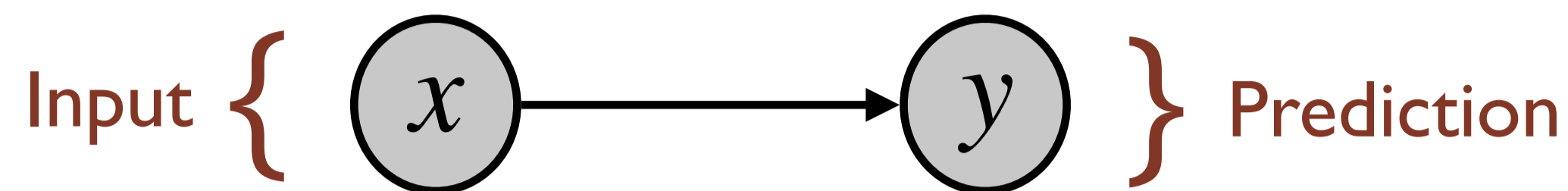# Learning Latent Trees with Stochastic Perturbations and Differentiable Dynamic Programming
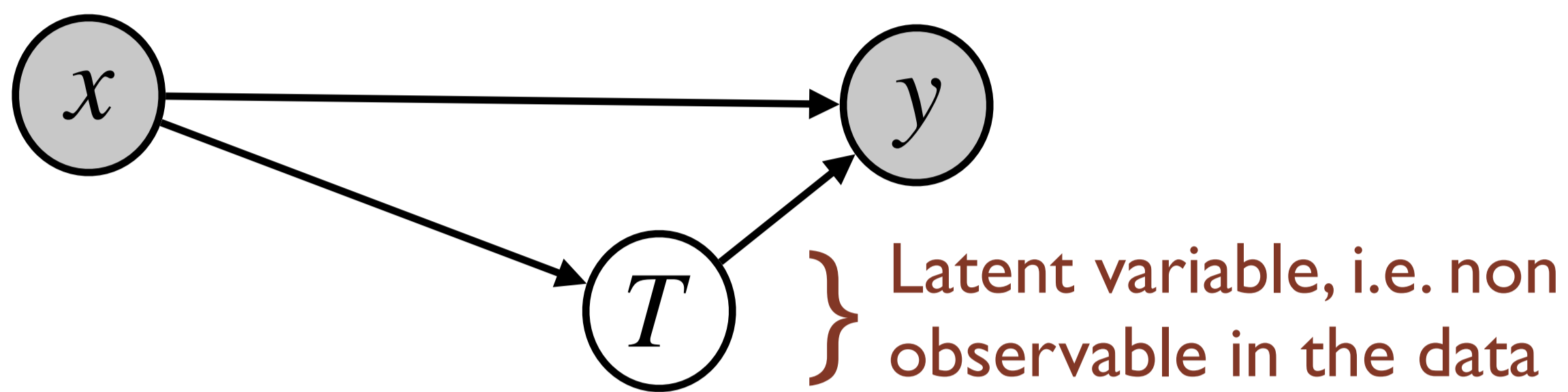
Caio Corro, Ivan Titov

ILCC, School of Informatics, University of Edinburgh    ILLC, University of Amsterdam

## Latent Variable Models

Supervised learning can be understood as inferring the probability distribution corresponding to a directed graphical model.
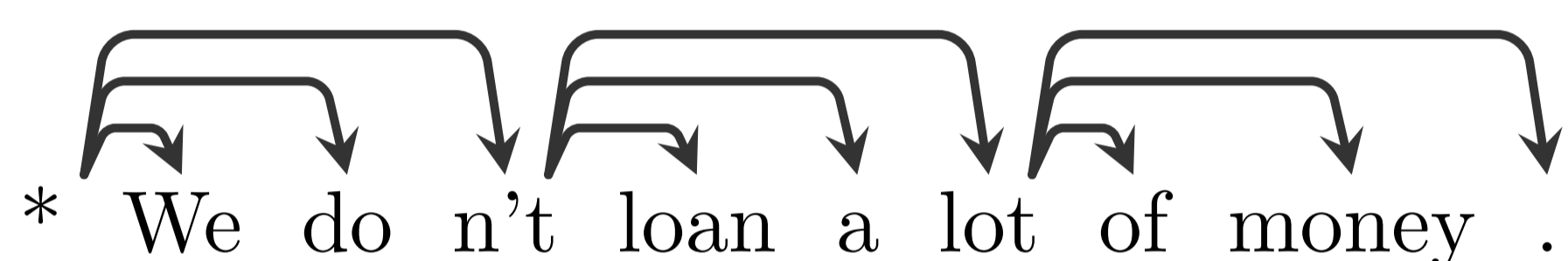
Input $\{$ $x$ → $y$ $\}$ Prediction

Latent variables can model unobserved inter-dependencies or introduce knowledge about the structure of a given problem.

