



Machine Reasoning: Technology, Dilemma and Future

Nan DUAN, Duyu TANG, Ming ZHOU
Microsoft Research Asia
{nanduan, dutang, mingzhou}@microsoft.com

EMNLP-2020 Tutorial



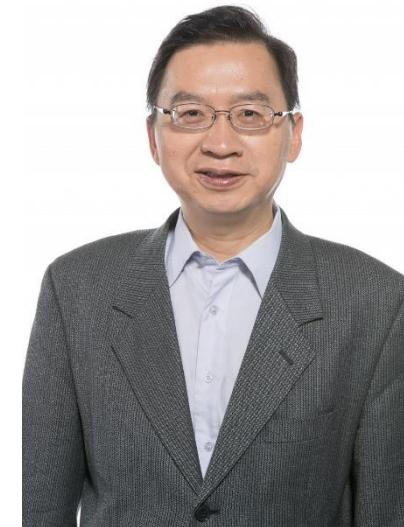
Nan DUAN

Principal Researcher
Microsoft Research Asia
nanduan@microsoft.com



Duyu TANG

Senior Researcher
Microsoft Research Asia
dutang@microsoft.com

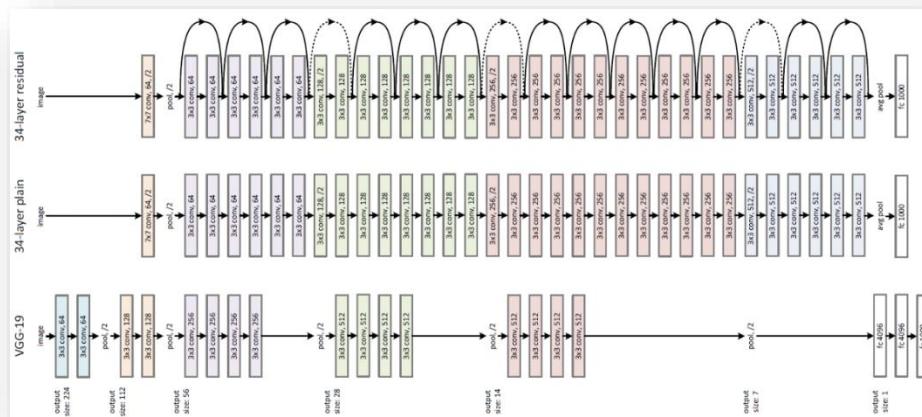
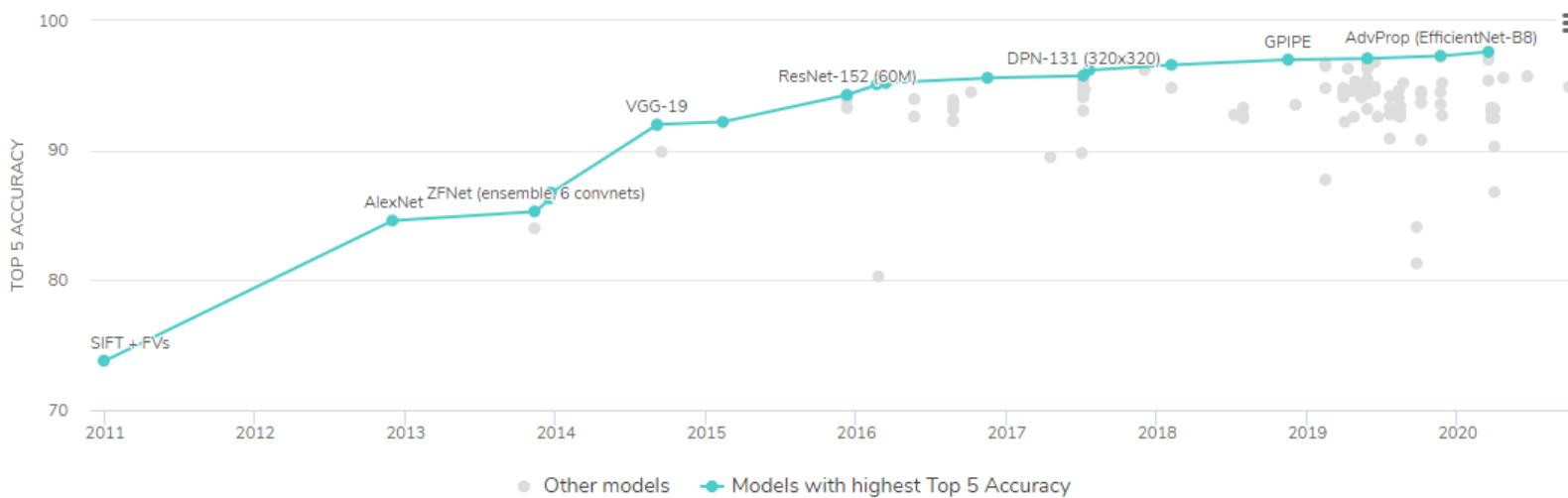


Ming ZHOU

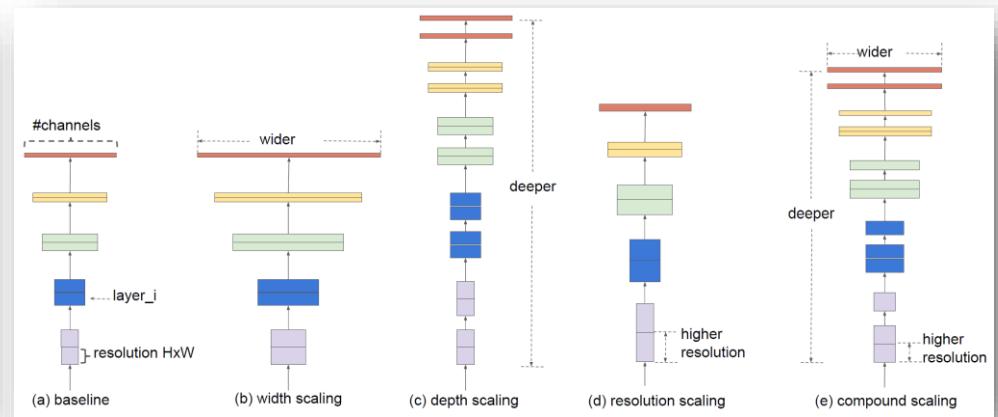
Assistant Managing Director
Microsoft Research Asia
mingzhou@microsoft.com

Current AI Paradigm: Pre-trained Deep Neural Networks (achieving SOTA results on ImageNet)

Image Classification on ImageNet



ResNet (He et al., 2015)



EfficientNet (Tan and Le, 2020)



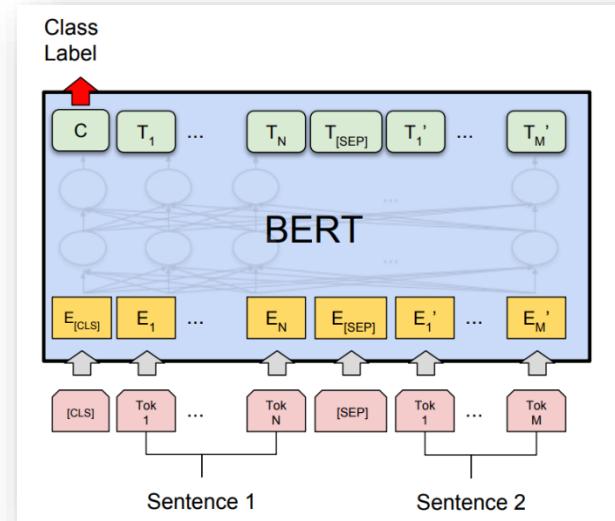
SimCLR (Chen et al., 2020)

Current AI Paradigm: Pre-trained Deep Neural Networks (achieving SOTA results on SuperGLUE)

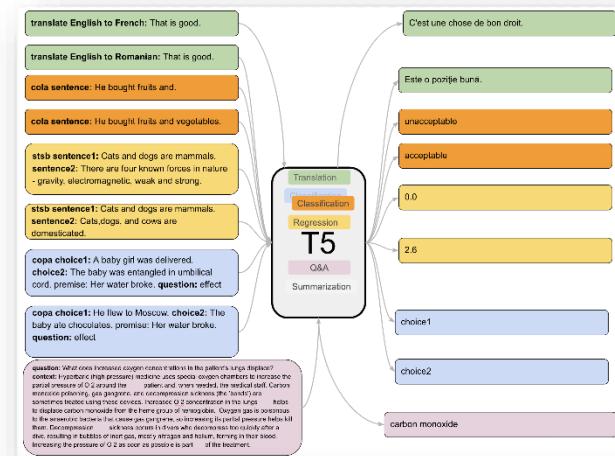
SuperGLUE GLUE Paper Code Tasks Leaderboard FAQ Diagnostics Submit Login

Leaderboard Version: 2.0

Rank	Name	Model	URL	Score	BoolQ	CB	COPA	MultiRC	ReCoRD	RTE	WiC	WSC	AX-b	AX-g
1	SuperGLUE Human Baselines	SuperGLUE Human Baselines		89.8	89.0	95.8/98.9	100.0	81.8/51.9	91.7/91.3	93.6	80.0	100.0	76.6	99.3/99.7
2	T5 Team - Google	T5		89.3	91.2	93.9/96.8	94.8	88.1/63.3	94.1/93.4	92.5	76.9	93.8	65.6	92.7/91.9
3	Huawei Noah's Ark Lab	NEZHA-Plus		86.7	87.8	94.4/96.0	93.6	84.6/55.1	90.1/89.6	89.1	74.6	93.2	58.0	87.1/74.4
4	Alibaba PAI&ICBU	PAI Albert		86.1	88.1	92.4/96.4	91.8	84.6/54.7	89.0/88.3	88.8	74.1	93.2	75.6	98.3/99.2
5	Tencent Jarvis Lab	RoBERTa (ensemble)		85.9	88.2	92.5/95.6	90.8	84.4/53.4	91.5/91.0	87.9	74.1	91.8	57.6	89.3/75.6
6	Zhuiyi Technology	RoBERTa-mtl-adv		85.7	87.1	92.4/95.6	91.2	85.1/54.3	91.7/91.3	88.1	72.1	91.8	58.5	91.0/78.1
7	Facebook AI	RoBERTa		84.6	87.1	90.5/95.2	90.6	84.4/52.5	90.6/90.0	88.2	69.9	89.0	57.9	91.0/78.1
8	Infosys : DAWN : AI Research	RoBERTa-ICETS		77.4	84.7	88.2/91.6	85.8	78.4/37.5	82.9/82.4	83.8	69.1	65.1	35.2	93.8/68.8
9	Timo Schick	iPET (ALBERT) - Few-Shot (32 Examples)		75.4	81.2	79.9/88.8	90.8	74.1/31.7	85.9/85.4	70.8	49.3	88.4	36.2	97.8/57.9
10	IBM Research AI	BERT-mtl		73.5	84.8	89.6/94.0	73.8	73.2/30.5	74.6/74.0	84.1	66.2	61.0	29.6	97.8/57.3
11	Ben Mann	GPT-3 few-shot - OpenAI		71.8	76.4	52.0/75.6	92.0	75.4/30.5	91.1/90.2	69.0	49.4	80.1	21.1	90.4/55.3
12	SuperGLUE Baselines	BERT++		71.5	79.0	84.8/90.4	73.8	70.0/24.1	72.0/71.3	79.0	69.6	64.4	38.0	99.4/51.4
		BERT		69.0	77.4	75.7/83.6	70.6	70.0/24.1	72.0/71.3	71.7	69.6	64.4	23.0	97.8/51.7
		Most Frequent Class		47.1	62.3	21.7/48.4	50.0	61.1/0.3	33.4/32.5	50.3	50.0	65.1	0.0	100.0/50.0
		CBoW		44.5	62.2	49.0/71.2	51.6	0.0/0.5	14.0/13.6	49.7	53.1	65.1	-0.4	100.0/50.0
		Outside Best		-	80.4	-	84.4	70.4/24.5	74.8/73.0	82.7	-	-	-	-
-	Stanford Hazy Research	Snorkel [SuperGLUE v1.9]		-	-	88.6/93.2	76.2	76.4/36.3	-	78.9	72.1	72.6	47.6	-



BERT (Devlin et al., 2018)



T5 (Raffel et al., 2020)

Current AI Paradigm: Pre-trained Deep Neural Networks (achieving human parity on SQuAD 2.0/CoQA)

SQuAD2.0

The Stanford Question Answering Dataset

What is SQuAD?

Stanford Question Answering Dataset (SQuAD) is a reading comprehension dataset, consisting of questions posed by crowdworkers on a set of Wikipedia articles, where the answer to every question is a segment of text, or *span*, from the corresponding reading passage, or the question might be unanswerable.

SQuAD2.0 combines the 100,000 questions in SQuAD1.1 with over 50,000 unanswerable questions written adversarially by crowdworkers to look similar to answerable ones. To do well on SQuAD2.0, systems must not only answer questions when possible, but also determine when no answer is supported by the paragraph and abstain from answering.

[Explore SQuAD2.0 and model predictions](#)

[SQuAD2.0 paper \(Rajpurkar & Jia et al. '18\)](#)

SQuAD 1.1, the previous version of the SQuAD dataset, contains 100,000+ question-answer pairs on 500+ articles.

[Explore SQuAD1.1 and model predictions](#)

[SQuAD1.0 paper \(Rajpurkar et al. '16\)](#)

Leaderboard

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1	SA-Net on Albert (ensemble) QIANXIN	90.724	93.011
2	SA-Net-V2 (ensemble) QIANXIN	90.679	92.948
2	Retro-Reader (ensemble) Shanghai Jiao Tong University http://arxiv.org/abs/2001.09694	90.578	92.978
3	ATRLP+PV (ensemble) Hithink RoyalFlush	90.442	92.877
3	ELECTRA+ALBERT+EntitySpanFocus (ensemble) SRCB_DML	90.442	92.839
4	ELECTRA+ALBERT+EntitySpanFocus (ensemble) SRCB_DML	90.420	92.799
4	EntitySpanFocus+AT (ensemble) RICOH_SRCB_DML	90.454	92.748

CoQA



A Conversational Question Answering Challenge

What is CoQA?

CoQA is a large-scale dataset for building Conversational Question Answering systems. The goal of the CoQA challenge is to measure the ability of machines to understand a text passage and answer a series of interconnected questions that appear in a conversation. CoQA is pronounced as **coca** .

[CoQA paper](#)

CoQA contains 127,000+ questions with answers collected from 8000+ conversations. Each conversation is collected by pairing two crowdworkers to chat about a passage in the form of questions and answers. The unique features of CoQA include 1) the questions are conversational; 2) the answers can be free-form text; 3) each answer also comes with an evidence subsequence highlighted in the passage; and 4) the passages are collected from seven diverse domains. CoQA has a lot of challenging phenomena not present in existing reading comprehension datasets, e.g., coreference and pragmatic reasoning.

Download

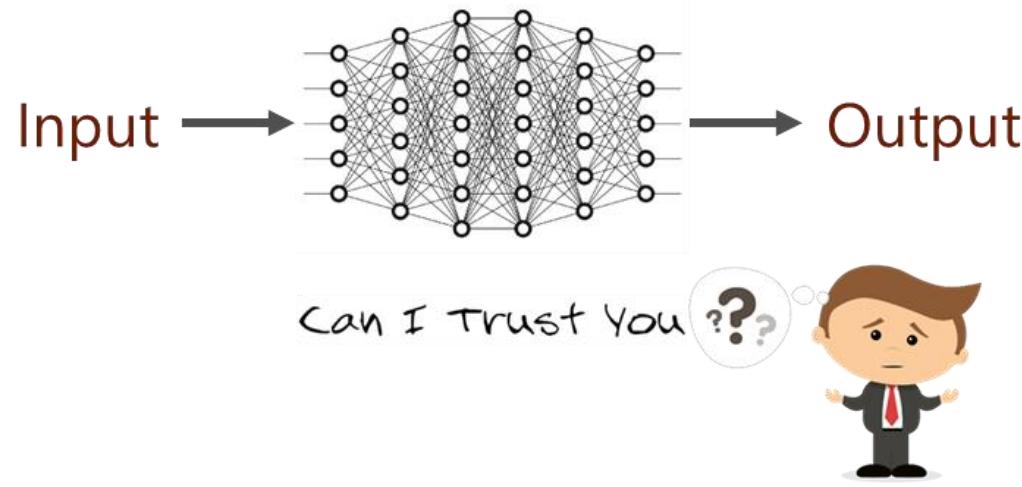
Browse the examples in CoQA:

[Browse CoQA](#)

Leaderboard

Rank	Model	In-domain	Out-of-domain	Overall
	Human Performance Stanford University (Reddy & Chen et al. TACL '19)	89.4	87.4	88.8
1	RoBERTa + AT + KD (ensemble) Zhuiyi Technology https://arxiv.org/abs/1909.10772	91.4	89.2	90.7
1	TR-MT (ensemble) WeChatAI	91.5	88.8	90.7
2	RoBERTa + AT + KD (single model) Zhuiyi Technology https://arxiv.org/abs/1909.10772	90.9	89.2	90.4
3	TR-MT (ensemble) WeChatAI	91.1	87.9	90.2
4	Google SQuAD 2.0 + MMFT (ensemble) MSRA + SDRG	89.9	88.0	89.4
5	TR-MT (single model) WeChatAI	90.4	86.8	89.3
6	XLNet + Augmentation (single model) Xiaoming	89.9	86.9	89.0

However, shortcomings still exist.

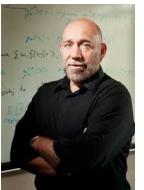


- The predictions of such models are not interpretable and trustable. } **interpretability issue**
- The working mechanisms of such models are hard to be analyzed or understood. } **knowledge issue**
- Such models are usually hard to be integrated with existing knowledge, such as rules, common sense and knowledge graphs. } **knowledge issue**
- ...

What's the POSSIBLE solution?

- a) machine reasoning as an interpretable and knowledge-based decision-making process
- b) ...
- c) ...
- d) ...

What is Machine Reasoning (defined by others)



- Reason is the capacity of consciously making sense of things, applying logic, and adapting or justifying practices, institutions, and beliefs based on new or existing information. —[Wikipedia](#)
- A program has common sense (reasoning) if it automatically deduces for itself a sufficiently wide class of immediate consequences of anything it is told and what it already knows. —[John McCarthy, 1958](#)
- Bayesian methods provide a formalism for reasoning about partial beliefs under conditions of uncertainty. In this formalism, propositions are given numerical parameters signifying the degree of belief accorded them under some body of knowledge, and the parameters are combined and manipulated according to the rules of probability theory. —[Judea Pearl, 1988](#)
- The Learning to Reason theory developed here is concerned with studying the entire process of learning a knowledge base representation and reasoning with it. —[Roni Khardon and Dan Roth, 1994](#)
- A plausible definition of "reasoning" could be "algebraically manipulating previously acquired knowledge in order to answer a new question". —[Léon Bottou, 2011](#)
- Such a nugget of information or knowledge seems to fit well as a conscious state. Combining such conscious states sequentially in order to make more complex predictions and inferences or actions is basically what reasoning is about. —[Yoshua Bengio, 2019](#)

Knowledge

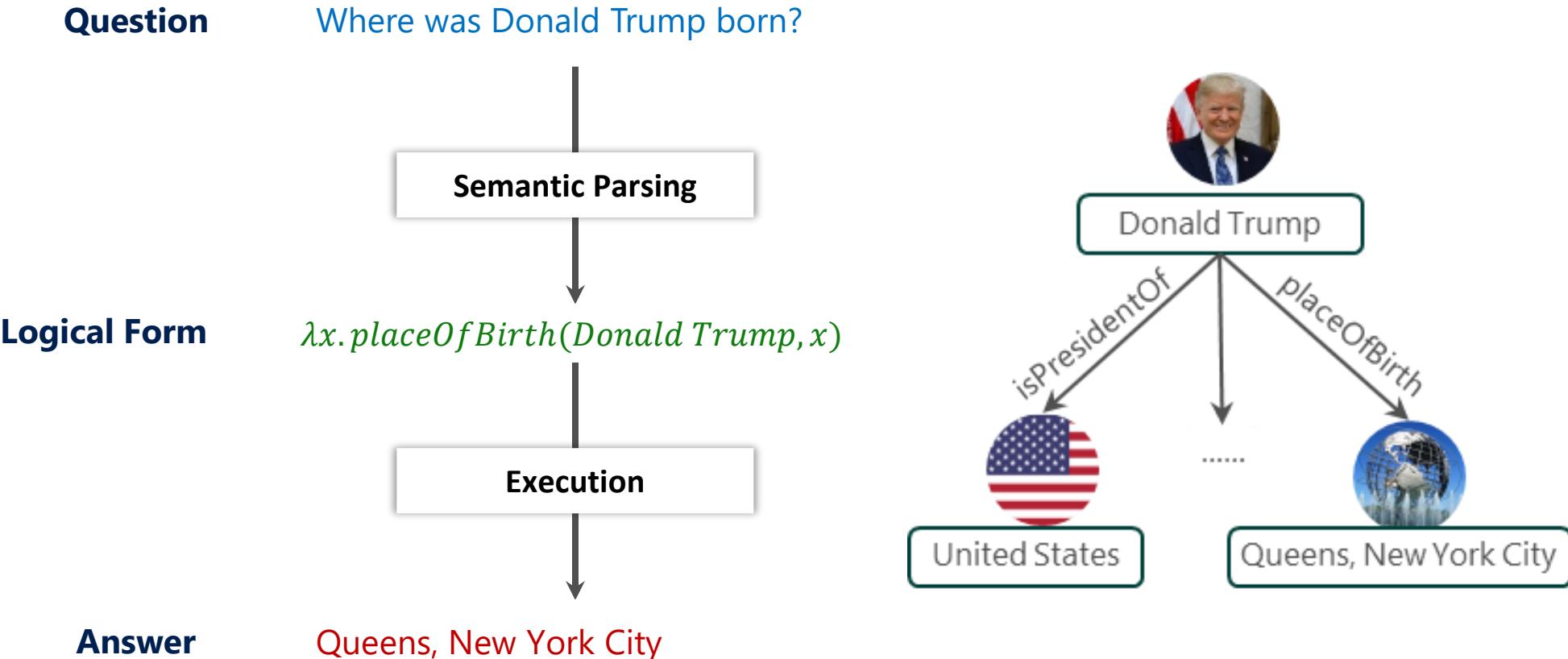
Inference Algorithms

Interpretability

What is Machine Reasoning (defined by this tutorial)

An **interpretable** decision-making process that can solve problems or draw conclusions from **what the system is told (i.e. facts and observations)** and **already knows (i.e. models, common sense and knowledge)** under certain constraints.

Single-turn Knowledge-based QA



Multi-turn Knowledge-based QA



Tell me the movies with Tom Hanks and Meg Ryan

$$\lambda x. \text{film_film_actor}(x, \text{Tom Hanks}) \wedge \text{film_film_actor}(x, \text{Meg Ryan})$$


Sleepless in Seattle, You've Got Mail,...



When was **he** born ?

$$\lambda x. \text{people_person_dateofbirth}(\text{Tom Hanks}, x)$$

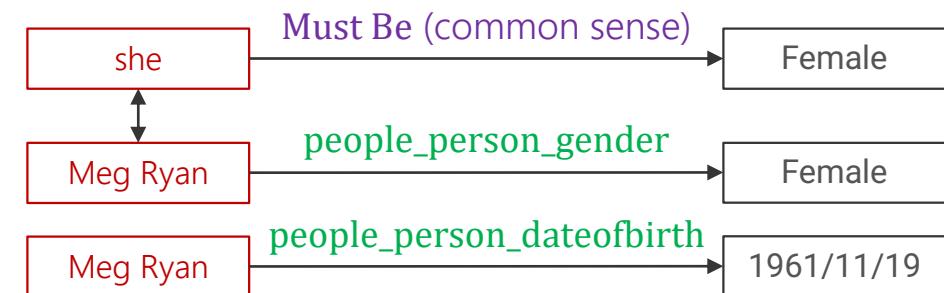
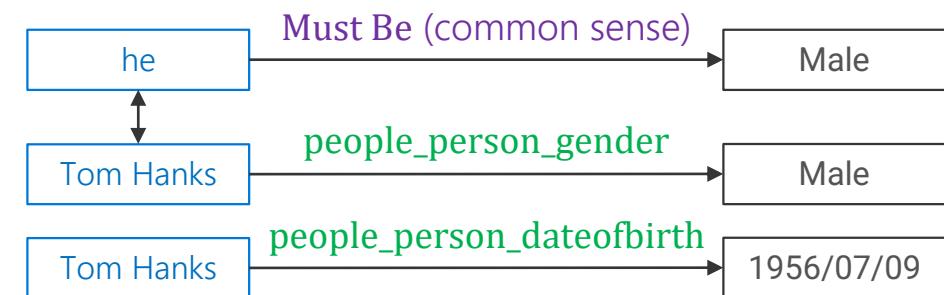

1956/07/09



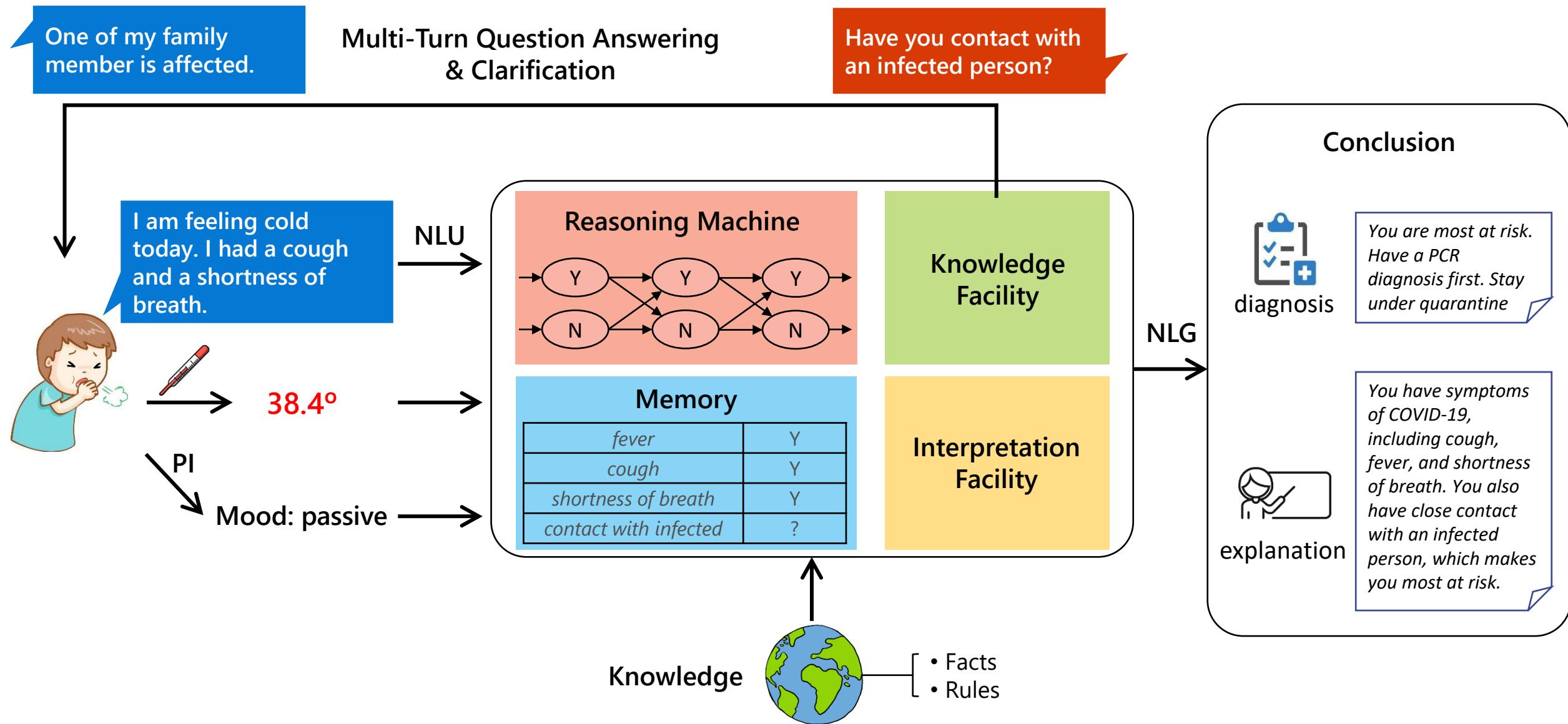
How about her ?

$$\lambda x. \text{people_person_dateofbirth}(\text{Meg Ryan}, x)$$

1961/11/19



Domain-Specific Expert System



Fact Verification

Claim: The Rodney King riots took place in the most populous county in the USA.

Reasoning needs fact knowledge extracted from evidence sentences

TRUE

Evidence #1:

The 1992 Los Angeles riots, also known as the Rodney King riots were a series of riots, lootings, arsons, and civil disturbances that occurred in Los Angeles County, California in April and May 1992.

Evidence #2:

Los Angeles County, officially the County of Los Angeles, is the most populous county in the USA.



4

3

1

2

5



Law School Admission Test (LSAT)

A university library budget committee must reduce exactly five of eight areas of expenditure—G, L, M, N, P, R, S, and W—in accordance with the following conditions:

1. If both G and S are reduced, W is also reduced.

2. If N is reduced, neither R nor S is reduced.

3. If P is reduced, L is not reduced.

4. Of the three areas L, M, and R, exactly two are reduced.

Question 1

If both M and R are reduced, which one of the following is a pair of areas neither of which could be reduced?

A. G, L

B. G, N

C. L, N

D. L, P

E. P, S

- Reduced: M, R
- Not reduced: L
- Not reduced: L, N

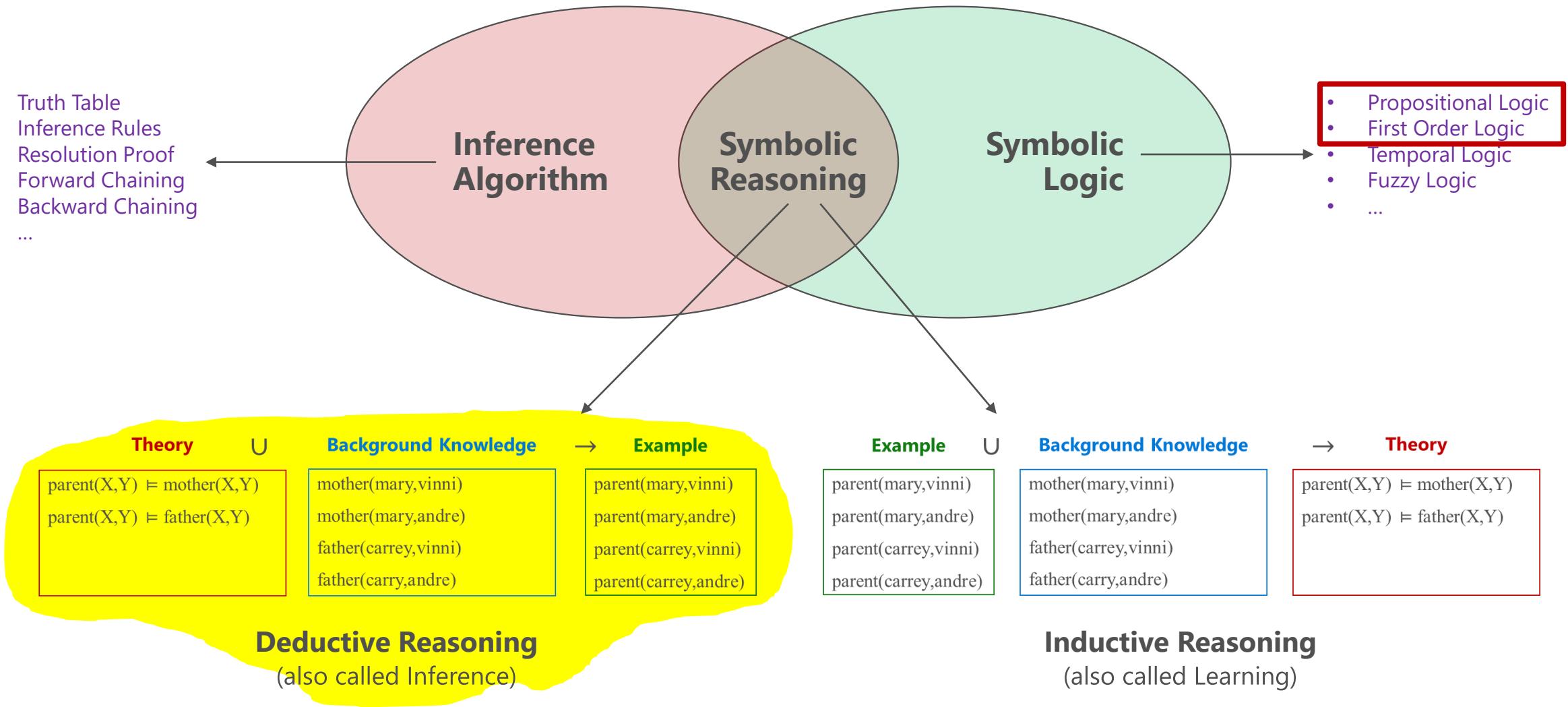
Tutorial Outline

- **Machine Reasoning Frameworks**
 - Symbolic Reasoning
 - Probabilistic Reasoning
 - Neural-Symbolic Reasoning
 - Neural-Evidence Reasoning
- **Dilemma: Interpretability vs. Performance**
 - Empirical Success of Pre-trained Models
 - Efforts on Interpretability AI
- **Summary**

Symbolic Reasoning

Symbolic Reasoning

- Truth Table
- Inference Rules
- Resolution Proof
- Forward Chaining
- Backward Chaining
- ...



Propositional Logic

- **Proposition**

- A proposition is defined as a declarative sentence that is either **True** or **False**.

- **Propositional Logic**

- The area of logic which deals with propositions.

- **Syntax of Propositional Logic**

- **Constants**
 - True, False
 - **Propositional Symbols**
 - Birds can fly, $1+1=2$, ...
 - **Connectives**
 - \neg (negation), \wedge (conjunction), \vee (disjunction), \Rightarrow (implication), \Leftrightarrow (equivalence)
 - **Atomic Sentences**
 - Constructed from constants and propositional symbols
 - **Composite Sentences**
 - Constructed from atomic sentences via connectives

Propositional Logic Reasoning: Entailment

Infer whether a knowledge base KB can entail a sentence α : $KB \models \alpha$.

KB_1 : *I will go to the beach (r) only if it is sunny today (p).* $r \Rightarrow p$

KB_2 : *If I don't go to the beach, then I will go to the cinema (s).* $\neg r \Rightarrow s$

KB_3 : *If I go the cinema, then I will be home before 6PM (t).* $s \Rightarrow t$



α : It is not sunny today, so I will be home before 6PM. $\neg p \Rightarrow t$

Logical Rules

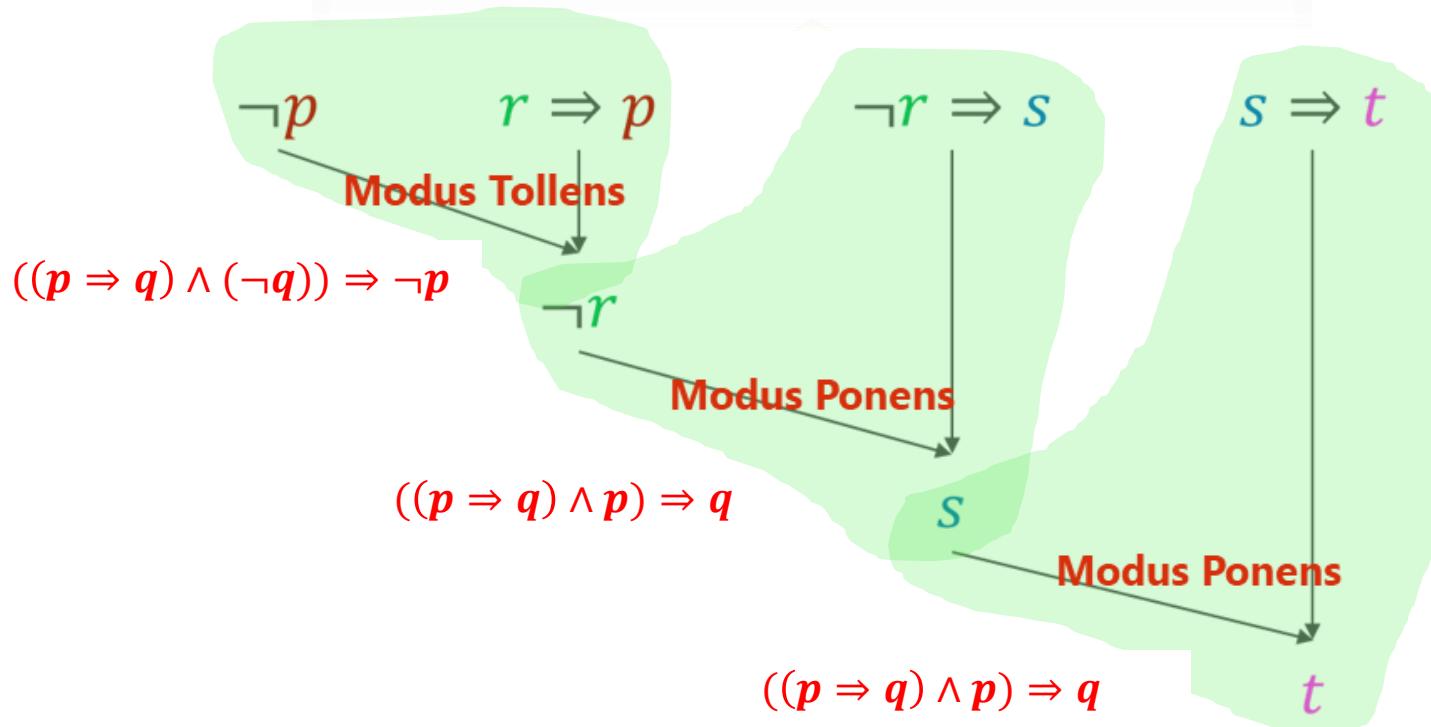
1.	$p \wedge q \equiv q \wedge p$	commutativity of \wedge
2.	$p \vee q \equiv q \vee p$	commutativity of \vee
3.	$(p \wedge q) \wedge r \equiv p \wedge (q \wedge r)$	associativity of \wedge
4.	$(p \vee q) \vee r \equiv p \vee (q \vee r)$	associativity of \vee
5.	$\neg(\neg p) \equiv p$	double-negation elimination
6.	$p \Rightarrow q \equiv \neg q \Rightarrow \neg p$	contraposition
7.	$p \Rightarrow q \equiv \neg p \vee q$	implication elimination
8.	$p \Leftrightarrow q \equiv (p \Rightarrow q) \wedge (q \Rightarrow p)$	biconditional elimination
9.	$\neg(p \wedge q) \equiv \neg p \vee \neg q$	De Morgan
10.	$\neg(p \vee q) \equiv \neg p \wedge \neg q$	De Morgan
11.	$(p \wedge (q \vee r)) \equiv ((p \wedge q) \vee (p \wedge r))$	distributivity of \wedge over \vee
12.	$(p \vee (q \wedge r)) \equiv ((p \vee q) \wedge (p \vee r))$	distributivity of \vee over \wedge
13.	$(p_1 \wedge \dots \wedge p_n) \Rightarrow p_i$	And-elimination
14.	$p_i \Rightarrow (p_1 \wedge \dots \wedge p_n)$	Or-introduction
15.	$((p \Rightarrow q) \wedge p) \Rightarrow q$	Modus Ponens
16.	$((p \Rightarrow q) \wedge (\neg q)) \Rightarrow \neg p$	Modus Tollens
17.	$((p \vee q) \wedge (\neg q \vee r)) \Rightarrow p \vee r$	Resolution

Inference with Logical Rules

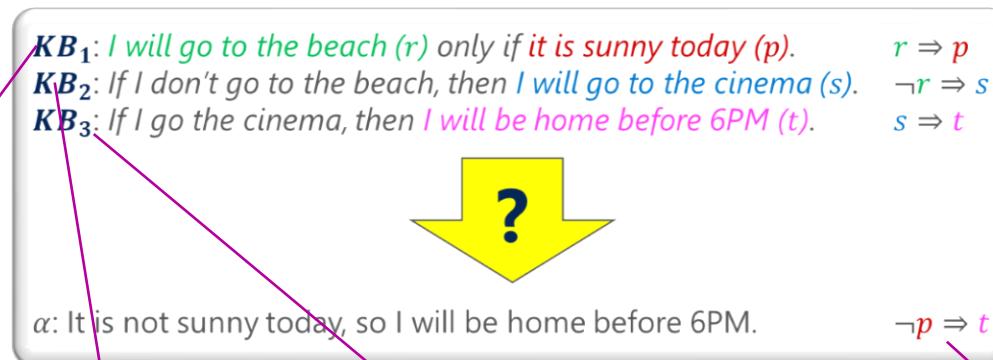
KB₁: I will go to the beach (r) only if it is sunny today (p). $r \Rightarrow p$
KB₂: If I don't go to the beach, then I will go to the cinema (s). $\neg r \Rightarrow s$
KB₃: If I go to the cinema, then I will be home before 6PM (t). $s \Rightarrow t$



α : It is not sunny today, so I will be home before 6PM. $\neg p \Rightarrow t$



Inference with Resolution Proof

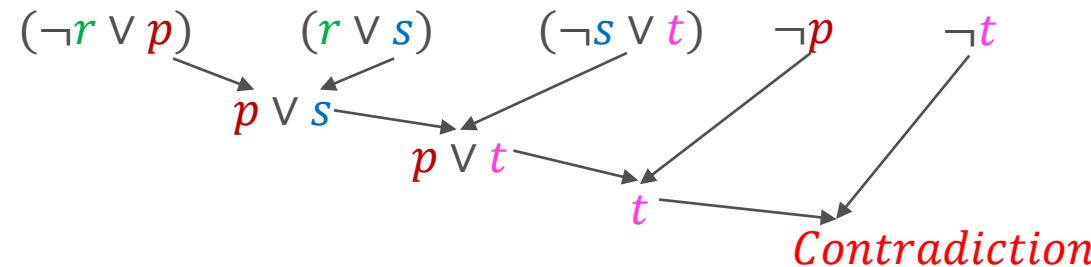


Resolution
 $((p \vee q) \wedge (\neg q \vee r)) \Rightarrow p \vee r$

1. Convert KB and the input sentence to Conjunctive Normal Forms (CNFs);

$$r \Rightarrow p \rightarrow (\neg r \vee p); \quad \neg r \Rightarrow s \rightarrow (r \vee s); \quad s \Rightarrow t \rightarrow (\neg s \vee t); \quad \neg(\neg p \Rightarrow t) \rightarrow \neg p \wedge \neg t;$$

2. Starting from KB and $\neg\alpha$, apply the resolution rule iteratively, until Contradiction is reached (proved), or no more sentences can be derived (disproved).



The Expressive Power of Propositional Logic is Limited.

To express “**every cat likes fish**”, propositional logic must include the statements about every cat.



Garfield is a cat.
Garfield likes fish.

...



Tom is a cat.
Toms likes fish.

...



Sylvester is a cat.
Sylvester likes fish.

...

From Propositional Logic to First Predicate Logic

1. *Introduce objects, their properties and relations.*
2. *Introduce variables to refer to arbitrary objects.*
3. *Introduce quantifiers to make statements over groups objects.*

every cat likes fish \leftrightarrow $\forall x \text{ cat}(x) \Rightarrow \text{like}(x, \text{fish})$

First Order Logic

▪ Syntax of First Order Logic

- Constants (a type of Term)
 - Represent specific entities, such as *China*, *Microsoft*
- Variables (a type of Term)
 - Represent entities of a certain type, such as x , y , z
- Functions (a type of Term)
 - Applied to one or more terms
- Predicates
 - Functions from one or more terms to a Boolean, such as *PlaceOfBirth(BillGates, Seattle) = true*
- Literal
 - Any predicate (or its negation) applied to any set of terms, such as *Female(Mary)*, $\neg\text{Female}(x)$
- Clause
 - Any disjunction of literals whose variables are universally quantified, such as $\forall x: \text{Female}(x) \vee \text{Male}(x)$
- Connectives
 - \neg (negation), \wedge (conjunction), \vee (disjunction), \Rightarrow (implication), \Leftrightarrow (equivalence)
- Quantifiers
 - \forall (Universal), \exists (Existential)
- Atomic Sentences
 - Constructed from a predicate followed by a parenthesis with a sequence of terms
- Composite Sentences
 - Constructed by combining atomic sentences using connectives

First Order Logic Reasoning: Entailment

Infer whether a knowledge base KB can entail a sentence α : $KB \models \alpha$.

KB_1 : $\forall x \text{ cat}(x) \Rightarrow \text{like}(x, \text{fish})$

//all cats like fish.

KB_2 : $\forall x \forall y (\text{cat}(x) \wedge \text{like}(x, y)) \Rightarrow \text{eat}(x, y)$

//cats eat everything they like.

KB_3 : $\text{cat}(\text{Tom})$

//tom is a cat.



α : $\text{eat}(\text{Tom}, \text{fish})$

Straightforward solution: propositionalizing a first order logic KB to a propositional logic KB.

Issue: with n k -ary predicates and m constants, the number of instantiations will be $n \cdot m^k$.

Alternative solution: Forward Chaining Algorithm and Backward Chaining Algorithm.

Substitution & Generalized Modus Ponens (GMP)

▪ Substitution

- Given $\theta = \{x_1/t_1, \dots, x_n/t_n\}$ and an expression f , substitution $f\theta = \text{SUBST}(\theta, f(x_1, \dots, x_n, x_{n+1}, \dots)) = f(t_1, \dots, t_n, x_{n+1}, \dots)$ replaces all occurrences of " x_i " with " t_i " in the expression f .
- An example
 - $\text{SUBST}(\{x/Tom\}, \text{like}(x, \text{fish})) = \text{like}(Tom, \text{fish})$

▪ Generalized Modus Ponens (GMP)

- Given $n + 1$ terms f_1, \dots, f_n , $(p_1 \wedge \dots \wedge p_n \Rightarrow q)$, if there exists a substitution θ such that $f_i\theta = p_i\theta$ for all $i = 1, 2, \dots, n$, then $(f_1 \wedge \dots \wedge f_n) \wedge (p_1 \wedge \dots \wedge p_n \Rightarrow q) \Rightarrow q\theta$.
- An example
 - Facts: $f_1 = \text{cat}(Tom)$ and $f_2 = \text{like}(Tom, \text{fish})$
 - Rule: $\text{cat}(x) \wedge \text{like}(x, y) \Rightarrow \text{eat}(x, y)$
 - Substitution: $\theta = \{x/Tom, y/fish\}$
 - Result of GMP: $\text{eat}(Tom, \text{fish})$

Inference with Forward Chaining

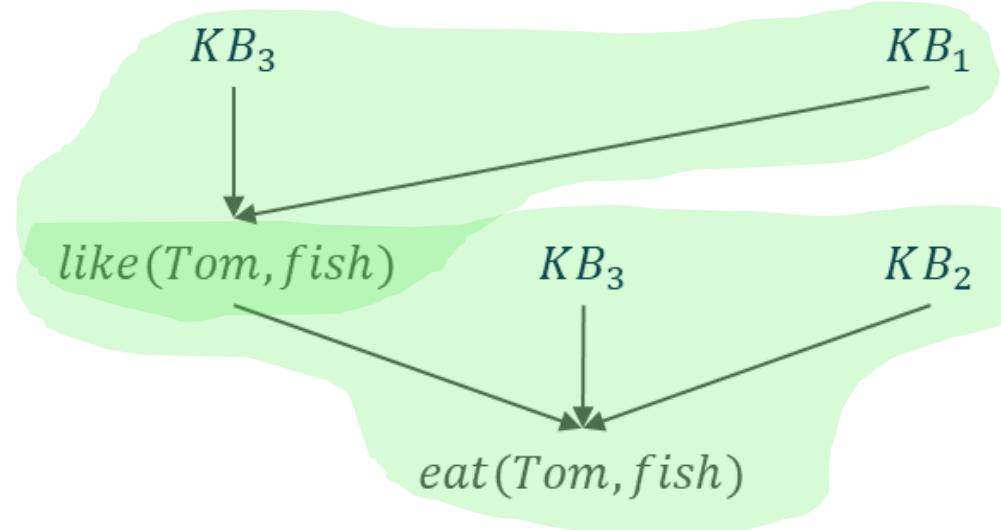
KB_1 : $\forall x \text{ cat}(x) \Rightarrow \text{like}(x, \text{fish})$ //all cats like fish.
 KB_2 : $\forall x \forall y (\text{cat}(x) \wedge \text{like}(x, y)) \Rightarrow \text{eat}(x, y)$ //cats eat everything they like.
 KB_3 : $\text{cat}(\text{Tom})$ //tom is a cat.



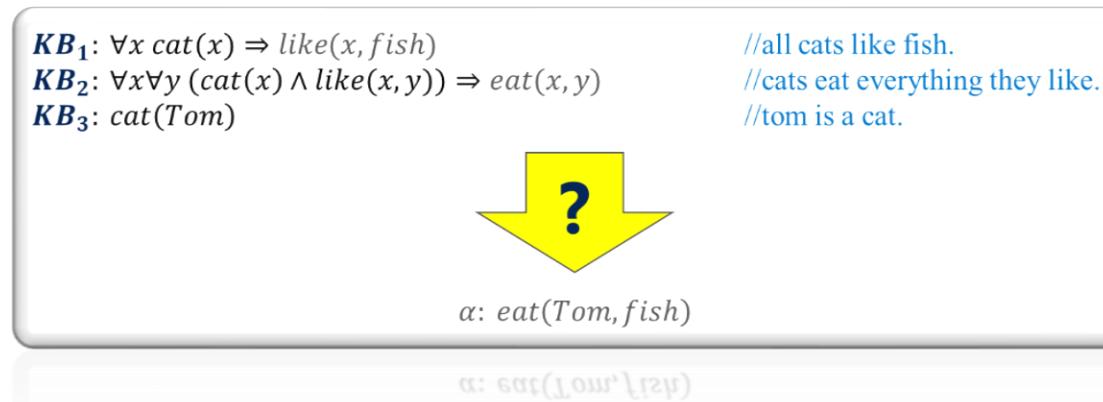
$\alpha: \text{eat}(\text{Tom}, \text{fish})$

$\alpha: \text{eat}(\text{Tom}, \text{fish})$

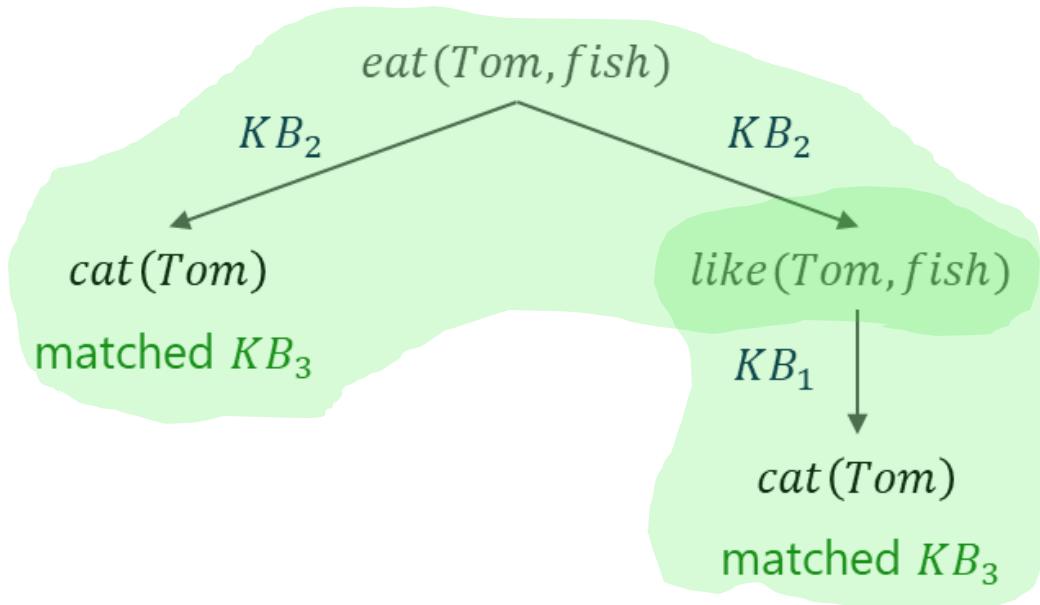
Start with atomic sentences in KB and apply GMP in the forward direction to extract more data until the goal is reached.



Inference with Backward Chaining



Start with the goal and work backward to find known facts that support the goal.



