Learning in Limited Data Settings Advancing Personalized Medicine in Cancer Treatment Planning*

Vaibhav Rajan Google DeepMind

*work funded at National University of Singapore



No sign-off yet | Foundation Medicine, Inc. | 1.888.988.3639

Can AI help us identify the right drug for such cancer patients?

Sample Preparation: 150 Second St., 1st Floor, Cambridge, MA 02141- CLIA: 22D2027531 Sample Analysis: 150 Second St., 1st Floor, Cambridge, MA 02141- CLIA: 22D2027531 PAGE 1. OÉ 25

NORMAL CELL AND CANCER CELL DEVELOPMENT



Image by brgfx on Freepik

Cancer is a *genetic disease*, i.e., it is caused by changes to genes (mutations) Cancer is a leading cause of death worldwide (one-in-six deaths, 2020)



Each cell in our body contains 23 pairs of chromosomes

Each chromosome is a sequence of "base pairs", bases are A, C, G, T

Gene: **subsequence of the chromosome** which has functional importance

~20,000 genes have been identified

https://www.cancer.gov/about-cancer/understanding/what-is-cancer

https://www.who.int/news-room/fact-sheets/detail/cancer

CANCER TREATMENT

- Treatment remains challenging
 - Complex disease: *Every cancer has an individual set of mutations*
 - A drug that works for one cancer patient, might have absolutely no effect on another
- > Treatment must be tailored to each patient: **personalized therapy**



https://www.worldwidecancerresearch.org/news-opinion/2021/march/why-havent-we-cured-cancer-yet/ https://en.wikipedia.org/wiki/Personalized_medicine





CANCER GENOMICS DATA

The Cancer Genome Atlas (TCGA)

Since 2006 > 11,000 patients 2.5 PetaBytes of Data 33 cancer types

Many similar data collection efforts to understand cancer



REPRESENTING GENOMIC DATA

Raw sequence (rarely used)

...ACCTTTCGGCCGGACCCCC...

Mutation Vector

Genes of interest

 G1	G2	G3	G4	G5	G6	G7	G8	G9	
1	0	1	0	0	0	0	1	0	
7 \									

Binary indicator: 1 \rightarrow mutation in gene, 0 \rightarrow no mutation

Gene Expression Vector

Genes of interest:

st:	G1	G2	G3	G4	G5	G6	G7	G8	G9	
	6	2.1	3	0	0	0	0	1	0	
		ount or	real va	alue: in	dicates	activit	y level	of gen	e	

Sequence of Mutations

G1: R273C, G1: S1372L, G2: L145V ...

In gene G1, at location 273 a mutation changed R to C in the protein

DRUG RESPONSE MEASUREMENTS

- 1. Response Evaluation Criteria In Solid Tumors (RECIST)
 - > Standard way to measure how well a cancer patient responds to treatment.

		RECIST
Cood response (label + 1)	CR	Complete Response
	PR	Partial Response
Bad response (label -1)	PD	Progressive Disease
bad response (laber - 1) 2	SD	Stable Disease

2. Progression-free Survival (PFS)

The length of time during and after the treatment (days/months/years), that a patient lives without the cancer getting worse.

DRUG RESPONSE PREDICTION (DRP)



Given:

- a patient's genomic profile and
- a drug

Will the response of the patient to the drug be good?

DRUG RESPONSE PREDICTION (DRP)

X: Patient's genomic data (e.g., mutation vector or gene expression)

- \succ Y: RECIST value after administering drug d
- > $Y \sim f_d(X) \rightarrow$ binary classification
- > Challenge
 - X: abundant, but...
 - *Y*: extremely limited for any drug *d*
- > Why?
 - Each patient is given one/few drugs, counterfactual unknown

CELL LINES: A RELATED "DOMAIN"



Extract cancer cells and clone them in lab (living cells, continue growing)
 Ensures each cell has same genomic data (X)

Administer multiple drugs on cell lines, measure response Y

DRUG RESPONSE MEASUREMENT IN CELL LINES

Area under the Dose Response Curve (AUDRC)

Real-valued [0,1]



- Administer progressively increasing concentration (X-axis) of drug and measure the amount of cancer cells (Y-axis) killed: Dose Response Curve (DRC)
- Lesser concentration kills more cells → more effective drug
 →Lower AUDRC
- E.g. efficacy of III > I > II

Vis, D. J. et al. Multilevel models improve precision and speed of IC50 estimates. Pharmacogenomics 2016.

CELL LINES: A RELATED "DOMAIN"



Extract cancer cells and clone them in lab (living cells, continue growing)
 Ensures each cell has same genomic data (X)

Administer multiple drugs on cell lines, measure response Y

> Can a Drug Response Prediction model on cell lines $Y \sim f_d(X)$ work for patients?

