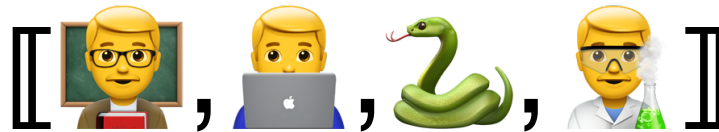


Lecture Notes for **Machine Learning in Python**



A History of Deep Learning And Introduction to Toolkits

Logistics and Agenda

- Logistics
 - Grading update
- Agenda
 - “Deep Learning” History
 - Deep Learning toolkits

Class Overview, by topic

Table Data
Visualization

Numpy, Pandas, Seaborn
Overviews with some in-depth discussion

Dimension
Reduction and
Image Processing

Scikit-learn, Scikit Image,
Intuition only, Some mathematics

Linear and
Logistic
Regression

Numpy, Recreate API for Scikit-learn
Detailed mathematics for simple optimization
intuition for advanced optimization

Neural Networks
and Back Prop.

Numpy
Detailed mathematics for NN operations

Wide and Deep
Networks

Convolutional
Networks

Recurrent
Networks

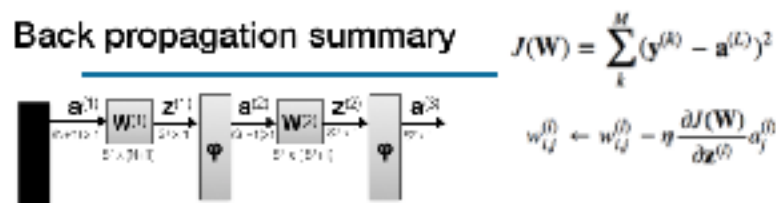
Keras, Tensorflow
Intuition, Detailed implement.

Ethics in
Language Models

ConceptNet
Case studies

Last time:

Back propagation summary



1. Forward propagate to get \mathbf{z}, \mathbf{a} for all layers
2. Get final layer gradient
3. Update back propagation variables
4. Update each $\mathbf{W}^{(l)}$

for each $y^{(i)}$

$$\frac{\partial J(W)}{\partial z^{(2)}} = -2(y^{(i)} - a^{(2)}) \times a^{(1)} + (1 - a^{(1)})$$

$$\frac{\partial J(W)}{\partial z^{(l)}} = \text{diag}[a^{(l+1)} * (1 - a^{(l+1)})] \cdot W^{(l+1)} \cdot \frac{\partial J(W)}{\partial z^{(l+1)}}$$

$$W^{(l)} \leftarrow W^{(l)} - \eta \frac{\partial J(W^{(l)})}{\partial z^{(l)}} \cdot a^{(l)}$$

Review

$$W_{k+1} = W_k - \rho_k$$

- Cross entropy

$$Y^{(2)} = A^{(2)} - Y$$

now final layer update

- Momentum

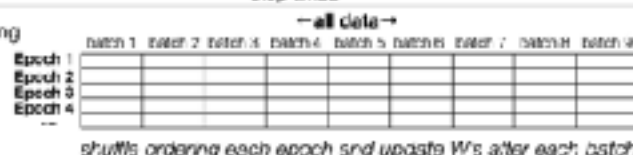
$$\rho_k = \alpha \nabla J(W_k) + \beta \nabla J(W_{k-1})$$

- Nesterov's Momentum

$$\rho_k = \beta \nabla J(W_k + \alpha \nabla J(W_{k-1})) + \alpha \nabla J(W_{k-1})$$

step twice

- Mini-batching

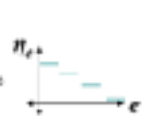


- Learning rate adjustment (eta)

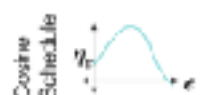
Star step

$$\eta_e = \eta_0 \cdot d^{\frac{1-e}{d}}$$

η_e epochs between reductions
 $0 < d < 1$

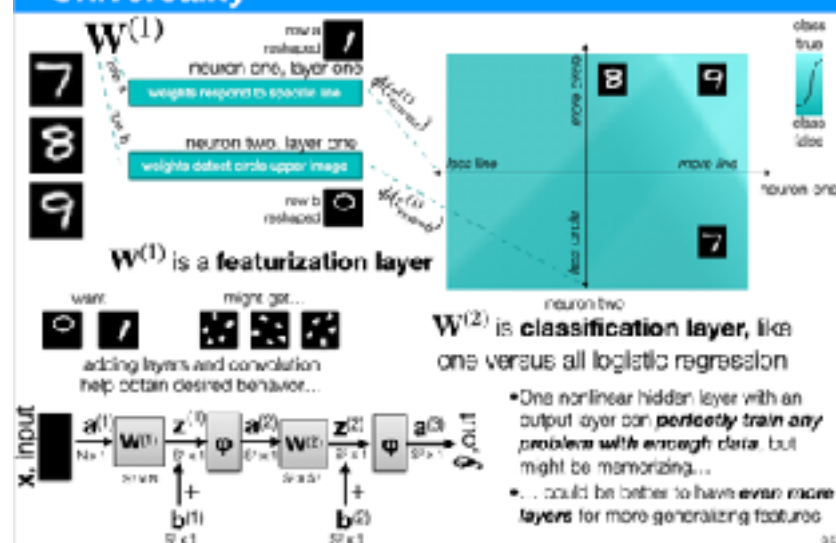


$$\eta_e = \eta_{\min} + \frac{1}{2}(\eta_{\max} - \eta_{\min}) \left(1 + \cos\left(\frac{e}{e_{\max}} \pi\right) \right)$$



36

Universality



McNemar Example

Model 1	Model 2	Label	Matrix
T-shirt	T-shirt	T-shirt	A
Sneaker	T-shirt	Sneaker	B
T-shirt	Pullover	Pullover	C
Sneaker	Sneaker	Sneaker	A
T-shirt	Sneaker	Sneaker	C
Pullover	Pullover	T-shirt	D
Pullover	T-shirt	Pullover	B
Sneaker	Sneaker	Sneaker	A
Sneaker	Sneaker	Sneaker	A

Model 1 correct	Model 2 wrong
4 ^A	2 ^B
2 ^C	1 ^D

McNemar and Edwards, 1946

$$\chi^2 \approx \frac{(|B - C| - 1)^2}{B + C}$$

$$\chi^2 = \frac{(|2 - 2| - 1)^2}{2 + 2} = 0.25$$

Confidence	0.90	0.95	0.99
1 DOF, Critical Value	2.706	3.841	6.635

<https://www.it-ebooks.info/book/view/5650440/4130>

Since $0.25 < 3.841$, we cannot reject the null hypothesis. This means **we should not say the models' performance are different** based on the evidence.

Self Test

- Modern day machine learning scientists have a rigorous background in neural biology or in processes that are biologically plausible.
 - **A. True.** Neural networks are windows into the mind.
 - **B. True.** PhD's in Machine Learning require advanced biology study.
 - **C. True.** They have made bold claims about human intelligence, and therefore must have this expertise.
 - **D. False.** They receive no formal training in this nor do they regularly consult with experts in this field.