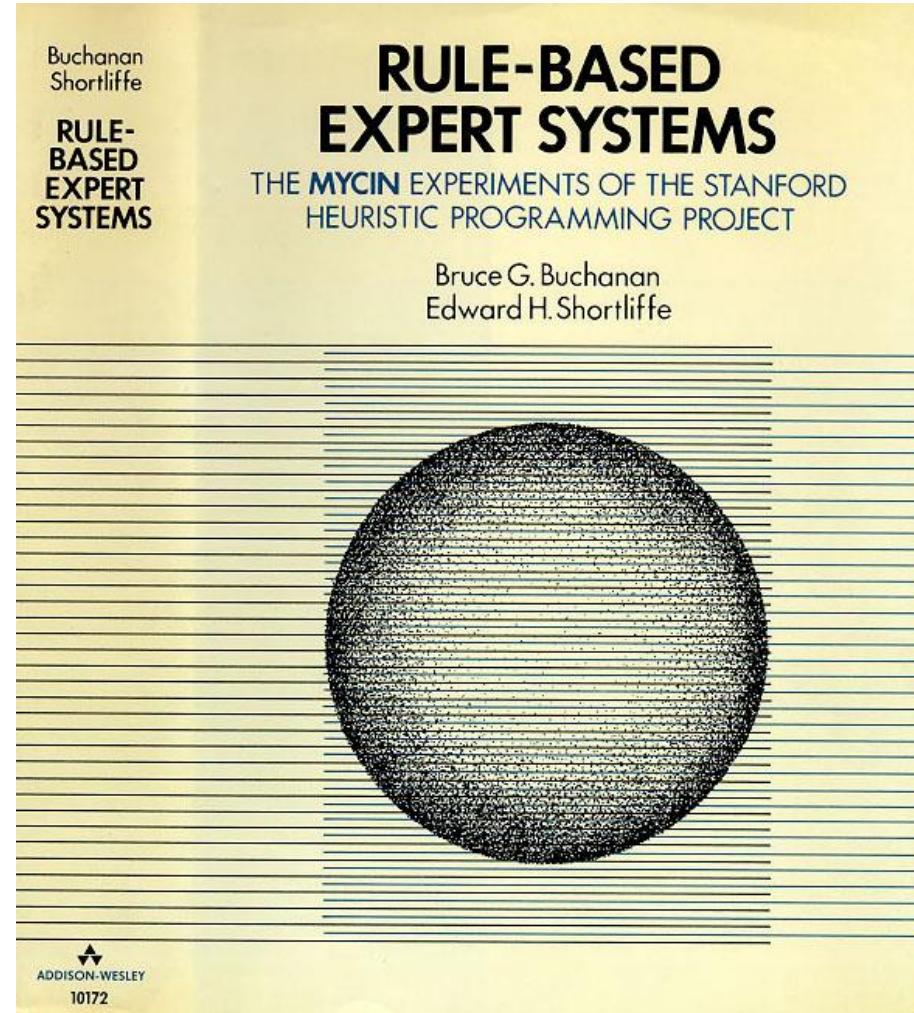
Successful Application: MYCIN (Expert System)

- An expert system developed at Stanford by Ed Feigenbaum, Bruce Buchanan, Edward Shortliffe in the 1970s, designed to diagnose blood infections.

- Hundreds of rules such as:

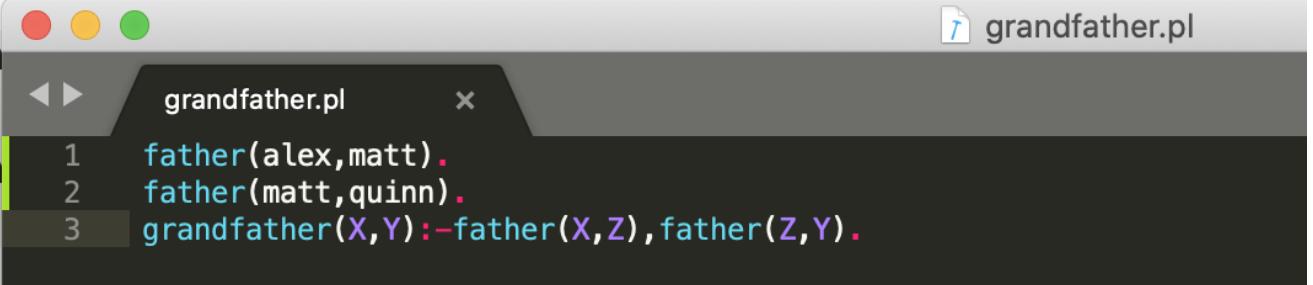
If: (1) the strain of the organism is gram-positive, and
(2) the morphology of the organism is coccus, and
(3) the growth conformation of the organism is clumps,
then there is suggestive evidence (0.7) that the identity of the organism is staphylococcus.

- "I said that MYCIN is the granddaddy of expert systems." -- Allen Newell (1984)



Successful Application: Prolog

- Prolog is a logic programming language widely used in first order logic-based scenarios.
- Syntax of Prolog
 - Predicates, objects and functions
 - sibling, parent_child, father_child, mother_child
 - trude, sally, tom, erica, mike
 - Variables
 - X, Y
 - Facts
 - mother_child(trude, sally)
 - father_child(tom, sally)
 - father_child(tom, erica)
 - father_child(mike, tom)
 - Rules (Knowledge)
 - sibling(X, Y) :- parent_child(Z, X), parent_child(Z, Y)
 - parent_child(X, Y) :- father_child(X, Y)
 - parent_child(X, Y) :- mother_child(X, Y)
 - Queries
 - ?- sibling(sally, erica) YES



```
grandfather.pl
1 father(alex,matt).
2 father(matt,quinn).
3 grandfather(X,Y):-father(X,Z),father(Z,Y).
```



```
?- father(K,matt).
K = alex.

?- grandfather(V,quinn).
V = alex .

?- grandfather(alex,quinn).
true.

?- |
```

Benefits & Limitations

- Benefits

- Good interpretability
- Easy to incorporate domain knowledge
- Doesn't require to be trained on huge amounts of data

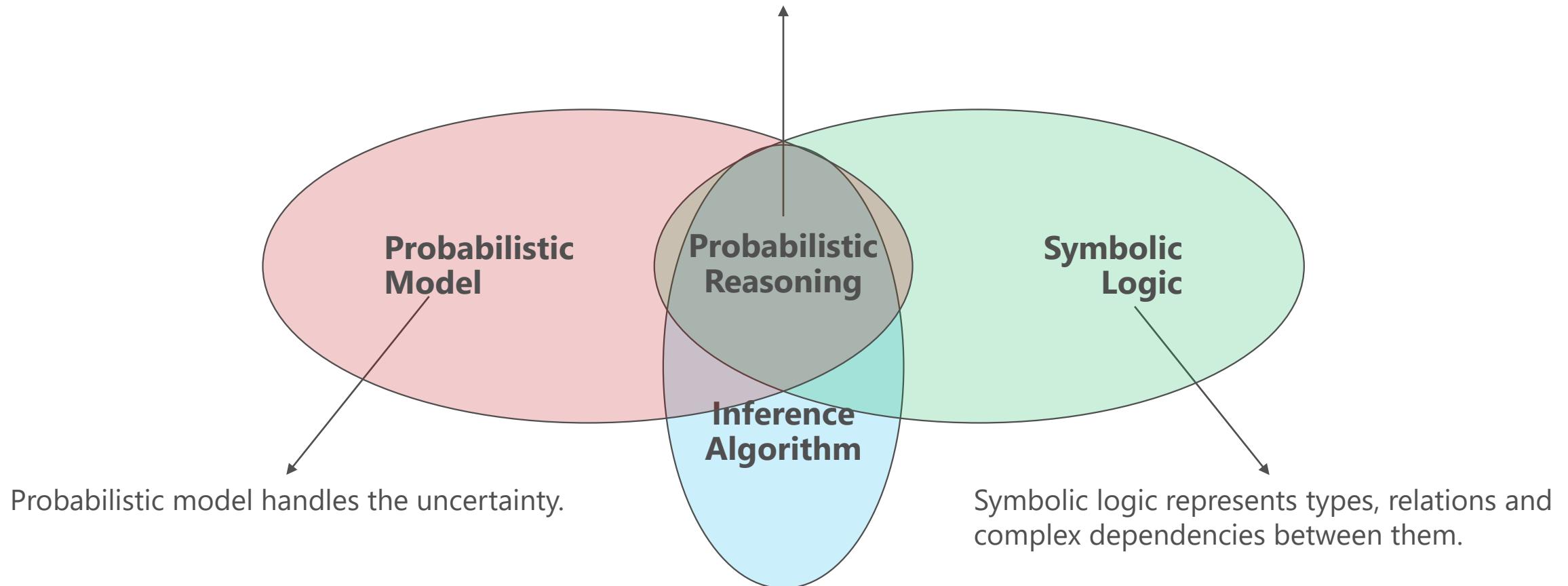
- Challenges

- Cannot deal with data uncertainty or noise
- Cannot make use of similarities between predicates or constants in training data
- Require hand-crafted logical rules

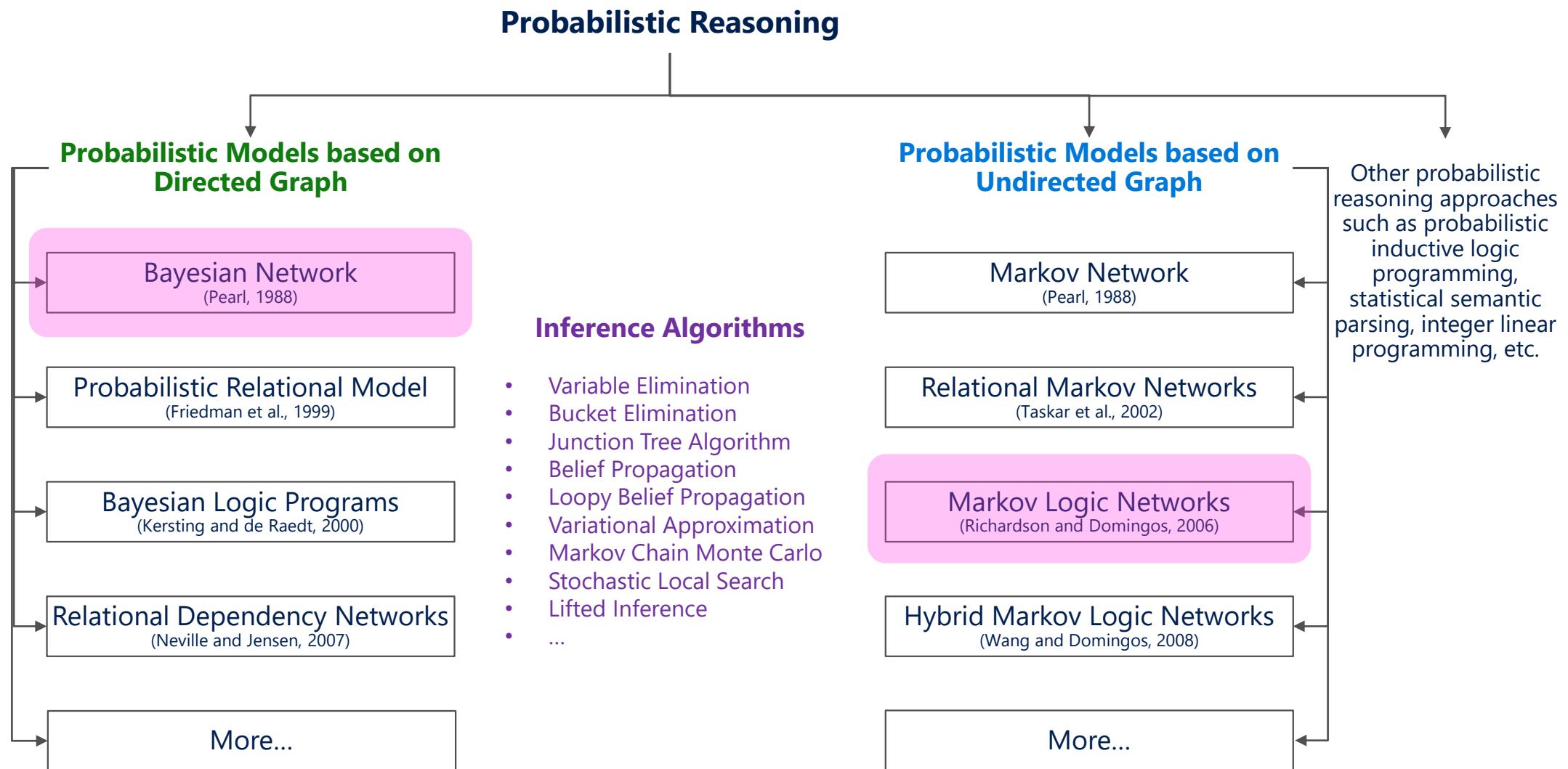
Probabilistic Reasoning

Probabilistic Reasoning

A probabilistic reasoning system integrates probabilistic models with symbolic logic.



Probabilistic Reasoning Models



Bayesian Network (BN)

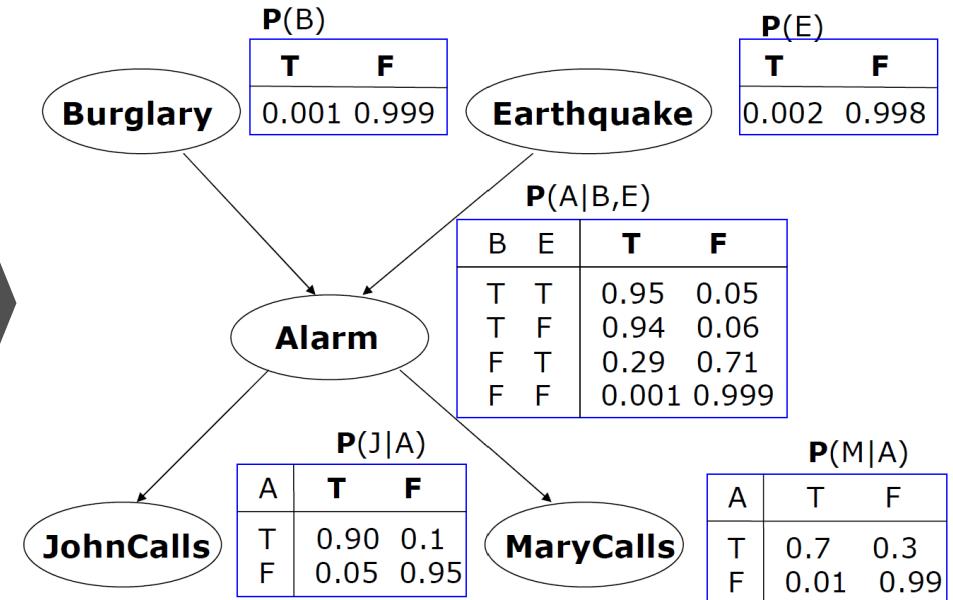
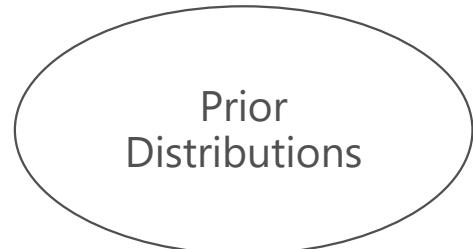
- A Bayesian Network $BN = \langle G, \Theta \rangle$ is a directed acyclic graph.
- A Bayesian Network represents a joint probability distribution over a set of random variables.
- The nodes represent the random variables.
- The edges represent direct dependencies between the nodes.
- $\theta_{x_i|\text{pa}(X_i)} = P(x_i|\text{pa}(X_i))$ is a parameter for each x_i of X_i conditioned on $\text{pa}(X_i)$, the parents of X_i .

$$P_{BN}(X_1, \dots, X_n) = \prod_{i=1}^n P(x_i|\text{pa}(X_i)) = \prod_{i=1}^n \theta_{X_i|\text{pa}(X_i)}$$

An Example (from Judea Pearl)

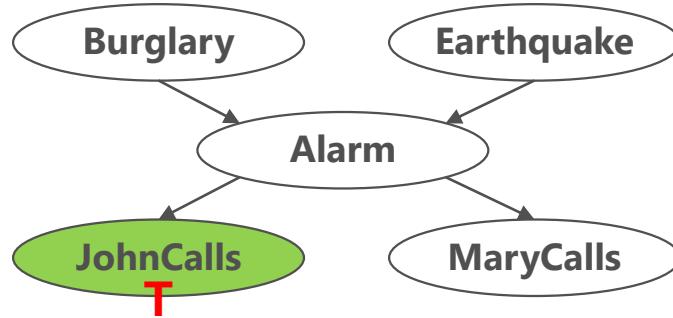
Assume your house has an **alarm system** against **burglary**. You live in the seismically active area and the alarm system can get occasionally set off by an **earthquake**. You have two neighbors, **Mary** and **John**, who do not know each other. If they hear the alarm, they call you, but this is not guaranteed.

- Burglary
- Earthquake
- Burglary \Rightarrow Alarm
- Earthquake \Rightarrow Alarm
- Alarm \Rightarrow JohnCalls
- Alarm \Rightarrow MaryCalls

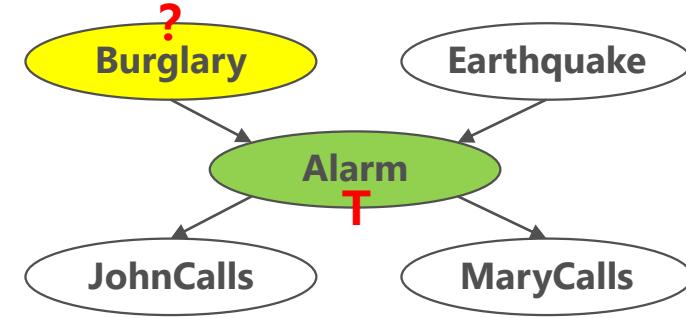


Typical Reasoning Tasks

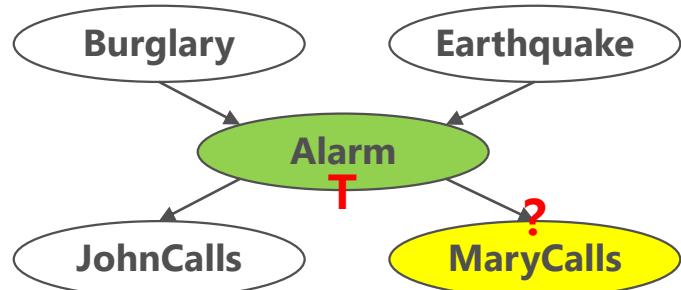
 observed variables
 unobserved variables
 of interests



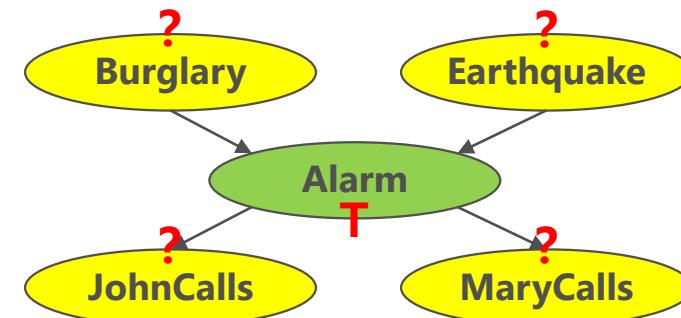
Likelihood: $P(J = T) ?$



Diagnosis: $P(B|A = T) ?$



Prediction: $P(M|A = T) ?$



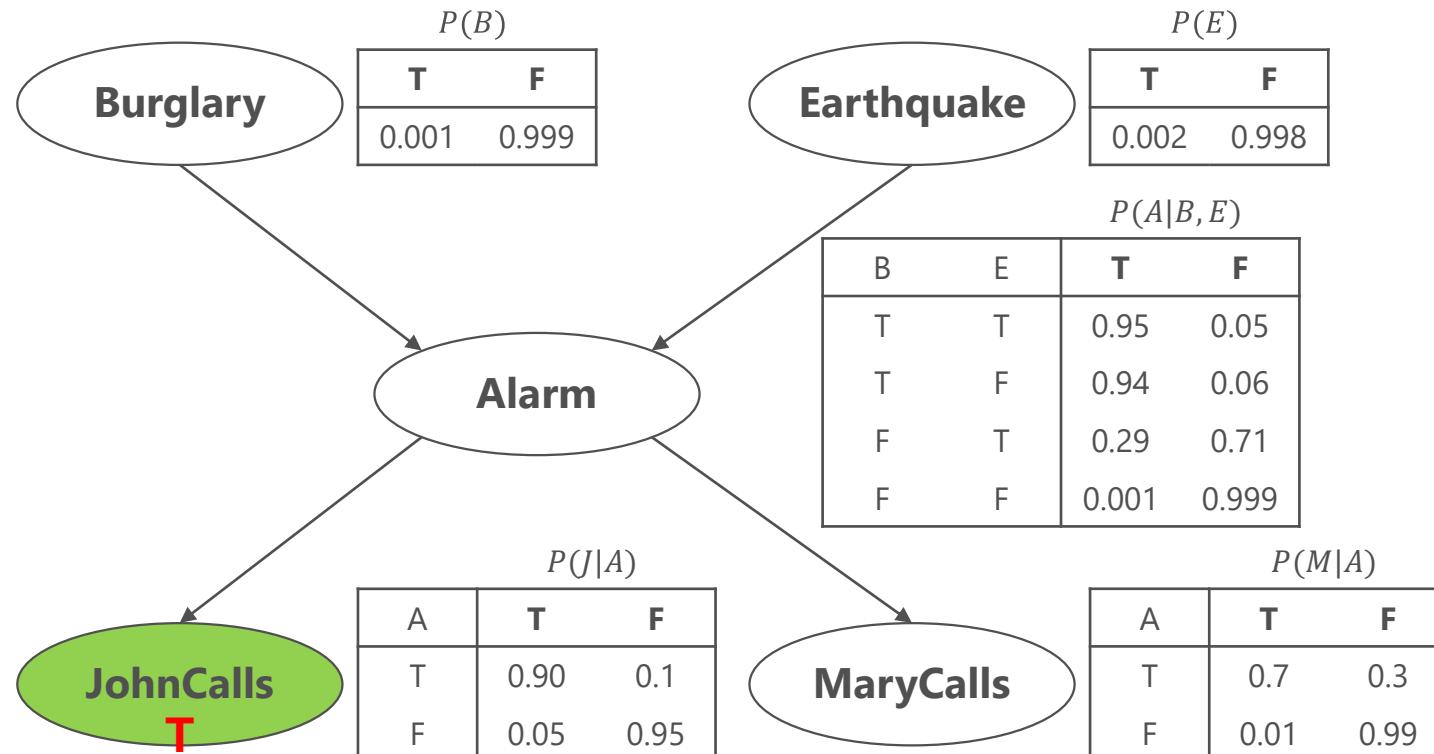
Maximum Probable Explanation:

$$\max_{B,E,J,M} P(B, E, J, M | A = T) ?$$

Reasoning Task (1): Likelihood

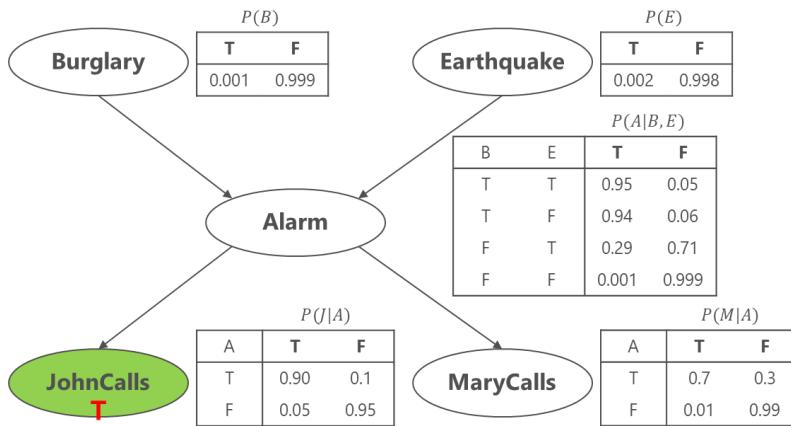
Compute the probability of the observed variables X .

$$P(X) = \sum_{\mathbf{Z}} P(\mathbf{Z}, X) = \sum_{z_1} \dots \sum_{z_k} P(z_1 = z_1, \dots, z_k = z_k, X)$$



What is the probability $P(J=T)$?

Inference with Variable Elimination



$$\begin{aligned}
 P(J = T) &= \sum_a \sum_b \sum_e \sum_m P(B = b)P(E = e)P(A = a | B = b, E = e)P(J = T | A = a)P(M = m | A = a) \\
 &= \sum_a \sum_b \sum_e P(B = b)P(E = e)P(A = a | B = b, E = e)P(J = T | A = a) \sum_m P(M = m | A = a) \\
 &= \sum_a \sum_b \sum_e P(B = b)P(E = e)P(A = a | B = b, E = e)P(J = T | A = a) f_m(A = a) \\
 &= \sum_a \sum_b P(B = b)P(J = T | A = a) \sum_e P(E = e)P(A = a | B = b, E = e) \\
 &= \sum_a \sum_b P(B = b)P(J = T | A = a) f_e(A = a, B = b) \\
 &= \sum_a P(J = T | A = a) \sum_b P(B = b) f_e(A = a, B = b) \\
 &= \sum_a P(J = T | A = a) f_b(A = a)
 \end{aligned}$$

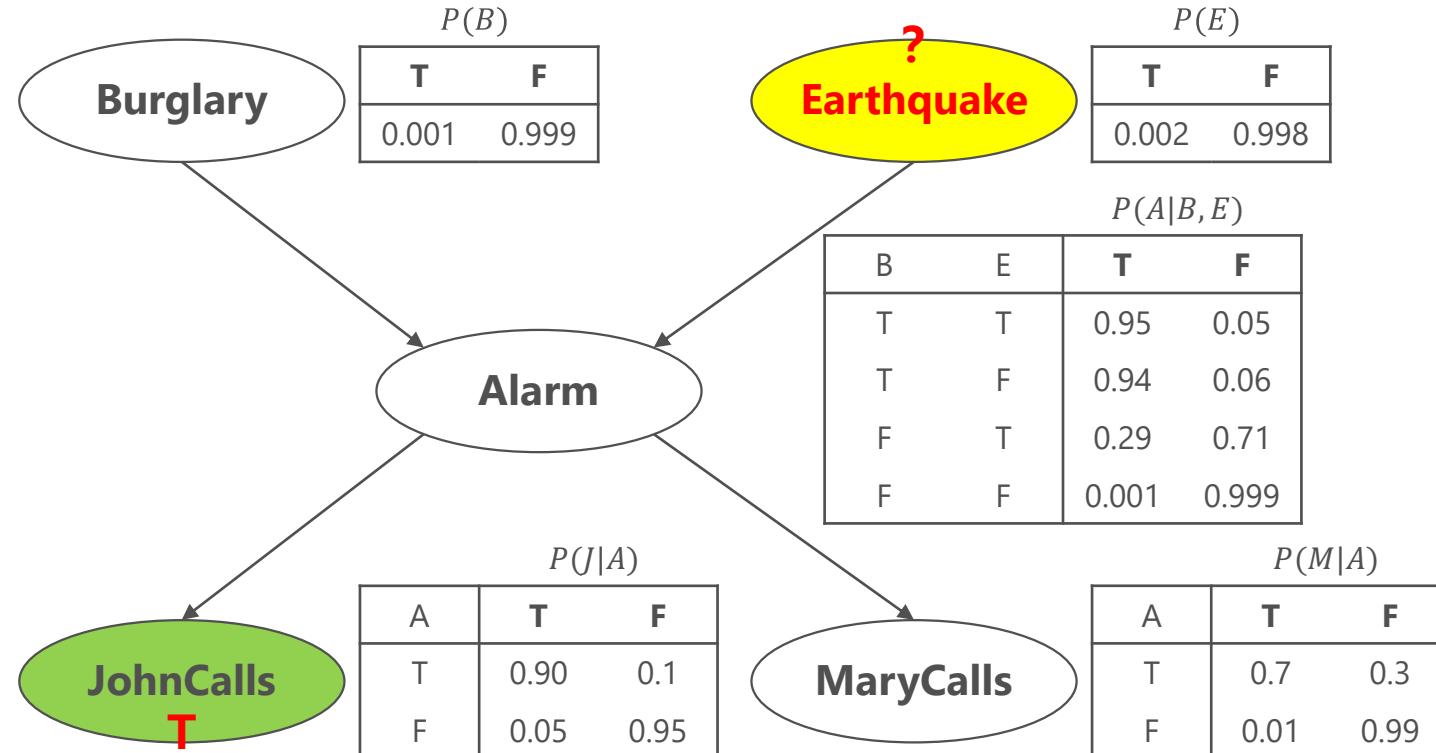
Annotations:

- A red oval highlights the term $\sum_a \sum_b \sum_e$.
- A purple arrow points from $f_m(A = a)$ to the term $\sum_m P(M = m | A = a)$.
- A blue arrow points from $f_e(A = a, B = b)$ to the term $\sum_e P(E = e)P(A = a | B = b, E = e)$.
- A red arrow points from $f_b(A = a)$ to the term $\sum_b P(B = b)$.

Reasoning Task (2): Belief Updating

Compute the conditional probability distribution of a hidden variable Z given the observed variables X .

$$P(Z = z_i | X) = \frac{P(Z = z_i, X)}{P(X)}$$



What is the conditional probability distribution $P(E|J=T)$?

Inference with Markov Chain Monte Carlo (MCMC)

initialize starting values for variables X_1, \dots, X_n in BN

do until convergence:

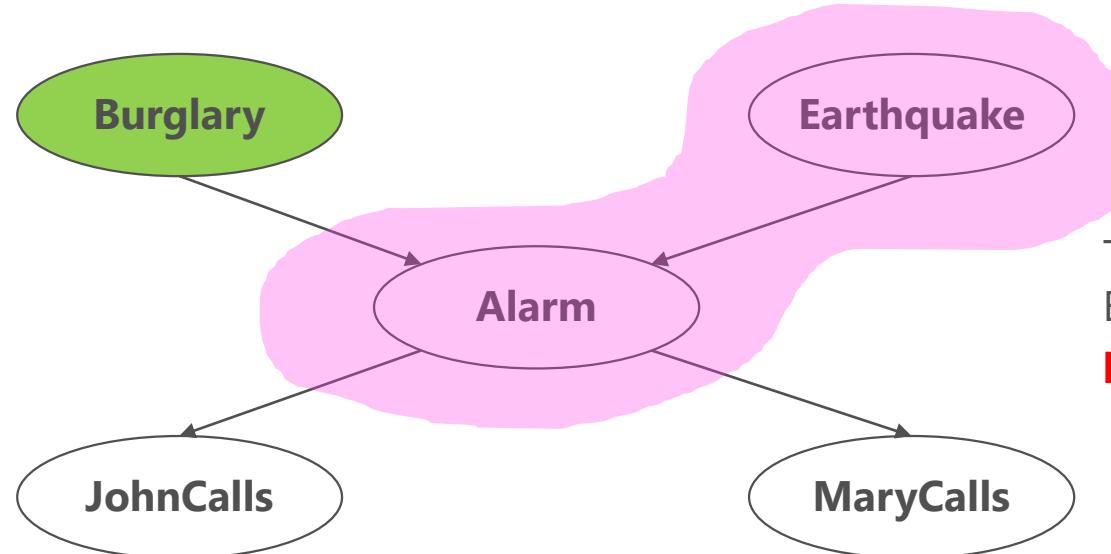
 pick an ordering of the variables

for each variable X_i in order **do**

 sample x_i from $p(x_i|x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) = p(x_i|\text{MB}(x_i))$

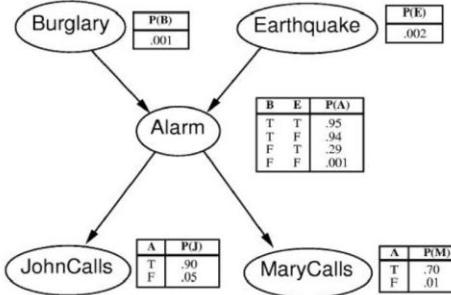
 update $X_i \leftarrow x_i$ and immediately use this new value for sampling other variables

estimate the marginal probability based on samples

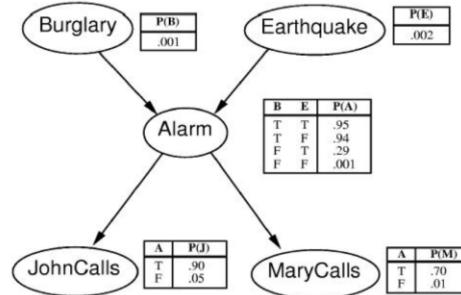


The **Markov blanket** of a node in a Bayesian Network is **the set of its parents, children and co-parents nodes**.

Inference with Markov Chain Monte Carlo (MCMC)



t	B	E	A	J	M
0	F	F	F	F	F
1	F				
2					
3					
4					



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T			
2					
3					
4					

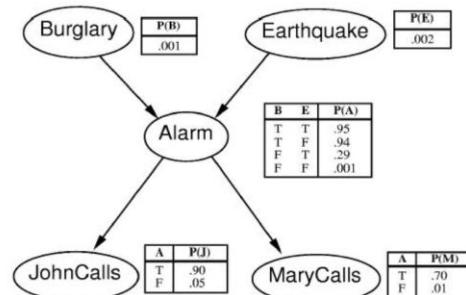
- Sampling $P(B|A,E)$ at $t = 1$: Using Bayes Rule,

$$P(B|A,E) \propto P(A|B,E)P(B)$$

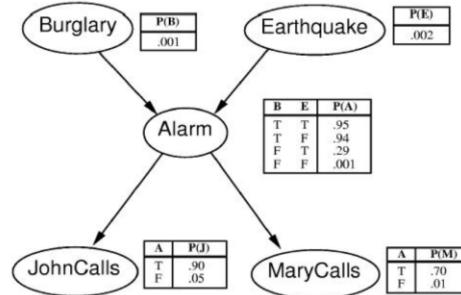
- $(A,E) = (F,F)$, so we compute the following, and sample $B = F$

$$P(B=T | A=F, E=F) \propto (0.06)(0.01) = 0.0006$$

$$P(B=F | A=F, E=F) \propto (0.999)(0.999) = 0.9980$$



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T	F	T	F
2	F	T	T	T	T
3	T	F	T	F	T
4	T	F	T	F	F



t	B	E	A	J	M
0	F	F	F	F	F
1	F	T	F	T	F
2					
3					
4					

- Now $t = 2$, and we repeat the procedure to sample new values of B,E,A,J,M

...

- And similarly for $t = 3, 4$, etc.

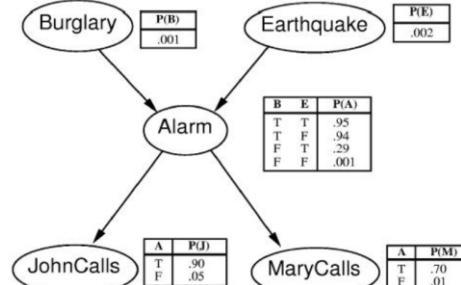
- Sampling $P(E|A,B)$: Using Bayes Rule,

$$P(E|A,B) \propto P(A|B,E)P(E)$$

- $(A,B) = (F,F)$, so we compute the following, and sample $E = T$

$$P(E=T | A=F, B=F) \propto (0.71)(0.02) = 0.0142$$

$$P(E=F | A=F, B=F) \propto (0.999)(0.998) = 0.9970$$



- Sampling $P(M|A)$: No need to apply Bayes Rule

- $A = F$, so we compute the following, and sample $M = F$

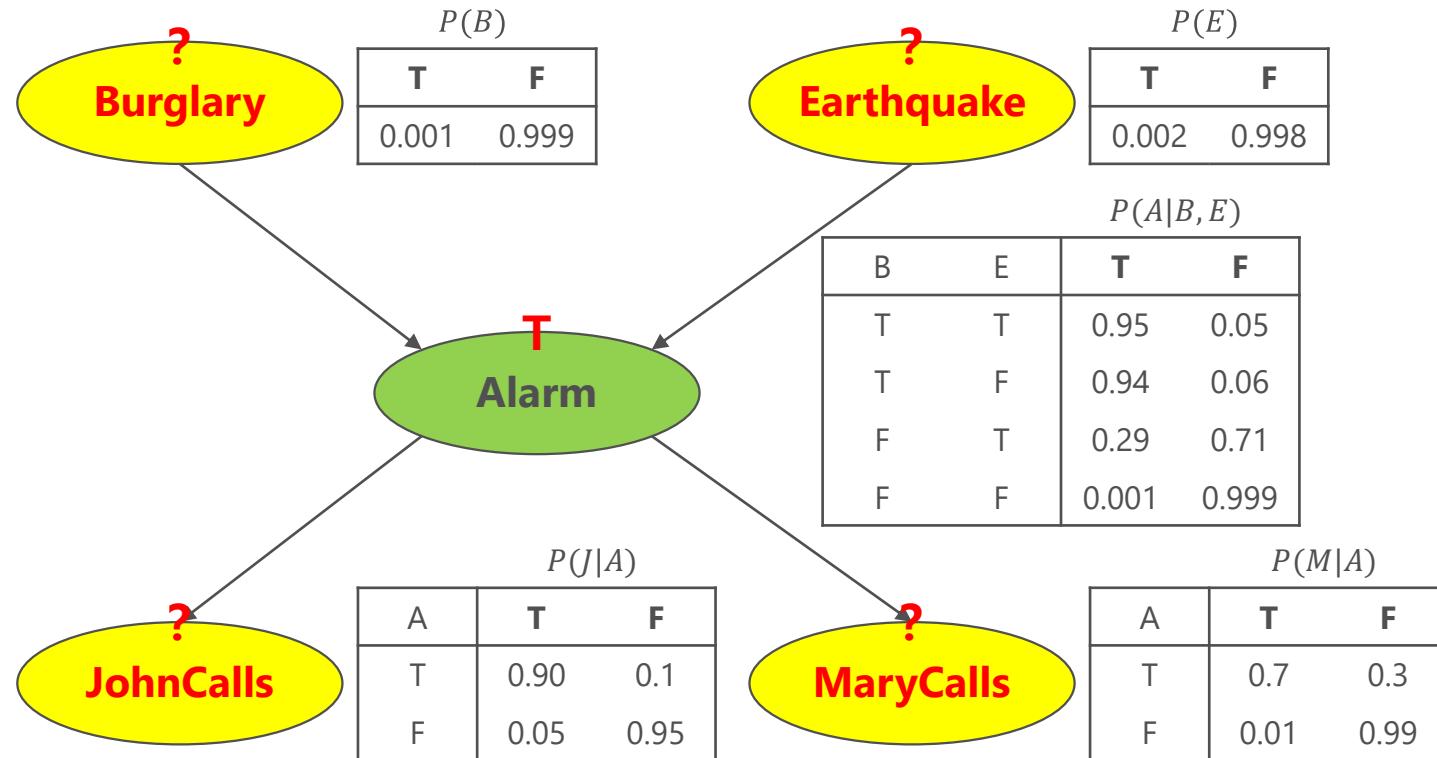
$$P(M=T | A=F) \approx 0.01$$

$$P(M=F | A=F) \approx 0.99$$

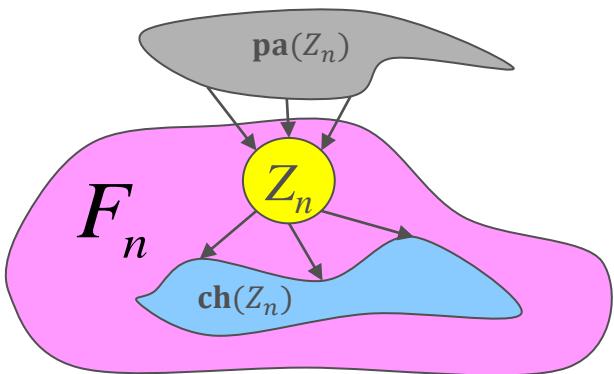
Reasoning Task (3): Maximum Probable Explanation

Find the most probable joint assignment of the unobserved variables Z given the observed variables X .

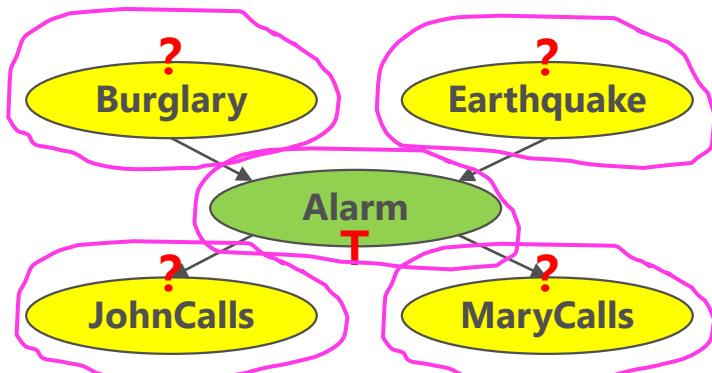
$$\text{MPA}(\mathbf{Z}|X) = \max_{\mathbf{Z}} P(\mathbf{Z}|X)$$



Inference with Bucket Elimination



$$\begin{aligned}
 \text{MPA}(\mathbf{Z}|X) &= \max_{\mathbf{Z}} P(\mathbf{Z}|X) = \max_{\mathbf{Z}} P(\mathbf{Z}, X) = \max_{\mathbf{Z}} \prod_{i=1}^n P(Z_i, X | \text{pa}(Z_i)) \\
 &= \max_{\mathbf{Z}_{n-1}} \left\{ \prod_{Z_i \in \mathbf{Z} - F_n} P(Z_i, X | \text{pa}(Z_i)) \right\} \cdot \max_{Z_n} \left\{ P(Z_n, X | \text{pa}(Z_n)) \prod_{Z_j \in \text{ch}(Z_n)} P(Z_j, X | \text{pa}(Z_j)) \right\} \\
 &= \boxed{\max_{\mathbf{Z}_{n-1}} \left\{ \prod_{Z_i \in \mathbf{Z} - F_n} P(Z_i, X | \text{pa}(Z_i)) \right\}} \cdot \boxed{h_{Z_n}}
 \end{aligned}$$

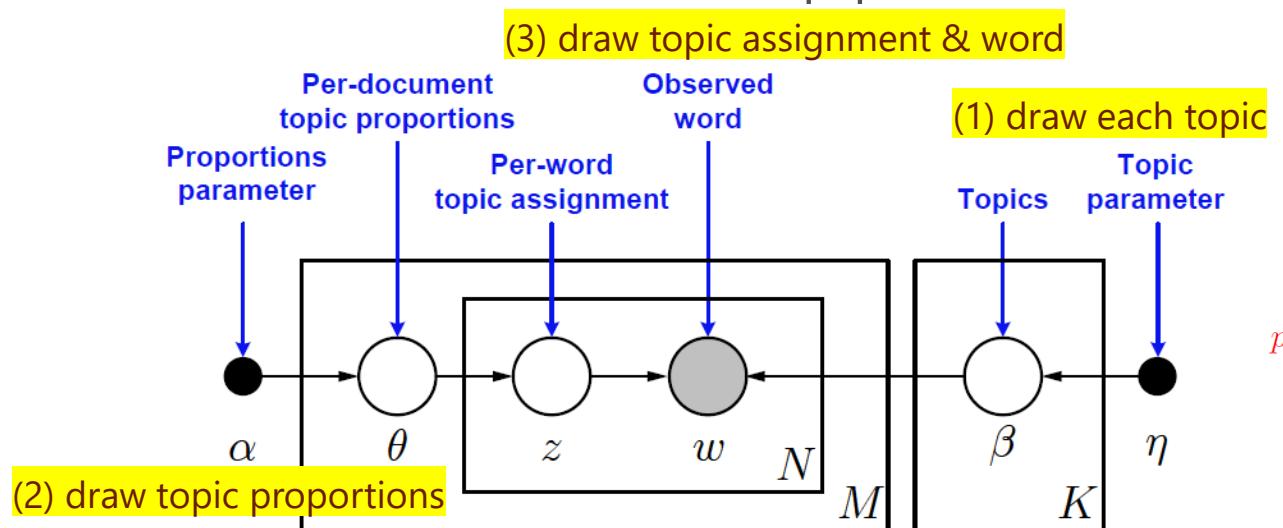


Consider the ordering: EBAMJ

- Bucket_J: $\max_J P(J|A = T) = h_J(A)$
- Bucket_M: $\max_M P(M|A = T) = h_M(A)$
- Bucket_A: $\max_A P(A = T|B, E) h_J(A) h_M(A) = h_A(B, E)$
- Bucket_B: $\max_B P(B) h_A(B, E) = h_B(E)$
- Bucket_E: $\max_E P(E) h_B(E) = \max_{E, B, M, J} P(E, B, M, J, A = T)$

Derive the most probable assignment.

Successful Application: Document Modeling

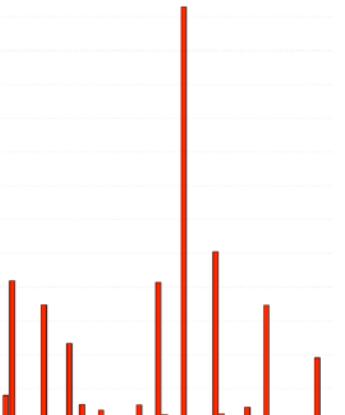


$$p(w, z, \theta, \beta) = \left(\prod_{d=1}^D p(\theta_d | \alpha) \prod_{n=1}^{N_d} p(w_{dn} | z_{dn}, \beta_{1:K}) p(z_{dn} | \theta_d) \right) \left(\prod_{k=1}^K p(\beta_k | \eta) \right)$$

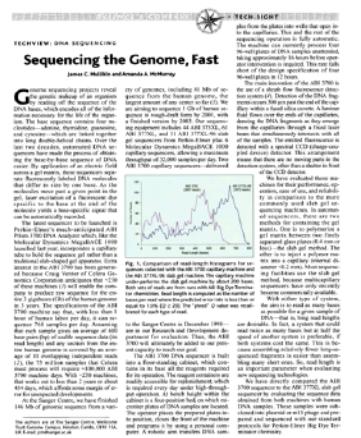
$$p(w_{dn} | z_{dn}, \beta_{1:K}) = \prod_{k=1}^K \text{Mult}(w_{dn} | \beta_k)^{z_{dn}^k} \quad p(\beta_k | \eta) = \text{Dir}(\beta_k | \eta)$$

$$p(z_{dn} | \theta_d) = \text{Mult}(z_{dn} | \theta_d) \quad p(\theta_d | \alpha) = \text{Dir}(\theta_d | \alpha)$$

Topic proportions



Original article



Most likely words from top topics

sequence	devices	data
genome	device	information
genes	materials	network
sequences	current	web
human	high	computer
gene	gate	language
dna	light	networks
sequencing	silicon	time
chromosome	material	software
regions	technology	system
analysis	electrical	words
data	fiber	algorithm
genomic	power	number
number	based	internet

Successful Application: Medical Diagnosis

nature communications

Explore our content ▾ Journal information ▾

nature > nature communications > articles > article

Article | Open Access | Published: 11 August 2020

Improving the accuracy of medical diagnosis with causal machine learning

Jonathan G. Richens Ciarán M. Lee & Saurabh Johri

Nature Communications 11, Article number: 3923 (2020) | Cite this article

15k Accesses | 886 Altmetric | Metrics

A Publisher Correction to this article was published on 16 September 2020

This article has been updated

Abstract

Machine learning promises to revolutionize clinical decision making and diagnosis. In medical diagnosis a doctor aims to explain a patient's symptoms by determining the diseases causing them. However, existing machine learning approaches to diagnosis are purely associative, identifying diseases that are strongly correlated with a patient's symptoms. We show that this inability to disentangle correlation from causation can result in sub-optimal or dangerous diagnoses. To overcome this, we reformulate diagnosis as a counterfactual inference task and derive counterfactual diagnostic algorithms. We compare our counterfactual algorithms to the standard associative algorithm and 44 doctors using a test set of clinical vignettes. While the associative algorithm achieves an accuracy placing in the top 48% of doctors in our cohort, our counterfactual algorithm places in the top 25% of doctors, achieving expert clinical accuracy. Our results show that causal reasoning is a vital missing ingredient for applying machine learning to medical diagnosis.

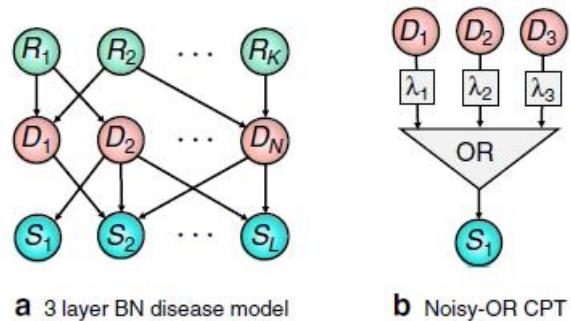
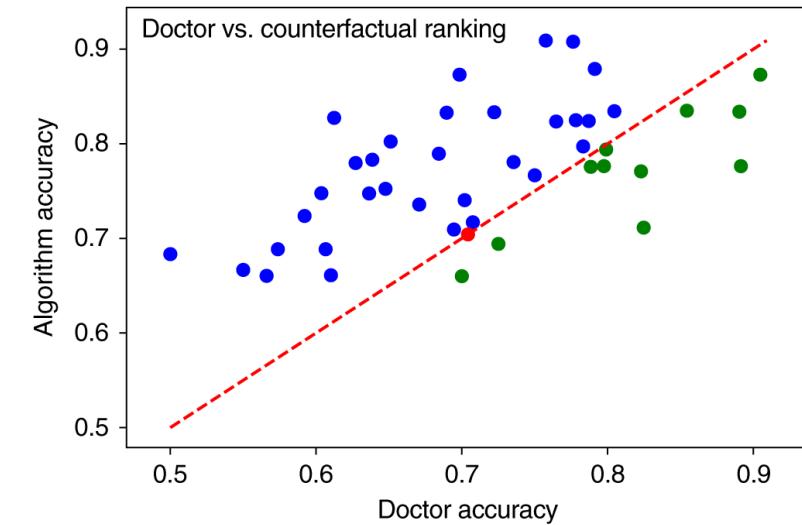
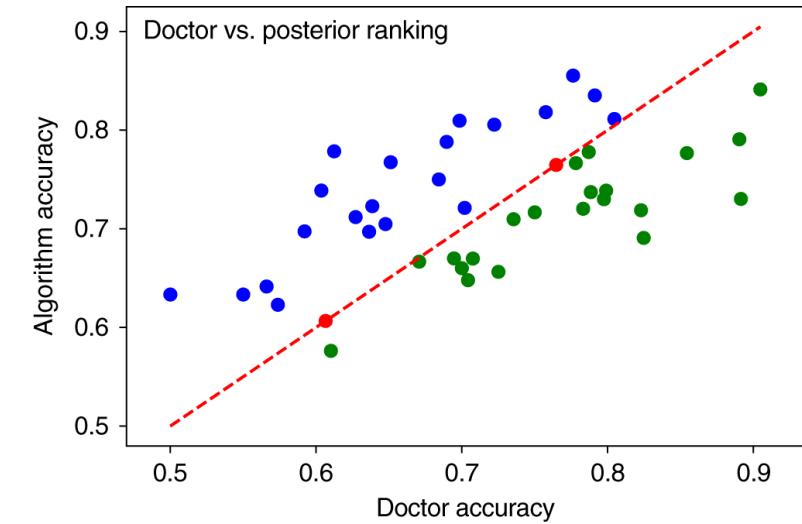


Fig. 2 Generative structure of our diagnostic Bayesian networks. **a** Three-layer Bayesian network representing risk factors R_i , diseases D_j and symptoms S_k . **b** noisy-OR CPT. S is the Boolean OR function of its parents, each with an independent probability λ_i of being ignored, removing them from the OR function.



