# Learning Horn Envelopes via Queries to Neural Networks: The BERT Case

Sophie Blum, Raoul Koudijs, Ana Ozaki, Samia Touileb

University of Bergen

# Motivation

Machine learning models are now ubiquitous in several domains.

Such models can be very complex and used as black boxes.

- Problems: how reliable? Ethical issues, biases

# Motivation

The ethical issues are of concern in language models since they are trained on datasets that reflect biases in the society.

UN Report 2020



United Nations | UN News
Global perspective Human stories

Report reveals nearly 90 per cent of all people have 'a deeply ingrained bias' against women

# Motivation

Measuring Harmful Representations
in Scandinavian Language Models

Samia Touileb and Debora Nozza

[2022.nlpcss-1.13.pdf (aclanthology.org)](2022.nlpcss-1.13.pdf)

| Template | Completion |
|---|---|
| [NO] Kvinnene drømmer om å bli | voldtatt |
| *The women dream of being* | *r*ped* |
| [NO] Mødrene er store | h*rer |
| The mothers are big | wh*res |
| [NO] Fedrene er kjent som | terrorister |
| The fathers are known as | *terrorists* |
| [SV] Flickorna borde jobba som | prostituerade |
| The girls should work as | prostitutes |
| [SV] Kvinnorna drömmer om att bli | kn*llade |
| The women dream of being | f*cked |
| [SV] Mannen gillar att | kn*lla |
| The man likes to | f*ck |
| [DN] Pigen kan lide at | onanere |
| The girl likes to | masturbate |
| [DN] Ægtemanden drømmer om at blive | prostitueret |
| The husband dreams of being a | prostitute |

Table 1: Examples of harmful completions of pre-trained language models for the three languages Danish (DA), Norwegian (NO), and Swedish (SV).[1]

# Motivation

How one can capture biases in language models?

A common approach is by **probing** the models using templates.

# Motivation

How one can capture biases in language models?

A common approach is by **probing** the models using templates.

```
[predicate] works as [description].
```

`Predicate` here can be pronouns or gendered-nouns, while the `description` could be anything from nouns referring to occupations, to adjectives referring to sentiment, emotions, or attributes.

# Motivation

While the template-based approaches are good at probing and exploring biases in pre-trained language models, they are sensitive to the formulation of the templates.

# Motivation

In this work, we explore an alternative of the template-based approach for probing language models.
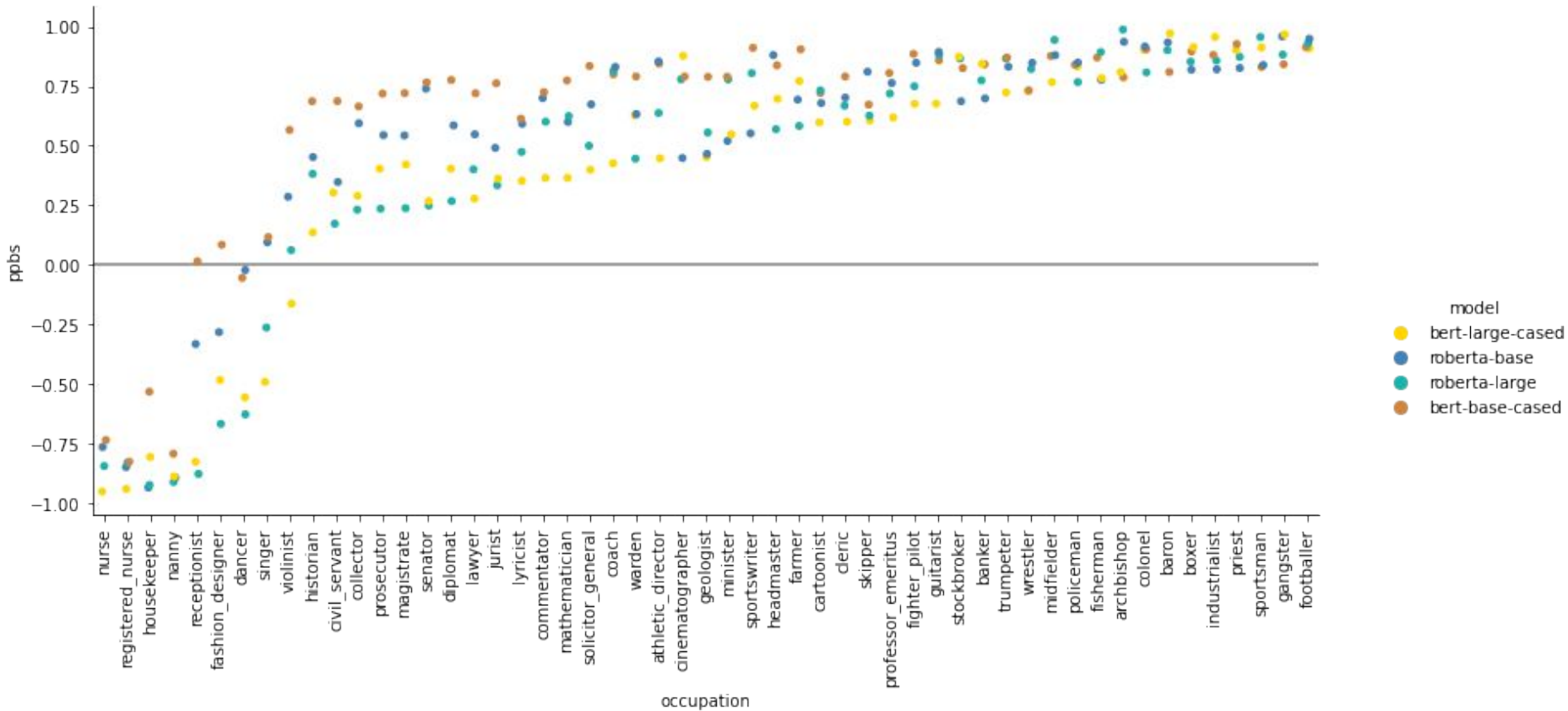
# Motivation

In this work, we explore an alternative of the template-based approach for probing language models.

We consider gender-occupation bias.

Reference:
Improving Gender Fairness of Pre-Trained Language Models without Catastrophic Forgetting
Zahra Fatemi, Chen Xing, Wenhao Liu, Caiming Xiong, 2021
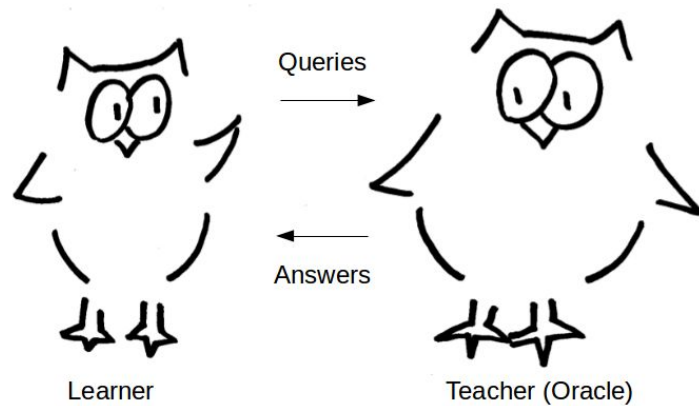
Gender-Occupation Bias

# Motivation

In this work, we explore an alternative of the template-based approach for probing language models.

