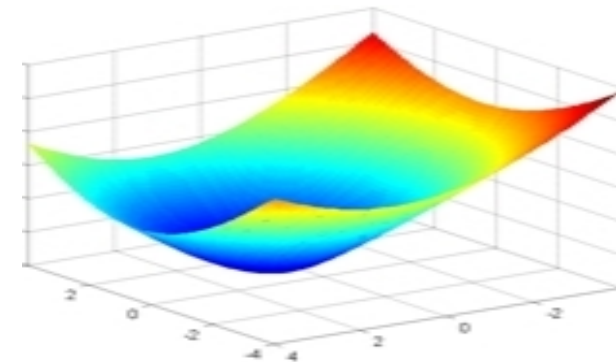# Deep Learning in NLP

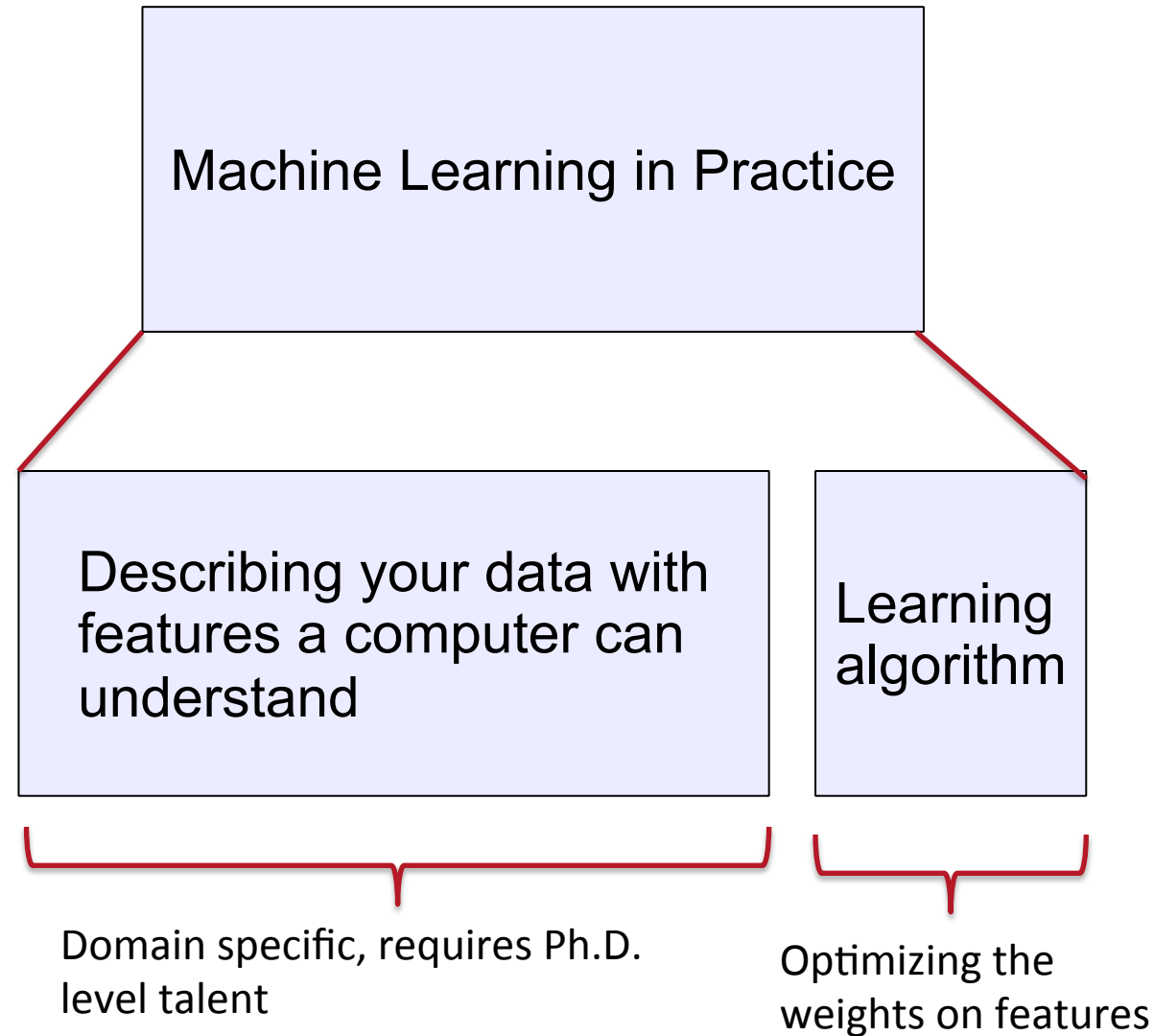Many slides adapted from Richard Socher, Tom Mitchell

# What's Deep Learning (DL)?

- Deep learning is a subfield of machine learning

- Most machine learning methods work well because of human-designed representations and input features
  - For example: features for finding named entities like locations or organization names (Finkel, 2010):

| Feature | NER |
|---|---|
| Current Word | ✓ |
| Previous Word | ✓ |
| Next Word | ✓ |
| Current Word Character n-gram | all |
| Current POS Tag | ✓ |
| Surrounding POS Tag Sequence | ✓ |
| Current Word Shape | ✓ |
| Surrounding Word Shape Sequence | ✓ |
| Presence of Word in Left Window | size 4 |
| Presence of Word in Right Window | size 4 |

- Machine learning becomes just optimizing weights to best make a final prediction
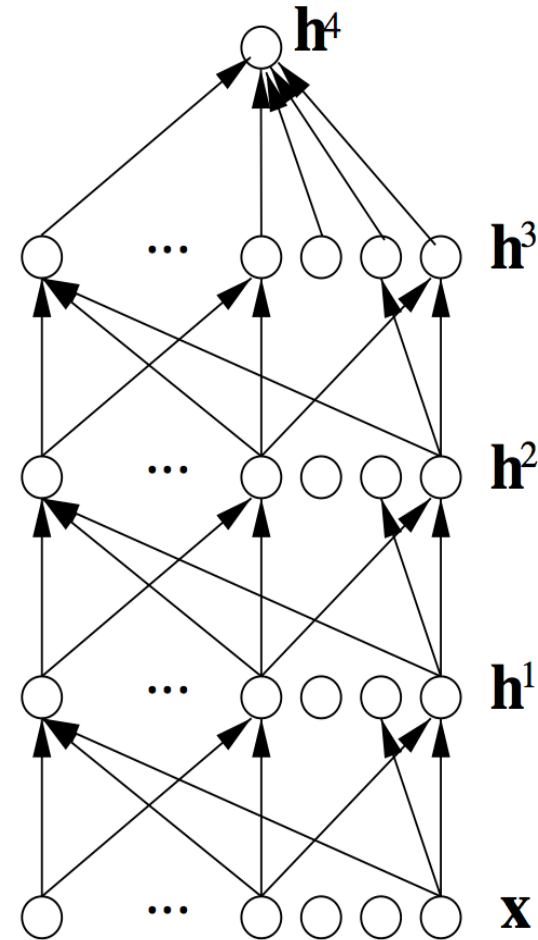
# Machine Learning vs Deep Learning

# What's Deep Learning (DL)?

- **Representation learning** attempts to automatically learn good features or representations

- **Deep learning** algorithms attempt to learn (multiple levels of) representation and an output

- From "raw" inputs **x** (e.g. words)

# On the history and term of "Deep Learning"

- We will focus on different kinds of **neural networks**
- The dominant model family inside deep learning

- Only clever terminology for stacked logistic regression units?
  - Somewhat, but interesting modeling principles and actual connections to neuroscience in some cases

# Reasons for Exploring Deep Learning

- Manually designed features are often over-specified, incomplete and take a long time to design and validate

- **Learned Features** are easy to adapt, fast to learn

- Deep learning provides a very flexible, (almost?) universal, learnable framework for **representing** world, visual and linguistic information.

- Deep learning can learn **unsupervised** (from raw text) and **supervised** (with specific labels like positive/negative)
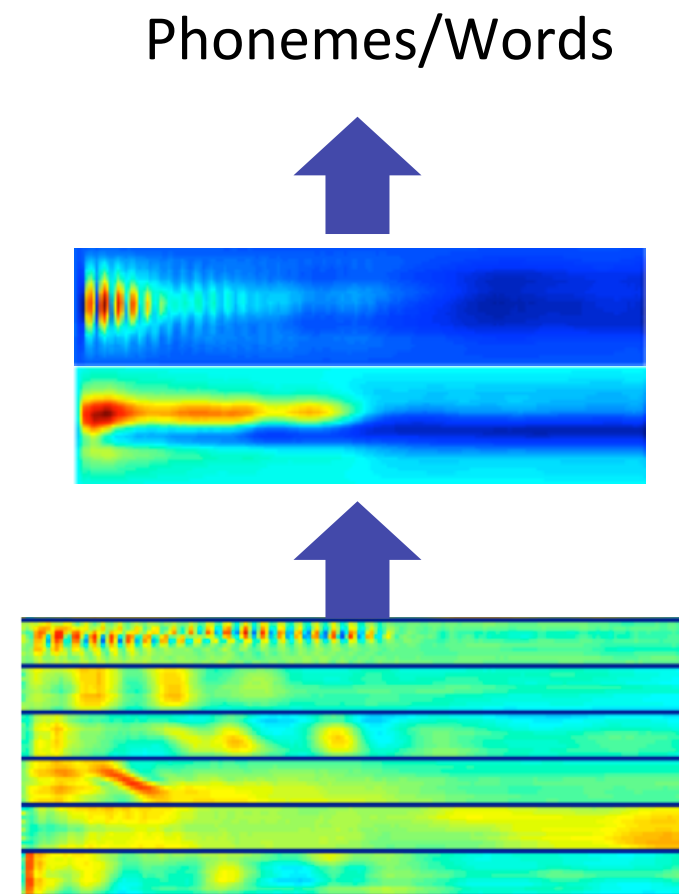
# Reasons for Exploring Deep Learning

- In 2006 **deep** learning techniques started outperforming other machine learning techniques. Why now?

- DL techniques benefit more from a lot of data

- Faster machines and multicore CPU/GPU help DL

- New models, algorithms, ideas

→ **Improved performance** (first in speech and vision, then NLP)
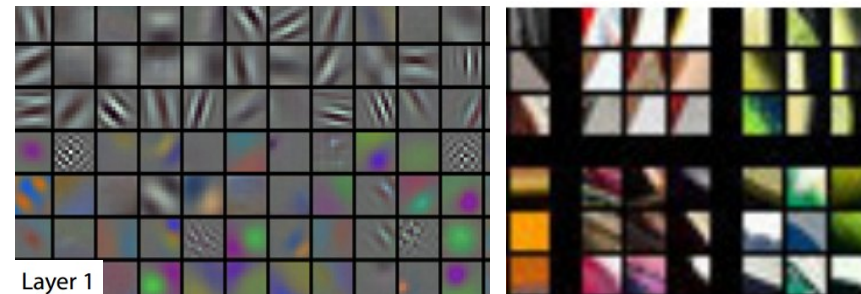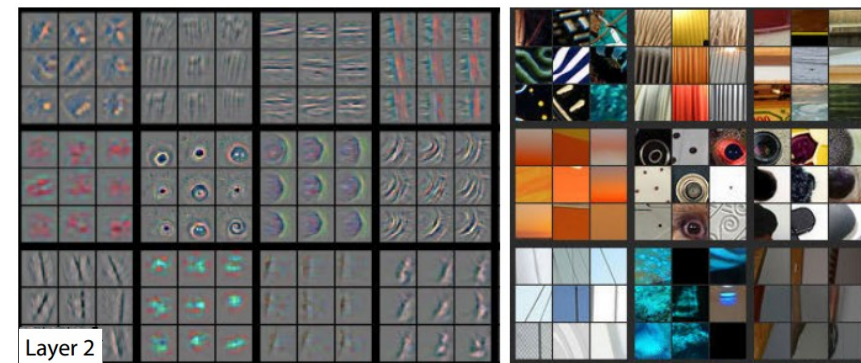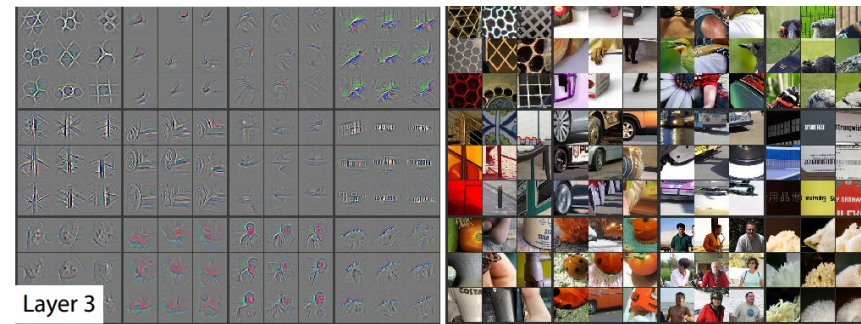
# Deep Learning for Speech

- The first breakthrough results of "deep learning" on large datasets happened in speech recognition

- Context-Dependent Pre-trained Deep Neural Networks for Large Vocabulary Speech Recognition Dahl et al. (2010)

| Acoustic model | Recog \ WER | RT03S FSH | Hub5 SWB |
|---|---|---|---|
| Traditional features | 1-pass –adapt | **27.4** | **23.6** |
| Deep Learning | 1-pass –adapt | **18.5** (−33%) | **16.1** (−32%) |

