

# Learning Feature Importance for a Deep Learning Cancer Classifier

*Jacob Chmura*

*Wednesday, July 31, 2019*

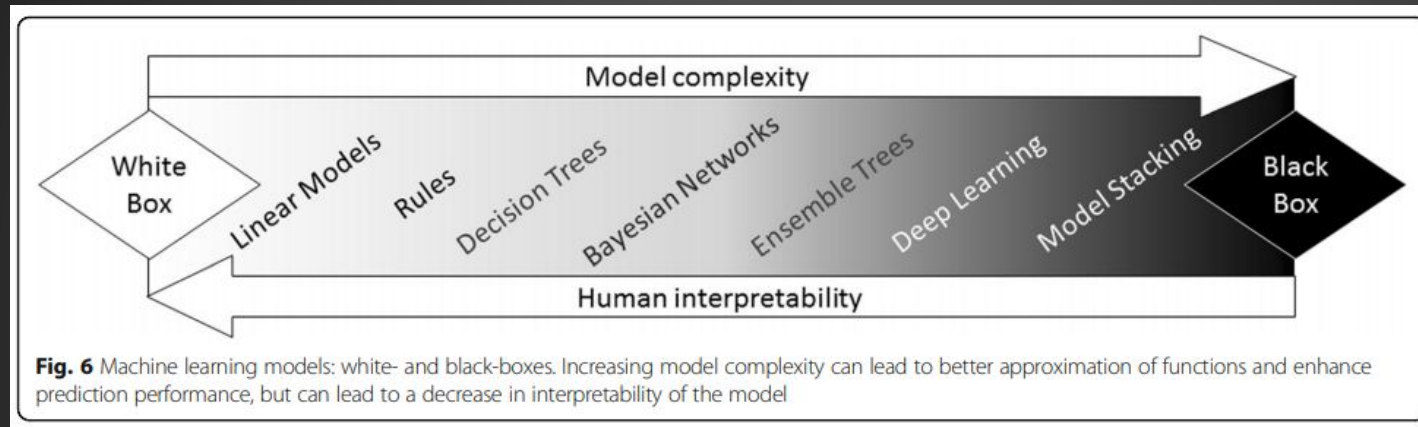
*Supervisor: Quaid Morris*

University of Toronto Undergraduate Summer Research Program

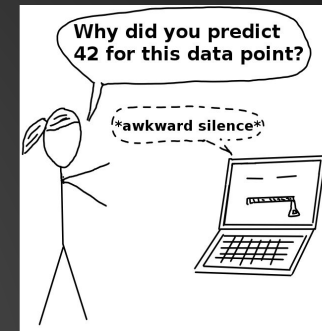


VECTOR INSTITUTE

# Lack of Interpretability in Deep Learning is a Limitation



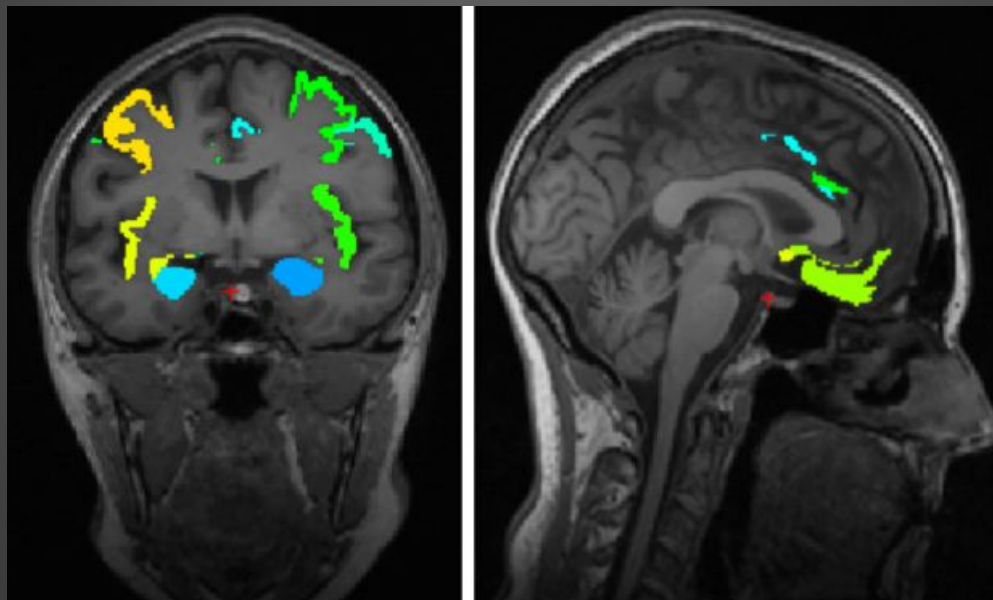
<https://onehealth.ifas.ufl.edu/media/onehealthifasufledu/pdfs/publications/Prosperi-article.pdf>



- *Object recognition network*: could tell us which pixels of the image responsible for a label being picked

## Lack of Interpretability in Deep Learning is a Limitation

- *Medical imaging model*: could help inform the doctor of the part of the image that resulted in the recommendation. Knowing the strengths and weaknesses of a model is essential in clinical settings.



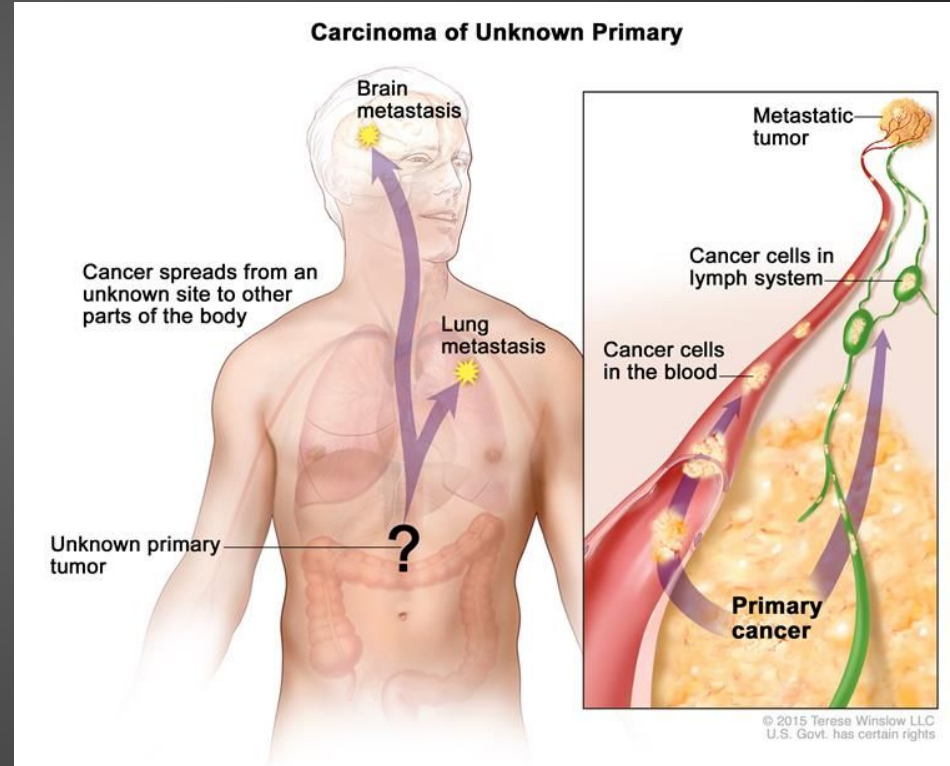
# Classifying Primary and Metastatic Cancers based on Mutation Patterns using Deep Learning Techniques

## *Situation:*

- pathologist can't identify primary cancer tumor well
- at autopsy, the primary cannot be identified roughly 70% of the time
- fourth most common cause of cancer death
- therapeutic options are driven by tissue of origin

## *Idea:*

- Different tumour types have dramatically different patterns of mutation
- Use this to train a deep network to classify cancer classes



# Features for the Model

~ 5000 length feature vector split into three main categories:

- (1) *Mutation Distribution* (~3000): encode information about cell type
- divided genome into bins and created features corresponding to the number of mutations per bin

(2) *Mutation Types* (~150): reflect environmental exposures of the cell of origin

- Ex. skin cancers have mutation types strongly correlated with UV light-induced DNA damage.