$x$, $y$, $T$ } Latent variable, i.e. non observable in the data

## Projective Dependency Tree

We are interested in **latent projective dependency trees** that implicitly encode hierarchical decomposition of a sentence into spans.

\* We do n't loan a lot of money .

## Distribution over Trees

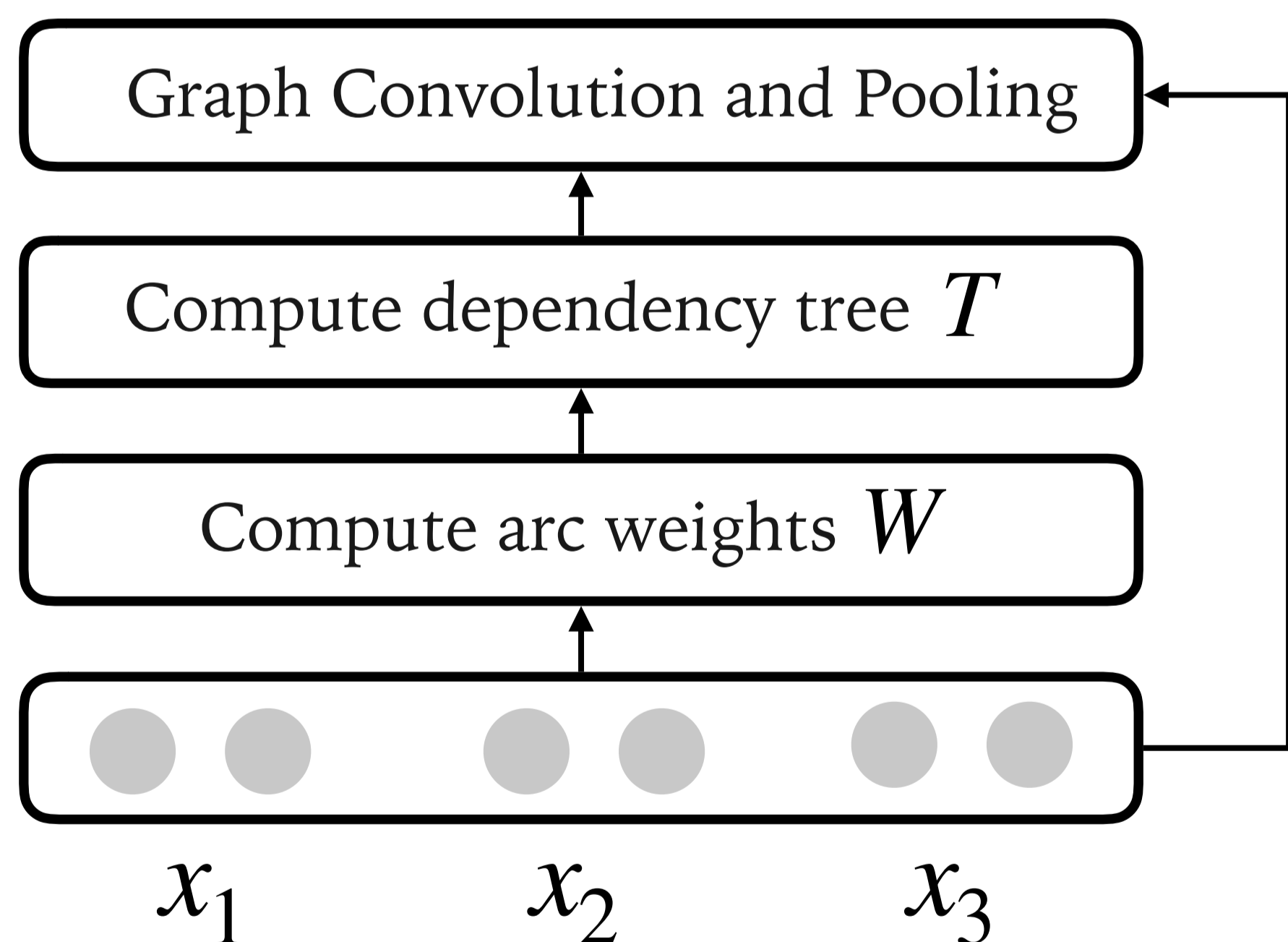The probability distribution over dependency trees is a log-linear model factored over arc weights.

