

344.063/163 KV Special Topic:

Natural Language Processing with Deep Learning

N-gram Representations with Convolutional Neural Networks



Navid Rekab-Saz

navid.rekabsaz@jku.at

Institute of Computational Perception

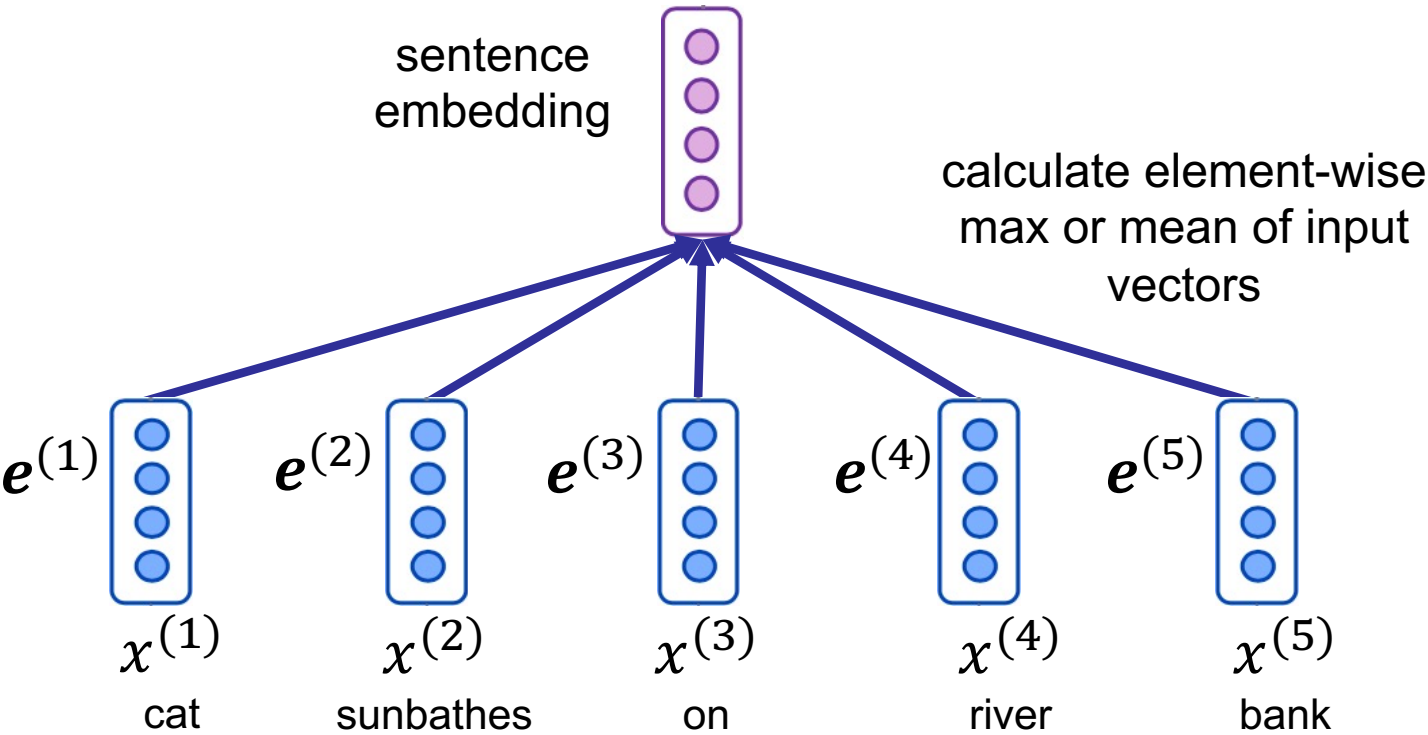
Agenda

- N-Gram Embeddings with CNN
- CNN in practice
 - Document classification
 - From characters to word embedding
 - CNN in information retrieval models

Agenda

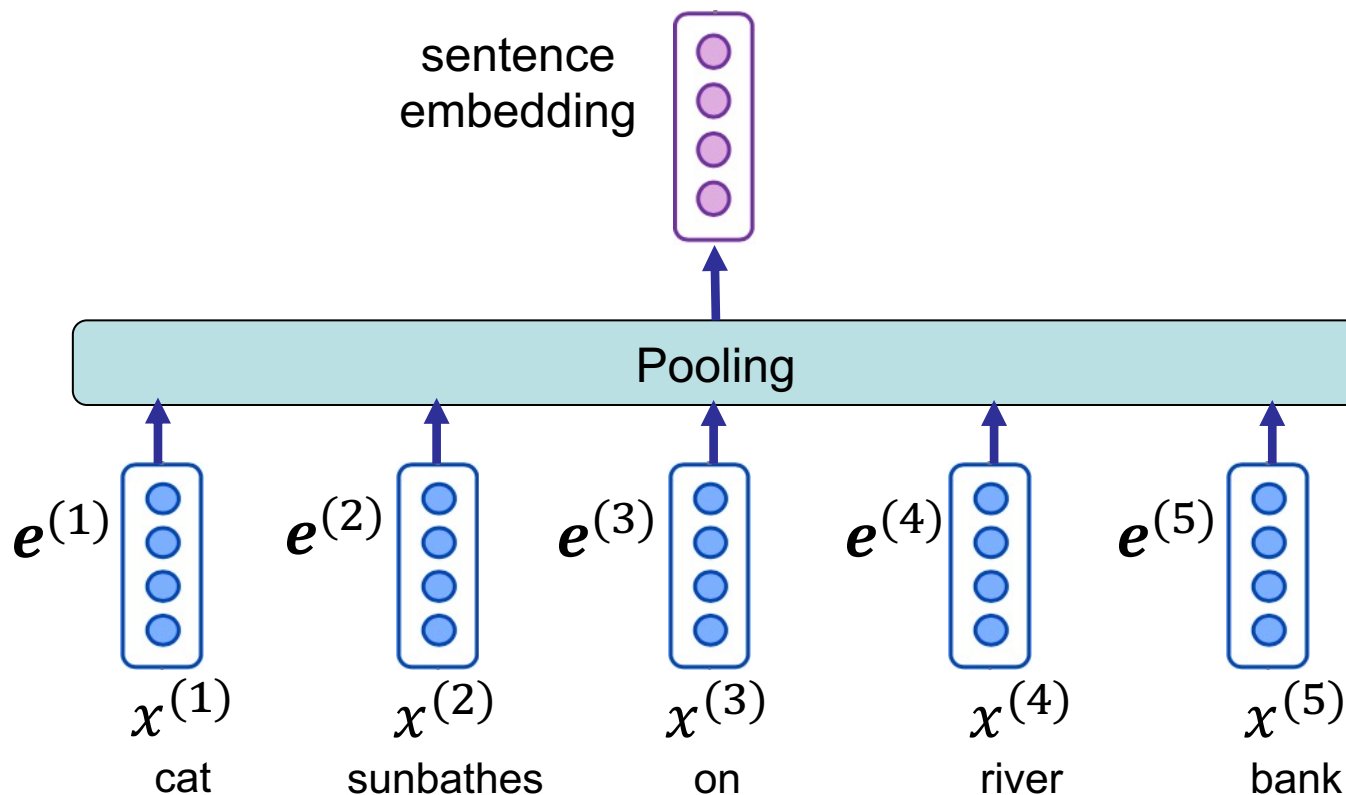
- **N-Gram Embeddings with CNN**
- CNN in practice
 - Document classification
 - From characters to word embedding
 - CNN in information retrieval models

Sentence embedding – from word embeddings



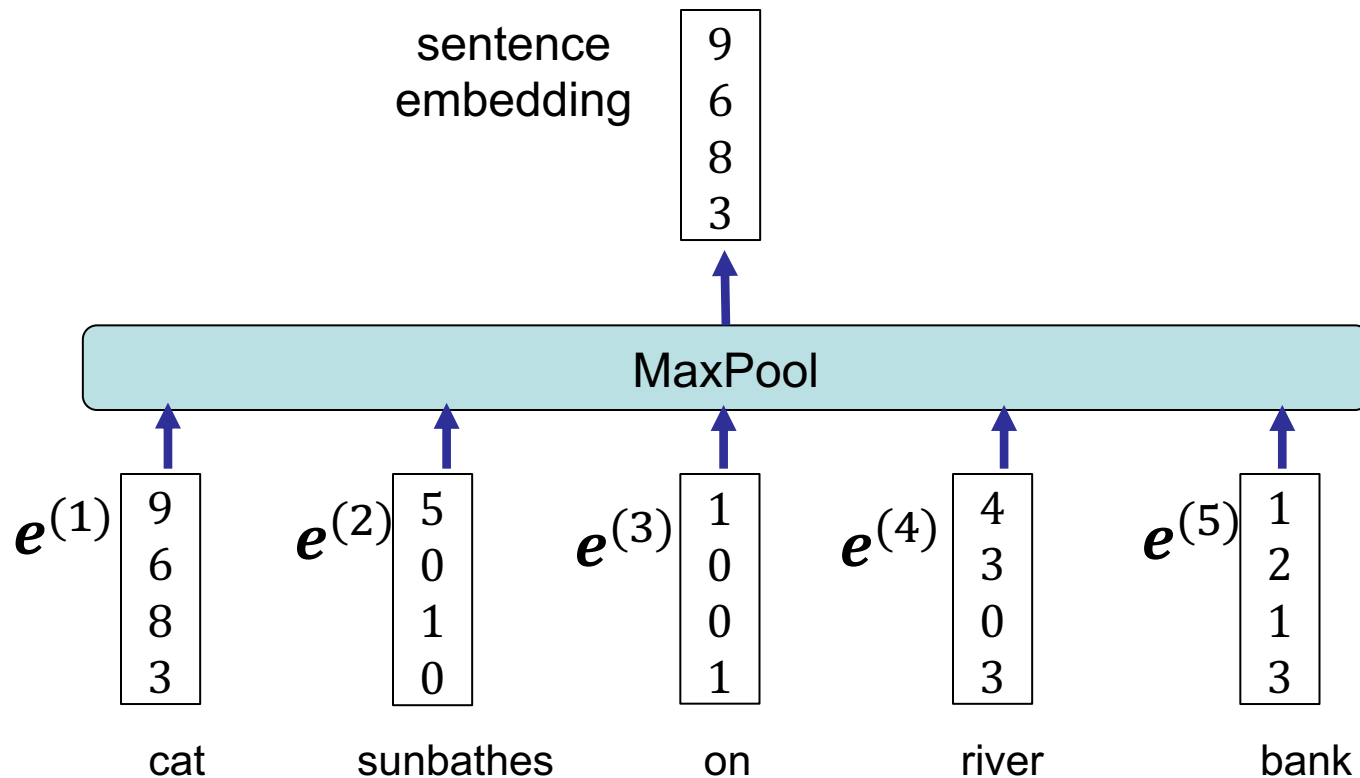
Sentence embedding – from word embeddings

- **Pooling:** element-wise operation on input vectors resulting to an output vector