(3) *Driver Genes* (~2000): tumour types are distinguished by high frequencies of alterations in particular driver genes and pathways

Wild type gene	C G A C T G G C T G A C
Transition (AT pair replaced by GC pair)	C G <u>G</u> C T G G C <u>C</u> G A C
Transversion (AT pair replaced by TA pair)	C G <u>T</u> C T G G C <u>A</u> G A C
Insertion (GC pair inserted)	C G A <u>G</u> C T G G C T <u>C</u> G A C
Deletion (AT pair deleted)	C G C T G G C G A C

[http://www.biochem.uthscsa.edu/med/06-Mechanisms-of-Mutation/PreqMechanismsofMutation\\_print.html](http://www.biochem.uthscsa.edu/med/06-Mechanisms-of-Mutation/PreqMechanismsofMutation_print.html)

# Results

True Class	Predicted Class																							
	Kidney-RCC	Skin-Melanoma	Liver-HCC	Breast-AdenoCA	ColoRect-AdenoCA	Ovary-AdenoCA	Lymph-BNHL	Panc-AdenoCA	Panc-AdenoCA	Myeloid-MPN	CNS-Medullo	CNS-GBM	Panc-Endocrine	Head-SCC	Lung-SCC	Lymph-CLL	Eso-AdenoCA	Thy-AdenoCA	Kidney-ChRCC	CNS-PiloAstro	Uterus-AdenoCA	Lung-AdenoCA	Bone-Osteosarc	Stomach-AdenoCA
Kidney-RCC (143)	99	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0
Skin-Melanoma (106)	0	98	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Liver-HCC (306)	1	0	98	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Breast-AdenoCA (198)	0	0	0	96	0	1	0	1	0	1	1	0	0	0	0	0	0	1	1	0	0	1	0	0
ColoRect-AdenoCA (52)	0	0	0	0	96	0	0	0	2	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0
Ovary-AdenoCA (112)	0	0	0	1	0	96	0	0	0	0	0	0	0	0	0	0	0	0	1	0	3	0	0	0
Lymph-BNHL (105)	0	0	0	0	0	0	95	0	0	0	0	0	0	0	0	5	0	0	0	0	0	0	0	0
Prost-AdenoCA (189)	1	0	0	1	0	0	0	94	0	0	2	0	0	0	0	0	0	0	0	4	0	0	0	0
Panc-AdenoCA (235)	0	0	0	1	0	0	0	0	94	0	0	0	0	0	0	0	0	0	0	0	0	1	0	2
Myeloid-MPN (46)	0	0	0	0	0	0	0	2	0	0	93	2	0	0	0	0	0	0	0	2	0	0	0	0
CNS-Medullo (146)	1	0	0	1	0	0	0	0	0	1	93	0	1	0	0	0	0	0	0	4	0	0	0	0
CNS-GBM (41)	2	0	0	0	0	0	0	0	0	5	93	0	0	0	0	0	0	0	0	0	0	0	0	0
Panc-Endocrine (85)	0	0	1	1	0	0	0	2	0	1	0	89	1	0	0	0	1	0	0	0	1	0	1	1
Head-SCC (57)	0	0	0	7	0	0	0	0	0	0	0	0	88	4	0	0	2	0	0	0	0	0	0	0
Lung-SCC (48)	0	0	0	2	0	0	0	0	0	0	0	0	2	88	0	0	0	0	0	0	8	0	0	0
Lymph-CLL (95)	0	0	0	0	0	0	12	0	0	0	0	0	1	0	0	87	0	0	0	0	0	0	0	0
Eso-AdenoCA (98)	0	0	0	0	1	0	0	0	3	0	0	0	0	1	0	0	84	0	0	1	0	1	0	9
Thy-AdenoCA (48)	0	0	0	4	0	0	0	2	0	0	0	0	8	0	0	0	0	81	2	2	0	0	0	0
Kidney-ChRCC (45)	2	0	0	2	0	0	0	0	2	0	0	0	9	0	0	0	0	2	80	0	0	0	0	2
CNS-PiloAstro (89)	0	0	0	0	0	0	0	3	0	3	11	0	1	0	0	1	0	0	0	80	0	0	0	0
Uterus-AdenoCA (40)	2	0	0	5	0	10	0	0	0	0	0	0	0	2	2	0	0	0	0	0	78	0	0	0
Lung-AdenoCA (38)	3	0	0	0	0	0	0	3	5	0	0	0	3	0	3	0	0	3	0	5	0	74	0	3
Bone-Osteosarc (44)	0	0	0	2	0	2	0	2	2	5	0	0	7	0	0	0	0	0	0	5	0	0	73	2
Stomach-AdenoCA (70)	0	0	0	1	4	0	0	3	9	0	0	0	1	1	0	0	14	0	1	1	0	0	1	61

- classifier achieves an accuracy of 91%
- most frequent classification errors for Stomach-AdenoCA samples were two other upper gastrointestinal tumours, (Eso-AdenoCA and Panc-AdenoCA)



## Goals

- *Aim 1*: take existing model and assess feature importance
- *Aim 2*: produce certainty estimates and extend to rare cancer types

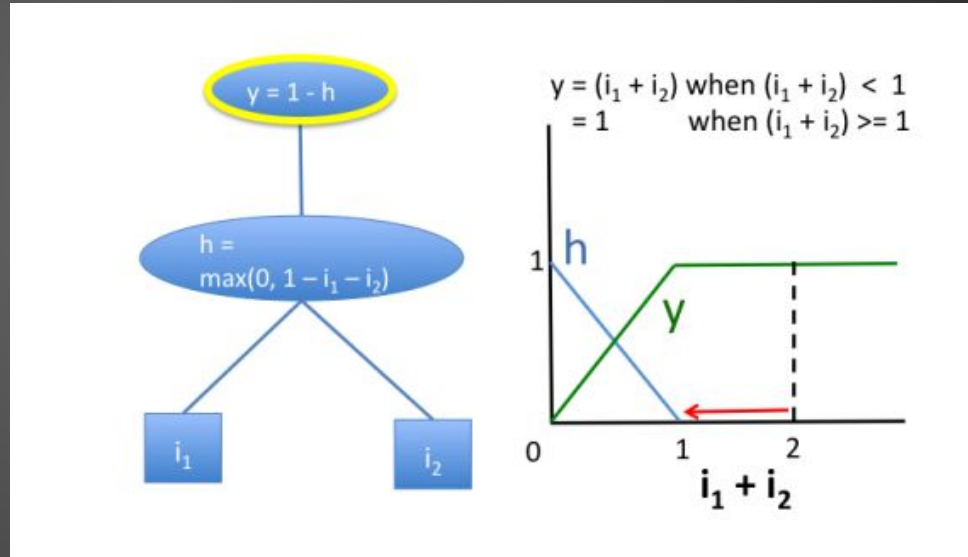
# How do we learn which features are important?

## Idea 1: Perturbations?

- Make small changes to individual inputs and observe the impact on later neurons in the network

## *Problems:*

- computationally inefficient
- **Saturation problem**



(Avanti Shrikumar, 2017)

**Saturation problem:** lack of local change need not imply zero importance



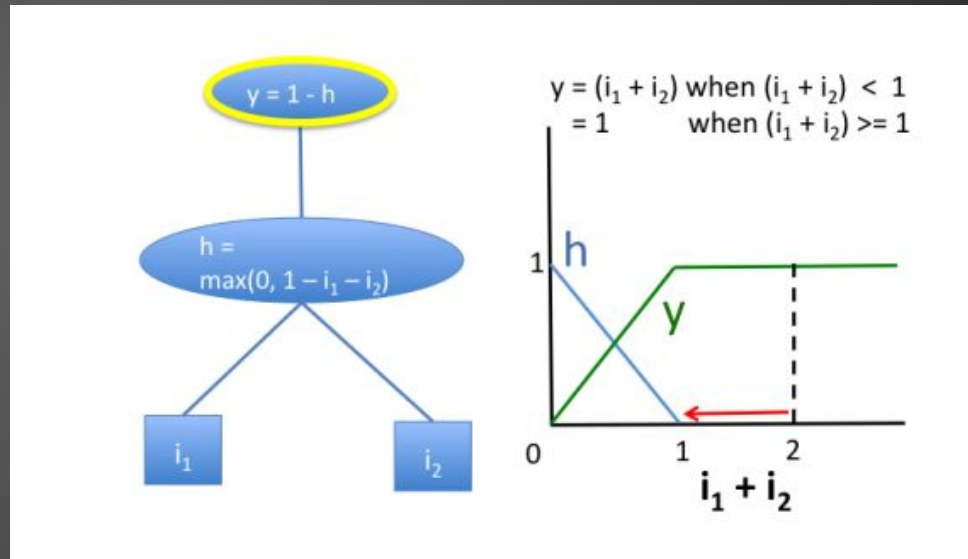
# How do we learn which features are important?