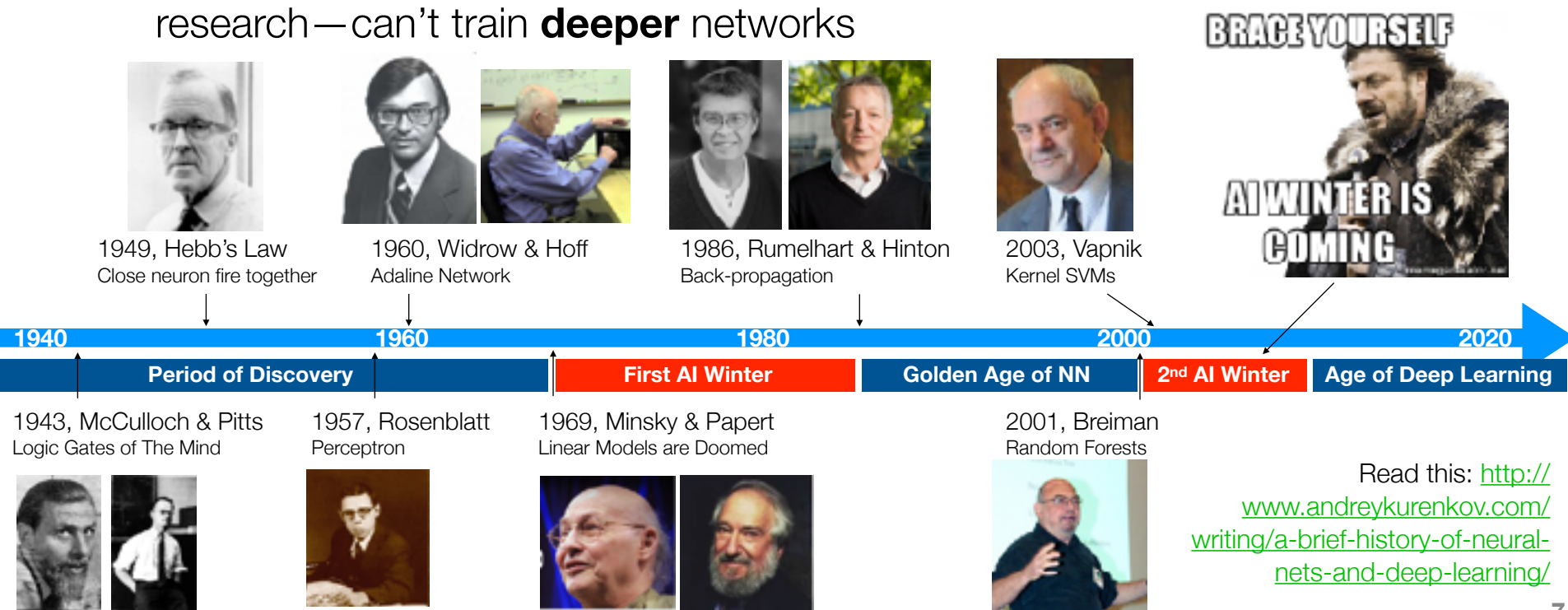
Some History of Deep Learning

When you move on to
Deep Learning



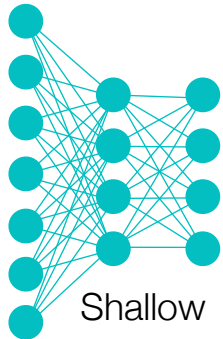
Machine Learning Timeline (Neural Nets)

- Up to this point: back propagation ended first AI winter
- 80's, 90's: neural networks for image processing start to get deeper
 - but back propagation no longer efficient for training
 - Back propagation gradient **stagnates** research—can't train **deeper** networks
- Second AI winter begins, research in NN plummets
- Funding for and accepted papers with Neural Networks asymptotically approaches zero

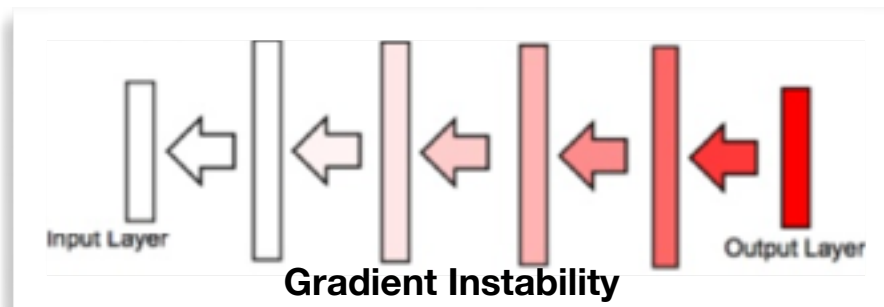
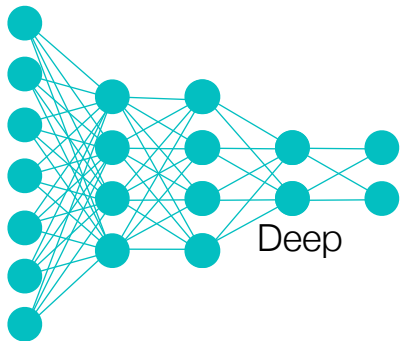


History of Deep Learning: Winter

- AI Winter is coming:



Easy to train, performs on par with other methods



Hard to train, performs worse than other methods
~chance (untrainable)

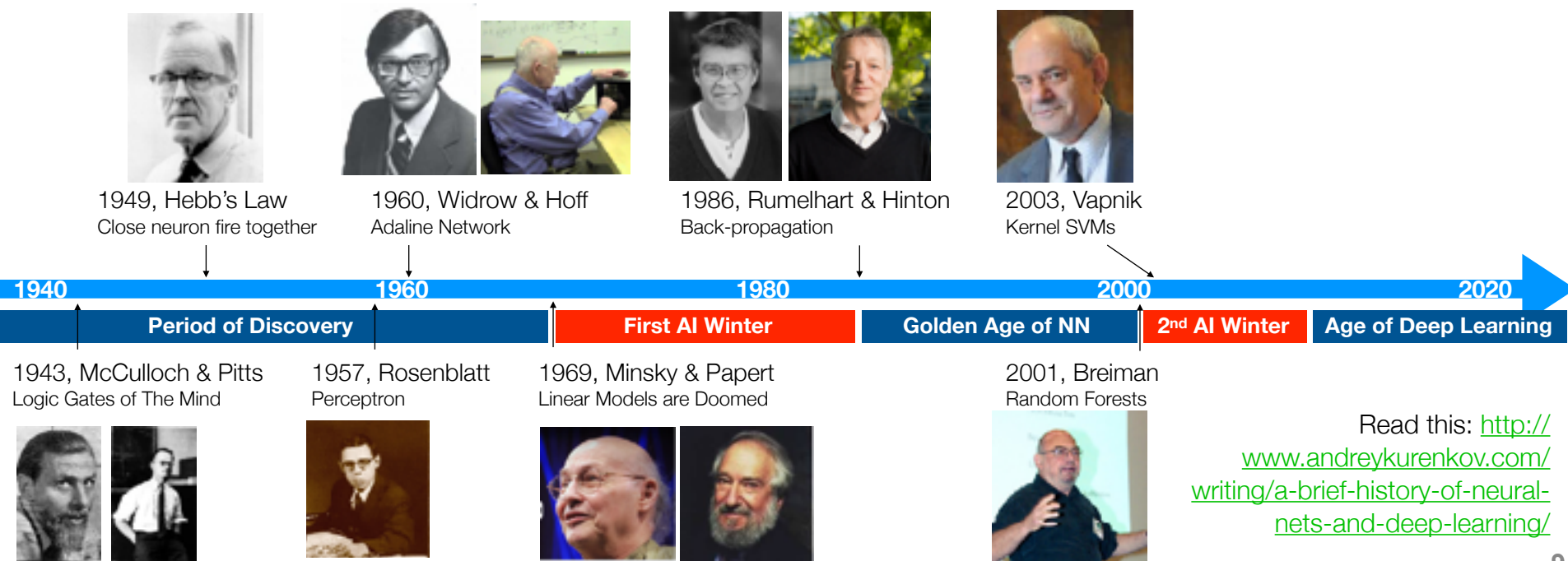
1990's: Many companies founded on the promise of advanced ML and the internet...

