

TOKENIZATION

Intro to neural networks

The need for segmentation

- Text as a stream of characters



- We need a way to understand the meaning of text

- Break into words (assign meaning to word) ← Tokenization
- Break into sentences (put word meanings back to sentence meaning)

sometimes
text is too long so
we divide it into chunk ↑
before doing NLP task we should
start from cutting the text into
chunk,
sentences,
words.
(cenglish)
Tokenizer .

Tokenization

→ to make into tokens
→ Segmentation. ayırma
aka ↴
1 token.

- A token should
 - 1. Linguistically significant (meaning)
 - 2. Methodologically useful (using with model)

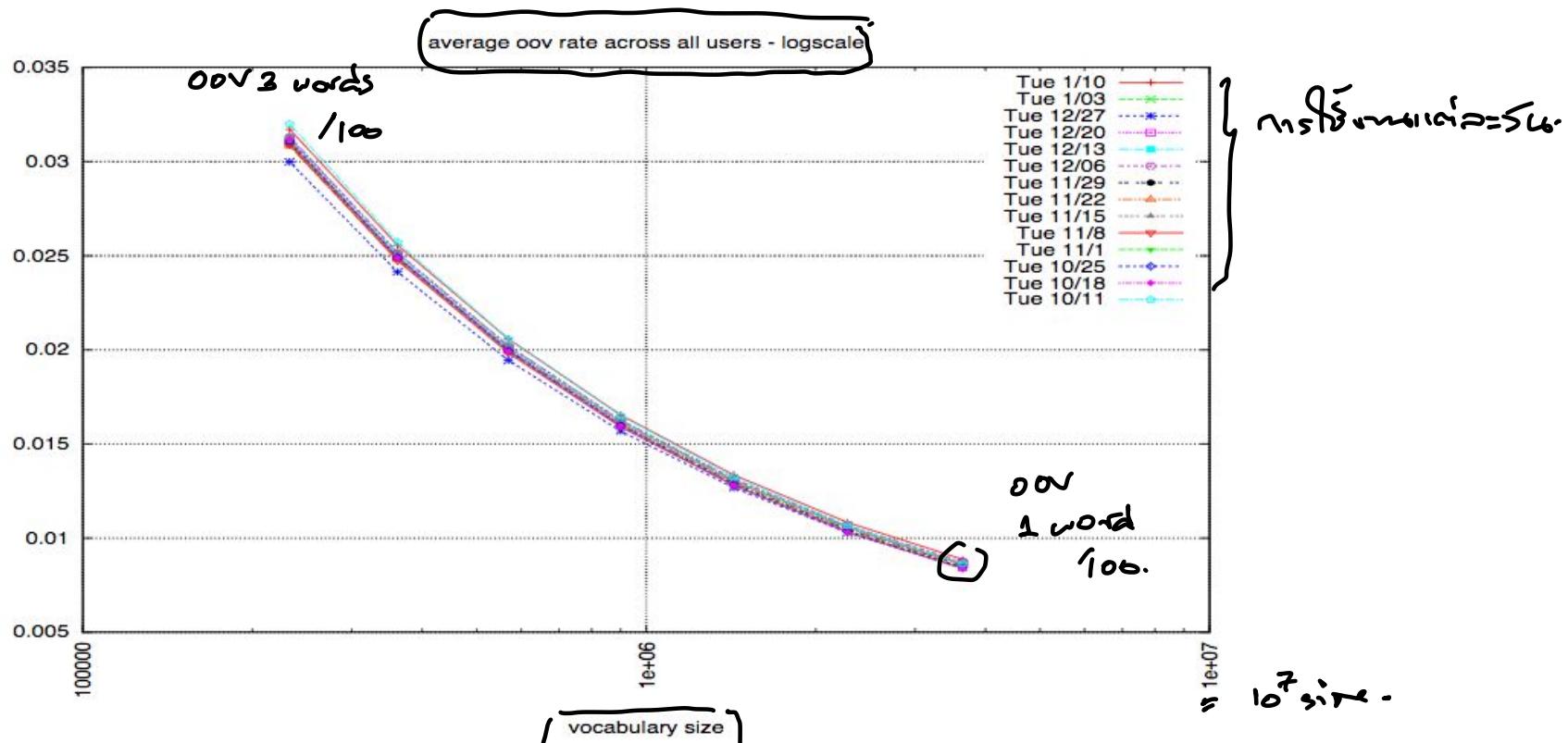
Dictionary-based vs Machine-learning-based

- Dictionary-based (last time)
 - Longest matching
 - Maximal matching
- Machine-learning-based
 - ** Today topic.*
 - *↳ training data.*

why we use ML-based

Dictionary-based drawbacks

- Cannot handle words outside of the dictionary (Out-of-Vocabulary, OOV words)
- Performs worse than machine-learning-based approach



NEURAL NETWORKS

Deep learning = Deep neural networks =
neural networks

DNNs (Deep Neural Networks)

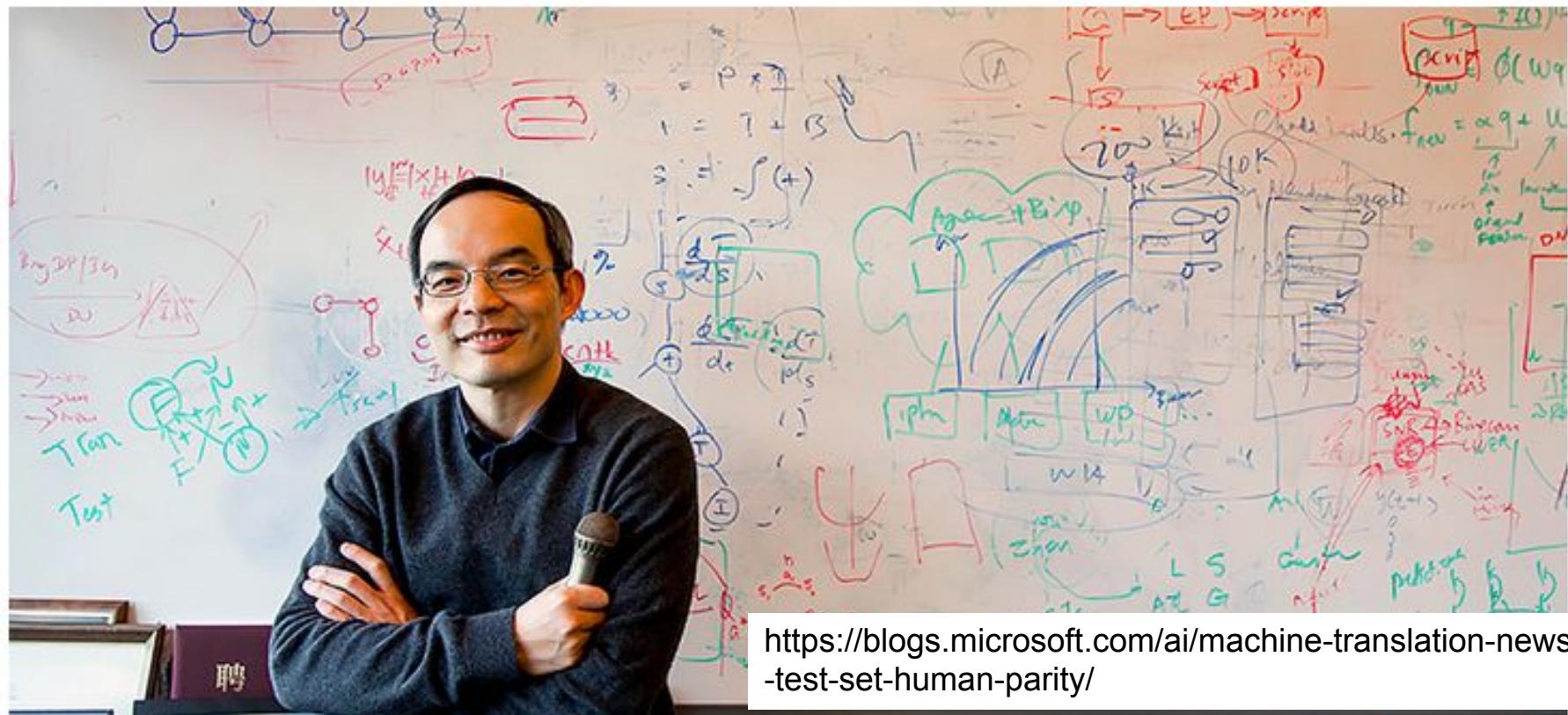
- Why deep learning?
- Greatly improved performance in ASR and other tasks (Computer Vision, Robotics, Machine Translation, NLP, etc.)
- Surpassed human performance in many tasks

Microsoft reaches a historic milestone, using AI to match “human performance” in translating news from Chinese to English

Mar 14, 2018 | Allison Linn



Hinge size model than gg.trans
and have better performance.



Deep learning in NLP

Easy task modest gains

	Traditional ML	Deep learning
Sentiment (th)	72%	$-4\% \rightarrow$ 76%
Topic classification (th)	67%	$-3\% \rightarrow$ 70%
PoS (th)	96%	$-1\% \rightarrow$ 97%

Harder task larger gains

	Traditional ML	Deep learning
QA	51%*	$-39\% \rightarrow$ 90%
Creating image from text	??? (<i>cannot do</i>)	<u>very good</u>



wangchanbert

<https://www.aclweb.org/anthology/>

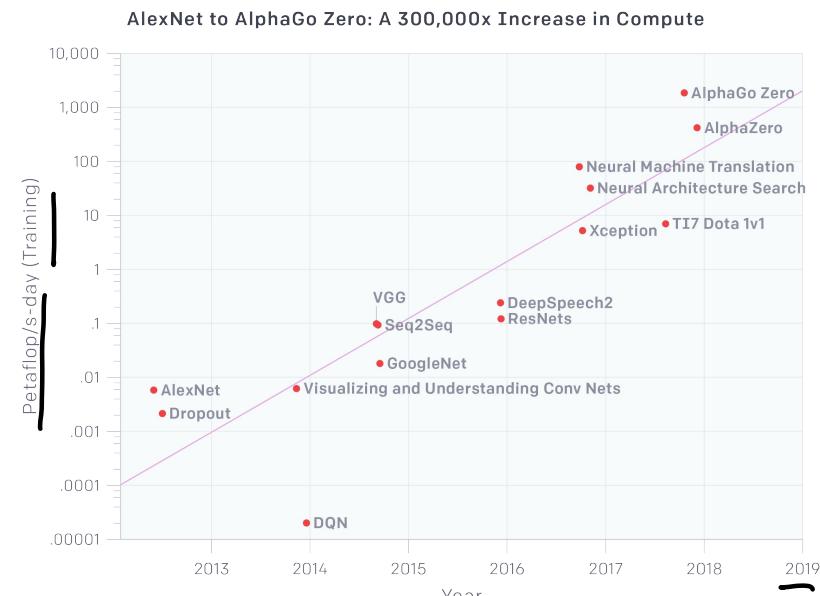
D16-1264/

<https://openai.com/blog/dall-e/>

Why now

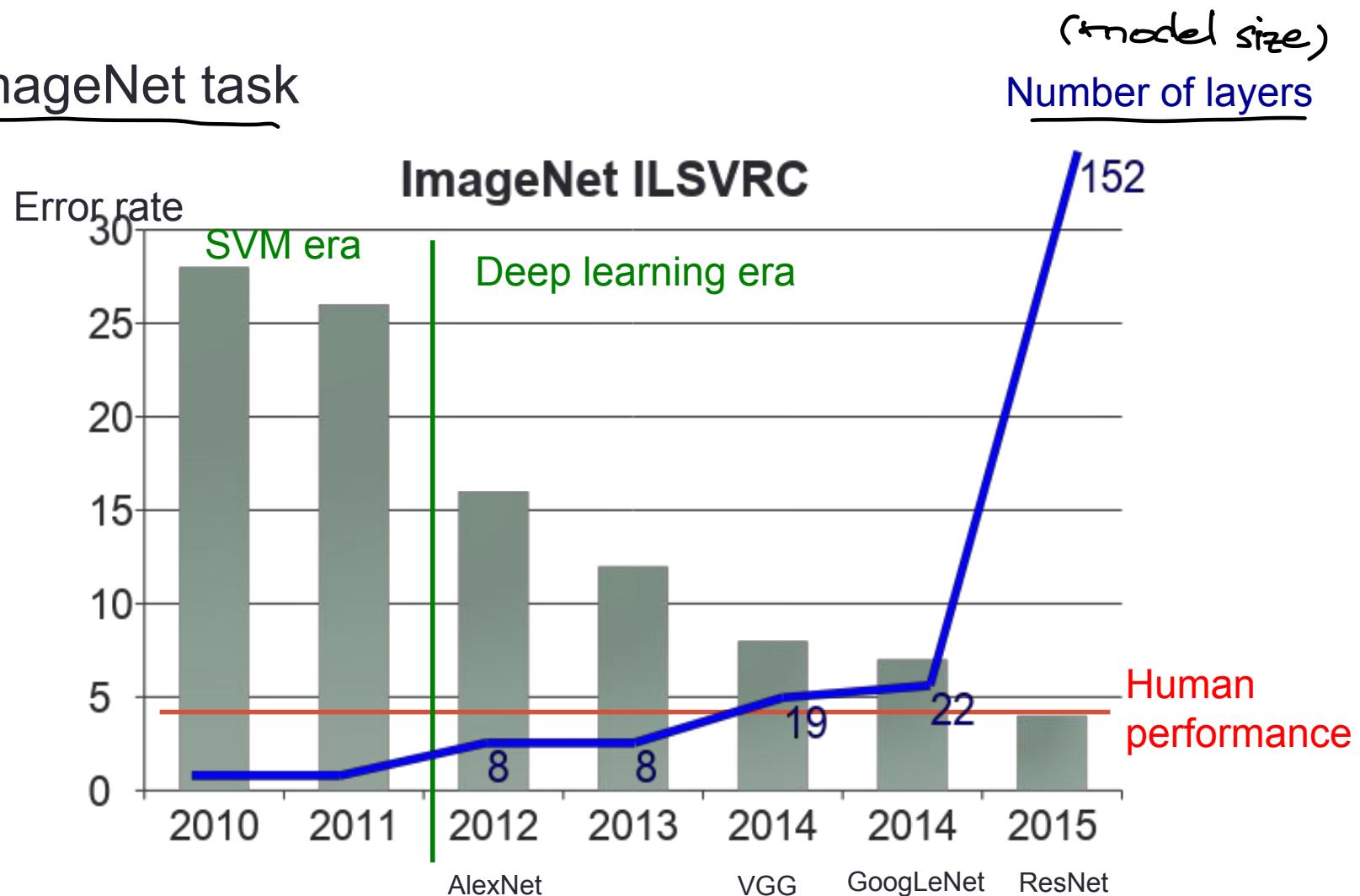
- Neural Networks has been around since 1990s
 - **Big data** – DNN can take advantage of large amounts of data better than other models
 - **GPU** – Enable training bigger models possible
 - **Deep** – Easier to avoid bad local minima when the model is large

\downarrow Petaflop/s-day
Tflops \downarrow Petaflop..
 \times day to train.

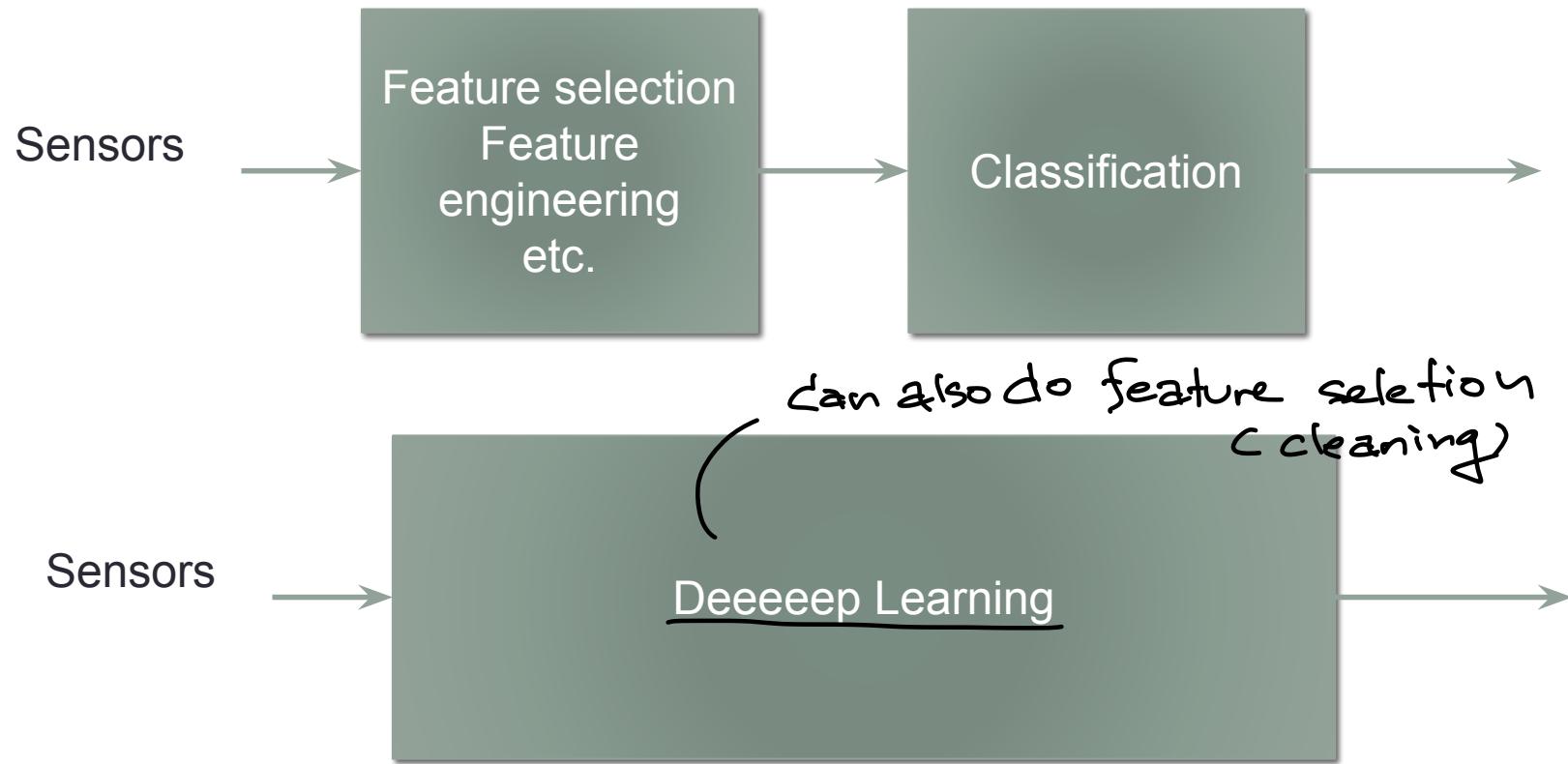


Wider and deeper networks

- ImageNet task



Traditional VS Deep learning

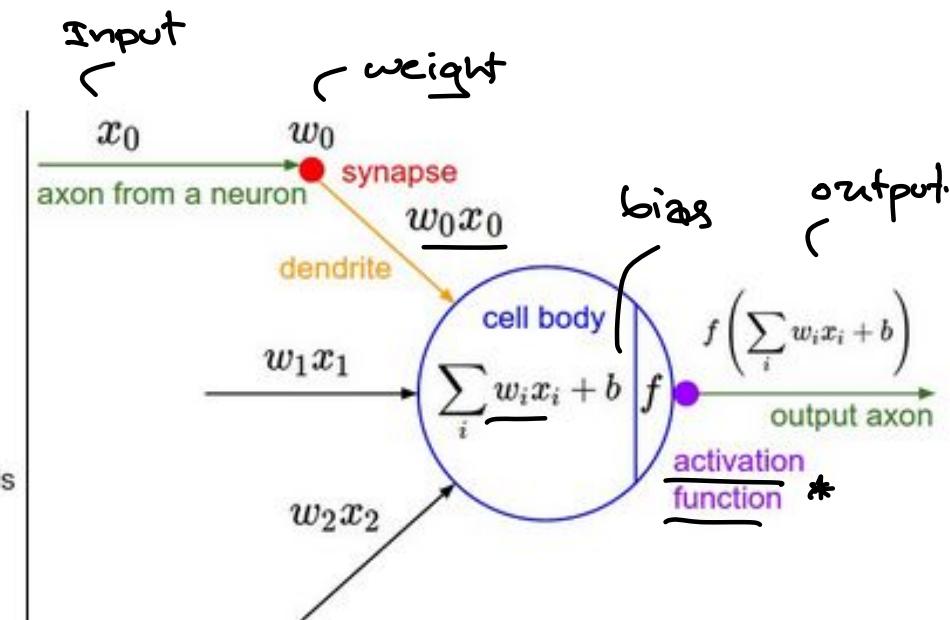
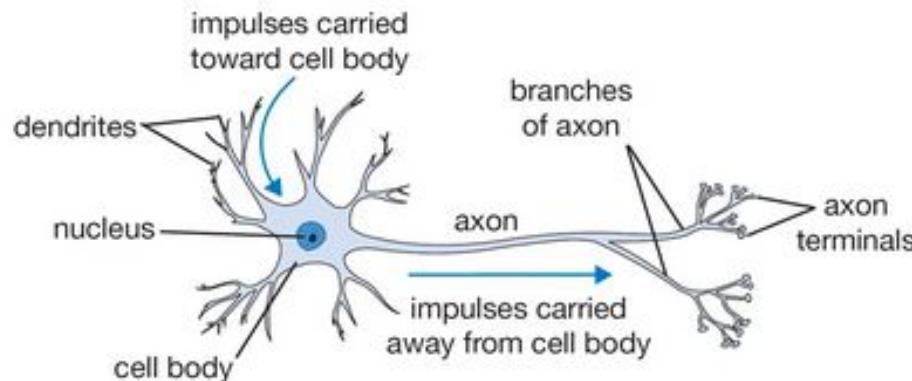


Neural networks

- Fully connected networks
 - Neuron
 - Non-linearity
 - Softmax layer
- DNN training
 - Loss function
 - SGD and backprop
 - Learning rate
 - Overfitting
- CNN, RNN, LSTM, GRU

Fully connected networks

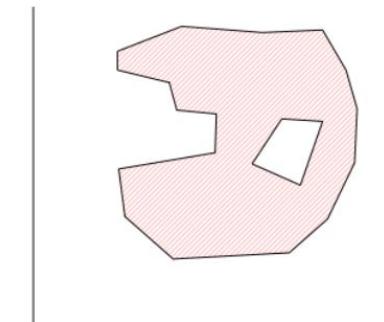
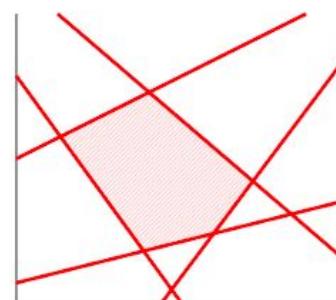
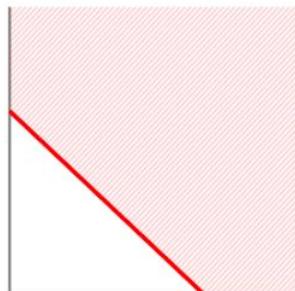
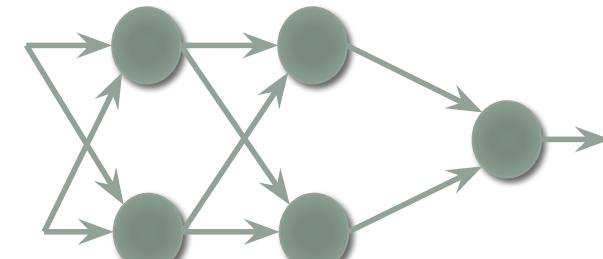
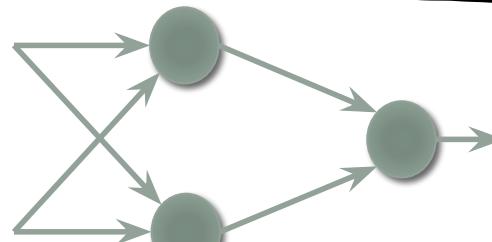
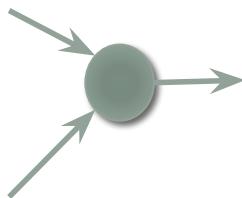
- Many names: feed forward networks or deep neural networks or multilayer perceptron or artificial neural networks
- Composed of multiple 'neurons'



Combining neurons

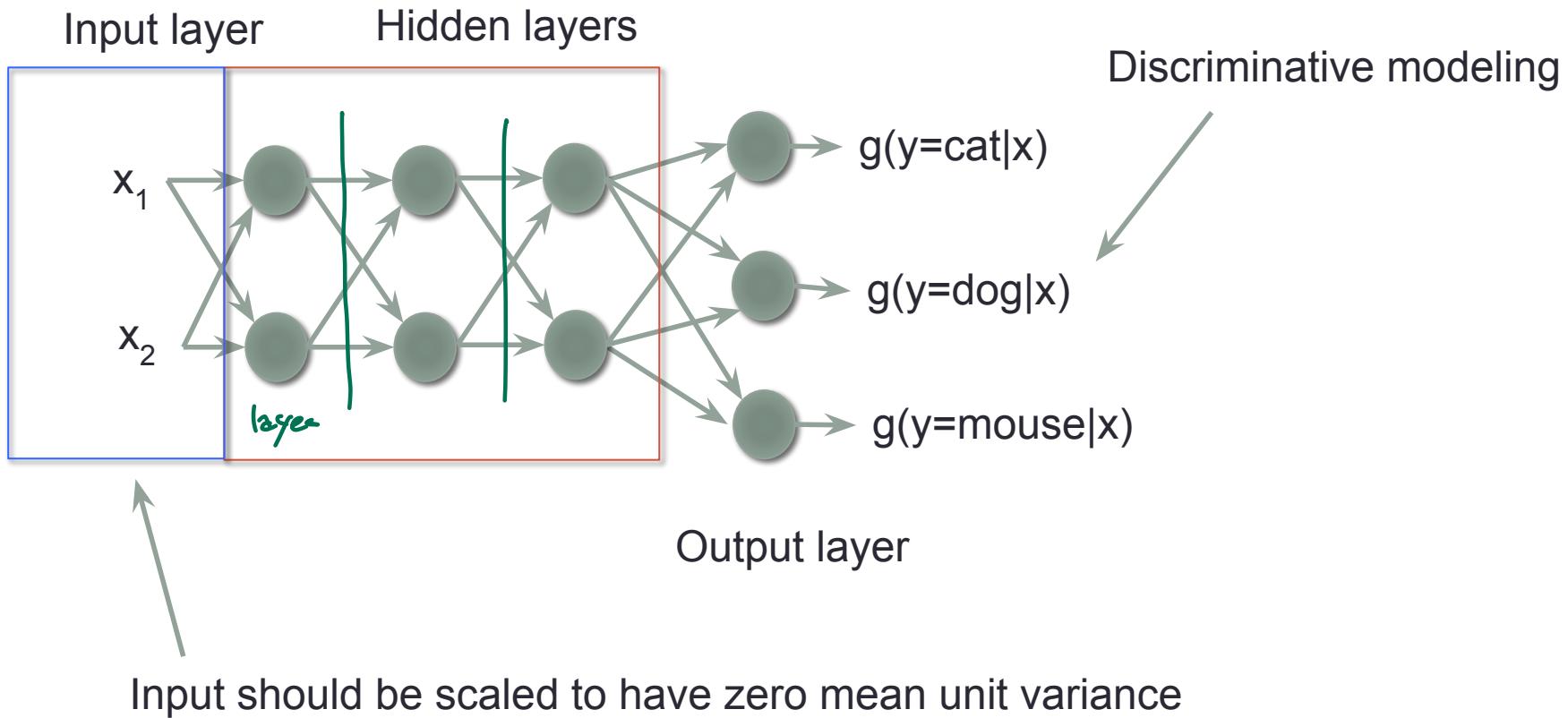
→ make model more complex.

- Each neuron splits the feature space with a hyperplane
- Stacking neuron creates more complicated decision boundaries
- More powerful but prone to overfitting



Terminology

Deep in Deep neural networks means many hidden layers

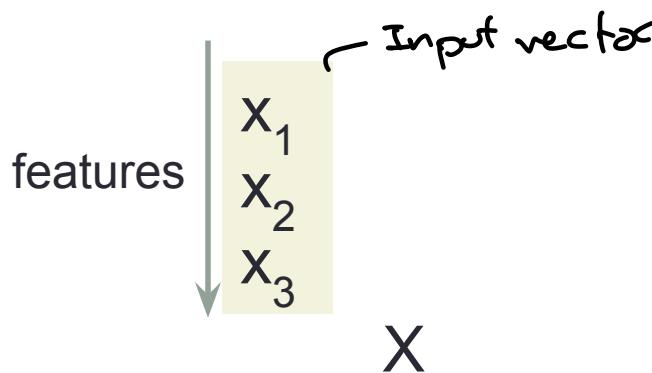


Neural Network as

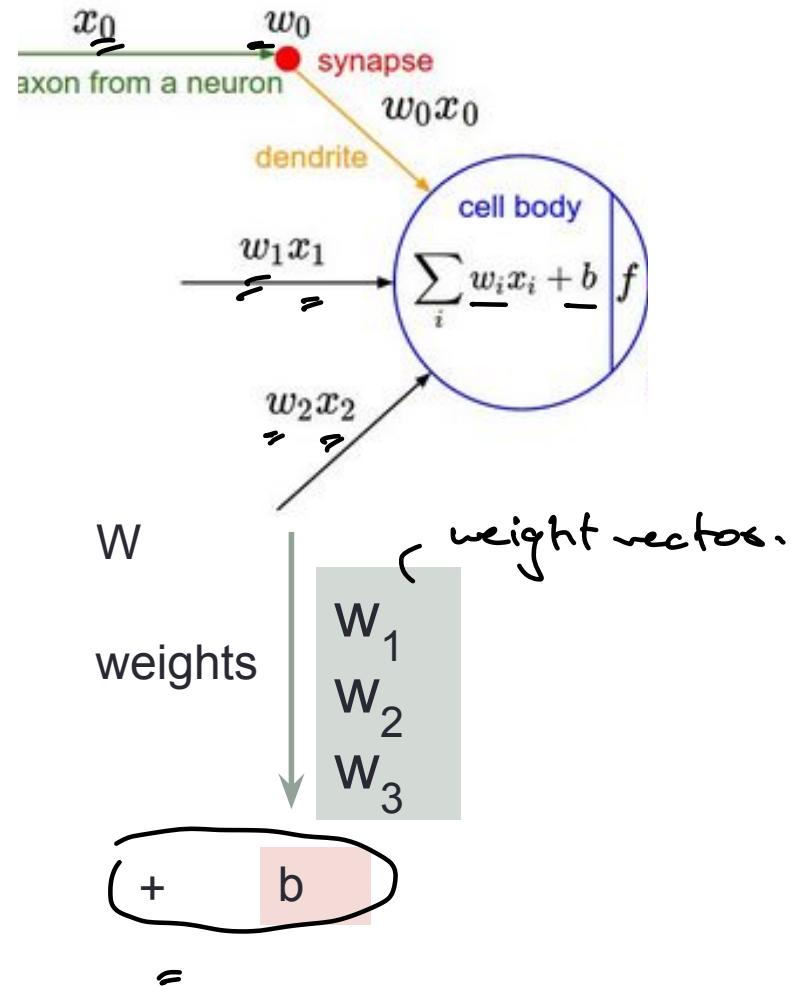
Matrices

$$\text{vector} = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}$$

- Inputs



$$W^T X + b = \boxed{W_1 \ W_2 \ W_3} \quad \boxed{x_1 \ x_2 \ x_3}$$



Matrices

This is why we use "jpn"

- Inputs

samples				
features	x_{11}	x_{12}	x_{13}	x_{14}
	x_{21}	x_{22}	x_{23}	x_{24}
	x_{31}	x_{32}	x_{33}	x_{34}

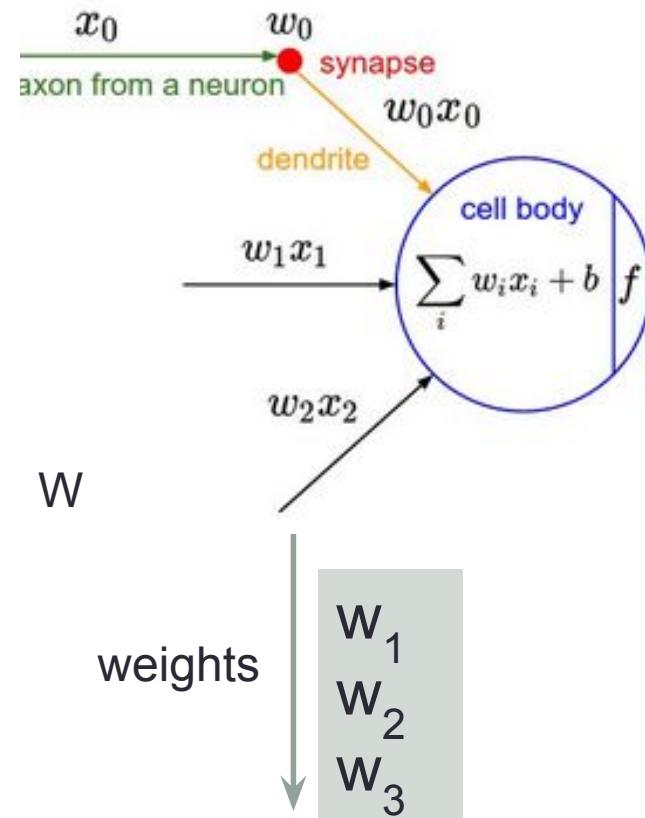
$$W^T X + b$$

$$\boxed{w_1 \ w_2 \ w_3}$$

$$\begin{pmatrix} x_{11} \\ x_{21} \\ x_{31} \end{pmatrix} \begin{pmatrix} x_{12} \\ x_{22} \\ x_{32} \end{pmatrix} \begin{pmatrix} x_{13} \\ x_{23} \\ x_{33} \end{pmatrix} \begin{pmatrix} x_{14} \\ x_{24} \\ x_{34} \end{pmatrix}$$

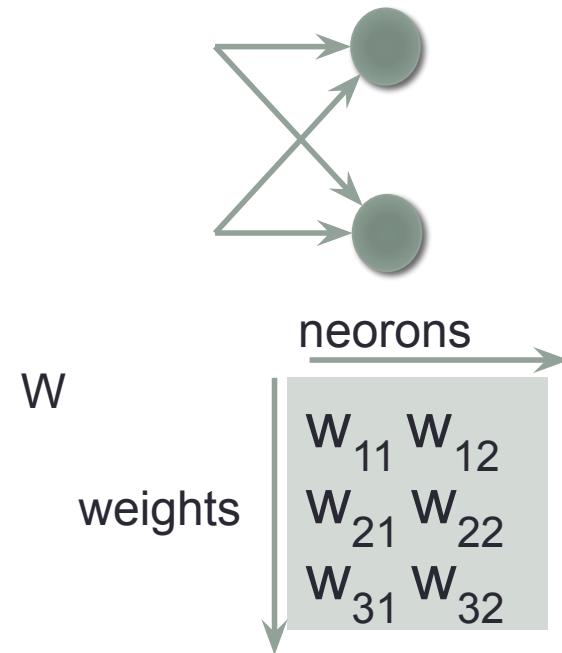
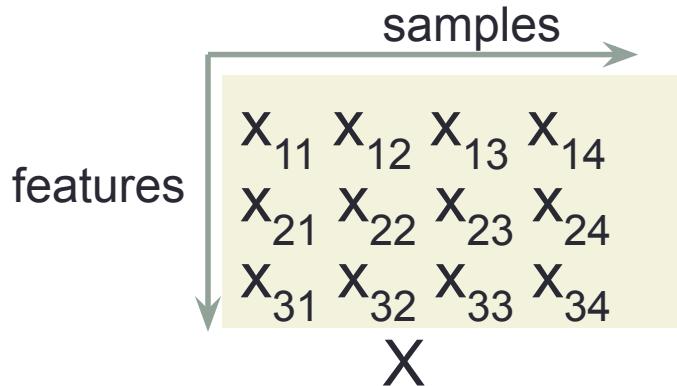
+

$$= \boxed{b}$$



Matrices

- Inputs



$$W^T X + b$$

$$\begin{matrix} w_{11} & w_{21} & w_{31} \\ w_{21} & w_{22} & w_{23} \end{matrix}$$

