

GCN for Dependency Parsing

Outline

1. Background
2. Task Definition
3. Graph-based Dependency Parsing with Graph Neural Networks

1. Syntax

- The word **syntax** comes from the Greek syntaxis , meaning “setting out together or arrangement”, and refers to **the way words are arranged together**.

一、词类: n、adj、adv、v、conj、prep、pron、art、num (红: 实词 黑: 虚词)
名 形容 副 动 连 介 代 冠 数

所有格 n { 专有: 人名, 地名, 国家名 例: Tom, Nan Chong, China 中国
the Great Wall 长城
普通 { 可数 pl 规则, 不规则 计量 { a book, two books 一, 两本书
two baskets of apples 两篮子苹果
不可数: 液体, 气体, 抽象 a bottle of water 一瓶水
two bottles of water 两瓶水

adj ...的

级 { v run fast 跑得快
修饰 { adj very good 非常好
adv very fast 非常快

时间、地点、原因: yesterday there why
(在) 昨天 在那儿 为什么

语法

1. 词类: n, 分类, 计量, 所有格, adj, adv, v, conj, prep, pron, art, num.

2. 四种 { 陈述: 抓主干 SVO, SLP, 理枝叶 DZB。5个简单句。
疑问 { 一般: 助, 情, 系。特殊: 几个W, 直接代替, 变一般。
反意: 翻译, 回答, 祈使句, shall, will。选择: or, which
祈使: (you) Don't Do + O/L + P 感叹: how, what, 四

3. 时态: 定义, 顺序。现, 过, 将, 过将; 一, 进, 完, 完进, 6种形式。否定
语态: be done / be being done / have been done

4. 并列句: 平2, 转3, 选3, 因3。

从句 { n, SOPT, 两种情况, that, whether, if, 9个W。
adj, 关系词, 代, that, who, whom, which, whose, 副, when, where, why。
adv, 时6, 地2, 伴1, 原1, 让7, 杂5, 目2, 结2, 方3。

5. 非: to do, doing, done。to+L+P n, adv Ling+P n, adv

6. 独: (with) S₁ + adj/adv/非 7. 特: there, it, 省, 倒, 插

8. 直间: 陈述, 一般, 特殊, 宾从。ask sb to do sth

9. 主谓一致: 并, 名, 代, 其, number, percent

10. 虚拟语气: if, as if, wish 往后推。332 (should) do

等效思想

介宾 → adj, adv

pron → n, adj

num { 基数 → 单 n, adj

序数 → 单 n, adj, adv

从句 { 单 n: S, O, P, 同

adj

adv

非谓语 { to do → 单 n, adj, adv

动词 { doing → 单 n, adj, adv

done → adj, adv

独立主格 → adv

1.1 Part-Of-Speech

- Parts of speech (POS): word classes or syntactic categories
 - noun, verb, pronoun, preposition, adverb, conjunction, participle, and article
- Example: The Penn TreeBank

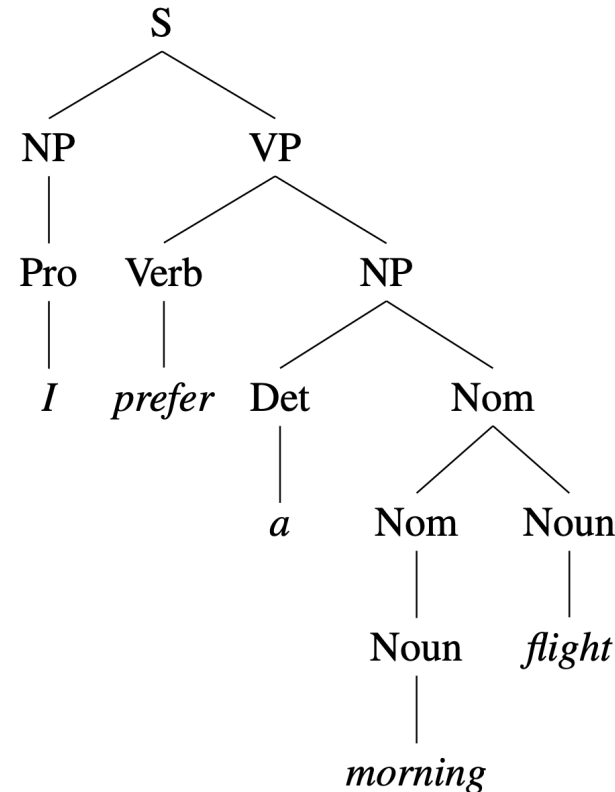
Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coordinating conjunction	<i>and, but, or</i>	PDT	predeterminer	<i>all, both</i>	VBP	verb non-3sg present	<i>eat</i>
CD	cardinal number	<i>one, two</i>	POS	possessive ending	<i>'s</i>	VBZ	verb 3sg pres	<i>eats</i>
DT	determiner	<i>a, the</i>	PRP	personal pronoun	<i>I, you, he</i>	WDT	wh-determ.	<i>which, that</i>
EX	existential 'there'	<i>there</i>	PRP\$	possess. pronoun	<i>your, one's</i>	WP	wh-pronoun	<i>what, who</i>
FW	foreign word	<i>mea culpa</i>	RB	adverb	<i>quickly</i>	WP\$	wh-possess.	<i>whose</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	RBR	comparative adverb	<i>faster</i>	WRB	wh-adverb	<i>how, where</i>
JJ	adjective	<i>yellow</i>	RBS	superlatv. adverb	<i>fastest</i>	\$	dollar sign	<i>\$</i>
JJR	comparative adj	<i>bigger</i>	RP	particle	<i>up, off</i>	#	pound sign	<i>#</i>
JJS	superlative adj	<i>wildest</i>	SYM	symbol	<i>+, %, &</i>	"	left quote	<i>' or "</i>
LS	list item marker	<i>1, 2, One</i>	TO	"to"	<i>to</i>	"	right quote	<i>' or "</i>
MD	modal	<i>can, should</i>	UH	interiection	<i>ah, oops</i>	(left paren	<i>[, (, {, <</i>
NN	sing or mass noun	<i>llama</i>	VB	verb base form	<i>eat</i>)	right paren	<i>],), }, ></i>
NNS	noun, plural	<i>llamas</i>	VBD	verb past tense	<i>ate</i>	,	comma	<i>,</i>
NNP	proper noun, sing.	<i>IBM</i>	VBG	verb gerund	<i>eating</i>	.	sent-end punc	<i>. ! ?</i>
NNPS	proper noun, plu.	<i>Carolinas</i>	VBN	verb past part.	<i>eaten</i>	:	sent-mid punc	<i>: ; ... --</i>

1.2 Constituency Tree

- Constituency: abstraction for groups of words -> hierarchical graph.

