

## ARTICLE OPEN

## Predicting ligand-dependent tumors from multi-dimensional signaling features

Helge Hass<sup>1,2</sup>, Kristina Masson<sup>1</sup>, Sibylle Wohlgemuth<sup>3</sup>, Violette Paragas<sup>1</sup>, John E. Allen<sup>1</sup>, Mark Sevecka<sup>1</sup>, Emily Pace<sup>1,4</sup>, Jens Timmer<sup>2,5</sup>, Joerg Stelling<sup>3</sup>, Gavin MacBeath<sup>1</sup>, Birgit Schoeberl<sup>1</sup> and Andreas Raue<sup>1</sup>

Targeted therapies have shown significant patient benefit in about 5–10% of solid tumors that are addicted to a single oncogene. Here, we explore the idea of ligand addiction as a driver of tumor growth. High ligand levels in tumors have been shown to be associated with impaired patient survival, but targeted therapies have not yet shown great benefit in unselected patient populations. Using an approach of applying Bagged Decision Trees (BDT) to high-dimensional signaling features derived from a computational model, we can predict ligand dependent proliferation across a set of 58 cell lines. This mechanistic, multi-pathway model that features receptor heterodimerization, was trained on seven cancer cell lines and can predict signaling across two independent cell lines by adjusting only the receptor expression levels for each cell line. Interestingly, for patient samples the predicted tumor growth response correlates with high growth factor expression in the tumor microenvironment, which argues for a co-evolution of both factors in vivo.

*npj Systems Biology and Applications* (2017)3:27; doi:10.1038/s41540-017-0030-3

## INTRODUCTION

The combination of Herceptin® with chemotherapy demonstrated a dramatically increased survival benefit for a subset of women with HER2 amplified advanced breast cancer, which ultimately led to FDA approval in 1998.<sup>1</sup> Since then, targeted cancer therapies have become an accepted therapeutic modality for the treatment of cancer and have contributed to a decrease in cancer related mortality.<sup>2</sup> However, the benefit of targeted therapies to date has been restricted to 5–10% of solid tumors addicted to oncogenes.<sup>3–5</sup> Identifying these relatively rare patients via predictive diagnostic tests relying on genomic biomarkers has created Precision Medicine.<sup>6–8</sup>

Retrospective analyses of several clinical studies of breast, gastric or lung adenocarcinoma identified increased receptor and/or growth factor expression as prognostic markers for patients with poor prognosis, which highlights the role of ligand-induced signaling as oncogenic drivers.<sup>9–12</sup> Here we aim to decipher what drives ligand-induced proliferation.

We present the first comprehensive proliferation screen across 58 cell lines comparing to which extent the growth factors EGF (epidermal growth factor), HRG (heregulin), IGF-1 (insulin growth factor 1) and HGF (hepatocyte growth factor) induce cell proliferation. We find that about half of the cell lines do not respond to any of the ligands whereas the other half of the cell lines respond to a least one ligand. We compare the observed ligand-induced proliferation with the response to treatment with antibodies targeting the ErbB receptor family members, a subfamily of four closely related receptor tyrosine kinases (RTKs): EGFR (ErbB1), HER2/c-neu (ErbB2), HER3 (ErbB3) and HER4 (ErbB4) as well as the insulin growth factor receptor (IGF-1R) and the hepatocyte growth factor receptor (Met). Not surprisingly, the

antibodies targeting the respective RTK inhibit ligand-induced proliferation. The antibodies also inhibited basal proliferation in some cell lines that do not respond to exogenous ligand addition, which could be driven by autocrine signaling.

The need has been recognized for computational approaches to deal with the complexity of signal transduction and its dysregulation in cancer to ultimately understand drug activity.<sup>13–17</sup> Large collections of genetic and genomic data led to efforts to disentangle the complex mechanisms using machine-learning algorithms.<sup>18–21</sup> It was previously shown that simulated patient-specific signaling responses derived from mechanistic signaling models using RNA sequencing data from patient biopsies can be robust biomarkers that are predictive of patient outcome.<sup>22</sup> Here, we combined machine learning and mechanistic modeling to predict which cell lines proliferate in the presence of ligand. We used RNA sequencing data as inputs into a comprehensive mechanistic model capturing the ErbB, IGF-1R and Met signaling pathways. Our novel approach uses simulated signaling features and mutation status of a specific cell line as inputs into a Bagged Decision Tree, which predicts whether tumor cells proliferate in the presence of a growth factor. We achieved a substantial gain in accuracy compared to predictions based on RNA sequencing data alone by inclusion of simulated signaling features such as the area under curve of distinct heterodimers and phosphorylated S6 for in vitro models.

Applying this approach to patient data, the prediction of ligand-dependent tumor samples based on mRNA data from The Cancer Genome Atlas (TCGA) revealed that colorectal and lung cancer are the two indications most responsive to EGF, which agrees with the approval of EGFR inhibitors in these indications. In addition, the prediction of responders in patient samples revealed a correlation

<sup>1</sup>Merrimack Pharmaceuticals, Inc., Cambridge, MA 02139, USA; <sup>2</sup>Institute of Physics, University of Freiburg, Freiburg, Germany; <sup>3</sup>Department of Biosystems Science and Engineering and SIB Swiss Institute of Bioinformatics, ETH Zuerich, Zuerich, Switzerland; <sup>4</sup>Celgene, San Francisco, CA 94158, USA and <sup>5</sup>BIOSS Centre for Biological Signalling Studies, University of Freiburg, Freiburg im Breisgau, Germany  
Correspondence: Andreas Raue (araue@merrimack.com)

Received: 2 May 2017 Revised: 23 August 2017 Accepted: 28 August 2017

Published online: 20 September 2017

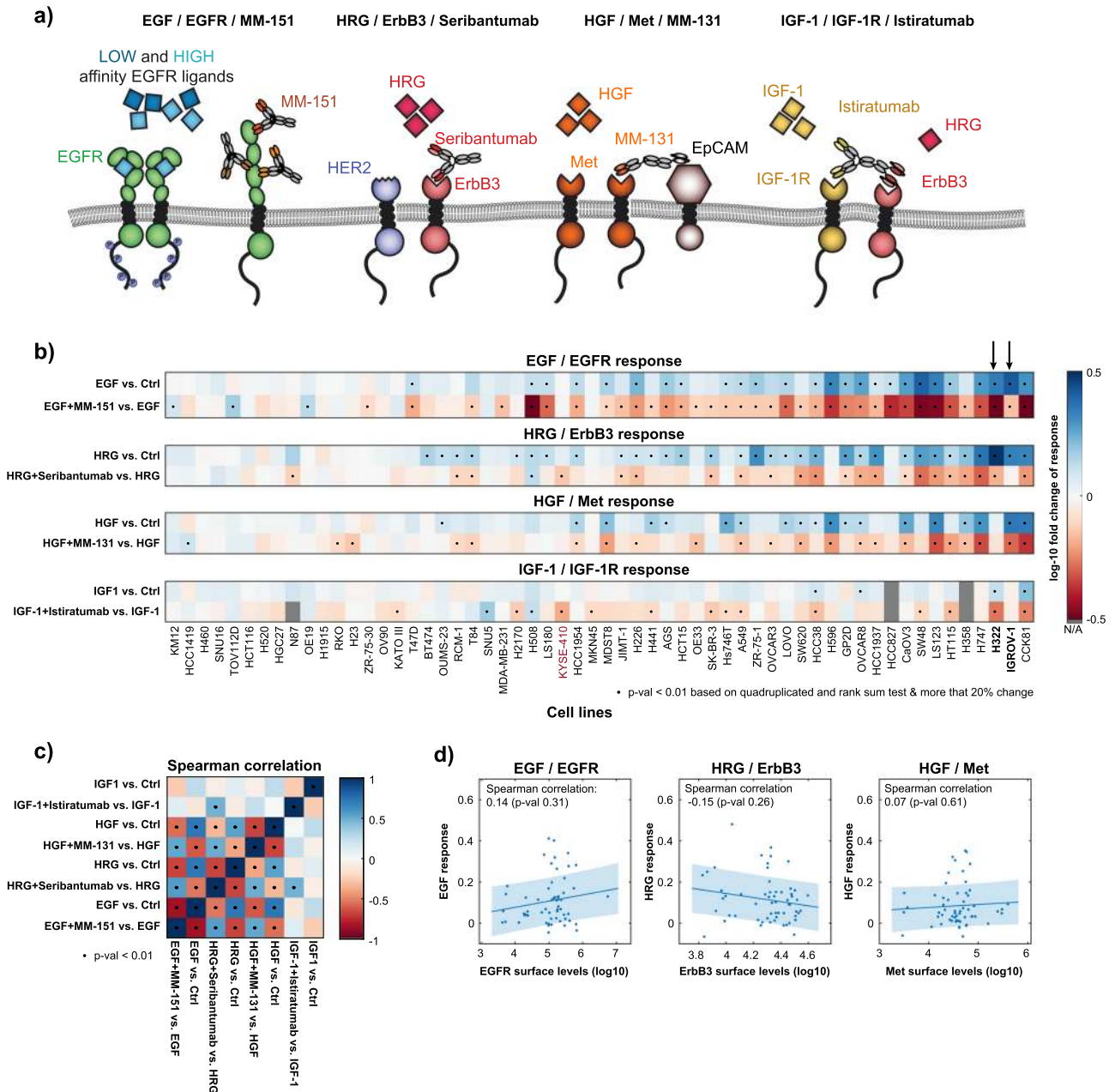
between predicted tumor growth and measured ligand expression in the tumor microenvironment, which argues for a co-evolution of ligand production and the ability of the tumor cells to respond to stimulation.

## RESULTS

### In vitro proliferation screen

To investigate growth factor-induced proliferation we screened a panel of 58 cancer cell lines (10 ovarian cancer, 11 breast cancer, 13 lung cancer, 11 gastric cancer, and 23 colorectal cancer cell

lines) for response to the exogenously added ligands EGF, HRG, HGF, and IGF-1 (Supplementary Fig. 1) that bind to EGFR, ErbB3, Met, and IGF-1R, respectively. In addition to ligand stimulation, cells were also treated with ligand blocking antibodies: MM-151, an oligoclonal therapeutic composed of three monoclonal antibodies targeting EGFR;<sup>23</sup> Seribantumab (MM-121), a monoclonal antibody targeting ErbB3;<sup>16</sup> MM-131, a bispecific antibody co-targeting Met and EpCAM;<sup>24</sup> and Istaratumab (MM-141), a bispecific antibody co-targeting IGF-1R and ErbB3.<sup>25</sup> Fig. 1a illustrates the RTKs, their corresponding ligands and the mechanism of action of the ligand blocking antibodies. Proliferation was quantified in a 3D spheroid formation assay at the 3-day time



**Fig. 1** Proliferation screen across 58 cell lines. **a** Ligand/Receptor and antagonistic antibodies used in the in vitro proliferation screen. **b** Results of the proliferation screen across 58 cell lines. Dots mark a significant increase in ligand induced proliferation or decrease in the presence of ligand plus antibody. The ligand effect is normalized to the medium control, whereas the antibody plus ligand effect is relative to ligand alone. The two cell lines marked with an arrow, as well as five additional cell lines that were not included in the proliferation screen, were used to train the computational model to signaling data. **c** Correlation pattern of ligand and antibody effects across all cell lines. **d** Linear correlation of receptor expression to ligand induced proliferation. The proliferation in response to ligand (y-axis) is displayed as log<sub>10</sub>-fold change with respect to day 0. The receptor surface levels (x-axis) are absolute measurements of receptors/cell by qFACS on a log<sub>10</sub>-scale

point (Fig. 1b) by measuring ATP content as surrogate for cell number (CellTiter-Glo® assay). Response was classified as positive if the signal at the 3-day time point was more than 20% above the respective control, plus being significant at a confidence level  $\alpha = 0.05$  (measured in quadruplicates, Wilcoxon rank-sum test). Per this screen approximately 45% of cell lines responded to EGF, 55% of cell lines responded to HRG, 33% of cell lines responded to HGF and 7% of cell lines responded to IGF-1. The low response rate to IGF-1 in this proliferation screen may reflect the presence of IGF-1 in the low-serum medium and the modest absolute inhibition point to the importance of IGF-1 mediated signaling for survival rather than for proliferation.<sup>26</sup> We and others observed a generally weaker MAPK activation via IGF-1R (see Fig. 3c) compared to the other growth factors in the screen.<sup>27,28</sup> Further, the CellTiter-Glo® assay relies on metabolic function and hence can be limited as readout for IGF-1 stimulation.<sup>29</sup>

Figure 1b shows the response to treatment with ligand in combination with the respective blocking antibody compared to the ligand effect alone. Depending on the ligand treatment, 5–17% of cell lines were ligand non-responsive, but the antibodies inhibited basal proliferation, which is indicative of autocrine driven proliferation. Even though IGF-1 did not induce a proliferative response in most cell lines, MM-141 inhibited proliferation in about 19% of the cell lines indicating that IGF-1 might be present in low-serum medium.