CELL LINES: A RELATED "DOMAIN"



Extract cancer cells and clone them in lab (living cells, continue growing)
 Ensures each cell has same genomic data (X)

Administer multiple drugs on cell lines, measure response Y

> Can a Drug Response Prediction model on cell lines $Y \sim f_d(X)$ work for patients?

No: drug responses differ across patients and cell lines



PROBLEM STATEMENT

Given:

Domain	Genomic Profile	Drug Response	#samples	$N_p \ll N_c \ll N_t$ $P(X_c) \neq P(X_t)$
Cell Lines	X _c	$Y_c^d \in R$ (AADRC)	N_c labeled	$f_c^d \not\sim f_t^d$
Patients	X _t	$Y_t^d \in \{0,1\}$ (RECIST)	N_p labeled N_t unlabeled	where $Y_c^d \sim f_c^d(X_c), Y_t^d \sim f_t^d(X_t)$

➤ Infer: Drug Response Prediction model $f_t^d: Y_t^d \sim f_t^d(X_t), \quad \forall \text{ drug } d \in \{d_1, d_2, \dots d_n\}$



No sign-off yet | Foundation Medicine, Inc. | 1.888.988.3639

Sample Preparation: 150 Second St., 1st Floor, Cambridge, MA 02141- CLIA: 22D2027531 Sample Analysis: 150 Second St., 1st Floor, Cambridge, MA 02141- CLIA: 22D2027531

PAGE 1 Of 25

DO PREVIOUS METHODS WORK IN THE CLINIC?

Input Data type

- Gene expression data
 - Assumed as inputs in previous methods
 - Not measured in FDA approved clinical panels

Mutation data

- Very sparse (typically a patient has ~10
 - mutations in panel out of ~million possible)
- Previous methods do not perform well with such inputs

Drug Repurposing

> Use of a drug for one cancer in another cancer
 > Need predictions on drugs not in training set



DRP Requirements

For clinical translation

- Training with mutations available in clinical sequencing reports
- Predict on drugs unseen during training
- Model varying length mutations

Utilise all available auxiliary patient response information (PFS)

From prior DRP literature

- Handle input discrepancy
- Handle output discrepancy
- Model patient mutation heterogeneity

Model Design

Model varying length mutations	Genes and mutations to be tokenized	Use transformers on gene and mutation level tokens

4

Model Design

Model varying length mutations	Genes and mutations to be tokenized	Use transformers on gene and mutation level tokens
Use all patient response-related information	Model survival information (PFS)	Pretraining transformers to predict survival

Model Training

Stage 1: Pretraining transformer with survival data

- Use all patient response-related information (PFS)
- Handles varying length mutation data
- Training on mutations

Stage 2: Training the Multi-task learning (MTL) model

- Handles input
 discrepancy
- Handles output
 discrepancy
- Allows drug repurposing



DISCRETE SEQUENTIAL DATA



DISCRETE SEQUENTIAL DATA



OUR APPROACH

1. Capture functional effects of mutations

• Features indicating harmfulness of each mutation

2. Transformers for sequential inputs

• Tokenization at gene and mutation level

3. Survival data for Supervised Representation Learning

• Pretraining and joint multi-task learning

FUNCTIONAL ANNOTATIONS

- Clinvar: 3 indicators pathogenic, benign, significance unknown
- GPD: 3 location indicators protein information unit, linker unit, non-coding unit
- Annovar: Predictions from 17 algorithms indicating deleteriousness (harmful/not)

23-dimensional feature vector per mutation



TRANSFORMERS FOR SEQUENTIAL INPUTS



MULTITASK LEARNING





MULTITASK LOGISTIC REGRESSION



PRETRAINING WITH SURVIVAL PREDICTION





EXPERIMENTS

Domain	Genomic Profile	Drug Response (RECIST)	Survival (PFS)	#(sample, drug) pairs	Drugs tested
Cell Lines	324 genes from	689		3632	
Patients	FoundationOne panel Mutation sequences	470	2512	15732	drugs with RECIST labels in > 80 patients

> Evaluation only on patients: 3 random 80-20 splits of (sample, drug) pairs

- > Binary classification task: RECIST category prediction
 - Metrics: AUROC and AUPRC

