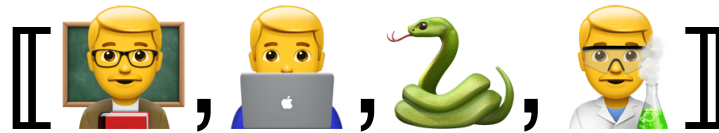


Lecture Notes for **Machine Learning in Python**



A “Not so Early” History of Deep Learning

Logistics and Agenda

- Logistics
 - Grading update
- Agenda
 - Town Hall
 - “Deep Learning” History

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

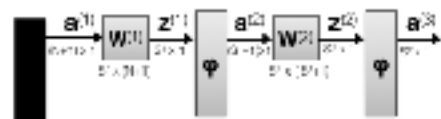
Keras, Tensorflow
Intuition, Detailed implement.

Ethics in
Language Models

ConceptNet
Case studies

Last time:

Back propagation summary



$$J(W) = \sum_k (y^{(k)} - a^{(L)})^2$$

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta \frac{\partial J(W)}{\partial z^{(l)}} a_j^{(l)}$$

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 $\mathbf{W}^{(l)}$

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

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

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

Practical Implementation of Architectures

- A new cost function: **Cross entropy**

$$J(W) = -[y^{(l)} \ln a^{(l)} + (1 - y^{(l)}) \ln(1 - a^{(l)})]$$

speeds up initial training

$$\frac{\partial J(W)}{\partial z^{(L)}} = (a^{(L+1)} - y^{(l)})$$

vectorized backpropagation
sigma1 = (A1 - Y_end) # <- this is only line
sigma2 = (W2.T @ sigma1) * A2 * (1 - A2)

$$\frac{\partial J(W)}{\partial z^{(2)}} = (a^{(3)} - y^{(l)})$$

grad1 = sigma2[1:, :] @ A1
grad2 = sigma3 @ A2.T

new update

$$\frac{\partial J(W)}{\partial z^{(l)}} = -2(y^{(l)} - a^{(l)}) * a^{(l)} * (1 - a^{(l)})$$

vectorized backpropagation
sigma1 = -2 * (Y_end - A1) * A2 * (1 - A1)
sigma2 = (W2.T @ sigma1) * A2 * (1 - A2)

$$\frac{\partial J(W)}{\partial z^{(l)}} = -2(y^{(l)} - a^{(l)}) * a^{(l)} * (1 - a^{(l)})$$

grad1 = sigma2[1:, :] @ A1
grad2 = sigma3 @ A2.T

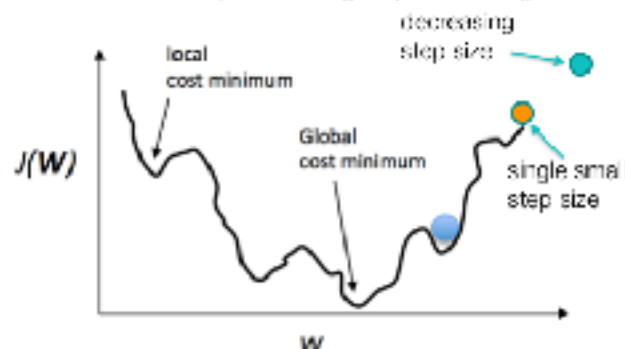
$$\frac{\partial J(W)}{\partial z^{(2)}} = -2(y^{(2)} - a^{(2)}) * a^{(1)} * (1 - a^{(1)})$$

old update

bp-5

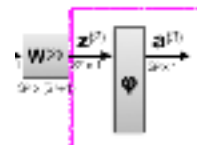
Problems with Advanced Architectures

- Space is no longer convex
 - One solution:
 - start with large step size
 - "cool down" by decreasing step size for higher iterations



Practical Implementation of Architectures

- A new nonlinearity: **rectified linear units**



$$\phi(z^{(l)}) = \begin{cases} z^{(l)}, & \text{if } z^{(l)} > 0 \\ 0, & \text{else} \end{cases}$$

it has the advantage of **large gradients** and **extremely simple** derivative

$$\frac{\partial \phi(z^{(l)})}{\partial z^{(l)}} = \begin{cases} 1, & \text{if } z^{(l)} > 0 \\ 0, & \text{else} \end{cases}$$