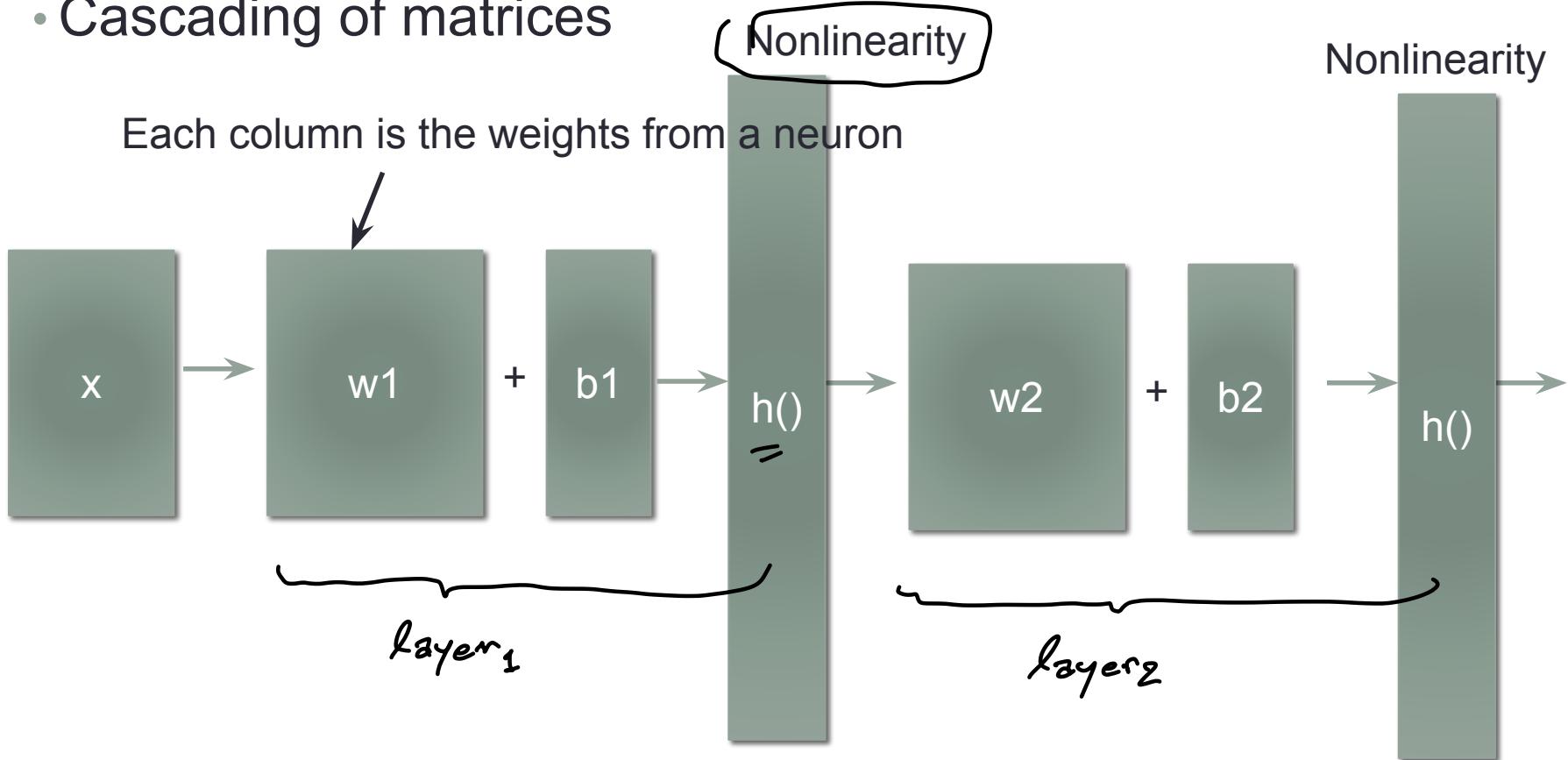
$$\begin{matrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{matrix}$$

+

$$\begin{matrix} b_1 \\ b_2 \end{matrix}$$

More linear algebra

- Cascading of matrices

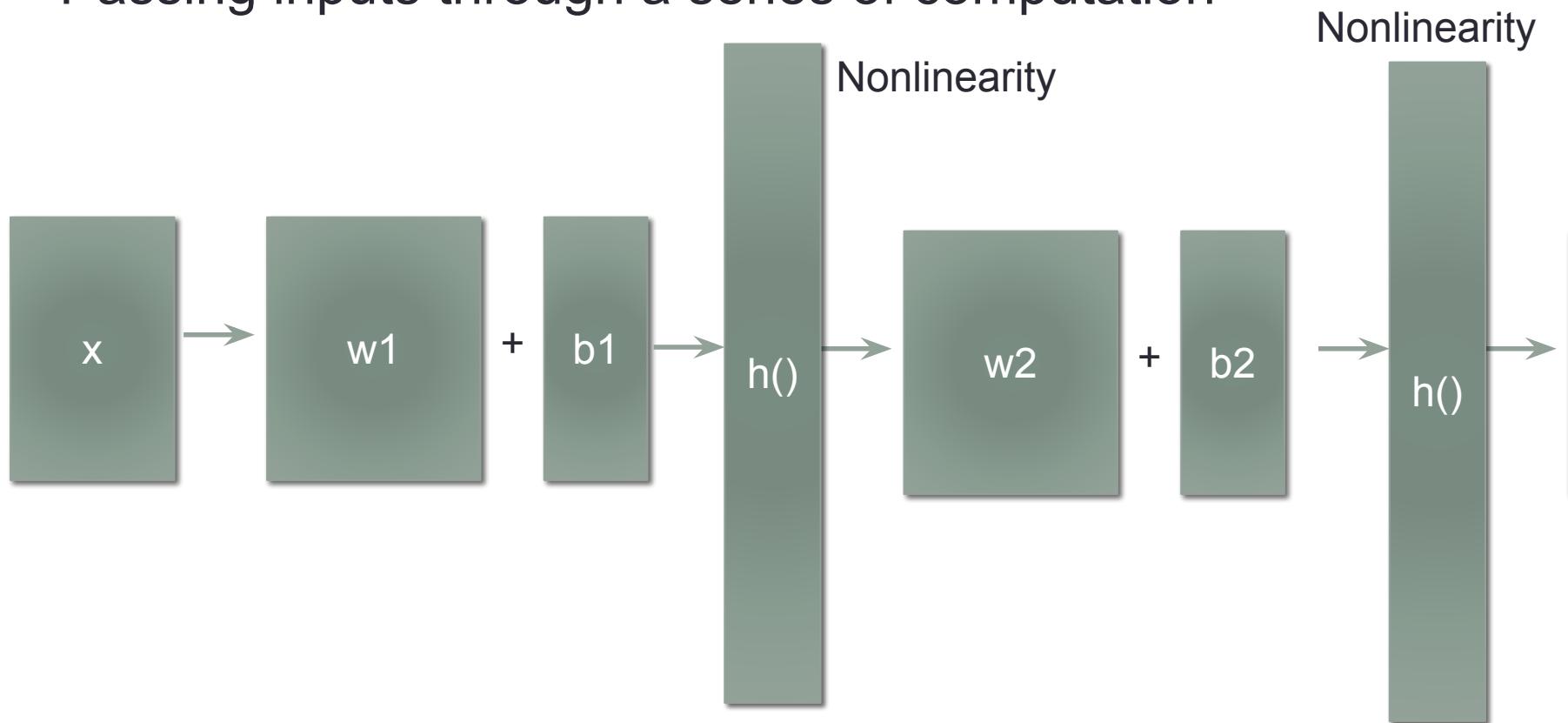


$$h(W_2^T \underbrace{h(W_1^T X + b_1)}_{\text{---}} + b_2)$$

Computation graph

visualize network as picture below.

- Passing inputs through a series of computation



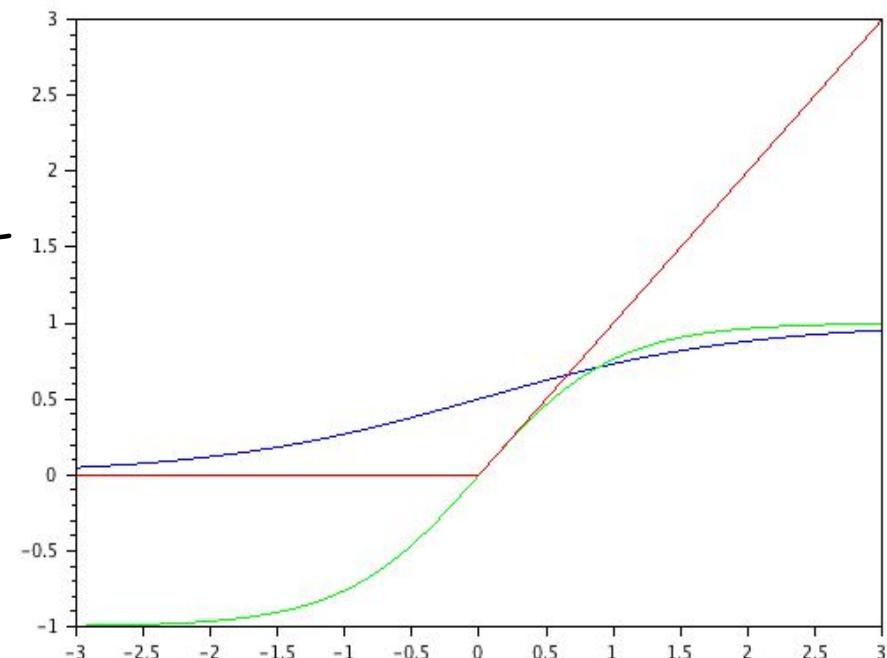
$$h(W_2^T h(W_1^T X + \mathbf{b}_1) + \mathbf{b}_2)$$

Non-linearity

if we use linear function. we can ignore MatrixA.

$w_2^T A w_1^T \rightarrow b$ \hookrightarrow non-linear 1 matrix
(layer 2 has no neurons)

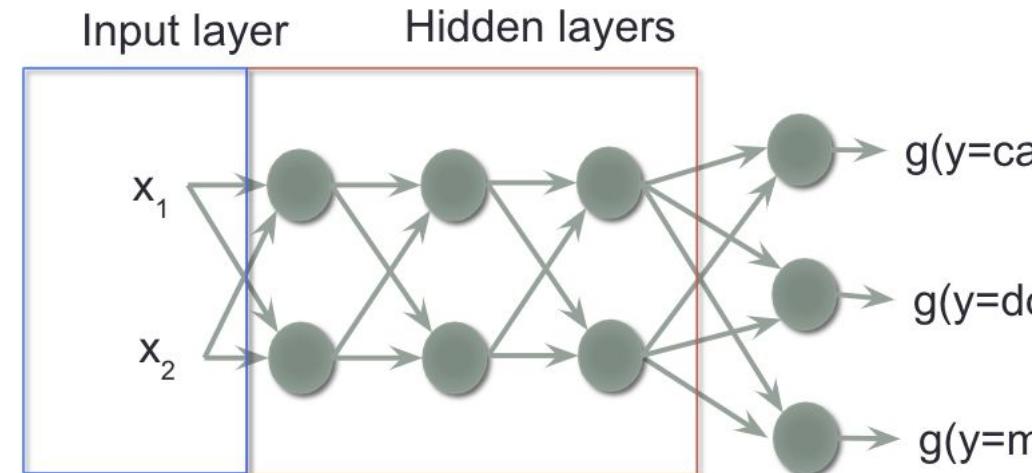
- The Non-linearity is important in order to stack neurons
 - If F is linear, a multi layered network can be collapsed as a single layer (by just multiplying weights together)
- Sigmoid or logistic function
- tanh
- Rectified Linear Unit (ReLU)
- Most popular is ReLU and its variants (Fast to train, and more stable)



Output layer – “Softmax” layer

- We usually want the output to mimic a probability function ($0 \leq P \leq 1$, sums to 1)
- Current setup has no such constraint
- The current output should have highest value for the correct class.
 - Value can be positive or negative number
- Takes the exponent
- Add a normalization

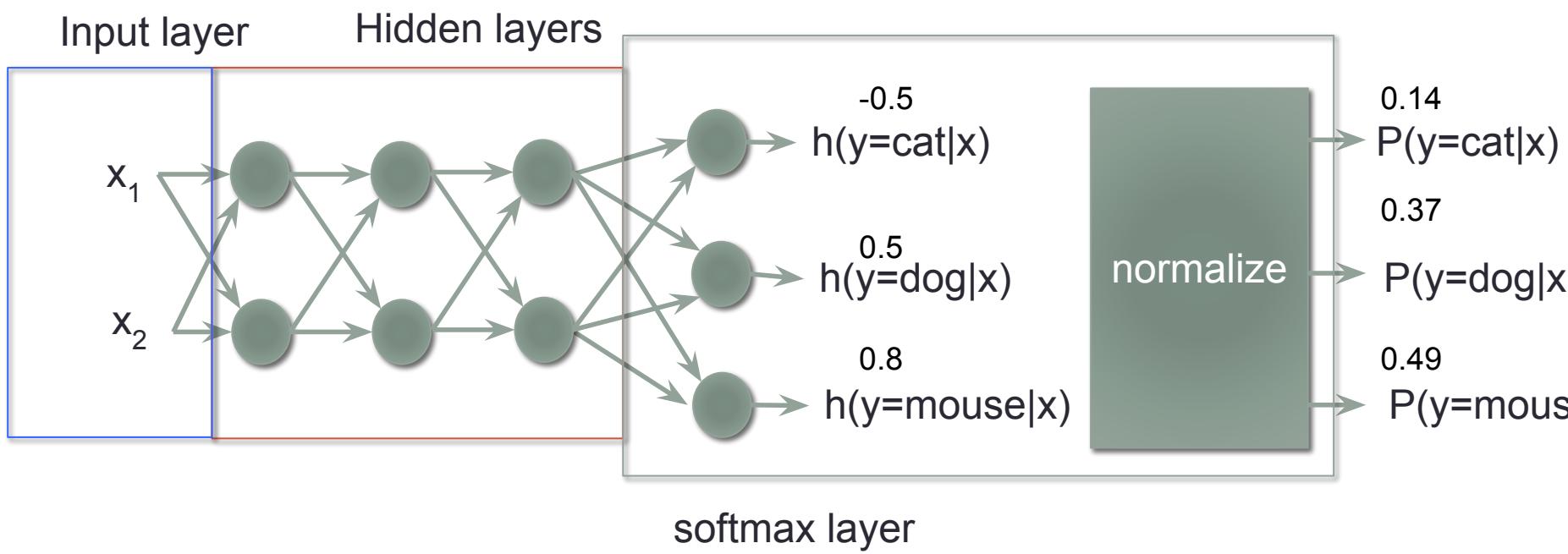
ກົດຕົວມາແລ້ວ
ສອບ.



Softmax layer

→ we want $y=j$ (whatever) to be prob.

$$P(y = j|x) = \frac{e^{h(y=j|x)}}{\sum_y e^{h(y|x)}} - \text{exponent } (\cancel{+,- \rightarrow > 0}) - \text{sum of them} \rightarrow \oplus \text{ must } \leq 1.$$

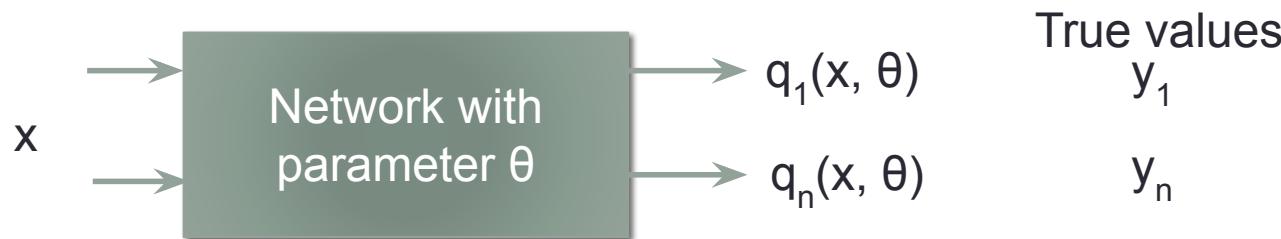


Neural networks

- Fully connected networks
 - Neuron
 - Non-linearity
 - Softmax layer
- DNN training
 - Loss function
 - SGD and backprop
 - Learning rate
 - Overfitting
- CNN, RNN, LSTM, GRU

Objective function (Loss function)

- Can be any function that summarizes the performance into a single number
- Cross entropy
- Sum of squared errors mse



goal \rightarrow minimize loss.

Cross entropy loss

(cost func)

✓ $y_n = 1 \Rightarrow$
✗ $y_n = 0 \Rightarrow$

- Used for softmax outputs (probabilities), or "classification" tasks

$$L = -\sum_n y_n \log q_n(x, \theta)$$

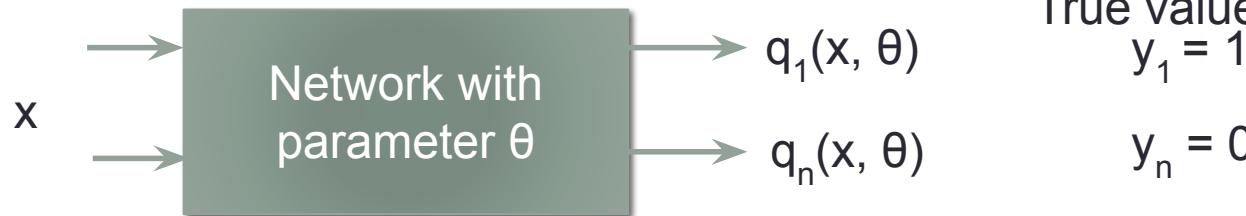
- Where y_n is 1 if data x comes from class n
0 otherwise

నుండి $y_3 \log q_3 = \underline{\log 0.3}$

నుండి: $y_1 \log q_1 = \log 1$
 $= 0.$

నుండి: $y_0 \log q_2 = \underline{\log 0}$

- L only has the term from the correct class
- L is non negative with highest value when the output matches the true values, a "loss" function



True values
 $y_1 = 1$

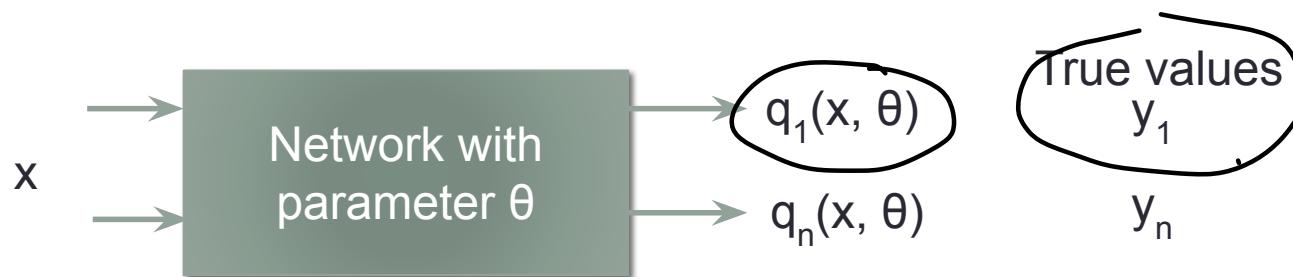
$y_n = 0$

Sum of squared errors

- Used for any real valued outputs such as regression

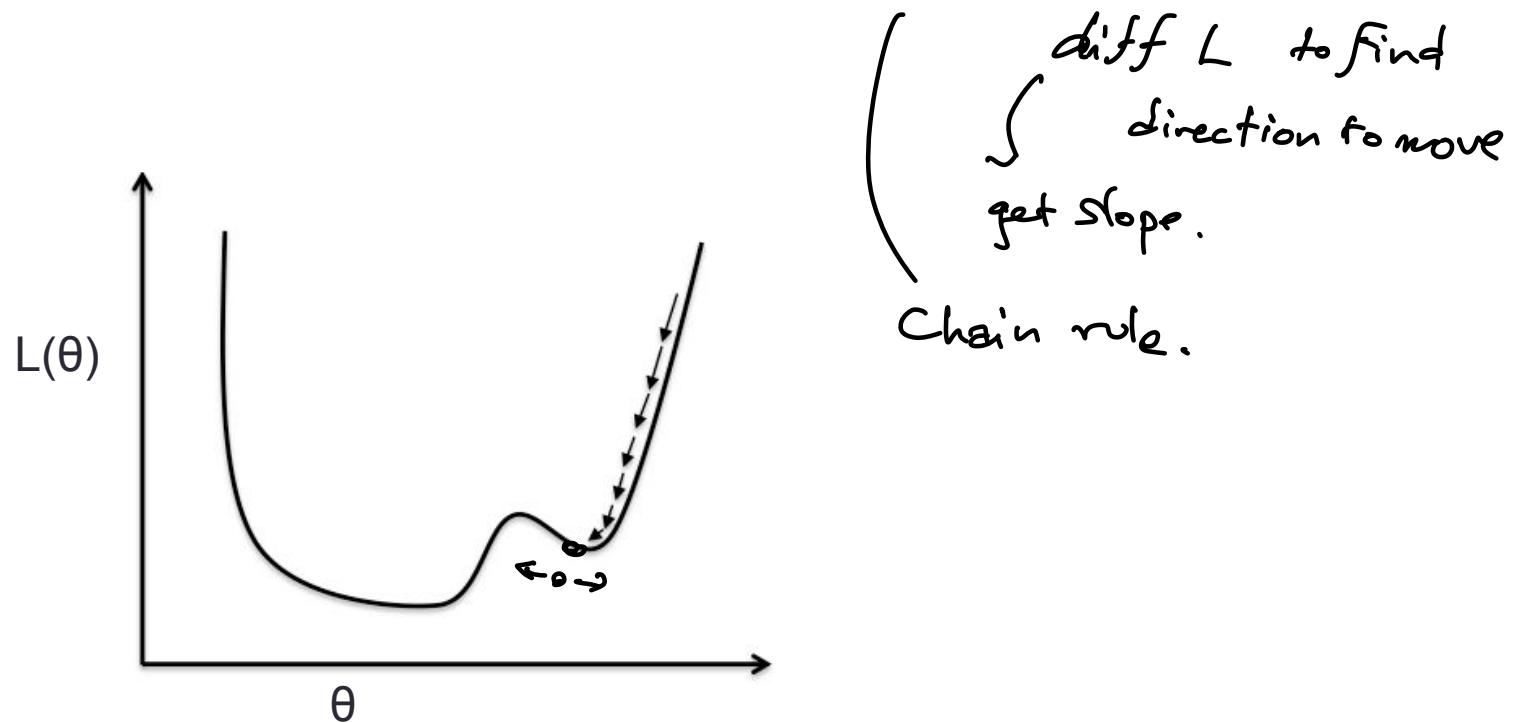
$$L = \frac{1}{2} \sum_n (y_n - q_n(x, \theta))^2$$

- Non negative, the better the lower the loss



Minimization using gradient descent

- We want to minimize L with respect to θ (^(Loss) weights and biases)
 - Differentiate with respect to θ
 - Gradients passes through the network by Back Propagation



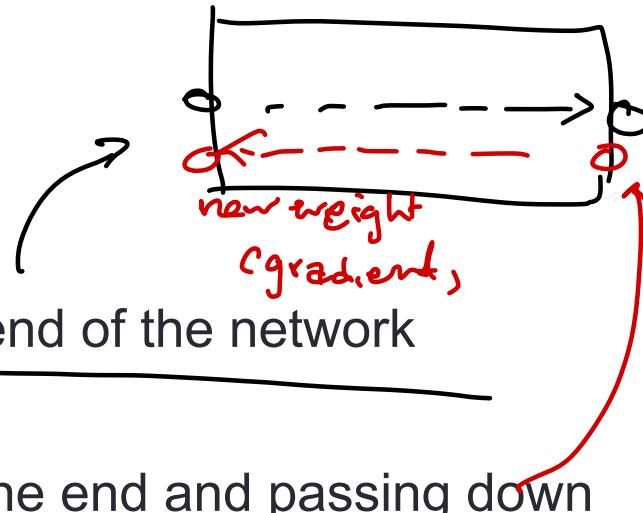
Differentiating a neural network model

- We want to minimize loss by gradient descent
- A model is very complex and have many layers! How do we differentiate this!!?



Back propagation

- Forward pass
 - Pass the value of the input until the end of the network
- Backward pass
 - Compute the gradient starting from the end and passing down gradients using chain rule



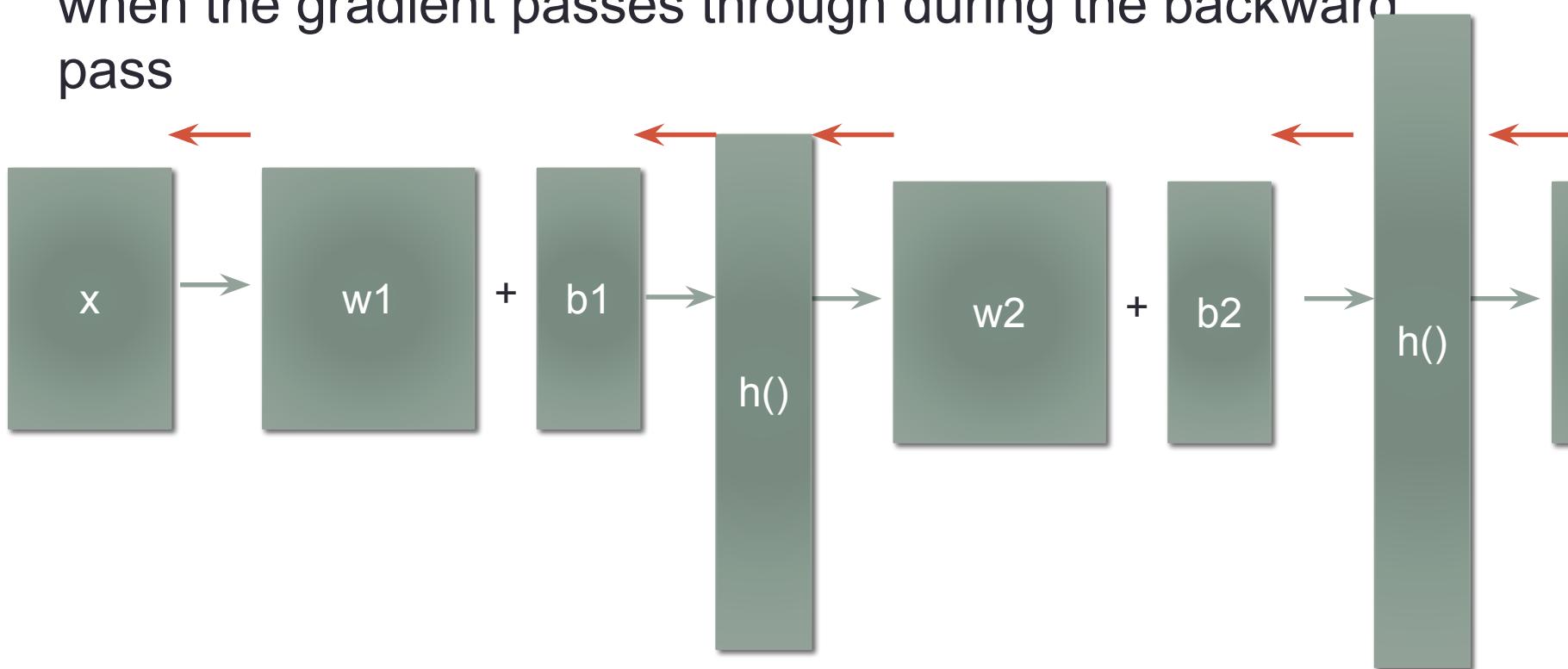
Examples to read

<https://alonalj.github.io/2016/12/10/What-is-Backpropagation/>

<https://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/>

Backprop and computation graph

- We can also define what happens to a computing graph when the gradient passes through during the backward pass



This lets us to build any neural networks without having to redo all the derivation as long as we define a forward and backward computation for the block.

Initialization

କେମାନ୍ତିରୁଣ୍ଡିଲୁବାକୁ ?
ପାରାମିଟ୍ରୋଫୁଲିଙ୍ଗି.

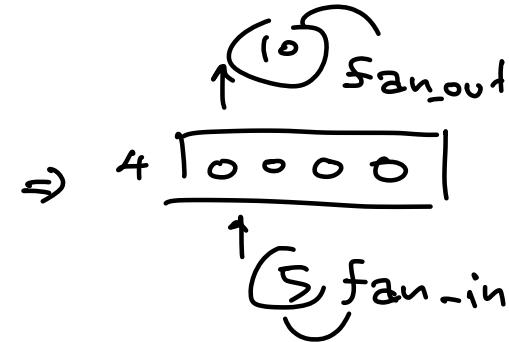
- The starting point of your descent
- Important due to local minimas
- Not as important with large networks AND big data
- Now usually initialized randomly

- One strategy (Xavier init) (Glorot)

$$\underline{\text{var}(w)} = \frac{2}{(\text{fan_in} + \text{fan_out})}$$

- For ReLUs (He init)

$$\underline{\text{var}(w)} = \frac{2}{\text{fan_in}}$$



- Or use a pre-trained network as initialization

X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. 2010

Kaiming He, Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. 2015

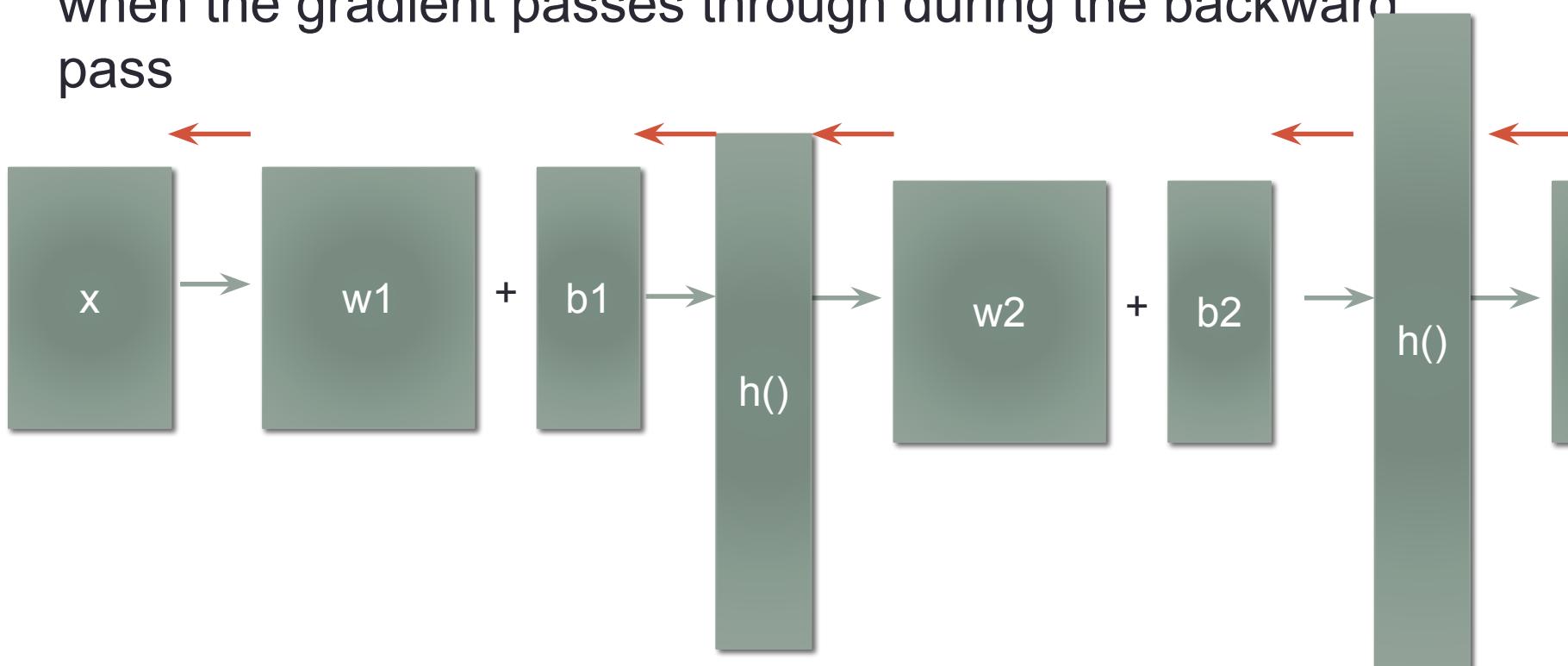
Stochastic gradient descent (SGD)

- Consider you have one million training examples *too big*
 - Gradient descent computes the objective function of all samples, then decide direction of descent
 - Takes too long
 - SGD computes the objective function on subsets of samples
 - The subset should not be biased and properly randomized to ensure no correlation between samples
 - The subset is called a mini-batch
 - Size of the mini-batch determines the training speed and accuracy
 - Usually somewhere between 32-1024 samples per mini-batch
 - Definition: 1 batch vs 1 epoch
-

Note: one can read relationship between batch size and learning rate here
<https://arxiv.org/abs/1711.00489>

Backprop and computation graph

- We can also define what happens to a computing graph when the gradient passes through during the backward pass

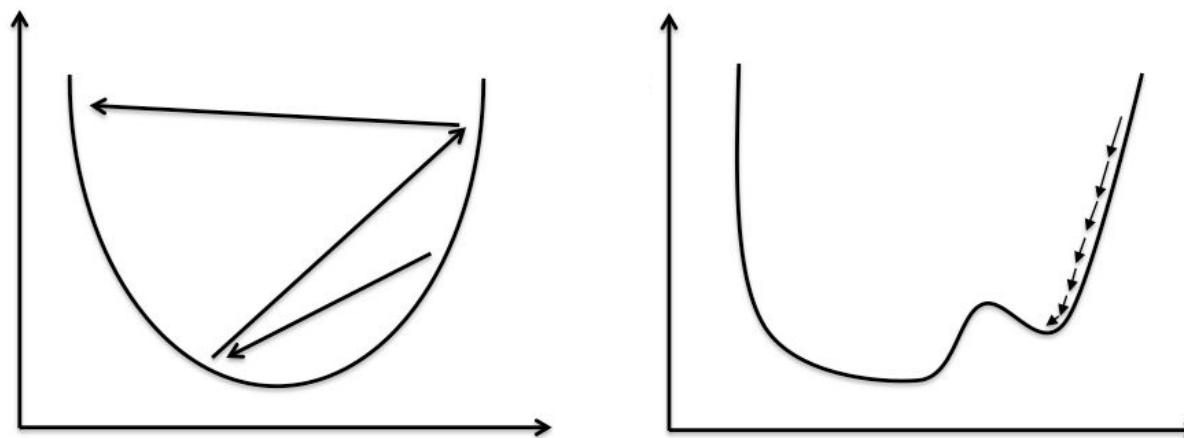


This lets us to build any neural networks without having to redo all the derivation as long as we define a forward and backward computation for the block.

Learning rate

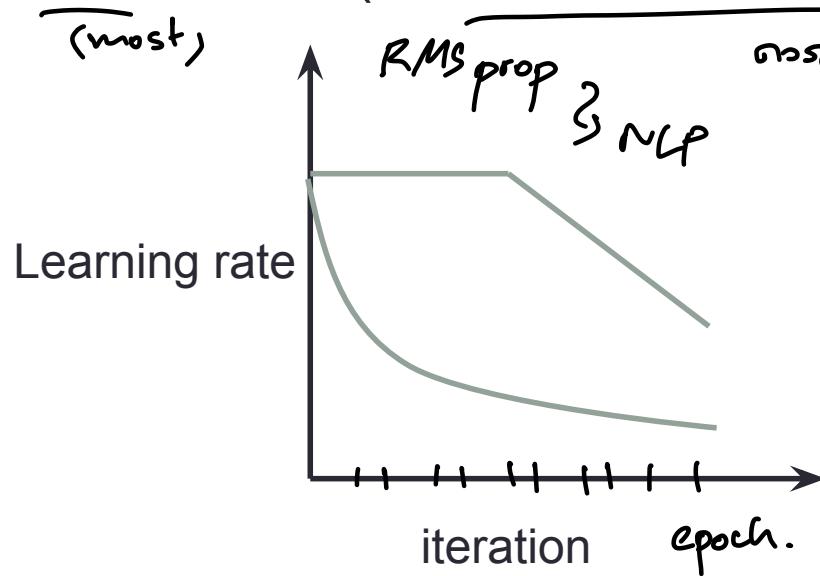
→ "significant."

- How fast to go along the gradient direction is controlled by the learning rate
- Too large models diverge
- Too small the model get stuck in local minimas and takes too long to train



Learning rate scheduling

- Usually starts with a large learning rate then gets smaller later
 - Depends on your task
 - Automatic ways to adjust the learning rate : Adagrad, Adam, etc. (still need scheduling still)
- for Lr \leftarrow after epoch
 which
- optimizer.
- $\theta = \text{init}$
 model
- \rightarrow train
- (N local
 minima)



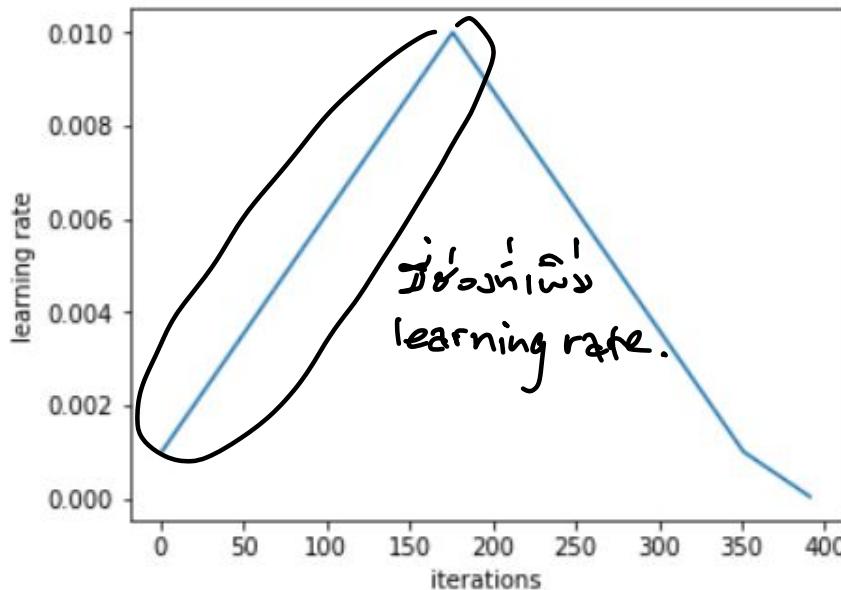
Learning rate warm up

ສ້າງລົງໃນ network ຂັງ
(Big Network)

Initial point of the network can be at a bad spot.

