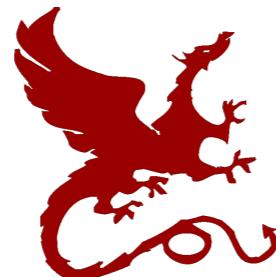


Neural Representation Learning in Linguistic Structured Prediction

Lingpeng Kong
Carnegie Mellon University



Linguistic Structures

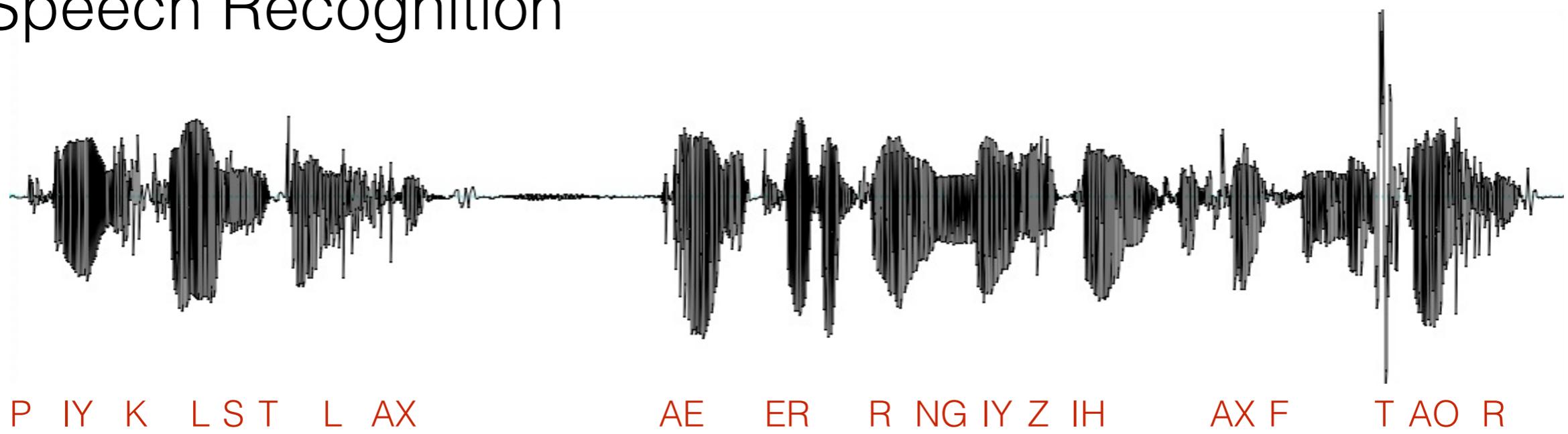
Word Segmentation

1997 年 继 续 主 办 。

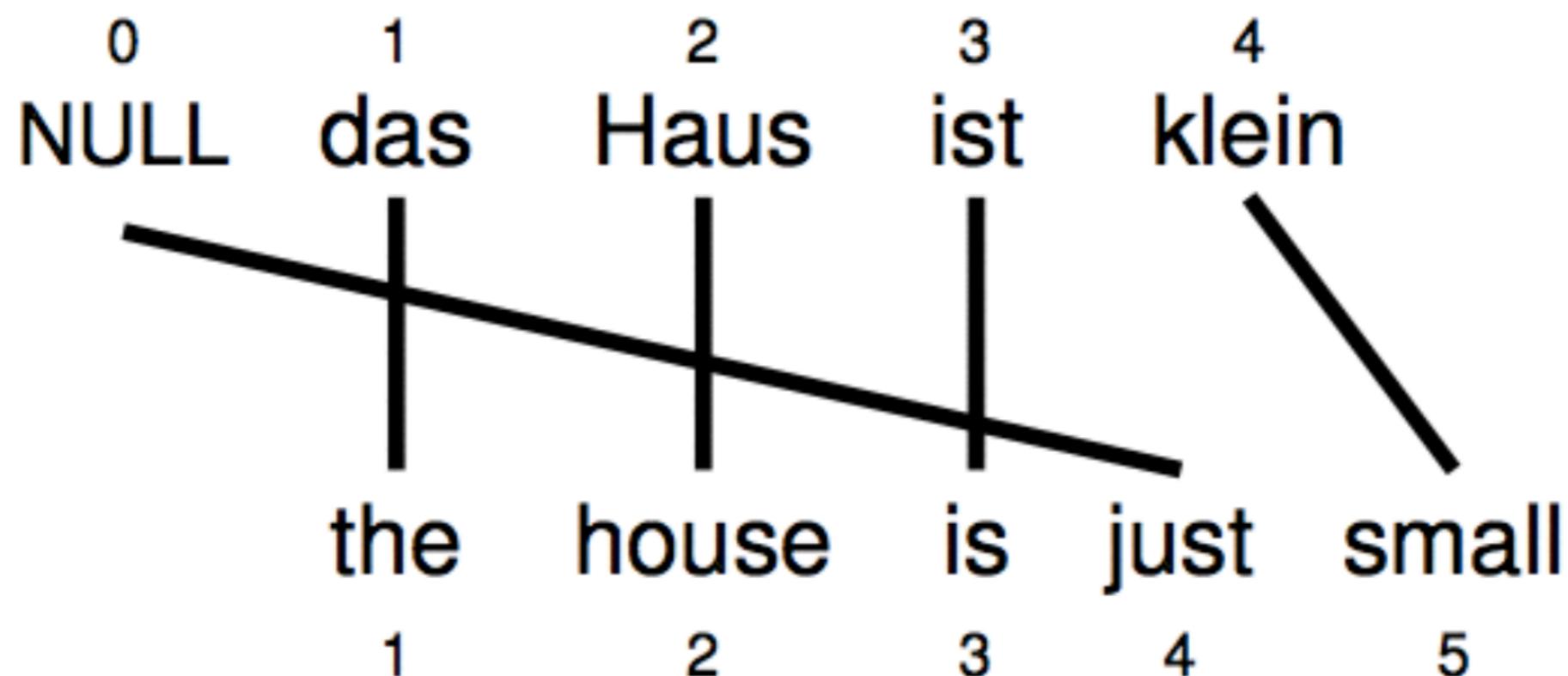
Chunking

NP PP NP PP
The angle of cats' ears is an important clue to their mood

Speech Recognition



Linguistic Structures

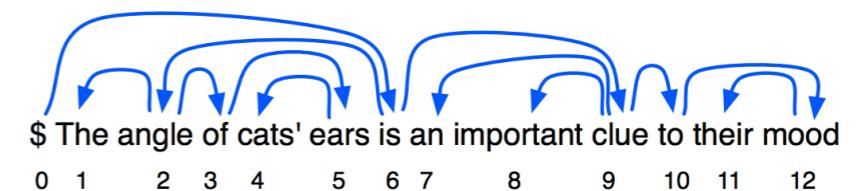


Neural Representation Learning

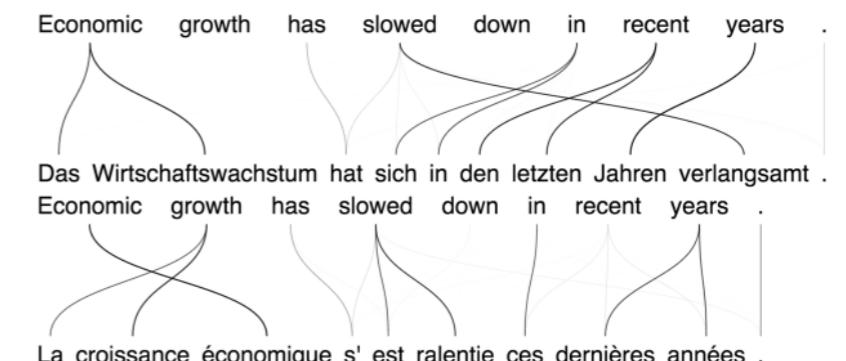
Speech Recognition (Graves et al, 2013)



Dependency Parsing (Andor et al, 2016)

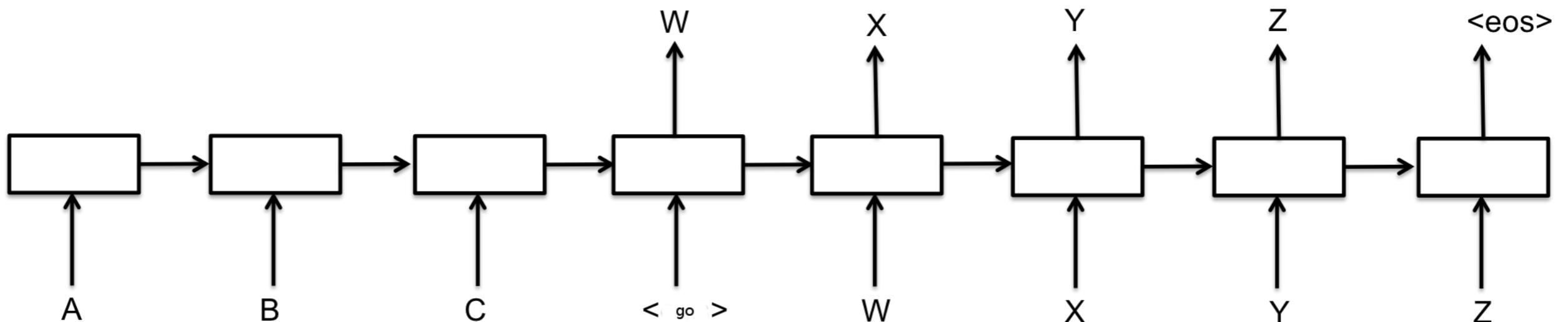


Machine Translation (Bahdanau et al, 2015)