I *prefer* *a* *morning* *flight*
| | | | |
PRP VBP DT NN NN

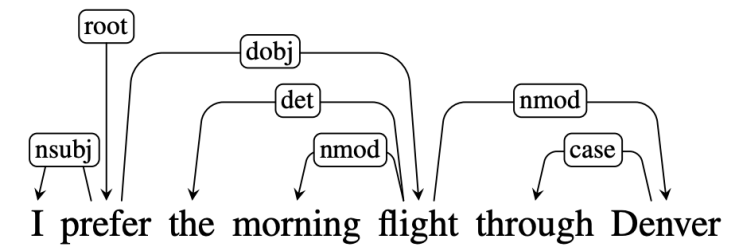
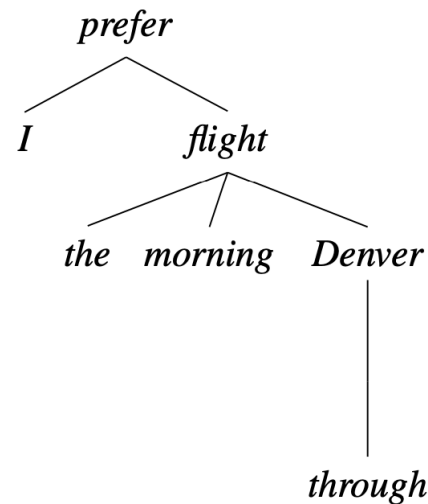
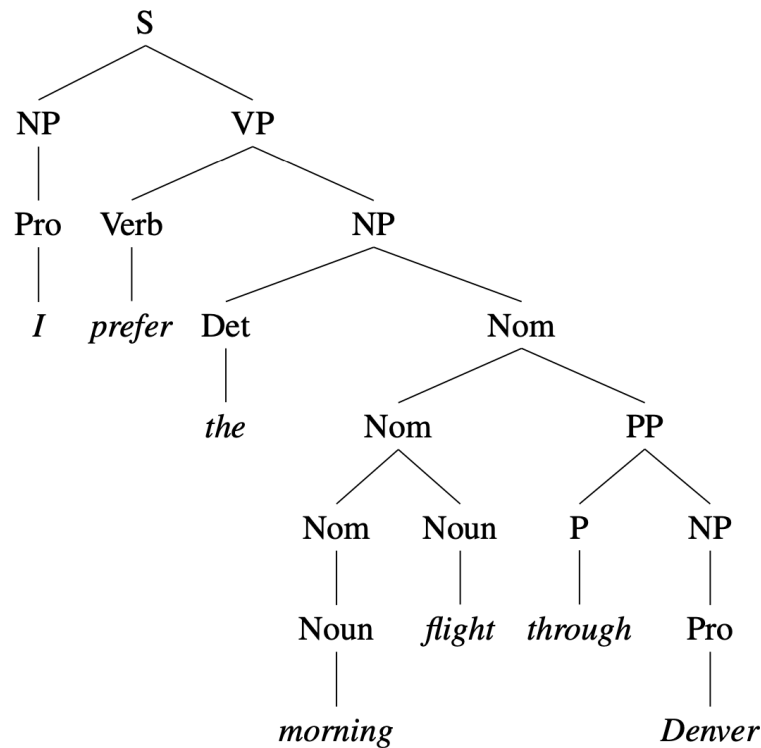
PRP: Personal Pronoun
VBP: Verb Present
DT: Determiner
NN: Single or Mass Noun



S: Sentence
NP: Noun Phrase
VP: Verb Phrase
Pro: Pronoun
Det: Determiner
Nom: Nominal Noun

1.3 Dependency Tree

- Dependency: relation between a pair of words -> dependency tree.



root

nsubj: nominal subject

dobj: direct object

det: determiner

nmod: nominal modifier

case: prepositions etc.

Important property: DAG constrain
(Directed Acyclic Graph)

2. Dependency Parsing via Graph Algorithm

- Given

$s = w_1, \dots, w_n$	A sentence containing n words
$G = (V, E)$	The complete graph over the sentence s .
V	The set of nodes (w_1, \dots, w_n and a root node)
E	The set of node edges $e = (i, j, r)$, where $w_i \xrightarrow{r} w_j$

- Find the best dependency tree \hat{T} which has the highest score.

$$\hat{T} = \operatorname{argmax}_{T \in G} \operatorname{score}(T, s)$$

$$\operatorname{score}(T, s) = \sum_{e \in T} \operatorname{score}(e, s)$$

2. Dependency Parsing via Graph Algorithm

- Two steps for solving the problem:
 - 1) Compute scores for each edge. (*main focus of the paper)
 - 2) Infer dependency tree based on the estimated edge scores.
 - i. Search for optimal tree (without edge labels)
 - ii. Predict edge labels

3. The paper

Graph-based Dependency Parsing with Graph Neural Networks

Tao Ji , Yuanbin Wu , and Man Lan

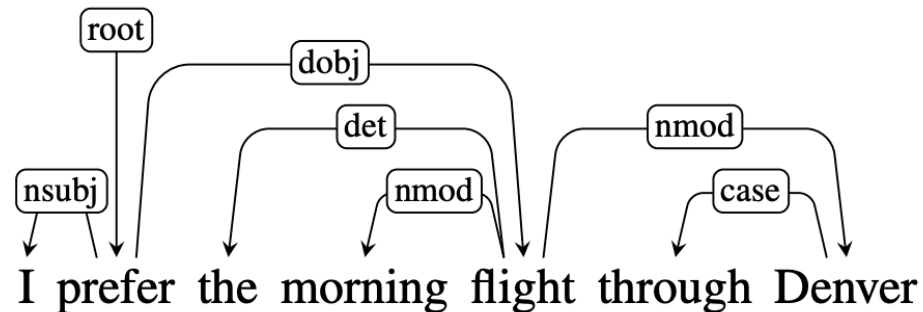
Department of Computer Science and Technology,

East China Normal University

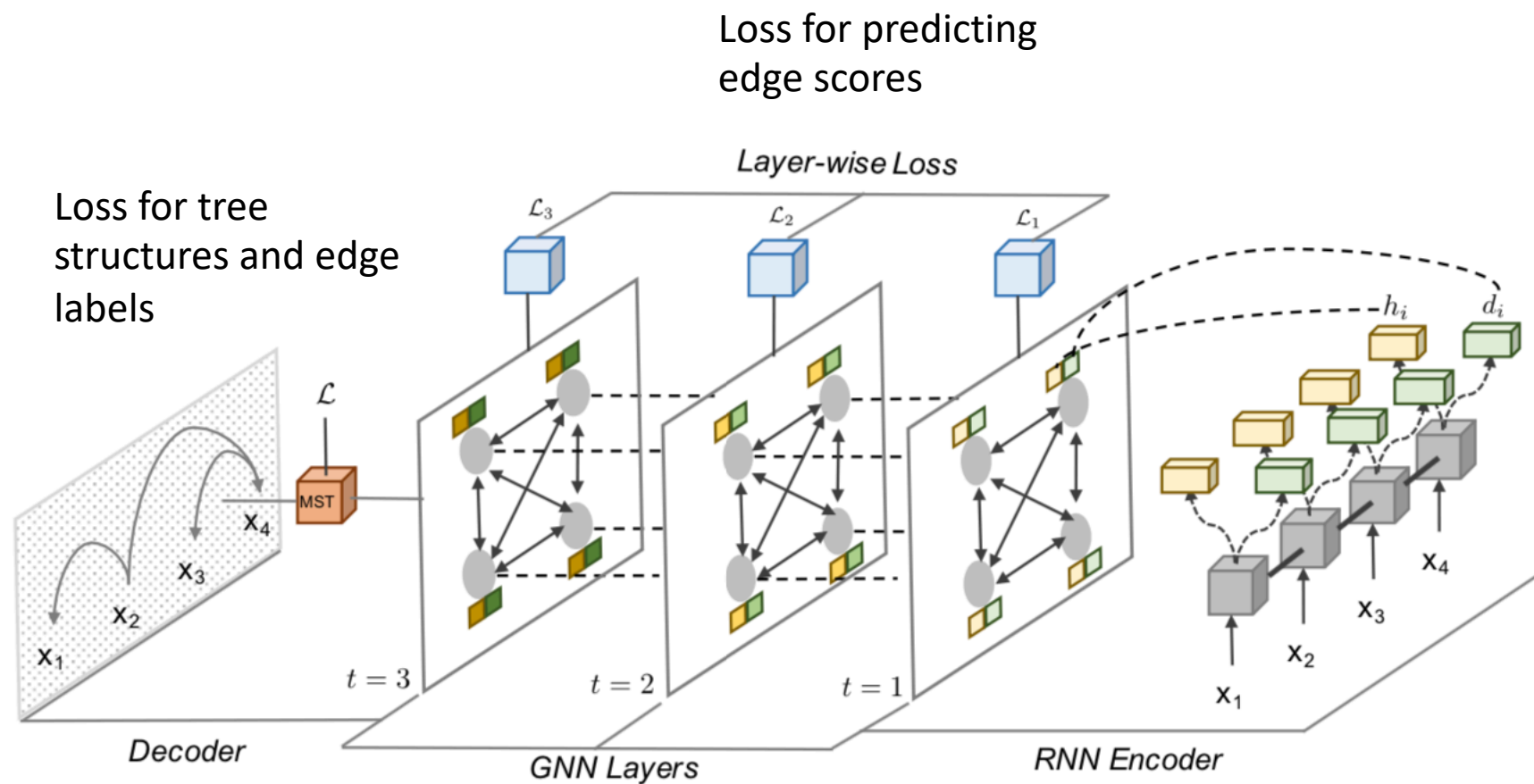
`{taoji.cs}@gmail.com {ybwu,mlan}@cs.ecnu.edu.cn`

