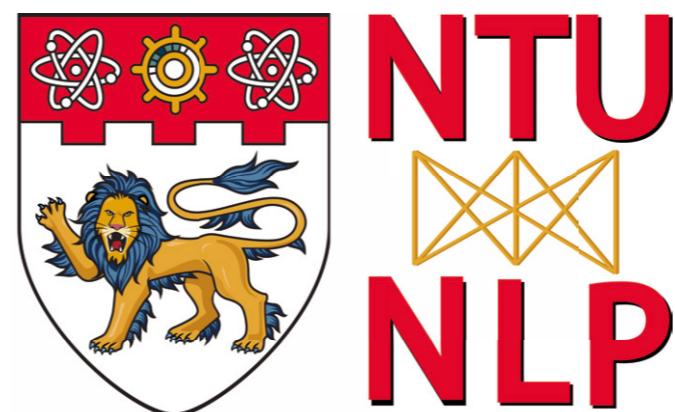


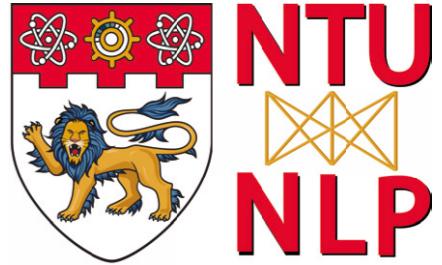
Deep Learning for Natural Language Processing

Shafiq Joty



Lecture 5: More on Word Vectors

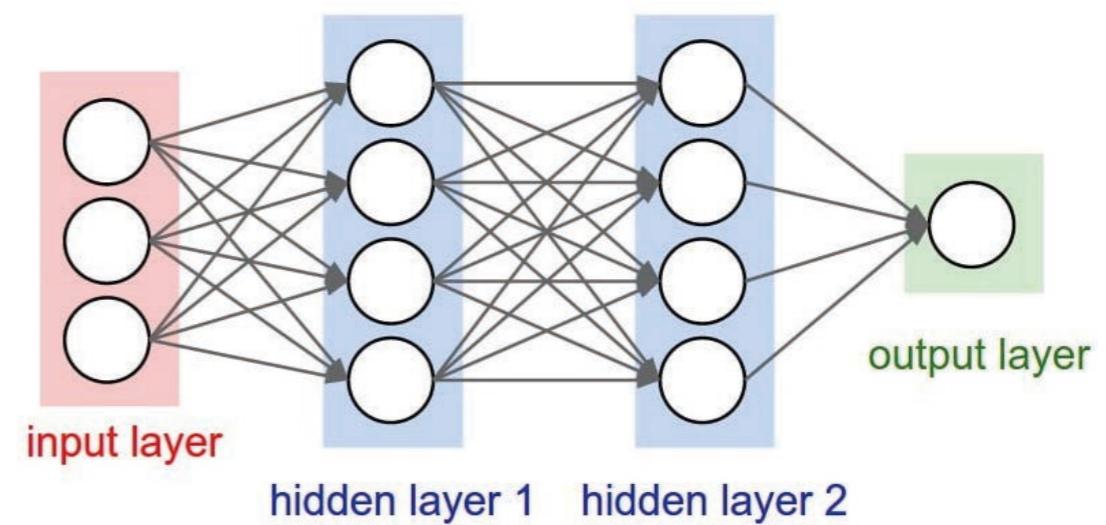
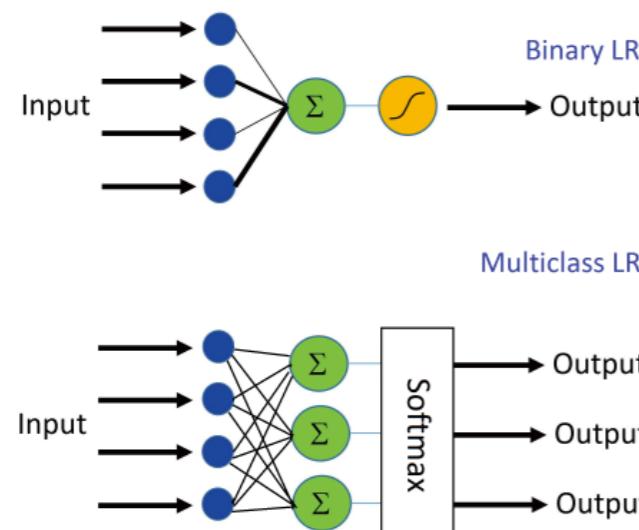
Announcements



- Assignment 1 submission deadline today
- Extension for ICML and KDD deadlines
- Start thinking about your course project

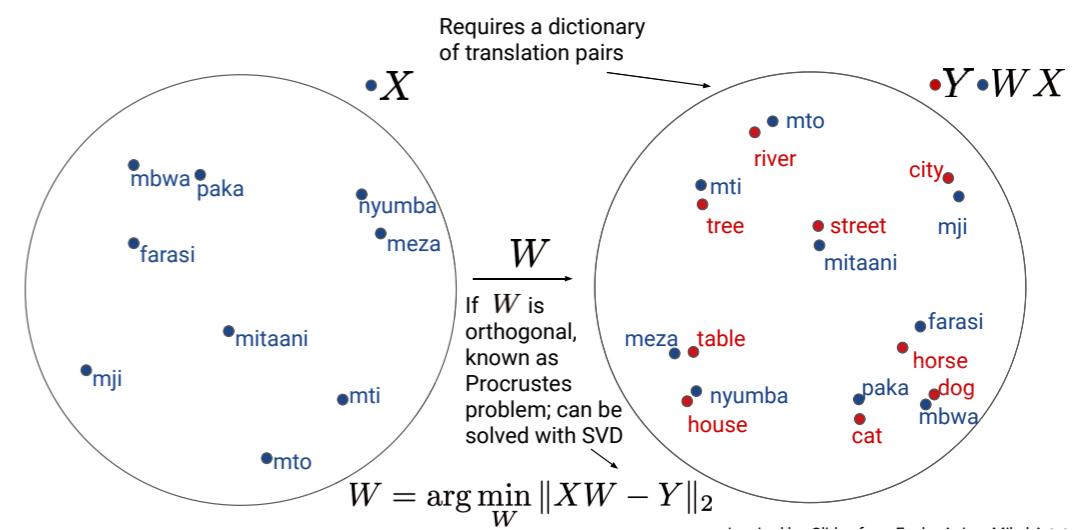
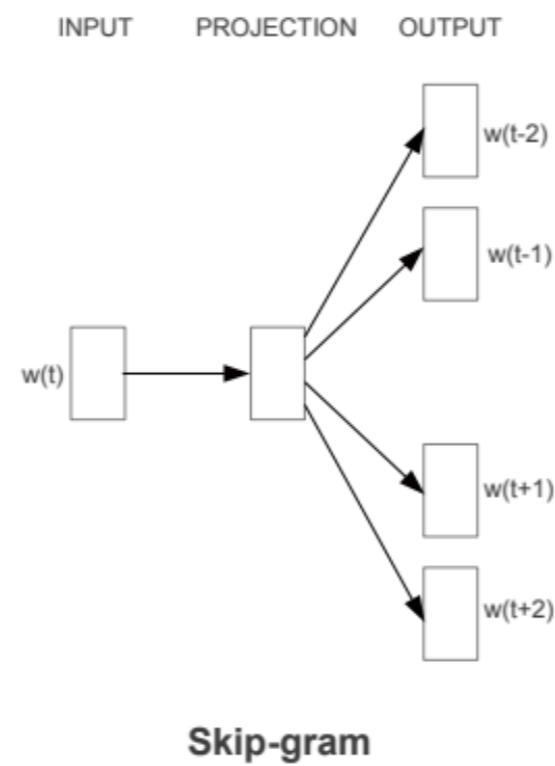
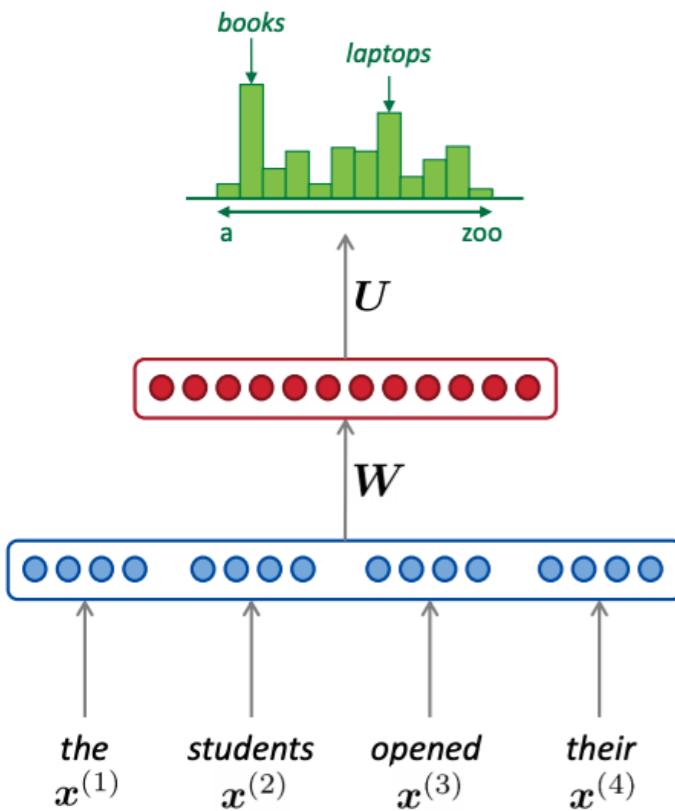
Things Covered So Far

- Lecture 2: Linear Models (Linear & Logistic Regression)
 - Training with gradient descent and SGD
- Lecture 3: Extended linear models to MLP/FNN
 - Training with SGD (+ Backdrop)

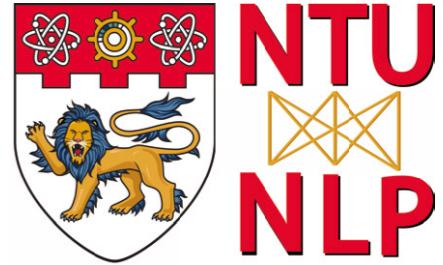


Things Covered So Far

- Lecture 4 & 5: Used simple FNN as
 - Simple Language Models
 - Word2Vec Models



Where we are



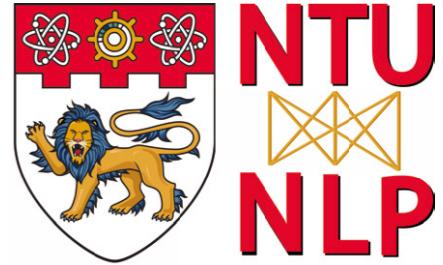
Models/Algorithms

- Linear models
- Feed-forward Neural Nets (FNN)

NLP tasks/applications

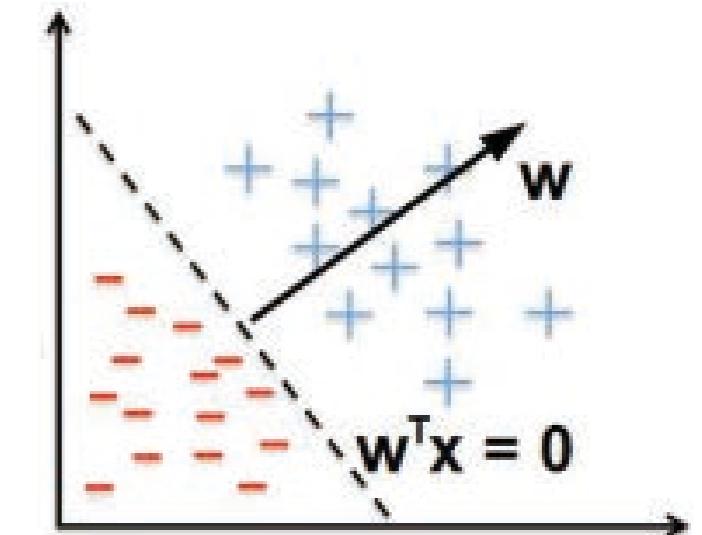
- Word meaning
- Language modelling

Linear Classifiers



- For linear models (logistic regression), we have a linear decision boundary (hyperplane).
- **Traditional approach:** assume x_i are fixed.

From lecture 2:



Instead of **maximizing** the log-likelihood, we can equivalently **minimize** the **negative log likelihood** or NLL.

The negative log-likelihood for logistic regression is given by

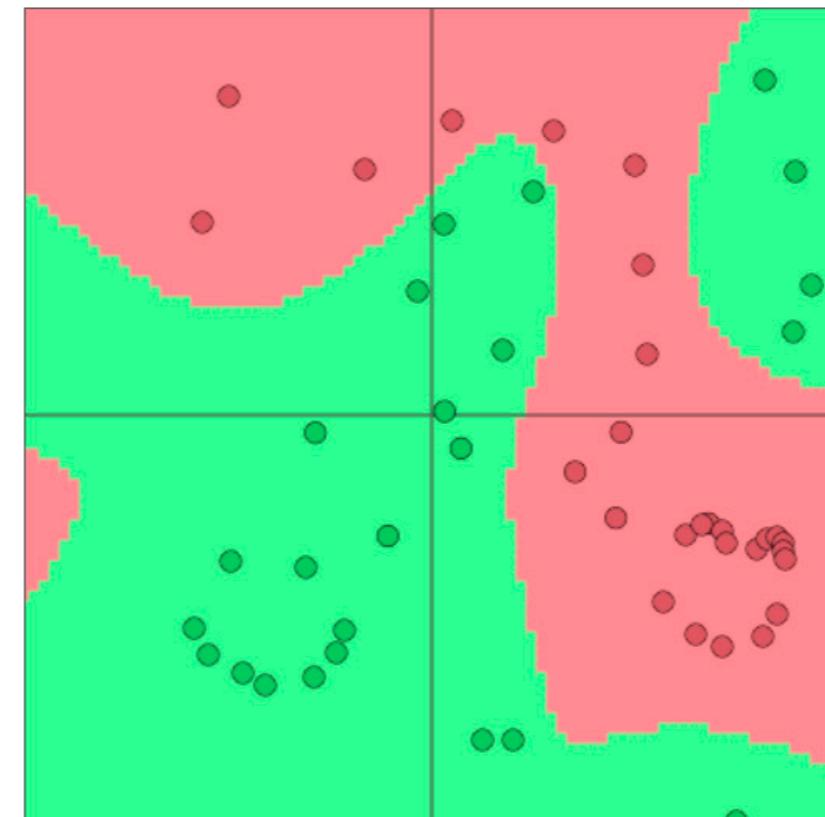
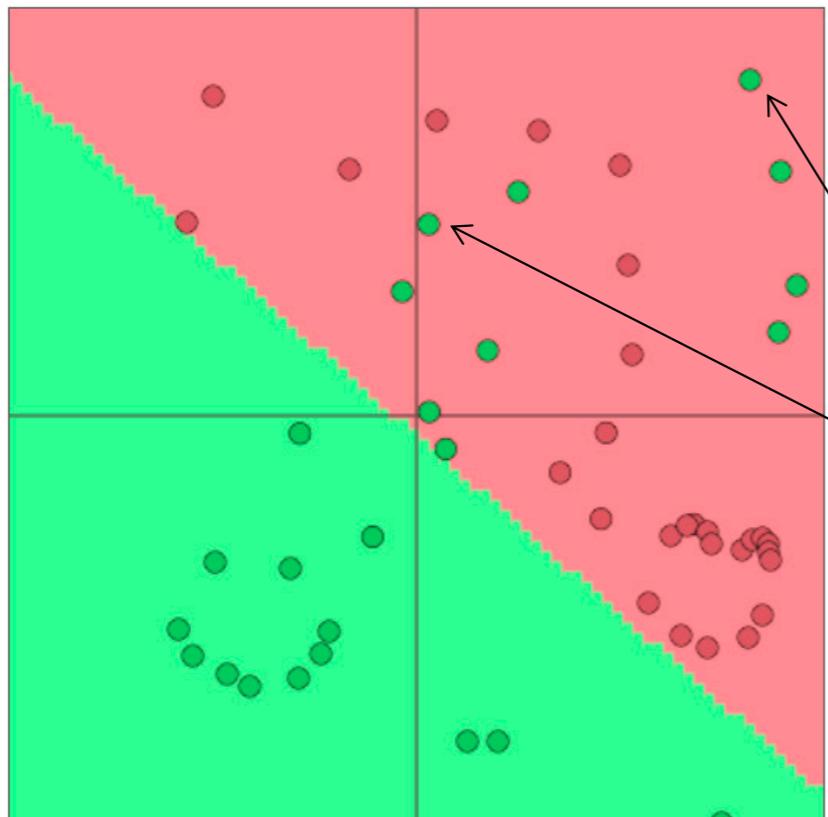
$$\begin{aligned} NLL(\mathbf{w}) &= - \sum_{i=1}^N \log[\mu_i^{1(y_i=1)} \times (1 - \mu_i^{1(y_i=0)})] \\ &= - \sum_{i=1}^N [y_i \log \mu_i + (1 - y_i) \log(1 - \mu_i)] \end{aligned}$$

This is also called the **Cross Entropy Error Function**.

