

15 **Abstract**

16 There is a paucity of molecules that progress through the drug development pipeline, making
17 the drug discovery process expensive and frustrating. Innovative approaches to drug
18 development are therefore required to maximise opportunities. Strategies like the Connectivity
19 Map (CMap), which compares >7,000 gene expression signatures generated from more than
20 1,000 drugs, can produce associations between currently unrelated therapeutics, unveiling new
21 mechanisms of action and favouring drug repositioning. Here, we discuss these opportunities
22 that could aid the drug development process and propose rigorous publication of 'omics' data
23 with open access and data sharing. We, pharmacologists of the third millennium, must aim
24 towards maximising knowledge in an unbiased and cost-effective manner, to deliver new drugs
25 for the global benefit of patients.

26

27 **Main text**

28 As learnt from Darwin's *Origin of Species*, it is not the strongest, nor the most intelligent of the
29 species that survives but the one that is the most adaptable to change. We could extrapolate
30 this statement to the current situation of the pharmaceutical industry, which seems unable to
31 sustain its own growth, due to the worldwide challenging economical climate and current
32 research strategies, perhaps too much seduced by technology and forgetting the unpredictable
33 nature of research discoveries [1, 2]. There is an unquestionable need for change and a re-
34 invention of the drug development process to guarantee, in a cost-effective manner, the
35 transition from basic research to patient benefit [3].

36

37 We now know that patients are not all the same, even if they receive the same diagnosis [4].
38 They may belong to a particular disease subtype that might require a specific therapy. The so-
39 called 'omics' (a suffix etymologically derived from the Greek, meaning the totality of something)
40 represent one of the best strategies to reveal differences between patients, as the study of the
41 totality of the genome, transcriptome, proteome, lipidome or metabolome does not require
42 previous knowledge on the nature of these differences.

43

44 Genomics, however, can contribute not only to patient stratification [5] but can also impact the
45 entire drug development process [6], including target identification, deciphering drugs
46 mechanisms of action, implementation of individualized medicines to seek optimal benefit for
47 each patient and to monitor drug response and toxicity. In this article we will discuss innovative
48 whole genome-based strategies that contribute to drug discovery and development by i)
49 identification of novel treatments for a specific disease, ii) discovery of mechanisms of action of
50 novel or known compounds and, finally, iii) for drug repositioning studies. We will also highlight
51 the need for more standardized methods and data-sharing policies to ensure full exploitation of
52 these findings into genuine clinical benefit.

53

54 ***Emerging strategies for drug discovery and drug repositioning***

55 The pharmaceutical industry needs to adapt according to the current economical situation. A re-
56 invention of the innovation process is necessary, as technological innovation has not been

57 proportionally translated into scientific innovation. Therefore, besides new instruments, new
58 concepts are needed to improve the efficiency of drug discovery [1, 2]. One of the main
59 consequences of any genome-wide study is the massive amount of information that is
60 generated. Whilst analyses of multiple hits can be more sophisticated than simple listing (up-
61 and down-regulated genes), current approaches tend to follow a more integrated interpretation
62 from a systems-oriented perspective [7-9].

63

64 A novel and powerful opportunity derives from the **connectivity map (CMap)** [10-12]. CMap is
65 an open-source software that allows a new interpretation of microarray data by comparing gene
66 expression profiles of interest with those obtained for hundreds of bioactive small molecules,
67 most of which are FDA-approved drugs. The most recent version (build 02,
68 <http://www.broadinstitute.org/cmap/>) of this database contains 7,056 gene expression profiles
69 from 1,309 bioactive compounds in 5 different human cell lines. The signatures contained in the
70 database can be compared with any gene-expression profile of interest following two
71 approaches: a *disease-centered* approach, when we use the gene expression profile of a
72 disease, and a *drug-centered* approach, when we use the gene expression profile of another
73 drug of interest. As a result, the 1,309 CMap drugs will be ranked according to the similarity with
74 the gene-signature of interest. Therefore, drugs with negative score (i.e. they present opposite
75 profiles to the signature of interest) might have the potential as new treatments for specific
76 diseases while drugs with positive score (i.e. they have similar gene expression profiles) could
77 be useful for identification of novel actions of existing drugs or to unravel drug mechanisms of
78 action [10] (Figure 1). Active efforts are currently being made to increase the capabilities of the
79 CMap. The new forthcoming version (<http://lincscloud.org/>) will represent a dramatic expansion
80 of the database and will contain almost one million of gene expression profiles. In addition to the
81 expansion in the number of pharmacological perturbagens (over 5,000 compounds), one of the
82 major novelties of the new CMap will be incorporation of genetic perturbations, that is gene
83 expression profiles obtained by up-regulation or down-regulation using shRNA of specific
84 genes, including drug targets and candidate disease genes.

