

Linguistically-Informed Neural Networks for Natural Language Processing

Joost Bastings

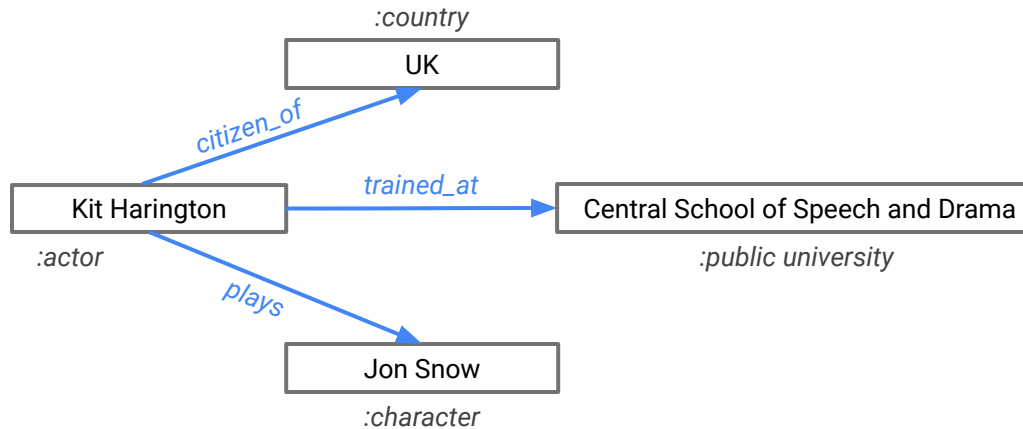
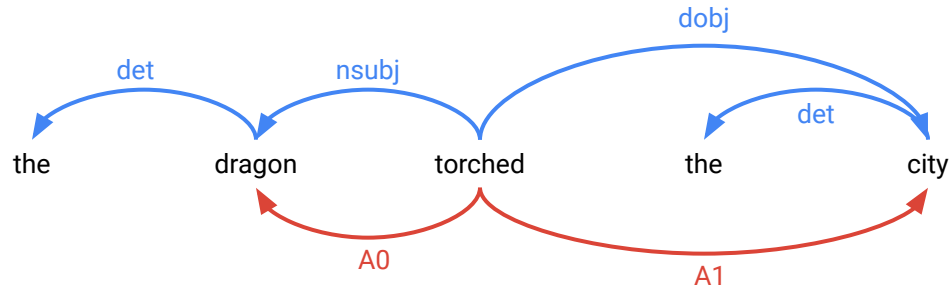
Institute for Logic, Language and Computation (ILLC), University of Amsterdam
bastings.github.io

NLP2 20-5-2019

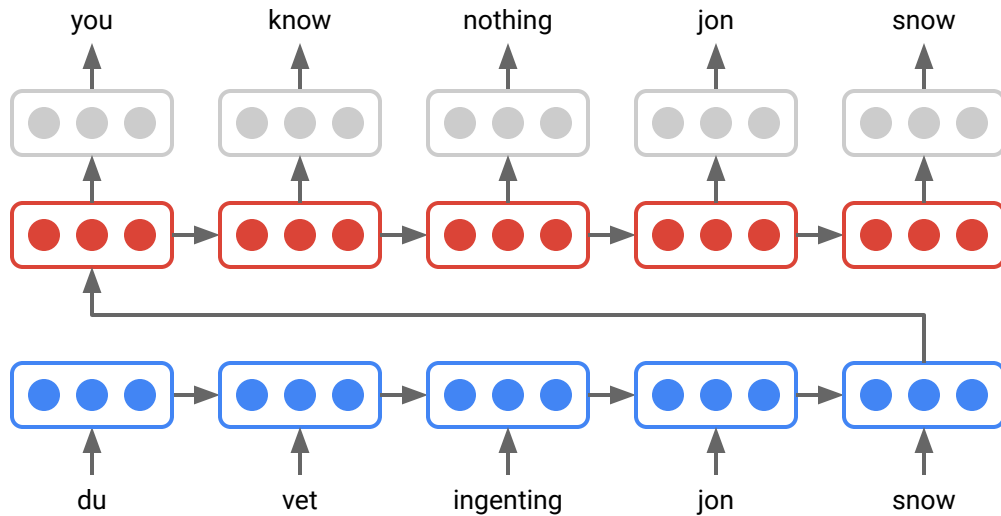
Outline

1. **Linguistic structures**
2. **Conditioning on graphs**
 - a. **Graph Convolutional Networks:** neural message passing
 - b. **Applications in NLP:** using GCNs for encoding prior knowledge and inference
3. **Inducing graphs**
 - a. **Inducing graphs for Neural Machine Translation**
 - b. **Transformer**
 - i. Self-Attention: fully-connected graphs
 - ii. Linguistically-Informed Self-Attention
 - iii. BERT

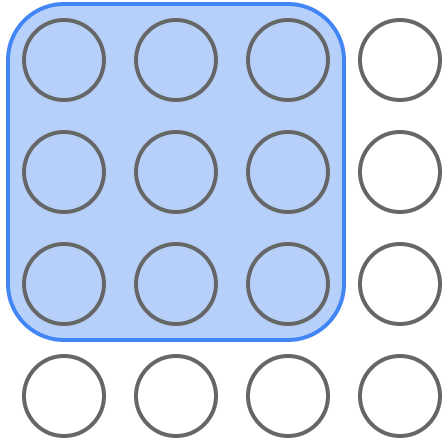
In NLP is natural to represent linguistic/prior knowledge as graphs



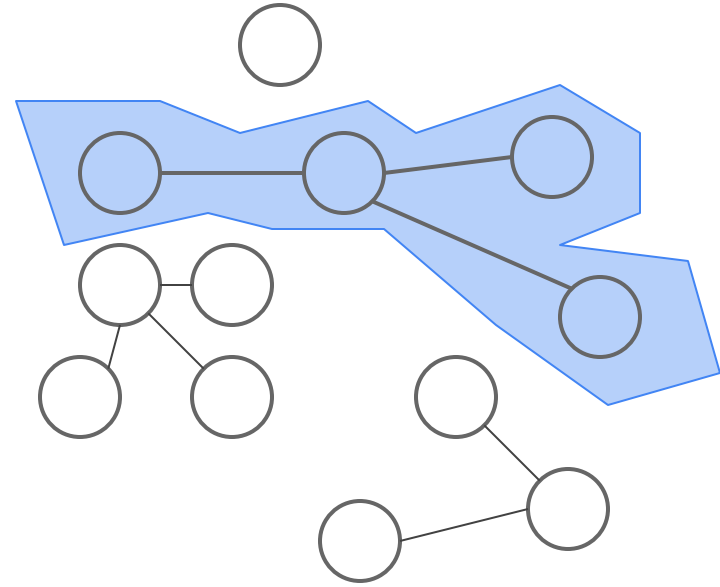
Encoder-Decoder models are agnostic to this



Convolution vs Graph Convolution

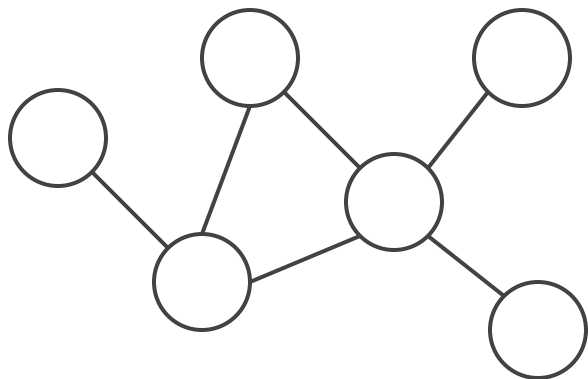


2D convolution
e.g. image filter

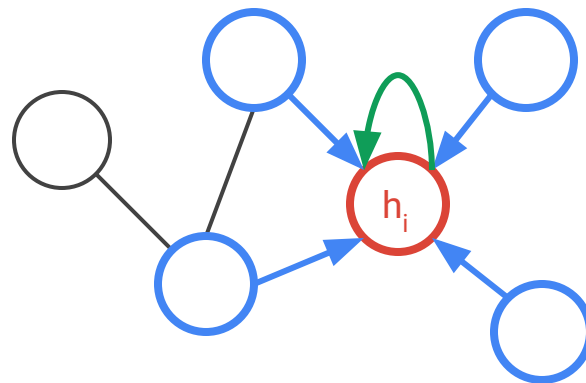


graph convolution
e.g. social network

Graph Convolutional Networks: message passing



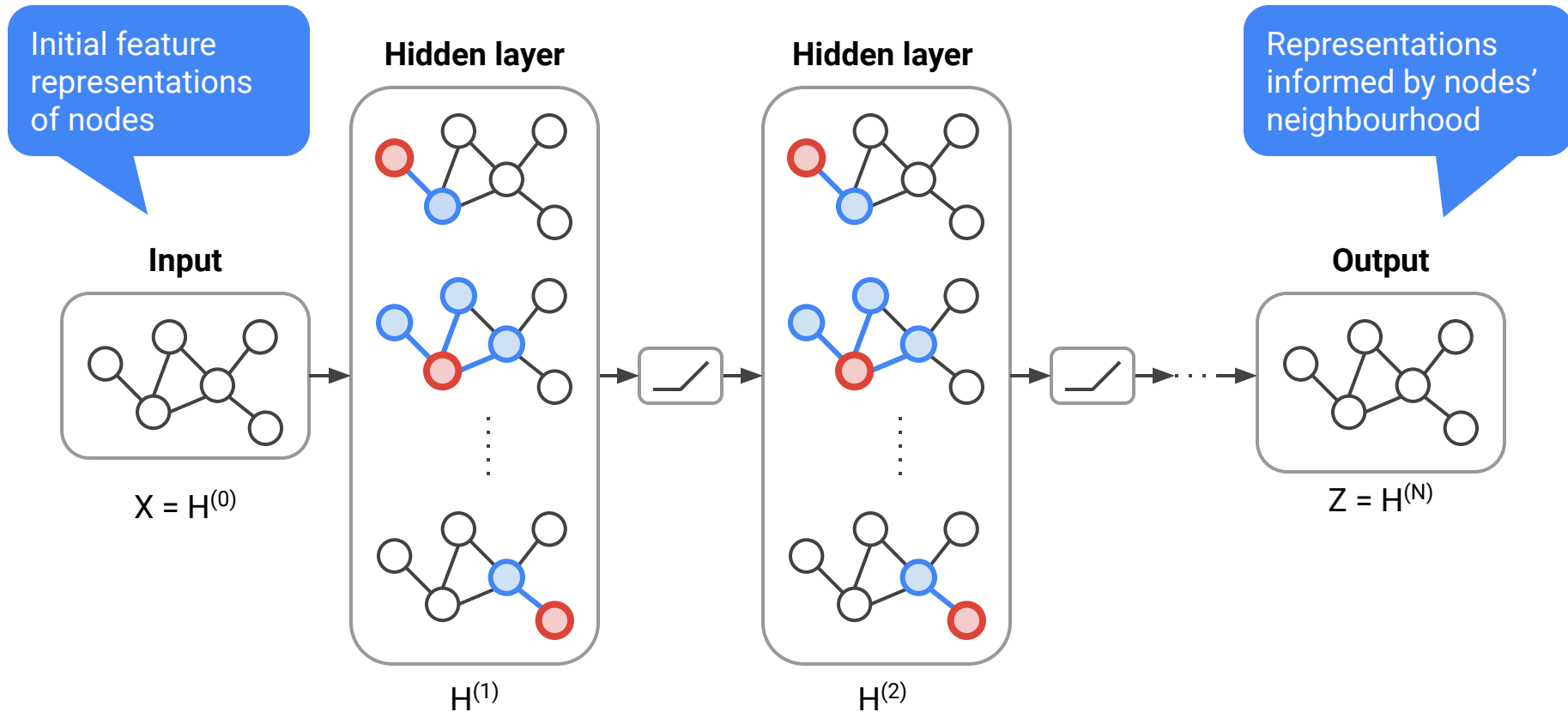
Undirected graph



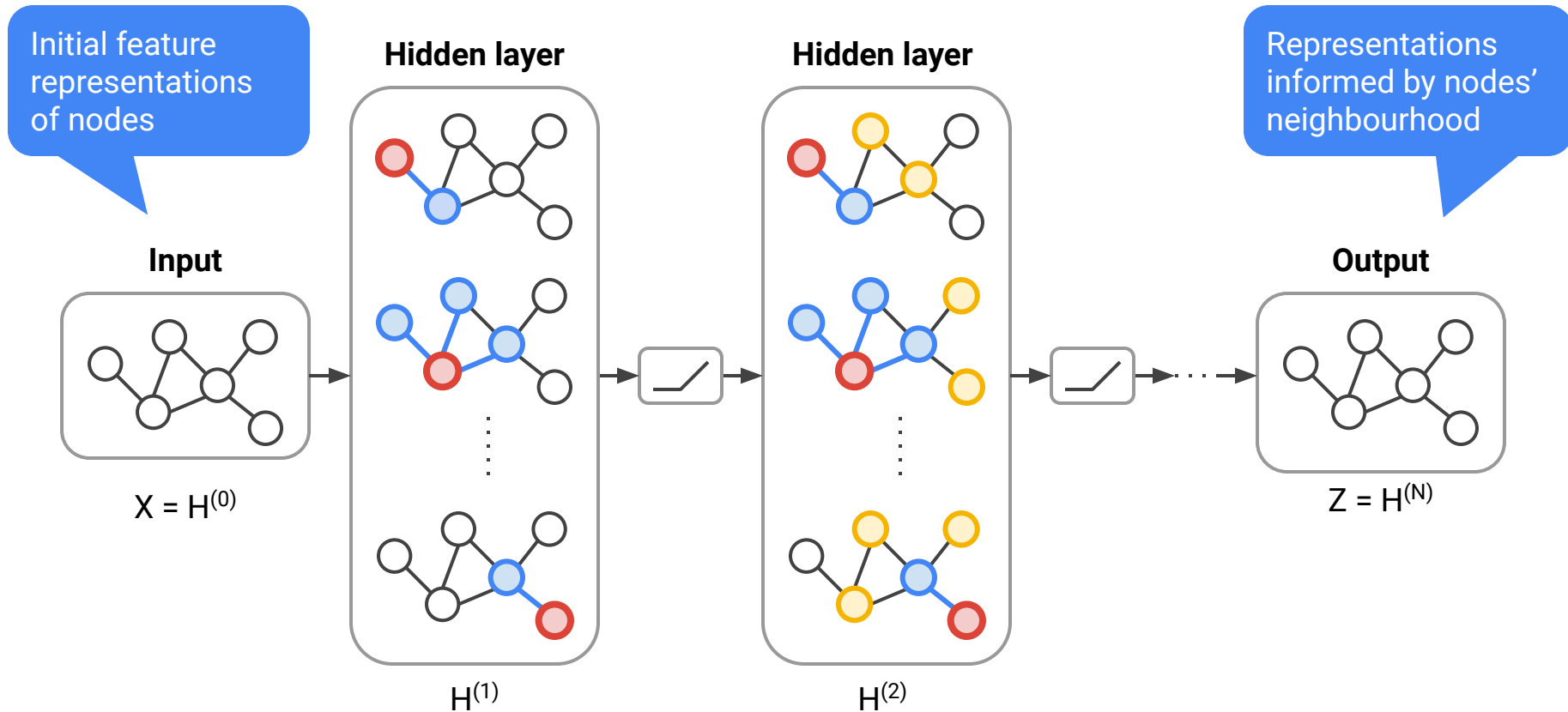
Update of red node

$$\mathbf{h}_i = \text{ReLU} \left(\mathbf{W}'\mathbf{h}_i + \sum_{j \in \text{neighbors}(i)} \mathbf{W}\mathbf{h}_j \right)$$

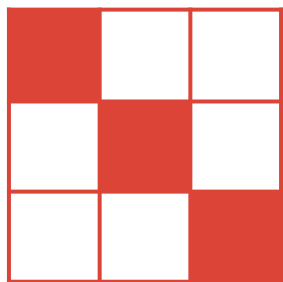
GCNs: multilayer convolution operation



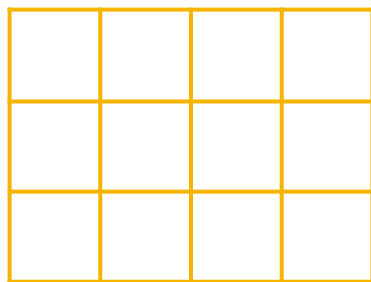
GCNs: multilayer convolution operation



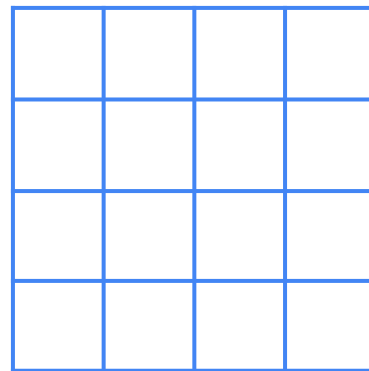
Implementation note



\hat{A}



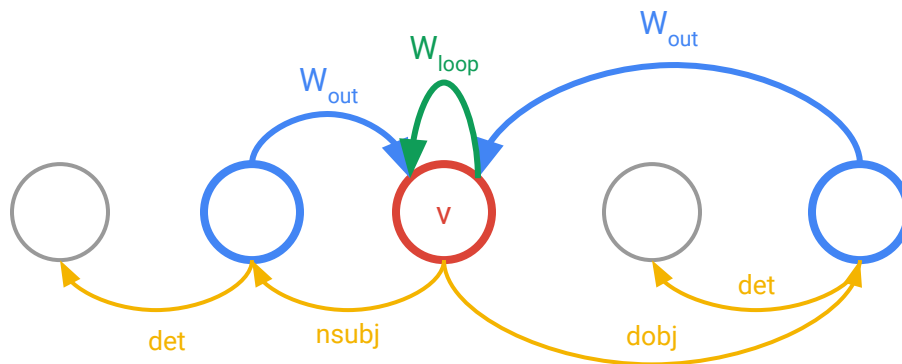
H



W

$$\hat{A} = A + I$$
$$H^{l+1} = \text{ReLU}(\hat{A}H^l W^l)$$

Syntactic GCNs: directionality and labels



Weight matrix for each direction:

$W_{out}, W_{in}, W_{loop}$

Bias for each label,
e.g. \mathbf{b}_{nsubj}

$$\mathbf{h}_v = \text{ReLU} \left(\sum_{u \in \mathcal{N}(v)} W_{\text{dir}(u,v)} \mathbf{h}_u + \mathbf{b}_{\text{lab}(u,v)} \right)$$

Syntactic GCNs: edge-wise gating

$$g_{u,v} = \sigma \left(\mathbf{h}_u \cdot \hat{\mathbf{w}}_{\text{dir}(u,v)} + \hat{b}_{\text{lab}(u,v)} \right)$$

$$\mathbf{h}_v = \text{ReLU} \left(\sum_{u \in \mathcal{N}(v)} g_{u,v} \left(W_{\text{dir}(u,v)} \mathbf{h}_u + \mathbf{b}_{\text{lab}(u,v)} \right) \right)$$

Conditioning on graphs: Syntax-aware Semantic Role Labeling

Semantic Role Labeling: “Who did what to whom?”