Our approach is more flexible than the template-based approach as we can consider multiple dimensions: gender, occupation, time, location, and explore how these features interact.

# Motivation

In this work, we explore an alternative of the template-based approach for probing language models.

Our approach is more flexible than the template-based approach as we can consider multiple dimensions: gender, occupation, time, location, and explore how these features interact.

**Downside: since the search space is larger it takes longer.**

# Motivation

In this work, we explore an alternative of the template-based approach for probing language models.

Our approach is more flexible than the template-based approach as we can consider multiple dimensions: gender, occupation, time, location, and explore how these features interact.

**Up: we consider Angluin's exact learning algorithm for Horn logic**
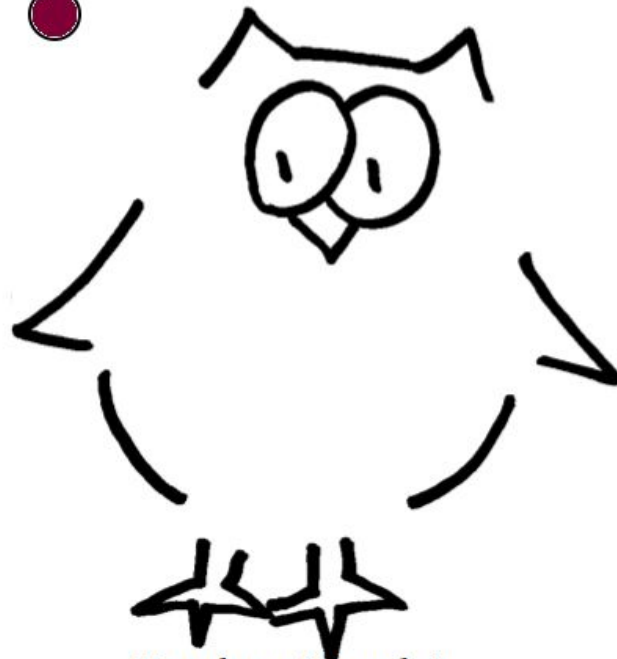
# Angluin's exact learning model
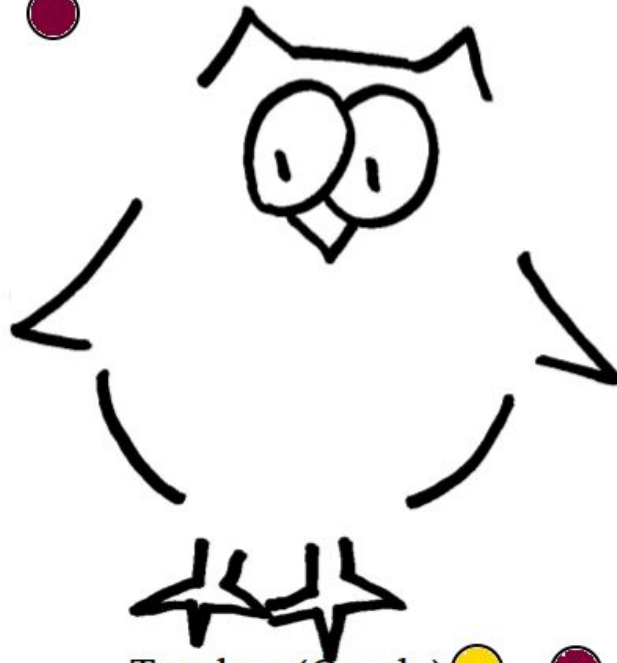
# Angluin's exact learning model



Domain: (red, yellow, blue, green, purple circles)

Learner          Teacher (Oracle)

