

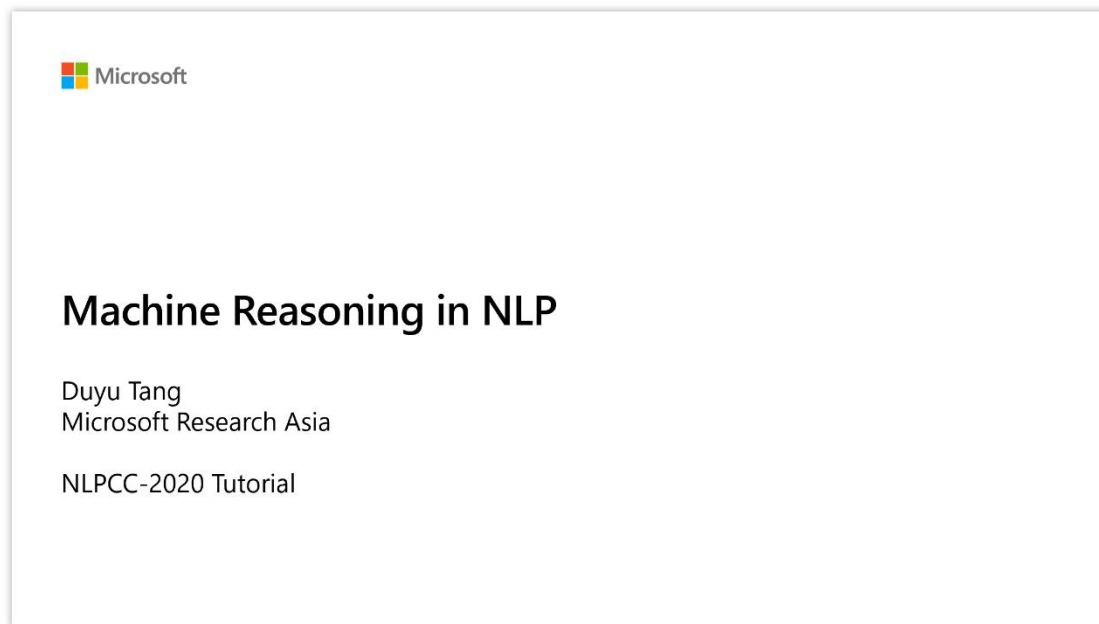


# Machine Reasoning in NLP

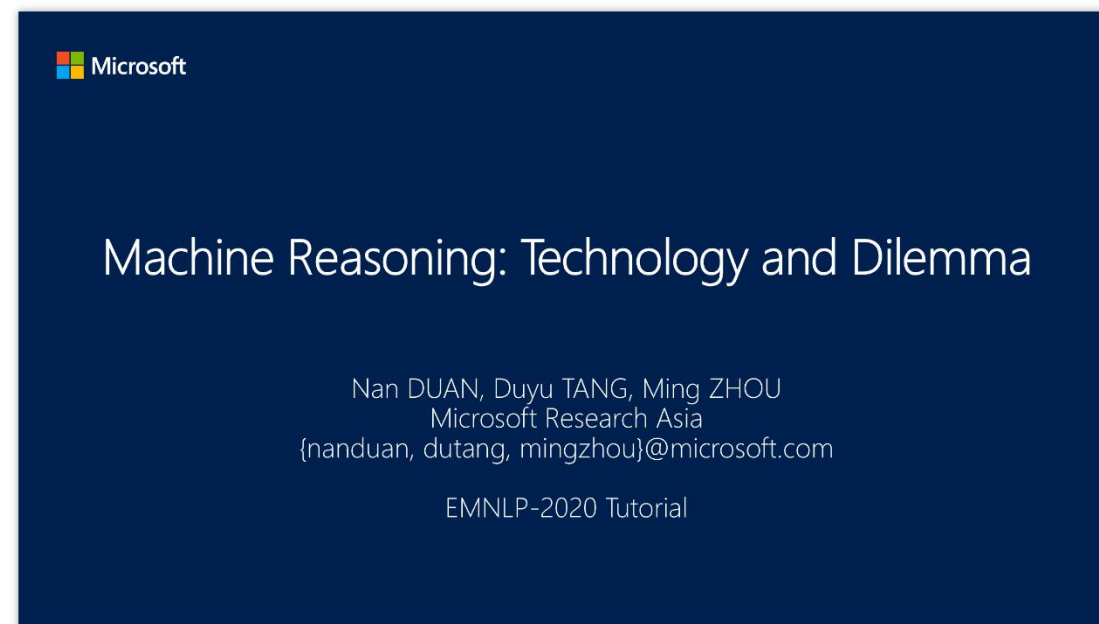
Duyu Tang  
Microsoft Research Asia

# Past Tutorials on Machine Reasoning

- NLPCC–2020 Tutorial (3 hours)



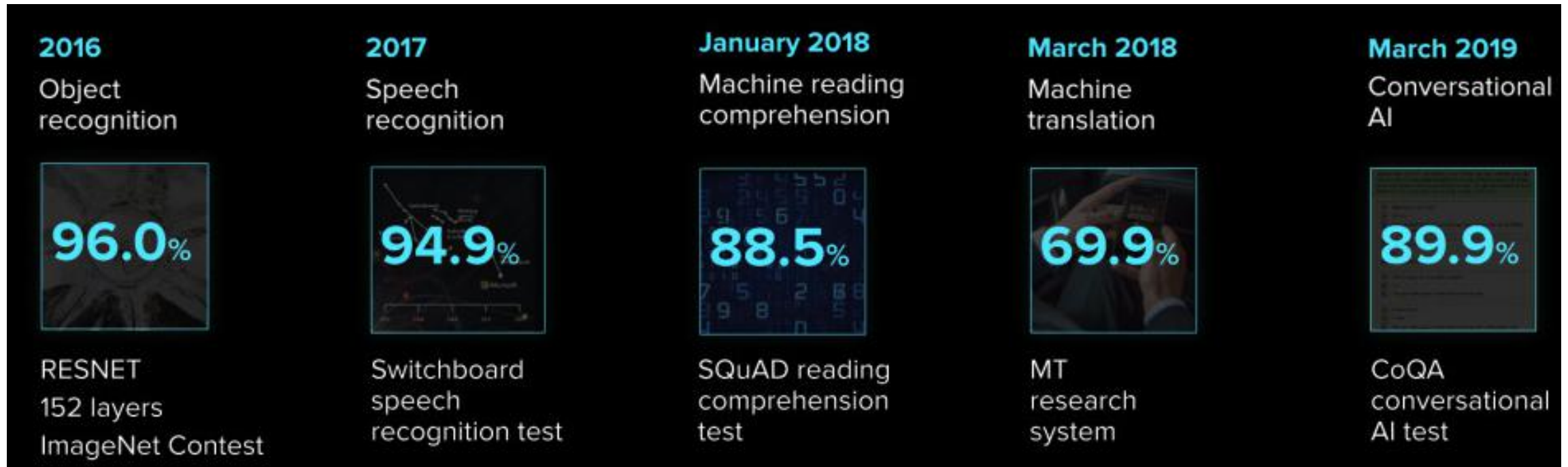
- EMNLP–2020 Tutorial (3 hours)

