



# TIME SERIES ANALYSIS & RECURRENT NEURAL NETWORKS

#13

- Bayesian RNN
- Generative Adversarial Networks (GAN)
- Attention & Transformers

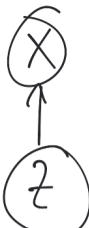
**Main lecture:** Daniel Durstewitz

**Exercises:** Leonard Bereska, Manuel Brenner,  
Daniel Kramer, Georgia Koppe

Heidelberg University

## Bayesian RNN

$$x = \{x_1, \dots, x_T\}$$



$$p_\theta(x, z), p_\theta(z|x)$$

$$p(\theta, z|x) = \frac{p(x|\theta, z)p(\theta, z)}{p(x)}$$

posterior

$$p(x) = \int_{\theta} \int_z p(x, z, \theta) dz d\theta$$

→ Variational Inference

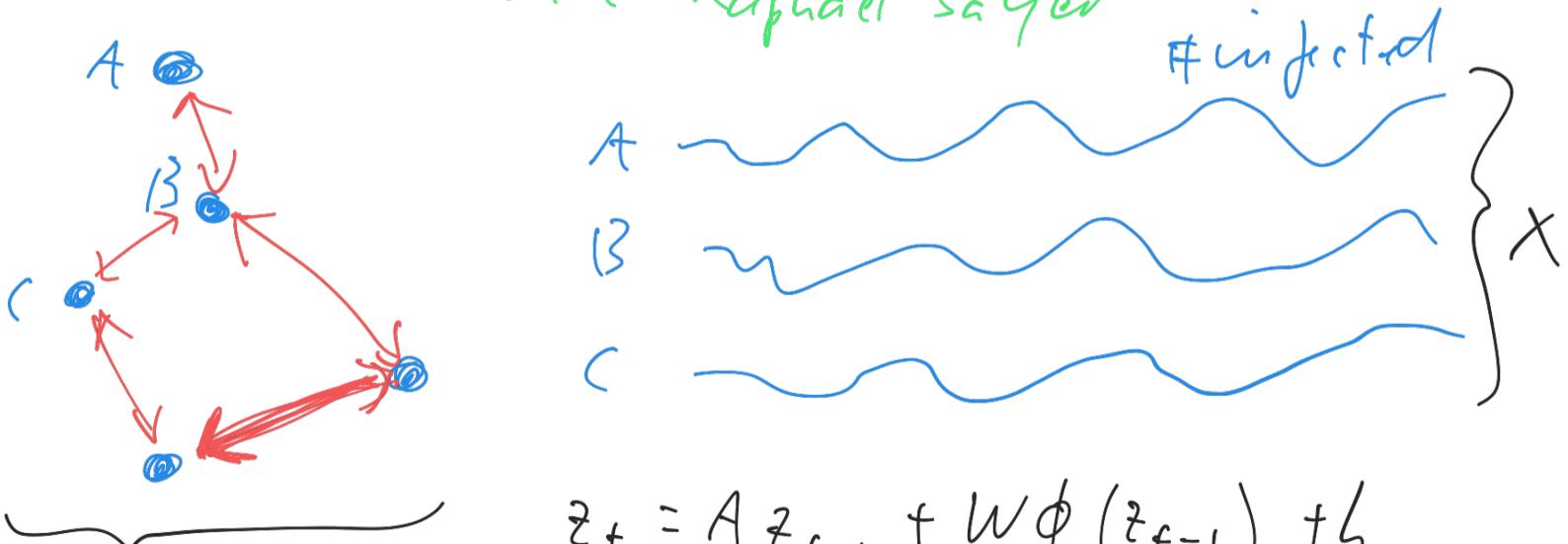
$$\begin{aligned} & \underset{\mathcal{Z}}{\operatorname{argmin}} \left\{ KL[q_Z(\theta, z|x) \| p(\theta, z|x)] \right\} \\ &= \int_{\theta, z} q_Z(\theta, z|x) \log \frac{q_Z(\theta, z|x)}{p(\theta, z|x)} d\theta dz \\ &= -H_q[q_Z(\theta, z|x)] - E_q[\log \frac{p(\theta, z|x)}{p(x)}] \\ &= -H_q[q] - E_q[\log p(\theta, z|x)] + E_q[\log p(x)] \\ &= -ELBO[\mathcal{Z}] + \log p(x) \end{aligned}$$

$$q_Z(\theta, z|x) = q_\phi(\theta|x) q_\psi(z|x) \text{ assumption!}$$

$$\mathcal{Z} = \{\phi, \psi\}$$

$$\rightarrow \underset{\phi, \psi}{\operatorname{argmax}} E_{q_\psi} \left[ ELBO_\phi(\phi, z) + \log p(\theta) - \log q_\psi(\theta|x) \right]$$

- appendix Kingma & Welling (2014)  
- MSc Raphael Sayer



$$z_t = Az_{t-1} + W\phi(z_{t-1}) + h$$

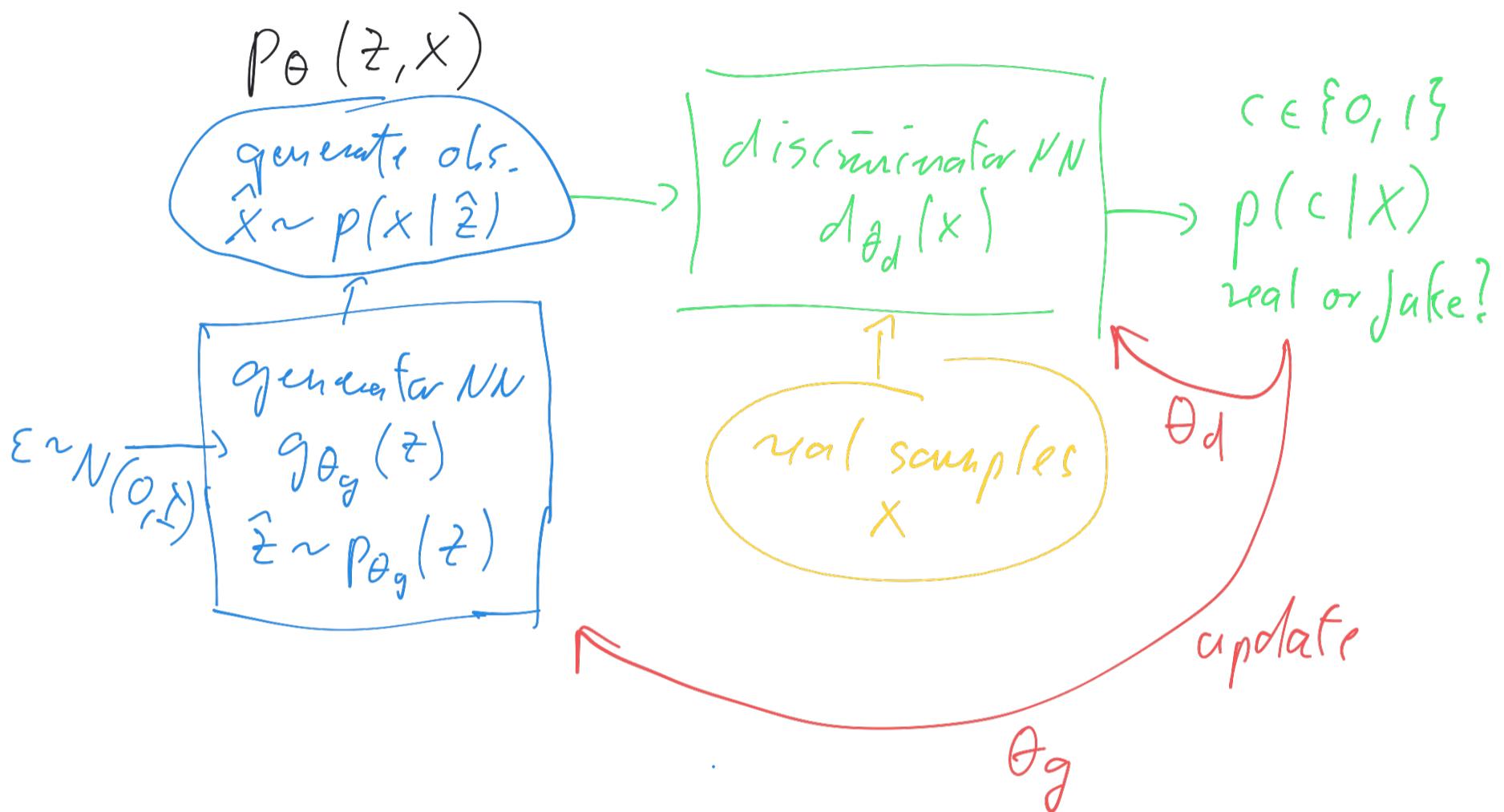
struct. prior

$$p(w) \longrightarrow p(w|z|x)$$

- Previous experiments →  $p(A, W, h)$

# Generative Adversarial Networks (GAN)

Goodfellow et al. (2014)



