

Introduction to Transformer and GANs for NLP

2/2565: FRA501 Introduction to Natural Language Processing with Deep learning
Week 09

Paisit Khanarsa, Ph.D.

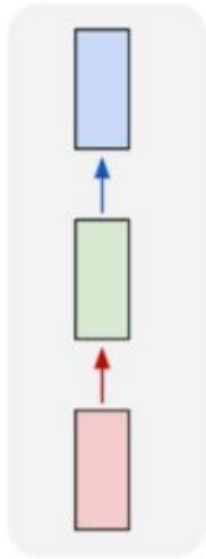
Institute of **Field Robotics** (FIBO), King Mongkut's University of Technology Thonburi

Outlines

- Recap of Attention Mechanism
- Transformer
- Beam Search
- GAN for NLP

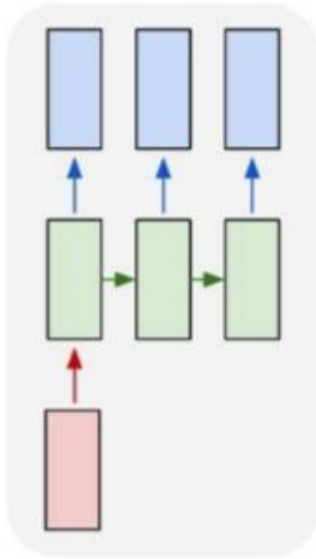
Recap of Attention Mechanism

one to one



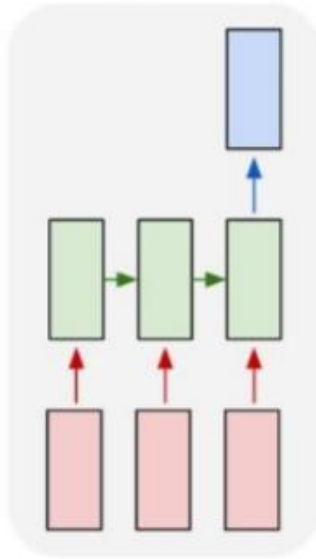
Fixed-sized input to fixed-sized output (e.g. image classification)

one to many



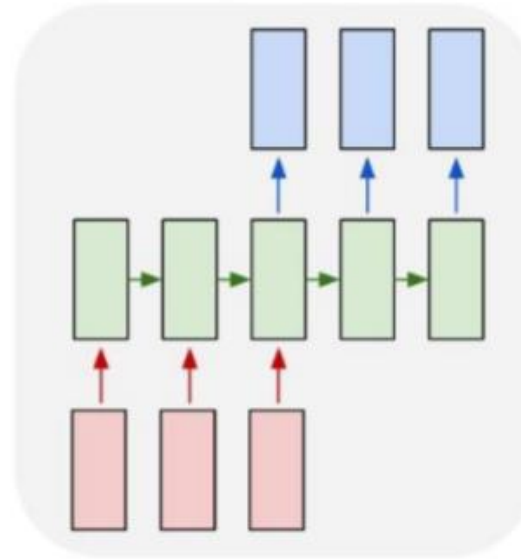
Sequence output (e.g. image captioning takes an image and outputs a sentence of words).

many to one



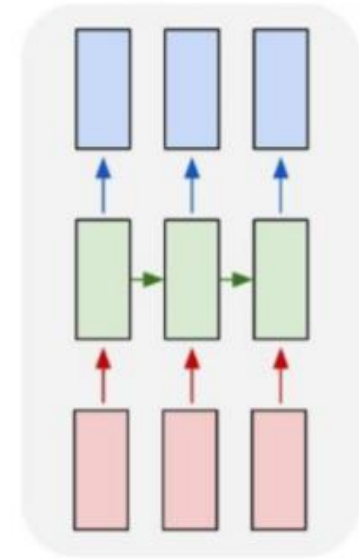
Sequence input (e.g. sentiment analysis where a given sentence is classified as expressing positive or negative sentiment).

many to many



Sequence input and sequence output (e.g. Machine Translation: an RNN reads a sentence in English and then outputs a sentence in French)

many to many

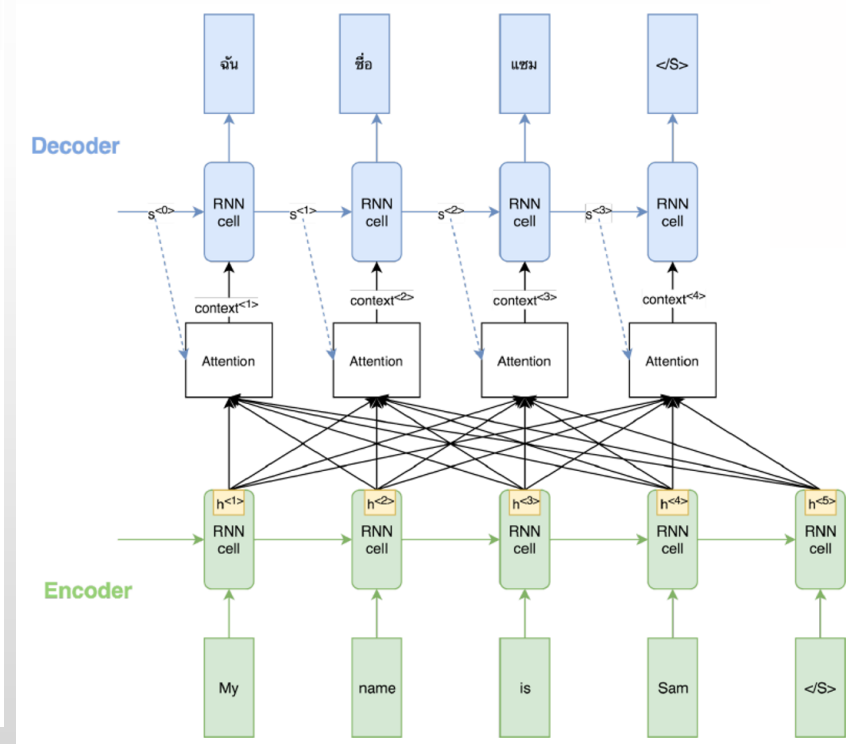
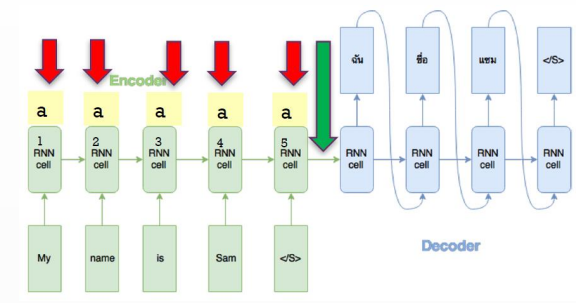
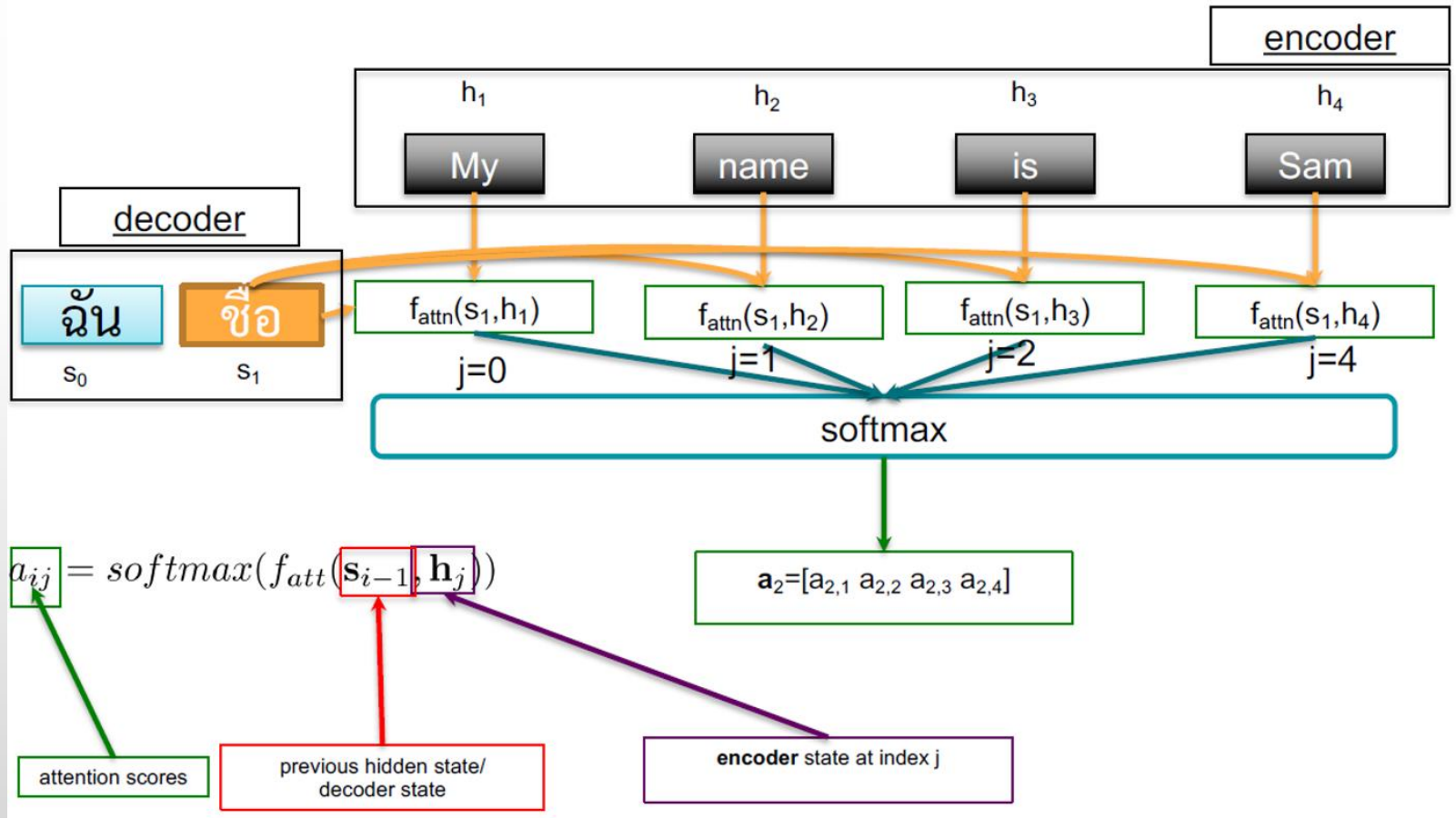


Synced sequence input and output (e.g. video classification where we wish to label each frame of the video)

i : decoder index
 j : encoder index

Attention Calculation: Attention Scores

- Example: Now we want to predict “ฉัน”



Type of Attention Mechanisms

	My	name	is	Sam	encoder
ฉัน					
ชื่อ					
แซน					
decoder					

$$a_{ij} = \text{softmax}(f_{att}(s_{i-1}, h_j))$$

- **Additive attention:** The original attention mechanism (Bahdanau et al., 2015) uses a one-hidden layer feed-forward network to calculate the attention alignment:

$$f_{att}(s_{i-1}, h_j) = \tanh(W_a[s_{i-1}; h_j])$$

- **Multiplicative attention:** Multiplicative attention (Luong et al., 2015) simplifies the attention operation by calculating the following function:

$$f_{att}(s_{i-1}, h_j) = S_{i-1}^T W_a h_j$$

- **Self-attention:** Without any additional information, however, we can still extract relevant aspects from the sentence by allowing it to attend to itself using self-attention (Lin et al., 2017)

$$a = \text{softmax}(W_{s_2} \tanh(W_{s_1} H^T))$$

- **Key-value attention:** key-value attention (Daniluk et al., 2017) is a recent attention variant that separates form from function by keeping separate vectors for the attention calculation.

Self Attention

	My	name	is	Sam	encoder
ฉัน					
ชื่อ					
แซม					
decoder					

$$a_{ij} = \text{softmax}(f_{\text{att}}(s_{i-1}, h_j))$$

- Without any additional information, we can still extract relevant aspects from the sentence by allowing it to attend to itself using self-attention (Lin et al., 2017)

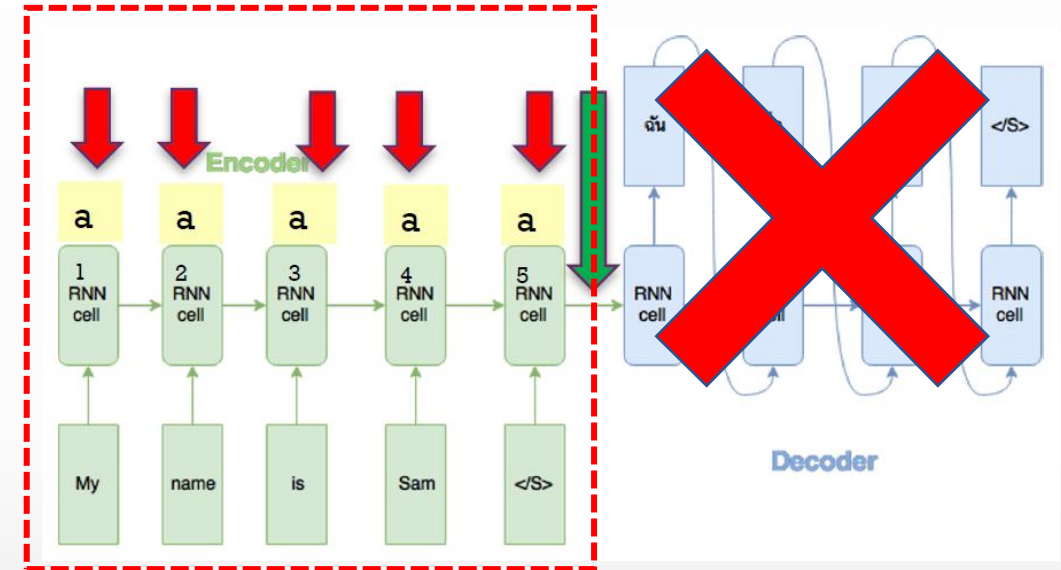
$$H = (h_1, h_2, \dots, h_n)$$

Fully connected layer

$$a = \text{softmax}(W_{s_2} \tanh(W_{s_1} H^T))$$

One-hidden layer (Dense)

- w_{s_1} is a weight matrix, w_{s_2} is a vector of parameters. Note that these parameters are tuned by the neural networks.
- The objective is to improve a quality of embedding vector by adding context information.



