



## Bonus Lecture: Deep Learning



### A brief history

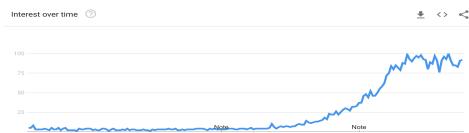
- 1943: neural networks  $\Leftrightarrow$  logical circuits (McCulloch/Pitts)
- 1949: "cells that fire together wire together" learning rule (Hebb)
- 1969: theoretical limitations of neural networks (Minsky/Papert)
- 1974: backpropagation for training multi-layer networks (Werbos)
- 1986: popularization of backpropagation (Rumelhardt, Hinton, Williams)

### A brief history

- 1980: Neocognitron, a.k.a. convolutional neural networks (Fukushima)
- 1989: backpropagation on convolutional neural networks (LeCun)
- 1990: recurrent neural networks (Elman)
- 1997: Long Short-Term Memory networks (Hochreiter/Schmidhuber)
- 2006: unsupervised layerwise training of deep networks (Hinton et al.)

## Google Trends

Query: deep learning

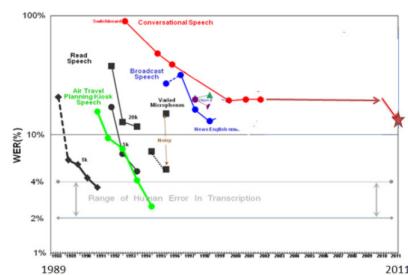


CS221

6

[figure from Li Deng]

## Speech recognition (2009-2011)



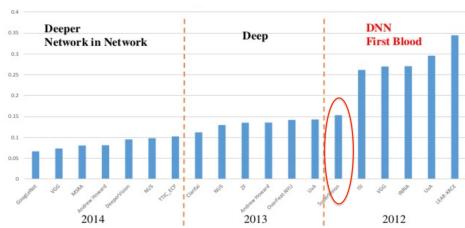
- Steep drop in WER due to deep learning
- IBM, Google, Microsoft all switched over from GMM-HMM

CS221

8

[Krizhevsky et al., 2012, a.k.a. AlexNet]

## Object recognition (2012)

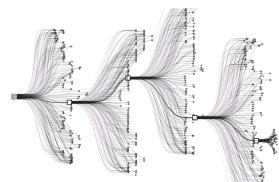
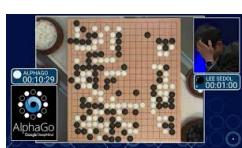


- Landslide win in ILSVRC object recognition competition
- Computer vision community switched to CNNs
- Simpler than hand-engineered features (SIFT)

CS221

10

## Go (2016)



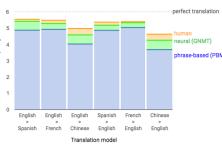
- Defeated world champion Le Sedol 4-1
- Simple architecture (in contrast, DeepBlue was search + hand-crafted heuristics)
- 2017: AlphaGoZero does not require human expert games as supervision

CS221

12

## Machine translation (2016)

Input sentence:	Translation (PHMT):	Translation (GNMT):	Translation (human):
李克强计划在月底同加拿大总理特鲁多会晤，这是两国领导人首次会面。	Li Keqiang will start the annual dialogue between Chinese and Canadian Prime Minister. This will be the first ministerial-level dialogue between the two premiers.	Li Keqiang will initiate the annual dialogue mechanism with Prime Minister Justin Trudeau of Canada during this visit, which will be the first ministerial-level dialogue with Premier Trudeau of Canada.	Li Keqiang will start the annual dialogue mechanism with Prime Minister Justin Trudeau of Canada during this visit, which will be the first ministerial-level dialogue with Premier Trudeau of Canada.



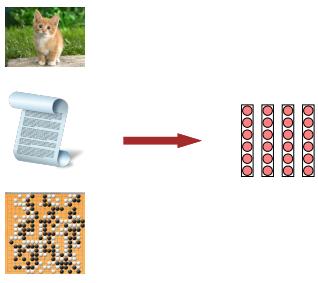
- Decisive wins have taken longer to achieve in NLP (words are meaningful in a way that pixels are not)
- Current state-of-the-art in machine translation
- Simpler architecture (throw out word alignment, phrases tables, language models)

CS221

14

## What is deep learning?

A family of techniques for learning compositional vector representations of complex data.