Neural Representation Learning

Sequential Recurrent Neural Networks (RNNs)



Dense Representation

Goal

- The combination of explicit structure representations and learned distributed representations

The importance of linguistic structures:

Automated language understanding ([Manning, 2016](#))

Inductive bias ([Mitchell, 1980](#))

Outline

- Introduction
- Part I — Segmental Recurrent Neural Networks
[ICLR 2016, INTERSPEECH 2016, 2017]
- Part II — Inference and Regularization in Sequence to Sequence Models with Stochastic Attention
- Conclusion

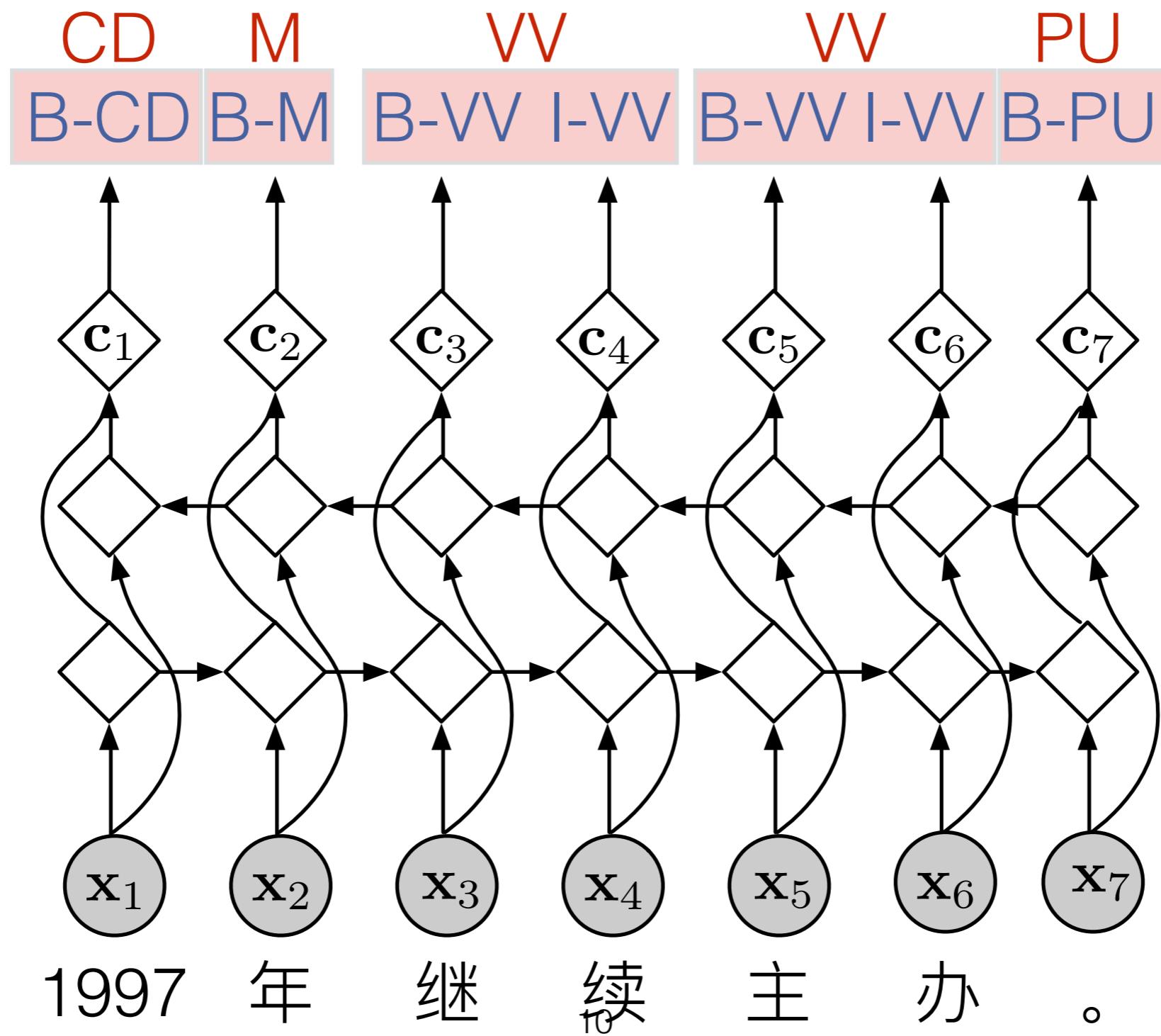
Outline

- Introduction
- Part I — Segmental Recurrent Neural Networks
[ICLR 2016, INTERSPEECH 2016, 2017]
- Part II — Inference and Regularization in Sequence to Sequence Models with Stochastic Attention
- Conclusion