Example goal: identify a **stock purchase** event by **Snow Ltd.**

Many different surface forms!

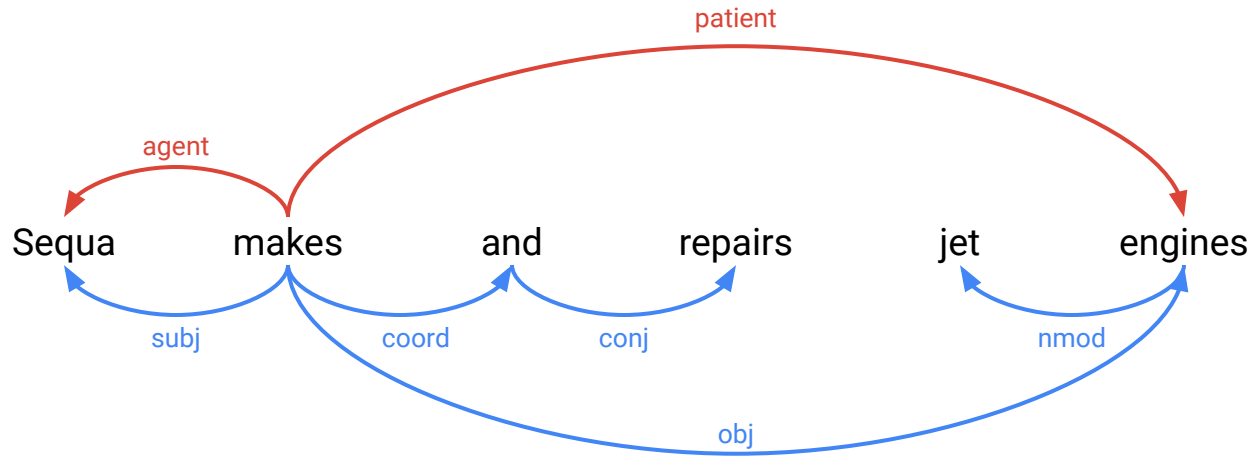
- Snow Ltd. bought the stock
- They sold the stock to Snow Ltd.
- The stock was bought by Snow Ltd.
- The purchase of the stock by Snow Ltd.
- The stock purchase by Snow Ltd. ...

Semantic roles are **abstract models** of the role an argument plays in the event described by the predicate

Roles can be predicate-specific
(A0, A1 are usually *agent* and *patient*)

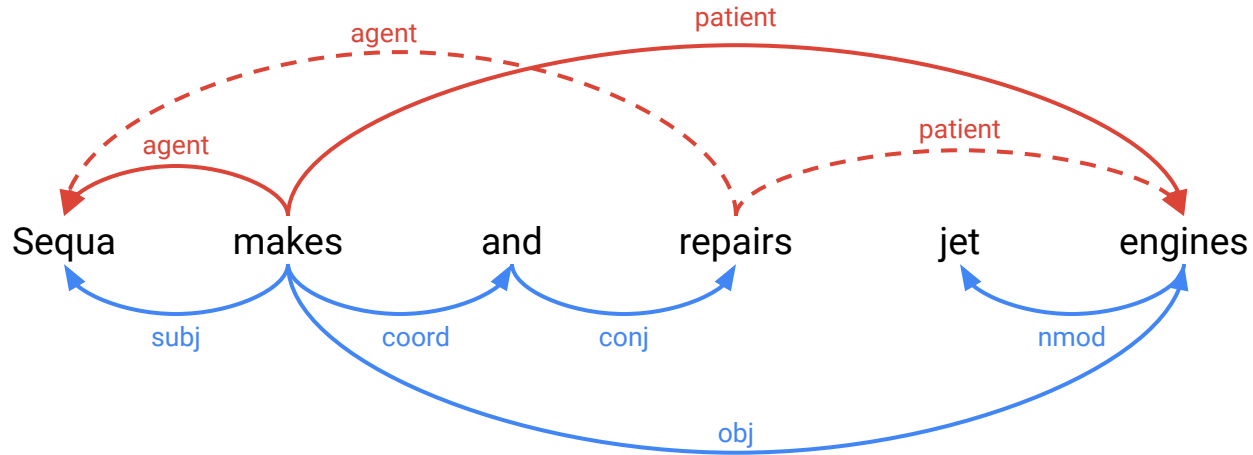
SRL is the task of assigning semantic role labels to the constituents of a sentence

Syntax/semantics interaction



Some syntactic dependencies are **mirrored** in the semantic graph

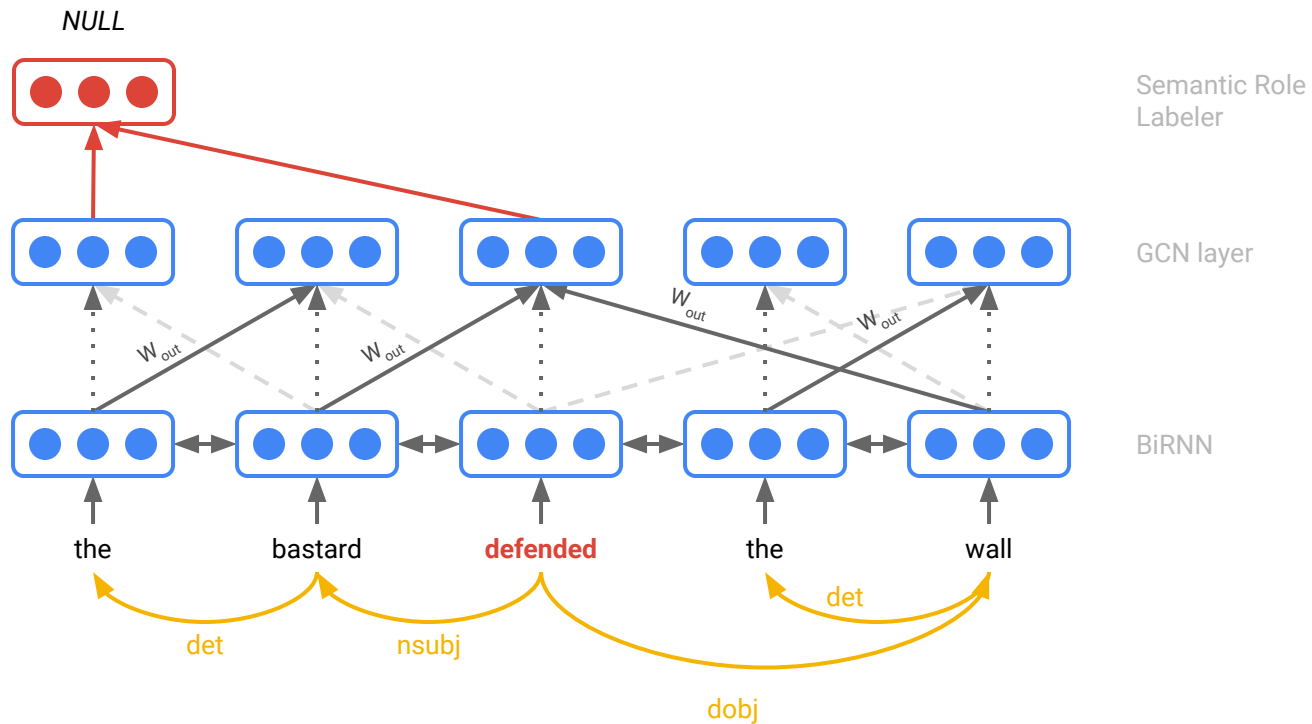
Syntax/semantics interaction



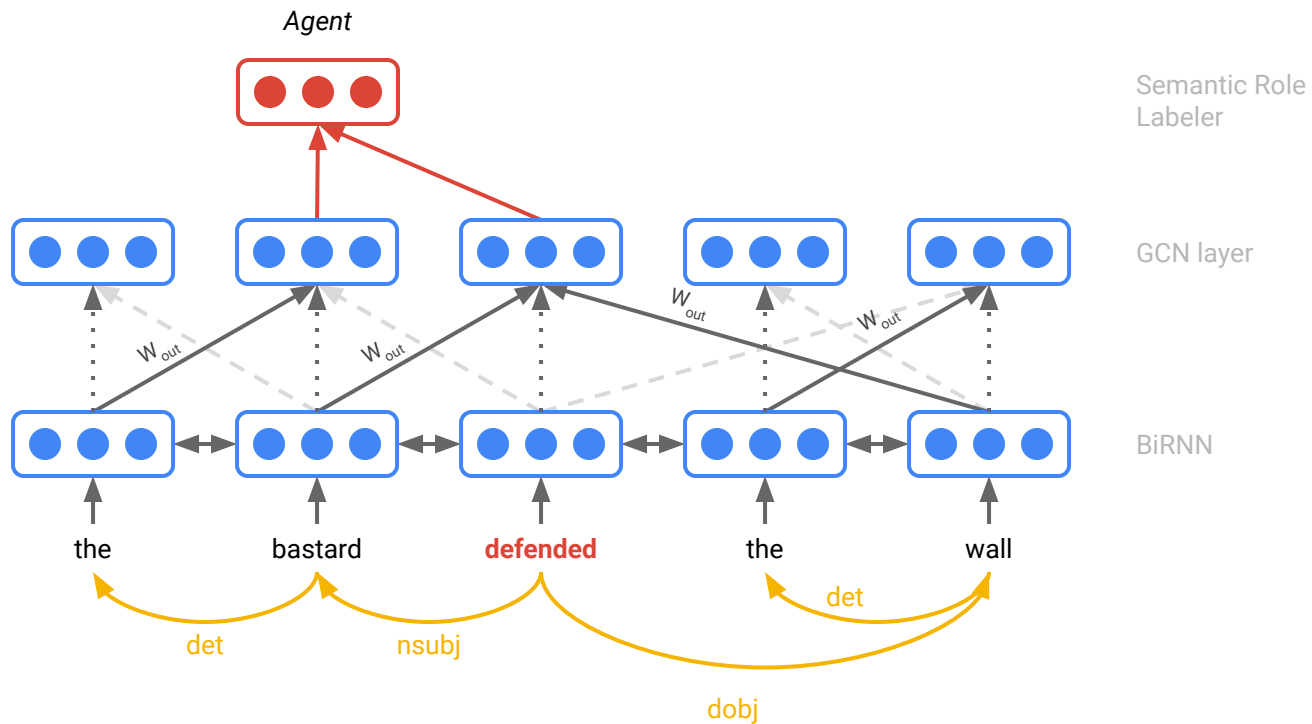
Some syntactic dependencies are **mirrored** in the semantic graph