85

86 Thus, the query of the CMap could be used for drug repositioning, that is, giving novel
87 indications for an existing drug [13, 14]. For example, the anticonvulsant drug topiramate was
88 linked (with a negative score) with the gene expression signature of IBD [15]. This prediction
89 was experimentally assessed using the trinitrobenzenesulfonic (TNBS) acid-induced colitis
90 model, in which the administration of topiramate significantly reduced intestinal inflammation.
91 Using a similar approach, the histone deacetylase inhibitor vorinostat was predicted as a
92 candidate therapeutic drug for gastric cancer, soliciting a series of *in vitro* investigations to
93 explore this functional association [16]. It is worth noting, that the CMap was proposed as a
94 'hypothesis generating tool', which means that confirmation studies are an absolute requirement
95 to validate initial predictions. Hassane *et al.* queried the CMap with the gene expression
96 signature produced by the drug parthenolide on acute myelogenous leukemia (AML) cells. This
97 drug was previously shown to ablate these cancer cells, and the predictions made with the
98 CMap led to the identification of novel agents (celastrol and 4-hydroxy-2-nonenal) that could
99 also markedly affect AML cells [17]. A CMap analysis also allowed Zhong *et al* to propose a
100 combination with angiotensin-converting enzyme inhibitors and histone deacetylase inhibitors
101 as a renoprotective therapy [18]

102

103 Interrogation of the CMap can also serve for the identification of novel mechanisms of action of
104 drugs. Hypoxia-inducible factor (HIF) 2a inhibitors were found by the CMap to be associated
105 (positive score) with the anti-inflammatory prostaglandin PGJ₂ [19]. This finding incited
106 subsequent experiments that showed how PGJ₂ was acting as an endogenous regulator of
107 HIF2a translation, suggesting this action as part of the anti-inflammatory effects of the
108 prostaglandin. The CMap approach has also facilitated identification of novel classes of drugs
109 including HSP90 inhibitors [20], and dissection of the mechanism of action of a traditional
110 Chinese medicinal herbal formula [21].

111

112 We have recently queried the CMap using the gene expression signature produced by the
113 endogenous pro-resolving mediator Annexin A1 (AnxA1) [22]; whilst this analysis produced
114 predictable associations, e.g. with non-steroidal anti-inflammatory drugs and glucocorticoids,
115 unexpected associations also emerged. In particular, the positive association with histone

116 deacetylase inhibitors (HDACIs) brought us to investigate whether a functional and mechanistic
117 link between AnxA1 and HDACIs could exist. Further experimentation made us conclude that
118 AnxA1 contribute to the anti-inflammatory mechanism of action of HDACIs [23].