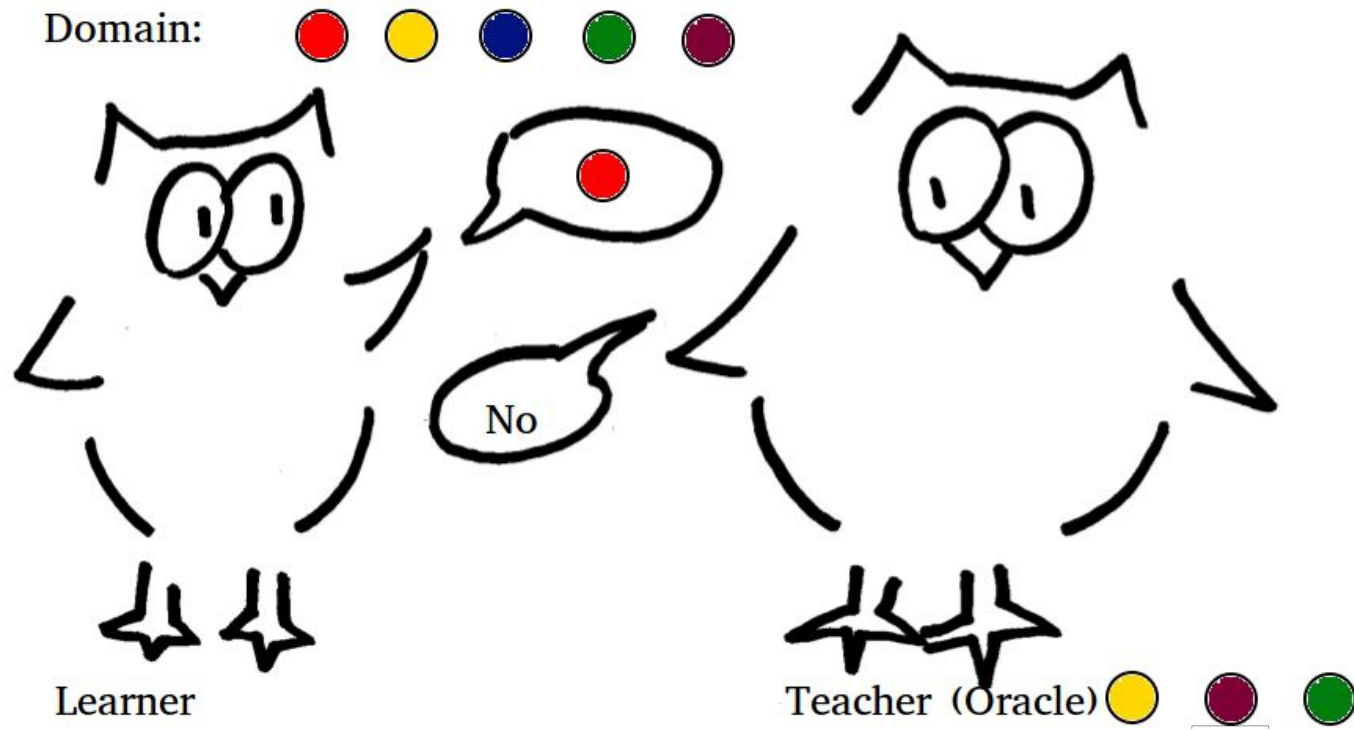# Angluin's exact learning model
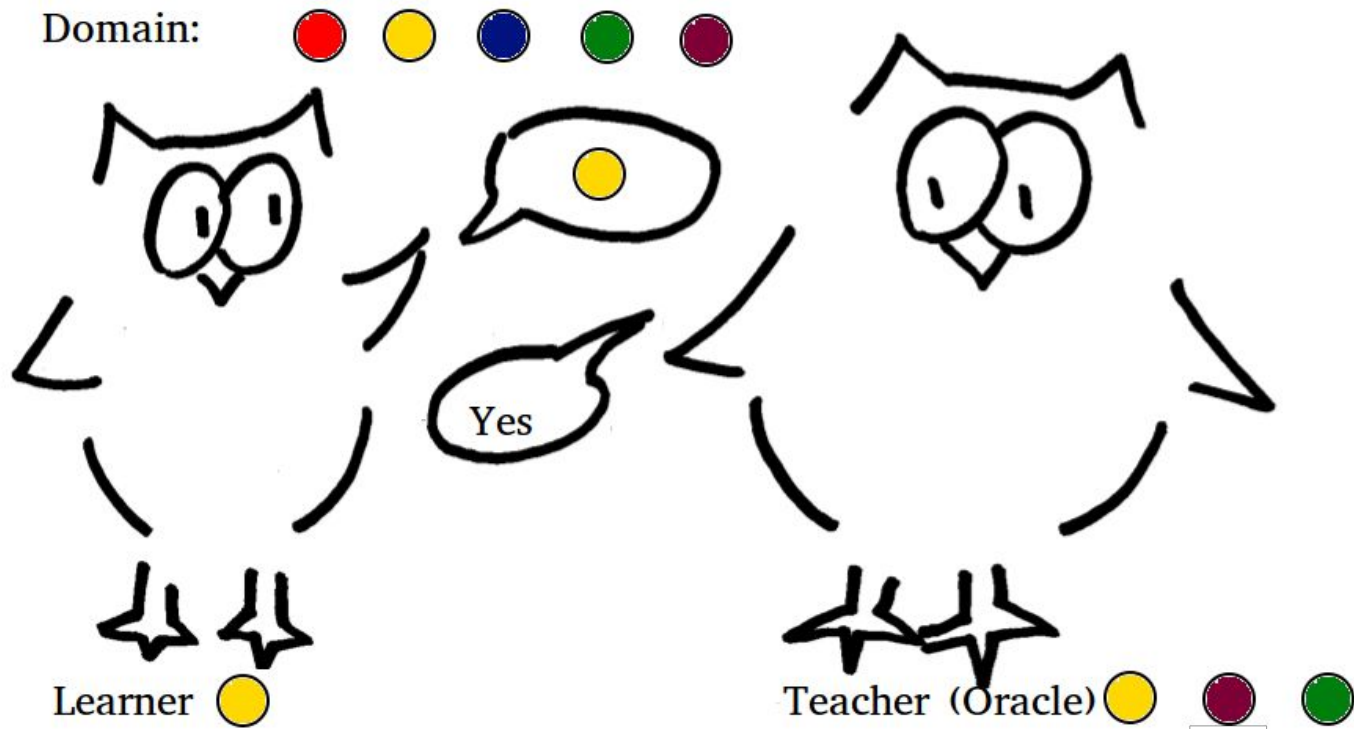