Markov Network

- A Markov Network (a.k.a. Markov Random Field) is an undirected graph.
- A Markov Network represents a joint probability distribution over a set of random variables.
- The graph has a node for each variable.
- The graph has a potential function $\phi_k(x_{\{k\}})$ for each clique $x_{\{k\}}$ in the graph.
- Markov Network is often represented as log-linear models, with each clique potential function replaced by an exponentiated weighted sum of features.

$$P(X = x) = \frac{1}{Z} \prod_k \phi_k(x_{\{k\}}) \approx \frac{1}{Z} \exp\left(\sum_j w_j f_j(x)\right)$$

Markov Logic Network (MLN)

- An MLN L is defined as a set of pairs (F_i, w_i) that combines probability and first-order logic in a unified model.
 - F_i is a first-order logic formula, which is denoted by predicates connected by operators.
 - w_i is the weight of F_i .
- An MLN L generates a Markov network $M_{L,C}$ by grounding all variables in formulas to a finite set of constants $C = \{c_1, \dots, c_{|C|}\}$.
 - One binary node for each grounding of each predicate in L , whose value is 1 if the ground atom is true, and 0 otherwise.
 - One feature for each grounding of each formula F_i in L , whose value is 1 if the ground atom is true, and 0 otherwise.
 - One weight w_i for each grounding of each formula F_i in L .
- The probability of a world x is defined as $P(x) = \frac{1}{Z} \exp(\sum_i w_i \cdot n_i(x))$

w_i is the weight
of formula i in x

$n_i(x)$ is the number of true
groundings of formula i in x

An Example (from Matthew Richardson and Pedro Domingos)

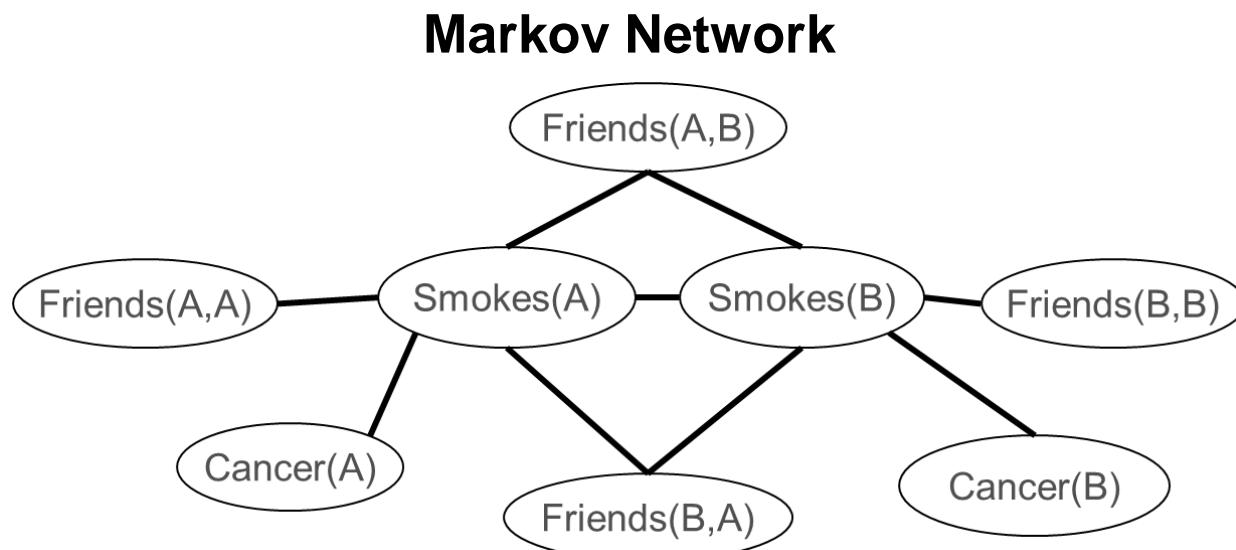
Smoking causes **cancer**.

Friends either both **smoke** or both don't smoke.



MLN

$$\begin{aligned}\forall x \text{ Smokes}(x) &\Rightarrow \text{Cancer}(x) \\ \forall x, y \text{ Friends}(x, y) &\Rightarrow (\text{Smokes}(x) \Leftrightarrow \text{Smokes}(y))\end{aligned}$$

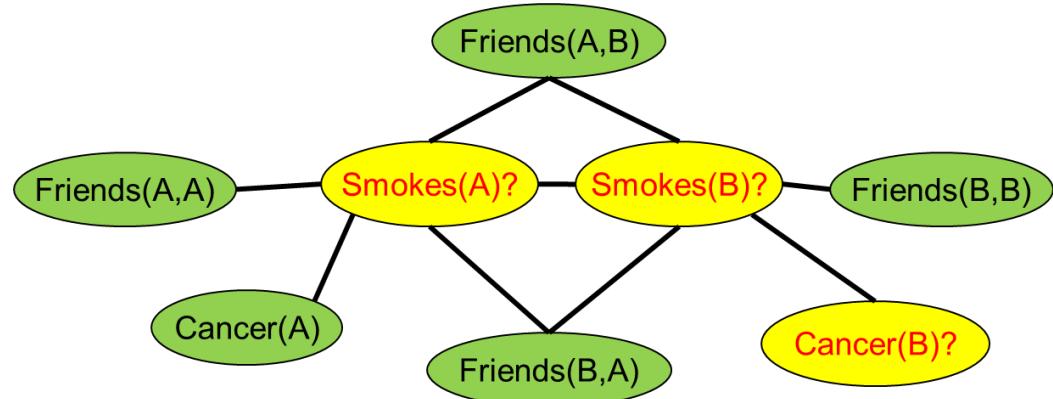


Constants

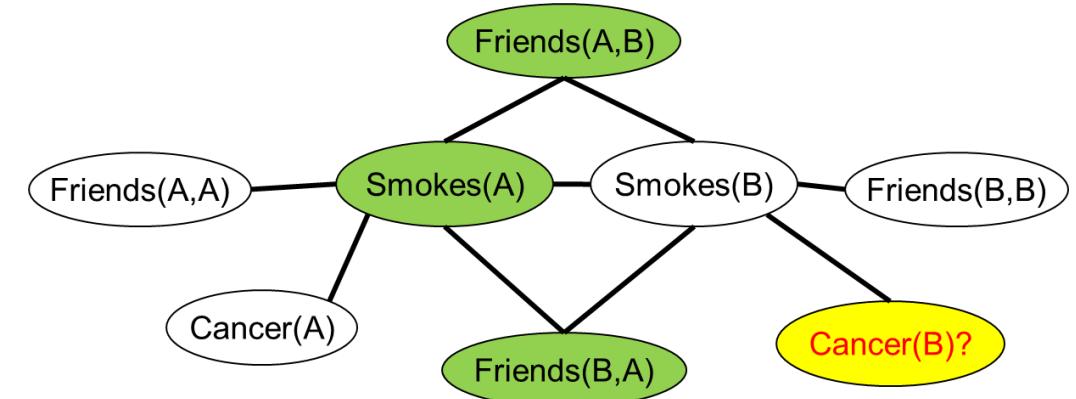
Anna (A), Bob (B)



Typical Reasoning Tasks



Maximum Probable Explanation



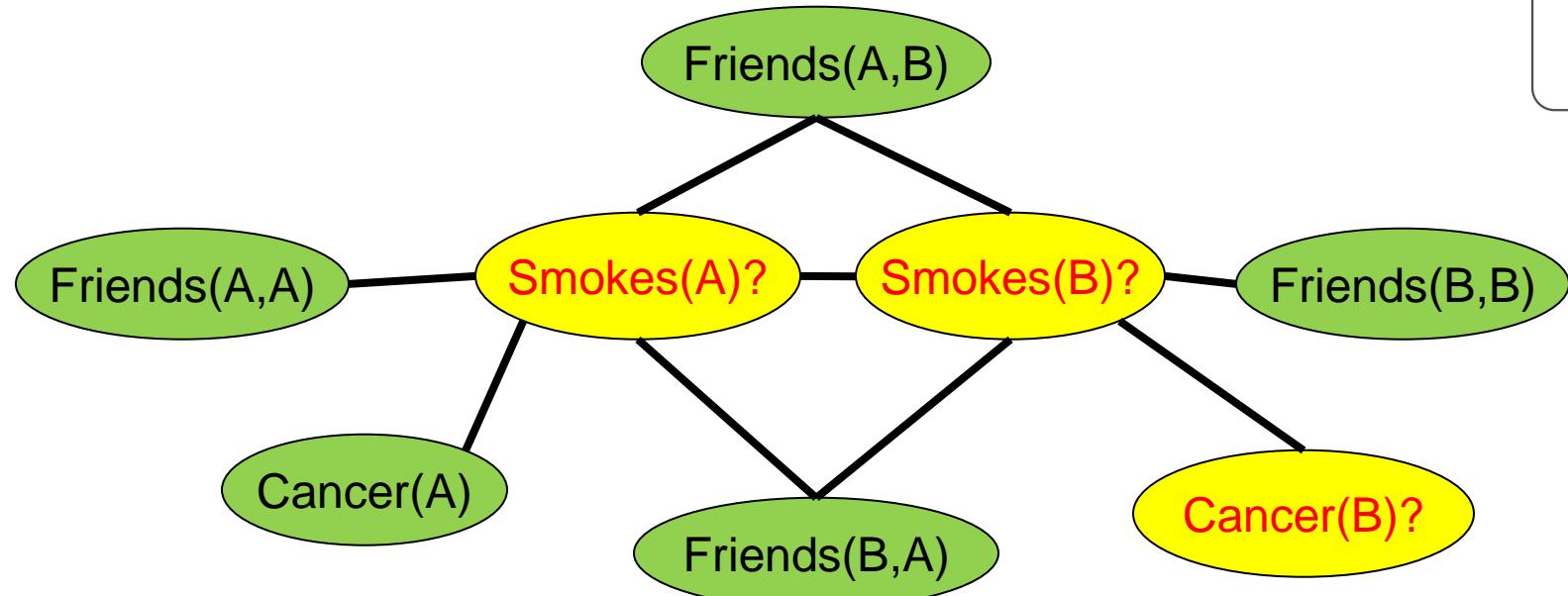
Conditional Probability

Reasoning Task (1): Maximum Probable Explanation

Find the most probable joint assignment of the unobserved variables y given the observed variables x

$$\hat{y} = \arg \max_y P(y|x) = \frac{1}{Z} \arg \max_y \sum_i w_i n_i(x, y) = \arg \max_y \sum_i w_i n_i(x, y)$$

No need to compute Z .



What are the most likely states of Smokes(A), Smokes(B) and Cancer(B)?

Inference with MaxWalkSAT

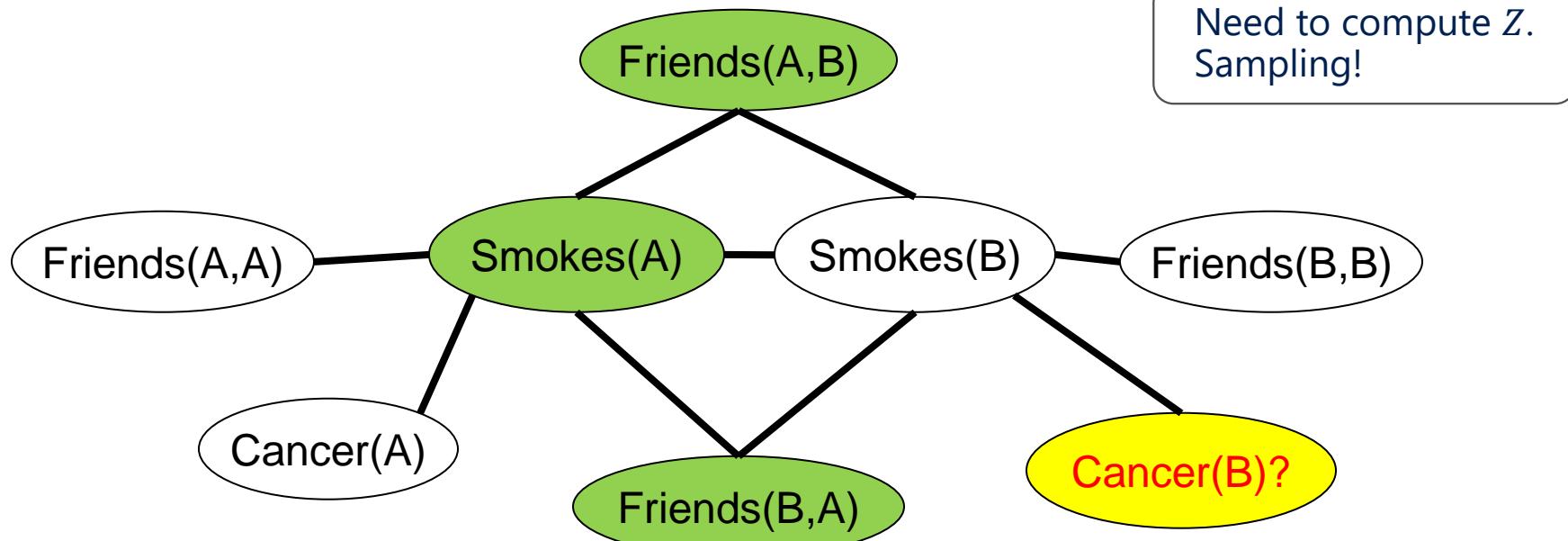
```
for  $i \leftarrow 1$  to  $max\_tries$  do
    solution = random truth assignment
    for  $j \leftarrow 1$  to  $max\_flips$  do
        if  $\sum weights(sat.\ clauses) > threshold$  then
            return solution
         $c \leftarrow$  random unsatisfied clause
        with probability  $p$ 
            flip a random variable in  $c$  //random
        else
            flip the random variable in  $c$  that maximizes  $\sum weights(sat.\ clauses)$  //greedy
    return failure, best solution found
```

Reasoning Task (2): Conditional Probability

Compute the marginal distribution of a node y given the evidence x .

$$\hat{y} = \arg \max_y P(y|x) = \frac{1}{Z} \arg \max_y \sum_i w_i n_i(x, y)$$

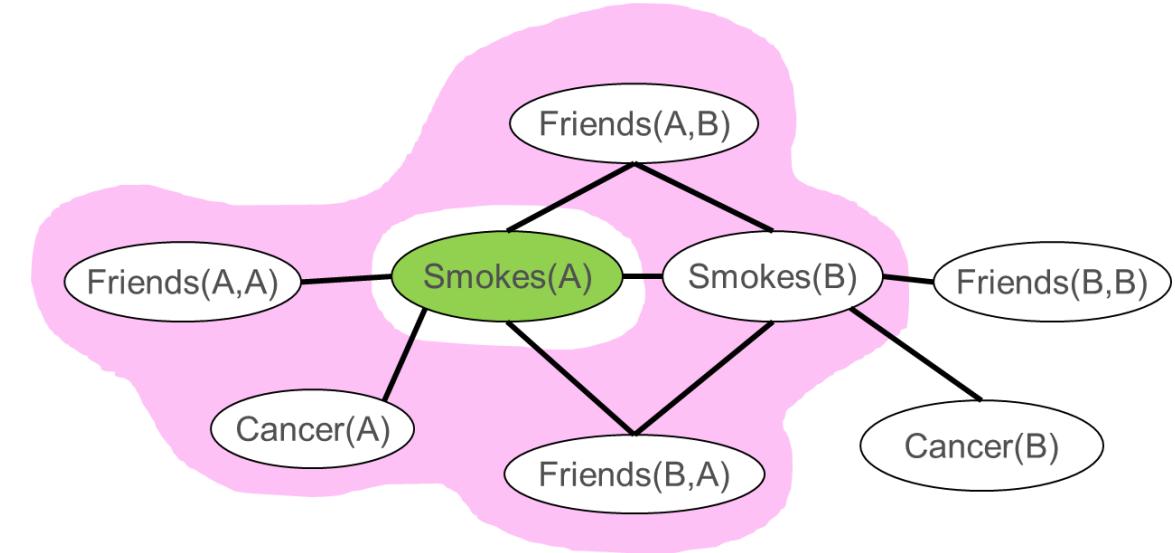
Need to compute Z .
Sampling!



What is the probability $P(\text{Cancer}(B) | \text{Smokes}(A), \text{Friends}(A,B), \text{Friends}(B,A))$?

Inference with MCMC

```
state ← random truth assignment
for  $i \leftarrow 1$  to  $num\_samples$  do
    for each variable  $x$  do
        sample  $x$  according to  $p(x|\text{MB}(x))$ 
        state ← state with new value of  $x$ 
    end for
end for
estimate the marginal probability based on samples
```



The **Markov blanket** of a node in a Markov Network is
the set of its neighboring nodes.

Successful Application: Information Extraction

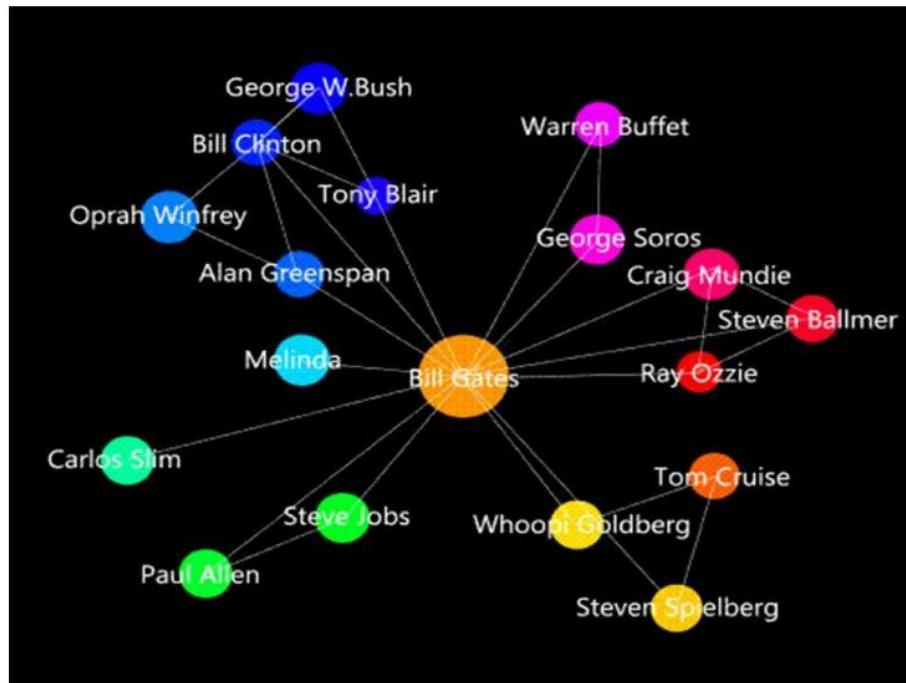


Figure 1: An entity relationship graph for the query “Bill Gates” generated by EntityCube (i.e. the English version of Renlifang).

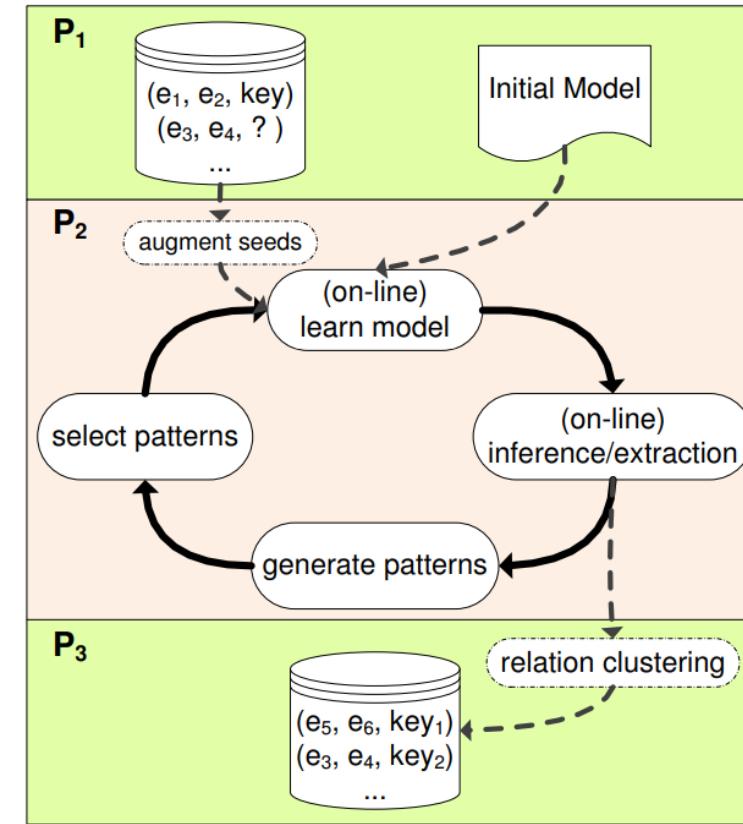


Figure 2: The framework of StatSnowball, with three parts – P_1 (input), P_2 (statistical extraction model), and P_3 (output).

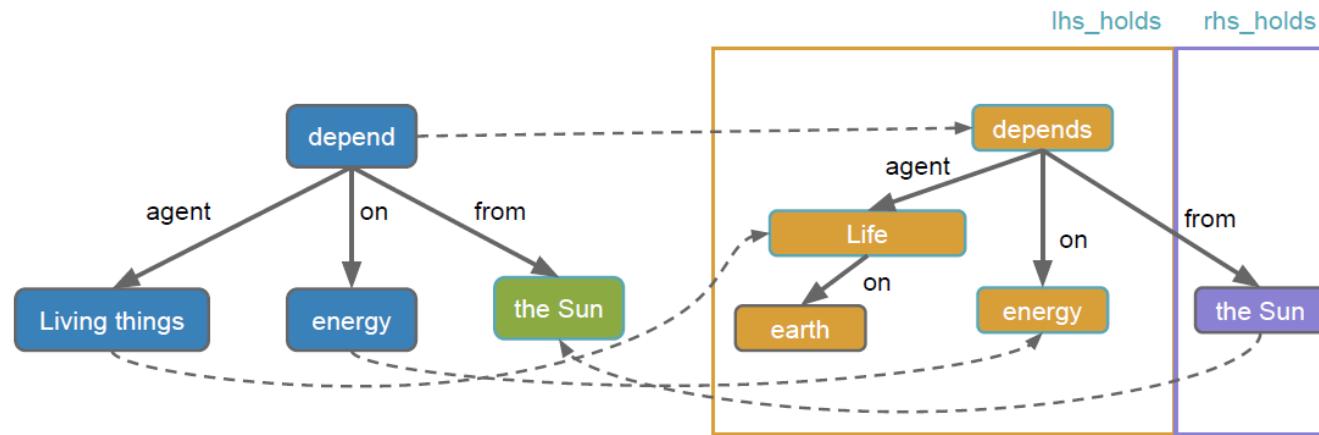
Successful Application: Semantic Parsing & QA

Microsoft buys Skype



BUY (n₁) \wedge MICROSOFT (n₂) \wedge SKYPE (n₃)
 \wedge BUYER (n₁, n₂) \wedge BOUGHT (n₁, n₃)

Hoifung Poon, Pedro Domingos. Unsupervised Semantic Parsing. EMNLP, 2009.



Tushar Khot, Niranjan Balasubramanian, Eric Gribkoff, Ashish Sabharwal, Peter Clark, Oren Etzioni. Exploring Markov Logic Networks for Question Answering. EMNLP, 2015.

Successful Application: MLN Toolkits

Alchemy: Open Source AI

Welcome to the Alchemy system! Alchemy is a software package providing a series of algorithms for statistical relational learning and probabilistic logic inference, based on the Markov logic representation. Alchemy allows you to easily develop a wide range of AI applications, including:

- Collective classification
- Link prediction
- Entity resolution
- Social network modeling
- Information extraction

Choose a version of Alchemy:

[Alchemy Lite](#)

Alchemy Lite is a software package for inference in Tractable Markov Logic (TML), the first tractable first-order probabilistic logic. Alchemy Lite allows for fast, exact inference for models formulated in TML. Alchemy Lite can be used in batch or interactive mode.

[Alchemy 2.0](#)

Alchemy 2.0 is a new version of Alchemy based on probabilistic theorem proving as the inference engine. It also includes the most widely used algorithms in the original Alchemy system.

[Alchemy 1.0](#)

The original version of Alchemy; no longer under active development.

<https://alchemy.cs.washington.edu/>

HAZY Project Tuffy

HOME PEOPLE PROJECTS PUBLICATIONS [✉](#) [Follow](#)

TUFFY DOWNLOAD DEMO DOCUMENTATION CHANGELOG PEOPLE CONTACT

Hazy > Tuffy

Meet Tuffy (Ver 0.4 released July 13, 2014!)

Tuffy is now released under Apache Version 2!

"We balance probabilities and choose the most likely. It is the scientific use of the imagination." -- The Master

Tuffy is an open-source Markov Logic Network inference engine, and part of Felix.

Check out our new demos built with Tuffy/Felix!

Markov Logic Networks (MLNs) is a powerful framework that combines statistical and logical reasoning; they have been applied to many data intensive problems including information extraction, entity resolution, text mining, and natural language processing. Based on principled data management techniques, Tuffy is an MLN inference engine that achieves scalability and orders of magnitude speedup compared to prior art implementations. It is written in Java and relies on PostgreSQL. For a brief introduction to MLNs and the technical details of Tuffy, please see our upcoming paper [\[1\]](#) or the technical report [\[2\]](#).

When designing and developing the user interface of Tuffy, we used Alchemy as a reference system. Thus, users who have experiences with Alchemy should be able to pick up Tuffy easily.

The current version (0.3) of Tuffy is capable of the following MLN tasks:

- **MRF partitioning**, a technique that can result in dramatically improved result quality (see our paper [\[1\]](#));
- **MAP inference**, where we want to find out the most likely possible world;
- **Marginal inference**, where we want to estimate marginal probabilities;
- **Weight learning**, where we want to learn the weights of MLN rules given training data.

<http://i.stanford.edu/hazy/tuffy/>

Benefits & Limitations

▪ Benefits

- Logic handles the complexity
- Probability handles the uncertainty
- Doesn't require to be trained on huge amounts of data

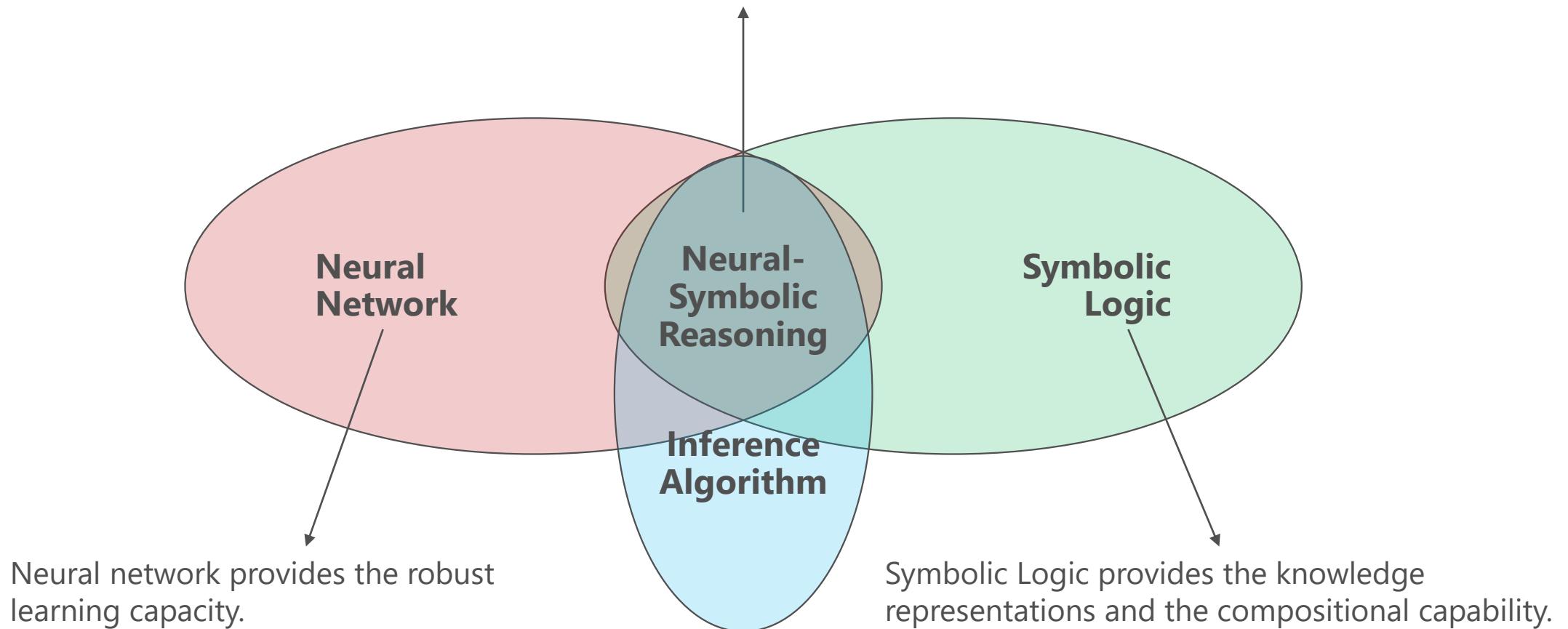
▪ Challenges

- Require hand-crafted logic formulas and are difficult to scale to big data sets
- Inference remains challenging due to the complicated relational structures
- Cannot make use of similarities between predicates or constants in training data

Neural-Symbolic Reasoning

Neural-Symbolic Reasoning

A neural-symbolic reasoning system integrates neural networks with symbolic logic.



4 Typical Neural-Symbolic Reasoning Approaches

- **Knowledge Graph Reasoning**

How to reason based on symbolic knowledge graphs using latest neural network techniques.

- **Neural Semantic Parsing**

How to convert natural language queries to programs based on symbolic knowledge graphs.

- **Neural Module Networks**

How to convert natural language queries to programs based on pre-defined differentiable modules.

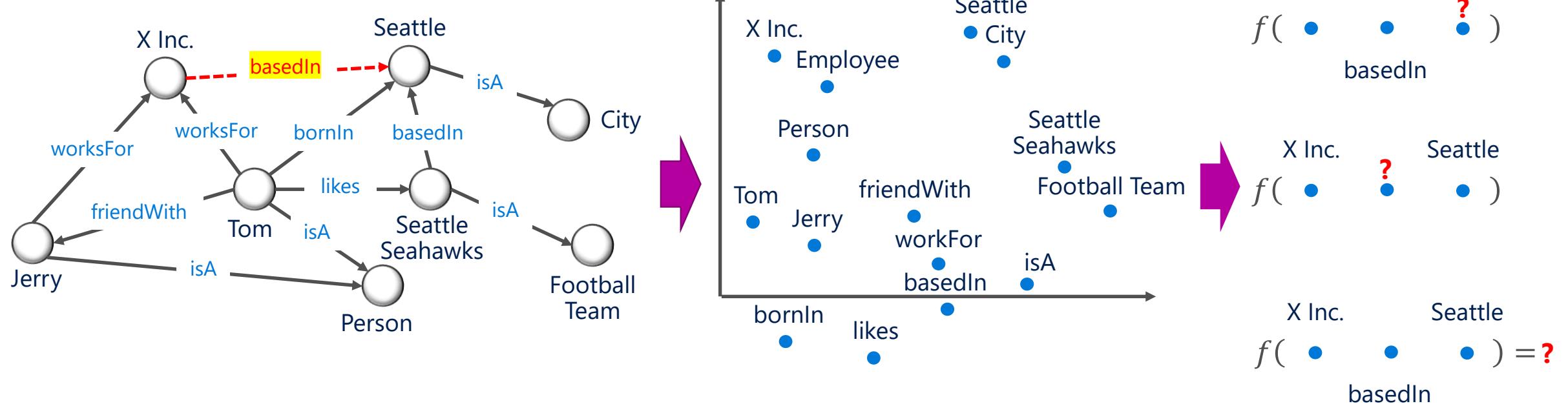
- **Symbolic Knowledge as Constraints**

How to learn better representations based on symbolic knowledge as constraints for various reasoning tasks .

4 Typical Neural-Symbolic Reasoning Approaches

- Knowledge Graph Reasoning ←
- Neural Semantic Parsing
- Neural Module Networks
- Symbolic Knowledge as Constraints

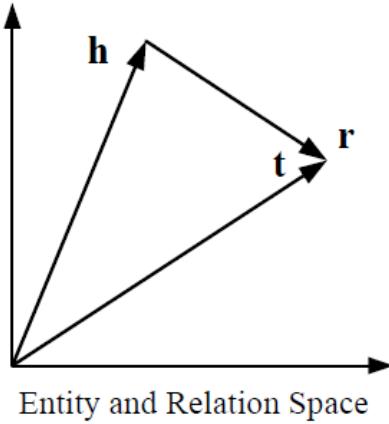
Knowledge Graph Reasoning



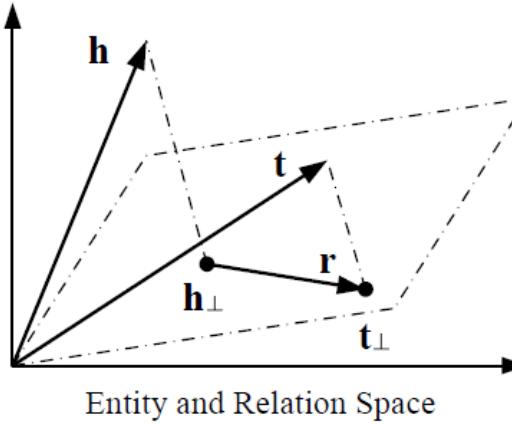
- knowledge graph completion
- natural language understanding
- question answering
- dialogue system
- recommendation system
- ...

(1): Translational Distance Model for Knowledge Graph Reasoning

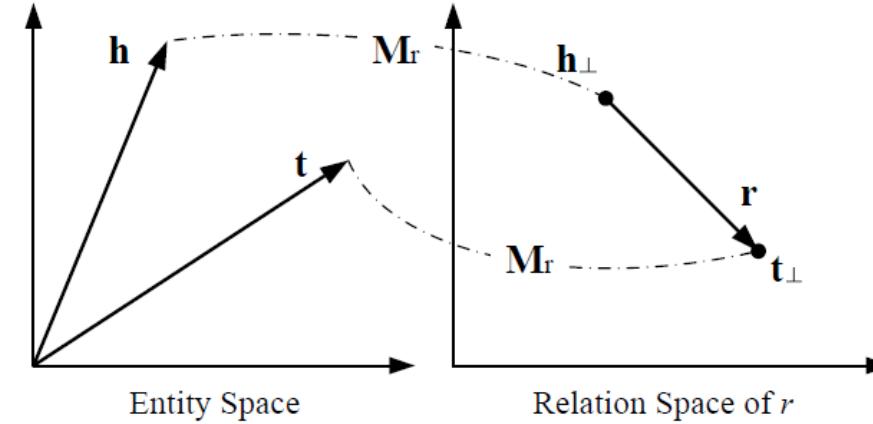
Measure the plausibility of a fact as the distance between the two entities.



(a) TransE.



(b) TransH.



(c) TransR.

$$f_r(h, t) = -\|h + r - t\|_{1/2}$$

$$f_r(h, t) = -\|h_\perp + r - t_\perp\|_2^2$$

$$h_\perp = h - w_r^\top h w_r$$

$$t_\perp = t - w_r^\top t w_r$$

$$f_r(h, t) = -\|h_\perp + r - t_\perp\|_2^2$$

$$h_\perp = M_r h$$

$$t_\perp = M_r t$$

Antoine Bordes, Nicolas Usunier, Alberto Garcia-Dur'an. Translating Embeddings for Modeling Multi-relational Data. NeurIPS, 2013.

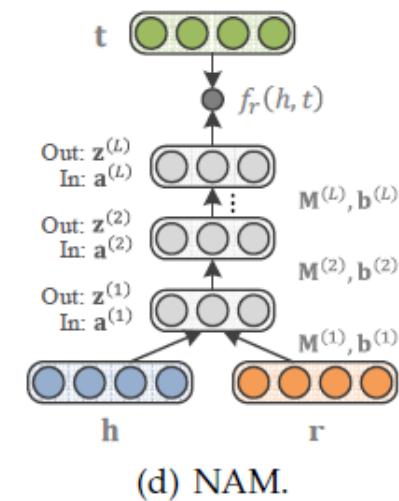
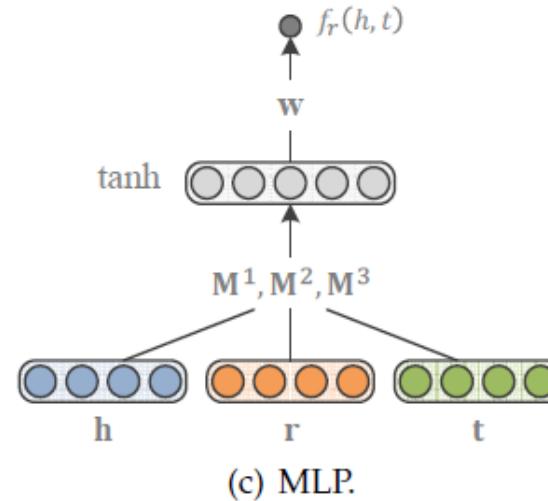
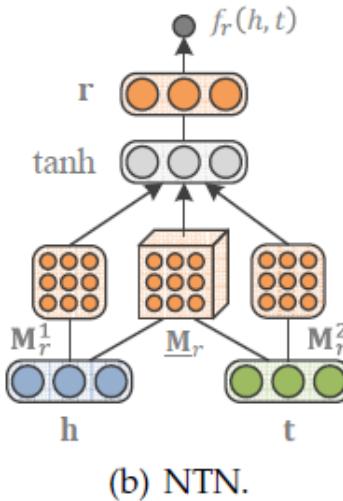
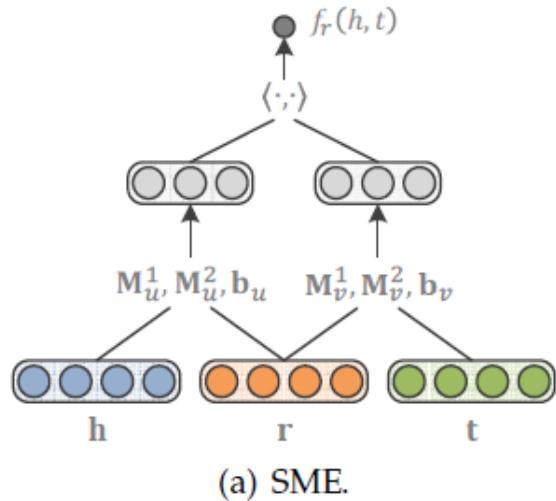
Zhen Wang, Jianwen Zhang, Jianlin Feng, Zheng Chen. Knowledge Graph Embedding by Translating on Hyperplanes. AAAI, 2014.

Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, Xuan Zhu. Learning Entity and Relation Embeddings for Knowledge Graph Completion. AAAI, 2015.

Quan Wang, Zhendong Mao, Bin Wang, Li Guo. Knowledge Graph Embedding: A Survey of Approaches and Applications. IEEE-TKDE, 2017.

(2): Semantic Matching Model for Knowledge Graph Reasoning

Measure plausibility of facts by matching latent semantics of entities and relations embodied in vector space representations.



$$f_r(h, t) = g_u(h, r)^\top g_v(t, r).$$

$$f_r(h, t) = \mathbf{r}^\top \tanh(\mathbf{h}^\top \underline{\mathbf{M}}_r \mathbf{t} + \mathbf{M}_r^1 \mathbf{h} + \mathbf{M}_r^2 \mathbf{t} + \mathbf{b}_r)$$

$$f_r(h, t) = \mathbf{w}^\top \tanh(\mathbf{M}^1 \mathbf{h} + \mathbf{M}^2 \mathbf{r} + \mathbf{M}^3 \mathbf{t})$$

$$f_r(h, t) = \mathbf{t}^\top \mathbf{z}^{(L)}$$

$$g_u(\mathbf{h}, \mathbf{r}) = \mathbf{M}_u^1 \mathbf{h} + \mathbf{M}_u^2 \mathbf{r} + \mathbf{b}_u$$

$$\mathbf{a}^{(\ell)} = \mathbf{M}^{(\ell)} \mathbf{z}^{(\ell-1)} + \mathbf{b}^{(\ell)}, \quad \ell = 1, \dots, L,$$

$$g_v(\mathbf{t}, \mathbf{r}) = \mathbf{M}_v^1 \mathbf{t} + \mathbf{M}_v^2 \mathbf{r} + \mathbf{b}_v$$

$$\mathbf{z}^{(\ell)} = \text{ReLU}(\mathbf{a}^{(\ell)}), \quad \ell = 1, \dots, L$$

$$\mathbf{z}^{(0)} = [\mathbf{h}; \mathbf{r}] \in \mathbb{R}^{2d}$$

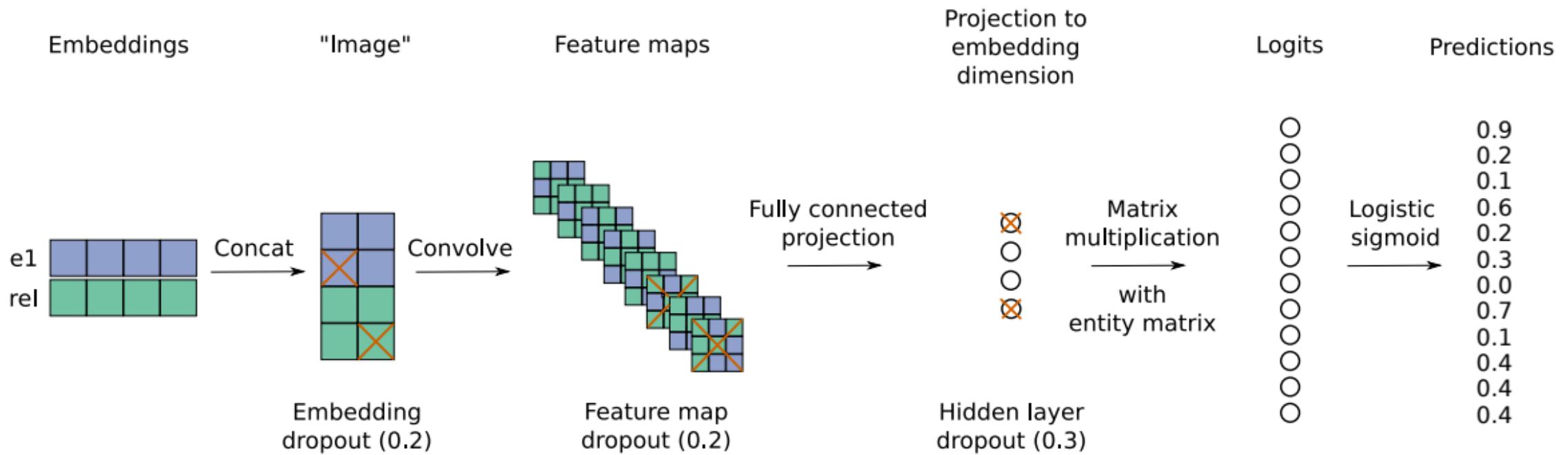
Xavier Glorot, Antoine Bordes, Jason Weston, Yoshua Bengio. A Semantic Matching Energy Function for Learning with Multi-relational Data. arXiv, 2013.

Richard Socher, Danqi Chen, Christopher D. Manning, Andrew Y. Ng. Reasoning With Neural Tensor Networks for Knowledge Base Completion. NeurIPS, 2013.

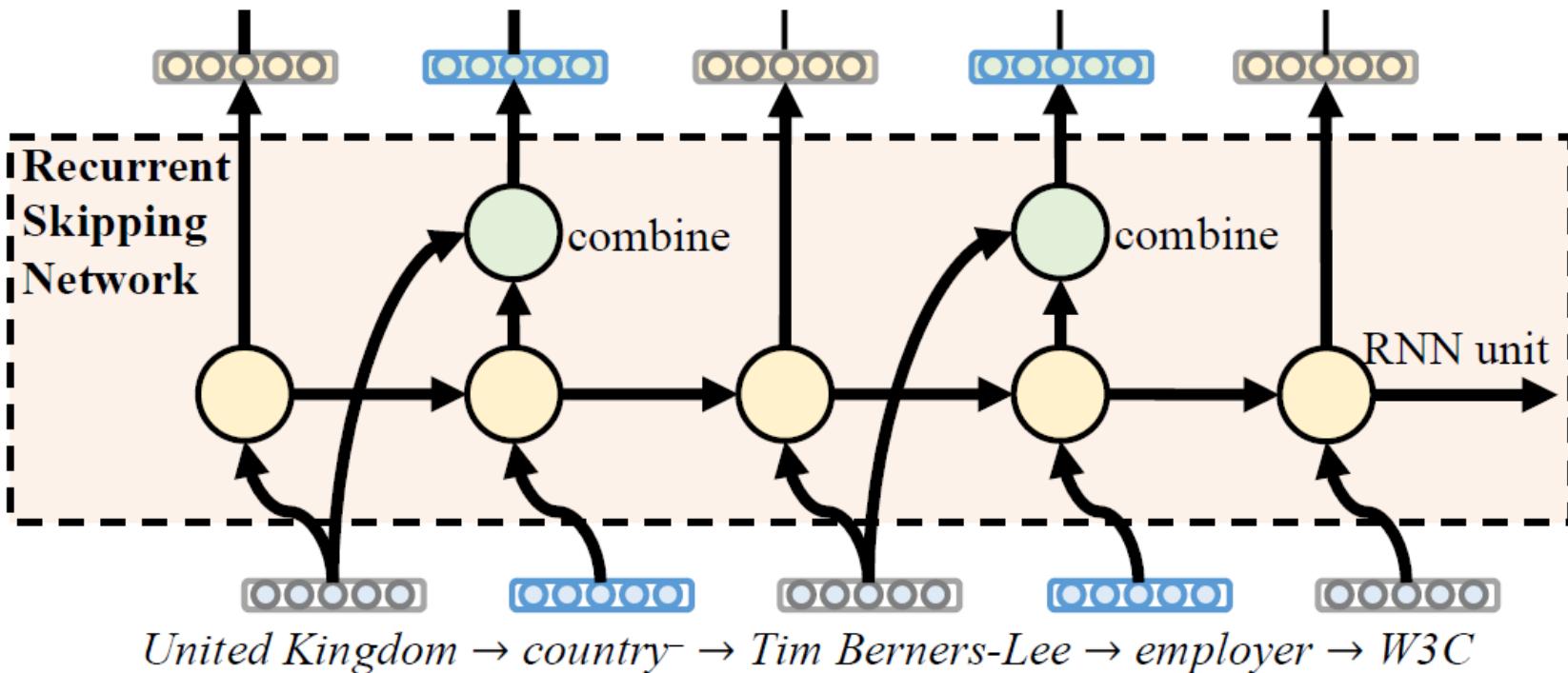
Xin Luna Dong, Evgeniy Gabrilovich, Jeremy Heitz, Wilko Horn, Ni Lao, Kevin Murphy, Thomas Strohmann, Shaohua Sun, Wei Zhang. Knowledge Vault: A Web-Scale Approach to Probabilistic Knowledge Fusion. KDD, 2014.

Quan Liu, Hui Jiang, Andrew Evdokimov, Zhen-Hua Ling, Xiaodan Zhu, Si Wei, Yu Hu. Probabilistic Reasoning via Deep Learning: Neural Association Models. arXiv, 2016.