119

120 Though innovative and promising, the CMap strategy is however not devoid of limitations,
121 although the new version **discussed above** might resolve some of them. Firstly,
122 pharmacologically relevant effects do not necessarily need to be reflected at the transcriptional
123 level. Secondly, the database was generated with a limited number of compounds and cell
124 lines. For example, the under-representation of certain drug classes, such as kinase inhibitors in
125 the current version (build 02) might bias the results. Thirdly, gene expression signatures of
126 interest are often not measured in the same cells/tissues as those used in the CMap. In
127 addition, different treatment durations can lead to different results due to feedback regulation of
128 the target, for example when studying G-protein coupled receptors. Other non-biological
129 phenomenon such as the "batch effect", which affects the microarrays, compounds and cell
130 used, can also impact the accuracy of the predictions [24]. Finally, as mentioned before, the
131 CMap has to be considered a hypothesis-generating tool where results need to be validated by
132 further experimentation. In any case, its potential could be significant and, indeed, similar
133 approaches for connecting drugs and genes are starting to emerge. For example, the tool
134 MANTRA (Mode of Action by NeTwoRk Analysis) allows analysis of the CMap data with an
135 innovative approach that takes into consideration the variability in the transcriptional responses
136 to the drug due to cell-line specific effects, different concentrations of drug applied and distinct
137 experimental conditions [25]. Another example is DvD (Drug versus Disease), a new tool that
138 combines together the data from the CMap, and the public microarray repositories Gene
139 Expression Omnibus and Array Express [26]. In addition to new analytical tools, new powerful
140 technologies such as next generation sequencing (NGS), currently generating data faster than
141 they can be analyzed, might be incorporated and applied to drug discovery and development
142 [27].

143

144 ***Successful translational research: importance of data-sharing and replication***

145 Despite the large number of studies using these powerful high-throughput 'omics' analysis
146 conducted over the last decade, it is striking and concerning the low number of discoveries that
147 have been translated into practice. To improve these odds, it is absolutely fundamental that
148 research discoveries are reproduced and validated in independent studies. A recent analysis of
149 18 microarray studies showed that only 2 were fully reproduced by independent researchers
150 [28]: the main reason for failure was the unavailability of the data necessary to reproduce the
151 published results. Similarly, analysis of the top 50 journals with highest impact factors revealed
152 that only 70% require a mandatory public deposition of microarrays data to guarantee
153 publication. More surprisingly, even if journals were subjected to data availability policies, 59%
154 of the articles analysed did not fully adhere to their requirements [29]. Scientific journals should
155 fully adhere to data-sharing policies to ensure reproducibility as a cornerstone of the scientific
156 process. Because CMap studies are based on a selection of a number of up- and down-
157 regulated genes obtained from previously conducted microarray analyses, the selection criteria
158 and the list of genes used for the analysis should be available to ensure transparency and
159 reproducibility.

160

161 Other publication practices might also be considered, such as the general tendency to publish
162 the more spectacular results, which might be not fully representative of the true 'real-life' result.
163 Journals should allow and promote publication of independent re-analysis and confirmation
164 studies, not only initial evidence, as replication is essential for the consolidation of scientific
165 knowledge and its eventual translation. In addition, underestimation and general refusal of
166 negative data also distorts the real picture [30, 31]. From the bench side, a more accurate
167 communication of microarray data is needed, although this aspect has improved thanks to
168 MIAME (minimum information about a microarray experiment), consisting of a number of
169 recommendations on the information that needs to be provided to enable the unambiguous
170 interpretation of microarray-based experimental results [32].

171

172 ***Challenges and future directions***

173 Despite its slow starting, we truly believe that integration of "omics' into the drug development
174 process and clinical practice will become a reality in future years. Innovative tools and

175 databases promoting the re-use of publicly available information provide new opportunities for
176 drug development at a low cost [33]. Initiatives like the Connectivity Map described here provide
177 publicly available tools to extract useful information from whole-genome studies, often not fully
178 exploited in part due to the difficulty associated to the analysis of large amount of information.
179 Addition of more gene expression signatures representing more drugs and more cell lines, as it
180 will happen with forthcoming CMap versions, would increase its usefulness. Data-sharing
181 policies should be fully implemented and Journals should encourage authors to submit sufficient
182 details to allow independent assessments of their findings. This transparency is of vital
183 importance for the performance of meta-analysis, which might help to overcome the variation
184 between individual studies.

