# Predicting Cancer Phenotypes based on Somatic Genomic Alterations via Genomic Impact Transformer

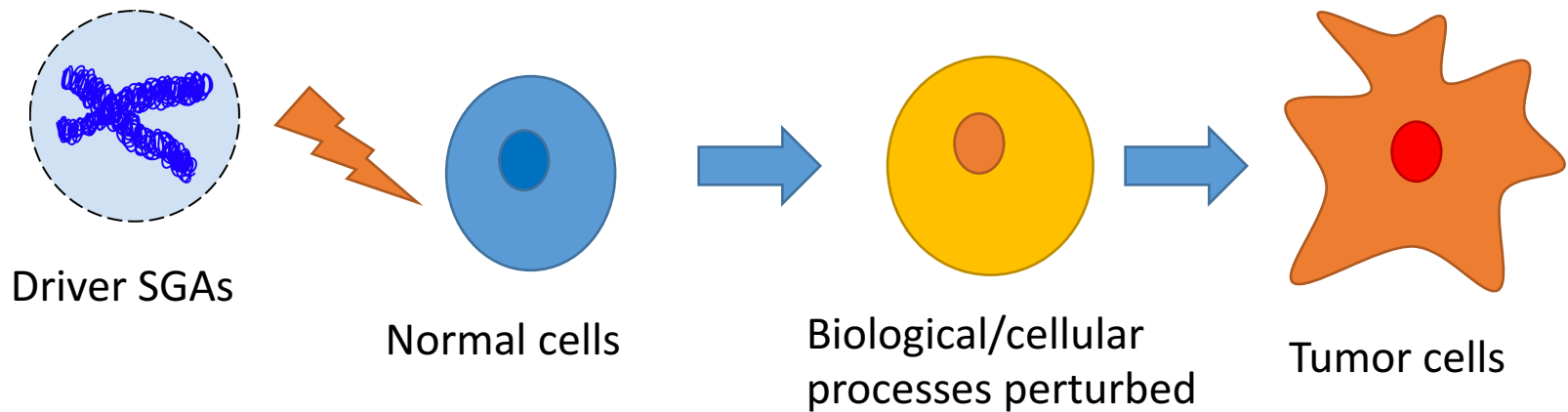Yifeng Tao[1], Chunhui Cai[2], William W. Cohen[1*], Xinghua Lu[2*]

[1]Carnegie Mellon University

[2]University of Pittsburgh

# Background

o Cancers are mainly caused by somatic genomic alterations (SGAs)
  - o Driver SGAs → causal to tumor development
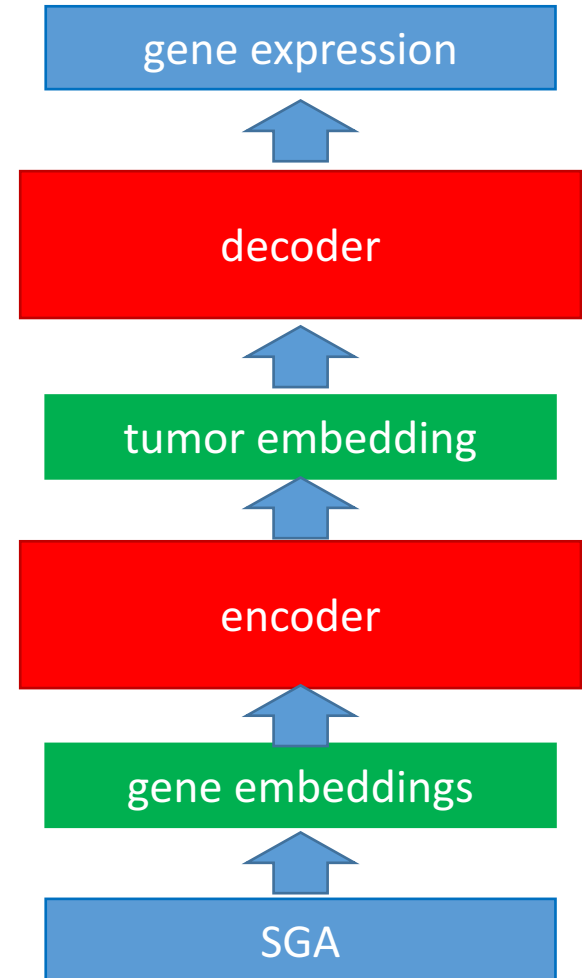  - o Passenger SGAs → neutral mutations



