

Supplements

Quality assessment of pharmacology data

The Quality Assessment (QA) flag is an output from Combenefit, which was used for the synergy score and monotherapy curve fitting (see Online Methods). Scores are defined as the following:

Flag	Meaning
0	No data was found in a combination folder or NaN was found in a combination file.
-1	At least one of the measured drug effects was above 125% of starting cell count. This is unlikely to be genuine and a major experimental issue is suspected.
-2	No flag '-1' but measured effects below -10% were found. By definition, effects should always be positive because cell viability is being measured. Very small negative values are sometime encountered due to quantification problems at high concentrations. These were tolerated up to -10% below which major issues were suspected.
-3	No flags '-1' or '-2' but combination dose-response showed very strong fluctuations. The combination dose-response was smoothed and compared to the original non-smoothed version. If differences above 25% were found the experiment was flagged as measurements likely to be unreliable.
1	None of the previous problems were encountered. Data is supposed to be ok.

All experiments were observed to have some level of synergy or antagonism and had non-zero synergy scores. Most of these were due to random variation in the experiments and had synergy within +/-1, and only 404 experiments with low variability non-zero experiments due to quality issues in the assay and were flagged accordingly (Fig. S1B). Only high quality data (QA=1) were included in the testset, while experiments also with low quality were made available for training.

Drug combination synergy is variable across cells, but reproducible across replicates

367 had a replicate experiment where the same drug combination, cell line and concentration ranges tested (Fig. S2A). The Spearman correlation between synergy scores across the replicates was 0.56, which is comparable to the correlation of 0.63 from the 315 replicate combinations screening experiments (Fig. S2B) completed by O'Neil *et al*¹. We observed that in all instances where the synergy and antagonism were not measured by both replicates, the quality of one or both replicates had been flagged as low (see above QA flags). Notably, the variance in synergy scores in DREAM (Fig. S2C) was larger than the variance the dataset from O'Neil et al (Fig. S2D).

Best ranked teams' Methods:

Best ranked teams' methods are below, who performed consistently well across all 3 sub-challenges, or outstanding well for at least one sub-challenge.

Yuanfang Guan

This method begins by limiting the feature space to those features (expression, CNV, methylation, mutations) mapping to genes that are putative targets of any drug in a given sub-challenge. If a drug is less specific and has multiple targets, all of them are included.

For each drug combination a separate classifier is built for predicting synergy. Each drug combination classifier is made of several random forests, and each using as distinct data type. For example, in sc1A three classifiers are created per drug combination, and with their predictions are averaged (table 1). This use of one classifier for each data type was motivated to improve the stability of the predictions, in the case that one of the feature sets contains outliers.

Sc1A Classifier	Data type
1	Mono-therapy data of Drug A and B
2	Drug A, B count data
3	All DC's data with either Drug A or B
Sc1B Classifier	Data type
1	CNV, mutation (and mean values for both), count on drug A,
2	CNV, mutation, count (and mean) on drug B
3.	Normalized CNV, mutation (and mean), count on drug A and B
Sc2 Classifier	Data type
1	Drug A, B, count,
2	Normalized All DC's data with either Drug A and B
3	Drug space as a simple parsing of original table provided in the Challenge
4	Mono-therapy data of Drug A and B

Table 1

The choice of classifier construct, randomForest, compared to using an SVM, regression, or boosting, was chosen for expedience, not any perceived advantage in accuracy, though the randomForest does facilitate prediction in a nonlinear context². Instead, greater prediction accuracy is achieved by creating new features that consist of scaling the features provided by AZ-Dream with a posterior probability from a predefined functional network¹.

Network Based Feature Scaling:

Feature scaling is motivated by the observation that when we use the original features as input, we found the prediction values are similar for the same cell line, regardless of drug perturbations. This is due to the fact the genomic and expression data are static for a given cell line across all drug combinations effectively leaving just three parameters for modeling drug synergy (drug A, drug B and cell line). Scaling cell line features by the functional network of a drug target creates a dynamic parameter space for each cell line that allows for better modeling synergy.

For gene expression, methylation, and copy number variation data, features are adjusted based on their probability of a functional relationship with drug target genes. Feature x_i , associated with gene i , is scaled to x'_i by:

$\{if\ g_i \in (DT_A, DT_B), x'_i = 0\ else, x'_i = x_i \times (1 - \max(e_{ij}), \forall\ g_j \in (DT_A, DT_B)\}$

where DT_A and DT_B are the genes targeted by drug A and B respectively, and e_{ij} is edge between genes i and j in the predefined functional relationship network¹.

For gene mutations, features are modified using the edge directly. Mutation feature x_i , associated with gene i , is changed to x'_i :

$\{if\ g_i \in (DT_A, DT_B), x'_i = x_i\ else, x'_i = x_i \times \max(e_{ij}), \forall\ g_j \in (DT_A, DT_B)\}$

The purpose of generating this new series of features is to simulate the effective values of expression, methylation, CNV and mutations post treatment. After scaling, predictive features are different for each cell line across different drug-combinations. We reduce the effective values of drug target in expression, methylation and CNV to zero, and reduce the values of other genes according to their connections to the drug target. For mutations, we assumed that the effects of drugs are equivalent to adding in new mutations to the system, with the effect values of drug targets being 1 (similar to mutated genes), while the other genes are increased in values according to their connections to the drug target.

A biological interpretation of this approach is that scaling cell line features with the functional network's edges, allows the RandomForest to model a drug's propagation through a targeted pathway cascade.

The weighting of different set of features was primarily done by cross-validation. However, we found that a single set of genomic features is often sufficient to achieve a similar performance as the entire set.

Mikhail Zaslavskiy

The model is an ensemble of three individual models trained exclusively on categorical features describing drug and cell line identities and one model trained on categorical identity features plus corresponding drug MonoTherapy results. There are two main ideas at the core of the proposed model. First, it is very easy to overfit when dealing with biological data, so the key factor is a proper design of the cross-validation scheme which covers not only meta parameter estimation but also model selection steps. In addition to avoiding overfitting pitfalls, the cross-validation design is important to address precisely the sub-challenge questions defined by corresponding training/test splits. Second, a rich sampling of the experimental space provides an excellent support for a competitive model even without additional features describing biological entities under consideration. When we have enough data on drug/drug and drug/cell line combinations (like in sc1A and sc1B) we can derive the information on drug and cell line similarities without using additional features. Of course, it is impossible to know in advance if the experimental results alone are enough to reach the maximum performance, so the model building process is to start with the set of baseline features (drug and cell line ids) and then add step-by-step more complex features (MonoTherapy results, drug features, cell line mutations, copy number variations et c.) verifying at each step if the addition of new features lead to an improvement of the cross-validation score.

The three individual models used to predict drug synergy from drug/cell line identities are a gradient boosting tree model (xgboost package) and an svm model (e1071 package) trained on original identity features represented by a binary matrix, and an elastic net model (R/glmnet package) model trained on average scores of drug-cell line combinations with the

same drug combination and a different cell line, average scores of drug-cell line combinations with the same cell line and one drug in common, average scores of drug-cell line combinations with the same cell line and no common drug. The fourth model is another gradient boosting tree model trained on drug MonoTherapy results and counts of categorical identity features. All four models were trained using 5-fold cross-validation, the final score was computed as a simple average of the individual models.

North Atlantic DREAM (NAD)

North Atlantic Dream team's solutions used different tree based models (Random Forest Regression and Extreme Gradient Boosting Trees, XGBoost³) to incorporate the presumed important interactions between cellular (mutations, copy number alterations etc.) and drug specific (drug targets, affected pathways etc.) features. For better representation of the similarities between cell lines and drug combinations, new sets of features were engineered using prior knowledge and also the monotherapy data.

The monotherapy data (IC50, Einf, Hill slope) for a drug combination was dependent on the ordering of the drugs in the combination, which seemed to be hard to represent in the machine learning model. To overcome this problem, North Atlantic Dream's model used monotherapy features that are independent of the ordering (such as min/max/absolute difference etc. type features from the original data). Also the expected volume under the dose-response surface (in case of additivity) was calculated using the original Loewe model⁴.

As the number of training examples for a given drug combination was relatively low (about a dozen for most combinations), it was crucial to find similarities between combinations beyond the trivial ones (same drug / same target). North Atlantic Dream created different feature sets based on Gene Ontology⁵ / KEGG Pathways⁶ and a directed signaling network⁷. For Gene Ontology based features a set of "cancer related" GO terms were selected (based on⁸), and for each drug a GO vector was created based on the association of GO terms with the target of the drug. For KEGG Pathway based features, KEGG Pathways containing the target genes were selected, and for each drug a KEGG vector was created, giving 1 values for pathways containing the target of drug, 0 otherwise. The GO/KEGG features for drug combinations were the sum of the two drug vectors of the combination. Based on the directed signaling network, for each drug combination the "similar" drug combinations were selected. Two drug combinations were defined similar, if the two targets of combination A are direct upward from the two targets of combination B. Based on this rule a similarity vector was created for each combination. To create these GO/KEGG/signaling networks based features, in case of "DNA targeting drugs" (i.e. chemotherapy drugs), the respective DNA damage response molecule was used as indirect target (based on Woods & Turchi, 2013⁹).

