# Attention
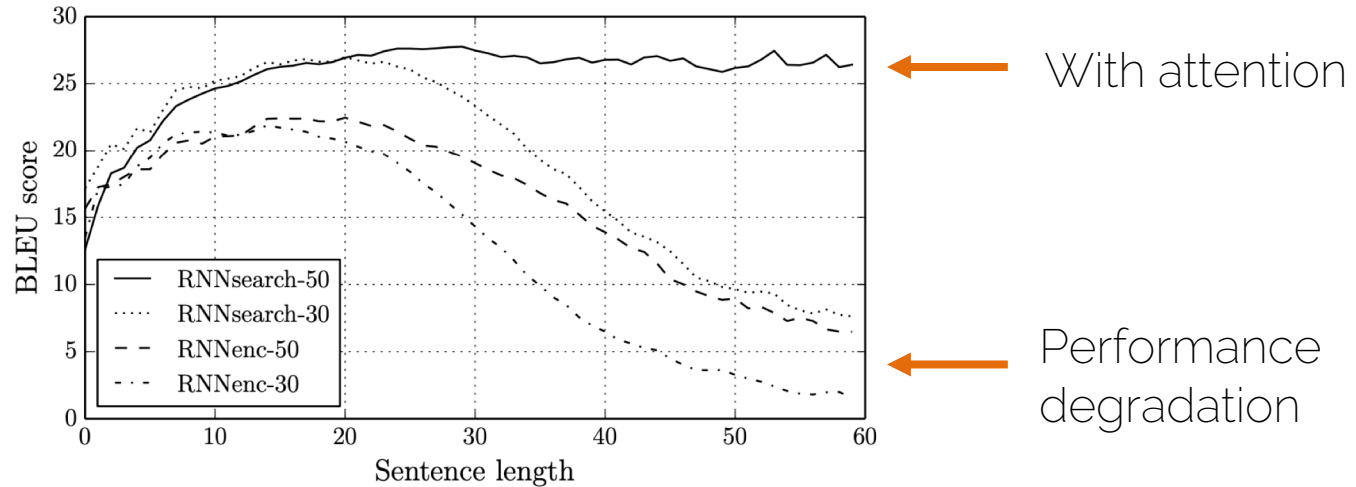
# The problem

- For very long sentences, the score for machine translation really goes down after 30-40 words.



With attention

Performance degradation

Bahdanau et al 2014. Neural machine translation by jointly learning to align and translate.

# Basic structure of a RNN

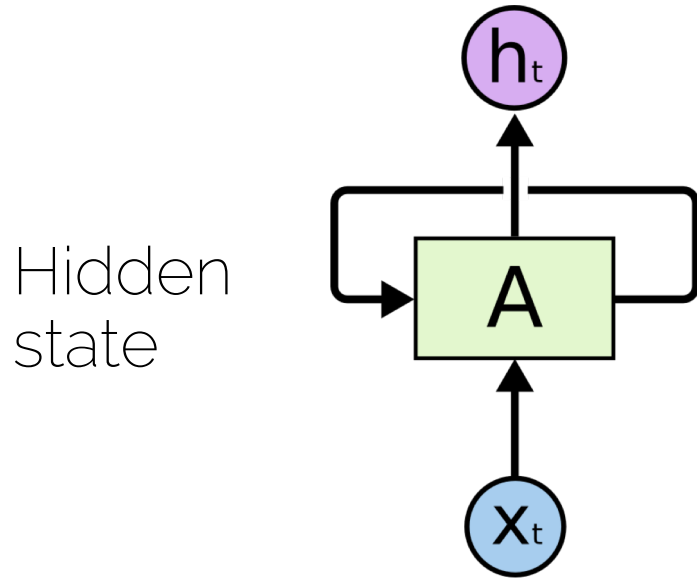- We want to have notion of "time" or "sequence"



$$\mathbf{A}_t = \boldsymbol{\theta}_c \mathbf{A}_{t-1} + \boldsymbol{\theta}_x \mathbf{x}_t$$

Hidden state

Previous hidden state

input

[Christopher Olah] Understanding LSTMs

# Basic structure of a RNN

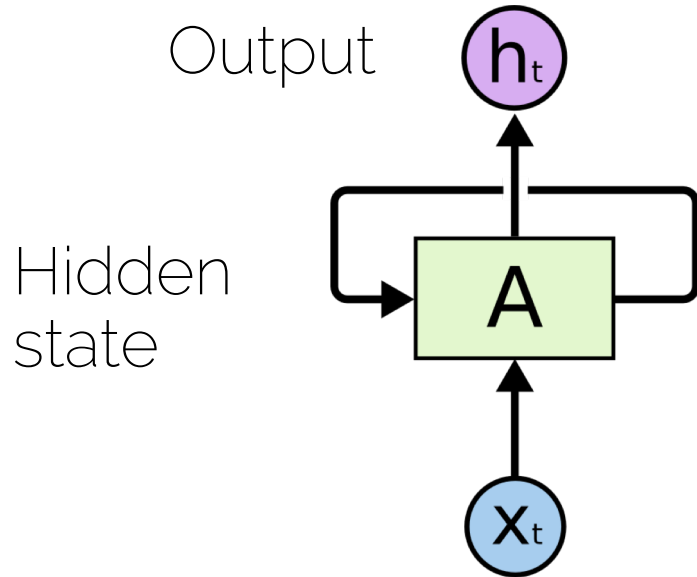- We want to have notion of "time" or "sequence"

