# Machine Learning for Data Generated from Next-Generation Sequencing

**Fangfang Xia**

fangfang@anl.gov

Computer Scientist

Argonne National Laboratory

2019-10-23

*DOE-NCI partnership to advance exascale development through cancer research*

**Jeremy Howard** @jeremyphoward · 22h

Seems like there is a real revolution going on in protein analysis thanks to deep learning language models. Here's an example that just got published



**Unified rational protein engineering with sequence-...**

UniRep learns fundamental protein features from millions of amino-acid sequences using a recurrent neural network. This summary of features can then be used ...

nature.com

💬 8          🔁 146          ♡ 629
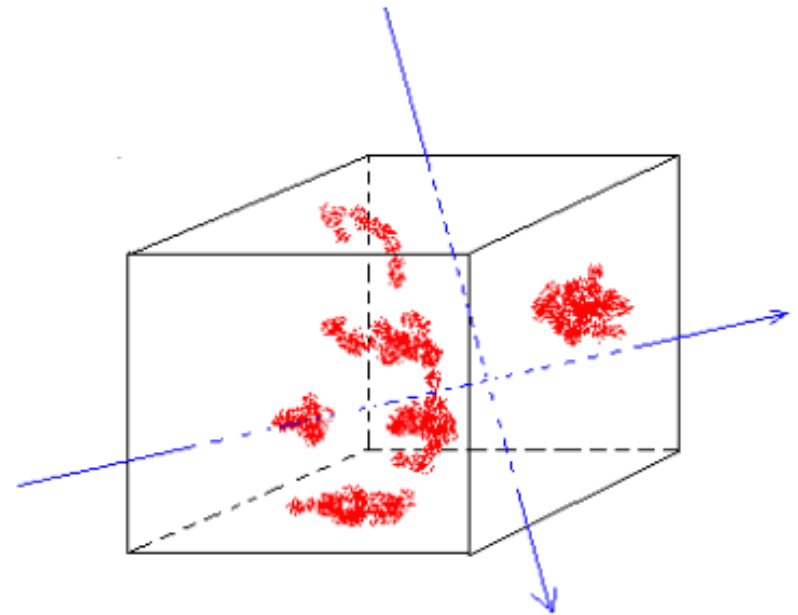
# Representation learning

dnaK
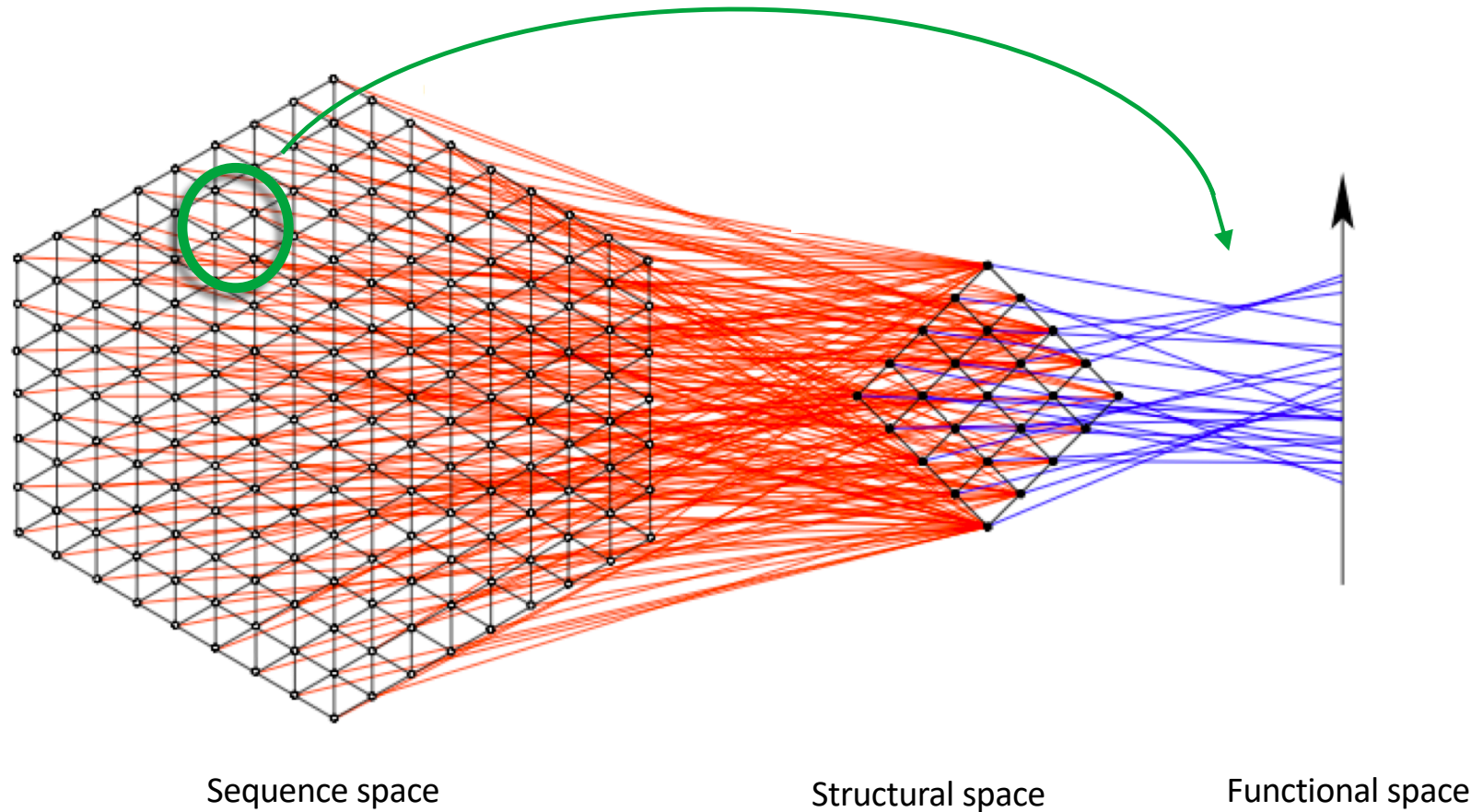
- atgggtaaaataattggtatcgacct

random

- tgaggtcgtagagtagacca

- $4^n$ possible sequences
  - Real sequences occupy a tiny fraction of all possible sequences
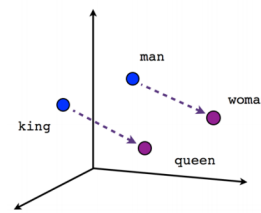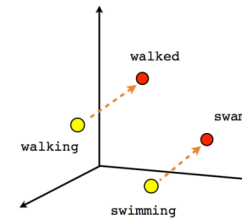
# Mapping the spaces: starting from raw seqs



Sequence space        Structural space        Functional space

# Embedding for categorical variables

| | |
|---|---|
| puppy | [1, 0, 0, 0] |
| dog | [0, 1, 0, 0] |
| kitten | [0, 0, 1, 0] |
| cat | [0, 0, 0, 1] |

| | |
|---|---|
| puppy | [0.9, 1.0, 0.0, 0.2] |
| dog | [1.0, 0.2, 0.0, 0.9] |
| kitten | [0.0, 1.0, 0.5, 0.1] |
| cat | [0.0, 0.2, 1.0, 1.0] |

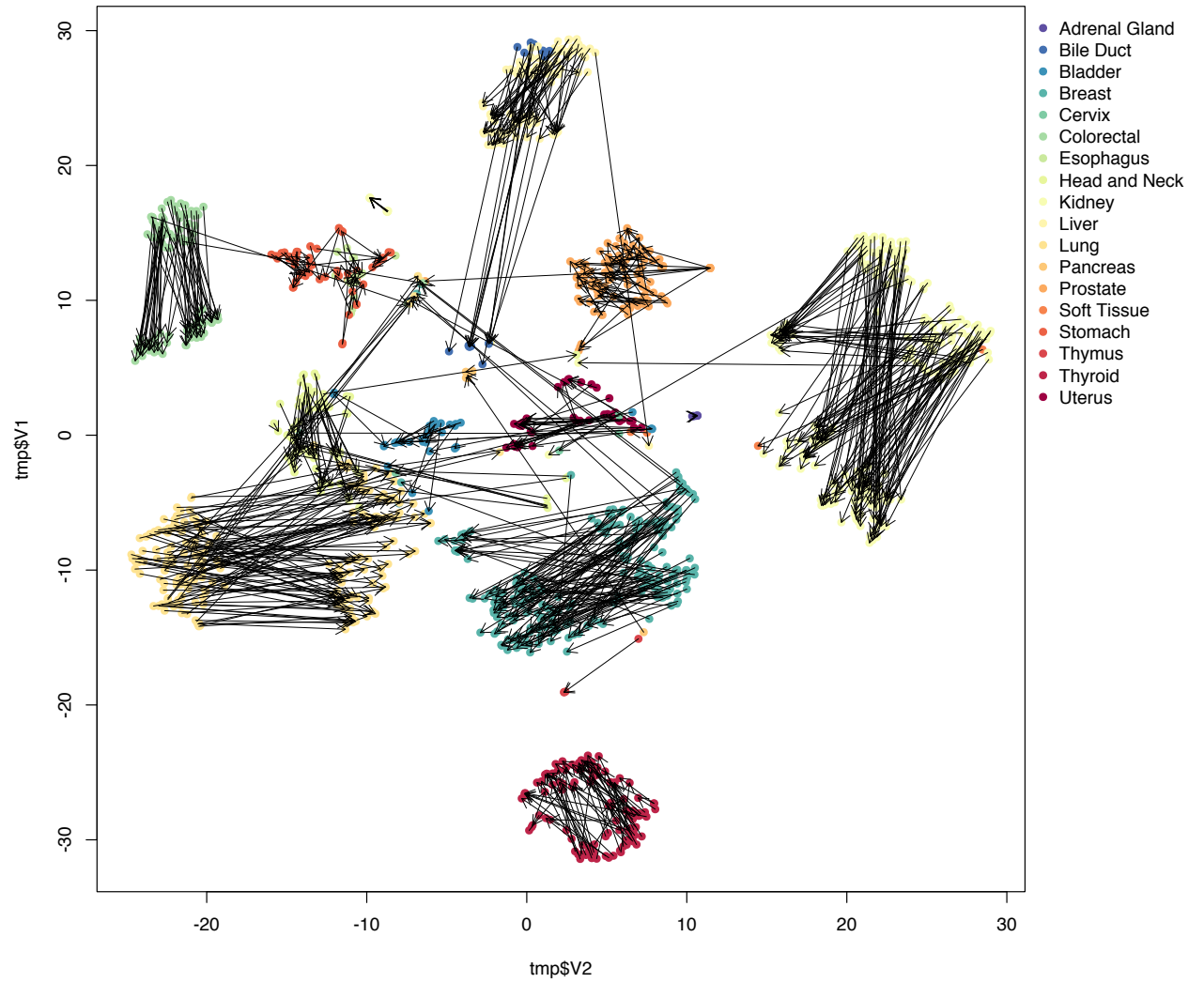*From the Instacart blog post 'Deep Learning with Emojis (not Math)'*

Normal => Tumor

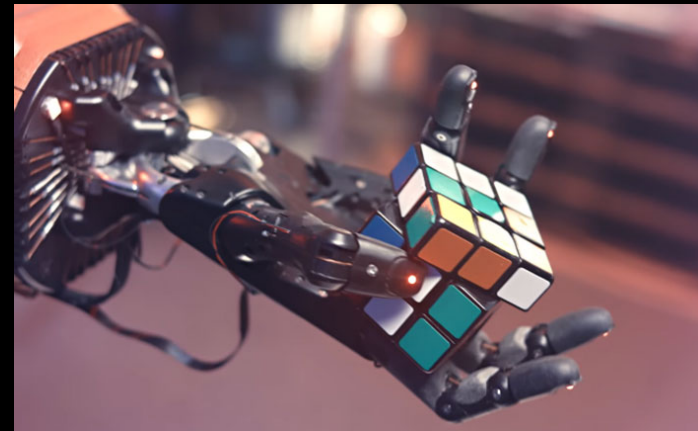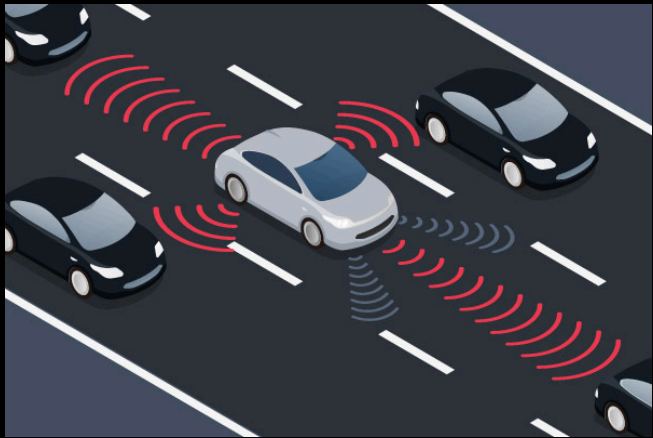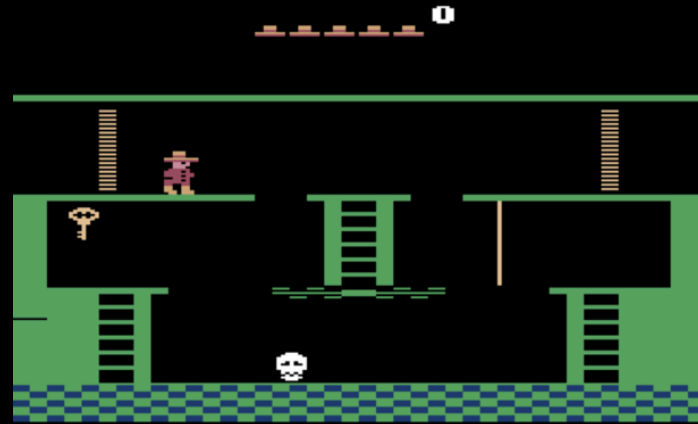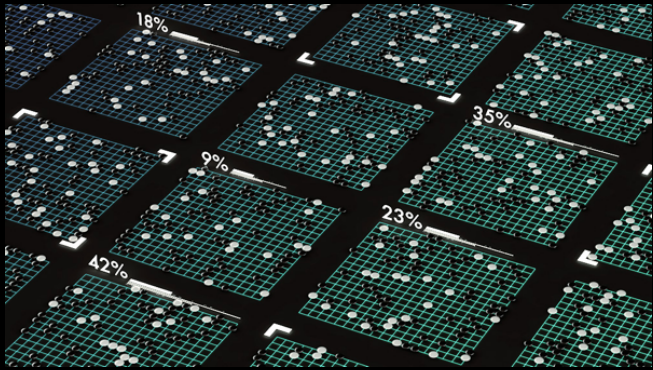# BIOLOGICAL STRUCTURE AND FUNCTION EMERGE FROM SCALING UNSUPERVISED LEARNING TO 250 MILLION PROTEIN SEQUENCES

**Alexander Rives** [*][†][‡]    **Siddharth Goyal** [*][§]    **Joshua Meier** [*][§]    **Demi Guo** [*][§]
**Myle Ott** [§]    **C. Lawrence Zitnick** [§]    **Jerry Ma** [†][§]    **Rob Fergus** [†][‡][§]

In the field of artificial intelligence, a combination of scale in data and model capacity enabled by unsupervised learning has led to major advances in representation learning and statistical generation. In biology, the anticipated growth of sequencing promises unprecedented data on natural sequence diversity. Learning the natural distribution of evolutionary protein sequence variation is a logical step toward predictive and generative modeling for biology. To this end we use unsupervised learning to train a deep contextual language model on 86 billion amino acids across 250 million sequences spanning evolutionary diversity. The resulting model maps raw sequences to representations of biological properties without labels or prior domain knowledge. The learned representation space organizes sequences at multiple levels of biological granularity from the biochemical to proteomic levels. Learning recovers information about protein structure: secondary structure and residue-residue contacts can be extracted by linear projections from learned representations. With small amounts of labeled data, the ability to identify tertiary contacts is further improved. Learning on full sequence diversity rather than individual protein families increases recoverable information about secondary structure. We show the networks generalize by adapting them to variant activity prediction from sequences only, with results that are comparable to a state-of-the-art variant predictor that uses evolutionary and structurally derived features.