Try not to go to fast - has a warm up period.

Useful for large datasets, or adaption (transfer learning)



Potentially leads to faster convergence and better accuracy

See links below for methods to select the shape of the triangle

<https://sgugger.github.io/the-1cycle-policy.html#the-1cycle-policy>

[Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour](#)

[Cyclical Learning Rates for Training Neural Networks](#)

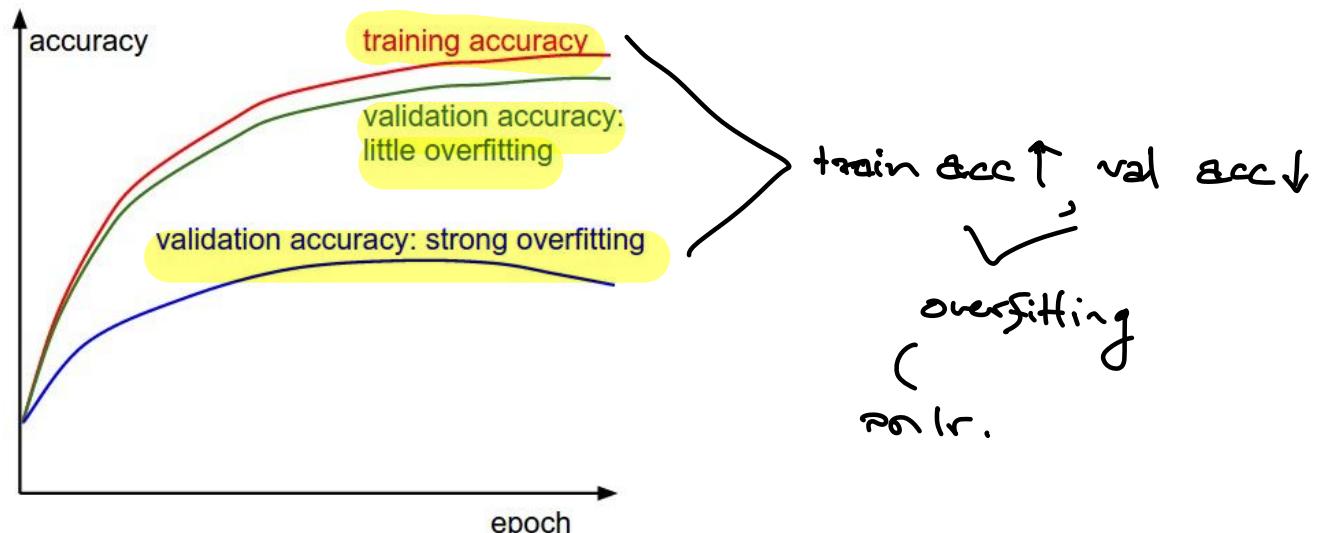
Optimizers

- Besides learning rate scheduling (coarse grain) we can do finer (and automatic) control of the learning rate
 - SGD + momentum (basic)
 - RMSprop
 - Faster than SGD but slower than ADAM
 - More stable than ADAM
 - ADAM + variants
 - Most popular for its ease of use

People find simple SGD with momentum and decay to perform better (with proper tuning)

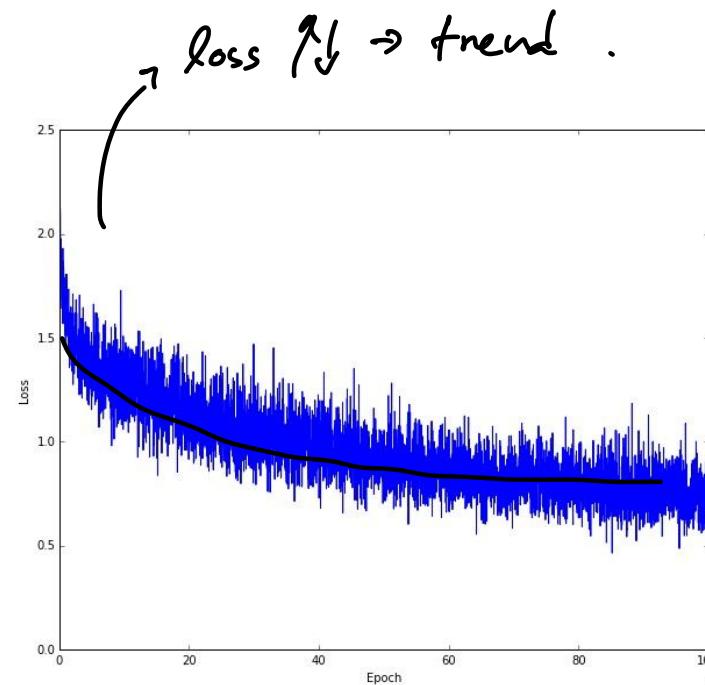
Overfitting

- You can keep doing back propagation forever!
- The training loss will always go down
- But it overfits
- Need to monitor performance on a held out set
- Stop or decrease learning rate when overfit happens



Monitoring performance

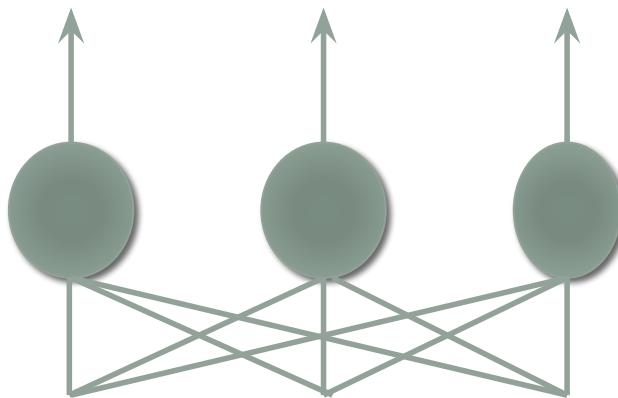
- Monitor performance on a “dev/validation set”
 - This is NOT the test set (cheat!)
- Can monitor many criterions
 - Loss function
 - Classification accuracy
- Sometimes these disagree
- Actual performance can be noisy, need to see the trend



Reducing overfitting - dropout

- A regularization technique for reducing overfitting
- Randomly turn off different subset of neurons during training
 - Network no longer depend on any particular neuron
 - Force the model to have redundancy – robust to any corruption in input data
 - A form of performing model averaging (ensemble of experts)
- Now a standard technique

Dropout visualized



Model

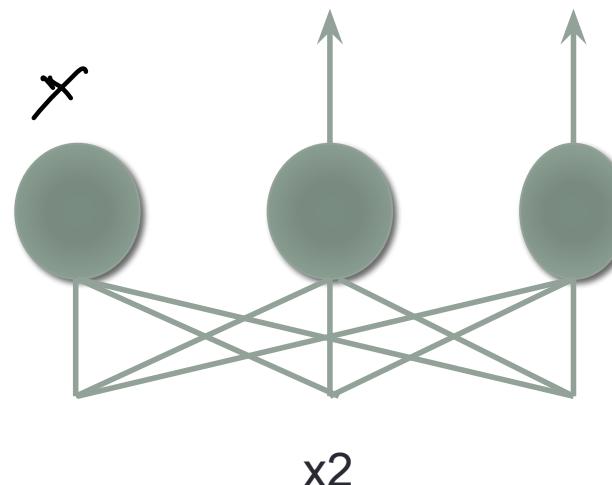
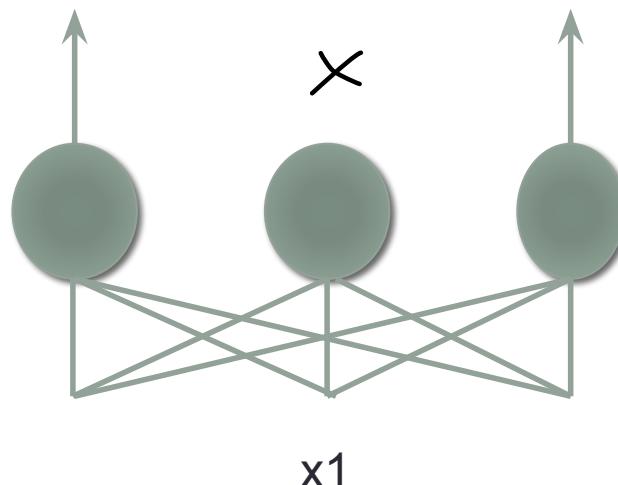
Dropout training time

as per convention

Drop out rate 0.33 $\frac{1}{3}$.

即 66% .

Randomly removed outputs
for each training sample

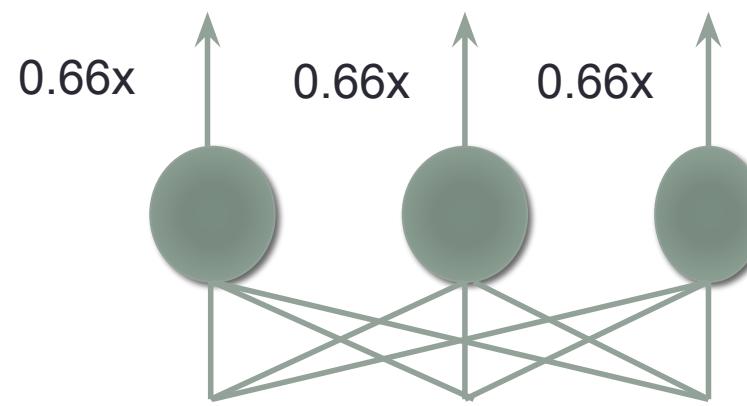
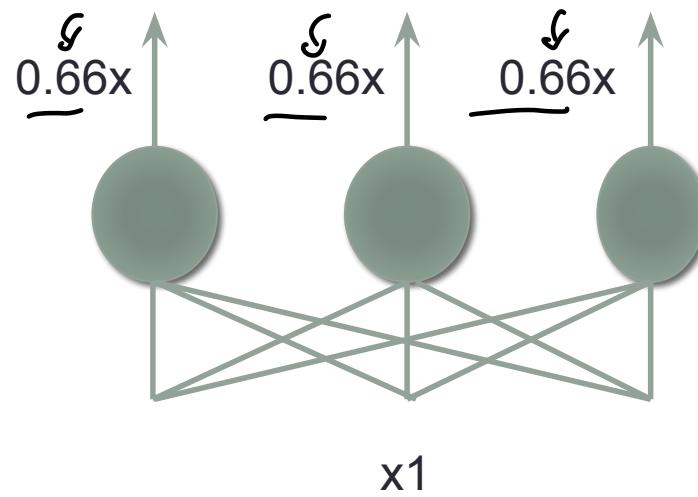


Dropout test time

Dropout rate 0.33

רְדוֹפָט → שְׁאַלְמִינְסְּ 66%

Scale outputs so the output contribution is around the same



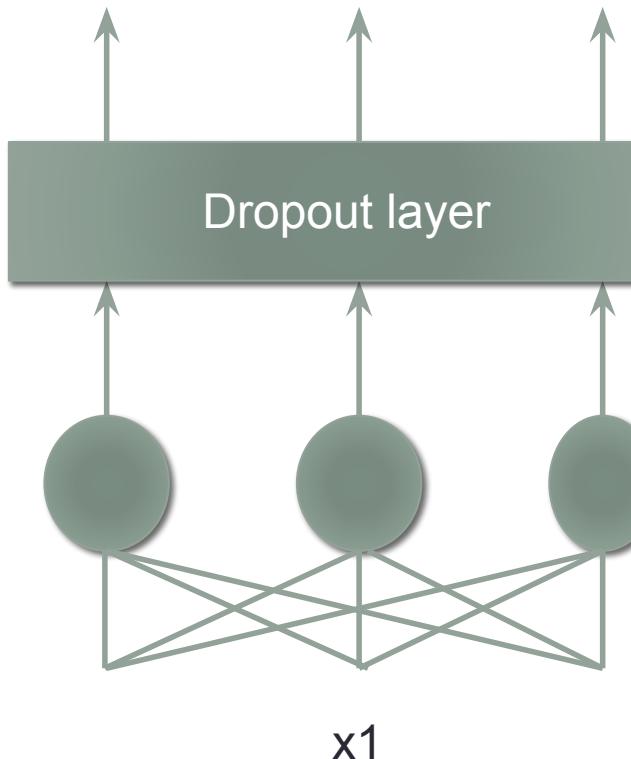
↳ why we open it all → because we use it

↳
random

✗ fd (no)
o
o
✗ time spent

↳ open all then scale.

Dropout implementation



Dropout layer

Just another layer that drops inputs

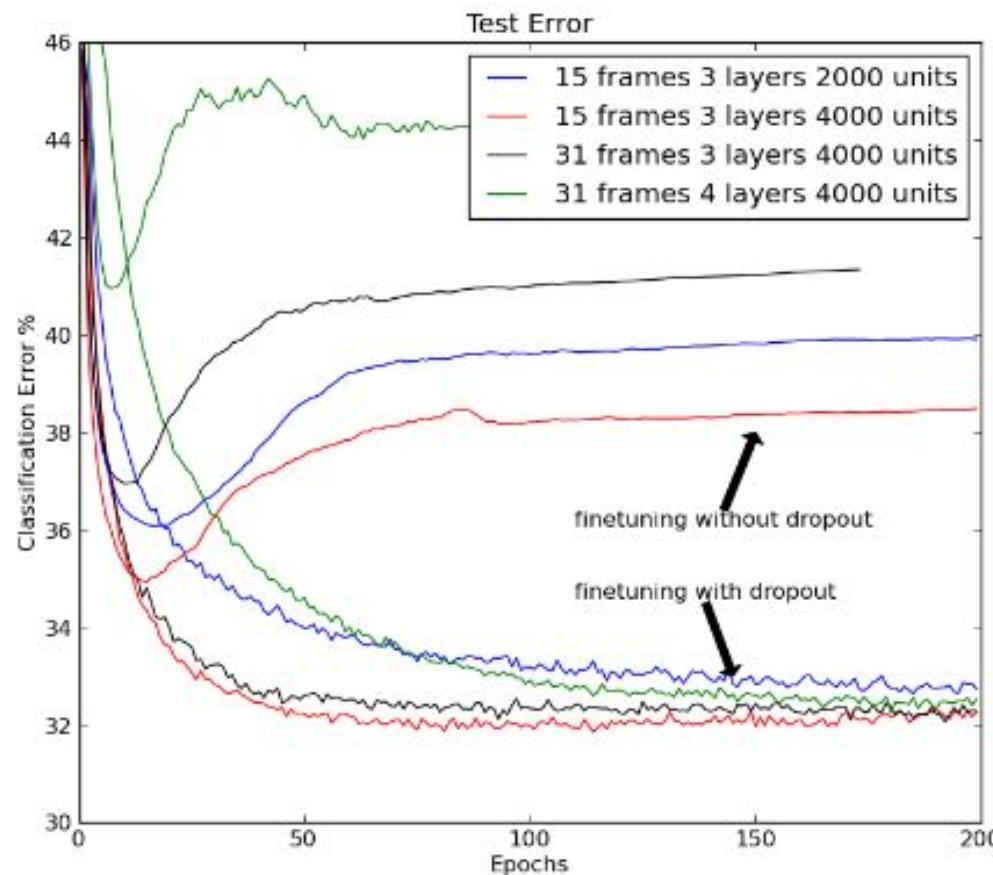
usually you $\times 0.66 \rightarrow \sigma_{\text{std}} = \times \text{original train}$,

Inverted dropout

A variant of dropout that scales the dropout at training time, so that you don't have to scale at test time.

Dropout on TIMIT

- A phoneme recognition task



→ reduce overfit

Batch normalization

- Recent technique for (implicit) regularization
- Normalize every mini-batch at various batch norm layers to standard Gaussian (different from global normalization of the inputs)
- Place batch norm layers before non-linearities
- Faster training and better generalizations

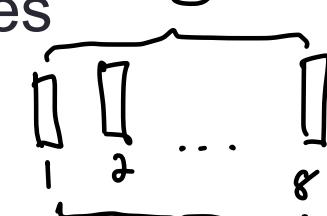
For each mini-batch that goes through batch norm

- Normalize by the mean and variance of the mini-batch for each dimension
- Shift and scale by learnable parameters

Replaces dropout in some networks

1. mean, std, mini batch
normalize
out.

$3 - 2, \dots, 5$



$$\hat{x} = \frac{x - \mu_b}{\sigma_b}$$

mini batch
mini batch
mini batch

$$y = \alpha \hat{x} + \beta$$

Dropout vs batchnorm

- You can add dropout in the hidden layers "(0-0.5)" *hidden layer.*
- Or input layers (0-0.2 is typical) *↓↓↓↓↓*
 - "Noising" the inputs, data augmentation
- Dropout in computer vision
 - use batchnorm instead
- Dropout in NLP *↖*
 - Usually works better than Batchnorm for NLP in simple architectures
 - Drop full words at the embedding
 - Recurrent dropout <http://arxiv.org/abs/1512.05287>

BERT • Recent works (transformers) find batchnorm to be better than dropout

- For seq2seq models, layer norm is popular

→ ນັບຕາມໄຕ້ມີຫົວ

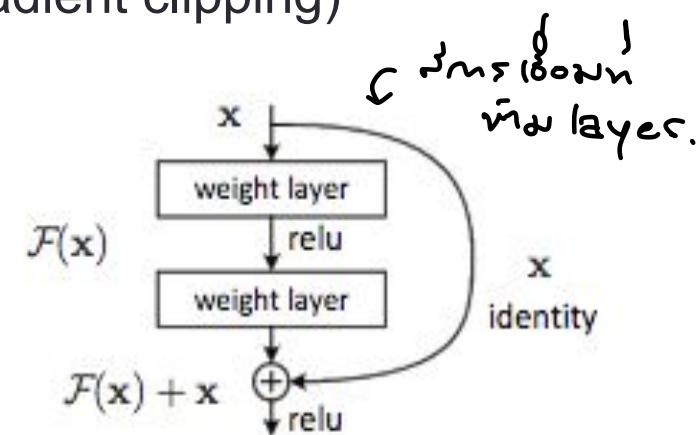
→ $100 = m \rightarrow$ not a number

Vanishing/Exploding gradient → ↑ 6 layers.

- Backprop introduces many multiplications down chain
- The gradient value gets smaller and smaller
 - The deeper the network the smaller the gradient in the lower layers
 - Lower layers changes too slowly (or not at all)
 - Hard to train very deep networks (>6 layers)
- The opposite can also be true. The gradient explodes from repeated multiplication
 - Put a maximum value for the gradient (Gradient clipping)

- How to deal with this?
 - Residual connection

<https://arxiv.org/pdf/1512.03385.pdf>



Neural networks

- Fully connected networks

- Neuron
- Non-linearity
- Softmax layer

- DNN training

- Loss function
- SGD and backprop
- Learning rate
- Overfitting

- CNN, RNN, LSTM, GRU

DENSE müssen. Select BatchSize, Epoch.



Theory: Biggest size that
match our "GPU"

* But if > 1000 → default in
Tensorflow
can be crash.

Do validation loss
→ overfit or no
over.

◦ DNN → do K-fold?

- K-fold → use in small data
- DNN usually not do (less different)

Convolutional Neural Networks (CNNs)

- Consider an image of a cat. DNNs need different neurons to learn every possible location a cat can be



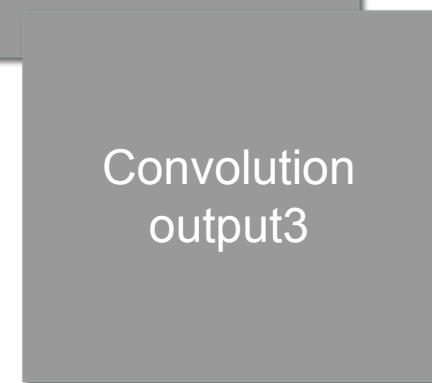
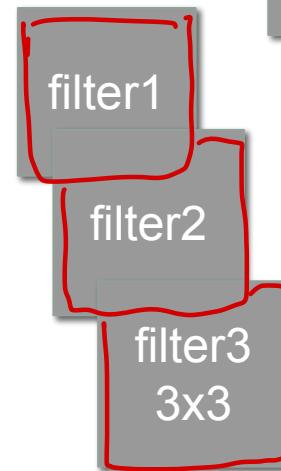
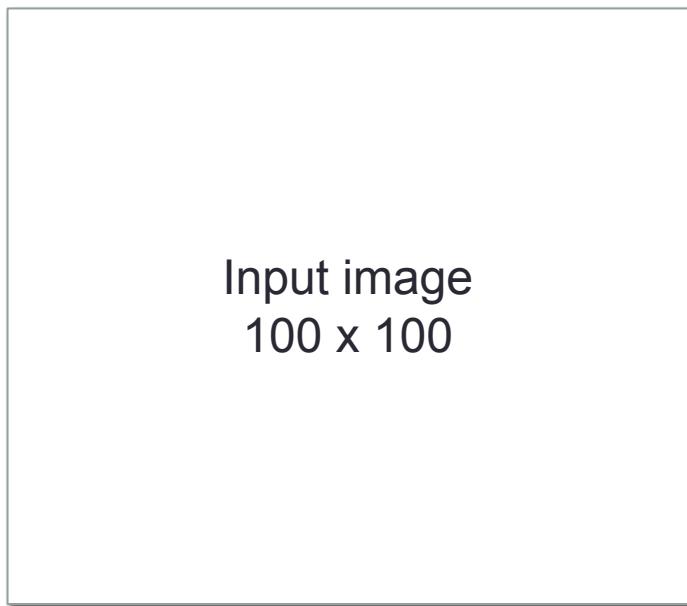
- Can we use the same parameters to learn that a cat exists regardless of location?

- 2 parts: convolutional layer and pooling layer

Convolutional filters

Multiply inputs with filter values

Output one feature map per filter

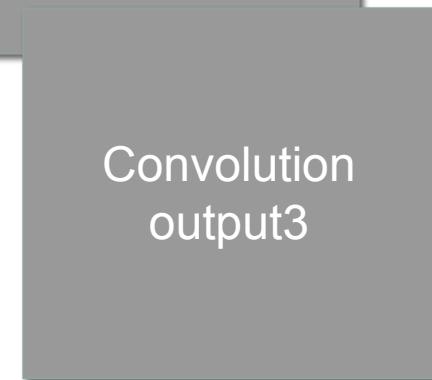
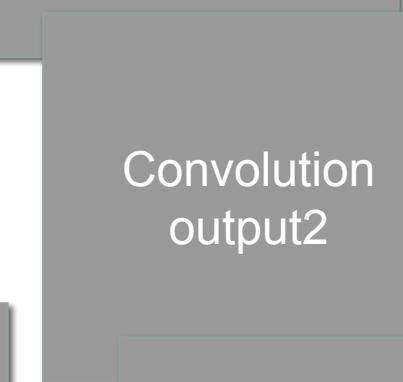
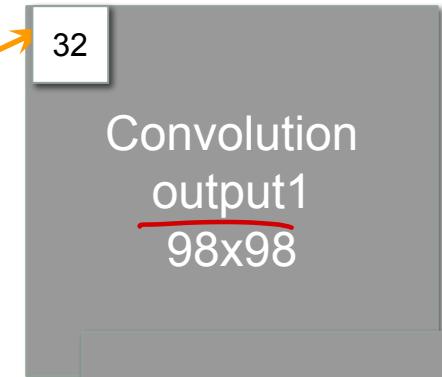
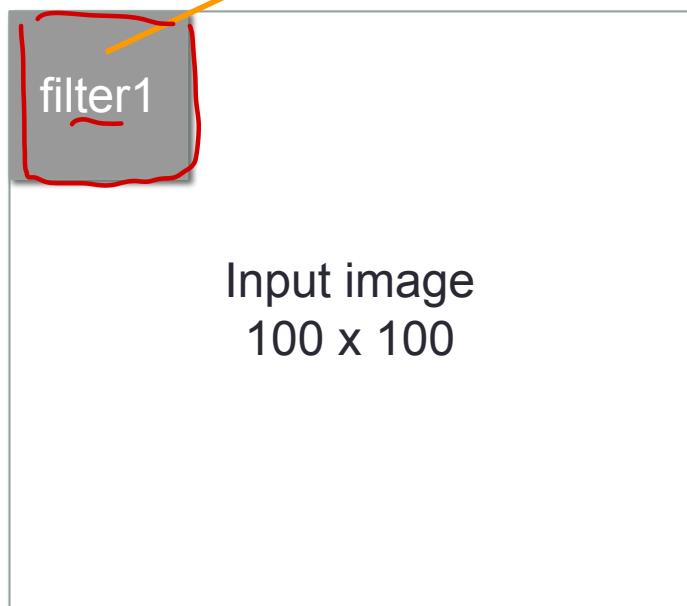


Convolutional filters

0	1	-1
1	0	1
1	2	0
1	2	3
4	5	6
7	8	9

}

$$1*2 + -1*3 + 1*4 + 1*6 + 1*7 + 2*8 = 32$$

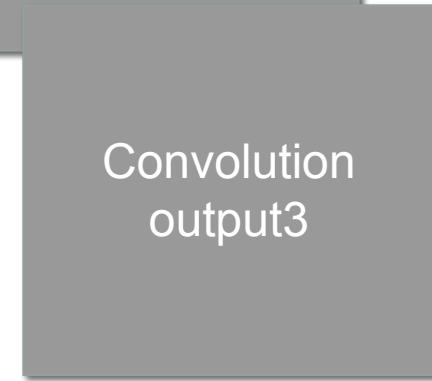
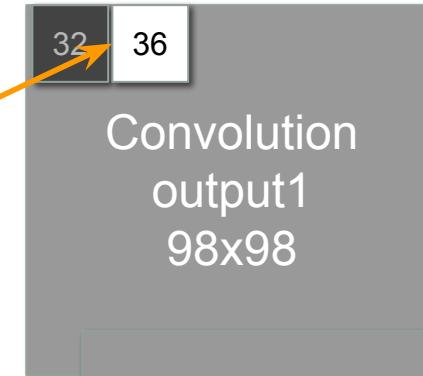
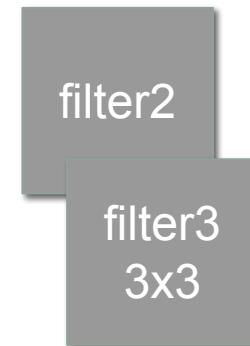
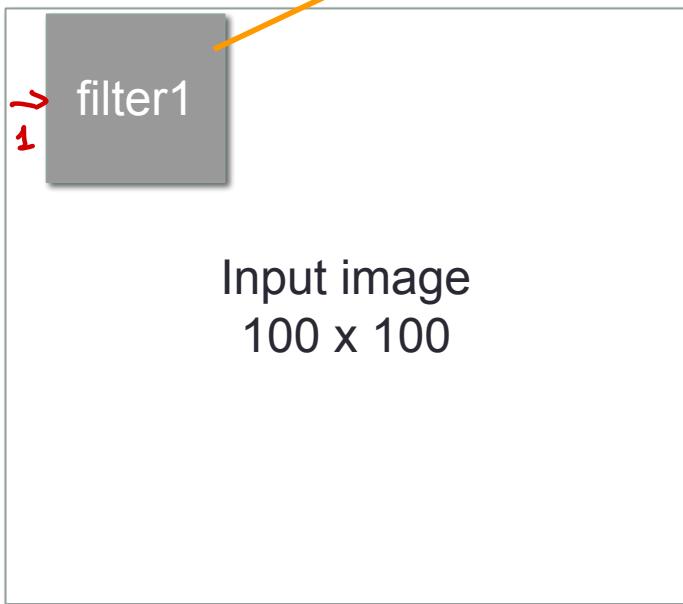


Convolutional filters

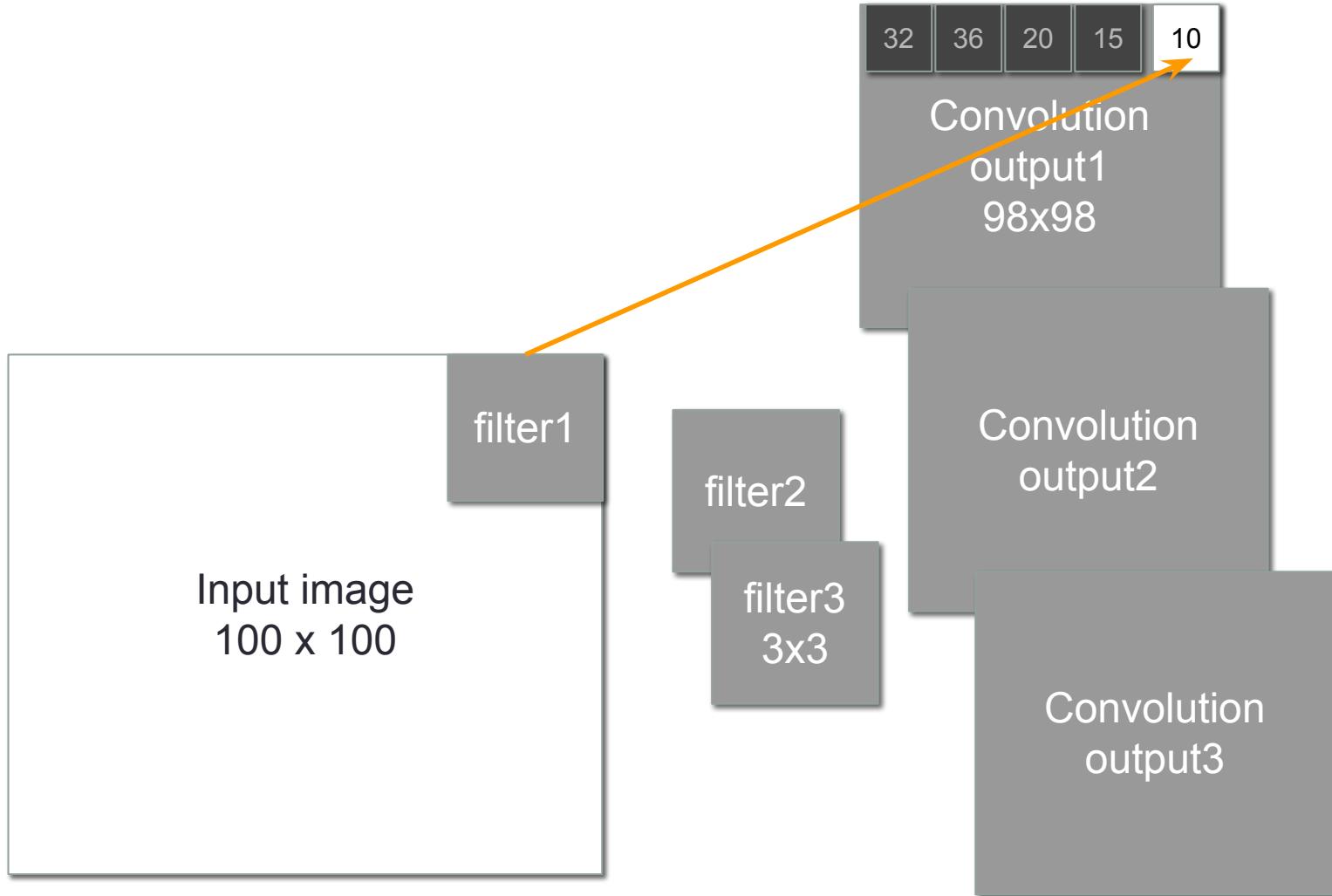
Stride of 1 *use 1 step*.

0	1	-1
1	0	1
1	2	0
2	3	1
5	6	3
8	9	8

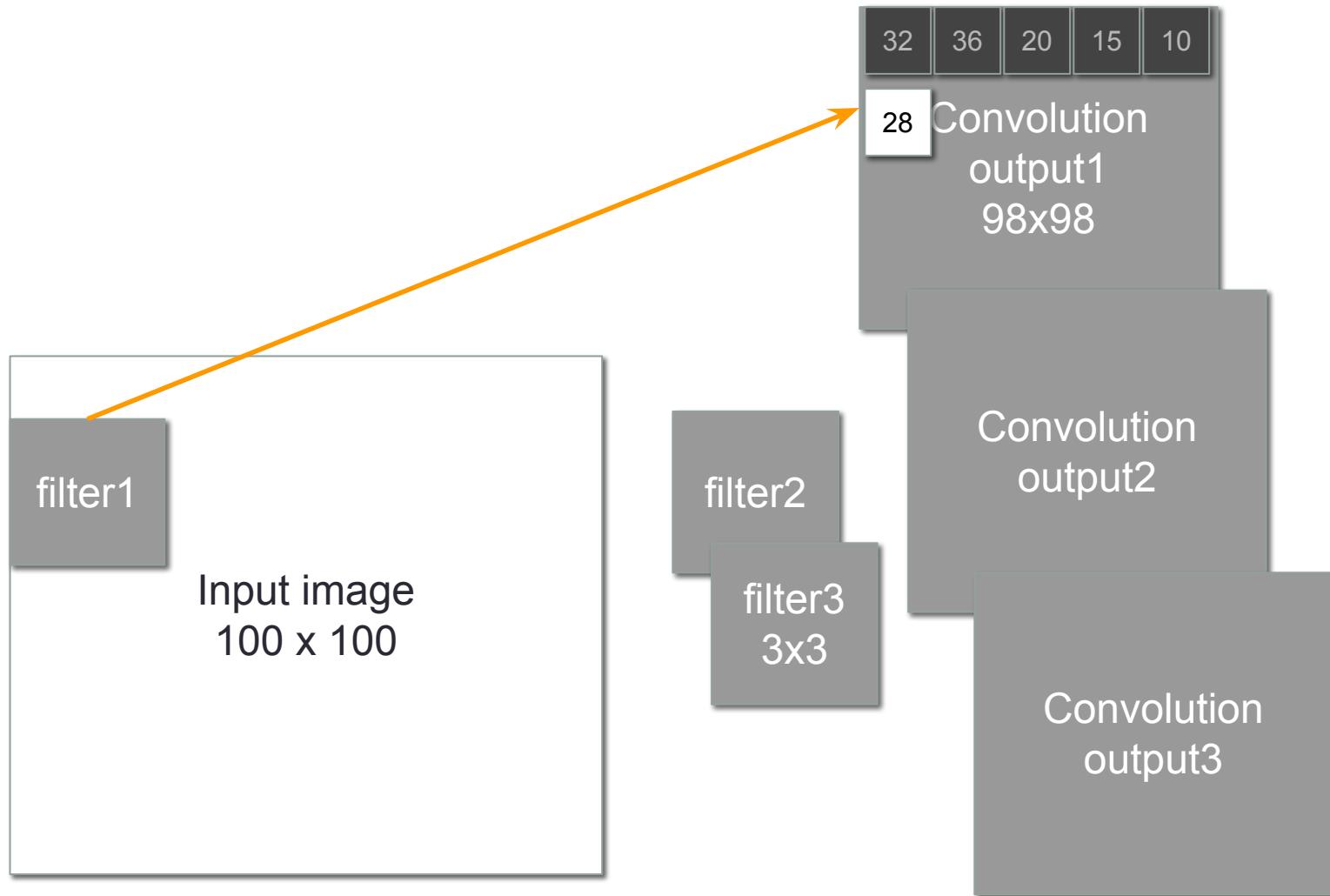
$$1*3 + -1*1 + 1*5 + 1*3 + 1*8 + 2*9 = 36$$



Convolutional filters

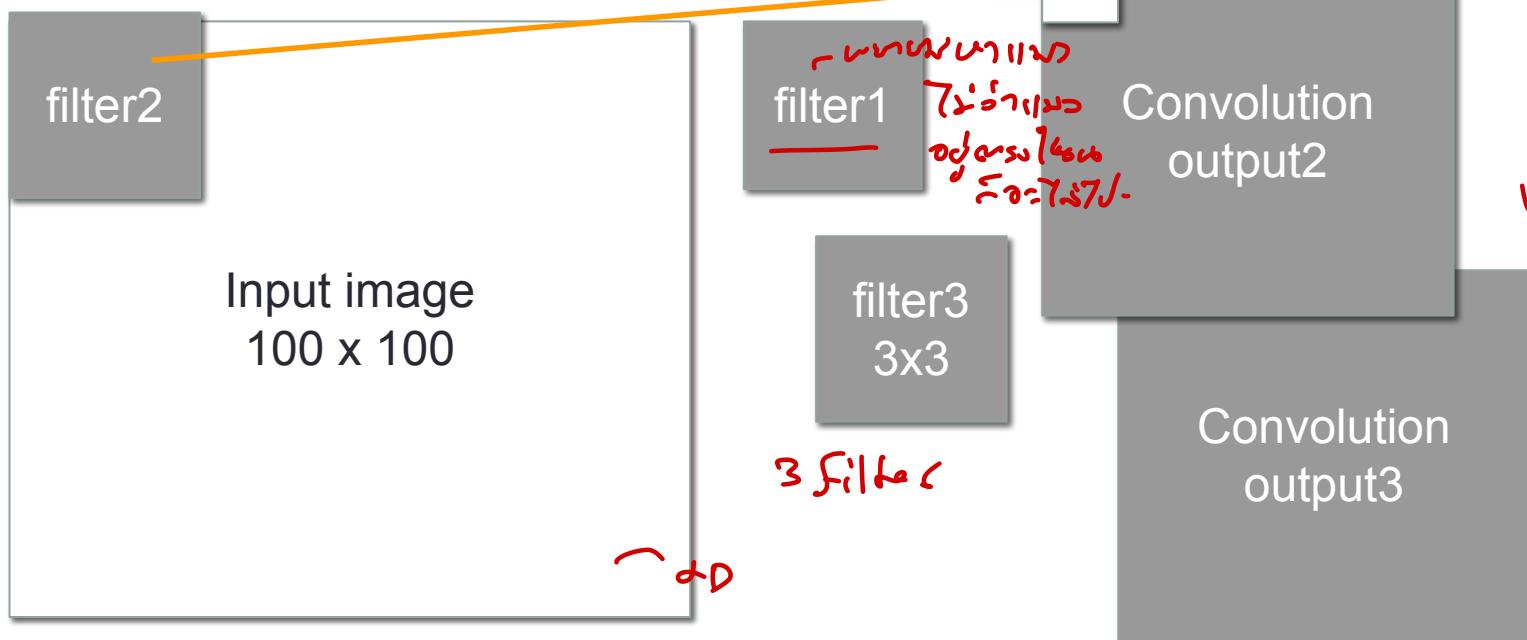
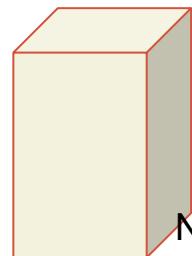