3.1 Motivation

- Previous work has shown first-order relations can help learning.
 - Head and dependent vector.
- Introduce structure knowledge into node representations.
 - Especially high-order relations (e.g. grandparents)
 - Graph neural network to capture long-term dependencies.



3.2 Overview



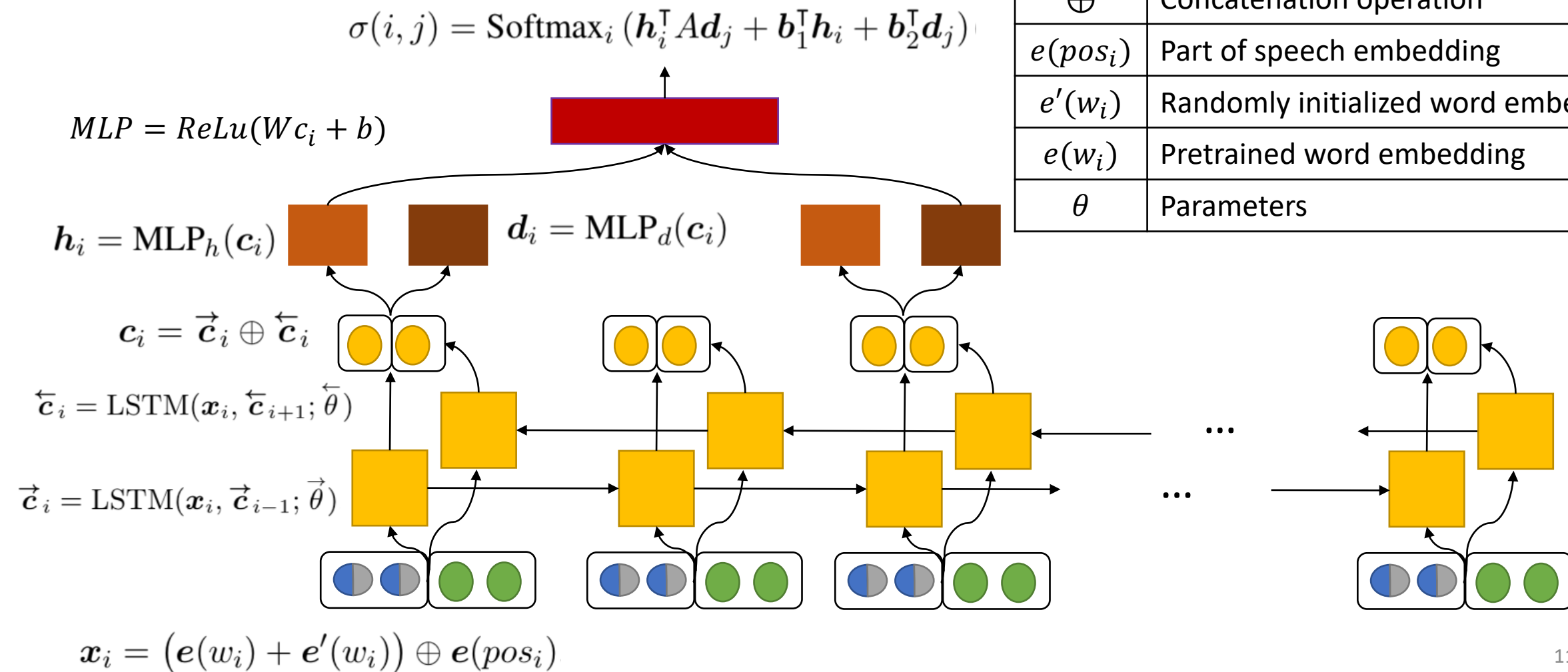
Infer optimal tree given a complete graph with edge scores.

