

1 Matching cell lines with cancer type and subtype of 2 origin via mutational, epigenomic and transcriptomic 3 patterns

4 Marina Salvadorés¹, Francisco Fuster-Tormo^{1,2}, Fran Supek^{1,3}

5 1 Institute for Research in Biomedicine (IRB Barcelona), The Barcelona Institute of Science and Technology, Barcelona, Spain

6 2 MDS Research Group, Institut de Recerca Contra la Leucèmia Josep Carreras, Institut Català d'Oncologia-Hospital Germans Trias
7 i Pujol, Universitat Autònoma de Barcelona, Badalona, Spain.

8 3 Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Spain

9

10 Abstract

11 Cell lines are commonly used as cancer models. Because the tissue and/or cell type of origin
12 provide important context for understanding mechanisms of cancer, we systematically examined
13 whether cell lines exhibit features matching the cancer type that supposedly originated them. To
14 this end, we aligned the mRNA expression and DNA methylation data between ~9,000 solid
15 tumors and ~600 cell lines to remove the global differences stemming from growth in cell culture.
16 Next, we created classification models for cancer type and subtype using tumor data, and applied
17 them to cell line data. Overall, the transcriptomic and epigenomic classifiers consistently identified
18 35 cell lines which better matched a different tissue or cell type than the one the cell line was
19 originally annotated with; we recommend caution in using these cell lines in experimental work.
20 Six cell lines were identified as originating from the skin, of which five were further corroborated
21 by the presence of a UV-like mutational signature in their genome, strongly suggesting
22 mislabelling. Overall, genomic evidence additionally supports that 22 (3.6% of all considered) cell
23 lines may be mislabelled because we predict they originate from a different tissue/cell type.
24 Finally, we cataloged 366 cell lines in which both transcriptomic and epigenomic profiles strongly
25 resemble the tumor type of origin, designating them as 'golden set' cell lines. We suggest these
26 cell lines are better suited for experimental work that depends on tissue identity and propose
27 tentative assignments to cancer subtypes. Finally, we show that accounting for the uncertain
28 tissue-of-origin labels can change the interpretation of drug sensitivity and CRISPR genetic
29 screening data. In particular, in brain, lung and pancreatic cancer cell lines, many novel
30 determinants of drug sensitivity or resistance emerged by focussing on the cell lines that are best
31 matched to the cancer type of interest.

32

33 **Key words:** cancer cell lines, human tumors, mRNA level, DNA methylation, mutational
34 signatures, cancer type, drug screen, genetic screen.

35

36 Introduction

37 Cell lines are an important research tool, often used in place of primary cells and intact organisms
38 to study biological processes. Cell lines are used for various applications such as testing drug
39 metabolism and cytotoxicity, study of gene function, generation of artificial tissues and synthesis
40 of biological compounds (1). In cancer research, cell lines derived from tumors are commonly
41 used as models because they are presumed to carry the genomic and epigenomic alterations that
42 arise in tumors (2). To understand the response of tumors to therapy, many studies have linked
43 genetic and/or epigenetic alterations with drug response across cell line panels, generating
44 datasets such as the Genomics of Drug Sensitivity in Cancer (GDSC) (3), the Cancer Cell Line
45 Encyclopedia (CCLE) (4), the Cancer Therapeutics Response Portal (5) and others. These efforts
46 have advanced our understanding of tumor biology by generating a massive resource of genomic,
47 transcriptomic, epigenomic and drug response data for hundreds of cell lines (2).

48

49 As a model for cancer, cell lines are cost effective, convenient and amenable to high-throughput
50 screening (1,2). However, a major question associated with the use of cell lines is whether they
51 are representative of the cancer they are meant to model, which may be complicated by issues
52 of misidentification (1,2,6).

53

54 Misidentified cell lines may lead to inconsistent conclusions across studies using the affected cell
55 lines. For instance, the cell lines referred to as HEp-2 and INT 407 in the literature are in fact
56 commonly cross-contaminated with HeLa (cervical cancer) cells, rather than being laryngeal
57 cancer and normal intestinal epithelium cells, respectively (7,8). Because of this, demonstrating
58 cell line identity via genetic markers is now a routine quality-control step. Current resources based
59 on large-scale cancer cell panels are therefore largely unaffected by this issue (4).

60

61 However, even if the identity of the cell line is correct, its properties may not match the cancer
62 type it is meant to model. One way in which this may happen is that tumors thought to originate
63 in a certain tissue might in fact be metastatic lesions originating from a distal site (9). Thus, cell
64 lines derived from such tumors would have a different tissue/cell type identity than that assigned
65 at isolation, constituting a case of mislabeling. It is conceivable that also in the case of primary

66 tumors, ambiguous histological or anatomical features may cause the cancer type or subtype to
67 be misdiagnosed for that tumor and therefore also for a cell line derived from it. Conceivably, the
68 process of establishing the culture might select for a rare cell type that is not representative of the
69 tumor isolate on the whole, meaning that the cell line would again effectively be mislabeled with
70 a different cell type (10). In addition to the initial changes upon adaptation to culture, cell lines
71 evolve over time due to selection and due to genetic drift, potentially diverging from the
72 characteristics of the originating tissue (1).

73
74 Tissue/cell type is a key determinant of response of cultured cells to a variety of experimental
75 conditions, including drug exposure and genetic perturbation (11,12). Therefore, having accurate
76 information on the tissue and cell type identity of a tumor cell line is important for interpreting the
77 experimental results obtained using these cell lines.

78
79 Recent work has examined cell line panels of individual cancer types, showing certain
80 discrepancies between the features of cell lines and corresponding tumor (sub)types. A gene
81 expression analysis of lung tumors and cell lines (10) suggested that some lung adenocarcinoma
82 cell lines did not resemble adenocarcinoma tumors but instead clustered with other lung tumor
83 subtypes (small-cell and squamous cell). A study of high-grade serous ovarian cancer (HGSOC)
84 cell lines using gene expression, driver gene mutations and copy number alteration (CNA) data
85 reported that two frequently used cell lines showed poor genetic similarity to profiles of this ovarian
86 cancer subtype (13). A study of a panel of renal cancer cell lines compared their CNA to kidney
87 tumors, finding that some cell lines used as models of the clear-cell carcinoma more closely
88 resemble papillary renal cancer (14). These examples highlight the need to systematically identify
89 the cell lines whose genotype and/or molecular phenotypes do not resemble the characteristics
90 of the matched human tumor type. A major challenge in the use of human tumor data to classify
91 cell lines are the widespread global changes in gene regulation between cell lines and tumors
92 that arise in cell culture conditions.

93
94 In this study, we performed a global analysis that aligned mRNA expression and DNA methylation
95 data between ~600 cancer cell lines and ~9,000 tumors from 22 different cancer types, adjusting
96 for global differences in transcriptomes and epigenomes. Classifiers trained on human tumor
97 mRNA and DNA methylation profiles were used to systematically identify those cell lines whose
98 genomic and epigenomic profiles are highly consistent with human tumors of their declared
99 cancer type of origin. Conversely, we used the same classifiers to identify those cell lines that

100 might be mislabeled with respect to cancer type or that might have diverged from their original
101 tissue and/or cell type identity. Our data suggests that tens of cell lines might be epigenetically
102 and/or genetically not consistent with their stated tissue or cell type of origin, which is an important
103 consideration for experiments that use these cell lines. We demonstrate this by reanalyzing
104 associations between drug sensitivity and genetic variation in a large panel of cell lines. After
105 explicitly accounting for putative cases of cell lines with mislabeled tissue identity, many novel
106 associations of genes with drug sensitivity or resistance were revealed.

107

108 **Results**

109 **1. Identification of tissue/cell type-of-origin for cell lines by a joint analysis with tumors**

110 During adaptation to cell culture, certain changes in the cell lines' physiology are inevitable, yet
111 ideally the cell lines should retain sufficient features of the tumor to be useful as experimental
112 models of tumor biology. Here, we systematically examined the global features of the
113 transcriptome and epigenome that reflect the tissue-of-origin of a tumor cell line. The tissue that
114 originated a tumor is well known to be a major determinant of drug responses -- including drugs
115 targeted to certain genetic mutations -- both *in vitro* (11,12) and also *in vivo* (15,16). Tissue of
116 origin is an important factor in shaping the networks of genetic interactions in cancer (17) and
117 also determines the phenotypes resulting from genetic perturbation (18). Therefore ascertaining
118 the tissue/cell type identity of cell lines is relevant for interpreting results of various experiments.

119

120 During the process of adaptation to cell culture, the cells undergo global changes in gene
121 regulation that affect many genes (19,20). In particular, the global patterns in transcriptomes and
122 epigenomes for cultured cells bear many similarities to other cultured cells, irrespective of the
123 originating tissue. Thus, there are commonalities in how culture affects gene regulation: for
124 example, proliferation genes in cultured cells have distinct DNA methylation and gene expression
125 patterns, when compared to tumor and normal tissues (19,20). These global alterations in gene
126 expression and DNA methylation mean that is not straightforward to directly compare cell line
127 transcriptomes/epigenomes with data obtained from actual tumors. Therefore, the cell culture-
128 induced shifts need to be carefully adjusted for in order to be able to track down tissue identity of
129 cell lines. To this end, we introduce a computational framework -- HyperTracker -- which can unify
130 transcriptome, epigenome and mutational data across tumors and cell lines, and provide robust
131 predictions of tissue/cell type and subtype identity.

132

133 In particular, we collected gene expression (RNA-Seq) and DNA methylation data (microarrays)
134 for 9,681 and 9,039 human tumors, respectively (TCGA), and additionally for 614 cell lines (CL)
135 of various solid cancer types. For gene expression data (GE), we examined transcript-per-million
136 (TPM) normalized counts for the 12,419 genes where RNA-Seq data was available for both cell
137 lines and tumors. For DNA methylation data (MET), we examined beta-values for 10,141 probes
138 from methylation arrays, after selecting a single probe per gene promoter with the highest
139 variance across the dataset. To align human tumor and cell line data, we quantile-normalized the
140 data and applied ComBat, a batch effect correction method (21), which is highly performant
141 compared to other related methods (22). In brief, ComBat estimates parameters for location and
142 scale adjustment of each batch (TCGA and CL in our case) for each gene. Then, it removes the
143 variability which is particular to the CL but not present in TCGA, while retaining the intra-dataset
144 variability of the tumors, which should presumably be evident in both the tumor and in the cell line
145 datasets.

146

147 A principal component analysis in the data (pre- and post-adjustment) suggests that there were
148 indeed strong global differences between TCGA and CL, and that they are largely removed by
149 our approach (Fig 1a; Fig S1ab). To quantify this, we trained a classification model that predicts
150 the CL *versus* TCGA origin of the data points based on GE and MET (Fig 1b). The model is able
151 to distinguish CL *versus* TCGA perfectly when using the pre-adjustment datasets (AUC=1), while
152 the post-adjustment datasets (AUC(GE) = 0.44; AUC(MET) = 0.42) do not perform better than
153 random (0.5; Fig 1b), suggesting the cell-type specific signal has been largely removed. Finally,
154 we tested the optimal number of features (genes/probes) using tumor classifiers and calculating
155 the accuracy in the cell line data (Fig S1c); we selected 5,000 features with the highest standard
156 deviation for later analyses.