## Idea 2: Gradients?

- Compute the gradient of the outputs with respect to each input feature

## *Problems:*

- nonlinearities
- **Saturation problem**



(Avanti Shrikumar, 2017)

**Saturation problem:** lack of local change need not imply zero importance

# DeepLIFT (Deep Learning Important FeaTures) (Avanti Shrikumar, 2017)

*Philosophy:* explain the difference in output from some 'reference' output in terms of the difference of the input from some 'reference' input.

reference input: represents default or *neutral* baseline that is chosen

- references for all neurons can be found by choosing a reference input and propagating activations through the net

- ❑  $t$  represents target output neuron of interest
- ❑  $x_1, x_2, \dots, x_n$  represent neurons in some intermediate layer
- ❑  $t_0$  represent the reference activation of  $t$ .
- ❑  $\Delta t$  is difference-from-reference, that is  $\Delta t = t - t_0$ .

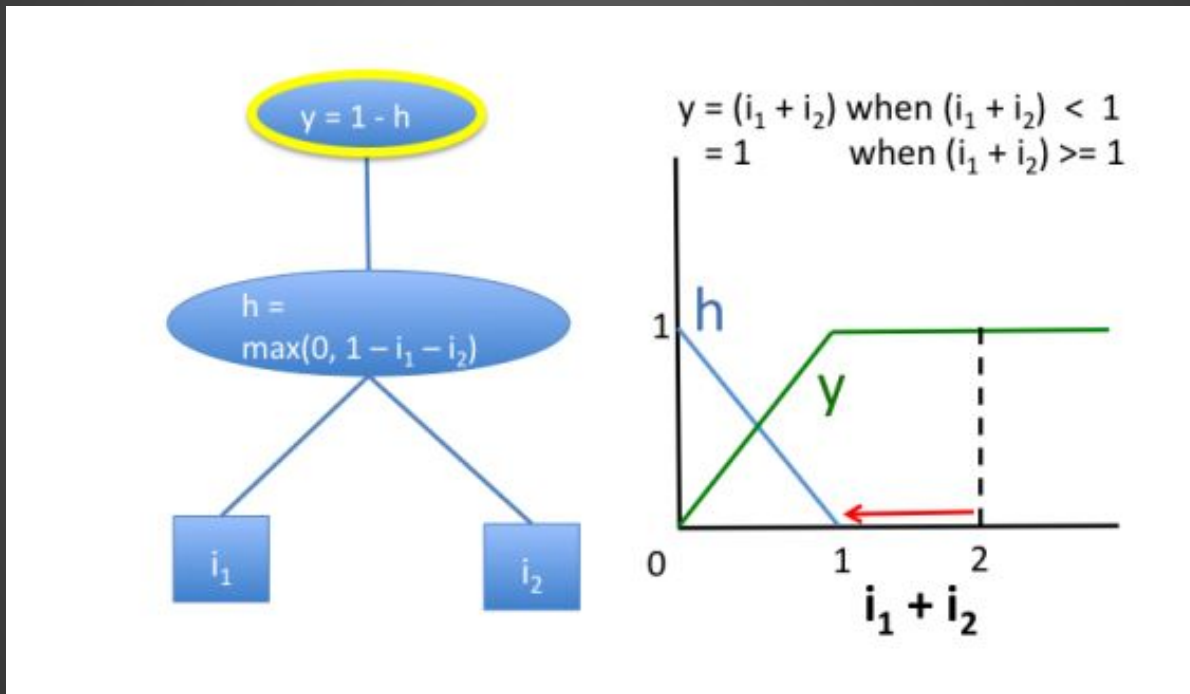
$$\sum_{i=1}^n C_{\Delta x_i \Delta t} = \Delta t$$

DeepLIFT assigns contribution scores:

- $C_{\Delta x_i \Delta t}$  is the amount in  $t$  that is *blamed* on difference-from-reference of  $x_i$ .

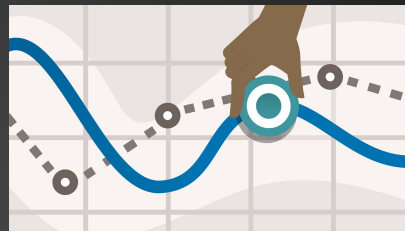
## Saturation Problem Revisited

- $C\Delta x_i\Delta t$  can be non-zero even when  $\partial t / \partial x_i$  is zero: a neuron can be signaling meaningful information even when its gradient is zero.



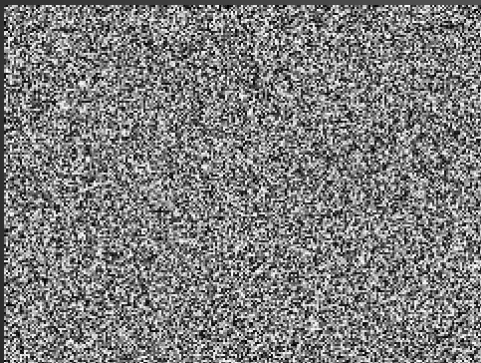
- DeepLIFT is sensitive to chosen reference

## Important Question: How do we choose a good reference?



*Intuition: ask ourselves what am I interested in measuring differences against?*

- needs to convey a complete absence of signal (allows us to interpret the attributions as a function of the input)



black image or noisy image signifies the absence of objects



zero embedding vector is a good baseline

# Choosing a reference for Cancer Classifier

~ 5500 length feature vector split into three main categories:

Mutation Distribution (~3000)

Mutation Type (~150)

Drivers (~2000)

★ Typical Feature Vector is extremely sparse

## Null Reference

- A vector of all zeros

0 0 0 ... 0 0 0 | 0 0 0 ... 0 0 0 | 0 0 0 ... 0 0 0

## Shuffle Reference

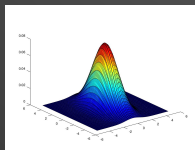
- permute within the three main sections of the input vector and averaging the results over multiple such references



