

11-9-2018

Predicting Response to Platin Chemotherapy Agents with Biochemically-inspired Machine Learning

Peter Rogan

The University of Western Ontario, progan@uwo.ca

Eliseos J. Mucaki

Western University, emucaki@uwo.ca

Dan Lizotte

Western University, dlizotte@uwo.ca

Follow this and additional works at: <https://ir.lib.uwo.ca/biochempub>

 Part of the [Biochemistry Commons](#)

Citation of this paper:

Rogan, Peter; Mucaki, Eliseos J.; and Lizotte, Dan, "Predicting Response to Platin Chemotherapy Agents with Biochemically-inspired Machine Learning" (2018). *Biochemistry Publications*. 190.

<https://ir.lib.uwo.ca/biochempub/190>

1 **Predicting Response to Platin Chemotherapy Agents with Biochemically-inspired**
2 **Machine Learning**

3 Eliseos J. Mucaki¹, Jonathan Z.L. Zhao¹, Daniel J. Lizotte^{2,3}, and [§]Peter K. Rogan^{1,2,3,4,5}

4

5 **Running Title:**

6 Predicting Responses to Platin Drugs by Machine Learning

7

8 **Author Affiliations**

9 ¹Department of Biochemistry, Schulich School of Medicine and Dentistry, Western University,
10 London, Canada, N6A 2C1

11 ²Department of Computer Science, Faculty of Science, Western University, London, Canada, N6A
12 2C1

13 ³Department of Epidemiology & Biostatistics, Faculty of Science, Western University, London,
14 Canada, N6A 2C1

15 ⁴Cytognomix Inc. London, Canada N5X 3X5

16 ⁵Department of Oncology, Schulich School of Medicine and Dentistry, Western University,
17 London, Canada, N6A 2C1

18

19 Author Emails: emucaki@uwo.ca, jzhao293@uwo.ca, dlizotte@uwo.ca, progan@uwo.ca

20 [§]**Correspondence to:** Peter K. Rogan (progan@uwo.ca), Department of Biochemistry, Schulich
21 School of Medicine and Dentistry, Western University, London, Ontario, Canada, N6A 2C1. 1
22 (519) 661-4255.

23 **ABSTRACT.**

24 Selection of effective genes that accurately predict chemotherapy response could
25 improve cancer outcomes. We compare optimized gene signatures for cisplatin,
26 carboplatin, and oxaliplatin response in the same cell lines, and respectively validate
27 each with cancer patient data. Supervised support vector machine learning was used to
28 derive gene sets whose expression was related to cell line GI₅₀ values by backwards
29 feature selection with cross-validation. Specific genes and functional pathways
30 distinguishing sensitive from resistant cell lines are identified by contrasting signatures
31 obtained at extreme vs. median GI₅₀ thresholds. Ensembles of gene signatures at
32 different thresholds are combined to reduce dependence on specific GI₅₀ values for
33 predicting drug response. The most accurate gene signatures for each platin are:
34 cisplatin: *BARD1, BCL2, BCL2L1, CDKN2C, FAAP24, FEN1, MAP3K1, MAPK13,*
35 *MAPK3, NFKB1, NFKB2, SLC22A5, SLC31A2, TLR4, TWIST1*; carboplatin: *AKT1,*
36 *EIF3K, ERCC1, GNGT1, GSR, MTHFR, NEDD4L, NLRP1, NRAS, RAF1, SGK1, TIGD1,*
37 *TP53, VEGFB, VEGFC*; oxaliplatin: *BRAF, FCGR2A, IGF1, MSH2, NAGK, NFE2L2,*
38 *NQO1, PANK3, SLC47A1, SLCO1B1, UGT1A1*. TCGA bladder, ovarian and colorectal
39 cancer patients were used to test cisplatin, carboplatin and oxaliplatin signatures
40 (respectively), resulting in 71.0%, 60.2% and 54.5% accuracy in predicting disease
41 recurrence and 59%, 61% and 72% accuracy in predicting remission. One cisplatin
42 signature predicted 100% of recurrence in non-smoking bladder cancer patients (57%
43 disease-free; N=19), and 79% recurrence in smokers (62% disease-free; N=35). This
44 approach should be adaptable to other studies of chemotherapy response, independent
45 of drug or cancer types.

46

47 **KEY WORDS.** Chemotherapy response, support vector machines, gene signatures,
48 cancer, cisplatin, oxaliplatin, carboplatin, machine learning, bladder cancer, breast
49 cancer, ovarian cancer

50 INTRODUCTION

51 Chemotherapy regimens are selected based on overall outcomes for specific
52 types and subtypes of cancer pathology, progression to metastasis, other high-risk
53 indications, and prognosis^{1,2}, and variability in tumor resistance has led to tiered
54 sequential strategies for selection of agents based on their overall efficacy³. We and
55 others have developed machine learning (ML)-based gene signatures (i.e. predictive
56 models) aimed at predicting response to specific chemotherapeutic agents and
57 minimizing chemoresistance based on inhibition of growth or drug targets (GI₅₀ or IC₅₀)⁴⁻
58 ⁶. In this study, we present integrated ML models of platin drug responses (cis-, carbo-
59 and oxaliplatin), and evaluate them on clinical outcomes data that were not used to
60 construct the signatures. Previous studies have reviewed the genes⁷, gene products⁸
61 and specific individual pathways that are activated and repressed by drugs⁹, but lack
62 comprehensive models of the global cellular response to drugs. We use integrated ML-
63 based signatures based on expression of multiple genes to predict key responses to
64 each of these platin agents, for the first time, at different resistance levels.

65 Cisplatin, carboplatin and oxaliplatin are each widely prescribed compounds for
66 their antineoplastic effects. While each contains platinum to form adducts with tumour
67 DNA, their effectiveness differs for specific types of cancers, such as bladder (cisplatin),
68 ovarian (cisplatin and carboplatin) and colorectal cancer (oxaliplatin). Carboplatin differs
69 in structure from cisplatin, exchanging the latter's dichloride ligands with a CBDCA
70 (cyclobutane dicarboxylic acid) group, while oxaliplatin is paired with both a DACH
71 (diaminocyclohexane) ligand and a bidentate oxalate group. These chelating ligands
72 have greater stability and solubility to aqueous solutions, which lead to differences in
73 drug toxicity compared to cisplatin¹⁰. Oxaliplatin can be up to two times as cytotoxic as
74 cisplatin, but it forms fewer DNA adducts¹¹. The large hydrophobic DACH ligand which

75 overlaps the major groove is thought to prevent binding of certain DNA repair enzymes
76 such as the POL polymerases, and may contribute to the low cross-resistance between
77 oxaliplatin and cisplatin and carboplatin¹⁰. While all three drugs can enter the cell via
78 copper transporters, organic cation transporters are oxaliplatin-specific and likely play a
79 role in its efficacy in colorectal cancer (CRC) cells where these transporters are
80 commonly overexpressed⁷. Oxaliplatin specifically plays a role in interfering with both
81 DNA and RNA synthesis, unlike cisplatin which only interferes with DNA¹². It is these
82 intrinsic properties between the platinum drugs which lead to differences in their activity
83 and resistance profiles, despite their similar mode of action.

84 We derived gene signatures to predict drug response at different sensitivity and
85 resistance levels for each of these agents. We and others have used supervised
86 learning algorithms, including random forest models¹³; support vector machine (SVM)
87 models⁶; neural networks¹⁴; and linear regression models⁵ to make these predictions.
88 Pathway and network analysis of gene expression have been used to indicate hundreds
89 of genes potentially up- and down-regulated upon cisplatin treatment¹⁵. Cisplatin-specific
90 gene signatures have been developed with integrative approaches such as elastic net
91 regression using inferred pathway activity of bladder cancer cell line data¹⁶. These
92 methods have implicated genes that have not been described previously. Supervised ML
93 with biochemically-relevant genes has also been useful for predicting drug response⁶. A
94 concern with each of these ML approaches is that an insufficient number of samples
95 coupled to a large number of features, i.e. gene expression changes, in each sample
96 can result in overfitting of the model affecting its generalizability with other sources of
97 data¹⁷. We therefore reduce the number of dimensions by selecting genes biologically
98 relevant to the drugs under observation^{6,17}. In this study, genes included in the final
99 signatures have well-defined roles in their corresponding drug responses

100 (Supplementary References, Section A). Additional selection criteria are necessary
101 when the number of genes implicated in peer-reviewed reports is still prohibitively large
102 compared to sample size.

103 Biochemically-inspired gene signatures have shown good performance in
104 predicting treatment response. A paclitaxel ML signature based on tumor gene
105 expression (GE) had a higher success predicting the pathological complete response
106 rate (pCR¹⁸) for sensitive patients (84% of patients with no / minimal residual disease)
107 than gene signatures based on differential GE analysis⁶. For gemcitabine, a signature
108 derived from both expression and copy number (CN) data from breast cancer cell lines
109 was derived, and subsequently applied to analysis of nucleic acids from patient archival
110 material. Multiple other outcome measures used to validate gene signatures include
111 prognosis⁵, Miller-Payne response¹⁹, and disease recurrence. Binary SVM classifiers
112 based on discrete time thresholds have been used to classify continuous outcome
113 measures such as prognosis and recurrence. By contrast, pCR is simpler to interpret
114 with binary SVM models. Nevertheless, differences in clinical recurrence have been
115 noted between patients demonstrated with pCR and those who do not exhibit disease
116 pathology¹⁸. This source of variability in defining patient response can confound
117 transferability of SVM models between different datasets.

118 We apply biochemically-inspired ML to predict and compare the cellular and
119 patient responses to cisplatin, carboplatin and oxaliplatin. We train models and perform
120 model selection for classification of platin resistance with cancer cell line data, and
121 validate using patient GE and clinical outcome data. Our previous gene signatures
122 derived from cell line data were based on median GI₅₀ for each drug⁶. Models (i.e. gene
123 signatures) learned and selected using the cell line data have not been re-trained prior to
124 application on the patient data, since GI₅₀ values are not available in patient samples.

125 This has been a necessary compromise; however, in this study, we derive different
126 signatures at the highest vs. the lowest levels of drug resistance. A series of candidate
127 gene signatures are derived by shifting the GI_{50} thresholds that distinguish sensitivity
128 from resistance. The frequency of genes selected at median vs. extreme thresholds
129 highlights pathways that most likely define these responses among different patient
130 subsets.

131 RESULTS

132 Selection of Platin Drug Related Genes

133 We documented genes in the peer-reviewed literature associated with drug
134 effectiveness or response (Supplemental References, Section B). For cisplatin,
135 carboplatin and oxaliplatin, this implicated 179, 90, and 288 genes, respectively (Suppl.
136 Table S1). Multiple factor analysis (MFA) was used to determine which genes were
137 correlated to GI_{50} in breast cancer cell lines through either GE and/or CN^{13} , significantly
138 reducing the sizes of the gene sets for cisplatin (N=39), carboplatin (N=28), and
139 oxaliplatin (N=55). Genes with significant relationships to GI_{50} and direction of correlation
140 (positive or inverse) are indicated in Figure 1. The diverse functions of these genes
141 included apoptosis, DNA repair, transcription, cell growth, metabolism, immune system,
142 signal transduction and membrane transport. Analysis of IC_{50} and GE levels for cisplatin-
143 treated bladder cancer cell lines confirmed these relationships evident from GI_{50} values
144 of different breast cancer lines. IC_{50} values were related to GE for *CFLAR*, *FEN1*,
145 *MAPK3*, *MSH2*, *NFKB1*, *PNKP*, *PRKAA2*, and *PRKCA*²⁰. Similarly, separate bladder cell
146 line IC_{50} values from the Genomics of Drug Sensitivity in Cancer project
147 (<http://www.cancerrxgene.org>; N=17) were correlated with GE for *CFLAR*, *FEN1*, and
148 *NFKB1*, in addition to *ATP7B*, *BARD1*, *MAP3K1*, *NFKB2*, *SLC31A2* and *SNAI1*.

149 We performed MFA on the GI_{50} values for cisplatin, carboplatin and oxaliplatin,
150 without consideration of either GE or CN. Responses to cis- and carboplatin were
151 directly correlated (a 6.2° separation between vectors), but neither was related to the
152 oxaliplatin response (Figure 2). Previous studies have shown that cisplatin-resistant cell
153 lines are generally sensitive to oxaliplatin^{21–23}.

