

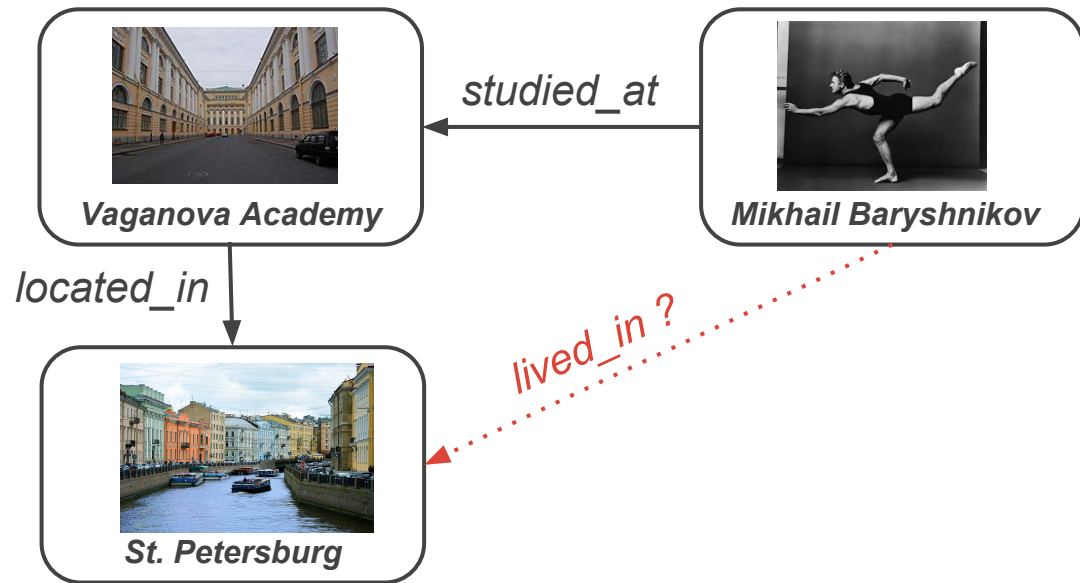
Extracting and Modeling Relations with Graph Convolutional Networks

Ivan Titov

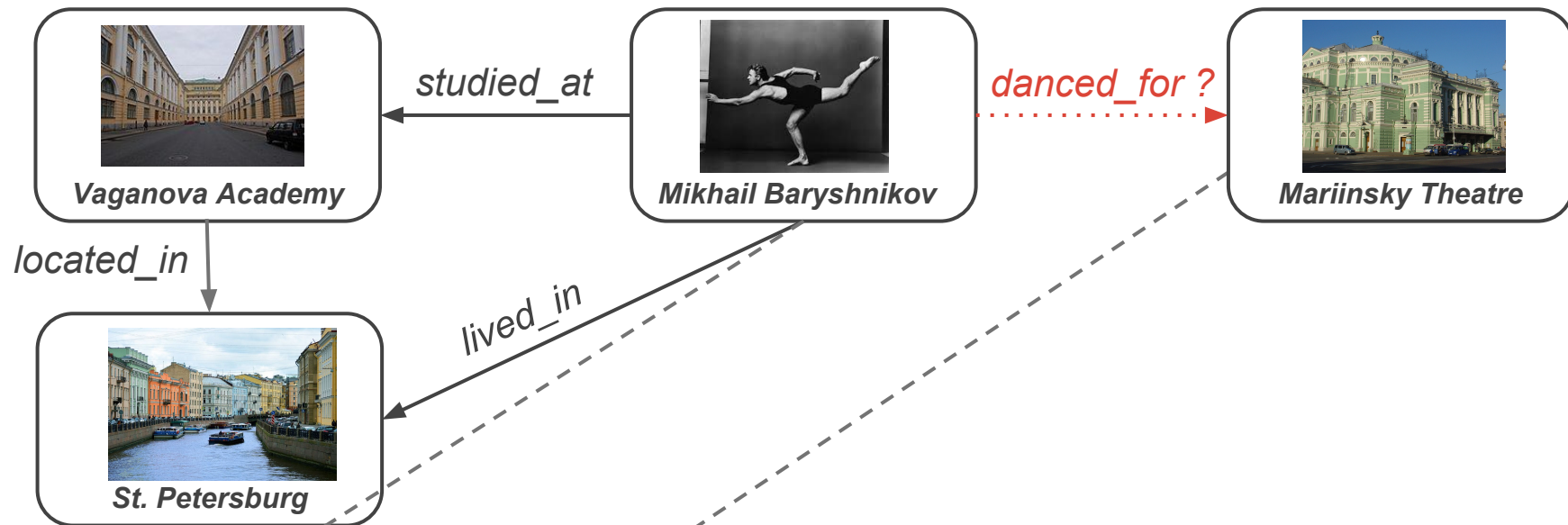
with Diego Marcheggiani, Michael Schlichtkrull, Thomas Kipf, Max Welling,
Rianne van den Berg and Peter Bloem



Inferring missing facts in knowledge bases: link prediction



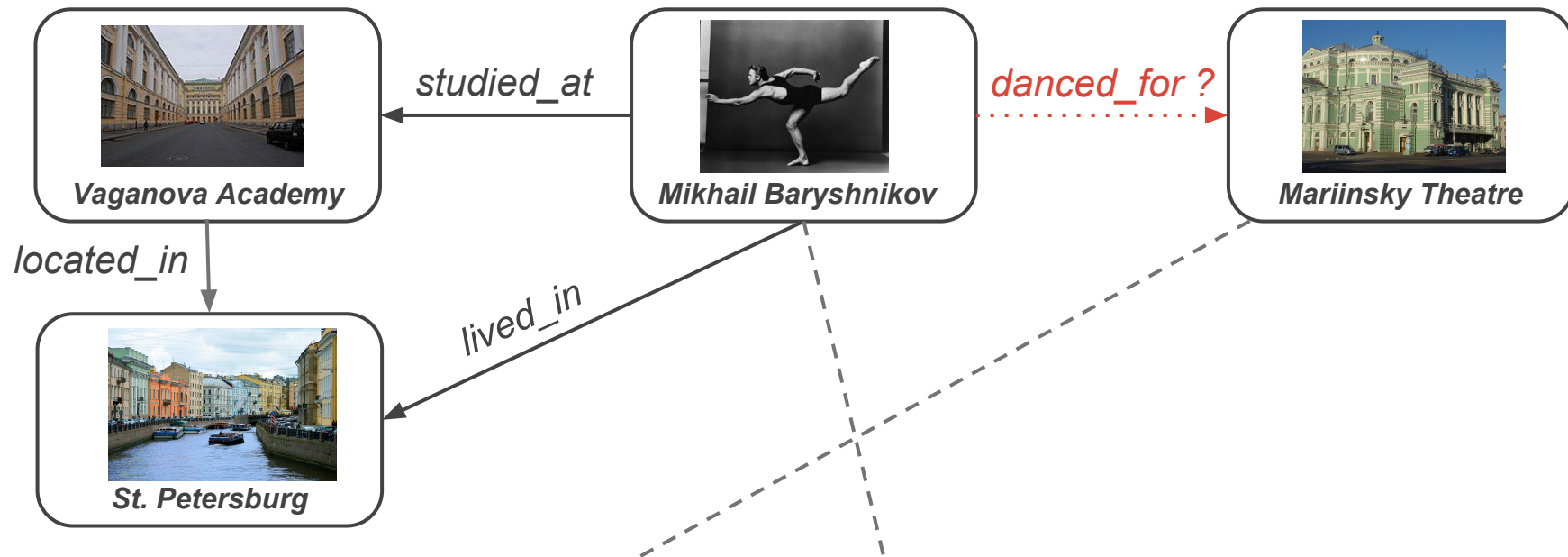
Relation Extraction



Baryshnikov danced for Mariinsky based in what was then Leningrad (now St. Petersburg)

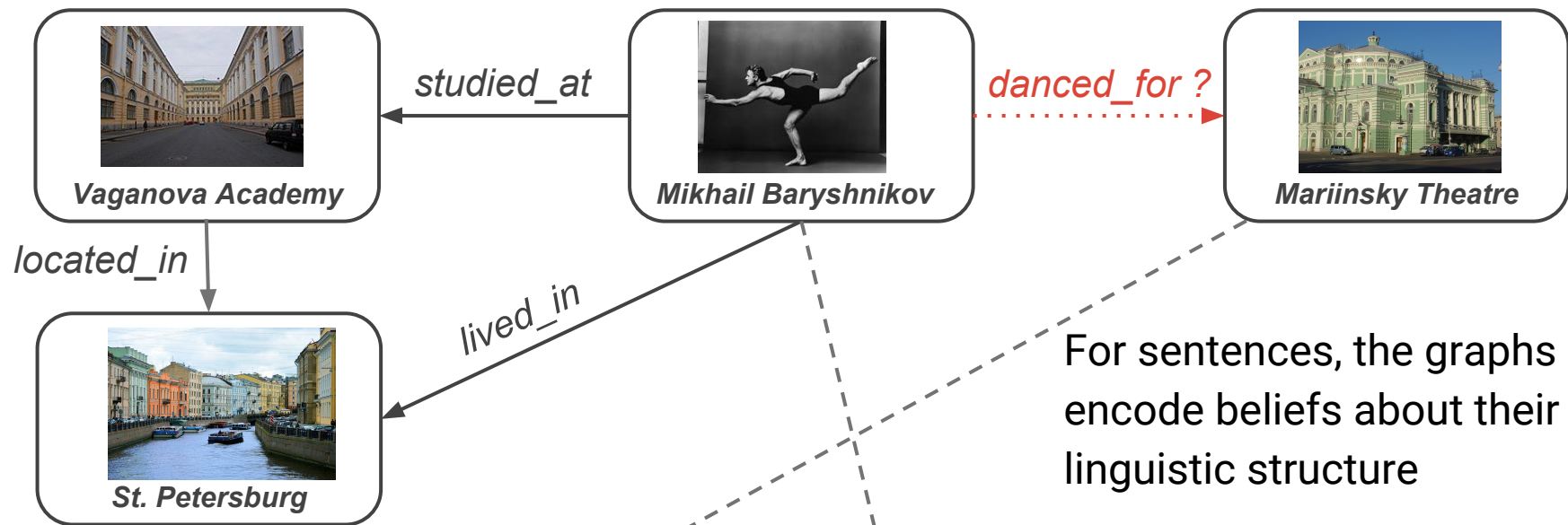
danced_for

Generalization of link prediction and relation extraction



After a promising start in Mariinsky ballet, Baryshnikov defected to Canada in 1974 ...

KBC: it is natural to represent both sentences and KB with **graphs**



For sentences, the graphs encode beliefs about their linguistic structure

After a promising start in Mariinsky ballet, Baryshnikov defected to Canada in 1974 ...

How can we model (and exploit) these graphs with **graph neural networks**?

Outline

Graph Convolutional Networks (GCNs)

Link Prediction with Graph Neural Networks

Relational GCNs

Denoising Graph Autoencoders for Link Prediction

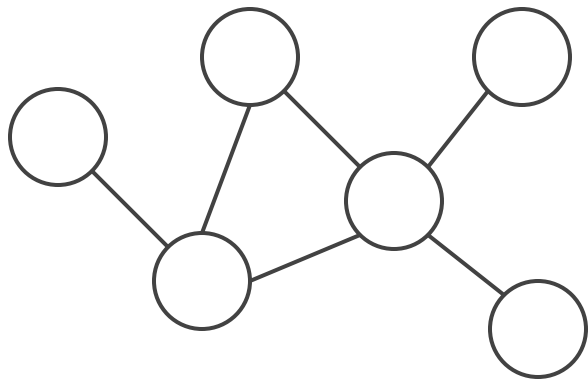
Extracting Semantic Relations: Semantic Role Labeling