2000's: Researchers have difficulty reconciling expressiveness with performance, cash in the market drops out, **no theory and no money**



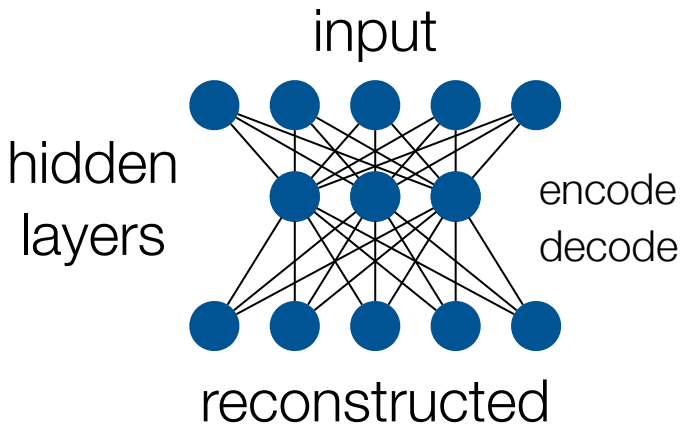
Machine Learning Timeline (Neural Nets)

- 2004: Hinton secures funding from CIFAR based on his reputation
 - *eventually*: Canada would be savior for neural networks
 - Hinton rebrands: **Deep Learning**
- 2006: Hinton publishes paper on using pre-training and Restricted Boltzmann Machines
- 2007: Another paper: Deep networks are more efficient when pre-trained
 - RBMs not really the important part

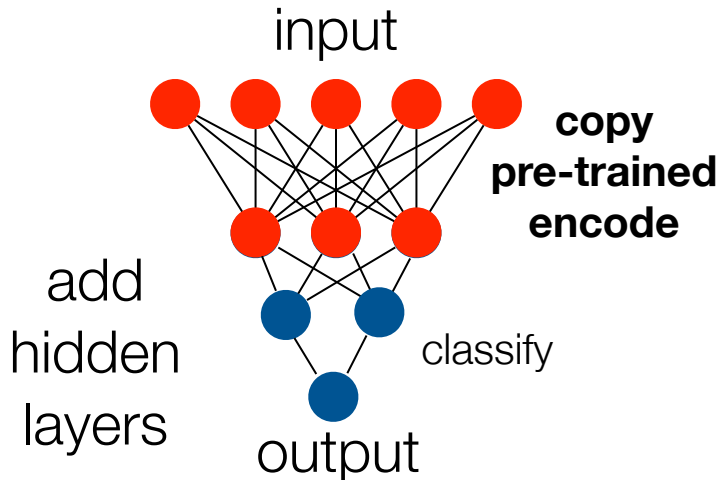
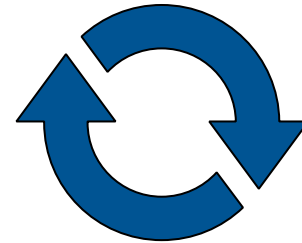


Pre-training: still in the long winter

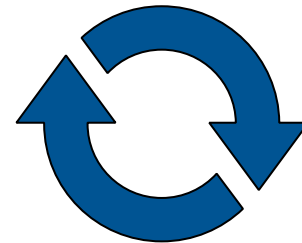
- auto-encoding (a form of pre-training)



train with lots of
unlabeled data

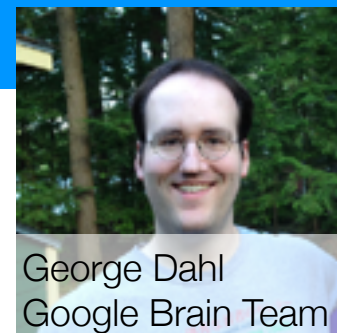
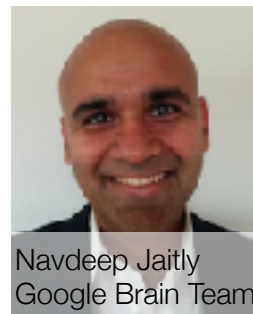


train with
labeled data



Still in the Long Winter

- 2009: Hinton's lab starts using GPUs, Also Andrew Ng
 - GPUs decrease training time by 70 fold...
- 2010: Hinton's and Ng's students go to internships with Microsoft, Google, IBM, and Facebook



Abdel-rahman Mohamed
Microsoft Research
Redmond, Washington | Computer Software
Current Microsoft
Previous University of Toronto, IBM, Microsoft
Education University of Toronto

- Xbox Voice
- Android Speech Recognition
- IBM Watson
- DeepFace
- All of Baidu



1949, Hebb's Law
Close neuron fire together



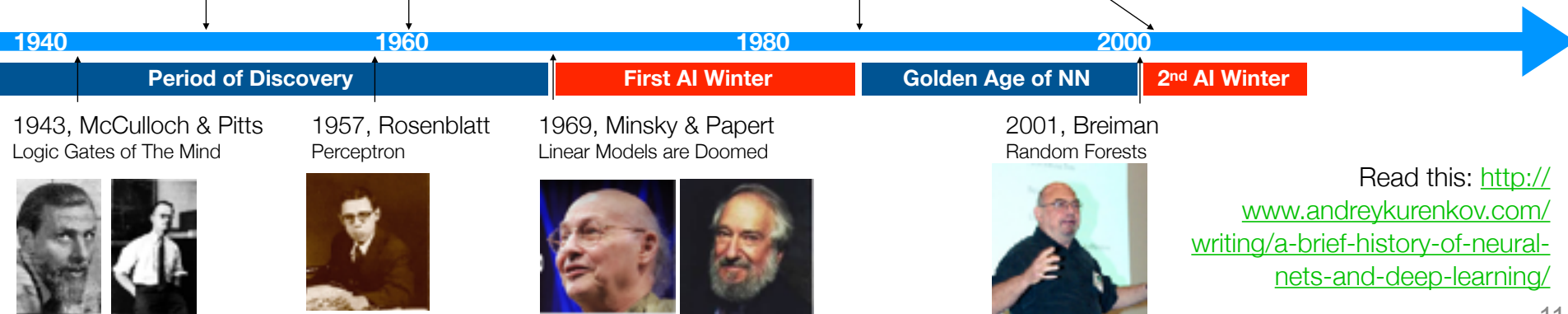
1960, Widrow & Hoff
Adaline Network



1986, Rumelhart & Hinton
Back-propagation



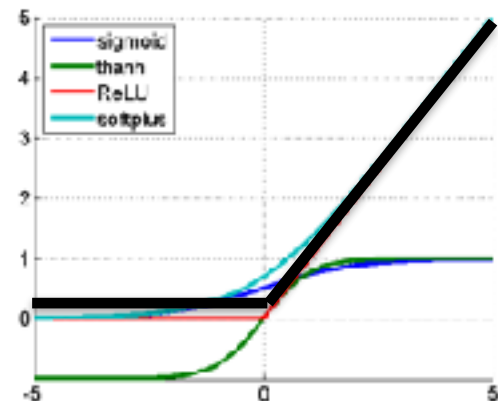
2003, Vapnik
Kernel SVMs



Read this: <http://www.andreykurenkov.com/writing/a-brief-history-of-neural-nets-and-deep-learning/>