79

Neural Networks and Deep Learning, Michael Nielsen, 2015

Town Hall



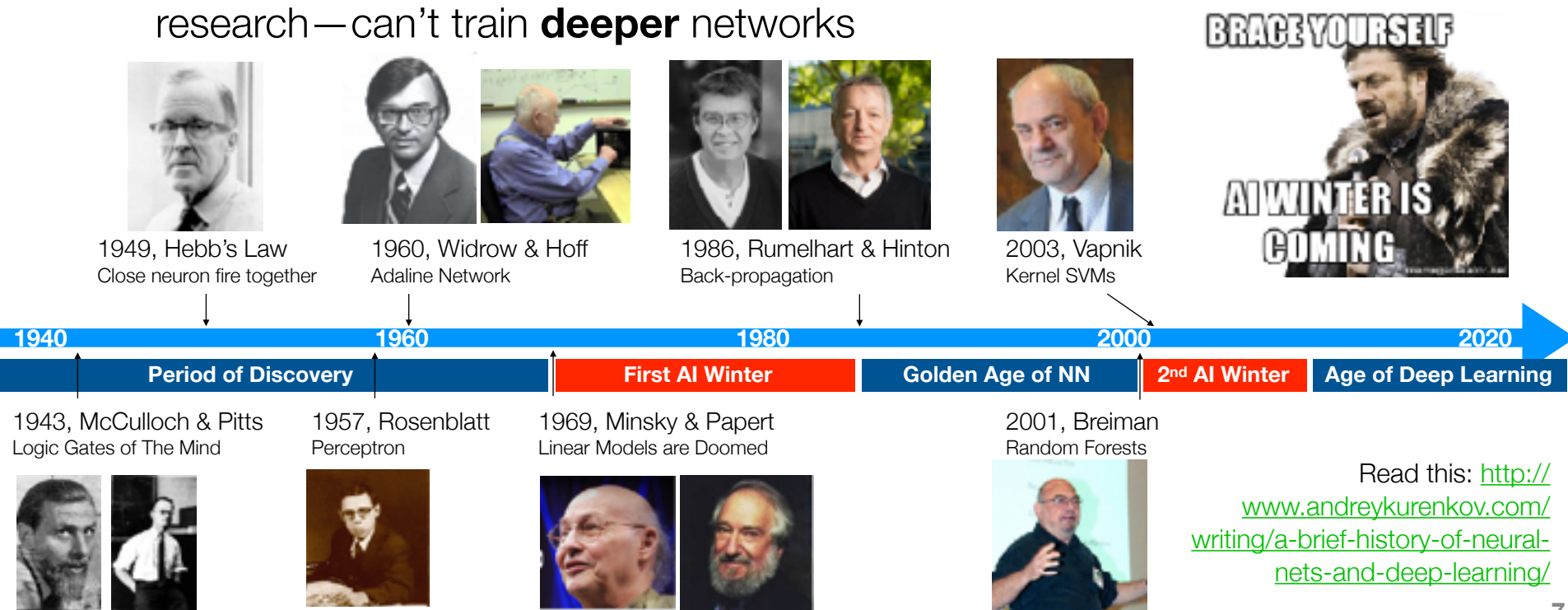
Some History of Deep Learning

When you move on to
Deep Learning



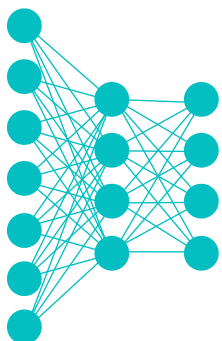
Machine Learning Timeline (Neural Nets)

- Up to this point: back propagation saved AI winter
- 80's, 90's, 2000'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