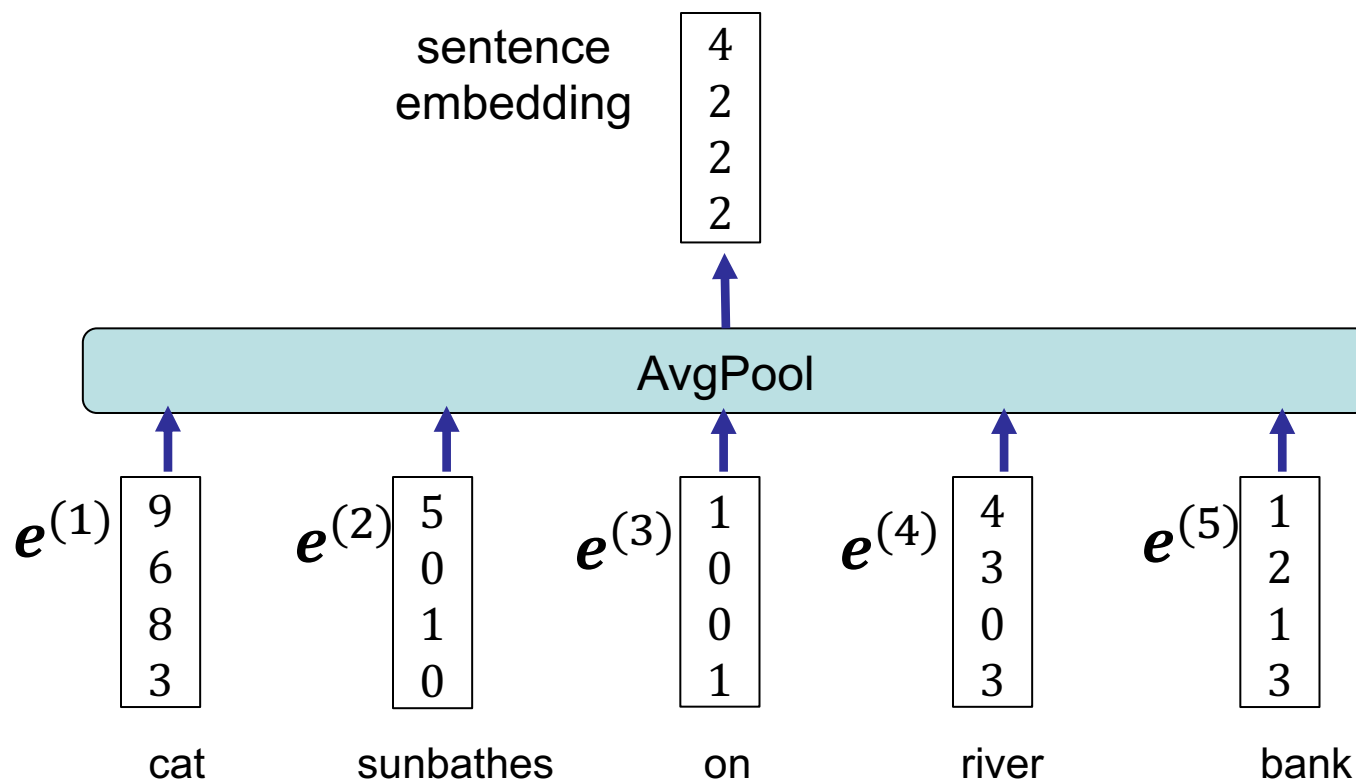
Sentence embedding – from word embeddings

- Pooling: element-wise operation on input vectors resulting to an output vector
- MaxPool: element-wise maximum of inputs

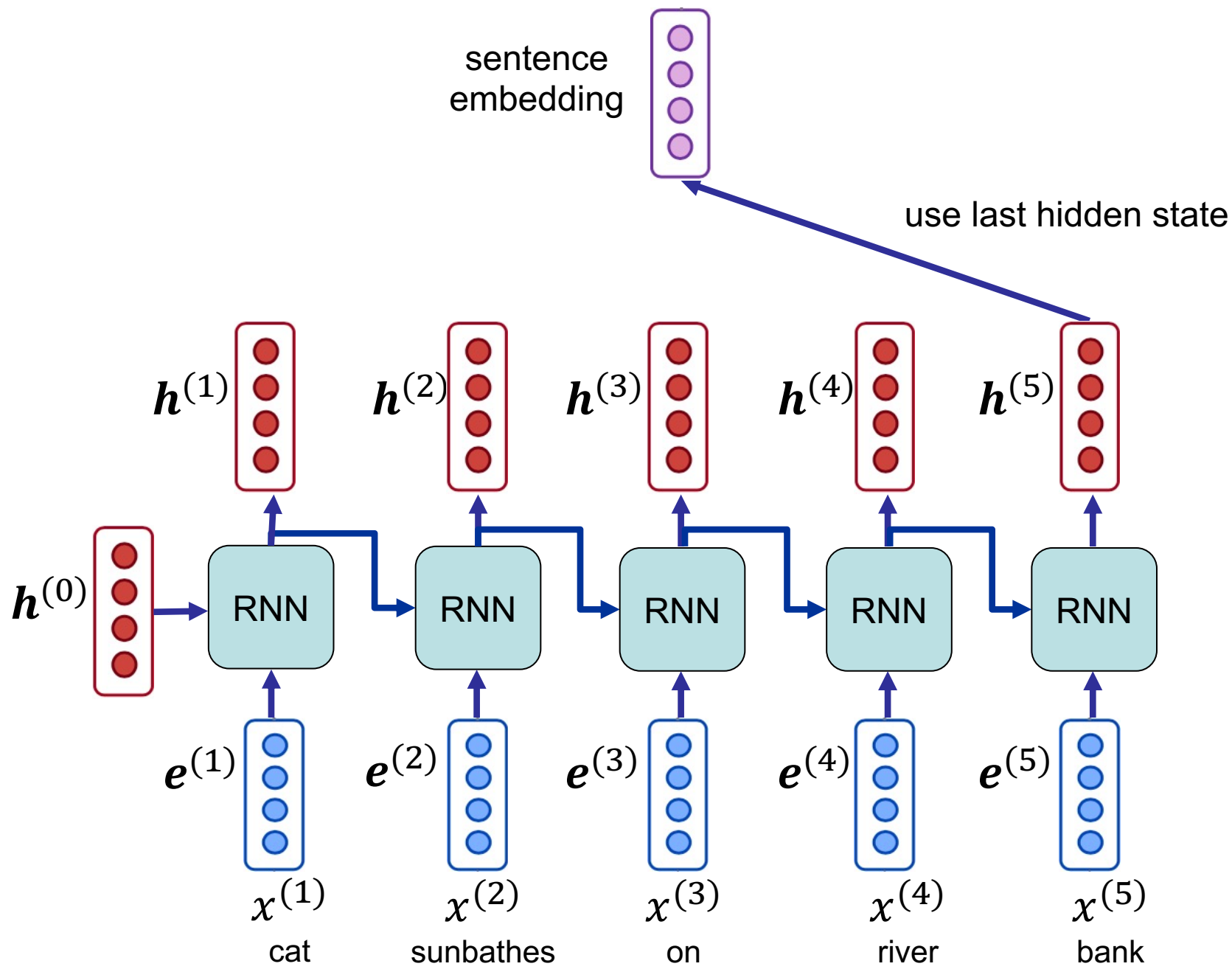


Sentence embedding – from word embeddings

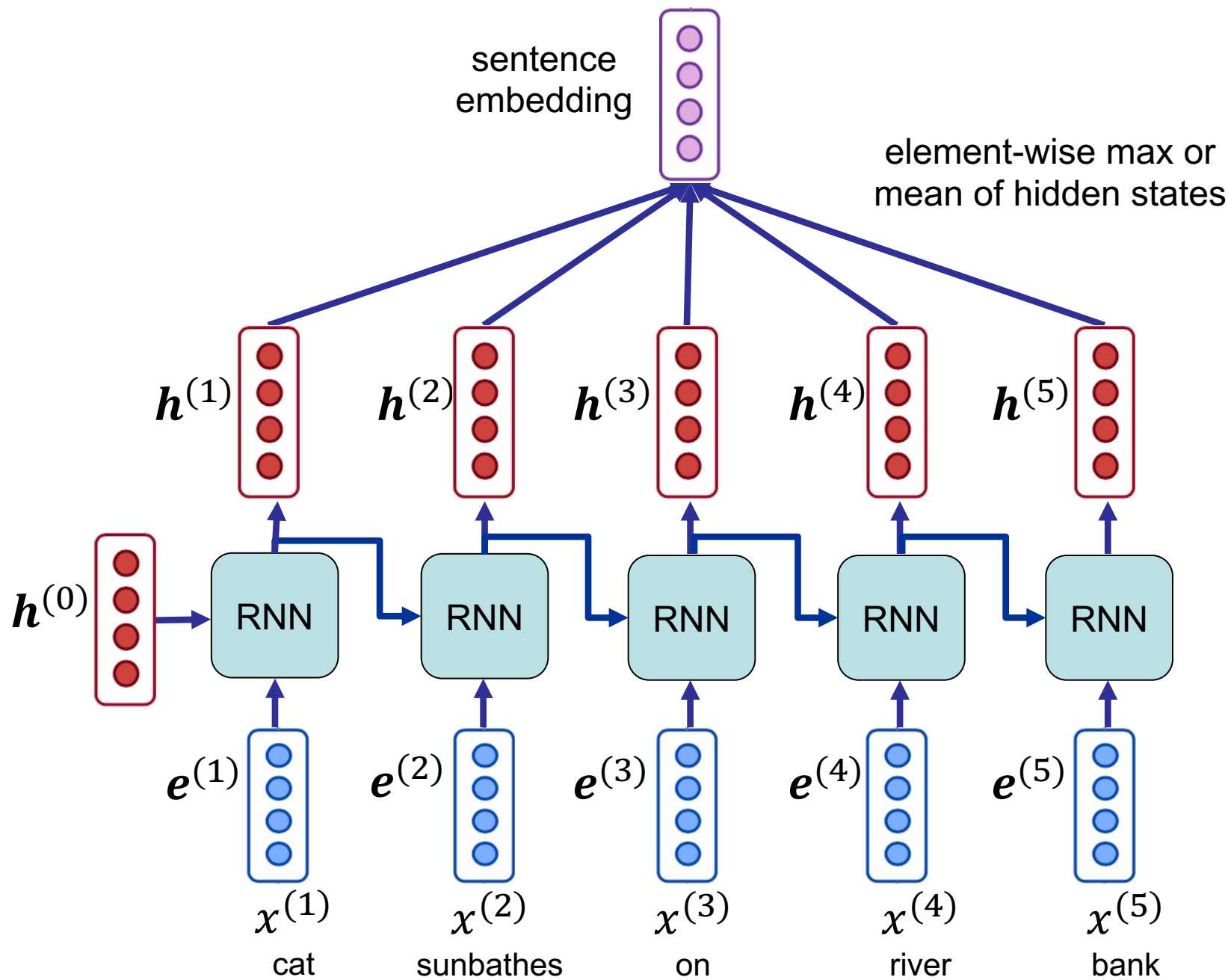
- **Pooling**: element-wise operation on input vectors resulting to an output vector
- **MaxPool**: element-wise maximum of inputs
- **AvgPool**: element-wise average of inputs



