u>https://github.com/KlausVigo/kknn</u>> (2016). 338 16 Zeileis, A. K. a. A. S. a. K. H. a. A. kernlab -- An {S4} Package for Kernel Methods in {R}. 339 Journal of Statistical Software 11, 1--20, doi:10.18637/jss.v011.i09 (2004). 340 17 Ripley, W. N. V. a. B. D. Modern Applied Statistics with S. Fourth edn, (Springer, 2002). 341

- 342 18 Majka, M. naivebayes: High Performance Implementation of the Naive Bayes Algorithm in R,
   343 <a href="https://CRAN.R-project.org/package=naivebayes">https://CRAN.R-project.org/package=naivebayes</a>> (2019).
- Wurm, M. J., Rathouz, P. J. & Hanlon, B. M. Regularized Ordinal Regression and the ordinalNet
  R Package. *J Stat Softw* 99, doi:10.18637/jss.v099.i06 (2021).
- T. Hastie, R. T., Balasubramanian Narasimhan, Gil Chu. *pamr: Pam: Prediction Analysis for Microarrays*, <<u>https://cran.r-project.org/package=pamr</u>> (2019).
- 348 21 Wiener, A. L. a. M. Classification and Regression by randomForest. *R News* 2, 18-22 (2002).
- 349 22 Terry Therneau [aut], B. A. a., cre], Brian Ripley [trl] (producer of the initial R port, maintainer
- **350** 1999-2017). *rpart: Recursive Partitioning and Regression Trees*,
- 351 <<u>https://github.com/bethatkinson/rpart</u>> (2018).
- Line Clemmensen, c. b. M. K. sparseLDA: Sparse Discriminant Analysis,
  <a href="http://www.imm.dtu.dk/~lhc">http://www.imm.dtu.dk/~lhc</a> (2016).
- Andrea Peters [aut], T. H. a., cre], Brian D. Ripley [ctb], Terry Therneau [ctb], Beth Atkinson
  [ctb]. ipred: Improved Predictors. (2021).
- 356 25 Kuhn, M. 8 Models Clustered by Tag Similarity, <<u>https://topepo.github.io/caret/models-clustered-</u>
   357 <u>by-tag-similarity.html</u>> (2014).
- Eden, E., Navon, R., Steinfeld, I., Lipson, D. & Yakhini, Z. GOrilla: a tool for discovery and
  visualization of enriched GO terms in ranked gene lists. *BMC Bioinformatics* 10, 48,
- doi:10.1186/1471-2105-10-48 (2009).
- 361 27 Assefa, A. T., Vandesompele, J. & Thas, O. SPsimSeq: semi-parametric simulation of bulk and
  362 single-cell RNA-sequencing data. *Bioinformatics* 36, 3276-3278,
- 363 doi:10.1093/bioinformatics/btaa105 (2020).
- Zhang, W. *et al.* Comparison of RNA-seq and microarray-based models for clinical endpoint
   prediction. *Genome Biol* 16, 133, doi:10.1186/s13059-015-0694-1 (2015).
- 29 Qian, H. *et al.* Dual roles of p62/SQSTM1 in the injury and recovery phases of acetaminophen-
- induced liver injury in mice. *Acta Pharm Sin B* 11, 3791-3805, doi:10.1016/j.apsb.2021.11.010
  (2021).

# 369 Methods

#### 370 Dataset composition

371 The original dataset used for this manuscript originated from a 3D liver cell culture made of spheroids by

- 372 InSphero Inc. 42 samples were taken, 21 per group: control DMSO and acetaminophen (APAP) exposure.
- 373 Those 21 samples were the accumulation of samples in triplicate at 7 different time points: 2, 8, 24, 72,
- 374 168, 240, and 336 hours. The applied dose was selected based on a PBPK profile, whose design is
- 375 described in Kuepfer et al<sup>10</sup>. Total RNA was extracted and ribo-depleted. The libraries were sequenced on
- an Illumina Hiseq 2000, where the average sequencing depth was 27.8 million reads. The nomenclature
- 377 of the samples was as follows: "Exposure"\_" Timepoint"\_" Replicate", where exposure could be either
- 378 the control group ("ConDMSO") or the treatment ("APAP"), followed by any of the seven time points
- 379 (002 336), and ending with any of the triplicates (1 3).

380 This dataset was used as input to the DESeq2 pipeline in three different configurations: 3 samples

- 381 (triplicates from a single time point: 24 hours) from each group contrasted against each other, 9 samples
- 382 (triplicates from 3-time points: 2, 8, and 24 hours) from each group, and 21 samples (triplicates from 7-
- time points: 2, 8, 24, 72, 168, 240, and 336 hours) from each group. All programming scripts that pertain
- to the analyses of this manuscript were written in the R language (major version 3, minor version 6.3)<sup>11</sup>.
- For all 3 comparisons, we used the R package 'DESeq2' (version 1.26.0)<sup>12</sup> to obtain both the quantitative
- 386 and statistical features that would be used as input for the model. The quantitative features were obtained
- 387 from the normalized counts. The normalization was performed using DESeq2's median of ratios. The
- 388 features used for our models are listed and described in the 'Supplementary Methods' document.

# 389 Model training

390 The targets to be predicted were annotated manually for the training dataset. The three labels assigned 391 were 'irrelevant', 'dubious', and 'relevant'. The 'irrelevant' label was assigned when a gene was not 392 estimated to present a biological effect, as a consequence of (but not limited to) both groups being almost identical, their distributions being so variable that no clear change can be identified, or not enough reads 393 394 are sequenced for most samples for that gene. 'Relevant' genes were those where a substantial change could be detected for most samples in both groups, so that a researcher would consider them for further 395 396 research for their biological effect, such as the distributions of both groups spanning a different range of high quantification levels. Observations that did not classify for either of the previously described classes 397 398 were labeled as 'Dubious'. To evaluate each gene in each comparison, we plotted the normalized counts for both groups using both barplots and boxplots. For each comparison, we labeled a random set of 399 observations to include a diverse set of distribution changes for the model to learn from. Thus, after each 400 401 comparison, we obtained a distinct dataset, which we named in an increasing number of replicates: datasets 3R, 9R, and 21R. 402

## 403 Pre-processing and modeling

After forming the complete dataset, it led to a total of 160 features. Aside from the 3 datasets derived 404 from the 3 different comparisons (3, 9, and 21 replicates), we generated a new dataset that included the 405 406 observations of all 3 datasets (dataset named MixR). For all four datasets, we performed several pre-407 processing steps to filter out redundant or uninformative features that might either worsen or slow down 408 the model training process. For the first 2 steps, we used an under-sampled version of the dataset. Under-409 sampling reduces the size of all groups to the size of the smallest one. In the first step, we used a (near) zero-variance filter using the 'nearZeroVar' function from the caret<sup>13</sup> library. Without under-sampling, a 410 411 feature with a common value that was predominant in the most common labels would have presented a reduced variance, and thus would have been more prone to be discarded due to over-representation of a 412 413 target label, and not due to low variance. The second step involved the removal of features highly correlated to each other (correlation > 0.99). The last step of pre-processing involved the exclusion of 414 415 features with a linear dependency on other features, using caret's 'findLinearCombos' function.

