

Machine Learning for Graphs and Sequential Data

Sequential Data – Neural Network Approaches

lecturer: Prof. Dr. Stephan Günnemann
www.daml.in.tum.de

Summer Term 2020

Data Analytics and
Machine Learning 

Roadmap

- Chapter: Temporal Data / Sequential Data
 1. Autoregressive Models
 2. Markov Chains
 3. Hidden Markov Models
 - 4. Neural Network Approaches**
 - a) **Word Vectors**
 - b) RNNs
 - c) Non-Recurrent Models (ConvNets, Transformer)
 5. Temporal Point Processes

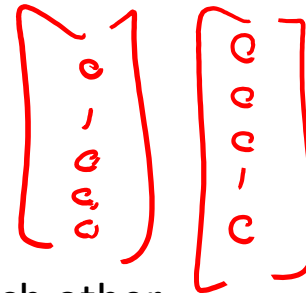
Introduction

- Text is everywhere
- Applying machine learning to textual data to solve machine translation, question answering, sentiment analysis etc.

- Example:

It's a brilliant, honest performance by Nicholson, but the film is an agonizing bore except when the fantastic Kathy Bates turns up.

- Goal: given text predict whether it is positive or negative
- Problem: how to represent words to input them into a subsequent model
- One solution: one-hot encoding
 - High dimensional
 - Too sparse
 - Assumes the words are **independent** of each other



Introduction

- Words as vectors while keeping the underlying language properties
- E.g. similar words should have vectors near each other
- **Distributional hypothesis** – words that appear in similar contexts have similar meanings

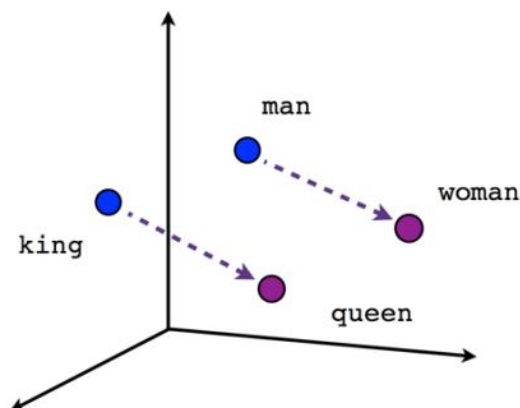
You shall know a word by the company it keeps.

J. R. Firth

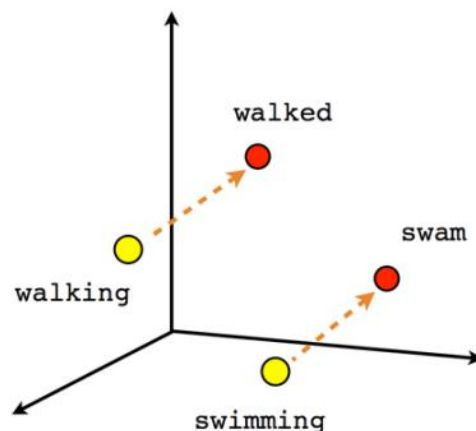
- Example: *hotel* and *motel*
 - Can be used interchangeably in many sentences while remaining meaningful
- However: *duck* – an animal vs. *duck* – to lower head quickly

Introduction

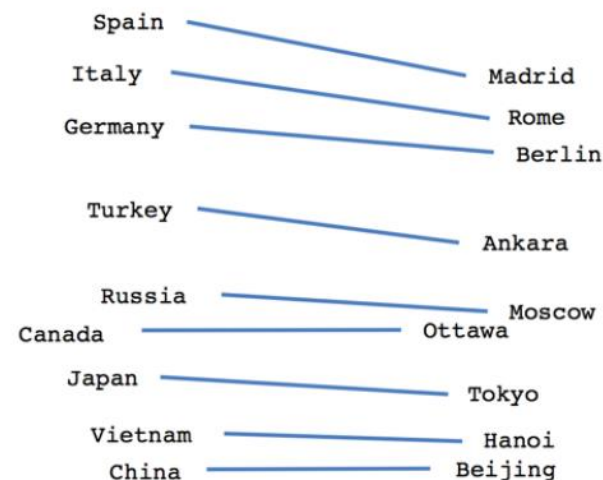
WORD EMBEDDINGS



Male-Female



Verb tense



Country-Capital

Illustration of how vectors can represent linguistic concepts

Figure from <https://www.tensorflow.org/tutorials/representation/word2vec>

Co-occurrence Matrix

- To be aware of **context** we can simply count how many times each word appeared beside other words
- If the text is given with words $\{x_1, \dots, x_N\}$, then a window of size l around a word x_i is $\{x_{i-l}, \dots, x_{i-1}, x_{i+1}, \dots, x_{i+l}\}$
- We slide this window over sentences and count the co-occurrences
- Example:

I like dogs. I like cats too. They hate each other.

- For the first sentence the windows ($l = 1$) are:
 - (\emptyset, like) (I, dogs) $(\text{like}, .)$ (dogs, \emptyset)

i like dogs .

Co-occurrence Matrix

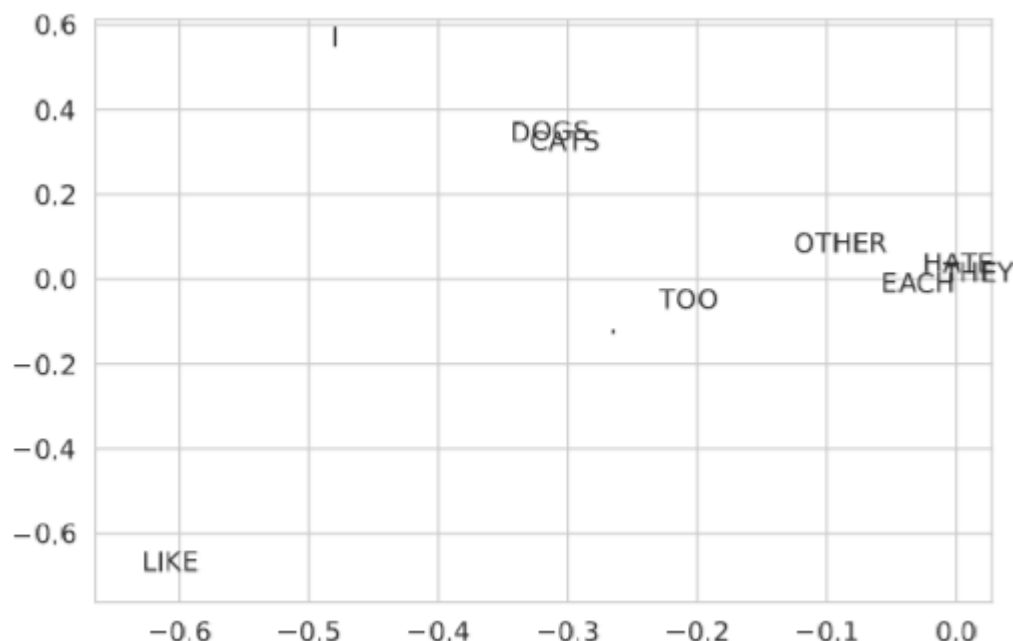
- After counting all the pairs we get a co-occurrence matrix M :