16



## Roadmap

Feedforward neural networks

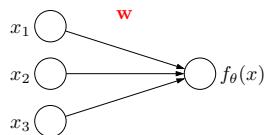
Convolutional neural networks

Recurrent neural networks

Unsupervised learning

Final remarks

## Review: linear predictors

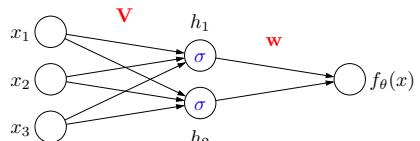


Output:

$$f_\theta(x) = \mathbf{w} \cdot \mathbf{x}$$

Parameters:  $\theta = \mathbf{w}$

## Review: neural networks



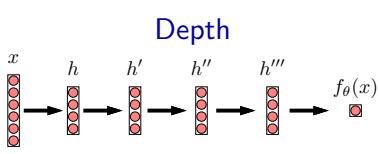
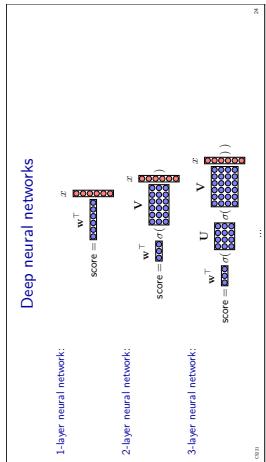
Intermediate hidden units:

$$h_j(x) = \sigma(\mathbf{v}_j \cdot \mathbf{x}) \quad \sigma(z) = (1 + e^{-z})^{-1}$$

Output:

$$f_\theta(x) = \mathbf{w} \cdot \mathbf{h}(x)$$

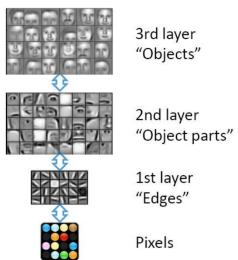
Parameters:  $\theta = (\mathbf{V}, \mathbf{w})$



Intuitions:

- Hierarchical feature representations
- Can simulate a bounded computation logic circuit (original motivation from McCulloch/Pitts, 1943)
- Learn this computation (and potentially more because networks are real-valued)
- Depth  $k + 1$  logic circuits can represent more than depth  $k$  (counting argument)
- Formal theory/understanding is still incomplete

## What's learned?



## Summary

- Deep networks learn hierarchical representations of data
- Train via SGD, use backpropagation to compute gradients
- Non-convex optimization, but works empirically given enough compute and data

CS221

30

## Review: optimization

Regression:

$$\text{Loss}(x, y, \theta) = (f_\theta(x) - y)^2$$

 **Key idea: minimize training loss**

$$\text{TrainLoss}(\theta) = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(x, y) \in \mathcal{D}_{\text{train}}} \text{Loss}(x, y, \theta)$$
$$\min_{\theta \in \mathbb{R}^d} \text{TrainLoss}(\theta)$$

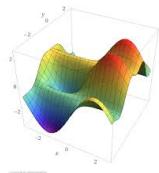
### Algorithm: stochastic gradient descent

```
For  $t = 1, \dots, T$ :  
  For  $(x, y) \in \mathcal{D}_{\text{train}}$ :  
     $\theta \leftarrow \theta - \eta_t \nabla_{\theta} \text{Loss}(x, y, \theta)$ 
```

CS221

32

## Training



- Non-convex optimization
- No theoretical guarantees that it works
- Before 2000s, empirically very difficult to get working

CS221

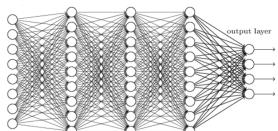
34

## What's different today

Computation (time/memory)      Information (data)

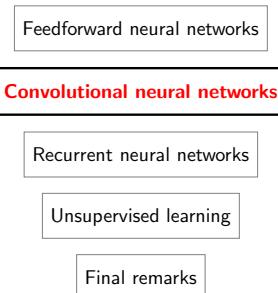


## How to make it work

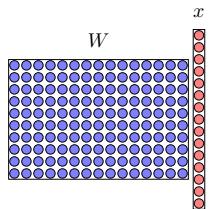


- More hidden units (over-provisioning)
- Adaptive step sizes (AdaGrad, ADAM)
- Dropout to guard against overfitting
- Careful initialization (pre-training)
- Batch normalization
- Model and optimization are tightly coupled