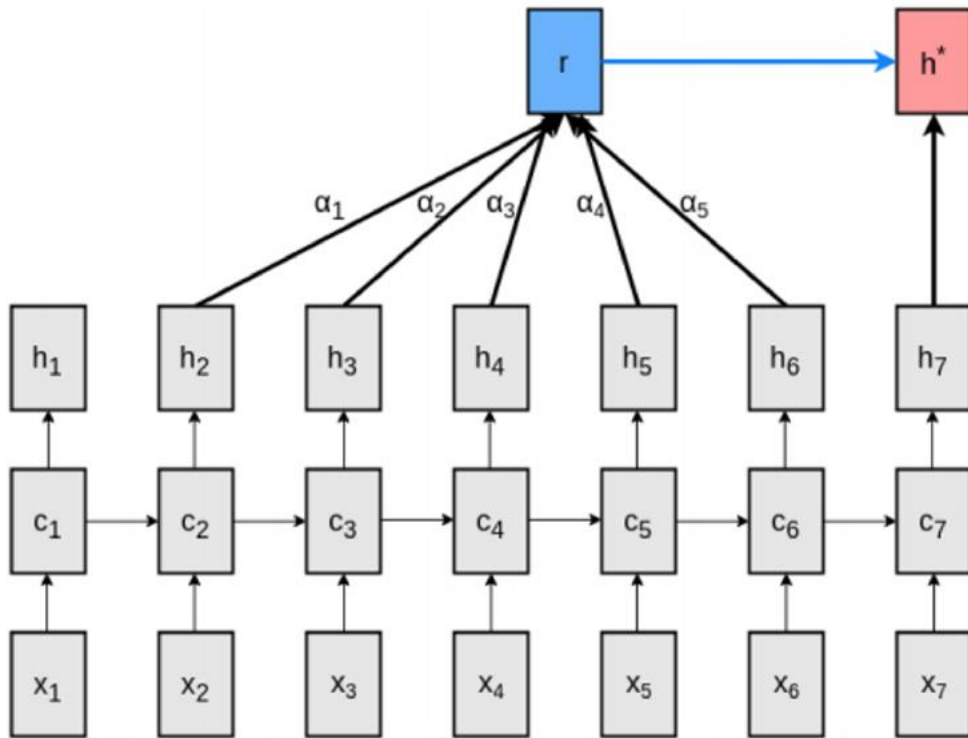
Additive Attention

$$a_{ij} = \text{softmax}(\tanh(W_a[s_{i-1}; h_j]))$$

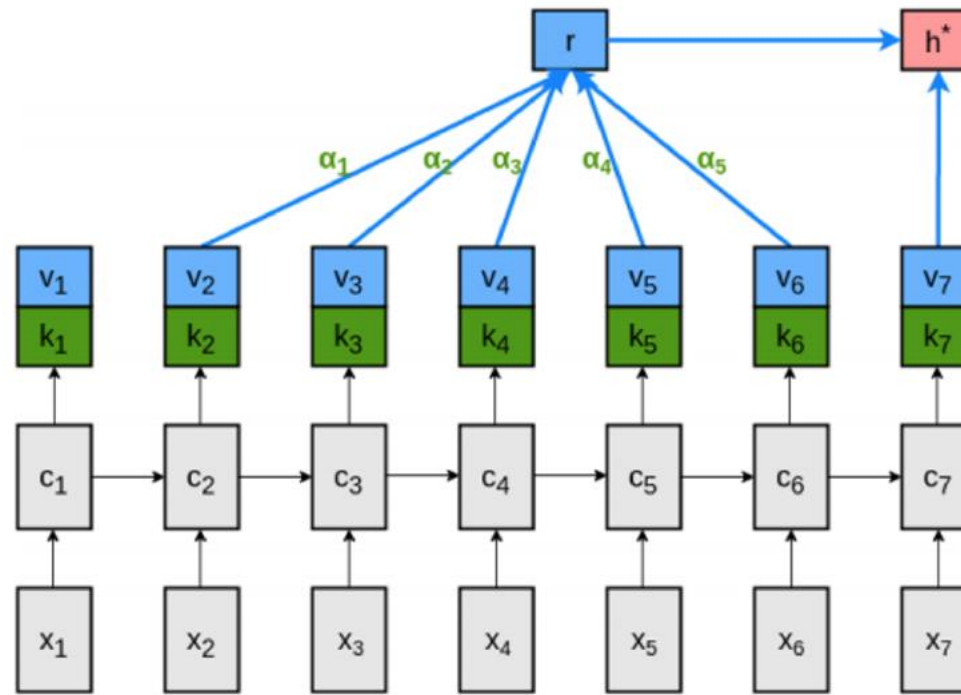
Key-value attention

$$a_{ij} = \text{softmax}(f_{\text{att}}(s_{i-1}, h_j))$$

$$c_i = \sum_j a_{ij} h_j$$



(a) Neural language model with attention.



(b) Key-value separation.

Value=encoded vector

Key=used for attention score calculation

Daniluk, M., Rocktäschel, T., Welbl, J., & Riedel, S. (2017). Frustratingly short attention spans in neural language modeling. arXiv preprint arXiv:1702.04521.

Self Attention

	My	name	is	Sam	encoder
ฉัน					
ชื่อ					
แซม					
decoder					

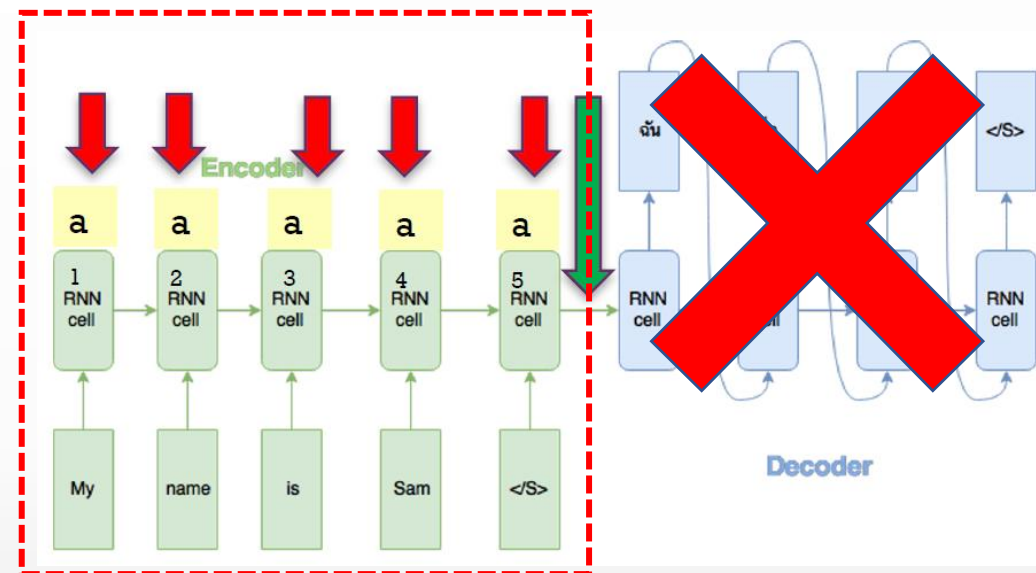
$$a_{ij} = \text{softmax}(f_{\text{att}}(s_{i-1}, h_j))$$

- Without any additional information, we can still extract relevant aspects from the sentence by allowing it to attend to itself using self-attention (Lin et al., 2017)

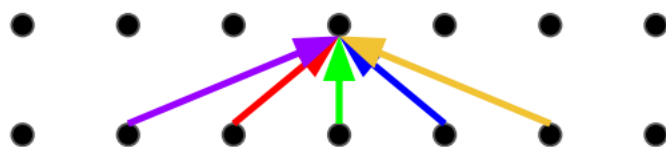
$$H = (h_1, h_2, \dots, h_n)$$

Fully connected layer

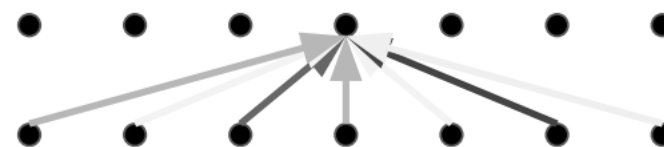
$$a = \text{softmax}(W_a \tanh(W_h H^T))$$



Convolution



Self-Attention

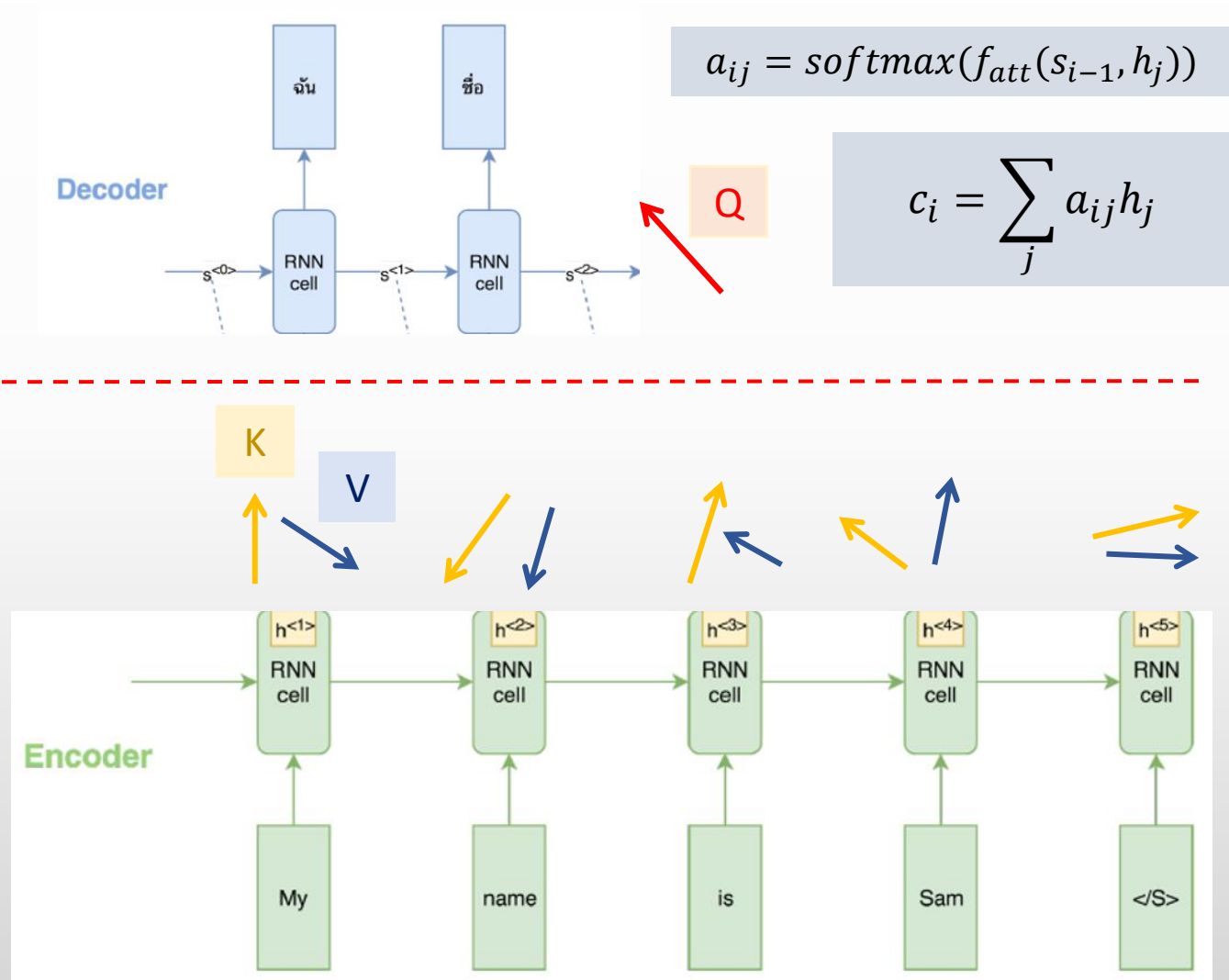
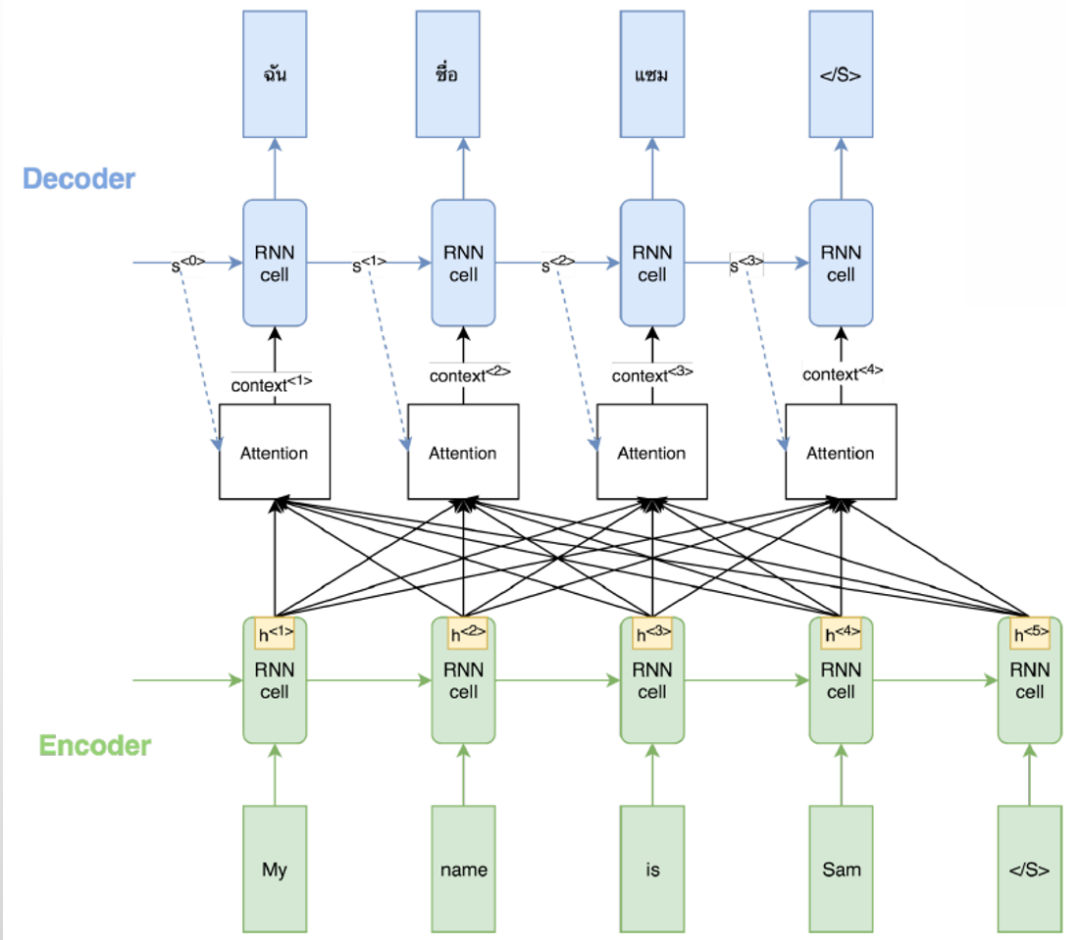


attention

$$W_a[s_{i-1}; h_j])$$

- w_{s_1} is that the
- The of adding

Query, Key and Value in Attention mechanism



Transformer

- Attention is all you need

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.

