

# SemSUM: Semantic Dependency Guided Neural Abstractive Summarization

**Hanqi Jin,<sup>1,2,3\*</sup> Tianming Wang,<sup>1,3\*</sup> Xiaojun Wan<sup>1,2,3</sup>**

<sup>1</sup>Wangxuan Institute of Computer Technology, Peking University

<sup>2</sup>Center for Data Science, Peking University

<sup>3</sup>The MOE Key Laboratory of Computational Linguistics, Peking University  
{jinhnqi, wangtm, wanxiaojun}@pku.edu.cn

# Task

Source	even some of his rivals in the roland garros locker-room are hoping that roger federer can create a bit of tennis history by winning all four grand slams .
Target	even fellow players keen for federer grand slam dream picture
Seq2seq	federer hoping to win at roland garros
SemSUM	even some rivals hope federer to make grand slam history

Table 1: Example summaries of a sentence with and without semantics guidance.

# Task

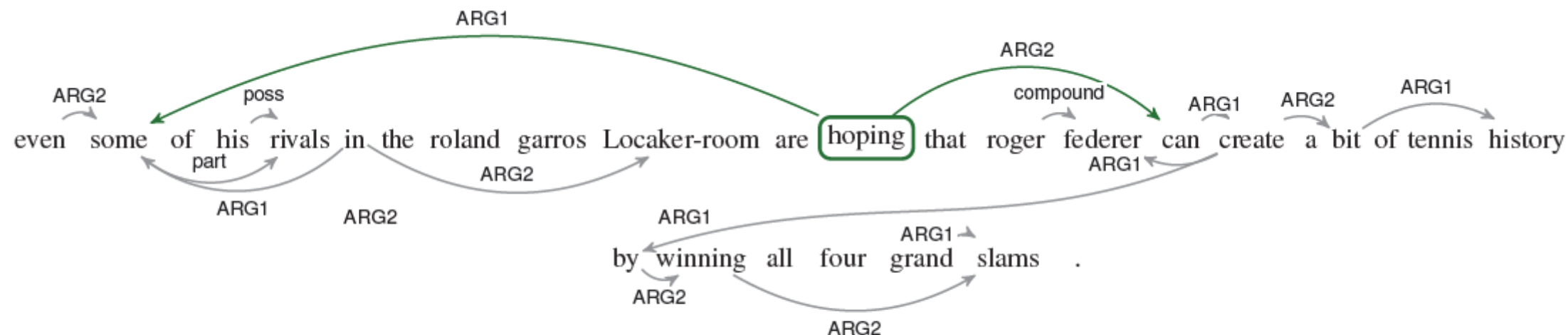


Figure 1: An example sentence annotated with a semantic dependency graphs. The green color represents the dependency of root node “hoping”. Some dependency edges are omitted for display.