Investigation of correlations between the ligand and antibody responses across all cell lines revealed a checkerboard pattern of significant positive correlations between EGF, HRG, and HGF as well as anti-correlations of those ligands and their respective antibody responses (Fig. 1c). This suggests a general trend that cell lines are either responsive to multiple ligands and their respective antibodies (right hand side of Fig. 1b), or are generally non-responsive to any given ligand or antibody (left hand side of Fig. 1b). For IGF-1/IGF-1R, the only significant correlation was observed between Istaratumab treatment and Seribantumab treatment. This can be attributed to both antibodies (co-)targeting ErbB3 and, therefore, some cell lines respond to both Istaratumab and Seribantumab independent of an IGF-1 effect (see, e.g., KYSE-410 cell line in Fig. 1b). The general lack of correlation patterns for IGF-1/IGF-1R responses as were observed for the ErbB family and HGF/Met can be explained by the lack of IGF-1 induced proliferation in this screen.

In the following, we will focus on the question of how ligand dependence can be predicted. A necessary condition for response to any given ligand is the presence of its respective receptor. First, we used a univariate analysis (Fig. 1d) and found that receptor surface levels measured by qFACS do not correlate significantly with the respective ligand response. Based on this data, a simple linear model cannot stratify responsiveness. Next, we investigated whether a multi-pathway signaling model featuring the complex receptor interactions as well as the cross-talk between the mitogen-activated protein (MAP) kinase and the phosphoinositide 3-kinase (PI3K) signaling pathways can be used to predict the phenotypic response. Specifically, signaling features like the area under curve (AUC), quasi steady-state and the signal amplitude of receptor homo- or heterodimers and downstream components were considered as inputs into a decision tree classification algorithm.

#### Multi-pathway computational model

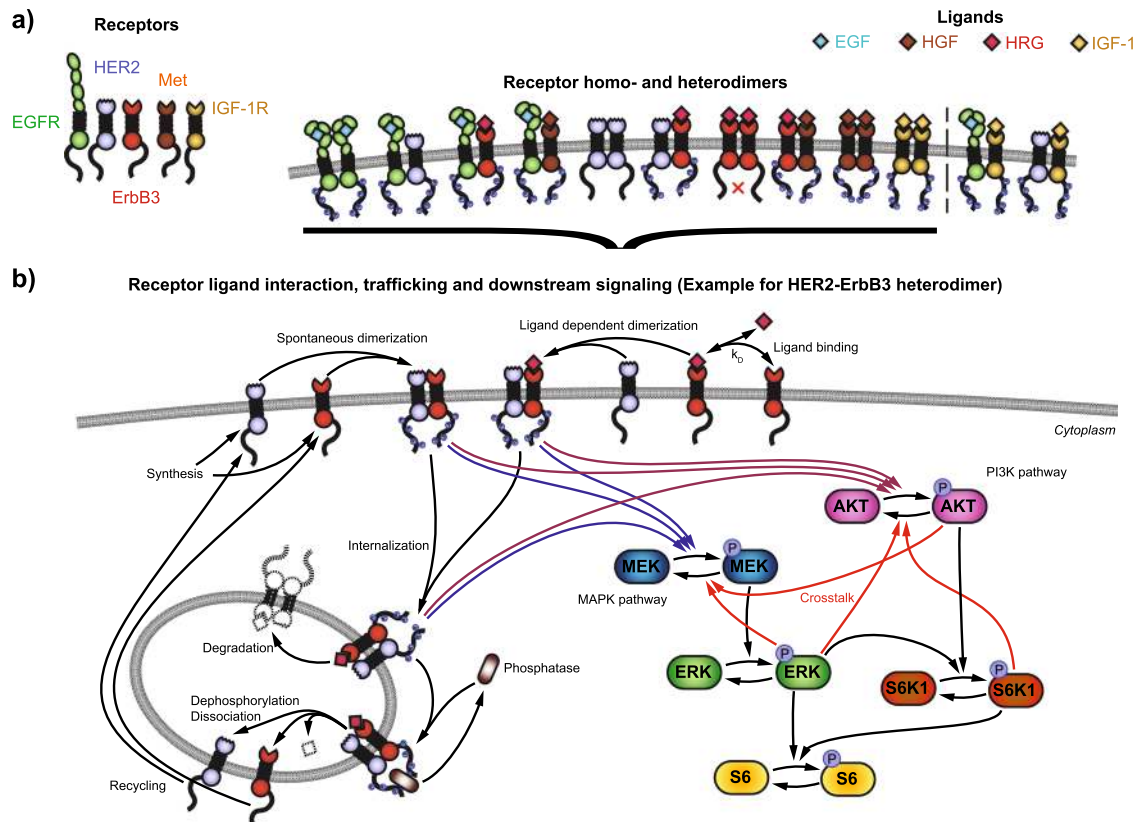
To construct a comprehensive signal transduction model that could be used to predict proliferation in response to growth factors for all 58 cell lines, we built on a previously published model of ErbB receptor signaling.<sup>16</sup> We extended the computational model to include IGF1-R and Met (Fig. 2a) as well as 12 homo and heterodimers for which biological evidence can be found.<sup>30–33</sup> Our analysis considers EGFR, HER2, ErbB3, Met and IGF-

1R homodimers as well as the heterodimers EGFR-HER2, EGFR-ErbB3, EGFR-Met, HER2-ErbB3, ErbB3-Met, IGF-1R-IGF-1R, EGFR-IGF-1R, and HER2-IGF-1R. The latter two were later removed from the computational model without impacting the model performance. Figure 2b depicts the structure of the model for the example of a signaling HER2-ErbB3 heterodimer. The complete model consists of 62 differential equations and replicates the model structure shown in Fig. 2b for each of the considered ten homo- and heterodimers. In short, receptors bind ligand with published dissociation constants ( $K_D$ ). Bound receptors can form homo and heterodimers and subsequently undergo endocytosis. After internalization, the receptor dimers can get either dephosphorylated and recycled to the cell surface or they get degraded in the lysosomes.<sup>34,35</sup> Downstream of the receptor, all homo- and heterodimers except the ErbB3-homodimer, which cannot transphosphorylate due to its lack of intrinsic kinase activity, can activate the MAP kinase cascade as well as the PI3K/AKT pathway. ERK and AKT phosphorylation converge in the phosphorylation of S6K1 and S6. Several known feedback mechanisms between the pathways<sup>27</sup> were implemented in the computational model. Mathematical details, executable code to simulate the model and instructions to replicate our findings are available in the [supplementary materials](#) and on [biomodels.org](#).

The computational model is constructed with the aim to capture the signaling dynamics of key components of the signaling pathway including receptor homo- and heterodimerization. It is not intended to be a complete compendium of all the known molecular interactions.<sup>33,36,37</sup> Size and complexity of the computational model were chosen to reflect the available experimental data, and to facilitate efficient computation. This is particularly important during model parameter calibration, which uses parameter estimation algorithms to match the available experimental data as closely as possible (see Methods section for details). Mutations were not implemented in the computational signaling model as they appear to increase the signaling baseline but not necessarily the signaling dynamics.<sup>38</sup> However, the mutation status for each cell line was used in the machine learning classification.

For model calibration, phosphoproteomic time course data from protein microarrays<sup>39</sup> for the receptor phosphorylation as well as for phospho-MEK, phospho-ERK, phospho-AKT, and phospho-S6 across all seven cancer cell lines (H322M, BxPc-3, A431, BT-20, ACHN, ADRr, and IGROV-1) were used. Only the two cell lines H322M and IGROV-1 were included in the cell line proliferation screen in Fig. 1. These seven cancer cell lines represent different cancer indications (lung adenocarcinoma, pancreatic, epidermoid, breast and ovarian cancer) and were selected based on the molecular diversity with respect to the mutation status and differences in receptor expression. A key challenge for building computational models that can describe and predict signaling dynamics of different cell lines is to limit the number of model parameters that are specific to one cell line.<sup>40</sup> In this case, it was possible to restrict all kinetic rate constants to the same value and to adjust only the receptor expression for individual cell lines. Due to the analytically calculated basal activation levels of all homo- and heterodimers as well as of the downstream components, which were derived from steady-state constraints,<sup>41</sup> the receptor expression impact the signaling response throughout the model. Therefore, the individual receptor expression of each cell line enables distinct model responses upon ligand stimulation. The receptor expression was measured using quantitative flow-cytometry (qFACS) in combination with RNA sequencing data (see Methods section). The model can accurately describe the time course data of seven training cell lines, with 85.7% of the data points within two standard deviations of the model uncertainty (see Fig. 3 for a selection of the data and Suppl. Figs. 11–40 for a comprehensive comparison of model simulations and experimental data).