→ g and d play a '0-sum' game

$v(\theta_d, \theta_g)$  : payoff for discriminator

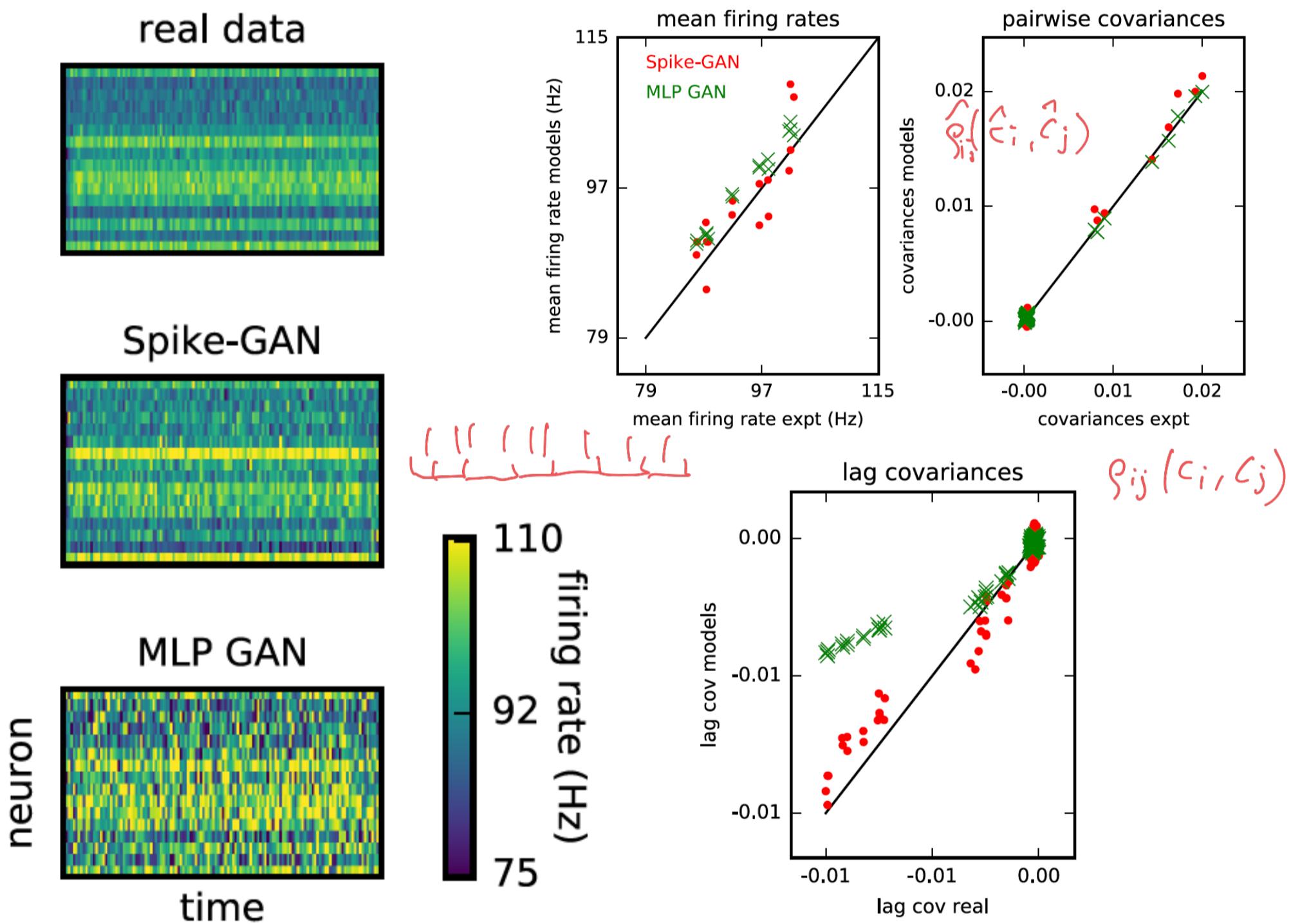
-  $v(\theta_d, \theta_g)$  : payoff for generator

$$\text{e.g. } v(\theta_d, \theta_g) = \bar{E}_{x \sim p_{\text{data}}} [\log d_{\theta_d}(x)] + \begin{cases} p(c=1|x) & c=1 : \text{real} \\ p(c=0|x) & c=0 : \text{fake} \end{cases}$$

$$\bar{E}_{\hat{x} \sim p_{\theta_g}} [\log (1 - d_{\theta_d}(\hat{x})]$$

$$\theta_g^* = \arg \min_{\theta_g} \left[ \max_{\theta_d} [v(\theta_d, \theta_g)] \right]$$

→ data augmentation



Molano-Mazon ... Panzeri (2018) ICLR

# Attention mechanisms

## Language models

- machine translation
- Query

## • Representation of language

- categ. distribs., "1-hot-encoding"
- vocabs.  $V$  of  $K$  words

$$c_t = (0 \ 0 \ 0 \ 1 \ 0 \ 0 \ . \ . \ . \ 0) \in \{0, 1\}^K$$

$\uparrow$   
i<sup>th</sup> word in  $V$

$$\sum c_{it} = 1, \quad \{c_1, \dots, c_T\}$$

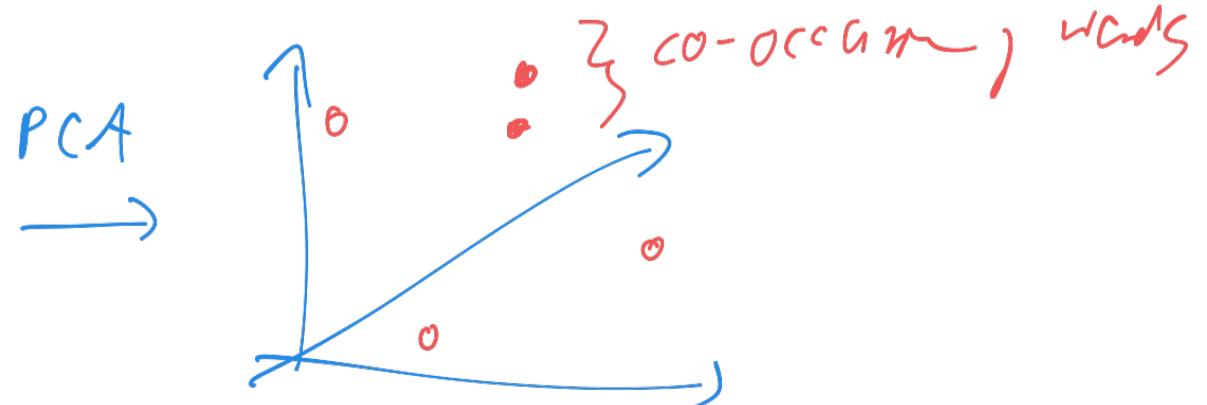
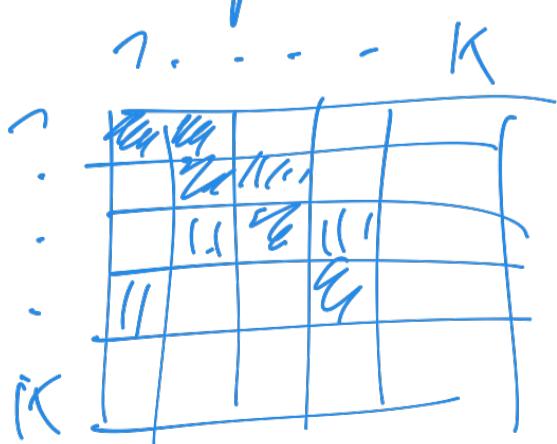
$$\rightarrow p(c_t) = \prod_{i=1}^K \pi_i^{c_{it}}$$

output from  
 $NN$

$$\pi_i = \text{Softmax}[x_i] = \frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}}$$