[†]Correspondence to <arives@cs.nyu.edu>, <maj@fb.com>, and <robfergus@fb.com>
[‡]Dept. of Computer Science, New York University, USA
[§]Facebook AI Research, USA

# What is Machine Learning?

*The complexity in traditional computer programming is in the code (programs that people write). In machine learning, algorithms (programs) are in principle simple and the complexity (structure) is in the data. Is there a way that we can automatically learn that structure? That is what is at the heart of machine learning.*

-- Andrew Ng

That is, machine learning is the about the construction and study of systems that can learn from data. This is very different than traditional computer programming.

# The Cartoon Form

## Traditional Programming

Data →

Program →

**Computer**

→ Output

## Machine Learning

Data →

Output →

**Computer**

→ Program

# Examples in applying ML to NGS data

- Classifying cancer type with gene expression profiles

- Removing study bias in tumor gene expression profiles

- Classifying cancer type with SNP data

- Drug response prediction - introduction
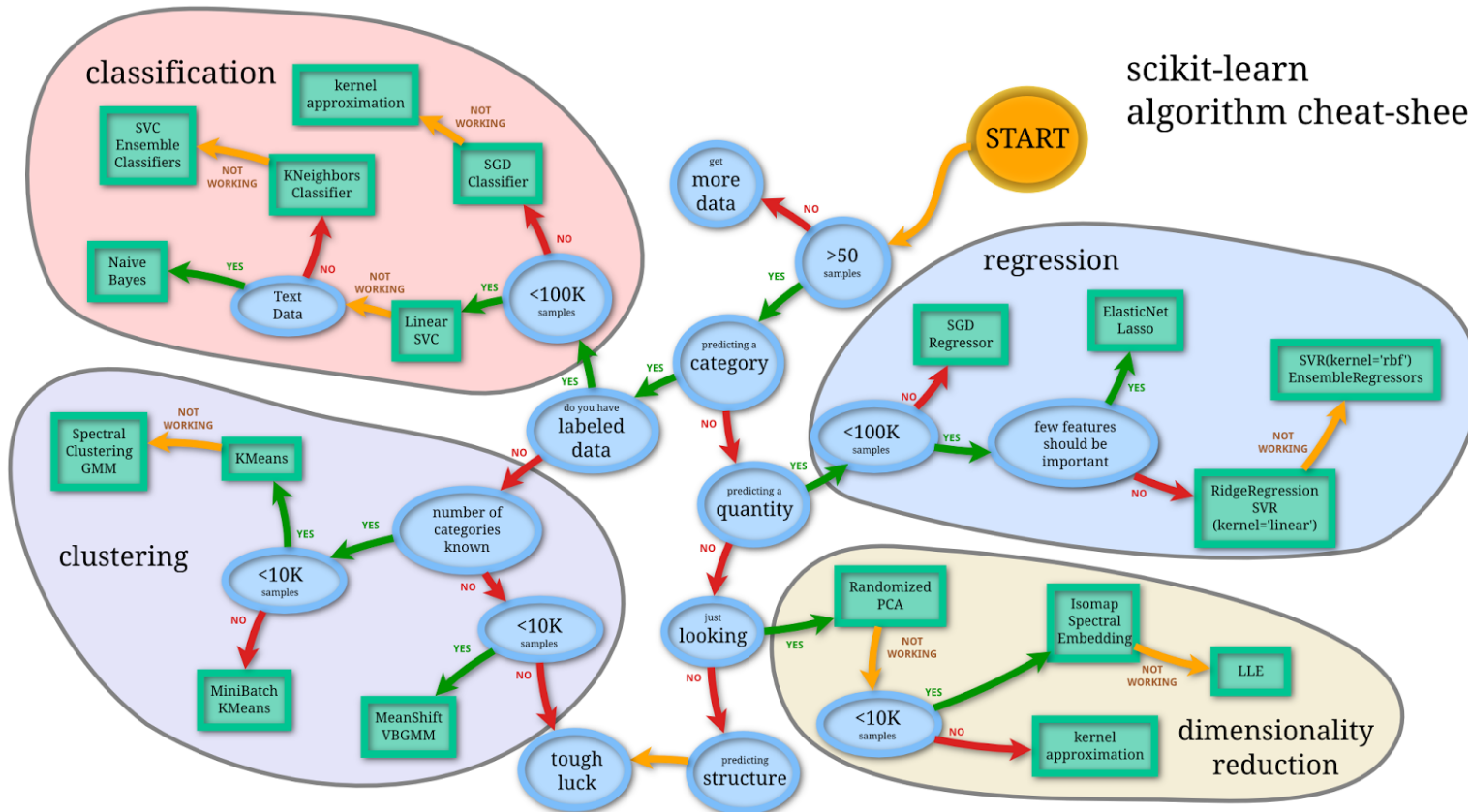
# Four Typical Problems



scikit-learn
algorithm cheat-sheet

**START**

**classification**

- kernel approximation
- SVC Ensemble Classifiers
- KNeighbors Classifier
- SGD Classifier
- Naive Bayes
- Text Data
- Linear SVC
- <100K samples
- >50 samples
- get more data
- predicting a category
- do you have labeled data

**regression**

- SGD Regressor
- ElasticNet Lasso
- SVR(kernel='rbf') EnsembleRegressors
- <100K samples
- few features should be important
- RidgeRegression SVR (kernel='linear')
- predicting a quantity

**clustering**

- Spectral Clustering GMM
- KMeans
- number of categories known
- <10K samples
- MiniBatch KMeans
- MeanShift VBGMM

**dimensionality reduction**

- Randomized PCA
- Isomap Spectral Embedding
- LLE
- <10K samples
- kernel approximation
- just looking
- predicting structure
- tough luck

# Deep learning in biology and medicine

*https://github.com/greenelab/deep-review*

# *Keras.js*

## *Tensorflow playground*

Epoch
000,564

Learning rate
0.03

Activation
Tanh

Regularization
None

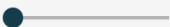Regularization rate
0

Problem type
Classification

## DATA

Which dataset do you want to use?

Ratio of training to test data: 50%
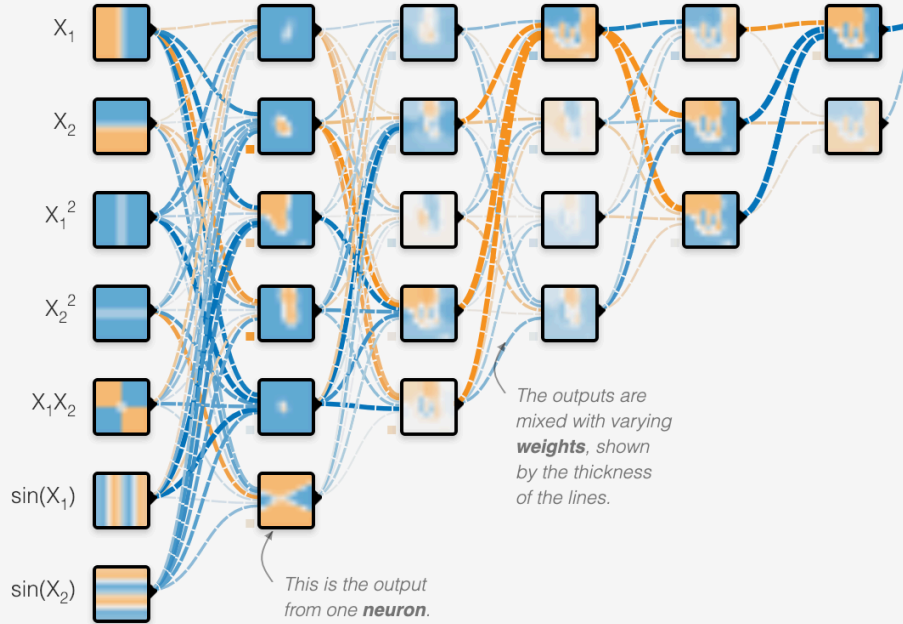
Noise: 0

Batch size: 10

REGENERATE

## FEATURES

Which properties do you want to feed in?

$X_1$

$X_2$

$X_1^2$

$X_2^2$

$X_1X_2$

$sin(X_1)$

$sin(X_2)$

This is the output from one **neuron**. Hover to see it larger.

## 5 HIDDEN LAYERS

6 neurons    5 neurons    4 neurons    3 neurons    2 neurons

The outputs are mixed with varying **weights**, shown by the thickness of the lines.

## OUTPUT

Test loss 0.223
Training loss 0.154

6
5
4
3
2
1
0
-1
-2
-3
-4
-5
-6

-6 -5 -4 -3 -2 -1 0 1 2 3 4 5 6

Colors shows data, neuron and weight values.

-1    0    1

☐ Show test data    ☐ Discretize output

# A few "good" runs

# *Deep Learning Basics*

# Neuron

**Dendrites**

**Cell Body**

**Connections from other Neurons**

**Axon**

**Neurotransmitter Molecules**

# Mathematical Model of a Neuron

inputs

weights

$x_1$

$w_{1j}$

$x_2$

$w_{2j}$

$x_3$

$w_{3j}$

$x_n$

$w_{nj}$

$\Sigma$

transfer
function

net input

$net_j$

activation
functon

$\varphi$

$o_j$

activation

$\theta_j$

threshold

z = w·x + b

o = $\varphi$(z)

# Activation

$$\sigma(z) \equiv \frac{1}{1 + e^{-z}}.$$



Sigmoid



ReLU

| Name | Plot | Equation | Derivative |
|---|---|---|---|
| Identity | | $f(x) = x$ | $f'(x) = 1$ |
| Binary step | | $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ | $f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$ |
| Logistic (a.k.a Soft step) | | $f(x) = \dfrac{1}{1 + e^{-x}}$ | $f'(x) = f(x)(1 - f(x))$ |
| TanH | | $f(x) = \tanh(x) = \dfrac{2}{1 + e^{-2x}} - 1$ | $f'(x) = 1 - f(x)^2$ |
| ArcTan | | $f(x) = \tan^{-1}(x)$ | $f'(x) = \dfrac{1}{x^2 + 1}$ |
| Rectified Linear Unit (ReLU) | | $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ | $f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ |
| Parameteric Rectified Linear Unit (PReLU) [2] | | $f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ | $f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ |
| Exponential Linear Unit (ELU) [3] | | $f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ | $f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ |
| SoftPlus | | $f(x) = \log_e(1 + e^x)$ | $f'(x) = \dfrac{1}{1 + e^{-x}}$ |

# Loss and Accuracy

Regression problem

y_true:  [   -50,    -10,    +80 ]
y_pred:  [   -30,    +10,    +60 ]

MSE = 0.04

Classification *view*

[ 1, 1, 0 ]
[ 1, 0, 1 ]

ACC = 0.67

y_true:  [   -50,    -10,    +80 ]
y_pred:  [   -50,    -10,    +45 ]

MSE = 0.04

[ 1, 1, 0 ]

ACC = 1.00

# Classification loss

- **Softmax**

- **Cross entropy**

How different are the two
probability distributions?

| scores | softmax $\longrightarrow$ | probabilities | true probabilities | cross entropy |
|---|---|---|---|---|
| 0.6 | | 0.059 | 0.1 | |
| -1 | $\dfrac{e^{z_j}}{\sum_{k=1}^{K} e^{z_k}}$ | 0.012 | 0.2 | 0.866 |
| 3.2 | | **0.797** | **0.6** | |
| 1.4 | | 0.132 | 0.1 | |

$$H(p,q) = -\sum_{x} p(x)\log q(x)$$

**c**

Output units

Hidden units H2

Hidden units H1

Input units

$y_l = f(z_l)$

$z_l = \sum_{k \,\varepsilon\, H2} w_{kl} \, y_k$

$y_k = f(z_k)$

$z_k = \sum_{j \,\varepsilon\, H1} w_{jk} \, y_j$

$y_j = f(z_j)$

$z_j = \sum_{i \,\varepsilon\, \text{Input}} w_{ij} \, x_i$

$w_{kl}$

$w_{jk}$

$w_{ij}$

**d**

Compare outputs with correct
answer to get error derivatives



$$\frac{\partial E}{\partial y_l} = y_l - t_l$$

$$\frac{\partial E}{\partial z_l} = \frac{\partial E}{\partial y_l} \frac{\partial y_l}{\partial z_l}$$

$$\frac{\partial E}{\partial y_k} = \sum_{l \, \varepsilon \, out} w_{kl} \frac{\partial E}{\partial z_l}$$

$$\frac{\partial E}{\partial z_k} = \frac{\partial E}{\partial y_k} \frac{\partial y_k}{\partial z_k}$$

$$\frac{\partial E}{\partial y_j} = \sum_{k \, \varepsilon \, H2} w_{jk} \frac{\partial E}{\partial z_k}$$