## Roadmap



## Motivation

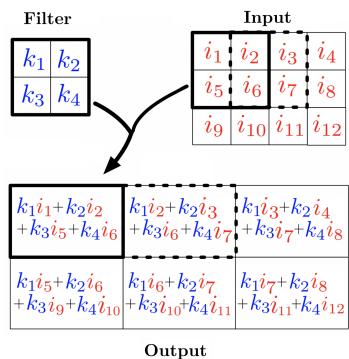


- **Observation:** images are not arbitrary vectors
- **Goal:** leverage spatial structure of images (translation invariance)

CS221

42

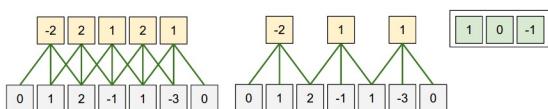
## Idea: Convolutions



CS221

44

## Prior knowledge

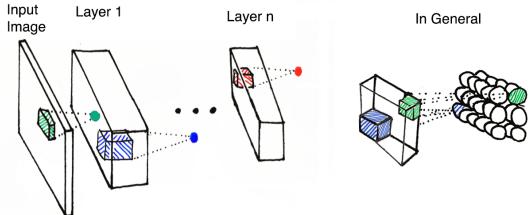


- **Local connectivity:** each hidden unit operates on a local image patch (3 instead of 7 connections per hidden unit)
- **Parameter sharing:** processing of each image patch is same (3 parameters instead of  $3 \cdot 5$ )
- **Intuition:** try to match a pattern in image

CS221

46

## Convolutional layers



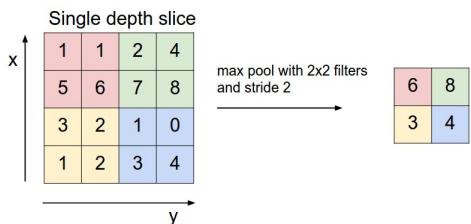
- Instead of vector to vector, we do volume to volume

CS221

48

[figure from Andrej Karpathy]

## Max-pooling

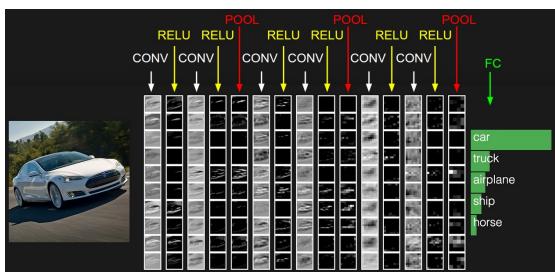


- Intuition: test if there exists a pattern in neighborhood
- Reduce computation, prevent overfitting

CS221

50

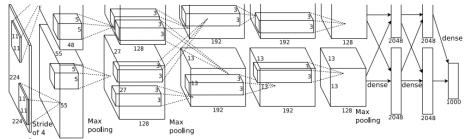
## Example of function evaluation



CS221

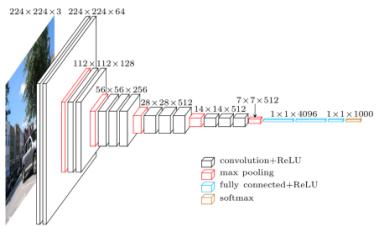
52

## AlexNet



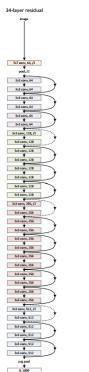
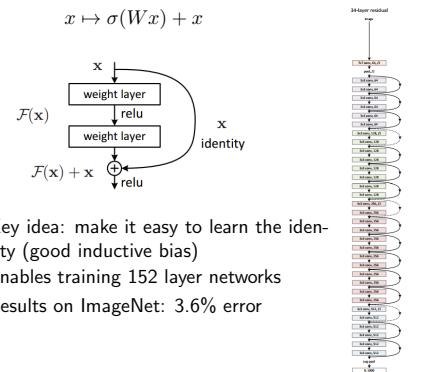
- **Non-linearity:** use ReLU ( $\max(z, 0)$ ) instead of logistic
- **Data augmentation:** translate, horizontal reflection, vary intensity, dropout (guard against overfitting)
- **Computation:** parallelize across two GPUs (6 days)
- **Results on ImageNet:** 16.4% error (next best was 25.8%)

## VGGNet



- **Architecture:** deeper but smaller filters; uniform
- **Computation:** 4 GPUs for 2-3 weeks
- **Results on ImageNet:** 7.3% error (AlexNet: 16.4%)

## Residual networks





## Summary

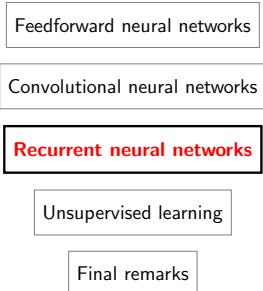
- Key idea: locality of connections, capture spatial structure
- Filters have parameter sharing; most parameters in last fully connected layers
- Depth really matters
- Applications to text, Go, drug design, etc.