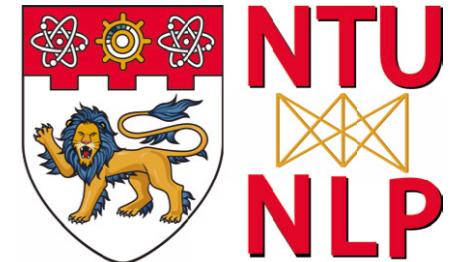
$$H(p, q) = - \sum_{c=1}^C p(c) \log q(c)$$

Neural Classifiers

- Linear models or Softmax (\approx logistic regression) alone not very powerful
- Real world classification problems are often too complex involving non-linear complex decision boundary
- Neural networks can learn more complex transformation functions and non-linear decision boundaries.



Neural vs. Linear Classifiers

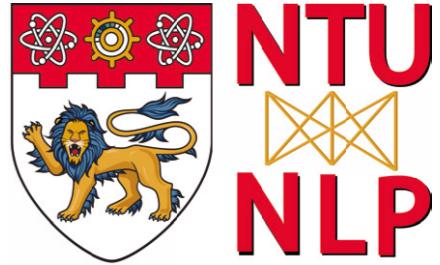


- We learn **both** W (weights, parameters) **and** representations (x & $f(x)$)

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \nabla_{W_{\cdot 1}} \\ \vdots \\ \nabla_{W_{\cdot d}} \\ \nabla_{x_{aardvark}} \\ \vdots \\ \nabla_{x_{zebra}} \end{bmatrix} \in \mathbb{R}^{Cd + Vd}$$

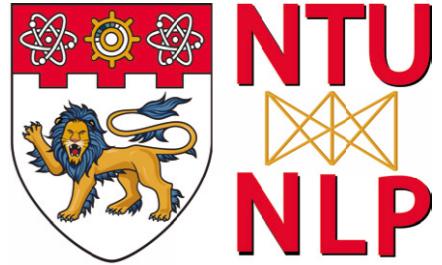
Very large number of parameters!

Lecture Plan



- Classification tasks in NLP
- Window-based Approach for Language Modeling
- Window-based Approach for NER, POS tagging, and Chunking
- Convolutional Neural Net for NLP
- Max-margin Training
- Scaling Softmax (for Language Modeling)
- Adaptive input and output.

Classification Problems



In a supervised setup, we have a training dataset

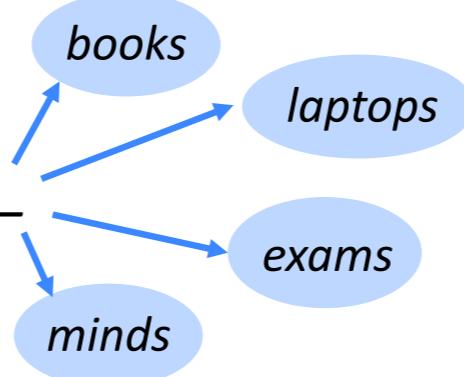
$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

- \mathbf{x}_i are **inputs**, e.g. words, sentences, documents, etc.
- y_i are **labels** (one of C classes) we try to predict, e.g.,
 - Other words: Language model
 - Classes: POS tag, sentiment, named entities
 - Multi-word sequences: Translation, Summarisation

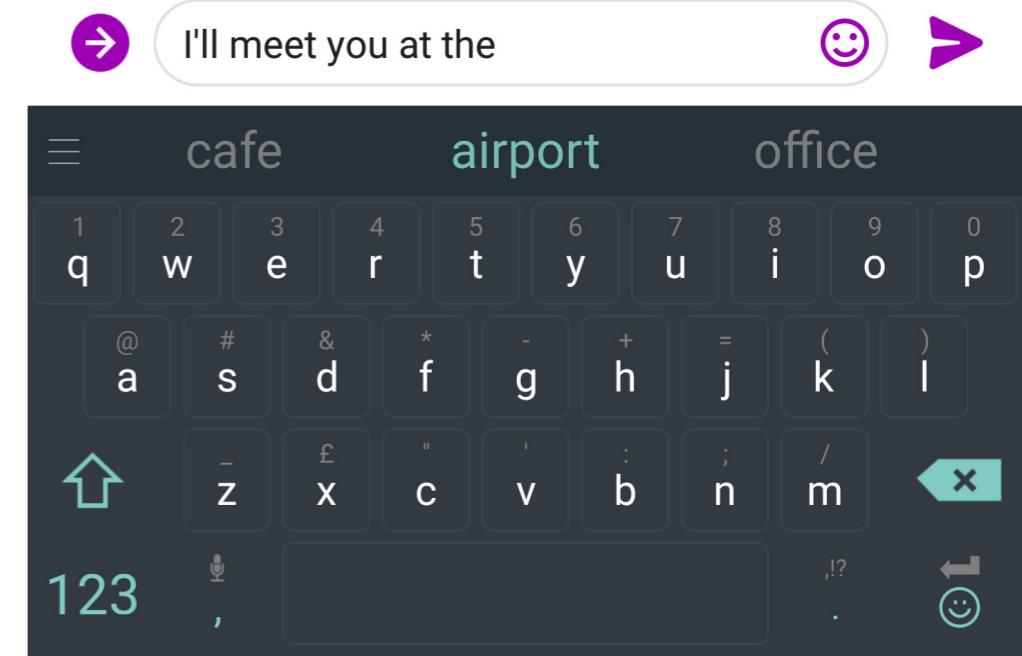
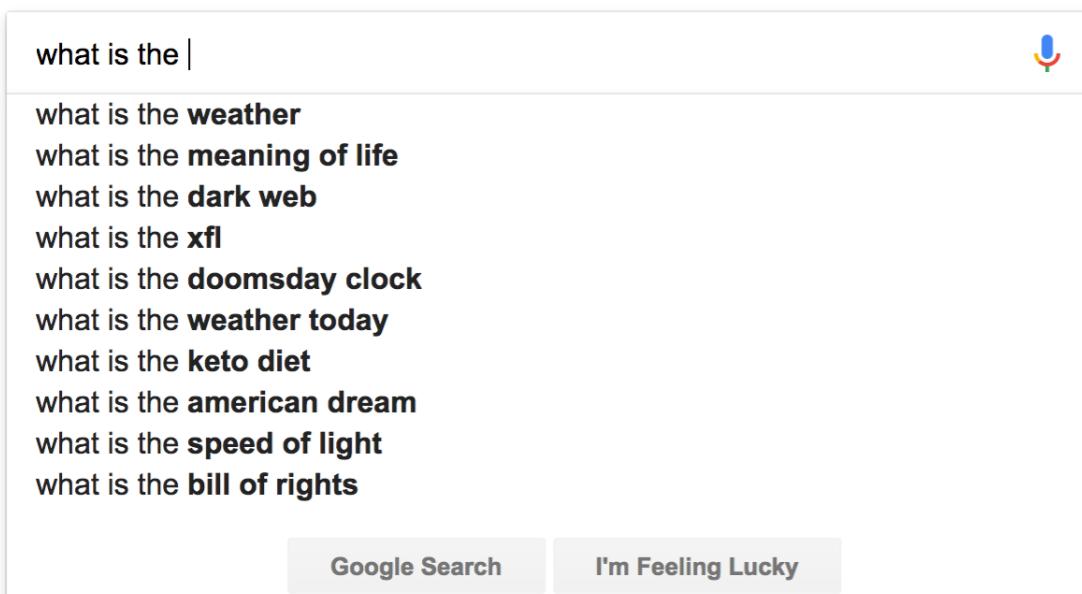
Language Modelling

A language model takes a list of words (history/context), and attempts to predict the word that follows them

the students opened their _____



Google



Language Modelling

(i) A language model takes a list of words (history/context), and attempts to predict the word that follows them

More formally: given a sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(t)}$, compute the probability distribution of the next word $x^{(t+1)}$:

$$P(x^{(t+1)} | x^{(t)}, \dots, x^{(1)})$$

where $x^{(t+1)}$ can be any word in the vocabulary $V = \{w_1, \dots, w_{|V|}\}$

Language Modelling

(ii) Language Models (LM) also assign a probability to a piece of text.

For example, if we have some text $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}$, then the probability of this text (according to the Language Model) is:

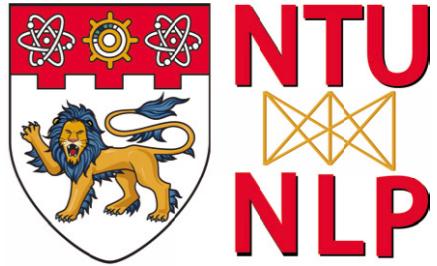
$$P(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(T)}) = P(\mathbf{x}^{(1)}) \times P(\mathbf{x}^{(2)} | \mathbf{x}^{(1)}) \times \dots \times P(\mathbf{x}^{(T)} | \mathbf{x}^{(T-1)}, \dots, \mathbf{x}^{(1)})$$

$$= \prod_{t=1}^T P(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(1)})$$



This is what our LM provides

Language Modelling - Ngram based Approach



Today: we will revisit window-based approach for language modelling

But, before that we will see how language modelling was done in the pre-neural era

Ngram based LM

A **n-gram** is a sequence of **n** consecutive words.

- **unigrams**: “the”, “students”, “opened”, “their”
- **bigrams**: “the students”, “students opened”, “opened their”
- **trigrams**: “the students opened”, “students opened their”
- **4-grams**: “the students opened their”

Idea: Collect statistics about how frequent different n-grams are, and use these to predict next word.

Ngram based LM

Markov assumption: make a simplifying assumption that a word's probability depends only on the preceding n-1 words

$$P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = P(\mathbf{x}^{(t+1)} | \overbrace{\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)}}^{n-1 \text{ words}})$$

How can we estimate this probability?

$$= \frac{P(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}{P(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}$$

$$\approx \frac{\text{count}(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}{\text{count}(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}$$

Ngram based LM

How can we estimate this probability?

$$= \frac{P(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}{P(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}$$

$$\approx \frac{\text{count}(\mathbf{x}^{(t+1)}, \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}{\text{count}(\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(t-n+2)})}$$

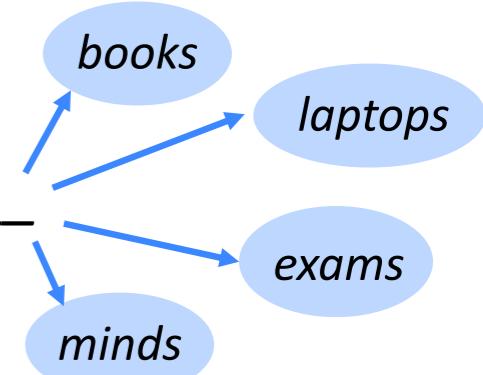
Count this in the corpus

For example, suppose that in the corpus:

- “students opened their” occurred 1000 times
- “students opened their books” occurred 400 times
 - $\rightarrow P(\text{books} \mid \text{students opened their}) = 0.4$
- “students opened their exams” occurred 100 times
 - $\rightarrow P(\text{exams} \mid \text{students opened their}) = 0.1$

Problem with the Markov Assumption

~~As the proctor started the clock, the students opened their~~ _____



For example, suppose that in the corpus:

- “students opened their” occurred 1000 times
- “students opened their **books**” occurred **400** times
 - $\rightarrow P(\text{books} \mid \text{students opened their}) = 0.4$
- “students opened their **exams**” occurred **100** times
 - $\rightarrow P(\text{exams} \mid \text{students opened their}) = 0.1$

Should we have discarded the “proctor” context?

Sparsity Issues with Ngrams

Sparsity Problem 1

Problem: What if “*students opened their w*” never occurred in data? Then w has probability 0!

(Partial) Solution: Add small δ to the count for every $w \in V$. This is called *smoothing*.

$$P(w|\text{students opened their}) = \frac{\text{count}(\text{students opened their } w)}{\text{count}(\text{students opened their})}$$

Sparsity Problem 2

Problem: What if “*students opened their*” never occurred in data? Then we can't calculate probability for any w !

(Partial) Solution: Just condition on “*opened their*” instead. This is called *backoff*.

Markov assumption was introduced in the first place to avoid sparsity

Storage Issues with Ngrams

Storage: Need to store count for all n -grams you saw in the corpus.

$$P(\mathbf{w}|\text{students opened their}) = \frac{\text{count(students opened their } \mathbf{w}\text{)}}{\text{count(students opened their)}}$$

Increasing n or increasing corpus increases model size!

Ngram LM in Practice

We can build a simple trigram language model over a 1.7 million word corpus (Reuters) in a few seconds on your laptop

- here is an example code: <https://nlpforhackers.io/language-models/>

today the _____

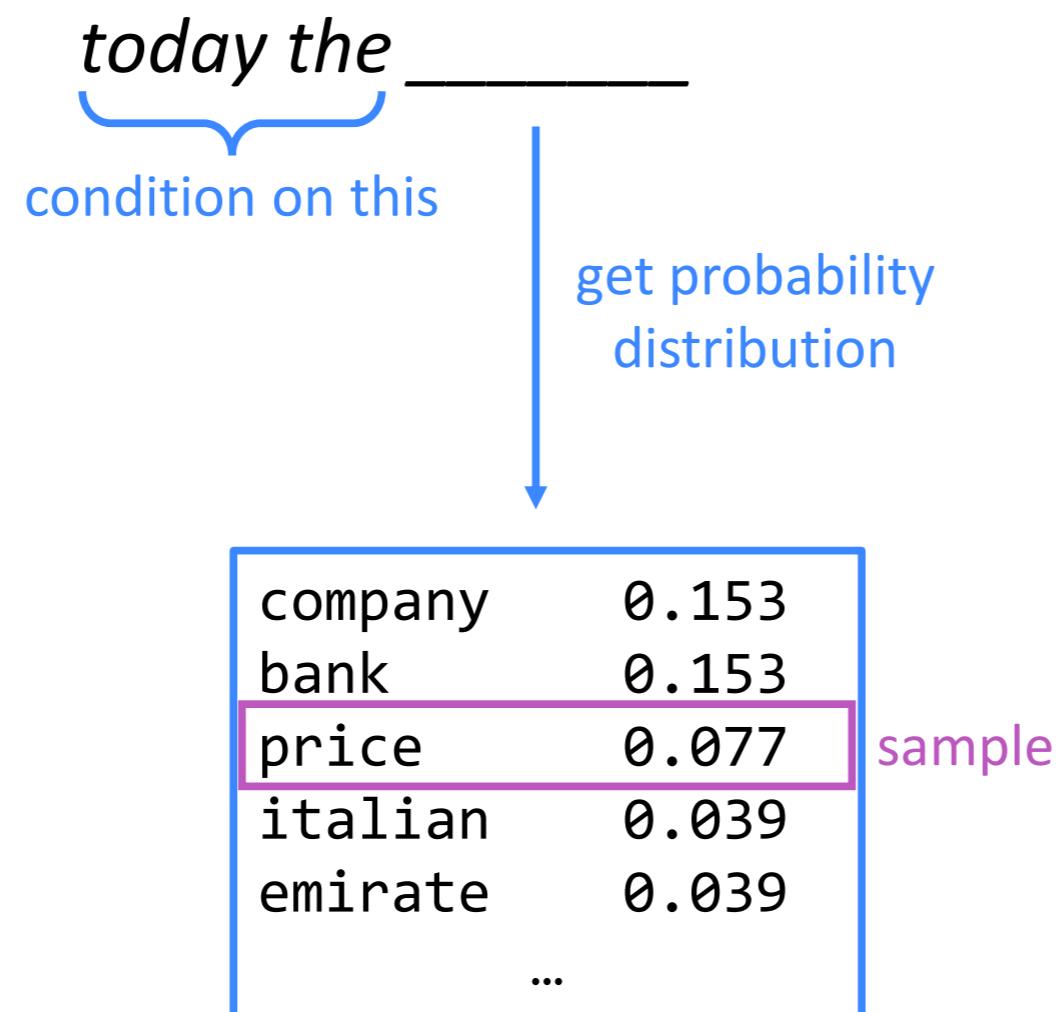
get probability distribution

company	0.153
bank	0.153
price	0.077
italian	0.039
emirate	0.039
...	