Convolutional filters



Convolutional filters

N filters means N feature maps
You get a 3 dimensional output



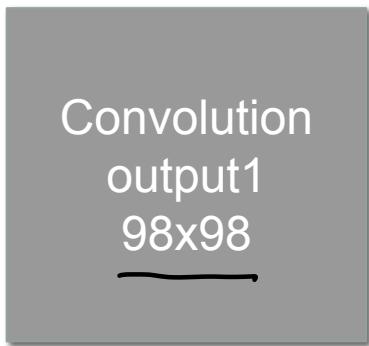
32	36	20	15	10
28	72	0	12	50
32	36	18	9	2
17	6	2	17	1

16

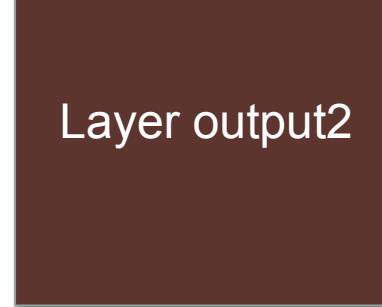
3D

Pooling/subsampling

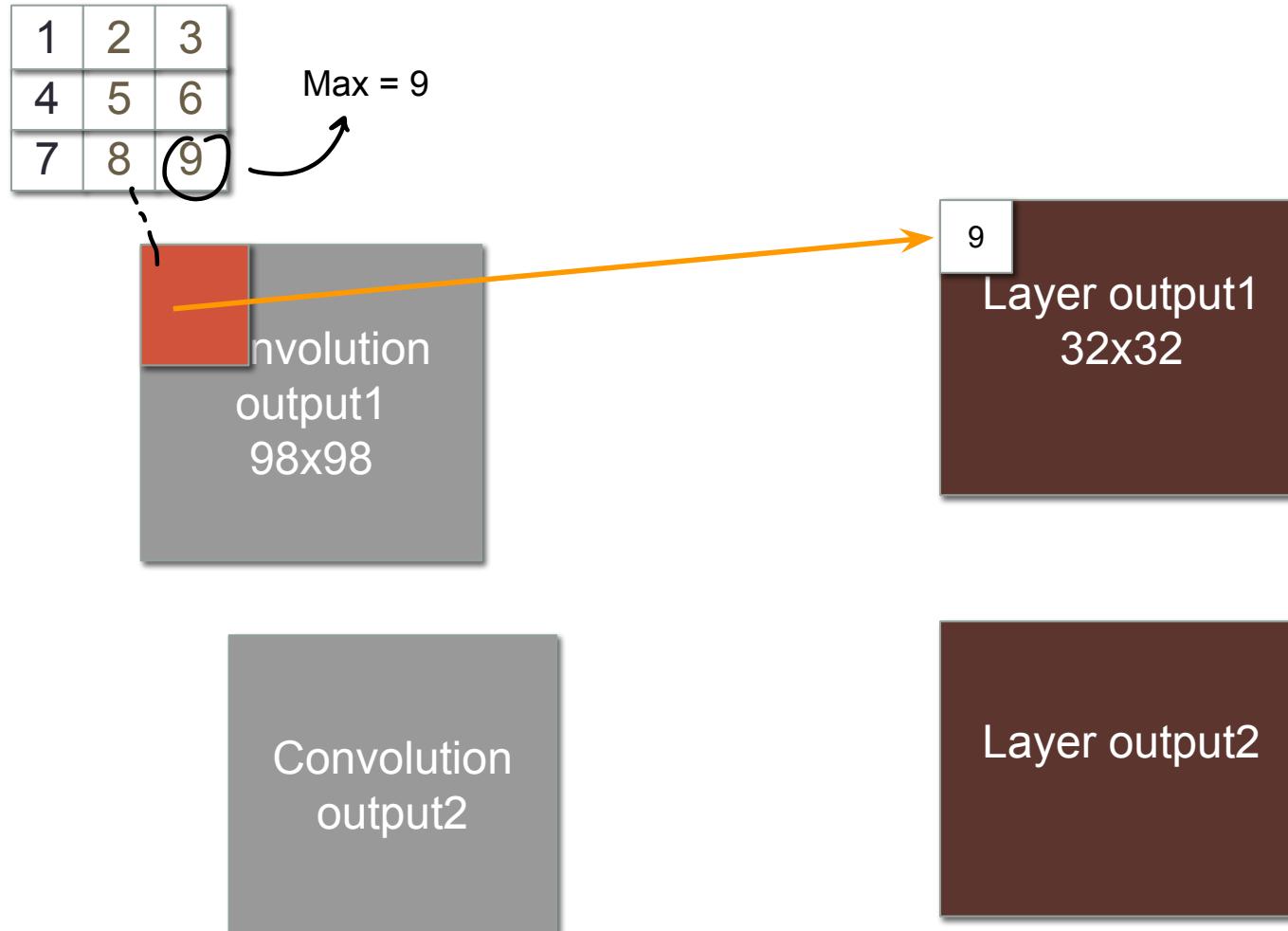
Reduce dimension of the feature maps



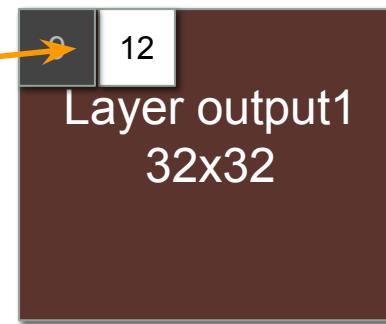
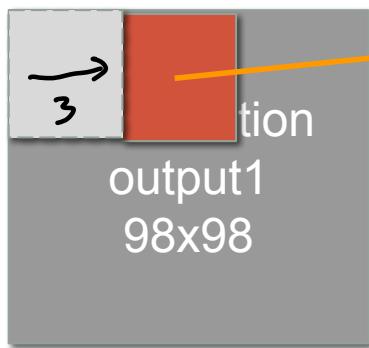
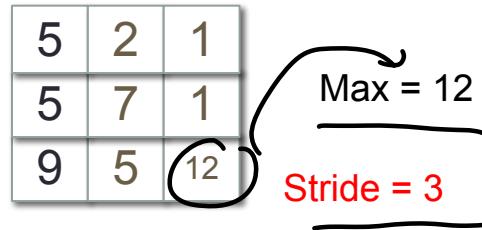
3x3 (Max) filter
with no overlap

A text label describing the pooling operation: "3x3 (Max) filter with no overlap". The word "Max" is circled in red.

Pooling/subsampling



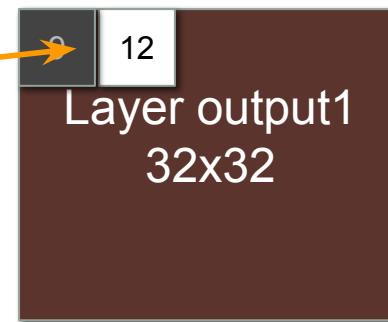
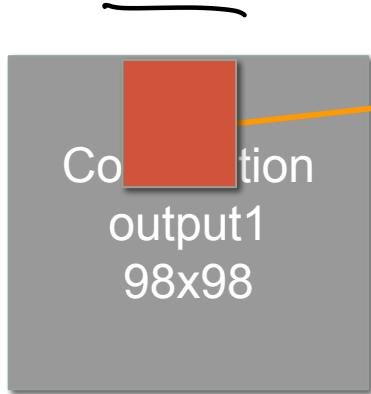
Pooling/subsampling



Pooling/subsampling

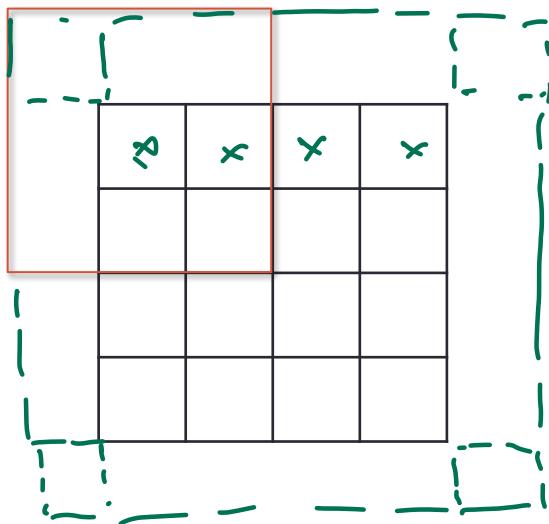
① = $\frac{\text{input size}}{\text{pool size}}$

Can use other functions besides max
Example, average



Convolution puzzle

5 filters, 3x3 filter (pad), stride 1, pad 1
What is the output size?

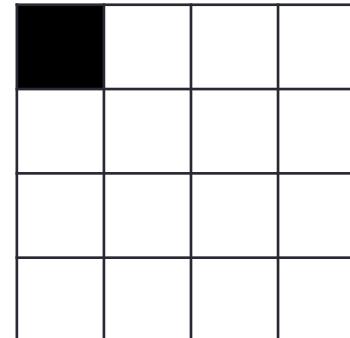
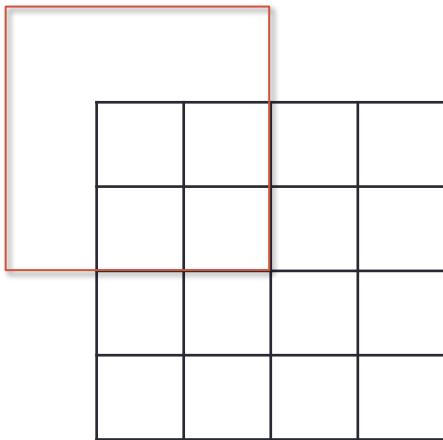


$$4 \times 4 \curvearrowleft 5 \text{ filters.}$$
$$4 \times 4 \times 5 =$$

Convolution puzzle

5 filters 3x3 filter pad, stride 1, pad 1

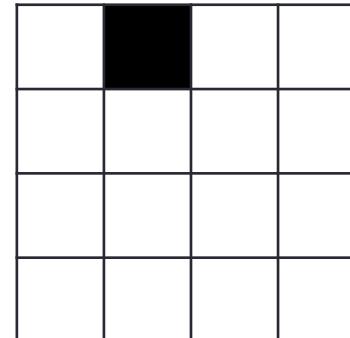
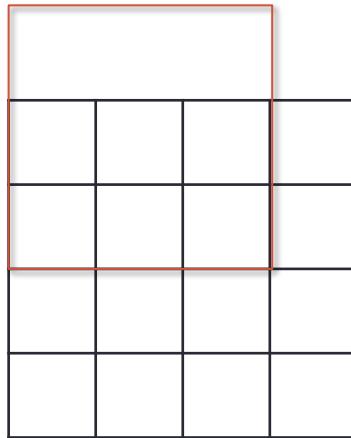
What is the output size?



Convolution puzzle

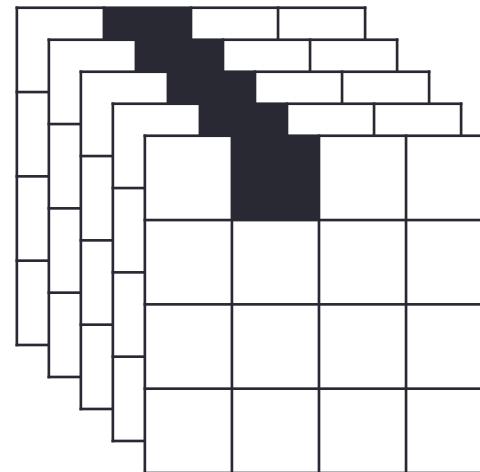
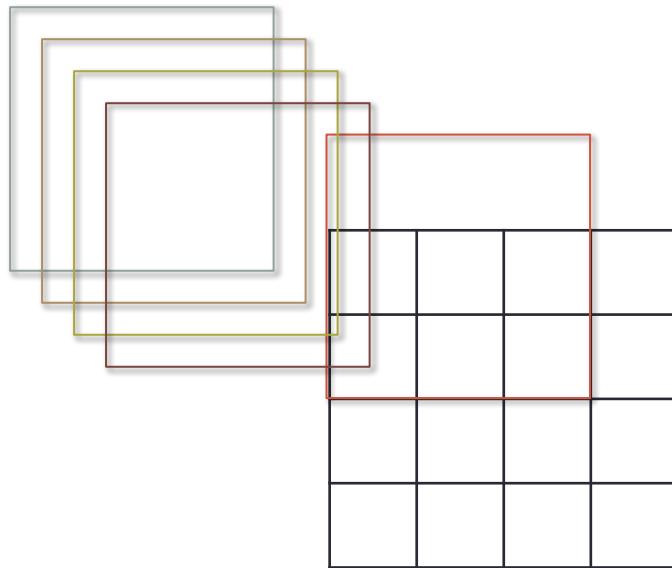
5 filters 3x3 filter pad, stride 1, pad 1

What is the output size?



Convolution puzzle

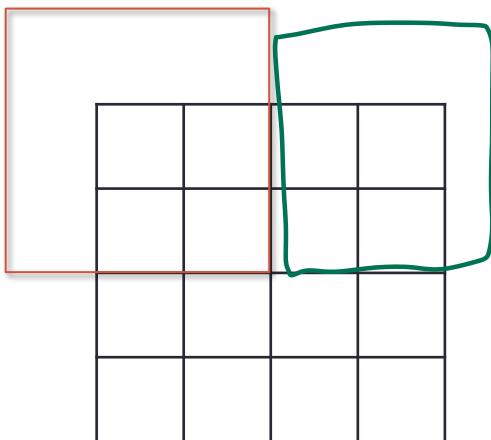
5 filters 3x3 filter pad, stride 1, pad 1



Convolution puzzle

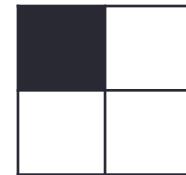
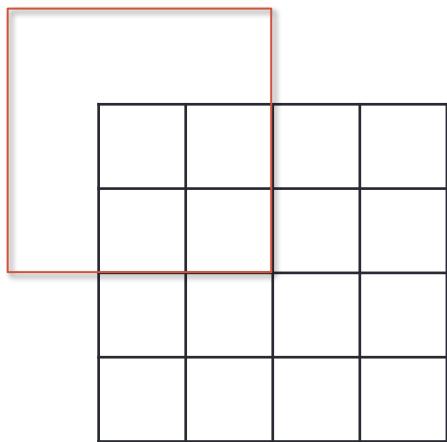
3x3 filter pad, stride 2, pad 1

→ 2×2



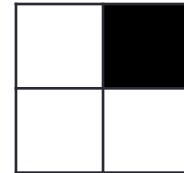
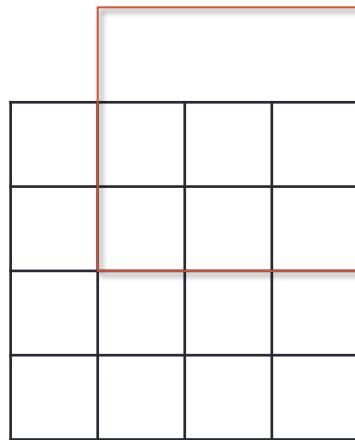
Convolution puzzle

3x3 filter pad, stride 2, pad 1



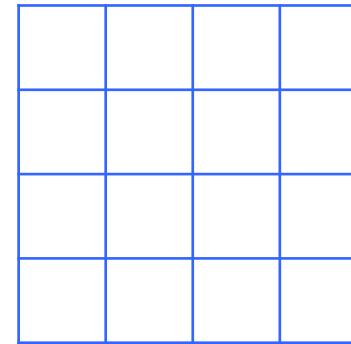
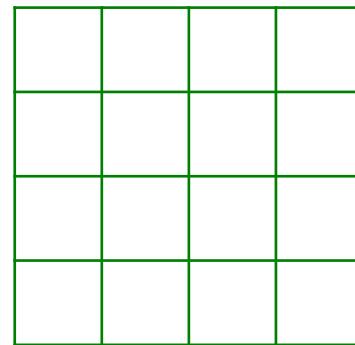
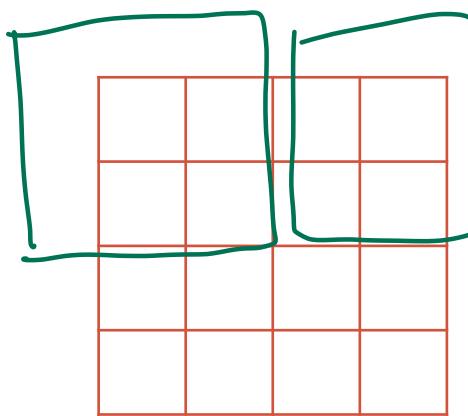
Convolution puzzle

3x3 filter pad, stride 2, pad 1



Convolution puzzle

RGB input (3 channels) 5 filters 3x3 filter pad, stride 2, pad 1

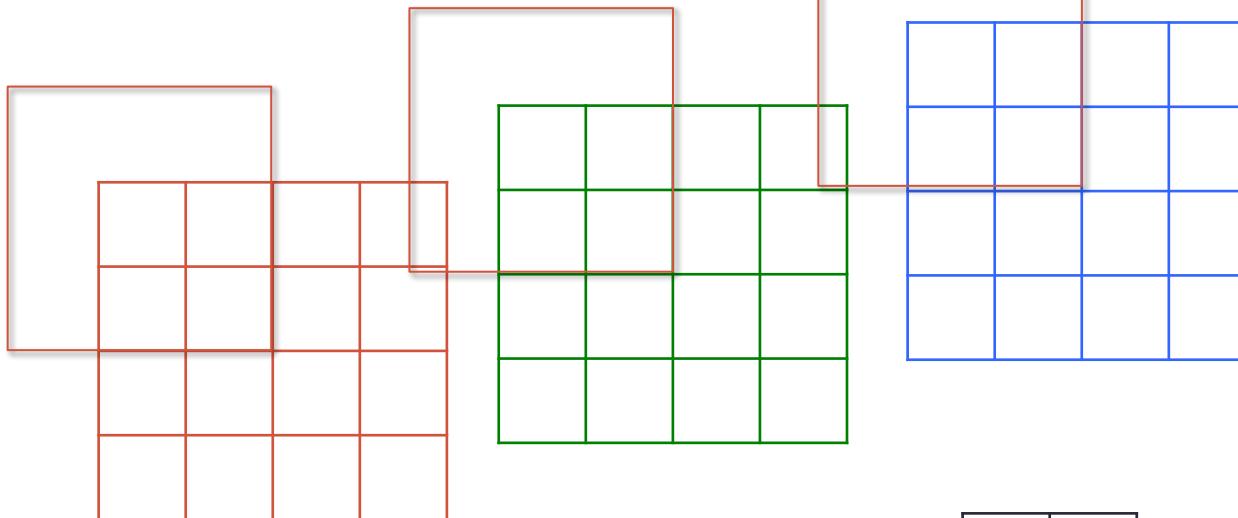


$2 \times 2 \times 5 \times 3$

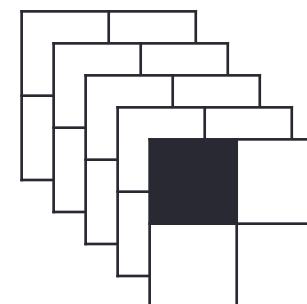
Convolution puzzle

filter views own input channel

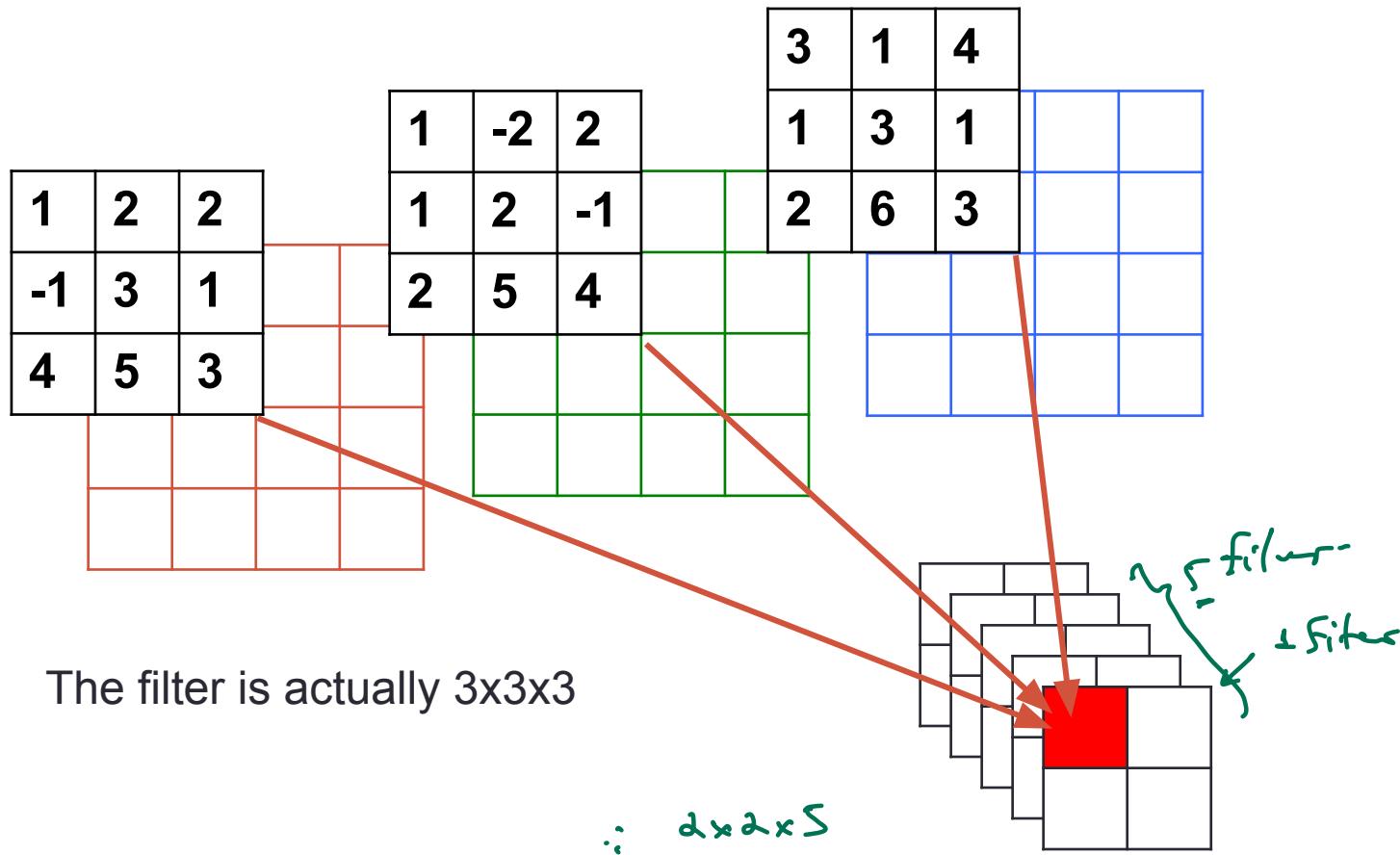
RGB input (3 channels) 5 filters 3x3 filter pad, stride 2, pad 1



The filter is actually 3x3x3



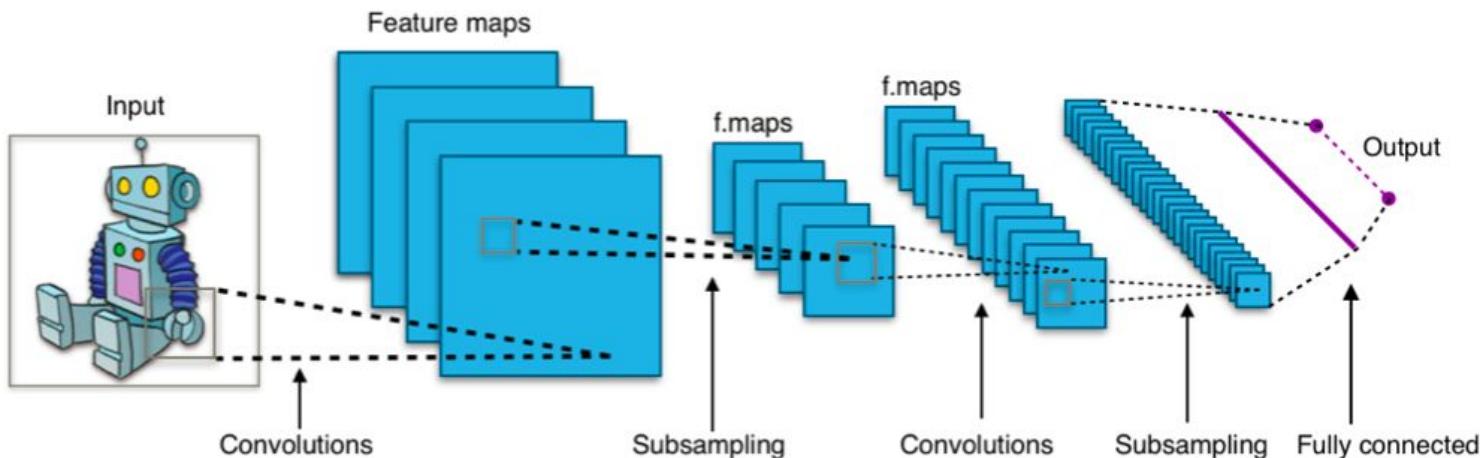
Convolution puzzle



CNN overview

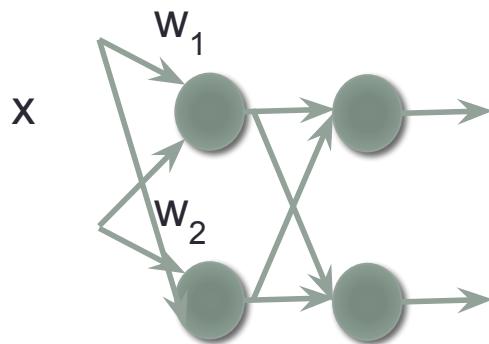
→ Toward's pattern, Fny, vector
Time series, NLP, CV.

- Filter size, number of filters, filter shifts, and pooling rate are all parameters
- Usually followed by a fully connected network at the end
 - CNN is good at learning low level features
 - DNN combines the features into high level features and classify

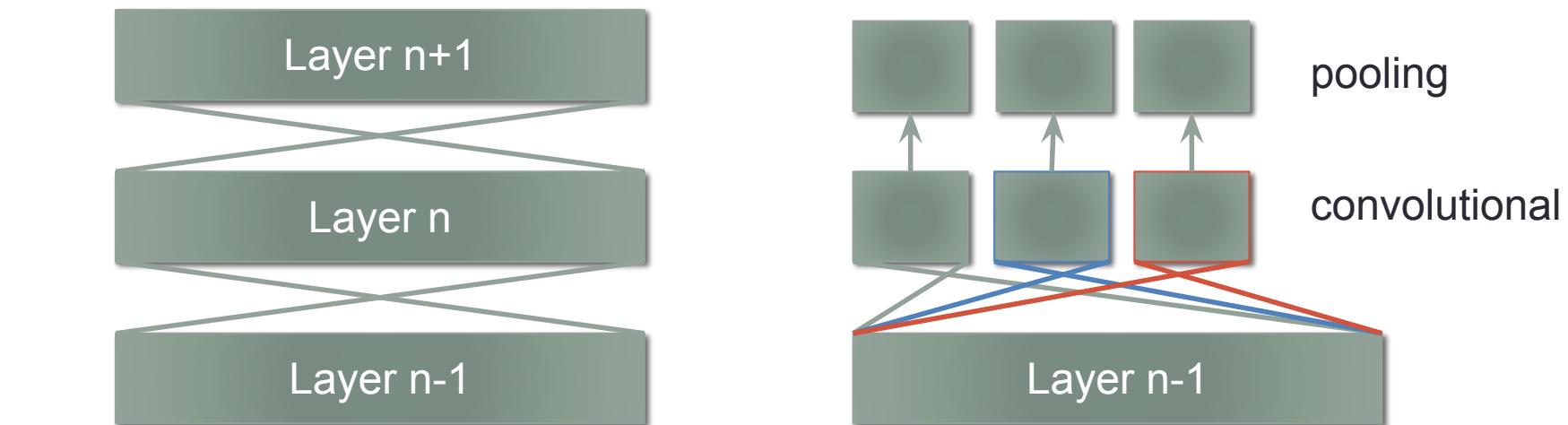


Parameter sharing in convolution neural networks

- $W^T x$

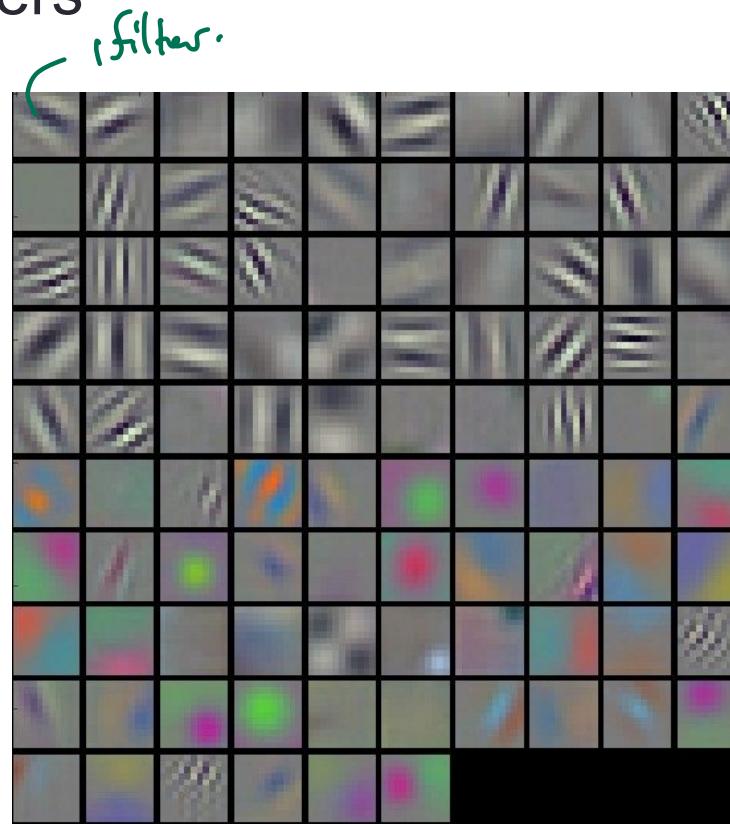


- Cats at different location might need two neurons for different locations in fully connect NNs.
- CNN shares the parameters in 1 filter
- The network is no longer fully connected



Visualizing convolutional layers

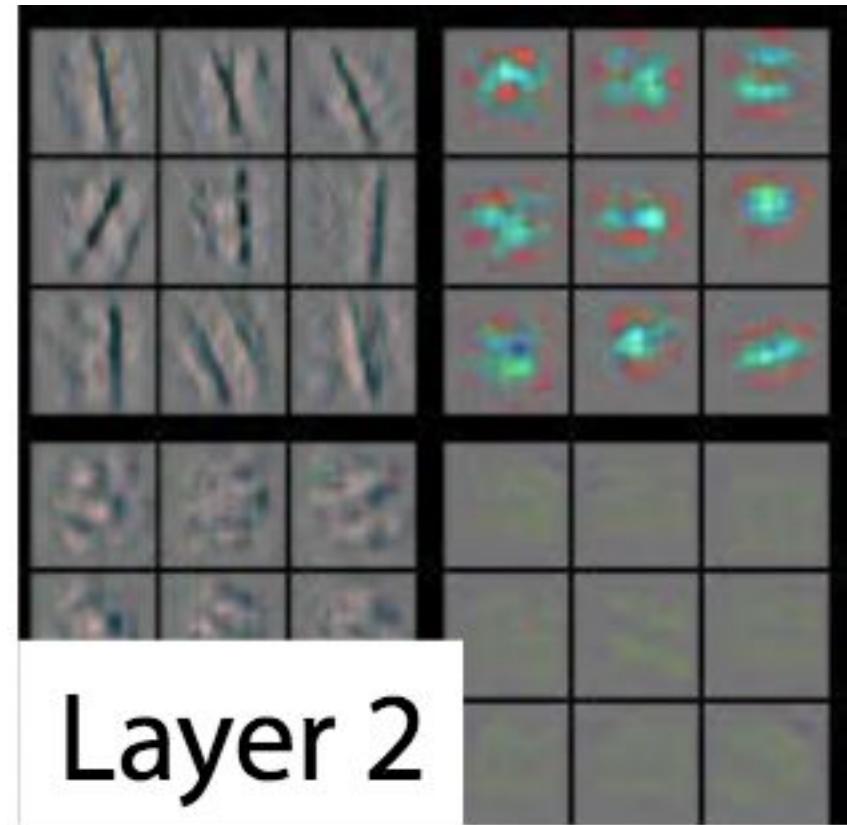
- We can visualize the weights of the filter
- “Matched filters”

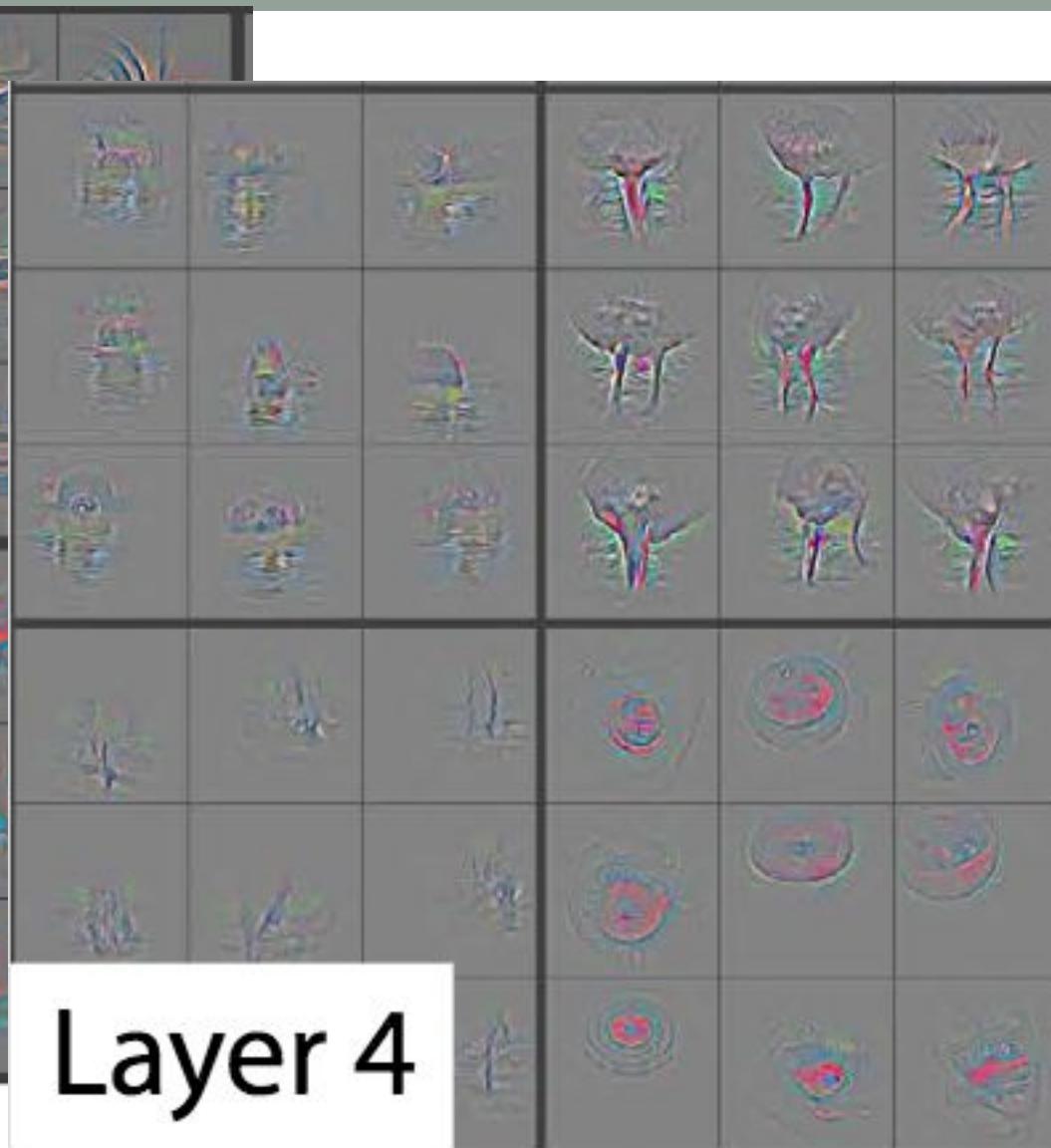


Higher layer captures higher-level
concepts



Layer 1

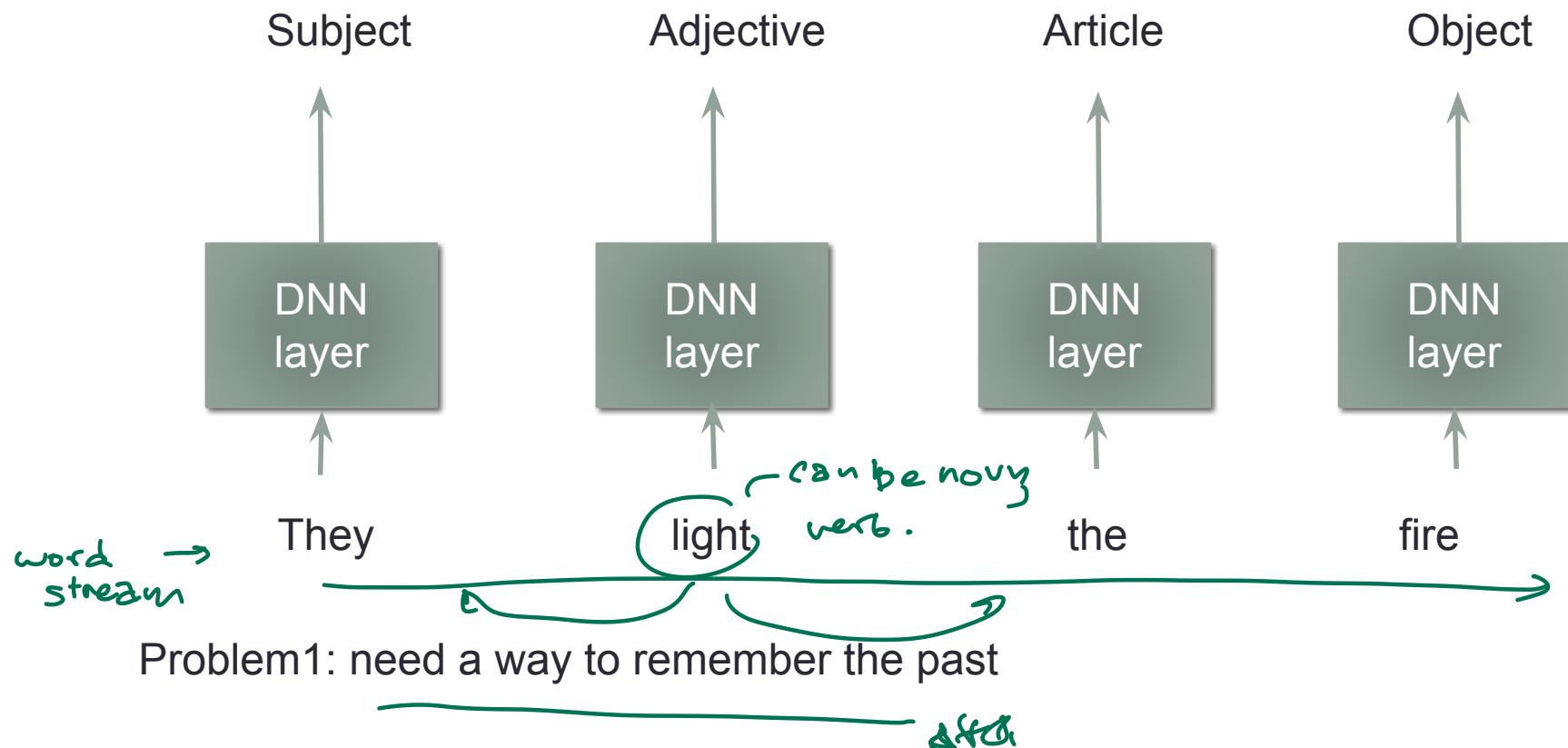




Recurrent neural network (RNN)