CS221

60



## Roadmap



CS221

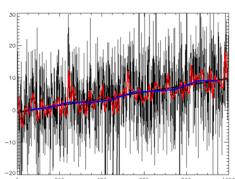
62

## Motivation: modeling sequences

Sentences:

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6 \quad x_7 \quad x_8 \quad x_9 \quad x_{10} \quad x_{11} \quad x_{12}$   
Paris Talks Set Stage for Action as Risks to the Climate Rise

Time series:

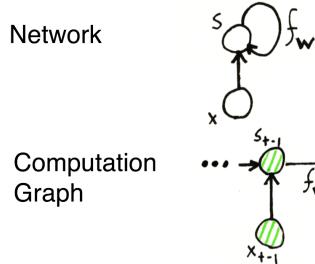


CS221

64

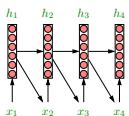
## Recurrent neural networks

Formula  $s_t = f_{\mathbf{W}}(s_{t-1}, x_t)$



66

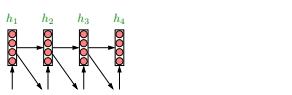
## Recurrent neural networks



$h_1 = \text{Encode}(x_1)$   
 $x_2 \sim \text{Decode}(h_1)$       Update context vector:  
 $h_2 = \text{Encode}(h_1, x_2)$        $h_t = \text{Encode}(h_{t-1}, x_t)$   
 $x_3 \sim \text{Decode}(h_2)$       Predict next character:  
 $h_3 = \text{Encode}(h_2, x_3)$        $x_{t+1} = \text{Decode}(h_t)$   
 $x_4 \sim \text{Decode}(h_3)$       context  $h_t$  compresses  $x_1, \dots, x_t$   
 $h_4 = \text{Encode}(h_3, x_4)$

68

## Simple recurrent network

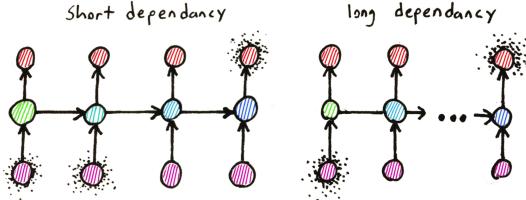


$$\text{Encode}(h_{t-1}, x_t) = \sigma(V h_{t-1} + W)$$

$$\text{Decode}(h_t) \sim \text{softmax}(W' h_t) = p(x_{t+1})$$

70

## Vanishing gradient problem

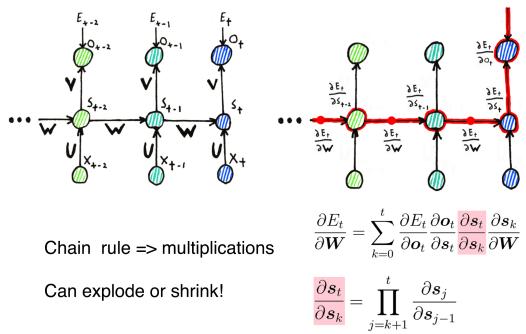


- RNNs can have long or short dependancies
- When there are long dependancies, gradients have trouble backpropagating through

CS221

72

## Vanishing gradient problem



CS221

74

## Long Short Term Memory (LSTM)

API:

$$(h_t, c_t) = \text{LSTM}(h_{t-1}, c_{t-1}, x_t)$$

Input gate:

$$i_t = \sigma(W_i x_t + U_i h_{t-1} + V_i c_{t-1} + b_i)$$

Forget gate (initialize with  $b_f$  large, so close to 1):

$$f_t = \sigma(W_f x_t + U_f h_{t-1} + V_f c_{t-1} + b_f)$$

Cell: additive combination of RNN update with previous cell

$$c_t = i_t \odot \tanh(W_c x_t + U_c h_{t-1} + b_c) + f_t \odot c_{t-1}$$

Output gate:

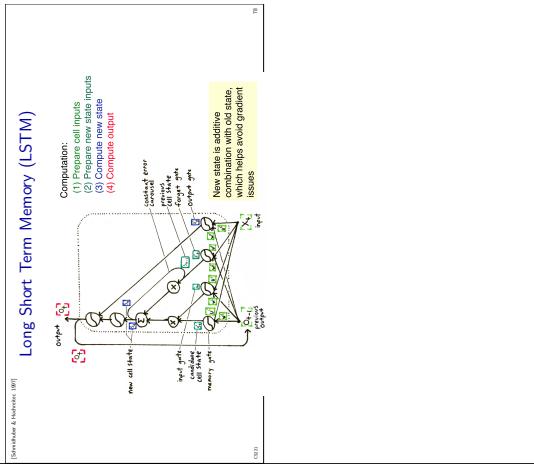
$$o_t = \sigma(W_o x_t + U_o h_{t-1} + V_o c_t + b_o)$$

Hidden state:

$$h_t = o_t \odot \tanh(c_t)$$

CS221

76



[from Andrej Karpathy's blog]

## Character-level language modeling

Sampled output:

*Naturalism and decision for the majority of Arab countries' capitalide was grounded by the Irish language by [[John Clair]], [[An Imperial Japanese Revolt]], associated with Guangzham's sovereignty. His generals were the powerful ruler of the Portugal in the [[Protestant Immineners]], which could be said to be directly in Cantonese Communication, which followed a ceremony and set inspired prison, training. The emperor travelled back to [[Antioch, Perth, October 25–21]] to note, the Kingdom of Costa Rica, unsuccessful fashioned the [[Thrales]], [[Cynth's Dajoard]], known in western [[Scotland]], near Italy to the conquest of India with the conflict.*

CS221

80

[from Andrej Karpathy's blog]

```
Cell sensitive to position in line:
The real importance of the crossing of the Berezina lies in the fact
that it plainly and indubitably proved the fallacy of all the plans of
the Russian general staff, and that it was the result of a single bold
line of action--the one Kutuzov, and the general mass of the army
had adopted. The French army had been marching for three days
at a continually increasing speed and all its energy was directed to
crossing the river. It had been compelled to do so because it had
to block its path. This was shown not so much by the arrangements it
had made for the crossing, but by the fact that the whole army
broke down, unarmed soldiers, people from Moscow and women with children
who had been driven from their homes, who had no clothes, pressed
forward into boats, and into the ice-covered water, and drowned.
```

Cell that turns on inside quotes:

```
You seem to imply that I have nothing to eat out of... On the
other hand, I have a good deal of time to eat, and I am invited to a
dinner party... I warmly repelled Chichagov, who tried by every word he
said to make me believe that I was a fool, and that I was incapable of
being animated by the same desire.
```

Kutuzov, however, replied with his subtle penetrating
glance, "I understand exactly what I said."

Cell that robustly activates inside if statements:

```
def __init__(self, max_length, mask):
    self.max_length = max_length
    self.mask = mask
    self._signals = []
    self._pending = []
    self._info = []

    self._decoder = Decoder(max_length, mask)

    self._decoder.set_start_token()
    self._decoder.set_end_token()

    self._decoder.set_max_length(max_length)
    self._decoder.set_min_length(1)
    self._decoder.set_stop_token(0)

    self._decoder.set_start_token()
    self._decoder.set_end_token()

    self._decoder.set_max_length(max_length)
    self._decoder.set_min_length(1)
    self._decoder.set_stop_token(0)
```

A large portion of cells are not easily interpretable. Here is a typical example:

```
def __init__(self, max_length, mask):
    self.max_length = max_length
    self.mask = mask
    self._decoder = Decoder(max_length, mask)
    self._decoder.set_start_token()
    self._decoder.set_end_token()
    self._decoder.set_max_length(max_length)
    self._decoder.set_min_length(1)
    self._decoder.set_stop_token(0)

    self._decoder.set_start_token()
    self._decoder.set_end_token()
    self._decoder.set_max_length(max_length)
    self._decoder.set_min_length(1)
    self._decoder.set_stop_token(0)

    self._decoder.set_start_token()
    self._decoder.set_end_token()
    self._decoder.set_max_length(max_length)
    self._decoder.set_min_length(1)
    self._decoder.set_stop_token(0)
```

CS221

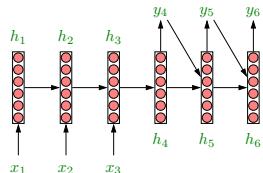
82

## Sequence-to-sequence model

Motivation: machine translation

$x$ : Je crains l'homme de un seul livre.