After the pre-processing of the dataset, we built different models using 11 different algorithms available to use via the 'caret' package: 'CSimca' (SIMCA<sup>14</sup>), 'kknn' (k-Nearest Neighbors<sup>15</sup>), 'lssvmRadial' (Least Squares Support Vector Machine with Radial Basis Function Kernel<sup>16</sup>), 'multinom' (Penalized Multinomial Regression<sup>17</sup>), 'naive\_bayes' (Naïve Bayes<sup>18</sup>), 'ordinalNet' (Penalized Ordinal Regression<sup>19</sup>), 'pam' (Nearest Shrunken Centroids<sup>20</sup>), 'rf' (Random Forest<sup>21</sup>), 'rpart2' (CART<sup>22</sup>), 'sparseLDA' (Sparse Linear Discriminant Analysis<sup>23</sup>), and 'treebag' (Bagged CART<sup>24</sup>). We selected these algorithms based on maximal dissimilarity sampling starting from the 'rf' algorithm<sup>25</sup>.

## 423 Model evaluation

424 To evaluate the models, we assessed the sensitivity, specificity, positive predictive value (PPV), and 425 negative predictive value (NPV) for the 3 labels (for each model). To summarize those 4 metrics, we 426 averaged their values at the label level, resulting in 3 metrics per model. We named each of these values 'equilibrated accuracies'. Next, to select the best model per comparison, the mean of these equilibrated 427 428 accuracies (average of the previous 3 values) was computed; and the model with the highest value was 429 selected. Even so, we verified in all cases that the metrics were not biased for labels that were relatively less important (i.e. a model with high performance for the 'dubious' label but low for the 'relevant' label). 430 Concretely, we first selected models that presented a near-maximum equilibrated accuracy for the 431 432 'relevant' label. A near maximum value was defined as a value with at least 98.5% of the maximum value 433 (i.e. if the maximum value was 1, values between 0.985 and 1 were considered near maximum). Out of 434 the selected models, we further selected those with near-maximum equilibrated accuracies of the 435 'irrelevant' label. And lastly, we selected models with near-maximum equilibrated accuracies of the 'dubious' label. If several models were selected in the last step, the model with the highest average 436 equilibrated accuracy was selected. 437

To calculate the variable importance for the models, we used the 'varImp' function from the 'caret' package. We used the default parameters to calculate and plot the variables' importance. The importance values were scaled between 0 and 100. This resulted in different calculations depending on the algorithm. For example, 'rf' and 'treebag' presented their variable importance calculation, while algorithms like 'sparseLDA' used ROC (Receiver Operating Characteristic) curve variable importance, where the variables were sorted by maximum importance across the classes.

## 444 Model comparison

To compare the models' accuracies across different datasets, for each model trained with a specific dataset, we tested each of those datasets: 3R, 9R, 21R, and MixR datasets. In case the same dataset would be used for training and testing, we split the full dataset into training and testing datasets to prevent an overfit. The accuracy metric used was Cohen's Kappa Coefficient, as it takes into account class imbalance while providing a straightforward metric.

## 450 Method comparison

- 451 We utilized the same initial datasets to compare the different methods to detect differential expression. 452 The p-values, FDR values, and log2 fold-change values were all obtained from the results table from 453 DESeq2 for each comparison. We use the '&' symbol to indicate that the conditions on both sides of the 454 symbol need to be true. The methods were as follows: FDR < 0.05, FDR < 0.05 & |log2 Fold Change 455 value| > 1.5, FDR < 0.01, FDR < 0.01 & |log2 Fold Change value| > 1.5, R-ODAF, and 'Relevant'
- 456 labelled genes by our model.

#### 457 Enrichment analysis

We used GORILLA<sup>26</sup> for the Gene Ontology enrichment analysis. The advanced parameters included a 'P-value threshold' of  $10^{-3}$ , and fast mode was enabled. The data used were derived from the 21R dataset. The FDR 0.05 gene set contained the significantly differentially expressed genes based on an FDR threshold of 0.05, which contained 5811 genes. The RELEVANT gene set contained the genes labeled as 'relevant' by our 'AutoRel' model, which contained 2686 genes.

## 463 Simulated dataset

The simulated datasets were generated using the SPsimSeq R package<sup>27</sup>. The Zhang bulk RNA-seq data<sup>28</sup> was used as a reference for the data generation process. Four comparisons were made using a different number of replicates: 50 (as performed in the package demonstration), 21, 9, and 3. The total number of genes was 3000, and the proportion of Differentially Expressed Genes was 10% (300 genes), with at least a 0.5 log-fold-change in the source data. Each of the 4 comparisons was performed 100-fold, generating a new/random dataset for each iteration. We applied the DESeq2 pipeline and used the consequent results as input for AutoRel's prediction process.

<sup>1</sup> Supplementary Methods (Chapter 5)

## 2 Model features

3 The features used in the machine learning model are derived from the normalized counts. The mean, 4 standard deviation (SD), and variance of all normalized counts per gene were added as features. 40 5 features were added (20 per group), each of them representing a quantile from 0% to 100% every 5% (i.e. quantile 55% control). 20 features were related to the normalized values: all replicates (per group) were 6 7 divided into 10 subgroups based on the order of their imputation. The median of each subgroup was obtained and saved as a feature (i.e. out of 21 samples, '1st subset median' would be the median between 8 the 1<sup>st</sup> and 2<sup>nd</sup> replicates; '2<sup>nd</sup> subset median' would represent the median between the 3<sup>rd</sup> and 4<sup>th</sup> replicates, 9 etc.). In addition, four features of N (absolute number) and proportion (N / number of samples) of zero 10 11 values per group (i.e. N nonexpressed samples treatment) were added. Next, 12 features for outlier 12 features (N and proportion per group for 3 groups of outliers: total, mild, and extreme) were included. For 13 each of the aforementioned features, an additional feature was added, where the feature of the treatment 14 was divided by the same feature for the control group, whose feature name would start with 'foldchange' 15 (foldchange N mild outliers).