Syntactic GCNs

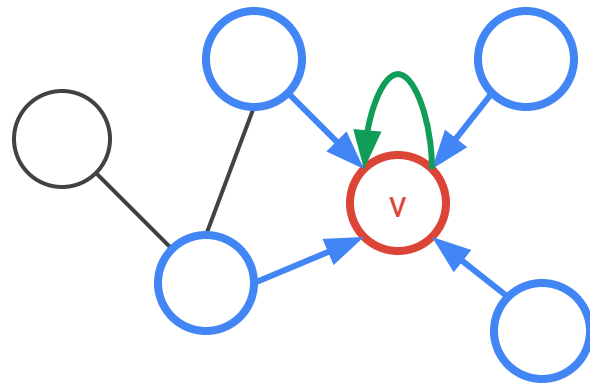
Semantic Role Labeling Model

Graph Convolutional Networks: Neural Message Passing

Graph Convolutional Networks: message passing

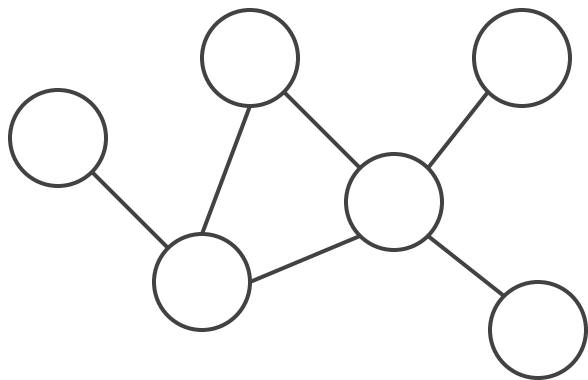


Undirected graph

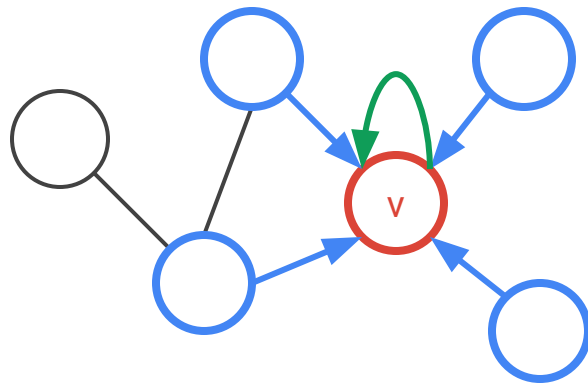


Update for node v

Graph Convolutional Networks: message passing



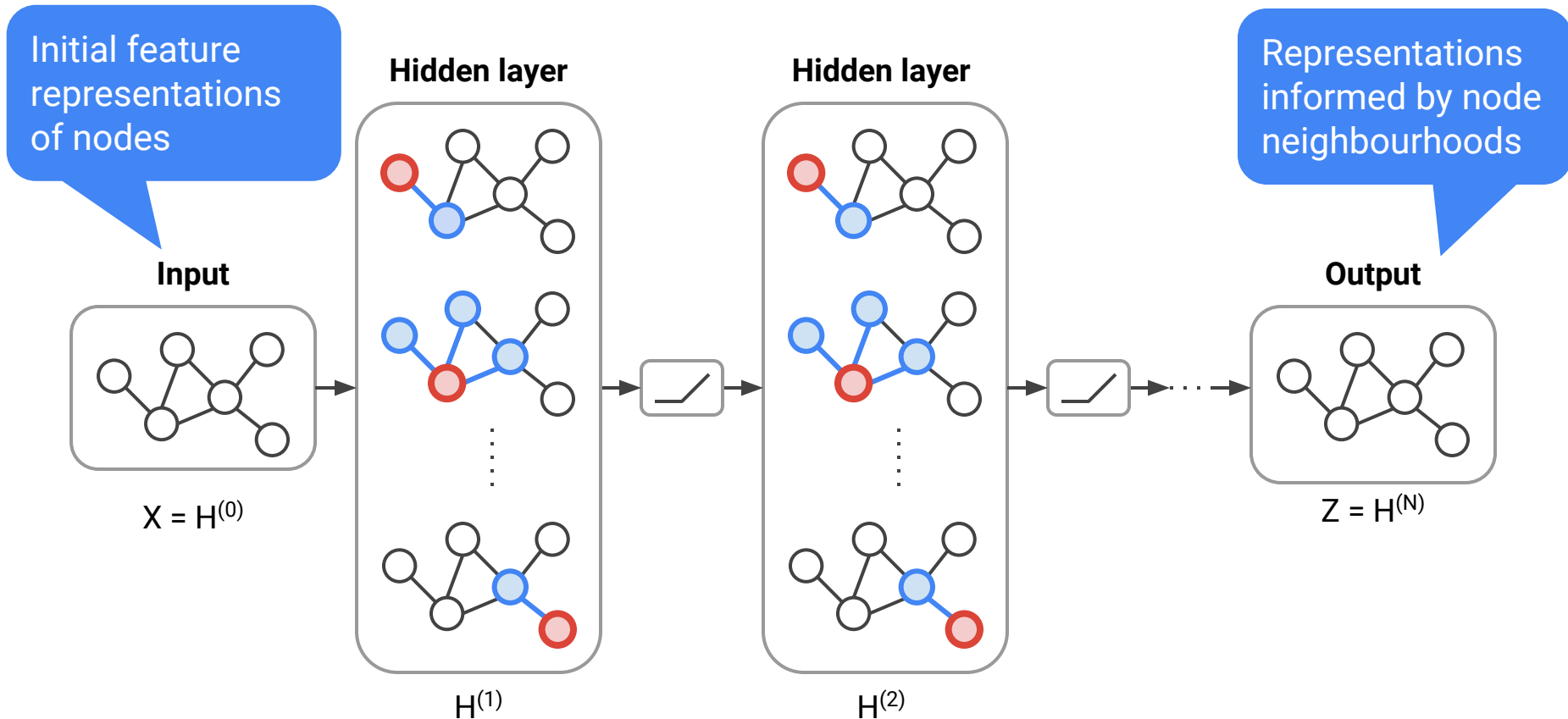
Undirected graph



Update for node v

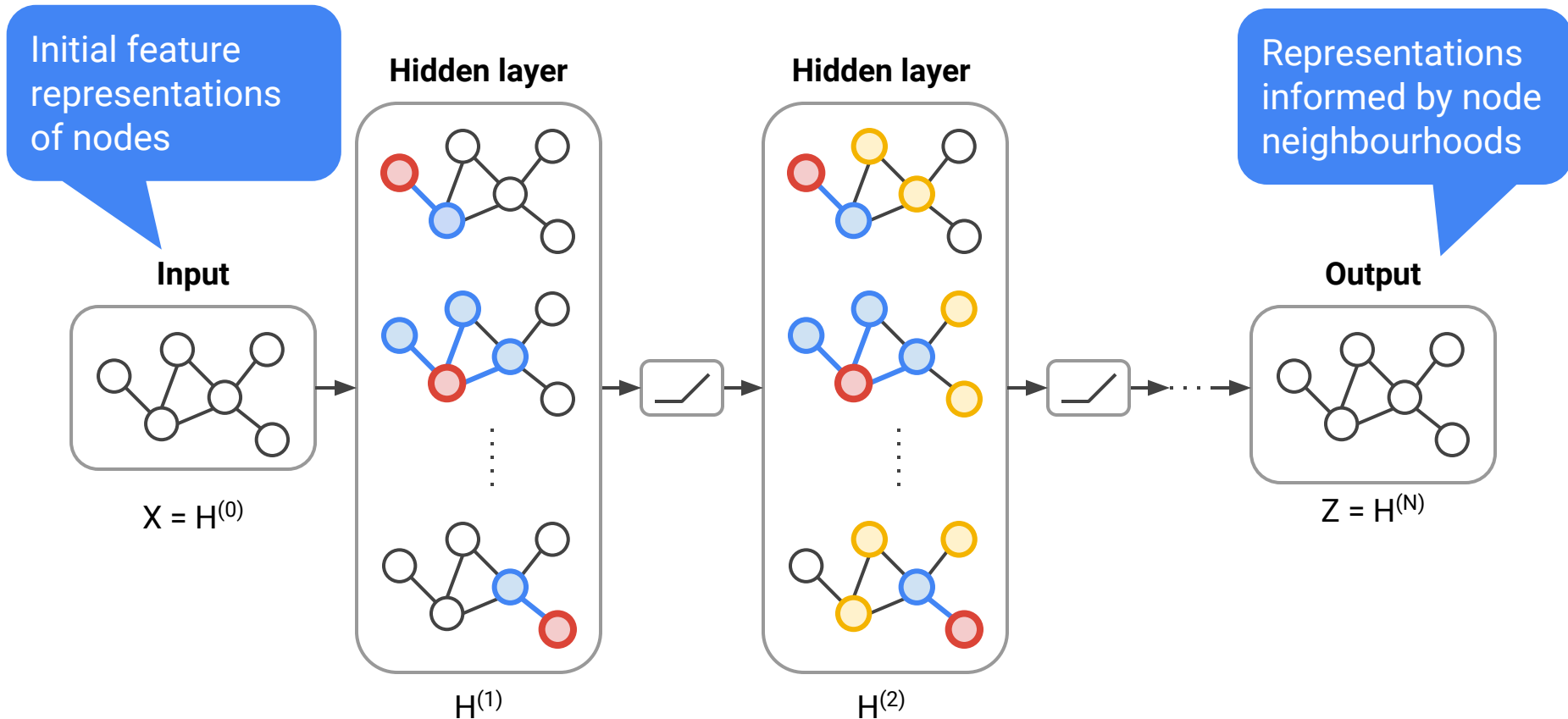
$$\mathbf{h}_v = \text{ReLU}(\mathbf{W}_{\text{loop}}\mathbf{h}_v + \sum_{u \in \mathcal{N}(v)} \mathbf{W}\mathbf{h}_u)$$

GCNs: multilayer convolution operation



Parallelizable computation, can be made quite efficient (e.g., Hamilton, Ying and Leskovec (2017)).

GCNs: multilayer convolution operation



Parallelizable computation, can be made quite efficient (e.g., Hamilton, Ying and Leskovec (2017)).

Graph Convolutional Networks: Previous work

Shown very effective on a range of problems - citations graphs, chemistry, ...

Mostly:

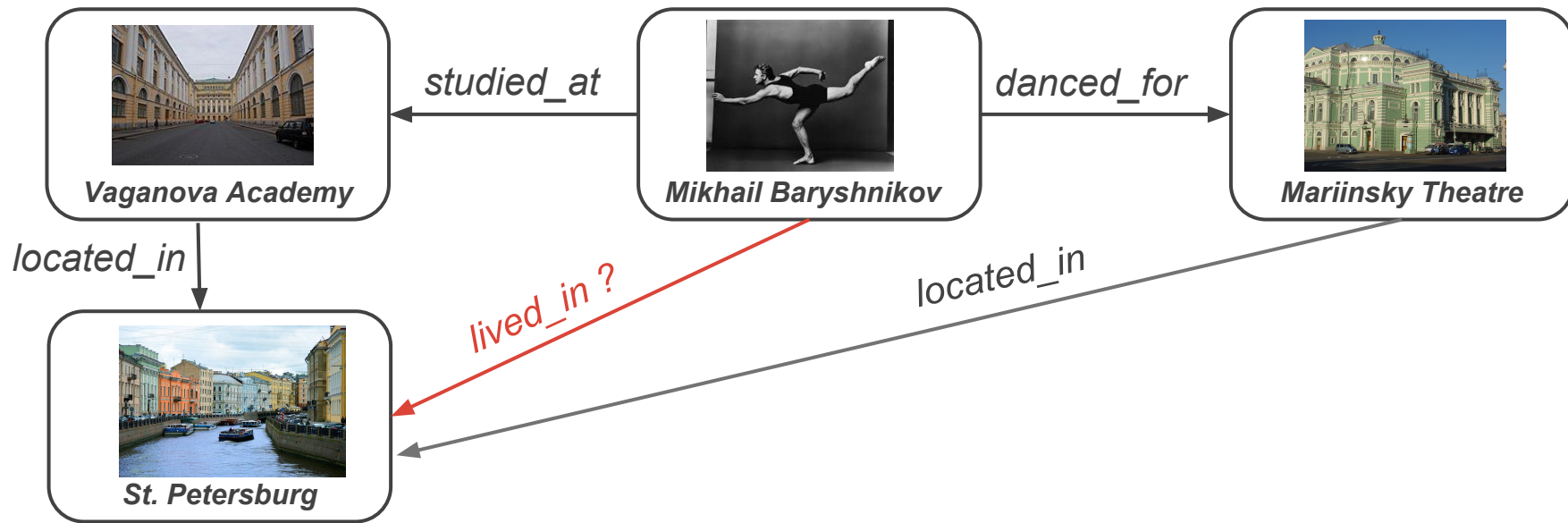
- Unlabeled and undirected graphs
- Node labeling in a single large graph (transductive setting)
- Classification of graphlets

How to apply GCNs to graphs we have in knowledge based completion / construction?

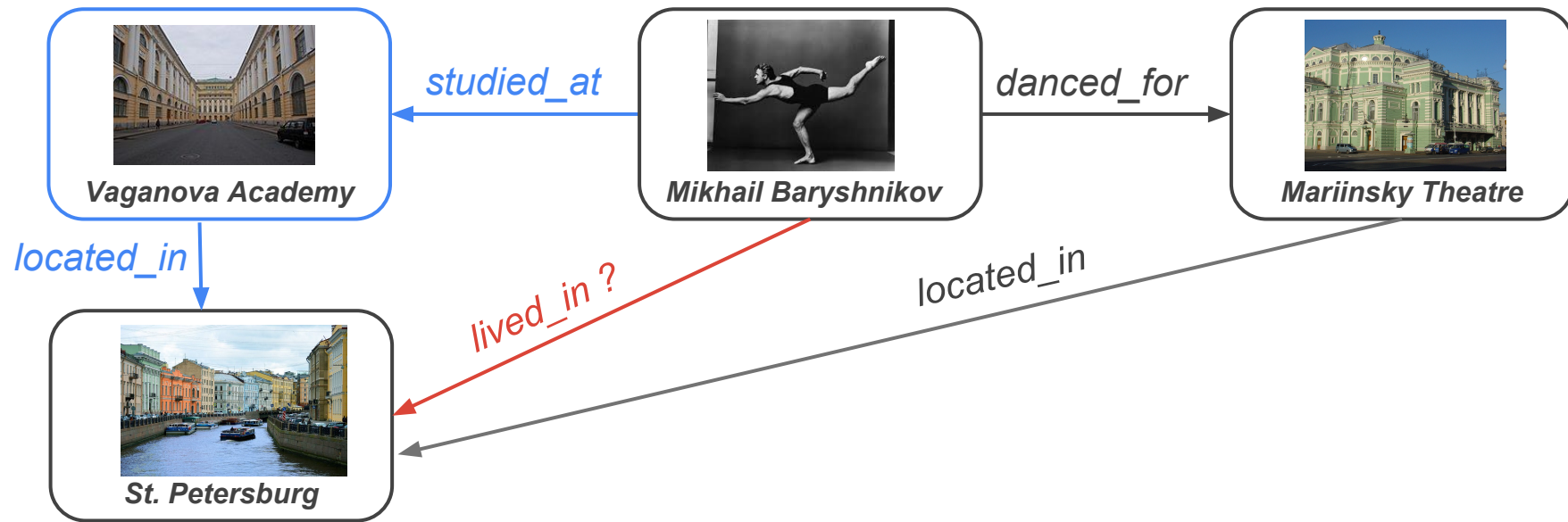
See Bronstein et al. (Signal Processing, 2017) for an overview

Link Prediction with Graph Neural Networks

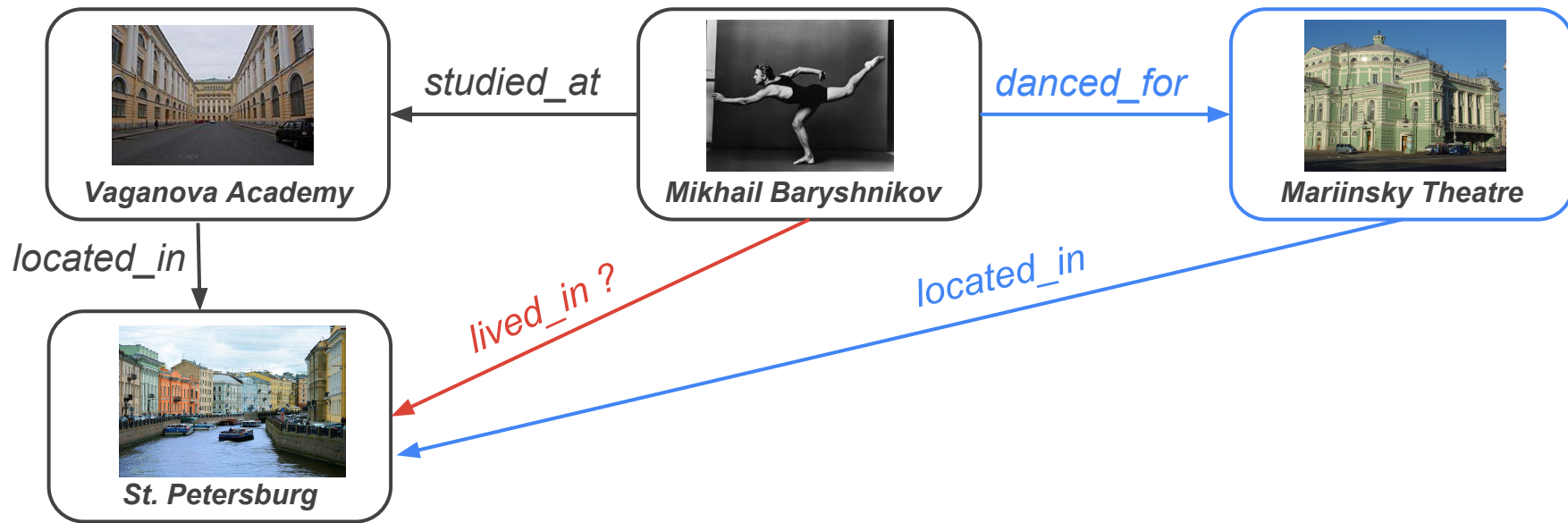
Link Prediction



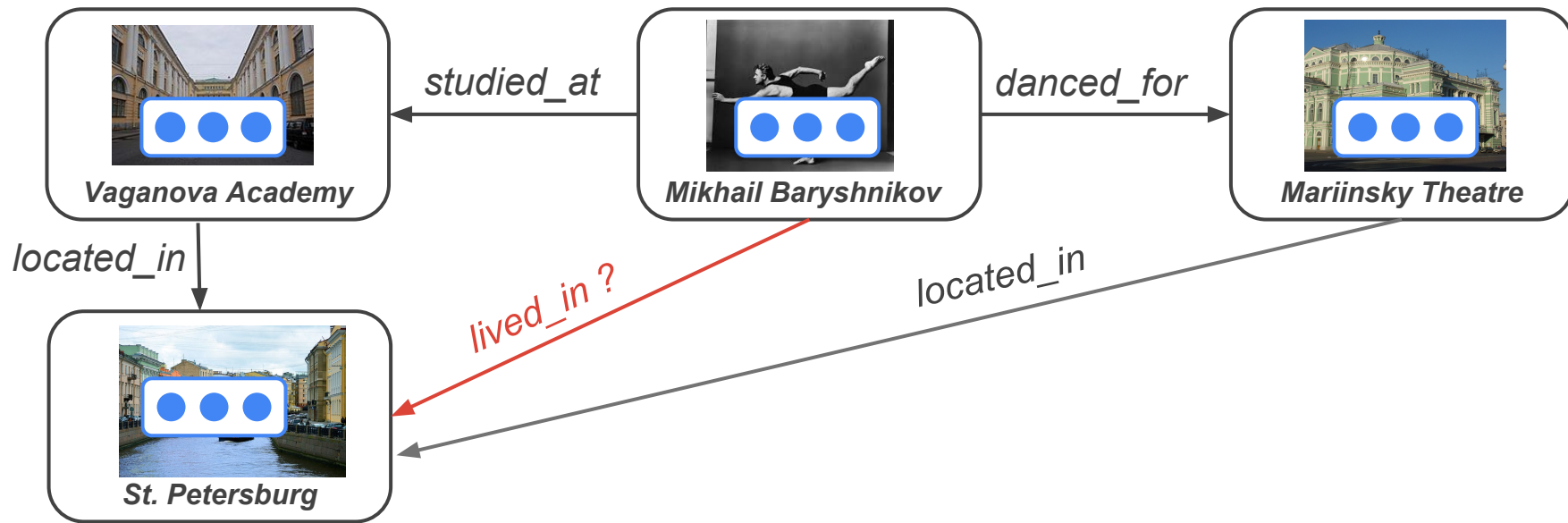
Link Prediction



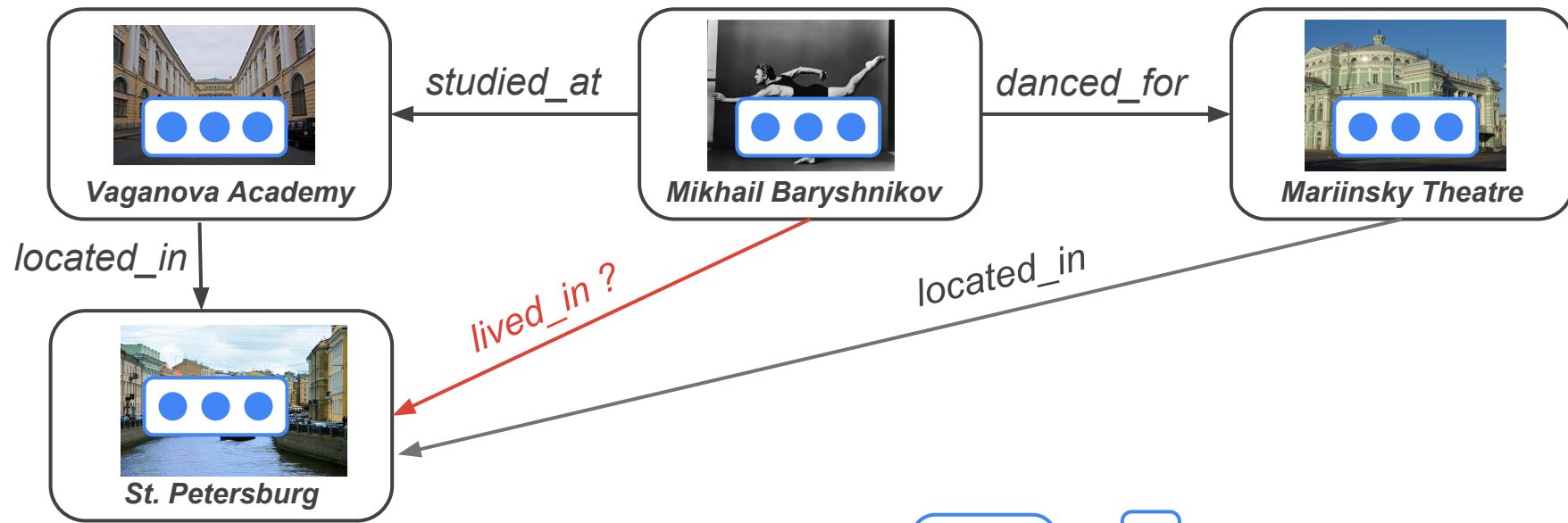
Link Prediction



KB Factorization



KB Factorization



A scoring function is used to predict whether a relation holds:

$$\begin{matrix} \text{Baryshnikov} & \times & \text{lived_in} & \times & \text{St. Petersburg} \end{matrix}$$