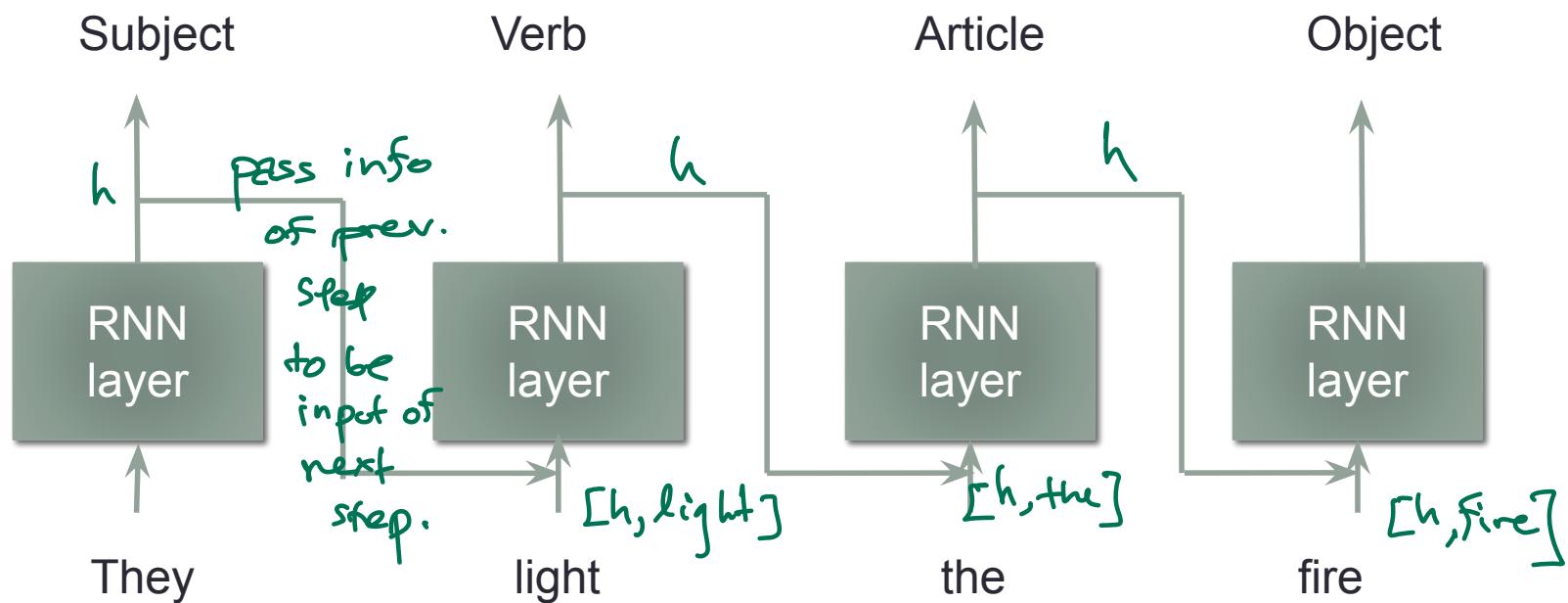
- DNN framework

Pos.



Recurrent neural network (RNN)

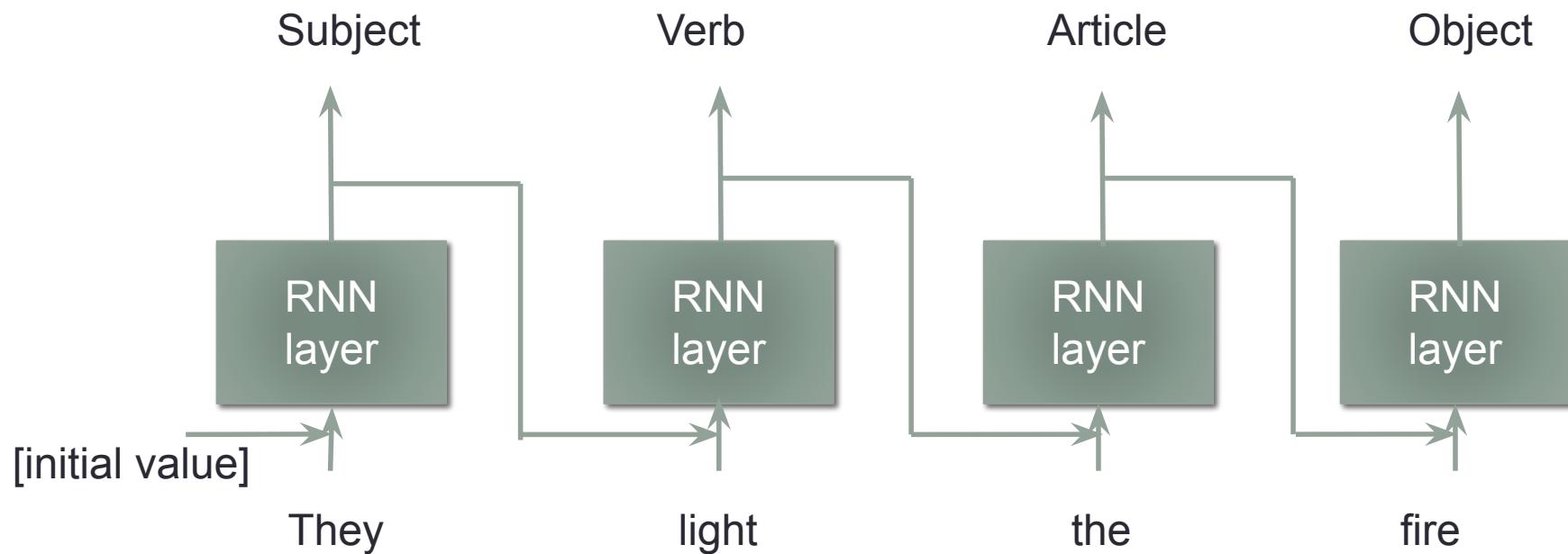
- RNN framework



Output of the layer encodes something meaningful about the past

Recurrent neural network (RNN)

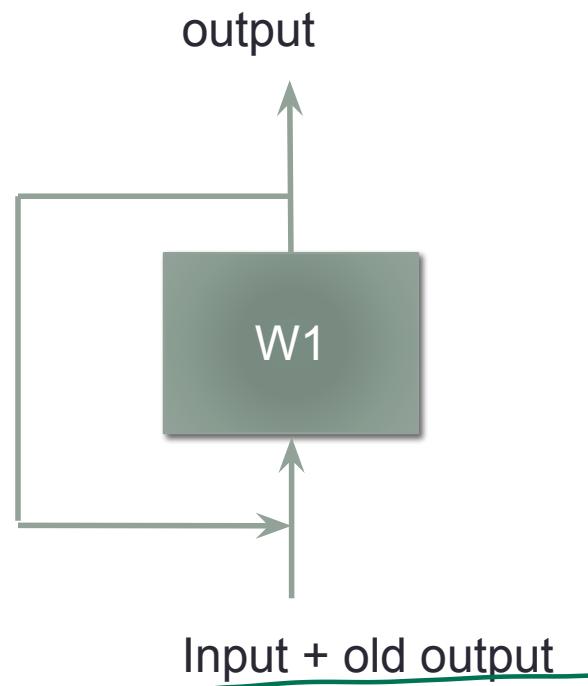
- RNN framework



New input feature = [original input feature, output of the layer at previous time step]

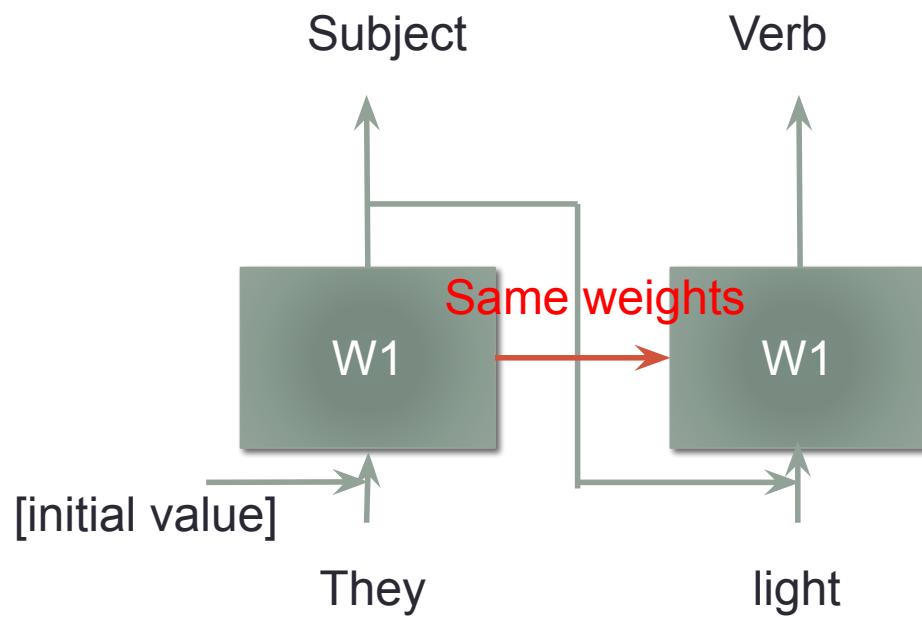
Recurrent neural network (RNN)

- Unrolling of a recurrent layer.



Recurrent neural network (RNN)

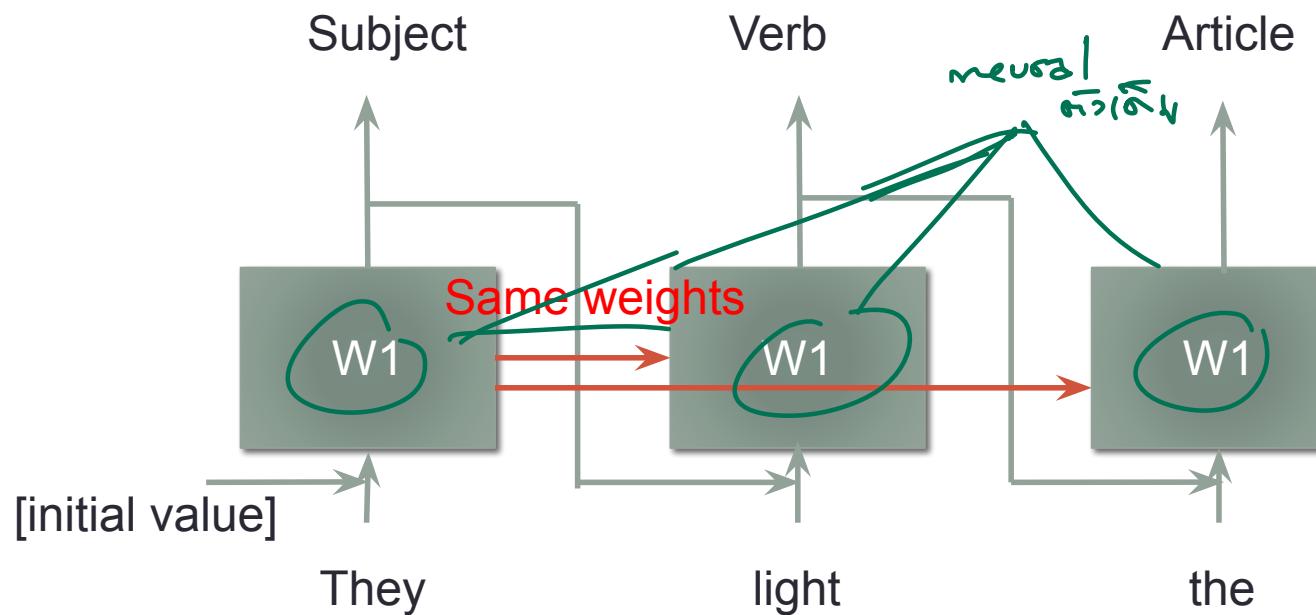
- Unrolling of a recurrent layer.



Parameter sharing

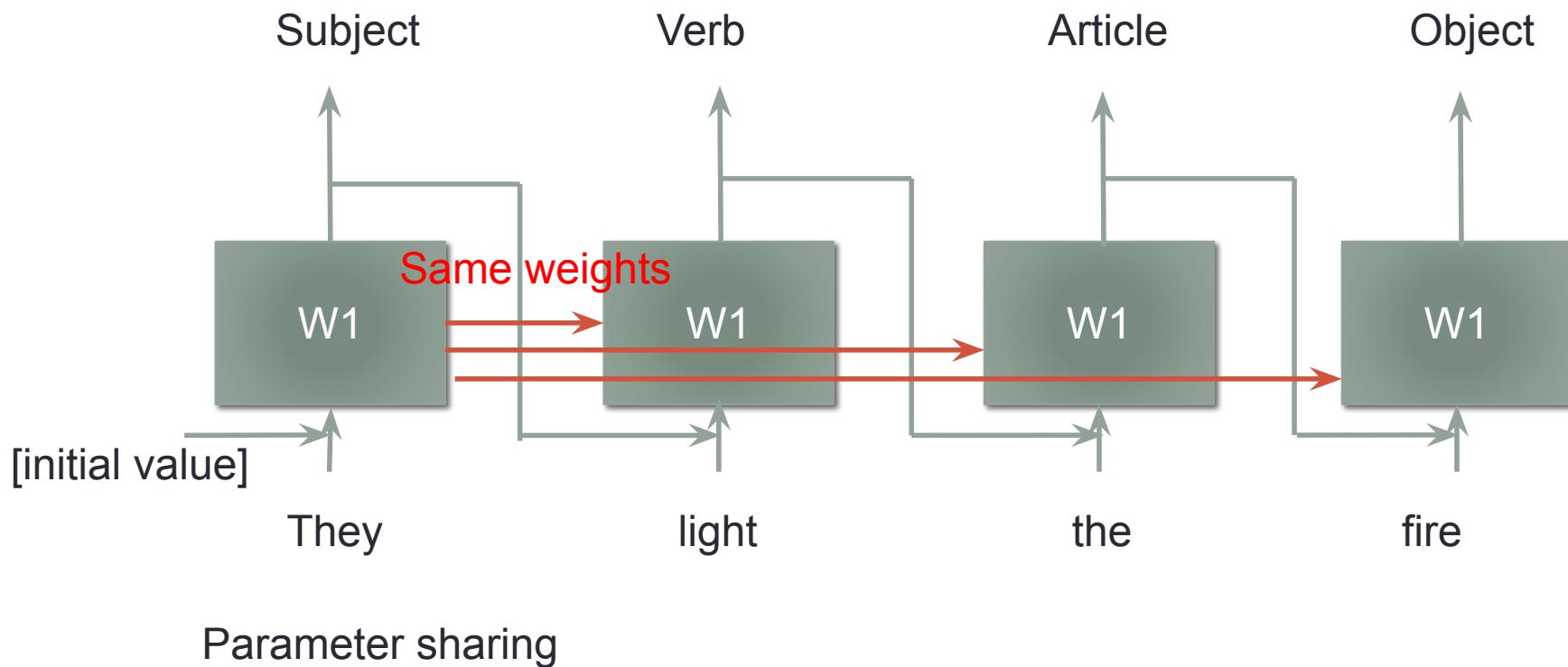
Recurrent neural network (RNN)

- Unrolling of a recurrent layer.



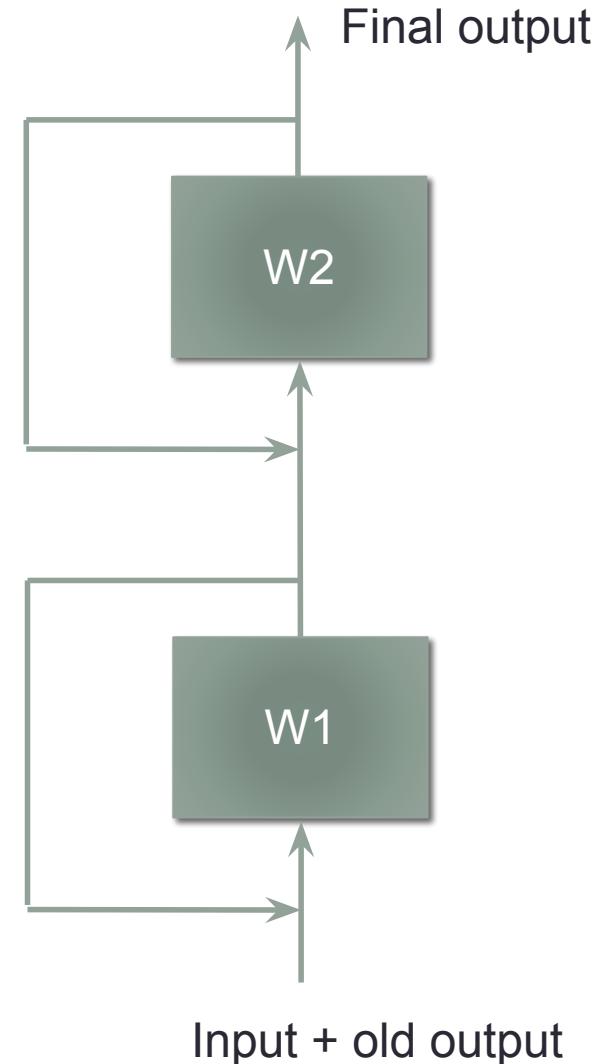
Recurrent neural network (RNN)

- Unrolling of a recurrent layer.



Recurrent neural network (RNN)

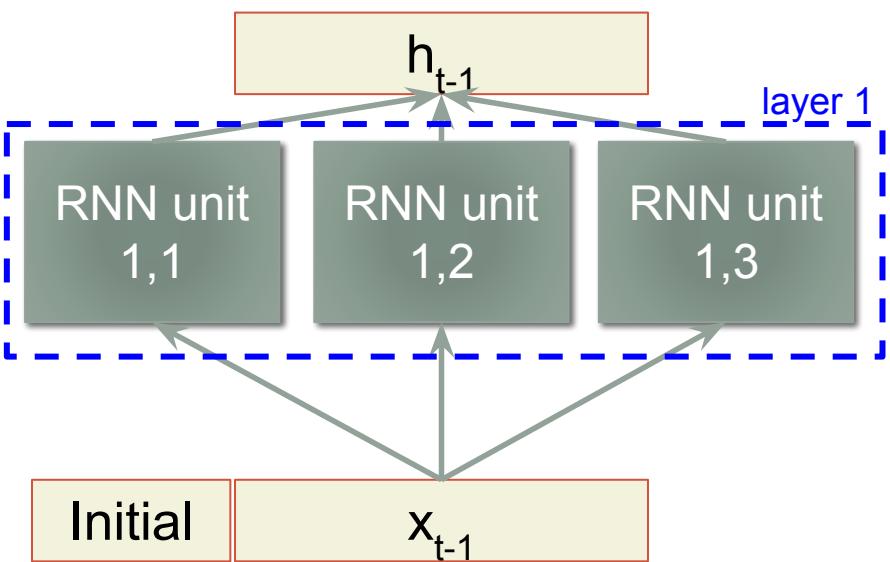
- Stacks of recurrent layer



RNN layers (expanded in time)

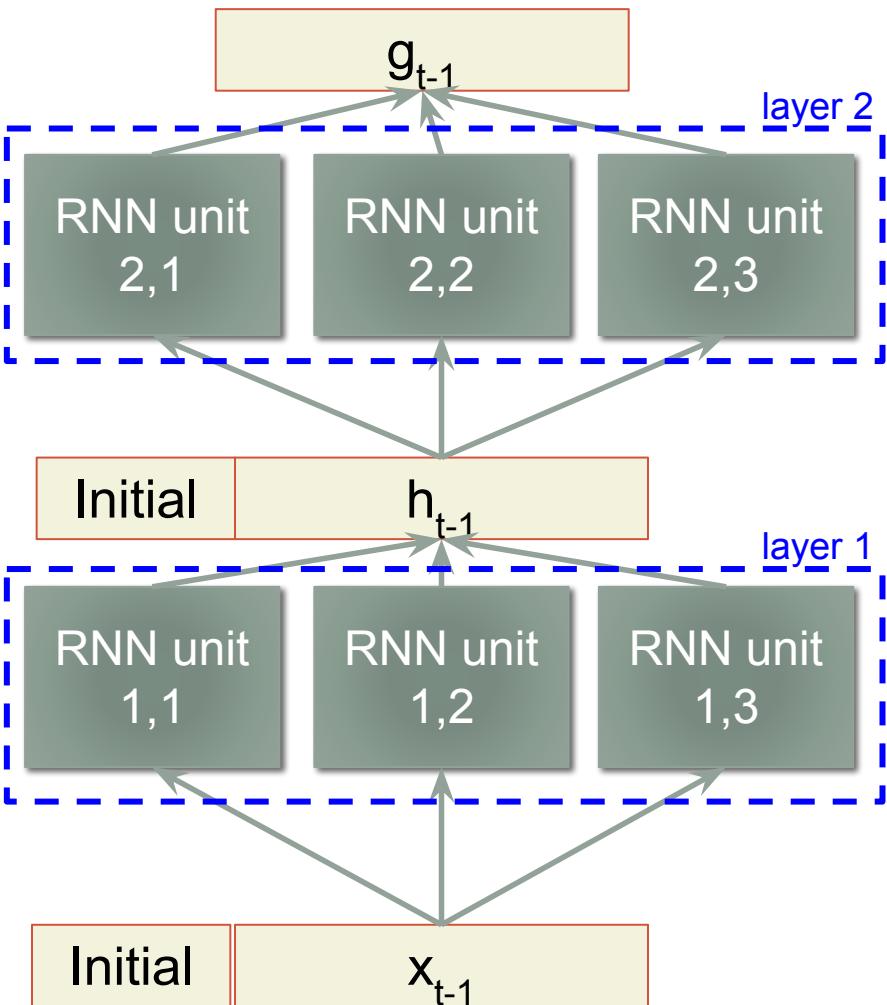
Time step 1

Time step 2



RNN layers (expanded in time)

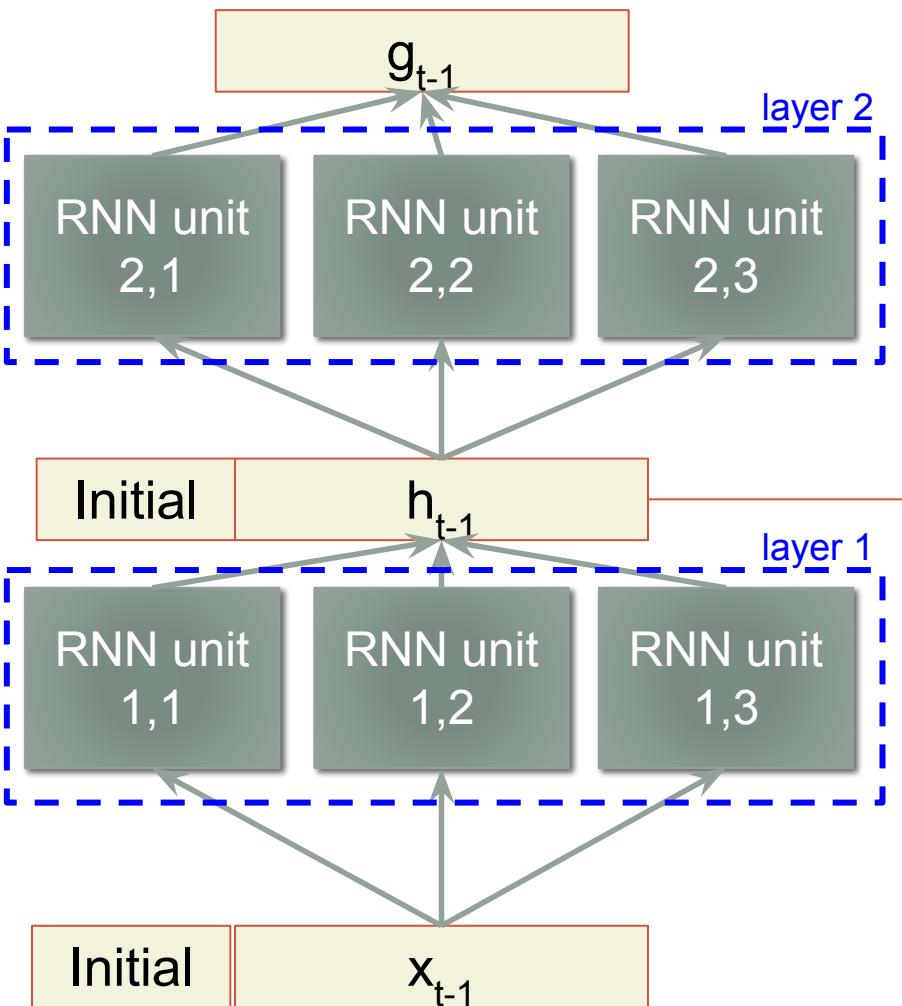
Time step 1



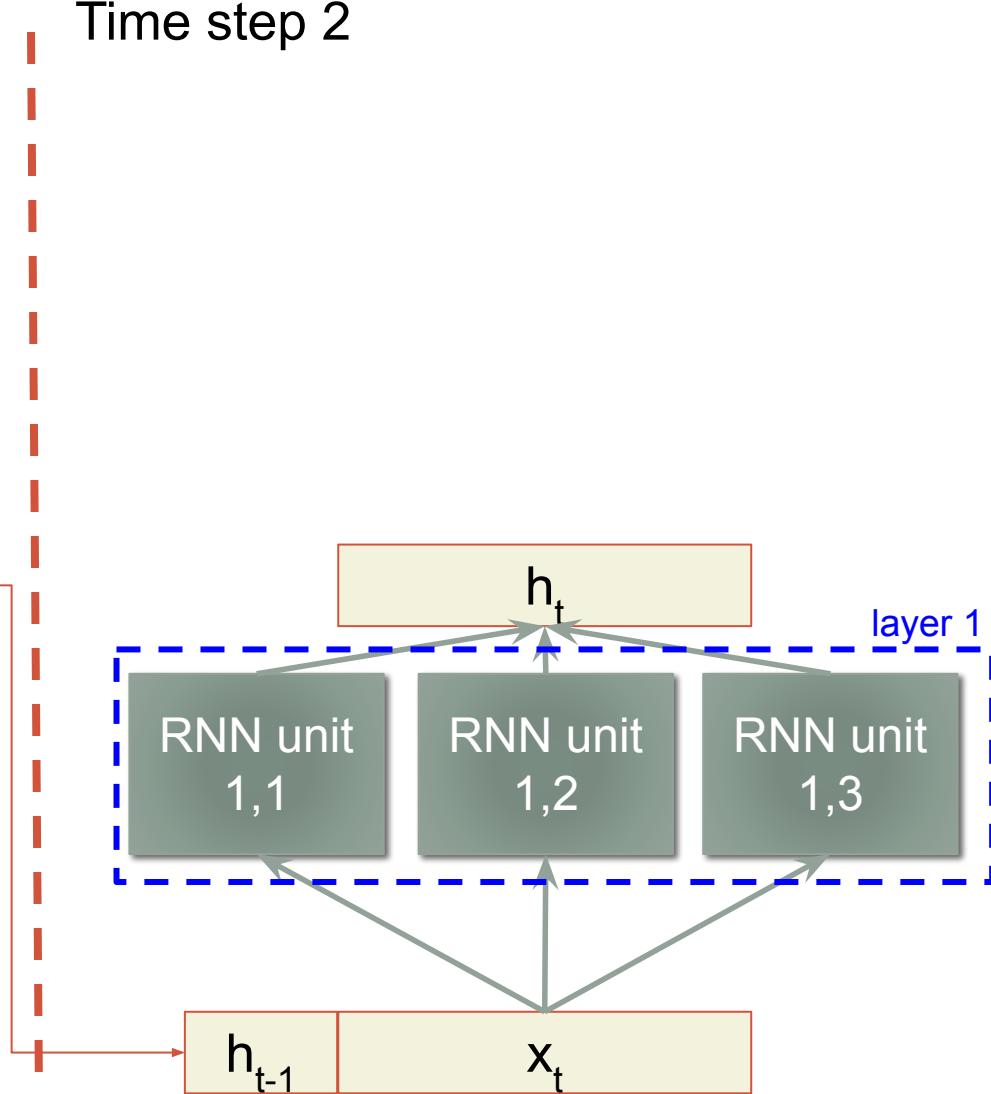
Time step 2

RNN layers (expanded in time)

Time step 1

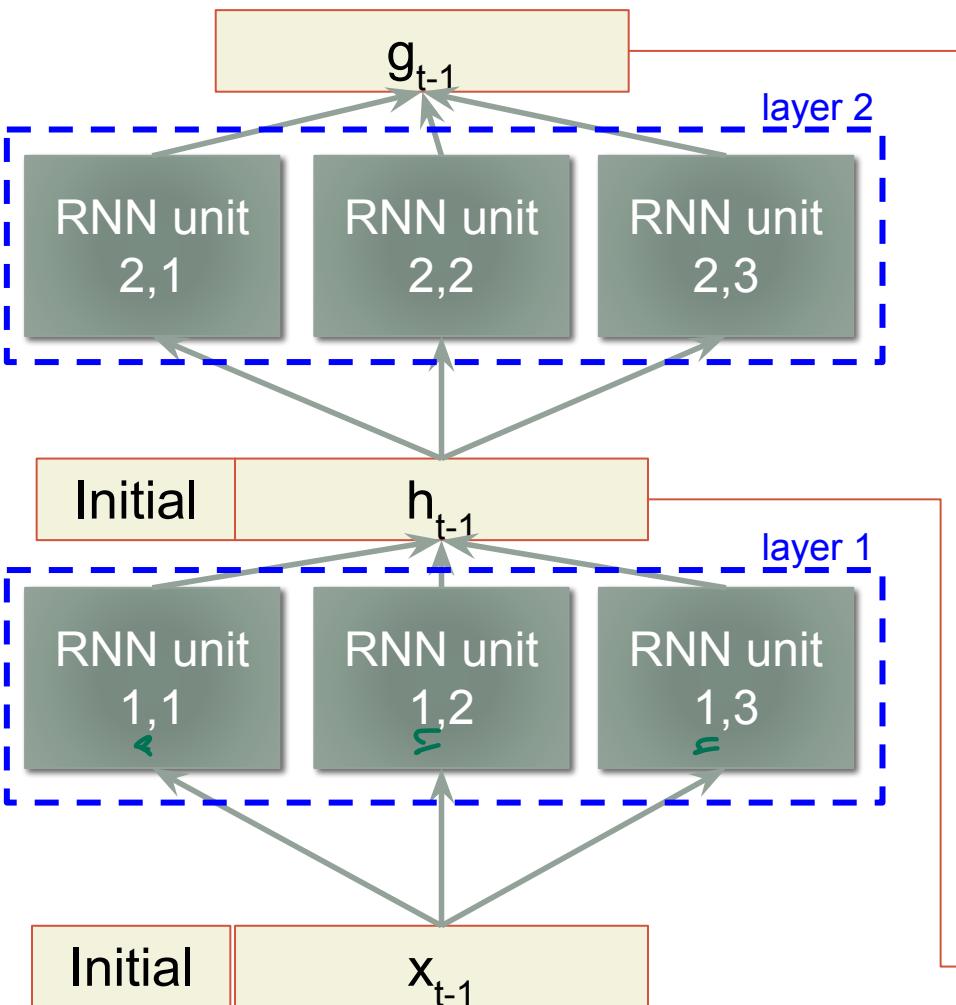


Time step 2

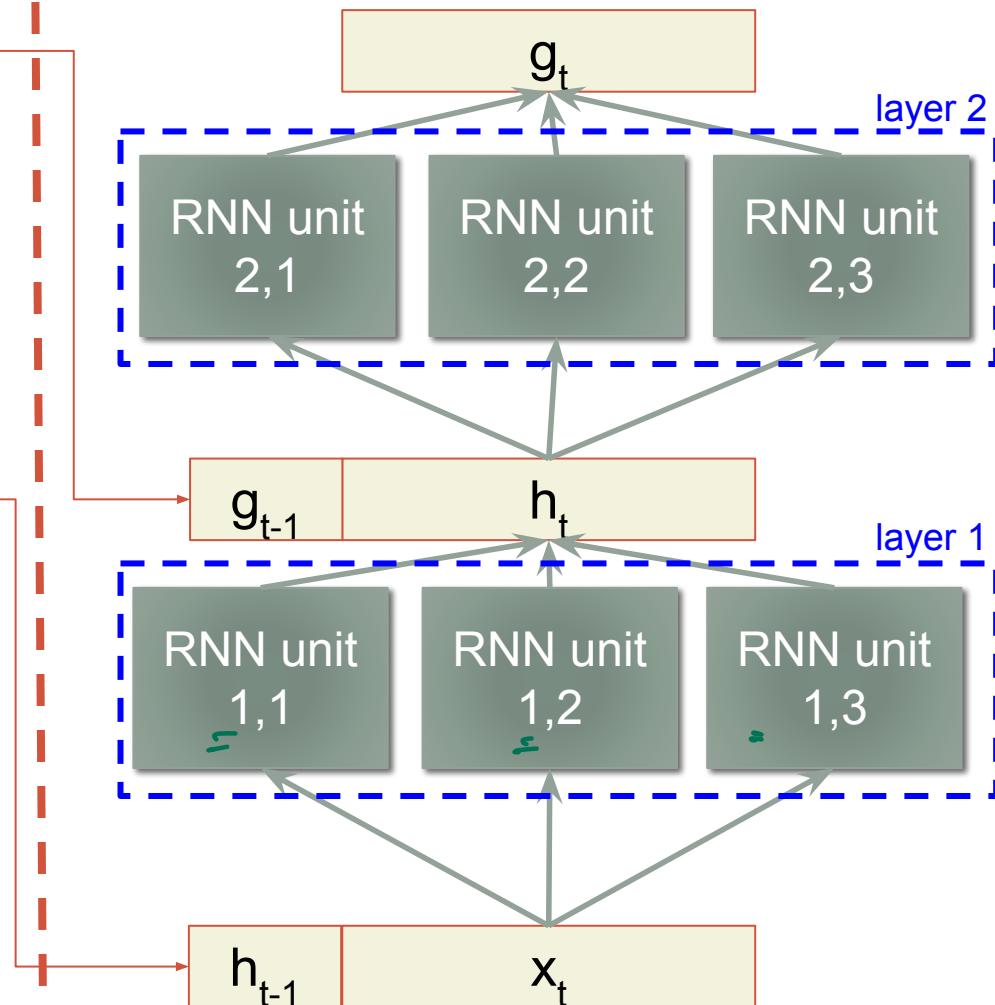


RNN layers (expanded in time)