Sparsity problem:
not much granularity
in the probability
distribution

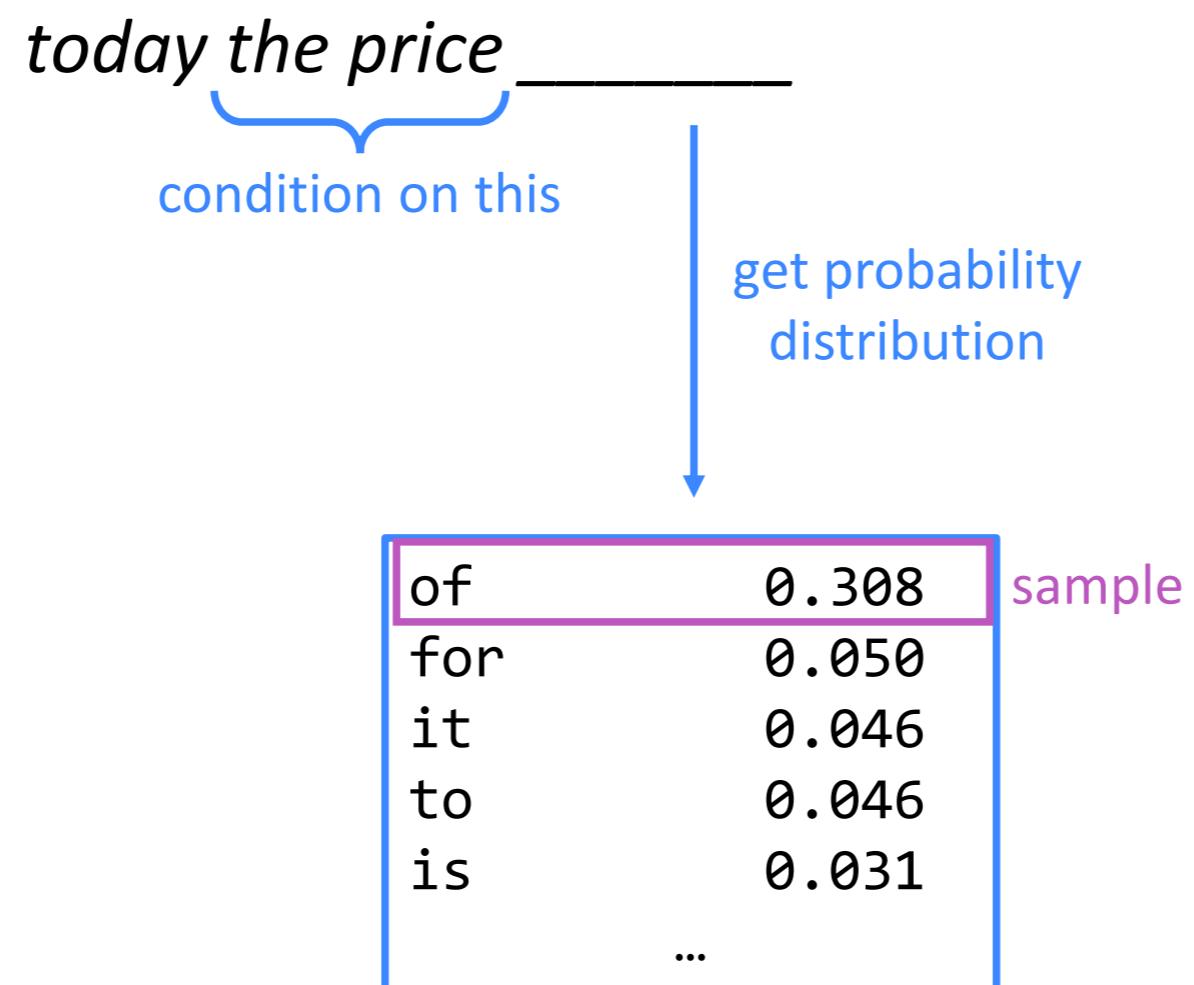
Ngram LM in Practice

We can also use a language model to generate text.



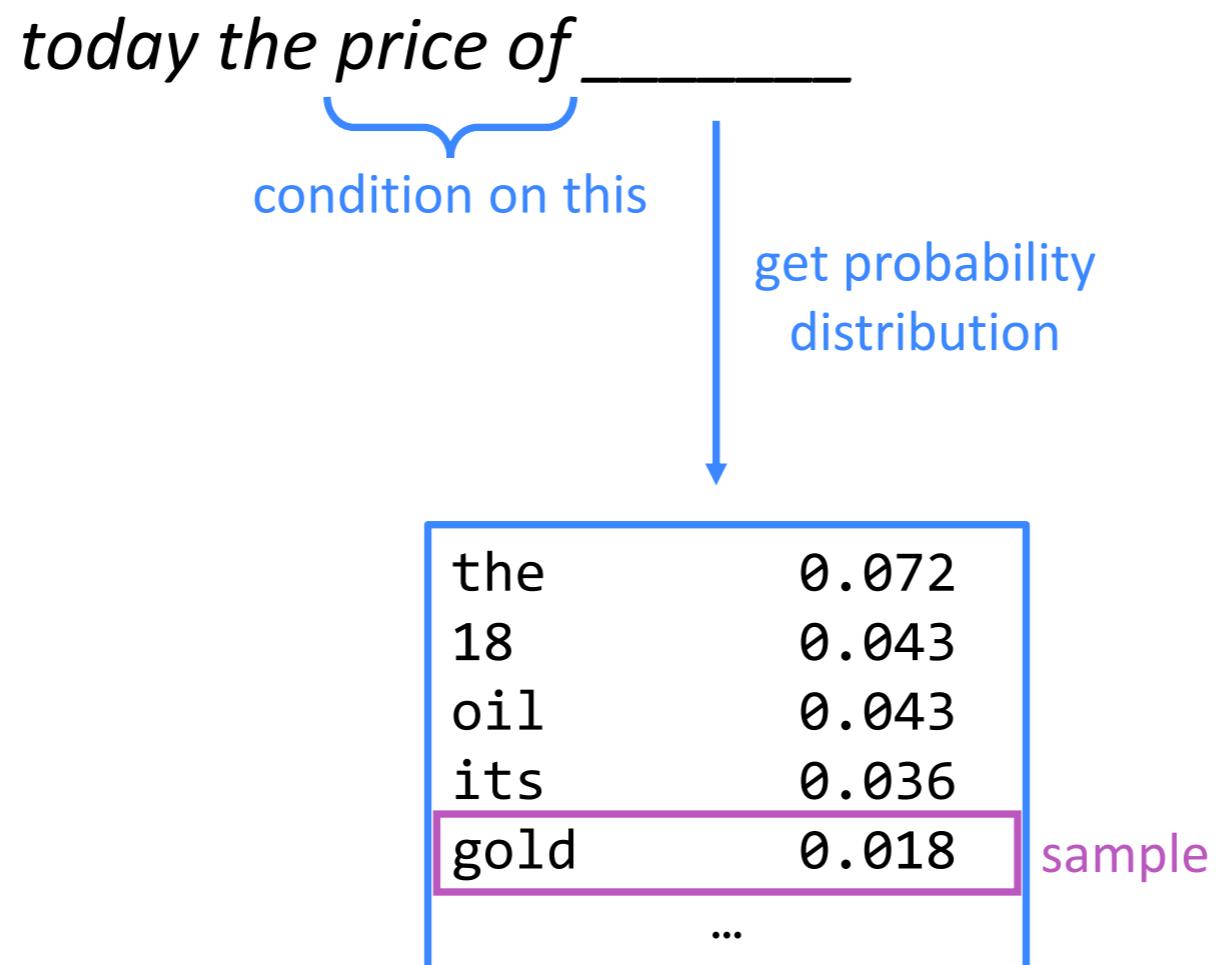
Ngram LM in Practice

We can also use a language model to generate text.



Ngram LM in Practice

We can also use a language model to generate text.



Ngram LM in Practice

If we keep going.....

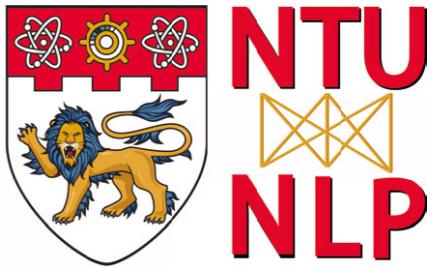
*today the price of gold per ton , while production of shoe
lasts and shoe industry , the bank intervened just after it
considered and rejected an imf demand to rebuild depleted
european stocks , sept 30 end primary 76 cts a share .*

Surprisingly grammatical!

...but **incoherent**. We need to consider more than three words at a time if we want to model language well.

But increasing n worsens sparsity problem, and increases model size...

Window-based Language Model (Revisited)



output distribution

$$\hat{y} = \text{softmax}(Uh + b_2) \in \mathbb{R}^{|V|}$$

hidden layer

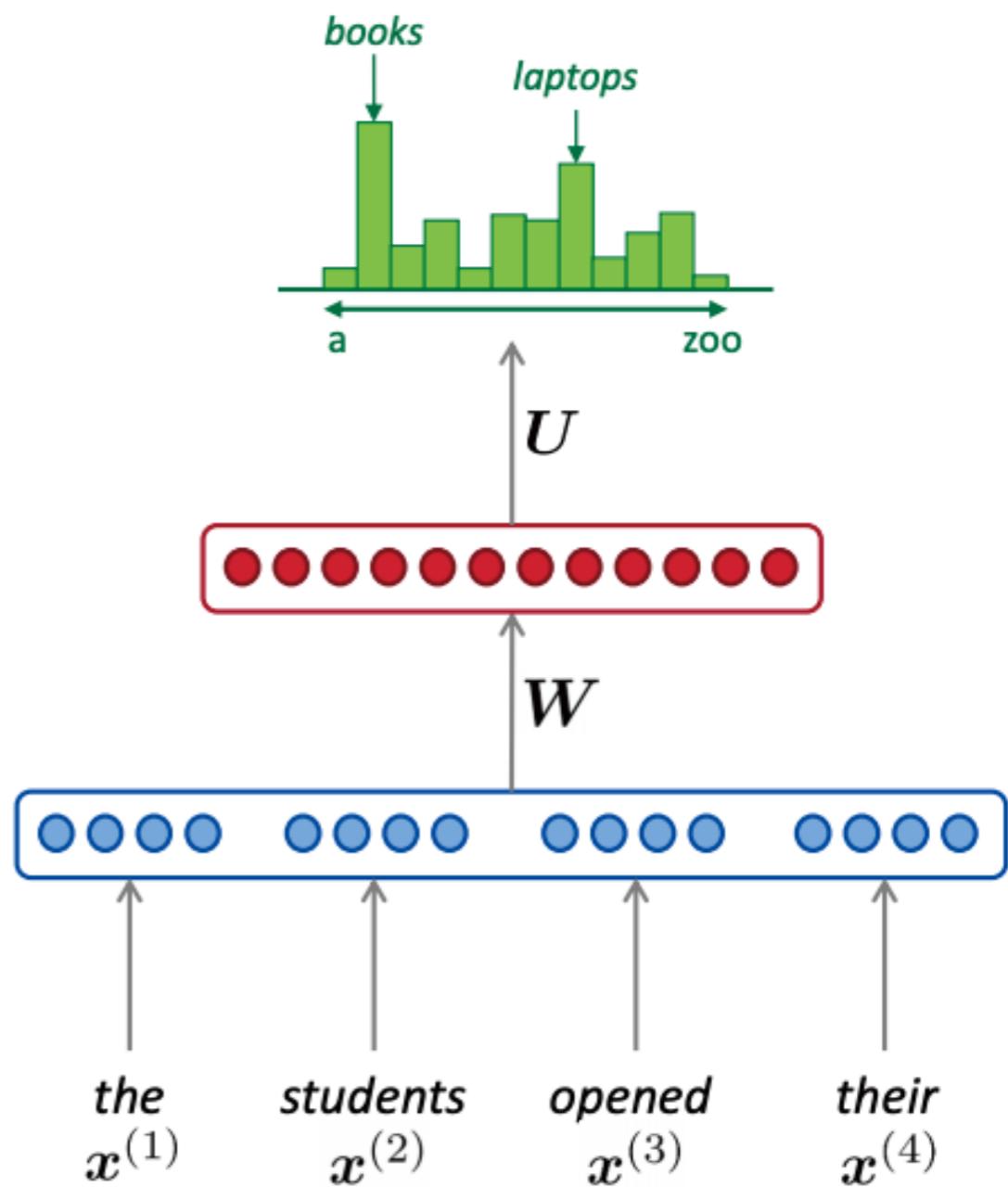
$$h = f(We + b_1)$$

concatenated word embeddings

$$e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$$

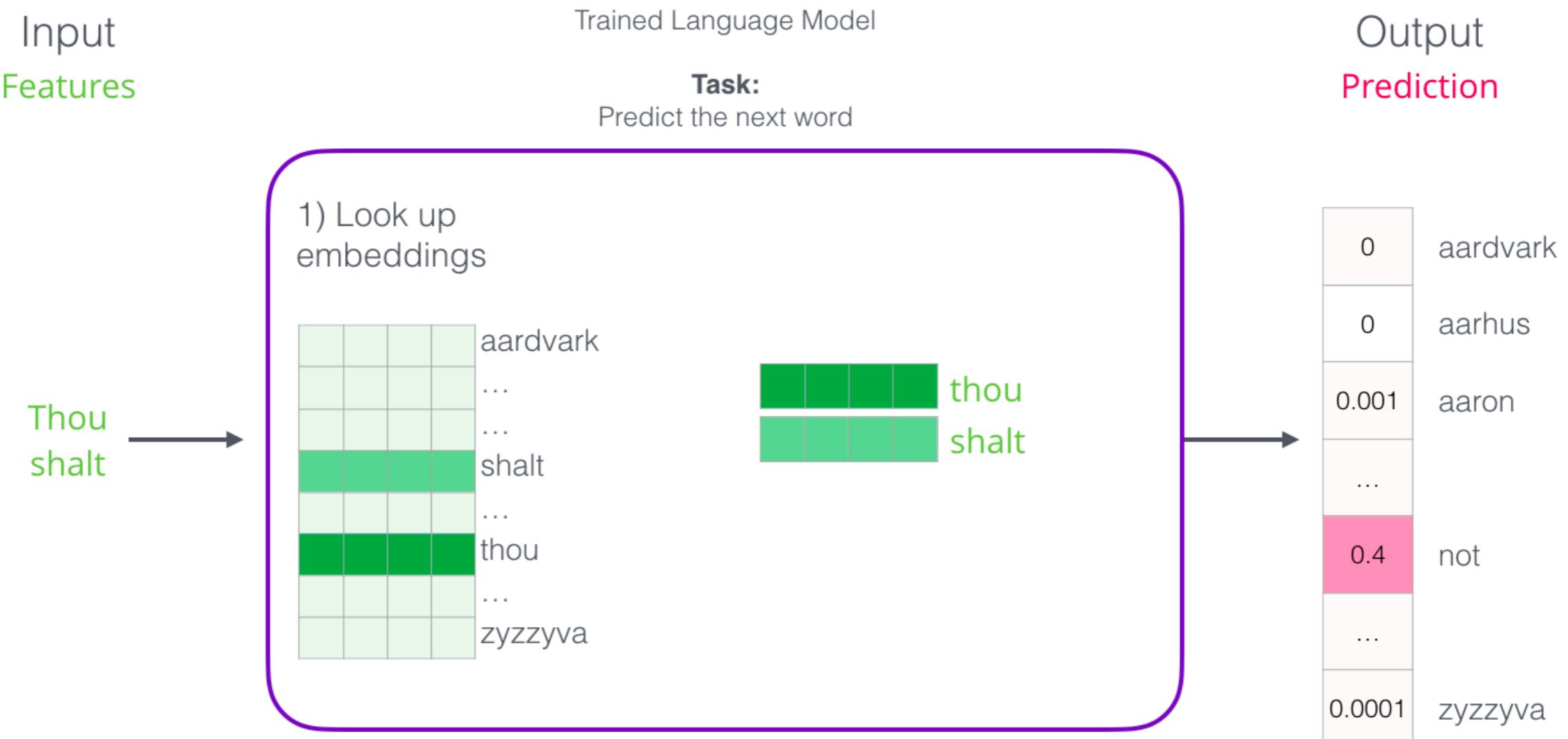
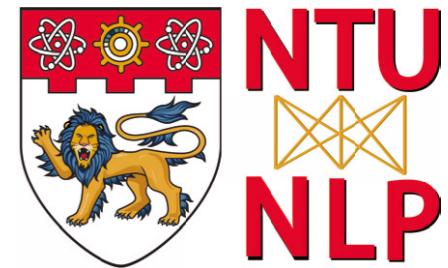
words / one-hot vectors

$$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$$

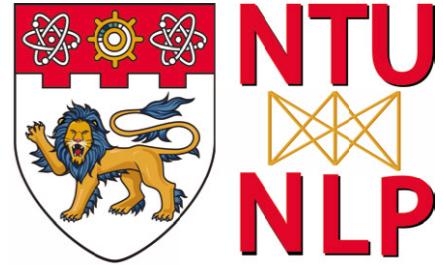


Source: Bengio, 2003

Window-based Language Model (Revisited)



Named Entity Recognition



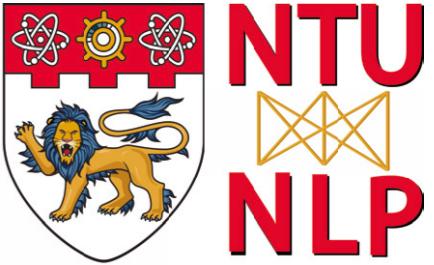
Named Entity Recognition (NER): Identify and classify entities such as persons, locations and organisations.