Driver SGAs     Normal cells     Biological/cellular processes perturbed     Tumor cells
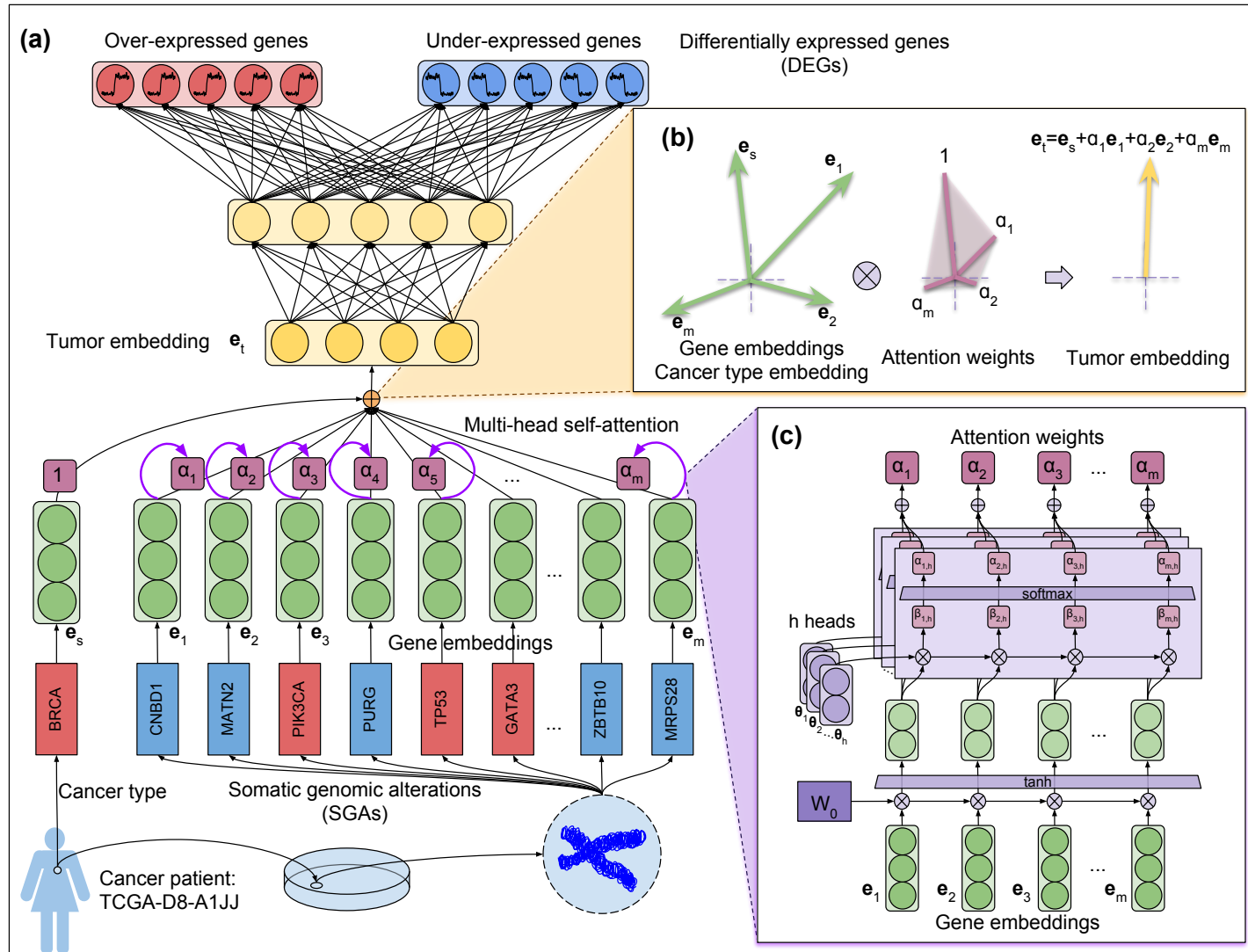
# Challenges