	.	I	cats	dogs	each	hate	like	other	they	too
.	0	0	0	1	0	0	0	1	0	1
I	0	0	0	0	0	0	2	0	0	0
cats	0	0	0	0	0	0	1	0	0	1
dogs	1	0	0	0	0	0	1	0	0	0
each	0	0	0	0	0	1	0	1	0	0
hate	0	0	0	0	1	0	0	0	1	0
like	0	2	1	1	0	0	0	0	0	0
other	1	0	0	0	1	0	0	0	0	0
they	0	0	0	0	0	1	0	0	0	0
too	1	0	1	0	0	0	0	0	0	0

- Pros: similar words have similar vectors
- Cons: still high dimensional and sparse
- Solution: reduce the dimension to get dense vector of fixed dimension

SVD

- We can reduce the dimension with an SVD decomposition: $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$
- If we take the first D columns of \mathbf{U} we get D -dimensional word vectors
- Applied to the previous example ($D = 2$):



- Problems: slow computation and hard to add new words

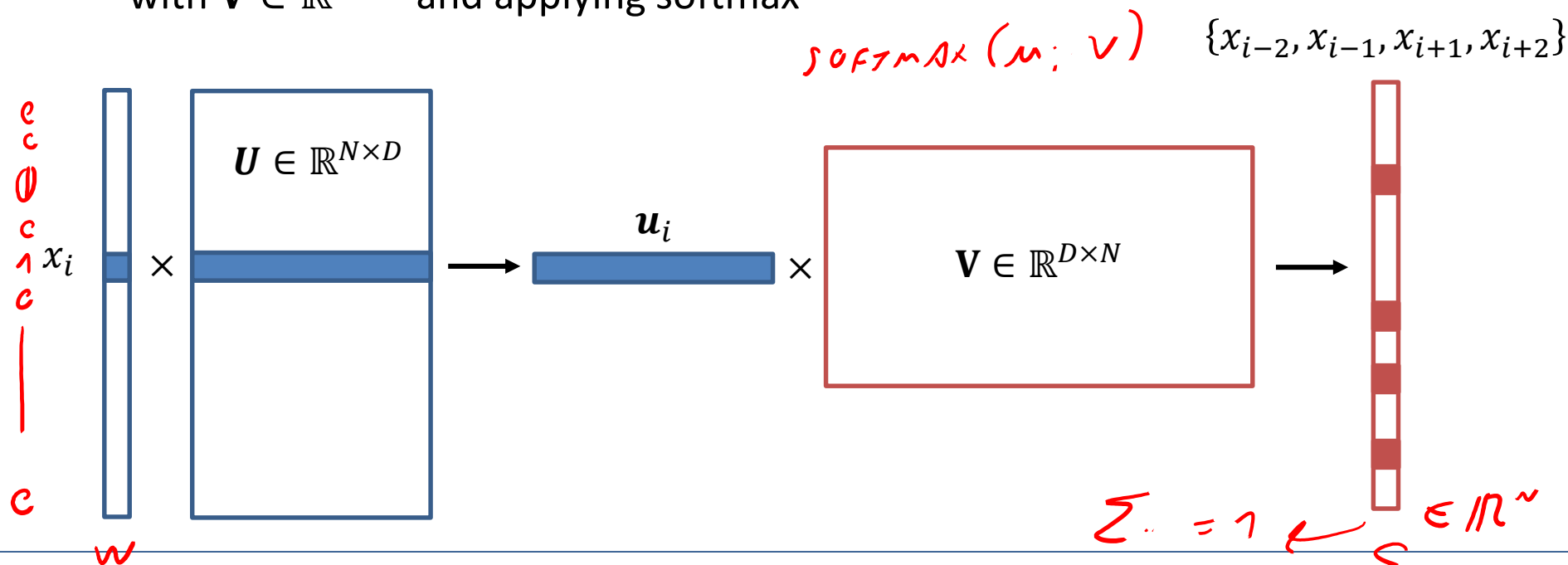
Word2Vec

- A different way to get word vectors is with a neural network
- Task: prediction of words based on context
- Two approaches:
 - **Continuous bag-of-words (CBOW)** C → W
 - Predicts a word from the words surrounding it (window)
 - Not good for rare words because the model might not predict them from context
 - **Skip-gram [1]** W → C
 - Predicts the surrounding context from the current word
 - Given a rare word it must *understand* it to predict the context
 - Slower to train but can work well with smaller amounts of data and with rare words

Skip-gram

$N = 1 \text{ m.e.}$
 $D = 32$
 $w \rightarrow c$

- Input: one-hot vector with dimension N
- Embedding: project the word to D -dimensional space with $U \in \mathbb{R}^{N \times D}$
 - Since input has zeros everywhere except on i th position, multiplication is equivalent to taking i th row of U
- Prediction: get probabilities of context words by multiplying embedding with $V \in \mathbb{R}^{D \times N}$ and applying softmax



Skip-gram

UNSUPERVISED
WORD
EMBEDDINGS

- Formally: if $S = \{x_{i-l}, \dots, x_{i-1}, x_{i+1}, \dots, x_{i+l}\}$ is a window of size l around the word x_i , and θ denotes model parameters, the objective is

$$\max_{\theta} \mathbb{E}[P(S|x_i, \theta)] = \min_{\theta} (-\mathbb{E}[P(S|x_i, \theta)])$$

$$\text{where } P(S|x_i, \theta) = \sum_{x_k \in S} P(x_k|x_i, \theta)$$

$$\text{and } P(x_k|x_i, \theta) = \text{softmax}(\mathbf{u}_i \mathbf{V})$$

$\uparrow f_{\psi}(\mathbf{u}_i)$

- The vector \mathbf{u}_i is the corresponding embedding
- We can choose to set $\mathbf{U} = \mathbf{V}$, giving less parameters to optimize but also less expressiveness

References

- [1] Mikolov, Tomas et al. (2013). “Efficient estimation of word representations in vector space”. In: arXiv preprint arXiv:1301.3781.
- [2] Morin, Frederic and Yoshua Bengio (2005). “Hierarchical probabilistic neural network language model.” In: Aistats. Vol. 5. Citeseer, pp. 246–252.

