

Embedding Techniques

MATRICES, TENSORS, AND NEURAL NETWORKS



Graphical Models: Downsides

Limitation to Logical Relations

- Representation restricted by manual design
 - Clustering? Assymmetric implications?
 - Information flows through these relations
- Difficult to generalize to unseen entities/relations

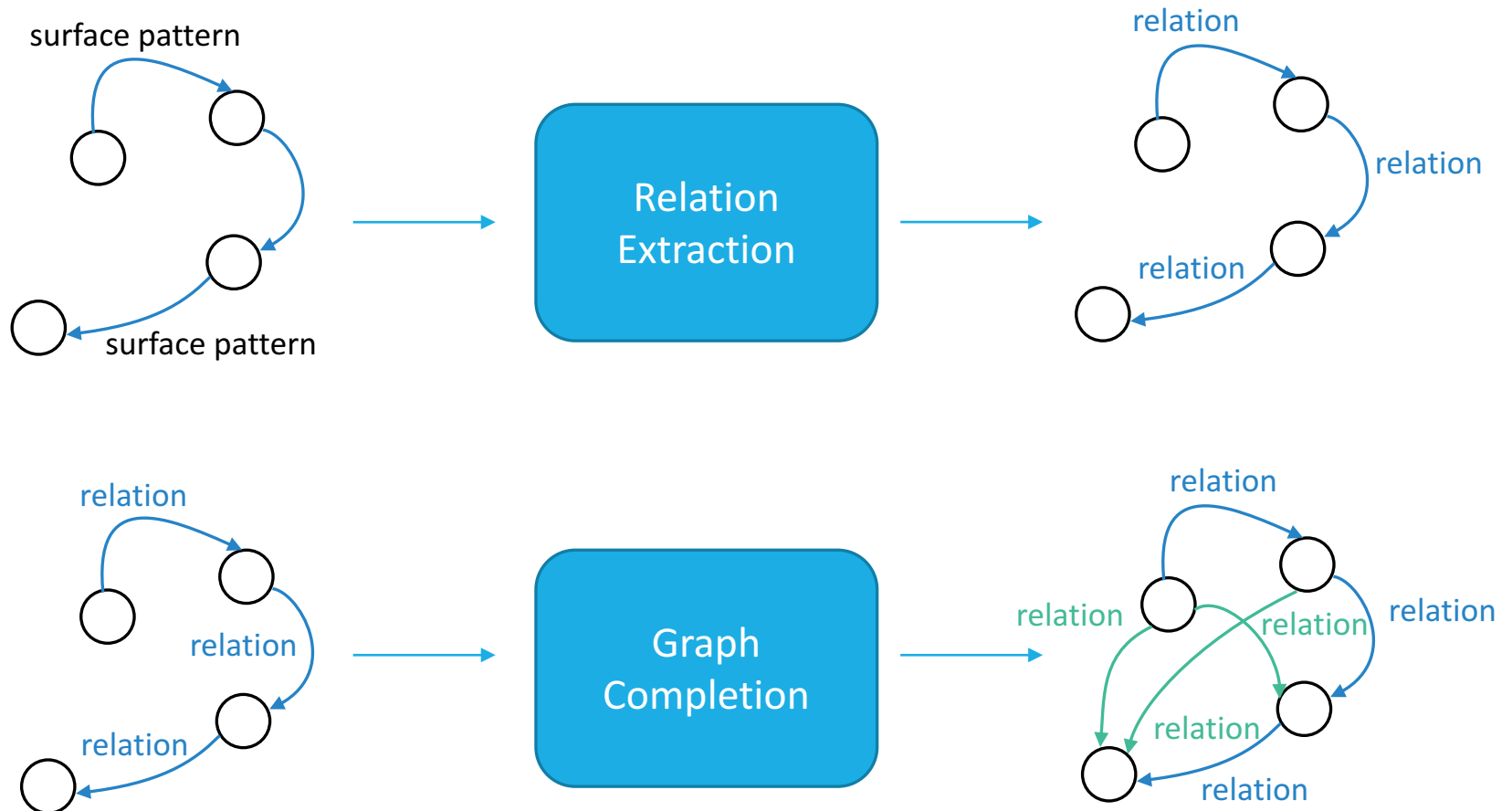
Computational Complexity of Algorithms

- Learning is NP-Hard, difficult to approximate
- Query-time inference is also NP-Hard
- Not easy to parallelize, or use GPUs
- Scalability is badly affected by representation