154 SVM-based signatures were initially derived for each platin drug from breast
155 cancer cell line GE data. A 13-gene signature for cisplatin that predicts whether
156 observed growth inhibition is above or below the median GI_{50} threshold (5.2% cross-
157 validation misclassification rate) consisted of *BARD1*, *BCL2L1*, *FAAP24*, *CFLAR*,
158 *MAP3K1*, *MAPK3*, *NFKB1*, *POLQ*, *PRKAA2*, *SLC22A5*, *SLC31A2*, *TLR4*, and *TWIST1*.
159 A similarly derived carboplatin signature included *AKT1*, *ATP7B*, *EGF*, *EIF3I*, *ERCC1*,
160 *GNGT1*, *HRAS*, *MTR*, *NRAS*, *OPRM1*, *RAD50*, *RAF1*, *SCN10A*, *SGK1*, *TIGD1*, *TP53*,
161 and *VEGFB* (10.4% cross-validation misclassification). For oxaliplatin, the final SVM
162 gene signature consisted of *AGXT*, *APOBEC2*, *BRAF*, *CLCN6*, *FCGR2A*, *IGF1*, *MPO*,
163 *MSH2*, *NAGK*, *NAT2*, *NFE2L2*, *NOTCH1*, *PANK3*, *PRSS1*, and *UGT1A1* (2.1% cross-
164 validation misclassification). A cisplatin SVM generated from 17 bladder cancer cell lines
165 in cancerRxgene resulted in 2 equally accurate signatures (with 11.8% cross-validation
166 misclassification) consisting of either *PNKP* and *PRKCA* or *ATP7B*, *CFLAR*, *FEN1*,
167 *MAPK3*, *NFKB1* and *SLC22A11*. These gene signatures were not useful for predicting
168 patient outcomes due to the limited size of the training set.

169 **GI_{50} -Threshold Independent Modeling**

170 In our previous studies, we set median GI_{50} value as the threshold to
171 distinguished drug resistance and sensitivity^{5,6}. An important question is whether the
172 genes contributing to drug response are consistent among different cell lines, each with

173 their own unique GI_{50} values. Different ML gene signatures were obtained by shifting the
174 GI_{50} threshold, which changed the labels of resistant vs. sensitive cell lines. After feature
175 selection, the compositions of the corresponding gene signatures for each threshold
176 were compared. Finally, ensemble averaging of all of these optimized SVMs with
177 Gaussian kernels were derived for different GI_{50} thresholds was used to create a single
178 aggregated, threshold-independent, ML-based predictive model, comprised of all genes
179 that were selected in any of the threshold-specific models (i.e. a composite gene
180 signature).

181 Kinase (*MAPK3*, *MAP3K1*) genes and apoptotic family members (*BCL2*,
182 *BCL2L1*) were most the common in the cisplatin signatures at different GI_{50} thresholds,
183 with consistent representation of error-prone and base-excision DNA repair genes as
184 well (Figure 3A; Supplementary Table S2A). The kinases are more concentrated in
185 signatures with lower drug sensitivity thresholds, whereas *BCL2* and *BCL2L1* are more
186 ubiquitous at all levels. The error prone polymerases, *POLD1* and *POLQ*, are more
187 frequent in gene signatures with lower sensitivity thresholds, while the flap endonuclease
188 *FEN1* tends to be present at high levels of resistance. Thresholded gene signatures for
189 carboplatin-related genes commonly contained the apoptotic family member *AKT1*,
190 transcription regulation genes *ETS2* and *TP53*, as well as cell growth factors *VEGFB*
191 and *VEGFC*, although the latter was less common at lower sensitivity thresholds (Figure
192 3B). Common oxaliplatin-related genes included transporters *SLCO1B1* and *GRTP1* (but
193 not *SLC47A1*), transcription genes *NFE2L2*, *PARP15* and *CLCN6*, as well as multiple
194 metabolism-related genes (Figure 3C).

195 GI_{50} thresholded ML models were also derived using the log-loss function to evaluate
196 whether an alternative loss function (for classification) would differ significantly to the
197 misclassification-based gene signatures (by both the distribution of selected genes and

198 by model accuracy to patient data). Log-loss penalizes false classifications, whose value
199 ranges from zero (or completely accurate), to 1 (or completely inaccurate;
200 Supplementary Table S3). The overall distribution of genes across GI_{50} thresholds has
201 many distinct similarities with the gene signatures derived by misclassification. For both
202 sets of cisplatin gene signatures, *BCL2*, *BCL2L1* and *FEN1* are common in low-to-
203 moderate GI_{50} thresholds, while *NFKB1* is enriched at high thresholds (Figure 3A; Suppl.
204 Figure 1A). For carboplatin, *AKT1*, *VEGFB* and *VEGFC* are similarly distributed across
205 GI_{50} thresholds with both methods, although *VEGFB* is less densely represented in log-
206 loss based gene signatures at low GI_{50} values (Figure 3B; Suppl. Figure 1B). In both
207 sets of oxaliplatin gene signatures, *SIAE* and *SLC47A1* are represented at high densities
208 across all GI_{50} thresholds, whereas *ABCG2* is present less frequently (<50% inclusion;
209 Figure 3C and Suppl. Figure 1C). There are differences between signatures selected by
210 minimizing log-loss and misclassification rates. *EGF* and *ERCC1* were selected at a
211 greater frequency at a moderate carboplatin GI_{50} with log-loss, rather than by
212 misclassification. Similarly, oxaliplatin signature genes, *APOBEC2*, *HLA-B*, *LTA*, and
213 *MPO*, were selected considerably more often by log-loss. Therefore, while the
214 misclassification and log-loss based gene signatures are not interchangeable, overall,
215 they are quite similar to one another.

216 Log-loss gene signatures were initially constructed either by (a) a modified
217 version of the misclassification-based method, or (b) using the backwards feature
218 selection (BFS) software described in Zhao *et al.* (2018)²⁵. Multiple signatures with low
219 log-loss values can have different compositions, consistent with the possibility that there
220 may be various diverse gene combinations that can give rise to signatures with
221 satisfactory performance. However, these signatures often contain a larger number of
222 gene features than the misclassification-based signatures, and raised concerns that they

223 might be more prone to overfitting. The log-loss minimized gene signatures generated by
224 both methods had comparable compositions. The median GI₅₀ thresholded cisplatin
225 gene signature generated by the log-loss modified software [*ATP7B*, *BCL2L1*, *CDKN2C*,
226 *CFLAR*, *ERCC2*, *ERCC6*, *FAAP24*, *FOS*, *GSTO1*, *GSTP1*, *MAP3K1*, *MAPK13*, *MAPK3*,
227 *MSH2*, *MT2A*, *PNKP*, *POLD1*, *POLQ*, *PRKAA2*, *PRKCA*, *PRKCB*, *SLC22A5*, *SLC31A2*,
228 *SNAI1*, *TLR4*, *TP63*] shares 15/19 genes with the signature generated by the BFS
229 software²⁵ [*ATP7B*, *BARD1*, *BCL2*, *BCL2L1*, *ERCC2*, *FAAP24*, *FEN1*, *FOS*, *MAP3K1*,
230 *MAPK13*, *MAPK3*, *MSH2*, *MT2A*, *NFKB1*, *PNKP*, *POLQ*, *PRKCB*, *SLC22A5*, *SNAI1*]).

231 *Impact of Features in Gene Signatures*

232 To determine the contribution of individual genes on overall cross-validation
233 accuracy of a gene signature, each gene was excluded (independently), and model
234 accuracy was reassessed within every SVM signature (Supplementary Tables S2A; S2B
235 and S2C contain cis-, carbo- and oxaliplatin gene signatures, respectively). Elimination
236 of *ERCC2*, *POLD1*, *BARD1*, *BCL2*, *PRKCA* and *PRKCB* consistently significantly
237 increase misclassification error (average > 16% increase) in moderate threshold cisplatin
238 SVMs (GI₅₀ thresholds: 5.1 to 5.5). *ERCC2* and *POLD1* perform critical functions in
239 nucleotide and base excision repair, respectively. *PRKCA* and *PRKCB* are paralogs with
240 significant roles in signal transduction. *BARD1* has been shown to reduce apoptotic
241 *BCL2* in the mitochondria²⁶, and has a key role in genomic stability through its
242 association with *BRCA1*. The genes *NFKB1*, *NFKB2*, *TWIST1*, *TP63*, *PRKAA2*, and
243 *MSH2* show a high variance in increased misclassification between different gene
244 signatures. The variance of these genes may be due to epistatic interactions with other
245 biological components, including the other genes in the SVM. For example, *NFKB1* and
246 *NFKB2* are jointly included in 7 SVMs generated at a moderate GI₅₀ threshold. There is
247 evidence of possible epistasis in that the removal of either of these genes, but not

248 necessary both, will have a large impact in model misclassification rates ($\geq 18.0\%$
249 increase). The misclassification variance of *NFKB1* with *NFKB2*, is significantly lower
250 than in SVM gene signatures lacking *NFKB2*.

251 *Derivation of Gene Signatures from Bladder Carcinoma Patient Data*

252 Gene signatures derived with cell line data are to be validated on cancer patient
253 data. To explore the similarities of the gene signatures to said patient data, we also
254 developed SVMs using the cisplatin and/or carboplatin-treated TCGA (The Cancer
255 Genome Atlas) bladder urothelial carcinoma patients, using post-treatment time to
256 relapse as a surrogate criterion for different GI_{50} resistance thresholds (as performed in
257 Mucaki *et al.* [2017]²⁴; Supplementary Table S4). Similar trends to cell line SVMs are
258 apparent: *POLQ* is frequently included in gene signatures with recurrence threshold of
259 longer duration, while *FEN1* is a marker of resistance, when time to relapse is shorter.
260 However *BCL2*, which is present in a majority of breast cancer cell line SVMs, is present
261 in only one gene signature derived from TCGA data. Similarly, *MSH2* was rarely
262 selected using cell lines, yet appears in nearly all patient derived SVMs with > 1 year
263 recurrence. However, independently derived patient SVMs could not be used for any
264 other analyses.

265 **Validation of Cell-line based Models against Cancer Patient Data**

266 GI_{50} -thresholded modeling for each platin drug, generated with the breast cancer
267 cell line data, produced 70 cisplatin, 83 carboplatin, and 83 oxaliplatin SVM gene
268 signatures, respectively. In order to understand how the choice of GI_{50} threshold for
269 training on cell line data impacts predictive accuracy when the resulting gene signatures
270 are applied to patient outcomes, each of the thresholded gene signatures was applied to
271 available platin-treated patient datasets²⁷⁻³¹. In this study, cisplatin gene signatures were

272 validated on bladder cancer patient data, carboplatin signatures were validated on
273 ovarian cancer patient data, and oxaliplatin signatures were validated on colorectal
274 cancer patient data. While the available data did contain the necessary GE information,
275 the clinical response metadata differed between studies. The response of bladder
276 cancer patients to cisplatin was provided as survival post-treatment by Als *et al.*³⁰,
277 whereas colorectal cancer patients treated with oxaliplatin were categorized as
278 responders and non-responders by Tsuji *et al.*³¹. TCGA provided two different measures
279 which were used to assess predictive accuracy in our gene signatures – clinical
280 response to chemotherapy and disease-free survival. Signature accuracy was found to
281 be similar using either measure (Supplementary Table S5A); however recurrence and
282 disease-free survival was used as the primary measure of response, as it was more
283 consistently recorded among the TCGA data sets tested. Patients from Als *et al.*³⁰ with
284 a ≥ 5 year survival post-treatment were labeled as sensitive to treatment. The
285 differences between these metadata may, in part, contribute to differences in the
286 prediction accuracy of the thresholded SVM gene signatures.

287 At higher resistance thresholds for any platin drug (low GI_{50}), where more cell
288 lines are labeled sensitive, the positive class (disease-free survival) is correctly
289 classified, while the negative class (recurrence) is highly misclassified (Suppl. Figures 2
290 and 3). The reverse is true for gene signatures derived using lower resistance thresholds
291 (high GI_{50}). For these reasons, SVMs generated at these extreme thresholds were not
292 very useful at predicting patient outcomes. When used to predict recurrence in the
293 TCGA datasets, sensitivity and specificity appears to be maximized in gene signatures
294 where the GI_{50} threshold for resistance was set near (but not necessarily at) the median
295 (Suppl. Figure 2; Suppl. Tables S5A to 5C). While this pattern holds true for data from
296 Tsuji *et al.*³⁰, oxaliplatin gene signatures where GI_{50} thresholds were set above the

297 median could better separate primary and metastatic CRC patients (best signature
298 predicting 92.6% metastatic and 60.7% primary cancers; Suppl. Table S5C). Although
299 less consistent, cisplatin gene signatures generated with thresholds above median GI_{50}
300 performed better when evaluating the patient dataset from Als *et al.*³⁰ (Suppl. Figure 3).

301 Gene signatures were individually evaluated for their accuracy in TCGA patients
302 using various recurrence times post-treatment to classify resistant and sensitive patients
303 (0.5 - 5 years; Supplemental Table S6A-C). The best performing cisplatin signature
304 (hereby identified as **Cis1**; Table 1) was able to accurately predict 71.0% of bladder
305 cancer patients who recurred after 18 mo. (N=31; 58.5% accurate for disease-free
306 patients [N=41]). The best performing carboplatin gene signature (designated **Car1**
307 [Table 1]) predicted recurrence of ovarian cancer after 4 years at an accuracy of 60.2%
308 (N=302; 61.0% accurate for disease-free patients [N=108]). For oxaliplatin, the best
309 performing gene signature (designated **Oxa1** [Table 1]) accurately predicted 71.6% of
310 the disease-free TCGA CRC patients after one year (N=88; 54.5% accuracy predicting
311 recurrence [N=11]). These gene signatures (based on GE measured by Affymetrix Gene
312 Chip Human Exon 1.0 ST arrays), TCGA sample expression data, as well as SVMs
313 based on bladder cell line data (based on expression measured by Affymetrix U133A
314 microarray), were added to the online web-based SVM calculator
315 (<http://chemotherapy.cytognomix.com>; introduced in Dorman *et al.* [2016]⁶) to predict
316 platin response.