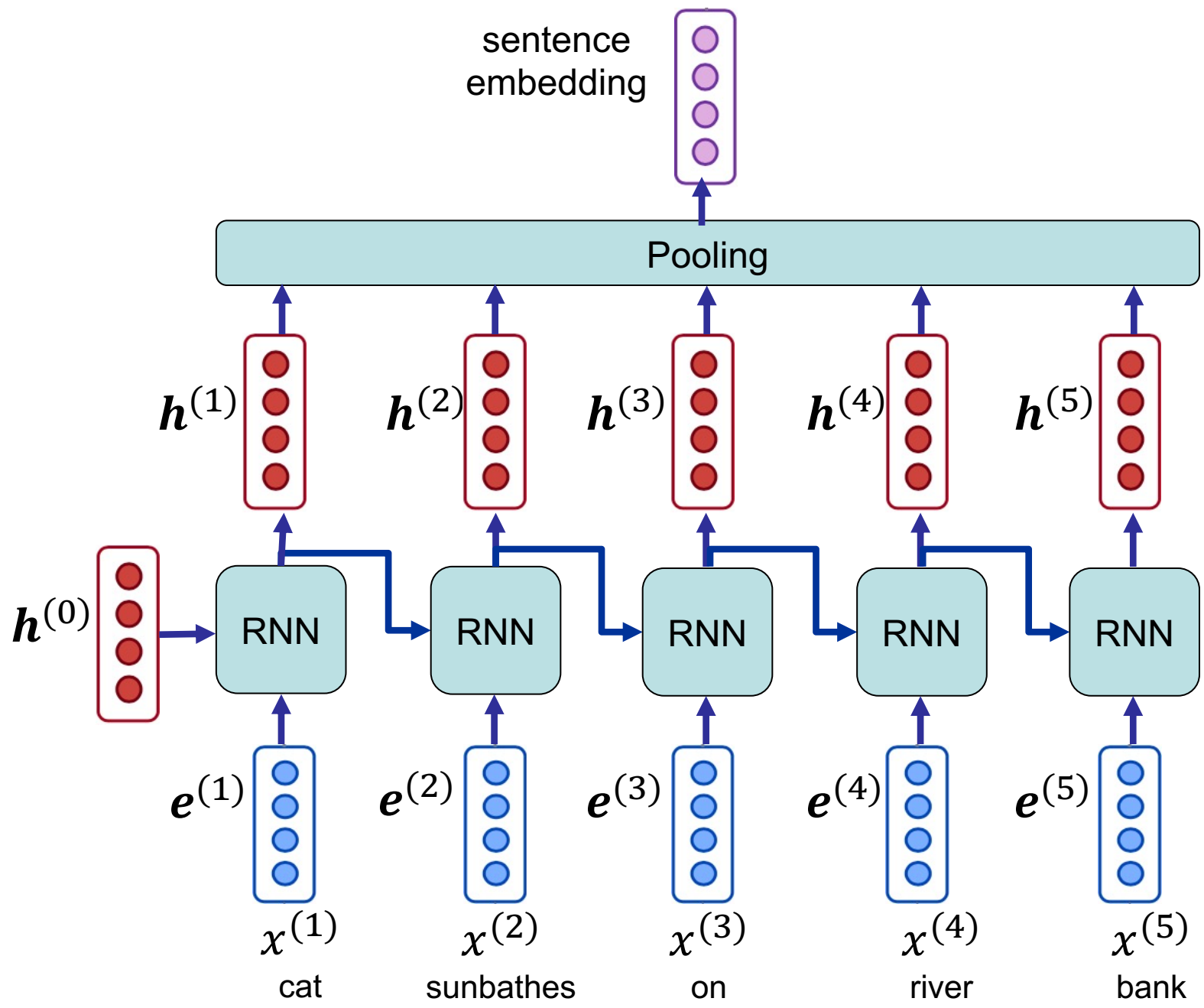
Sentence embedding – RNN



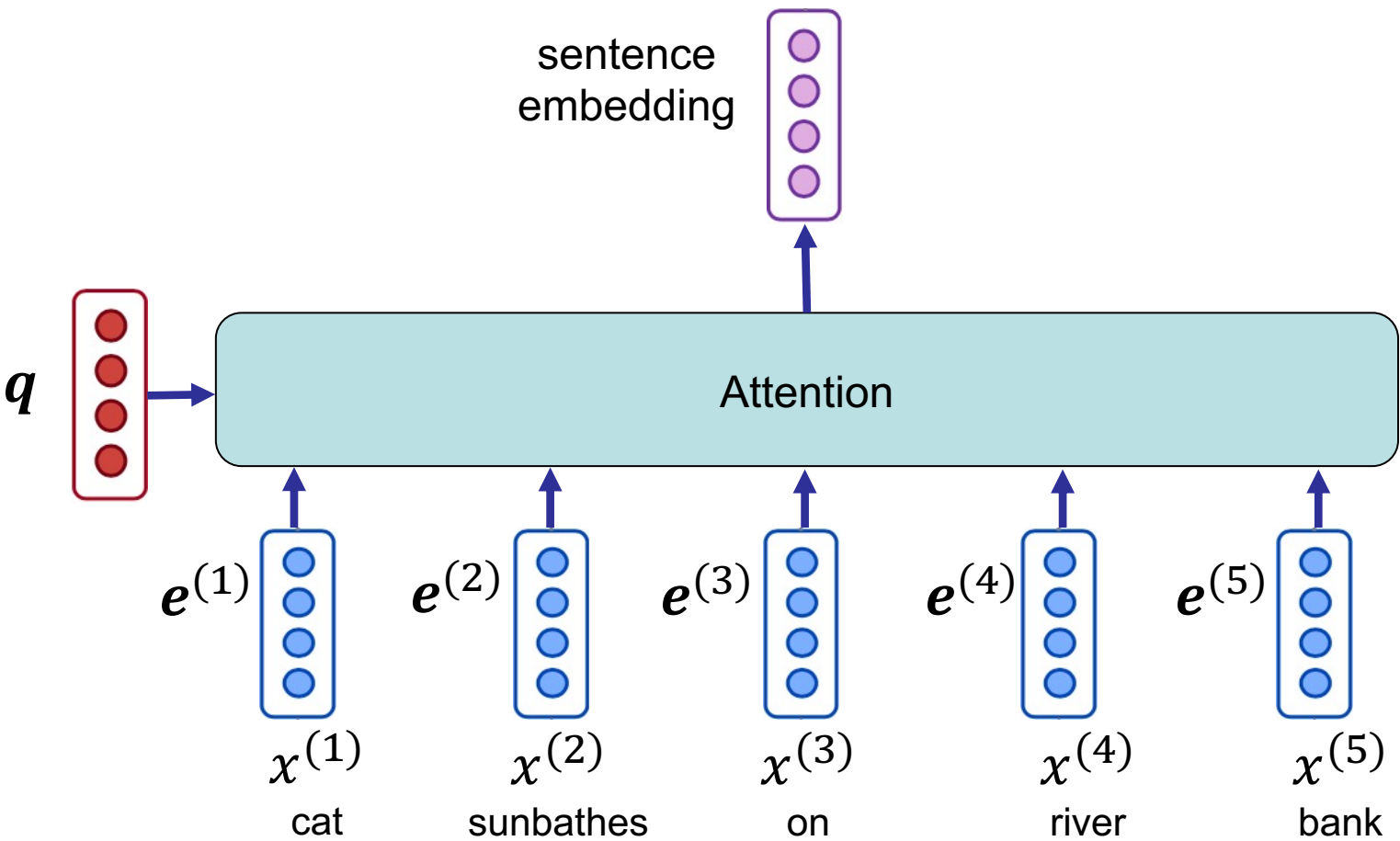
Sentence embedding – RNN



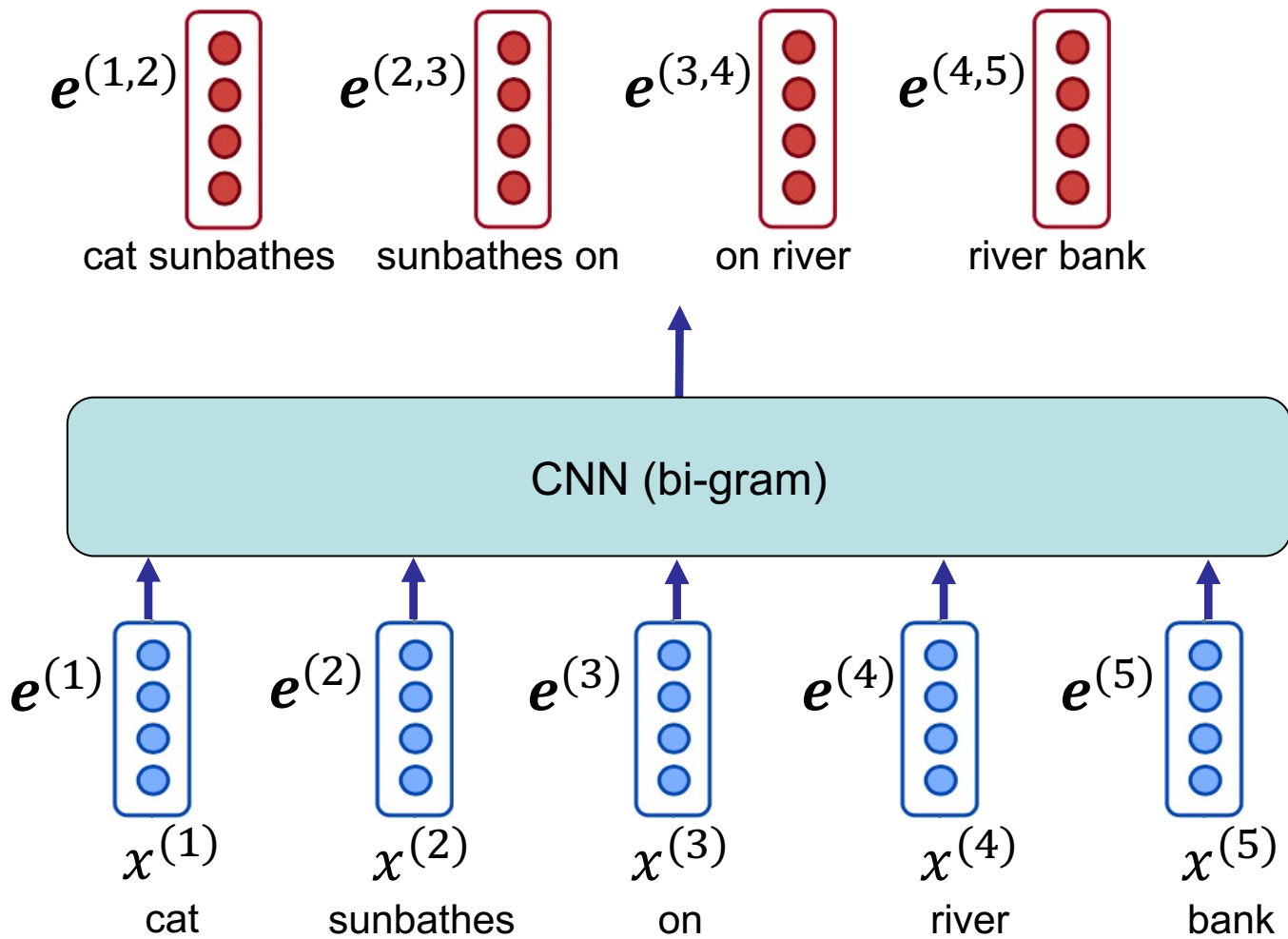
Sentence embedding – RNN



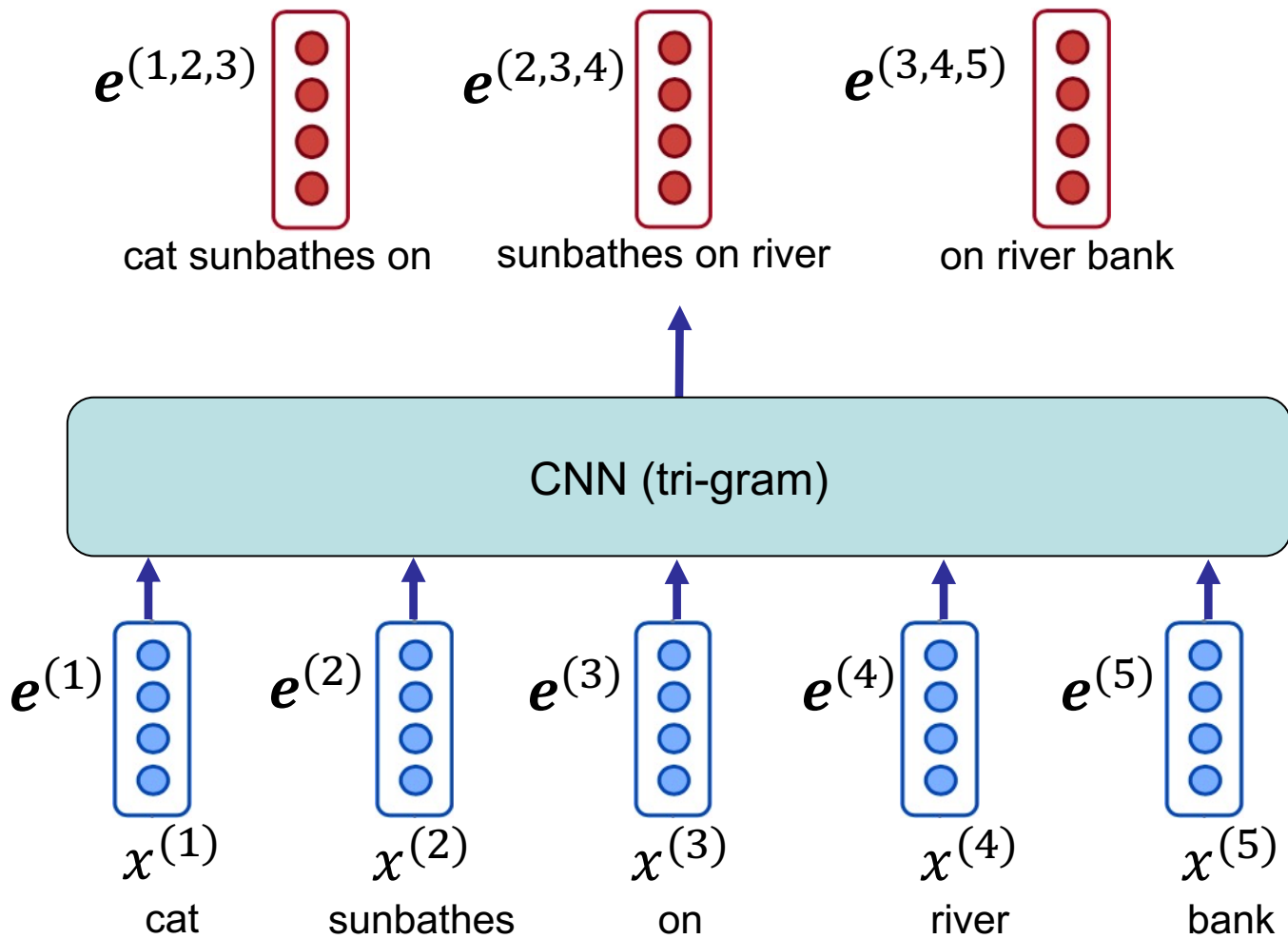
Sentence embedding – attention networks



N-gram embeddings



N-gram embeddings



Convolutional Neural Networks for NLP

- In many NLP models, we can benefit from the vectors which correspond to every sequence of input with a certain length
 - Like bi-gram, tri-gram, 4-gram embeddings

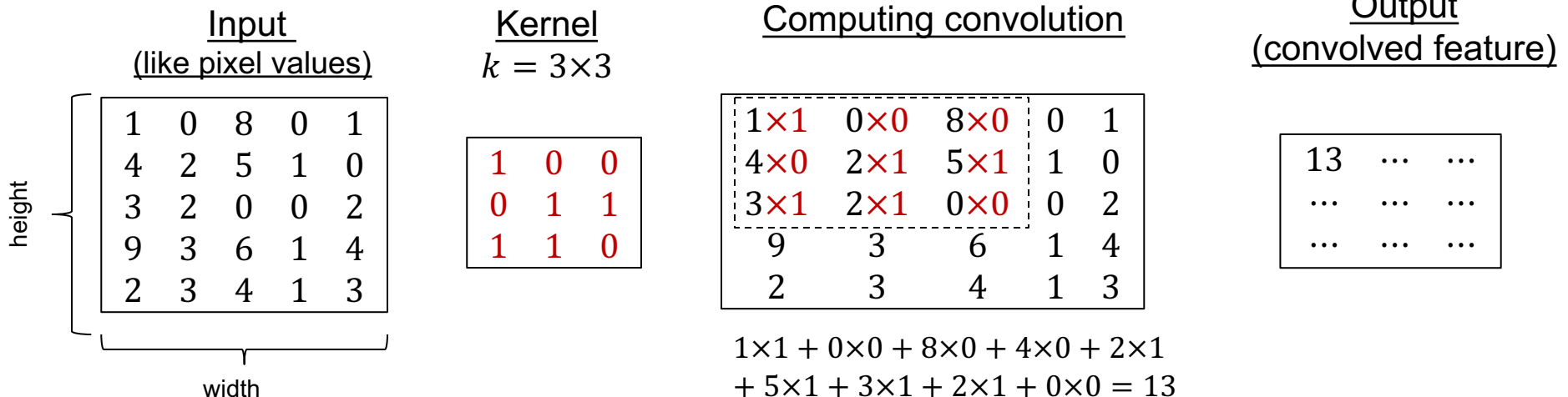
This lecture

- First part: How to create n-gram embeddings using Convolutional Neural Nets (CNNs)
- Second part: How to use these embeddings in different NLP models

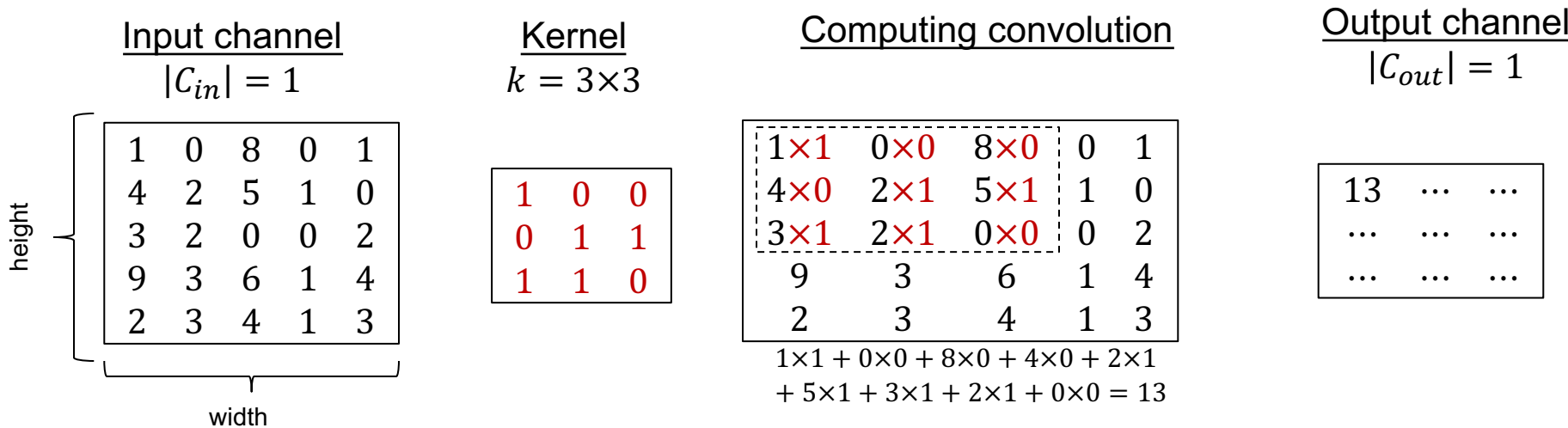
CNNs

- CNNs are widely used to extract features from images
 - CNNs capture **position-invariant** patterns from the input data, where ...
 - the patterns are captured by a set of **kernels**
- Kernel (or filter)
 - A kernel is a set of parameters, ...
 - applied to **every sequence** of input values of a certain length ...