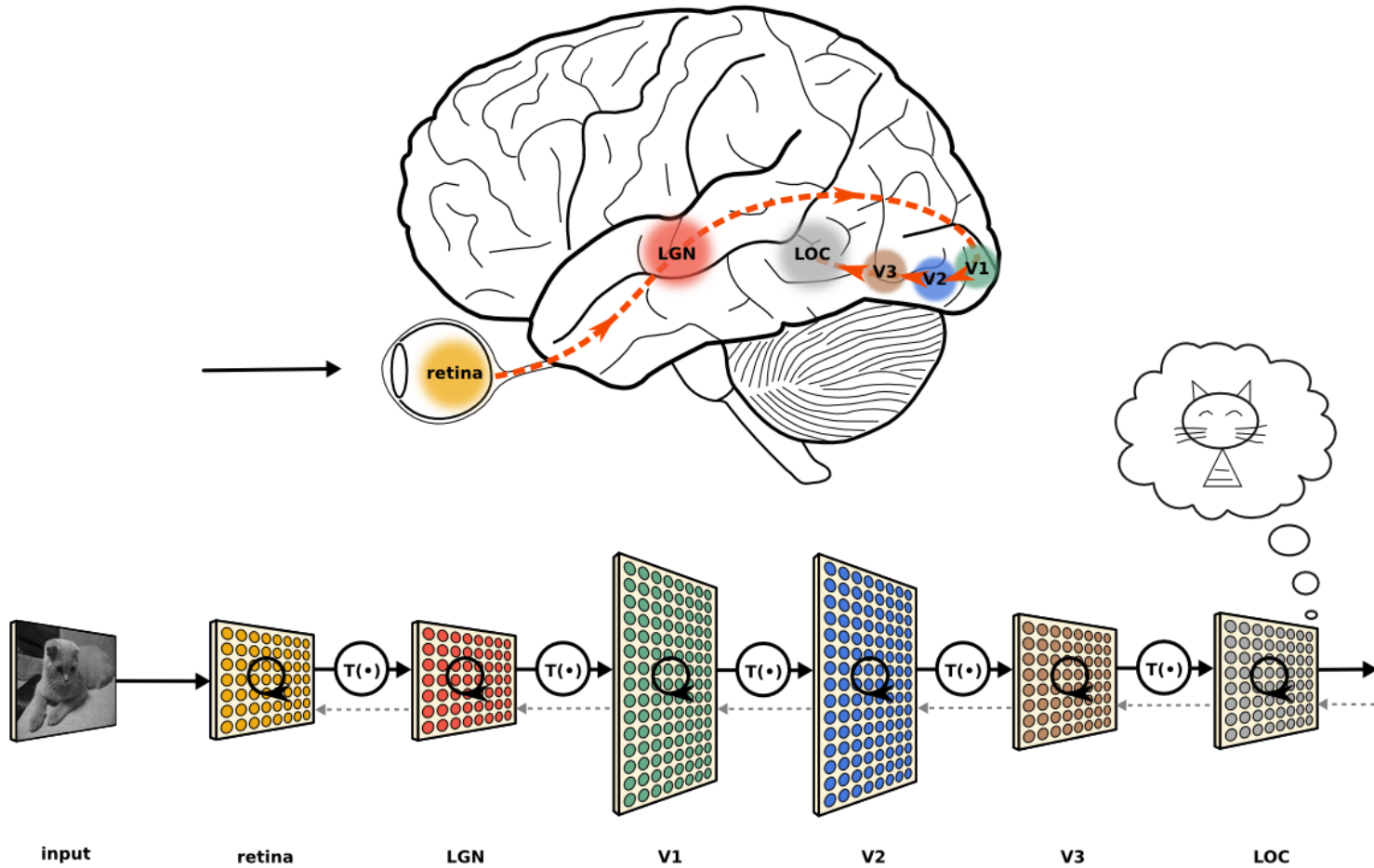$$\frac{\partial E}{\partial z_j} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial z_j}$$

$w_{kl}$

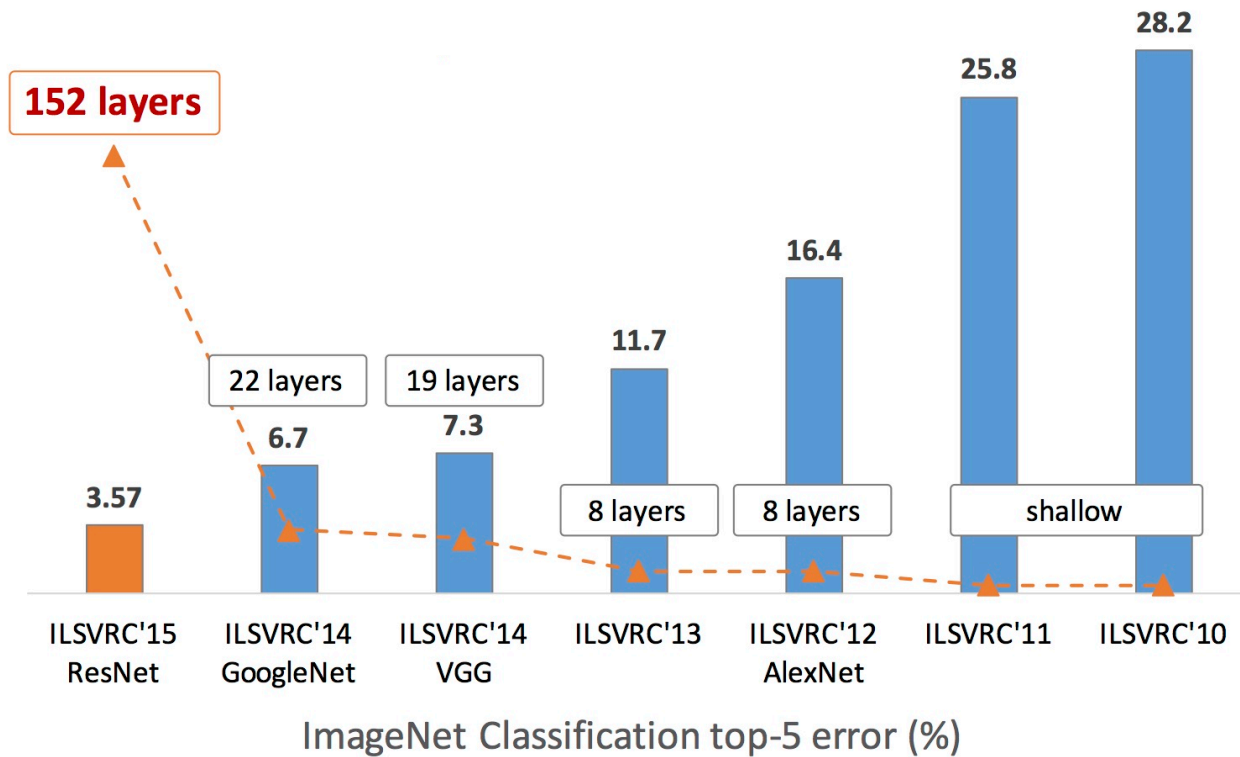$w_{jk}$
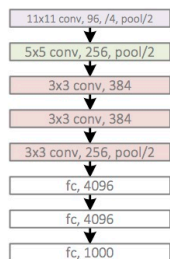
$w_{ij}$

# Human Vision System



input   retina   LGN   V1   V2   V3   LOC

# Increasing Depth Works...



ImageNet Classification top-5 error (%)

# Example CNNs Structures from ILSVRC Winners



**AlexNet, 8 layers (ILSVRC 2012)**

| |
|---|
| 11x11 conv, 96, /4, pool/2 |
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 3x3 conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

**VGG, 19 layers (ILSVRC 2014)**

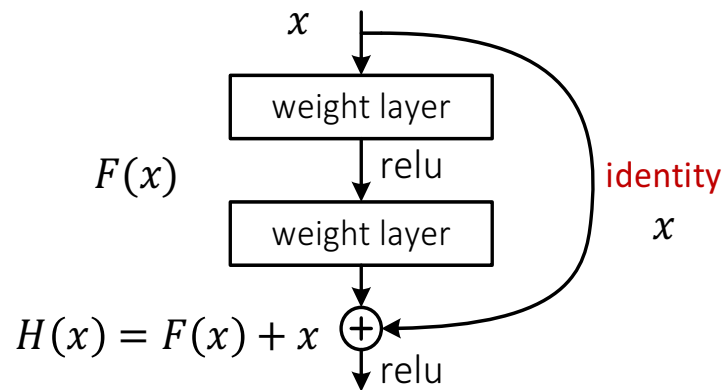| |
|---|
| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

**GoogleNet, 22 layers (ILSVRC 2014)**

Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 2016.

# Deep Residual Learning

- $F(x)$ is a residual mapping w.r.t. identity



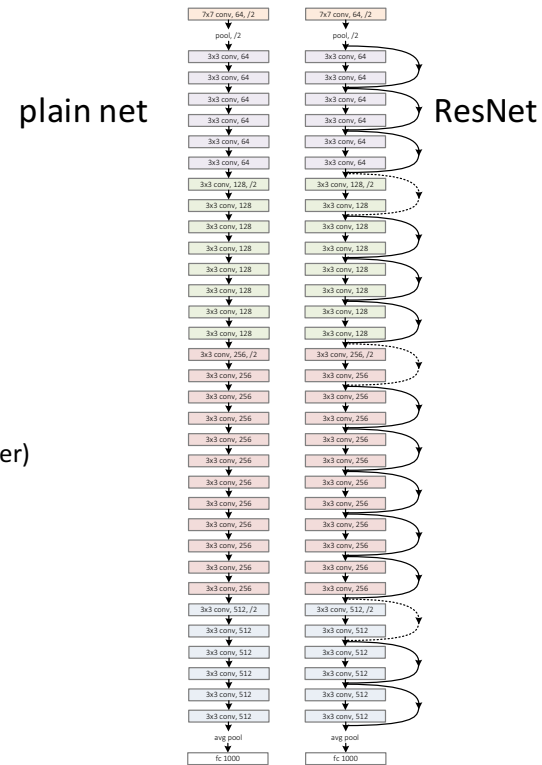- If identity were optimal, easy to set weights as 0

- If optimal mapping is closer to identity, easier to find small fluctuations

# Deep Residual Learning

## Network "Design"

plain net    ResNet

- Keep it simple

- Our basic design (VGG-style)
  - all 3x3 conv (almost)
  - spatial size /2 => # filters x2 (~same complexity per layer)
  - Simple design; just deep!

- Other remarks:
  - no hidden fc
  - no dropout

# Dropout

SRIVASTAVA, HINTON, KRIZHEVSKY, SUTSKEVER AND SALAKHUTDINOV

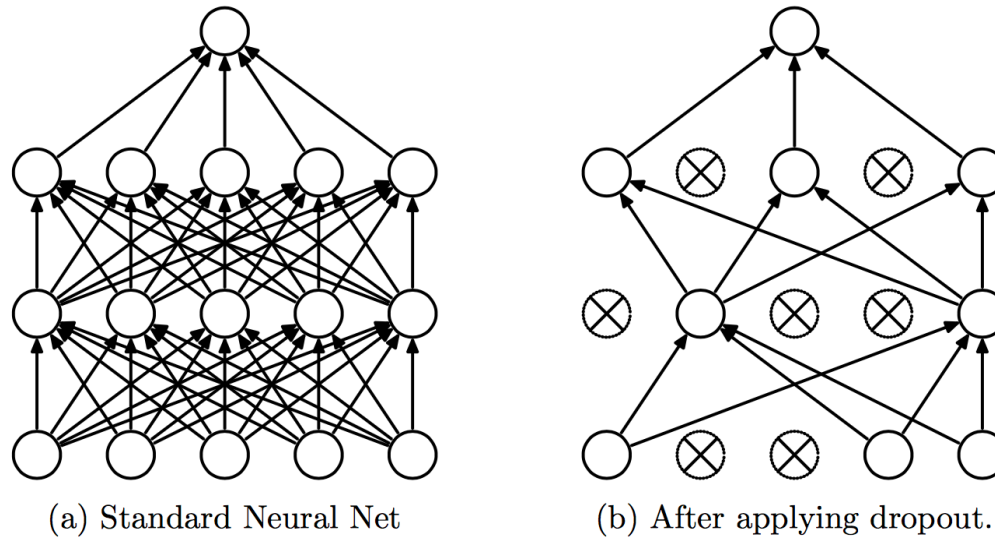(a) Standard Neural Net        (b) After applying dropout.

Figure 1: Dropout Neural Net Model. **Left**: A standard neural net with 2 hidden layers. **Right**: An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

**Andres Torrubia** @antor · Oct 18

Replying to @radekosmulski

Sometimes it feels like 👇

**QUESTION 8/10**

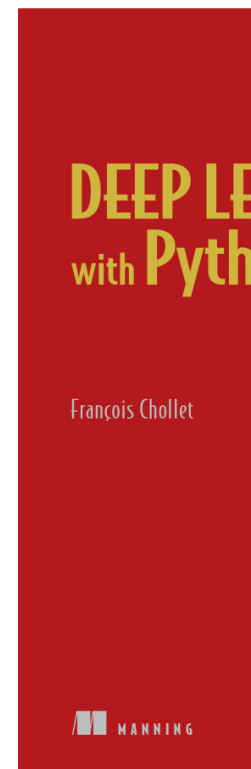As demonstrated in research, what is the trick that enables training a 10_000 layer deep neural network?