For cellular features mutations, copy number variations and gene expression were used. The main problem with cellular features was their large number. To overcome this problem, North Atlantic Dream team used a pre-assembled gene list (including target genes, known oncogenes and tumor suppressors, genes related to drug monotherapy resistance etc.^{10 11 12}). Genes with mutations / copy number alterations were selected from this gene list. The low number of training examples for a given drug combination made it hard to use traditional feature selection/reduction methods. However, based on the drug similarities defined above,

it was possible to select molecular features associated with the observed synergy scores for a given, similar set of drug combinations. During the original Challenge gene expression features were selected (from the expression of target genes and their direct neighbors in signaling network) using this method by Randomized Lasso Feature selection. In the later, collaborative phase of the Challenge, a similar method was used for mutation and copy number variation features.

For the final prediction different XGBoost models were created using subsets of the above defined features and/or different model parameters. The final submitted predictions were the ensembles of these models, either as simple averages or using hillclimbing¹³ on out-of-fold predictions. For the various tasks through the Challenge R, Python, SAS, and JMP were used.

NAD feature layer importance:

NAD's Random Forest Regression model was trained using different pairs of cell line and drug combination specific features. The tested cell line features included cell line label, mutations (pre-filtered for 469 cancer related gene) and CNV (pre-filtered for 292 gene). Combination related features tested here were drug label, drug target, Gene Ontology and KEGG pathway based features (feature size: 407 and 140, respectively) and signalling network based features (601). Baseline model used cell line and drug label as features, while in the other models the respective feature was either swapped with the corresponding baseline feature (e.g.: in *CNV model* CNV and drug label, in *target model* cell line label and drug target was used as features) or added to the features of baseline model (e.g. in *+target model* cell line label, drug label and drug target was used). Ensemble model in this case is the simple average of the prediction of these models. With all the used models 10 random cross-validation was performed, and the mean weighted Pearson correlation was calculated for each cross-validation run. For the cross-validations all the training and leaderboard data of the Challenge was used, and the size of the cross-validation set for each combination resembled the size of the test set of the Challenge.

NAD biomarker selection:

NAD ranked their biomarkers for each drug combination based on the number of times the feature was used for predicting the given combination in the Random Forest Regressor models. For each combination a separate Random Forest model was built (using mutation, CNV, drug target and KEGG features) where all the Challenge training and leaderboard data was used without the data of the actual combination. With this model the left out combination was predicted for all of the cell lines. The number of times a given cellular feature (mutation or CNV) was used (based on the internal structure of the Random Forest trees) was recorded for each combination - feature pair. For each drug combination - feature pair Mann-Whitney U test was performed to calculate the probability that the given feature is used more often for the given combination than other features, and that the given feature is used more often for the given combination than for other combinations. The final score of the feature for the combination was the product of these two probabilities. For each combination the used cellular features (mutation and CNV) was sorted decreasing order, and the top 5 features was used for further analysis.

DMIS

Support Vector Regression (SVR) were used as prediction model. The main difficulty was the high dimensionality problem of our feature space. To address this, novel literature-based approach was used. 200 genes were identified that most frequently occur in the context of cancer in the literature, and used only the mutations in those 200 genes as the features. To perform this gene selection task, the Biomedical Entity Search Tool (BEST) (<http://best.korea.ac.kr>)¹⁴ was applied. BEST finds an entity relevant to a query based on the number of co-occurrences between the query terms and the entity in the PubMed corpus, the authority of journals, the recency of articles, and the term frequency inverse document frequency (TF-IDF) weighting. BEST was queried by using the query term “cancer” and the top 200 cancer-related genes were collected. For CNV features, cBioPortal was used to collect 13 gene sets of cancer-related pathways. During the creation of various types of features, it was of essence to identify the best combinations of feature groups was a necessary part. However, testing all possible combinations of feature groups would require a considerable amount of computing resources. To address this problem, a high performance computing pipeline using HTCondor was constructed. HTCondor is an open-source computing framework for coarse-grained distributed parallelization of computationally intensive tasks. We ran our pipeline using 1,764 cores from Amazon Web Service, and selected the best combination of the feature groups. Through this process, the following features for sub-challenge 1B were selected: 118 drug IDs, 99 drug targets, 94 CNVs, 241 mutations, and maximum concentrations of the dosages for each sample.

The SVR model showed a good performance on AZ dataset. However, because the original model represents target and mutation features as sparse binary vectors, it is not appropriate to apply to O'Neil et al dataset with unseen cell lines and drugs. For translatability, a dense vector was created to capture and generalize characteristics of cell lines and drugs, and constructed a deep learning model which could utilize these vectors as input.

For the post-hoc analysis, an additional deep learning model was generated, which was composed of 6 layers including a preprocessing layer as the first layer. The second layer had 4 modules and the first module gets mutation features, and the second module gets target feature. These two modules embed sparse feature vectors as dense vectors, and generate a single vector using a convolutional neural network. Pre-trained mutation vectors were used to leverage mutation information from TCGA and Mikolov's Word2Vec algorithm¹⁵ for mutation embedding. In addition, Asgari's public protein dense vector¹⁶ for target embedding was used. The third module gets drug or monotherapy-related features, and the last module gets cell line-related features. Each module generates a single vector and the vectors are concatenated and are inputted into the next layer. The rest of layers are fully-connected layer. The output layer generates a single value, which is the predicted synergy score.

DMIS feature layer importance:

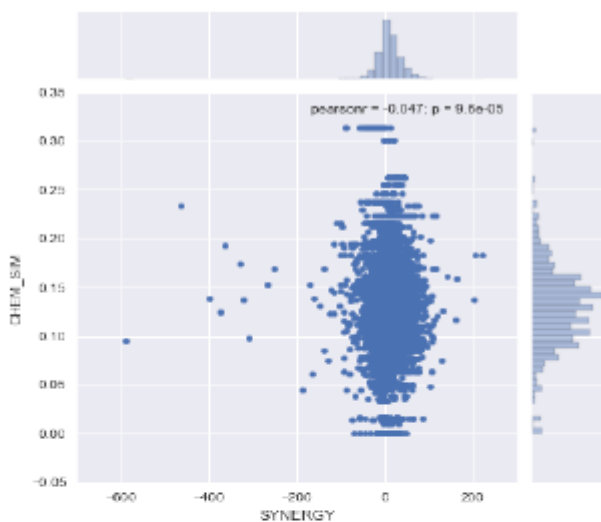
In the main Challenge, the Support Vector Regression (SVR) was used as machine learning model. For extracting feature layer importance, the accuracy of the model was estimated after randomly permuting the values of the feature. How much the permutation decreases the primary score of the SVR model is an estimate of feature layer importance.

DMIS biomarker selection:

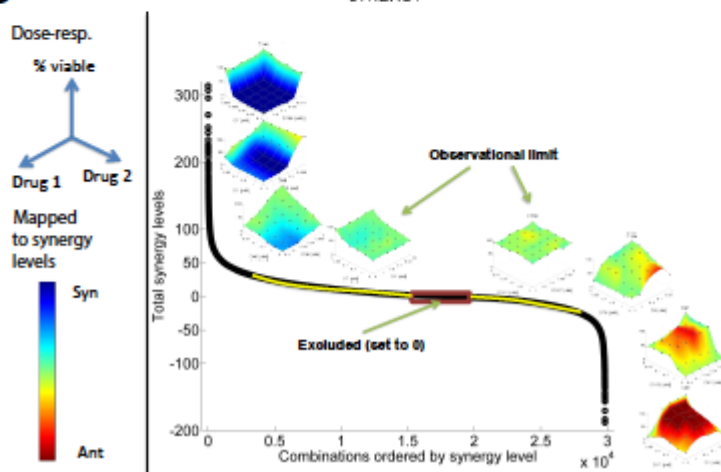
In order to get biomarker indications, important features were extracted for each drug combination individually. Therefore, samples were grouped by drug combination and for each combination a random forest model was built. The method called “mean decrease impurity” was applied to obtain feature rankings. Random forest consists of multiple decision trees. Each node in a tree corresponds to a feature and the same feature can appear in multiple trees. Each node in a tree splits the samples so that similar synergy scores group together. At each node, we can compute how much the split reduces the variance in the sample by taking the difference between the variance-before-split and the variance-after-split (this can be measured by weighted average of the two split groups). Finally, features were ranked based on their mean variance reduction, i.e., mean decrease impurity. The number of trees was limited to 200 and each tree randomly selected up to 204 features (a third of total 612 features available).