<https://arxiv.org/pdf/1706.03762.pdf>

Model Architecture of Transformer

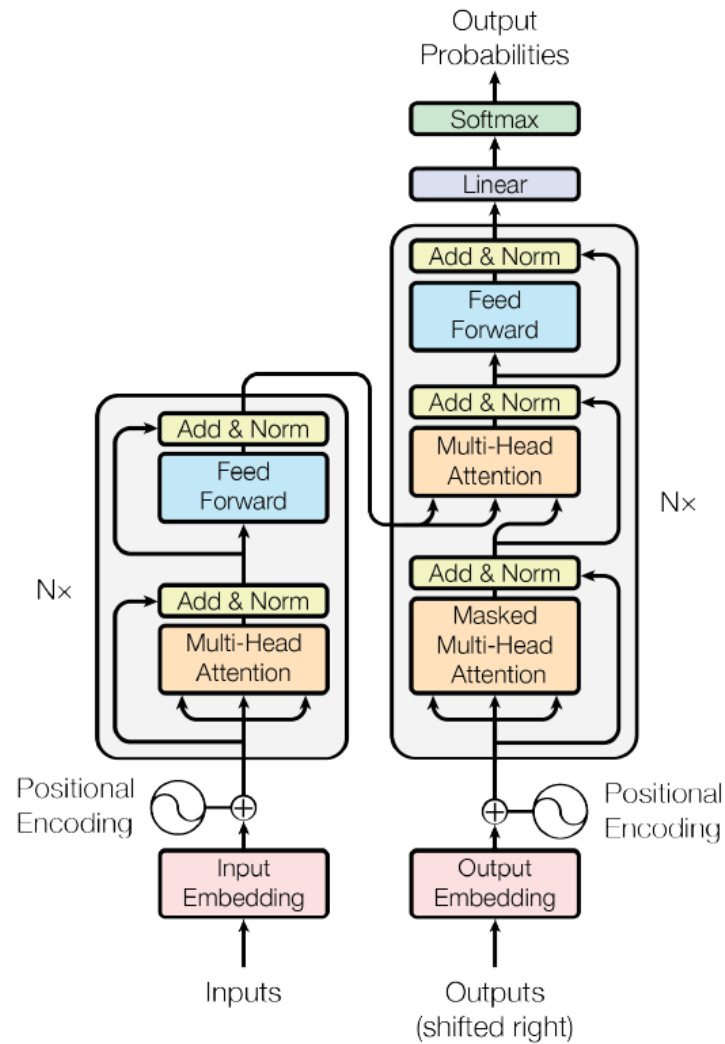
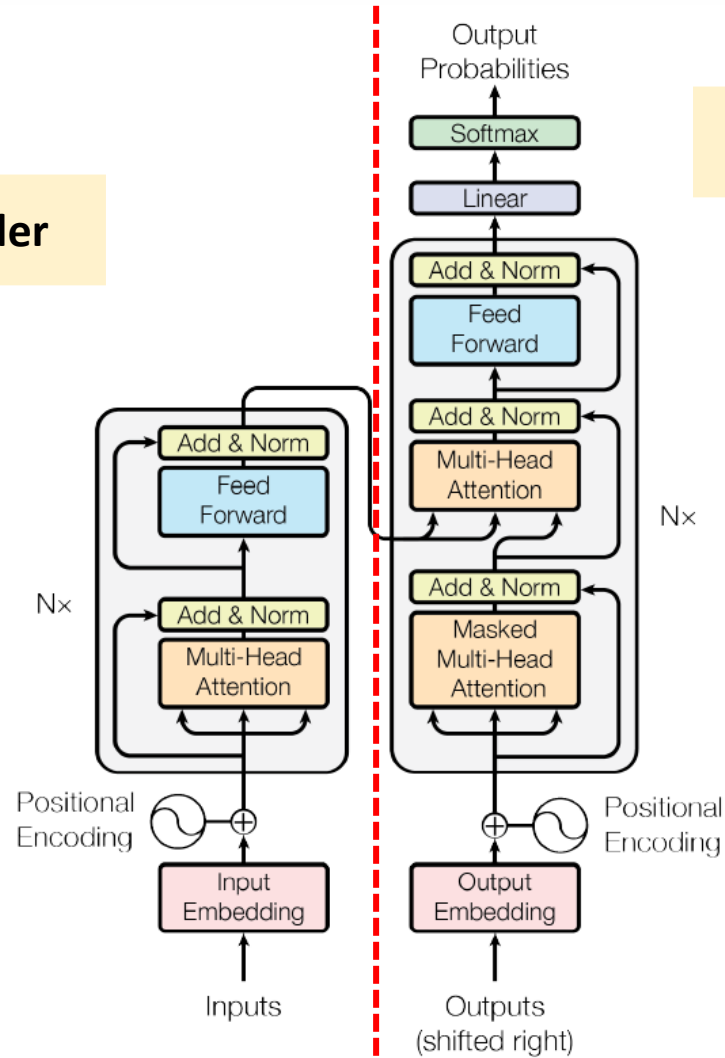


Figure 1: The Transformer - model architecture.

Model Architecture of Transformer

Encoder



Decoder

Figure 1: The Transformer - model architecture.

Model Architecture of Transformer

Example:

The black cat is on the black chair

Encoder

Decoder

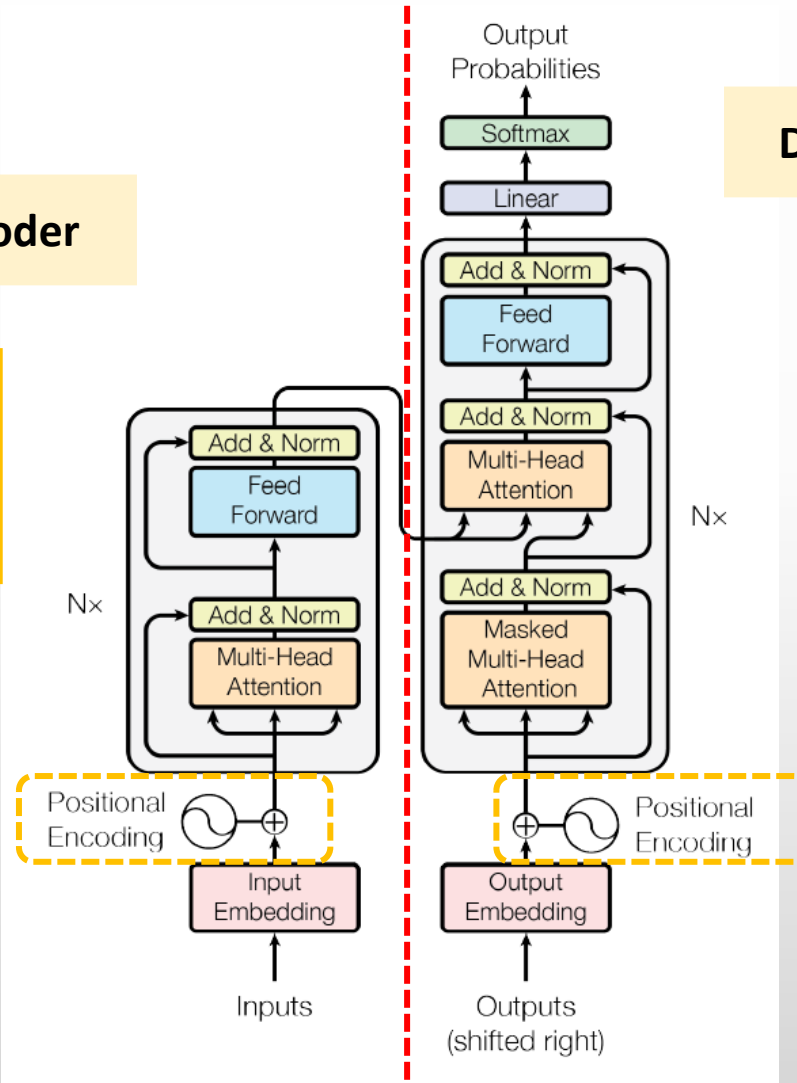


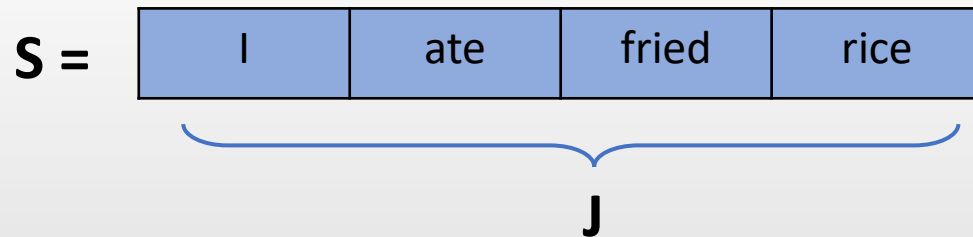
Figure 1: The Transformer - model architecture.

Recap: Positional Encoding

j =positional index
 d =dimensional index

- Positional Encoding (PE)
 - Position is modeled by a multiplicative term on each word vector with weights depending on the position in the sentence

$$f(j, d) = (1 - j/J) - (d/D)(1 - 2j/J)$$



$PE(I) =$

	1	0.13	○	0.70	=	0.09
	2	0.04	○	0.65		0.03
	○
	○
	D-1	0.91	○	0.30		0.27
	D	0.82	○	0.25		0.20

○ is element-wise multiplication

Model Architecture of Transformer

Encoder

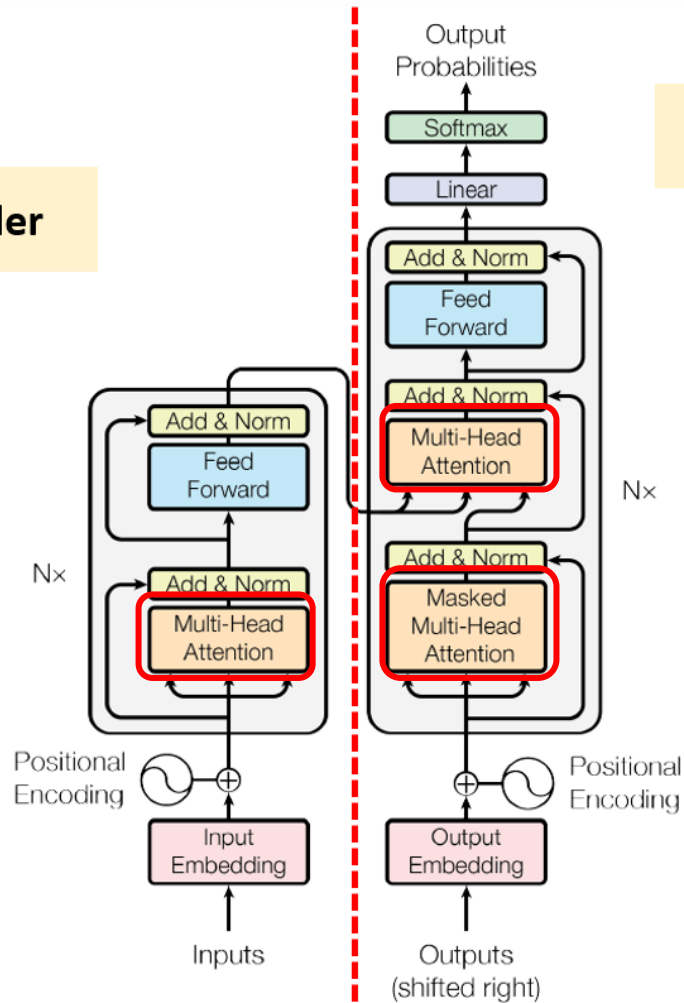
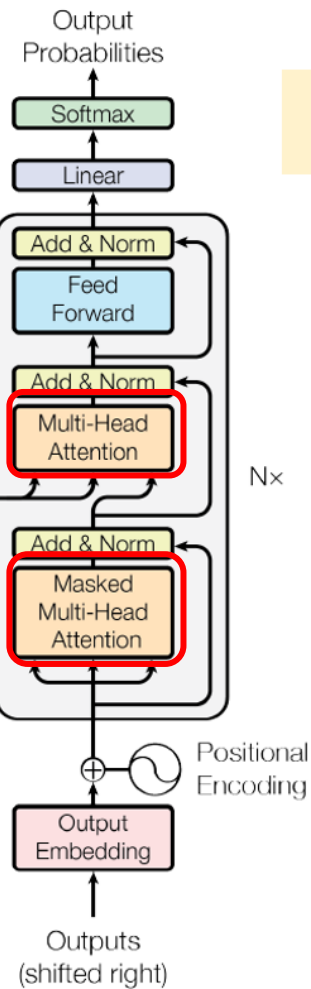
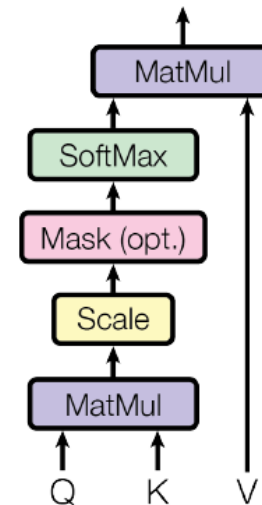


Figure 1: The Transformer - model architecture.

