

ORIGINAL ARTICLE

Pathway-extended gene expression signatures integrate novel biomarkers that improve predictions of patient responses to kinase inhibitors

Ashis J. Bagchee-Clark¹ | Eliseos J. Mucaki¹ | Tyson Whitehead² |Peter K. Rogan^{1,3} 

¹ Department of Biochemistry, Schulich School of Medicine and Dentistry, University of Western Ontario, London, Canada N6A 2C8, Canada

² SHARCNET, University of Western Ontario, London, Ontario N6A 5B7, Canada

³ Cytognomix Inc., 60 North Centre Road, Box 27052, London, Canada N5X 3X5, Canada

Correspondence

Peter K. Rogan, Department of Biochemistry, Schulich School of Medicine and Dentistry, University of Western Ontario, Siebens-Drake Research Institute, Room 201A, London, Ontario N6A 5C1, Canada. Email: progan@uwo.ca

Ashis Bagchee-Clark and Eliseos Mucaki contributed equally to this work.

Funding information

Compute Canada; SHARCNET

Abstract

Cancer chemotherapy responses have been related to multiple pharmacogenetic biomarkers, often for the same drug. This study utilizes machine learning to derive multi-gene expression signatures that predict individual patient responses to specific tyrosine kinase inhibitors, including erlotinib, gefitinib, sorafenib, sunitinib, lapatinib and imatinib. Support vector machine (SVM) learning was used to train mathematical models that distinguished sensitivity from resistance to these drugs using a novel systems biology-based approach. This began with expression of genes previously implicated in specific drug responses, then expanded to evaluate genes whose products were related through biochemical pathways and interactions. Optimal pathway-extended SVMs predicted responses in patients at accuracies of 70% (imatinib), 71% (lapatinib), 83% (sunitinib), 83% (erlotinib), 88% (sorafenib) and 91% (gefitinib). These best performing pathway-extended models demonstrated improved balance predicting both sensitive and resistant patient categories, with many of these genes having a known role in cancer aetiology. Ensemble machine learning-based averaging of multiple pathway-extended models derived for an individual drug increased accuracy to >70% for erlotinib, gefitinib, lapatinib and sorafenib. Through incorporation of novel cancer biomarkers, machine learning-based pathway-extended signatures display strong efficacy predicting both sensitive and resistant patient responses to chemotherapy.

KEYWORDS

biochemical pathways, gene signatures, machine learning, systems biology, tyrosine kinase inhibitors

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2020 The Authors. *MedComm* published by Sichuan International Medical Exchange & Promotion Association (SCIMEA) and John Wiley & Sons Australia, Ltd.

1 | INTRODUCTION

Selection of a chemotherapy regimen is largely determined by efficacy of a drug in eligible subjects for a specific type and stage of cancer, and considers duration, location and magnitude of responses.¹ Individuals progress to second-line chemotherapeutic agents after demonstrating or developing limited efficacy to or after relapse from first-line chemotherapeutics.^{2,3} It is feasible to consider personal differences in genomic responses as a means of differentiating between acceptable chemotherapies with otherwise similar response rates across populations of eligible patients.⁴

Previously, we developed gene signatures that predict patient responses to specific chemotherapies from gene expression (GE) and copy number (CN) levels in a set of distinct breast and/or bladder cancer cell lines,⁵ with each line characterized by the drug concentration that inhibited growth by half (GI_{50}).^{6,7} Support vector machine (SVM) and random forest machine learning (ML) models were built for each drug using expression and/or CN values from ‘curated genes’ with evidence from published cancer literature of a contribution to the function or response to said drug in cell lines or patients. This paper develops signatures for tyrosine kinase inhibitors (TKIs),⁸ for which literature on genes associated with response is somewhat more limited.

We developed a novel technique for generating biochemically inspired gene signature models by expanding the pool of genes for ML to include genes both possessing and lacking literature support. The premise for including novel genes or gene products in these models is that these candidates could be related to genes supported by documented evidence through biochemical pathways or interactions that also contribute to drug response. We then compare conventional ML-based gene signatures to corresponding pathway-extended (PE) versions for these TKIs.

Abnormal expression levels or mutations in tyrosine kinases are often causally related to tumour angiogenesis⁹ and metastasis¹⁰ in certain cancers.^{11,12} TKIs have emerged as effective anti-cancer therapies, owing to their activity by ATP-competitive inhibition of the catalytic binding site of these kinases.¹³ Despite a conserved mechanism of action, sorafenib, sunitinib, erlotinib, gefitinib, imatinib and lapatinib preferentially inhibit different tyrosine kinase targets and exhibit distinct pharmacokinetic profiles.^{13–15} Sorafenib and sunitinib both inhibit VEGFRs, PDGFRs, FLT3R, RET and c-Kit.^{15,16} However, structural differences produce different binding profiles. For example, in binding VEGFR, sorafenib stabilizes the DFG-out inactive conformation of the enzyme, which allows it to bind within an allosteric pocket,¹⁷ whereas sunitinib binds in and around the ATP-binding region, imparting lower kinase selectivity

and faster off-rates.¹⁸ Similarly, erlotinib and gefitinib are both preferential inhibitors of EGFR, and share analogous chemical structure^{19,20}; but post-absorption, gefitinib is localized to a greater extent in tumour tissue, whereas erlotinib preferentially accumulates in plasma.²¹ Imatinib is particularly selective for the ABL kinase,^{8,22,23} whereas lapatinib binds to both EGFR and ERBB2.²⁴ The specificities of TKIs for different tyrosine kinase targets and the relative activities of those targets in different tumour types largely determine which of these drugs are recommended to treat individual clinical indications. These include renal cell carcinoma (sunitinib, sorafenib), hepatocellular carcinoma (sorafenib), pancreatic cancer (erlotinib), lung cancer (erlotinib, gefitinib), breast cancer (lapatinib) and chronic myelogenous leukaemia (imatinib).

Tumour cells can exhibit intrinsic or acquired resistance to chemotherapy. Intrinsic responses refer to an inherent capability to suppress the effects of treatment or render treatment cytostatic to functional characteristics of these cells. In acquired resistance, the tumour mutates or undergoes epigenetic changes after an initial period of clinical success that renders it impervious to treatment.^{25,26} Cytostasis is often achieved by inhibition of glycolytic activity with signal transduction, with the largest group of drugs targeting tyrosine kinases.²⁷ On average, tumours initially responsive to TKI treatments such as erlotinib and gefitinib will progress again within a year of treatment.^{28,29} Intrinsic resistance to these TKI drugs tends to be uncommon in EGFR-positive tumors.³⁰

Recent studies have revealed novel pathways of resistance and sensitivity to chemotherapeutic drugs.^{31,32} This study aimed to generate models that comprehensively represent global drug responses by inclusion of novel genes or gene products discoverable through their interactions with gene products known to influence these responses. We modify supervised ML-based models to systematically identify novel biomarkers whose expression is related to GI_{50} . GE changes in cancer cell lines that expand conventional gene signatures beyond an initial curated set of genes are utilized, including or replacing the initial set with other genes that interact with them. The resulting signatures aim to improve accuracy of prediction of individual patient responses to chemotherapies targeted towards tyrosine kinases.

2 | METHODS

2.1 | Data and pre-processing of cell line and cancer patient data sets

Microarray GE, CN, and GI_{50} values of breast cancer cell lines treated with erlotinib, gefitinib, imatinib, lapatinib,

sorafenib and sunitinib (obtained from Daemen et al⁵) were used to derive ML-based gene signatures that predict drug responses. The median GI_{50} values for these cell lines were applied as the threshold distinguishing sensitivity from resistance during ML. The median (range) of GI_{50} values for erlotinib was 4.71 (4.18-6.54), gefitinib was 5.03 (4.48-6.45), imatinib was 4.69 (3.82-5.81), sorafenib was 4.27 (3.0-5.83) and sunitinib was 5.23 (4.70-5.98).^{5,6} For lapatinib, the threshold was set at the GI_{50} value with the maximum difference relative to adjacent cell lines (4.94 [ranges from 4.78 to 6.40]), because the GI_{50} of multiple cell lines was equal to the median value.

Performance of these gene signatures was assessed using published studies of cancer patients treated with these drugs. NCBI Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) sourced data sets contained GE data and linked clinical outcomes of each patient with non-small cell lung carcinoma (NSCLC; GSE61676, N = 43)³³ treated with erlotinib (in combination with bevacizumab), hepatocellular carcinoma (GSE109211, N = 67)³⁴ treated with sorafenib, breast cancer (GSE66399, N = 31)³⁵ treated with lapatinib ('Arm B' patient set only, which received lapatinib in combination with paclitaxel, fluorouracil, epirubicin and cyclophosphamide), chronic myelogenous leukaemia (GSE14671, N = 23)³⁶ treated with imatinib, breast cancer patients (GSE33658, N = 11)³⁷ treated with gefitinib (in combination with anastrozole and fulvestrant) and gliomas (GSE51305, N = 18)³⁸ treated with sunitinib. Each of these studies provided clinical information that included a treatment outcome measure that could then be utilized as a binary outcome measure for comparison with predictions made by various models. These outcome measurements vary from study to study. For patients treated with sorafenib or imatinib, a chemotherapy response biomarker was used to distinguish sensitive from resistant patients. For patients treated with erlotinib or lapatinib, outcome (i.e. survival vs death) was used as a surrogate for response. Cancer cell migration data distinguished patients sensitive versus resistant to sunitinib (where those with 'moderate induction' or 'moderate inhibition' were defined as resistant, and those with 'strong inhibition' were considered sensitive to the drug). Responses to gefitinib were classified based on Response Evaluation Criteria In Solid Tumors (RECIST) guidelines (where those with progressive disease are considered TKI resistant).³⁹

Patient selection criteria differed between studies. In the GSE61676 study (erlotinib), patient data were acquired from the SAKK 19/05 trial, where selection criteria consisted of patients with newly diagnosed or recurrent Stage IIIB or Stage IV NSCLC.³³ In the sorafenib study (GSE109211), tumour tissue was collected from the STORM trial, which enrolled patients with hepatocellular

carcinoma with complete radiological response after surgical resection or local ablation.³⁴ The lapatinib study (GSE66399) utilized data from the CHER-LOB study, where female adults with HER2+ breast cancer were selected.³⁵ In the GSE33658 patient cohort, CD34+ cells were isolated from peripheral blood collected from newly diagnosed chronic-phase chronic myelogenous leukaemia patients treated with imatinib.³⁶ In the gefitinib study (GSE33658), biopsies were taken from postmenopausal women with newly diagnosed ER+ breast cancer receiving anastrozole, fulvestrant and gefitinib.³⁷ In the sunitinib study (GSE51305), native glioma tissue samples were collected from patients with a diagnosis of high-grade glioma (WHO [World Health Organization] grade III or IV) who underwent surgical resection.³⁸

Different expression microarray platforms were used in these GEO data sets, for example GSE66399, GSE61676 and GSE51305 each measures GE values with distinct vendor and gene sets. To minimize batch effects and apply the cell line-based signatures to these patient data sets, the data were first normalized on a common scale using quantile normalization, according to our previously published approach.⁴⁰ If multiple microarray probes existed for the same gene, the mean of all probe measurements was determined.