Supplemental figures

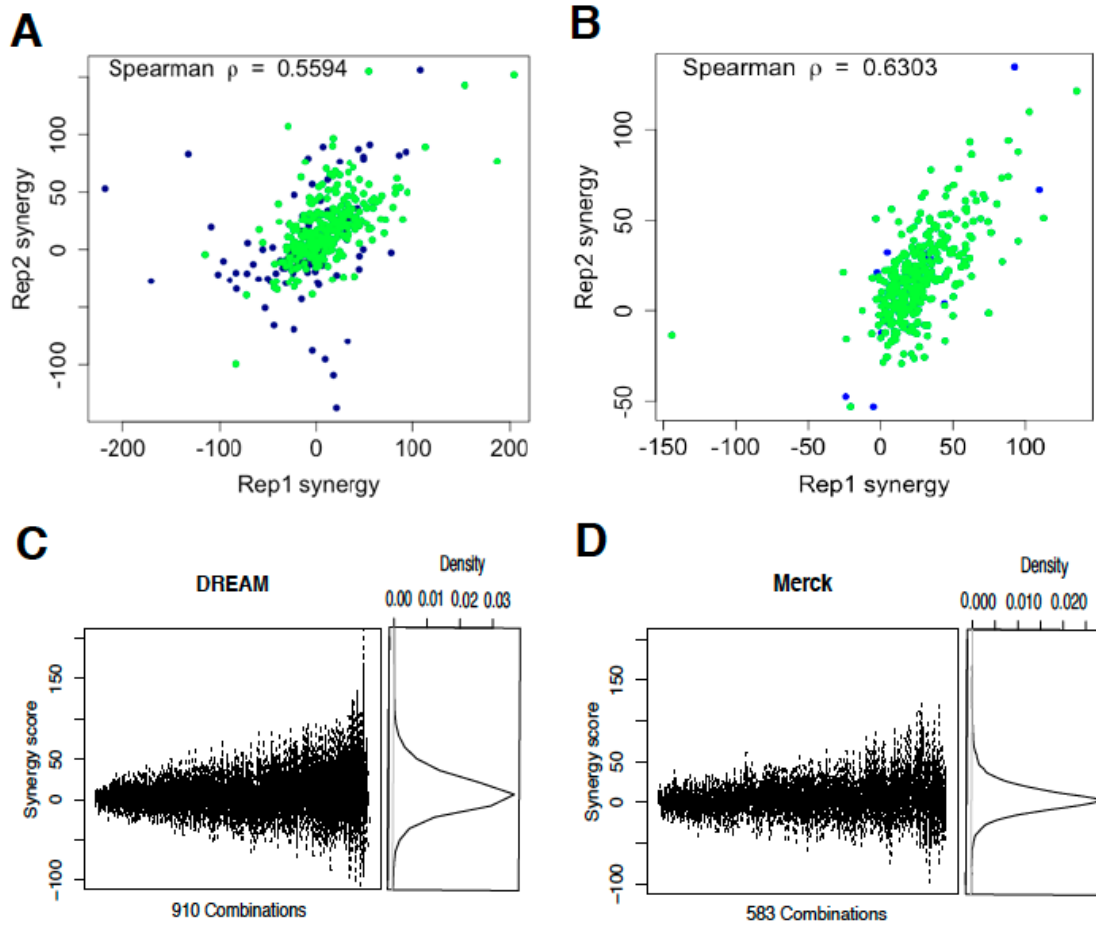
A



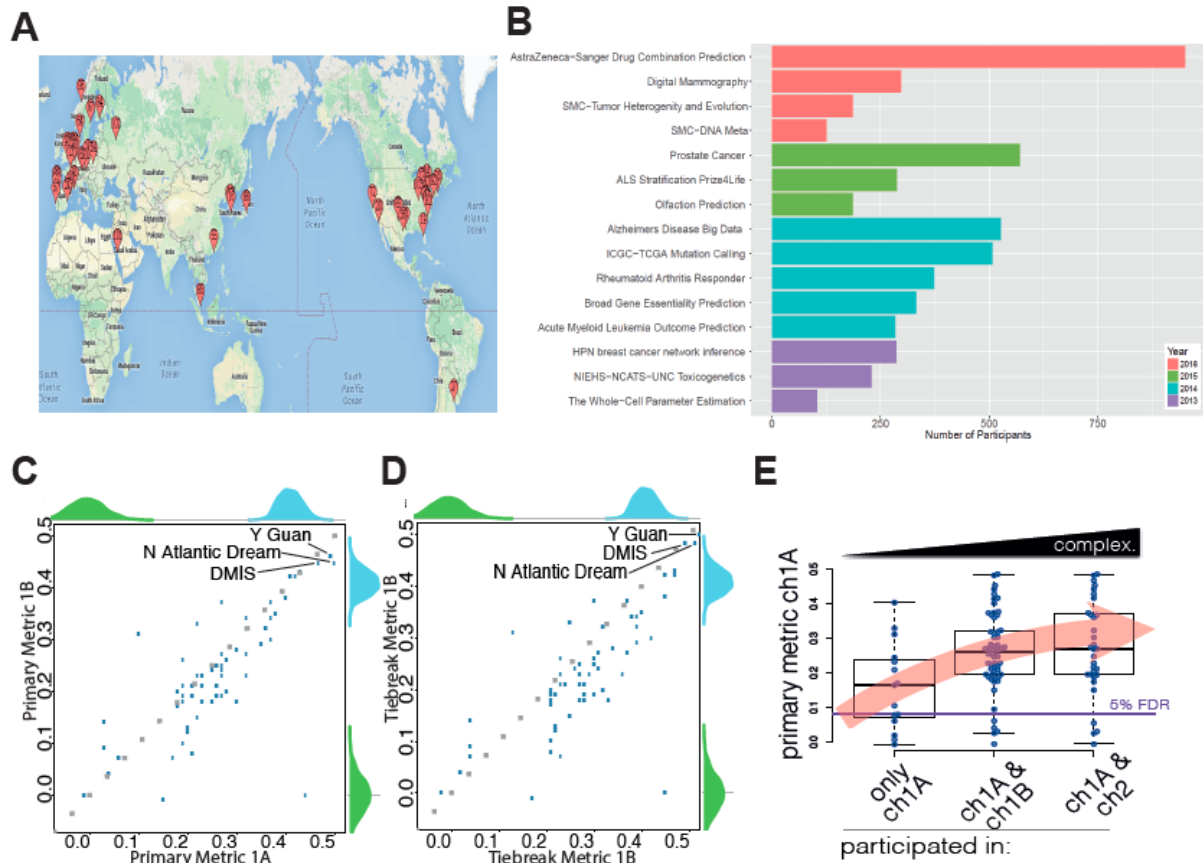
B



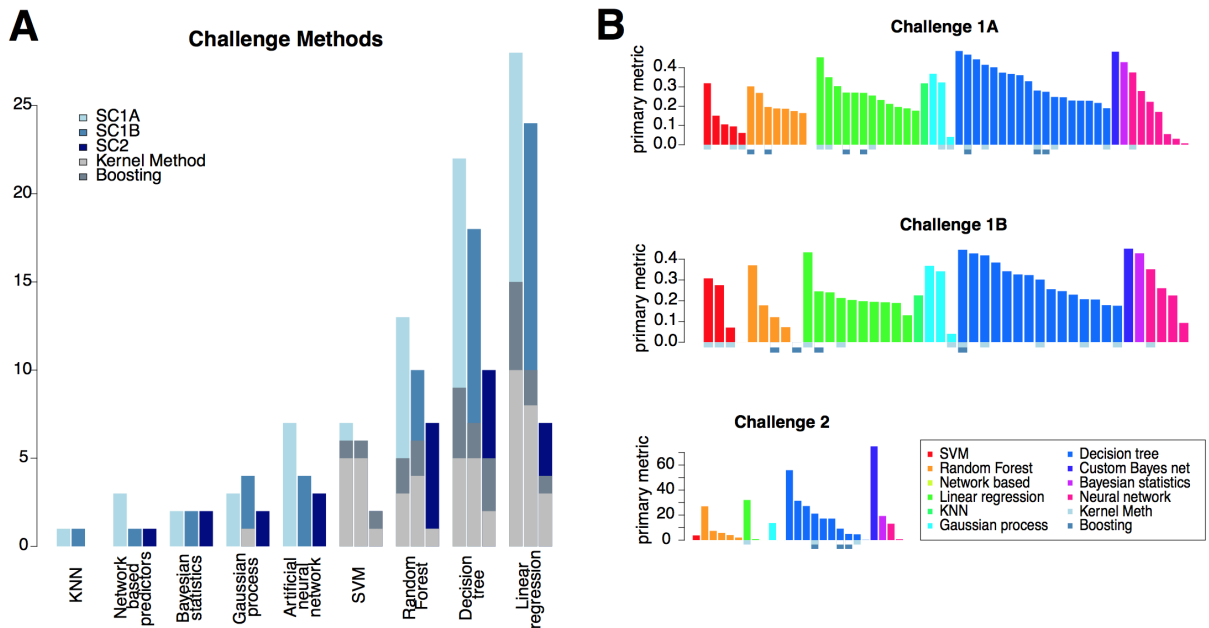
Supp. Figure 1: Distribution of synergy scores across all combinations. (A) Type of compounds in each combination is plotted against their synergy scores. (B) Chemical similar synergy scores of combinations are ordered from lowest to highest. 3D synergy heatmaps show additional cells killed (Syn) or not killed (Ant) beyond the additive effect of the two drugs at each dose. Two examples of combinations with total synergy scores of +/-20 show the limit at which synergy and antagonism can be visually confirmed. Experiments where non-zero total synergy was due to random variation were set to zero.