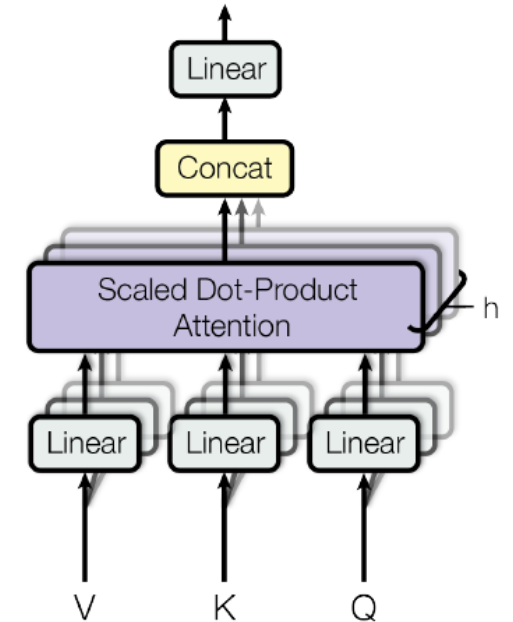
Decoder



Scaled Dot-Product Attention

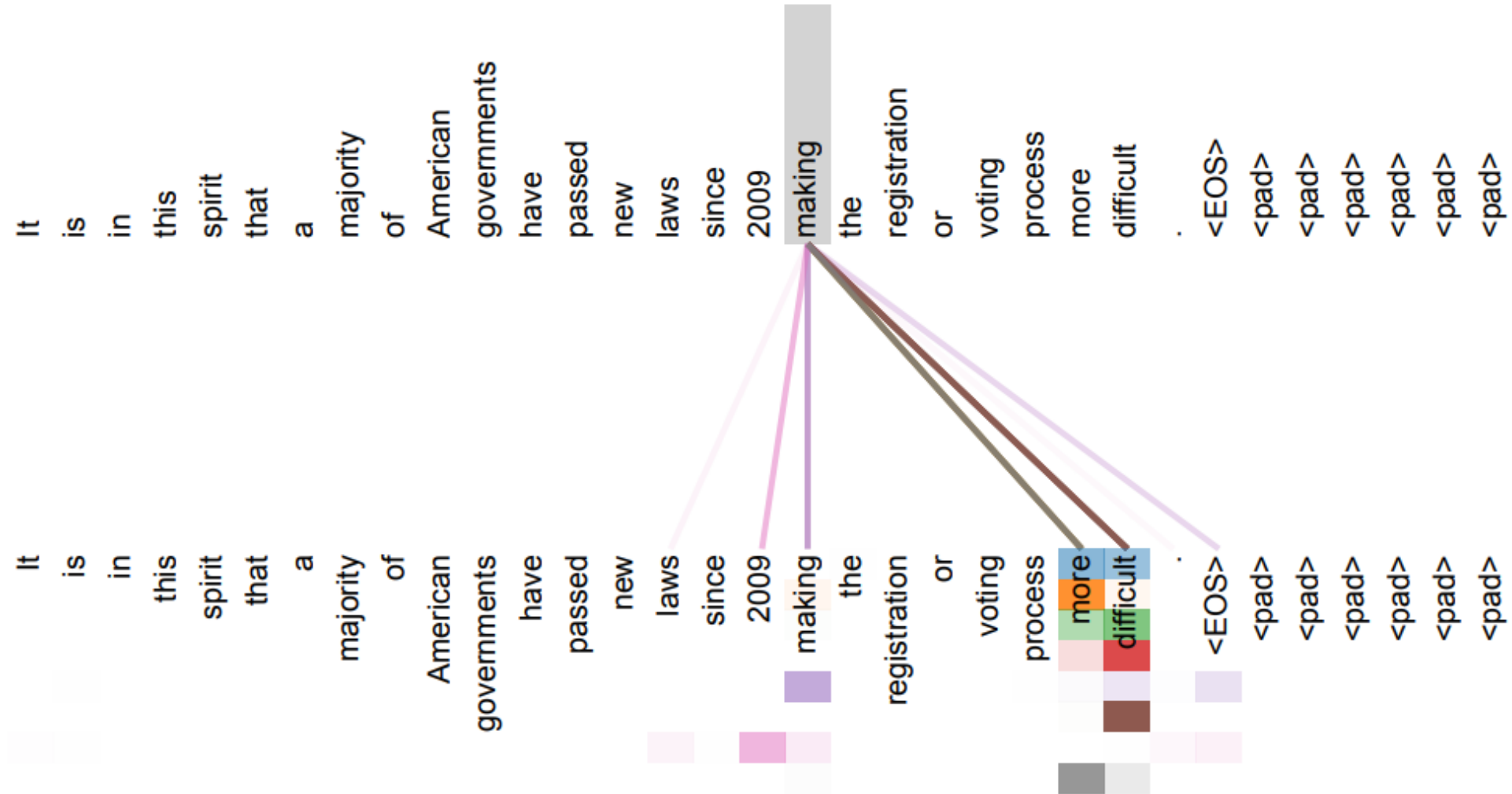


Multi-Head Attention

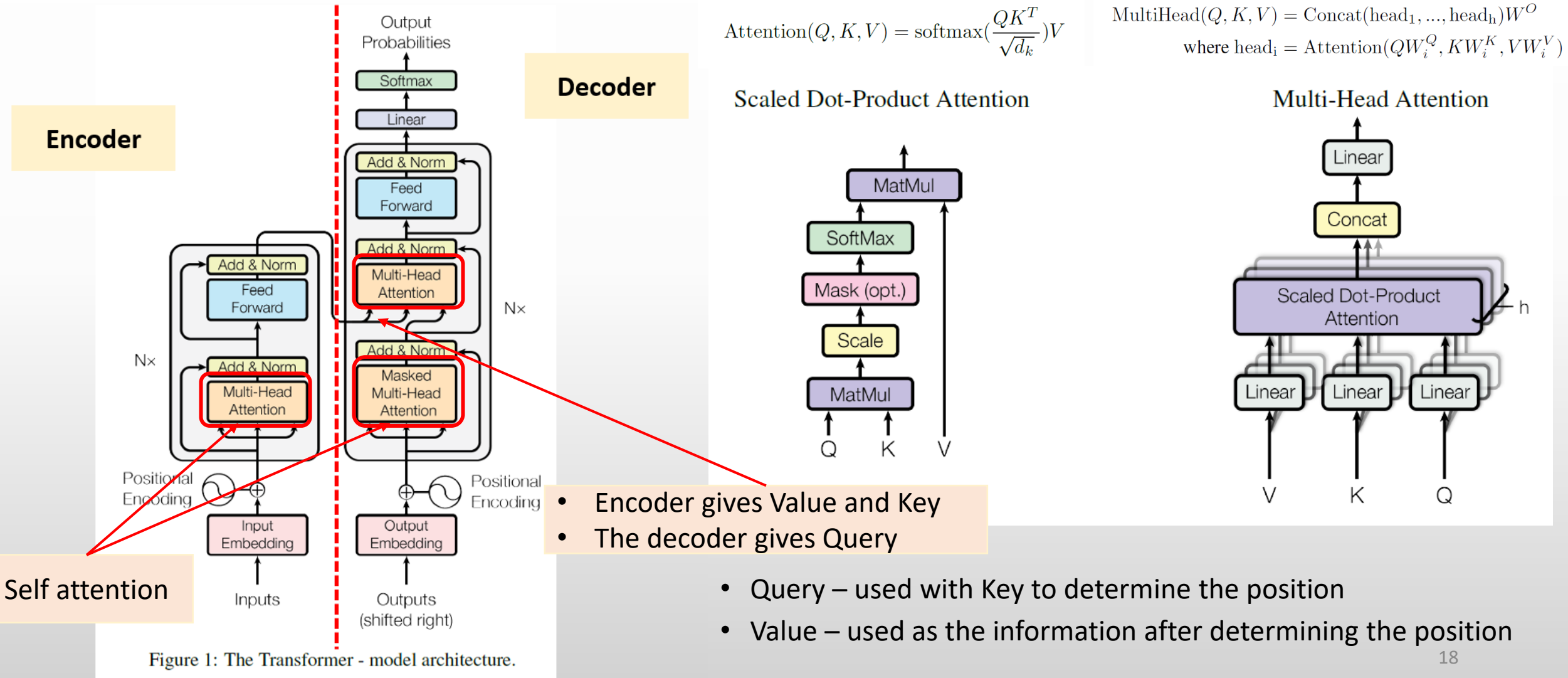


- Query – used with Key to determine the position
- Value – used as the information after determining the position

Attention Visualizations



Model Architecture of Transformer



- Encoder gives Value and Key
- The decoder gives Query

- Query – used with Key to determine the position
- Value – used as the information after determining the position

Model Architecture of Transformer

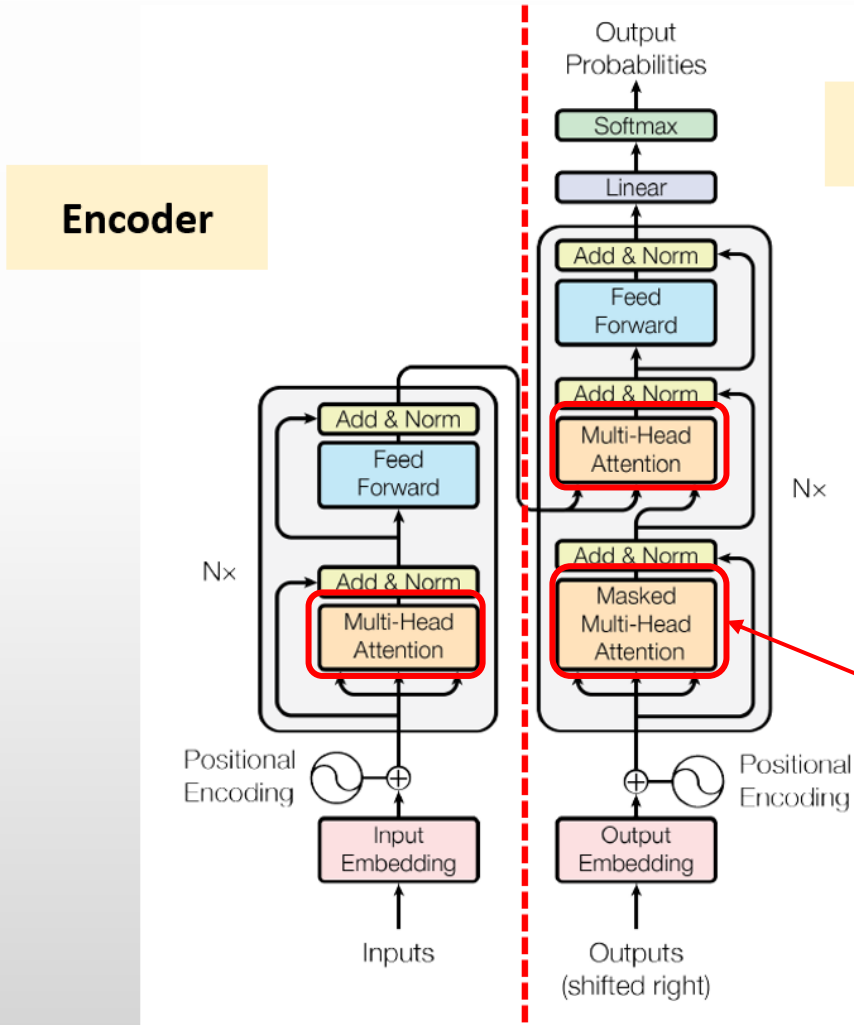


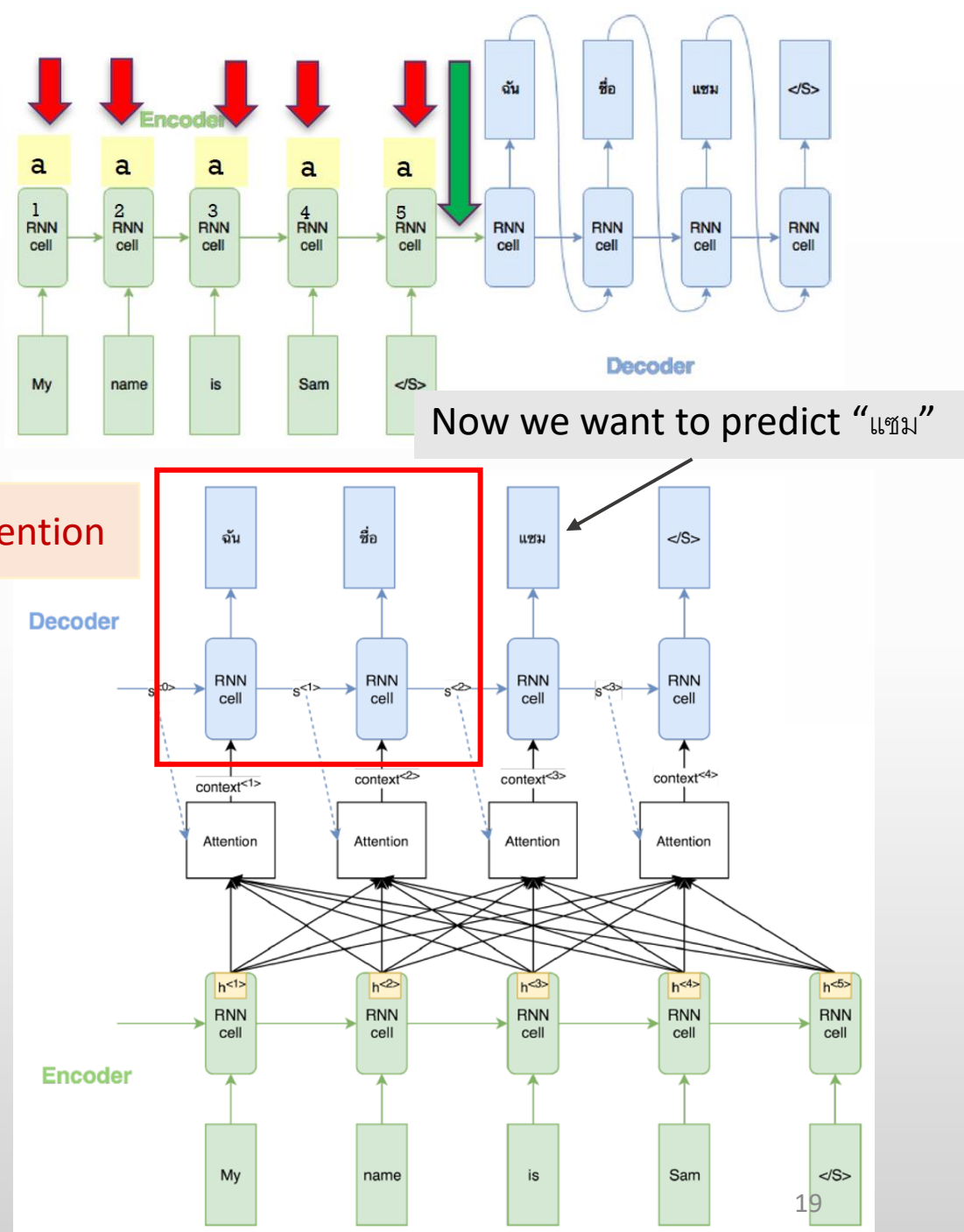
Figure 1: The Transformer - model architecture.