EXPERIMENTAL RESULTS

	5-Fluor	ouracil	Cisp	latin	Pacli	taxel	Ove	rall
	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC
CODE-AE	61.77(±11.63)	88.02(±4.85)	38.09(±16.34)	71.88(±8.46)	44.83(±15.08)	74.68(±5.11)	47.15(±7.82)	73.74(±6.90)
TCRP	48.39(±11.96)	$\overline{77.87(\pm 9.49)}$	50.86(±16.75)	79.60(±10.73)	60.66(±24.22)	78.69(±16.73)	47.36(±3.26)	75.70(±4.03)
TUGDA	58.39(±16.94)	83.40(±7.69)	37.47(±4.36)	72.07(±3.75)	41.13(±18.28)	71.67(±1.03)	46.18(±3.18)	75.42(±2.22)
Velodrome	50.91(±19.54)	79.08(±15.04)	48.61(±9.73)	78.80(±7.97)	63.26(±26.45)	80.02(±16.34)	52.56(±4.93)	77.62(±0.48)
DruID	64.73(±8.73)	85.55(±6.43)	67.38(±10.63)	86.30(±7.35)	63.43(±4.97)	82.55(±6.83)	62.36(±3.60)	82.06(±5.02)
PREDICT-AI	71.30(±3.87)	88.74(±2.82)	72.37(±10.10)	90.26 (±4.68)	62.19(±9.42)	81.08(±8.02)	64.96(±4.50)	84.85(±4.02)

EXPERIMENTAL RESULTS

	5-Fluor	ouracil	Cispl	atin	Pacli	taxel	Ove	rall
	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC
CODE-AE	61.77(±11.63)	88.02(±4.85)	38.09(±16.34)	71.88(±8.46)	44.83(±15.08)	74.68(±5.11)	47.15(±7.82)	73.74(±6.90)
TCRP	48.39(±11.96)	77.87(±9.49)	50.86(±16.75)	79.60(±10.73)	60.66(±24.22)	78.69(±16.73)	47.36(±3.26)	75.70(±4.03)
TUGDA	58.39(±16.94)	83.40(±7.69)	37.47(±4.36)	72.07(±3.75)	41.13(±18.28)	71.67(±1.03)	46.18(±3.18)	75.42(±2.22)
Velodrome	50.91(±19.54)	79.08(±15.04)	48.61(±9.73)	78.80(±7.97)	63.26(±26.45)	80.02(±16.34)	52.56(±4.93)	77.62(±0.48)
DruID	64.73(±8.73)	85.55(±6.43)	67.38(±10.63)	86.30(±7.35)	63.43(±4.97)	82.55(±6.83)	62.36(±3.60)	82.06(±5.02)
PREDICT-AI	71.30(±3.87)	88.74(±2.82)	72.37(±10.10)	90.26 (±4.68)	62.19(±9.42)	81.08(±8.02)	64.96(±4.50)	84.85(±4.02)

More details/results: https://arxiv.org/abs/2402.10551

Aishwarya Jayagopal, Hansheng Xue, Ziyang He, Robert J Walsh, Krishna Kumar H, David SP Tan, Tuan Z Tan, Jason J Pitt, Anand D Jeyasekharan, Vaibhav Rajan

Personalised Drug Identifier for Cancer Treatment with Transformers using Auxiliary Information KDD 2024

One Last Requirement...

DRP Requirements

For clinical translation

- Training with mutations available in clinical sequencing reports
- Predict on drugs unseen during training
- Model varying length mutations
- Utilise all available auxiliary patient response information (PFS) From prior DRP literature
- Handle input discrepancy
- Handle output discrepancy
- Model patient mutation heterogeneity





GANDALF: Generative AtteNtion based Data Augmentation and predictive modeLing Framework

ICLR 2025

Jayagopal, A., Zhang, Y., Walsh, R.J., Tan, T.Z., Jeyasekharan, A.D. and Rajan, V., 2025. GANDALF: Generative AttentioN based Data Augmentation and predictive modeLing Framework for personalized cancer treatment.

Model Design

Requirement	Considerations	Design Choice
Model patient mutation	Generate more patient-like data,	Generate samples from cell lines,
heterogeneity	from cell line profiles	using diffusion models

Treatment Recommendation System

MOLECULAR TUMOR BOARD (MTB)

- Treatment planning for complex cancer cases is increasingly done by a Molecular Tumor Board
- Several expert clinicians collectively decide on the most suitable treatment