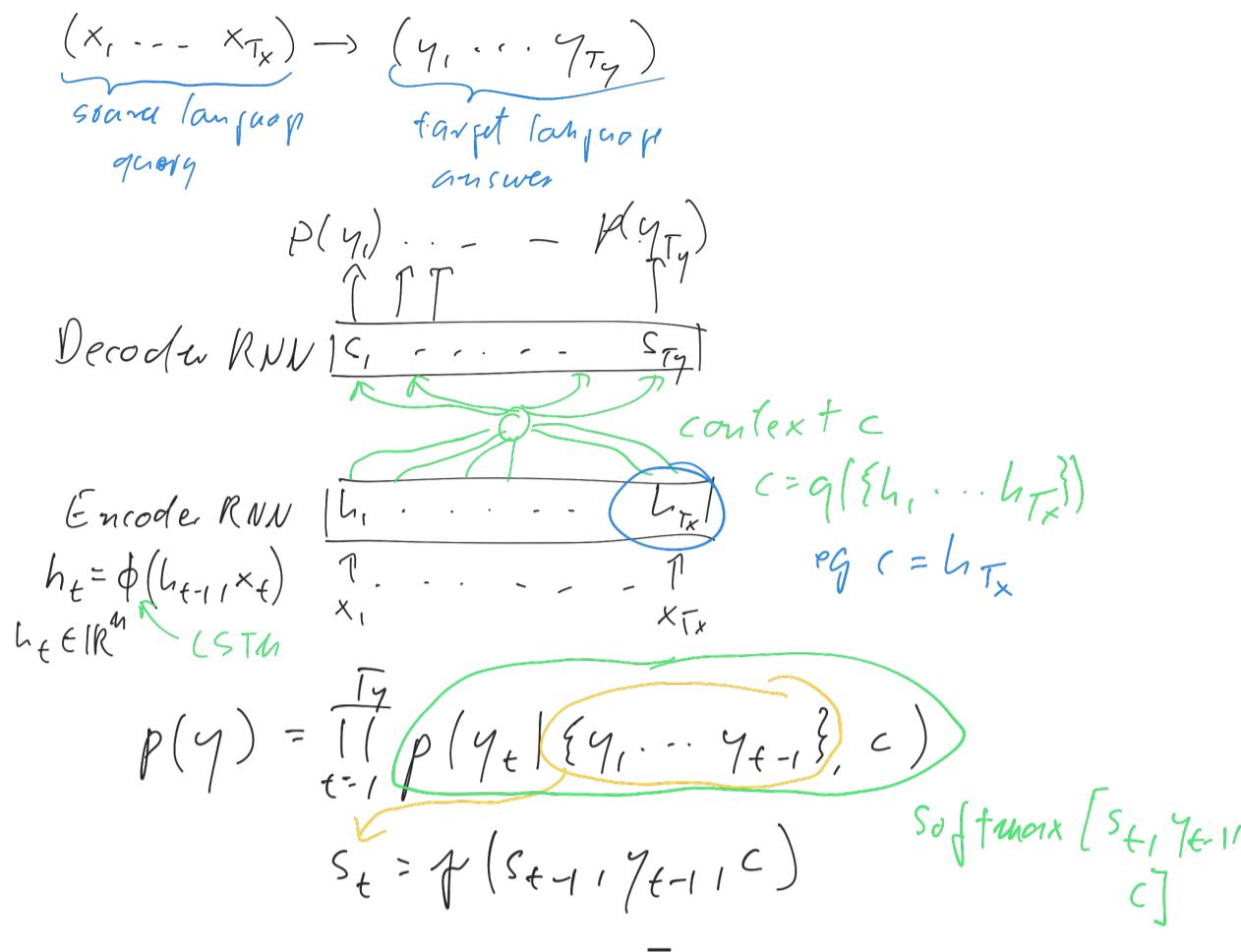
- ## • word embedding:
- contin. repr.
  - lower-dim.
  - acknowledge context

e.g. PCA on word cov. matrix



$$y_i = \underbrace{Ax_i}_{d \times k, d \ll k}.$$

"classical approach" (Sutskever, Graves, Bengio)



- Bahdanau, Cho, Bengio (2015) ICLR  
→ Attu.-mech.

$$s_t = \text{softmax}(s_{t-1}, y_{t-1}, c_t)$$

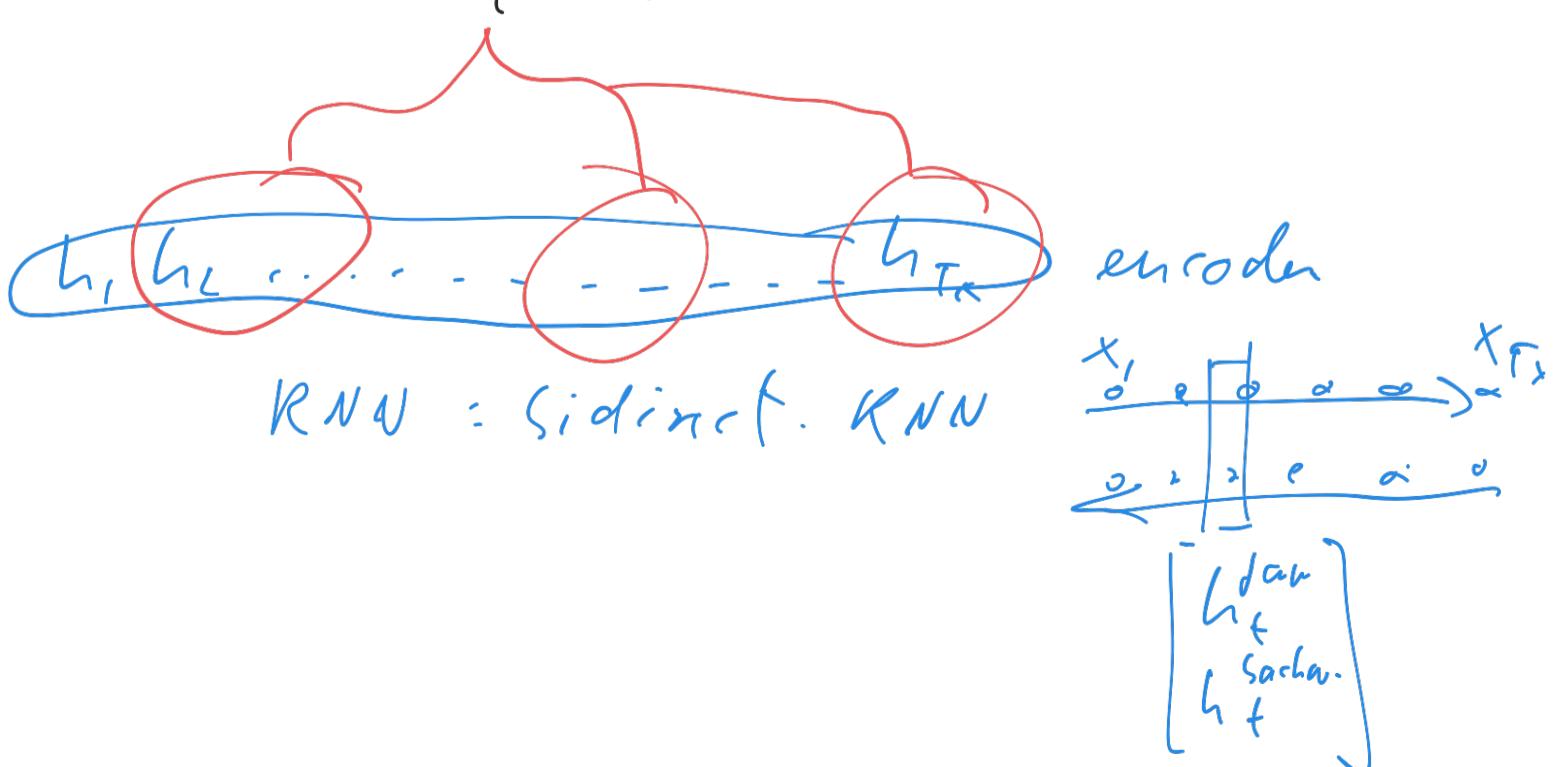
$$c_t = \sum_{j=1}^{T_x} \alpha_{tj} h_j \quad \text{"alignment scores, atten. weights"}$$

