# DiSCoVER:
# Evaluating An Algorithm for Cancer-Related Medicine Recommendations

Gabriel Diaz

UCSD STARS

Mesirov Lab
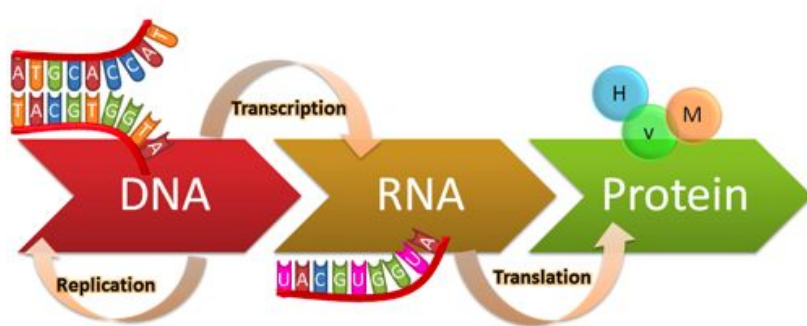
Summer Research Conference 2019

# Overview

- Background Information
- Proposed Solution
- Methods
- Results
- Future Steps

# Background Info

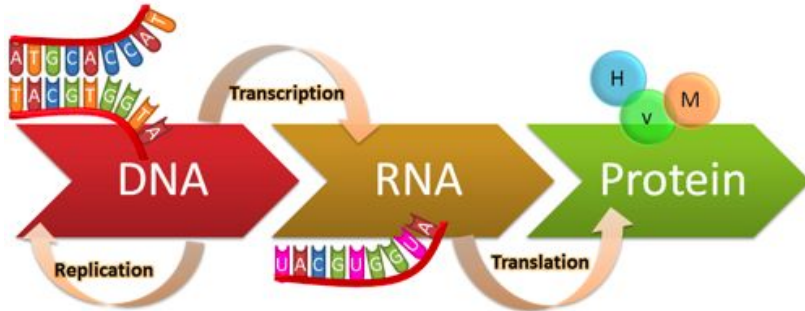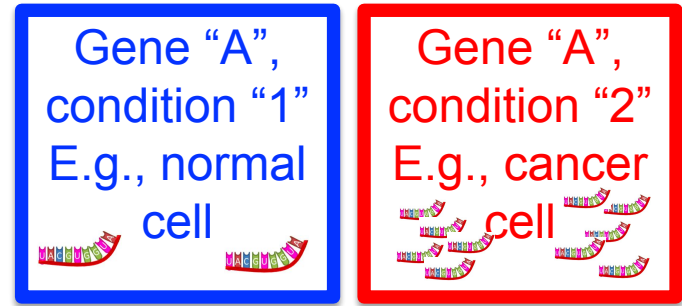- Precision Medicine
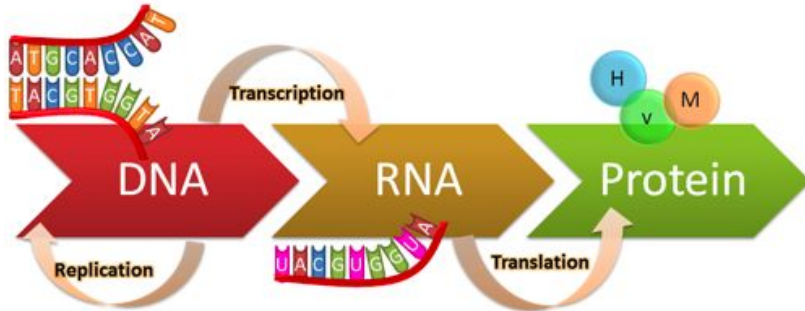- Protein Synthesis
- Gene Expression and Cancer



https://genius.com/Biology-genius-the-central-dogma-annotated

# Background Info

- Precision Medicine
- Protein Synthesis
- Gene Expression and Cancer
- GE Lists/Datasets



Gene "A", condition "1" E.g., normal cell
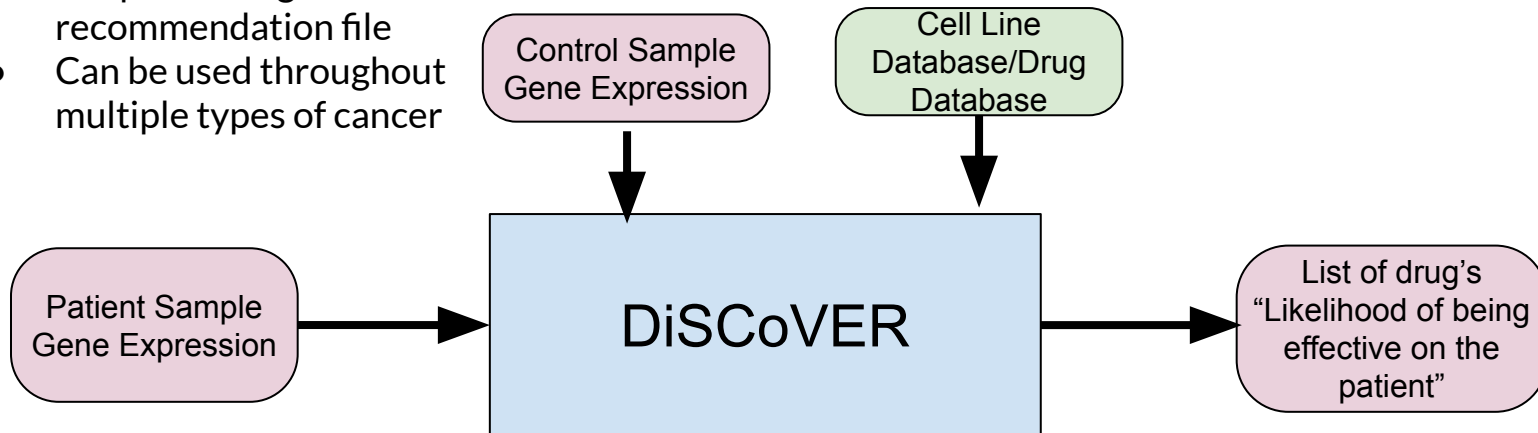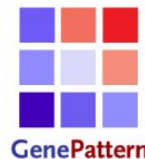
Gene "A", condition "2" E.g., cancer cell



