

Graph Convolutional Encoders for Syntax-aware NMT

Joost Bastings

Institute for Logic, Language and Computation (ILLC)
University of Amsterdam
<http://joost.ninja>

Collaborators



Ivan Titov



Wilker Aziz

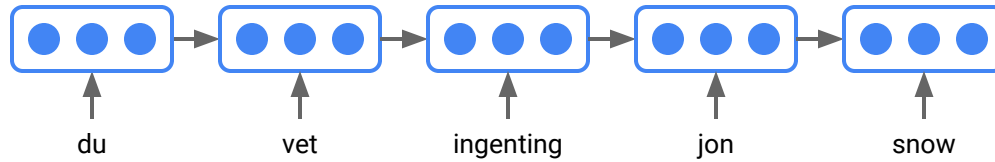


Diego Marcheggiani

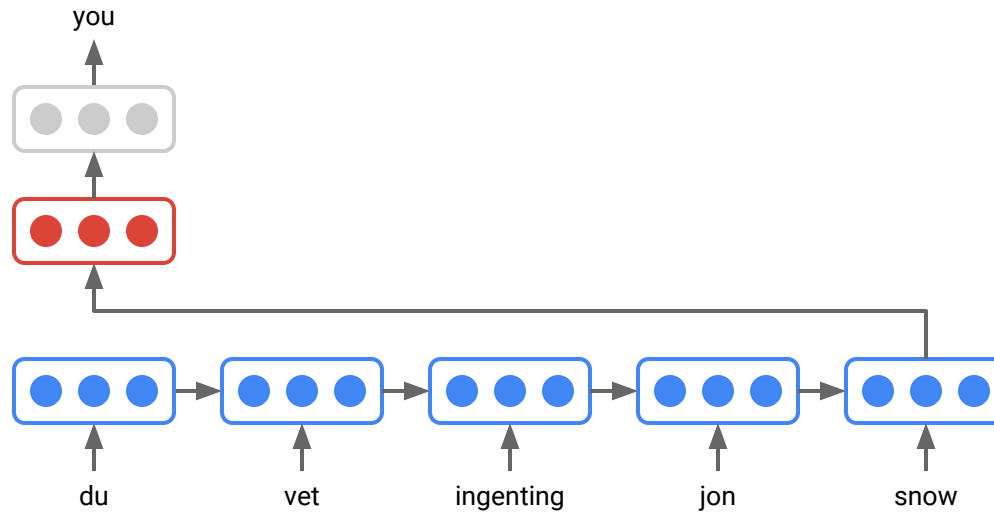


Khalil Sima'an

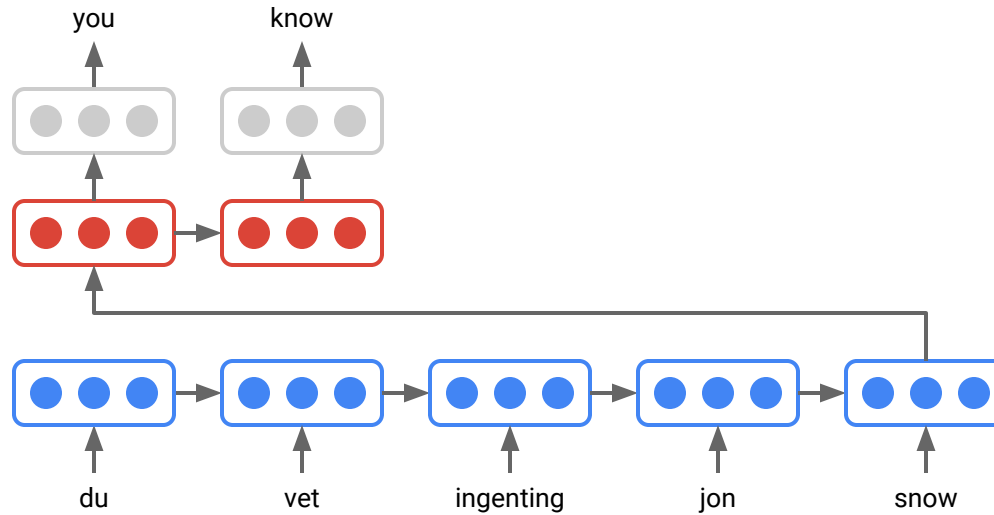
Encoder-Decoder models: no syntax, no hierarchical structure



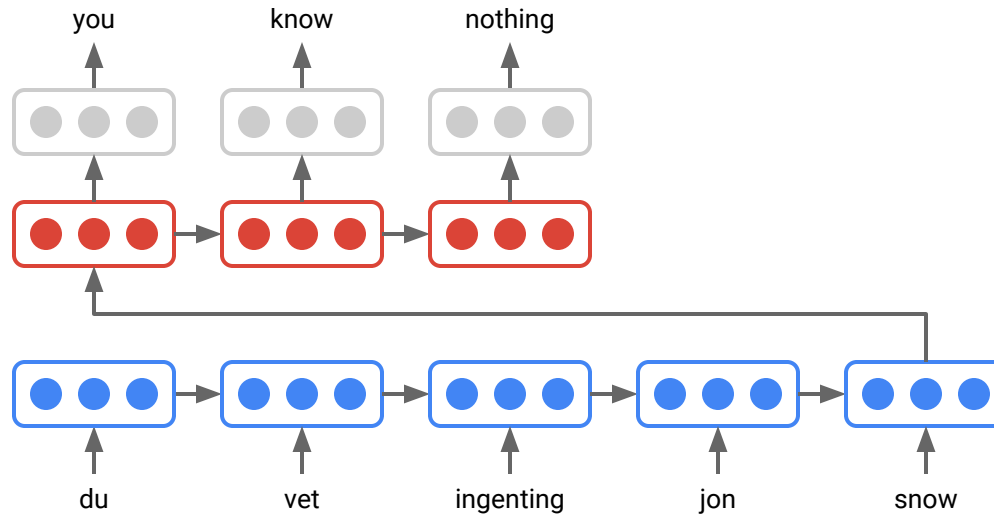
Encoder-Decoder models: no syntax, no hierarchical structure



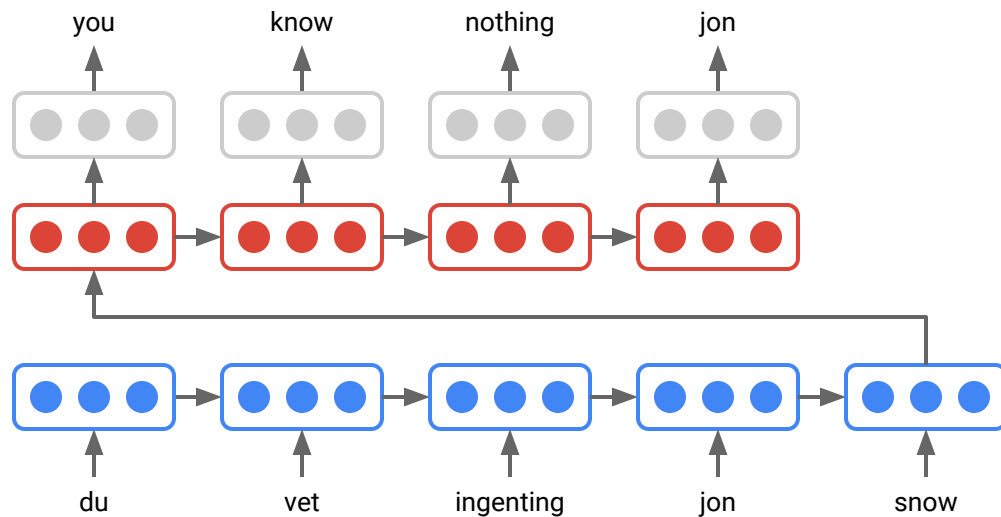
Encoder-Decoder models: no syntax, no hierarchical structure