Quartile rules: all potential differences between the quartiles of the 2 groups were described with different 16 17 features. Four degrees of changes were recorded as features depending on the shift between both group 18 distributions: mild changes with only a 25% difference (i.e. the 50% quantile of the treatment group above the 75% quantile of the control group), 50% difference (i.e. 50% quantile of the treatment group 19 20 above the 100% quantile of the control group), 75% difference (i.e. 25% quantile of the treatment group 21 above the 100% quantile of the control group), and 100% difference (i.e. 0% quantile of the treatment 22 group above the 100% quantile of the control group). The names for the rules were, respectively: 'onequartilediff rule', 'twoquartilediff rule', 'threequartilediff rule', and 'fourquartilediff rule'. In 23 addition, all possible differences were added as features (i.e. 'q1belmin' stood for 25% quantile of the 24 25 control group below 0% quantile of the treatment group). The last feature regarding quartiles was named 'quartilediff score'. The score represented an overall description of all potential changes between the two 26 27 distributions. Any quartile rule where the treatment group values were above the control ones contributed +1 to the score. Any quartile rule where the treatment group values were below the control ones 28 29 contributed -1 to the score. This resulted in a score that ranged between -10 and +10. For example, if the 30 quantile 50% control group was below the quantile 0% treatment group, this resulted in two rules being true: 'q2belmin' and 'q1belmin'; thus the resulting score would be +2. If simultaneously the 75% quantile 31 32 control group was above the 100% quantile treatment group, a negative rule would also apply ('q3abomax'), and thus the score would result in +1. 33

Aside from the third quartile rule, which inspired the aforementioned rules as features, we also incorporated additional rules inspired by the R-ODAF method<sup>1</sup>. The "expression consistency" rule uses values transformed to counts per million (CPM), and for an observation to be true, 75% of the values on either the control or treatment groups need to be above 1 CPM. A "spurious spike" rule was designed to detect single samples with most of either group's reads, where the threshold limit for the proportion of reads depends on the number of replicates: *Spike threshold*=1.4  $N^{-0.66}$ , where N is the number of replicates of the group).

Another source of features for the dataset was the result table from each DESeq analysis. The result table included base means across samples ('baseMean'), standard errors ('lfcSE'), log2 fold changes ('log2FoldChange'), p-values ('pvalue'), adjusted p-values ('padj'), and test statistics ('stat'). An additional feature was added from R-ODAF, where an FDR threshold limit of 0.01 was established ('fdrlowerthan0.01'). All NA values from the result table for 'log2FoldChange' and 'stat' were transformed to zero, while all NA values for 'lfcSE', 'pvalue', and 'padj' were set to one.

# <sup>1</sup> Supplementary Results (Chapter 5)

# 2 Dataset 3R: 3 replicates per group

3 The 'relevant' class was generally the best performing across the three classes, followed closely by 4 'irrelevant' (Figure 1). The 'dubious' class was generally the worst-performing class, where a clear 5 difference could be seen for all models but one ('CSimca'). We selected 'treebag' as the best-performing 6 model on average.



8 Figure 1: Equilibrated accuracies for the models built with 11 machine learning algorithms using the 3R dataset. The models
9 were sorted by the average of equilibrated accuracies in descending order.

To further investigate the best-performing model, we calculated the variable importance of the 'treebag' model ('treebag\_3R') to evaluate which features were crucial for a gene to be classified. Only the 20 most important features were plotted (Figure 2). Seven features distinguished as predominantly important, from most to least importance: 'pvalue', 'padj', 'rule\_cpm\_0.75\_above\_1', 'fdrlowerthan0.01', 'lfcSE', 'threequartilediff rule', and 'twoquartilediff rule'.



16 Figure 2: Variable importance for the 10 most important features sorted in decreasing order for the 'treebag' model trained with17 the 3R dataset.

18 The first two relate to the statistical significance, which is confirmed to be of primordial importance in a high number of observations even with only three replicates. The CPM rule described whether a gene was 19 20 considered expressed in at least 75% of the samples of either group. This rule was essential in our 21 decision-making during the labeling process, and thus it was also reflected in the model's variable importance. 'fdrlowerthan0.01' was able to label the most extreme cases of statistical significance. 22 'lfcSE' (Log Fold Change Standard Error) reflected the confidence of the log2 Fold Change, while the 23 feature that represented the latter 'log2FoldChange' was situated in the lower half of the variable 24 importance ranking. The last two features of the most important seven were related to the quartile 25 differences between groups: these "rules" showed how different both groups were by comparing their 26 quartiles. 27

# 28 Dataset 9R: 9 replicates per group

Even though the 'treebag' algorithm also performed correctly in the contrast with nine replicates per group (Figure 3), the best performing model was 'sparseLDA'. 'multinom' and 'treebag' presented a higher equilibrated accuracy for the 'relevant' class, but the difference was minimal. On the other hand, the equilibrated accuracy for the 'dubious' class was substantially higher for sparseLDA, leading to a higher mean equilibrated accuracy.



Figure 3: Equilibrated accuracies for the models built with 11 machine learning algorithms using the 9R dataset. The models
were sorted by the average of equilibrated accuracies in descending order.

The variable importance calculation was obtained via the use of the ROC curve (Figure 4). As there were 37 38 variable importances for each class, the variables were sorted by maximum importance across the classes. 39 Focusing on the 20 most important variables, we were able to detect some similarities with the most important ones from the previous contrast. All the seven most important features for the 'treebag 3R' 40 41 model were also present: the quartile difference rules ('twoquartilediff rule', 'onequartilediff'), the statistical features ('padj', 'pvalue', 'fdrlowerthan0.01', and 'lfcSE'), and the CPM rule 42 43 ('rule cpm 0.75 above 1'). Other features related to those groups were also detected ('threequartilediff rule' and 'baseMean'). In addition, we observed important features related to the 44 45 ('X8th subset median APA The', 'var APA The', treatment group 'sd APA The', 'Proportion nonexpressed samples APA The'), and up-regulation ('q3belq2', 'q2belq1', 'q1belmin', 46 47 'q3belq1', 'q2belmin', and 'q3belmin').



50 Figure 4: Variable importance for the 20 most important features for the 'sparseLDA' model trained with the 9R dataset. The

51 *features were sorted in a decreasing order based on the maximum importance value across the three classes.* 

52

# 53 Dataset 21R: 21 replicates per group

54 When analyzing the results for the equilibrated accuracies of this contrast (Figure 5), we noted that most 55 models presented an equilibrated accuracy for the 'dubious' class of approximately 0.5. This was due to 56 the lack of prediction of this class for these models, even though the training dataset included

57 observations from that class. The best performing model was thus random forest ('rf'), with high metrics

for 'relevant' and 'irrelevant' classes, but especially for the 'dubious' class, being the only model above0.6.



*Figure 5: Equilibrated accuracies for the models built with 11 machine learning algorithms using the 21R dataset. The models*were sorted by the average of equilibrated accuracies in descending order.

As performed for the previous contrasts, we evaluated the most important variables for the 'rf' model 63 64 (Figure 6). The 'threequartilediff\_rule' was predominantly the most important feature, having an 65 importance value 400% higher than the second most important feature. Again, we observed similarities 66 with the previous datasets: quartile difference features ('threequartilediff rule', 'quartilediff score', and CPM 67 'maxbelq1'), statistical features ('pvalue', 'padj', and 'lfcSE'), and the rule 68 ('rule cpm 0.75 above 1').



*Figure 6: Variable importance for the 10 most important features sorted in decreasing order for the 'rf' model trained with the* 

*21R dataset.* 

# 73 50 replicates



Figure 7: Boxplot of the ratio of several metrics between relevant and significant genes with 50 replicates per group. For every simulated dataset, a metric is derived for both relevant and significant genes compared to the population/original difference. The result of each metric from the relevant list of genes is divided by the result of each metric from the significant list of genes. If they perform equally, the value will be 1. If the relevant list of genes performs better, the value of the ratio will be above 1. If the relevant list of genes, the value of the ratio will be below 1. A vertical line at 1 visualizes where the equilibrium is situated.