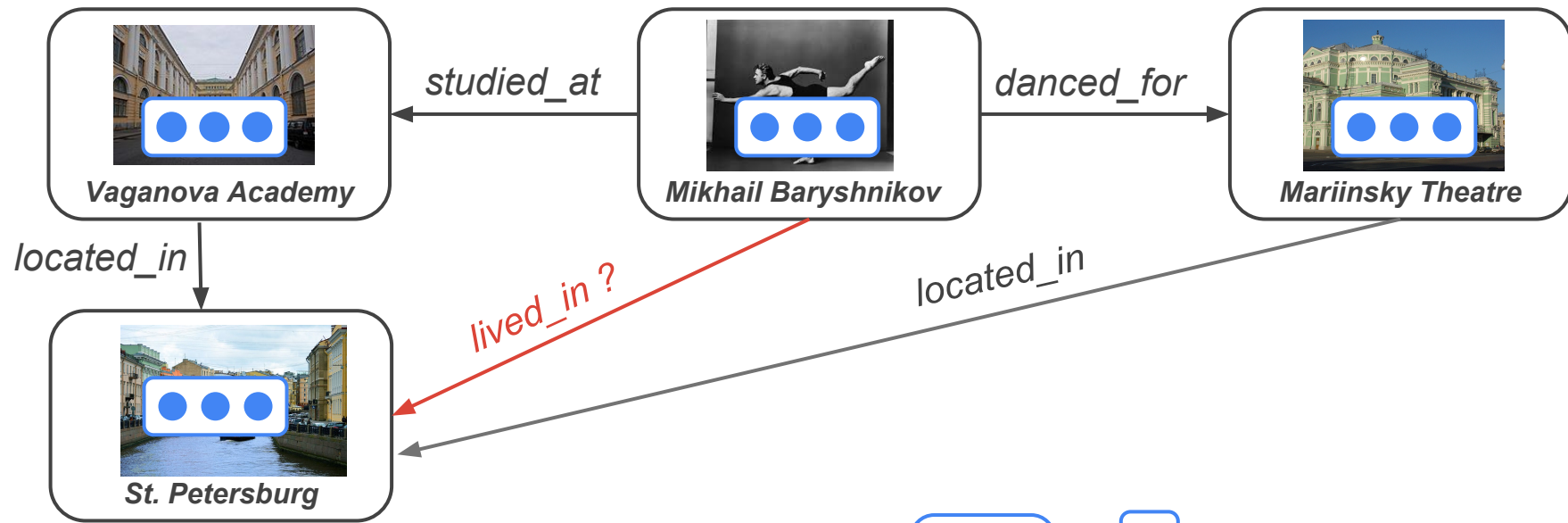
The diagram shows the matrix factorization process for the query *lived_in* between *Baryshnikov* and *St. Petersburg*. Each entity and relation is represented by a matrix of blue dots:

- Baryshnikov*: A 1x3 matrix.
- lived_in*: A 3x3 matrix.
- St. Petersburg*: A 3x1 matrix.

The matrices are multiplied together to score the relation.

RESICAL
(Nickel et al., 2011)

KB Factorization



A scoring function is used to predict whether a relation holds:

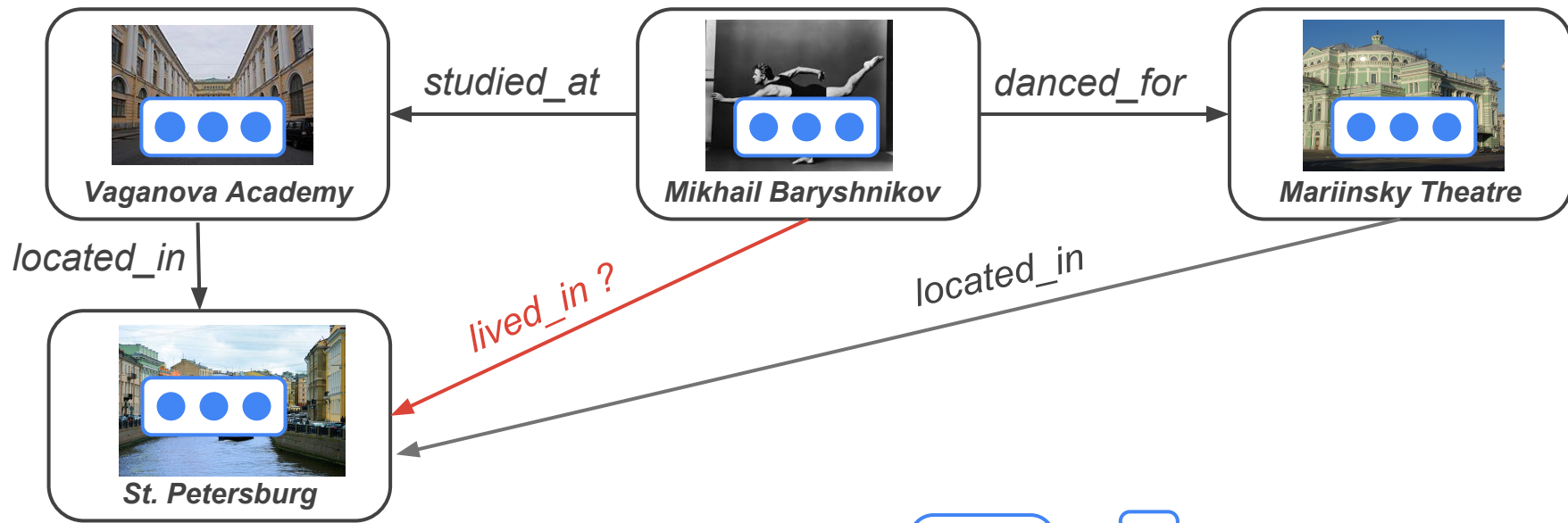
$$\begin{matrix} \text{Baryshnikov} & \times & \text{lived_in} & \times & \text{St. Petersburg} \end{matrix}$$

The diagram shows three blue boxes representing vectors:

- Baryshnikov**: A horizontal box with three blue dots.
- lived_in**: A square box with three blue dots.
- St. Petersburg**: A vertical box with three blue dots.

DistMult
(Yang et al., 2014)

KB Factorization



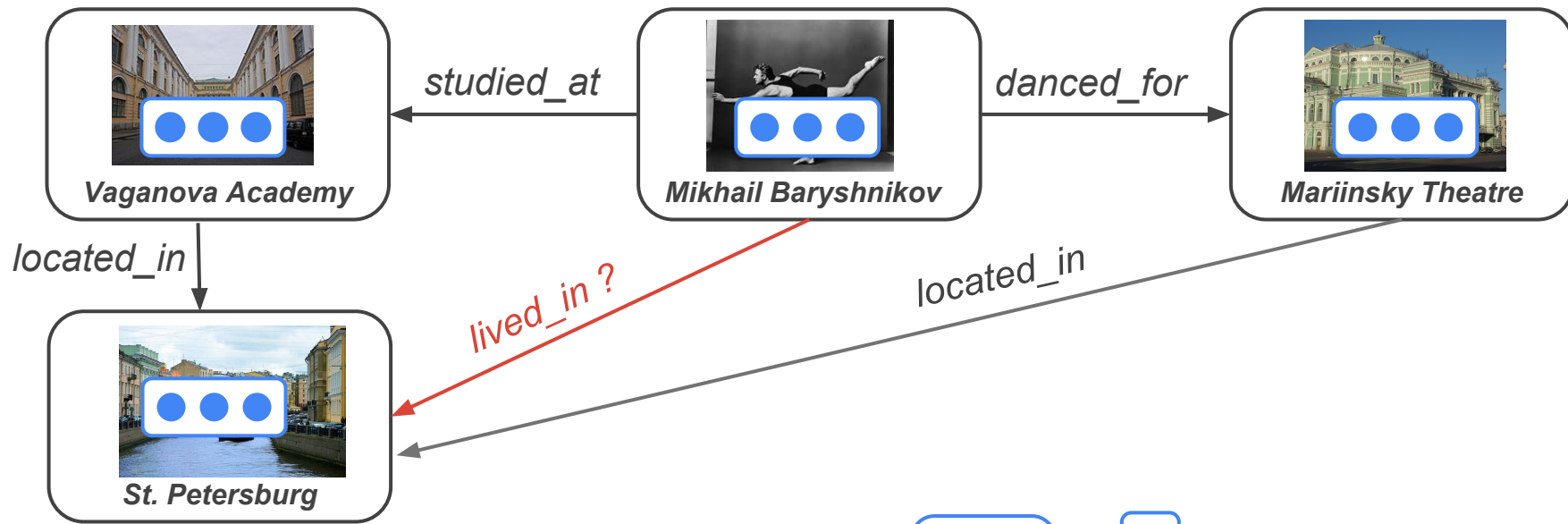
A scoring function is used to predict whether a relation holds:

$$\begin{matrix} \text{Baryshnikov} & \times & \text{lived_in} & \times & \text{St. Petersburg} \end{matrix}$$

DistMult
(Yang et al., 2014)

Relies on SGD to propagate information across the graph

Relational GCNs



A scoring function is used to predict whether a relation holds:

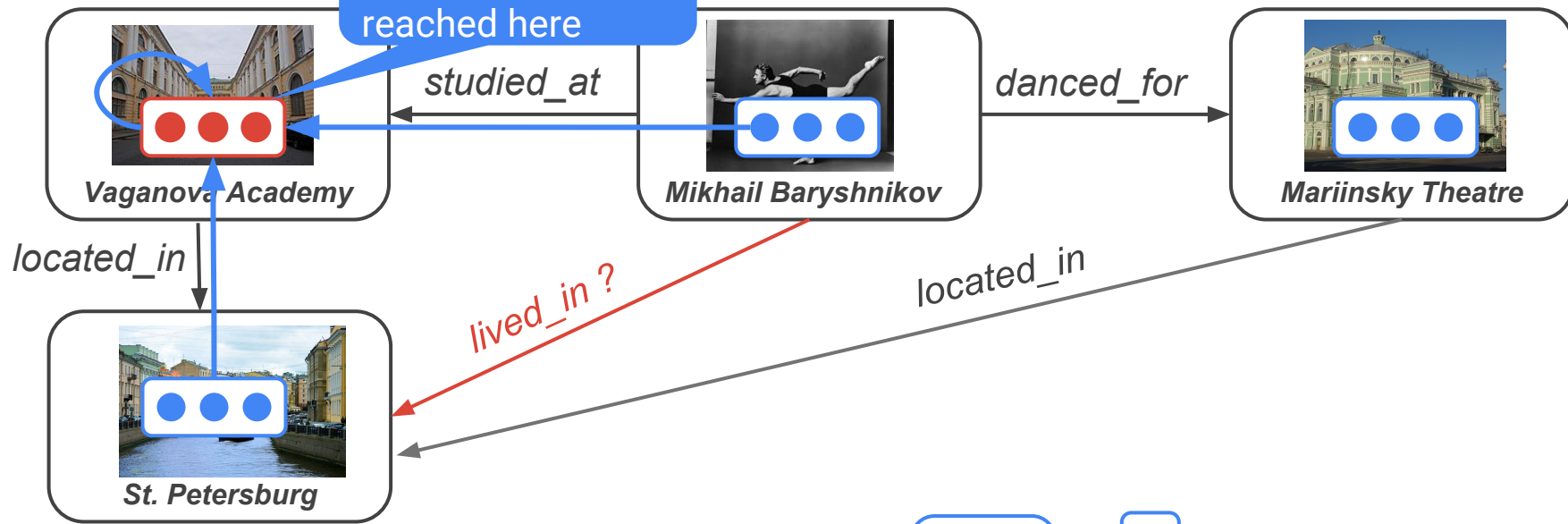
$$\begin{matrix} \text{Baryshnikov} \\ \text{[3 dots]} \end{matrix} \times \begin{matrix} \text{[3 dots]} \\ \text{lived_in} \end{matrix} \times \begin{matrix} \text{[3 dots]} \\ \text{St. Petersburg} \end{matrix}$$

DistMult
(Yang et al., 2014)

Use the same scoring function but with GCN node representations rather than parameter vectors

Relational GCNs

Info about
St. Petersburg
reached here



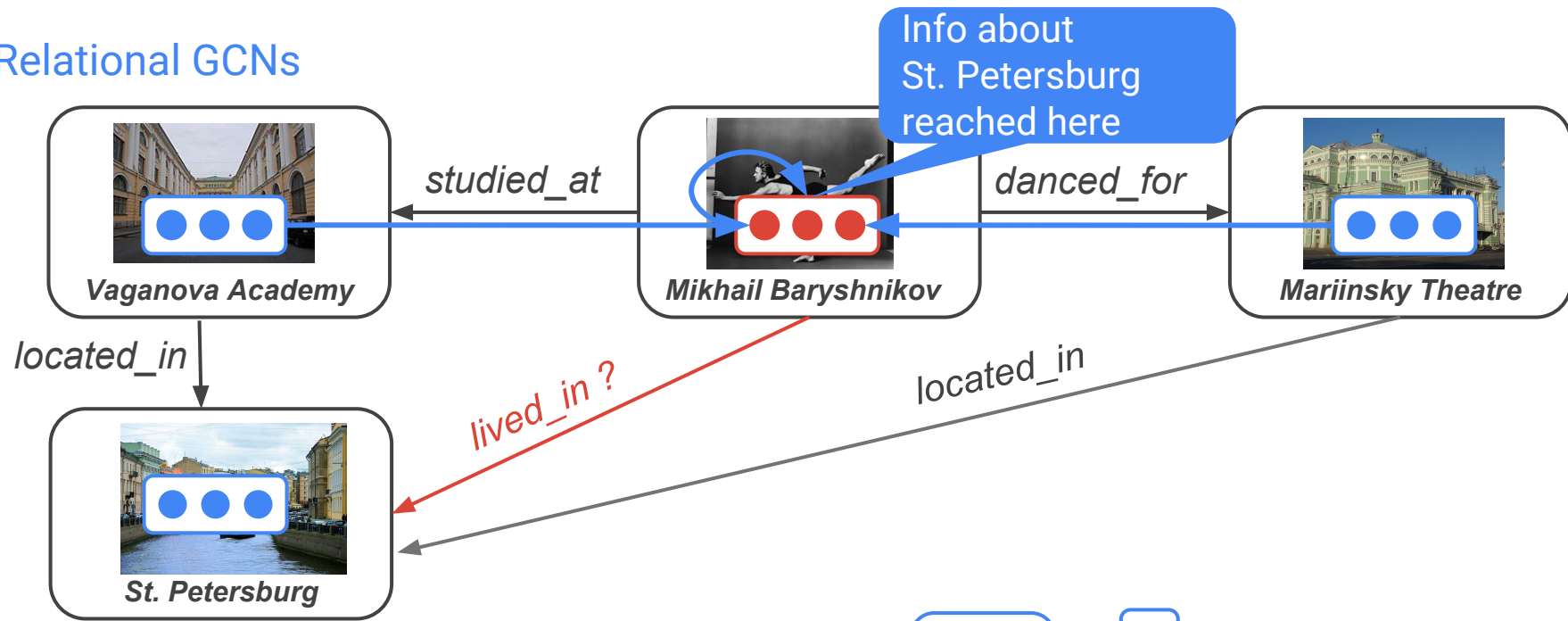
A scoring function is used to predict whether a relation holds:

$$\text{Baryshnikov} \times \text{lived_in} \times \text{St. Petersburg}$$

DistMult
(Yang et al., 2014)