2.2 | Multiple factor analysis and gene set expansion

Genes associated with therapeutic response or function were curated from previous peer-reviewed publications for each TKI (refer to References in the Supporting Information). Inclusion criteria were based on evidence of the gene or protein contributing to pharmacokinetic or pharmacodynamic response, or were established biomarkers of sensitivity or resistance. Multiple factor analysis (MFA) was performed using cell line expression and GI_{50} (concentration of drug inhibiting 50% growth) data⁵ for each curated gene using the *MFAPreselection* software we have developed (available in a Zenodo archive⁴¹). The archive describes the algorithm used by *MFAPreselection* to traverse pathway networks, dataflow within the program, and software code. MFA determines the relationship between GI_{50} and GE and/or CN data for all expressed genes as an angle that indicates the degree to which expression or CN correlates either directly (~ 0 degrees) or inversely (~ 180 degrees) with the GI_{50} of the set of cell lines.^{42,43} Circular plots, generated by *MFAPreselection*, indicate this correlation angle (Figure 1).⁴³ *MFAPreselection* searches for known gene pseudonyms and substitutes the correct alias (from www.genecards.org [downloaded July 2016]). In the Daemen et al's data set⁵ used for training in SVM

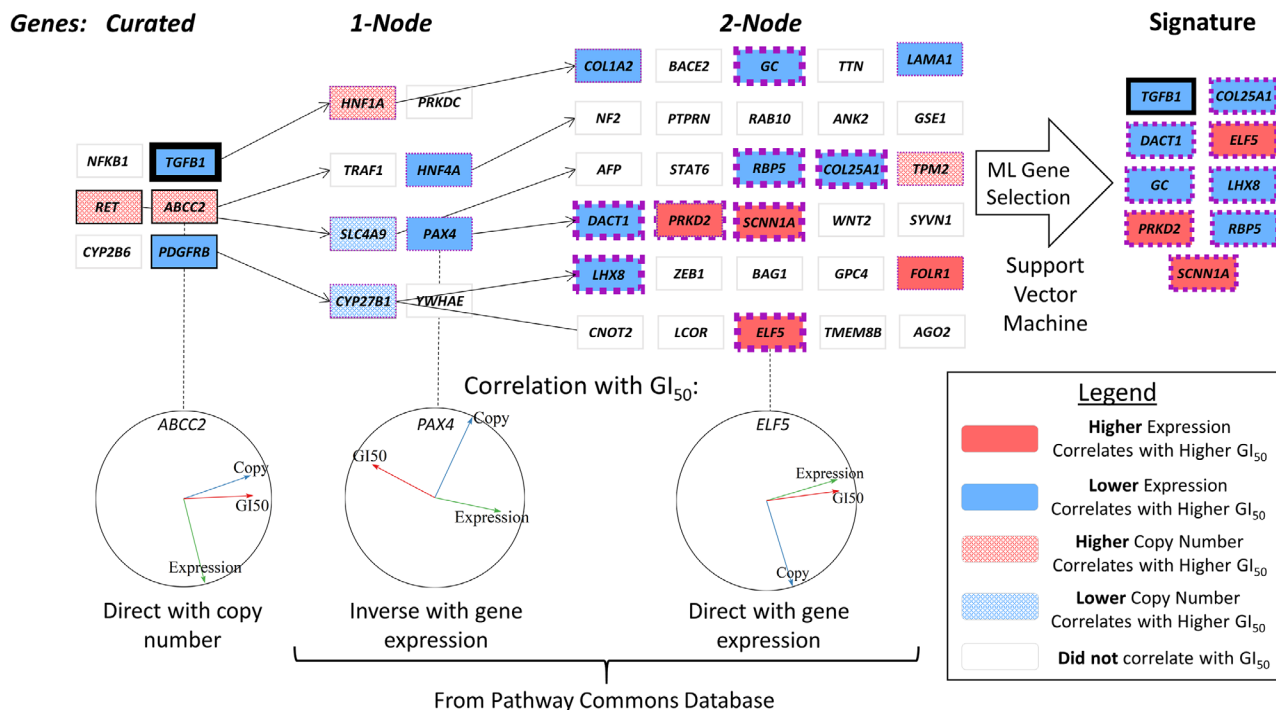


FIGURE 1 Procedure for pathway gene selection. An initial set of genes with known associations to a particular TKI (here, we show a subset of sorafenib-related genes) are selected and then evaluated by MFA, which was used to find a correlation between cell line drug sensitivity (GI_{50}) and the GE or CN of these genes in those cell lines (*left*). MFA correlation circles visualize these relationships (*bottom*). The gene list is extended, using pathway and interaction databases (i.e. PathwayCommons) to find genes related to curated genes which showed MFA correlation to GI_{50} (one-node distant genes; *middle-left*). The list is extended again from the MFA-correlating one-node distant genes (two-node distant genes; *middle-right*). All curated and extended genes which showed an MFA correlation were then used as features to generate a final predictive SVM gene signature for the evaluated TKI (*right*). Genes within the best performing sorafenib signature are indicated in thick borders (black for curated genes; purple for pathway-extended genes)

learning (see below), the microarray platform data were in some instances labelled with conflicting gene names. During pathway extension, associated genes were related to older gene aliases that have been deprecated and re-assigned by HUGO (Human Genome Organization) Gene Nomenclature Committee to other unrelated genes. This led to some spurious associations between genes during pathway extension. Examples include *PPY* which was mismatched due to its former alias ‘PNP’ as well as *DDR2* due to incorrect association with its former alias ‘TKT’. For the sorafenib model **PE-Sor**, associations of *GC* to *CNN1* and *CA3* were eliminated due to its original designation as ‘DBP’; however, its associations between *HNFI1A*, *CYP11B1*, *CYP27B1* and *PIK3R3* remained valid (Figure S1). This issue was addressed using a program script that removed these unsupported associations from the output of *MFAPreselection*.⁴¹

This script eliminated these spurious matches by confirming relations reported by *MFAPreselection* with the PathwayCommons Interaction SIF (Simple Interaction Format) file (‘Parentage-MFA-Path-Source-Program.Simple-Output-Version.pl’; provided in a Zenodo archive⁴¹). If corrected labels were not found or a gene

was absent from a microarray platform, then this cell line or gene is not included in the analysis.

ML signatures were expanded by *MFAPreselection* to include genes associated with curated genes by extension using components of adjacent biochemical pathways (pathway extension [PE]; Figure S2). To identify these relationships, *MFAPreselection* relied on the PathwayCommons database (version 8 [downloaded April 2016]) to assess expanded gene lists by inclusion of genes addressable from the curated set (one node distant from a curated gene), followed by a second iteration (two nodes distant from a curated gene; illustrated in Figure 1). During this process, genes that did not meet minimally defined levels of MFA correlation to drug GI_{50} (either positive or negative) were discarded and additional gene expansion steps also ignored these genes. These levels were determined using six different conditions set for the *MFAPreselection* software: maximum thresholds up to 10° and 20° from either full direct or inverse correlation for curated genes only (conditions #1 and 2, respectively); up to a 10° and 20° threshold for both curated genes and directly related genes (one-node distant; conditions #3 and 4, respectively) and up to a 10° and 20° threshold for curated genes and genes

up to two nodes distant from the curated gene set (conditions #5 and 6, respectively). Genes in which GI_{50} was correlated with CN (Tables S1A-F) were not considered for SVM analyses due to unavailability of CN data in patient data sets.

2.3 | SVM learning

Genes with expression levels correlated with GI_{50} were qualified for SVM analysis. SVMs were used to train GE data sets against GI_{50} data using the MATLAB statistics toolbox (similar to the procedure described in Mucaki et al⁴⁴ using SVM software developed in Zhao et al⁴⁰; software available at <https://doi.org/10.5281/zenodo.1170572>). Instead of using the 'fitsvm' function (as in Mucaki et al⁴⁴), a multiclass-compatible 'fitcecoc' function was used to generate SVM signatures, with both misclassification rate⁴⁴ and log loss⁴⁰ value used as performance metrics to derive optimal signature models. A forward feature selection (FFS) algorithm was used to generate these gene signatures (program from Zhao et al⁴⁰: 'FFS_strat_kfold_gridsearch.m'). FFS tests each gene at random from the qualified gene set by training a cross-validated Gaussian kernel SVM on the training data to determine the individual gene that produces the lowest misclassification rate or log loss value. Subsequent genes are then added to determine whether model performance is improved, until the performance criterion converges to a minimum value. Models were built using a range of C and sigma values (from 1 to 100,000, in multiples of 10 for each variable [where $C \geq \sigma$]; 21 total combinations). Because the goal of pathway extension was to expand and improve these models beyond curated signatures with more than two genes, PE-derived gene signatures with fewer than two genes were excluded from proceeding to the validation step.

2.4 | Validation of cell line-derived gene signatures using patient data

All derived multi-gene SVMs were validated against clinical patient data using traditional validation (MatLab program 'regularValidation_multiclassSVM.m' from Zhao et al⁴⁰). Performance was indicated by both overall predictive accuracy and by Matthews Correlation Coefficient (MCC; which assesses overall quality of a binary classifier by considering the balance of true and false positives and negatives). Overall, the best-performing gene signature for each drug was selected by MCC, as it is a metric not skewed by imbalanced data. Once the best performing SVM for each drug was established, leave-one-out cross-validation⁷ was used to determine the overall impact of each individ-

ual gene to the model itself (change in misclassification or log loss), as well as its impact on the accuracy of the model to predict chemotherapy response. Top-performing PE TKI models can be accessed to predict responses based on expression in individual patients with our web-based SVM calculator (<http://chemotherapy.cytogenomix.com>).⁶

Ensemble averaging of multiple SVM models involved weighting patient predictions from highest performing models derived for a particular TKI by the area under the curve (AUC) of each corresponding model (computed using the MATLAB function 'perfcurve'). MCC itself was also evaluated as a potential source of weights for ensemble averaging; however, AUC-weighted predictions were superior in overall performance. The number of models included in the ensemble varied, as the number of highest performance models for each TKI differed (4 for sorafenib; 2 for erlotinib, sorafenib, imatinib, sunitinib and gefitinib). A patient was considered resistant to a drug if the sum of all AUC-weighted predictions was > 0 and sensitive if this sum was < 0 .