Phonemes/Words

# Deep Learning for Computer Vision

- Most deep learning groups have (until recently) largely focused on computer vision

- Break through paper: ImageNet Classification with Deep Convolutional Neural Networks by Krizhevsky et al. 2012
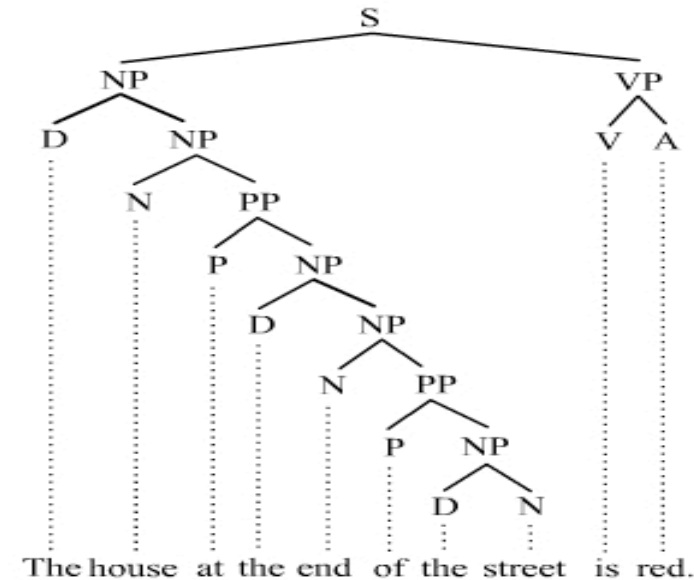


flamingo    cock    ruffed grouse    quail    partridge    …

Egyptian cat    Persian cat    Siamese cat    tabby    lynx    …



Layer 3

Layer 2

Layer 1

Zeiler and Fergus (2013)

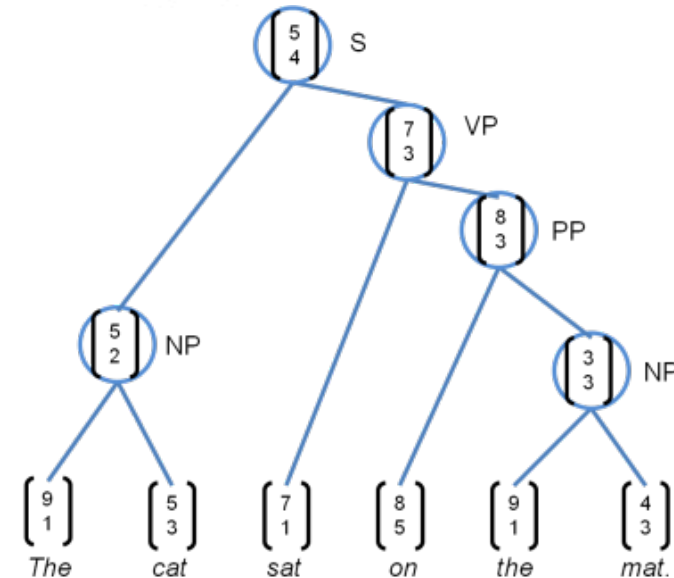# Neural word vectors - visualization

# Representations at NLP Levels: Syntax

- Traditional: Phrases
  Discrete categories like NP, VP

- DL:
  - Every word and every phrase is a vector
  - a neural network combines two vectors into one vector
  - Socher et al. 2011

# Machine Translation

- Many levels of translation have been tried in the past:

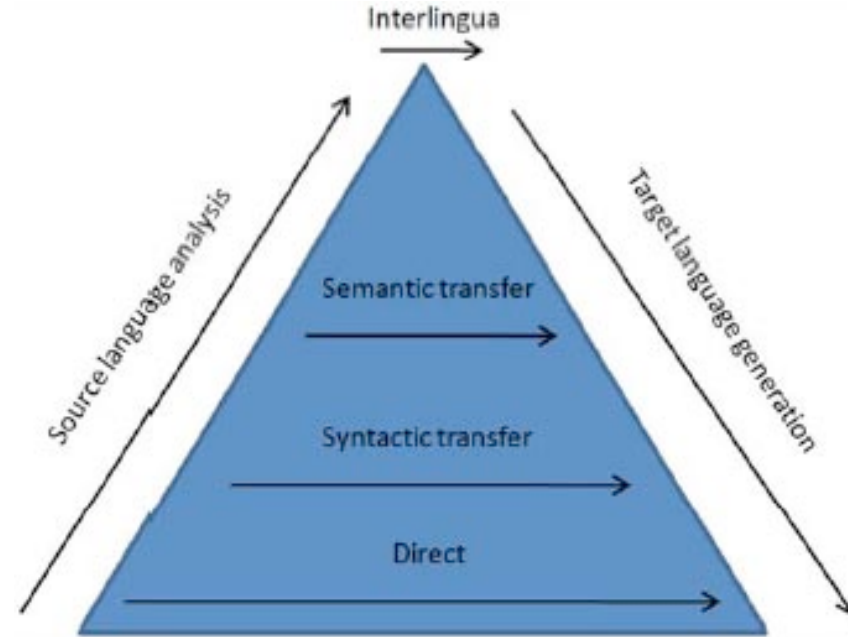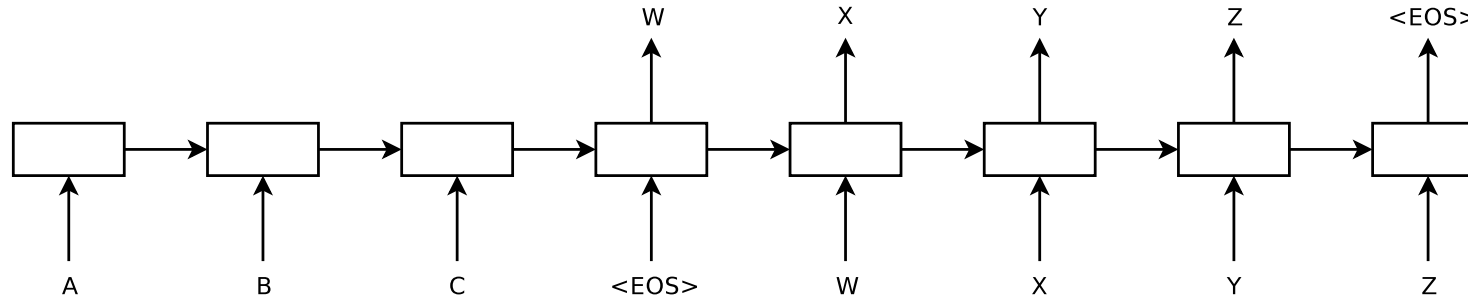- Traditional MT systems are very large complex systems



Figure 1: The Vauquois triangle

- What do you think is the interlingua for the DL approach to translation?

# Machine Translation

- Source sentence mapped to vector, then output sentence generated.



- Sequence to Sequence Learning with Neural Networks by Sutskever et al. 2014

- Very new but could replace very complex architectures!

# Neural Networks

# Connectionist Models

**Consider humans:**

- Neuron switching time $\sim .001$ second

- Number of neurons $\sim 10^{10}$

- Connections per neuron $\sim 10^{4-5}$

- Scene recognition time $\sim .1$ second

- 100 inference steps doesn't seem like enough

$\Rightarrow$ Much parallel computation

**Properties of neural nets:**

- Many neuron-like threshold switching units

- Many weighted interconnections among units

- Highly parallel, distributed process

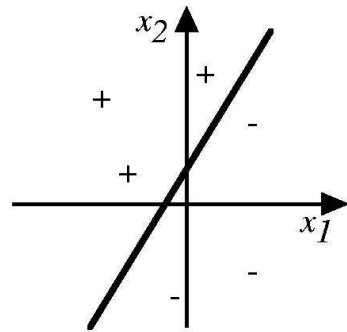- Emphasis on tuning weights automatically

# Perceptron



$$o(x_1, \ldots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \cdots + w_n x_n > 0 \\ -1 & \text{otherwise.} \end{cases}$$

Sometimes we'll use simpler vector notation:

$$o(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w} \cdot \vec{x} > 0 \\ -1 & \text{otherwise.} \end{cases}$$

# Decision Surface of a Perceptron



Represents some useful functions

- What weights represent $g(x_1, x_2) = AND(x_1, x_2)$?

But some functions not representable

- All not linearly separable
- Therefore, we'll want networks of these...

# Neural Nets for the Win!

- Neural networks can learn much more complex functions and nonlinear decision boundaries!

# Perceptron Training Rule

$$w_i \leftarrow w_i + \Delta w_i$$

where

$$\Delta w_i = \eta(t - o)x_i$$

Where:

- $t = c(\vec{x})$ is target value

- $o$ is perceptron output

- $\eta$ is small constant (e.g., 0.1) called *learning rate*

# Perceptron Training Rule

Can prove it will converge if

- Training data is linearly separable
- $\eta$ sufficiently small

# Gradient Descent

To understand, consider simpler *linear unit*, where

$$o = w_0 + w_1 x_1 + \cdots + w_n x_n$$

Let's learn $w_i$'s that minimize the squared error

$$E[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

Where $D$ is set of training examples

# Gradient Descent

Gradient:

$$\nabla E[\vec{w}] \equiv \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \cdots \frac{\partial E}{\partial w_n} \right]$$

Training rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

I.e.:

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

# Gradient Descent

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i}\frac{1}{2}\sum_d (t_d - o_d)^2$$

$$= \frac{1}{2}\sum_d \frac{\partial}{\partial w_i}(t_d - o_d)^2$$

$$= \frac{1}{2}\sum_d 2(t_d - o_d)\frac{\partial}{\partial w_i}(t_d - o_d)$$

$$= \sum_d (t_d - o_d)\frac{\partial}{\partial w_i}(t_d - \vec{w}\cdot\vec{x_d})$$

$$\frac{\partial E}{\partial w_i} = \sum_d (t_d - o_d)(-x_{i,d})$$

# Gradient Descent

GRADIENT-DESCENT$(training\_examples, \eta)$

Initialize each $w_i$ to some small random value

Until the termination condition is met, Do

- Initialize each $\Delta w_i$ to zero.

- For each $\langle \vec{x}, t \rangle$ in $training\_examples$, Do
  - Input instance $\vec{x}$ to unit and compute output $o$
  - For each linear unit weight $w_i$, Do

$$\Delta w_i \leftarrow \Delta w_i + \eta(t - o)x_i$$

- For each linear unit weight $w_i$, Do

$$w_i \leftarrow w_i + \Delta w_i$$

# Summary

Perceptron training rule guaranteed to succeed if

- Training examples are linearly separable
- Sufficiently small learning rate $\eta$

Linear unit training rule uses gradient descent

- Guaranteed to converge to hypothesis with minimum squared error
- Given sufficiently small learning rate $\eta$
- Even when training data contains noise
- Even when training data not separable by $H$

# Batch vs. Incremental Gradient Descent

**Batch Mode** Gradient Descent:

Do until convergence

1. Compute the gradient $\nabla E_D[\vec{w}]$

2. $\vec{w} \leftarrow \vec{w} - \eta \nabla E_D[\vec{w}]$

**Incremental Mode** Gradient Descent:

Do until convergence

    For each training example $d$ in $D$

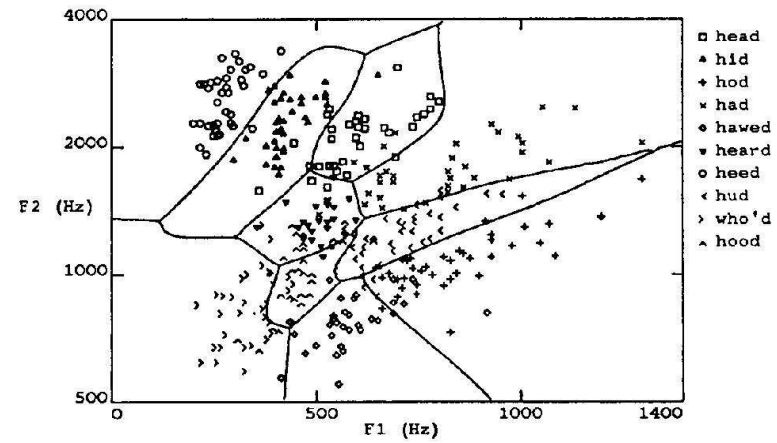        1. Compute the gradient $\nabla E_d[\vec{w}]$

        2. $\vec{w} \leftarrow \vec{w} - \eta \nabla E_d[\vec{w}]$

$$E_D[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$E_d[\vec{w}] \equiv \frac{1}{2} (t_d - o_d)^2$$

*Incremental Gradient Descent* can approximate *Batch Gradient Descent* arbitrarily closely if $\eta$ made small enough

# Multilayer Networks of Sigmoid Units

# Sigmoid Unit



$\sigma(x)$ is the sigmoid function

$$\frac{1}{1 + e^{-x}}$$

Nice property: $\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$

We can derive gradient descent rules to train

- One sigmoid unit

- *Multilayer networks* of sigmoid units $\rightarrow$ Backpropagation

# Error Gradient for a Sigmoid Unit

$$\frac{\partial E}{\partial w_i} \;=\; \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$=\; \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2$$

$$=\; \frac{1}{2} \sum_d 2(t_d - o_d) \; \frac{\partial}{\partial w_i} (t_d - o_d)$$

$$=\; \sum_d (t_d - o_d) \left( -\frac{\partial o_d}{\partial w_i} \right)$$

$$=\; -\sum_d (t_d - o_d) \; \frac{\partial o_d}{\partial net_d} \; \frac{\partial net_d}{\partial w_i}$$

