# Semi-Supervised Dependency Parsing with Variational Autoencoder
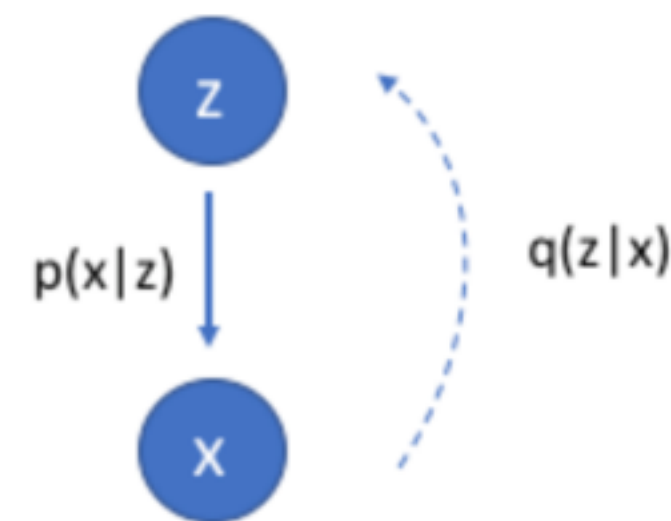
Wang Ge
2019/3/20

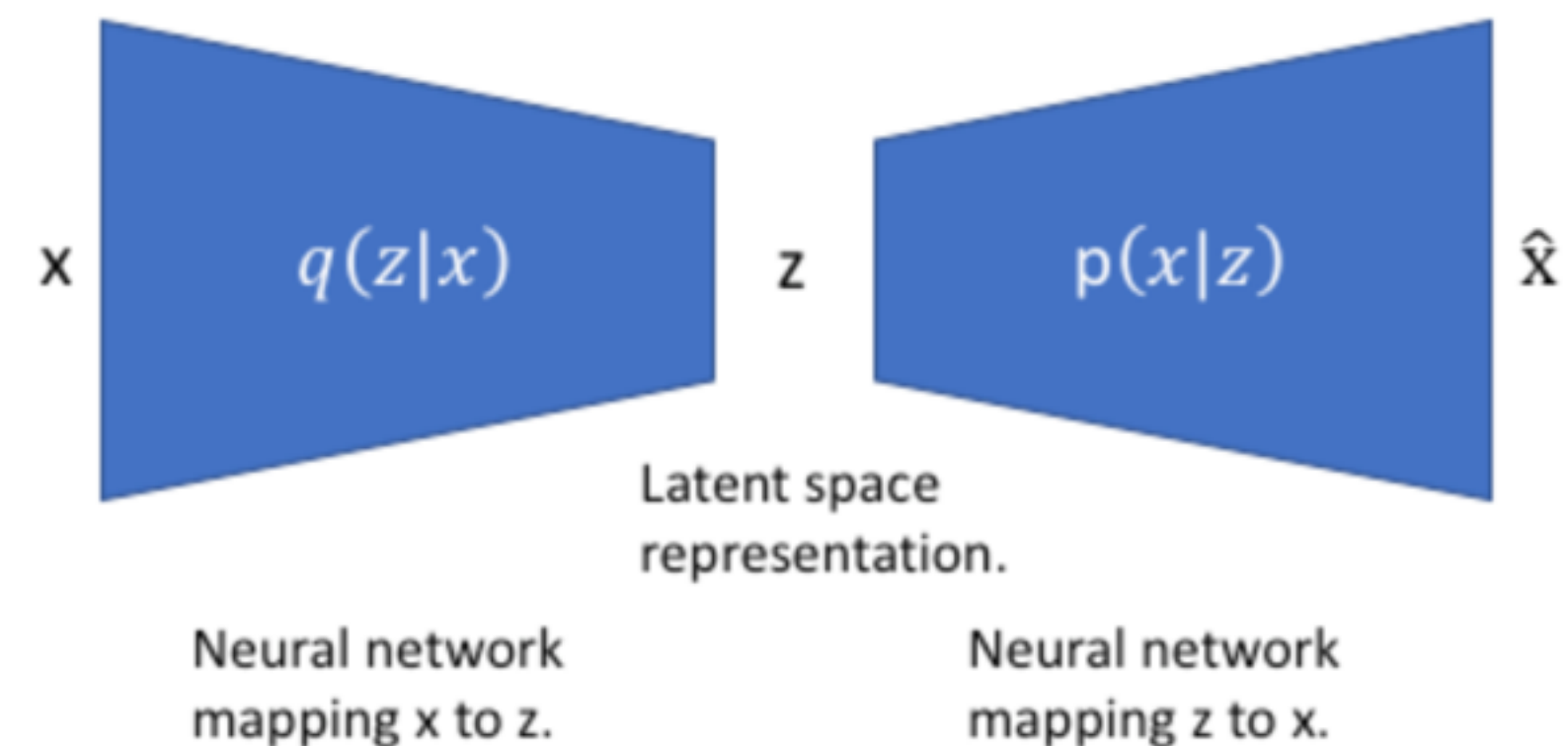# Why is Semi-Supervised Learning Necessary

- Learning a reliable dependency parser requires a large amount of labelled data, e.g. Penn Treebank (39000+ sentence for training).

