

# Use of a possibilistic logic-like framework for statistical relational learning

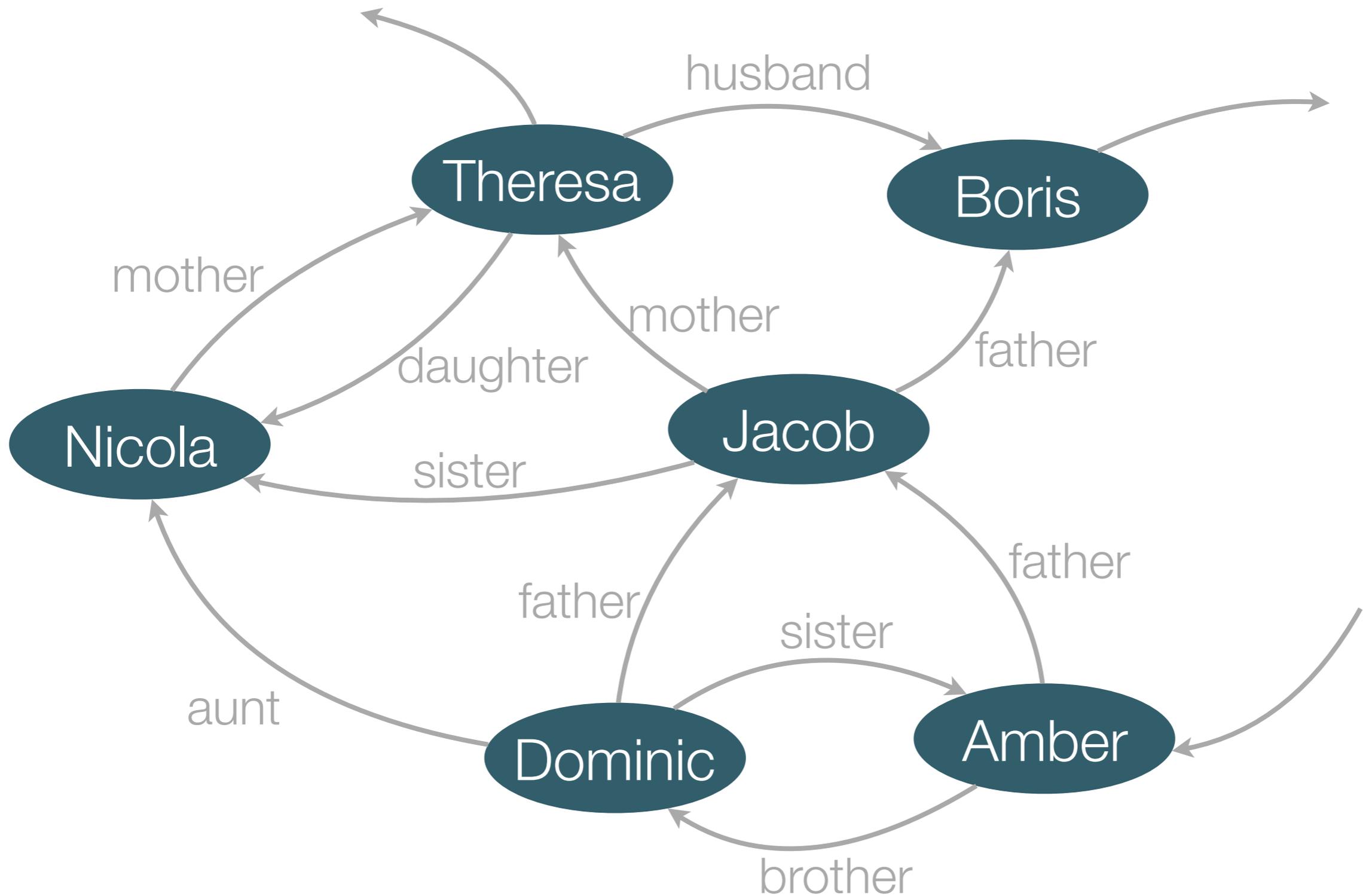
Steven Schockaert

Joint work with: Ondrej Kuzelka, Jesse Davis, Martin Svatos

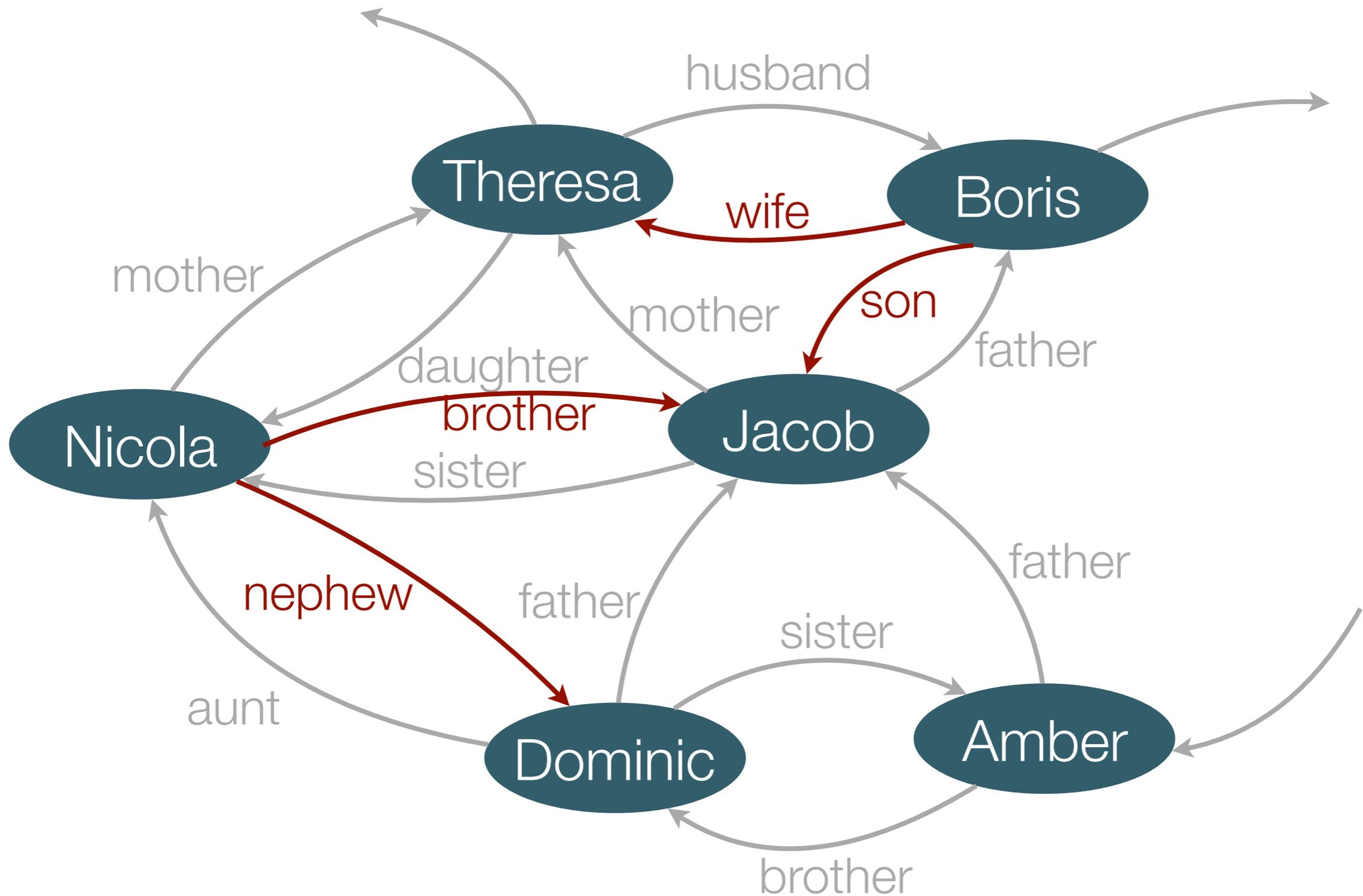
School of Computer Science & Informatics  
Cardiff University, Cardiff, UK  
SchockaertS1@cardiff.ac.uk  
<http://users.cs.cf.ac.uk/S.Schockaert>