But we know:

$$\frac{\partial o_d}{\partial net_d} = \frac{\partial \sigma(net_d)}{\partial net_d} = o_d(1 - o_d)$$

$$\frac{\partial net_d}{\partial w_i} = \frac{\partial(\vec{w} \cdot \vec{x}_d)}{\partial w_i} = x_{i,d}$$

So:

$$\frac{\partial E}{\partial w_i} = -\sum_{d \in D}(t_d - o_d)o_d(1 - o_d)x_{i,d}$$

Let: $\delta_k = -\frac{\partial E}{\partial net_k}$

$$\frac{\partial E}{\partial net_j} = \sum_{k \in Outs(j)} \frac{\partial E}{\partial net_k} \frac{\partial net_k}{\partial net_j}$$

$$= \sum_{k \in Outs(j)} -\delta_k \frac{\partial net_k}{\partial net_j}$$

$$= \sum_{k \in Outs(j)} -\delta_k \frac{\partial net_k}{\partial o_j} \frac{\partial o_j}{\partial net_j}$$

$$= \sum_{k \in Outs(j)} -\delta_k \, w_{kj} \frac{\partial o_k}{\partial net_j}$$

$$= \sum_{k \in Outs(j)} -\delta_k \, w_{kj} \, o_j(1 - o_j)$$

$$\delta_j = -\frac{\partial E}{\partial net_j} = o_j(1 - o_j) \sum_{k \in Outs(j)} \delta_k \, w_{kj}$$

# Backpropagation Algorithm

Initialize all weights to small random numbers
Until convergence, Do
    For each training example, Do

      1. Input it to network and compute network outputs
      2. For each output unit $k$

$$\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$$

      3. For each hidden unit $h$

$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in outputs} w_{h,k}\delta_k$$

      4. Update each network weight $w_{i,j}$

$$w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j}$$

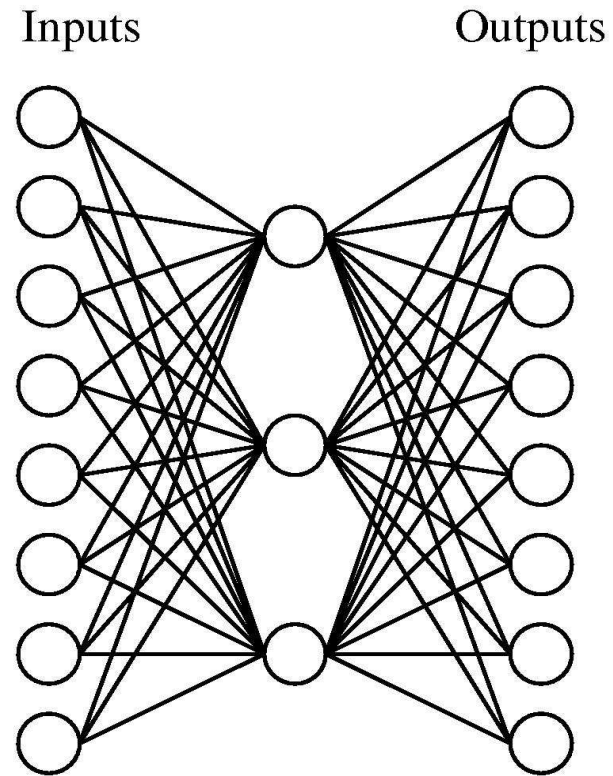    where $\Delta w_{i,j} = \eta \delta_j x_{i,j}$

# More on Backpropagation

- Gradient descent over entire *network* weight vector

- Easily generalized to arbitrary directed graphs

- Will find a local, not necessarily global error minimum
  - In practice, often works well
    (can run multiple times)

- Often include weight *momentum* $\alpha$

$$\Delta w_{i,j}(n) = \eta \delta_j x_{i,j} + \alpha \Delta w_{i,j}(n-1)$$

- Minimizes error over *training* examples
  - Will it generalize well to subsequent examples?

- Training can take thousands of iterations $\rightarrow$ slow!

- Using network after training is very fast
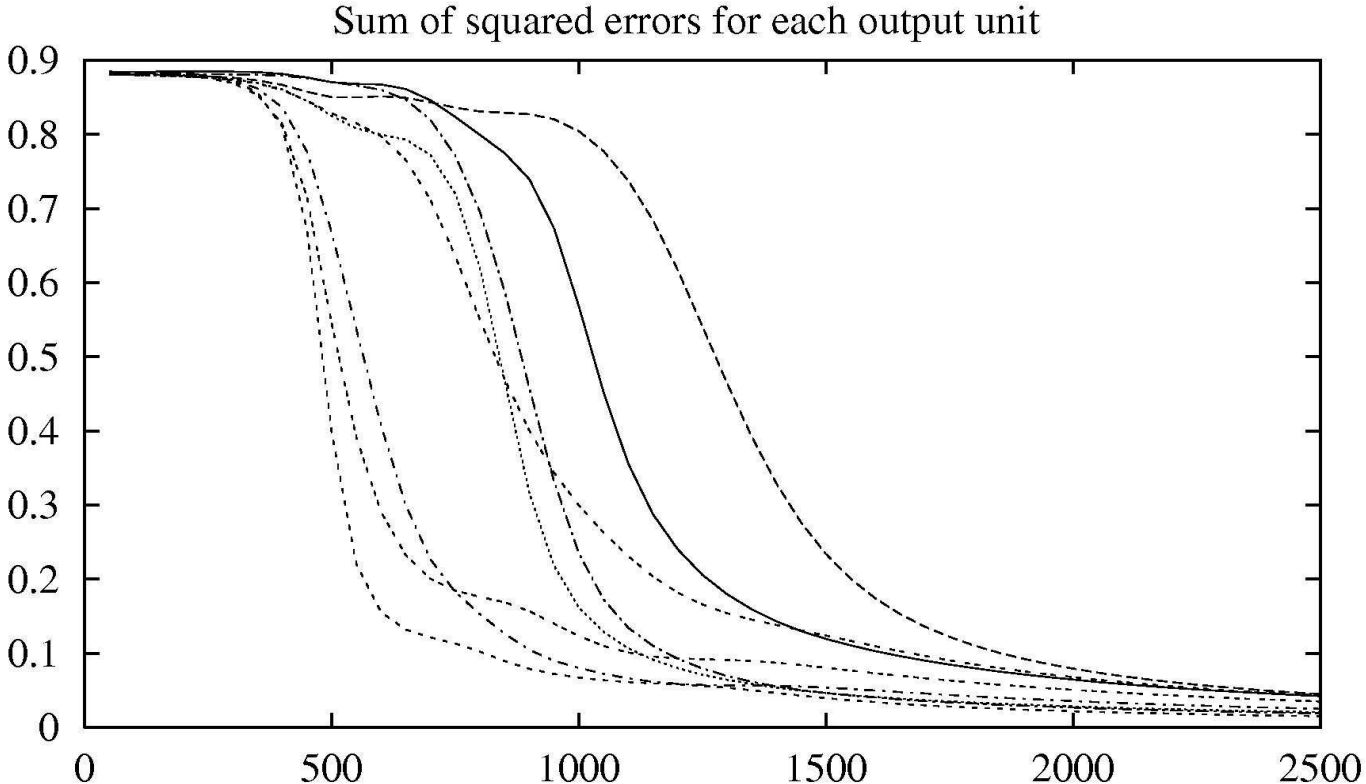
# Learning Hidden Layer Representations



Inputs                    Outputs

A target function:

| Input | | Output |
|-------|-----|--------|
| 10000000 | → | 10000000 |
| 01000000 | → | 01000000 |
| 00100000 | → | 00100000 |
| 00010000 | → | 00010000 |
| 00001000 | → | 00001000 |
| 00000100 | → | 00000100 |
| 00000010 | → | 00000010 |
| 00000001 | → | 00000001 |

Can this be learned?

Learned hidden layer representation:

| Input | | Hidden | | | | Output |
|---|---|---|---|---|---|---|
| | | | Values | | | |
| 10000000 | $\rightarrow$ | .89 | .04 | .08 | $\rightarrow$ | 10000000 |
| 01000000 | $\rightarrow$ | .01 | .11 | .88 | $\rightarrow$ | 01000000 |
| 00100000 | $\rightarrow$ | .01 | .97 | .27 | $\rightarrow$ | 00100000 |
| 00010000 | $\rightarrow$ | .99 | .97 | .71 | $\rightarrow$ | 00010000 |
| 00001000 | $\rightarrow$ | .03 | .05 | .02 | $\rightarrow$ | 00001000 |
| 00000100 | $\rightarrow$ | .22 | .99 | .99 | $\rightarrow$ | 00000100 |
| 00000010 | $\rightarrow$ | .80 | .01 | .98 | $\rightarrow$ | 00000010 |
| 00000001 | $\rightarrow$ | .60 | .94 | .01 | $\rightarrow$ | 00000001 |

# Training

Sum of squared errors for each output unit

# Training

Hidden unit encoding for input 01000000
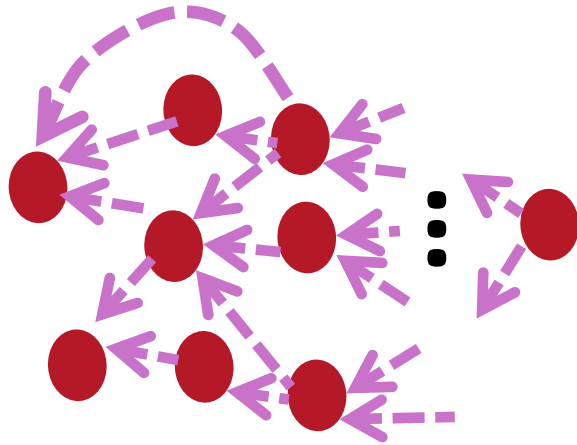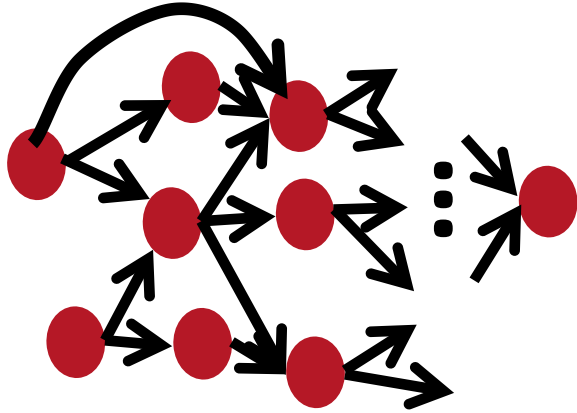
# Automatic Differentiation



- The gradient computation can be **automatically inferred** from the symbolic expression of the fprop.

- Each node type needs to know how to compute its output and how to compute the gradient wrt its inputs given the gradient wrt its output.
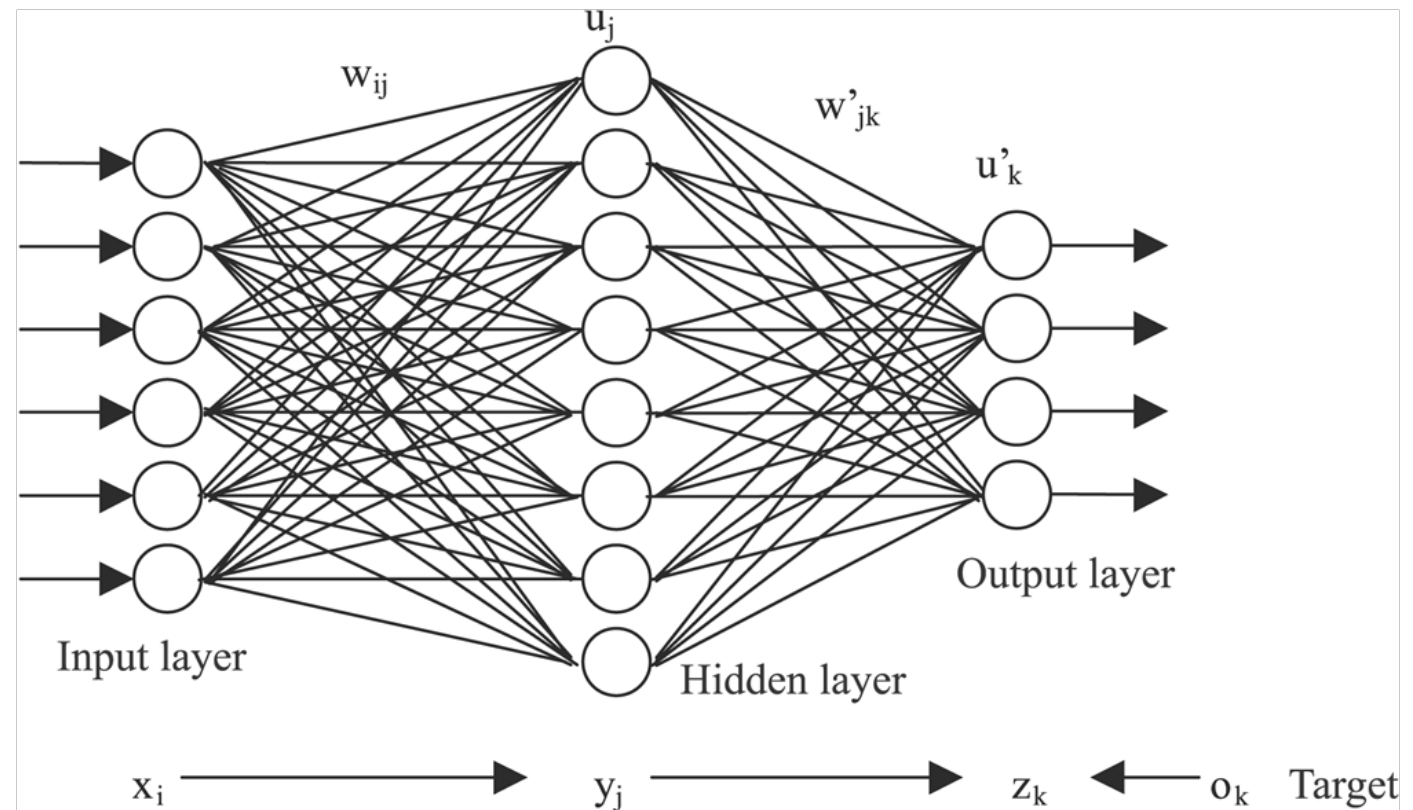
- Easy and fast prototyping

# Neural Network Language Models

# Review

- Deep Learning
  - Learning Representations of Inputs
- Neural Networks
  - Layers of Logistic Regression
  - Can represent any nonlinear function (with a large enough network)
  - Training with backpropagation
- Recent breakthroughs in predictive tasks
  - Speech Recognition
  - Object Recognition (computer vision)

# Q: How to model sequences with neural Networks?

- Fixed number of inputs.