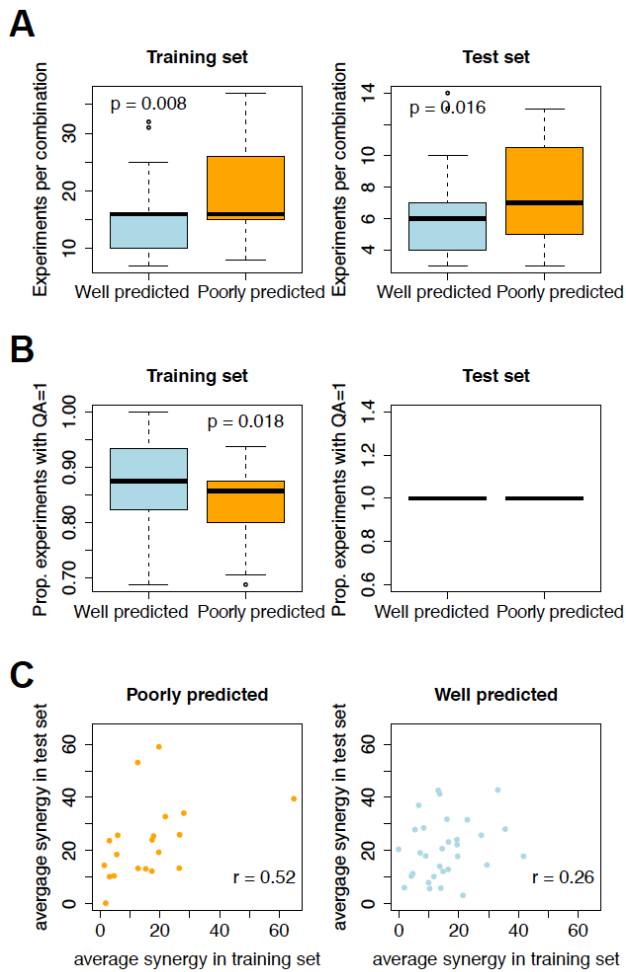
Supp. Figure 2: Reproducible of synergy scores. Shows correlation of drug combinations replicates from (A) DREAM Challenge and (B) external combination screen (O’Neil et al. 2016). In green are high quality data points, while in blue are points not passing the QC from CombeneFit (QC score=1). (C, D) Distribution of synergy scores across combinations screened for the DREAM Challenge and independent set of combinations screened by Merck.



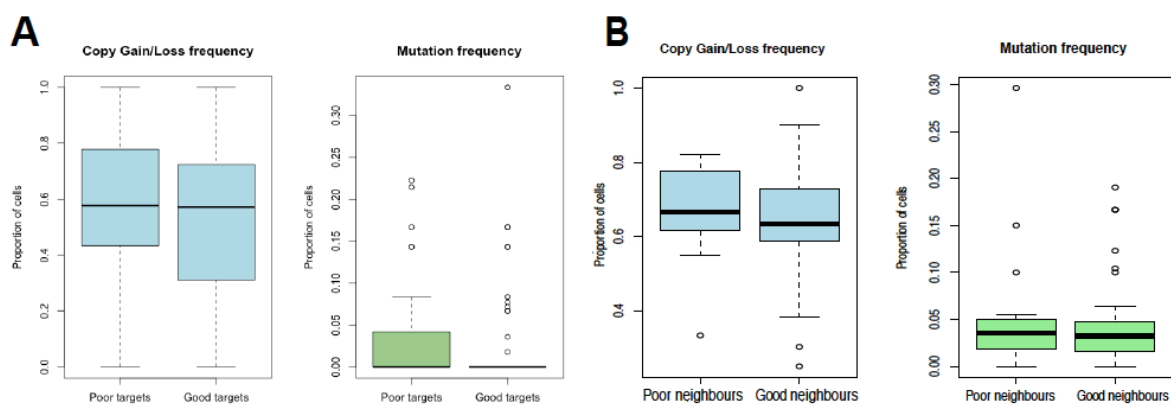
Supp. Figure 3: Participation in the Combinations Prediction DREAM Challenge. (A) Nearly 800 participants were located across five continents. (B) Comparison of participants across different DREAM Challenges. (C) Performances of sub-challenge 1A plotted against 1B based on the primary metric, average weighted Pearson correlation. (D) Performances of sub-challenge 1A plotted against 1B based on the tie-break metric, average weighted Pearson correlation of combinations with cases of synergy > 20. Dotted grey line shows 1:1 relationship between 1A and 1B performance. (E) Performance of participants in sub-challenge 1A grouped by whether they participated in one or more sub-challenges.



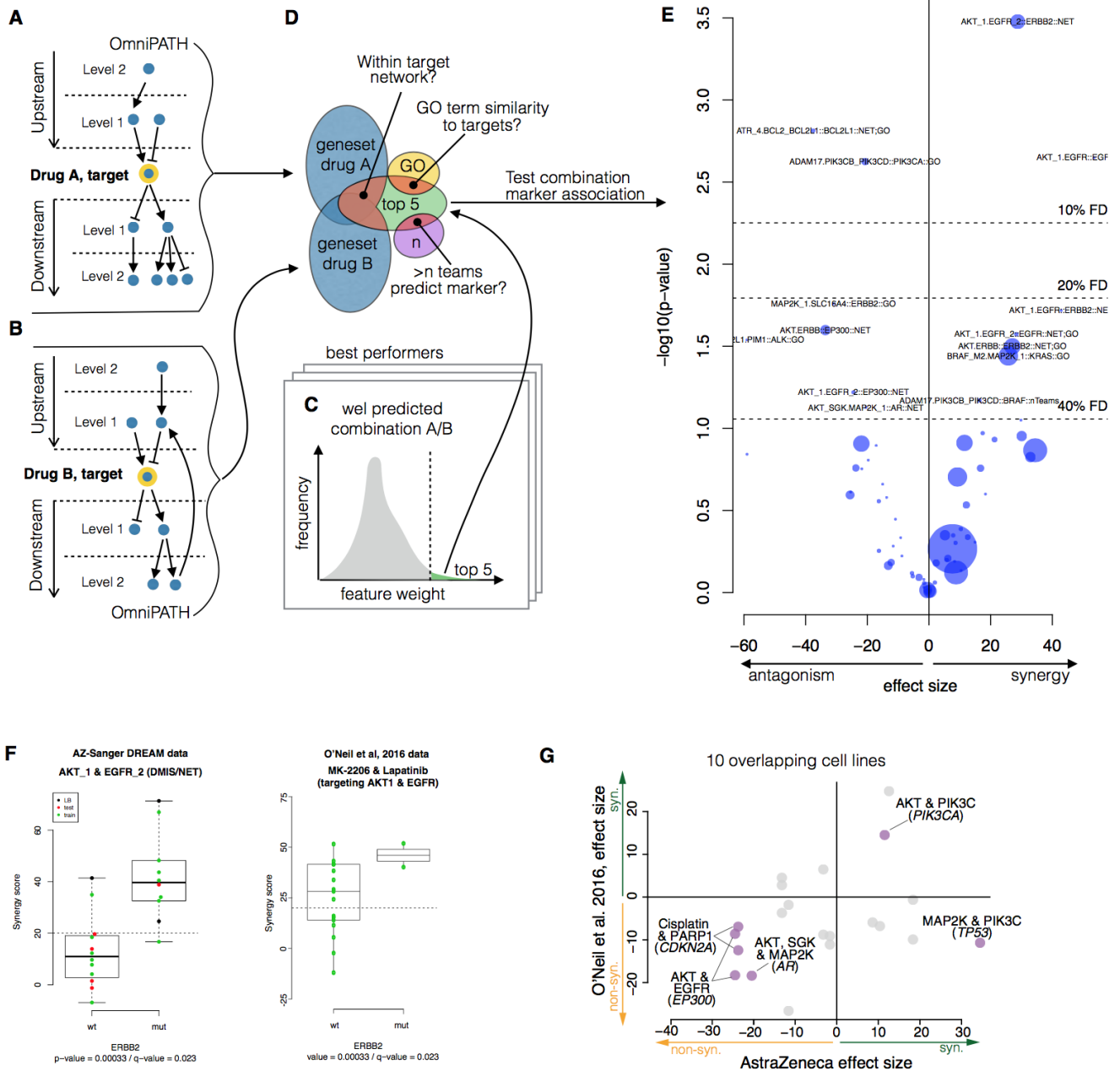
Supp. Figure 4: (A) Performance across different types of methods used by participants in each sub-challenge. Each bar represents occurrence of the method in subch 1A, 1B, and 2. (B) Performance of individual teams coloured by the primary type of machine learning method used in each sub-challenge. Indicators below each bar show cases where kernel and boosting techniques were used with the primary method.



Supp. Figure 5: Training and test set differences between well and poorly predicted combinations based on (A) number of experiments for each combination, (B) proportion of experiments with high quality (QA=1), and (C) average synergy score for each combination.

