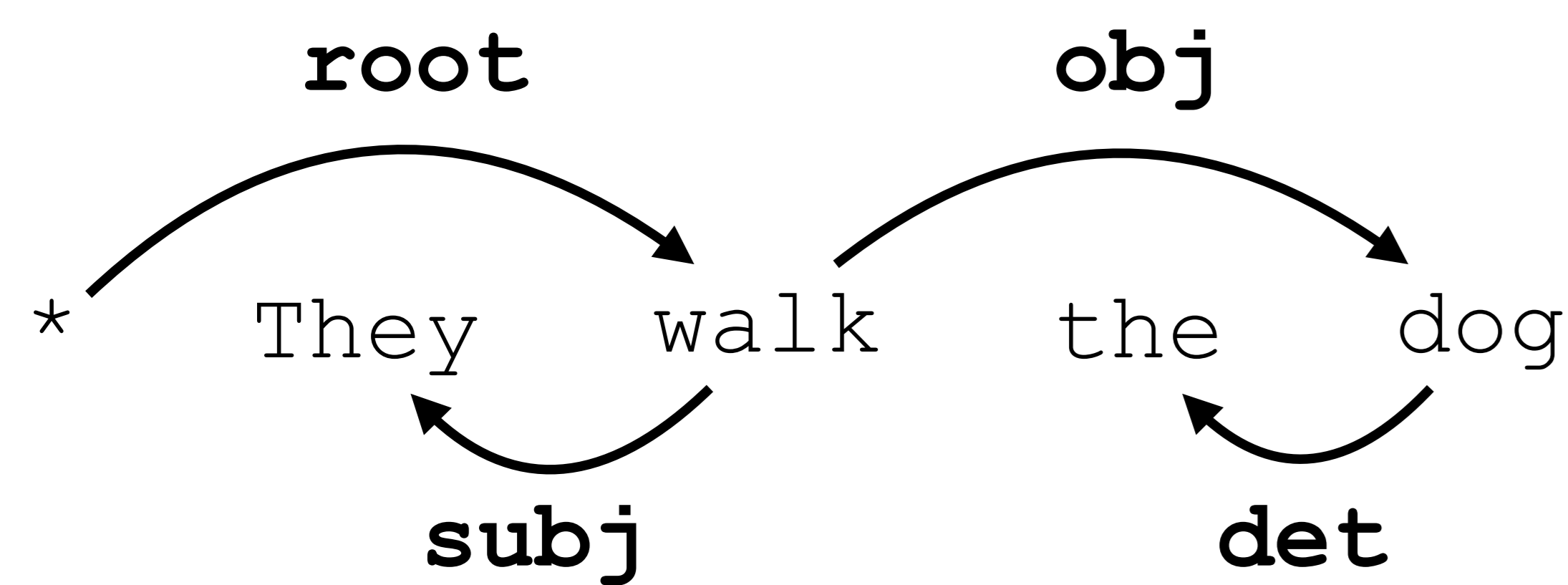


Dependency Tree

Syntactic structure that is useful in downstream tasks but annotation is expensive: datasets are small for many major languages (e.g. Vietnamese).



