

Embedding Probabilistic Logic for Machine Reading

aka Towards Two-Way Interaction with Reading Machines

▶ **Sebastian Riedel (University College London)**



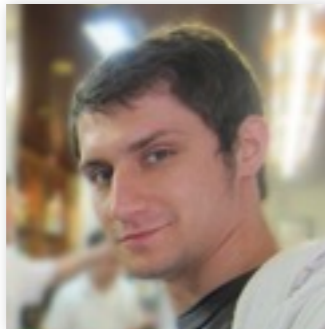
Collaborators



Tim Rocktäschel
UCL



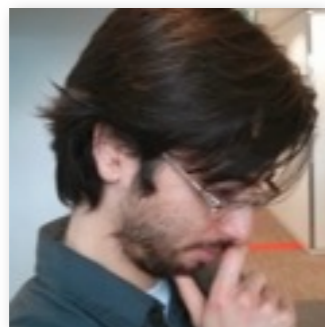
Limin Yao
UMass Amherst (Twitter)



Matko Bosnjak
UCL



Andrew McCallum
UMass Amherst



Ivan Sanchez
UCL



Ben Marlin
UMass Amherst



Sameer Singh
UWash

Machine Reading

“Who **works** in London and is **interested** in NLP?”

*interest(x, NLP),
worksFor(x, y),
in(y, London)*

[Kwiatkowski et al., 2013]



Narrow domain-specific schema

[Mintz et al., 2009]

Semantics

Syntax

Coreference

Statistical NLP

“**Sebastian Riedel** works in the area of **NLP** and is now Lecturer at **UCL**”

Machine Reading

[Riedel et al., 2013]

“Who **works** in London and is **interested** in NLP?”

*interest(x, NLP),
worksFor(x, y),
in(y, London)*

*in(UCL, London)
works-in-area-of(Seb, NLP)
lecturer-at(Seb, UCL)*

Relational DB

Semantics

Wide universal schema

Syntax

Coreference

Statistical NLP

“**Sebastian Riedel** works in the area of **NLP** and is now Lecturer at **UCL**”

Semantics as Reasoning

[Riedel et al., 2013]

“Who **works** in London and is **interested** in NLP?”

*interest(x, NLP),
worksFor(x, y),
in(y, London)*

*in(UCL, London)
works-in-area-of(Seb, NLP)
lecturer-at(Seb, UCL)
worksFor(x, y):
 faculty-at(x, y)

interest(x, y):
 works-in-area-of(x, y) [0.9]

faculty-at(x, y):
 lecturer-at(x, y)*

**Statistical Relational
Learner and Reasoner**

Wide universal schema

Syntax

Coreference

Statistical NLP

“**Sebastian Riedel** works in the area of **NLP** and is now Lecturer at **UCL**”

Benefit: Transitive Reasoning

“Who works in London and is interested in NLP?”

*interest(x, NLP),
worksFor(x, y),
in(y, London)*

in(UCL, London)
works-in-area-of(Seb, NLP)
lecturer-at(Seb, UCL)

*worksFor(x, y):
faculty-at(x, y)*

*interest(x, y):
works-in-area-of(x, y) [0.9]*

*faculty-at(x, y):
lecturer-at(x, y)*

**Statistical Relational
Learner and Reasoner**

Wide universal schema

Syntax

Coreference

Statistical NLP

“**Sebastian Riedel** works in the area of **NLP** and is now Lecturer at **UCL**”

Benefit: More Coverage

“Who is **faculty** in London and interested in NLP?”

*interest(x, NLP),
worksFor(x, y),
in(y, London)*

in(UCL, London)
works-in-area-of(Seb, NLP)
lecturer-at(Seb, UCL)
*worksFor(x, y):
faculty-at(x, y)*
*interest(x, y):
works-in-area-of(x, y) [0.9]*
*faculty-at(x, y):
lecturer-at(x, y)*

**Statistical Relational
Learner and Reasoner**

Wide universal schema

Syntax

Coreference

Statistical NLP

“**Sebastian Riedel** works in the area of **NLP** and is now Lecturer at **UCL**”

Benefit: Code Reuse

“Who **lives in** London and is interested in NLP?”

*interest(x, NLP),
worksFor(x, y),
in(y, London)*

*in(UCL, London)
works-in-area-of(Seb, NLP)
lecturer-at(Seb, UCL)
worksFor(x, y):
 faculty-at(x, y)

interest(x, y):
 works-in-area-of(x, y) [0.9]*

*livesIn(x, z):
 worksFor(x, y),
 locatedIn(y, z) [0.6]*

Statistical Relational Learner and Reasoner

[Lao et al., 2011]

Wide universal schema

Syntax

Coreference

Statistical NLP

“**Sebastian Riedel** works in the area of **NLP** and is now Lecturer at **UCL**”

Joint Inference

“Who **lives in** London and is interested in NLP?”

*interest(x, NLP),
worksFor(x, y),
in(y, London)*

Wide universal schema

in(UCL, London)
works-in-area-of(Seb, NLP)
lecturer-at(Seb, UCL)
*worksFor(x, y):
 faculty-at(x, y)*
*interest(x, y):
 works-in-area-of(x, y) [0.9]*

*livesIn(x, z):
 worksFor(x, y),
 locatedIn(y, z) [0.6]*

Syntax

Coreference

**Statistical Relational
Learner and Reasoner**

Statistical NLP

“**Sebastian Riedel** works in the area of **NLP** and is now Lecturer at **UCL**”

Reasoner and Learner

*Statistical Relational
Learner and Reasoner*

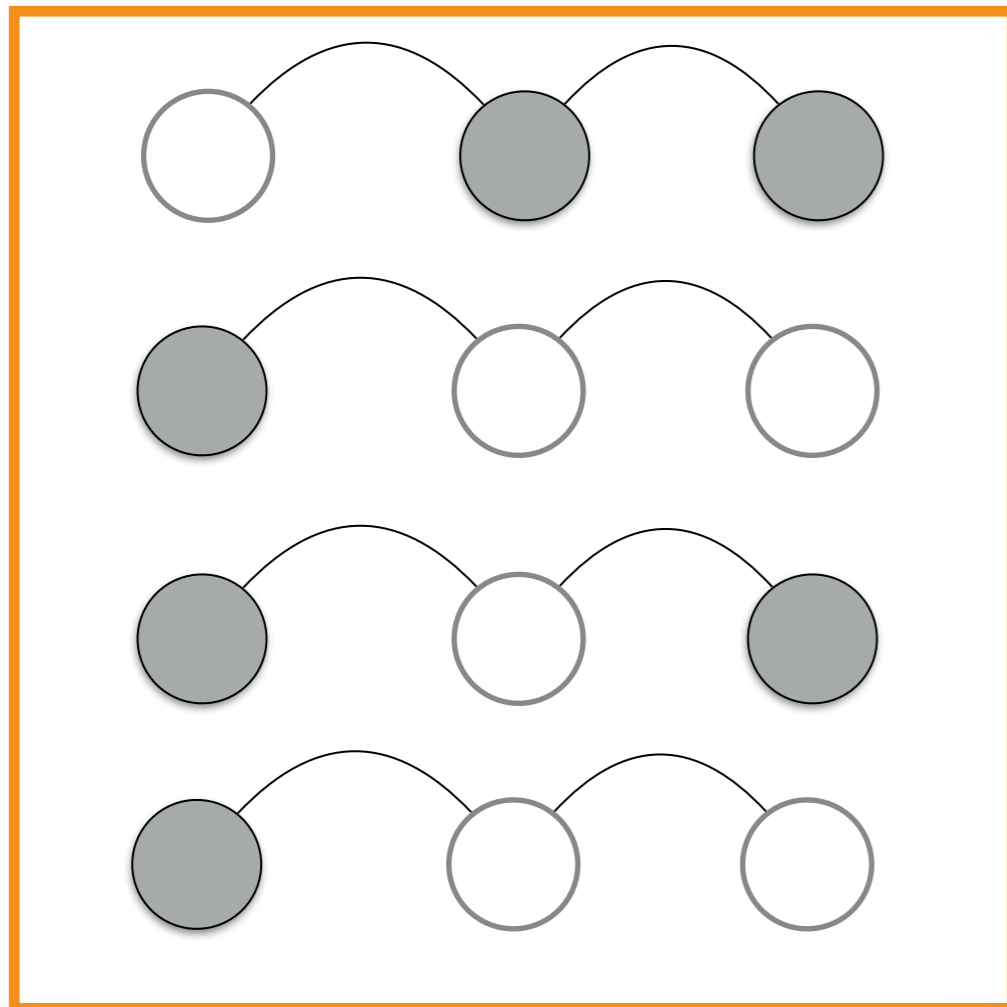


?

Probabilistic Logics

Use (weighted) logics to define graphical models

lecturer-at *prof-at* *works-for*



Examples

▶ Markov Logic

[Richardson and Domingos, 2006]

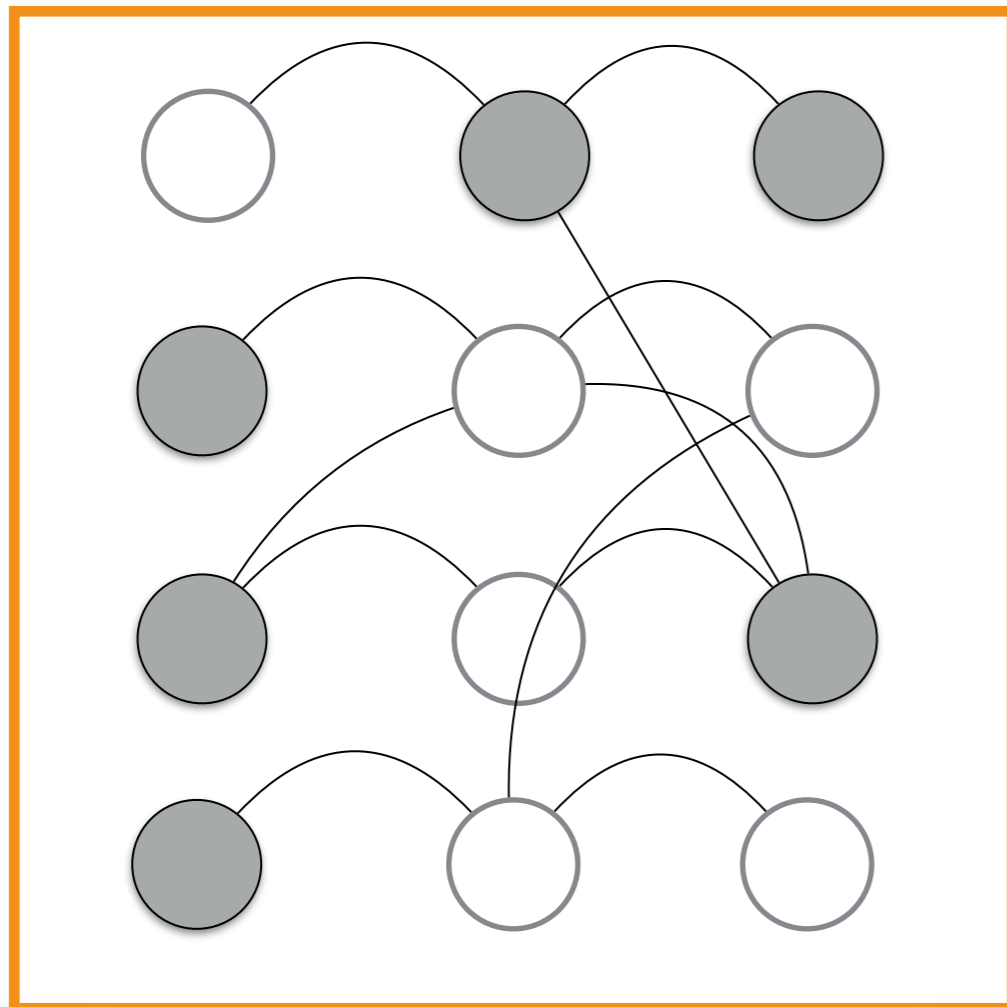
▶ Bayesian Logic Programs

[Kersting, 2007]

Probabilistic Logics

Use (weighted) logics to define graphical models

lecturer-at prof-at works-for



- Problems
- ▶ Inference
 - ▶ Rule Learning

Matrix Factorization

Think of database as a matrix or tensor

lecturer-at prof-at works-for

	1	1
1		
1		1
1		



Matrix Factorization

Embed entity (pairs) in low dimensional vector spaces

lecturer-at *prof-at* *works-for*

	1	1	?	?
1			?	?
1		1	?	?
1			?	?

Matrix Factorization

Embed relations in low dimensional vector spaces

		1	1
1			
1			1
1			

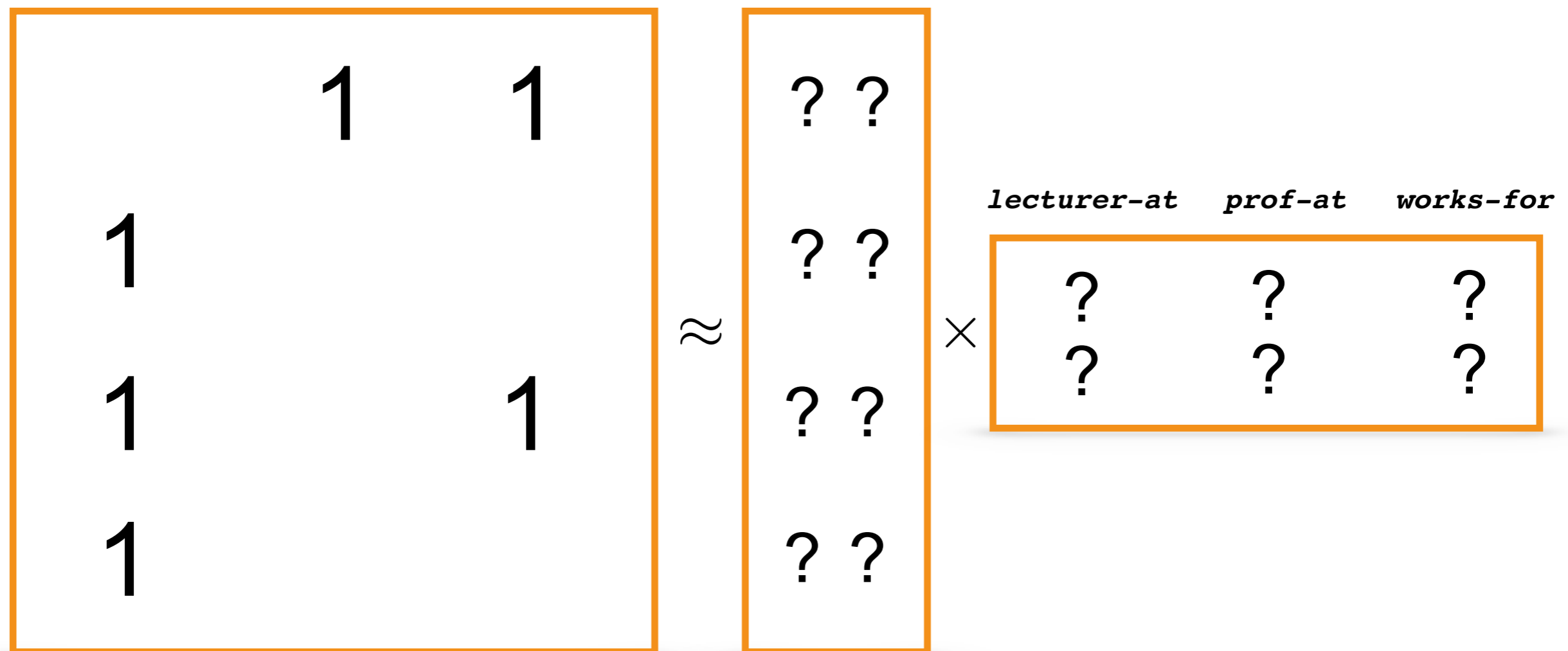
?	?
?	?
?	?
?	?

lecturer-at prof-at works-for

?	?	?
?	?	?