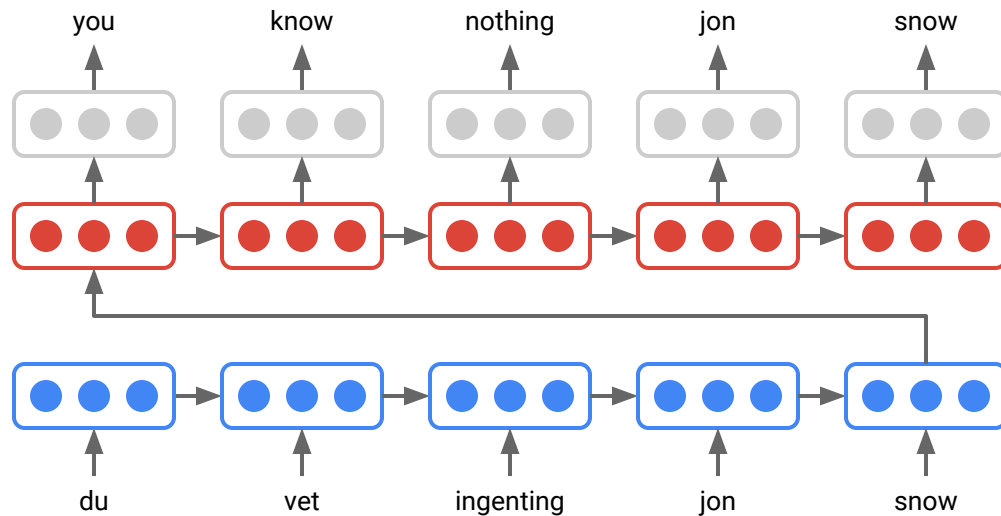
Encoder-Decoder models: no syntax, no hierarchical structure



Encoder-Decoder models: no syntax, no hierarchical structure



Encoder-Decoder models: no syntax, no hierarchical structure

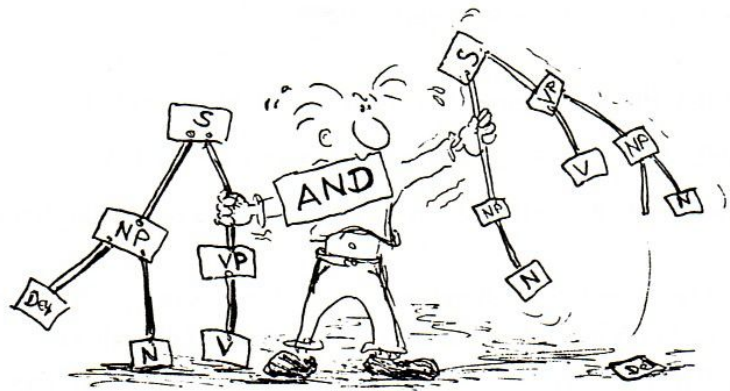


Syntax can help!

Syntax:

- captures **non-local dependencies** between words
- captures **relations** between words

How to incorporate syntax into NMT?



Related Work

Multi-task learning / latent-graph:

- **Eriguchi et al.** (ACL, 2017). Learning to parse and translate improves NMT.
- **Hashimoto et al.** (EMNLP, 2017). NMT with Source-Side Latent Graph Parsing.
- **Nadejde et al.** (WMT, 2017). Predicting Target Language CCG Supertags Improves NMT.

Tree-RNN:

- **Eriguchi et al.** (ACL, 2016). Tree-to-Sequence Attentional NMT.
- **Chen et al.** (ACL, 2017). Improved NMT with a Syntax-Aware Encoder and Decoder.
- **Yang et al.** (EMNLP, 2017). Towards Bidirectional Hierarchical Representations for Attention-based NMT.

Serialized trees:

- **Aharoni & Goldberg** (ACL, 2017). Towards String-to-Tree NMT.
- **Li et al.** (ACL, 2017). Modeling Source Syntax for NMT.

Sequence-to-Dependency:

- **Wu et al.** (ACL, 2017). Sequence-to-Dependency NMT.



$(((\lambda())(\lambda()) 'yoav))))$ @yoavgo

18 Apr

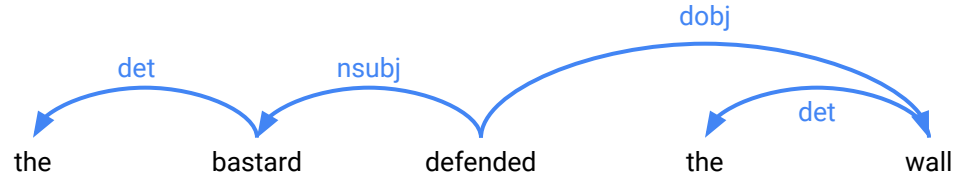
well, 2017 seems to be the year of
syntax-in-NMT :)

We have it on target side, but seems
effective also at source.

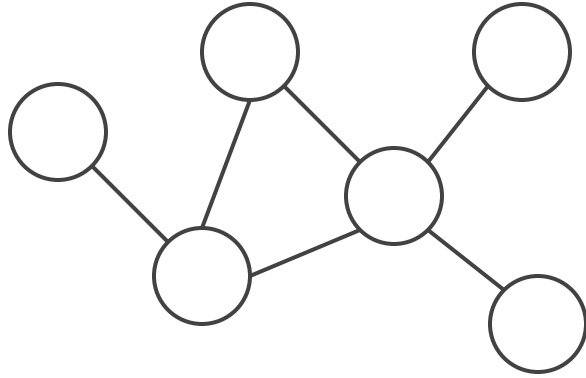
[syntax is good!]

Graph Convolutional Networks: Neural Message Passing

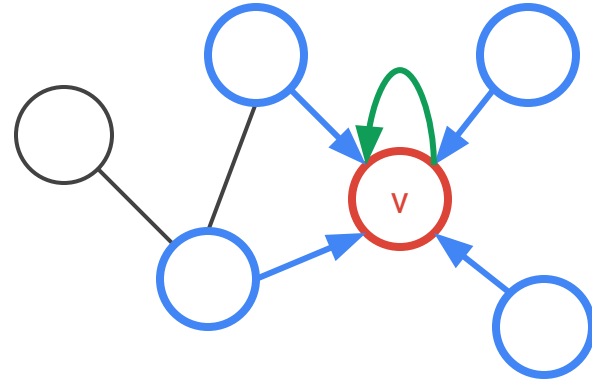
Many linguistic representations can be expressed as a graph