# How about just predicting the next word in the input?

- Q: what about just predicting the next word?

- From the context?
  - No longer a language model

- Word2Vec!

# Word2vec

- Learn continuous word embedding for each word
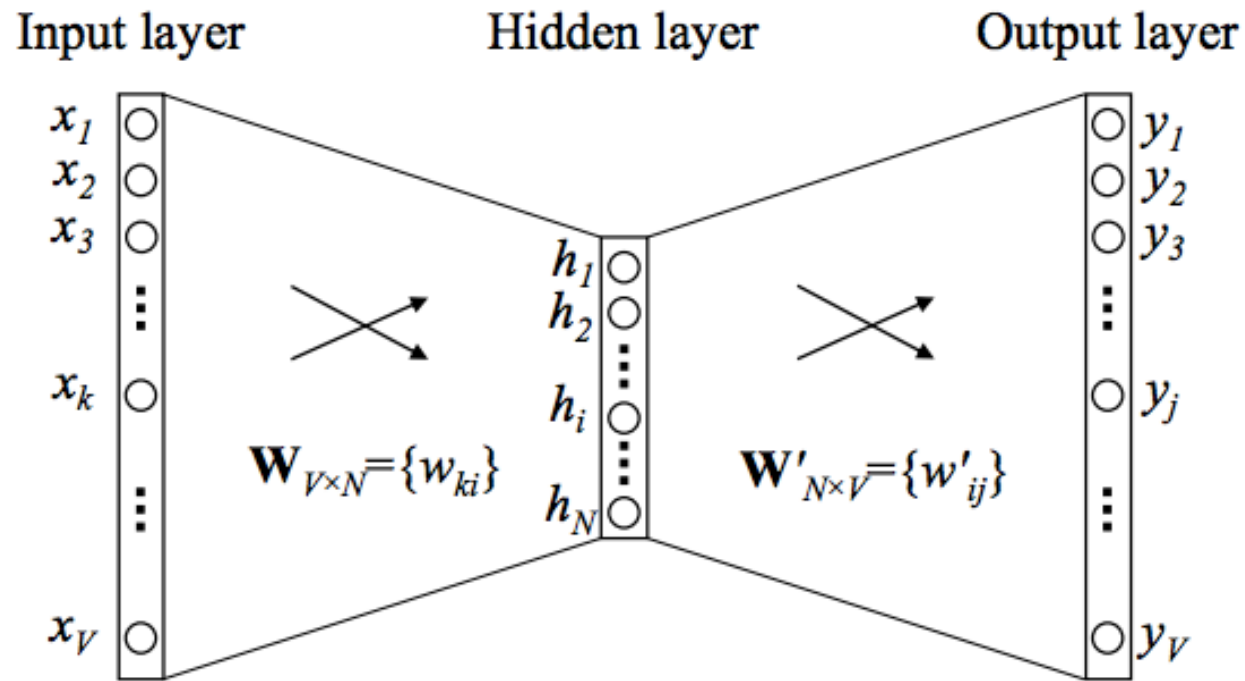  - Each word represented by a vector



Figure 1: A simple CBOW model with only one word in the context

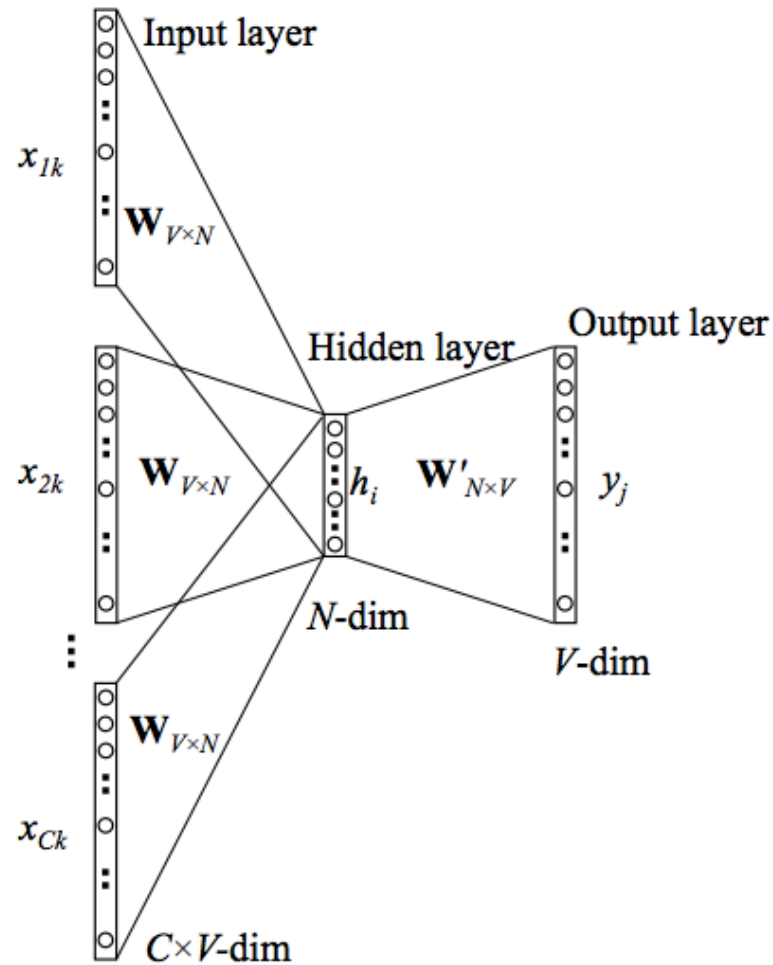# Using more than one word of context
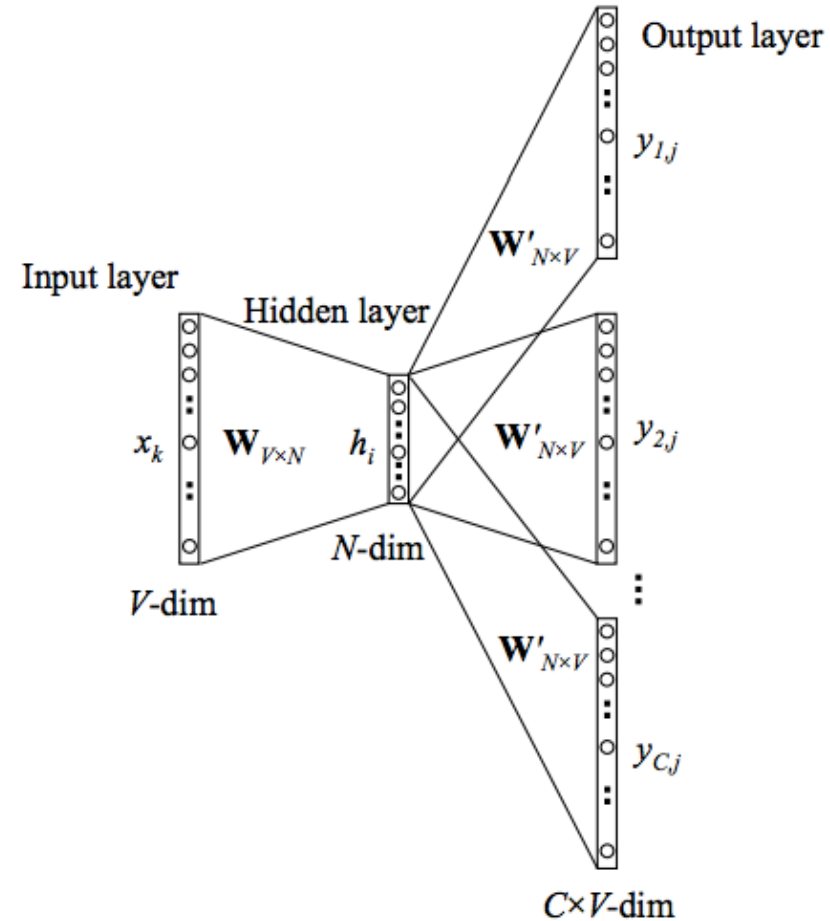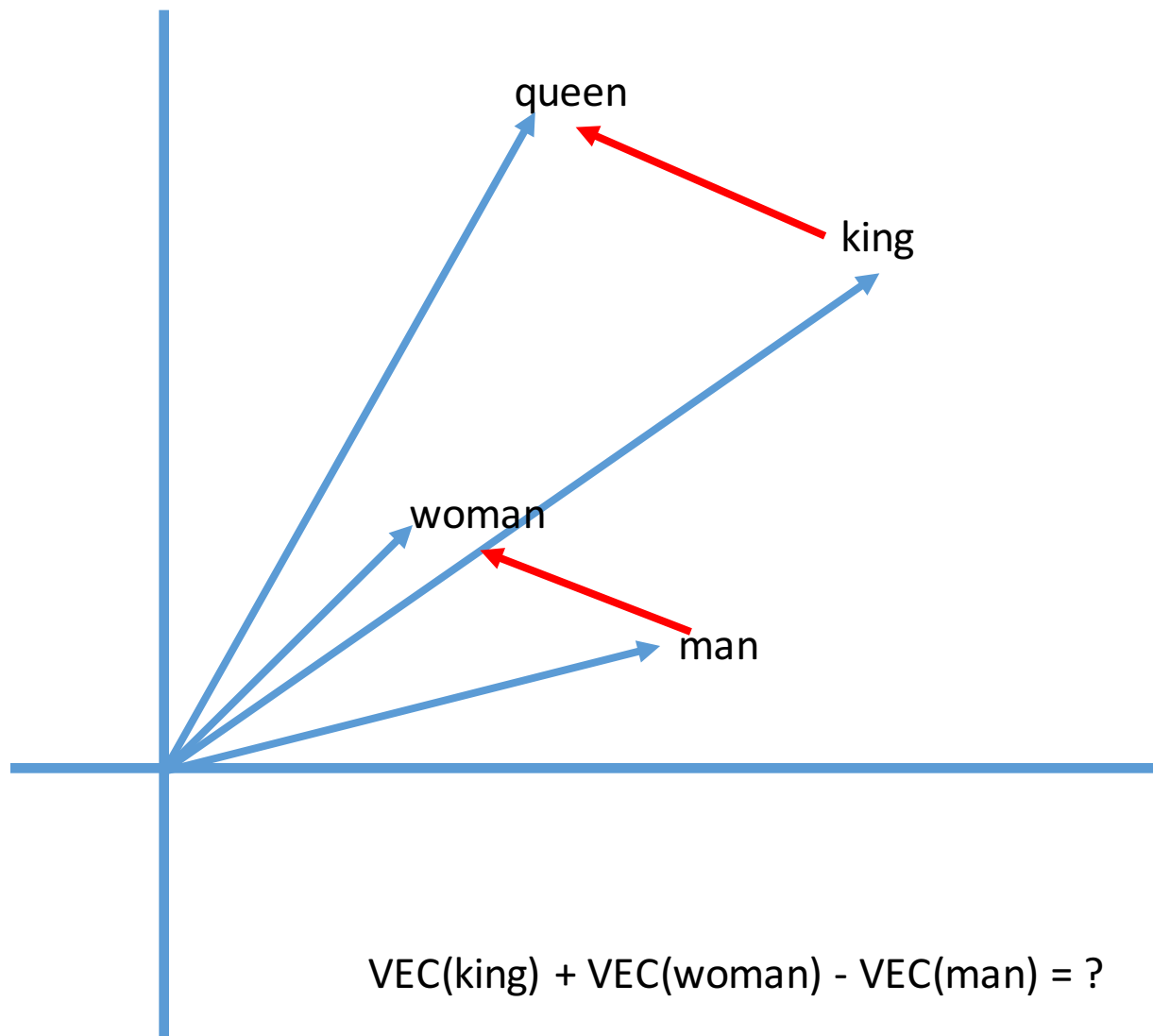


Figure 2: Continuous bag-of-word model

Figure 3: The skip-gram model.