157

158 Once the data was aligned, we set out to determine which cell lines have tissue identity not
159 matching the declared tissue-of-origin (henceforth: 'suspect set'), and conversely, which cell lines
160 have largely retained their tissue identity (henceforth: 'golden set'), by comparing against a large
161 set of tumors from 17 tissues in the TCGA (Fig 1c). Using TCGA data, we derived one-*versus*-
162 rest classification models (using Ridge regression), separately for the GE and the MET data. Some
163 pairs of cancer types were considered jointly in this analysis, based on their overall similarity, for
164 example stomach adenocarcinoma (TCGA code: STAD) and esophageal adenocarcinoma
165 (subset of samples from TCGA code: ESAD); see Methods for a full list. Our study focuses on

16

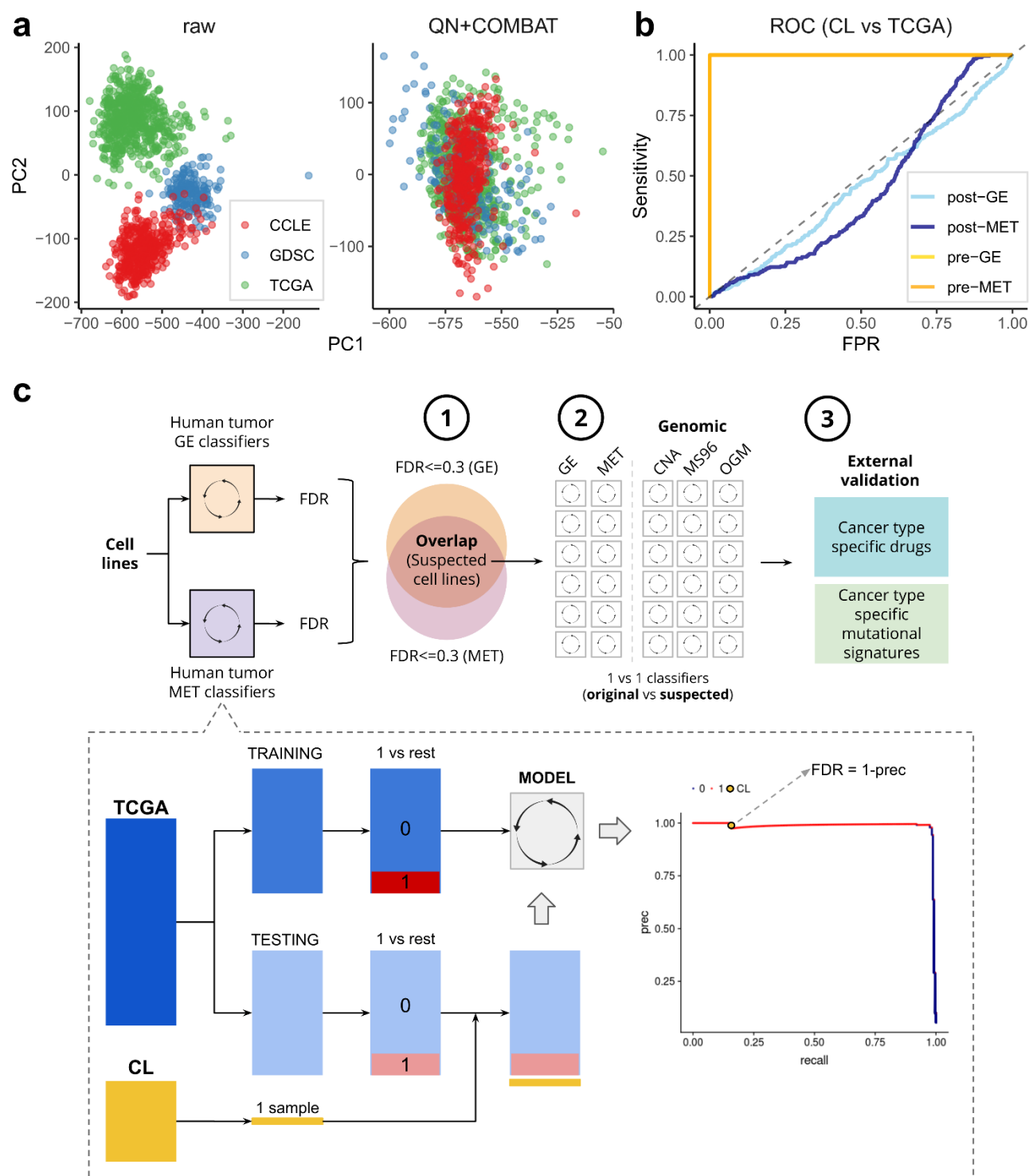


Fig. 1. Data alignment and methodology for classification. (a) Principal component (PC) 1 and PC2 of a PC analysis, in the gene expression (GE) data pre-adjustment for batch effects (raw) and post-adjustment (quantile normalization+COMBAT) (see Fig S1 for DNA methylation data (MET)). Colors represent the dataset sources (GDSC and CCLE are two sources for the cell lines, and TCGA is the source for the tumors). (b) ROC curves for classifying TCGA versus cell lines in the data pre-adjustment (orange) and post-adjustment (blue) for GE and MET. (c) An overview of the HyperTracker methodology applied in the manuscript. First, we systematically identified possible mislabeled cell lines using GE and MET data, independantly. Second, we used various types genomic data to corroborate the hits. Third, we further validate the cell lines suspected to originate from skin using independent data.

167 solid cancer types and does not examine blood cancers. In a crossvalidation test, TCGA models
168 had very high AUPRC scores: 0.98 and 0.97 for GE and MET respectively (average across cancer
169 types). This means that transcript level data and DNA methylation data are largely sufficient to
170 accurately distinguish those cancer types.

171
172 Next, we obtained predictions of cancer type identity for each cell line. For every cancer type, we
173 split TCGA data randomly into training and testing sets, and we used the calculated precision-
174 recall curve of the testing data to obtain the False Discovery Rate (FDR) score for every cell line
175 (details in Methods; all FDR values are listed in [Table S1](#)). The smaller the FDR, the more likely
176 the cell line is to belong to that particular cancer type. As expected, most of the cancer type labels
177 of the cell lines match the declared tissue of origin of that cell line -- they tend to cluster at low
178 FDR values for the cognate cancer type (red dots in [Fig 2a](#), [Fig S2](#)). However, among these many
179 correctly classified cell lines (red dots), there are some with similarly low FDR scores, but which
180 were originally annotated as belonging to another cancer type ([Fig 2a](#); blue dots with label shown).
181 A clustering analysis of the GE and MET values for the genes that had the highest weight in the
182 classification models ([Fig 2b](#), [Fig S3](#)) showed that in most cases, the samples clearly cluster by
183 cancer type, but not by CL *versus* TCGA label. Moreover, we observed that the suspected cell
184 lines (cell lines with highly confident FDR scores to a different cancer type) tend to cluster with
185 the newly-assigned cancer type by the classifier, rather than with the original one ([Fig 2b](#)).

186
187 In further analyses, we designated as the 'golden set' those cell lines that have $FDR \leq 0.3$ for
188 both GE and, independently, for MET in their originally declared cancer type ($n=366$ out of 614
189 examined cell lines, 60%). For these cell lines, two independent types of evidence --
190 transcriptomes and epigenomes -- support that they match their expected cancer type well,
191 suggesting these cell lines would be preferred as experimental models. Further, we designated
192 as the 'silver set' those cell those cell lines that have $FDR \leq 0.3$ for only one classifier (either
193 GE or MET but not both) ($n=131$ out of 614 examined cell lines, 21%). From the remaining 117
194 cell lines, we selected as 'suspect set' those CL which exhibit an $FDR \leq 30\%$ for both GE and
195 for MET, but in a different cancer type than declared for that cell line ($n=43$ out of 614, 7% of
196 analyzed cell lines) ([Fig 1c](#)). This set of cell lines may consist either of mislabeled cell lines, where
197 the cancer type of origin is different than it was thought, or of heavily diverged cell lines, where
198 the genomic and/or epigenomic alterations accumulating during cell culture have overridden the
199 original cancer type identity. Of note, cell line cross-contamination issues (23) cannot underlie the
200 trends we observe, because the repositories that provided GE and MET data have used genetic

201

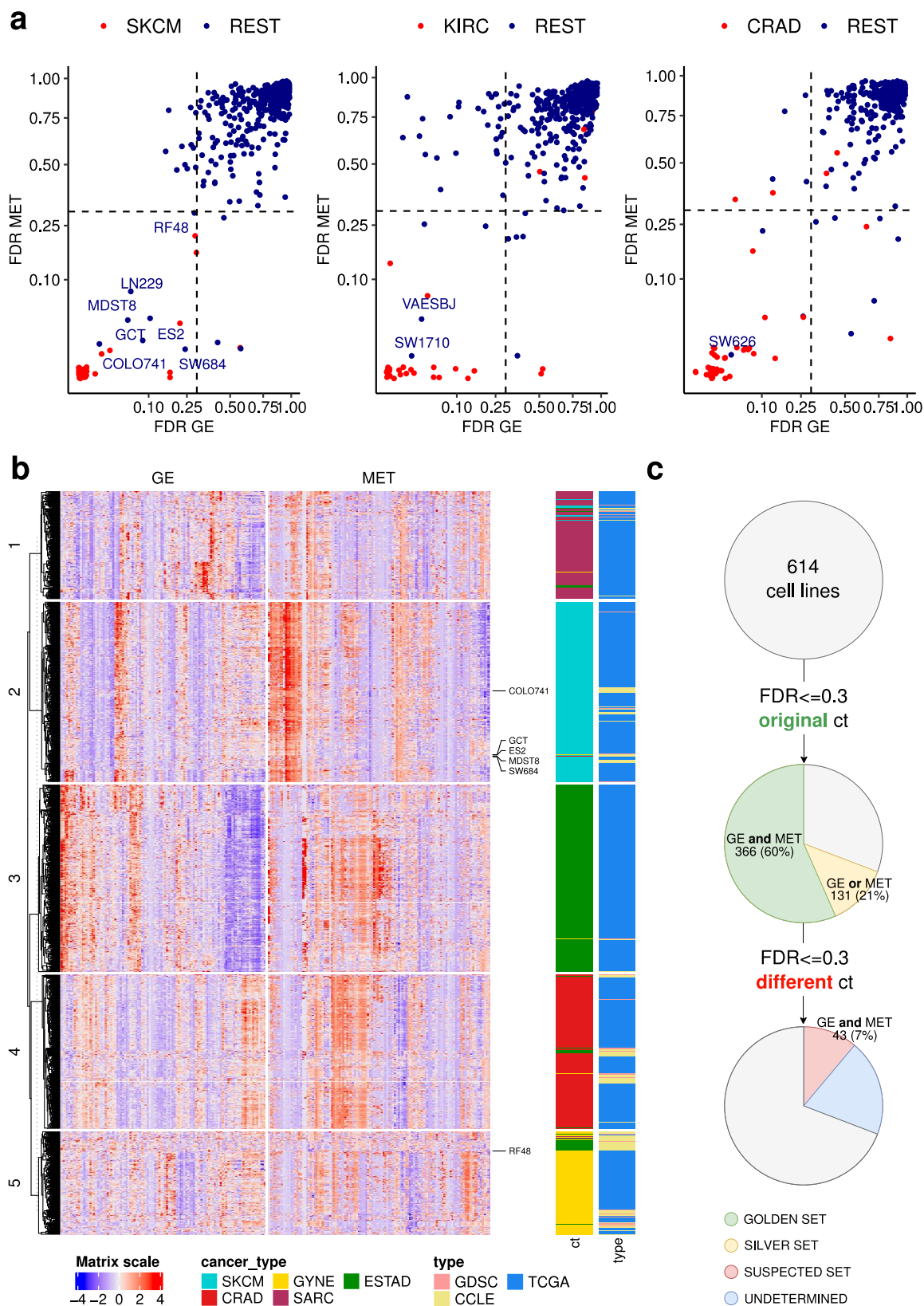


Fig. 2. Detection of cell lines mislabelled to a different cancer type. (a) False Discovery Rate (FDR) scores for 614 cell lines were calculated in MET and GE cancer type classifiers (one-versus-rest). The lower the FDR, the higher the confidence that the sample belongs to that particular cancer type (here, to SKCM, KIRC and CRAD from left to right, see Fig S2 for the other cancer types). The cell lines that were originally annotated as the cancer type that is being tested are shown in red, the rest in blue. (b) Heatmap for the 25 genes (GE) and CpG probes (MET) whose Ridge regression coefficients had the highest absolute values for SKCM (skin cancer) versus rest classifiers. The suspected skin cell lines are labeled in the right side of the heatmap. The cancer types shown are the suspected one (SKCM in this case) and additionally the originally declared cancer types of the suspected cell lines (here, ESTAD, SARC, CRAD and GYNE). See Fig S3 for the Heatmaps for the rest of the suspected cell lines. (c) Overview of the results from the systematic mislabelling testing of all cell lines. Cell lines with an $FDR \leq 0.3$ to its original cancer type in (i) GE and MET are assigned to the 'golden set' group and (ii) either GE or MET are assigned to the 'silver set'. If however the $FDR \leq 0.3$ to a different cancer type in GE and in MET, the cell line is assigned to the suspected set.