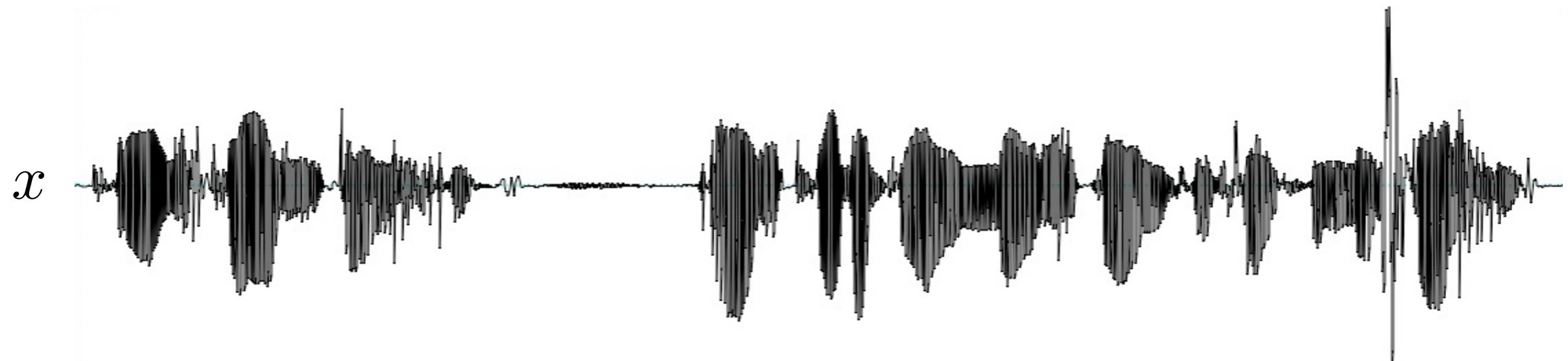
Joint Word Segmentation and POS Tagging

x	1997	年	继 续	主 办	。
y	CD	M	VV	VV	PU

Bi-LSTMs Tagger



Speech Recognition

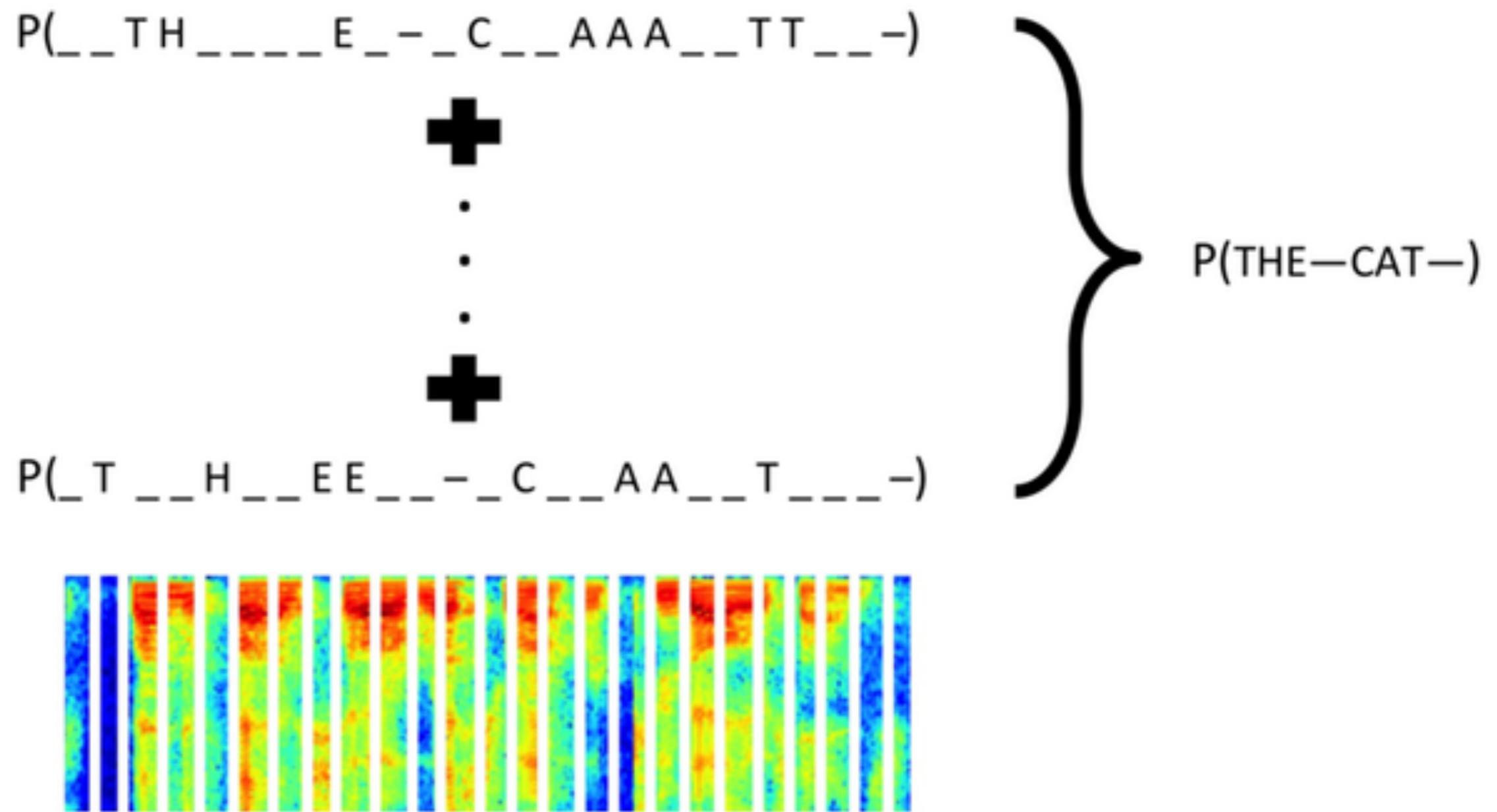


y P IY K L ST L AX AE ER R NG IY Z IH AX F T AO R

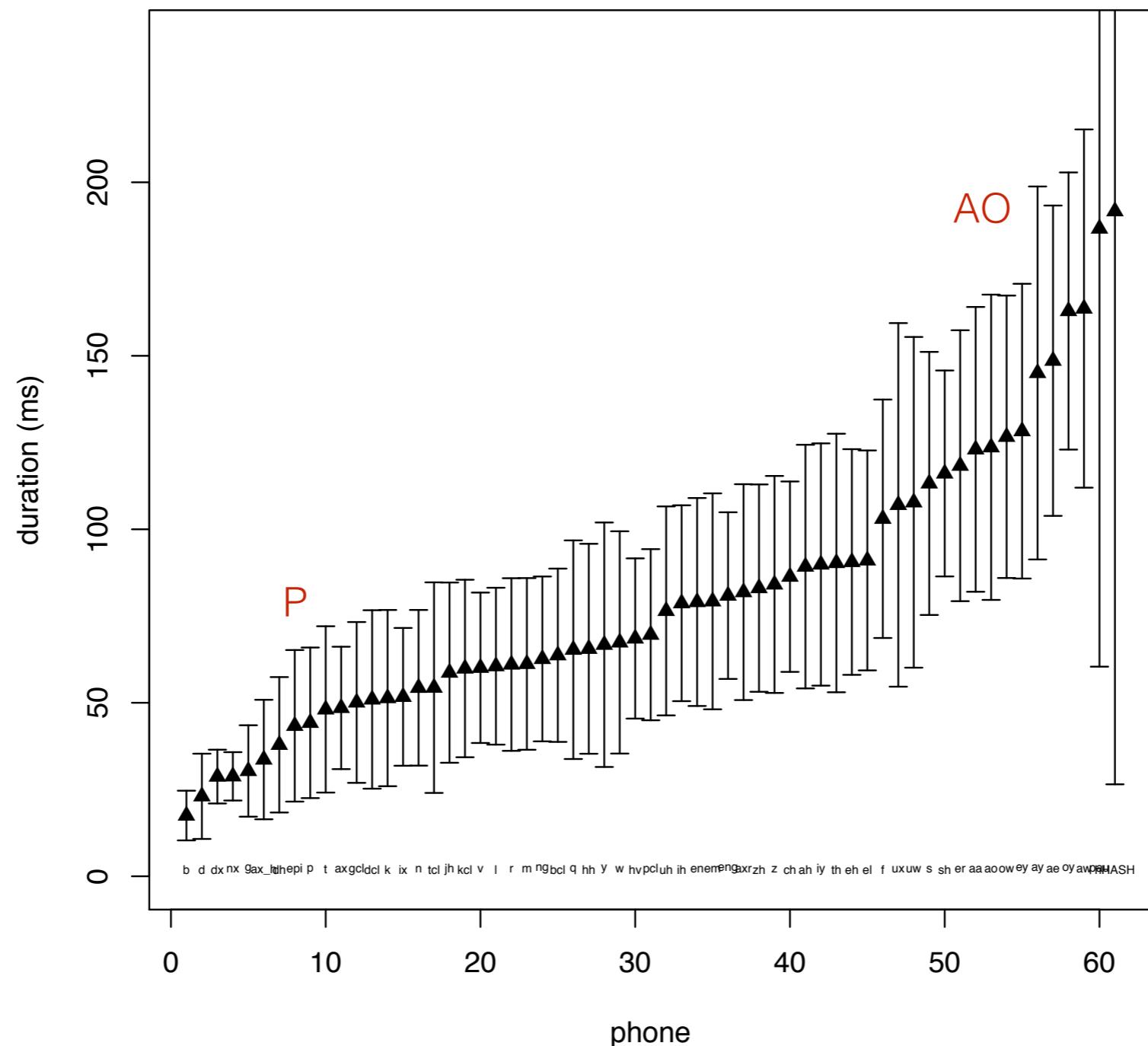
Text: *Please call Stella. Ask her to bring these things with her from the store.*

http://groups.linguistics.northwestern.edu/documentation/images/praat_aligned.jpg

Connectionist Temporal Classification (CTC)

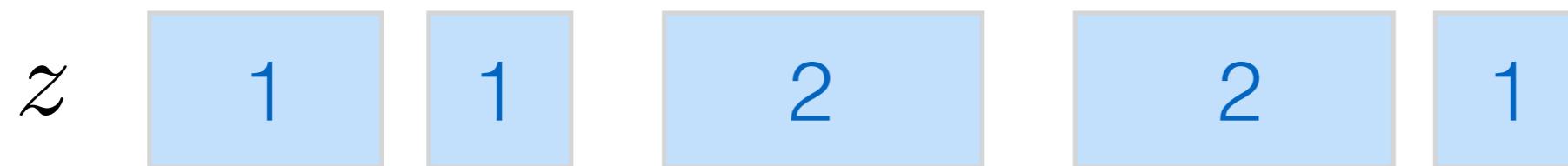


Duration Features



Segmental Recurrent Neural Networks (SRNNs)