... **but not all of them** – the syntax-semantics interface is far from trivial

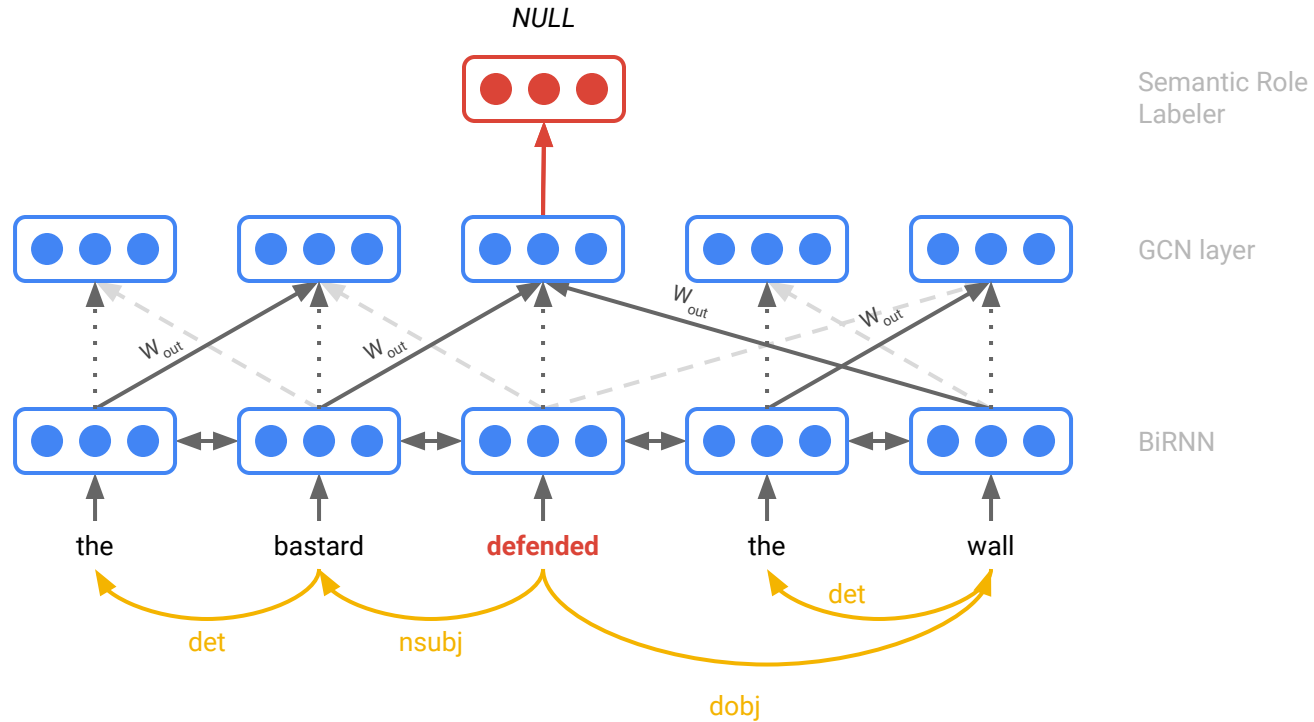
GCNs for SRL



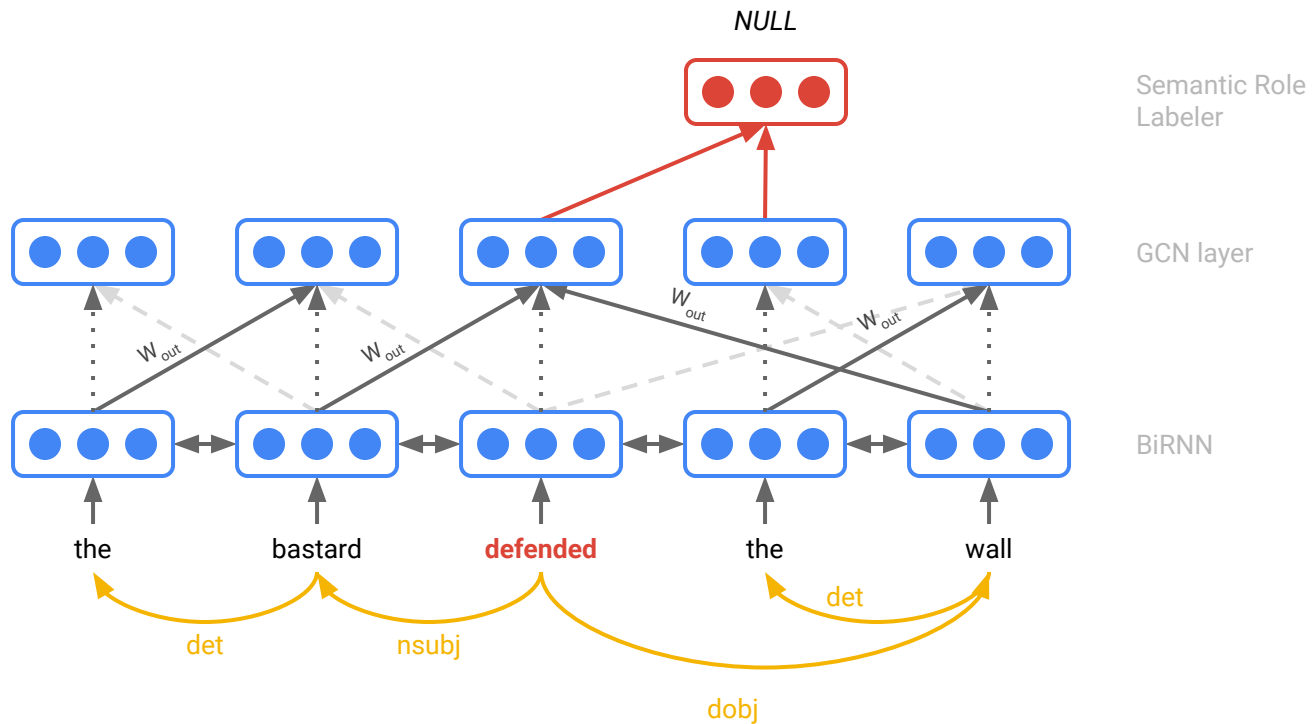
GCNs for SRL



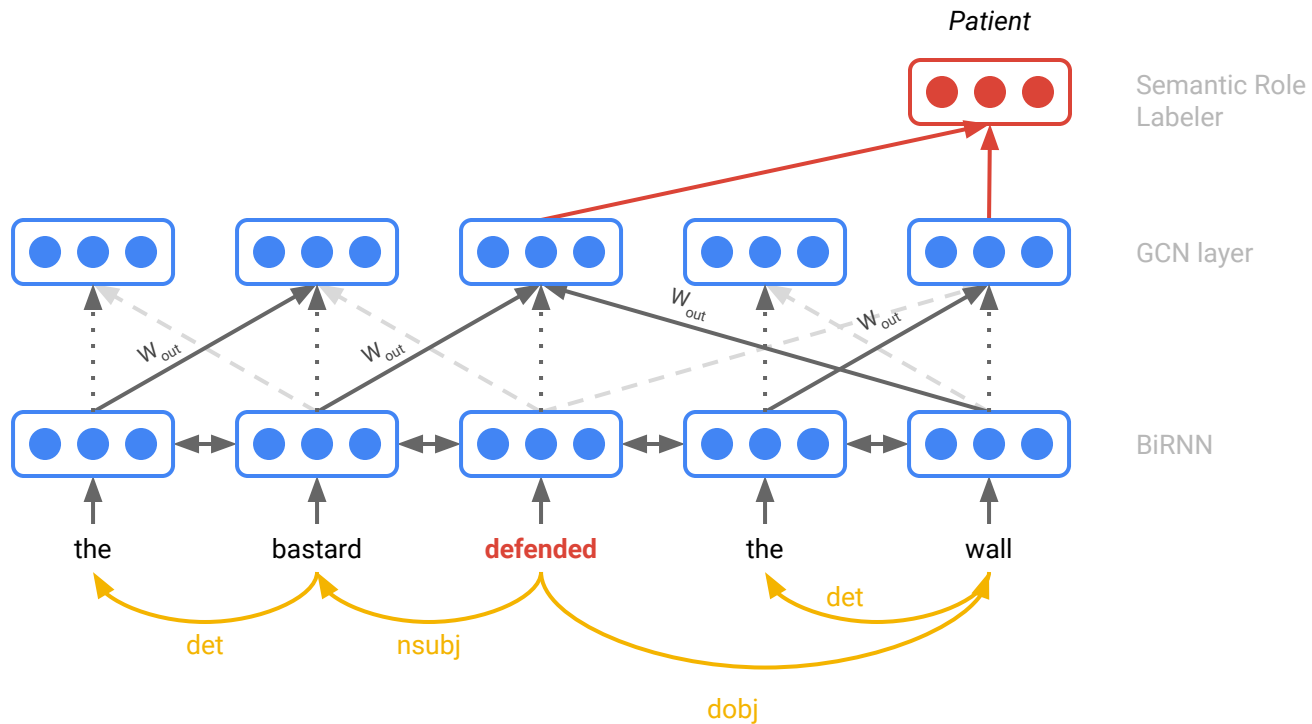
GCNs for SRL



GCNs for SRL

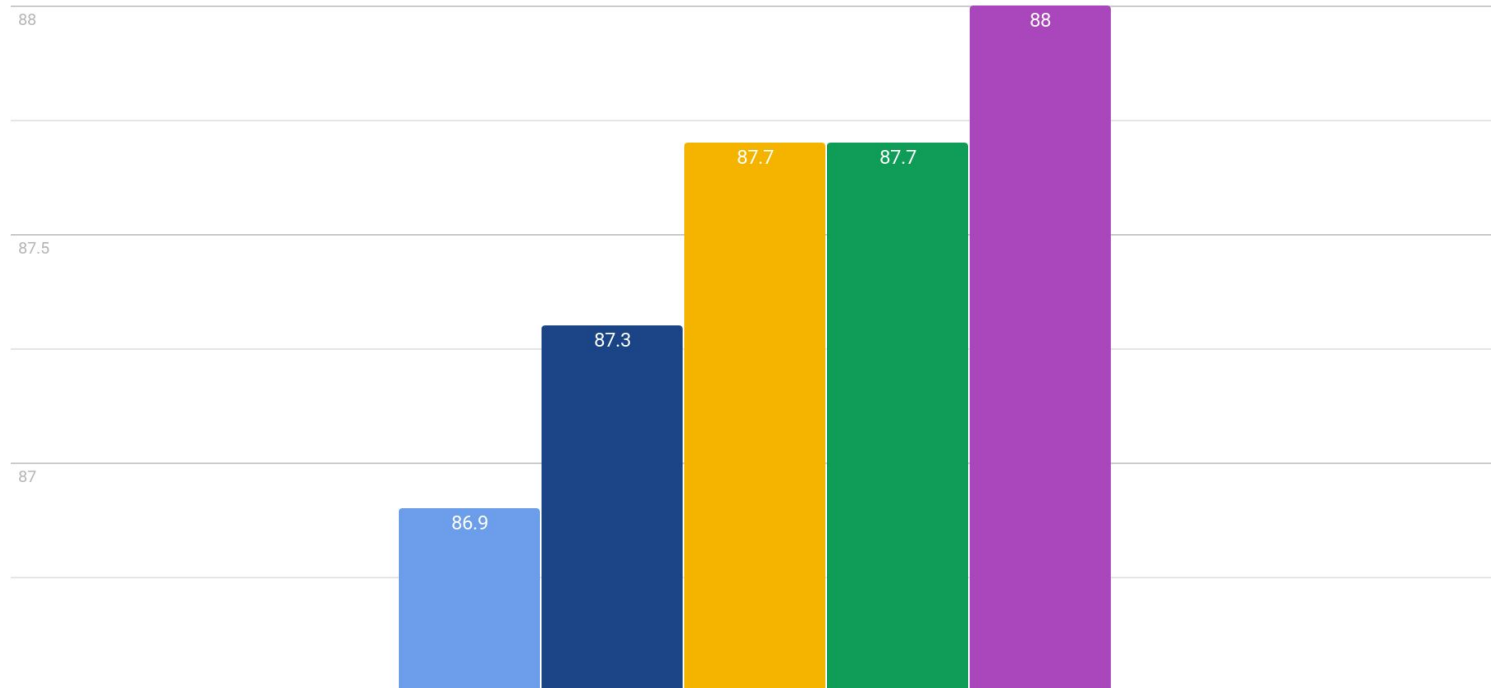


GCNs for SRL

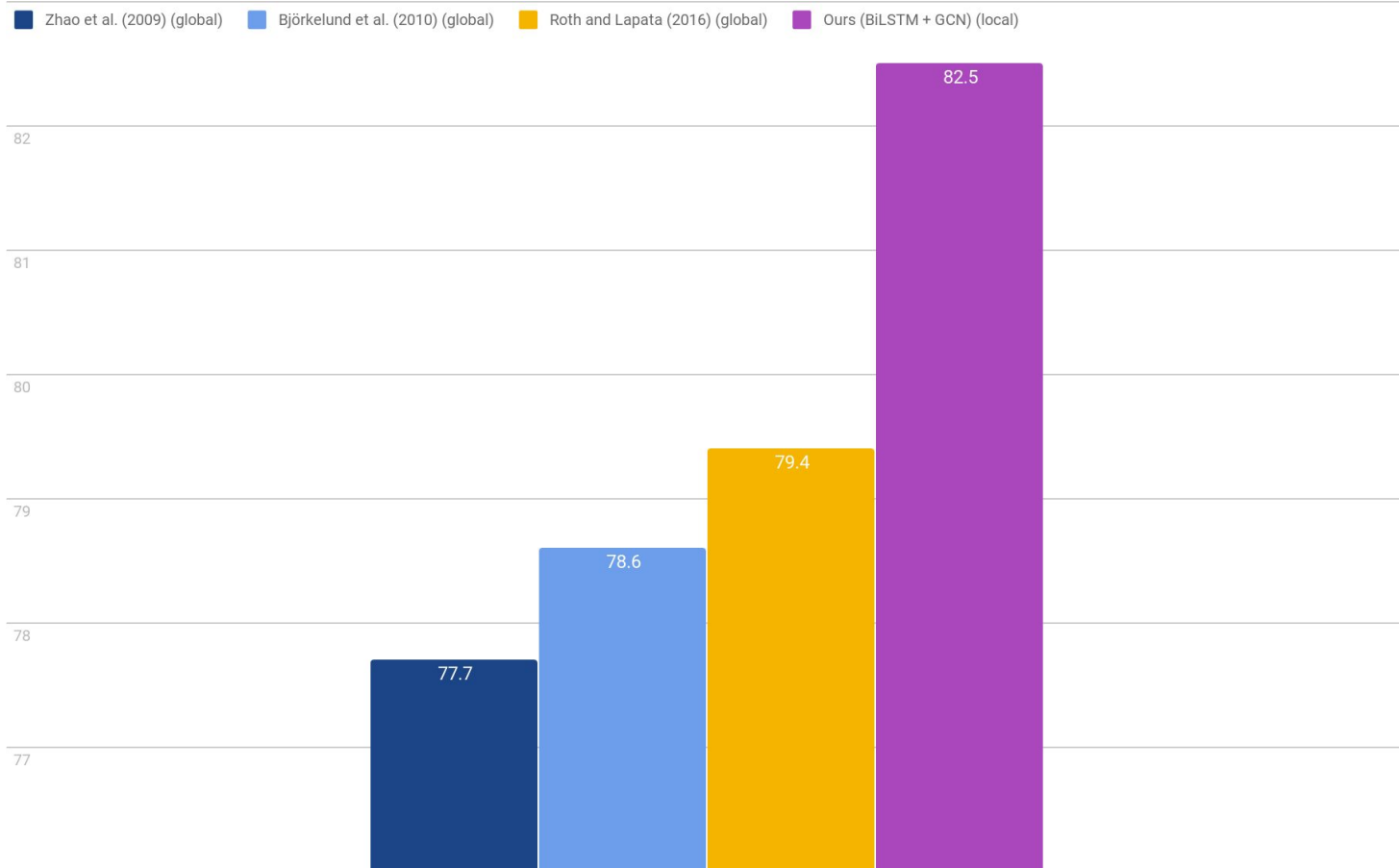


SRL Results (English)

■ Björkelund et al. (2010) (global) ■ FitzGerald et al. (2015) (global) ■ Roth and Lapata (2016) (global)
■ Marcheggiani et al. (2017) (local) ■ Ours (BiLSTM + GCN) (local)

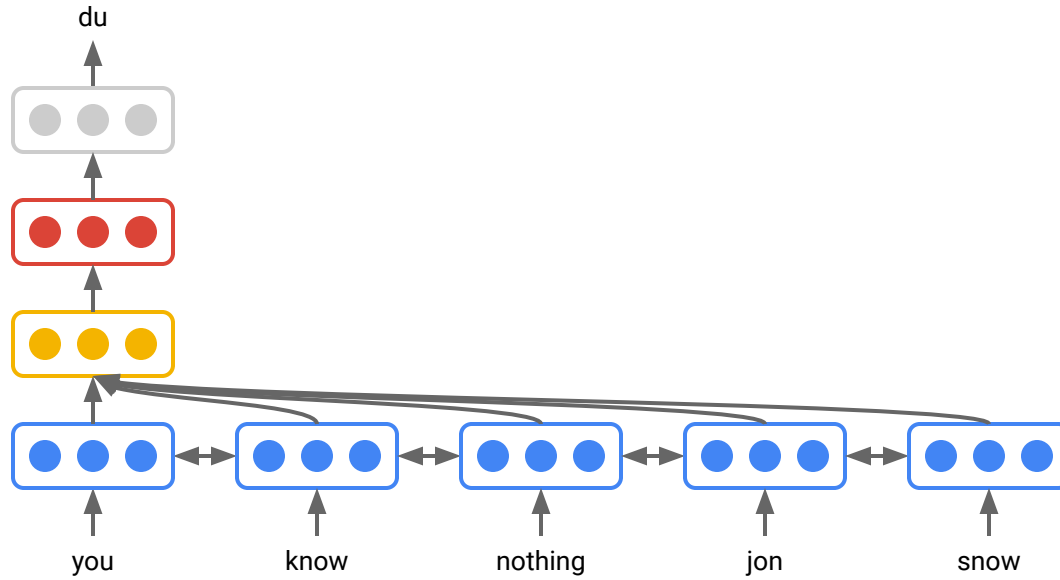


SRL Results (Chinese)