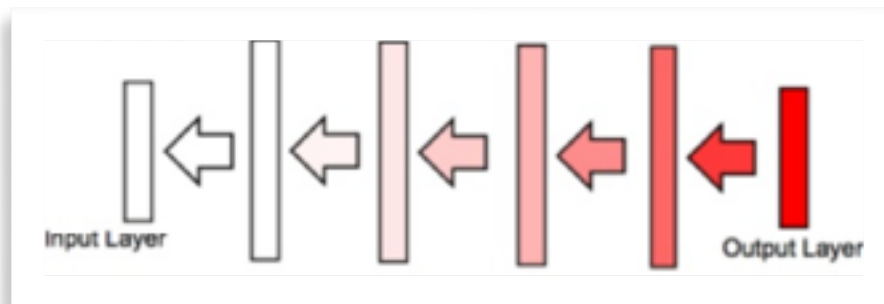


History of Deep Learning: Winter

- AI Winter is coming:



Easy to train, performs on par with other methods



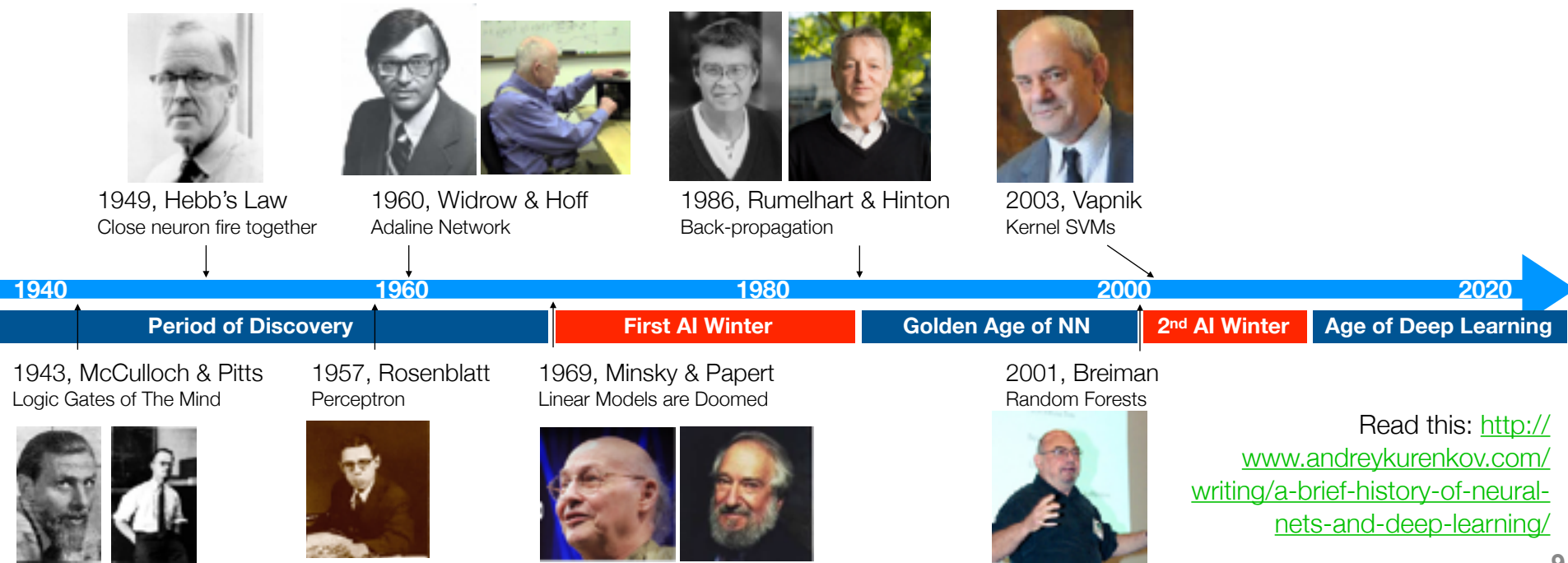
Hard to train, performs worse than other methods
~chance (untrainable)