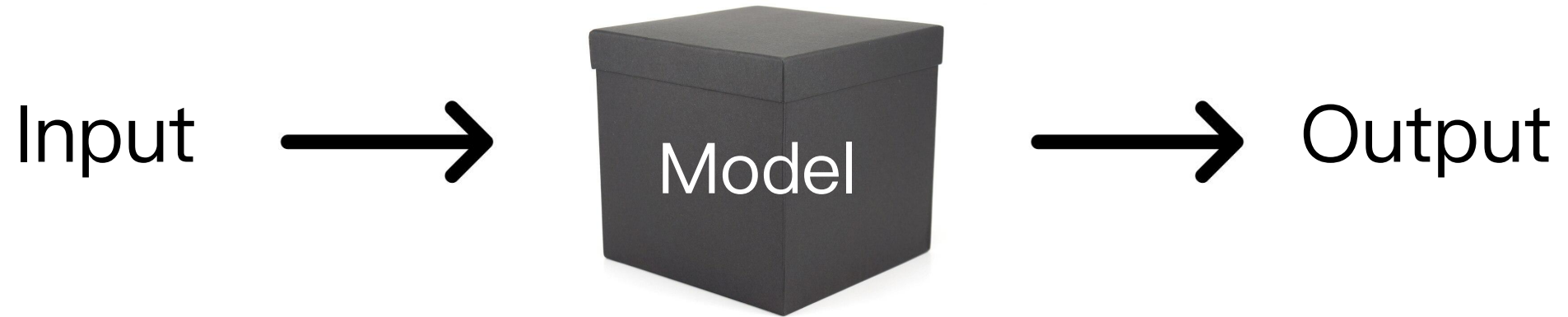


Both tutorials are available at  
<https://tangduyu.github.io/>

# Microsoft AI Breakthroughs

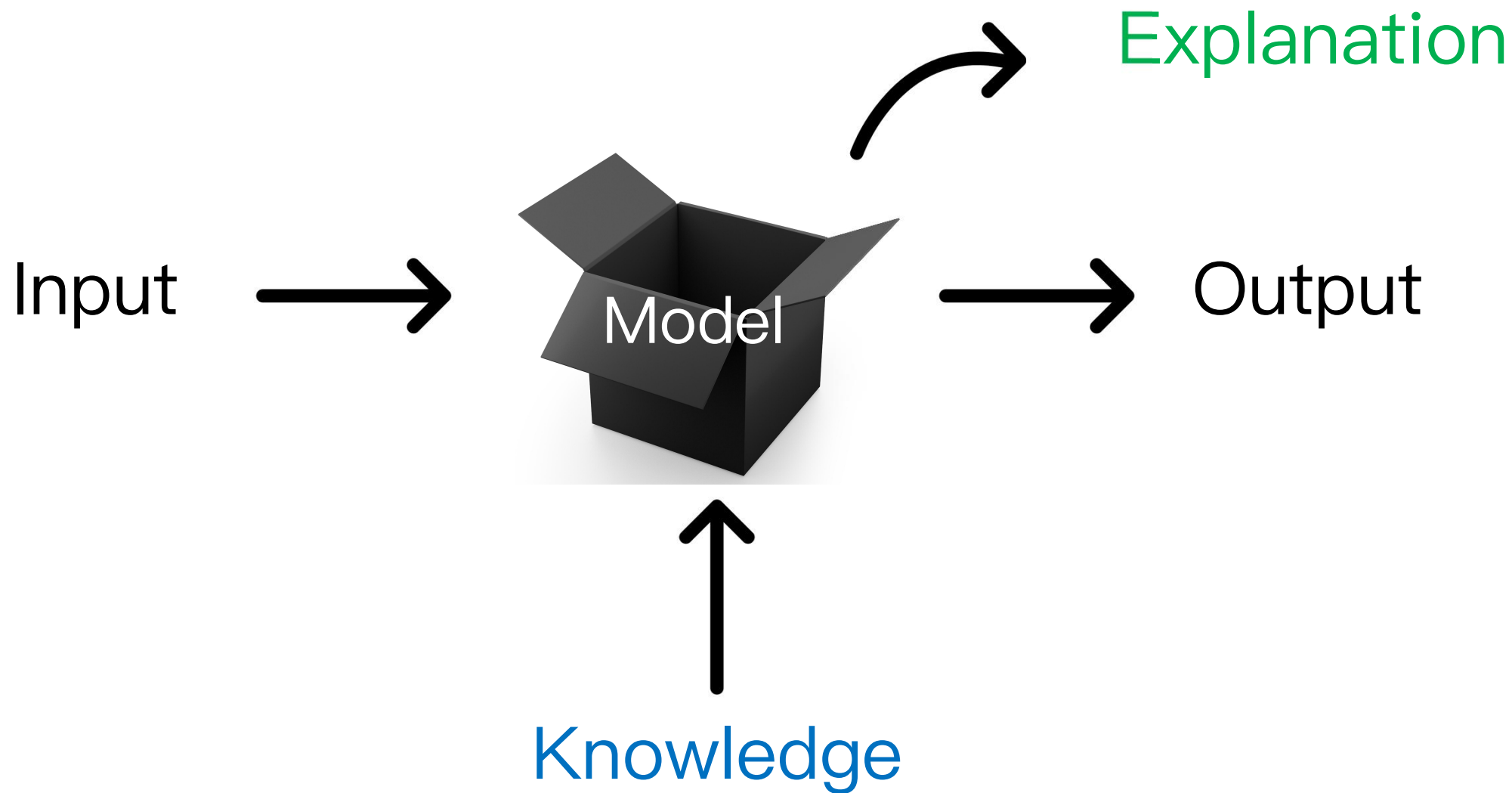
- Gradually approaching human parity



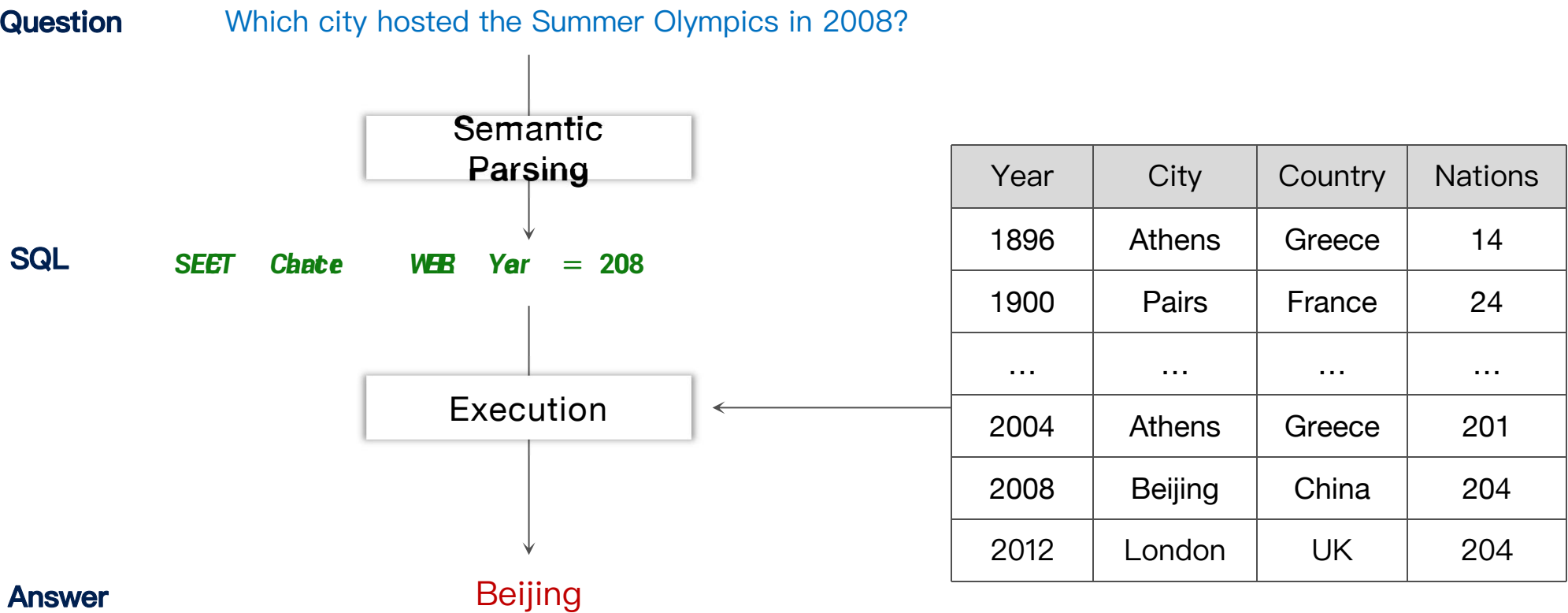


**Limitations:**

1. Lack of transparency of the decision-making process
2. Highly rely on annotated data, ignore human/expert knowledge



# Example #1: Simple Question Answering



# Example #2: Multi-Turn Question Answering

**Q1:** Which city hosted the Summer Olympics in 2008?  
*SELECT Character WHERE Year = 2008*    **A1:**    **Beijing**

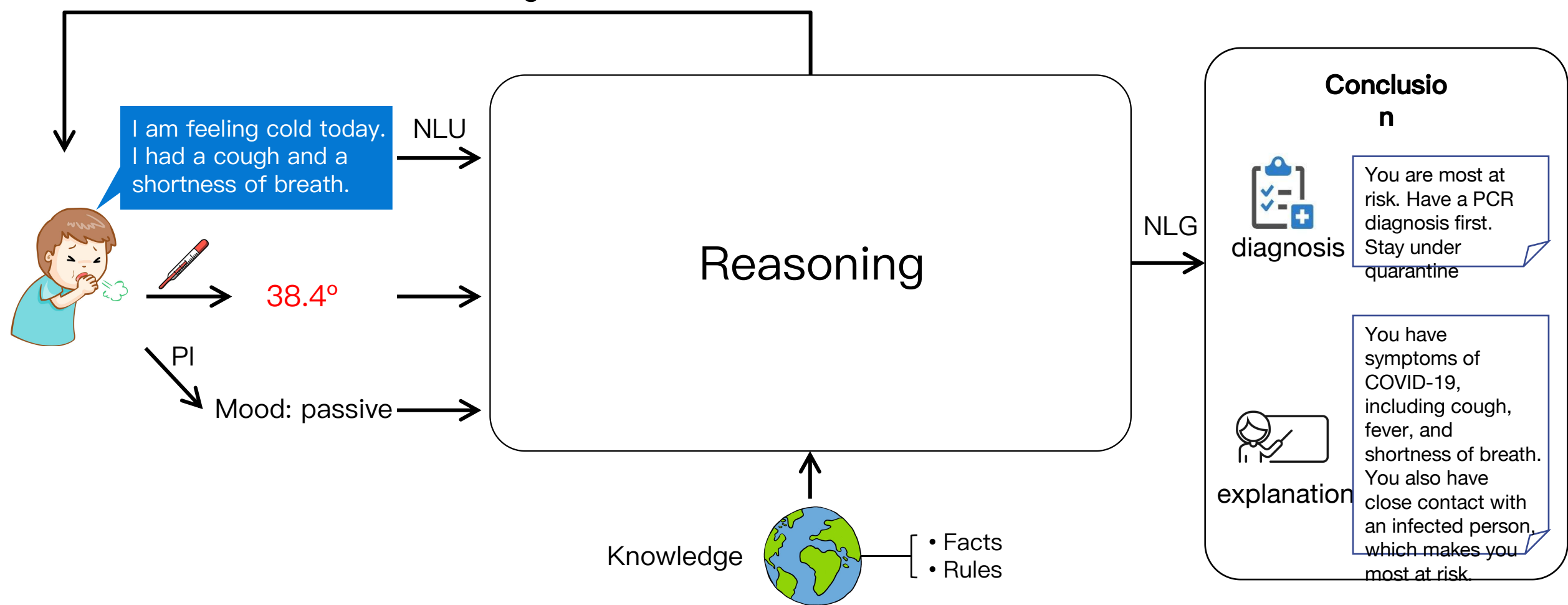
**Q2:** How many nations participate that year?  
*SELECT Nations WHERE Year = 2008*    **A2:**    **204**

**Q3:** How about 2004?  
*SELECT Nations WHERE Year = 2004*    **A3:**    **201**

Year	City	Country	Nations
1896	Athens	Greece	14
1900	Paris	France	24
...	...	...	...
2004	Athens	Greece	201
2008	Beijing	China	204
2012	London	UK	204

# Example #3: Medical Diagnosis

## Multi-Turn Question Answering & Clarification



**Required Knowledge:** Symbolic + Conversation + **Domain Knowledge**



# Example #3: Fact



WHO Coronavirus Disease (COVID-19) Dashboard  
Data last updated: 2020/10/5, 3:54pm CEST

i

←

Covid-19 Response Fund

Donate

Overview Data Table Explore

## Situation by Country, Territory & Area

Name	Cases - cumulative total	⇅	Cases - newly reported in last 24 hours	Deaths - cumulative total	Deaths - newly reported in last 24 hours	Transmission Classification