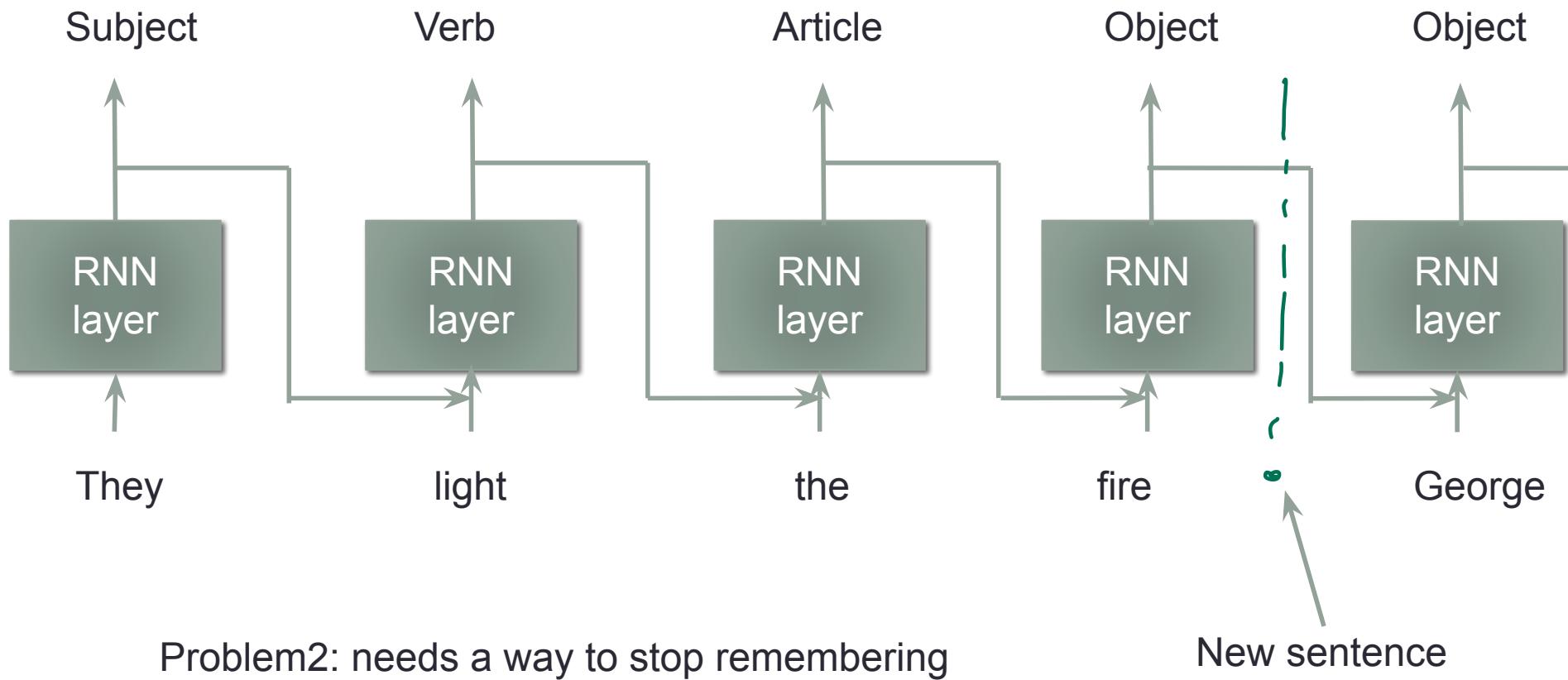
Time step 1



Time step 2



Recurrent neural network (RNN)

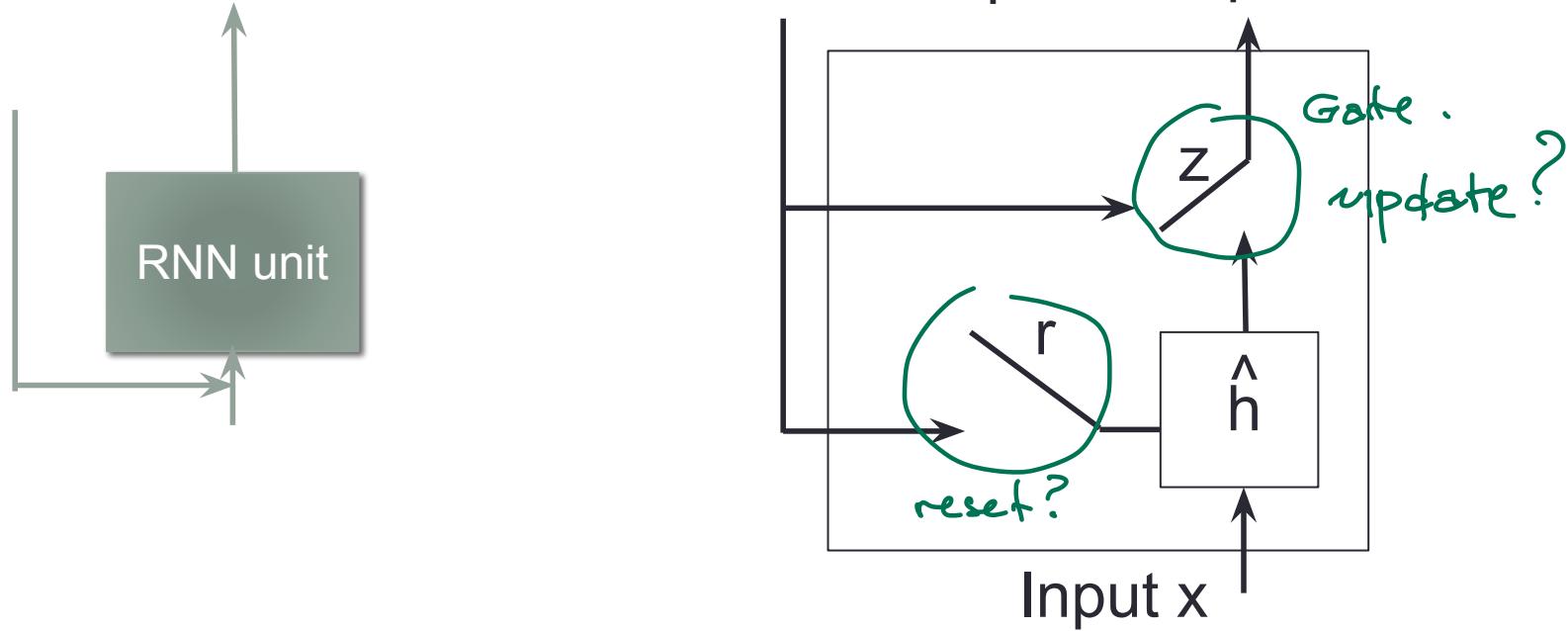


Can the network learn when to start and stop remembering things? 

ឧបនគរណីក្នុងការបង្កើតរបស់ខ្លួន

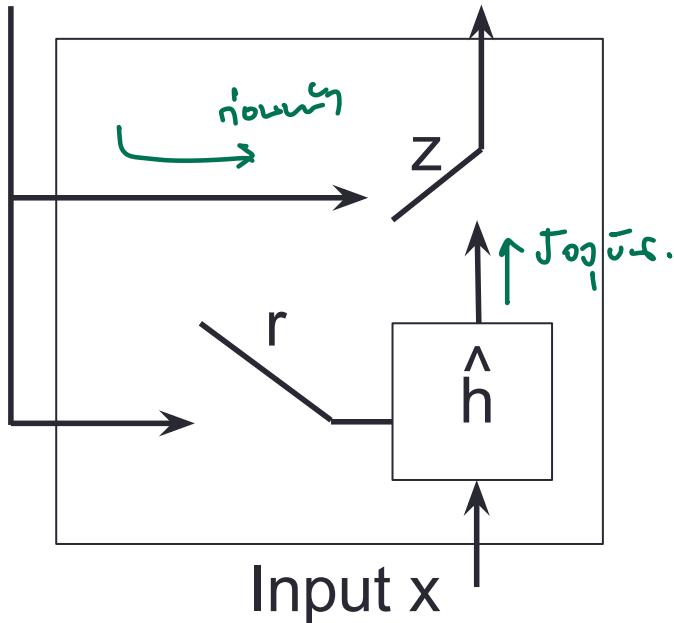
Gated Recurrent Unit (GRU)

- Forms a Gated Recurrent Neural Networks (GRNN)
- Add gates that can choose to reset (r) or update (z)



Gated Recurrent Unit (GRU)

Previous output Output h

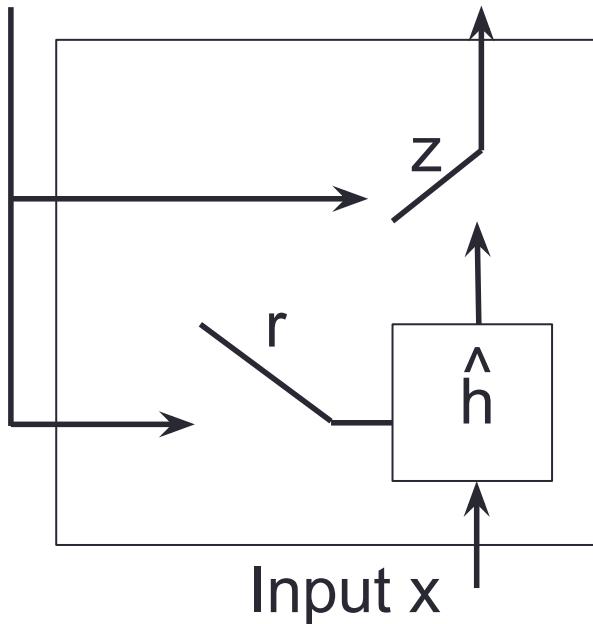


Neuron index j time index t

$$h_t^j = (1 - z_t^j)h_{t-1}^j + z_t^j \hat{h}_t^j$$

Gated Recurrent Unit (GRU)

Previous output Output history



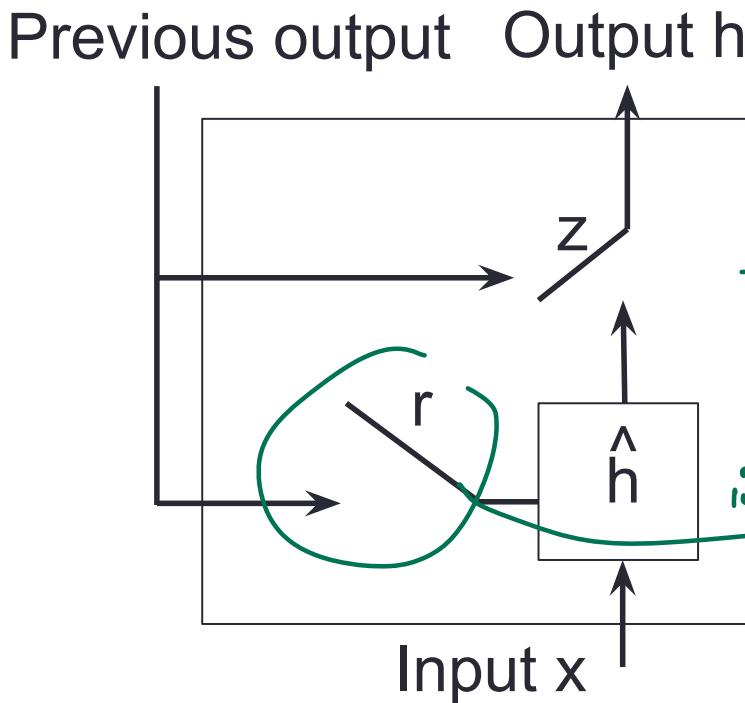
$$h_t^j = (1 - z_t^j)h_{t-1}^j + z_t^j \hat{h}_t^j$$

One GRU neuron output (scalar)

Gated Recurrent Unit (GRU)

$h(wx)$

$w[x]$

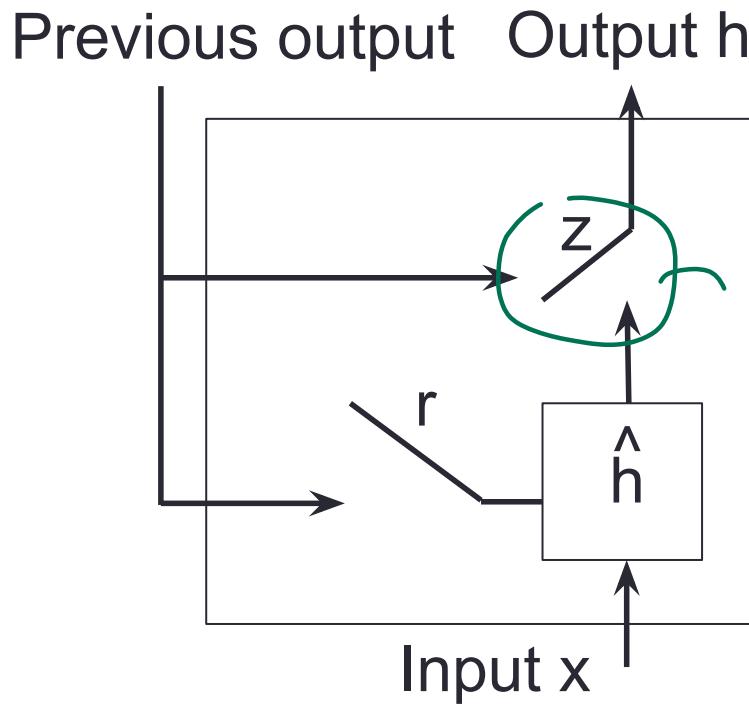


$$h_t^j = (1 - z_t^j)h_{t-1}^j + z_t^j \hat{h}_t^j$$
$$\hat{h}_t^j = \frac{\tanh^j(W_s x_t + U(\text{non-lin. Linear transform with matrix multiply})(r_t \odot h_{t-1})))}{\text{Element-wise product}}$$

Annotations:

- z_t^j : neural net.
- \tanh : $[-1, 1]$
- W_s : x_t
- U : $(r_t \odot h_{t-1})$
- $\text{non-lin. Linear transform with matrix multiply}$
- $\text{Element-wise product}$
- non-lin.
- $x_t^j = h_t^j$
- $\text{Vector (each value from each GRU unit in the previous layer)}$

Gated Recurrent Unit (GRU)



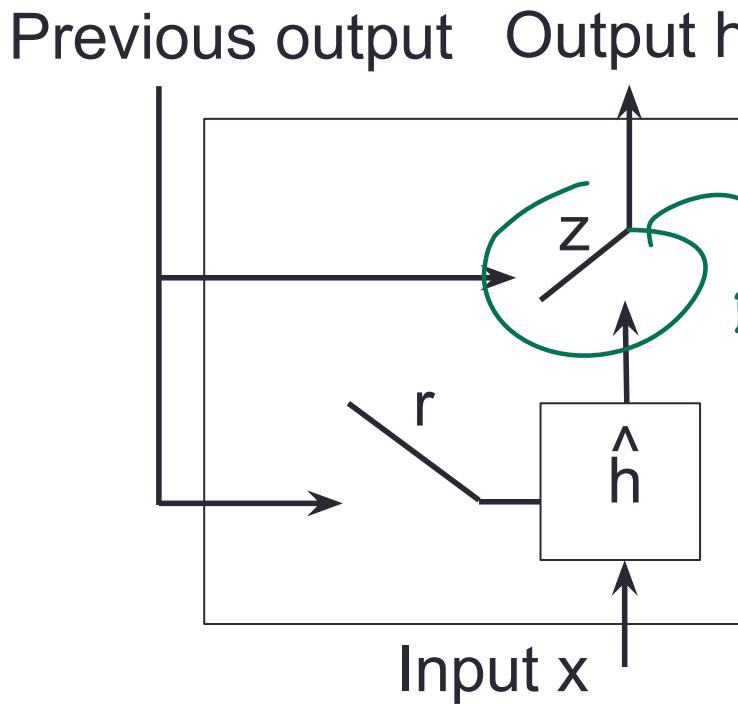
$$h_t^j = (1 - z_t^j)h_{t-1}^j + z_t^j \hat{h}_t^j$$

$$\hat{h}_t^j = \underline{\tanh^j}(W\mathbf{x}_t + U(\mathbf{r}_t \odot \mathbf{h}_{t-1}))$$

Takes the j-th element

linear combination → [0, 1]

Gated Recurrent Unit (GRU)



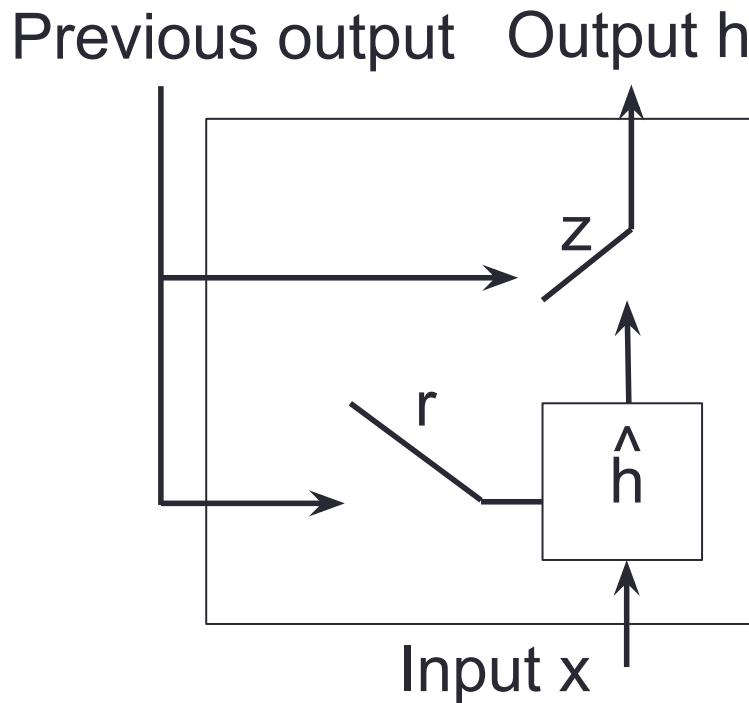
$$h_t^j = (1 - z_t^j)h_{t-1}^j + z_t^j \hat{h}_t^j$$

$$\hat{h}_t^j = \tanh^j(W\mathbf{x}_t + U(\mathbf{r}_t \odot \mathbf{h}_{t-1}))$$

$$z_t^j = \underline{\text{sigmoid}}^j(W_z \mathbf{x}_t + U_z \mathbf{h}_{t-1})$$

Indicates a different set of weights

Gated Recurrent Unit (GRU)



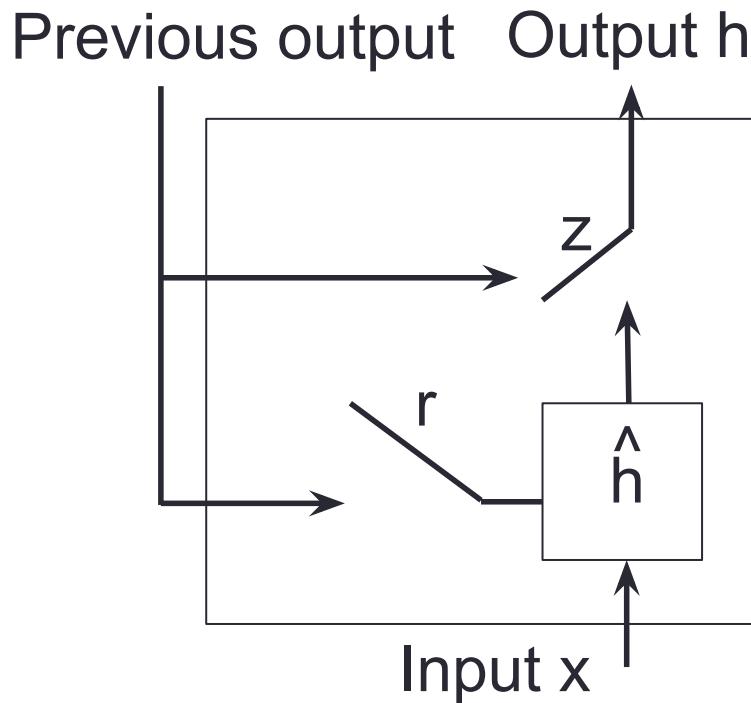
$$h_t^j = (1 - z_t^j)h_{t-1}^j + z_t^j \hat{h}_t^j$$

$$\hat{h}_t^j = \tanh^j(W\mathbf{x}_t + U(\mathbf{r}_t \odot \mathbf{h}_{t-1}))$$

$$z_t^j = \text{sigmoid}^j(W_z \mathbf{x}_t + U_z \mathbf{h}_{t-1})$$

Bounds the output to 0 to 1 for interpolation

Gated Recurrent Unit (GRU)



$$h_t^j = (1 - z_t^j)h_{t-1}^j + z_t^j \hat{h}_t^j$$

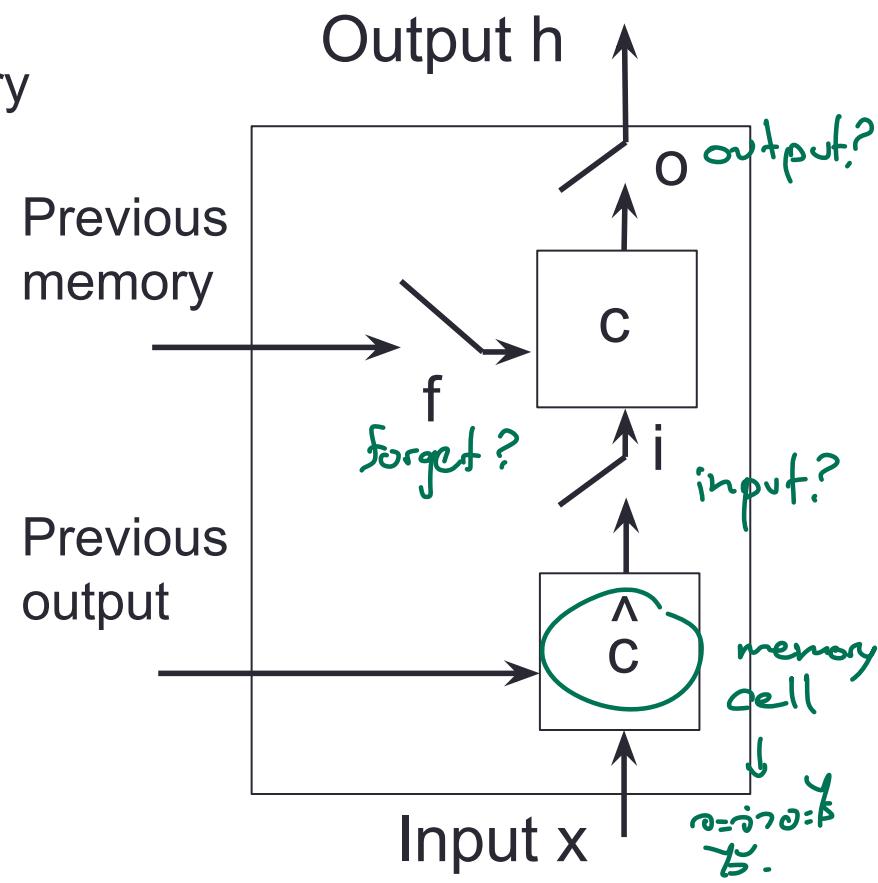
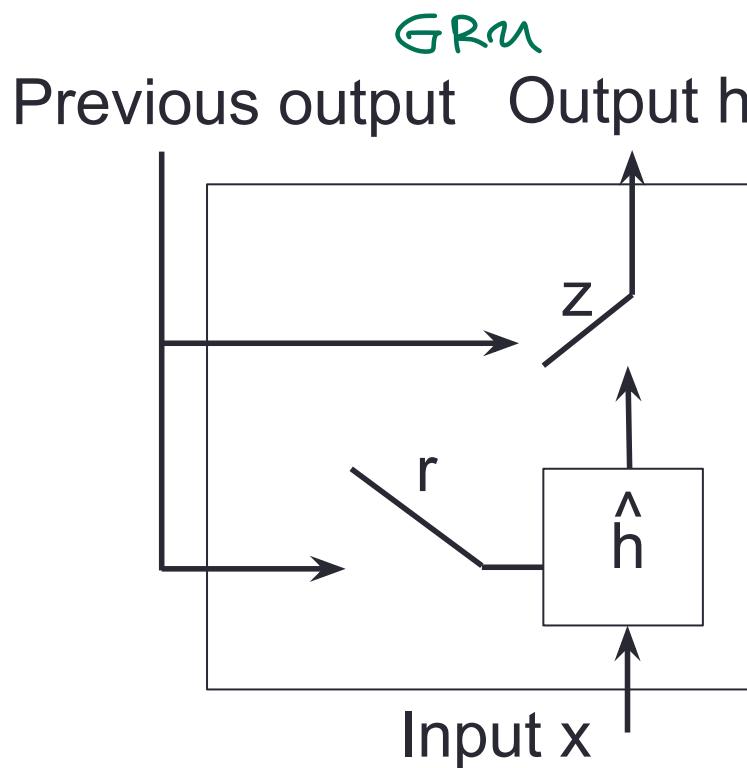
$$\hat{h}_t^j = \tanh^j(W\mathbf{x}_t + U(\mathbf{r}_t \odot \mathbf{h}_{t-1}))$$

$$z_t^j = \underline{\text{sigmoid}}^j(W_z \mathbf{x}_t + U_z \mathbf{h}_{t-1})$$

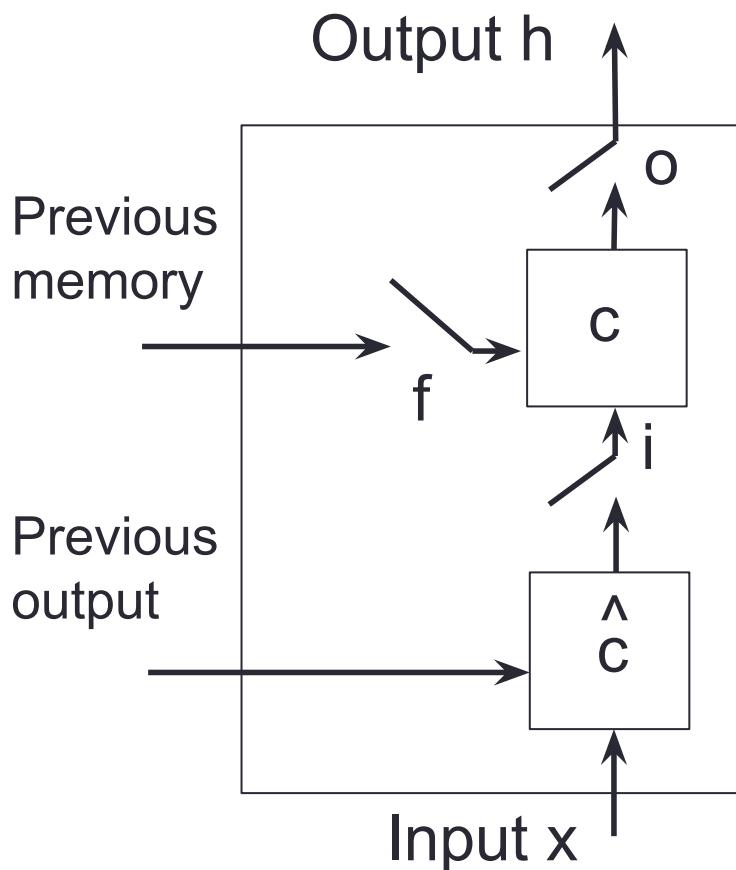
$$r_t^j = \underline{\text{sigmoid}}^j(W_r \mathbf{x}_t + U_r \mathbf{h}_{t-1})$$

Long Short-Term Memory (LSTM)

- Have 3 gates, forget (f), input (i), output (o)
- Has an **explicit memory cell** (c)
 - Does not have to output the memory



Long Short-Term Memory (LSTM)



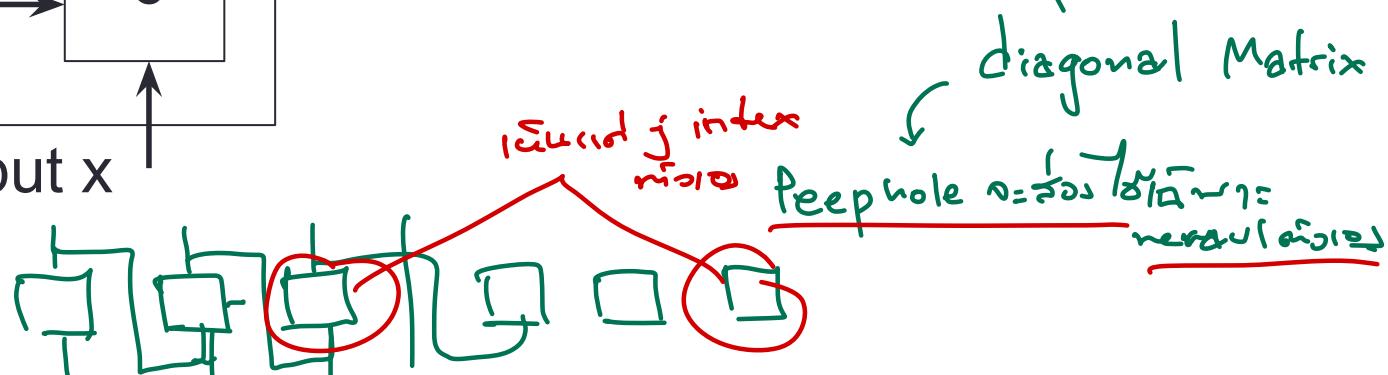
$$i_t^j = F^j(W_i \mathbf{x}_t + U_i \mathbf{h}_{t-1} + V_i \mathbf{c}_{t-1})$$

$$o_t^j = F^j(W_o \mathbf{x}_t + U_o \mathbf{h}_{t-1} + V_o \mathbf{c}_t)$$

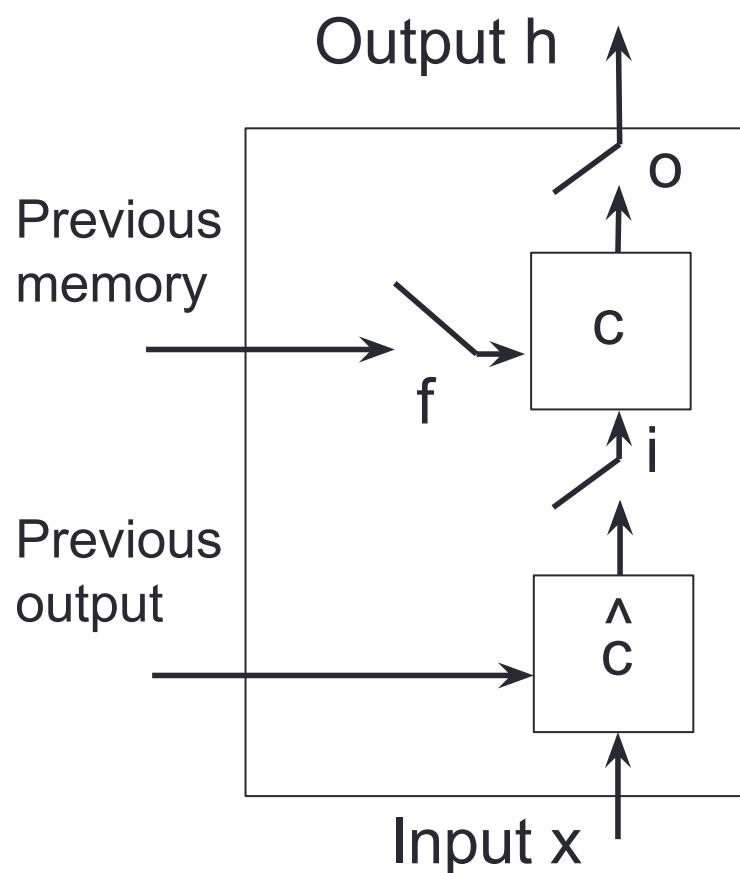
$$f_t^j = F^j(W_f \mathbf{x}_t + U_f \mathbf{h}_{t-1} + V_f \mathbf{c}_{t-1})$$

Contribution from memory “Peephole connection”

V s are diagonal matrices (Each cell can only see its own memory)



Long Short-Term Memory (LSTM)



$$i_t^j = F^j(W_i \mathbf{x}_t + U_i \mathbf{h}_{t-1} + V_i \mathbf{c}_{t-1})$$

$$o_t^j = F^j(W_o \mathbf{x}_t + U_o \mathbf{h}_{t-1} + V_o \mathbf{c}_t)$$

$$f_t^j = F^j(W_f \mathbf{x}_t + U_f \mathbf{h}_{t-1} + V_j \mathbf{c}_{t-1})$$

$$h_t^j = o_t^j \tanh(c_t^j)$$

$$c_t^j = f_t^j c_{t-1}^j + i_t^j \hat{c}_t^j$$

$$\hat{c}_t^j = \underbrace{\tanh^j}_{\text{non-linear}} \left(\underbrace{W_c \mathbf{x}_t + U_c \mathbf{h}_{t-1}}_{\text{neural}} \right)$$

GRU vs LSTM

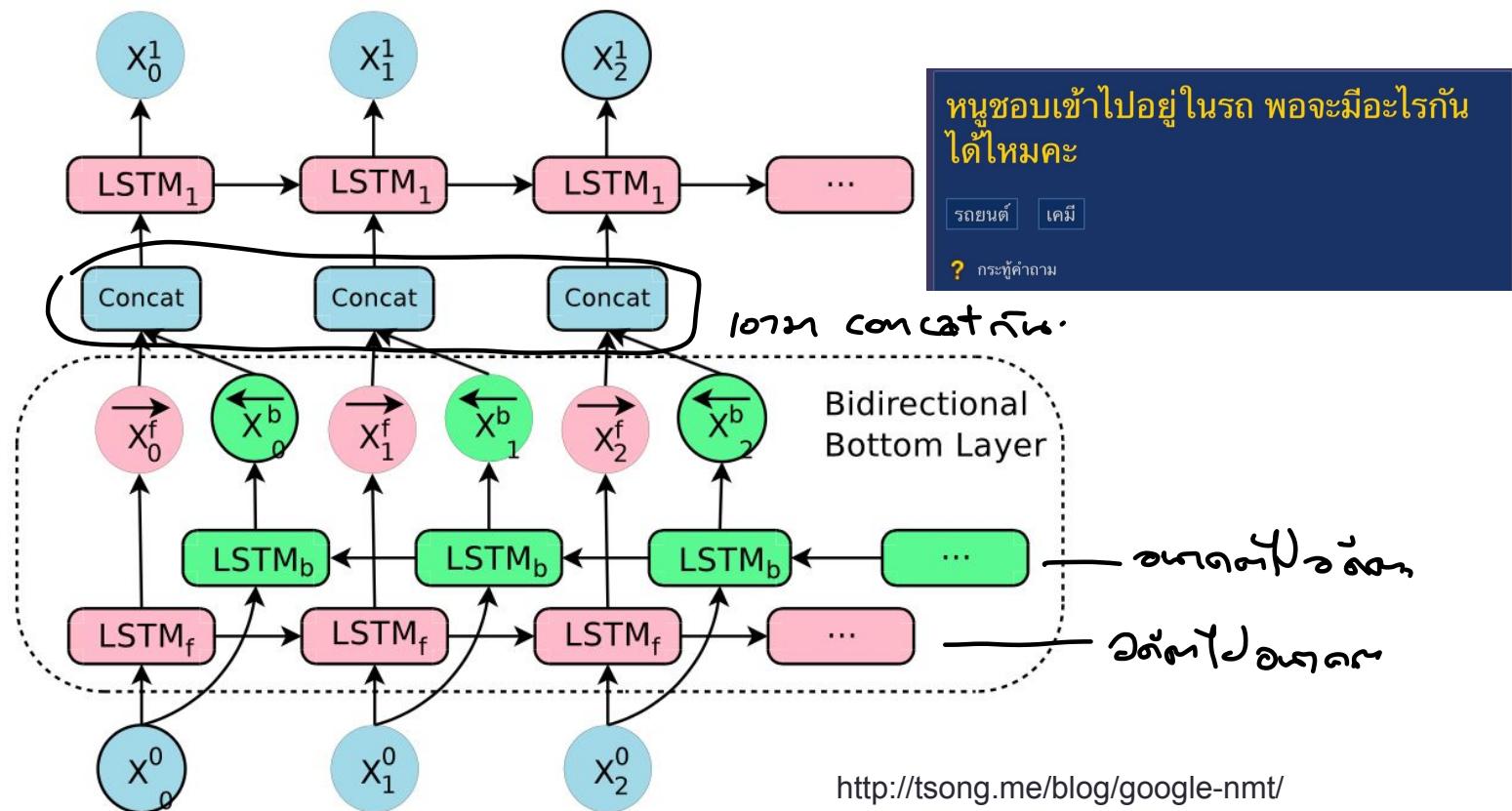
- GRU and LSTM offers the same performance with large dataset
 - GRU better for smaller dataset (less parameters)
 - GRU faster to train and faster runtime (smaller model)
- Use GRUs!

except condition that u want "C" → some model
need to use C value

↓
Seq2Seq

Bi-directional LSTM

- The previous GRU/LSTM only goes backward in time (uni-directional)
- Most of the time information from the future is useful for predicting the current output



LSTM remembers meaningful things

(cf value)

សុវត្ថិភាព. → និរាយការណ៍រាយ
 ↑ ↑
 T h e s □

1000 1 0 0 0 0 - 0
 neural next → និរាយការណ៍រាយ
 layer

Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.