202 markers to ascertain the identity of the cell lines (4). The fact that two classifiers based on
203 independent data types -- one genomic and one epigenomic -- reached the same predictions
204 adds confidence that these are *bona fide* cases of mistaken tissue/cell-type identity.

205

206 **2. Validation of individual examples of suspected mislabeled cell lines using genomic** 207 **classifiers**

208 We detected 43 cell lines that bear transcriptomic and also epigenomic features of a different
209 cancer type than the one they were originally annotated to. We next turned to support individual
210 examples of cell lines with reassigned tissue identity by analyzing independent data. In particular,
211 we used genomic sequence-based classifiers, which are able to predict the tissue of origin based
212 on mutation patterns (24,25). As in our recent work (24), we used the trinucleotide mutation
213 spectra and the oncogenic mutations. In this validation setting, we applied such genomic
214 classifiers to a problem of 'one-versus-one' classification, where we contrasted the originally
215 assigned cancer type *versus* the newly-proposed cancer type for each reassignment. We found
216 that such one-versus-one classifiers based on genomic data had satisfactory accuracy with our
217 whole-exome sequencing data sets (Fig S4; our past work (24) suggests whole genome
218 sequences are more powerful). Finally, we included an additional classifier based on copy number
219 alteration (CNA) profiles, which were also shown to yield accurate predictive models of tissue
220 specificity (24,25).

221

222 For the 43 examples of suspected cell lines tissue identity, we first derived one-*versus*-one
223 classification models separately for GE and MET. If a cell line is truly mislabelled when testing
224 the original *versus* the suspected cancer type, we should observe the same reassignment of the
225 cell line to be robustly observed across multiple runs of the classification algorithm, which use
226 different random initializations. Out of 20 iterations of the algorithm, a score of 20 indicates that
227 the cell line is consistently predicted as the suspected cancer type, and a score of 0 means that
228 the cell line is consistently assigned to the original cancer type. We randomized the labels to
229 obtain a background model of expected values (Fig 3b; Fig S5a). From the 43 suspected cell
230 lines, 35 are consistently reassigned to the other tissue (score>10), irrespective of the variability
231 in the predictive models introduced by resampling the data (Fig 3a; Fig S5b). Next, we calculated
232 the same score for the genomic classifiers (based on mutations and CNA, as described above)
233 on these 35 suspected cell lines (Fig 3a).

234

235 Of these, approximately two-thirds ($n=22$ cell lines) received high support for the new tissue label
236 by one or more genomic classifiers (Fig 3a; score ≥ 15 , corresponding to FDRs of 0%, 0% and
237 18% for the CNA, OGM and MS96 respectively, based on randomized data; Fig 3b). This data
238 suggests 22 cell lines are candidates for assignment to another cancer type, based on converging
239 evidence from the levels of the genome, epigenome and transcriptome, which provides
240 confidence. Reassuringly, this list contains two cell lines which have been previously shown to be
241 misclassified: SW626 which was initially annotated as ovarian cancer but later discovered to be
242 derived from colon cancer (26), and COLO741 which was originally thought to be a colon
243 adenocarcinoma cell line but later shown to originate from a melanoma (27). The fact that these
244 two known examples were detected and reassigned to the correct cancer type provides evidence
245 that our method is overall reliable.

246
247 The two plausible reasons why a cell line thought to originate from one cell type would need to be
248 reassigned to a different cell type are (i) that at the time of isolation, the cell line was not of the
249 type that it was thought to be (mislabeling), or (ii) that during prolonged cell culture, the cell line
250 diverged greatly and now resembles another cell type (transdifferentiation). Our data allows to
251 examine how prevalent each case is: mislabelling is expected to be reflected equally in both the
252 epigenome and the genome, while transdifferentiation is expected to be reflected more strongly
253 in the (presumably more malleable) epigenome, and less so in the genome, which retains the
254 mutations from the original tumor. We suggest that mislabelling at isolation is a much more
255 common scenario (Fig 3c, many reassigned cell lines are in the upper-right corner). However, it
256 is possible that there exist individual examples of cell lines that have effectively transdifferentiated
257 during culture, because their genomic features are consistent with the original tissue identity while
258 the epigenomic features are consistent with another tissue (Fig 3c, lower left corner, e.g. the
259 RPMI2650 and OACM51 cell lines are possible candidates).

260

261 **4. Validation of cell lines suspected to originate from the skin**

262 From the previous analysis, we identified a total of six cell lines which are reassigned from various
263 cancer types to skin cancer. We note that, of skin cancers, the TCGA study contains only
264 melanoma but not the non-melanoma skin cancers, so we are currently not able to distinguish
265 between cell type identities of different types of skin cancer.

266

267 To further support that these cells are indeed skin cancer cells, we performed an independent
268 analysis based on mutational signatures to confirm the mislabelling. Large-scale analyses of
269 trinucleotide mutation spectra across human tumors have revealed at least 30 different types of
270 mutational signatures (28). Of these, Signature 7 (C>T changes in CC and TC contexts) was
271 associated with exposure to UV light and is highly abundant in sun-exposed melanoma tumors
272 (29). The same signatures were recently estimated in cancer cell lines by two related methods
273 (30,31), which enabled us to use existence UV-linked Signature 7 to examine whether these cell
274 lines originated from the skin. Based on mutational burden of Signature 7, the known melanoma
275 cell lines (turquoise dots) are clearly separated from the rest (Fig 3e), meaning the approach can
276 distinguish skin-derived cells. Among the melanoma cell lines with high mutational burden of
277 Signature 7, we found four out of five of the suspected cell lines (Fig 3e), in particular GCT,
278 SW684, ES2 and MDST8 are very likely skin cells, and not sarcoma, sarcoma, ovarian cancer or
279 colorectal cancer, respectively, as originally thought. For the sixth suspected cell line COLO741,
280 the mutational signature data is not available, however COLO741 has been previously reported
281 of being melanoma based on the expression of skin-specific genes (27).

282
283 The RF48 cell line (originally considered stomach, here putatively reassigned to skin) does not
284 exhibit the UV signature nor the DNA methylation patterns of skin, therefore a highly confident
285 call cannot be made. Nonetheless, a pattern of cancer driver mutations in RF48 suggests it is
286 indeed not a stomach cell line (Fig 3a). Past work based on gene expression suggested that RF48
287 is indeed not representative of stomach -- instead, a lymphoid origin was proposed for RF48 (32).

288
289 Next, we sought to substantiate these findings using drug sensitivity data. In particular, two drugs
290 (dabrafenib and trametinib) that target mutant BRAF are approved for treating melanoma in the
291 clinic. These drugs are known to have poor efficacy in other cancer types bearing BRAF
292 mutations, such as in colon cancers (33) and therefore sensitivity to these drugs adds confidence
293 we are in fact looking at a melanoma cell line; (note that the converse does not necessarily hold
294 here: resistance does not imply it is not a melanoma). Therefore, we compared the IC50 of these
295 two drugs for all cell lines (Fig 3d). As expected, many melanoma cell lines cluster at low values
296 of IC50 for the two drugs, meaning these cells are sensitive to the drugs. Among this cluster we
297 observed two out of five of our suspected cell lines (ES2 and MDST8) providing further supporting
298 evidence these are of skin, likely melanoma skin cancer origin.