317 The TCGA bladder cancer data set contained 19 patients treated with carboplatin
318 (but not cisplatin), which enabled evaluation of the specificity of cisplatin models relative
319 to patients not treated with this drug. The cisplatin model which best predicted outcomes
320 of carboplatin-treated TCGA bladder patients was not **Cis1** (the best performing cisplatin
321 model) but rather **Cis12** at two years post-treatment (80% accurate for responding

322 patients [N=5]; 93% for recurrent patients [N=14]). **Cis12** contains 9 genes not present in
323 **Cis1** including *ATP7B*, which is a gene found in many of our carboplatin models. The
324 presence of this gene may have a significant impact on the overall accuracy of **Cis12** to
325 the carboplatin-treated bladder cancer patients. We also evaluated these 19 patients to
326 the carboplatin-specific gene signatures, and found the signature which best predicted
327 the response of these patients (**Car73**) was 84% accurate for patients after 1 year of
328 treatment (100% for responding patients [N=11]; 62.5% accuracy for recurrent [N=8]).
329 Interestingly, **Car73** shares the same *ATP7B* gene with **Cis12**. Two additional
330 carboplatin gene signatures are tied for overall accuracy (84%; **Car9** and **Car51**), but
331 more successfully predict non-responsive patients (87.5%; 82% accuracy for responding
332 patients). *AKT1*, *ETS2*, *GNGT1*, and *VEGFB* were shared among these carboplatin
333 gene signatures.

334 To evaluate the consistency in the response prediction of TCGA bladder cancer
335 patients treated with cisplatin, distances from the hyperplane for all SVMs generated
336 were determined for patients with short recurrence time (<6 mo., N=10; Supplementary
337 Figure 4). Despite showing similar levels of resistance to treatment, distances differed
338 between patients. While these patients would be expected to be indicated as highly
339 cisplatin resistant (hyperplane distance < 0), two patients (TCGA-XF-A9SU and TCGA-
340 FJ-A871) were predicted sensitive across nearly all SVM gene signatures. Similar
341 variation was also seen in patients with either a long recurrence time (>4 years) or no
342 recurrence at all after 6 years (Suppl. Figure 5).

343 An aggregate, threshold-independent model was generated for each individual
344 platin drug at different GI_{50} thresholds through ensemble ML, which involves the
345 averaging of hyperplane distances for each model to generate a composite score for
346 each TCGA patient tested (i.e. a composite gene signature). Hyperplane distances

347 across all 70 cisplatin gene signatures were similar, with a mean score of -0.22 and a
348 standard deviation of 3.5 hyperplane units (hu) across the set of patient data. The
349 ensemble model classified disease-free bladder cancer patients with 59% accuracy
350 and those with recurrent disease with 47% accuracy. Limiting ensemble averaging to
351 only cisplatin gene signatures generated at a moderate GI_{50} threshold (ranging from 5.10
352 to 5.50) did not significantly improve accuracy (44% for disease-free and 66% for
353 recurrent patients; Suppl. Table S7A). For carboplatin, ensemble ML did not produce
354 significantly better predictions than random, regardless of the GI_{50} threshold interval
355 selected (Suppl. Table S7B) or the similar mean hyperplane distances (-0.11 +/- 3.9 hu).
356 For oxaliplatin, the ensemble ML model (mean = -0.12 +/- 2.7 hu) was most accurate
357 after 1 year (60% accuracy for disease-free and 73% for recurrent patients; Suppl. Table
358 S7C). As in cisplatin, limiting this analysis to oxaliplatin SVM gene signatures with
359 moderate GI_{50} thresholds did not significantly increase accuracy.

360 *K-Fold Cross-Validation*

361 The misclassification-based cisplatin, carboplatin and oxaliplatin gene signatures
362 were also evaluated by k-fold cross-validation on TCGA bladder, ovarian and colorectal
363 cancer patient data, respectively. This cross-validation is independent of cell line data;
364 that is, the genes and hyper-parameters of signatures are used, but the GE data used is
365 exclusively from patients. Patients were evenly distributed in 5 groups with an equal (or
366 near-equal) ratio of disease-free and recurrent patients. The majority of the cisplatin
367 gene signatures showed an overall accuracy > 50%. The cisplatin gene signature which
368 performed best under the k-fold analysis (6-resistance level; *BARD1*, *BCL2*, *BCL2L1*,
369 *PRKAA2*, *PRKCA*, *PRKCB*, *TWIST1*) showed an overall accuracy of 71.2% (84.4%
370 accurate for sensitive and 53.9% accurate for resistant patients). The accuracy of the
371 carboplatin and oxaliplatin gene signatures did not exceed 60%. In general, treating the

372 patient data as a held-out test set yielded higher performance estimates than training
373 and evaluating the models on the patient data using k-fold cross-validation.

374 **Predicting cisplatin response in patients based on smoking history**

375 Tobacco smoking is known as the highest risk factor for the development of
376 bladder cancer³². Head and neck cancer patients who smoke while undergoing cisplatin
377 and radiotherapy treatment have been shown to have a shorter overall survival rate³³.
378 We therefore subdivided the patients based on their smoking history and tested the
379 thresholded gene signatures (Supplementary Tables S8 and S9). When testing patients
380 who were lifelong non-smokers, the prediction accuracy of **Cis1** predicted all non-
381 smoking patients who were recurrent after 18 months as cisplatin-resistant (N=5).
382 Prediction accuracy for disease-free patients was 57.1% (N=14). Another gene signature
383 (**Cis18**; Suppl. Table S8) had performed equally as well for non-smokers, and these two
384 gene signatures share the genes *BCL2*, *BCL2L1*, *FAAP24*, *MAP3K1*, *MAPK13*, *MAPK3*,
385 and *SLC31A2*. Threshold independent analysis predicted disease-free equally well, but
386 recurrence was less accurate (66.7%). Note that non-smokers make up a small subset
387 of the patients tested (N=19). Threshold-independent prediction of recurrence in patients
388 with a smoking history was 46% accurate (N=13), while disease-free patients were
389 correctly predicted at a rate of 58% (N=19). Recurrence in these patients was best
390 predicted by a gene signature built at the median GI₅₀ threshold (**Cis2**). Accuracy
391 improved for both disease-free (57.7% -> 61.9%) and recurrent patients (76.0% ->
392 78.6%) when excluding patients who quit smoking more than 15 years before diagnosis.
393 This SVM includes *CFLAR* and *PRKAA2*, genes which are not present in the two gene
394 signatures which performed well for non-smokers.

395 To determine which genes in these gene signatures led to discordant predictions
396 of patient outcome, we gradually altered the expression of each signature gene until the
397 misclassification was corrected. Expression of *MAP3K1*, *MAPK3*, *SLC22A5* and
398 *SLC31A2*, when altered, corrected discordant predictions of patient outcome. Altering
399 *BCL2L1* expression was more likely to correct the discordant predictions of **Cis1** (4 out
400 of 5) than with **Cis2** (2 out of 4). If the change exceeded ≥ 3 -fold the highest/lowest
401 expression of that gene and the prediction was still unchanged between different
402 patients, the effect of that gene was considered to be minor. Expression of *PRKAA2*,
403 *NFKB1*, *NFKB2* and *TWIST1* could not be altered sufficiently to correct a discordant
404 prediction.

405 *Cytosine Methylation Levels of Genes in Cisplatin Models*

406 Tobacco smoking has a significant impact on cytosine methylation levels in the
407 genome³⁴. CpG island methylation has been associated with smoking pack years in a
408 subset of the TCGA bladder urothelial carcinoma patients²⁶. We suspected that the level
409 of methylation measured in the SVMs which performed best for smoking and non-
410 smoking patients might differ, and with possible concomitant effects on GE. When
411 ranking each gene from **Cis1** by highest methylation and GE, 88 of 1080 patient: gene
412 combinations showed the expected inverse correlation between methylation levels and
413 GE (i.e. high methylation and low GE). Inverse correlation of methylation and GE was
414 more common than direct correlation (i.e. high methylation and high GE; N=17).
415 However, direct correlation was more common in patients with a recent smoking history
416 (70.5%). This pattern was also observed for **Cis2**, which best predicted recurrence in
417 smokers. In cases where methylation and GE are directly correlated, we propose that
418 smoking may alter expression by other effects, e.g. mutagenic, rather solely than by
419 epigenetic inactivation through methylation.

420 **DISCUSSION**

421 Using gene expression signatures, we derived both GI_{50} threshold-dependent
422 and -independent ML models which predict the chemotherapy responses for cisplatin,
423 carboplatin and oxaliplatin, respectively. The cisplatin gene signature **Cis1**
424 (Supplementary Table S6A) most accurately predicted response in bladder cancer
425 patients after 18 months, and **Car1** (Suppl. Table S6B) best predicted response in
426 ovarian cancer patients after 4 years. **Oxa1** (Suppl. Table S6C) more accurately
427 predicted disease-free patients than recurrent disease at the one year treatment
428 threshold. The thresholds which best represented time-to-recurrence differed between
429 the platin drugs in each cancer type. Cisplatin gene signatures had noticeably improved
430 performance when smoking history was taken into account.

431 The three platin drugs produce distinctly different gene signatures. Initial gene
432 sets exhibited some overlap between platin drugs (N=67 between any two platins), but
433 very few of these were correlated by MFA of GI_{50} with multiple platin drugs (*ATP7B*,
434 *BCL2* and *MSH2*). *BCL2L1*, *GSTP1*, *MAP3K1*, *MAPK3*, *MT1A*, and *MT2* were genes
435 common to multiple platin drugs whose expression was correlated with cisplatin GI_{50}
436 values but not with carboplatin and/or oxaliplatin values. Similarly, genes correlating only
437 to carboplatin GI_{50} included *AKT1*, *EGF*, *ERCC1*, *KRAS*, *LIG3*, *MTHFR*, *MTR*, *RAD50*,
438 *TP53*, while genes correlating to only oxaliplatin GI_{50} included *ATM*, *BCL2*, *CLCN6*,
439 *ERCC2*, *ERCC6*, and *UGT1A1*. Despite the close similarity between cisplatin and
440 carboplatin GI_{50} response (see Figure 2), only one gene (*ATP7B*) was related by MFA to
441 GI_{50} levels of both drugs. *BCL2* and *MSH2* correlated with both cisplatin and oxaliplatin
442 GI_{50} (*BCL2* did not correlate with carboplatin GI_{50}). The increase in misclassification
443 caused by the elimination of *MSH2* from any gene signature in which it was present was
444 significant; for example, misclassification of **Cis14** and **Oxa21** (Table 1) were increased

445 by 28.2% and 19.1%, respectively (Suppl. Tables S2A and S2C). These differences may
446 reflect the spectrum of activity, sensitivity, and toxicity of these signature genes^{21–23,35,36}.

447 Previous validation of patient data for other drugs validated with other datasets^{6,24}
448 using biochemically inspired machine learning have had better performance than those
449 reported here. We investigated the possibility that disease and molecular heterogeneity
450 in platin-treated patients may have affected the accuracy of our results. Model
451 predictions were reevaluated after stratifying clinical features such as time-to-disease
452 recurrence, cancer stage, and metastatic lymph node count. Breast cancer patients with
453 advanced disease (stage III and IV) were analyzed separately from those with earlier
454 stage diagnoses (stage I and II). Cisplatin gene signature **Cis1** performed best on stage
455 IV patients (overall accuracy 72.4% at a 2 year recurrence threshold), while **Oxa1**
456 similarly performed best in predicting late stage cancers (74.5% accurate for stage III
457 and 71.4% accurate for stage IV at a 2 year recurrence threshold). **Cis5** was also more
458 accurate for later stage cancer patients (72.4% overall accuracy at 18 months). The
459 accuracies of gene signatures were similar across all stages (e.g. **Car1** ranged from 58-
460 74%). Cisplatin-treated, TCGA bladder cancer patients and oxaliplatin-treated TCGA
461 colorectal cancer patients were also stratified by Lymph Node status (N0, N1, and N2
462 [bladder cancer patient data set comprised of only two N3 patients, which were included
463 with the analysis of N2 patients; N3 was not represented in colorectal cancer]). In TCGA
464 bladder cancer patients, **Cis1** exhibited ~60% accuracy across all categories; however it
465 performed better in sensitive N0 and N1 patients relative to N2. **Cis2** was less accurate
466 for N2 patients than for N0 and N1. Sensitive N2 patients were more likely to be
467 misclassified (<40%) than relapsed N2 patients. In TCGA colorectal cancer patients,
468 **Oxa1** was 88% accurate in N2 patients (95% accurate for sensitive N2 patients [n=19],
469 and 67% accurate for relapsed N2 patients [n=6]). Oxaliplatin gene signatures were less

470 accurate for N1 patients compared to N0 and N2. Thus, heterogeneity in disease stage
471 as well as metastatic phenotypes adversely confounds the overall accuracies of our
472 predictions.