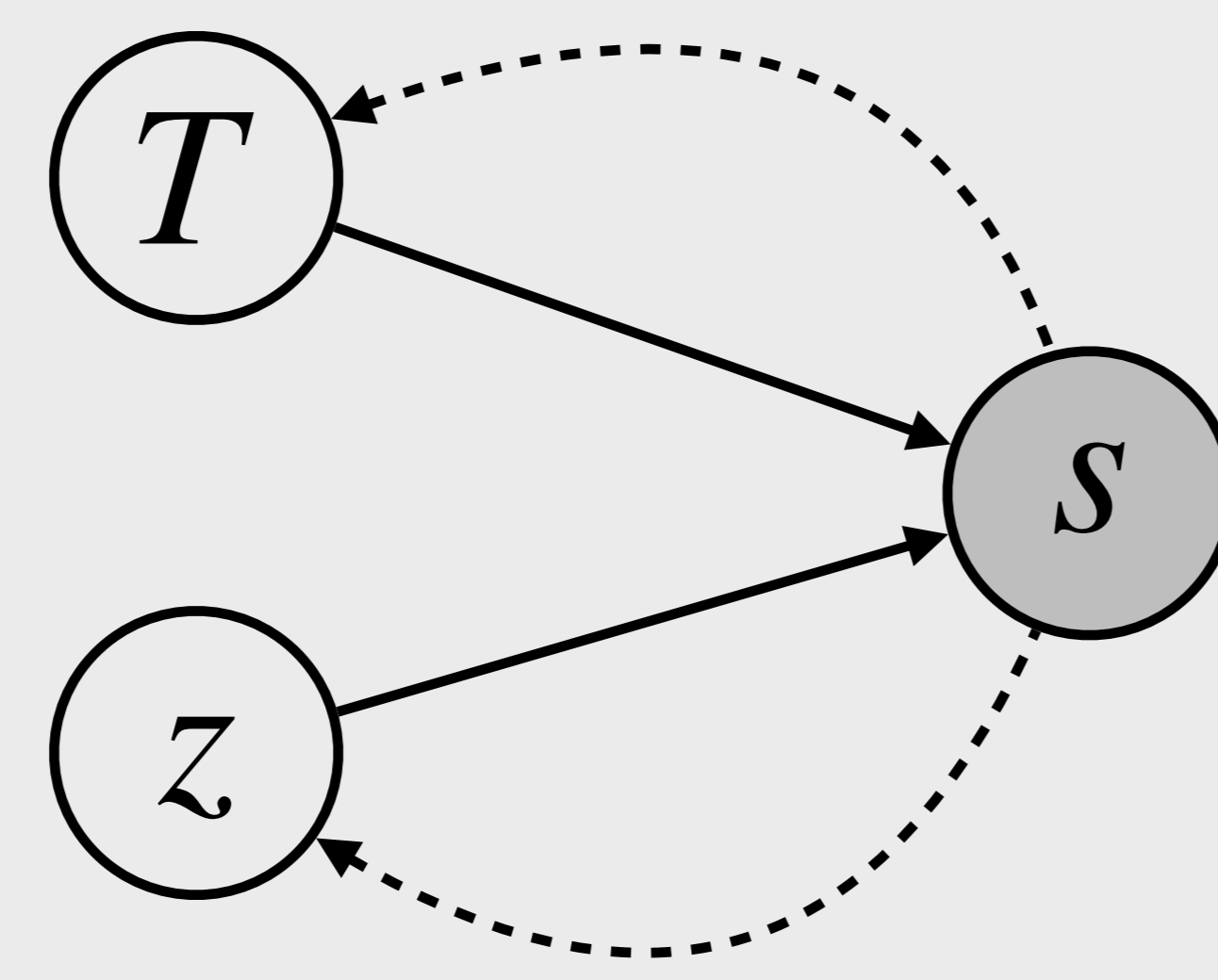
T : Tree described as an adjacency matrix

$\mathcal{T}(s)$: set of trees compatible with sentence s

W : matrix of arc weights computed with a NN

$$\operatorname{argmax}_{T \in \mathcal{T}(s)} \sum_{h,m} T_{h,m} \times W_{h,m}$$

Semi-Supervised Variational Auto-Encoder



We assume a sentence is generated from a latent tree and a latent embedding:

$$p(s) = \sum_T \int p(s, T, z) dz$$

We perform variational inference by jointly learning a distribution which is intended to be close to the posterior:

$$KL [q(T, z | s) || p(T, z | s)] \simeq 0$$

$$\mathcal{L} = - \sum_j \log p(s_j) \quad \left. \begin{array}{l} \text{Unsupervised auto-encoder loss} \\ \text{[Kingma & Welling, 2013]} \end{array} \right\}$$

$$- \sum_i \log q(T_i | s_i) \quad \left. \begin{array}{l} \text{Supervised loss} \\ \text{[Kingma et al., 2014]} \end{array} \right\}$$