# Predicting missing relational facts



# Predicting missing relational facts



# Knowledge graph embedding

Learn vector representation for each entity and scoring function for each relation such that:

$$P(e, r, f) = \text{score}_r(\mathbf{e}, \mathbf{f})$$

Drawbacks:

- ▶ Poor performance if “test graph” is sparse
- ▶ Only able to capture relatively simple dependencies
- ▶ Black box

Advantages:

- ▶ Scalable (millions of entities)
- ▶ Easy to incorporate text descriptions, images, etc.
- ▶ Only able to capture relatively simple dependencies

# Markov logic networks

---

0	<i>FacultyEmeritus(V1)</i>
-5.64478	<i>AdvisedBy(V1, V2)</i>
-5.10868	<i>Student(V1)</i>
-2.06003	<i>FacultyAffiliate(V1)</i>
-3.18786	<i>PostQuals(V1)</i>
-3.64133	<i>PostGenerals(V1)</i>
-0.80596	<i>FacultyAdjunct(V1)</i>
-0.0391718	<i>TempAdvisedBy(V1, V1)</i>
-4.12274	<i>Publication(V1, V2)</i>
-0.0551981	<i>AdvisedBy(V1, V1)</i>
-2.95067	<i>Faculty(V1)</i>
-6.56008	<i>ta(V1, V2, V3)</i>
-4.90068	<i>ProjectMember(V1, V2)</i>
-4.44769	<i>TempAdvisedBy(V1, V2)</i>
5.10868	<i>Professor(V1)</i>
-6.31879	<i>TaughtBy(V1, V2, V3)</i>
-4.69704	<i>PreQuals(V1)</i>
1.67999e-06	$\neg TaughtBy(V1, V2, V3) \vee \neg FacultyAdjunct(V2) \vee \neg TempAdvisedBy(V2, V4)$
1.08275	$\neg TaughtBy(V1, V2, V3) \vee \neg FacultyAdjunct(V2) \vee \neg Publication(V4, V2)$
6.07986	$Faculty(V1) \vee \neg TempAdvisedBy(V2, V1) \vee FacultyAffiliate(V1)$
2.48337	$PreQuals(V1) \vee PostQuals(V1) \vee ProjectMember(V2, V1) \vee \neg TempAdvisedBy(V1, V3)$
9.33506	$PreQuals(V1) \vee PostQuals(V1) \vee \neg AdvisedBy(V1, V2) \vee PostGenerals(V1)$
-12.0071	$\neg Student(V1) \vee FacultyAdjunct(V1) \vee Professor(V1) \vee Faculty(V1) \vee FacultyAffiliate(V1)$
-2.05975	$PostQuals(V1) \vee FacultyAdjunct(V1) \vee \neg Faculty(V1) \vee \neg AdvisedBy(V2, V1) \vee FacultyAffiliate(V1)$

---

# Markov logic networks

Learn domain knowledge in the form of a weighted set of rules  $w:F$

$$p_{\mathcal{M}}(\omega) = \frac{1}{Z} \exp \left( \sum_{w:F \in \mathcal{M}} w n_F(\omega) \right)$$

normalisation constant

number of groundings of  $F$  satisfied in  $\omega$

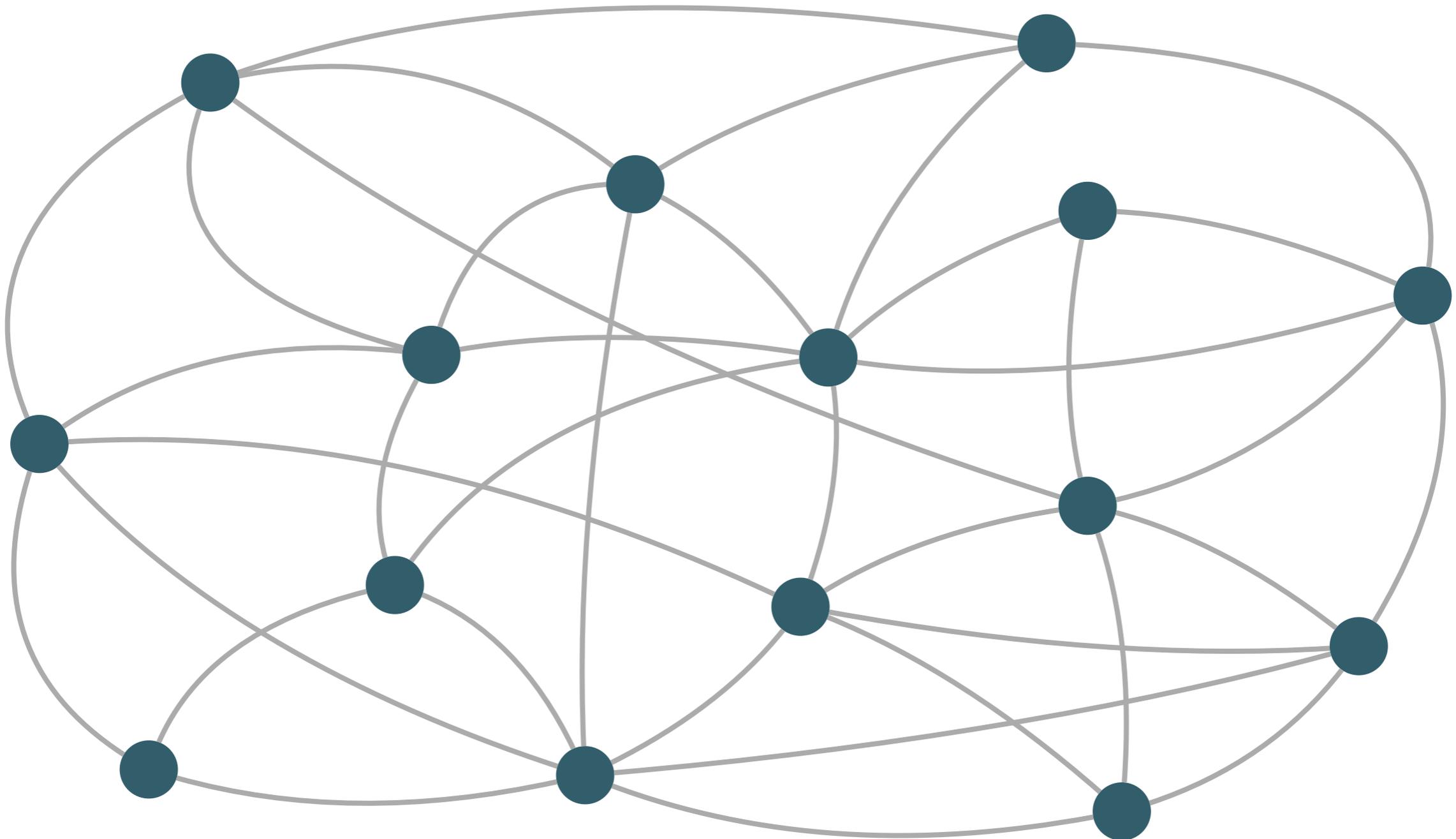
Given some set of evidence literals  $E$ , find the most probable world  $\omega$  satisfying all literals in  $E$ ; predict all literals satisfied in  $\omega$  as true

Drawbacks:

- ▶ Weights dependent on domain size
- ▶ Interaction of the rules hard to interpret
- ▶ Computationally rather expensive