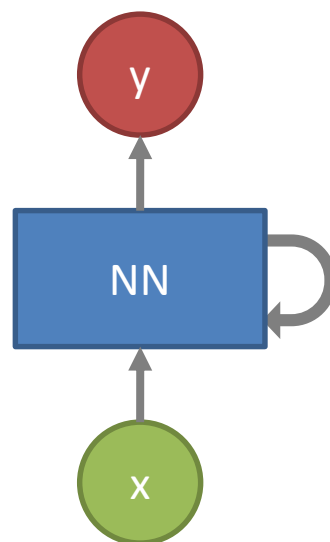
Roadmap

- Chapter: Temporal Data / Sequential Data
 1. Autoregressive Models
 2. Markov Chains
 3. Hidden Markov Models
 - 4. Neural Network Approaches**
 - a) Word Vectors
 - b) RNNs**
 - c) Non-Recurrent Models (ConvNets, Transformer)
 5. Temporal Point Processes

Introduction

$$\text{doc} = \begin{bmatrix} 1 \\ 23 \\ 7 \end{bmatrix} \quad \text{cap} = \begin{bmatrix} \vdots \end{bmatrix}$$

- In word embeddings, we learn a representation for every individual word
- How to process an entire sequence with neural networks?
 - In particular if the sequences have varying length?
- We can use **Recurrent Neural Networks (RNNs)**



$$x_1 \ x_2 \ \dots \ x_k$$

MLP

$$\sigma(w_x \cdot x + b_x)$$

$$f_1(f_1(f_1(x)))$$

PADDING

Definition

SEQUENCE-LEVEL
 $x^{(1)}, \dots, x^{(N)}$

TOKEN-LEVEL
 WORD-LEVEL

POS, NEG, NEUTRAL

- Given a sequence of inputs $\{x^{(1)}, \dots, x^{(N)}\}$ and outputs $\{y^{(1)}, \dots, y^{(N)}\}$ we want to know the probability $P(y^{(t)} | x^{(1)}, \dots, x^{(t)})$
- Represent a sequence $\{x^{(1)}, \dots, x^{(t-1)}\}$ with a hidden state $h^{(t-1)}$
- Neural network takes $h^{(t-1)}$ and current input and maps them to a new hidden state $h^{(t)}$ from which we can predict the output at step t
 - Also use $h^{(t)}$ in the next step