Getting out of the long Winter

- 2011: Glorot and Bengio investigate more systematic methods for why past deep architectures did not work
 - **discover some interesting, simple fixes:** the type of neurons chosen and the selection of initial weights
 - do not require pre-training to get deep networks properly trained, just sparser representations and less complicated derivatives



$$\text{ReLU: } f(x) = \max(0, x) \\ f'(x) = 1 \text{ if } x > 0 \text{ else } 0$$



1949, Hebb's Law
Close neuron fire together



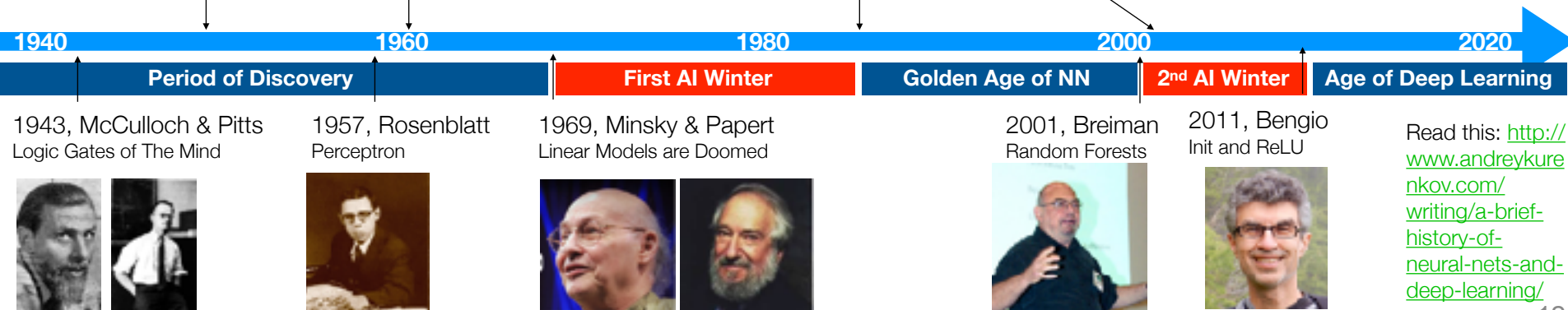
1960, Widrow & Hoff
Adaline Network



1986, Rumelhart & Hinton
Back-propagation



2003, Vapnik
Kernel SVMs



Machine Learning Timeline

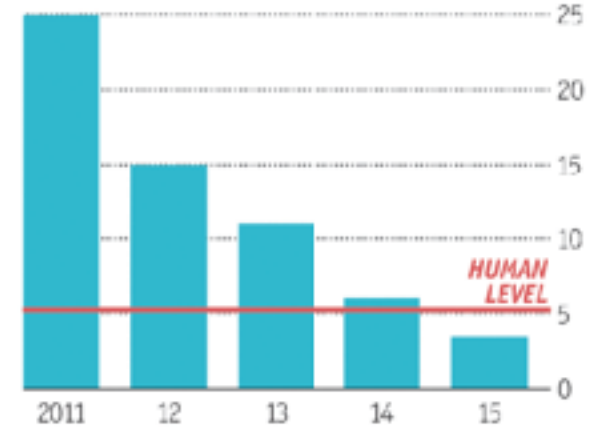
- **ImageNet competition occurs**
- **Second place:** 26.2% error rate
- **First place:**
 - From Hinton's lab, uses convolutional network with ReLU and dropout
 - 15.2% error rate
- Computer vision adopts deep learning with convolutional neural networks en masse



"I have had a hard time last fall so I really as she come pacin happ from skip a co Visio

Ever cleverer

Error rates on ImageNet Visual Recognition Challenge, %



Sources: ImageNet; Stanford Vision Lab

Economist.com



1949, Hebb's Law
Close neuron fire together



1960, Widrow & Hoff
Adaline Network



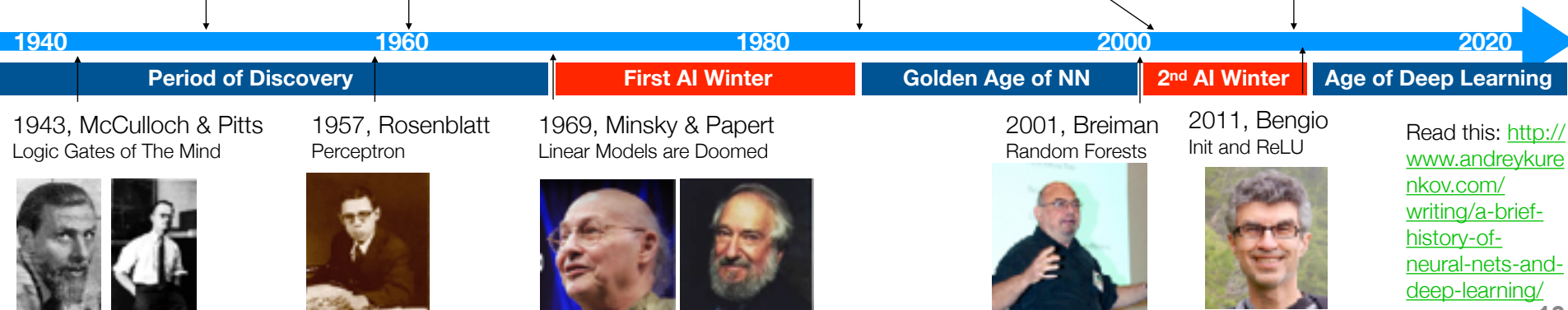
1986, Rumelhart & Hinton
Back-propagation



2003, Vapnik
Kernel SVMs



2012, Hinton, Fei-Fei Li
CNNs win ImageNet



Machine Learning Timeline (Acoustics)