Hidden
state

$$\mathbf{A}_t = \boldsymbol{\theta}_c \mathbf{A}_{t-1} + \boldsymbol{\theta}_x \mathbf{x}_t$$

Parameters to be learned

# Basic structure of a RNN

- We want to have notion of "time" or "sequence"

Output

$h_t$

Hidden state

$A$

$x_t$

$$\mathbf{A}_t = \boldsymbol{\theta}_c \mathbf{A}_{t-1} + \boldsymbol{\theta}_x \mathbf{x}_t$$
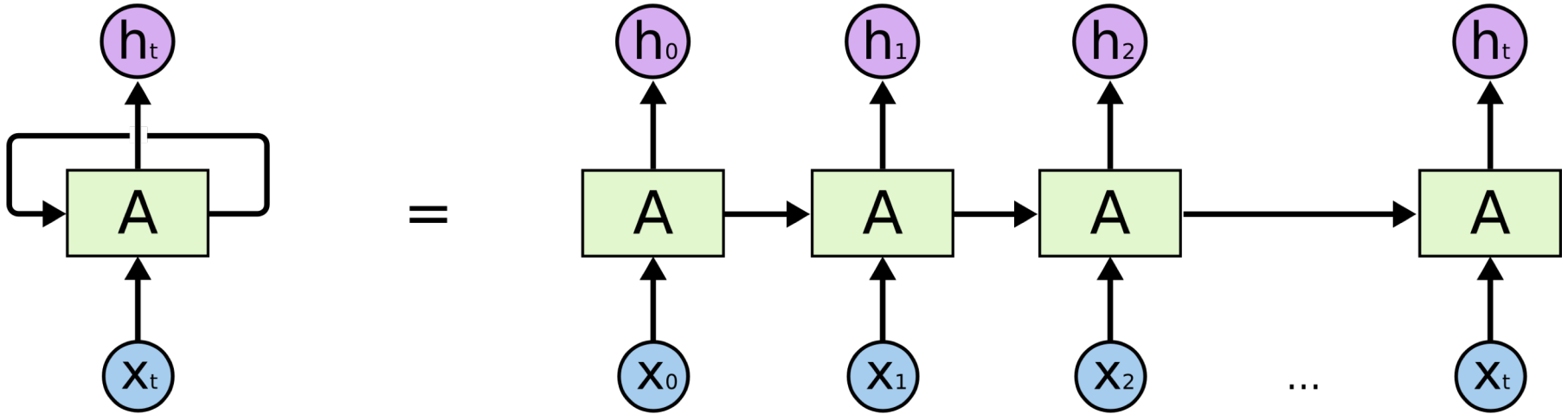
$$\mathbf{h}_t = \boldsymbol{\theta}_h \mathbf{A}_t$$

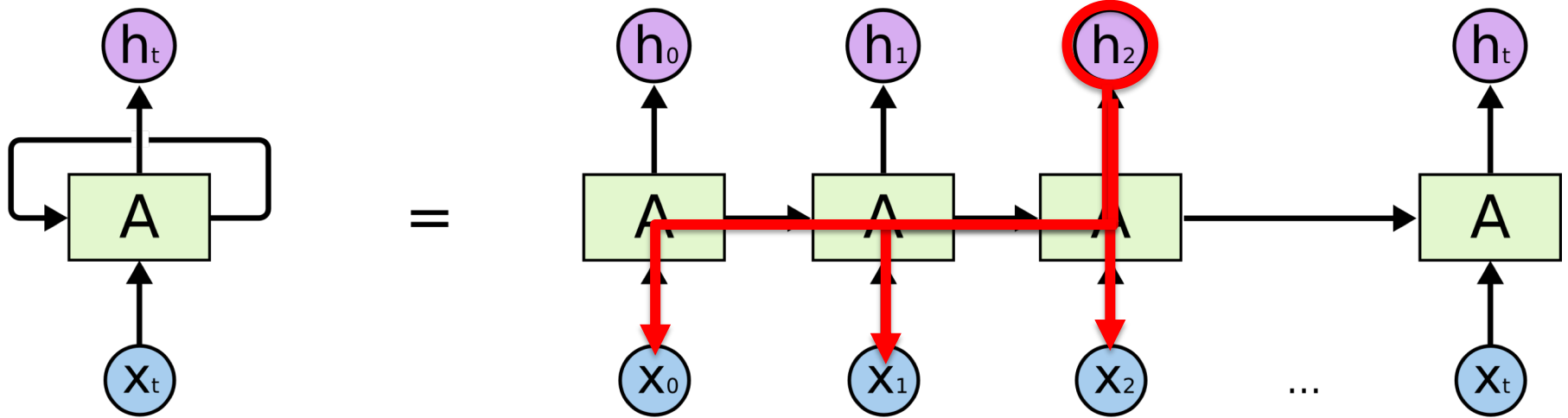Same parameters for each time step = generalization!