1.1 -1.6 ... 0.7 0.2 | 0.0 0.5 ... 1.1 -0.7 | 2.9 1.3 ... 2.1 0.2

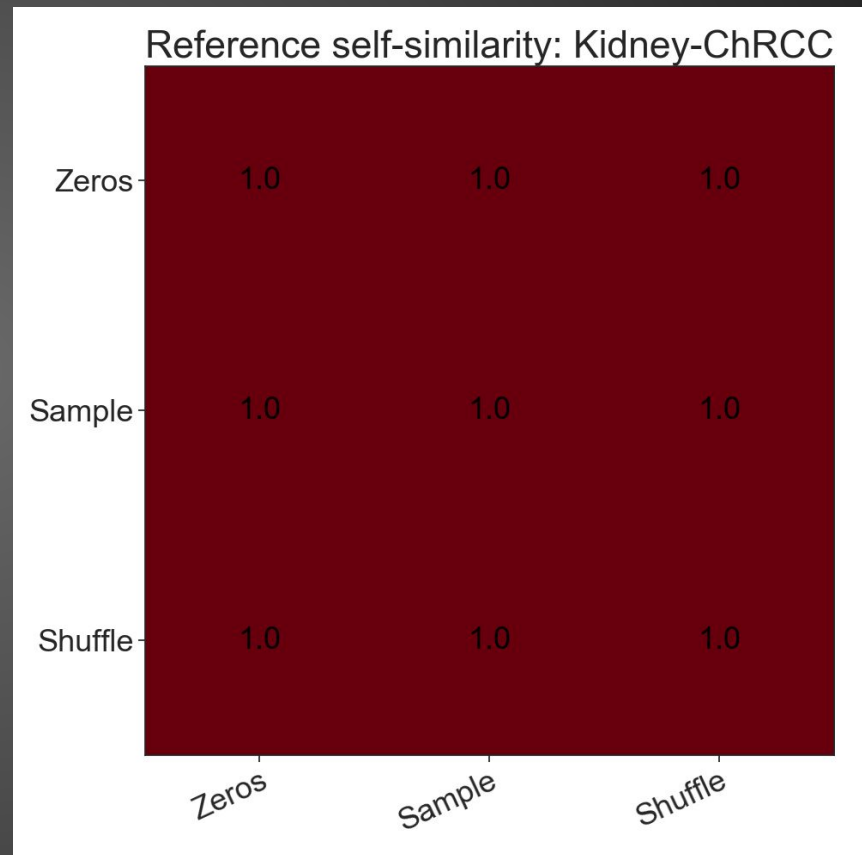
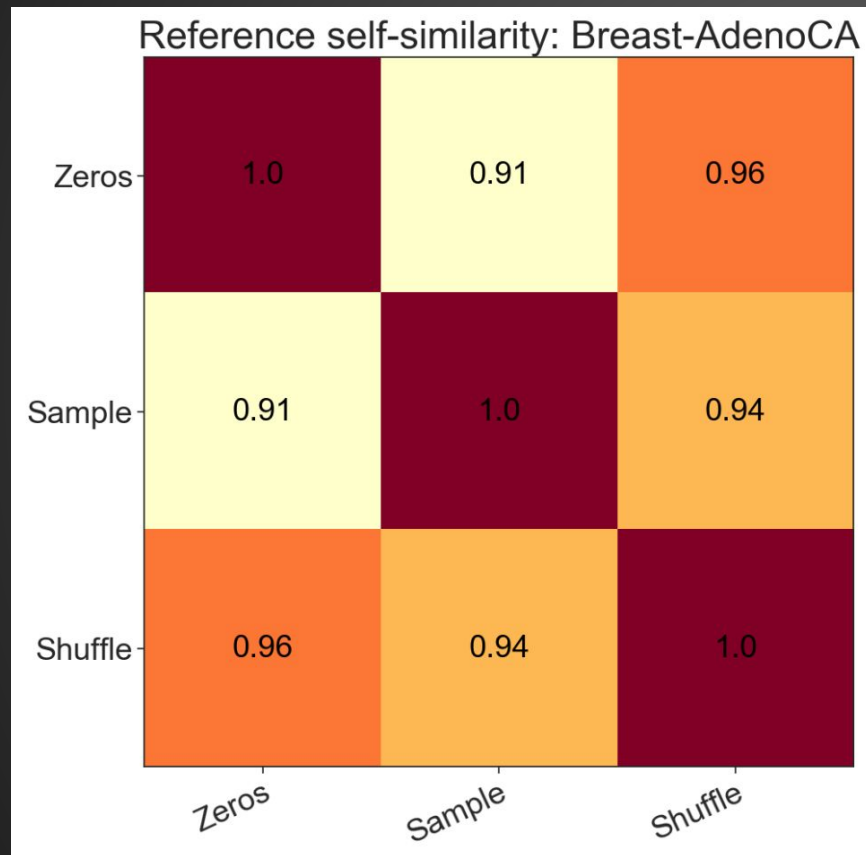
## Sample Reference

- For each z score feature, simply sample from a Gaussian
- For count features, create a distribution of frequencies and sample from it