Decoder

Encoder

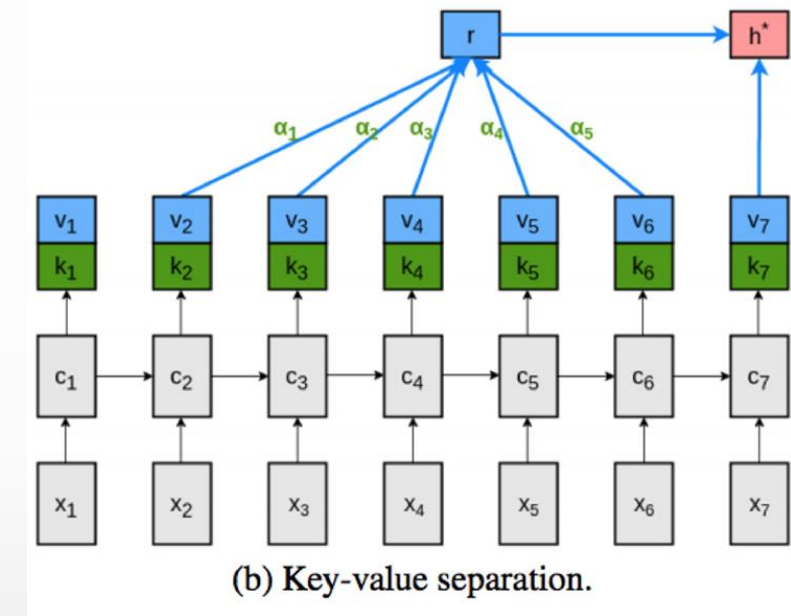
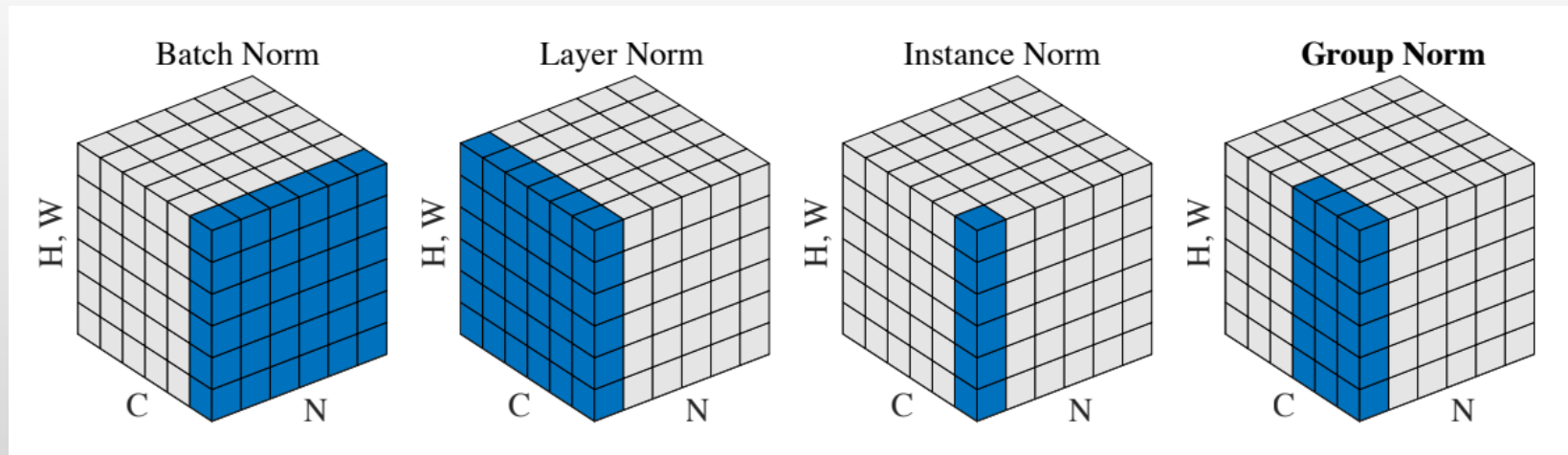
A mask prevents looking ahead (cheating)

Self attention



Model Architecture of Transformer

- **Layer norm**
- Normalize the mean and SD
- BN and GN are usually best, GN is better when batch size is small (Vision task)
- N – example in mini batch
- C – Channel output
- H,W – spatial coordinates (x,y)



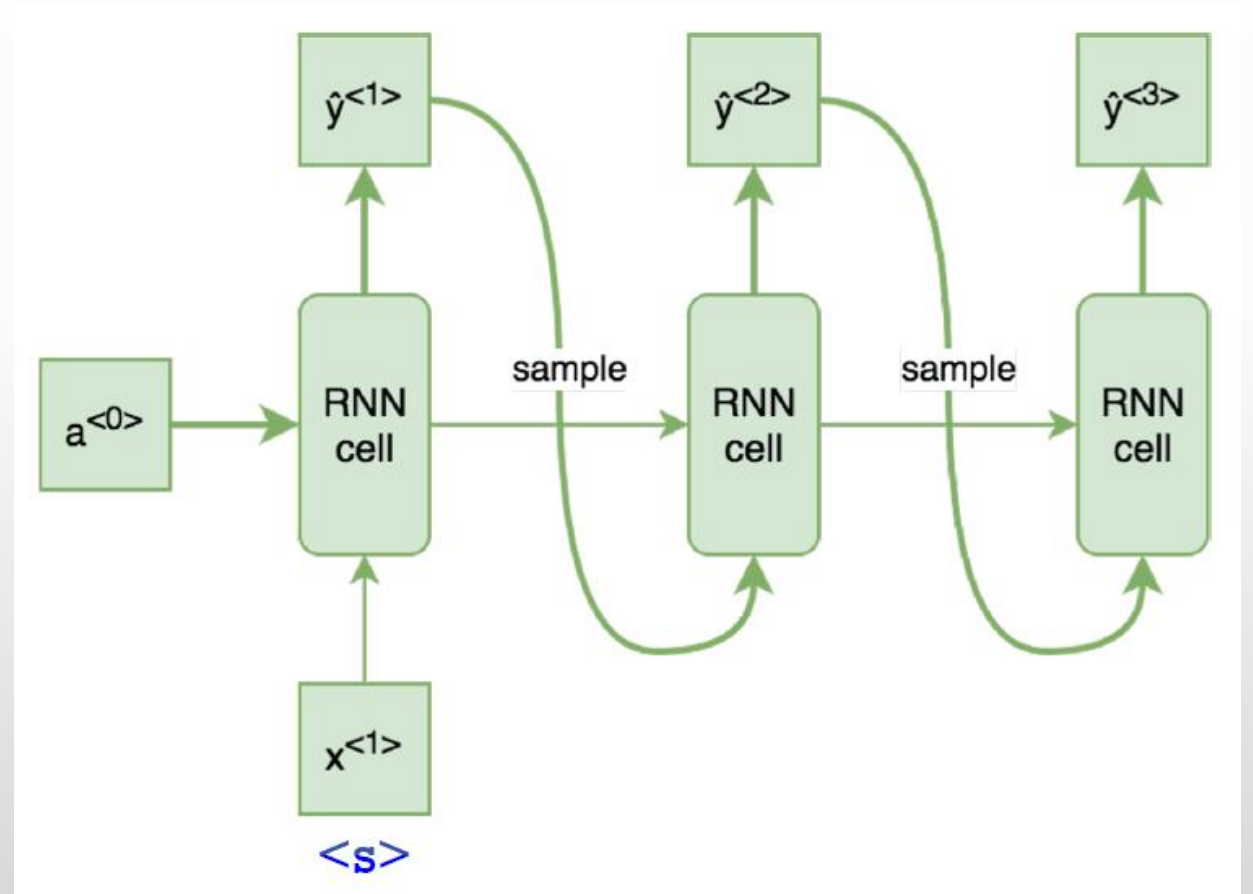
MT results with Transformer

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

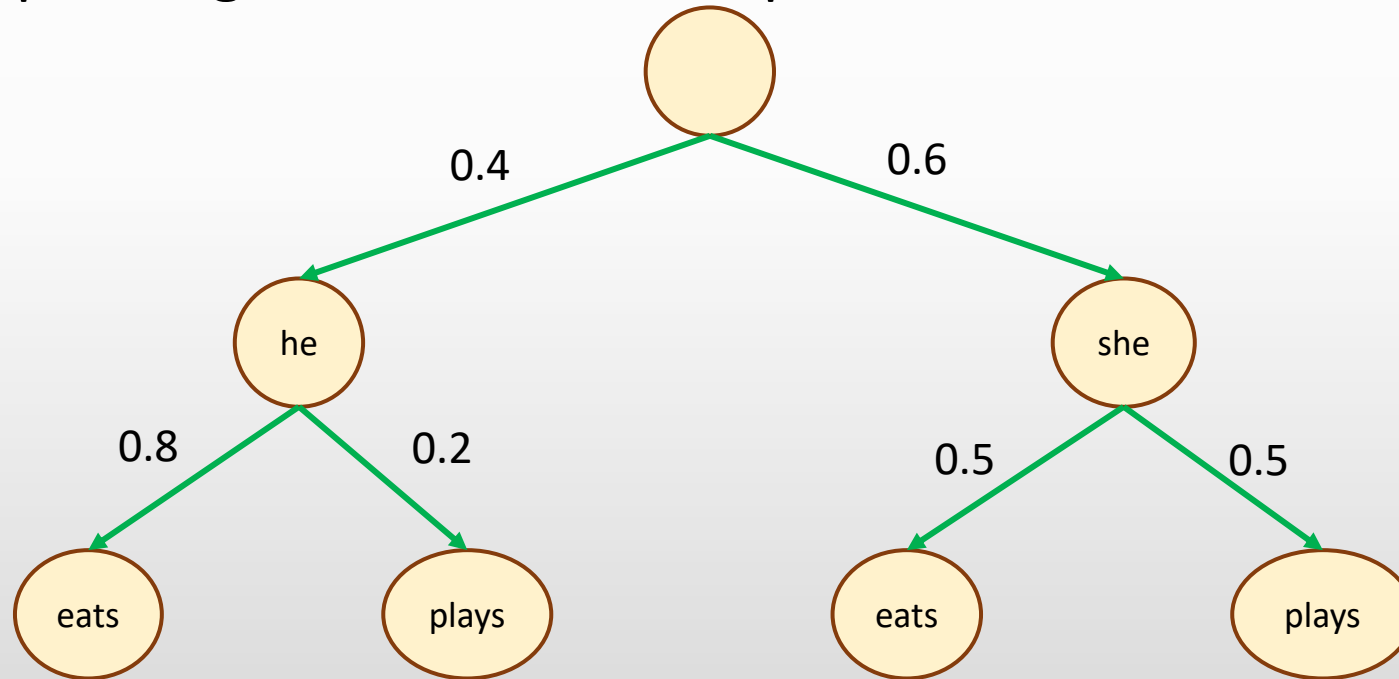
Text generation model (inference; testing)

- To generate a novel sequence, the inference model (testing phase) randomly samples an output from a softmax distribution.