473 Gene signature models derived from cell lines and tested on patients differ in
474 their outcome measures. The exact GI_{50} cell line threshold that is most predictive of
475 patient outcome is not known, and different groups use different methods to discretize
476 GI_{50} values^{37,38}. Therefore, we developed ML models for platin drugs which predict drug
477 response without relying on arbitrary GI_{50} thresholds. For cisplatin, SVM ensemble
478 averaging generated on different resistance thresholds shows a small increase in
479 accuracy over most gene signatures, better representing the sensitive, disease-free
480 class (59% accuracy). Interestingly, ensemble averaging of only the gene signatures
481 built using a moderate GI_{50} thresholds yielded results which better represented the
482 resistance class. This result closer matches the accuracy of **Cis1**, and may be due to
483 **Cis1** having a greater overall impact on the ensemble prediction. When limiting
484 ensemble averaging to only those gene signatures with the highest area under the curve
485 (AUC) at each resistance threshold, differences in predictions were negligible. Ensemble
486 ML can potentially avoid problems with poor performance and overfitting by combining
487 gene signatures that individually perform slightly better than chance³⁹.

488 It is difficult to reconcile gene signatures without features known to be related to
489 chemoresistance with tumor biology. Our thresholding approach may reveal potentially
490 important genes and pathways associated with platin resistance. It would be preferable
491 to explore pathways related to signature genes to improve accuracy, identify potential
492 targets for further study of chemoresistance, and expand the model parameters to take
493 into account alternate states besides those captured in the original signature⁴⁰.
494 Signatures for resistance may be useful for developing targeted intervention to re-

495 sensitize tumours. For example, the mismatch repair (MMR) gene *MSH2* is commonly
496 present in gene signatures at high resistance levels for oxaliplatin, which is of interest,
497 as MMR deficiency has been shown to be predictive for oxaliplatin resistance³⁶. Indeed,
498 *MLH1*, *MSH2* and *MSH6*-deficient cells are more susceptible to oxaliplatin, despite
499 MMR-deficiency being associated with cisplatin resistance³⁵. The autoimmune disease-
500 associated gene *SIAE*, which has been previously shown to have a strong negative
501 correlation to oxaliplatin response in advanced CRC patients⁴¹, was selected in the
502 majority of thresholded oxaliplatin gene signatures (Supplementary Table S2C). The
503 gene *BCL2*, which was commonly selected for cisplatin (Figure 3A), was rarely selected
504 for oxaliplatin (Figure 3C). At the highest levels of resistance to cisplatin, gene
505 signatures were enriched for genes belonging to DNA repair, anti-oxidative response,
506 apoptotic pathways and drug transporters (Figure 3A). These gene pathways are known
507 to be involved in cisplatin resistance^{42,43} and these specific genes may be explored in
508 subsequent work to identify the contribution to chemotherapy response in a biochemical
509 context.

510 Log-loss evaluates the accuracy of a classifier by penalizing erroneous
511 classifications, and is relevant in cases where data is imbalanced and/or have an
512 unequally distributed error cost. We assessed whether ML gene signatures based on
513 log-loss minimization could improve accuracy to predicting patient response
514 (Supplementary Table S3) and compared them to gene signatures generated by
515 minimizing cell line misclassification. When gene signatures generated by both methods
516 were highly similar (generated at the same GI₅₀ threshold, consist of a similar number of
517 genes and consist of ≥ 80% shared genes), prediction accuracy of TCGA cancer patient
518 outcomes were nearly indistinguishable, as accuracy can vary over different relapse
519 thresholds. Where significant differences in predictions were seen, the misclassification-
520 based gene signatures were more accurate overall (**Cis1**, **Cis17** and the “12-Resistant”

521 carboplatin gene signature were +8.3%, +5.6% and +3.9% more accurate compared to
522 the log-loss gene signature, respectively). Oxaliplatin gene signatures were dissimilar
523 across all GI_{50} thresholds, as the log-loss minimized ML gene signatures often contain
524 increased numbers of genes compared to the misclassification-based gene signatures.
525 Many of these larger gene signatures were less accurate in patients compared to gene
526 signatures which minimized misclassification rates consistent that this evaluation and
527 model selection method is more prone to overfitting. This pattern was also noted for
528 gene signatures generated at extreme GI_{50} thresholds for all three platin drugs in which
529 response was, by definition, somewhat imbalanced.

530 It may be feasible to predict responses to combination chemotherapy with the
531 gene signatures described here. Not included in the present analysis were signatures for
532 methotrexate, vinblastine, and doxorubicin, which comprise the MVAC cocktail used to
533 treat bladder cancer. This was due primarily to a lack of patients treated with this drug
534 combination in the TCGA bladder dataset (N=11). Individual signatures for several of
535 these drugs have been derived and analyzed using the patient data from METABRIC
536 (Molecular Taxonomy of Breast Cancer International Consortium)²⁴. A reasonable
537 approach to predicting combination chemotherapy would first determine the probability
538 of sensitivity or resistance to individual drugs, accounting for the misclassification rate by
539 each (defined as d_1, \dots, d_k). The ML classifiers output these probabilities, analogous to
540 their misclassification rates in a set of patients treated identically. If the model predicts
541 that the patient is sensitive to drug d_1 with 90% probability, and sensitive to drug d_2 with
542 5% probability, and the errors are independent, then the probability of sensitivity to the
543 combination is $1 - (1 - 0.9)(1 - 0.05) = 90.5\%$, and the probability of resistance is 9.5%,
544 assuming no synergistic effects between drugs. If interaction or dependence among
545 errors is suspected, the combined probability of resistance to the pair d_{12} could be
546 estimated based on the features that are shared by the signatures of both drugs. The

547 probability of sensitivity to the combination would then be given by $1 - (1 - d_{12}) \cdot (1 -$
548 $d_3) \cdot \dots \cdot (1 - d_k)$.

549 The predictive accuracy for the same gene signature could differentiate highly
550 between the two datasets. **Cis3** (Supplemental Table S6A) had an AUC of 0.64 when
551 validated against TCGA bladder cancer patients. However, the AUC was lower when
552 applied to the Als *et al.*²⁹ dataset (AUC=0.18). Patient metadata in the latter study only
553 indicated patient survival times, while we base the expected TCGA patient outcome on
554 time to disease recurrence. As the basis of our expected outcome differs between
555 datasets, these differences may be acting as a confounding factor to determine accuracy
556 of gene signatures. The datasets also differ in how expression was measured
557 (microarray vs. RNA-seq). The relevance of gene signatures based on training and
558 testing data from different platforms can affect the accuracy of validation, which might
559 not be improved by data normalization. In this study, datasets were subjected to z-score
560 normalization. In subsequent studies, other techniques to correct for some of these
561 effects have been described and could be applied⁴⁴.

562 In summary, we describe GI₅₀- or IC₅₀-threshold-independent ML gene signatures
563 to predict chemotherapy response to platin agents in cancer patients. Ensemble
564 machine learning produced combined signatures that were more accurate than most
565 individual gene signatures generated with different thresholds. Genes associated
566 cisplatin response included those which exacerbate resistance in patients with a history
567 of smoking. The methodology described here should be adaptable to other drugs and
568 cancer types. With a range of gene signatures for multiple drugs, it may be possible to
569 improve the efficacy of treatment by tailoring treatment to a patient's specific tumour
570 biology, and reduce treatment duration by limiting the number of different therapeutic
571 regimens prescribed before achieving a successful response⁴⁵.

572 MATERIALS AND METHODS

573 Data and preprocessing

574 *Cell-line Data Sets*

575 Microarray GE and data from breast cancer cell lines were used to train
576 ML-based gene signatures of drug response based on respective growth or target
577 inhibition data (GI_{50} or IC_{50}). Cell lines were treated with either cisplatin (N=39),
578 carboplatin (N=46), or oxaliplatin (N=47)¹³. Bladder cancer cell line GE and IC_{50}
579 measurements for cisplatin were obtained from cancerRxgene (N=17). However, all
580 models (gene signatures) used to evaluate patient data were trained on breast cancer
581 cell line data, because the number of bladder cancer cell lines was insufficient to
582 produce accurate signatures.

583 *Cancer Patient Data Sets*

584 RNA-seq GE and survival measurements were downloaded from TCGA for
585 bladder urothelial carcinoma (N=72 patients treated with cisplatin)²⁶, ovarian epithelial
586 tumor (N=410 treated with carboplatin)²⁷ and colorectal adenocarcinoma (N=99 treated
587 with oxaliplatin)²⁸. GE of cisplatin-treated patients of cell carcinoma of the urothelium
588 (N=30)²⁹ and for oxaliplatin-treated CRC patients (N=83)³⁰ were obtained from the Gene
589 Expression Omnibus. Clinical metadata and GE for TCGA patients were obtained from
590 Genomic Data Commons (<https://gdc.cancer.gov/>), while methylation HM450 (Illumina)
591 data for these patients was downloaded from cBioPortal⁴⁶.

592 *Development and Pre-Processing of Biochemically-Inspired Gene Sets*

593 Initial gene sets for developing signatures for each drug were identified from
594 previously published literature (see Supplemental References, Section B) and

595 databases, such as PharmGKB and DrugBank^{47,48}. The evidence supporting each gene
596 contained in the final signatures is independent scientific evidence that the genes
597 selected are not the result of spurious associations. The final gene sets were chosen
598 using MFA with the breast cancer cell-line data to analyze interactions between GE, CN,
599 and GI₅₀ data for the drug of interest⁴⁹. Genes whose GE and/or CN showed a direct or
600 inverse correlation with GI₅₀ were selected for SVM training. Because the number of
601 genes related to GI₅₀ for oxaliplatin exceeded the number of cell lines available for
602 training, we limited the input to the oxaliplatin ML model to those genes whose GE were
603 related to GI₅₀. Similarly, the number correlating genes in cisplatin treated cells
604 exceeded the number of cell lines. For cisplatin, genes whose expression correlated with
605 GI₅₀ were eliminated if they showed no or little expression in TCGA bladder cancer
606 patients (i.e. RNA-seq counts by Expectation Maximization [RSEM] were < 5.0 for
607 majority of individuals). This reduces the overall number of genes for SVM analysis, and
608 thus helps to avoid a data to size sample imbalance. For cisplatin, MFA was repeated
609 using IC₅₀ values for 17 bladder cancer cell lines; however, the available CN data
610 generally showed a lack of variation in the cell lines for these genes. Instead, the
611 available IC₅₀ values for three other cancer drugs (doxorubicin, methotrexate and
612 vinblastine) were compared with the IC₅₀ of cisplatin by MFA.

613 Applying an SVM model directly to patient data without a normalization approach
614 is imprecise when training and testing data are not obtained using similar methodology
615 (i.e. different microarray platforms). To compare the cell line GE microarray data and the
616 patient RNA-seq GE datasets, expression values were normalized by conversion to z-
617 scores using MATLAB⁵⁰. Although Log₂ intensity values from microarray data were not
618 available for TCGA samples, RNA-seq based GE and log₂ intensities from microarray
619 data are highly correlated⁵¹.

620 **Machine Learning**

621 SVMs were trained with breast cancer cell line GE datasets¹³ with the Statistics
622 Toolbox in MATLAB⁵⁰ similar to Dorman et al (2016)⁶. Rather than a linear kernel, we
623 used a Gaussian kernel function (fitcsvm), and then tested with leave-one-out cross-
624 validation (using the options 'crossval' and 'leaveout'). A greedy backwards feature
625 selection algorithm was used to improve classification accuracy⁵². BFS leaves out
626 individual genes from the initial MFA-qualified gene set, then trains a cross validated
627 Gaussian kernel SVM on the training samples, removing the gene with the highest
628 misclassification rate. The procedure is repeated until all genes have been evaluated.
629 The gene subset with the lowest misclassification rate⁶ or log-loss statistic²⁵ based on
630 cross-validation is selected as the gene signature for subsequent testing with patient GE
631 and clinical data. K-fold cross validation of the misclassification-based gene signatures
632 was performed using MATLAB software described in Zhao et al. (2018)²⁵.

633 SVMs minimized according to the log-loss classification function were also
634 generated with both software described in Zhao *et al.* (2018; uses multiclass compatible
635 'fitcecoc' function)²⁵, and with a modified version of the software described above (using
636 'fitSVMPosterior' to compute posterior probabilities). Computed probabilities differ
637 between 'fitSVMPosterior' and 'fitcecoc' (range: 0.02-0.04), thus the resultant gene
638 signatures will differ between the two programs. When given unbalanced data (e.g.
639 lower resistance thresholds), 'fitSVMPosterior' will warn that some classes are not
640 represented, and thus those folds will not predict the labels for those missing classes.
641 The log-loss gene signatures described in this manuscript were generated with the
642 multiclass compatible 'fitcecoc' function software²⁵.