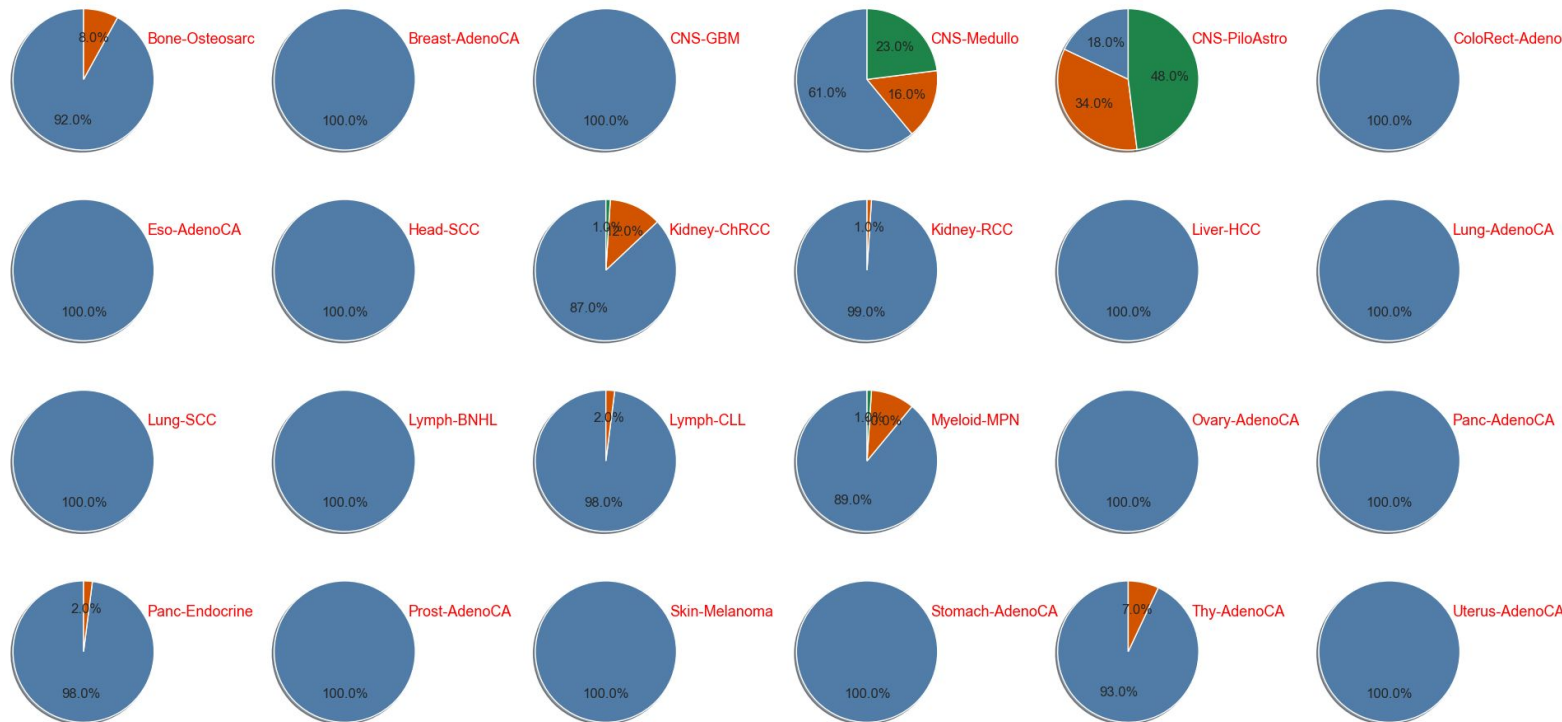


2.1 -0.6 ... -1.7 0.1 | 0.0 3.5 ... 0.1 0.7 | 1.9 1.3 ... 2.1 0.2

### Reference Self-Similarity Matrix

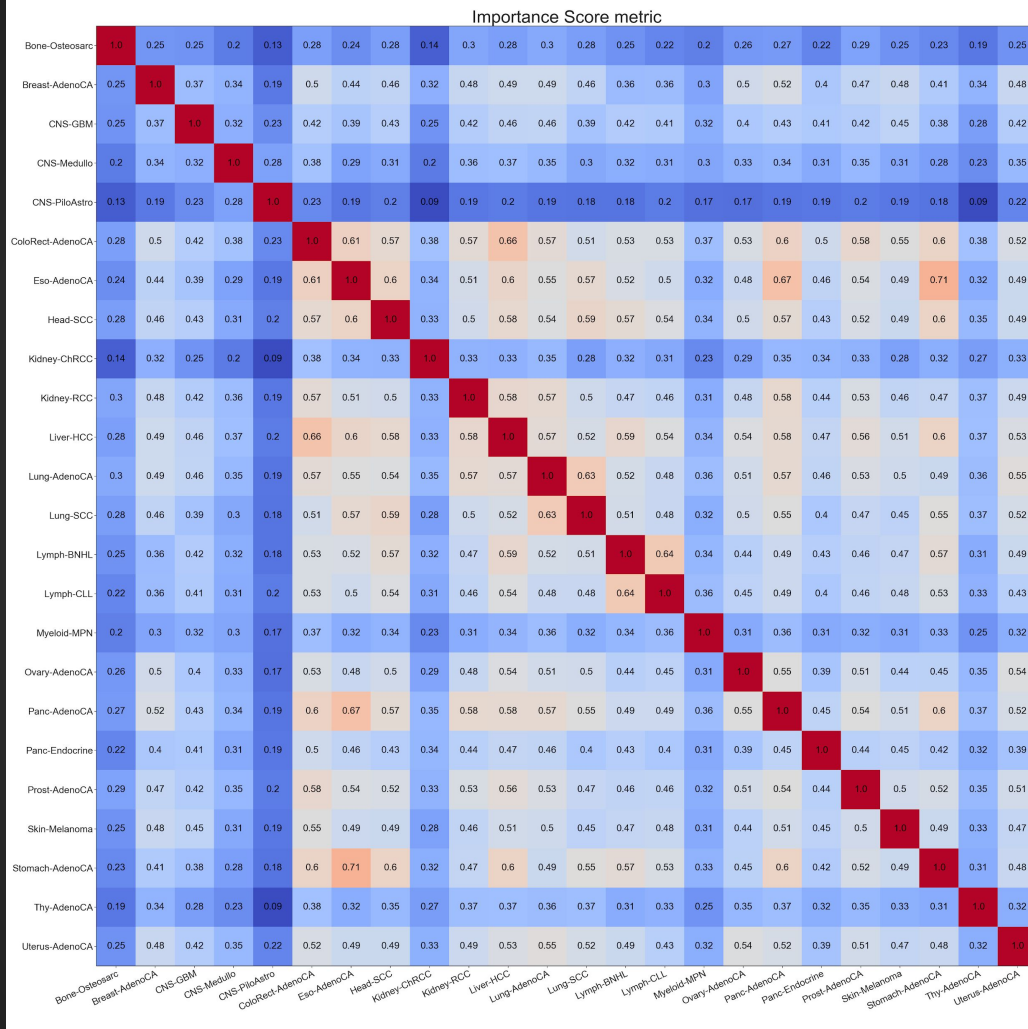


# Feature Type Distribution over 100 Most Significant Inputs



Retraining without driver features, the classifier does just as good

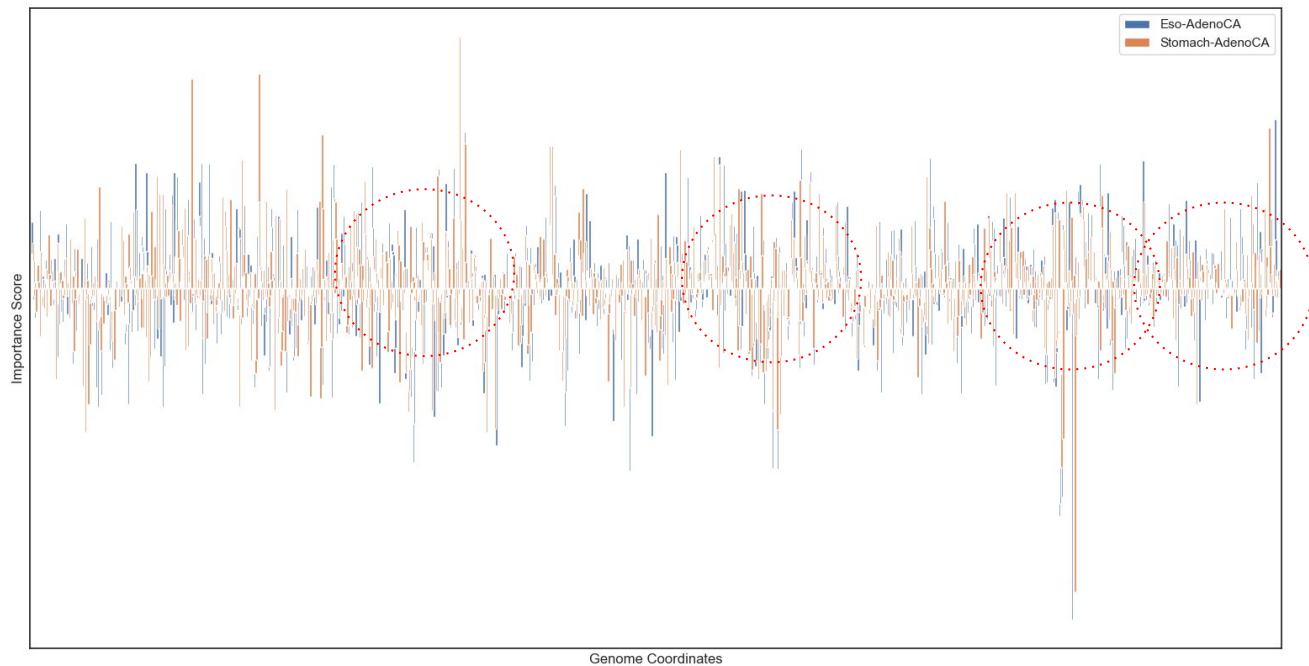




# Cancer class Self-Similarity over 100 Most Significant Input Features

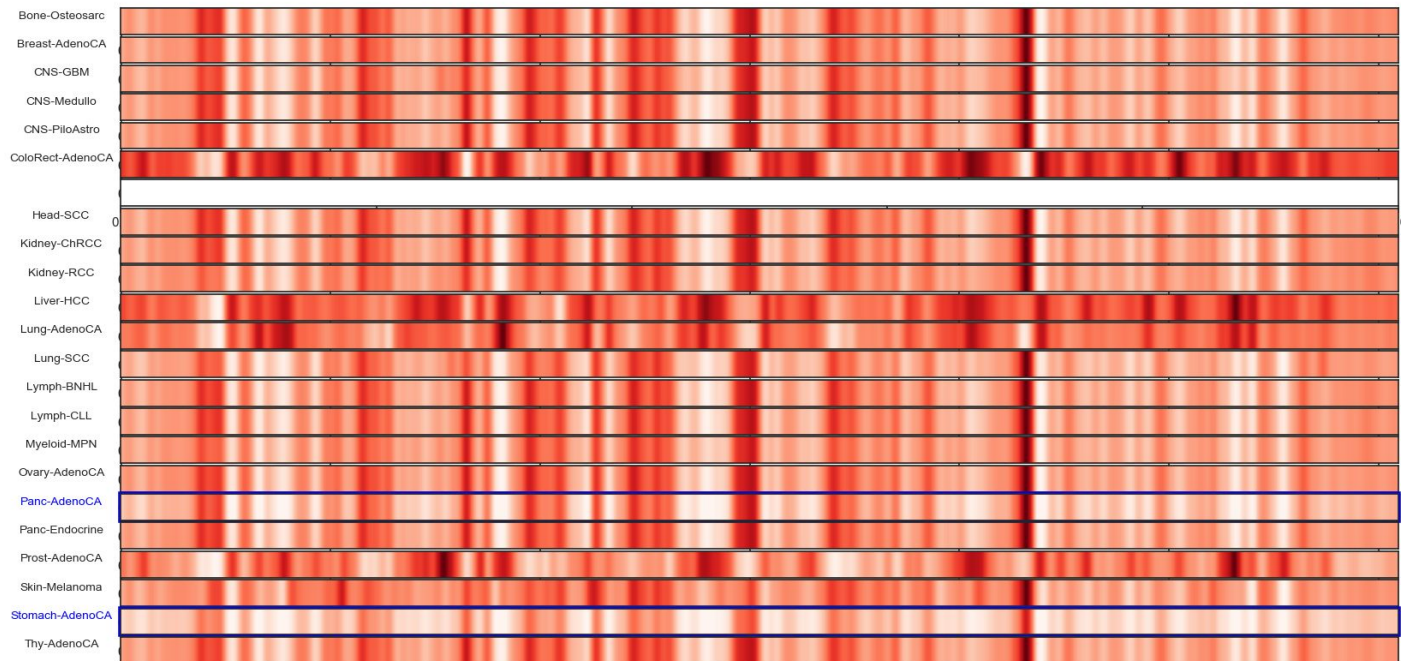
- Similarity score computed using adaptation of Jaccard Index over feature importance
- No drivers
- Average of all 3 references

