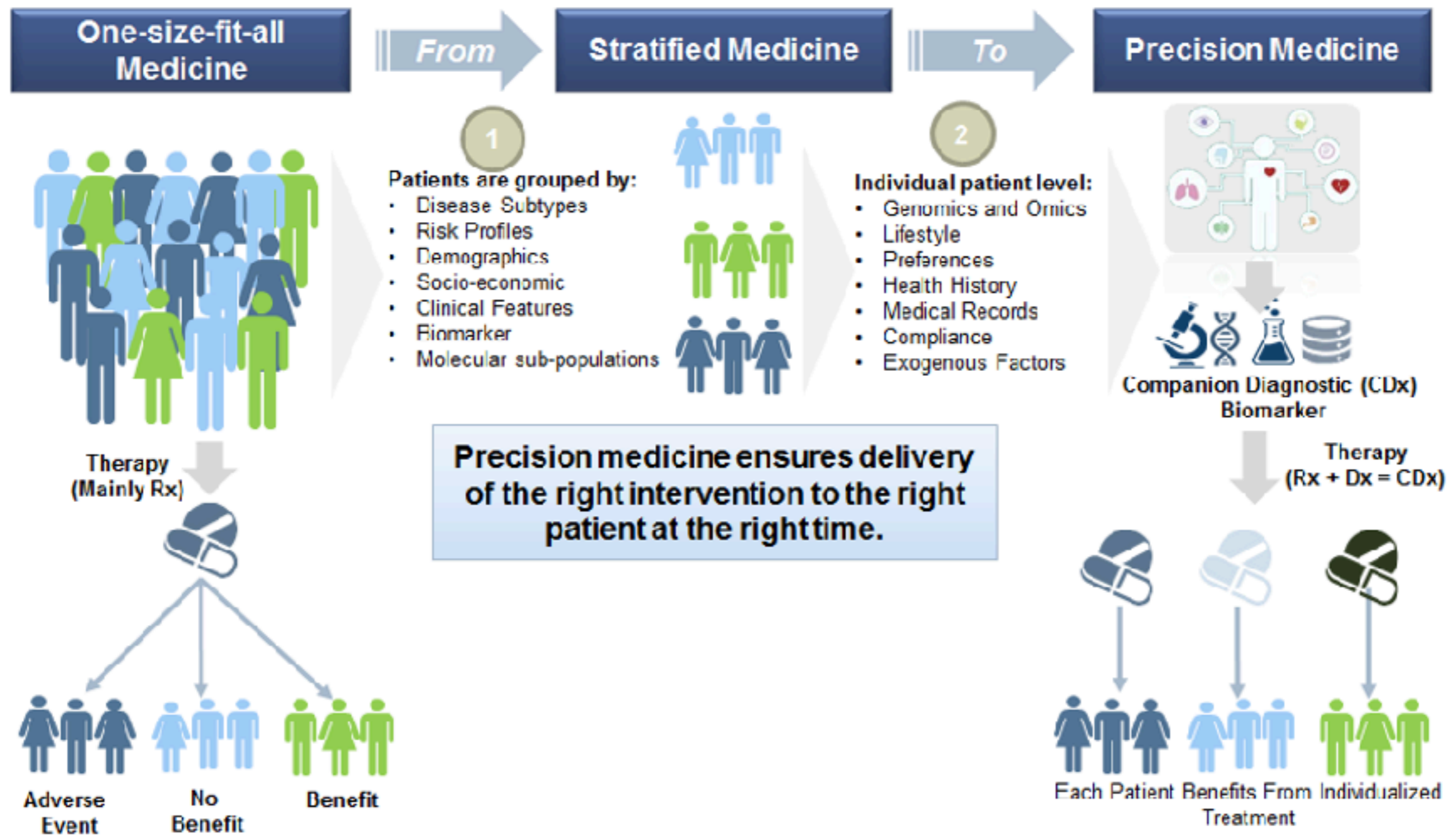


Precision medicine

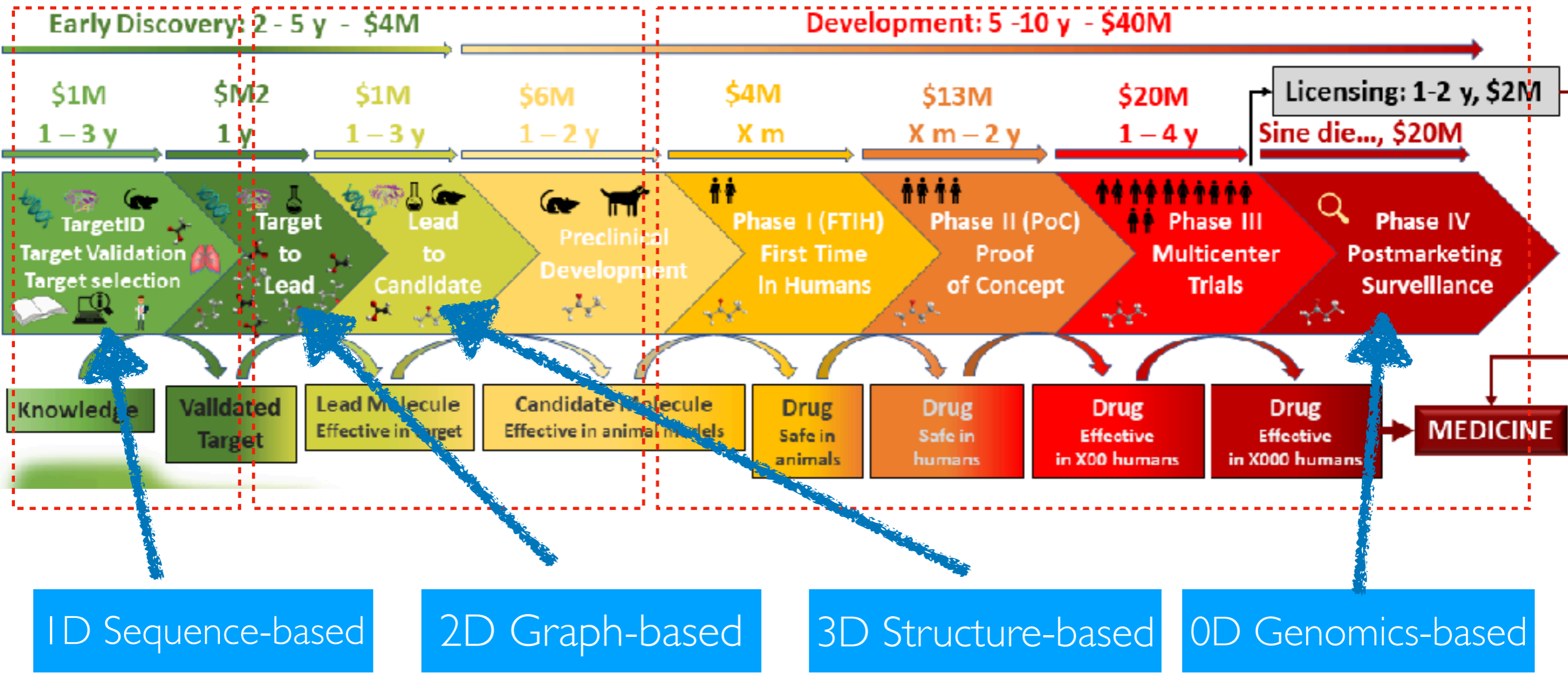
Sheng Wang

Precision medicine:

the right patient, the right drug, the right time, the right dose



Discover a new drug is very time-consuming



Sequence: understand target function using protein sequence. NLP to find targets (word sequence).

Graph: generate compound graph 2D structure (deep generative model)

Structure: modify structure according to 3D structure (geometric deep learning)

Genomics: side effects, personalized efficacy, repurposing, etc. (multi-modality)

- how to reuse an old drug

We don't have so many “drugs”

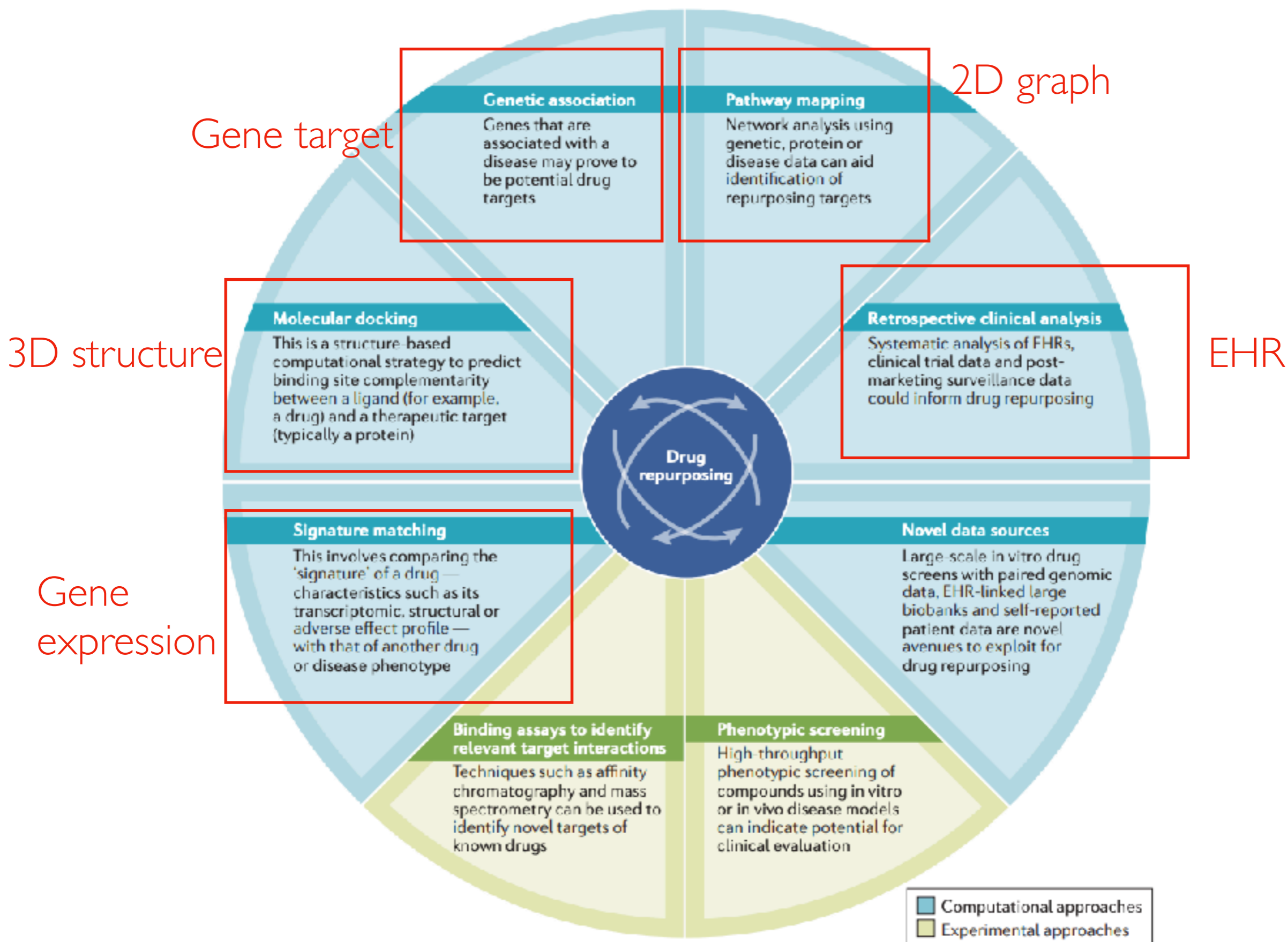
- Discovery new drug?
 - Often not in the scope of precision medicine
 - New patient cannot wait for a new drug
- Drug repurposing
 - Drug A, which is used to treat disease X, is later used to treat disease Y
 - Well-documented side effects and less restriction from FDA
- Drug combination
 - Drug A is not effective. Drug B is not effective. Drug A and B used together is effective.
- Personalized dosage
 - Widely used in clinics. Use genomics data to determine dosage (regression).

Drug repurposing

Table 1 | Selected successful drug repurposing examples and the repurposing approach employed

Drug name	Original indication	New indication	Date of approval	Repurposing approach used	Comments on outcome of repurposing
Zidovudine	Cancer	HIV/AIDS	1987	In vitro screening of compound libraries	Zidovudine was the first anti-HIV drug to be approved by the FDA
Minoxidil	Hypertension	Hair loss	1988	Retrospective clinical analysis (identification of hair growth as an adverse effect)	Global sales for minoxidil were US\$860 million in 2016 (Questec minoxidil sales report, 2017 ; see Related links)
Sildenafil	Angina	Erectile dysfunction	1998	Retrospective clinical analysis	Marketed as Viagra, sildenafil became the leading product in the erectile dysfunction drug market, with global sales in 2012 of \$2.05 billion ⁸
Thalidomide	Morning sickness	Erythema nodosum leprosum and multiple myeloma	1998 and 2006	Off-label usage and pharmacological analysis	Thalidomide derivatives have achieved substantial clinical and commercial success in multiple myeloma
Celecoxib	Pain and inflammation	Familial adenomatous polyps	2000	Pharmacological analysis	The total revenue from Celebrex (Pfizer) at the end of 2014 was \$2.69 billion (Pfizer 2014 financial report ; see Related links)
Atomoxetine	Parkinson disease	ADHD	2002	Pharmacological analysis	Strattera (Eli Lilly) recorded global sales of \$855 million in 2016
Duloxetine	Depression	SUI	2004	Pharmacological analysis	Approved by the EMA for SUI. The application was withdrawn in the US. Duloxetine is approved for the treatment of depression and chronic pain in the US
Rituximab	Various cancers	Rheumatoid arthritis	2006	Retrospective clinical analysis (remission of coexisting rheumatoid arthritis in patients with non-Hodgkin lymphoma treated with rituximab ¹⁴³)	Global sales of rituximab topped \$7 billion in 2015 (RFF ¹⁴⁶)
Raloxifene	Osteoporosis	Breast cancer	2007	Retrospective clinical analysis	Approved by the FDA for invasive breast cancer. Worldwide sales of \$237 million in 2015 (see Related links)
Fingolimod	Transplant rejection	MS	2010	Pharmacological and structural analysis ¹⁴⁴	First oral disease-modifying therapy to be approved for MS. Global sales for fingolimod (Gilenya) reached \$3.1 billion in 2017 (see Related links)
Dapoxetine	Analgesia and depression	Premature ejaculation	2012	Pharmacological analysis	Approved in the UK and a number of European countries; still awaiting approval in the US. Peak sales are projected to reach \$750 million
Topiramate	Epilepsy	Obesity	2012	Pharmacological analysis	Qsymia (Vivus) contains topiramate in combination with phentermine
Ketoconazole	Fungal infections	Cushing syndrome	2014	Pharmacological analysis	Approved by the EMA for Cushing syndrome in adults and adolescents above the age of 12 years (see Related links)
Aspirin	Analgesia	Colorectal cancer	2015	Retrospective clinical and pharmacological analysis	US Preventive Services Task Force released draft recommendations in September 2015 regarding the use of aspirin to help prevent cardiovascular disease and colorectal cancer ⁵²