○ dropout
○ batchnorm
○ initialization
● thoughts and prayers

**Submit answer**

💬 2     🔁 17     ♡ 72

# Essential Deep Learning Resources

- Books
  - Michael Nielsen's free online book on deep learning
  - Deep Learning textbook (Goodfellow et al, 2016)
  - Deep Learning with Python (Keras book)
- Online courses
  - FastAI
  - Stanford CS231n
  - Andrew Ng
- Python
  - Python Cookbook, 3rd Edition
- Other resources
  - PyTorch tutorials
  - Kaggle kernels / discussions; No free hunch interviews
  - Twitter
  - arxiv-sanity

# Improving Models with Domain Knowledge

- Books
  - Michael Nielsen's free online book on deep learning
  - Deep Learning textbook (Goodfellow et al, 2016)
  - Deep Learning with Python (Keras book)
- Online courses
  - FastAI
  - Stanford CS231n
  - Andrew Ng
- Python
  - Python Cookbook, 3rd Edition
- Other resources
  - PyTorch tutorials
  - Kaggle kernels / discussions; No free hunch interviews
  - Twitter
  - arxiv-sanity

About 10,000 deep learning papers have been written about "hard-coding priors about a specific task into a NN architecture works better than a lack of prior" -- but they're typically being passed as "architecture XYZ offers superior performance for [overly generic task category]"

| 749 Likes | 168 Retweets |
|---|---|
| Jun 1, 2019 at 5:19 PM | via **Twitter Web Client** |

**François Chollet** @fchollet  36d
You can always "buy" performance by either training on more data, better data, or by injecting task information into the architecture or the preprocessing. However, this isn't informative about the generalization power of the techniques used (which is the only thing that matters)

# *Cancer Type Classification with RNAseq*

🏠 Home    📑 Projects    Data    📊 Analysis    🔍 Quick Search    Login    🛒 Cart 0    GDC Apps

## Harmonized Cancer Datasets
# Genomic Data Commons Data Portal

*Get Started by Exploring:*

📑 **Projects**    ⊙ **Data**

*Perform Advanced Search Queries, such as:*

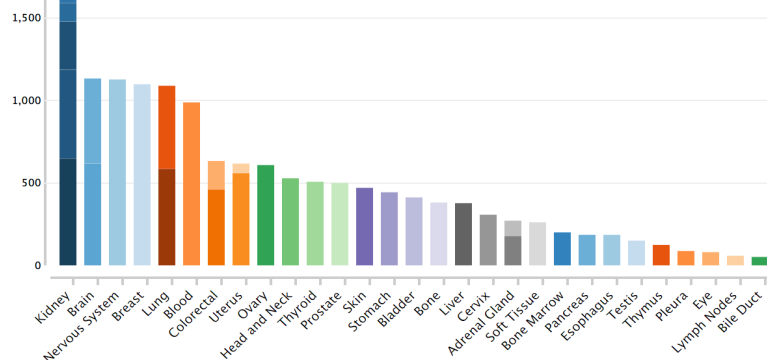| | | |
|---|---|---|
| Cases of kidney cancer diagnosed at the age of 20 and below | 736 Cases | 1,519 Files |
| CNV data of female brain cancer cases | 459 Cases | 1,788 Files |
| Gene expression quantification data in TCGA-GBM project | 166 Cases | 522 Files |

### Cases by Primary Site



| DATA PORTAL SUMMARY | PROJECTS | PRIMARY SITE | CASES | FILES |
|---|---|---|---|---|
| *Data Release 6.0 - May 9, 2017* | 📑 39 | 29 | 👤 14,551 | 📄 274,724 |

## Infrastructure

*Data is continuously being processed and harmonized by the GDC.*
*View GDC system statistics:*

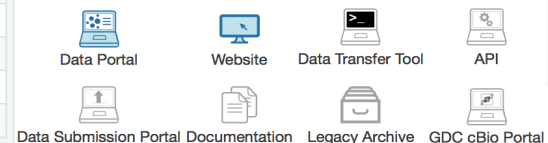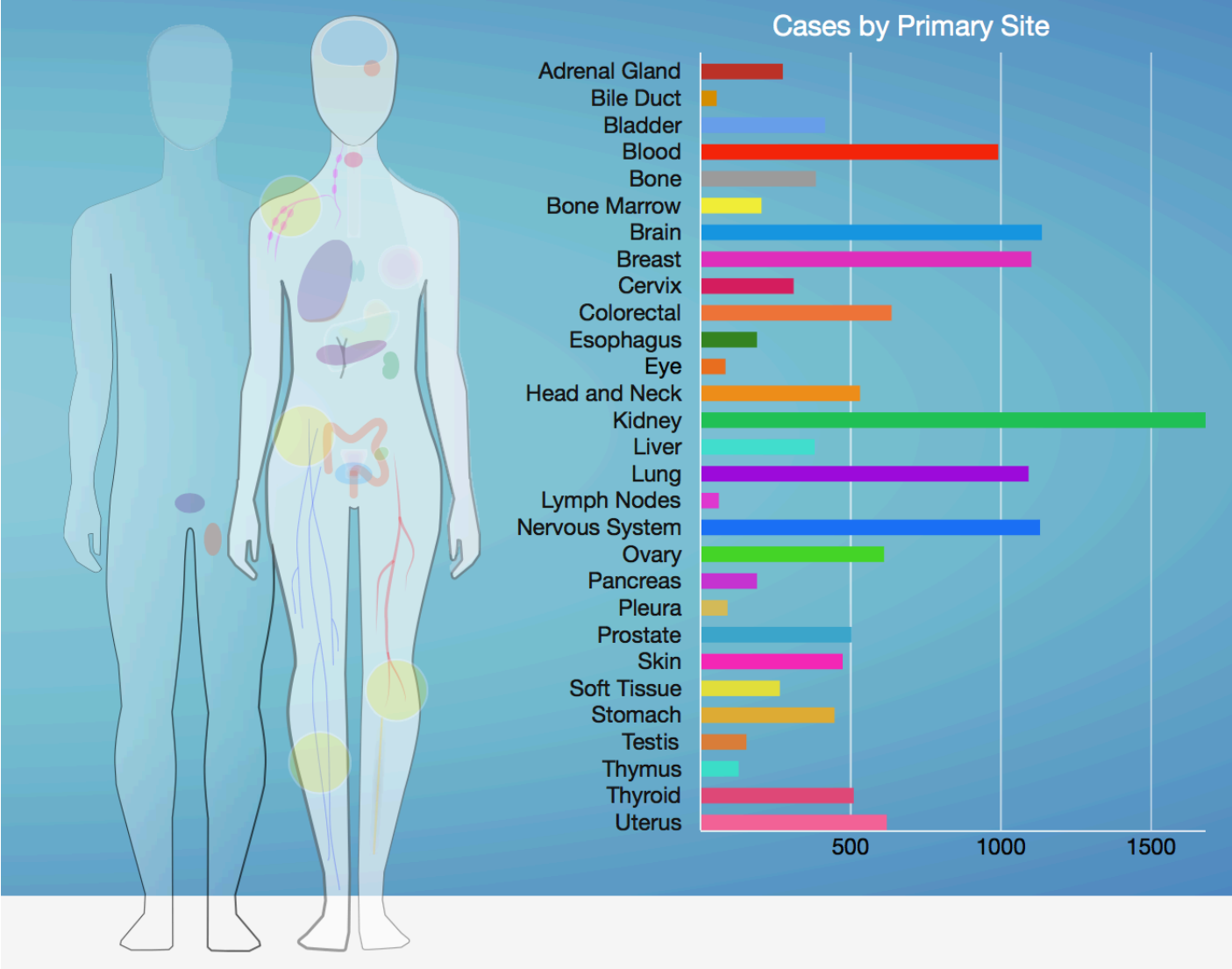| Compute Infrastructure | 12,800 Cores | 87.96 TB RAM |
|---|---|---|
| Storage Infrastructure | 4.98 PB Used | 5.42 PB Total |

## Documentation

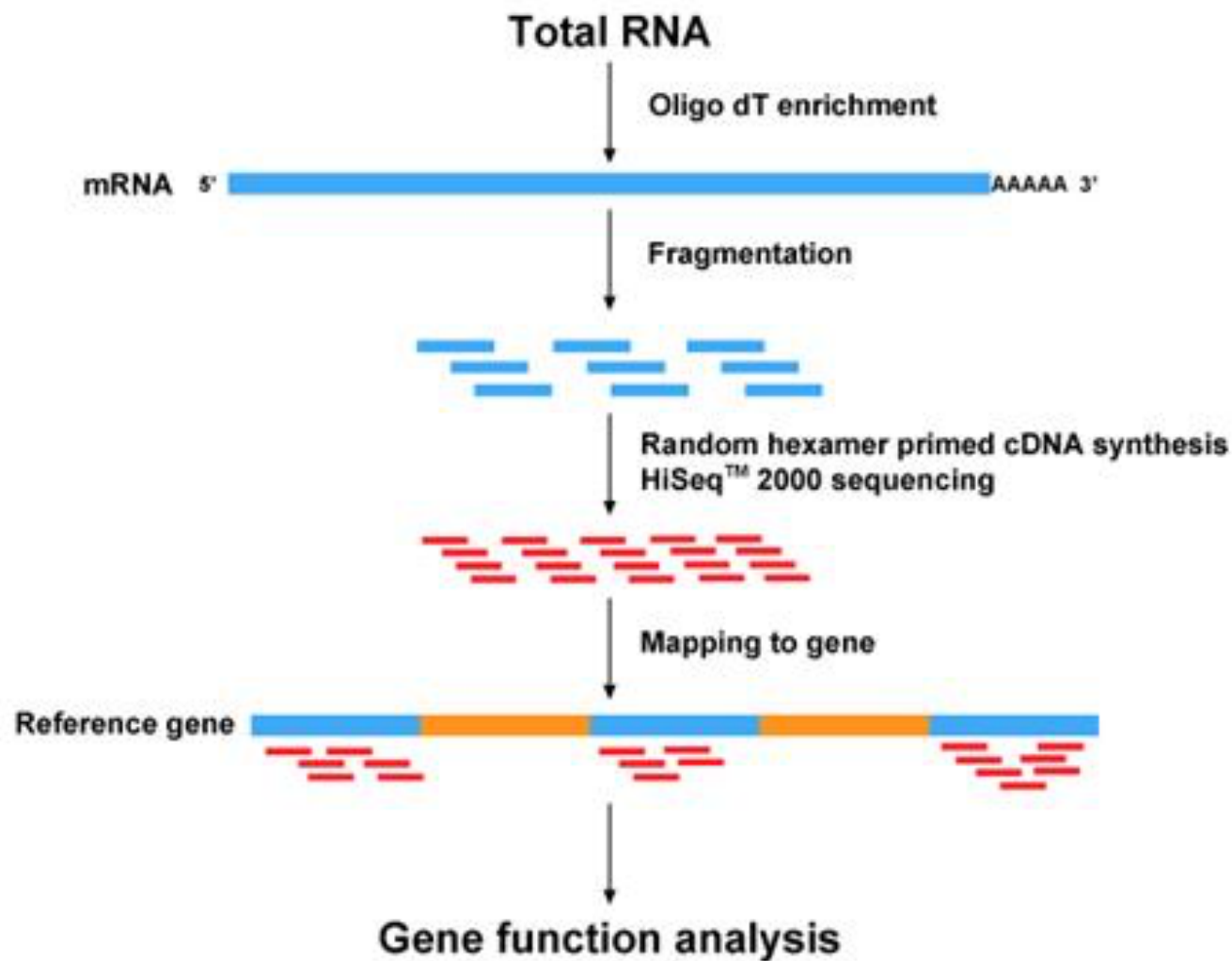Learn how to use the GDC Data Portal to its full potential with common topics such as:

Browse Data using Facet Search

Search Data with Advanced Search Technology

Project Based Data Availability

Controlled Access Data

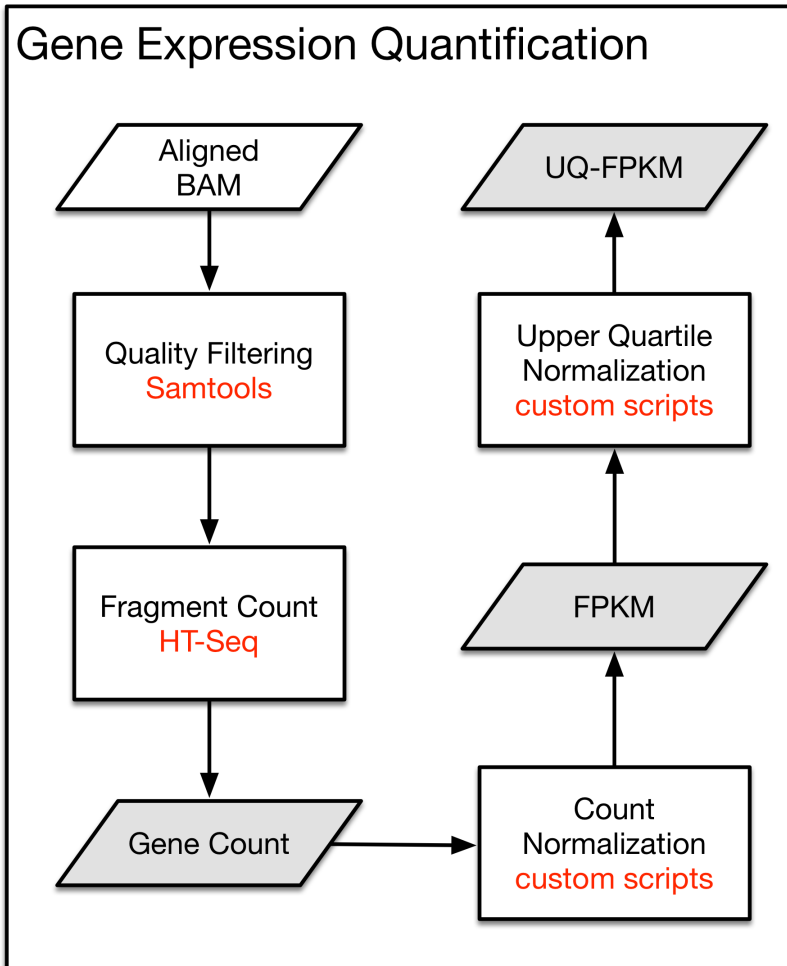**Visit the Documentation Website »**

## GDC Applications

The GDC Data Portal is a robust data-driven platform that allows cancer researchers and bioinformaticians to search and download cancer data for analysis. The GDC applications include:

Data Portal    Website    Data Transfer Tool    API

Data Submission Portal    Documentation    Legacy Archive    GDC cBio Portal

Cases by Primary Site

**Total RNA**

Oligo dT enrichment

mRNA  5'  AAAAA  3'

Fragmentation

Random hexamer primed cDNA synthesis
HiSeq™ 2000 sequencing

Mapping to gene

Reference gene

**Gene function analysis**

# Gene Expression Quantification



RPKM (reads per kilobase per million mapped reads)
Upper Quantile (UQ)

FPKM

The Fragments per Kilobase of transcript per Million mapped reads (FPKM) calculation normalizes read count by dividing it by the gene length and the total number of reads mapped to protein-coding genes.

Upper Quartile FPKM

The upper quartile FPKM (FPKM-UQ) is a modified FPKM calculation in which the total protein-coding read count is replaced by the 75th percentile read count value for the sample.