**Fig. 2** Structure of computational signaling model. **a** The receptors EGFR, HER2, ErbB3, Met, and IGF-1R can form several homo and heterodimers after ligand binding. **b** In the model, receptors are synthesized and either dimerize spontaneously or bind a ligand to form homo- and hetero-dimers, which results in trans-phosphorylation of the receptors. Activated receptors signal downstream and are prone for internalization, which leads to either degradation or dephosphorylation by a phosphatase followed by recycling to the cell surface. Downstream, the MAPK and PI3K cascade activate S6K1 and ultimately converge in the phosphorylation of S6. The MAPK and PI3K signaling pathways are interconnected via multiple crosstalk mechanisms

### Validation of the computational signaling model

Based on the trained model, predictions were generated for two independent validation cell lines (BT-474 M3, MDA-MB-231) and compared to the experimental data. The goodness of the predictions for the validation cell lines was equivalent to the goodness of fit of the training cell lines (Fig. 3b, Supplementary Figs. 35–40 for model fits to the available data). These simulation results validate that receptor expression is sufficient to predict signaling features of independent cell lines that were not used for model training. In addition, we generated model predictions that were based on random receptor surface levels, by taking non-matching values from randomly selected cell lines used in the cell viability assay (see Table 3). The decline in goodness of fit was on average 30% and statistically significant ( $p = 8.5 \times 10^{-9}$ , see Supplementary Fig. 2). These results illustrate the importance of receptor expression and their ratios to capture the distinct signaling features observed for each cell line.<sup>42</sup>

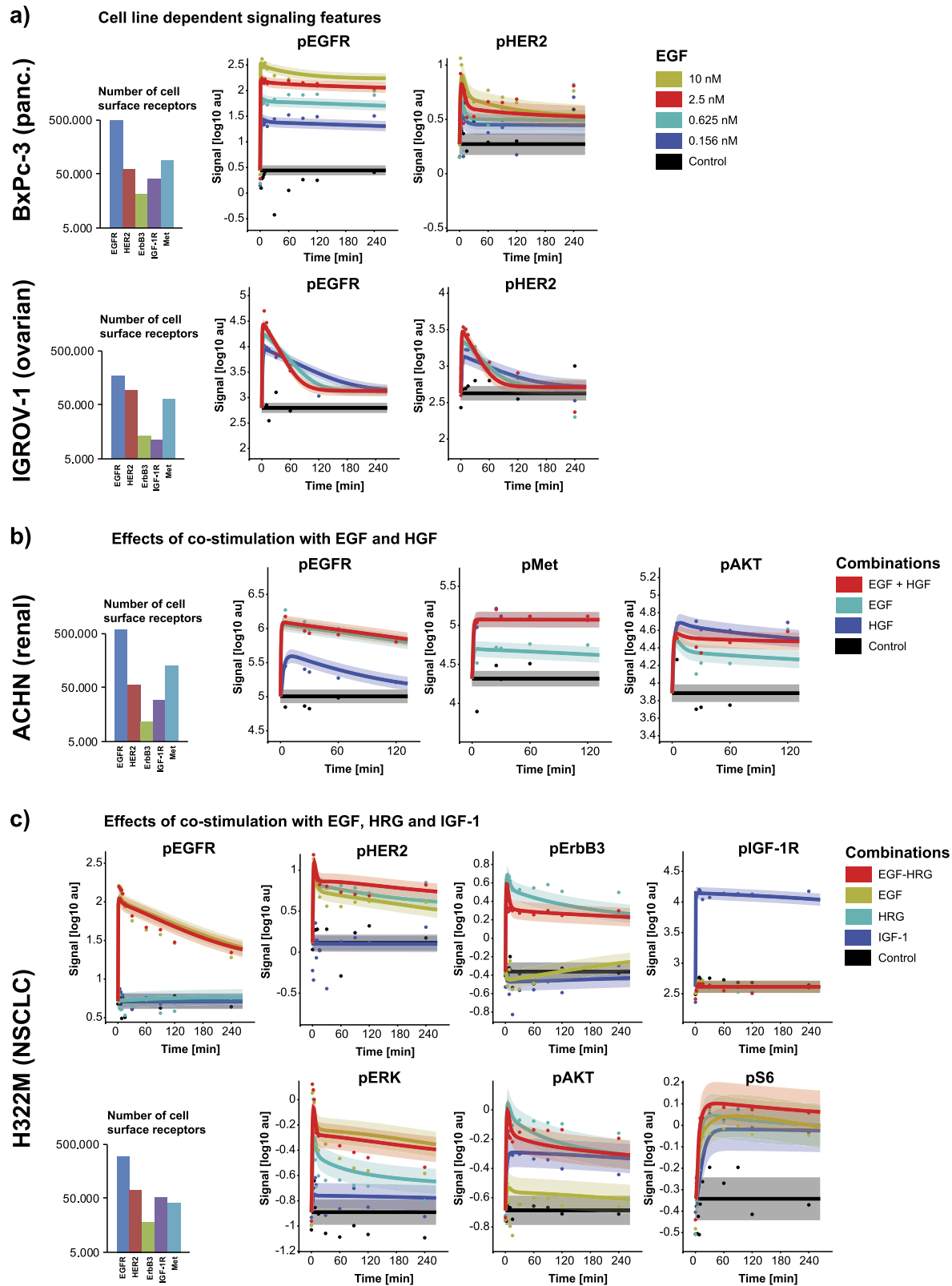
To further validate the presented model structure with its multiple receptor heterodimers, a simplified model lacking any heterodimerization capabilities was trained to the experimental data. In this setting, all receptors could signal downstream through homodimerization, disregarding the non-functional kinase unit of the ErbB3 receptor. Even with 39 parameters less, the reduced model had a goodness of fit impediment with associated  $p$ -value  $< 1.e-15$  in the corresponding likelihood-ratio test, showing the significant improvement of the computational model by including receptor heterodimerization. Besides, concordance of the basal receptor levels obtained via analytic steady state equations, reflecting the proposed receptor trafficking, was

assessed through extensive measurements of basal total and phosphorylation levels in 39 breast cancer cell lines.<sup>28</sup> A good correlation, especially for the ErbB receptor family, was found (see Supplementary Fig. 3), confirming the calibrated model parameters constituting cell-dependent steady states.

To further test the robustness and applicability of the model to ligands not included in the original training-set, we compared the predicted receptor activation patterns in response to different ligands of the EGFR-ligand family, such as Betacellulin (BTC). To this end, previously published<sup>16</sup> time-resolved data of the ADRr cell line for EGF and BTC with ligand concentration range between 0.1 nM and 10 nM was reanalyzed with the current model (Suppl. Fig. 4a). Differences in the ligand binding affinities of each ligand to the EGF receptor as well as different homo- as well as heterodimerization kinetics were sufficient to describe the experimental data (Supplementary Figs. 4b, c). BTC induces a stronger EGFR homodimerization compared to the stronger EGFR-HER2 heterodimerization induced by EGF (Supplementary Fig. 4d). These differences in EGFR homo and heterodimerization with HER2 were previously described.<sup>43,44</sup>

### Importance of receptor homo and heterodimers

Further insights into growth factor signaling and signal processing by the cancer cells can be gained by analyzing the computational model and why it can capture the distinct signaling dynamics across cell lines. This analysis revealed the importance of different homo and heterodimers in encoding information as a function of the ligand(s) present. The largest effect of heterodimerization on signal output is seen within the ErbB family. The interplay



**Fig. 3** Importance of receptor surface levels for model response, shown for a selection of calibration cell lines. **a** Cell line dependent signaling features: Model response to EGF stimulation of two different cell lines resulting in sustained or transient receptor phosphorylation in the BxPc-3 and IGROV-1 cells. Their respective receptor surface levels are shown on the left. The model fits are represented by the colored lines with respective uncertainties (67% confidence intervals) as shades. Data points are shown as dots in the same color. **b** Model fits for the cell line ACHN stimulated with HGF, EGF and the combination. **c** Model response to co-stimulation of EGF plus HRG in comparison to the stimulation with EGF, HRG or IGF-1 alone in H322M cells

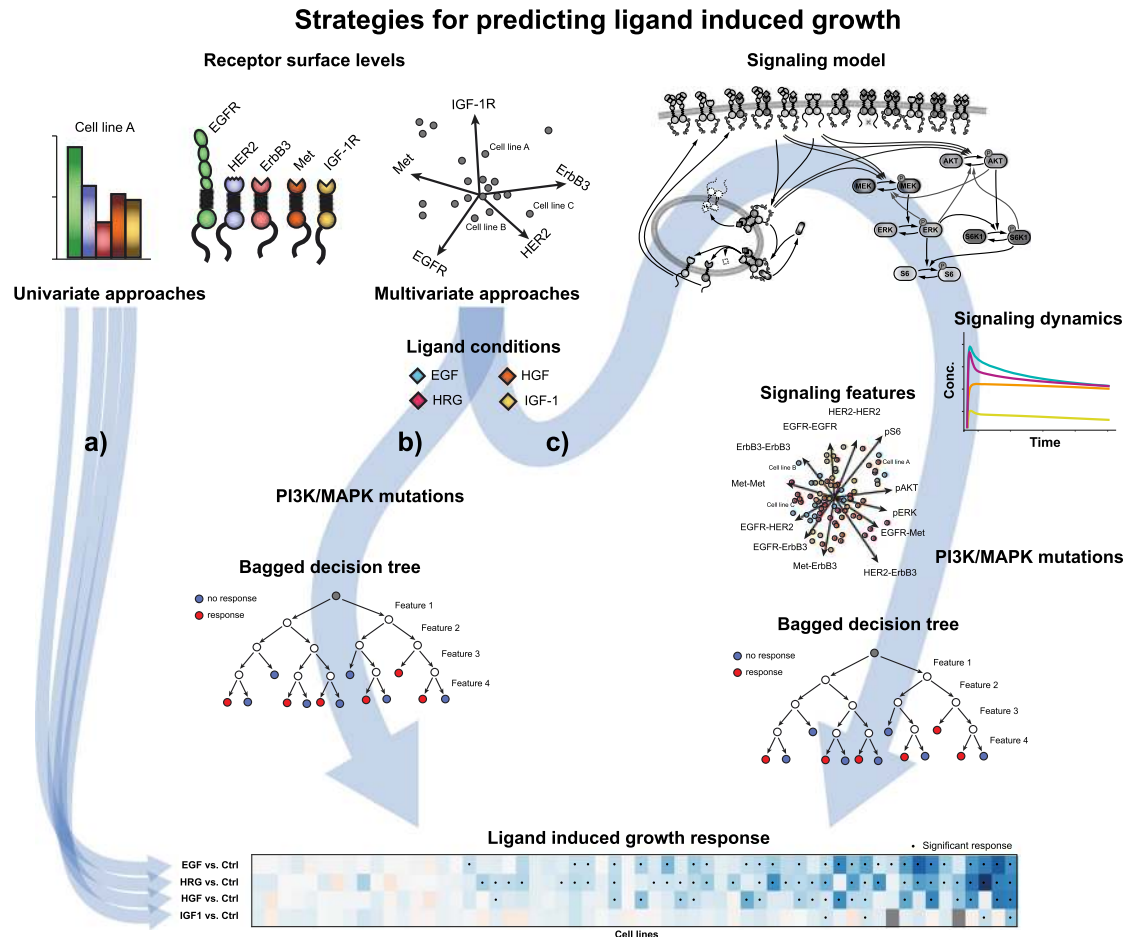
between the receptors explains the slower, more sustained receptor activation in response to EGF in the BxPc-3 cells, which are characterized by a high ratio of EGFR to other receptors (Fig. 3a). In contrast to the BxPc-3 cells, the IGROV-1 cells are characterized by low EGFR expression compared to other receptors leading to the observed transient and early activation of EGFR and HER2. For the ACHN cell line, we generated signaling data in response to EGF, HGF as well as to the combination of EGF and HGF. The computational model reveals sophisticated feedback regulation between the MAPK and PI3K pathways, e.g., reduced AKT activation comparing EGF and HGF co-stimulation to HGF only (Fig. 3b). In Fig. 3c another example is depicted: when EGF and HRG are present, EGFR and ErbB3 compete for HER2. The ligand combination results in reduced phospho-ErbB3 levels due to a dominant binding of HER2 to EGFR in the presence of EGF (Fig. 3c). We argue that mechanistic understanding of changes in receptor stoichiometry based on individual ligands or ligand mixtures can cause non-obvious signaling responses and is required to understand the ultimate phenotypic response.

IGF-1 is distinct from the other growth factors in our screen. IGF-1 displayed a much weaker ability to induce proliferation and similarly the effect of IGF-1 signaling in co-stimulation experiments is distinct to the HRG/EGF or HGF/EGF co-stimulation

experiments. The time-course data for co-stimulation of IGF-1 with either EGF or HRG did not show deviations from the respective stimulation with IGF-1 alone (see Supplementary Fig. 14). Consequently, the model parameters referring to the heterodimerization of IGF-1R (see heterodimers involving IGF-1R in Fig. 2a) with other receptors could be set to zero without a significant decline in the goodness of fit.