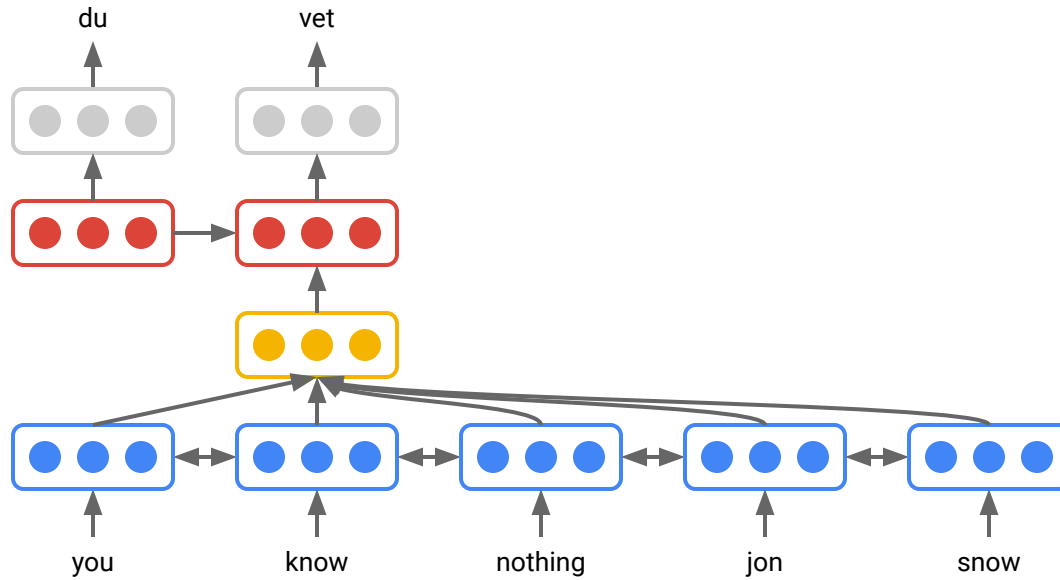


Conditioning on graphs: Linguistically-informed Neural Machine Translation

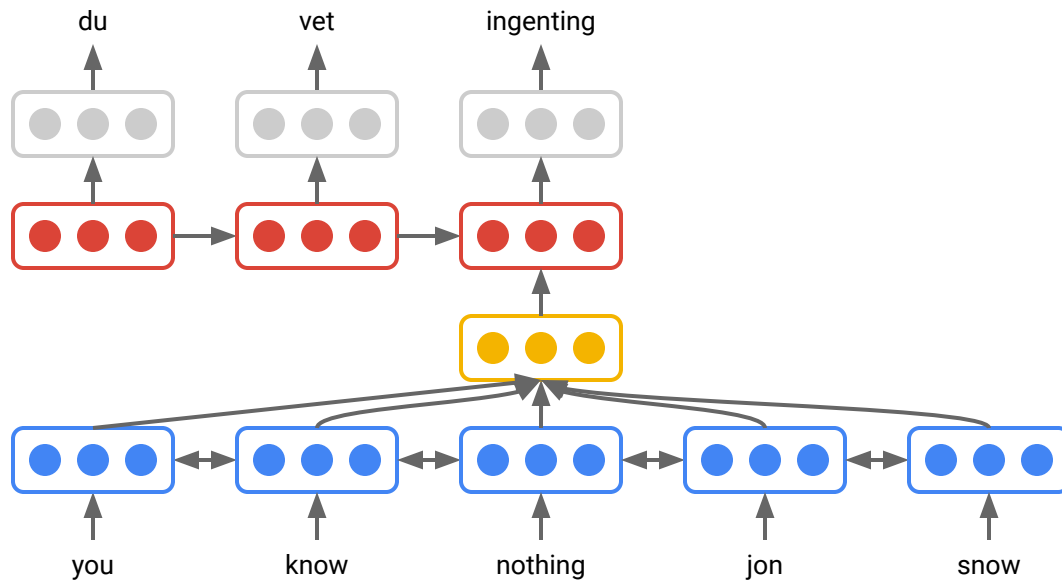
Encoder-decoder with Attention



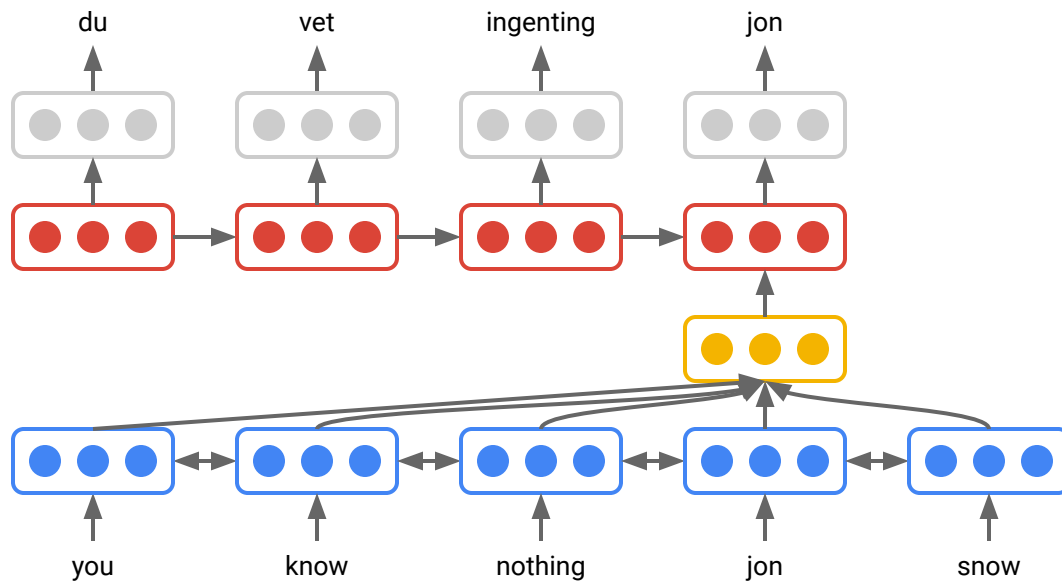
Encoder-decoder with Attention



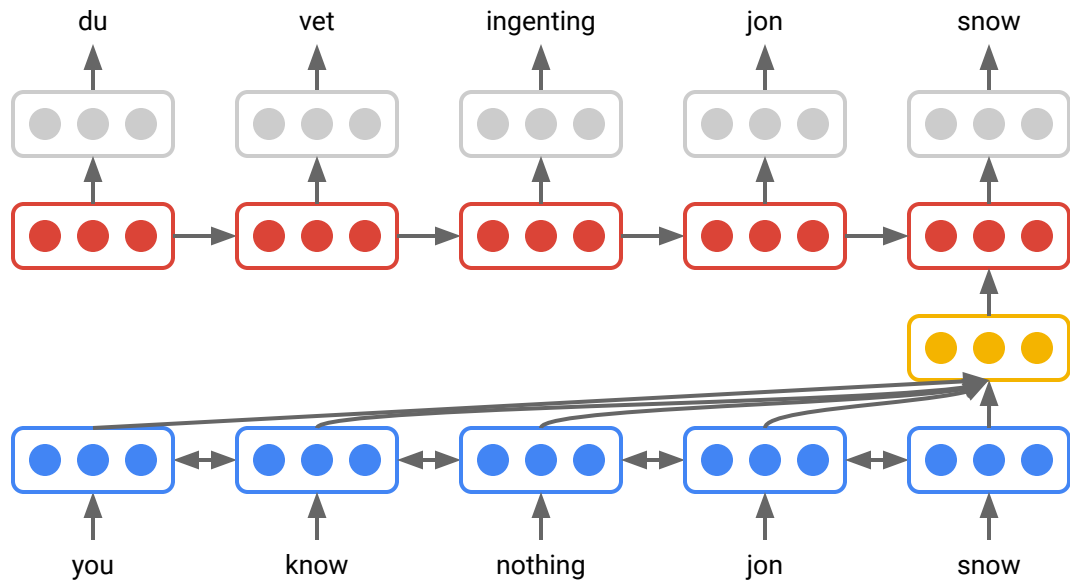
Encoder-decoder with Attention

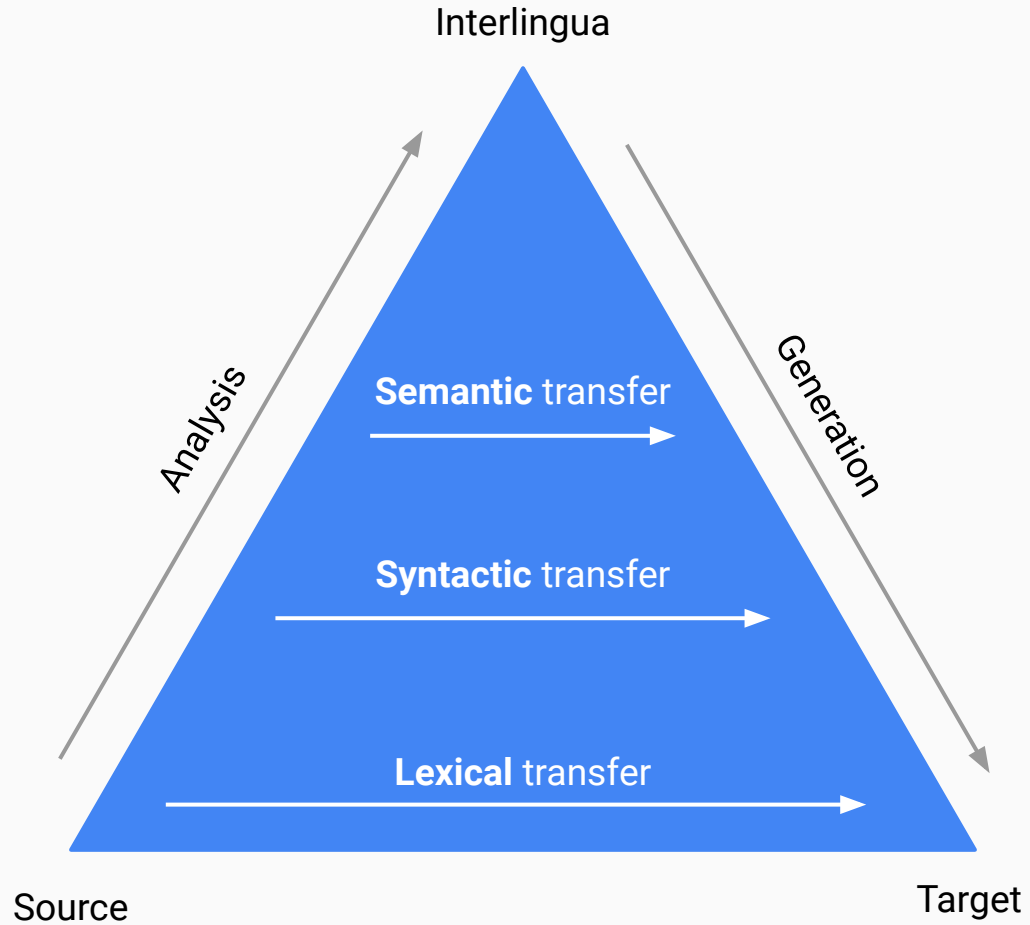


Encoder-decoder with Attention

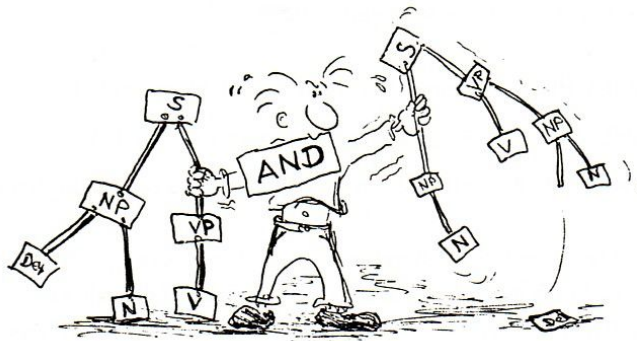


Encoder-decoder with Attention





Motivation

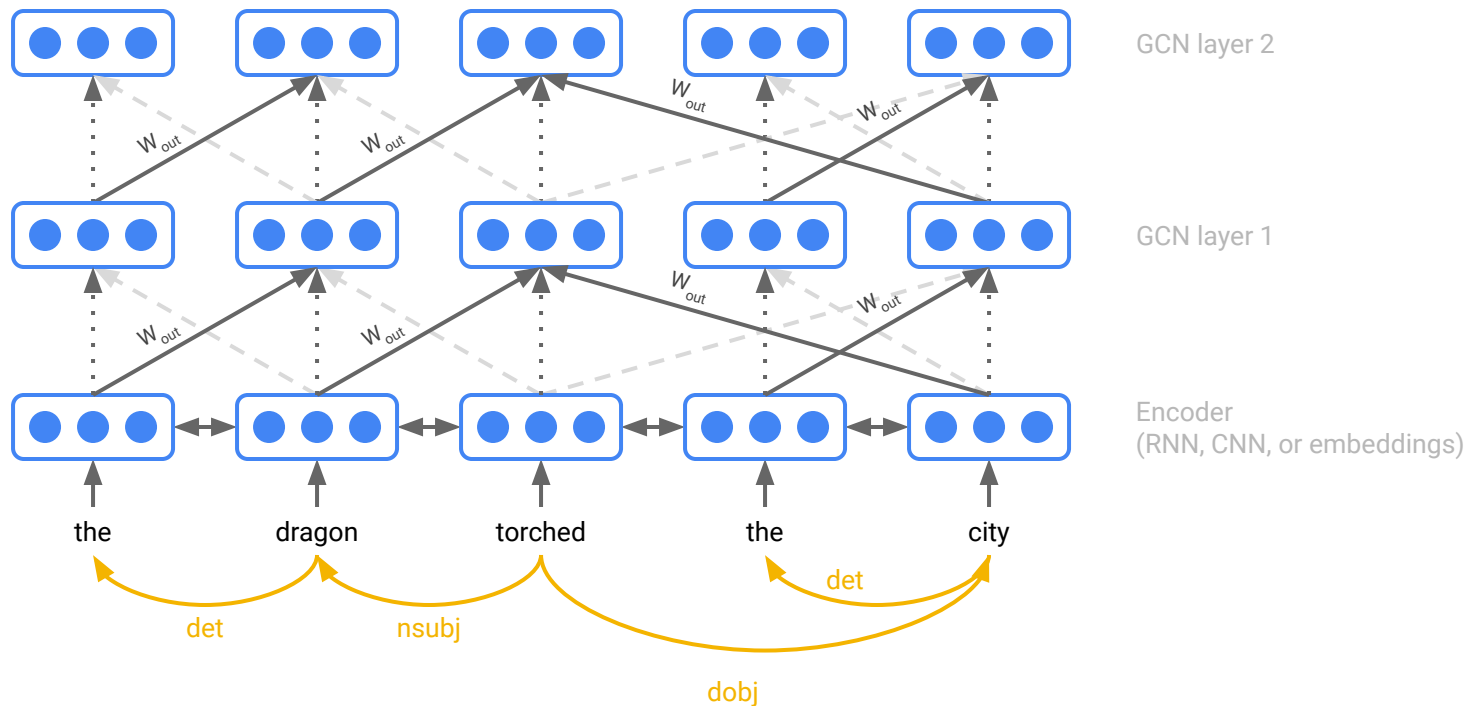


State-of-the-art NMT systems **lack** any explicit modeling of **syntax**.

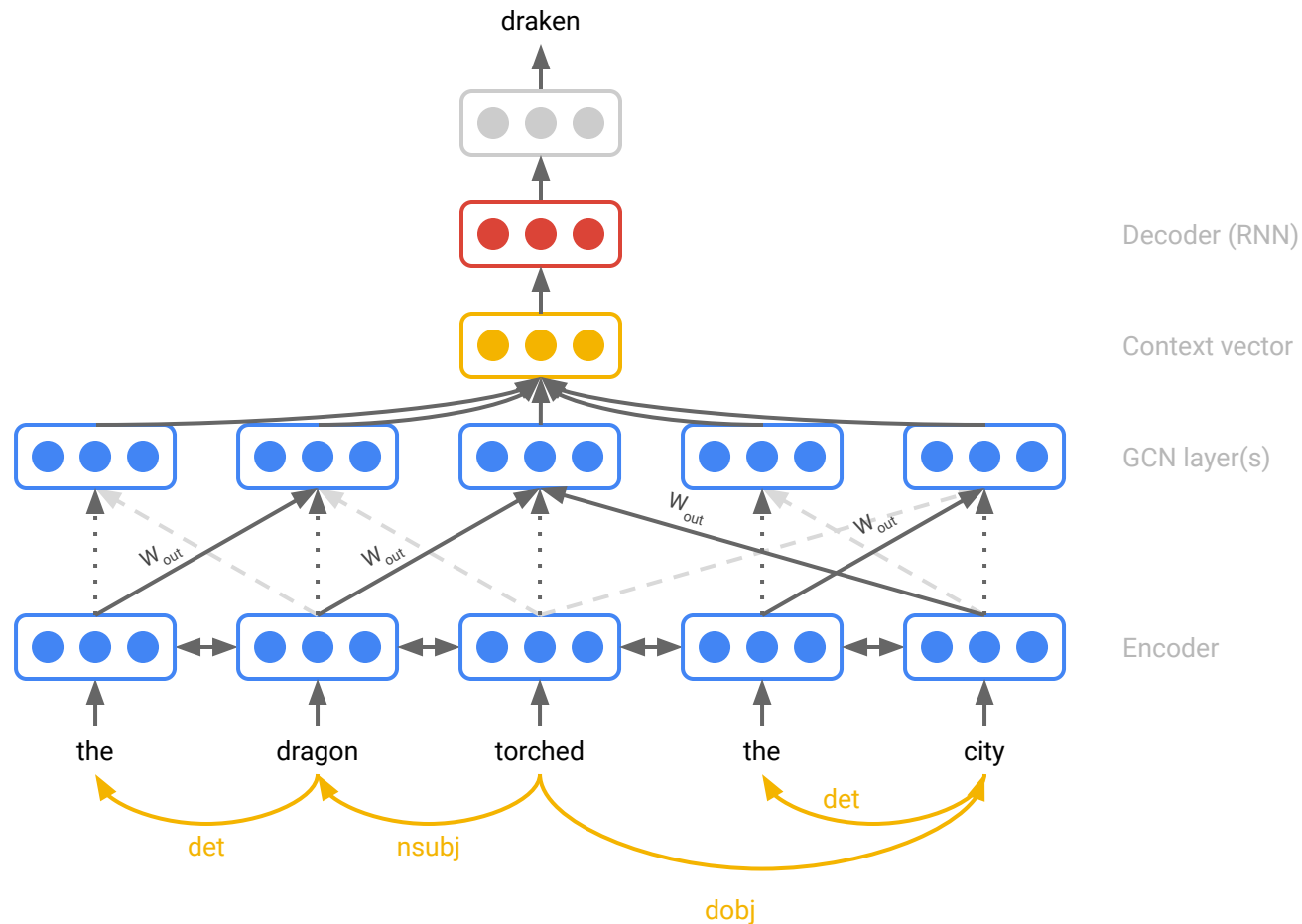
They fully rely on an **LSTM (or self-attention)** to capture syntactic phenomena.

We can use GCNs to condition on syntax and/or semantics

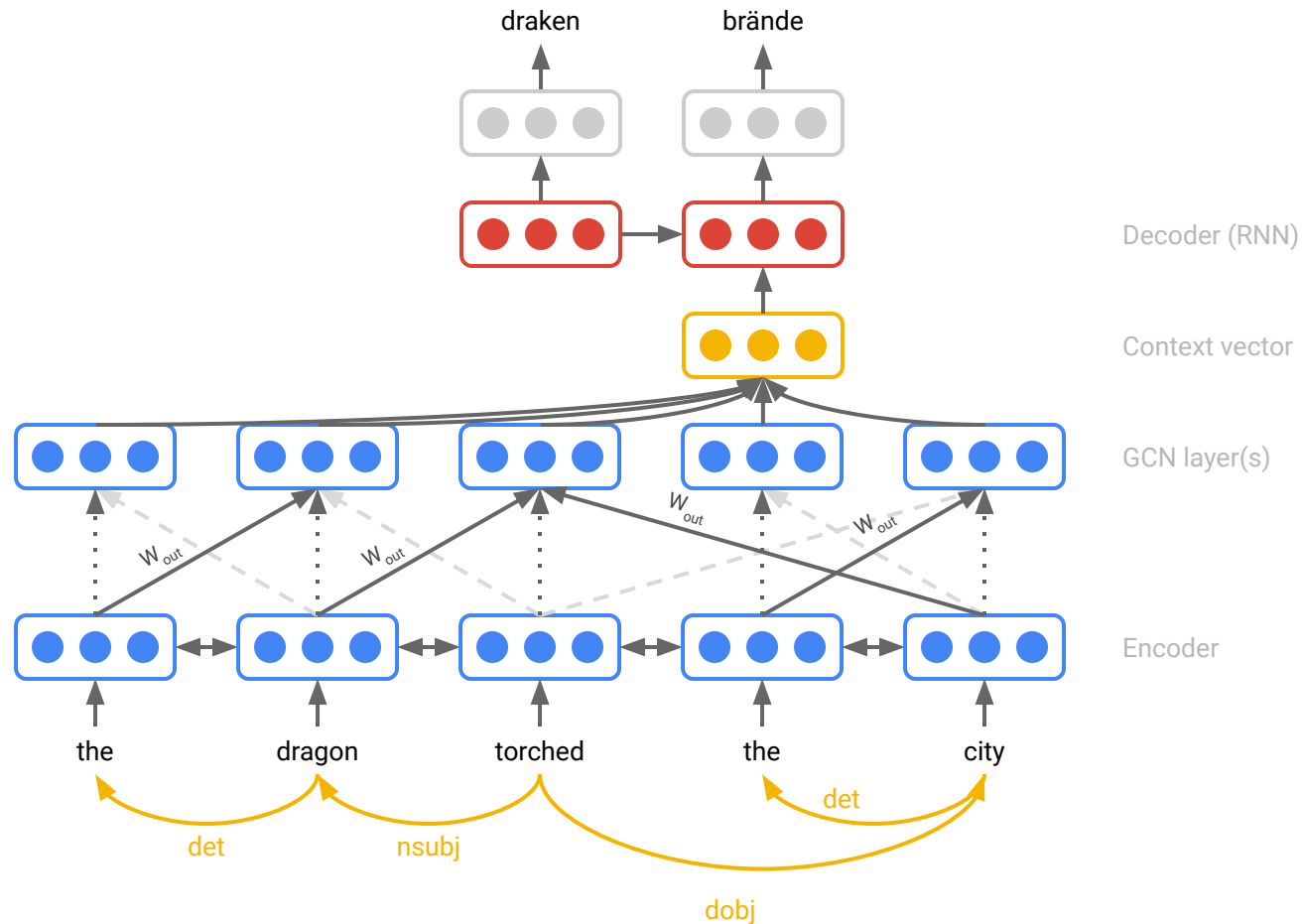
Graph Convolutional Encoders for NMT



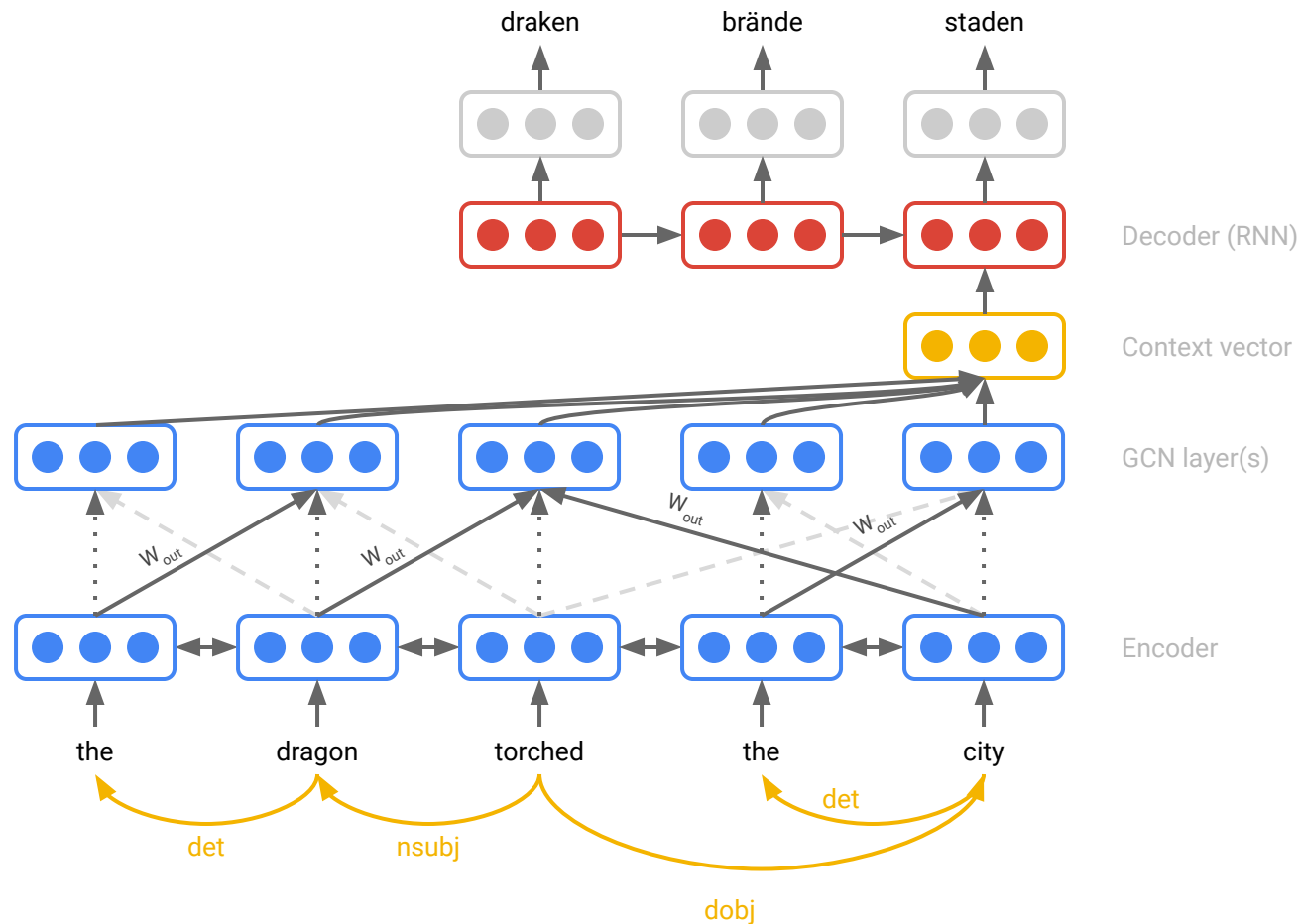
GCNs for NMT



GCNs for NMT



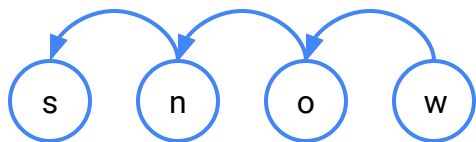
GCNs for NMT



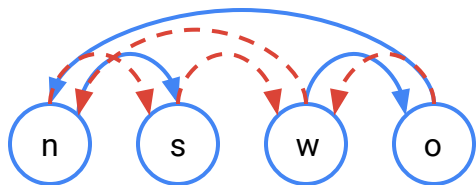
Artificial reordering experiment

Random sequences from **26 types**

Each token is linked to its **predecessor**



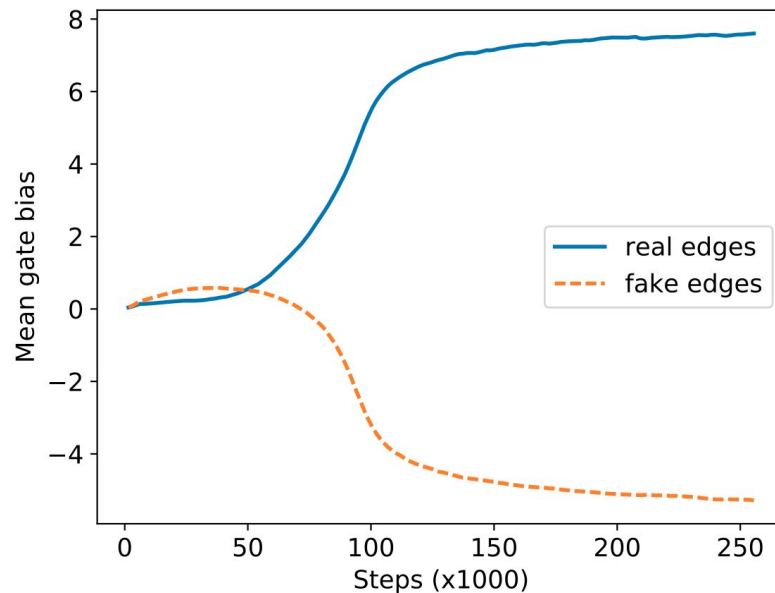
Sequences are randomly **permuted**



Each token is also linked to a **random** token with a “fake edge” using a different label set

A BiRNN+GCN model is able to learn how to put the permuted sequences **back into order**.