x 1997 年 继 续 主 办 。



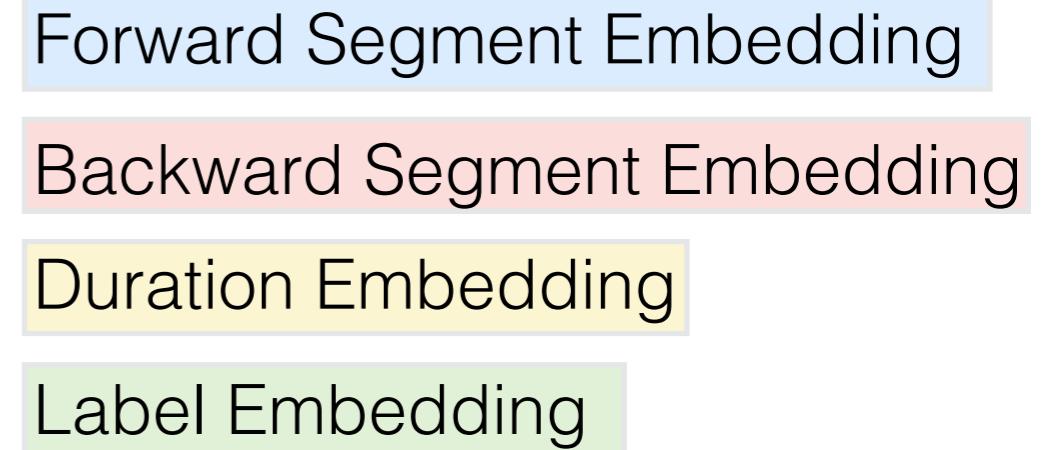
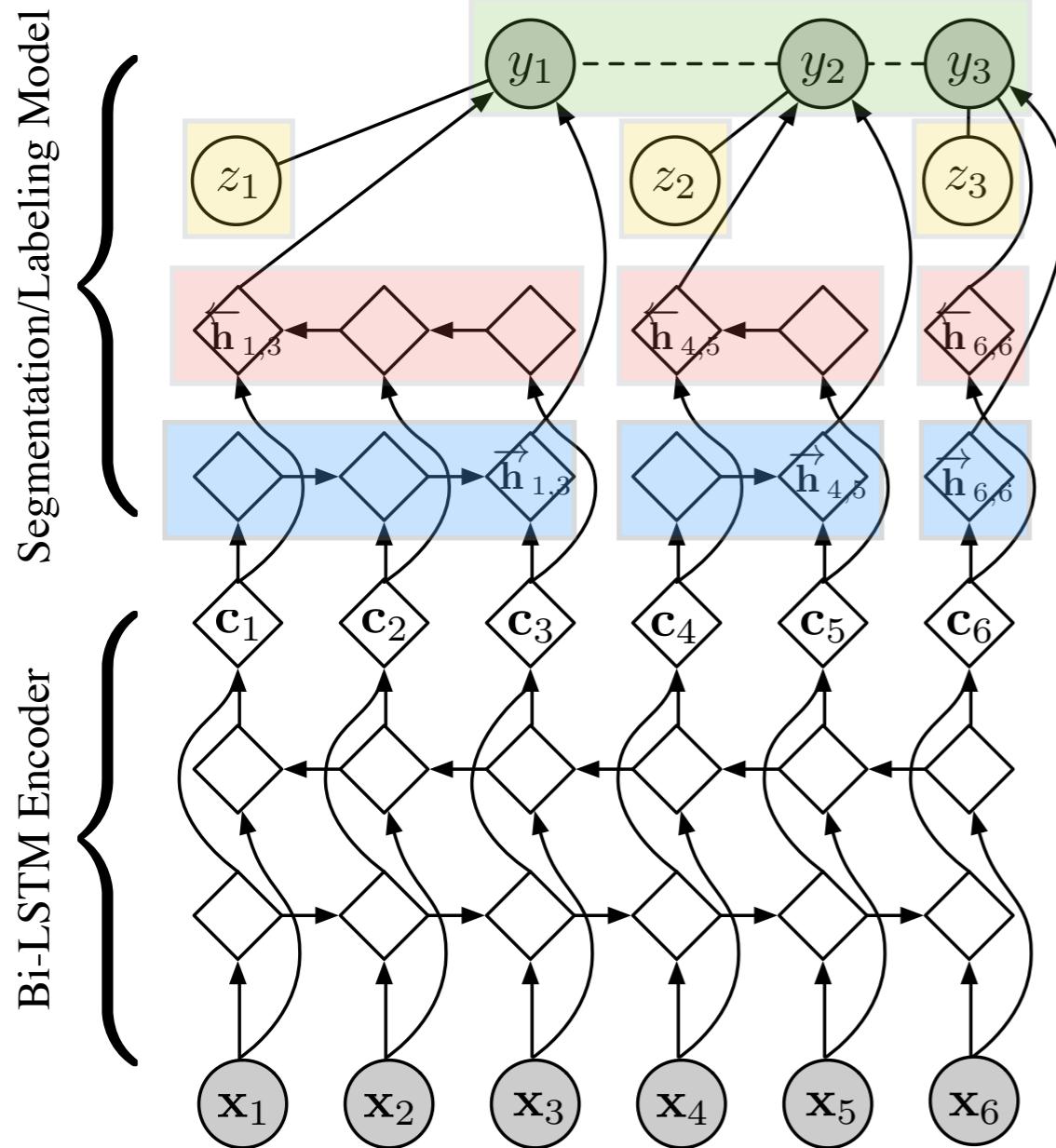
y CD M VV VV PU

SRNNs — $p(y, z|x)$

$$y^* = \arg \max_y \sum_z p(y, z | x)$$

$$\approx \arg \max_y \max_z p(y, z | x)$$

Segmental Recurrent Neural Networks (SRNNs)



$$p(\mathbf{y}, \mathbf{z} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{i=1}^{|\mathbf{y}|} \exp f(y_{i-k:i}, z_i, \mathbf{x})$$

$$f(y_{i-k:i}, z_i, \mathbf{x}_{s_i:s_i+z_i-1}) = \mathbf{w}^\top \phi(\mathbf{V}[\mathbf{g}_y(y_{i-k}); \dots; \mathbf{g}_y(y_i); \mathbf{g}_z(z_i); \overrightarrow{\text{RNN}}(\mathbf{c}_{s_i:s_i+z_i-1}); \overleftarrow{\text{RNN}}(\mathbf{c}_{s_i:s_i+z_i-1})] + \mathbf{a}) + b$$

Parameter Learning

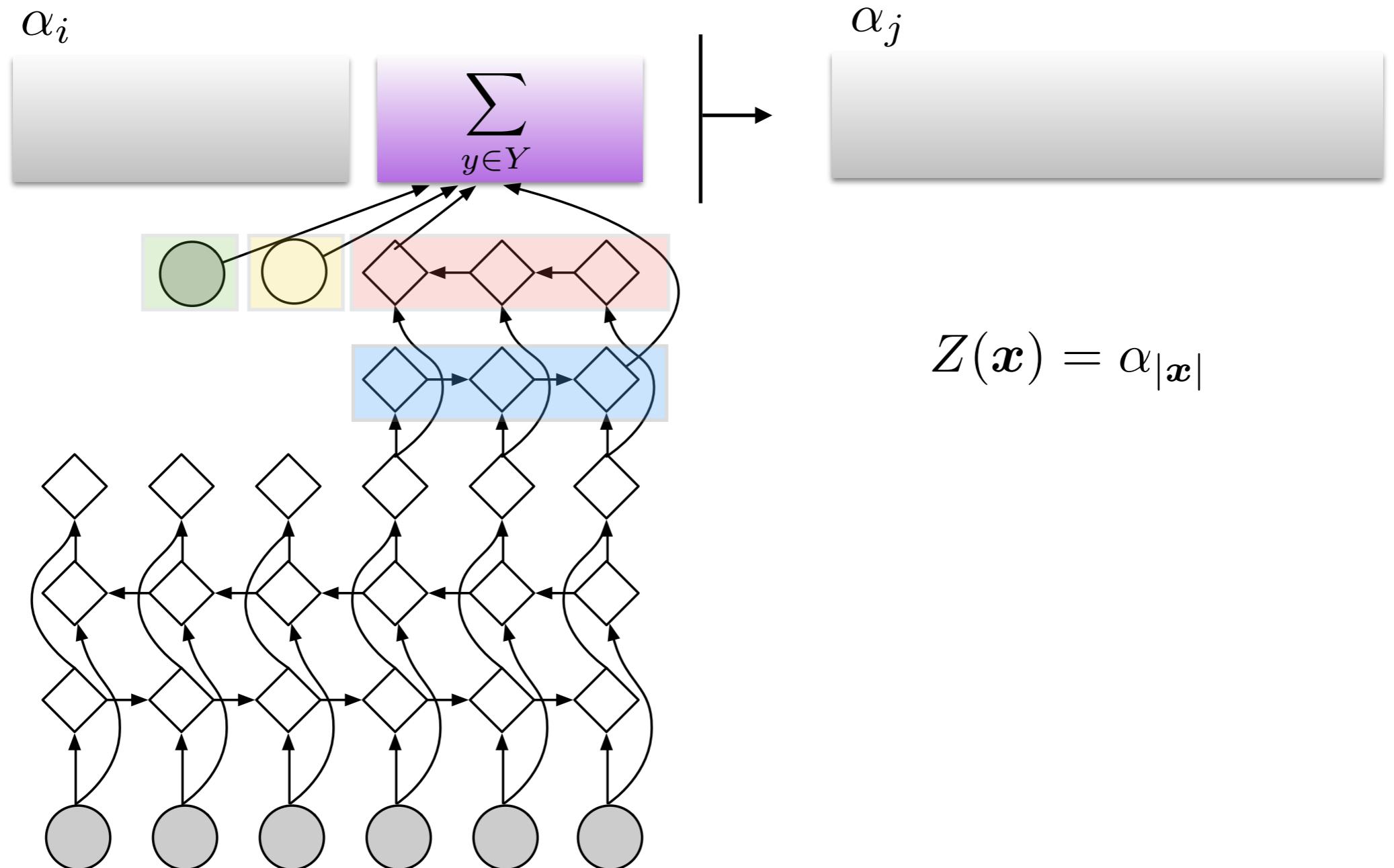
Fully Supervised

$$\begin{aligned}\mathcal{L} &= \sum_{(\mathbf{x}, \mathbf{y}, \mathbf{z}) \in \mathcal{D}} -\log p(\mathbf{y}, \mathbf{z} \mid \mathbf{x}) \\ &= \sum_{(\mathbf{x}, \mathbf{y}, \mathbf{z}) \in \mathcal{D}} \log Z(\mathbf{x}) - \log Z(\mathbf{x}, \mathbf{y}, \mathbf{z})\end{aligned}$$

Partially Supervised

$$\begin{aligned}\mathcal{L} &= \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} -\log p(\mathbf{y} \mid \mathbf{x}) \\ &= \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \sum_{\mathbf{z} \in \mathcal{Z}(\mathbf{x}, \mathbf{y})} -\log p(\mathbf{y}, \mathbf{z} \mid \mathbf{x}) \\ &= \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \overline{\log Z(\mathbf{x}) - \log Z(\mathbf{x}, \mathbf{y})}\end{aligned}$$