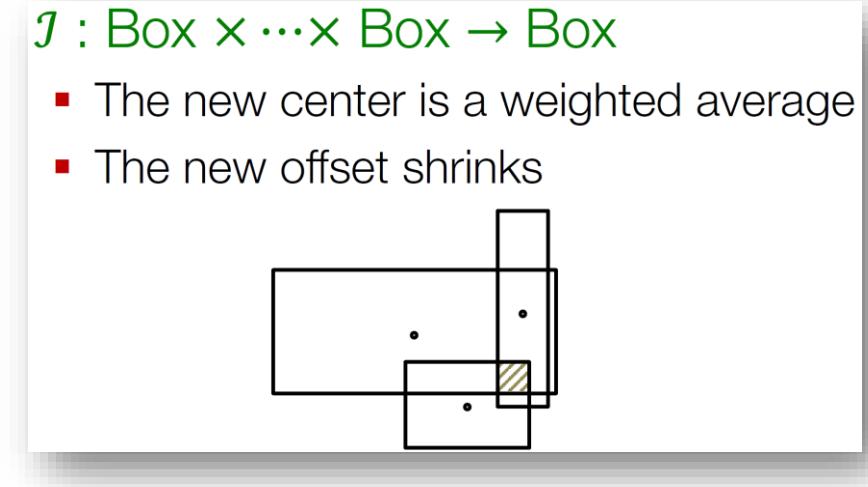
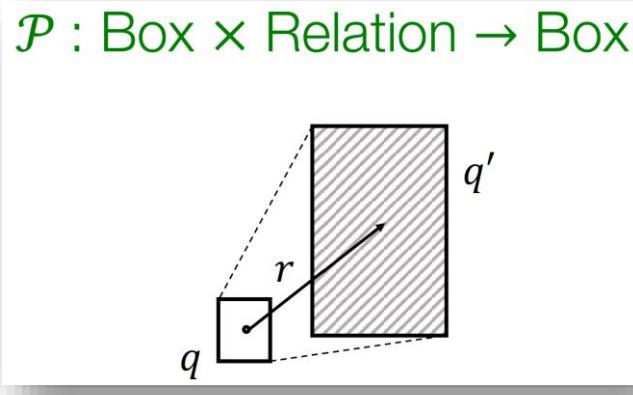
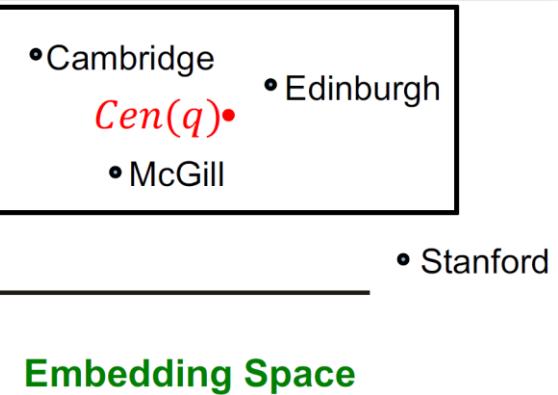
(3): Convolutional Neural Networks for Knowledge Graph Reasoning



(4): Recurrent Neural Networks for Knowledge Graph Reasoning



(5): Query2Box for Knowledge Graph Reasoning



$$\mathbf{q} = (\text{Cen}(\mathbf{q}), \text{Off}(\mathbf{q}))$$

$$\begin{aligned}\mathbf{q}' &= \mathbf{q} + \mathbf{r} \\ &= (\text{Cen}(\mathbf{q}) + \text{Cen}(\mathbf{r}), \text{Off}(\mathbf{q}) + \text{Off}(\mathbf{r}))\end{aligned}$$

$$\mathbf{p}_{\text{inter}} = (\text{Cen}(\mathbf{p}_{\text{inter}}), \text{Off}(\mathbf{p}_{\text{inter}}))$$

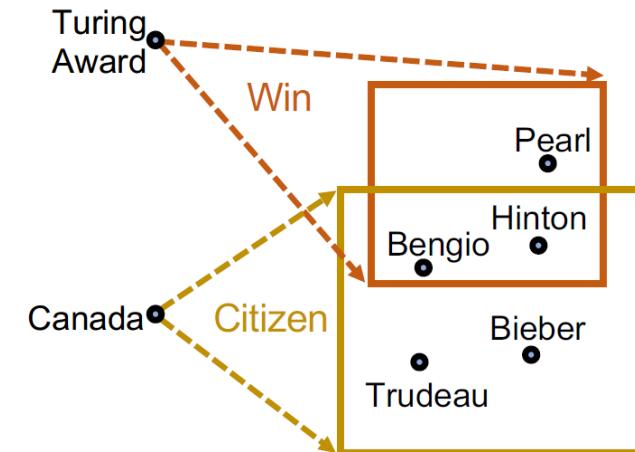
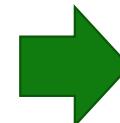
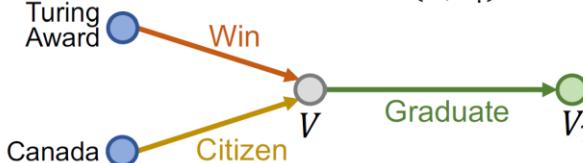
$$\text{Cen}(\mathbf{p}_{\text{inter}}) = \sum_i \mathbf{a}_i \odot \text{Cen}(\mathbf{p}_i), \quad \mathbf{a}_i = \frac{\exp(\text{MLP}(\mathbf{p}_i))}{\sum_j \exp(\text{MLP}(\mathbf{p}_j)),}$$

$$\text{Off}(\mathbf{p}_{\text{inter}}) = \text{Min}(\{\text{Off}(\mathbf{p}_1), \dots, \text{Off}(\mathbf{p}_n)\}) \odot \sigma(\text{DeepSets}(\{\mathbf{p}_1, \dots, \mathbf{p}_n\}))$$

(5): Query2Box for Knowledge Graph Reasoning

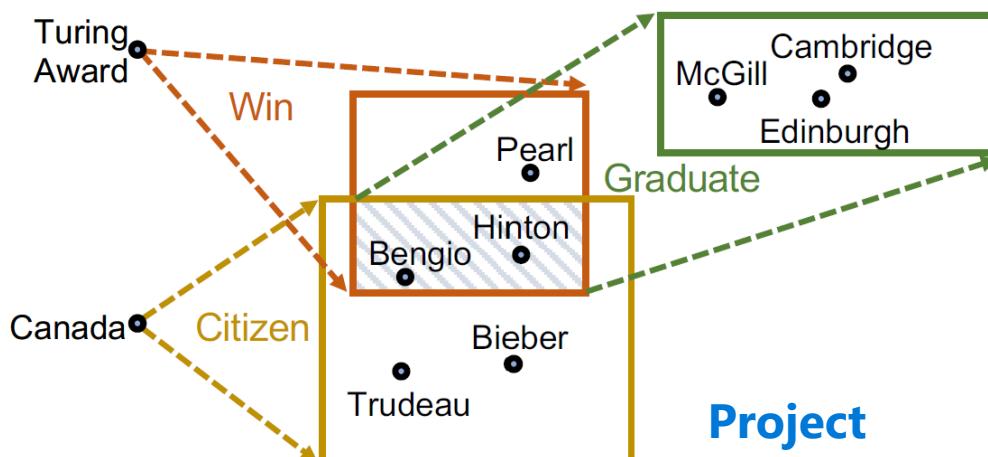
“Where did all Canadian citizens with Turing Award graduate?”

$$q = V_? . \exists V : Win(TuringAward, V) \wedge Citizen(Canada, V) \wedge Graduate(V, V_?)$$

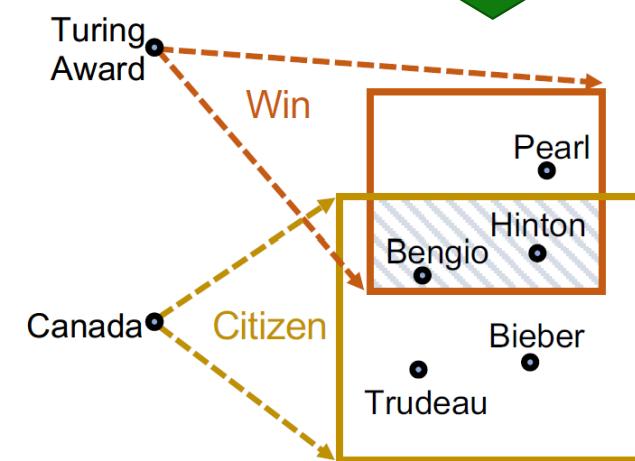


McGill Cambridge
Edinburgh

Project



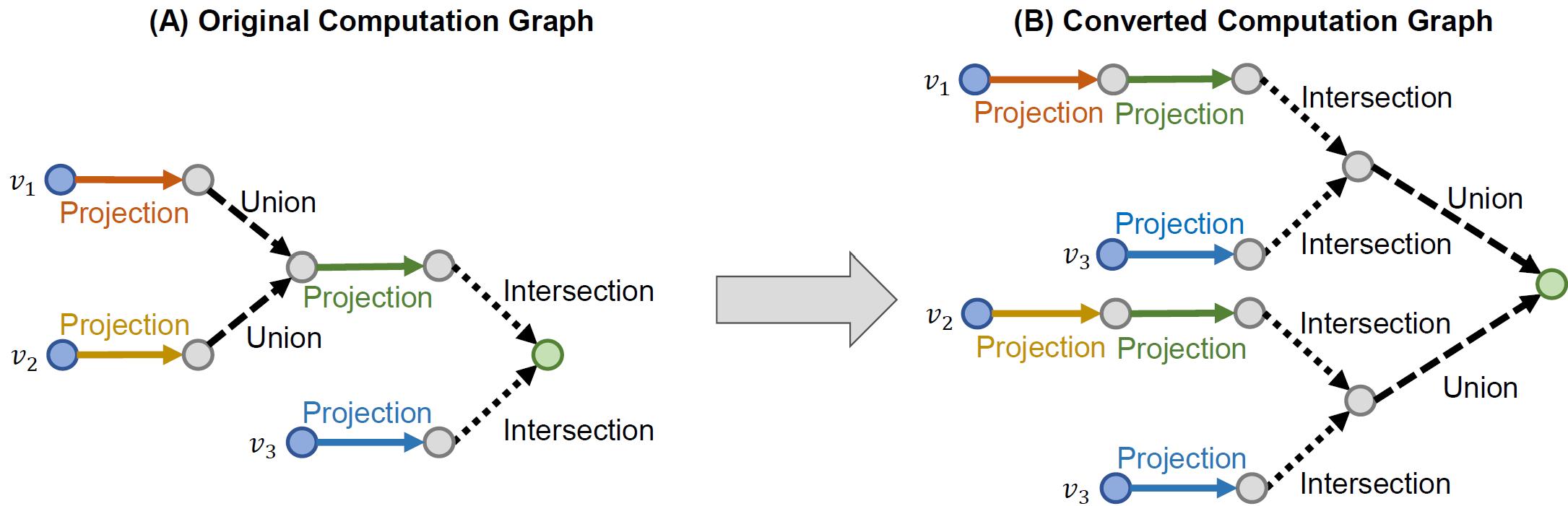
Project



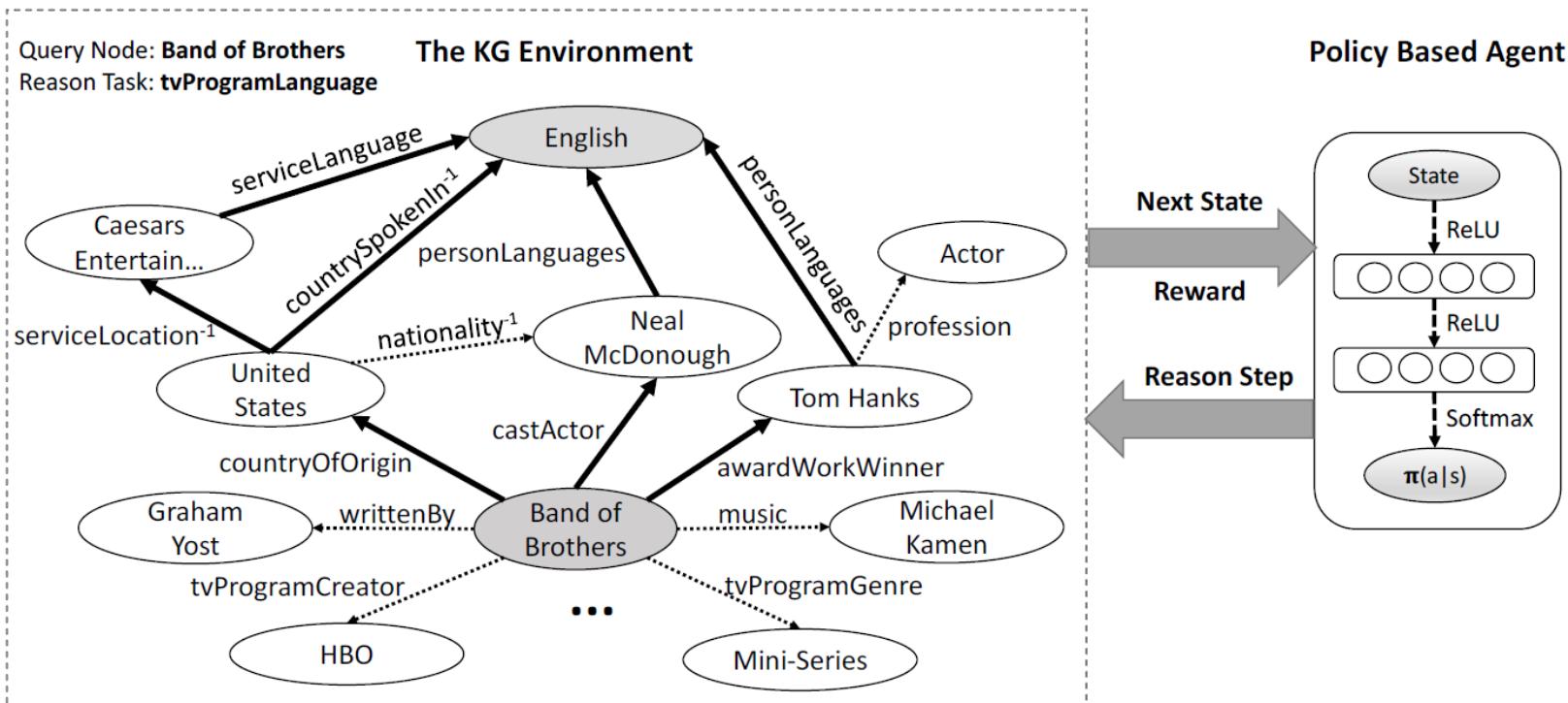
McGill Cambridge
Edinburgh

Intersection

(5): Query2Box for Knowledge Graph Reasoning



(6): Reinforcement Learning for Knowledge Graph Reasoning



- **Actions:** all relations in the KG.
- **States:** the state vector at step t is $s_t = (e_t, e_{target} - e_t)$, where e_t is the embedding of the current entity node and e_{target} is the embedding of the target entity node.

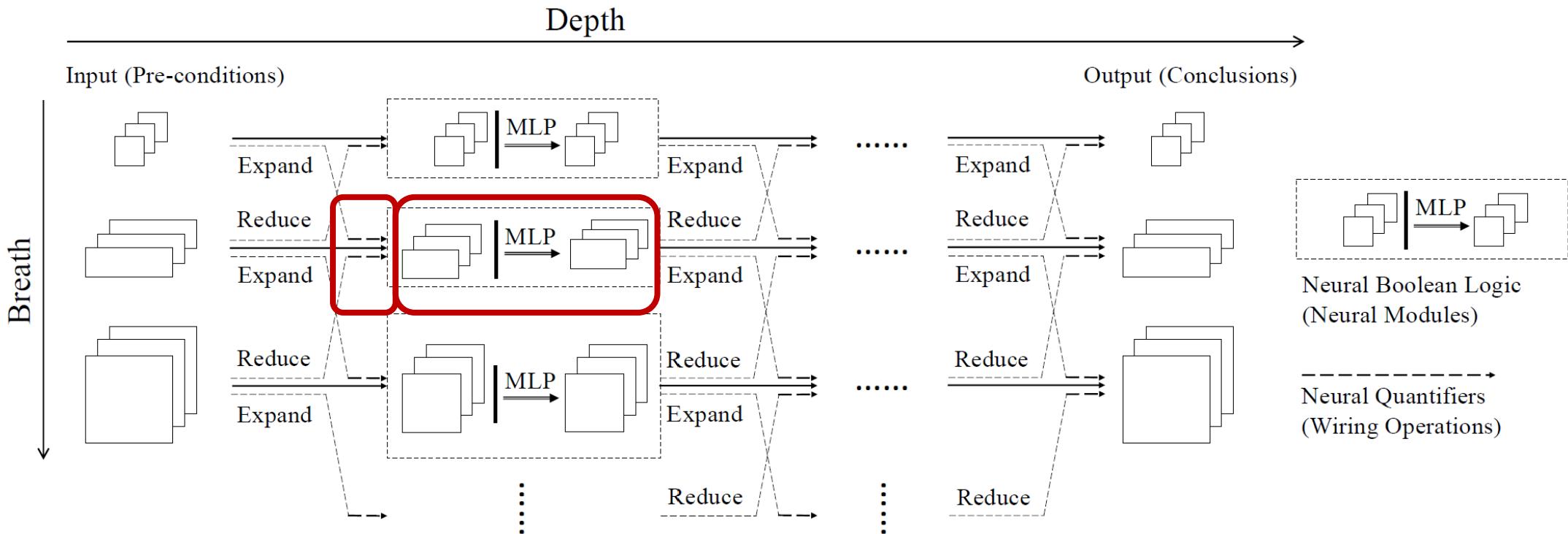
- **Rewards:**

$$r_{\text{GLOBAL}} = \begin{cases} +1, & \text{if the path reaches } e_{target} \\ -1, & \text{otherwise} \end{cases}$$

$$r_{\text{EFFICIENCY}} = \frac{1}{\text{length}(p)}$$

$$r_{\text{DIVERSITY}} = -\frac{1}{|F|} \sum_{i=1}^{|F|} \cos(p, p_i)$$

(7): Neural Forward Chaining for Knowledge Graph Reasoning



Expansion: $\forall x_{r+1} q(x_1, x_2, \dots, x_r, x_{r+1}) \leftarrow p(x_1, x_2, \dots, x_r)$

Reduction: $q(x_1, x_2, \dots, x_r) \leftarrow \forall x_{r+1} p(x_1, x_2, \dots, x_r, x_{r+1})$

Boolean logic: $O_i^{(r)} = \sigma \left(\text{MLP} \left(\text{Permute} \left(I_i^{(r)} \right); \theta_i^{(r)} \right) \right)$

(8): Neural Backward Chaining for Knowledge Graph Reasoning

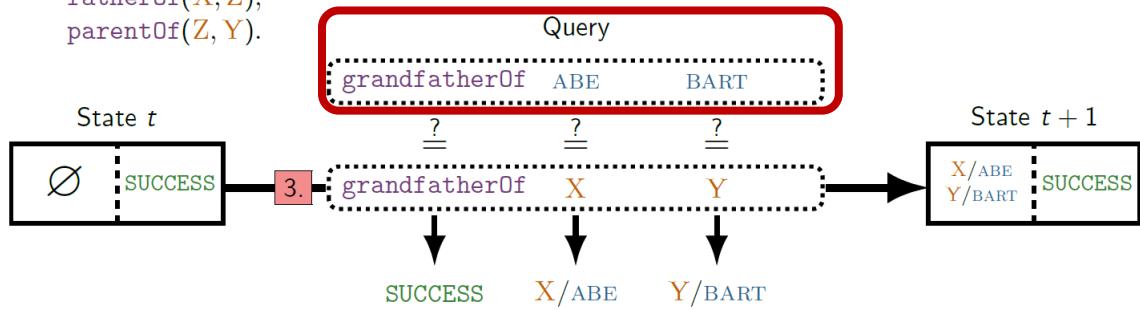
Example Knowledge Base:

```

1. fatherOf(ABE, HOMER).
2. parentOf(HOMER, BART).
3. grandfatherOf(X, Y) :-  

    fatherOf(X, Z),  

    parentOf(Z, Y).
  
```



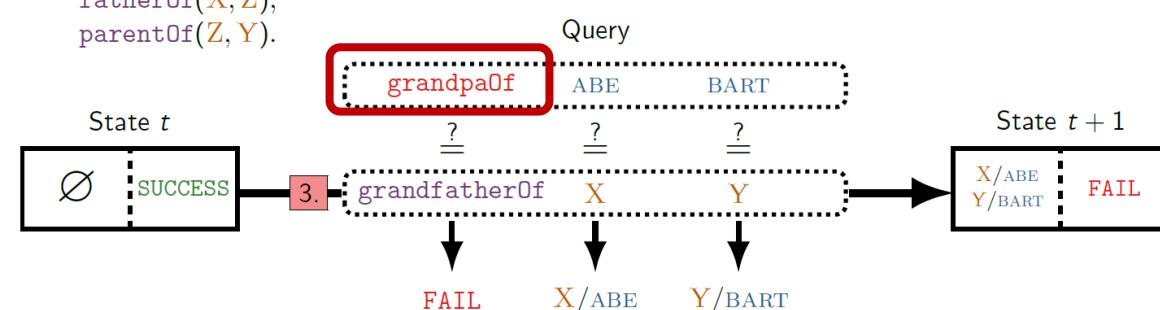
Example Knowledge Base:

```

1. fatherOf(ABE, HOMER).
2. parentOf(HOMER, BART).
3. grandfatherOf(X, Y) :-  

    fatherOf(X, Z),  

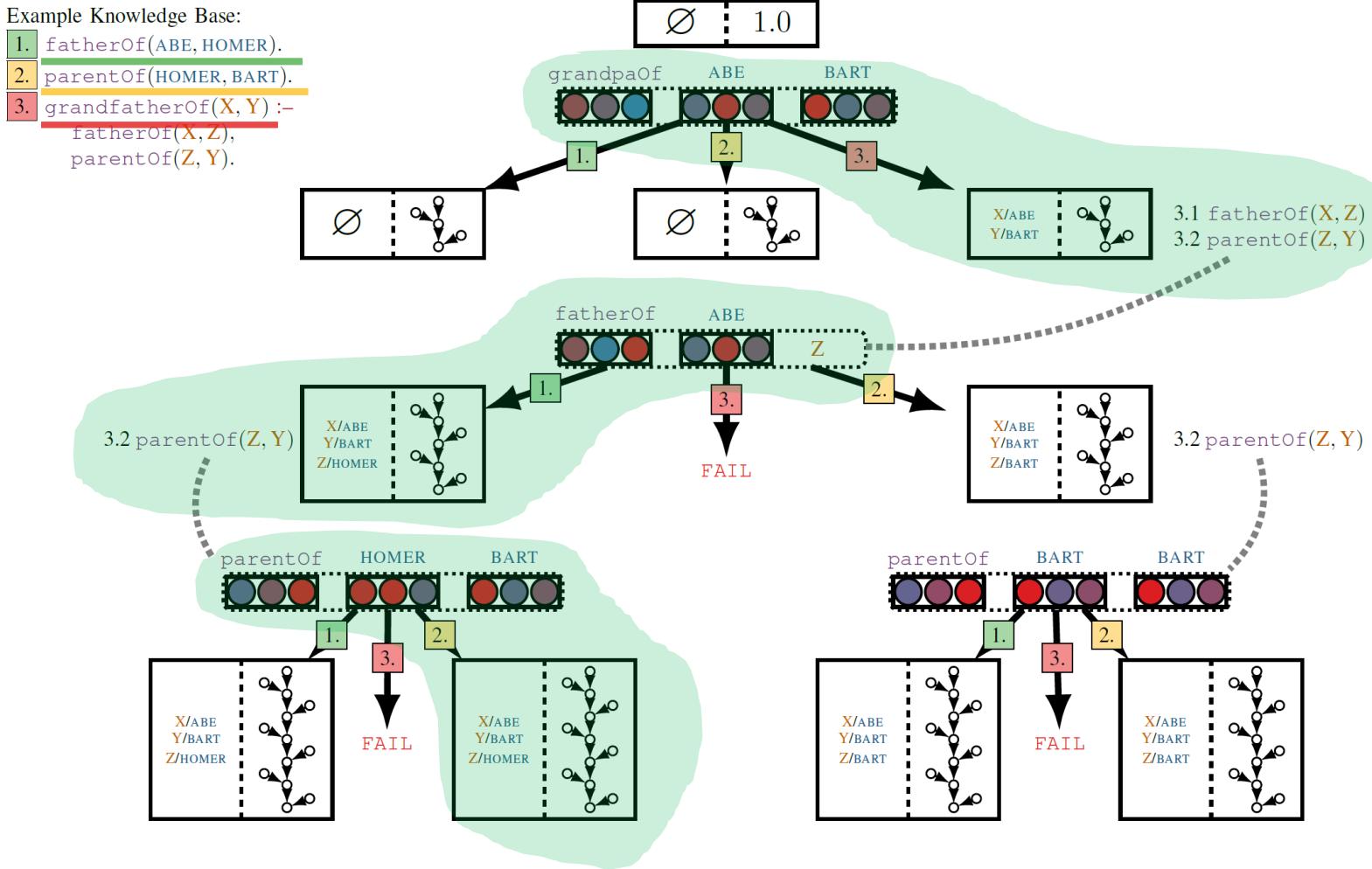
    parentOf(Z, Y).
  
```



Problem: $grandpaOf \neq grandfatherOf$ in backward chaining.

Solution: $v_{grandpaOf} \approx v_{grandfatherOf}$ in neural backward chaining.

(8): Neural Backward Chaining for Knowledge Graph Reasoning

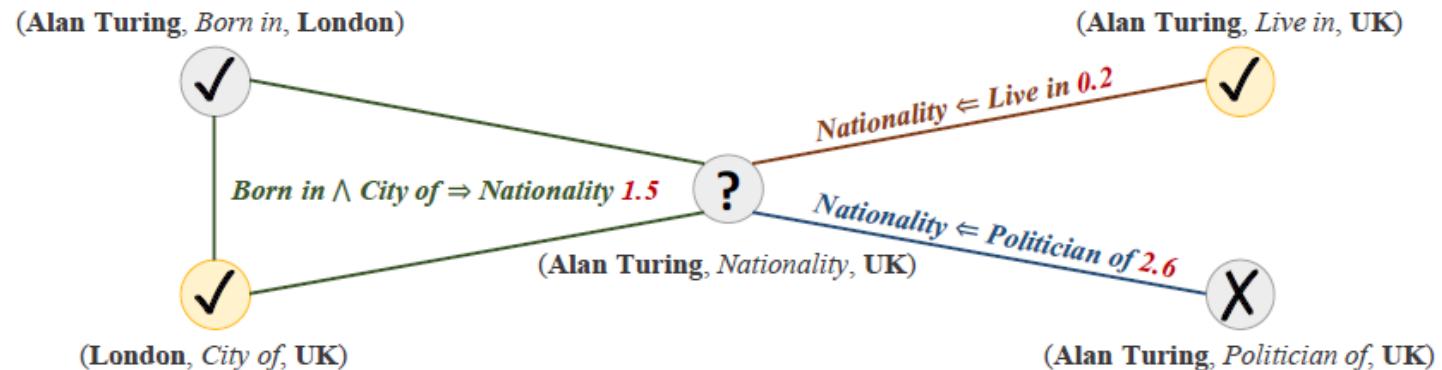


(9): Probabilistic Logic Neural Networks (pLogicNet) for Knowledge Graph Reasoning

0.2 $\text{Live}(X, Y) \Rightarrow \text{Nationality}(X, Y)$

2.6 $\text{Politician_of}(X, Y) \Rightarrow \text{Nationality}(X, Y)$

1.5 $\text{Born}(X, Y) \wedge \text{City_of}(Y, Z) \Rightarrow \text{Nationality}(X, Z)$



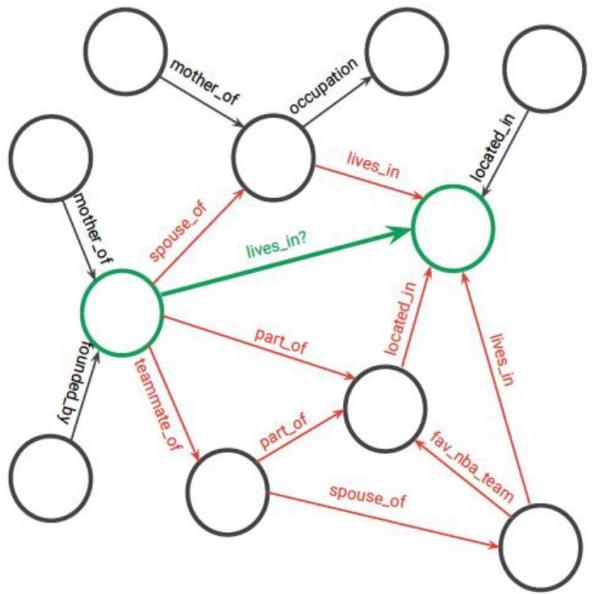
Given a set of first-order logic rules, pLogicNet uses an MLN to model the joint distribution of observed and unobserved triples:

$$\log p_w(\mathbf{v}_O) \geq \mathcal{L}(q_\theta, p_w) = \mathbb{E}_{q_\theta(\mathbf{v}_H)} [\log p_w(\mathbf{v}_O, \mathbf{v}_H) - \log q_\theta(\mathbf{v}_H)]$$

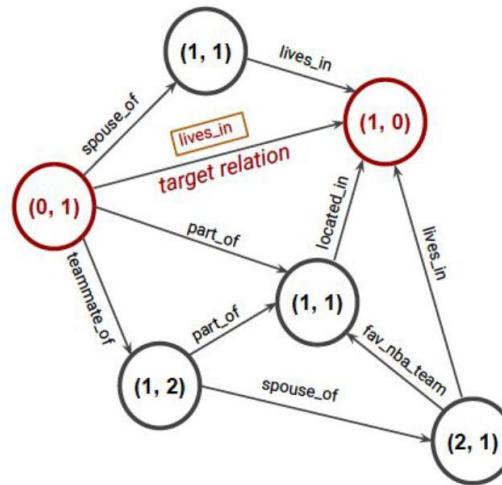
In the E-step: $q_\theta(\mathbf{v}_H) = \prod_{(h,r,t) \in H} q_\theta(\mathbf{v}_{(h,r,t)}) = \prod_{(h,r,t) \in H} \text{Ber}(\mathbf{v}_{(h,r,t)} | f(\mathbf{x}_h, \mathbf{x}_r, \mathbf{x}_t))$

In the M-step: $\ell_{PL}(w) \triangleq \mathbb{E}_{q_\theta(\mathbf{v}_H)} [\sum_{h,r,t} \log p_w(\mathbf{v}_{(h,r,t)} | \mathbf{v}_O \cup H \setminus (h,r,t))] = \mathbb{E}_{q_\theta(\mathbf{v}_H)} [\sum_{h,r,t} \log p_w(\mathbf{v}_{(h,r,t)} | \mathbf{v}_{MB(h,r,t)})]$

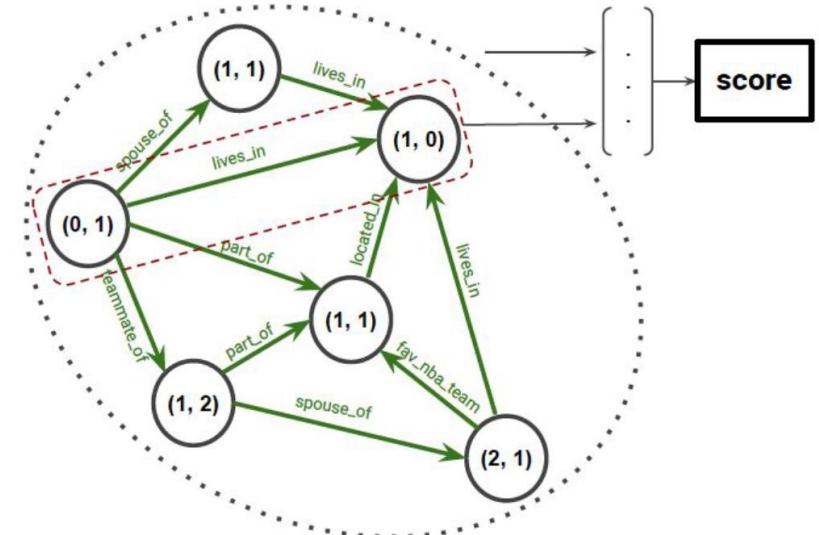
(10) Graph Inductive Learning (Grall) for Knowledge Graph Reasoning



1. Extract subgraph around candidate edge

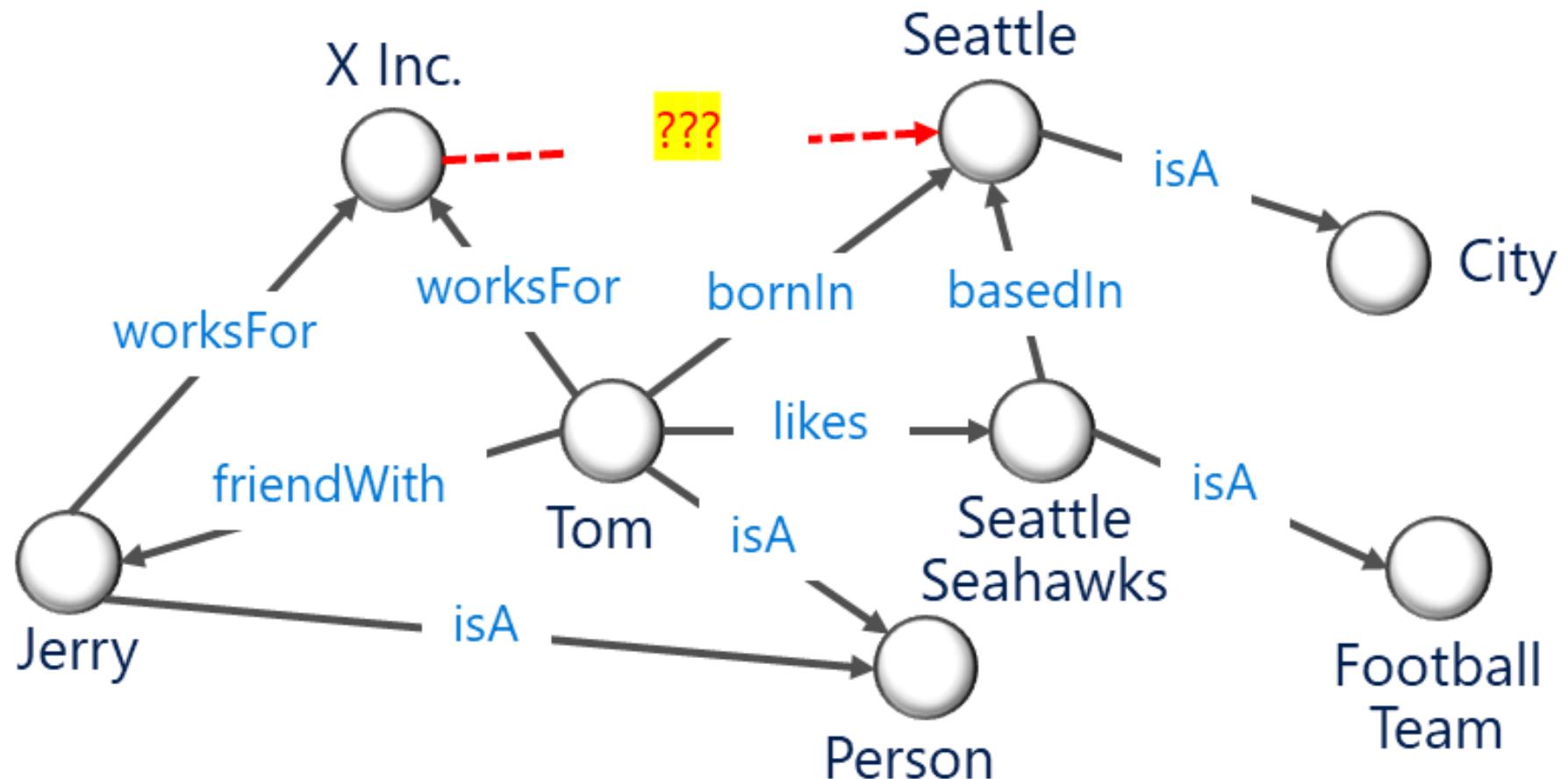


2. Assign structural labels to nodes

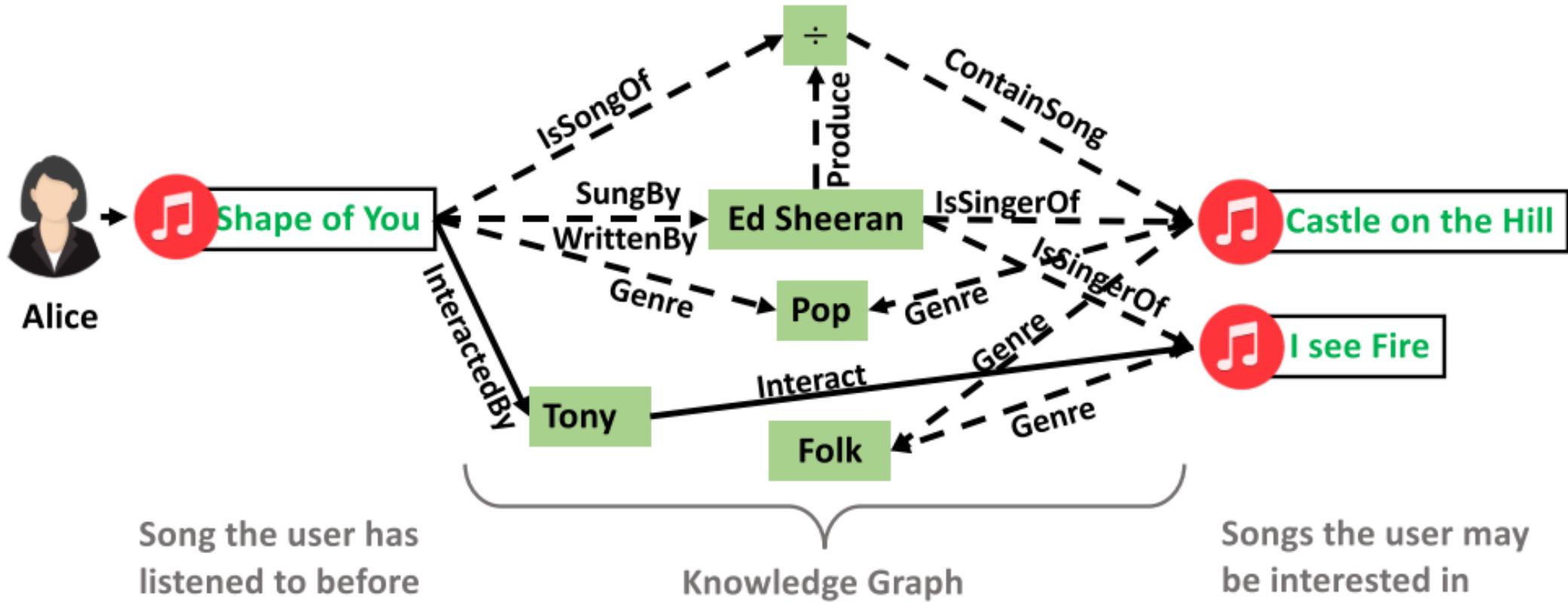


3. Run GNN on the extracted subgraph

Successful Application: Knowledge Graph Completion



Successful Application: Recommendation System

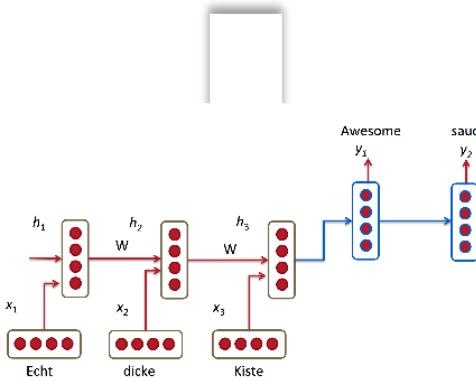


4 Typical Neural-Symbolic Reasoning Approaches

- Knowledge Graph Reasoning
- Neural Semantic Parsing 
- Neural Module Networks
- Symbolic Knowledge as Constraints

Neural Semantic Parsing

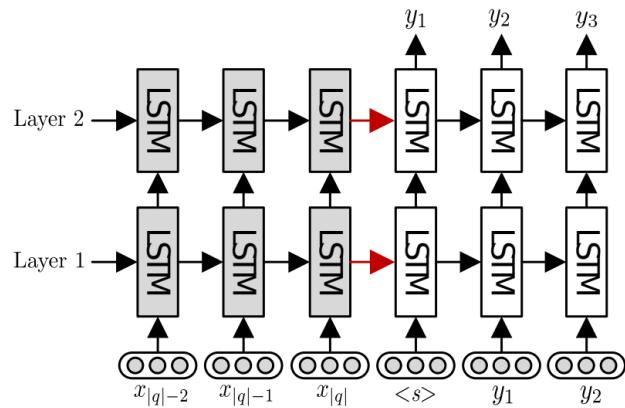
What city was Donald Trump born ?



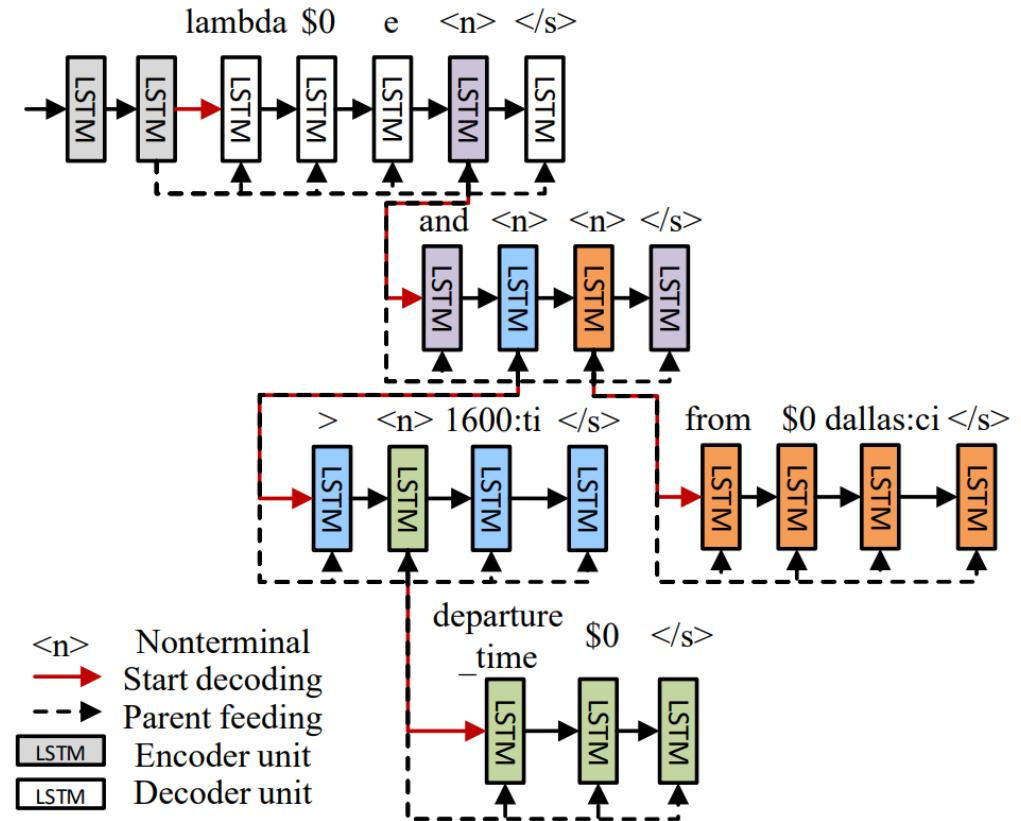
Neural Semantic Parser

$\lambda x. \text{Type}(\text{City}, x) \wedge \text{Place_of_Birth}(\text{Donald Trump}, x)$

(1): Seq2Seq & Seq2Tree

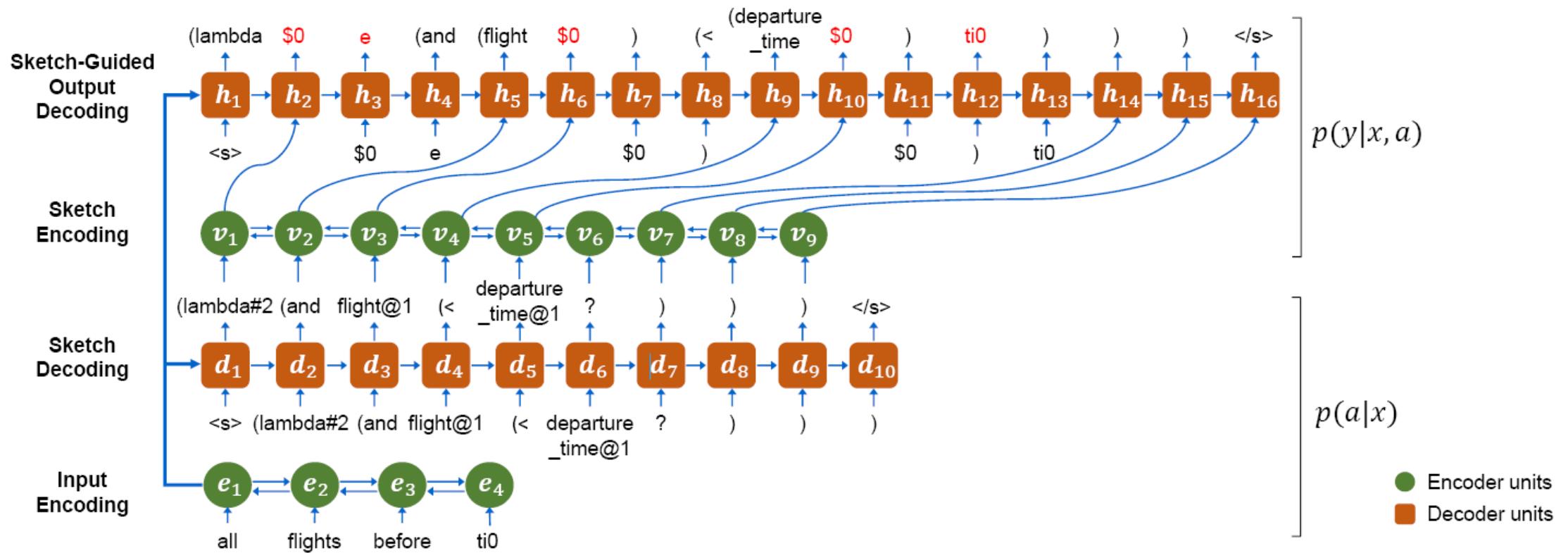


Seq2Seq

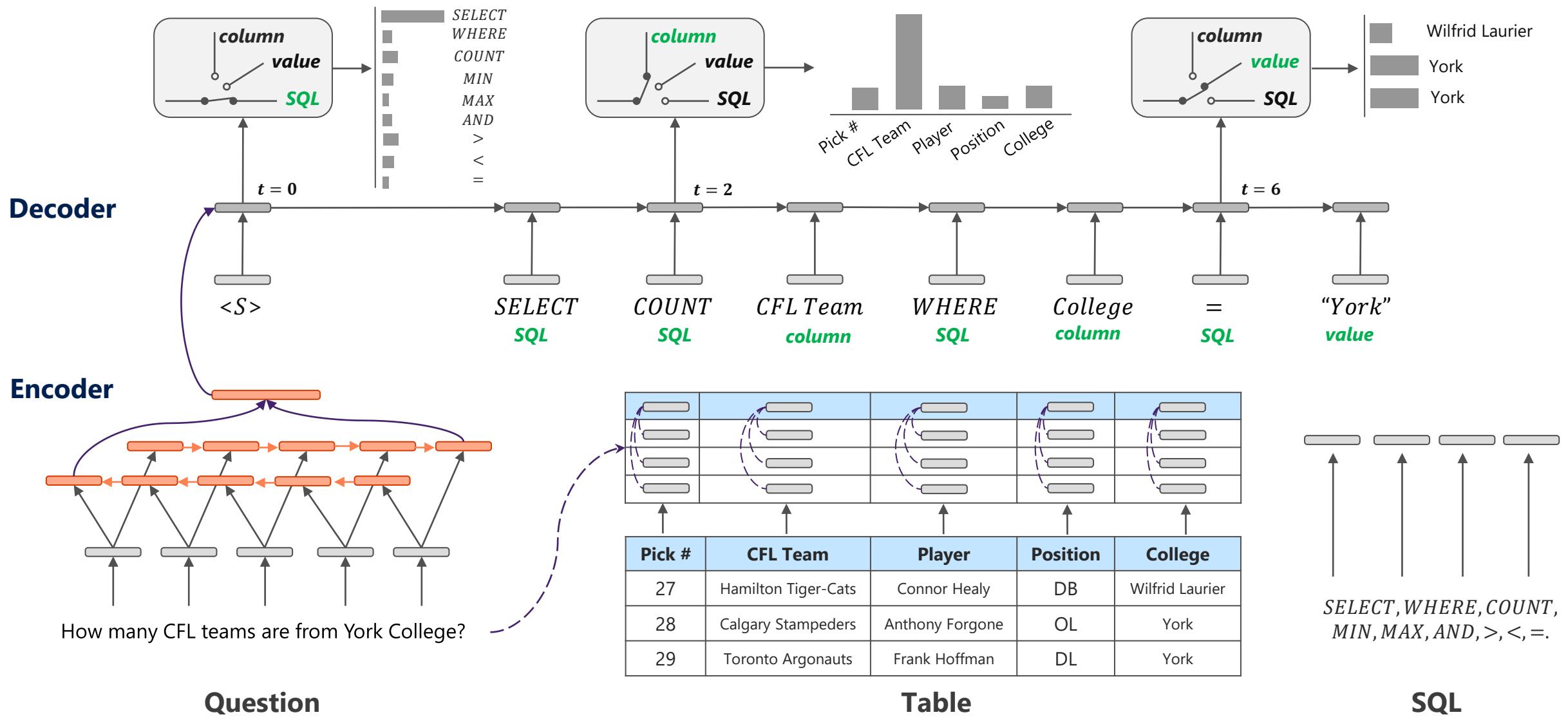


Seq2Tree

(2): Coarse-to-Fine Decoding

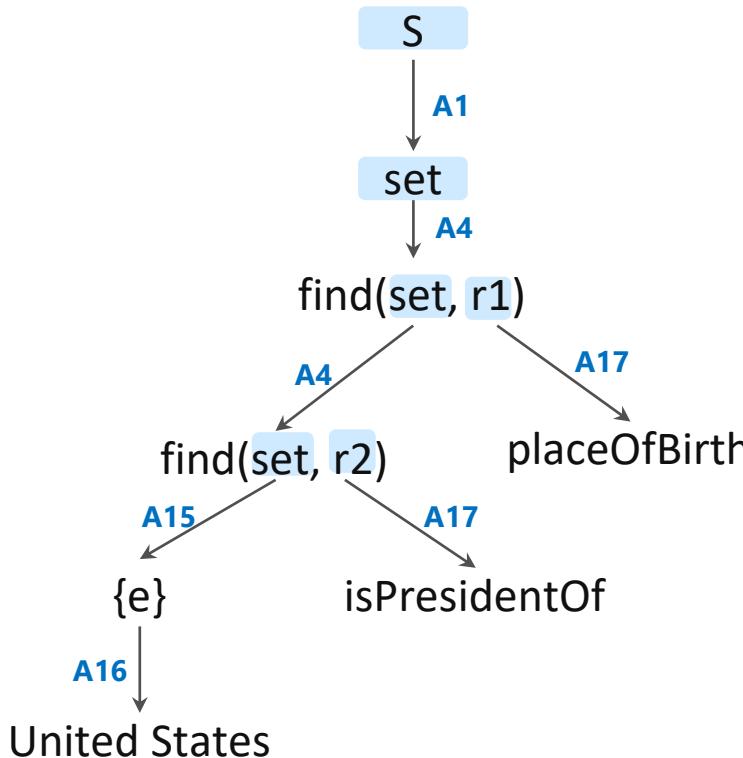
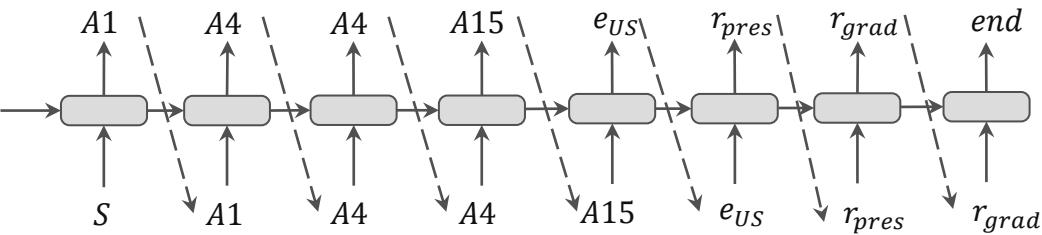


(3): Syntax-aware SQL Generation



(4): Dialog-to-Action (single-turn)

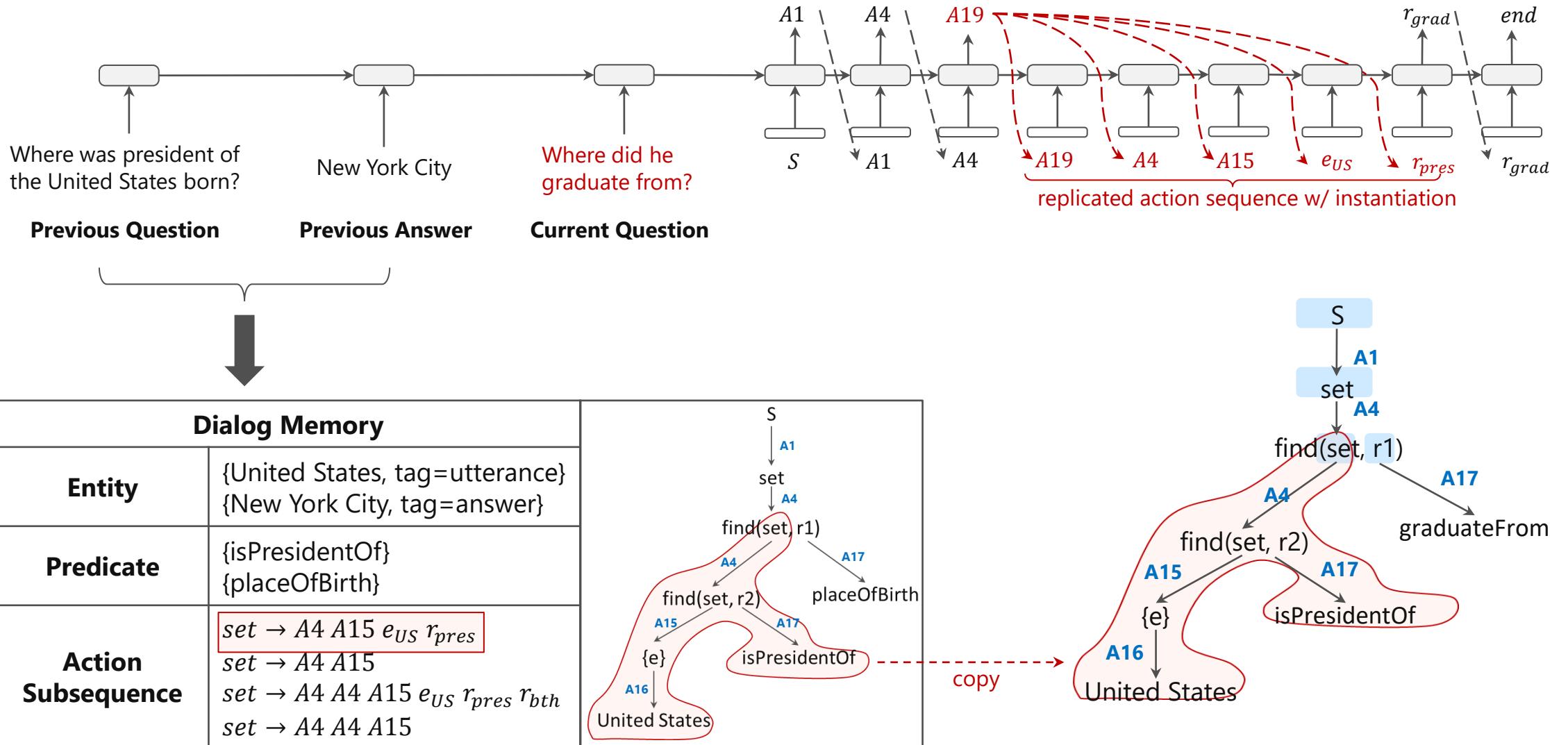
Where was the president of the United States born?