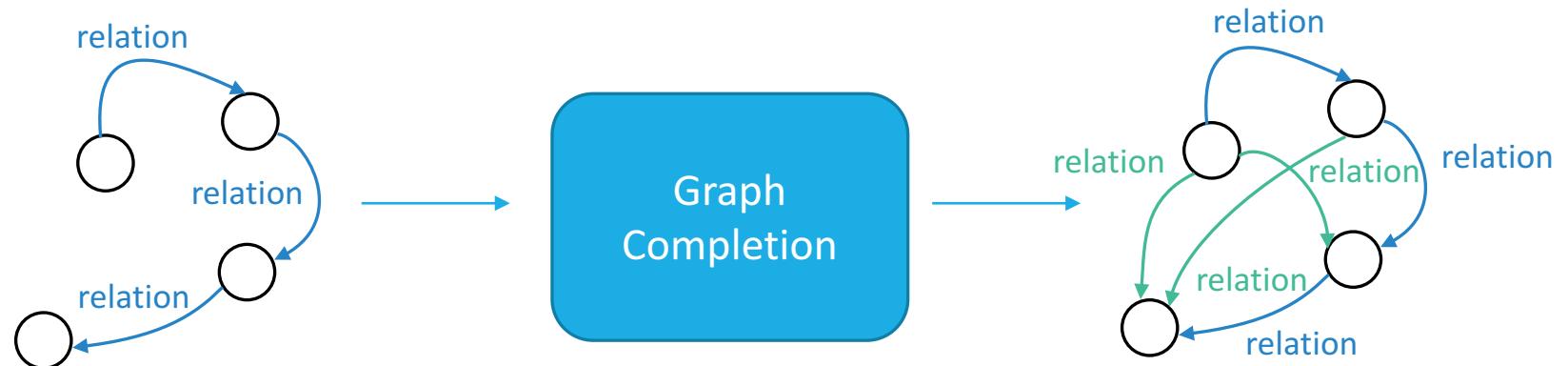
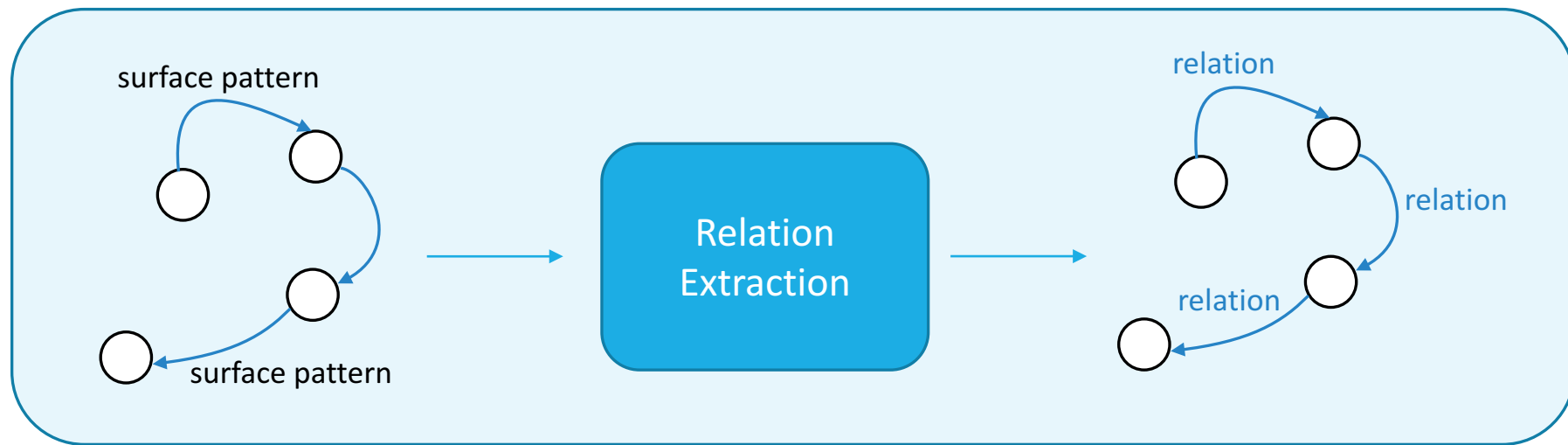
Embeddings

- Everything as dense vectors
 - Captures many relations
 - Learned from data
-
- Learning using stochastic gradient, back-propagation
 - Querying is often cheap
 - GPU-parallelism friendly