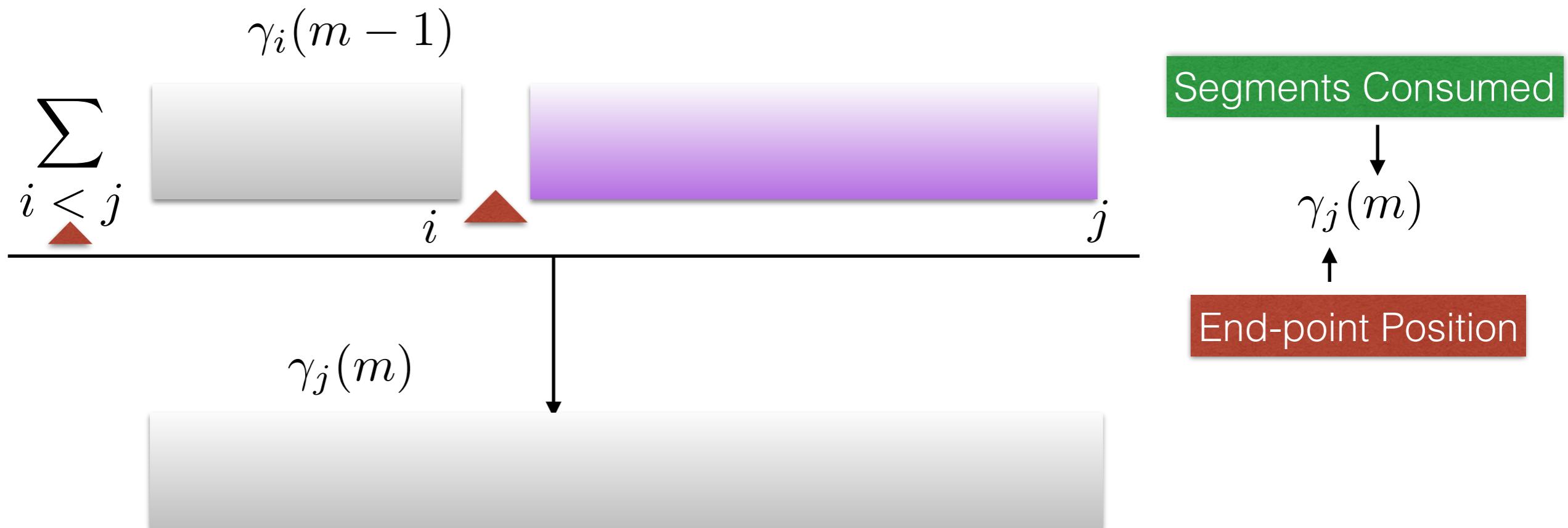
Dynamic Programming



$$\alpha_0 = 1$$

$$\alpha_j = \sum_{i < j} \alpha_i \times \sum_{y \in Y} \exp \mathbf{w}^\top \phi(\mathbf{V}[\mathbf{g}_y(y); \mathbf{g}_z(z_i)]; \overrightarrow{\text{RNN}}(\mathbf{c}_{s_i:s_i+z_i-1}); \overleftarrow{\text{RNN}}(\mathbf{c}_{s_i:s_i+z_i-1}) + \mathbf{a}) + b$$

Dynamic Programming

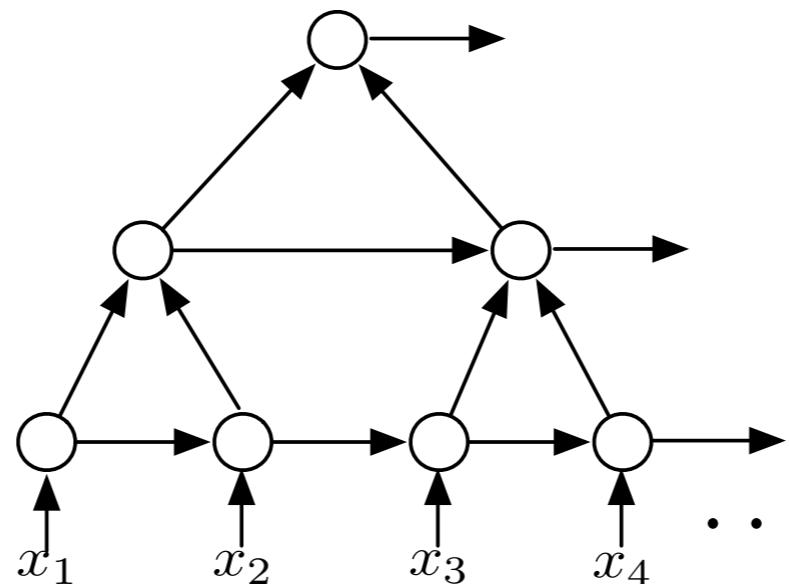


$$Z(\mathbf{x}, \mathbf{y}) = \gamma_{|\mathbf{x}|}(|\mathbf{y}|)$$

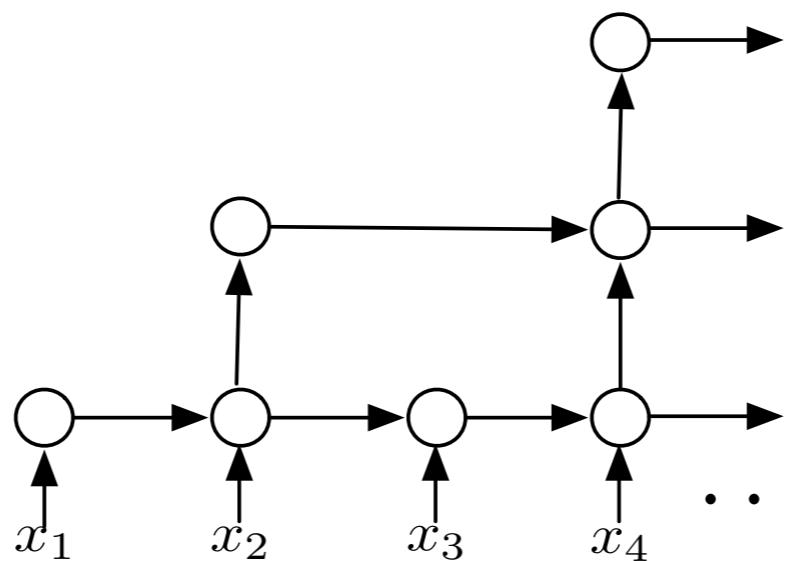
$$\gamma_j(m) = \sum_{i < j} \gamma_i(m-1) \times \left(\exp \mathbf{w}^\top \phi(\mathbf{V}[g_y(y_i); g_z(z_i); \overrightarrow{\text{RNN}}(\mathbf{c}_{s_i:s_i+z_i-1}); \overleftarrow{\text{RNN}}(\mathbf{c}_{s_i:s_i+z_i-1})] + \mathbf{a}) + b \right)$$

Further Speedup

Subsampling



a) concatenate / add



b) skip

Experiments

Online Hand Writing Recognition



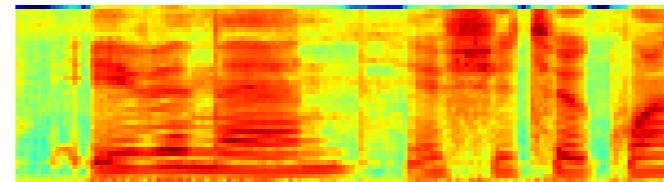
	P (seg)	R (seg)	F (seg)	Error
CTC	-	-	-	13.8%
SRNNs(Full)	98.9%	98.6%	98.6%	5.4%
SRNNs (Partial)	99.2%	99.1%	99.2%	2.7%

(Kassel, 1995)

(Taskar et al. 2004)