Researcher have difficulty reconciling expressiveness with performance



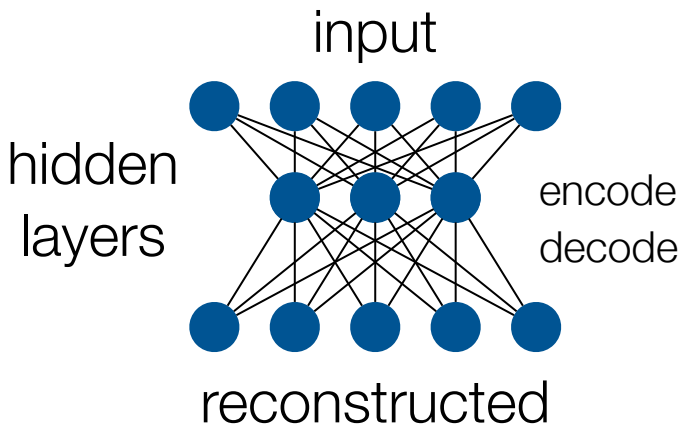
Machine Learning Timeline (Neural Nets)

- 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

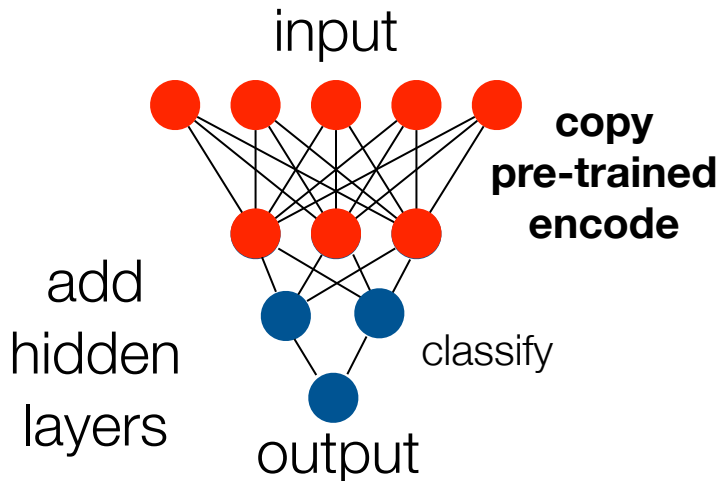
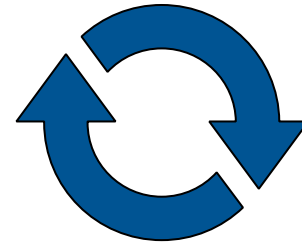


Pre-training: still in the long winter

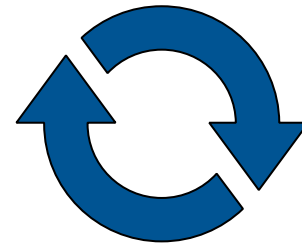
- auto-encoding (a form of pre-training)



train with lots of
unlabeled data



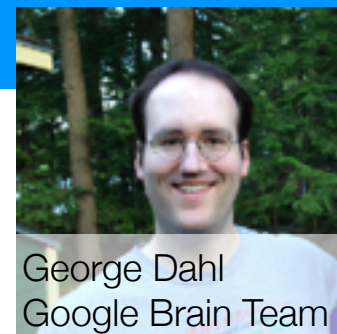
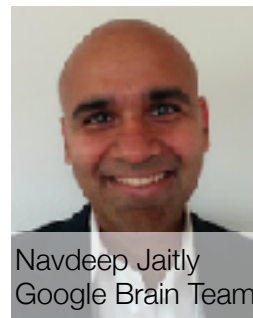
train with
labeled data



