# **Closing the Target Gap** A Computational Approach to Optimizing Therapeutic Selection for Cancer Patients

# Mikhail Grushko, Jeremy Goldstein, Zakary ElSeht, Alfonso Alarcon, Nichole Jones, Mahta Samizadeh, Yuhan Zhu, Johanne Kaplan, Katherine Arline

## **OBJECTIVES**

- . To identify optimal therapies for specific cancers and subtypes of cancer which may not have previously been clinically explored
- II. To identify optimal therapies for specific patients based on their -omic information
- III. To match new and existing therapies with the specific patient populations that are most likely to benefit from drug treatment

## **BACKGROUND**

While significant progress has been made in developing new therapies for cancer patients, many patients lack treatments that result in good outcomes. Existing patient-therapy matching algorithms frequently rely on mutations and well-studied targets for which only a limited number of FDA-approved therapies exist. In contrast, SHEPHERD's approach, called DELVE, uses computational and mathematical tools informed by transcriptomic data to match therapies with the models, cancers, and specific patients that will be most impacted by drug treatment, regardless of mutational status.

The DELVE platform, which is comprised of dozens of separate tools, was constructed to inform intelligent drug design, development, and clinical use by identifying the simultaneous mechanisms of action for any specific drug, and matching those mechanisms with -omic-level patient data. DELVE operates in contrast to the "target-based" development methodology which excludes the therapeutic utility of off-target effects.

For drug developers, DELVE is able to match specific molecules with predicted clinical utility in 161 cancers, and to identify potential utility in comparison with

standard of care. For point-of-care use, DELVE incorporates transcriptomic data and can match patient-level data with 396 drugs agnostic of the patient's mutational status. DELVE differs from transcriptomic direction reversal techniques in that, while it can incorporate direction reversal, it also can identify the potential for a drug to leverage dysregulations, such as gene overexpression, to boost efficacy.

For each therapy under study, DELVE requires canonical SMILES representing molecular structure and 100 or fewer post treatment IC50 values that measure the sensitivity of cell lines to the therapy under study and which include multiple cell lines that a) clearly respond to therapy and b) are resistant to therapy.

#### **Poster Overview**

A major component of DELVE is a broad repository of cancer -omic data, which has been aggregated over three years from public sources [Table 1]. Table 2 contains an overview of some of the dozens of tools that compose the DELVE platform. One of them, GCVA, is highlighted in figure 1. Figure 2 describes how GCVA-output genes interact at the protein level for a specific therapy. Figure 3 maps the computational pipeline which knits together DELVE tools to generate reports. Figure 4 contains example outputs from GCVA on generic drugs with high and low probability for therapeutic utility in additional patients. Figure 5 describes the overall accuracy achieved using DELVE tools in predicting outcomes on pre-clinical models as well as on retrospective clinical trial data.

## METHODS

DELVE leverages bioinformatics, chemoinformatics, proprietary algorithms, deep learning neural networks, random tree forests, and other tools to generate transcriptomic-level drug response-resistance signatures. DELVE integrates over 75,000 patient samples representing 161 cancers as well as healthy tissues, and was deployed to predict drug response and resistance across thousands of *in vivo*, *ex vivo*, and *in* vitro cancer models.

#### Cancer Naming

A significant challenge in pan-cancer analysis is heterogeneity in naming convention as we as misclassifications of cancer models. DELVE's tools operate on cancer patient omic data and cancer model data repositories which have been curated, cleaned, and when necessary renamed to conform with a consolidated data dictionary. Cancer naming conventions have been derived from prior work documented by Arline, et al.<sup>1,2</sup>

Patient Database

Microarray data have been processed using the crossmeta package and hand curation, RNAseq data has been normalized as log(TPM+10<sup>n</sup>) transformed transcripts per million and hand curated. Doppelgänger is used to remove duplicate patient samples by looking for suspicious similarities in expression between all samples.

Classification Subclassifications of cancers and classifications of specific patient samples as noted in figure 3a are accomplished with several tools. Random forest classification models are use to classify patient samples based both on mutational and transcriptomic characteristics. Unsupervised non-negative matrix factorization and pathway analysis models are used to generate subclassifications of cancers where they do not exist in literature.