In 1917, Einstein applied the general theory of relativity to model the large-scale structure of the universe. He was visiting the United States when Adolf Hitler came to power in 1933 and did not go back to Germany, where he had been a professor at the Berlin Academy of Sciences. He settled in the U.S., becoming an American citizen in 1940. On the eve of World War II, he endorsed a letter to President Franklin D. Roosevelt alerting him to the potential development of "extremely powerful bombs of a new type" and recommending that the U.S. begin similar research. This eventually led to what would become the Manhattan Project. Einstein supported defending the Allied forces, but largely denounced using the new discovery of nuclear fission as a weapon. Later, with the British philosopher Bertrand Russell, Einstein signed the Russell-Einstein Manifesto, which highlighted the danger of nuclear weapons. Einstein was affiliated with the Institute for Advanced Study in Princeton, New Jersey, until his death in 1955.

Tag colours:

LOCATION TIME PERSON ORGANIZATION MONEY PERCENT DATE

Named Entity Recognition



NER in other Language and Domains

Spanish

B-PER I-PER O O O B-ORG I-ORG I-ORG
Jaime Matos , de los Rebeldes de Moca

BIO encoding

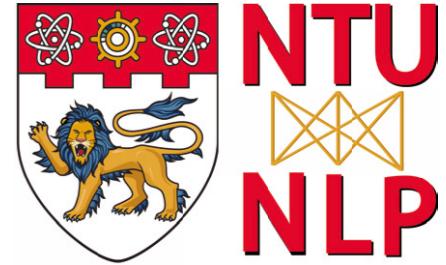
Fin.

B-ORG I-ORG I-ORG O O B-LOC O B-ORG
Silicon Valley bank , a California chartered bank ..

NER has many applications

- Knowledge graph
- Search
- Translation

Why NER Is Hard?



First National Bank Donates 2 Vans To Future School Of Fort Smith

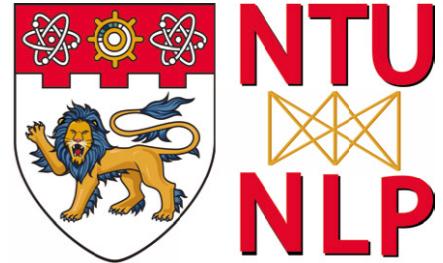
POSTED 3:43 PM, JANUARY 11, 2019, BY SNEWS WEB STAFF

- Hard to know if something is a name; Is there a school called “Future School” or is it a future school
- Hard to find name boundaries; Is the first entity “First National Bank” or “National Bank”
- Hard to know class of unknown/novel name; what class is “Zig Ziglar”? (A person.)

To find out more about Zig Ziglar and read features by other Creators Syndicate writers and

- Entity class is ambiguous and depends on context; Charles Schwab is a Person or an Org? where Larry Ellison and Charles Schwab can live discreetly amongst wooded estates. And

Other Similar NLP Tasks



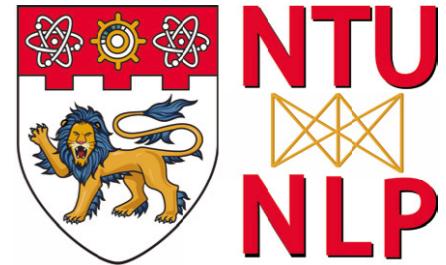
Chunking

Tokens	Most	large	cities	in	the	US	had	morning	and	afternoon	newspapers	.
POS Tags	JJS	JS	NNS	IN	DT	NNP	VBD	NN	CC	NN	NNS	.
Chunk IDs	B-NP	I-NP	I-NP	B-NP	B-NP	I-NP	B-NP	B-NP	I-NP	I-NP	I-NP	B-NP

The diagram illustrates the six chunks identified by the chunk IDs. Dashed green lines connect the start of each chunk to its corresponding token in the sequence:

- 1st chunk:** Most
- 2nd chunk:** large cities
- 3rd chunk:** in the US
- 4th chunk:** had morning
- 5th chunk:** and afternoon
- 6th chunk:** newspapers.

Other Similar NLP Tasks



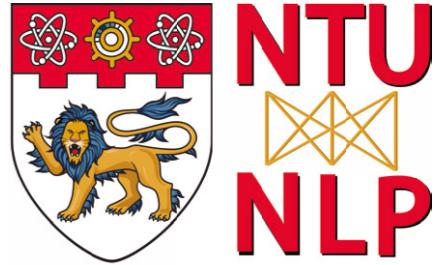
Part-of-Speech (POS) Tagging

The_DT first_JJ time_NN he_PRP was_VBD shot_VBN in_IN the_DT
hand_NN as_IN he_PRP chased_VBD the_DT robbers_NNS outside_RB ...

first	time	shot	in	hand	as	chased	outside
JJ	NN	NN	IN	NN	IN	JJ	IN
RB	VB	VBD	RB	VB	RB	VBD	JJ
		VBN	RP			VBN	NN
							RB

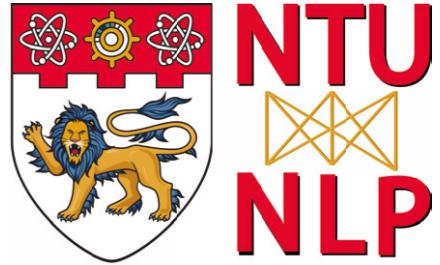
Similar to NER, both POS tagging and Chunking have many applications in other NLP tasks and downstream applications

Window-based Approach



- In general, classification of one token depends in its context.
- Example: resolving ambiguous named entities:
 - Paris → Paris, France vs. Paris Hilton vs. Paris, Texas
 - Hathaway → Berkshire Hathaway vs. Anne Hathaway
 - Washington → Washington DC vs. Denzel Washington
- A window-based approach involves:
 - Extracting a sequence of words from the text.
 - Computing a vector representation for this window.
 - And to classify the average vector.

Window-based Approach



- Train a classifier to classify a center word by taking words from its neighborhood.
 - We can put hidden layers to make the model deep.
 - Example: Classify “Paris” in with window length of 2:

... museums in Paris are amazing ...

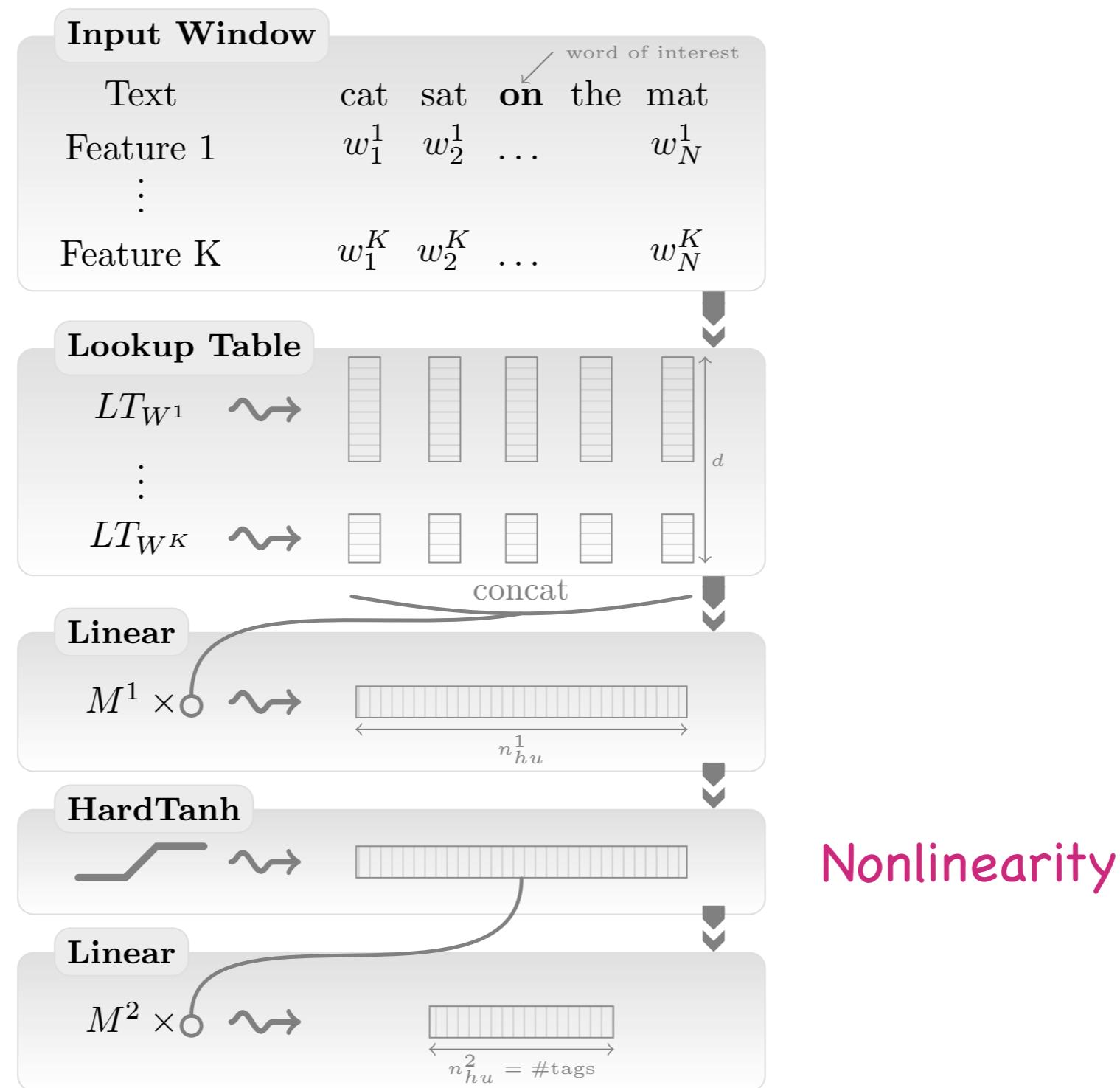
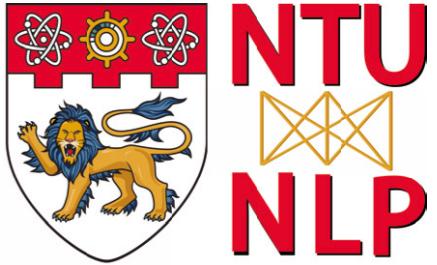
A horizontal row consisting of five separate groups. Each group contains four solid red circles arranged in a horizontal line. The groups are evenly spaced and aligned horizontally.

$$X_{\text{window}} = [x_{\text{museums}} \quad x_{\text{in}} \quad x_{\text{Paris}} \quad x_{\text{are}} \quad x_{\text{amazing}}]^T$$

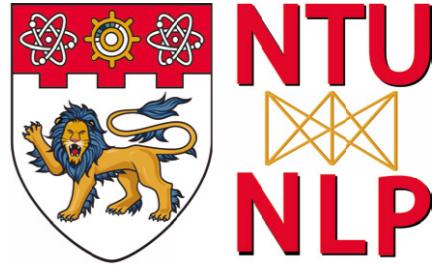
$$x_{\text{window}} = x \in \mathbb{R}^{5d}$$

- Similar to word2vec, we go over all positions in a sentence. But this time, it is supervised.

Window-based Approach



Window-based Approach



- We can use the same **Softmax** as the final classification layer:

$$\mathcal{P}(y = c | \mathbf{x}, \mathbf{W}) = \frac{\exp(\mathbf{w}_c^\top \mathbf{x})}{\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x})}$$

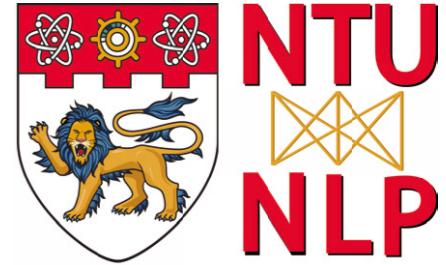
W is final layer
weights

- Use the same **Cross Entropy** loss:

$$\ell(\mathbf{W}) = \sum_{i=1}^N \sum_{c=1}^C y_{ic} \log \mu_{ic}$$

$$= \sum_{i=1}^N \left[\left(\sum_{c=1}^C y_{ic} \mathbf{w}_c^\top \mathbf{x}_i \right) - \log \left(\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x}_i) \right) \right]$$

Where we are



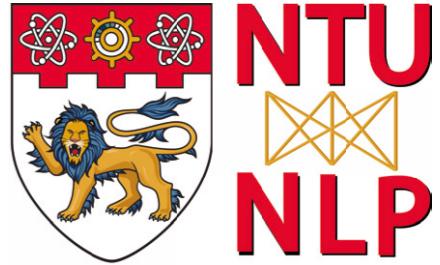
Models/Algorithms

- Linear models
- Feed-forward Neural Nets (FNN)
- Window-based methods
- Convolutional Nets

NLP tasks/applications

- Word meaning
- Language modelling
- Sequence tagging

Convolutional Neural Nets



Convolutional neural network (CNN, or ConvNet) is a type of feed-forward network where the individual neurons are tiled in such a way that they respond to overlapping regions (context windows) in the input.

In other words, CNN is a type of feed-forward neural network with

- Local connectivity (with the input layer)
- Share parameters across spatial/temporal positions

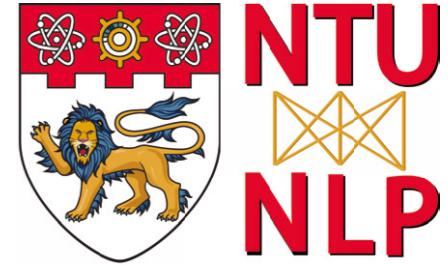
Advantage over FNN

- Requires less parameters (sharing)

Advantage over Recurrent Neural Net (next lecture)

- Parallelizable

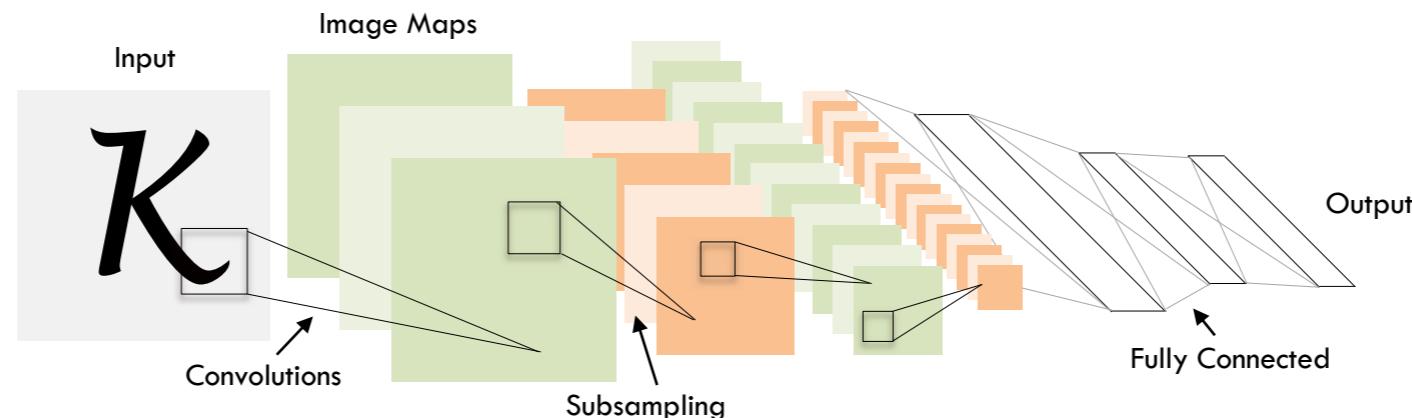
CNNs in Computer Vision



- CNNs were responsible for major breakthroughs in Image Classification (e.g., face recognition) and are the core of most Computer Vision systems today.