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

 Research at Google



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



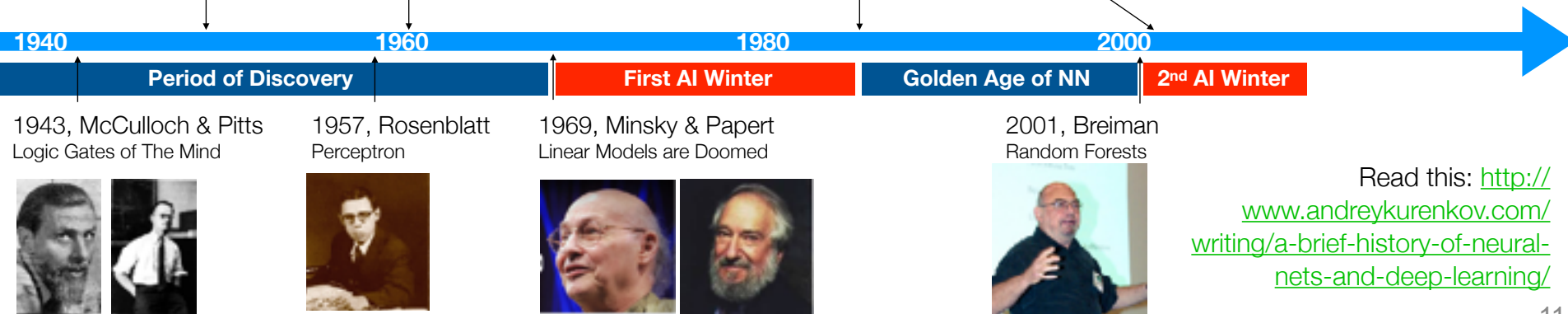
1960, Widrow & Hoff
Adaline Network



1986, Rumelhart & Hinton
Back-propagation



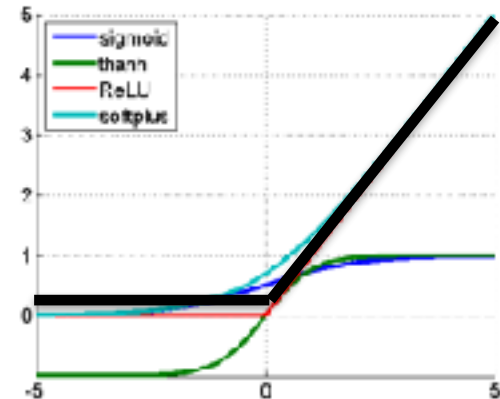
2003, Vapnik
Kernel SVMs



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



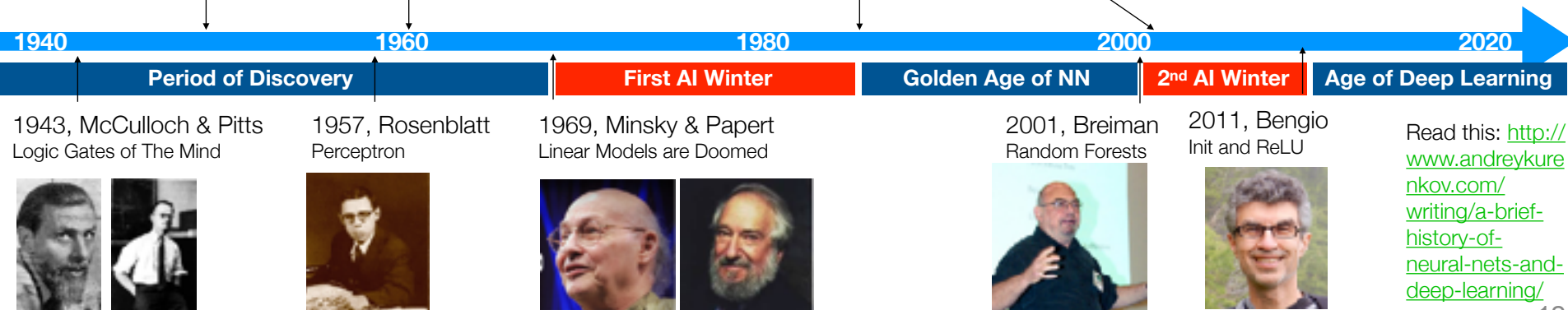
1960, Widrow & Hoff
Adaline Network



1986, Rumelhart & Hinton
Back-propagation



2003, Vapnik
Kernel SVMs



Machine Learning Timeline

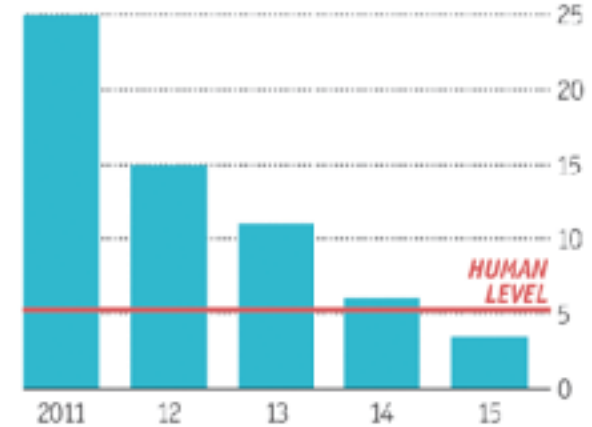
- **ImageNet competition occurs**
- **Second place:** 26.2% error rate
- **First place:**
 - From Hinton's lab, uses convolutional network with ReLU and dropout
 - 15.2% error rate