# Basic structure of a RNN

- Unrolling RNNs

Hidden state is the same

[Christopher Olah] Understanding LSTMs
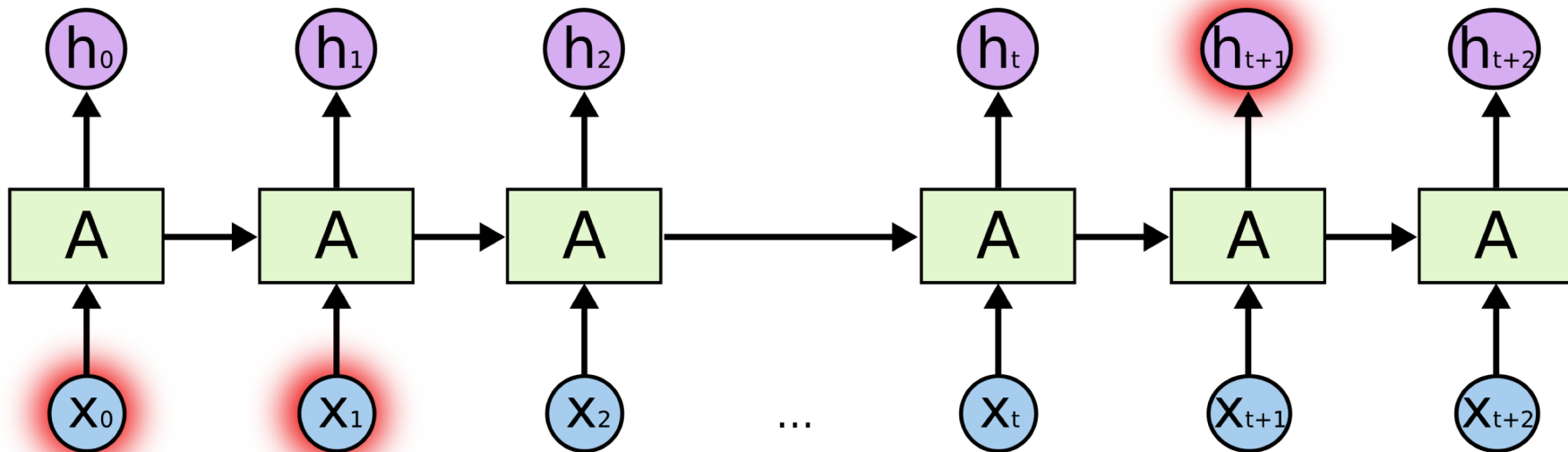
# Basic structure of a RNN

- Unrolling RNNs
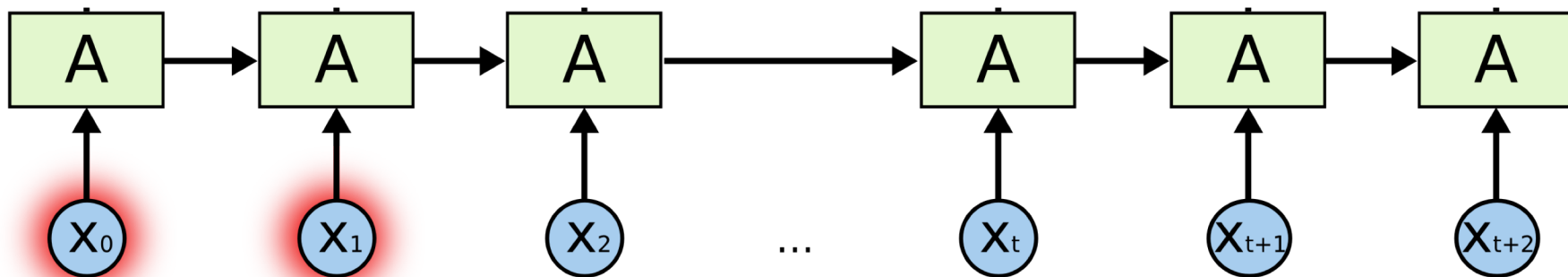
# Long-term dependencies



I moved to Germany ...                    so I speak German fluently
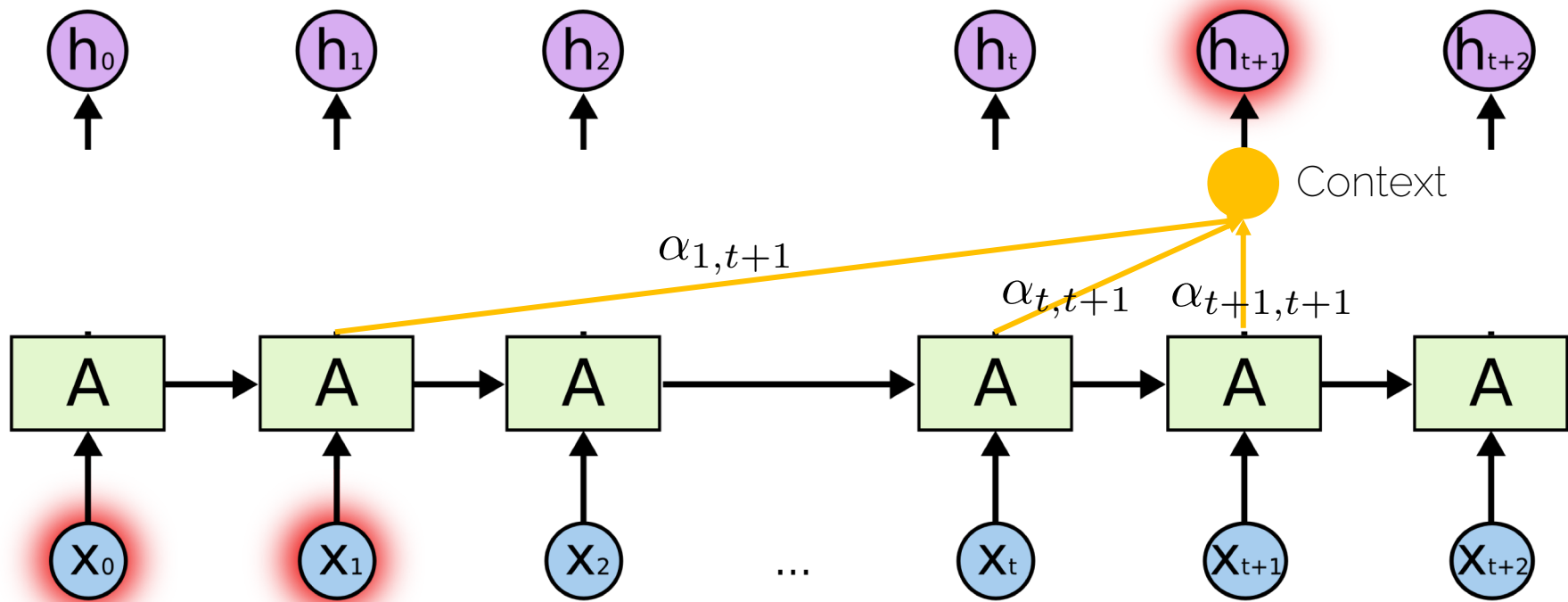
# Attention: intuition