3 | RESULTS

3.1 | Generating SVM signatures using breast cancer cell line-training data

Genes associated with drug response or function were curated for gefitinib (N = 113), sunitinib (N = 90), erlotinib (N = 71), imatinib (N = 157), sorafenib (N = 73) and lapatinib (N = 91) (curated genes are provided in Table S1 and labelled as '0' node distant genes). In general, MFA was performed using 48 breast cancer cell lines using GE, CN and GI_{50} values for each gene.⁵ Biochemically inspired ML-based signatures for each TKI, derived from curated genes, were obtained according to our previously described approach.⁶ MFA analysis was also performed on genes encoding proteins related to these curated genes (through interaction or as neighbours in the same biochemical pathway) to identify those that also correlated, either directly or inversely, with GI_{50} (all GI_{50} -correlated PE genes are provided in Table S1 [labelled as 1-node and 2-node distant genes]). This expanded set of GI_{50} -correlating genes were then used to derive SVMs containing combinations of curated and PE genes. The derived signatures for each TKI minimized either misclassification or log loss to generate the best performing models. The best performing curated and PE SVM signatures for erlotinib (**C-Erl**, **PE-Erl**), sorafenib (**C-Sor**, **PE-Sor**), gefitinib (**C-Gef**, **PE-Gef**), lapatinib (**C-Lap**, **PE-Lap**), imatinib (**C-Ima**, **PE-Ima**) and sunitinib (**C-Sun**, **PE-Sun**) are summarized in Table 1, whereas the performance of all models is indicated in Table S2.

TABLE 1 Performance of curated SVMs and PE models on training and patient testing data

TKI (Patient tumour type; GEO data set)	Model	Gene signature (SVM: C; σ)	Validation			Training	
			MCC	Sensitive	Resistant	Overall	Misclassification
Erlotinib (NSC Lung Carcinoma; GSE61676)	Cur	BAX, FOXO1 (100; 1)	0.08	100%	3%	23%	0.38
	PE	NEK7, SLCO3A1, RELB, FRMD4A, HSD17B2, CDK6, PALM, IL1RN, SMYD1, BAG2, GNG3, SULF1 (1000; 100)	0.41	42%	93%	83%	0.01
Sorafenib (Hepatocellular Carcinoma; GSE109211)	Cur	PDGFRB, TGFB1, SLCO1B1 (10,000; 1)	0.28	96%	29%	50%	0.72
	PE	ELF5, RBP5, GC, PRKD2, SCNN1A, COL25A1, TGFB1, DACT1, LHX8 (100,000; 1000)	0.72	72%	95%	88%	0.05
Gefitinib (Breast Cancer; GSE33658)	Cur	GRP (10,000; 10)	0.16	13%	100%	29%	0.53
	PE	CNTN1, CXCL2, NTRK3, GCG (10,000; 10)	0.67	100%	50%	91%	0.58
Lapatinib (Breast Cancer; GSE66399)	Cur	ERBB2 (10,000; 1)	0.31	13%	100%	77%	0.72
	PE	FBP1, ITGAI1, TRIM68, BCAT1, ZNF780A, UTP20, GRB7 (10; 10)	0.33	63%	74%	71%	0.01
Imatinib (Leukaemia; GSE14671)	Cur	IL3, ABL2, CDKN1A (10,000; 10)	0.23	41%	83%	52%	0.42
	PE	LIF, MRGPRF, GRM3, TNNT1, CACNA2D1 (100,000; 100)	0.18	84%	33%	70%	0.06
Sunitinib (Glioma; GSE51305)	Cur	HGF, VEGFC, TSCI, AXL, ENPP2, NFKB1 (100,000; 10,000)	0.31	87%	45%	59%	0.14
	PE	EPHA2, NR4A1, SIAF (100,000; 100,000)	0.61	67%	92%	83%	0.20

Curated and PE SVMs were derived for each TKI based on ability sorting cancer cell lines. The C (box-constraint), σ (kernel-scale) and features comprising the best-performing model are indicated. Models listed are those which exhibited optimal performance, defined as the model with the highest MCC against the patient data set. 'Validation' indicates the predicted drug response of patients made by each curated and PE model as compared to the observed response provided by these studies. 'Training' indicates either percent misclassification or overall log-loss of the cell line-based model by cross-validation, depending on which minimization metric was used in said model derivation. Abbreviations: Cur, Curated models derived from genes associated with drug in literature; PE, Pathway-extended models; MCC, Matthews Correlation Coefficient; Sensitive, Accuracy to Drug Responsive Patients; Resistant, Accuracy to Non-Responsive Patients; Overall, Combined accuracy.

3.2 | Validation of cell line-based SVM signatures using cancer patient data

Cell line-derived SVMs for TKIs were initially evaluated on patient data sets where patients were treated with the same agent.⁴⁰ Erlotinib signatures were validated using patients with NSCLC (GSE61676; N = 9 survived, 34 died), sorafenib signatures were validated using patients with hepatocellular carcinoma (GSE109211; N = 21 sensitive, 46 resistant), sunitinib signatures were validated using outcomes of patients with high-grade gliomas (GSE51305; N = 6 sensitive, 12 resistant), imatinib signatures were validated using outcomes of patients with chronic myelogenous leukaemia (GSE14671; N = 17 sensitive, 6 resistant) and lapatinib and gefitinib signatures were validated based on breast cancer outcomes (GSE66399 [N = 8 survived, 23 died] and GSE33658 [N = 10 sensitive, 2 with resistant], respectively).

MCC (range: -1 to +1) was the primary determinant of model performance, as it measures overall accuracy (OA) while accounting for representation between binary prediction categories.⁴⁵ This was necessary, as patient data sets available exhibited imbalances in the ratios of responsive to non-responsive patients in terms of their respective observed clinical outcomes. Models based on features generated under relaxed constraints (condition #6) generated the best performing SVM on patient data for every TKI, except sorafenib. The best-performing PE model was **PE-Sor**, which accurately predicted patient responses with 0.72 MCC (and 88% OA). The best performing curated model was **Cur-Lap**, with 0.31 MCC (and 77% OA). In comparison to curated SVMs, PE SVMs predicted patient response with 0.26 higher MCC and 33% higher OA (13% increase in accuracy predicting sensitive patients; 13% increase in accurately predicting resistant patients). Except for imatinib, the best-performing PE model outperformed their curated counterpart. This difference in performance is evident in Figure 2, as predictive accuracy for PE models is consistently higher for both resistant and sensitive patient outcomes.

The erlotinib (GSE61676) and gefitinib (GSE33658) studies utilized for model testing provide patient GE data both pre- and post-treatment. This provided an opportunity to determine whether to determine whether short-term drug exposure altered GE and model accuracy. For erlotinib, blood samples were obtained prior to and 24 h post-treatment. For gefitinib, biopsies were taken prior to and 3 weeks post-treatment. Both **PE-Erl** and **PE-Gef** exhibited slightly lower performance for the pre-treatment samples (Table S3), with five additional patients misclassified with **PE-Erl** (73% OA with N = 43 total patients) and two additional misclassified individuals with **PE-Gef** (73% OA; N = 12 patients). MCC for **PE-Gef** is

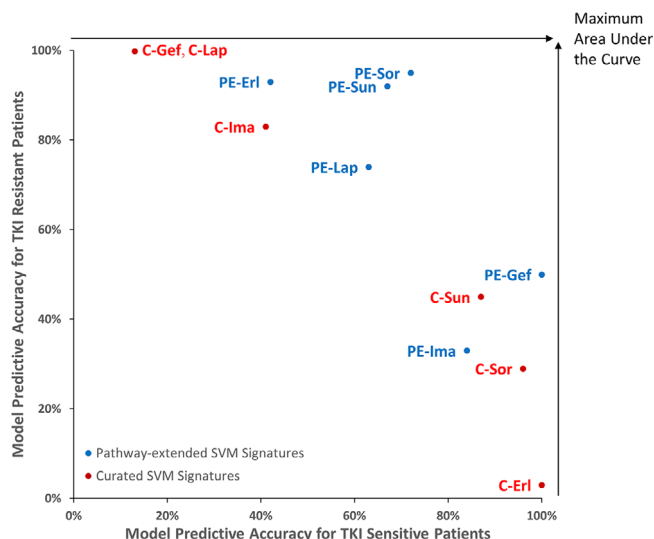


FIGURE 2 Accuracy of curated and pathway-extended SVMs on TKI sensitive and resistant patients. The predictive accuracy of the best-performing curated (C-) and pathway-extended (PE-) models for each TKI was arranged based on their accuracy in classification of drug-sensitive and drug-resistant tumour patients. This illustrates how curated models are often only accurate towards one patient class (sensitive or resistant) but not both (red), which is an issue as the patient data were often imbalanced (number of sensitive | resistant patients in each study: lapatinib ['Lap'; n = 8 | 23], imatinib ['Ima'; n = 17 | 6], sunitinib ['Sun'; n = 6 | 12], erlotinib ['Erl'; n = 9 | 34], gefitinib ['Gef'; n = 10 | 2] and sorafenib ['Sor'; n = 21 | 46]). Conversely, predictions by PE SVMs were often more balanced (blue), possessing moderate to high accuracy for both sensitive and resistant patients, and consequently greater accuracy as a whole

significantly lower (-0.15), because the model misclassifies all untreated individuals as resistant. Treatment with these drugs perturbs predictions, but to a limited extent.

3.3 | Composition of PE SVM signatures and contributions of individual features

PE SVM signatures contain either genes from peer-reviewed literature about the drug response ('initial' or 'curated' genes), those related to these genes through direct interactions or as neighbours within the same pathways (one-node distant genes), or genes associated with these one-node distant genes (two-node distant genes). To better comprehend the composition of and relationships between genes in the best-performing PE SVM signatures, we analysed the connection networks for each model (see Table S4 for connection network for all other top performing PE models). For example, although **PE-Sor** consists of one curated gene and eight two-node distant genes, there are an additional 6 curated and 10 one-node genes that

connect the genes in **PE-Sor** by pathway extension (Figure 3A shows a two-dimensional visualized connection network for this drug; see Figure 3B-F for lapatinib, gefitinib, sunitinib, imatinib and erlotinib, respectively). Due to the complexity of the relationships between gene products for erlotinib, it was not feasible to create an unequivocal two-dimensional network diagram for this drug response, and is instead presented in tabular form (Figure 3F). Nevertheless, it is apparent from the majority of these network diagrams that genes that were two-nodes distant from the curated gene set were most commonly selected in the best performing PE models. Furthermore, the two-node distant genes selected interacted with multiple curated or one-node distant genes.

To determine the degree to which each gene in a signature contributed to the accuracy of the overall model prediction, we performed leave-one-out cross-validation for each gene in the best-performing model for each drug. We then reassessed the predictions of the resultant signature for the observed responses in the cell lines used for model training (Table S5) and for the patient data used for testing (Figure 4). Based on patient data, the gene features eliminated from models that had the highest impact on performance were as follows: *CDK6*, *BAG2*, *SULT1E1* and *IL1RN* (**PE-Erl**); *CNTN1*, *GCG* and *NTRK3* (**PE-Gef**); *GRB7* and *BCAT1* (**PE-Lap**); *ELF5*, *TGFB1*, *PRKD2*, *RBP5* and *GC* (**PE-Sor**); *EPHA2* and *SIAE* (**PE-Sun**) and *CACNA2D1* and *GRM3* (**PE-Ima**). Genes removed that improved predictive performance on patient data included *FBP1* (**PE-Lap**), *PLAT* (**PE-Sor**) and *LHX8* (**PE-Sor**).