Encode high-order information and predict edge scores

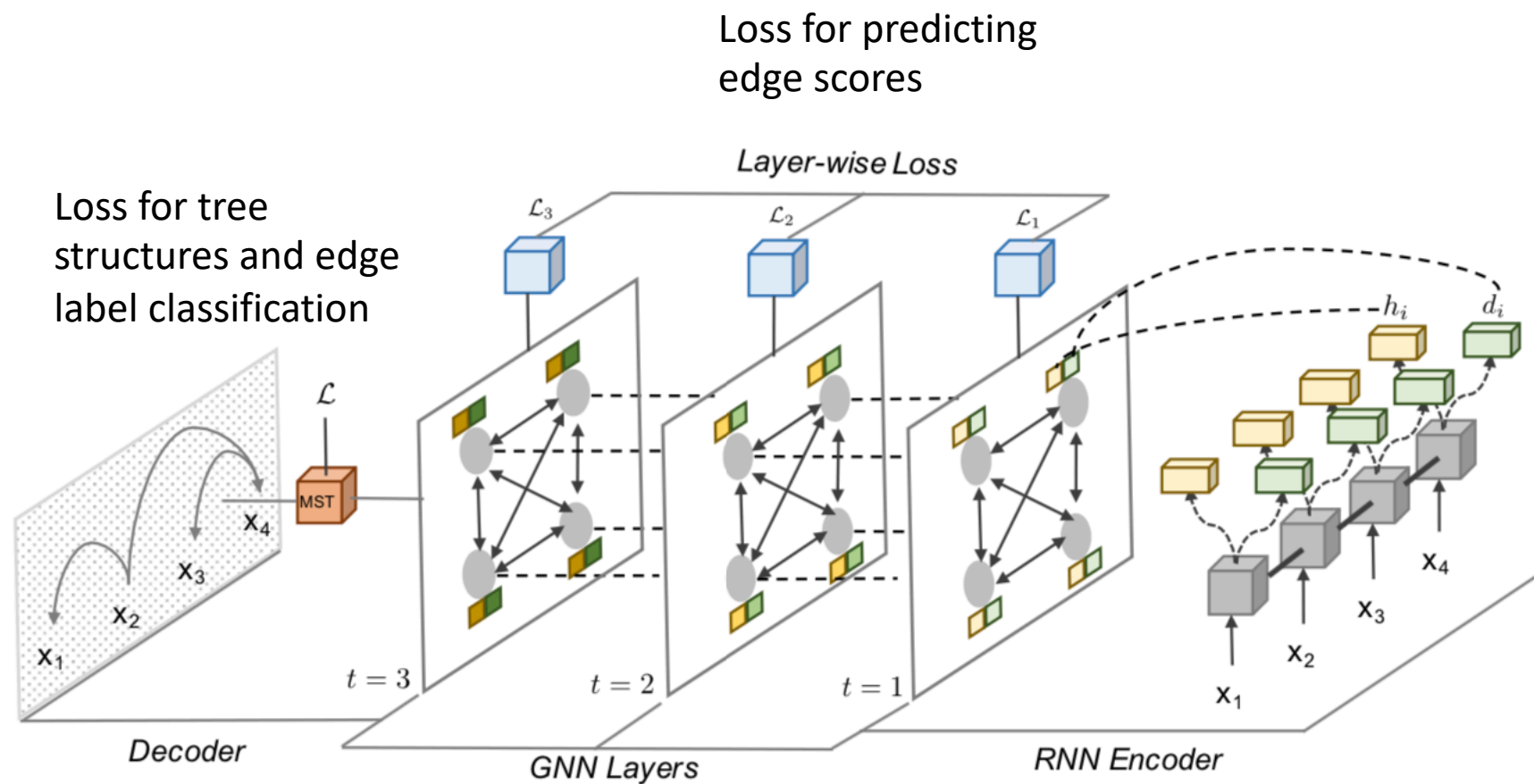
Basic node representations

3.3 RNN Encoder

$\sigma(i, j)$	Score for edge $w_i \rightarrow w_j$
W, A, b_*	Weights
\oplus	Concatenation operation
$e(pos_i)$	Part of speech embedding
$e'(w_i)$	Randomly initialized word embedding
$e(w_i)$	Pretrained word embedding
θ	Parameters



Overview



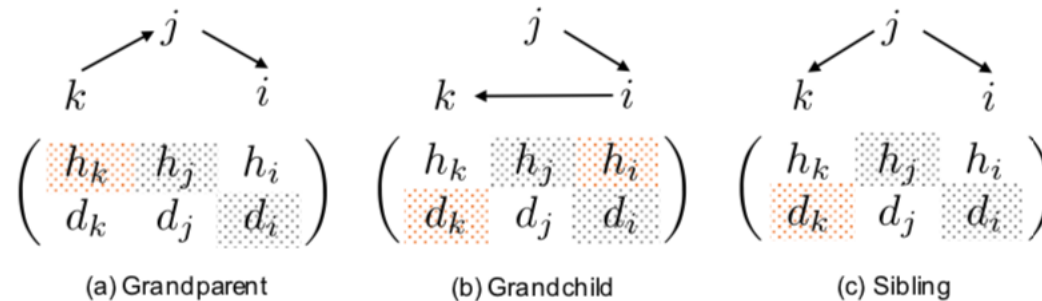
Infer optimal tree given a complete graph with edge scores.

Encoding high-order information

Basic node representations

3.4 GNN Layers

- Encode three types of high-order information for edge (j, i)

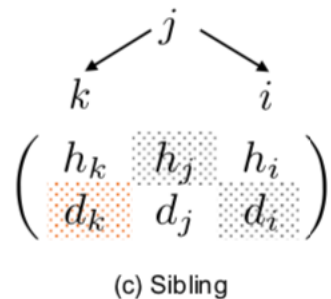
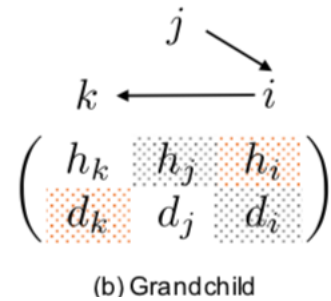
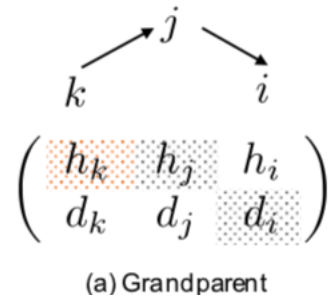


- 1) The grey shadows indicate first-order information: parent \rightarrow children
- 2) The orange shadows are the second-order information.

- Multi-layer GCN is able to explicitly capture this high-order information.

3.4 GNN Layers (cont')

- Encode three types of high-order information for edge (j, i)



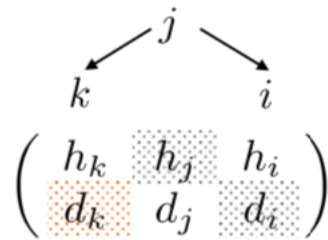
$$\begin{cases} \mathbf{h}_i^t = g\left(W_1 \sum_{j \in \mathcal{N}(i)} \alpha_{ji}^t \mathbf{h}_j^{t-1} + B_1 \mathbf{h}_i^{t-1}\right) \\ \mathbf{d}_i^t = g\left(W_2 \sum_{j \in \mathcal{N}(i)} \alpha_{ij}^t \mathbf{d}_j^{t-1} + B_2 \mathbf{d}_i^{t-1}\right). \end{cases}$$

$$\begin{cases} \mathbf{h}_i^t = g\left(W_1 \sum_{j \in \mathcal{N}(i)} \alpha_{ij}^t \mathbf{d}_j^{t-1} + B_1 \mathbf{h}_i^{t-1}\right) \\ \mathbf{d}_i^t = g\left(W_2 \sum_{j \in \mathcal{N}(i)} \alpha_{ji}^t \mathbf{h}_j^{t-1} + B_2 \mathbf{d}_i^{t-1}\right). \end{cases}$$

W_1, B_1	Weights for updating the head vector
W_2, B_2	Weights for updating the dependent vector
α_{ij}^t	Score for edge (i, j) at the t^{th} layer.
\mathbf{h}_i^t	The head vector for node i at the t^{th} layer.
\mathbf{d}_i^t	The dependent vector for node i at the t^{th} layer.
g	Non-linear activation function

3.4 GNN Layers – My understanding

- How to gather information of neighbors?
 - 1) Head nodes: gather information from neighbors' dependent nodes.
 - 2) Dependent nodes: gather information from neighbors' head nodes.

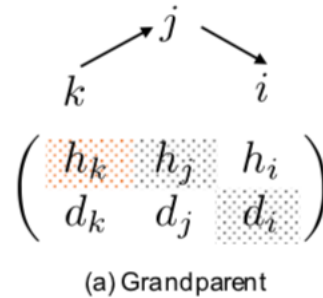


$$h_j^t = g(W_1 \sum_{i \in \mathcal{N}(j)} \alpha_{ji}^t d_i^{t-1} + B_1 h_j^{t-1})$$

$$d_i^t = g(W_2 \sum_{j \in \mathcal{N}(i)} \alpha_{ji}^t h_j^{t-1} + B_2 d_i^{t-1})$$

3.4 GNN Layers – My understanding (cont')

- Now, we want to capture the 2nd order information, how to propagate h_k to d_i via node h_j ?