 - to create the output vector in respect to that sequence

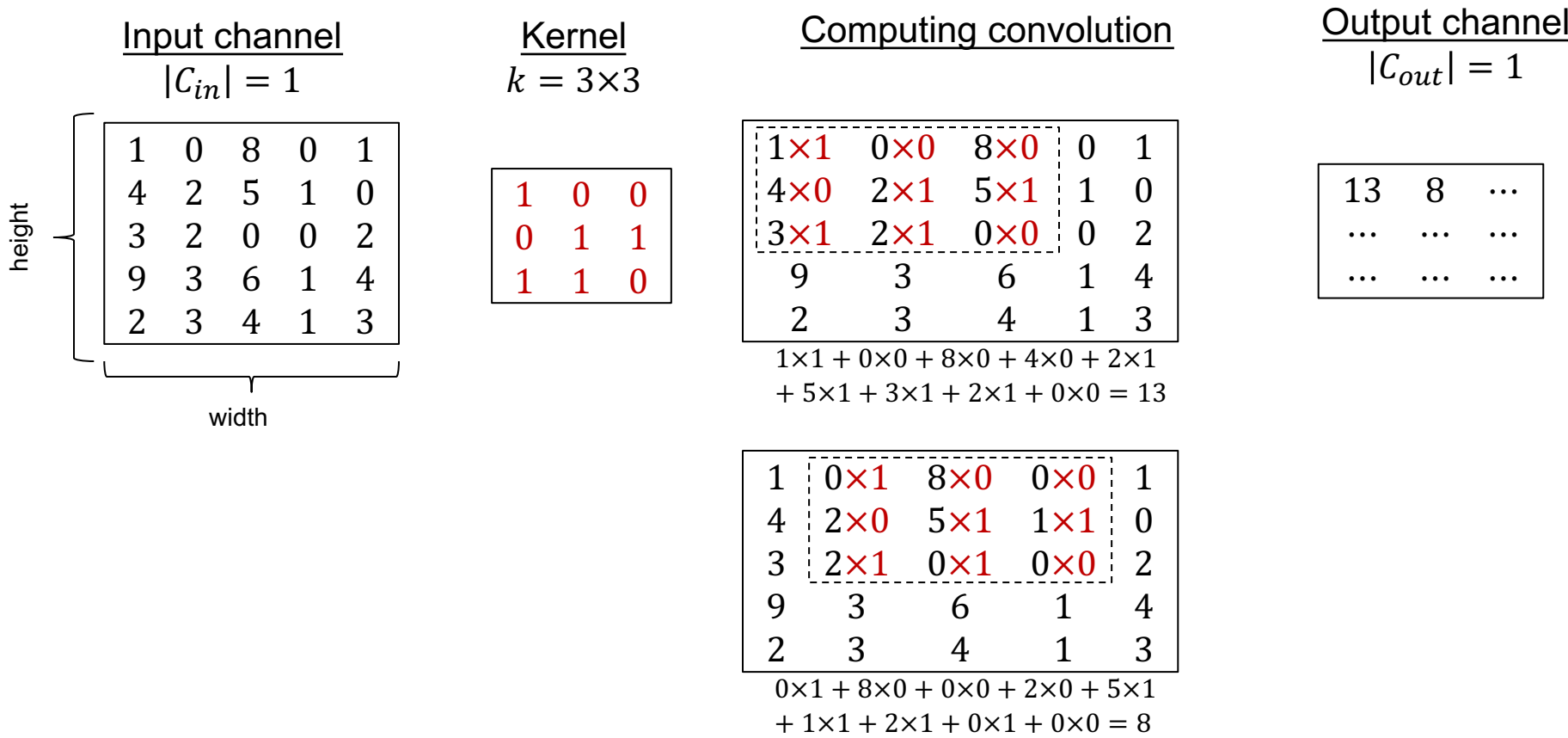
Example: 2d Image data with Conv2d



2-dimensional CNN (CONV2D)



2-dimensional CNN (CONV2D)



2-dimensional CNN (CONV2D)

Input channel
 $|C_{in}| = 1$

height

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

width

Kernel
 $k = 3 \times 3$

1	0	0
0	1	1
1	1	0

Computing convolution

1×1	0×0	8×0	0	1
4×0	2×1	5×1	1	0
3×1	2×1	0×0	0	2
9	3	6	1	4
2	3	4	1	3

$1 \times 1 + 0 \times 0 + 8 \times 0 + 4 \times 0 + 2 \times 1 + 5 \times 1 + 3 \times 1 + 2 \times 1 + 0 \times 0 = 13$

1	0×1	8×0	0×0	1
4	2×0	5×1	1×1	0
3	2×1	0×1	0×0	2
9	3	6	1	4
2	3	4	1	3

$0 \times 1 + 8 \times 0 + 0 \times 0 + 2 \times 0 + 5 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 1 + 0 \times 0 = 8$

1	0	8	0	1
4	2×1	5×0	1×0	0
3	2×0	0×1	0×1	2
9	3×1	6×1	1×0	4
2	3	4	1	3

$2 \times 1 + 5 \times 0 + 1 \times 0 + 2 \times 0 + 0 \times 1 + 0 \times 1 + 3 \times 1 + 6 \times 1 + 1 \times 0 = 11$

Output channel
 $|C_{out}| = 1$

13	8	...
...	11	...
...

Calculate other values!