## Genome Wide Mutation Importance: Chromosome 11-15

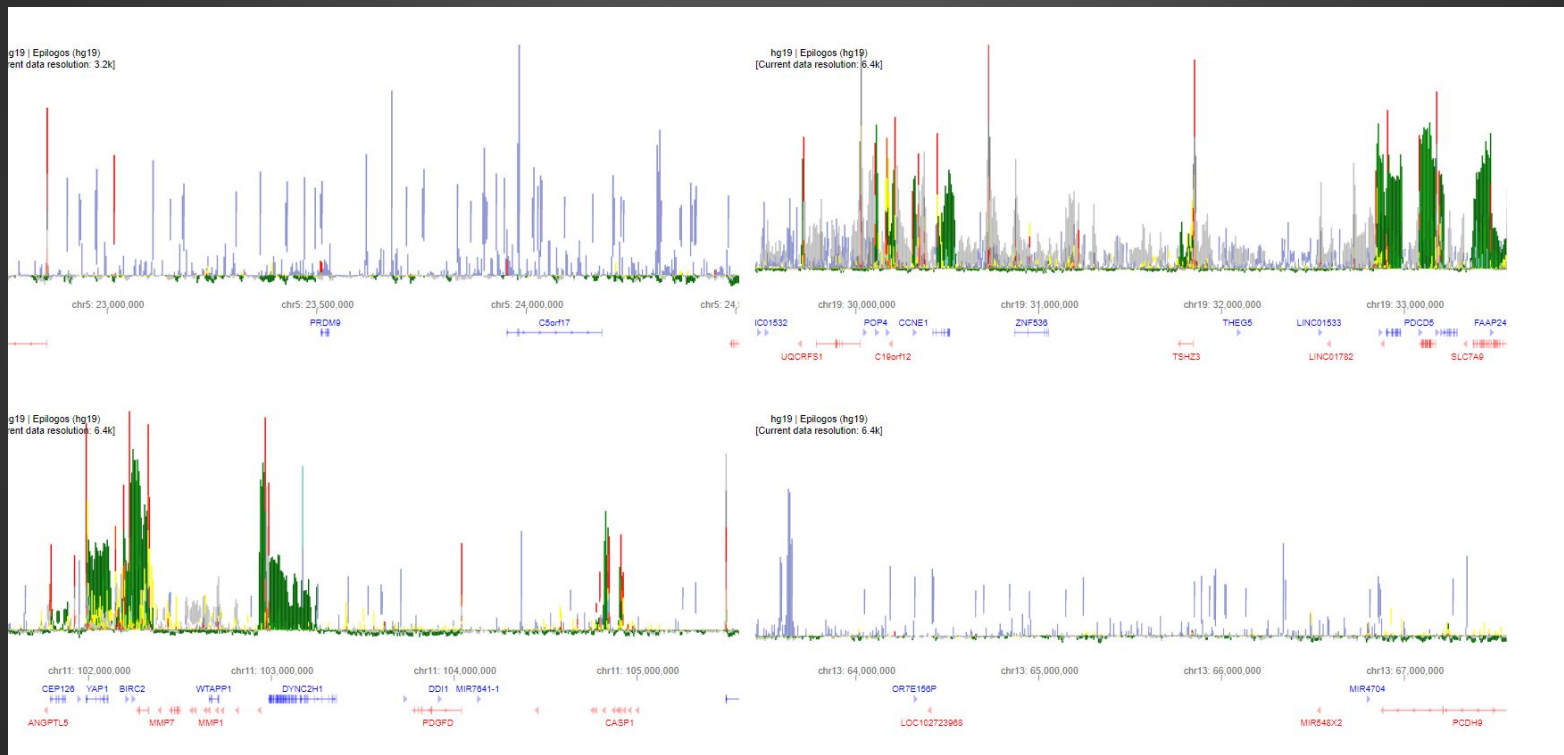


# Convolution of Importance Scores Across Genome

Convolution over the Genome Eso-AdenoCA



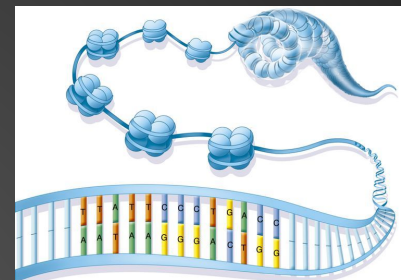
# Genome Annotation Process: *Extract genes from most important DeepLIFT scoring genome regions*



*Bone-Osteosarc Important Regions*

# Chromatin States Extraction Controlled for Cell Type

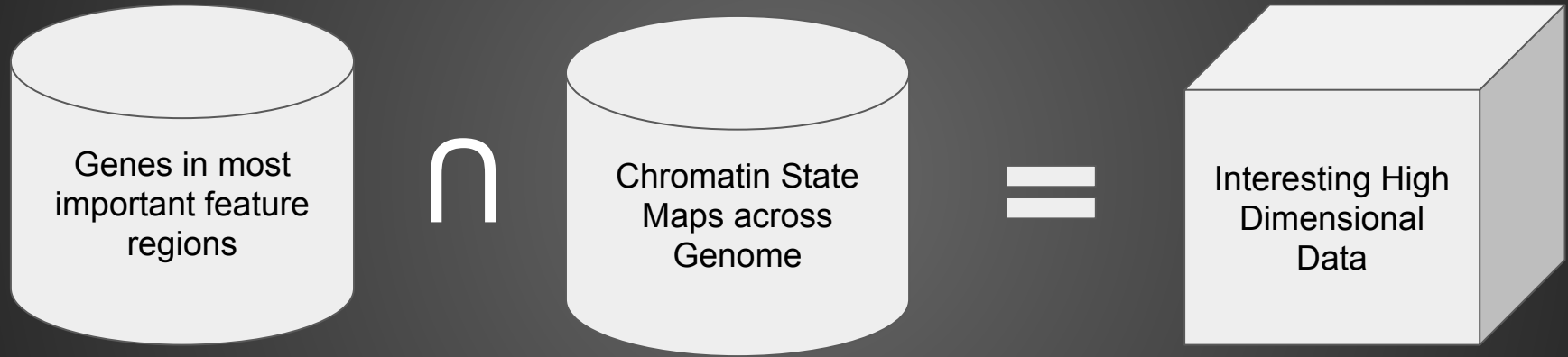
- Chromatin is understood to be a complex genome organizer
- Map chromatin function along the genome to discrete states



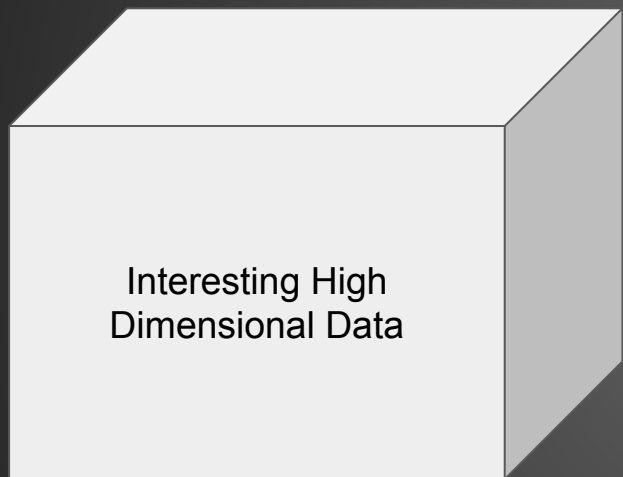
*Does it transcribe?*

STATE NO.	MNEMONIC	DESCRIPTION	COLOR NAME	COLOR CODE
1	TssA	Active TSS	Red	255,0,0
2	TssAFlnk	Flanking Active TSS	Orange Red	255,69,0
3	TxFlnk	Transcr. at gene 5' and 3'	LimeGreen	50,205,50
4	Tx	Strong transcription	Green	0,128,0
5	TxWk	Weak transcription	DarkGreen	0,100,0
6	EnhG	Genic enhancers	GreenYellow	194,225,5
7	Enh	Enhancers	Yellow	255,255,0
8	ZNF/Rpts	ZNF genes & repeats	Medium Aquamarine	102,205,170
9	Het	Heterochromatin	PaleTurquoise	138,145,208
10	TssBiv	Bivalent/Poised TSS	IndianRed	205,92,92
11	BivFlnk	Flanking Bivalent TSS/Enh	DarkSalmon	233,150,122
12	EnhBiv	Bivalent Enhancer	DarkKhaki	189,183,107
13	ReprPC	Repressed PolyComb	Silver	128,128,128
14	ReprPCWk	Weak Repressed PolyComb	Gainsboro	192,192,192
15	Quies	Quiescent/Low	White	255,255,255