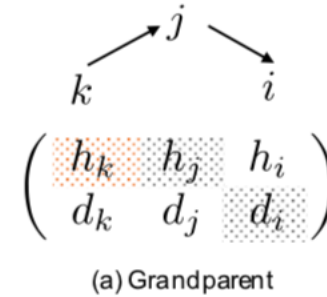
- 1) $h_k \rightarrow d_m$
- 2) $d_m \rightarrow h_j$
- 3) $h_j \rightarrow d_i$ for some word m



Conclusion: unable to directly capture the 2nd order information.

3.4 GNN Layers – My understanding (cont')

- How to solve this problem?
 - Shortcut from h_k to h_j .
 - 1) h_k to h_j . (shortcut update)
 - 2) h_j to d_i . (standard update)



3.4 GNN Layers – Put updating functions together

- Synchronous updates:
 - update h_i and d_i at the same time

$$\begin{cases} \mathbf{h}_i^t = g\left(W_1 \sum_{j \in \mathcal{N}(i)} (\alpha_{ji}^t \mathbf{h}_j^{t-1} + \alpha_{ij}^t \mathbf{d}_j^{t-1}) + B_1 \mathbf{h}_i^{t-1}\right) \\ \mathbf{d}_i^t = g\left(W_2 \sum_{j \in \mathcal{N}(i)} (\alpha_{ij}^t \mathbf{h}_j^{t-1} + \alpha_{ji}^t \mathbf{d}_j^{t-1}) + B_2 \mathbf{d}_i^{t-1}\right). \end{cases}$$

- Asynchronous updates:
 - updates h_i first, then update d_i

$$\begin{cases} \mathbf{h}_i^{t-\frac{1}{2}} = g\left(W_1 \sum_{j \in \mathcal{N}(i)} (\alpha_{ji}^t \mathbf{h}_j^{t-1} + \alpha_{ij}^t \mathbf{d}_j^{t-1}) + B_1 \mathbf{h}_i^{t-1}\right) \\ \mathbf{d}_i^t = g\left(W_2 \sum_{j \in \mathcal{N}(i)} (\alpha_{ij}^t \mathbf{h}_j^{t-\frac{1}{2}} + \alpha_{ji}^t \mathbf{d}_j^{t-1}) + B_2 \mathbf{d}_i^{t-1}\right), \end{cases}$$

3.4 GNN Layers – Edge scores

1. Soft scores:

$$\alpha_{ij}^t = \sigma^t(i, j) = P^t(i|j)$$
$$\sigma(i, j) = \text{Softmax}_i (\mathbf{h}_i^\top A \mathbf{d}_j + \mathbf{b}_1^\top \mathbf{h}_i + \mathbf{b}_2^\top \mathbf{d}_j)$$

2. $\{0, 1\}$ scores:

$$\alpha_{ij}^t = \begin{cases} 1, & i = \arg \max_{i'} P^t(i'|j) \\ 0, & \text{otherwise} \end{cases}$$

3. Top k scores:

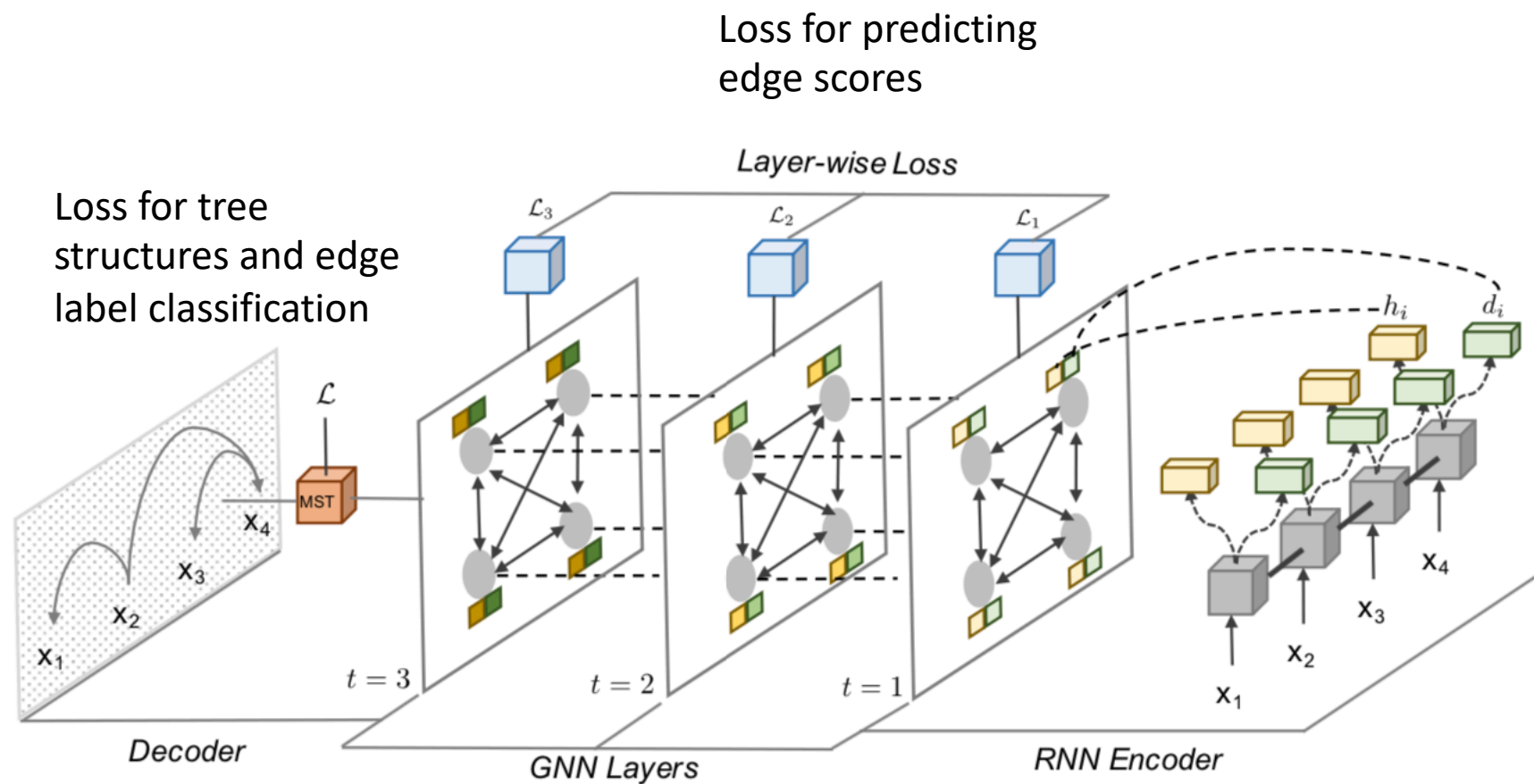
$$\alpha_{ij}^t = \begin{cases} \text{Softmax}_i (\mathbf{h}_i^\top A \mathbf{d}_j + \mathbf{b}_1^\top \mathbf{h}_i + \mathbf{b}_2^\top \mathbf{d}_j), & i \in \mathcal{N}_k^t(j) \\ 0, & \text{otherwise} \end{cases}$$