ATTENTION: Which hidden states are more important to predict my output?

I moved to Germany … so I speak German fluently

# Attention: intuition



$\alpha_{1,t+1}$

$\alpha_{t,t+1}$   $\alpha_{t+1,t+1}$
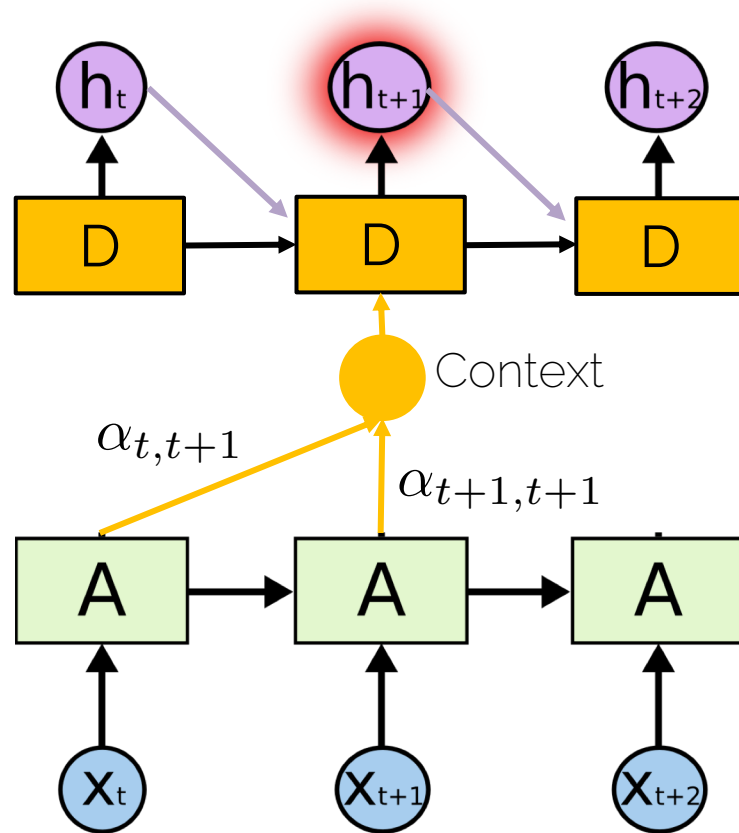
Context

I moved to Germany ...          so I speak German fluently

# Attention: architecture

- A decoder processes the information

- Decoders take as input:
  - Previous decoder hidden state
  - Previous output
  - Attention

# Attention

- $\alpha_{1,t+1}$ indicates how much the word in the position $1$ is important to translate the work in position $t+1$

- The context aggregates the attention

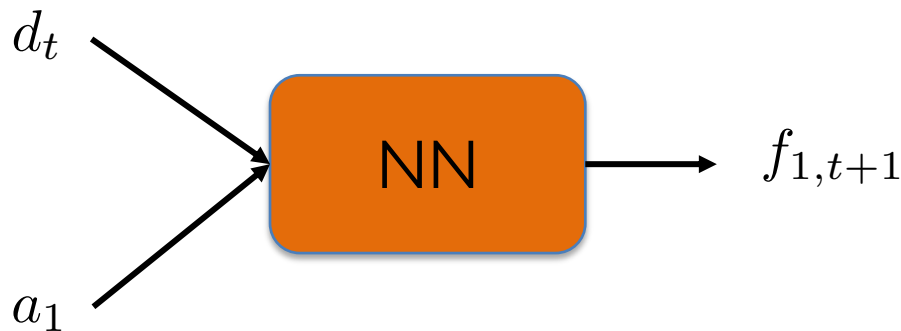$$c_{t+1} = \sum_{k=1}^{t+1} \alpha_{k,t+1} a_k$$