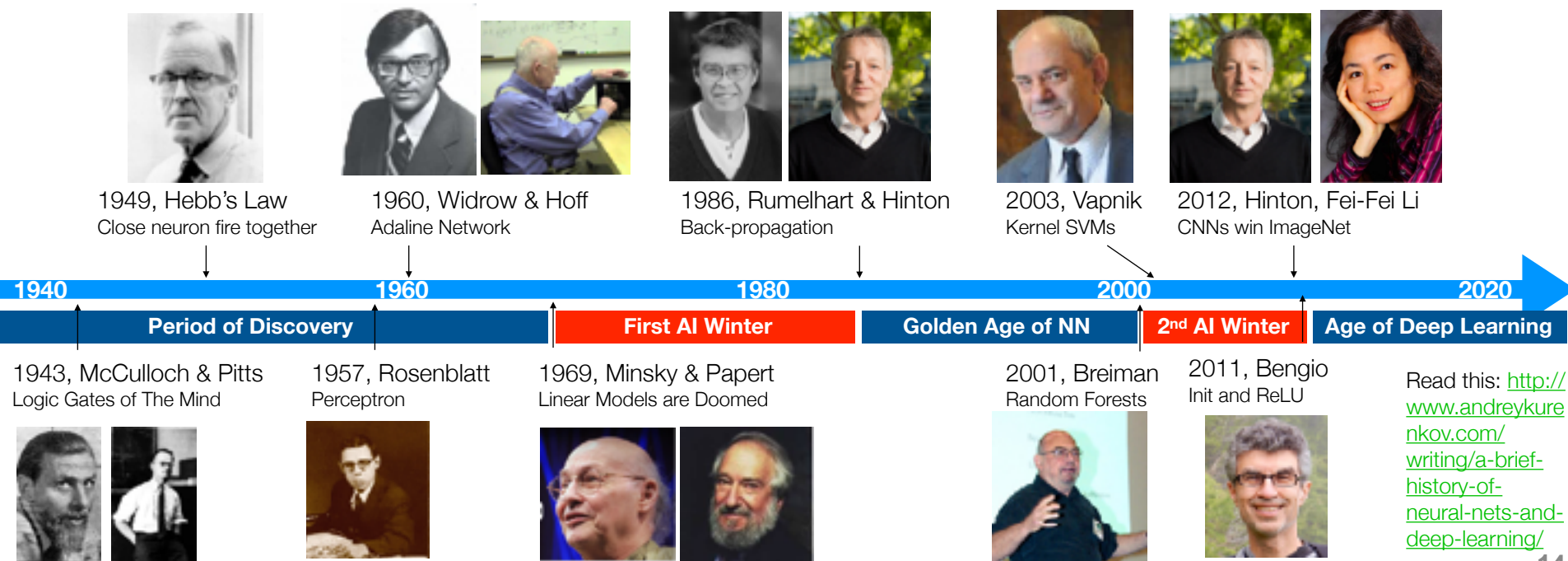
- 2012: Hinton Lab, Google, IBM, and Microsoft jointly publish paper, popularity for deep learning methods increases

Deep Neural Networks for Acoustic Modeling in Speech Recognition

[The shared views of four research groups]

[Geoffrey Hinton, Li Deng, Dong Yu, George E. Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N. Sainath, and Brian Kingsbury]

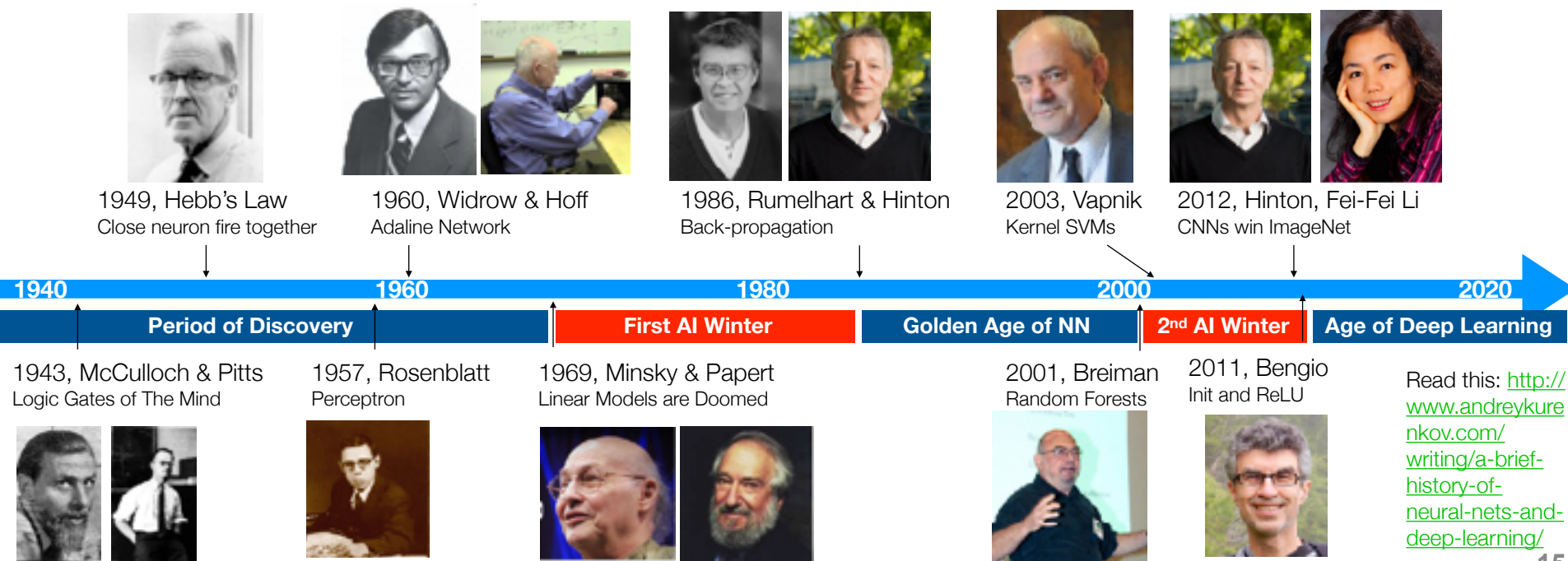
<https://www.cs.toronto.edu/~gdahl/papers/deepSpeechReviewSPM2012.pdf>









Machine Learning Timeline (Neural Nets)


- 2013: Andrew Ng and Google (BrainTeam)
 - run unsupervised feature creation on YouTube videos (becomes computer vision benchmark)

The work resulted in unsupervised neural net learning of an unprecedented scale - 16,000 CPU cores powering the learning of a whopping 1 billion weights. The neural net was trained on Youtube videos, entirely without labels, and learned to recognize the most common objects in those videos.



A summary of the Deep Learning people:

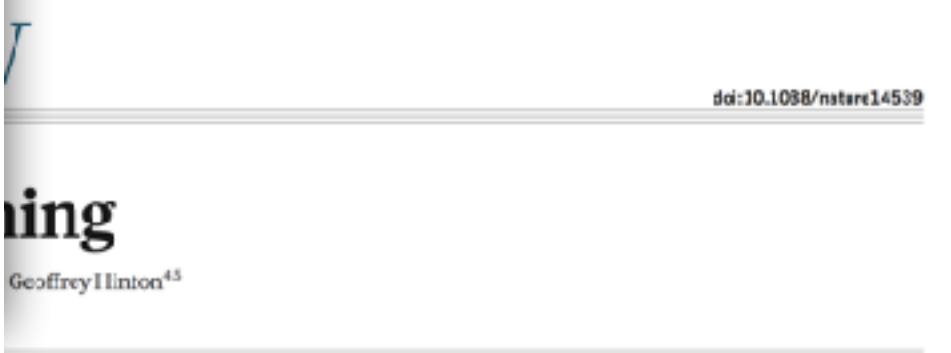
					
Yoshua Bengio	Yann LeCun	Geoffrey Hinton	FeiFei Li	Andrew Ng	Daphne Koller
Stayed at Univ. Montreal Advises IBM	Heads Facebook AI Team	Univ. Toronto Google	Stanford (HAI) Former Chief Scien., AI/ML Google Cloud	Coursera Baidu Google	Stanford Founded Coursera MacArthur Genius



- Hinton: Restricted Boltzmann Machine, Deep autoencoder
- Bengio: neural language modeling.
- LeCun: Convolutional Neural Network
- NIPS, ICML, CVPR, ACL
- Google Brain, Deep Mind.
- FaceBook AI.