PE-Gef consists of four pathway-extended genes (*CNTN1*, *CXCL2*, *NTRK3* and *GCG*) and one curated gene, *GCG*. *GCG* encodes a hormone preprotein which is cleaved into four peptides, including glucagon-like peptide 2, which has been found to reduce gefitinib-induced intestinal atrophy in mice.⁴⁶ Removal of *NTRK3* from **PE-Gef** had the largest impact on model performance, reducing MCC to 0. *NTRK3* has a critical role in secretory breast cancer gene, with the *EVT6-NTRK3* fusion oncogene being considered a primary initiating event.^{47,48}

PE-Sun, which consists of three pathway-extended genes, *SIAE*, *NR4A1* and *EPHA2*, was evaluated in gliomas. *NR4A1* is essential for colony formation of glioblastoma cells on soft agar.⁴⁹ Of 14 glioblastoma specimens, 13 possessed elevated *EPHA2* levels.⁵⁰ Removal of *NR4A1* from **PE-Sun** did not alter overall accuracy or MCC of the model, whereas removal of *EPHA2* decreased overall accuracy by 55% and MCC by 0.94. Regarding *SIAE*, alterations in cell surface sialylation by glucocorticosteroids have been suggested to promote glioma formation.⁵¹

PE-Sor (*COL25A1*, *TGFB1*, *DACT1*, *RBP5*, *PRKD2*, *GC*, *ELF5*, *LHX8* and *SCNN1A*) was used to predict sorafenib response in hepatocellular carcinoma (HCC) patients.

Removal of *RBP5*, *PRKD2*, *GC* and *ELF5* significantly reduced overall accuracy (>50%) and MCC (>0.7) (Figure 4A). *RBP5* is linked to aggressive tumour features in HCC,⁵² *PRKD2* is upregulated in HCC and correlated with metastasis,⁵³ and decreased actin-free GC levels have been found to relate with disease severity in HCC.^{54,55} Vitamin D₃, which is bound by GC, lowers the effective dose of sorafenib required for its cytostatic effect in melanoma and differentiated thyroid carcinoma.⁵⁶ *ELF5* has not been directly connected to HCC, but has been associated with a wide range of cancers.^{57,58} Genes in **PE-Sor** that have not been as strongly linked to cancer (*COL25A1* and *LHX8*) did not change model accuracy to the same extent (<10%) when removed (Figure 4A). Removal of the curated gene *TGFB1*, which enhances the apoptotic activity and sensitizes cells to sorafenib,⁵⁹ decreased overall accuracy by 60% in HCC patients. The respective contexts of the curated Sorafenib-related genes juxtaposed with the PE genes in **PE-Sor** are indicated in a cellular schematic of the roles and functions of these genes (Figure 5).

PE-Ima (*LIF*, *MRGPRF*, *GRM3*, *TNNI1* and *CACNA2D1*) predicted imatinib response in chronic myelogenous leukaemia patients. *LIF* encodes a protein which prevents continued growth of myeloid leukaemia cells by inducing terminal differentiation,⁶⁰ although independent removal of *LIF* did not notably affect model performance. Downregulation of *CACNA2D1* from **PE-Ima** is associated with erythroid differentiation of K562 and KCL-22 chronic myeloid leukaemia cells.⁶¹ Removal of *CACNA2D1* decreased both classification accuracy and MCC (−44% and −0.18, respectively; Figure 4E).

A second PE model (indicated in green in Figure 4E) exhibited comparable performance to **PE-Ima**: *TNNI1* and *WASF3* ($C = 10,000$, $\sigma = 10,000$), with an OA of 57% (47% accurate with sensitive and 83% with resistant patients; MCC = 0.27). *WASF3* has been implicated in breast cancer metastasis.⁶² *TNNI1*, a gene that is shared by both this model and **PE-Ima**, is one of the three inhibitory subunits of smooth muscle troponin that are all overexpressed in breast cancer.⁶³ Interestingly, the kinase, *TNNI3K*, that phosphorylates this protein is essential for proliferation of mononuclear diploid cardiomyocytes during heart muscle repair due to injury.⁶⁴ Phosphorylation of troponin would appear to have a previously uncharacterized moonlighting function in tumour development.⁶⁵ If imatinib inhibits *TNNI3K* through an off-target effect, this may modulate *TNNI1* activation and possibly, an associated proliferative phenotype.

PE-Lap (*FBP1*, *ITGA11*, *TRIM68*, *BCAT1*, *ZNF780A*, *UTP20* and *GRB7*) predicted outcomes of breast cancer patients treated with lapatinib. Independent removal of *BCAT1* reduced accuracy in predicting sensitive patients. Silencing or knockdown of *BCAT1* has been associated

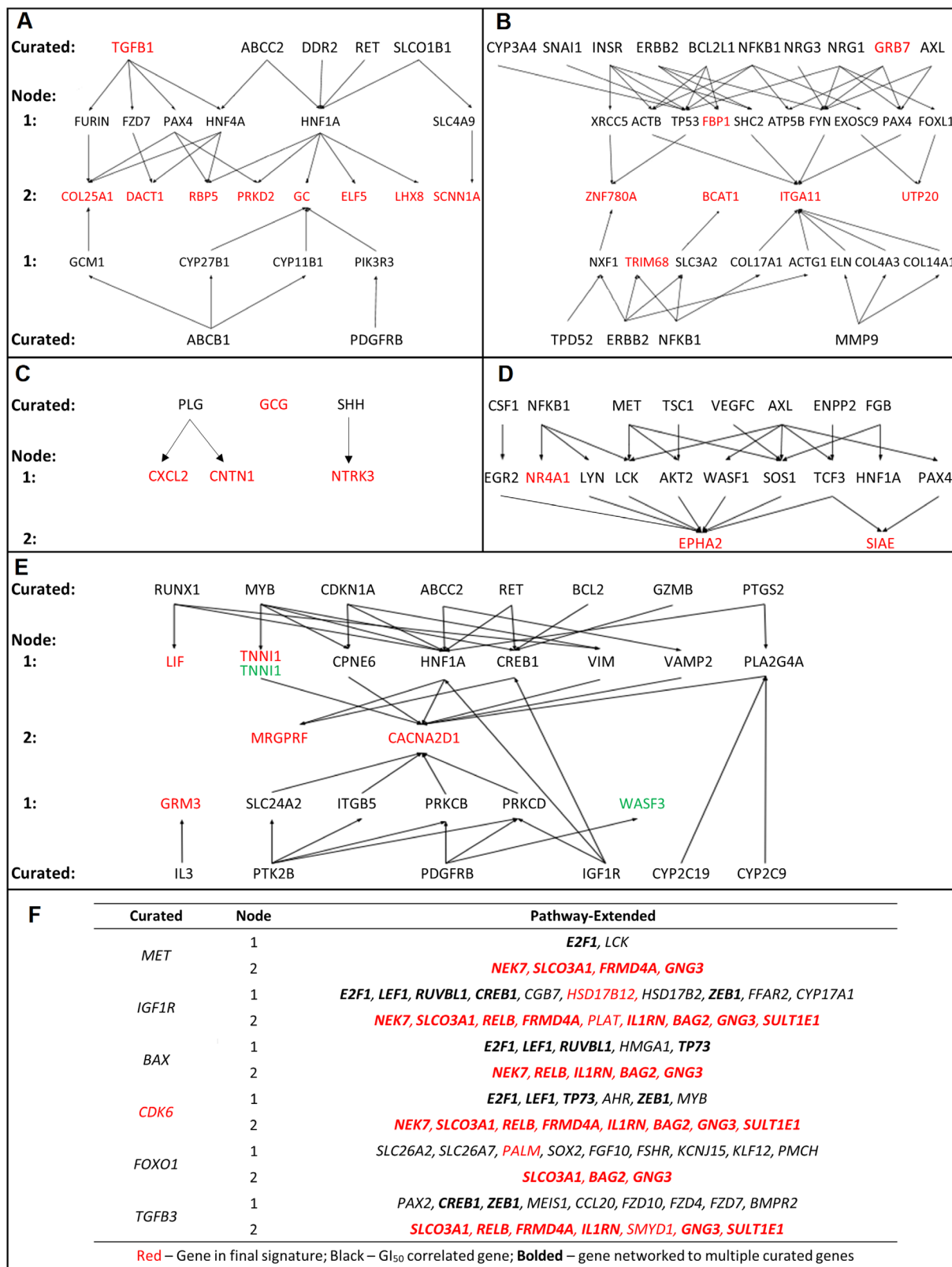


FIGURE 3 Connection network for pathway-extended TKI SVMs. Schematic relationships outlining the pathway connections for the best-performing PE model for each drug in panels (A) sorafenib, (B) lapatinib, (C) gefitinib, (D) sunitinib, (E) imatinib and (F) erlotinib. All symbols indicated are gene names. The erlotinib model was highly interconnected and is represented as a table. Genes in red are features selected for the final **PE-Sor** gene signature, whereas genes coloured green were chosen in a separate PE gene signature with comparable performance. Genes in black were not part of the final signature themselves but correlated with efficacy to sorafenib by MFA and expanded the gene pool through biochemical connections they possessed to one-node or two-node distant genes

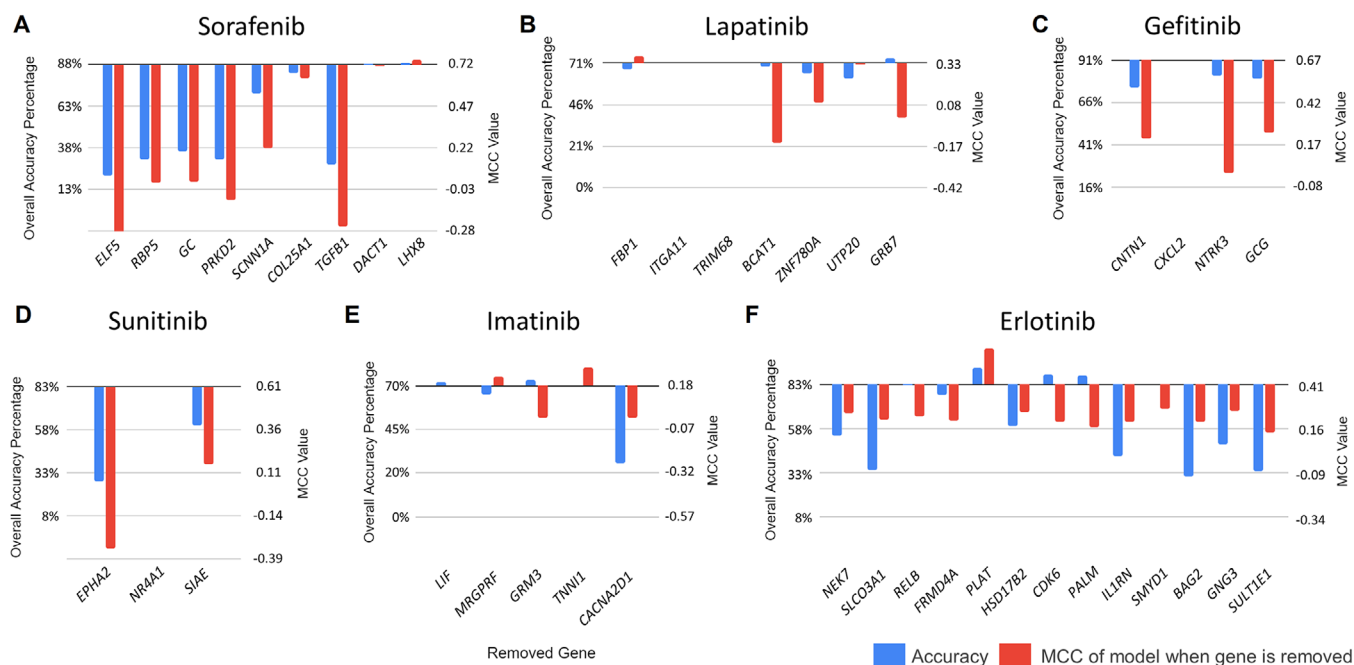


FIGURE 4 Effect of removal of individual genes from signature on overall accuracy using patient tumour data. The patient classification accuracy and MCC of the strongest performing PE models are altered upon the removal of each component gene listed. These PE TKI gene signatures are (A) sorafenib [PE-Sor], (B) lapatinib [PE-Lap], (C) gefitinib [PE-Gef], (D) sunitinib [PE-Sun], (E) imatinib [PE-Ima] and (F) erlotinib [PE-Erl]. Blue and red bars denote the overall accuracy and MCC of the model after gene removal, respectively