- **Soft** attention: All attention masks alpha sum up to 1

# Computing the attention mask

- We can train a small neural network

Previous state of
the decoder $\quad d_t$

Hidden state of
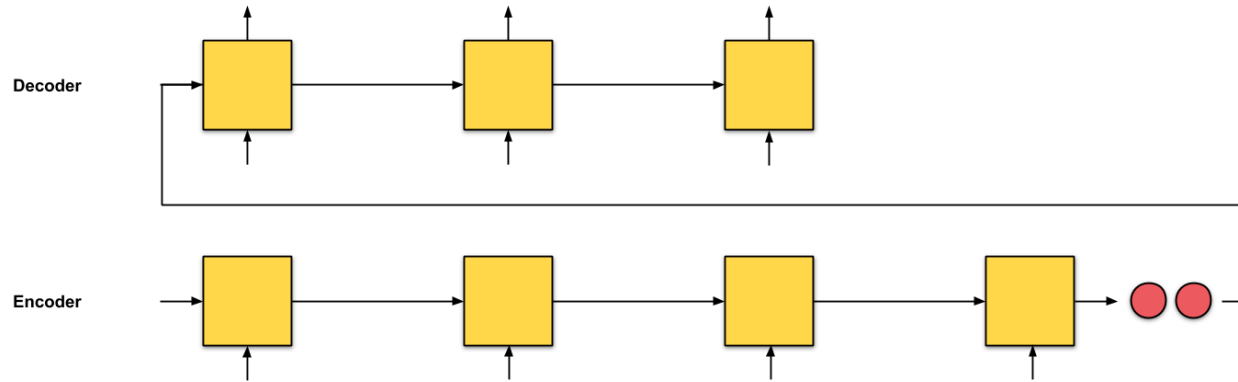the encoder $\quad a_1$

NN

$f_{1,t+1}$

- Normalize $\quad \alpha_{1,t+1} = \dfrac{\exp^{f_{1,t+1}}}{\sum_{k=1}^{t+1} \exp^{f_{k,t+1}}}$

- Animations? As a summary with the se2seq example in here [https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3](https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3)

# Seq2Seq

- How do we translate?
- First read *the whole* sentence in language 1.
- *Afterwards*, translate the whole sentence in language 2.



Sutskever et al. „Sequence to Sequence Learning with Neural Networks". NIPS 2014
Picture from: https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3

# Seq2Seq + Attention?

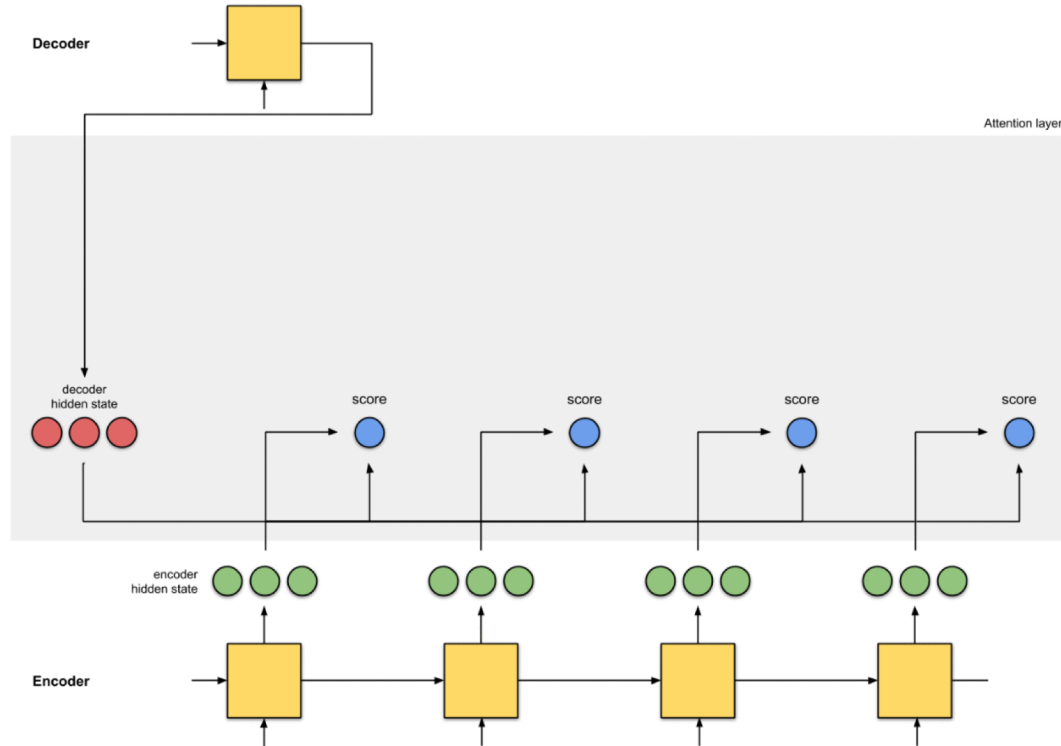- If the sentence is very long, we might have forgotten what was said at the beginning.

- Solution: take "notes" of keywords as we read the sentence in language 1.

- Use attention!

# Seq2Seq + Attention



Animation from: https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3

# Seq2Seq + Attention



Animation from: https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3

# Seq2Seq + Attention



Animation from: https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3

# Seq2Seq + Attention



Animation from: https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3

# Seq2Seq + Attention



Animation from: https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3

# Seq2Seq + Attention



Animation from: https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3

# Attention for vision

# Why do we need attention?

- We use the whole image to make the classification



- Are all pixels equally important?