Figure 8: Boxplot of the ratio of several metrics between relevant-dubious and significant genes with 50 replicates per group.
For every simulated dataset, a metric is derived for both relevant and significant genes compared to the population/original
difference. The result of each metric from the relevant list of genes is divided by the result of each metric from the significant list
of genes. If they perform equally, the value will be 1. If the relevant list of genes performs better, the value of the ratio will be
above 1. If the relevant list of genes performs worse than the significant list of genes, the value of the ratio will be below 1. A
vertical line at 1 visualizes where the equilibrium is situated.

# 88 21 replicates



# 89

Figure 9: Boxplot of the ratio of several metrics between relevant and significant genes with 21 replicates per group. For every
simulated dataset, a metric is derived for both relevant and significant genes compared to the population/original difference. The
result of each metric from the relevant list of genes is divided by the result of each metric from the significant list of genes. If they
perform equally, the value will be 1. If the relevant list of genes performs better, the value of the ratio will be above 1. If the
relevant list of genes performs worse than the significant list of genes, the value of the ratio will be below 1. A vertical line at 1

95 visualizes where the equilibrium is situated.



97 Figure 10: Boxplot of the ratio of several metrics between relevant-dubious and significant genes with 21 replicates per group.
98 For every simulated dataset, a metric is derived for both relevant and significant genes compared to the population/original
99 difference. The result of each metric from the relevant list of genes is divided by the result of each metric from the significant list
100 of genes. If they perform equally, the value will be 1. If the relevant list of genes performs better, the value of the ratio will be
101 above 1. If the relevant list of genes performs worse than the significant list of genes, the value of the ratio will be below 1. A
102 vertical line at 1 visualizes where the equilibrium is situated.

# 103 9 replicates



# 104

Figure 11: Boxplot of the ratio of several metrics between relevant and significant genes with 9 replicates per group. For every
simulated dataset, a metric is derived for both relevant and significant genes compared to the population/original difference. The
result of each metric from the relevant list of genes is divided by the result of each metric from the significant list of genes. If they

108 perform equally, the value will be 1. If the relevant list of genes performs better, the value of the ratio will be above 1. If the

109 relevant list of genes performs worse than the significant list of genes, the value of the ratio will be below 1. A vertical line at 1

110 *visualizes where the equilibrium is situated.* 



Figure 12: Boxplot of the ratio of several metrics between relevant-dubious and significant genes with 9 replicates per group.
For every simulated dataset, a metric is derived for both relevant and significant genes compared to the population/original difference. The result of each metric from the relevant list of genes is divided by the result of each metric from the significant list of genes. If they perform equally, the value will be 1. If the relevant list of genes, the value of the ratio will be above 1. If the relevant list of genes, the value of the ratio will be below 1. A

*vertical line at 1 visualizes where the equilibrium is situated.* 

# 118 3 replicates



# 119

Figure 13: Boxplot of the ratio of several metrics between relevant and significant genes with 3 replicates per group. For every simulated dataset, a metric is derived for both relevant and significant genes compared to the population/original difference. The result of each metric from the relevant list of genes is divided by the result of each metric from the significant list of genes. If they perform equally, the value will be 1. If the relevant list of genes performs better, the value of the ratio will be above 1. If the relevant list of genes performs worse than the significant list of genes, the value of the ratio will be below 1. A vertical line at 1

125 *visualizes where the equilibrium is situated.*


#### 126

127 Figure 14: Boxplot of the ratio of several metrics between relevant-dubious and significant genes with 3 replicates per group.
128 For every simulated dataset, a metric is derived for both relevant and significant genes compared to the population/original
129 difference. The result of each metric from the relevant list of genes is divided by the result of each metric from the significant list
130 of genes. If they perform equally, the value will be 1. If the relevant list of genes performs better, the value of the ratio will be
131 above 1. If the relevant list of genes performs worse than the significant list of genes, the value of the ratio will be below 1. A

132 vertical line at 1 visualizes where the equilibrium is situated.

## Chapter 6:

# **General Discussion and Summary**

### 1 Chapter 6: General Discussion and Summary

OMICs data has become the new tool to study all kinds of areas related to molecular biology: from cancer 2 3 and its prognosis<sup>1</sup> to microbiology discovery<sup>2</sup>. Its use has become so universal that the data generation process has led to an enormous library of molecular data. Even though the size of data is crucial for 4 5 research, its management and analysis are at least as important. In the case of this thesis, the data used 6 originated in the HeCaToS project. This project spanned tens of compounds, with different doses, 7 timepoints, and cell models. Even so, the abundance of data was not free of important limitations. The 8 number of replicates was minimal (three). In addition, these were technical replicates, i.e., repeated 9 simultaneously and originating from the same batch of cells. A great limitation of such design is that the 10 results obtained are at best indicative of potential changes, but cannot be used as evidence per se. The full 11 investment in the data generation process, while helpful, undermines its own results by not being able to 12 have access to further resources to validate the results. A different approach might be to increase the 13 number of both technical and biological replicates, even at the cost of decreasing the number of 14 conditions tested. The danger of the latter approach is that, while being more scientifically sound, it might 15 hinder its financing, as it may be evaluated at face value to be less ambitious than other broader-scoped 16 projects.

17 As evidence of our intention to be biologically sound and more data skeptical, one can observe that in Chapter 2 we applied *in vitro* techniques to validate the existence of the circular RNAs studied. Ideally, 18 19 we would have not only validated their existence but ratified their effect at the toxicological level by either overexpressing or knocking down each of the transcripts described in the post-transcriptional 20 21 regulation axes. This is especially relevant for the circRNAs themselves, which, even though they have 22 been frequently studied since their hypothesized function as miRNA sponges<sup>3</sup>, some bioinformatic analyses showed that this is rather the exception than the rule<sup>4,5</sup>. In addition, artificially modifying the 23 24 expression levels of the studied circRNAs would help us discern whether those changes are simply a 25 result or also a mechanism of the toxicity exerted by those compounds. Other modifications would have 26 introduced more specificity in our novel method to detect and quantify circRNAs. For example, one of 27 the main weaknesses of the remapping strategy is that circRNAs are not mapped based on their most 28 characteristic property: the back-spliced junction (BSJ). This is because the sequences introduced to be 29 mapped were structured in the classical or genomic 5' to 3' form. To this effect, no read that mapped the BSJ would be assigned to their circRNA sequence. One solution would be to truncate or split the original 30 31 sequences in half and join the 3' end of the second half to the 5' end of the first one. In this manner, the 32 characteristical BSJ would lay in the middle of the sequence. The reads mapping to the BSJ would not only increase the read count for each circRNA but would help to more correctly assess their expression. 33

This is due to the algorithms used by quantification software, which assign shared reads (reads that map to more than one location) by looking at the proportion of non-shared reads (reads that uniquely map to a genomic locus). This change in the transcriptome file would still allow quantifying circRNAs more broadly than CIRI2<sup>6</sup>, as circRNAs without a BSJ read would still be identified, but simultaneously improve the quantification step for the circRNAs whose BSJ reads were sequenced.