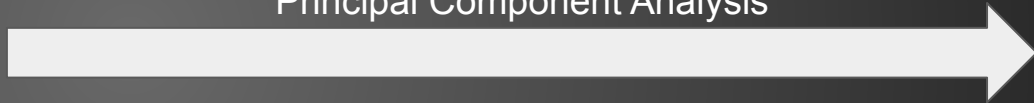
Vector Space Embedding: *We want to control for Suppressed Chromatin Markers*



## Vector Space Embedding: *I want to see what this looks like ...*

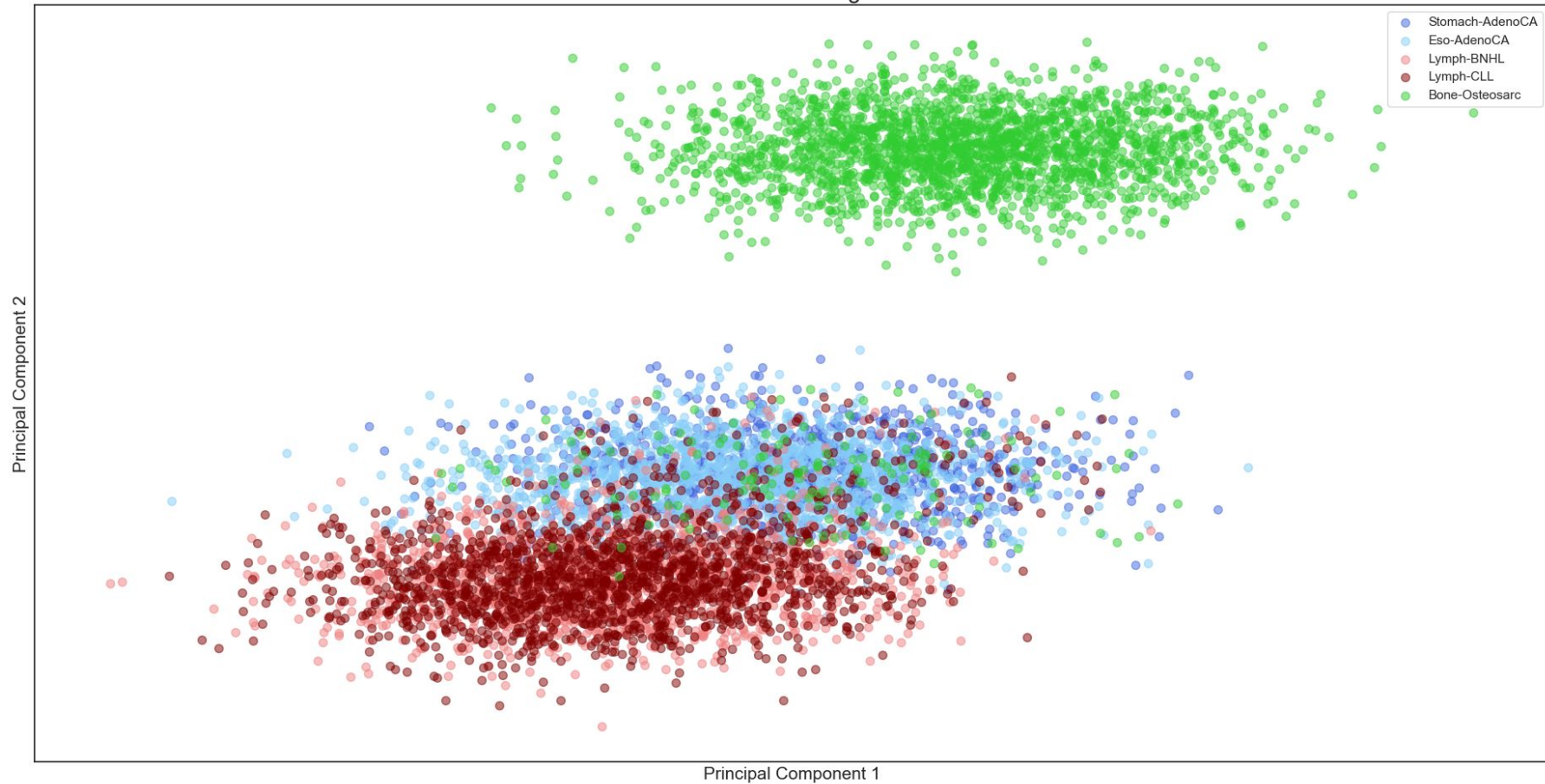


Principal Component Analysis





Vector Embedding PCA



## Conclusions

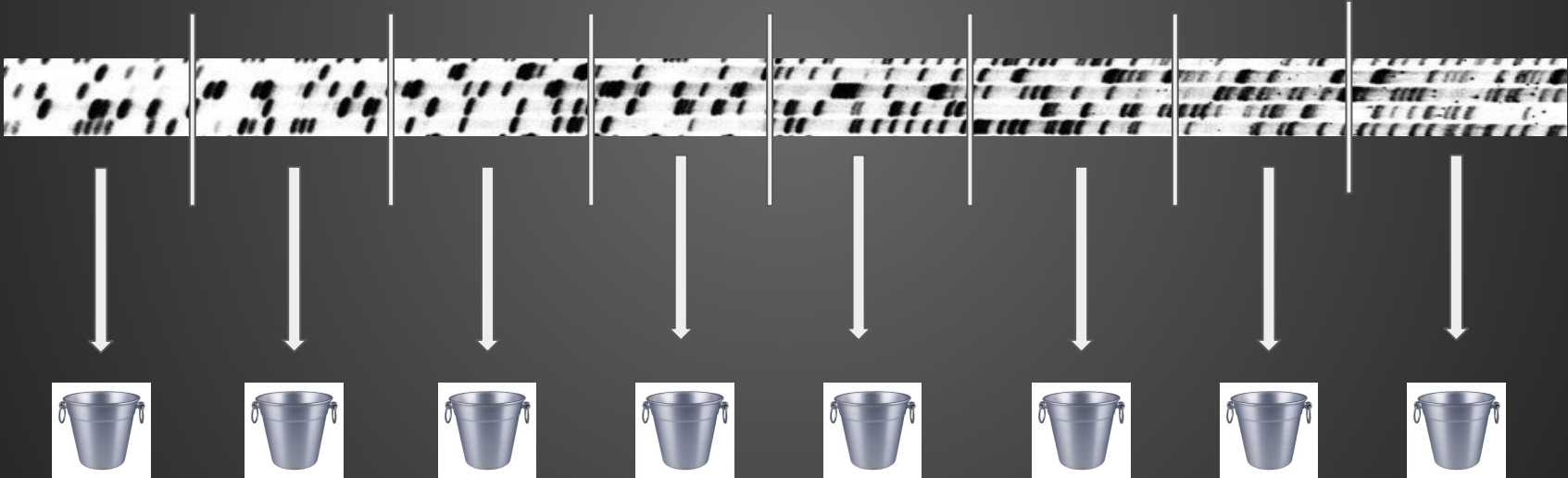
- Driver gene and pathway features don't provide additional information for cancer classification. Training with mutation distribution and type is sufficient.
- Model misclassifications reflect shared biological characteristics in mutation topology.
- Discovered patterns in chromatin marks in misclassified cancers
- Research is fun
- I learned a lot of biology

## Next Steps

### Better genome segmentation

- Explore genome segmentations based on mutational process activity
- increase information content between mutational density and cell type