1998

LeCun et al.



of transistors



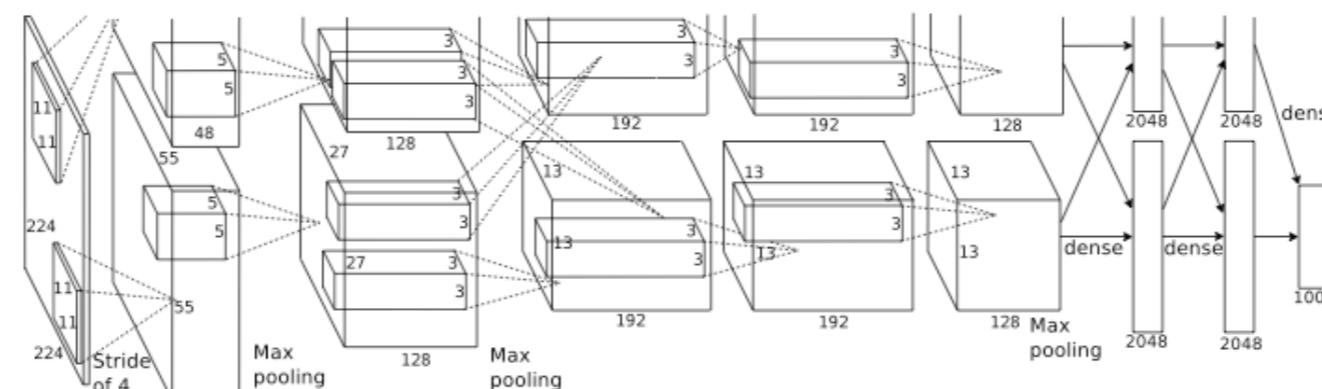
10^6

of pixels used in training

10^7 NIST

2012

Krizhevsky et al.



of transistors



10^9

GPUs



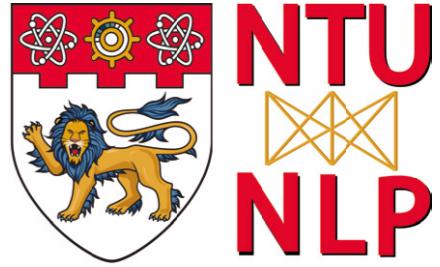
of pixels used in training

10^{14} IMAGENET

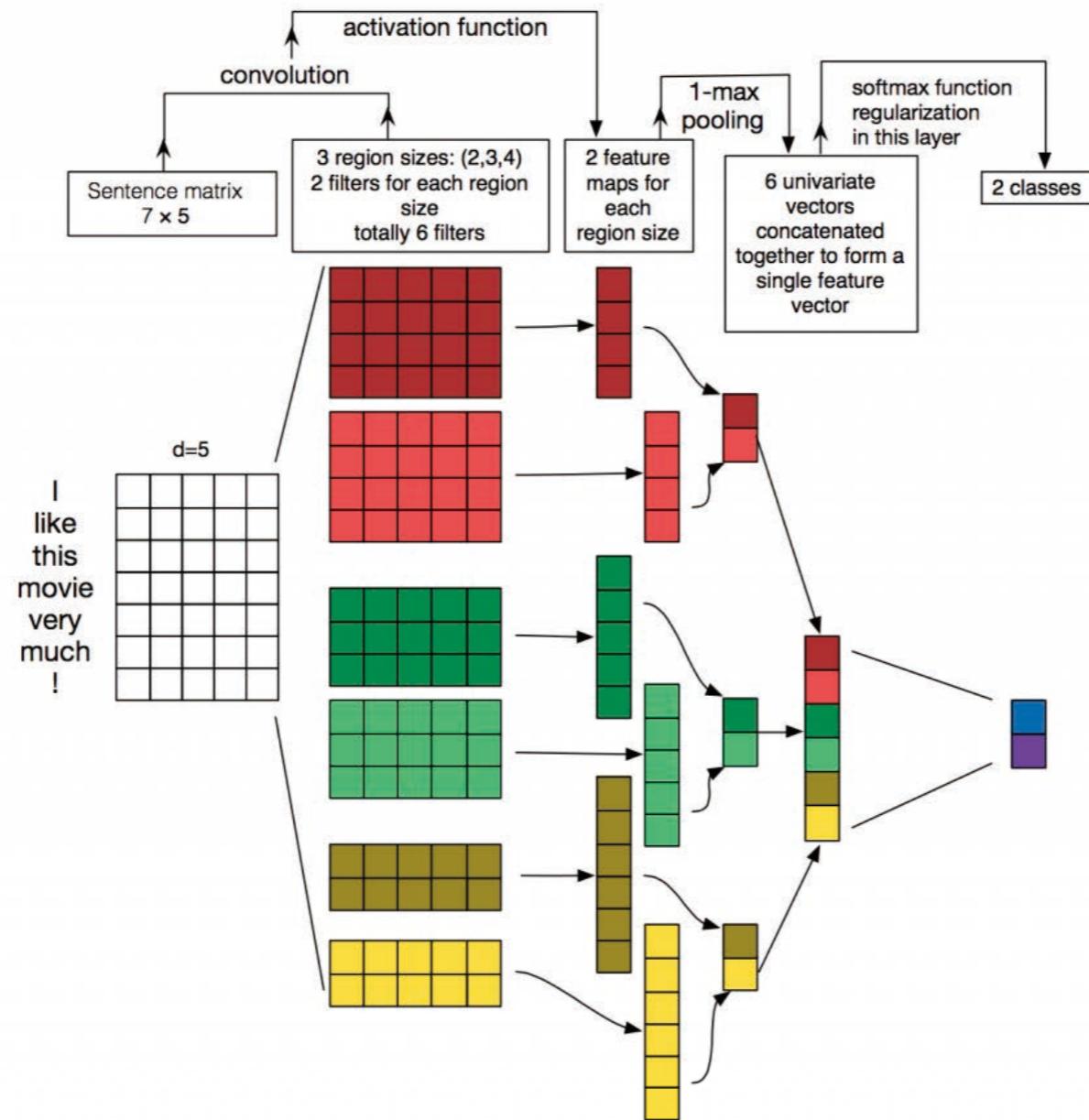
Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Source: cs231n

CNNs in NLP



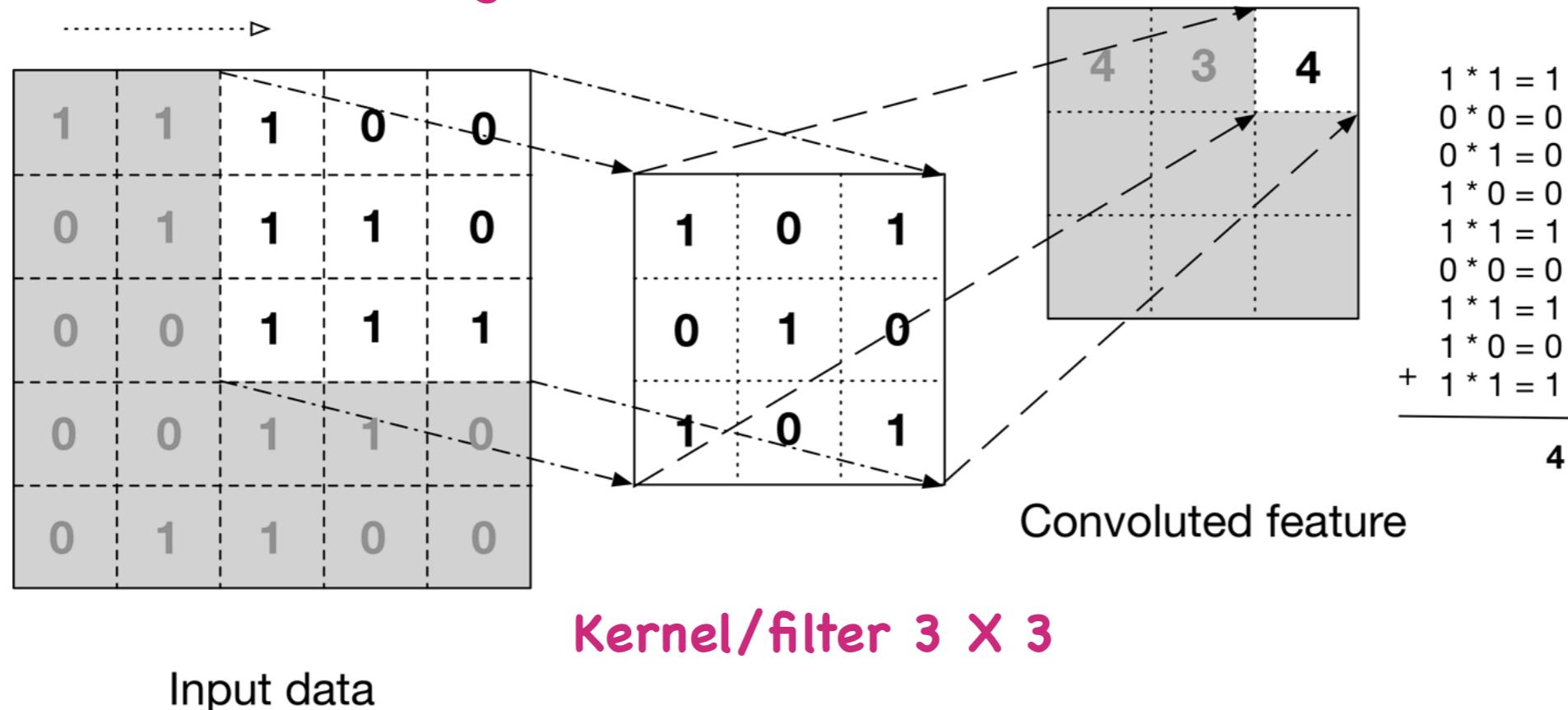
- CNNs have in the past years shown break-through results in some NLP tasks, especially classifying phrases and sentences.



Convolution in Image

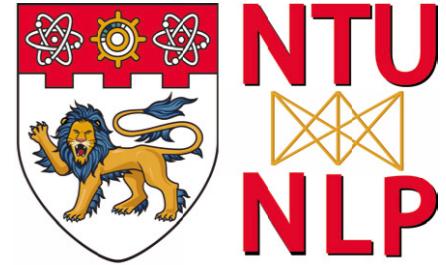
A CNN typically involves two operations: **convolution** and **pooling**.

● (2D) Convolution in image



- The kernel is moved **s** pixels to the right/down, where **s** is the **stride length**
- The output of this process is a matrix with all its entries filled, called the **convolved feature** or **input feature map**
- An input image can be convolved with multiple convolution kernels at once, creating one output (feature map) for each kernel.

Convolution in NLP



Convolution in NLP

- Generally 1D (temporal)
- For each window, it is same as our window-based FFN method
- One Conv. operation

$$c_i = f(\mathbf{w} \cdot \mathbf{x}_{i:i+h-1} + b)$$

filter $\mathbf{w} \in \mathbb{R}^{hk}$

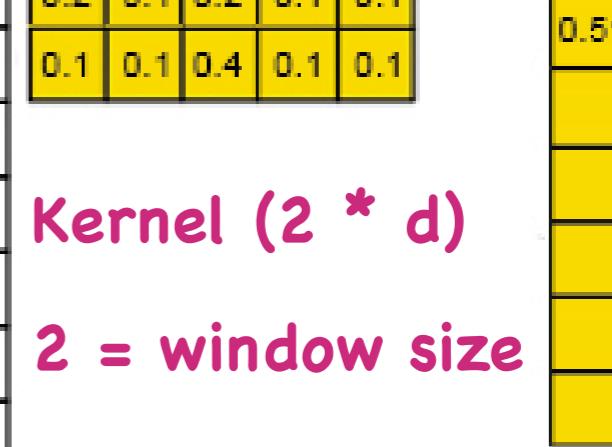
- One feature map

$$\mathbf{c} = [c_1, c_2, \dots, c_{n-h+1}]$$

I
like
this
movie
very
much
!

0.6	0.5	0.2	-0.1	0.4
0.8	0.9	0.1	0.5	0.1
0.4	0.6	0.1	-0.1	0.7
---	---	---	---	---
---	---	---	---	---
---	---	---	---	---
---	---	---	---	---

0.2	0.1	0.2	0.1	0.1
0.1	0.1	0.4	0.1	0.1



Word vectors

0.6	0.5	0.2	-0.1	0.4
0.8	0.9	0.1	0.5	0.1
0.4	0.6	0.1	-0.1	0.7
---	---	---	---	---
---	---	---	---	---
---	---	---	---	---
---	---	---	---	---

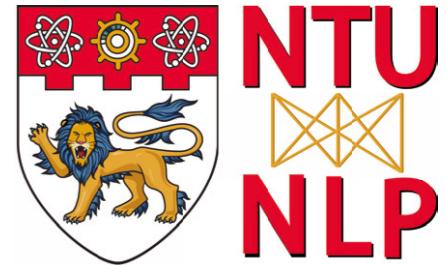
0.2	0.1	0.2	0.1	0.1
0.1	0.1	0.4	0.1	0.1

One Feature map



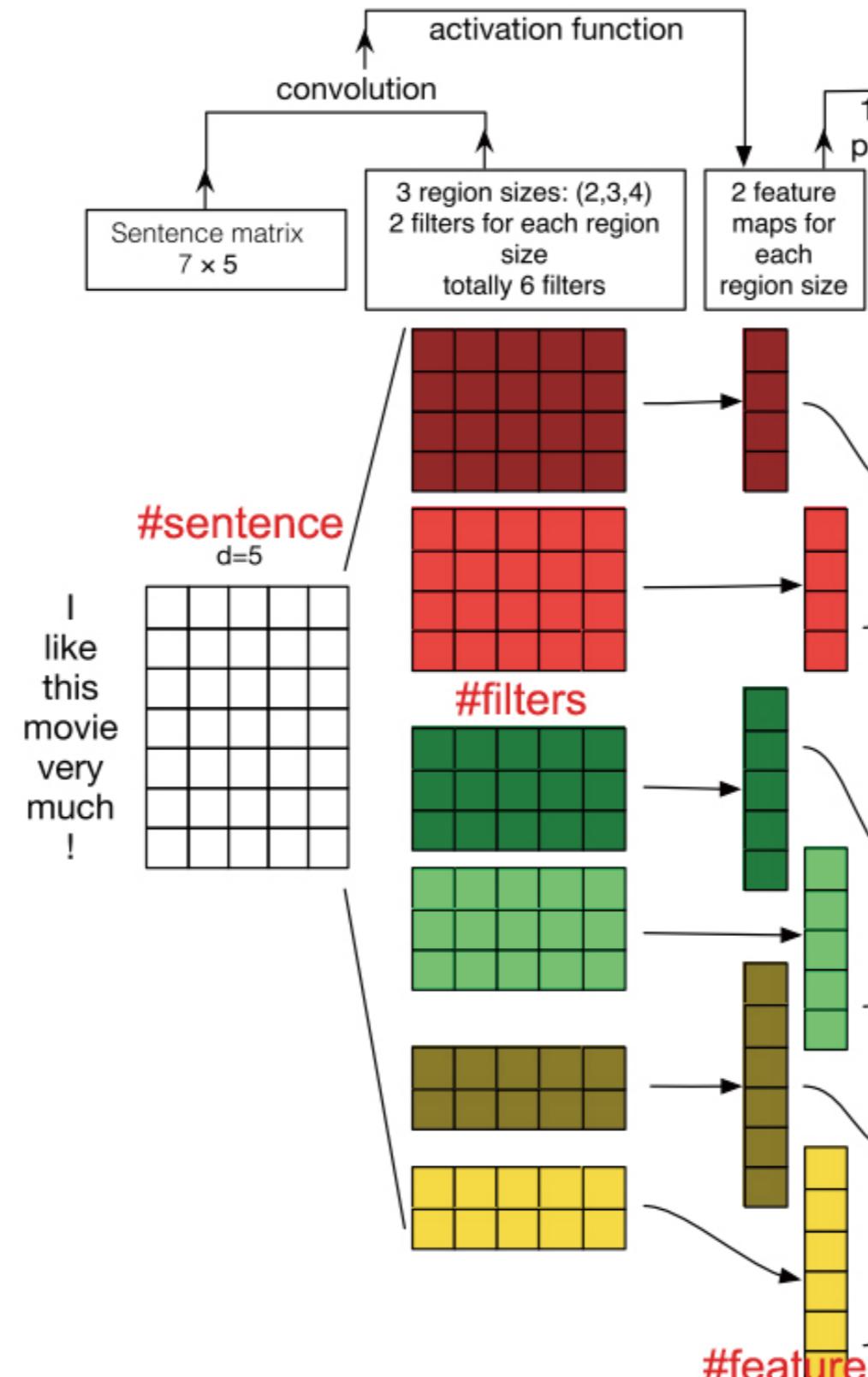
Convolution with 1 filter

Convolution in NLP



Convolution with multiple filters

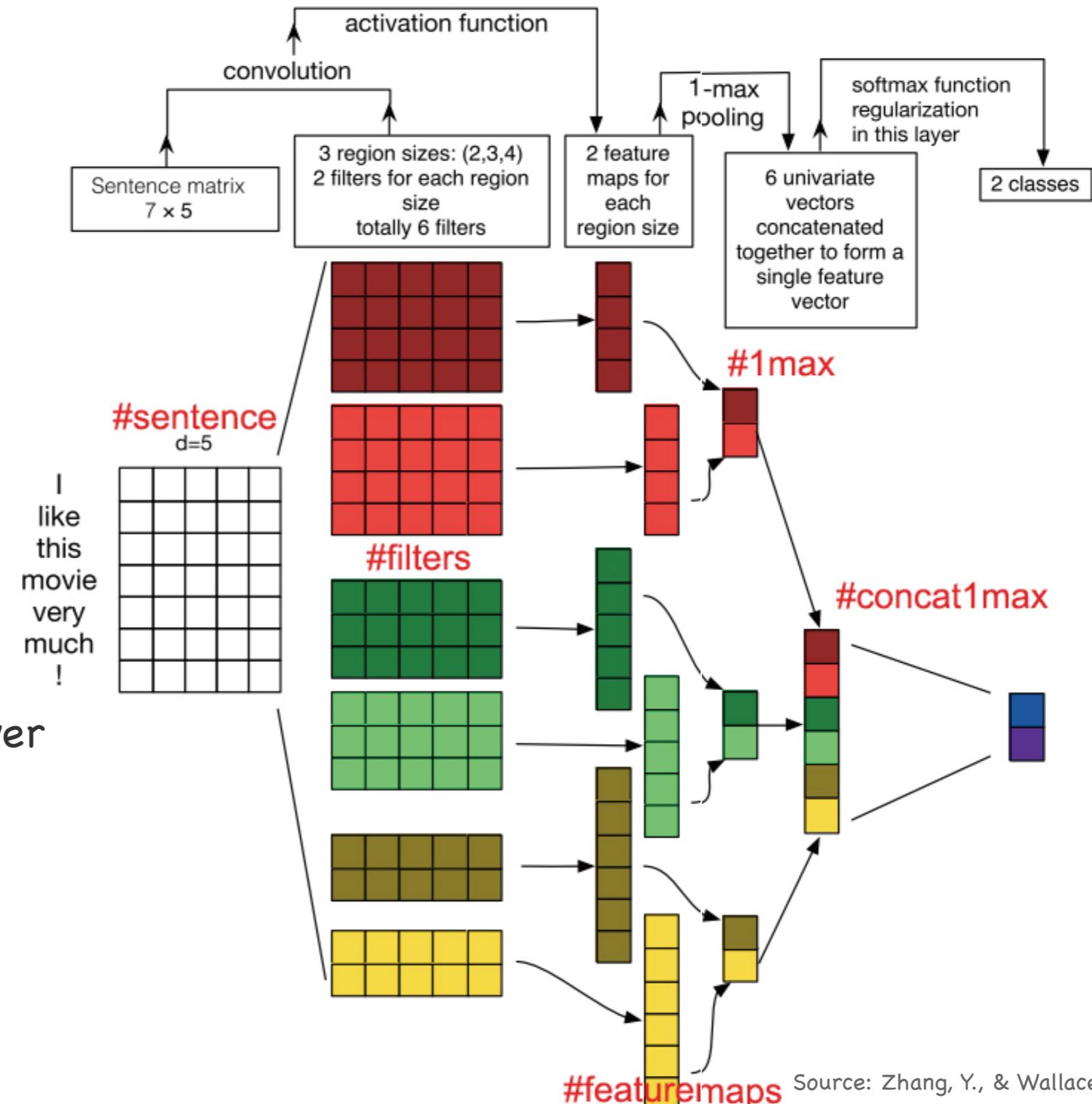
Multiple filters (with varying window sizes) to obtain multiple features