$y$ : Fear the man of one book.

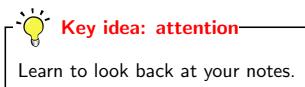


Read in a sentence first, output according to RNN:

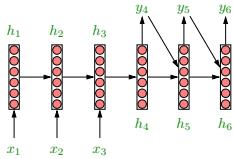
$$h_t = \text{Encode}(h_{t-1}, x_t \text{ or } y_{t-1}), \quad y_t = \text{Decode}(h_t)$$

## Attention-based models

Motivation: long sentences — compress to finite dimensional vector?



## Attention-based models



Distribution over input positions:

$$\alpha_t = \text{softmax}([\text{Attend}(h_1, h_{t-1}), \dots, \text{Attend}(h_L, h_{t-1})])$$

Generate with attended input:

$$h_t = \text{Encode}(h_{t-1}, y_{t-1}, \sum_{j=1}^L \alpha_t h_j)$$

Transformer models: attention only – no RNN!

## Machine translation

L'accord sur la zone économique européenne a été signé en août 1992.

The agreement on the European Economic Area was signed in August 1992.

<end>

## Image captioning

A woman is throwing a frisbee in a park.A dog is standing on a hardwood floor.A little girl sitting on a bed with a teddy bear.A group of people sitting on a boat in the water.

## Summary

- Recurrent neural networks: model sequences (non-linear version of Kalman filter or HMM)
- Logic intuition: learning a program with a for loop (reduce)
- LSTMs mitigate the vanishing gradient problem
- Attention-based models: when only part of input is relevant at a time
- Newer models with "external memory": memory networks, neural Turing machines



## Roadmap

Feedforward neural networks

Convolutional neural networks

Recurrent neural networks

**Unsupervised learning**

Final remarks

## Motivation

- Deep neural networks require lot of data
- Sometimes not very much labeled data, but plenty of unlabeled data (text, images, videos)
- Humans rarely get direct supervision; can learn from raw sensory information?

## Autoencoders

Analogy:

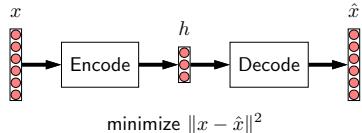
A A A A B B B B B B  $\rightarrow$  4 A's, 5 B's  $\rightarrow$  A A A A B B B B B B



**Key idea: autoencoders**

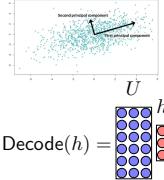
If we can compress a data point and still reconstruct it, then we have learned something generally useful.

General framework:



## Principal component analysis

**Input:** points  $x_1, \dots, x_n$



$$\text{Encode}(x) = U^\top x \quad \text{Decode}(h) = U h$$

(assume  $x_i$ 's are mean zero and  $U$  is orthogonal)

**PCA objective:**

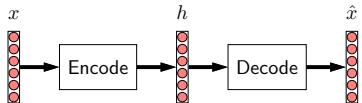
$$\text{minimize } \sum_{i=1}^n \|x_i - \text{Decode}(\text{Encode}(x_i))\|^2$$

CS221

102

## Autoencoders

Increase dimensionality of hidden dimension:



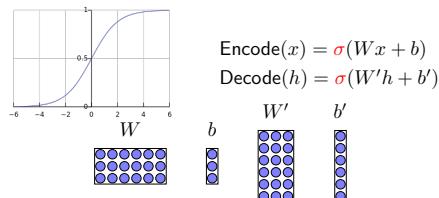
- **Problem:** learning nothing — just set Encode, Decode to identity function!
- Need to control complexity of Encode and Decode somehow...

CS221

104

## Non-linear autoencoders

Non-linear transformation (e.g., logistic function):



**Loss function:**

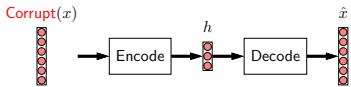
$$\text{minimize } \|x - \text{Decode}(\text{Encode}(x))\|^2$$

**Key:** non-linearity makes life harder, prevents degeneracy

CS221

106

## Denoising autoencoders



Types of noise:

- Blankout:  $\text{Corrupt}([1, 2, 3, 4]) = [0, 2, 3, 0]$
- Gaussian:  $\text{Corrupt}([1, 2, 3, 4]) = [1.1, 1.9, 3.3, 4.2]$