RNAseq and signaling features are predictive of phenotype

The calibrated mechanistic signaling model can be used to simulate signaling features for the stimulation with EGF, HRG, HGF or IGF-1 for all cell lines of the cell viability screen only using their receptor expression as inputs. To connect signaling features derived from the computational model to the phenotypic response observed in the cell proliferation screen, we applied a machine learning approach. Based on different sets of input features that were selected based on their prediction ability (see below), we trained bootstrap-aggregating (bagged) decision trees (BDTs). BDTs are highly efficient for multivariate analysis and allow for a comprehensive interpretation of the chosen features.<sup>45,46</sup> Therein, a multitude of trees are trained, with each single tree aiming at discriminating growing from non-growing cells based



**Fig. 4** Strategies for predicting ligand-induced phenotypic response. Based on the receptor expression of individual cancer cell lines, either a univariate or multivariate approach can be used to predict the phenotypic response to ligand stimulation. **a** Univariate approaches relate the respective receptor expression to the observed ligand induced proliferation for each of the four ligands separately. **b–c** Multivariate approaches such as bagged decision trees (BDTs) relate high-dimensional feature sets to the observed phenotype. **b** In this case the feature set consists of the five receptor surface levels as well as information about the respective ligand stimulation and mutation status. **c** The calibrated and validated signaling model allows to simulate the expected signaling dynamics for each individual cell line based on its receptor expression and ligands present. Based on the mechanistic knowledge that the signaling model incorporates, it can expand the initial five-dimensional feature set to a 12-dimensional feature set. This expanded feature set, together with information about mutation status is now connected to the observed growth responses by a bagged decision tree

on the provided feature space. At each node, a tree divides the data through selected features that yield the best improvement in signal-to-noise. Combining the ensemble of trees, high-dimensional and non-linear regions in feature space prone for cell growth are obtained and can be used for prediction. More details can be found in the [supplementary materials](#) and in Supplementary Fig. 5.

To investigate the contribution of dynamic signaling features derived from the computational model on the predictive performance of the machine learning approach, we considered two different sets of input features (additional sets are reported in Supplementary Fig. 6). The first feature set contained the receptor expression, KRAS and PI3K mutation status and ligand treatment as binary input (Fig. 4b).<sup>47,48</sup> The second feature set did not contain the receptor expression explicitly but only signaling features derived from the computational model based on the receptor expression and the ligand stimulation (Fig. 4c) as well as the mutation status (see Table 3). The cell line specific signaling features consist of the AUC of all phosphorylated receptor homo- and heterodimers in addition to AKT, ERK, and S6 phosphorylation. Inclusion of the fold-change of these features as well as their quasi steady-state levels were tested but did not yield substantial benefits on top of the information given by the area under curve (see Supplementary Fig. 6). Figure 4 illustrates both multivariate prediction strategies as well as the univariate analysis shown in Fig. 1d.

To evaluate the accuracy of BDT predictions, the cell lines were randomly split 500 times into training and testing sets. For each ligand, BDT training was performed on the training data for all available ligands, while efficiency was calculated on the testing cell lines for the chosen ligand only. By leaving out whole cell lines as opposed to a fraction of the total data, possible bias due to correlated responses to different ligands in the same cell line is avoided. We monitored the fraction of true predictions as a metric for the prediction accuracy. Both feature sets resulted in a better prediction of proliferation compared to random data, which results in 50% true predictions and serves as control (Fig. 5a). Exceptions were the predictions for IGF-1 and HGF stimulation using the receptor expression only, where the performance drops insignificantly below the control. Training on features derived from the computational model improved the prediction of cell proliferation significantly compared to control ( $p$ -value of  $<1.e-2$ , see Fig. 5b) for the combination of all ligands except IGF-1, while BDT predictions based on receptor expression alone did not result in a statistically significant improvement ( $p$ -value of 0.15, see Fig. 5b). The respective distributions for single ligand induced proliferation predictions can be found in Suppl. Fig. 7. The BDT training was robust with respect to the relative amount of training and testing data and to the significance threshold upon which a cell is labeled as proliferating (Supplementary Fig. 8). For IGF-1, the low number of responders resulted in a low correlation within the events and a statistical bias of their relative amount in training or testing data. These circumstances rendered robust prediction impossible and the IGF-1 data set was excluded, e.g. in the combination of all ligands (see Fig. 5a).

One of the advantages of mechanistic computational models is that it is possible to gain insights into cellular signal processing. Thus, the importance of different model features during BDT training can be traced back to develop hypotheses about what ultimately drives proliferation. The training features are ranked by their impact on data classification, which is measured by the average gain in signal-to-noise ratio over all trees. As illustrated in Table 1, the most important features rely on the homo- and heterodimerization stoichiometry of the ErbB receptor family as well as on the downstream signaling. A more detailed overview is given in Fig. 5c, which illustrates the data from the proliferation screen and the model-derived features that proved important during training of the BDT (see Table 1). This coincides with prior

biological knowledge that ErbB receptors induce proliferation.<sup>13</sup> It can be observed that EGF, HGF or HRG stimulation induce very specific homo and heterodimerization patterns of EGFR and EGFR/HER2 respectively. Moreover, phospho-S6 is an important feature to predict proliferation in the presence of HRG and HGF, both mainly activating the PI3K pathway. Its importance might be a result of the crosstalk between the MAPK and PI3K pathways and the convergence of both pathways in S6 phosphorylation.<sup>49</sup> Moreover, PI3K mutation status and EGFR heterodimerization patterns help to identify clusters of EGF dependent cell lines. RAS mutation status and differences in activation of downstream targets further predict proliferation after HRG and HGF stimulation. Independent of the KRAS mutation status, HRG induces increased levels of AKT phosphorylation while inducing the same amount of S6 phosphorylation as HGF.

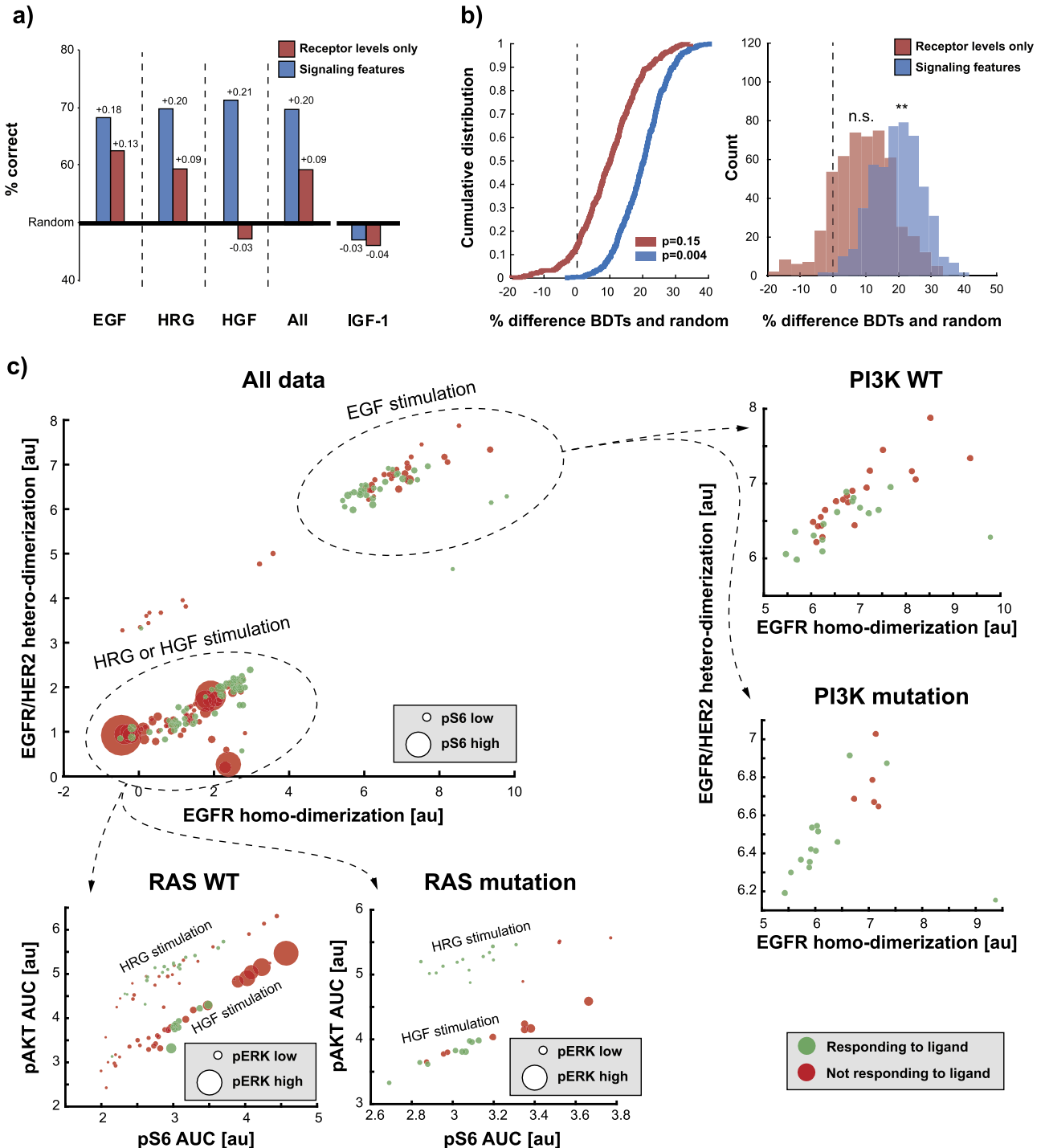
#### Application to patient data

In the previous section, we showed that a mechanistic computational model in combination with decision tree classification can predict *in vitro* proliferation. Next, we applied this novel approach to patient data. Using the data from the TCGA Research Network (<http://cancergenome.nih.gov/>), we use our model to predict if an individual patient tumor would show a proliferation response if stimulated by the ligand of interest. The patient data set includes 2909 samples from patients with breast, colorectal, lung, and ovarian cancer. The input to our model is receptor RNA expression measured by RNA sequencing for each tumor sample. The measured RNA expressions between our *in vitro* cancer cell line data and the data from patient tumors were on different scales. Therefore, we normalized the expressions to their respective means. Subsequently, the expression of EGFR, Her2, ErbB3, Met, and IGF-1R were used to perform model simulations and to extract the signaling features needed to predict ligand-dependent tumors using the previously described BDT algorithm. The number of ligand-dependent tumors differed within indications and ligand (EGF, HGF or HRG). The number of predicted ligand responsive tumors was highest for HRG followed by EGF and lowest for HGF (Fig. 6a). Lung and colorectal cancer seem to be most responsive to EGF, which is congruent with the high prevalence of EGFR mutations and overexpression in these indications.<sup>50–52</sup> In Fig. 1b we observed that ligand induced proliferation is correlated with treatment response to an antibody targeting the respective receptor. Therefore, the approvals of EGFR inhibitors in non-small cell lung cancer (NSCLC) and CRC confirm the predicted dependence on EGFR signaling.<sup>53,54</sup> In contrast, the low dependence of NSCLC on HGF signaling might explain the failure of Onartuzumab (MetMab), a Met blocking antibody in a Phase 3 study in NSCLC.<sup>55</sup> Similarly, EGFR inhibitors have not yet proven to result in clinical benefit in breast cancer,<sup>56</sup> which is also in agreement with the predicted low EGF dependence of breast cancer. The predicted high responsiveness of breast cancer and lung cancer to HRG seems to agree with the retrospective analysis of two clinical studies with Seribantumab and the finding that HRG expression appears to be predictive of patients responding to therapy.<sup>57</sup> Unfortunately, the precision of the predictions cannot be assessed as no outcome data to treatment with ligand blocking antibodies is available for the TCGA data set. The predicted ligand-dependence only considers the molecular makeup of individual tumors. However, a tumor that is predicted to be ligand-dependent would respond only if the respective ligands were also present in sufficient amounts in the tumor microenvironment. The local concentration of ligands, however, cannot be inferred from our analysis and it is difficult to match to the *in vitro* data that was used for model training. This impacts the predicted number of ligand-dependent tumors in this data set.