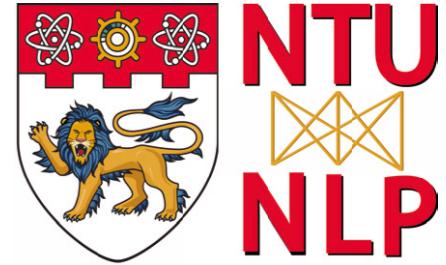
Two Basic Operations in CNNs

Pooling

- Mean, Max
- Performed over convoluted features for
 - Dimension reduction
 - Selecting important features
- No parameter is added
- Often followed by a Fully connected layer to ensure fixed length input to final layer

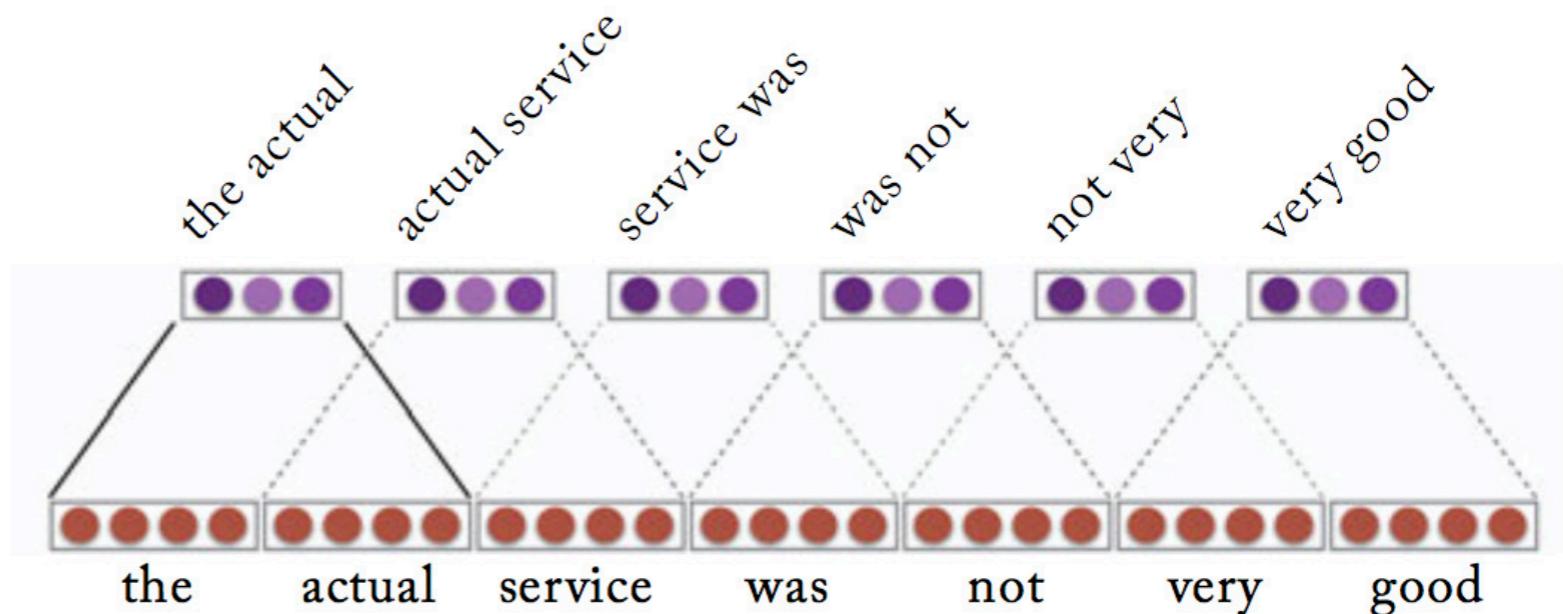


CNN for Tagging Tasks



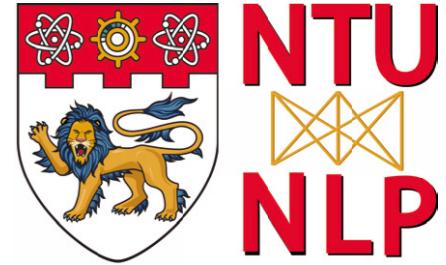
CNN applied to phrases in a sentence

- Apply to POS tagging, Chunking, NER
- We can add more hidden (CNN, FNN) layers
- Final layer is a softmax layer over the possible tags
- Extract meaningful sub-structures (n-grams) that are useful for the overall prediction task
- Also can be applied to character level to capture morphological features



Convolution with **window size=2** and dimensional output (number of filters/output channel) = 3.

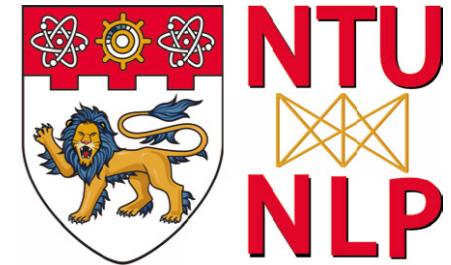
CNN for Tagging Tasks



- Apply to POS tagging, Chunking, NER
- We can add more hidden (CNN, FNN) layers
- Final layer is a softmax layer over the possible tags
- Extract meaningful sub-structures (n-grams) that are useful for the overall prediction task
- Also can be applied to character level to capture morphological features

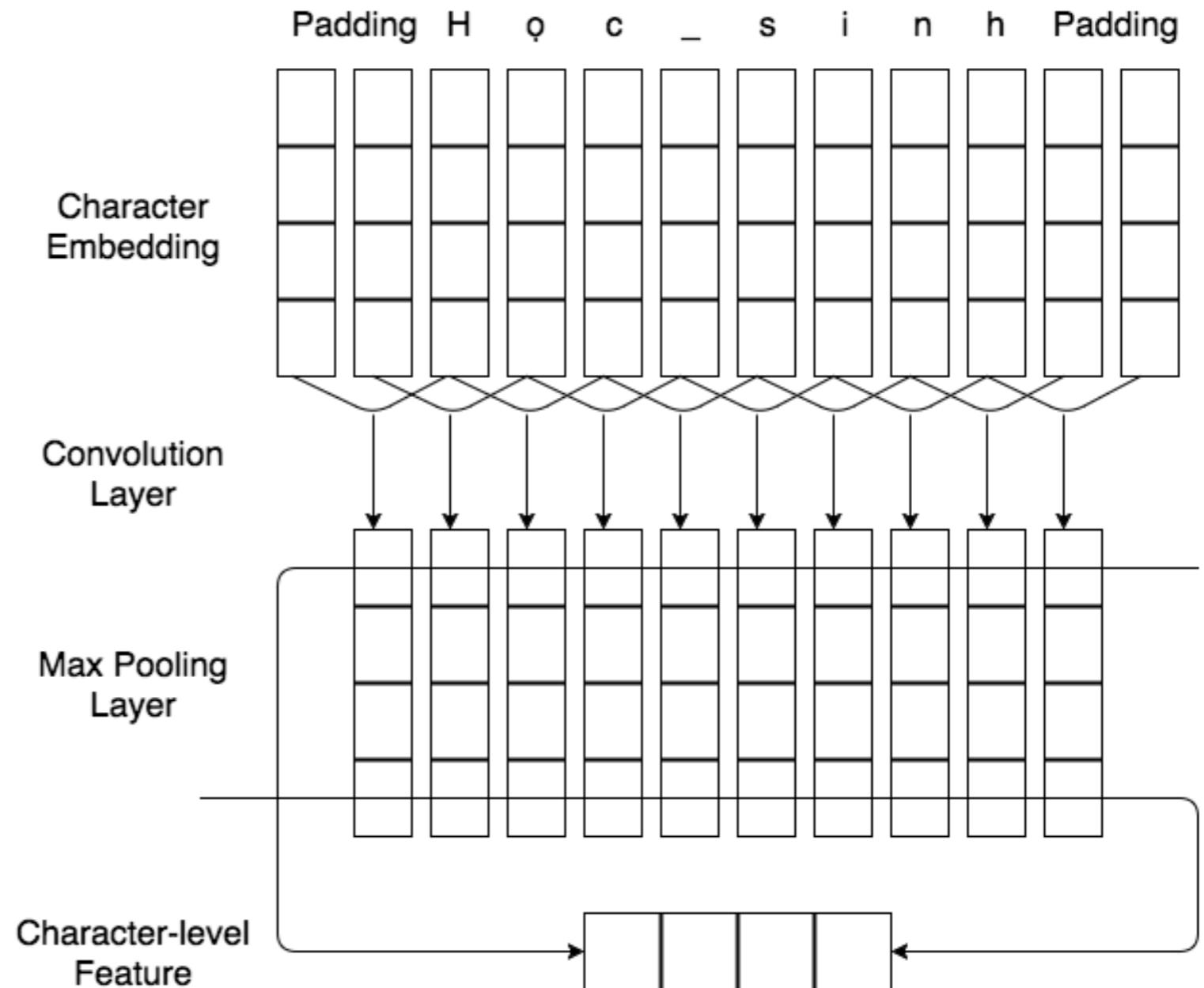
Approach	POS (PWA)	CHUNK (F1)	NER (F1)	SRL (F1)
Benchmark Systems	97.24	94.29	89.31	77.92
NN+WLL	96.31	89.13	79.53	55.40
NN+SLL	96.37	90.33	81.47	70.99
NN+WLL+LM1	97.05	91.91	85.68	58.18
NN+SLL+LM1	97.10	93.65	87.58	73.84
NN+WLL+LM2	97.14	92.04	86.96	58.34
NN+SLL+LM2	97.20	93.63	88.67	74.15

CNN for Encoding

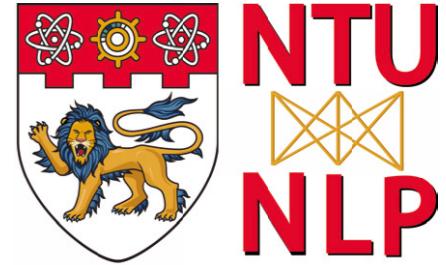


Character-level CNN

- CNNs can be applied to characters of a word
- Model morphology/subword
- Handle unknown (out of vocabulary) words
- Commonly used in Neural MT

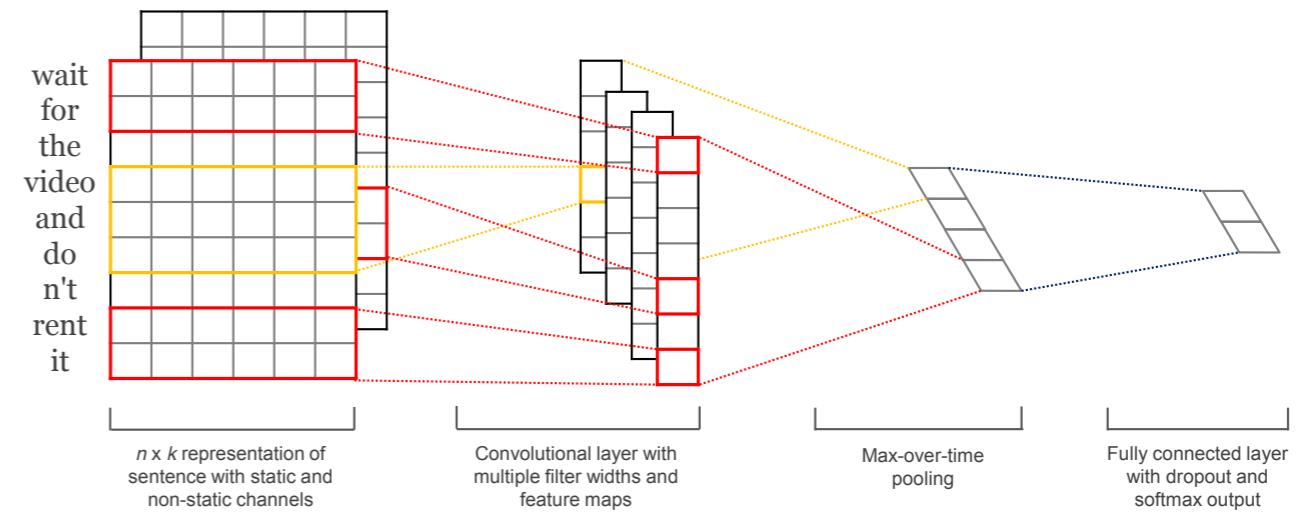


CNN for Tagging Tasks

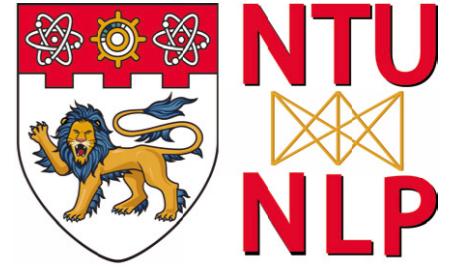


CNN for sentence modelling

- Apply to sentence classification tasks (e.g., sentiment)
- Final layer is a softmax layer over the possible classes
- Convolution extracts meaningful local (n-gram) features that are useful for the overall sentence level prediction task



Other CNN Types



Dilated Convolution in Image

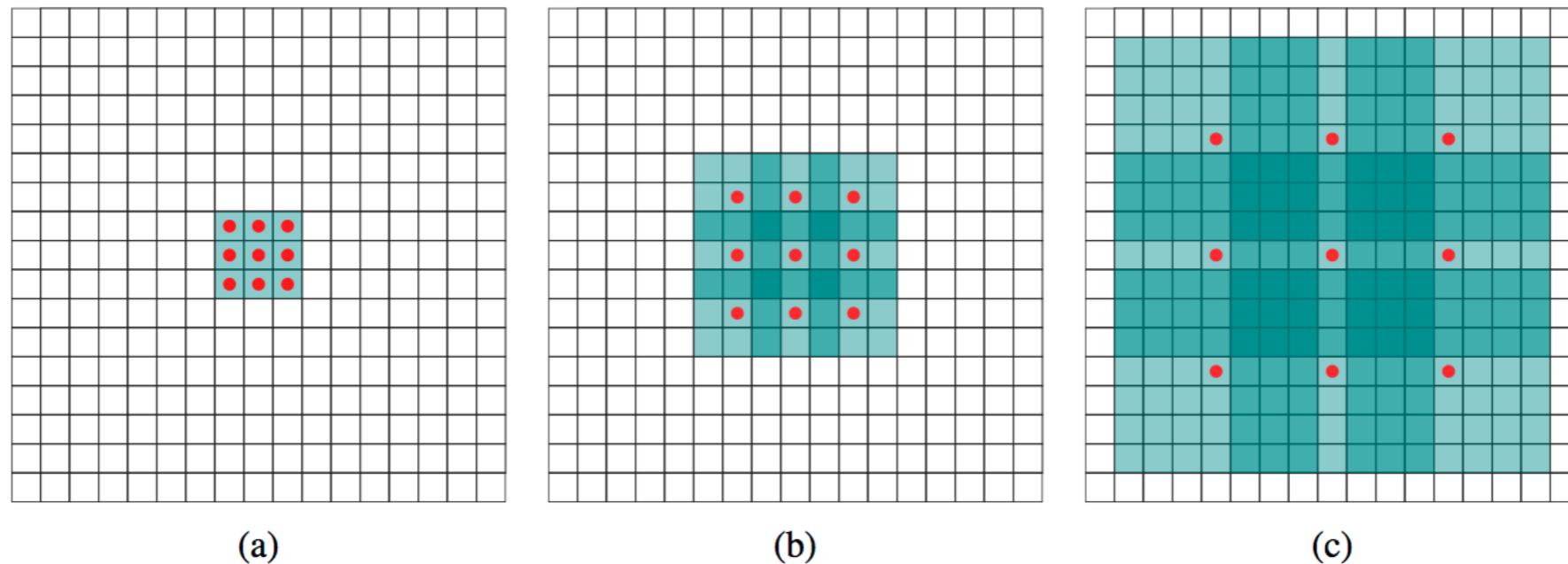
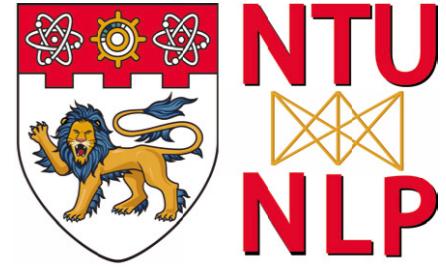


Figure 1: Systematic dilation supports exponential expansion of the receptive field without loss of resolution or coverage. (a) F_1 is produced from F_0 by a 1-dilated convolution; each element in F_1 has a receptive field of 3×3 . (b) F_2 is produced from F_1 by a 2-dilated convolution; each element in F_2 has a receptive field of 7×7 . (c) F_3 is produced from F_2 by a 4-dilated convolution; each element in F_3 has a receptive field of 15×15 . The number of parameters associated with each layer is identical. The receptive field grows exponentially while the number of parameters grows linearly.