with reduced growth of triple negative breast cancer.⁶⁶ Removal of *ITGA11* or *TRIM68* did not alter **PE-Lap** accuracy (Figure 4B).

PE-Erl consisted of *NEK7*, *SLCO3A1*, *RELB*, *FRMD4A*, *HSD17B2*, *CDK6*, *PALM*, *IL1RN*, *SMYD1*, *BAG2*, *GNG3* and *SULT1E1* and was used to predict chemotherapy response in NSCLC patients. *BAG2* and *SULT1E1* are novel biomarkers of erlotinib efficacy, as removal of either gene led to imbalanced predictions of sensitive patients by this signature. Overexpression of *BAG2* has been associated with poor disease-specific survival in lung cancer,⁶⁷ whereas the *SULT1E1* polymorphism rs4149525 has been associated with shortened overall survival in NSCLC.⁶⁸ This model originally contained *PLAT*, which when eliminated from the erlotinib data set of 43 patients significantly increased in overall accuracy (+10%) and MCC (+0.21) of the model predictions (Figure 4F). *PLAT* was therefore considered a false positive result from ML, and therefore eliminated from gene signature. Our post hoc analysis demonstrated that the majority of genes (75%) in **PE-Erl** were associated with the NSCLC phenotype.

3.4 | Performance of PE SVM signatures on sex-stratified patients

Previous studies have suggested that females may be more sensitive to TKI treatment than males.^{69,70} We therefore

stratified TKI model performance by sex in the GSE61676 data set, which provided patient sex information along with response (19 male [3 sensitive] and 24 female [6 sensitive] patients). Considering all patients, **PE-Erl** predicted patient response with an MCC of 0.41 and 83% overall accuracy (42% and 93% accurate in patients sensitive and resistant to this drug, respectively). In males alone, **PE-Erl**'s overall accuracy was lower (76%), with MCC notably decreased to 0.11, as **PE-Erl** did not predict individuals who were sensitive or resistant to the drug as accurately (27% and 85%, respectively). In females, **PE-Erl** performed better than for the full data set, with 85% OA (MCC = 0.56), of which resistance was predicted with 99% accuracy and sensitivity was predicted with 42% accuracy (Table S6). This indicates that the **PE-Erl** signature more precisely captures factors that contribute to greater sensitivity in females.

The predictive performance of erlotinib PE model **PE-Erl** to the NSCLC data set GSE61676 was higher in female patients than male patients (0.45 greater MCC; 9% greater OA). This was consistent with the possibility that **PE-Erl** contains gene(s) distinguishing sex-differentiated sensitivity to the drug. Of the 12 genes comprising **PE-Erl**, independent removal of *RELB* and *CDK6* features from the model notably reduced accuracy of the predicted response in female patients who were sensitive to the drug. *RELB* has previously been identified as a sex-discriminatory candidate gene in trichostatin A-treated chronic lymphocytic

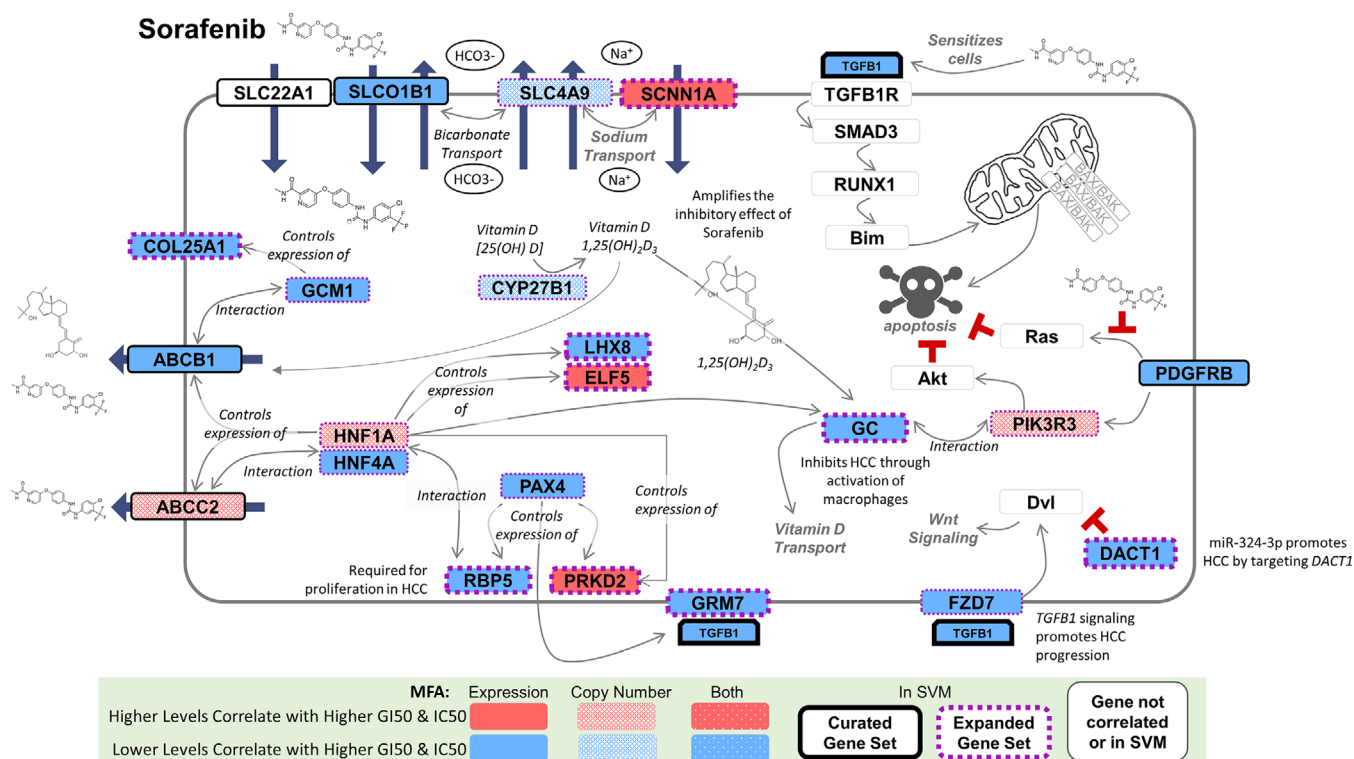


FIGURE 5 Schematic of the pathway-extended genes in the sorafenib model PE-Sor. The best performing sorafenib model **PE-Sor** is a nine-gene model consists of a single curated gene (*TGFB1*) and eight genes selected by pathway extension (*ELF5*, *RBP5*, *GC*, *PRKD2*, *SCNN1A*, *COL25A1*, *DACT1* and *LHX8*). This cell schematic provides context of the cellular mechanisms of action and/or known relationships between genes with a documented impact on sorafenib activity ('curated' genes; black borders) and those genes selected by pathway extension (purple borders). Genes with grey borders are neither curated nor pathway-extended genes and are simply present to give context between genes and their known cellular functions. Thicker borders specify those genes in the **PE-Sor** model, whereas gene colour coding indicates how GE and/or CN correlated to sorafenib GI_{50} by MFA

leukaemia cells due to repressed expression in resistant male cells, but upregulation in resistant female cells.⁷¹ *RELB* also possesses pro-survival functions across multiple cancer types^{72–74} and has been identified as a prognostic biomarker for NSCLC patients.⁷⁵ Overall, *RELB* is a top candidate gene to explain the improved accuracy of **PE-Erl** in female NSCLC patients.

3.5 | AUC-weighted ensemble model predictions

Ensemble learning consolidates hypotheses of multiple models to potentially improve predictive performance.⁷⁶ For ensemble learning, each model's AUC was computed and used to weigh predictions made for each model within the ensemble.⁷⁷ There were four SVMs for sorafenib possessing strong predictive accuracy with patient-derived expression data. Therefore, all were used for ensemble averaging. For the other TKIs, ensemble learning combined the top- and second-best performing SVMs (Table 2). Ensemble averaging improved both

OA and MCAQ4C for erlotinib (OA: 84% [+1%]; MCC: 0.45 [+0.04]) and sorafenib (OA: 91% [+3%]; MCC: 0.79 [+0.07]). For patients with the same predicted outcome in $\geq 75\%$ of cases after ensemble learning, overall accuracy exceeded 80% for all TKIs except lapatinib. Discordant consensus predictions between multiple signatures for the same drug (majority outcome occurred $< 75\%$ for each patient) exhibited lower overall accuracy.