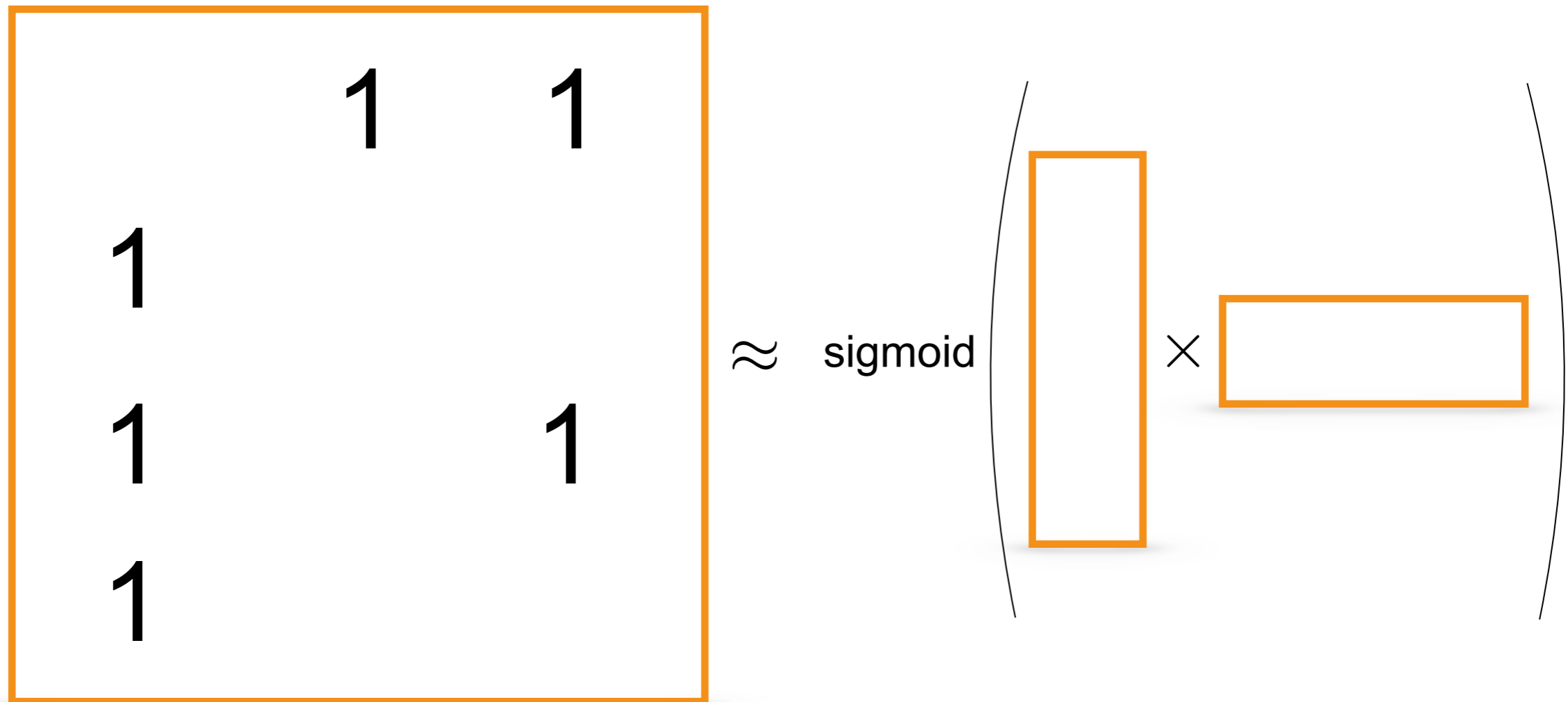
Matrix Factorization

Find a matrix-matrix product that approximates observed DB



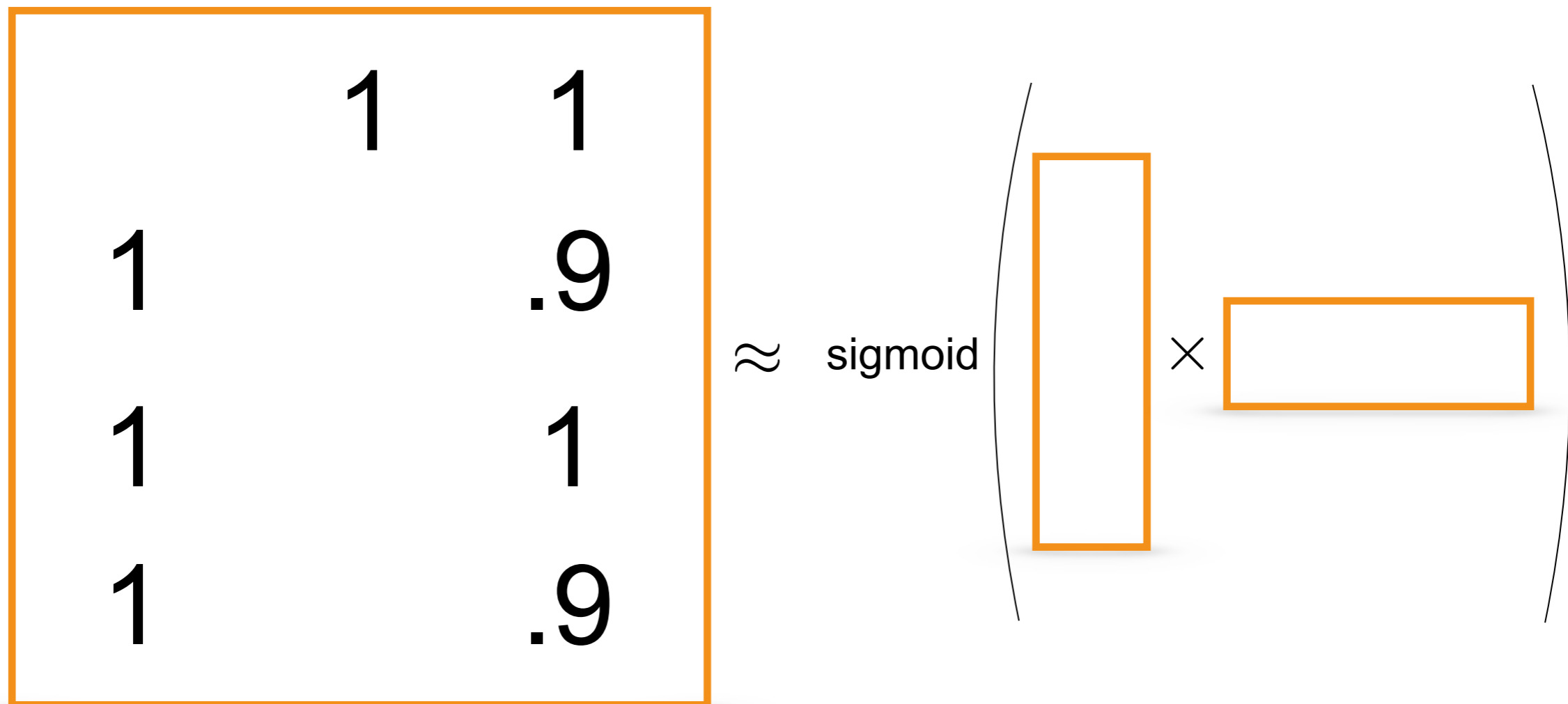
Matrix Factorization

Or a non-linear function of this product



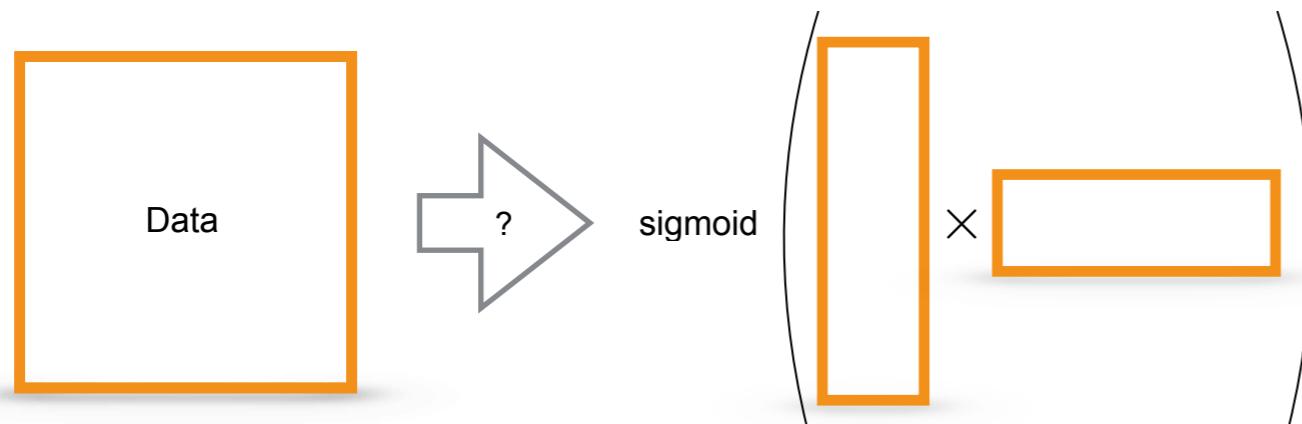
Matrix Factorization

Low rank forces some 0 cells to become non-zero => prediction

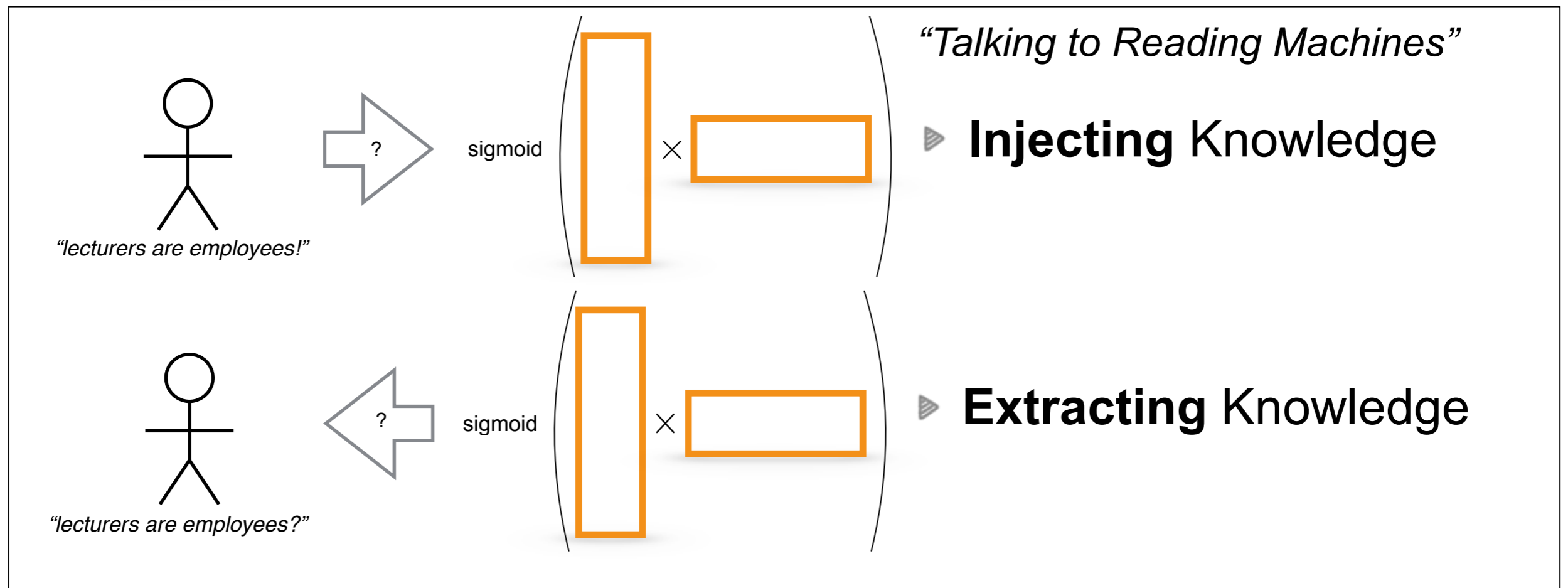


[Nickel, Bordes, ...]

Overview



► **Matrix Factorization Models**



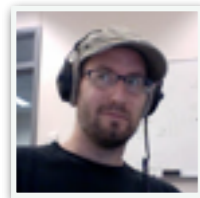
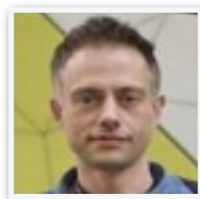
"Talking to Reading Machines"

► **Injecting Knowledge**

► **Extracting Knowledge**

Universal Schema Matrix

Schema contains structured and unstructured (~OpenIE) relations



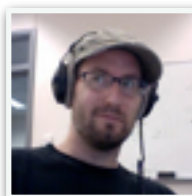







<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
1	1		1
1			
1		1	

Goal: Learn to *Complete*

Schema contains structured and unstructured (~OpenIE) relations

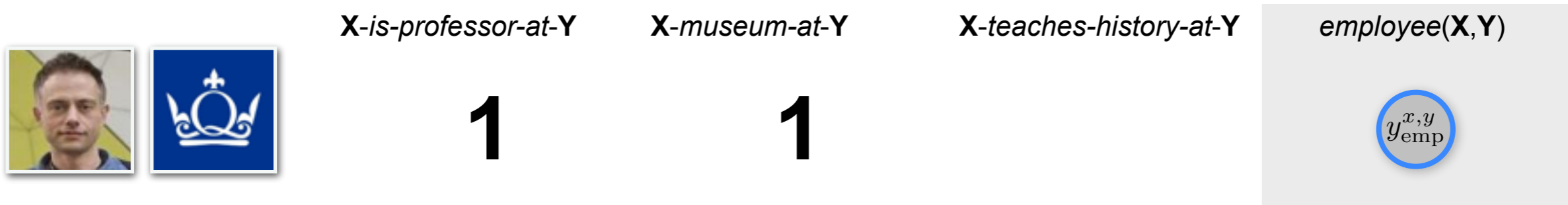


	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
 	1	1	?	1
 	1	?	?	?
 	1	?	1	?
 	?	?	?	?

Model N: Baseline Classifier

[Mintz et al 2009,...]

Standard supervised relation extractor ...



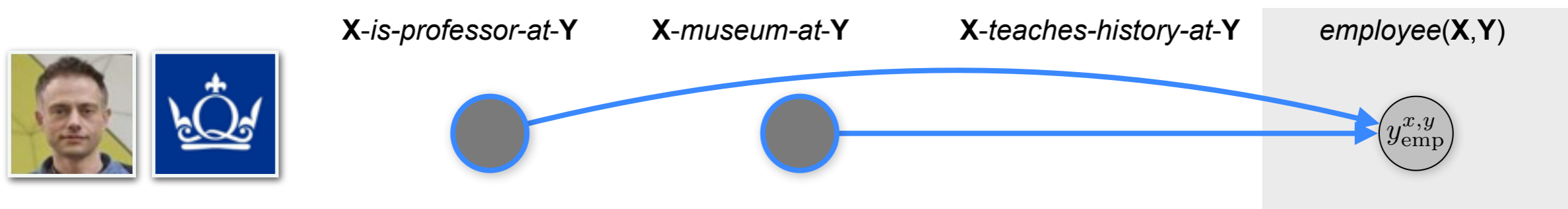
○ training data

$$p(\underline{y_{emp}^{x,y}} = 1 | \quad)$$

Model N: Classifier

[Mintz et al 2009,...]

Standard supervised relation extractor ...



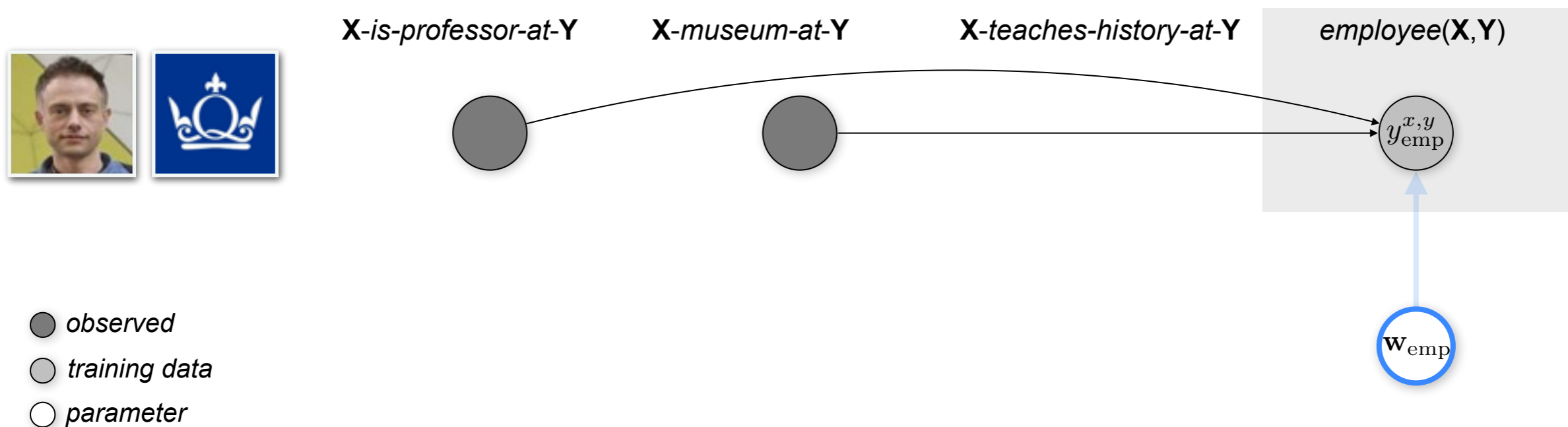
- *observed*
- *training data*

$$p(y_{\text{emp}}^{x,y} = 1 | \underline{\mathbf{f}}_{\text{emp}}^{x,y})$$

Model N: Classifier

[Mintz et al 2009,...]

Standard supervised relation extractor ...

