



# RNN-based encoder-decoder models

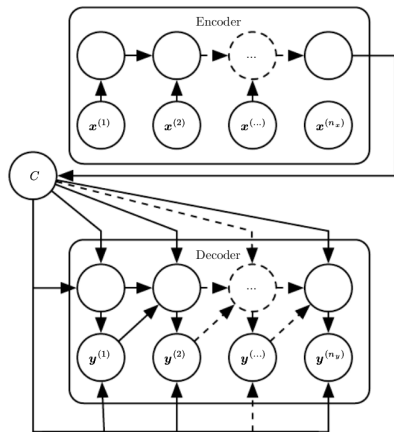
For a pair  $(x, y)$ , we model the conditional distribution:

$$p_{\theta}(y|x) = \prod_t p_{\theta}(y^t|x, y^{<t})$$

- $(h_{enc}^t)$  a context to condition the decoder.
- $(h_{dec}^t)$  a predictor for the emission prob.

$$h^t = f_{\theta}(h^{t-1}, y^t) \quad (\text{recursion})$$

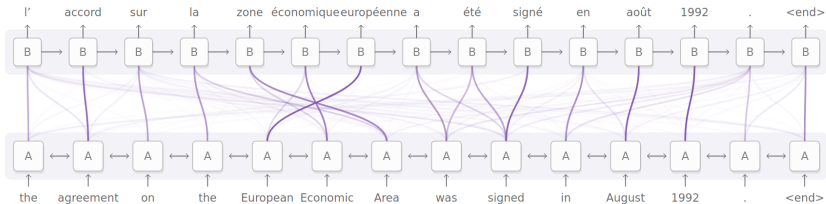
$$p_{\theta}(y^{t+1}|h^t) = \sigma(\mathbf{W}h^t) \quad (\text{prediction})$$



source: I. Goodfellow, Y. Bengio, and A. Courville,  
Deep Learning. MIT Press.

# Applications

## Machine translation



source: distill.pub

## Image captioning



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



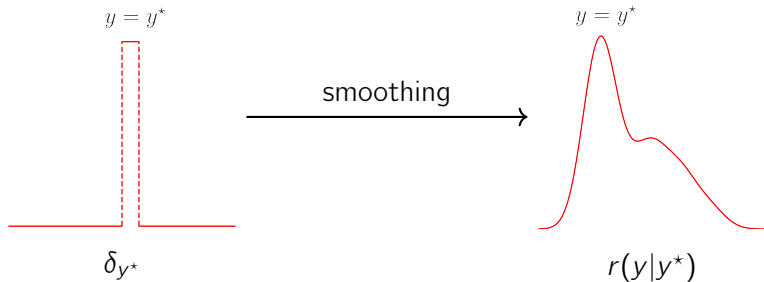
A stop sign is on a road with a mountain in the background.



OCR, Speech recognition, Times series...

## Training objectives: Smoothing

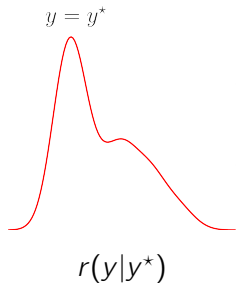
$$\ell_{\text{ML}}(y^*, x) = -\ln p_{\theta}(y^*|x) = D_{\text{KL}}(\delta(y|y^*)||p_{\theta}(y|x)) = \sum_{t=1}^T D_{\text{KL}}(\delta(y_t|y_t^*)||p_{\theta}(y_t|h_t))$$



- Zero-one Loss, all the outputs  $y \neq y^*$  are treated equally.
- Discrepancy at the sentence level between the training and evaluation reward.

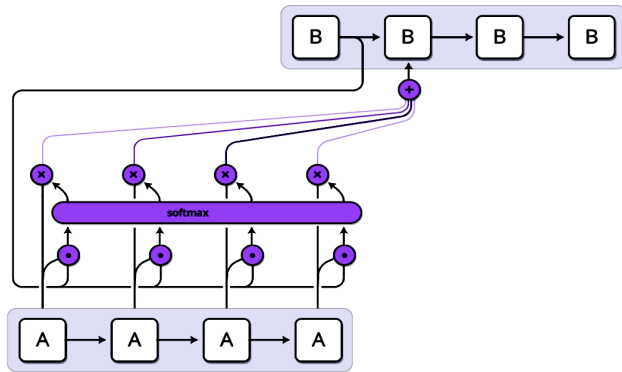
# Training objectives: Smoothing

$$\ell_{\text{ML}}(y^*, x) = -\ln p_{\theta}(y^*|x) = D_{\text{KL}}(\delta(y|y^*)||p_{\theta}(y|x)) = \sum_{t=1}^T D_{\text{KL}}(\delta(y_t|y^*_t)||p_{\theta}(y_t|h_t))$$



- Token-level: smooths w.r.t tokens similarity as assessed by a given word embedding (Word2vec, GloVe,...)
- Sequence-level: smooths w.r.t an evaluation reward via importance sampling.

# Encoder-decoder with attention



source: distill.pub

- The attending RNN generates a query describing what it wants to focus on and matches it with the source codes.

# Our 2D-Convolutional coder

Target sequence

Source sequence

	<start>	Alice	told	Bob	that	Charlie	told
Alice							
a							
dit							
a				•			
Bob							
que							
Charlie							

- Every target token has its *own interpretation* of the source tokens.
- Each convolution has the future target tokens masked  $\Rightarrow$  Autoregressive.
- The emission prob. is generated by a **DenseNet**.

Under review.

Thank you for your attention.

**Questions?**