However, apart from the relative number of ligand-dependent tumors, we observed a significant correlation between ligand



## BDT prediction accuracy and efficiency



**Fig. 5** Prediction of ligand-induced proliferation using BDTs. **a** Ratio of true predictions after BDT training with simulated signaling features or receptor expression only, compared to random predictions in the presence of EGF, HRG, IGF or HGF. **b** For 500 random splits of training and testing cell lines, the BDT outcome is compared to random growth assessment as histogram and cumulative density function, showing the significant improvement due to mechanistic modeling. **c** Data of in-vitro cell viability screen showing proliferation response (green) or no significant response (red) in different 2D representations of the feature space

expression and the predicted response to ligands (Fig. 6b, additional data in Supplementary Fig. 9). The predicted ligand-dependent tumor samples from patients with breast and colorectal cancer display statistically significant higher (*t*-test) amounts of the corresponding ligand compared to the predicted ligand

independent tumors, if we compare the mean expression. This suggests that tumors that express ligands evolved to be sensitive to ligands, or vice versa. In our opinion, this is an indirect proof that the model predictions can be applied to data from patient tumors and that they could be clinically relevant.



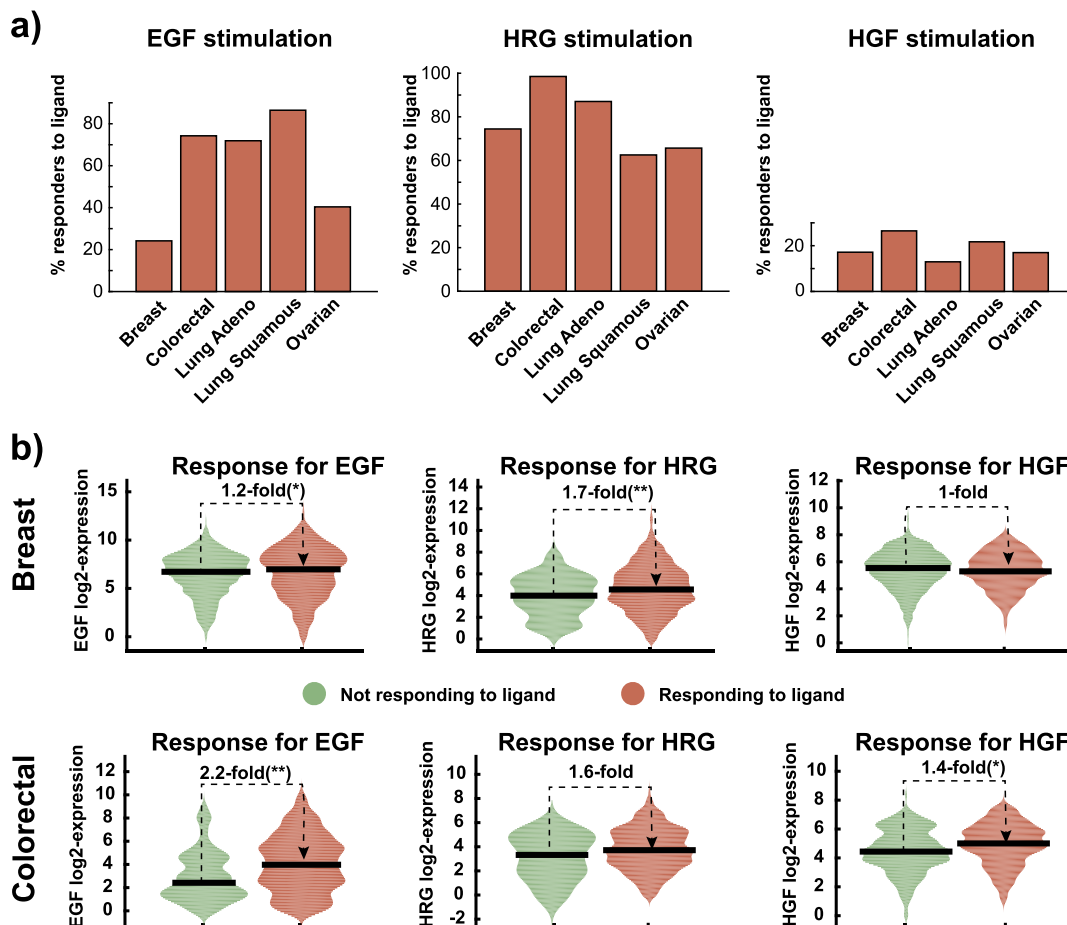
## DISCUSSION

The research to understand and therapeutically battle various cancer types has made significant progress over the last two decades, decreasing the overall mortality by roughly 2% per year since 2001.<sup>58</sup> Targeted therapy and combinations thereof have become important areas of drug development with rising FDA approval rates in the last years.<sup>59,60</sup> However, by studying cellular fate across multiple cell lines and indications, we and others have learned that complex interactions between receptors as well as positive and negative feedback regulation between signaling pathways can diminish drug efficacy. To obtain a deeper

understanding of cell response to exogenous stimuli, phenotypic responses need to be studied in the context of multiple signaling pathways as well as mutation status. In this work, we developed a computational model describing multiple signaling pathways and show that a BDT algorithm using simulated signaling features can accurately predict ligand-dependent proliferation in vitro.

The signaling model incorporates the ErbB receptor family as well as the Met and IGF-1R receptors. Parameters of the model were estimated based on a variety of time-resolved data from seven different cell lines including a wide range of ligand concentrations with comprehensive single ligand and co-stimulations. The cell lines cover a broad range of ratios of receptor expression. The ligand concentrations used in the proliferation screen were in the range of concentrations used for the signaling experiments. While retaining a good fit to the experimental data, we could keep all kinetic parameters in the signaling model constant and just vary the receptor expression to describe the experimental signaling data for each cell line. It is important to note that this is not a direct proof that the reaction rate constants are identical between the cell lines. However, our model with all rate constants set to the same values, is in line with the receptor phosphorylation and downstream signaling data. We do not argue that the model presented here is correct in all of its aspect, but we could show that the signaling dynamics predicted by the model are useful for predicting the cellular responses to ligand stimulation and that this presents an approach that should be explored further by incorporating clinical patient response

Feature	BDT importance score
pS6 AUC	0.65
EGFR homodimerization	0.56
EGFR-HER2 heterodimerization	0.56
Met-ErbB3 heterodimerization	0.55
pAKT AUC	0.43
pERK AUC	0.37
PI3K mutation status	0.35
RAS mutation status	0.30



**Fig. 6** Predicting ligand dependent tumors from the TCGA data set. **a** Predicted percentage of tumors that would respond to ligand exposure. The predictions were obtained by using the receptor RNA expression measured in breast, colorectal, lung, and ovarian cancers as inputs to our model. **b** The measured RNA expression of the ligands in predicted responders (red) vs. non-responders (green). The mean expression (black horizontal lines) and statistical significance of differences is indicated. The receptor mRNA expression is measured in transcripts per million and is displayed on a log<sub>2</sub> scale

data. The computational model not only describes the data for the seven training cell lines but also predicts the signaling responses of two additional, independent validation cell lines. We showed that the goodness of fit is dependent on the absolute receptor expression and the formation of different homo- and heterodimers. The observed variability in model response was achieved by differences in internalization, degradation and recycling rates for different receptor homo- and heterodimers. Saturation of different downstream model components is sensitive to the receptor expression. The analytically solved steady states for all cell lines were important to implement the complex receptor dimerization properties in the signaling model. Thus, ligand stimulation results in a complex re-distribution of receptor homo- or heterodimerization and facilitated distinct downstream activation patterns. The calculated steady states, especially for the ErbB receptor family, were found in concordance with measured basal total and phosphorylation levels of 39 breast cancer cell lines utilized from.<sup>28</sup> The fact that the developed computational model can accurately describe the signaling responses across multiple cell lines enables the prediction of signaling dynamics for cell lines of the proliferation screen that were not used to construct the computational signaling model.

To predict proliferation in response to ligand stimulation, we linked simulated signaling features to the data of the cell proliferation screen using a supervised machine learning approach. Tree-based classification algorithms are widely used for machine learning<sup>61–63</sup> and carve out regions in feature space that best distinguish between different data classifications, here proliferation vs. stasis upon ligand stimulation. BDTs were trained on the in vitro proliferation screen across 58 cell lines. The algorithm was trained with either receptor expression and ligand stimulation conditions as Boolean columns or with signaling features extracted from the computational model. The signaling features included the integrated area under curve of all receptor complexes as well as the phosphorylation of AKT, ERK, and S6. Both feature sets lead to a better prediction of proliferation compared to control, where response was predicted based on random data. However, the signaling features allowed for a more robust and statistically improved prediction of proliferation. In addition, the computational model allows us to gain insights into the underlying processes driving ligand-dependent proliferation. In all cases homo or heterodimerization of the ErbB receptors was important. In the case of EGF, the PI3K mutation status mattered in addition. For HRG and HGF possible RAS mutations together with AKT and S6 phosphorylation were important features to predict cell proliferation. However, the importance of features in the tree-based approach is very sensitive to the utilized data, and additional measurements are needed to infer the role of MAPK and PI3K signaling in inducing growth.

Simulated signaling features are advantageous over using receptor expression directly as input features for two reasons: First, the dynamic range of receptor activation as well as of the downstream components is described quantitatively and renders the model outputs more robust to receptor expressions, which span multiple orders of magnitude. Second, the interplay between receptors and the included feedback mechanisms adds a source of information on top of the receptor expression and ligand information alone, resulting in a non-linear input transformation that improves the detection of regions governing proliferation.

To demonstrate the applicability of this novel approach to patient samples, data from 2909 patients with breast, colorectal, lung or ovarian cancer were analyzed. For these samples, model simulations were conducted to extract signaling features required for BDT prediction of ligand-dependence. Interestingly, we observed a significant correlation between the measured ligand expression and the predicted ligand-dependence for breast and colorectal cancer. This may be a consequence of evolutionary adaption of these tumor cells due to the growth advantage from

ligand-mediated signaling. Therefore, the presence of ligands in the tumor micro-environment may be a favorable biomarker for RTK-directed drug treatment. This rationale together with possible mutations, which contributed substantially in the BDT training, are currently being explored in the clinic.<sup>57</sup>

However, both the prediction of proliferation and the mechanistic computational model have limitations. For one, machine learning in feature-space regions that are not covered by many cell lines is not efficient as illustrated by IGF-1 in the proliferation screen data-set. The mechanistic model is limited in its ability to fully reproduce data in cases of either receptor overexpression (see Supplementary Fig. 19), which probably transfers to the presence of activating mutations, as well as in cell lines harboring PI3K or RAS mutations (see Supplementary Figs. 22, 23, and 37 to 40). We also encountered computational limitations since the complexity of the mechanistic computational model is on the verge of what is currently computationally feasible. This said, more emphasis on model selection, e.g. profile likelihood-based model reduction<sup>64</sup> with additional prior knowledge may allow us to better bridge between quantitative time-resolved data and large-scale genomic and phenotypic data. To understand ligand mixtures and pathway redundancy a greater variety of single ligands, e.g., FGF and PDGF, or stimulations with ligand mixtures might aid to more accurately determine model parameters for receptor dimerization, trafficking and downstream activation in the future. Further, our model is currently limited to MAPK and PI3K pathways. Future work, should expand on this approach by including other signaling pathways, e.g., data from the JAK-STAT signaling pathway. An expanded model may improve the accuracy of the predictions as well as additional perturbations like specific gene knockdown etc. would help to improve the signaling model.

