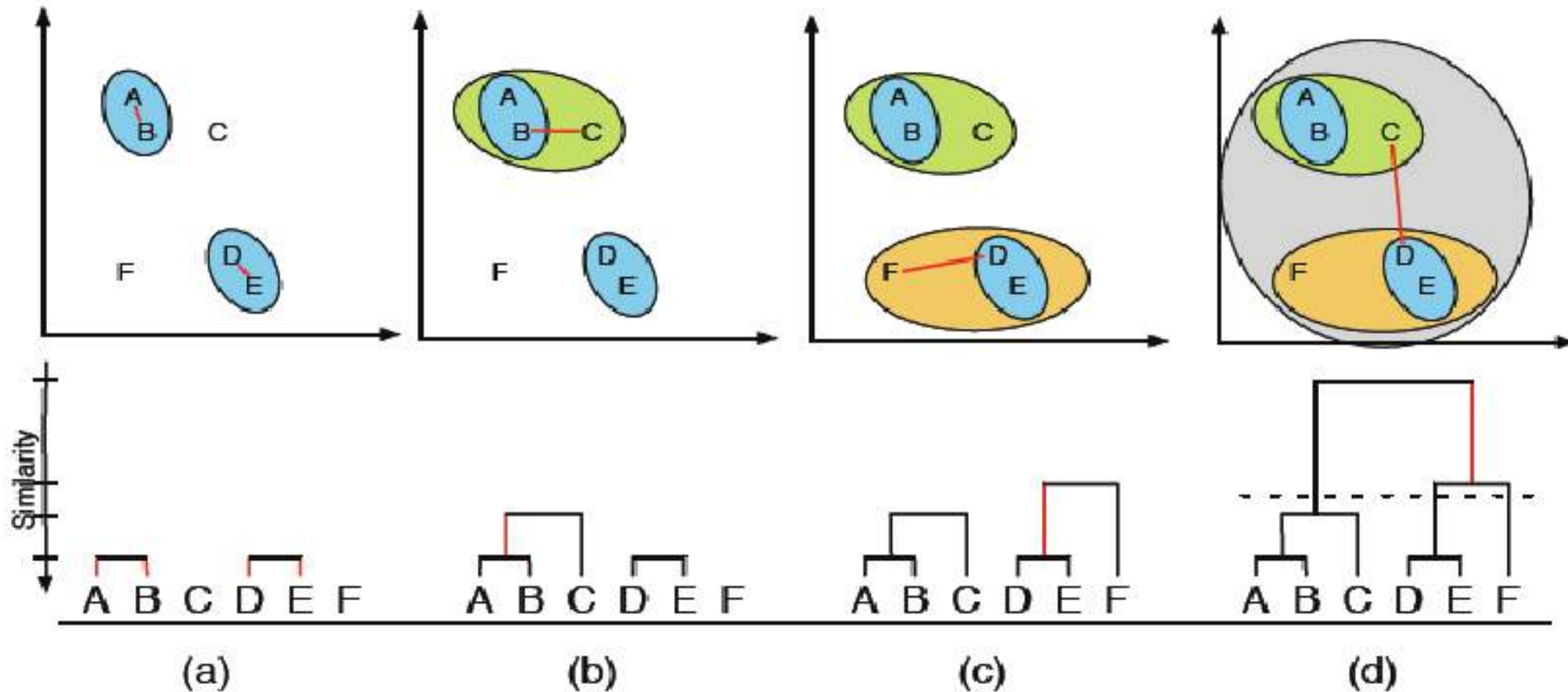


# Hierarchical clustering

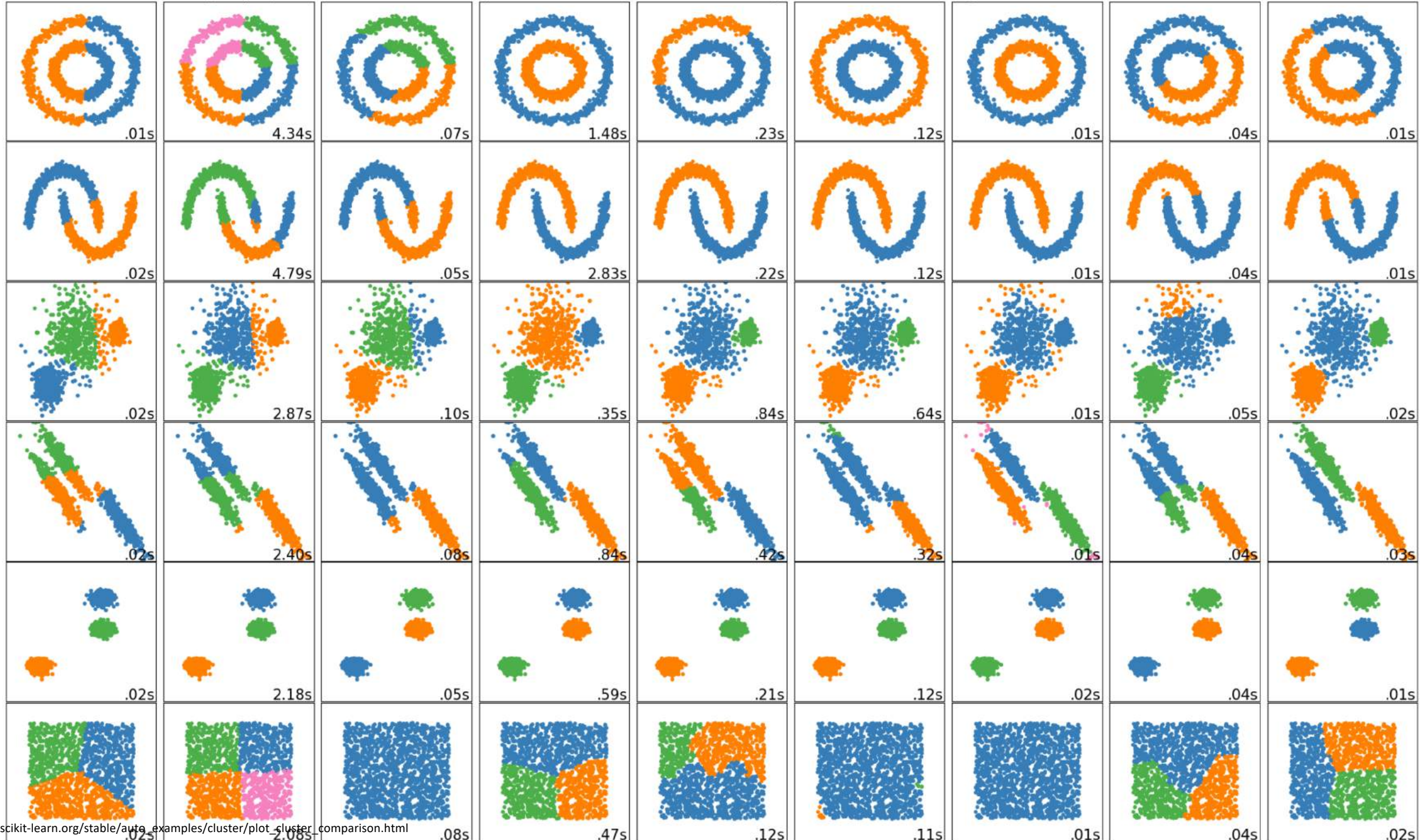
## Example: Hierarchical Agglomerative Clustering