299
300

301

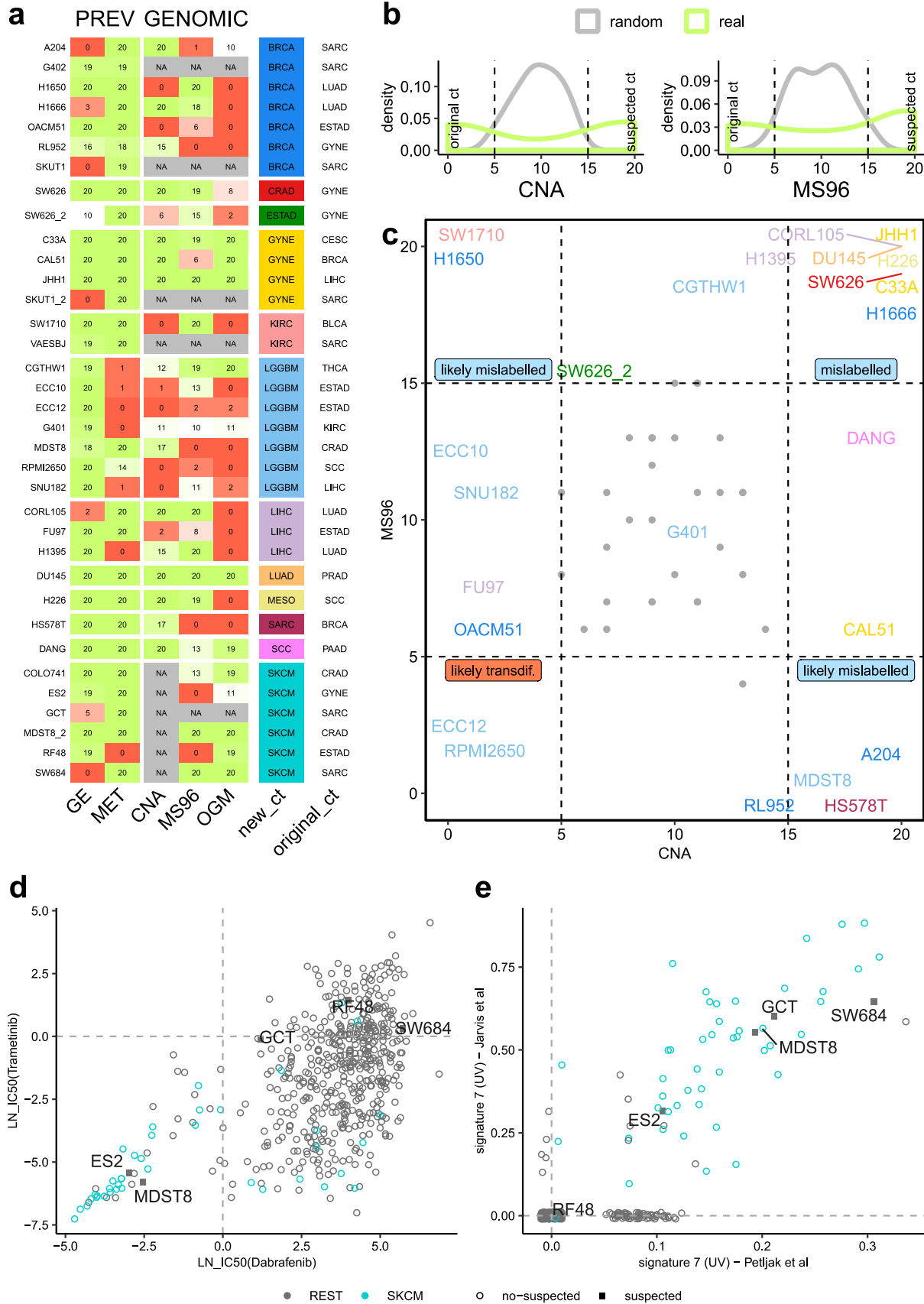


Fig. 3. Further evidence supporting tissue identity of the suspected cell lines. (a) Prediction score (0-20) for each suspected cell line for 20 runs of one-versus-one classifiers that predicted suspected versus original cancer type in GE, MET, copy number alteration (CNA), mutational spectrum (MS96) and oncogenic mutations (OGM). A value of 20 means that the cell line is predicted as suspected consistently in the 20 runs of the algorithm, and a value of 0 means it is predicted as original cancer type the 20 runs. **(b)** Histograms of the prediction scores for CNA and MS96 for the models based on actual data, and a baseline on randomized data (shuffling the labels). **(c)** Prediction scores for MS96 and CNA for the suspected cell lines. Colors represent the suspected cancer type (see column "new_ct" in panel a). Grey dots represent the random values. **(d)** Drug sensitivity (IC50) for mutant BRAF inhibitors dabrafenib and trametinib for 614 cell lines. Cell lines originally labelled as skin cancer shown in red, and skin-suspected cell lines are marked with a square and their sample id. **(e)** Burden of UV-associated mutation Signature 7 (estimated from two different sources) in 614 cell lines. Cell lines originally labelled as skin cancer are shown in red and skin-suspected cell lines are marked with a square and the sample id.

302 In conclusion, from the six cell lines suspected of originating from skin, four of them are confirmed
303 by the UV mutational signatures and two of those are additionally confirmed by the drug sensitivity
304 to BRAF inhibitors. This striking example demonstrates how the transcriptome and epigenome-
305 based tissue/cell type classifiers are able to link cultured human cell lines with their correct cancer
306 type of origin.

307
308 In addition to these examples of skin cell lines, we have further supported several other cancer
309 type reassignments using drug sensitivity data (34) (results summarized in [Table S2](#)). This
310 provided evidence that the DANG cell line is consistent with squamous cell carcinoma of the lung
311 or of the head and neck (SCC), rather than with its original assignment of pancreatic
312 adenocarcinoma (this reassignment is also observed with multiple genomic classifiers; ([Fig 3a](#)).
313 Similarly, SW1710 may be a kidney, rather than a bladder cell line, based on the original
314 reassignment via transcriptome and epigenome, based on mutational patterns ([Fig 3a](#)) and
315 additionally supported in the global analysis of drug responses ([Table S2](#)). We note that such
316 analyses of drug screening data can be applied to distinguish only certain pairs of tissues and not
317 all reassignments can be reliably validated in this test (see AUC scores in [Table S2](#)).

318

319 **5. Identification of subtypes for cell lines using multi-omics analyses**

320 Tumors are heterogeneous and major differences exist between tumor samples of the same
321 cancer type. To manage this variability, researchers have attempted to subdivide each cancer
322 type based on their molecular characteristics, including global patterns in gene expression and
323 DNA methylation (35–37). However, with the exception of a few tumor types, in particular breast
324 cancer, molecular subtypes are still being established or refined, in order to better predict disease
325 progression in response to particular treatments.

326 Since drug screens and genetic screens performed in cell lines have the intent of serving as
327 models for actual tumors, it is important to establish a method that can transfer the subtype
328 assignments from tumors to cell lines, thereby establishing which cell line(s) are the most
329 appropriate model for which cancer subtype.

330 Previously, molecular subtypes from tumors have been transferred to cell lines using different
331 strategies. For breast cancer, cell lines subtypes have been assigned mainly using Prediction
332 Analysis for Microarrays (PAM) analysis, which is based on a restricted set of gene expression
333 markers (38). For colorectal cancer, the cell lines were stratified into the consensus molecular

334 subtypes (CMS) integrating transcriptomic and genomic data (39). For renal cancer, subtypes
335 were assigned to the cell lines using gene expression data (14). In a recent pan-cancer study,
336 subtypes have been assigned to a set of 600 cell lines (40). It has been proposed that the cell
337 lines do not usually represent all subtypes of a particular cancer type, possibly due to a bias
338 introduced during the process of immortalization (38,40).

339 Our approach to assign subtypes to cell lines herein is to apply the same strategies that have
340 allowed us to get accurate cancer type classifiers: first, the integration of transcriptomic and
341 epigenomic data to boost confidence in the predictions, and second, careful adjustment of the
342 two data types to make them comparable between TCGA tumors and cell lines (Fig S1).

343 An important consideration in the task of inferring the cell lines' subtypes is the absence of true
344 labels needed for systematic validation, thus assignments should be treated as tentative.
345 However, for breast cancer cell lines the subtype labels are available (38) and can be used as a
346 benchmark of our multi-omics based methodology.

347 We examined proposed subtypes for 15 cancer types in TCGA and generated subtype classifiers
348 (Methods) for each cancer type. The combination of both data types (GE and MET) achieved a
349 higher cross-validation accuracy in the TCGA (median AUPRC across cancer types: 0.81) than
350 GE (0.76) or MET (0.72) separately. Therefore, we used the combined datasets to generate
351 subtype classifiers and propose assignments of the cell lines to cancer subtypes. Since we are
352 using one-versus-rest classifiers each cell line can be assigned to more than one subtype.
353 However, the majority of them are assigned to only one subtype (Fig S6a); we used only those in
354 further analysis. As a benchmark, we calculated the accuracy for the breast cancer cell lines with
355 subtypes available (Fig S6b): the median AUPRC (across breast cancer subtypes) for CL is 0.83.
356 This suggests acceptable performance in obtaining tentative subtype assignments for cell lines in
357 all 15 cancer types, which we provide as a resource in Table S3. This resource is complementary
358 to a recent set of subtype predictions for 9 cancer types based on transcriptomes (40).

359 Next, we examined if the relative prevalence of subtypes is similar between tumors and cell line
360 panels of the same cancer type. Cell line panels of some cancer types have good representation
361 of subtypes, for instance lung squamous cell cancer, head and neck squamous cell cancer, lung
362 adenocarcinoma, and gastric/esophageal cancers (Fig S6c). However, the converse is the case
363 for liver, skin and thyroid cancer cell lines, in which a single subtype predominates in cell line
364 panels but not in tumors (statistics listed in Supplementary Table 4). Additionally, we observe
365 suboptimal representation (where half of the tumor subtypes are not represented) in the kidney,

366 bladder and brain cancer cell line panels, when considering the 463 cell lines we analyzed. This
367 suggests that -- in some cancer types more than others -- the commonly used cell line panels do
368 not represent the diversity of molecular subtypes in tumors, which should be taken into account
369 when interpreting experimental data. One possible reason is the relative ease of culturing certain
370 subtypes, compared to others (2).

371

372 **6. Accounting for mislabelled cell lines reveals new associations in drug screening data**

373 We detected 35 cell lines that may have a tissue or cell type identity different than the one
374 originally assigned to them. Because the cell type is an important determinant of drug response
375 in cancer cell lines and in tumors (11), we hypothesized that the inclusion of this new tissue
376 information into analyses of genetic determinants of drug sensitivity may change the results. In a
377 comprehensive study, Iorio *et al.* searched for associations between drug response and Cancer
378 Functional Events (CFEs): the recurrent mutations, CNA and hypermethylation events present in
379 human tumors (11). Here, we used GDCSTools (27) to reproduce the results of that study,
380 however after filtering the cell lines to those that better represent the cancer type in question. In
381 particular, we repeated the same analysis using for each tissue (i) all the cell lines; (ii) only the
382 cell lines in the 'golden set' (G); (iii) as a less stringent filtering criterion, only the cell lines in the
383 'golden and silver set' (G&S). Additionally, as controls we included a random subset of cell lines
384 that matches (iv) the number of cell lines in 'golden set' (r_G) and (v) the number in 'golden and
385 silver set' combined (r_G&S).

386

387 For the majority of the cancer types, we observed that one of the filtered subsets recovered a
388 higher number of significant (at $FDR \leq 25\%$) associations of CFE with drug sensitivity or
389 resistance, than were recovered using all cell lines (Fig S7). For instance, for glioblastoma, using
390 the 'golden set' cell lines we found 23 new associations, which were not recovered from the entire
391 cell line panel nor from the random-subset controls (Fig 4b). For example, this recovers the
392 positive association of CDKN2A loss with camptothecin sensitivity (Fig 4c), which was previously
393 reported in an independent analysis of the NCI-60 cell line panel screening data (41). Similarly,
394 for pancreatic adenocarcinoma, benefits were observed by focusing on cell lines that resemble
395 the corresponding cancer type better: using only the 'golden set' plus 'silver set' cell lines, 10 new
396 significant associations were found (Fig 4b). For instance, we detected that SMAD4-mutant cell
397 lines are more resistant to piperlongumine, a natural product claimed to have antitumor properties