2-dimensional CNN (CONV2D)

Input channels (like RGB)
 $|C_{in}| = 3$

$C_{in}^{(1)}$

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

$C_{in}^{(2)}$

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

$C_{in}^{(3)}$

3	1	0	0	6
4	2	2	0	7
2	1	0	0	1
6	2	0	2	2
4	1	0	3	6

Kernel
 $k = 3 \times 3$

1	0	0
0	1	1
1	1	0

0	0	0
0	0	0
1	0	0

0	1	1
1	0	1
1	1	0

Computing convolution

1×1	0×0	8×0	0	1
4×0	2×1	5×1	1	0
3×1	2×1	0×0	0	2
9	3	6	1	4
2	3	4	1	3

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

3×0	1×1	0×1	0	6
4×1	2×0	2×1	0	7
2×1	1×1	0×0	0	1
6	2	0	2	2
4	1	0	3	6

Output channel
 $|C_{out}| = 1$

$C_{out}^{(1)}$

28
...
...

$(1 \times 1 + 0 \times 0 + 8 \times 0 + 4 \times 0 + 2 \times 1 + 5 \times 1 + 3 \times 1 + 2 \times 1 + 0 \times 0)$
 $+ (1 \times 0 + 7 \times 0 + 4 \times 0 + 3 \times 0 + 1 \times 0 + 3 \times 0 + 5 \times 1 + 0 \times 0 + 9 \times 0)$
 $+ (3 \times 0 + 1 \times 1 + 0 \times 1 + 4 \times 1 + 2 \times 0 + 2 \times 1 + 2 \times 1 + 1 \times 1 + 0 \times 0)$
 $= 28$

Parameters are shown in red

19

1-dimensional CNN (CONV1D)

Input channels
 $|C_{in}| = 4$

Kernel
 $k = 3$

Computing convolution

Output channel
 $|C_{out}| = 1$

$w_i^{(j)}$
kernel weights for
 j th input channel and
 i th output channel

$C_{in}^{(4)}$	1	0	8	0	1	$w_1^{(4)}$	1	0	0
$C_{in}^{(3)}$	1	7	4	6	0	$w_1^{(3)}$	0	0	0
$C_{in}^{(2)}$	3	1	0	0	6	$w_1^{(2)}$	0	1	1
$C_{in}^{(1)}$	4	2	1	0	2	$w_1^{(1)}$	1	0	1
input sequence length = 5						kernel size = 3			

$C_{out}^{(1)}$
output sequence length = 3

1-dimensional CNN (CONV1D)

Input channels

$$|C_{in}| = 4$$

Kernel

$$k = 3$$

Computing convolution

1×1	0×0	8×0	0	1
1×0	7×0	4×0	6	0
3×0	1×1	0×1	0	6
4×1	2×0	1×1	0	2

Output channel

$$|C_{out}| = 1$$

$w_i^{(j)}$
kernel weights for
 j th input channel and
 i th output channel

$$(1 \times 1 + 0 \times 0 + 8 \times 0) + (1 \times 0 + 7 \times 0 + 4 \times 0) + (3 \times 0 + 1 \times 1 + 0 \times 1) + (4 \times 1 + 2 \times 0 + 1 \times 1) = 7$$

$C_{in}^{(4)}$	1	0	8	0	1	$w_1^{(4)}$	1	0	0
$C_{in}^{(3)}$	1	7	4	6	0	$w_1^{(3)}$	0	0	0
$C_{in}^{(2)}$	3	1	0	0	6	$w_1^{(2)}$	0	1	1
$C_{in}^{(1)}$	4	2	1	0	2	$w_1^{(1)}$	1	0	1
input sequence length = 5						kernel size = 3			

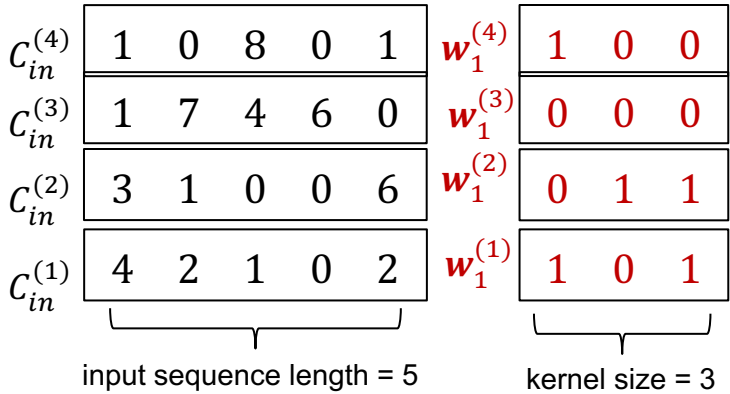
$C_{out}^{(1)}$ 7
 └──────────┘
 output sequence length = 3

1-dimensional CNN (CONV1D)

Input channels
 $|C_{in}| = 4$

Kernel
 $k = 3$

$w_i^{(j)}$
kernel weights for
 j th input channel and
 i th output channel



Computing convolution

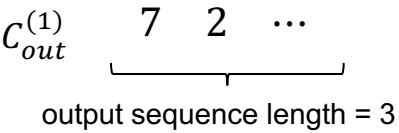
1×1	0×0	8×0	0	1
1×0	7×0	4×0	6	0
3×0	1×1	0×1	0	6
4×1	2×0	1×1	0	2

$(1 \times 1 + 0 \times 0 + 8 \times 0) + (1 \times 0 + 7 \times 0 + 4 \times 0) + (3 \times 0 + 1 \times 1 + 0 \times 1) + (4 \times 1 + 2 \times 0 + 1 \times 1) = 7$

1	0×1	8×0	0×0	1
1	7×0	4×0	6×0	0
3	1×0	0×1	0×1	6
4	2×1	1×0	0×1	2

$(0 \times 1 + 8 \times 0 + 0 \times 0) + (7 \times 0 + 4 \times 0 + 6 \times 0) + (1 \times 0 + 0 \times 1 + 0 \times 1) + (2 \times 1 + 1 \times 0 + 0 \times 1) = 2$

Output channel
 $|C_{out}| = 1$



1-dimensional CNN (CONV1D)

Input channels
 $|C_{in}| = 4$

Kernel
 $k = 3$

Computing convolution

Output channel
 $|C_{out}| = 1$

$w_i^{(j)}$
kernel weights for
 j th input channel and
 i th output channel

1×1	0×0	8×0	0	1
1×0	7×0	4×0	6	0
3×0	1×1	0×1	0	6
4×1	2×0	1×1	0	2

$(1 \times 1 + 0 \times 0 + 8 \times 0) + (1 \times 0 + 7 \times 0 + 4 \times 0) + (3 \times 0 + 1 \times 1 + 0 \times 1) + (4 \times 1 + 2 \times 0 + 1 \times 1) = 7$

$C_{in}^{(4)}$	1	0	8	0	1	$w_1^{(4)}$	1	0	0
$C_{in}^{(3)}$	1	7	4	6	0	$w_1^{(3)}$	0	0	0
$C_{in}^{(2)}$	3	1	0	0	6	$w_1^{(2)}$	0	1	1
$C_{in}^{(1)}$	4	2	1	0	2	$w_1^{(1)}$	1	0	1
input sequence length = 5						kernel size = 3			

1	0×1	8×0	0×0	1
1	7×0	4×0	6×0	0
3	1×0	0×1	0×1	6
4	2×1	1×0	0×1	2

$(0 \times 1 + 8 \times 0 + 0 \times 0) + (7 \times 0 + 4 \times 0 + 6 \times 0) + (1 \times 0 + 0 \times 1 + 0 \times 1) + (2 \times 1 + 1 \times 0 + 0 \times 1) = 2$

1	0	8×1	0×0	1×0
1	7	4×0	6×0	0×0
3	1	0×0	0×1	6×1
4	2	1×1	0×0	2×1

$(8 \times 1 + 0 \times 0 + 1 \times 0) + (4 \times 0 + 6 \times 0 + 0 \times 0) + (0 \times 0 + 0 \times 1 + 6 \times 1) + (1 \times 1 + 0 \times 0 + 2 \times 1) = 17$

$C_{out}^{(1)}$ 7 2 17
 └────────┘
 output sequence length = 3

1-dimensional CNN applied to word embeddings

Input channels

$|C_{in}| = 4$

Number of input channels $|C_{in}|$
=
dimension of word embedding.

Conv1d sees every dimension
as a channel

Kernel

$k = 3$

W_i
kernel weights for
 i th output channel

W_1

1	0	0
0	0	0
0	1	1
1	0	1

Computing convolution

1×1	0×0	8×0	0	1
1×0	7×0	4×0	6	0
3×0	1×1	0×1	0	6
4×1	2×0	1×1	0	2

Output channel

$|C_{out}| = 1$

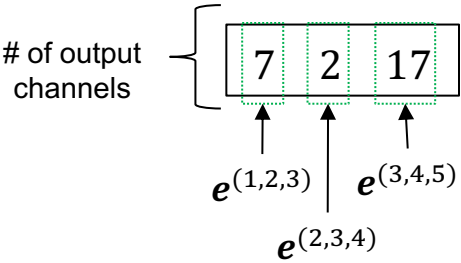
Number of output channels $|C_{out}|$
=
dimension of n-gram embeddings

$C_{in}^{(4)}$	1	0	8	0	1
$C_{in}^{(3)}$	1	7	4	6	0
$C_{in}^{(2)}$	3	1	0	0	6
$C_{in}^{(1)}$	4	2	1	0	2

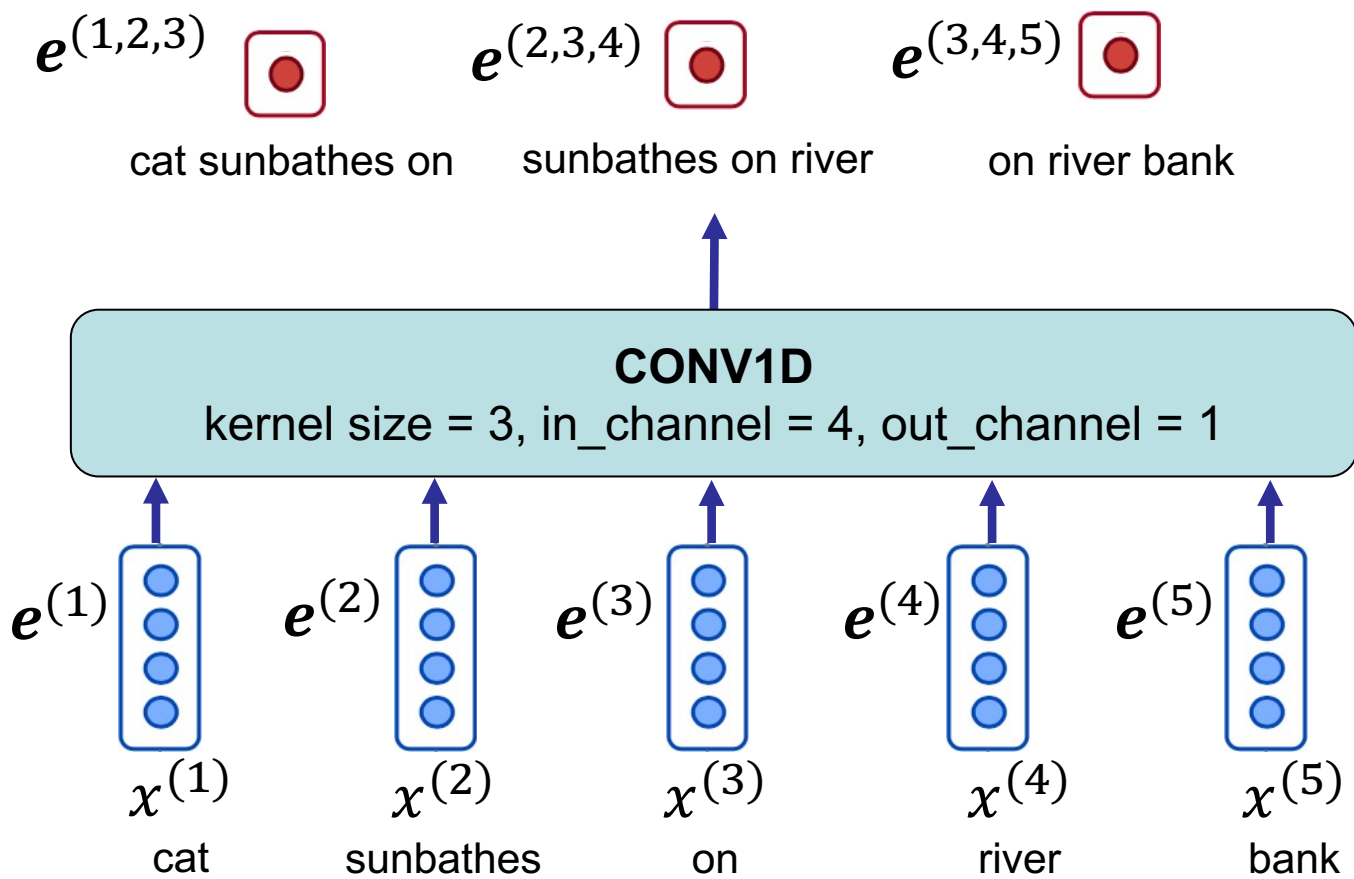
$e^{(1)} e^{(2)} e^{(3)} e^{(4)} e^{(5)}$
 $x^{(1)} x^{(2)} x^{(3)} x^{(4)} x^{(5)}$