Objective:

$$\text{minimize } \|x - \text{Decode}(\text{Encode}(\text{Corrupt}(x)))\|^2$$

Algorithm: pick example, add fresh noise, SGD update

Key: noise makes life harder, prevents degeneracy

CS221

108

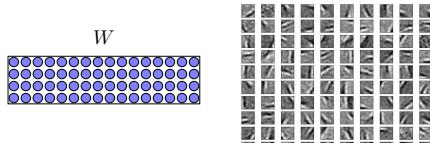
[Figure 7 of Vincent et al. (2010)]

## Denoising autoencoders

MNIST: 60,000 images of digits (784 dimensions)



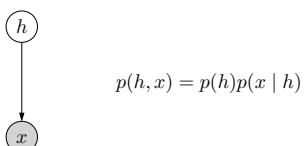
200 learned filters (rows of W):



110

## Variational autoencoders

Motivation: learn a latent-variable model



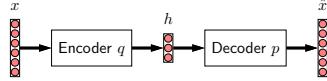
E-step in EM: computing  $p(h | x)$  is intractable

Solution: approximate using a neural network  $q(h | x)$

CS221

112

## Variational autoencoders



**Objective:** maximize

$$\log p(x) \geq \mathbb{E}_{q(h|x)}[\log p(x | h)] - \text{KL}(q(h | x)||p(h))$$

**Algorithm:**

- Sample  $h$  from encoder  $q$ , gradient update on  $q$  and  $p$
- Reparametrization trick [Kingma/Welling, 2014]

CS221

114

[Rajpurkar+ 2016]

## Reading comprehension (SQuAD)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **graupel** and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain in scattered locations are called "showers".

What causes precipitation to fall?  
**gravity**

100K examples

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?  
**graupel**

Where do water droplets collide with ice crystals to form precipitation?  
**within a cloud**

CS221

116

## Raw text

Stanford University (officially Leland Stanford Junior University)<sup>[1]</sup> (colloquially "The Farm") is a private research university in Stanford, California. Stanford is known for its academic strength, wealth, proximity to Silicon Valley, and ranking as one of the world's top universities.<sup>[2][3][4][5][6]</sup> The university was founded in 1891 by Leland and Jane Stanford in memory of their only child, Leland Stanford Jr., who had died of typhoid fever at age 15 the previous year. Stanford was a U.S. Senator and former Governor of California who made his fortune as a railroad tycoon. The school admitted its first students on October 1, 1891.<sup>[3][7]</sup> as a coeducational and non-denominational institution. Stanford University struggled financially after the death of Leland Stanford in 1893 and again after much of the campus was damaged by the 1906 San Francisco earthquake.<sup>[8][9]</sup> Following World War II, Provost Frederick Terman supported faculty and graduate entrepreneurs to build a self-sufficient local industry in what would later be known as Silicon Valley.<sup>[10]</sup> The university is also one of the top fundraising institutions in the United States.<sup>[11]</sup> The university is a member of the Association of American Universities and is often cited as a leading research institution.

The university is organized around three traditional schools consisting of 40 academic departments at the undergraduate and graduate level and four professional schools that focus on graduate programs in Law, Medicine, Education and Business. Stanford's undergraduate program is one of the top three most selective in the United States by acceptance rate.<sup>[12][13][14]</sup> Students compete in 36 varsity sports, and the university is one of two private institutions in the Division I Pac-12 Conference. It has gained 117 NCAA team championships<sup>[15]</sup> the most for a university. Stanford athletes have won 512 individual championships,<sup>[16]</sup> and Stanford has won the NACDA Directors' Cup for 23 consecutive years, beginning in 1994–1995.<sup>[17]</sup> In addition, Stanford students and alumni have won 270 Olympic medals including 139 gold medals.<sup>[18]</sup>

...

3.3 billion words

CS221

118

## Unsupervised pre-training

labeled

unlabeled

120

[Devlin+ 2018]

## BERT



Paris Talks ... Stage for \_\_\_\_\_ as Risks to ... Climate Rise



Paris Talks Set Stage for Action as Risks to the Climate Rise

- Tasks: fill in words, predict whether is next sentence

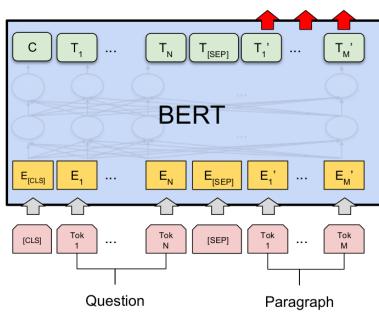
- Trained on 3.3B words, 4 days on 64 TPUs