o Driver SGA detection
  o Solution 1: frequency
  o Solution 2: conserved domain of protein
  o Problem: downstream effect of SGAs

o SGA/tumor representation
  o Solution: a higher dimensional one-hot/sparse vector
  o Problem: little information/knowledge
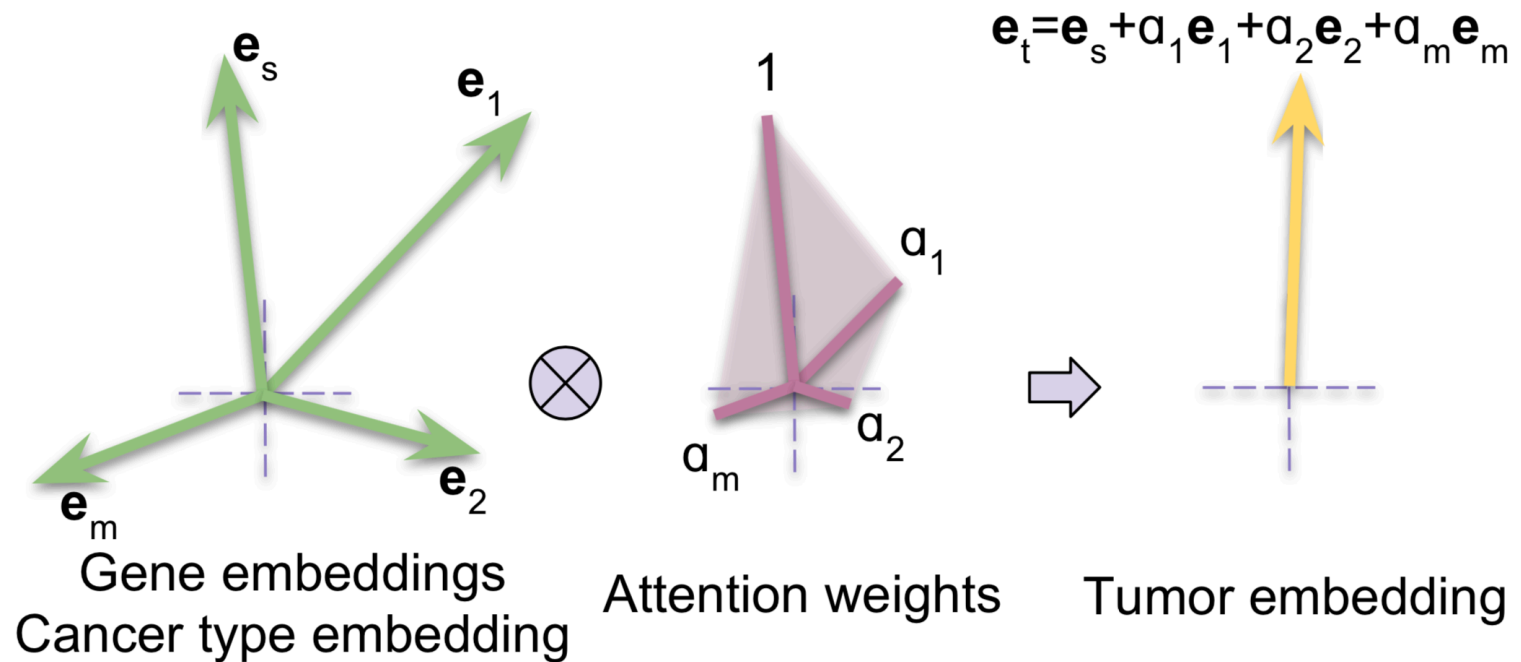
# Genomic Impact Transformer (GIT)