# Model

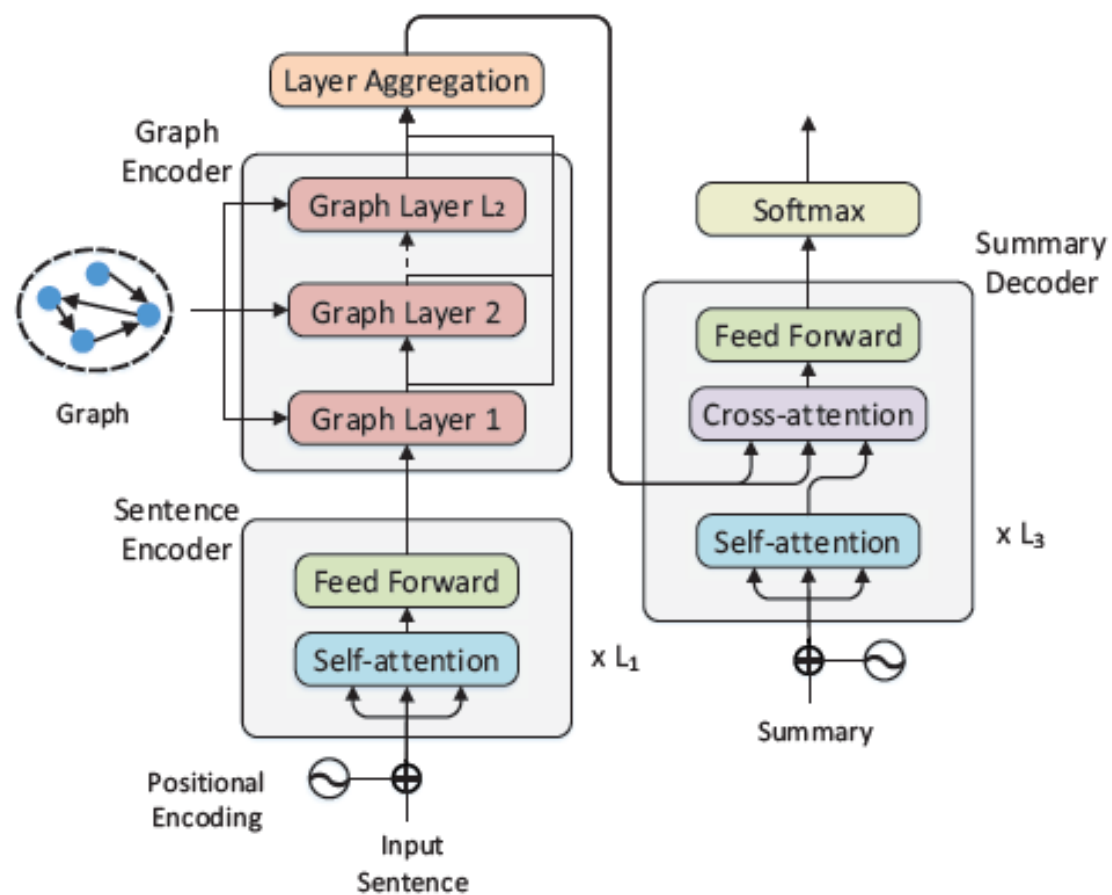


Figure 2: The overview of our SemSUM model

# Graph Encoder

- **Graph attention mechanism**

$$u_{out} = \text{ReLU}([g_{head} \parallel e_{type} \parallel g_{tail}]W_{out} + b_{out})$$

$$u_{in} = \text{ReLU}([g_{head} \parallel e_{type} \parallel g_{tail}]W_{in} + b_{in})$$

$$\hat{g}_v = \sum_{u \in \mathcal{N}(v)} \alpha(u, g_v) u W^V$$

$$\alpha(u, g_v) = \frac{\exp\left((uW^K)^\top g_v W^Q\right)}{\sum_{z \in \mathcal{N}(v)} \exp\left((zW^K)^\top g_v W^Q\right)}$$

$$\text{MHGAT}(g_v) = \left( \begin{array}{c} H \\ \parallel \\ \hat{g}_v^j \\ j=1 \end{array} \right) W^O$$

- **Layer aggregation**

- we adopt a bidirectional LSTM to aggregate the outputs of all graph layers:

$$o = \overleftarrow{h}_1 + \overrightarrow{h}_N$$

# Dataset

- **Gigaword**

- 3.8M sentence-summary pairs for training and 189K pairs for development
- 1951 sentence-summary pairs for test

- **DUC2004**

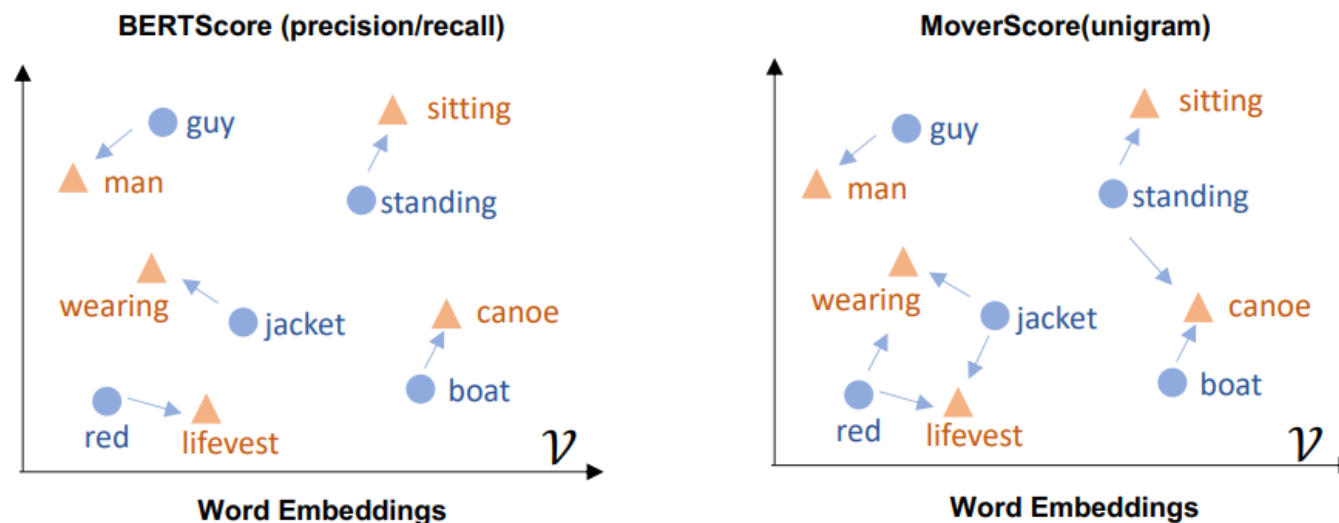
- 500 input sentence with each sentence paired with 4 different human-written reference summaries for test.

- **MSR-ATC**

- 785 input sentences with each sentence paired with 3-5 summaries for test

# Metrics

- Rough
- BERTScore
- MoverScore



● System x: A **guy** with a **red jacket** is **standing** on a **boat**

▲ Ref y: A **man wearing a lifevest** is **sitting** in a **canoe**

Figure 1: An illustration of MoverScore and BERTScore.