Calculations

$$FPKM = \frac{RC_g * 10^9}{RC_{pc} * L} \qquad FPKM - UQ = \frac{RC_g * 10^9}{RC_{g75} * L}$$

- **RC$_g$:** Number of reads mapped to the gene
- **RC$_{pc}$:** Number of reads mapped to all protein-coding genes
- **RC$_{g75}$:** The 75th percentile read count value for genes in the sample
- **L:** Length of the gene in base pairs

**Note:** The read count is multiplied by a scalar ($10^9$) during normalization to account for the kilobase and 'million mapped reads' units.

# Each Sample has > 60,000 columns

# Encode numeric columns

- Log transformation: log (1+x)

- Standard scalar: z score

- Rank => Gaussian

- Discretization

- Mark missing values

# One Hot Encoding of Categories

| State | Binary | One-Hot | Hamming 2 | Hamming 3 |
|-------|--------|---------|-----------|-----------|
| S0 | 000 | 00000001 | 0000 | 000000 |
| S1 | 001 | 00000010 | 0011 | 000111 |
| S2 | 010 | 00000100 | 0101 | 011001 |
| S3 | 011 | 00001000 | 0110 | 011110 |
| S4 | 100 | 00010000 | 1001 | 101010 |
| S5 | 101 | 00100000 | 1010 | 101101 |
| S6 | 110 | 01000000 | 1100 | 110011 |
| S7 | 111 | 10000000 | 1111 | 110100 |

# Open Source Framework Comparison

| | Languages | Tutorials and training materials | CNN modeling capability | RNN modeling capability | Architecture: easy-to-use and modular front end | Speed | Multiple GPU support | Keras compatible |
|---|---|---|---|---|---|---|---|---|
| Theano | Python, C++ | ++ | ++ | ++ | + | ++ | + | + |
| Tensor-Flow | Python | +++ | +++ | ++ | +++ | ++ | ++ | + |
| Torch | Lua, Python (new) | + | +++ | ++ | ++ | +++ | ++ | |
| Caffe | C++ | + | ++ | | + | + | + | |
| MXNet | R, Python, Julia, Scala | ++ | ++ | + | ++ | ++ | +++ | + |
| Neon | Python | + | ++ | + | + | ++ | + | |
| CNTK | C++ | + | + | +++ | + | ++ | + | + |

# Keras

- https://keras.io/
- Minimalist, highly modular neural networks library
- Written in Python
- Capable of running on top of either TensorFlow/Theano and CNTK
- Developed with a focus on enabling fast experimentation

```python
from keras.layers import Input, Dense
from keras.models import Model


input_layer = Input(shape=(1000,))
fc_1 = Dense(512, activation='relu')(input_layer)
fc_2 = Dense(256, activation='relu')(fc_1)
output_layer = Dense(10, activation='softmax')(fc_2)


model = Model(input=input_layer, output=output_layer)
model.compile(optimizer='rmsprop',
              loss='categorical_crossentropy',
              metrics=['accuracy'])


model.fit(bow, newsgroups.target)
predictions = model.predict(features).argmax(axis=1)
```

# DNN hyperparameter examples

- Data preprocess
    - Positive features: logarithmic transformation y = log(1+x)
    - Mixed features: standard scaler
- Number of hidden layers: 4
- Number of neurons in hidden layers: 4000-2000-1000-1000
- Activation function: ReLU
- Dropouts: Input: 0%; layers 1,2,3: 25%; layer 4: 10%
- Initialization: no unsupervised pretraining
- Optimization: learning rate = 0.05, momentum = 0.9, and weight decay = 0.0001
- Training epochs: as large as possible (dropout can prevent overfitting)

# *Code Examples*

# Cancer Type Classification

```
4320/4320 [==============================] - 87s - loss: 3.2885 - acc: 0.0537 - val_loss: 2.9542 - val_acc: 0.0556
Epoch 2/400
4320/4320 [==============================] - 76s - loss: 2.9777 - acc: 0.0752 - val_loss: 2.8273 - val_acc: 0.1083
Epoch 3/400
4320/4320 [==============================] - 78s - loss: 2.8117 - acc: 0.1176 - val_loss: 2.5971 - val_acc: 0.2194
Epoch 4/400
4320/4320 [==============================] - 77s - loss: 2.5094 - acc: 0.2060 - val_loss: 2.1191 - val_acc: 0.3306
Epoch 5/400
4320/4320 [==============================] - 78s - loss: 2.0385 - acc: 0.3442 - val_loss: 1.6411 - val_acc: 0.4648
Epoch 6/400
4320/4320 [==============================] - 75s - loss: 1.4995 - acc: 0.5079 - val_loss: 0.9846 - val_acc: 0.7704
Epoch 7/400
4320/4320 [==============================] - 77s - loss: 1.0688 - acc: 0.6481 - val_loss: 0.5628 - val_acc: 0.8796
Epoch 8/400
4320/4320 [==============================] - 76s - loss: 0.7657 - acc: 0.7461 - val_loss: 0.4952 - val_acc: 0.8509
Epoch 9/400
4320/4320 [==============================] - 76s - loss: 0.5729 - acc: 0.8123 - val_loss: 0.2803 - val_acc: 0.9287
Epoch 10/400
4320/4320 [==============================] - 79s - loss: 0.4389 - acc: 0.8620 - val_loss: 0.1962 - val_acc: 0.9398
Epoch 11/400
```

Model Loss

Cancer Type Classification
18 types each with ~300 RNAseq profiles

https://github.com/ECP-CANDLE/Benchmarks/tree/frameworks/Pilot1/TC1

Model Accuracy

Cancer Type Classification
18 types each with ~300 RNAseq profiles

https://github.com/ECP-CANDLE/Benchmarks/tree/frameworks/Pilot1/TC1

# VAE Latent Representation of GDC Expression



Color legend:
- 9 — Uterine Corpus Endometrial Carcinoma
- 8 — Thyroid Carcinoma
- 7 — Skin Cutaneous Melanoma
- 6 — Prostate Adenocarcinoma
- 5 — **Other Cancer Types**
- 4 — Lung Squamous Cell Carcinoma
- 3 — Lung Adenocarcinoma
- 2 — Head and Neck Squamous Cell Carcinoma
- 1 — Breast Invasive Carcinoma
- 0 — Brain Lower Grade Glioma

https://github.com/ECP-CANDLE/Benchmarks/tree/frameworks/Pilot1/P1B1

# How did we know it might work?

- Build autoencoders first with the features you are going to work with
- If you get reasonable accuracy then the model can learn a representation and that is a good sign
- Class balance seems to matter
- Number of training examples matters > 1000 is good > 10,000 better, > 100,000 much better
- Hyper parameter search is also important once you get something that basically works

# Generate Compact Molecular Signatures

- **For each agent or class of agents we will apply feature selection methods to the models to generate where possible a compact molecular signature that retains prediction performance**
  - Typical reduced signatures include O(10)-O(100) features from >> 50,000 starting features
  - Features may be genes, SNPs, µRNA etc.

- **Developed and applied multiple feature selection methods**
  - Selection criterion: Chi2, Anova, mutual info, ensemble ML, deep neural networks
  - Algorithms: ranking, intersection, iterat~~~maximization
  - Supervised recursive binning

- **Extracted compact features**
  - Features from cancer type prediction
  - 50 features: 0.981 accuracy
  - 20 features: 0.976 accuracy
  - 14 features: 0.973 accuracy
  - RNAseq is more informative than miRN~



LOOCV prediciton accuracy (JMIM method)

- i: miRNA
- ii: RNAseq
- iii: all
- iv: (i) and (ii) interspersed

# Analyze Molecular Signatures to Provide Insight to Potential Mechanisms

- **Started mapping gene features to pathways**
  - Enrichment analysis will be applied to the signatures to identify associated pathways
  - Pathways will be identified that associate with both sensitive and resistant response phenotypes

- **Identified co-located or known interacting pairs of gene and microRNA signatures**

  - Top miRNA feature hsa.mir.10a is co-located with ENSG00000120075.5
  - It has also been experimentally verified that this miRNA downregulates the corresponding HOX genes

| Rank | MIR | RNA | MIR&RNA |
|---|---|---|---|
| 1 | hsa.mir.10a | ENSG00000119888.9 | ENSG00000119888.9 |
| 2 | hsa.mir.205 | ENSG00000170370.11 | ENSG00000170370.11 |
| 3 | hsa.mir.181a.2 | ENSG00000157551.16 | ENSG00000157551.16 |
| 4 | hsa.mir.135a.1 | ENSG00000124664.9 | ENSG00000124664.9 |
| 5 | hsa.mir.203a | ENSG00000102554.12 | ENSG00000102554.12 |
| 6 | hsa.mir.196b | ENSG00000009765.13 | hsa.mir.10a |
| 7 | hsa.mir.194.1 | ENSG00000275410.3 | ENSG00000009765.13 |
| 8 | hsa.mir.9.3 | ENSG00000274173.1 | ENSG00000275410.3 |
| 9 | hsa.mir.196a.2 | ENSG00000120075.5 | ENSG00000274173.1 |
| 10 | hsa.mir.429 | ENSG00000124466.8 | hsa.mir.205 |
| 11 | hsa.mir.375 | ENSG00000204385.9 | ENSG00000204385.9 |
| 12 | hsa.mir.584 | ENSG00000103449.10 | ENSG00000104447.10 |
| 13 | hsa.mir.135b | ENSG00000104447.10 | ENSG00000103449.10 |
| 14 | hsa.mir.10b | ENSG00000255794.5 | ENSG00000078399.14 |
| 15 | hsa.let.7i | ENSG00000189334.7 | ENSG00000189334.7 |
| 16 | hsa.mir.125b.2 | ENSG00000165215.6 | ENSG00000165215.6 |
| 17 | hsa.mir.30a | ENSG00000137203.9 | ENSG00000137203.9 |
| 18 | hsa.mir.200c | ENSG00000078399.14 | ENSG00000255794.5 |
| 19 | hsa.mir.203b | ENSG00000103942.11 | ENSG00000103942.11 |
| 20 | hsa.mir.944 | ENSG00000046653.13 | ENSG00000046653.13 |
| 21 | hsa.mir.1301 | ENSG00000151322.17 | ENSG00000151322.17 |
| 22 | hsa.mir.138.1 | ENSG00000123892.10 | ENSG00000123892.10 |

1.00 Mb
48.6Mb                48.7M
Chromosome bands        q21.32
Contigs        < AC103702.3
Genes
(Comprehensive set ...

2 >   < HOXB2   < HOXB5   CTD-2377D24.(
  < HOXB1   < HOXB4   < HOXB9   < C
  HOXB-AS1 >   < HOXB7   < F
  < HOXB3   < HOXB6
  HOXB-AS3 >   < HOXB8
  HOXB-AS2 >   < HOXB-AS4 >
    < MIR10A   < MIR196A1
      < RP11-357H14.1

# Close examination of prediction error

- Confusion matrix
- Local feature importance
- Force plots

# *RNAseq Bias Removal*

*Alex Partin*

# *Cancer Type Classification with SNPs*

# GDC: 10K samples with 10M mutations

## Somatic Mutations

### Overall Survival Plot

**17,424 Cases with Survival Data**

drag to zoom



Survival Rate (y-axis: 0.0, 0.2, 0.4, 0.6, 0.8, 1.0)
Duration (years) (x-axis: 0, 5, 10, 15, 20, 25, 30)

Showing **1 - 10** of **3,142,246** somatic mutations

JSON | TSV | Save/Edit Mutation Set

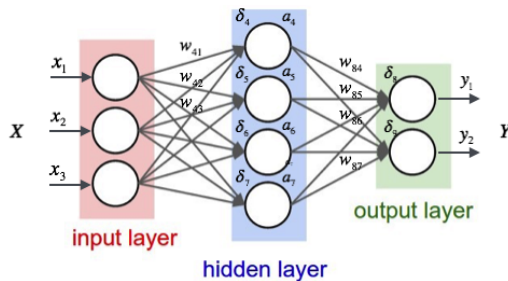| DNA Change | Type | Consequences | # Affected Cases in Cohort | # Affected Cases Across the GDC | Impact | Survival |
|---|---|---|---|---|---|---|
| chr7:g.140753336A>T | Substitution | **Missense** BRAF V600E | 565 / 10,202  5.54% | 565 / 10,202 ◄ | MO DH PR | 📉 |
| chr2:g.208248388C>T | Substitution | **Missense** IDH1 R132H | 388 / 10,202  3.80% | 388 / 10,202 ◄ | MO DL PO | 📉 |
| chr3:g.179218303G>A | Substitution | **Missense** PIK3CA E545K | 258 / 10,202  2.53% | 258 / 10,202 ◄ | MO DH PR | 📉 |
| chr3:g.179234297A>G | Substitution | **Missense** PIK3CA H1047R | 234 / 10,202  2.29% | 234 / 10,202 ◄ | MO TO PO | 📉 |
| chr12:g.25245350C>T | Substitution | **Missense** KRAS G12D | 208 / 10,202  2.04% | 208 / 10,202 ◄ | MO DH BE | 📉 |
| chr12:g.25245350C>A | Substitution | **Missense** KRAS G12V | 176 / 10,202  1.73% | 176 / 10,202 ◄ | MO DH PO | 📉 |
| chr3:g.179218294G>A | Substitution | **Missense** PIK3CA E542K | 167 / 10,202  1.64% | 167 / 10,202 ◄ | MO DH PR | 📉 |
| chr17:g.7675088C>T | Substitution | **Missense** TP53 R175H | 156 / 10,202  1.53% | 156 / 10,202 ◄ | MO TO BE | 📉 |
| chr17:g.7673803G>A | Substitution | **Missense** TP53 R273C | 125 / 10,202  1.23% | 125 / 10,202 ◄ | MO DH PR | 📉 |