Use the same scoring function but with GCN node representations rather than parameter vectors

Relational GCNs



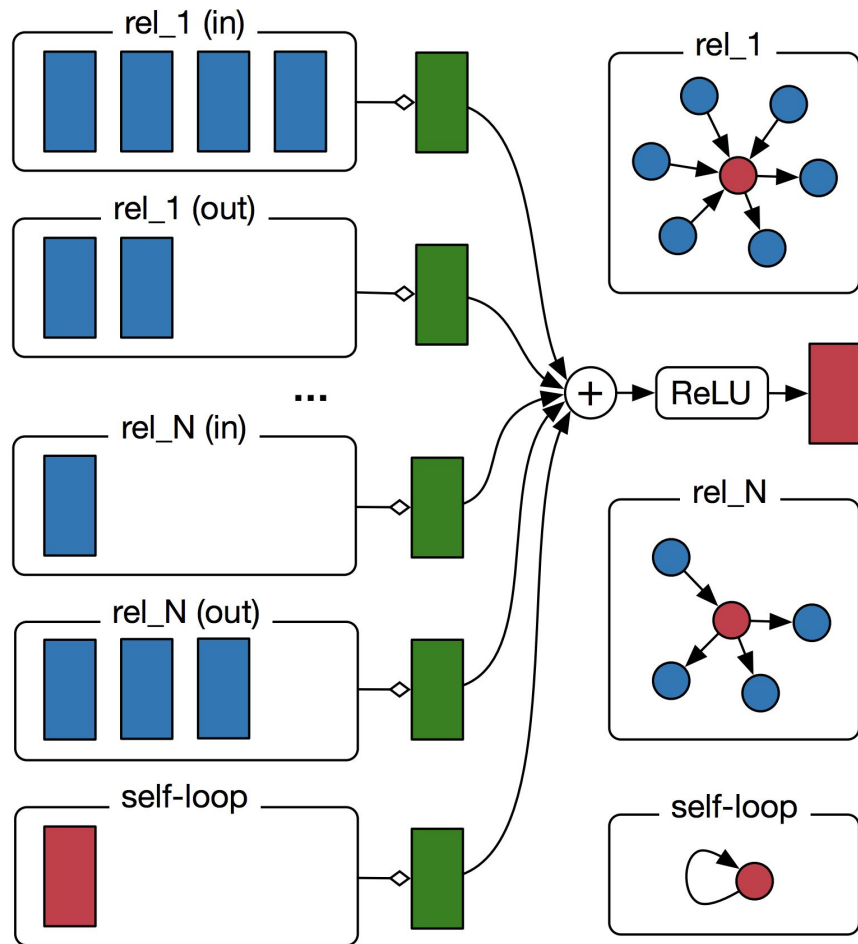
A scoring function is used to predict whether a relation holds:

$$\begin{matrix} \text{Baryshnikov} & \times & \text{lived_in} & \times & \text{St. Petersburg} \end{matrix}$$

DistMult
(Yang et al., 2014)

Use the same scoring function but with GCN node representations rather than parameter vectors

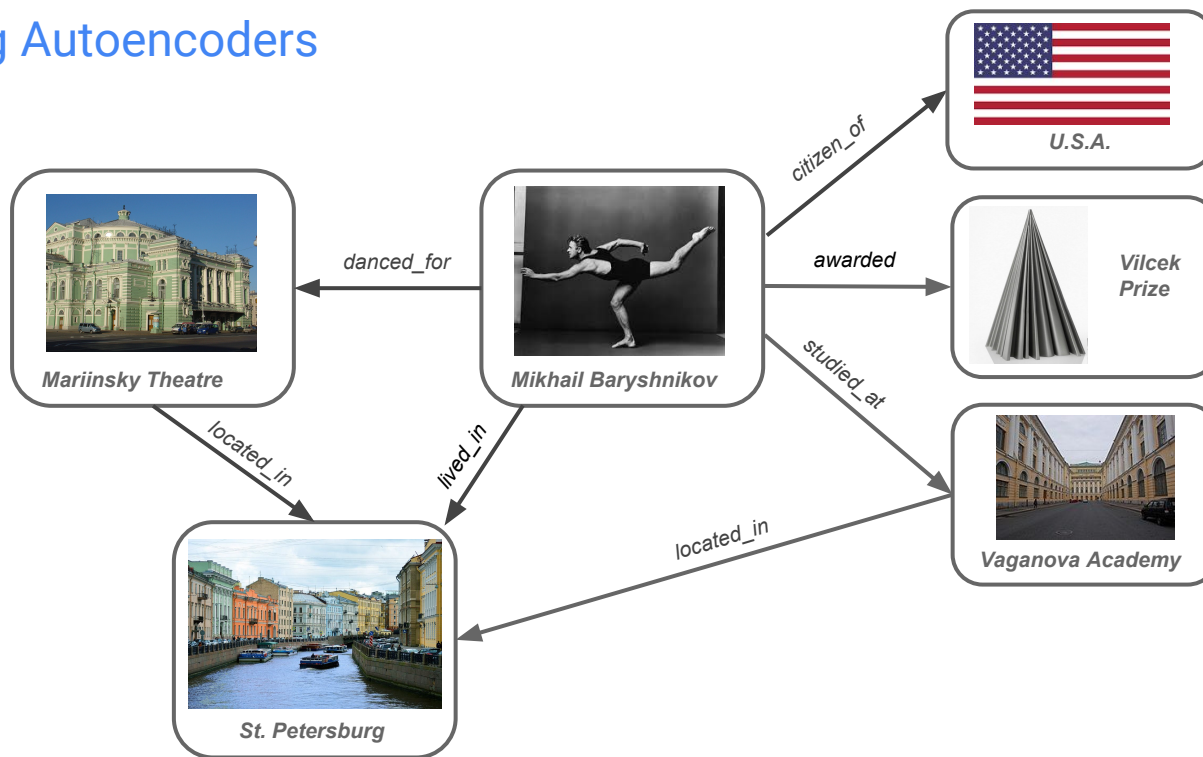
Relational GCNs



How do we train Relational GCNs?

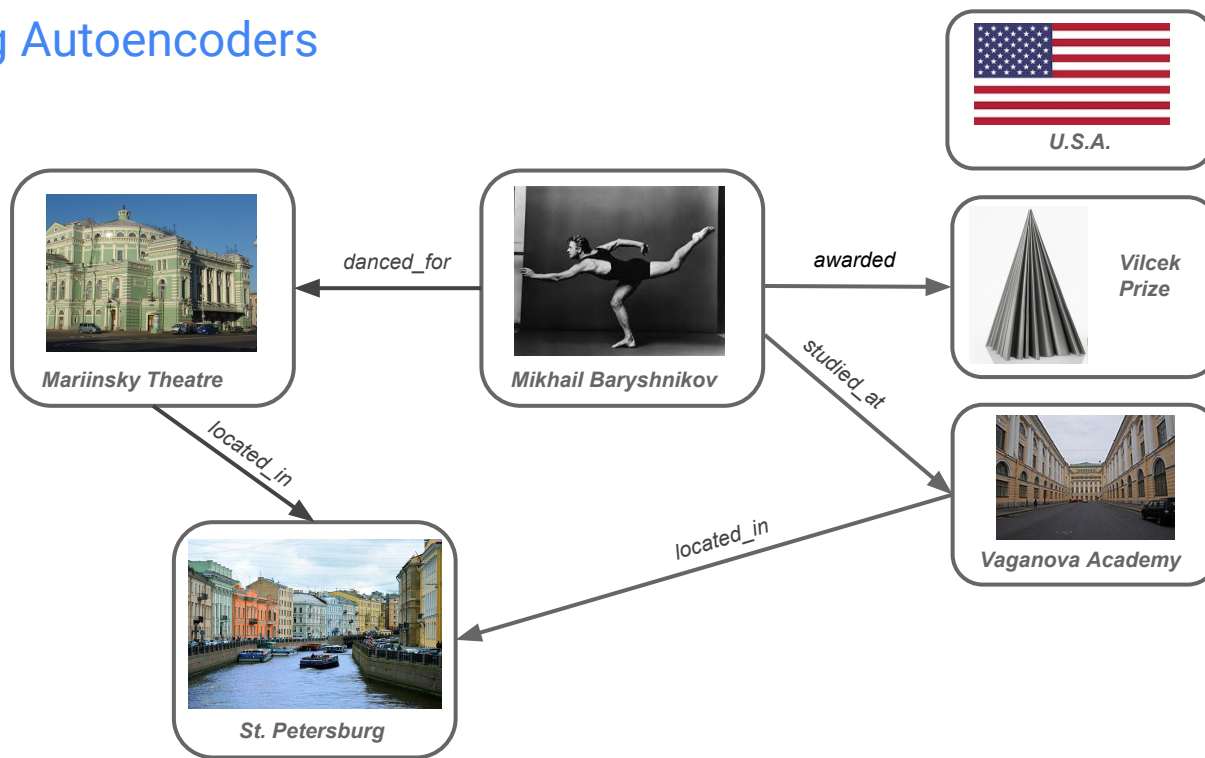
How do we compactly parameterize Relational GCNs?

GCN Denoising Autoencoders



Take the training graph

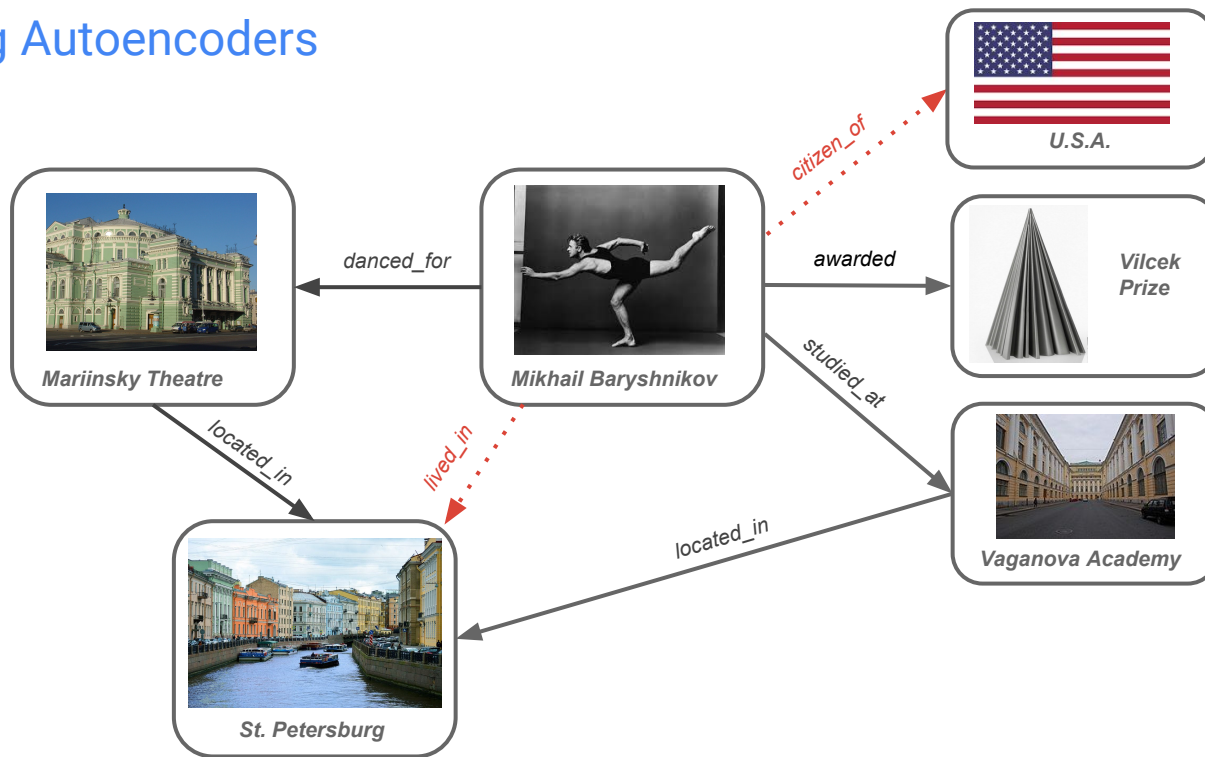
GCN Denoising Autoencoders



Produce a noisy version: drop some random edges

Use this graph for encoding nodes with GCNs

GCN Denoising Autoencoders

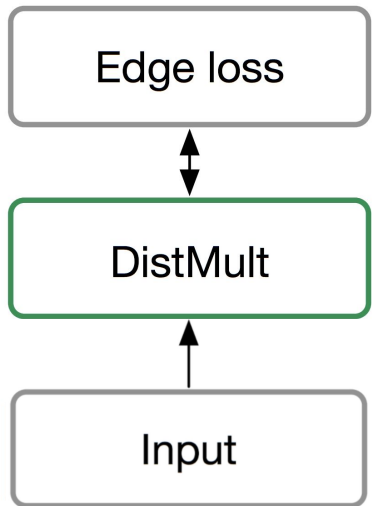


Force the model to reconstruct the original graph (including dropped edges)

(a ranking loss on edges)

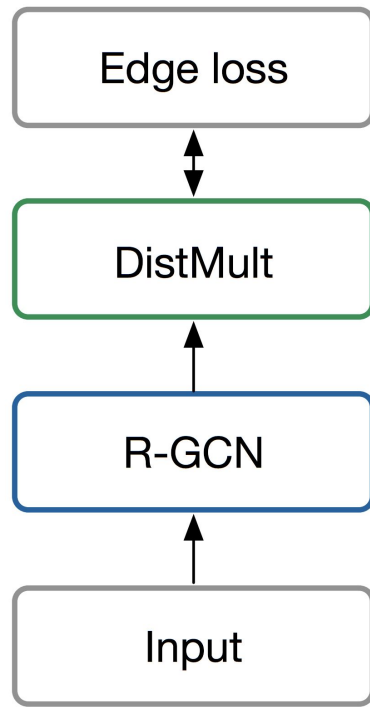
Training

Classic DistMult



(e.g., node embeddings)

Our R-GCN



Decoder

Encoder

(e.g., node embeddings)

GCN Autoencoders: Denoising vs Variational

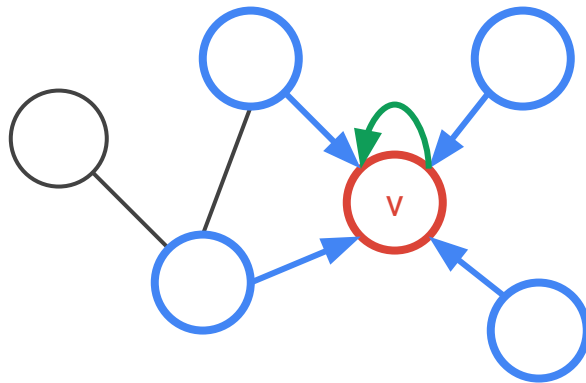
Instead of denoising AEs, we can use variational AEs to train R-GCNs

VAE R-GCN can be regarded as **an inference network** performing amortized variational inference

Intuition:

R-GCN AEs are amortized versions of factorization models

Relational GCN



$$\mathbf{h}_v = \text{ReLU}\left(\frac{1}{|\mathcal{N}(v)|} \sum_{u \in \mathcal{N}(v)} W_{r(u,v)} \mathbf{h}_u\right)$$

There are too many relations in realistic KBs, we cannot use full rank matrices W_r

Naive logic:

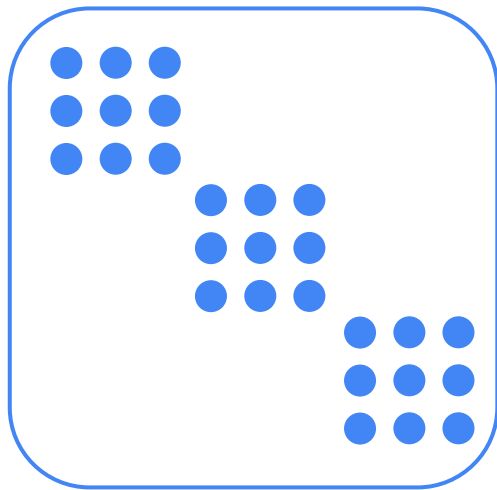
We score with a diagonal matrix (DistMul), let's use a diagonal one in GCN

$$W_r = \begin{matrix} \bullet & & \\ & \bullet & \\ & & \bullet \end{matrix}$$

Block diagonal assumption:

Latent features can be grouped into sets of tightly inter-related features, modeling dependencies across the sets is less important

$$W_r =$$



Basis / Dictionary learning:

Represent every KB relation as a linear combination of basis transformations

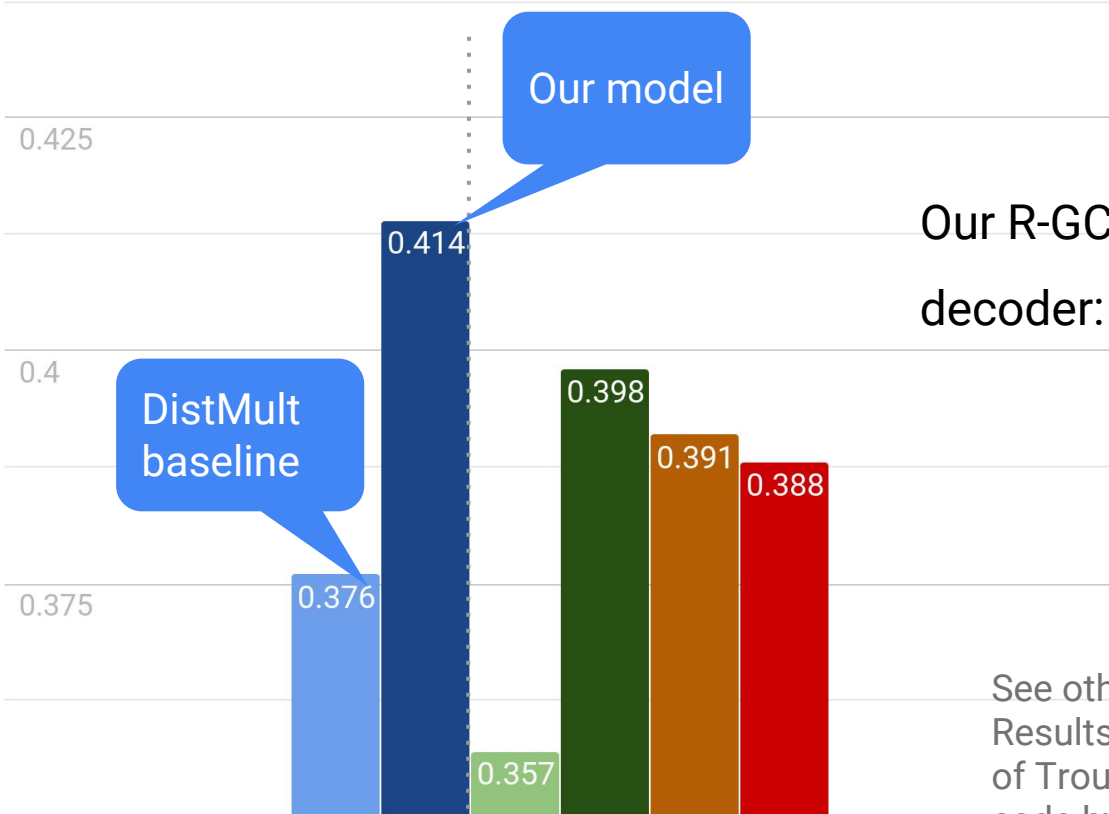
$$W_r = \sum_{b=1}^B a_{rb} V_b$$

Diagram illustrating the equation $W_r = \sum_{b=1}^B a_{rb} V_b$ with annotations:

- The term a_{rb} is labeled "coefficients" with an arrow pointing to it.
- The term V_b is labeled "basis transformations" with an arrow pointing to it.

Results on FB15k-237 (hits@10)

■ DistMult ■ R-GCN (block diagonal) ■ CP ■ TransE ■ HolE
■ ComplEX



Our R-GCN relies on DistMult in the decoder: DistMult is its natural baseline

See other results and metrics in the paper.
Results for ComplEX, TransE and HolE from code of Trouillon et al. (2016). Results for HolE using code by Nickel et al. (2015)

Relational GCNs

Fast and simple approach to Link Prediction

Captures multiple paths without the need to explicitly marginalize over them

Unlike factorizations, can be applied to **subgraphs unseen in training**

FUTURE WORK:

R-GCNs can be used in combination with **more powerful factorizations / decoders**

Objectives favouring **recovery of paths** rather than edges

Gates and memory may be effective

Extracting Semantic Relations

Semantic Role Labeling

Closely related to the relation extraction task

Discovering the predicate-argument structure of a sentence

Sequa makes and repairs jet engines

Semantic Role Labeling

Closely related to the relation extraction task

Discovering the predicate-argument structure of a sentence

- Discover predicates

Sequa

makes

and

repairs

jet

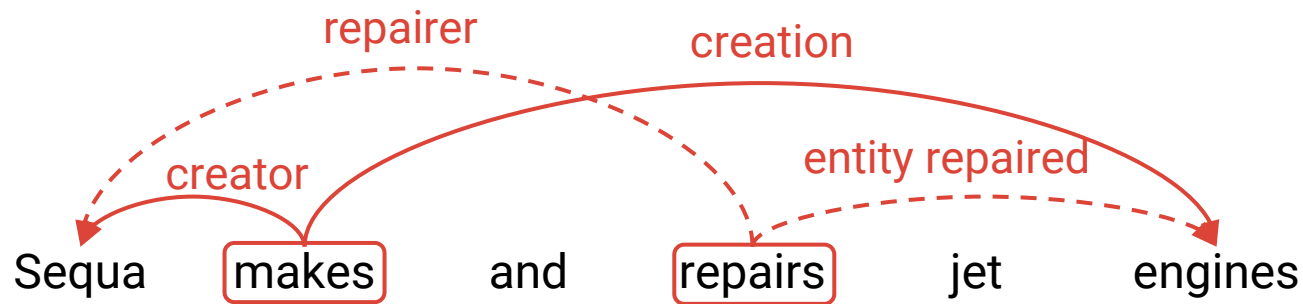
engines

Semantic Role Labeling

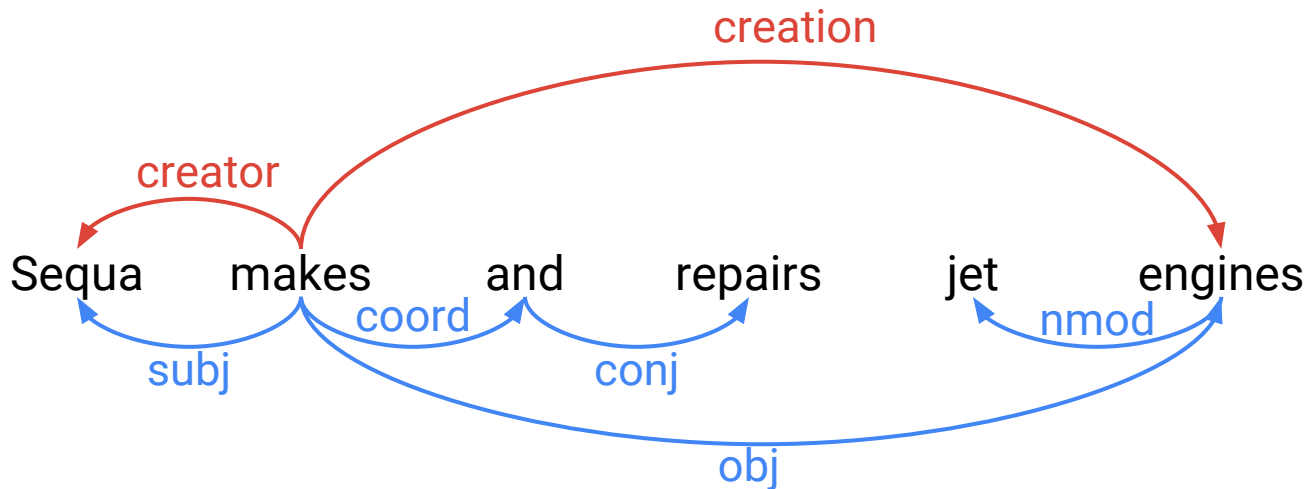
Closely related to the relation extraction task

Discovering the predicate-argument structure of a sentence

- Discover predicates
- Identify arguments and label them with their semantic roles

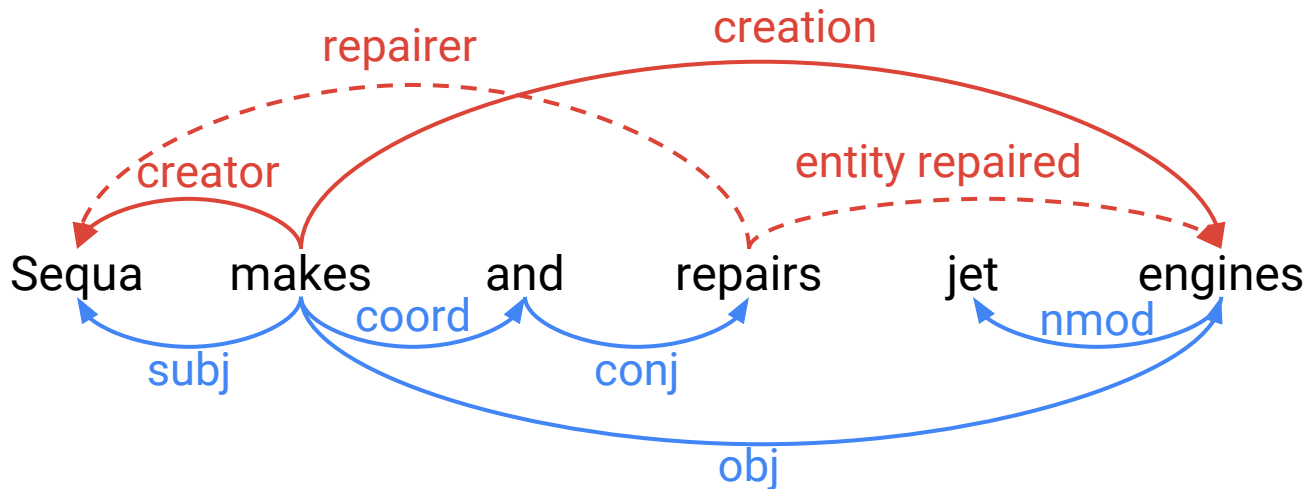


Syntax/semantics interaction



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

Syntax/semantics interaction

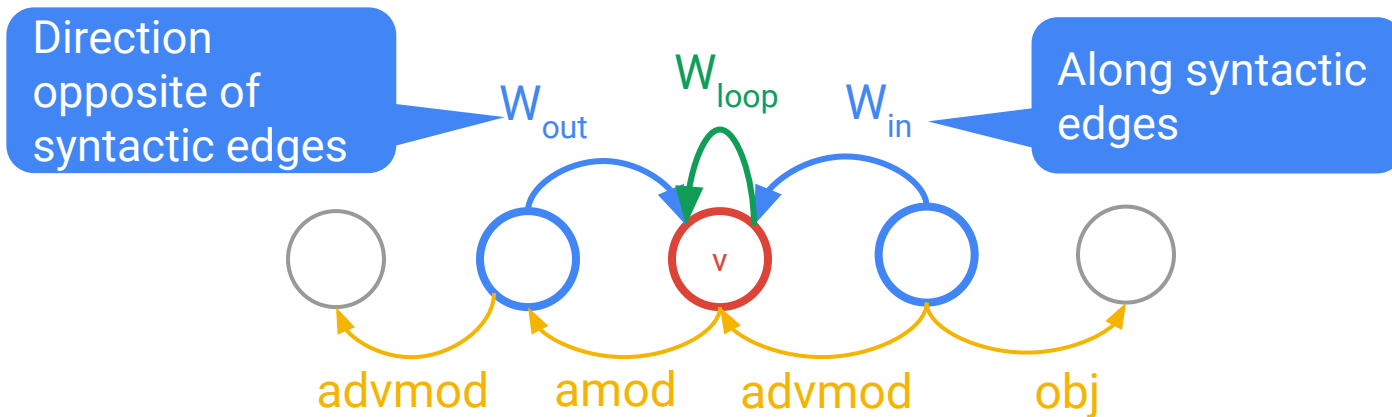


Some syntactic dependencies are **mirrored** in the semantic graph

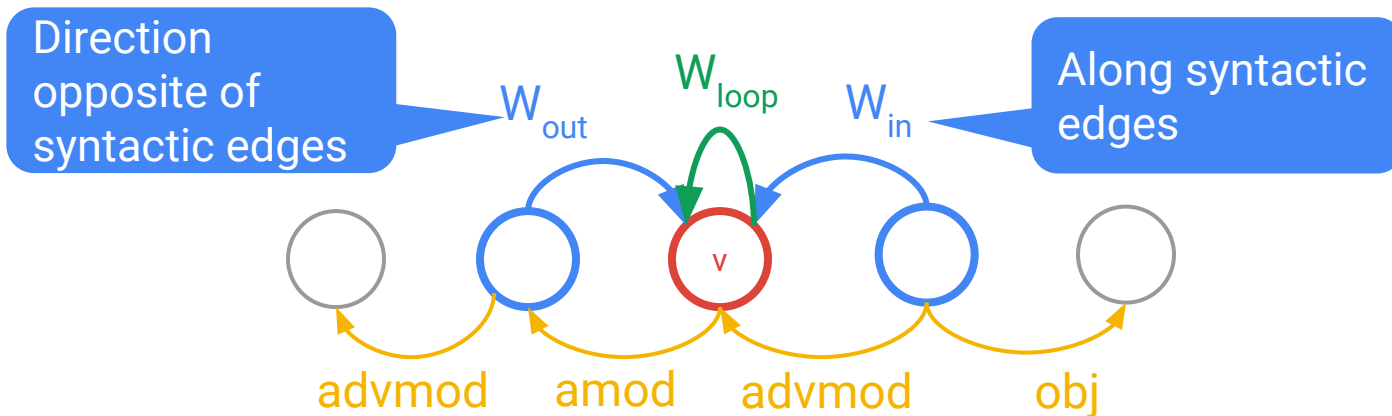
... **but not all of them** – the syntax-semantics interface is far from trivial

GCNs provide a flexible framework for capturing interactions between the graphs

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 + direction, e.g. $\mathbf{b}_{in-subj}$

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

Syntactic GCNs: edge-wise gating

We use parsers to predict syntax

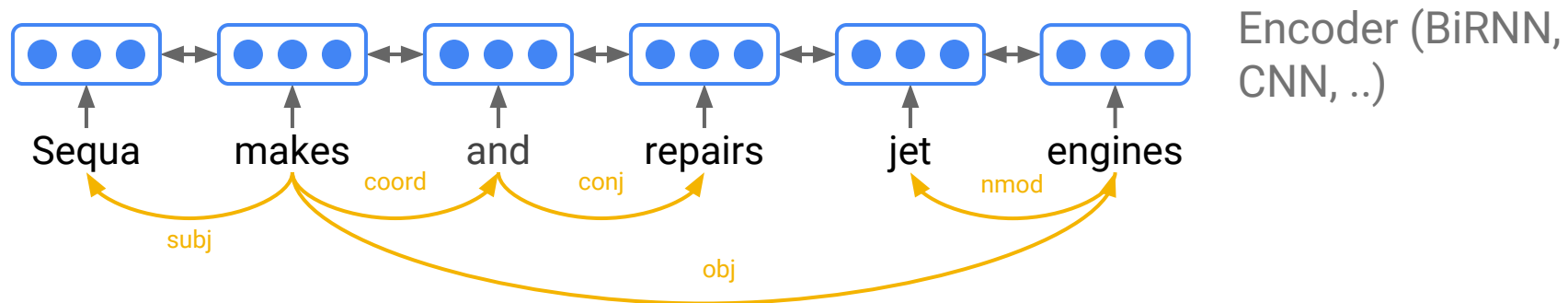
Not all edges are equally informative for the downstream task or reliable