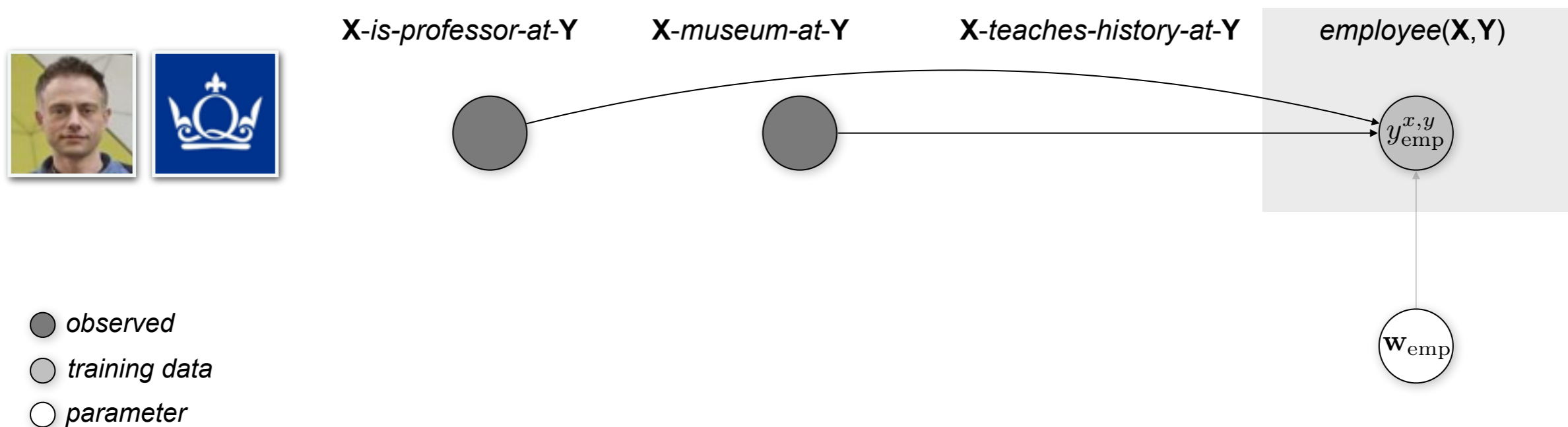


$$p(y_{\text{emp}}^{x,y} = 1 | \mathbf{f}_{\text{emp}}^{x,y}, \underline{\mathbf{w}}_{\text{emp}})$$

Model N: Classifier

[Mintz et al 2009,...]

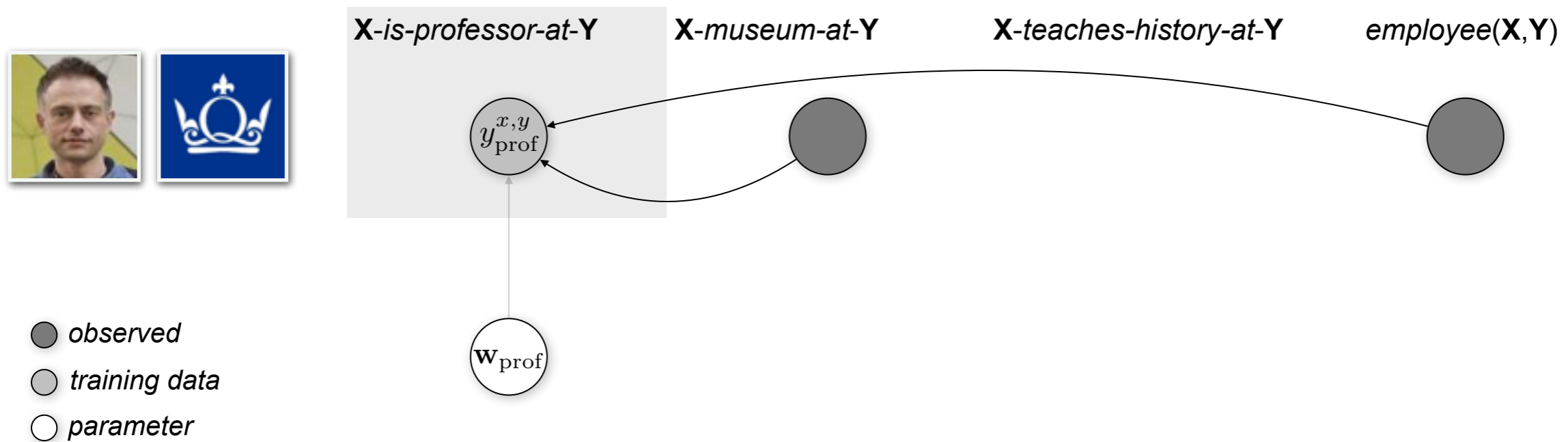
Standard supervised relation extractor ...



$$p(y_{\text{emp}}^{x,y} = 1 | \mathbf{f}_{\text{emp}}^{x,y}, \mathbf{w}_{\text{emp}}) \propto \exp[\langle \mathbf{f}_{\text{emp}}^{x,y}, \mathbf{w}_{\text{emp}} \rangle]$$

Model N: Classifier

... for each pattern

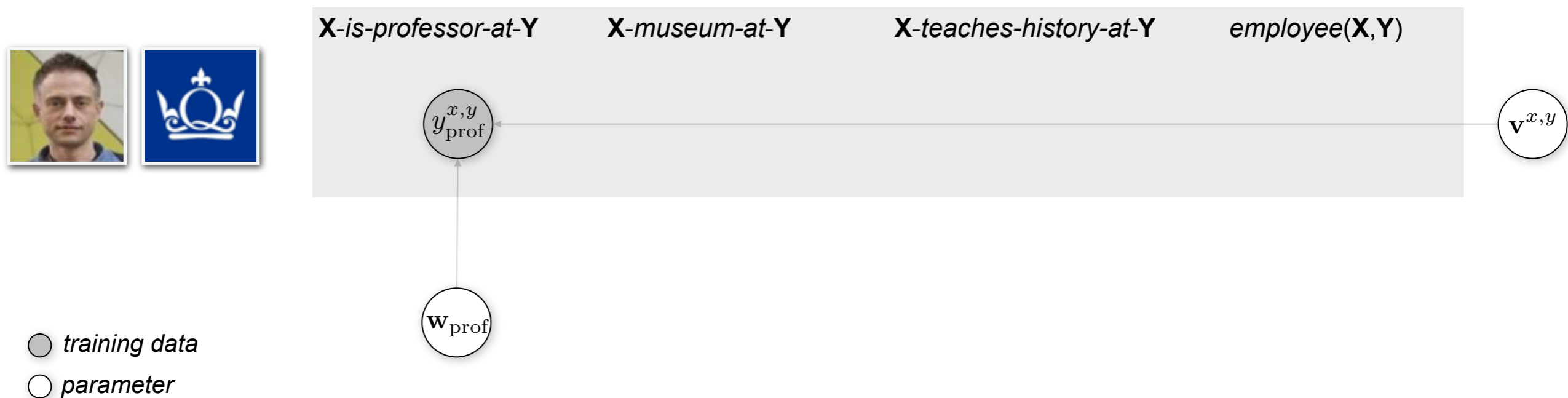


$$p(y_{\text{prof}}^{x,y} = 1 | \mathbf{f}_{\text{prof}}^{x,y}, \mathbf{w}_{\text{prof}}) \propto \exp[\langle \mathbf{f}_{\text{prof}}^{x,y}, \mathbf{w}_{\text{prof}} \rangle]$$

Model F: Latent Feature (Factorization)

[Collins et al, 2001]

Model the probability of a pair (x,y) being in relation “prof”

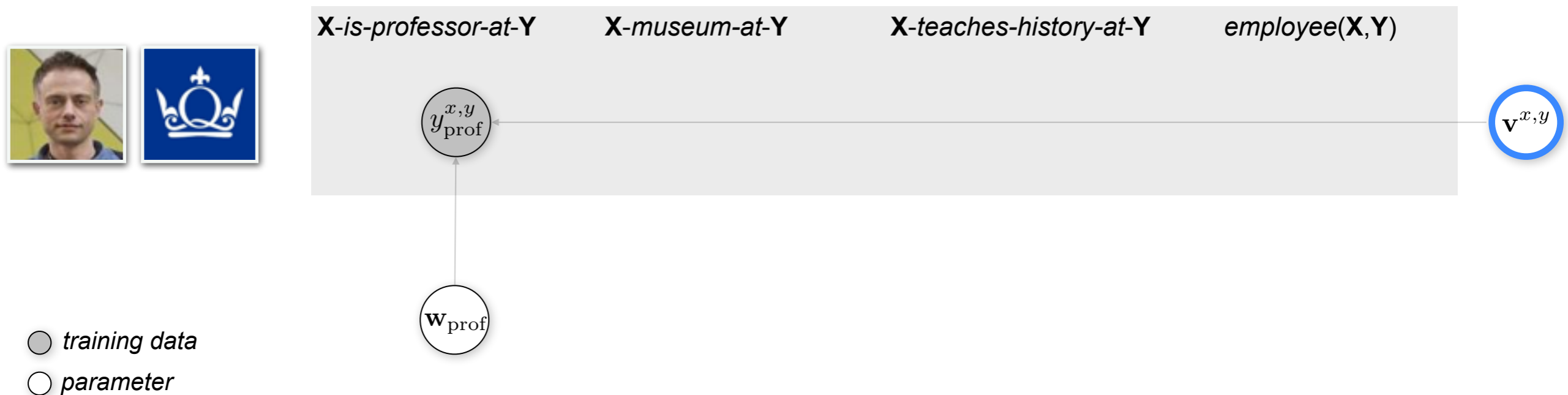


$$p(\underline{y_{\text{prof}}^{x,y}} = 1 | \mathbf{v}^{x,y}, \mathbf{w}_{\text{prof}}) \propto \exp[\langle \mathbf{v}^{x,y}, \mathbf{w}_{\text{prof}} \rangle]$$

Model F: Latent Feature (Factorization)

[Collins et al, 2001]

Per tuple **latent feature** vector

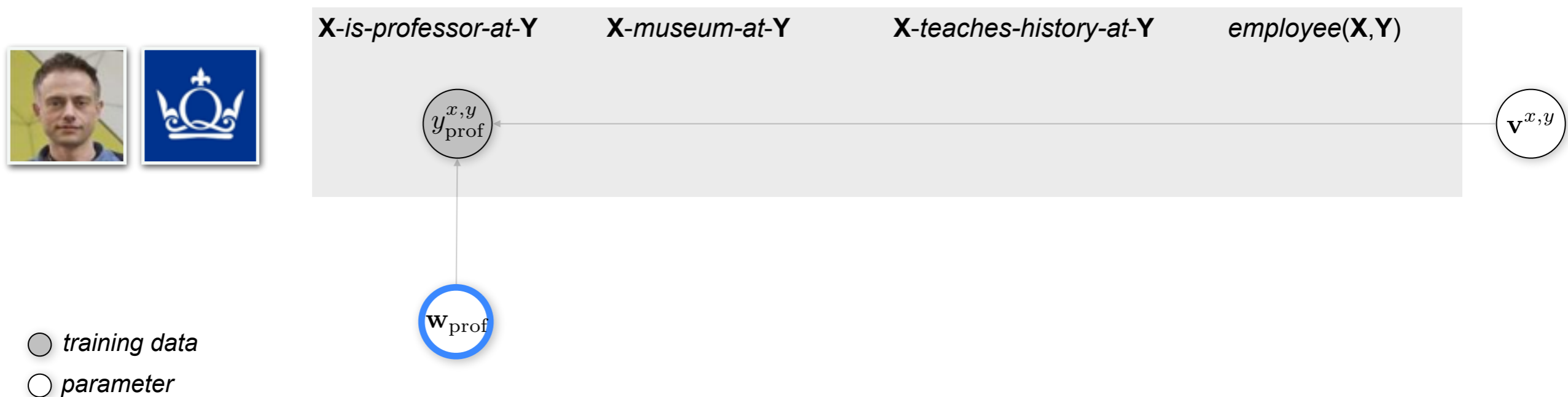


$$p(y_{prof}^{x,y} = 1 | \underline{\mathbf{v}^{x,y}}, \mathbf{w}_{prof}) \propto \exp[\langle \underline{\mathbf{v}^{x,y}}, \mathbf{w}_{prof} \rangle]$$

Model F: Latent Feature (Factorization)

[Collins et al, 2001]

Per tuple **latent feature** vector

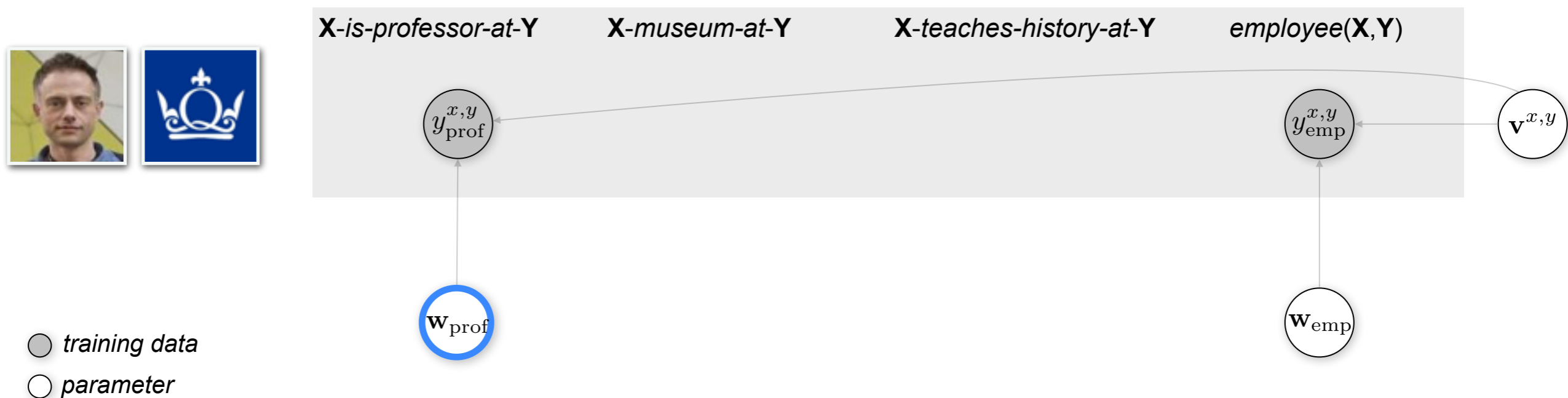


$$p(y_{prof}^{x,y} = 1 | \mathbf{v}^{x,y}, \underline{\mathbf{w}}_{prof}) \propto \exp[\langle \mathbf{v}^{x,y}, \underline{\mathbf{w}}_{prof} \rangle]$$

Model F: Latent Feature (Factorization)

[Collins et al, 2001]

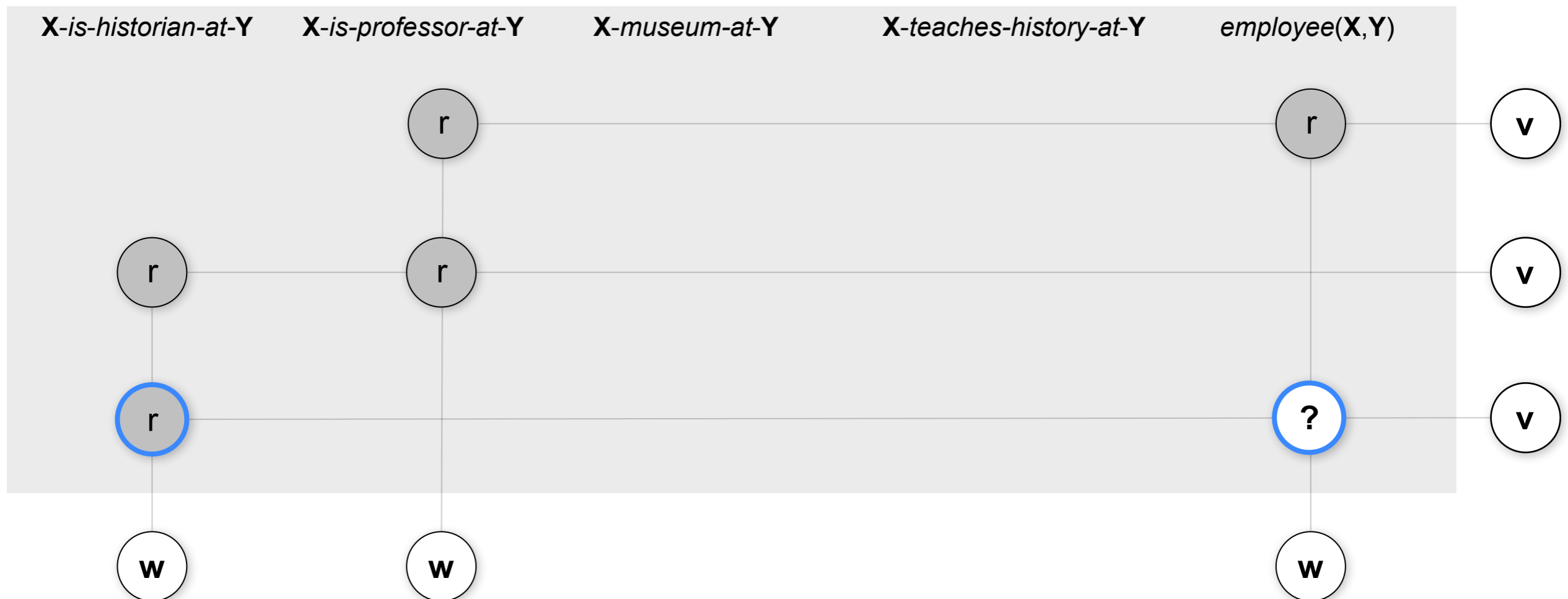
Per tuple latent feature vector



$$\begin{aligned}
 p(y_{\text{prof}}^{x,y} = 1 | \mathbf{v}^{x,y}, \mathbf{w}_{\text{prof}}) &\propto \exp[\langle \mathbf{v}^{x,y}, \mathbf{w}_{\text{prof}} \rangle] \\
 &= \textit{sigmoid}(\langle \mathbf{v}^{x,y}, \mathbf{w}_{\text{prof}} \rangle)
 \end{aligned}$$