GCVA is a proprietary algorithm which intakes cell lines transcriptome data and associated IC50s measuring sensitivity to drug treatment. GCVA utilizes the direct comparison between the highest and lowest responding samples in this training set. This approach ensures the detection of transcriptomic variation that is most directly responsible for the differences in drug response across the panel of sensitivity data at hand. In general, the more sensitivity data are available the more reliable the predictions

For instances where data are inadequate due to data set size, sampling bias, noise, or imbalance, a separate variation of GCVA exists. This variation institutes sampling adjustments to avoid Bayesian imbalance and reshuffles sensitivity order within extreme samples. Such

algorithm adjustments allow DELVE to process data for under-researched therapeutics and to deliver highly predictive sensitivity projections, as quantified by the measures in the results section and in figure 5. **Machine Learning** Machine learning, and deep learning in particular, are used in a variety of DELVE tools due to their versatility and variety of computational methods. DELVE's Blood Brain Barrier Deep Learning Network, Deep Learning IC50 Network, Deep Learning Synergy Networks and Chemical Analytics all utilize a variety of network or classifier based machine learning algorithms. The most important feature made possible through machine learning is ability to make predictions just from the known chemical structure alone in the absence of sensitivity data. ML methods utilized in DELVE include convolutional neural networks (CNN), feedforward neural networks, random tree classifiers, support vector classifiers, non-negative matrix factorization, natural language transformers, and chemoinformatic featurizers.

## RESULTS

DELVE was able to correctly classify the highest and lowest responding drug-cell line pairs with 96% sensitivity [95% CI +/- 0.95%] and 88% specificity [95% CI +/- 1.0%]. Across published in vivo studies related to 20 FDA-approved cancer therapies, drug-model pairs which achieved a high DELVE score showed greater than or equal to 90% tumor growth inhibition in vivo 78% of the time. Drug-model pairs with low DELVE scores showed less than or equal to 60% tumor growth inhibition in vivo 77% of the time. Across 71 FDA-approved drugs, using chemical structure alone, the platform was able to predict at least one approved solid tumor indication 84% of the time. Random chance indicationtherapy pairing was correct only 21% of the time. Across 10 failed therapies, the platform was able to predict failure with 82% accuracy.

1 Arline, K., Rare Isn't Rare; (Abstract #7739). Presented at The American Association of Cancer Researchers Annual Proceedings, April 2018, Chicago, Illinois 2 Treuting, R.L., Rare Cancer's 'Valley of Death'; (Abstract #2505). Presented at The American

Association of Cancer Researchers Annual Proceedings, March 2019, Atlanta, Georgia.

#### TABLE 1: DATA SOURCES

TABLE 1A: PATIENT DATABASE

	Count
Cancer Samples	77,615
Represented Cancers	161
Healthy Tissue Samples	13,284

ERD's patient database has been cted from 1,215 public sources. These sist of RNAsea and microarray data c y and tumor tissue. Data were cleaned sed, and analyzed for quality.

#### **TABLE 1B: MODEL REPOSITORY**

	Model Count	Omic Data Count	<b>Cancers Represented</b>
In Vitro	1,898	1,408	165
<i>Ex Vivo /</i> PDX	2,029	1,950	149
In Vivo	940	279	158

Cancer models, including *in vitro*, *ex vivo*, PDX, and *in vivo*, have been curated by hand to share naming conventions across all data sources. Cancer naming is determined via SHEPHERD's research identifying specific forms of cancer.<sup>1,2</sup> -Omic data on these models has been aggregated when available or generated internally by SHEPHERD.

The availability of this information allows DELVE to identify models that best represent specific cancers, and to identify those models most likely to be highly affected by a specific drug treatment from *in vitro* research tools through to *in vivo*.

### FIGURE 1: GCVA



### FIGURE 2: PROTEIN LEVEL CONNECTIONS

GCVA Predicted Genes	Literature and Bioinformatic-Identified Interactions	
MSX1	EP300, TP53, NEFM, SHH, PTCH1, GAS1, GLI1	
GUCY1B3	KDR, MAPK1, MAPK3	
GLI2	MYB, EP300, RB1, MAPK1, TP53, CDKN1A, CCNE1, CCNB1, CCNA1, ABL1, PSM9, KDR, TGFB1, MAPK3, TNF, CDH1, SHH, GLI1, PTCH1, SUFU, GAS1, SMO, HHIP	
FOXG1	GLI1, MK167, CCNA1, CDKN1A, TP53, CDK7, SHH, MYB, EP300, MAPK3, MAPK1	
FOXD4L6	CDK7, EP300	
ETV4	ABL1, MAPK1, MAPK3, CCND1, MKI67, TP53, EP300, CCNB1, GLI1	
EPS8L1	ABL1, CDH1, MAPK1, MAPK3	

While GCVA is powered by transcriptomic data, proteins associated with GCVA-derived genes for a specific therapy frequently have documented interconnections with genes and proteins known by other methods to interact with the molecule's canonical and non-canonical targets. These targets and interactions are identified via literature review and bioinformatic analysis, including pre- and post-treatment cell line sequencing, as part of DELVE. The example in this figure is for the repurposed therapy mebendazole. The GCVA-predicted gene GLI2 has 23 protein-level interactions with proteins that have a known relationship with the repurposed therapeutic mebendazole