$$\alpha_{tj} = \text{softmax}(e_{tj}) = \frac{\exp(e_{tj})}{\sum_{k=1}^{T_x} \exp(e_{tk})}$$

"relevance / attu. func."  $e_{tj} = g(h_j, s_{t-1})$

$$p(y_t) = \underbrace{v_a^T \tanh(W_a s_{t-1} + U_a h_j)}_{1.\text{-Lidd-Lage FNN}}$$

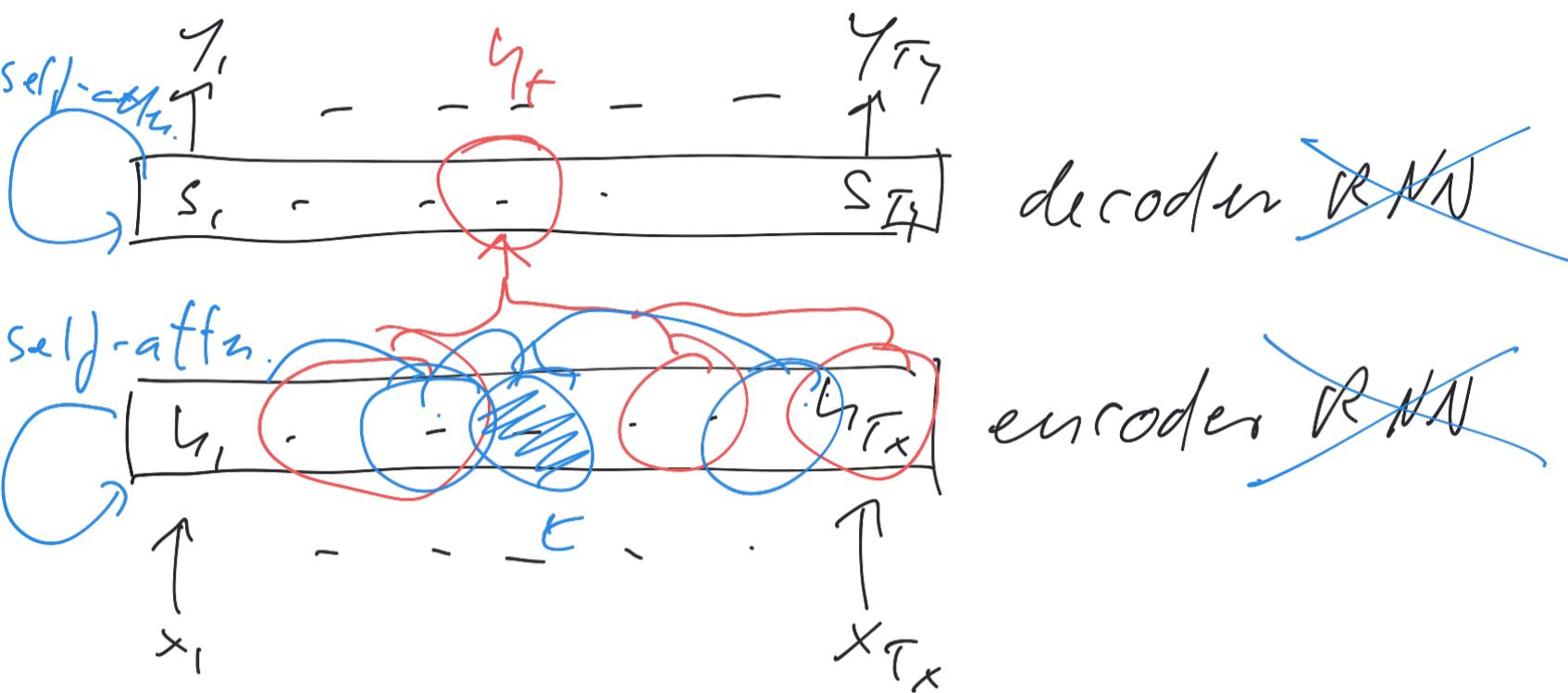
$s_t$  decoder



The  
agreement  
on  
the  
European  
Economic  
Area  
was  
signed  
in  
August  
1992  
. <end>  
  
L'  
accord  
sur  
la  
zone  
économique  
européenne  
a  
été  
signé  
en  
août  
1992  
. <end>

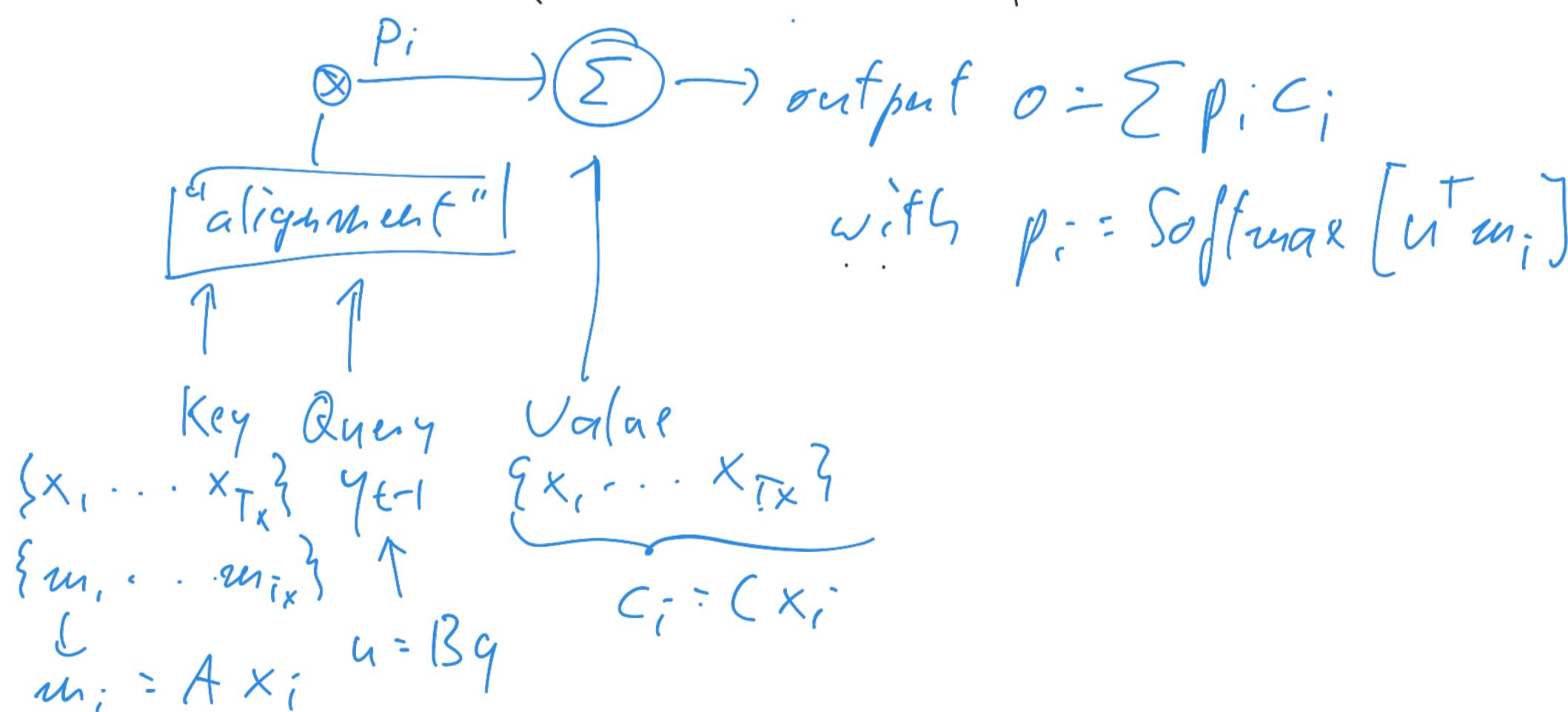
Destruction  
of  
the  
equipment  
means  
that  
Syria  
can  
no  
longer  
produce  
new  
chemical  
weapons  
. <end>  
  