$\mathcal{N}_k(j)$ top k neighbors
with highest scores

4. Average:

$$\alpha_{ij}^t = \frac{1}{n}, \quad \forall j \in V, i \in V/\{j\}$$

Overview



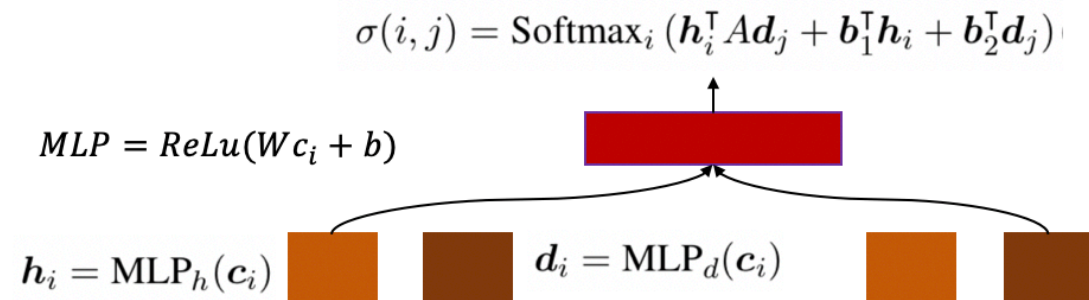
Infer optimal tree given a complete graph with edge scores.

Encoding high-order information

Basic node representations

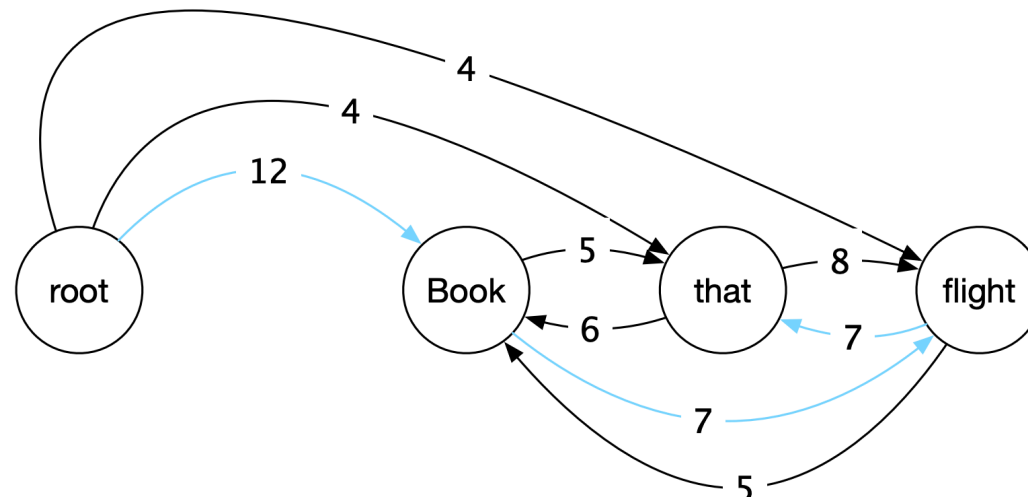
3.5 Decoding Optimal Tree

- Decoding optimal tree structure
 - Maximum Spanning Tree (MST) Algorithm– Chu-Liu Edmonds Algorithm
- Label prediction (classification):
 - Similar as computing edge scores.
 - Multi-Layer Perceptron to model $P(r|i, j)$ for edge (i, j)
 - Input: representations for node i and j
 - Output: $r = \operatorname{argmax} P(r|i, j)$

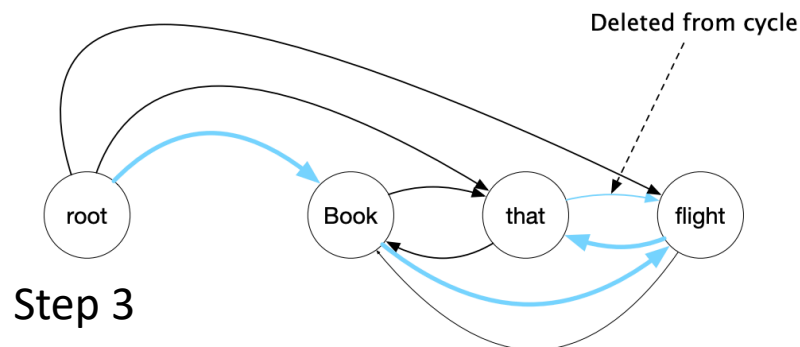
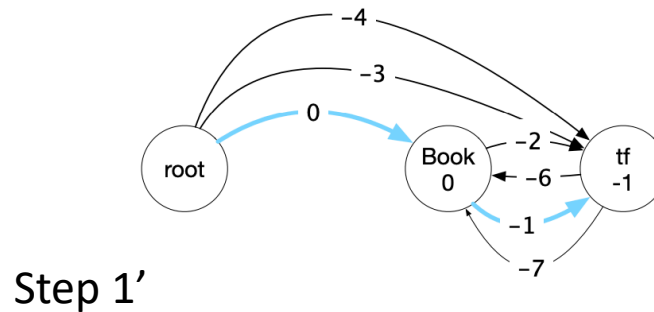
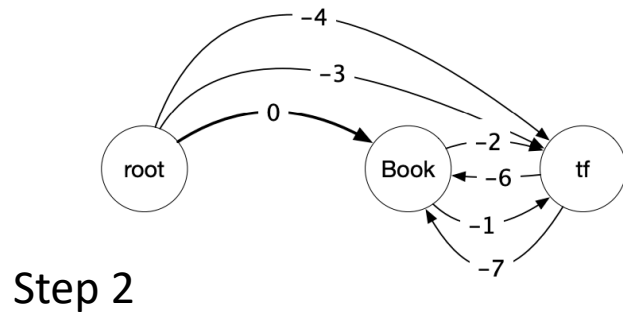
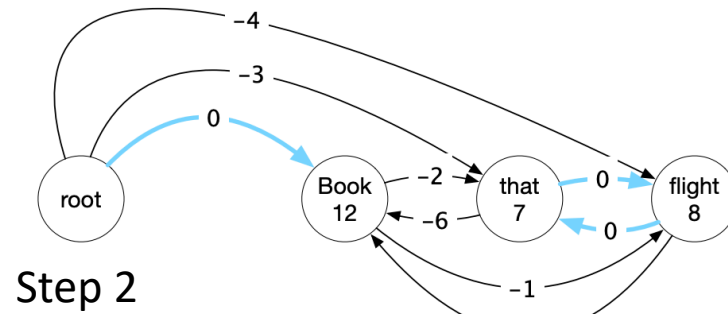
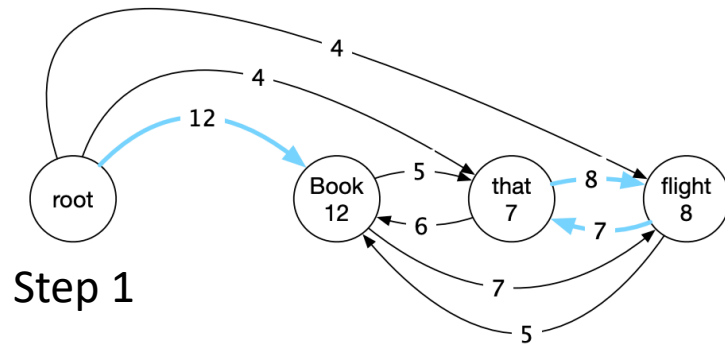


3.5.1 Maximum Spanning Tree (MST)

- Spanning Tree: A spanning tree $T(V, E')$ of a directed graph $G = (V, E)$ is a tree over G , which contains all of the nodes in G with minimum number of edges.
- MST: A MST is a spanning tree $T(V, E')$ with highest edge weights.