4 | DISCUSSION

Pathway-extended GE signatures generally improved accuracy of predicted patient responses to specific TKIs. Compared to signatures composed solely of literature curated genes, PE signatures revealed previously unknown gene loci that contributed to drug response and, on average, had consistently better predictive performance. Aside from higher OA, the prediction accuracy for both sensitive and resistant patient groups (measured by MCC) was consistently more balanced. For example, **Cur-Lap** was the sole curated model with higher OA than

TABLE 2 Models used in the ensemble averaging analysis of patient data

TKI	Gene signatures (SVM: C; σ)	AUC	MCC	Sensitive	Resistant	Overall	Consensus [†]	Non-consensus [†]
Erlotinib	1. NEK7, SLC03A1, RELB, FRMD4A, HSD17B2, CDK6, PALM, IL1RN, SMYD1, BAG2, GNG3, SULF1E1 (1000; 100)	0.61	0.45	44%	94%	84%	89%	57%
	2. RET, HDGF, B3GNT5, BAG2, DSP, CAPN1, MAF, BCL2, MAP2K6, RPL13A, PTPRZ1, OLIG2 (10,000; 100)	0.47						
Sorafenib	1. ELF5, RBP5, GC, PRKD2, SCNN1A, COL25A1, TGFBI, DACT1, LHX8 (100,000; 1000)	0.88	0.79	81%	96%	91%	94%	83%
	2. ELF5, RBP5, GC, PRKD2, SCNN1A, GRM7, COL25A1, TGFBI, DACT1, LHX8 (100,000; 1000)	0.85						
	3. CPTIC, DOPEY2, KRT26, DGUOK, DLCL, CYP11B1, CALCA, MAPK1, ANK3, KRAS, FURIN, OR2A14 (10,000; 10,000)	0.86						
	4. FURIN, HTR3D, LAMAI, STMN2, SLITRK3, CACNAIS, SCN2A, CCL5, TRIM32 (100,000; 100,000)	0.93						
Lapatinib	1. FBP1, ITGAI1, TRIM68, BCAT1, ZNF780A, UTP20, GRB7 (10; 10)	0.70	0.33	63%	74%	71%	55%	80%
	2. S100A12, API5, GRHL1, TASIR1, TUBB1, CORO1A (100,000; 100)	0.56						
Sunitinib	1. EPHA2, NR4A1, SIAE (100,000; 100,000)	0.78	0.40	83%	58%	67%	91%	29%
	2. SCN3B, MED29, MPST, TSCI, AHR, CARD9, RPL3 (100,000; 100,000)	0.79						
Imatinib	1. LIF, MRGPRF, GRM3, TNNI7, CACNA2D1 (100,000; 100)	0.55	0.27	47%	83%	57%	85%	20%
	2. WASF3, TNNI7 (10,000; 10,000)	0.65						
Gefitinib	1. CNTN1, CXCL2, NTRK3, GCG (10,000; 10)	0.5	0.39	89%	50%	82%	100%	33%
	2. BDNF, PRKCB (10,000; 100)	0.69						

[†]Consensus and non-consensus predictions are when the ensemble predicts the same outcome for a patient \geq and $<75\%$ of the time, respectively. Ensemble averaging amalgamates predictions from numerous SVMs for an individual TKI, weighted by AUC (indicated). Each SVM signature included within ensemble predicted the response of each patient treated with its associated TKI, and the majority prediction was used of that of the ensemble. Overall, the accuracy of the ensemble prediction was equivalent to or greater than any individual model within it.

its PE counterpart; however, its predictions were more skewed resulting in lower MCC. Furthermore, both MCC and overall accuracy were increased by AUC-weighted ensemble averaging of multiple PE models for sorafenib and erlotinib. Except for lapatinib, the highest OAs were evident in patients receiving a 'consensus' prediction (where $\geq 75\%$ of predictions made by the models in the ensemble predicted the same outcome for a patient). The improved predictive performance of PE SVMs, both individual and as ensembles of models, suggests that the genes within these signatures may refine the predominant mechanisms of both sensitivity and resistance to TKI therapy. PE gene models may be more useful in selecting chemosensitivity regimens for patients compared to models solely consisting of previously implicated genes known to respond to a specific chemotherapy.

Pathway extension and the inclusion of pathway-related genes allowed for a larger pool of genes involved in ML. We avoided overfitting⁷⁸ by pre-filtering these genes based on correlation with GI_{50} . Furthermore, independent validation was determined by the identity and expression level of these features in patients treated with these drugs. Signatures containing pathway-related genes produced higher performing SVM signatures, consistent with the possibility that optimal molecular indicators of chemo-response may identify genes upstream or downstream of, or are interactors with, previously known cancer biomarkers. Generating SVMs from curated genes assures that features selected do not arise from statistical association alone. Generating PE SVMs required systematic selection of genes with established relationships to curated genes. For the best-performing PE gene signatures, most signature genes validated in the present study had been independently associated with abnormalities of expression, CN or mutation in these tumour types (see Additional References, Supporting Materials). Expanded signatures could potentially assist in the identification of novel biomarkers of chemo-response in these tissues.

Primary and secondary genes in PE gene signatures can offer context for drug responses without predicate literature support. The relationships between curated genes and genes selected through pathway extension for sorafenib are illustrated in Figure 5. The vitamin D transporter encoded by *GC* is a major determinant of the response to this drug, as overall prediction accuracy is decreased by 52% upon its removal from **PE-Sor** (Figure 4A). In fact, *GC* is two nodes distant from multiple curated genes (*ABCB1*, *ABCC2* and *HNF1A*, among others [Figure 3A]). The *ABCB1* transporter has been implicated in sorafenib-related toxicities based on efflux efficiency.^{79,80} *ABCB1* also carries out efflux of Vitamin D₃,⁸¹ and the 1,25-dihydroxy-vitamin D₃ isoform (or 1.25D) activates *ABCB1* expression.⁸² Vitamin D is converted to this 1.25D isoform by *CYP27B1*, which

is one-node distant from *ABCB1*. Similarly, *GC* binds specifically to 1.25D, which puts *GC* one-node distant from *ABCB1*. The growth inhibitory effect of sorafenib has been shown to be amplified by 1.25D.⁵⁶ Together, these network connections provide context that integrates functions and roles of individual genes of the tumour response to sorafenib. The PE signatures will be useful for understanding drug toxicity, although it was not explicitly a goal of this study. The importance of *GC* in **PE-Sor** may explain why a lower sorafenib dose is effective for treatment. Supplemental vitamin D₃ reduces toxicity to sorafenib at this lower dose in differentiated thyroid carcinoma that is non-responsive to iodine therapy.⁵⁶

The best performing SVMs for TKIs shared several common genetic pathways. Multiple PE models contained genes related to NOD-like receptor signalling (erlotinib: *NEK7*, *RELB*), PI3K-AKT signalling pathway (erlotinib: *CDK6*, *GNG3*; lapatinib: *ITGAI1*; sunitinib: *EPHA2*, *NR4A1*) and Ras-Raf-MEK-ERK pathway (erlotinib: *CDK6*, *RELB*; sorafenib: *TGFB1*; sunitinib: *EPHA2*, *NR4A1*). Aberrant NOD-like receptor signalling drives carcinogenesis,⁸³ whereas numerous cancer therapies target either or both of *PI3K* and *AKT*.^{84–86} The Ras-Raf-MEK-ERK pathway involves several protein kinases activated by tyrosine kinase receptors, with oncogenic mutations most prominently affecting Ras and B-Raf within the pathway.⁸⁷ These pathways, which are disrupted broadly among different cancers, are implicated across numerous high-performing ML models predicting TKI response.

Several pathway-extended (*EPHA2*, *PRKD2* and *PDGFRB*) and curated (*CDK6* and *ABL2*) gene products extrapolated from the highest performing signatures were bound to kinases based on a proteomic analysis of target selectivity for 243 kinase inhibitors on 259 distinct tyrosine kinases.⁸⁸ Few tyrosine kinase target genes from this proteomic analysis for TKIs in the current study exhibited correlations between GI_{50} and either GE or CN (<20° threshold; Table S7). *RET* was the only SVM gene implicated in the response to a TKI for both GE and protein (sorafenib; Concentration- and Target-Dependent Selectivity of 0.515; Klaeger et al.⁸⁸). Therefore, expression of genes that are either positively or inversely correlated with drug response is generally unrelated to quantification of proteins that directly interact with the kinases themselves. If absence of signature genes from those corresponding to proteomic analysis is not attributable to either experimental or specific cell lines used, then signature GE is more likely indirectly regulated by gene products that are selective for most TKIs. Many genes in the PE SVMs were two nodes distant from curated biomarkers, which is consistent with the possibility that these represent common control points in the regulation

of drug responses. In this regard, such control points exhibit behaviour similar to state-cycle attractors of self-organizing systems.⁸⁹ From a ML perspective, the dimensionality of the SVM model is reduced, avoiding overfitting, by substituting these control point genes for curated genes. Improvement in the prediction accuracy for both the sensitive and resistant patient categories might also be a consequence of these biomarkers being control points for *multiple* curated genes. Consider two curated genes that are ‘controlled’ or regulated by the same two node biomarker, where inclusion of one of these improves accuracy for detecting drug sensitivity, and the other improves detection of resistance. Substituting the controlling gene for both curated genes in the PE signature might improve accuracy of detection of both outcomes.

Transferability of these cell line-based models to other independent cell line data sets was also evaluated.⁹⁰ PE TKI models were analysed using data from the Sanger Genomics of Drug Sensitivity in Cancer Project (GDSC), including RNA-seq-derived cancer cell line-derived GE data (E-MTAB-3983; ArrayExpress) based on IC_{50} values of cell lines in CancerRxGene.⁹¹ Using median IC_{50} to distinguish sensitivity from resistance, the top SVM that we derived for each TKI could not significantly separate cell lines sensitive and resistant to the same drug in GDSC (MCC from 0 to 0.19; OA ranging from 50 to 58%). Altering the IC_{50} thresholds did not significantly change these results. When applying this analysis to cell lines from specific tissues used in the derivation of the specific TKI signatures, **PE-Ima** was more accurate for seven imatinib-treated cell lines derived from intestinal tumours (OA of 69%; MCC of -0.41). The disparity in performance between the training and testing data sets may be related to differences in the expression patterns in different tissue types, or batch effects. IC_{50} measurements for the same cell line and drug are known to vary significantly between studies, especially when the cell line is drug insensitive,⁹² which may contribute to the poor correlation between results of both data sets.

Transferability of SVMs to different patient data sets may also be confounded by several other limitations of applying ML models derived from cell line expression to predict responses to the same drugs using patient GE data. By contrast with tumours, cancer cell lines tend to have a stable genetic profile when grown under controlled culturing conditions. Consequently, they tend to lack the genetic heterogeneity present in many tumour types,⁹³ particularly during progression, which often occurs concomitant with evolution of acquired chemotherapy resistance.⁹⁴ Cancer cell lines also lack extracellular matrix, which contributes to tumour growth, migration and invasion in vivo. These differences may challenge prediction accuracy of cell line-based SVMs using patient GE and/or CN. Clinical

outcome measures within patient data sets were not consistent between different studies of the same tumour type. Finally, the cell line GE data used for training originated in this study solely from breast cancer, whereas patient tumour GE data were also derived from other cancer types.

5 | CONCLUSIONS

The enhanced performance of chemotherapy response models developed using pathway extension (over curated-only models) suggests that an interaction between a drug and its target may not directly relate with drug response; sensitivity could also be caused by a cellular event downstream of the drug-target interaction. PE models derived in this study demonstrated strong efficacy in selecting relevant genes, identifying novel molecular biomarker candidates, and predicting patient responses to TKIs. Strong-performing PE models appear to predict chemotherapy response in a cancer type-specific fashion, as many pathway-related genes selected by SVM software as novel candidate biomarkers of TKI efficacy were already prognostic biomarkers for the cancer type patients within the testing set were afflicted with. Ensemble averaging of multiple PE SVMs improved predictive accuracy in most cases and was found to be most commonly correct when predictions were highly consistent across each model constituting the ensemble. **PE-Erl** was also shown to have greater accuracy when considering solely female NSCLC patients. Interestingly, *RELB*, a feature in this signature, had previously demonstrated sexually dimorphic expression upon cancer treatment. The process of including pathway-related genes in biochemically inspired gene signatures can produce highly specific and accurate SVMs. PE models may have practical value, both in identifying novel biomarkers of chemosensitivity and in selecting effective chemotherapeutic agents.