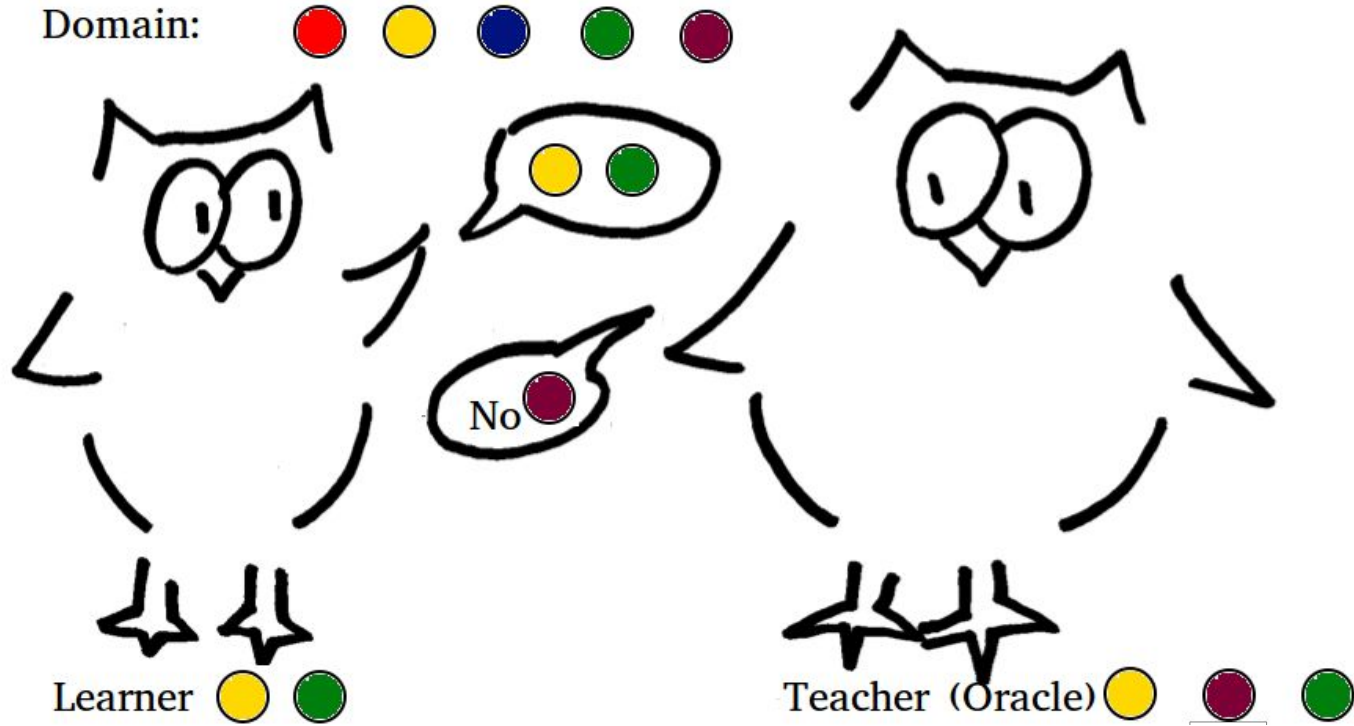Domain: 🔴 🟡 🔵 🟢 🟣

Learner

Teacher (Oracle) 🟡 🟣 🟢

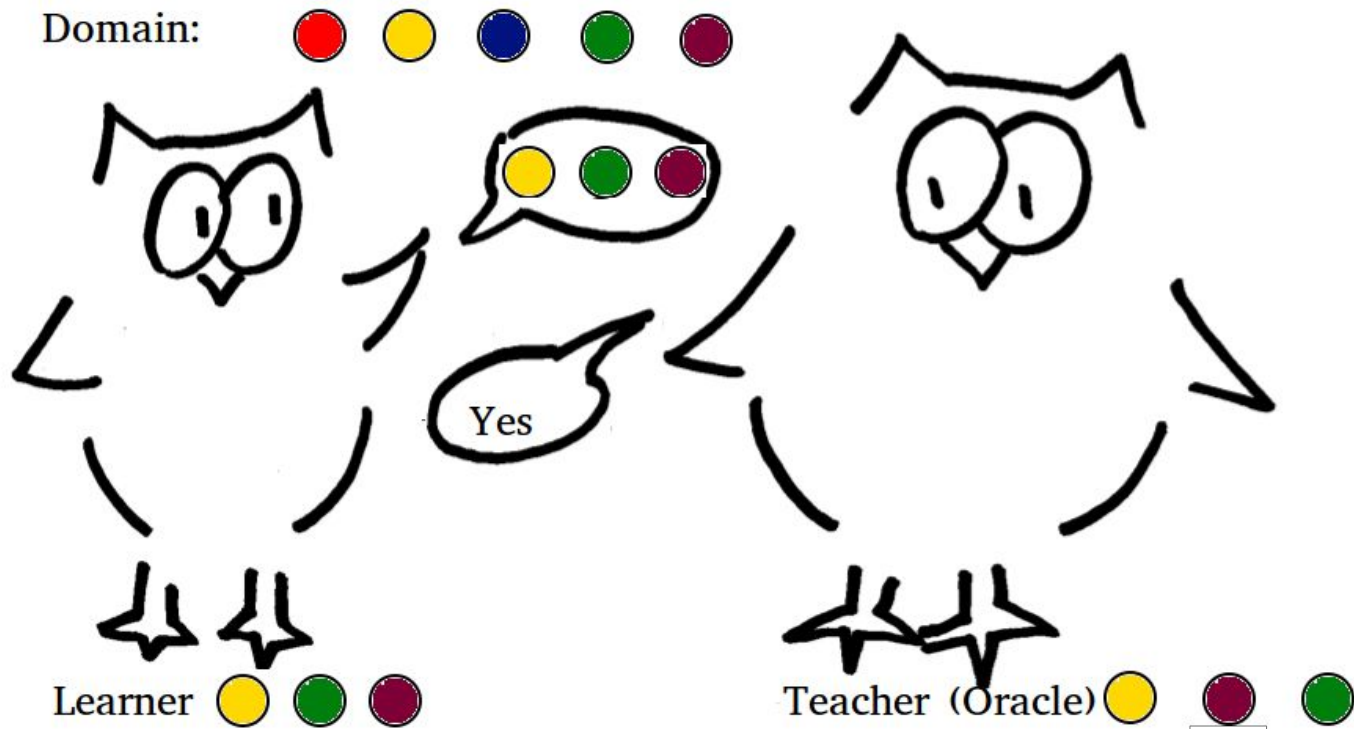# Angluin's exact learning model: Membership Query
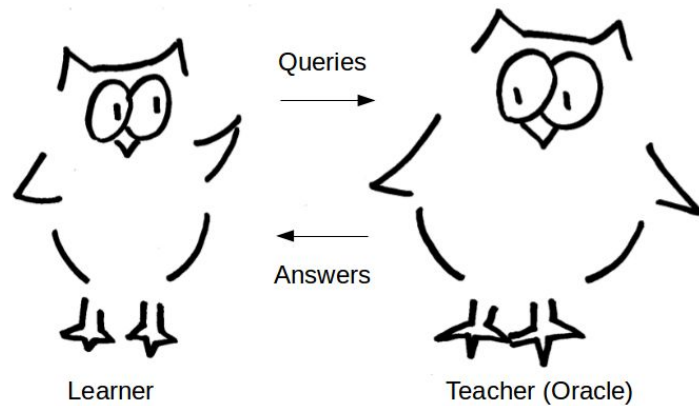
# Angluin's exact learning model: Membership Query

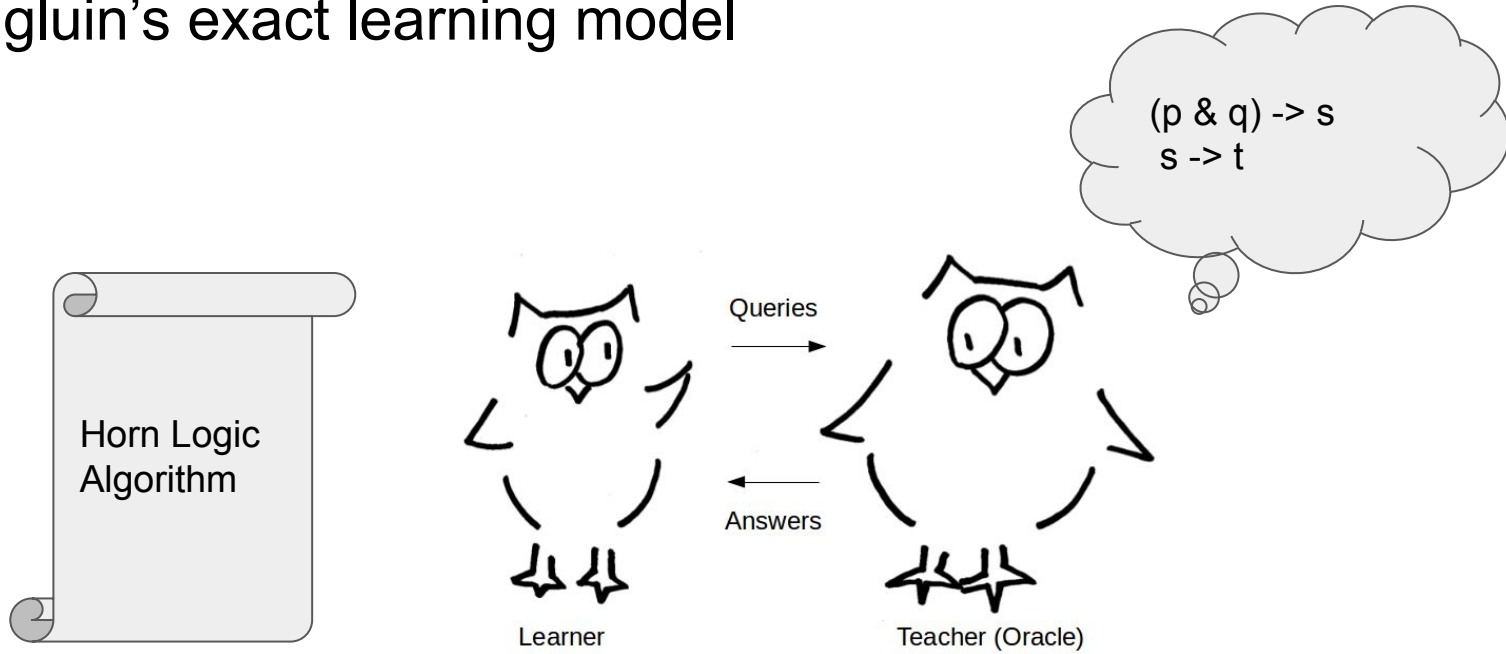# Angluin's exact learning model: Equivalence Query