# Hierarchical clustering

- Usually most computationally efficient
- Truly deterministic(?)
  - Are other forms of clustering deterministic?

MiniBatchKMeans AffinityPropagation MeanShift SpectralClustering Ward AgglomerativeClustering DBSCAN Birch GaussianMixture

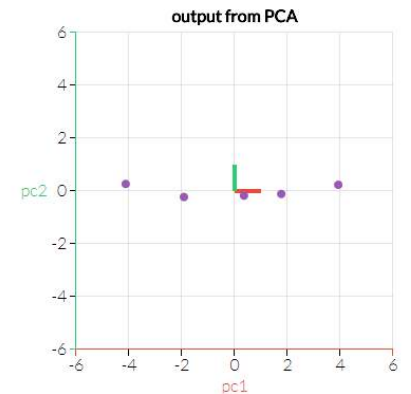
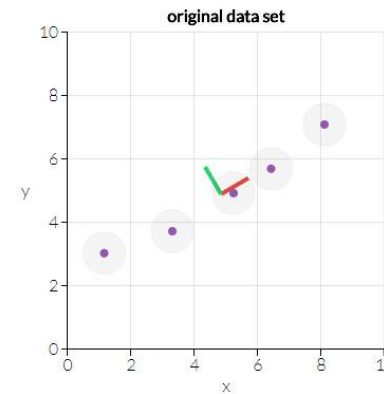


# Pitfalls of clustering

- Choice of method affects clusters identified
- How many clusters? How do you determine #?
- Non-deterministic (depending on seeding)
- Do all datapoints share similar profile?
- Check percent of data points in each cluster (e.g.5%-35%)
- Outlier treatment