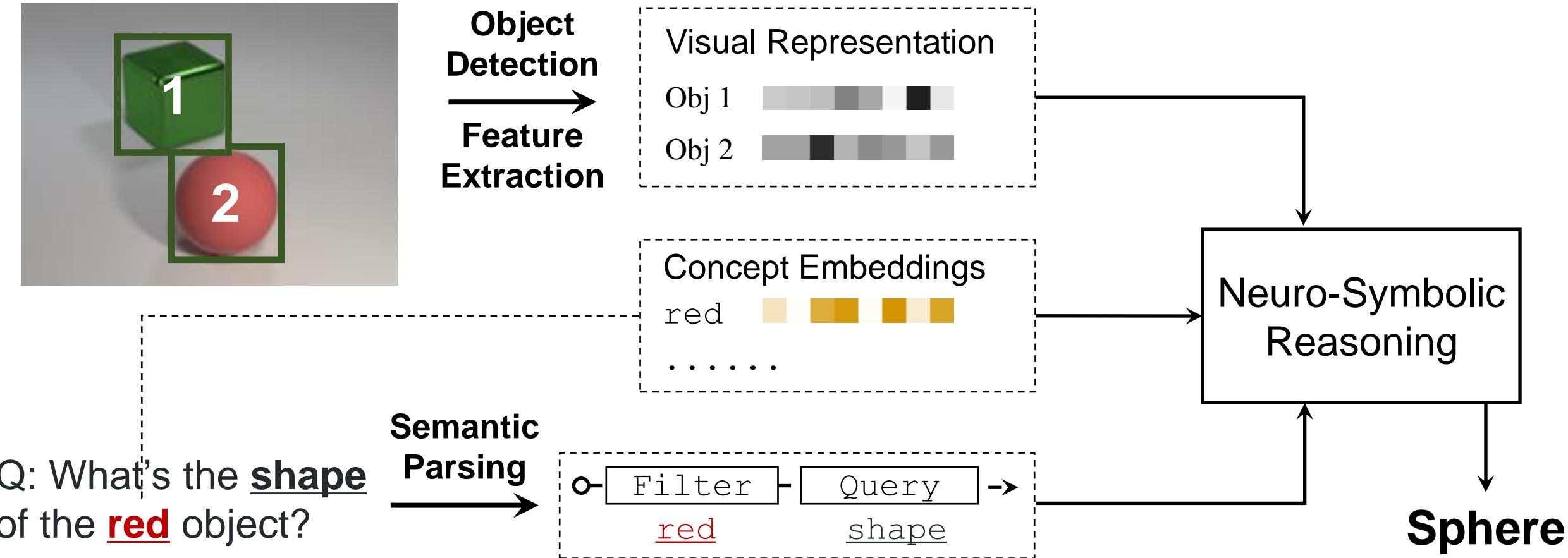
- A1: $S \rightarrow set$
- A4: $set \rightarrow find(set, r_1)$
- A4: $set \rightarrow find(set, r_2)$
- A15: $set \rightarrow \{e\}$
- A16: $e \rightarrow United\ States$
- A17: $r_2 \rightarrow isPresidentOf$
- A17: $r_1 \rightarrow placeOfBirth$

Action	Operation	Description
A1-A3	$S \rightarrow Set \mid Num \mid Bool$	S is start symbol
A4	$Set \rightarrow Find(R, E)$	Set of entities with a r edge to e
A5	$Num \rightarrow Count(Set)$	Total number of set
A6	$Bool \rightarrow (\epsilon, E, Set)$	Whether $\epsilon \in Set$
A7	$Set \rightarrow Set \cup Set$	Union of Sets
A8	$Set \rightarrow Set \cap Set$	Intersection of Sets
A9	$Set \rightarrow Set - Set$	Difference of Sets
A10	$Set \rightarrow larger(set, r, num)$	Entity from set linking to more than num entities with relation r
A11	$Set \rightarrow less(set, r, num)$	Entity from set linking to less than num entities with relation r
A12	$Set \rightarrow equal(set, r, num)$	Entity from set linking to num entities with relation r
A13	$Set \rightarrow argmax(set, r, num)$	Entity from set linking to most entities with relation r
A14	$Set \rightarrow argmin(set, r, num)$	Entity from set linking to least entities with relation r
A15	$Set \rightarrow \{e\}$	
A16-A18	$e \mid r \mid num \rightarrow constant$	instantiation for entity e, predicate r or number num
A19-A21	$Set \mid Num \mid Bool \rightarrow action(i-1)$	Replicate previous operation sequence

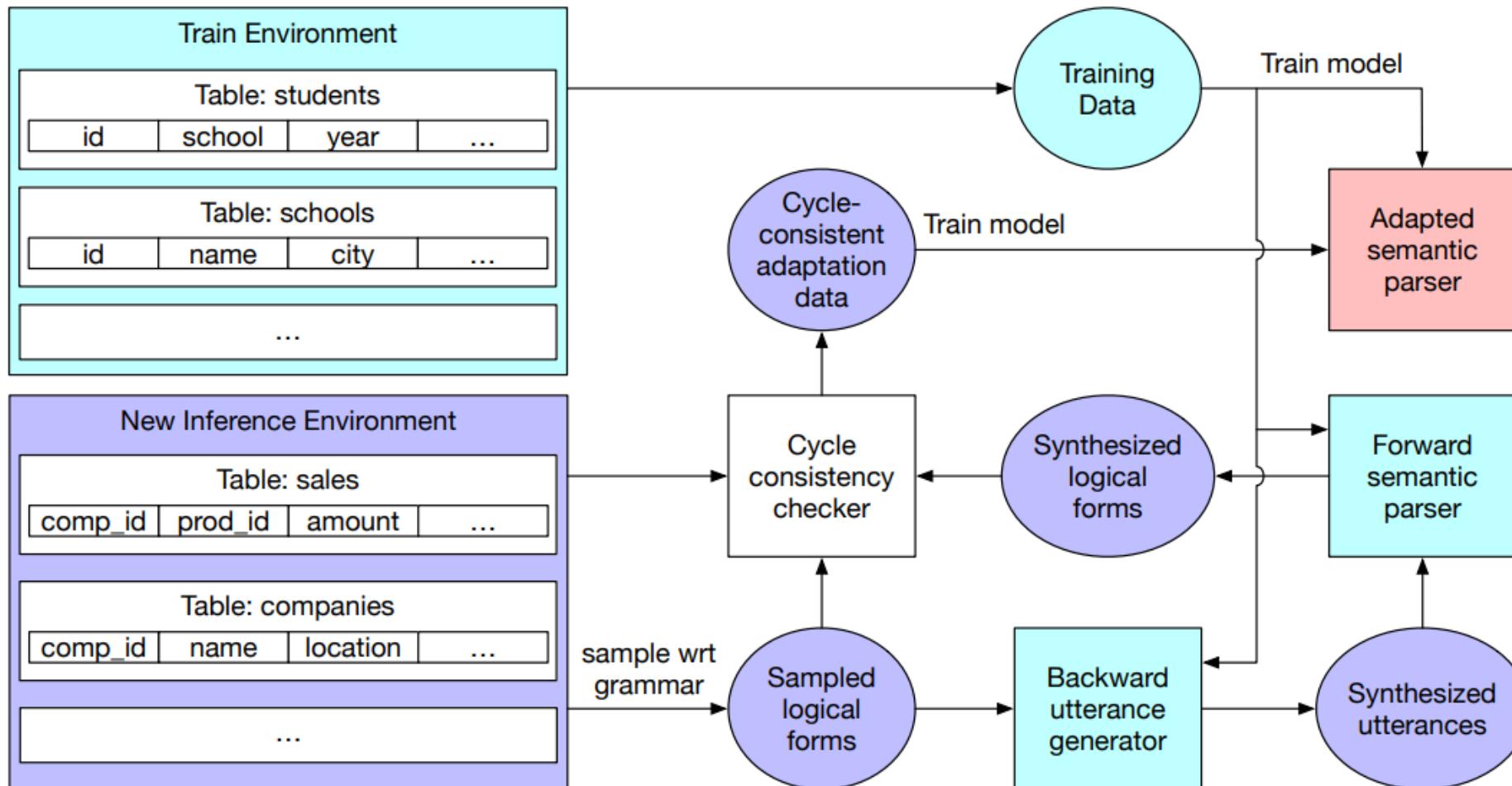
(4): Dialog-to-Action (multi-turn)



(5): Neuro-Symbolic Concept Learner



(6): Zero-shot Executable Semantic Parsing



Successful Application: Knowledge QA in Search Engine

The image displays two screenshots of the Bing search engine interface, illustrating its implementation of Knowledge Graph technology.

Search Query 1: "who is president of united states in 2000"

The search bar shows the query "who is president of united states in 2000". The results page, titled "Microsoft Show results from Microsoft", indicates 18,300,000 results. It features a large image of Bill Clinton with the caption "President in 2000" and the name "Bill Clinton" highlighted with a red box.

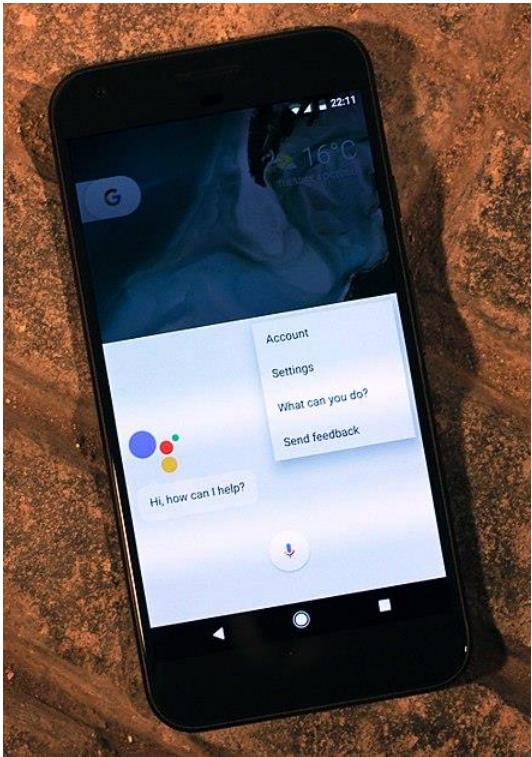
Search Query 2: "how high is yao ming's wife"

The search bar shows the query "how high is yao ming's wife". The results page, also titled "Microsoft Show results from Microsoft", indicates 4,920,000 results. It features a large image of Ye Li with the caption "Ye Li · Height" and the height "6' 3\"" highlighted with a red box. Below this, a grid of images and names for other basketball players is shown:

- Yao Ming**: 7' 6"
- Yi Jianlian**: 6' 11"
- LeBron James**: 6' 8"
- Wang Zhizhi**: 7' 1"

Bing Knowledge QA

Successful Application: Intelligent Virtual Assistant



Google Assistant running
on a Pixel XL smartphone



Apple TV remote control, with
which users can ask Siri the virtual
assistant to find content to watch



Amazon Echo smart
speaker running the
Alexa virtual assistant

4 Typical Neural-Symbolic Reasoning Approaches

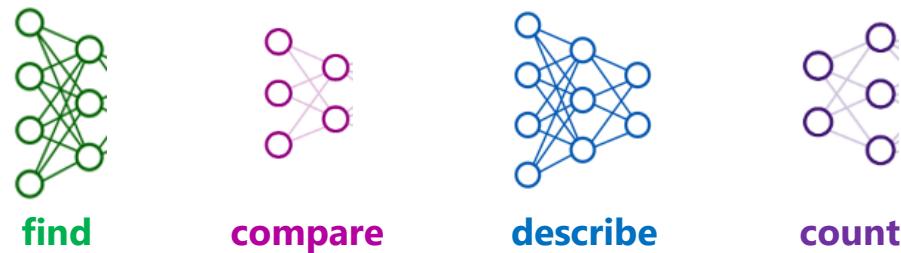
- Knowledge Graph Reasoning
- Neural Semantic Parsing
- Neural Module Networks 
- Symbolic Knowledge as Constraints

Neural Module Networks

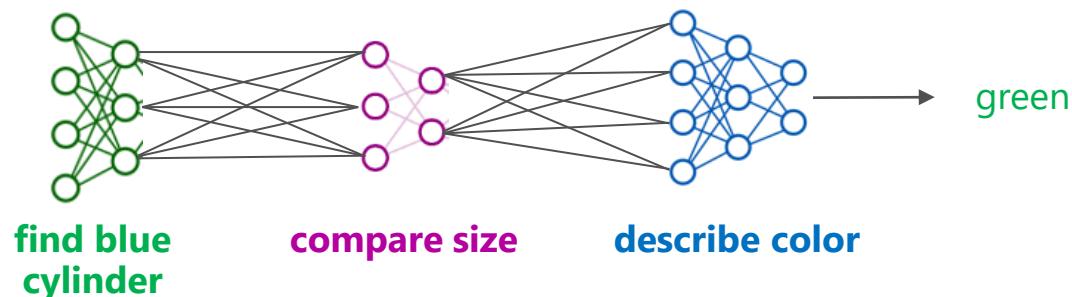
What color is the thing with the same size as the blue cylinder?



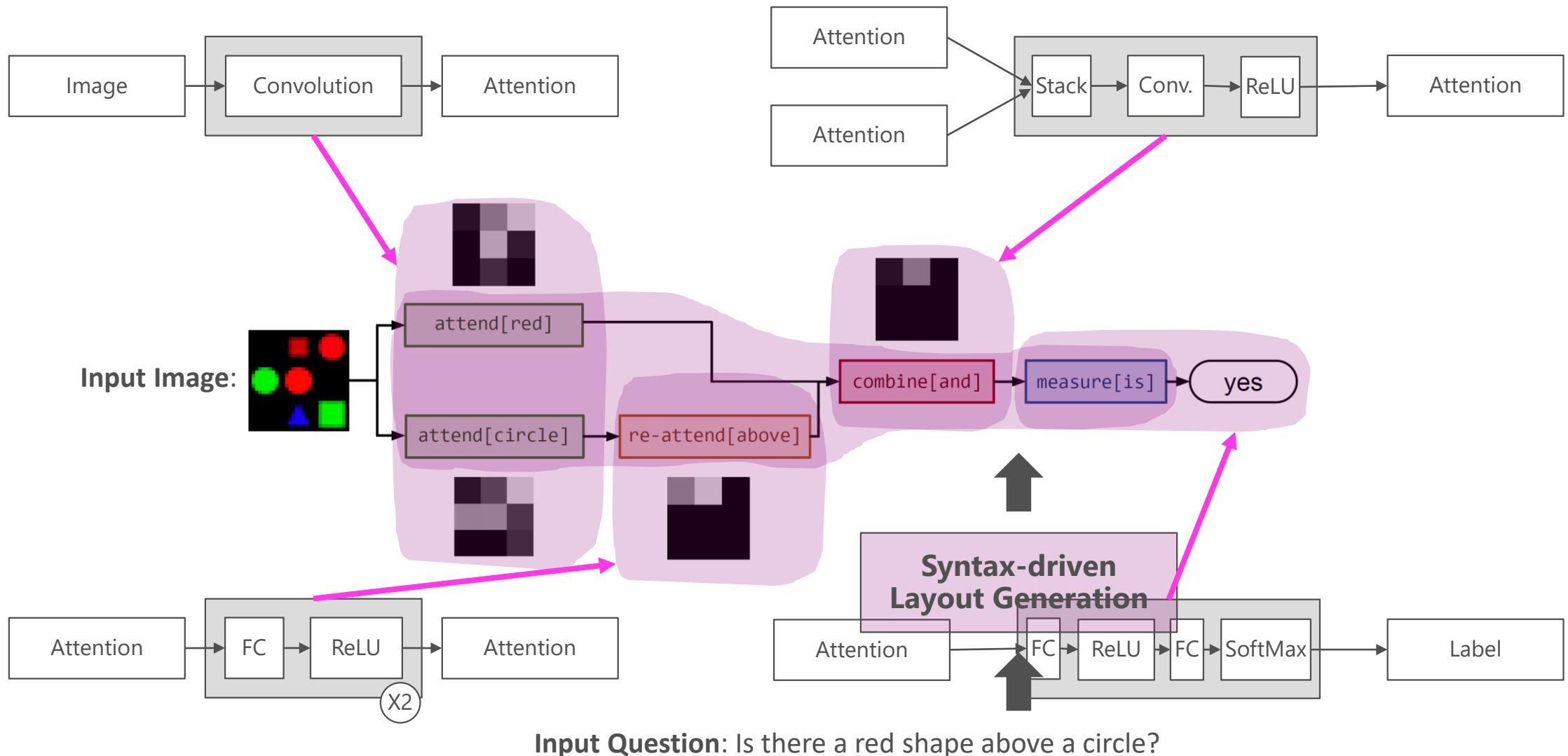
(1) define a collection of neural "modules", each of which implements a single step of reasoning



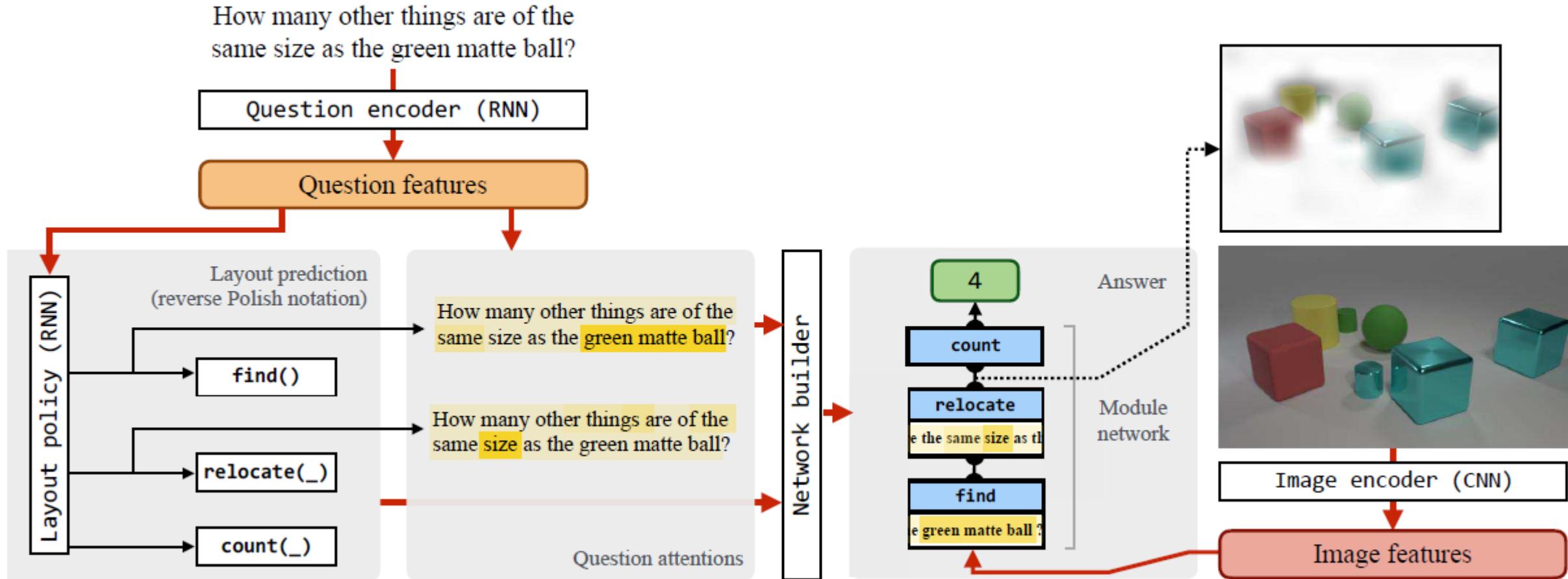
(2) compose neural modules to form different neural networks for different inputs



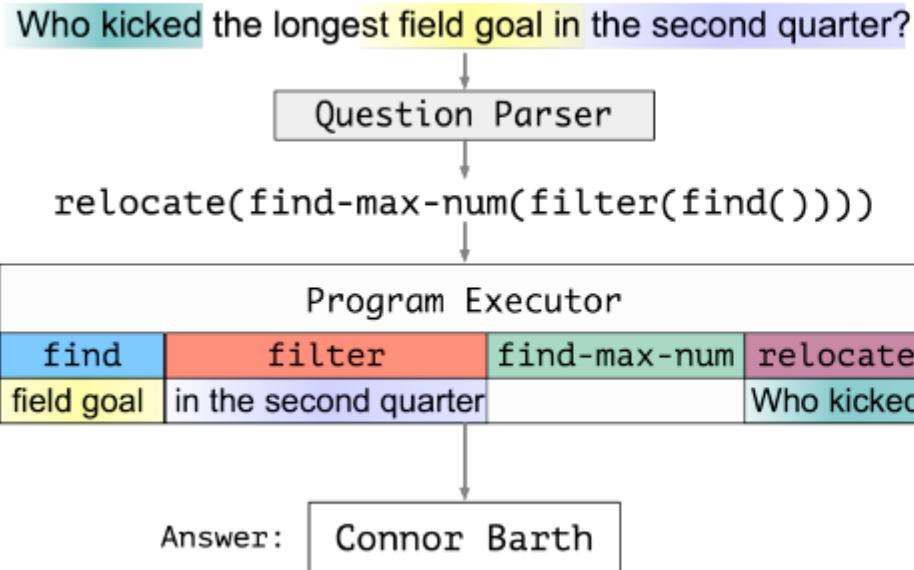
(1): Neural Module Networks for Visual Reasoning



(2): End-to-End Neural Module Networks for Visual Reasoning



(3): Neural Module Networks for Text Reasoning

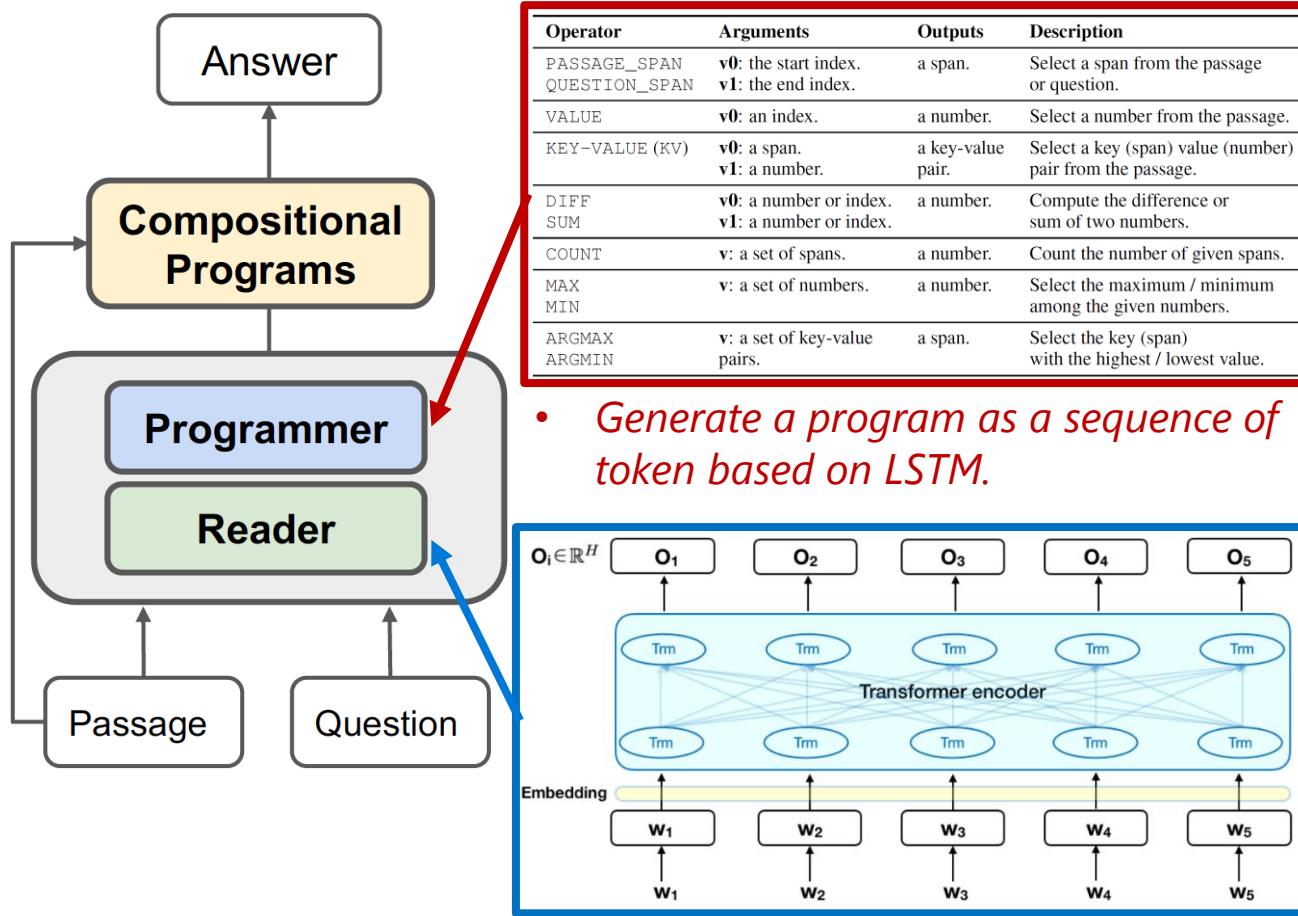


In the first quarter, Buffalo trailed as Chiefs QB Tyler Thigpen completed a 36-yard TD pass to RB Jamaal Charles. The Bills responded with RB Marshawn Lynch getting a 1-yard touchdown run. In the second quarter, Buffalo took the lead as kicker Rian Lindell made a 21-yard and a 40-yard field goal. Kansas City answered with Thigpen completing a 2-yard TD pass. Buffalo regained the lead as Lindell got a 39-yard field goal. The Chiefs struck with kicker Connor Barth getting a 45-yard field goal, yet the Bills continued their offensive explosion as Lindell got a 34-yard field goal, along with QB Edwards getting a 15-yard TD run. In the third quarter, Buffalo continued its poundings with Edwards getting a 5-yard TD run, while Lindell got himself a 48-yard field goal. Kansas City tried to rally as Thigpen completed a 45-yard TD pass to WR Mark Bradley, yet the Bills replied with Edwards completing an 8-yard TD pass to WR Josh Reed. In the fourth quarter, Edwards completed a 17-yard TD pass to TE Derek Schouman.

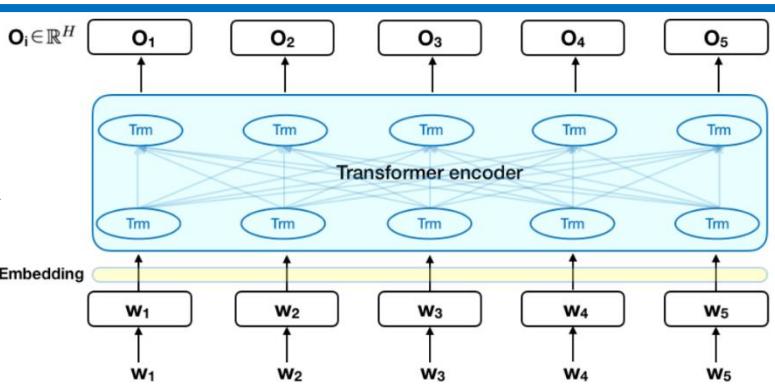
Reasoning Steps

- 1) find all instances of “field goal” in the paragraph,
- 2) select the ones “in the second quarter”,
- 3) find their lengths, compute the “longest” of them,
- 4) and then find “who kicked” it.

(4): Neural Symbolic Reader



- *Generate a program as a sequence of token based on LSTM.*



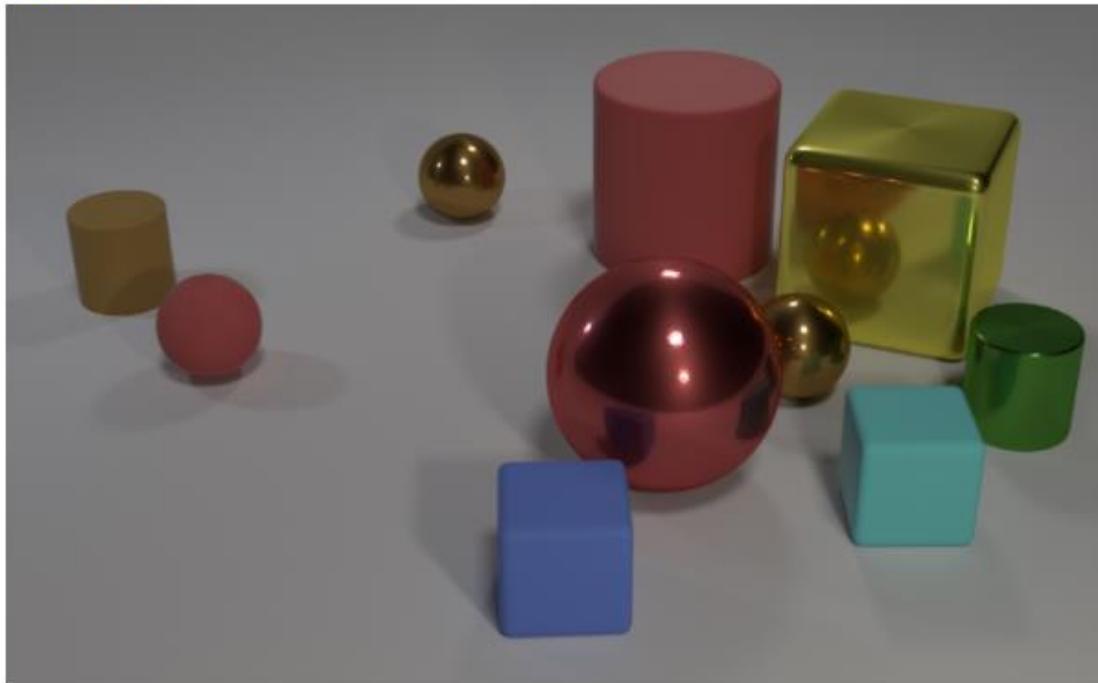
- *Encode input tokens into vector space representations based on BERT.*

Operator	Arguments	Outputs	Description
PASSAGE_SPAN	v_0 : the start index.	a span.	Select a span from the passage or question.
QUESTION_SPAN	v_1 : the end index.		
VALUE	v_0 : an index.	a number.	Select a number from the passage.
KEY-VALUE (KV)	v_0 : a span. v_1 : a number.	a key-value pair.	Select a key (span) value (number) pair from the passage.
DIFF	v_0 : a number or index.	a number.	Compute the difference or sum of two numbers.
SUM	v_1 : a number or index.		
COUNT	v : a set of spans.	a number.	Count the number of given spans.
MAX	v : a set of numbers.	a number.	Select the maximum / minimum among the given numbers.
MIN			
ARGMAX	v : a set of key-value pairs.	a span.	Select the key (span) with the highest / lowest value.
ARGMIN			

Passage	Question & Answer
	Multiple spans ...the population was spread out with 26.20% under the age of 18, 9.30% from 18 to 24, 26.50% from 25 to 44 , 23.50% from 45 to 64 , and 14.60% who were 65 years of age or older...
	Question: Which groups in percent are larger than 16%? Program: PASSAGE_SPAN(26,30), PASSAGE_SPAN(46,48), PASSAGE_SPAN(55,57) Result: 'under the age of 18', '25 to 44', '45 to 64'
	Date When major general Nathanael Greene took command in the south, Marion and lieutenant colonel Henry Lee were ordered in January 1781 ... On August 31 , Marion rescued a small American force trapped by 500 British soldiers...
	Question: When did Marion rescue the American force? Program: PASSAGE_SPAN(71,71), PASSAGE_SPAN(72,72), PASSAGE_SPAN(32,32) Result: 'August', '31', '1781'
	Numerical operations ...Lassen county had a population of 34,895 . The racial makeup of Lassen county was 25,532 (73.2%) white (U.S. census), 2,834 (8.1%) African American (U.S. census)... Program: DIFF(SUM(10,12)) Result: $34895 - (25532 + 2834) = 6529$

Successful Application: Visual Reasoning & QA

Questions in CLEVR test various aspects of visual reasoning including **attribute identification**, **counting**, **comparison**, **spatial relationships**, and **logical operations**.



Q: Are there an **equal number** of **large things** and **metal spheres**?

Q: **What size** is the **cylinder** that is **left of** the **brown metal** thing that is **left of** the **big sphere**?

Q: There is a **sphere** with the **same size as** the **metal cube**; is it **made of the same material** as the **small red sphere**?

Q: **How many** objects are **either small cylinders or red** things?

Successful Application: Text Reasoning & QA

DROP: A Reading Comprehension Benchmark Requiring Discrete Reasoning Over Paragraphs

Dheeru Dua, Yizhong Wang, Pradeep Dasigi,
Gabriel Stanovsky, Sameer Singh and Matt Gardner
NAACL 2019.

With system performance on existing reading comprehension benchmarks nearing or surpassing human performance, we need a new, hard dataset that improves systems' capabilities to actually *read* paragraphs of text. DROP is a crowdsourced, adversarially-created, 96k-question benchmark, in which a system must resolve references in a question, perhaps to multiple input positions, and perform discrete operations over them (such as addition, counting, or sorting). These operations require a much more comprehensive understanding of the content of paragraphs than what was necessary for prior datasets.

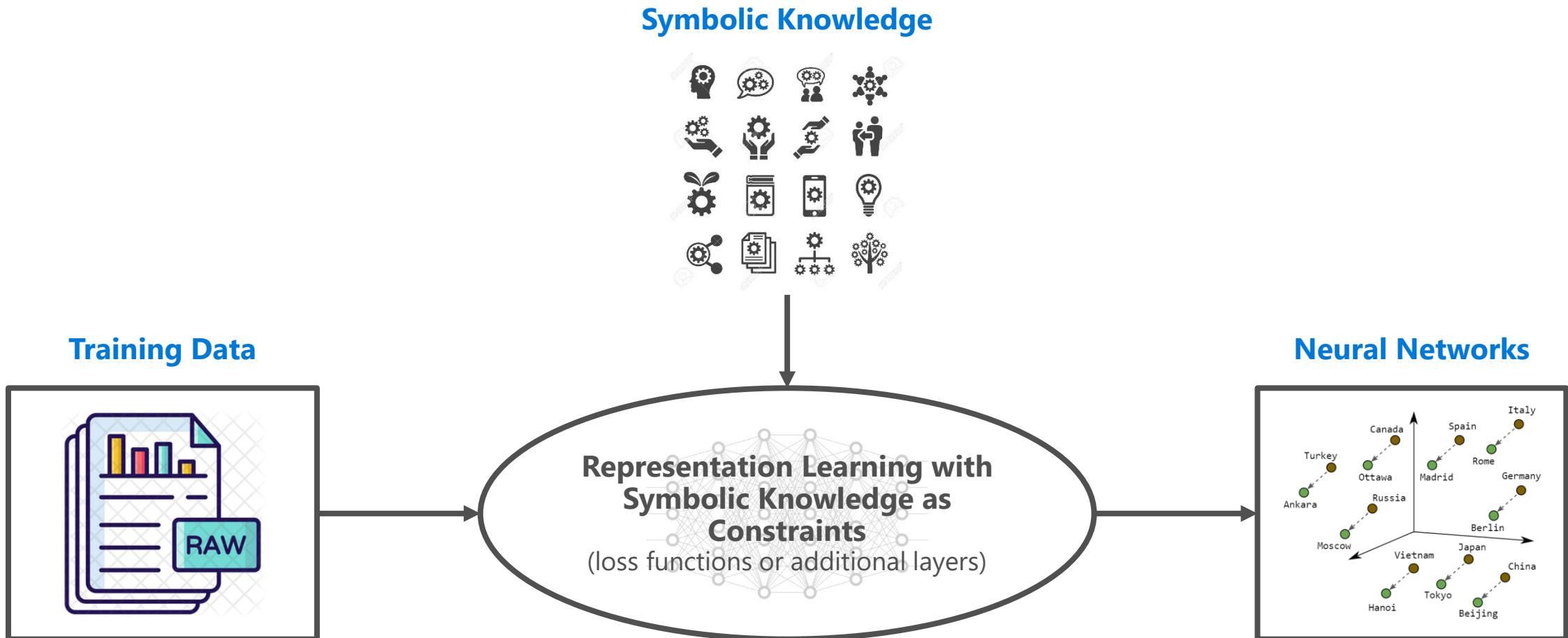
<https://allennlp.org/drop.html>

4 Typical Neural-Symbolic Reasoning Approaches

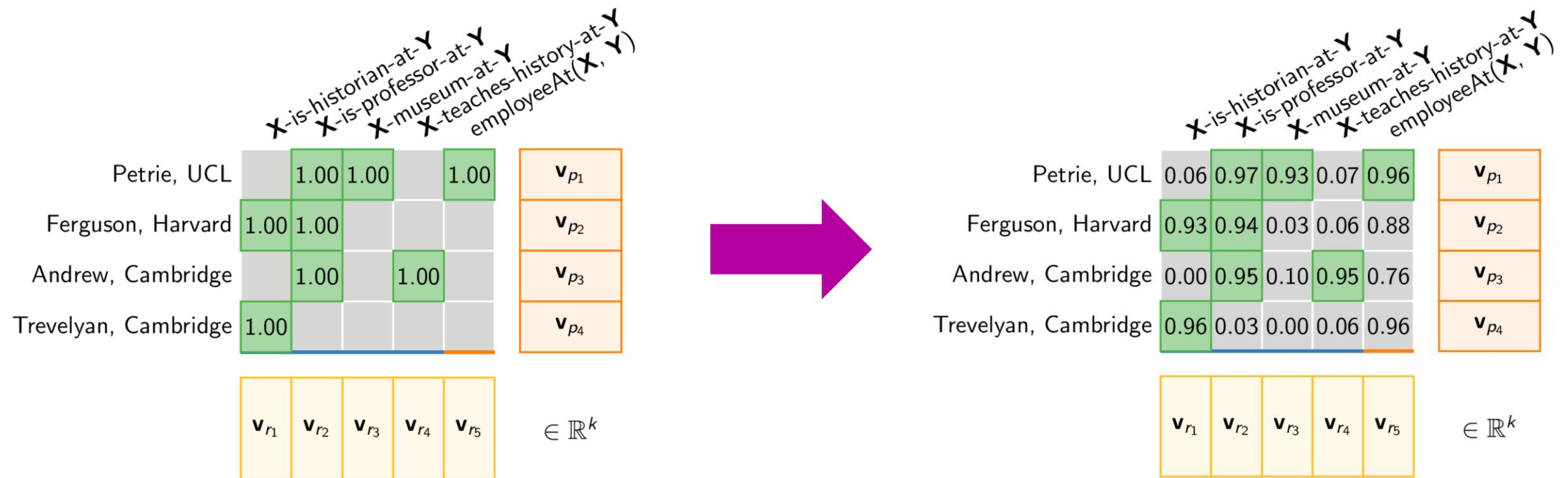
- Knowledge Graph Reasoning
- Neural Semantic Parsing
- Neural Module Networks
- Symbolic Knowledge as Constraints



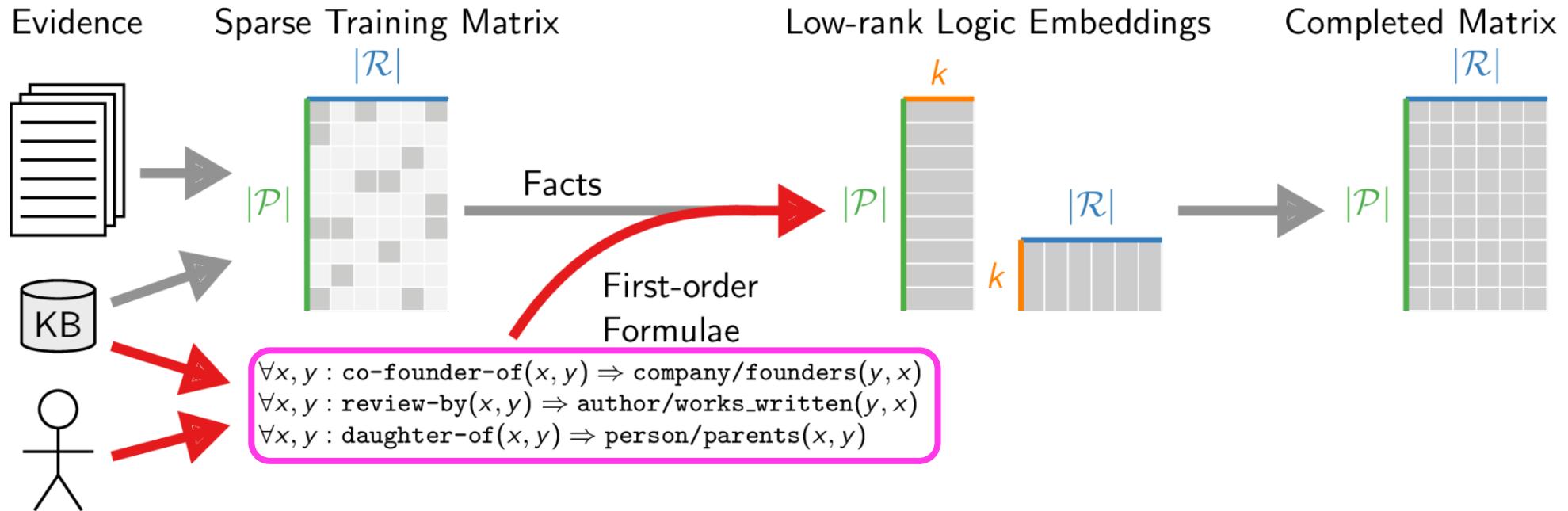
Symbolic Knowledge as Constraints



(1): First Order Logic Formulas for Relation Extraction



(1): First Order Logic Formulas for Relation Extraction

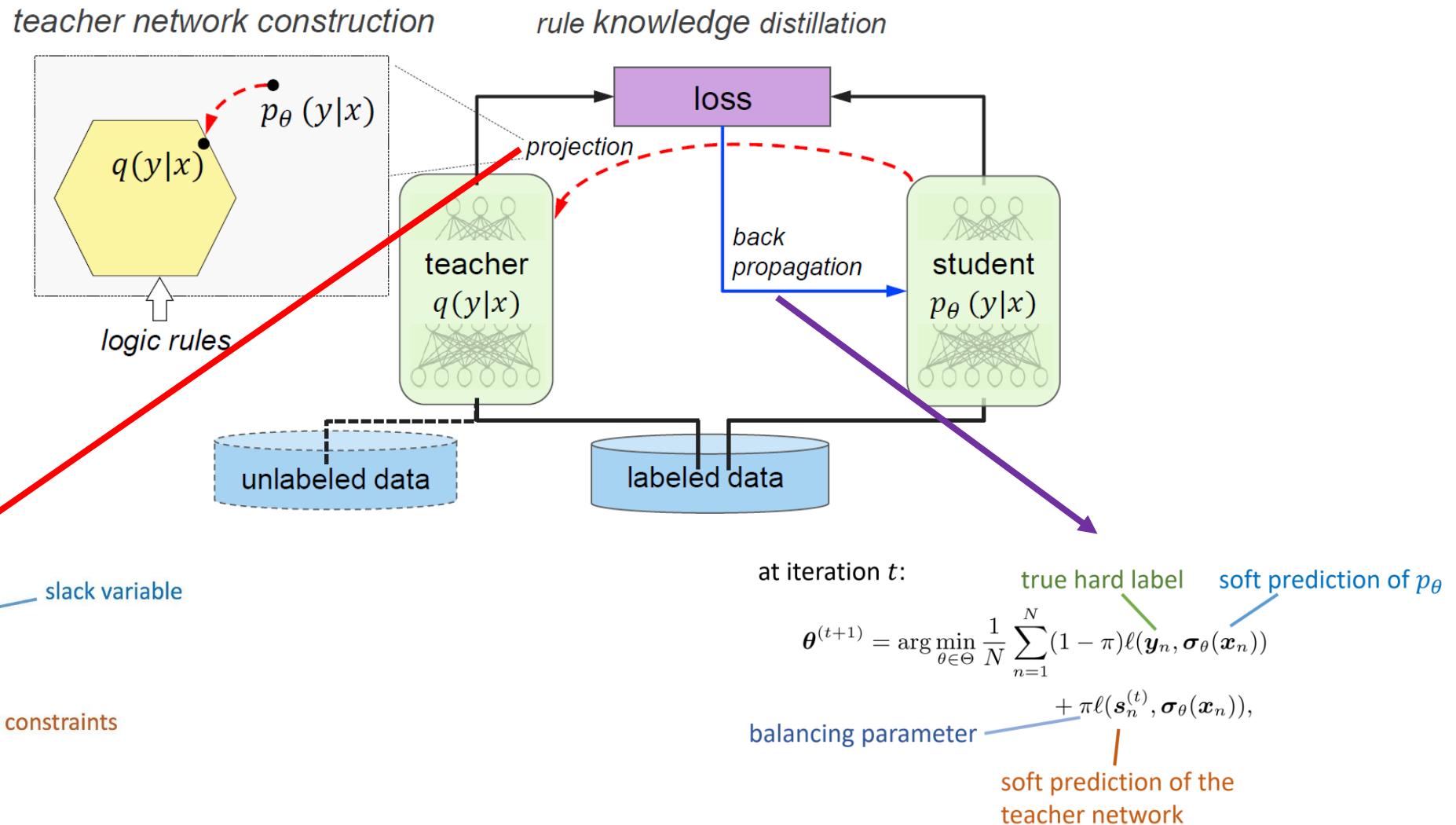


- **Joint-likelihood training objective:** $\min_{\{\mathbf{v}_s\}, \{\mathbf{v}_{ij}\}} - \sum_{\mathcal{F} \in \mathfrak{F}} \log([\mathcal{F}])$

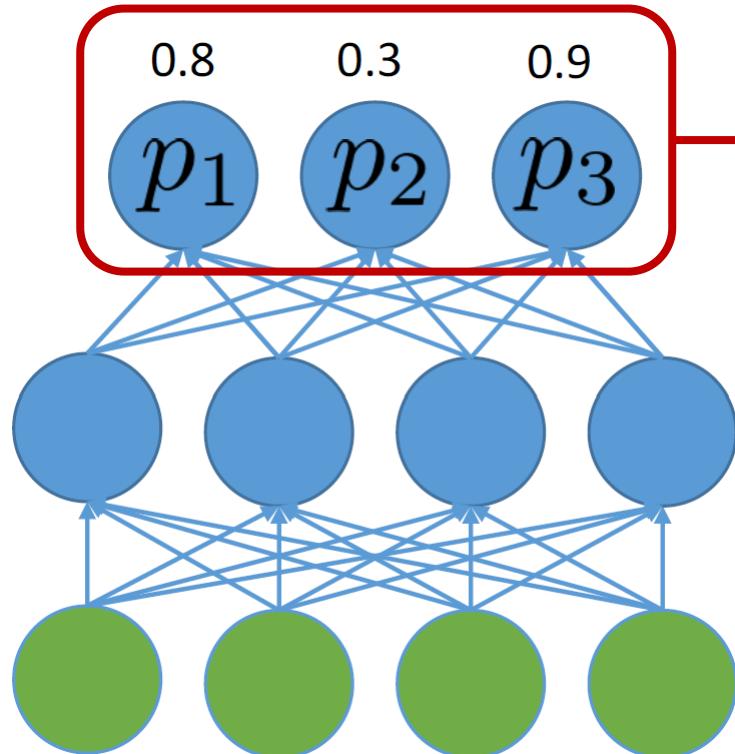
- If $\mathcal{F} = r_s(e_i, e_j)$ is a ground fact, then $[\mathcal{F}] = p(r_s(e_i, e_j) = \text{true}) = \sigma(\mathbf{v}_s \cdot \mathbf{v}_{ij})$;
- If \mathcal{F} is a first order logic formula, then $[\mathcal{F}]$ is defined based on T-norm.

- $[\neg \mathcal{A}] = 1 - [\mathcal{A}]$
- $[\mathcal{A} \wedge \mathcal{B}] = [\mathcal{A}][\mathcal{B}]$
- $[\mathcal{A} \vee \mathcal{B}] = [\mathcal{A}] + [\mathcal{B}] - [\mathcal{A}][\mathcal{B}]$
- $[\mathcal{A} \Rightarrow \mathcal{B}] = [\mathcal{A}]([B] - 1) + 1$
- ...

(2): Knowledge-Regularized Teacher Network



(3) Semantic Loss with Symbolic Knowledge



existing loss + $w \cdot$ semantic loss

$$L^s(\text{exactly-one}, p) \propto -\log \sum_{i=1}^n p_i \prod_{j=1, j \neq i}^n (1 - p_j)$$

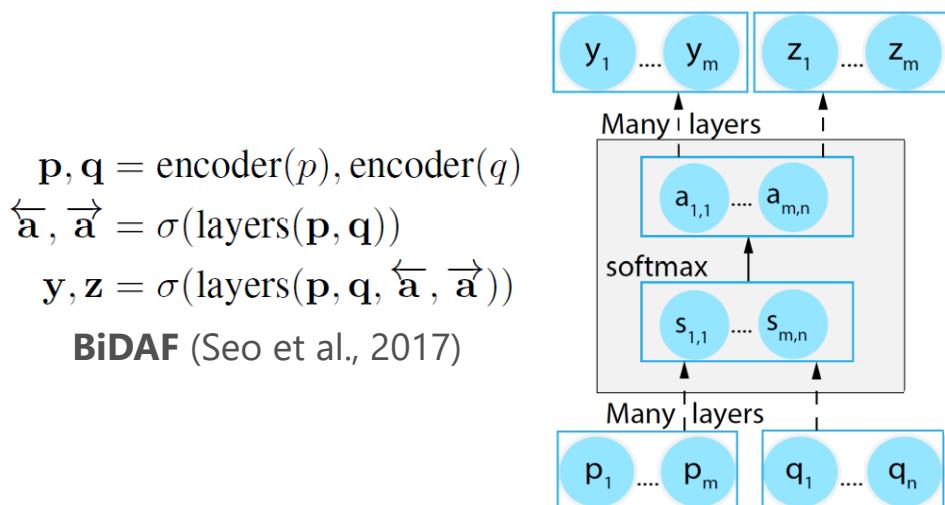
Table 2: FASHION. Test accuracy comparison between MLP with semantic loss and ladder nets.

Accuracy % with # of used labels	100	500	1000	ALL
Ladder Net (Rasmus et al., 2015)	81.46 (± 0.64)	85.18 (± 0.27)	86.48 (± 0.15)	90.46
Baseline: MLP, Gaussian Noise	69.45 (± 2.03)	78.12 (± 1.41)	80.94 (± 0.84)	89.87
MLP with Semantic Loss	86.74 (± 0.71)	89.49 (± 0.24)	89.67 (± 0.09)	89.81

(4) Augmenting Neural Networks with First-order Logic

Paragraph: Gaius Julius Caesar (July 100 BC – 15 March 44 BC), [Roman general], statesman, Consul and notable author of Latin prose, played a critical role in the events that led to the demise of the Roman Republic and the rise of the Roman Empire through his various military campaigns.

Question: Which [Roman general] is known for writing prose?



%Train	BiDAF	+ R_1	+ R_2	+ELMo	+ELMo, R_1
10%	57.5	61.5	60.7	71.8	73.0
20%	65.7	67.2	66.6	76.9	77.7
40%	70.6	72.6	71.9	80.3	80.9
100%	75.7	77.4	77.0	83.9	84.1

Constrained layer: $\mathbf{y} = g(W\mathbf{x} + \rho \cdot d(\mathbf{z}))$

FOL rule: $\forall i, j \in C, k_{i,j} \rightarrow a'_{i,j}$

- $k_{i,j}$: p_i is related to q_j in ConceptNet via edges {Synonym, DistinctFrom, IsA, Related};
- $a'_{i,j}$: constrained model decision for the above alignment.

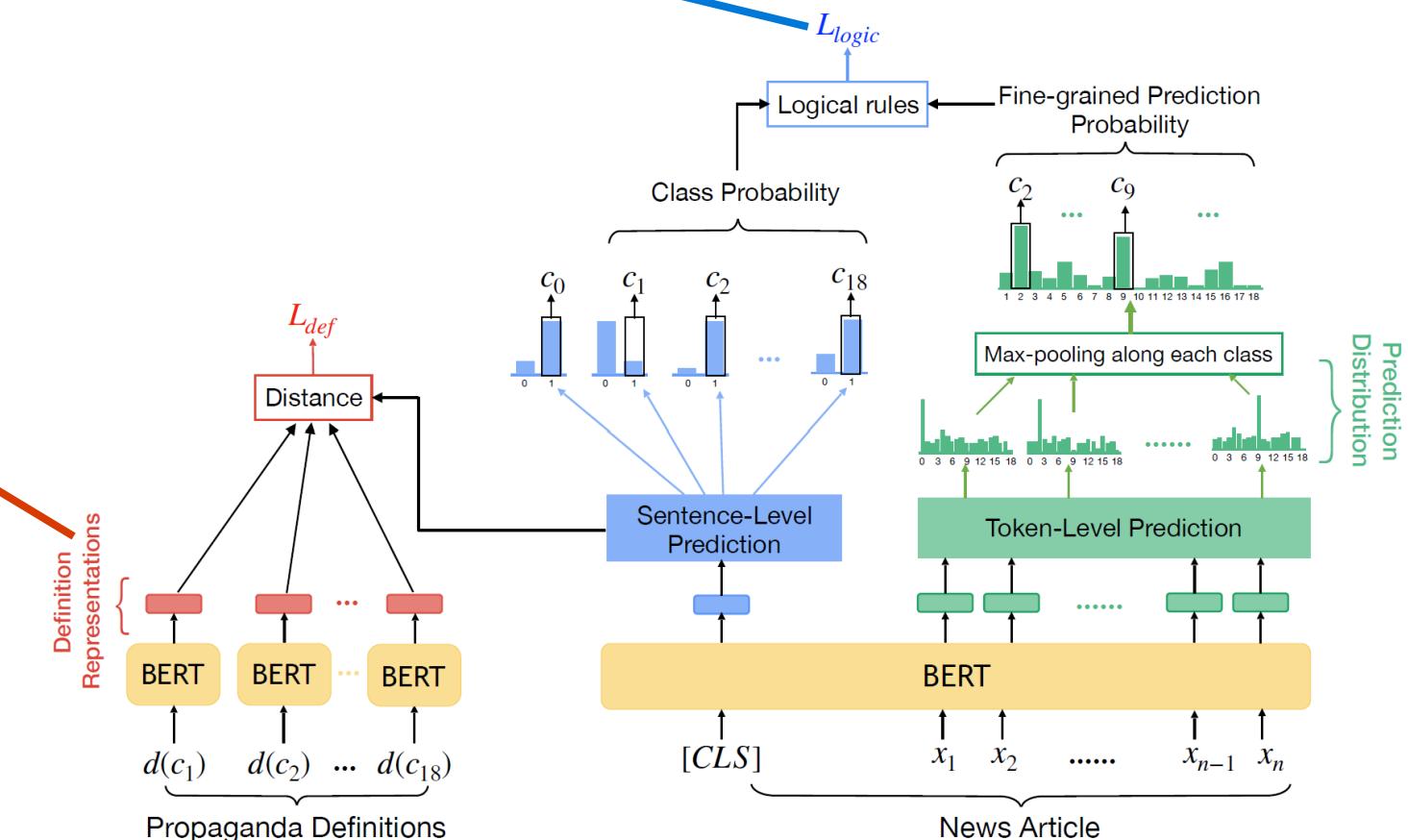
T-norm: $\bigvee_i Z_i = \min(1, \sum_i z_i)$

(5): Inject Logical Consistency Knowledge for Propaganda Detection

$$P(f_c(x) \Rightarrow g_c(x)) = P(f_c(x))(g_c(x) - 1) + 1$$

Propaganda Technique	Instances		
	Train	Dev	Test
Loaded Language	1,811	127	177
Name Calling,Labeling	931	68	86
Repetition	456	35	80
Doubt	423	23	44
Exaggeration,Minimisation	398	37	44
Flag-Waving	206	13	21
Appeal to fear-prejudice	187	32	20
Causal Oversimplification	170	24	7
Slogans	120	3	13
Black-and-White Fallacy	97	4	8
Appeal to Authority	91	2	23
Thought-terminating Cliches	70	4	5
Whataboutism	55	1	1
Reductio ad hitlerum	44	5	5
Red Herring	24	0	9
Straw Men	11	0	2
Obfus.,Int. Vagueness,Confusion	10	0	1
Bandwagon	10	2	1
Total	5,114	380	547

Table 1: The statistics of all 18 propaganda techniques.