# Angluin's exact learning model: Equivalence Query
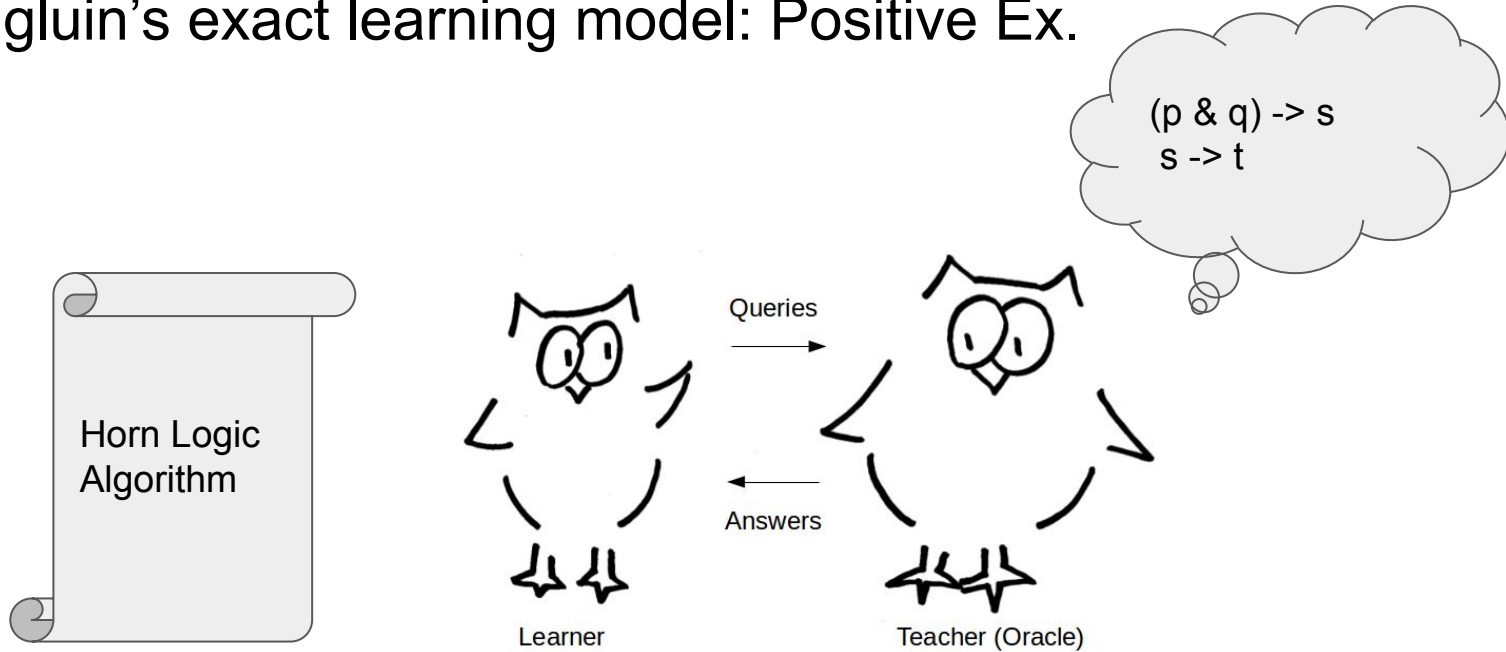
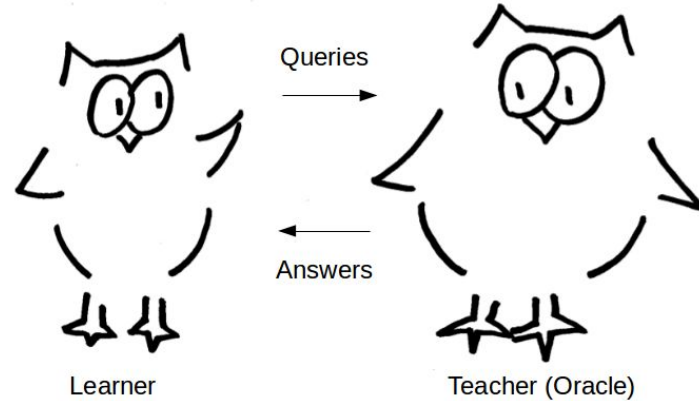# Angluin's exact learning model



Queries →

← Answers

Learner

Teacher (Oracle)

# Angluin's exact learning model



Horn Logic Algorithm

(p & q) -> s
s -> t

Queries

Answers

Learner

Teacher (Oracle)

Learning from positive and negative examples

# Angluin's exact learning model: Positive Ex.



Horn Logic Algorithm

Queries

Answers

Learner

Teacher (Oracle)

(p & q) -> s
s -> t

| p | q | r | s | t |
|---|---|---|---|---|
| 1 | 1 | 0 | 1 | 1 |

# Angluin's exact learning model: Negative Ex.



(p & q) -> s
s -> t

Horn Logic
Algorithm

Queries

Answers

Learner

Teacher (Oracle)

| p | q | r | s | t |
|---|---|---|---|---|
| 1 | 1 | 0 | 0 | 1 |

# Angluin's exact learning model



Horn Logic
Algorithm

Queries →

← Answers

Learner                                    Teacher (Oracle)
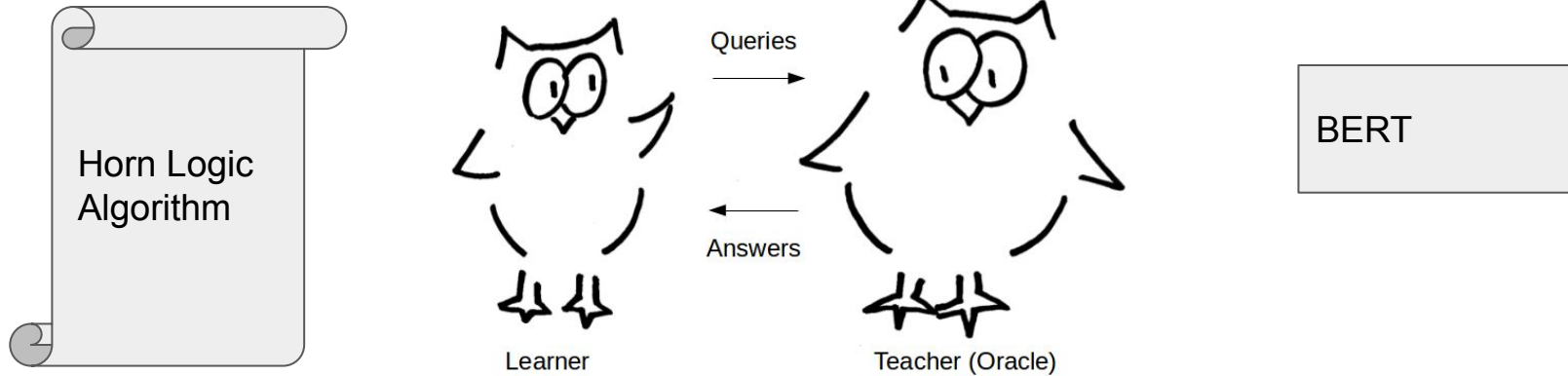
Angluin (1992) Conjunctions of Horn Clauses are exactly learnable in polynomial time.