Beam search

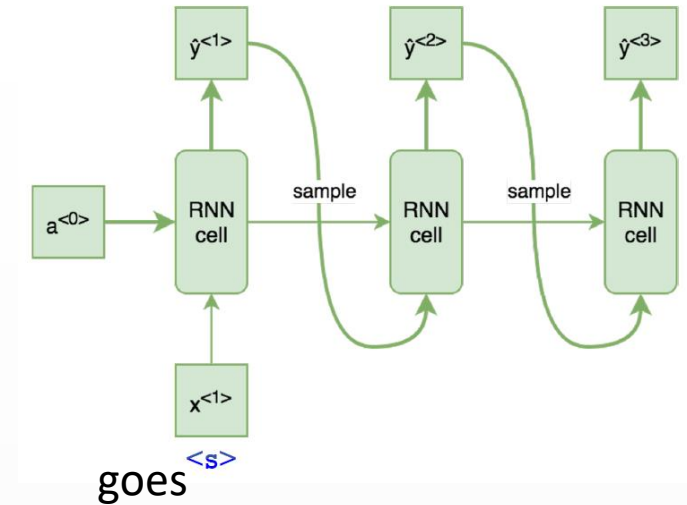
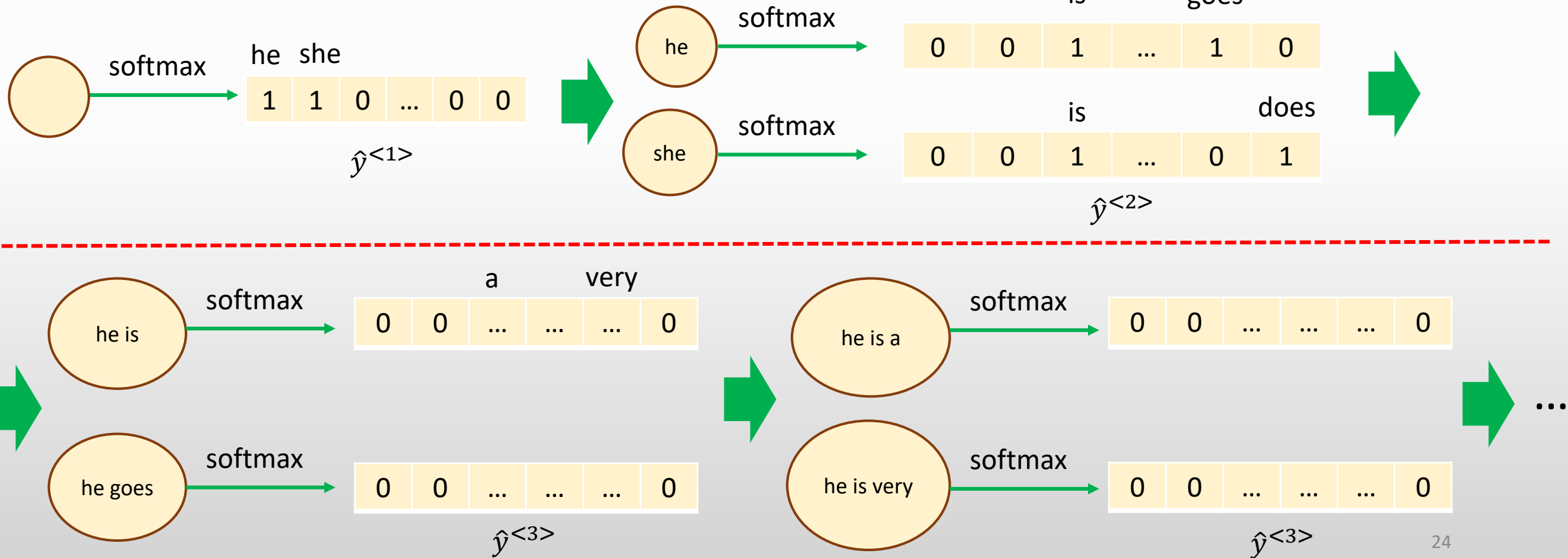
- Autoregressive requires making a decision even though it might not be optimal for the whole sequence
- However, expanding the tree will be exponential



- $P(\text{he eats}) = 0.32$
- $P(\text{she eats}) = 0.30$

Beam search

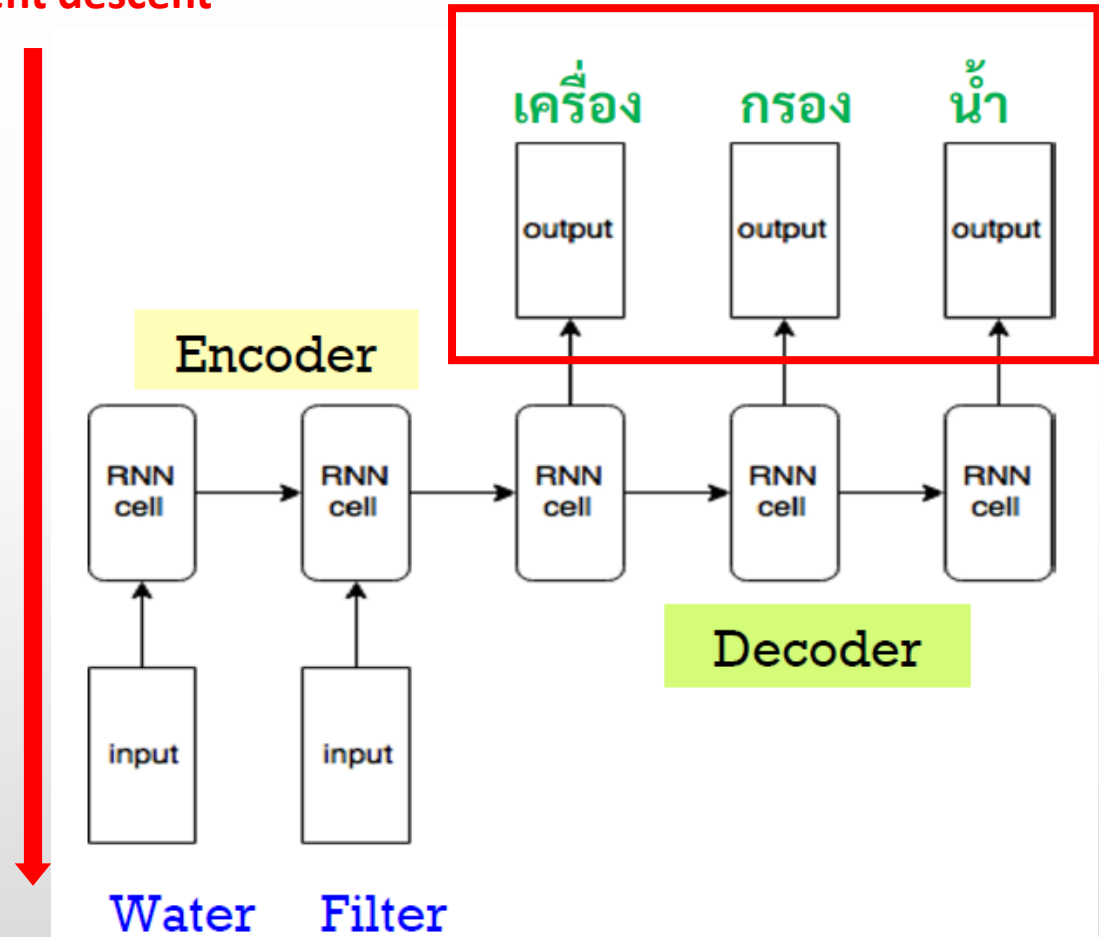
- Keep K active states at each step ; K = 2



RNN architectures in NPL

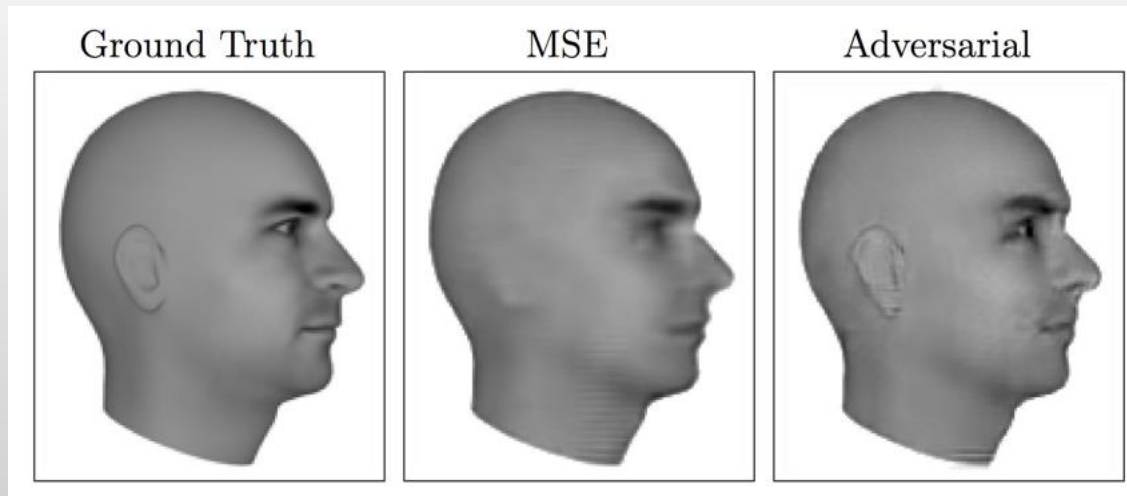
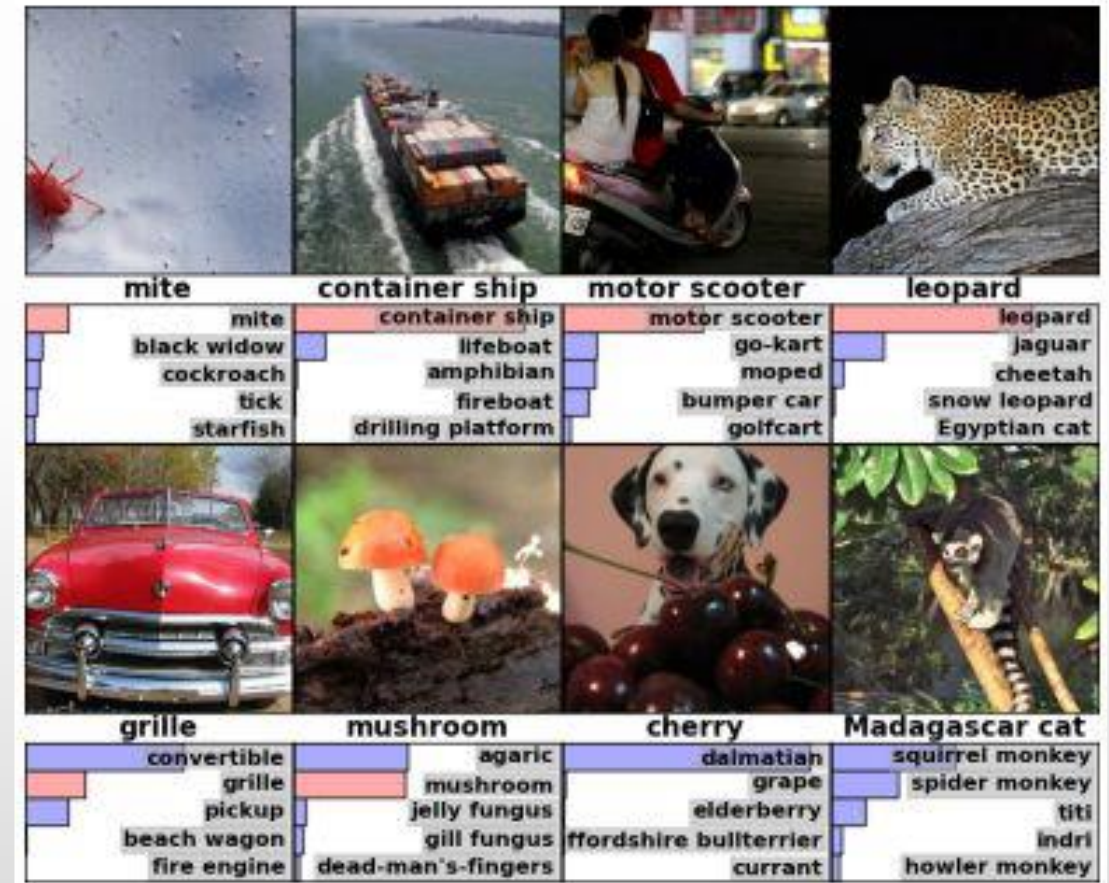
- **Many-to-many (encoder-decoder)**
- Sequence Input, Sequence output
- These two sequences can be of different length
- E.g. Machine Translation
 - Input: English Sentence
 - Output: Thai Sentence
- Machine Translation is also a text generation task

Gradient descent

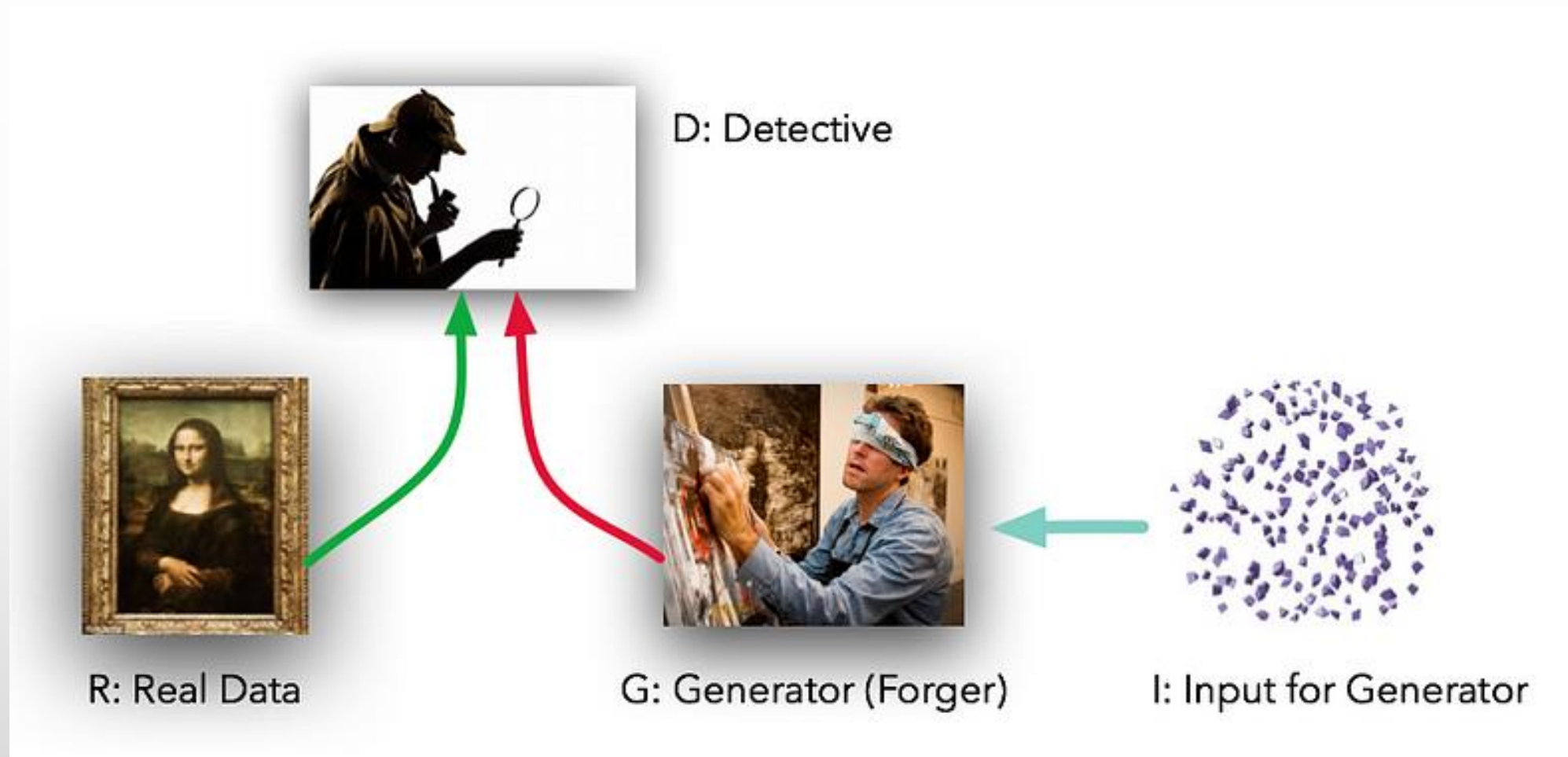


Generative Adversarial Networks (GANs)