39 As seen in Chapter 2, the integration of all the multi-omics is not a straightforward process. Even the 40 study of two circular RNAs leads to a complex network of interactions that affect the translation of messenger RNAs, upon which the whole structure and function of the cell are based. For this reason, we 41 42 proposed in Chapter 3 to formulate an equation that would help quantify the number of transcripts available for translation. Even though from the toxicological view the dataset used (HeCaToS) involved 43 certain complexities due to the various factors to consider (time, dose, or batch effects), it included 44 different omics data extracted simultaneously from the same samples, thus a special design for a big 45 46 dataset that is truly beneficial when aiming for a new crossomics (integration of multiple omics data) 47 strategy. In addition, the fact that the RNA-Sequencing data was ribo-depleted (instead of using poly (A) capture) permitted the quantification of circRNAs, thus expanding the level of RNA types analyzed. 48 49 Including important factors that affect the levels of free transcripts can become endless, not only for the 50 calculations necessary, but also due to data that is not accessible at that moment in time (such as 51 epigenetics or translatomics), and the different number of factors that affect each gene: some will be 52 limited by transcriptional factors while others might be mostly regulated by histone acetylation. Thus, 53 even to the best of our efforts, the results showed that the prediction of transcript availability was much more complex than a formulation derived from our current knowledge of the main drivers of post-54 transcriptional regulation. 55

This issue brought about our published work<sup>7</sup> encapsulated in **Chapter 4**, where we addressed the issue 56 57 from a different perspective. First, we selected the proteomics expression as the target to be predicted. 58 Second, we switched our method from human formulation to machine learning. These changes were 59 sounder for several reasons. First, having proteomics data allowed us to directly evaluate the accuracy of 60 our research. Second, machine learning is mainly limited by the data it uses. So, if we fed the data in the most comprehensive way possible, it would output the most optimized prediction, as was the case. Even 61 62 so, one would acknowledge that the information available is still far from completely describing the molecular landscape of a cell. And so, even though the metrics showed a good prediction accuracy in 63 terms of R<sup>2</sup>, it is fair to wonder how much of the model accuracy is due to being trained in untreated 64 conditions: if the model has been trained with the same protein on a different sample, outputting the same 65 66 value that was trained on is a safe bet which would also lead to a potentially good correlation score. In

67 any case, if the conditions were met, the ideal situation to generate the best model would be to make a 68 model for each protein (to be predicted), where the variables encapsulate all other molecules. In this 69 manner, all potential effects would be considered, from miRNAs to transcription factors, and thus if such an abundant source of data was available, one could in principle predict any molecule's quantity based on 70 all other variables. Of course, such a model would only be useful if one would also have access to the 71 huge number of variables that trained the model, and thus it would probably be more of a proof of 72 73 concept of whether we can completely map and understand the molecular biology of the cell. In addition, 74 such a model would allow us to modify the abundance of specific molecules in order to study the effects 75 on all the components of the cell. A model with these characteristics is, in the author's opinion, still not possible nowadays due to technological limitations, and so one would caution in favor of skepticism 76 77 whenever projects aim to develop a model of similar characteristics to predict all toxic effects inside a 78 cell.

79 This is though not a leeway as a scientist to presume that using machine learning is per se a universal 80 method to solve current scientific problems as long as there is enough data. As an example, just feeding a model with all published transcriptomics data and assumed results derived from that will mostly result, in 81 82 the author's opinion, in a big model with very stochastic results. This is because the human component in 83 the decisions taken, be it due to different methods, thresholds, or prejudices used, will ultimately lead to 84 different conclusions, and these decisions would all be mixed into a single model. For this reason, even 85 though one of the main limitations of Chapter 5 is that most of the observations fed into the AutoRel 86 were labeled by the author of this thesis, it also ensures a more replicable model of the reasoning behind those decisions, that is, that the model can more accurately label observations in concordance to how it 87 has been trained. As mentioned in the chapter, this would represent a further step in the direction of 88 automatizing manual filtering steps in RNA-Sequencing analyses, as was the case for our collaboration in 89 R-ODAF<sup>8</sup>. In any case, it can be argued that different strategies could have been applied to represent more 90 91 extensively the expert decision-making. The most direct manner would have been to introduce all 92 possible interpretations from the maximal number of scientists in a single model. One could also instead 93 have built a model per researcher, and a global algorithm would be used to assign a label based on the 94 majority vote from all those distinct models. A different strategy could have been to only label as relevant 95 genes that have been previously proven to be differentially expressed via several distinct confirmation experiments, or at least train these models more extensively by using simulation datasets, where the 96 97 differentially expressed genes (DEGs) are known beforehand. A limitation of the use of simulated datasets, though, is that it selects its DEGs by fold change alone, which is a limited metric to use. It is 98 99 based on the comparison of means, which are easily tilted by the presence of outliers. In addition, 100 simulated datasets do not take into account low quantification values, and thus a difference of a few reads

is considered a DEG, as small values have a higher chance at random to present a higher fold changevalue. In practice, this can easily be the result of an irrelevant small variation of expression.

103 It is of interest to note the use of previous machine learning approaches to evaluate differentially expressed genes. For example, LASSO regression analysis is usually utilized after a differential 104 expression analysis to further filter the number of DEGs<sup>9</sup>. The use of random forest has been recently 105 suggested for the selection of important genes in microarray data<sup>10</sup>. The theory is based on the hypothesis 106 that if two conditions (control vs treatment) are able to be separated by the data available (gene 107 108 expression data), an artificial intelligence model should be able to differentiate between them. In addition, 109 the commonplace existence of variable importance values would help identify which variables (genes) are the most responsible for the differences/splitting between the two conditions. We consider the evaluation 110 of several classification algorithms (including LASSO and random forest) of great interest to discern to 111 which extent these algorithms can correctly identify such genes, by the use of validated gene expression 112 113 results.

### 114 Summary

Novel bioinformatics approaches have been applied to elucidate the function of circular RNAs in the mechanism of action of several cardiotoxicants through their importance in the post-transcriptional regulation. We further our investigation of this regulation via a formulation to predict the quantity of coding transcripts that are available for translation. The complexity of this task makes necessary the introduction of machine learning strategies, which help us predict and impute proteomics values in untreated samples. Finally, we make use of the artificial intelligence methods to classify genes according to their biological relevance based on expert labeling.