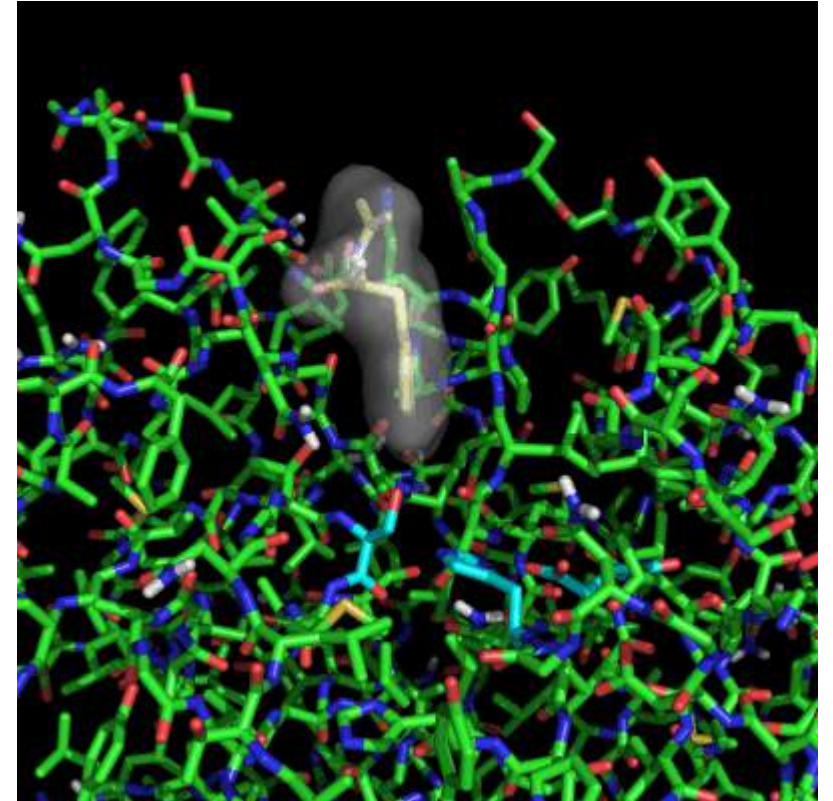
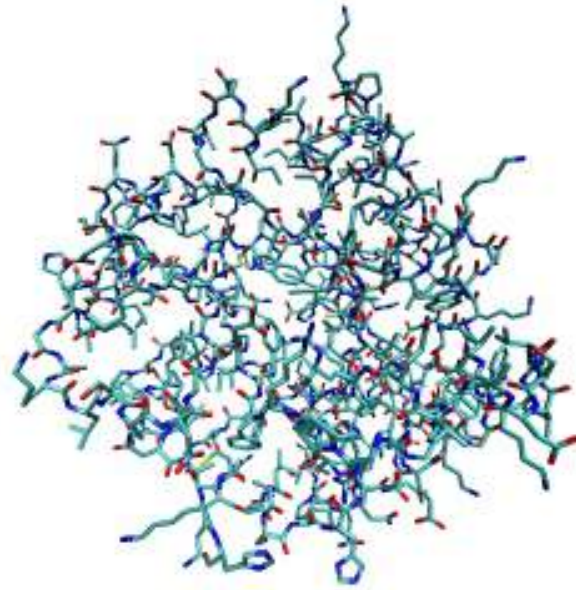
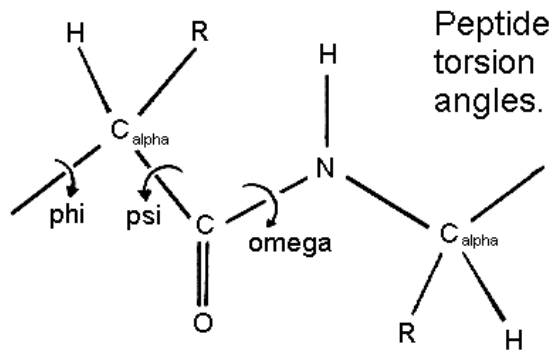
# Decomposition/Dimensionality Reduction

- Useful to perform as preprocessing
- Can tell you about the structure of your data
- **PCA** is one of the most common methods
  - ICA, TICA, etc.

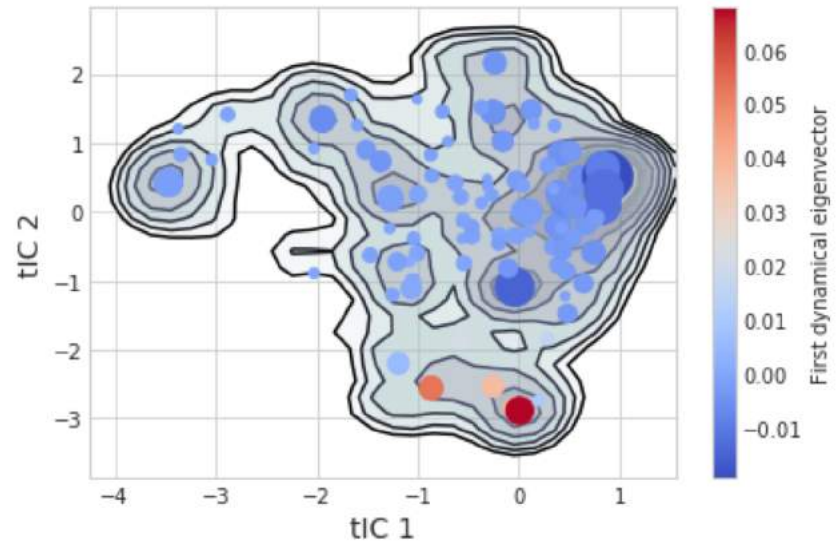
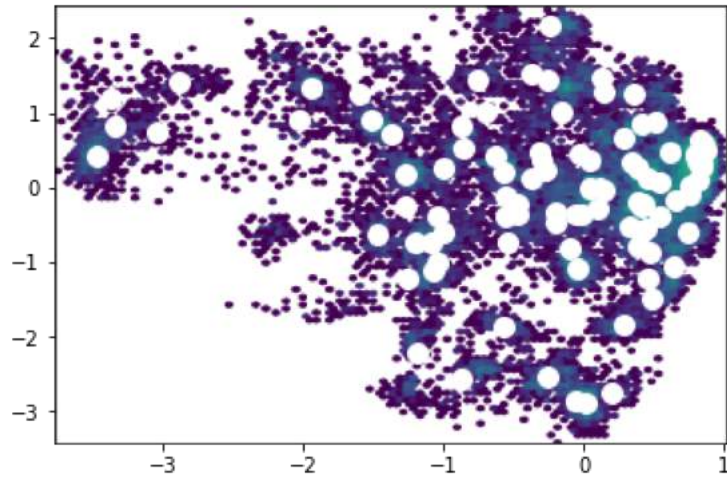