Global	35,109,317		294,763	1,035,341	4,526	
United States...	7,305,270	<div></div>	49,036	208,064	698	Community transmission
India	6,623,815	<div></div>	74,442	102,685	903	Clusters of cases
Brazil	4,906,833	<div></div>	26,310	145,987	599	Community transmission
Russian Fed...	1,225,889	<div></div>	10,888	21,475	117	Clusters of cases
Colombia	848,147	<div></div>	6,616	26,556	159	Community transmission
Peru	824,985	<div></div>	3,421	32,665	56	Community transmission
Argentina	790,818	<div></div>	11,129	20,795	196	Community transmission
Spain	789,932	<div></div>	0	32,086	0	Clusters of cases
Mexico	757,953	<div></div>	4,863	78,880	388	Community transmission
South Africa	681,289	<div></div>	1,573	16,976	38	Community transmission

<https://covid19.who.int/table>

# Example #3: Rule

## SYMPTOM

IF <a person is infected >  
THEN <he/she may have  
fever>



mild

IF <a person is infected >  
THEN <he/she may have  
cough and shortness of  
breath>



severe

## DIAGNOSIS & TREATMENT

IF < disease == COVID-19>  
THEN  
<diagnosis(disease)=PCR  
(Polymerase Chain



IF < disease == COVID-19>  
THEN <treatment(disease)=none  
AND vaccine(disease)=none >  
(by Jan 31, 2020)



## GROUPS AT RISK

IF <a person has close contact with  
animals>  
THEN <the person is at risk>



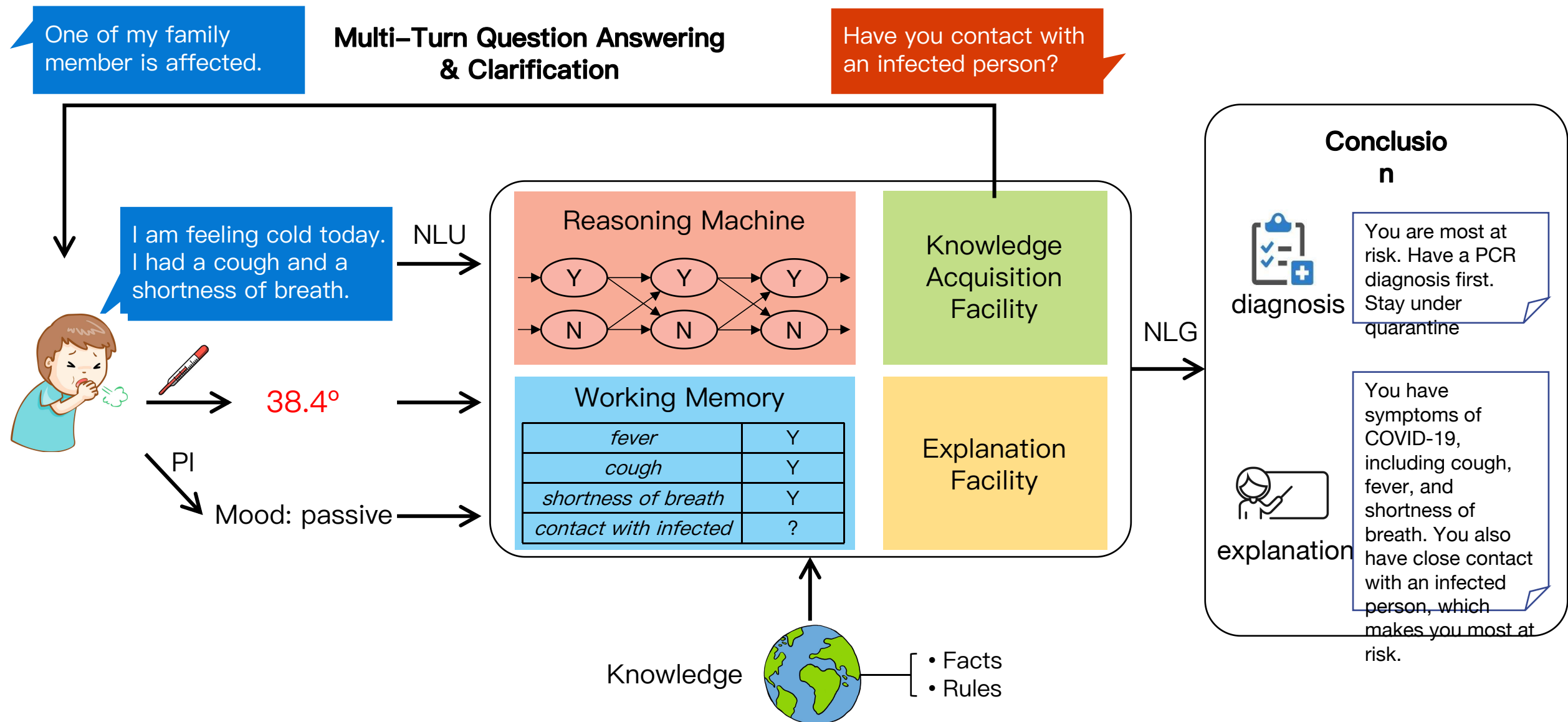
IF <a person is a live animal market  
worker>  
THEN <he/she has close contact with  
animals>

IF <a person has close contact with an  
infected person>  
THEN <the person is most at risk>