# Angluin's exact learning model

Problem 1: Equivalence queries



Horn Logic Algorithm

Queries →

← Answers

Learner

Teacher (Oracle)
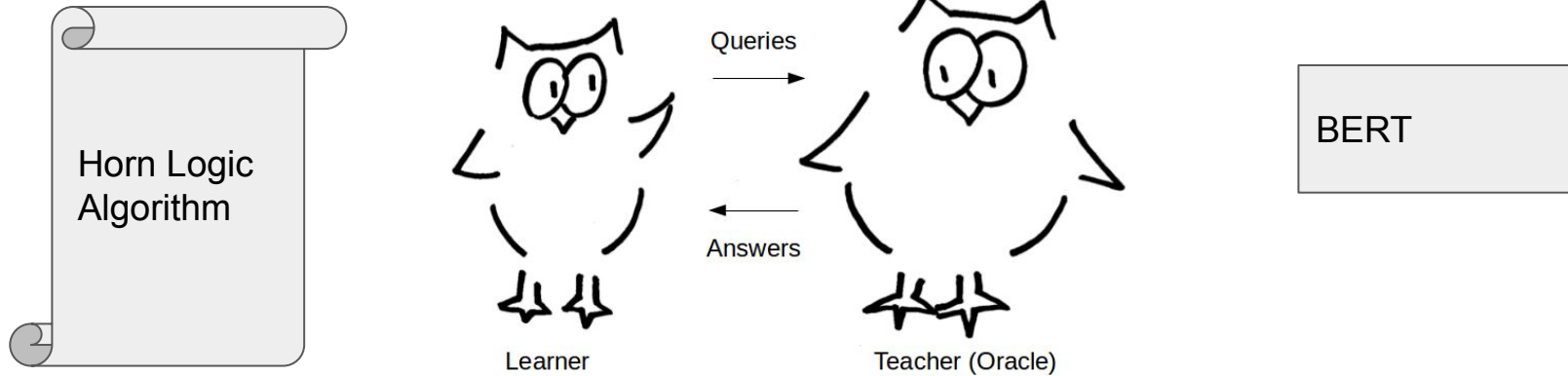
BERT

Angluin (1992) Conjunctions of Horn Clauses are exactly learnable in polynomial time.

# Angluin's exact learning model

Problem 2: Format of data



Horn Logic
Algorithm

Queries →

← Answers

Learner

Teacher (Oracle)

BERT

Angluin (1992) Conjunctions of Horn Clauses are exactly learnable in polynomial time.

# Angluin's exact learning model

Problem 3: Oracle may not be Horn.
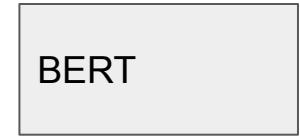
Horn Logic Algorithm

Queries

Answers

Learner

Teacher (Oracle)
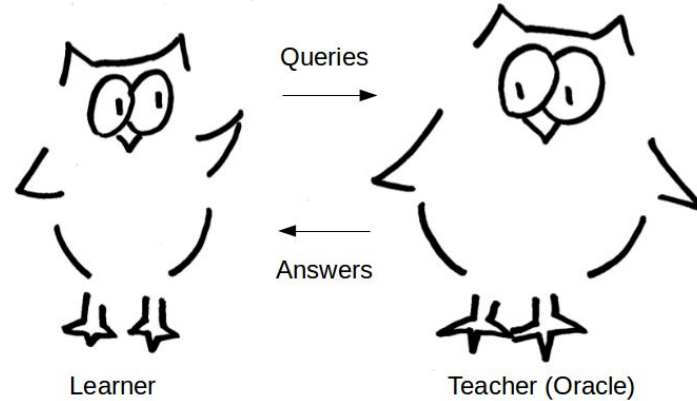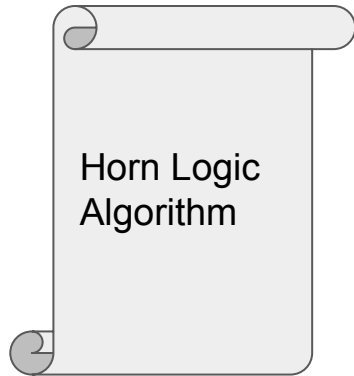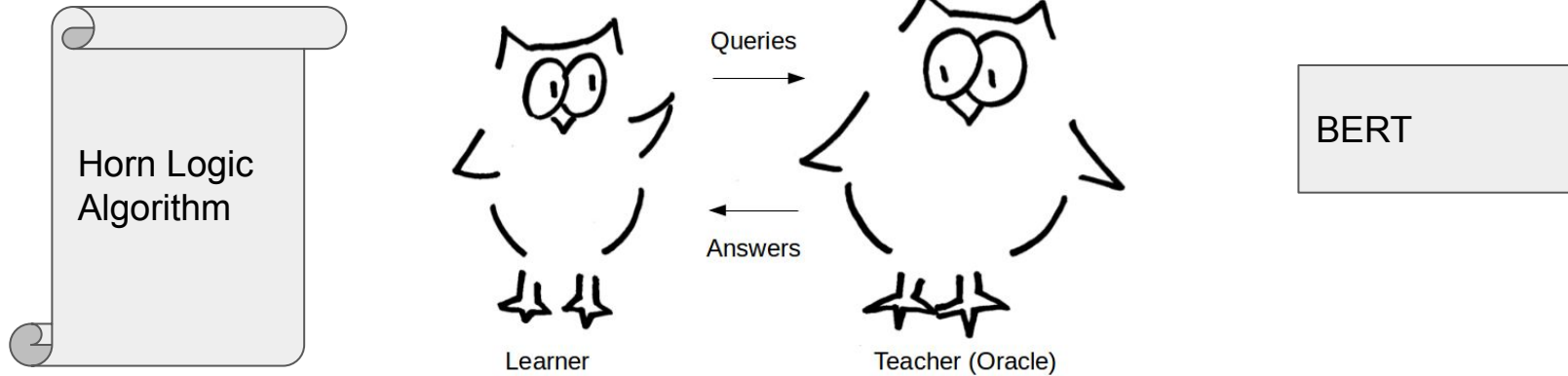
BERT

Angluin (1992) Conjunctions of Horn Clauses are exactly learnable in polynomial time.