# Word2Vec: fast to train

- Word2Vec is a fairly simple model,
- But Can efficiently train word vectors on really big corpora
- This is probably the main advantage of Word2vec over other approaches...
  - Principal Component Analysis
  - Recurrent Neural Network Language Models
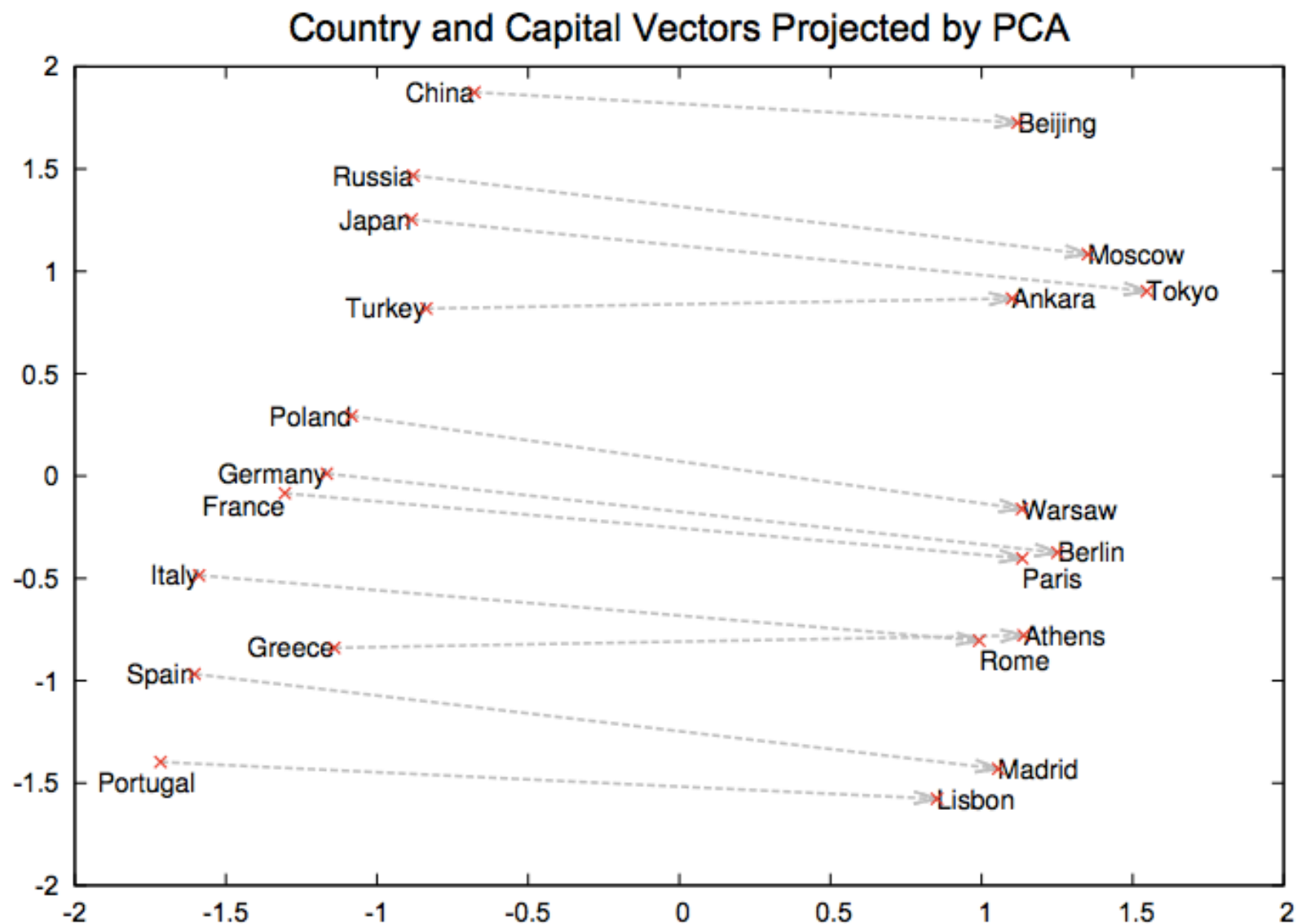
VEC(king) + VEC(woman) - VEC(man) = ?

Figure 2: Two-dimensional PCA projection of the 1000-dimensional Skip-gram vectors of countries and their capital cities. The figure illustrates ability of the model to automatically organize concepts and learn implicitly the relationships between them, as during the training we did not provide any supervised information about what a capital city means.

# The Unreasonable Effectiveness of Word Representations for Twitter Named Entity Recognition

**Colin Cherry** and **Hongyu Guo**
National Research Council Canada
first.last@nrc-cnrc.gc.ca

| System | Fin10Dev | Rit11 | Fro14 | Avg |
|---|---|---|---|---|
| CoNLL | 27.3 | 27.1 | 29.5 | 28.0 |
| + Brown | 38.4 | 39.4 | 42.5 | 40.1 |
| + Vector | 40.8 | 40.4 | 42.9 | 41.4 |
| + Reps | 42.4 | 42.2 | 46.2 | 43.6 |
| Fin10 | 36.7 | 29.0 | 30.4 | 32.0 |
| + Brown | 59.9 | 53.9 | 56.3 | 56.7 |
| + Vector | 61.5 | 56.4 | 58.4 | 58.8 |
| + Reps | 64.0 | 58.5 | 60.2 | 60.9 |
| CoNLL+Fin10 | 44.7 | 39.9 | 44.2 | 42.9 |
| + Brown | 54.9 | 52.9 | 58.5 | 55.4 |
| + Vector | 58.9 | 55.2 | 59.9 | 58.0 |
| + Reps | 58.9 | 56.4 | 61.8 | 59.0 |
| + Weights | 64.4 | 59.6 | 63.3 | 62.4 |

Table 5: Impact of our components on Twitter NER performance, as measured by F1, under 3 data scenarios.

# Language Modeling

- $x$ is a "history" $w_1, w_2, \ldots w_{i-1}$, e.g.,

  *Third, the notion "grammatical in English" cannot be identified in any way with the notion "high order of statistical approximation to English". It is fair to assume that neither sentence (1) nor (2) (nor indeed any part of these sentences) has ever occurred in an English discourse. Hence, in any statistical*

- $y$ is an "outcome" $w_i$
- Example features:

$$f_1(x, y) = \begin{cases} 1 & \text{if } y = \texttt{model} \\ 0 & \text{otherwise} \end{cases}$$

$$f_2(x, y) = \begin{cases} 1 & \text{if } y = \texttt{model} \text{ and } w_{i-1} = \texttt{statistical} \\ 0 & \text{otherwise} \end{cases}$$

$$f_3(x, y) = \begin{cases} 1 & \text{if } y = \texttt{model}, w_{i-2} = \texttt{any}, w_{i-1} = \texttt{statistical} \\ 0 & \text{otherwise} \end{cases}$$

$$f_4(x, y) = \begin{cases} 1 & \text{if } y = \texttt{model}, w_{i-2} = \text{any} \\ 0 & \text{otherwise} \end{cases}$$

$$f_5(x, y) = \begin{cases} 1 & \text{if } y = \texttt{model}, w_{i-1} \text{ is an adjective} \\ 0 & \text{otherwise} \end{cases}$$

$$f_6(x, y) = \begin{cases} 1 & \text{if } y = \texttt{model}, w_{i-1} \text{ ends in "ical"} \\ 0 & \text{otherwise} \end{cases}$$

$$f_7(x, y) = \begin{cases} 1 & \text{if } y = \texttt{model}, \text{author} = \text{Chomsky} \\ 0 & \text{otherwise} \end{cases}$$

$$f_8(x, y) = \begin{cases} 1 & \text{if } y = \texttt{model}, \text{"model" is not in } w_1, \ldots w_{i-1} \\ 0 & \text{otherwise} \end{cases}$$

$$f_9(x, y) = \begin{cases} 1 & \text{if } y = \texttt{model}, \text{"grammatical" is in } w_1, \ldots w_{i-1} \\ 0 & \text{otherwise} \end{cases}$$

# Defining Features in Practice

▶ We had the following "trigram" feature:

$$f_3(x, y) = \begin{cases} 1 & \text{if } y = \texttt{model}, \ w_{i-2} = \texttt{any}, \ w_{i-1} = \texttt{statistical} \\ 0 & \text{otherwise} \end{cases}$$

▶ In practice, we would probably introduce one trigram feature for every trigram seen in the training data: i.e., for all trigrams $(u, v, w)$ seen in training data, create a feature

$$f_{N(u,v,w)}(x, y) = \begin{cases} 1 & \text{if } y = w, \ w_{i-2} = u, \ w_{i-1} = v \\ 0 & \text{otherwise} \end{cases}$$

where $N(u, v, w)$ is a function that maps each $(u, v, w)$ trigram to a different integer

# Language Modeling

▶ $x$ is a "history" $w_1, w_2, \ldots w_{i-1}$, e.g.,

*Third, the notion "grammatical in English" cannot be identified in any way with the notion "high order of statistical approximation to English". It is fair to assume that neither sentence (1) nor (2) (nor indeed any part of these sentences) has ever occurred in an English discourse. Hence, in any statistical*

▶ Each possible $y$ gets a different score:

$$v \cdot f(x, model) = 5.6 \qquad v \cdot f(x, the) = -3.2$$
$$v \cdot f(x, is) = 1.5 \qquad v \cdot f(x, of) = 1.3$$
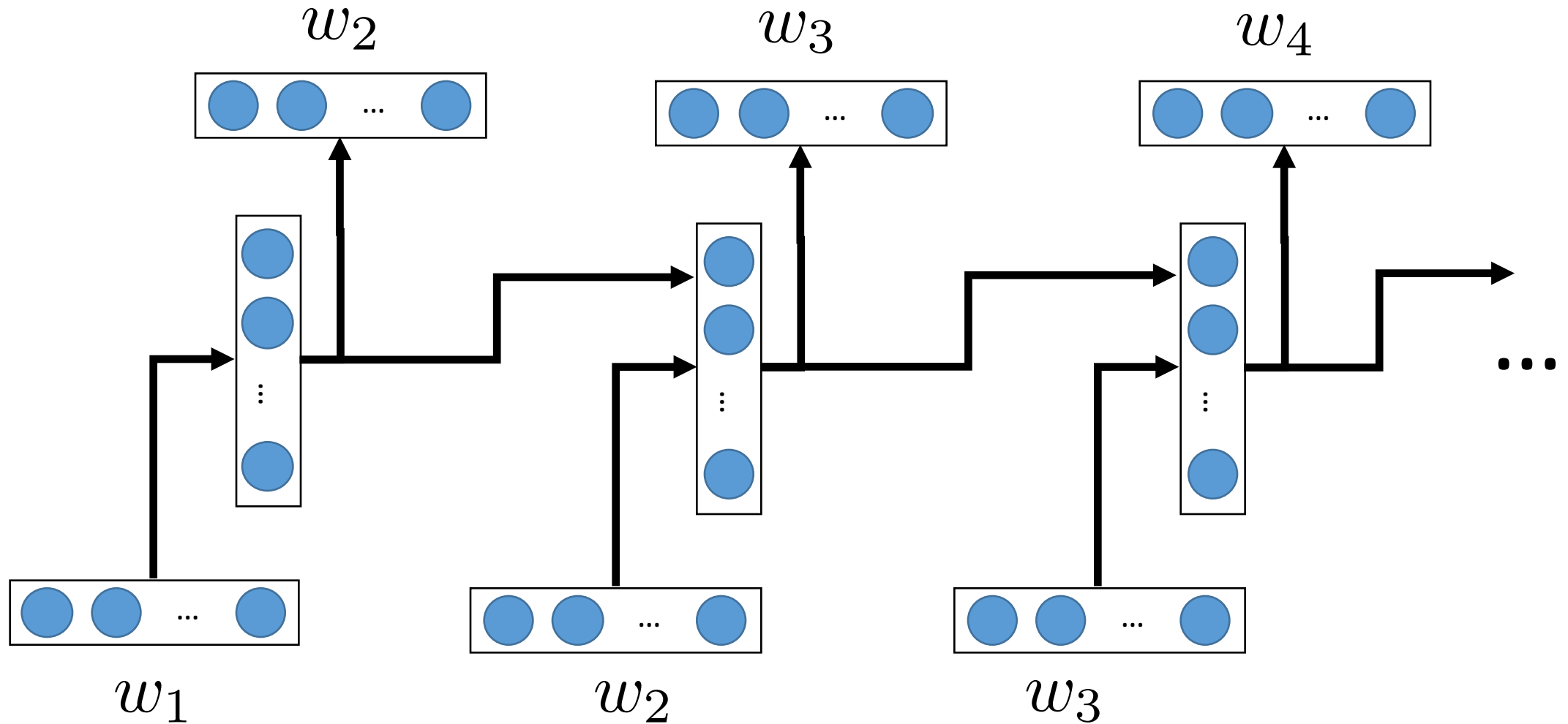$$v \cdot f(x, models) = 4.5 \qquad \ldots$$

# Log-Linear Models

- We have some input domain $\mathcal{X}$, and a finite label set $\mathcal{Y}$. Aim is to provide a conditional probability $p(y \mid x)$ for any $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.

- A feature is a function $f : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$
  (Often binary features or indicator functions
  $f_k : \mathcal{X} \times \mathcal{Y} \to \{0, 1\}$).