643 *Derivation of gene signatures for different drug resistance thresholds*

644 We have previously set a conventional GI_{50} threshold distinguishing sensitivity
645 from resistance at the *median* of the range of drug concentrations that inhibited cell
646 growth by 50%⁶. We hypothesized that different gene signatures could be derived for
647 different levels of drug resistance by varying this threshold. ML experiments for
648 classifying resistance or sensitivity at GI_{50} values generated a series of optimized
649 Gaussian SVM gene signatures whose performance were assessed with patient
650 expression data for each signature. A heat map which illustrates the frequencies of
651 genes appearing in these gene signatures was created with the R language *hist2d*
652 function.

653 A composite gene signature was created by ensemble averaging of all gene
654 signatures generated at each resistance threshold. Ensemble averaging combines
655 signatures through averaging the weighted accuracy of a set of related models³⁹. The
656 decision function for the ensemble classifier is the mean of the decision function scores
657 of the component classifiers, weighted by the AUC.

658 *Significance of cell line-derived gene signatures*

659 The significance of the derived SVMs (whether the observed performance of the
660 gene signatures could have arisen by chance) was first assessed by permutation
661 analysis with randomized cell line labels and with random sets of genes, as described
662 previously⁶. Using the median cisplatin GI_{50} as the resistance threshold, 10,000 gene
663 signatures based on random gene selection (15 genes) had higher rates of
664 misclassification than the best median SVM gene signatures (2 signatures with 7.7%
665 misclassification). Cisplatin, carboplatin and oxaliplatin GE data for random cell line label
666 combinations (n=10,000) generated only 8, 1 and 1 signatures, respectively, with lower
667 error rates than the best biochemically-inspired signatures. When minimizing for log-loss

668 (rather than misclassification), random gene analysis (10,000 iterations; median cisplatin
669 GI_{50} threshold) resulted only in gene signatures with a higher log-loss than the signature
670 generated with the initial cisplatin gene set. Log-loss based random label analysis
671 ($n=2000$ combinations) resulted in 3.4% of random label gene signatures resulted in a
672 lower log-loss than the cisplatin signature at the same GI_{50} threshold (5.27). This was not
673 entirely surprising, since this result depends on the GI_{50} threshold used for labeling. The
674 differences between GI_{50} values for cell lines close to the median GI_{50} used in this
675 analysis are almost negligible (e.g. 5.11 vs 5.12) and likely within the measurement error
676 for these values.

677 Regarding the specificity of the cisplatin gene signatures, the best performing
678 cisplatin gene signatures (**Cis1** and **Cis2**) were used to evaluate participants who were
679 treated with other drugs (using an 18 months post-treatment threshold). In such patients,
680 36.5% of those who were disease-free were predicted accurately with the **Cis1** signature
681 ($N=178$; 22% less accurate than platin-treated patients), and 62.9% accurate for those
682 with recurrent disease ($N=70$; 8.1% less accurate). **Cis2** was 43.8% accurate for
683 disease-free non-platin treated patients ($N=178$; 12.3% lower accuracy), and 60.0% of
684 those who relapsed ($N=70$; 2.9% less accurate). GE changes in patients treated with
685 platin drugs are better modeled by cancer cell-line based predictors than in patients
686 receiving other drug treatments.

687 **ACKNOWLEDGEMENTS**

688 Katherina Baranova contributed to the cisplatin gene signatures and Dimo Angelov
689 developed automated feature selection. We thank Murray Junop for commenting on the
690 manuscript. Compute Canada and Shared Hierarchical Academic Research Computing
691 Network (SHARCNET) provided high performance computing and storage facilities.

692 **CONFLICTS OF INTEREST**

693 PKR cofounded CytoGnomix Inc., which hosts the interactive resource described in this
694 study for prediction of responses to chemotherapy agents. The other authors have no
695 conflicts of interest.

696 **AUTHOR CONTRIBUTIONS**

697 PKR and DL designed the methodology. EJM and JZ performed analyses. EJM and
698 PKR wrote the manuscript.

699 **FUNDING**

700 PKR is supported by NSERC (RGPIN-2015-06290), Canadian Foundation for
701 Innovation, Canada Research Chairs, and CytoGnomix.

702

703 Supplementary information accompanies the manuscript on the *Signal Transduction and*
704 *Targeted Therapy* website <http://www.nature.com/sigtrans>.

705 **REFERENCES**

- 706 1. Cardoso, F. *et al.* Locally recurrent or metastatic breast cancer: ESMO Clinical Practice
707 Guidelines for diagnosis, treatment and follow-up. *Ann. Oncol.* **23**, vii11–vii19 (2012).
- 708 2. Oostendorp, L. J., Stalmeier, P. F., Donders, A. R. T., van der Graaf, W. T. & Ottevanger, P. B.
709 Efficacy and safety of palliative chemotherapy for patients with advanced breast cancer
710 pretreated with anthracyclines and taxanes: a systematic review. *Lancet Oncol.* **12**, 1053–
711 1061 (2011).

- 712 3. Alfarouk, K. O. *et al.* Resistance to cancer chemotherapy: failure in drug response from
713 ADME to P-gp. *Cancer Cell Int.* **15**, 71 (2015).
- 714 4. Gąsowska-Bodnar, A. *et al.* Survivin Expression as a Prognostic Factor in Patients With
715 Epithelial Ovarian Cancer or Primary Peritoneal Cancer Treated With Neoadjuvant
716 Chemotherapy: *Int. J. Gynecol. Cancer* **24**, 687–696 (2014).
- 717 5. Hatzis, C. *et al.* A genomic predictor of response and survival following taxane-anthracycline
718 chemotherapy for invasive breast cancer. *JAMA* **305**, 1873–1881 (2011).
- 719 6. Dorman, S. N. *et al.* Genomic signatures for paclitaxel and gemcitabine resistance in breast
720 cancer derived by machine learning. *Mol. Oncol.* **10**, 85–100 (2016).
- 721 7. Zhang, S. *et al.* Organic Cation Transporters Are Determinants of Oxaliplatin Cytotoxicity.
722 *Cancer Res.* **66**, 8847–8857 (2006).
- 723 8. Poisson, L. M. *et al.* A metabolomic approach to identifying platinum resistance in ovarian
724 cancer. *J. Ovarian Res.* **8**, (2015).
- 725 9. Cadoná, F. C. *et al.* Guaraná a Caffeine-Rich Food Increases Oxaliplatin Sensitivity of
726 Colorectal HT-29 Cells by Apoptosis Pathway Modulation. *Anticancer Agents Med. Chem.* **16**,
727 1055–1065 (2016).
- 728 10. Kasparkova, J., Vojtiskova, M., Natile, G. & Brabec, V. Unique Properties of DNA Interstrand
729 Cross-Links of Antitumor Oxaliplatin and the Effect of Chirality of the Carrier Ligand. *Chem. –*
730 *Eur. J.* **14**, 1330–1341 (2008).
- 731 11. Woynarowski, J. M. *et al.* Oxaliplatin-Induced Damage of Cellular DNA. *Mol. Pharmacol.* **58**,
732 920–927 (2000).
- 733 12. Tashiro, T., Kawada, Y., Sakurai, Y. & Kidani, Y. Antitumor activity of a new platinum
734 complex, oxalato (trans-1,2-diaminocyclohexane)platinum (II): new experimental data.
735 *Biomed. Pharmacother.* **43**, 251–260 (1989).

- 736 13. Daemen, A. *et al.* Modeling precision treatment of breast cancer. *Genome Biol.* **14**, R110
737 (2013).
- 738 14. Yuan, Y. *et al.* Identification of the biomarkers for the prediction of efficacy in first-line
739 chemotherapy of metastatic colorectal cancer patients using SELDI-TOF-MS and artificial
740 neural networks. *Hepatogastroenterology.* **59**, 2461–2465 (2012).
- 741 15. L'Espérance, S., Bachvarova, M., Tetu, B., Mes-Masson, A.-M. & Bachvarov, D. Global gene
742 expression analysis of early response to chemotherapy treatment in ovarian cancer
743 spheroids. *BMC Genomics* **9**, 99 (2008).
- 744 16. Nickerson, M. L. *et al.* Molecular analysis of urothelial cancer cell lines for modeling tumor
745 biology and drug response. *Oncogene* (2016). doi:10.1038/onc.2016.172
- 746 17. Yuryev, A. Gene expression profiling for targeted cancer treatment. *Expert Opin. Drug*
747 *Discov.* **10**, 91–99 (2015).
- 748 18. Sataloff, D. M. *et al.* Pathologic response to induction chemotherapy in locally advanced
749 carcinoma of the breast: a determinant of outcome. *J. Am. Coll. Surg.* **180**, 297–306 (1995).
- 750 19. Ogston, K. N. *et al.* A new histological grading system to assess response of breast cancers
751 to primary chemotherapy: prognostic significance and survival. *Breast Edinb. Scotl.* **12**, 320–
752 327 (2003).
- 753 20. Earl, J. *et al.* The UBC-40 Urothelial Bladder Cancer cell line index: a genomic resource for
754 functional studies. *BMC Genomics* **16**, 403 (2015).
- 755 21. Rixe, O. *et al.* Oxaliplatin, tetraplatin, cisplatin, and carboplatin: Spectrum of activity in drug-
756 resistant cell lines and in the cell lines of the national cancer institute's anticancer drug
757 screen panel. *Biochem. Pharmacol.* **52**, 1855–1865 (1996).
- 758 22. Mehmood, R. K. Review of Cisplatin and oxaliplatin in current immunogenic and monoclonal
759 antibody treatments. *Oncol. Rev.* **8**, 256 (2014).

- 760 23. Kweekel, D. M., Gelderblom, H. & Guchelaar, H.-J. Pharmacology of oxaliplatin and the use
761 of pharmacogenomics to individualize therapy. *Cancer Treat. Rev.* **31**, 90–105 (2005).
- 762 24. Mucaki, E. J. *et al.* Predicting Outcomes of Hormone and Chemotherapy in the Molecular
763 Taxonomy of Breast Cancer International Consortium (METABRIC) Study by Biochemically-
764 inspired Machine Learning. *F1000Research* **5**, 2124 (2017).
- 765 25. Zhao, J. Z. L., Mucaki, E. J. & Rogan, P. K. Predicting ionizing radiation exposure using
766 biochemically-inspired genomic machine learning. *F1000Research* **7**, 233 (2018).
- 767 26. Tembe, V. *et al.* The BARD1 BRCT domain contributes to p53 binding, cytoplasmic and
768 mitochondrial localization, and apoptotic function. *Cell. Signal.* **27**, 1763–1771 (2015).
- 769 27. Robertson, A. G. *et al.* Comprehensive Molecular Characterization of Muscle-Invasive
770 Bladder Cancer. *Cell* **171**, 540–556.e25 (2017).
- 771 28. Cancer Genome Atlas Research Network. Integrated genomic analyses of ovarian
772 carcinoma. *Nature* **474**, 609–615 (2011).
- 773 29. Cancer Genome Atlas Network. Comprehensive molecular characterization of human colon
774 and rectal cancer. *Nature* **487**, 330–337 (2012).
- 775 30. Als, A. B. *et al.* Emmprin and survivin predict response and survival following cisplatin-
776 containing chemotherapy in patients with advanced bladder cancer. *Clin. Cancer Res. Off. J.*
777 *Am. Assoc. Cancer Res.* **13**, 4407–4414 (2007).
- 778 31. Tsuji, S. *et al.* Potential responders to FOLFOX therapy for colorectal cancer by Random
779 Forests analysis. *Br. J. Cancer* **106**, 126–132 (2012).
- 780 32. Freedman, N. D., Silverman, D. T., Hollenbeck, A. R., Schatzkin, A. & Abnet, C. C. Association
781 between smoking and risk of bladder cancer among men and women. *JAMA* **306**, 737–745
782 (2011).

- 783 33. Fortin, A., Wang, C. S. & Vigneault, E. Influence of smoking and alcohol drinking behaviors
784 on treatment outcomes of patients with squamous cell carcinomas of the head and neck.
785 *Int. J. Radiat. Oncol. Biol. Phys.* **74**, 1062–1069 (2009).
- 786 34. Joehanes, R. *et al.* Epigenetic Signatures of Cigarette Smoking. *Circ. Cardiovasc. Genet.*
787 (2016). doi:10.1161/CIRCGENETICS.116.001506
- 788 35. Raymond, E., Faivre, S., Chaney, S., Woynarowski, J. & Cvitkovic, E. Cellular and Molecular
789 Pharmacology of Oxaliplatin1. *Mol. Cancer Ther.* **1**, 227–235 (2002).
- 790 36. Alex, A. K. *et al.* Response to Chemotherapy and Prognosis in Metastatic Colorectal Cancer
791 With DNA Deficient Mismatch Repair. *Clin. Colorectal Cancer* (2016).
792 doi:10.1016/j.clcc.2016.11.001
- 793 37. Sos, M. L. *et al.* Predicting drug susceptibility of non-small cell lung cancers based on genetic
794 lesions. *J. Clin. Invest.* **119**, 1727–1740 (2009).
- 795 38. Laderas, T. G., Heiser, L. M. & Sönmez, K. A Network-Based Model of Oncogenic
796 Collaboration for Prediction of Drug Sensitivity. *Front. Genet.* **6**, (2015).
- 797 39. Clemen, R. T. Combining forecasts: A review and annotated bibliography. *Int. J. Forecast.* **5**,
798 559–583 (1989).
- 799 40. Airley, R. *Cancer chemotherapy*. (Wiley-Blackwell, 2009).
- 800 41. Li, X.-X. *et al.* RNA-seq identifies determinants of oxaliplatin sensitivity in colorectal cancer
801 cell lines. *Int. J. Clin. Exp. Pathol.* **7**, 3763–3770 (2014).
- 802 42. Borst, P., Rottenberg, S. & Jonkers, J. How do real tumors become resistant to cisplatin? *Cell*
803 *Cycle Georget. Tex* **7**, 1353–1359 (2008).
- 804 43. Wernyj, R. & Morin, P. Molecular mechanisms of platinum resistance: still searching for the
805 Achilles heel. *Drug Resist. Updat.* **7**, 227–232 (2004).