$W$ : matrix of arc weights computed with a NN

$T$ : boolean adjacency matrix,
    i.e $T_{h,m} = 1$ iff arc $x_h \rightarrow x_m$ is in the tree

$$p(T\,|\,x) = \frac{\sum_{h,m} T_{h,m} \times W_{h,m}}{Z(T,x)}$$

## Neural Architecture

A **Graph Convolutional Network** [Kipf and Welling, 2017] is used to compute the sentence representation w.r.t. the dependency tree.

Graph Convolution and Pooling

Compute dependency tree $T$

Compute arc weights $W$

$x_1$    $x_2$    $x_3$

## Training Loss

We maximise the log-likelihood of training data via SGD:

$$\log p(y\,|\,x) = \log \mathbb{E}_{T \sim p(T|x)}[\,p(y\,|\,T,x)\,]$$

Unfortunately, exact marginalisation is intractable:

$$= \log \sum_T p(T\,|\,x) \times p(y\,|\,T,x)$$

Therefore, we derive a bound using Jensen's inequality:

$$\geq \mathbb{E}_{T \sim p(T|x)}[\,\log p(y\,|\,T,x)\,]$$

Which can be approximated via Monte-Carlo method.

## Contributions

1. We show that a **latent tree model** can be estimated by drawing global approximate samples via **Gumbel perturbation and differentiable dynamic programming**

2. We demonstrate that constraining the structures to be projective dependency trees is beneficial

3. We show the effectiveness of our approach on two standard tasks and on a synthetic dataset

## Perturb-and-MAP

Approximate sampling method for log-linear models:

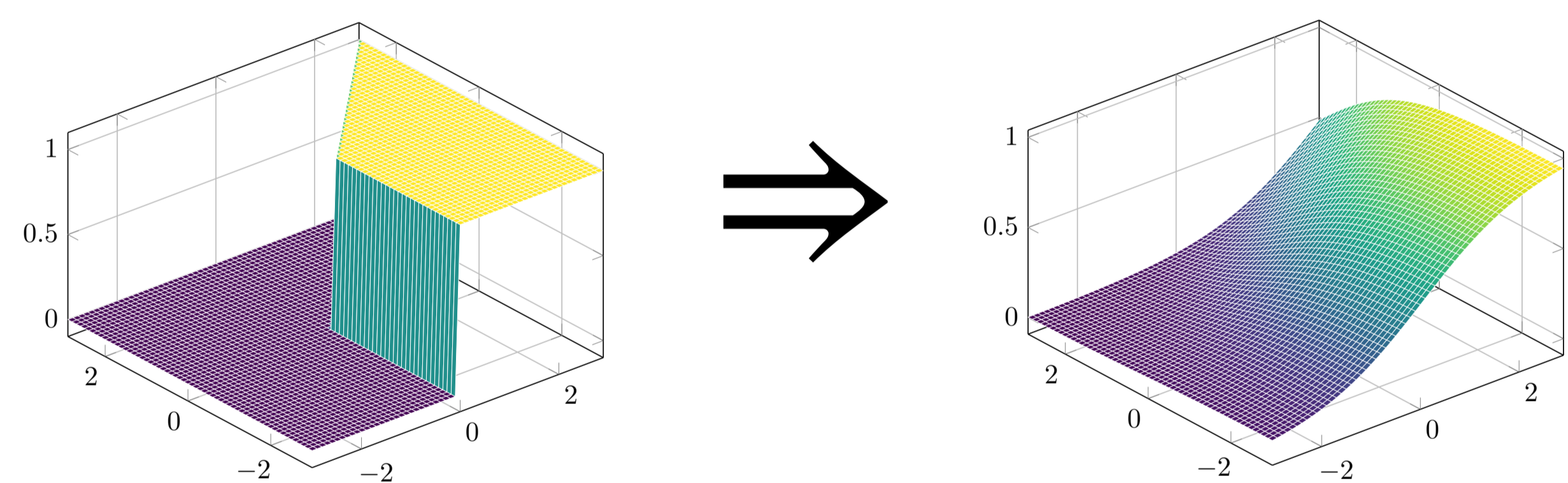$$G \sim \mathscr{G}(0,1)$$

$$\widetilde{W} = W + G$$
} Arc weight perturbation with Gumbel noise [Papandreou & Yuille, 2011]