Ligand-dependence or addiction to growth factors (ligands) might prove to be far more prevalent than oncogene addiction given the high ligand prevalence in solid tumors and potentially as much more complex since multiple ligands are expressed in the tumor microenvironment. To successfully treat patients with ligand-dependent tumors with targeted inhibitors (small molecules or monoclonal antibodies) or rational combinations of targeted inhibitors, a better understanding of ligand-dependence is crucial, especially given the redundancy of signaling pathways within a tumor cell and tumor heterogeneity. We argue that the mechanistic understanding of changes in receptor stoichiometry based on individual ligands or ligand mixtures result in non-obvious signaling responses that are relevant to the ultimate phenotypic response. The work presented here demonstrates that for targeted therapies to be successful in the clinic the ligand hierarchy as well as co-dependence need to be understood and most likely require the measurement of multiple ligands and respective receptor expression in tumor biopsies. New methodologies like single cell RNAseq<sup>65</sup> may allow us in the future to characterize the clonal composition of tumors and to determine which cellular fraction is ligand-dependent and which drug combination is best suited to eliminate all ligand-dependent tumor cells.

The presented novel approach of using BDTs in conjunction with simulated signaling features is the beginning of how complex mechanistic models and large data sets can be combined to understand cell-specific complexity but also heterogeneous tumors better. We demonstrated that mechanistic computational models of signaling pathways can help bridge between large scale in vitro observations and clinical hypotheses. In the future, selected in vivo studies should be used to validate rational combination regimens. Previous efforts predicting drug sensitivity based on large and diverse data sets found that gene expression data proved most valuable, together with exploiting non-linear relationships and addition of prior knowledge of biological pathways.<sup>18,66,67</sup> Yet, significant improvement in predictions

proved to be challenging across multiple approaches and data sets. With the presented approach, mechanistic knowledge can be easily combined with known datasets from RNA sequencing, copy-number alterations and mutation information to improve the prediction of patient-individual drug response and unravel the interplay between complex signaling and cellular fate.

## METHODS

### Cell lines and reagents

All cell lines used in the viability screen were purchased from American Type Culture Collection. In brief, cells were cultured in RPMI 1640 medium (Life Technologies) supplemented with 10% fetal bovine serum (FBS, Life Technologies) and 1% penicillin-streptomycin (pen/strep, Life Technologies) at 37 °C and 5% CO<sub>2</sub>. Recombinant human heregulin-1-β1 (NRG1b) EGF domain and HGF were obtained from PeproTech, and recombinant human IGF-1 and EGF were obtained from RD Systems. Seribantumab, MM-131, Istitutumab and MM-151 was manufactured in-house by the Merrimack Pharmaceuticals Protein Engineering Department and stored at 4 °C. CellTiter-Glo® was obtained from Promega and reconstituted fresh for each experiment.

### Spheroid formation assay and in vitro screening conditions

To measure cell viability in a three-dimensional spheroid culture, cells were seeded into 384-well low-binding multi-spheroid culture plates (Scivax, USA) in the relevant growth medium supplemented with 4% FBS and 1% pen/strep at a density of 1500 cells/well. To allow for spheroid formation, plates were incubated for 24 h, after which cells were treated with ligands (5 nM HRG1, 5 nM EGF, 50 nM IGF-1, 1 nM HGF) and/or inhibitors (1 μM of Seribantumab, MM-131 and MM-151 and 0.5 μM of Istitutumab) in 4% FBS containing medium. Following 72 h of incubation, cell viability was determined by incubation with CellTiterGlo® reagent for 10 min and measuring well luminescence on an Envision (Perkin Elmer) plate reader.

### Experimental time-course data and selected cell lines

The available data consists of three different data sets, which include time-resolved concentration measurements of activated and total receptors as well as of various phosphorylated downstream targets. The largest data set comprises nine cell lines with measurements of the EGFR, HER2 and ErbB3 receptors of the ErbB family and the IGF1-receptor, together with the downstream targets ERK, AKT, S6K1 and S6. Four different ligand concentrations of EGF, HRG, and IGF-1 ranging from 0.156 to 10 nM are used and 12 measurement time points up to 240 min are taken. In addition, co-stimulations of the respective ligands are available for two of the nine cell lines. Out of the nine cell lines, six cell lines are used for model calibration, the remaining three to validate the model. Cell lines used for calibration include H322M (non-small cell lung cancer), BxPc-3 (pancreatic cancer), A431 (epidermoid cancer), BT-20 (breast cancer), ADRr (ovarian cancer) and IGROV-1 (ovarian cancer). BT474 (breast cancer), MDA-MB-231 (breast cancer) and ACHN (renal cancer) are utilized to validate the model.

Measurements after either EGF or BTC stimulation are available for one of the calibration cell lines, ADRr, with ligand concentrations between 0.11 and 9.26 nM, and 12 measurement time points up to 240 min. These include phosphorylation of EGFR, HER2, ErbB3, ERK, and AKT. Apart from that, measurements with HGF and EGF as well as their co-stimulation is available for ACHN, which is used for validation with respect to HRG and IGF-1, spanning the same concentrations and measurements up to 120 min. Therein, phosphorylated EGFR and Met phosphorylation as well as phospho-ERK and phospho-AKT are measured. The receptor concentrations in all experiments are measured by ELISA whereas the downstream components are measured by lysate microarray.

### Quantification of receptor expression in cell lines using qFACS

Cells were trypsinized, washed, and stained using fluorescently labeled antibodies, see list below. Antibodies were labeled as previously described (16). Receptor numbers were determined by assessing the antibody-binding capacity of the fluorescently labeled antibody via quantitative fluorescence-activated cell sorting. Antibody-binding capacity was determined using Simply CellularQuantumBeads (BangsLabs, Fishers, IN), per the manufacturer's instructions. List of antibodies and target receptors used in qFACS method: Erbitux (EGFR), Trastuzumab (HER2), Anti-human

HER3 Ab generated by Merrimack, A12 IgG (anti-human IGF1-R), Anti-human Met, Mouse Anti-Human EpCAM-APC (BD Biosciences, Cat# EBA-1).

### ELISA and Lysate microarray (reverse phase protein array) measurements

Matching high density lysate matrices from nine cells lines (BxPc-3, H322M, ACHN, IGROV-1, A-431, BT-20, ADRr, BT-474-M3, and MDA-MB-231) were generated for both the ELISA and lysate microarray studies. Lysate matrices were developed by treating each cell line with EGF, HRG1b1, IGF1, individually at four different doses (10 nM, 2.5 nM, 0.625 nM, and 0.156 nM), as well as in all two-way combinations of these three ligands at a single 2.5 nM dose of each ligand (a total of 15 conditions) in 96-well plates (Corning). Antibodies directed at pAKT (S473), pMEK (S217/S221), pERK (T202/Y204), p-S6K1 (T389), and pS6 (S235/S236) were obtained from Cell Signaling Technology. Cells were cultured using standard tissue culture techniques in RPMI supplemented with 10% FBS and penicillin/streptomycin. Cells were counted using a hemocytometer and seeded at 10,000 cells/well into 40x 96-well tissue culture plates (Corning). After 24–48 h, once cells had reached approximately 50% confluency, medium was aspirated and replaced with low- in medium (RPMI supplemented with 0.5% FBS and penicillin/streptomycin) (Gibco). 16–24 h later, 2x ligand solutions prepared in low-serum medium were added simultaneously to all 96 wells of each plate. Plates were returned to the incubator for the prescribed incubation times, then placed on ice to stop ligand stimulation. Medium was aspirated from all wells, cells were washed once with ice-cold PBS, and then lysed in 30 μL/well of lysis buffer. Twelve time points as well as an untreated time zero were collected for all conditions (0, 2, 4, 6, 8, 10, 15, 30, 60, 90, 120, and 240 min). Lysates were then collected and frozen. For ELISA lysates, M-PER buffer (ThermoFisher) was supplemented with protease and phosphatase inhibitors (Roche) and NaCl to a final concentration of 150 mM. For lysate microarrays, lysis buffer was prepared as previously described.<sup>68</sup>

Total and phospho ELISAs of EGFR, ErbB2, ErbB3, and IGF-1R were performed as previously described.<sup>69</sup> Protein specific antibodies were used for capture in all cases, while an anti-phospho tyrosine antibody was used for detection in all phosphor assays (see Supplementary Table 1). Antibody screening for this study and well as lysate preparation and printing were performed as previously described.<sup>68,70</sup> Arrays were printed with a 7-point dilution series for each lysate onto 16-pad slides (GraceBioLabs OnCyte Avid, Bend, OR) using the Aushon 2170 micro-arraying robot in 2x4 8-pin mode (Aushon Biosystems, Billerica, MA). Slides were washed in 100 mM Tris-HCl pH 9.0 at RT for several days, followed by 3x 5 min PBST, after which slides were spun dry. ProPlate 16-well slide modules (GraceBioLabs) were then attached. Arrays were blocked in Odyssey Blocking Buffer (LiCor, Lincoln, NE) at 4 °C for 1 h, after which blocking solution was replaced with 1:500 primary antibodies (all rabbit, see list below) mixed with 1:1000 anti-beta-actin antibody (mouse) (Sigma A1978) for 24 h at 4 °C with agitation. Slides were then washed 3x 5 min with PBST, then incubated in secondary antibodies (1:1000 anti-rabbit-800 and 1:1000 anti-mouse-680 (LiCor) in 5% BSA/PBST) at 4 °C for 5 h. Arrays were washed briefly in PBST, ProPlate modules were removed, and whole-slides were washed 3x 5 min in PBST. Slides were then spun dry at room temperature. Slides were scanned on the LiCor Odyssey scanner at 21 μm resolution and under the "highest" quality setting in both the 700 and 800 nm channels. Spot intensities were extracted using the LiCor ImageStudio software with manual spot alignment.

### Mechanistic modeling

Mechanistic models based on ordinary differential equations (ODEs) are frequently used for the description of biochemical reaction networks. They are composed of kinetic rate equations and every component  $x$  of the model has a biological counterpart. The time evolution  $x(t)$  of the model concentrations is obtained by integration of the corresponding system of ODEs

$$\dot{x}(t, u, \theta) = f(x, \theta_d), \quad (1)$$

depending on initial and kinetic rate parameters comprised in  $\theta_d$ . These are linked to measured concentrations of the involved constituents  $y(t)$  by an observational function

$$y(t_i) = g(x(t_i, \theta_d), \theta_o) + \epsilon(t_i) \quad (2)$$

with the assumption of Gaussian errors  $\epsilon \sim N(0, \sigma)$  that is often achieved via log transformation. In addition, the observation function includes e.g.

scaling and offset parameters, summarized in  $\theta_0$ . Both observational and dynamic parameters are comprised in  $\theta$ . To compare the model response to measured data at time points  $t_i$ , the scaled log-likelihood is calculated via

$$-2 \log(\mathcal{L}) = \chi^2(\theta) = \sum_i \left( \frac{y_i - g(x(t_i, \theta))}{\sigma_i} \right)^2 + \text{const.} \quad (3)$$

Within the maximum likelihood framework, the optimized parameter set  $\hat{\theta}$  is estimated through minimization of  $\chi^2(\theta)$ .