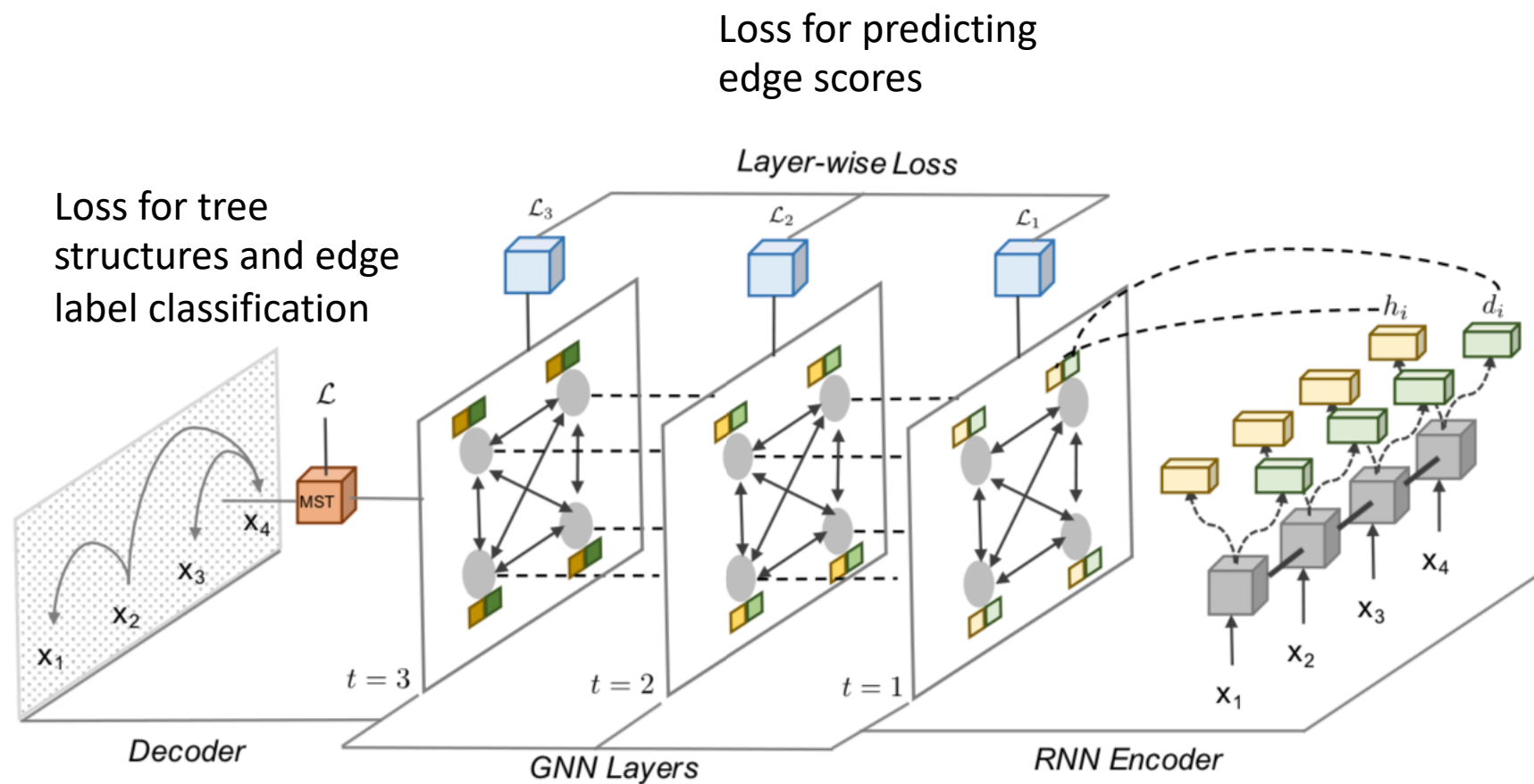


3.5.2 Chu-Liu Edmonds Algorithm



1. Initialization: for each v , choose the incoming edge with highest score.
2. Break cycles: subtract highest scores for all incoming edges, and merge nodes in cycles.
3. Split nodes, and delete edges according to Tree constrain.

Overview



Infer optimal tree given a complete graph with edge scores.

Encoding high-order information

Basic node representations

3.5.3 Learning Objective

- Loss for the output tree structure

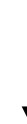
- Output from the model:

- tree structure τ
 - probability of edge labels $P(r|i, j)$
 - Probability of edge $P^\tau(i|j)$

- Ground truth: T

n : the number of edges in T

$$\mathcal{L}_0 = -\frac{1}{n} \sum_{(i,j,r) \in T} (\log P^\tau(i|j) + \log P(r|i, j))$$



$$L = \lambda_1 \mathcal{L}_0 + \lambda_2 \mathcal{L}'$$

- Layer-wise loss

- Edge scores $P^t(i|j)$ for each GCN layer

$$\mathcal{L}' = \sum_{t=1}^{\tau} \mathcal{L}_t = \sum_{t=1}^{\tau} -\frac{1}{n} \sum_{(i,j,r) \in T} \log P^t(i|j)$$

3.6 Experimental Results

- Unlabeled Attachment Score (UAS):
 - The percentage of words that have the corrected **head**.
- Labeled Attachment Score (LAS):
 - The percentage of words that have the corrected **head and label**.
- Unlabeled Complete Match (UCM):
 - The percentage of predicted **trees** that are completely corrected.
- Labeled Complete Match (LCM):
 - The percentage of predicted **trees and labels** are completely corrected.

3.6.1 Main Results

Parser		Test	
		UAS	LAS
(Chen and Manning, 2014)	T	91.8	89.6
(Dyer et al., 2015)		93.1	90.9
(Ballesteros et al., 2016)		93.56	92.41
(Weiss et al., 2015)		94.26	91.42
(Andor et al., 2016)		94.61	92.79
(Ma et al., 2018) §		95.87	94.19
(Kiperwasser and Goldberg, 2016a) §	G	93.0	90.9
(Kiperwasser and Goldberg, 2016b)		93.1	91.0
(Wang and Chang, 2016)		94.08	91.82
(Cheng et al., 2016)		94.10	91.49
(Kuncoro et al., 2016)		94.26	92.06
(Zheng, 2017) §		95.53	93.94
(Dozat and Manning, 2017)		95.74	94.08
Baseline	G	95.68	93.96
Our Model §		95.97	94.31

Table 1: Results on the English PTB dataset. The § indicates parsers using high-order features. “T” represents transition-based parser, and “G” represents a graph-based parser.

3.6.2 Different Layers and Updating Methods

GNN Layer	GNN Model	Dev		Test	
		UAS	LAS	UAS	LAS
$l = 0$	Baseline	95.58	93.74	95.68	93.96
$l = 1$	$d \triangleright h$	95.75	93.84	95.83	94.15
	$h \triangleright h$	95.78	93.80	95.91	94.12
	$hd \triangleright h$	95.77	93.87	95.88	94.23
$l = 2$	$d \triangleright h$	95.80	93.85	95.88	94.17
	$h \triangleright h$	95.77	93.83	95.85	94.13
	$hd \triangleright h$	95.79	93.90	95.92	94.24
$l = 3$	$d \triangleright h$	95.74	93.78	95.87	94.14
	$h \triangleright h$	95.75	93.80	95.90	94.15
	$hd \triangleright h$	95.71	93.82	95.93	94.22

Table 2: Impact of l and different high-order information integration methods on PTB dataset. “ $d \triangleright h$ ” corresponds to the Equation 7, “ $h \triangleright h$ ” corresponds to the Equation 6, “ $hd \triangleright h$ ” corresponds to the Equation 8.