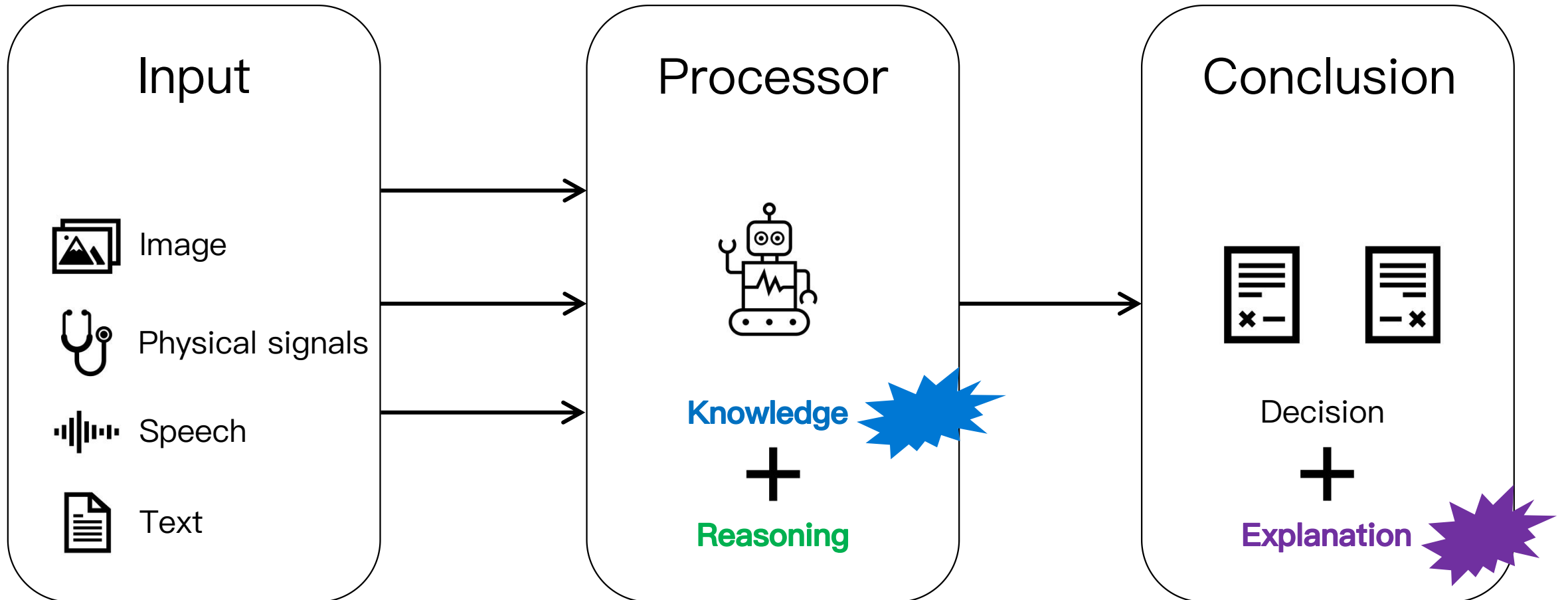


IF <a person is a healthcare worker or a  
family member of infected person>  
THEN <the person has close contact  
with infected person >

# Example #3: Machine Reasoning Pipeline



# Features of Machine Reasoning



+ Propositional/First-Order Logic  
+ Neuro-symbolic

$(\alpha \wedge \beta) \equiv (\beta \wedge \alpha)$  commutativity of  $\wedge$   
 $(\alpha \vee \beta) \equiv (\beta \vee \alpha)$  commutativity of  $\vee$   
 $((\alpha \wedge \beta) \wedge \gamma) \equiv (\alpha \wedge (\beta \wedge \gamma))$  associativity of  $\wedge$   
 $((\alpha \vee \beta) \vee \gamma) \equiv (\alpha \vee (\beta \vee \gamma))$  associativity of  $\vee$   
 $\neg(\neg\alpha) \equiv \alpha$  double-negation elimination  
 $(\alpha \Rightarrow \beta) \equiv (\neg\beta \Rightarrow \neg\alpha)$  contraposition  
 $(\alpha \Rightarrow \beta) \equiv (\neg\alpha \vee \beta)$  implication elimination  
 $(\alpha \Rightarrow \beta) \equiv ((\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha))$  biconditional elimination

First-Order Logic

**SEET** **CON** **CE** **Tam** **WEB**

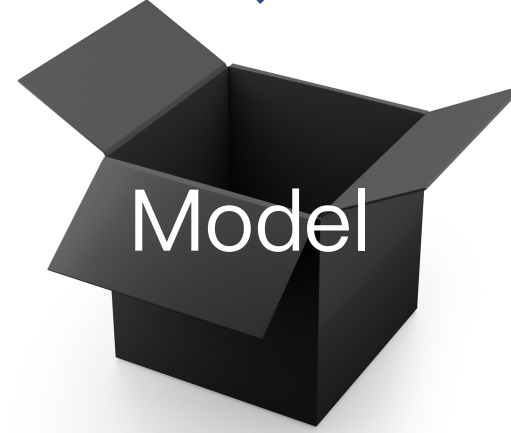
*Cbèg* = "Yok "

$\lambda x. \text{pepe} . \text{peon} . \text{pae} \_ \text{of\_bit} \text{ h(}$

*Dakd* *Tupn* , *x*)

Symbolic Operations

Input



Model



Output

+ ~~Pre-trained models (e.g. ELMo, BERT, GPT)~~  
+ Evidence (e.g. retrieved docs from Wikipedia/web,  
retrieved facts from Wikidata/ConceptNet)



# Agenda

Opening

Logic-based Models in NLP

Neural-Symbolic Models in NLP

Evidence-based Models in NLP

Summary

# Logic-based Models in NLP

# Outline

- Propositional Logic and First–Order Logic
- Inference/Theorem Proving: Forward and Backward Chaining
- Application in NLP





# Propositional Logic

- **Logical constants:** true, false
- **Propositional symbols:** P, Q, S, ... (**atomic sentences**)
- Wrapping **parentheses:** ( ... )
- Sentences are combined by **connectives:**
  - $\wedge$  ...and [conjunction]
  - $\vee$  ...or [disjunction]
  - $\rightarrow$  ...implies [implication / conditional]
  - $\leftrightarrow$  ..is equivalent [biconditional]
  - $\neg$  ...not [negation]
- **Literal:** atomic sentence or negated atomic sentence

# Propositional Logic Examples

- P means “It is hot”
- Q means “It is humid”
- R means “It is raining”
  
- $(P \wedge Q) \rightarrow R$   
“If it is hot and humid, then it is raining”
- $Q \rightarrow P$   
“If it is humid, then it is hot”
- Q  
“It is humid.”

# Propositional Logic Syntax

- Given: a set of proposition symbols  $\{X_1, X_2, \dots, X_n\}$ 
  - (we often add **True** and **False** for convenience)
- $X_i$  is a sentence
- If  $\alpha$  is a sentence then  $\neg\alpha$  is a sentence
- If  $\alpha$  and  $\beta$  are sentences then  $\alpha \wedge \beta$  is a sentence
- If  $\alpha$  and  $\beta$  are sentences then  $\alpha \vee \beta$  is a sentence
- If  $\alpha$  and  $\beta$  are sentences then  $\alpha \Rightarrow \beta$  is a sentence
- If  $\alpha$  and  $\beta$  are sentences then  $\alpha \Leftrightarrow \beta$  is a sentence
- And p.s. there are no other sentences!

# Logical Equivalence

$(\alpha \wedge \beta) \equiv (\beta \wedge \alpha)$	commutativity of $\wedge$
$(\alpha \vee \beta) \equiv (\beta \vee \alpha)$	commutativity of $\vee$
$((\alpha \wedge \beta) \wedge \gamma) \equiv (\alpha \wedge (\beta \wedge \gamma))$	associativity of $\wedge$
$((\alpha \vee \beta) \vee \gamma) \equiv (\alpha \vee (\beta \vee \gamma))$	associativity of $\vee$
$\neg(\neg\alpha) \equiv \alpha$	double-negation elimination
$(\alpha \Rightarrow \beta) \equiv (\neg\beta \Rightarrow \neg\alpha)$	contraposition
$(\alpha \Rightarrow \beta) \equiv (\neg\alpha \vee \beta)$	implication elimination
$(\alpha \Leftrightarrow \beta) \equiv ((\alpha \Rightarrow \beta) \wedge (\beta \Rightarrow \alpha))$	biconditional elimination
$\neg(\alpha \wedge \beta) \equiv (\neg\alpha \vee \neg\beta)$	De Morgan
$\neg(\alpha \vee \beta) \equiv (\neg\alpha \wedge \neg\beta)$	De Morgan
$(\alpha \wedge (\beta \vee \gamma)) \equiv ((\alpha \wedge \beta) \vee (\alpha \wedge \gamma))$	distributivity of $\wedge$ over $\vee$
$(\alpha \vee (\beta \wedge \gamma)) \equiv ((\alpha \vee \beta) \wedge (\alpha \vee \gamma))$	distributivity of $\vee$ over $\wedge$