Other CNN Types



Dilated Convolution in NLP

Regular convolution,

$$c_t = W_c \bigoplus_{k=0}^r x_{t \pm k}$$

Dilated Convolution,

$$c_t = W_c \bigoplus_{k=0}^r x_{t \pm k\delta}$$

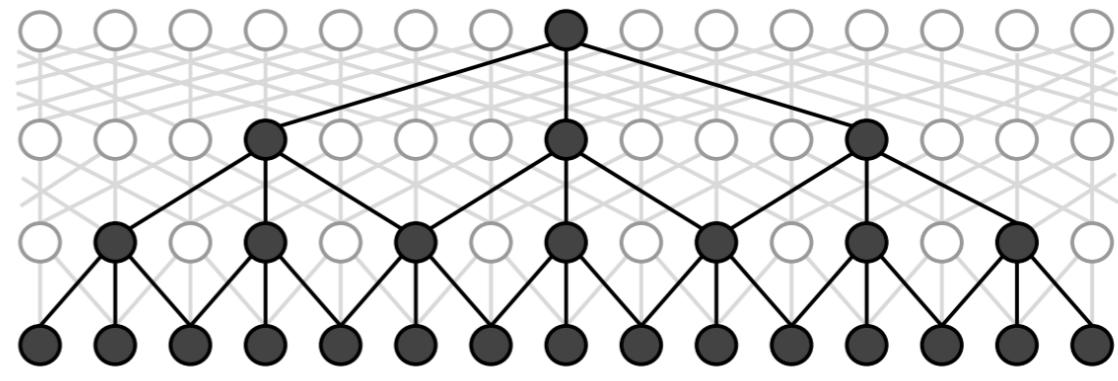
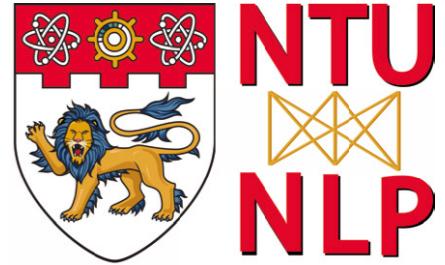


Figure: A dilated CNN block with maximum dilation width 4 and filter width 3.

$$r = l(w - 1) + 1$$

where , r number of node covered at layer l with filter width w .
total number of node at layer l is, $2^{l+1} - 1$ (**exponential growth**).

Other CNN Types



Dilated Convolution in NLP

Regular convolution,

$$c_t = W_c \bigoplus_{k=0}^r x_{t \pm k}$$

Dilated Convolution,

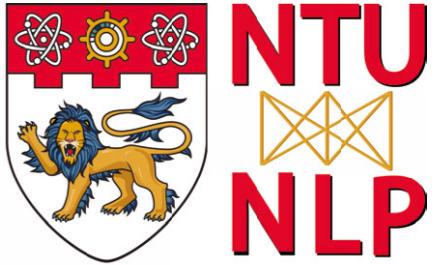
$$c_t = W_c \bigoplus_{k=0}^r x_{t \pm k\delta}$$

Model	Speed
Bi-LSTM-CRF	1×
Bi-LSTM	9.92×
ID-CNN-CRF	1.28×
5-layer CNN	12.38×
ID-CNN	14.10×

Model	F1
Ratinov and Roth (2009)	86.82
Collobert et al. (2011)	86.96
Lample et al. (2016)	90.33
Bi-LSTM	89.34 ± 0.28
4-layer CNN	89.97 ± 0.20
5-layer CNN	90.23 ± 0.16
ID-CNN	90.32 ± 0.26
Collobert et al. (2011)	88.67
Passos et al. (2014)	90.05
Lample et al. (2016)	90.20
Bi-LSTM-CRF (re-impl)	90.43 ± 0.12
ID-CNN-CRF	90.54 ± 0.18

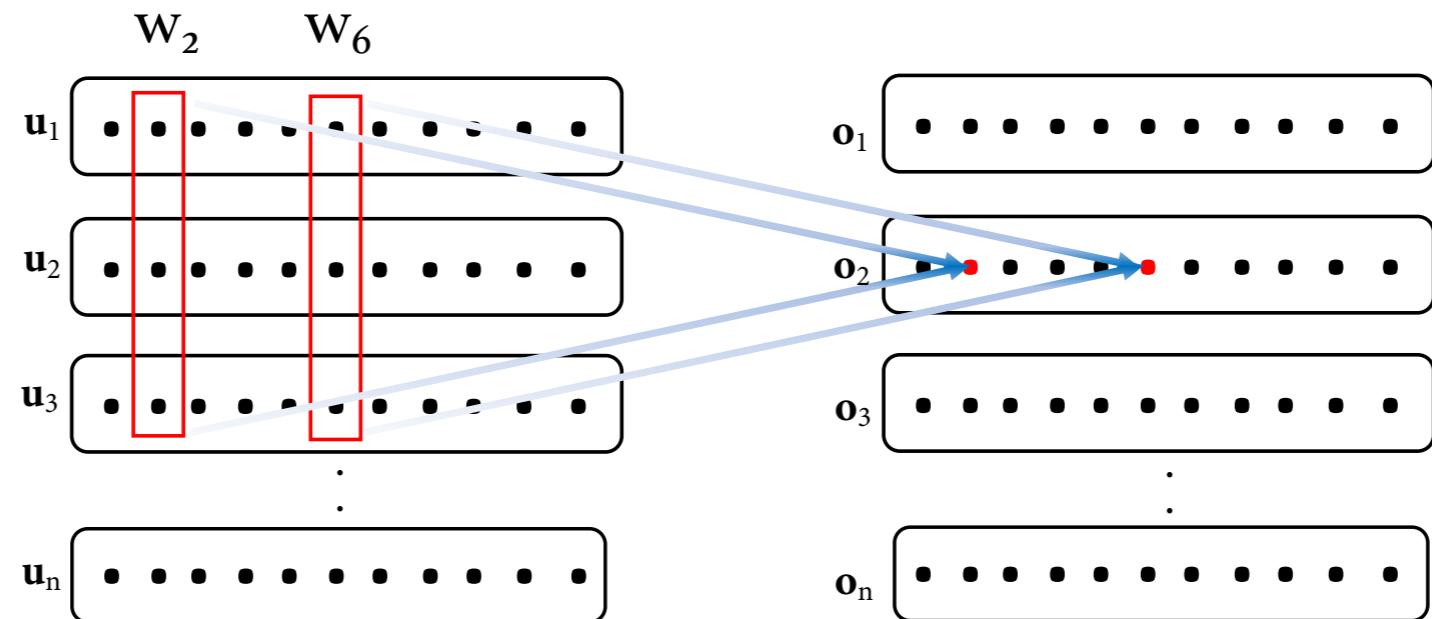
NER Results

Other CNN Types



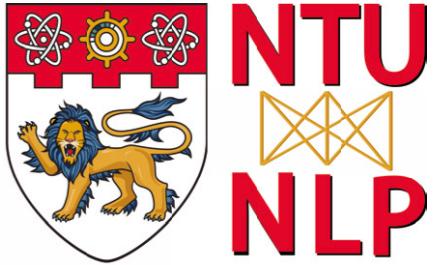
Depth-wise Convolution in NLP

The convolutions are done over the input dimensions.



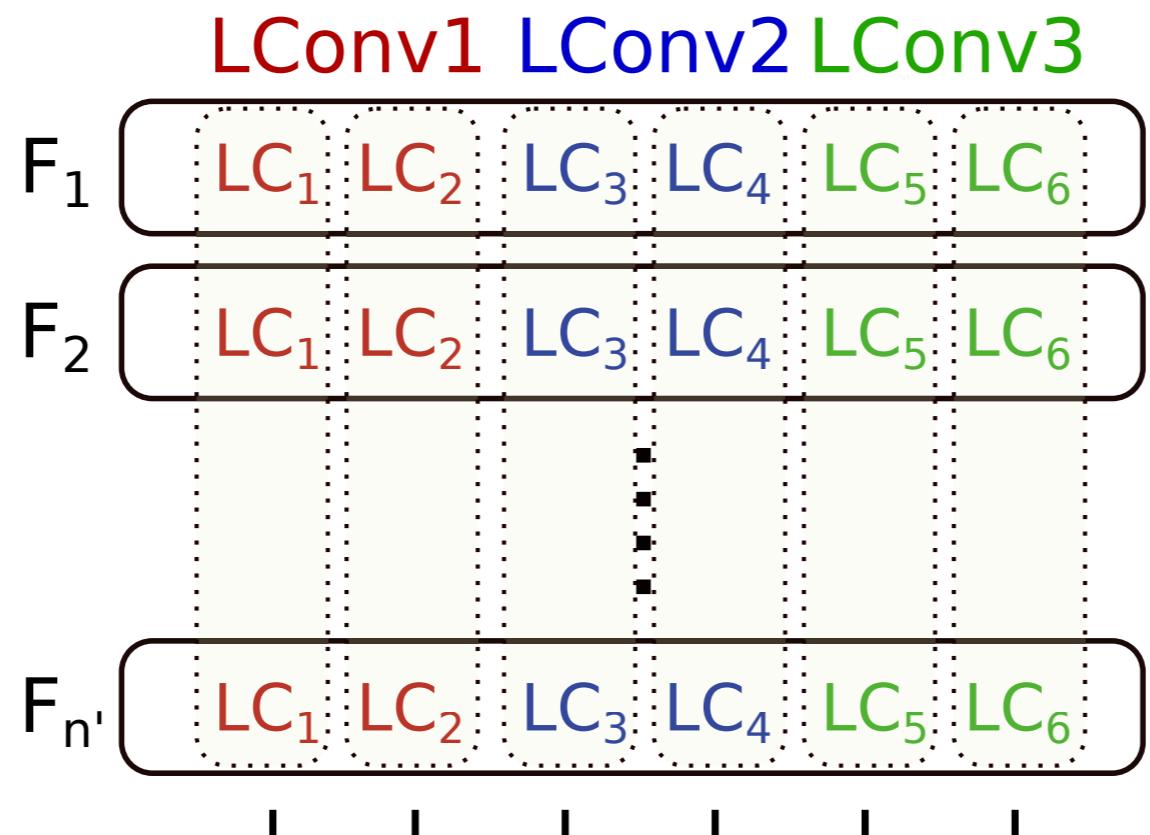
Depth-wise convolution for kernel size $k = 3$.

Other CNN Types

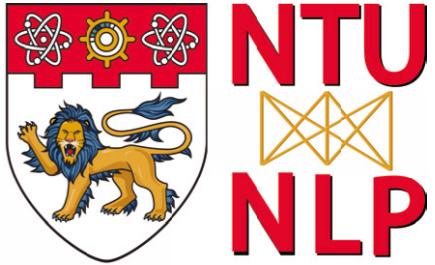


Depth-wise Lightweight Convolution

- The convolutions are done over the input dimensions.
- Consecutive dimensions share parameters

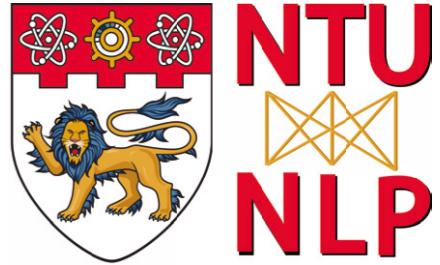


Lecture Plan



- Classification tasks in NLP
- Window-based Approach for Language Modeling
- Window-based Approach for NER, POS tagging, and Chunking
- Convolutional Neural Net for NLP
- Max-margin Training
- Scaling Softmax (for Language Modeling)
- Adaptive input and output.

Pairwise (Contrastive) Training



Training losses covered so far

- Mean-Squared Error (MSE) for regression

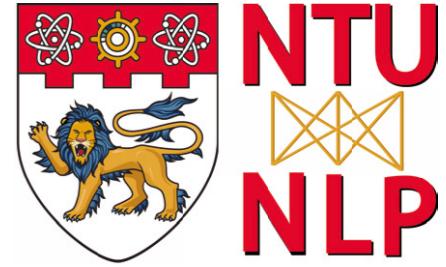
$$\sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

- Cross entropy

binary CE
$$-\sum_{i=1}^N [y_i \log \mu_i + (1 - y_i) \log(1 - \mu_i)]$$

For multi-class
$$-\sum_{i=1}^N \sum_{c=1}^C y_{ic} \log \mu_{ic}$$

Pairwise (Contrastive) Training



- Cross entropy

$$\text{Binary CE} = - \sum_{i=1}^N [y_i \log \mu_i + (1 - y_i) \log(1 - \mu_i)]$$

Sigmoid
 $\text{sigm}(\mathbf{w}^\top \mathbf{x})$

$$\text{For multi-class} = - \sum_{i=1}^N \sum_{c=1}^C y_{ic} \log \mu_{ic}$$

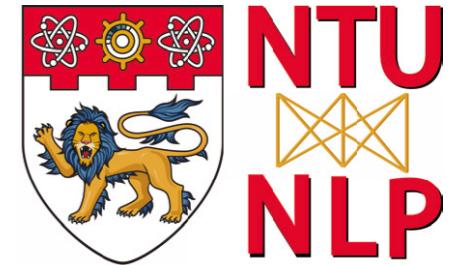
Softmax
 $\mathcal{P}(y = c | \mathbf{x}, \mathbf{W}) = \frac{\exp(\mathbf{w}_c^\top \mathbf{x})}{\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x})}$

$$\text{Negative sampling} = - \log \sigma(u_o^\top v_c) + \sum_{i=1}^k \mathbb{E}_{j \sim P(w)} [\log \sigma(-u_j^\top v_c)]$$

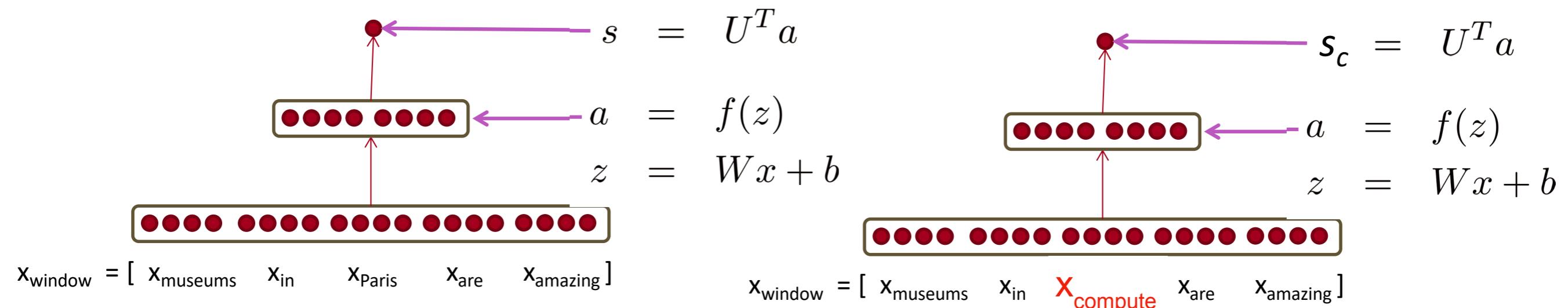
Sigmoid

- Negative sampling is a type of Noise Contrastive Estimation (NCE)
- Pairwise ranking loss is also similar in spirit

Pairwise (Contrastive) Training



- Pairwise ranking loss (Hinge loss, Margin loss, Contrastive loss)



Positive example (middle word
Paris is original)

Negative example (middle word
Paris is a random word)

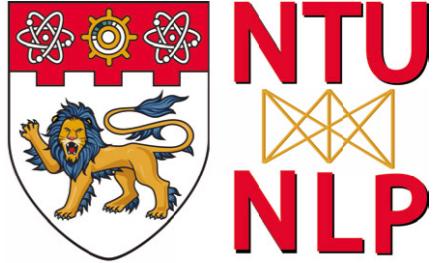
Minimize

$$J = \max(0, 1 - s + s_c)$$

xxx |← 1 →| ooo

See Collobert et al (2011)

Scaling Softmax: Hierarchical Softmax



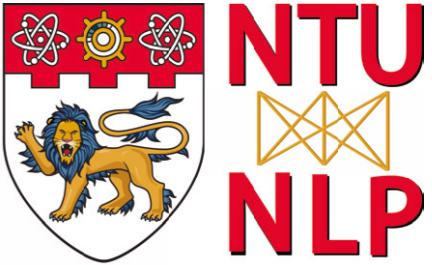
If h is the last hidden layer of size d and V is the output vocabulary of k items, Softmax layer requires $d \times k$ (e.g., 1000×200000) computations for every example during training and testing.