- Computer vision adopts deep learning with convolutional neural networks en masse



"I have had a hard time last fall so I need to do as much as I can to come back to the Pacific. I have been happy from the time I skip a class to a class in Vision."

Ever cleverer

Error rates on ImageNet Visual Recognition Challenge, %



Sources: ImageNet; Stanford Vision Lab

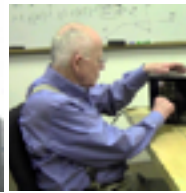
Economist.com



1949, Hebb's Law
Close neuron fire together



1960, Widrow & Hoff
Adaline Network



1957, Rosenblatt
Perceptron



1986, Rumelhart & Hinton
Back-propagation



2003, Vapnik
Kernel SVMs



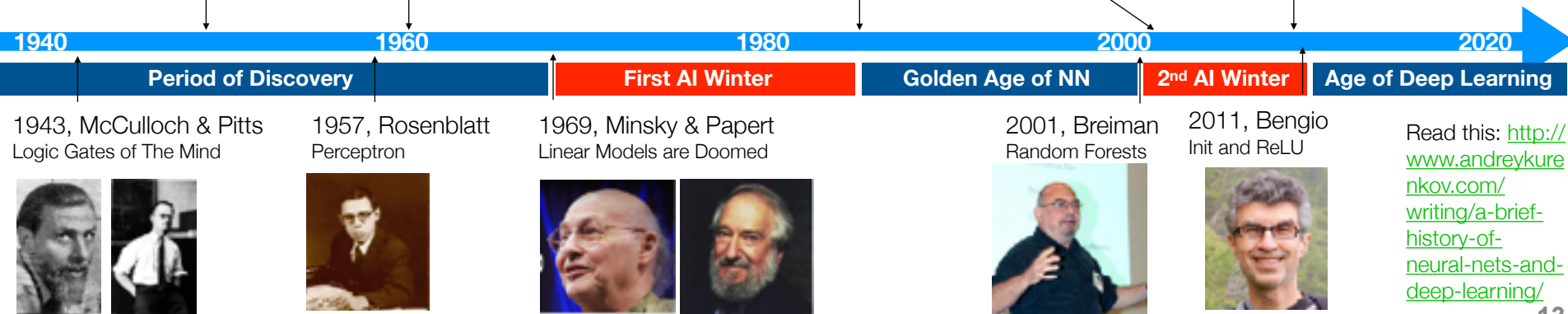
2012, Hinton, Fei-Fei Li
CNNs win ImageNet



2011, Bengio
Init and ReLU



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



Machine Learning Timeline (Neural Nets)