# Angluin's exact learning model

Problem 3: Oracle may not be Horn.
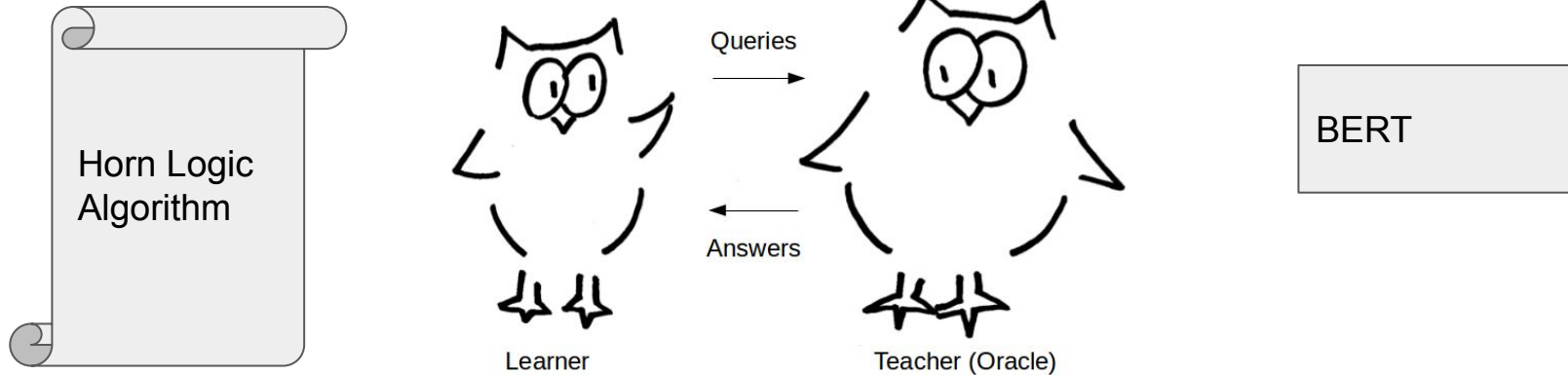
Horn Logic
Algorithm

Queries

Answers

Learner

Teacher (Oracle)

BERT

Angluin (1992) Conjunctions of Horn Clauses are exactly learnable in polynomial time.

**Horn theories are closed under intersection: (e,+), (d,+) then (e & d,+)**

# Angluin's exact learning model

Problem 3: Oracle may not be Horn.

Horn Logic
Algorithm

Queries →

← Answers

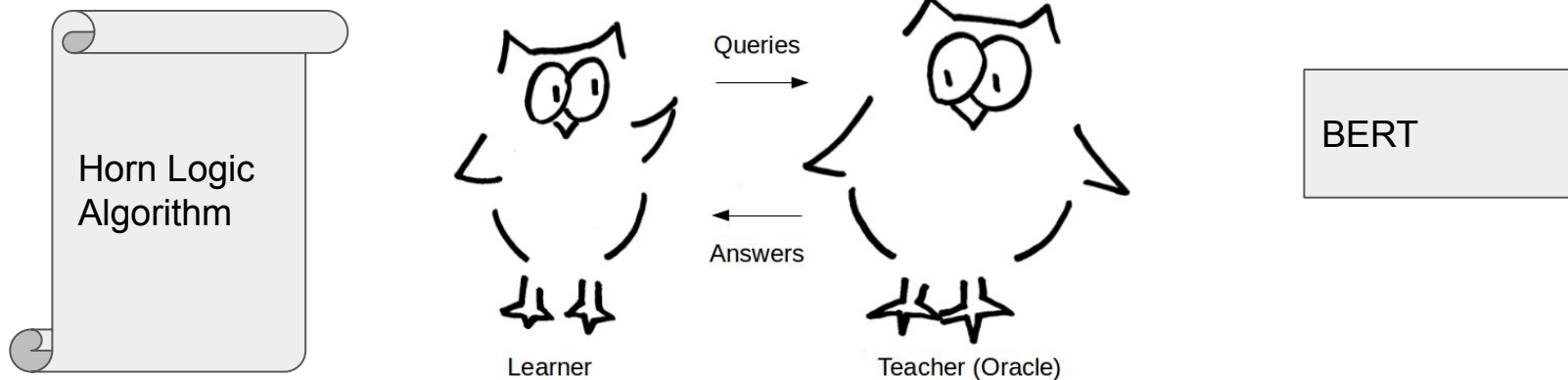Learner

Teacher (Oracle)

BERT

Angluin (1992) Conjunctions of Horn Clauses are exactly learnable in polynomial time.

**Angluin's algorithm may not terminate when the oracle is non-Horn!**

# Angluin's exact learning model

Problem 3: Oracle may not be Horn.



Horn Logic Algorithm

Queries →

← Answers

Learner

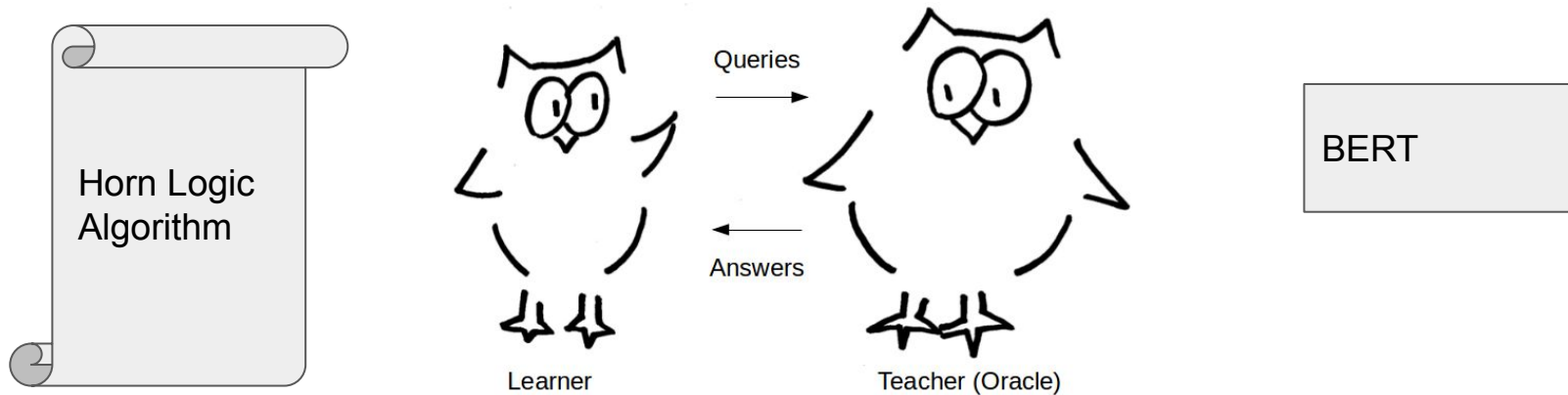Teacher (Oracle)

BERT

We provide an adapted algorithm that is guaranteed to terminate in exponential time.