La  
destruction  
de  
l'  
équipement  
signifie  
que  
la  
Syrie  
ne  
peut  
plus  
produire  
de  
nouvelles  
armes  
chimiques  
. <end>

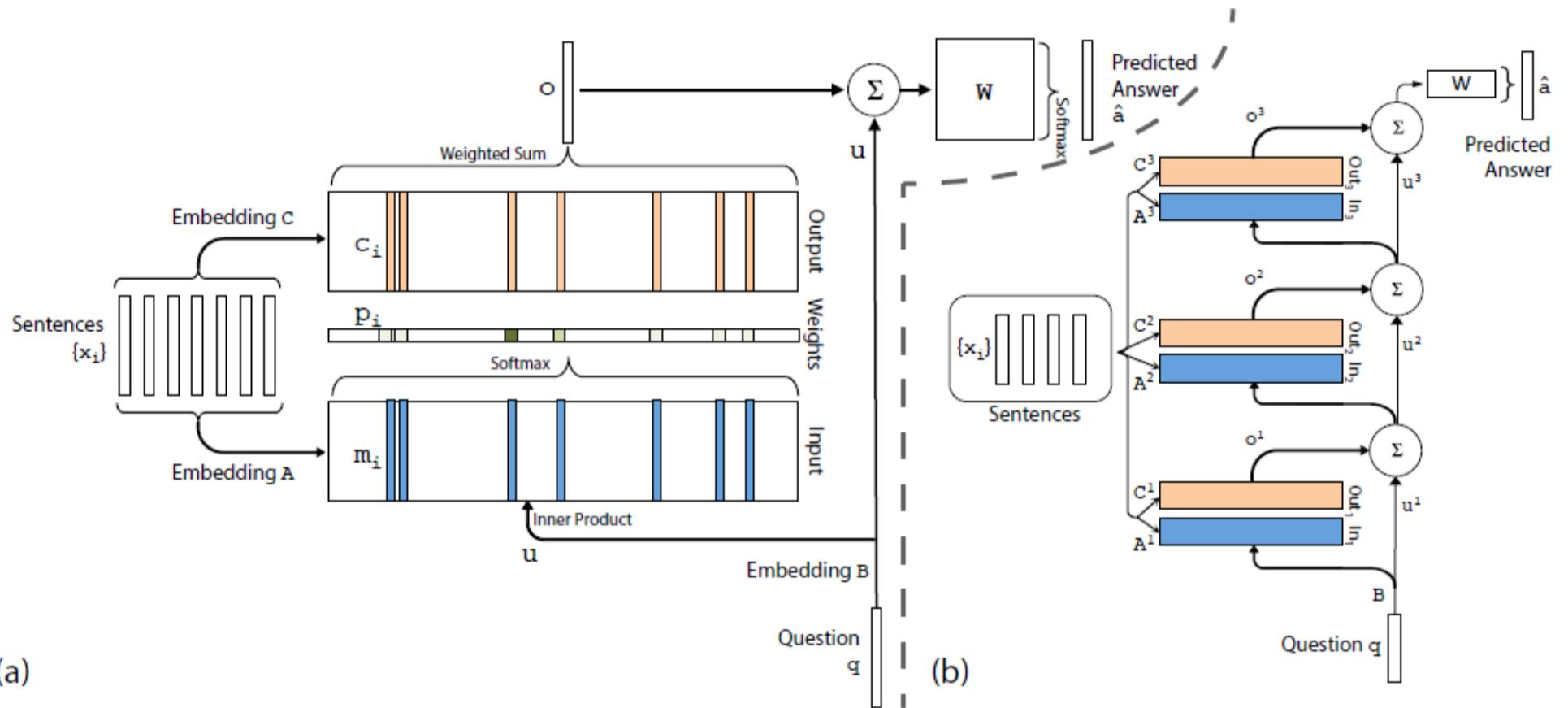
Bahdanau ... Bengio (2015) ICLR



- strongly parallelizable
- Transformer, Vaswani et al. (2017)  
NIPS

Atttn. mech. (Sukhbaatar et al. 2015,



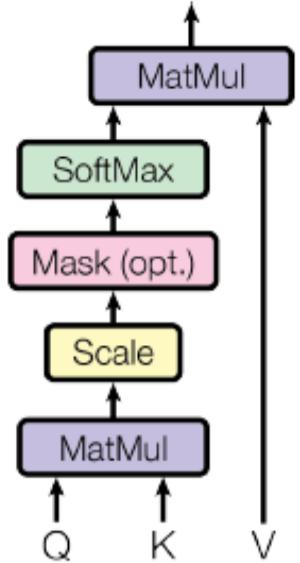


(a)

(b)