- The update equations are

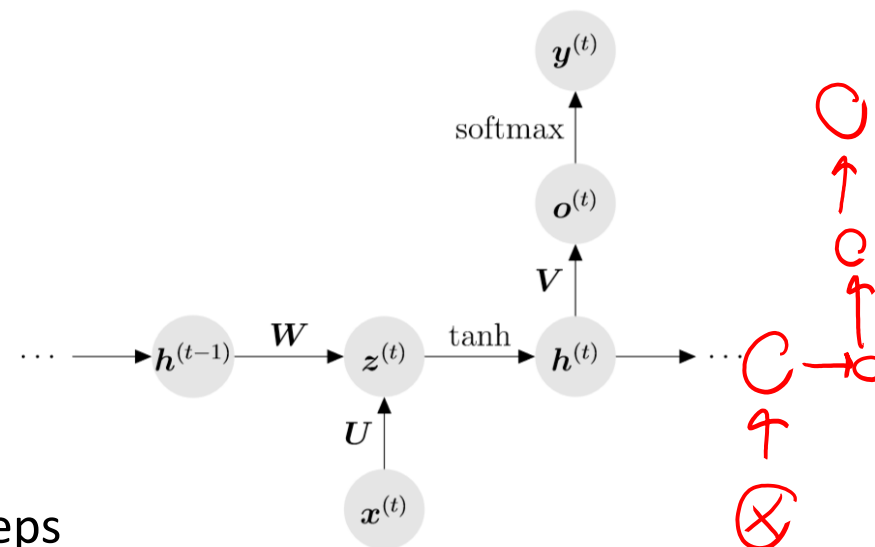
$$z^{(t)} = Wh^{(t-1)} + Ux^{(t)} + b$$

$$h^{(t)} = \tanh(z^{(t)})$$

$$o^{(t)} = Vh^{(t)}$$

$$\hat{y}^{(t)} = \text{softmax}(o^{(t)})$$

- The weights are shared over all time steps



Objective

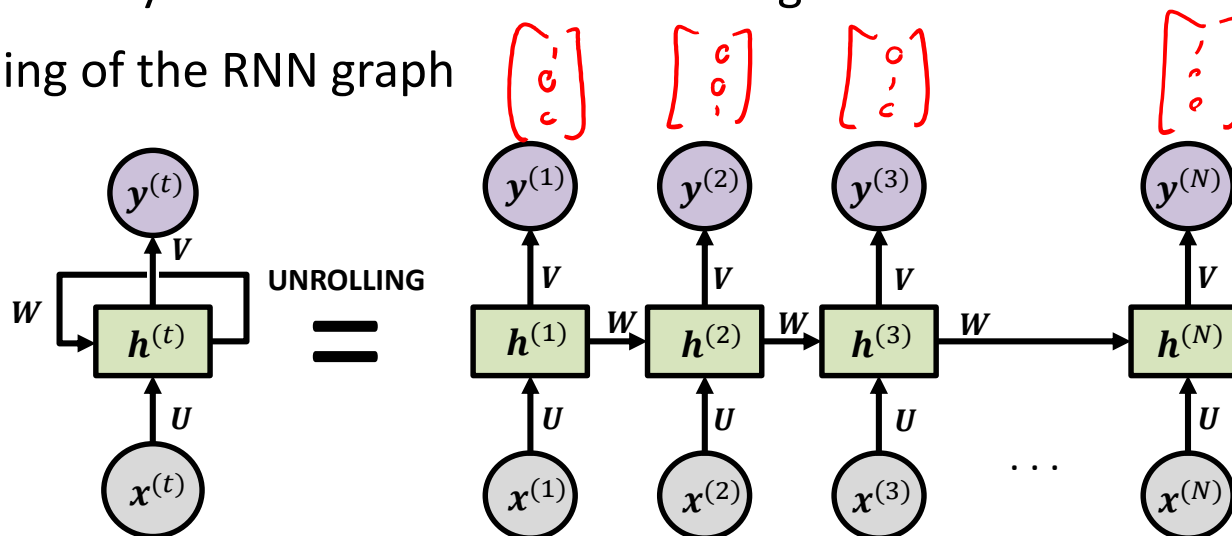
$$\prod_{i=1}^n p(y_i | x_i)$$

- The negative log-likelihood is

$$\Theta = \{U, V, W\}$$

$$\begin{aligned} L &= -\log \prod_t p_{\text{model}}(y^{(t)} | x^{(1)}, \dots, x^{(t)}) \\ &= -\sum_t \log p_{\text{model}}(y^{(t)} | x^{(1)}, \dots, x^{(t)}) = -\sum_t L^{(t)} \end{aligned}$$

- Network fully differentiable – train with a gradient based method
- Unrolling of the RNN graph



$$x \rightarrow h_1 \rightarrow h_2 \rightarrow y$$

Backpropagation through time

- All functions used in the update equations are differentiable (linear, tanh, softmax) → We can compute the derivative w.r.t the parameters:

$$\frac{\partial L}{\partial \mathbf{V}} = \sum_t (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}) (\mathbf{h}^{(t)})^T$$

$$\frac{\partial L}{\partial \mathbf{W}} = \sum_t \text{diag}(1 - (\mathbf{h}^{(t)})^2) \frac{\partial L}{\partial \mathbf{h}^{(t)}} (\mathbf{h}^{(t-1)})^T$$

$$\frac{\partial L}{\partial \mathbf{U}} = \sum_t \text{diag}(1 - (\mathbf{h}^{(t)})^2) \frac{\partial L}{\partial \mathbf{h}^{(t)}} (\mathbf{x}^{(t)})^T$$

$$\mathbf{z}^{(t)} = \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} + \mathbf{b}$$

$$\mathbf{h}^{(t)} = \tanh(\mathbf{z}^{(t)})$$

$$\mathbf{o}^{(t)} = \mathbf{V}\mathbf{h}^{(t)}$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)}) \text{ or } \mathbf{o}^{(t)}$$

- Since parameters are shared over the steps, final derivative are a sum of all the contributions at every step t .

Backpropagation through time

- The hidden state $\mathbf{h}^{(t)}$ recursively depends on all previous hidden states $\mathbf{h}^{(t-1)}, \dots, \mathbf{h}^{(0)}$ i.e.

$$\mathbf{h}^{(t)} = \tanh(\mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} + \mathbf{b})$$

- The gradient $\frac{\partial L}{\partial \mathbf{h}^{(t)}}$ depends on future times

$$\frac{\partial L}{\partial \mathbf{h}^{(t)}} = \mathbf{v}^T (\hat{\mathbf{y}}^{(t)} - \mathbf{y}^{(t)}) + \mathbf{W}^T \text{diag}(1 - (\mathbf{h}^{(t+1)})^2) \frac{\partial L}{\partial \mathbf{h}^{(t+1)}}$$

- The impact of future times might **vanish** or **explode** (e.g. 1-D example: $W > 1$ or $W < 1$) → RNN cannot retain information for many steps.

$$\frac{\partial L}{\partial \mathbf{h}^{(t)}} = \sum_{s=t}^N \frac{\partial L}{\partial \mathbf{h}^{(s)}} \frac{\partial \mathbf{h}^{(s)}}{\partial \mathbf{h}^{(t)}} = \sum_{s=t}^N \frac{\partial L}{\partial \mathbf{h}^{(s)}} \prod_{t+1 \leq k \leq s} \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{h}^{(k-1)}} = \sum_{s=t}^N \frac{\partial L}{\partial \mathbf{h}^{(s)}} \prod_{t \leq k \leq s} W (1 - (h^{(k)})^2)$$

- Solution: change the RNN architecture so it can keep information longer
- Main idea: not every input should be fully taken into account when updating the hidden state – update partially with a *gating* mechanism
- **Gated Recurrent Unit (GRU)** [2]

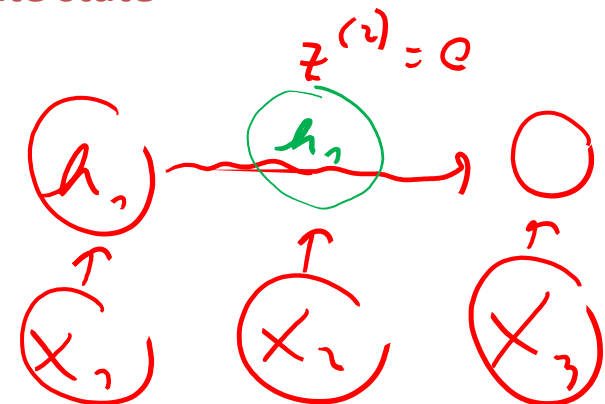
Gates

 $\in [0, 1]$

Simple RNN update – gives candidate state

How much to take from previous state vs. candidate state

$$7(x) = \begin{bmatrix} c.9 \\ c.8 \\ 0 \\ c.21 \end{bmatrix}$$



LSTM

- More powerful architecture: **Long Short-Term Memory (LSTM)** [3]
- Introduces a cell state $\mathbf{c}^{(t)}$ in addition to $\mathbf{h}^{(t)}$ – we have two states

Previous
cell state

$$\mathbf{f}^{(t)} = \sigma \left(\mathbf{W}_f \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right)$$

Forget gate

$$\mathbf{i}^{(t)} = \sigma \left(\mathbf{W}_i \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right)$$

Input gate

$$\mathbf{o}^{(t)} = \sigma \left(\mathbf{W}_o \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right)$$

Output gate

$$\mathbf{c}^{(t)} = \mathbf{f}^{(t)} \odot \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \odot \tanh \left(\mathbf{W} \left[\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)} \right] \right)$$

$$\mathbf{h}^{(t)} = \mathbf{o}^{(t)} \odot \tanh(\mathbf{c}^{(t)})$$

$\in [0, 1]$

Simple RNN update
– LSTM treats it as
an input

Update hidden state (now the output)
using a cell state

Summary

- LSTM and GRU are two examples of improvements to the basic RNN
- Gating enables *skipping* some inputs to capture long-term dependencies
 - Actually, since it uses an element-wise product, it can *remember* or *forget* per individual dimension of a hidden state
 - Avoids gradient problems that RNN has
- It is fully differentiable so we can derive gradients for all the parameters as in the RNN and train it with, e.g., gradient descent
- Many variations on LSTM architecture
 - E.g. *peephole LSTM* – replaces $\mathbf{h}^{(t)}$ with $\mathbf{c}^{(t)}$ in all the equations