1	0×1	8×0	0×0	1
1	7×0	4×0	6×0	0
3	1×0	0×1	0×1	6
4	2×1	1×0	0×1	2

1	0	8×1	0×0	1×0
1	7	4×0	6×0	0×0
3	1	0×0	0×1	6×1
4	2	1×1	0×0	2×1



N-gram embeddings



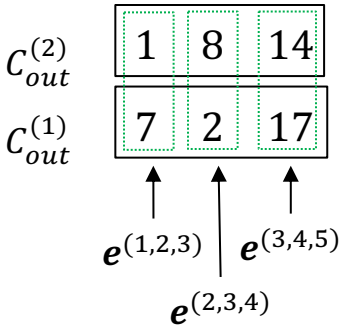
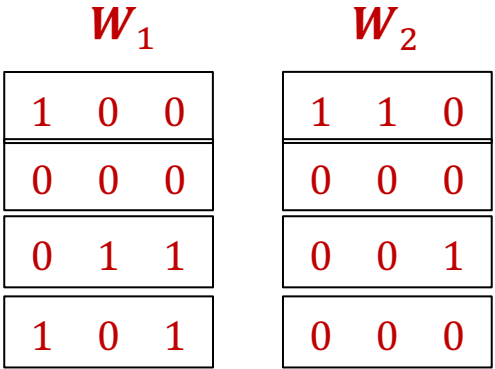
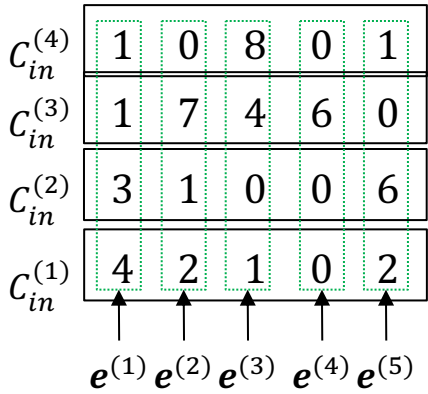
1-dimensional CNN in NLP

Input channels
 $|C_{in}| = 4$

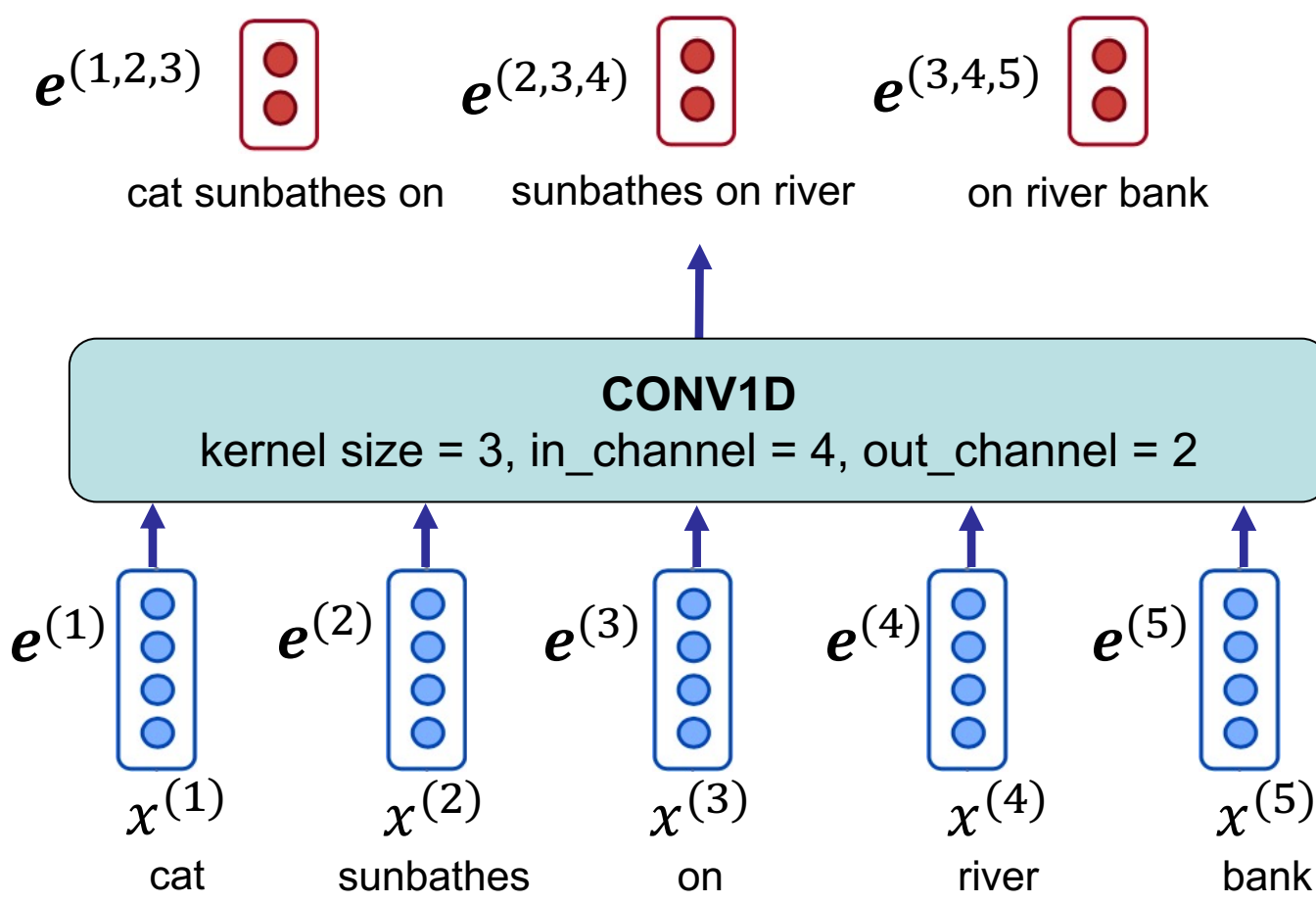
Kernels
 $k = 3$

Output channels
 $|C_{out}| = 2$

W_i : kernel weights for i th output channel



N-gram embeddings



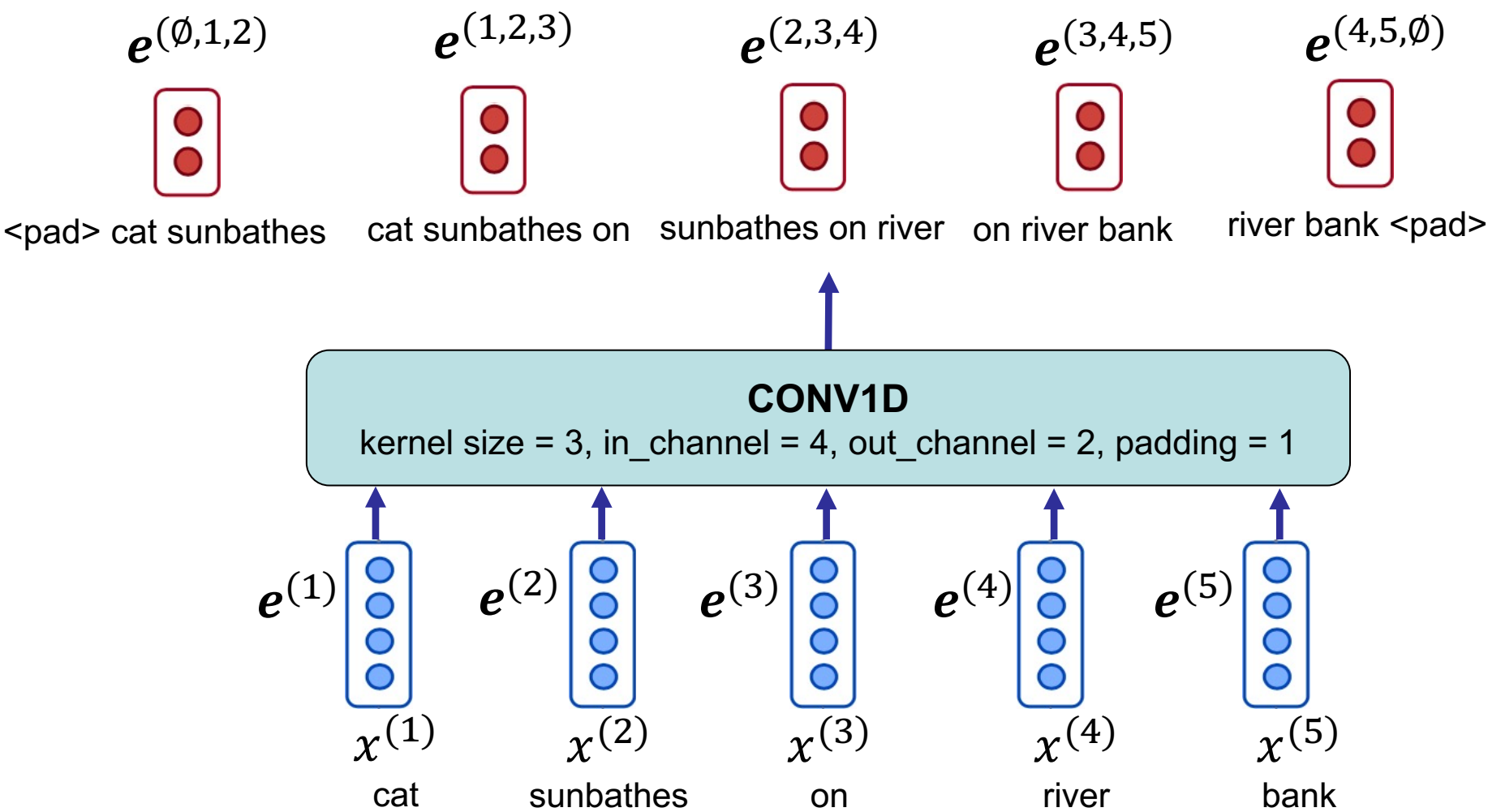
Other notions

- Padding:
 - adds zero vectors to the beginning and end of the sequence
- Stride:
 - The length of the steps over the sequence on which the convolutions are applied
 - Default is 1