Graph Convolutional Networks: message passing

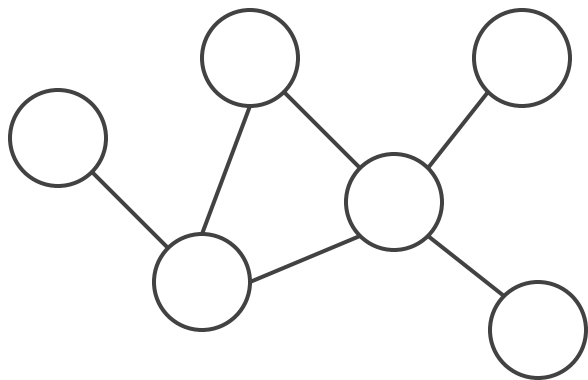


Undirected graph

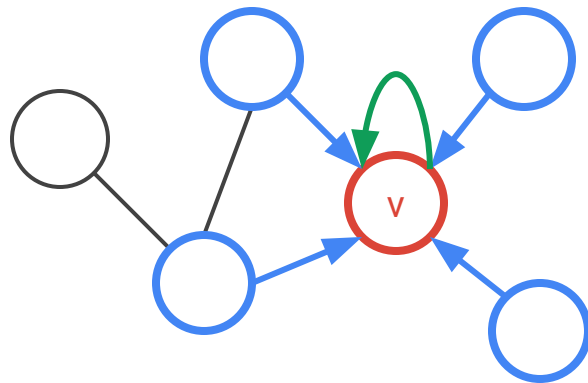


Update of node v

Graph Convolutional Networks: message passing



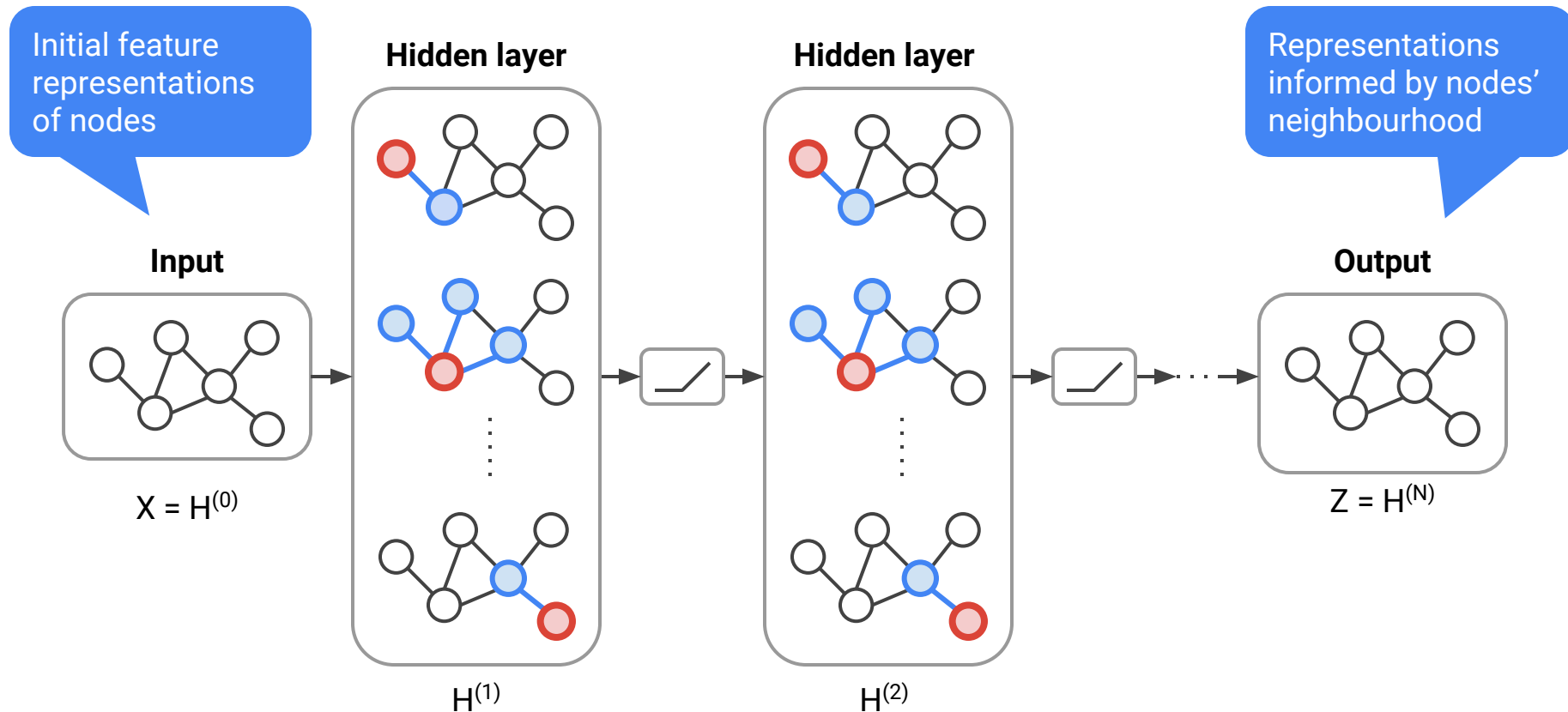
Undirected graph



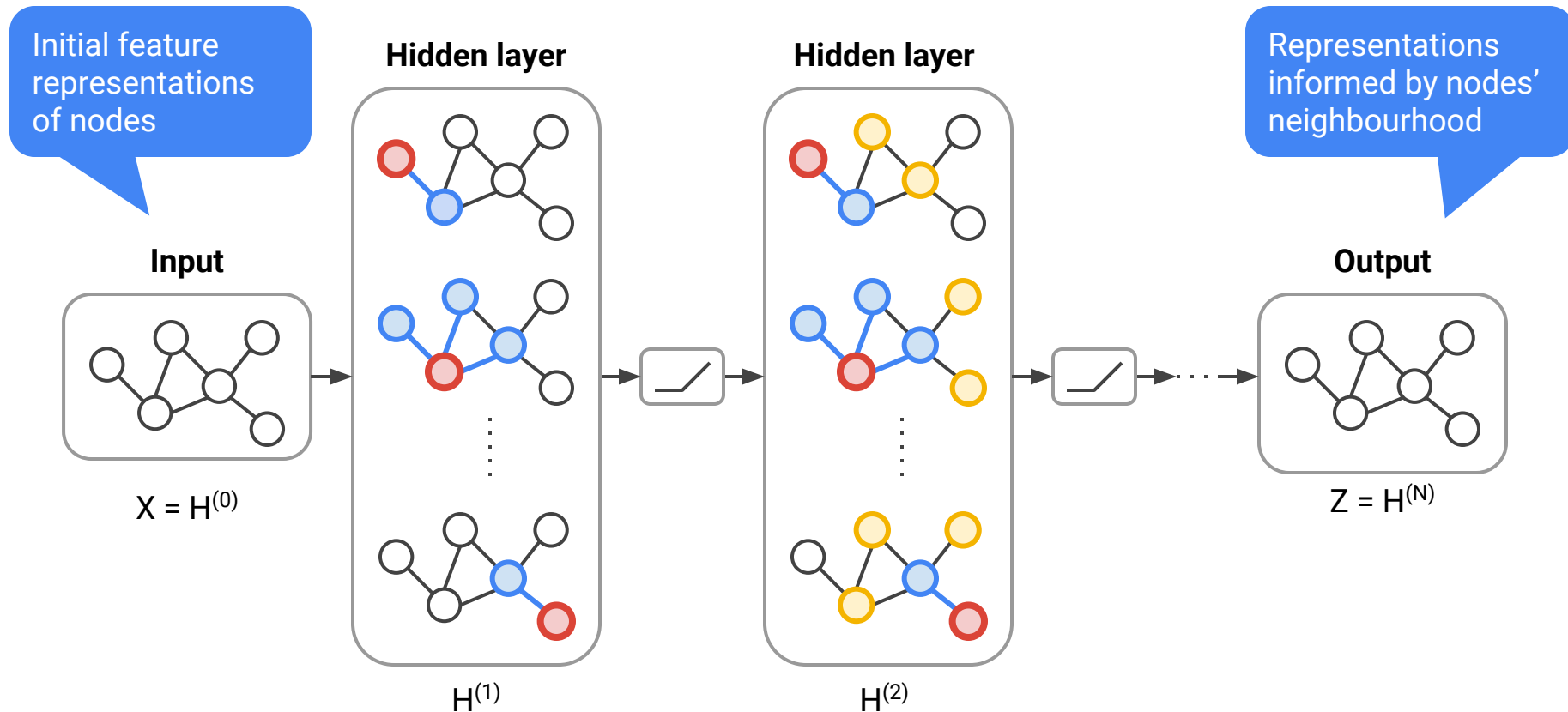
Update of node v

$$\mathbf{h}_v = \text{ReLU}(W_{\text{loop}}\mathbf{h}_v + \sum_{u \in \mathcal{N}(v)} W\mathbf{h}_u)$$