- <http://setosa.io/ev/principal-component-analysis/>

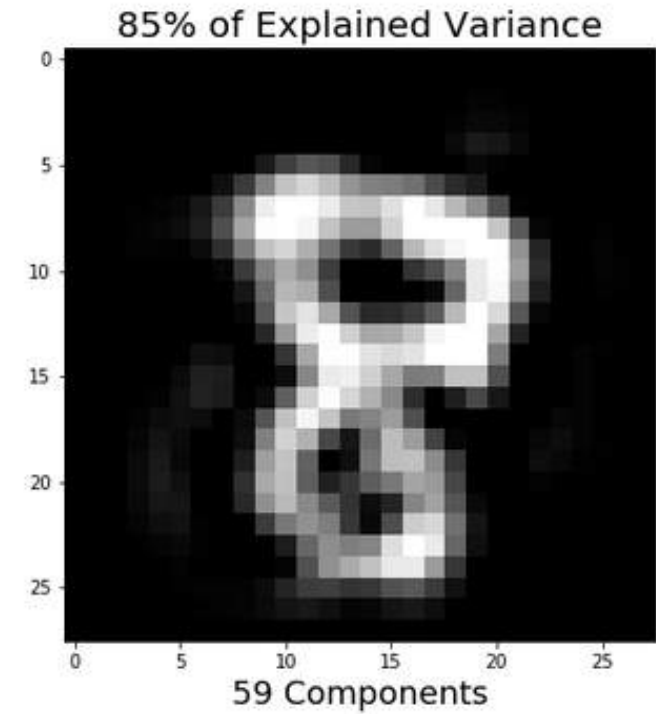
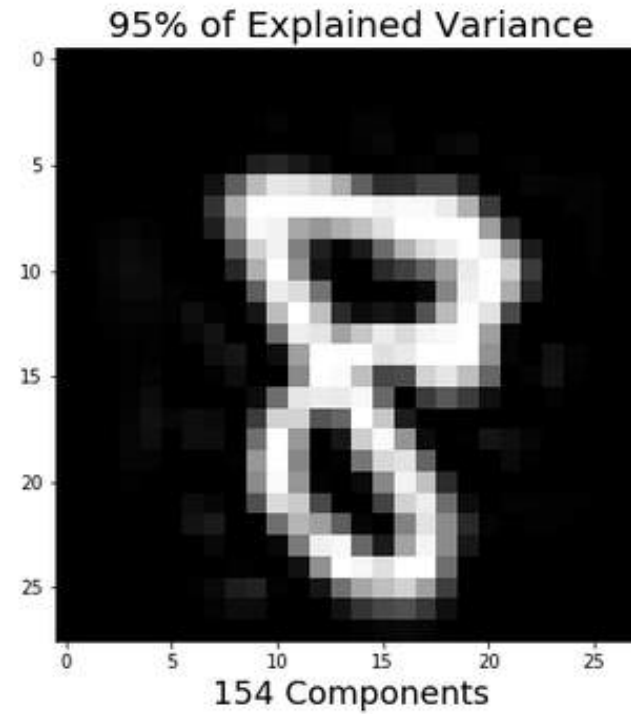
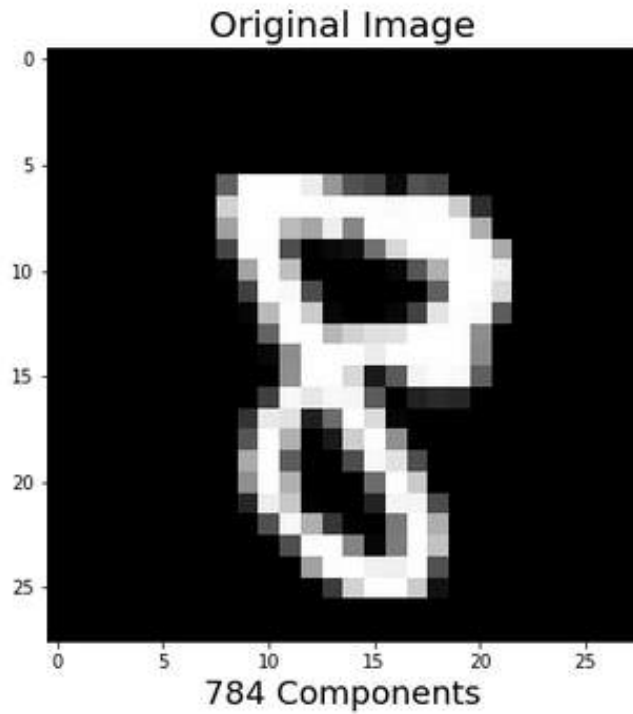
# Personal Example



# Personal Example



# PCA – Image Compression





# Sklearn pipelines

```
pipeline = Pipeline([
    ('ordinal_to_nums', DFTransform(_ordinal_to_nums, copy=True)),
    ('union', DFFeatureUnion([
        ('categorical', Pipeline([
            ('select', DFTransform(lambda X: X.select_dtypes(include=['object']))),
            ('fill_na', DFTransform(lambda X: X.fillna('NA'))),
            ('one_hot', DFTransform(_one_hot_encode)),
        ])),
        ('numerical', Pipeline([
            ('select', DFTransform(lambda X: X.select_dtypes(exclude=['object']))),
            ('fill_median', DFTransform(lambda X: X.fillna(X.median()))),
            ('add_features', DFTransform(_add_features, copy=True)),
            ('remove_skew', DFTransform(_remove_skew, copy=True)),
            ('find_outliers', DFTransform(_find_outliers, copy=True)),
            ('normalize', DFTransform(lambda X: X.div(X.max())))
        ])),
    ])),
])
```