# Result

Model	RG-1	RG-2	RG-L	WMD	BERT
ABS(Rush, Chopra, and Weston 2015)	29.55	11.32	26.42	-	-
SEASS(Zhou et al. 2017)	36.15	17.54	33.63	-	-
Re3Sum(Cao et al. 2018a)	37.04	19.03	34.46	-	-
FTSum(Cao et al. 2018b)	37.27	17.65	34.24	-	-
PostEnsemble(Kobayashi 2018)	37.52	18.55	34.86	-	-
Sun-Attention(Niu et al. 2019)	38.27	16.45	36.08	-	-
MASS(Song et al. 2019)	38.73	19.71	35.96	34.28	<b>61.56</b>
BiSET(Wang, Quan, and Wang 2019b)	<b>39.11</b>	<b>19.78</b>	<b>36.87</b>	33.79	61.24
Transformer	36.69	18.08	34.22	32.50	60.48
TFM&GCN	37.51	19.03	34.89	33.67	61.02
SemSUM	38.78	<b>19.75</b>	36.09	<b>34.39</b>	<b>61.56</b>

Table 2: ROUGE F1, WMD unigram and BERTScore F1 evaluation results on the Gigaword test set.

Model	RG-1	RG-2	RG-L	WMD	BERT
ABS(Rush, Chopra, and Weston 2015)	26.55	7.06	22.05	-	-
SEASS(Zhou et al. 2017)	29.21	9.56	25.51	-	-
ERAML(Li et al. 2018)	29.33	10.24	25.24	-	-
WACNNs(Yuan et al. 2019)	30.54	10.87	26.94	-	-
Transformer	29.78	9.61	25.85	22.95	56.36
TFM&GCN	30.24	10.44	26.32	24.80	57.21
SemSUM	<b>31.00</b>	<b>11.11</b>	<b>26.94</b>	<b>26.71</b>	<b>57.99</b>

Table 3: ROUGE recall, WMD unigram and BERTScore F1 evaluation results on the DUC2004 test set.

Model	RG-1	RG-2	RG-L	WMD	BERT
ABS(Rush, Chopra, and Weston 2015)	20.27	5.26	17.10	-	-
SEASS(Zhou et al. 2017)	25.75	10.63	22.90	-	-
Transformer	29.29	12.45	25.93	13.76	54.31
TFM&GCN	32.53	15.41	29.26	15.61	55.48
SemSUM	<b>33.82</b>	<b>17.08</b>	<b>30.62</b>	<b>17.14</b>	<b>56.19</b>

Table 4: ROUGE F1, WMD unigram and BERTScore F1 evaluation results on the MSR-ATC test set.



# Result

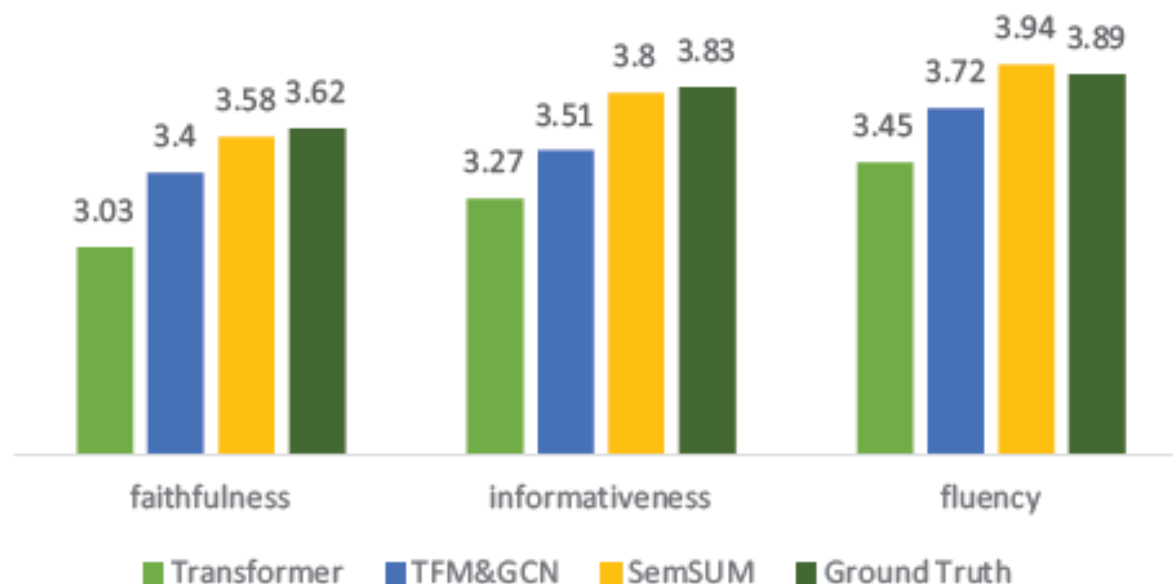


Figure 4: Human evaluation. They are rated on a Likert scale of 1(worst) to 5(best).

# Ablation study

Model	RG-1	RG-2	RG-L
SemSUM	48.25	26.54	45.14
only sentence encoder	47.14	24.94	44.00
only graph encoder	47.28	25.21	44.17
only graph encoder*	47.50	25.51	44.39
without layer aggregation	48.02	25.83	44.82
without separated mappings	47.82	25.57	44.59

Table 5: ROUGE F1 evaluation results on the development set of ablation study.  $\star$  denotes taking the adjacent relation as an edge type.