185

186 In conclusion, costs and objective difficulties associated with the drug discovery process require
187 innovative approaches, where the benefits of available information is maximised. In this sense,
188 drug repositioning and identification of new mechanisms represent a low-cost process since
189 making use of already developed therapeutics: these have often been used in humans,
190 therefore facilitating rapid testing in clinical settings and rapid completion of drug repositioning.
191 The CMap can be of great help for this, even more if potentiated with more meaningful protocols
192 (e.g. use of primary cells). On the other hand, an organized multi-disciplinary effort is needed,
193 from basic scientists, clinicians, research journals and regulatory bodies, to make the concept of
194 translational medicine a reality and not a future perspective. An effort by bio-informatics to make
195 these powerful tools easy to use and to interpret by basic scientists (biologists,
196 pharmacologists...) will also be desirable. This must be our priority considering that the ultimate
197 goal of drug development is improvement of the quality of life of patients. And sooner or later,
198 we all will be patients!

199

200

201 **Conflict of interest**

202 The authors declare no conflict of interest.

203

204

205 **References**

- 206 1 Kaitin, K.I. (2010) Deconstructing the drug development process: the new face of
207 innovation. *Clinical pharmacology and therapeutics* 87, 356-361
- 208 2 Kaitin, K.I. (2010) The Landscape for Pharmaceutical Innovation: Drivers of Cost-
209 Effective Clinical Research. *Pharmaceutical outsourcing* 2010
- 210 3 Milne, C.P. and Kaitin, K.I. (2009) Translational medicine: an engine of change for
211 bringing new technology to community health. *Science translational medicine* 1, 5cm5
- 212 4 Loscalzo, J., *et al.* (2007) Human disease classification in the postgenomic era: a
213 complex systems approach to human pathobiology. *Molecular systems biology* 3, 124
- 214 5 Alizadeh, A.A., *et al.* (2000) Distinct types of diffuse large B-cell lymphoma
215 identified by gene expression profiling. *Nature* 403, 503-511
- 216 6 Roses, A.D. (2004) Pharmacogenetics and drug development: the path to safer and
217 more effective drugs. *Nature reviews. Genetics* 5, 645-656
- 218 7 Dudley, J.T., *et al.* (2010) Drug discovery in a multidimensional world: systems,
219 patterns, and networks. *Journal of cardiovascular translational research* 3, 438-447
- 220 8 Al-Shahrour, F., *et al.* (2007) From genes to functional classes in the study of
221 biological systems. *BMC bioinformatics* 8, 114
- 222 9 Joyce, A.R. and Palsson, B.O. (2006) The model organism as a system: integrating
223 'omics' data sets. *Nature reviews. Molecular cell biology* 7, 198-210
- 224 10 Lamb, J., *et al.* (2006) The Connectivity Map: using gene-expression signatures to
225 connect small molecules, genes, and disease. *Science* 313, 1929-1935
- 226 11 Lamb, J. (2007) The Connectivity Map: a new tool for biomedical research. *Nature*
227 *reviews. Cancer* 7, 54-60
- 228 12 Gullans, S.R. (2006) Connecting the dots using gene-expression profiles. *The New*
229 *England journal of medicine* 355, 2042-2044
- 230 13 Lussier, Y.A. and Chen, J.L. (2011) The emergence of genome-based drug
231 repositioning. *Science translational medicine* 3, 96ps35
- 232 14 Sirota, M., *et al.* (2011) Discovery and preclinical validation of drug indications
233 using compendia of public gene expression data. *Science translational medicine* 3,
234 96ra77
- 235 15 Dudley, J.T., *et al.* (2011) Computational repositioning of the anticonvulsant
236 topiramate for inflammatory bowel disease. *Science translational medicine* 3, 96ra76
- 237 16 Claerhout, S., *et al.* (2011) Gene expression signature analysis identifies vorinostat
238 as a candidate therapy for gastric cancer. *PloS one* 6, e24662
- 239 17 Hassane, D.C., *et al.* (2008) Discovery of agents that eradicate leukemia stem cells
240 using an in silico screen of public gene expression data. *Blood* 111, 5654-5662
- 241 18 Zhong, Y., *et al.* (2013) Renoprotective effect of combined inhibition of angiotensin-
242 converting enzyme and histone deacetylase. *Journal of the American Society of*
243 *Nephrology : JASN* 24, 801-811
- 244 19 Zimmer, M., *et al.* (2010) The connectivity map links iron regulatory protein-1-
245 mediated inhibition of hypoxia-inducible factor-2a translation to the anti-inflammatory
246 15-deoxy-delta12,14-prostaglandin J2. *Cancer research* 70, 3071-3079
- 247 20 Hieronymus, H., *et al.* (2006) Gene expression signature-based chemical genomic
248 prediction identifies a novel class of HSP90 pathway modulators. *Cancer cell* 10, 321-
249 330
- 250 21 Wen, Z., *et al.* (2011) Discovery of molecular mechanisms of traditional Chinese
251 medicinal formula Si-Wu-Tang using gene expression microarray and connectivity
252 map. *PloS one* 6, e18278
- 253 22 Renshaw, D., *et al.* (2010) Downstream gene activation of the receptor ALX by the
254 agonist annexin A1. *PloS one* 5