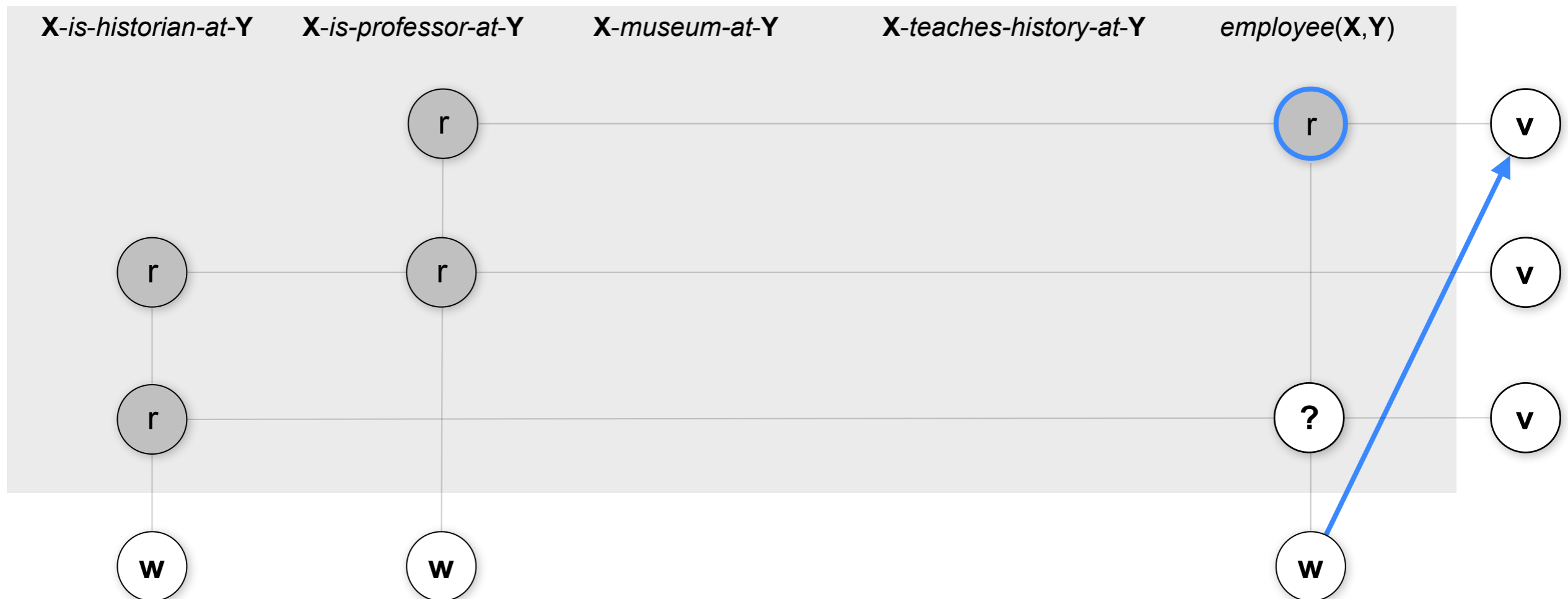
Model F: Latent Feature (Factorization)

Transitive Reasoning



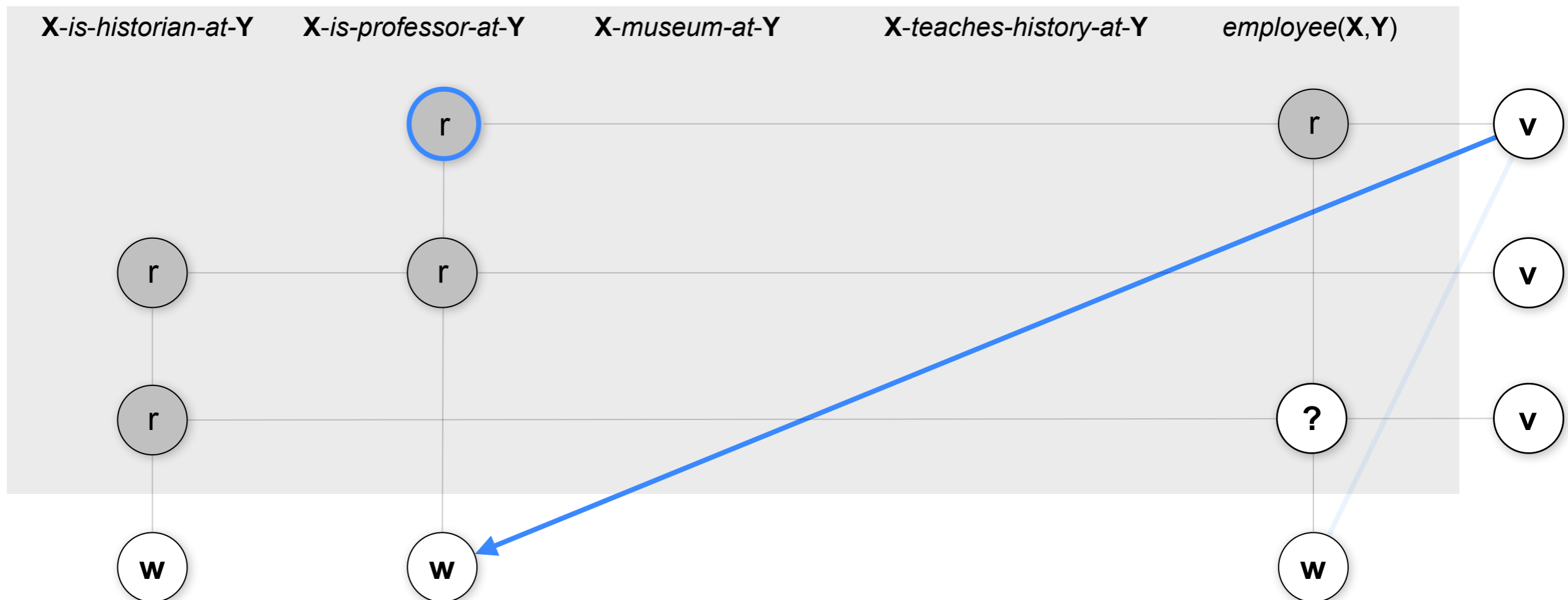
Model F: Latent Feature (Factorization)

Transitive Reasoning



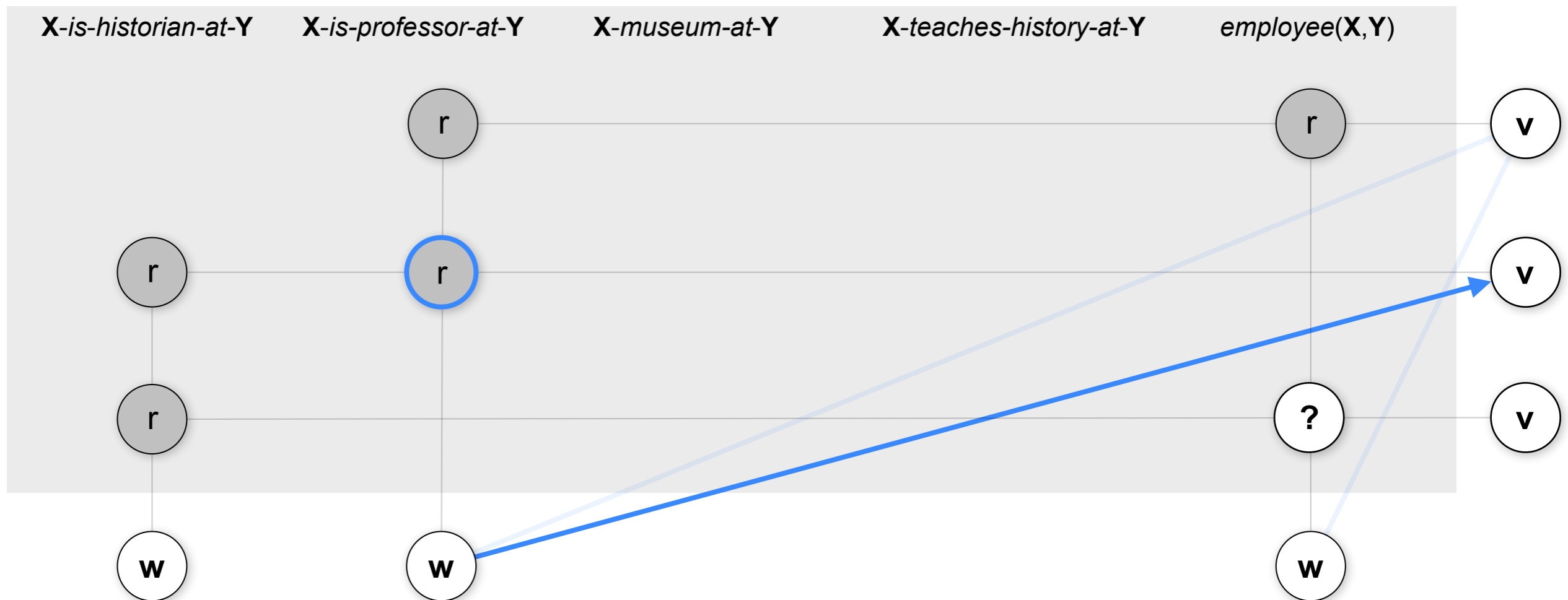
Model F: Latent Feature (Factorization)

Transitive Reasoning



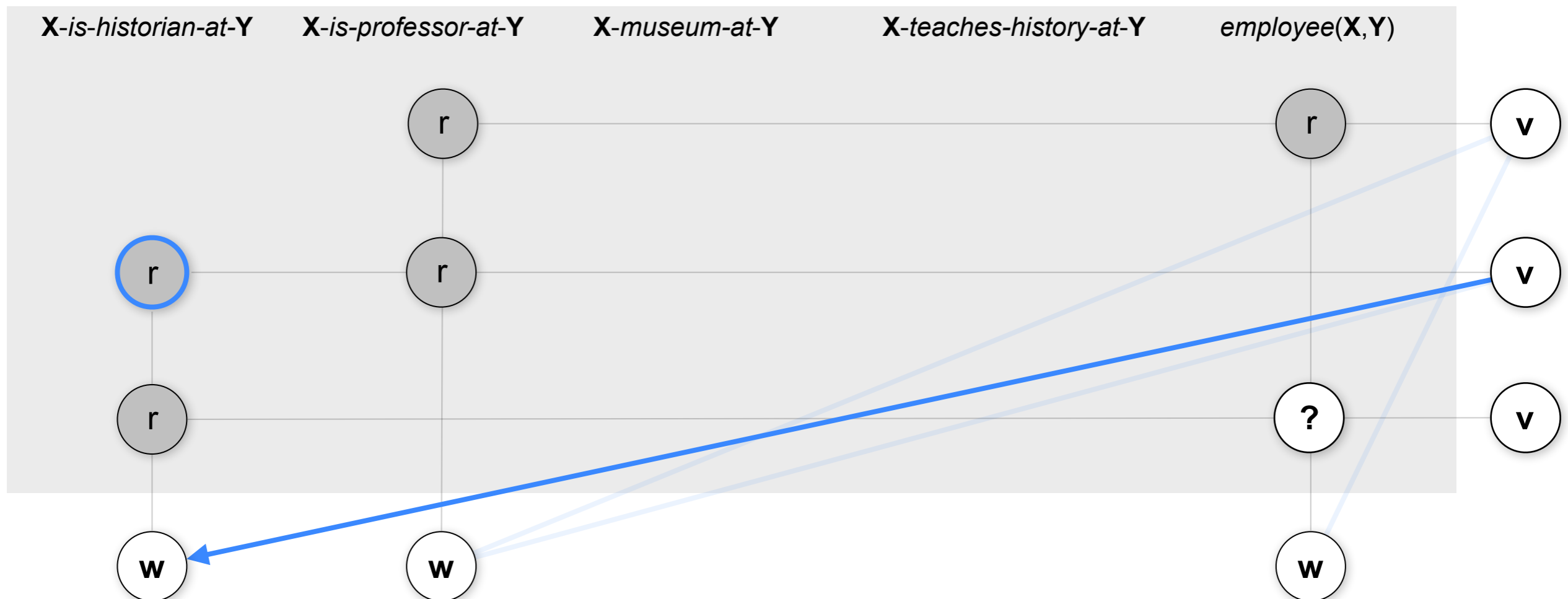
Model F: Latent Feature (Factorization)

Transitive Reasoning



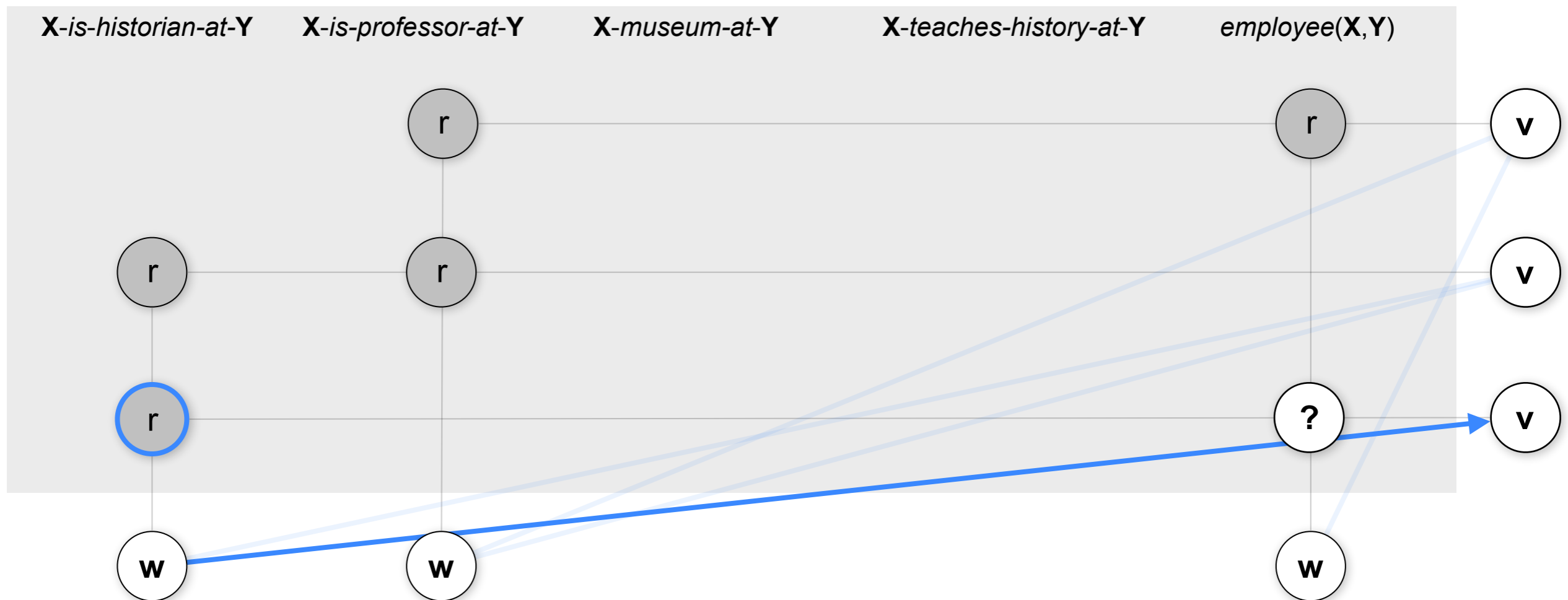
Model F: Latent Feature (Factorization)

Transitive Reasoning



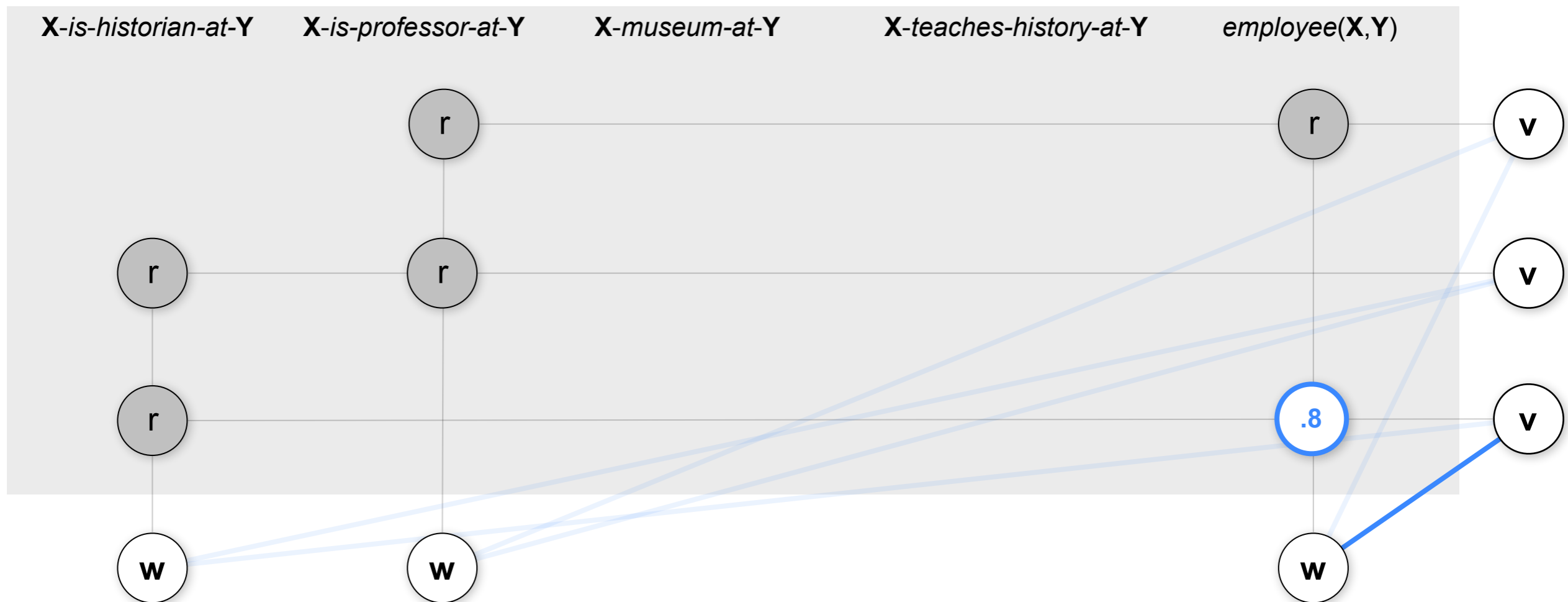
Model F: Latent Feature (Factorization)