# Case study

Source	german parliament called on the international olympic committee on thursday to do more for women in sport .
Target	olympics told to help women
Transformer	german parliament calls for more women in sport
SemSUM	german parliament urges ioc to do more for women
Source	democrats in georgia and alabama , borrowing an idea usually advanced by conservative republicans , are promoting bible classes in the public schools .
Target	democrats in southern states push bills on bible study
Transformer	bible classes in public schools are promoting bible classes
SemSUM	democrats promote bible classes in public schools

Table 6: Case Study.

# **Encoding Syntactic Constituency Paths for Frame-Semantic Parsing with Graph Convolutional Networks**

**Emanuele Bastianelli, Andrea Vanzo, and Oliver Lemon**

Interaction Lab, MACS, Heriot-Watt University, Edinburgh, UK

`{e.bastianelli, a.vanzo, o.lemon}@hw.ac.uk`

# Task

- **Target identification**
- **Frame identification**
- **Argument identification**
- **Argument classification**

# Task

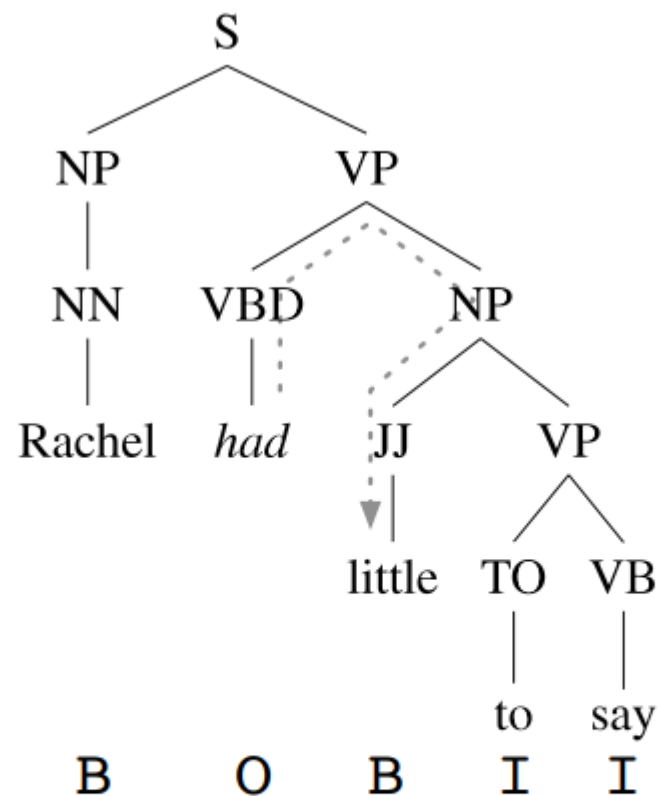


Figure 1: Constituency path between the target predicate *had* and the word *little*. At the bottom, the IOB tagging of the corresponding argument spans.

# Encoder

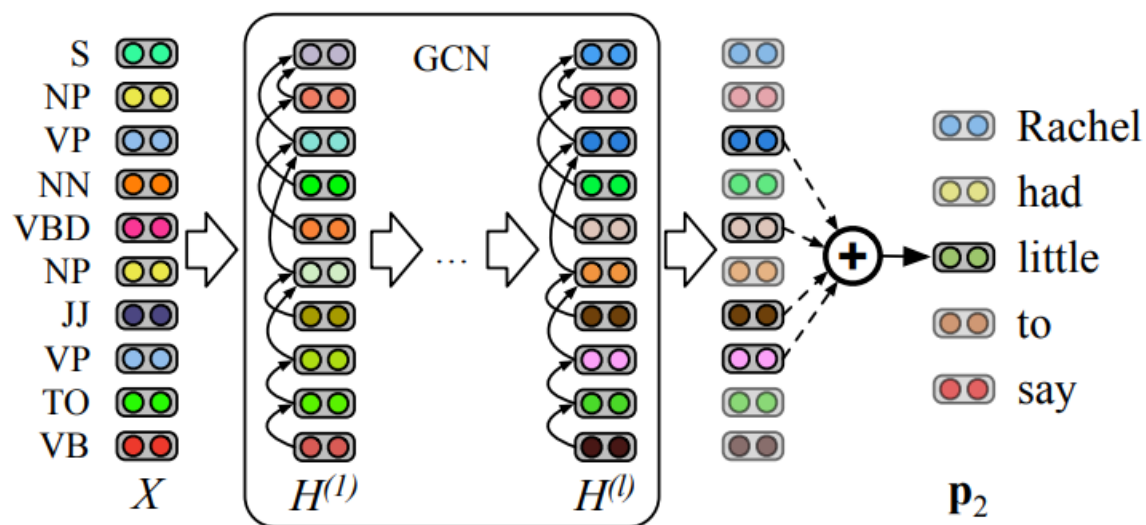
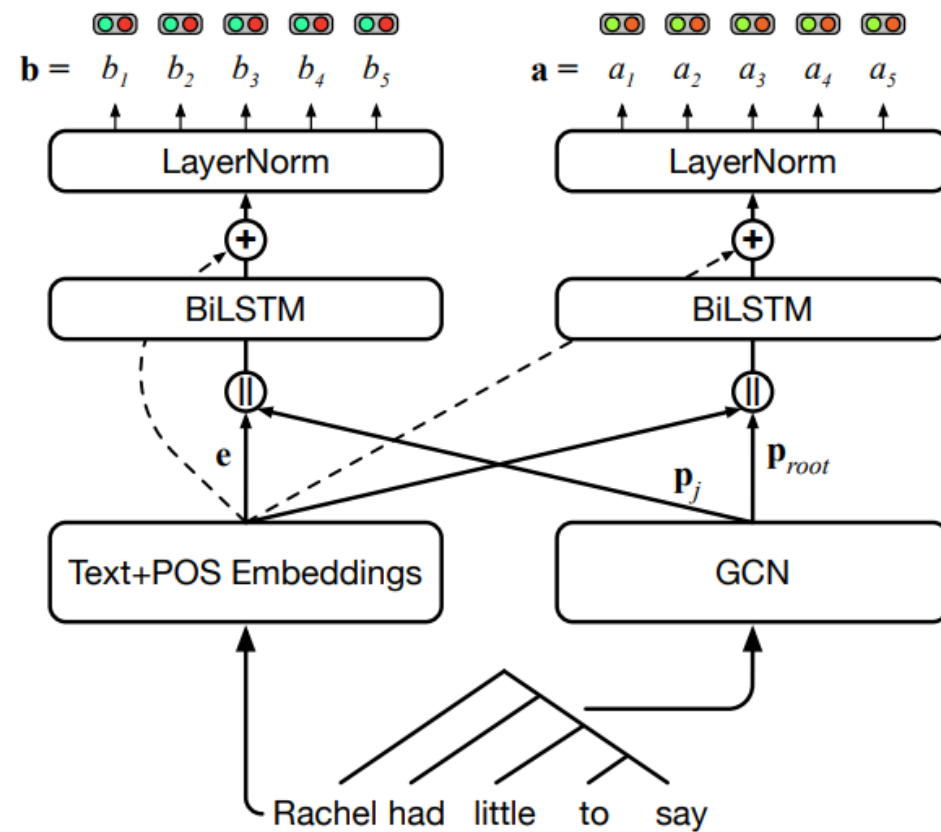