Made Deep Learning Instruction Accessible



deep learning

Deep learning allows computational models that are composed of multiple processing layers to learn representations of data with multiple levels of abstraction. These methods have dramatically improved the state-of-the-art in speech recognition, visual object recognition, object detection and many other domains such as drug discovery and genomics. Deep learning discovers intricate structure in large data sets by using the backpropagation algorithm to indicate how a machine should change its internal parameters that are used to compute the representation in each layer from the representation in the previous layer. Deep convolutional nets have brought about breakthroughs in processing images, video, speech and

Credit for Deep Learning

Official ACM @TheOfficialACM

Yoshua Bengio, Geoffrey Hinton and Yann LeCun, the fathers of #DeepLearning, receive the 2018 #ACMTuringAward for conceptual and engineering breakthroughs that have made deep neural networks a critical component of computing today. bit.ly/2HVJtdV



Yoshua Bengio

Geoffrey Hinton

Yann LeCun



Machine learning is the science of credit assignment. The machine learning community itself profits from proper credit assignment to its members. The inventor of an important method should get credit for inventing it. She may not always be the one who popularizes it. Then the popularizer should get credit for popularizing it (but not for inventing it). Relatively young research areas such

Review of Deep Learning History

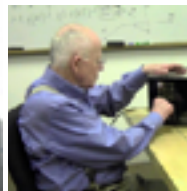
- Up to this point: back propagation saved AI winter for NN (Hinton and others!)
- 80's, 90's, 2000's: convolutional networks for image processing start to get deeper
 - but back propagation no longer does great job at training them
- SVMs and Random Forests gain traction...
 - The second AI winter begins, research in NN plummets
- 2004: Hinton secures funding from CIFAR in 2004 Hinton rebrands: Deep Learning
- 2006: Auto-encoding and Restricted Boltzmann Machines
- 2007: Deep networks are more efficient when pre-trained
- 2009: GPUs decrease training time by 70 fold...
- 2010: Hinton's students go to internships with Microsoft, Google, and IBM, making their speech recognition systems faster, more accurate and deployed in only 3 months...
- 2012: Hinton Lab, Google, IBM, and Microsoft jointly publish paper, popularity sky-rockets for deep learning methods
- 2011-2013: Ng and Google run unsupervised feature creation on YouTube videos (becomes computer vision benchmark)
- 2012+: Pre-training is not actually needed, just solutions for vanishing gradients (like ReLU, SiLU, initializations, more data, GPUs)



1949, Hebb's Law
Close neuron fire together



1960, Widrow & Hoff
Adaline Network



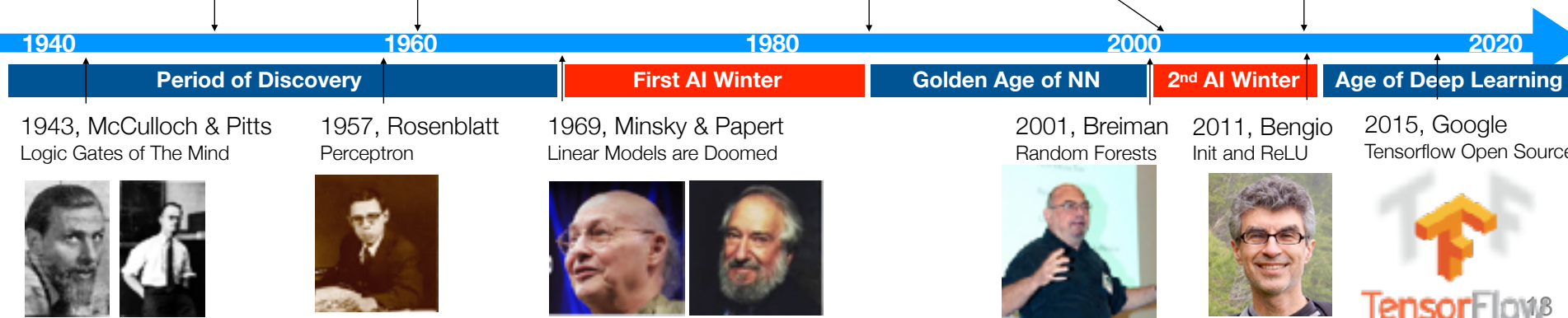
1986, Rumelhart & Hinton
Back-propagation



2003, Vapnik
Kernel SVMs



2012, Hinton, Fei-Fei Li
CNNs win ImageNet



The New Paradigm



Tensor Packages for Deep Learning

"Further discussion of it merely incumbers the literature and befogs the mind of fellow students."

- 2007: **NIPS** program committee rejects a paper on deep learning by *al. et.* Hinton because they already accepted a paper on deep learning and two papers on the same topic would be excessive.
- ~2009: A reviewer tells Yoshua Bengio that papers about neural nets have no place in **ICML**.
- ~2010: A **CVPR** reviewer rejects Yann LeCun's paper even though it beats the state-of-the-art. The reviewer says that it tells us nothing about computer vision because everything is learned.