ACKNOWLEDGEMENT

The authors acknowledge Compute Canada and Shared Hierarchical Academic Research Computing Network (SHARCNET) for a high-performance programming grant and computing facilities.

ORCID

Peter K. Rogan  <https://orcid.org/0000-0003-2070-5254>

REFERENCES

1. Pazdur R. Response rates, survival, and chemotherapy trials. *J Natl Cancer Inst.* 2000;92(19):1552-1553.
2. Thigpen JT, Vance RB, Khansur T. Second-line chemotherapy for recurrent carcinoma of the ovary. *Cancer.* 1993;71(4 Suppl):1559-1564.

3. Huisman C, Smit EF, Giaccone G, Postmus PE. Second-line chemotherapy in relapsing or refractory non-small-cell lung cancer: a review. *J Clin Oncol*. 2000;18(21):3722-3730.
4. Schiller JH, Harrington D, Belani CP, et al. Comparison of four chemotherapy regimens for advanced non-small-cell lung cancer. *N Engl J Med*. 2002;346(2):92-98.
5. Daemen A, Griffith OL, Heiser LM, et al. Modeling precision treatment of breast cancer. *Genome Biol*. 2013;14(10):R110.
6. Dorman SN, Baranova K, Knoll JHM, et al. Genomic signatures for paclitaxel and gemcitabine resistance in breast cancer derived by machine learning. *Mol Oncol*. 2016;10(1):85-100.
7. Mucaki EJ, Zhao JZL, Lizotte DJ, Rogan PK. Predicting responses to platin chemotherapy agents with biochemically-inspired machine learning. *Signal Transduct Target Ther*. 2019;4:1.
8. Kannaiyan R, Mahadevan D. A comprehensive review of protein kinase inhibitors for cancer therapy. *Expert Rev Anticancer Ther*. 2018;18(12):1249-1270.
9. Jeltsch M, Leppänen V-M, Saharinen P, Alitalo K. Receptor tyrosine kinase-mediated angiogenesis. *Cold Spring Harb Perspect Biol*. 2013;5(9).
10. Butti R, Das S, Gunasekaran VP, Yadav AS, Kumar D, Kundu GC. Receptor tyrosine kinases (RTKs) in breast cancer: signaling, therapeutic implications and challenges. *Mol Cancer*. 2018;17(1):34.
11. Paul MK, Mukhopadhyay AK. Tyrosine kinase - role and significance in cancer. *Int J Med Sci*. 2004;1(2):101-115.
12. Biscardi JS, Ishizawa RC, Silva CM, Parsons SJ. Tyrosine kinase signalling in breast cancer: epidermal growth factor receptor and c-Src interactions in breast cancer. *Breast Cancer Res*. 2000;2(3):203-210.
13. Hartmann JT, Haap M, Kopp H-G, Lipp H-P. Tyrosine kinase inhibitors - a review on pharmacology, metabolism and side effects. *Curr Drug Metab*. 2009;10(5):470-481.
14. Haouala A, Zanolari B, Rochat B, et al. Therapeutic drug monitoring of the new targeted anticancer agents imatinib, nilotinib, dasatinib, sunitinib, sorafenib and lapatinib by LC tandem mass spectrometry. *J Chromatogr B Analyt Technol Biomed Life Sci*. 2009;877(22):1982-1996.
15. Ferguson FM, Gray NS. Kinase inhibitors: the road ahead. *Nat Rev Drug Discov*. 2018;17(5):353-377.
16. Kumar R, Crouthamel M-C, Rominger DH, et al. Myelosuppression and kinase selectivity of multikinase angiogenesis inhibitors. *Br J Cancer*. 2009;101(10):1717-1723.
17. Wu P, Nielsen TE, Clausen MH. FDA-approved small-molecule kinase inhibitors. *Trends Pharmacol Sci*. 2015;36(7):422-439.
18. Aziz MA, Serya RAT, Lasheen DS, et al. Discovery of potent VEGFR-2 inhibitors based on furopyrimidine and thienopyrimidine scaffolds as cancer targeting agents. *Sci Rep*. 2016;6:24460.
19. Juan O, Popat S. Treatment choice in epidermal growth factor receptor mutation-positive non-small cell lung carcinoma: latest evidence and clinical implications. *Ther Adv Med Oncol*. 2017;9(3):201-216.
20. Yun C-H, Boggon TJ, Li Y, et al. Structures of lung cancer-derived EGFR mutants and inhibitor complexes: mechanism of activation and insights into differential inhibitor sensitivity. *Cancer Cell*. 2007;11(3):217-227.
21. McKillop D, Partridge EA, Kemp JV, et al. Tumor penetration of gefitinib (Iressa), an epidermal growth factor receptor tyrosine kinase inhibitor. *Mol Cancer Ther*. 2005;4(4):641-649.
22. Hantschel O, Rix U, Superti-Furga G. Target spectrum of the BCR-ABL inhibitors imatinib, nilotinib and dasatinib. *Leuk Lymphoma*. 2008;49(4):615-619.
23. Lin Y-L, Meng Y, Jiang W, Roux B. Explaining why Gleevec is a specific and potent inhibitor of Abl kinase. *Proc Natl Acad Sci USA*. 2013;110(5):1664-1669.
24. Johnston SRD, Leary A. Lapatinib: a novel EGFR/HER2 tyrosine kinase inhibitor for cancer. *Drugs Today*. 2006;42(7):441-453.
25. Lovly CM, Shaw AT. Molecular pathways: resistance to kinase inhibitors and implications for therapeutic strategies. *Clin Cancer Res*. 2014;20(9):2249-2256.
26. Kim YR, Kim SY. Machine learning identifies a core gene set predictive of acquired resistance to EGFR tyrosine kinase inhibitor. *J Cancer Res Clin Oncol*. 2018;144(8):1435-1444.
27. Serkova NJ, Eckhardt SG. Metabolic imaging to assess treatment response to cytotoxic and cytostatic agents. *Front Oncol*. 2016;6:152.
28. Mok TS, Wu Y-L, Thongprasert S, et al. Gefitinib or carboplatin-paclitaxel in pulmonary adenocarcinoma. *N Engl J Med*. 2009;361(10):947-957.
29. Rosell R, Carcereny E, Gervais R, et al. Erlotinib versus standard chemotherapy as first-line treatment for European patients with advanced EGFR mutation-positive non-small-cell lung cancer (EORTC): a multicentre, open-label, randomised phase 3 trial. *Lancet Oncol*. 2012;13(3):239-246.
30. Gainor JF, Shaw AT. Emerging paradigms in the development of resistance to tyrosine kinase inhibitors in lung cancer. *J Clin Oncol*. 2013;31(31):3987-3996.
31. He L, Zhu H, Zhou S, et al. Wnt pathway is involved in 5-FU drug resistance of colorectal cancer cells. *Exp Mol Med*. 2018;50(8):101.
32. Ye Q, Liu K, Shen Q, et al. Reversal of multidrug resistance in cancer by multi-functional flavonoids. *Front Oncol*. 2019;9:487.
33. Baty F, Joerger M, Früh M, Klingbiel D, Zappa F, Brutsche M. 24h-gene variation effect of combined bevacizumab/erlotinib in advanced non-squamous non-small cell lung cancer using exon array blood profiling. *J Transl Med*. 2017;15(1):66.
34. Pinyol R, Montal R, Bassaganyas L, et al. Molecular predictors of prevention of recurrence in HCC with sorafenib as adjuvant treatment and prognostic factors in the phase 3 STORM trial. *Gut*. 2019;68(6):1065-1075.
35. Guarneri V, Dieci MV, Frassoldati A, et al. Prospective biomarker analysis of the randomized CHER-LOB study evaluating the dual anti-HER2 treatment with trastuzumab and lapatinib plus chemotherapy as neoadjuvant therapy for HER2-Positive Breast Cancer. *The Oncologist*. 2015;20(9):1001-1010.
36. McWeeney SK, Pemberton LC, Loriaux MM, et al. A gene expression signature of CD34+ cells to predict major cytogenetic response in chronic-phase chronic myeloid leukemia patients treated with imatinib. *Blood*. 2010;115(2):315-325.
37. Massarweh S, Tham YL, Huang J, et al. A phase II neoadjuvant trial of anastrozole, fulvestrant, and gefitinib in patients with newly diagnosed estrogen receptor positive breast cancer. *Breast Cancer Res Treat*. 2011;129(3):819-827.
38. Moeckel S, Meyer K, Leukel P, et al. Response-predictive gene expression profiling of glioma progenitor cells in vitro. *PloS One*. 2014;9(9):e108632.
39. Eisenhauer EA, Therasse P, Bogaerts J, et al. New response evaluation criteria in solid tumours: revised RECIST guideline (version 1.1). *Eur J Cancer*. 2009;45(2):228-247.