- 806 44. Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data
807 using empirical Bayes methods. *Biostat. Oxf. Engl.* **8**, 118–127 (2007).
- 808 45. Akamatsu, N., Nakajima, H., Ono, M. & Miura, Y. Increase in acetyl CoA synthetase activity
809 after phenobarbital treatment. *Biochem. Pharmacol.* **24**, 1725–1727 (1975).
- 810 46. Gao, J. *et al.* Integrative analysis of complex cancer genomics and clinical profiles using the
811 cBioPortal. *Sci. Signal.* **6**, pl1 (2013).
- 812 47. Whirl-Carrillo, M. *et al.* Pharmacogenomics Knowledge for Personalized Medicine. *Clin.*
813 *Pharmacol. Ther.* **92**, 414–417 (2012).
- 814 48. Law, V. *et al.* DrugBank 4.0: shedding new light on drug metabolism. *Nucleic Acids Res.* **42**,
815 D1091–D1097 (2014).
- 816 49. Abdi, H. & Williams, L. J. Principal component analysis. *Wiley Interdiscip. Rev. Comput. Stat.*
817 **2**, 433–459 (2010).
- 818 50. MATLAB and Statistics Toolbox Release 2012b, The MathWorks, Inc., Natick, Massachusetts,
819 United States.
- 820 51. Marioni, J. C., Mason, C. E., Mane, S. M., Stephens, M. & Gilad, Y. RNA-seq: an assessment
821 of technical reproducibility and comparison with gene expression arrays. *Genome Res.* **18**,
822 1509–1517 (2008).
- 823 52. Bermingham, M. L. *et al.* Application of high-dimensional feature selection: evaluation for
824 genomic prediction in man. *Sci. Rep.* **5**, 10312 (2015).

825

827 **Table 1: Gene Signatures Which Best Predicted Response in TCGA Cancer Patients**

Gene Signature ID	Cancer Type Tested	GI50 Threshold	Signature (C, σ^*)
Cis1 (Cisplatin)	Bladder	5.11	<i>BARD1, BCL2, BCL2L1, CDKN2C, FAAP24, FEN1, MAP3K1, MAPK13, MAPK3, NFKB1, NFKB2, SLC22A5, SLC31A2, TLR4, TWIST1</i> (100000, 100)
Cis2 (Cisplatin)	Bladder	5.12	<i>BARD1, BCL2L1, CFLAR, FAAP24, MAP3K1, MAPK3, NFKB1, POLQ, PRKAA2, SLC22A5, SLC31A2, TLR4, TWIST1</i> (10000, 100)
Cis3 (Cisplatin)	Bladder	5.60	<i>BCL2, CFLAR, ERCC2, ERCC6, FAAP24, FEN1, MAP3K1, NFKB1, NFKB2, PNKP, POLQ, PRKCB, SLC22A5, SNAI1, TLR4</i> (100000, 100)
Cis12 (Cisplatin)	Bladder	5.40	<i>ATP7B, BCL2, BCL2L1, CDKN2C, ERCC2, FAAP24, GSTO1, MAP3K1, MAPK3, MT2A, NFKB1, NFKB2, POLD1, POLQ, PRKCB, SNAI1, TLR4, TP63</i> (10000, 100)
Cis14 (Cisplatin)	Bladder	5.16	<i>BARD1, BCL2, BCL2L1, CDKN2C, FAAP24, FEN1, FOS, GSTP1, MAP3K1, MAPK13, MAPK3, MSH2, NFKB1, POLD1, POLQ, PRKAA2, PRKCB, SLC22A5, SLC31A2, SNAI1, TWIST1</i> (10000, 100)
Cis17 (Cisplatin)	Bladder	5.10	<i>ATP7B, BCL2, BCL2L1, FEN1, GSTP1, MAP3K1, MAPK3, MT2A, NFKB1, PNKP, POLQ, PRKAA2, PRKCB, SLC31A2, TLR4, TP63</i> (100000, 100)
Car1 (Carboplatin)	Ovarian	4.22	<i>AKT1, EIF3K, ERCC1, GNGT1, GSR, MTHFR, NEDD4L, NLRP1, NRAS, RAF1, SGK1, TIGD1, TP53, VEGFB, VEGFC</i> (100000, 100)
Car9 (Carboplatin)	Ovarian	4.32	<i>AKT1, ATP7B, EIF3I, ETS2, GNGT1, HRAS, KRAS, LIG3, MTHFR, MTR, NRAS, RAD50, SCN10A, TIGD1, TP53, VEGFB</i> (10000, 100)
Car51 (Carboplatin)	Ovarian	4.34	<i>AKT1, EGF, EIF3I, ERCC1, ETS2, GNGT1, KRAS, MTHFR, MTR, NEDD4L, NLRP1, NRAS, RAD50, RAF1, SGK1, TIGD1, TP53, VEGFB, VEGFC</i> (10000, 100)
Car73 (Carboplatin)	Ovarian	4.09	<i>AKT1, ATP7B, ETS2, GNGT1, HRAS, NLRP1, SCN10A, VEGFB</i> (100000, 1000)
Oxa1 (Oxaliplatin)	Colorectal	5.10	<i>BRAF, FCGR2A, IGF1, MSH2, NAGK, NFE2L2, NQO1, PANK3, SLC47A1, SLCO1B1, UGT1A1</i> (10, 10)
Oxa21 (Oxaliplatin)	Colorectal	5.10	<i>BRAF, IGF1, IGF1R, KLF3, MSH2, NAT2, NFE2L2, NQO1, PANK3, PRSS1, SIAE, SLC47A1, SLCO1B1, UGT1A1</i> (1000, 100)
*C - The box-constraint. σ - the kernel-scale ("sigma"). Bolded gene signatures are those that best overall performance against TCGA cancer patient gene expression data.			

829

830 **FIGURE LEGENDS**

831 **Figure 1.** Schematic of platinum drug sensitivity and resistance genes which showed
832 MFA correlation for GI_{50} of A) cisplatin, B) carboplatin, and C) oxaliplatin. The genes
833 used to derive the SVM are shown in context of their effect in the cell and role in cisplatin
834 mechanisms of action. GE and CN correlation with inhibitory drug concentration by MFA
835 of breast (GI_{50}) and bladder (IC_{50}) cancer cell line data.

836 **Figure 2:** GI_{50} values for cell lines treated with the three platin drugs were plotted in
837 order of ascending oxaliplatin GI_{50} . For most cell lines, there is a visible trend between
838 the GI_{50} for cisplatin and carboplatin, reflecting the correlation between the two drugs
839 seen by MFA. Despite this correlation, carboplatin shows a much smaller variance (0.22)
840 compared to cisplatin (0.37; oxaliplatin variance is 0.34).

841 **Figure 3.** The variation in gene composition of misclassification-based SVMs at different
842 GI_{50} thresholds for A) cisplatin, B) carboplatin, and C) oxaliplatin. GI_{50} intervals are
843 indicated on the left, with the number of cell lines with GI_{50} values within said intervals in
844 brackets. Each box represents the density of genes appearing in optimized Gaussian
845 SVM gene signatures in those functional categories, with darker grey indicating frequent
846 genes in indicated GI_{50} threshold intervals, while lighter grey indicates less commonly
847 selected genes. The number of thresholded gene signatures used to derive the density
848 plot within each interval is equal (or greater, in the case of multiple equally performing
849 gene signatures) to the number of cell lines within that GI_{50} interval.

850 **Supplementary Figure 1.** The variation in gene composition of log-loss based SVMs at
851 different GI_{50} thresholds for A) cisplatin, B) carboplatin, and C) oxaliplatin. Each box
852 represents the density of genes appearing in optimized Gaussian log-loss SVM gene

853 signatures in those functional categories, with darker grey indicating frequent genes in
854 indicated GI_{50} threshold intervals, while lighter grey indicates less commonly selected
855 genes.

856 **Supplementary Figure 2.** Classification accuracy of gene signatures on TCGA bladder
857 cancer patients treated with cisplatin and/or carboplatin as the resistance threshold is
858 varied. Recurrence and disease-free survival are used as a binary measure to assess
859 performance. The x-axis indicates movement of the resistance threshold, with more cell
860 lines labeled sensitive on the left and more labeled resistant on the right. Maximal AUC
861 is indicated by the downward arrows.

862 **Supplementary Figure 3.** Classification accuracy of SVM gene signatures for cisplatin,
863 at a range of response thresholds, were assessed using gene expression data for
864 cisplatin-treated bladder cancer patients from Als *et al.*²⁹. Patients with a ≥ 5 year
865 survival post-treatment were labeled sensitive. Red arrows indicate the SVM gene
866 signatures with the highest positive predictive value (PPV) in the accuracy of
867 classification of patient outcome.

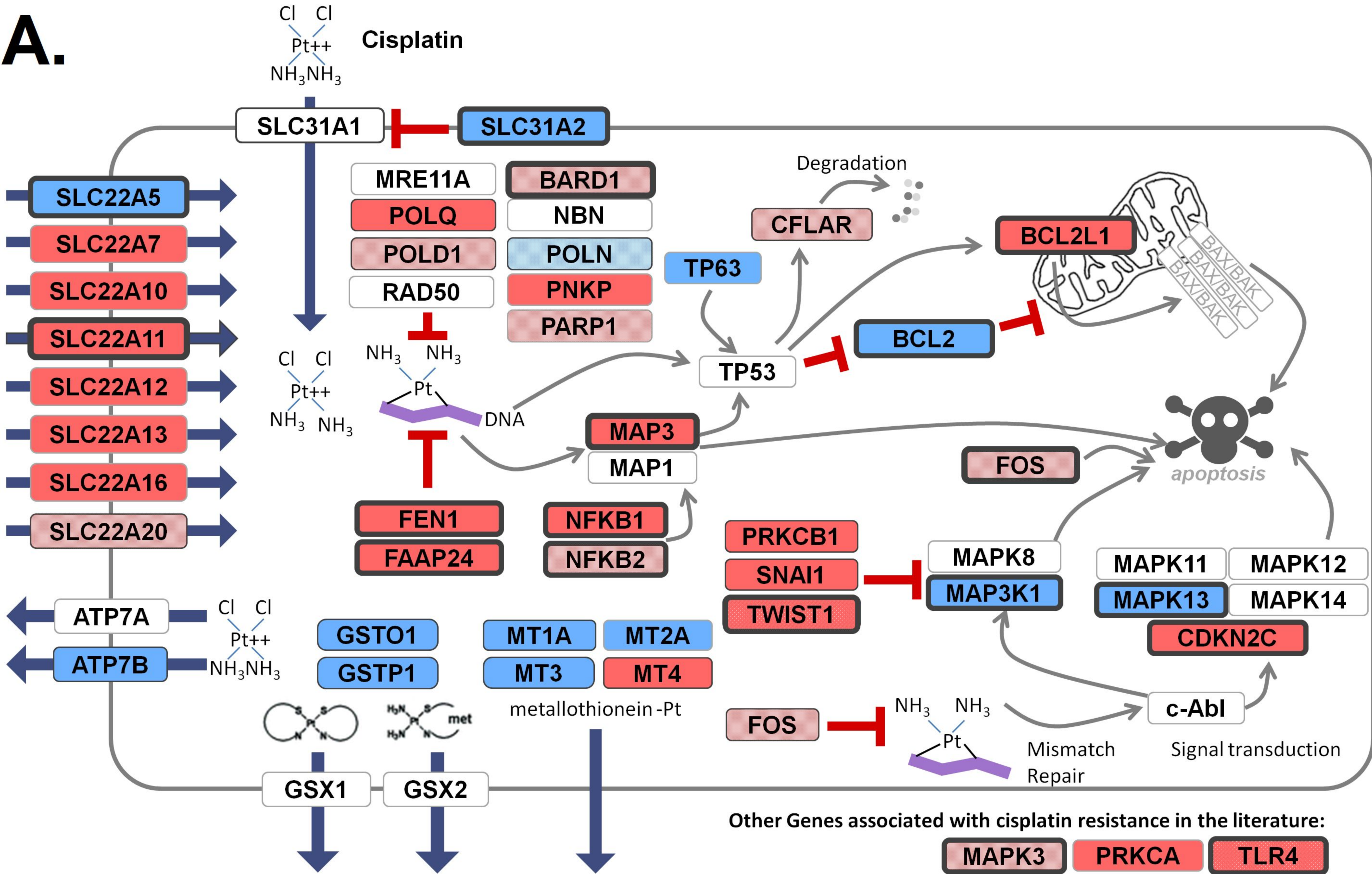
868 **Supplementary Figure 4.** Hyperplane distance calculated by all thresholded SVMs for
869 recurrent (<6 months) TCGA patients. Each diagram represents the predictions of all
870 SVMs for all patients who had recurrence less than 6 months after treatment (N=10).
871 Each point represents an SVM, where the x-axis represents the number of cell lines set
872 to resistant (in order of lowest to highest GI_{50}), and the y-axis represents the calculated
873 hyperplane distance. A negative hyperplane distance would represent a prediction of
874 resistance to cisplatin. Despite this, some patients show a strong preference towards
875 predictions of sensitivity (i.e. TCGA-XF-A9SU).