Experiments

Speech Recognition



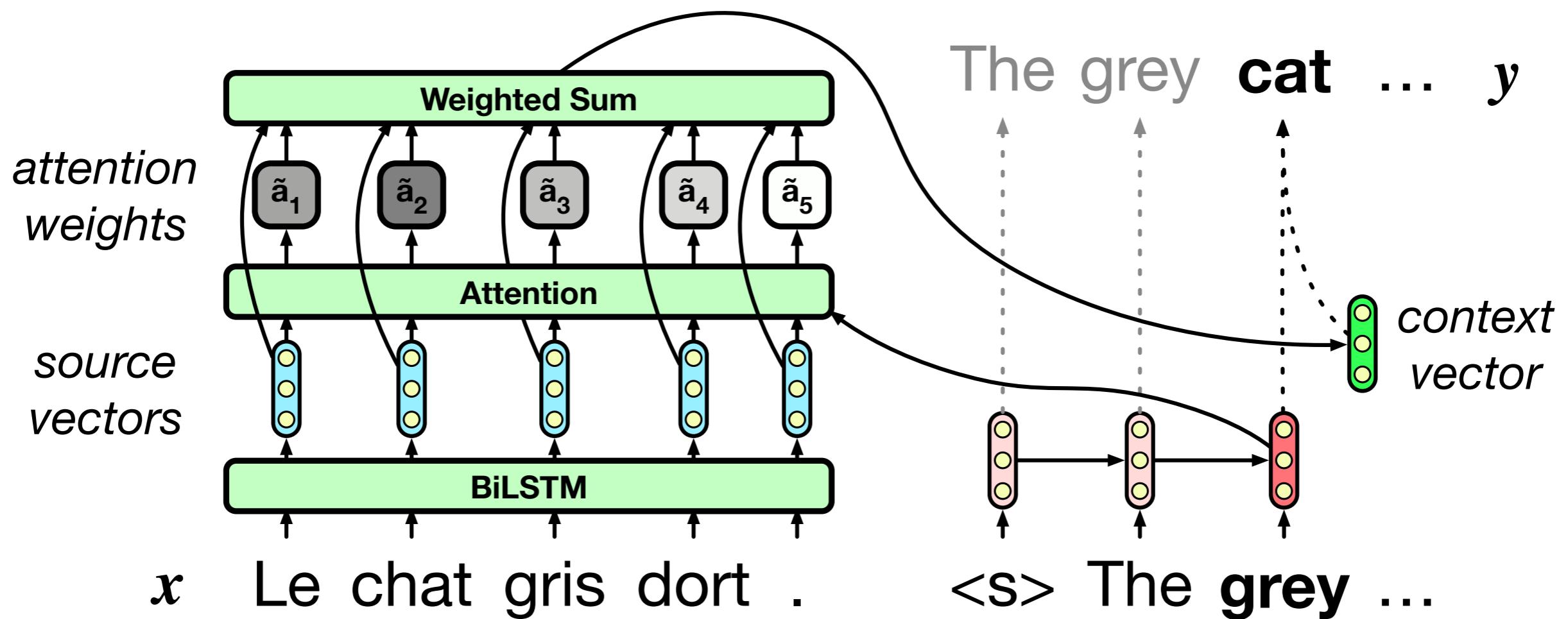
*Multi-task Learning with CTC and Segmental CRF
for Speech Recognition* [[INTERSPEECH 2016](#)]

Segmental Recurrent Neural Networks for End-to-end Speech Recognition [[INTERSPEECH 2017](#)]

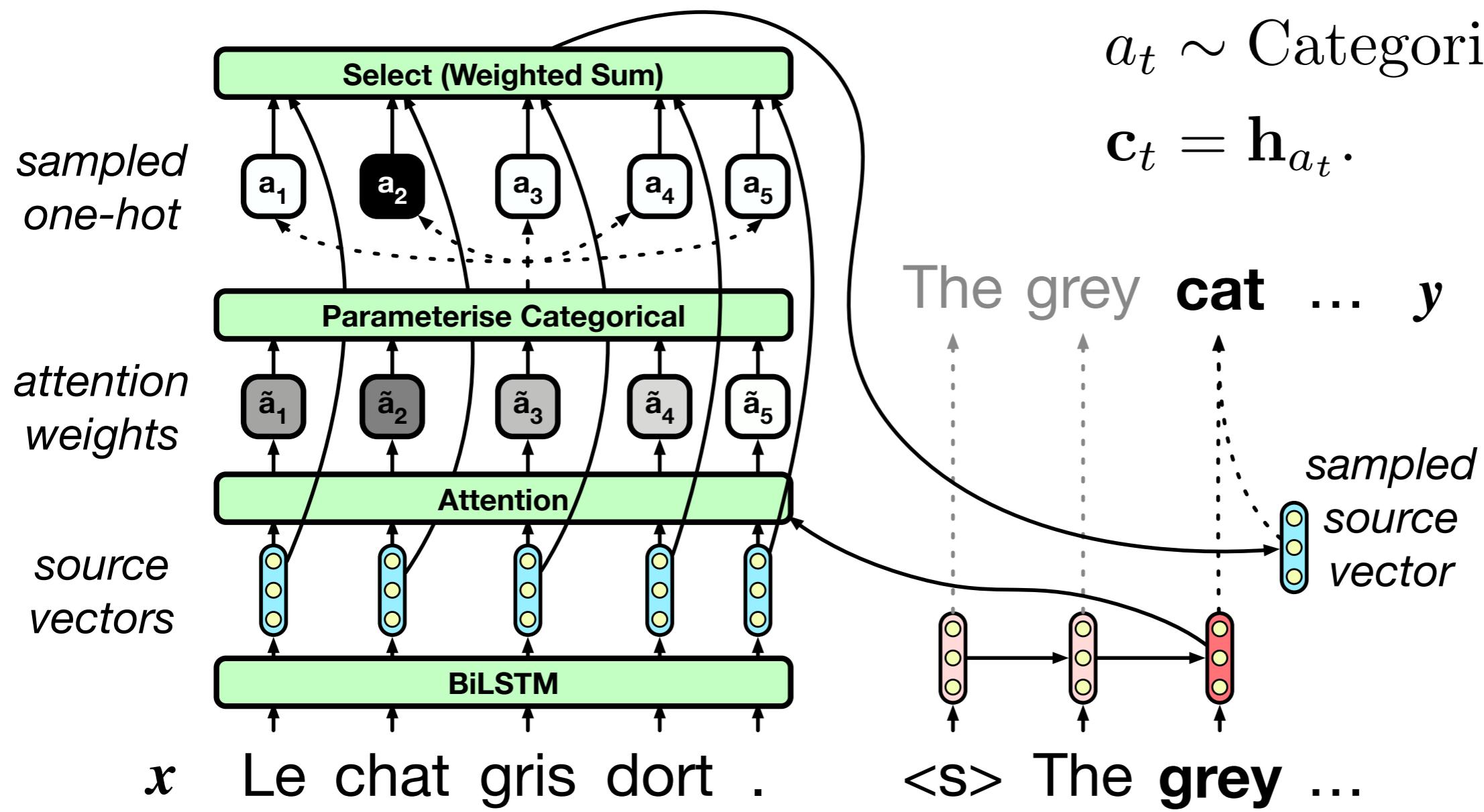
Outline

- Introduction
- Part I — Segmental Recurrent Neural Networks
[ICLR 2016, INTERSPEECH 2016, 2017]
- Part II — Inference and Regularization in Sequence to Sequence Models with Stochastic Attention
- Conclusion