# How to deal with n >> p ?

- Dropout

- Regularization

- Locally connected networks
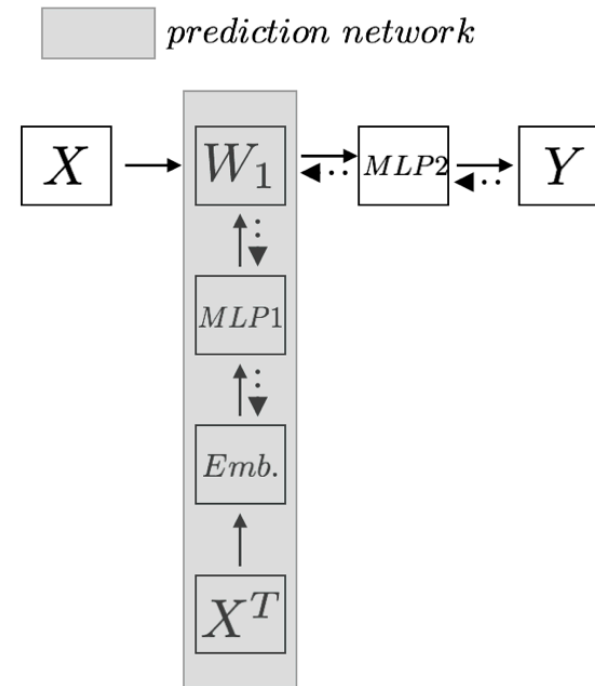
**Diet Networks: Thin Parameters for Fat Genomics**

Adriana Romero, Pierre Luc Carrier, Akram Erraqabi, Tristan Sylvain, Alex Auvolat, Etienne Dejoie, Marc-André Legault, Marie-Pierre Dubé, Julie G. Hussin, Yoshua Bengio

# Diet Network

- ## Suppose we have

  - 1000 samples
  - 1,000,000 features
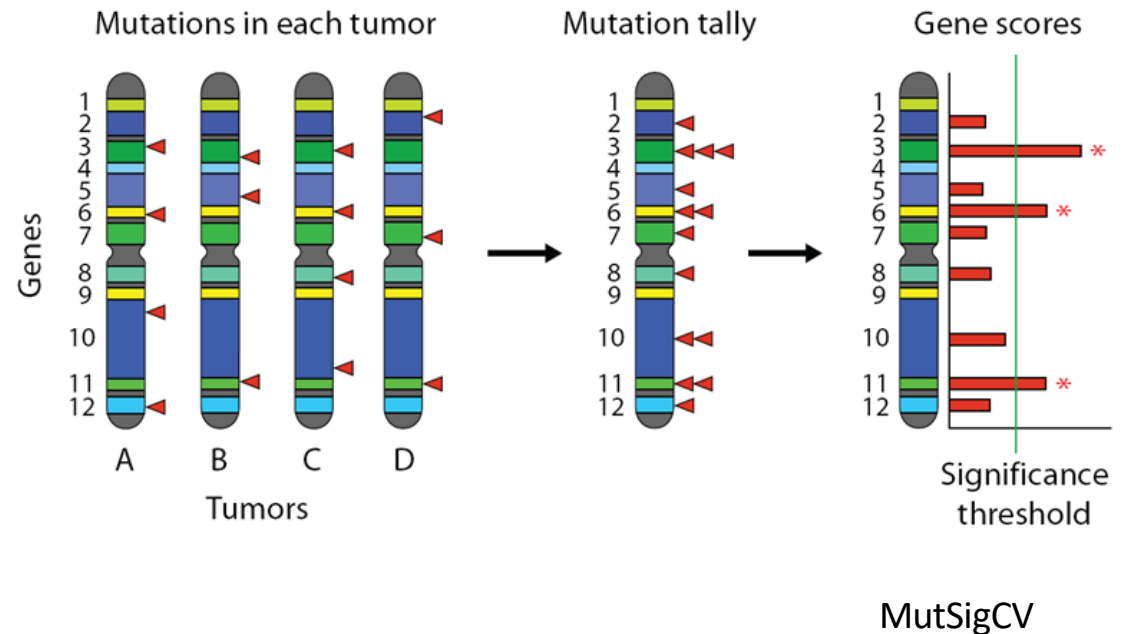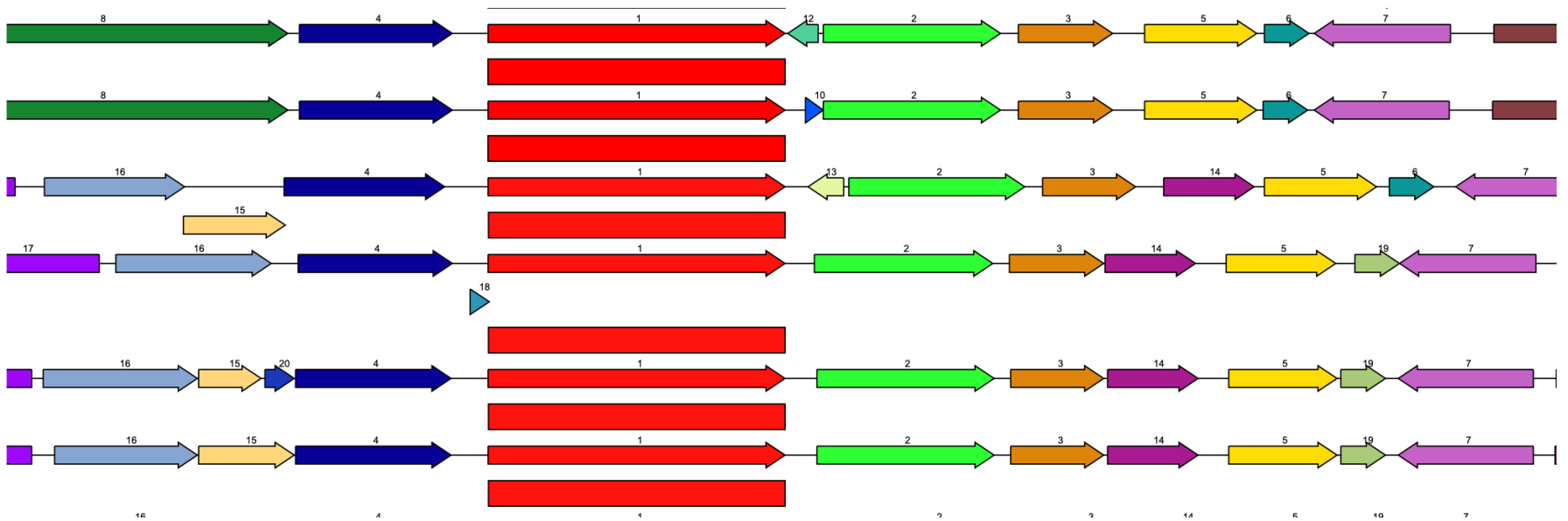  - 100 neurons in the hidden layer



  - Parameters in first layer = 100M

# How do we represent the sparse mutations?

- Gene level / pathway level

- Weighting by impact

- Filtering by significance

- Convert to images
  - Variant calling
  - Annotation



MutSigCV

# Deep annotation with compare region images

# Identification of genomic islands and operons

**Table 3.** Quality of tools predictions: % predictions made with GI features | % predictions missed with GI features, over the testing genomes dataset.

| Predictor \ Target | ShutterIsland | IslandViewer | AlienHunter | IslandPick | SIGI | **Average** |
|---|---|---|---|---|---|---|
| ShutterIsland | N/A | 91% \| 64% | 87% \| 47% | 89% \| 31% | 87% \| 36% | **89% \| 45%** |
| IslandViewer | 94% \| 67% | N/A | 89% \| 45% | 80% \| n/a | 87% \| n/a | **88% \| 56%** |
| AlienHunter | 74% \| 70% | 66% \| 60% | N/A | 73% \| 21% | 71% \| 42% | **71% \| 48%** |
| IslandPick | 69% \| 76% | 34% \| 86% | 49% \| 53% | N/A | 54% \| 44% | **52% \| 65%** |
| SIGI | 67% \| 75% | 45% \| 77% | 48% \| 51% | 50% \| 35% | N/A | **53% \| 60%** |
| Dimob | n/a \| 66% | n/a \| 28% | n/a \| 43% | n/a \| 25% | n/a \| 23% | **n/a \| 37%** |
| Phispy | n/a \| 68% | n/a \| 70% | n/a \| 50% | n/a \| 33% | n/a \| 39% | **n/a \| 52%** |
| PhageFinder | n/a \| 68% | n/a \| 70% | n/a \| 50% | n/a \| 34% | n/a \| 39% | **n/a \| 52%** |
| Islander | n/a \| 75% | n/a \| 71% | n/a \| 51% | n/a \| 33% | n/a \| 39% | **n/a \| 54%** |
| Phaster | n/a \| 75% | n/a \| 71% | n/a \| 51% | n/a \| 33% | n/a \| 39% | **n/a \| 54%** |

Assaf et al. 2019

# WORD2VEC

WINDOW

THE QUICK BROWN FOX JUMPS OVER THE LAZY DOG

CLASSIFIERS

Output Layer
Softmax Classifier

Hidden Layer
Linear Neurons

Input Vector

Probability that the word at a randomly chosen, nearby position is "**abandon**"

... "**ability**"

... "**able**"

... "**zone**"

A '1' in the position corresponding to the word "ants"

10,000 positions

300 neurons

10,000 neurons

http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/

"Skip-Gram"
With one-hot encoded centre word, we can predict context words.

Hidden layer creates embeddings

# Gene sets from MSigDB



22,596 gene sets

**H** — **hallmark gene sets** are coherently expressed signatures derived by aggregating many MSigDB gene sets to represent well-defined biological states or processes.

**C1** — **positional gene sets** for each human chromosome and cytogenetic band.

**C2** — **curated gene sets** from online pathway databases, publications in PubMed, and knowledge of domain experts.

**C3** — **motif gene sets** based on conserved cis-regulatory motifs from a comparative analysis of the human, mouse, rat, and dog genomes.

**C4** — **computational gene sets** defined by mining large collections of cancer-oriented microarray data.

**C5** — **GO gene sets** consist of genes annotated by the same GO terms.

**C6** — **oncogenic gene sets** defined directly from microarray gene expression data from cancer gene perturbations.

**C7** — **immunologic gene sets** defined directly from microarray gene expression data from immunologic studies.

# Gene2vec



# Sample2image



Alena Harley, The Mystery of the Origin
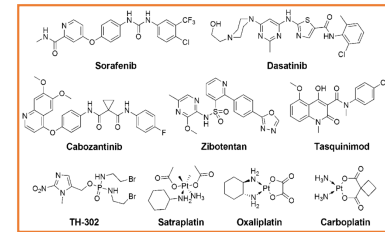
Confusion matrix

78% accuracy

Alena Harley, The Mystery of the Origin

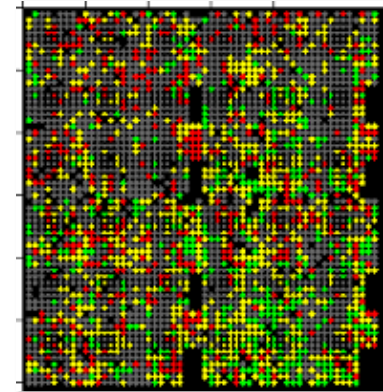# Drug Response Prediction

# Modeling Drug Response

**Drug (s)**
descriptors
fingerprints
structures
SMILES
dose



$$\mathcal{R} = f(\mathcal{T}, \mathcal{D}_1, \mathcal{D}_2)$$

IC50
GI50
% growth
Z-score
AUC
**Response**

gene expression levels
SNPs
protein abundance
microRNA
methylation
**Tumor**

# Cell Line Features

- NCI-60: 60 cell lines

- Molecular Assays: 20
  - Gene expression array
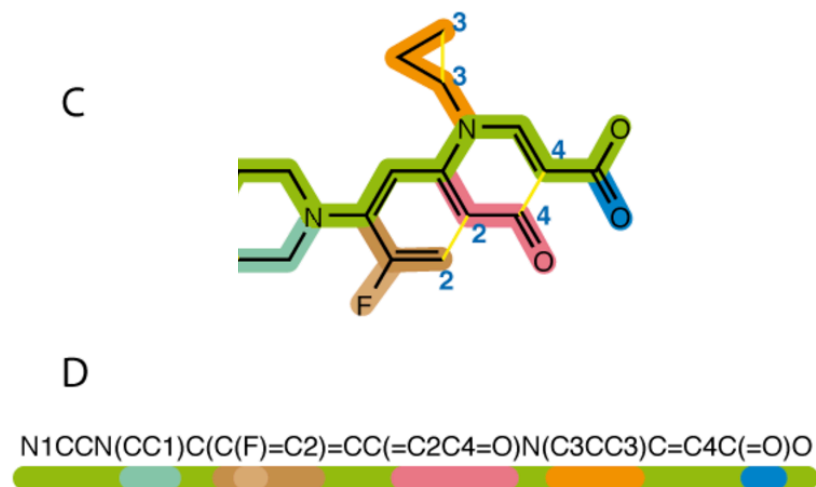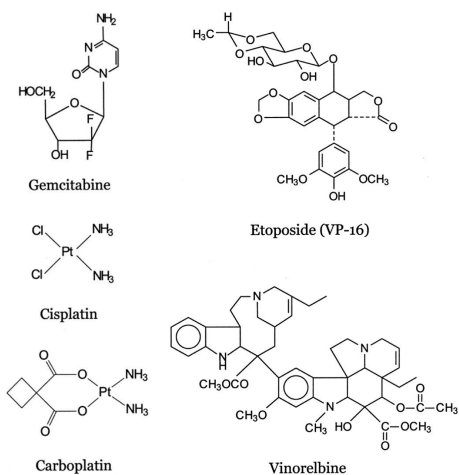  - RNA-seq
  - Mutations
  - Protein abundance
  - microRNA



Figure adapted from Kundaje et al. Nature 2015

# Drug Features

- SMILES strings
- 2D or 3D structures
- Graph convolutions
- Descriptors
- Fingerprints

C

D

N1CCN(CC1)C(C(F)=C2)=CC(=C2C4=O)N(C3CC3)C=C4C(=O)O

Gemcitabine

Etoposide (VP-16)

Cisplatin

Carboplatin

Vinorelbine

Paclitaxel

Docetaxel

Ifosfamide

| 0D | 1D | 2D | 3D | 4D |
|----|----|----|----|----|
| atom count | fragment counts | topological descriptors | geometrical | combination of atomic coordinates and sampling of conformations |
| molecular weight | e.g. # of OH # of NH | e.g. Weiner index Harrary index | atomic coordinates | |
| sum of atomic properties | | | energy grid | |