Since analytical solutions of non-linear ODE systems are in general not available, a numerical integration must be performed. In this work, the dynamical system and its sensitivities were integrated by the CVODES integrator of the SUNDIALS suite.<sup>71</sup> Therein, an implicit BDF integration method<sup>72</sup> with attached KLU sparse solver was chosen.<sup>73</sup> The inner derivatives of the likelihood needed in gradient-based parameter estimation were computed via forward sensitivities supplied to the integration algorithm.<sup>74</sup> Numerical optimization was conducted using a trust-region based, large scale nonlinear optimization algorithm implemented in the MATLAB function LSQNONLIN.<sup>75</sup> For the mathematical

modeling and visualization, the open-source and freely available d2d framework,<sup>76</sup> based on MATLAB, was used.

#### Calculation of receptor surface levels

We established a relationship between receptor levels on the cell surface measured by qFACS (see Methods section) and receptor mRNA expression from the Cancer Cell Line Encyclopedia (CCLE) for 124 cancer cell lines.<sup>77</sup> A good correlation between mRNA and protein expression was previously shown in an independent study.<sup>78</sup> By fitting a linear model (see Suppl. Fig. 10) we could calculate receptor surface levels or receptor mRNA expression also for the cell lines where data was missing (see Table 2, Table 3, and Supplementary Fig. 41).

#### Data availability

The signaling model and code for machine learning analysis from this publication have been deposited to [BioModels.org](https://www.ebi.ac.uk/biomodels/) with the identifier MODEL1708210000. In addition, the computational model including all proteomic data, the phenotypic data from the in vitro viability screen

**Table 2.** Receptor surface levels for cell lines used in model calibration and validation estimated from mRNA expression values of the CCLE database

Cell-line	EGFR	ErbB2	ErbB3	IGF-1R	Met
<b>Cell lines used for model training</b>					
H322M	305047	69526	17984	52974	41903
BxPc-3	491421	58492	20827	40703	86574
BT-20	1365875	86856	20909	43231	43604
A431	1446255	121004	19918	87181	24296
ADRR	252245	54253	34412	117175	41684
IGROV-1	171058	93007	12897	10825	62346
ACHN	656461	49346	11883	33311	107201
-	-	-	-	-	-
<b>Cell lines used for model validation</b>					
MDA	381279	44926	9284	19446	49020
BT-474-M3	31423	1325386	28118	35556	19861

**Table 3.** Cell lines used for machine learning on an in vitro cell viability screen with their respective mutations and receptor surface levels

Cell line	KRAS mutation	PIK3CA mutation	Receptor surface levels (in thousand)			
			EGFR	HER2	ErbB3	Met
CCK-81	wt	C420R, C472Y	18.2	1325.4	28.1	19.9
GP2D	G12D	H1047L	4.2	2359.7	35.0	3.0
H508	wt	E545K	185.1	63.4	20.1	42.0
H747	G13D	Wt	193.1	2281.8	21.6	57.0
HCT116	G13D	H1047R	118.2	121.0	21.2	58.1
HCT15	G13D	E545K, D549N	290.1	382.8	23.4	40.3
HT115	wt	p.R88Q, p.E321D	381.3	44.9	9.3	49.0
KM12	wt	Wt	180.3	1852.1	29.2	15.5
LOVO	G13D	Wt	23.9	105.7	33.8	3.8



**Table 3** continued

Cell line	KRAS mutation	PIK3CA mutation	Receptor surface levels (in thousand)			
			EGFR	HER2	ErbB3	Met
LS123	G12S	Wt	23.5	121.2	31.8	11.5
LS180	G12D	H1047R	5.3	2364.7	39.5	3.2
MDST8	wt	Wt	47.7	69.7	32.9	34.1
OUMS23	wt	Wt	149.2	117.2	32.1	34.3
RCM-1	G12V	Wt	98.7	139.9	31.4	41.4
RKO	wt	H1047R	632.2	84.5	21.3	68.9
SW48	wt	G914R	174.5	46.4	16.5	67.7
SW620	G12V	Wt	97.3	88.9	23.1	29.4
T84	G13D	E542K	117.8	99.7	32.5	36.6
OVCAR-8	wt	Wt	61.0	47.8	28.7	34.1
MKN-45	wt	Wt	121.6	66.3	30.3	51.1
SNU-5	wt	Wt	213.2	55.9	22.6	57.1
H441	G12V	Wt	96.8	93.5	28.6	43.4
HCC827	wt	Wt	409.4	58.8	7.0	43.4
A549	G12S	Wt	102.7	57.3	25.4	19.8
H322M	wt	Wt	96.8	82.3	29.0	49.1
H358	G12C	Wt	21.2	44.1	11.7	18.8
ZR-75-1	wt	Wt	543.0	70.0	20.6	24.8
MDA-MB-231	G13D	Wt	4.4	34.3	26.4	36.7
BT-474	wt	K111N	149.4	55.1	29.6	63.4
HCC1419	wt	Wt	80.1	95.0	23.1	26.7
HCC1937	wt	Wt	5.2	57.0	8.5	18.0
HCC1954	wt	H1047R	182.2	33.6	6.7	395.2
HCC38	wt	Wt	188.3	78.6	20.2	76.1
JIMT-1	wt	Wt	270.8	1287.4	21.3	56.1
SK-BR-3	wt	Wt	167.7	65.6	26.4	302.2
T47D	wt	H1047R	228.5	2393.0	24.2	25.2
ZR-75-30	wt	Wt	116.5	3263.7	32.1	35.5
AGS	G12D	E453K	287.6	479.3	23.1	246.8
HGC27	wt	Wt	138.4	74.1	30.0	34.0
Hs746T	wt	Wt	374.0	69.0	18.4	297.7
KATO III	wt	Wt	366.9	42.1	11.1	73.2
KYSE-410	wt	Wt	84.0	13.8	9.5	21.4
N87	wt	Wt	20.7	2060.9	22.5	38.2
OE19	wt	Wt	662.5	34.4	21.0	87.5
OE33	wt	Wt	140.6	24.5	8.3	49.7
SNU-16	wt	Wt	305.0	69.5	18.0	41.9
H1915	wt	Wt	98.6	53.7	18.2	55.0
H2170	wt	Wt	165.4	92.7	27.7	158.3
H226	wt	Wt	71.3	36.5	7.3	30.6
H23	G12C	Wt	5.8	11.4	9.9	16.2
H460	Q61H	E545K	634.6	28.6	10.6	53.8
H520	wt	Wt	5501.4	55.1	19.3	82.3
H596	wt	E545K	276.8	40.7	21.0	58.5
CaOV3	wt	Wt	171.1	93.0	12.9	62.3
IGROV-1	wt	Wt	19.0	47.9	17.1	30.1
OV90	wt	Wt	121.3	45.5	20.1	13.4
OVCAR-3	wt	Wt	221.8	55.0	20.7	22.4
TOV-112D	Wt	Wt	12.9	64.9	9.8	3.2

together with the RNA sequencing data obtained from the CCLE<sup>77</sup> and the TCGA Research Network (<http://cancergenome.nih.gov/>), respectively, have been deposited within freely available modeling toolbox Data2Dynamics (<http://data2dynamics.org>; repository folder Examples/Hass\_npjSysBio2017). It includes MATLAB code and documented script files to readily perform all analysis steps outlined in this publication. The code folder is also available in the Figshare repository (10.6084/m9.figshare.5331544). In addition, the mechanistic signaling model was also implemented in the open-source R package dMod<sup>71</sup> (<https://github.com/dkaschek/dMod>; Figshare repository: 10.6084/m9.figshare.5336338). The package contains the experimental data for all calibration cell lines and allows to simulate model trajectories.

## ACKNOWLEDGEMENTS

We thank Tim Heinemann, Jeffrey Kearns, Sergio Iadevaia, Yasmin Hashambhoy-Ramsay, and Tim Maiwald for their constructive feedback and proof reading the manuscript, and Daniel Kaschek for implementing the model in R.

## AUTHOR'S CONTRIBUTIONS

HH build the mechanistic signaling model, built BDT and TCGA predictions, performed the computational analysis and wrote manuscript. KM performed the cell viability screen and wrote the manuscript. MS and JA performed the ELISA and lysate microarray assays. VP performed the qFACS assay. SW helped built the mechanistic signaling model. JT and JS revised the manuscript. BS and GM helped plan the study and wrote the manuscript. AR planned the study, supervised computational work and wrote the manuscript.

## ADDITIONAL INFORMATION

**Supplementary information** accompanies the paper on the *npj Systems Biology and Applications* website (<https://doi.org/10.1038/s41540-017-0030-3>).