Two Related Tasks



Two Related Tasks



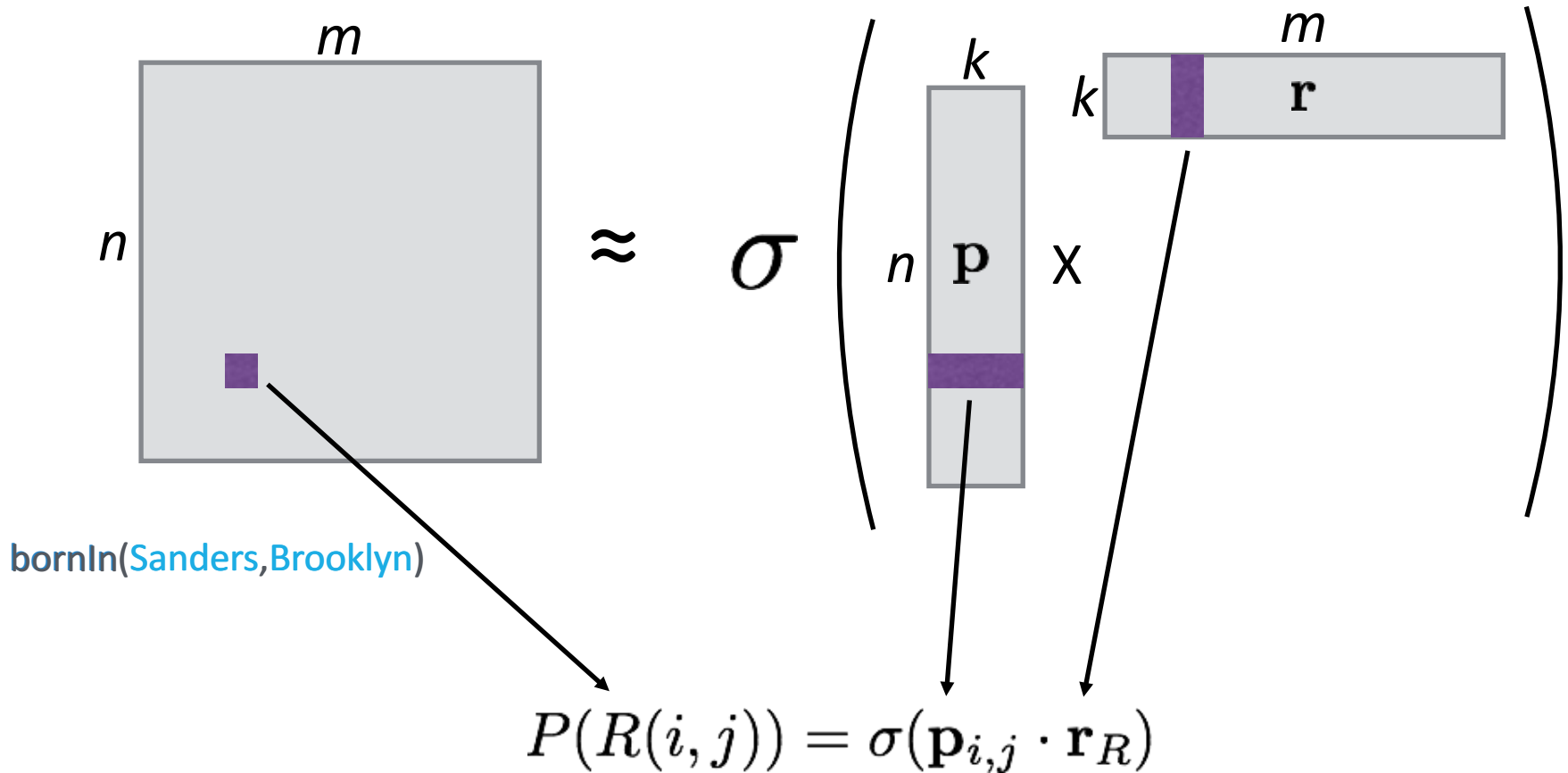
Relation Extraction as a Matrix

Sanders was born in Brooklyn, to Dorothy and Eli Sanders.

Entity Pairs

	<i>was born in</i> $\leftarrow n\text{subjpas-born} \leftarrow n\text{mod:in-}$	<i>was born to</i>	<i>and</i>	<i>birthplace(x,y)</i>	<i>spouse(x,y)</i>
Bernie Sanders, Brooklyn	1			?	
Bernie Sanders, Dorothy Sanders		1			
Bernie Sanders, Eli Sanders		1			
Dorothy Sanders, Eli Sanders			1		?
Barack Obama, Hawaii	1			1	
Barack Obama, Michelle Obama			1		1

Matrix Factorization



Training

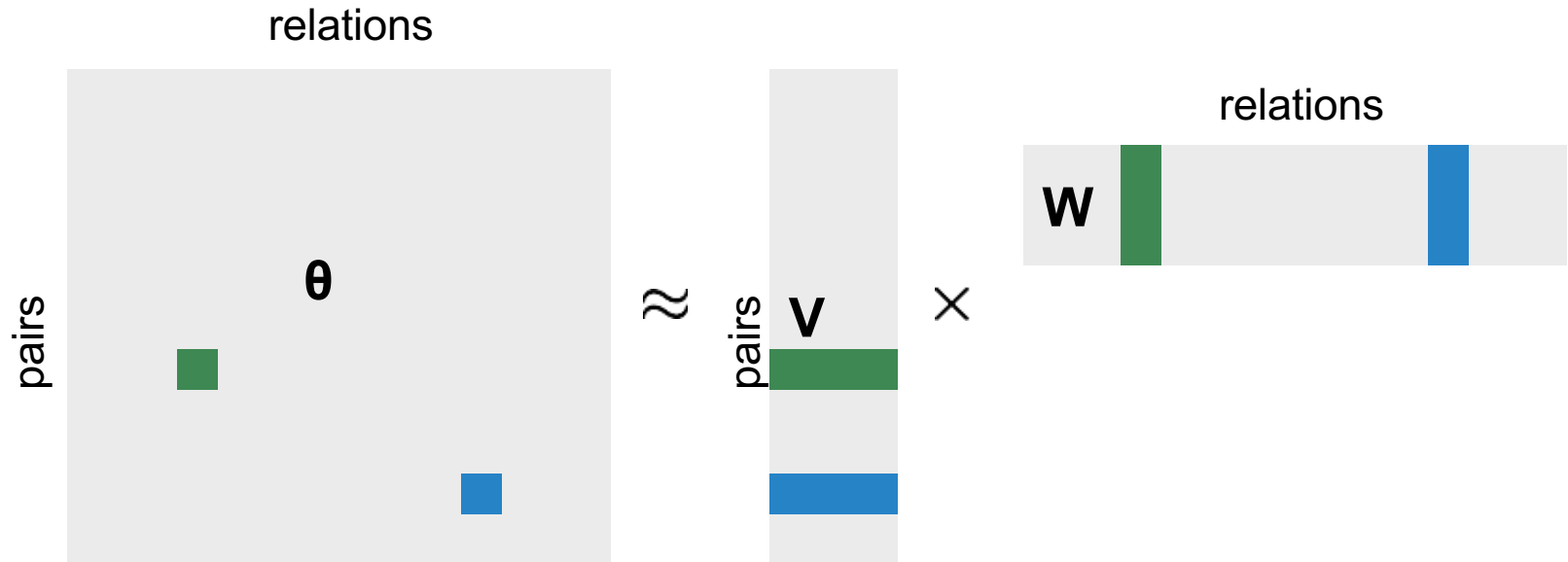
Loss Function:
$$\max_{\mathbf{v}, \mathbf{w}} \log \prod_{x,y,r} \exp \langle \mathbf{v}^{x,y}, \mathbf{w}_r \rangle - \lambda (\|\mathbf{v}\|_2^2 + \|\mathbf{w}\|_2^2)$$

Desiderata from the training algorithm:

- Do not instantiate the whole matrix!
- Do not hold all the observed cells in memory
- Each iteration linear in the no. of observations

Solution: Stochastic Gradient Descent!