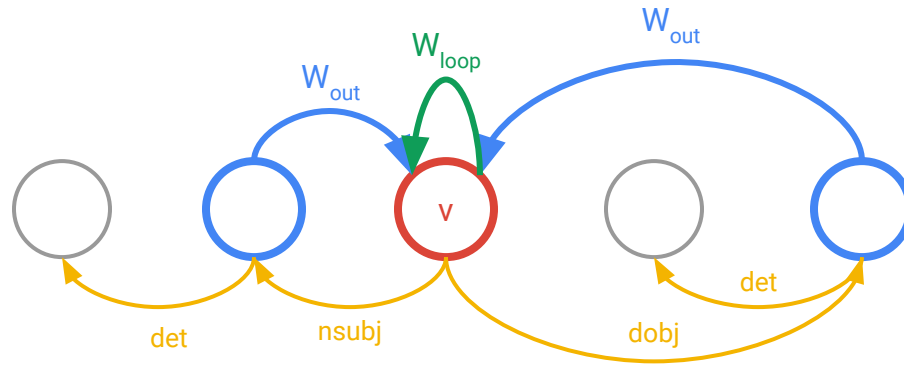
GCNs: multilayer convolution operation



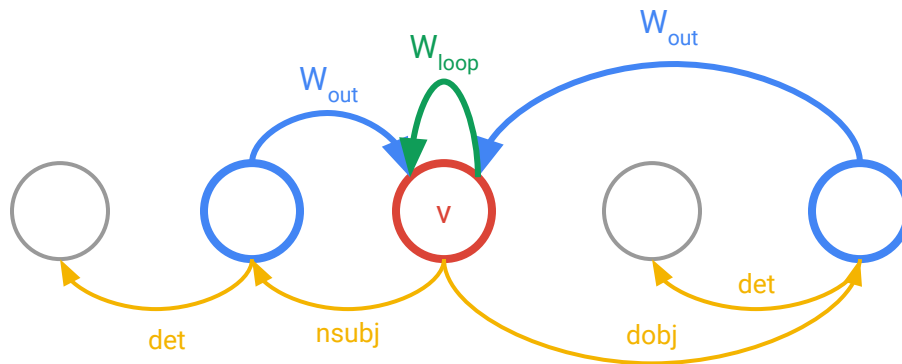
GCNs: multilayer convolution operation



Syntactic GCNs: directionality and labels



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}

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

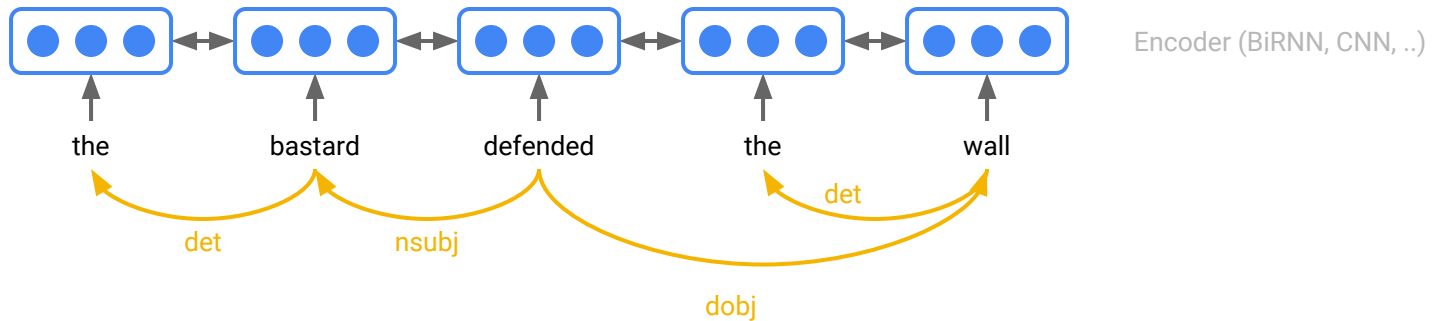
Syntactic GCNs: edge-wise gating

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

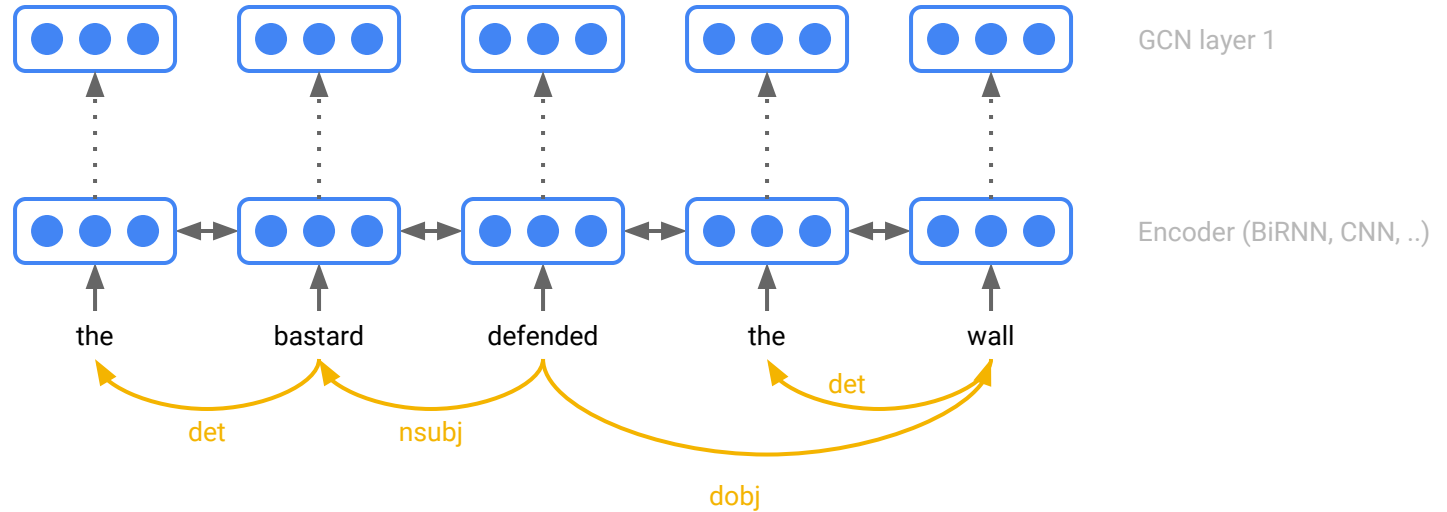
$$\mathbf{h}_v = \text{ReLU} \left(\sum_{u \in \mathcal{N}(v)} g_{u,v} \left(W_{\text{dir}(u,v)} \mathbf{h}_u + \mathbf{b}_{\text{lab}(u,v)} \right) \right)$$