Transitive Reasoning



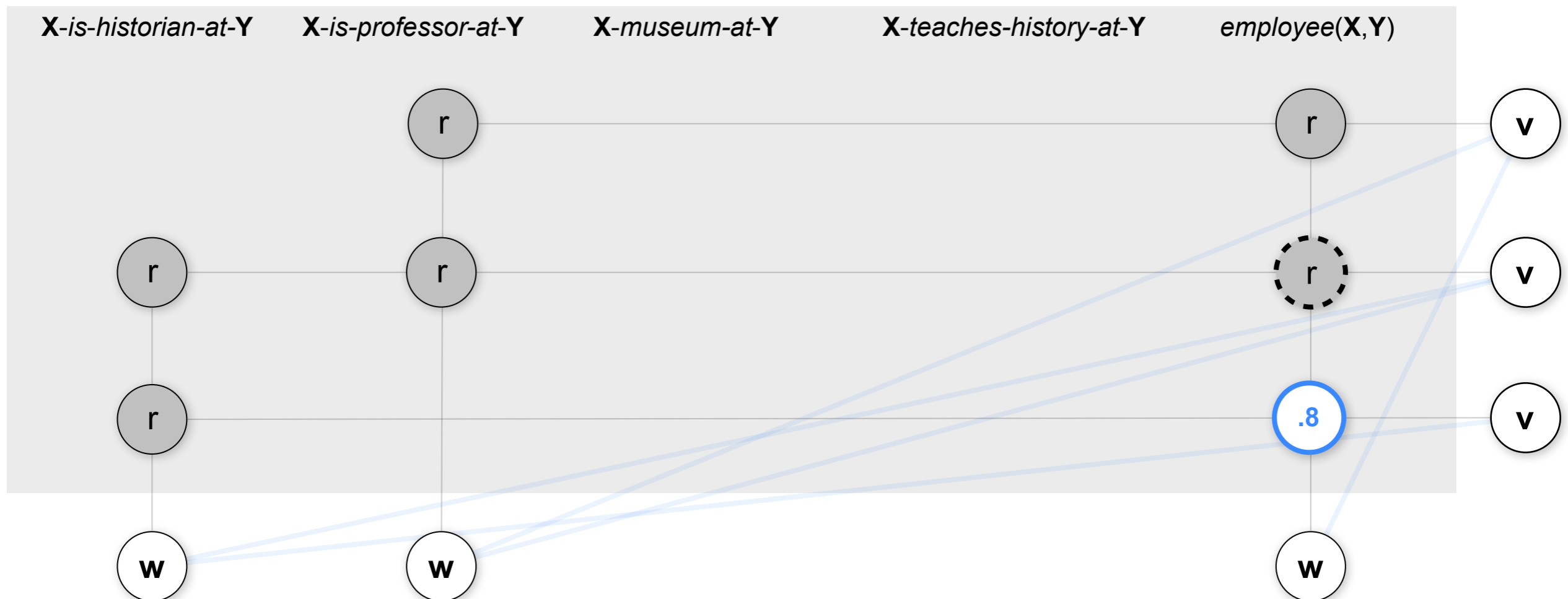
Model F: Latent Feature (Factorization)

Transitive Reasoning



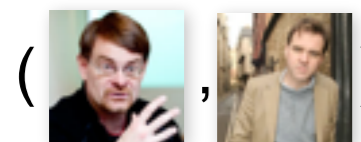
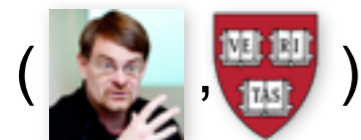
Model F: Latent Feature (Factorization)

Bootstrapping without fantasy



Model E: Selectional Preferences

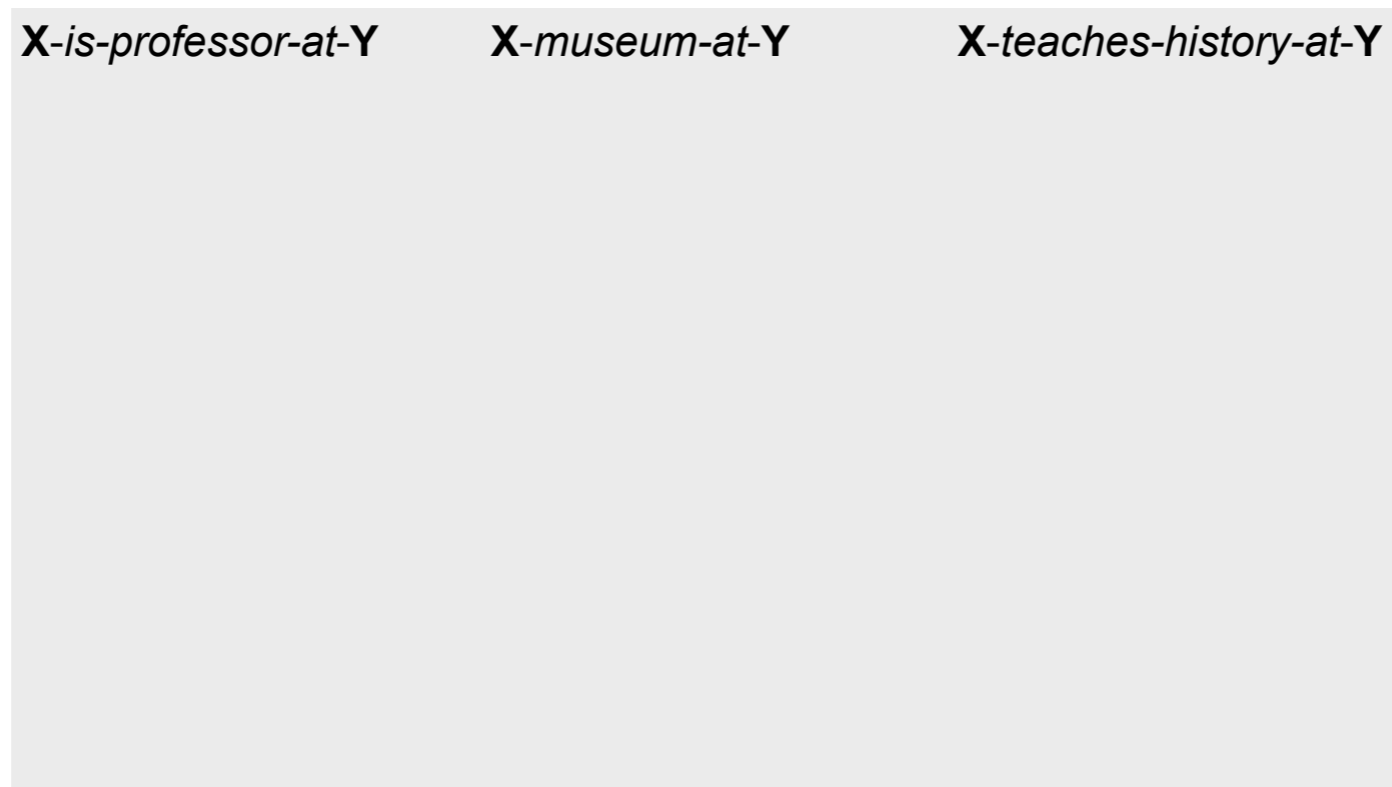
Relations have **entity type restriction**



X-is-professor-at-Y

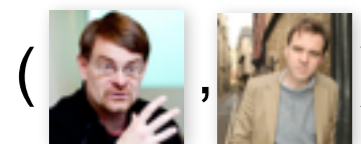
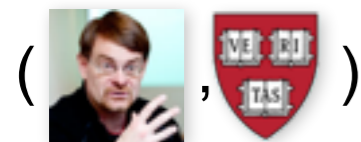
X-museum-at-Y

X-teaches-history-at-Y



Model E: Selectional Preferences

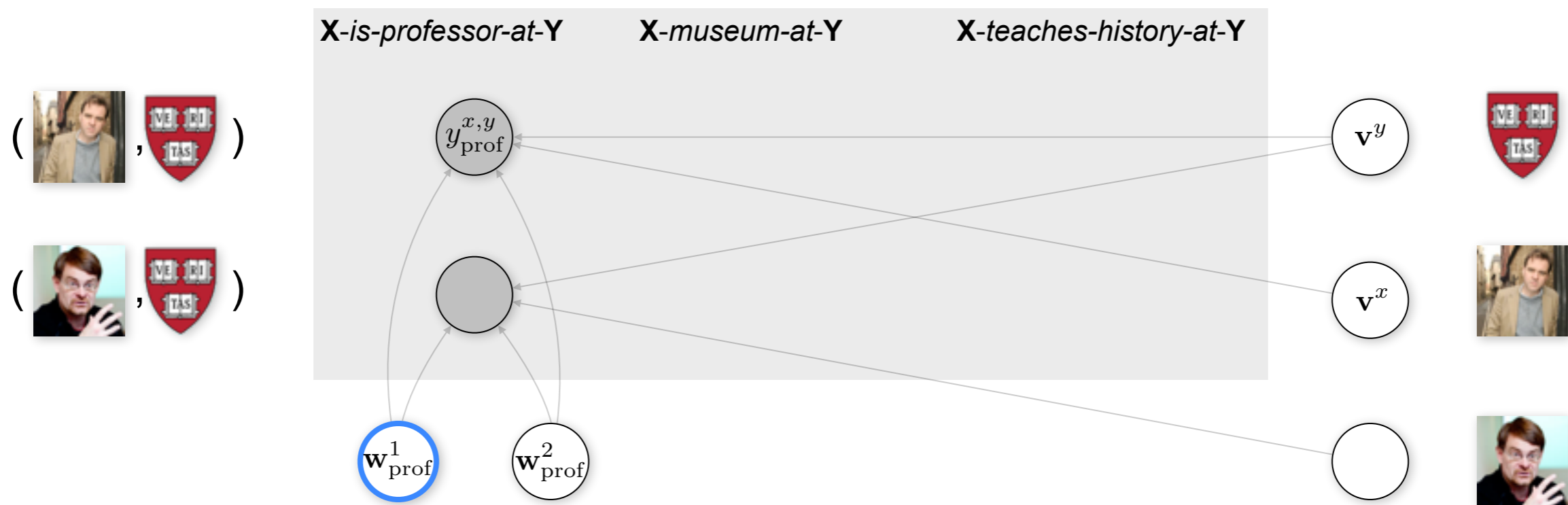
Relations have **entity type restriction**



<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>
0	0	0

Model E: Selectional Preferences

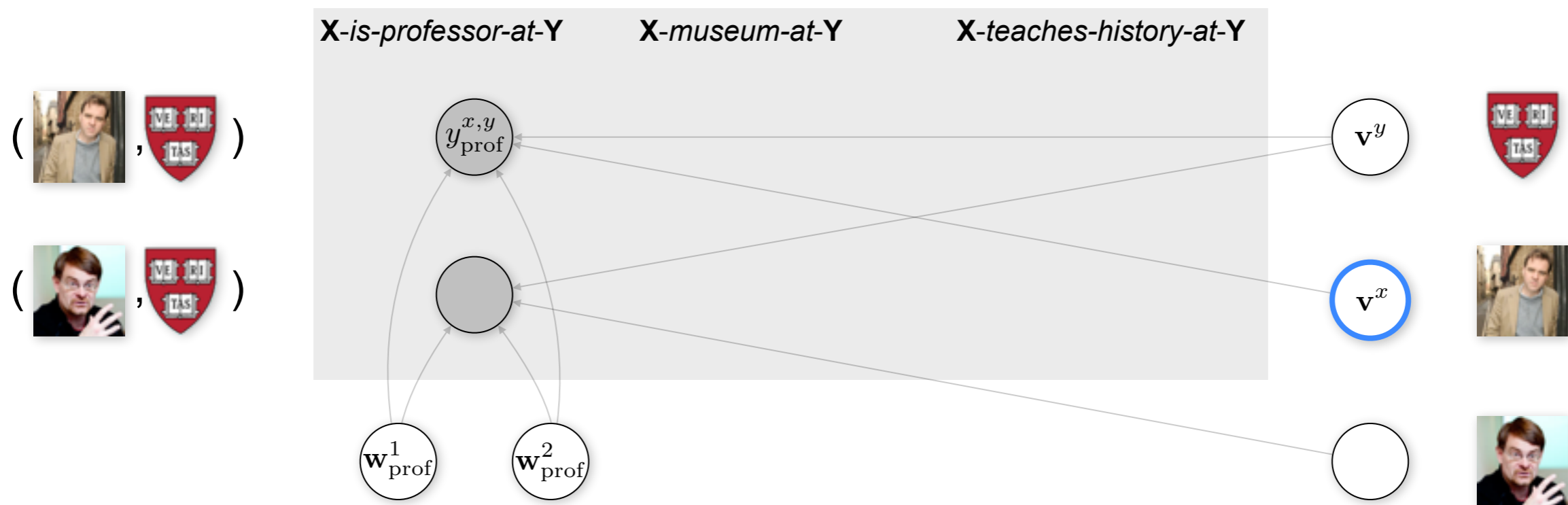
Argument Slot 1 weight vector ...



$$p(y_{\text{prof}}^{x,y} = 1 | \dots) \propto \exp[\langle \mathbf{v}^x, \mathbf{w}_{\text{prof}}^1 \rangle + \langle \mathbf{v}^y, \mathbf{w}_{\text{prof}}^2 \rangle]$$

Model E: Selectional Preferences

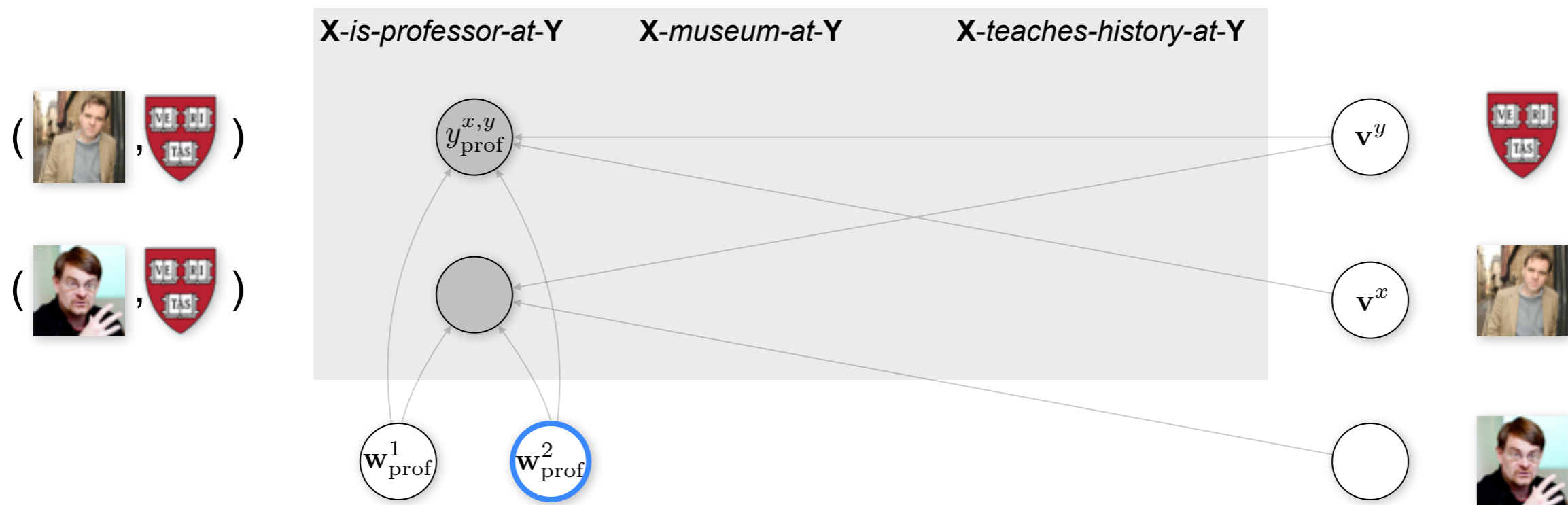
... dot-product with feature vector of **entity 1**



$$p(y_{prof}^{x,y} = 1 | \dots) \propto \exp[\langle \underline{v}^x, w_{prof}^1 \rangle + \langle v^y, w_{prof}^2 \rangle]$$