o GIT: encoder-decoder architecture

o Mimic cellular signaling process

o Driver SGA detection

   o Problem: downstream effect of SGAs

   o Solution: supervised by gene expressions

o SGA/tumor representation

   o Problem: little information/knowledge

   o Solution: gene/tumor embedding

gene expression

↑

decoder

↑

tumor embedding

↑

encoder

↑

gene embeddings

↑

SGA

# Genomic Impact Transformer (GIT)

# Encoder: Attention



$$e_t = e_s + \alpha_1 e_1 + \alpha_2 e_2 + \alpha_m e_m$$

Gene embeddings
Cancer type embedding

Attention weights

Tumor embedding

# Encoder: Attention

$$\boxed{\alpha_1, \alpha_2, ..., \alpha_m} = \boxed{\text{Function}_{\text{Attention}}} (\boxed{\mathbf{e}_1, \mathbf{e}_2, ..., \mathbf{e}_m})$$

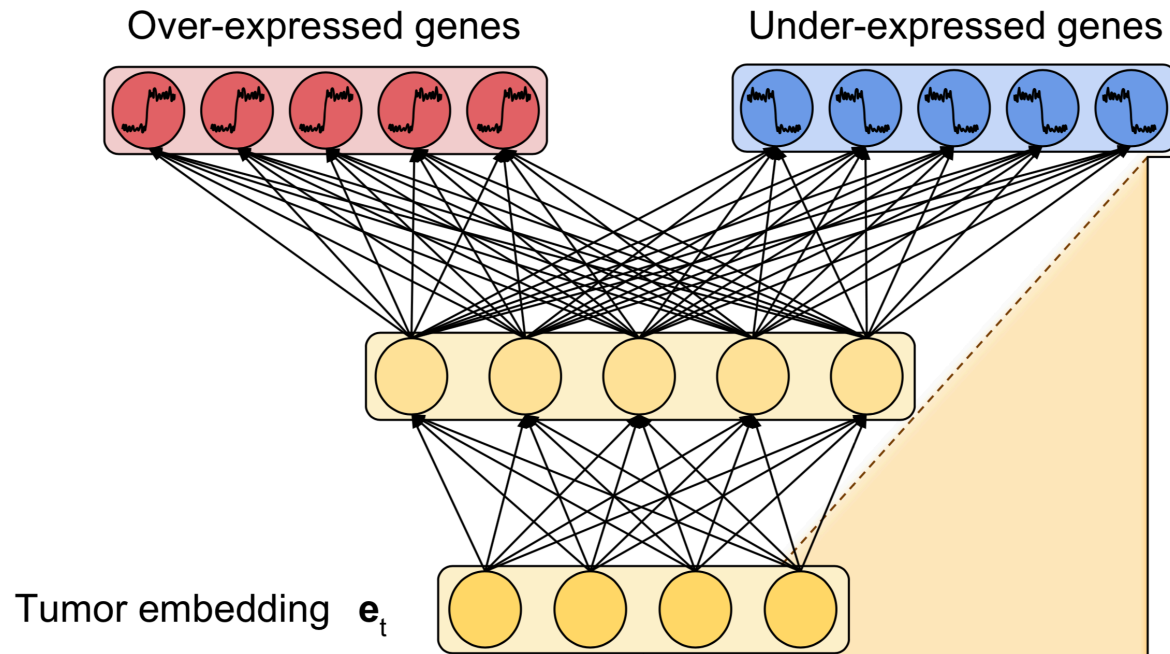$$\alpha_g = \sum_{j=1}^{h} \alpha_{g,j} = \alpha_{g,1} + \alpha_{g,2} + ... + \alpha_{g,h}, \; g = 1, 2, ..., m$$

$$\alpha_{1,j}, \alpha_{2,j}, ..., \alpha_{m,j} = \text{softmax}(\beta_{1,j}, \beta_{2,j}, ..., \beta_{m,j})$$

$$\beta_{g,j} = \boldsymbol{\theta}_j^{\mathsf{T}} \cdot \tanh(W_0 \cdot \mathbf{e}_g), \; g = 1, 2, ..., m$$