Approaches used in drug repurposing

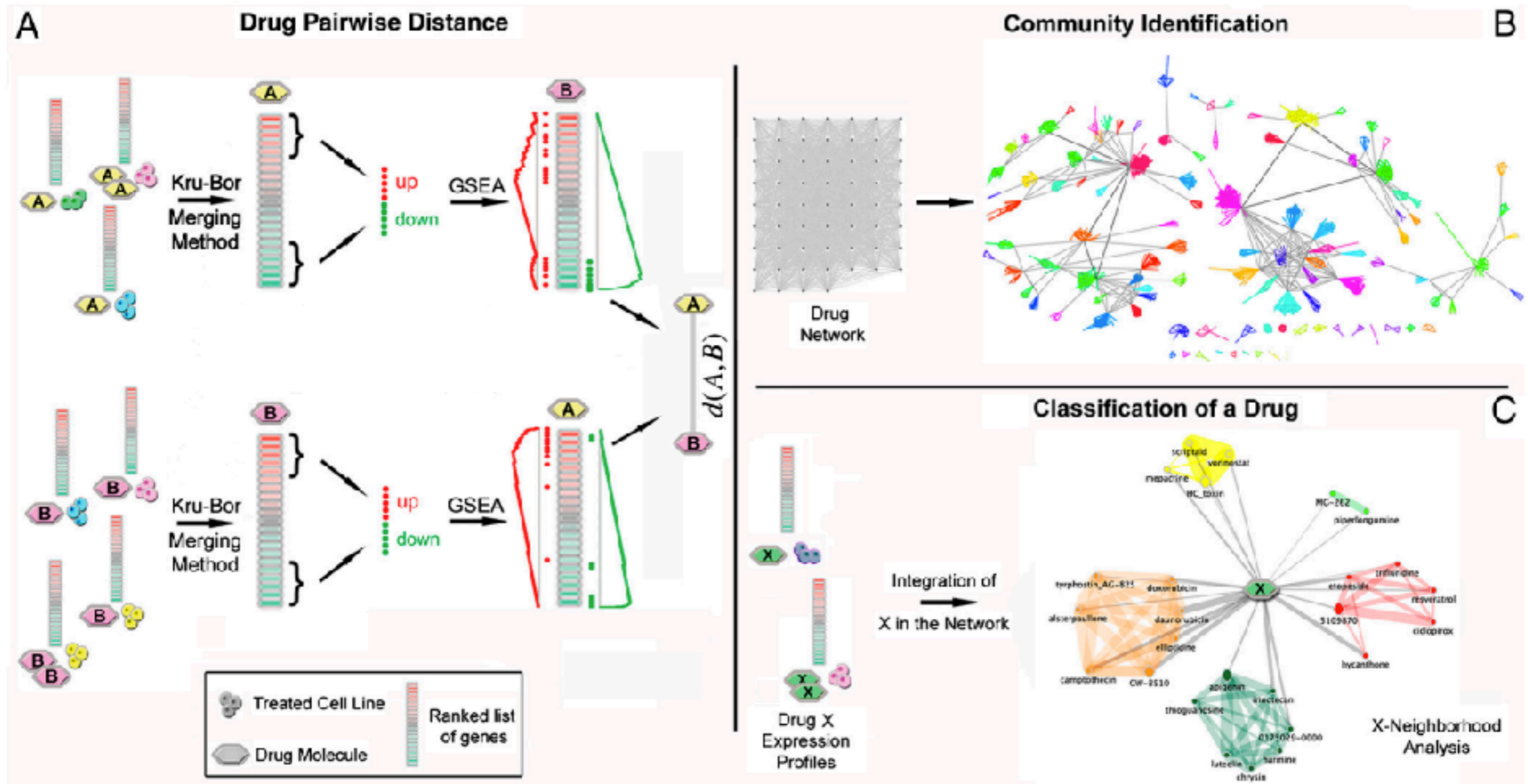


Drug repurposing strategy

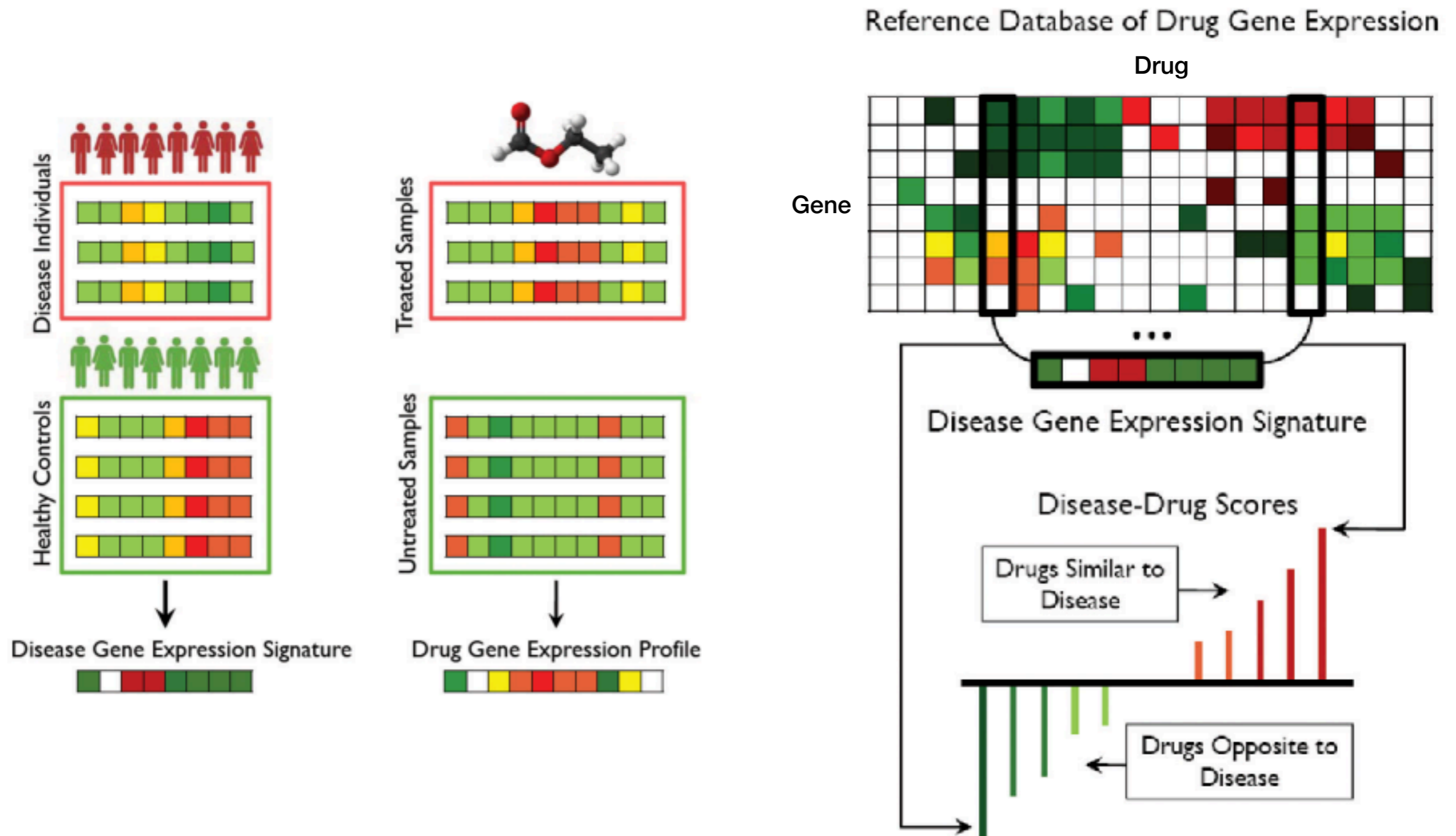
- Drug-based
 - If drug A can cure disease X and is similar to drug B, then B might be also treat X
- Disease-based
 - If disease X and Y have similar profiles and indications, and drug R can cure X, then R can also cure Y.

Use gene expression after treatment

Drugs target on similar proteins or have similar Mode of Actions have similar (after treatment) expression.



Compare disease expression and drug expression



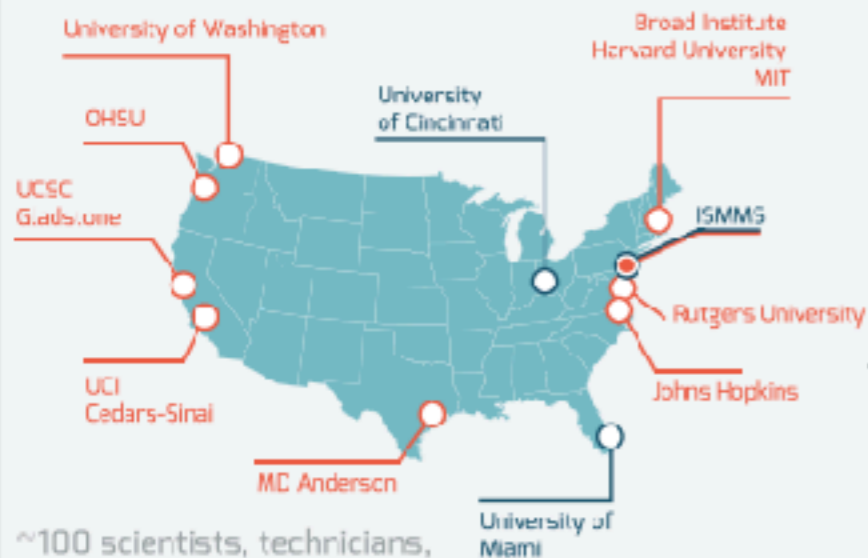
Expression-based drug repurposing

- People realized that the performance (accuracy, coverage) depends on the data, rather than the model
- How about we just generate the expression of X drugs on Y tissues
 - LINCS: Library of Integrated Network-based Cellular Signatures
 - 15 institutions, >1000 cell lines, >5000 drugs, 1000 genes
 - 1.3 million after treatment gene expression vectors
 - cMAP: 3 cell lines, but 20k genes

LINCS

BY THE NUMBERS

15 INSTITUTIONS



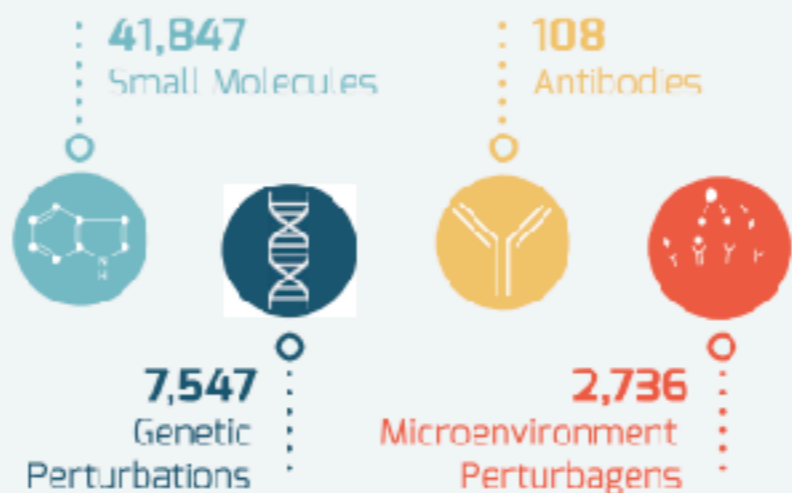
Data Coordination and Integration Center

BD2K-LINCS | ISMMS | University of Miami | University of Cincinnati

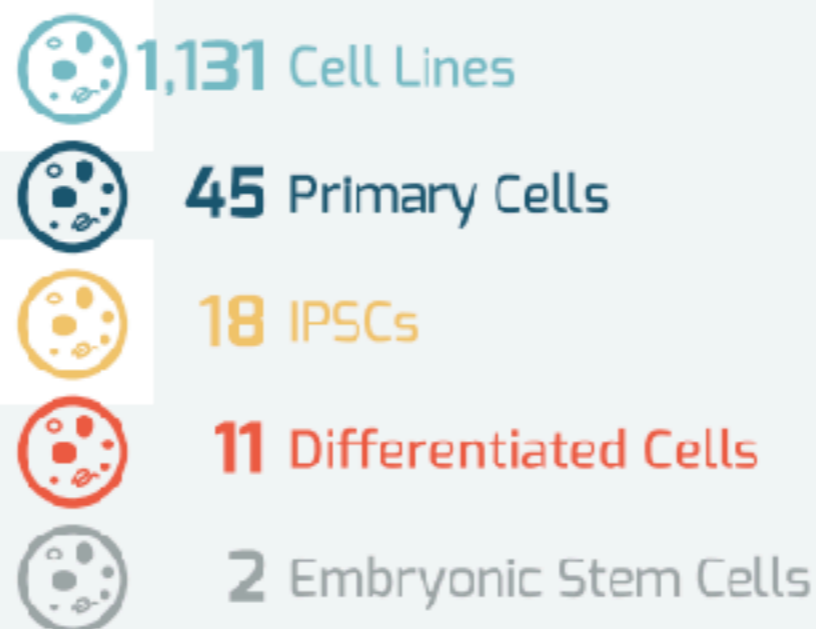
Data and Signature Generation Centers

NeuroLINCS	MEP LINCS	DToxS	FCCE
UC Cedars-Sina Glastone Johns Hopkins MIT	OHSU MD Anderson	ISMMS Rutgers University	Broad Institute Harvard University UCSC
	HMS LINCS	Broad Transcriptomics	