Figure 2: Process of learning constituent encodings and evaluating path features  $\mathbf{p}_2$  from the node *had* over the tree in Figure 1. Detail of feature evaluation for the word *little*.



# Decoder

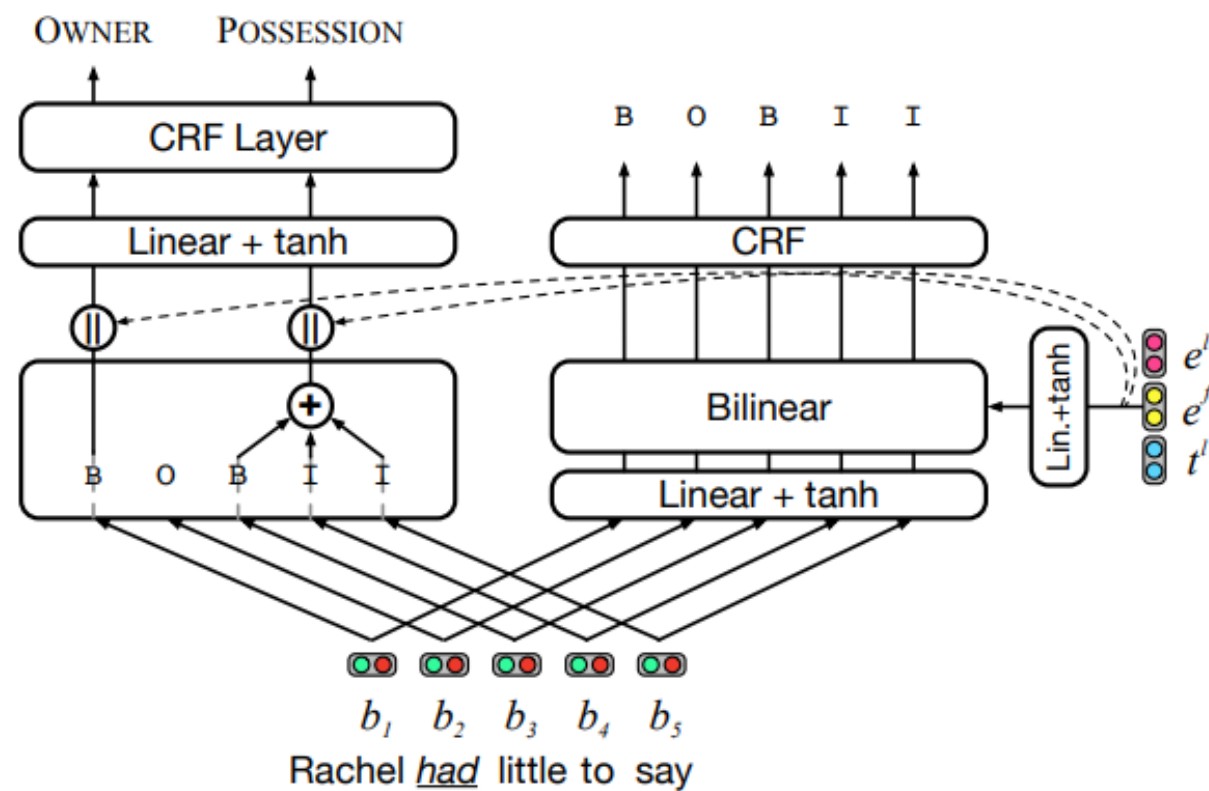


Figure 4: AI and AC sub-networks on the right and on the left, respectively.