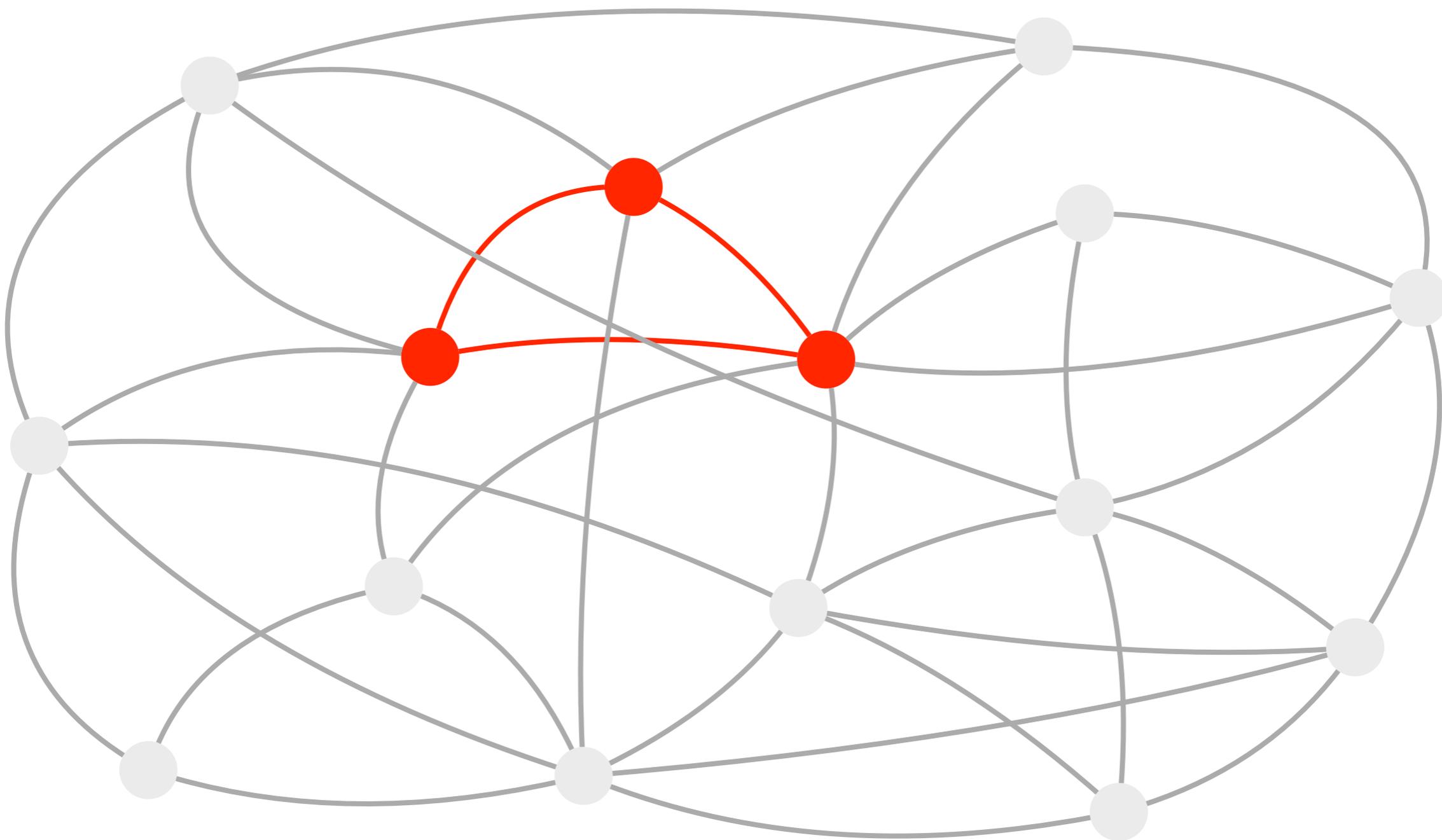
# Relational marginals

We only have one relational graph, i.e. one training example. How can we estimate probabilities from one example? What does the probability of a (universally quantified) clause represent?



# Relational marginals

Relational marginal distribution  $\approx$  distribution over size-k subgraphs



# Relational marginals

Relational marginal distribution  $\approx$  distribution over size-k subgraphs

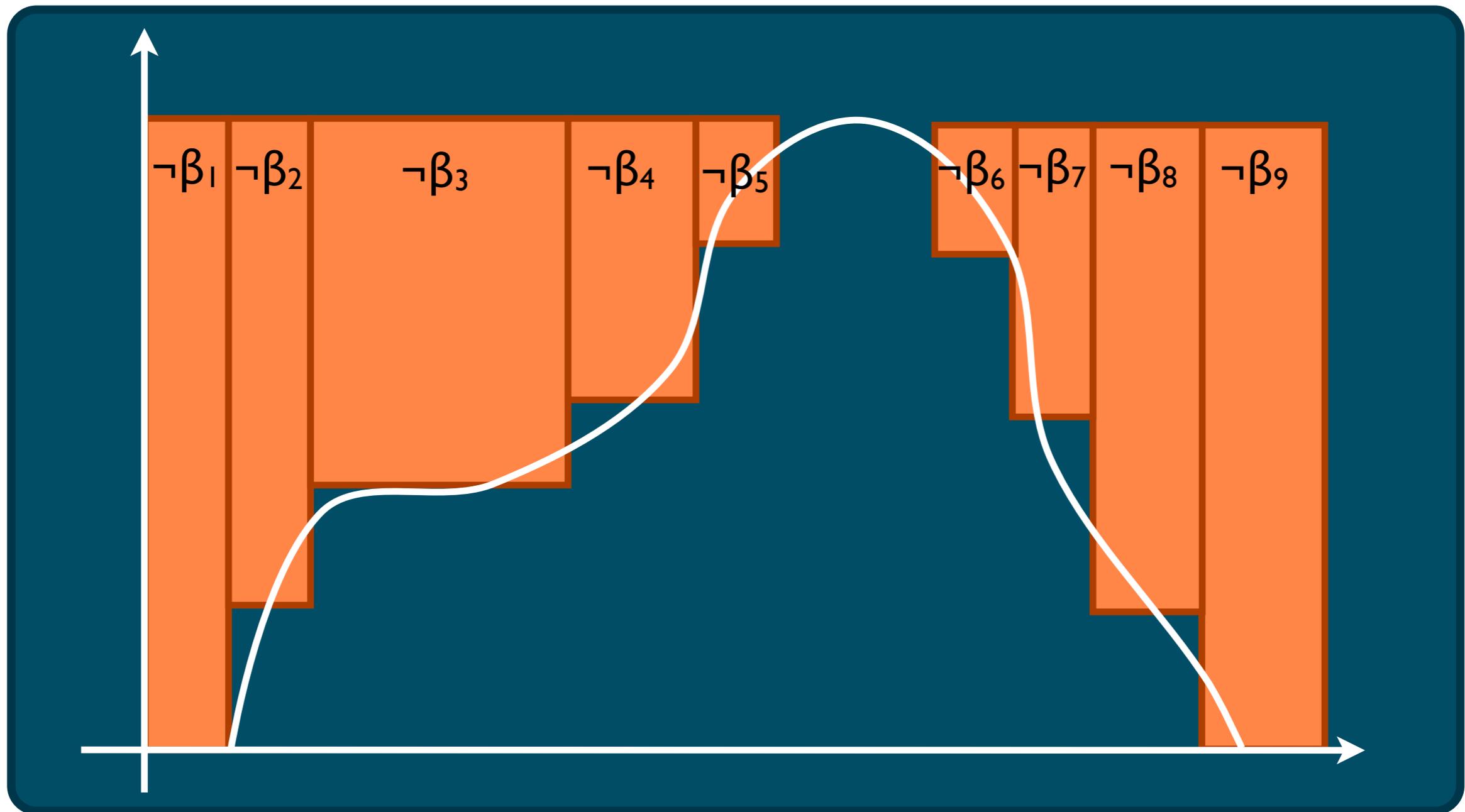
Probability of a clause given by the relational marginal distribution:

$$P_{\Upsilon, k}(\alpha) = \sum_{\omega: \omega \models \alpha} P_{\Upsilon, k}(\omega)$$

**Main idea:** encode the relational marginal distribution using possibilistic logic

- ▶ Needs first-order formulas, but these are merely templates for sets of ground formulas
- ▶ Exact encoding exponentially large, so need heuristic method for selecting set of rules

# Possibilistic encoding of probability distributions



Possibilistic logic:  $\pi(\omega) = \min\{1 - \lambda : \omega \not\models \beta, (\beta, \lambda) \in \Theta\}$

# Relational possibilistic logic

$\text{son}(X, Y) \wedge \text{male}(Y) \rightarrow \text{father}(Y, X)$  1

---

$\text{son}(X, Y) \wedge \text{son}(Z, Y) \rightarrow \text{sibling}(X, Z)$  0.5

---

$\text{son}(X, Y) \wedge \text{married}(Y, Z) \rightarrow \text{son}(X, Z)$  0.25

# Relational possibilistic logic

$$E = \{\text{son}(a, b), \text{married}(b, c)\}$$

$$\text{son}(a, b) \wedge \text{male}(b) \rightarrow \text{father}(b, a)$$

...