40. Zhao JZL, Mucaki EJ, Rogan PK. Predicting ionizing radiation exposure using biochemically-inspired genomic machine learning. *FI000Res*. 2018;7:233.
41. Bagchee-Clark AJ, Mucaki EJ, Whitehead T, Rogan PK. Pathway-extended multigene expression signatures of chemotherapy responses to tyrosine kinase inhibitors: supporting data and program code. Published online August 10, 2020. <https://zenodo.org/record/3843516>
42. Abdi H, Williams LJ. Principal component analysis. *Wiley Interdiscip Rev Comput Stat*. 2010;2(4):433-459.
43. Lê S, Josse J, Husson F. FactoMineR: an R package for multivariate analysis. *J Stat Softw*. 2008;25(1).
44. Mucaki EJ, Baranova K, Pham HQ, et al. Predicting outcomes of hormone and chemotherapy in the molecular taxonomy of breast cancer international consortium (METABRIC) study by biochemically-inspired machine learning. *FI000Res*. 2016;5:2124.
45. Chicco D, Jurman G. The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC Genomics*. 2020;21(1):6.
46. Hare KJ, Hartmann B, Kissow H, Holst JJ, Poulsen SS. The intestinotrophic peptide, glp-2, counteracts intestinal atrophy in mice induced by the epidermal growth factor receptor inhibitor, gefitinib. *Clin Cancer Res*. 2007;13(17):5170-5175.
47. Tognon C, Knezevich SR, Huntsman D, et al. Expression of the ETV6-NTRK3 gene fusion as a primary event in human secretory breast carcinoma. *Cancer Cell*. 2002;2(5):367.
48. Li Z, Tognon CE, Godinho FJ, et al. ETV6-NTRK3 fusion oncogene initiates breast cancer from committed mammary progenitors via activation of AP1 complex. *Cancer Cell*. 2007;12(6):542-558.
49. Eger G, Papadopoulos N, Lennartsson J, Heldin C-H. NR4A1 promotes PDGF-BB-induced cell colony formation in soft agar. *PLoS ONE*. 2014;9(9):e109047.
50. Wykosky J, Gibo DM, Stanton C, Debinski W. EphA2 as a novel molecular marker and target in glioblastoma multiforme. *Mol Cancer Res*. 2005;3(10):541-551.
51. Wielgat P, Trofimiuk E, Czarnomysy R, Braszko JJ, Car H. Sialic acids as cellular markers of immunomodulatory action of dexamethasone on glioma cells of different immunogenicity. *Mol Cell Biochem*. 2019;455(1-2):147-157.
52. Ho JCY, Cheung ST, Poon WS, Lee YT, Ng IOL, Fan ST. Down-regulation of retinol binding protein 5 is associated with aggressive tumor features in hepatocellular carcinoma. *J Cancer Res Clin Oncol*. 2007;133(12):929-936.
53. Zhu Y, Cheng Y, Guo Y, et al. Protein kinase D2 contributes to TNF- α -induced epithelial mesenchymal transition and invasion via the PI3K/GSK-3 β /catenin pathway in hepatocellular carcinoma. *Oncotarget*. 2016;7(5):5327-5341.
54. Gressner OA, Gao C, Siluschek M, Kim P, Gressner AM. Inverse association between serum concentrations of actin-free vitamin D-binding protein and the histopathological extent of fibrogenic liver disease or hepatocellular carcinoma. *Eur J Gastroenterol Hepatol*. 2009;21(9):990-995.
55. Nonaka K, Onizuka S, Ishibashi H, et al. Vitamin D binding protein-macrophage activating factor inhibits HCC in SCID mice. *J Surg Res*. 2012;172(1):116-122.
56. Izhakov E, Sharon O, Knoll E, et al. A sorafenib-sparing effect in the treatment of thyroid carcinoma cells attained by co-treatment with a novel isoflavone derivative and 1,25 dihydroxyvitamin D3. *J Steroid Biochem Mol Biol*. 2018;182:81-86.
57. Piggan CL, Roden DL, Gallego-Ortega D, Lee HJ, Oakes SR, Ormandy CJ. ELF5 isoform expression is tissue-specific and significantly altered in cancer. *Breast Cancer Res*. 2016;18(1):4.
58. Gallego-Ortega D, Ledger A, Roden DL, et al. ELF5 drives lung metastasis in luminal breast cancer through recruitment of Gr1+ CD11b+ myeloid-derived suppressor cells. *PLoS Biol*. 2015;13(12):e1002330.
59. Fernando J, Sancho P, Fernández-Rodríguez CM, et al. Sorafenib sensitizes hepatocellular carcinoma cells to physiological apoptotic stimuli. *J Cell Physiol*. 2012;227(4):1319-1325.
60. Hoffman-Liebermann B, Liebermann DA. Interleukin-6- and leukemia inhibitory factor-induced terminal differentiation of myeloid leukemia cells is blocked at an intermediate stage by constitutive c-myc. *Mol Cell Biol*. 1991;11(5):2375-2381.
61. Ruan J, Liu X, Xiong X, et al. miR-107 promotes the erythroid differentiation of leukemia cells via the downregulation of Ccna2d1. *Mol Med Rep*. 2015;11(2):1334-1339.
62. Teng Y, Pi W, Wang Y, Cowell JK. WASF3 provides the conduit to facilitate invasion and metastasis in breast cancer cells through HER2/HER3 signaling. *Oncogene*. 2016;35(35):4633-4640.
63. Mertins P, Yang F, Liu T, et al. Ischemia in tumors induces early and sustained phosphorylation changes in stress kinase pathways but does not affect global protein levels. *Mol Cell Proteomics*. 2014;13(7):1690-1704.
64. Patterson M, Barske L, Van Handel B, et al. Frequency of mononuclear diploid cardiomyocytes underlies natural variation in heart regeneration. *Nat Genet*. 2017;49(9):1346-1353.
65. Johnston JR, Chase PB, Pinto JR. Troponin through the looking-glass: emerging roles beyond regulation of striated muscle contraction. *Oncotarget*. 2018;9(1):1461-1482.
66. Thewes V, Simon R, Hlevnjak M, et al. The branched-chain amino acid transaminase 1 sustains growth of antiestrogen-resistant and ER α -negative breast cancer. *Oncogene*. 2017;36(29):4124-4134.
67. Yue X, Zhao Y, Liu J, et al. BAG2 promotes tumorigenesis through enhancing mutant p53 protein levels and function. *eLife*. 2015;4:e08401.
68. Yamamoto N, Sato Y, Tamura T, et al. Genetic polymorphisms correlate with overall survival (OS) in advanced non-small cell lung cancer (NSCLC) treated with carboplatin (CBDCA) and paclitaxel (PTX). *J Clin Oncol*. 2008;26(15_suppl):8034.
69. Zeng Z, Chen H-J, Yan H-H, Yang J-J, Zhang X-C, Wu Y-L. Sensitivity to epidermal growth factor receptor tyrosine kinase inhibitors in males, smokers, and non-adenocarcinoma lung cancer in patients with EGFR mutations. *Int J Biol Markers*. 2013;28(3):249-258.
70. Schmetzer O, Flörcken A. Sex differences in the drug therapy for oncologic diseases. In: Regitz-Zagrosek V, ed. *Sex and Gender Differences in Pharmacology. Handbook of Experimental Pharmacology*. Berlin, Germany: Springer; 2012:411-442.
71. Marteau J-B, Rigaud O, Brugat T, et al. Concomitant heterochromatinisation and down-regulation of gene expression unveils epigenetic silencing of RELB in an aggressive subset of chronic lymphocytic leukemia in males. *BMC Med Genomics*. 2010;3:53.

72. Roy P, Mukherjee T, Chatterjee B, Vijayaragavan B, Banoth B, Basak S. Non-canonical NF κ B mutations reinforce pro-survival TNF response in multiple myeloma through an autoregulatory RelB:p50 NF κ B pathway. *Oncogene*. 2017;36(10):1417-1429.
73. Wang M, Tang J. RelB facilitates cell migration and invasion in breast cancer via MMP1 upregulation. *Ann Oncol*. 2018;29:ix19.
74. Mineva ND, Wang X, Yang S, et al. Inhibition of RelB by 1,25-dihydroxyvitamin D3 promotes sensitivity of breast cancer cells to radiation. *J Cell Physiol*. 2009;220(3):593-599.
75. Qin H, Zhou J, Zhou P, et al. Prognostic significance of RelB overexpression in non-small cell lung cancer patients. *Thorac Cancer*. 2016;7(4):415-421.
76. Polikar R. Ensemble learning. In: Zhang C, Ma Y, eds. *Ensemble Machine Learning: Methods and Applications*. New York, NY: Springer; 2012:1-34.
77. LeDell E, van der Laan MJ, Petersen M. AUC-maximizing ensembles through metalearning. *Int J Biostat*. 2016;12(1):203-218.
78. Pham HNA, Triantaphyllou E. The impact of overfitting and overgeneralization on the classification accuracy in data mining. In: Maimon O, Rokach L, eds. *Soft Computing for Knowledge Discovery and Data Mining*. New York, NY: Springer; 2008:391-431.
79. Qin C, Cao Q, Li P, et al. The influence of genetic variants of sorafenib on clinical outcomes and toxic effects in patients with advanced renal cell carcinoma. *Sci Rep*. 2016;6:20089.
80. Agarwal S, Elmquist WF. Insight into the cooperation of P-glycoprotein (ABCB1) and breast cancer resistance protein (ABCG2) at the blood-brain barrier: a case study examining sorafenib efflux clearance. *Mol Pharm*. 2012;9(3):678-684.
81. Margier M, Collet X, le May C, et al. ABCB1 (P-glycoprotein) regulates vitamin D absorption and contributes to its transintestinal efflux. *FASEB J*. 2019;33(2):2084-2094.
82. Tachibana S, Yoshinari K, Chikada T, Toriyabe T, Nagata K, Yamazoe Y. Involvement of vitamin D receptor in the intestinal induction of human ABCB1. *Drug Metab Dispos Biol Fate Chem*. 2009;37(8):1604-1610.
83. Velloso FJ, Trombetta-Lima M, Anschau V, Sogayar MC, Correa RG. NOD-like receptors: major players (and targets) in the interface between innate immunity and cancer. *Biosci Rep*. 2019;39(4).
84. Liu Q, Yu S, Zhao W, Qin S, Chu Q, Wu K. EGFR-TKIs resistance via EGFR-independent signaling pathways. *Mol Cancer*. 2018;17(1):53.
85. Sadeghi N, Gerber DE. Targeting the PI3K pathway for cancer therapy. *Future Med Chem*. 2012;4(9):1153-1169.
86. Morgensztern D, McLeod HL. PI3K/Akt/mTOR pathway as a target for cancer therapy. *Anticancer Drugs*. 2005;16(8):797-803.
87. Burotto M, Chiou VL, Lee J-M, Kohn EC. The MAPK pathway across different malignancies: a new perspective. *Cancer*. 2014;120(22):3446-3456.
88. Klaeger S, Heinzlmeir S, Wilhelm M, et al. The target landscape of clinical kinase drugs. *Science*. 2017;358(6367):eaan4368.
89. Kauffman S. At home in the universe: the search for laws of self-organization and complexity. *Choice Rev Online*. 1996;33(06):33-3294-3233-3294.
90. Brooks EA, Galarza S, Gencoglu MF, Cornelison RC, Munson JM, Peyton SR. Applicability of drug response metrics for cancer studies using biomaterials. *Philos Trans R Soc Lond B Biol Sci*. 2019;374(1779):20180226.
91. Yang W, Soares J, Greninger P, et al. Genomics of drug sensitivity in cancer (GDSC): a resource for therapeutic biomarker discovery in cancer cells. *Nucleic Acids Res*. 2013;41(D1):D955-D961.
92. Brooks EA, Gencoglu MF, Corbett DC, Stevens KR, Peyton SR. An omentum-inspired 3D PEG hydrogel for identifying ECM-drivers of drug resistant ovarian cancer. *APL Bioeng*. 2019;3(2):026106.
93. Schmitt MW, Prindle MJ, Loeb LA. Implications of genetic heterogeneity in cancer. *Ann N Y Acad Sci*. 2012;1267:110-116.
94. Lee S-C, Xu X, Lim Y-W, et al. Chemotherapy-induced tumor gene expression changes in human breast cancers. *Pharmacogenet Genomics*. 2009;19(3):181-192.

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of the article.

How to cite this article: Bagchee-Clark AJ, Mucaki EJ, Whitehead T, Rogan PK. Pathway-extended gene expression signatures integrate novel biomarkers that improve predictions of patient responses to kinase inhibitors. *MedComm*. 2020;1–17.
<https://doi.org/10.1002/mco.2.46>