Benefits & Limitations

▪ Benefits

- Integrate neural networks with symbolic knowledge to make them more interpretable
- Support robust inference based on vector representations
- Support compositional problem solving

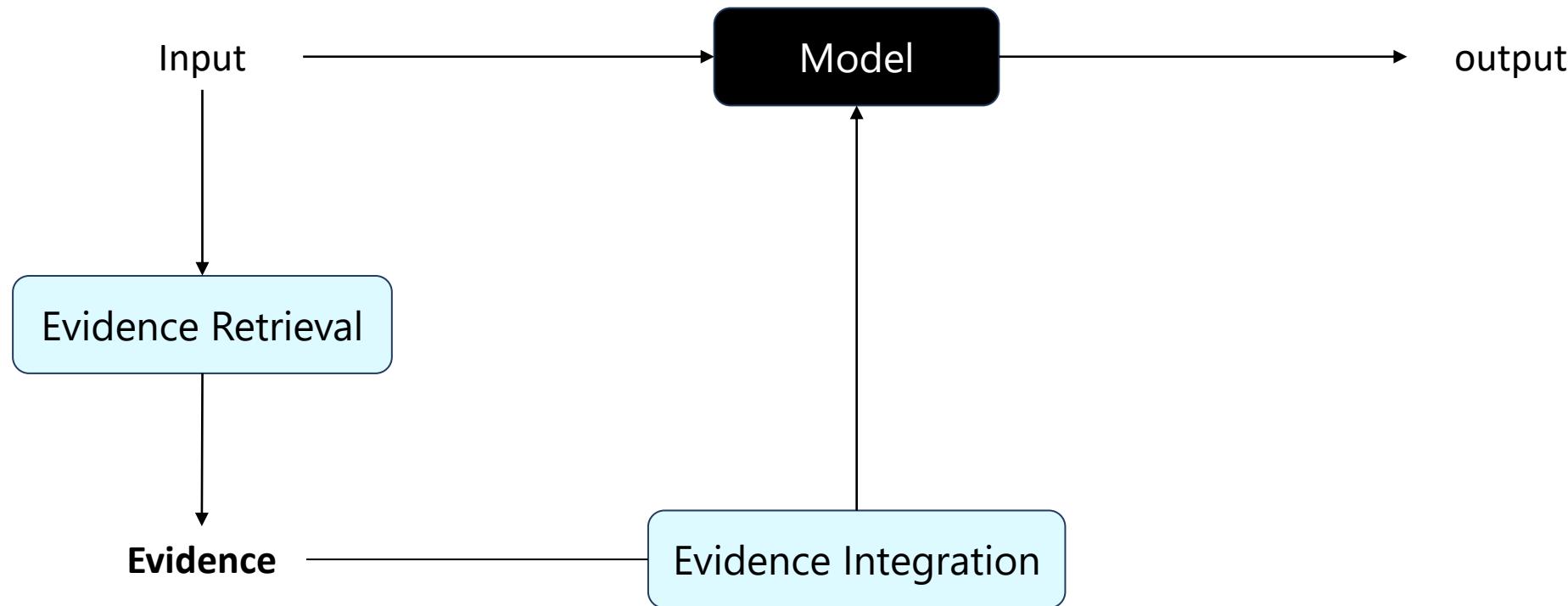
▪ Challenges

- Require large-scale labeled data for model training
- Require symbolic knowledge to have enough coverage on tasks of interest
- (Sometimes) Require hand-crafted rules or symbols, which are usually hard to scale

Neural-Evidence Reasoning

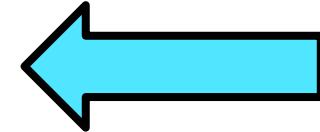
General Framework

- Consider evidence as an additional input of the model



Outline

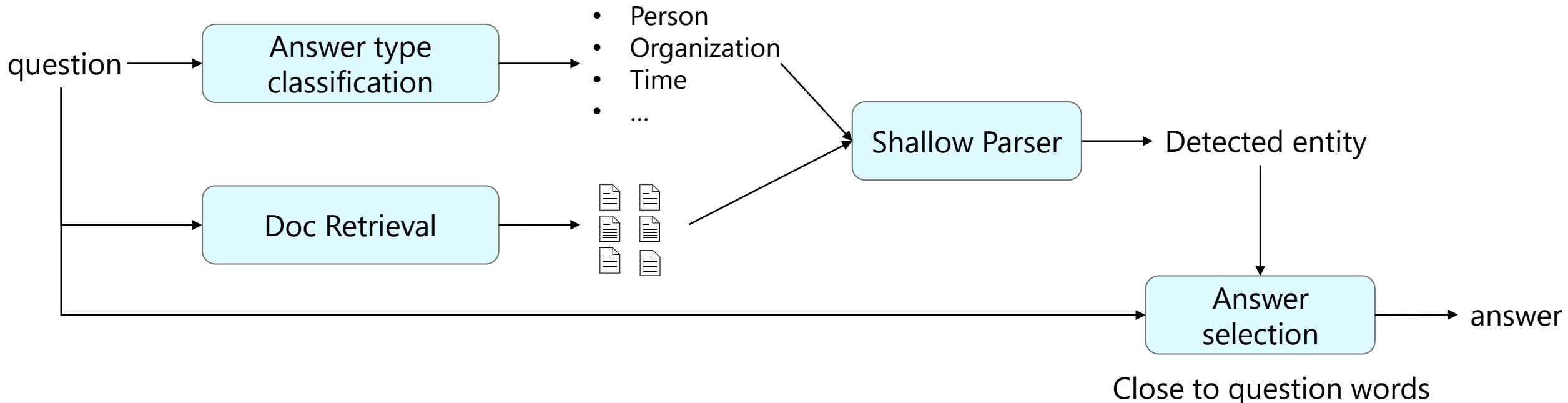
- **Text Evidence**
 - Applications: Open QA, Inferential Text Generation



- Fact Evidence
 - Applications: CommonsenseQA, Fact Checking
- Iterative Evidence
 - Multi-hop QA

Task #1: Open QA

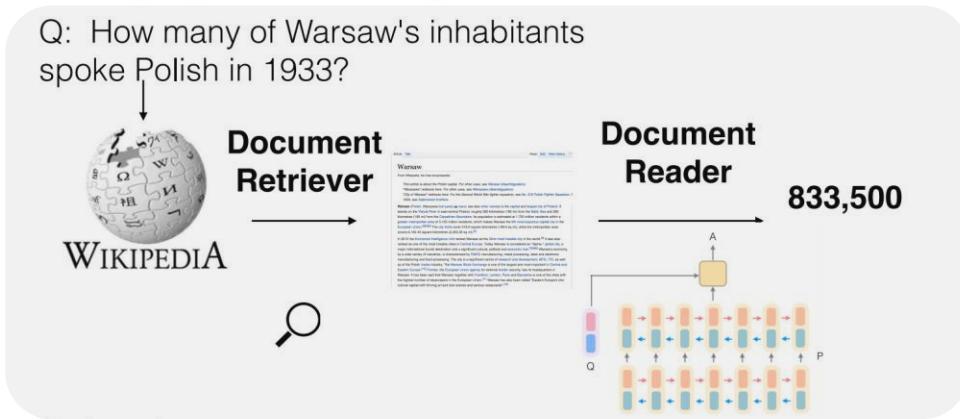
- First large-scale evaluation of domain-independent QA systems.
- Participants were given 200 fact-based, short-answer questions
- Each question was guaranteed to have at least one document in the collection that explicitly answered the question.
- Participants returned a ranked list of [document-id, answer-string] pairs per question such that each answer string was believed to contain an answer to the question.



Sparse Retrieval Model (DrQA)

Document Retriever + Document Reader

- Document retriever: finding relevant articles from 5 million Wikipedia articles
- Document reader (reading comprehension system): identifying the answer spans from those articles



- Datasets:
 - SQuAD ([Rajpurkar et al, 2016](#))
 - TREC ([Baudiš and Šedivý, 2005](#))
 - WebQuestions ([Berant et al, 2013](#))
 - WikiMovies ([Miller et al, 2016](#))

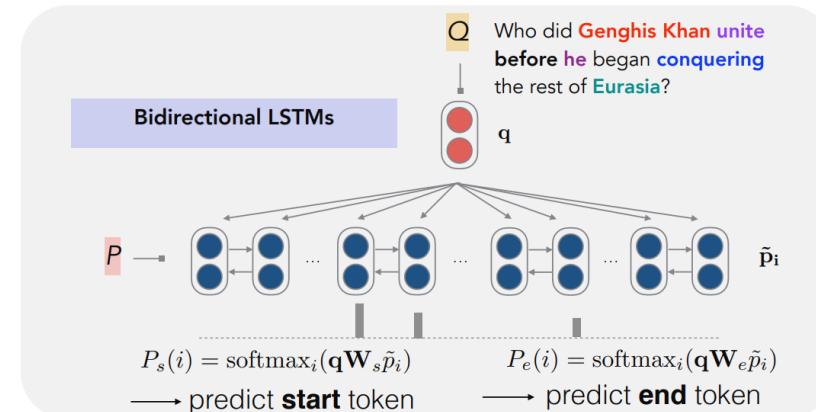
Document Retriever

TF-IDF bag-of-words vectors + efficient bigram hashing
(Weinberger et al., 2009)

Document Reader

Task: given paragraph P and question Q, the goal is to find a span A in the paragraph which answers the question.

Model: similar to AttentiveReader ([Hermann et al, 2015](#); [Chen et al, 2016](#)). We aim to keep it **simple!**



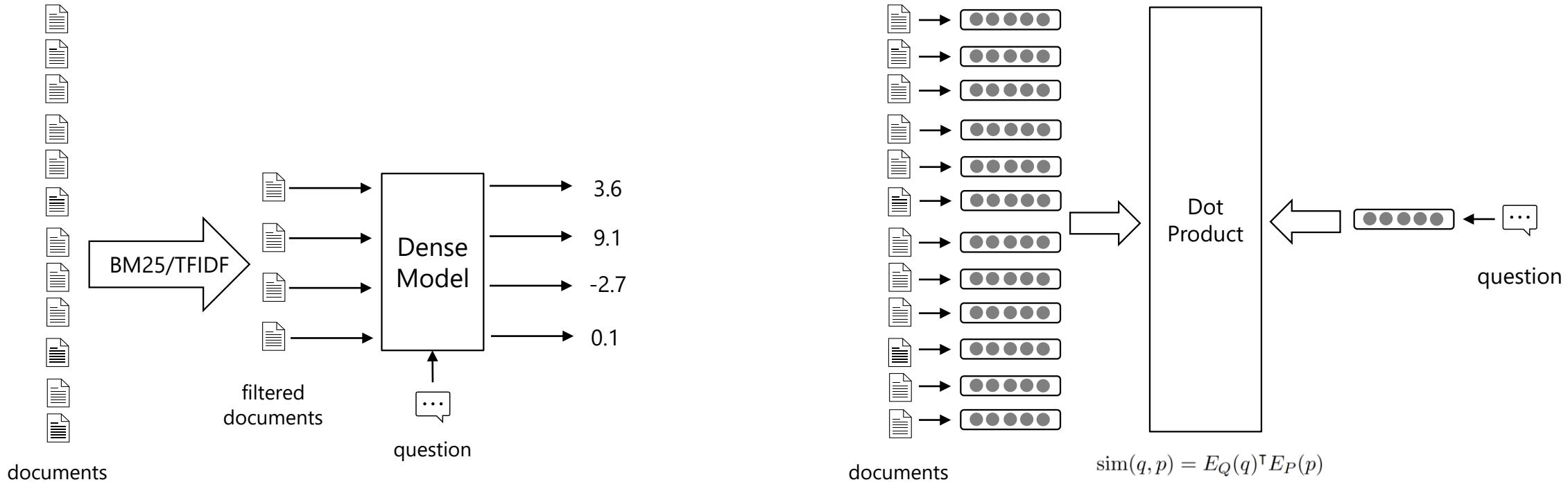
The input vectors consist of:

- Word embeddings
- Exact match features: whether the word appears in question
- Token features: POS, NER, term frequency
- Aligned question embedding

Data: SQuAD + **Distantly Supervised** Data

$(Q, A) \longrightarrow (P, Q, A)$ if P is retrieved and A can be found in P

Dense Retrieval Model (DPR, Dense Passage Retrieval)



$$L(q_i, p_i^+, p_{i,1}^-, \dots, p_{i,n}^-) = -\log \frac{e^{\text{sim}(q_i, p_i^+)}}{e^{\text{sim}(q_i, p_i^+)} + \sum_{j=1}^n e^{\text{sim}(q_i, p_{i,j}^-)}}$$

Who is the **bad guy** in lord of the rings?

Sala Baker is an actor and stuntman from New Zealand. He is best known for portraying the **villain** Sauron in the Lord of the Rings trilogy..

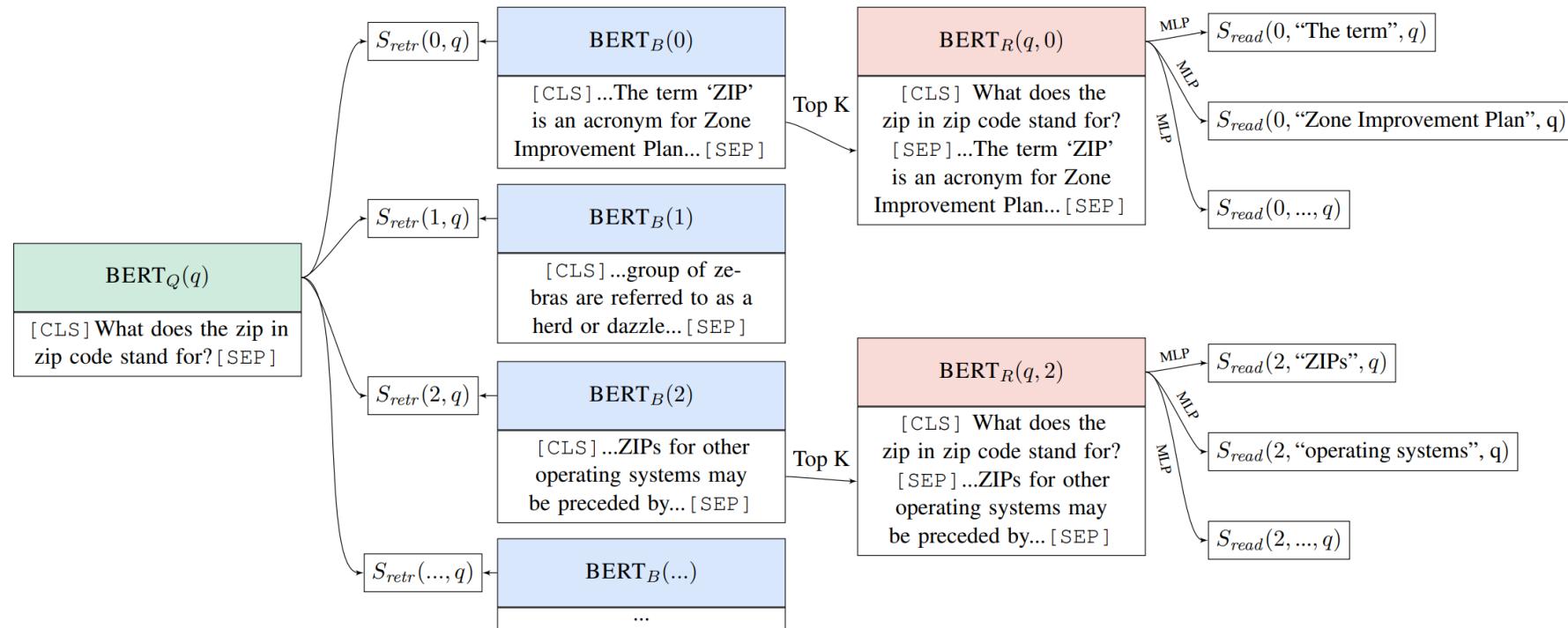
DPR Results

Training	Model	NQ	TriviaQA	WQ	TREC	SQuAD
Single	BM25+BERT (Lee et al., 2019)	26.5	47.1	17.7	21.3	33.2
Single	ORQA (Lee et al., 2019)	33.3	45.0	36.4	30.1	20.2
Single	HardEM (Min et al., 2019a)	28.1	50.9	-	-	-
Single	GraphRetriever (Min et al., 2019b)	34.5	56.0	36.4	-	-
Single	PathRetriever (Asai et al., 2020)	32.6	-	-	-	56.5
Single	REALM _{Wiki} (Guu et al., 2020)	39.2	-	40.2	46.8	-
Single	REALM _{News} (Guu et al., 2020)	40.4	-	40.7	42.9	-
Single	BM25	32.6	52.4	29.9	24.9	38.1
	DPR	41.5	56.8	34.6	25.9	29.8
	BM25+DPR	39.0	57.0	35.2	28.0	36.7
Multi	DPR	41.5	56.8	42.4	49.4	24.1
	BM25+DPR	38.8	57.9	41.1	50.6	35.8

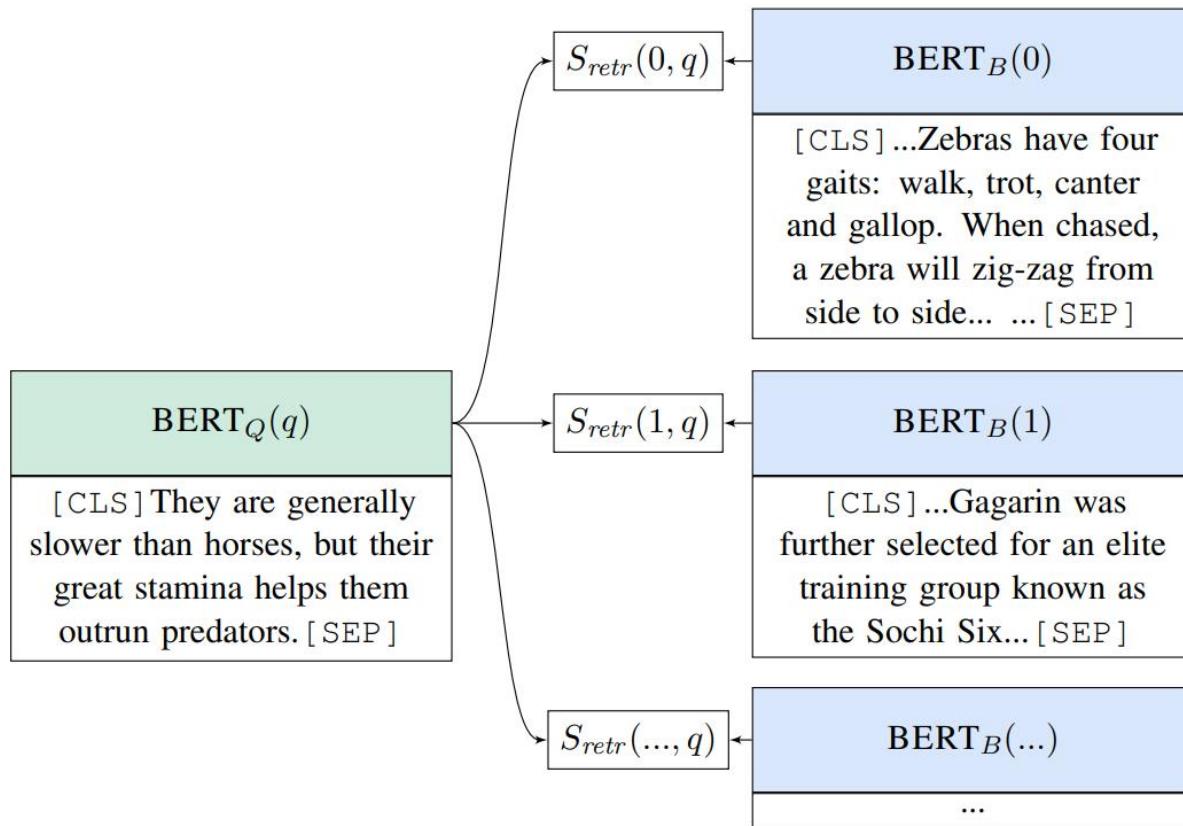
Table 4: End-to-end QA (Exact Match) Accuracy. The first block of results are copied from their cited papers. REALM_{Wiki} and REALM_{News} are the same model but pretrained on Wikipedia and CC-News, respectively. *Single* and *Multi* denote that our Dense Passage Retriever (DPR) is trained using individual or combined training datasets (all except SQuAD). For WQ and TREC in the *Multi* setting, we fine-tune the reader trained on NQ.

Joint Retrieval and Reader

- ORQA: Open-Retriever Question Answering
 - jointly learn the retriever and reader from question-answer string pairs
 - pre-train the retriever with an Inverse Cloze Task.



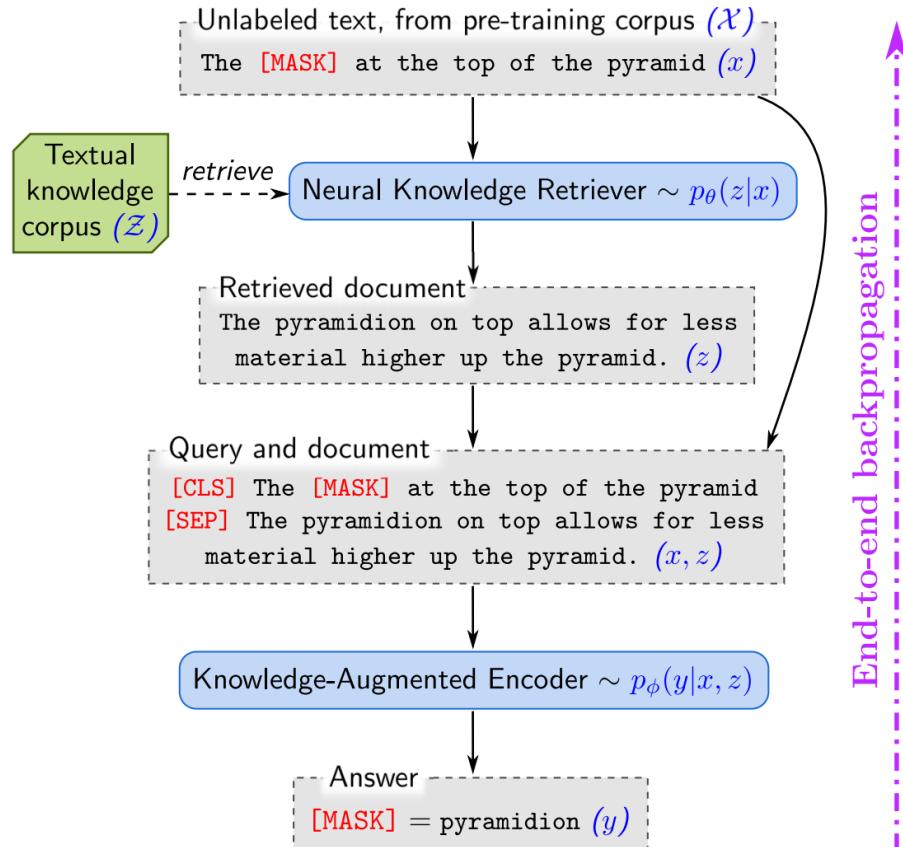
Pre-train Retrieval Model with Inverse Cloze Task



- In ICT, a sentence is treated as a pseudo-question, and its context is treated as pseudo-evidence.
- Given a pseudo-question, ICT requires selecting the corresponding pseudo-evidence out of the candidates in a batch.

Kenton Lee, Ming-Wei Chang, Kristina Toutanova.
"Latent Retrieval for Weakly Supervised Open Domain Question Answering." ACL-2019

Pre-train Retrieval Model with REALM



Knowledge Retriever: dense inner product model

$$p(z|x) = \frac{\exp f(x, z)}{\sum_{z'} \exp f(x, z')},$$
$$f(x, z) = \text{Embed}_{\text{input}}(x)^{\top} \text{Embed}_{\text{doc}}(z),$$

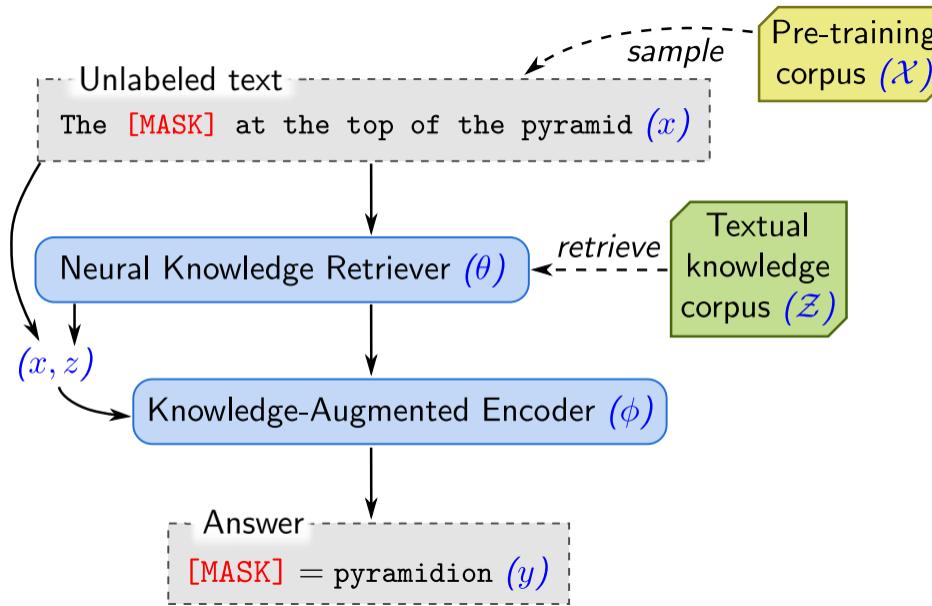
$$\text{Embed}_{\text{input}}(x) = \mathbf{W}_{\text{input}} \text{BERT}_{\text{CLS}}(\text{join}_{\text{BERT}}(x))$$

$$\text{Embed}_{\text{doc}}(z) = \mathbf{W}_{\text{doc}} \text{BERT}_{\text{CLS}}(\text{join}_{\text{BERT}}(z_{\text{title}}, z_{\text{body}}))$$

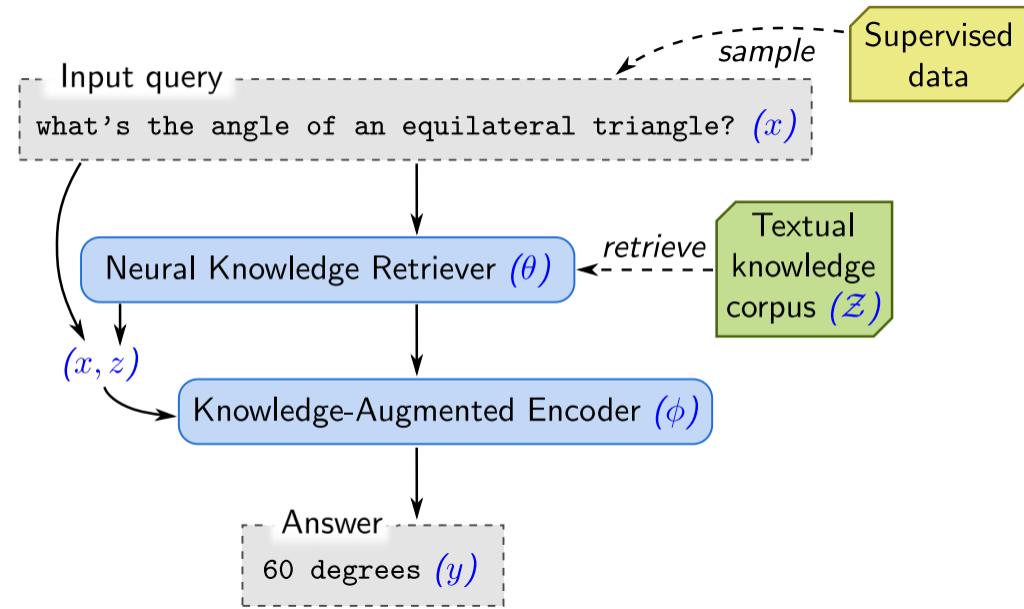
$$\text{join}_{\text{BERT}}(x) = [\text{CLS}] x [\text{SEP}]$$

$$\text{join}_{\text{BERT}}(x_1, x_2) = [\text{CLS}] x_1 [\text{SEP}] x_2 [\text{SEP}]$$

Pre-train Retrieval Model with REALM



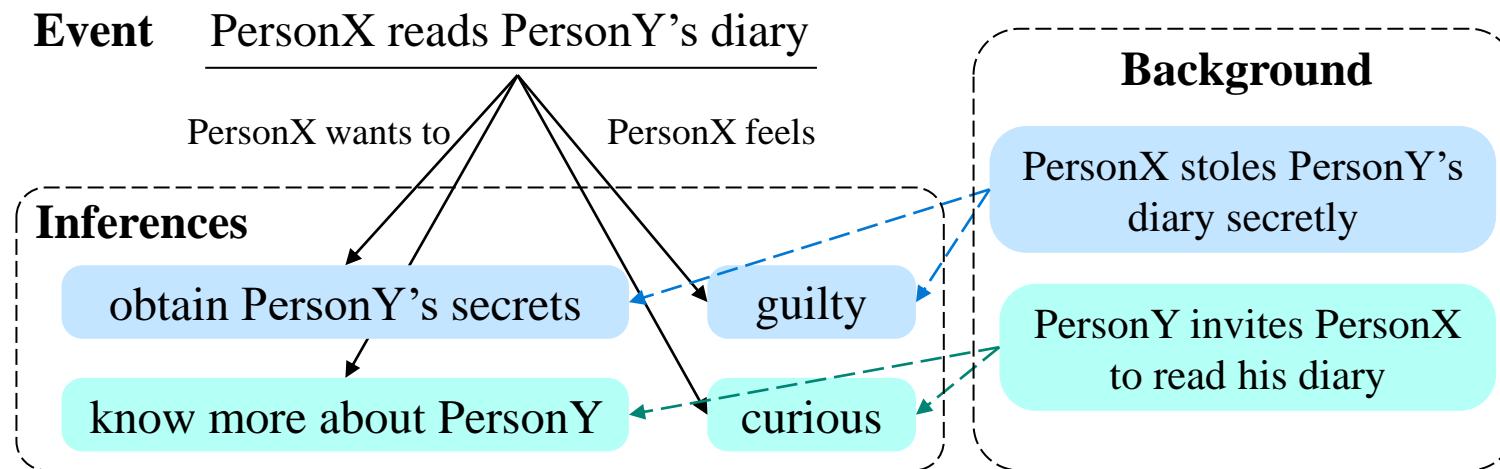
Unsupervised pre-training



Supervised fine-tuning

Task #2: Inferential Text Generation

- ATOMIC & Event2Mind

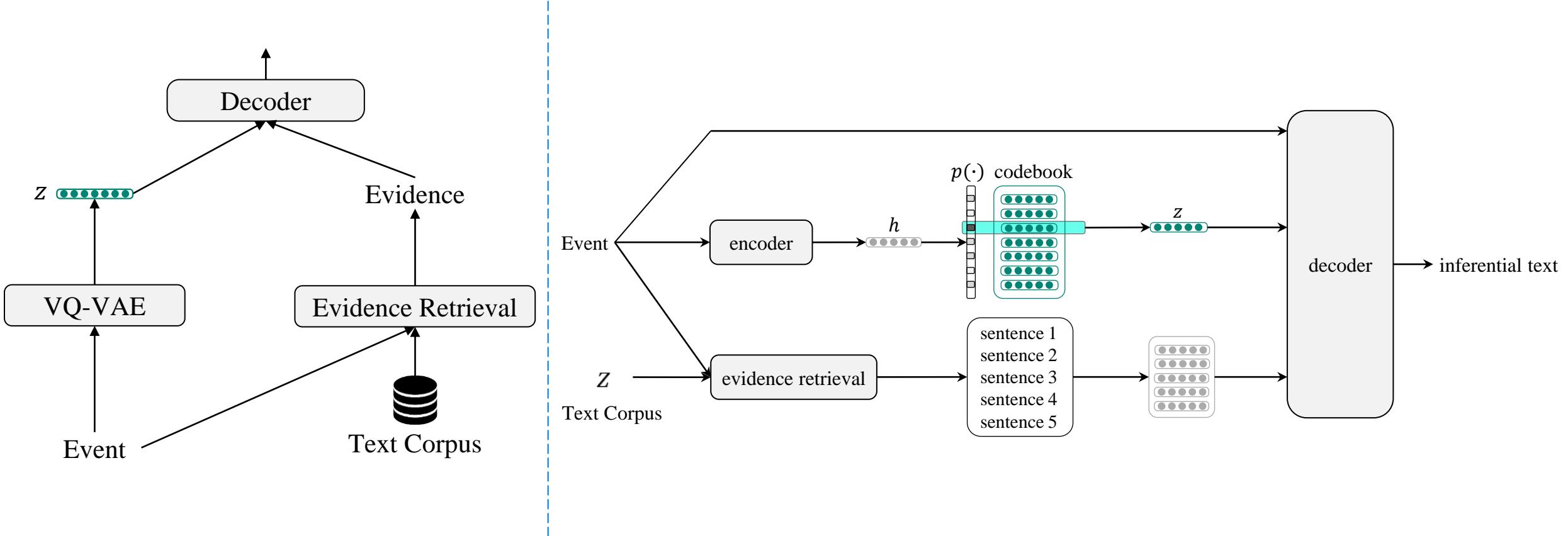


Maarten Sap, Ronan LeBras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith and Yejin Choi .

"**ATOMIC: An Atlas of Machine Commonsense for If-Then Reasoning .**" AAAI-2019.

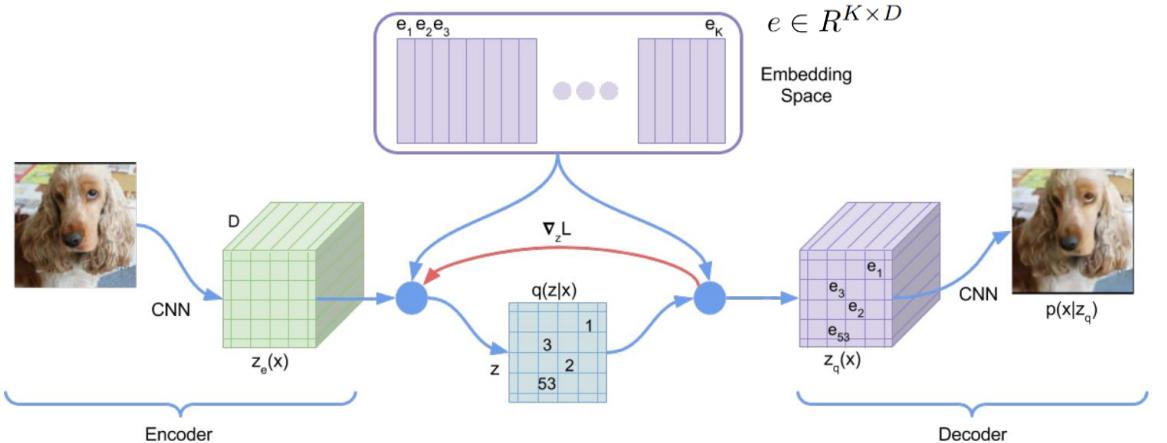
Hannah Rashkin, Maarten Sap, Emily Allaway, Noah A. Smith and Yejin Choi. "**Event2Mind: Commonsense Inference on Events, Intents, and Reactions**". ACL-2018

Evidence-Aware VQ-VAE (EA-VQ-VAE)



BM25 and select top K sentences from BooksCorpus as the evidence.

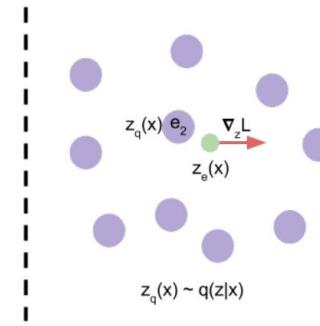
VQ-VAE (Vector Quantised Variational AutoEncoder)



$$q(z = k|x) = \begin{cases} 1 & \text{for } k = \operatorname{argmin}_j \|z_e(x) - e_j\|_2, \\ 0 & \text{otherwise} \end{cases}$$

$$z_q(x) = e_k, \quad \text{where } k = \operatorname{argmin}_j \|z_e(x) - e_j\|_2$$

- One of the simplest dictionary learning algorithms, Vector Quantisation (VQ).
- The VQ objective uses the L2 error to move the embedding vectors e_i towards the encoder outputs $z_e(x)$ as shown in the **second term** of equation 3.



- The gradient $\nabla_z L$ (in red) will push the encoder to change its output, which could alter the configuration in the next forward pass.
- Since the output representation of the encoder and the input to the decoder share the same D dimensional space, the gradients contain useful information for how the encoder has to change its output to lower the reconstruction loss.

$$L = \log p(x|z_q(x)) + \|\operatorname{sg}[z_e(x)] - e\|_2^2 + \beta \|z_e(x) - \operatorname{sg}[e]\|_2^2$$

- The decoder optimises the first loss term only, the encoder optimises the first and the last loss terms, and the embeddings are optimised by the middle loss term.

add a **commitment loss** to make sure the encoder commits to an embedding and its output does not grow

State-of-the-art on ATOMIC

Methods	xIntent	xNeed	xAttr	xEffect	xReact	xWant	oEffect	oReact	oWant	Overall
Single Task										
S2S	8.17	12.35	2.96	5.26	3.43	13.44	6.42	4.09	7.08	7.02
VRNMT	9.52	13.35	4.87	4.42	7.64	9.80	13.71	5.28	10.79	8.82
CWVAE	12.12	15.67	5.63	14.64	8.13	15.01	11.63	8.58	13.83	11.69
Multi Task										
S2S*	24.53	23.85	5.06	9.44	5.38	24.68	7.93	5.60	21.30	14.20
COMET*	25.82	25.54	5.39	10.39	5.36	26.41	8.43	5.65	21.96	15.00
COMET	-	-	-	-	-	-	-	-	-	15.10
EA-VQ-VAE	26.89	25.95	5.72	10.96	5.68	25.94	8.78	6.10	22.48	15.40

- Avg BLEU score with multiple references
- for example, 10 predictions with 3 references. First, for each prediction, calculate the BLEU score with each reference separately, and choose the highest one. Then, Average the 10 BLEU scores.

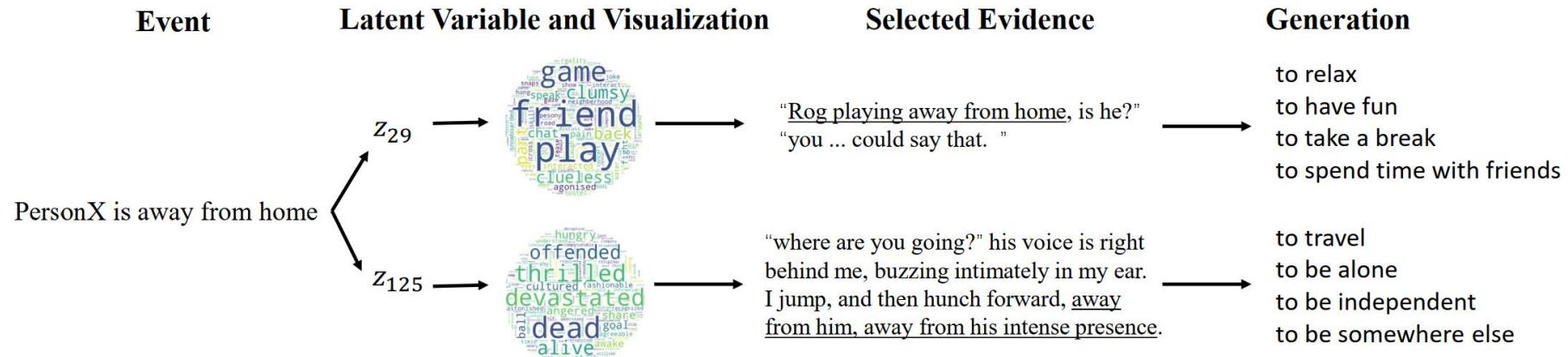


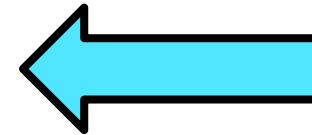
Figure 4: An examples of Event2Mind dataset on the xIntent dimension (i.e. “*PersonX wants*”).

Outline

- Text Evidence
 - Applications: Open QA, Inferential Text Generation

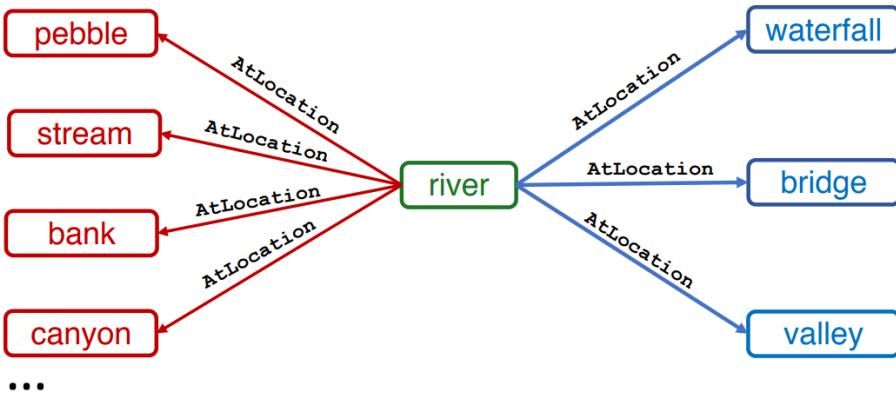
- Fact Evidence
 - Applications: CommensenseQA, Fact Checking

- Iterative Evidence
 - Multi-hop QA



Task #1: CommonsenseQA

- a) Sample ConceptNet for specific subgraphs



- A source concept (in green) and three target concepts (in blue) are sampled from CONCEPTNET
- Crowd-workers generate three questions, each having one of the target concepts for its answer (✓), while the other two targets are not (X). Then, for each question, workers choose an additional distractor from CONCEPTNET (in red), and author one themselves (in purple).

- b) Crowd source corresponding natural language questions and two additional distractors

Where on a **river** can you hold a cup upright to catch water on a sunny day?

✓ **waterfall**, X **bridge**, X **valley**, X **pebble**, X **mountain**

Where can I stand on a **river** to see water falling without getting wet?

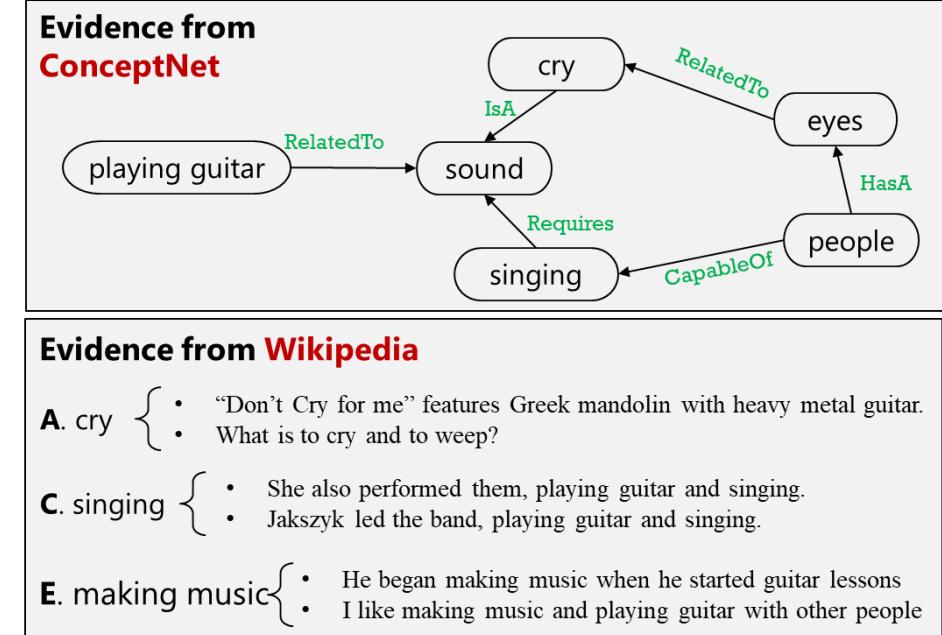
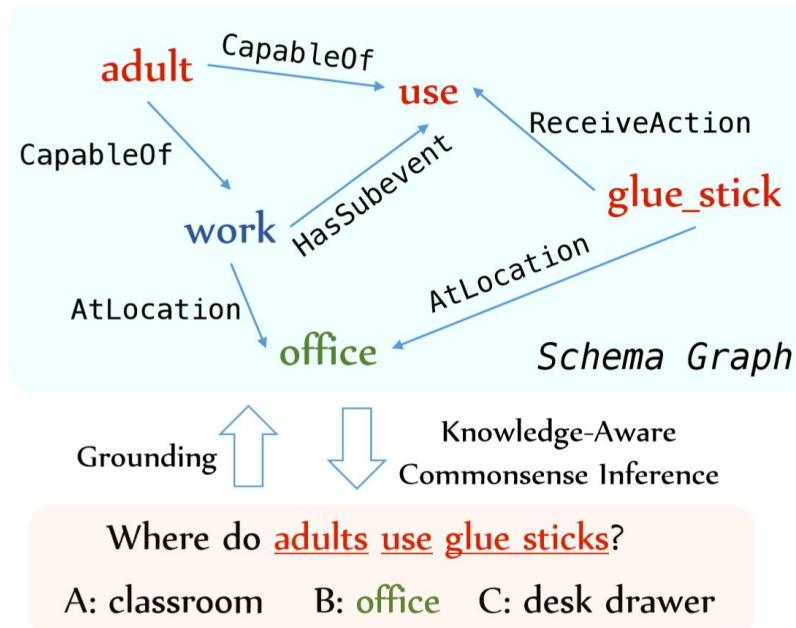
X **waterfall**, ✓ **bridge**, X **valley**, X **stream**, X **bottom**

I'm crossing the **river**, my feet are wet but my body is dry, where am I?