$$\underset{T \in \mathscr{T}(s)}{argmax} \sum_{h,m} T_{h,m} \times \widetilde{W}_{h,m}$$
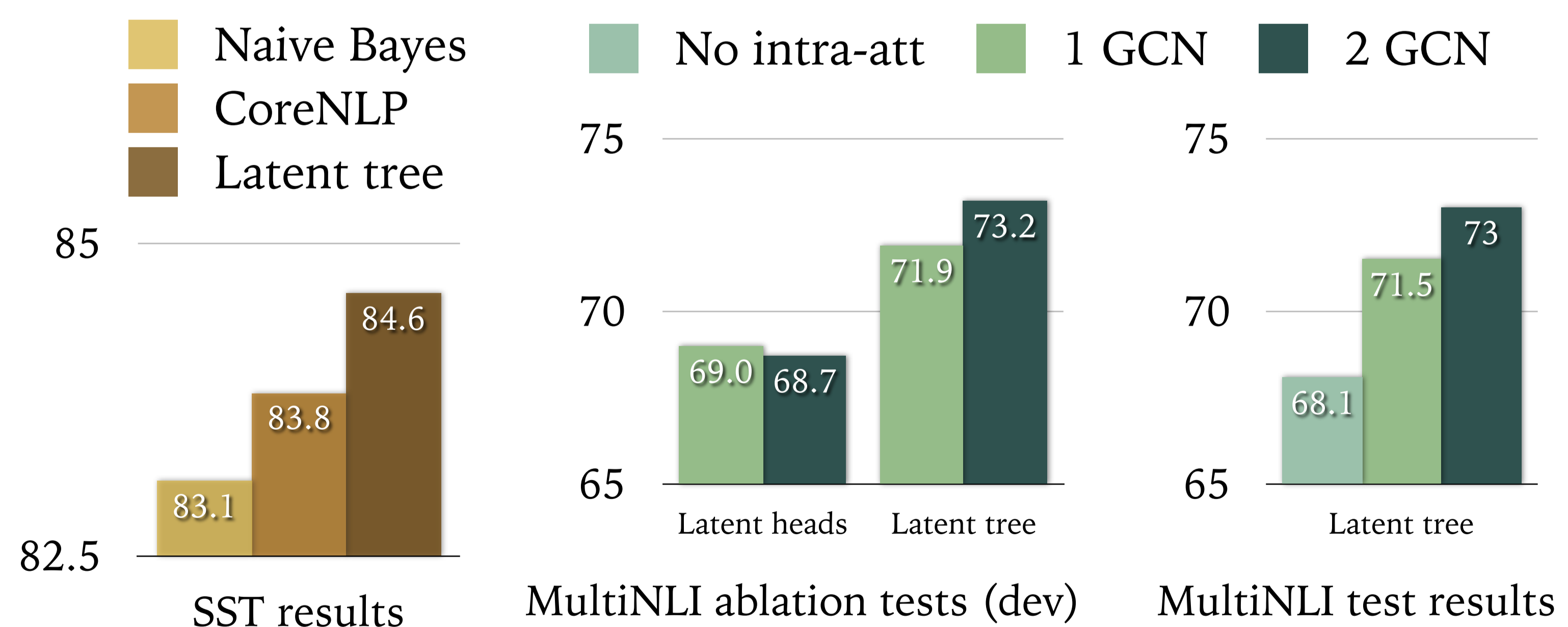} Solved with dynamic programming [Eisner, 1996]

## 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.



## Experimental Results

Experimentally, we observe that our **Latent Tree (LT)** model improves comparable baselines on **sentiment analysis** with syntactic trees predicted by CoreNLP and on **Natural Language Inference** datasets.



Naive Bayes, CoreNLP, Latent tree — SST results (83.1, 83.8, 84.6)

No intra-att, 1 GCN, 2 GCN — MultiNLI ablation tests (dev): Latent heads (69.0, 68.7), Latent tree (71.9, 73.2)

MultiNLI test results: Latent tree (68.1, 71.5, 73)