https://genius.com/Biology-genius-the-central-dogma-annotated

# Background Info

- Precision Medicine
- Protein Synthesis
- Gene Expression and Cancer
  GE lists/datasets



https://genius.com/Biology-genius-the-central-dogma-annotated

Gene "A", condition "1" E.g., normal cell

Gene "A", condition "2" E.g., cancer cell

Rows: Genes

Columns: Cell Sample From Patients

| Name | TCGA-A7-A0CE-01 |
|---|---|
| TSPAN6 | 7.023572 |
| DPM1 | 5.210623 |
| SCYL3 | 4.545135 |
| C1orf112 | 4.982679 |
| FGR | 2.144122 |
| CFH | 4.67671 |
| FUCA2 | 7.331436 |
| GCLC | 5.163128 |
| NFYA | 5.415756 |
| STPG1 | 4.159719 |
| NIPAL3 | 6.384105 |
| LAS1L | 6.187223 |
| ENPP4 | 3.863163 |
| SEMA3F | 5.744268 |
| ANKIB1 | 6.885191 |
| KRIT1 | 5.720907 |
| RAD52 | 3.838566 |
| BAD | 3.730756 |
| LAP3 | 6.832838 |
| CD99 | 6.452564 |
| HS3ST1 | 0.959698 |
| MAD1L1 | 5.675192 |
| LASP1 | 7.668393 |
| SNX11 | 4.388745 |

# Background Info: DiSCoVER

- "Disease-model Signature vs. Compound-Variety Enriched Response"
- Input: 2 gene expression files
- Output: 1 drug recommendation file
- Can be used throughout multiple types of cancer

# Issue At Hand

Scientists sometimes need to preprocess data before running a computer analysis. For DiSCoVER, this may result in inaccurate recommendations.

# Solution

Find relationships between the outputs of DiSCoVER, where one output comes from a raw dataset and the other comes from a preprocessed dataset. If both outputs are not closely related, the accuracy of DiSCoVER's recommendations can may be questionable.

# Methods

Collect Open-Source Data From TCGA

Parse/Upload each Dataset to DiSCoVER

Rank/Compare Scores

# Methods

| Collect Open-Source Data From TCGA | Parse/Upload each Dataset to DiSCoVER | Rank/Compare Scores |

Columns: Cancer Patients

Rows: Genes



Raw Dataset



Preprocessed Dataset

- We are using RNA-Seq samples from 20 breast cancer (BRCA) patients.
- Rows represent genes, Columns represent cancerous or normal cells from each patient
- 755 columns by 40 rows

NIH ▶ NATIONAL CANCER INSTITUTE

# Methods

# Methods

Collect Open-Source Data From TCGA → Parse/Upload each Dataset to DiSCoVER → Rank/Compare Scores

TCGA Gene Expression Dataset (1 of 2)



Cancerous Tissue Sample ("01")

Normal/Control Tissue Sample ("11")

40 Columns

X 20

# Methods

# Methods

# Methods

| Collect Open-Source Data From TCGA | Parse/Upload each Dataset to DiSCoVER | Rank/Compare Scores |

## Example Dataset

| drug | TCGA-BH-A0B3-01_DiSCoVER_result | TCGA-BH-A0BC-01_DiSCoVER_result | TCGA-A7-A0CH-01_DiSCoVER_result | TCGA-BH-A0B7-01_DiSCoVER_result | TCGA-A7-A13G-01_DiSCoVER_result |
|---|---|---|---|---|---|
| gdsc_Motesanib | 0.597 | -0.596 | -0.611 | 0.602 | -0.686 |
| gdsc_PD173074 | 0.590 | -0.528 | -0.548 | -0.528 | -0.591 |
| gdsc_GW441756 | 0.573 | -0.568 | 0.576 | 0.515 | 0.653 |
| gdsc_Cetuximab | 0.570 | 0.539 | 0.606 | 0.551 | 0.706 |
| gdsc_Afatinib | 0.554 | 0.577 | 0.584 | 0.617 | 0.767 |