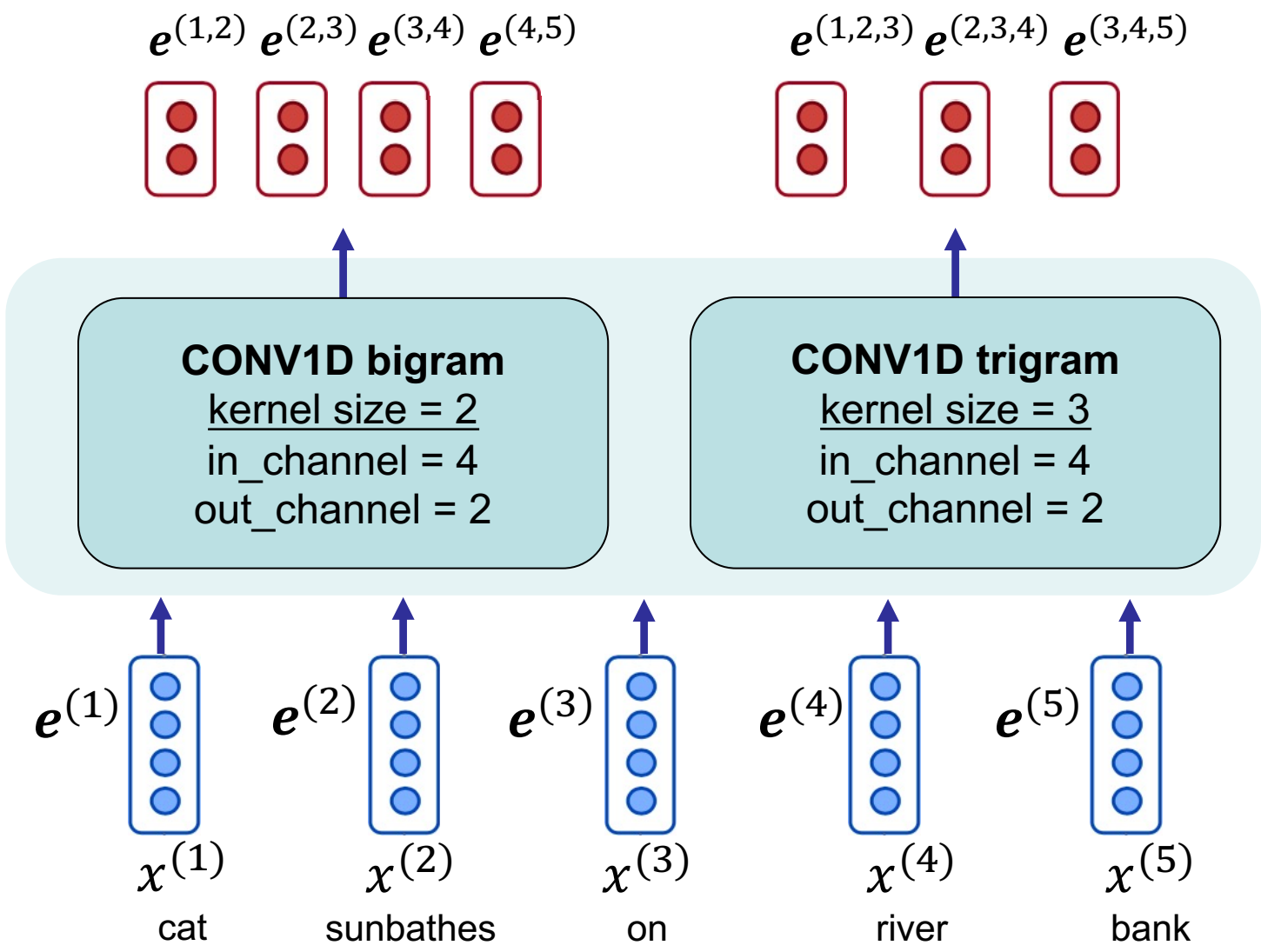
More notions with graphic:

https://github.com/vdumoulin/conv_arithmetic/blob/master/README.md

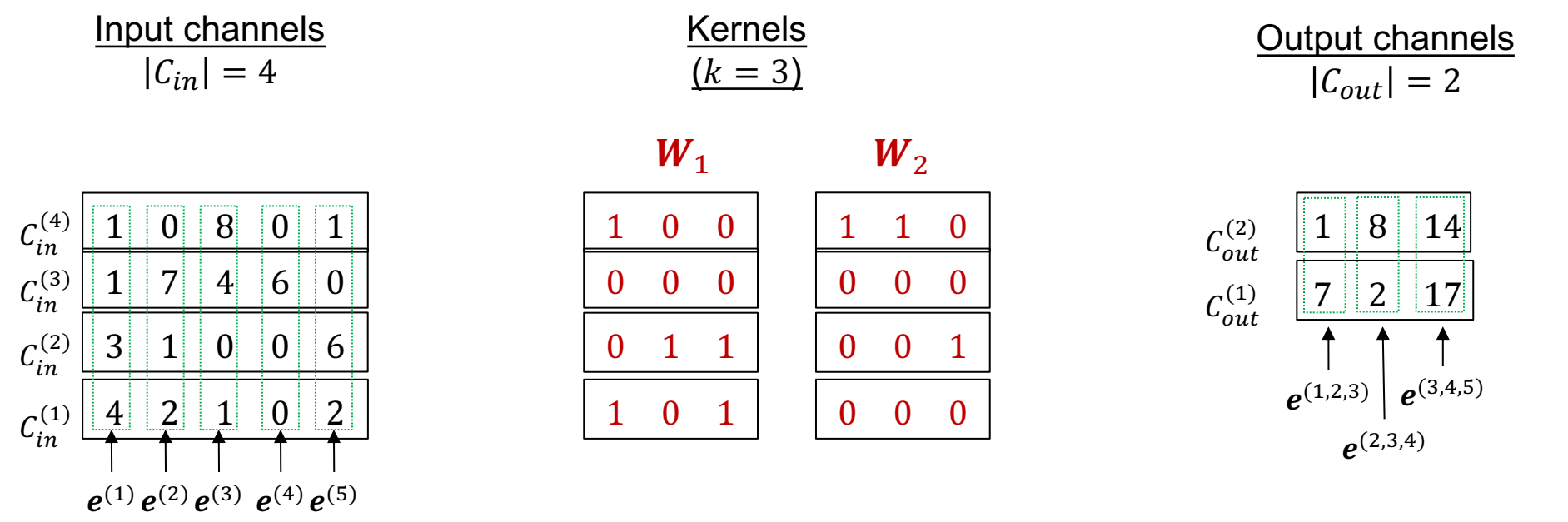
N-gram embeddings



N-gram embeddings



1-dimensional CNN in NLP



Informal formulation of the calculation in Conv1D:

$e_i^{(x,...,x+k)}$

$=$

$\text{torch.sum}([e^{(x)} ; \dots ; e^{(x+k)}] \odot W_i) + b_i$

Position i th of the output embedding corresponding to inputs x till $x + k$

Input embedding x

Element-wise multiplication

Bias term of the i th output channel

CNN – summary

- A model to capture patterns in local proximities, learnt through many (linear) kernels
 - Output embeddings are position-invariant
- NLP mostly uses Conv1D
 - `in_channels` is the dimension of input embeddings
 - `out_channels` is the dimension of output embeddings
 - `kernel_size` is the length of n-gram

CONV1D

```
CLASS torch.nn.Conv1d(in_channels, out_channels, kernel_size, stride=1, padding=0,  
dilation=1, groups=1, bias=True, padding_mode='zeros')
```

[SOURCE]

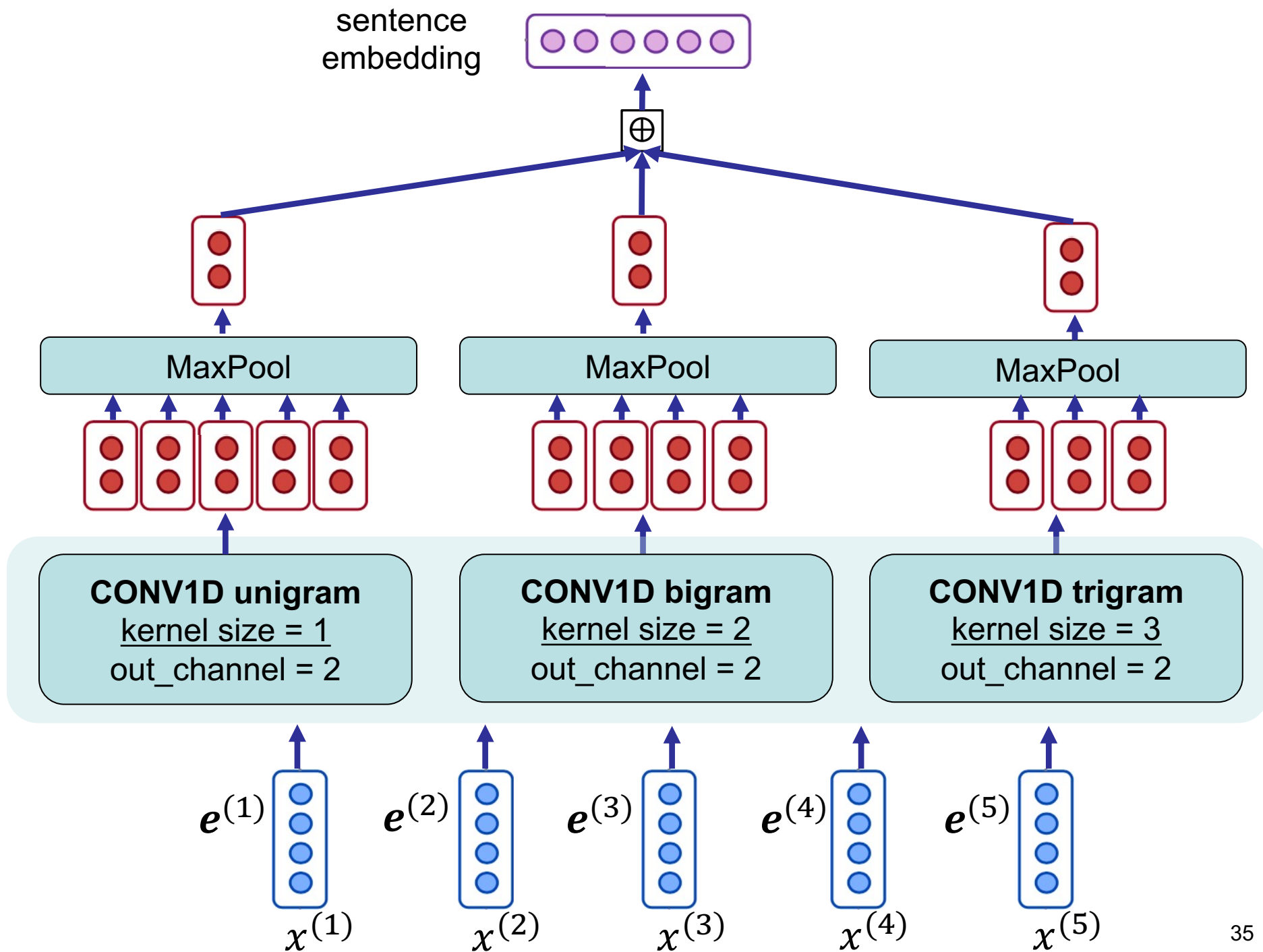
Agenda

- N-Gram Embeddings with CNN
- **CNN in practice**
 - Document classification
 - From characters to word embedding
 - CNN in information retrieval models

Document classification

Document classification with CNNs

1. Create unigram, bigram, trigram, etc. embeddings
2. Apply pooling to merge embeddings of each n-gram over whole the sequence, resulting in several n-gram features
3. Concatenate n-gram features as the final document feature (document embedding)



Document classification

- Unigram embeddings ($k = 1$)? ... can't we just use the original word embeddings?
 - Unigram CNN adds an extra neural network layer with very few additional parameters
 - CNN with $k = 1$ applies the same parameters to all word embeddings (position invariant)
 - Unlike fully connected a feed forward layer which is position variant and adds a lot more parameters

Composing word embeddings from character embeddings

- Instead of predefined word vectors (**static** word embeddings), **compose** the embedding of a word from the embeddings of its characters
 1. Define one vector for every character
 - The embedding matrix will be much smaller in comparison with the ones of word embeddings
 2. Use CNNs to create a word embedding from its character embeddings
 - In the same way that we created a document embedding from word embeddings
 - Each CNN results in a character N-gram embedding