Training: Stochastic Updates



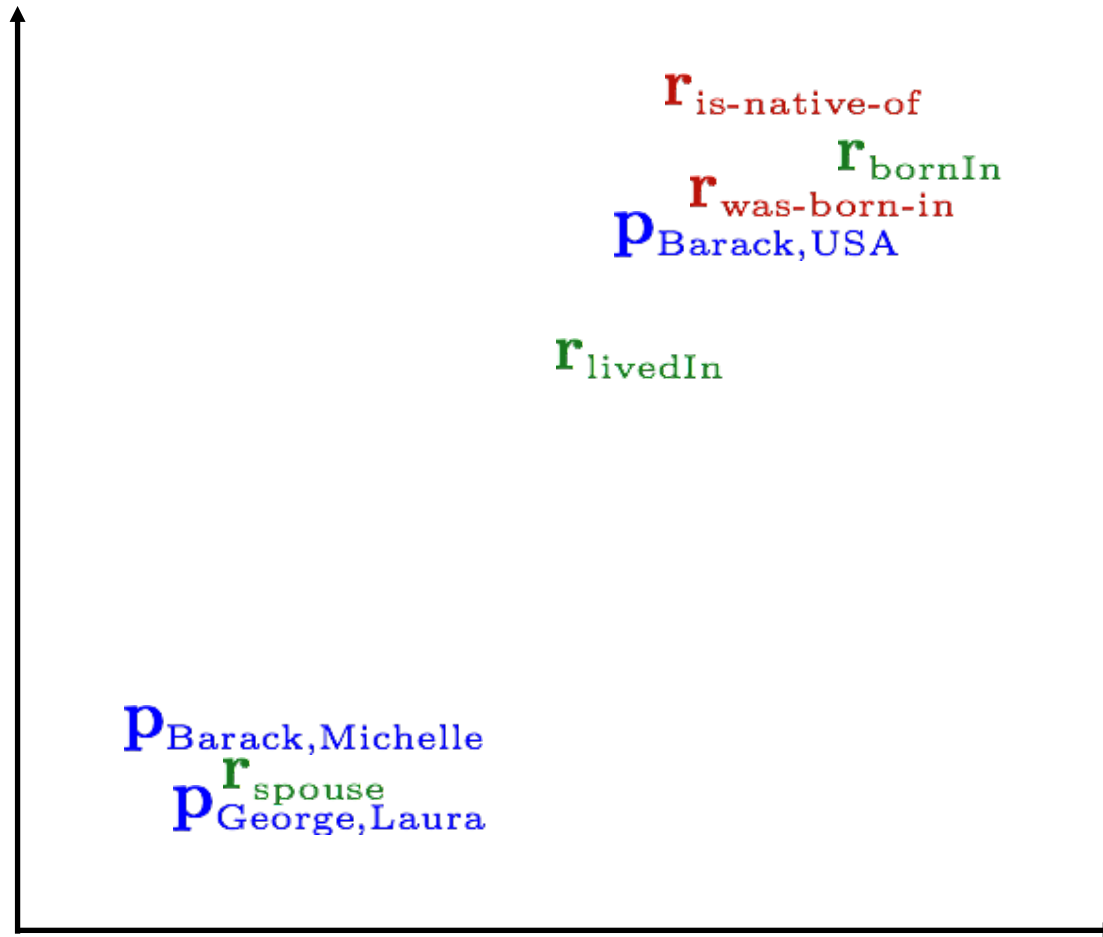
Pick an **observed** cell, $\theta_{x,y}^r$:

- Update $\mathbf{v}^{x,y}$ & \mathbf{w}^r such that $\theta_{x,y}^r$ is higher

Pick any random cell, assume it is **negative**:

- Update $\mathbf{v}^{x,y}$ & \mathbf{w}^r such that $\theta_{x,y}^r$ is lower

Relation Embeddings



Embeddings \sim Logical Relations

Relation Embeddings, w

- Similar embedding for 2 relations denote they are paraphrases
 - is married to, spouseOf(X,Y), /person/spouse
- One embedding can be contained by another
 - $w(\text{topEmployeeOf}) \subset w(\text{employeeOf})$
 - $\text{topEmployeeOf}(X,Y) \rightarrow \text{employeeOf}(X,Y)$
- Can capture logical patterns, without needing to specify them!

Entity Pair Embeddings, v

- Similar entity pairs denote similar relations between them
- Entity pairs may describe multiple “relations”
 - independent **foundedBy** and **employeeOf** relations

Similar Embeddings

similar underlying embedding

X own percentage of Y X buy stake in Y

similar embedding

Time, Inc Amer. Tel. and Comm.	1	1
Volvo Scania A.B.		1
Campeau Federated Dept Stores		
Apple HP		

Successfully predicts “Volvo owns percentage of Scania A.B.”
from “Volvo bought a stake in Scania A.B.”

Implications

$X \text{ historian at } Y \rightarrow X \text{ professor at } Y$

		X professor at Y	X historian at Y
(Freeman,Harvard) → (Boyle,OhioState)	Kevin Boyle Ohio State		1
	R. Freeman Harvard	1	