- Assign each word w to a unique class $c(w)$.
- For a given input, first predict the class c (using a softmax)
- You do a softmax over all the words for just that cluster.
- Factorize the conditional probability of word w_t given final hidden state h_t as

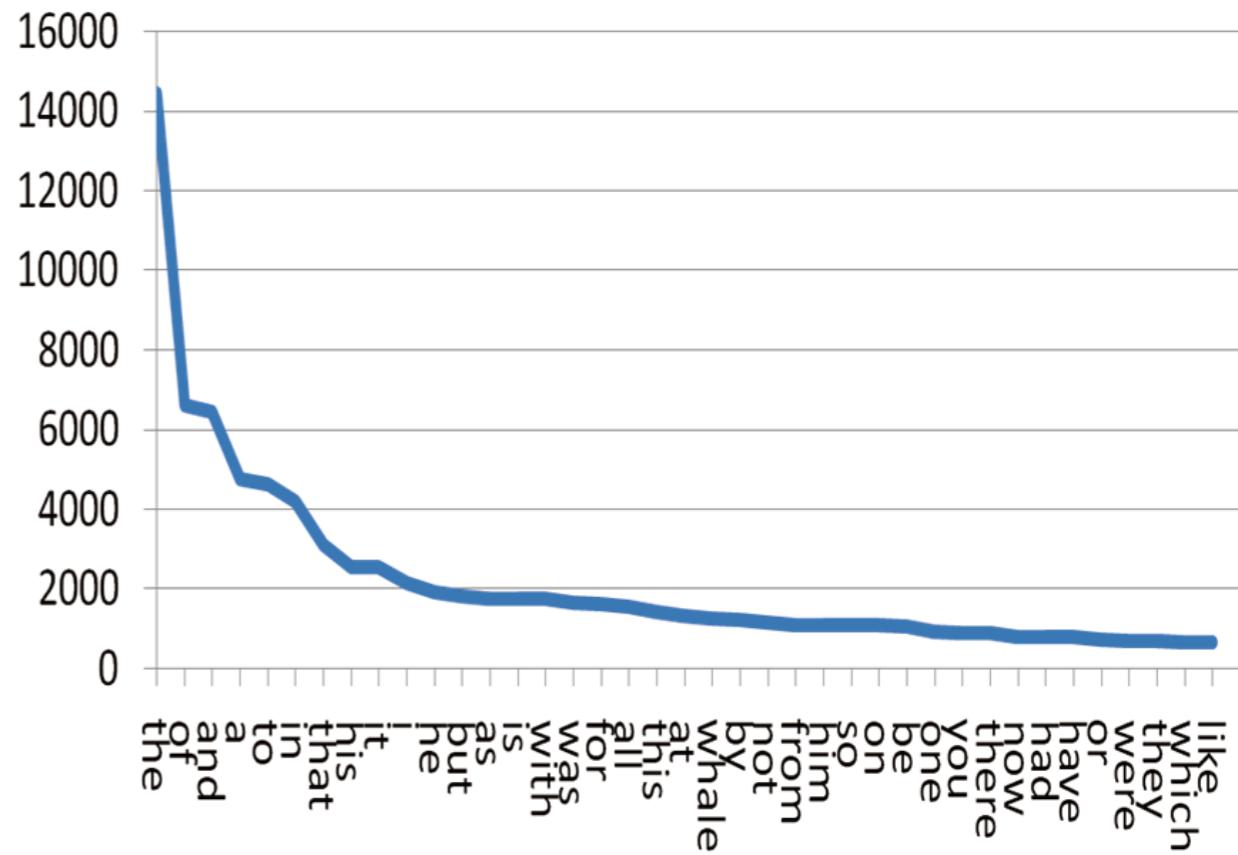
$$p(w_t | h_t) = p_1(c(w_t) | h_t) \times p_2(w_t | c(w_t), h_t)$$

- If each class contains \sqrt{k} words, the computation cost $O(d\sqrt{k})$

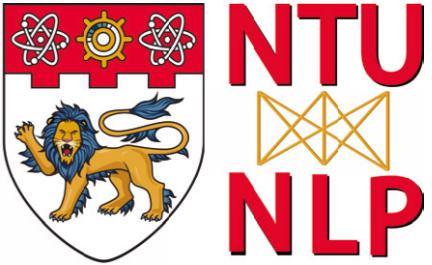
Scaling Softmax: Adaptive Softmax



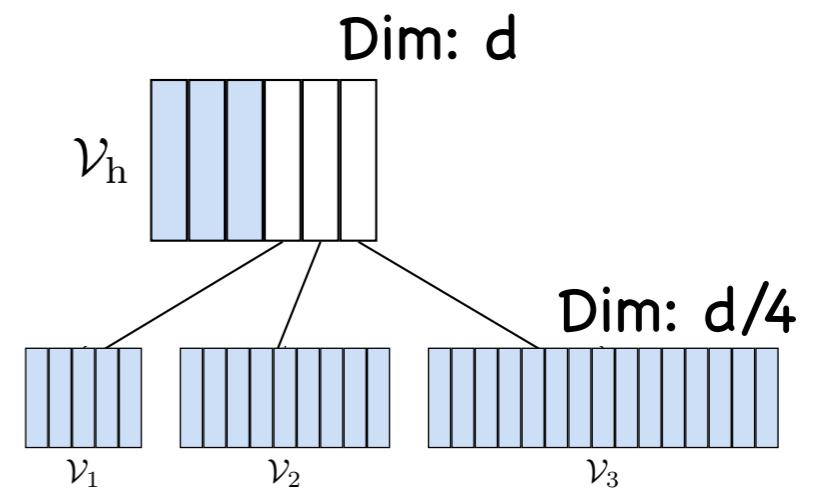
- In natural languages, the distribution of the words follows a Zipf law (Zipf, 1949). Most of the probability mass is covered by a small fraction of the dictionary.
- Partition the dictionary V into two clusters as V_h and V_t , where V_h denotes the head of the distribution consisting of the most frequent words, and V_t is the tail containing a large number of rare words.
- The partitions have the property: $|V_h| \ll |V_t|$ and probabilities $P(V_h) \gg P(V_t)$



Scaling Softmax: Adaptive Softmax

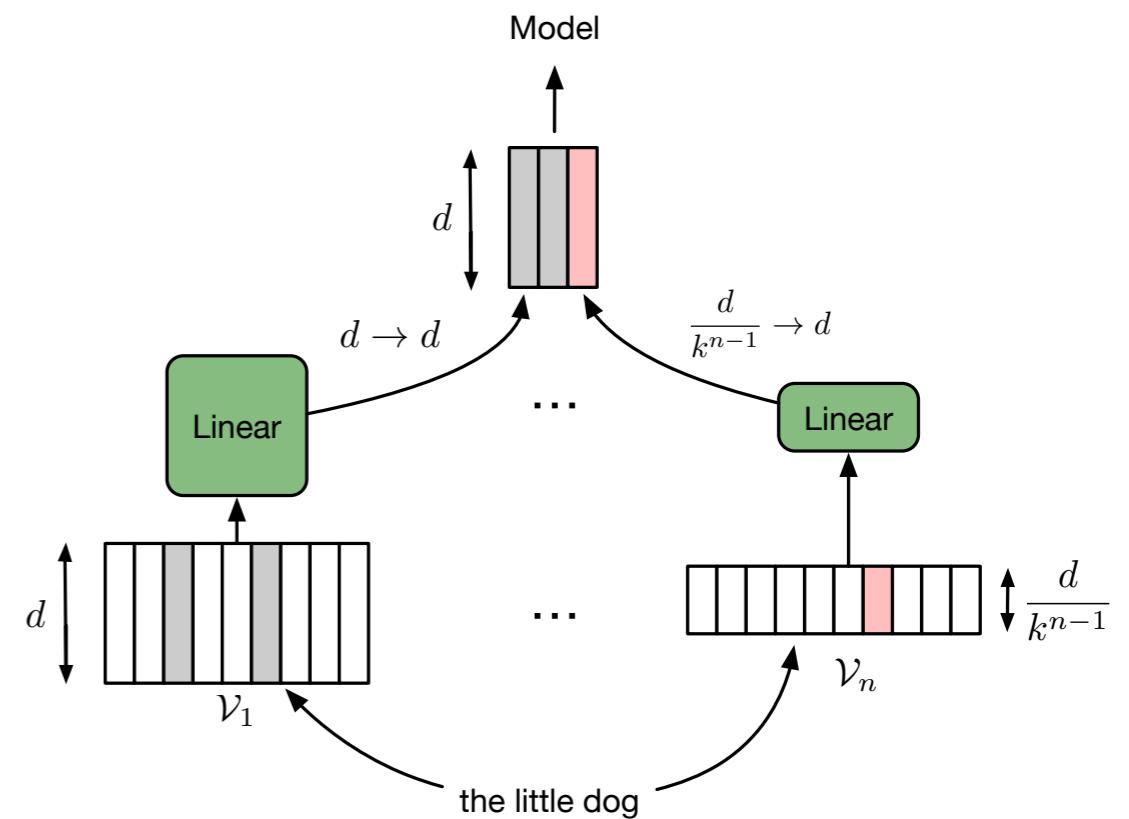


- Put the words in V_h in the root node along with the “other class” ids.
- The tree has two levels.
- If the word w comes from V_h , the probability is simply $p(w|h)$. If it comes from one of the other clusters then the probability is $P(V_i|h) * P(w|V_i, h)$
- Reduce complexity (and parameters) further by having a variable capacity scheme for output word embeddings, assigning more parameters (d for V_h) to frequent words and fewer parameters ($d/4$ for others) to rare words. This is done by projecting h into h_i of dimensions $d/4$ (for other classes)

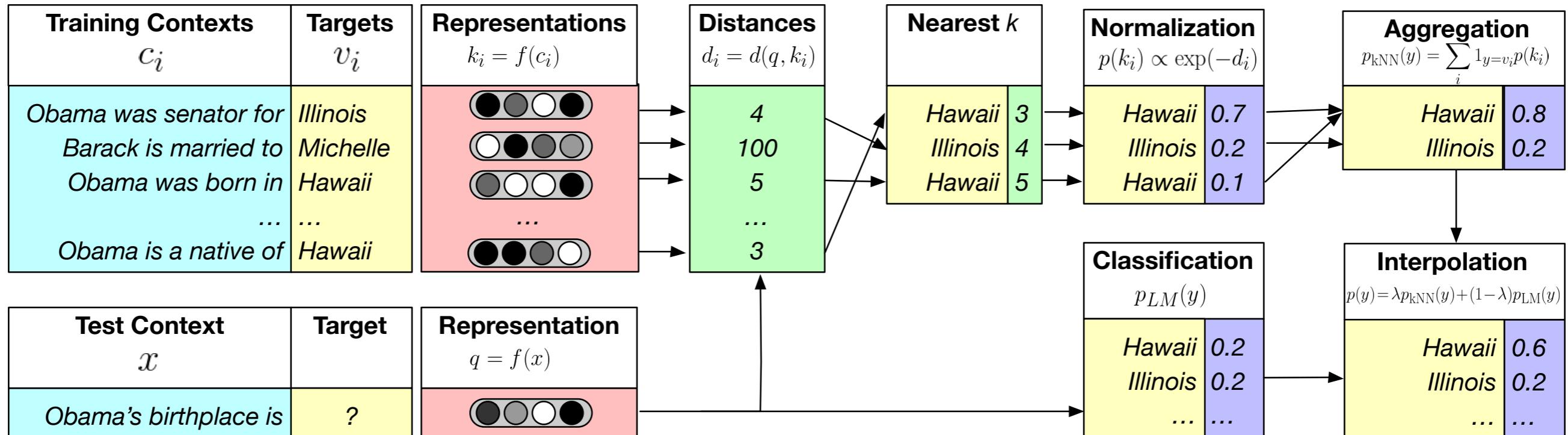
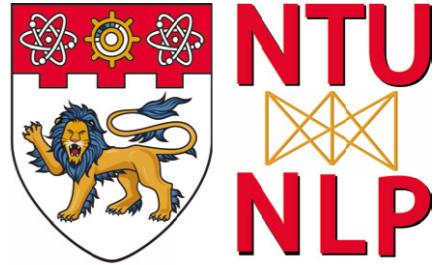


Adaptive Input

- The variable capacity idea can be applied to input word embeddings as well
- Assign more capacity to frequent words and reduce the capacity for less frequent words with the benefit of reducing overfitting to rare words.
- Define a number of clusters that partitions the frequency ordered vocabulary $V = V_1 \cup V_2, \dots, V_{n-1} \cup V_n$ such that $V_i \cap V_j = \emptyset$ for $\forall i, j$, and $i \neq j$, where V_1 contains the most frequent words and V_n the least frequent words.
- Reduce the capacity for each cluster by a factor of k . That is, if words in V_1 have dimension d , then words in V_n have dimension $\frac{d}{k^{n-1}}$

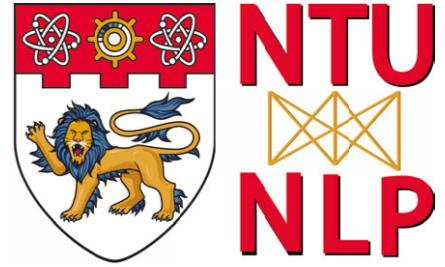


KNN-LM



- Construct a datastore with an entry for each training token, and an encoding of its context.
- For inference, encode a test context, and retrieve the k most similar training contexts from the datastore, along with the corresponding targets.
- A distribution over targets is computed based on the distance of the corresponding context from the test context. This distribution is then interpolated with the original model's output distribution

Language Model Results



	Test	Train Time (hours)	Parameters
Dauphin et al. (2017)	31.9	-	428M
Józefowicz et al. (2016)	30.0	-	1,040M
Shazeer et al. (2017)	28.0	-	4,371M [†]
Char-CNN	25.88	79	366M
Adaptive inputs	25.22	55	331M
Adaptive inputs (large)	23.91	72	465M
Adaptive inputs (very large)	23.02	145	1026M
10 LSTMs + SNM10-SKIP (Shazeer et al., 2016)	23.7	-	-

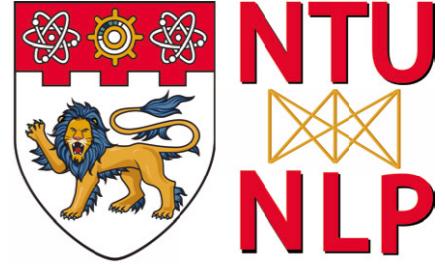
Table 1: Test perplexity on BILLION WORD. Adaptive inputs share parameters with an adaptive softmax. Training times of Char-CNN and Adaptive input models are measured when training with 64 GPUs.

[†]does not include embedding and softmax layers

	Test	Train Time (hours)	Parameters
Grave et al. (2016)	40.8	-	
Dauphin et al. (2017)	37.2	-	229M
Merity et al. (2018)	33.0	-	151M
Rae et al. (2018)	29.2	-	
Adaptive inputs	18.7	67	247M

Table 2: Test perplexity on WIKITEXT-103 (cf. Table 1). Training time is based on 8 GPUs.

Language Model Results



Model	Perplexity (\downarrow)		# Trainable Params
	Dev	Test	
Baevski & Auli (2019)	17.96	18.65	247M
+Transformer-XL (Dai et al., 2019)	-	18.30	257M
+Phrase Induction (Luo et al., 2019)	-	17.40	257M
Base LM (Baevski & Auli, 2019)	17.96	18.65	247M
+ k NN-LM	16.06	16.12	247M
+Continuous Cache (Grave et al., 2017c)	17.67	18.27	247M
+ k NN-LM + Continuous Cache	15.81	15.79	247M

Table 1: Performance on WIKITEXT-103. The k NN-LM substantially outperforms existing work. Gains are additive with the related but orthogonal continuous cache, allowing us to improve the base model by almost 3 perplexity points with no additional training. We report the median of three random seeds.

Model	Perplexity (\downarrow)		# Trainable Params
	Dev	Test	
Base LM (Baevski & Auli, 2019)	14.75	11.89	247M
+ k NN-LM	14.20	10.89	247M

Table 2: Performance on BOOKS, showing that k NN-LM works well in multiple domains.