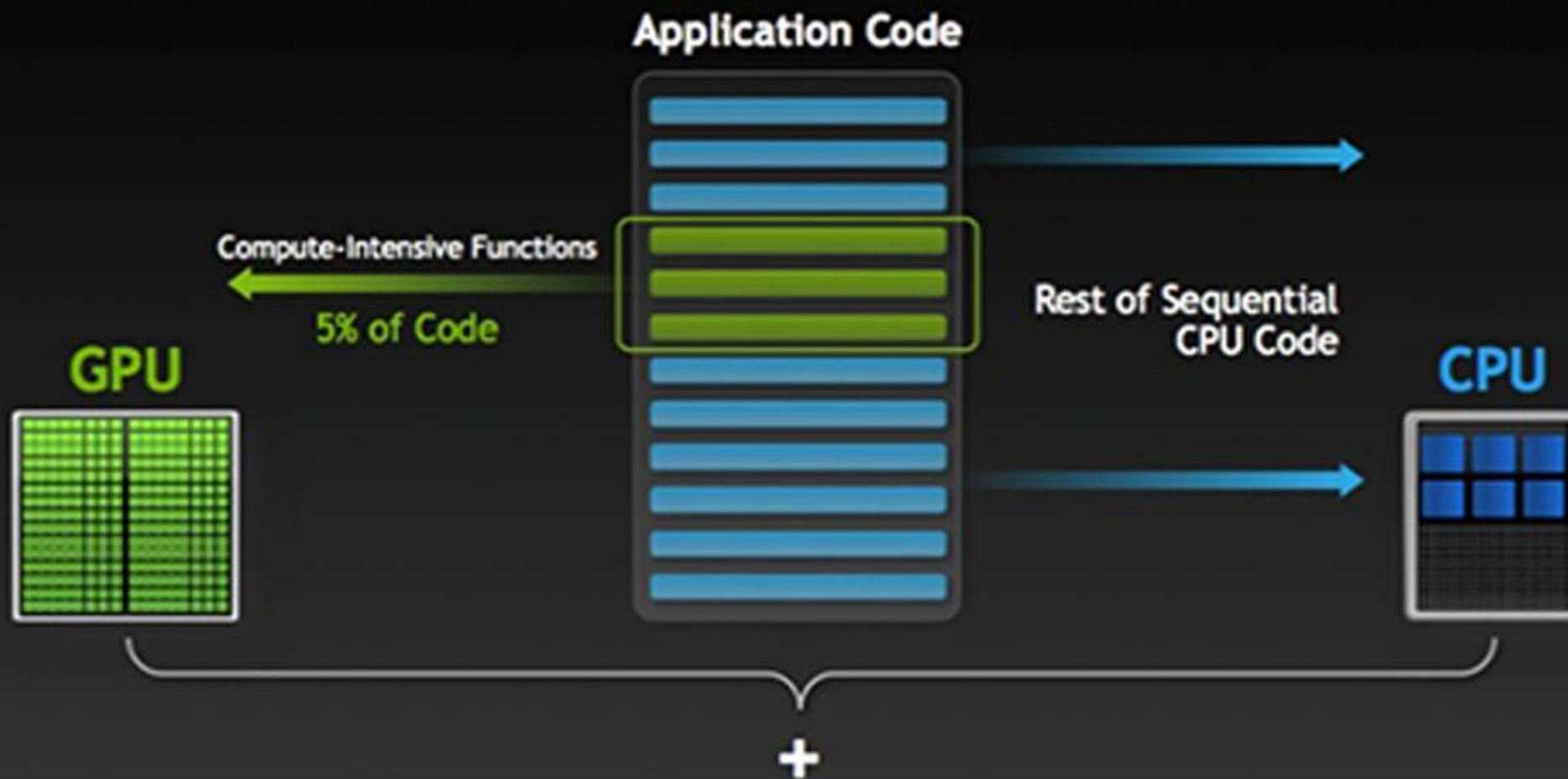
# Natural Language Processing

- SpaCy <https://nicschrading.com/project/Intro-to-NLP-with-spaCy/>
- Harry Potter <http://botnik.org/content/harry-potter.html>

# Advances that enable ML

- Computational Architecture
- Parallel processing
- High Performance Computing
- GPU-acceleration

# How GPU Acceleration Works



# Troubleshooting ML

- Overfitting?
  - Reduce # of features
    - Manually select
    - Model selection (sklearn has options for this)

```
>>> from sklearn.feature_selection import VarianceThreshold
>>> X = [[0, 0, 1], [0, 1, 0], [1, 0, 0], [0, 1, 1], [0, 1, 0], [0, 1, 1]]
>>> sel = VarianceThreshold(threshold=(.8 * (1 - .8)))
>>> sel.fit_transform(X)
array([[0, 1],
       [1, 0],
       [0, 0],
       [1, 1],
       [1, 0],
       [1, 1]])
```