# Model

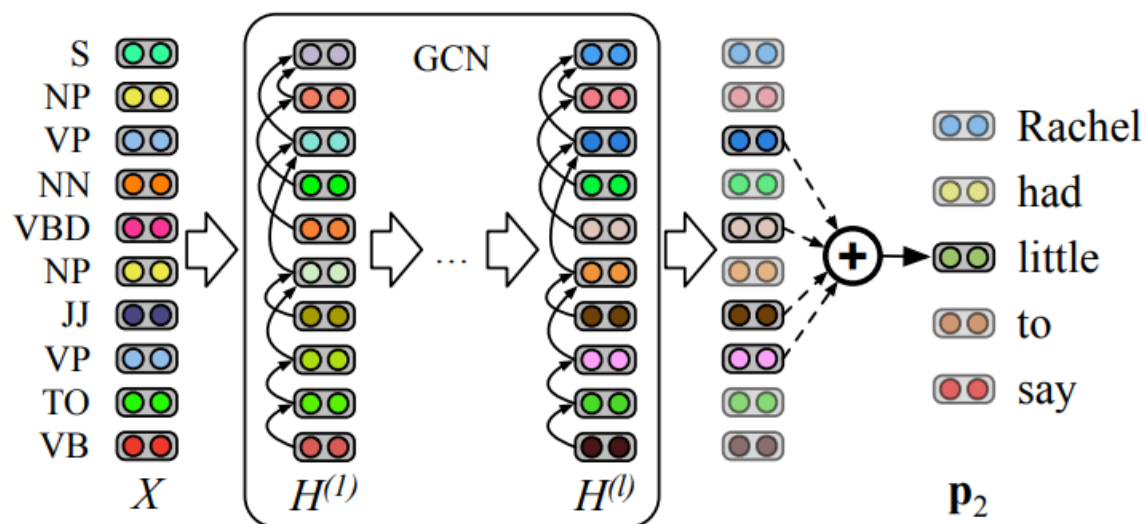
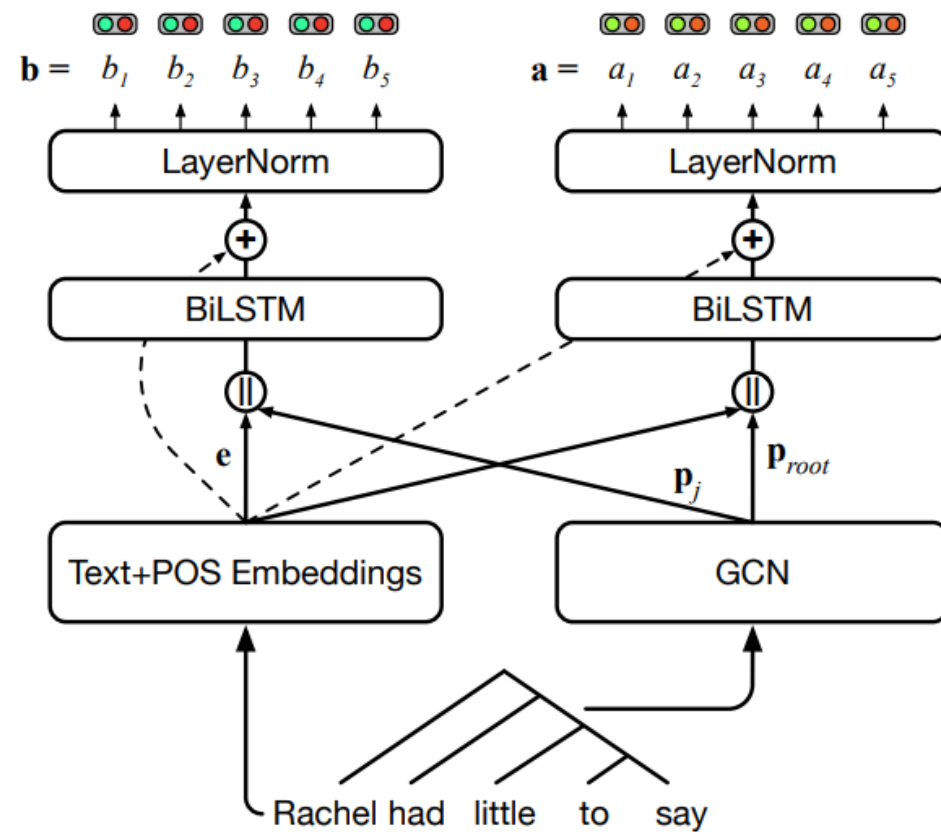


Figure 2: Process of learning constituent encodings and evaluating path features  $\mathbf{p}_2$  from the node *had* over the tree in Figure 1. Detail of feature evaluation for the word *little*.