Real edges are distinguished from fake edges.



Machine Translation experiments

		Train	Validation <i>newstest2015</i>	Test <i>newstest2016</i>
English-German	NCv11	227 k	2169	2999
	WMT'16	4.5 M		
English-Czech	NCv11	181 k	2656	2999

We parse English source sentences using Google's **SyntaxNet**

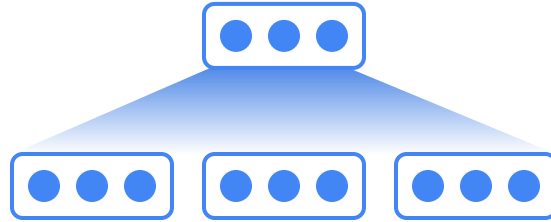
BPE on **target side** (8k merges, 16k for WMT'16)

Embeddings: 256 units, GRUs/CNNs: 512 units (800 for WMT'16)

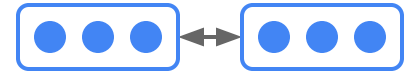
Three baselines



Bag of words

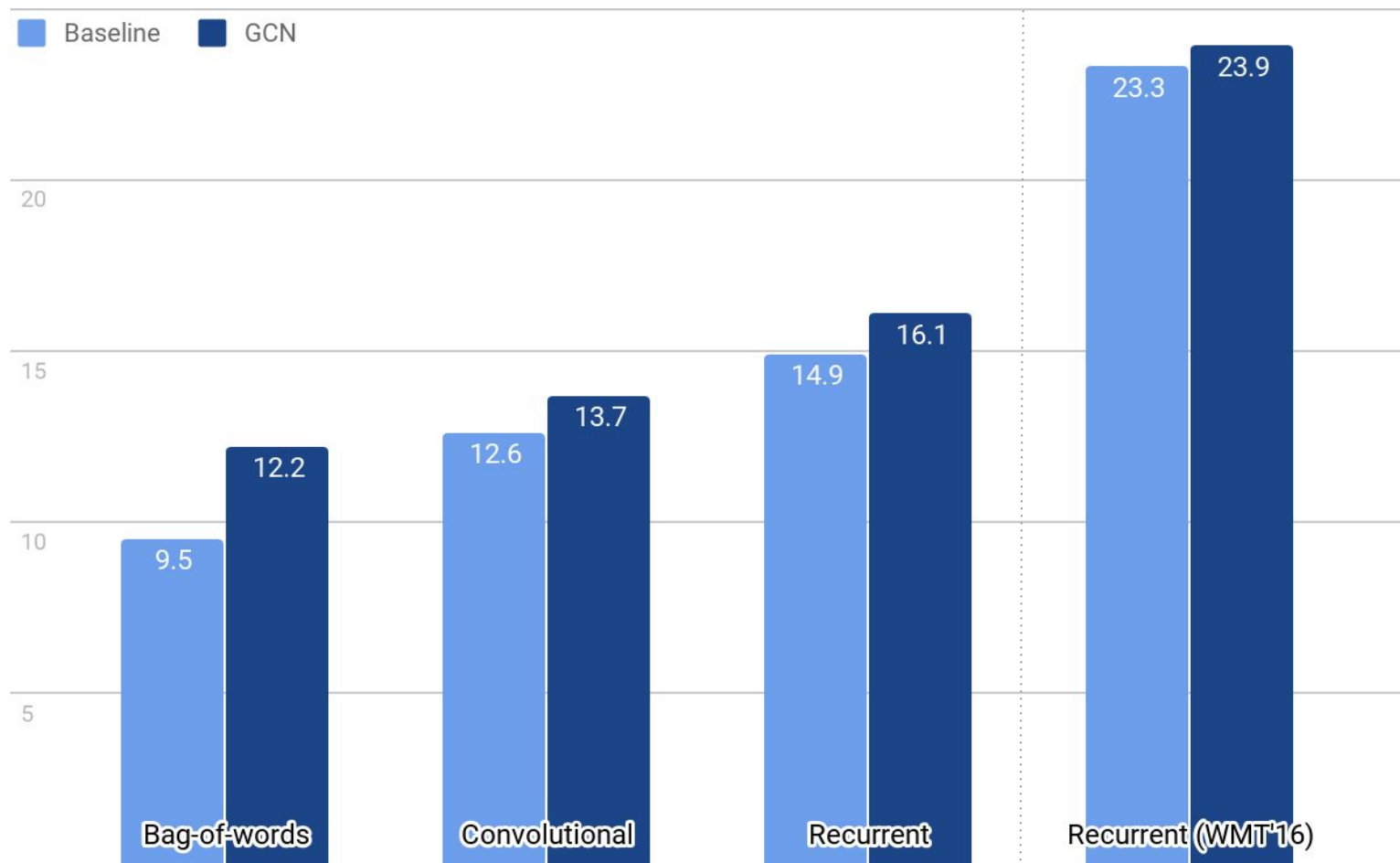


Convolutional

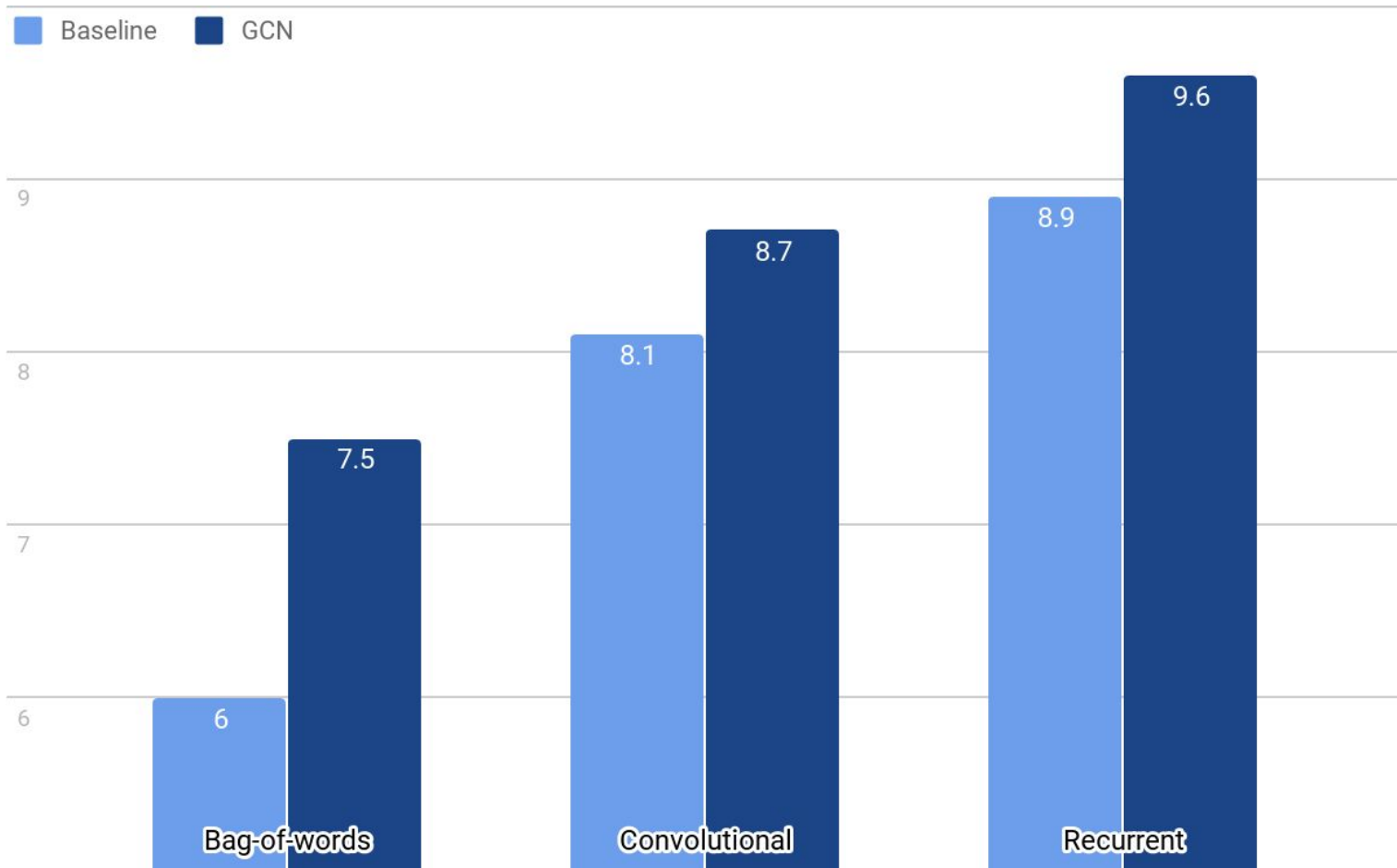


Recurrent

Results English-German (BLEU)



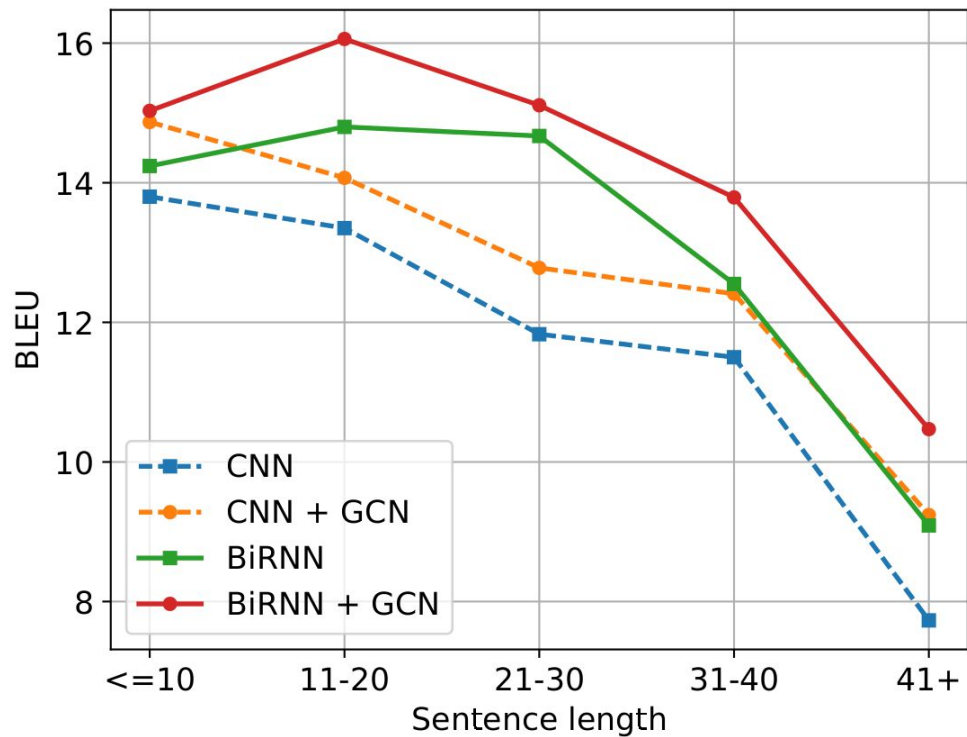
Results English-Czech (BLEU)



Effect of GCN layers

	English-German		English-Czech	
	BLEU ₁	BLEU ₄	BLEU ₁	BLEU ₄
BiRNN	44.2	14.1	37.8	8.9
+ GCN 1L	45.0	14.1	38.3	9.6
+ GCN 2L	46.3	14.8	39.6	9.9

Effect of sentence length



Flexibility of GCNs

In principle we can condition on any kind of graph-based linguistic structure:

Semantic Role Labels

Co-reference chains

AMR semantic graphs

NMT with Syntax *and* Semantics

We can condition on *both* syntax and semantics:

- Dependency graph
- Semantic role structures

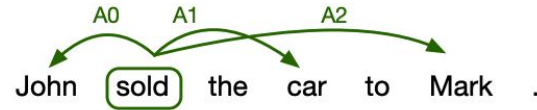
	BiRNN	CNN
Baseline (Bastings et al., 2017)	14.9	12.6
+Sem	15.6	13.4
+Syn (Bastings et al., 2017)	16.1	13.7
+Syn + Sem	15.8	14.3

Table 1: Test BLEU, En–De, News Commentary.

	BiRNN
Baseline (Bastings et al., 2017)	23.3
+Sem	24.5
+Syn (Bastings et al., 2017)	23.9
+Syn + Sem	24.9

Table 2: Test BLEU, En–De, full WMT16.

Semantics can help to get the right translation



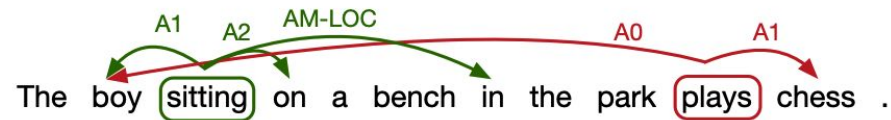
BiRNN John verkaufte das Auto nach Mark .

Sem John verkaufte das Auto an Mark .



BiRNN Der Junge zu Fuß die staubige Straße ist ein Bier trinken .

Sem Der Junge , der die staubige Straße hinunter geht , trinkt ein Bier .



BiRNN Der Junge auf einer Bank im Park spielt Schach .

Sem Der Junge sitzt auf einer Bank im Park Schach .

Inducing graphs

Inducing graphs

Instead of conditioning on a graph, we can try to **induce** one

The idea is the following:

1. **Predict** the graph for a sentence
2. Process that sentence **using** the predicted graph, instead of conditioning on an externally provided graph

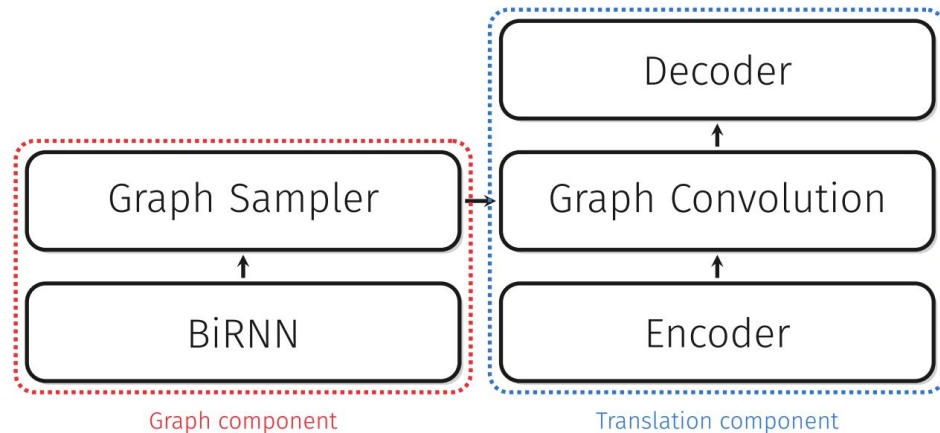
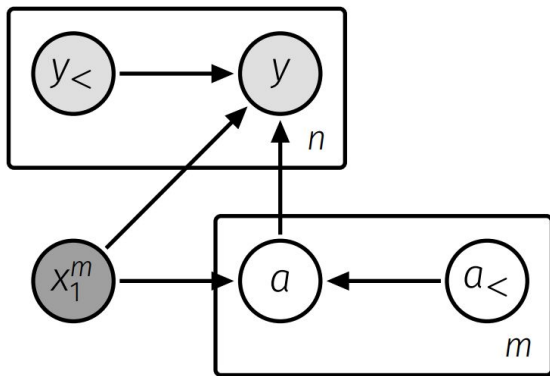
There is lots of related work here on inducing **trees**, usually with natural language inference (NLI) as a downstream task.

See e.g. Williams et al. "Do latent tree learning models identify meaningful structure in sentences?." TACL, 2018.