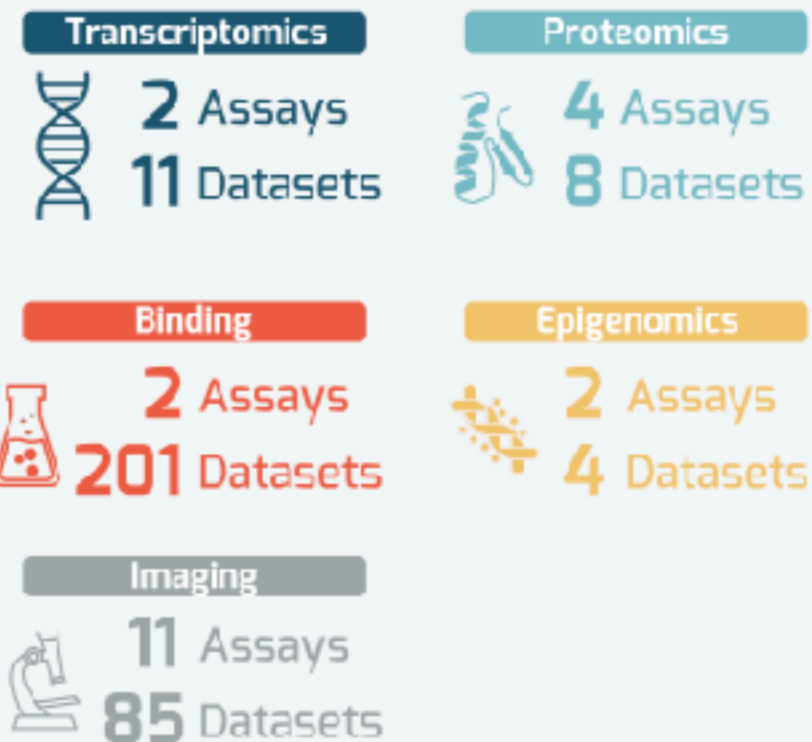
4 PERTURBATION TYPES



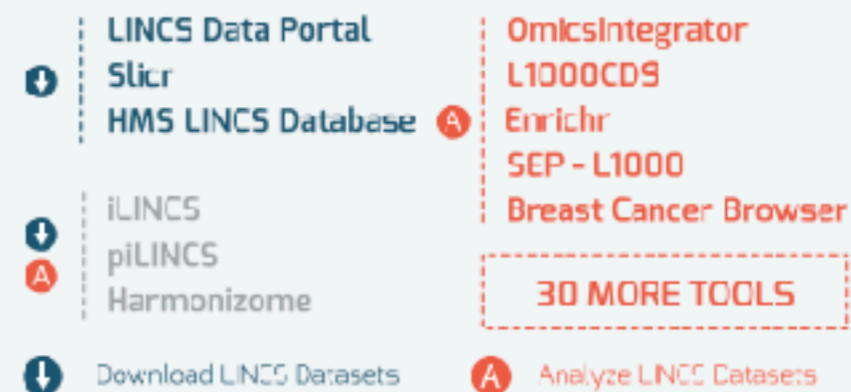
5 CELL TYPES



5 SIGNATURE TYPES



FEATURED LINCS TOOLS

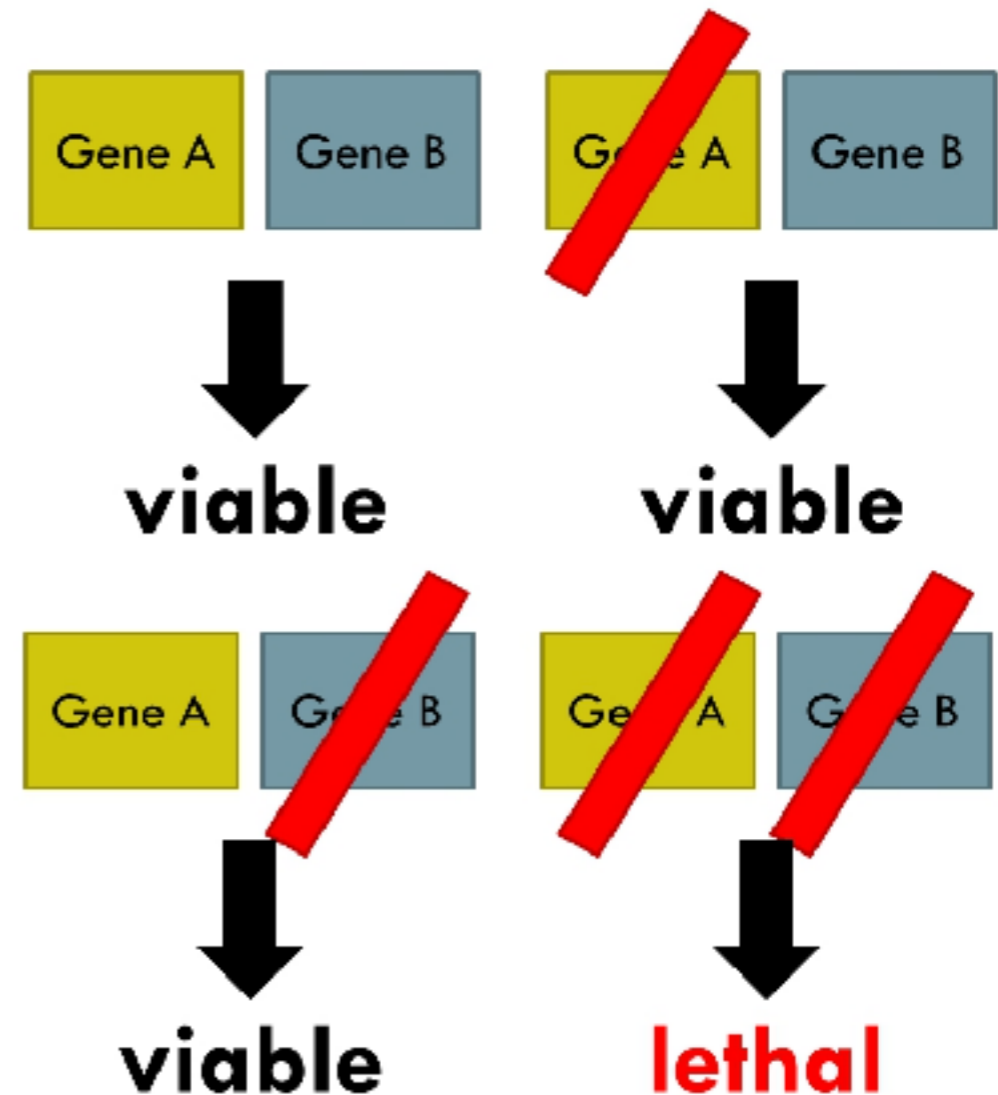
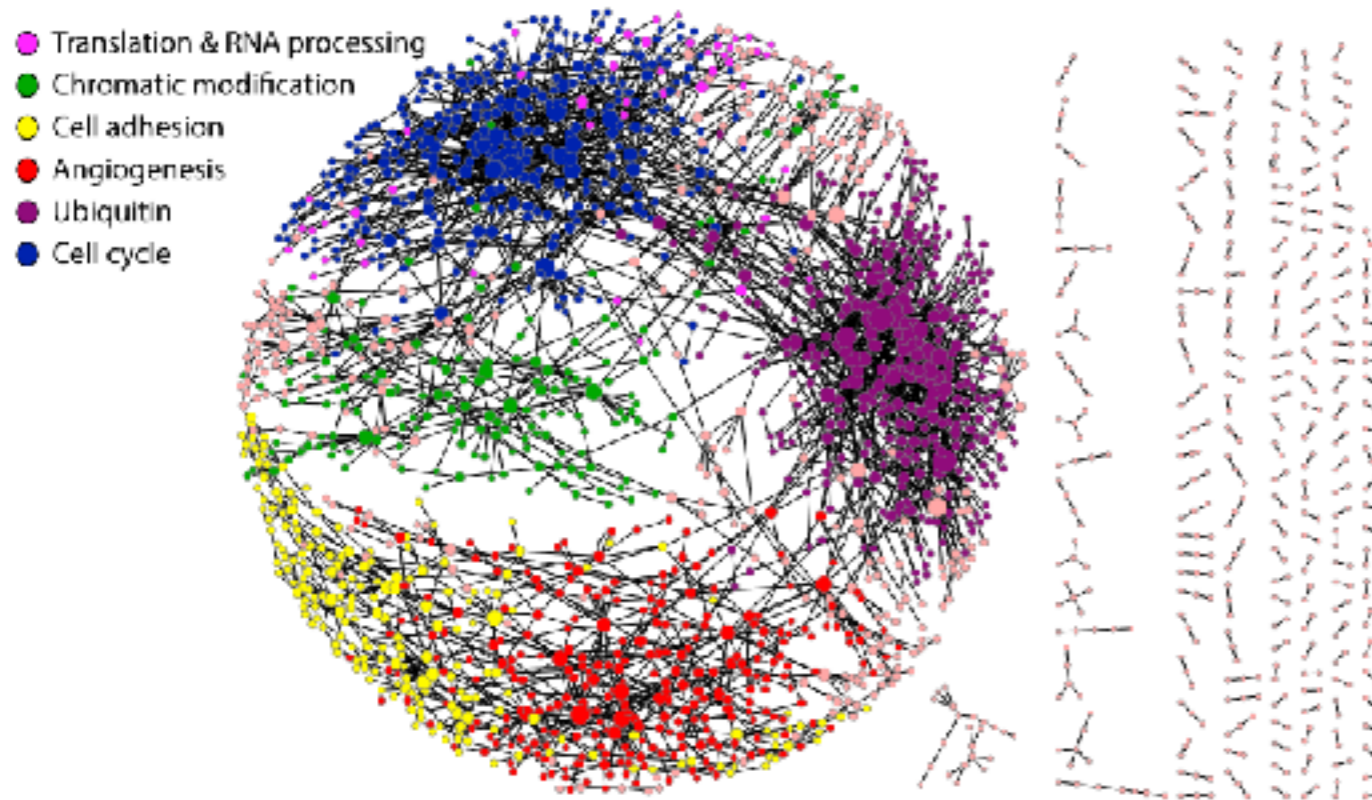


- **Adverse drug reaction prediction** (Wang et al. Drug-induced adverse events prediction with the LINCS L1000 data)
- **Drug target identification** (Xia et al. Target Predictions using LINCS Data)
- **Expression signature comparison** (Xiao et al. SigMat: A Classification Scheme for Gene Signature Matching)
- **Drug response prediction** (Lu et al. Drug-induced cell viability prediction from LINCS-L1000 through WRFEN-XGBoost algorithm)

We don't have so many “drugs”

- Discovery new drug?
 - Often not in the scope of precision medicine
 - New patient cannot wait for a new drug
- Drug repurposing
 - Drug A, which is used to treat disease X, is later used to treat disease Y
 - Well-documented side effects and less restriction from FDA
- Drug combination
 - Drug A is not effective. Drug B is not effective. Drug A and B used together is effective.
- Personalized dosage
 - Widely used in clinics. Use genomics data to determine dosage (regression).

Synthetic lethality: Gene A **OR** Gene B

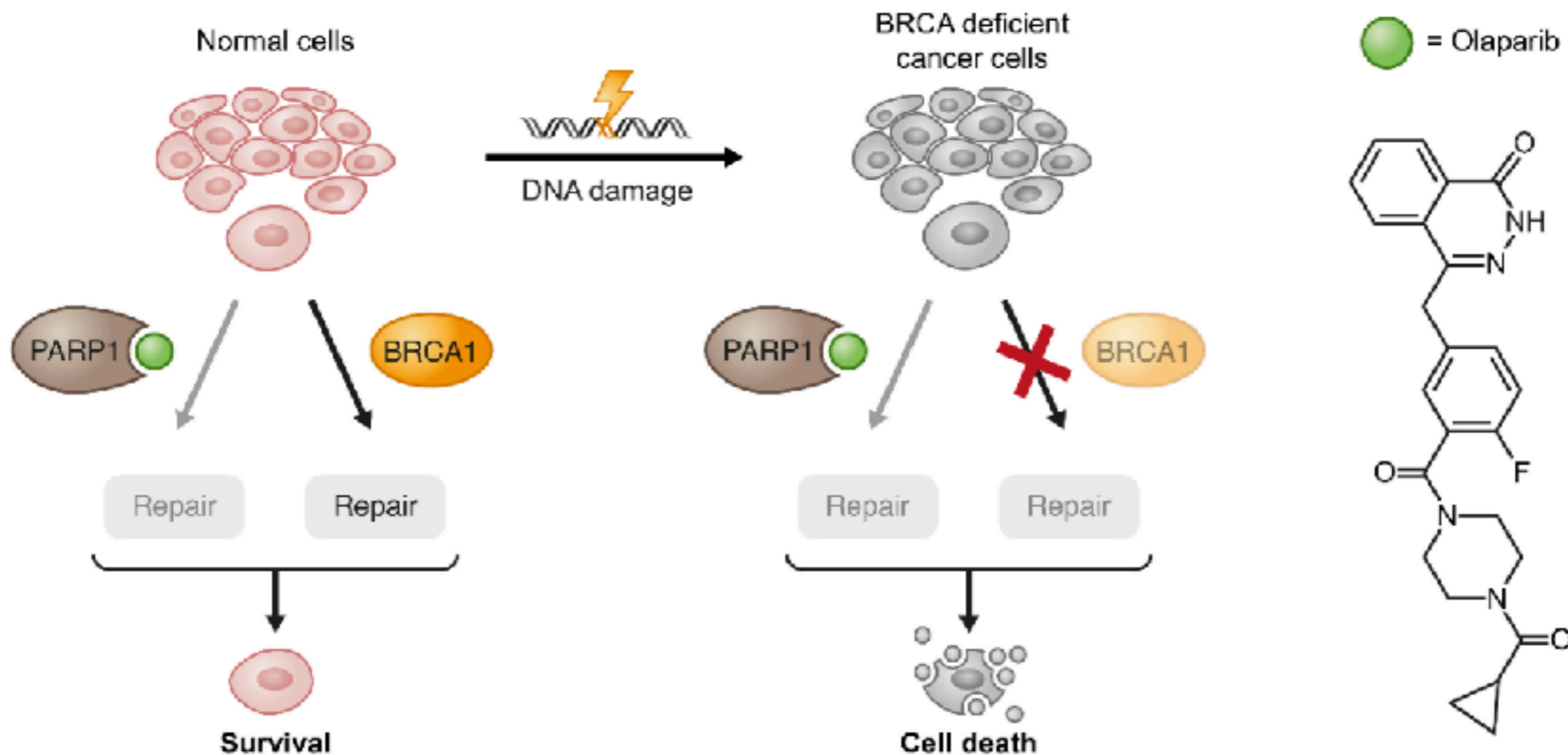


Question: how to leverage SL in drug combination discovery?

Drug combination therapy

- **Breast cancer**
 - an alkylating agent (cyclophosphamide) and antimetabolites (methotrexate and 5-fluorouracil)
- **Anti-HIV cocktail**
 - Use of three or more antiretroviral medicines
- We don't have so many single drug candidate
- Drug combinations ($k \geq 2$) offer us more treatment plans

Drug treatment based on synthetic lethality



Goal: We want to make normal cells survive and kill cancer cells (BRCA deficient cancer cells)

Prior knowledge: PARP1 (off) + BRCA1 (off) → cell death

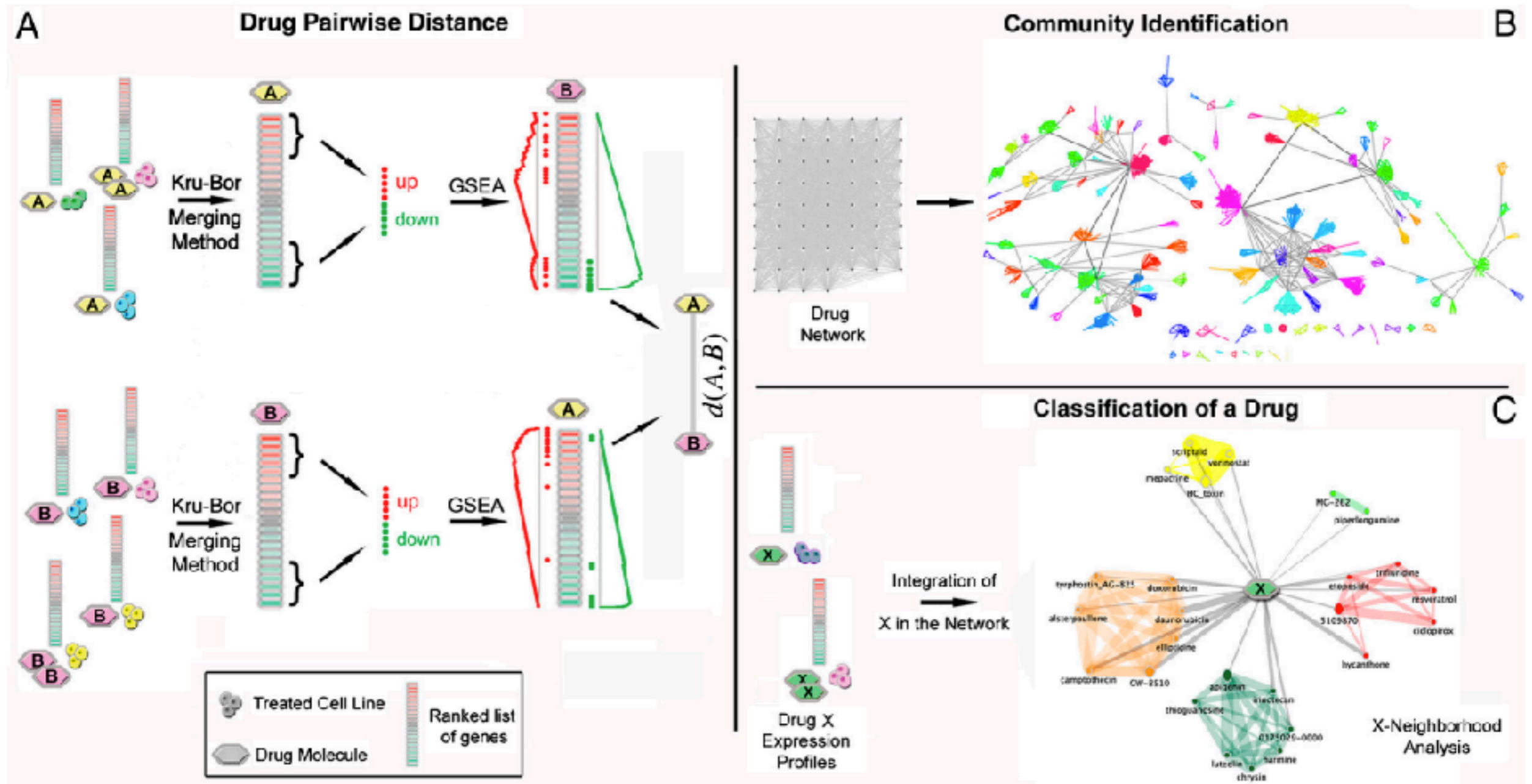
Solution: Turn off PARP1 using Olaparib

Results:

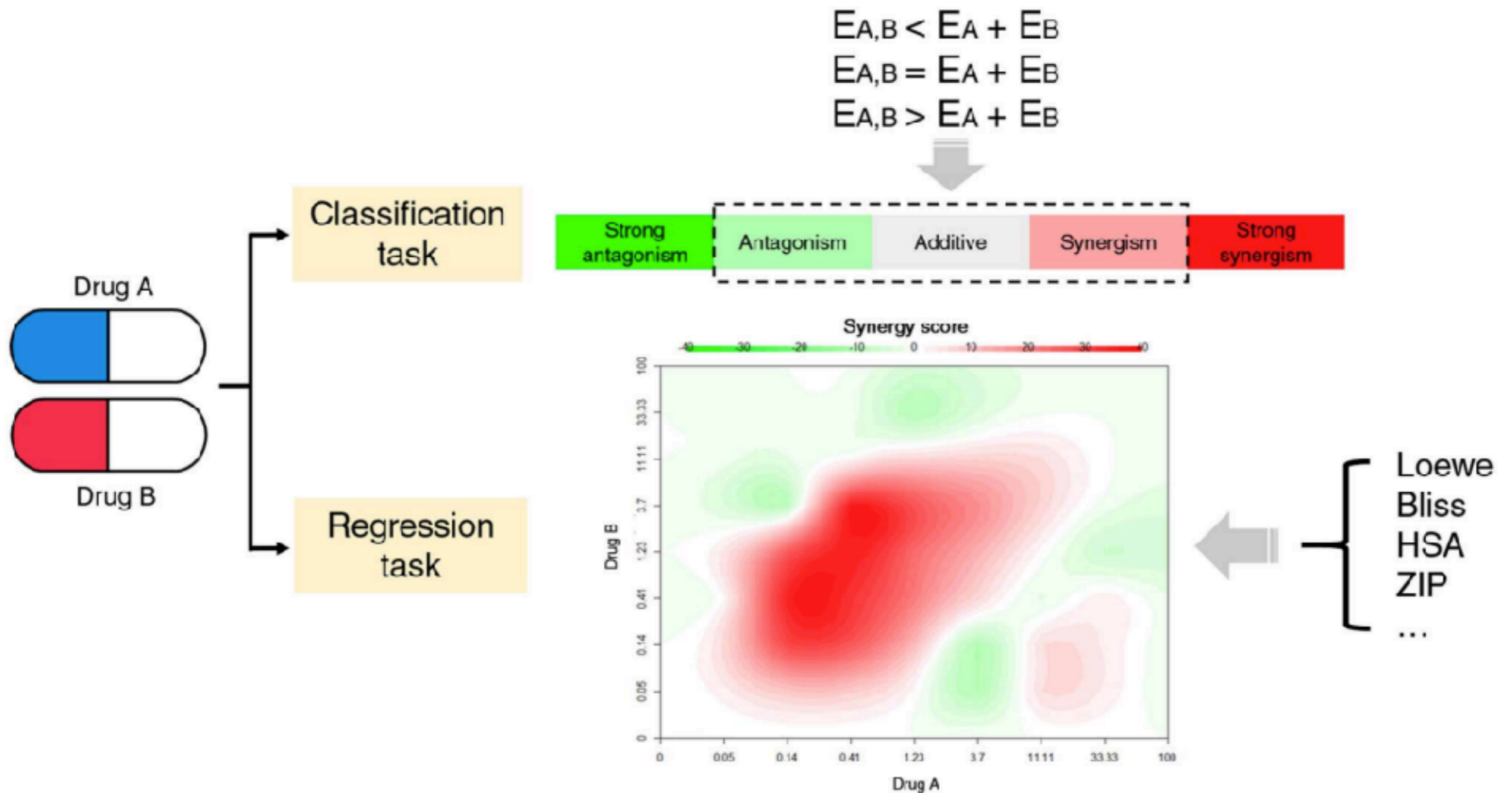
- Normal cells: PARP1 (off) + BRCA1 (on) → cell survive
- Cancer cells: PARP1 (off) + BRCA1 (off) → cell death

Use gene expression after treatment

Drugs target on similar proteins or have similar Mode of Actions have similar (after treatment) expression.

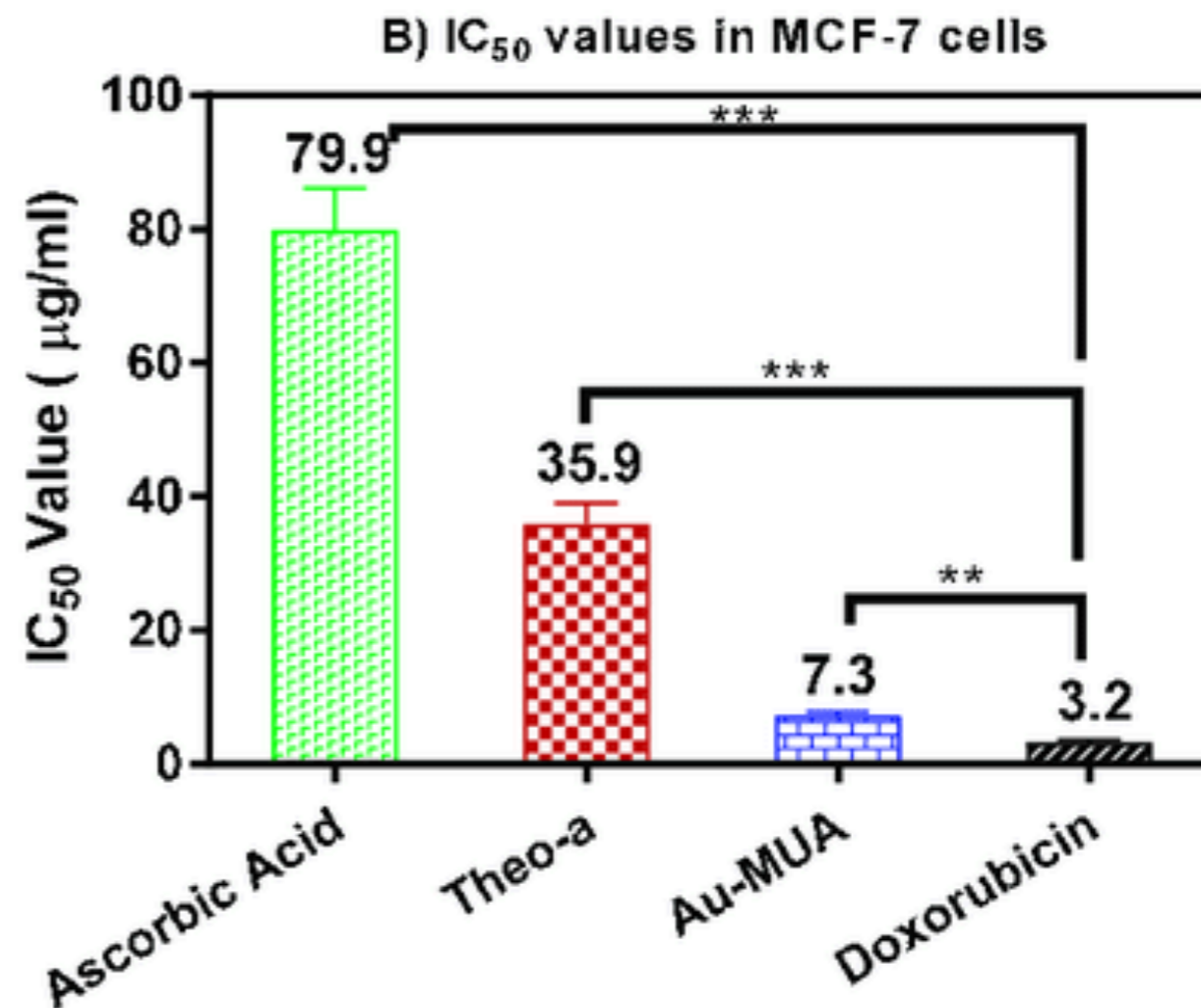
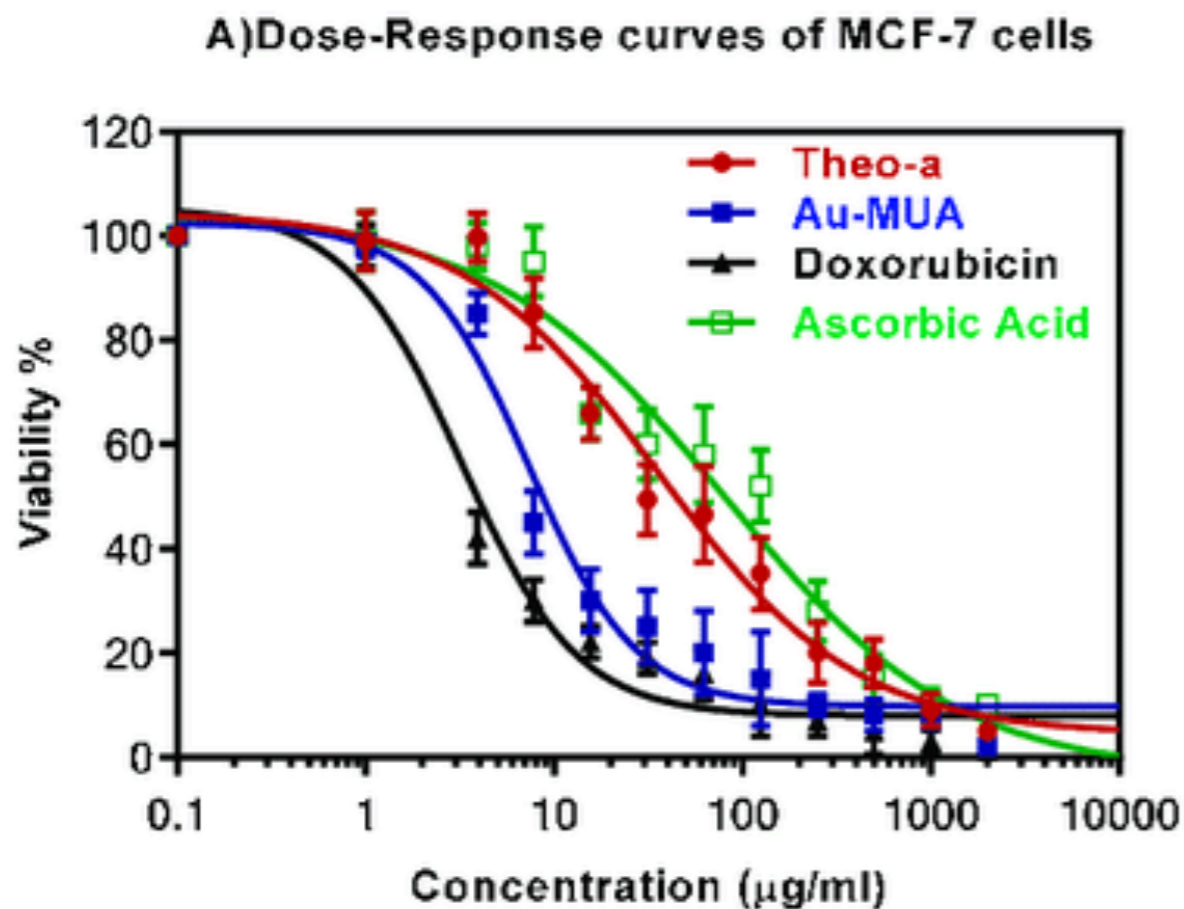


Drug combination prediction



$E(A)$ is the efficacy of using drug A (e.g., IC50)

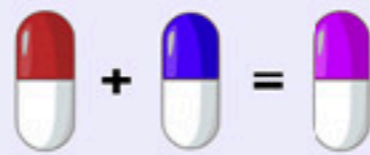
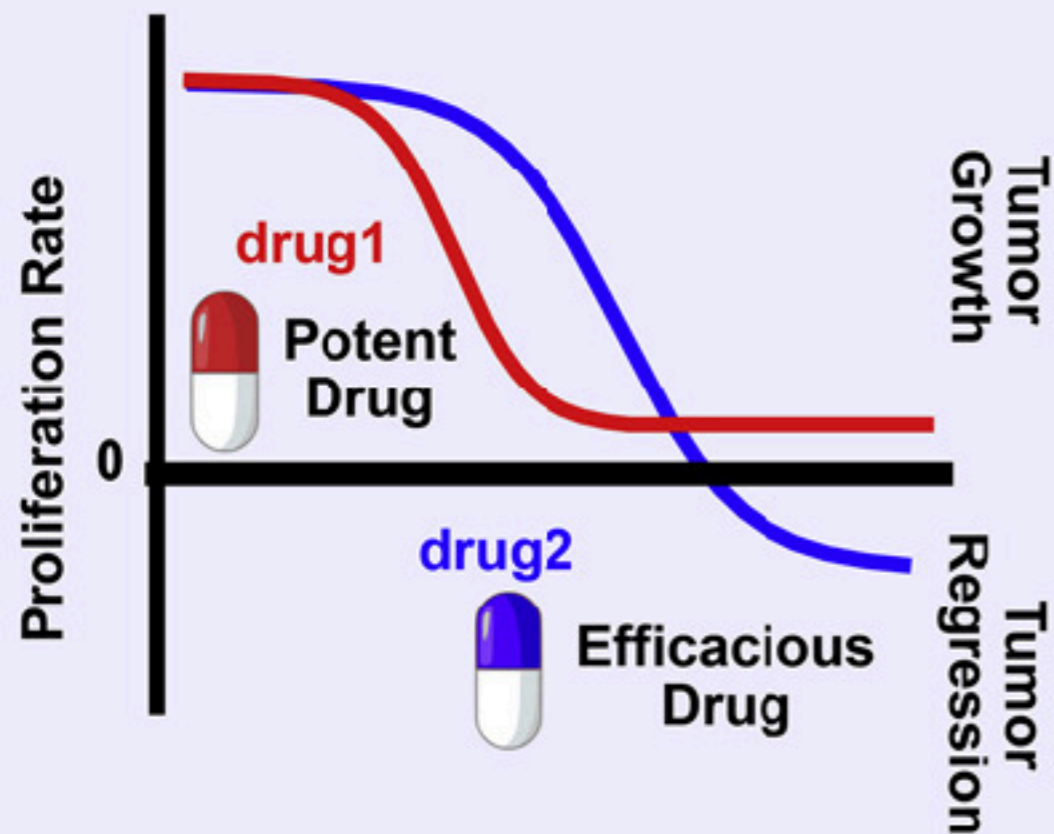
Dose-response curve



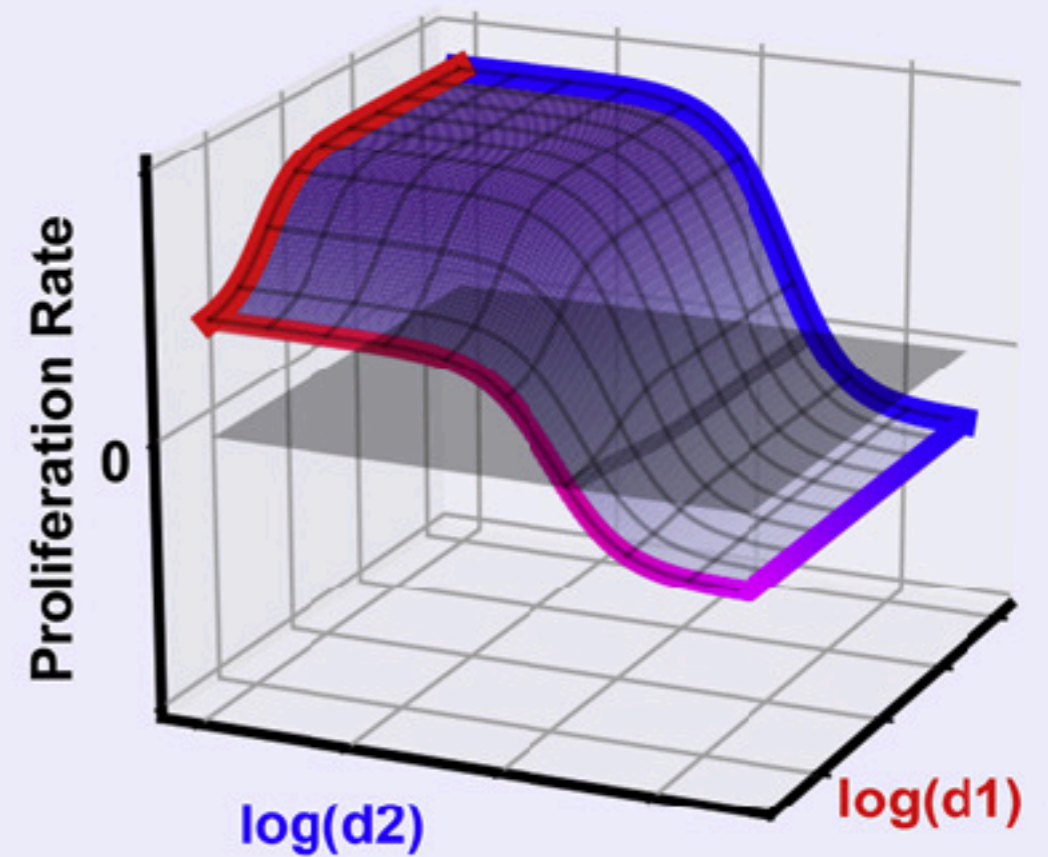
$\text{IC}(50)$: concentration of 50% viability