### 122 Bibliography

- Lin, E. *et al.* Integrative Analysis of the Genomic and Immune Microenvironment Characteristics
   Associated With Clear Cell Renal Cell Carcinoma Progression: Implications for Prognosis and
   Immunotherapy. *Front Immunol* 13, 830220, doi:10.3389/fimmu.2022.830220 (2022).
   Grujcic, V., Taylor, G. T. & Foster, R. A. One Cell at a Time: Advances in Single-Cell Methods
- and Instrumentation for Discovery in Aquatic Microbiology. *Front Microbiol* 13, 881018,
  doi:10.3389/fmicb.2022.881018 (2022).
- 129 3 Ma, C. *et al.* circRNA CDR1as Promotes Pulmonary Artery Smooth Muscle Cell Calcification by
- 130 Upregulating CAMK2D and CNN3 via Sponging miR-7-5p. *Mol Ther Nucleic Acids* 22, 530-
- 131 541, doi:10.1016/j.omtn.2020.09.018 (2020).

- Guo, J. U., Agarwal, V., Guo, H. & Bartel, D. P. Expanded identification and characterization of 132 4 133 mammalian circular RNAs. Genome Biol 15, 409, doi:10.1186/s13059-014-0409-z (2014). 134 5 You, X. et al. Neural circular RNAs are derived from synaptic genes and regulated by development and plasticity. Nat Neurosci 18, 603-610, doi:10.1038/nn.3975 (2015). 135 6 Gao, Y., Zhang, J. & Zhao, F. Circular RNA identification based on multiple seed matching. 136 Brief Bioinform 19, 803-810, doi:10.1093/bib/bbx014 (2018). 137 7 Ochoteco Asensio, J., Verheijen, M. & Caiment, F. Predicting missing proteomics values using 138 machine learning: Filling the gap using transcriptomics and other biological features. Comput 139 Struct Biotechnol J 20, 2057-2069, doi:10.1016/j.csbj.2022.04.017 (2022). 140 8 Verheijen, M. C. et al. R-ODAF: Omics data analysis framework for regulatory application. 141 Regul Toxicol Pharmacol 131, 105143, doi:10.1016/j.yrtph.2022.105143 (2022). 142 9 143 Miduo Tan, G. H., Jingjing Chen, Jiansheng Yi, Xi Liu, Ni Liao, Yi Hu, Wei Zhou, Qiong Guo. Construction and validation of an eight pyroptosis-related lncRNA risk model for breast cancer. 144 American Journal of Translational Research 14(5), 2779-2800 (2022). 145 146 10 Hediger, S., Michel, L. & Naf, J. On the use of random forest for two-sample testing. Comput Stat Data An 170, doi:ARTN 107435 147
- 148 10.1016/j.csda.2022.107435 (2022).

149

## Impact

## <sup>1</sup> Scientific and Social Impact

Currently, regulatory agencies rely on the use of animals for testing drugs for several reasons. First, the 2 3 most accurate way of testing drugs would be testing them in humans, which is not ethically possible. Second, modeling everything that happens in our bodies in a Petri dish is nowadays unfeasible: the 4 complexity of a human body requires 30 trillion cells to keep functioning. In addition, each of those cells 5 6 also involves its level of complexity: even a simple yeast cell contains 42 million proteins<sup>1</sup>, which does 7 not include other essential molecules like sugars and fatty acids. Even so, recent incentives (such as the regulation (EC) No 1223/2009 of the European parliament<sup>2</sup>) have pushed the scientific community to 8 search for alternative testing methods without the use of animals. Specifically, the area of Toxicology is 9 affected, as it studies the potentially toxic effects of drugs both before and after being released on the 10 market. A popular method for studying the effects a compound can have in humans is by testing them in 11 12 *vitro*, that is, by exposing human cells. Doing so helps to narrow the bridge between what is being tested (human cells outside the body) and the actual goal of the study (human cells inside the body) when 13 compared to animal testing, where the model and the end goal belong to different species. Nowadays, the 14 15 development of induced pluripotent stem cell (iPSC) technology allows reverting any human skin cell to another tissue type (such as cardiac cells that contain the same DNA as the donor of the skin cells), 16 17 without the need for surgery or invasive biopsies. Although some drugs may kill cells by simply 18 destroying the membrane that encapsulates the cell, most of them disturb the cell in more subtle ways. One of these ways is the deregulation of the number of proteins synthesized by a cell. Proteins are 19 essential molecules, as they perform most of the cell functions. For this reason, disturbing or blocking 20 21 their production can lead to the disruption or death of a cell and/or the ones that depend on it. The processes that lead to the making of proteins involve mainly RNAs, molecules that work as messengers 22 23 from DNA to proteins. Therefore, in Toxicogenomics, studying how specific treatments can affect these 24 molecules can help understand their mechanisms of toxicity.

25 In Chapter 2, we assessed how a recently discovered class of RNA, called circular RNAs (circRNAs), 26 are disturbed in heart cells by known toxicants. These circRNAs have been hypothesized to regulate microRNAs (miRNAs) by letting the latter bind to the former. When messenger RNAs (mRNAs) are not 27 28 bound to miRNAs, they can provide the instructions to produce proteins. When miRNAs are occupied binding circRNAs, they are not able to bind to mRNAs. For this reason, by assessing changes in the 29 number of circRNAs, miRNAs, mRNAs, and proteins; we helped better understand how these 30 compounds (that are still in use) are being toxic in the human heart without the use of animal testing. 31 Going further into the thesis, we realized that knowing how many proteins are in a cell is crucial to 32

understanding how a drug (or disease, or any other perturbation) affects a cell. Unfortunately, the
technology used for doing it, mass spectrometry, does not measure all proteins in a cell. Instead,
researchers tend to use transcriptomics, which exhaustively measures the quantity of RNAs. Nevertheless,
the number of RNAs and the number of proteins do not always perfectly correlate with each other.

37 In Chapter 3, we designed an equation to estimate how many RNAs are available to produce proteins. We did that by counting the total number of mRNAs and subtracting the ones that will be affected by 38 39 miRNAs, but only those miRNAs that are not binding to other RNAs (like another mRNA or circRNA). 40 Nonetheless, the formula, which is focused mainly on RNA molecules, demonstrated an added value for 41 only a subset of proteins. As a result, in Chapter 4, we went a step further. We built a large dataset with RNAs and their corresponding proteins and trained a machine learning model to predict the latter. 42 Machine learning algorithms "learn" how to predict values by looking at how other similar values behave. 43 Using our data, our model predicted well the increases and decreases of proteins, which can help others to 44 45 predict how many proteins there are in a sample of cells based on how many there are in similar ones.