Model E: Selectional Preferences

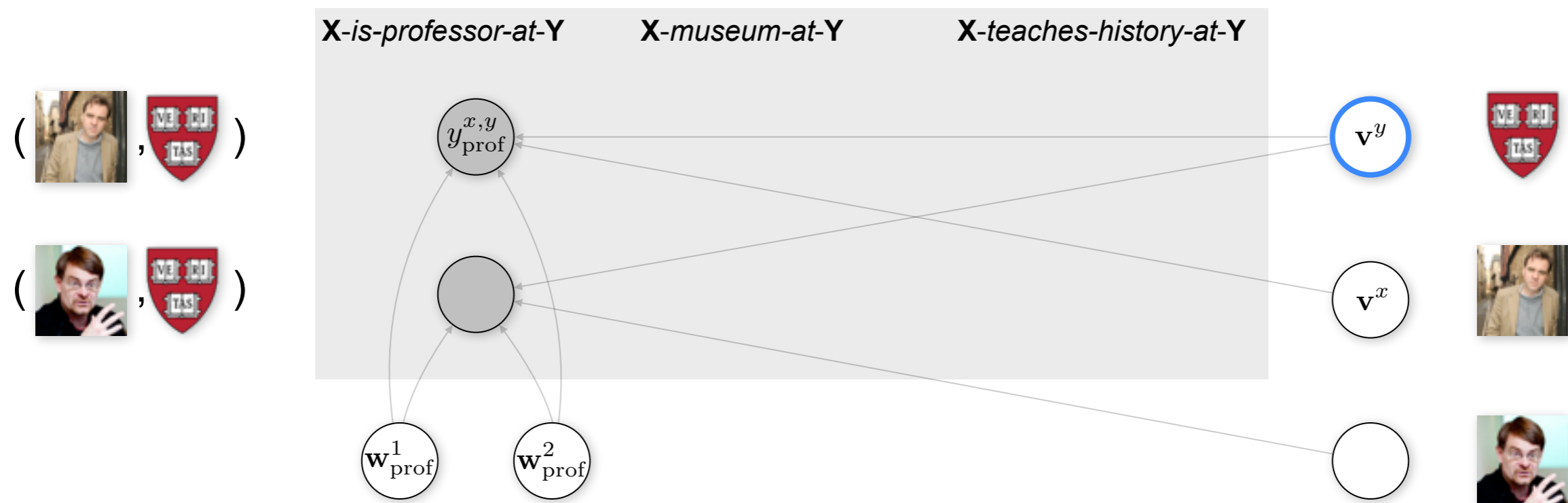
Argument Slot 2 weight vector ...



$$p(y_{\text{prof}}^{x,y} = 1 | \dots) \propto \exp[\langle \mathbf{v}^x, \mathbf{w}_{\text{prof}}^1 \rangle + \langle \mathbf{v}^y, \mathbf{w}_{\text{prof}}^2 \rangle]$$

Model E: Selectional Preferences

... dot-product with feature vector of **entity 2**



$$p(y_{\text{prof}}^{x,y} = 1 | \dots) \propto \exp[\langle \mathbf{v}^x, \mathbf{w}_{\text{prof}}^1 \rangle + \langle \underline{\mathbf{v}}^y, \mathbf{w}_{\text{prof}}^2 \rangle]$$

Combinations

models capture different aspects of the data, combine them (e.g., NF)

$$p(y_{\text{emp}}^{x,y} = 1 | \dots) \propto \exp[\langle \mathbf{f}_{\text{emp}}^{x,y}, \mathbf{w}_{\text{emp}}^{\text{N}} \rangle + \langle \mathbf{v}^{x,y}, \mathbf{w}_{\text{emp}}^{\text{F}} \rangle]$$

Evaluation (Structured)

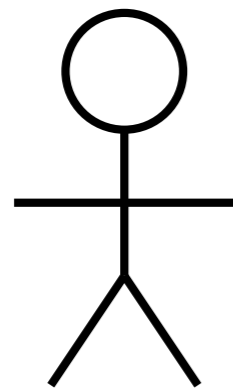
[Mintz 09; Yao 11; Surdenau 12]

Evaluate **average precision** per **Freebase** relation.

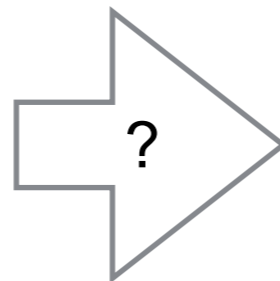
Relation	MI09	YA11	SU12	N+F+E
<i>employee</i>	0.67	0.64	0.7	0.79
<i>containedby</i>	0.48	0.51	0.54	0.69
<i>parents</i>	0.39			0.39
...
Weighted MAP	0.48	0.52	0.57	0.69
MAP	0.32	0.42	0.56	0.63

~45 minutes to train our models on 4000 relations, ~50k entity pairs

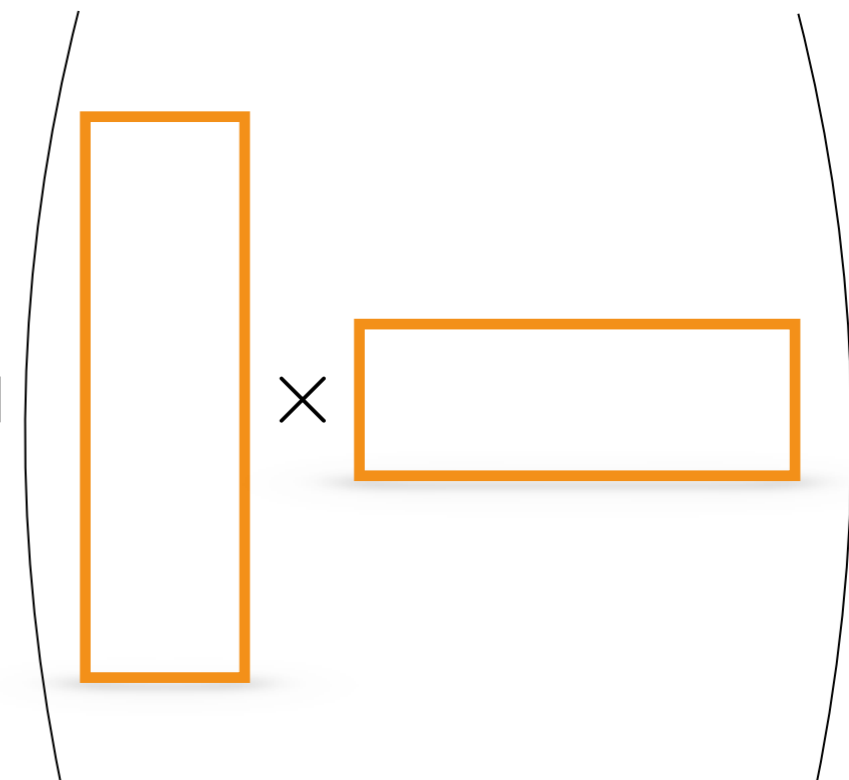
Injecting Knowledge



“lecturers are employees!”



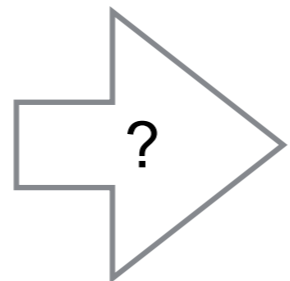
sigmoid



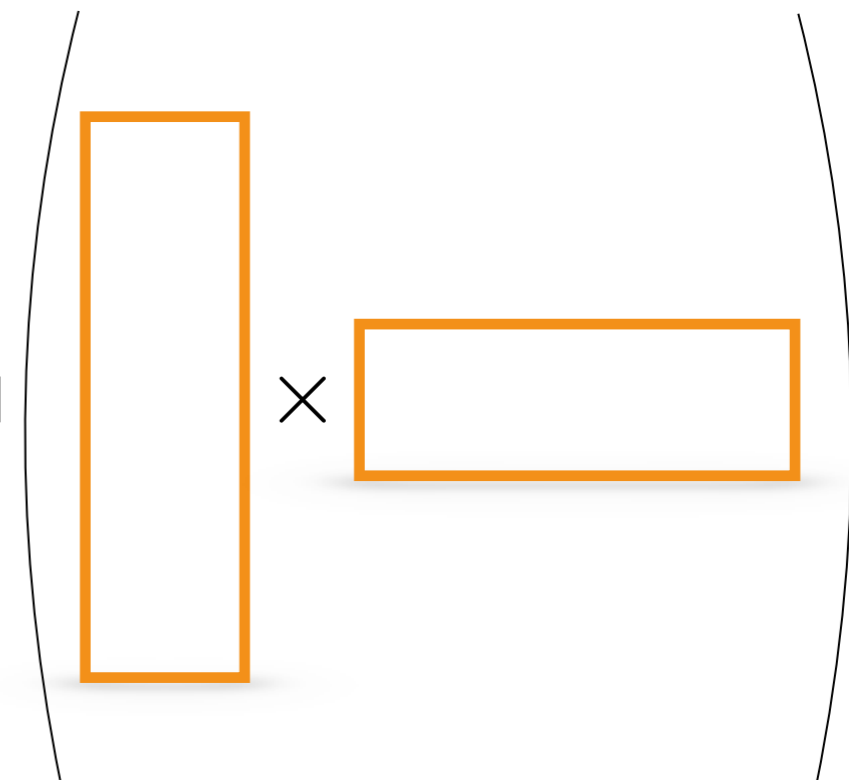
Injecting Knowledge



“a liquid turns into a solid when its temperature is lowered below its freezing point

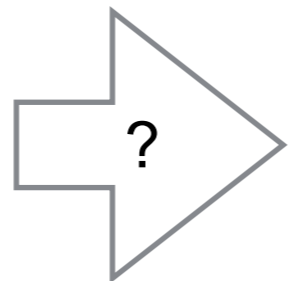


sigmoid

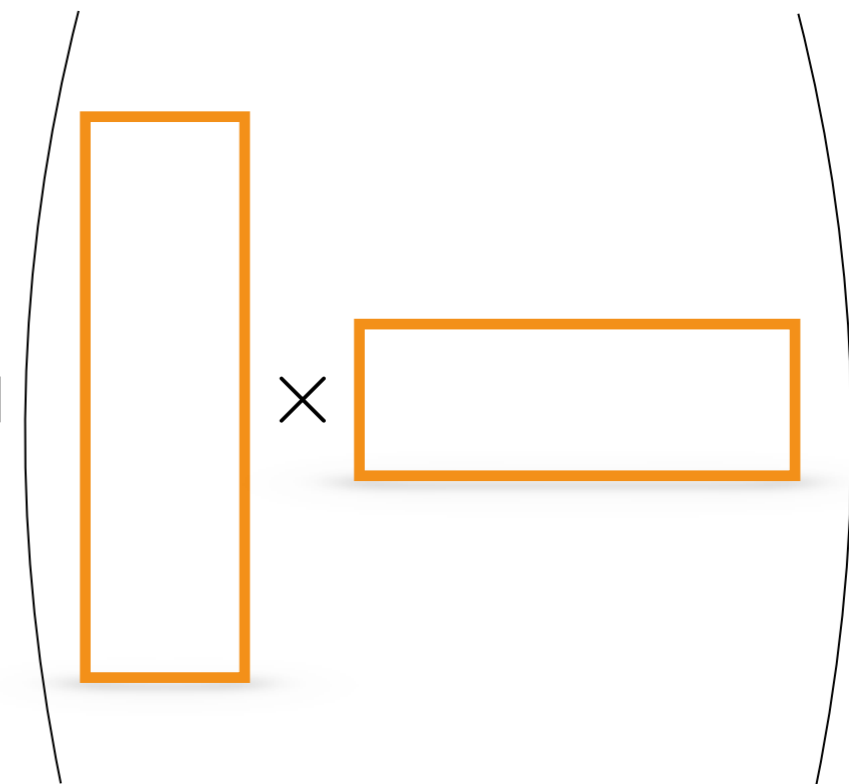


Injecting Knowledge: Rules

$\forall x,y: \text{birthplace}(x,y)$
 $\Rightarrow \text{bornIn}(x,y)$



sigmoid



Goal: Predict Unseen Cells ...

native-of 's birthplace bornIn livesIn

		1	1
1	1		
	1		
1			?

... By Using Rules and Data

native-of 's birthplace bornIn livesIn

		<u>1</u>	<u>1</u>
<u>1</u>	<u>1</u>		
	<u>1</u>	1	.8
<u>1</u>		1	.8

$\forall x,y: \text{birthplace}(x,y) \Rightarrow \text{bornIn}(x,y)$

Baselines

native-of 's birthplace *bornIn* *livesIn*

		1	1
1	1	1	1
	1	1	.8
1	.8	.8	.6

- Rules only
- Rules after learning
- Rules before learning

$$\forall x,y: \text{birthplace}(x,y) \Rightarrow \text{bornIn}(x,y)$$

Pre-Injection may not add data at all

native-of 's birthplace bornIn livesIn

1	1	1	1
1			
	1		
1			

- Rules before learning

$\forall x,y: \text{birthplace}(x,y) \Rightarrow \text{bornIn}(x,y)$

Pre-Injection may not add data at all

native-of 's birthplace bornIn livesIn

1			1
1			
	1		
1			

- Rules before learning

$$\forall x,y: \text{birthplace}(x,y) \Rightarrow \text{bornIn}(x,y)$$

Idea: Iterate

native-of 's birthplace bornIn livesIn

1	.8	1	1
1	.8	1	
	1		
1			

- Inference with model
- Apply rules

... and learn again

$\forall x,y: \text{birthplace}(x,y) \Rightarrow \text{bornIn}(x,y)$

Our approach

[Rocktaeschel et al 15]

- ▶ **Directly optimise to fulfil formulae in expectation**
- ▶ formulae have **compositional expectations**

$$E_{\mathbf{v}, \mathbf{w}}[\textit{birthplace}(\textit{Seb}, \textit{HH})] = E_{\mathbf{v}, \mathbf{w}}[y_{\textit{birthplace}}^{\textit{Seb}, \textit{HH}}] = \textit{sigm}(\langle \mathbf{v}^{\textit{Seb}, \textit{HH}}, \mathbf{w}_{\textit{birthplace}} \rangle)$$

$$E_{\mathbf{v}, \mathbf{w}}[r(X_1, X_2)] = \textit{sigm}(\langle \mathbf{v}^{X_1, X_2}, \mathbf{w}_r \rangle)$$

$$E_{\mathbf{v}, \mathbf{w}}[A \wedge B] = E_{\mathbf{v}, \mathbf{w}}[A] \times E_{\mathbf{v}, \mathbf{w}}[B]$$

$$E_{\mathbf{v}, \mathbf{w}}[\neg A] = 1 - E_{\mathbf{v}, \mathbf{w}}[A]$$

$$E_{\mathbf{v}, \mathbf{w}}[A \Rightarrow B] = 1 - (E_{\mathbf{v}, \mathbf{w}}[A] \times (1 - E_{\mathbf{v}, \mathbf{w}}[B]))$$

Our approach

[Rocktaeschel et al 15]

- ▶ **Directly** optimise to **fulfil formulae in expectation**
- ▶ formulae have **compositional expectations**
- ▶ **quantification** through **grounding**

$$\begin{aligned} E_{\mathbf{v},\mathbf{w}}[\forall x.f(x)] &= E_{\mathbf{v},\mathbf{w}}[f(X_1) \wedge \dots \wedge f(X_n)] \\ &= E_{\mathbf{v},\mathbf{w}}[f(X_1)] \times \dots \times E_{\mathbf{v},\mathbf{w}}[f(X_n)] \end{aligned}$$

General Framework

[Rocktaeschel et al 15]

- ▶ Find embeddings \mathbf{v} and \mathbf{w} that...
- ▶ Maximize **log expectation** of a set of formulae f

$$\arg \max_{\mathbf{v}, \mathbf{w}} \sum_f \log(E_{\mathbf{v}, \mathbf{w}}[f])$$

- ▶ **Generalizes** regular (binary) matrix factorization with logistic loss
- ▶ Get gradients by **back-propagation** through $\log(E[.])$ tree
- ▶ Optimize via **SGD** / Adagrad etc.

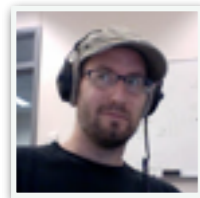
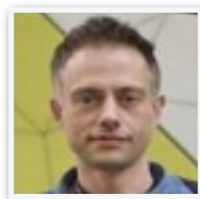
Experiments

[Rocktaeschel et al 15]

- ▶ “Zero-shot” learning
 - ▶ Given: a lot of surface form data, but **no Freebase relations**
 - ▶ Goal: given few (36) Freebase rules, learn to Freebase relations

Experiments: Zero-Shot Learning

Remove Freebase data from training set ...