Over 4000 descriptors can be calculated by Dragon software

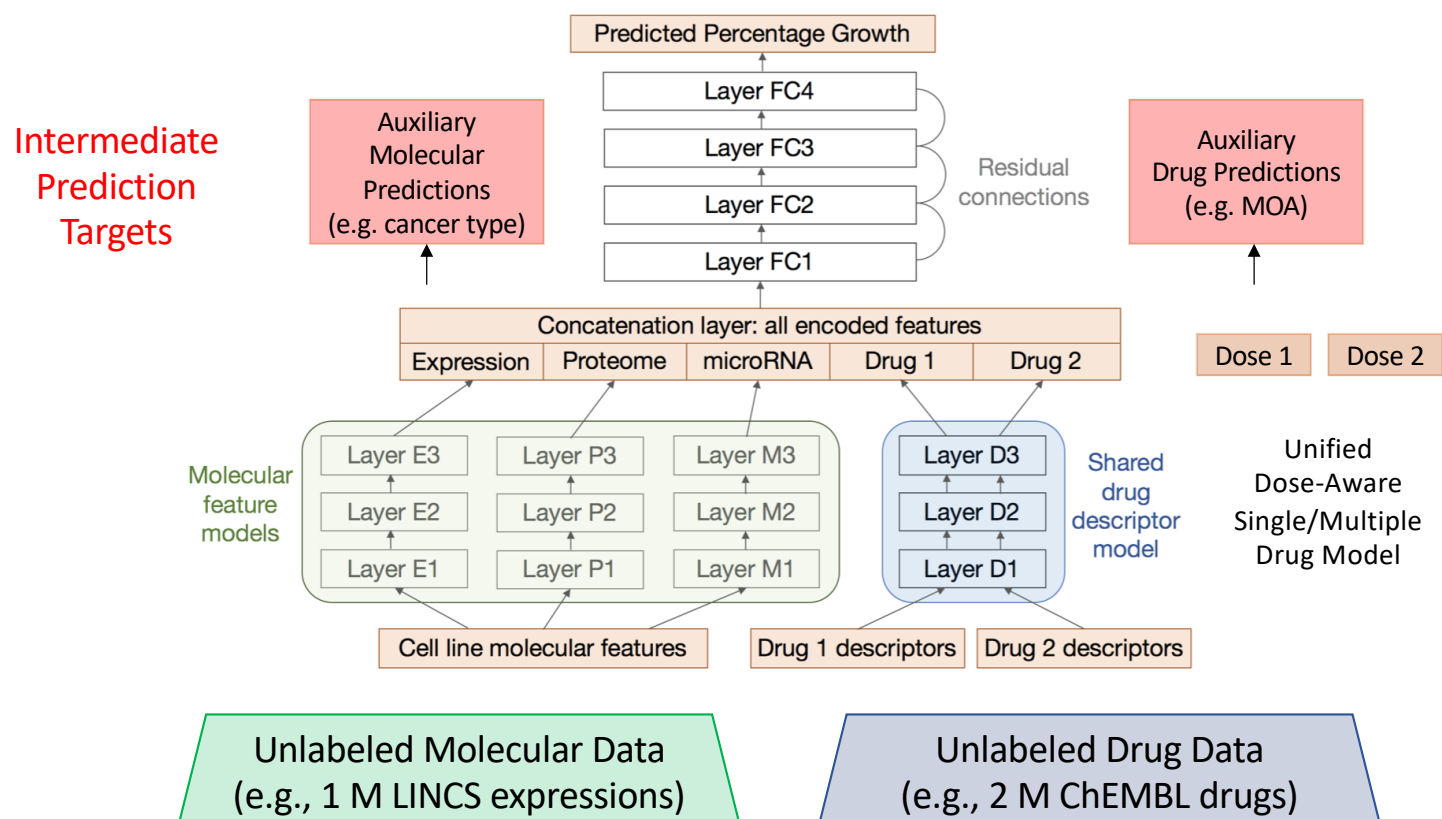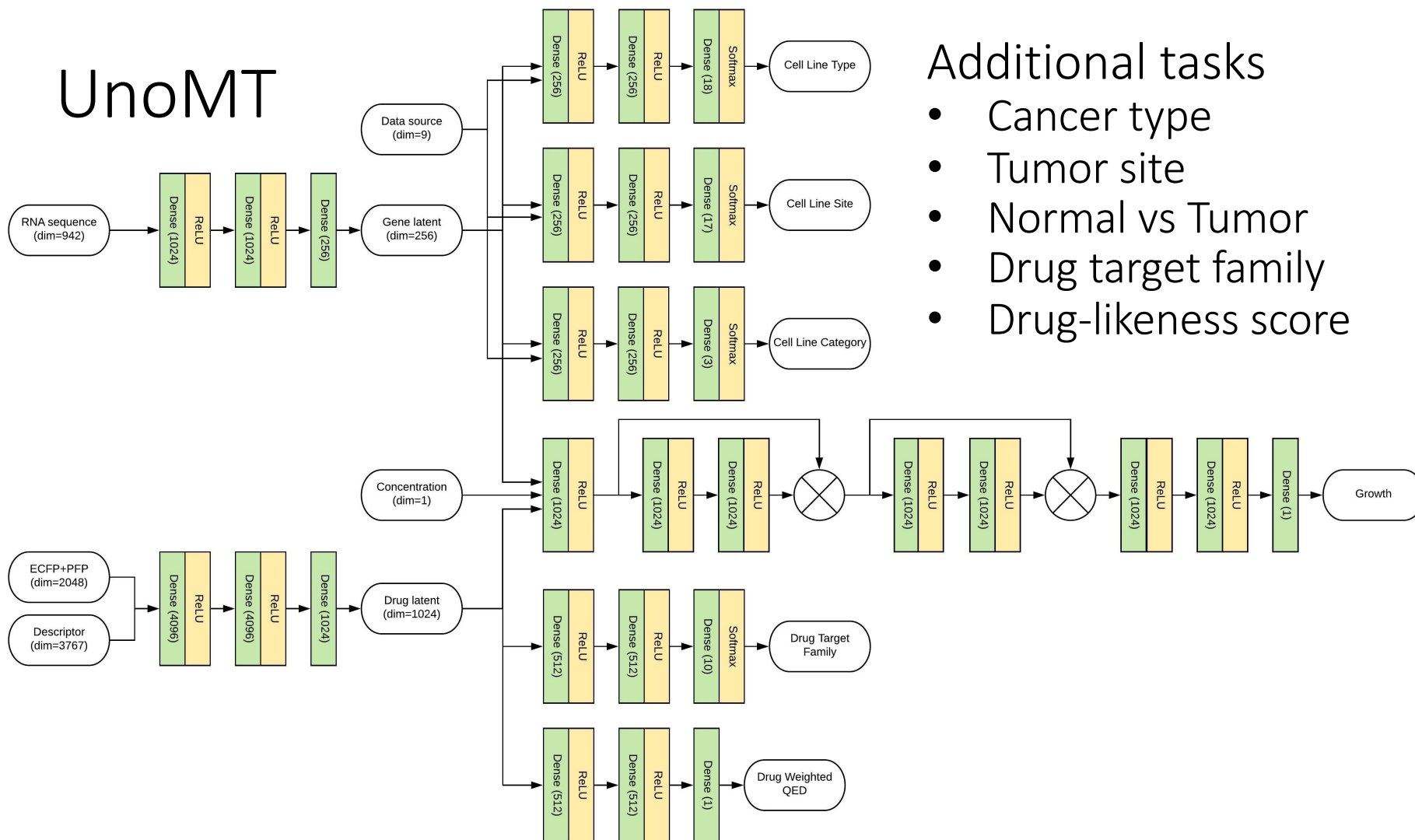# Deep Learning Model for Drug Pair Response



Fig. 2. **Neural network architecture.** The orange square boxes, from bottom to top, represent input features, encoded features, and output growth values. Feature models are denoted by round shaded boxes: green for molecular features and blue for drug features. There are multiple types of molecular features that are fed into submodels for gene expression, proteome, and microRNA. The descriptors for the two drugs share the same descriptor model. All encoded features are then concatenated to form input for the top fully connected layers. Most connecting layers are linked by optional residual skip connections if their dimensions match.
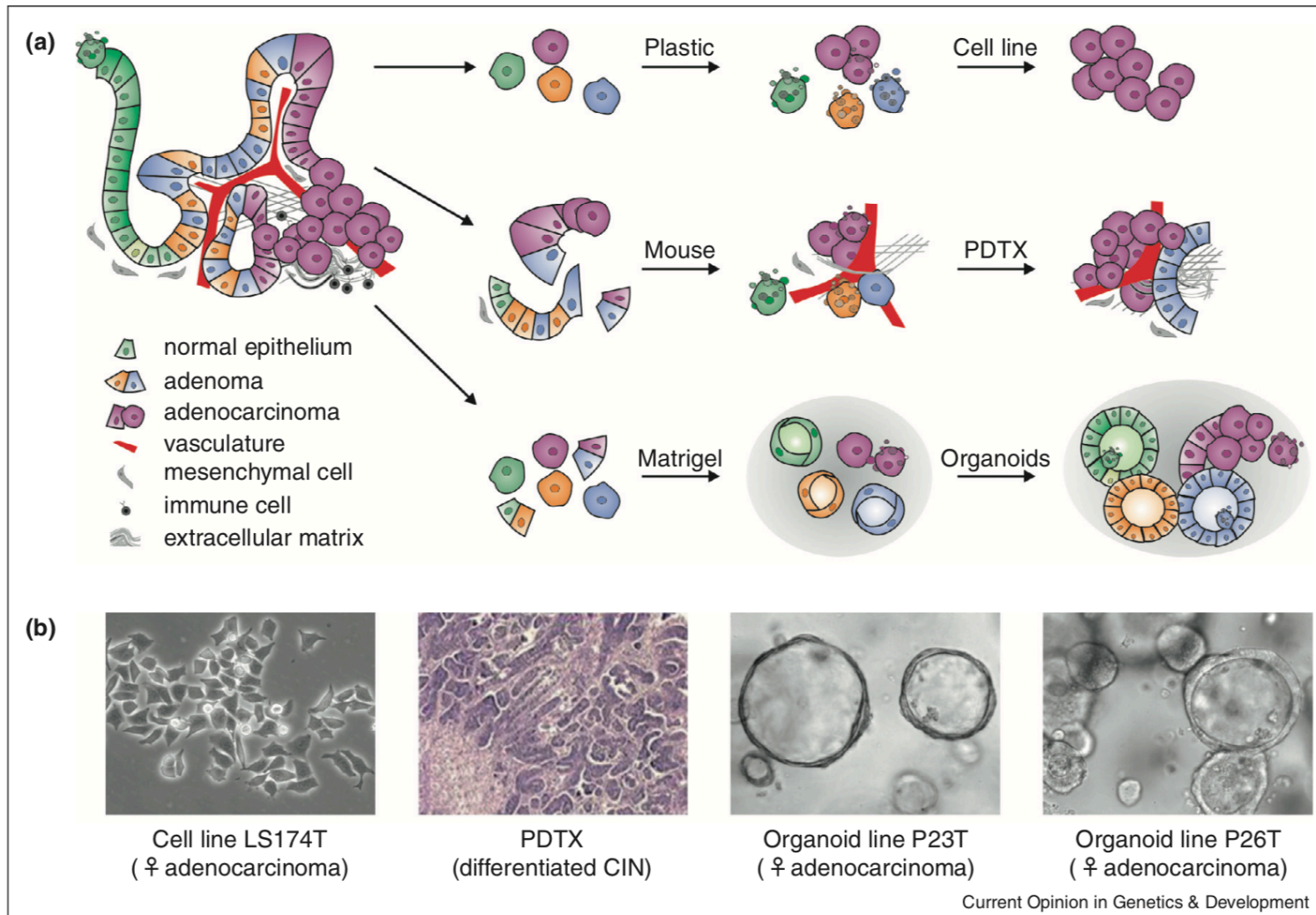
# Uno: Predicting Single/Paired Drug Response

# UnoMT

Additional tasks
- Cancer type
- Tumor site
- Normal vs Tumor
- Drug target family
- Drug-likeness score

(a)

| | Plastic | | Cell line |
| | Mouse | | PDTX |

normal epithelium
adenoma
adenocarcinoma
vasculature
mesenchymal cell
immune cell
extracellular matrix

Matrigel — Organoids

(b)

Cell line LS174T
(♀ adenocarcinoma)

PDTX
(differentiated CIN)

Organoid line P23T
(♀ adenocarcinoma)

Organoid line P26T
(♀ adenocarcinoma)

Current Opinion in Genetics & Development

Organoid cultures for the analysis of cancer phenotypes, Sachs and Clevers, 2014

# Data in place for model training and testing

**Table 1.** Integrating cell line, PDX, and real tumor samples across multiple studies

**Dose Independent**

| Data Source | # Tumor Samples | # Drugs | # Dose Response Samples | Treatment Type |
|---|---|---|---|---|
| NCI-ALMANAC | 60 | 104 | 3,686,475 | Drug pair |
| CCLE | 504 | 24 | 93,251 | Single drug |
| CTRPv2 | 887 | 544 | 6,171,005 | Single drug |
| gCSI | 409 | 16 | 58,094 | Single drug |
| GDSC | 1,075 | 249 | 1,894,212 | Single drug |
| NCI | 60 | 52,671 | 18,862,308 | Single drug |
| GDC | 11,081 | N/A | N/A | N/A |
| NCI-PDM | 1,198 | 12 | 518* | Single and paired drugs |

Dose Independent values (left column):
- 11,671
- 395,264
- 6,456
- 225,481
- 3,780,150 (60,000)

\* PDM drug response were measured differently from cell line dose response data.