- Labelling data with dependency parse trees is a challenging and laborious work. For many languages, enough labelled data are not available.

- Unsupervised methods require POS tag information, lacking robustness over different annotation criteria such as WSJ and UD

# Variational Autoencoder
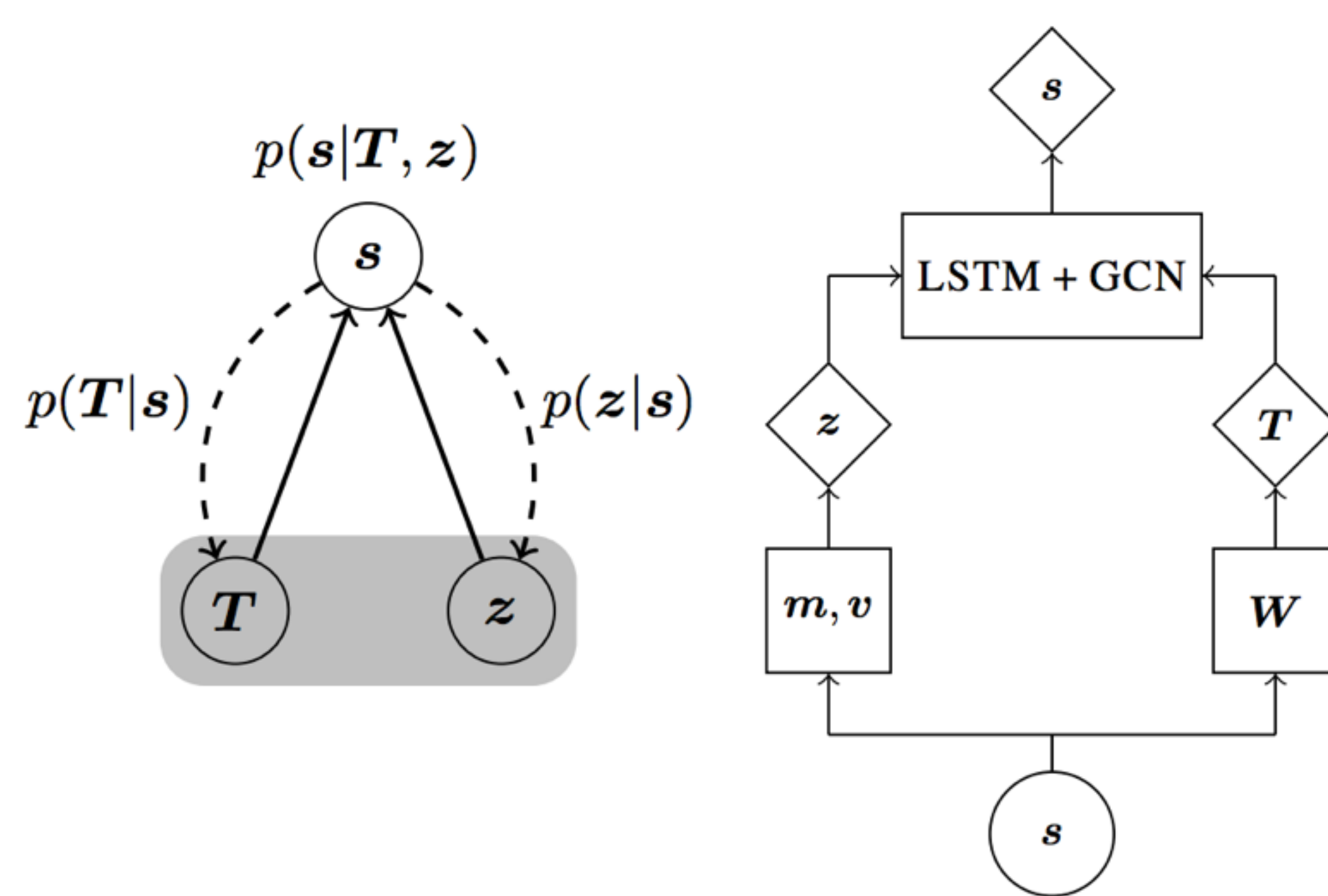
- An elegant generative model with latent variables



q(z|x)

p(x|z)

We'd like to use our observations to understand the hidden variable.