Drug combination dose-response curve

1D Dose-Response Curve



2D Dose-Response Surface



DeepSynergy: deep learning-based drug synergistic prediction

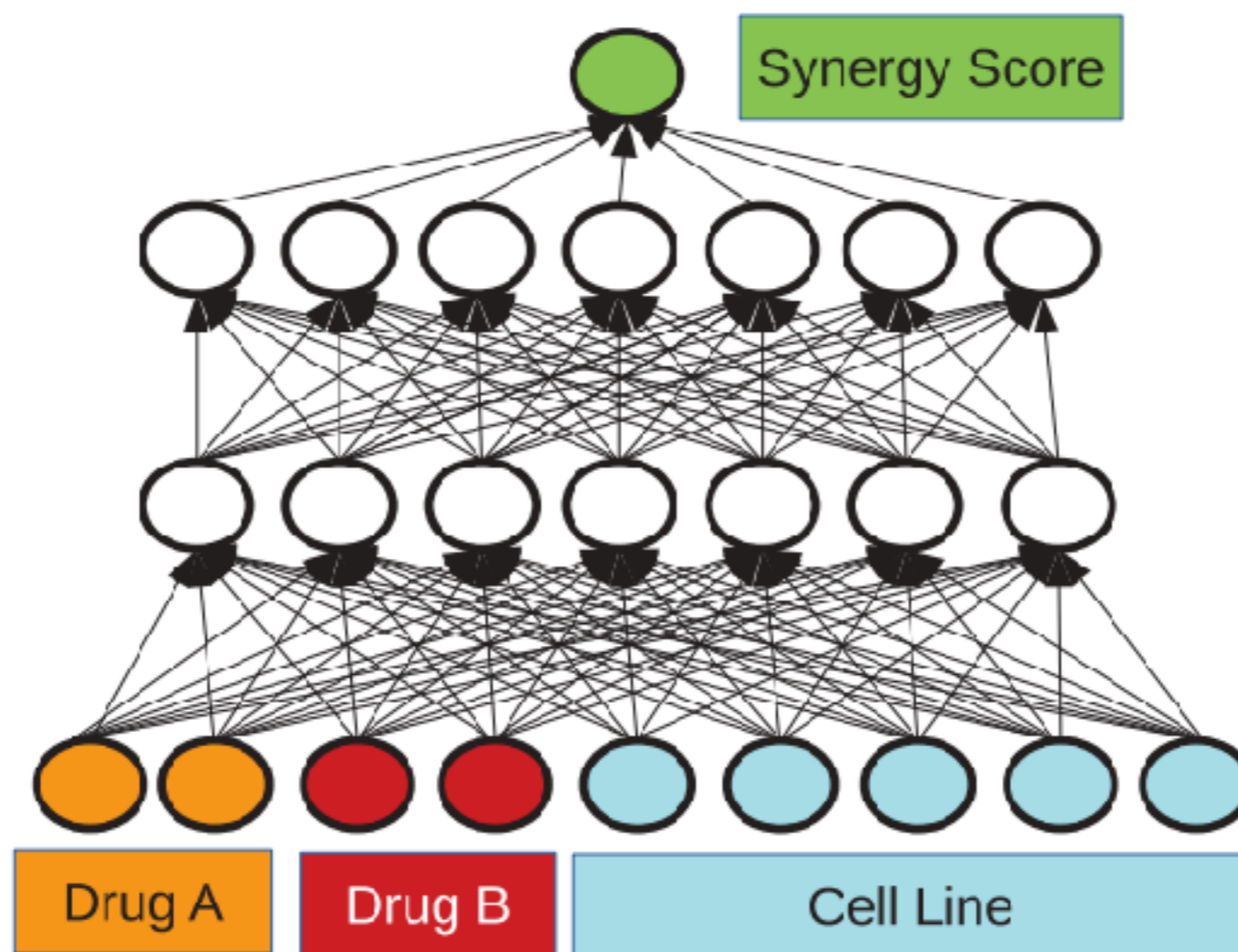
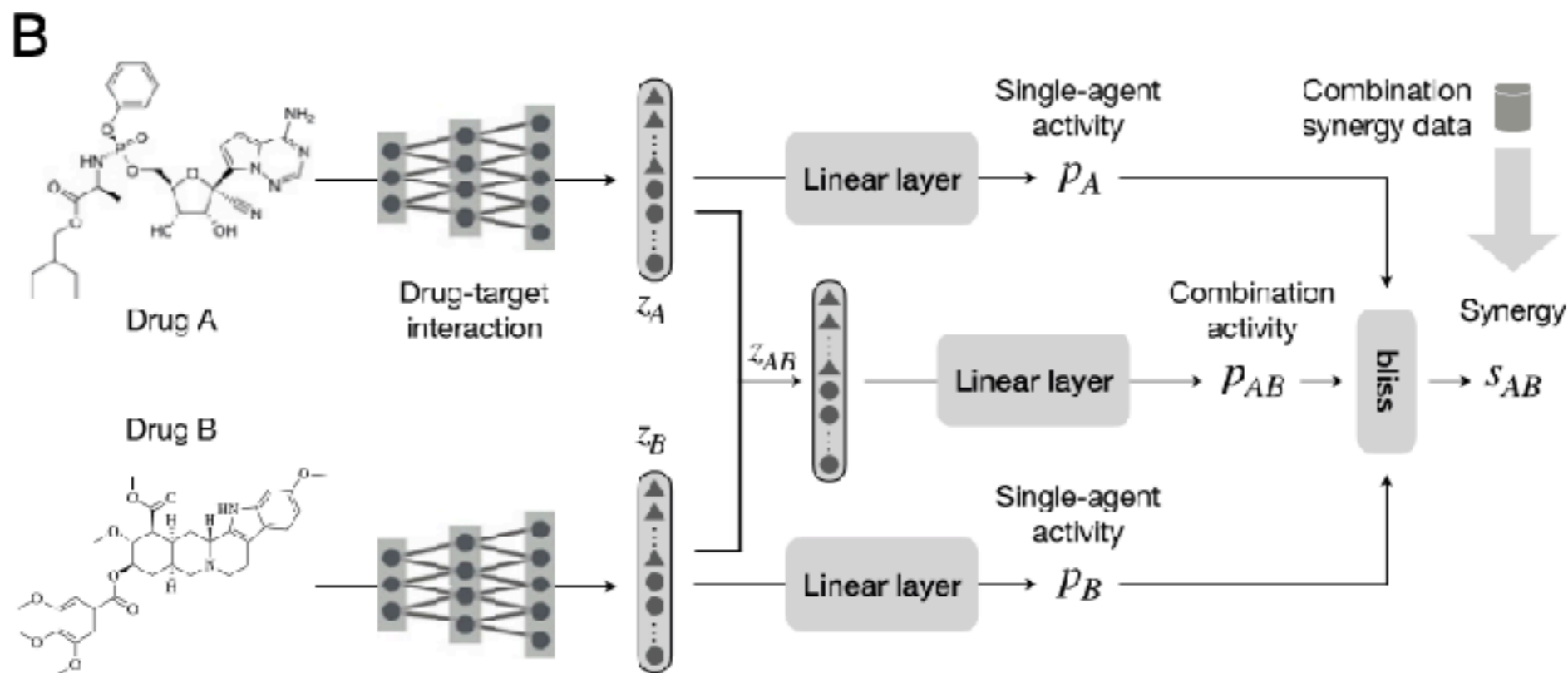
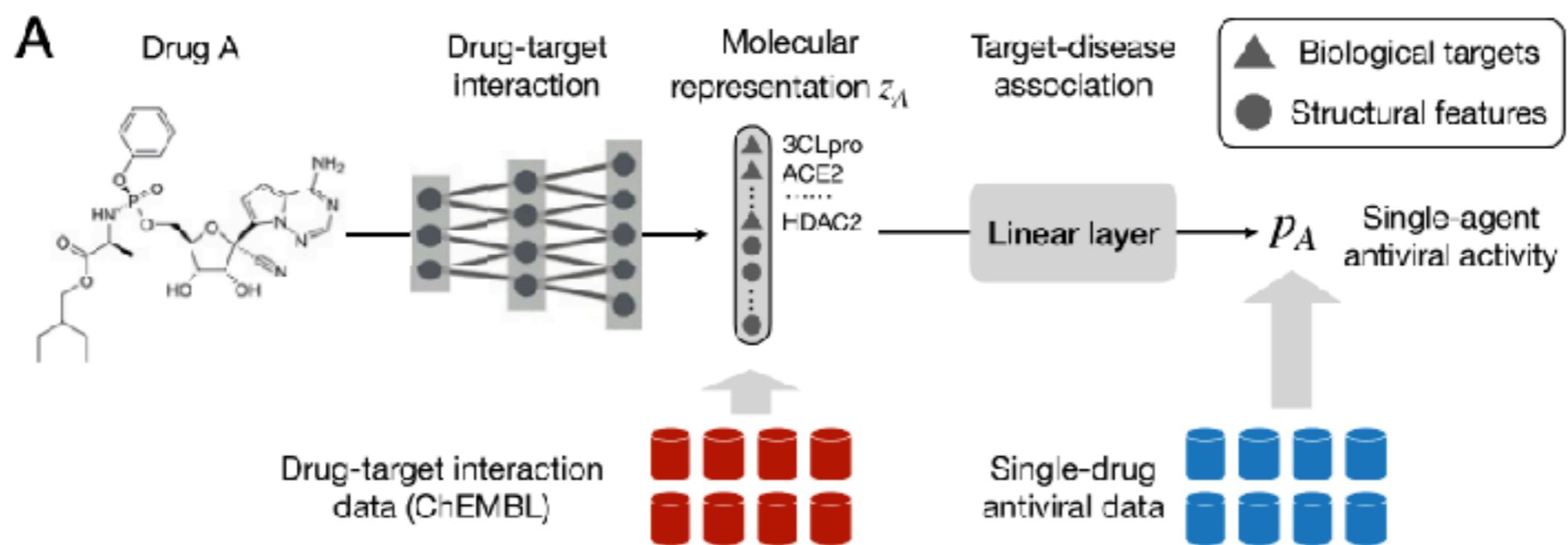


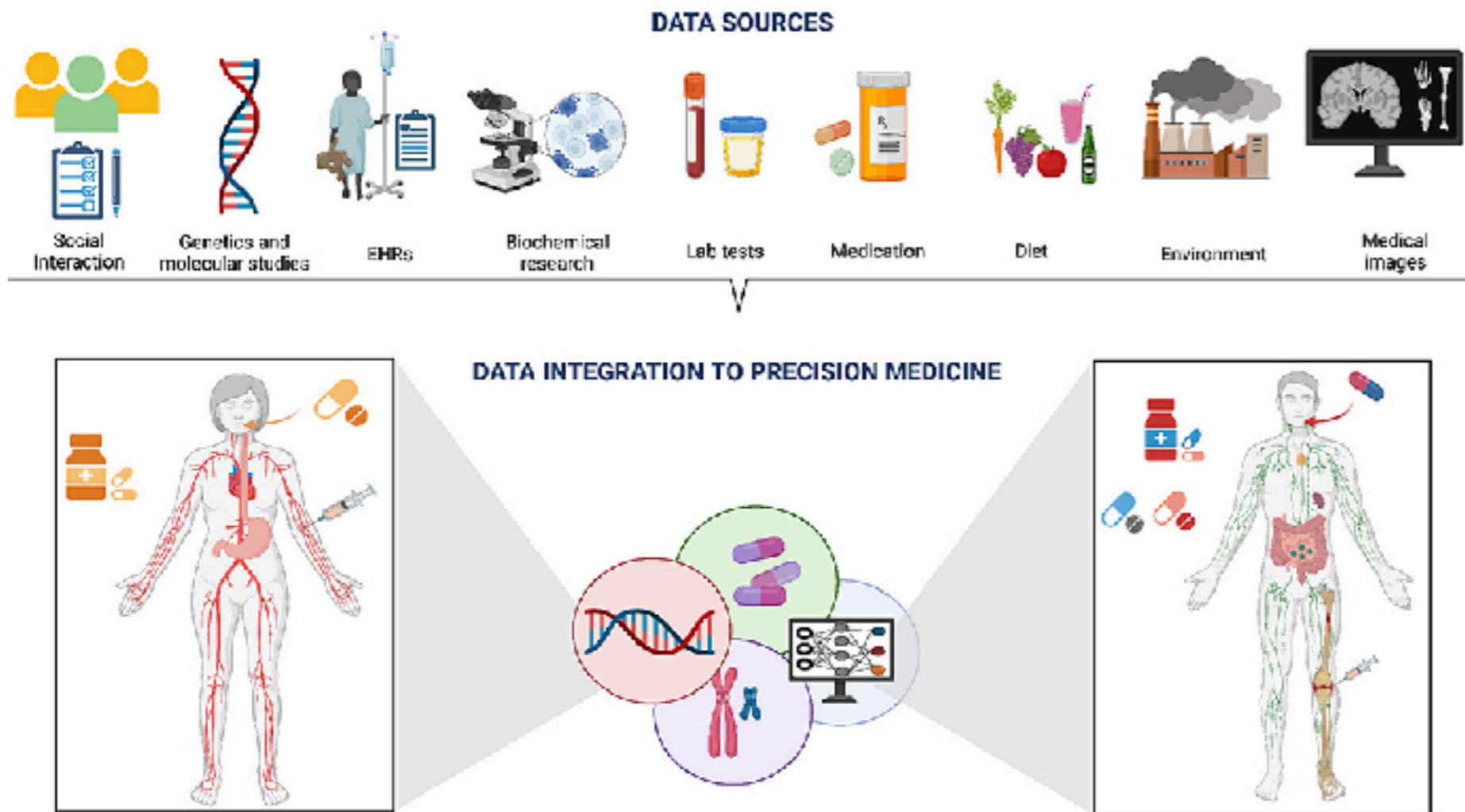
Table 3. Performance metrics for the classification task

Performance Metric	ROC AUC	PR AUC	ACC	BACC	PREC	TPR	TNR	Kappa
Deep Neural Networks	0.90 ± 0.03	0.59 ± 0.06	0.92 ± 0.03	0.76 ± 0.03	0.56 ± 0.11	0.57 ± 0.09	0.95 ± 0.03	0.51 ± 0.04
Gradient Boosting Machines	0.89 ± 0.02	0.59 ± 0.04	0.87 ± 0.01	0.80 ± 0.03	0.38 ± 0.04	0.71 ± 0.05	0.89 ± 0.01	0.43 ± 0.03
Random Forests	0.87 ± 0.02	0.55 ± 0.04	0.92 ± 0.01	0.73 ± 0.04	0.57 ± 0.04	0.49 ± 0.08	0.96 ± 0.01	0.48 ± 0.04
Support Vector Machines	0.81 ± 0.04	0.42 ± 0.08	0.76 ± 0.06	0.73 ± 0.03	0.23 ± 0.04	0.69 ± 0.08	0.77 ± 0.07	0.24 ± 0.05
Elastic Nets	0.78 ± 0.04	0.34 ± 0.10	0.75 ± 0.05	0.71 ± 0.02	0.21 ± 0.03	0.65 ± 0.07	0.76 ± 0.06	0.22 ± 0.03
Baseline (Median Polish)	0.77 ± 0.04	0.32 ± 0.09	0.76 ± 0.04	0.70 ± 0.03	0.22 ± 0.03	0.62 ± 0.06	0.78 ± 0.04	0.22 ± 0.04

Drug combinations for treating COVID-19



Two key problems



- How to cluster patients
- We don't have so many "drugs"