# Patient Derived Xenograft Models

## Cancer Cell Lines



**Patient-derived xenografts (PDX) & conditionally reprogrammed cell lines**

Create reprogrammed cell lines

Tumorigenesis

Transplantation into NSG mice
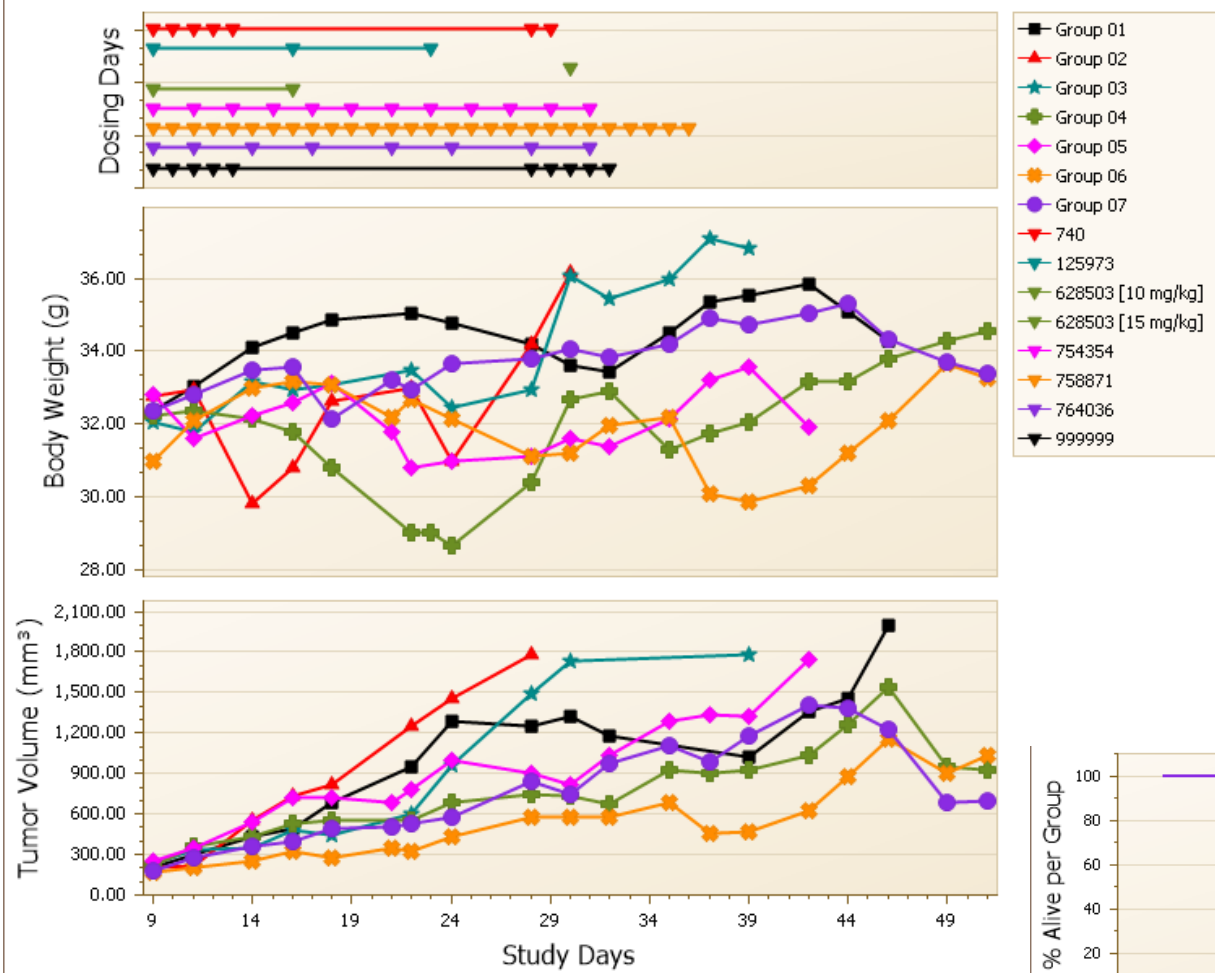
Tumor/patient heterogeneity

Nature Rev. Clin. Oncol. 11: 649-662, 2014.

## CL and PD Xenografts

ZFZJ2-1, 146476-266-R, Urothelial/bladder ca

AZD8055 ~2.2x delay. No regression

| Paclitaxel | 125973 |
|------------|--------|
| AZD8055 | 758871 |
| Dinaciclib | 764036 |
| Docetaxel | 628503 |
| GSK-461364 | 754354 |
| Methotrexate | 740 |

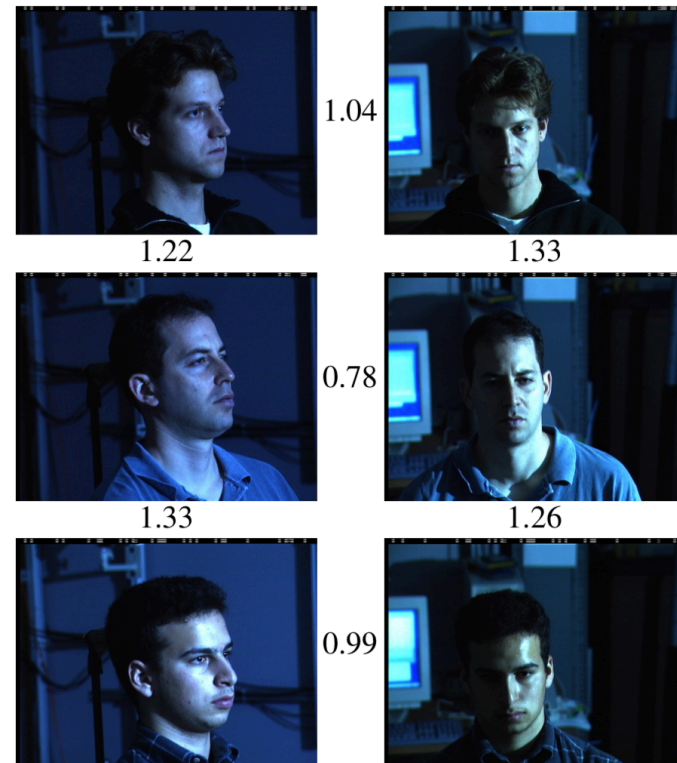# ZFZJ2-1, 146476-266-R, Urothelial/bladder ca



Op Meeting:

# *CANDLE prediction analysis notebook*

# Perceptual distance vs data distance

- Metric learning
- Representation learning
- Feature encoding
- Embedding
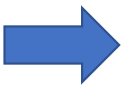- WordVec, ProteinVec

# Siamese network



Pairs could be gene expression replicates or samples from the same cancer type

# Focusing on the difficult parts

- ### Mine the difficult samples



Figure 2. **Model structure.** Our network consists of a batch input layer and a deep CNN followed by $L_2$ normalization, which results in the face embedding. This is followed by the triplet loss during training.
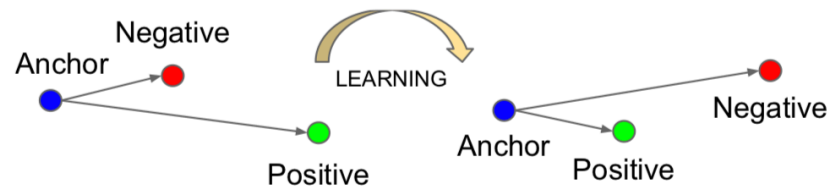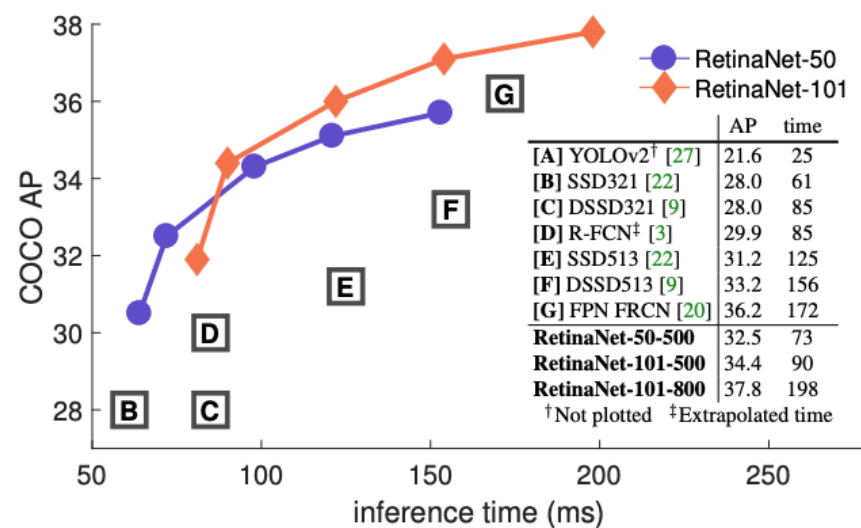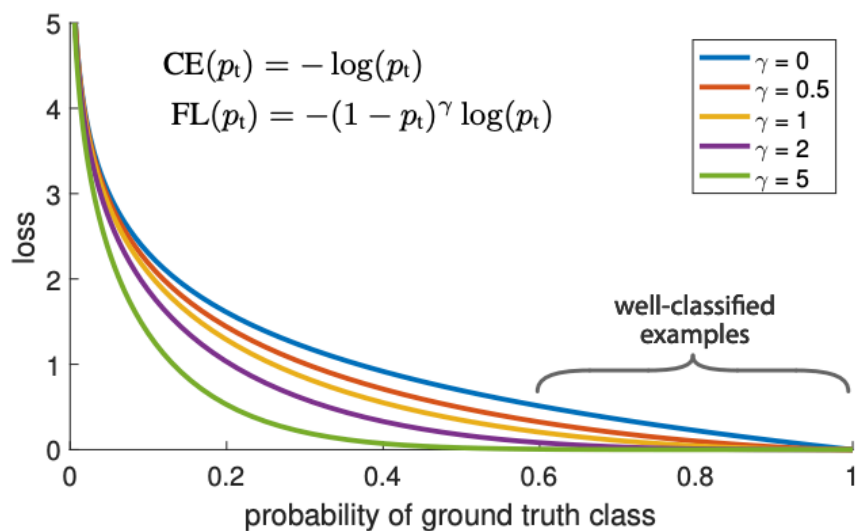
- ### Change loss function



Figure 3. The **Triplet Loss** minimizes the distance between an *anchor* and a *positive*, both of which have the same identity, and maximizes the distance between the *anchor* and a *negative* of a different identity.

# Focal Loss for Dense Object Detection

Tsung-Yi Lin    Priya Goyal    Ross Girshick    Kaiming He    Piotr Dollár

Facebook AI Research (FAIR)

*DL Challenge:*

*from surface pattern recognition*
*to deep mechanistic understanding*

# Just fancy regression?

Sequence to sequence models

Function approximation

| Input | | Output |
|---|---|---|
| 1+1 | = | 2 |
| 19+28 | = | 47 |
| 577+45 | = | 622 |
| 👍+👍 | = | ✌️ |
| 1️⃣9️⃣+2️⃣8️⃣ | = | 4️⃣7️⃣ |
| 五七七+四五 | = | 六二二 |

```c
#include<stdio.h>
void main(int argc, char *argv[])
{
    int k,r;
    int i=0,j=1,f;
    int sum=1;

    r=10;
    for (k=2; k<r; k++) {
        f=i+j;
        i=j;
        j=f;
        sum=sum+j;
    }

    printf("%d\n", sum);
}
```
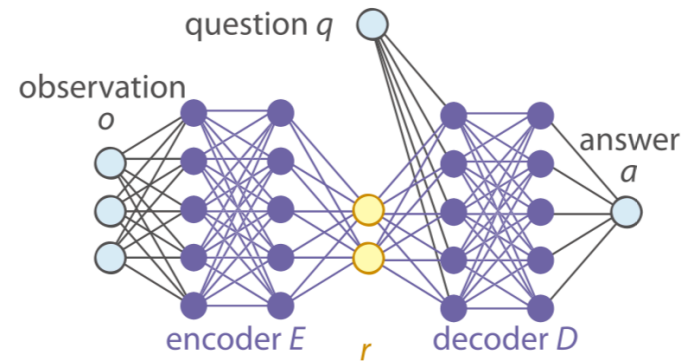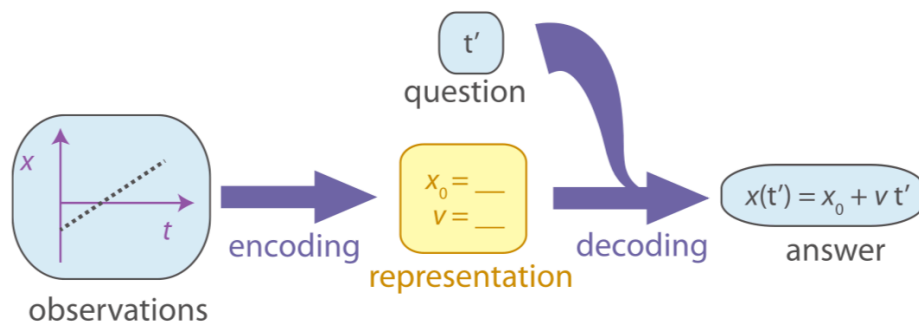
=>   88

# Discovering physical concepts with neural networks

Raban Iten, Tony Metger, Henrik Wilming, Lidia del Rio, Renato Renner

We introduce a neural network architecture that models the physical reasoning process and that can be used to extract simple physical concepts from experimental data without being provided with additional prior knowledge. We apply the neural network to a variety of simple physical examples in classical and quantum mechanics, like damped pendulums, two-particle collisions, and qubits. The network finds the physically relevant parameters, exploits conservation laws to make predictions, and can be used to gain conceptual insights. For example, given a time series of the positions of the Sun and Mars as observed from Earth, the network discovers the heliocentric model of the solar system – that is, it encodes the data into the angles of the two planets as seen from the Sun. Our work provides a first step towards answering the question whether the traditional ways by which physicists model nature naturally arise from the experimental data without any mathematical and physical pre-knowledge, or if there are alternative elegant formalisms, which may solve some of the fundamental conceptual problems in modern physics, such as the measurement problem in quantum mechanics.

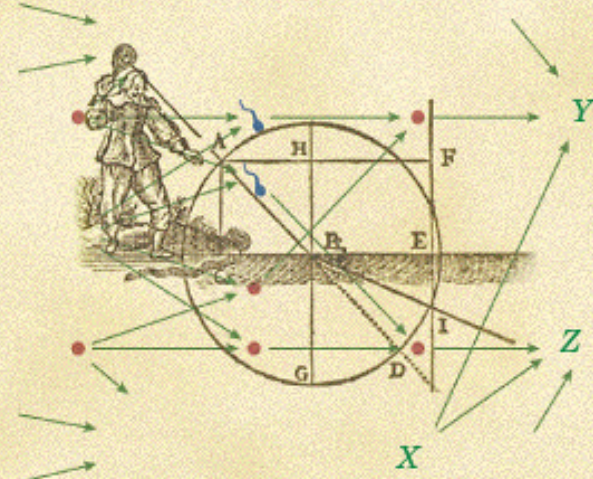JUDEA PEARL
*WINNER OF THE TURING AWARD*
AND DANA MACKENZIE

# THE
# BOOK OF
# WHY

α ➤ β

THE NEW SCIENCE
OF CAUSE AND EFFECT



# CAUSALITY

MODELS, REASONING,
AND INFERENCE

## JUDEA PEARL

# Thank you