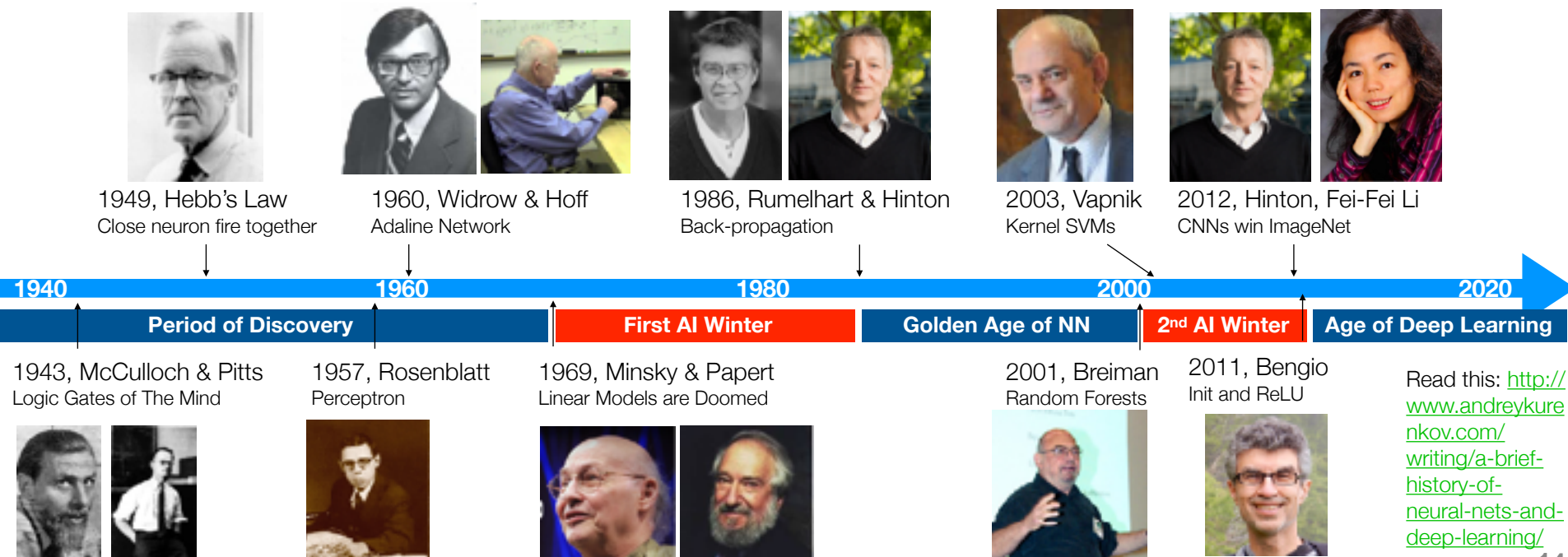
- 2012: Hinton Lab, Google, IBM, and Microsoft jointly publish paper, popularity for deep learning methods increases

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]