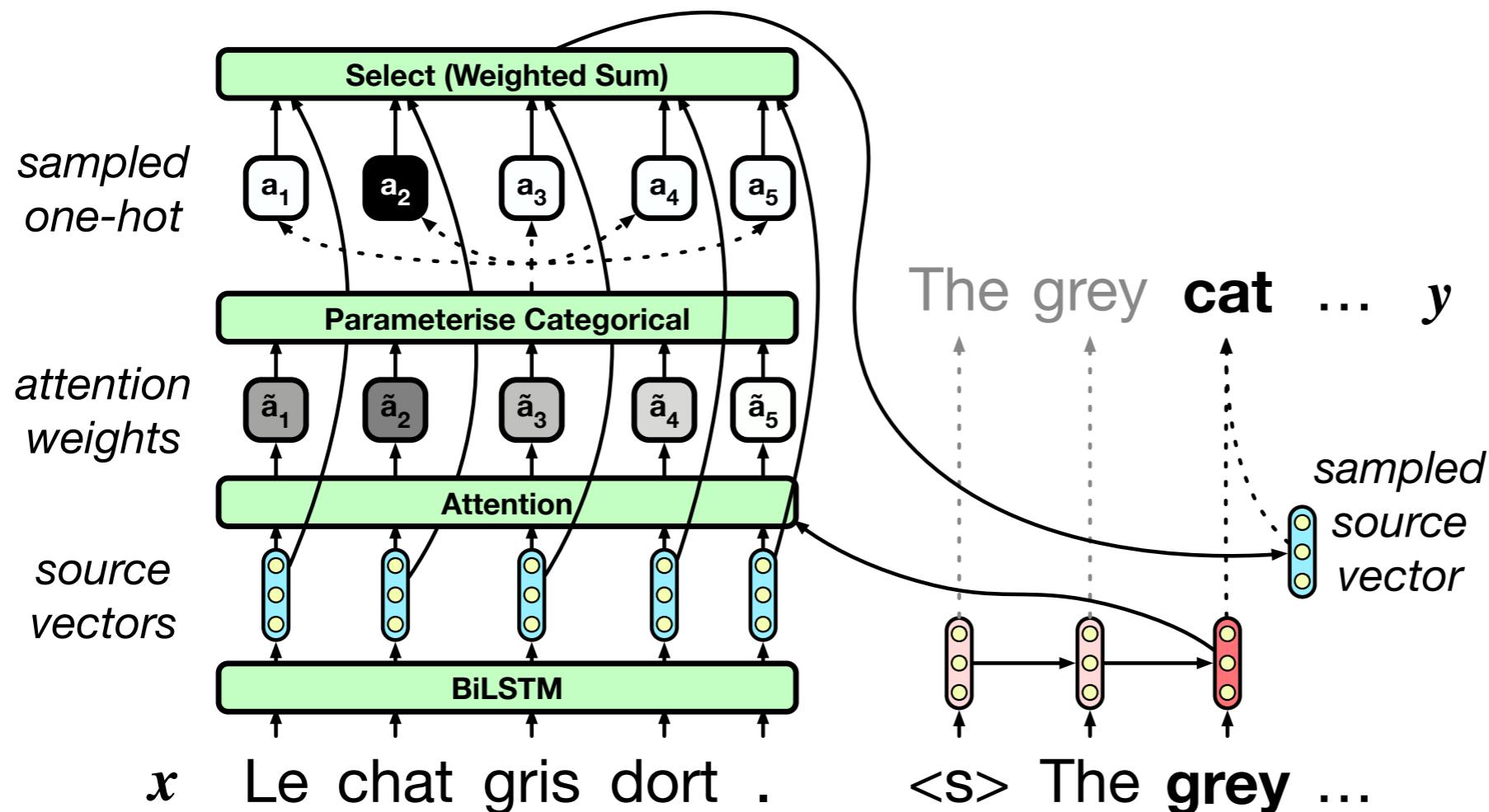
Deterministic v.s. Stochastic Attention



Deterministic v.s. Stochastic Attention



Marginal Likelihood and Training Objective



$$p(\mathbf{y} \mid \mathbf{x}) = \prod_{t=1}^M \sum_{a_t=1}^N p(a_t \mid \mathbf{x}, \mathbf{y}_{<t}) p(y_t \mid a_t, \mathbf{x}, \mathbf{y}_{<t})$$

Approximating the Marginal Likelihood

Variational lower bound:

$$\begin{aligned}\log p(\mathbf{y}|\mathbf{x}) &= \sum_t \log \sum_a p(a|\mathbf{x}, \mathbf{y}_{<t}) p(y_t|a, \mathbf{x}, \mathbf{y}_{<t}) \\ &= \sum_t \log \sum_a p(y_t, a|\mathbf{x}, \mathbf{y}_{<t}) \\ &\leq \sum_t \mathbb{E}_{q(a)} \log \frac{p(y_t, a|\mathbf{x}, \mathbf{y}_{<t})}{q(a)} \\ &= \sum_t \mathbb{E}_{q(a)} \log p(y_t, a|\mathbf{x}, \mathbf{y}_{<t}) - H(q)\end{aligned}$$

Approximating the Marginal Likelihood

Variational lower bound:

$$\begin{aligned}\log p(\mathbf{y}|\mathbf{x}) &= \sum_t \log \sum_a p(a|\mathbf{x}, \mathbf{y}_{<t}) p(y_t|a, \mathbf{x}, \mathbf{y}_{<t}) \\ &= \sum_t \log \sum_a p(y_t, a|\mathbf{x}, \mathbf{y}_{<t}) \\ &\leq \sum_t \mathbb{E}_{q(a)} \log \frac{p(y_t, a|\mathbf{x}, \mathbf{y}_{<t})}{q(a)} \\ &= \sum_t \mathbb{E}_{q(a)} \log p(y_t, a|\mathbf{x}, \mathbf{y}_{<t}) - H(q)\end{aligned}$$

REINFORCE:

$$\log p(\mathbf{y}|\mathbf{x}) = \sum_t \log \sum_a p(a|\mathbf{x}, \mathbf{y}_{<t}) p(y_t|a, \mathbf{x}, \mathbf{y}_{<t})$$

One-sample approximation

Experiments: Deterministic vs. Stochastic Attention

Hypothesis: stochasticity in memory access (rather than mean-field-like approximations) results in easier learning problems.

Model	Inference	BLEU	PPL
Deterministic	-	31.87	5.25
Stochastic	exact	31.91	4.65
Stochastic	variational	30.10	5.40
Stochastic	REINFORCE	29.85	5.31

Punchline: stochasticity doesn't really help us.

Stochastic Attention

- **Let's not give up yet!**
 - Neural nets can fit anything ([Zhang, Bengio, Hardt, Recht, Vinyals, ICLR 2017](#)).
 - (Stochastic) attention should be sensible, not just a random fit
 - Let's **regularize** the posterior distributions so they look more like what we expect posteriors to be ([Ganchev, Graça, Gillenwater, Taskar, JMLR 2010](#))
- **Strategy:**
 - Apply KL penalty (true PR penalty)
 - Use variants of IS (biased estimator) using the expected posterior as the instrumental distribution

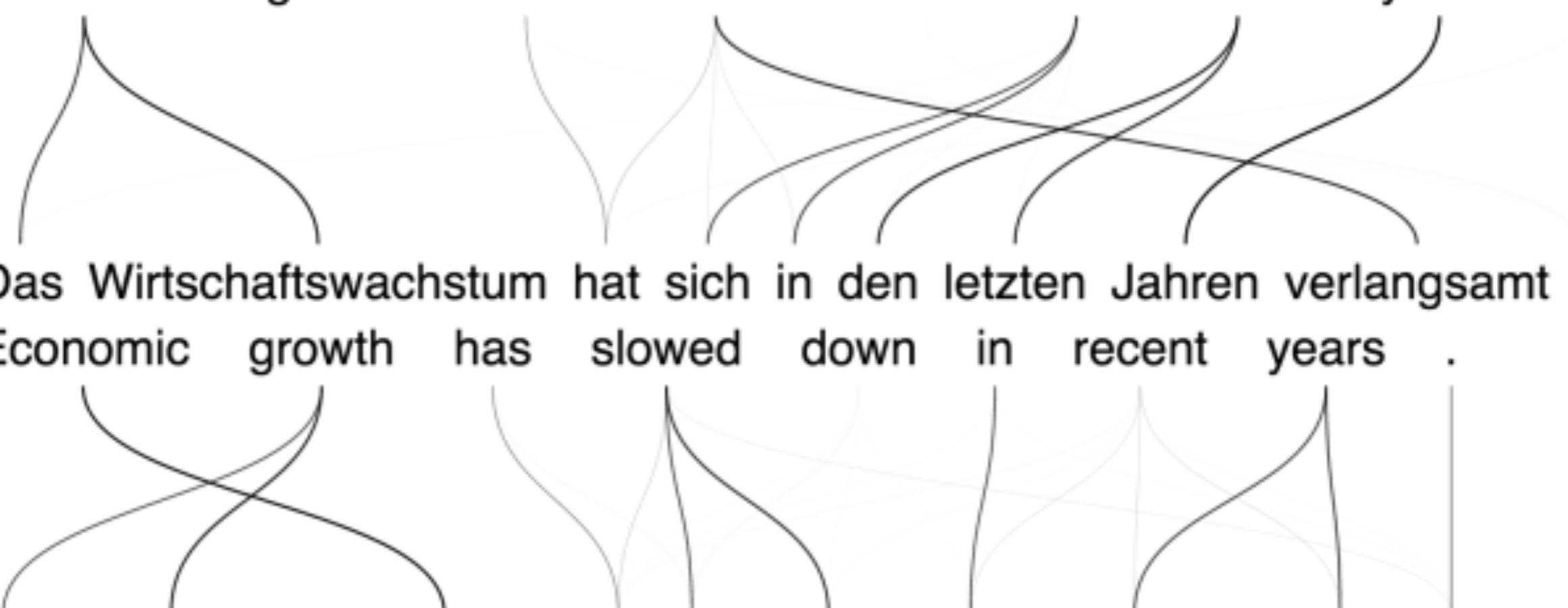
IBM Models

Economic growth has slowed down in recent years .

Das Wirtschaftswachstum hat sich in den letzten Jahren verlangsamt .

Economic growth has slowed down in recent years .

La croissance économique s' est ralentie ces dernières années .



Great q distribution!

Posterior Regularization

Exact:

$$\sum_t \log \sum_a \underbrace{p(a|x, y_{<t})}_{+ \gamma \times \text{KL}[p(a|y_t, x, y_{<t}) || q(a)]} p(y_t|a, x, y_{<t})$$

Posterior Regularization

Importance Sampling:

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}) &= \prod_t \sum_a p(a|\mathbf{x}, \mathbf{y}_{<t}) p(y_t|a, \mathbf{x}, \mathbf{y}_{<t}) \\ &= \prod_t \sum_a p(y_t, a|\mathbf{x}, \mathbf{y}_{<t}) \\ &= \prod_t \sum_a q(a) w(a, \mathbf{x}, \mathbf{y}) \\ &= \prod_t \mathbb{E}_{q(a)} w(a, \mathbf{x}, \mathbf{y}) \quad \text{Monte Carlo approximation} \end{aligned}$$

$$w(a, \mathbf{x}, \mathbf{y}) = \frac{p(y_t, a|\mathbf{x}, \mathbf{y}_{<t})}{q(a)}$$

Posterior Regularization

Jensen IS:

$$\begin{aligned}\log p(\mathbf{y}|\mathbf{x}) &= \sum_t \log \sum_a p(a|\mathbf{x}, \mathbf{y}_{<t}) p(y_t|a, \mathbf{x}, \mathbf{y}_{<t}) \\ &= \sum_t \log \sum_a p(y_t, a|\mathbf{x}, \mathbf{y}_{<t}) \\ &\leq \sum_t \mathbb{E}_{q(a)} \log \frac{p(y_t, a|\mathbf{x}, \mathbf{y}_{<t})}{q(a)} \\ &= \sum_t \mathbb{E}_{q(a)} \log p(y_t, a|\mathbf{x}, \mathbf{y}_{<t}) - H(q)\end{aligned}$$