GNN Layer	GNN Model	Dev		Test	
		UAS	LAS	UAS	LAS
$l = 2$	Synch	95.79	93.90	95.92	94.24
	H-first	95.88	93.94	95.97	94.31
	D-first	95.78	93.91	95.95	94.27

Table 3: Impact of different GNN update methods on PTB dataset. “Synch” is our default synchronized setting (Equation 8). “H-first” is an asynchronous update method that first updates head word representation (Equation 9). Similarly, the “D-first” model first updates dependent word representation.

3.6.2 Different Layers and Updating Methods

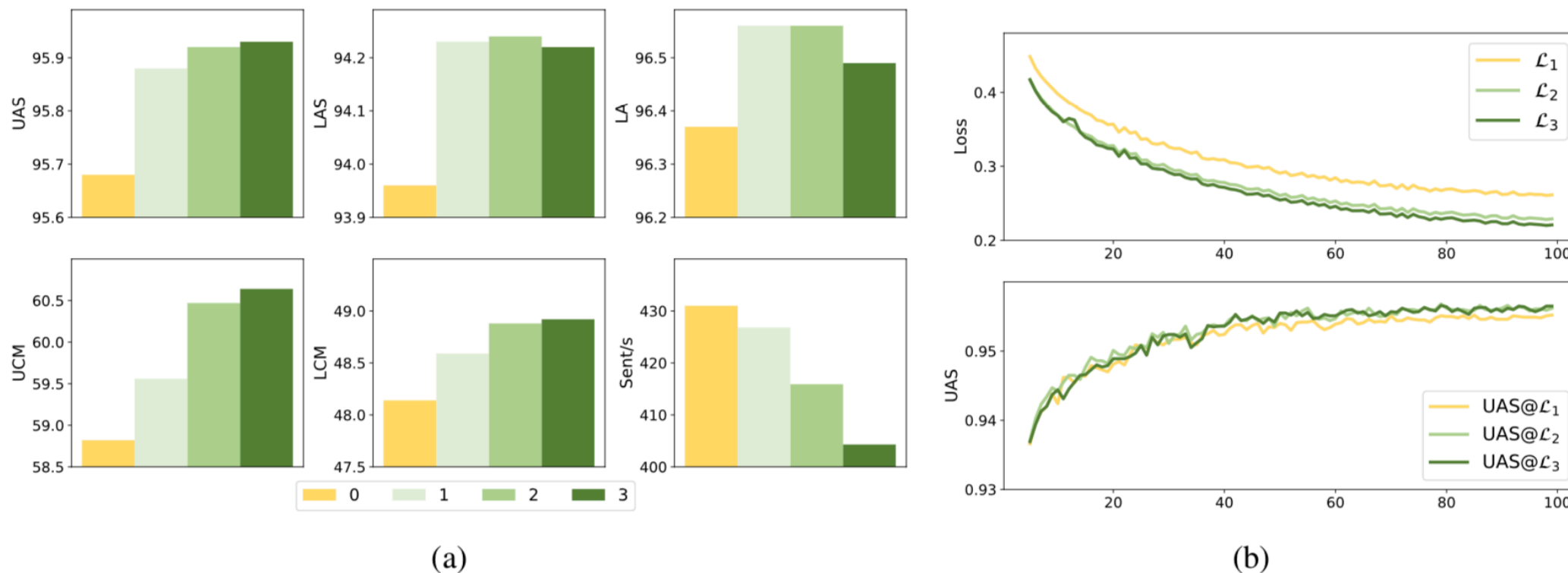


Figure 3: (a) Parsing performance and speed of different layers of our $hd \triangleright h$ model on the test set. (b) Layer-wise training loss and development set's UAS of our 3-layer $hd \triangleright h$ model.

3.6.3 Different Scoring Functions

GNN Layer	GNN Model	Dev		Test	
		UAS	LAS	UAS	LAS
$l = 2$	All=1	95.71	93.73	95.76	94.07
	Hard-1	95.69	93.70	95.80	94.13
	Hard-2	95.73	93.78	95.90	94.20
	Hard-3	95.81	93.88	95.88	94.20
$l = 2$	Soft	95.88	93.94	95.97	94.31

Table 4: Impact of different kinds of graph weights on PTB dataset. “All=1” means setting all weights to 1 (Equation 12), “Hard-k” means renormalization at the top-k weights of each node (Equation 11), “Soft” is our default model setting (Equation 8).

3.6.4 Error Analysis

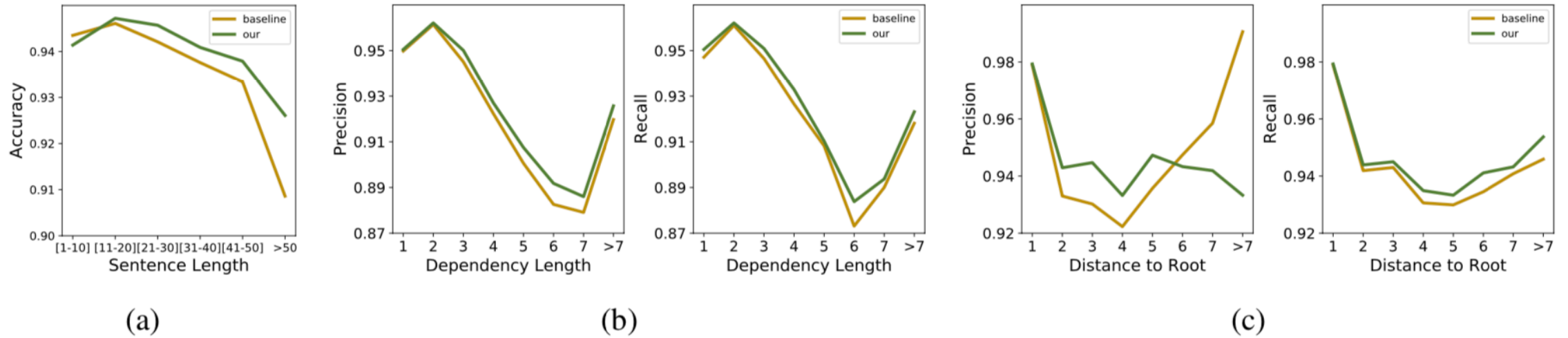


Figure 4: Parsing performance of baseline and our best parser relative to length and graph factors.

References

1. Jurafsky, Daniel, and James H. Martin. "Speech and Language Processing: An introduction to speech recognition, computational linguistics and natural language processing." *Upper Saddle River, NJ: Prentice Hall* (2008).
2. Marcus, Mitchell, et al. "The Penn Treebank: annotating predicate argument structure." *Proceedings of the workshop on Human Language Technology*. Association for Computational Linguistics, 1994.
3. Ji, Tao, Yuanbin Wu, and Man Lan. "Graph-based dependency parsing with graph neural networks." *Proceedings of ACL*. 2019.