X **waterfall**, X **bridge**, ✓ **valley**, X **bank**, X **island**

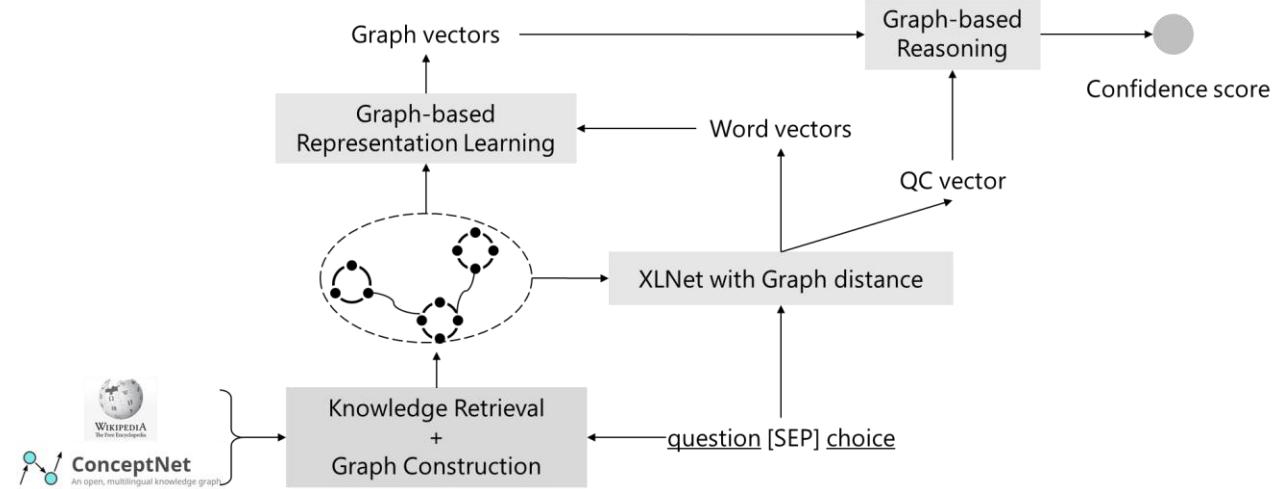
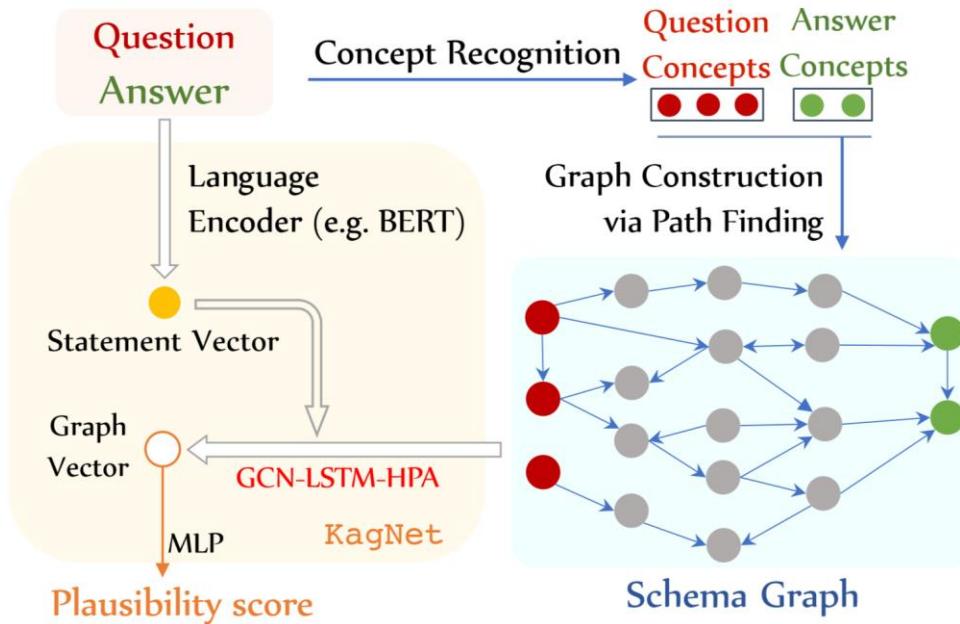
Evidence Retrieval

- Evidence from ConceptNet
 - Match ConceptNet Vocab
 - Find paths between QA-concept pair
 - Path pruning by length (≤ 5 nodes) and embedding-based metric
- Evidence from both ConceptNet (< 3 hops) and Wikipedia (via elasticsearch)



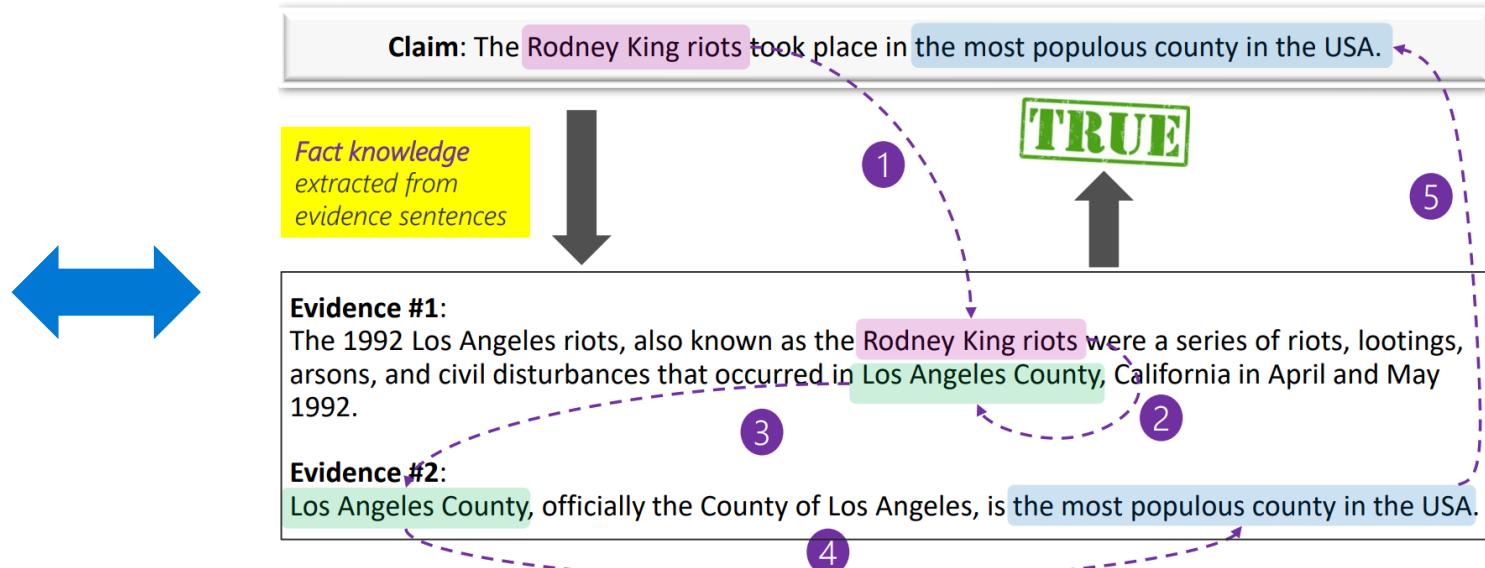
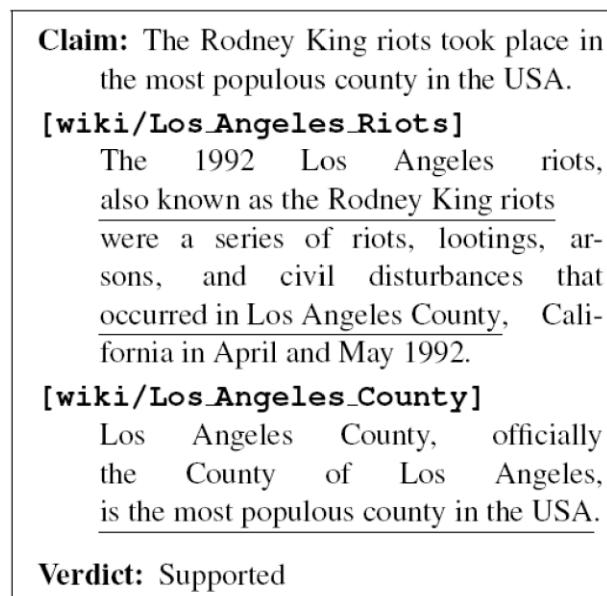
Graph Reasoning

- BERT + Graph Neural Net
- XLNet + Graph Neural Net

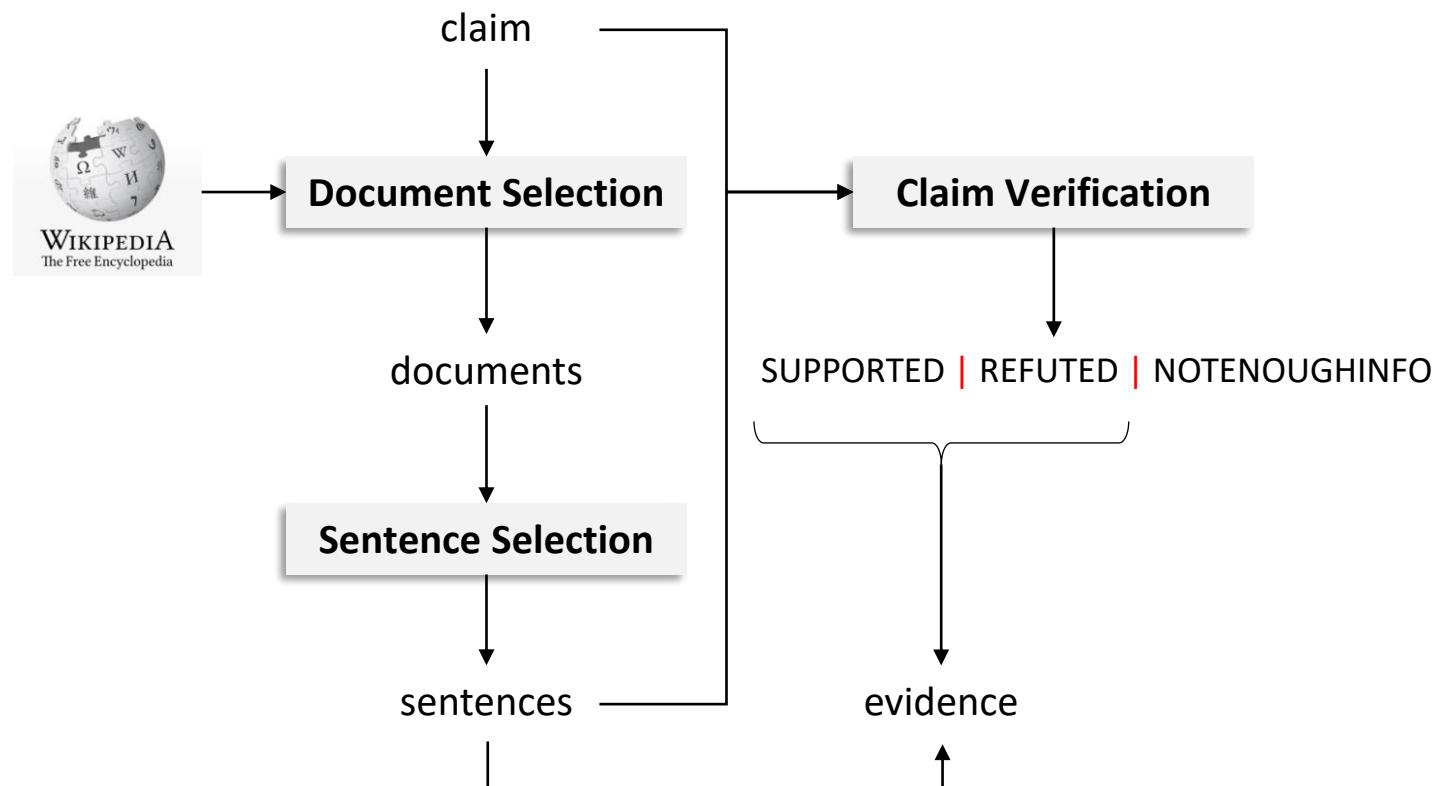


Task #2: FEVER (Fact Extraction and VERification)

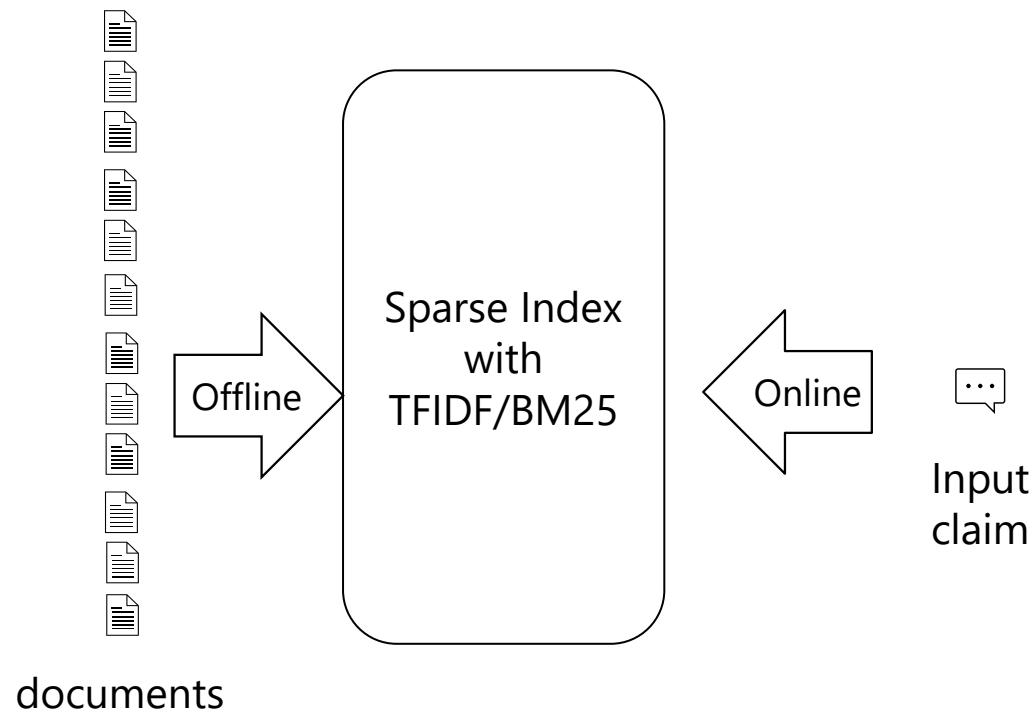
- Input: A claim + A collection of Wikipedia documents
- Output:
 - Classify the given claim as SUPPORTED, REFUTED or NOTENOUGHINFO
 - For the first two classes, provide the supporting sentences as evidence



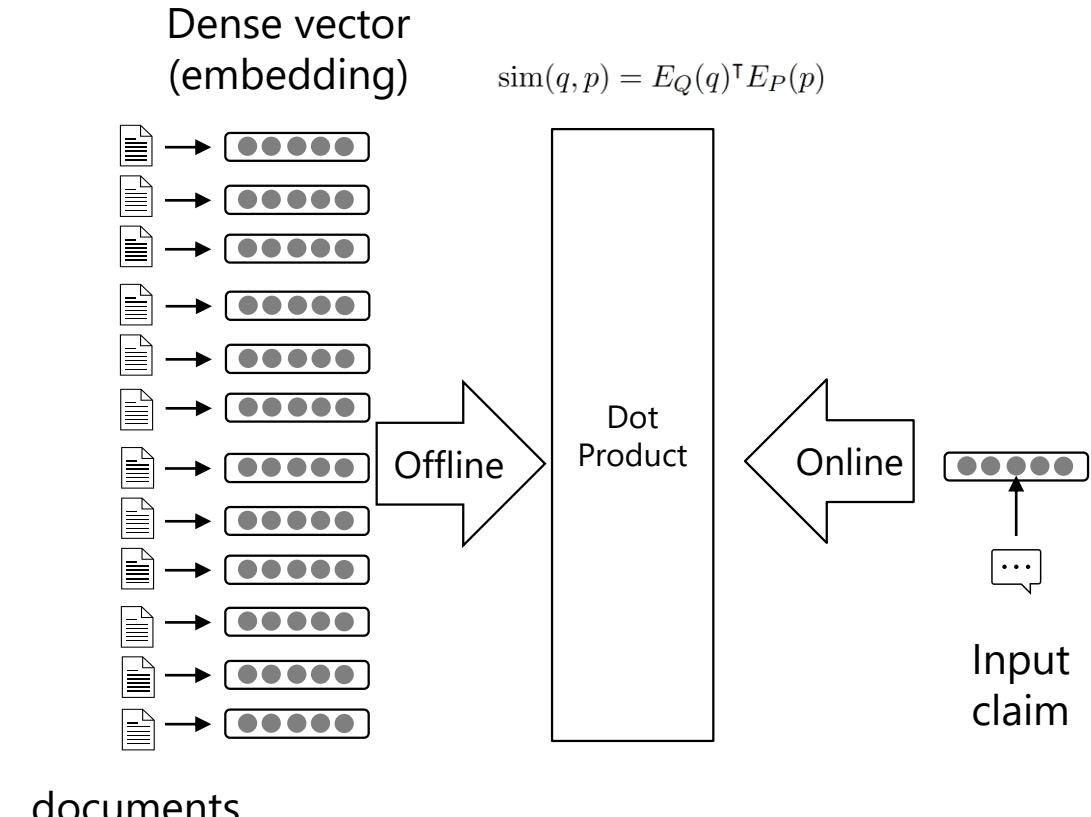
Pipeline



Document Retrieval



Sparse Retrieval Models

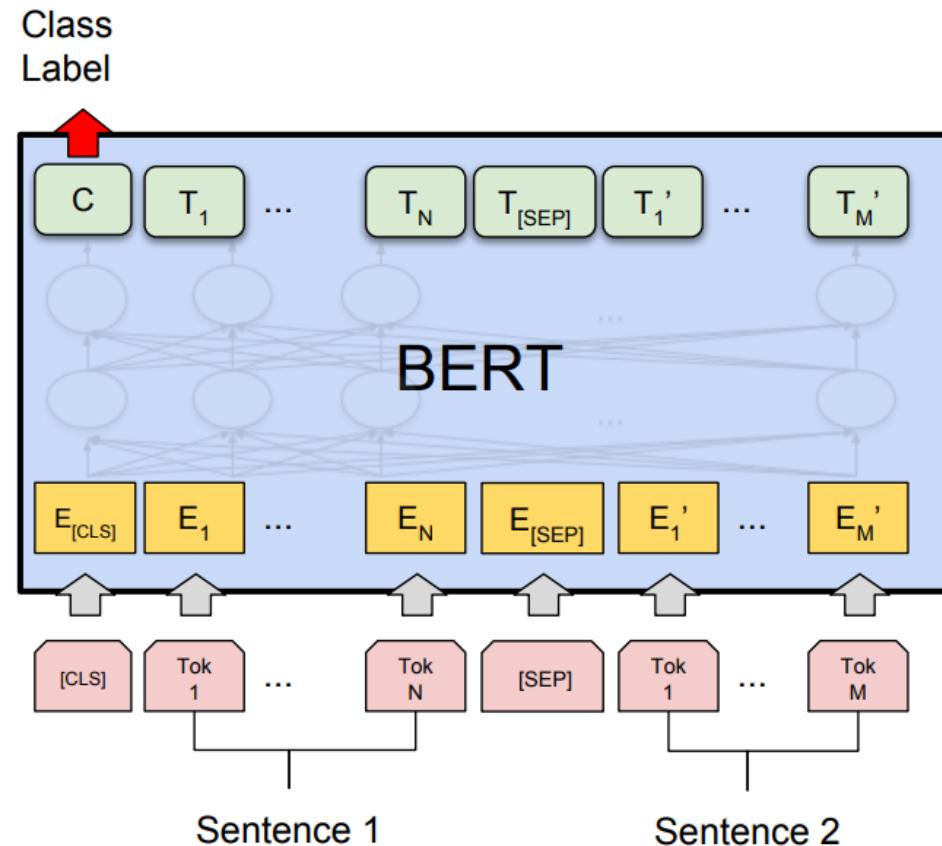


Dense Retrieval Models

Evidence Selection

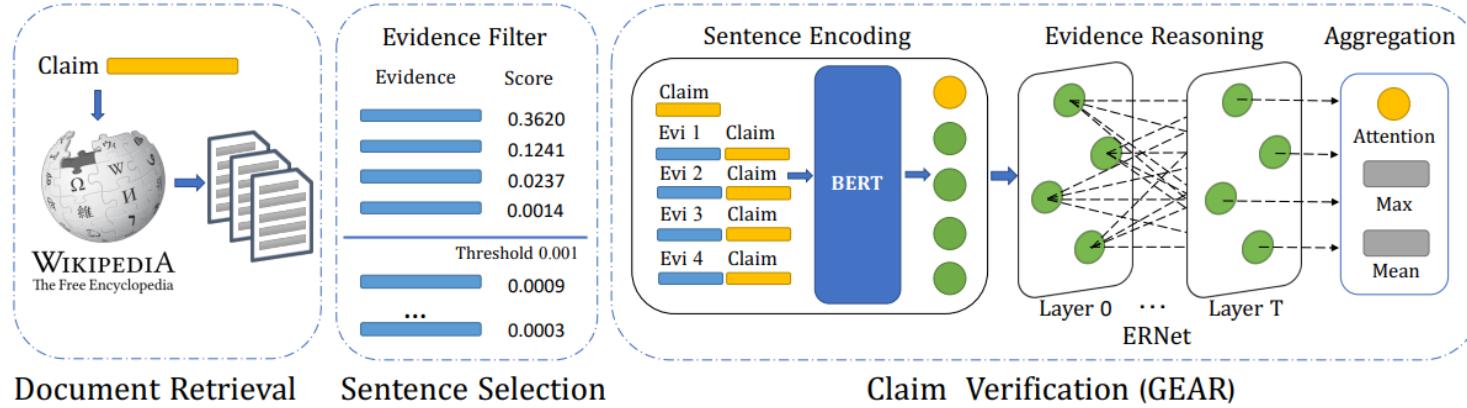
- Natural Language Inference

Model	Matched	Mismatched	Paper / Source	Code
RoBERTa (Liu et al., 2019)	90.8	90.2	RoBERTa: A Robustly Optimized BERT Pretraining Approach	Official
XLNet-Large (ensemble) (Yang et al., 2019)	90.2	89.8	XLNet: Generalized Autoregressive Pretraining for Language Understanding	Official
MT-DNN-ensemble (Liu et al., 2019)	87.9	87.4	Improving Multi-Task Deep Neural Networks via Knowledge Distillation for Natural Language Understanding	Official
Snorkel MeTaL(ensemble) (Ratner et al., 2018)	87.6	87.2	Training Complex Models with Multi-Task Weak Supervision	Official
Finetuned Transformer LM (Radford et al., 2018)	82.1	81.4	Improving Language Understanding by Generative Pre-Training	
Multi-task BiLSTM + Attn (Wang et al., 2018)	72.2	72.1	GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding	
GenSen (Subramanian et al., 2018)	71.4	71.3	Learning General Purpose Distributed Sentence Representations via Large Scale Multi-task Learning	



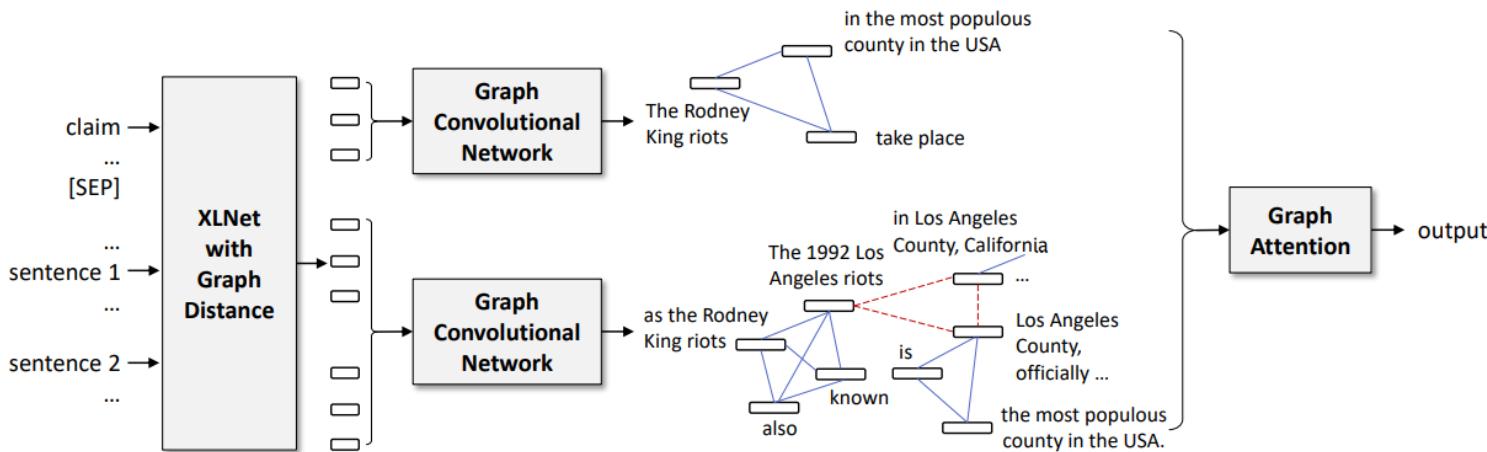
Jacob Devlin, Ming-Wei Chang, Kent Lee, Kristina Toutanova. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding." NAACL-2019.

Claim Verification



- Evidence aggregation with graph neural network, where each node represents a claim or claim/evidence pair.

Jie Zhou, Xu Han, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, Maosong Sun. "GEAR: Graph-based Evidence Aggregating and Reasoning for Fact Verification." ACL-2019.



- Evidence aggregation with graph neural network, where graph is constructed with fine-grained info extracted with SRL.

Wanjun Zhong, Jingjing Xu, Duyu Tang, Zenan Xu, Nan Duan, Ming Zhou, Jiahai Wang, Jian Yin. "Reasoning Over Semantic-Level Graph for Fact Checking" ACL-2020.

Outline

- Text Evidence
 - Applications: Open QA, Inferential Text Generation
- Fact Evidence
 - Applications: CommonsenseQA, Fact Checking
- Iterative Evidence
 - Multi-hop QA



Task: HotpotQA

- Multi-hop Reasoning across Multiple Documents

SQuAD dataset

*When was the **Millwall Football Club** founded?*



HotpotQA dataset

*When was the football club founded **in which Walter Otto Davis played at centre forward**?*

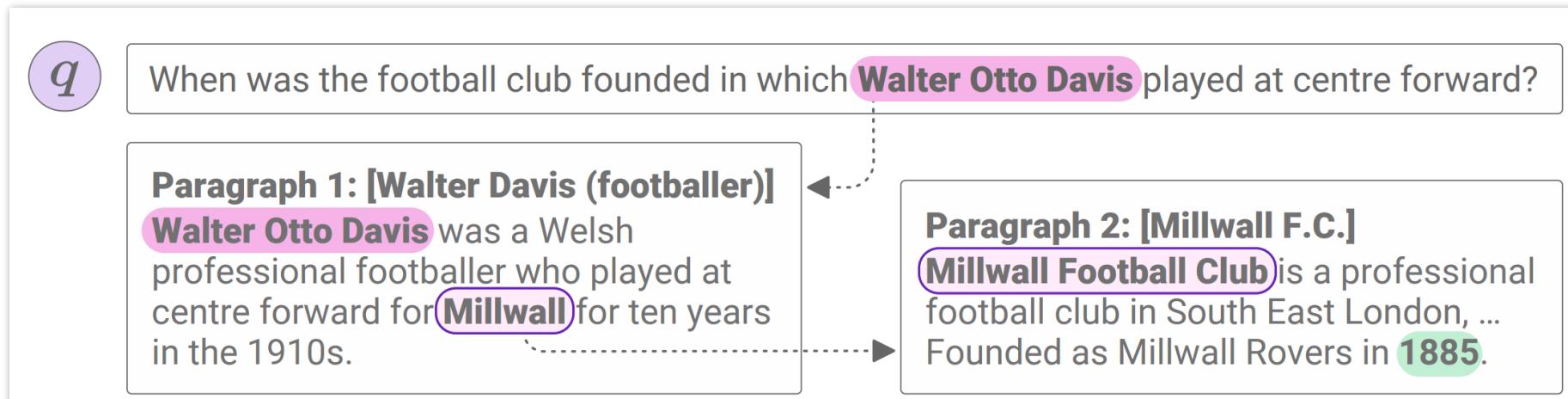


Image courtesy [Asai, et al. 2020]

Learning to Retrieve Reasoning Paths

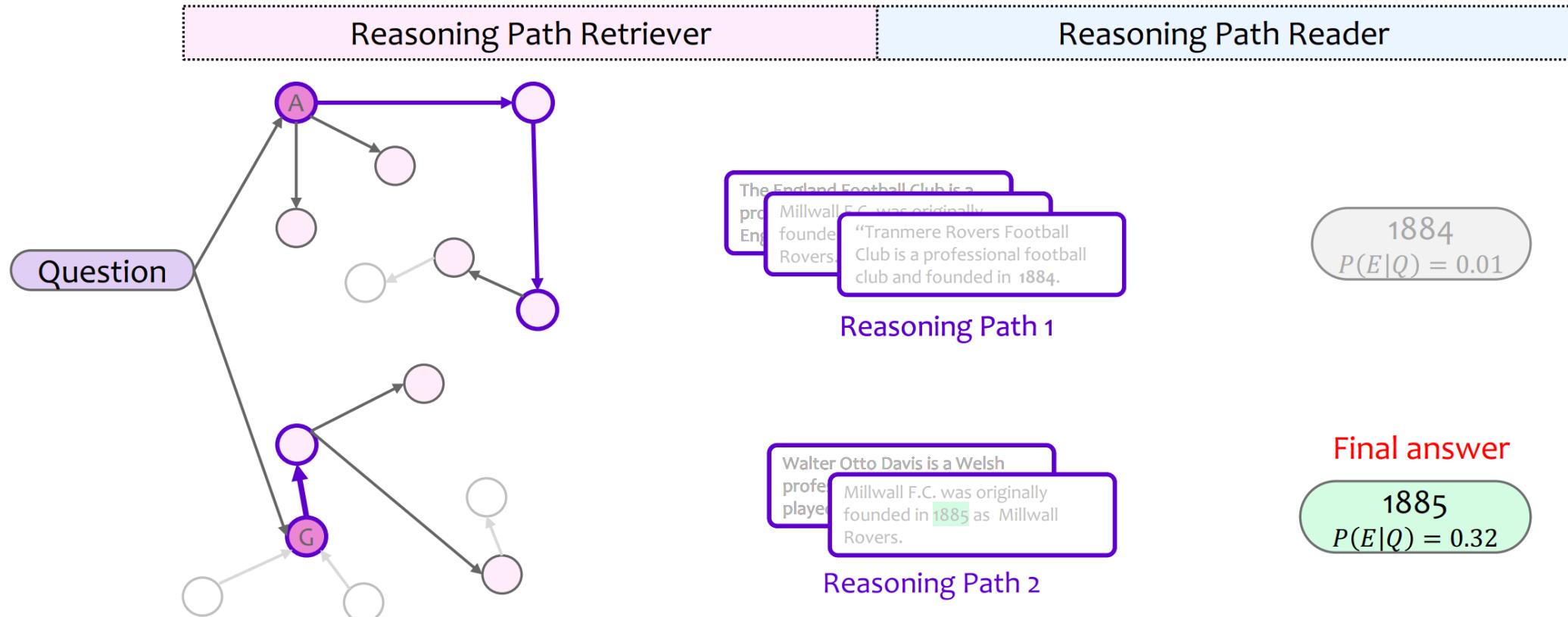
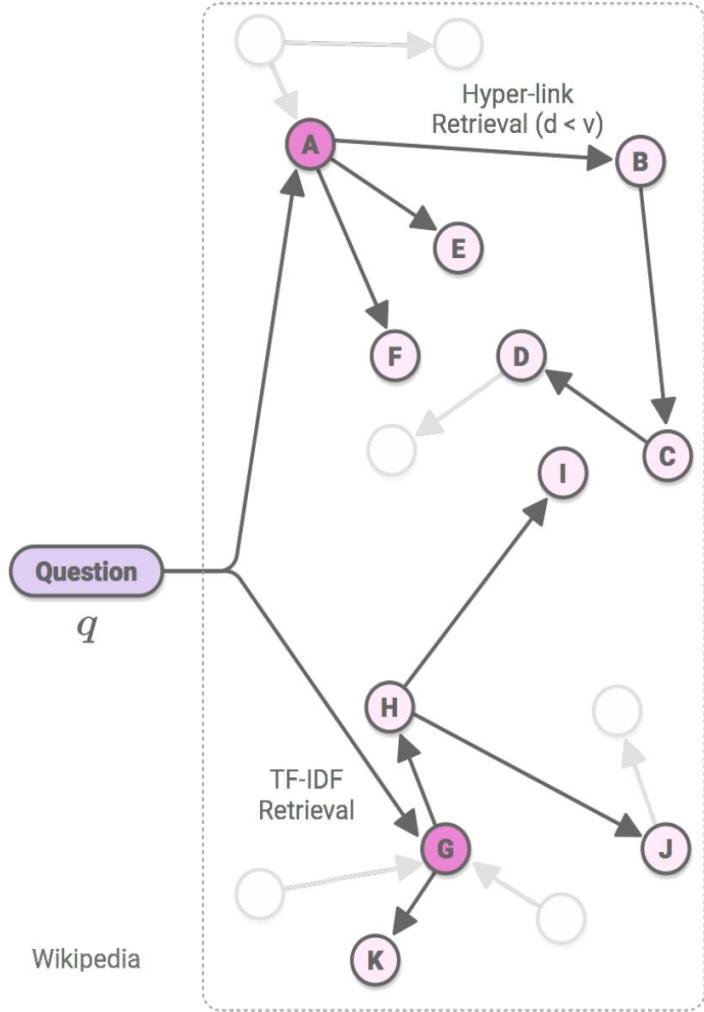
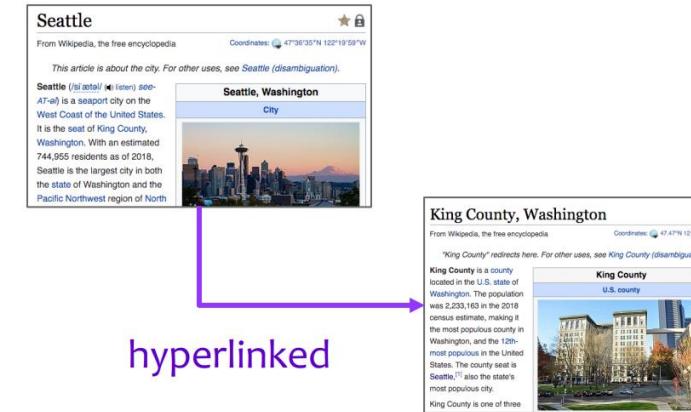


Image courtesy [Asai, et al. 2020]

Graph Construction



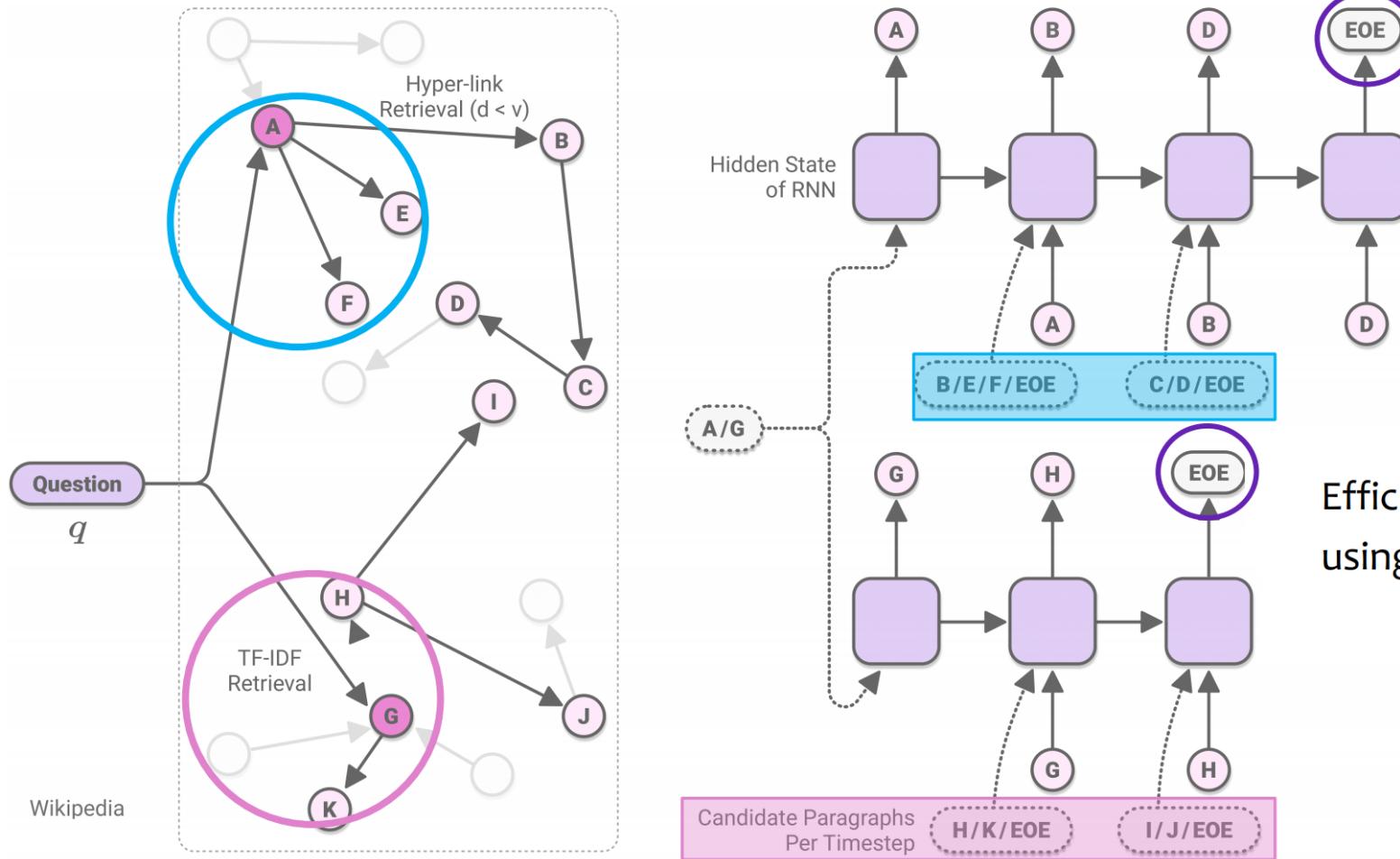
The Wikipedia Graph covers all articles.



- Each article as a node of the graph.
 - TF-IDF top articles (represented by filled pink circles)
 - Hyper-linked articles (represented by white circles with black outlines)
- Each hyperlink as an edge of the graph.

Image courtesy [Asai, et al. 2020]

Evidence Retrieval Model

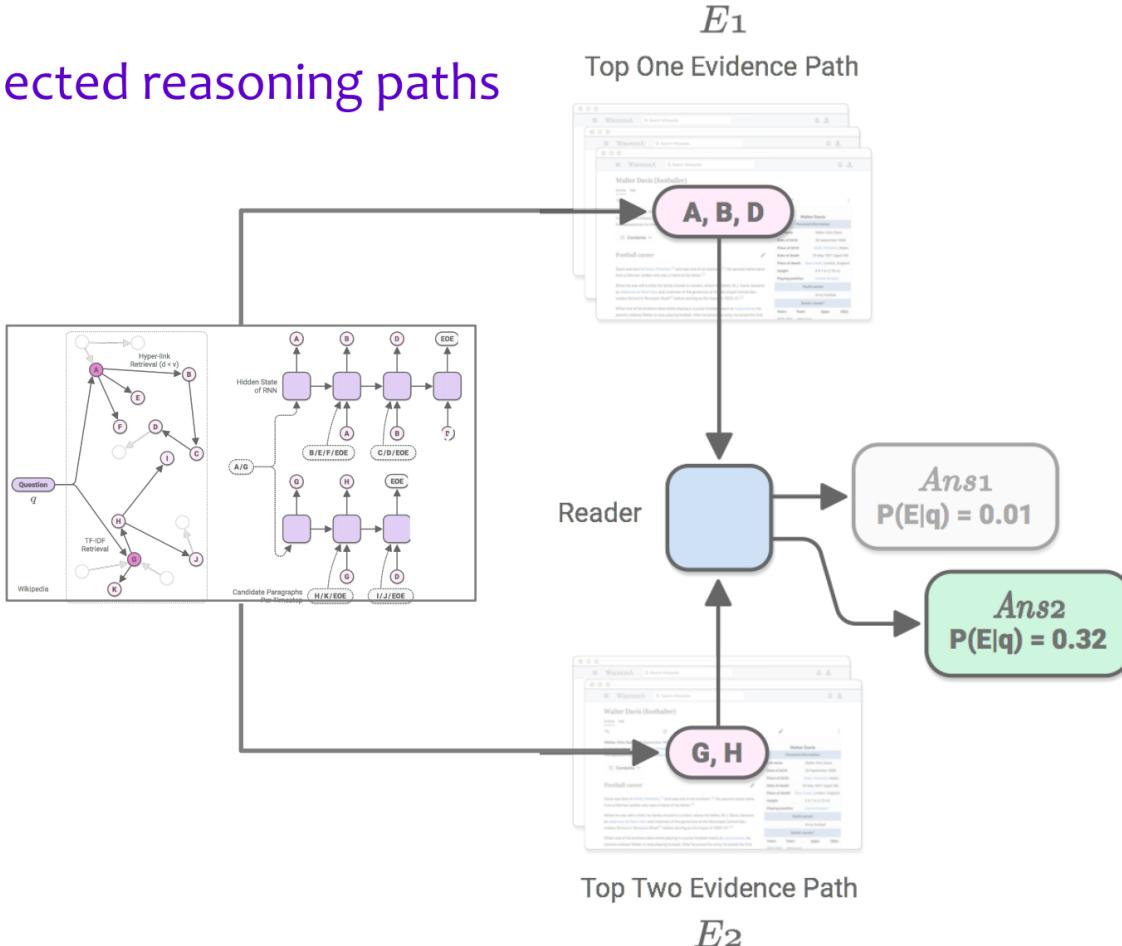


Efficiently finding top k paths
using beam search

Image courtesy [Asai, et al. 2020]

Reader

Selected reasoning paths



1. extracts answers from reasoning paths
2. estimates scores of the selected reasoning paths.

Final answer: the best span extracted from the best reasoning path

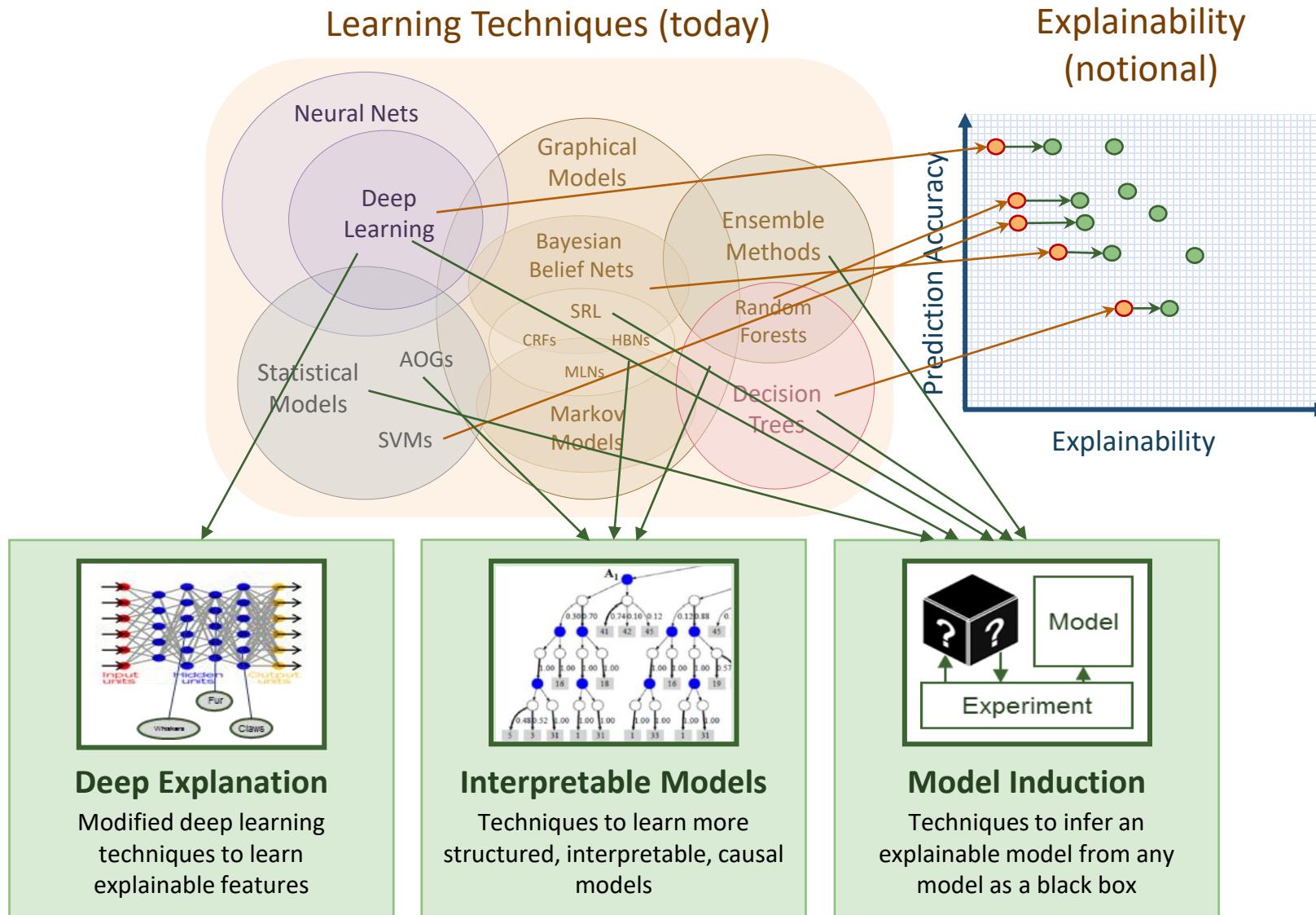
Image courtesy [Asai, et al. 2020]

Benefits and Challenges

- Benefits
 - Easy to be combined with neural models (e.g. pre-trained models), good performance
 - Broad applications as the definition of evidence can be broad
- Challenges
 - Evidence retrieval is a fundamental challenge
 - Interpretability is limited by the retrieved evidence, instead of the inference process

Dilemma: Interpretability vs. Performance

Performance vs. Interpretability



Explainable Artificial Intelligence (XAI),
David Gunning, DARPA/I2O

Empirical Success of Pre-trained Model (1/4): Grade 8 New York Regents Science Exam

1. Which equipment will best separate a mixture of iron filings and black pepper? (1) magnet (2) filter paper (3) triplebeam balance (4) voltmeter
2. Which form of energy is produced when a rubber band vibrates? (1) chemical (2) light (3) electrical (4) sound
3. Because copper is a metal, it is (1) liquid at room temperature (2) nonreactive with other substances (3) a poor conductor of electricity (4) a good conductor of heat
4. Which process in an apple tree primarily results from cell division? (1) growth (2) photosynthesis (3) gas exchange (4) waste removal

Example questions from the NY Regents Exam (8th Grade), illustrating the need for both scientific and commonsense knowledge.

• AristoBERT

- **Background Knowledge:** use up to 10 of the top sentences found by the IR solver, truncated to fit into the BERT max tokens setting (256).
- **Curriculum Fine-Tuning:** first fine-tune on the RACE training set, and then fine-tune on current corpus

Test Set	Num Q	IR	PMI	ACME	TupInf	Multee	AristoBERT	AristoRoBERTa	ARISTO
Regents 4th	109	64.45	66.28	67.89	63.53	69.72	86.24	88.07	89.91
Regents 8th	119	66.60	69.12	67.65	61.41	68.91	86.55	88.24	91.60
Regents 12th	632	41.22	46.95	41.57	35.35	56.01	75.47	82.28	83.54
ARC-Easy	2376	74.48	77.76	66.60	57.73	64.69	81.78	82.88	86.99
ARC-Challenge	1172	n/a [†]	n/a [†]	20.44	23.73	37.36	57.59	64.59	64.33

Empirical Success of Pre-trained Model (2/4): Discrete Reasoning over Natural Language

Reasoning	Passage (some parts shortened)	Question	Answer
Subtraction (28.8%)	That year, his Untitled (1981) , a painting of a haloed, black-headed man with a bright red skeletal body, depicted amid the artists signature scrawls, was sold by Robert Lehrman for \$16.3 million, well above its \$12 million high estimate.	How many more dollars was the Untitled (1981) painting sold for than the 12 million dollar estimation?	4300000
Comparison (18.2%)	In 1517 , the seventeen-year-old King sailed to Castile. There, his Flemish court In May 1518 , Charles traveled to Barcelona in Aragon.	Where did Charles travel to first, Castile or Barcelona?	Castile

QDGAT is based on QANet and ALBERT model.

AI2 Allen Institute for AI

 Leaderboard

DROP

DROP is a QA dataset which tests comprehensive understanding of paragraphs. In this crowdsourced, adversarially-created, 96k question-answering benchmark, a system must resolve multiple references in a question, map them onto a paragraph, and... [more](#)

[Public Submissions](#) [Getting Started](#) [About](#)

 Human Performance

Rank	Submission	Created	F1
1	QDGAT - ALBERT AntGroup KG & NLP	09/08/2020	0.9010
2	Numeric Transformer - Albert OneConnect GammaLab NYC	03/17/2020	0.8911
3	QDGAT Ensemble AntGroup KG & NLP	12/16/2019	0.8838
4	sna_albert+ Ensemble OneConnect GammaLab	12/03/2019	0.8795
5	QDGAT - RoBERTa AntGroup KG & NLP	06/01/2020	0.8779
6	Numeric Transformer - RoBERTa OneConnect GammaLab NYC	03/03/2020	0.8759
	QDGAT		

Empirical Success of Pre-trained Model (3/4): Reason over Rules in Natural Language

(Input Facts:) Alan is blue. Alan is rough. Alan is young.

Bob is big. Bob is round.

Charlie is big. Charlie is blue. Charlie is green.

Dave is green. Dave is rough.

(Input Rules:) Big people are rough.

If someone is young and round then they are kind.

If someone is round and big then they are blue.

All rough people are green.

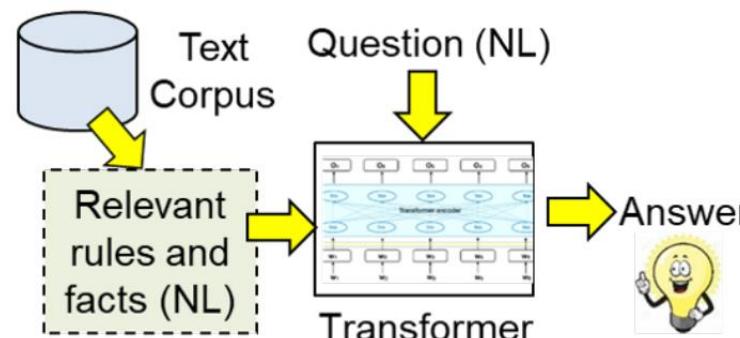
Q1: Bob is green. True/false? [Answer: T]

Q2: Bob is kind. True/false? [F]

Q3: Dave is blue. True/false? [F]

Training	Num Q	Mod0 $D = 0$	Mod1 $D \leq 1$	Mod2 $D \leq 2$	Mod3 $D \leq 3$	MMax D_{Max}
Test (own)	~ 20000	100	99.8	99.5	99.3	99.2
Test (DMax)	20192	53.5	63.5	83.9	98.9	99.2
Depth=0	6299	100	100	100	100	100
Depth=1	4434	57.9	99.0	98.8	98.5	98.4
Depth=2	2915	34.3	36.8	98.8	98.8	98.4
Depth=3	2396	20.4	23.1	71.1	98.5	98.8
Depth=4	2134	10.2	11.4	43.4	98.8	99.2
Depth=5	2003	11.2	12.3	37.2	97.6	99.8

Out-of-distribution tests (reasoning depth unseen in training)



RuleTakers: A linguistic analog, where a transformer serves as a “soft theorem prover” over knowledge expressed linguistically.

Empirical Success of Pre-trained Model (4/4): Logical Reasoning

Context:

In jurisdictions where use of headlights is optional when visibility is good, drivers who use headlights at all times are less likely to be involved in a collision than are drivers who use headlights only when visibility is poor. Yet Highway Safety Department records show that making use of headlights mandatory at all times does nothing to reduce the overall number of collisions.

Question: Which one of the following, if true, most helps to resolve the apparent discrepancy in the information above?

Options:

- A. In jurisdictions where use of headlights is optional when visibility is good, one driver in four uses headlights for daytime driving in good weather.
- B. Only very careful drivers use headlights when their use is not legally required.
- C. The jurisdictions where use of headlights is mandatory at all times are those where daytime visibility is frequently poor.
- D. A law making use of headlights mandatory at all times is not especially difficult to enforce.

Answer: B

ReClor	
construction method	exams
context type	written text
# of options	4
# of context	6,138
# of questions	6,138
Vocab size	26,576
Context Len	73.6
Question Len	17.0
Option Len	20.6

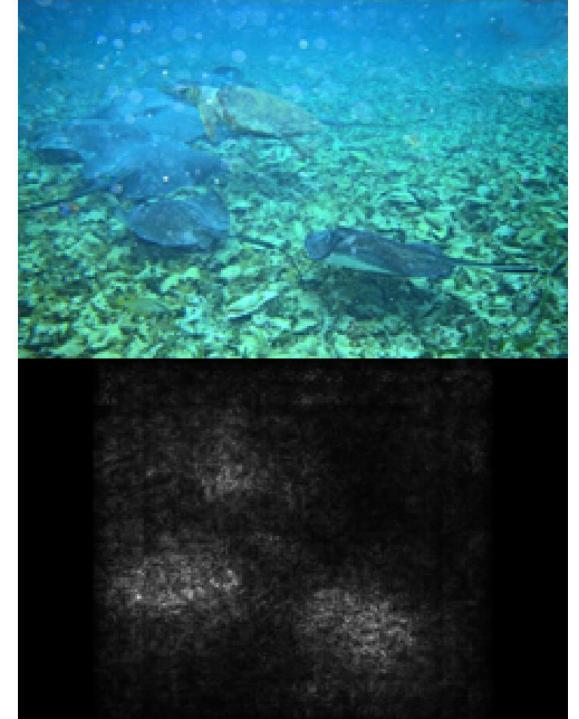
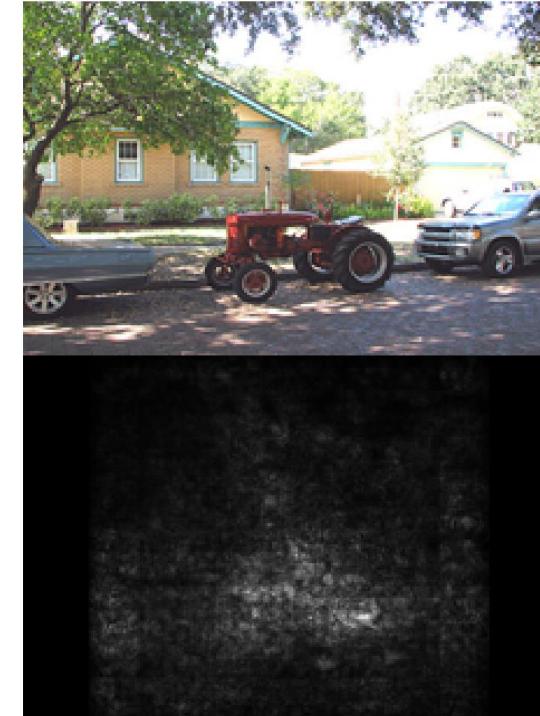
Model	Input	RACE	Val	Test	Test-E	Test-H
Chance	(C, Q, A)		25.0	25.0	25.0	25.0
fastText			32.0	30.8	40.2	23.4
Bi-LSTM			27.8	27.0	26.4	27.5
GPT			47.6	45.4	73.0	23.8
GPT-2			52.6	47.2	73.0	27.0
BERT _{BASE}	(C, Q, A)		54.6	47.3	71.6	28.2
	(C, Q, A)	✓	55.2	49.5	68.9	34.3
BERT _{LARGE}	(A)		46.4	42.4	69.3	21.3
	(Q, A)		48.8	43.4	72.7	20.4
	(C, Q, A)		53.8	49.8	72.0	32.3
	(C, Q, A)	✓	55.6	54.5	73.9	39.3
XLNet _{BASE}	(C, Q, A)		55.8	50.4	75.2	30.9
	(C, Q, A)	✓	62.0	55.5	76.1	39.3
XLNet _{LARGE}	(A)		45.0	42.9	66.1	24.6
	(Q, A)		47.8	43.4	68.6	23.6
	(C, Q, A)		62.0	56.0	75.7	40.5
	(C, Q, A)	✓	70.8	62.4	77.7	50.4
RoBERTa _{BASE}	(C, Q, A)		55.0	48.5	71.1	30.7
	(C, Q, A)	✓	56.8	53.0	72.5	37.7
RoBERTa _{LARGE}	(A)		48.8	43.2	69.5	22.5
	(Q, A)		49.8	45.8	72.0	25.2
	(C, Q, A)		62.6	55.6	75.5	40.0
	(C, Q, A)	✓	68.0	65.1	78.9	54.3
Graduate Students	(C, Q, A)		—	63.0	57.1	67.2
Ceiling Performance	(C, Q, A)		—	100	100	100

Interpretable Techniques

- Post-hoc methods
 - Saliency Maps
 - LIME (Local Interpretable Model-Agnostic Explanations)
 - TCAV (Testing with Concept Activation Vectors)
 - Interpretation Generation
 - ...
- Intrinsic methods
 - Attention
 - Interpretable CNN Filters
 - Neural Module Network
 - ...