Graph Convolutional Encoders

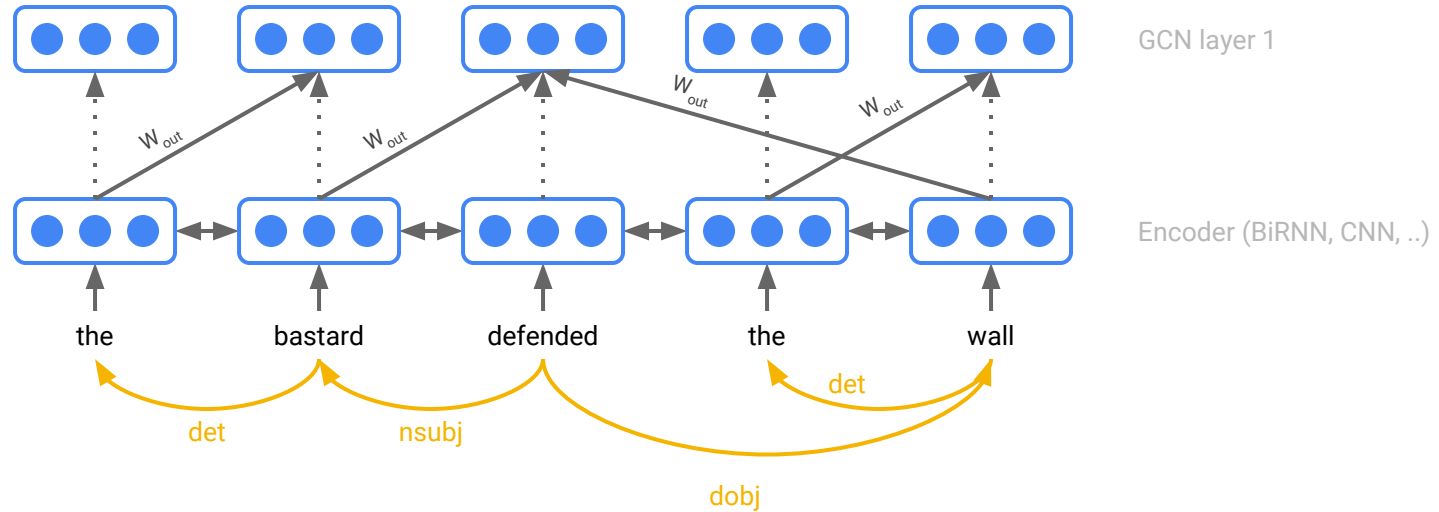
Graph Convolutional Encoders



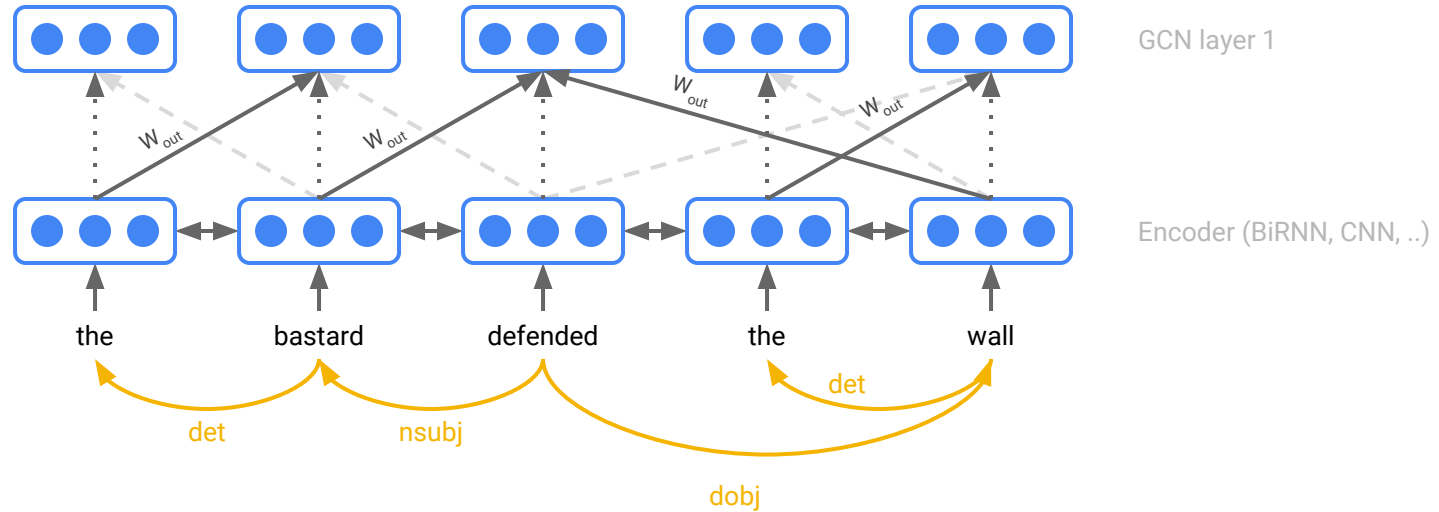
Graph Convolutional Encoders



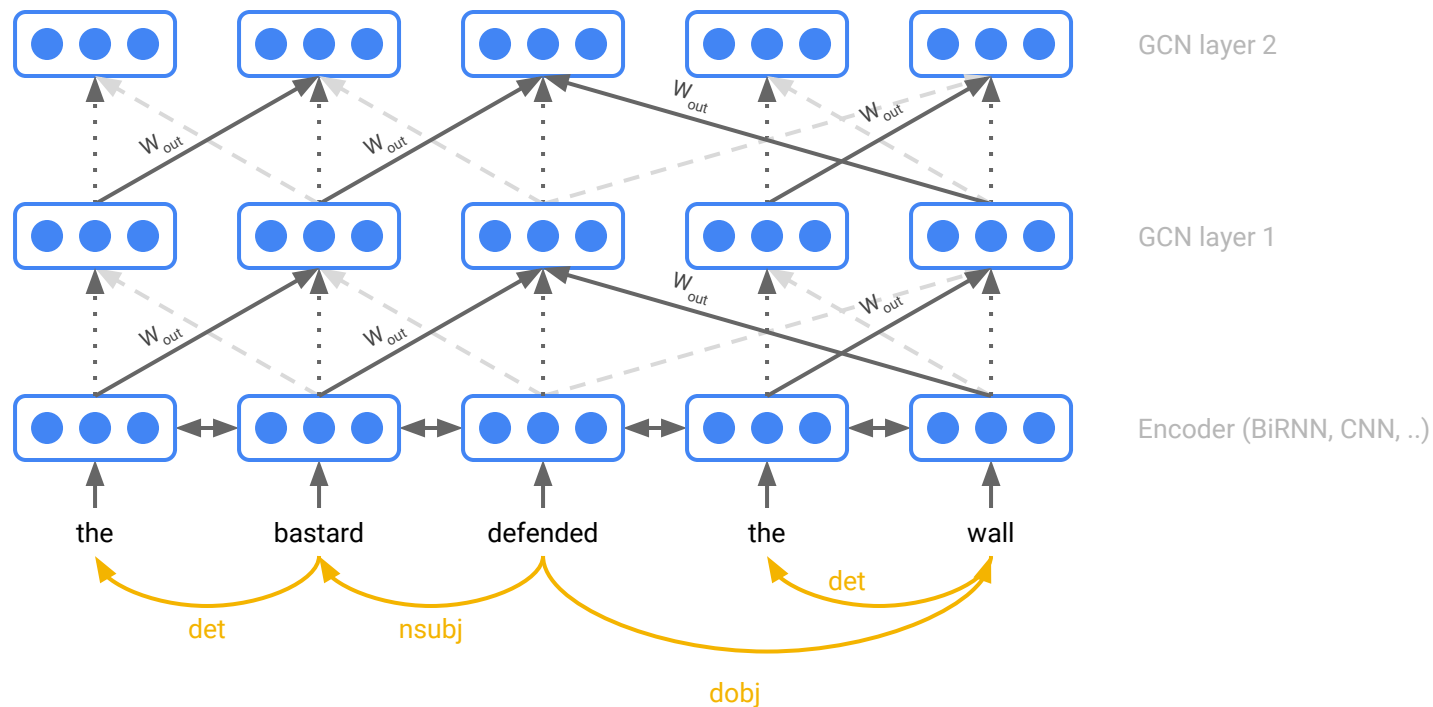
Graph Convolutional Encoders



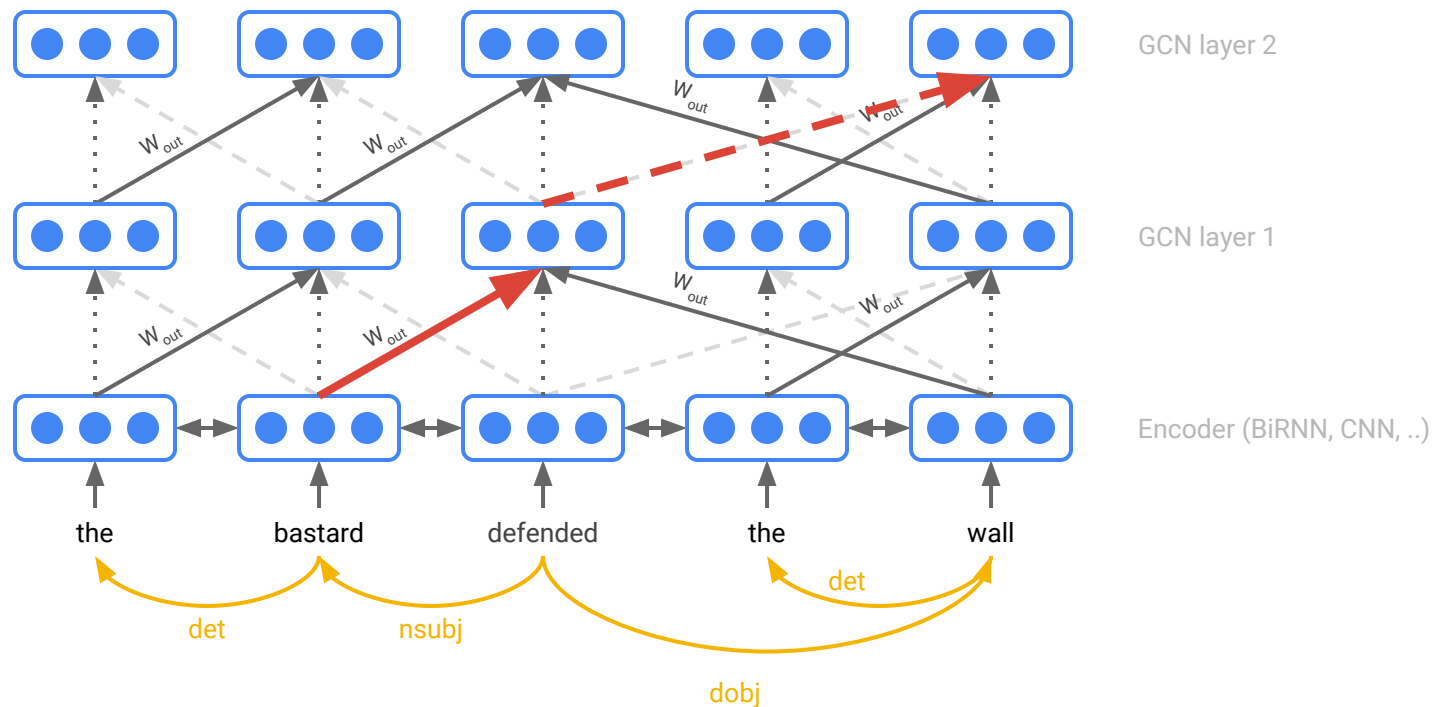
Graph Convolutional Encoders



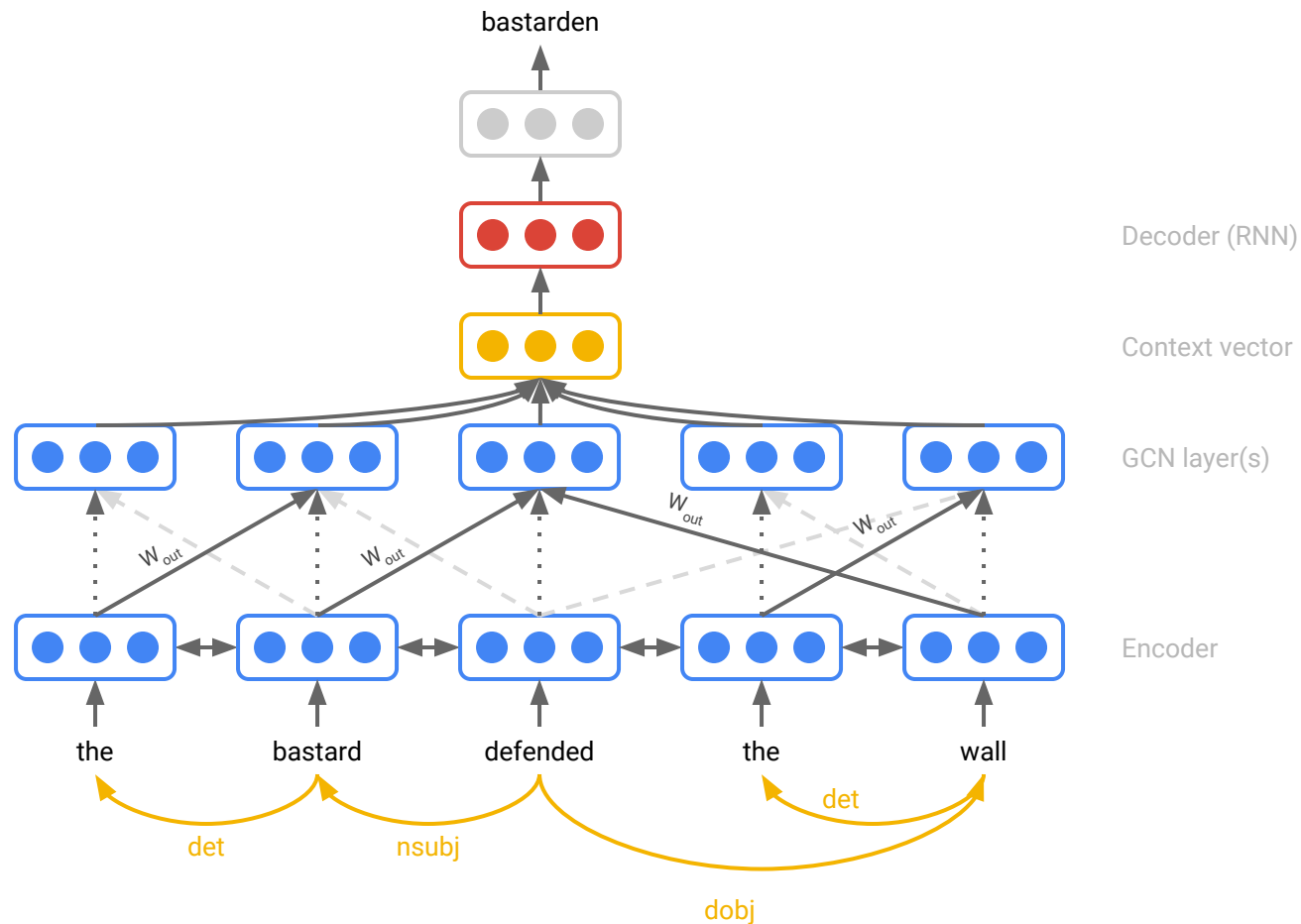
Graph Convolutional Encoders



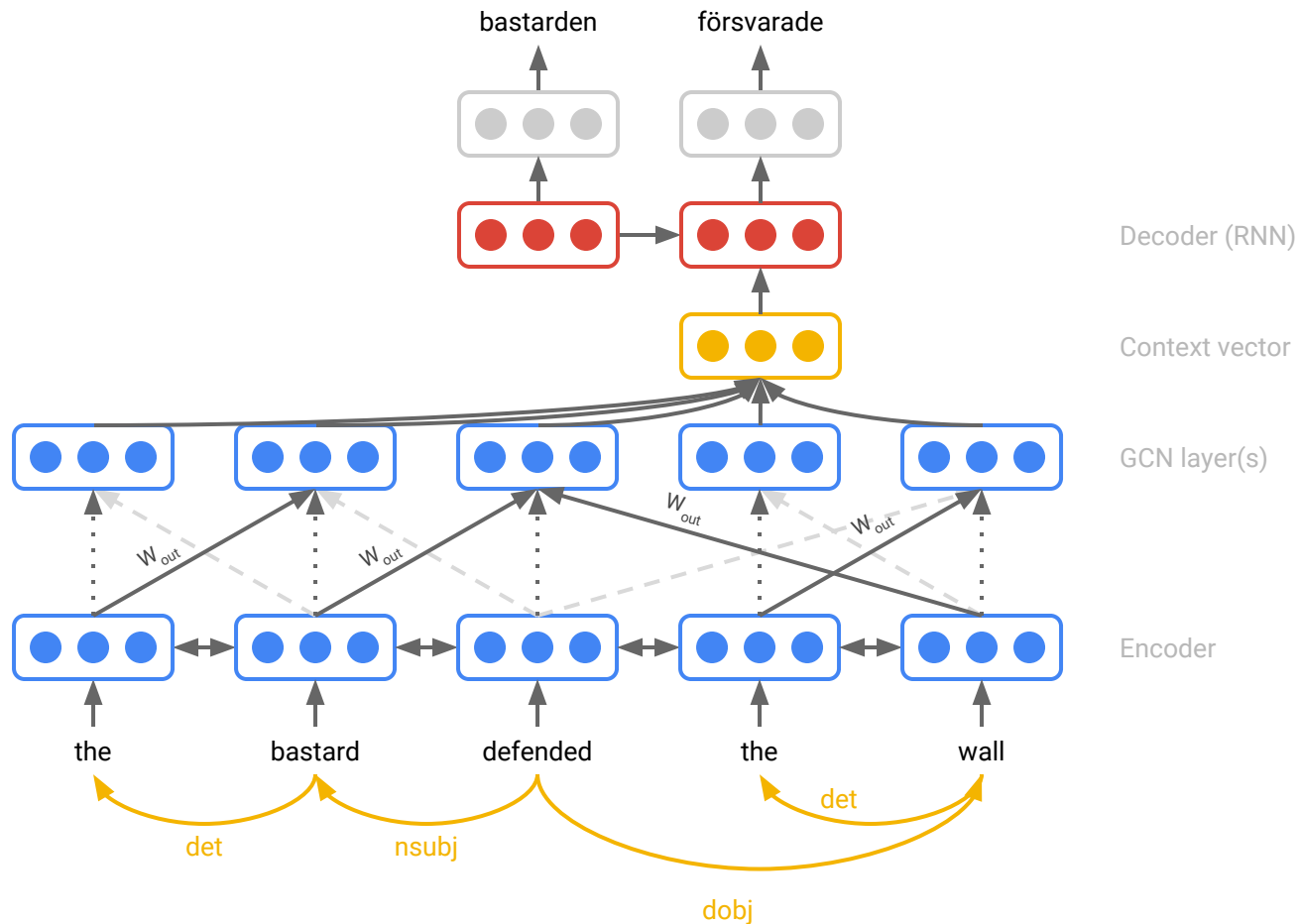
Graph Convolutional Encoders



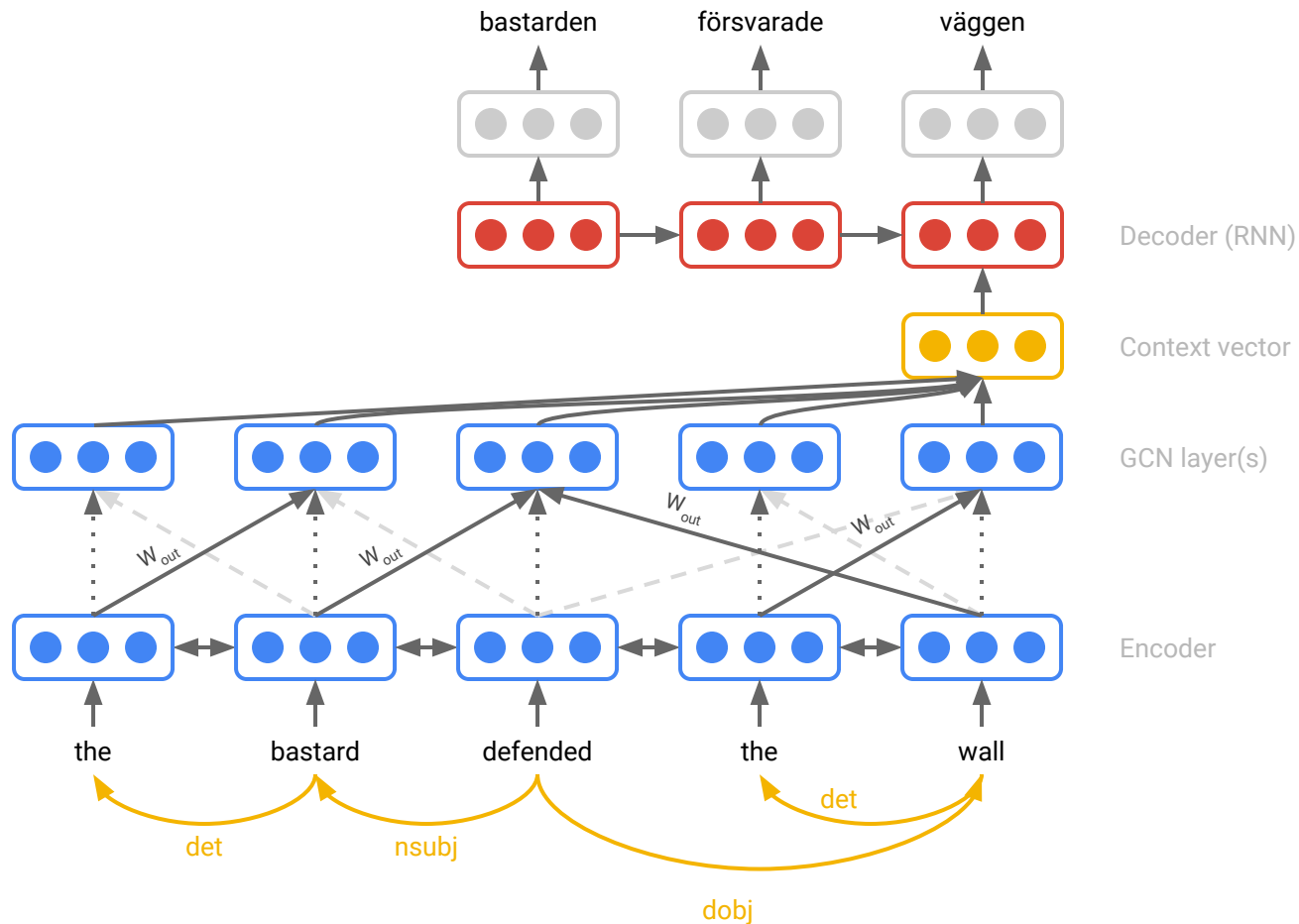
GCNs for NMT



GCNs for NMT



GCNs for NMT