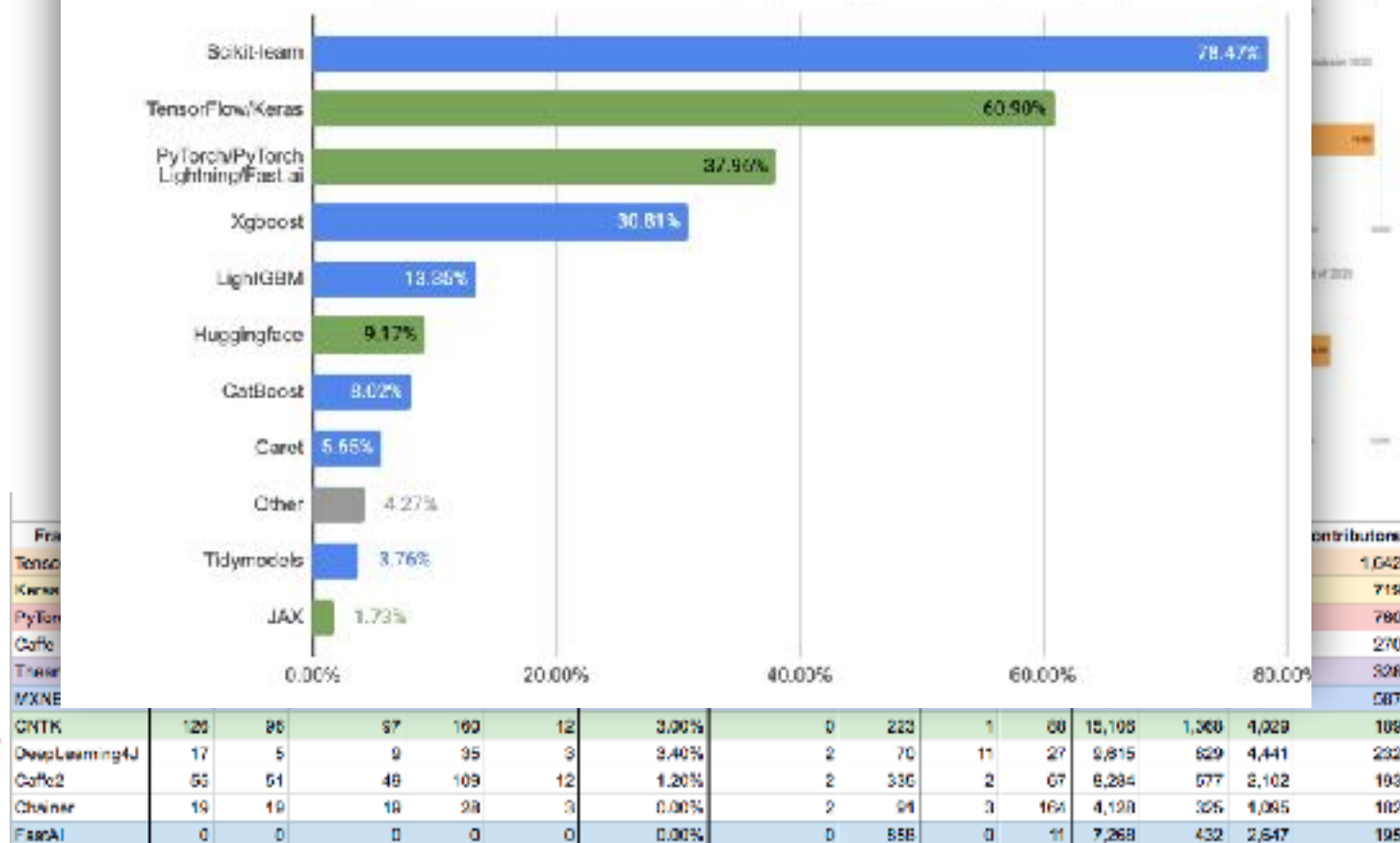


Options for Deep Learning Toolkits

1.  TensorFlow
2.  Keras
3.  PyTorch
4.  Caffe
5.  theano
6.  Apache MXNet
7.  Microsoft CNTK
8.  DL4J
9.  Caffe2
10.  Chainer
11.  fast.ai

Overview of Deep Learning frameworks adoption metrics over 2020

2022 Machine Learning & Data Science Survey by Kaggle: library usage (N=14,531)



Tensorflow

- Open sourced library from Google
- Second generation release from Google Brain
 - supported for Linux, Unix, Windows
 - Also works on Android/iOS
- Released November 9th, 2015
 - (this class first offered January 2016)
- PyTorch released January 2017



Programmatic creation

- Most toolkits use python to build a **computation graph** of operations
 - Build up computations
 - Execute computations
- **Most Toolkits Support:**
 - tensor creation
 - functions on tensors
 - automatic differentiation
- Tensors are just multidimensional arrays
 - like in Numpy
 - scalars (biases and constants)
 - vectors (e.g., input arrays)
 - 2D matrices (e.g., images)
 - 3D matrices (e.g., color images)
 - 4D matrices (e.g., batches of color images)

Tensor basic functions

- Easy to define operations on tensors

```
a = tf.constant(5.0)
```

```
b = tf.constant(6.0)
```

```
c = a * b
```

Numpy	TensorFlow
<code>a = np.zeros((2,2)); b = np.ones((2,2))</code>	<code>a = tf.zeros((2,2)), b = tf.ones((2,2))</code>
<code>np.sum(b, axis=1)</code>	<code>tf.reduce_sum(a, reduction_indices=[1])</code>
<code>a.shape</code>	<code>a.get_shape()</code>
<code>np.reshape(a, (1,4))</code>	<code>tf.reshape(a, (1,4))</code>
<code>b * 5 + 1</code>	<code>b * 5 + 1</code>
<code>np.dot(a,b)</code>	<code>tf.matmul(a, b)</code>
<code>a[0,0], a[:,0], a[0,:]</code>	<code>a[0,0], a[:,0], a[0,:]</code>

Also supports convolution: `tf.nn.conv2d`, `tf.nn.conv3D`

Tensor neural network functions

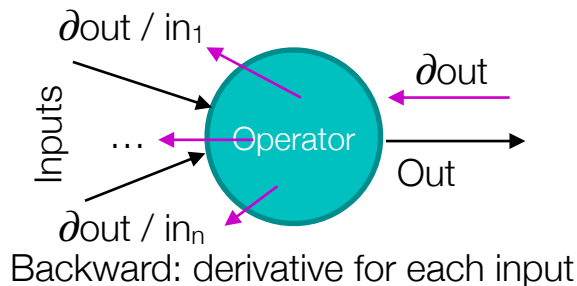
- Easy to define operations on layers of networks
 - `relu(features, name=None)`
 - `bias_add(value, bias, data_format=None, name=None)`
 - `sigmoid(x, name=None)`
 - `tanh(x, name=None)`
 - `conv2d(input, filter, strides, padding)`
 - `conv1d(value, filters, stride, padding)`
 - `conv3d(input, filter, strides, padding)`
 - `conv3d_transpose(value, filter, output_shape, strides)`
 - `sigmoid_cross_entropy_with_logits(logits, targets)`
 - `softmax(logits, dim=-1)`
 - `log_softmax(logits, dim=-1)`
 - `softmax_cross_entropy_with_logits(logits, labels, dim=-1)`
- Each function created *knows its gradient*
- **Automatic Differentiation** is just **chain rule**
- But... lets start simple...

Auto Differentiation (under the hood)

```
class add(Operator):
    """Binary addition operation."""
    def __init__(self, a, b):
        super().__init__()
        self.inputs=[a, b]

    def forward(self, a, b):
        return a+b

    def backward(self, a, b, dout):
        return dout, dout
```



```
class matmul(Operator):
    """Binary multiplication"""
    def __init__(self, a, b):
        super().__init__()
        self.inputs=[a, b]

    def forward(self, a, b):
        return a@b

    def backward(self, a, b, dout):
        return dout@b.T, a.T@dout
```

```
def forward_pass(order, feed_dict={}):
    """ Performs the forward pass, returning output of graph.
    Args:
        order: a topologically sorted array of nodes
        feed_dict: a dictionary values for placeholders.
    Returns:
        1. the final result of the forward pass.
        2. edits the graph to fill in its values.
    """
    for node in order:
        tmp = [prev_node.value for prev_node in node.inputs]
        node.value = node.forward(*tmp)

    return order[-1].value
```

Could also have GPU implementations

```
def backward_pass(order):
    """ Perform the backward pass to retrieve gradients.
    gradients of nodes as listed in same order as input """
    vis = set()
    order[-1].gradient = 1
    for node in reversed(order):
        inputs = node.inputs
        grads = node.backward(*[x.value for x in inputs],
                               dout=node.gradient)
        for inp, grad in zip(inputs, grads):
            if inp not in vis:
                inp.gradient = grad
            else:
                inp.gradient += grad
            vis.add(inp)
    return [node.gradient for node in order]
```

Tensor function evaluation

```
import tensorflow as tf

a = tf.constant(5.0)
b = tf.constant(6.0)

c = a*b

// pre tensorflow 2
with tf.Session() as sess:
    print(sess.run(c))
    print(c.eval())

// post tensorflow 2
print(c)
```

output = 30

- Easy to define operations on tensors
 - constant
 - variables
 - placeholders
- Nothing evaluated until the output is needed
- Tensorflow now run in default eager execution