AI-IN-THE-LOOP MTB



TREATMENT RECOMMENDATION SYSTEM @ NUH, SINGAPORE



https://pharmacope.ai/

TREATMENT RECOMMENDATION SYSTEM @ NUH, SINGAPORE

Case list / Patient 4 Patient information Patient 4 Age: NA Age: NA Sex: NA Cenomic Data Variant's of Unknown significance UDENTIFIED VINORELBINE ERLOTINIB VINORELBINE VINORELBI	Feedback
Patient information Patient 4 Patient 4 Patient 4 Sex: NA Sex: NA Sex: NA Recommendation Recommendation Recommendation Sex: NA Recommendation Sex: NA Recommendation Recommenda	Feedback
Are: NA Sex: NA Genomic Data VARIANTS OF UNKNOWN SIGNIFICANCE GENE ALTERATIONS VARIANTS OF UNKNOWN SIGNIFICANCE IDENTIFIED None found. WYC amplification-equivocal	Feedback
Patient 4 Age: NA Sex: NA Genomic Data Sex: NA Genomic Cata Variants of UNKNOWN SIGNIFICANCE IDENTIFIED GENE ALTERATION IDENTIFIED None found. MYC amplification-equivocal Drug Recommended Internatives VINORELBINE ERLOTINIB O.8974 O.8974	Feedback
Age: NA Sex: NA Genomic Data CETUXIMAB GENOMIC ALTERATIONS IDENTIFIED VARIANTS OF UNKNOWN SIGNIFICANCE IDENTIFIED GENE ALTERATION None found. ERLOTINIB MYC amplification-equivocal	~
Genomic Data CETUXIMAB 0.8980 GENOMIC ALTERATIONS VARIANTS OF UNKNOWN SIGNIFICANCE ALTERNATIVES IDENTIFIED IDENTIFIED VINORELBINE GENE ALTERATION None found. WYC amplification-equivocal ERLOTINIB	~
GENOMIC ALTERATIONS VARIANTS OF UNKNOWN SIGNIFICANCE VINORELBINE 0.8974 IDENTIFIED IDENTIFIED ERLOTINIB 0.7594 MYC amplification-equivocal IDENTIFIED IDENTIFIED	
GENOMIC ALTERATIONS VARIANTS OF UNKNOWN SIGNIFICANCE VINORELBINE O.8974 IDENTIFIED IDENTIFIED ERLOTINIB 0.7594 MYC amplification-equivocal IDENTIFIED IDENTIFIED	
GENE ALTERATION None found. ERLOTINIB 0.7594 MYC amplification-equivocal	~
MYC amplification-equivocal	~
DOCETAXEL 0.7531	~
APC K534* BELINOSTAT DO.7317	~
TP53 splice site 93-1_96delTCTGG 0.7226	~
APC E129Q CARMUSTINE 0.7216	~
NPP4B CYTARABINE 0.7194	~
ARIDIB Q121_Q128>Q PALBOCICLIB 0.7027	~
MLL S2319T TEMOZOLOMIDE 0.7023	

https://pharmacope.ai/

SUPPORTING EVIDENCE: WHY?

Clinical decision of prescribing a specific drug

- Clinical guidelines
- Evidence of efficacy from:
 - 1. Previous clinical trials
 - 2. Mechanism of action

Completely or partially unknown

Mostly unknown

>Multiple incomplete or indirect sources of evidence needed

• Even to evaluate a DRP model in a clinical trial

SUPPORTING EVIDENCE

Model Explainability

• XAI algorithms to highlight genes most important for prediction

Auxiliary Drug Databases

• Gene-drug associations, clinical trial information (mostly limited to effects of single gene mutations)

Drug-level validation across patients

• Difference in the prediction (for the input patient) to the predictions on a reference set

Patient-level validation across drugs

• Distribution of predictions for all drugs

SUPPORTING EVIDENCE



Difference in the prediction (for the input patient) to the predictions on a reference set

Clinical trial



https://clinicaltrials.gov/study/NCT05719428

CONCLUSION

Advanced state-of-the-art in Cancer Drug Response Prediction (DRP) literature

> Al for Personalized Cancer Treatment

First* clinical trial where patients are being treated in a Molecular Tumor Board with our DRP-based recommendations

https://clinicaltrials.gov/study/NCT05719428

First* (& current best) DRP model on clinically available genomic data for personalized cancer treatment

WISER: ICML 2024

GANDALF: ICLR 2025

DRUID: Cell iScience PREDICT-AI: KDD 2024

TEAM

Aishwarya Jayagopal Hansheng Xue Ziyang He Krishna Kumar Ragunathan Mariappan Debabrata Mahapatra **Vaibhav Rajan**

Tuan Zea Pitt Jason J Pitt

Renie Ravin

Robert J Walsh Patrick Williams Jaynes Diana Lim David Shao Peng Tan Anand D Jeyasekharan

ML Team at School of Computing, NUS Bioinformatics Team at Cancer Science Institute, NUS

Primal Rave Software Labs, India Clinical Team at National University Hospital

Principal Investigators: Vaibhav Rajan and Anand D Jeyasekharan

FUTURE DIRECTIONS

Techniques

- LLMs
- Gene representations from Knowledge Graphs

Model Improvements

- Drug combinations
- Additional inputs clinical data, genomic rearrangements
- Temporality

Decision Support Systems

- Supporting Evidence from External Knowledge Bases
- Integrate supporting evidence in ranking recommendations
- Improved AI-in-the-loop system for collective decision making in MTBs

Thank you!

Thanks to Aishwarya Jayagopal for most of the slides