References

- [1] Cho, Kyunghyun et al. (2014). “Learning phrase representations using RNN encoder-decoder for statistical machine translation”. In: arXiv preprint arXiv:1406.1078.
- [2] Pascanu, Razvan, Tomas Mikolov, and Yoshua Bengio (2013). “On the difficulty of training recurrent neural networks”. In: International conference on machine learning, pp. 1310–1318.
- [3] Hochreiter, Sepp and Jürgen Schmidhuber (1997). “Long short-term memory”. In: Neural computation 9.8, pp. 1735–1780.
- [4] Peters, Matthew E., Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. "Deep contextualized word representations." *arXiv preprint arXiv:1802.05365* (2018).

Roadmap

- Chapter: Temporal Data / Sequential Data
 1. Autoregressive Models
 2. Markov Chains
 3. Hidden Markov Models
 - 4. Neural Network Approaches**
 - a) Word Vectors
 - b) RNNs
 - c) Non-Recurrent Models (ConvNets, Transformer)**
 5. Temporal Point Processes

Introduction

- Sometimes when modeling a sequence we do not need the complete history to produce the output
- Example: **generating speech**
 - Raw audio has many data points (16000 per second)
 - Important relations on many time scales

- Recall an autoregressive model:

$$X_t = c + \sum_{i=1}^p \varphi_i X_{t-i} + \varepsilon_t$$

- Uses a fixed window of p previous inputs and performs regression
- Can we use neural networks to capture more complex behavior?
 - RNNs share the parameters across time steps, but **depend on full history**
 - We can instead use **Convolutional Neural Networks (ConvNets)**

Recap: Definition

- The convolution $f * g$ of functions f and g is

$$(f * g)(t) = \int_{-\infty}^{+\infty} f(\tau)g(t - \tau)d\tau$$

- In image processing, given an image I and a kernel K , both 2-D matrices, the convolution can be written as:

$$(K * I)(i, j) = \sum_m \sum_n I(i - m, j - n)K(m, n)$$

- Output is again a 2-D matrix (transformed image), where an element (pixel) is a sum of its neighbors, weighted by a kernel

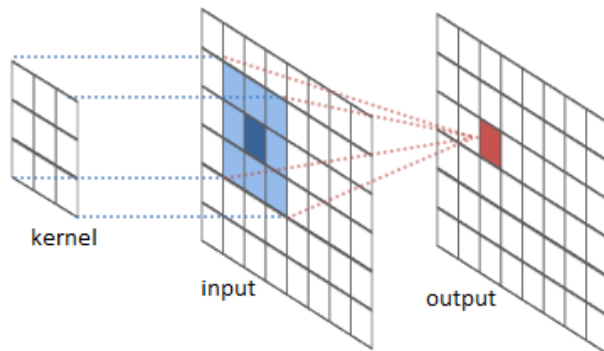


Figure from
<https://intellabs.github.io/RiverTrail/tutorial/>

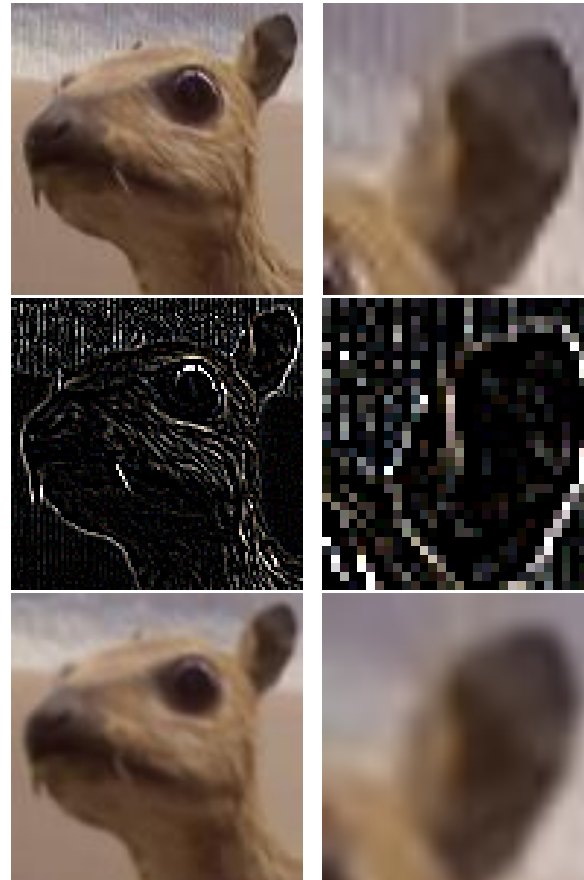
Recap: Examples

Kernel

Output

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

$$\begin{bmatrix} 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \end{bmatrix}$$


Source: [https://en.wikipedia.org/wiki/Kernel_\(image_processing\)](https://en.wikipedia.org/wiki/Kernel_(image_processing))