How to cluster patients

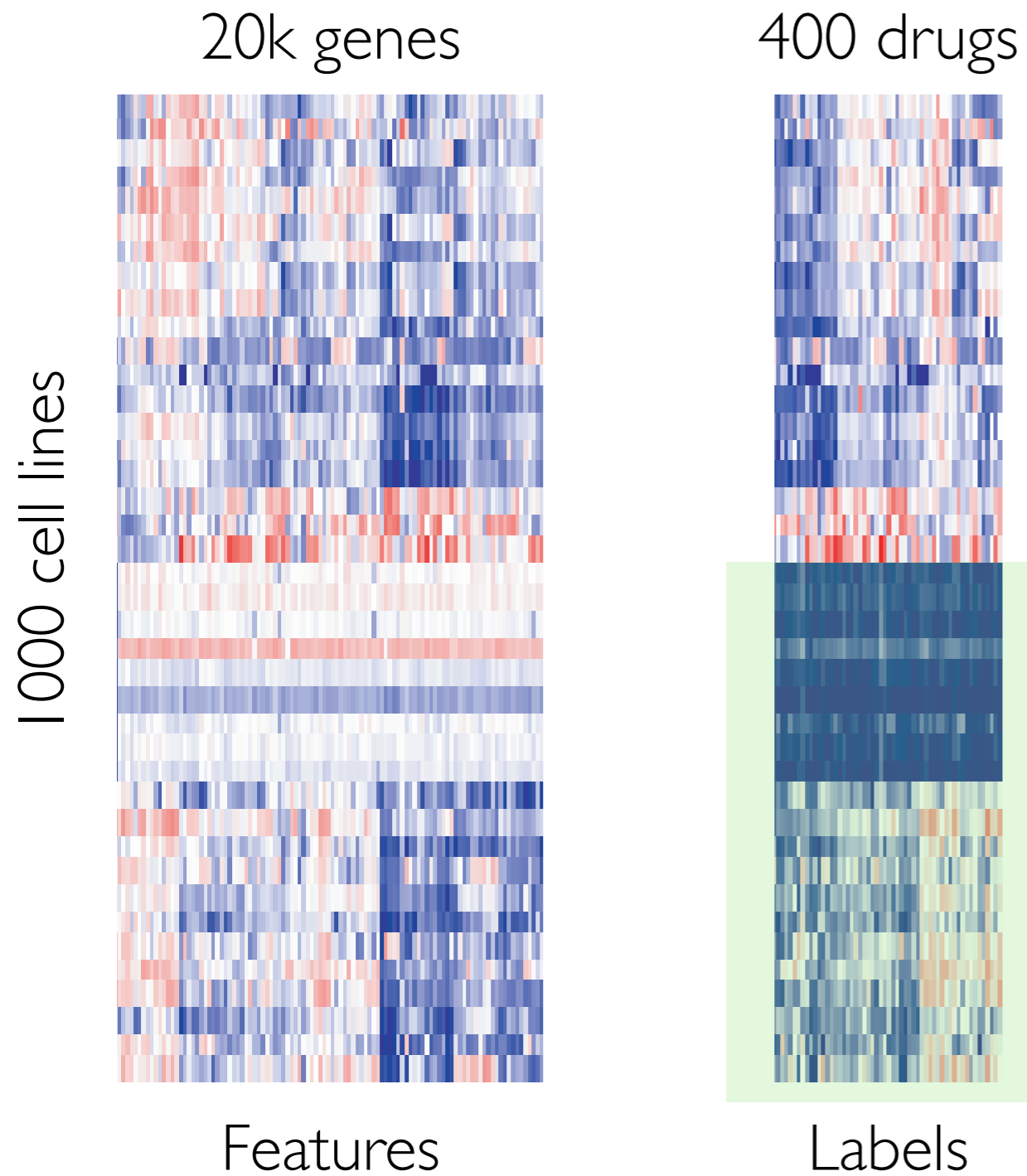


- **Patient clustering = data integration**
 - Find a “signature” vector for each patient
 - Signature is integrated from different data sources
- **Heterogeneous data integration**
 - General challenges: Heterogeneous, missing values, noise, privacy
- **Precision medicine specific data integration challenges:**
 - Batch effects (different preprocessing pipelines, sequencing techniques, reference ranges)
 - Unpaired data (some patients only have genomics data, some patients only have EHR data, very few patients have both)

Public biomedical databases

DATA REPOSITORY	WEB LINK	DISEASE	TYPES OF MULTI-OMICS DATA AVAILABLE
The Cancer Genome Atlas (TCGA)	https://cancergenome.nih.gov/	Cancer	RNA-Seq, DNA-Seq, miRNA-Seq, SNV, CNV, DNA methylation, and RPPA
Clinical Proteomic Tumor Analysis Consortium (CPTAC)	https://cptac-data-portal.georgetown.edu/cptacPublic/	Cancer	Proteomics data corresponding to TCGA cohorts
International Cancer Genomics Consortium (ICGC)	https://icgc.org/	Cancer	Whole genome sequencing, genomic variations data (somatic and germline mutation)
Cancer Cell Line Encyclopedia (CCLE)	https://portals.broadinstitute.org/ccle	Cancer cell line	Gene expression, copy number, and sequencing data; pharmacological profiles of 24 anticancer drugs
Molecular Taxonomy of Breast Cancer International Consortium (METABRIC)	http://molonc.bccrc.ca/aparicio-lab/research/metabric/	Breast cancer	Clinical traits, gene expression, SNP, and CNV
TARGET	https://ocg.cancer.gov/programs/target	Pediatric cancers	Gene expression, miRNA expression, copy number, and sequencing data
Omics Discovery Index	https://www.omicsdi.org	Consolidated data sets from 11 repositories in a uniform framework	Genomics, transcriptomics, proteomics, and metabolomics

Personalized drug response prediction: multi-label regression problem

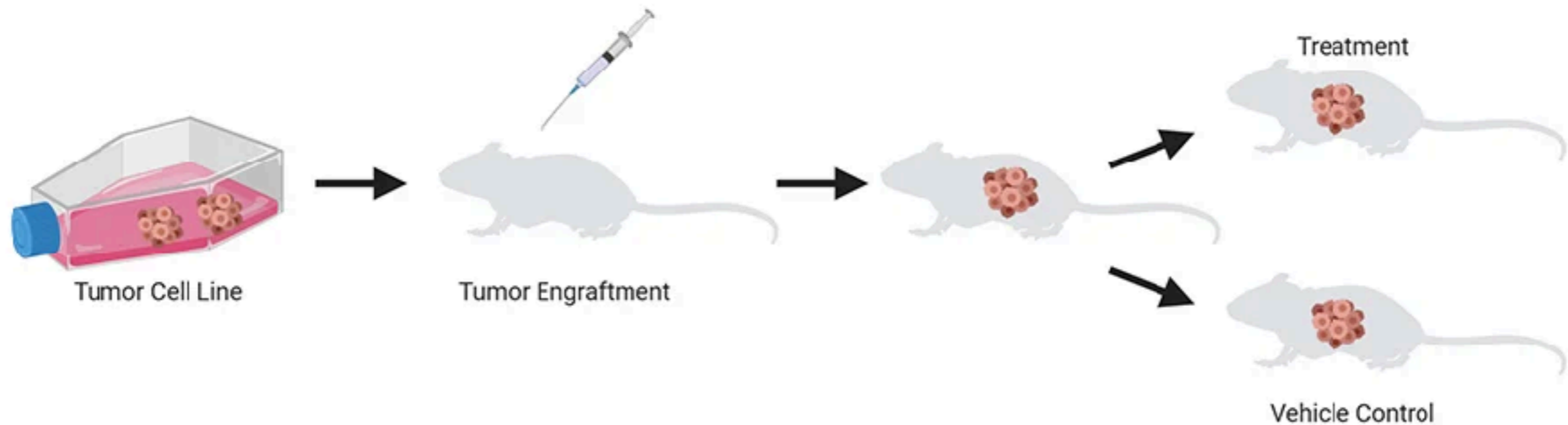


CCLC data: ~1000 cell lines, 20k genes, 400 drugs

Three settings

- Test patients: no drugs are observed for this patient
- Test drugs: no patients are observed for this drug
- Test <patient, drug> pairs

Cell line, xenograft, tumor, patient



- Cell line is a “copy” of a patient. We cannot test one patient with many drugs. But we can copy a cell line many times.
- Cell line is cheaper than xenograft. Xenograft is cheaper than patient data
- Xenograft data: Gao et al. High-throughput screening using patient-derived tumor xenografts to predict clinical trial drug response
- TCGA has some patient data
- ML question: how to integrate cell line data, xenograft data and patient data

Batch effects: inconsistency or consistency?

ANALYSIS

doi:10.1038/nature12831

Inconsistency in large pharmacogenomic studies

Benjamin Haibe-Kains^{1,2}, Nehme El-Hachem¹, Nicolai Juul Birkbak³, Andrew C. Jin⁴, Andrew H. Beck^{4*}, Hugo J. W. L. Aerts^{5,6,7*} & John Quackenbush^{5,8*}

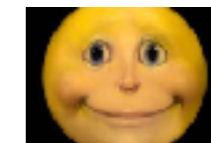
Two large-scale pharmacogenomic studies were published recently in this journal. Genomic data are well correlated between studies; however, the measured drug response data are highly discordant. Although the source of inconsistencies remains uncertain, it has potential implications for using these outcome measures to assess gene-drug associations or select potential anticancer drugs on the basis of their reported results.

ANALYSIS

doi:10.1038/nature15736

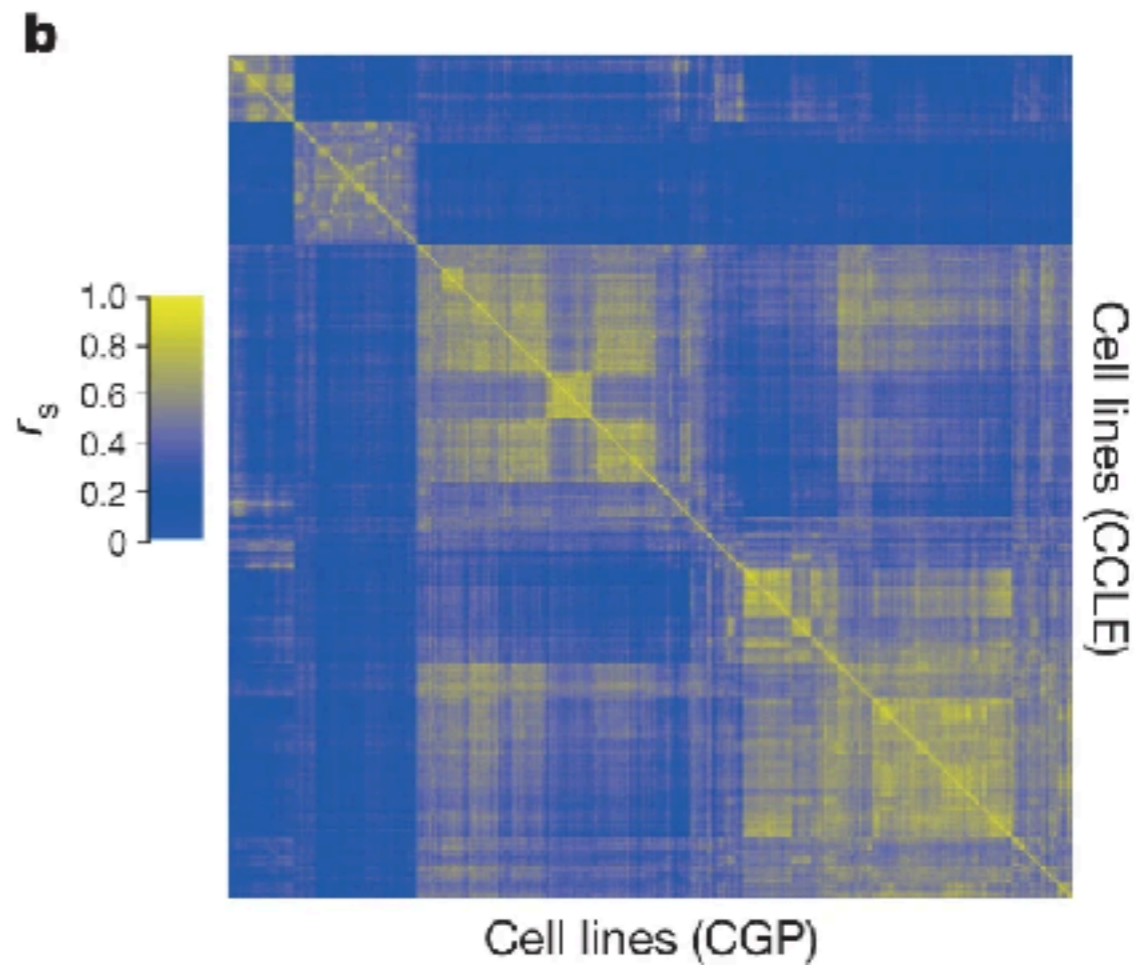
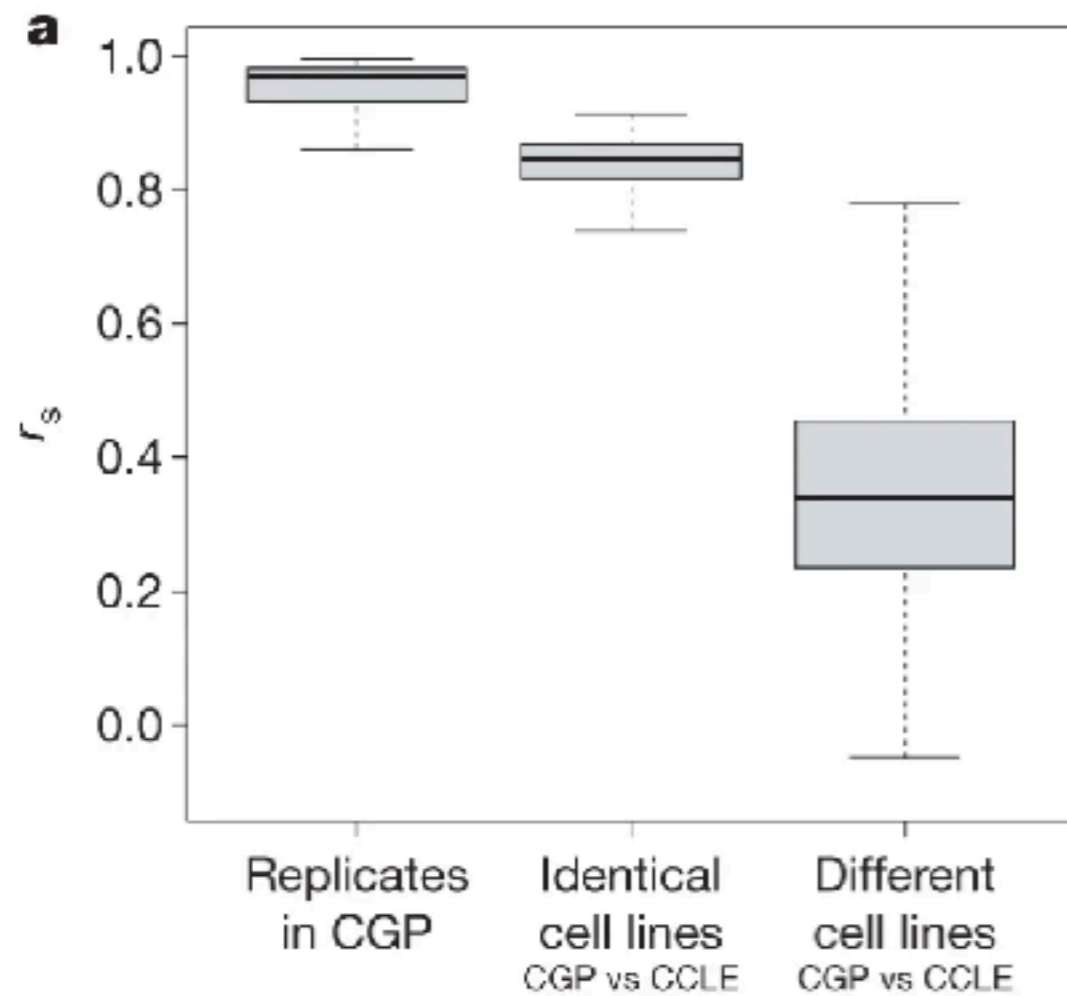
Pharmacogenomic agreement between two cancer cell line data sets