876 **Supplementary Figure 5.** Hyperplane distance calculated by all thresholded SVMs for
877 sensitive TCGA patients. Each diagram represents the predictions of all SVMs for all
878 patients who had recurrence > 4 years after treatment (top; N=3), or patients who
879 showed no recurrence after 6 years (bottom; N=6). Each point represents an SVM,
880 where the x-axis represents the number of cell lines set to resistant (in order of lowest to
881 highest GI₅₀), and the y-axis represents the calculated hyperplane distance. A positive
882 hyperplane distance would represent a prediction of sensitivity to cisplatin.

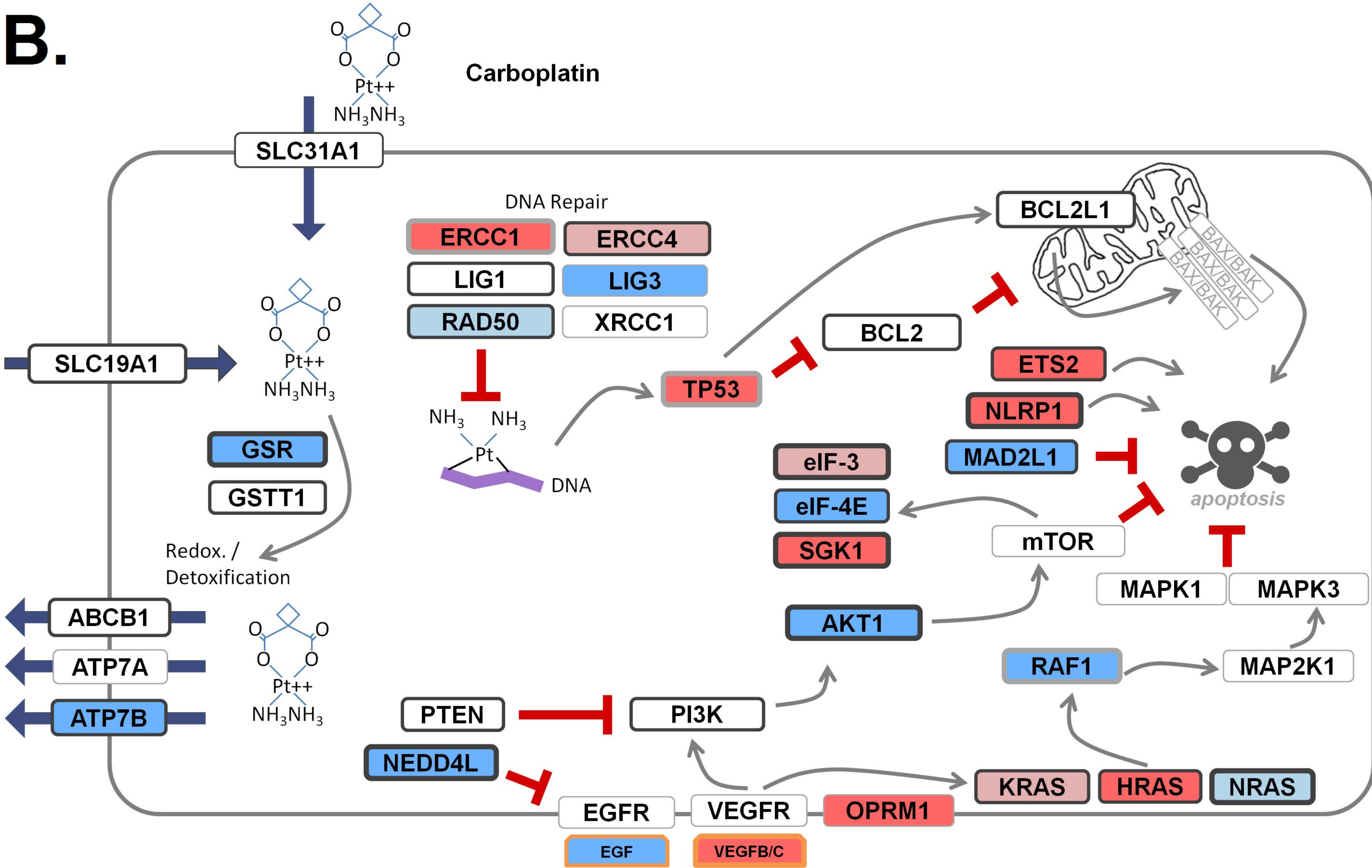
883 **Supplementary References.** A) Experimental evidence supporting inclusion of genes
884 using references relating expression to platinum drug efficacy. A subset of genes are
885 shown with consistent significant increase in misclassification error. B) The initial peer-
886 reviewed literature used to develop gene signatures associated with cis-, carbo- and
887 oxaliplatin response.

888 **Supplementary Tables.** Details of gene signatures, validation, and accuracy are
889 indicated in Tables: S1A) Genes Selected for MFA of Expression/Copy Number to
890 Cisplatin GI₅₀; S1B) Genes Selected for MFA of Expression/Copy Number to
891 Carboplatin GI₅₀; S1C) Genes Selected for MFA of Expression/Copy Number to
892 Oxaliplatin GI₅₀; S2A) SVM Models by Varying Resistance Thresholds and Impact on
893 Misclassification, Categorized by Gene Function for Cisplatin; S2B) SVM Models by
894 Varying Resistance Thresholds and Impact on Misclassification, Categorized by Gene
895 Function for Carboplatin; S2C) SVM Models by Varying Resistance Thresholds and
896 Impact on Misclassification, Categorized by Gene Function for Oxaliplatin; S3A)
897 Cisplatin SVM Models Derived Over a Range of Response Thresholds Using Log-Loss
898 Minimization; S3B) Carboplatin SVM Models Derived Over a Range of Response
899 Thresholds Using Log-Loss Minimization; S3C) Oxaliplatin SVM Models Derived Over a
900 Range of Response Thresholds Using Log-Loss Minimization; S4) SVM Models

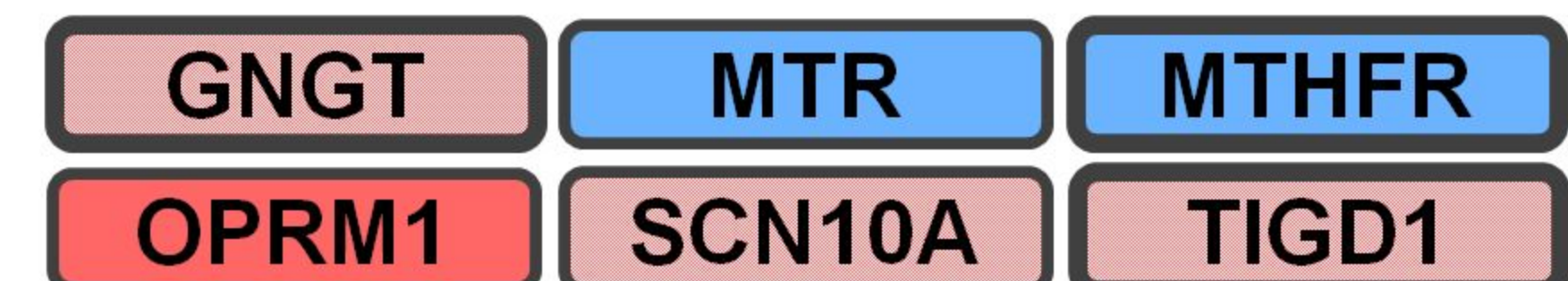
901 Generated by Bladder Cancer Patient Data at Various Time to Recurrence thresholds,
902 Categorized by Gene Function for Cisplatin; S5A) Cisplatin SVM Models Derived Over a
903 Range of Response Thresholds for TCGA Bladder Cancer Patients using
904 Misclassification; S5B) Carboplatin SVM Models Derived Over a Range of Response
905 Thresholds for TCGA ovarian epithelial tumor patients using Misclassification; S5C)
906 Oxaliplatin SVM Models Derived Over a Range of Response Thresholds for TCGA
907 colorectal adenocarcinoma patients using Misclassification; S6A) Accuracy of SVMs for
908 TCGA Bladder Cancer Patients, in relation to Time to Recurrence Post-Treatment With
909 Cisplatin; S6B) Accuracy of SVMs for TCGA Bladder Cancer Patients, in relation to Time
910 to Recurrence Post-Treatment With Carboplatin; S6C) Accuracy of SVMs for TCGA
911 Bladder Cancer Patients, in relation to Time to Recurrence Post-Treatment With
912 Oxaliplatin; S7A) Accuracy of Threshold Independent Analysis (Ensemble Averaging) for
913 Cisplatin; S7B) Accuracy of Threshold Independent Analysis (Ensemble Averaging) for
914 Carboplatin; S7C) Accuracy of Threshold Independent Analysis (Ensemble Averaging)
915 for Oxaliplatin; S8) Accuracy of SVMs for Non-Smoking TCGA Bladder Cancer Patients;
916 S9) Accuracy of SVMs for Smoking TCGA Bladder Cancer Patients (within 15 years of
917 diagnosis).

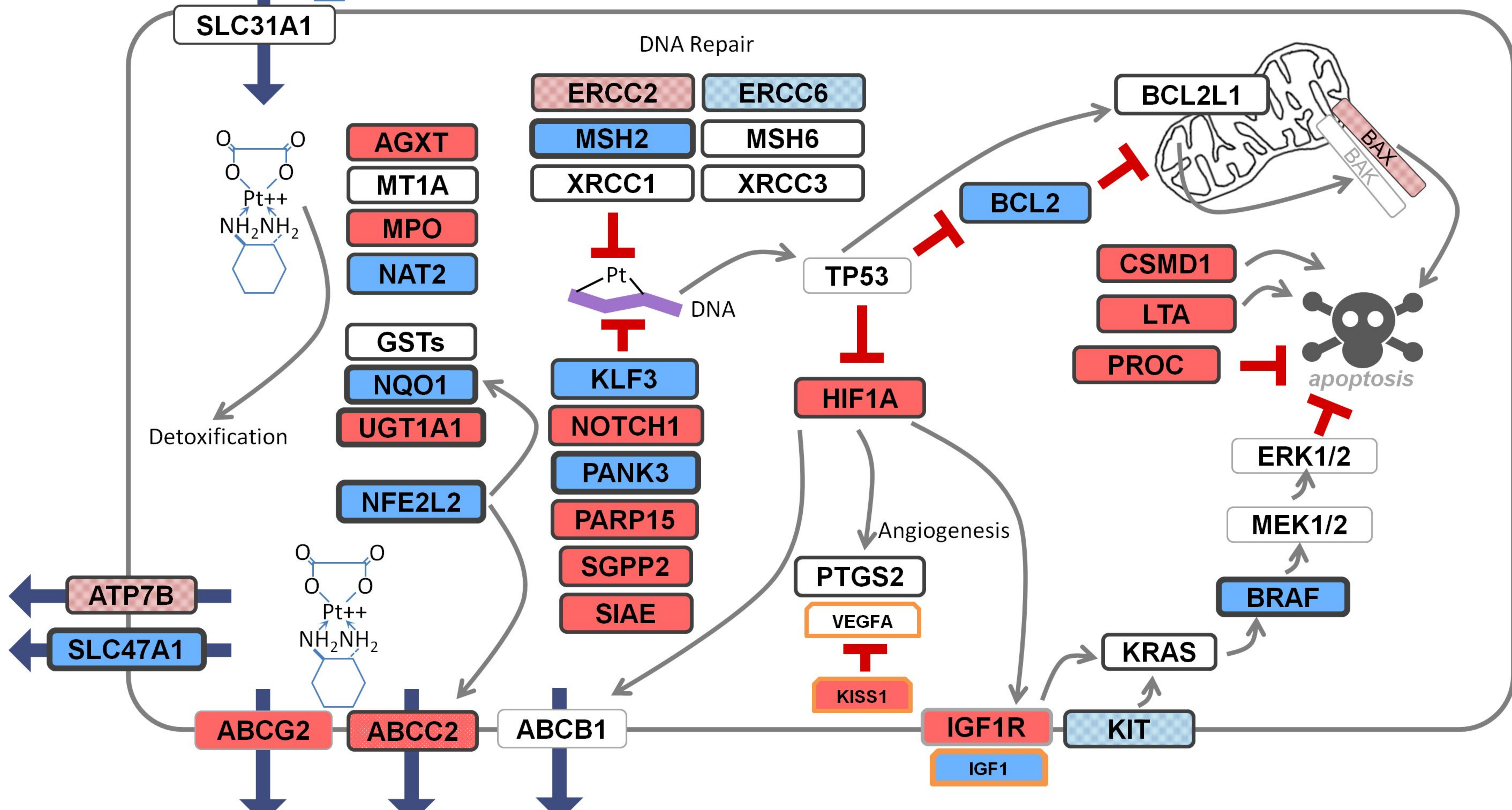
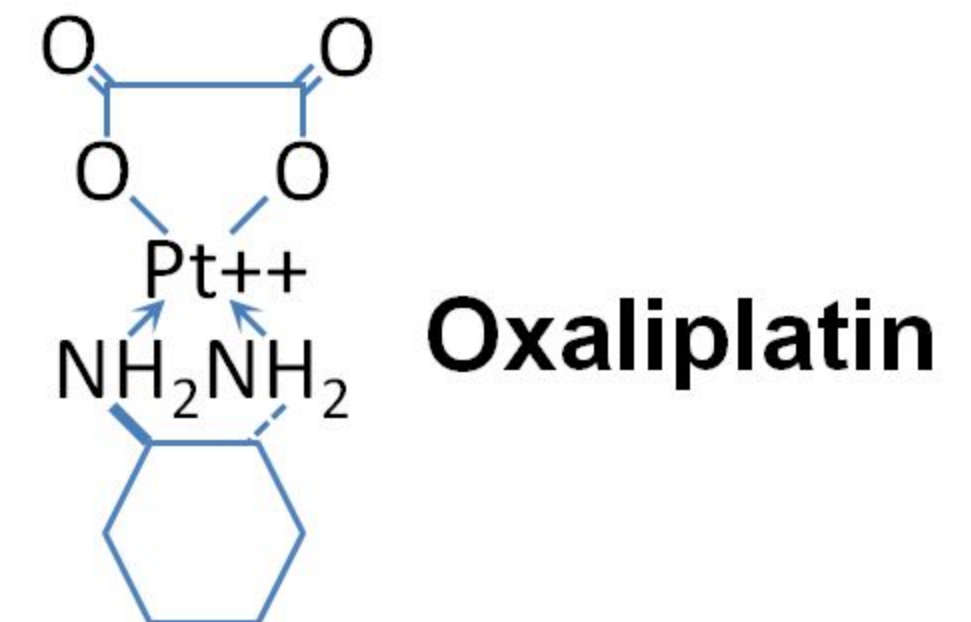