Recap: Definition

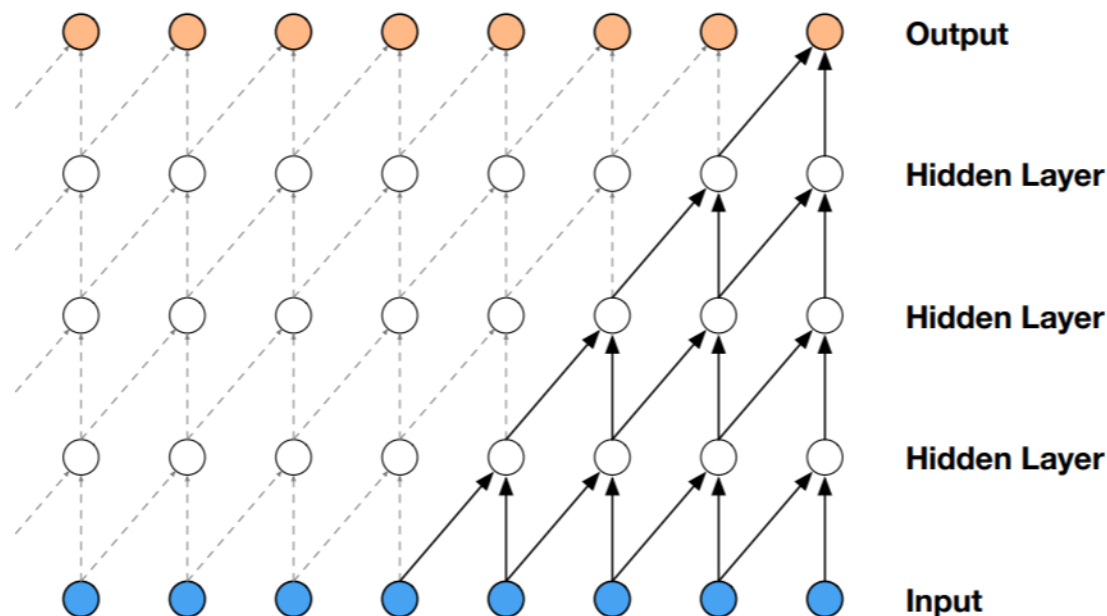
- A **Convolutional Neural Network** contains convolutions as its layers
 - **Kernel** is often called **filter**
 - The parameters/weights of the filter are learned



Example of learned filters [1]

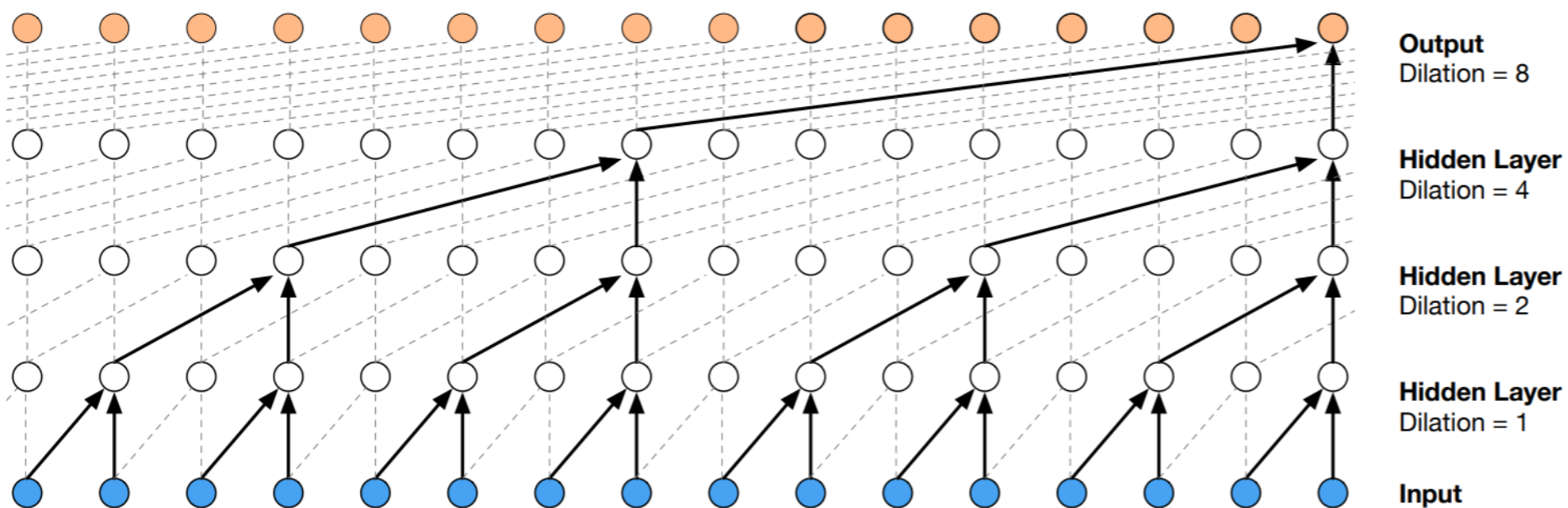
WaveNet

- Sequences are 1-D so we can use a 1-D version of ConvNets
- WaveNet** [2] is an architecture that uses 1-D ConvNets to model speech
 - In addition, it uses special convolutions to ensure causality and increase receptive field
- Causal convolutions** – ensure that the output only depends on the past



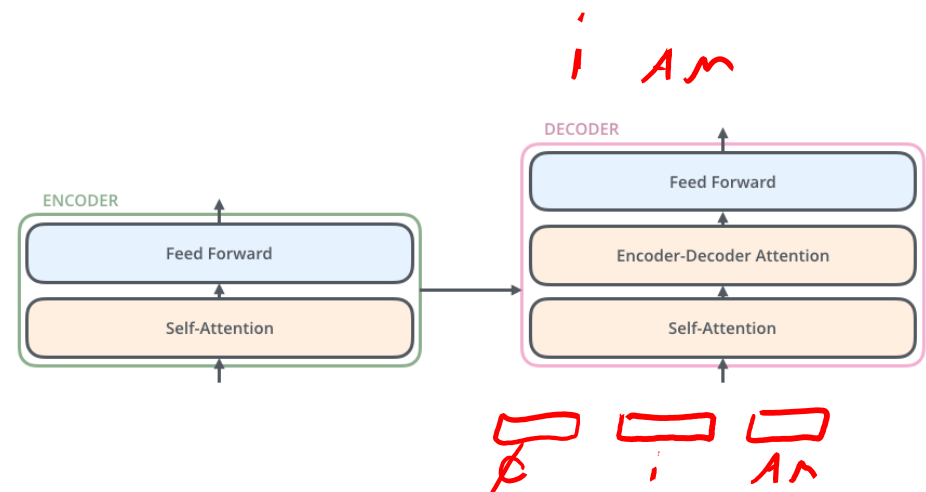
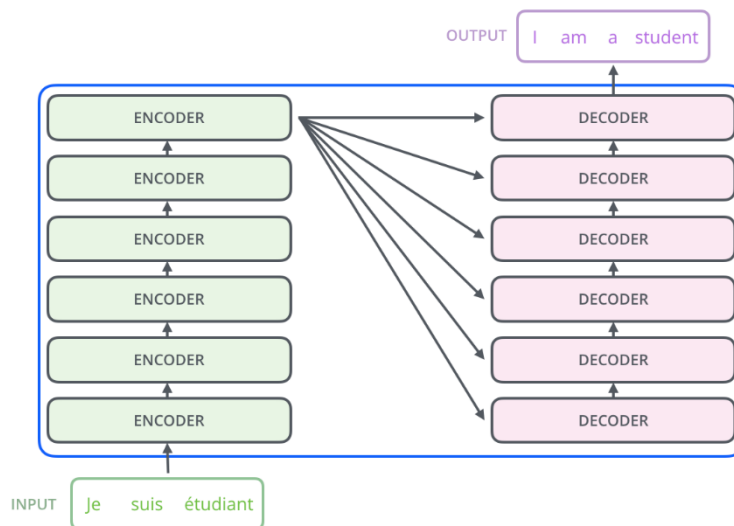
WaveNet