Cell that turns on inside quotes:

"You mean to imply that I have nothing to eat out of... On the contrary, I can supply you with everything even if you want to give dinner parties" warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

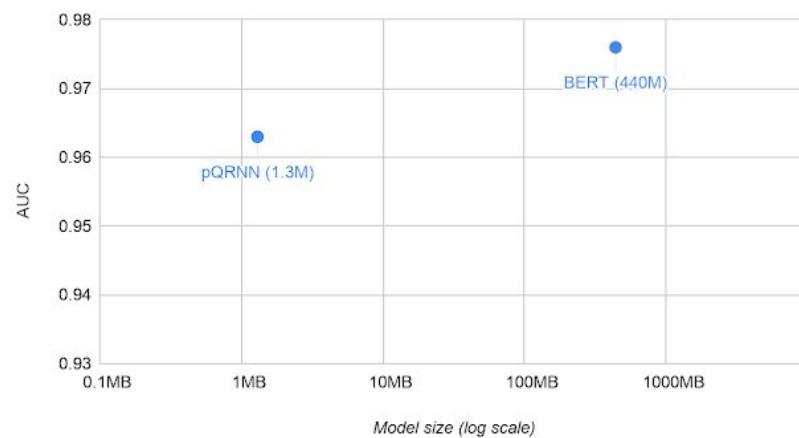
Notes on RNNs

→ "_rnn"

Considered obsolete by some due to speed and performance

- Improved speed using QRNN <https://arxiv.org/abs/1611.01576>
- Can be as good as BERTs on simple tasks

Model Comparison



with less size.



<https://ai.googleblog.com/2020/09/advancing-nlp-with-efficient-projection.html>

Easier to train than attention-based models (easier to tune and faster convergence). Less memory requirements.

Embeddings

- A way to encode information to a lower dimensional space
 - PCA
 - We learn about this lower dimensional space through data

Hw.

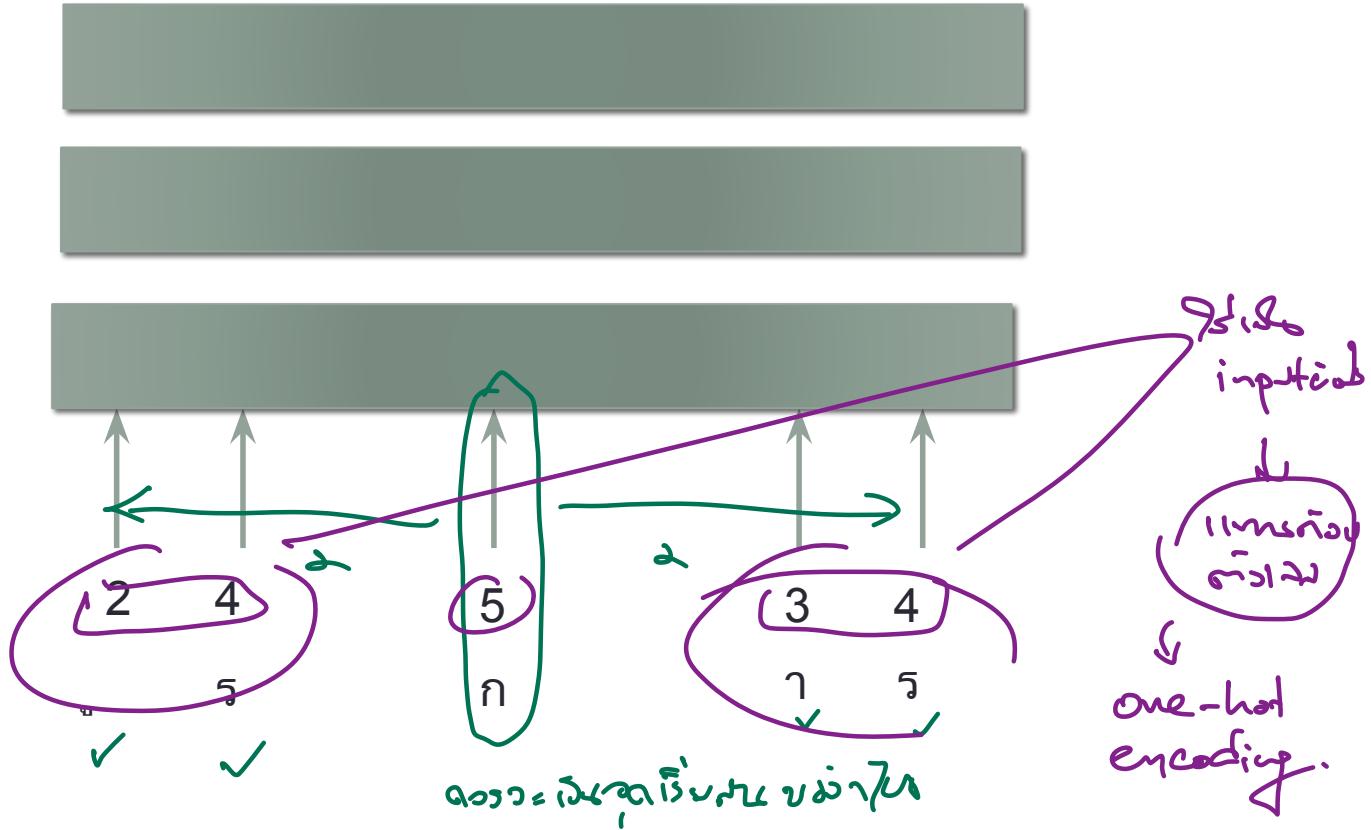
Word segmentation with fully connected networks



1 = word beginning, 0 = word middle



Logistic function $o_{j,1}$



One hot encoding

- Categorical representation is usually represented by one hot encoding
- Categorical representations examples:
 - Words in a vocabulary, characters in Thai language

Apple -> 1 -> [1, 0, 0, 0, ...]

Bird -> 2 -> [0, 1, 0, 0, ...]

Cat -> 3 -> [0, 0, 1, 0, ...]

} change word vector.
1 in one position

- Sparse representation

- Spare means most dimensions are zero

invis.

One hot encoding

- Sparse – but lots of dimension
 - Curse of dimensionality
- Does not represent meaning.

Apple -> 1 -> [1, 0, 0, 0, ...]

Bird -> 2 -> [0, 1, 0, 0, ...]

Cat -> 3 -> [0, 0, 1, 0, ...]

$$|\text{Apple} - \text{Bird}| = |\text{Bird} - \text{Cat}|$$

To encode/unt distance
between!
Bird ~ cat more
proximity.

Getting meaning into the feature vectors

- You can add back meanings by hand-crafted rules
 - Old-school NLP is all about feature engineering
 - Word segmentation example:

- Cluster Numbers
 - Cluster letters
 - Concatenate the

میں جسے

- Concatenate them

$$\bullet \cdot 1 = [0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0, 1, 0]$$

$$\bullet \eta = [0\ 0\ 0\ 1\ 0\ 0\ 0\ 0, 0, 1]$$

- Which rules to use?

- Try as many as you can think of, and do feature selection or use models that can do feature selection

Dense representation

- We can encode sparse representation into a lower dimensional space
 - $F: \mathbb{R}^N \rightarrow \mathbb{R}^M$, where $N > M$

meaningful

Apple $\rightarrow 1 \rightarrow [1, 0, 0, 0, \dots] \rightarrow [2.3, 1.2]$

Bird $\rightarrow 2 \rightarrow [0, 1, 0, 0, \dots] \rightarrow [-1.0, 2.4]$

Cat $\rightarrow 3 \rightarrow [0, 0, 1, 0, \dots] \rightarrow [-3.0, 4.0]$

Θ Susan \rightarrow neural network
layer.

one-hot \rightarrow dense representation

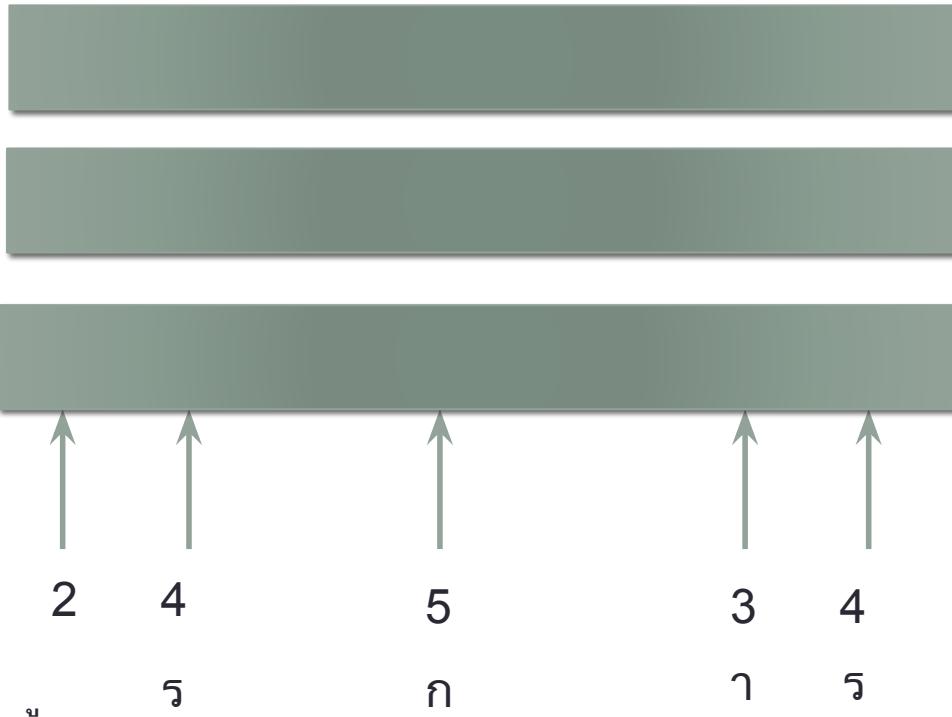
```
graph LR; A[Apple] --> 1[1]; 1 --> S1["[1, 0, 0, 0, ...]"]; S1 --> D1["[2.3, 1.2]"]; B[Bird] --> 2[2]; 2 --> S2["[0, 1, 0, 0, ...]"]; S2 --> D2["[-1.0, 2.4]"]; C[Cat] --> 3[3]; 3 --> S3["[0, 0, 1, 0, ...]"]; S3 --> D3["[-3.0, 4.0]"]; S1 --- S2 --- S3; D1 --- D2 --- D3; S1 --- E["embedding layer"]; D1 --- NN["neural network layer."]; S1 --- OR["one-hot → dense representation"];
```

Word segmentation with fully connected networks

1 = word beginning, 0 = word middle



Logistic function

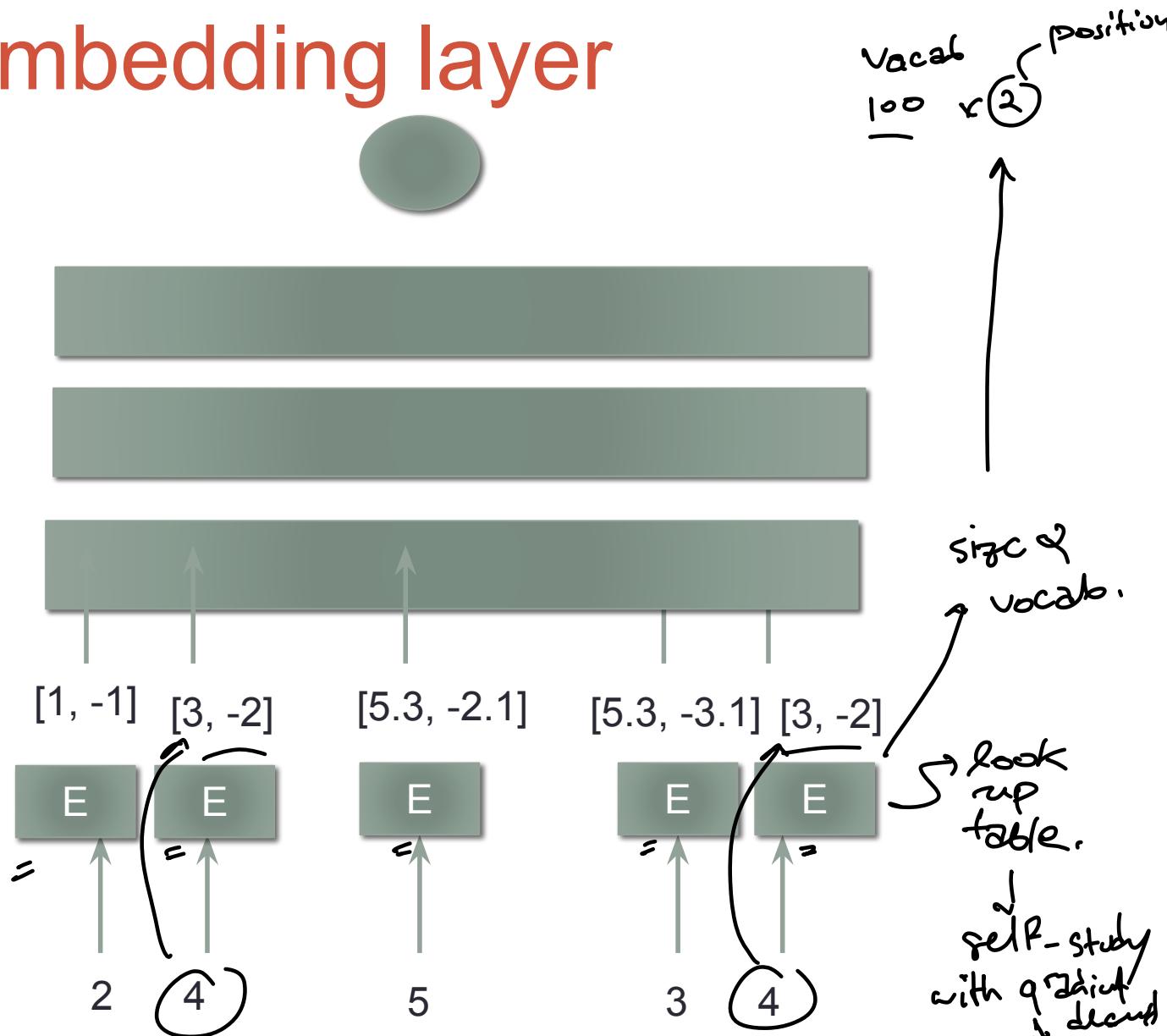


Adding embedding layer

Embedding layer
shares the same
weights

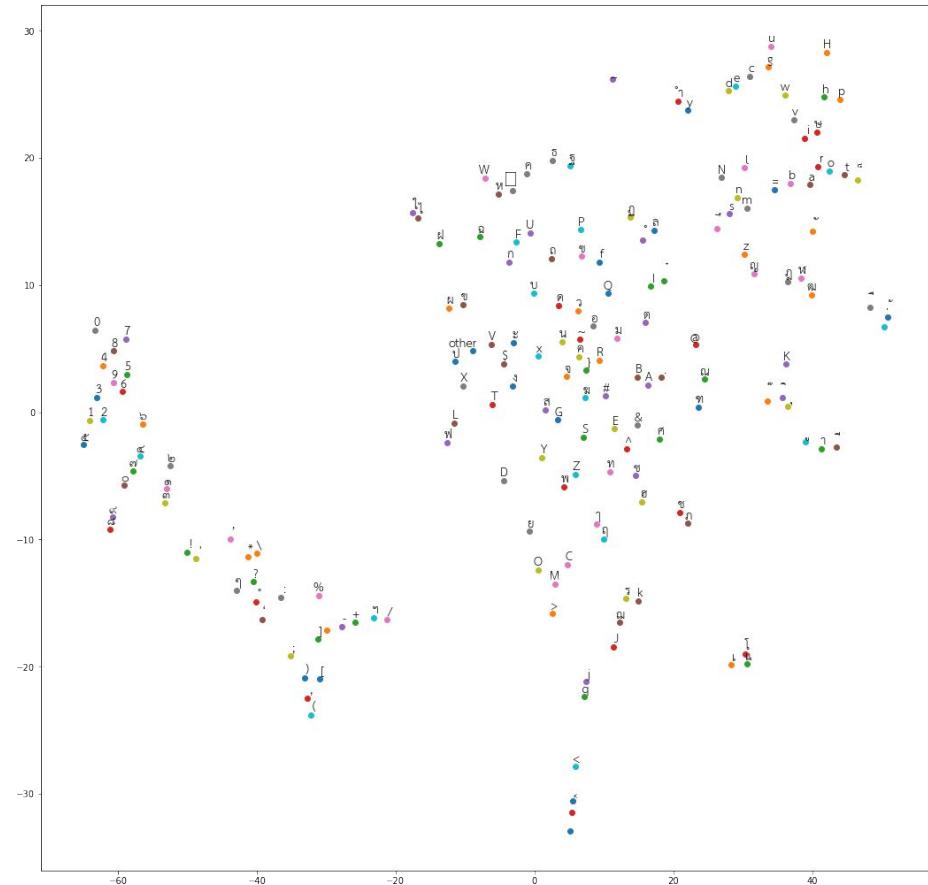
Parameter sharing!

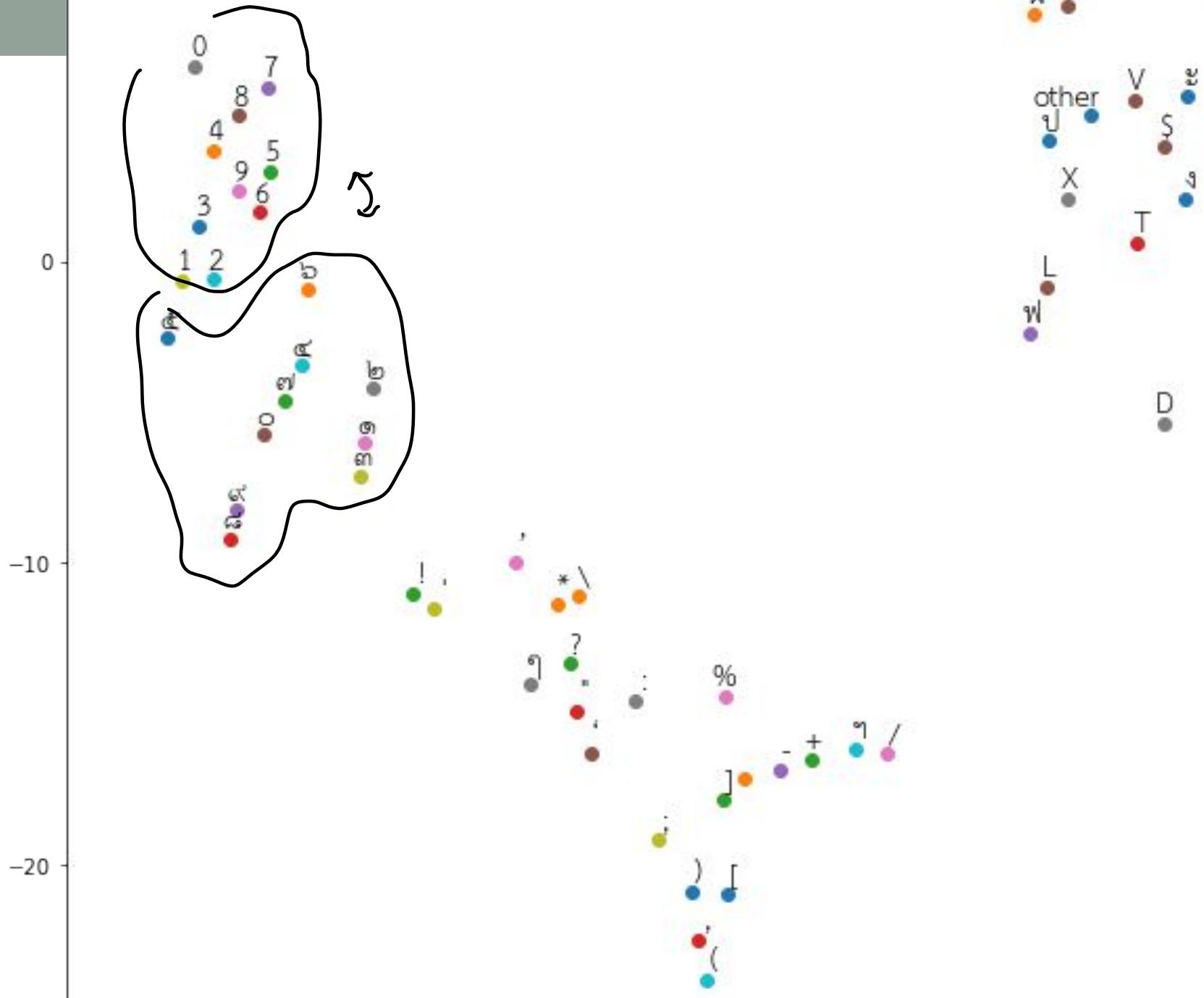
More on embeddings
in the next two
lectures!



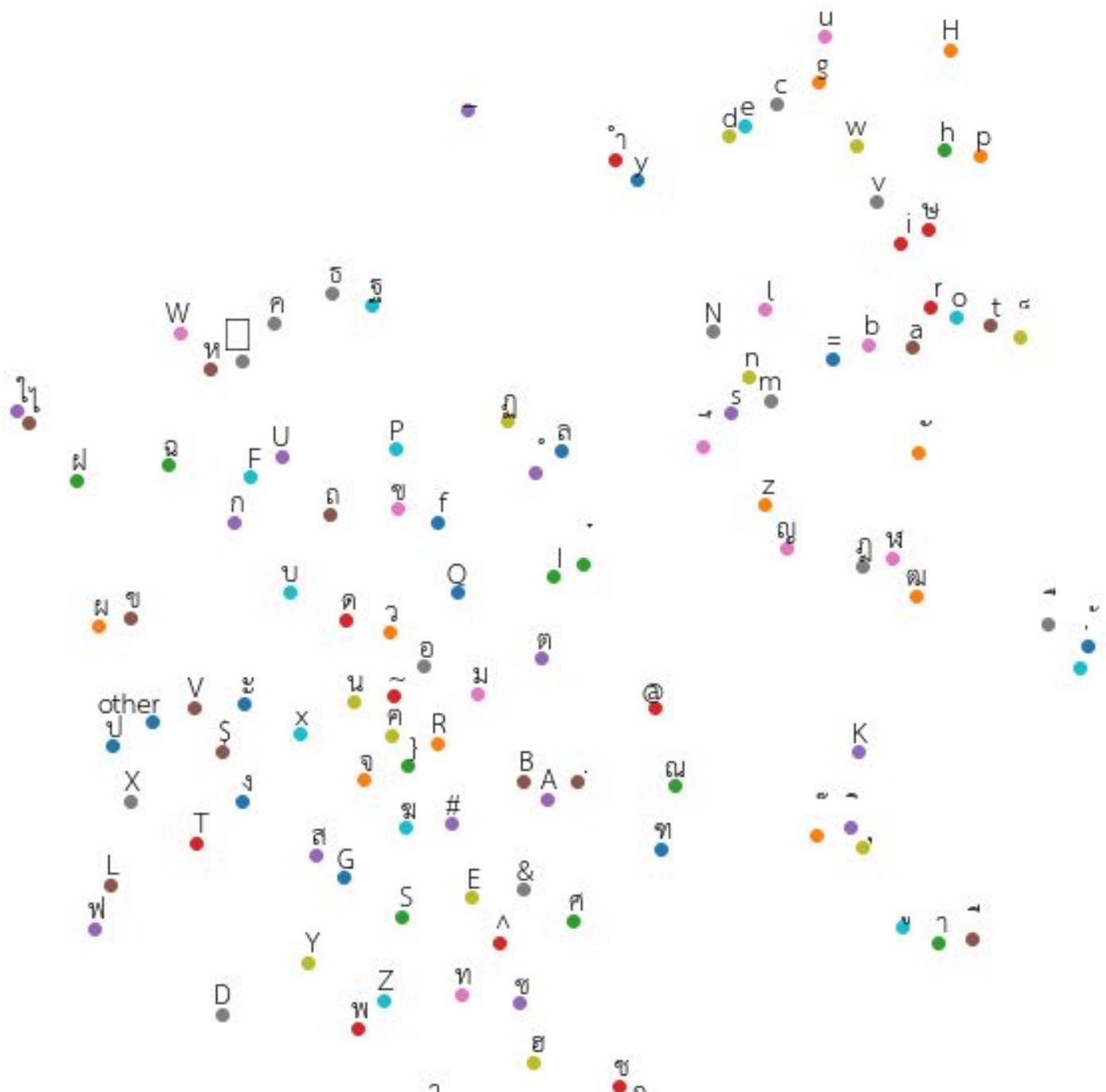
Embedding and meaning (semantics)

- Meaning is inferred from the task *visualizat-*
- Embedding of 32 dimensions -> t-SNE into 2 dimension for visualization
- Automatically!



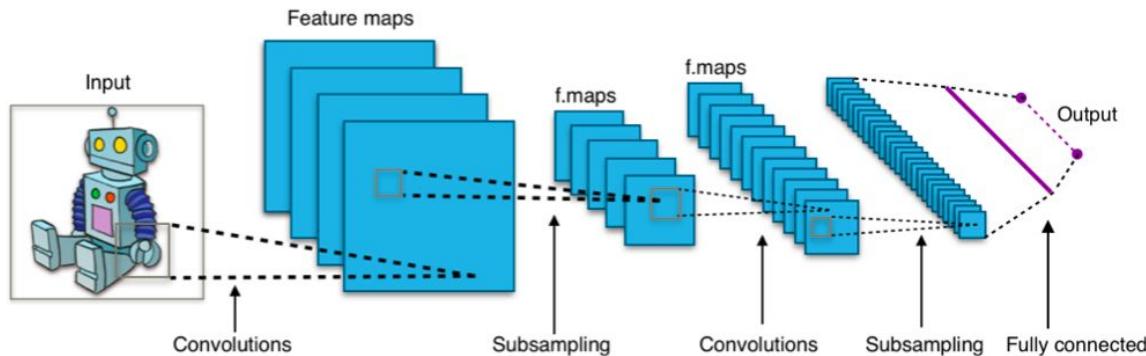
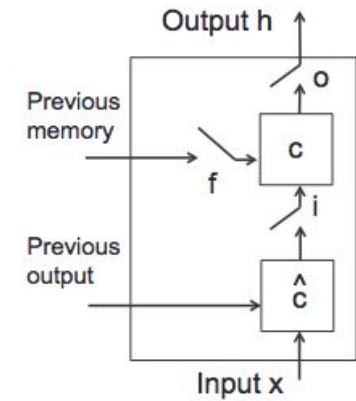
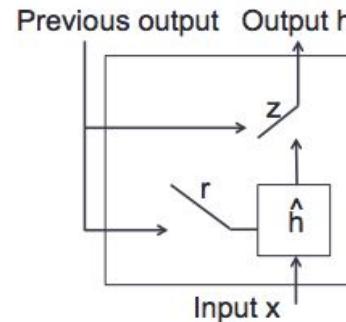






Neural networks

- Fully connected networks
 - SGD, backprop
- CNN
- RNN, LSTM, GRU

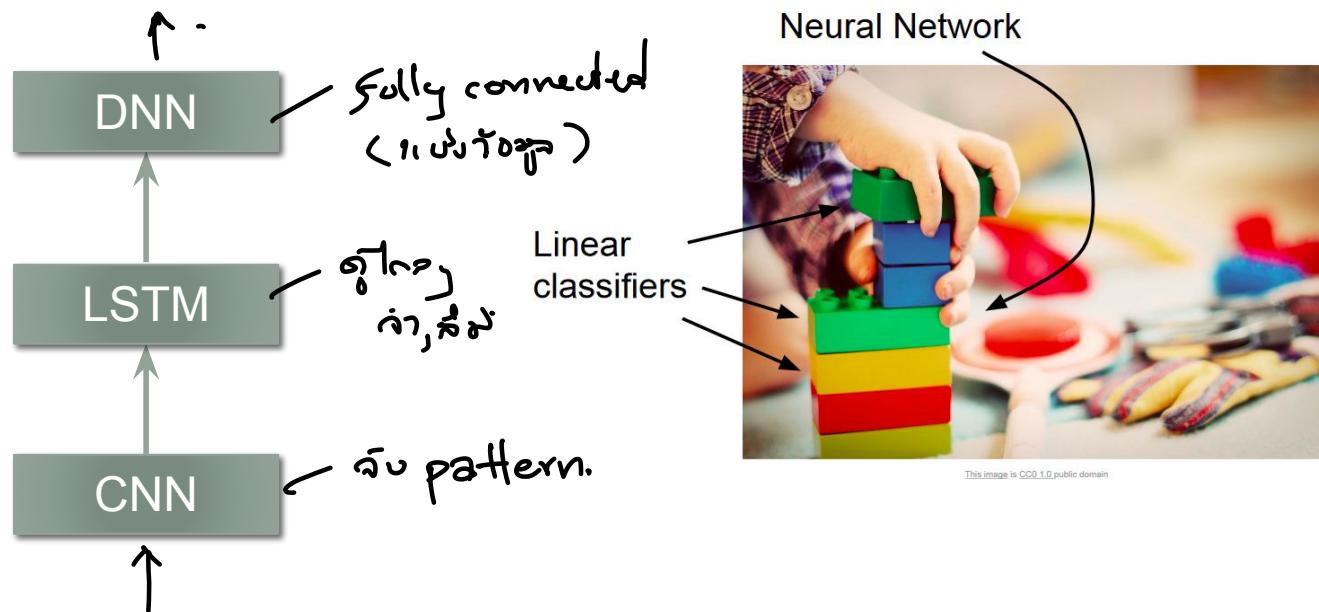


Attention modeling
Recursive neural networks

← Future lectures

DNN Legos

- Typical models now consists of all 3 types
 - CNN: local structure in the feature. Used for feature learning.
 - LSTM: remembering longer term structure or across time
 - DNN: Good for mapping features for classification. Usually used in final layers



Back to tokenization...

TABLE II
RESULTS OF THE SIX BEST TEAMS

Type of participants	F-Measure (%)	Time (mm:ss)
<i>Non-Students^a</i>	97.94937	00:47
<i>Non-Students</i>	97.84097	02:46
<i>Non-Students</i>	97.18822	00:26
<i>Bachelor Students^b</i>	95.78162	01:08
<i>Master Students</i>	95.56670	12:14
<i>PhD+Master Students</i>	92.02067	02:28

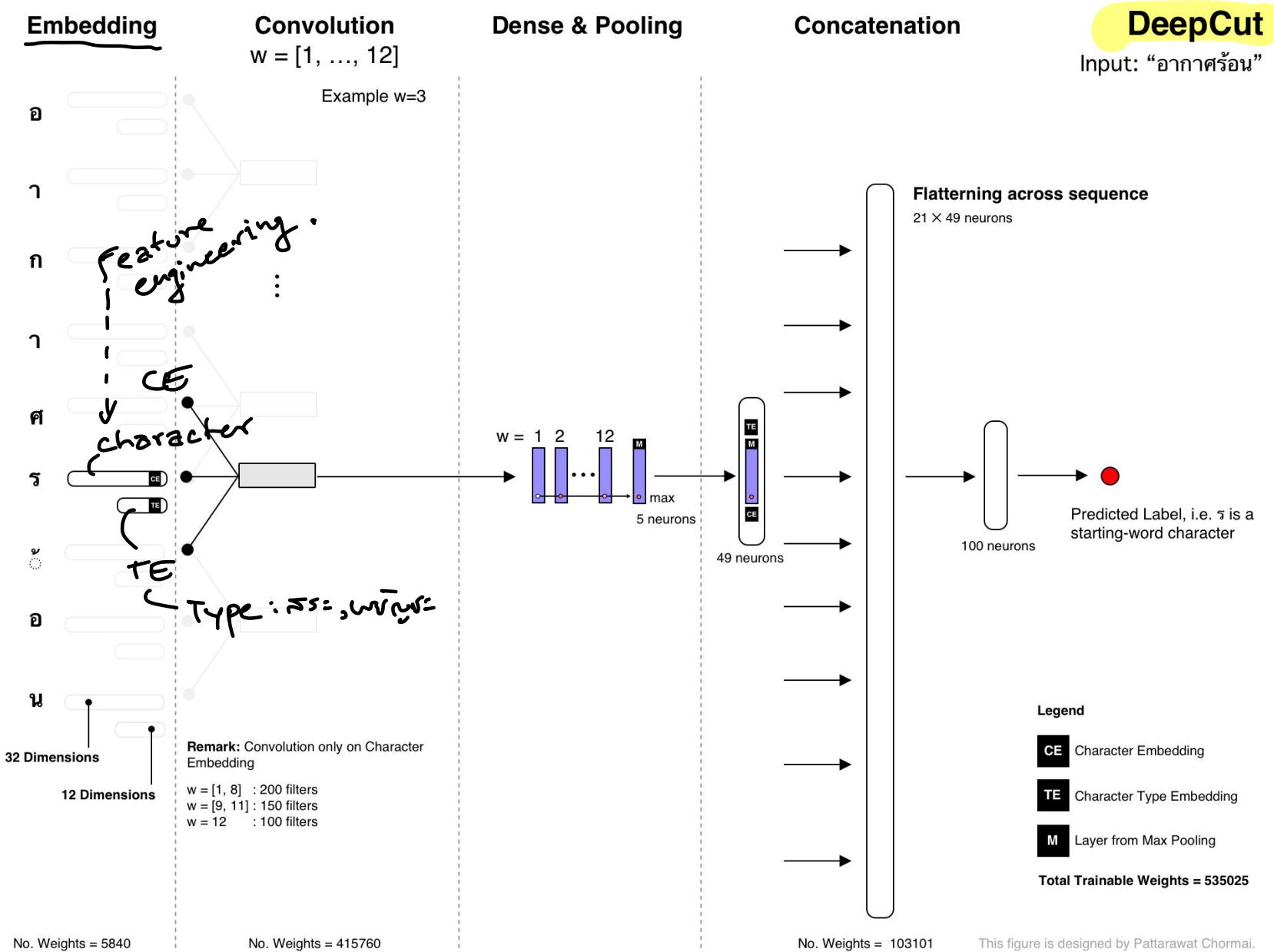
^aBest of the BEST 2009 Award Winner

^bBEST Student 2009 Award Winner

BEST 2009 : Thai word segmentation software contest

<http://ieeexplore.ieee.org/document/5340941/>

https://sertiscorp.com/thai-word-segmentation-with-bi-directional_rnn/



Input sequence



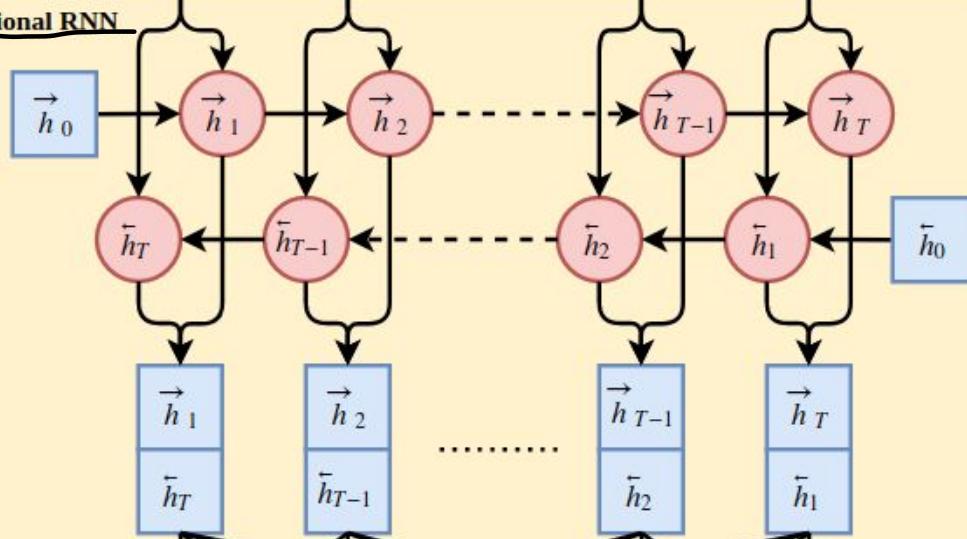
<https://www.sertiscorp.com/november-20-2017>