46 As mentioned before, in the area of Toxicogenomics it is of great interest to study the changes happening 47 in a cell. Transcriptomics is exceptionally good at counting how many RNAs there are, consequently 48 using this technology helps us understand which molecules change in quantity due to a specific cause. 49 The statistical tools used to detect changes, though, do not work without fault. For this reason, experts in 50 this technology can manually detect these errors. On the flip side, this manual curation is pretty timeintensive when taking into account the number of genes to be evaluated, and requires specialist 51 52 knowledge to do so. That is why, in **Chapter 5**, we again trained a machine learning model. In this case, 53 we taught the model to recognize the profile of genes that are typically of interest to the researcher. We built several of them with different characteristics and selected the best one, which we named 'AutoRel'. 54 Even though AutoRel was not flawless, it showed improvements by removing genes that were not of 55 interest. 56

57 In an era of increasing societal pressure against animal testing, added to the inherent shortcomings of 58 animal assays, regulatory agencies need to reevaluate their historical procedure of risk assessment. With 59 the rapid development of methodologies that allow the analysis of the complete set of biological entities 60 in a cell exposed to any substance, regulators will need both a better understanding of all the complex 61 interactions behind molecular biology and powerful data analysis tools to integrate them. The work of this 62 thesis contributes to this necessary transition toward a next-generation risk assessment. This is achieved 63 by the discovery of new changes that happen when a toxic compound affects a cell, predicting protein 64 measures that are usually unknown with artificial intelligence, and filtering results in an automated way to have a better understanding of the changes that occur in a cell. In aggregate, this contributes to assessmore accurately how toxic a drug is, making the use of treatments more safe and reliable.

### 67 Bibliography

68

Ho, B., Baryshnikova, A. & Brown, G. W. Unification of Protein Abundance Datasets Yields a
Quantitative Saccharomyces cerevisiae Proteome. *Cell Syst* 6, 192-205 e193,
doi:10.1016/j.cels.2017.12.004 (2018).

Regulation (EC) No 1223/2009 of the European Parliament and of the Council of 30 November
 2009 on cosmetic products (Text with EEA relevance),
 <<u>http://data.europa.eu/eli/reg/2009/1223/oj</u>> (2009).

75

Juan Ochoteco Asensio was born on the 29<sup>th</sup> of March of 1994 in Barcelona, Spain. In 2012, he started his Bachelor in Biochemistry at the University of Navarra. During his bachelor's, he took an internship in preanalytical and analytical phases in a clinical laboratory at the Biochemistry Lab of the Hospital Clinic de Barcelona. In addition, he attended a basic proteomics course at the Centre of Applied Medical Research in Pamplona (University Clinic of Navarra). For the bachelor's thesis, he worked on the analysis of the astroglia and microglia activation in the aging brain for the study of neurodegenerative diseases at the Centre of Applied Medical Research in Pamplona (University Clinic of Navarra).



After finishing his BSc degree in 2016, he started his Master's degree in Bioinformatics at the Faculty of Biosciences of the Autonomous University of Barcelona. For his master's thesis, he worked on a project that focused on re-annotating the genome of *Drosophila buzzatii*.

In March 2018, he started working as a PhD student at the Toxicogenomics department of Maastricht University under the supervision of Prof. Dr. Jos Kleinjans and Dr. Florian Caiment. His PhD project took advantage of the big datasets generat ed in the HeCaToS project, the latter of which aimed at developing integrative in silico tools for predicting human liver and heart toxicity. The results achieved throughout these four years are presented in this thesis.

## List of papers

### Published

- Verheijen MC, Meier MJ, Asensio JO, et al. R-ODAF: **Omics data analysis framework** for regulatory application. *Regul Toxicol Pharmacol*. 2022;131:105143. doi:10.1016/j.yrtph.2022.105143
- Ochoteco Asensio J, Verheijen M, Caiment F. Predicting missing proteomics values using machine learning: Filling the gap using transcriptomics and other biological features. *Comput Struct Biotechnol J*. 2022;20:2057-2069. Published 2022 Apr 22. doi:10.1016/j.csbj.2022.04.017

In submission process

• AutoRel: Machine Learning model for selecting differentially expressed genes. Juan Ochoteco Asensio, Jelmer Faber, Twan van den Beucken, Marcha Verheijen, Florian Caiment.

### Preprints

• Quantifying the number of translatable transcripts through the use of OMICs involved in post-transcriptional regulation. Juan Ochoteco Asensio, Jos Kleinjans, Florian Caiment. bioRxiv 2022.06.20.496876; doi: https://doi.org/10.1101/2022.06.20.496876

Other contributions

- Evaluation of transcriptomic quantification tools **Bioinformatics & Systems Biology** Conference 2019
- P11-06 Unraveling novel gene networks in anthracycline-induced cardiotoxicity. JG Faber, JO Asensio, F Caiment, T van den Beucken. Toxicology Letters 368, S164. 2022
- SOC-IV-09 Investigating the role of nuclear receptors in valproic acid-induced liver steatosis. K Guo, J Faber, JO Asensio, F Caiment, T van den Beucken. Toxicology Letters 368, S54-S55. 2022

# Acknowledgments

Let's be honest. Most people (or at least I would) will first come to this section to see if their names are here. Indeed, this is one of the lessons you learn in academia: we may like science, but that curiosity is also present in the human affairs. What is this section, you may ask? This is where I acknowledge that all these pages would not have been possible without so many people that have brought me to this moment, and for that I would like to thank them.

First, I would like to thank **Florian Caiment**. Thanks for teaching me so many valuable lessons about how to do research, which involves creativity (which we both share) and deadlines. I am grateful to have met a rigorous scientist in so many aspects, which has helped me being open and critical about the work we do. Interestingly enough, you also have fun while trying to make the most out of your work. All while not having sacrificed everything for your work. That sounds a pretty good future if I ever get there. And of course, it was so fun to always learn new French idioms with you, such as "Couper l'herbe sous le pied de quelqu'un".

**Prof. Dr. Jos Kleinjans**, thank you for giving me the opportunity to be a part of the last year of the HeCaToS project and providing the funding for my doctoral contract. Ever since your departure from the department, I have missed your productive motivation to move forward in delivering results.

I would like to acknowledge all the members of my assessment committee: **Prof. dr. H.J.M. Smeets**, **Prof. dr. J. Quackenbush**, **Prof. dr. R. Peeters**, **Dr. S. Hayat**, and **Dr. L. Eijssen**.

I also want to thank **Dr. T. van den Beucken** for his abundant help and contribution in the discussion of my work, giving always a realist criticism from a biological perspective, so necessary in the Bioinformatics field. Needless to say, I need to add also **Dr. M. Verheijen** here, as she has been there from the start helping me getting to understand the HeCaToS project, while also giving a friendly environment to the department. You literally celebrated my birthday on my second day of work! Even if our party attendance ratios are completely opposite. **Dr. D. Jennen** for your suggestions and hall conversations. **Dr. J. Briede** for your educational role. **Dr. J. Krauskopf** for your bioinformatic tips and bits. **Christa Graus** for so much help in everything needed in my day to day work. For our conversations at the office. Ik wens u nog success verder. Misschien zie ik je in Catalonië. **Duncan** for being the wisest of all of us, and make us wonder whether we are actually worth a doctoral degree. **Rosella**, la postdoc italiana più forte, grazie per la tua energia e immediatezza. Quei due sono sicuramente una combinazione potente. Saluta il ragazzo di Nazaret. E, come sempre, Forza Italia! **Margot**, ik heb je voor een korte tijd ontmoet, maar ik zie dat je heel energetisch bent, en dat is heel goed voor deze afdeling. Bedankt voor al onze gesprekken.