- Say we have $m$ features $f_k$ for $k = 1 \ldots m$
  $\Rightarrow$ A feature vector $f(x, y) \in \mathbb{R}^m$ for any $x \in \mathcal{X}$ and $y \in \mathcal{Y}$.

- We also have a **parameter vector** $v \in \mathbb{R}^m$

- We define

$$p(y \mid x; v) = \frac{e^{v \cdot f(x,y)}}{\sum_{y' \in \mathcal{Y}} e^{v \cdot f(x,y')}}$$

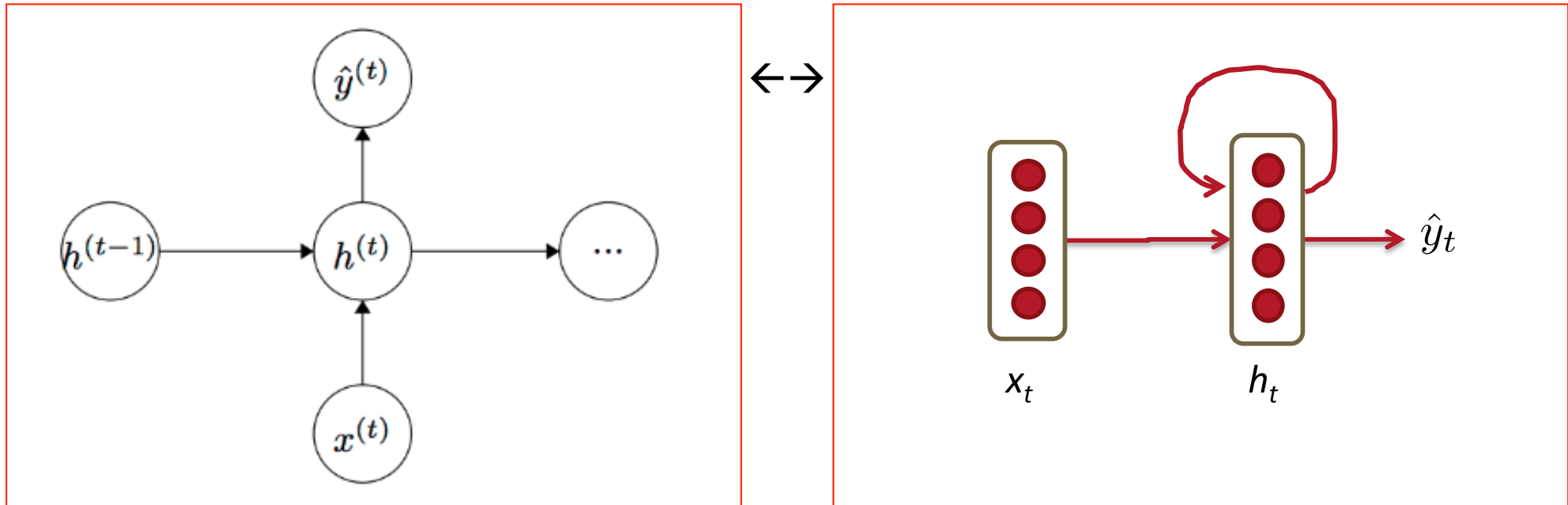# Recurrent Neural Network Language Models

# Recurrent Neural Network language model

Given list of word **vectors**: $x_1, \ldots, x_{t-1}, x_t, x_{t+1}, \ldots, x_T$

At a single time step:

$$h_t = \sigma\left(W^{(hh)} h_{t-1} + W^{(hx)} x_{[t]}\right)$$

$$\hat{y}_t = \text{softmax}\left(W^{(S)} h_t\right)$$

$$\hat{P}(x_{t+1} = v_j \mid x_t, \ldots, x_1) = \hat{y}_{t,j}$$

# Recurrent Neural Network language model

Main idea: we use the same set of W weights at all time steps!

Everything else is the same:

$$h_t = \sigma\left(W^{(hh)}h_{t-1} + W^{(hx)}x_{[t]}\right)$$

$$\hat{y}_t = \text{softmax}\left(W^{(S)}h_t\right)$$

$$\hat{P}(x_{t+1} = v_j \mid x_t, \ldots, x_1) = \hat{y}_{t,j}$$
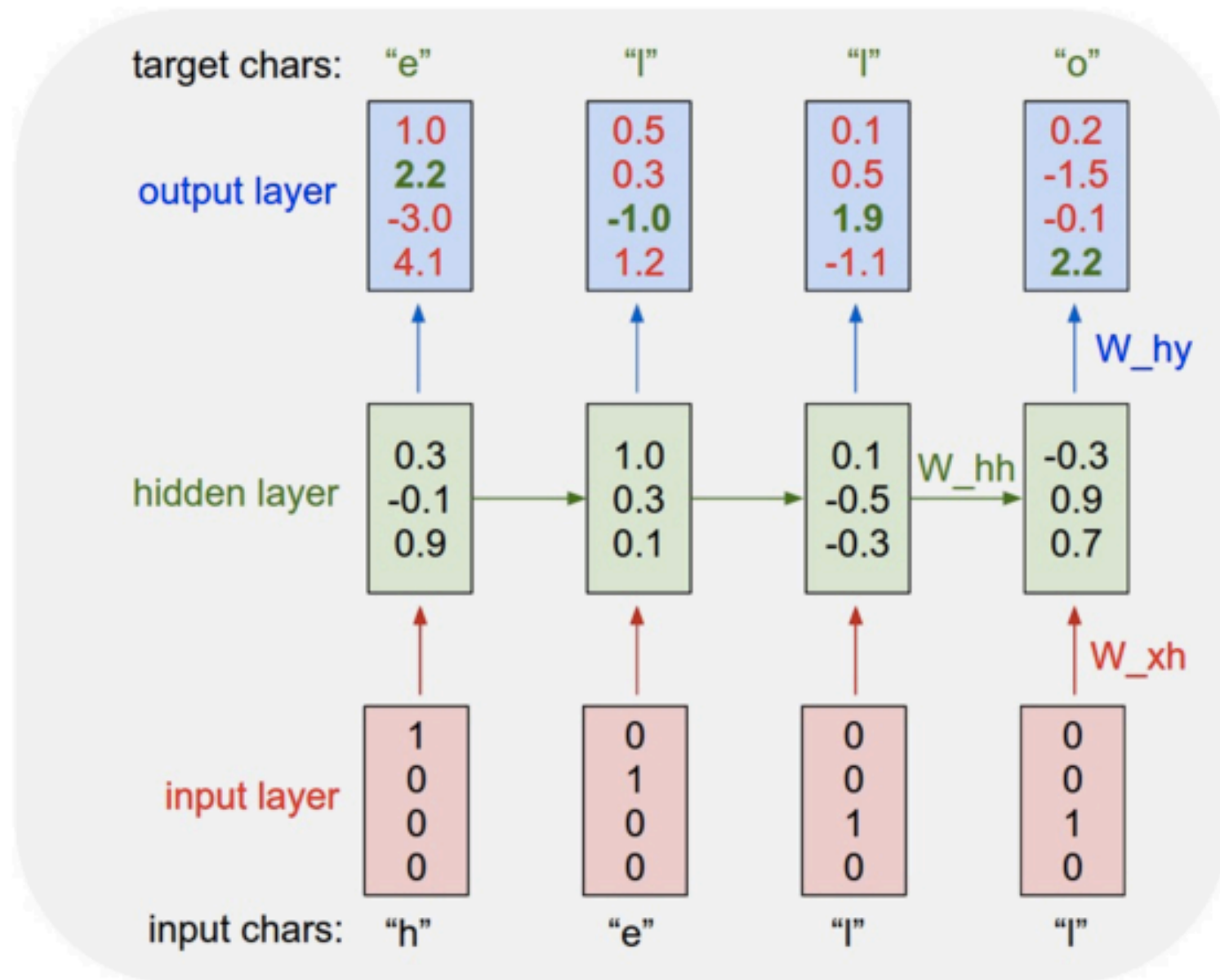
$h_0 \in \mathbb{R}^{D_h}$ is some initialization vector for the hidden layer at time step 0

$x_{[t]}$ is the column vector of L at index [t] at time step t

$$W^{(hh)} \in \mathbb{R}^{D_h \times D_h} \qquad W^{(hx)} \in \mathbb{R}^{D_h \times d} \qquad W^{(S)} \in \mathbb{R}^{|V| \times D_h}$$

# Training RNN language models

- It's just a neural network!
  - With a slightly different / more complicated structure
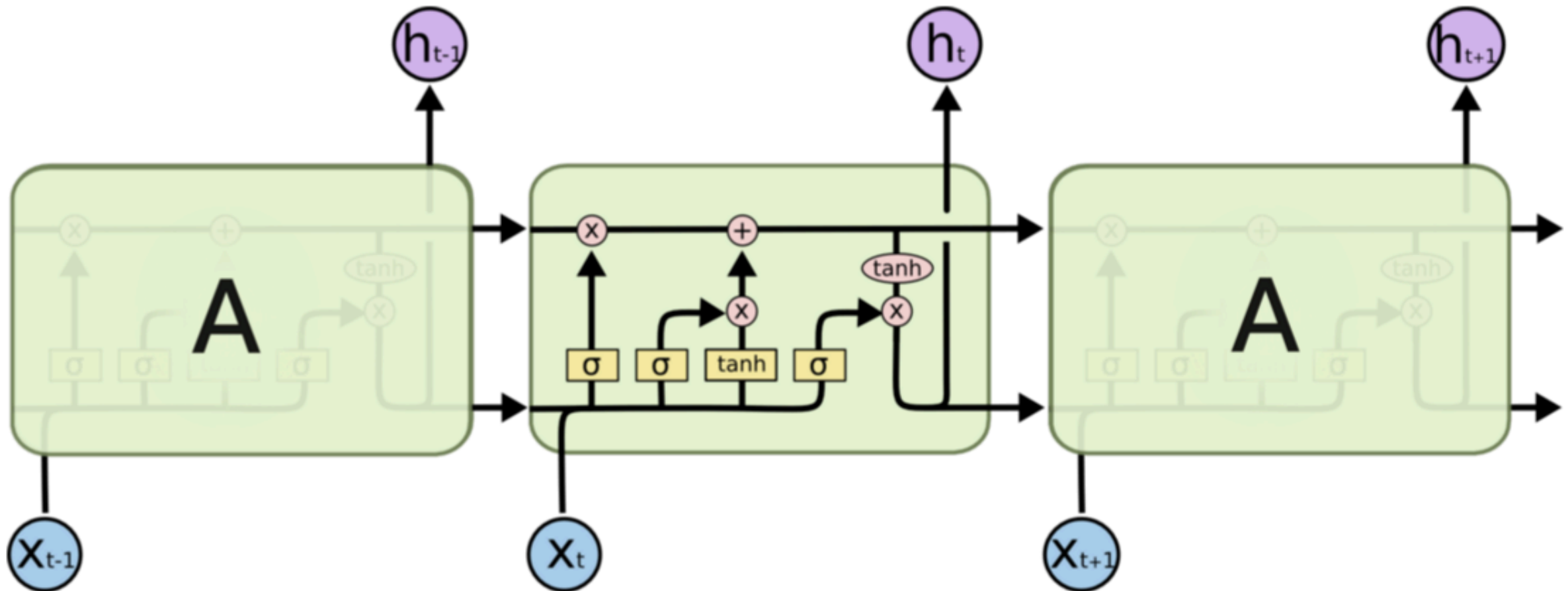  - Straightforward to compute gradients of all parameters wrt. Outputs
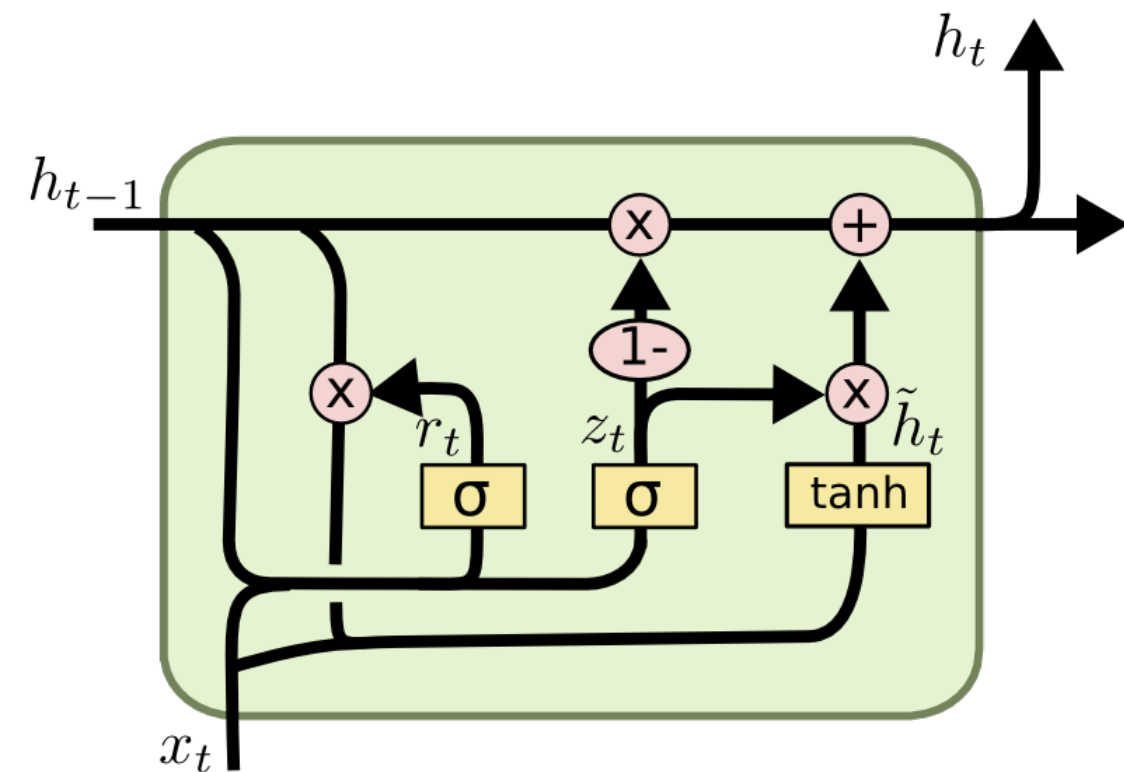
# More complicated stuff…