X    q(z|x)    z    p(x|z)    X̂

Latent space representation.

Neural network mapping x to z.
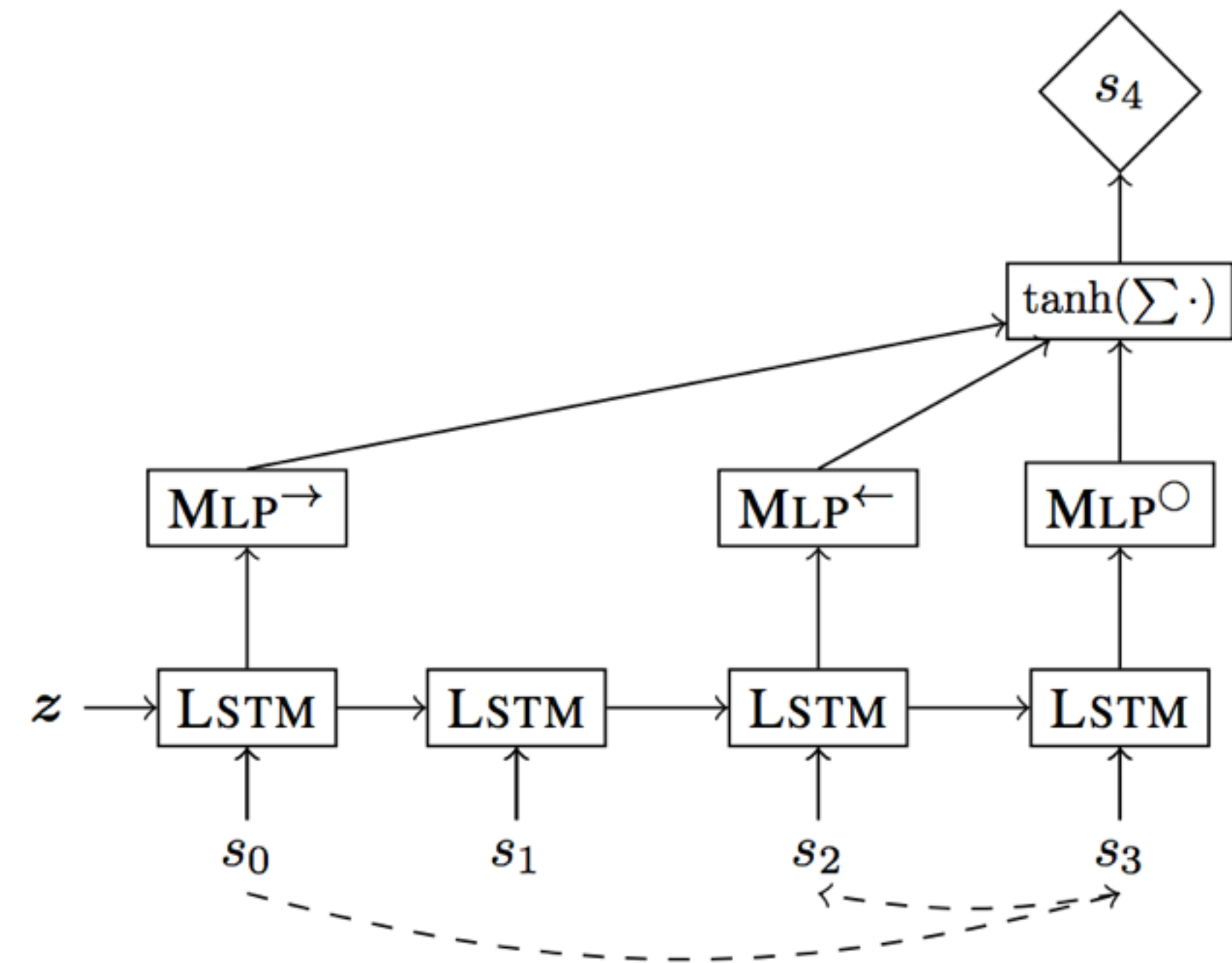
Neural network mapping z to x.