- **Dilated convolutions** – skip some inputs to increase the receptive field
 - Dilation of 1 gives standard convolution
 - If we start with dilation of 1 in the first layer and double it with every layer (2,4,8...) the receptive field will be the exponential of the number of layers
 - E.g. with 4 layers we use 16 inputs in the first layer



Transformers

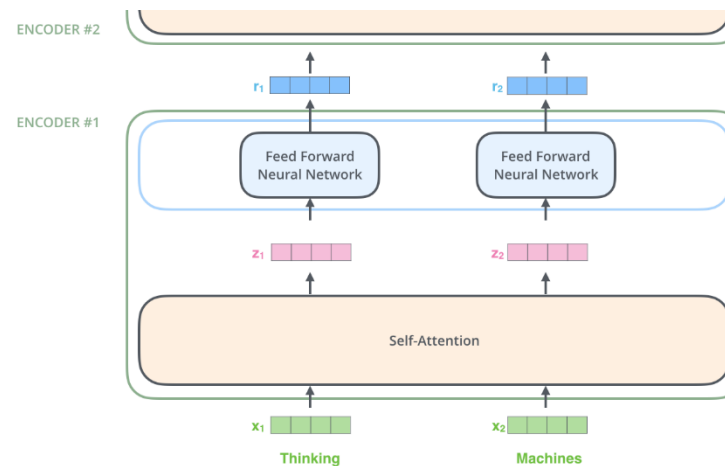
- **Transformers** [3] are fast models using attention mechanisms
 - Like WaveNet, it is not a recurrent neural network → parallelizable and fast
 - It is composed of a stack of encoders and decoders using self-attention
 - Achieves high performances in NLP translations



The following images are taken from [4] Jalammar blog

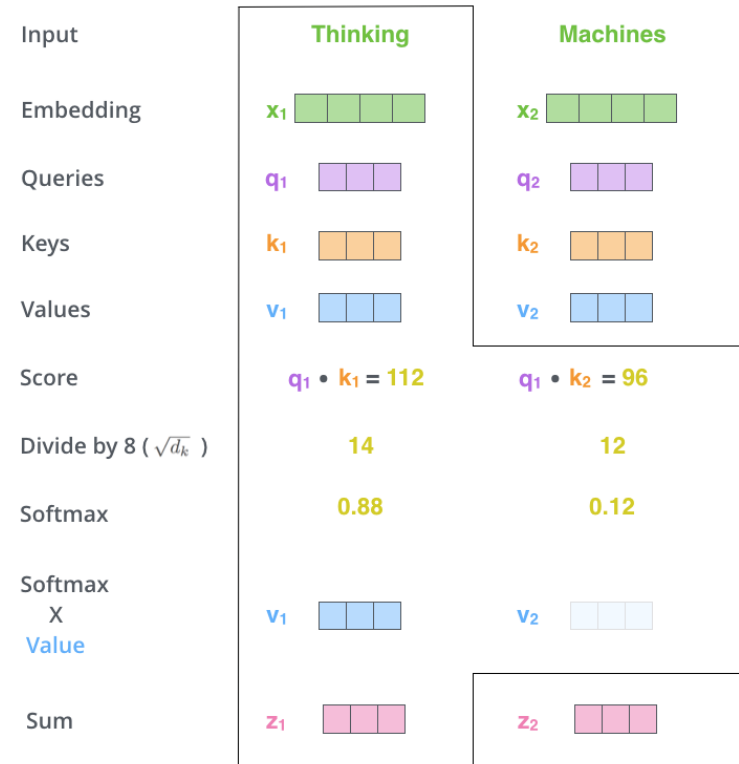
Transformers

- Words are represented with **embeddings** and flow all at once during training.
- Below illustrated by an encoder block
 - The self-attention layer “couples” the embeddings
 - The rest handles the embeddings independently



Transformers – Attention

- Attention is a learned weighting over the elements x_j (given element x_i)
 - The weighting is computed by applying softmax to query/key scores
 - Query depends on x_i ; key on x_j
 - The weight indicates how much of v_j we use (the “value” of x_j)
- Self-attention: the attention is on the input signal itself
- It is easy computable in a matrix formulation.