- Technical Problem: vanishing / exploding gradients
  - (sort of) Solutions: LSTM units or GRUs
- Encoder-decoder RNN
- Sequence-to-sequence
  - Machine translation
  - Conversation generation
- Image captioning

These are all just end-to-end neural networks trained with backpropagation!
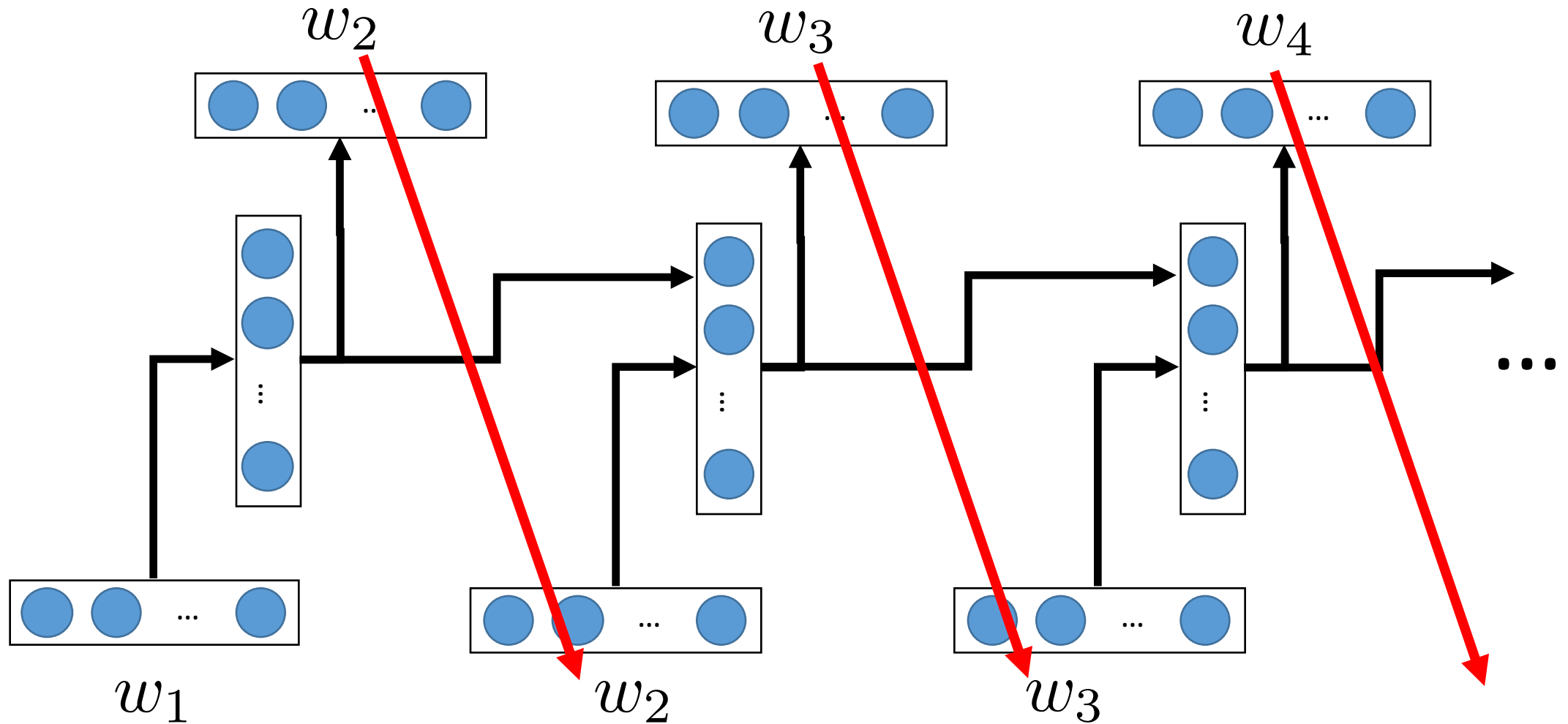
# LSTM Unit...

# GRU Unit...



$$z_t = \sigma\left(W_z \cdot [h_{t-1}, x_t]\right)$$

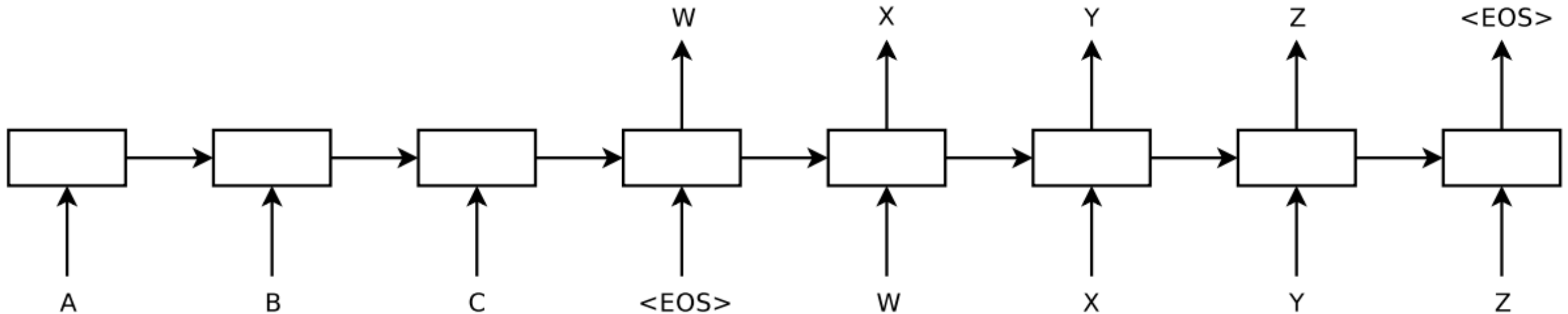$$r_t = \sigma\left(W_r \cdot [h_{t-1}, x_t]\right)$$

$$\tilde{h}_t = \tanh\left(W \cdot [r_t * h_{t-1}, x_t]\right)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Generating text from an RNN LM

# Neural Translation Models (sequence to sequence)

# Conversation Generation

**Human**: *what is the purpose of life?*
**Machine**: *to serve the greater good.*
**Human**: *what is the purpose of living?*
**Machine**: *to live forever.*
**Human**: *what is the purpose of existence?*
**Machine**: *to find out what happens when we get to the planet earth.*
**Human**: *where are you now?*
**Machine**: *i'm in the middle of nowhere.*
**Human**: *what is the purpose of dying?*
**Machine**: *to have a life.*
**Human**: *what is the purpose of being intelligent?*
**Machine**: *to find out what it is.*
**Human**: *what is the purpose of emotions?*
**Machine**: *i don't know.*
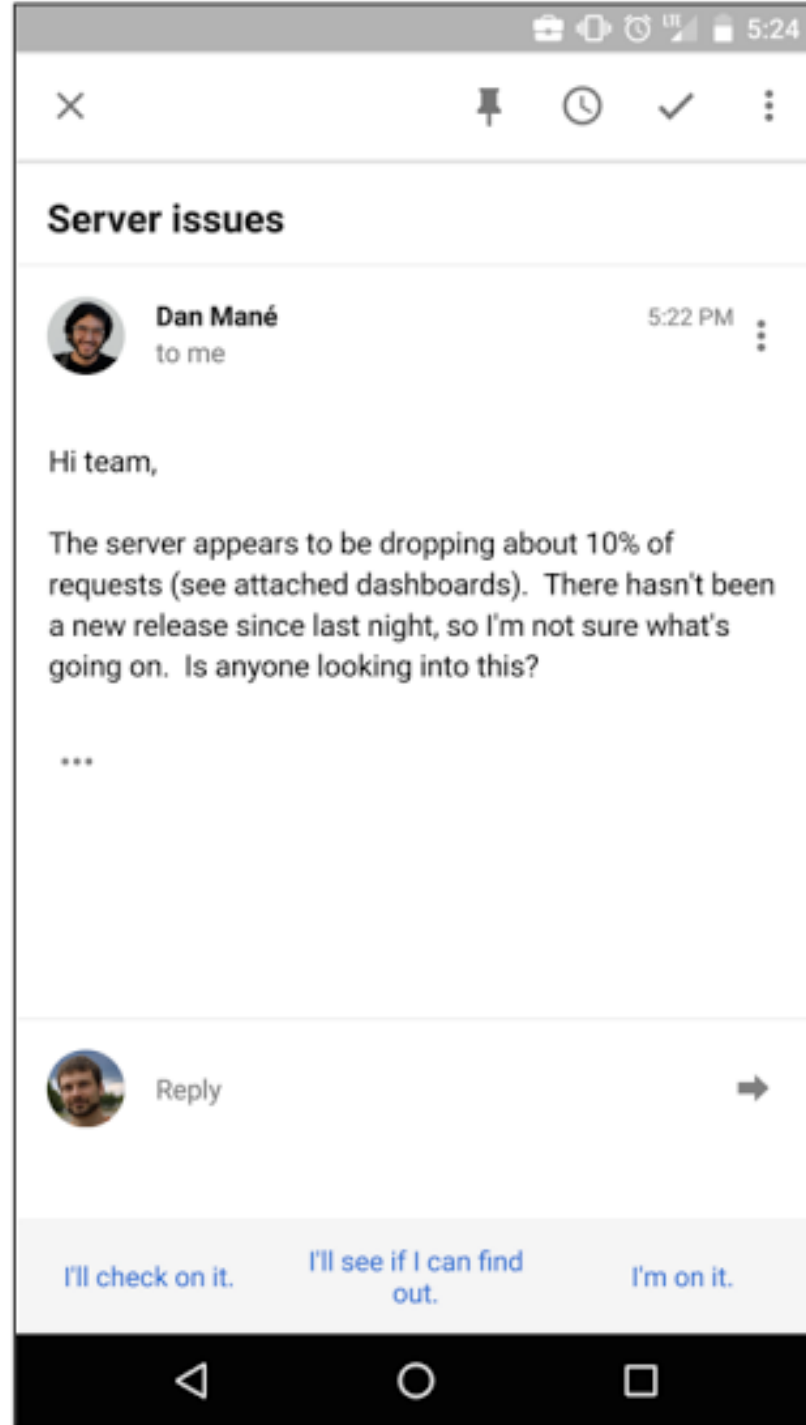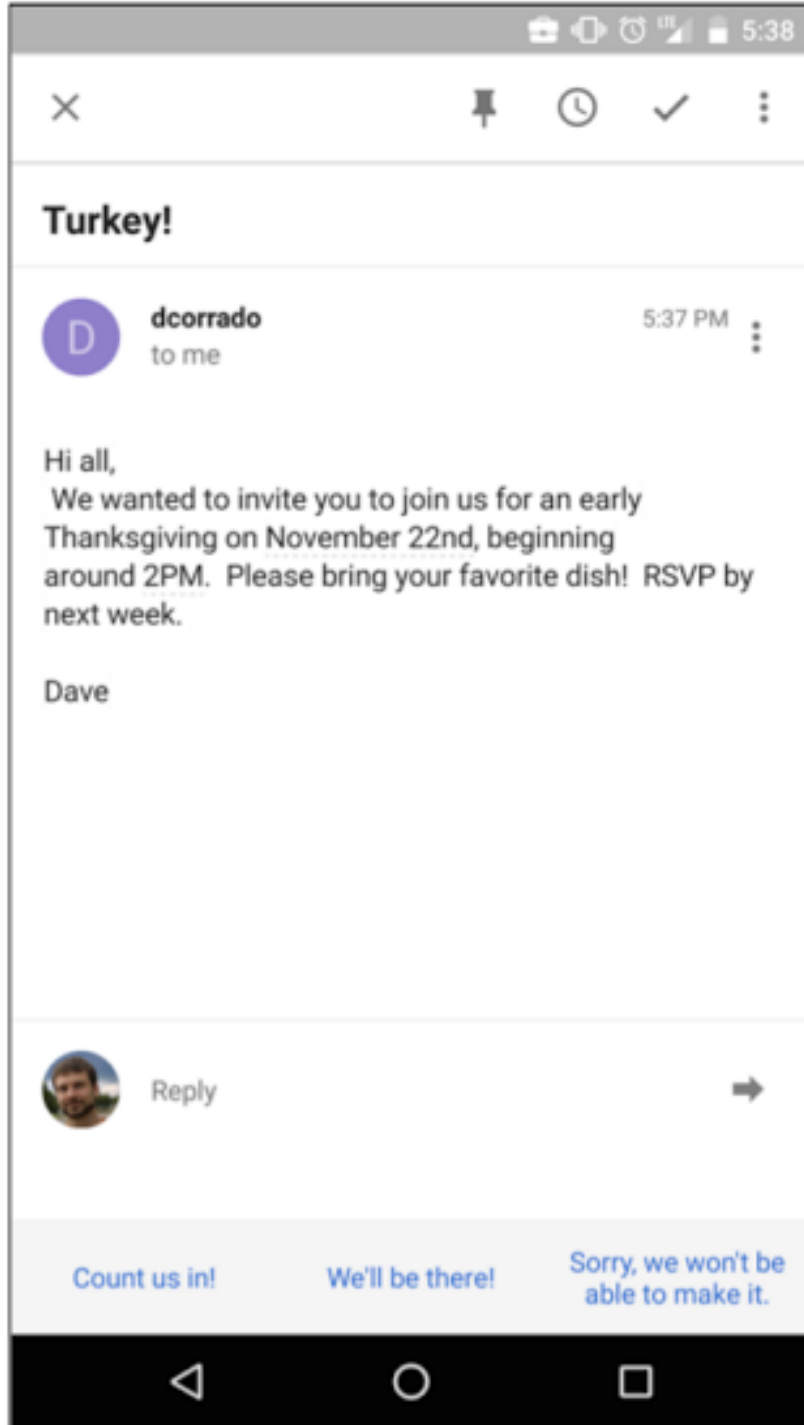
# Google Research Blog