Inducing graphs with Neural Machine Translation

We define a deep generative model with:

1. a **graph component**
samples graph a conditioned on x
2. a **translation component**
given x and a , predicts translation y



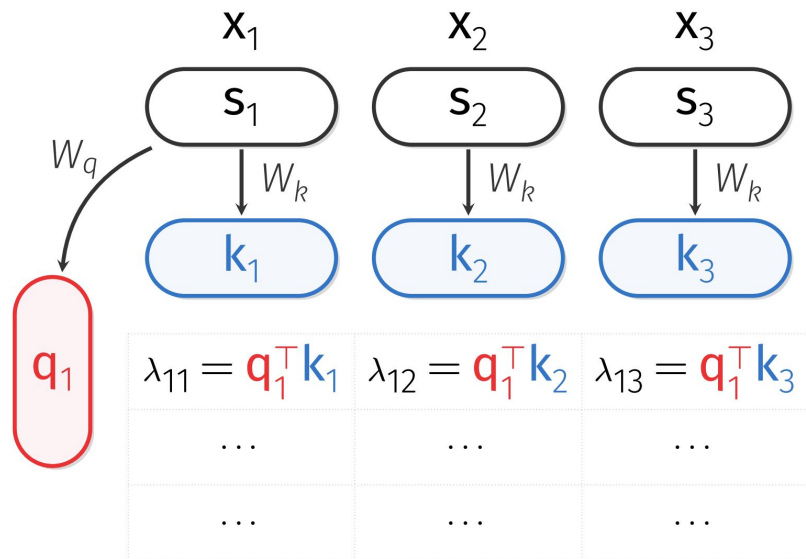
Graph component

Samples for each source position i an m -dimensional probability vector A_i :

$$A_i | x_{1..m} \sim \text{Concrete}(\tau, \lambda_i)$$

We interpret A_i as a **head distribution**

- a_{ik} is the relative strength of the edge $x_i \rightarrow x_k$
- 'Head potentials' λ_i are computed with self-attention

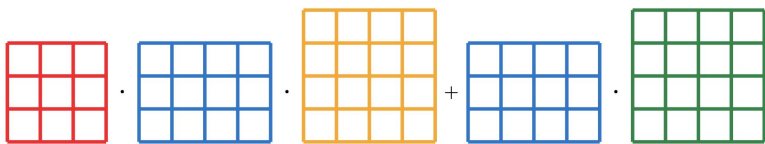


Self-attention

Translation component

Attentive encoder-decoder that incorporates graph $A = a_{m..1}$ using graph convolution

$S = \text{encode}(x_1, \dots, x_m)$



$$\text{GCN}(S, A) = \text{ReLU}(ASW_{\text{IN}} + SW_{\text{LOOP}} + \mathbf{b})$$

We then sample target words one step a a time:

$$Y_j | X_1^m, a_1^m, y_{<j} \sim \text{Cat}(\pi_j)$$

$$\pi_j = f_{\theta}(x_1^m, a_1^m, y_{<j})$$

Objective

We optimize the following objective:

$$\begin{aligned}\log p(y | x) &\geq \mathbb{E} \left[\log p(y | x, a) \right] \\ &\approx \log p(y | x, a) \\ &a \sim p(a | x_1^m)\end{aligned}$$

Experiments

German-English (IWSLT14) and Japanese-English (ASPEC)

	Train	Dev	Test	Vocabulary
De-En	153K	7282	6750	32010/22823
Ja-En	2M	1790	1812	16384 (SPM)

Results

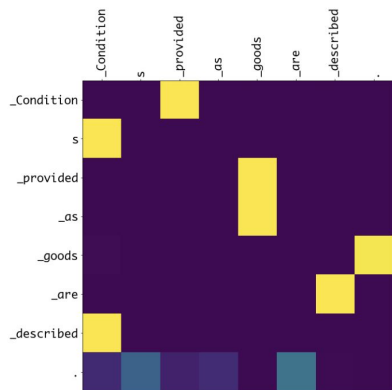
		IWSLT14		WAT17	
ENCODER		DE-EN	EN-DE	JA-EN	EN-JA
Ext. baseline	RNN	27.6	-	-	28.5
Baseline	Emb.	22.7	17.9	18.1	18.1
Baseline	CNN	23.6	19.1	23.0	24.6
Baseline	RNN	27.6	22.4	26.0	28.7
Latent Graph	Emb.	24.0	18.7	23.2	24.3
Latent Graph	CNN	24.6	20.3	24.6	26.7
Latent Graph	RNN	27.2	22.4	26.0	29.1

ENCODER	MEAN HEAD DISTANCE		MEAN ENTROPY	
	JA-EN	EN-JA	JA-EN	EN-JA
Emb.	4.0 \pm 6.9	3.8 \pm 5.6	0.49 \pm 0.18	0.42 \pm 0.18
CNN	6.1 \pm 6.5	6.7 \pm 7.1	1.21 \pm 0.28	1.47 \pm 0.30
RNN	4.3 \pm 6.5	2.0 \pm 5.4	0.51 \pm 0.20	0.00 \pm 0.01

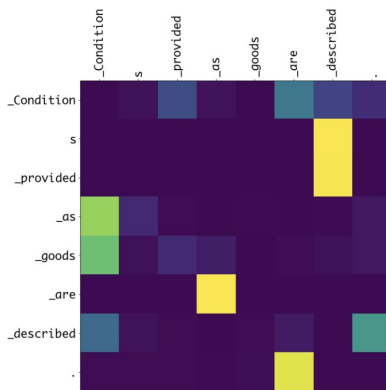
The mean **head distance** shows that the graphs may capture **non-local** dependencies

The mean **entropy** shows that the graphs have a high degree of **sparsity**

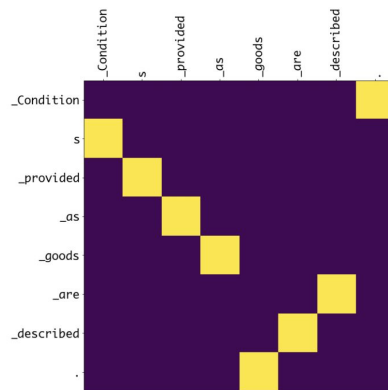
Examples



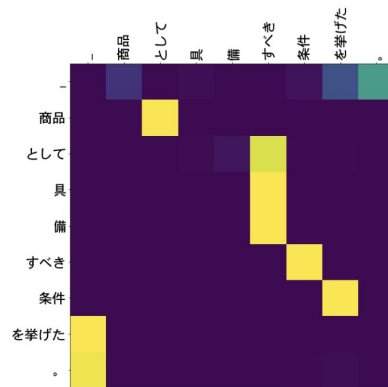
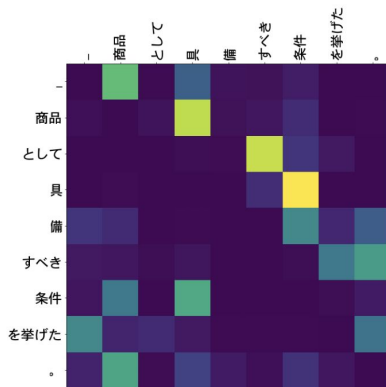
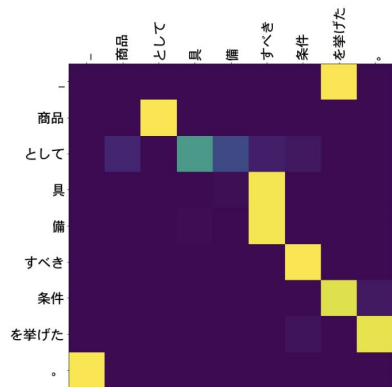
(a) En-Ja Emb



(b) En-Ja CNN



(c) En-Ja RNN



Further readings

Not a deep generative model, but with restrictions on the induced graph:

Tran & Bisk. Inducing Grammars with and for Neural Machine Translation.

ACL NMT workshop 2018

Transformer

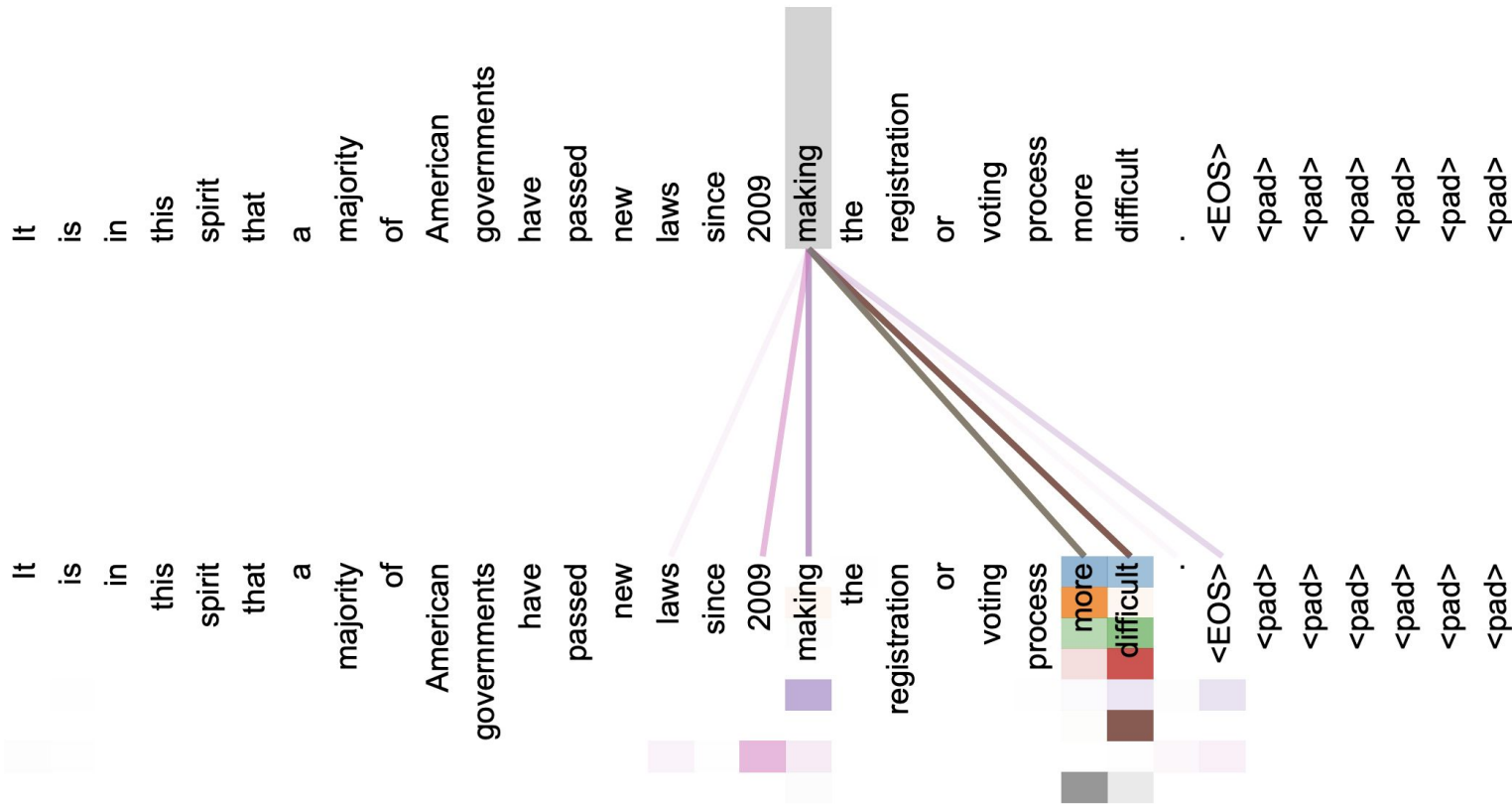
Transformer (Vaswani et al., 2017)