# Why do we need attention?

- Wouldn't it be easier and computationally more efficient to just run our classification network on the patch?

# Soft attention for captioning

# Image captioning



14x14 Feature Map

1. Input Image
2. Convolutional Feature Extraction
3. RNN with attention over the image
4. Word by word generation

A bird flying over a body of water

Xu et al 2015. Show attention and tell: neural image caption generation with visual attention.

# Image captioning

- Input: image
- Output: a sentence describing the image.
- **Encoder**: a classification CNN (VGGNet, AlexNet). This computes a feature maps over the image.
- **Decoder**: an attention-based RNN
  - In each time step, the decoder computes an attention map over the entire image, effectively deciding which regions to focus on.
  - It receives a context vector, which is the weighted average of the conv net features.

# Conventional captioning



LSTM only sees the image once!

# Attention mechanism

A girl is throwing a frisbee in the park

# Attention mechanism



A girl is throwing a frisbee in the park

# Attention mechanism



A girl is throwing a frisbee in the park

# Attention mechanism



A girl is throwing a frisbee in the park

# Attention mechanism



$y_i$: Output of encoder are the image features which still retain spatial information (no FC layer!)

$Z_i$: Output of attention model

$h_i$: Hidden state of LSTM

# Attention mechanism



How does the attention model look like?

# Attention model

- Attention architecture



Output attention

Any past hidden state

Visual features

Image: https://blog.heuritech.com/2016/01/20/attention-mechanism/

# Attention model

- Inputs = feature descriptor for each image patch

# Attention model

- Inputs = feature descriptor for each image patch



Still related to the spatial location of the image

# Attention model

- We want an bounded output

$$m_i = \tanh(W_{cm}c + W_{ym}\, y_i)$$

# Attention model

- Softmax to create the attention values between 0 and 1

# Attention model

- Multiplied by the image features → ranking by importance

# Hard attention model

- Choosing one of the features by sampling with probabilities $s_i$

# Types of attention

- **Soft attention**: deterministic process that can be backproped

- **Hard attention**: stochastic process, gradient is estimated through Monte Carlo sampling.

- Soft attention is the most commonly used since it can be incorporated into the optimization more easily

# Types of attention

- Soft vs hard attention



Soft

Hard

A    bird    flying    over    a    body    of    water    .

# Types of attention: soft



Image:
H x W x 3

CNN

Grid of features
(Each D-
dimensional)

a | b
c | d

Attention
module

Distribution over
grid locations
$p_a + p_b + p_c + p_c = 1$

$p_a$ | $p_b$
$p_c$ | $p_d$

Final context

- Can be backproped
- Uses all the image

Image: Stanford CS231n lecture

# Types of attention: hard



Input image:
H x W x 3

Box Coordinates:
(xc, yc, w, h)

Cropped and
rescaled image:
X x Y x 3

Gradient is 0 almost everywhere
Gradient is undefined at x = 0

- You can view it as an image cropping!

- If we cannot use gradient descent, what alternative could we use to train this function?

Reinforcement Learning

Image: Stanford CS231n lecture

# Image captioning with attention



A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

A little <u>girl</u> sitting on a bed with a teddy bear.

A group of <u>people</u> sitting on a boat in the water.

A giraffe standing in a forest with <u>trees</u> in the background.

Xu et al 2015. Show attention and tell: neural image caption generation with visual attention.

# Interesting works on attention

- Luong et al, "Effective Approaches to Attentionbased Neural Machine Translation," EMNLP 2015
- Chan et al, "Listen, Attend, and Spell", arXiv 2015
- Chorowski et al, "Attention-based models for Speech Recognition", NIPS 2015
- Yao et al, "Describing Videos by Exploiting Temporal Structure", ICCV 2015
- Xu and Saenko, "Ask, Attend and Answer: Exploring Question-Guided Spatial Attention for Visual Question Answering", arXiv 2015
- Zhu et al, "Visual7W: Grounded Question Answering in Images", arXiv 2015
- Chu et al. „Online Multi-Object Tracking Using CNN-based Single Object Tracker with Spatial-Temporal Attention Mechanism". ICCV 2017

# Conditioning

# When do we need conditioning?

- Scene understanding from an image and an audio source. Both need to be processed!

# When do we need conditioning?

- Visual Question and Answering: the sentence (question) needs to be understood, the image is needed to create the answer.



Are there an equal number of large things and metal spheres?

# When do we need conditioning?

- Visual Question and Answering: the sentence (question) needs to be understood, the image is needed to  create the answer.



Are there an equal number of large things and metal spheres?

# When do we need conditioning?

- We have two sources, can we process one **in the context** of the other?

- **Conditioning**: the computation carried out by a model is conditioned or *modulated* by information extracted from an auxiliary input.

- Note: a similar thing can be obtained with attention

# When do we need conditioning?

- Generate images based on a word
- Do we need to retrain a model for each word?