Word embeddings from character embeddings

Task: Language modeling

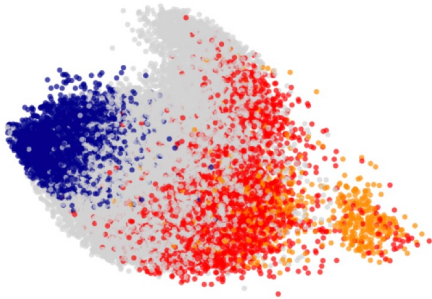
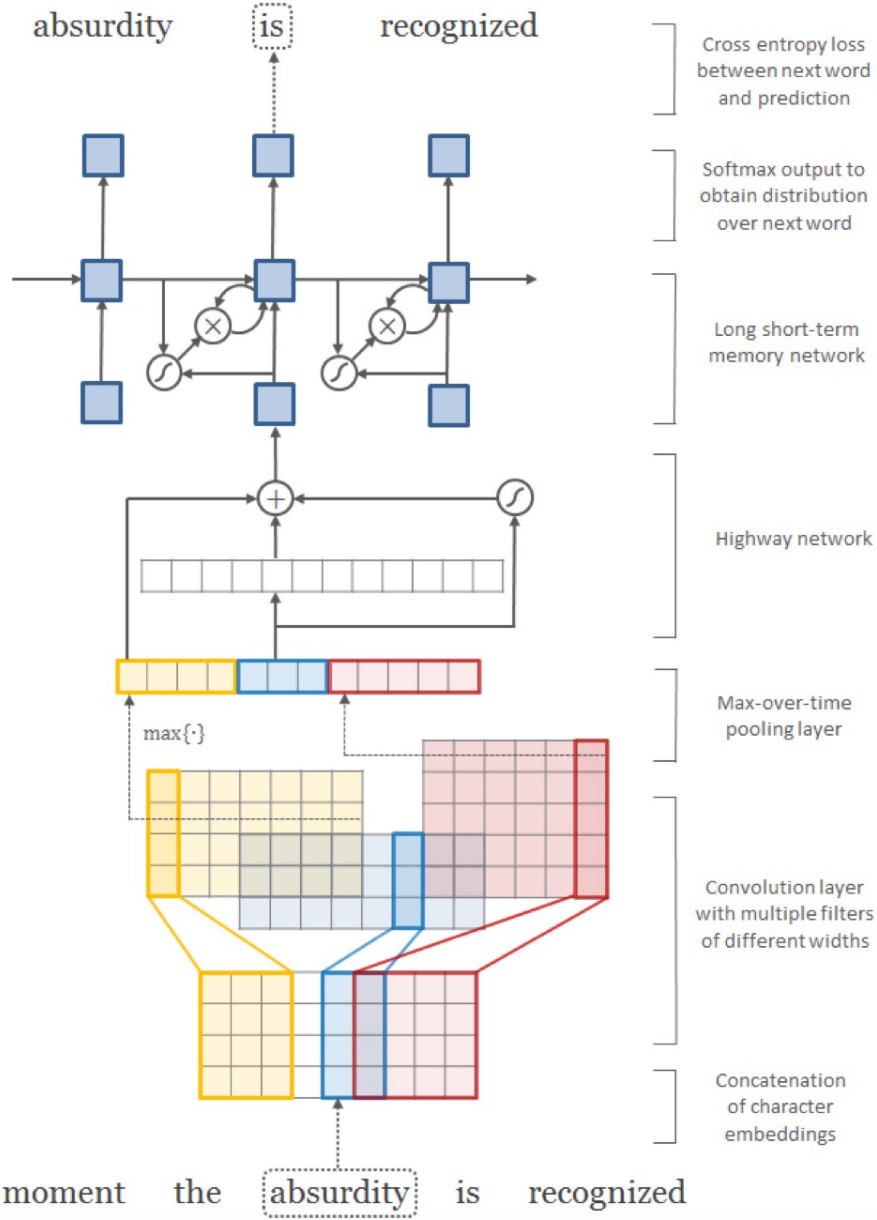
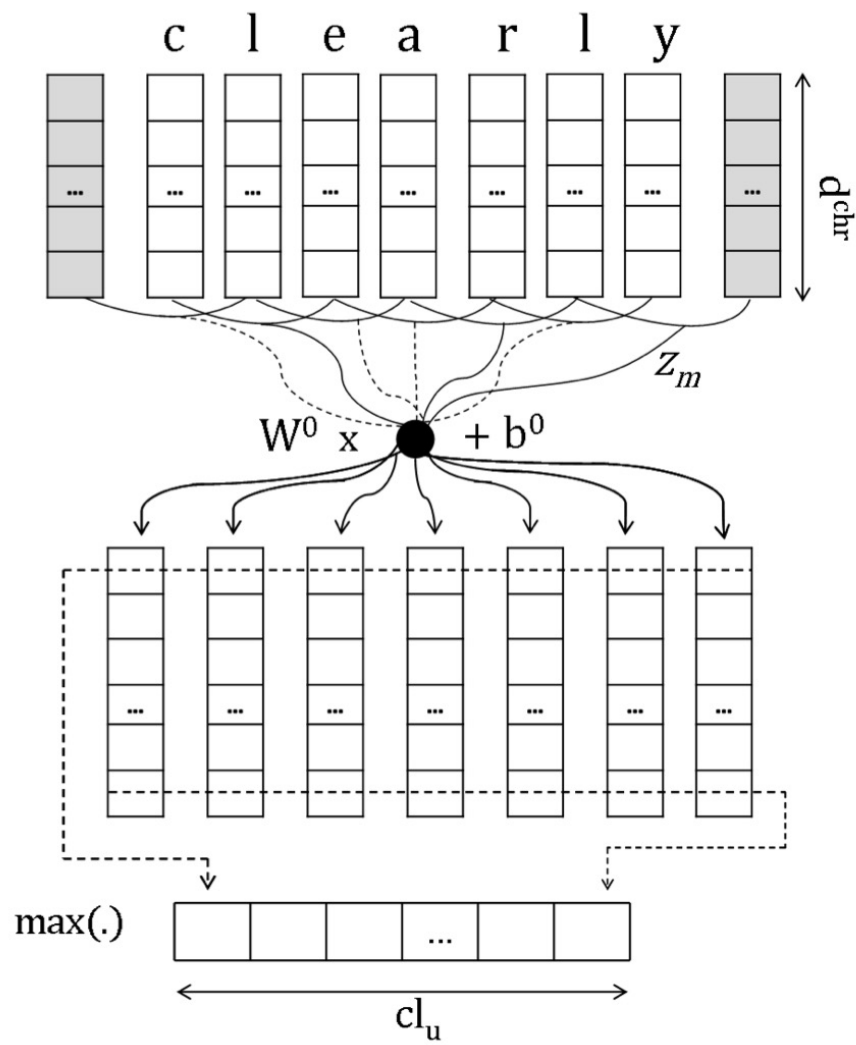


Figure 2: Plot of character n -gram representations via PCA for English. Colors correspond to: prefixes (red), suffixes (blue), hyphenated (orange), and all others (grey). Prefixes refer to character n -grams which start with the start-of-word character. Suffixes likewise refer to character n -grams which end with the end-of-word character.



Word embeddings from character embeddings

Task: part-of-speech tagging



Dos Santos, C., & Zadrozny, B. (2014, June). Learning character-level representations for part-of-speech tagging. In *International Conference on Machine Learning* (pp. 1818-1826). PMLR.

Kim, Y., Jernite, Y., Sontag, D., & Rush, A. (2016, March). Character-aware neural language models. In *Proceedings of the AAAI conference on artificial intelligence* (Vol. 30, No. 1).

CNN word embeddings from character embeddings

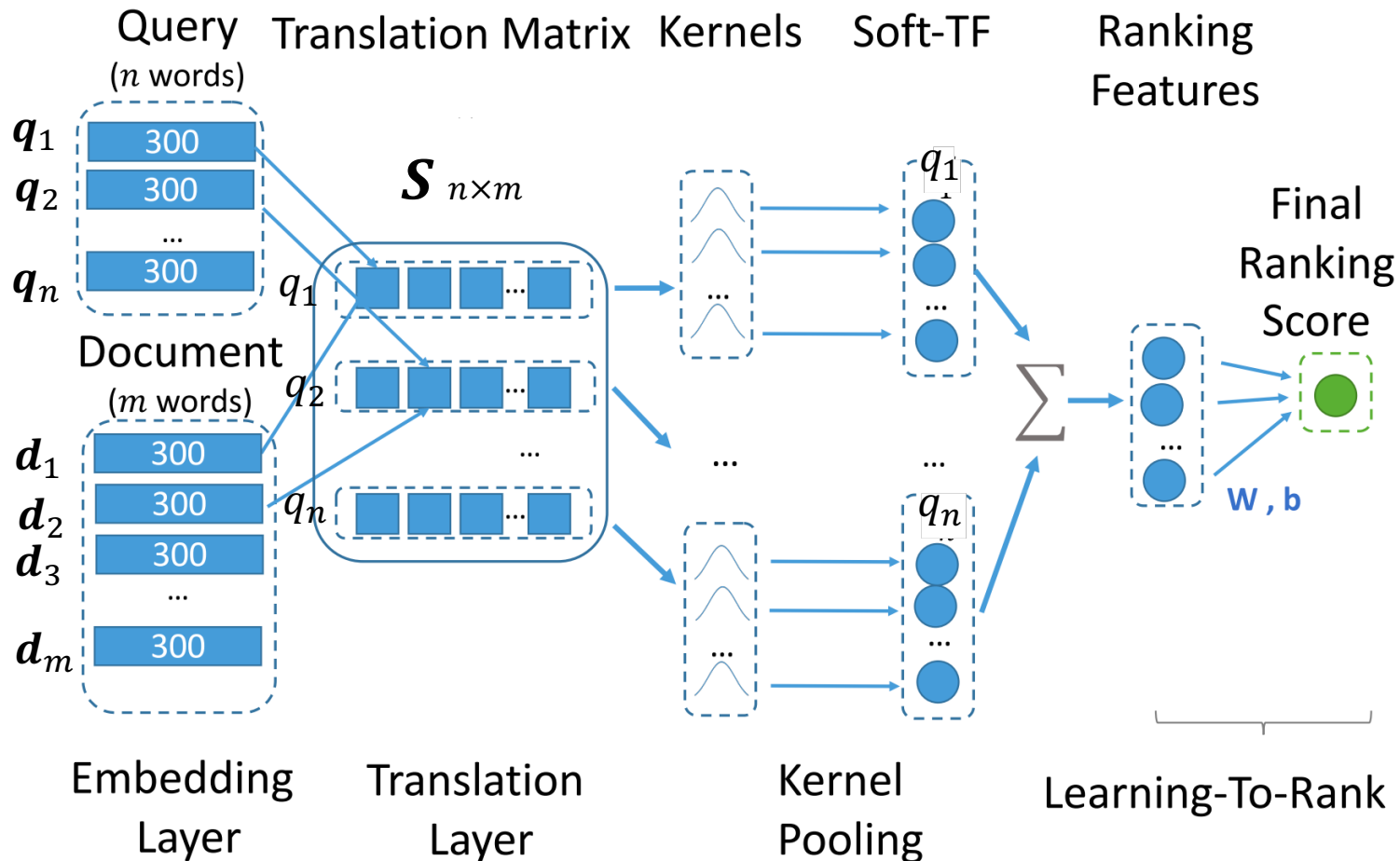
Pros:

- Overall, less parameters in comparison with static word embeddings
- This method resolves the difficulties of handling out-of-vocabularies (OOV)
- Semantic and syntactic regularities are transferred across words, which can benefit some words by providing better generalization

Cons:

- Achieving word embeddings require some computation (feedforward through the CNNs)
- Since every word is composed solely from character embeddings, the quality of some word embeddings might not be as good as static word embeddings

A neural information retrieval model – recap



For details look at Natural Language Processing course - Lecture 6: Information Retrieval with Neural Networks:
<https://www.jku.at/en/institute-of-computational-perception/teaching/alle-lehrveranstaltungen/natural-language-processing/>
Reference: Xiong, C., Dai, Z., Callan, J., Liu, Z., & Power, R. (2017). End-to-end neural ad-hoc ranking with kernel pooling. In *Proceedings of the 40th International ACM SIGIR conference on research and development in information retrieval*

The same model enhanced with n-gram embeddings