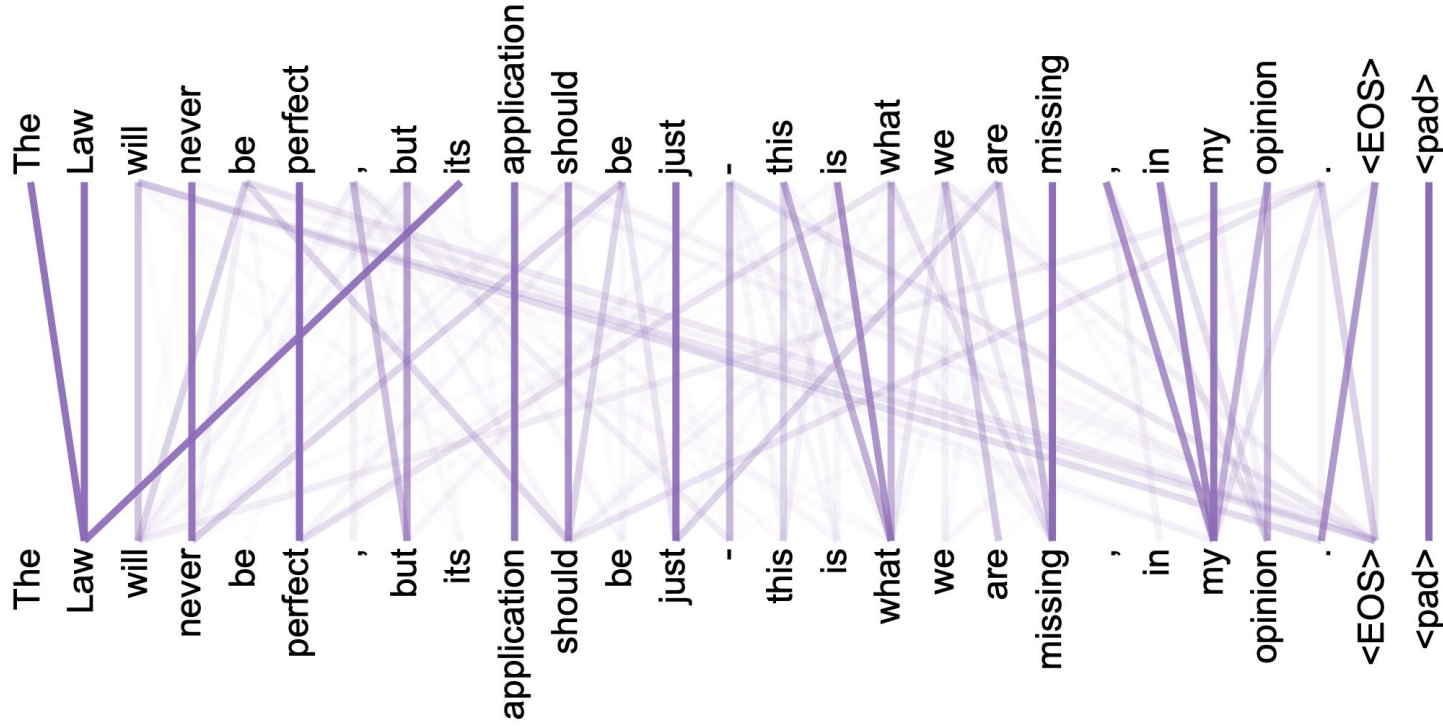


imgflip.com

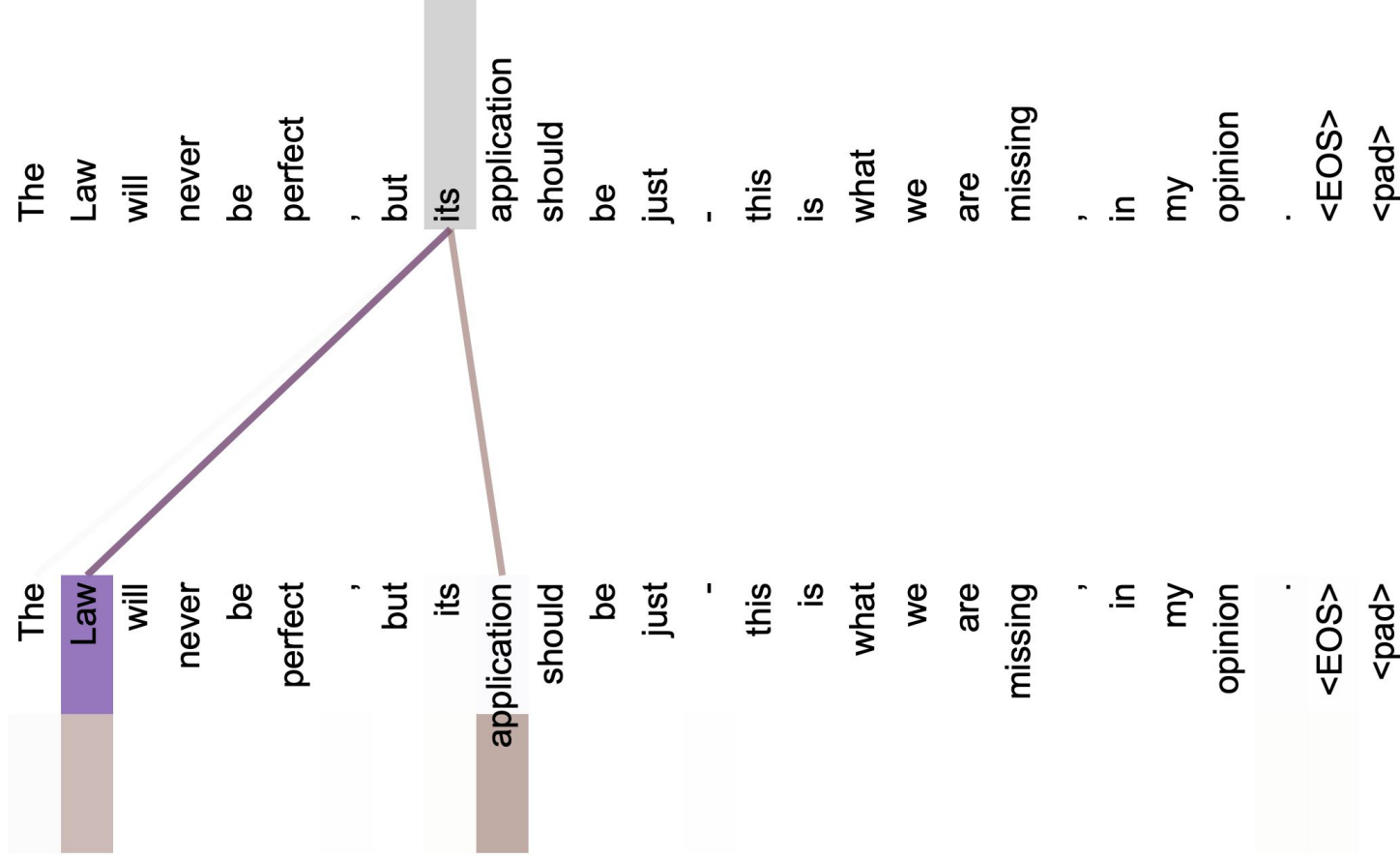
Attention visualization



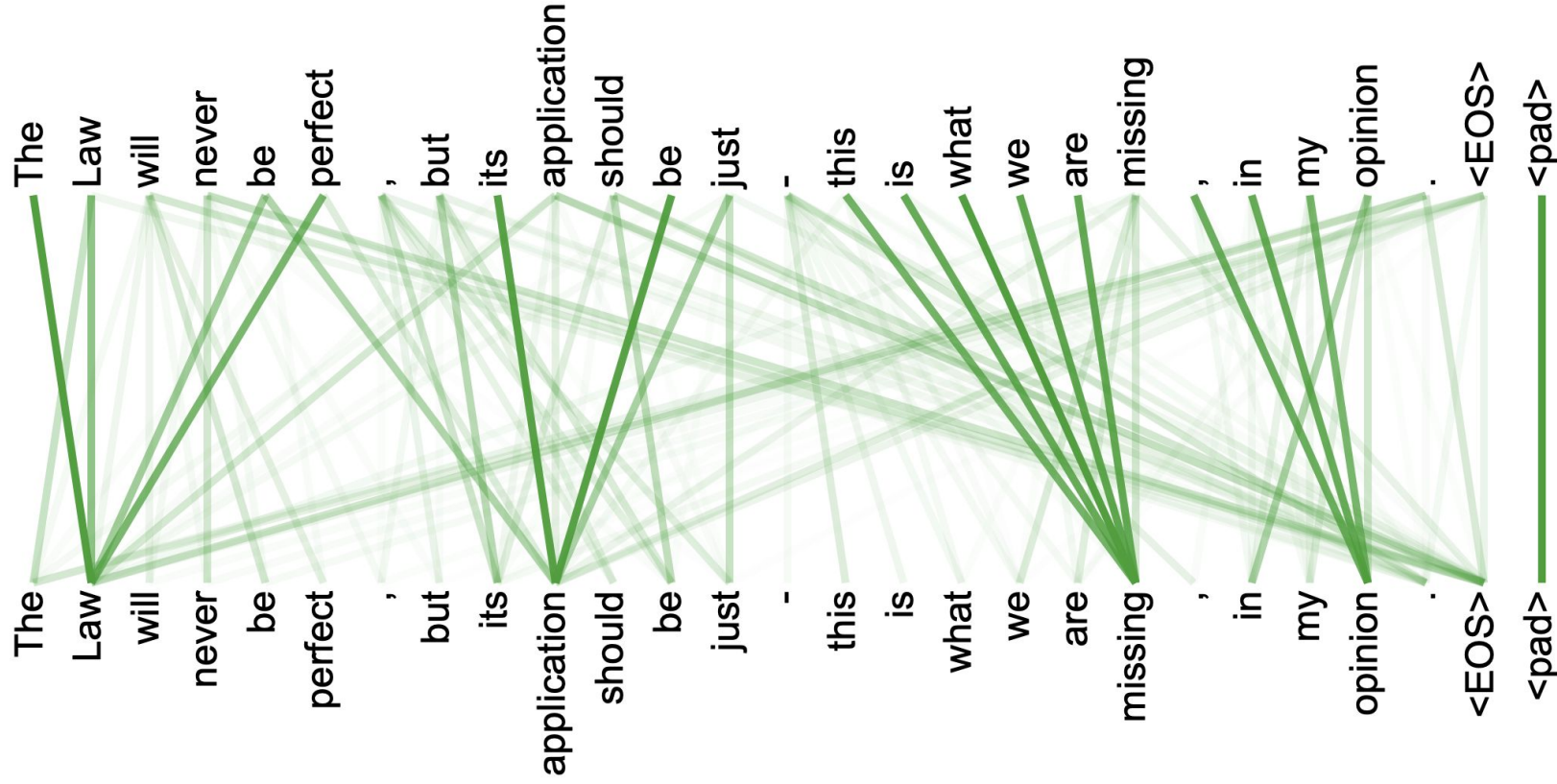
Attention visualization



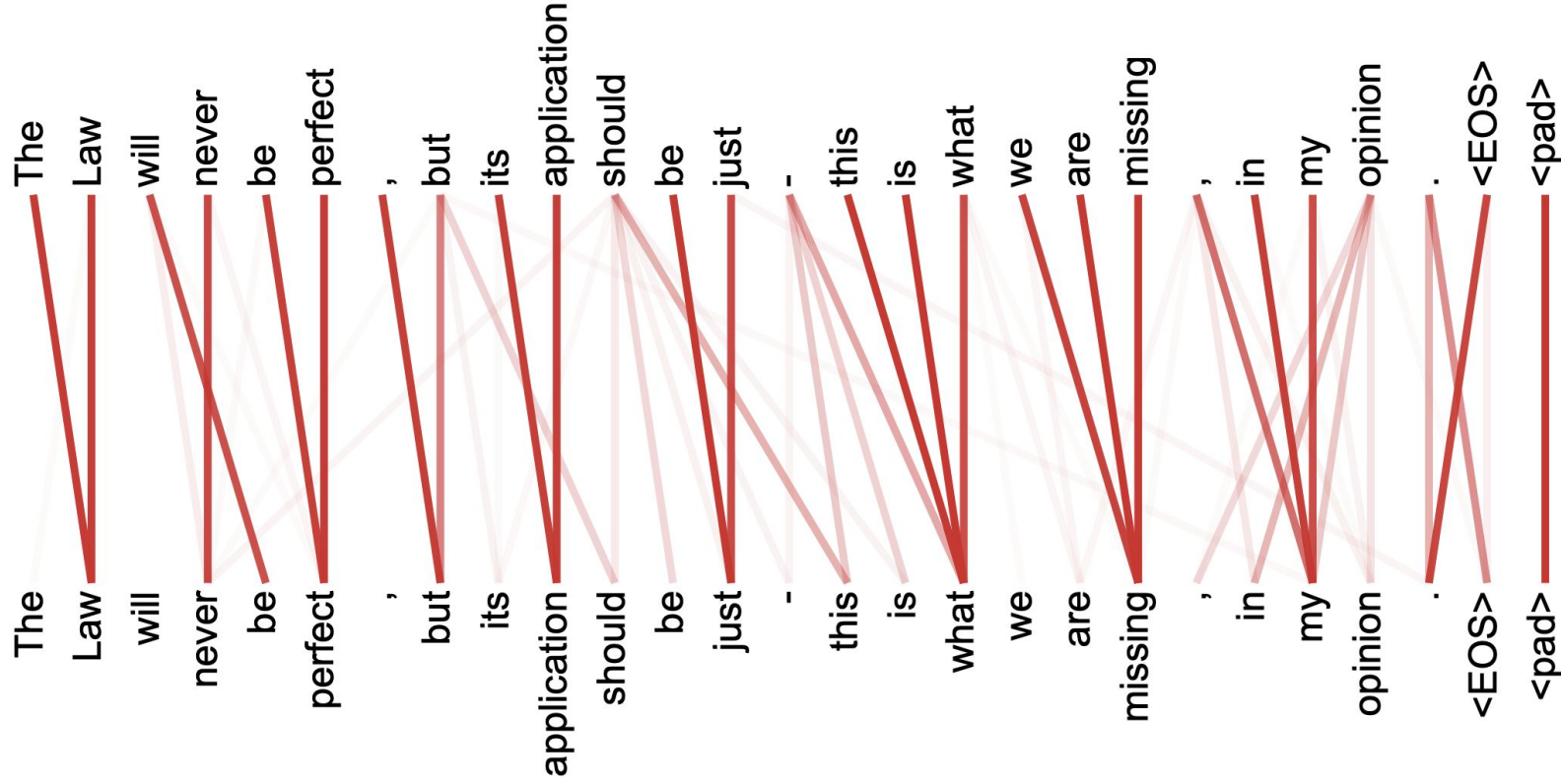
Attention visualization



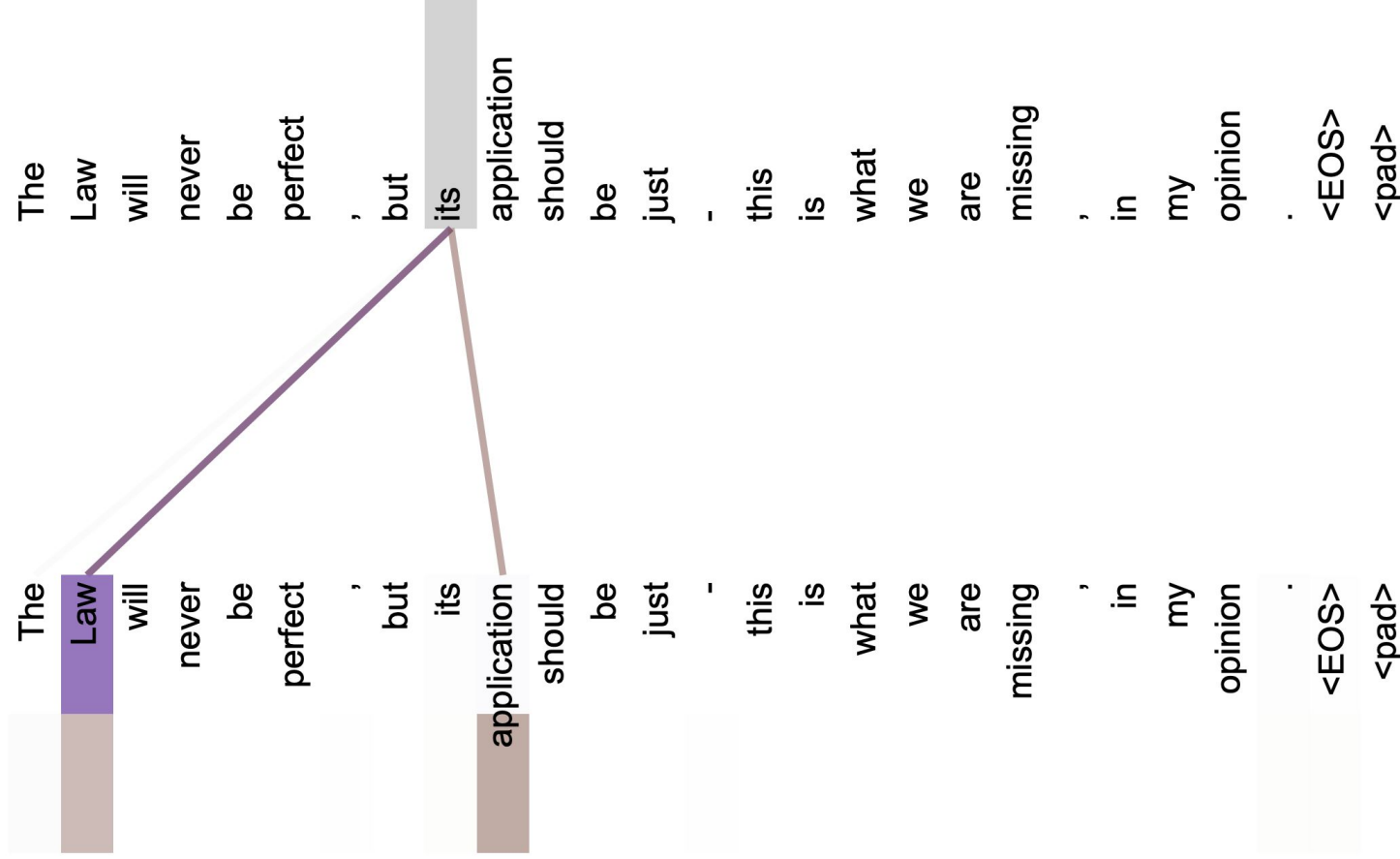
Attention visualization



Attention visualization



Attention visualization



Transformer (Vaswani et al., 2017)

Start with input sequence w_1, w_2, \dots, w_n

Map to sequence of vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$
each $\mathbf{x}_i \in \mathbb{R}^d$

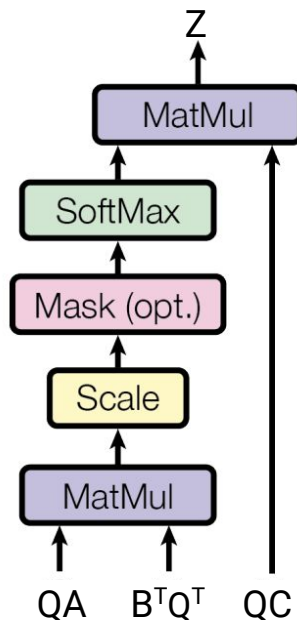
Stack the vectors to make matrix $Q \in \mathbb{R}^{n \times d}$

Define the parameters of the transformation:

$$A \in \mathbb{R}^{d \times l} \quad B \in \mathbb{R}^{d \times l} \quad C \in \mathbb{R}^{d \times o}$$

Then we can define:

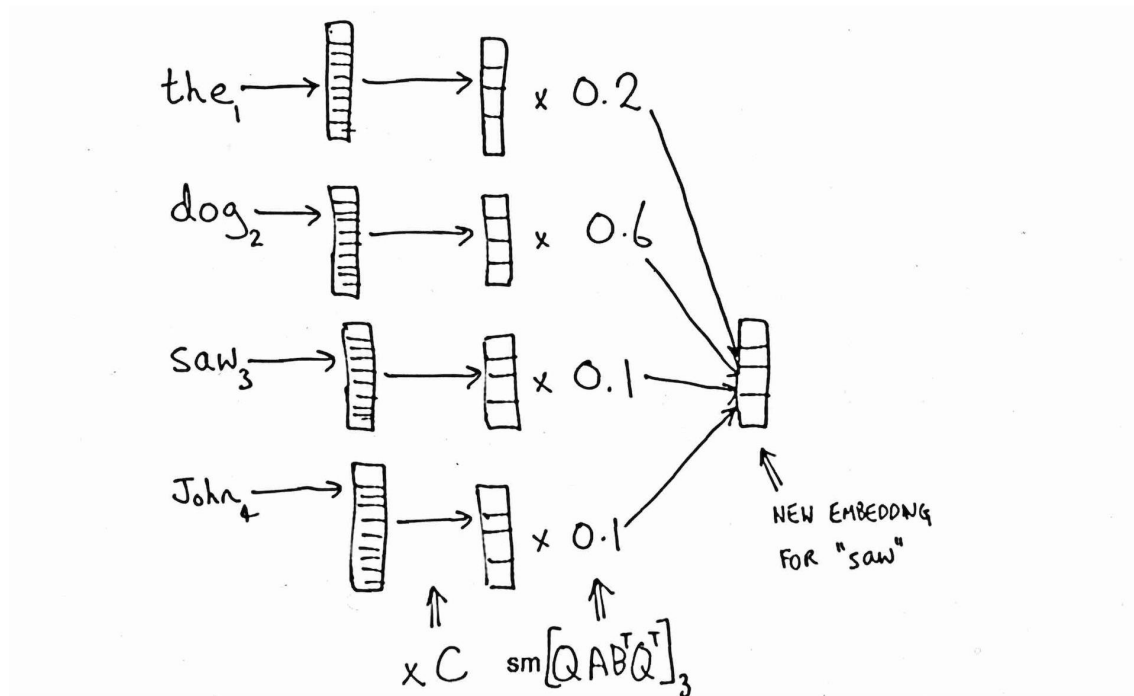
$$Z = \text{softmax}(QAB^TQ^T) QC$$



Intuition:

- QC is a matrix of new word embeddings
- QAB^TQ^T is an $n \times n$ matrix of inner products in a new l -dimensional space
- $\text{softmax}(\cdot)$ makes the matrix positive and the rows sum to 1 (*self-attention*)

Encoder example



Multi-Headed Self-Attention

Say $d = 512$, $o = 64$, and $h = 8$

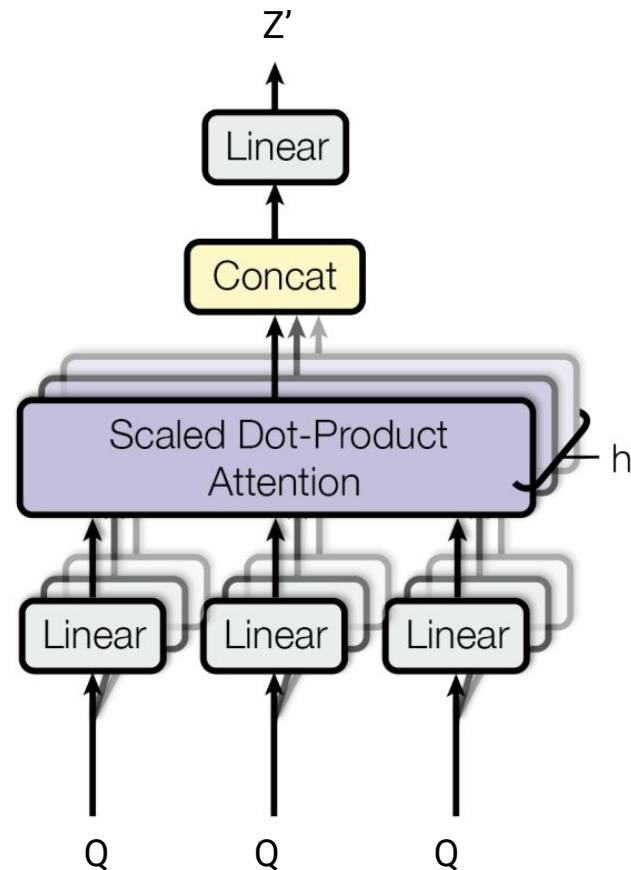
Define parameters $A^j, B^j \in \mathbb{R}^{d \times l}$ $C^j \in \mathbb{R}^{d \times o}$ for $j = 1 \dots h$

For $j = 1 \dots h$:

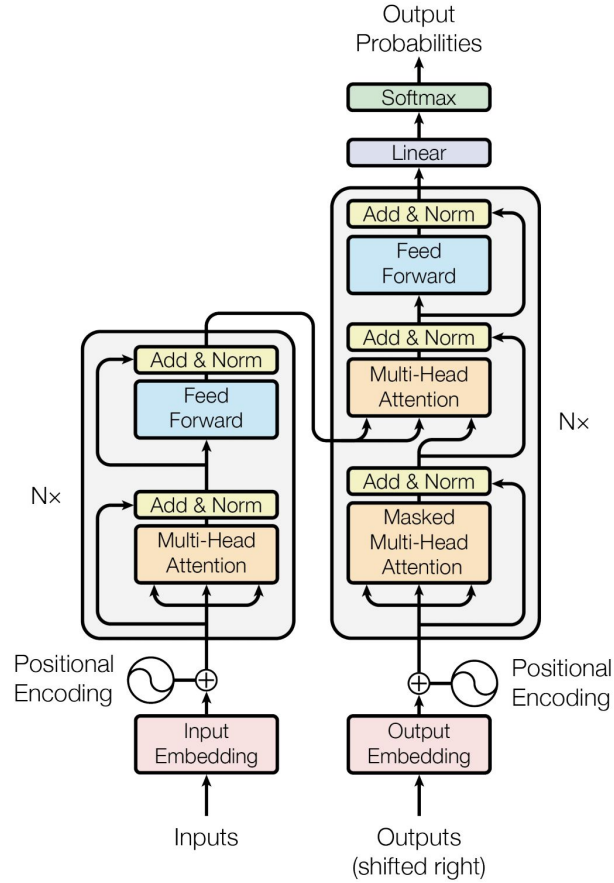
$$Z^j \in \mathbb{R}^{n \times 64} = \text{softmax}(QA^jB^jTQ^T) QC^j$$

$$Z \in \mathbb{R}^{n \times 512} = \text{concat}(Z^1, Z^2, \dots, Z^h)$$

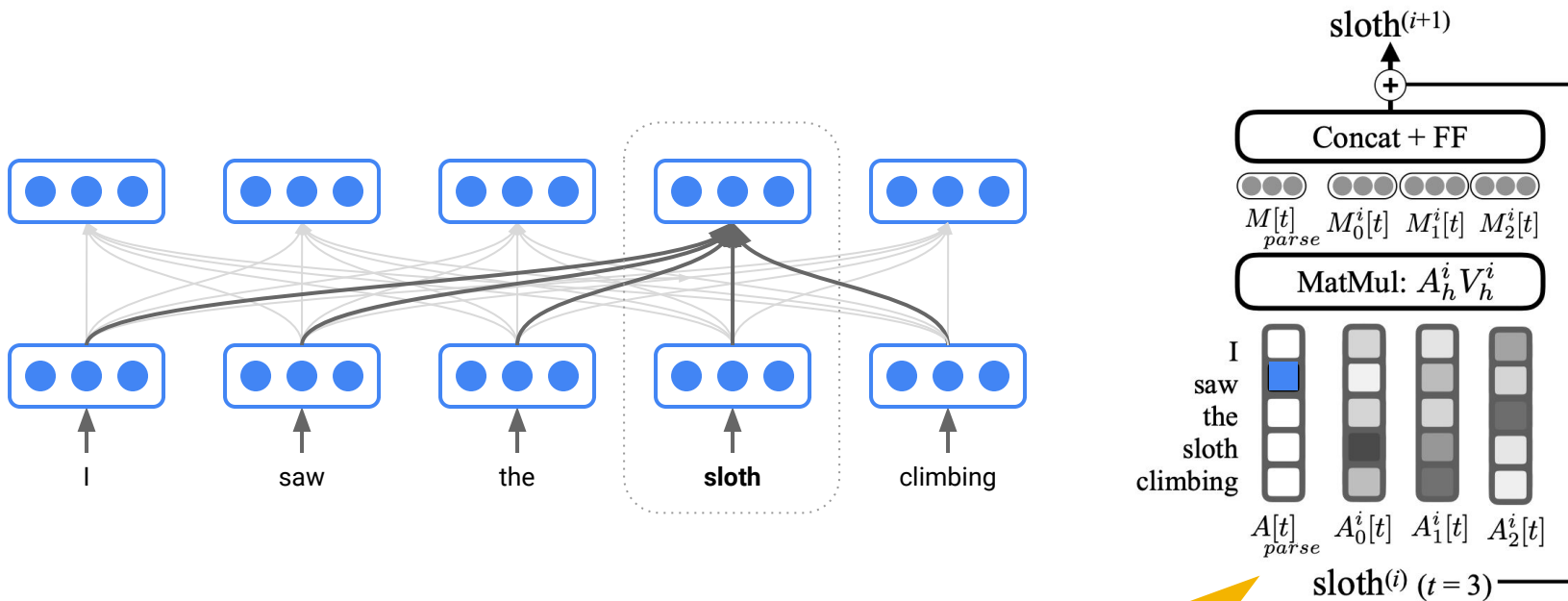
$$Z' = \text{feed-forward}(\text{layer-norm}(Q + Z))$$