255 23 Montero-Melendez, T., *et al.* (2013) Gene expression signature-based approach
256 identifies a pro-resolving mechanism of action for histone deacetylase inhibitors. *Cell*
257 *death and differentiation* 20, 567-575
258 24 Iskar, M., *et al.* (2010) Drug-induced regulation of target expression. *PLoS*
259 *computational biology* 6
260 25 Iorio, F., *et al.* (2010) Discovery of drug mode of action and drug repositioning from
261 transcriptional responses. *Proceedings of the National Academy of Sciences of the*
262 *United States of America* 107, 14621-14626
263 26 Pacini, C., *et al.* (2013) DvD: An R/Cytoscape pipeline for drug repurposing using
264 public repositories of gene expression data. *Bioinformatics* 29, 132-134
265 27 Woollard, P.M., *et al.* (2011) The application of next-generation sequencing
266 technologies to drug discovery and development. *Drug discovery today* 16, 512-519
267 28 Ioannidis, J.P., *et al.* (2009) Repeatability of published microarray gene expression
268 analyses. *Nature genetics* 41, 149-155
269 29 Alsheikh-Ali, A.A., *et al.* (2011) Public availability of published research data in
270 high-impact journals. *PloS one* 6, e24357
271 30 Young, N.S., *et al.* (2008) Why current publication practices may distort science.
272 *PLoS medicine* 5, e201
273 31 Ioannidis, J.P. (2005) Why most published research findings are false. *PLoS*
274 *medicine* 2, e124
275 32 Brazma, A., *et al.* (2001) Minimum information about a microarray experiment
276 (MIAME)-toward standards for microarray data. *Nature genetics* 29, 365-371
277 33 Rung, J. and Brazma, A. (2013) Reuse of public genome-wide gene expression data.
278 *Nature reviews. Genetics* 14, 89-99
279
280

281

282

283 **Figure Legends**

284 **Figure 1. The Connectivity Map concept.** The Connectivity Map (CMap build 02) is a
285 database that contains the gene expression signatures (obtained with the Affymetrix Genechip
286 HG-U133A) of more than 1,300 bioactive molecules. Differentially expressed genes were
287 identified by comparing cells treated with each distinct drug with untreated cells. A gene
288 expression signature of interest (e.g. of a drug on a particular cell type (A) or a disease (B)) can
289 be compared with those contained in the CMap database. If the signatures compared are
290 similar (that will be identified by a 'positive' score), this could potentially be used to predict novel
291 actions or suggest mechanism of actions of known or novel compounds. On the other hand,
292 comparisons with a disease signature and identification of a 'negative' score (i.e. the gene
293 signatures are the opposite) could be used for drug repositioning studies or to suggest new

294 treatments for that disease. Experimental validation is further required to confirm hypothesis or
295 predictions furnished by the CMap.
296