# Problems with Propositional Logic

- Hard to identify “individuals” (e.g., Mary, 3)
- Can’t directly talk about properties of individuals or relations between individuals (e.g., “Bill is tall”)
- Generalizations, patterns, regularities can’t easily be represented (e.g., “all triangles have 3 sides”)
- Lack of variables prevents stating more general rules
  - We need a set of similar rules for each cell
- First–Order Logic is expressive enough to concisely represent this kind of information
  - FOL adds relations, variables, and quantifiers, e.g.,
    - “*Every elephant is gray*”:  $\forall x (\text{elephant}(x) \rightarrow \text{gray}(x))$
    - “*There is a white alligator*”:  $\exists x (\text{alligator}(X) \wedge \text{white}(X))$

# First–Order Logic

- First–order logic models the world in terms of
  - **Objects**, which are things with individual identities
  - **Properties** of objects that distinguish them from other objects
  - **Relations** that hold among sets of objects
  - **Functions**, which are a subset of relations where there is only one “value” for any given “input”
- Examples:
  - Objects: Students, lectures, companies, cars ...
  - Relations: Brother–of, bigger–than, outside, part–of, has–color, occurs–after, owns, visits, precedes, ...
  - Properties: blue, oval, even, large, ...
  - Functions: father–of, best–friend, second–half, one–more–than ...

# First–Order Logic Examples

- Universal quantification
  - $(\forall x)P(x)$  means that  $P$  holds for **all** values of  $x$  in the domain associated with that variable
  - E.g.,  $(\forall x) \text{dolphin}(x) \rightarrow \text{mammal}(x)$
- Existential quantification
  - $(\exists x)P(x)$  means that  $P$  holds for **some** value of  $x$  in the domain associated with that variable
  - E.g.,  $(\exists x) \text{mammal}(x) \wedge \text{lays-eggs}(x)$
  - Permits one to make a statement about some object without naming it

Language	Propositional logic	First–order logic
Syntax	The world contains facts	The world contains objects, relations, and functions
Semantics	$\alpha \wedge \beta$ is true in a world iff $\alpha$ is true and $\beta$ is true (etc.)	$\phi(\sigma)$ is true in a world if $\sigma = o_j$ and $\phi$ holds for $o_j$ ; etc.

# Outline

- Propositional Logic and First–Order Logic
- Inference/Theorem Proving: Forward and Backward Chaining
- Application in NLP



**We are Here**



# Forward Chaining

- Start with given proposition symbols (atomic sentence)
  - e.g., A and B
- Iteratively try to infer truth of additional proposition symbols
  - e.g.,  $A \wedge B \Rightarrow C$ , therefor we establish C is true
- Continue until
  - — no more inference can be carried out, or
  - — goal is reached



# Forward Chaining Algorithm

```
function PL-FC-ENTAILS?(KB, q) returns true or false
    count ← a table, where count[c] is the number of symbols in c's
    premise
    inferred ← a table, where inferred[s] is initially false for all s
    agenda ← a queue of symbols, initially symbols known to be true
    in KB
    while agenda is not empty do
        p ← Pop(agenda)
        if p = q then return true
        if inferred[p] = false then
            inferred[p] ← true
            for each clause c in KB where p is in c.premise do
                decrement count[c]
                if count[c] = 0 then add c.conclusion to agenda
    return false
```

# Backward Chaining

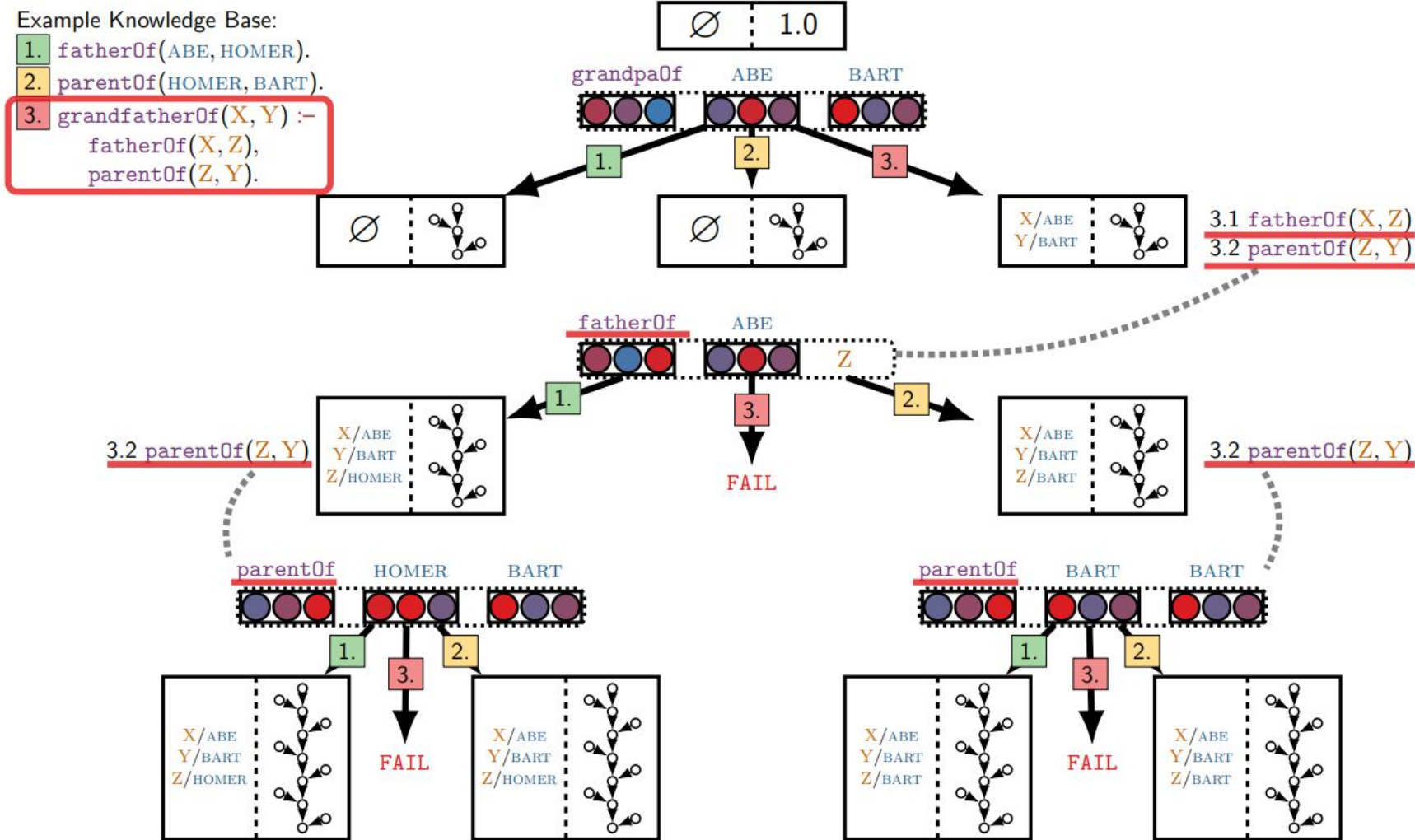
- Idea: work backwards from the query  $Q$ :
  - to prove  $Q$  by BC,
    - check if  $Q$  is known already, or
    - prove by BC all premises of some rule concluding  $q$
- Avoid loops: check if new subgoal is already on the goal stack
- Avoid repeated work: check if new subgoal
  - 1. has already been proved true, or
  - 2. has already failed