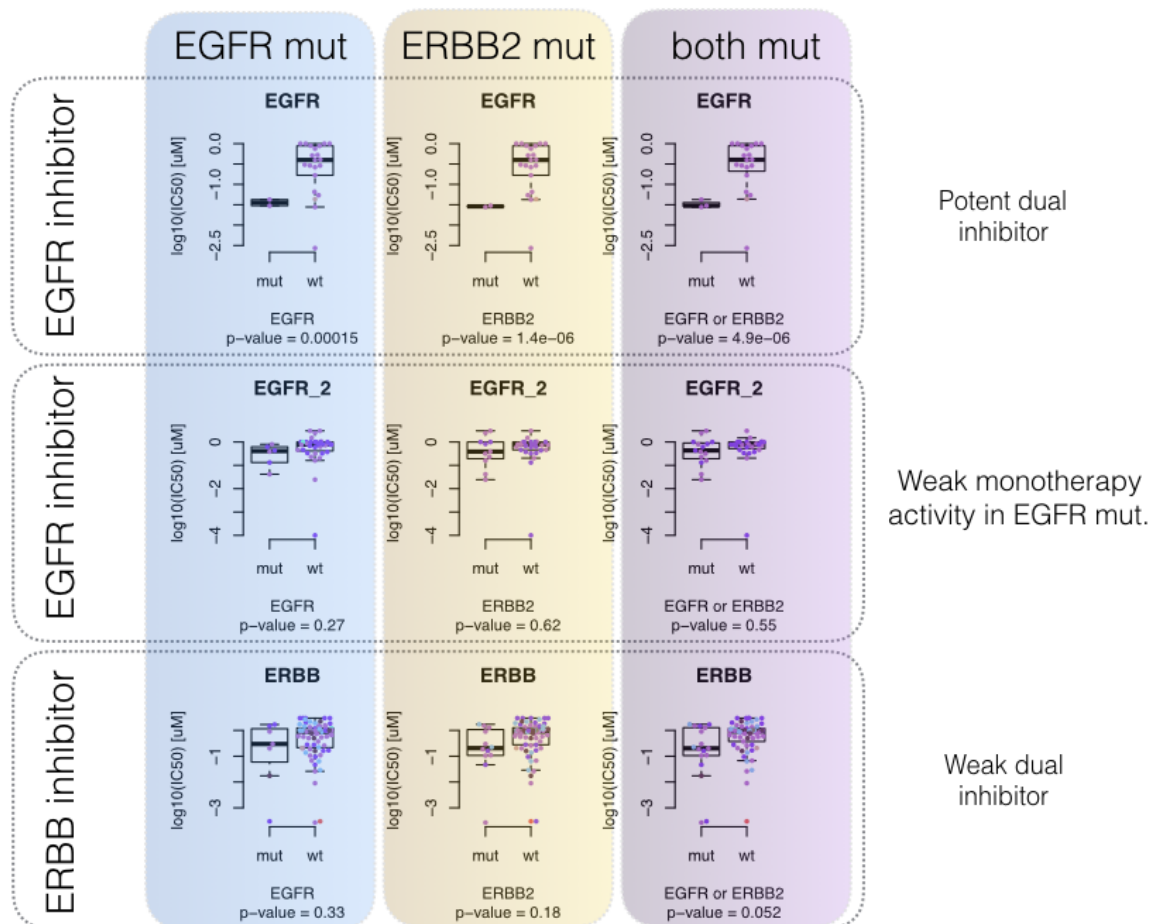


Supp. Figure 6: Alteration frequency in drug targets and nearest neighbours. (A) Copy number and somatic mutation frequency in the gene targets of drug combinations. (B) Copy number and somatic mutation frequency in the genes that are nearest interacting neighbors of the targets, as determined by *OmniPath*.

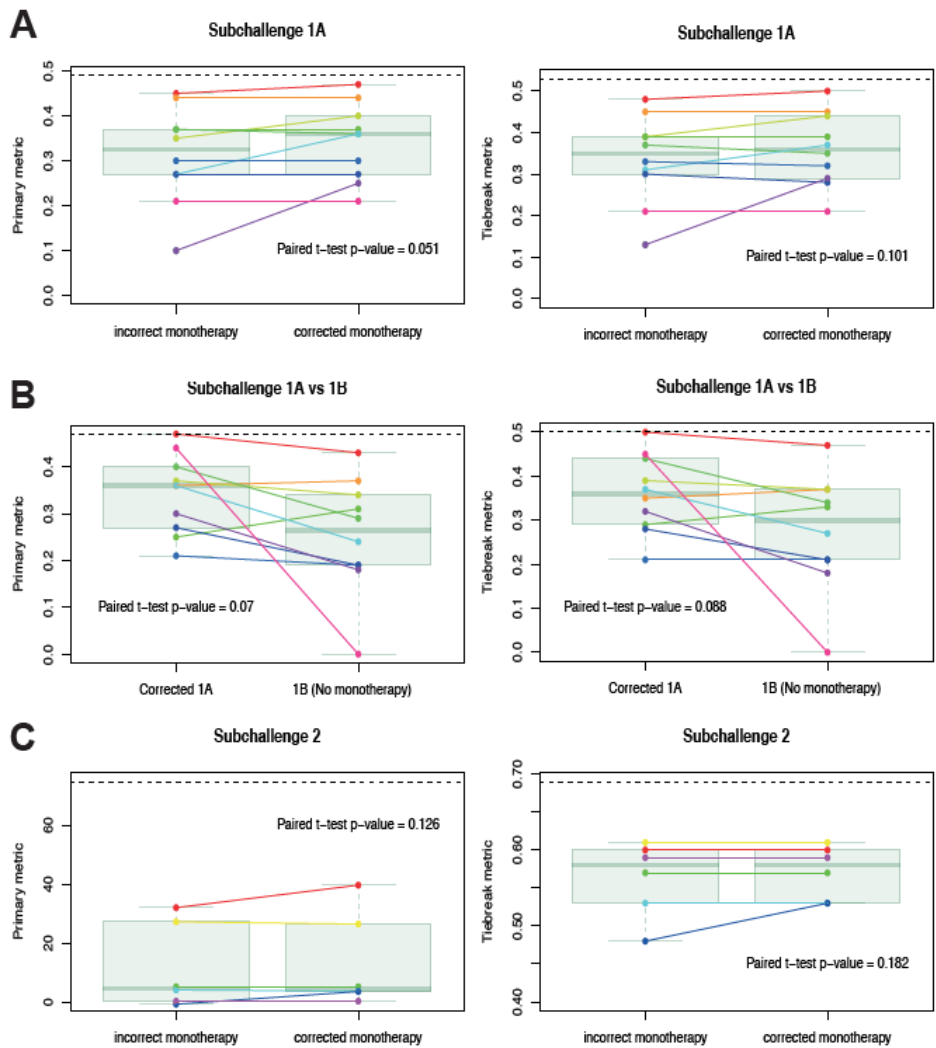


Supp. Figure 7: Post-hoc analysis of synergy biomarkers. (A) and (B) shows the target centric exploration of 2 levels up- and downstream of the putative drug targets from each combination. (C) The top 5 ranked features from well predicted models were chosen for exploring their target enrichment. (D) Additionally, the top 5 features were further investigated if they had GO term similarity to the putative target larger than 0.5, or two independent teams relied on the same features. (E) This putative gene-to-combination association set was tested with an ANOVA model. Here exemplified with (F) AKT combined with EGFR in ERBB2 mutants showed synergy and independently validated with O'Neil et al.

2016. (G) General validation of synergy biomarkers in only overlapping cell lines between AZ-DREAM and O’Neil et al. 2016.



Supp. Figure 8: Monotherapy markers of EGFR and EBB2 inhibitors. Exploration of EGFR, ERBB mutations and copy number changes alone and in combinations.



Supp. Figure 9: Value of monotherapy on prediction performance. (A) Prediction performance of teams (dots) in sub-challenge 1A when given correct and incorrect monotherapy data. (B) Comparison of performance between sub-challenge 1A and 1B for teams that had correct monotherapy data but could use it in 1A and could not in 1B. (C) Prediction performance of teams in sub-challenge 2 when given correct and incorrect monotherapy data. Horizontal dashed line indicates the level of the top performing team.

Supplemental tables

Supplementary Table 4: Prediction performance on independent combinations screen

		Performance: Average weighted Pearson Correlation			
Team	Method	All experiments	Same Cells	Similar Drug	Similar Combination
Mikhail	1A model	0.04	0.06	0.05	0.1
Mikhail	1B model	-0.05	-0.07	-0.02	0
NorthAtlanticDream	1A model	0.05	0.05	0.05	0.03
NorthAtlanticDream	1B model	0.03	0.07	0.05	0.07
DMIS	new 1A model (Deep Learning)	0.11	0.13	0.11	0.1
DMIS	1A model	0.08	0.12	0.08	0.03
DMIS	1B model	-0.03	0.01	0	-0.05
Ensemble	Average 1A model	0.13	0.17	0.13	0.11

Supplemental source code

1. Winning method (code freeze)
2. Scoring code is available online at <https://www.synapse.org/#!/Synapse:syn4991619>

Supplementary references

1. O’Neil, J. *et al.* An Unbiased Oncology Compound Screen to Identify Novel Combination Strategies. *Mol. Cancer Ther.* **15**, 1155–1162 (2016).
2. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001).
3. Chen, T. & Guestrin, C. XGBoost: A Scalable Tree Boosting System. in *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* 785–794 (ACM, 2016).
4. Yadav, B., Wennerberg, K., Aittokallio, T. & Tang, J. Searching for Drug Synergy in Complex Dose–Response Landscapes Using an Interaction Potency Model. *Comput. Struct. Biotechnol. J.* **13**, 504–513 (2015).
5. Gene Ontology Consortium. Gene Ontology Consortium: going forward. *Nucleic Acids Res.* **43**, D1049–56 (2015).

6. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference resource for gene and protein annotation. *Nucleic Acids Res.* **44**, D457–62 (2016).
7. Babur, Ö. *et al.* Systematic identification of cancer driving signaling pathways based on mutual exclusivity of genomic alterations. *Genome Biol.* **16**, 45 (2015).
8. Sun, Y. *et al.* Combining genomic and network characteristics for extended capability in predicting synergistic drugs for cancer. *Nat. Commun.* **6**, 8481 (2015).
9. Woods, D. & Turchi, J. J. Chemotherapy induced DNA damage response: convergence of drugs and pathways. *Cancer Biol. Ther.* **14**, 379–389 (2013).
10. An, O., Dall’Olio, G. M., Mourikis, T. P. & Ciccarelli, F. D. NCG 5.0: updates of a manually curated repository of cancer genes and associated properties from cancer mutational screenings. *Nucleic Acids Res.* **44**, D992–9 (2016).
11. Wagner, A. H. *et al.* DGIdb 2.0: mining clinically relevant drug–gene interactions. *Nucleic Acids Res.* **44**, D1036–D1044 (2016).
12. Yang, W. *et al.* Genomics of Drug Sensitivity in Cancer (GDSC): a resource for therapeutic biomarker discovery in cancer cells. *Nucleic Acids Res.* **41**, D955–61 (2013).
13. Caruana, R., Niculescu-Mizil, A., Crew, G. & Ksikes, A. Ensemble Selection from Libraries of Models. in *Proceedings of the Twenty-first International Conference on Machine Learning* 18– (ACM, 2004).
14. Lee, S. *et al.* BEST: Next-Generation Biomedical Entity Search Tool for Knowledge Discovery from Biomedical Literature. *PLoS One* **11**, e0164680 (2016).
15. Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S. & Dean, J. Distributed Representations of Words and Phrases and their Compositionality. in *Advances in Neural Information Processing Systems 26* (eds. Burges, C. J. C., Bottou, L., Welling, M., Ghahramani, Z. & Weinberger, K. Q.) 3111–3119 (Curran Associates, Inc., 2013).
16. Asgari, E. & Mofrad, M. R. K. Continuous Distributed Representation of Biological Sequences for Deep Proteomics and Genomics. *PLoS One* **10**, e0141287 (2015).

Consortium Authors

Jordi Abante, Barbara Schmitz Abecassis⁴⁶, Nanne Aben^{86,92}, Delasa Aghamirzaie⁹⁵, Tero Aittokallio¹¹⁴, Farida S. Akhtari¹⁴, Bissan Al-lazikani¹⁵⁹, Tanvir Alam³⁵, Amin Allam¹²⁰, Chad Allen²⁸, Mariana Pelicano de Almeida⁴⁶, Doaa Altarawy^{4,63}, Vinicius Alves⁹¹, Alicia Amadoz⁴⁷, Benedict Anchang⁸⁰, Albert A. Antolin¹⁵⁹, Jeremy R. Ash¹², Victoria Romeo Aznar¹²⁴, Wail Balawi³⁵, Moeen Bagheri¹⁷¹, Vladimir Bajic³⁵, Gordon Ball¹⁶³, Pedro J. Ballester^{2,32,107,112}, Delora Baptista²⁵, Christopher Bare¹⁴⁴, Mathilde Bateson¹¹¹, Andreas Bender²⁸, Denis Bertrand³⁴, Bhagya¹⁰⁶, Keith A. Boroevich²³, Evert Bosdriesz^{17,86}, Salim Bougouffa³⁵, Gergana Bounova⁸⁶, Thomas Brouwer³⁶, Barbara Bryant³⁷, Krishna Bulusu⁶, Manuel Calaza³⁰, Alberto Calderone⁹, Stefano Calza^{71,75}, Stephen Capuzzi⁹¹, Jose Carbonell-Caballero²⁶, Daniel Carlin⁷³, Hannah Carter⁷³, Luisa Castagnoli⁴⁸, Remzi Celebi⁵³, Gianni Cesareni⁴⁸, Hyeokyoon Chang⁵⁵, Guocai Chen¹⁶⁹, Haoran Chen^{125,156}, Lijun Cheng¹⁰⁶, Ariel Chernomoretz¹²⁴, Davide Chicco¹³⁸, sunghwan cho⁴², Kwang-Hyun Cho⁴², Daeseon Choi, Jaejoon Choi⁸, Kwanghun Choi⁵⁵, Minsoo Choi⁴², Martine De Cock¹⁷⁴, Elizabeth Coker¹⁵⁹, Isidro Cortes-Ciriano⁴¹, Miklos Cserzo^{79,122}, Cankut Cubuk³¹, Christina Curtis¹⁵⁴, Dries Van Daele¹²¹, Cuong C. Dang^{2,32,107,112}, Tjeerd Dijkstra¹²⁷, Joaquin Dopazo³¹, Sorin Draghici^{64,78}, Anastasios Drosou²⁹, Jonathan Dry⁶, Michel Dumontier¹¹⁶, Friederike Ehrhart⁴⁶, Fatma-Elzahraa Eid^{3,62}, Mahmoud ElHefnawi^{132,133}, Haitham Elmarakeby^{3,62}, Bo van Engelen⁴⁵, Hatice Billur Engin⁷³, Iwan de Esch⁸⁵, Chris Evelo⁴⁶, Andre O Falcao¹²³, Sherif Farag⁹¹, Carlos Fernandez-Lozano¹⁶⁵, Kathleen Fisch²², Asmund Flobak^{134,151}, Chiara Fornari, Amir B K Foroushani¹⁵⁷, Donatien Chedom Fotso, Denis Fourches¹³, Stephen Friend¹⁴⁴, Arnaldo Frigessi¹³⁵, Feng Gao⁵⁰, Xiaoting Gao, Jeffrey M. Gerold¹⁴⁰, Pierre Gestraud¹⁰⁸, Samik Ghosh¹⁶¹, Jussi Gillberg⁵⁴, Antonia Godoy-Lorite⁴⁰, Lizzy Godynyuk⁴⁶, Adam Godzik¹⁴⁵, Anna Goldenberg¹⁰¹, David Gomez-Cabrero^{130,163}, Mehmet Gonen^{33,147}, Chris de Graaf⁸⁵, Harry Gray¹²⁹, Maxim Grechkin¹⁷³, Yuanfang Guan¹⁶⁷, Roger Guimera^{40,104}, Justin Guinney¹⁴⁴, Emre Guney¹¹⁹, Benjamin Haibe-Kains¹³⁹, Younghyun Han⁴², Takeshi Hase¹⁶¹, Di He¹³⁷, Liye He¹¹⁴, Lenwood S. Heath⁶², Kristoffer H. Hellton⁷⁰, Manuela Helmer-Citterich⁴⁸, Marta R. Hidalgo³¹, Daniel Hidru¹⁷¹, Steven M. Hill¹²⁹, Seungpyo Hong, Eivind Hovig^{69,82}, Ya-Chih Hsueh¹⁷¹, Zhiyuan Hu, Zhiyuan Hu, Justin K Huang¹⁰, R. Stephanie Huang¹⁶⁸, Laszlo Hunyady^{79,122}, Jinseub Hwang⁵⁷, Tae Hyun Hwang, Woochang Hwang⁸, Yongdeuk Hwang⁶¹, Olexandr Isayev⁹¹, Oliver Bear Don't Walk IV⁴⁹, John Jack¹⁴, Samad Jahandideh¹⁴⁵, Minji Jeon⁵⁵, Jiadong Ji^{51,59,148}, Yousang Jo⁴³, Piotr J. Kamola²³, Georgi K. Kanev⁸⁵, Jaewoo Kang^{55,117}, Loukia Karacosta⁸⁰, Mostafa Karimi^{65,125}, Samuel Kaski⁵⁴, Marat Kazanov^{1,150}, Abdullah M Khamis³⁵, Suleiman Ali Khan¹¹⁴, Narsis A. Kiani⁵, Allen Kim¹⁷², Jinhan Kim¹⁵³, Juntae Kim⁴², Kiseong Kim⁴³, Kyung Kim⁴⁴, Sunkyu Kim⁵⁵, Yongsoo Kim^{86,88}, Yunseong Kim⁴², Paul D. W. Kirk¹²⁹, Hiroaki Kitano¹⁶¹, Gunter Klambauer¹¹⁵, David Knowles^{68,80}, Melissa Ko⁵², Alvaro Kohn-Luque¹³⁶, Albert J. Kooistra⁸⁵, Melaine A. Kuenemann¹³, Martin Kuiper¹³⁴, Christoph Kurz⁹⁸, Mijin Kwon⁴³, Astrid L_greid¹³⁴, Twan van Laarhoven³⁹, Simone Lederer³⁹, Heewon Lee¹¹⁷, Jeon Lee⁸³, Yun Woo Lee, Eemeli Lepp_aho⁵⁴, Richard Lewis²⁸, Lang Li¹⁶⁰, James Liley⁷⁴, Weng Khong Lim^{27,149}, Chieh Lin¹⁹, Yiyi Liu¹⁷⁶, Yosvany Lopez^{23,38,72}, Joshua Low⁴⁶, Artem Lysenko²³, Daniel