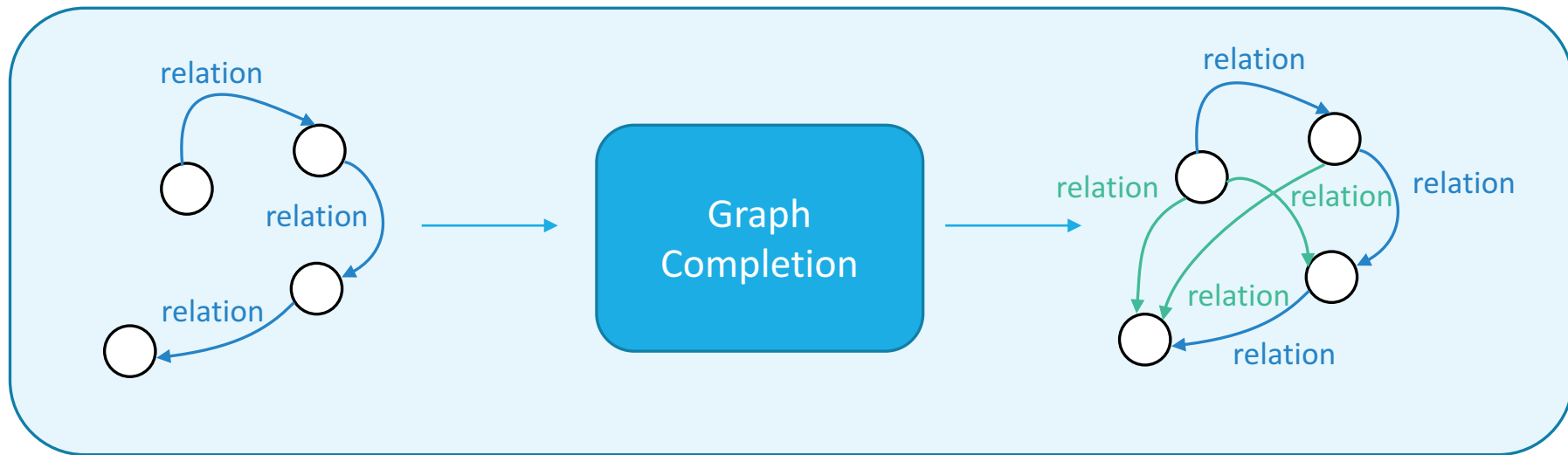
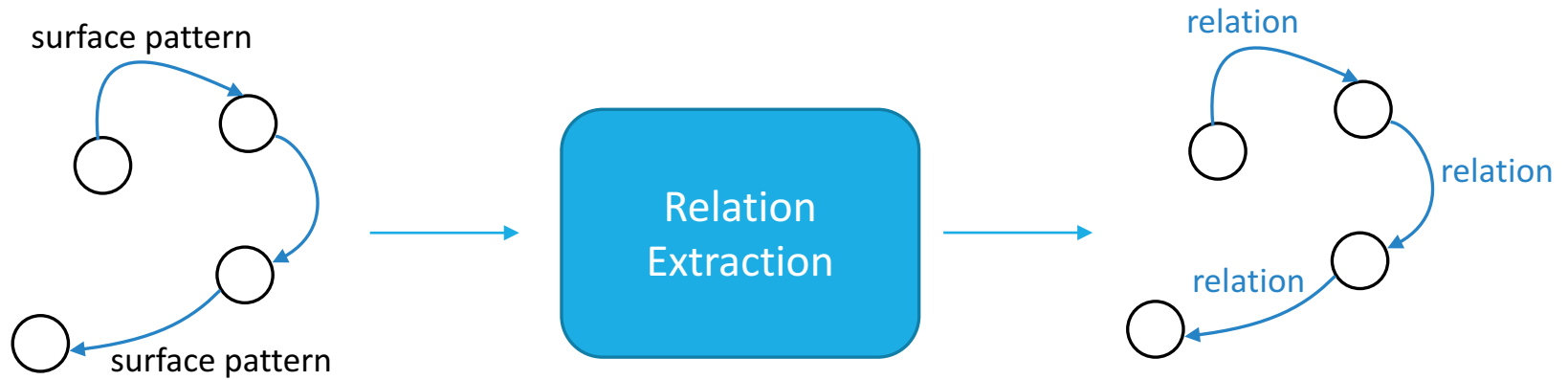
Learns asymmetric entailment:

$\text{PER historian at UNIV} \rightarrow \text{PER professor at UNIV}$

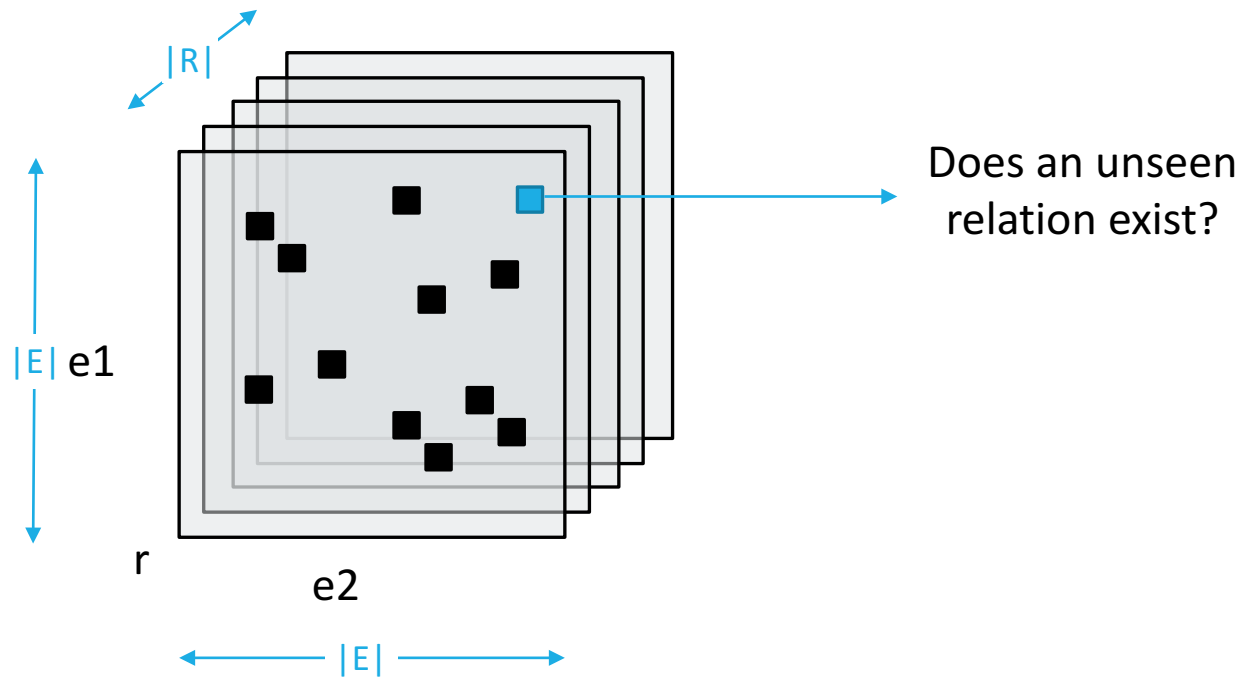
But,

$\text{PER professor at UNIV} \not\rightarrow \text{PER historian at UNIV}$

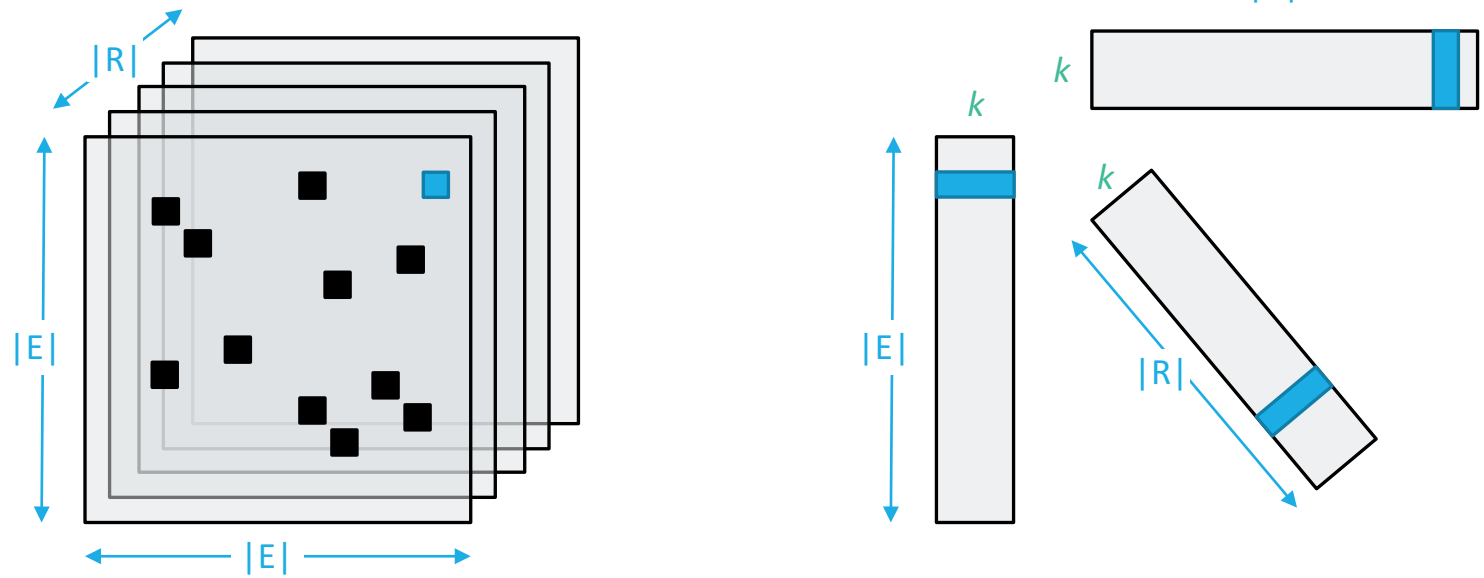
Two Related Tasks



Tensor Formulation of KG



Factorize that Tensor



Many Different Factorizations

CANDECOMP/PARAFAC-Decomposition

$$S(r(a, b)) = \sum_k R_{r,k} \cdot e_{a,k} \cdot e_{b,k}$$

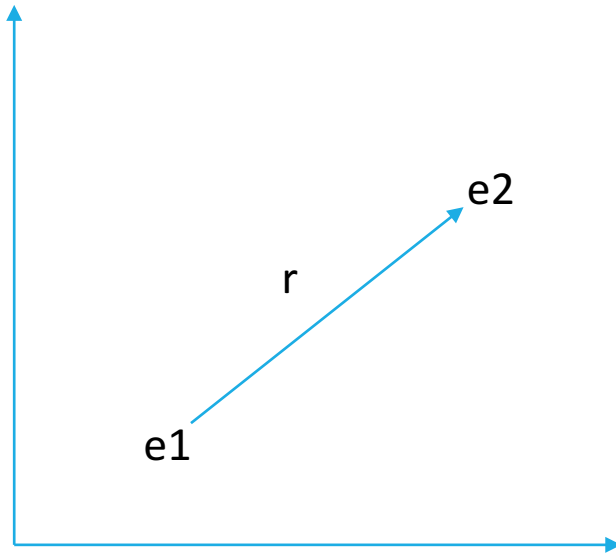
Tucker2 and RESCAL Decompositions

$$S(r(a, b)) = (\mathbf{R}_r \times \mathbf{e}_a) \times \mathbf{e}_b$$

Model E

$$S(r(a, b)) = \mathbf{R}_{r,1} \cdot \mathbf{e}_a + \mathbf{R}_{r,2} \cdot \mathbf{e}_b$$

Translation Embeddings



TransE

$$S(r(a, b)) = -\|\mathbf{e}_a + \mathbf{R}_r - \mathbf{e}_b\|_2^2$$

Parameter Estimation: SGD

Training Objective

$$\theta = \operatorname{argmax}_{\theta} \sum_{r_{ab} \in \mathcal{P}} \sum_{r'_{a'b'} \in \mathcal{N}} \mathcal{L}()$$