Image: https://distill.pub/2018/feature-wise-transformations/

# Concatenation-based conditioning



"puppy"

conditioning representation

**Concatenation-based conditioning** simply concatenates the conditioning representation to the input.

input

Image: https://distill.pub/2018/feature-wise-transformations/

# Concatenation-based conditioning



"puppy"

conditioning
representation

**Concatenation-based conditioning** simply concatenates the conditioning representation to the input.

input

concatenate

Image: https://distill.pub/2018/feature-wise-transformations/

# Concatenation-based conditioning

"puppy"

conditioning
representation

**Concatenation-based conditioning** simply concatenates the conditioning representation to the input.

input

concatenate

The result is passed through a linear layer to produce the output.

linear

output

Image: https://distill.pub/2018/feature-wise-transformations/

# Concatenation-based conditioning

- Source: image (high-dimensional) and pose (low-dimensional) → expressed as an image (same dimensionality)



Condition image + Target pose → Generated image

L. Ma et al. „Pose guided person image generation". NIPS 2017

# Concatenation-based conditioning

- Source: image (high-dimensional) and pose (low-dimensional)
  → expressed as an image (same dimensionality)



Wait for the GAN intro in a few weeks!

L. Ma et al. „Pose guided person image generation". NIPS 2017

# Concatenation-based conditioning

- Sources: image (high-dimensional) and measurements (low-dimensional)



A. Dosovitskiy and V. Koltun. Learning to act by predicting the future. ICLR 2017

# Conditional biasing

Think about the similarities with concatenation -based conditioning



**Conditional biasing** first maps the **conditioning representation** to a bias vector.

The bias vector is then added to the input.

Image: https://distill.pub/2018/feature-wise-transformations/

# Conditional scaling



conditioning representation → [linear] → [scaling vector]

**Conditional scaling** first maps the **conditioning representation** to a scaling vector.

input → ⊙ → output

The scaling vector is then multiplied with the input.

Image: https://distill.pub/2018/feature-wise-transformations/

# Conditional scaling

- Reminds you of…. Gating
  - Long-Short Term Memory units

- Gating allows you to learn which inputs are more related between e.g. the two sources

- All conditioning so far is on a feature level → efficient and effective → number of parameters to be learned scales linearly with the number of features of the NN

# Conditional scaling

- Can one do both conditional scaling and biasing?

Conditional Affine Transformation

The **FiLM generator** processes the conditioning information and produces parameters that describe how the target network should alter its computation.

Here, the **FiLM-ed network**'s computation is conditioned by two FiLM layers.

input

sub-network

conditioning

**FiLM generator**

FiLM parameters

FiLM

Information coming from e.g. the other source

sub-network

FiLM

sub-network

output

Image: https://distill.pub/2018/feature-wise-transformations/

E. Perez et al. „FiLM: Visual Reasoning with a General Conditioning Layer". AAAI 2018.  65

In a **fully-connected** network, FiLM applies a different affine transformation to each feature.

In a **convolutional** network, FiLM applies a different affine transformation to each channel, consistent across spatial locations.

First, each feature (or channel) is scaled by the corresponding γ parameter.

γ

Then, each feature (or channel) is shifted by the corresponding β parameter.

β

γ

β

Image: https://distill.pub/2018/feature-wise-transformations/

E. Perez et al. „FiLM: Visual Reasoning with a General Conditioning Layer". AAAI 2018.  66

# What can we do with conditioning?

- Visual Reasoning with Multi-hop Feature Modulation Strub et al. ECCV 2018.

- GuessWhat?! Visual object discovery through multi-modal dialogue. de Vries et al CVPR 2017
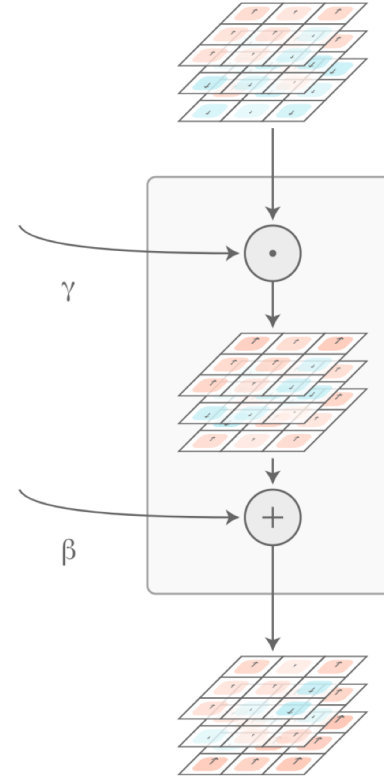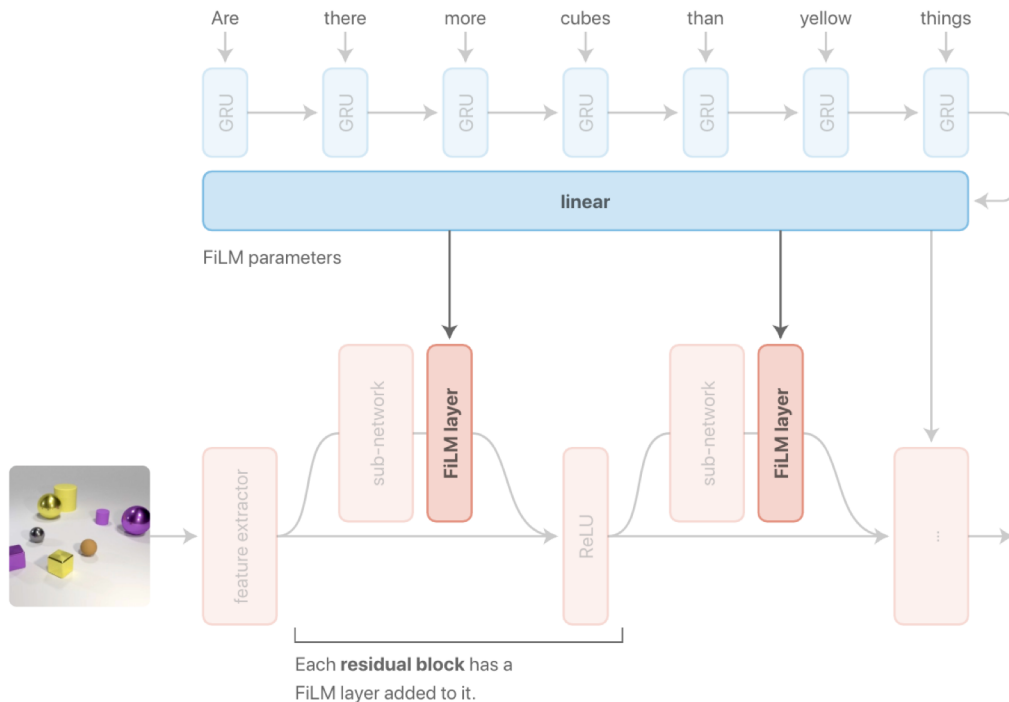
- A learned representation for artistic style. Dumoulin et al ICLR 2017