The Cancer Cell Line Encyclopedia and Genomics of Drug Sensitivity in Cancer Investigators*



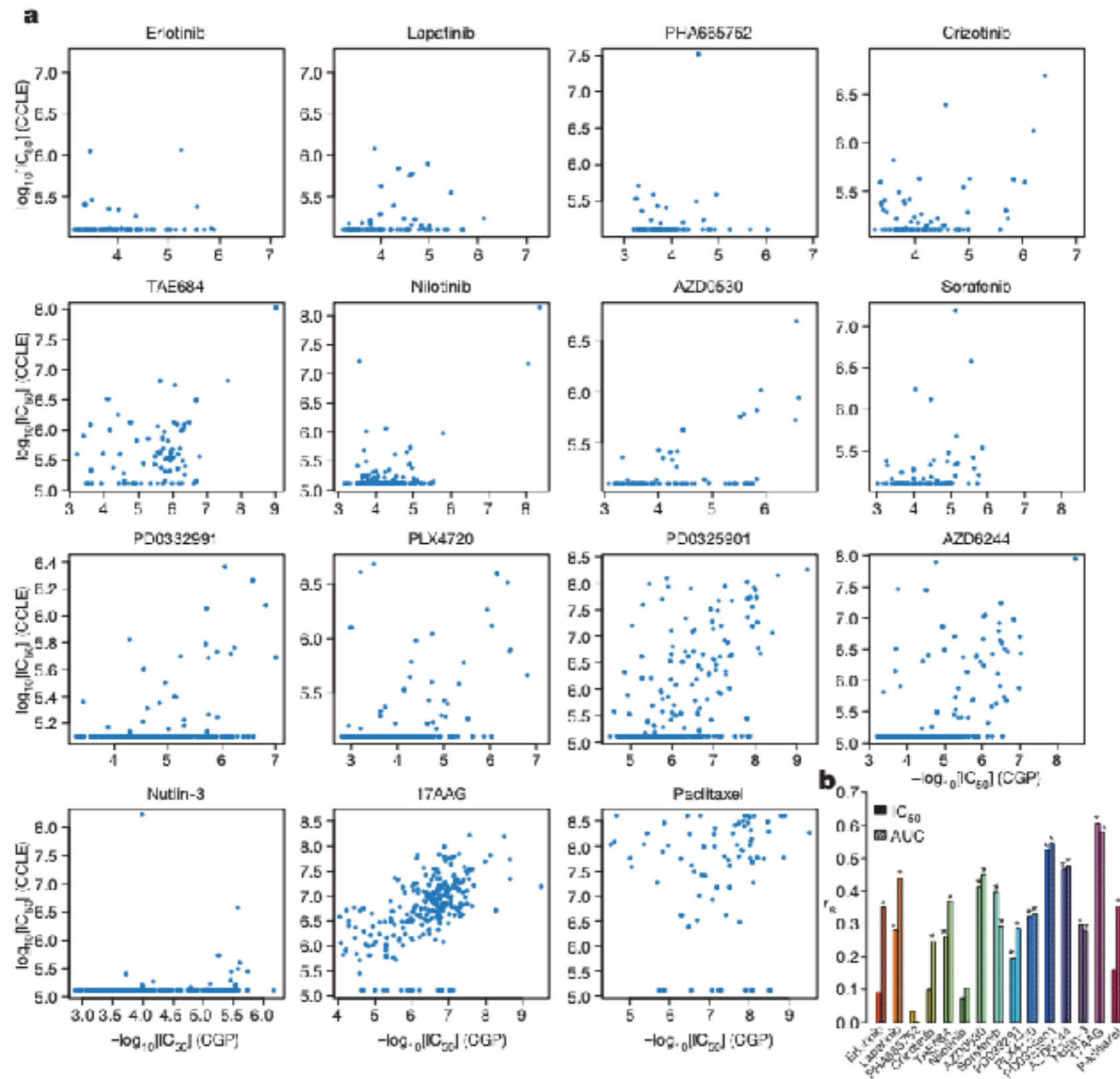
Large cancer cell line collections broadly capture the genomic diversity of human cancers and provide valuable insight into anti-cancer drug response. Here we show substantial agreement and biological consistency between drug sensitivity measurements and their associated genomic predictors from two publicly available large-scale pharmacogenomics resources: The Cancer Cell Line Encyclopedia and the Genomics of Drug Sensitivity in Cancer databases.

Low correlation between drug response data

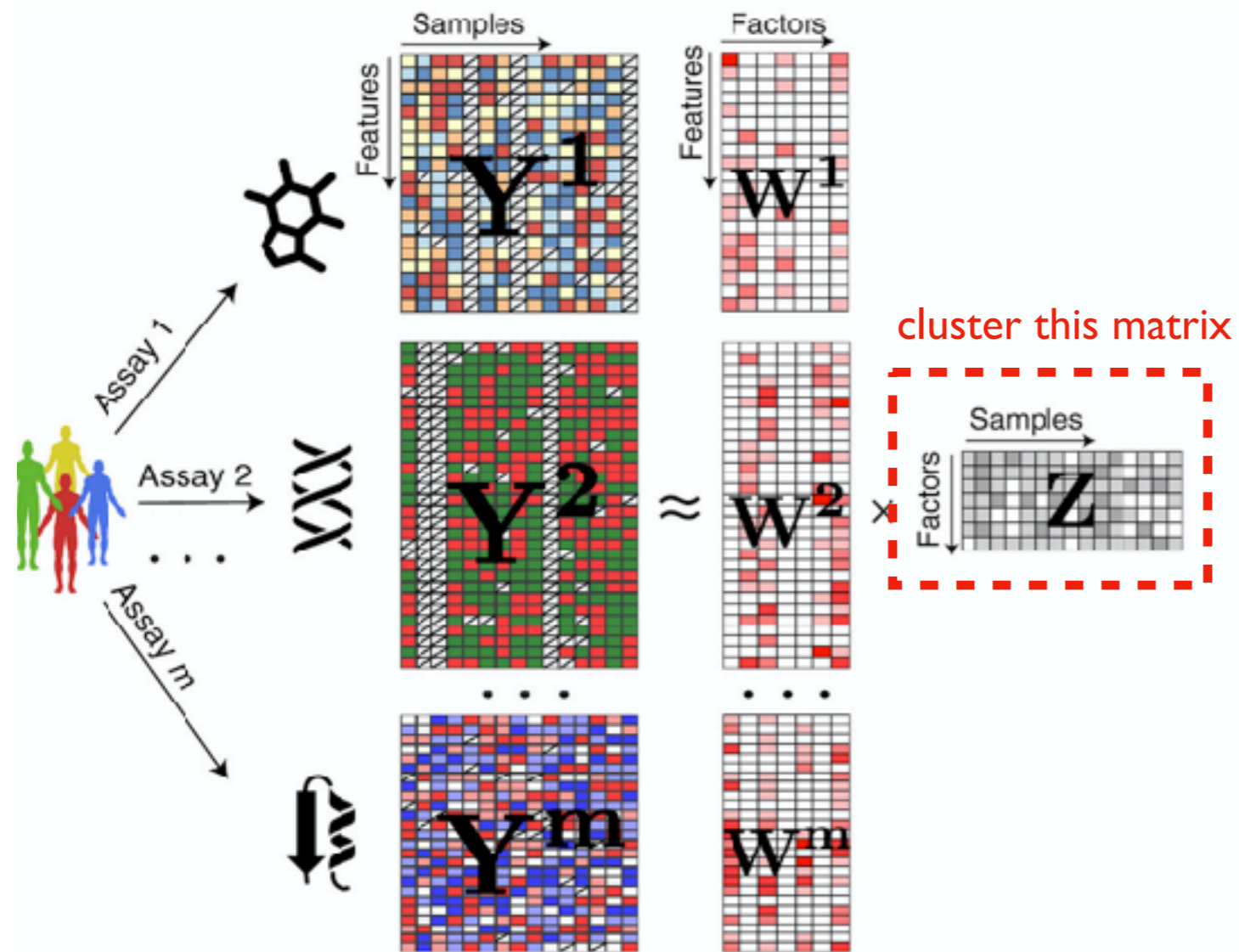


Integrating two datasets

High correlation between gene expression (features)

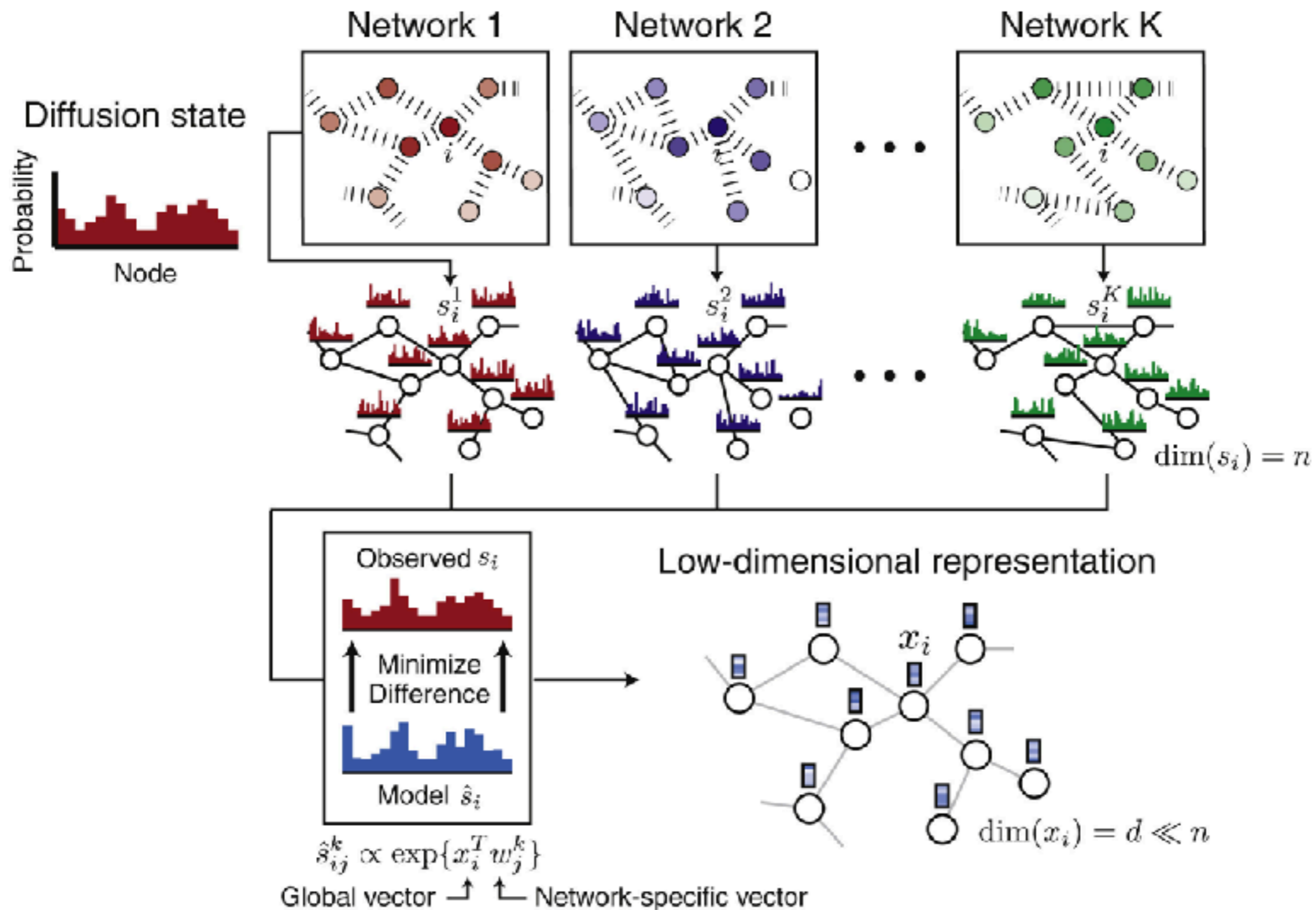


General framework: jointly decompose multiple data matrices

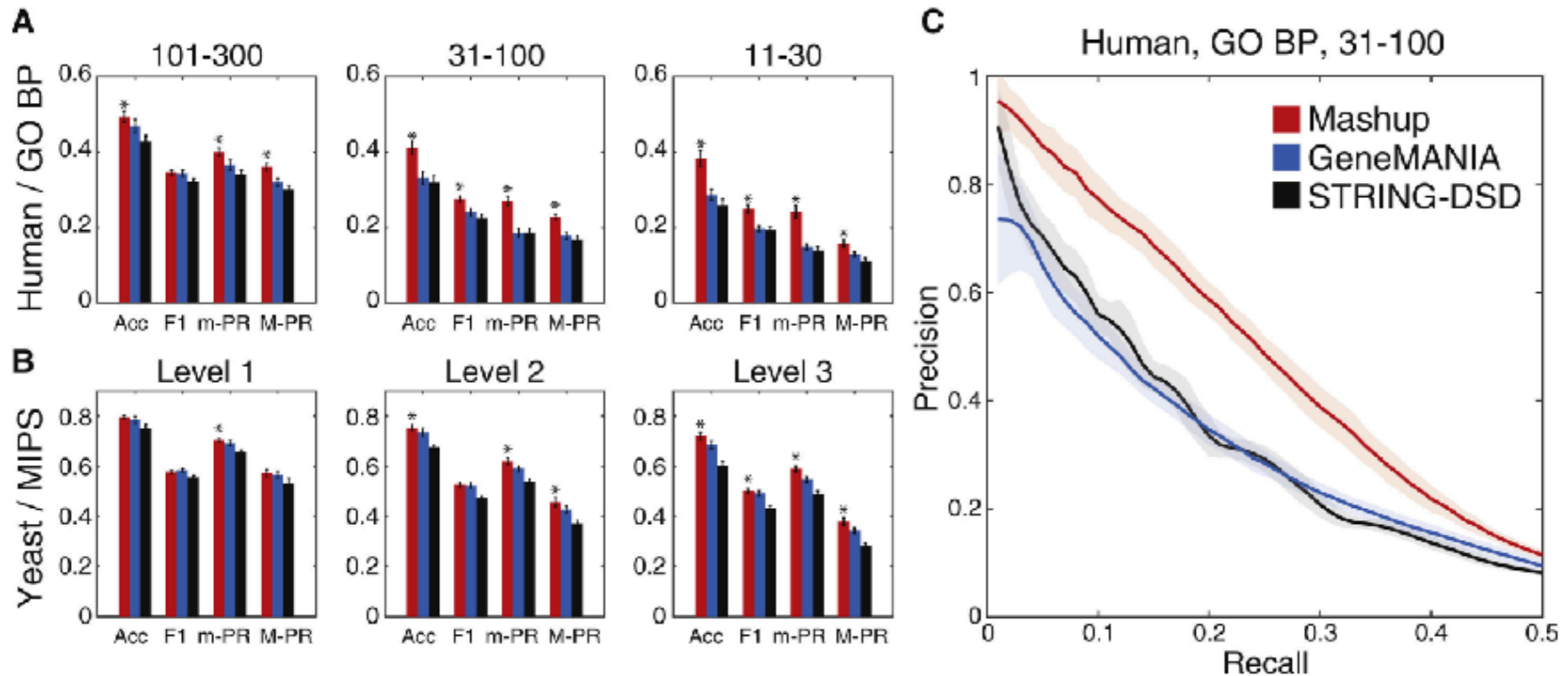


- Key ideas
 - One matrix capture batch effects
 - One matrix capture common patterns
- Detail implementations
 - What distribution?
 - Mutation (Bernoulli)
 - Expression count (Poisson)
 - How to decompose?
 - NMF, SVD, NN, MF