1

---

$$\text{son}(c, b) \wedge \text{male}(b) \rightarrow \text{father}(b, c)$$

$$\text{son}(a, b) \wedge \text{son}(c, b) \rightarrow \text{sibling}(a, c)$$

...

0.5

$$\text{son}(c, b) \wedge \text{son}(a, b) \rightarrow \text{sibling}(c, a)$$

---

$$\text{son}(a, b) \wedge \text{married}(b, c) \rightarrow \text{son}(a, c)$$

...

0.25

$$\text{son}(c, b) \wedge \text{married}(b, a) \rightarrow \text{son}(c, a)$$

# Rule learning

Learn a set of hard constraints

- ▶ Exhaustive search over all small clauses
- ▶ Keep all clauses which are never violated in the given knowledge graph

Construct training examples for rule learner

- ▶ **Positive examples:** literals from the knowledge graph
- ▶ **Negative examples:** literals which are missing from the knowledge graph but which are consistent with the learned constraints

Rule learner

- ▶ Use beam search to find a rule which maximises accuracy
- ▶ Repeat process several times for each predicate

# Weight learning

Main principle for weight learning is maximising the likelihood of a given set of sampled sub-graphs

Assume ordering  $\alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_n$  is given

Number of worlds whose probability is  $1 - \lambda_i$ '

Maximize:  $\prod_{\omega \in E} P(\omega) = \prod_{i=1}^n (1 - \lambda'_i)^{|\mathcal{E}_{i+1}| - |\mathcal{E}_i|}$

where  $\mathcal{E}_i = \{\omega \in \mathcal{E} \mid \omega \models \alpha_i \wedge \dots \wedge \alpha_n\}$ .

Subject to:

$$\lambda'_1 \leq \lambda'_2 \leq \dots \leq \lambda'_n$$

Weights to be learned have to satisfy the given ordering

Ensure probabilities sum to 1

$$\sum_{i=1}^k (1 - \lambda'_i) \cdot (|M_{i+1}| - |M_i|) = 1$$

Number of models of  $\alpha_i \wedge \dots \wedge \alpha_n$

Can be solved using geometric programming

Ordering of the rules is done in a greedy fashion

# Learned theories

---

... (116 hard constraints not shown here)

$(Publication(V0, V1) \leftarrow TempAdvisedBy(V3, V1) \wedge Publication(V0, V3), \lambda_{16})$

$(PostQuals(V0) \leftarrow TempAdvisedBy(V0, V3) \wedge TempAdvisedBy(V4, V3) \wedge PostQuals(V4), \lambda_{15})$

$(AdvisedBy(V0, V1) \leftarrow FacultyAffiliate(V1), Publication(V3, V1) \wedge Publication(V3, V0), \lambda_{15})$

$(PostGenerals(V0) \leftarrow Publication(V2, V0) \wedge AdvisedBy(V0, V4), \lambda_{14})$

$(FacultyAdjunct(V0) \leftarrow AdvisedBy(V2, V0), FacultyAdjunct(V4) \wedge AdvisedBy(V2, V4), \lambda_{13})$

$(PostGenerals(V0) \leftarrow TaughtBy(V2, V0, V4) \wedge Student(V0), \lambda_{13})$

$(PostQuals(V0) \leftarrow AdvisedBy(V0, V3) \wedge ProjectMember(V4, V3), \lambda_{12})$

$(AdvisedBy(V0, V1) \leftarrow Publication(V3, V1) \wedge FacultyAdjunct(V1) \wedge Publication(V3, V0), \lambda_{11})$

$(PostQuals(V0) \leftarrow AdvisedBy(V0, V3) \wedge PostQuals(V4) \wedge AdvisedBy(V4, V3), \lambda_{10})$

$(Publication(V0, V1) \leftarrow FacultyAffiliate(V4), FacultyAffiliate(V1) \wedge Publication(V0, V4), \lambda_9)$

$(PostGenerals(V0) \leftarrow AdvisedBy(V0, V3) \wedge PostGenerals(V4) \wedge AdvisedBy(V4, V3), \lambda_8)$

$(Faculty(V0) \leftarrow TempAdvisedBy(V2, V0), \lambda_7)$

$(PostQuals(V0) \leftarrow AdvisedBy(V0, V3) \wedge FacultyAdjunct(V3), \lambda_6)$

$(Faculty(V0) \leftarrow AdvisedBy(V2, V0), \lambda_5)$

$(Publication(V0, V1) \leftarrow FacultyAdjunct(V4), FacultyAdjunct(V1) \wedge Publication(V0, V4), \lambda_4)$

$(PreQuals(V0) \leftarrow TempAdvisedBy(V0, V3), \lambda_3)$

$(FacultyAdjunct(V0) \leftarrow Publication(V2, V4), PreQuals(V4) \wedge Publication(V2, V0), \lambda_2)$

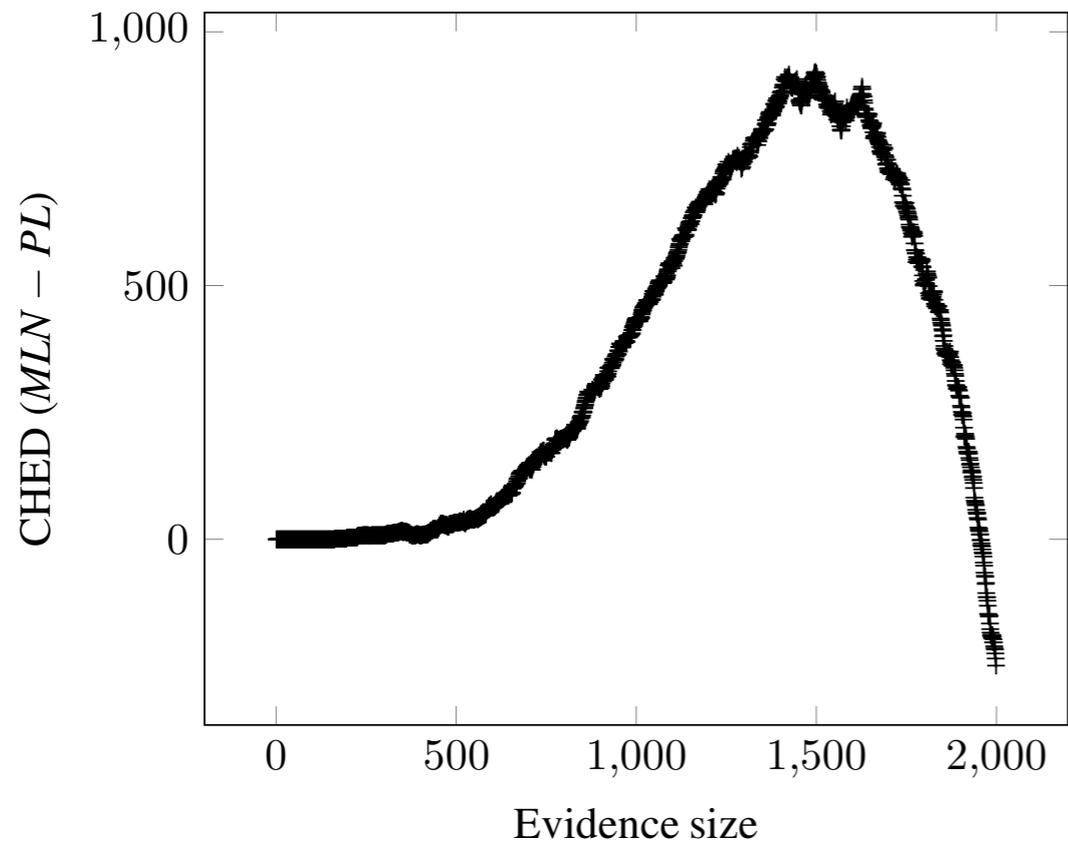
$(PostQuals(V0) \leftarrow AdvisedBy(V0, V3), \lambda_1)$