Computation Graph with Code

```
import tensorflow as tf
X = tf.placeholder()
y = tf.placeholder()
```

$$J(\mathbf{W}) = \frac{1}{N} \sum_i (y^{(i)} - \underbrace{(\mathbf{W} \cdot \mathbf{x}^{(i)} + \mathbf{b})}_{y_{pred}})^2$$

1. **Setup** Variables and computations

```
W = tf.Variable("weights", (1,num_features),
               initializer=tf.random_normal_initializer())
b = tf.Variable("bias", (1,),
               initializer=tf.constant_initializer(0.0))
```

```
def feedforward(X,y):
    y_pred = tf.matmul(X,W) + b
    loss = tf.reduce_sum((y-y_pred)**2)
    return loss
```

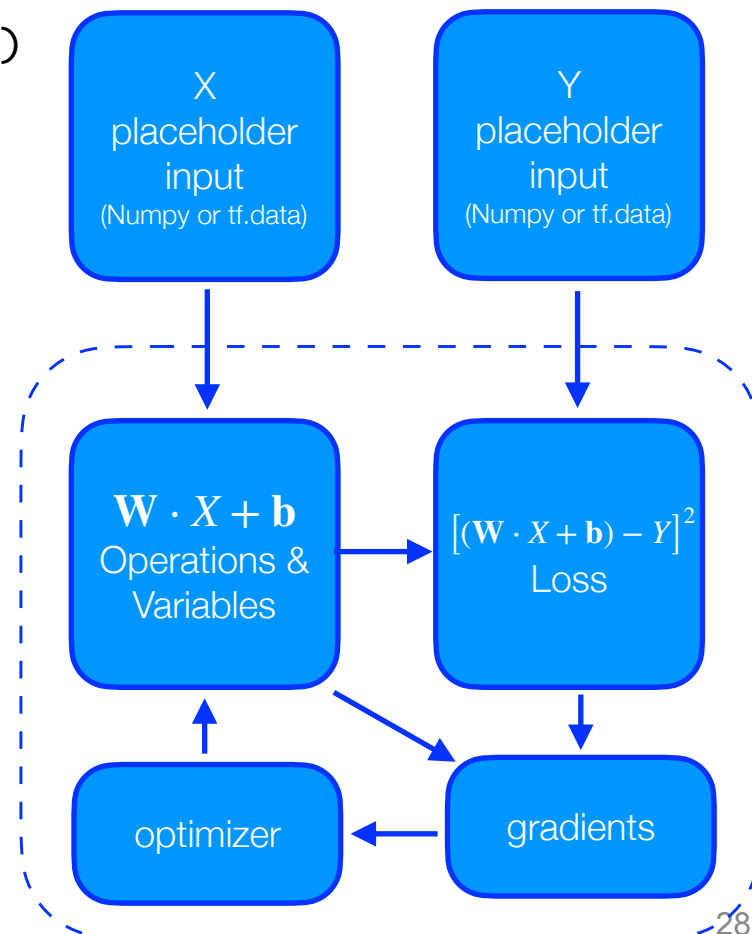
2. Add **optimization** operation to computation graph
Adjusts variables (W, b) to minimize loss with auto differentiation

```
opt = tf.train.AdamOptimizer()

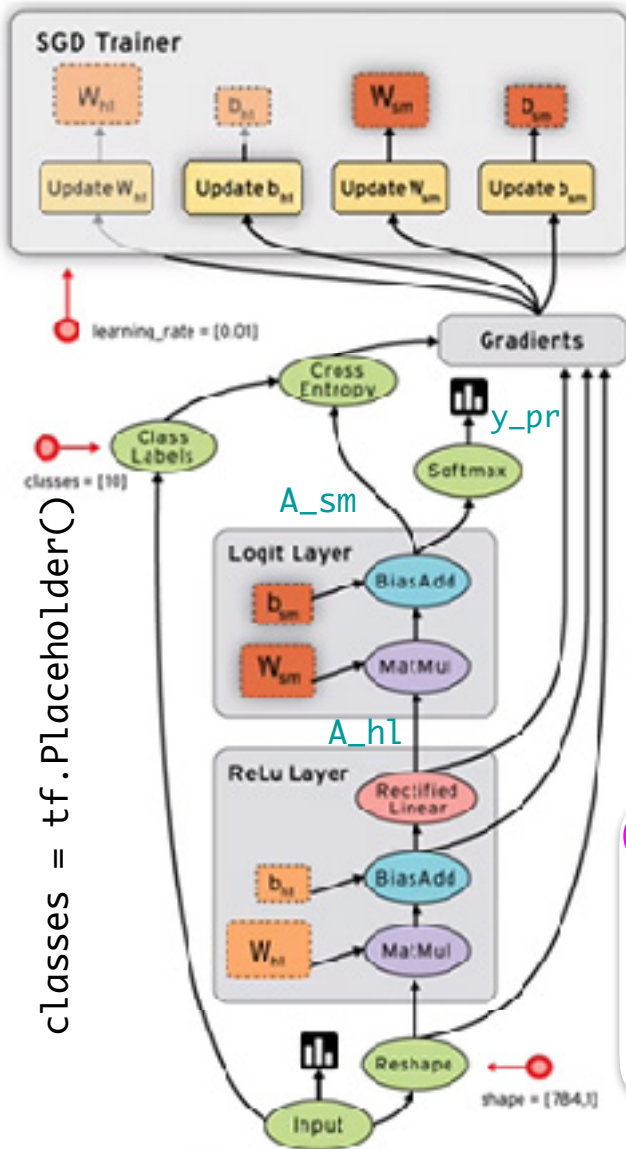
tf.initialize_all_variables()
feedforward(X_numpy, y_numpy)
... track gradients and variables ...
opt.apply_gradients(zip(grads, train_vars))
```

3. **Run graph operation** once, → one optimization update
on all variables

<http://www.datasciencecentral.com/profiles/blogs/google-open-source-tensorflow>



Computation Graph, Two Layer Network



```
Input = tf.placeholder() # size is 28x28
Input = tf.reshape(Input, [784,1])
classes = tf.placeholder()
```

```
W_sm = tf.Variable(...)
b_sm = tf.Variable(...)
W_hl = tf.Variable(...)
b_hl = tf.Variable(...)
```

```
trainable_variables =
    [W_sm, b_sm, W_hl, b_hl]
```

```
def model_forward(Input):
    A_hl = tf.relu( tf.matmul(Input, W_hl) + b_hl )
    A_sm = tf.matmul(A_hl, W_sm) + b_sm
    return A_sm
```

```
y_pr = tf.softmax(A_sm)
loss = tf.sparse_softmax_cross_entropy_with_logits
opt = tf.train.SGDOptimizer(learning_rate=0.01)
```

```
for features, labels in train_data:
    with tf.GradientTape( ) as tape:
        yhat = model_forward(features)
        loss_val = loss(labels, yhat)
        grads = tape.gradient(loss_val, trainable_variables)
        opt.apply_gradients(zip(grads, trainable_variables))
```

Tensorflow Simplification

- **Self Test:** Can the syntax be simplified?
 - (A) **Yes**, we could write a generic mini-batch optimization computation graph, then use it for arbitrary graph instructions
 - (B) **Yes**, but we need to learn the Keras API, which can be mixed with tensorflow operations
 - (C) **Yes**, but we need to understand how to access the gradients to apply them, a lot like PyTorch
 - (D) **All of the above**

End of Session

- Next Time:
 - Introduction to TensorFlow
 - Wide and Deep Networks