# Result

Model	P	R	F1
Das et al. (2014)	37.5	57.5	45.4
Swayamdipta et al. (2017)	—	—	73.2
<b>OURS</b> -ELMo	69.84	78.88	74.08
<b>OURS</b> -ELMo-GCN	70.08	78.72	74.15
<b>OURS</b> -BERT	71.71	78.51	74.96
<b>OURS</b> -BERT-GCN	70.66	<b>84.12</b>	<b>76.80</b>
<b>OURS</b> -BERT-GCN-JL	<b>72.59</b>	79.43	75.86

Table 1: Precision, Recall and F1 of the Target Identification task.

Model	Acc
Das et al. (2014)	83.6
Hermann et al. (2014)	88.4
Hartmann et al. (2017)	87.6
Yang and Mitchell (2017)	88.2
Peng et al. (2018)	89.9
<b>OURS</b> -ELMo	88.89
<b>OURS</b> -ELMo-GCN	88.82
<b>OURS</b> -BERT	89.90
<b>OURS</b> -BERT-GCN	89.83
<b>OURS</b> -BERT-GCN-JL	<b>90.10</b>

Table 2: Frame Identification results using gold targets, in terms of Accuracy.

Model	P	R	F1
Das et al. (2014) <sup>†</sup>	65.6	53.8	59.1
Kshirsagar et al. (2015) <sup>†</sup>	66.0	60.4	63.1
Yang and Mitchell (2017) <sup>†</sup>	70.2	60.2	65.5
Swayamdipta et al. (2018)	69.2	69.0	69.1
Marcheggiani and Titov (2019)	69.8	68.8	69.3
<b>OURS</b> -ELMo	63.89	67.36	65.58
<b>OURS</b> -ELMo-GCN	72.02	73.70	72.85
<b>OURS</b> -BERT	71.19	74.26	72.69
<b>OURS</b> -BERT-GCN	74.23	<b>76.94</b>	<b>75.56</b>
<b>OURS</b> -BERT-GCN-JL	<b>74.56</b>	74.43	74.50

Table 3: Semantic Role Labeling results with gold targets and frames.

Input	SRL		
	No-GCN	GCN	gain
GloVe	49.13	68.23	+19.1
ELMo	65.58	72.85	+7.27
BERT	72.69	75.56	+2.87

Table 4: Ablation on the use of constituency features with different input layers for the SRL on FN15.