$(\perp, \lambda_{\perp})$

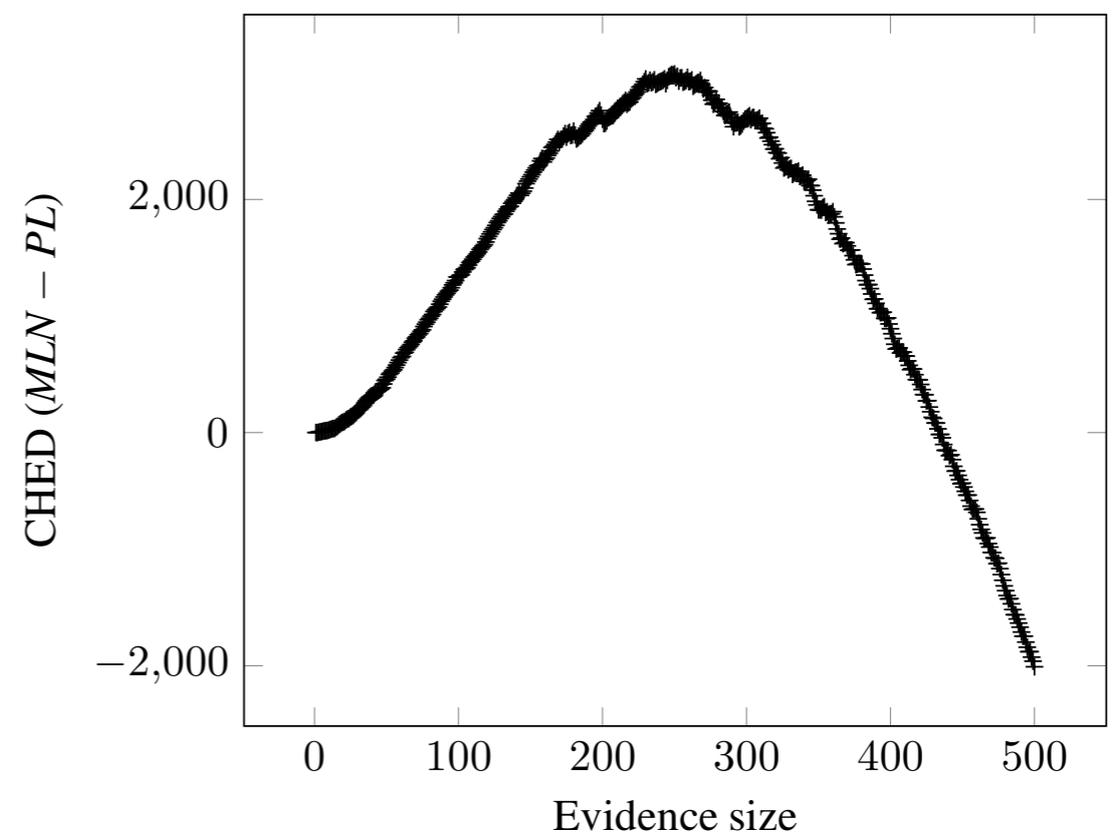
---

# Results

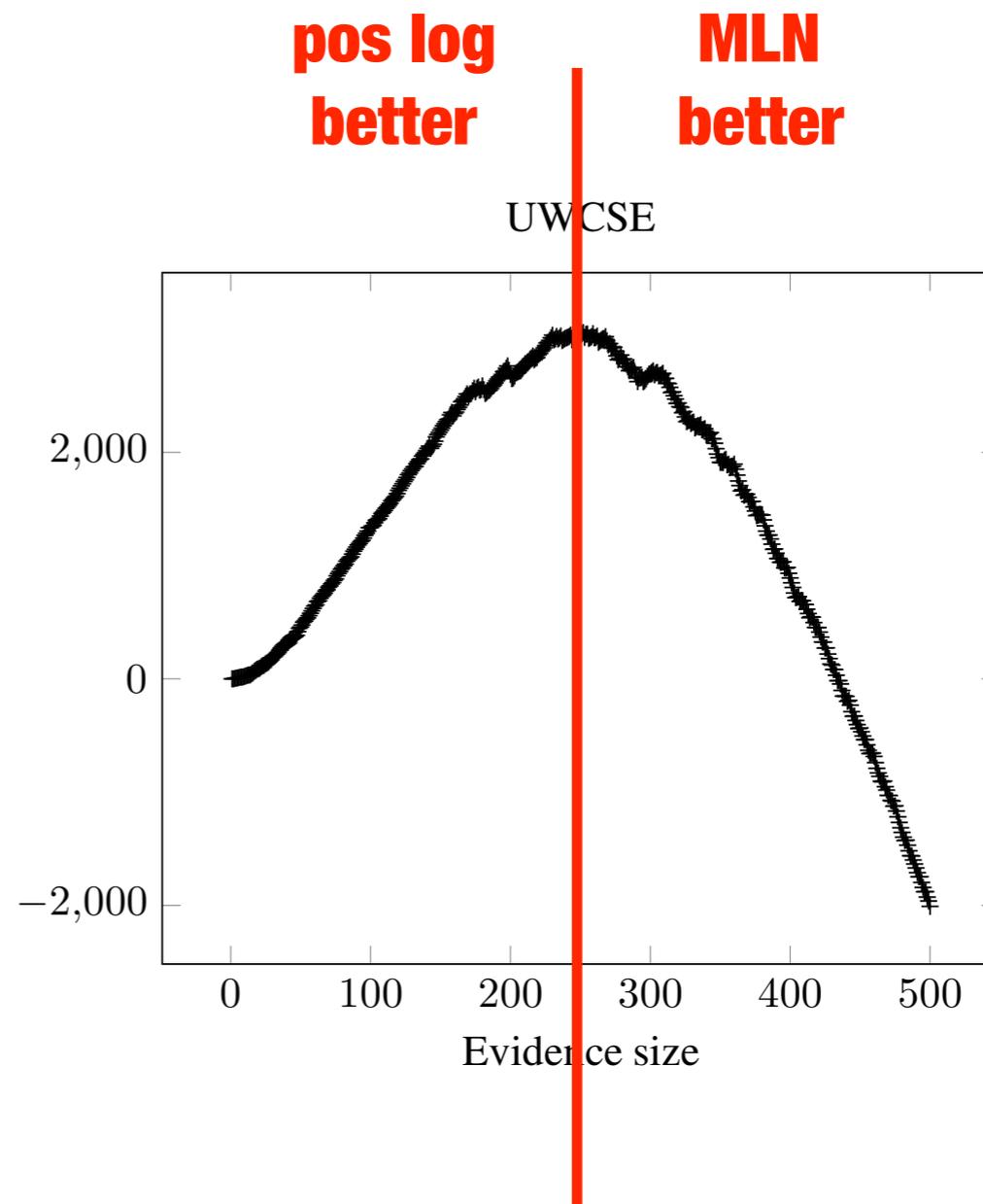
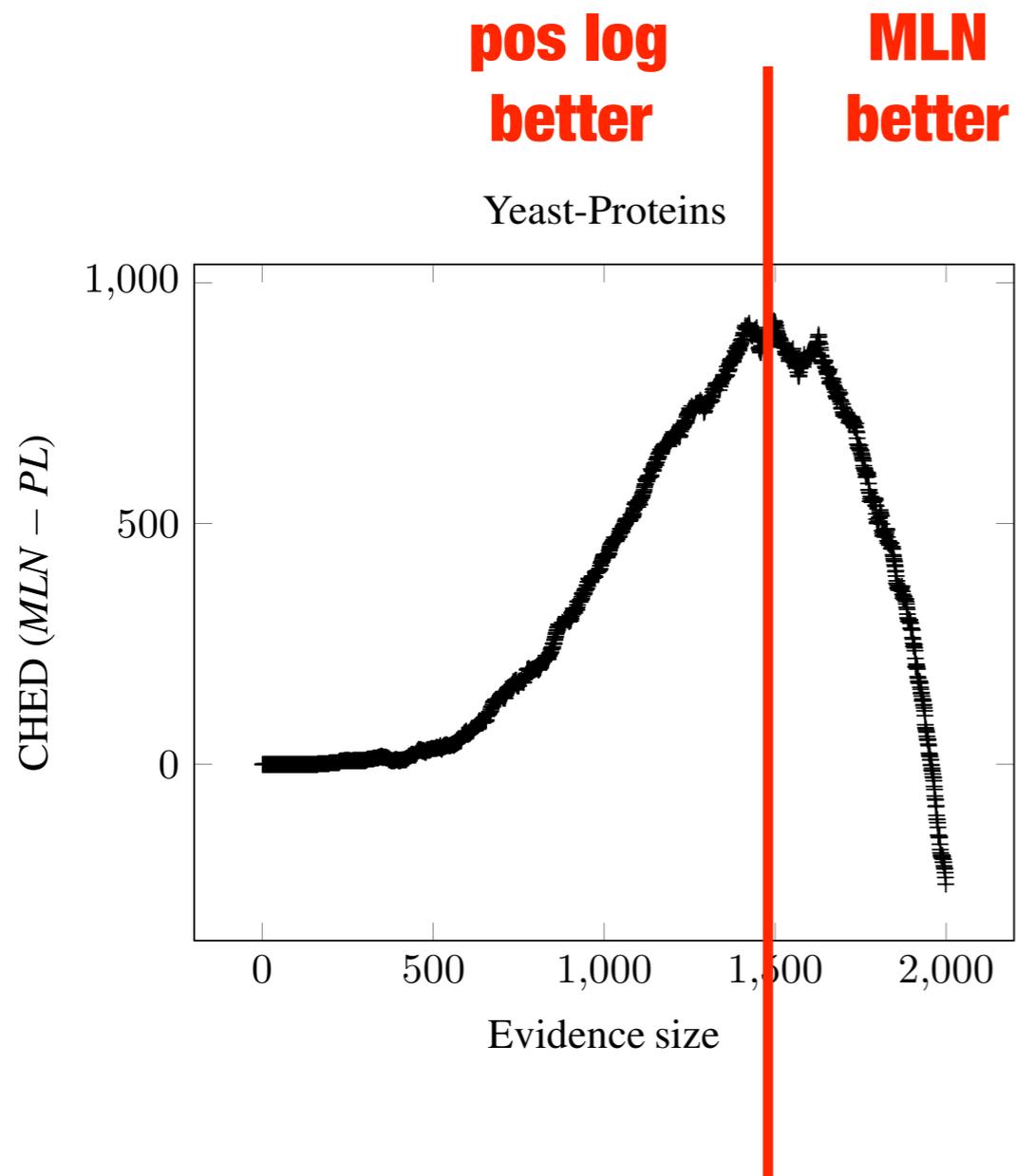
Yeast-Proteins