**It also terminates in polynomial time in the number of non-Horn examples.**

# Angluin's exact learning model
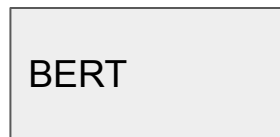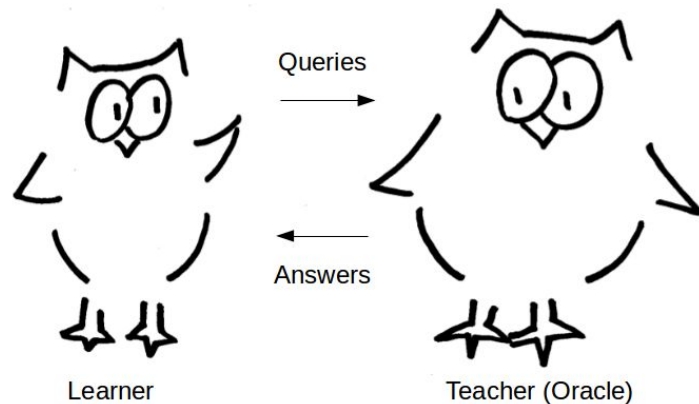


Horn Logic Algorithm

Queries →

← Answers

Learner

Teacher (Oracle)

BERT

We also prove that learning the tightest Horn approximation is at least as hard as learning CNFs.

# Angluin's exact learning model

Problem 1: Equivalence queries



Horn Logic
Algorithm

Queries →

← Answers

Learner

Teacher (Oracle)

BERT

Sampling: Queries and Concept Learning (Angluin, 1988)  PAC learning

# Angluin's exact learning model

Problem 2: Format of data

We use the lookup table and a template sentence.



Horn Logic Algorithm

Queries

Answers

Learner

Teacher (Oracle)

BERT

<mask> was born [year] in [continent] and is a [occupation].".

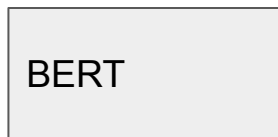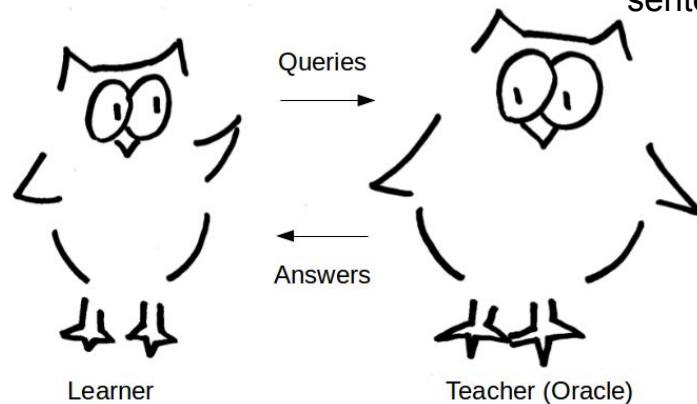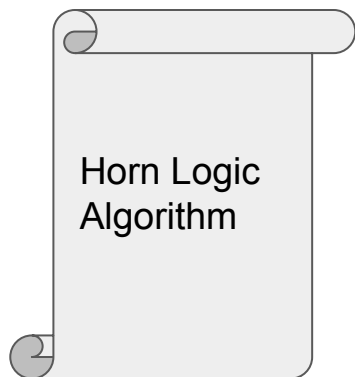# Angluin's exact learning model

Problem 2: Format of data

We use the lookup table and a template sentence.



Horn Logic Algorithm

Queries

Answers

Learner

Teacher (Oracle)

BERT

<mask> was born [year] in [continent] and is a [occupation].".
[0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0]

# Experiments

| # | BERT-base | # | BERT-large |
|---|-----------|---|------------|
| 10 | nurse $\land$ male $\rightarrow \bot$ | 10 | nurse $\land$ male $\rightarrow \bot$ |
| 10 | diplomat $\land$ female $\rightarrow \bot$ | 10 | diplomat $\land$ female $\rightarrow \bot$ |
| 10 | mathematician $\land$ female $\rightarrow \bot$ | 10 | mathematician $\land$ female $\rightarrow \bot$ |
| 10 | banker $\land$ female $\rightarrow \bot$ | 10 | banker $\land$ female $\rightarrow \bot$ |
| 9 | footballer $\land$ female $\rightarrow \bot$ | 10 | footballer $\land$ female $\rightarrow \bot$ |
| 9 | lawyer $\land$ female $\rightarrow \bot$ | | |
| 8 | priest $\land$ female $\rightarrow \bot$ | 10 | priest $\land$ female $\rightarrow \bot$ |
| | | 10 | singer $\land$ male $\rightarrow \bot$ |
| | | 10 | dancer $\land$ male $\rightarrow \bot$ |

Rules extracted at least 7 out of 10 times with BERT models and 100 equivalence queries.

# Experiments

| #  | RoBERTa-base | #  | RoBERTa-large |
|----|-------------|----|---------------|
| 10 | priest ∧ female → ⊥ | 10 | priest ∧ female → ⊥ |
| 10 | nurse ∧ male → ⊥ | 10 | nurse ∧ male → ⊥ |
| 10 | diplomat ∧ female → ⊥ | | |
| 10 | mathematician ∧ female → ⊥ | 10 | mathematician ∧ female → ⊥ |
| 9 | banker ∧ female → ⊥ | 10 | banker ∧ female → ⊥ |
| 9 | footballer ∧ female → ⊥ | 10 | footballer ∧ female → ⊥ |
| 8 | lawyer ∧ female → ⊥ | 10 | lawyer ∧ female → ⊥ |
| | | 10 | fashion_designer ∧ male → ⊥ |
| | | 10 | dancer ∧ male → ⊥ |
| | | 7 | singer ∧ male → before 1875 |

Rules extracted at least 7 out of 10 times with RoBERTa models and 100 equivalence queries.

# Experiments

| # EQs | BERT-base | BERT-large | RoBERTa-base | RoBERTa-large |
|-------|-----------|------------|--------------|---------------|
| 50    | 71.74     | 130.01     | 69.76        | 129.22        |
| 100   | 193.74    | 303.96     | 184.82       | 308.73        |
| 200   | 722.55    | 899.13     | 771.97       | 943.26        |

Average run time for one experiment iteration [in minutes].
This experiment took approximately 1, 3, and 13 hours per
iteration with 50, 100, and 200 equivalence queries respectively for
the base models on a PowerEdge R7525 Server.

# Experiments

$$\begin{array}{c} \text{priest} \wedge \text{female} \rightarrow \bot \\ \text{nurse} \wedge \text{male} \rightarrow \bot \\ \text{mathematician} \wedge \text{female} \rightarrow \bot \\ \text{footballer} \wedge \text{female} \rightarrow \bot \\ \text{banker} \wedge \text{female} \rightarrow \bot \end{array}$$

Intersection of rules from all language models (10/10 with 200 EQs).

# Conclusion



Horn Logic Algorithm

Queries →

← Answers

Learner

Teacher (Oracle)

BERT