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

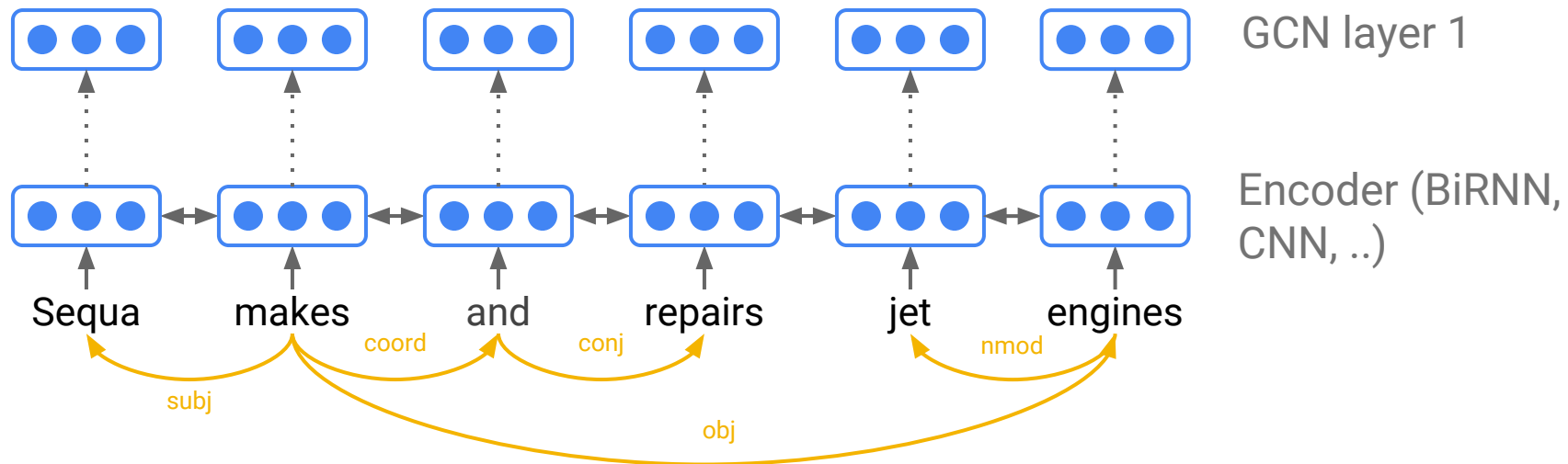
$$\mathbf{h}_v = \text{ReLU} \left(\sum_{u \in \mathcal{N}(v)} g_{u,v} \left(W_{\text{dir}(u,v)} \mathbf{h}_u + \mathbf{b}_{\text{lab}(u,v)} \right) \right)$$