# Outline

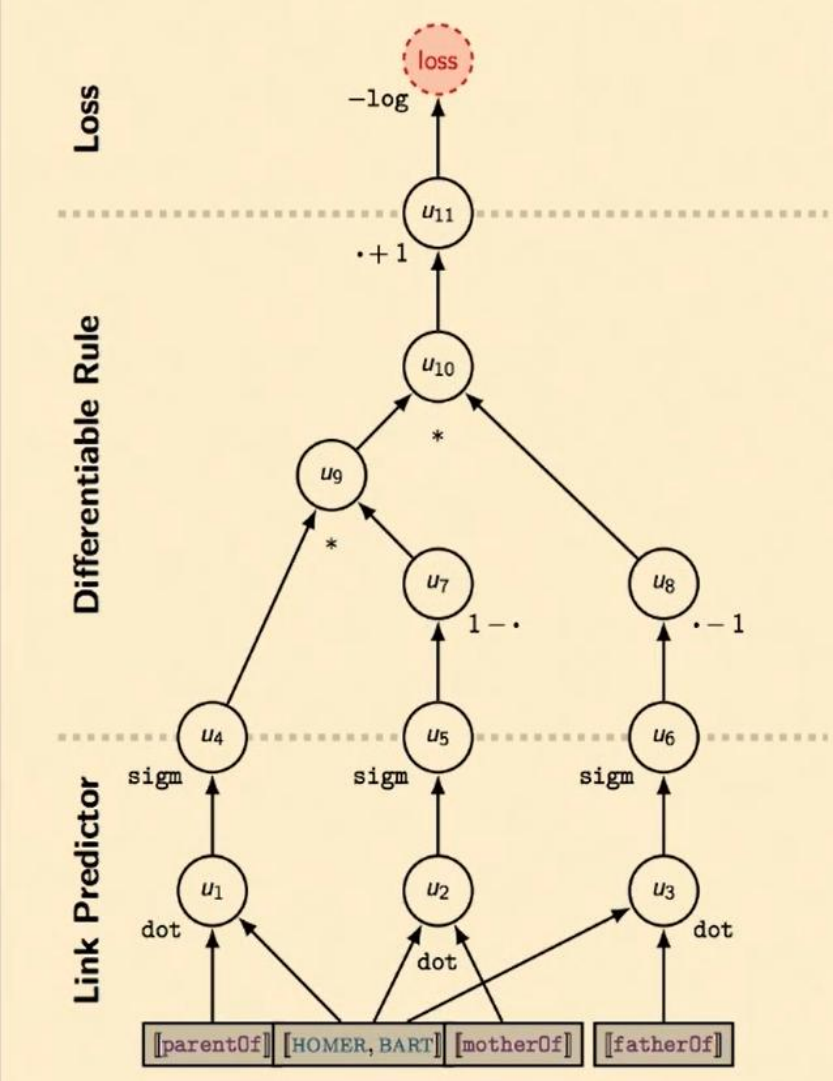
- Propositional Logic and First–Order Logic
- Inference/Theorem Proving: Forward and Backward Chaining
- Application in Knowledge Base Completion
  - Neural Backward Chaining
  - Logic as constraints



# Neural Backward Chaining

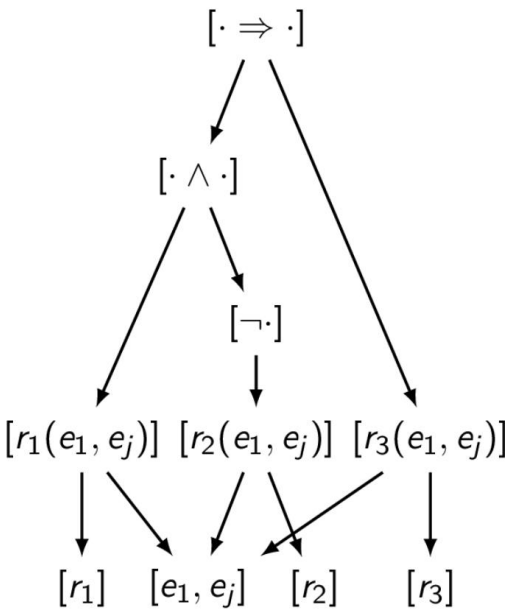


# Logic as Constraints



$$[\mathcal{F}] = \begin{cases} \sigma(\mathbf{v}_s \cdot \mathbf{v}_{ij}) & \text{if } \mathcal{F} = r_s(e_i, e_j), \text{ i.e., facts} \\ 1 - [\mathcal{A}] & \text{if } \mathcal{F} = \neg \mathcal{A} \\ [\mathcal{A}] * [\mathcal{B}] & \text{if } \mathcal{F} = \mathcal{A} \wedge \mathcal{B} \end{cases}$$

$$r_1(e_i, e_j) \wedge \neg r_2(e_i, e_j) \Rightarrow r_3(e_i, e_j)$$



# Performance V.S. Interpretability

- Neural Backward Chaining
  - Good interpretability, limited scope of application (e.g., completion on structured KB)
- Regularizing Neural Models
  - Good performance with neural models as backbone, limited interpretability



# Neural–Symbolic Models in NLP

# Symbolic Language

<b>First–Order Logic</b>	<i>Every prime greater than two is odd.</i>	$\forall x.\text{prime}(x) \wedge \text{more}(x, 2) \rightarrow \text{odd}(x)$
<b>Lambda Calculus</b>	<i>How many primes are less than 10?</i>	$\text{count}(\lambda x.\text{prime}(x) \wedge \text{less}(x, 10))$
<b>Lambda DCS</b>	<i>How many primes are less than 10?</i>	$\text{count}(\text{prime} \sqcap (\text{less}.10))$

- First–order logic fails to construct a set and manipulating it.
- The  $\lambda$  operator can be thought of as constructing a set of all  $x$  that satisfy the condition; in symbols,  $[\lambda x. f(x)]_c = \{x : [f(x)]_c = \text{true} \}$ .

$\text{prime}$

$\text{less than}$

$10$

$$\frac{\overline{N[\lambda x.\text{prime}(x)]} \quad \frac{\overline{(N \setminus N)/NP[\lambda y.\lambda f.\lambda x.f(x) \wedge \text{less}(x, y)]} \quad \overline{NP[10]}}{\overline{N \setminus N[\lambda f.\lambda x.f(x) \wedge \text{less}(x, 10)]}}_{(>)} \quad \overline{N[\lambda x.\text{prime}(x) \wedge \text{less}(x, 10)]}}_{(<)}$$

Lambda Calculus

$\text{prime}$

$\text{less than}$

$10$

$$\frac{\overline{N[\text{prime}]} \quad \frac{\overline{N[N[\text{less}]]} \quad \overline{N[10]}}{\overline{N[\text{less}.10]}}_{(\text{join})} \quad \overline{N[\text{prime} \sqcap \text{less}.10]}}_{(\text{intersect})}$$