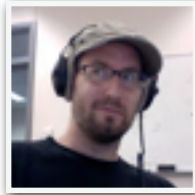







<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
1	1		1
1			
1		1	

Experiments: Zero-Shot Learning

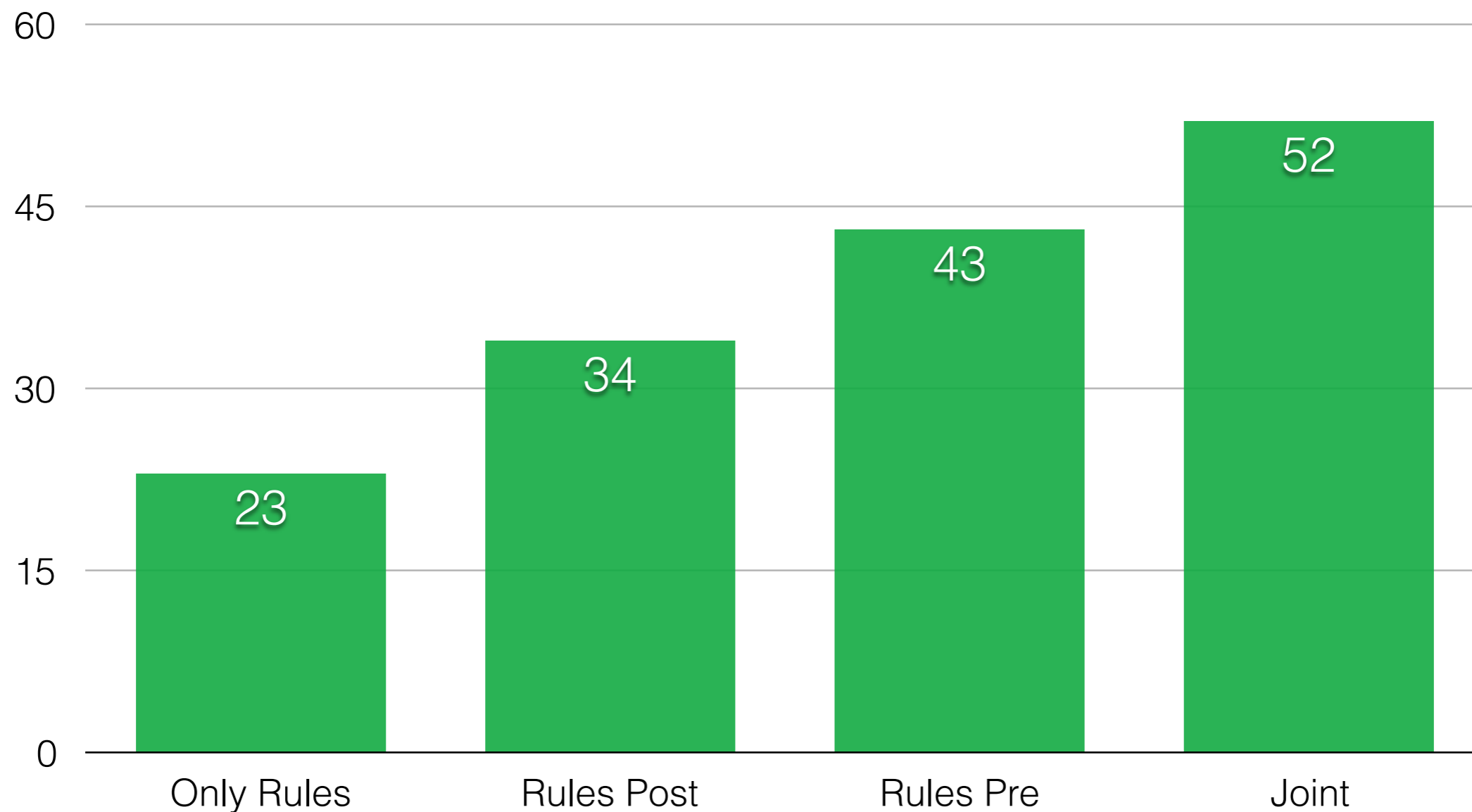
and learn only from surface form relations, and rules



	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
 	1	1		
 	1			
 	1		1	
 				

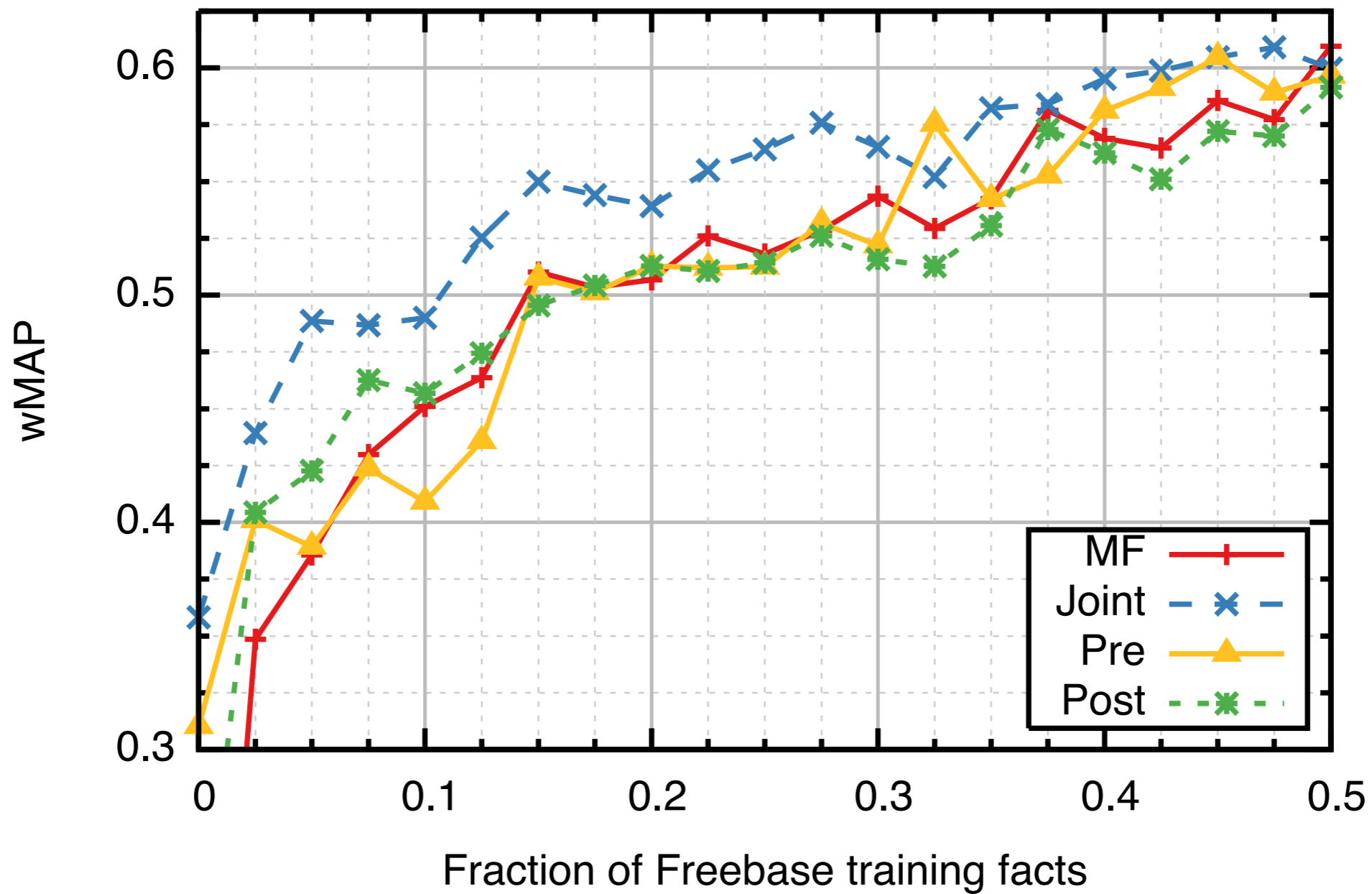
Zero-Shot Learning Results (MAP)

[Rocktaeschel et al 15]



Learning Curve

[Rocktaeschel et al 15]



Generating Data?

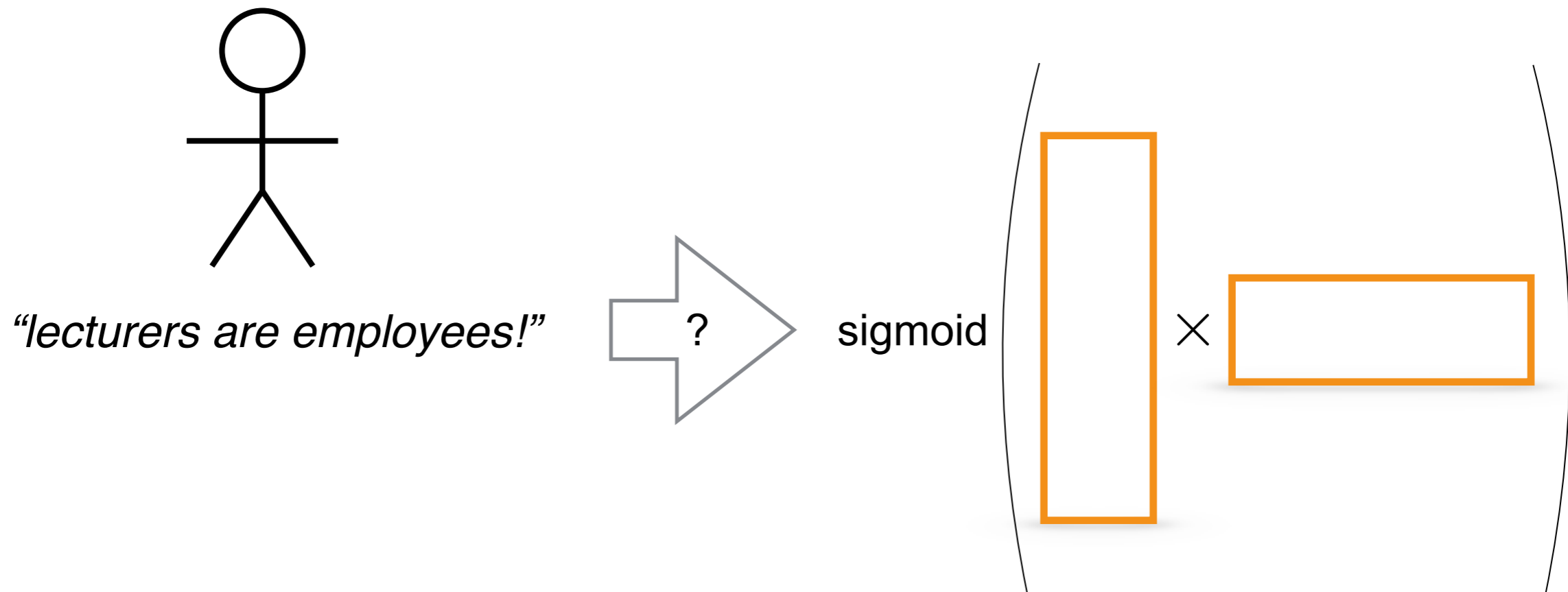
native-of 's birthplace bornIn livesIn

		1	1
1	1		
	1		?
1			
	1	1	

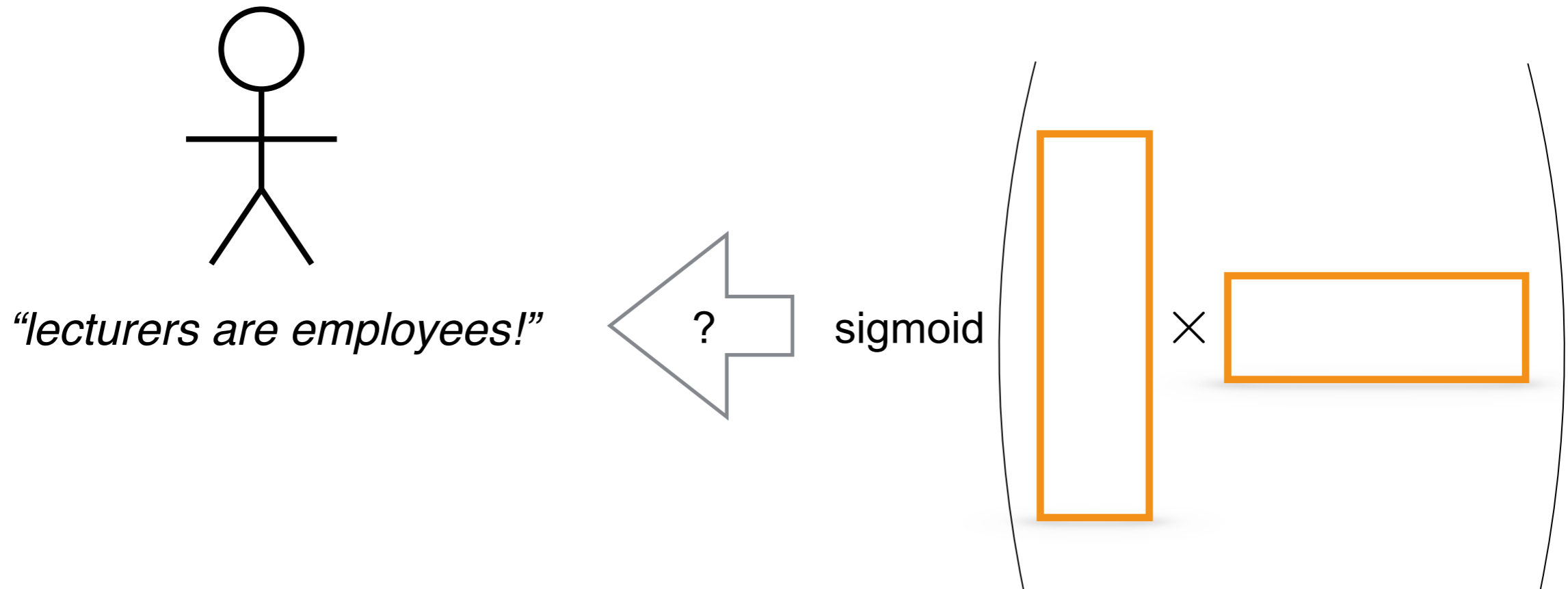
Hasn't worked yet

- ▶ **Row embeddings overtrain**
- ▶ **At test time premise appears with other relations**

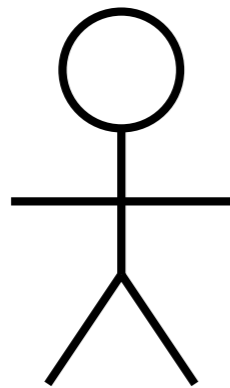
Challenge 1: Injecting Symbolic Rules



Challenge 2: Extracting Explanations

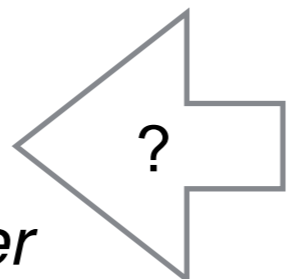


Challenge 2: Extracting Explanations

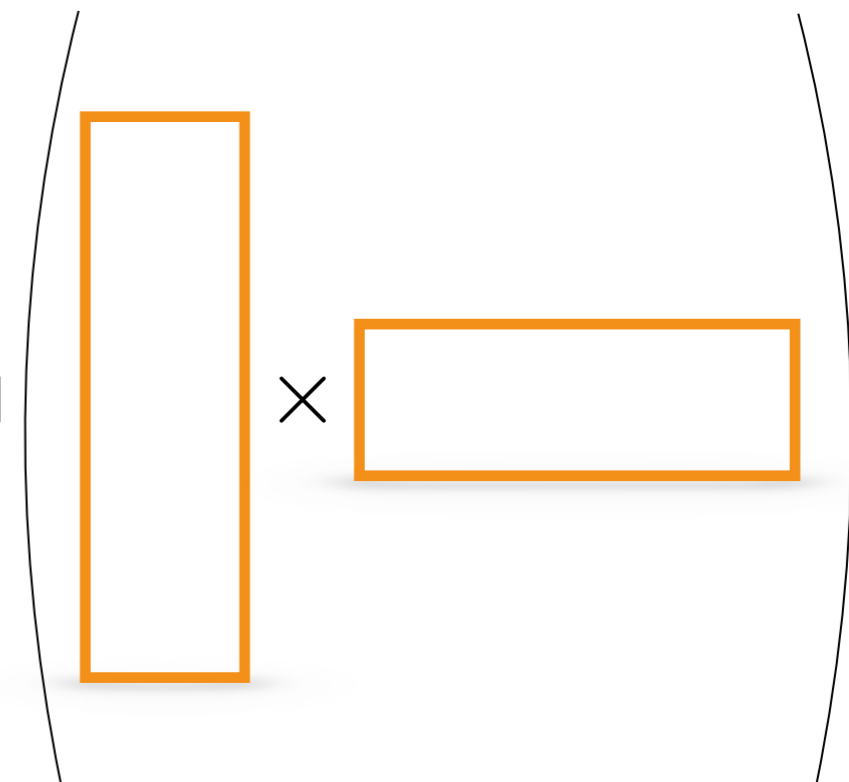


“I returned Sebastian because we know he is a lecturer at UCL, which is in London, so he most likely lives in London

...