	MFA: Expression	MFA: Copy Number	MFA: Both	
Higher Levels Correlate with Higher GI50 & IC50				No MFA Correlation
Lower Levels Correlate with Higher GI50 & IC50				

B.

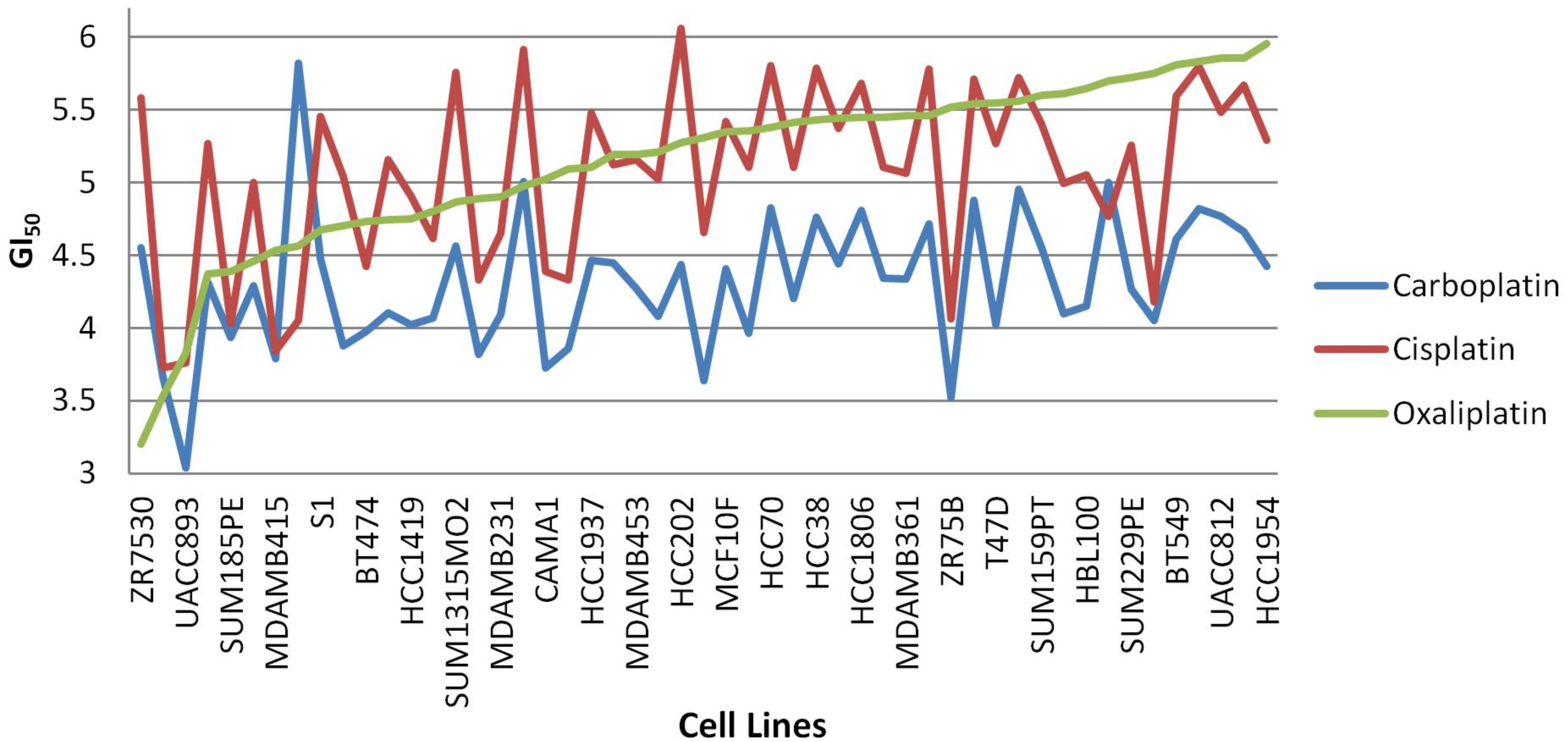
Other Genes associated with carboplatin resistance in the literature:

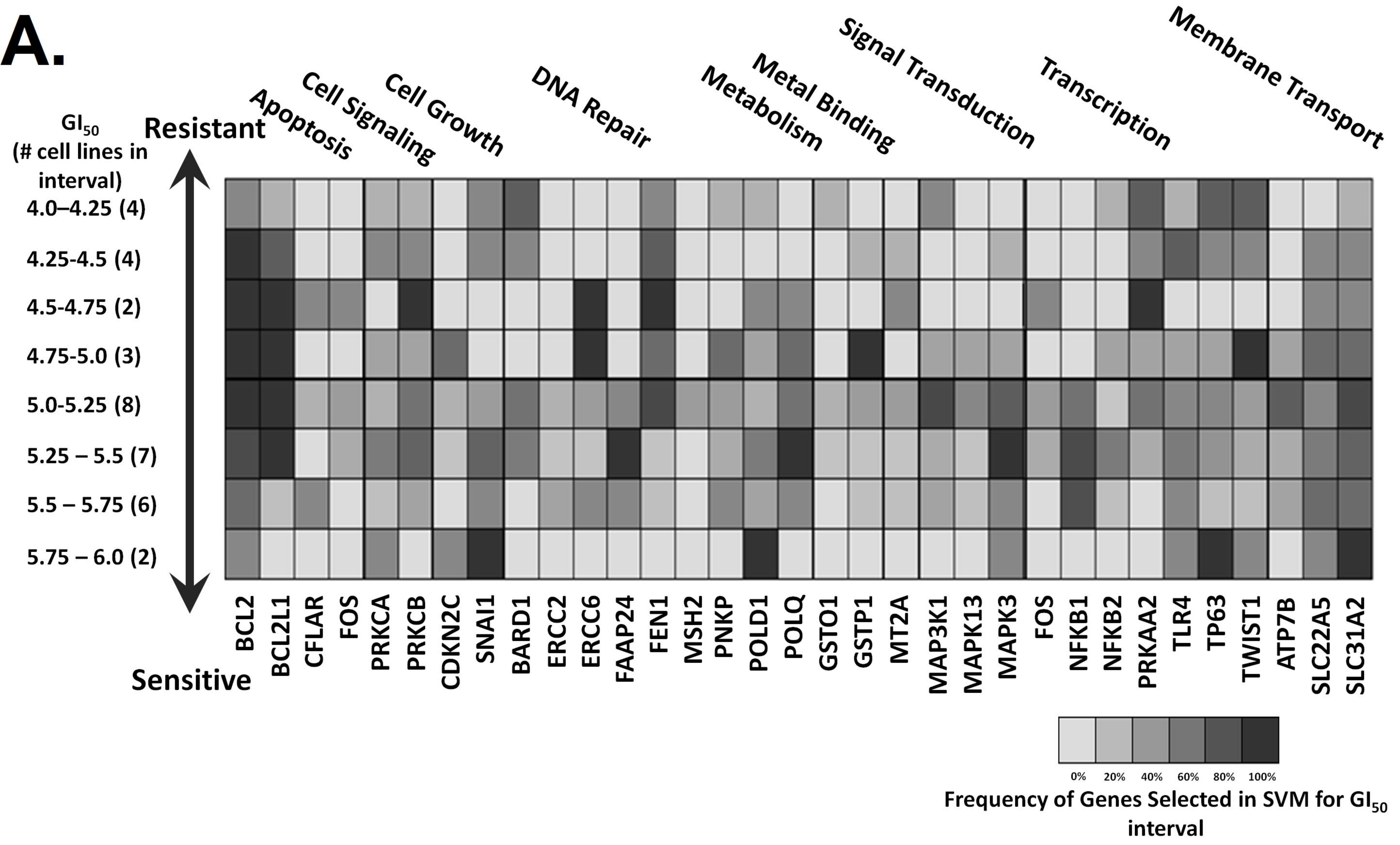


C.

Other Genes associated with oxaliplatin resistance in the literature:

- | | | | | |
|---------|-------|-------|--------|---------|
| APOBEC2 | CLCN6 | GRTP1 | FCGR2A | ICAM5 |
| HLA-B | KLF3 | NAGK | PRSS1 | SLCO1B1 |





B.

GI_{50}
(# cell lines in
interval)

Resistant

3.50 – 3.75 (4)

3.76 – 4.00 (7)

4.01 – 4.25 (8)

4.26 – 4.50 (11)

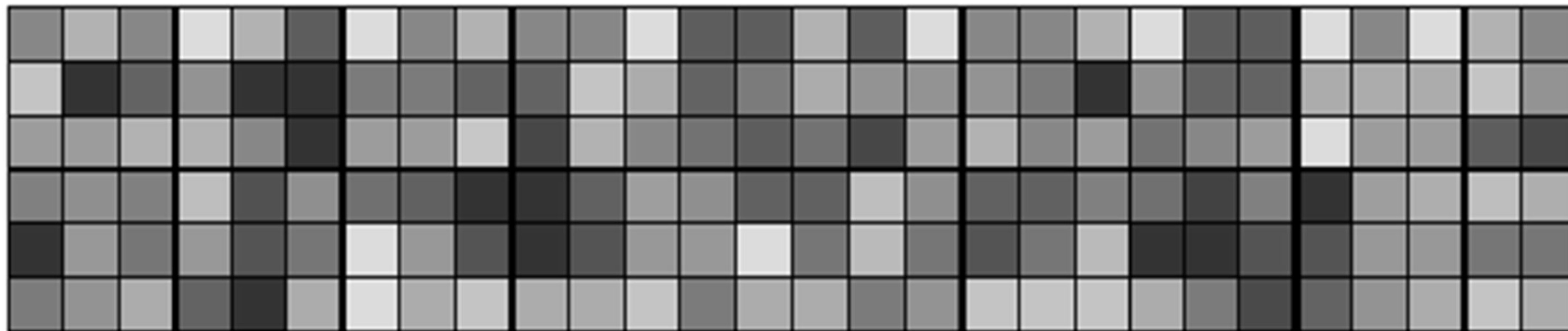
4.51 – 4.75 (5)

4.76 – 5.00 (8)

DNA Repair

Cell Growth

Metabolism

Signal Transduction /
ApoptosisTranscription /
DNA BindingMembrane Transport
Translation

ERCC1

LIG3

RAD50

EGF

VEGFB

VEGFC

GSR

MTHFR

MTR

AKT1

GNGT1

HRAS

MAD2L1

NLRP1

NRAS

OPRM1

RAF1

KRAS

NEDD4L

SGK1

ETS2

TIGD1

TP53

EIF3I

EIF3K

EIF4E2

ATP7B

SCN10A



0% 20% 40% 60% 80% 100%

Sensitive

Frequency of Genes Selected in SVM for GI_{50}
interval