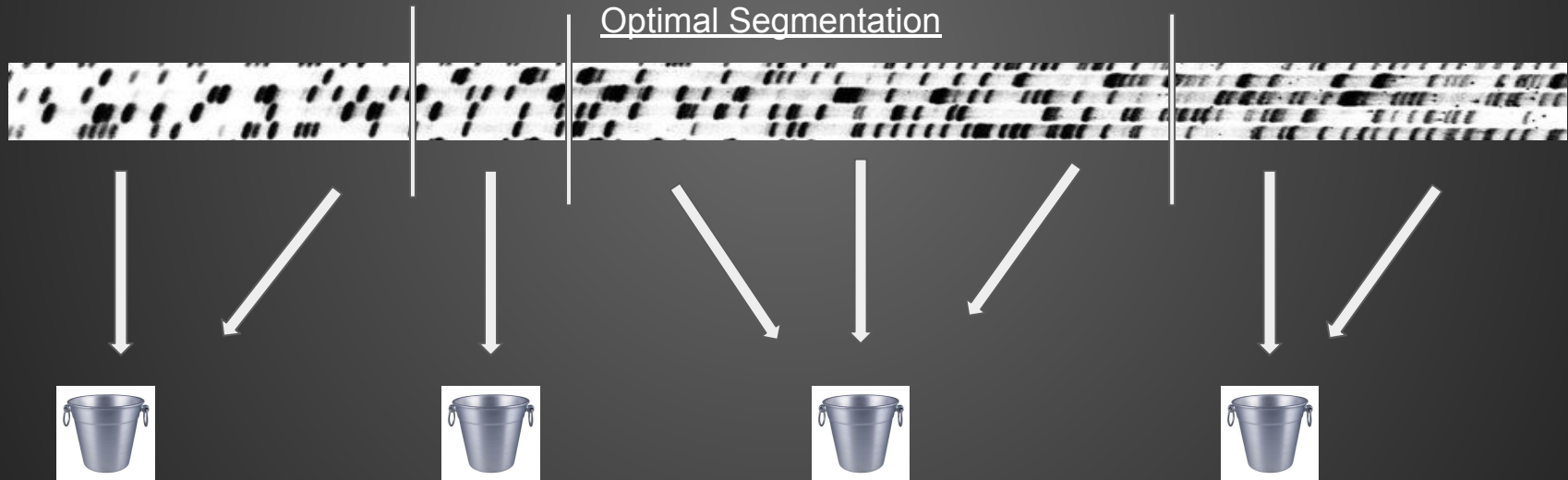
### Naive Segmentation



## Next Steps

### Better genome segmentation

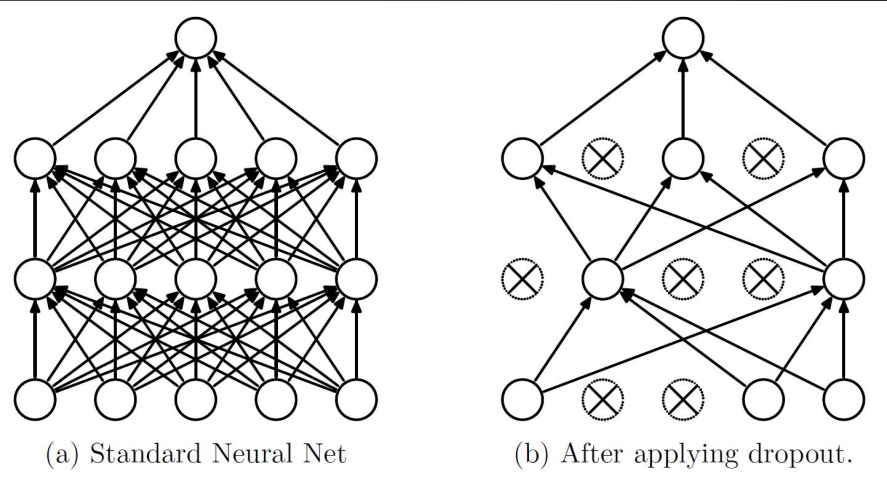
- Explore genome segmentations based on mutational process activity
- increase information content between mutational density and cell type



## Next Steps

### Extend classifier to rare cancer types

- Adopt a bayesian approach by implementing MC-dropout
- Certainty Estimates provide interpretable results

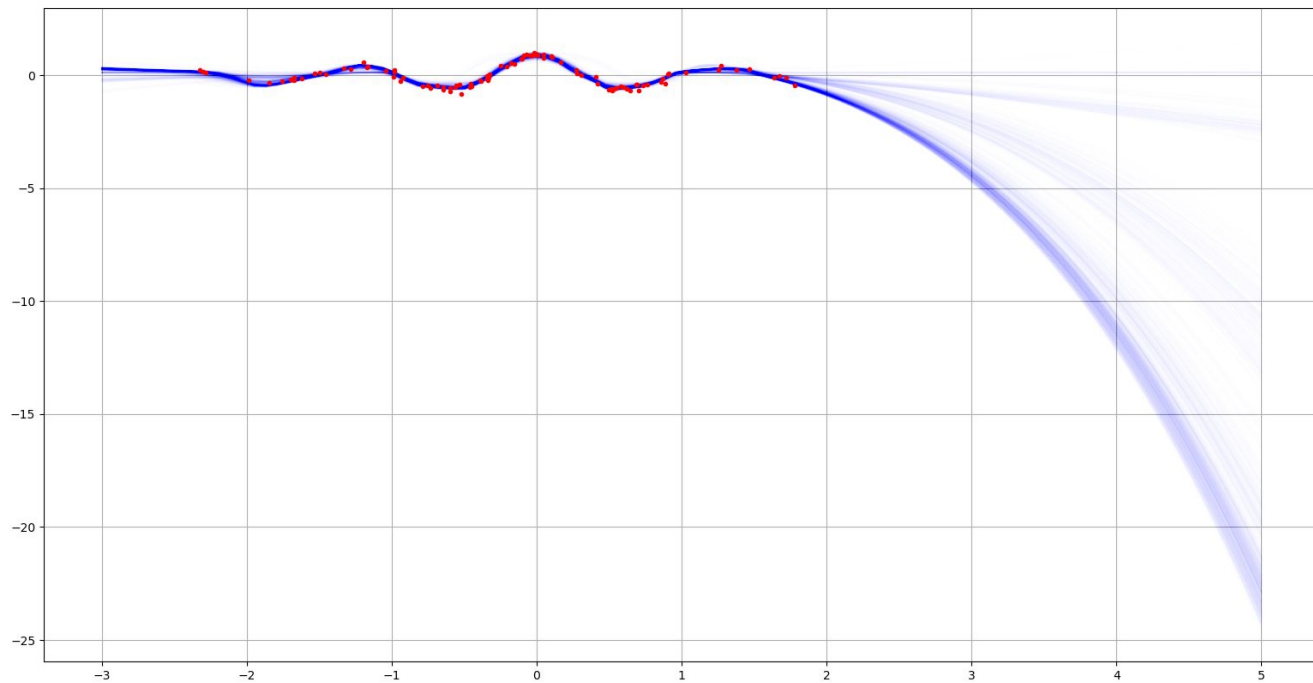


Source: Srivastava, Hinton, Krizhevsky, Sutskever and Salakhutdinov (2014)

MC-dropout (Monte Carlo Dropout): *dropping neuron activations during test time ~ Non-deterministic*

- Generates predictions that one can interpret as samples from a probability distribution

## Regression Example



# Acknowledgements

Thanks to everyone in *The Morris Lab*:

*Adamo Young*  
*Alexander Sasse*  
*Alinda Selega*  
*Amir*  
*Amit Deshwar*  
*Audrina Zhou*  
*Caitlyn*  
*Chris*  
*Chris Cole*

*Farzan Taj*  
*Gurnit Atwal*  
*Haoran Zhang*  
*Ilyes Baali*  
*Jarry Barber*  
*Jeff Wintersinger*  
*Jingping Qiao*  
*Kaitlin Iaverty*  
*Kate Nie*

*Kathryn Yu*  
*Kim Skead*  
*Linda*  
*Nik Krosigk*  
*Nil Sahin*  
*Quaid Morris*  
*Rozy Razavi*  
*Yoonsik Park*  
*Yulia Rubanova*

Special thanks to Gurnit Atwal and Quaid Morris