Importance Sampling

$$w(a, \mathbf{x}, \mathbf{y}) = \frac{p(y_t, a|\mathbf{x}, \mathbf{y}_{<t})}{q(a)}$$

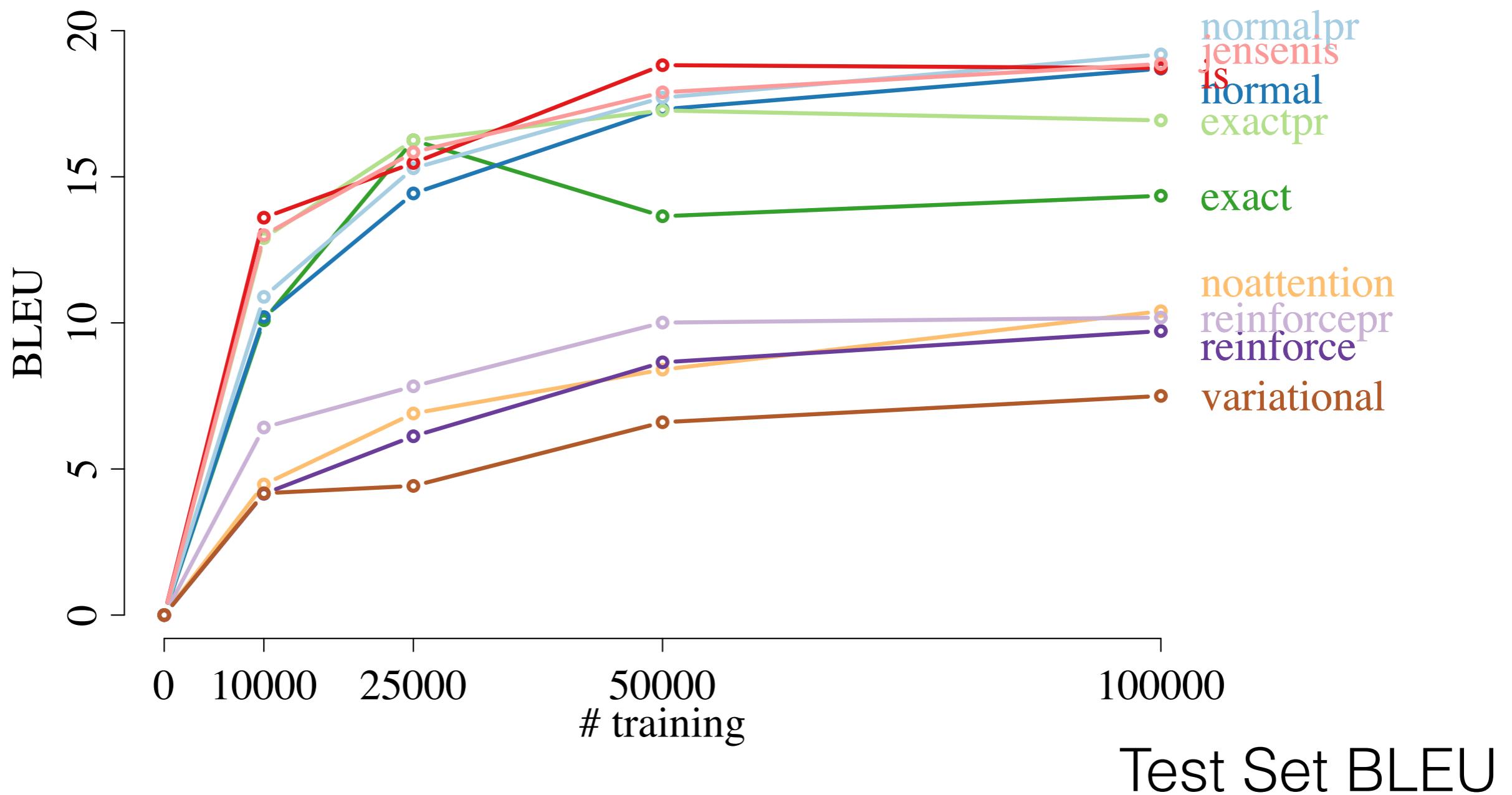
sample from a fixed $q(a)$

Experiments: Posterior regularization

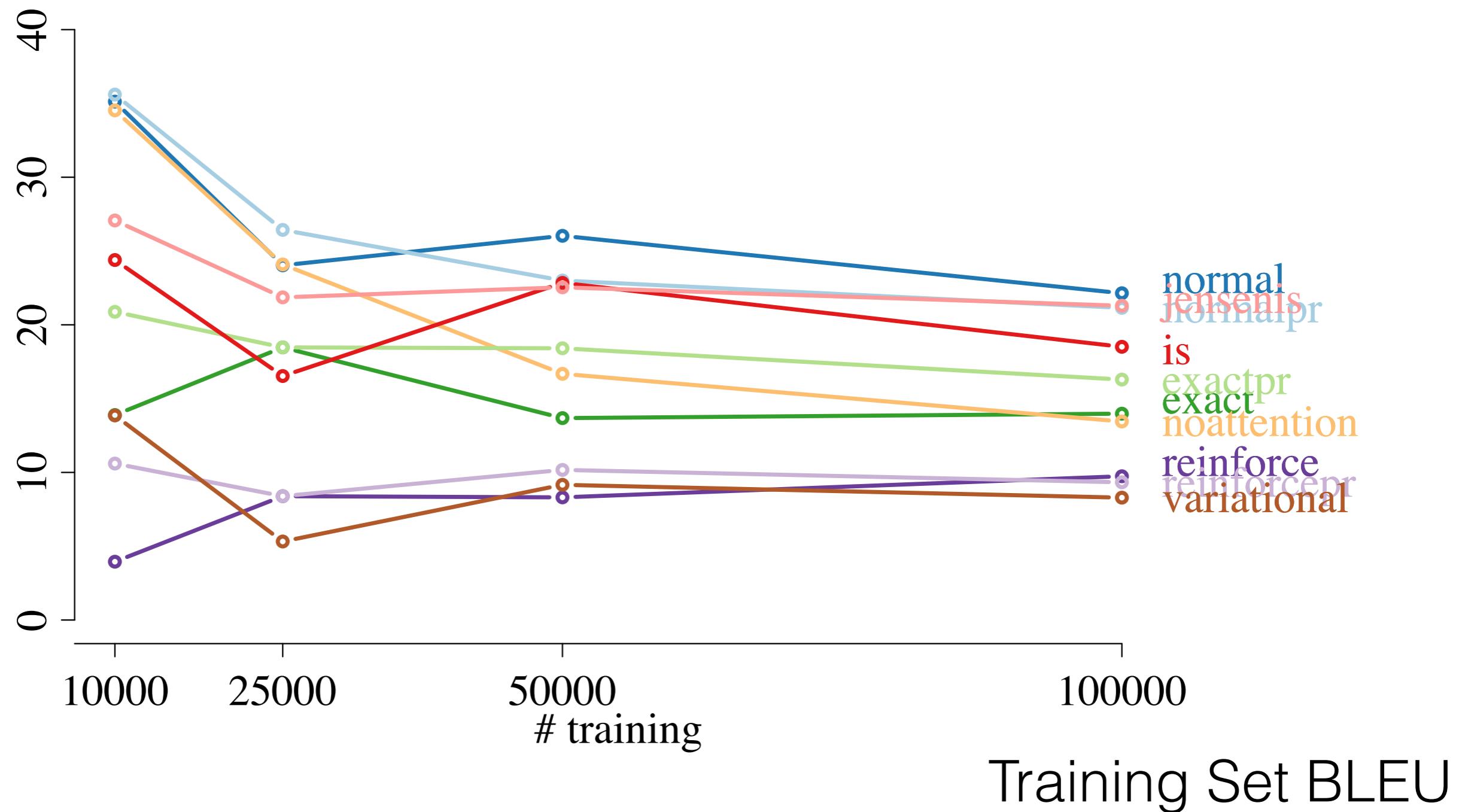
Hypothesis: neural nets can fit anything (Zhang, Bengio, Hardt, Recht, Vinyals).
Let's regularize posterior distributions (Ganchev, Graca, Gillenwater, Taskar).

Model	Inference	PR	BLEU	PPL
Deterministic	exact	none	31.87	5.25
Stochastic	exact	none	31.91	4.65
Deterministic	exact	full	32.48	5.20
Stochastic	exact	full	35.17	4.03
Stochastic	IS with q	approximate	34.68	4.04
Stochastic	Jensen bound IS with q	approximate	35.40	3.94

Experiments: Different training set size



Experiments: Different training set size



Conclusion

- The combination of explicit structure representations and learned distributed representations
 - Part I — Segmental Recurrent Neural Networks [ICLR 2016, INTERSPEECH 2016, 2017]
 - Part II — Inference and Regularization in Sequence to Sequence Models with Stochastic Attention

Thank you!