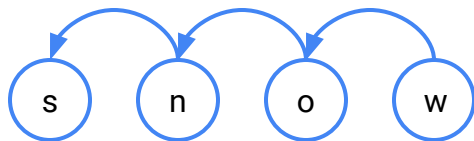


Experiments

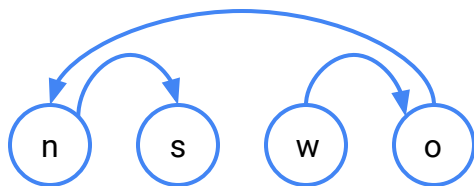
Artificial reordering experiment

Random sequences from **26 types**

Each token is linked to its **predecessor**



Sequences are randomly **permuted**

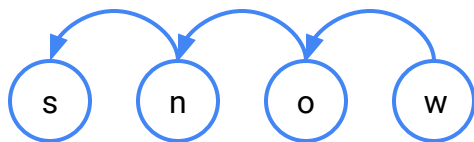


Each token is also linked to a **random** token with a “fake edge” using a different label set

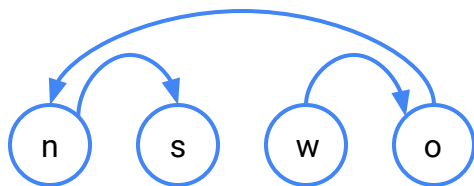
Artificial reordering experiment

Random sequences from **26 types**

Each token is linked to its **predecessor**