Lambda DCS

# Semantic Parsing

map an utterance  $x$  in a context  $c$  to an action  $y$

$x$ : *What is the largest prime less than 10?*

$c$ : **primes** : {2, 3, 5, 7, 11, ...}



$y$ : 7

# A Derivation

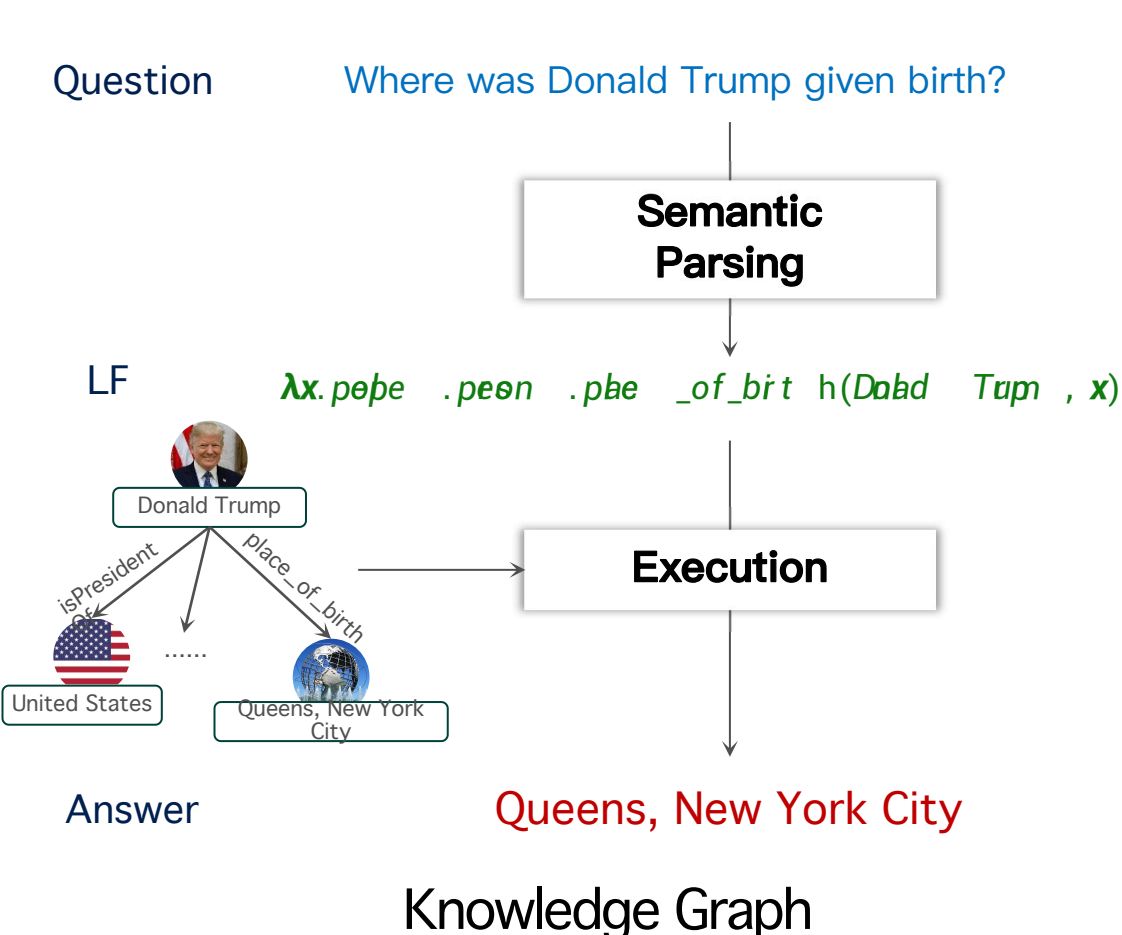
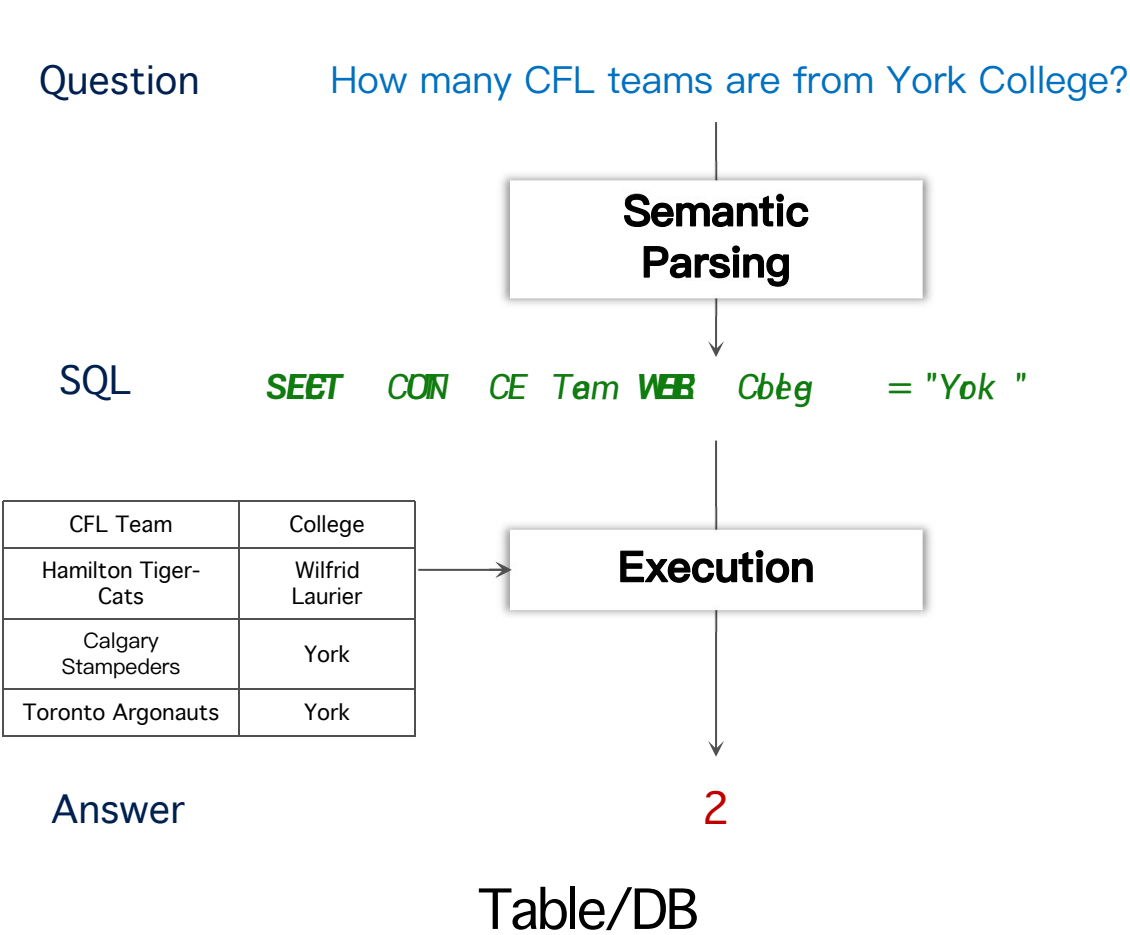
*What is the largest prime less than 10 ?*

$\text{NP}[\mathbf{primes}]$		$\text{NP}[10]$
$\text{QP}[(-\infty, 10)]$		
$\text{NP}[\mathbf{primes} \cap (-\infty, 10)]$		
$\text{NP}[\max(\mathbf{primes} \cap (-\infty, 10))]$		
$\text{ROOT}[\max(\mathbf{primes} \cap (-\infty, 10))]$		

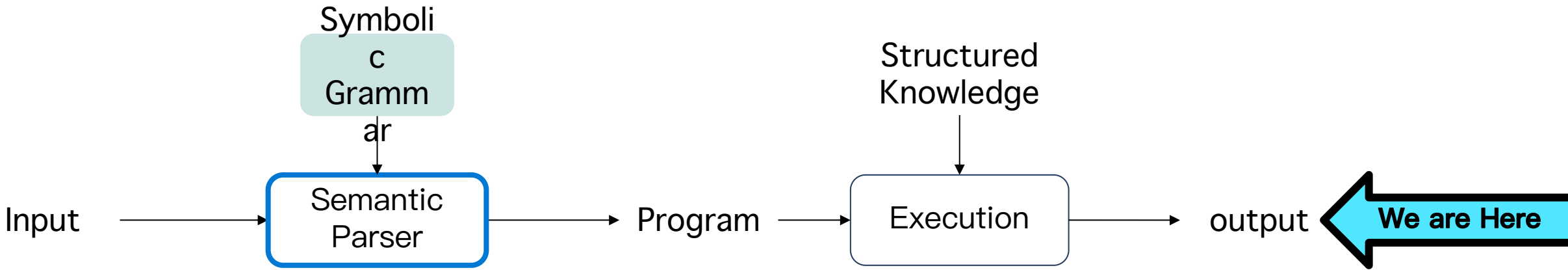
(R1)	<i>prime</i>	$\Rightarrow$	$\text{NP}[\mathbf{primes}]$
(R2)	<i>10</i>	$\Rightarrow$	$\text{NP}[10]$
(R3)	<i>less than</i> $\text{NP}[z]$	$\Rightarrow$	$\text{QP}[(-\infty, z)]$
(R4)	$\text{NP}[z_1]$ $\text{QP}[z_2]$	$\Rightarrow$	$\text{NP}[z_1 \cap z_2]$
(R5)	<i>largest</i> $\text{NP}[z]$	$\Rightarrow$	$\text{NP}[\max(z)]$
(R6)	<i>largest</i> $\text{NP}[z]$	$\Rightarrow$	$\text{NP}[\min(z)]$
(R7)	<i>What is the</i> $\text{NP}[z]?$	$\Rightarrow$	$\text{ROOT}[z]$

# Semantic Parsing

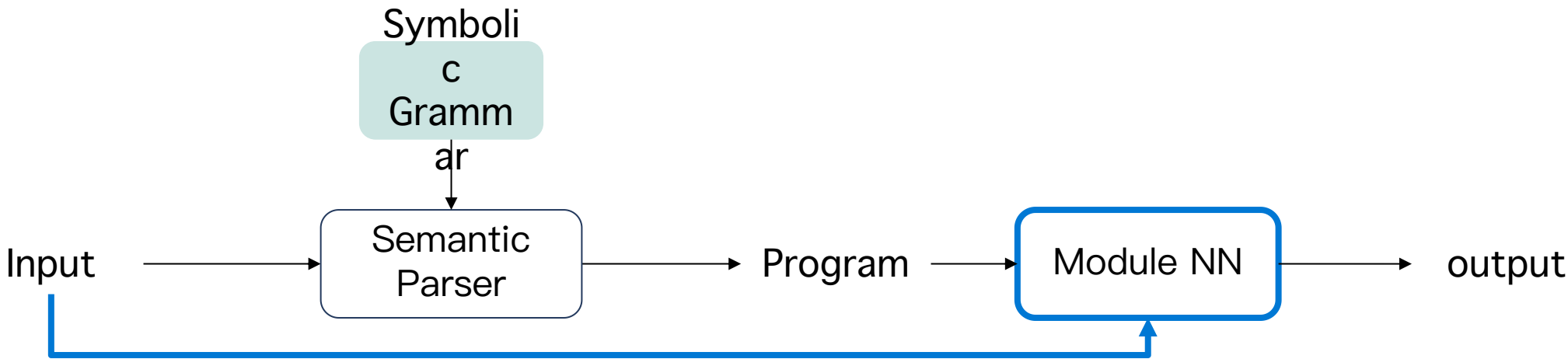
- Map Natural Language into machine executable logical forms