$$\text{softmax} \left(\frac{Q \times K^T}{\sqrt{d_k}} \right) V$$

$$= Z$$

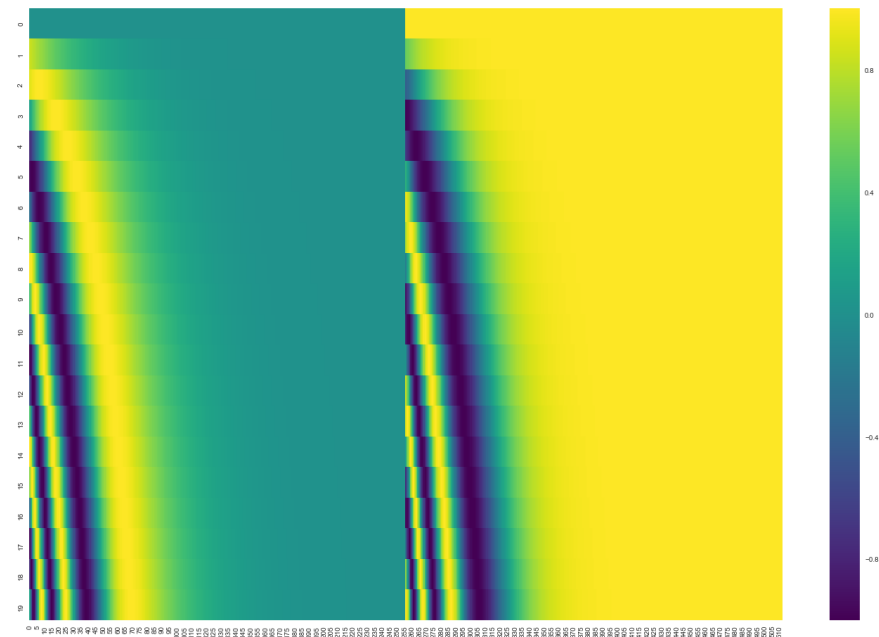
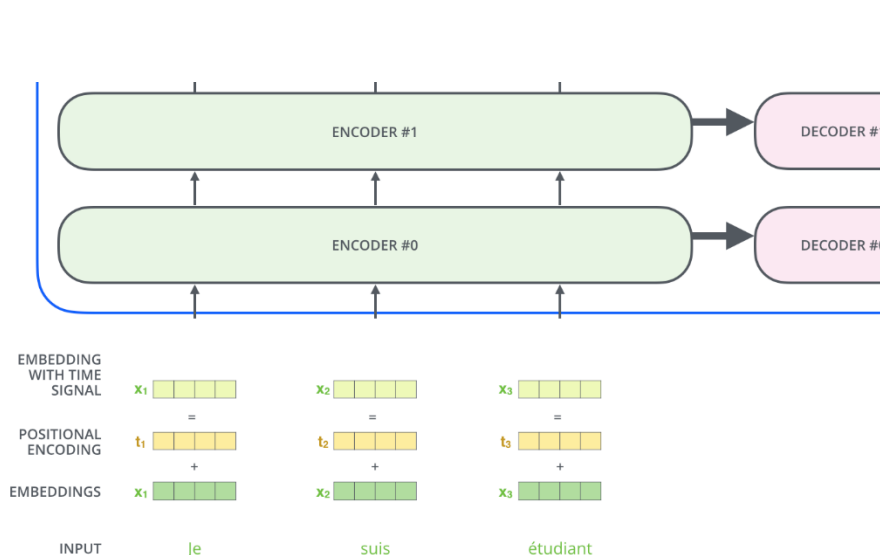
Hand-drawn red annotations on the matrix equation:

- A red box highlights the top-left element of the $Q \times K^T$ product.
- Red arrows point from the elements of the Q matrix to the corresponding elements in the product matrix.
- Red arrows point from the elements of the K^T matrix to the corresponding elements in the product matrix.
- Red arrows point from the elements of the V matrix to the final result Z .

“Self-attention allows the model to look at other positions in the input sequence for clues that can help lead to a better encoding for this word” [4]

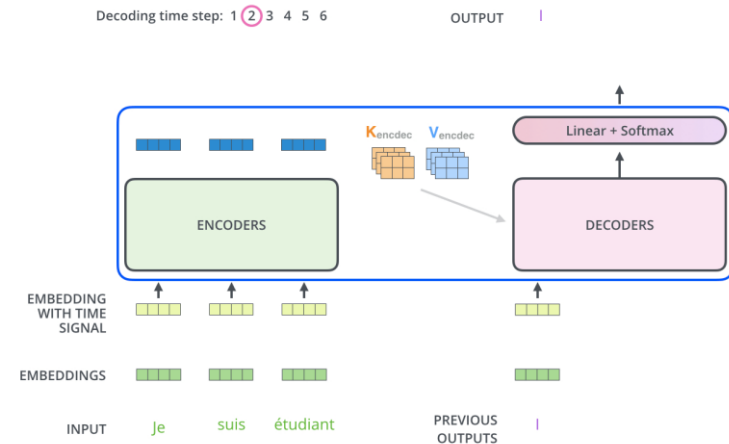
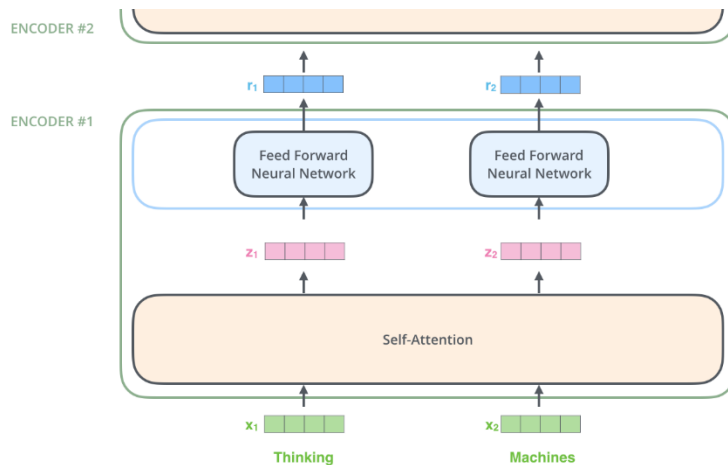
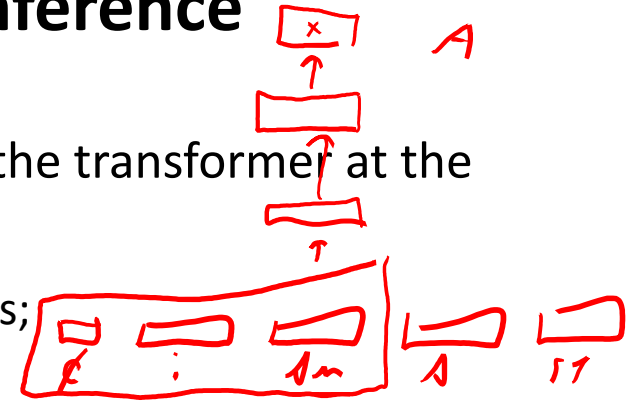
Transformers – Positional Encoding

- A “plain” transformer actually does not care about the order, i.e. the architecture itself is not aware of the non-i.i.d. nature
- Transformers use a **positional encoding** to represent the word order. Positional encoding: meaningful static vectors which are concatenated with the word embeddings.



Transformers: Notes on Training and Inference

- During training, the embeddings flow all through the transformer at the same time/in parallel
 - This enables efficient training on very large datasets; crucial for the success of recent models
- At inference time, decoding is done one step after the other until the end of sentence symbol is reached.



Questions – NN

1. In an RNN, the hidden state at a given time influences all hidden states into the future. However, an RNN cannot model long-term dependencies. Why?
2. What is the receptive field of a causal convolution and dilated convolution with n layers ?

References

- [1] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." In *Advances in NIPS*, pp. 1097-1105. 2012.
- [2] Van Den Oord, Aäron, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. "WaveNet: A generative model for raw audio." *SSW 125* (2016).
- [3] Ashish Vaswani et al., "Attention Is All You Need" NIPS 2017
- [4] Jalammar blog, <http://jalammar.github.io/illustrated-transformer/>