Final Transformer model



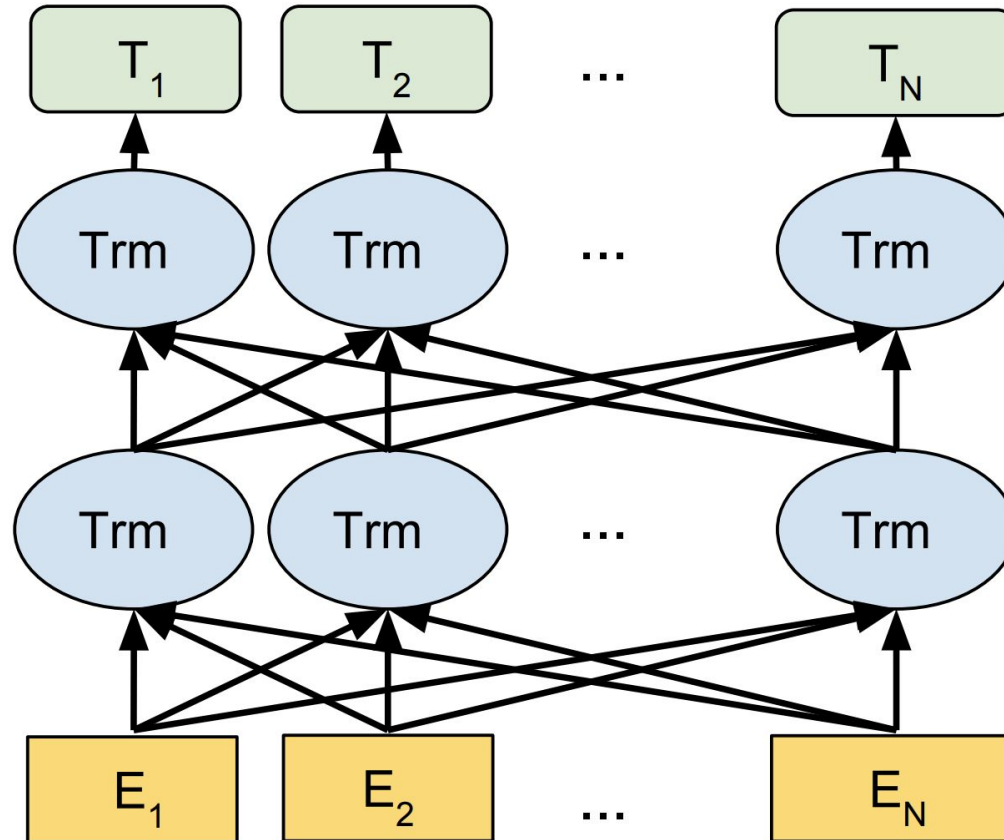
Linguistically-Informed Self-Attention



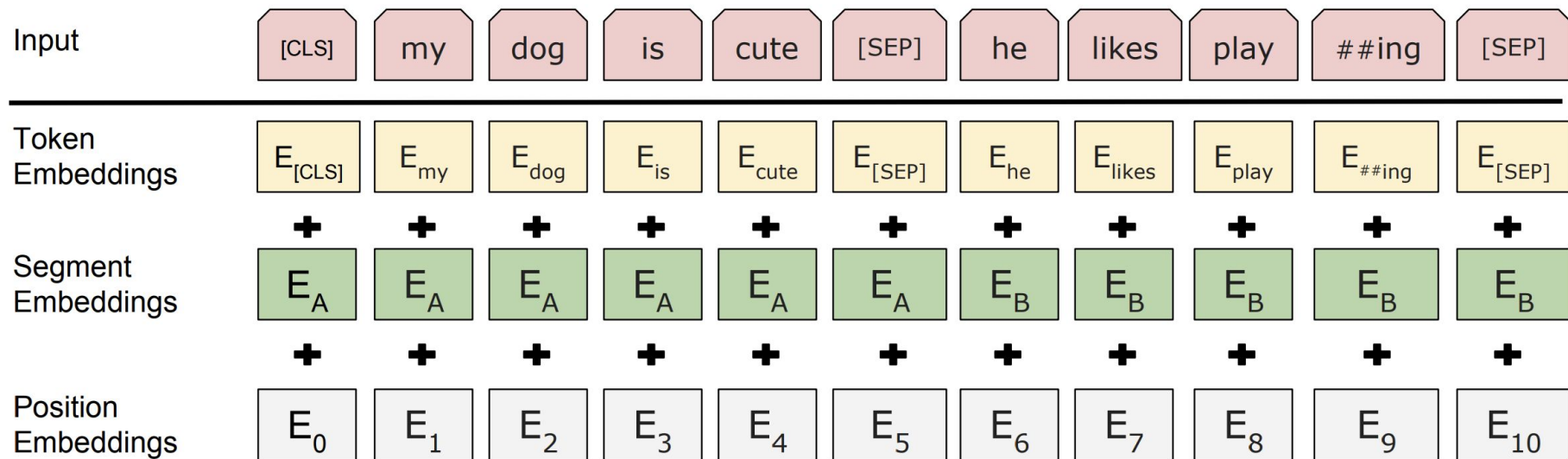
Train first attention head to predict the syntactic head

Pre-trained Transformers

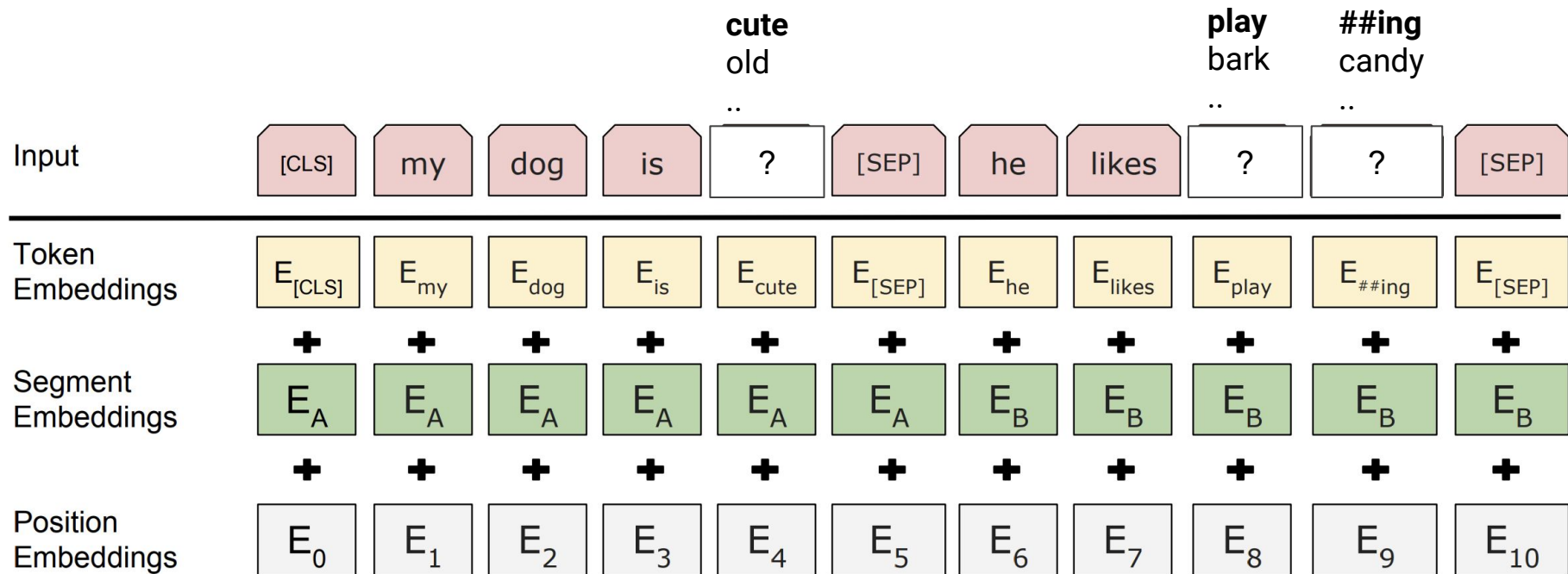
BERT: Pre-trained Transformer (Devlin et al., 2018)



BERT input



Training: with masked input, predict gaps

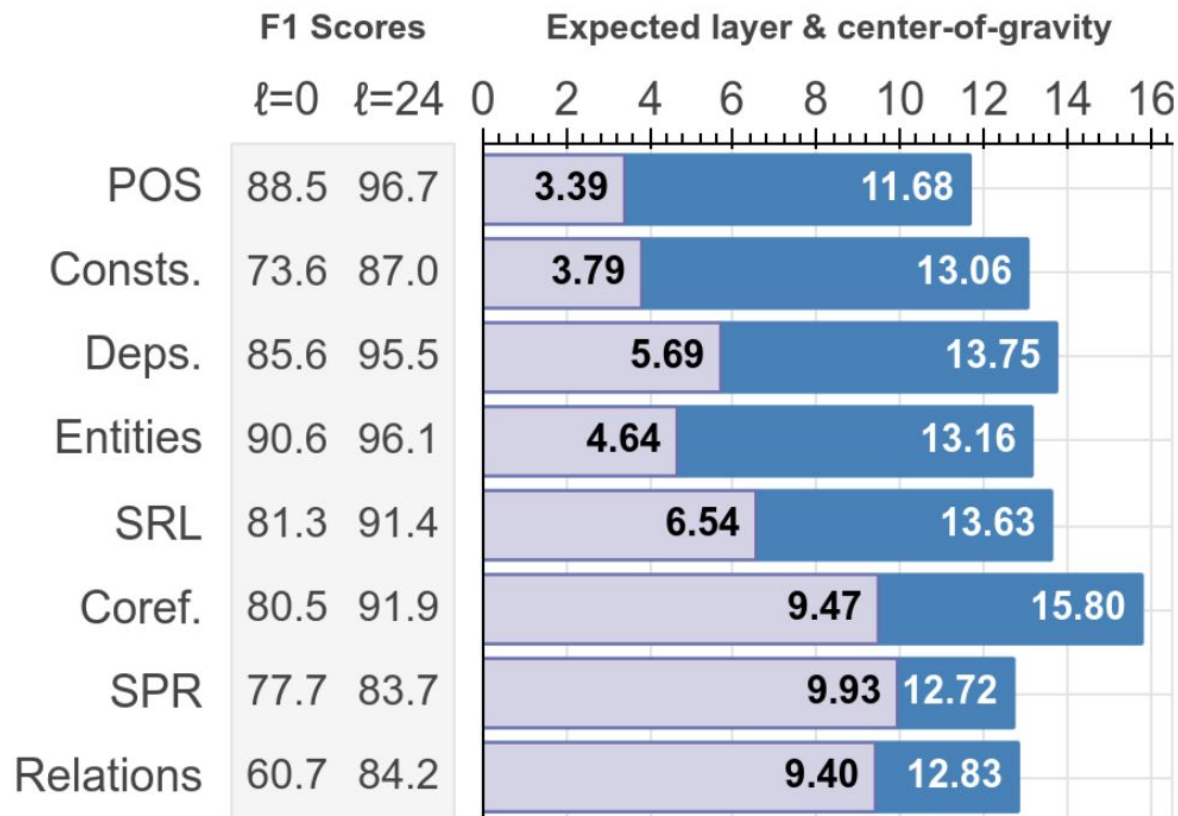


Training: next sentence prediction

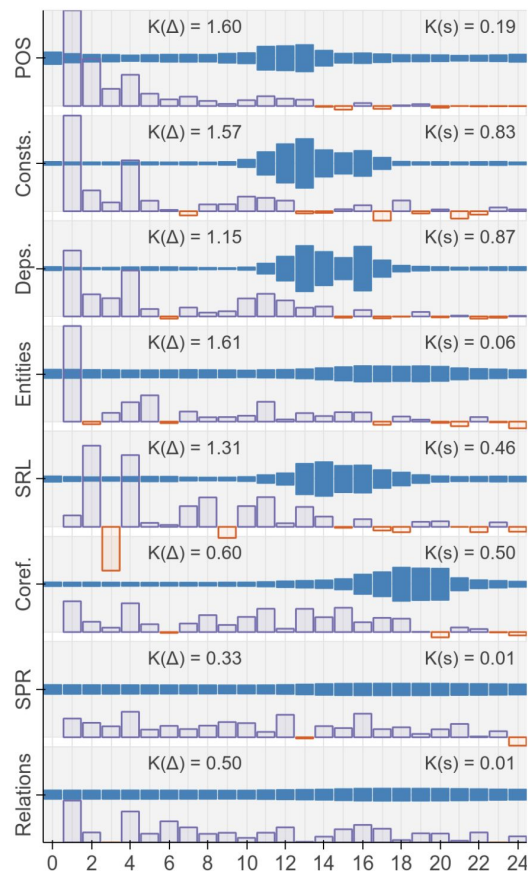
Sentence A = The man went to the store.
Sentence B = He bought a gallon of milk.
Label = IsNextSentence

Sentence A = The man went to the store.
Sentence B = Penguins are flightless.
Label = NotNextSentence

What can pre-training give us? (Tenney et al., 2019)

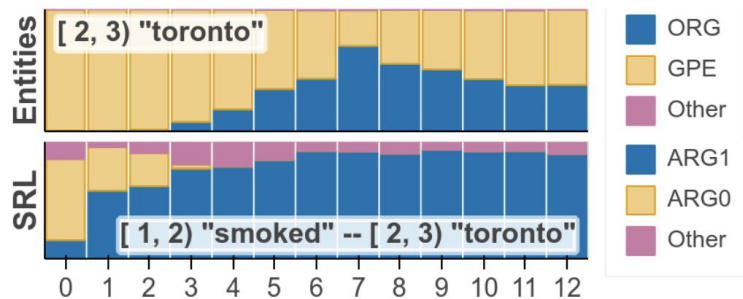


What can pre-training give us? (Tenney et al., 2019)

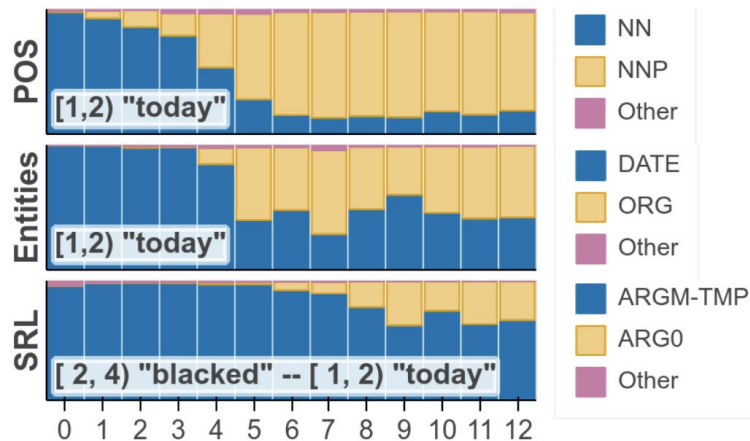


What can pre-training give us? (Tenney et al., 2019)

(a) he smoked **toronto** in the playoffs with six hits, ...



(b) china **today** blacked out a cnn interview that was ...



Summary

Many structures in NLP take the form of graphs

We can condition on them using graph convolutional networks, GCNs

It becomes more and more interesting to see how much linguistics is already present in models, especially now that pre-trained models are becoming very powerful (e.g. BERT)

A recent BERT analysis shows that the hierarchy of BERT layers resembles a classic NLP pipeline, where POS-tagging is done first, then dependency parsing, NER, etc.

However, semantic information seems to be spread out a lot; currently it is unclear why

Thank you!

GCN resources:

- Syntactic GCN (simple PyTorch version)::
<https://tinyurl.com/syngcn>
- Undirectional GCN from Thomas Kipf with blogpost:
<https://github.com/tkipf/pygcn>
<https://tkipf.github.io/graph-convolutional-networks/>
- The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables
<https://arxiv.org/abs/1611.00712>

NMT resources:

- Annotated Encoder-Decoder with Attention blog post (NMT tutorial in PyTorch):
https://bastings.github.io/annotated_encoder_decoder/
- Joey NMT - a simple NMT toolkit in PyTorch:
<https://github.com/joeynmt/joeynmt>

Graph Convolution Surveys:

- Graph Neural Networks: A Review of Methods and Applications <https://arxiv.org/abs/1812.08434>
- Deep Learning on Graphs: A Survey
<https://arxiv.org/abs/1812.04202>
- A Comprehensive Survey on Graph Neural Networks
<https://arxiv.org/abs/1901.00596>

Related paper in NLP:

- Beck et al. (2018) Graph-to-Sequence Learning using Gated Graph Neural Networks <https://arxiv.org/abs/1806.09835> (AMR and NMT with dependency input)