UWCSE



# Results



# Cautious inference with learned rules

Learned rules are noisy, which means that we cannot simply use classical logic

Possibilistic logic addresses this by prioritising the most certain rules in case of conflict

Another possibility is to weaken the classical entailment relation

- ▶ AMIE: only allow 1-step inferences
- ▶ We want something in between

# k-entailment

Given a set of rules  $\Phi$ , a set of constraints  $\Gamma$  and a set of facts  $E$ , we say that a fact  $\alpha$  is k-entailed if there exists a set of at most  $k$  constants  $C$  such that:

Restricting the facts in  $E$  to those that only contain the constants from  $C$  restores consistency

Restricting the facts in  $E$  still allows us to derive  $\alpha$

# Stratified k-entailment (STRiKE)

$\alpha_1$

$\alpha_2$

$\alpha_3$

...

$\alpha_n$

# Stratified k-entailment (STRiKE)

$\alpha_1$

---

$\alpha_2$

$\alpha_3$

...

$\alpha_n$

Use k-entailment to find all facts that can be derived from  $\alpha_1$

# Stratified k-entailment (STRiKE)

$\alpha_1$

$\alpha_2$

---

$\alpha_3$

...

$\alpha_n$

Use k-entailment to find all facts that can be derived from  $\{\alpha_1, \alpha_2\}$ , but assign a lower certainty level to these facts

# Stratified k-entailment (STRiKE)

$\alpha_1$

$\alpha_2$

$\alpha_3$

...