# Methods

## Example Dataset

5]:

| | drug | TCGA-BH-A0B3-01_DiSCoVER_result | TCGA-BH-A0BC-01_DiSCoVER_result | TCGA-A7-A0CH-01_DiSCoVER_result | TCGA-BH-A0B7-01_DiSCoVER_result | TCGA-A7-A13G-01_DiSCoVER_result |
|---|---|---|---|---|---|---|
| 0 | ccle_17-AAG | 515.0 | 186.0 | 186.0 | 350.0 | 183.0 |
| 1 | ccle_AEW541 | 251.0 | 390.0 | 322.0 | 373.0 | 269.0 |
| 2 | ccle_AZD0530 | 378.0 | 174.0 | 158.0 | 223.0 | 123.0 |
| 3 | ccle_AZD6244 | 570.0 | 305.0 | 297.0 | 601.0 | 369.0 |
| 4 | ccle_Erlotinib | 565.0 | 92.0 | 64.0 | 150.0 | 42.0 |
| 5 | ccle_Irinotecan | 242.0 | 642.0 | 722.0 | 293.0 | 712.0 |

# Methods

| Collect Open-Source Data From TCGA | Parse/Upload each Dataset to DiSCoVER | Rank/Compare Scores |
|---|---|---|

## Ranking of Raw Output



## Ranking of Preprocessed Output

# Methods

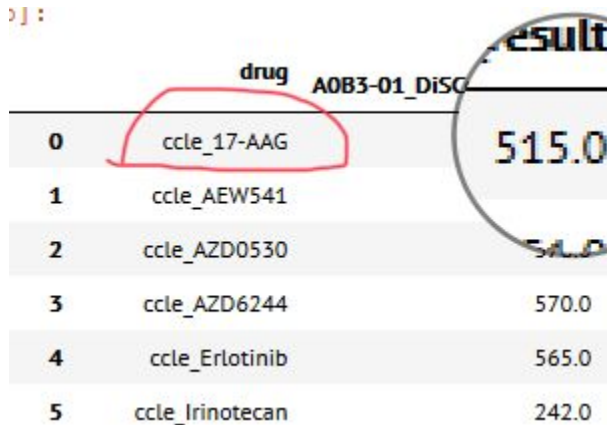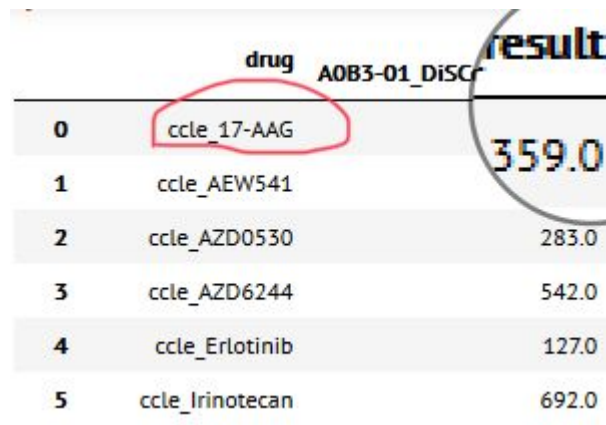| | | |
|---|---|---|
| Collect Open-Source Data From TCGA | Parse/Upload each Dataset to DiSCoVER | Rank/Compare Scores |

## Kendall's Tau

- Represents degree of correlation between two columns of ranked data
- Can range from -1.0 to 1.0
- Based on concordant and discordant pairs

| Drug | Patient1 | Patient2 | C | D |
|---|---|---|---|---|
| A | 1 | 1 | 5 | 0 |
| B | 2 | 2 | 4 | 0 |
| C | 3 | 4 | 2 | 1 |
| D | 4 | 3 | 2 | 0 |
| E | 5 | 6 | 0 | 1 |
| F | 6 | 5 | 0 | 0 |
| | | Total: | 13 | 2 |
| | | KT | 0.733 | |

# Preliminary Results:



Kendall Coefficients for Patient Drug Rankings

# Possible Areas of Further Research

- Analyze possible data preprocessing methods
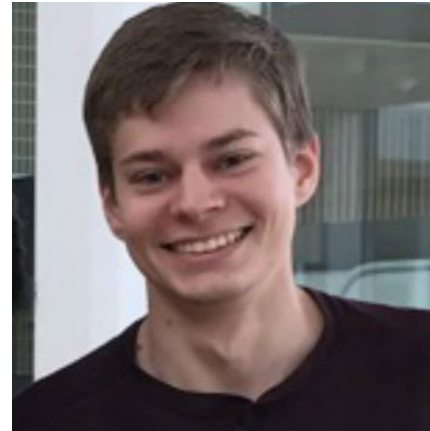- Determine potential sources of "noisy data" for recommendation scores

# Acknowledgments - Mesirov Lab



**Dr. Jill Mesirov**
PI/Lab Director

**Dr. Edwin Juarez**
Bioinformatics Programmer
/Research Guide

**Alex Wenzel**
PhD Candidate

**Owen Chapman**
PhD Candidate

# Thank You!