The gate weights the message

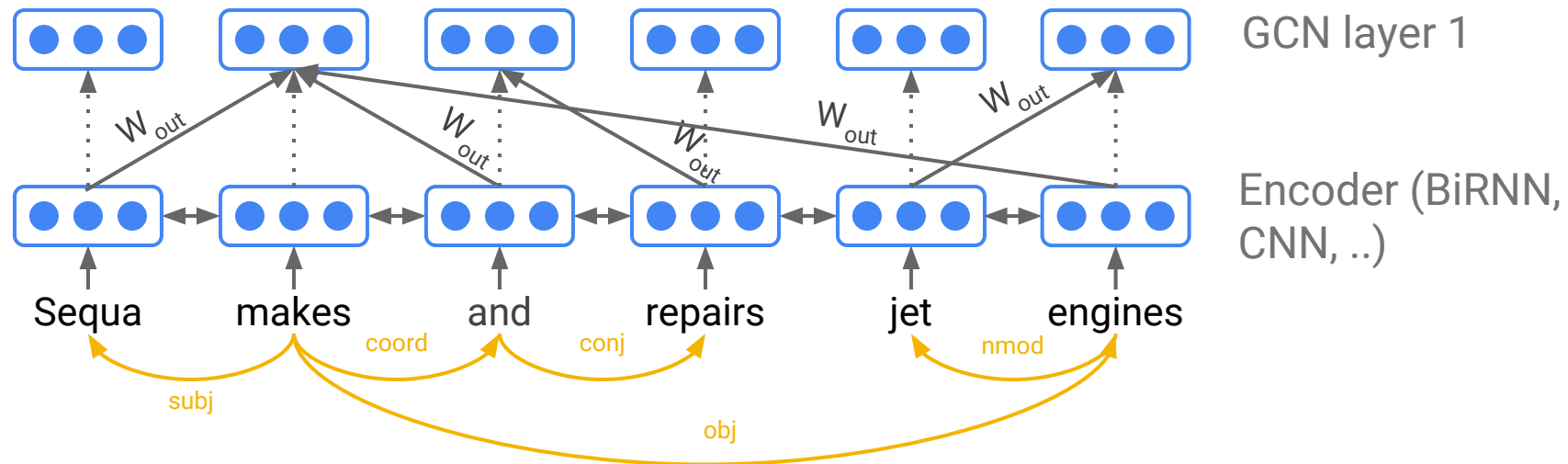
Graph Convolutional Encoders



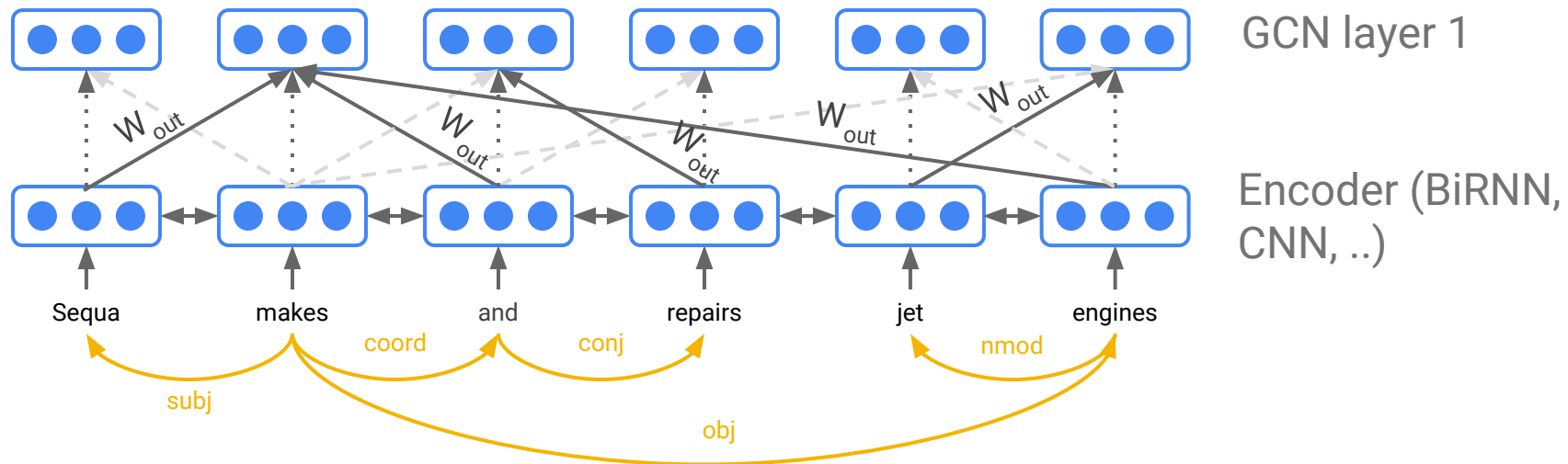
Graph Convolutional Encoders



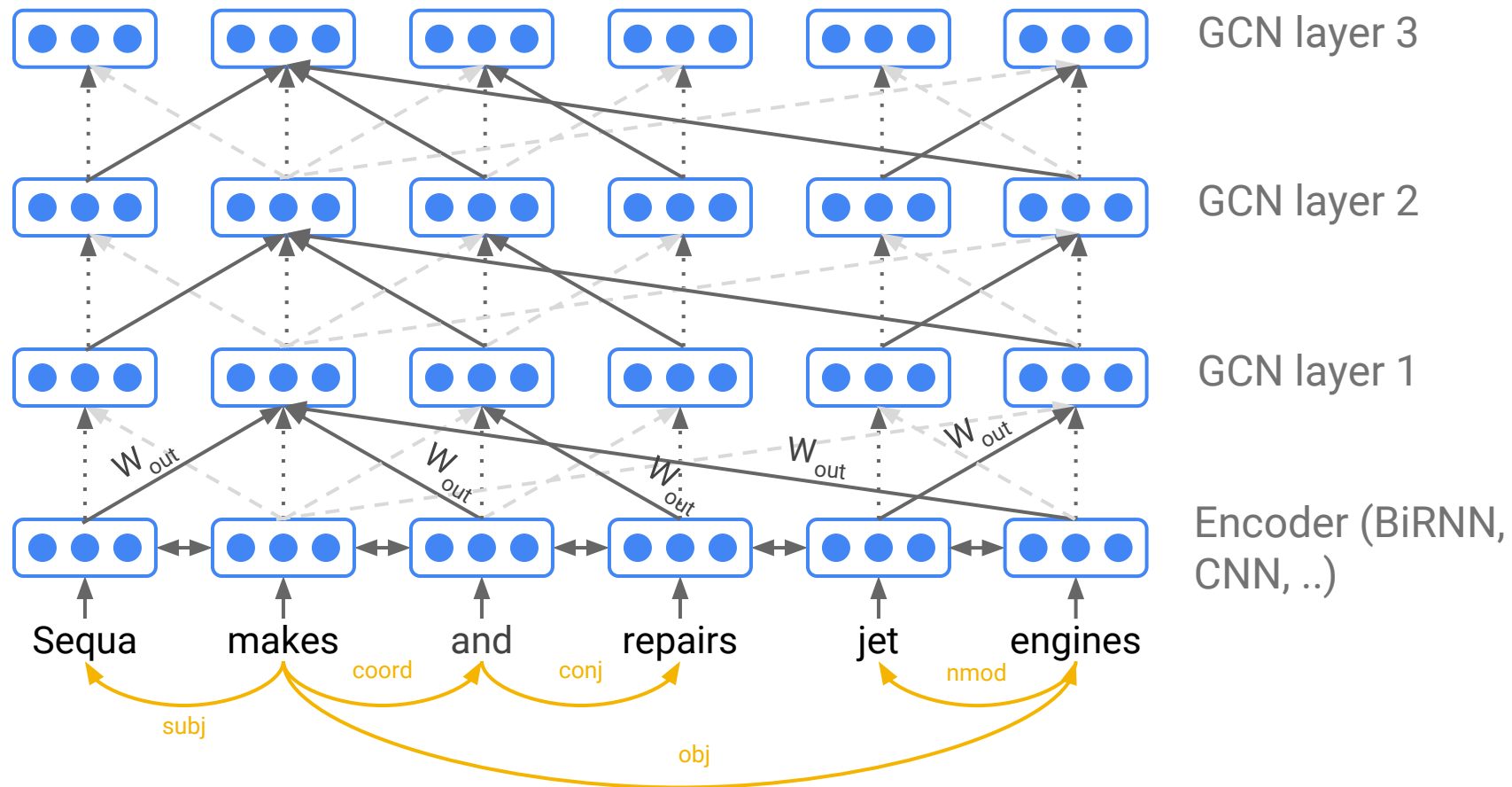
Graph Convolutional Encoders



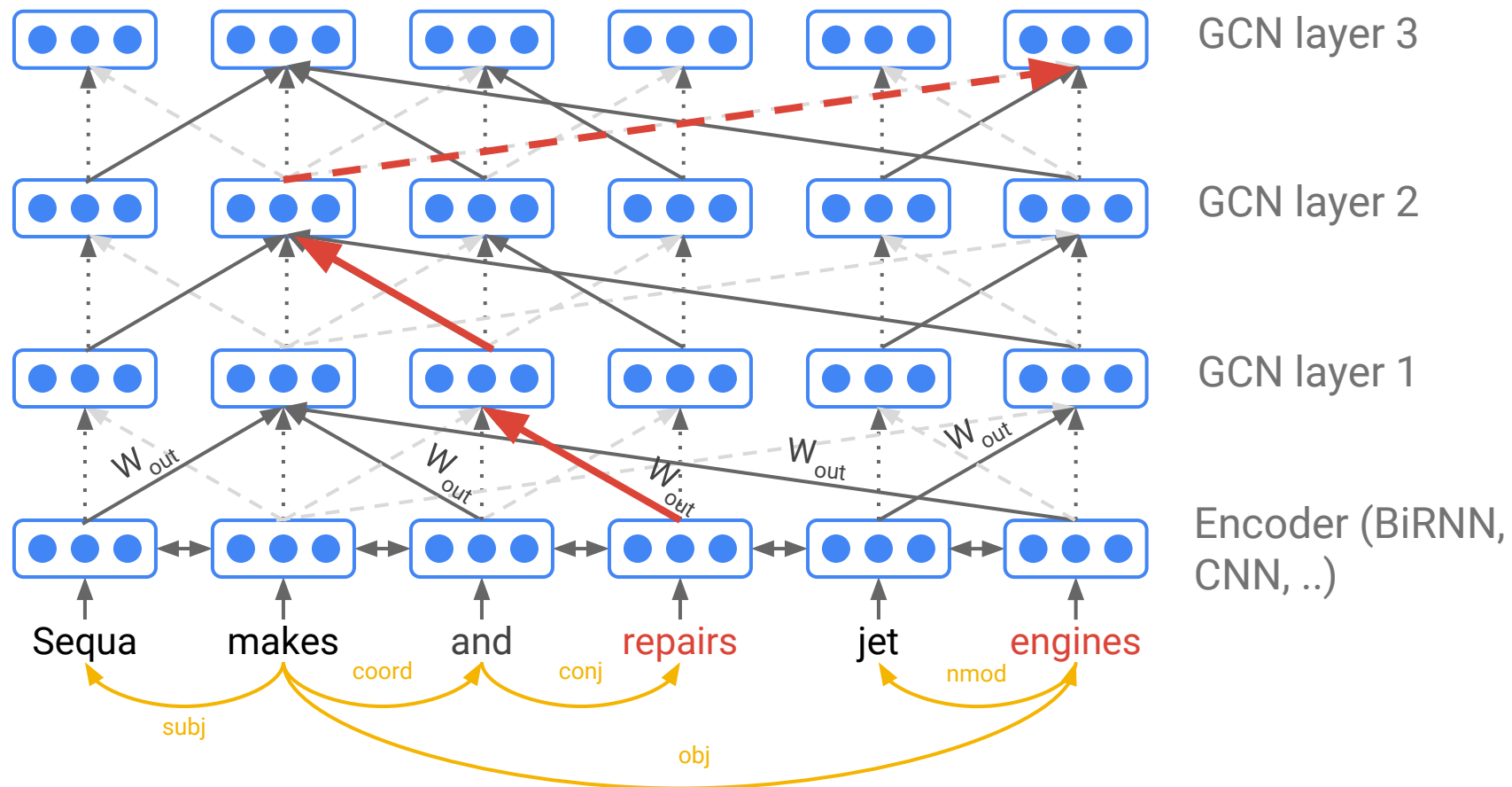
Graph Convolutional Encoders



Graph Convolutional Encoders

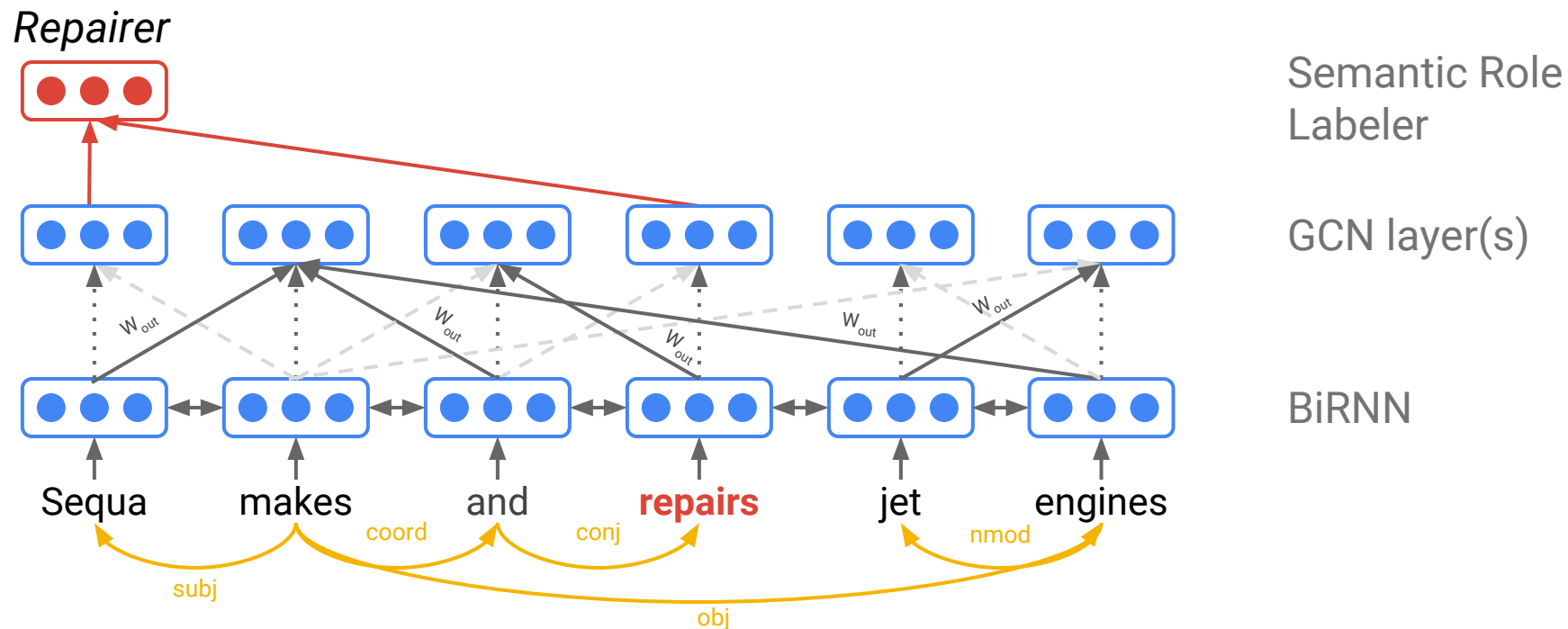


Graph Convolutional Encoders

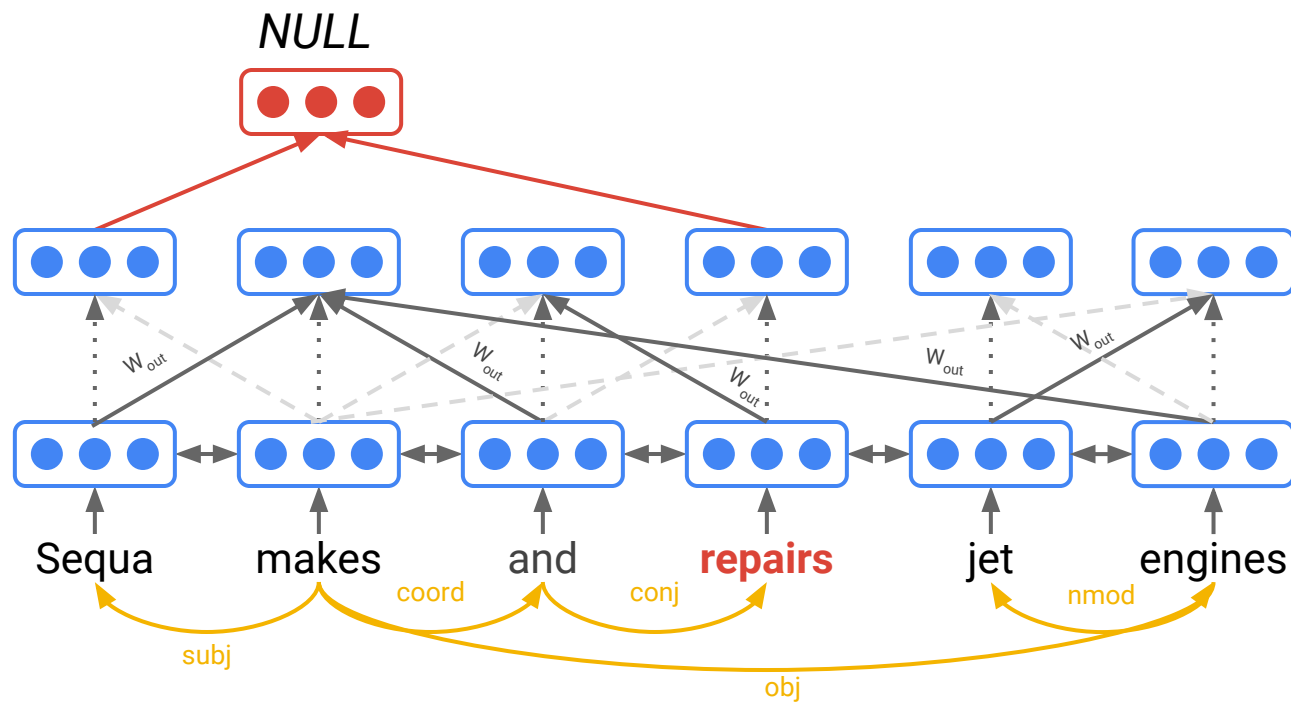


How do we construct a GCN-based semantic role labeler?

GCNs for Semantic Role Labeling



GCNs for Semantic Role Labeling

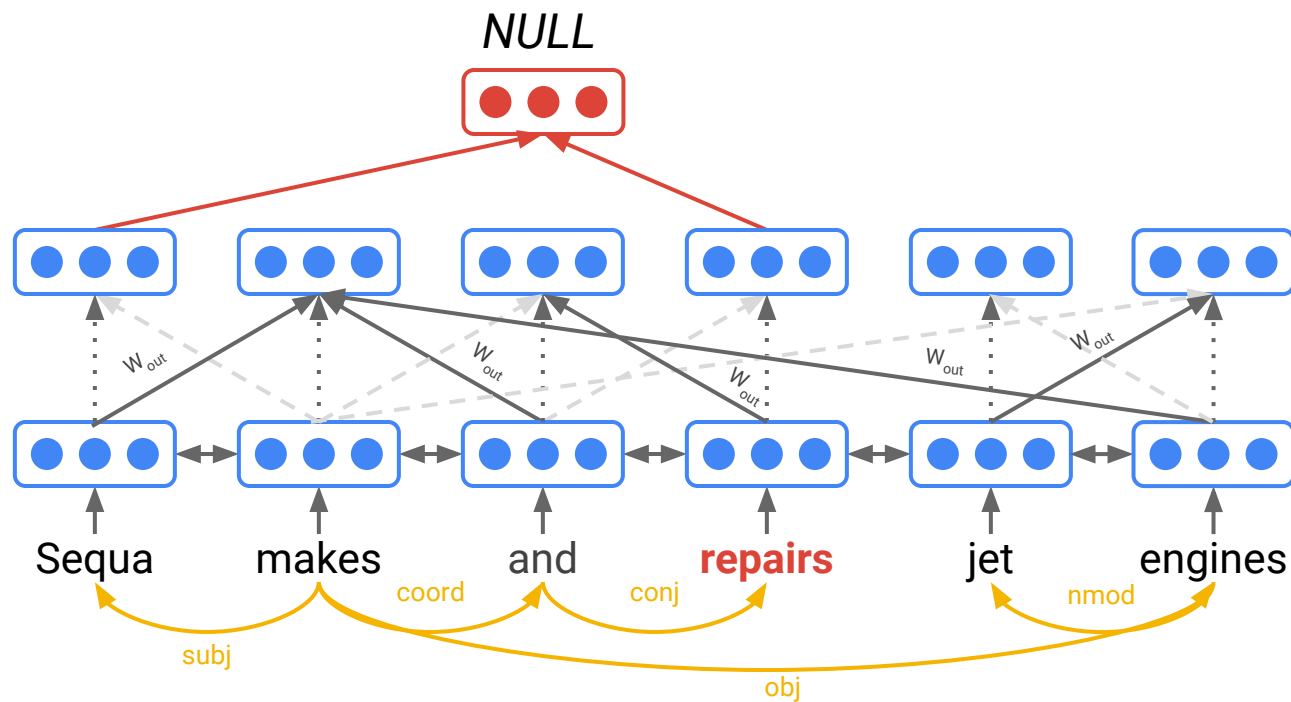


Semantic Role
Labeler

GCN layer(s)

BiRNN

GCNs for Semantic Role Labeling

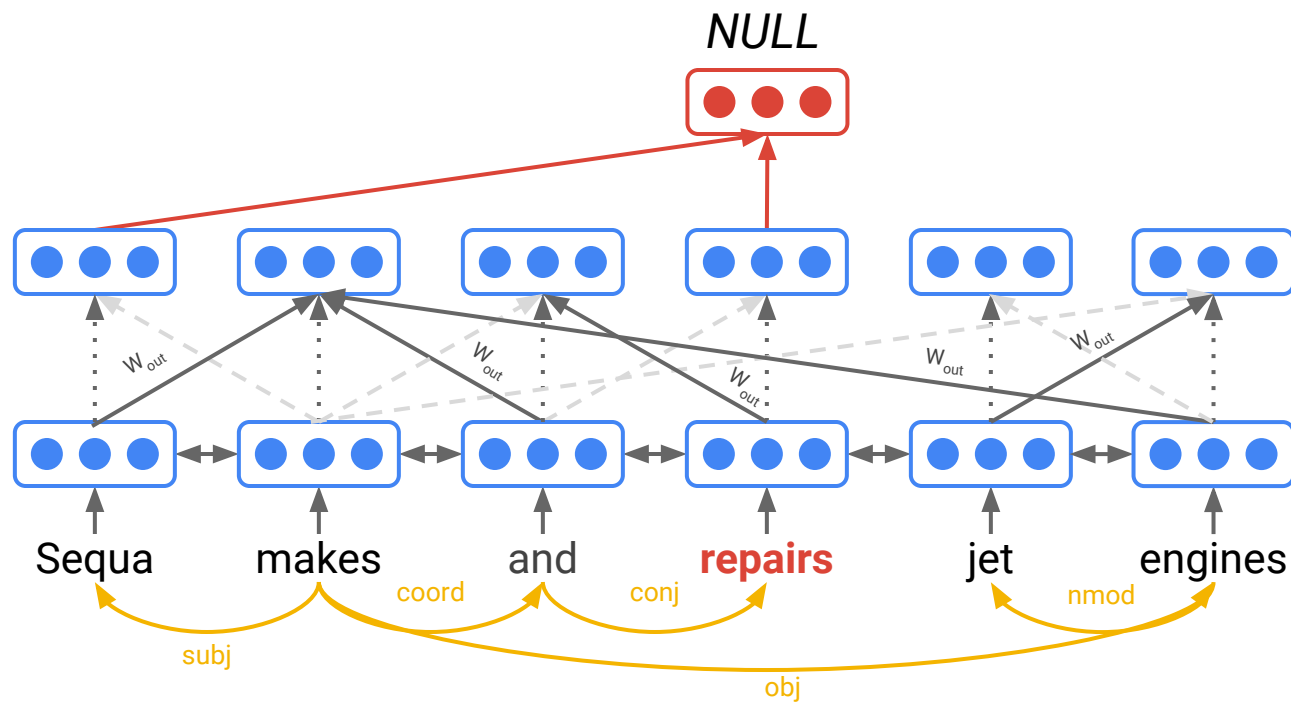


Semantic Role
Labeler

GCN layer(s)

BiRNN

GCNs for Semantic Role Labeling

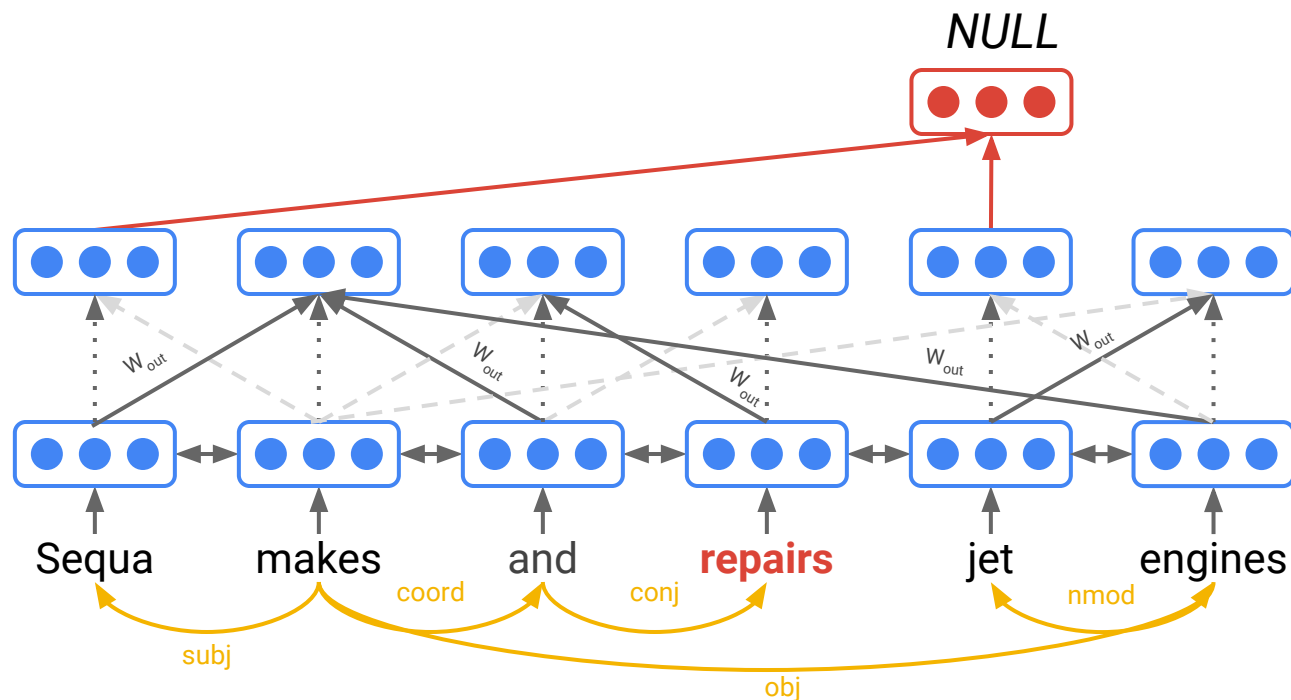


Semantic Role
Labeler

GCN layer(s)