CS221

122

## BERT

Start/End Span



124

CS221

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar et al. '16)	82.304	91.221
1 <small>Oct 05, 2018</small>	BERT (ensemble) Google A.I.	87.433	93.160
2 <small>Oct 10, 2018</small>	BERT (single model) Google A.I.	85.083	91.835
2 <small>Sep 09, 2018</small>	nlnet (ensemble) Microsoft Research Asia	85.356	91.202
2 <small>Sep 26, 2018</small>	nlnet (ensemble) Microsoft Research Asia	85.954	91.677
3 <small>Jul 11, 2018</small>	QANet (ensemble) Google Brain & CMU	84.454	90.490
4 <small>Jul 06, 2018</small>	r-net (ensemble) Microsoft Research Asia	84.003	90.147
5 <small>Mar 19, 2018</small>	QANet (ensemble) Google Brain & CMU	83.877	89.737
5 <small>Sep 09, 2018</small>	nlnet (single model) Microsoft Research Asia	83.468	90.133
5 <small>Jun 20, 2018</small>	MARS (ensemble) YUANFUDAO research NLP	83.982	89.796
6 <small>Sep 01, 2018</small>	MARS (single model) YUANFUDAO research NLP	83.185	89.547

126

CS221

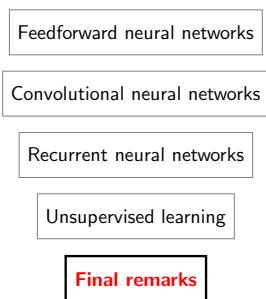
## Unsupervised learning

- Principle: make up prediction tasks (e.g.,  $x$  given  $x$  or context)
- Hard task → pressure to learn something
- Loss minimization using SGD
- Discriminatively fine tune: initialize feedforward neural network and backpropagate to optimize task accuracy
- How far can one push this?

128

CS221

## Roadmap



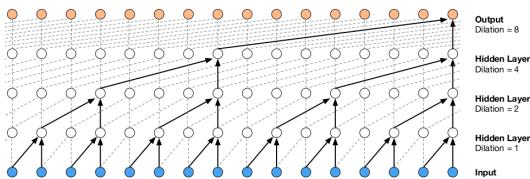
130

CS221

## WaveNet for audio generation

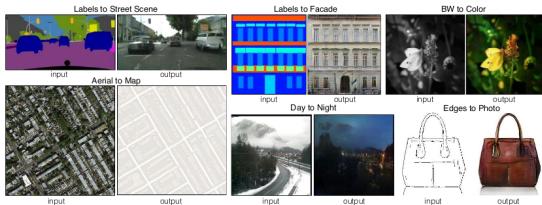


- Work with **raw** audio (16K observations / second)



- Key idea: **dilated convolutions** captures multiple scales of resolution, not recurrent

## Conditional adversarial networks



Key idea: game between

- **Generator:** generates fake images
- **Discriminator:** distinguishes between fake/real images

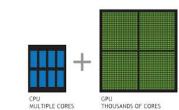
## Getting things to work

**Better optimization algorithms:** SGD, SGD+momentum, AdaGrad, AdaDelta, momentum, Nesterov, Adam

**Tricks:** initialization, gradient clipping, batch normalization, dropout

**More hyperparameter tuning:** step sizes, architectures

**Better hardware:** GPUs, TPUs



...wait for a long time...

## Theory: why does it work?

Two questions:

- Approximation: why are neural networks good hypothesis classes?
- Optimization: why can SGD optimize a high-dimensional non-convex problem?

Partial answers:

- 1-layer neural networks can approximate any continuous function on compact set [Cybenko, 1989; Barron, 1993]
- Generate random features works too [Rahimi/Recht, 2009; Andoni et. al, 2014]
- Use statistical physics to analyze loss surfaces [Choromanska et al., 2014]

CS221

138



## Summary

Phenomena	Ideas
Fixed vectors	Feedforward NNs
Spatial structure	convolutional NNs
Sequence	recurrent NNs LSTMs
Sequence-to-sequence	encoder-decoder attention-based models
Unsupervised	autoencoders variational autoencoders any auxiliary task

CS221

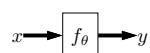
140

## Outlook

Extensibility: able to compose modules



Learning programs: think about analogy with a computer



CS221

142