398

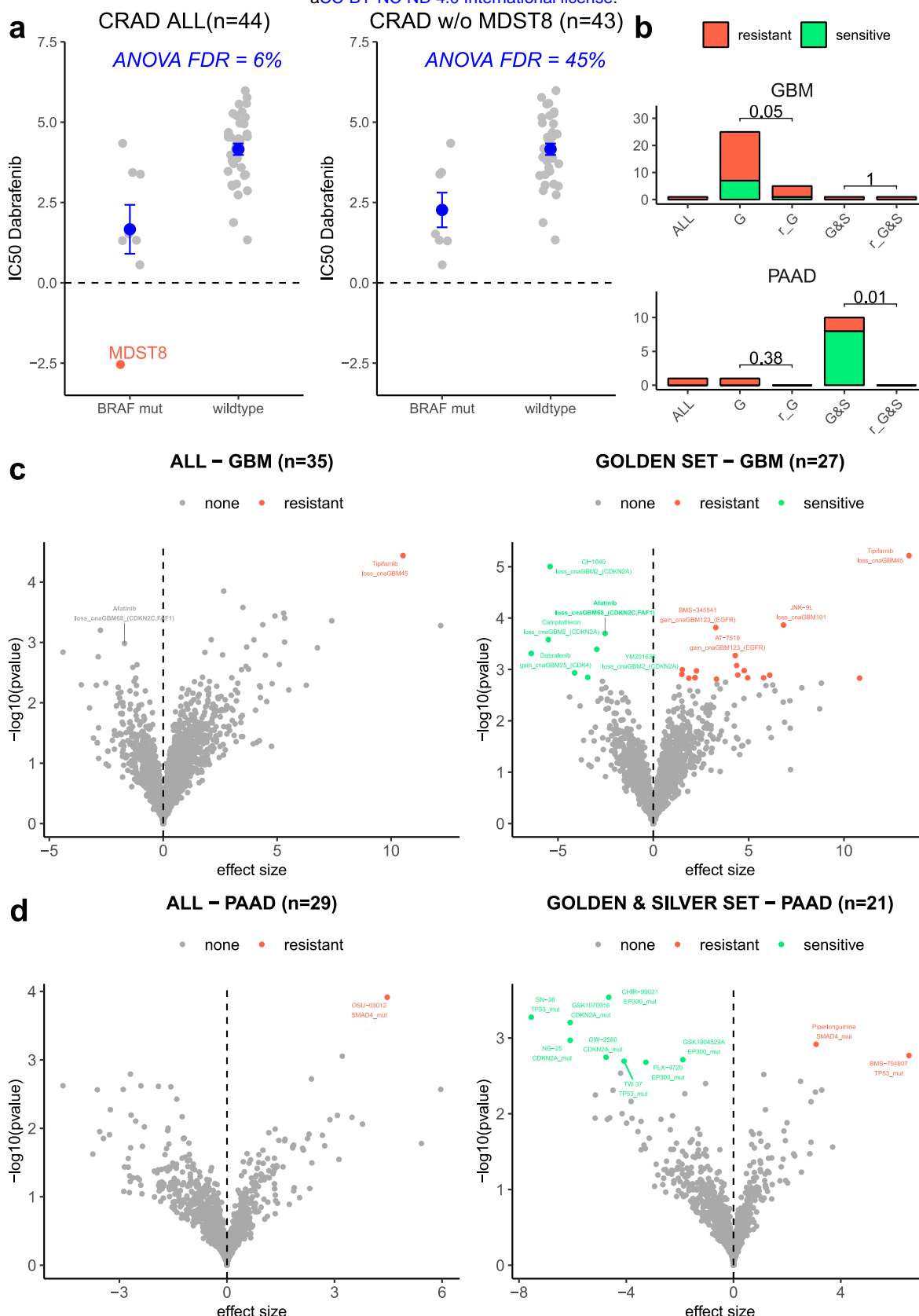


Fig. 4. Drug sensitivity association testing using high-confidence sets of cell lines. (a) Drug sensitivity (IC₅₀) to dabrafenib in all CRAD cell lines (left) and all CRAD cell lines except MDST8, which is suspected of being skin cancer (right). Two groups are compared: cell lines with BRAF mutation and without (wild-type). ANOVA FDR for this association (dabrafenib and BRAF mutation) shown in blue for both datasets. Horizontal line is shown at 0, because score <0 implies sensitivity to the drug. (b) Number of significant associations between Cancer Functional Events (CFEs) and drugs detected in the ANOVA test for all cell lines (ALL), cell lines in the golden set (G), cell lines in the golden plus silver set (G&S), random subset of cell lines that match the number of cell lines in the golden set (r_G) and in the golden plus silver set (r_G&S). For the random subsets, the number of significant associations is calculated 10 times (with different random selection) and the median of the 10 runs is shown. P-values for a sign test (one-tailed, alternative = “less”) between the number of associations in the G/G&S versus the number of associations in r_G/r_G&S are shown. See Fig S7 for remaining cancer types. (c) Differential sensitivity of drugs were analysed by ANOVA for all brain cancer cell lines (left) and the brain cancer cell lines in the golden set only (right). Each point is an association between the sensitivity of a drug and a genetic feature (CFE). (d) Differential sensitivity of drugs were analysed by ANOVA for all pancreatic (PAAD) cell lines (left) and PAAD cell lines in the golden and silver set only (right). Each point is an association between the sensitivity of a drug and a genetic feature (CFE).

399 exerted via multiple pathways (42,43). Mutations of the tumor suppressor gene EP300 were
400 associated with higher sensitivity to three drugs in pancreatic cancer cell lines (Fig 4d).

401
402 The observation that more associations were found despite using a somewhat lower number of
403 cell lines (thus less statistical power) emphasizes the importance of using the cell lines that more
404 closely model the tissue and/or cell type-of-origin of the cognate tumor.

405
406 Colorectal cancers provides an illustrative example of how important is to remove
407 nonrepresentative cell lines from drug screening efforts. In the Iorio *et al.* (2016) study, 50
408 colorectal adenocarcinoma (CRAD) cell lines were tested. Of those, we strongly suspected that
409 MDST8 derives from skin. To test the influence of this individual mislabelled cell line we have
410 performed association testing with all CRAD cell lines, and after excluding MDST8. For the
411 association of the drug dabrafenib with BRAF mutation status, we observed that all CRAD cell
412 lines (irrespective of BRAF mutation) are not sensitive, except for MDST8 which is strongly
413 sensitive (Fig 4a). The FDR of the ANOVA analysis when using all cell lines is 6%, while when
414 removing MDST8 the FDR is 45%. Therefore, in this case, the presence of a single mislabelled
415 cell line is sufficient to cause the appearance of a false association between a drug and a feature.
416 This is fully consistent with clinical responses: in contrast to the good response of patients with
417 BRAF-mutant melanoma to dabrafenib, colorectal tumors with the same BRAF V600 mutation are
418 not sensitive to BRAF or MEK inhibitor monotherapy (33).

419
420 **7. Accounting for mislabelled cell lines reveals new associations in genetic screening data**

421 Motivated by the many novel associations revealed by reanalyzing the drug screening data, we
422 asked if the same extends to genetic screening data in cancer cell lines, because results in genetic
423 screens may also depend on cell lineage (12). To further investigate, we analyzed CRISPR
424 screening data from Project Score and Project Achilles (see Methods), from which 347 cell lines
425 overlap our tested cell lines. Then, we applied the same association testing method, which was
426 however underpowered because the number of available overlapping cell lines was smaller.
427 Nonetheless, in colorectal and ovarian cancer, we observed that focussing only on the ‘golden
428 set’ and/or ‘silver set’, the number of associations recovered increased (as a control, there were
429 no increases in the random cell line subsets of the same size; Fig 5a, Fig S8).

430

431 To illustrate the importance of removing suspect cell lines in gene dependency screenings, we
432 provide two examples of associations that were originally not detected as significant due to the
433 presence of a mislabelled cell line. For ovarian cancer, the presence of SW626 (mislabelled cell
434 line confirmed by the literature (26)) prevents finding the association between MED8 dependency
435 and a copy number gain in the region containing ASXL1 (cnaOV72) as significant (Fig 5b).
436 Similarly, for colorectal cancer the presence of MDST8 (mislabelled cell line confirmed by the UV
437 mutational signature) prevents finding the association between TUBB4B dependency and a copy
438 number gain in STK4 (cnaCOREAD32) (Fig 5c). Finally, a significant association between WRN
439 dependency and MLL2 (also known as KMT2D) gene mutation is recovered only with the filtered
440 cell lines in ovarian cancer (Fig 5d). This WRN-MLL2 association has been recently reported
441 using a somewhat different set of cell lines (from Project Score) (44) that partially overlap our set.

442
443 Finally, our re-analyses of drug screening and genetic screening data revealed an interesting
444 association independently supported in both drug and genetic data. The drug afatinib inhibits the
445 EGFR protein and is clinically indicated for EGFR-mutated lung cancer, however in EGFR-altered
446 glioblastoma afatinib is generally not considered to elicit a response (45). Consistently, afatinib
447 sensitivity was associated with EGFR alterations in lung cancer previously (11), as well as in our
448 re-analysis (FDR G&S = 0.6%), but not in the brain cell line panel (all FDR \geq 25%). However,
449 using the focussed (golden set) of brain cancer cell lines revealed a significant association
450 (ANOVA FDR = 15%, Fig S9a) between afatinib sensitivity and a different genetic alteration: copy
451 number loss in a region at 1p32.3 containing the CDKN2C and FAF1 genes (id: cnaGBM68).
452 Remarkably, the same loss at 1p32.3 is associated with sensitivity to genetic knockout of EGFR
453 in brain cell line panels in two independent large-scale genetic screens (Project Scores and
454 Project Achilles, Fig S9cd) and another drug screen (PRISM, Fig S9b). The meta-analysis of the
455 two drug screens and two genetic screens suggests high strength of combined evidence
456 ($p=0.00094$, Fisher's method of combining p-values) linking the loss at 1p32.3 (chr1: 51169045-
457 51472178) with sensitivity to pharmacological or genetic EGFR inhibition in brain cells, suggesting
458 a strong candidate for follow-up work.

459
460 In summary, the presence of cell lines with dubious or incorrect labels of tissue identity may
461 strongly impact association studies of drug or CRISPR screening data in two different ways. First,
462 the presence of mislabelled cell lines can cause the appearance of spurious associations that do
463 not reflect the biology of the cancer type of interest. Second, the presence of mislabelled or
464 divergent cell lines can prevent the recovery of true associations.