Perturb-and-Parse

As exact marginalisation over dependency trees is intractable, we introduce a reparametrization for differentiable Monte-Carlo estimation.

$$G \sim \mathcal{G}(0,1)$$

$$\tilde{W} = W + G$$

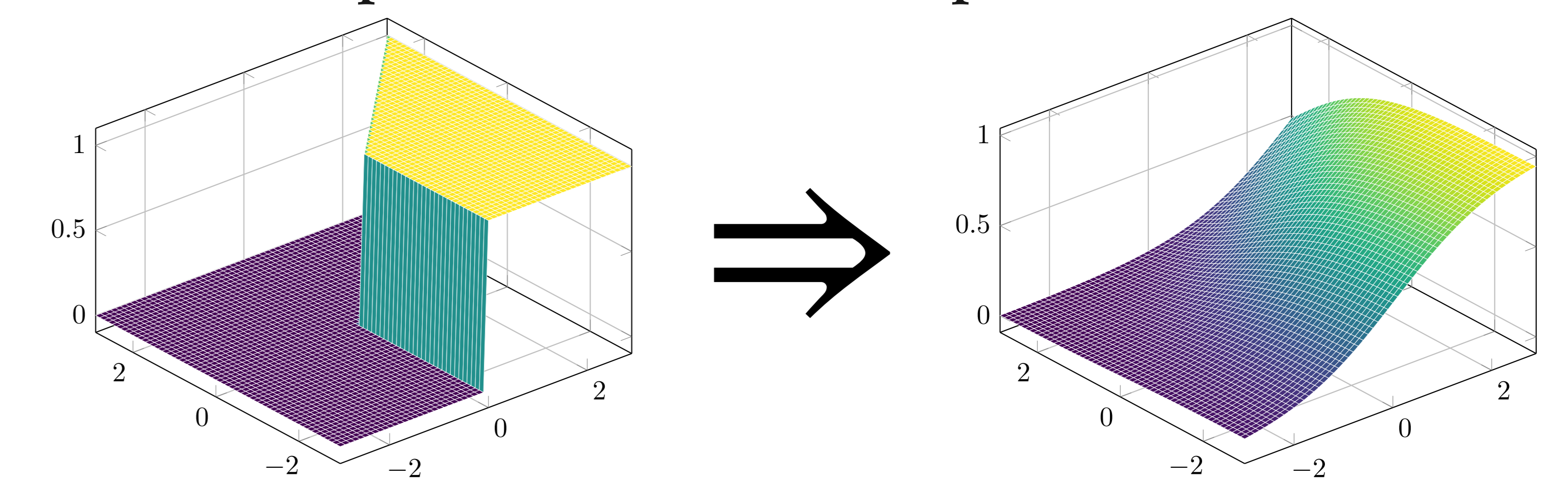
$$\operatorname{argmax}_{T \in \mathcal{T}(s)} \sum_{h,m} T_{h,m} \times \tilde{W}_{h,m}$$

} Arc weight perturbation with Gumbel noise [Papandreou & Yuille, 2011]

} Solved with dynamic programming [Eisner, 1996]

Differentiable Dynamic Programming

The dynamic programming approach for parsing relies on recursive calls to the *one-hot-argmax* op, which introduces ill-defined derivatives during the backward pass. We replace *one-hot-argmax* ops with *softmax* ops to smooth the optimization landscape.



[Maddison et al., 2017; Goyal et al., 2018]