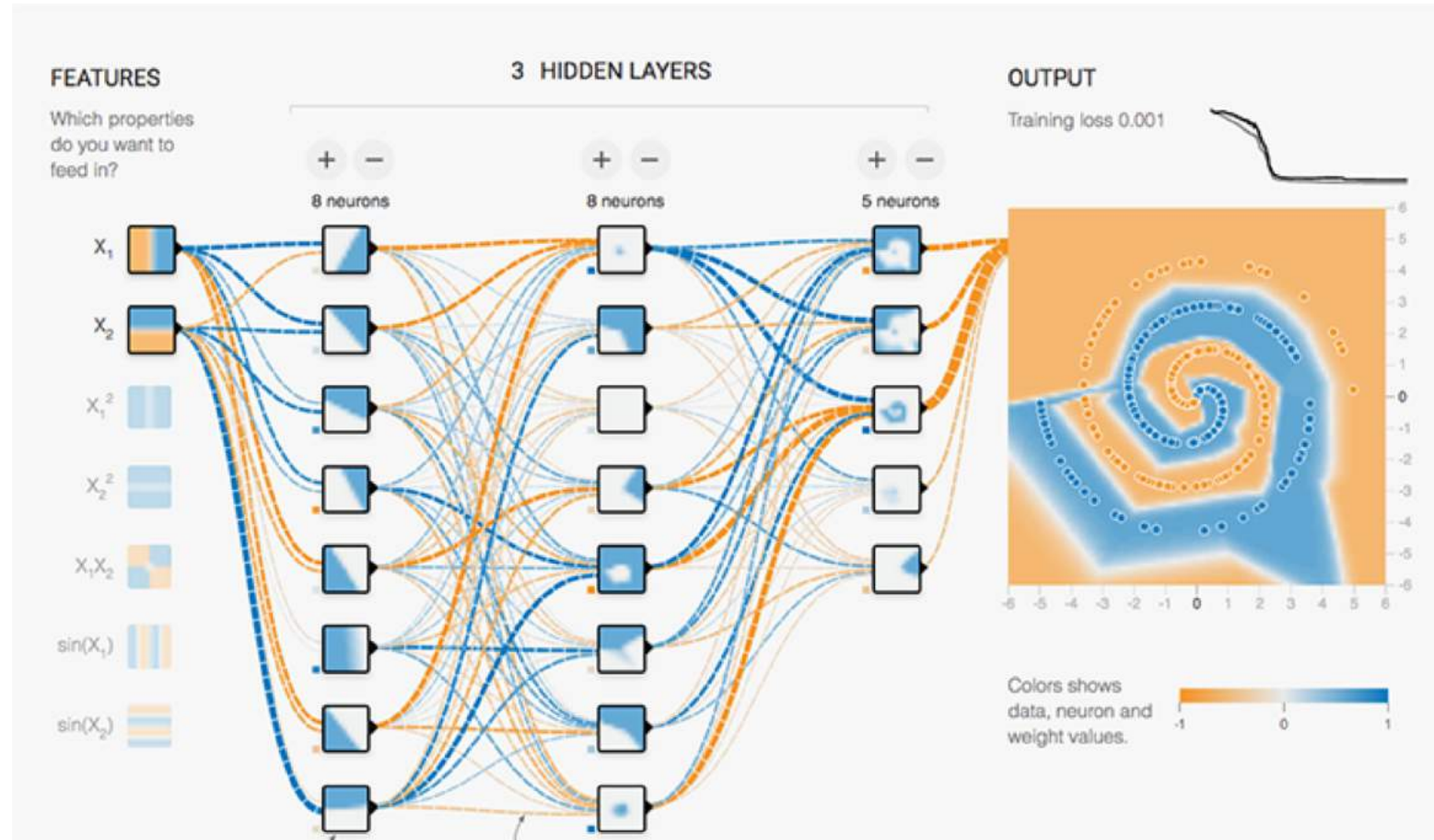
# Troubleshooting ML

- Overfitting?
  - Reduce # of features
    - Manually select
    - Model selection (sklearn has options for this)
  - Regularization
  - Revisit scaling of features

# Troubleshooting ML

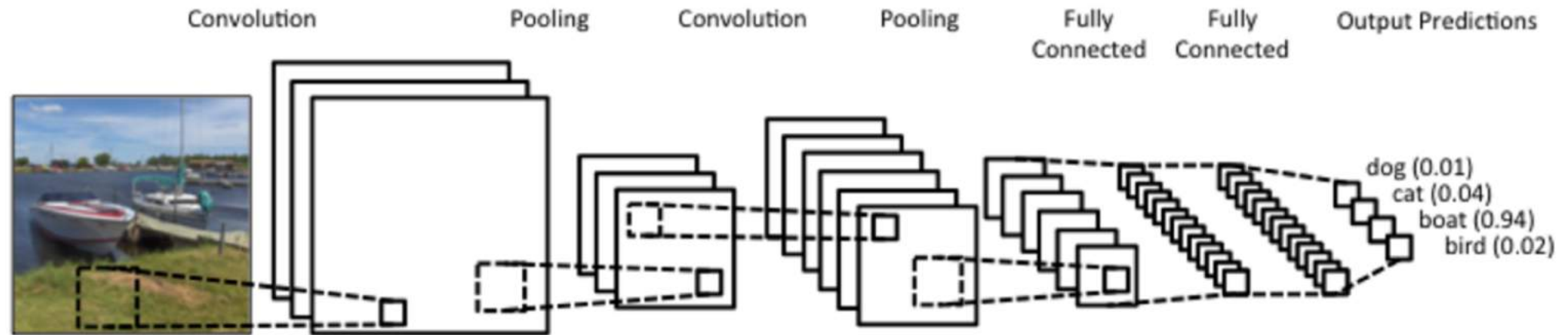
- Large errors
  - More training examples? (addressing high variance)
  - Smaller sets of features (addressing high variance)
  - Additional features (addressing high bias)
  - Adding polynomial features  $x^2$ ,  $x^3$ , etc. (sklearn has options for this as well) (addressing high variance)