Distance

$$\mathcal{L}(x, y) = -\|x - y\|_2^2$$

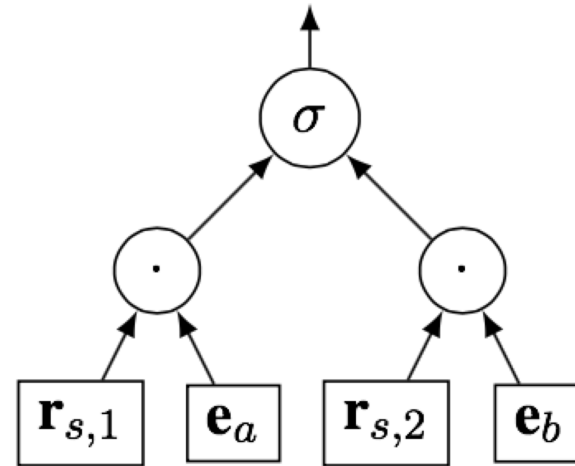
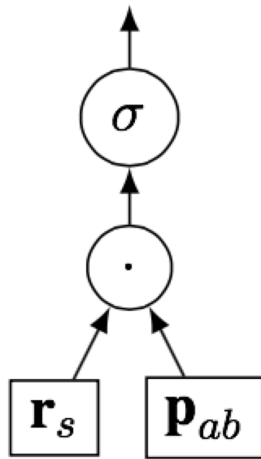
Likelihood

$$\mathcal{L}(x, y) = p(x)^{p(y)} (1 - p(x))^{(1-p(y))}$$

Stochastic Gradient Descent

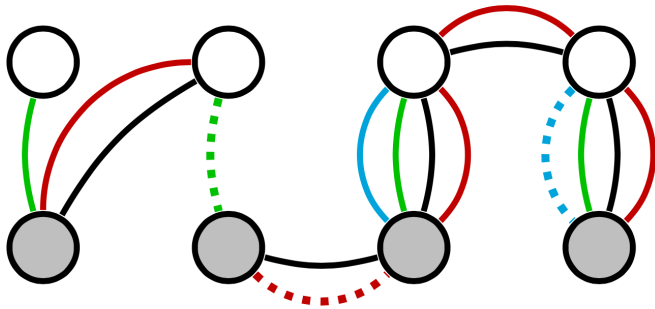
Negative Sampling ...

Matrix vs Tensor Factorization

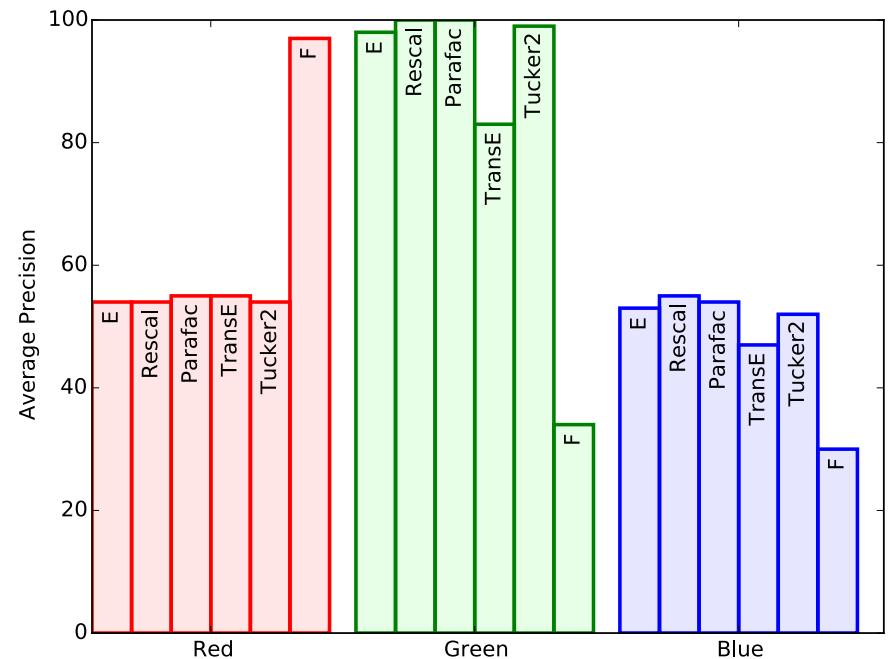


- No encoding of type information
 - Can only predict for entity pairs that appear in text together
 - Sufficient evidence has to be seen for each entity pair
- Assume low-rank for pairs
 - But many relations are not!
 - Spouse: you can have only ~ 1
 - Cannot learn pair specific information

What they can, and can't, do..



- **Red**: deterministically implied by **Black**
 - needs *pair-specific* embedding
 - Only **F** is able to generalize
- **Green**: needs to estimate entity types
 - needs *entity-specific* embedding
 - Tensor factorization generalizes, **F** doesn't
- **Blue**: implied by **Red** and **Green**
 - Nothing works much better than random



Compositional Neural Models

So far, we're learning vectors for each entity/surface pattern/relation..

But learning vectors independently ignores “composition”

Composition in Surface Patterns

- Every surface pattern is not unique
- Synonymy: A is B's spouse.
A is married to B.
- Inverse: X is Y's parent.
Y is one of X's children.
- Can the representation learn this?

Composition in Relation Paths

- Every relation path is not unique
- Explicit: A parent B, B parent C
A grandparent C
- Implicit: X bornInCity Y, Y cityInState Z
X “bornInState” Z
- Can the representation capture this?

Composing Dependency Paths

... was born to ...



... 's parents are ...



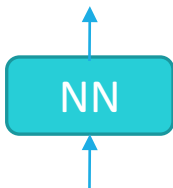
(never appears in
training data)

\birthplace

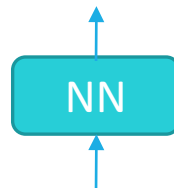


But we don't need linked data to know they mean similar things...

Use neural networks to produce the embeddings from text!



... was born to ...