- Easy to be applied to semi-supervised learning, not suitable for structured latent variable like parse tree

# Recent Work Combining Semi-Supervised Parsing with VAE



Graphical Model and Computational Graph

Decoder Implementation

# Model Detail

- Generative Story

$$T \sim p(T|n) \qquad z \sim p(z|n) \qquad s \sim p(s|T,z,n)$$

- Objective Function

$$\log p_\theta(s) = \mathbb{E}_{q_\phi(T,z|s)}[\log p_\theta(s|T,z)] - \mathrm{KL}[q_\phi(T,z|s)|p(T,z)] + \mathrm{KL}\left[q_\phi(T,z|s)\|p_\theta(T,z|s)\right]$$

$$\log p_\theta(s) \geq \mathbb{E}_{q_\phi(T,z|s)}[\log p_\theta(s|T,z)] - \mathrm{KL}[q_\phi(T,z|s)|p(T,z)] = \tilde{\mathcal{E}}_{\theta,\phi}(s)$$

$$\bar{\mathcal{E}}_{\theta,\phi}(s,T) = \mathbb{E}_{q_\phi(z|s)}[\log p_\theta(s|T,z)] - \mathrm{KL}[q_\phi(z|s)|p(z)]$$

$$\bar{\mathcal{E}}_{\theta,\phi}(s,T) = \mathbb{E}_{q_\phi(z|s)}[\log p_\theta(s|T,z)] - \mathrm{KL}[q_\phi(z|s)|p(z)]$$

# Encoder and Decoder

- For the encoder:

$$\boldsymbol{W} = \text{DEPWEIGHTS}(\boldsymbol{s}) \qquad\qquad \boldsymbol{m}, \log \boldsymbol{v}^2 = \text{EMBPARAMS}(\boldsymbol{s})$$

$$q_\phi(\boldsymbol{T}|\boldsymbol{s}) = \frac{\exp(\sum_{i,j} W_{i,j} T_{i,j})}{\sum_{\boldsymbol{T}'} \exp(\sum_{i,j} W_{i,j} T'_{i,j})} \qquad\qquad q_\phi(\boldsymbol{z}|\boldsymbol{s}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{m}, \boldsymbol{v})$$

- For the decoder:

$$\boldsymbol{g}^i = \tanh\left( \text{MLP}^{\circ}(\boldsymbol{o}^i) + \sum_{h=0}^{i-1} T_{h,i} \times \text{MLP}^{\frown}(\boldsymbol{o}^h) + \sum_{m=0}^{i-1} T_{i,m} \times \text{MLP}^{\frown}(\boldsymbol{o}^m) \right)$$

# Differentiable Perturb-and-Parse

- Sampling dependency tree by perturbing weight matrix

$$\boldsymbol{W} = \textsc{EmbParams}(\boldsymbol{s})$$
$$\boldsymbol{P} \sim \mathcal{G}(0, 1)$$
$$\boldsymbol{T} = \textsc{Eisner}(\boldsymbol{W} + \boldsymbol{P})$$

where the perturbation is sampled from Gumbel distribution

- Sampling dependency tree by perturbing weight matrix

Replace the one-hot-argmax peaked-softmax

$$o_i = \mathbb{1}[\forall 1 \leq j \leq k, j \neq i : v_i > v_j]$$ with $$o_i = \frac{\exp(^1\!/_\tau \, v_i)}{\sum_{1 \leq j \leq k} \exp(^1\!/_\tau \, v_j)}$$

# Experiment Results

## (a) Parsing results

|  | English | French | Swedish |
|---|---|---|---|
| **Supervised** | 88.79 / 84.74 | 84.09 / 77.58 | 86.59 / 78.95 |
| **VAE w. $z$** | 89.39 / 85.44 | 84.43 / 77.89 | 86.92 / 80.01 |
| **VAE w/o $z$** | 89.50 / 85.48 | 84.69 / 78.49 | 86.97 / 79.80 |
| **Kipperwasser & Goldberg** | 89.88 / 86.49 | 84.30 / 77.83 | 86.93 / 80.12 |

## (b) Dependency length analysis

| **Distance** | **Supervised** <br> **Re / Pr** | **Semi-sup.** <br> **Re / Pr** |
|---|---|---|
| **(to root)** | 93.46 / **89.30** | 93.84 / **92.41** |
| **1** | 95.61 / 94.07 | 95.33 / 94.57 |
| **2** | 93.01 / 90.88 | 92.50 / 92.09 |
| **3 . . . 6** | 85.95 / 88.13 | 87.31 / 87.93 |
| **> 7** | **72.47** / 83.26 | **78.72** / 83.11 |

## (c) Dependency label analysis

| **Label** | **Supervised** <br> **Re / Pr** | **Semi-sup.** <br> **Re / Pr** |
|---|---|---|
| **mwe** | 75.58 / 81.25 | 90.70 / 84.78 |
| **advmod** | 87.27 / 85.95 | 87.32 / 87.51 |
| **appos** | 77.49 / 80.27 | 81.39 / 81.03 |