**Competing interests:** The authors declare no competing financial interests.

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## REFERENCES

- Slamon, D. J. et al. Use of chemotherapy plus a monoclonal antibody against HER2 for metastatic breast cancer that overexpresses HER2. *N. Engl. J. Med.* **344**, 783–792 (2001).
- Howlader, N. et al. *SEER Cancer Statistics Review, 1975-2013*. (National Cancer Institute, Bethesda, MD, 2016).
- Luo, J., Solimini, N. L. & Elledge, S. J. Principles of cancer therapy: oncogene and non-oncogene addiction. *Cell* **136**, 823–837 (2009).
- Wilson, T. R., Longley, D. B. & Johnston, P. G. Chemoresistance in solid tumours. *Ann. Oncol.* **17**, 315–324 (2006).
- Zahreddine, H. & Borden, K. L. B. Mechanisms and insights into drug resistance in cancer. *Front. Pharmacol.* **4**, 28 (2013).
- Ledford, H. Ways to fix the clinical trial. *Macmillan Publ. Ltd. Nat.* **477**, 526–528 (2011).
- Nelson, M. R. et al. The genetics of drug efficacy: opportunities and challenges. *Nat. Rev. 1Genet.* **17**, 197–206 (2016).
- Paez, J. G. et al. EGFR mutations in lung cancer: correlation with clinical response to gefitinib therapy. *Science* **304**, 1497–1500 (2004).
- Arteaga, C. L. Epidermal growth factor receptor dependence in human tumors: more than just expression? *Oncologist* **7**, 31–39 (2002).
- Liu, F., Wang, L., Perna, F. & Nimer, S. D. Beyond transcription factors: how oncogenic signalling reshapes the epigenetic landscape. *Nat. Rev. Cancer* **16**, 359–372 (2016).
- Tateishi, M., Ishida, T., Mitsudomi, T., Kaneko, S. & Sugimachi, K. Immunohistochemical evidence of autocrine growth factors in adenocarcinoma of the human lung. *Cancer Res.* **50**, 7077–7080 (1990).
- Umekita, Y., Ohi, Y., Sagara, Y. & Yoshida, H. Co-expression of epidermal growth factor receptor and transforming growth factor- $\alpha$  predicts worse prognosis in breast-cancer patients. *Int. J. Cancer* **89**, 484–487 (2000).
- Arteaga, C. L. & Engelman, J. A. ERBB receptors: from oncogene discovery to basic science to mechanism-based cancer therapeutics. *Cancer Cell.* **25**, 282–303 (2014).
- Chong, C. R. & Jänne, P. The quest to overcome resistance to EGFR-targeted therapies in cancer. *Nat. Med.* **19**, 1389–1400 (2013).
- Holohan, C., Van Schaeybroeck, S., Longley, D. B. & Johnston, P. G. Cancer drug resistance: an evolving paradigm. *Nat. Rev. Cancer* **13**, 714–726 (2013).
- Schoeberl, B. et al. Therapeutically targeting ErbB3: a key node in ligand-induced activation of the ErbB receptor-P13K axis. *Sci. Signal.* **2**, ra31 (2009).
- Yarden, Y. & Pines, G. The ERBB network: at last, cancer therapy meets systems biology. *Nat. Rev. Cancer* **12**, 553–563 (2012).
- Altman, R. B. Predicting cancer drug response: advancing the dream. *Cancer Discov.* **5**, 237–238 (2015).
- Hill, S. M. et al. Inferring causal molecular networks: empirical assessment through a community-based effort. *Nat. Method.* **13**, 310–318 (2016).
- Menden, M. P. et al. Machine learning prediction of cancer cell sensitivity to drugs based on genomic and chemical properties. *PLoS. ONE.* **8**, e61318 (2013).
- Radivojac, P. et al. A large-scale evaluation of computational protein function prediction. *Nat. Method.* **10**, 221–227 (2013).
- Fey, D. et al. Signaling pathway models as biomarkers: patient-specific simulations of JNK activity predict the survival of neuroblastoma patients. *Sci. Signal.* **8**, ra130 (2015).
- Kearns, J. D. et al. Enhanced targeting of the EGFR network with MM-151, an oligoclonal anti-EGFR antibody therapeutic. *Mol. Cancer Ther.* **14**, 1625–1636 (2015).
- Abu-Yousif, A. O. et al. Mechanistic characterization of MM-131, a bispecific antibody that blocks c-Met signaling through concurrent targeting of EpCAM. *Cancer Res.* **75**, 1690 (2015).
- Fitzgerald, J. B. et al. MM-141, an IGF-1R- and ErbB3-directed bispecific antibody, overcomes network adaptations that limit activity of IGF-1R inhibitors. *Mol. Cancer Ther.* **13**, 410–425 (2014).
- Luey, B. C. & May, F. E. B. Insulin-like growth factors are essential to prevent anoikis in oestrogen-responsive breast cancer cells: importance of the type I IGF receptor and PI3-kinase/Akt pathway. *Mol. Cancer* **15**, 8 (2016).
- Mendoza, M. C., Er, E. E. & Blenis, J. The Ras-ERK and PI3K-mTOR pathways: cross-talk and compensation. *Trends Biochem. Sci.* **36**, 320–328 (2011).
- Niepel, M. et al. Profiles of basal and stimulated receptor signaling networks predict drug response in breast cancer lines. *Sci. Signal.* **6**, ra84 (2013).
- Endo, H., Okuyama, H., Ohue, M. & Inoue, M. Dormancy of cancer cells with suppression of AKT activity contributes to survival in chronic hypoxia. *PLoS. One.* **9**, e98858 (2014).
- Engelman, J. A. et al. MET amplification leads to gefitinib resistance in lung cancer by activating ERBB3 signaling. *Sci. (80-.)* **316**, 1039–1043 (2007).
- Jin, Q. & Esteva, F. J. Cross-talk between the ErbB/HER family and the type I insulin-like growth factor receptor signaling pathway in breast cancer. *J. Mammary Gland. Biol. Neoplasia.* **13**, 485–498 (2008).
- Lai, A. Z., Abella, J. V. & Park, M. Crosstalk in Met receptor oncogenesis. *Trends Cell. Biol.* **19**, 542–551 (2009).
- Oda, K., Matsuoka, Y., Funahashi, A. & Kitano, H. A comprehensive pathway map of epidermal growth factor receptor signaling. *Mol. Syst. Biol.* **1**, E1–E17 (2005).
- Avraham, R. & Yarden, Y. Feedback regulation of EGFR signalling: decision making by early and delayed loops. *Nat. Rev. Mol. Cell. Biol.* **12**, 104–117 (2011).
- Waterman, H. & Yarden, Y. Molecular mechanisms underlying endocytosis and sorting of ErbB receptor tyrosine kinases. *FEBS Lett.* **490**, 142–152 (2001).
- Citri, A. & Yarden, Y. EGF-ERBB signalling: towards the systems level. *Nat. Rev. Mol. Cell. Biol.* **7**, 505–516 (2006).
- Desbois-Mouthon, C. et al. Insulin-like growth factor-1 receptor inhibition induces a resistance mechanism via the epidermal growth factor receptor/HER3/AKT signaling pathway: rational basis for cotargeting insulin-like growth factor-1 receptor and epidermal growth factor receptor. *Clin. Cancer Res.* **15**, 5445–5456 (2009).
- Yarar, D., Lahdenranta, J., Kubasek, W., Nielsen, U. B. & MacBeath, G. Heregulin-ErbB3-driven tumor growth persists in PI3 kinase mutant cancer cells. *Mol. Cancer Ther.* **14**, 2072–2080 (2015).
- Sevecka, M., Wolf-Yadlin, A. & MacBeath, G. Lysate microarrays enable high-throughput, quantitative investigations of cellular signaling. *Mol. Cell. Proteom.* **10**, M110.005363 (2011).
- Kirouac, D. C. et al. Computational modeling of ERBB2-amplified breast cancer identifies combined ErbB2/3 blockade as superior to the combination of MEK and AKT inhibitors. *Sci. Signal.* **6**, ra68 (2013).
- Rosenblatt, M., Timmer, J. & Kaschek, D. Customized steady-state constraints for parameter estimation in non-linear ordinary differential equation models. *Front. Cell Dev. Biol.* **4**, 41 (2016).
- Shi, T. et al. Conservation of protein abundance patterns reveals the regulatory architecture of the EGFR-MAPK pathway. *Sci. Signal.* **9**, rs6 (2016).
- Macdonald-Obermann, J. L. & Pike, L. J. Different epidermal growth factor (EGF) receptor ligands show distinct kinetics and biased or partial agonism for homodimer and heterodimer formation. *J. Biol. Chem.* **289**, 26178–26188 (2014).
- Yarden, Y. The EGFR family and its ligands in human cancer: signalling mechanisms and therapeutic opportunities. *Eur. J. Cancer* **37**, 3–8 (2001).

45. Breiman, L. Random forests. *Mach. Learn.* **45**, 5–32 (2001).
46. Schapire, R. E. The strength of weak learnability. *Mach. Learn.* **5**, 197–227 (1990).
47. Kingsford, C. & Salzberg, S. L. What are decision trees? *Nat. Biotechnol.* **26**, 1011–1013 (2008).
48. Rokach, L. & Maimon, O. *Data mining with decision trees: theory and applications*. (World scientific, 2014).
49. Annovazzi, L. et al. mTOR, S6 and AKT expression in relation to proliferation and apoptosis/autophagy in glioma. *Anticancer. Res.* **29**, 3087–3094 (2009).
50. Gazdar, A. F., Shigematsu, H., Herz, J. & Minna, J. D. Mutations and addiction to EGFR: the achilles ‘heal’ of lung cancers? *Trends Mol. Med.* **10**, 481–486 (2004).
51. Moroni, M. et al. Gene copy number for epidermal growth factor receptor (EGFR) and clinical response to antiEGFR treatment in colorectal cancer: a cohort study. *Lancet Oncol.* **6**, 279–286 (2005).
52. Sharma, S. V., Bell, D. W., Settleman, J. & Haber, D. A. Epidermal growth factor receptor mutations in lung cancer. *Nat. Rev. Cancer* **7**, 169–181 (2007).
53. Laurent-Puig, P. et al. Analysis of PTEN, BRAF, and EGFR status in determining benefit from cetuximab therapy in wild-type KRAS metastatic colon cancer. *J. Clin. Oncol.* **27**, 5924–5930 (2009).
54. Kris, M. et al. Efficacy of gefitinib, an inhibitor of the epidermal growth factor receptor tyrosine kinase, in symptomatic patients with non-small cell lung cancer: a randomized trial. *JAMA* **290**, 2149–2158 (2003).
55. Pérol, M. Negative results of METLung study: an opportunity to better understand the role of MET pathway in advanced NSCLC. *Transl. Lung Cancer Res.* **3**, 392–394 (2014).
56. Masuda, H. & Zhang, D. Role of epidermal growth factor receptor in breast cancer. *Breast Cancer Res.* **136**, 1–21 (2012).
57. Schoeberl, B. et al. Systems biology driving drug development: from design to the clinical testing of the anti-ErbB3 antibody seribantumab (MM-121). *npj Syst. Biol. Appl.* **3**, 16034 (2017).
58. Ryerson, A. B. et al. Annual report to the nation on the status of cancer, 1975–2012, featuring the increasing incidence of liver cancer. *Cancer* **122**, 1312–1337 (2016).
59. FDA. *Accelerating the Development of New Pharmaceutical Therapies*. (Silver Spring, MD, 2015).
60. Mullard, A. 2014 FDA drug approvals. *Nat. Publ. Gr.* **14**, 77–81 (2015).
61. ATLAS Collaboration. Evidence for the Higgs-boson Yukawa coupling to tau leptons with the ATLAS detector. *J. High. Energy Phys.* **2015**, 117 (2015).
62. Dieterich, T. G. Experimental comparison of three methods for constructing ensembles of decision trees: bagging, boosting, and randomization. *Mach. Learn.* **40**, 139–157 (2000).
63. Lessmann, S., Baesens, B., Seow, H. V. & Thomas, L. C. Benchmarking state-of-the-art classification algorithms for credit scoring: an update of research. *Eur. J. Oper. Res.* **247**, 124–136 (2015).
64. Maiwald, T. et al. Driving the model to its limit: profile likelihood based model reduction. *PLoS. ONE.* **11**, e0162366 (2016).
65. Patel, A. P. et al. Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science* **344**, 1396–1401 (2014).
66. Bansal, M. et al. A community computational challenge to predict the activity of pairs of compounds. *Nat. Biotech.* **32**, 1213–1222 (2014).
67. Costello, J. C. et al. A community effort to assess and improve drug sensitivity prediction algorithms. *Nat. Biotechnol.* **32**, 1–103 (2014).
68. Sevecka, M. & MacBeath, G. State-based discovery: a multidimensional screen for small-molecule modulators of EGF signaling. *Nat. Methods* **3**, 825–831 (2006).
69. Schoeberl, B. et al. A Data-Driven Computational Model of the ErbB Receptor Signaling Network. in *2006 International Conference of the IEEE Engineering in Medicine and Biology Society* 53–54 (IEEE, 2006). <https://doi.org/10.1109/IEMBS.2006.259754>.
70. Kaushansky, A. et al. Quantifying protein-protein interactions in high throughput using protein domain microarrays. *Nat. Protoc.* **5**, 773–790 (2010).
71. Hindmarsh, A. C. et al. SUNDIALS: suite of nonlinear and differential/algebraic equation solvers. *ACM Trans. Math. Softw.* **31**, 363–396 (2005).
72. Gear, C. Simultaneous Numerical Solution of Differential-Algebraic Equations. *IEEE Trans. Circuit Theory* **18**, 89–95 (1971).
73. Davis, T. A. & Natarajan, E. P. Algorithm 907. *ACM Trans. Math. Softw.* **37**, 1–17 (2010).
74. Leis, J. R. & Kramer, M. A. The simultaneous solution and sensitivity analysis of systems described by ordinary differential equations. *ACM Trans. Math. Softw.* **14**, 45–60 (1988).
75. Coleman, T. F. & Li, Y. An interior, trust region approach for nonlinear minimization subject to bounds. *SIAM J. Optim.* **6**, 418–445 (1996).
76. Raue, A. et al. Data2Dynamics: a modeling environment tailored to parameter estimation in dynamical systems. *Bioinformatics* **31**, 3558–3560 (2015).
77. Barretina, J. et al. The cancer cell line encyclopedia enables predictive modelling of anticancer drug sensitivity. *Nature* **483**, 603–607 (2012).
78. Edfors, F. et al. Gene-specific correlation of RNA and protein levels in human cells and tissues. *Mol. Syst. Biol.* **12**, 883 (2016).



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2017