Post-Hoc #1: Saliency Maps

- Given an image I_0 , a class c , and a classification ConvNet with the class score function $S_c(I)$, we would like to rank the pixels of I_0 based on their influence on the score $S_c(I_0)$.



Post-Hoc #1: Saliency Maps

- Given an image I_0 , a class c , and a classification ConvNet with the class score function $S_c(I)$, we would like to rank the pixels of I_0 based on their influence on the score $S_c(I_0)$.

$$S_c(I) = w_c^T I + b_c$$



$$S_c(I) \approx w^T I + b$$

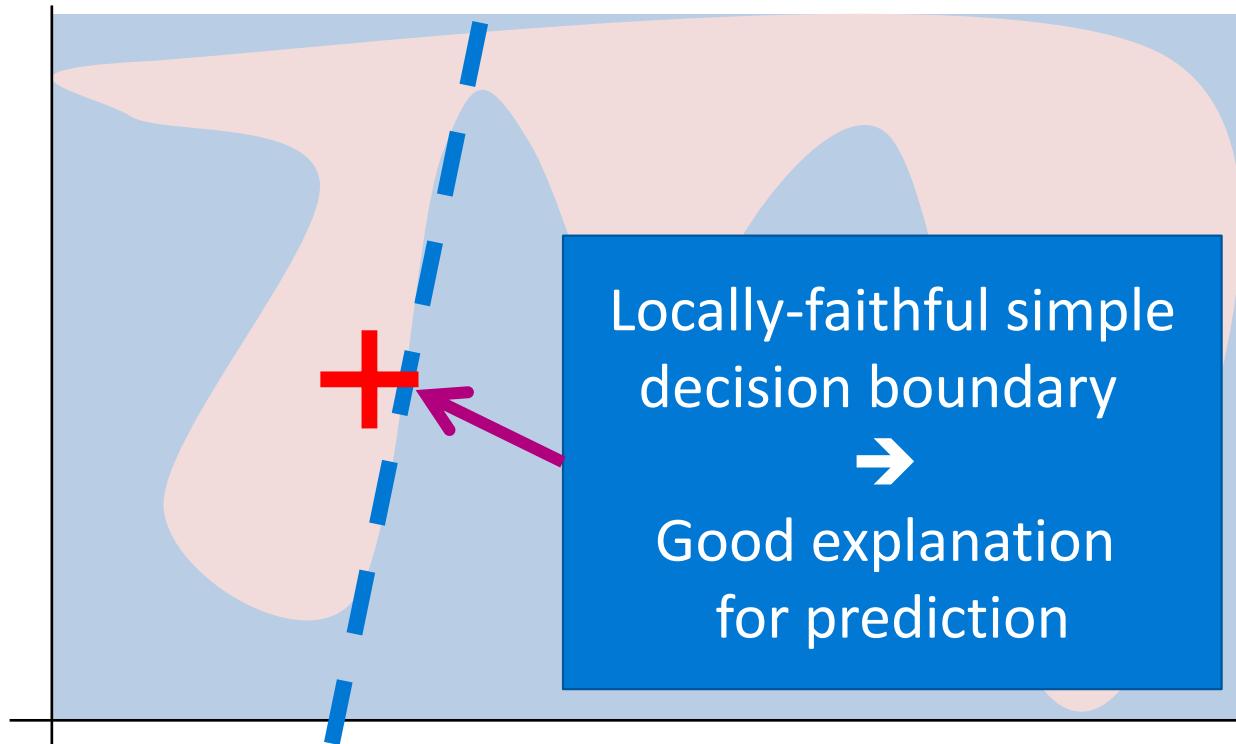
$$w = \frac{\partial S_c}{\partial I} \Big|_{I_0}$$

- first-order Taylor expansion
- the magnitude of the derivative indicates which pixels need to be changed the least to affect the class score the most.

Post-Hoc #2: Local Interpretable Model-Agnostic Explanations (LIME)

- Key Idea
 - Pick a model class interpretable by humans
 - Locally approximate global (blackbox) model. Simple model globally bad, but locally good

1. Sample points around x_i
2. Use complex model to predict labels for each sample
3. Weigh samples according to distance to x_i
4. Learn new simple model on weighted samples
5. Use simple model to explain



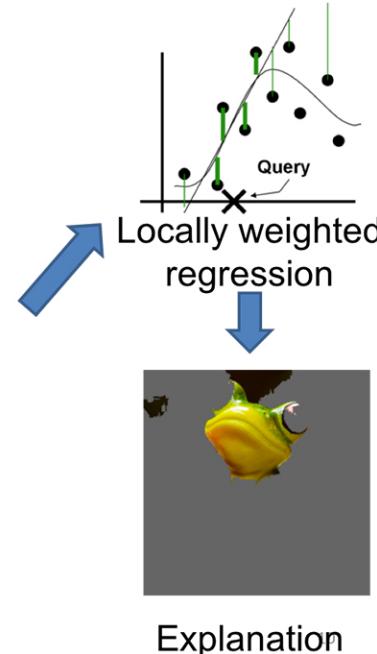
Post-Hoc #2: Local Interpretable Model-Agnostic Explanations (LIME)



Original Image
 $P(\text{tree frog}) = 0.54$

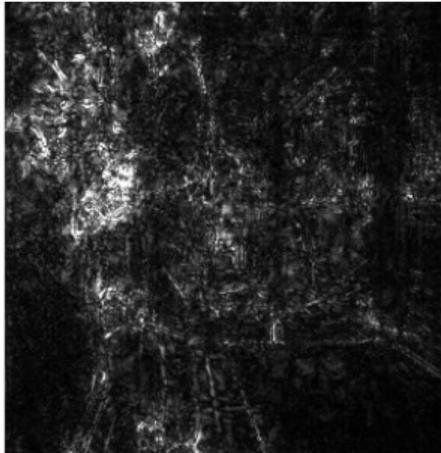
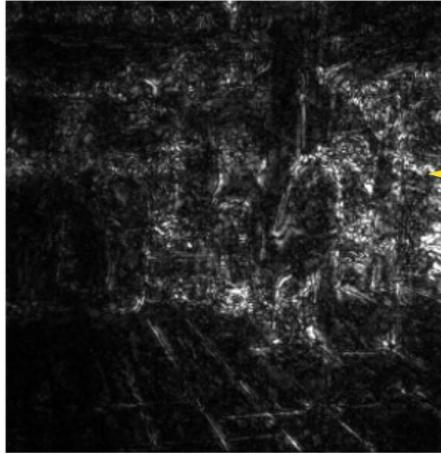


Perturbed Instances	$P(\text{tree frog})$
A photograph of a tree frog with a large red rectangular mask over its eyes and mouth.	0.85
A photograph of a tree frog with its body mostly yellow.	0.00001
A photograph of a tree frog with its eyes red.	0.52



- generate a data set of perturbed instances by turning some of the interpretable components “off” (in this case, making them gray).
- For each perturbed instance, get the probability that a tree frog is in the image according to the model.
- then learn a simple (linear) model on this data set, which is locally weighted—that is, care more about making mistakes in perturbed instances that are more similar to the original image.
- In the end, present the superpixels with highest positive weights as an explanation, graying out everything else.

Post-Hoc #3: Testing with Concept Activation Vectors (TCAV)



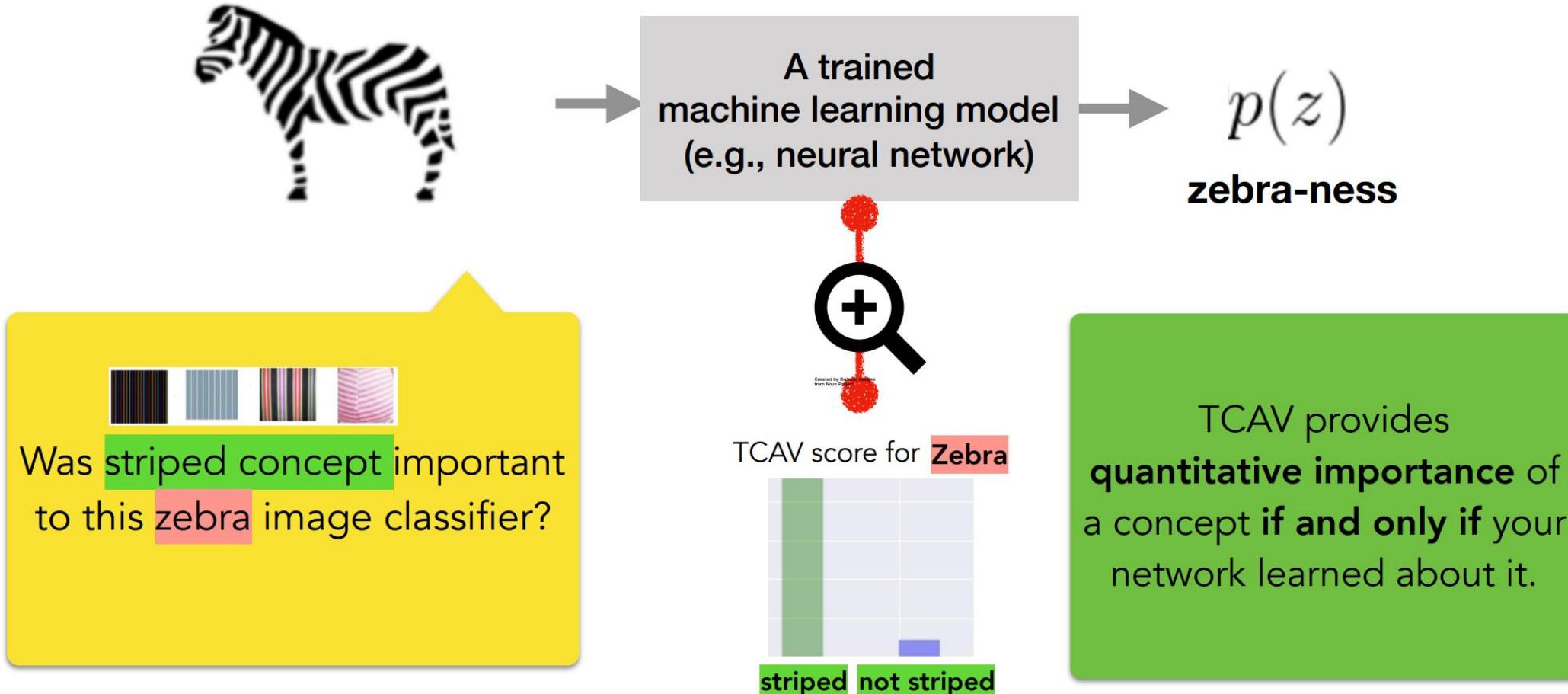
Were there more pixels on the cash machine than on the person?

Did the 'human' concept matter?
Did the 'glasses' or 'paper' matter?

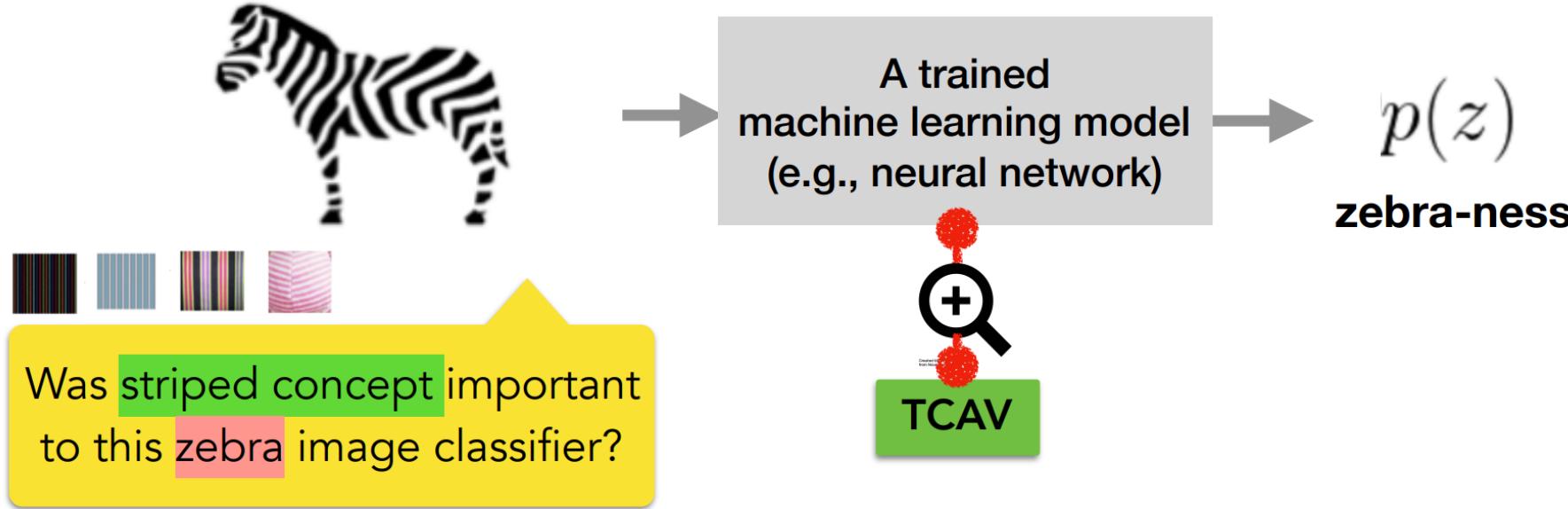
Which concept mattered more?

Is this true for all other cash machine predictions?

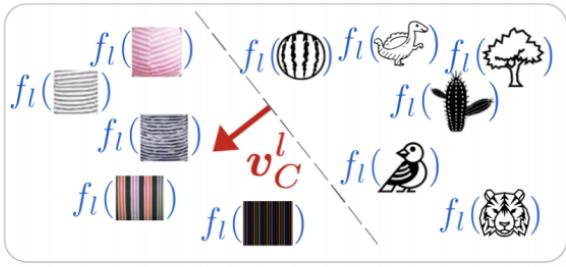
Post-Hoc #3: Goal of TCAV



Post-Hoc #3: TCAV Pipeline



1. Learning CAVs



2. Getting TCAV score

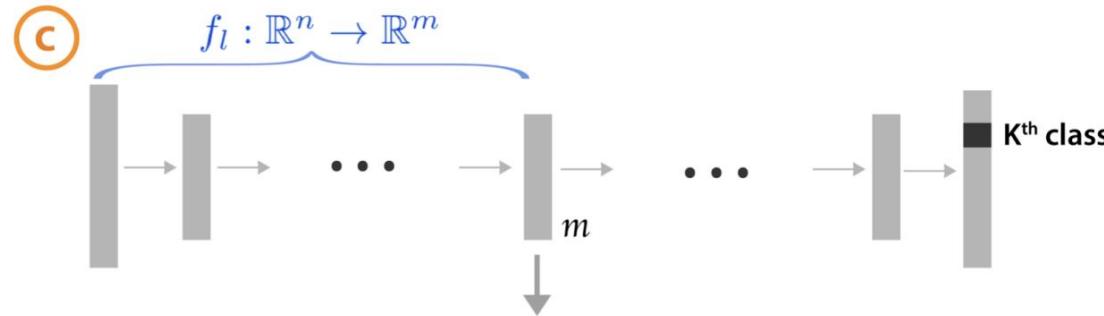
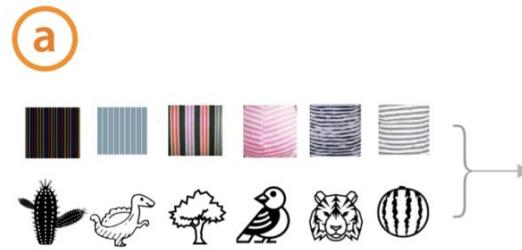
$$S_{C,k,l}(\text{zebra})$$
$$S_{C,k,l}(\text{zebra}) \quad \left. \right\} \rightarrow \text{TCAV}_{Q_{C,k,l}}$$
$$S_{C,k,l}(\text{zebra})$$

3. CAV validation

Qualitative
Quantitative

Post-Hoc #3: Concept Activation Vector (CAV)

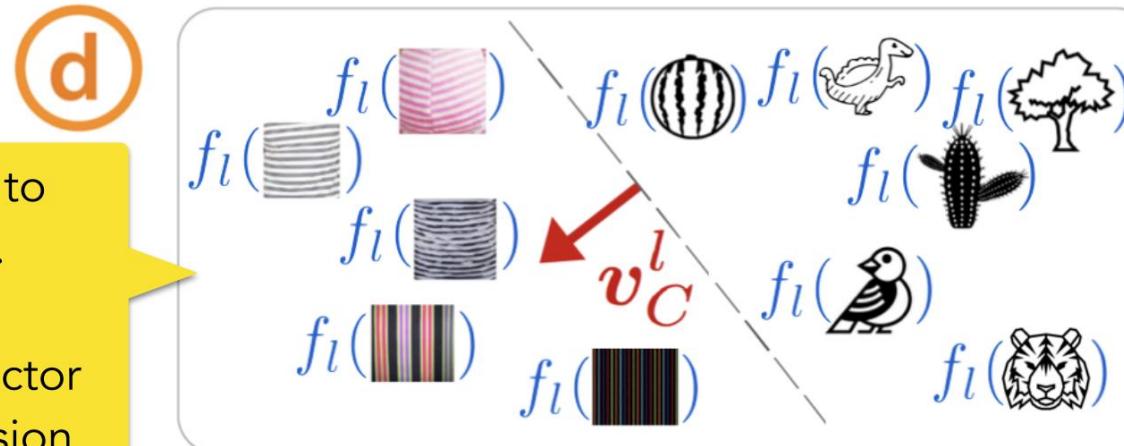
Inputs:



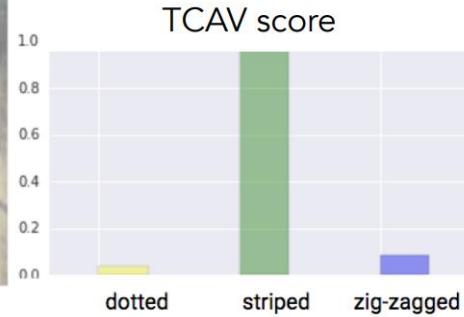
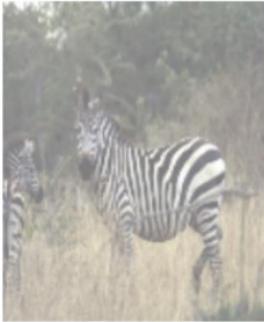
Train a linear classifier to separate activations.

CAV (v_C^l) is the vector **orthogonal** to the decision boundary.

[Smilkov '17, Bolukbasi '16, Schmidt '15]



Post-Hoc #3: Derivative with CAV to get prediction sensitivity



$$\begin{aligned} \text{zebra-ness} &\rightarrow \frac{\partial p(z)}{\partial \mathbf{v}_C^l} = S_{C,k,l}(\mathbf{x}) \\ \text{striped CAV} &\rightarrow \frac{\partial \mathbf{v}_C^l}{\partial \mathbf{v}_C^l} = S_{C,k,l}(\mathbf{x}) \end{aligned}$$

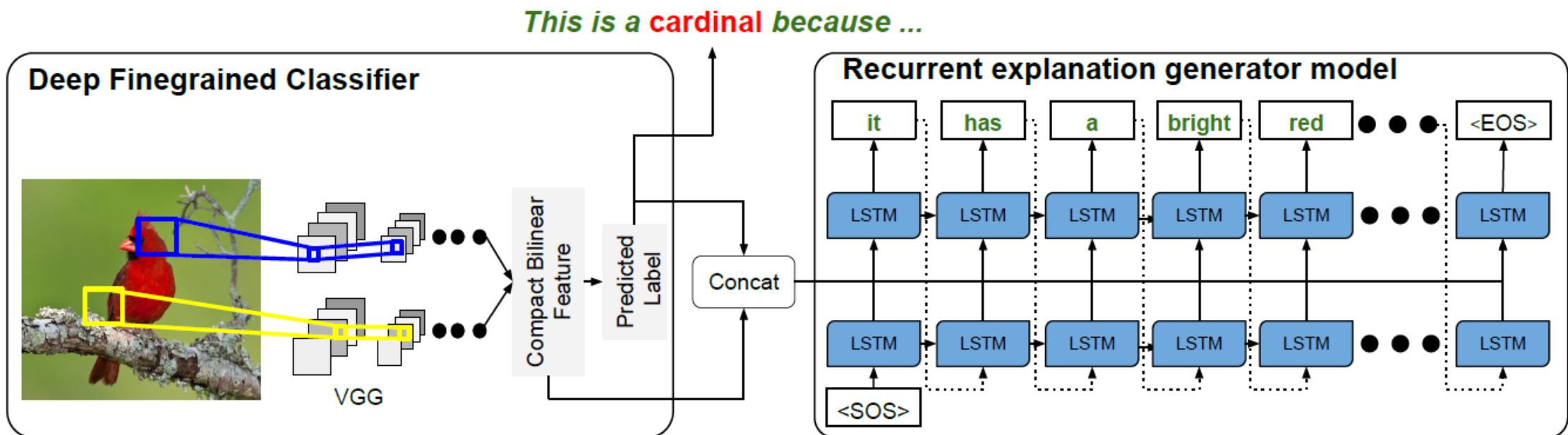
Directional derivative with CAV

$$S_{C,k,l}(\begin{array}{c} \text{zebra} \\ \text{dotted} \end{array})$$
$$S_{C,k,l}(\begin{array}{c} \text{zebra} \\ \text{striped} \end{array})$$
$$S_{C,k,l}(\begin{array}{c} \text{zebra} \\ \text{zig-zagged} \end{array})$$
$$S_{C,k,l}(\begin{array}{c} \text{zebra} \\ \text{striped} \end{array})$$

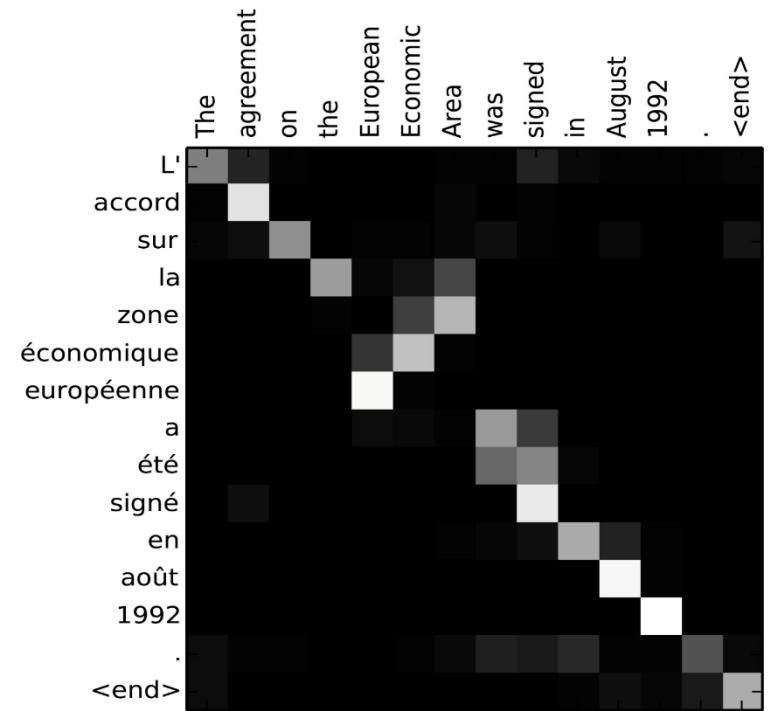
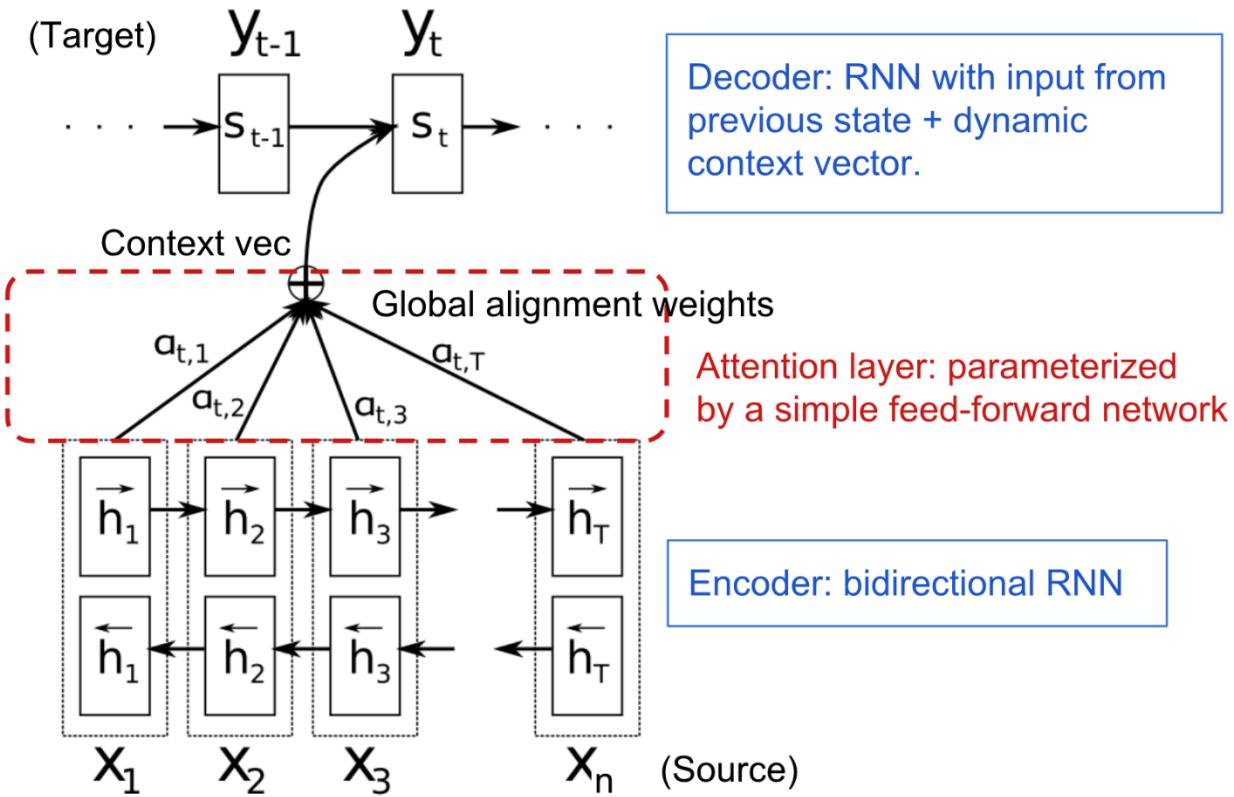
$$\text{TCAV}_{Q_{C,k,l}} = \frac{|\{\mathbf{x} \in X_k : S_{C,k,l}(\mathbf{x}) > 0\}|}{|X_k|}$$

- the “conceptual sensitivity” of class k to concept C in neural activation layer l can be computed as the directional derivative $S_{C,k,l}(x)$.
- X_k denote all inputs with that given label.

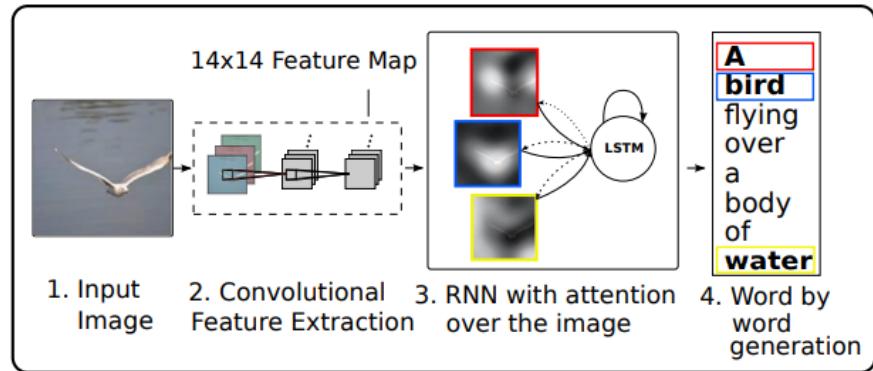
Post-Hoc #4: Interpretation Generation



Intrinsic #1: Attention



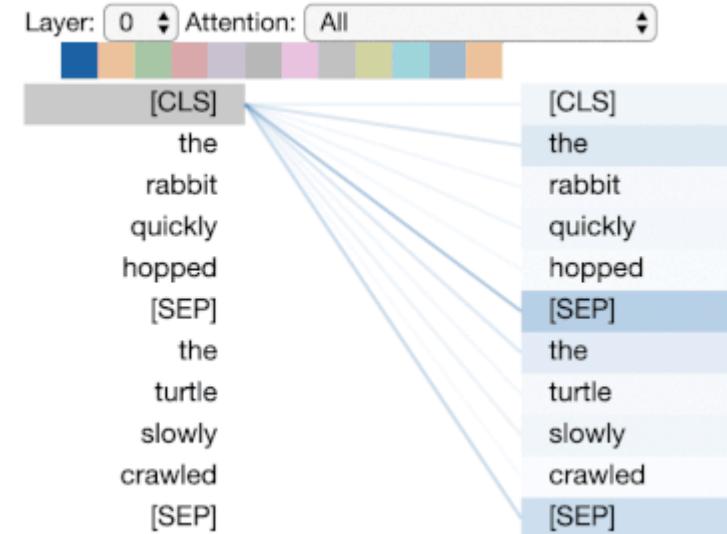
Intrinsic #1: Attention



A woman is throwing a frisbee in a park.

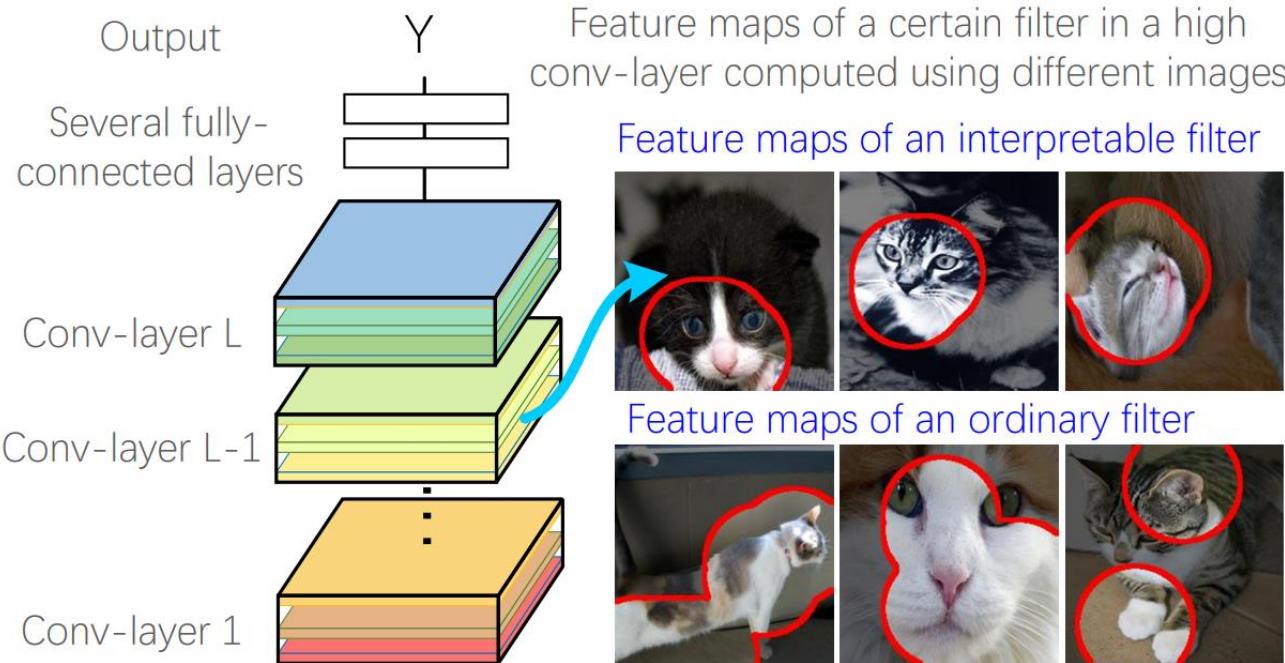


A dog is standing on a hardwood floor.



The *attention-head view* visualizes the attention patterns produced by one or more attention heads in a given transformer layer.

Intrinsic #2: Interpretable CNNs

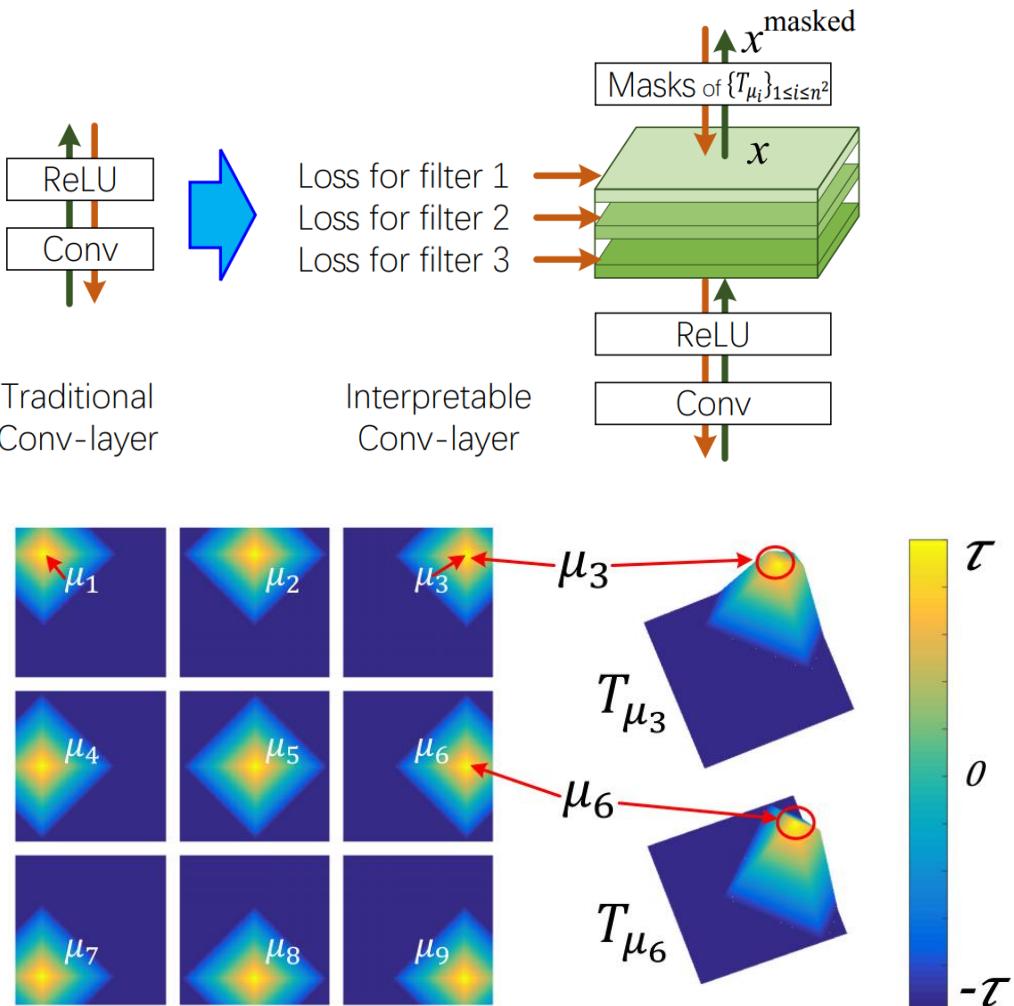


Goal: without any additional human supervision, modify a CNN to obtain interpretable knowledge representations in its conv-layers.

Aim to force each filter in a high conv-layer to represent an object part without using any additional supervision.

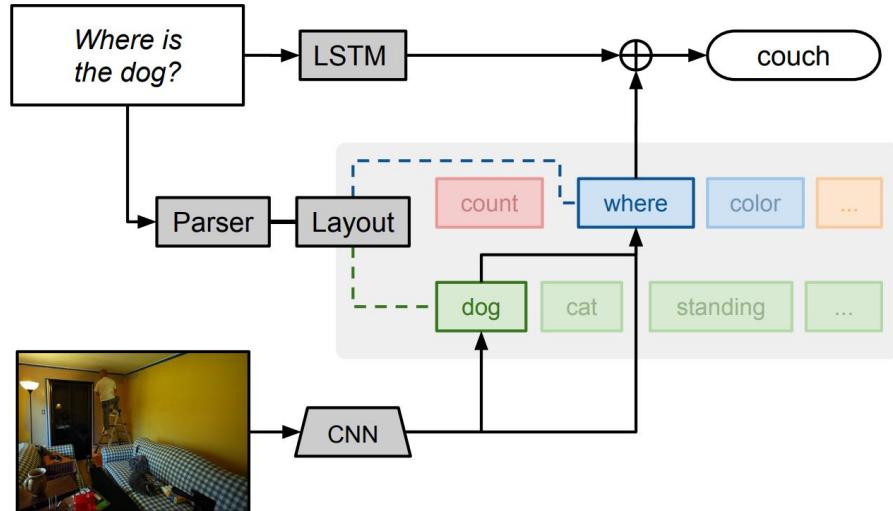
In a traditional CNN, a high-layer filter may describe a mixture of patterns, i.e., the filter may be activated by both the head part and the leg part of a cat. In contrast, the filter in our interpretable CNN is activated by a single part.

Intrinsic #2: Interpretable CNNs



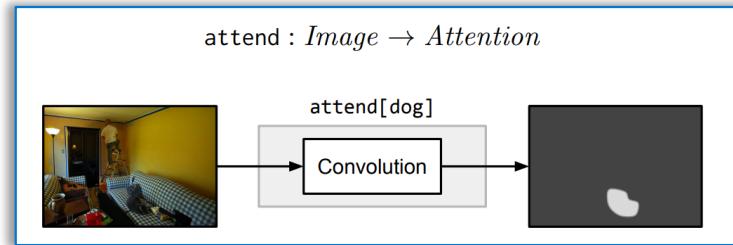
- Add additional losses to force each convolutional filter in our interpretable CNN to represent a specific object part.
- The filter loss is formulated as the minus mutual information between a set of feature maps and a set of pre-defined templates
- In comparisons, a filter in ordinary CNNs usually represents a mixture of parts and textures.
- Clear semantic meanings of middle-layer filters are of significant values in real applications.

Intrinsic #3: Neural Module Network

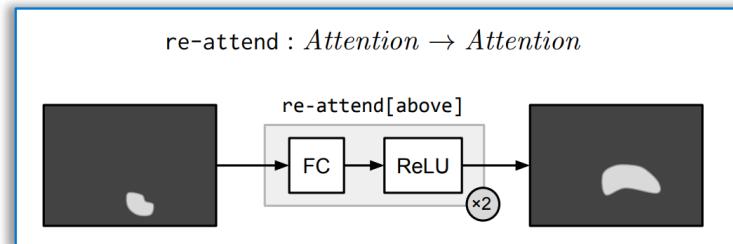


**Reusable neural modules
with different architectures**

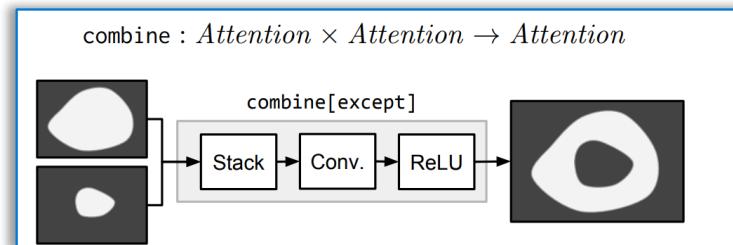
Attention



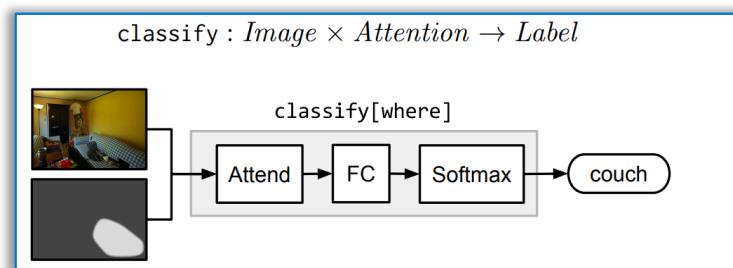
Re-Attention



Combination



Classification



Benefits and Challenges

- Post-hoc Interpretable Models
 - Benefits: remain existing high-performance model unchanged
 - Challenges: to what extent was the interpretable model to mimic the original model
- Intrinsic Interpretable Models
 - Benefits: good interpretability
 - Challenges: performance may be discounted as non-interpretable strategies which could bring performance boost may not be easily utilized
- Dilemma of “Performance versus Interpretability”
 - How to develop models to take the best of both worlds
 - How to introduce evaluation metrics to foster research towards both dimensions

Summary

Machine Reasoning Frameworks

(covered by this tutorial)

Symbolic Reasoning

Symbolic Knowledge
+ Inference based on Logic Rules

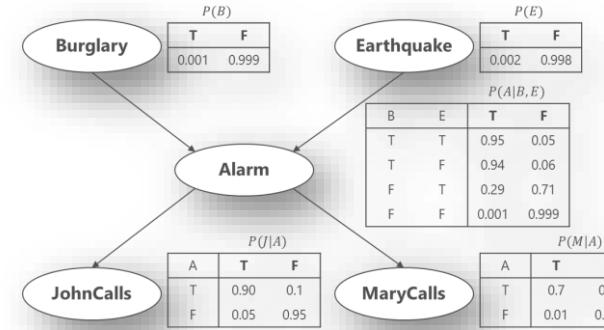
- KB₁:** $\forall x \text{ cat}(x) \Rightarrow \text{like}(x, \text{fish})$
KB₂: $\forall x \forall y (\text{cat}(x) \wedge \text{like}(x, y)) \Rightarrow \text{eat}(x, y)$
KB₃: $\text{cat}(\text{Tom})$



$\alpha: \text{eat}(\text{Tom}, \text{fish})$

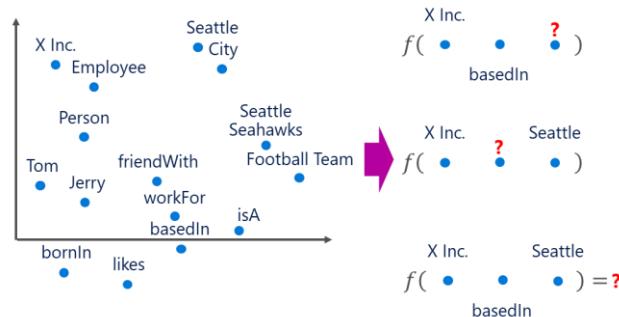
Probabilistic Reasoning

Probabilistic Symbolic Knowledge
+ Inference based on Probabilistic Graphical Models



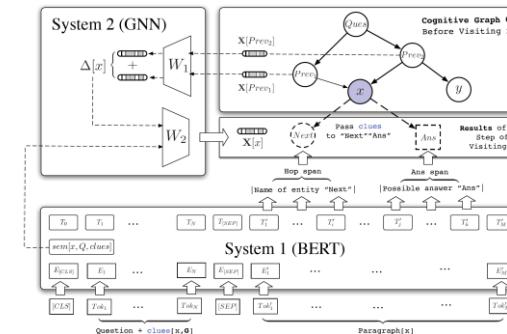
Neural-Symbolic Reasoning

Vector Representation of Symbolic Knowledge
+ Inference based on Neural Networks

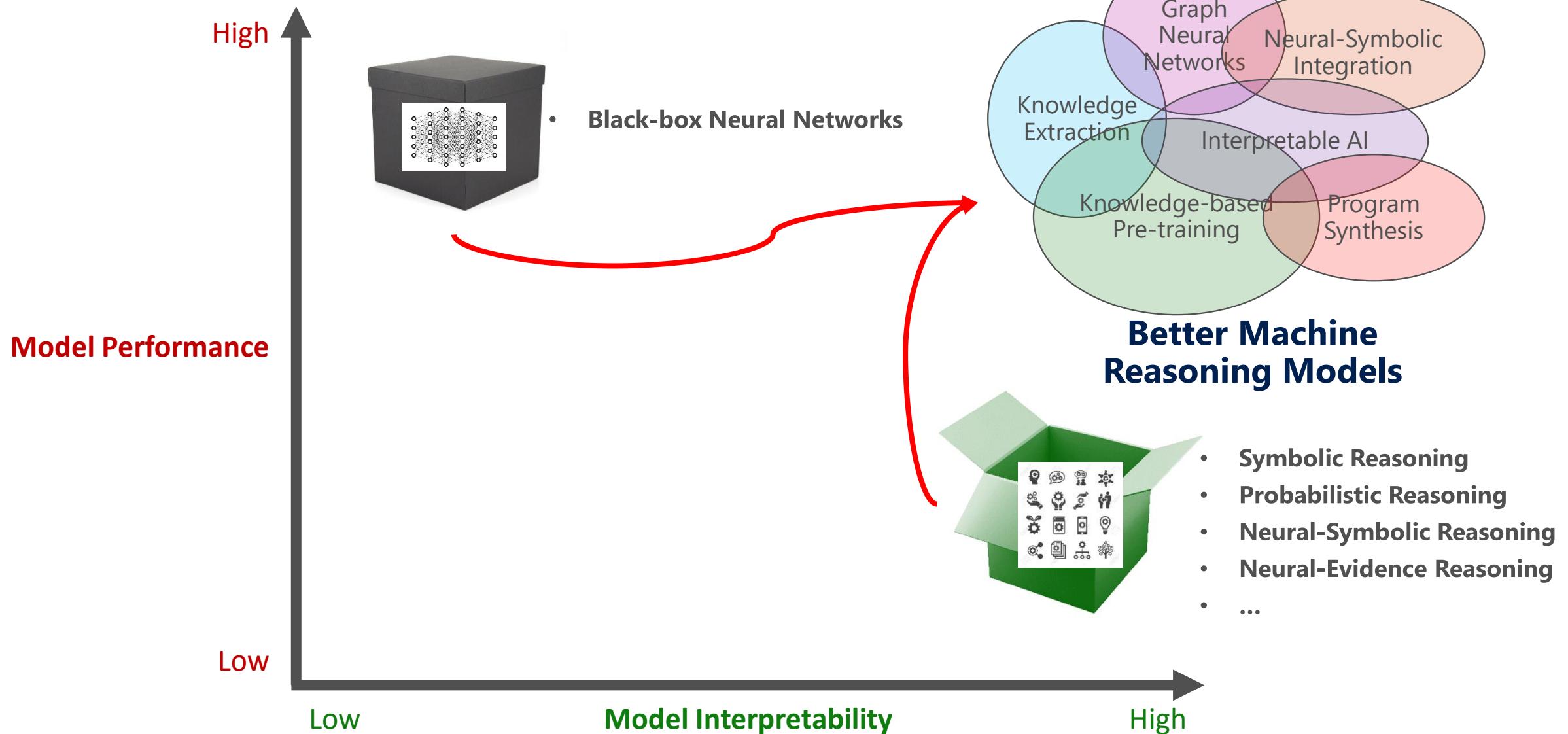


Neural-Evidence Reasoning

Vector Representation of Non-Symbolic Evidence Knowledge
+ Inference based on Neural Networks



Usefulness, Dilemma and Bright Future 😊



Further Readings



Commonsense Knowledge Representation and Reasoning in Natural Language processing



Yejin Choi



Vered Shwartz



Maarten Sap



Antoine Bosselut



Dan Roth

<https://homes.cs.washington.edu/~msap/acl2020-commonsense/>

Neuro-Symbolic Visual Reasoning and Program Synthesis

Virtual CVPR 2020 Tutorial

<http://nscv.csail.mit.edu/>

The Third AI Summer AAAI Robert S. Engelmore Memorial Lecture

Henry Kautz

Professor, Department of Computer Science
University of Rochester

<https://www.cs.rochester.edu/u/kautz/talks/index.html>

Scalable Construction and Reasoning of Massive Knowledge Bases

Xiang Ren¹ Nanyun Peng¹ William Yang Wang²

University of Southern California¹

University of California, Santa Barbara²

<https://www.aclweb.org/anthology/N18-6003/>

Statistical Relational Learning

Pedro Domingos

Dept. of Computer Science & Eng.
University of Washington



<https://homes.cs.washington.edu/~pedrod/cikm13.html>

Probabilistic Logic Programming and its Applications

Luc De Raedt

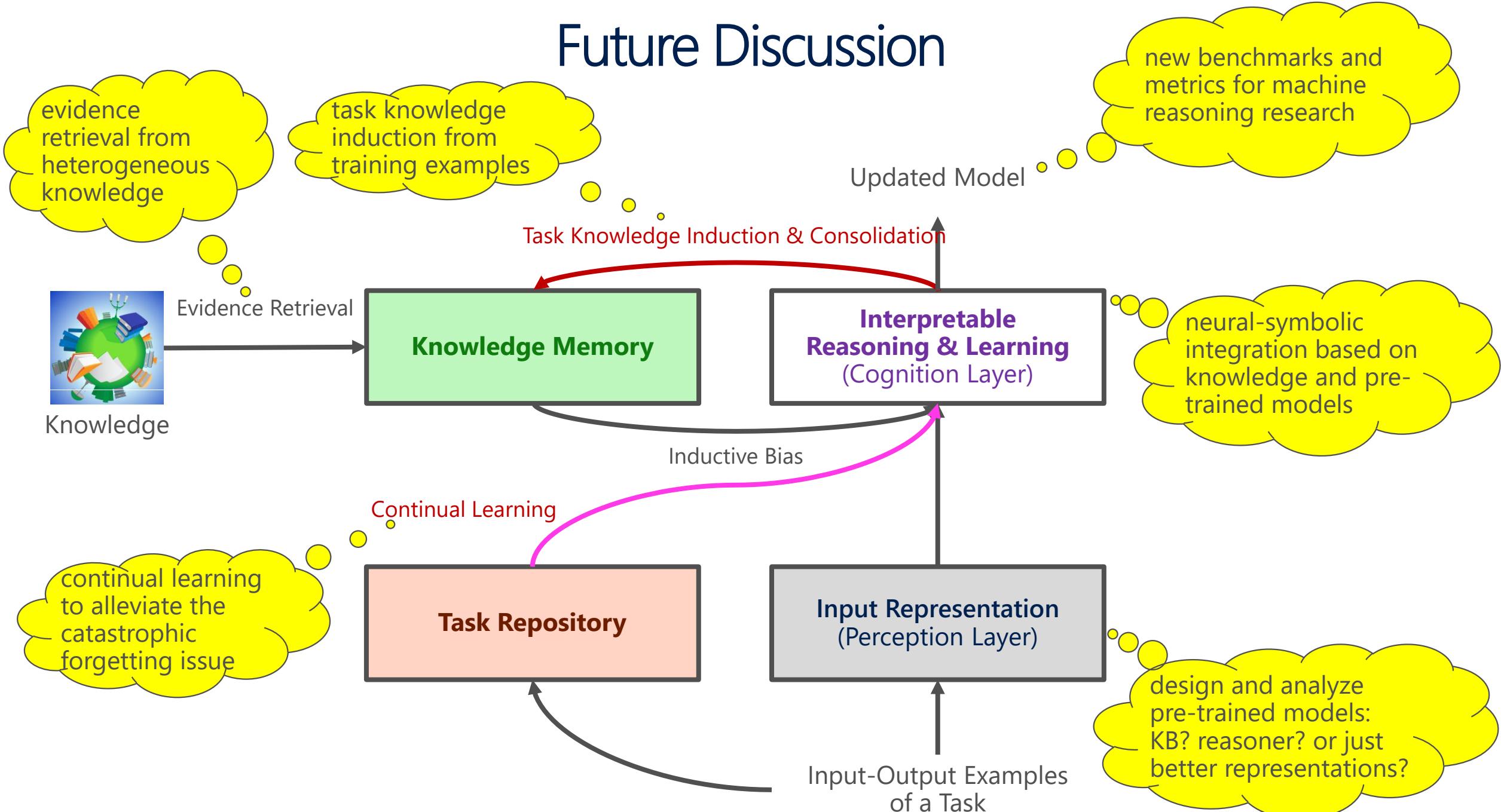
with many slides from Angelika Kimmig



The Turing, London, September 11, 2017

<https://logic-data-science.github.io/Slides/DeRaedt.pdf>

Future Discussion



Thank you/谢谢!