Sukhbataar ... Fergus (2015) NIPS

### Scaled Dot-Product Attention



Vaswani et al. (2017) NIPS

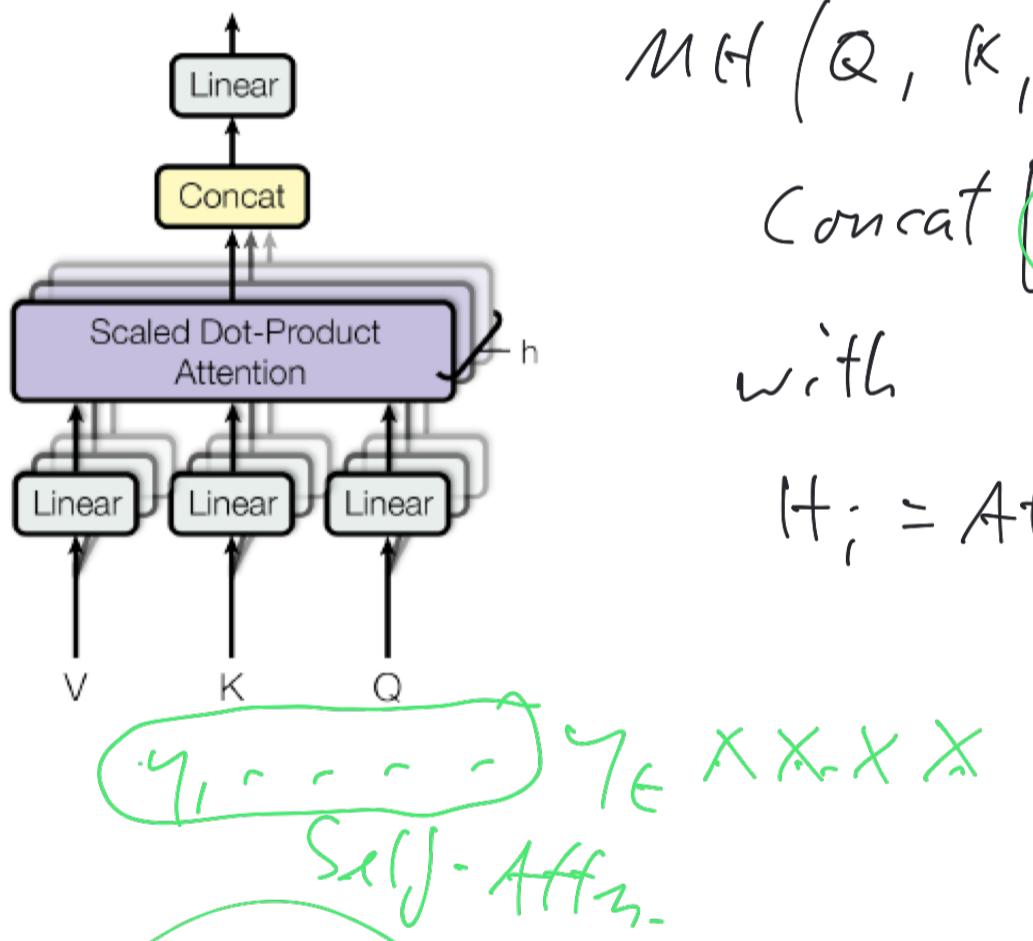
$$Q = \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_{T_q} \end{bmatrix}, K = \begin{bmatrix} k_1 \\ k_2 \\ \vdots \\ k_{T_K} \end{bmatrix}, V = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_{T_V} \end{bmatrix}$$

$$\text{Attn}(Q, K, V) =$$

$$\text{softmax} \left[ \frac{Q K^T}{\sqrt{d_k}} \right] V$$

$x_1 \dots x_T$   
 $y_1 \dots y_T$

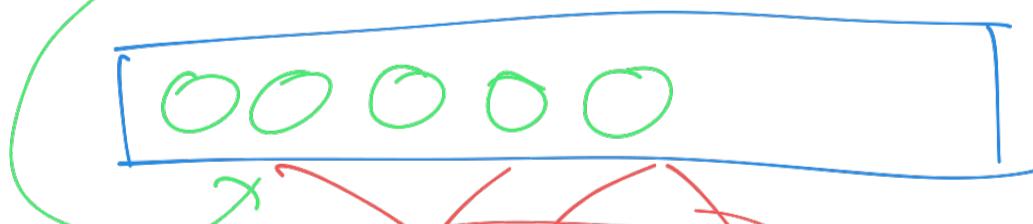
### Multi-Head Attention



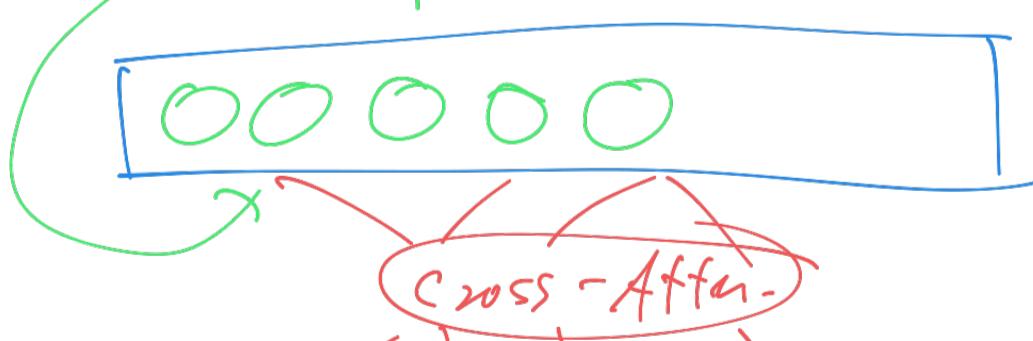
$$MHA(Q, K, V) = \text{concat}[H_1, \dots, H_n] W_0$$