Machado²⁵, Neel Madhukar¹⁷⁵, Dries De Maeyer⁹⁶, Ana Belen Malpartida⁴⁶, Hiroshi Mamitsuka¹¹, Francesco Marabita¹⁶³, Kathleen Marchal⁹⁶, Pekka Marttinen⁹⁹, Daniel Mason²⁸, Mike J Mason¹⁴⁴, Alireza Mazaheri⁵, Arfa Mehmood, Ali Mehreen, Michael P. Menden⁶, Magali Michaut⁸⁶, Ryan A. Miller⁴⁶, Costas Mitsopoulos¹⁵⁹, Dezso Modos²⁸, Marijke Van Moerbeke, Keagan Moo¹⁷², Alison Motsinger-Reif¹⁴, Rajiv Movva⁶⁸, Sebastian Muraru⁴⁶, Eugene Muratov⁹¹, Mushthofa Mushthofa⁹⁷, Niranjana Nagarajan³⁴, Sigve Nakken⁸¹, Aritro Nath¹⁶⁸, Elias Neto¹⁴⁴, Pierre Neuvial¹¹⁰, Richard Newton, Tin Nguyen⁵⁶, Zheng Ning⁷¹, Carlos De Niz¹⁵⁸, Thea Norman¹⁴⁴, Baldo Oliva¹⁵⁵, Catharina Olsen^{118,126}, Antonio Palmeri⁴⁸, Bhawan Panesar¹⁷¹, Stavros Papadopoulos²⁹, seonyeong park, Jaesub Park⁴³, Sungjoon Park⁵⁵, Yudi Pawitan⁷¹, Daniele Peluso⁴⁸, Sriram Pendyala, Jian Peng¹⁶⁶, Livia Perfetto⁴⁸, Stefano Pirro⁴⁸, Sylvia Plevritis⁸⁰, Regina Politi⁹¹, Hoifung Poon¹²⁸, Eduard Porta¹⁴⁵, Isak Prellner¹⁶, Kristina Preuer¹¹⁵, Miguel Angel Pujana²¹, Ricardo Ramnarine¹⁷¹, John E. Reid¹²⁹, Fabien Reyal¹⁰⁹, Sylvia Richardson¹²⁹, Camir Ricketts¹⁷⁵, Linda Rieswijk^{46,84}, Miguel Rocha²⁵, Carmen Rodriguez-Gonzalvez¹⁵⁹, Kyle Roell¹⁴, Daniel Rotroff¹⁴, Julian R. de Ruiter^{86,87}, Rukawa, Benjamin Sadacca¹⁰⁹, Julio Saez-Rodriguez^{90,143}, Zhaleh Safikhani¹³⁹, Fita Safitri, Francisco Salavert, Marta Sales-Pardo⁴⁰, Sebastian Sauer⁹⁴, Moritz Schlichting⁴⁵, Jose A. Seoane¹⁵⁴, Jordi Serra^{7,20}, Ming-Mei Shang¹⁶³, Alok Sharma^{23,113,170}, Hari Sharma¹⁷¹, Yang Shen^{65,125}, Motoki Shiga⁶⁶, Moonshik Shin⁴³, Ziv Shkedy, Kevin Shopsowitz⁹³, Sam Sinal¹⁴⁰, Dylan Skola¹⁰, Petr Smirnov¹³⁹, Izel Fourie Soerensen²⁴, Peter Soerensen²⁴, Je-Hoon Song⁴², Sang Ok Song¹⁵³, Othman Soufan³⁵, Andreas Spitzmueller, Boris Steipe¹⁷¹, Gustavo Stolovitzky^{102,103}, Chayaporn Suphavilai^{34,60}, Bence Szalai^{79,122,142}, Sergio Pulido Tamayo¹²¹, David Tamborero¹⁴¹, Jing Tang¹¹⁴, Zia-ur-Rehman Tanoli¹¹⁴, Marc Tarres-Deulofeu⁴⁰, Jesper Tegner^{15,163}, Liv Thommesen¹³⁴, Seyed Ali Madani Tonekaboni¹³⁹, Hong Tran⁶², Ewoud De Troyer, Amy Truong¹⁴⁴, Tatsuhiko Tsunoda^{23,38,72}, Gabor Turu^{79,122}, Guang-Yo Tzeng¹⁷¹, Lieven Verbeke⁹⁶, Giovanni Di Veroli⁶, Santiago Videla¹²⁴, Daniel Vis⁸⁶, Andrey Voronkov⁸¹, konstantinos votis²⁹, Ashley Wang¹⁷¹, Dennis Wang¹⁶², Hong-Qiang Horace Wang¹⁸, Po-Wei Wang¹⁹, Sheng Wang¹⁶⁶, Wei Wang⁵⁰, Xiaochen Wang¹⁷⁶, Xin Wang⁵⁰, Krister Wennerberg¹¹⁴, Lorenz Wernisch¹²⁹, Lodewyk Wessels^{17,86,92}, Gerard JP van Westen⁸⁹, Bart A. Westerman⁷⁷, Simon Richard White¹²⁹, Egon Willighagen⁴⁶, Russ Wolfinger¹⁴⁶, Tom Wurdinger⁷⁶, Lei Xie⁵⁸, Shuilian Xie^{125,156}, Hua Xu¹⁶⁹, Bhagwan Yadav¹⁰⁰, Christopher Yau, Christopher Yau, Huwate Yeerna, Jia Wei Yin¹⁷¹, Michael Yu¹⁰, MinHwan Yu⁵⁵, Thomas Yu¹⁴⁴, So Jeong Yun^{152,153}, So Jeong Yun^{152,153}, Alexey Zakharov¹³¹, Alexandros Zamichos²⁹, Massimiliano Zanin¹⁶⁴, Mikhail Zaslavskiy, Li Zeng¹⁷⁶, Hector Zenil⁵, Frederick Zhang¹⁷¹, Pengyue Zhang¹⁰⁵, Wei Zhang⁷³, Hongyu Zhao¹⁷⁶, Lan Zhao⁶⁷, Wenjin Zheng¹⁶⁹, Azedine Zoufir²⁸, Manuela Zucknick¹³⁶

1. A.A. Kharkevich Institute for Information Transmission Problems Moscow Russia
2. Aix-Marseille University UM105 France
3. Al-Azhar University Cairo Egypt
4. Alexandria University Alexandria Egypt
5. Algorithmic Dynamics Lab Unit of Computational Medicine Department of Medicine Solna SciLifeLab Center for Molecular Medicine Karolinska Institute

6. AstraZeneca
7. Bellvitge Biomedical Biomedical Research Institute (IDIBELL)
8. Bio-Synergy Research Center Daejeon Republic of Korea.
9. Bioinformatics and Computational Biology Unit Department of Biology University of Rome Tor Vergata Rome 00133 Italy
10. Bioinformatics and Systems Biology Program University of California San Diego La Jolla CA
11. Bioinformatics Center Institute for Chemical Research Kyoto University Japan
12. Bioinformatics Research Center Department of Chemistry Department of Statistics North Carolina State University Raleigh NC
13. Bioinformatics Research Center Department of Chemistry North Carolina State University Raleigh NC
14. Bioinformatics Research Center North Carolina State University Raleigh NC
15. Biological and Environmental Science and Engineering Division KAUST
16. Cambridge Rindge and Latin High School
17. Cancer Genomics Netherlands
18. Cancer Hospital of Chinese Academy of Sciences Hefei China
19. Carnegie Mellon University
20. Catalan Institute of Oncology (ICO)
21. Catalan Institute of Oncology (ICO) Bellvitge Biomedical Biomedical Research Institute (IDIBELL)
22. Center for Computational Biology and Bioinformatics University of California San Diego La Jolla CA
23. Center for Integrative Medical Sciences RIKEN Japan
24. Center for Quantitative Genetics and Genomics Department of Biology and Genetics Aarhus University
25. Centre Biological Engineering (CEB) University of Minho
26. Centre de Regulacio Genomica (CRG) Barcelona Institute for Science and Technology Barcelona Spain
27. Centre for Computational Biology Duke-NUS Medical School Singapore
28. Centre for Molecular Informatics Department of Chemistry University of Cambridge
29. CERTH-ITI
30. CIMUS University of Santiago de Compostela
31. Clinical Bioinformatic Area Fundacion Progreso y Salud CDCA Hospital Virgen del Rocio Sevilla Spain
32. CNRS UMR7258 Marseille France.
33. College of Engineering Koc University Istanbul Turkey
34. Computational and Systems Biology Genome Institute of Singapore
35. Computational Bioscience Research Center (CBRC) KAUST
36. Computer Laboratory University of Cambridge
37. Constellation Pharmaceuticals
38. CREST JST Japan
39. Data Science Radboud University Netherlands
40. Departament d'Enginyeria Quimica Universitat Rovira i Virgili
41. Departement de Biologie Structurale et Chimie Institut Pasteur Unite de Bioinformatique Structurale CNRS UMR 3825 Paris France
42. Department of Bio and Brain Engineering Korea Advanced Institute of Science and Technology (KAIST)
43. Department of Bio and Brain Engineering Korea Advanced Institute of Science and Technology Daejeon Republic of Korea.