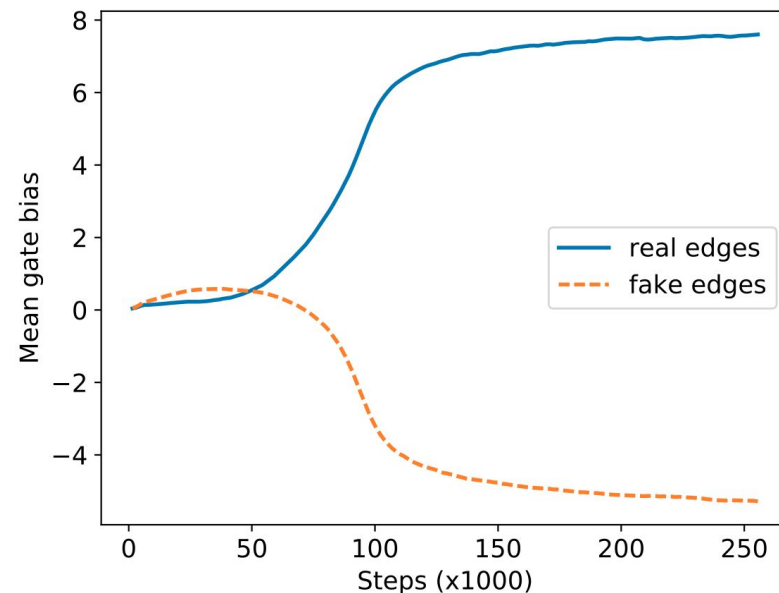
Sequences are randomly **permuted**



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**.

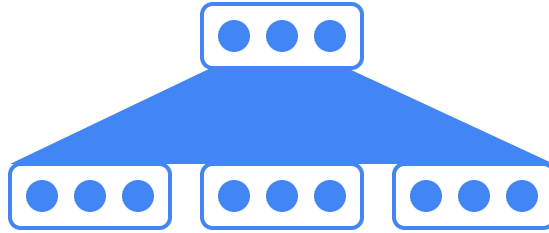
Real edges are distinguished from fake edges.



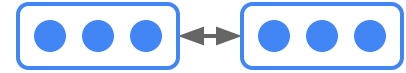
Three baselines



Bag of words



Convolutional



Recurrent

Machine Translation experiments

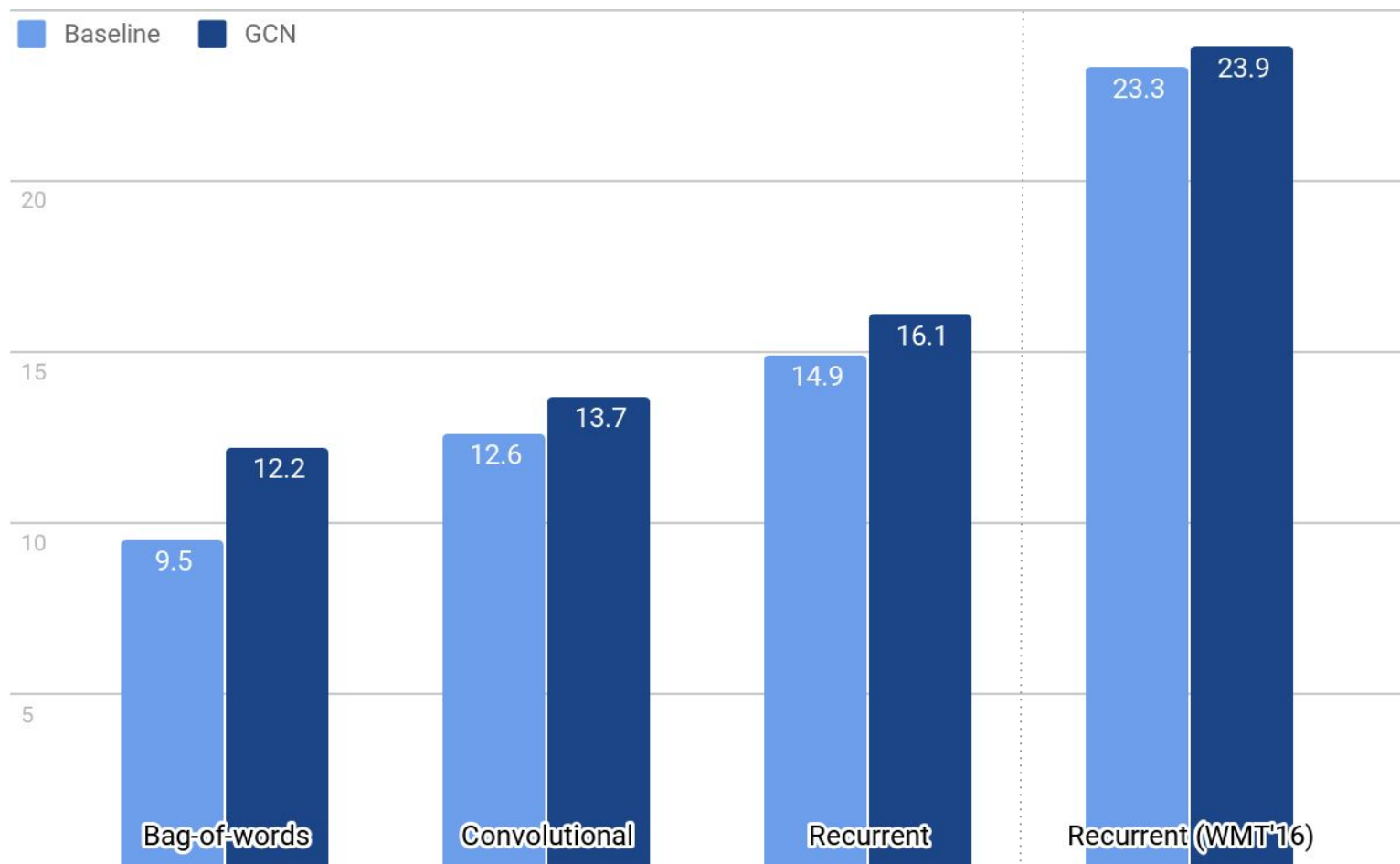
		Train	Validation <i>newstest2015</i>	Test <i>newstest2016</i>
English-German	NCv11	227 k		
	WMT'16	4.5 M	2169	2999
English-Czech	NCv11	181 k	2656	2999