# Decoder: MLP

$$\hat{y} = \sigma(W_2 \cdot \mathrm{ReLU}(W_1 \cdot \mathrm{ReLU}(\mathbf{e}_t) + b_1) + b_2)$$
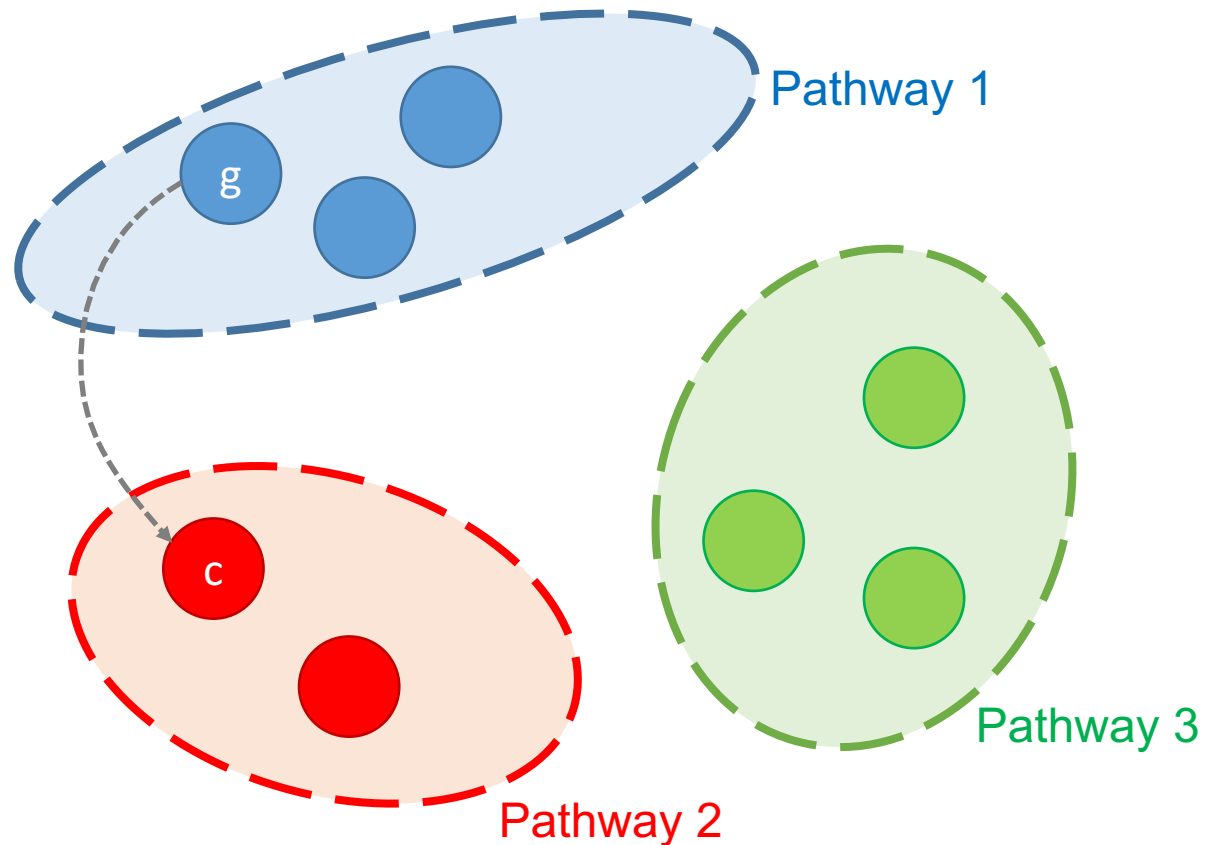