BiRNN

GCNs for Semantic Role Labeling

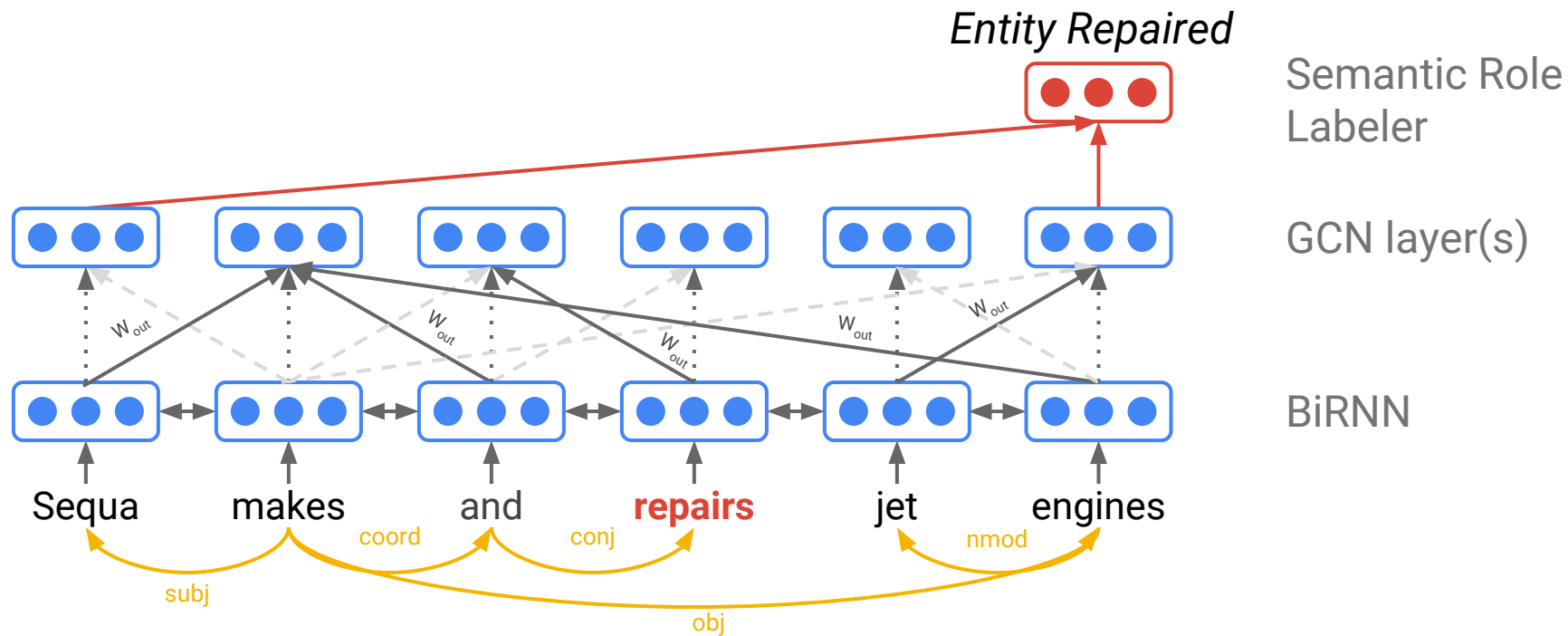


Semantic Role
Labeler

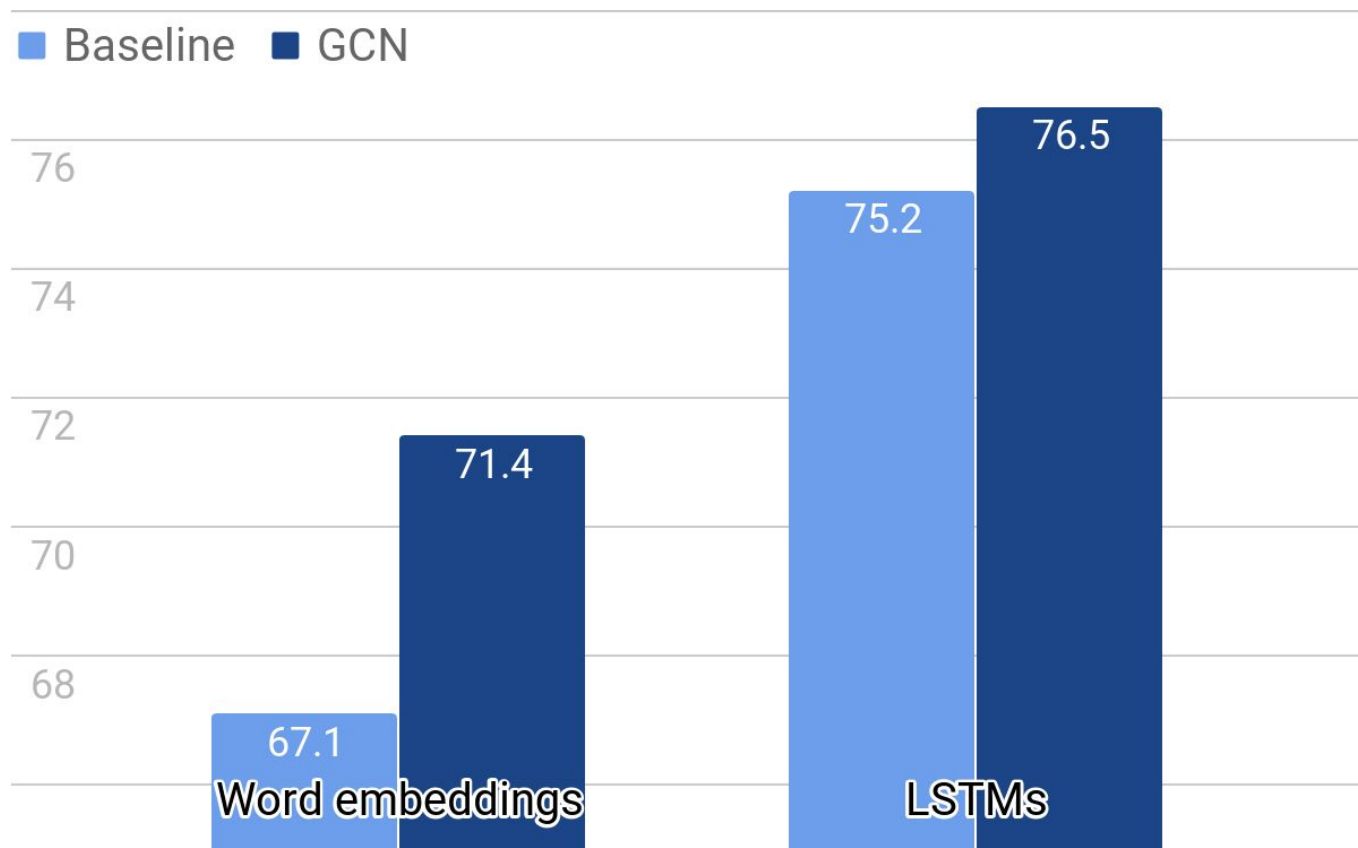
GCN layer(s)

BiRNN

GCNs for Semantic Role Labeling



Results (F1) on Chinese (CoNLL-2009, dev set)

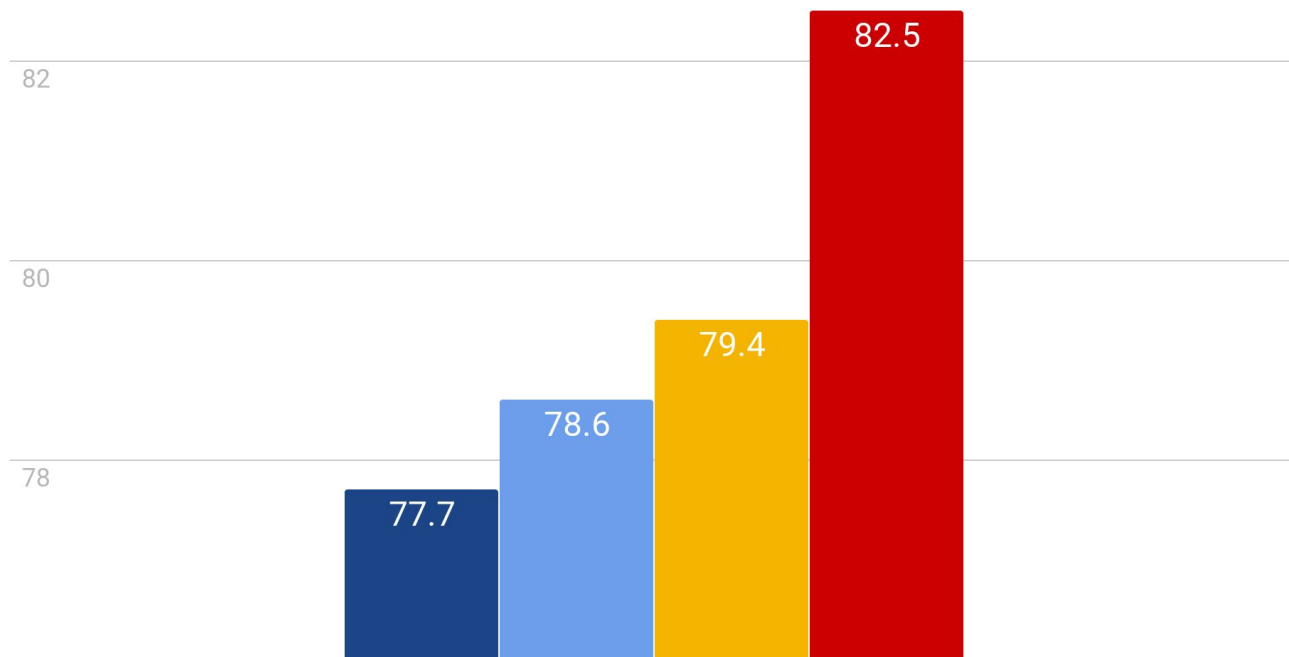


Marcheggiani & Titov (EMNLP, 2017)

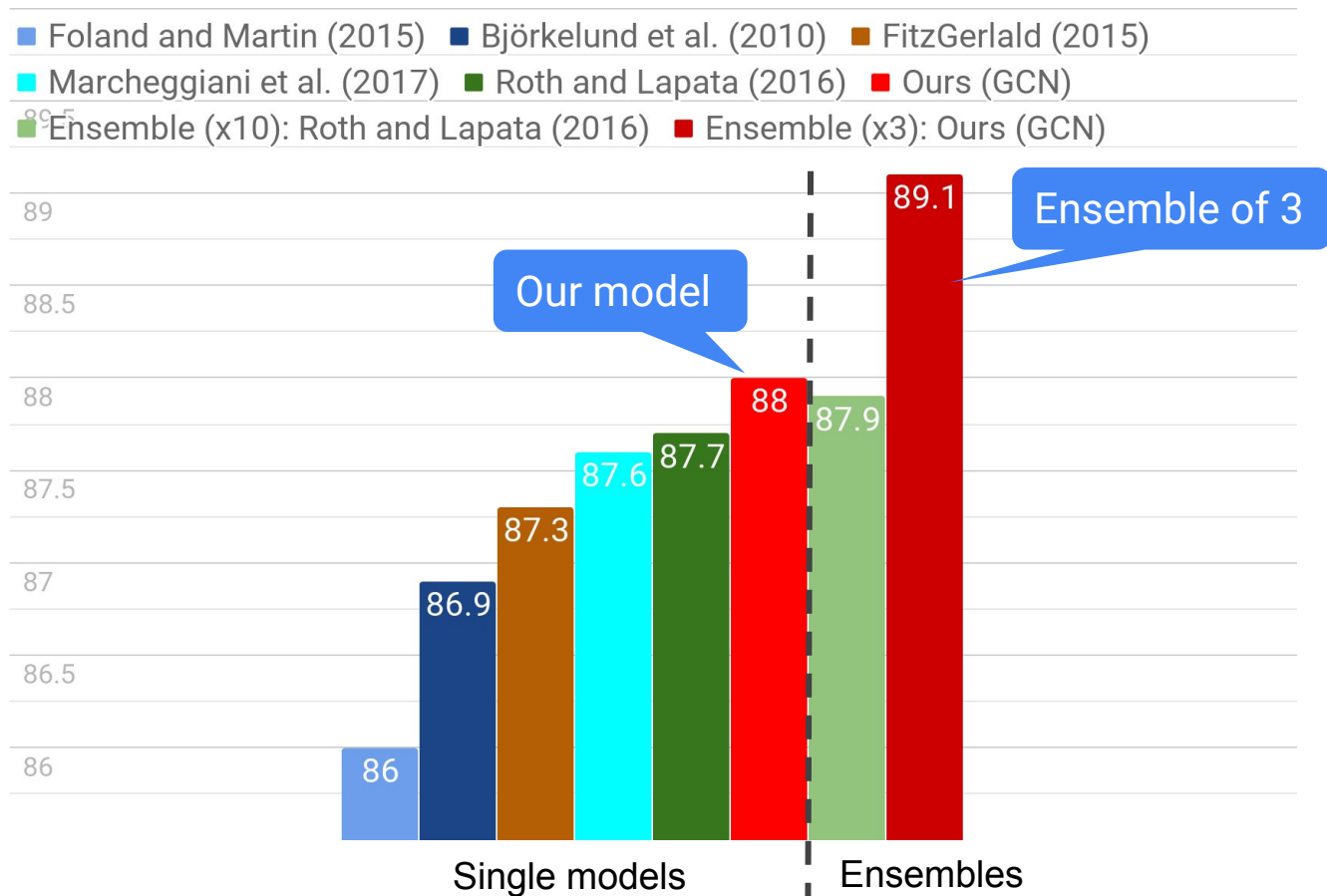
Predicate disambiguation is excluded from the F1 metric

Results (F1) on Chinese (CoNLL-2009, test set)

■ Zhao et al. (2009) ■ Björkelund et al. (2010) ■ Roth and Lapata (2016)
■ Ours (GCN)



Results (F1) on English (CoNLL-2009)



Flexibility of GCN encoders

Simple and fast approach to integrating linguistic structure into encoders

In principle we can exploit almost **any kind** of linguistic structure:

Semantic role labeling structure

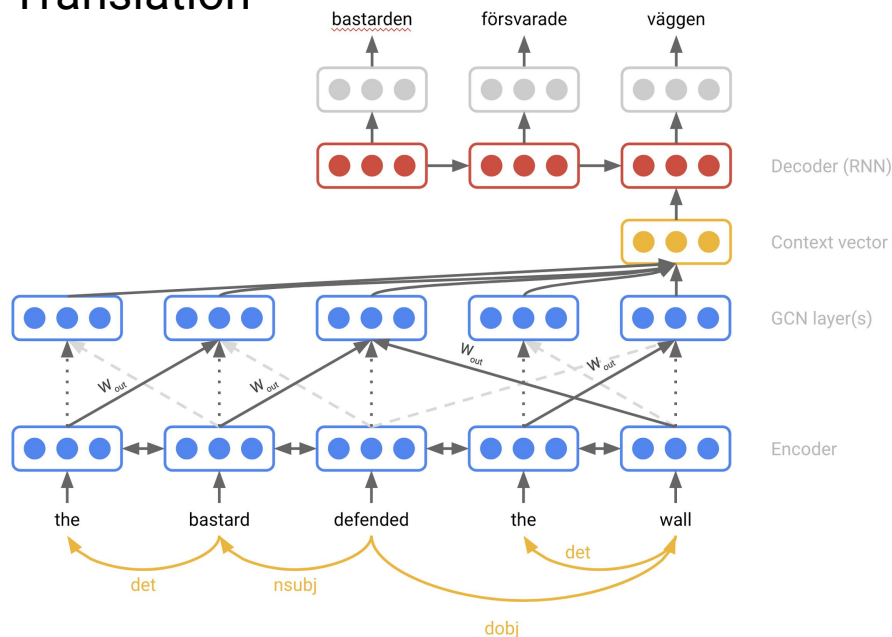
Co-reference chains

AMR semantic graphs

Their combination

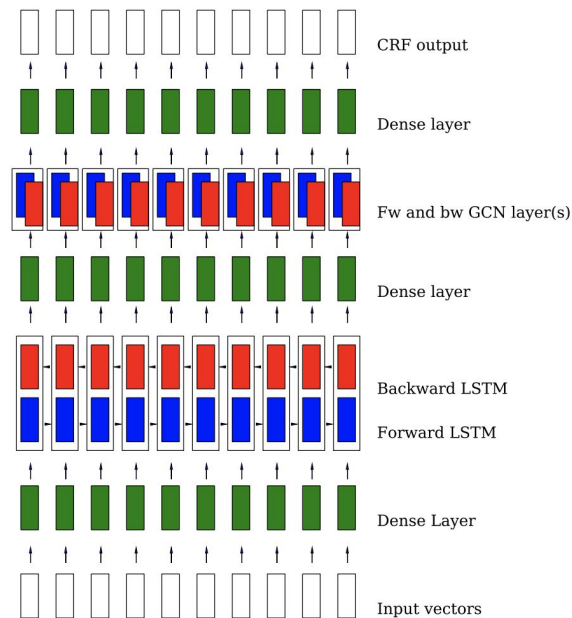
Other applications of syntactic GCN encoders

We also showed them effective as encoders in Neural Machine Translation



Bastings et al. (EMNLP, 2017)

Others recently applied them to NER



Cetoli et al. (arXiv:1709.10053)

Conclusions

GCNs are in subtasks of KBC (and in NLP beyond KBC):

- Semantic Roles: we proposed GCNs for encoding linguistic knowledge
- Link prediction: GCNs for link prediction (and entity classification) in multi-relational knowledge bases

Code available

We are hiring! (PhD students / postdocs)



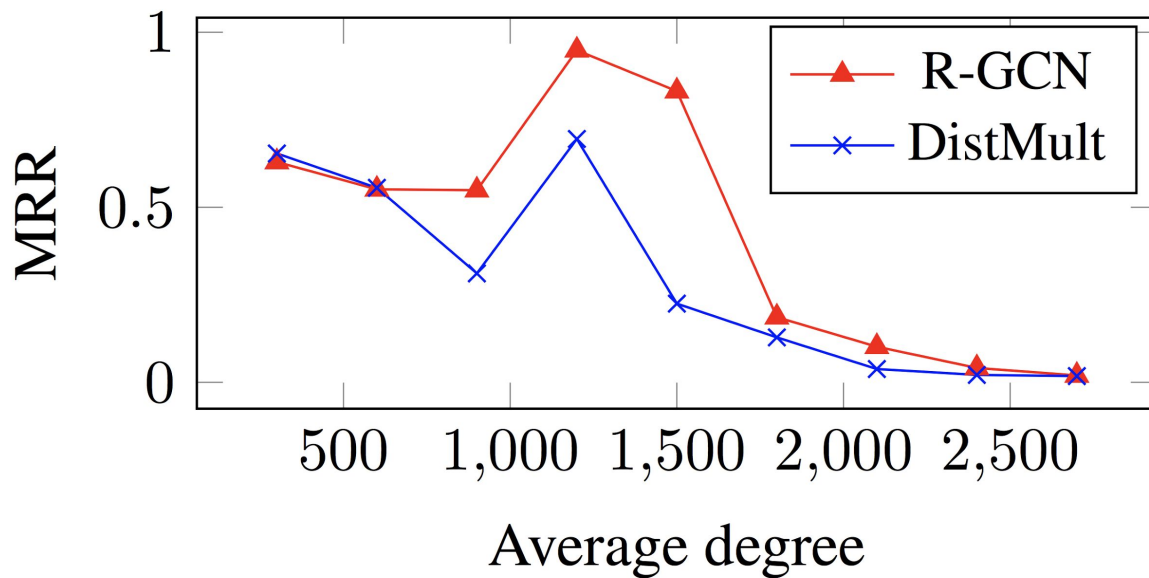
European
Research
Council



Nederlandse Organisatie voor
Wetenschappelijk Onderzoek



Analysis / Discussion



- Improvement across the board, especially in the middle of the range

Effect of Distance between Argument and Predicate (English)