Acknowledgments

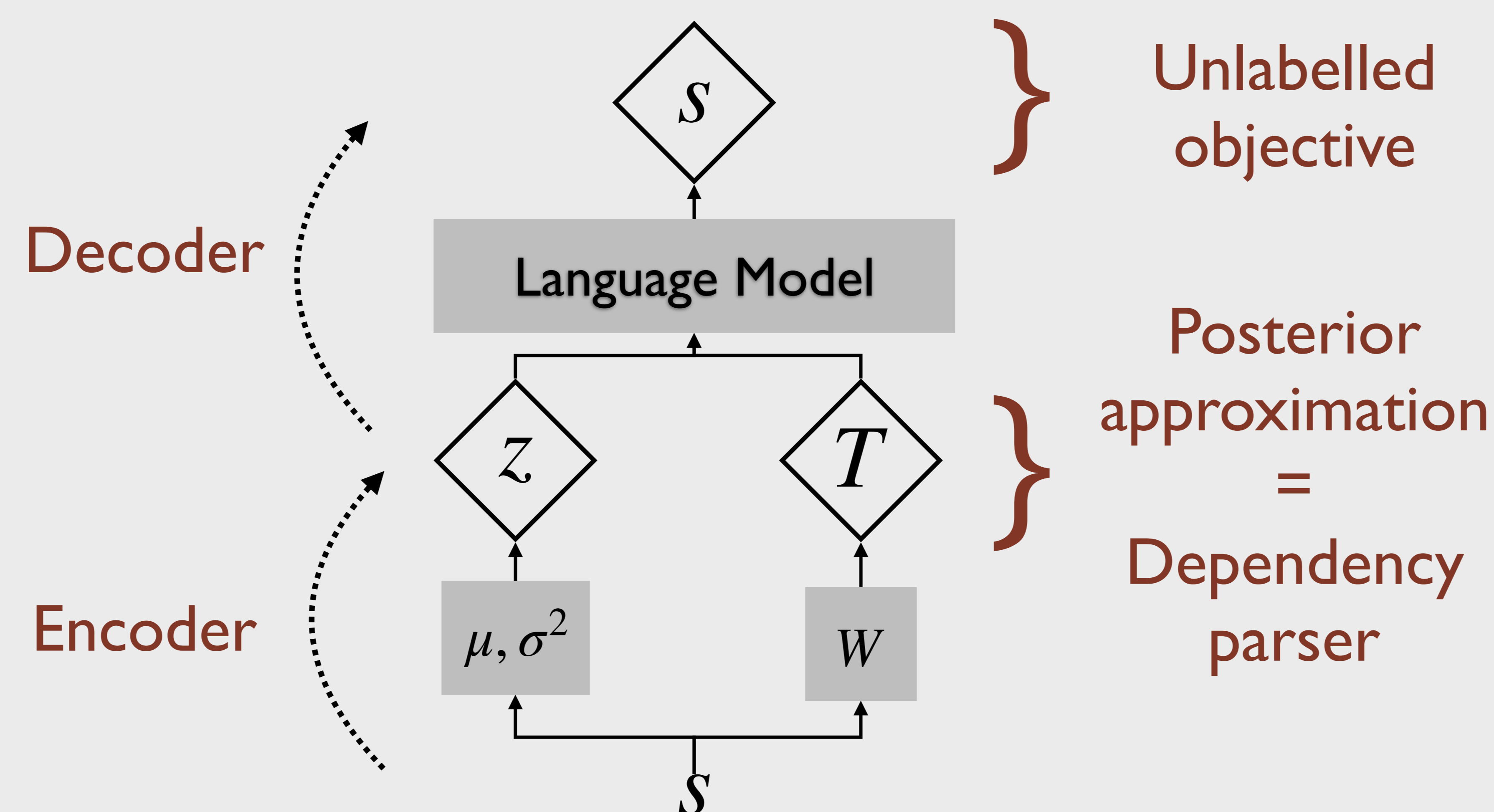
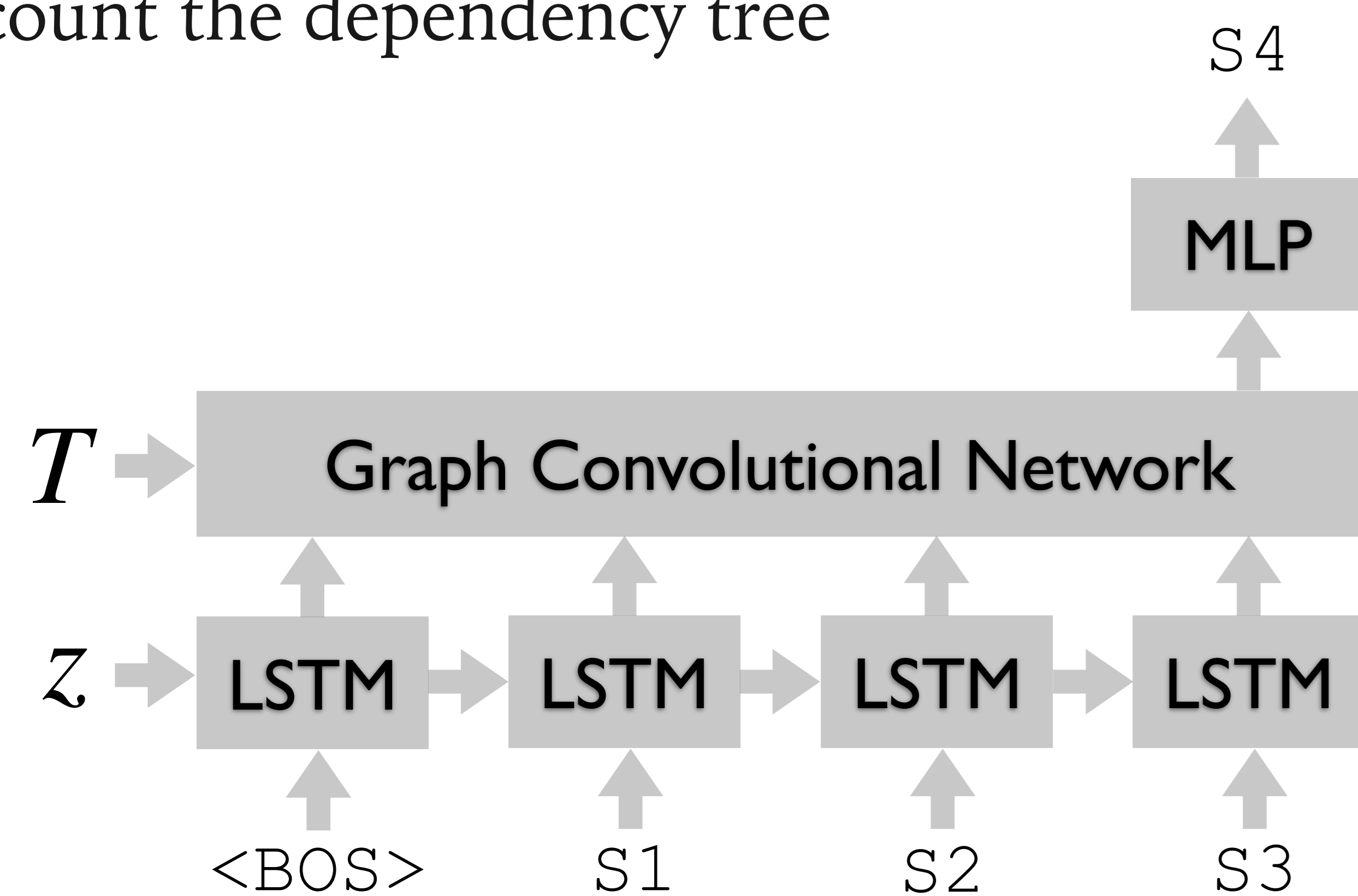


NWO VIDI 639.022.518
ERC BroadSem 678254



Syntactic language Model

Auto-regressive model that takes into account the dependency tree



Experimental results

	English	French	Swedish
Supervised	88.79 / 84.74	84.09 / 77.58	86.59 / 78.95
VAE with z	89.39 / 85.44	84.43 / 77.89	86.92 / 80.01
VAE without z	89.50 / 85.48	84.69 / 78.49	86.97 / 79.80

Unlabeled/labeled attachment scores. Only 10% of the data is labeled (~4000 sentences for English)

	Supervised prec/recall	Semi-sup. prec/recall
(root)	93.46 / 89.30	93.84 / 92.41
> 7	72.47 / 83.26	78.72 / 83.11

The main improvement is observed on root word identification and long distance dependencies (arcs crossing at least 7 words)