Over-expressed genes        Under-expressed genes

Tumor embedding  $\mathbf{e}_t$

# Pre-training Gene Embedding: Gene2Vec

○ Co-occurrence pattern

$$\Pr\left(c \in \mathrm{Context}(g) \mid g\right) = \frac{\exp\left(\mathbf{e}_g^{\mathsf{T}} \mathbf{v}_c\right)}{\sum_{c' \in \mathcal{G}} \exp\left(\mathbf{e}_g^{\mathsf{T}} \mathbf{v}_{c'}\right)}$$
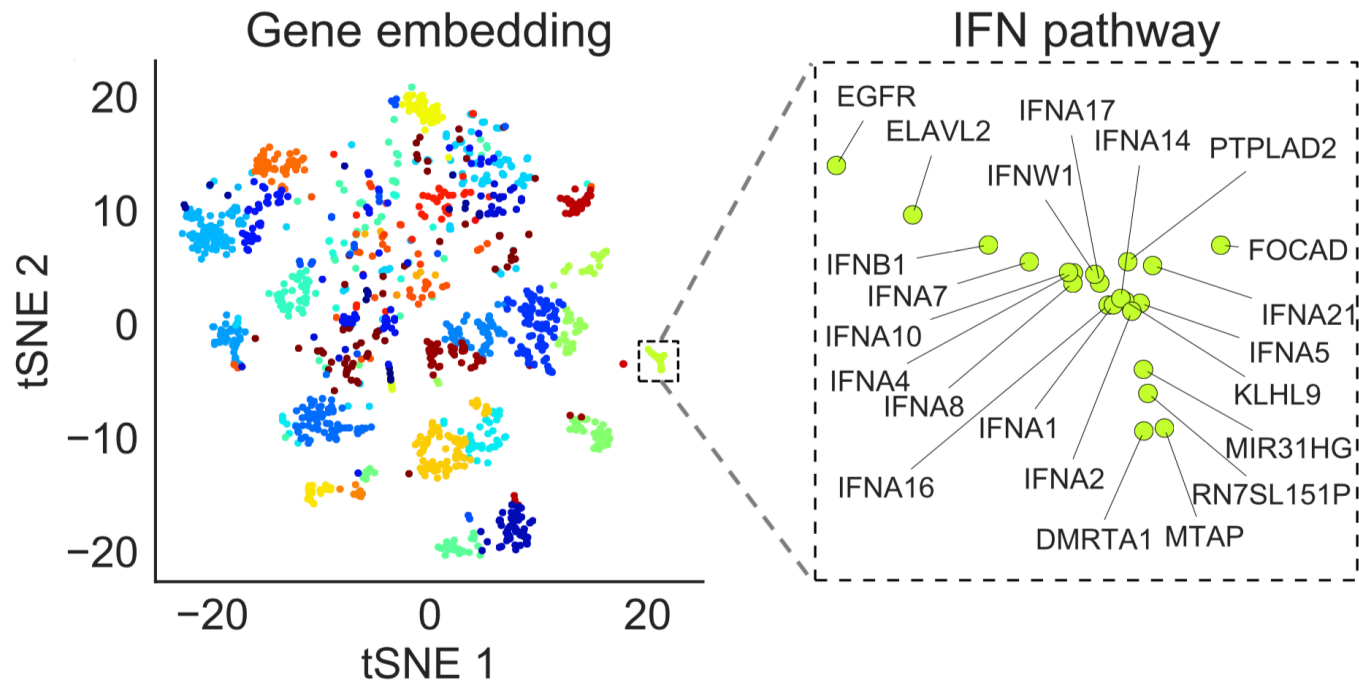