# Deep Learning

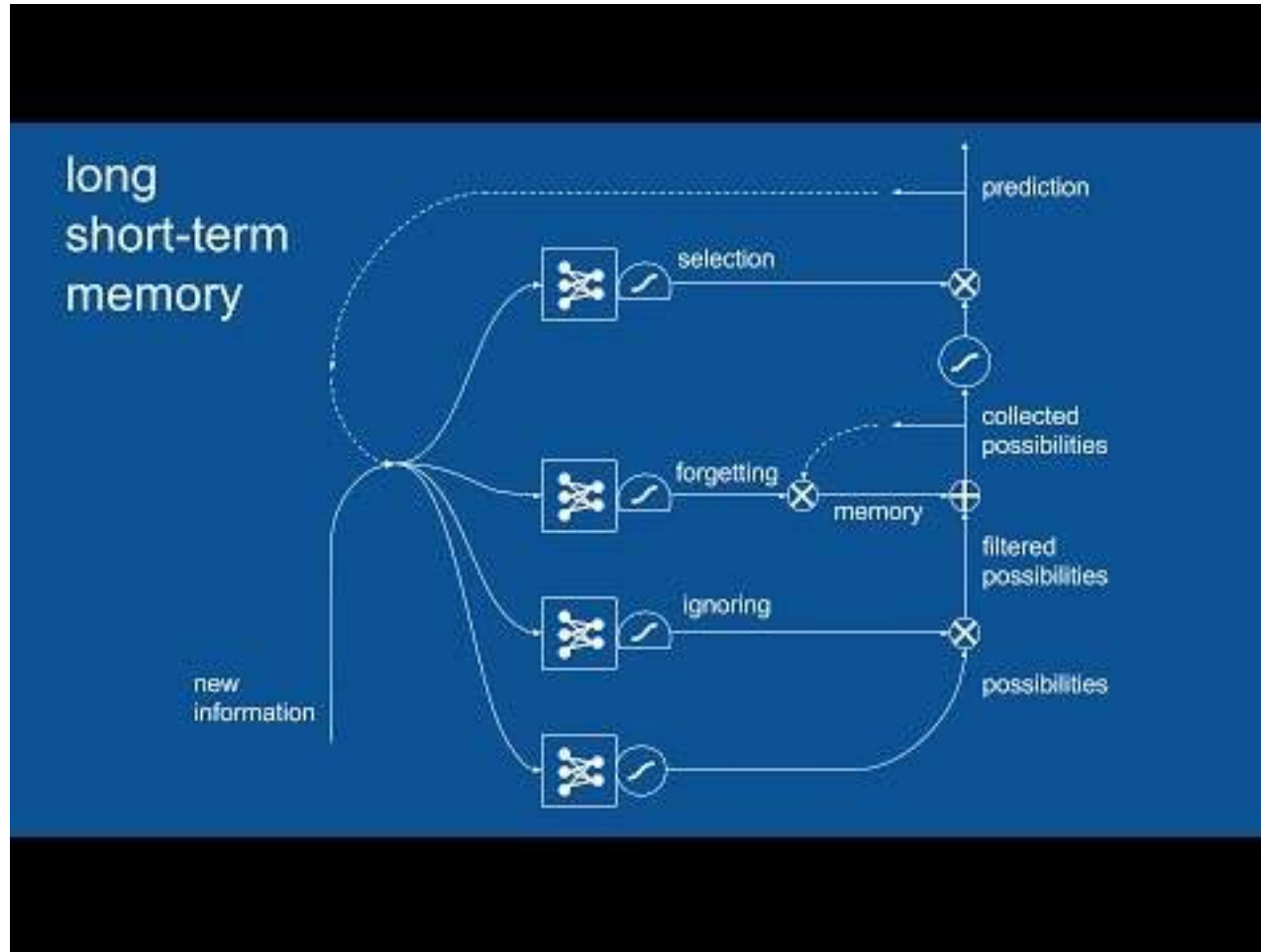




# Convolutional Neural Networks (CNNs)



# Long Short Term Memory



# Tensorflow



# TensorBoard

**TensorBoard** EVENTS IMAGES GRAPH HISTOGRAMS

**Fit to screen**

**Run** `cifar-train`

**Upload**

**Color** Structure  
color: same substructure  
gray: unique substructure

**Main Graph**

**Auxiliary nodes**

- Variable
- inc
- GradientDescent
- save
- avg
- gradients
- ExponentialDecay

**Graph** (\* = expandable)

- Namespace\*
- OpNode
- Unconnected series\*
- Connected series\*
- Constant
- Summary
- Dataflow edge
- Control dependency edge
- Reference edge