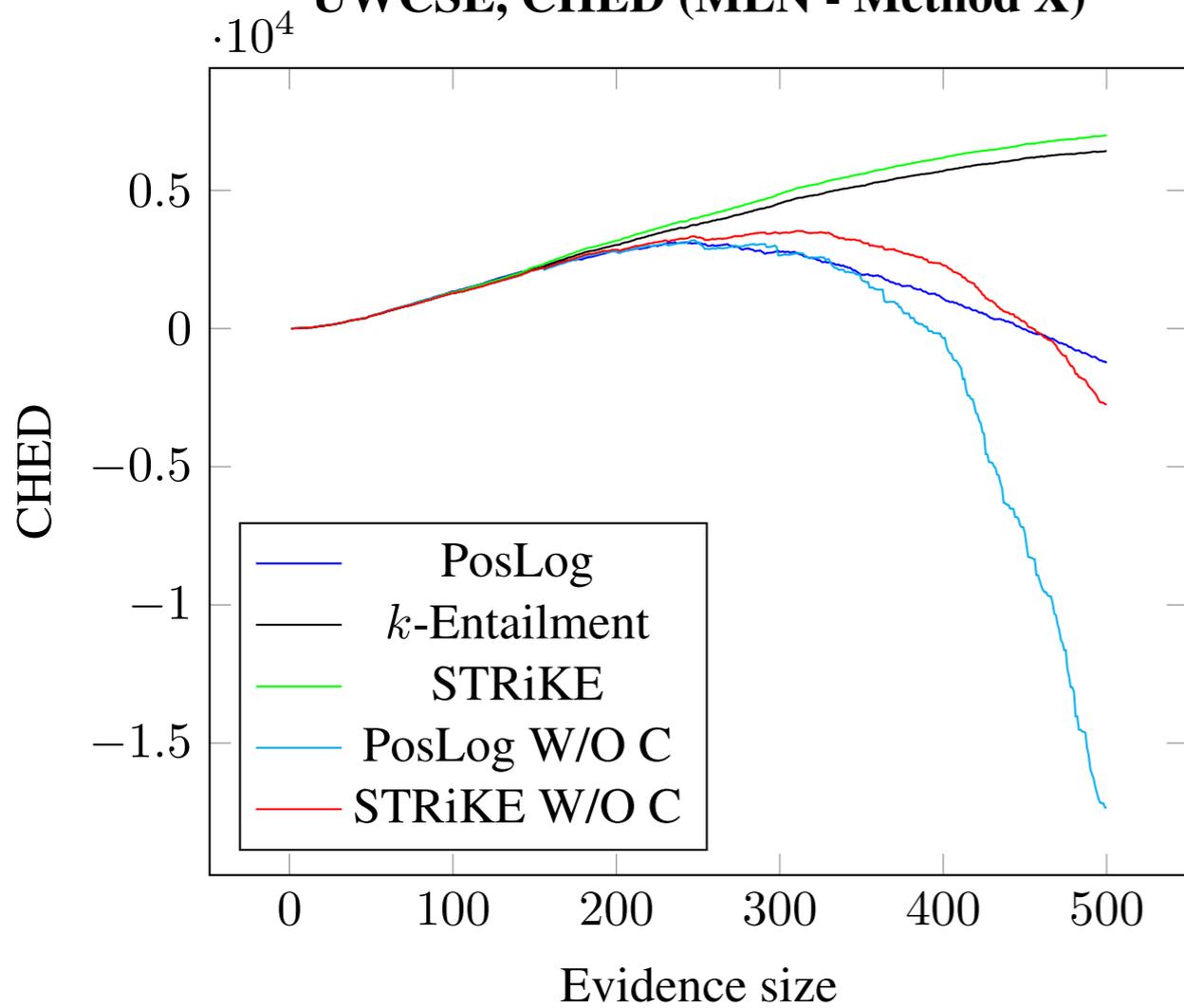
$\alpha_n$

Use k-entailment to find all facts that can be derived from  $\{\alpha_1, \alpha_2, \dots, \alpha_n\}$ , and assign the lowest certainty level to these facts

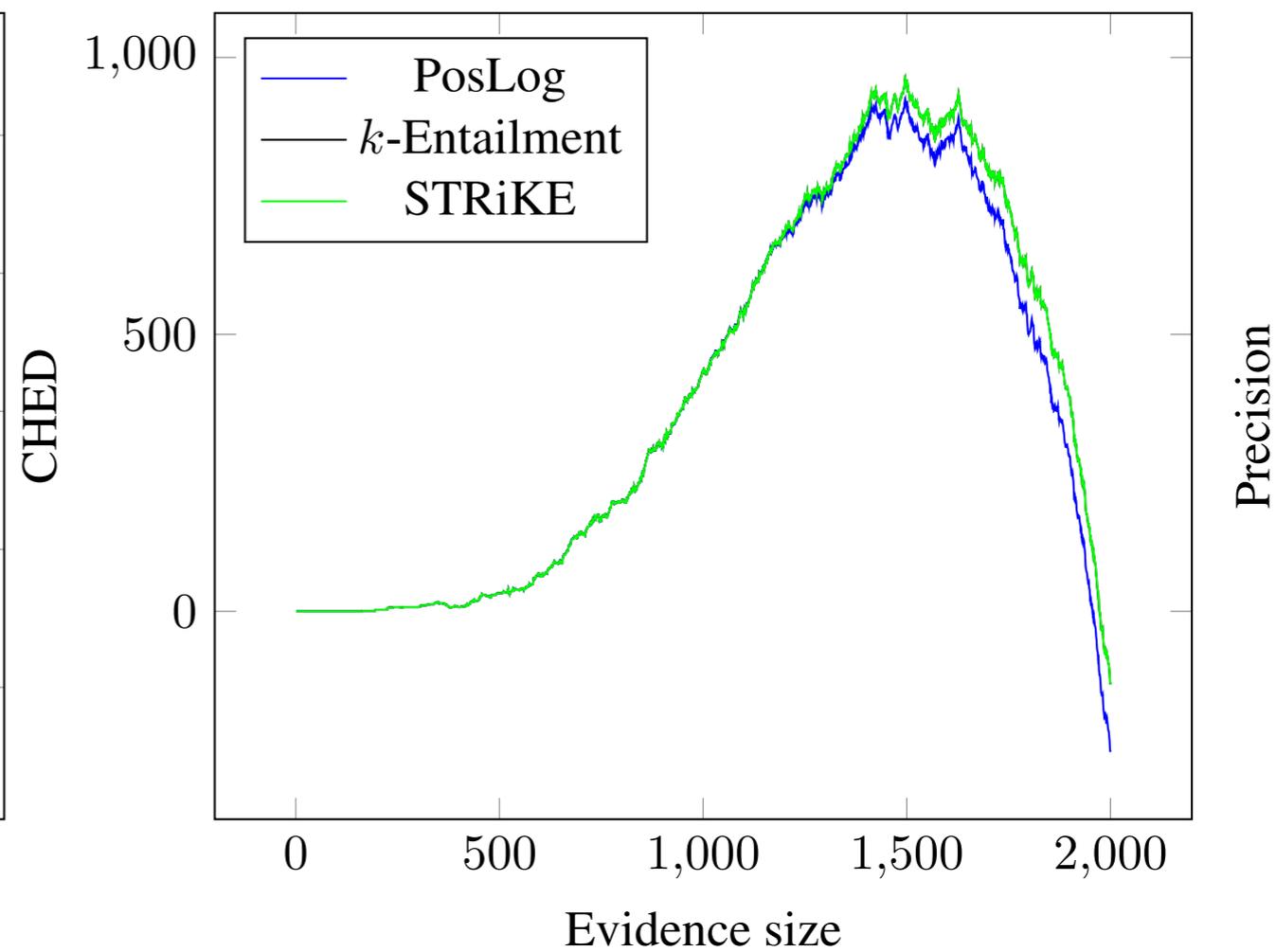
**Note:** set of entailed facts may not be consistent

# Results

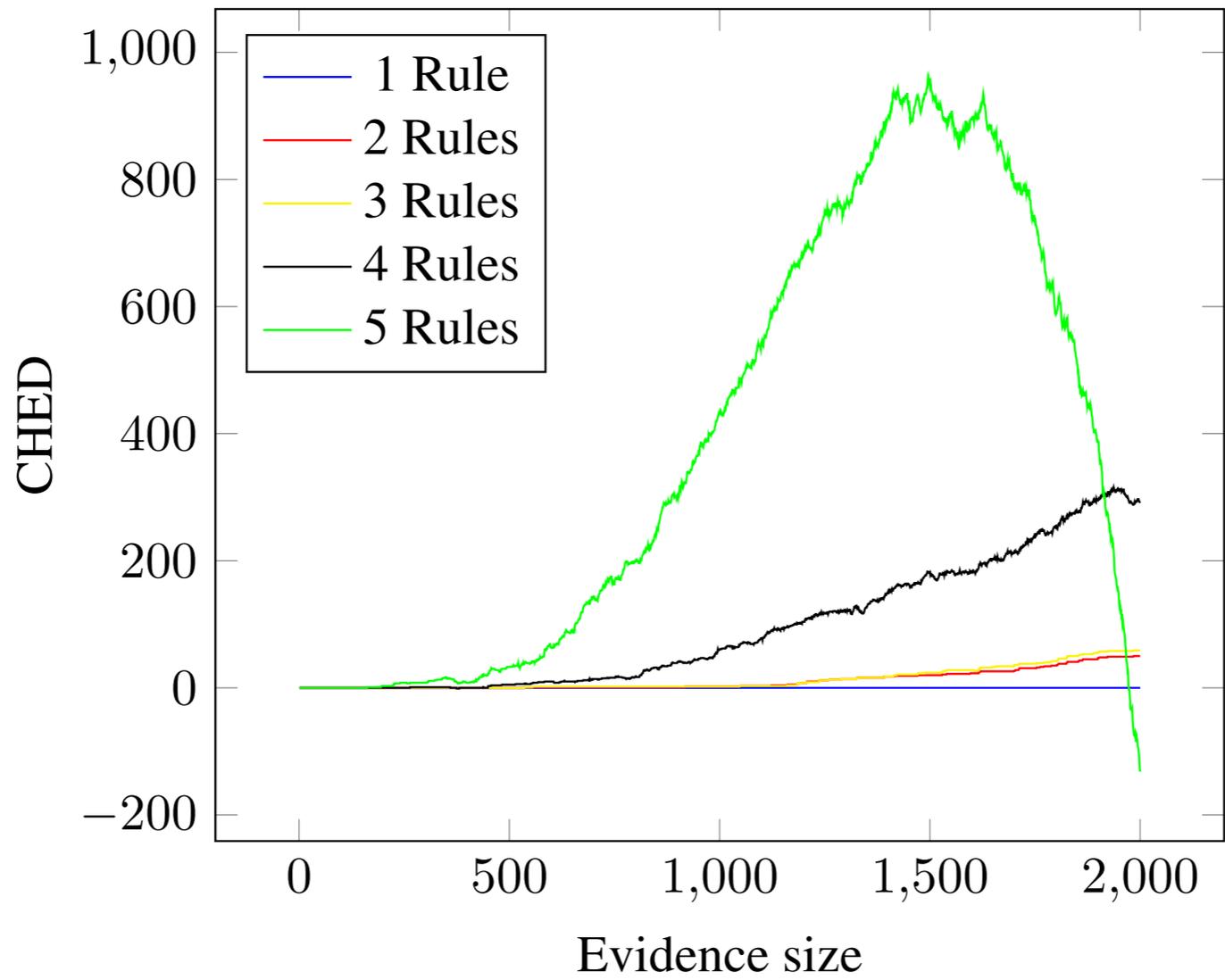
UWCSE, CHED (MLN - Method X)