44. Department of Bioengineering University of Washington
45. Department of Bioinformatics - BiGCaT NUTRIM Maastricht University Maastricht 6229 ER Maastricht The Netherlands
46. Department of Bioinformatics - BiGCaT NUTRIM Maastricht University The Netherlands
47. Department of Bioinformatics Igenomix SL Valencia Spain.
48. Department of Biology University of Rome Tor Vergata Rome Italy
49. Department of Biomedical Informatics Columbia University in the City of New York
50. Department of Biomedical Sciences City University of Hong Kong Hong Kong
51. Department of Biostatistics School of Public Health Shandong University China.
52. Department of Cancer Biology Stanford University Stanford
53. Department of Computer Engineering Ege University Turkey
54. Department of Computer Science Aalto University
55. Department of Computer Science and Engineering Korea University Seoul Korea
56. Department of Computer Science and Engineering University of Nevada Reno NV
57. Department of Computer Science and Statistics Daegu University South Korea
58. Department of Computer Science Hunter College and The Graduate Center The City University of New York
59. Department of Computer Science Hunter College The City University of New York NY
60. Department of Computer Science National University of Singapore
61. Department of Computer Science The University of Suwon Suwon Republic of Korea.
62. Department of Computer Science Virginia Tech Blacksburg VA
63. Department of Computer Science Virginia Tech Blacksburg VA.
64. Department of Computer Science Wayne State University Michigan
65. Department of Electrical and Computer Engineering Texas A
66. Department of Electrical Electronic and Computer Engineering Gifu University
67. Department of Electronic Engineering City University of Hong Kong Hong Kong
68. Department of Genetics Stanford University
69. Department of Informatics University of Oslo Norway
70. Department of Mathematics University of Oslo Norway
71. Department of Medical Epidemiology and Biostatistics Karolinska Institute Stockholm Sweden
72. Department of Medical Science Mathematics Medical Research Institute Tokyo Medical and Dental University Japan
73. Department of Medicine University of California San Diego La Jolla CA
74. Department of Medicine University of Cambridge UK
75. Department of Molecular and Translational Medicine University of Brescia Brescia Italy
76. Department of Neurosurgery Cancer Center Amsterdam HZ Amsterdam The Netherlands
77. Department of Neurosurgery Cancer Center Amsterdam The Netherlands
78. Department of Obstetrics and Gynecology Wayne State University Michigan
79. Department of Physiology Faculty of Medicine Semmelweis University Budapest Hungary
80. Department of Radiology Stanford University
81. Department of Tumor Biology Institute for Cancer Research Oslo University Hospital Norway
82. Department of Tumor Biology Institute for Cancer Research Oslo University Hospital Norway and
83. Dept. Bioinformatics University of Texas Southwestern Medical Center
84. Division of Environmental Health Sciences School of Public Health University of California Berkeley CA United States
85. Division of Medicinal Chemistry AIMMS Vrije Universiteit Amsterdam The Netherlands
86. Division of Molecular Carcinogenesis NKI Amsterdam The Netherlands

87. Division of Molecular Pathology NKI Amsterdam The Netherlands
88. Division of Oncogenomics NKI Amsterdam The Netherlands
89. Drug Discovery and Safety Leiden Academic Centre for Drug Research Leiden University
Einsteinweg 55 2333 CC Leiden The Netherlands
90. EMBL-EBI
91. Eshelman School of Pharmacy University of North Carolina at Chapel Hill
92. Faculty of EEMCS Delft University of Technology The Netherlands
93. Faculty of Medicine University of British Columbia
94. FOM University of Applied Sciences
95. Genetics Bioinformatics and Computational Biology Virginia Tech Blacksburg VA
96. Ghent University
97. Ghent University Bogor Agricultural University
98. Helmholtz Zentrum Munchen Institute of Health Economics and Health Care
Management
99. Helsinki Institute for Information Technology HIIT Department of Computer Science Aalto
University Finland
100. Hematology Research Unit Helsinki Department of Clinical Chemistry and
Hematology University of Helsinki and Helsinki University Hospital Comprehensive
Cancer Center Helsinki Finland
101. Hospital for Sick Children Toronto Ontario Canada
102. IBM Research
103. Icahn School of Medicine at Mt. Sinai
104. ICREA
105. Indiana University - Purdue University Indianapolis
106. Indiana University School of Medicine
107. INSERM U1068 Marseille France
108. Institut Curie Inserm U900 Paris France
109. Institut Curie RT2 Lab Paris France
110. Institut de Mathematiques de Toulouse Universite Paul Sabatier Toulouse France
111. Institut HyperCube Paris France
112. Institut Paoli-Calmettes Marseille France
113. Institute for Integrated and Intelligent Systems Griffith University Australia
114. Institute for Molecular Medicine Finland FIMM University of Helsinki Finland
115. Institute of Bioinformatics Johannes Kepler University Linz Austria
116. Institute of Data Science Maastricht University Maastricht Netherlands
117. Interdisciplinary Graduate Program in Bioinformatics Korea University Seoul
Korea
118. Interuniversity Institute of Bioinformatics in Brussels (IB) Belgium
119. Joint IRB-BSC-CRG Program in Computational Biology Institute for Research in
Biomedicine Barcelona Spain
120. King Abdullah University of Science and Technology (KAUST)
121. KU Leuven
122. Laboratory of Molecular Physiology Hungarian Academy of Sciences and
Semmelweis University (MTA-SE) Budapest
123. LaSIGE Faculdade de Ciencias Universidade de Lisboa Portugal
124. Leloir Institute Buenos Aires Argentina
125. M University
126. Machine Learning Group (MLG) Department d' Informatique Universite libre de
Bruxelles (ULB) Brussels 1050 Belgium
127. Max Planck Institute for Developmental Biology Tuebingen Germany

128. Microsoft Research Redmond
129. MRC Biostatistics Unit University of Cambridge UK
130. Mucosal and Salivary Biology Division King's College London Dental Institute
London UK
131. National Center for Advancing Translational Sciences National Institutes of
Health Rockville MD
132. National Research Centre Cairo Egypt
133. Nile University Cairo Egypt
134. Norwegian University of Science and Technology
135. Oslo Centre for Biostatistics and Epidemiology University of Oslo and Oslo
University Hospital Norway
136. Oslo Centre for Biostatistics and Epidemiology Department of Biostatistics
University of Oslo Norway
137. Ph.D. Program in Computer Science The Graduate Center The City University of
New York
138. Princess Margaret Cancer Centre Toronto Canada
139. Princess Margaret Cancer Centre University Health Network Toronto Ontario
Canada
140. Program for Evolutionary Dynamics Harvard University
141. Research Unit on Biomedical Informatics University Pompeu Fabra UPF
Barcelona Spain
142. RWTH Aachen University Faculty of Medicine Joint Research Center for
Computational Biomedicine Aachen Germany
143. RWTH-Aachen
144. Sage Bionetworks
145. Sanford Burnham Prebys Medical Discovery Institute
146. SAS Institute Inc. Cary NC
147. School of Medicine Koc University Istanbul Turkey
148. School of Statistics Shandong University of Finance and Economics Jinan China
149. SingHealth Duke-NUS Institute of Precision Medicine Singapore
150. Skolkovo Institute of Science and Technology Moscow Russia
151. St. Olav's University Hospital
152. Standigm Inc.
153. Standigm Inc. Seoul Korea
154. Stanford University
155. Structural Bioinformatics Group GRIB IMIM Department of Experimental and Life
Sciences Universitat Pompeu Fabra Barcelona Catalonia Spain
156. Texas A
157. Texas State University San Marcos Texas
158. Texas Tech University
159. The Institute of Cancer Research London UK
160. The Ohio State University College of Medicine Department of Biomedical
Informatics
161. The Systems Biology Institute Tokyo
162. The University of Sheffield
163. Unit of Computational Medicine Department of Medicine Solna SciLifeLab Center
for Molecular Medicine Karolinska Institute
164. Universidad Politecnica de Madrid
165. University of A Coruna
166. University of Illinois at Urbana-Champaign Urbana

167. University of Michigan Department of Computational Medicine and Bioinformatics
168. University of Minnesota
169. University of Texas Health Science Center
170. University of the South Pacific Fiji
171. University of Toronto Toronto Canada
172. University of Washington
173. University of Washington Seattle
174. University of Washington Tacoma
175. Weill Cornell Medicine
176. Yale University