## Computer, respond to this email.
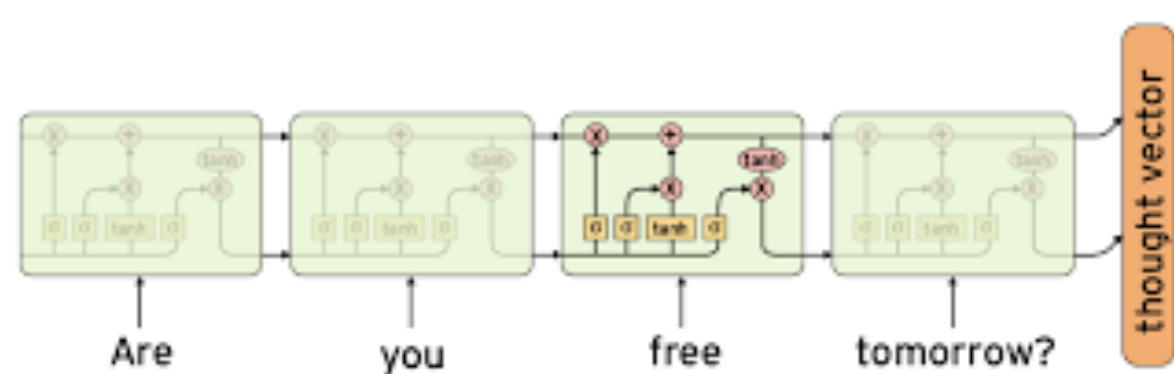
Tuesday, November 03, 2015

Posted by Greg Corrado*, Senior Research Scientist

**Machine Intelligence for You**

What I love about working at Google is the opportunity to harness cutting-edge machine intelligence for users' benefit. Two recent Research Blog posts talked about how we've used machine learning in the form of deep neural networks to improve voice search and YouTube thumbnails. Today we can share something even wilder -- Smart Reply, a deep neural network that writes email.

ENCODER

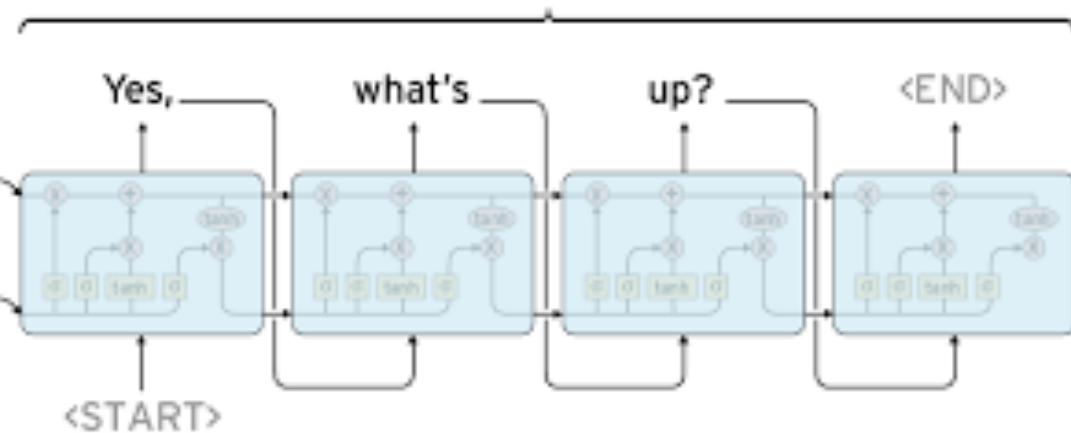Reply

thought vector

Yes,     what's     up?     <END>

Are     you     free     tomorrow?

<START>

Incoming Email

DECODER

# Show and Tell: A Neural Image Caption Generator

Oriol Vinyals
Google
vinyals@google.com

Alexander Toshev
Google
toshev@google.com

Samy Bengio
Google
bengio@google.com
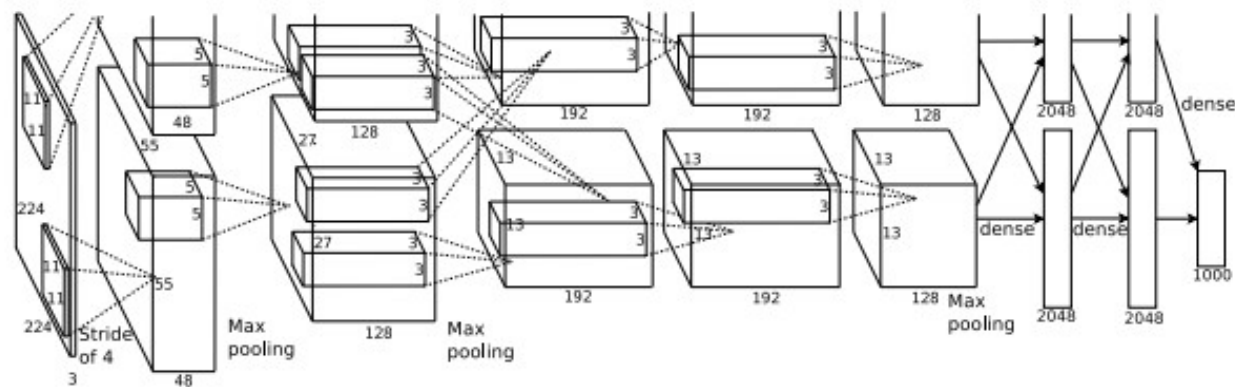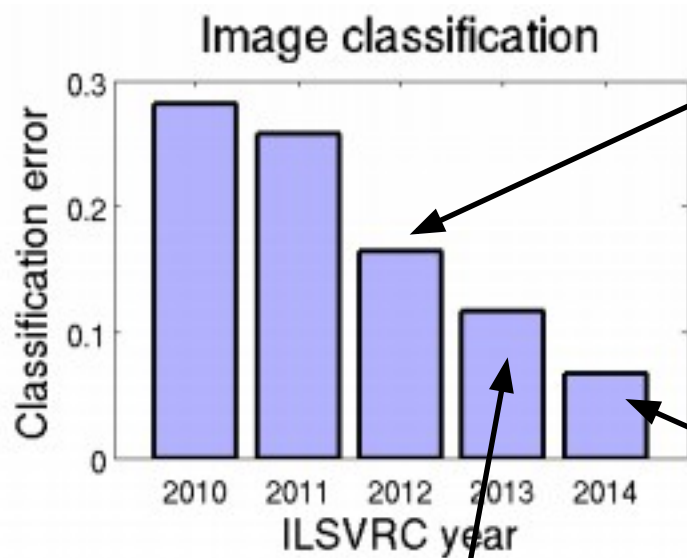
Dumitru Erhan
Google
dumitru@google.com

[Krizhevsky, Sutskever, Hinton. 2012] **16.4% error**

Image classification

[Szegedy et al., 2014] **6.6% error**
[Simonyan and Zisserman, 2014] **7.3% error**

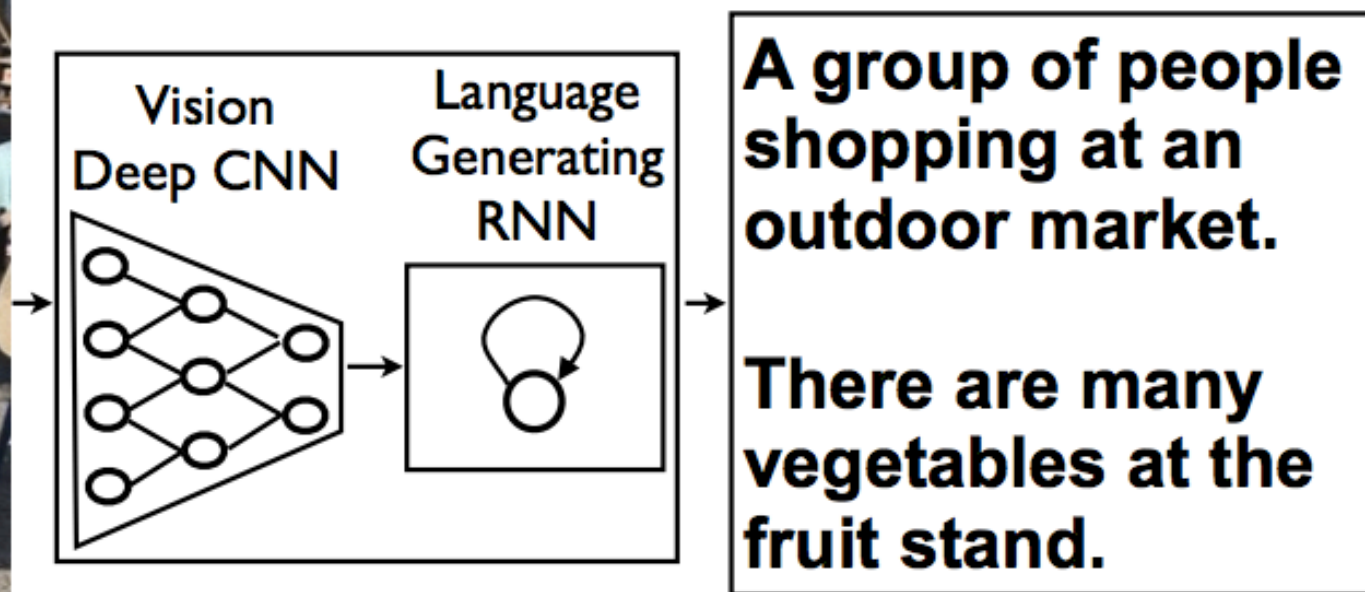[Zeiler and Fergus, 2013] **11.1% error**

Figure 1. NIC, our model, is based end-to-end on a neural network consisting of a vision CNN followed by a language generating RNN. It generates complete sentences in natural language from an input image, as shown on the example above.

# Image Sentence Datasets

a man riding a bike on a dirt path through a forest.
bicyclist raises his fist as he rides on desert dirt trail.
this dirt bike rider is smiling and raising his fist in triumph.
a man riding a bicycle while pumping his fist in the air.
a mountain biker pumps his fist in celebration.



Microsoft COCO
*[Tsung-Yi Lin et al. 2014]*
mscoco.org


currently:
~120K images
~5 sentences each

# + Transfer Learning

use weights
pretrained from
ImageNet

| image |
| conv-64 |
| conv-64 |
| maxpool |
| conv-128 |
| conv-128 |
| maxpool |
| conv-256 |
| conv-256 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| FC-4096 |
| FC-4096 |



"straw hat"

training example

y0   y1   y2

h0 → h1 → h2

x0
<START>   x1 "straw"   x2 "hat"

<START>   straw   hat

use word vectors
pretrained with
word2vec [1]

*[1] Mikolov et al., 2013*

# Summary of the approach

We wanted to describe images with sentences.

1. Define a single function from input -> output
2. Initialize parts of net from elsewhere if possible
3. Get some data
4. Train with SGD

# Wow I can't believe that worked



a group of people standing around a room with remotes
logprob: -9.17

a young boy is holding a baseball bat
logprob: -7.61

a cow is standing in the middle of a street
logprob: -8.84

# Wow I can't believe that worked



a cat is sitting on a toilet seat
logprob: -7.79



a display case filled with lots of different types of donuts
logprob: -7.78



a group of people sitting at a table with wine glasses
logprob: -6.71

# Well, I can kind of see it



a man standing next to a clock on a wall
logprob: -10.08



a young boy is holding a
baseball bat
logprob: -7.65



a cat is sitting on a couch with a remote control
logprob: -12.45

# Training an RNN/LSTM...

- **Clip** the gradients (important!). 5 worked ok
- **RMSprop** adaptive learning rate worked nice
- Initialize softmax **biases** with log word frequency distribution
- Train for **long time**

# Summary

- Deep learning is a popular area in machine learning recently
  - Very successful in speech recognition and computer vision
- Becoming very popular in NLP these days
- Main motivation:
  - Learn feature representations from data
  - Alternative to hand-engineered features
- Neural networks:
  - Primary deep learning approach
  - Layers of logistic regressions – can directly calculate gradients from outputs
  - Nonlinear decision boundaries