- Conditional image generation with PixelCNN decoders. van den Oord et al. NIPS 2016

# Visual Question and Answering



Are there more cubes than yellow things

GRU → GRU → GRU → GRU → GRU → GRU → GRU

**linear**

FiLM parameters
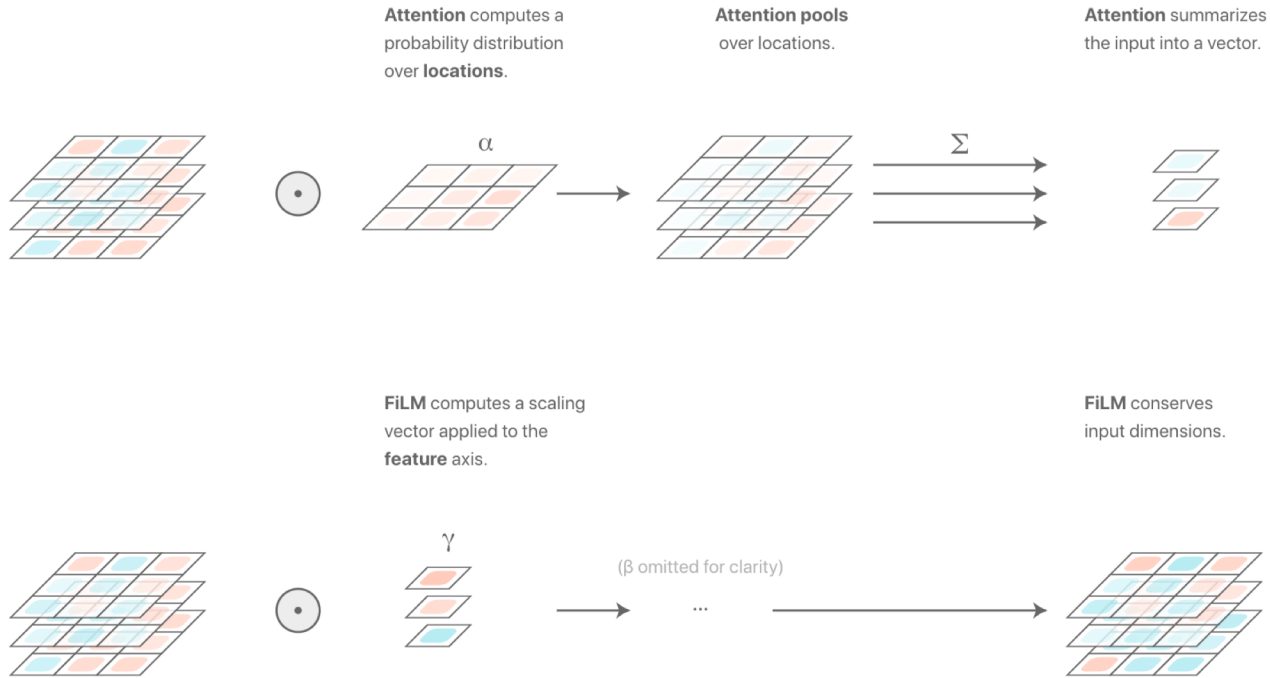
The **linguistic pipeline** acts as the FiLM generator.

FiLM layers in each residual block modulate the **visual pipeline**.

feature extractor → sub-network → **FiLM layer** → ReLU → sub-network → **FiLM layer** → ...

Each **residual block** has a FiLM layer added to it.

# Attention vs Conditioning



**Attention** computes a probability distribution over **locations**.

**Attention pools** over locations.

**Attention** summarizes the input into a vector.

$\alpha$

$\Sigma$

**FiLM** computes a scaling vector applied to the **feature** axis.

**FiLM** conserves input dimensions.

$\gamma$

(β omitted for clarity)

...

Image: https://distill.pub/2018/feature-wise-transformations/

# Attention vs Conditioning

- Attention: assumes that specific **locations** contain the most useful information

- Conditioning: assumes that specific **feature maps** contain the most useful information

Image: https://distill.pub/2018/feature-wise-transformations/

# Next lecture

- Next Monday: lecture on visualization