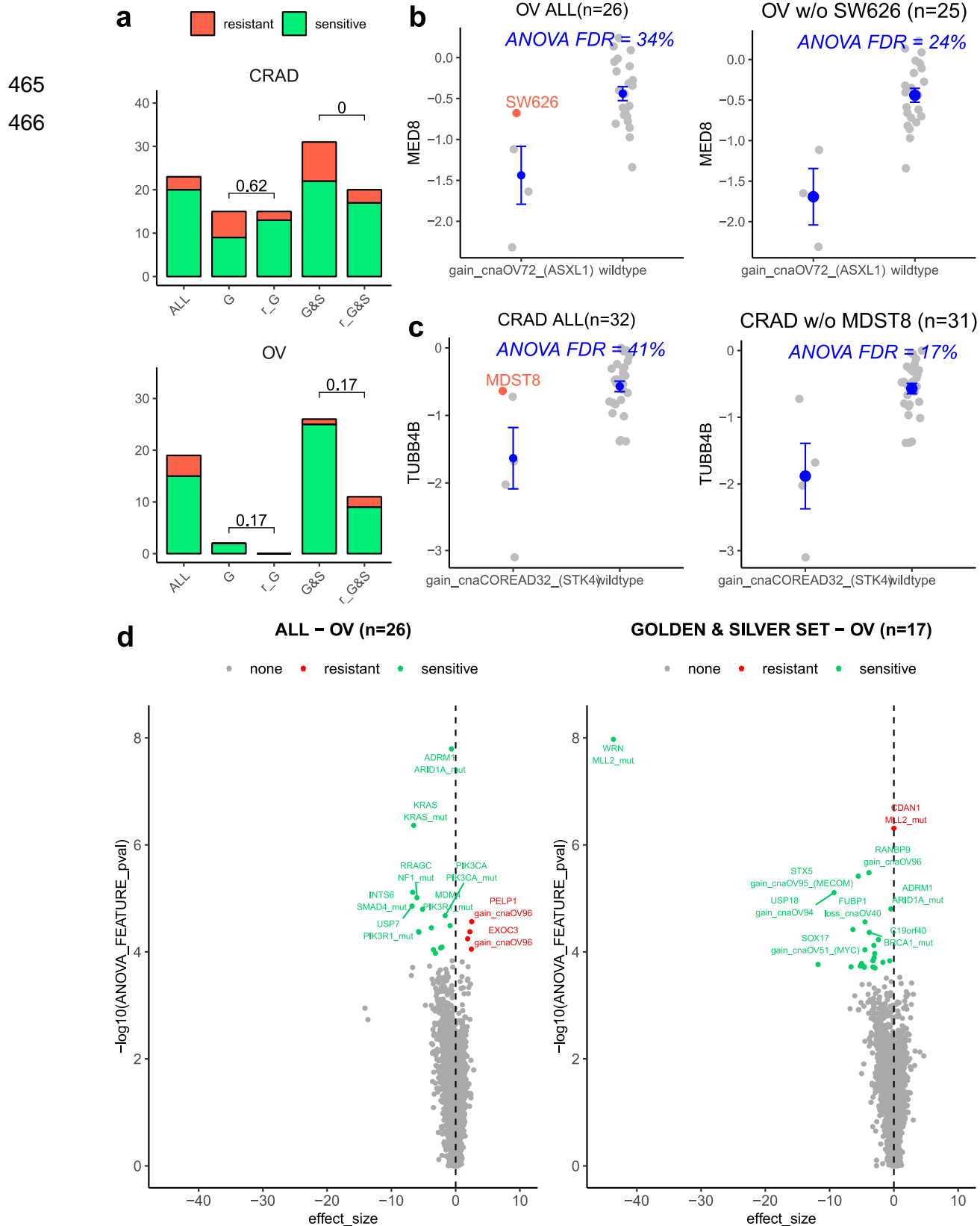


Fig. 5. Analysis of genetic screening data using high-confidence cell lines. (a) Number of significant associations between Cancer Functional Events (CFEs) and gene dependencies (in CRISPR knockout screens) detected in the ANOVA test for all cell lines (ALL), cell lines in the golden set (G), cell lines in the golden and silver set (G&S), random subset of cell lines that match the number of cell lines in the golden set (r_G) and in the golden and silver set (r_G&S). For the random subsets, the number of significant associations is calculated 10 times and median. P-value for a sign test (one-tailed) between the associations in the G/G&S and the associations in the 10 runs of r_G/r_G&S are shown. See Fig S8 for remaining cancer types. (b) Fitness effect (fold change) for MED8 k.o. in all OV cell lines (left) and all OV cell lines except SW626, which is suspected of originating from CRAD (right). Two groups are compared: cell lines with copy number gain in a region containing ASXL1 (gain_cnaOV72), and without (wild-type). ANOVA FDR for this association (MED8 k.o. and gain_cnaOV72) is shown in blue for both datasets. (c) Fitness effect (fold change) for TUBB4B k.o. in all CRAD cell lines (left) and all CRAD cell lines except MDST8, which is suspected to originate from skin (right). Two groups are compared: cell lines with copy number gain in region containing STK4 (gain_cnaCOREAD32) and without (wild-type). ANOVA FDR for this association (TUBB4B k.o. and gain_cnaCOREAD32) shown in blue for both datasets. (d) Differential dependency biomarkers were analysed by ANOVA for all ovarian cancer (OV) cell lines (left) and OV cell lines in the golden and silver set only (right). Each point is an association between the fitness effect of a gene and a genetic feature (CFE).

467 Discussion

468 Cell lines are commonly used as models for tumors, however it is an open question how to best
469 apply the available cell line panels to learn about cancer biology. The availability of genomic data
470 from large tumor cohorts and from cell line panels has spurred multiple efforts to find which cell
471 line(s) are closer to tumors by their transcriptomic (10,39,40,46) and/or genomic features (13,14),
472 presumably making better models, and which are more distant from examples of actual tumors,
473 presumably making less good models of tumor biology.

474
475 Our work addresses a different question: we attempt to detect the cancer type (i.e. tissue and/or
476 cell type) that originated the cell line, in order to ascertain if this matches the declared origin of
477 the cell line. A mismatch may conceivably stem from the sampling step, for instance a metastasis
478 might have a different tissue-of-origin than thought at time collection. The work-up after collecting
479 the tumor sample may have inaccurately assigned the cancer type based on unclear histological
480 or anatomical features. Another possibility is that the mismatch might stem from the step of
481 adaptation to cell culture, where a minority cell type not representative of the tumor prevails over
482 other tumoral cells. We consider these to be cases of cell line mislabeling during isolation. In
483 addition, we would also detect cases where the cell line might have acquired some properties of
484 a different tissue/cell type during culture, however our analyses (Fig 3c) suggest this is a less
485 common occurrence, although individual examples of this cannot be ruled out.

486
487 Importantly, this phenomenon of tissue / cell type mislabeling is distinct from well-known and
488 widespread cell line misidentification issues (23,47), where one cell line (often HeLa) was
489 mistakenly used in place of another cell line, commonly due to cross-contamination. The cell line
490 panels that provided data used in our analyses (GDSC and CCLE) have authenticated their cell
491 lines (4,44), thus misidentification/cross-contamination cannot underlie our observation that the
492 mislabeling of the cancer type of origin is not uncommon. (We note there were rare cases of
493 misidentified cells reported in these panels (44) however these do not overlap our results.)

494
495 Methodologically improving over previous work, we introduce the HyperTracker framework that
496 performs global analyses that independently examine transcriptomic, epigenomic and several
497 types of mutational features. Additionally, while carefully adjusting for the known bulk differences
498 between cell lines and tumors, which might have resulted e.g. from impurities in tumors, or from
499 altered expression of cell-cycle-related genes in cell lines (40,48). Parallel analyses of different

500 omics data provides increased confidence in our inferences, which suggested, remarkably, that
501 5.7% (35 of total 614 considered cell lines) exhibit significant transcriptomic and epigenomic
502 features of a different tissue/cell type than the declared cell type of origin. For 3.6% (22 cell lines)
503 these reassignments to a different cancer type were additionally supported in at least one type of
504 genomic evidence. This increased confidence that these were indeed examples of cell lines with
505 mislabeled (or, less likely, diverged) tissue/cell-type identity. A striking example are cell lines GCT,
506 SW684, ES2 and MDST8 that we predict to originate from the skin, based on the presence of the
507 UV mutational signature, in addition to transcriptome/epigenome data. These cases are
508 reminiscent of the recent reports of UV mutational signatures found in some cases of presumed
509 lung cancers, suggesting they may instead be metastases originating from a sun-exposed area
510 of skin (49).

511
512 In interpreting our data, an important consideration is that the cancer samples types in TCGA may
513 not necessarily reflect the full diversity of rarer subtypes within a cancer type, which may cause
514 some ambiguous predictions. For instance, the ECC10 and ECC12 cell lines are assigned to
515 STAD cancer type (stomach adenocarcinoma) when matched with TCGA tumors. However, these
516 cell lines originate from gastric small cell neuroendocrine carcinomas. This may explain why, in
517 our analysis, gene expression patterns point towards brain tissues, while mutational features
518 suggest stomach cancer. In such cases of disagreement between different types of features, a
519 future use of a more exhaustive set of reference tumor data may help resolve the ambiguity and
520 improve confidence in predictions.

521
522 The genomic classifiers we employed here were based on whole-exome sequences and were
523 overall less powerful than the transcriptome/DNA methylation classifiers in our data (Fig S4).
524 Recent work by us and others (24,25) suggests that analyzing whole-genome sequences of these
525 cell lines would permit use of additional, highly predictive features based on regional mutation
526 density of chromosomal domains. This may provide further genomic evidence for the identity of
527 the cell-of-origin for the 35 suspected cell lines. Experimental work on these cell lines will provide
528 further evidence to support or refute our predictions based on global analyses of omics data.

529
530 Knowing the correct tissue-of-origin label for a cell line is important, because this has a strong
531 bearing on the response of the cell line to drug treatment and to genetic perturbation. We
532 demonstrate the implications of this general principle to analyses of drug and genetic screening
533 data: by accounting for suspect cell lines, the power to discover new determinants of sensitivity

534 to drug/genetic perturbation may increase substantially for some cancer types, such as brain, lung
535 and pancreatic cancers. Therefore, when designing future screening efforts, it is not only
536 important to increase the number of cell lines to gain more power, but it is also important to focus
537 on the cell lines that are most consistent with the tissue and/or cell type of interest.

538
539

540 **Methods**

541 **Omics data collection and preparation**

542 DNA methylation data. We downloaded DNA methylation data as beta values (platform Illumina
543 Human Methylation 450) from GDC Data Portal (50) for TCGA samples and from Genomics of
544 Drug Sensitivity in Cancer (GDSC) (3) for CL samples. We filtered out all probes outside promoter
545 regions and probes with NA values in more than 100 samples. For the probes in promoter regions,
546 we selected only one probe per gene, keeping the probe with the highest standard deviation (sd)
547 across samples. We transformed the beta-values to m-values (log₂ ratio of the intensities of
548 methylated probe versus unmethylated probe). In total, we have 10,141 features for 942 CL
549 samples and 8,453 TCGA samples.

550

551 Gene expression data. We downloaded gene expression data as transcripts per million (TPM)
552 from GDC Data Portal (50) for TCGA samples and from GDSC (3) and the Cancer Cell Line
553 Encyclopedia (CCLE) (4) for CL samples. We filtered out genes with NA values in more than 100
554 samples and selected the overlapping genes between the 3 sources. We removed low expressed
555 genes (TMP<1 in 90% of the samples). We applied square root transformation to the data. In
556 total, we have 12,419 features for 942 CL samples and 9,149 TCGA samples.

557

558 Finally, for both DNA methylation (MET) and gene expression (GE), we created datasets of
559 different sizes: 1,000; 2,000; 3,000; 5,000; and 8,000 features by selecting the genes/probes with
560 the highest standard deviation across TCGA samples only.

561

562 Copy Number Alteration data. We downloaded Copy Number Alteration data (computed by gene)
563 from GDC Data Portal (50) for TCGA samples and from DepMap (51) for CL samples. In total, we
564 have 20,491 features for 942 CL samples and 9,188 TCGA samples. To reduce the dataset, we
565 selected 299 cancer driver genes (52) and filtered out the rest.

566

567 Mutation data. For human tumors, we downloaded mutation data as whole exome sequencing
568 (WES) MC3 dataset (53) from the GDC Data Portal for TCGA samples. For cell lines, bam files
569 were obtained from European Genome-phenome Archive (EGA) (ID number:
570 EGAD00001001039). Variant calling was performed using Strelka (version 2.8.4) with default
571 parameters. Variant annotation was performed using ANNOVAR (version 2017-07-16). In
572 samples where Strelka was unable to run, a re-alignment was performed using Picard tools
573 (version 2.18.7) to convert the bams to FASTQ and, following that, the alignment was performed
574 by executing bwa sampe (version 0.7.16a) with default parameters. The resulting bam files were
575 sorted and indexed using Picard tools. To account for germline variants, we removed all mutations
576 that were present in the gnomAD database (54) at an allele frequency ≥ 0.001 in any of the
577 populations. Finally, using the filtered somatic mutations we calculated three set of mutational
578 features: Regional Mutation Density (RMD), Mutation Spectra (MS96) and Oncogenic Mutations
579 (OGM) as described in Salvadores et. al (24). RMD features did not exhibit high accuracy when
580 applied to exome-sequencing data and so were not considered further in this analysis.

581

582 For the cell line samples, we matched their cancer types to the TCGA cancer types using the cell
583 line metadata from GDSC (3) and manually annotated those that did not have TCGA label using
584 cellosaurus (55). Next, we selected the cell lines from solid tumor that had a matching cancer type
585 in TCGA, ending up with a total of 614 cell lines from 22 cancer types. Blood cancers (LAML and
586 DLBC) are not tested because they are commonly growth in suspension, therefore their confusion
587 with solid tumors is less likely to occur. For further analysis, we merged the cancer types that
588 were overall similar: HNSC with LUSC and ESCC (SCC), GBM with LGG (LGGBM), STAD with
589 ESAD (ESTAD) and OV with UCEC (GYNE).