# Keras

## Keras: The Python Deep Learning library



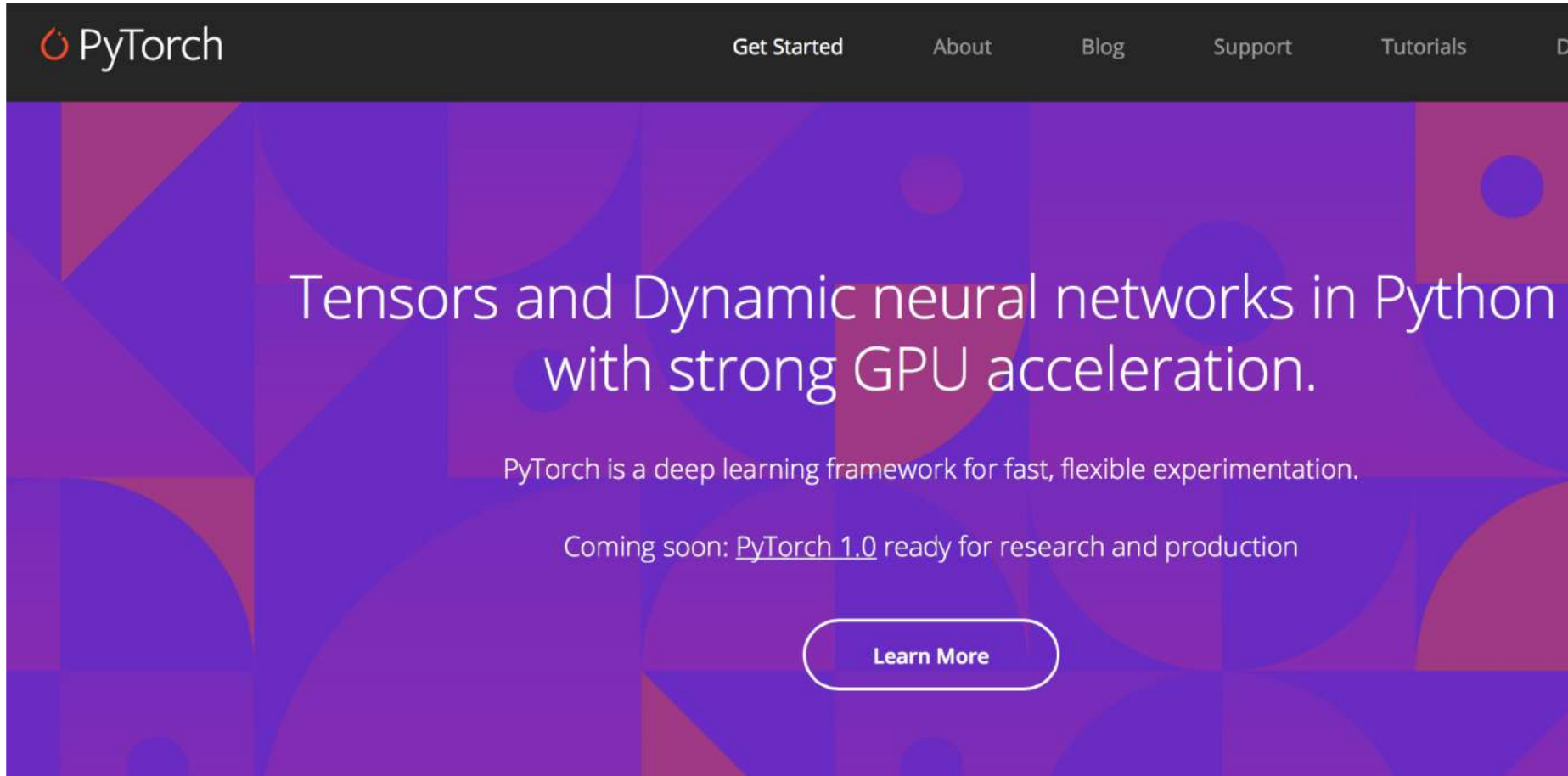
### You have just found Keras.

Keras is a high-level neural networks API, written in Python and capable of running on top of **TensorFlow**, **CNTK**, or **Theano**. It was developed with a focus on enabling fast experimentation. *Being able to go from idea to result with the least possible delay is key to doing good research.*

Use Keras if you need a deep learning library that:

- Allows for easy and fast prototyping (through user friendliness, modularity, and extensibility).
- Supports both convolutional networks and recurrent networks, as well as combinations of the two.
- Runs seamlessly on CPU and GPU.

# PyTorch



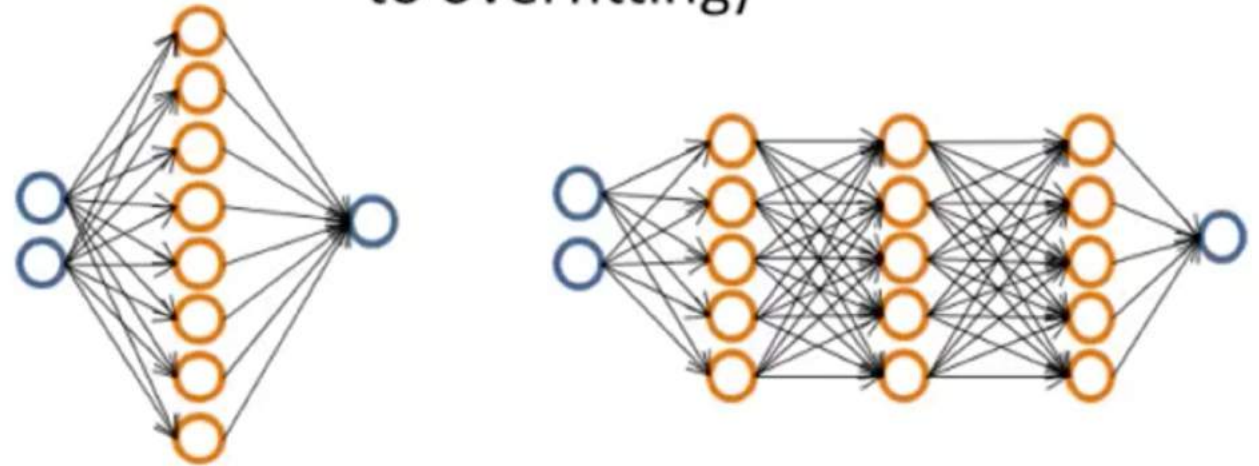
# Neural networks and overfitting

“Small” neural network  
(fewer parameters; more  
prone to underfitting)



Computationally cheaper

“Large” neural network  
(more parameters; more prone  
to overfitting)



Computationally more expensive.

Use regularization ( $\lambda$ ) to address overfitting.

# When is machine learning or deep learning appropriate?

- Consider the problem at hand
- How much data do you have available?
- How computationally intensive is your approach?
- Can your dataset be loaded into memory entirely?
- Do you have a GPU accessible?
- Interoperability (sklearn pickling)
- Deep learning hype



# Data Sources & Learning

- UCI
- Kaggle
- 538, NYT
  
- Coursera ML Andrew Ng
- Google machine learning course
- Fast.ai