- Learning distributions
- Supervised learning tasks usually have one correct answer
- Sometimes there are more than one possibility
 - What is the next frame of a video?
 - What is the missing pixels in an image?
 - What word is missing from the blank?
 - I eat ____



Generative Adversarial Networks (GANs)



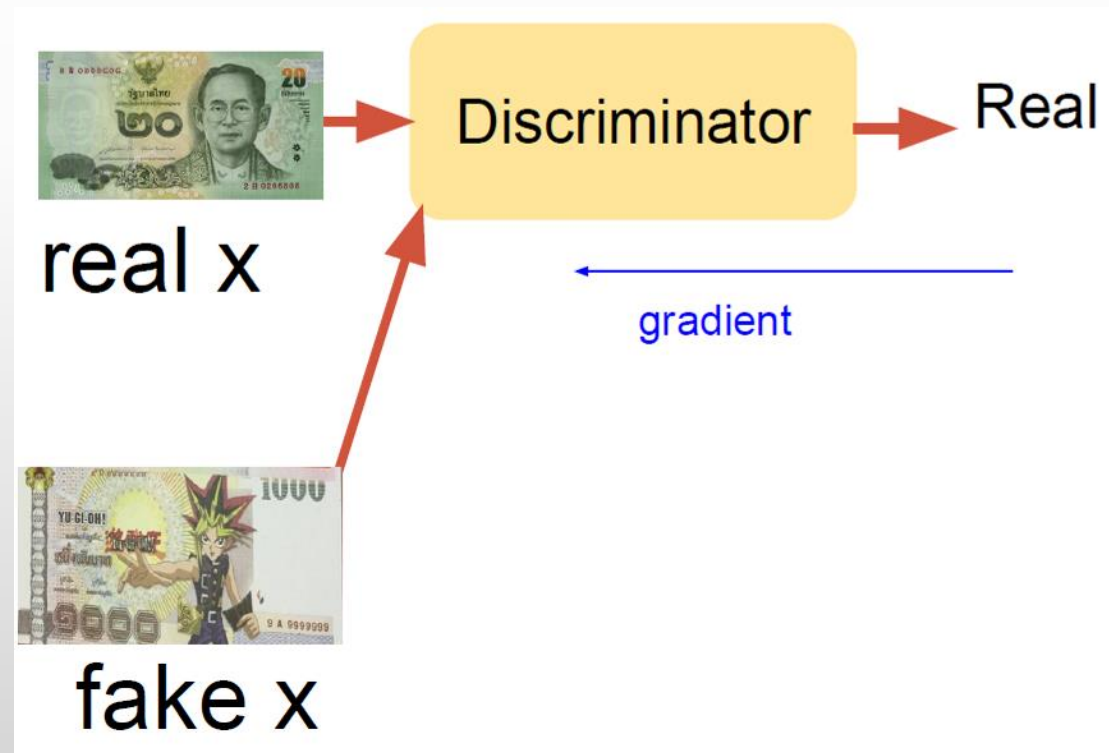
Generative Adversarial Networks (GANs)

- Consider a money counterfeiter
 - He wants to make fake money that looks real
 - There's a police that tries to differentiate fake and real money.
- The counterfeiter is the adversary and is generating fake inputs. – Generator network
- The police is try to discriminate between fake and real inputs. – Discriminator network



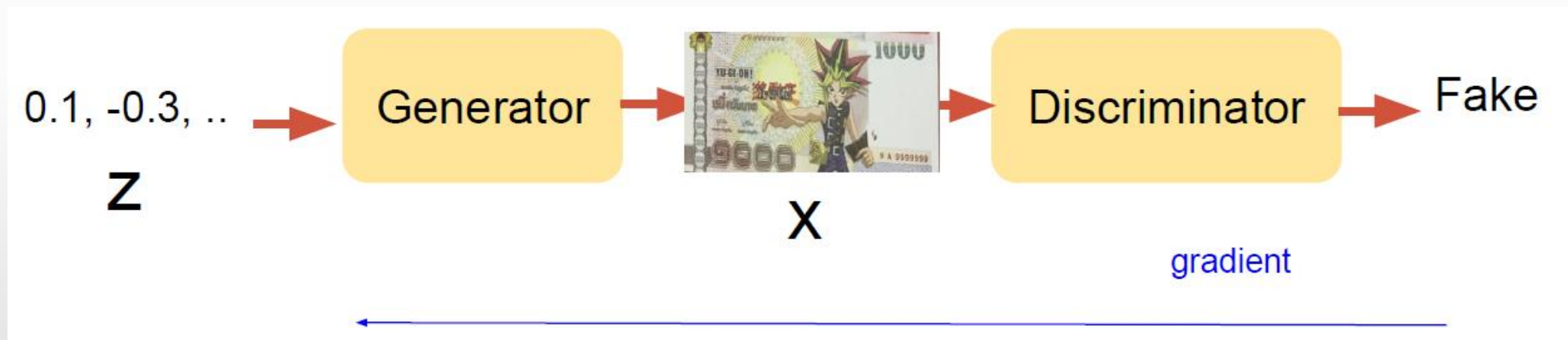
Generative Adversarial Networks (GANs)

- Train Discriminator
 - The discriminator learns to differentiate real and fake data
 - Learns by backpropagation



Generative Adversarial Networks (GANs)

- Train Generator
 - The generator learns to be better by the gradient given by the discriminator



Generative Adversarial Networks (GANs)

- Generator (Money Faker):

- Maximize \mathcal{Y}

$$\min_G \mathbb{E}_{\mathbf{z}} [\log(1 - D(G(\mathbf{z})))]$$

$$\mathcal{L}_G = -\frac{1}{m} \sum_{i=1}^m \log(D(G(z^{(i)})))$$

- Discriminator (Police):

- For real images => Maximize \mathcal{Y}
- For generated images from the faker => Minimize \mathcal{Y}

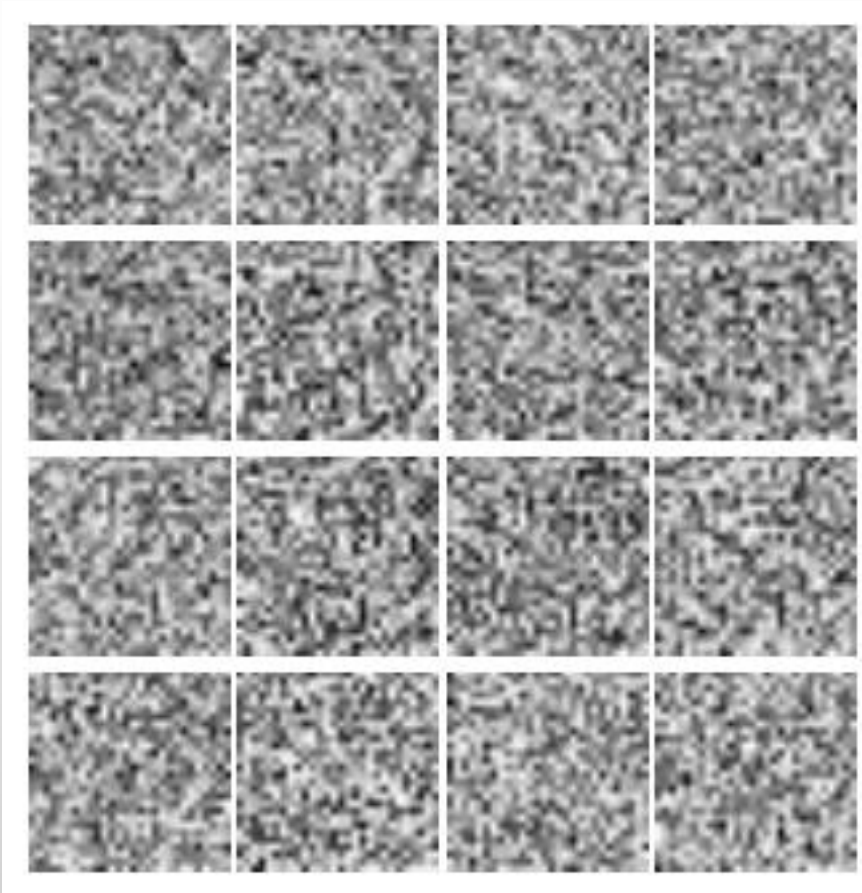
$$\max_D \mathbb{E}_{\mathbf{z}} [\log(1 - D(G(\mathbf{z})))] + \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(D(\mathbf{x}))]$$

$$\mathcal{L}_D = -\frac{1}{m} \sum_{i=1}^m (\log(D(x^{(i)})) + \log(1 - D(G(z^{(i)}))))$$



Generative Adversarial Networks (GANs)

- Generator output starts from random noise and gets better as we train.



GANs Loss Formulations

Discriminator

$$\max_D \mathbb{E}_{\mathbf{z}} [\log(1 - D(G(\mathbf{z}))) + \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(D(\mathbf{x}))]]$$

Generator

$$\min_G \mathbb{E}_{\mathbf{z}} [\log(1 - D(G(\mathbf{z})))]$$

GAN	DISCRIMINATOR LOSS	GENERATOR LOSS
MM GAN	$\mathcal{L}_D^{\text{GAN}} = -\mathbb{E}_{x \sim p_d} [\log(D(x))] - \mathbb{E}_{\hat{x} \sim p_g} [\log(1 - D(\hat{x}))]$	$\mathcal{L}_G^{\text{GAN}} = \mathbb{E}_{\hat{x} \sim p_g} [\log(1 - D(\hat{x}))]$
NS GAN	$\mathcal{L}_D^{\text{NSGAN}} = -\mathbb{E}_{x \sim p_d} [\log(D(x))] - \mathbb{E}_{\hat{x} \sim p_g} [\log(1 - D(\hat{x}))]$	$\mathcal{L}_G^{\text{NSGAN}} = -\mathbb{E}_{\hat{x} \sim p_g} [\log(D(\hat{x}))]$
WGAN	$\mathcal{L}_D^{\text{WGAN}} = -\mathbb{E}_{x \sim p_d} [D(x)] + \mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})]$	$\mathcal{L}_G^{\text{WGAN}} = -\mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})]$
WGAN GP	$\mathcal{L}_D^{\text{WGANGP}} = \mathcal{L}_D^{\text{WGAN}} + \lambda \mathbb{E}_{\hat{x} \sim p_g} [(\nabla D(\alpha x + (1 - \alpha \hat{x})) _2 - 1)^2]$	$\mathcal{L}_G^{\text{WGANGP}} = -\mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})]$
LS GAN	$\mathcal{L}_D^{\text{LSGAN}} = -\mathbb{E}_{x \sim p_d} [(D(x) - 1)^2] + \mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})^2]$	$\mathcal{L}_G^{\text{LSGAN}} = -\mathbb{E}_{\hat{x} \sim p_g} [(D(\hat{x}) - 1)^2]$
DRAGAN	$\mathcal{L}_D^{\text{DRAGAN}} = \mathcal{L}_D^{\text{GAN}} + \lambda \mathbb{E}_{\hat{x} \sim p_d + \mathcal{N}(0, c)} [(\nabla D(\hat{x}) _2 - 1)^2]$	$\mathcal{L}_G^{\text{DRAGAN}} = \mathbb{E}_{\hat{x} \sim p_g} [\log(1 - D(\hat{x}))]$
BEGAN	$\mathcal{L}_D^{\text{BEGAN}} = \mathbb{E}_{x \sim p_d} [x - \text{AE}(x) _1] - k_t \mathbb{E}_{\hat{x} \sim p_g} [\hat{x} - \text{AE}(\hat{x}) _1]$	$\mathcal{L}_G^{\text{BEGAN}} = \mathbb{E}_{\hat{x} \sim p_g} [\hat{x} - \text{AE}(\hat{x}) _1]$

GANs Loss Formulations

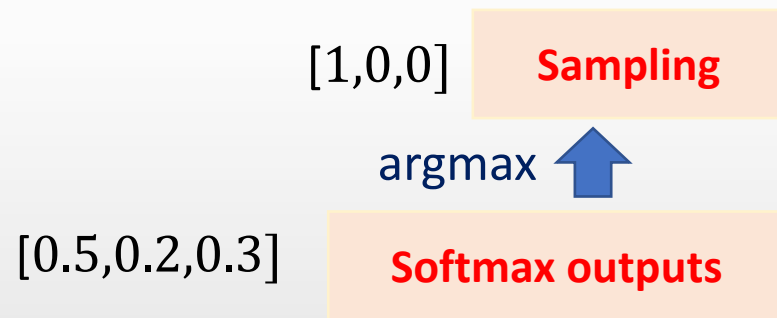
DCGAN	LSGAN	WGAN (clipping)	WGAN-GP (ours)
Baseline (G : DCGAN, D : DCGAN)			
			