with

$$H_i = \text{Attn}(QW_i^Q, KW_i^K, VW_i^V)$$



Decoder



Encoder



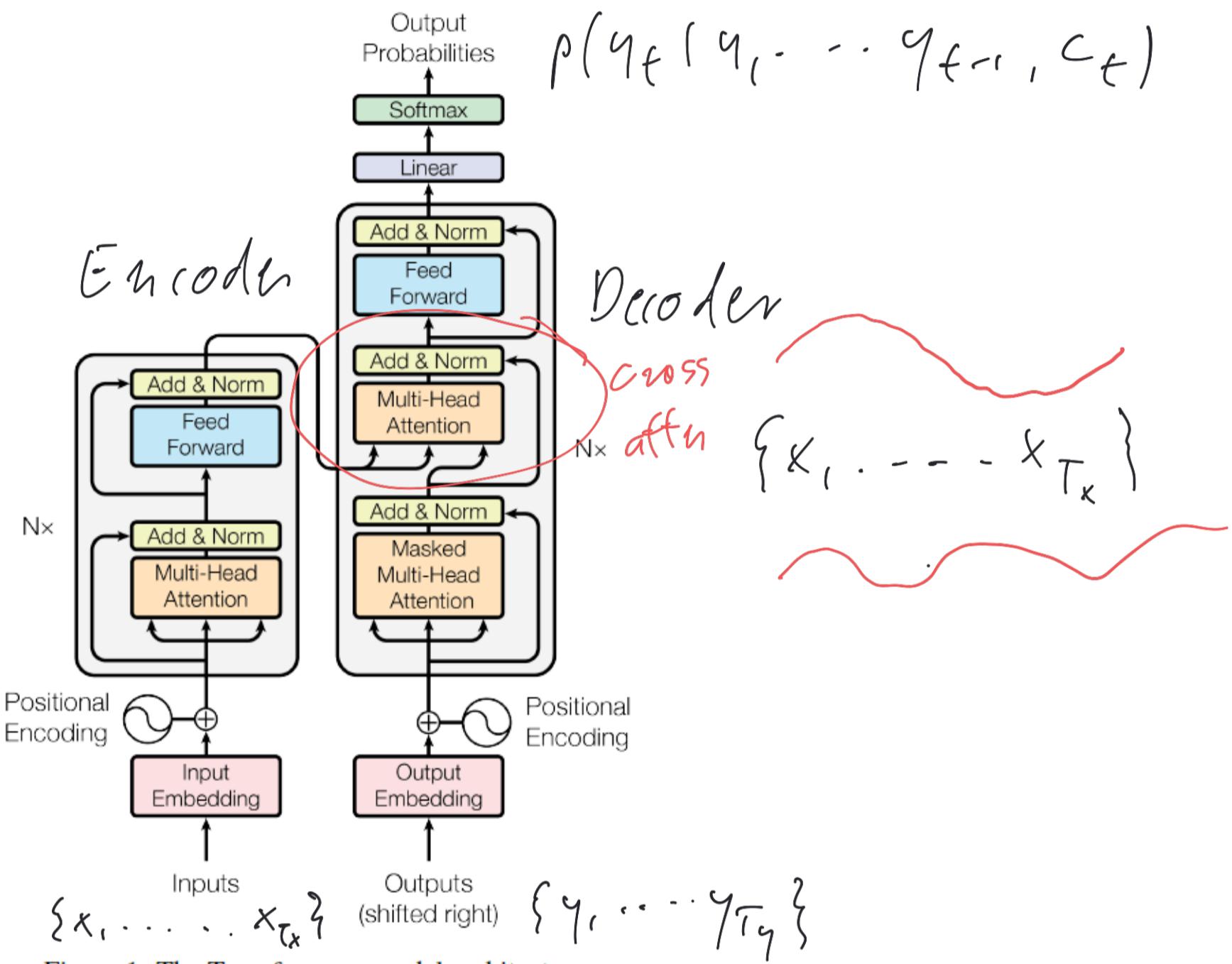
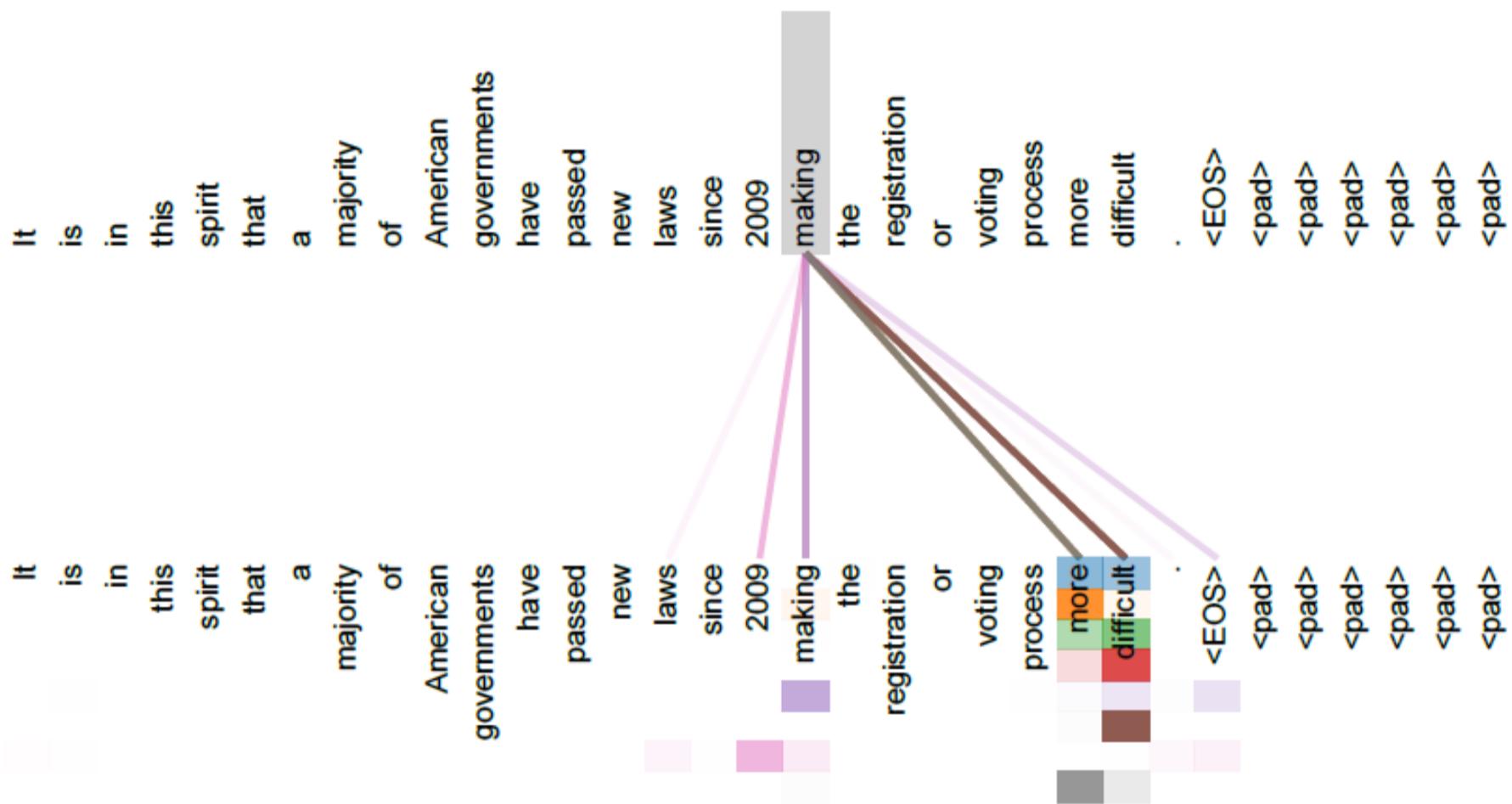


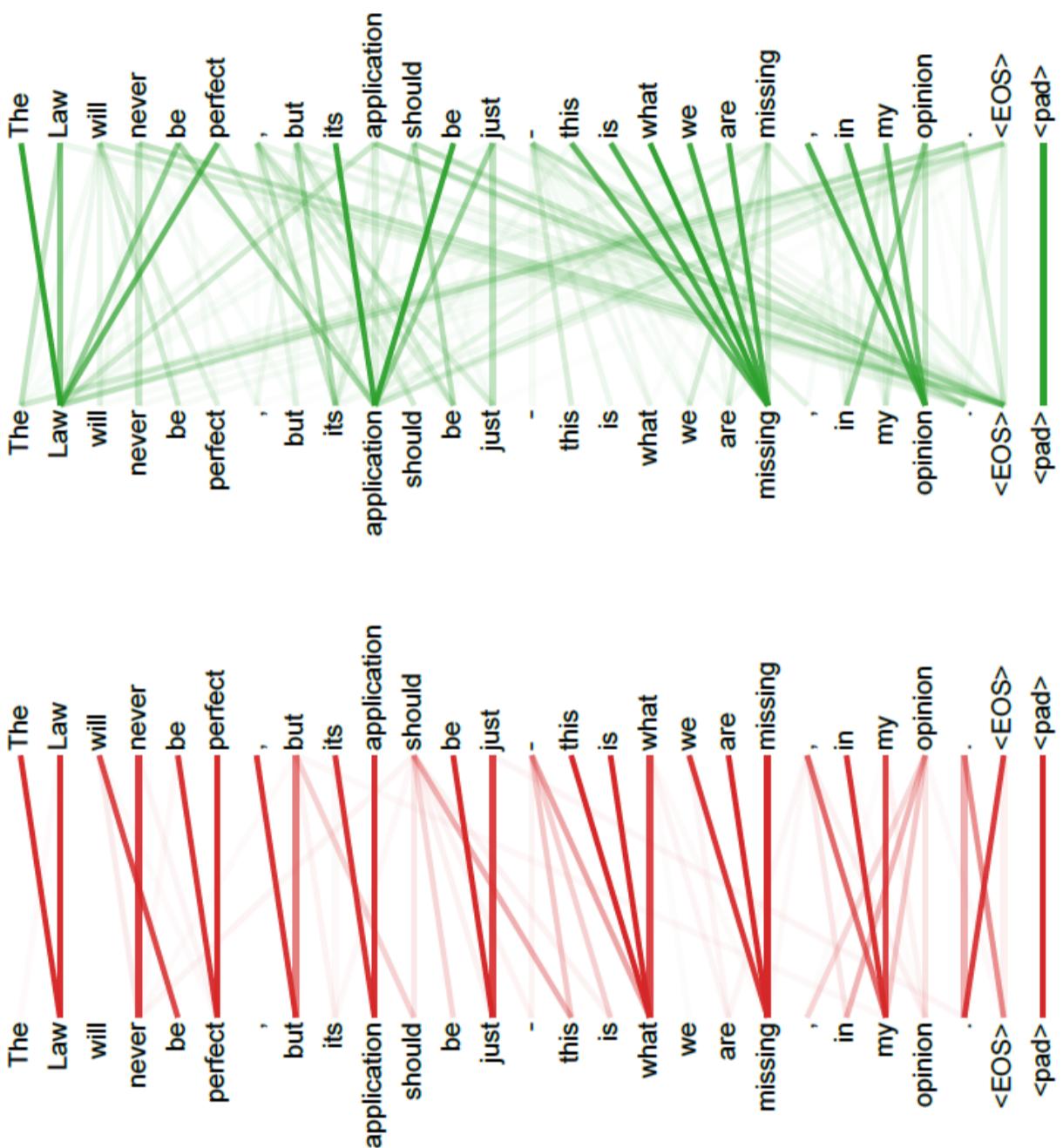
Figure 1: The Transformer - model architecture.