Mashup: integrating multiple networks



Mashup improves protein function prediction

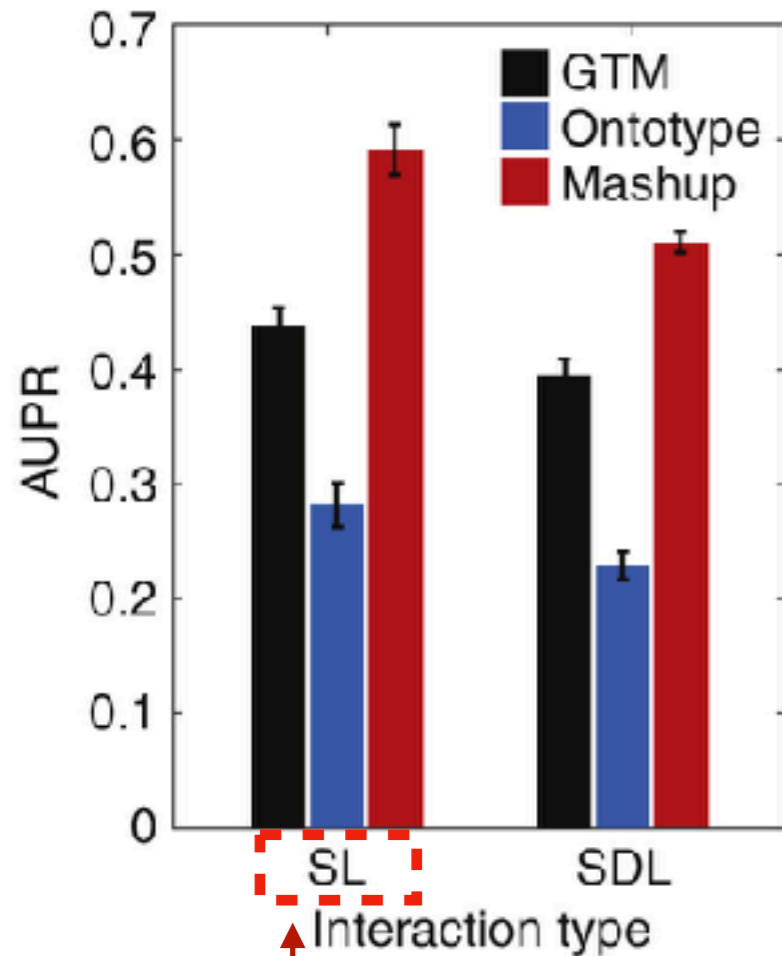


Protein function prediction is a good benchmark for machine learning algorithms because of it is high-quality and has many annotations. It can be used to evaluate:

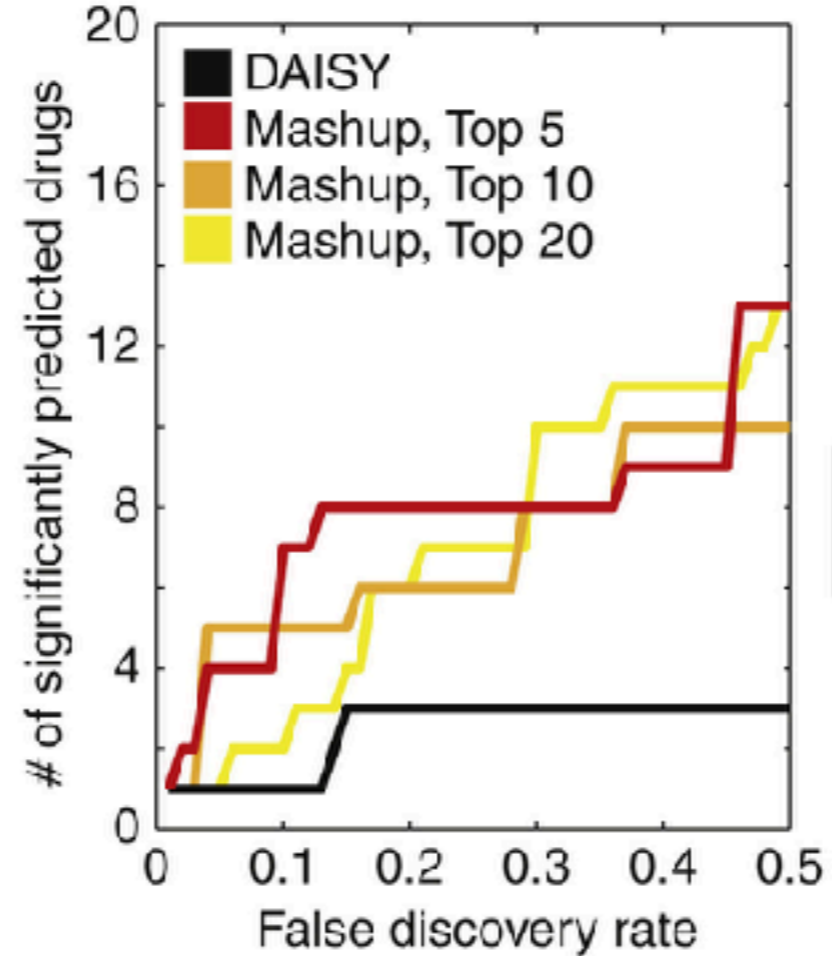
- Network-based approach
- Sequence-based approach
- Few-shot/zero-shot learning

Mashup enables genetic interaction prediction

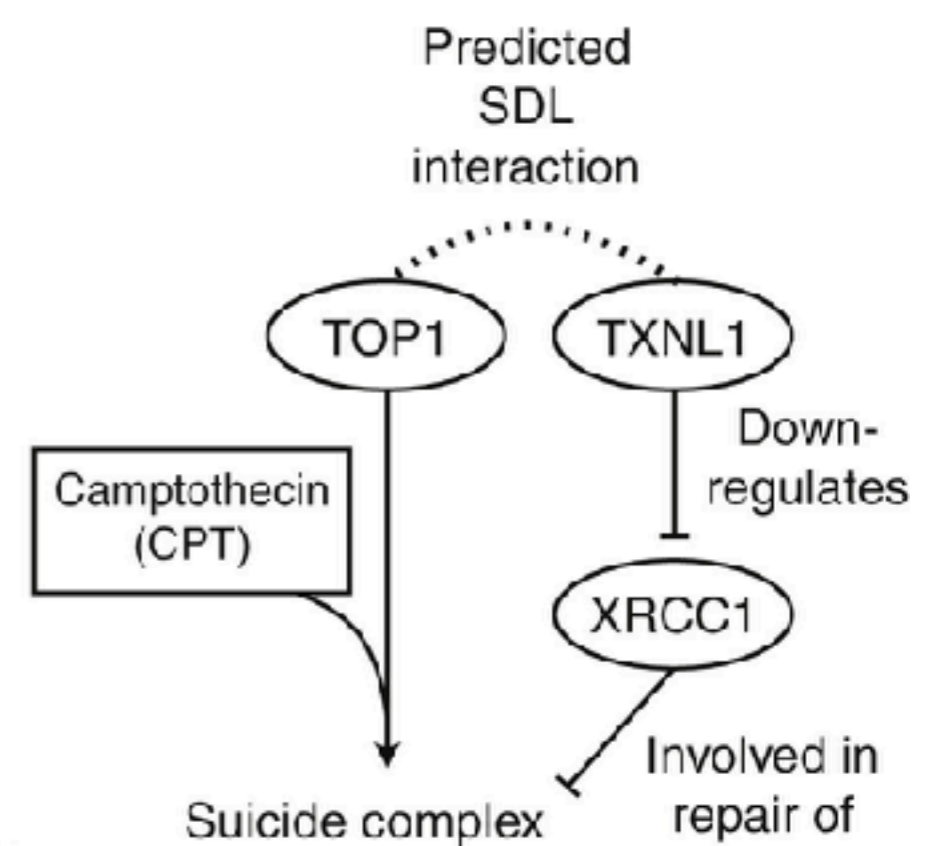
A



B



C

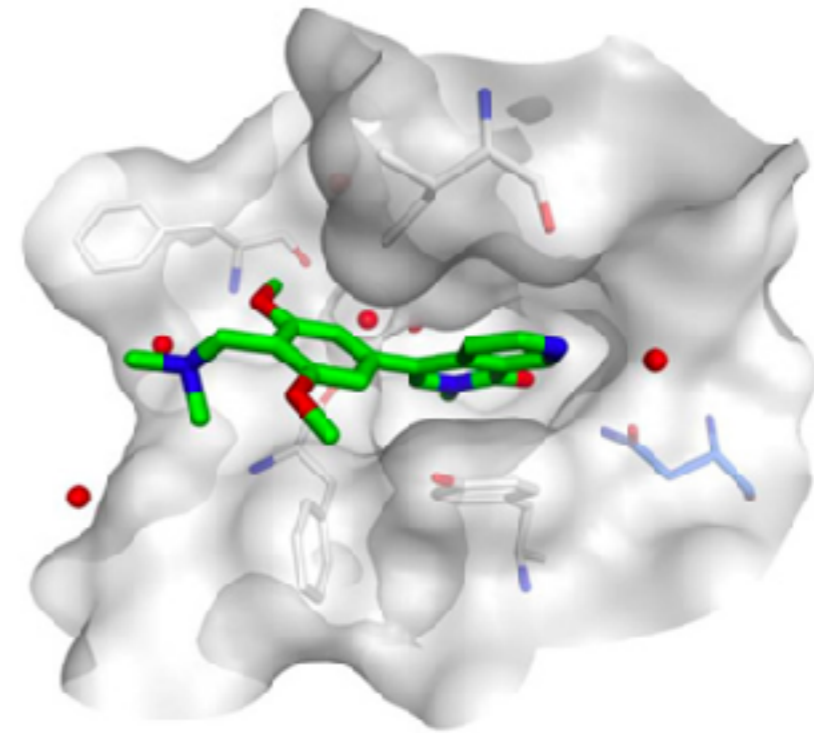
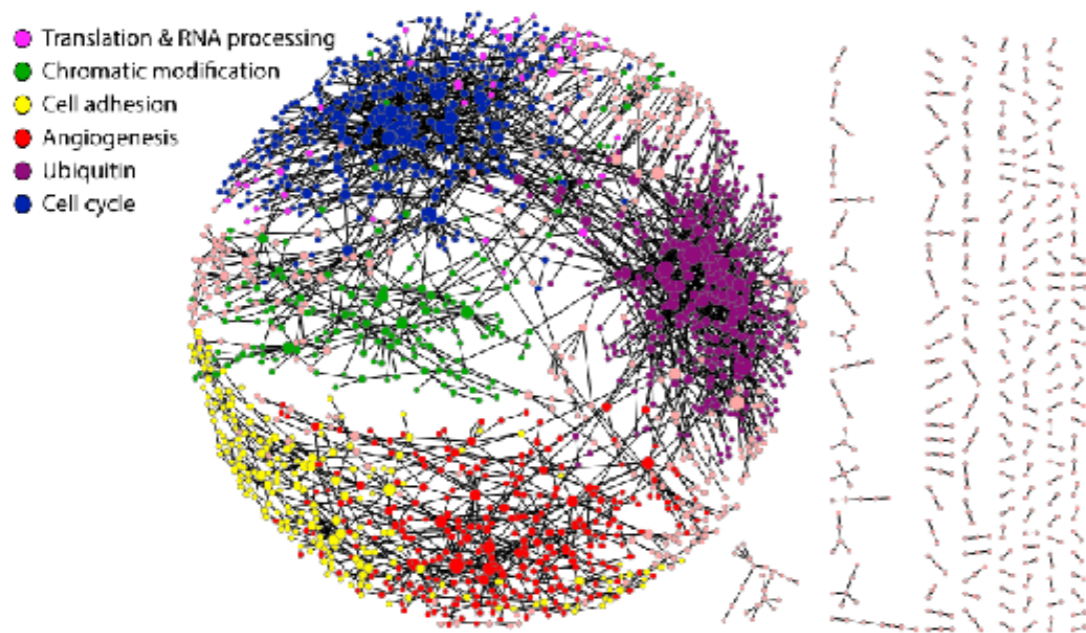


Synthetic Lethality

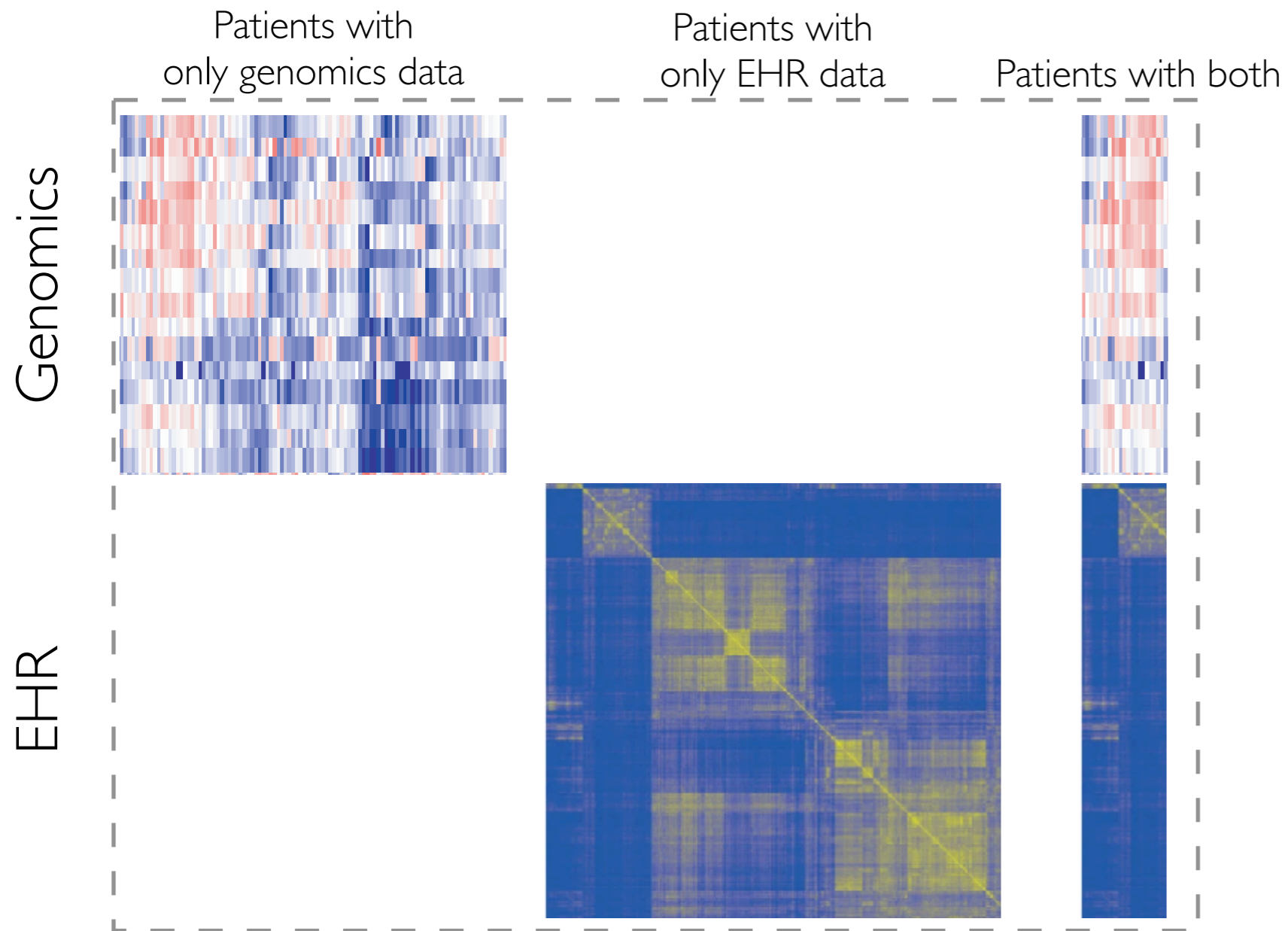
How to leverage SL to develop (personalized) drug (combination) therapy?

Integrate three sources:

- Mutation data of the patient (mutation A)
- SL network (Gene B has SL effect with Gene A)
- Drug target information (Drug X inhibits Gene B)



How to handle unpaired data



ML question: How to integrate all these patients?

Translation between features: RNA to ATAC translation