# Graph Convolutions over Constituent Trees for Syntax-Aware Semantic Role Labeling

**Diego Marcheggiani<sup>1\*</sup>**

**Ivan Titov<sup>2,3</sup>**

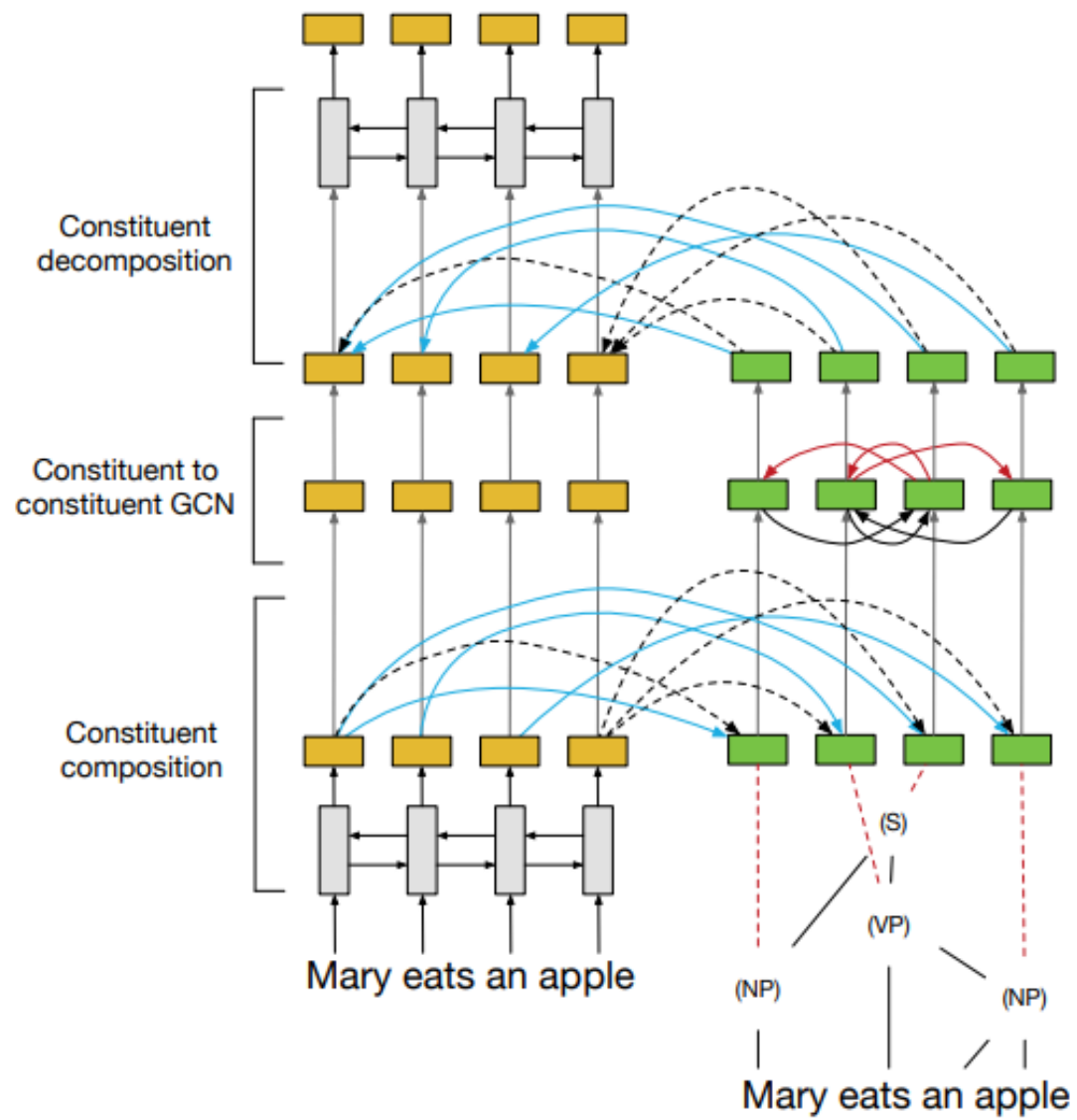
<sup>1</sup>Amazon

<sup>2</sup>ILCC, School of Informatics, University of Edinburgh

<sup>3</sup>ILLC, University of Amsterdam

marchegg@amazon.es    ititov@inf.ed.ac.uk

# Model



# Result

Model	$P$	$R$	$F_1$
Yang and Mitchell (2017) (SEQ)	63.4	66.4	64.9
Yang and Mitchell (2017) (ALL)	70.2	60.2	65.5
Swayamdipta et al. (2018) $\dagger\dagger$	69.2	69.0	69.1
SpanGCN $\dagger$	69.8	68.8	69.3

Table 4: Results on FrameNet 1.5 test set using gold frames.  $\dagger$  indicates syntactic models and  $\dagger\dagger$  indicates multi-task learning models.

Single	WSJ Test		
	$P$	$R$	$F_1$
He et al. (2017)	83.1	83.0	83.1
He et al. (2018a)	84.2	83.7	83.9
Tan et al. (2018)	84.5	85.2	84.8
Ouchi et al. (2018)	84.7	82.3	83.5
Strubell et al. (2018)(LISA) $\dagger\dagger$	84.7	84.6	84.6
SpanGCN $\dagger$	85.8	85.1	85.4

Single / Context. Emb.			
He et al. (2018a)(ELMo)	-	-	87.4
Li et al. (2019)(ELMo)	87.9	87.5	87.7
Ouchi et al. (2018)(ELMo)	88.2	87.0	87.6
Wang et al. (2019)(ELMo) $\dagger$	-	-	88.2
SpanGCN (ELMo) $\dagger$	87.5	87.9	87.7
SpanGCN (RoBERTa) $\dagger$	87.7	88.1	87.9

Single	Brown Test		
	$P$	$R$	$F_1$
He et al. (2017)	72.9	71.4	72.1
He et al. (2018a)	74.2	73.1	73.7
Tan et al. (2018)	73.5	74.6	74.1
Ouchi et al. (2018)	76.0	70.4	73.1
Strubell et al. (2018)(LISA) $\dagger\dagger$	74.8	74.3	74.6
SpanGCN $\dagger$	76.2	74.7	75.5

Single / Context. Emb.			
He et al. (2018a)(ELMo)	-	-	80.4
Li et al. (2019)(ELMo)	80.6	80.4	80.5
Ouchi et al. (2018)(ELMo)	79.9	77.5	78.7
Wang et al. (2019)(ELMo) $\dagger$	-	-	79.3
SpanGCN(ELMo) $\dagger$	79.4	79.6	79.5
SpanGCN(RoBERTa) $\dagger$	80.5	80.7	80.6

Table 2: Precision, recall and  $F_1$  on the CoNLL-2005 test sets.  $\dagger$  indicates syntactic models and  $\dagger\dagger$  indicates multi-task learning models.

Single	Test		
	$P$	$R$	$F_1$
He et al. (2017)	81.7	81.6	81.7
He et al. (2018a)	-	-	82.1
Tan et al. (2018)	81.9	83.6	82.7
Ouchi et al. (2018)	84.4	81.7	83.0
Swayamdipta et al. (2018) $\dagger\dagger$	85.1	81.2	83.8
SpanGCN $\dagger$	84.5	84.3	84.4

Single / Context. Emb.			
Peters et al. (2018a)(ELMo)	-	-	84.6
He et al. (2018a)(ELMo)	-	-	85.5
Li et al. (2019)(ELMo)	85.7	86.3	86.0
Ouchi et al. (2018)(ELMo)	87.1	85.3	86.2
Wang et al. (2019)(ELMo) $\dagger$	-	-	86.4
SpanGCN (ELMo) $\dagger$	86.3	86.8	86.5
SpanGCN (RoBERTa) $\dagger$	86.5	87.1	86.8

Table 3: Precision, recall and  $F_1$  on the CoNLL-2012 test set.  $\dagger$  indicates syntactic models and  $\dagger\dagger$  indicates multi-task learning models.