Vaswani et al. (2017) NIPS

## Attention Visualizations



Vaswani et al. (2017) NIPS



Vaswani et al. (2017) NIPS

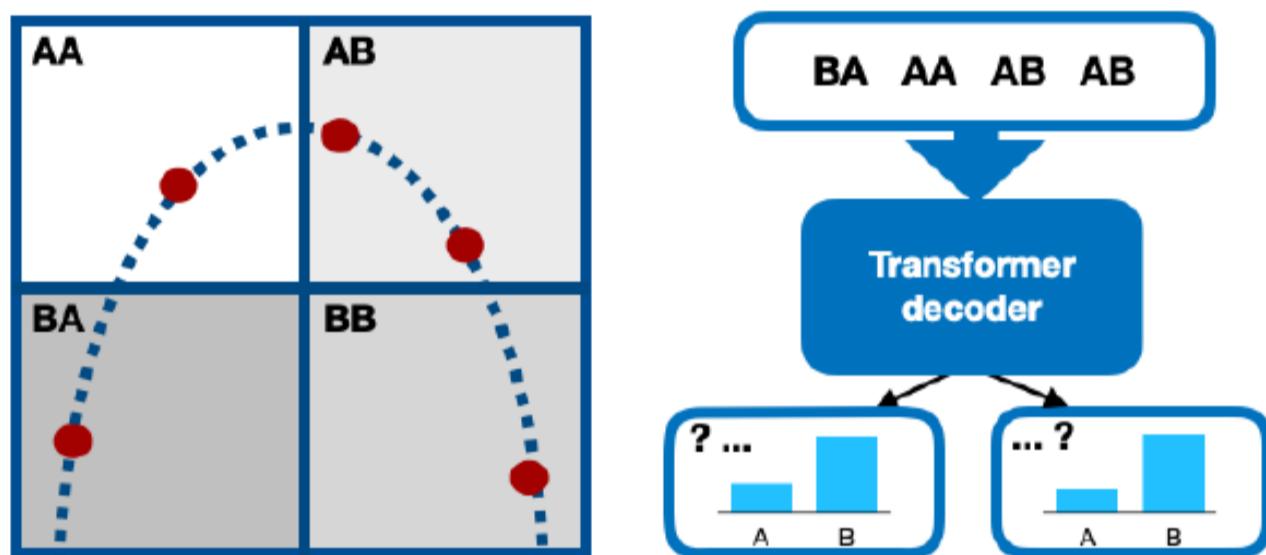


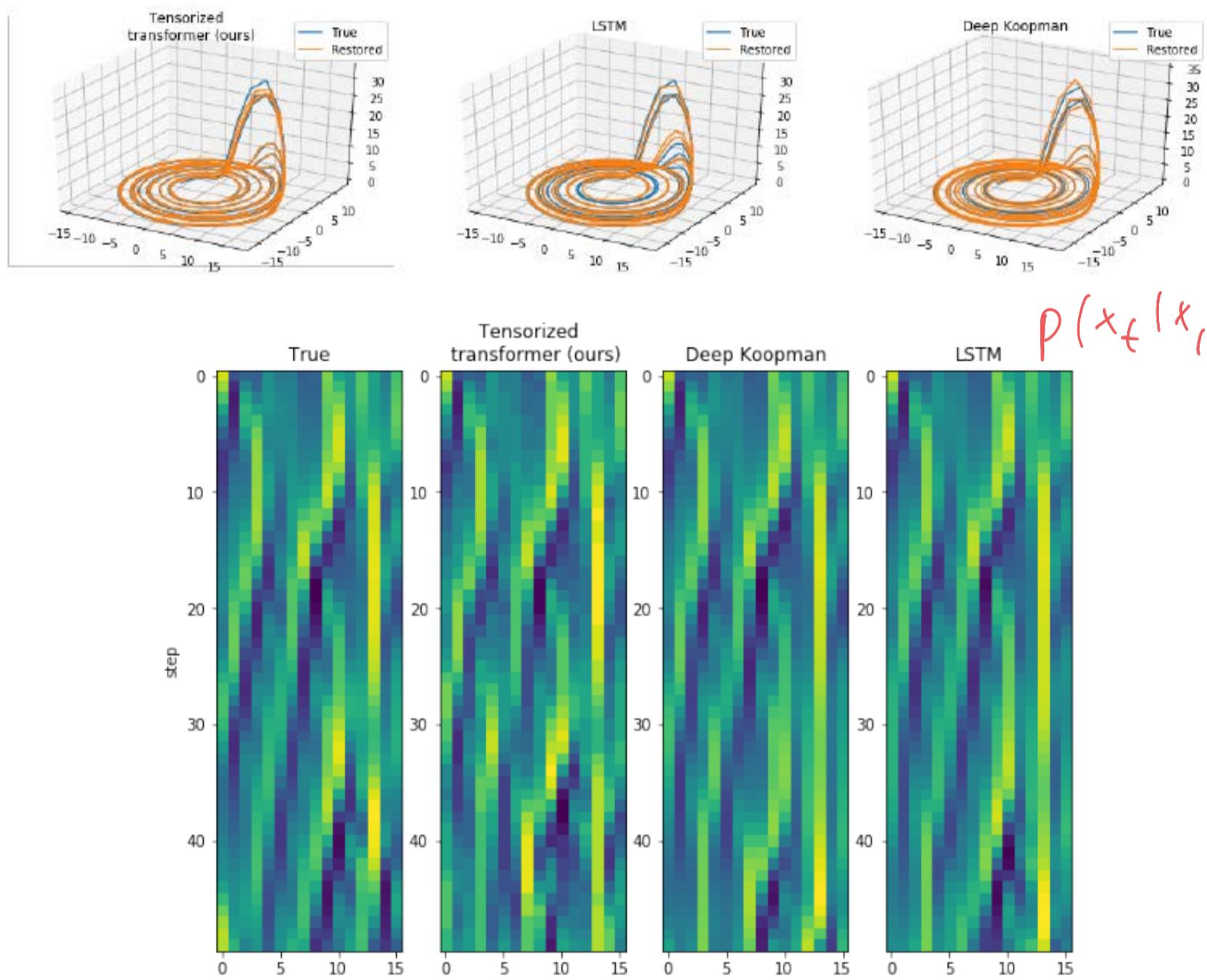
Figure 1: An example of the proposed method for the 2-dimensional trajectory. The initial states are discretized along each dimension; then transformer decoder is used to predict conditional probabilities along each dimension.

Shalova & Oseledets (2020) arXiv

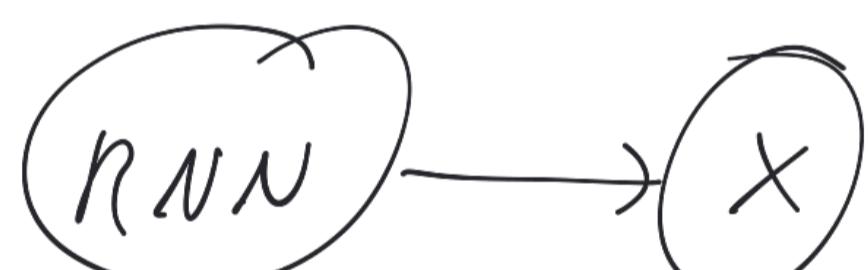
$$p(x_{k+1} | x_1, \dots, x_k)$$

$$\mathcal{L} = \text{cross-entropy} = -\frac{1}{N} \sum_t \log p(x_t | x_1, \dots, x_{t-1})$$

Softmax [Transform]  
 evaluated at time  $x_t$



Shalova & Oseledets (2020) arXiv

- basic dist. cov, cov } 
- ARMA
- Poisson - AR
- latent var., lat. DS 
- EM, KFS  $p(z|x)$
- DS
- RNN, long-term dep.
-  , EM: EKF, PF,  
Caplace approx.
- seq-VAE / VAE
- Bayesian RNN, GAN
- Affin., Transforms