I would also like to thank my paranymphs. First, **Yannick Schrooders**. What can I say? I cannot imagine how my junior years could have been without you being part of the Gossip Guys group. It would have been then difficult to think how many things we would have done together, from you coming to my place in Barcelona, to me marrying you soon. In any case, you know I will there to give some of my energy for the Spirit Bomb. Weest altijd jezelf, en doe je ding! Mem, **Antoni Vallbona Garcia**, que

daixonare de tu? M'has vist amb xubec i fins i tot mesclar ous amb cargols. A vegades et sembla que cerc na Maria per sa cuina, i vols que a la defensa vagi gat perque mira que en soc de poma! De totes maneres, gracies per estar amb mi tot aquest temps. La nostra amistat demostra que persones molt diferents poden viure juntes, i esper que pugui saber de tu fins al final en aquest mon terrenal. Porta't bé, i un dels millors consells: digues la veritat (o com a minim no menteixis), i estima (o com a minim no odiïs).

I would like to thank Jarno for letting me have his place in the office. Thank you Evelyn for being so direct with everyone, to the point of leaving the office because of me. Thank you Rajinder for offering your house for parties (Raji shots!) and cooking such delicious foods (that I was able to try), in addition to let us know that cricket is a thing. Thanks to **Tim**, who helped me so much in the beginning of the doctorate with programming, and for sharing so many memories together. Hopefully your hockey skills are unmatched, and you keep reading as much as you do. Manon Marjoleen Jacqueline, what a fun time was to have you around in the office! So much energy and feminine taste in our otherwise male environment. And so much coffee and food of mine in your table! Marta, maggari, sei l'ultima collega che avevo in ufficio prima di essere estradato in isolamento e reclusione. Spero che le tue piante crescano alte e sane, e apprezzo così tanto le nostre lunghe conversazioni su così tante cose importanti. Mantieni intatta la tua curiosità e cerca sempre la verità! Sean, wat een gast hoor! Zo open om over alles te praten (en grappen makken). Ik wil jou danken voor al de tijd die je hebt in onze vriendschaap geinvesteerd, en ik wardeer je openheid. Je bent echt slim (maar wees nederig). De expert in belegging en motors! Gedraag je, en ik hoop dat we nog meer kunnen elkaar leren kennen in de toekomst. Thanks to all the interns that I have also shared a room with. Unfortunately, my memory is really bad, so thanks to Irene, Wendy, the Destiny 2 guy, and so many others that were fascinated with my madness. Also I remember now Nkoli, the tallest Brittish woman I had ever met, and almost as unhinged as me.

Now, time to thank other (ex)PhD students. Heloïse, enchanté de vous rencontrer dame française! Et maintenant maman apparemment. Nhan, the expert judo and colleague of HeCaToS. We started and will finish (God willing) the same year. Thanks for all the help throughout this time. Thanks for the Vietnamese transparent tacos and tiếng mèo kêu. Bạn là một người phụ nữ mạnh mẽ. Tiếp tục theo cách đó! Daniela, from the Spanish region of Portugal. Esta linda senhora é na verdade uma PhD agora! Espero que você seja o melhor na Suíça e continue rindo desse seu jeito característico. E lembre-se, o Amor e a Verdade prevalecerão! E eles realmente existem! E como você sempre disse, tudo ficará bem no final. Jelmer, je bent echt lief. Bedankt voor al je samenwerk, zonder jou zou ik deze thesis niet beeindigt. We zijn niet altijd eens met elkaar, maar nog heb ik echt genieten van deze tijd met jou. Wees altijd grappig and raar zoals altijd, met jou dansen in de gang en expert racen in Mario Kart. Nikki, de meest briljant Nederlans meisje die nog een PhD doet in ons afdeling. Bedankt dat je de meest normale van ons bent, en tegelijkertijd ook zo grappig en leuk het is om met jou te zijn. Zeg aan Mickey dat ik nog hem herinner (en dat hij raketten schoon maakt). Wijchen!. Nicolaj, der neue Deutsche in der Stadt. Schön, Sie auch kennenzulernen, und hoffentlich können Sie in einem so interessanten Büro überleben. باشد که این جوان ها را تا پایان دکترا شاد کنید. مواظب خودت باش من برای مادران .TGX جوان ترین دختر در .Somaieh احترام زبادي قائلم. Julia (& Frank), thanks for everything! So many deep conversations, your sweetness, your acrobatic skills, your work ethic, your wedding invitation, your heart blessings, your wine/knife selling, and so many other things. Thanks for your yeehaw time. Qingfeng, Zhiling, Yueqin, Kaidi, and Na:向所有中国女孩致敬。不幸的是,我总是很难记住你所有的名字。但是感谢中国的糖果! Milena, dziękuję za informacje o sytuacji bezrobocia w Holandii. Brian, nimekujua kwa muda mfupi tu, lakini ninathamini tabasamu lako na ukimya wako, ingawa kama Mhispania ningekuhimiza kuzungumza zaidi! **Saad**, میں بھی آپ کو صرف تھوڑی دیر کے لیے جانتا ہوں، لیکن کسی بھی بایو انفار میٹیشن کا ہمیشہ خیر مقدم ہوتا ہے۔ میں بھی آپ کو صرف تھوڑا گرم پائیں گے۔ امید ہے کہ آپ اس ملک کو برفانی سویڈن کے مقابلے میں کم از کم تھوڑا گرم پائیں گے۔

Salutations to other PhD's I have met from other departments: **Grzegorz** (jedyny inny katolicki dr. uczennica, która była na przyjęciach i mszy świętej), **Jana** (das klügste Mädchen, das fortschrittliche Einzelzell-Transkriptomik verwendet, um Menschen mit der Schönheit (ihrer Handlungen) zu beeindrucken), **Niccolo** (sei pazzo come me, goditi la vita negli Stati Uniti e continua a essere il ragazzo più sicuro di sé che abbia mai incontrato finora), **Delia** (arriba la gente del acueducto, eres genial y muy divertida, acuerdate de pórtarte bien y deja a la gente con un poco de sangre querida, espero que mis plaquetas te traten bien), **Julia** (che energica ragazza italiana, spero che tu possa continuare a scalare pareti sia fisicamente che nella vita), and **Celia** (bona damisela de Barcelona, espero que el teu temps a Maastricht te'n recordi encara que sigui una mica a Barna).

Thanks to the Twitch streamers **JPStudyTime** and **Sarayoposita** for their digital company in my working from home time.

Gracias a mi familia por estar siempre ahí. Sé que siempre podré volver a casa a descansar y disfrutar de la buena vida hasta que me sienta inútil y tenga que trabajar de nuevo. Para que quede patente, ahí van los nombres: Juan Ignacio, María Montserrat, Marcos, María, Isabel, Luis, Beatriz, Ignacio, Sofía, y los tres hermanos que nos esperan en el cielo. Os conozco desde hace ya tiempo, así que son demasiadas anécdotas y de todas formas no hace falta que os coma las orejas.

Bedankt **Melanie** om zo liefdevol met mij altijd te zijn. Je herinnert me van de liefde van God: bemminen zoals ik ben.