Pathway 1

Pathway 2
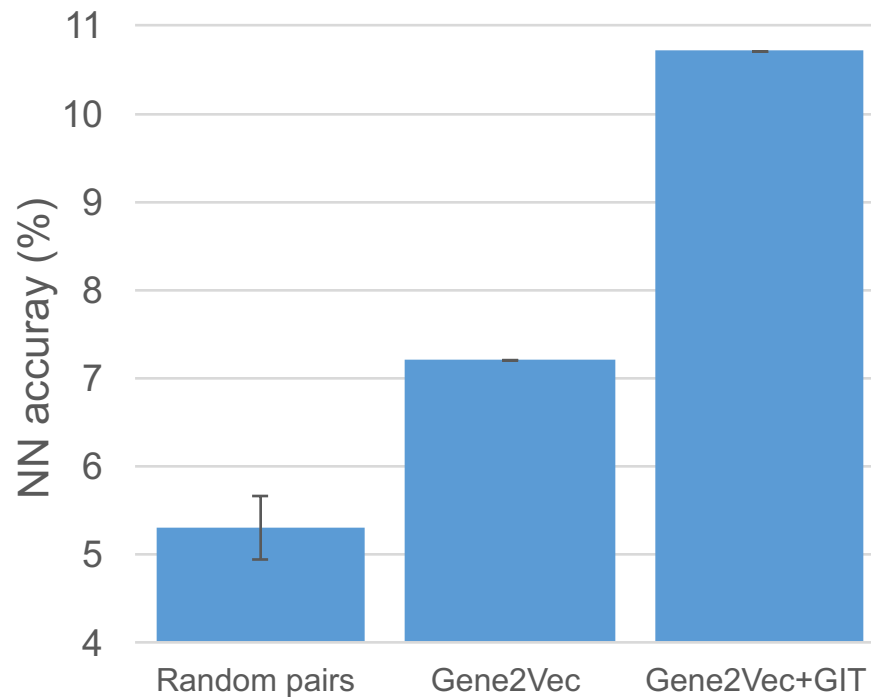
Pathway 3

# Performance

o Predicting gene expression using SGAs

# Gene Embedding
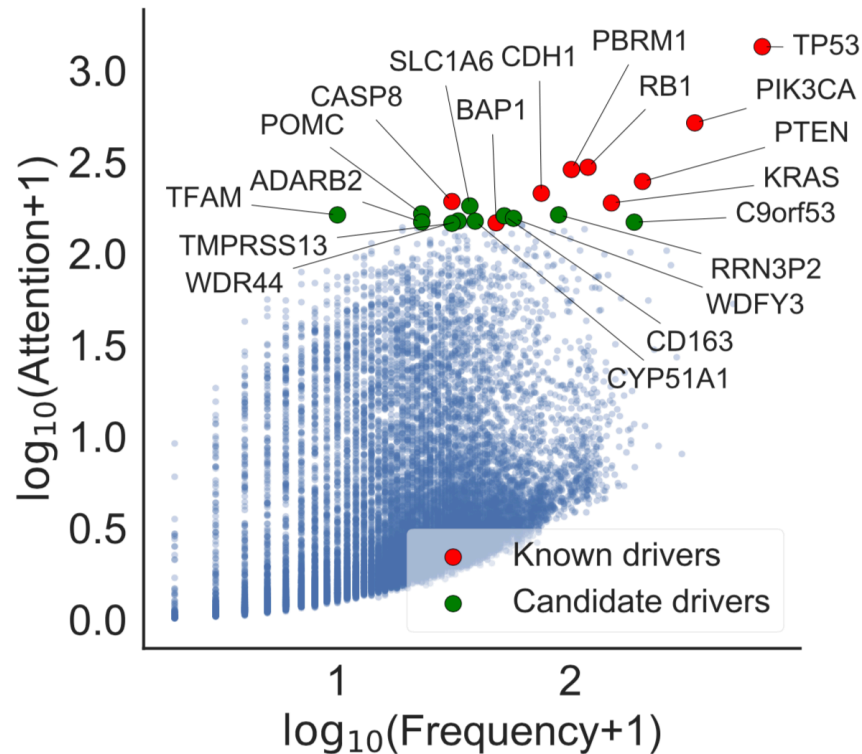
# Gene Embedding

o NN accuracy: functional similar genes → closer in embedding space

$$\text{NN accuracy} = \mathbb{E}_{\mathbf{e}_c \in \text{NN}(\mathbf{e}_g)} \left[ \mathbb{1} \left( \text{GO}(g) \cap \text{GO}(c) \neq \emptyset \right) \right]$$
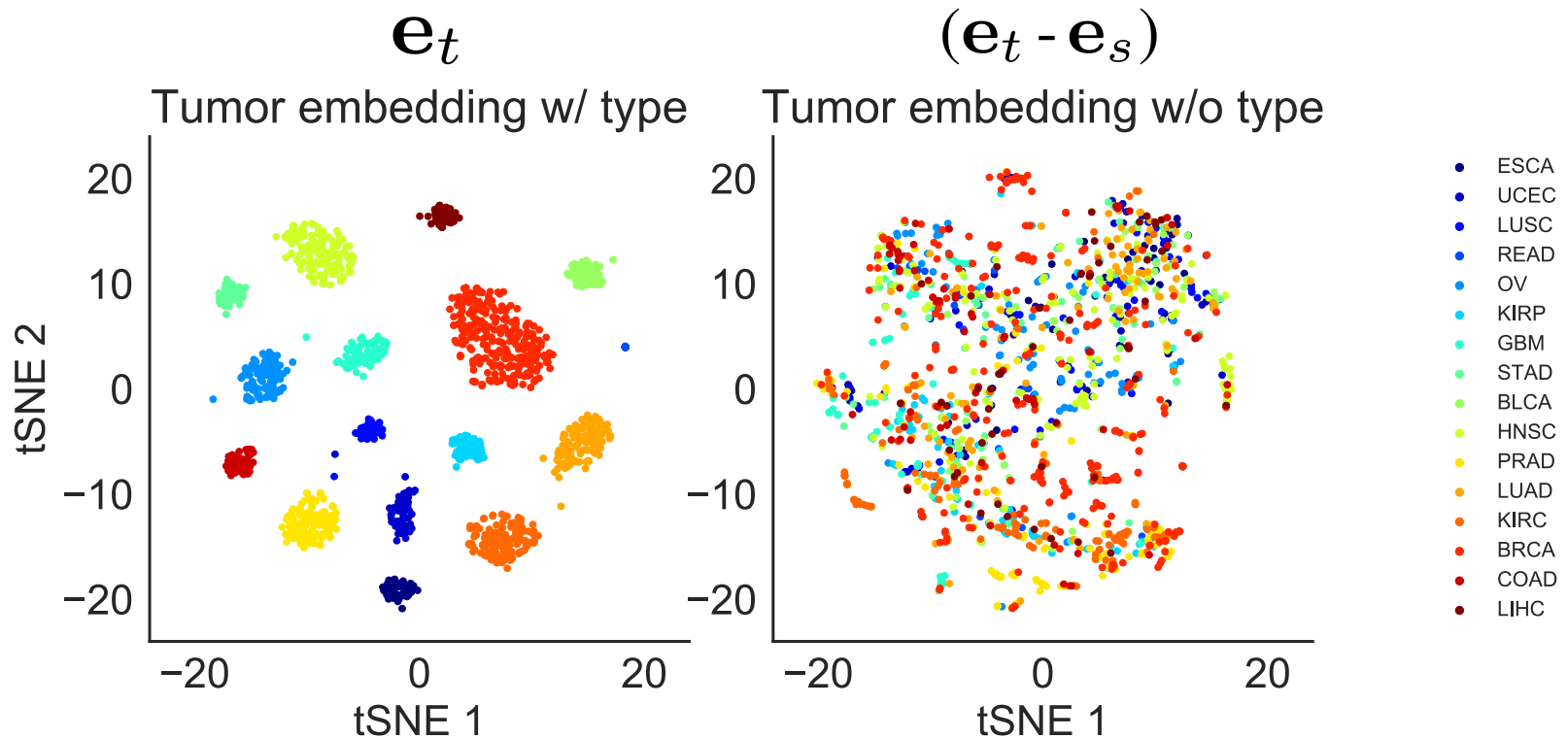
# Candidate Drivers via Attention

# Tumor Embedding

o Common cellular signaling process across cancer types

$$\mathbf{e}_t \qquad (\mathbf{e}_t - \mathbf{e}_s)$$



Tumor embedding w/ type

Tumor embedding w/o type

Legend: ESCA, UCEC, LUSC, READ, OV, KIRP, GBM, STAD, BLCA, HNSC, PRAD, LUAD, KIRC, BRCA, COAD, LIHC

# Application - Survival Analysis

# Application – Drug Response

# Summary

o Biological-inspired neural network to mimic cellular signaling

o Distinguish drivers from passengers with supervision of expression

o Gene embedding: informative of gene functions

o Tumor embedding: transferable to other phenotype prediction tasks

o Gene2Vec: speed up training and alleviate overfitting

# Acknowledgements