G : No BN and a constant number of filters, D : DCGAN			
			
G : 4-layer 512-dim ReLU MLP, D : DCGAN			
			
No normalization in either G or D			
			
Gated multiplicative nonlinearities everywhere in G and D			
			
tanh nonlinearities everywhere in G and D			
			
101-layer ResNet G and D			
			

GANs for text generation

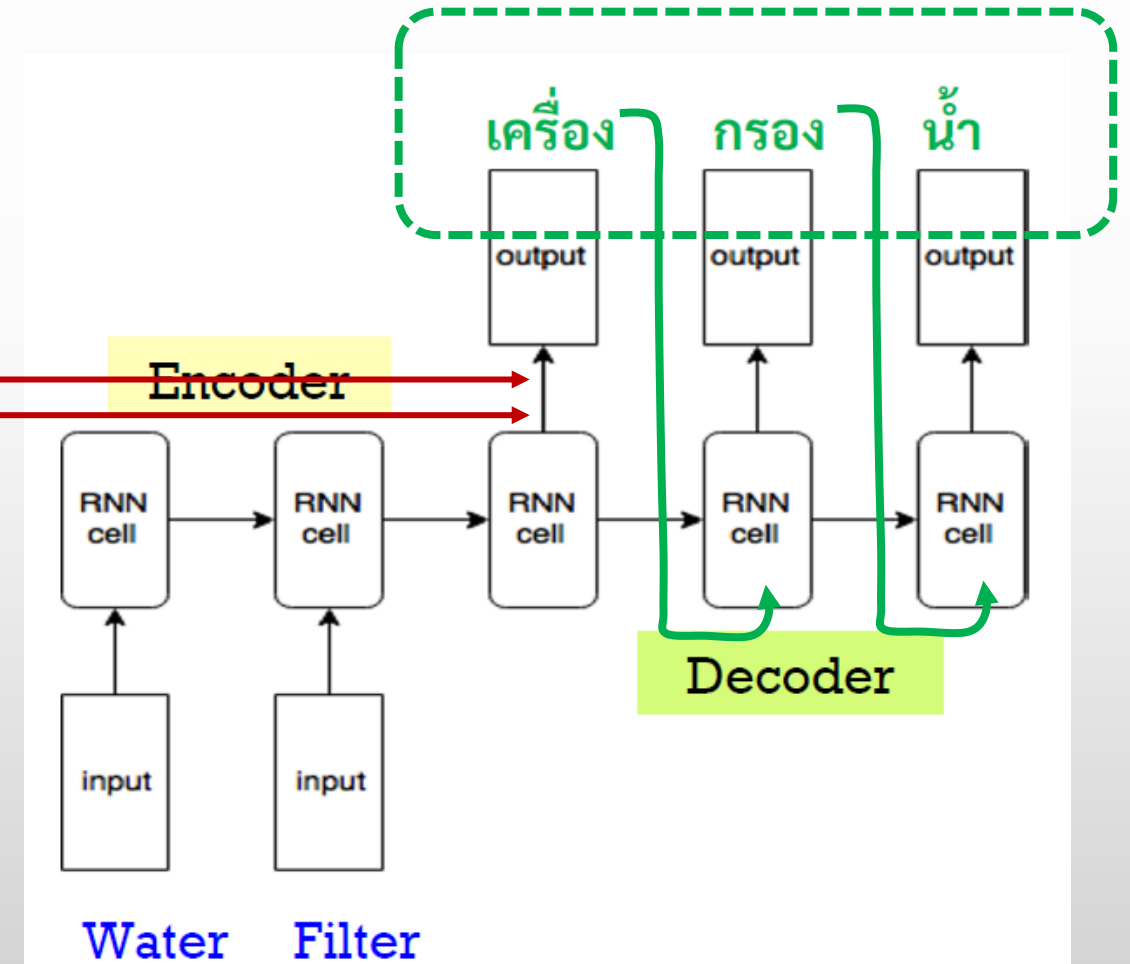
- Autoregressive decoding includes a sampling process
- Cannot gradient descent



Discriminator sees the sampled sentence.
Says it's good or bad

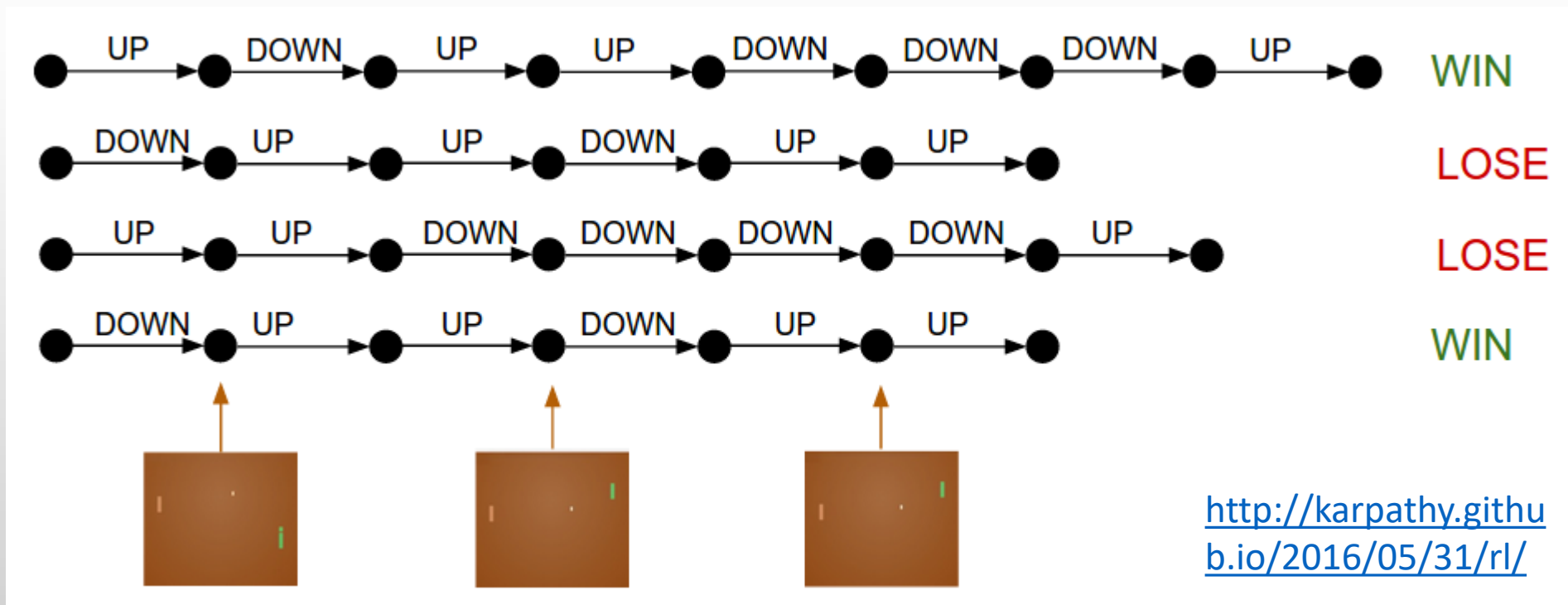


- Two popular methods
 - REINFORCE
 - Gumbel-Softmax approximation



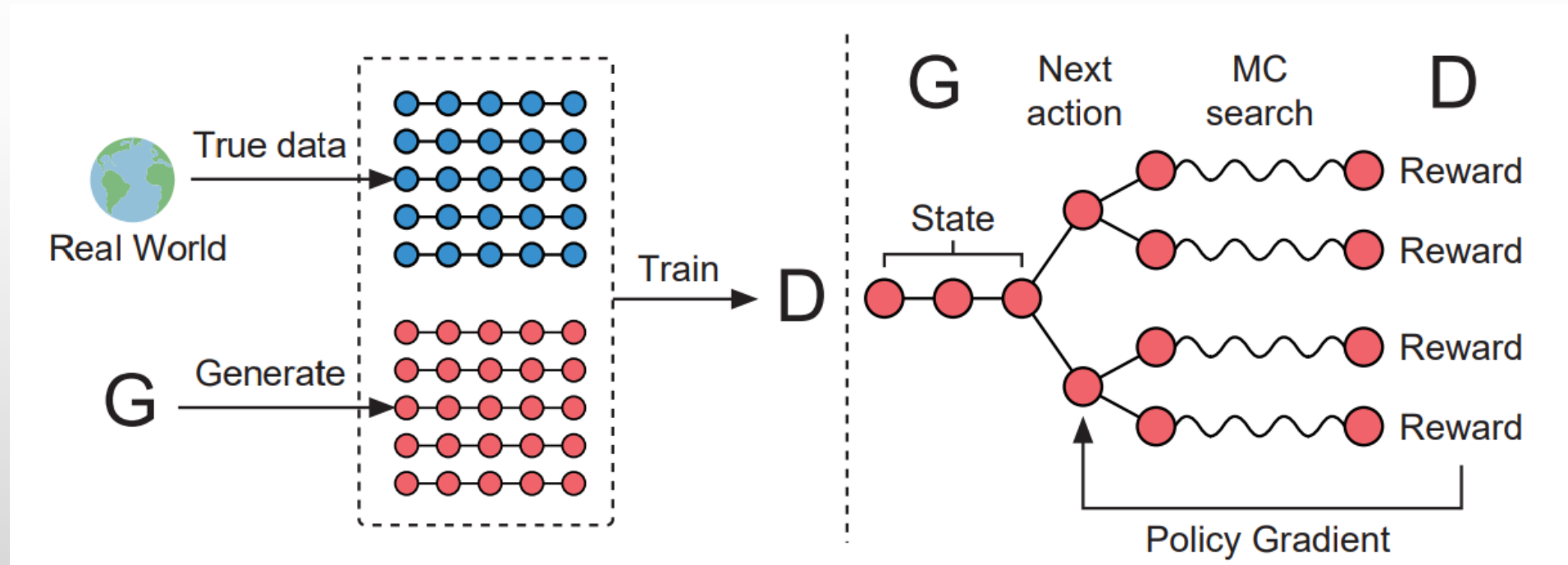
RL and policy gradients

- Credit assignment problem in reinforcement learning
- Which move makes you win?
- For RL with policy gradients, we increase the probability of every move that results in a win (REINFORCE algorithm)



GAN with text generation(SeqGAN)

- Use policy gradient to update the generator (the agent in RL setting)
- The discriminator (critic) gives the reward



Gumbel Softmax

- From softmax prob [0.3 0.4 0.1 0.2] -> word index z
- Pick randomly based on prob (sampling)

random value generated from Gumbel dist.

$$z = \text{one_hot} \left(\arg \max_i [g_i + \log \pi_i] \right)$$

index for each word

prob values from softmax

Gumbel Softmax

- From softmax prob [0.3 0.4 0.1 0.2] -> word index z
- Pick randomly based on prob (sampling)

random value generated from Gumbel dist.

$$z = \text{one_hot} \left(\arg \max_i [g_i + \log \pi_i] \right)$$

estimate using Gumbel trick

index for each word

prob values from softmax

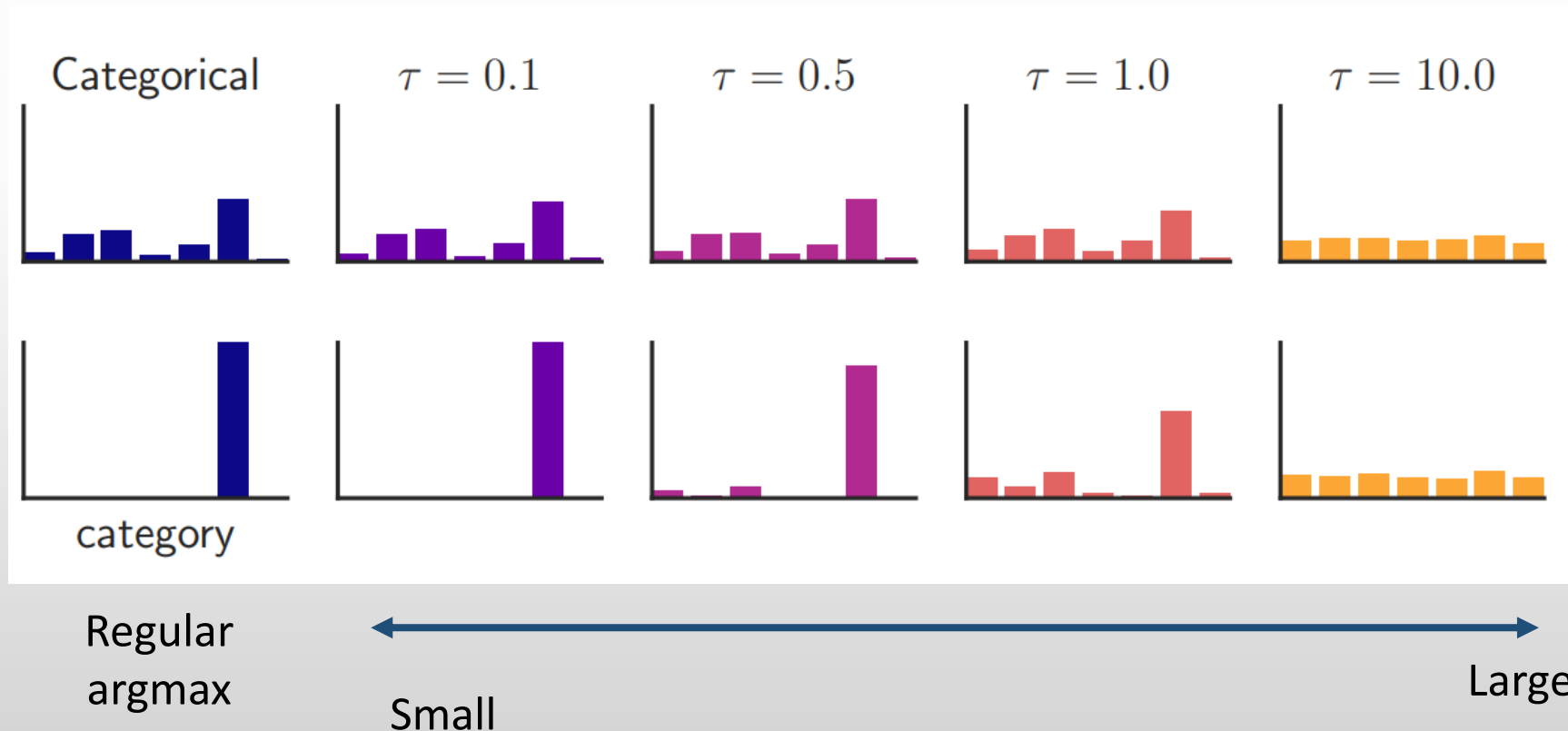
$$y_i = \frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{j=1}^k \exp((\log(\pi_j) + g_j)/\tau)} \quad \text{for } i = 1, \dots, k.$$

Temperature parameter

Gumbel Softmax

- Temperature can be scaled
- y can be back propagated through

$$y_i = \frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{j=1}^k \exp((\log(\pi_j) + g_j)/\tau)} \quad \text{for } i = 1, \dots, k.$$



HW4

[https://drive.google.com/drive/folders/13nW1B1DGkwRrUjZfsEc_rKDiG8MCaE6mU?usp=share link](https://drive.google.com/drive/folders/13nW1B1DGkwRrUjZfsEc_rKDiG8MCaE6mU?usp=share_link)