Embedding lookup



$$e_t = W_c c_t$$

Bi-directional RNN



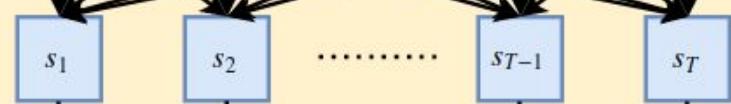
$$z_t = \text{sigmoid}(W_z e_t + U_z h_{t-1} + b_z)$$

$$r_t = \text{sigmoid}(W_r e_t + U_r h_{t-1} + b_r)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \tanh(W_h e_t + U_h (r_t \odot h_{t-1}) + b_h)$$

$$H_t = [\vec{h}_t, \tilde{h}_{T-t+1}]$$

Output score

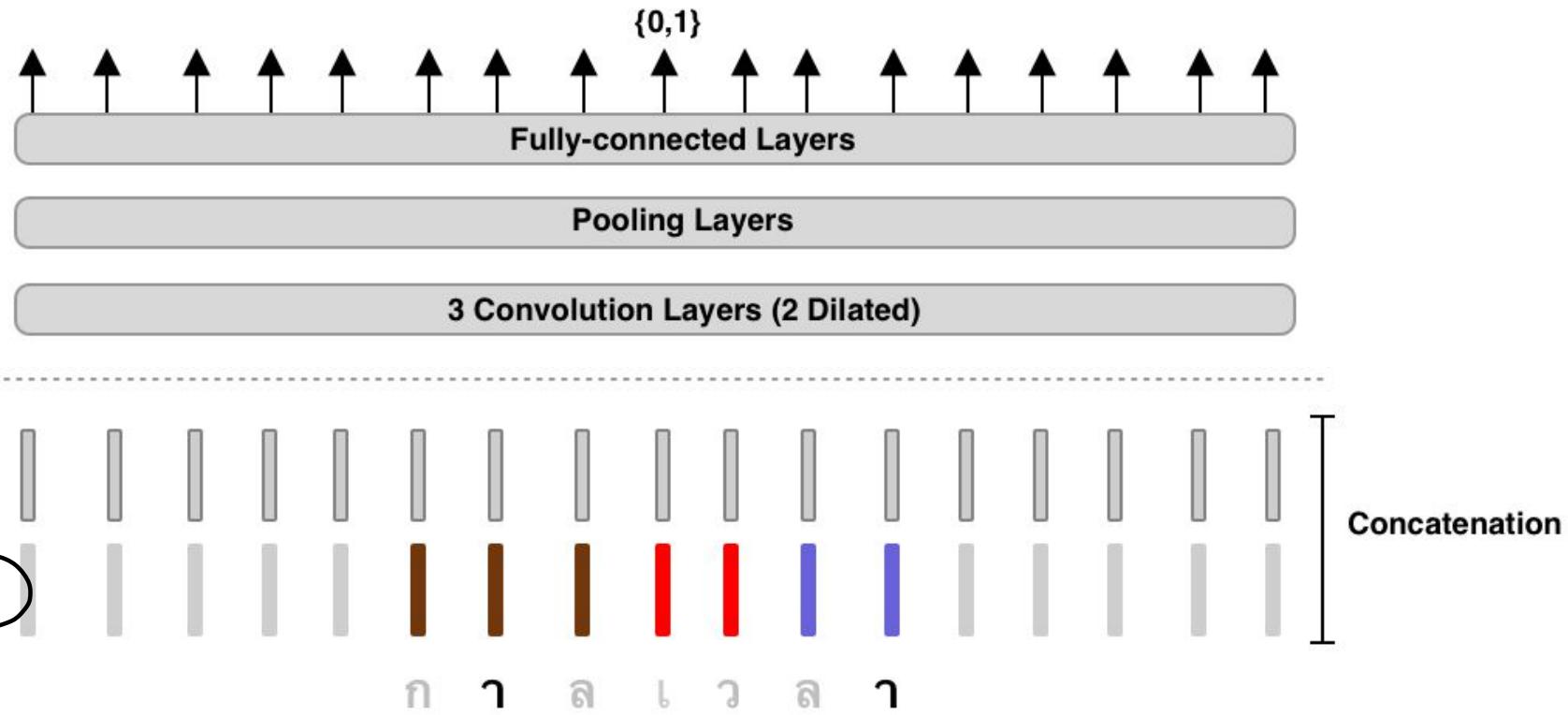


$$s_t = H_t W_s + b_s^\top$$

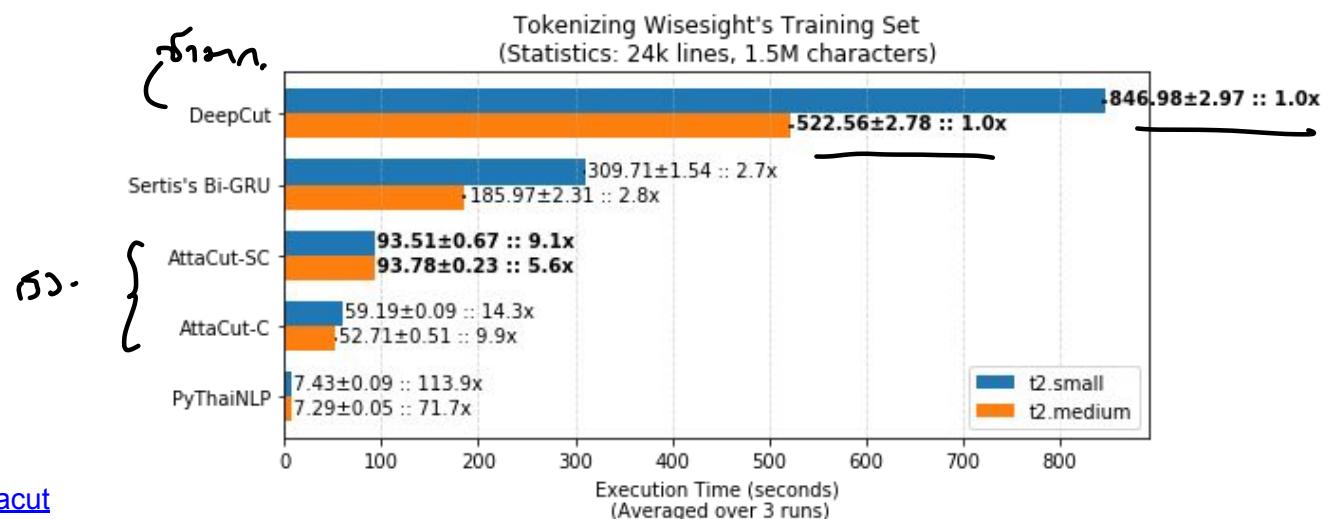
Softmax output



$$p_t = \frac{\exp(s_t)}{\exp(s_{t1}) + \exp(s_{t2})}$$



Last Updated: 29/08/2019		PyThaiNLP newmm	Sertis Bi-GRU	DeepCut	<u>AttaCut-C</u>	<u>AttaCut-SC</u>
BEST Validation Set						
Character-Level	precision	0.94±0.11	0.95±0.10	0.99±0.05	0.97±0.07	0.98±0.05
	recall	0.83±0.09	0.99±0.02	0.99±0.03	0.98±0.04	0.99±0.03
	f1	0.88±0.08	0.97±0.07	0.99±0.04	0.98±0.05	0.99±0.04
Word-Level	precision	0.73±0.16	0.91±0.14	0.97±0.07	0.94±0.10	0.96±0.08
	recall	0.65±0.16	0.94±0.10	0.97±0.07	0.94±0.09	0.97±0.08
	f1	0.68±0.15	0.93±0.12	0.97±0.07	0.94±0.10	0.97±0.08
BEST Test Set						
Character-Level	precision	0.91±0.15	0.92±0.11	0.96±0.08	0.94±0.10	0.95±0.09
	recall	0.85±0.09	0.98±0.04	0.98±0.04	0.98±0.04	0.98±0.04
	f1	0.86±0.11	0.95±0.08	0.97±0.06	0.96±0.07	0.96±0.07
Word-Level	precision	0.70±0.19	0.85±0.18	0.92±0.14	0.88±0.17	0.91±0.15
	recall	0.64±0.18	0.90±0.14	0.93±0.12	0.91±0.14	0.92±0.13
	f1	0.67±0.19	0.87±0.16	0.93±0.13	0.89±0.16	0.91±0.14



ผมเห็นคนวางแผนการนี้เมื่อ 20-30 ปีที่แล้วทำเรื่องตัดคำ วงการ
นี้มันไม่ไปไหนเลยใช่ไหมเนี่ย

มิตรสหาย Business Development ท่านนึง



Words of caution

Statistical tokenizers fail on mismatched data

A tokenizer trained on social text might not be able to cut simple words like

ມະນຸງ ມະລະກອ

ຕົກໆງໆຈູ້ · ຈະ ໄປເກີ

ໜັດໝາຍ

ກໍານົດ ພະ

<https://www.aclweb.org/anthology/2020.emnlp-main.315/>

	WS160	Twitter	TNHC	ຄະນະກົມດົມດົມ
Deepcut	BEST ຈົກໆນົດ 99	93.8	93.5	
Attacut	BEST	93.5	80.8	ອກກົມ

speed
fde
off

witten
robustness.

* Statistical tokenizers fails unpredictably

ໜູກຮອບ => |ໜູກ|ຮອບ|

ຂ້າວຜັດຄະນຳໜູກຮອບໜຶ່ງຈານ => |ຂ້າວຜັດ|ຄະນຳ|ໜູກ|ຮອບ|ໜຶ່ງ|ຈານ|

Might need rule-based override (Deepcut has this)

For speed, maximal matching (newmm) is reliable.

- drawbacks?

Words of caution

Tokenization performance effects downstream task performance *depend on task.*

Can be small (1%) or large (10%)

Specialized tokenizer can help your downstream task

Example: ecommerce search | ห្ម|ຟັງ| |ຕ່າງ| ห្ម|



 Ekapol Chuangsuwanich
 October 25, 2020 · 
...

Month 1 in ຕັດຄຳ: the best solution is deep learning
 Month 6 in ຕັດຄຳ: the best solution is no segmentation
 Month 24 in ຕັດຄຳ: the best solution is sentencepiece
 Month 36 in ຕັດຄຳ: the best solution is dict-based
 Month infinity in ຕັດຄຳ: the best solution is a new language 😊
 The thing I find the most fascinating is that for *all* the changes in Thai NLP in the last 3 years, this problem still causes other shit to break.
 #ThaiNLP #WordSegmentation

Tokenization - English

- Even English has tokenization issues!
 - Space is usually not enough
 - aren't
 - are + n't
 - aren't
 - arent
 - aren t
 - are + not
 - San Francisco
- Usually includes the text normalization step
- Again depends on application
 - “aren't” might be different from “are not” for sentiment analysis

Other English issues - hyphens

ñuussnō

- “the *New York-based* co-operative was *fine-tuning* forty-two *K-9-like* models. “
- Lexical vs Sentential hyphens

Tokenization of non-standard text

- Twitter

@SentimentSymp: can't wait for the Nov 9 #Sentiment talks! YAAAAAAAY!!! >:-D <http://sentimentsymposium.com/>.

Needs to correctly tokenize

Emoticons

Twitter markup (# and @)

Capitalization (and html tags for bold, etc.)

Lengthening

Dates

Hand-crafted for tweets

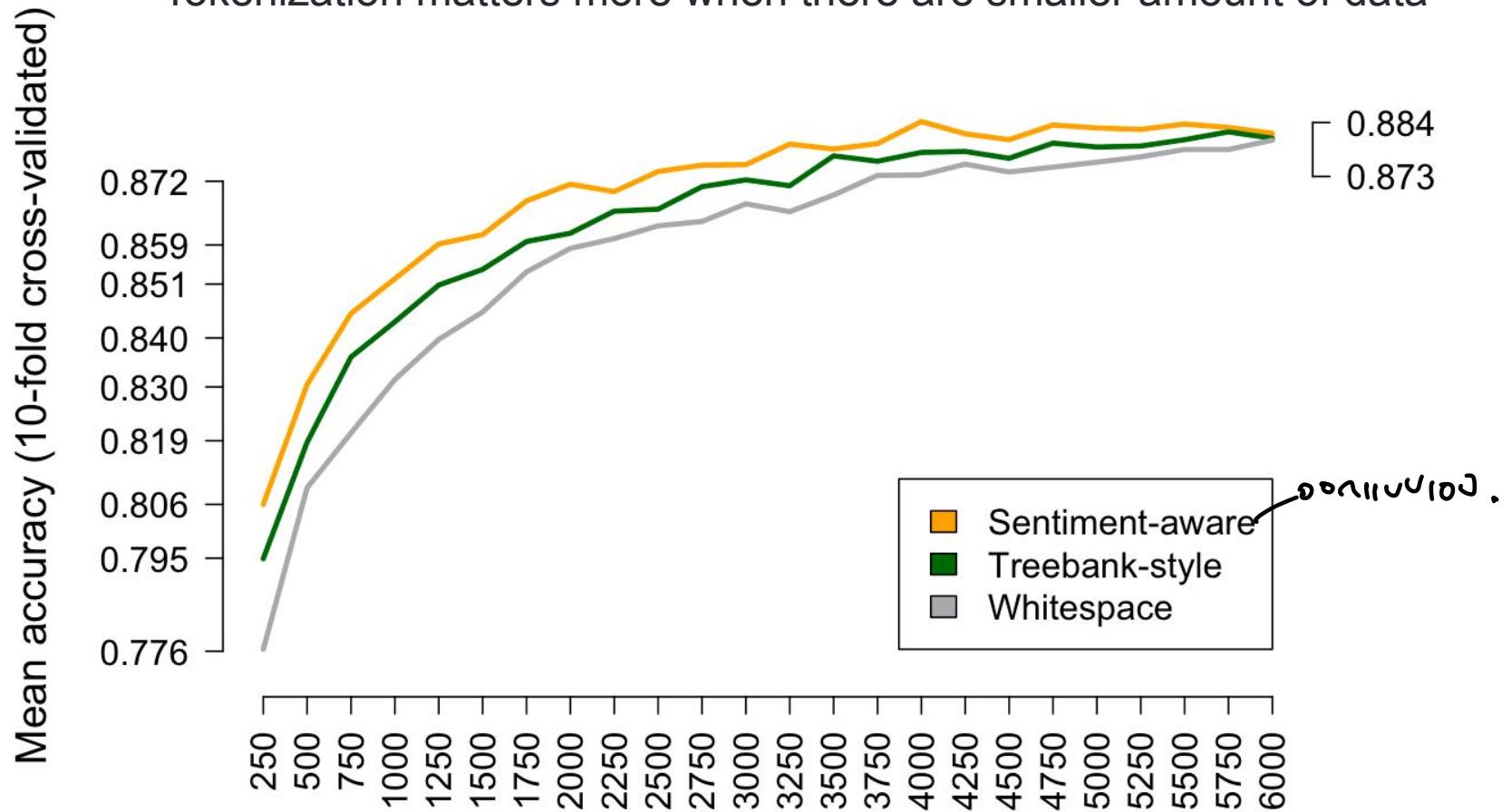
```
@sentimentsymp
:
can't
wait
for
the
Nov_09
#sentiment
talks
!
YAAAAAY
!
!
!
>:-D
http://sentimentsymposium.com/
.
```

@
SentimentSymp
:
ca
n't Standard tokenizer
wait (Stanford tokenizer)
for
the
Nov
9

Sentiment
talks
!
YAAAAAY
!
!
!
>
;
:
-D
http
:
//sentimentsymposium.com/
.

Sentiment analysis (in-domain)

Tokenization matters more when there are smaller amount of data



Sentiment analysis (out-of-domain)

Better tokenizer gives better portability



End-to-end models

End-to-end learning



- Classical machine learning systems usually break the problem into smaller subtasks
 - Self-driving:
 - Image -> objects detection -> path finding -> steering
 - Speech2speech translation:
 - Speech A -> text A -> text B -> Speech B
- End-to-end models use one large neural networks process the input and generate the desired output
 - Image -> steering
 - Speech A -> Speech B

End-to-end NLP?

Discourse

Semantics

CommunicationEvent(e)
Agent(e, Alice)
Recipient(e, Bob)
SpeakerContext(s)
TemporalBefore(e, s)

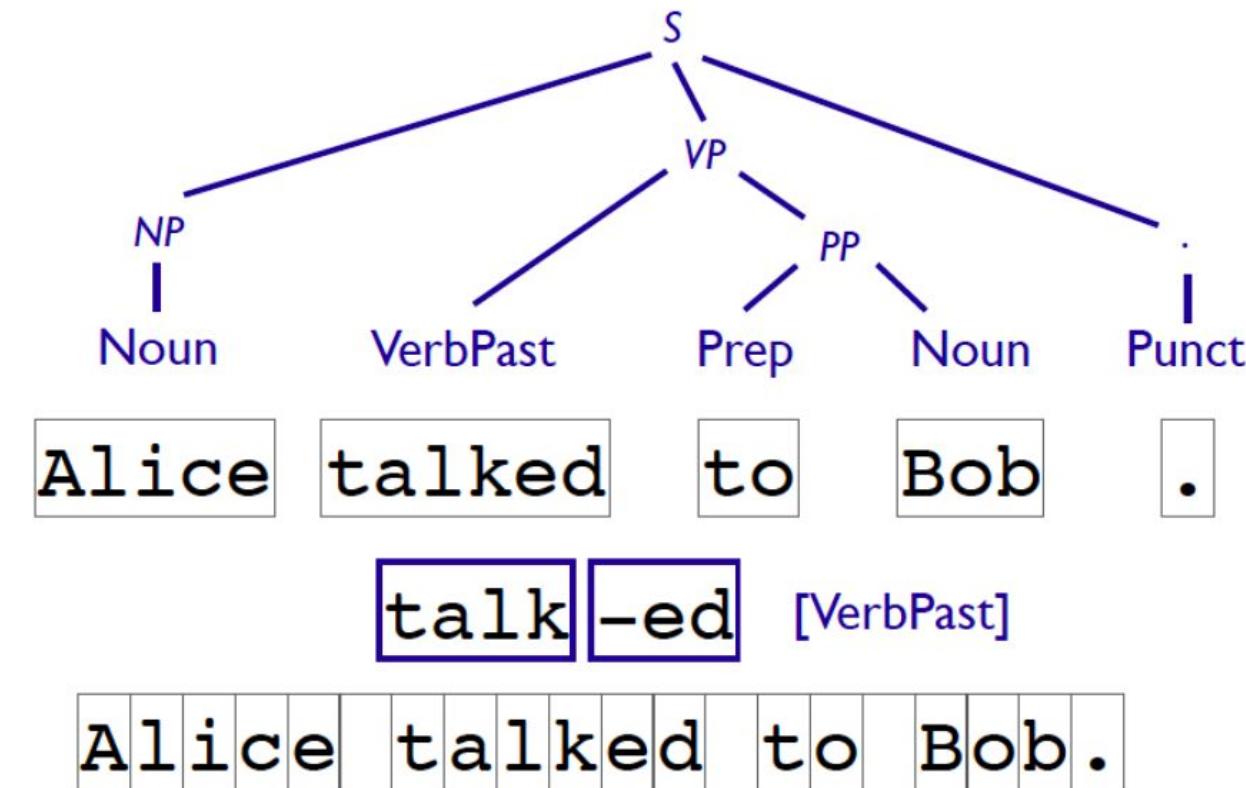
Syntax: Constituents

Syntax: Part of Speech

Words

Morphology

Characters



Towards no tokenization? :Character-aware neural language models

- Input: previous characters
- Output: next word

CNN over characters capture character sequence patterns

Fully connected

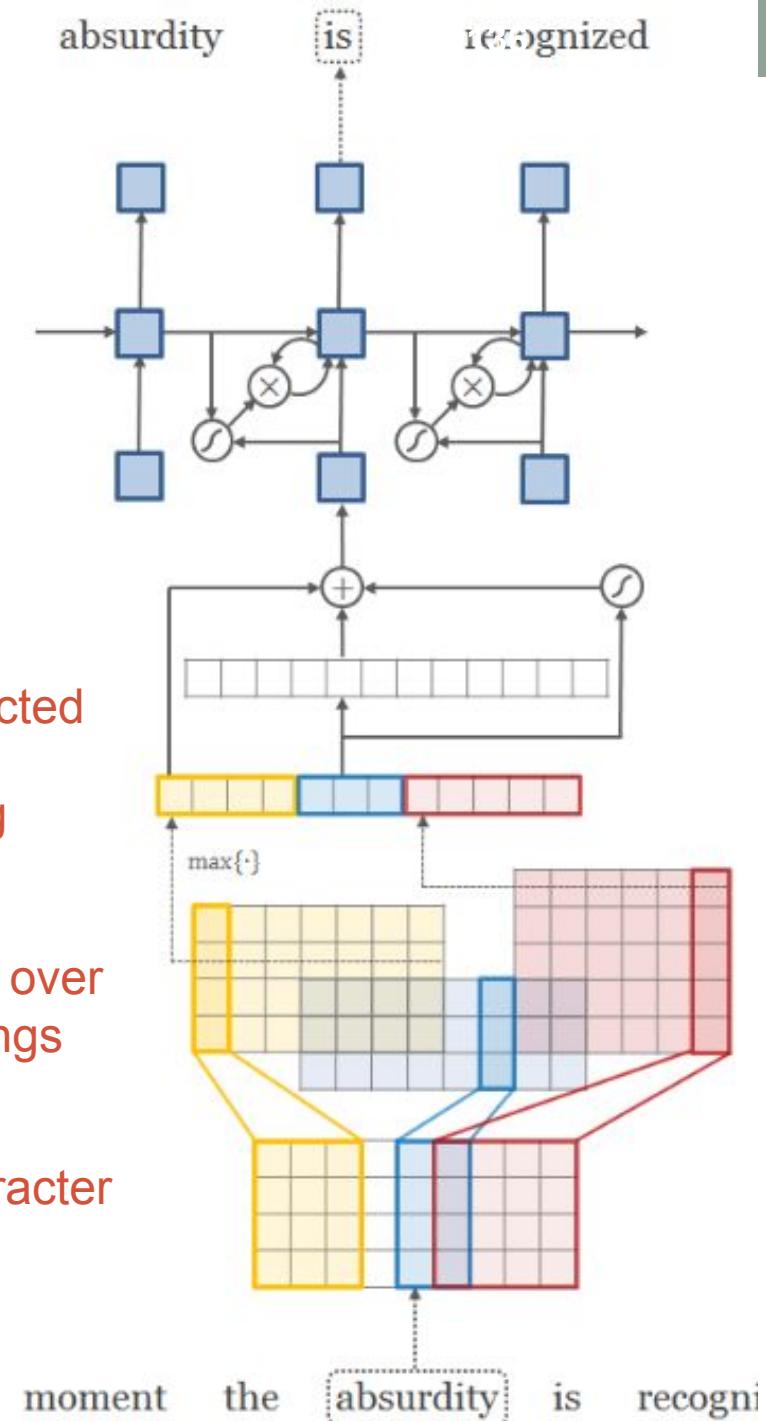
LSTM

Fully connected

Max pooling

Convolutional layer over character embeddings

Character to character embeddings



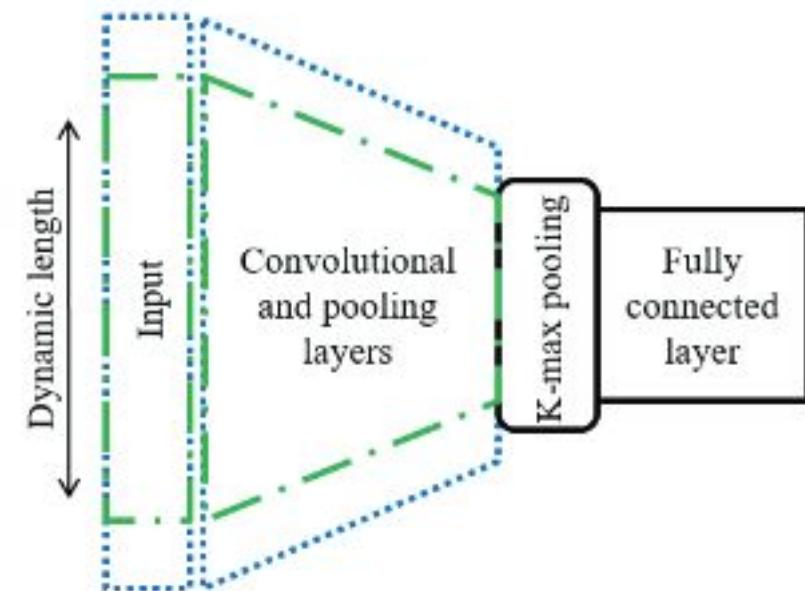
Towards no tokenization

- Text classification using charCNN on Thai

Method	Accuracy (%)	F_1 (%)
Naïve Bayes, BoW	87.2	87.1
Naïve Bayes, TF-IDF	89.0	88.9
Logistic Regression, BoW	94.8	94.8
Logistic Regression, TF-IDF	94.7	94.7
SVM, BoW	93.7	93.7
SVM, TF-IDF	95.2	95.2
DCNN (Kalchbrenner et al., 2014)	95.9	95.9
Proposed Char-CNN	95.4	95.4

Characters

Word-based
methods



A character-level convolutional neural network with dynamic input length for Thai text categorization
<http://ieeexplore.ieee.org/document/7886102/>

Caveats of end-to-end models

- Requires lots of data for the specific task
- Hard to fix specific mistakes by the model *

text ~~survive~~ → text ~~survive~~ → text classify →

text ~~survive~~ → text classify →

↑
Add data → not guarantee diagnosis..: end2end ~~with~~ ~~without~~

Things to consider when thinking about tokenization

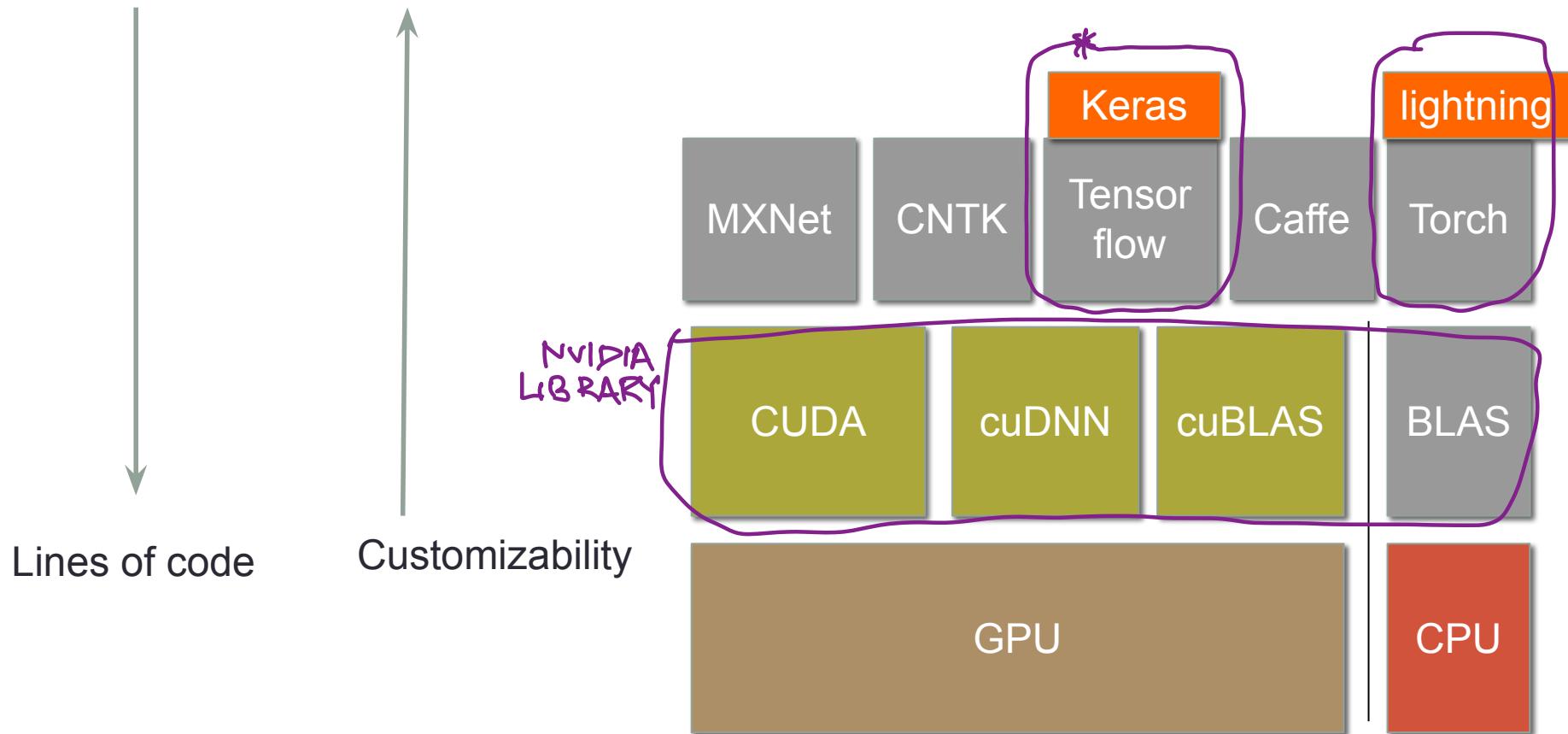
Know your use cases

Syllable (word)

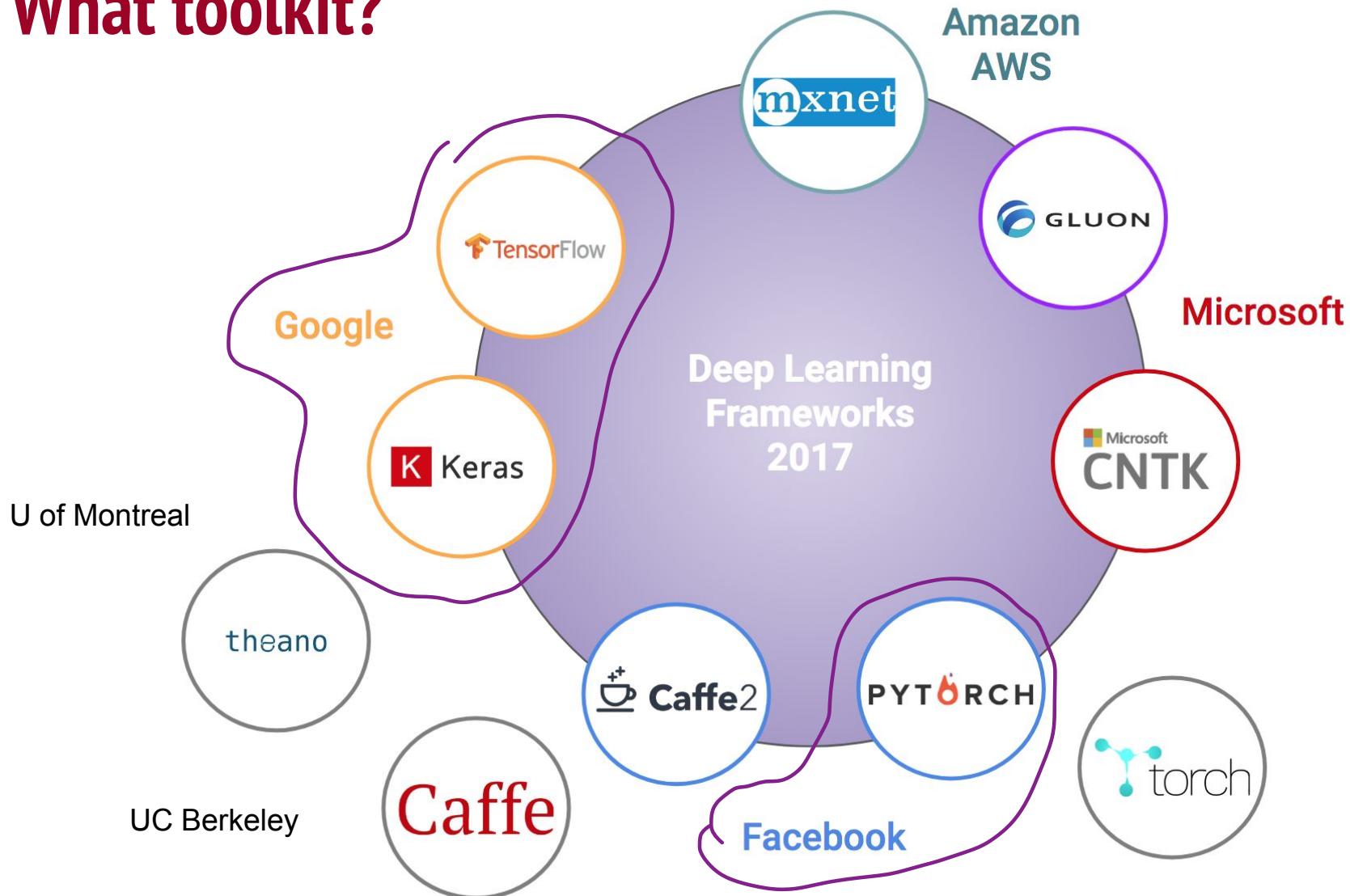
<u>Segment text into characters</u>	<u>Segment text into word-pieces</u>	<u>Segment text into words</u>
<p><u>Embedding table size.</u></p> <ul style="list-style-type: none"> → <u>Small embedding table</u> → <u>Needs a powerful model to learn very long range influences</u> → Individual tokens are not very meaningful <p><u>To character</u> ຈະກົດຕົວ.</p>	<ul style="list-style-type: none"> → <u>Medium-sized embedding table</u> → Needs a model that can learn long range dependencies → Individual tokens are partially meaningful <p><u>From character to word-piece</u></p>	<ul style="list-style-type: none"> → <u>Very large embedding table</u> → Smaller number of segments means the model need not learn long range influences → Individual tokens are meaningful <p><u>From character to word</u>.</p>

What toolkit

Tradeoff between customizability and ease of use

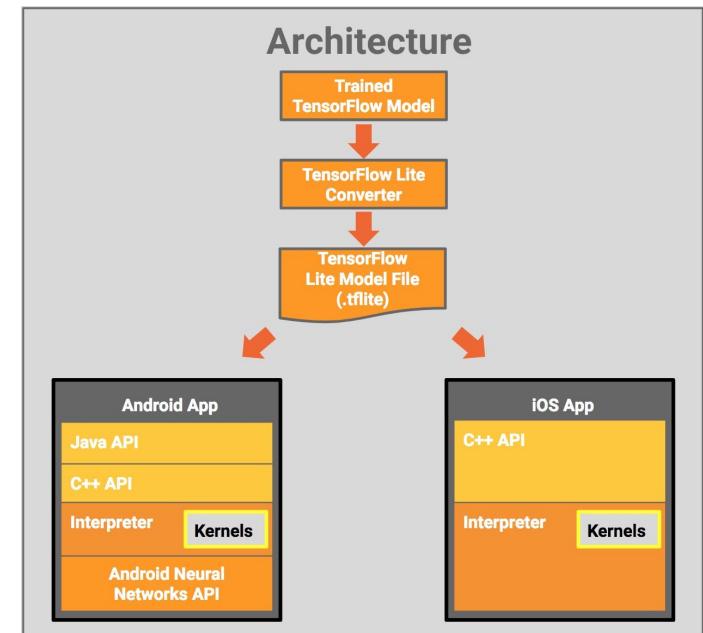


What toolkit?



Which?

- Easiest to use and play with deep learning: Keras
- Easiest to use and tweak: pytorch
- Easiest to deploy: tensorflow
 - Tensorflow lite for mobile
 - TensorRT support
 - tf-serving
 - Kubeflow



Which?

- Easiest to use and play with deep learning: Keras
- Easiest to use and tweak: pytorch
- Easiest to deploy:
 - Tensorflow lite for mobile
 - TensorRT support
 - tf-serving
 - Kubeflow
- Community: TensorFlow?

Demos

- Keras
- You should read before doing the homework
 - https://keras.io/guides/functional_api/

Keras steps

- 1 • Define the network -
 - 2 • Compile the network -
 - 3 • Fit the network ✓

```
def get_feedforward_nn():
    input1 = Input(shape=(21,))
    x = Dense(100, activation='relu')(input1)
    x = Dense(100, activation='relu')(x)
    x = Dense(100, activation='relu')(x)
    out = Dense(1, activation='sigmoid')(x)

    model = Model(inputs=input1, outputs=out)
    model.compile(optimizer=Adam(),
                  loss='binary_crossentropy',
                  metrics=['acc'])

    return model
```

Keras is easy!

Dense layer

Dense class

```
tf.keras.layers.Dense(  
    units,  
    activation=None, must put  
    use_bias=True,  
    kernel_initializer="glorot_uniform",  
    bias_initializer="zeros",  
    kernel_regularizer=None,  
    bias_regularizer=None,  
    activity_regularizer=None,  
    kernel_constraint=None,  
    bias_constraint=None,  
    **kwargs  
)
```

Just your regular densely-connected NN layer.

`Dense` implements the operation: `output = activation(dot(input, kernel) + bias)` where `activation` is the element-wise activation function passed as the `activation` argument, `kernel` is a weights matrix created by the layer, and `bias` is a bias vector created by the layer (only applicable if `use_bias` is `True`).

Dropout layer

Dropout class

```
tf.keras.layers.Dropout(rate, noise_shape=None, seed=None, **kwargs)
```

Applies Dropout to the input.

The Dropout layer randomly sets input units to 0 with a frequency of `rate` at each step during training time, which helps prevent overfitting. Inputs not set to 0 are scaled up by $1/(1 - \text{rate})$ such that the sum over all inputs is unchanged.

Note that the Dropout layer only applies when `training` is set to True such that no values are dropped during inference. When using `model.fit`, `training` will be appropriately set to True automatically, and in other contexts, you can set the kwarg explicitly to True when calling the layer.

(This is in contrast to setting `trainable=False` for a Dropout layer. `trainable` does not affect the layer's behavior, as Dropout does not have any variables/weights that can be frozen during training.)

Conv2D layer

Conv2D class

```
tf.keras.layers.Conv2D(  
    filters,  
    kernel_size,  
    strides=(1, 1),  
    padding="valid",  
    data_format=None,  
    dilation_rate=(1, 1),  
    groups=1,  
    activation=None,  
    use_bias=True,  
    kernel_initializer="glorot_uniform",  
    bias_initializer="zeros",  
    kernel_regularizer=None,  
    bias_regularizer=None,  
    activity_regularizer=None,  
    kernel_constraint=None,  
    bias_constraint=None,  
    **kwargs  
)
```

2D convolution layer (e.g. spatial convolution over images).

Deep learning demo

Red boundaries Indicate output and input tensors