We parse English source sentences using **SyntaxNet**

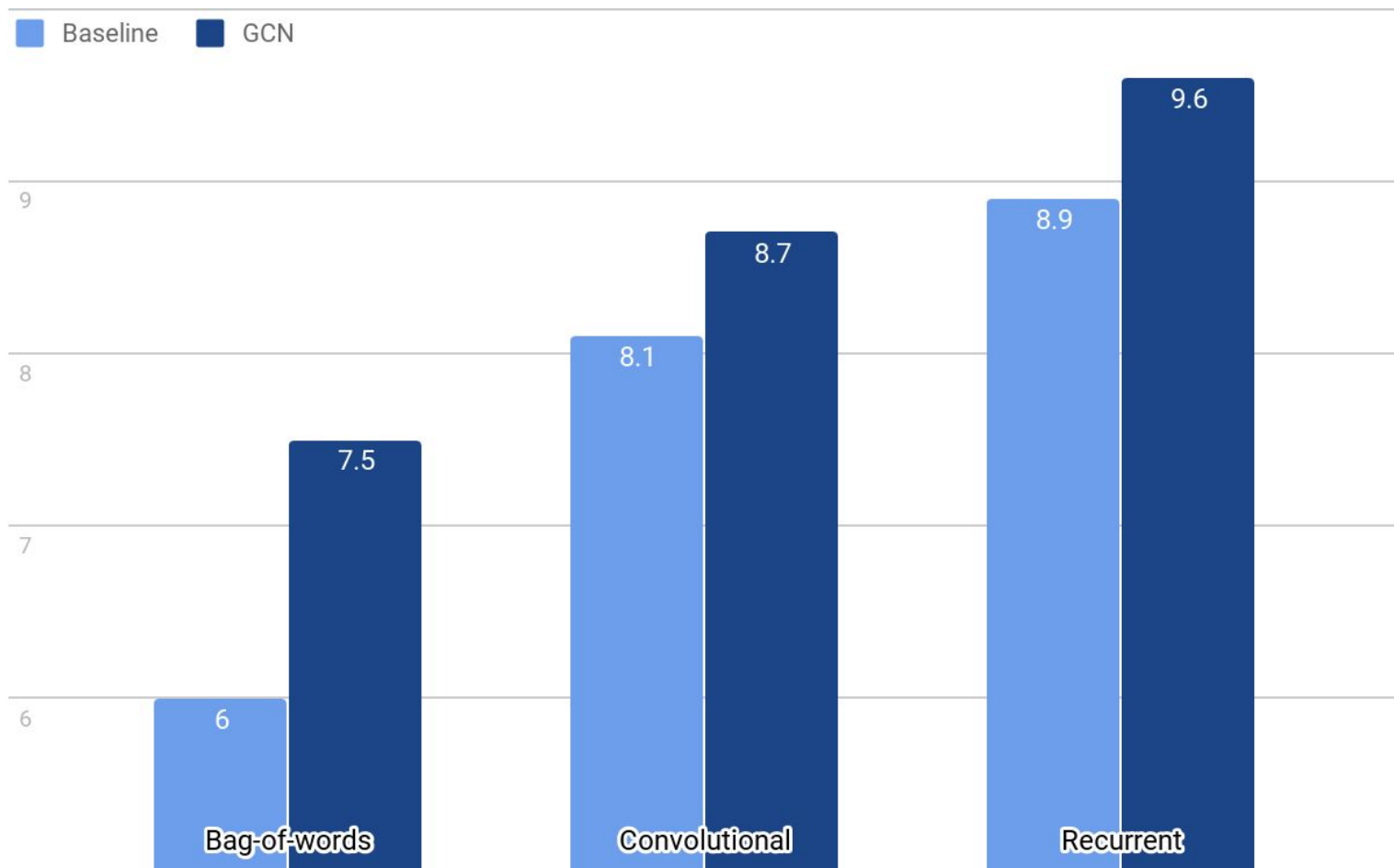
BPE on **target side** (8k merges, 16k for WMT'16)

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

Results English-German (BLEU)



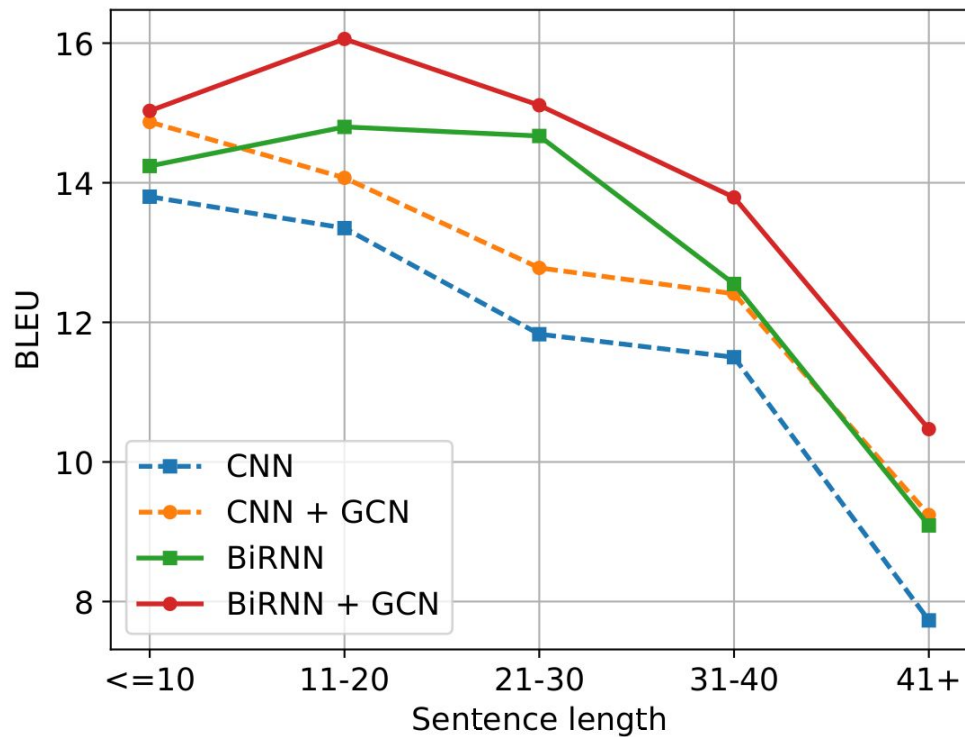
Results English-Czech (BLEU)



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

Conclusion

Simple and effective approach to integrating syntax into NMT models

Consistent BLEU improvements for two challenging language pairs

GCNs are capable of encoding any kind of graph-based structure

Thank you!

Code: Neural Monkey + GCN (TensorFlow)

joost.ninja

Code: Semantic Role Labeling (Theano)

diegma.github.io

Ivan is hiring PhDs/Postdocs. Deadline tomorrow!

This work was supported by the European Research Council (ERC StG BroadSem 678254) and the Dutch National Science Foundation (NWO VIDI 639.022.518, NWO VICI 277-89-002).