sigmoid



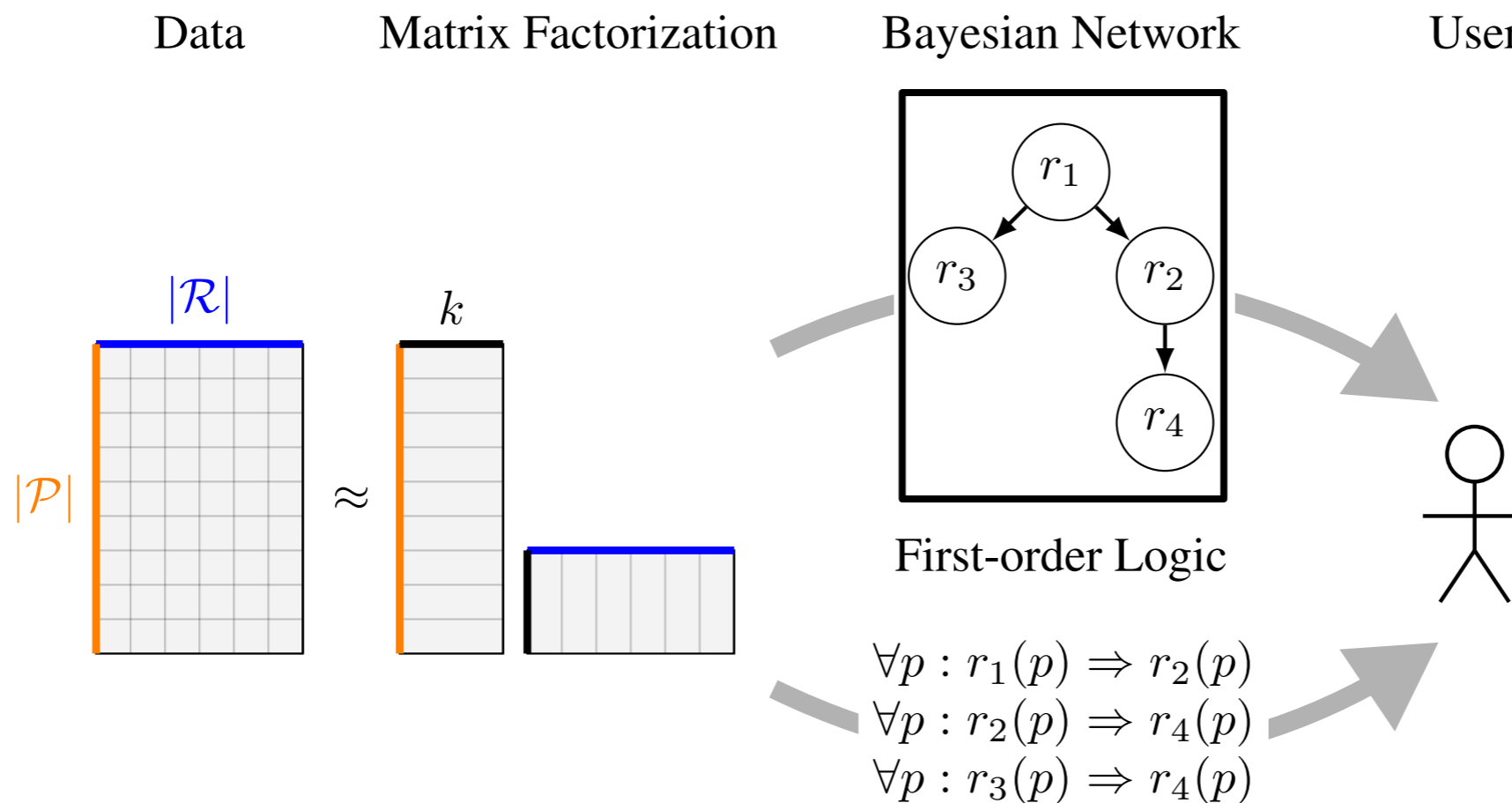
[Thrun 1995, NIPS, Craven 1996, NIPS]

“Knowledge Extraction”

- ▶ Learn a more interpretable model from distributed representations (for interpretation, not for use)
 - ▶ Neural Networks => **if-then rules** (Thrun, 95)
 - ▶ Neural Networks => **Decision Trees** (Craven, 96)
- ▶ Open Questions
 - ▶ Go **beyond classification**: joint models
 - ▶ Use to provide **proofs** of complex predictions
 - ▶ Integrate into a **dialog** between human and machine

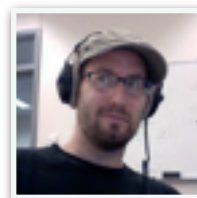
Explaining Matrix Factorization

[Sanchez et al. 2015, KRR]



Extracting Bayesian Networks

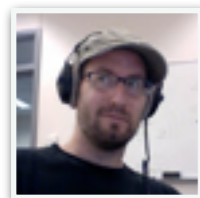
Learn Embeddings from Data



<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
1	1		1
1			
1		1	

Extracting Bayesian Networks

Generate data from embeddings (threshold or sample)

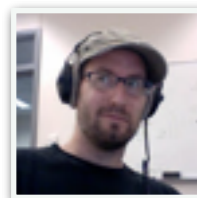


<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
1	1	0	1
1	0	0	1
1	0	1	1

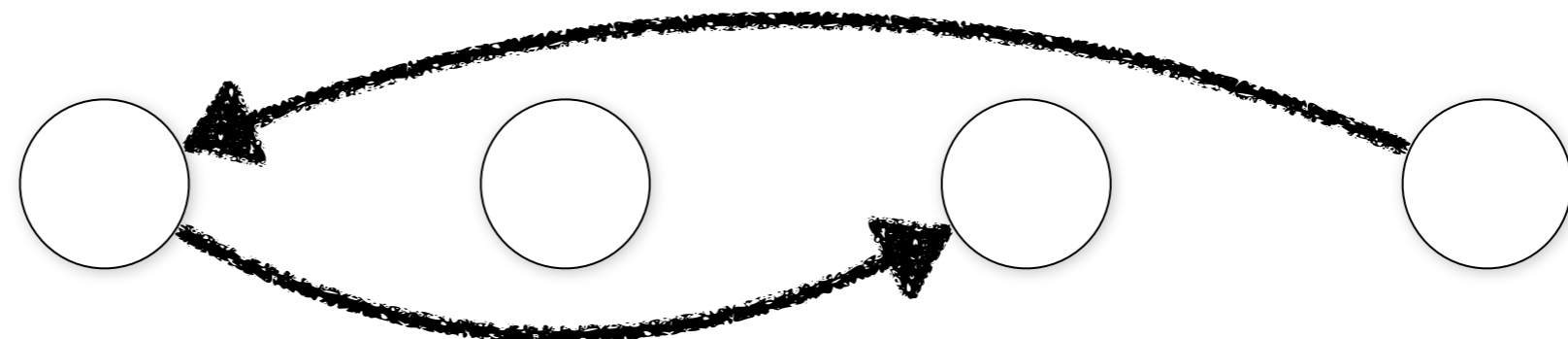


Extracting Bayesian Networks

Learning a tree shaped Bayesian Network



<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
1	1	0	1
1	0	0	1
1	0	1	1

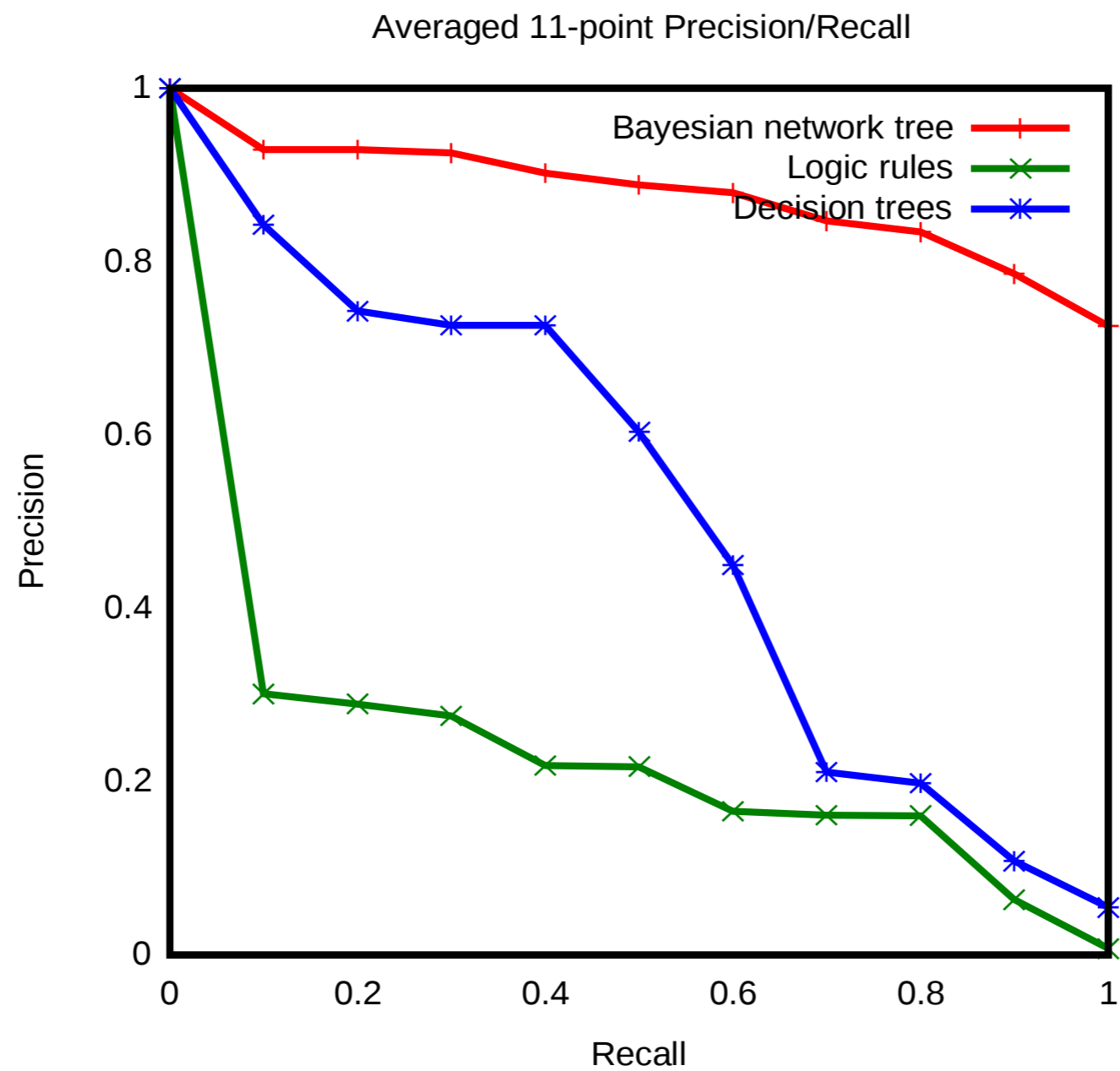


Benefits of Bayesian Network Trees

[Sanchez et al. 2015, KRR]

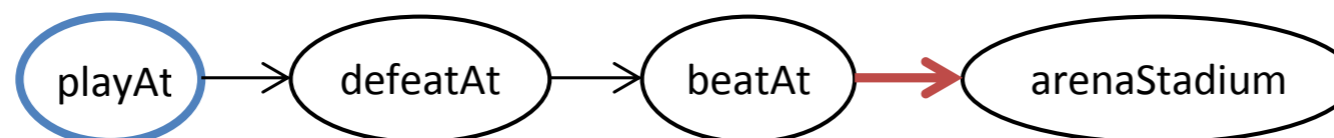
- ▶ Provide a **joint model** over all relations
 - ▶ more **compact** (than one decision tree per relation)
 - ▶ more **faithful** to the **joint** MF model
- ▶ **Probabilistic** interpretation, captures probabilistic nature of MF
- ▶ **Very scalable**
 - ▶ Learning: **Prim Algorithm** to find Maximum Spanning Tree over mutual information
 - ▶ Inference: Belief Propagation in **non-loopy** graph

Faithfulness



A “Proof”

- ▶ Model observed *playAt*(**Eagles**, **Canton**)
- ▶ Model wrongly predicted *arenaStadium*(**Eagles**, **Canton**)
- ▶ The Bayesian Network can provide this “proof”



- ▶ Todo: evaluate this in a downstream “debugging” task

Summary

- ▶ Do semantics in a **probabilistic relational reasoner**
- ▶ Reasoner: **matrix/tensor factorization** (or other LV models)
- ▶ Models itself don't need to be interpretable if we know ...
- ▶ ... how to **Interact with uninterpretable models**
 - ▶ **inject** explanations and logical rules
 - ▶ Approach: **optimize** embeddings to fulfil formulae
 - ▶ **extract** explanations
 - ▶ for example: by using an interpretable BN **proxy** model

Thanks

Training

Negative Data

Usually **unavailable** or **sparse**, so...

<i>X-is-historian-at-Y</i>	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
	1	1		1
1	1			
	1		1	
1				

Negative Data

...subsample, which can work...

<i>X-is-historian-at-Y</i>	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
	1	1		1
1	1	0		
	1		1	
1				

Negative Data

but often **does not**

<i>X-is-historian-at-Y</i>	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
	1	1		1
1	1			
0	1		1	
1				

Negative Data

and you need to sample a lot (wasting resources)

<i>X-is-historian-at-Y</i>	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
0	1	1	0	1
1	1	0		0
0	1	0	1	
1	0		0	

Implicit Feedback

Often users only **click/view/buy** items, or not, but **no rating**

User 1	User 2	User 3	User 4	User 5	
	1	1		1	Item 1
1	1				Item 2
	1		1		Item 3
1					Item 4

Ranking

[Rendle et al.,09]


for all (observed, not observed) pairs in column: $\text{prob}(o) > \text{prob}(n)$

<i>X-is-historian-at-Y</i>	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
	1	1		1
0.9	1			
0.95	1		1	
1				

Ranking

[Rendle et al.,09]

for all (observed, not observed) pairs in a column: $\text{prob}(o) > \text{prob}(n)$

<i>X-is-historian-at-Y</i>	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
	1	1		1
0.9	1			
	1		1	
0.85				
1				

Training: Stochastic Gradient Descent

[Rendle et al.,09]

Sample observed fact...

<i>X-is-historian-at-Y</i>	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
	1	1		1
1	1			
	1		1	
1				

Training: Stochastic Gradient Descent

[Rendle et al.,09]

Sample unobserved cell for same relation

<i>X-is-historian-at-Y</i>	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
	1	1		1
1	1			
	1		1	
1				

Training: Stochastic Gradient Descent

[Rendle et al.,09]

Estimate current beliefs and gradient, **update** parameters accordingly

<i>X-is-historian-at-Y</i>	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
	0.8	1		1
1	1			
	1		1	
1	0.9			

Training: Stochastic Gradient Descent

[Rendle et al.,09]

Estimate current beliefs and gradient, **update** parameters accordingly

<i>X-is-historian-at-Y</i>	<i>X-is-professor-at-Y</i>	<i>X-museum-at-Y</i>	<i>X-teaches-history-at-Y</i>	<i>employee(X,Y)</i>
	0.85	1		1
1	1			
	1		1	
1	0.8			

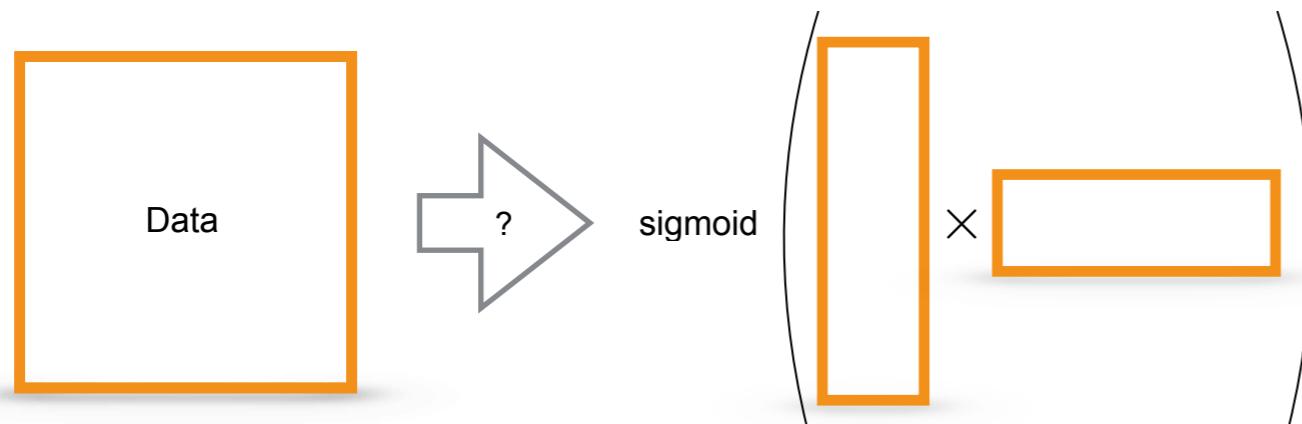
How can we do this?

native-of 's birthplace bornIn livesIn

		1	1
1	1		
		1	
1			?

birthplace(x,y) => bornIn(x,y)

Overview: Embeddings and ...



► **Learning from Data**
[NAACL 2013]

"lecturers are employees!"

sigmoid

"Talking to Uninterpretable Models"

► **Injecting Knowledge**
[SP 2014, NAACL 2015]

"lecturers are employees?"

sigmoid

► **Extracting Knowledge**
[KRR 2015]