TABLE 2: DELVE TOOLS		
ΤοοΙ	Application	
Subclassification Algorithms	Expression and mutation based methods for finding subgroups in pa- tient data and identifying potential new targets in understudied cancers	
Deep Learning Synergy Networks	Prediction of enhanced drug efficacy in combination with standard-of-care or other oncology assets, both across indications and for specific cell lines	
Deep Learning IC50 Network	Ability to predict drug efficacy across 1,300 cell lines based on molecular structure, or structure and seed data	
GCVA	Proprietary drug response-resistance algorithm predicts efficacy across cell lines, <i>ex vivo</i> and <i>in vivo</i> models, aggregate patient datasets, and on individual patient data	
GCVA - Toxicity	GCVA deployed on healthy tissue samples to predict toxicity on organs	
RNA sequencing Pipeline	Internally developed sequencing pipeline including quality control, align- ment, fusion detection, somatic variant calling, and multiple differential gene expression, pathway, and gene-set analysis techniques.	
Model Repository	-Omic data on <i>in vitro, ex vivo</i> , and <i>in vivo</i> models, cleaned and classi- fied according to SHEPHERD cancer classifications	
BBB Deep Learning Network	Prediction of the ability of an asset to cross the blood-brain barrier, built on hand-curated chemical interactions	
Chemical Analytics	Ability to predict lipophilicity and hydrophilicity	
Homogeneity Analysis	Ability to predict homogeneity of cancers based on patient samples, and to compare cancers with each other	
High Throughput In Silico Pipeline	Ability to automatically assess the utility of hundreds of drugs on tens of thousands of patient and model data points	



data are derived from characteristics of the drug as revealed by GCVA analysis of cell line IC50 data. The area between the green curve and the red line quantifies GCVA's ability to improve on population response rates for populations selected without DELVE at any proportion of "true" population responder share.





For patient-specific analysis, DELVE is executed simultaneously on 396 therapies which include oncology-specific targeted therapies, chemotherapies, chemotherapy-related agents, current an discontinued therapies in development, and repurposed drugs with known activity in cancer. The average time for report generation is 8 minutes leveraging Amazon AWS infrastructure.

#### FIGURE 4: EXAMPLE THERAPEUTIC REPOSITIONING ANALYSIS

DELVE analysis of public data on existing cancer therapies as well as failed therapies These charts are derived from DELVE's GCVA algorithm and *in vitro* screening data. reveals variable opportunities for benefiting additional patients and redeployment i They plot the difference between predicted response rates in a DELVE-selected patient additional cancers. Therapies which appear to have more potential for repositioning population (green arc) versus the response rate that otherwise would be achieved in a typically show efficacy in a broader range of pre-clinical models than therapies which are clinical trial. unlikely to succeed beyond a specific cancer, or which fail in cancer clinical trials overall

#### FIGURE 4A: HIGH PROBABILITY FOR REPOSITIONING – TEMSIROLIMUS



#### FIGURE 5: ACCURACY

As outlined in results, DELVE accuracy has been characterized on preclinical tools and on retrospective clinical trial data. In vitro efficacy accuracy is DELVE's ability to predict highest and lowest performing drug-cell line pairs within an appropriate IC50 range for each drug. Clinical approvals accuracy is DELVE's ability to predict at least one successful drug approval for each therapy. Clinical failure accuracy is DELVE's ability to identify that specific cancers under study would not successfully be treated by a specified drug.





## **Presentation 1146P** • 2021 ESMO Congress

Fifty-six separate processes are executed for a therapeutic focused workflow, resulting in a list of cancers on which the therapy is predicted to be effective and associated statistics. The average time for report generation is 22 minutes.

#### FIGURE 4B: LOW PROBABILITY FOR REPOSITIONING – PEMETREXED



#### DATA HIGHLIGHTS

• DELVE can predict therapeutic efficacy across 161 cancers, including estimates of comparative efficacy to standard of care.

- DELVE can identify high responding pre-clinical models from a repository of 4,867 models representing 165 cancers.
- DELVE can predict the comparative efficacy of 396 compounds on specific patient samples.
- Single patient reports are generated in approximately 8 minutes.
- DELVE can operate on small molecules and *in vitro*-assessable biologics.

## CONCLUSIONS

. DELVE is able to predict drug efficacy on cancer models and to correctly select indications for existing therapies, supporting its utility in predicting new indications for cancer therapies. Because the platform can operate on any single transcriptome, it is possible to use this tool to match patients with therapies from which they may benefit without reliance on a previously classified target.

II. DELVE has the ability to predict new indications for in-development and existing oncology therapies across 161 cancer types. DELVE is also able to offer more precise patient population selection guidance. which can increase the overall clinical trial response rate.

## delve@shepherd.bio • shepherd.bio

## Boston, MA • USA