590

591 The identification of the cell line samples were performed by the databases providing the data
592 using short tandem repeat (STR) analysis (4,44). Of note, they reported a few commonly
593 misidentified cell lines: Ca9-22, RIKEN, MKN28, KP-1N, OVMIU and SK-MG-1 (44). These cell
594 lines do not overlap with our suspected samples and additionally the misidentification does not
595 impact tissue or cancer type of origin.

596

597 **Data alignment between tumors and cell lines**

598 For the alignment of TCGA and CL data we first applied quantile normalization (*R package*
599 *preprocessCore 1.46.0*) and second applied ComBat (*R package sva 3.32.1*), a batch effect
600 correction method. We used ComBat as if our dataset was the TCGA and CL data combined, and
601 the batch effects were whether a sample belongs to TCGA or CL (for MET) or a sample belongs
602 to TCGA, GDSC or CLLE (for GE). We applied this method for GE, MET, CNA, MS96 and RMD.
603 For validation, we calculated a principal component analysis (PCA) subsampling TCGA data to
604 match the CL samples (stratified by cancer types). Additionally, we calculated Elastic Net (EN)
605 classifiers to predict (in the processed dataset) TCGA *versus* CL and calculated the AUC and
606 AUPRC to check whether the process of alignment is being successful or not.

607
608 In addition to the chosen adjustment method, we tested other approaches based on Canonical
609 Correlation Analysis, Partial Least Squares and principal component analysis, which did not
610 exceed accuracy of ComBat (data not shown) and therefore were not examined further.

611

612 **Cancer type classifiers**

613 For the TCGA dataset we generated Ridge regression model for predicting the cancer type in a
614 One-vs-Rest manner (*using cv.glmnet function with alpha=0 and family = binomial, R package*
615 *glmnet 2.0.18*). To calculate the accuracy, we trained classifiers in the TCGA dataset and tested
616 in the CL dataset. In particular, we calculated the Area Under the Receiver Operating
617 Characteristic curves (AUC) and the Area Under the Precision Recall curve (AUPRC) for each
618 cancer type vs the rest (all the rest of cancer types combined).

619
620 FDR Score. For each cell line, we calculated an FDR score of belonging to a particular cancer
621 type. For this, we divided the TCGA data into two datasets (training and testing) of the same size
622 keeping the cancer type proportions. For each cancer type, we trained classifiers in the TCGA
623 training dataset and we introduced the cell lines one by one with the testing data and calculated
624 the precision recall (PR) curve (TCGA testing + 1CL). We set the cell line FDR score for that
625 specific cancer type as (1 - precision) at the threshold where the cell line is situated in the PR
626 curve. Overall, for every cell line we obtained 17 FDR scores, 1 for each possible cancer type.
627 We repeated this procedure 5 times and calculated the median FDR for every case to get more
628 robust values. In addition, when training for 1 cancer type (label = 1) versus the rest of cancer
629 types combined (label = 0) we made some exceptions and removed those cancer types which
630 are similar and therefore the classifier is not good at separating them (e.g. when we calculated

631 FDR for ESTAD, we removed from the rest CRAD and PAAD, all hidden cases in [Table S7](#)). This
632 is conservative with respect to reassigning cell lines to another cancer type, however some
633 resolution is traded off because the more closely related cancer types are, by design, not
634 distinguished. We have further attempted to reclassify cell lines within these hidden tissues and
635 the combined ones. However, when using One-vs-One classifiers the accuracy is not good
636 enough for distinguishing the two cancer types in the cell lines (data not shown).

637
638 Once we have a list of suspected cell lines, we have an “original” cancer type and a “suspected”
639 cancer type. Therefore, we generated One-vs-One classifiers (original *versus* suspected) using
640 TCGA dataset (balancing the classes) and for each suspected cell line we checked if it is predicted
641 as “original” or “suspected”. We repeated this prediction 20 times and counted the number of
642 times a cell line is predicted as suspected. Therefore, we defined a prediction score (range
643 between 0 and 20) for every cell line, where 0 means never predicted as suspected and 20 always
644 predicted as suspected. As a control, we repeated the same procedure with randomized cancer
645 type labels. We calculated this prediction score for GE, MET, CNA, OGM and MS96 datasets.
646 For calculating the FDR at a score ≥ 15 , we applied this formula: $FDR = FP/(FP+TP)$ where FP
647 was the number of cell lines with score ≥ 15 in the randomized data and TP the number of cell
648 lines with score ≥ 15 in the actual data.

649

650 **Independent validation**

651 We downloaded drug sensitivity for the CL from GDSC (3). From all the drugs we selected
652 *trametinib* and *dabrafenib*, FDA-approved drugs for melanoma treatment. We compared IC50
653 values for these two drugs for all cancer types.

654

655 We downloaded mutational signatures from cell lines available in Jarvis *et al* (31) and Petljak *et*
656 *al* (30) and we compared the exposures of all cell lines for Signature 7 (UV light). In Petljak *et al*
657 dataset the signature 7 is divided into Signature 7a, b, c and d. Therefore, we used the sum of
658 exposures across all four subtypes of Signature 7.

659

660 We downloaded another set of drug screening data (PRISM 19Q3) (34) for the CL dataset. For
661 the suspected cell lines, we generated One-vs-One classifiers (*using cv.glmnet function with*
662 *alpha=1 and family = binomial, R package glmnet 2.0.18*) for predicting original *versus* suspected
663 cancer type based on the drug sensitivity data. We performed 20 runs of each case and count

664 how many times it is predicted as suspected (prediction score 0-20). Additionally, we calculated
665 the AUC for each classifier.

666

667 **Subtype classifiers**

668 We downloaded cancer subtypes for the TCGA samples from the *R package TCGAbiolinks 2.12.6*
669 (56). We combined the GE and MET datasets. For this data, we generated Ridge regression
670 model for predicting the subtypes in a One-vs-Rest manner (*using cv.glmnet function with alpha=0*
671 *and family = binomial, R package glmnet 2.0.18*) within each cancer type. We trained models in
672 TCGA and we predicted subtypes for the cell lines. Additionally, we used cell line's subtypes for
673 breast cancer from a previous paper (38) to calculate the confusion matrix and the AUPRC.

674

675 We performed a chi-square test (*R package stats 3.6.0*) and calculated the cramer's V statistic
676 (*R package lsr 0.5*) for checking whether the proportion of subtypes between TCGA and CL is
677 maintained for each cancer type.

678

679 **Drug and CRISPR screening data**

680 We downloaded drug sensitivity and cancer functional event (CFE) data from the Iorio *et al.* study
681 (11). Cancer Functional Events (CFEs) are a collection of recurrent mutations, CNA and
682 hypermethylation events present in human tumors (11). We used GDSCTools (57) to search for
683 associations between the drugs and the CFEs in every cancer as they did. In particular, we
684 performed this analysis using for each tissue (i) all the cell lines; (ii) only the cell lines in the golden
685 set (G); (iii) only the cell lines in the golden and silver set (G&S). Additionally, as controls we
686 included a random subset of all cell lines matching (iv) the number of cell lines in goldenset (r_G)
687 and (v) the number in golden and silver set combined (r_G&S). We counted the number of
688 significant hits ($FDR \leq 25\%$) for each of the cancer types. For the controls, we repeated the
689 subsampling 10 times and took the median of significant hits. We compared the number of hits
690 for all the cell lines (same as in Iorio *et al.* study) with the number of hits for the different subsets
691 of cell lines according to our grouping. Additionally, we performed a sign test (*R package BSDA*
692 *1.2.0*) comparing the significant hits in the G/G&S subsets versus the significant hits over 10 runs
693 in the random G/ random G&S and calculated the p-value for all cancer types (alternative = "less").

694

695 Similarly, we downloaded gene dependency data from Project Score (44) and Project Achilles
696 (58) processed with the Project Score pipeline and combined them. From a total of uniquely 696

697 cell lines, 357 overlap with the 600 cell lines tested with our method. For those 357 tested cell
698 lines, we repeated the same procedure as described above for the drug sensitivity data.

699

700 References

- 701 1. Kaur G, Dufour JM. Cell lines: Valuable tools or useless artifacts. *Spermatogenesis*. 2012 Jan 1;2(1):1–5.
- 702 2. Goodspeed A, Heiser LM, Gray JW, Costello JC. Tumor-Derived Cell Lines as Molecular Models of Cancer
703 Pharmacogenomics. *Mol Cancer Res MCR*. 2016 Jan;14(1):3–13.
- 704 3. Yang W, Soares J, Greninger P, Edelman EJ, Lightfoot H, Forbes S, et al. Genomics of Drug Sensitivity in Cancer (GDSC): a
705 resource for therapeutic biomarker discovery in cancer cells. *Nucleic Acids Res*. 2013 Jan;41(Database issue):D955–961.
- 706 4. Barretina J, Caponigro G, Stransky N, Venkatesan K, Margolin AA, Kim S, et al. The Cancer Cell Line Encyclopedia enables
707 predictive modelling of anticancer drug sensitivity. *Nature*. 2012 Mar;483(7391):603–7.
- 708 5. Basu A, Bodycombe NE, Cheah JH, Price EV, Liu K, Schaefer GI, et al. An interactive resource to identify cancer genetic and
709 lineage dependencies targeted by small molecules. *Cell*. 2013 Aug 29;154(5):1151–61.
- 710 6. Geraghty RJ, Capes-Davis A, Davis JM, Downward J, Freshney RI, Knezevic I, et al. Guidelines for the use of cell lines in
711 biomedical research. *Br J Cancer*. 2014 Sep 9;111(6):1021–46.
- 712 7. Vaughan L, Glänzel W, Korch C, Capes-Davis A. Widespread Use of Misidentified Cell Line KB (HeLa): Incorrect Attribution
713 and Its Impact Revealed through Mining the Scientific Literature. *Cancer Res*. 2017 Jun 1;77(11):2784–8.
- 714 8. Neimark J. Line of attack. *Science*. 2015 Feb 27;347(6225):938–40.
- 715 9. Campbell BB, Light N, Fabrizio D, Zatzman M, Fuligni F, Borja R de, et al. Comprehensive Analysis of Hypermutation in
716 Human Cancer. *Cell*. 2017 Nov 16;171(5):1042–1056.e10.
- 717 10. Virtanen C, Ishikawa Y, Honjoh D, Kimura M, Shimane M, Miyoshi T, et al. Integrated classification of lung tumors and cell
718 lines by expression profiling. *Proc Natl Acad Sci*. 2002 Sep 17;99(19):12357–62.
- 719 11. Iorio F, Knijnenburg TA, Vis DJ, Bignell GR, Menden MP, Schubert M, et al. A Landscape of Pharmacogenomic Interactions in
720 Cancer. *Cell*. 2016 Jul 28;166(3):740–54.
- 721 12. Stewart A, Coker EA, Pölsterl S, Georgiou A, Minchom AR, Carreira S, et al. Differences in Signaling Patterns on PI3K
722 Inhibition Reveal Context Specificity in *KRAS*-Mutant Cancers. *Mol Cancer Ther*. 2019 Aug;18(8):1396–404.
- 723 13. Domcke S, Sinha R, Levine DA, Sander C, Schultz N. Evaluating cell lines as tumour models by comparison of genomic
724 profiles. *Nat Commun*. 2013 Jul 9;4:2126.
- 725 14. Sinha R, Winer AG, Chevinsky M, Jakubowski C, Chen Y-B, Dong Y, et al. Analysis of renal cancer cell lines from two major
726 resources enables genomics-guided cell line selection. *Nat Commun*. 2017 Aug;8(1):15165.
- 727 15. Jonsson P, Bandlamudi C, Cheng ML, Srinivasan P, Chavan SS, Friedman ND, et al. Tumour lineage shapes BRCA-mediated
728 phenotypes. *Nature*. 2019 Jul;571(7766):576–9.
- 729 16. Kopetz S, Desai J, Chan E, Randolph Hecht J, J O'Dwyer P, Maru D, et al. Phase II Pilot Study of Vemurafenib in Patients
730 With Metastatic BRAF-Mutated Colorectal Cancer. *J Clin Oncol Off J Am Soc Clin Oncol*. 2015 Oct 13;33.
- 731 17. Park S, Lehner B. Epigenetic epistatic interactions constrain the evolution of gene expression. *Mol Syst Biol*. 2013;9:645.
- 732 18. Kim E, Dede M, Lenoir WF, Wang G, Srinivasan S, Colic M, et al. A network of human functional gene interactions from
733 knockout fitness screens in cancer cells. *Life Sci Alliance*. 2019 Apr 1;2(2):e201800278.
- 734 19. Lukk M, Kapushesky M, Nikkilä J, Parkinson H, Goncalves A, Huber W, et al. A global map of human gene expression. *Nat*
735 *Biotechnol*. 2010 Apr;28(4):322–4.
- 736 20. Varley KE, Gertz J, Bowling KM, Parker SL, Reddy TE, Pauli-Behn F, et al. Dynamic DNA methylation across diverse human
737 cell lines and tissues. *Genome Res*. 2013 Mar;23(3):555–67.
- 738 21. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical Bayes methods.
739 *Biostatistics*. 2007 Jan 1;8(1):118–27.
- 740 22. Chen C, Grennan K, Badner J, Zhang D, Gershon E, Jin L, et al. Removing Batch Effects in Analysis of Expression Microarray