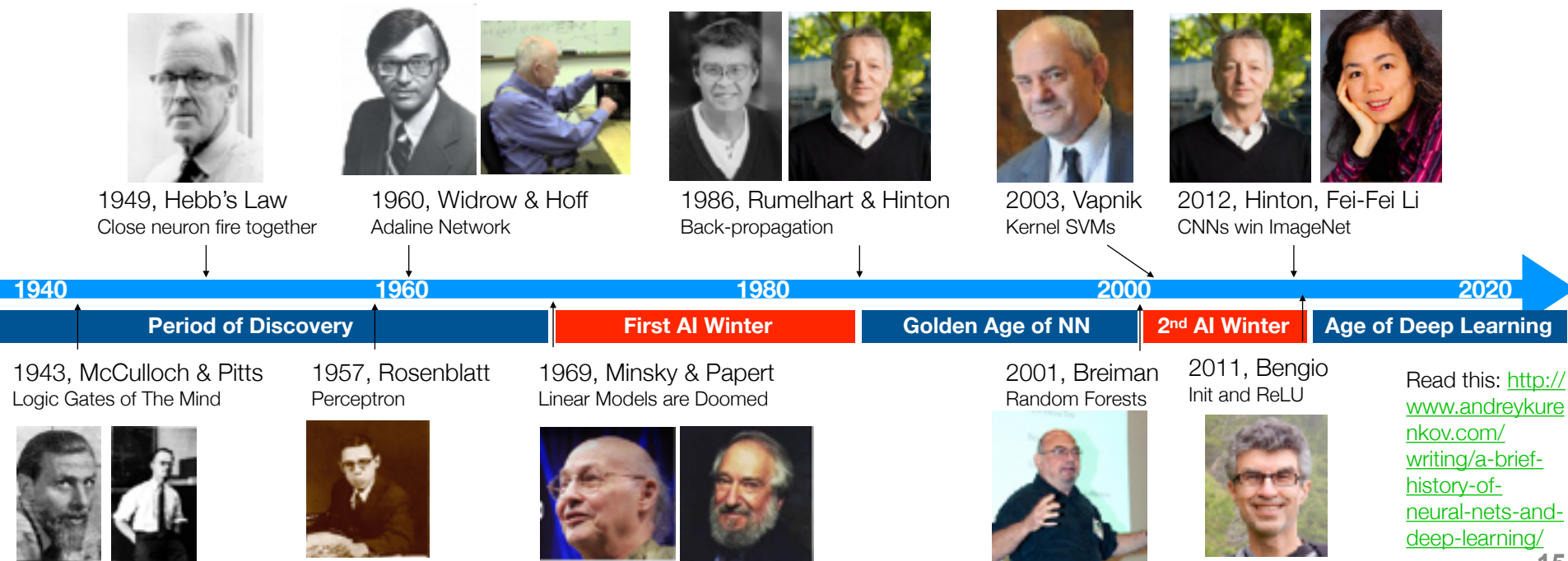
<https://www.cs.toronto.edu/~gdahl/papers/deepSpeechReviewSPM2012.pdf>









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

doi:10.1088/nature14539

ing

Geoffrey Hinton⁴⁵

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



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

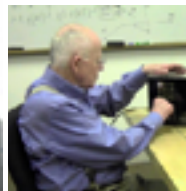
- Up to this point: back propagation saved AI winter for NN (Hinton and others!)
- 80's, 90's, 2000's: convolutional networks for image processing start to get deeper
 - but back propagation no longer does great job at training them
- SVMs and Random Forests gain traction...
 - The second AI winter begins, 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



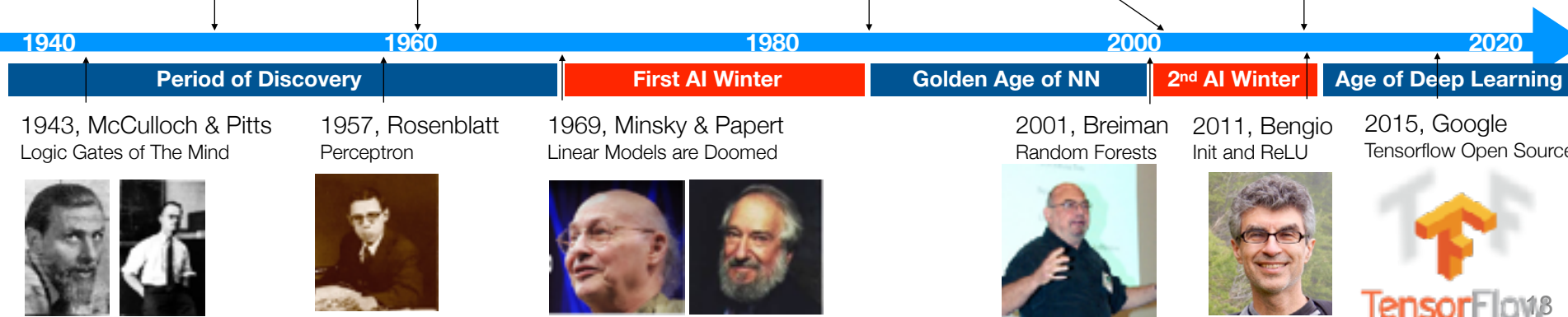
1986, Rumelhart & Hinton
Back-propagation



2003, Vapnik
Kernel SVMs



2012, Hinton, Fei-Fei Li
CNNs win ImageNet



End of Session

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