Proteins, CHED (MLN - Method X)



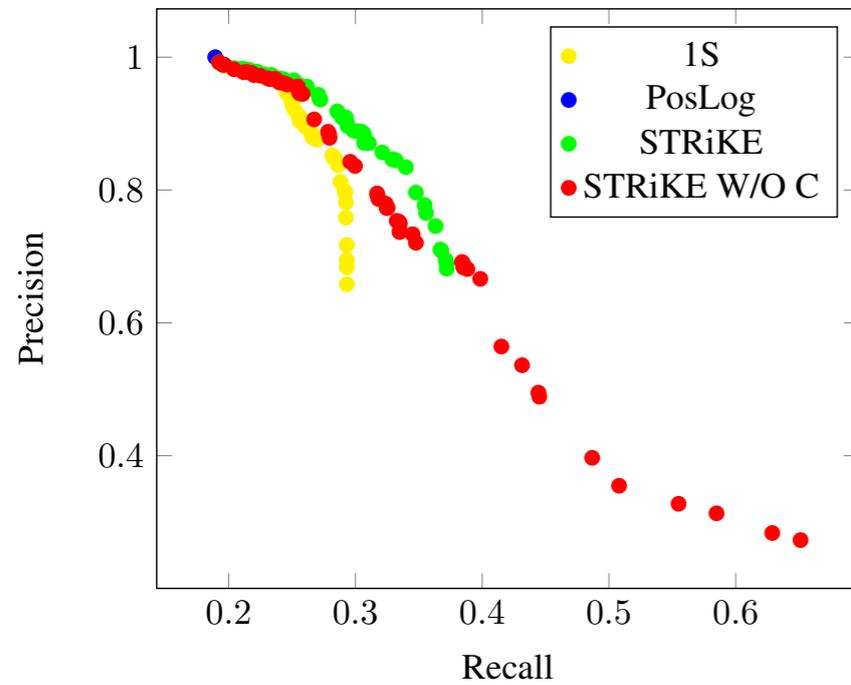
# Results



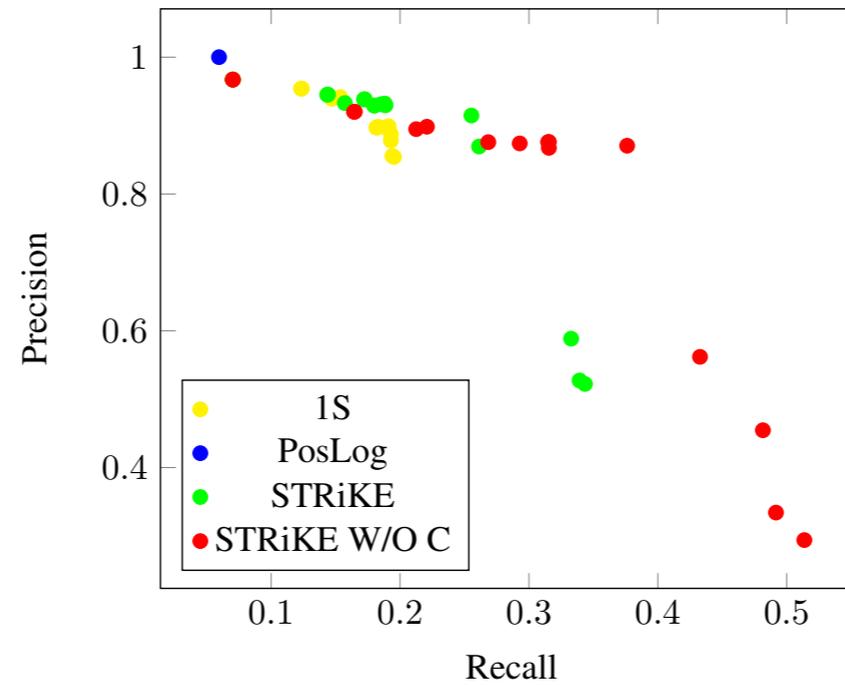
**Yeast-proteins**

# Results

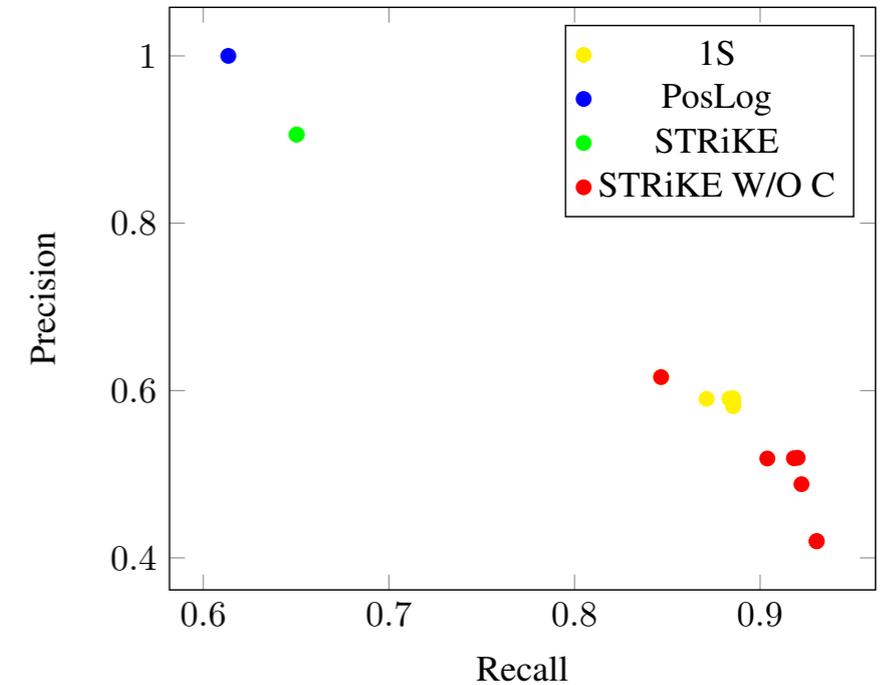
Kinship, Precision-Recall, Evidence Size = 500



UMLS, Precision-Recall, Evidence Size = 100



Nations, Precision-Recall, Evidence Size = 500



**Note:** a different rule learner was used for these experiments

# Conclusions

Possibilistic logic (variants) provide us with a useful framework for encoding and exploiting learned relational knowledge

Compared to Markov logic networks, possibilistic logic theories are much easier/faster to learn, and perform better in most knowledge base completion setting (but they are far less suitable for marginal inference)

Best results were obtained with a variant based on k-entailment, to avoid the drowning issue

Different viewpoints are possible:

- ▶ Possibilistic logic theory encodes a probability distribution
- ▶ Rank rules according to how “safe” they are (PAC guarantees)