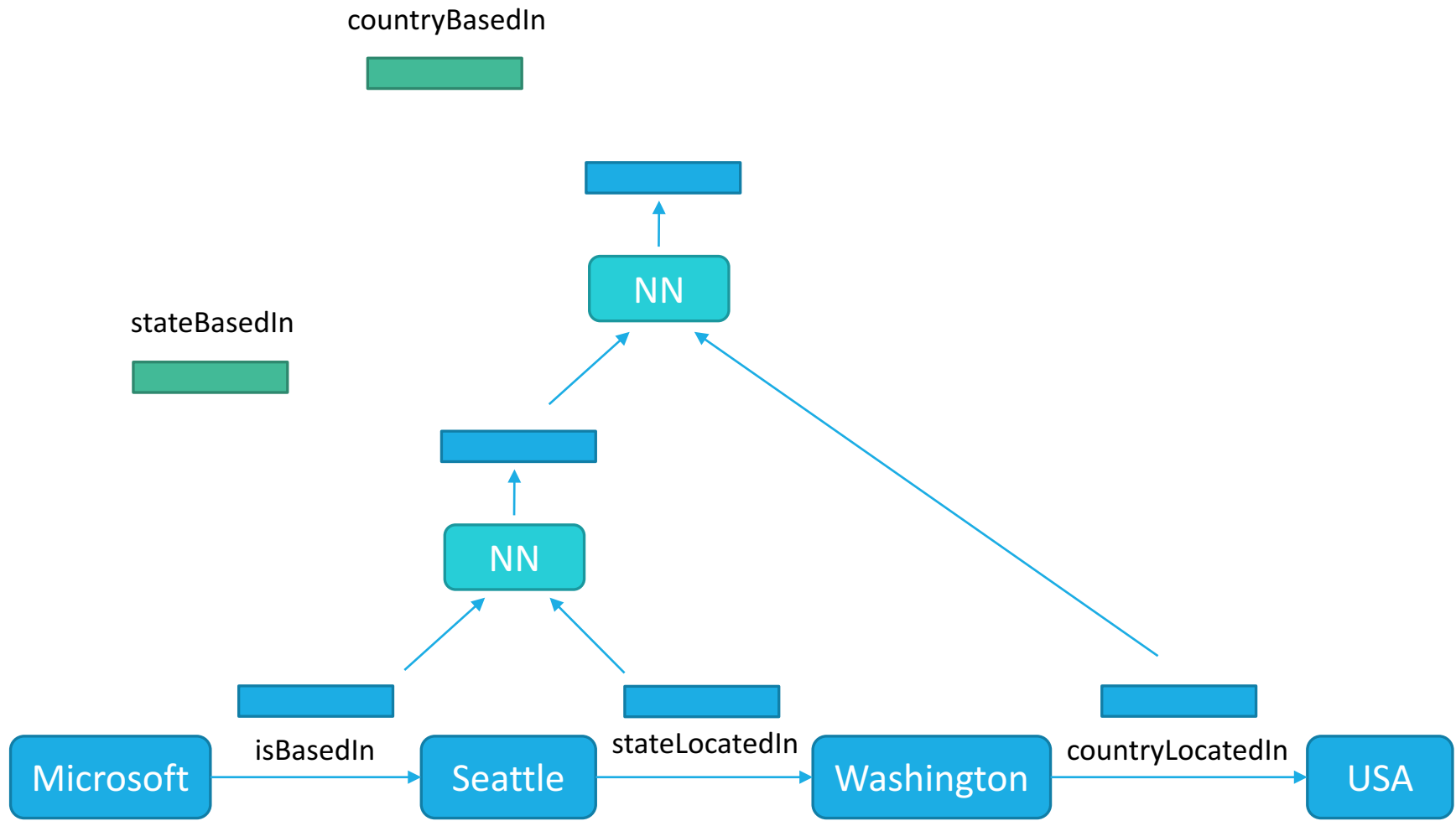


... 's parents are ...



\birthplace

Composing Relational Paths



Review: Embedding Techniques

Two Related Tasks:

- Relation Extraction from Text
- Graph (or Link) Completion

Relation Extraction:

- Matrix Factorization Approaches

Graph Completion:

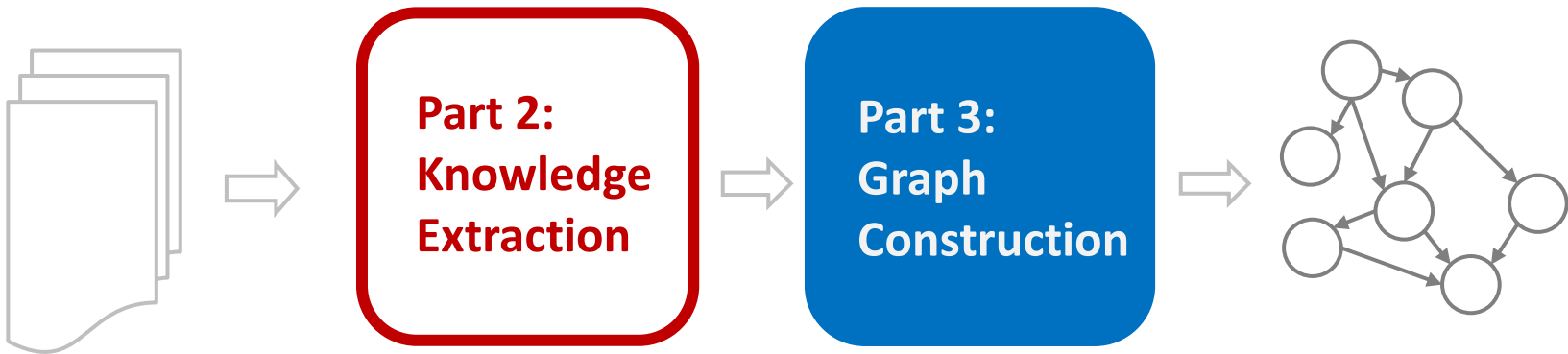
- Tensor Factorization Approaches

Compositional Neural Models

- Compose over dependency paths
- Compose over relation paths

Tutorial Overview

Part 1: Knowledge Graphs



Part 4: Critical Analysis