- 741 Data: An Evaluation of Six Batch Adjustment Methods. PLoS ONE [Internet]. 2011 Feb 28 [cited 2019 Jul 17];6(2). Available
742 from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3046121/>
- 743 23. Horbach SPJM, Halffman W. The ghosts of HeLa: How cell line misidentification contaminates the scientific literature. PLOS
744 ONE. 2017 Oct 12;12(10):e0186281.
- 745 24. Salvadores M, Mas-Ponte D, Supek F. Passenger mutations accurately classify human tumors. PLOS Comput Biol. 2019 Apr
746 15;15(4):e1006953.
- 747 25. Jiao W, Atwal G, Polak P, Karlic R, Cuppen E, Danyi A, et al. A deep learning system can accurately classify primary and
748 metastatic cancers based on patterns of passenger mutations. bioRxiv. 2019 Jan 22;214494.
- 749 26. Furlong MT, Hough CD, Sherman-Baust CA, Pizer ES, Morin PJ. Evidence for the colonic origin of ovarian cancer cell line
750 SW626. J Natl Cancer Inst. 1999 Aug 4;91(15):1327–8.
- 751 27. Medico E, Russo M, Picco G, Cancelliere C, Valtorta E, Corti G, et al. The molecular landscape of colorectal cancer cell lines
752 unveils clinically actionable kinase targets. Nat Commun. 2015 Apr 30;6:7002.
- 753 28. Alexandrov LB, Jones PH, Wedge DC, Sale JE, Campbell PJ, Nik-Zainal S, et al. Clock-like mutational processes in human
754 somatic cells. Nat Genet. 2015 Dec;47(12):1402–7.
- 755 29. Hayward NK, Wilmott JS, Waddell N, Johansson PA, Field MA, Nones K, et al. Whole-genome landscapes of major melanoma
756 subtypes. Nature. 2017 May;545(7653):175–80.
- 757 30. Petljak M, Alexandrov LB, Brummel JS, Price S, Wedge DC, Grossmann S, et al. Characterizing Mutational Signatures in
758 Human Cancer Cell Lines Reveals Episodic APOBEC Mutagenesis. Cell. 2019 Mar 7;176(6):1282–1294.e20.
- 759 31. Jarvis MC, Ebrahimi D, Temiz NA, Harris RS. Mutation Signatures Including APOBEC in Cancer Cell Lines. JNCI Cancer
760 Spectr [Internet]. 2018 Jan 1 [cited 2019 Sep 2];2(1). Available from:
761 <https://academic.oup.com/jncics/article/2/1/pky002/4942295>
- 762 32. Ji J, Chen X, Leung SY, Chi J-TA, Chu KM, Yuen ST, et al. Comprehensive analysis of the gene expression profiles in human
763 gastric cancer cell lines. Oncogene. 2002 Sep 19;21(42):6549–56.
- 764 33. Corcoran RB, Atreya CE, Falchook GS, Kwak EL, Ryan DP, Bendell JC, et al. Combined BRAF and MEK Inhibition With
765 Dabrafenib and Trametinib in BRAF V600–Mutant Colorectal Cancer. J Clin Oncol. 2015 Dec 1;33(34):4023–31.
- 766 34. Corsello SM, Nagari RT, Spangler RD, Rossen J, Kocak M, Bryan JG, et al. Non-oncology drugs are a source of previously
767 unappreciated anti-cancer activity. bioRxiv. 2019 Aug 9;730119.
- 768 35. The Cancer Genome Atlas Network. Comprehensive molecular portraits of human breast tumours. Nature. 2012
769 Oct;490(7418):61–70.
- 770 36. Ceccarelli M, Barthel FP, Malta TM, Sabedot TS, Salama SR, Murray BA, et al. Molecular Profiling Reveals Biologically
771 Discrete Subsets and Pathways of Progression in Diffuse Glioma. Cell. 2016 Jan 28;164(3):550–63.
- 772 37. Liu Y, Sethi NS, Hinoue T, Schneider BG, Cherniack AD, Sanchez-Vega F, et al. Comparative Molecular Analysis of
773 Gastrointestinal Adenocarcinomas. Cancer Cell. 2018 Apr 9;33(4):721–735.e8.
- 774 38. Neve RM, Chin K, Fridlyand J, Yeh J, Baehner FL, Fevr T, et al. A collection of breast cancer cell lines for the study of
775 functionally distinct cancer subtypes. Cancer Cell. 2006 Dec;10(6):515–27.
- 776 39. Ronen J, Hayat S, Akalin A. Evaluation of colorectal cancer subtypes and cell lines using deep learning [Internet].
777 Bioinformatics; 2018 Nov [cited 2019 Aug 23]. Available from: <http://biorxiv.org/lookup/doi/10.1101/464743>
- 778 40. Yu K, Chen B, Aran D, Charalel J, Yau C, Wolf DM, et al. Comprehensive transcriptomic analysis of cell lines as models of
779 primary tumors across 22 tumor types. Nat Commun. 2019 Aug 8;10(1):1–11.
- 780 41. Ikediobi ON, Reimers M, Durinck S, Blower PE, Futreal AP, Stratton MR, et al. In vitro differential sensitivity of melanomas to
781 phenothiazines is based on the presence of codon 600 BRAF mutation. Mol Cancer Ther. 2008 Jun 1;7(6):1337–46.
- 782 42. Yamaguchi Y, Kasukabe T, Kumakura S. Piperlongumine rapidly induces the death of human pancreatic cancer cells mainly
783 through the induction of ferroptosis. Int J Oncol. 2018 Mar 1;52(3):1011–22.
- 784 43. Dhillon H, Mamidi S, McClean P, Reindl KM. Transcriptome Analysis of Piperlongumine-Treated Human Pancreatic Cancer
785 Cells Reveals Involvement of Oxidative Stress and Endoplasmic Reticulum Stress Pathways. J Med Food. 2016 Apr
786 27;19(6):578–85.
- 787 44. Behan FM, Iorio F, Picco G, Gonçalves E, Beaver CM, Migliardi G, et al. Prioritization of cancer therapeutic targets using

- 788 CRISPR–Cas9 screens. *Nature*. 2019 Apr;568(7753):511–6.
- 789 45. Westphal M, Maire CL, Lamszus K. EGFR as a Target for Glioblastoma Treatment: An Unfulfilled Promise. *CNS Drugs*.
790 2017;31(9):723–35.
- 791 46. Vincent KM, Findlay SD, Postovit LM. Assessing breast cancer cell lines as tumour models by comparison of mRNA
792 expression profiles. *Breast Cancer Res*. 2015 Aug 20;17(1):114.
- 793 47. Capes- Davis A, Theodosopoulos G, Atkin I, Drexler HG, Kohara A, MacLeod RAF, et al. Check your cultures! A list of cross-
794 contaminated or misidentified cell lines. *Int J Cancer*. 2010;127(1):1–8.
- 795 48. Ertel A, Verghese A, Byers SW, Ochs M, Tozeren A. Pathway-specific differences between tumor cell lines and normal and
796 tumor tissue cells. *Mol Cancer*. 2006 Nov 2;5(1):55.
- 797 49. Campbell JD, Alexandrov A, Kim J, Wala J, Berger AH, Pedamallu CS, et al. Distinct patterns of somatic genome alterations in
798 lung adenocarcinomas and squamous cell carcinomas. *Nat Genet*. 2016;48(6):607–16.
- 799 50. Grossman RL, Heath AP, Ferretti V, Varmus HE, Lowy DR, Kibbe WA, et al. Toward a Shared Vision for Cancer Genomic
800 Data. *N Engl J Med*. 2016 Sep 22;375(12):1109–12.
- 801 51. Cancer Cell Line Encyclopedia Consortium, Genomics of Drug Sensitivity in Cancer Consortium. Pharmacogenomic
802 agreement between two cancer cell line data sets. *Nature*. 2015 Dec 3;528(7580):84–7.
- 803 52. Bailey MH, Tokheim C, Porta-Pardo E, Sengupta S, Bertrand D, Weerasinghe A, et al. Comprehensive Characterization of
804 Cancer Driver Genes and Mutations. *Cell*. 2018 05;173(2):371-385.e18.
- 805 53. Ellrott K, Bailey MH, Saksena G, Covington KR, Kandoth C, Stewart C, et al. Scalable Open Science Approach for Mutation
806 Calling of Tumor Exomes Using Multiple Genomic Pipelines. *Cell Syst*. 2018 Mar;6(3):271-281.e7.
- 807 54. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. Variation across 141,456 human exomes and
808 genomes reveals the spectrum of loss-of-function intolerance across human protein-coding genes. *bioRxiv*. 2019 Aug
809 13;531210.
- 810 55. Bairoch A. The Cellosaurus, a Cell-Line Knowledge Resource. *J Biomol Tech JBT*. 2018 Jul;29(2):25–38.
- 811 56. TCGAAbiolinks: an R/Bioconductor package for integrative analysis of TCGA data [Internet]. [cited 2019 Oct 7]. Available from:
812 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4856967/>
- 813 57. Cokelaer T, Chen E, Iorio F, Menden MP, Lightfoot H, Saez-Rodriguez J, et al. GDSCTools for mining pharmacogenomic
814 interactions in cancer. *Bioinformatics*. 2018 Apr 1;34(7):1226–8.
- 815 58. Meyers RM, Bryan JG, McFarland JM, Weir BA, Sizemore AE, Xu H, et al. Computational correction of copy number effect
816 improves specificity of CRISPR-Cas9 essentiality screens in cancer cells. *Nat Genet*. 2017 Dec;49(12):1779–84.
- 817
- 818
- 819