# Outline



**Part 1: Semantic Parsing**



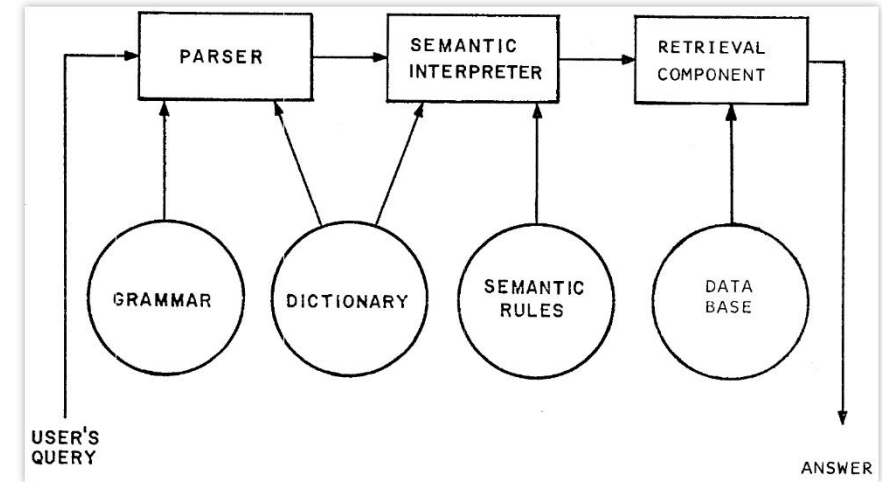
**Part 2: Module Network**

# LSNLIS (Lunar Sciences Natural Language Information System)

A question–answering system to enable a lunar geologist to conveniently access, compare, and evaluate the chemical analysis data on lunar rock and soil composition that is accumulating as a result of the Apollo moon missions.

Two DB files. One is a 13,000 line table of chemical and age analysis of the Apollo 11 samples extracted from the reports of a the First Annual Lunar Science Conference, and the second is a keyphrase index to those reports.

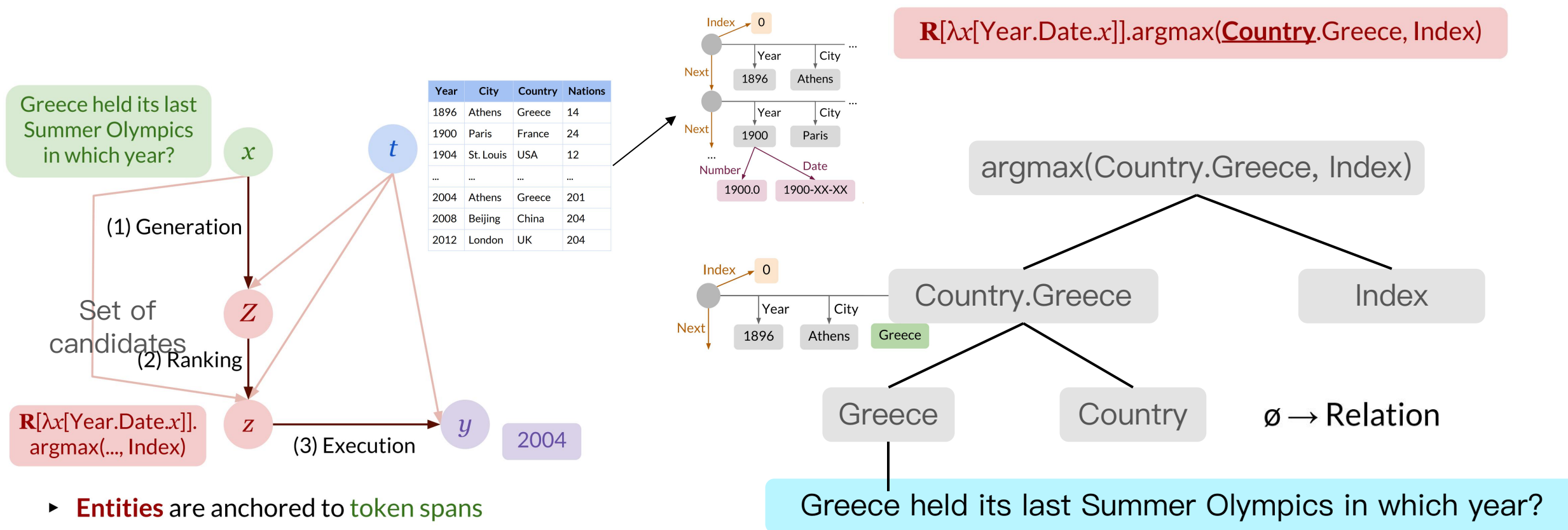
A dictionary about 3500 words.



1. List the rocks which contain chromite and ulvospinel.
2. Give me all references on fayalitic olivine.
3. What minerals have been identified in the lunar samples?
4. What analyses of olivine are there?
5. What is the average analysis of Ir in rock S10055?
6. List the modes for all low Rb rocks.
7. Give me the K / Rb ratios for all lunar samples.
8. Has the mineral analcite been identified in any lunar sample?
9. What is the concentration of La in rock S10034?
10. Identify all samples in which glass was found.
11. Give me all modal analyses of lunar fines.
12. In what samples has apatite been identified?

# Floating Parser

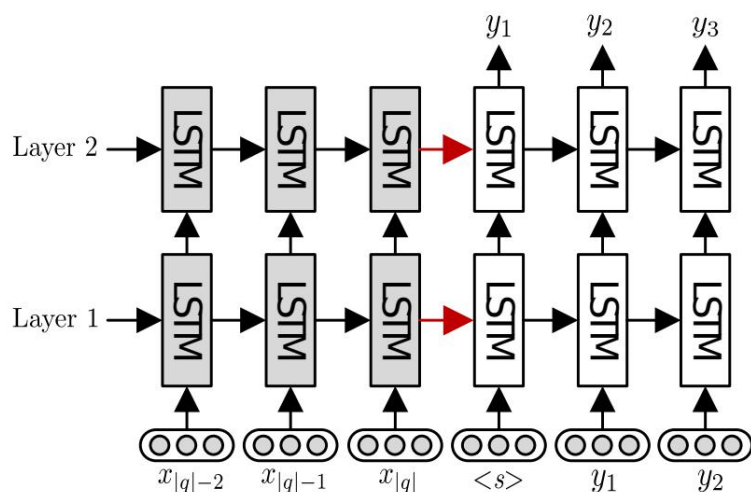
- Build formulas bottom-up according to a set of deduction rules
- Allow formulas to be created from nothing ("floating")



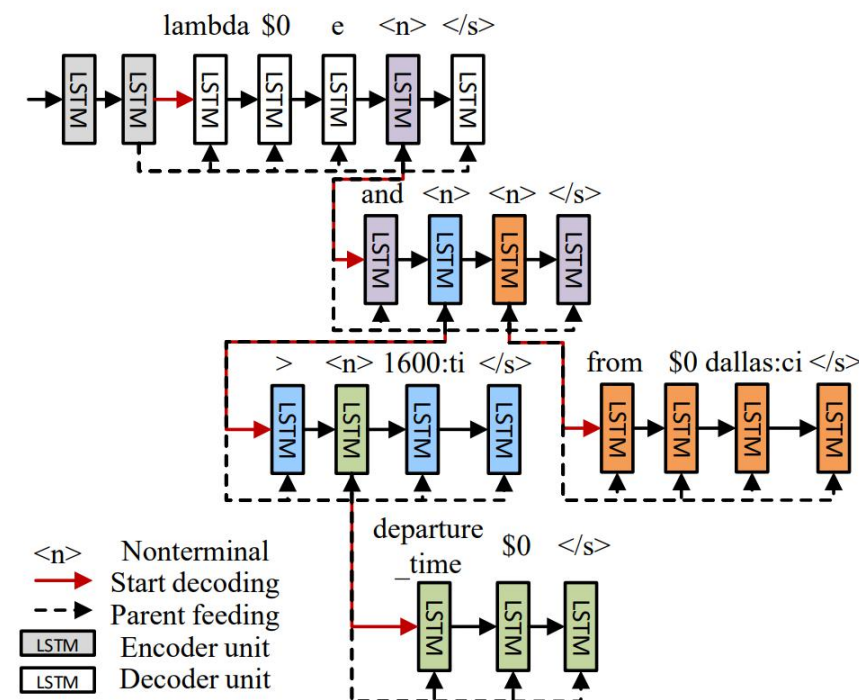
- **Entities** are anchored to **token spans**
- **Relations** and **Operations** are not



# Seq2Seq and Seq2Tree

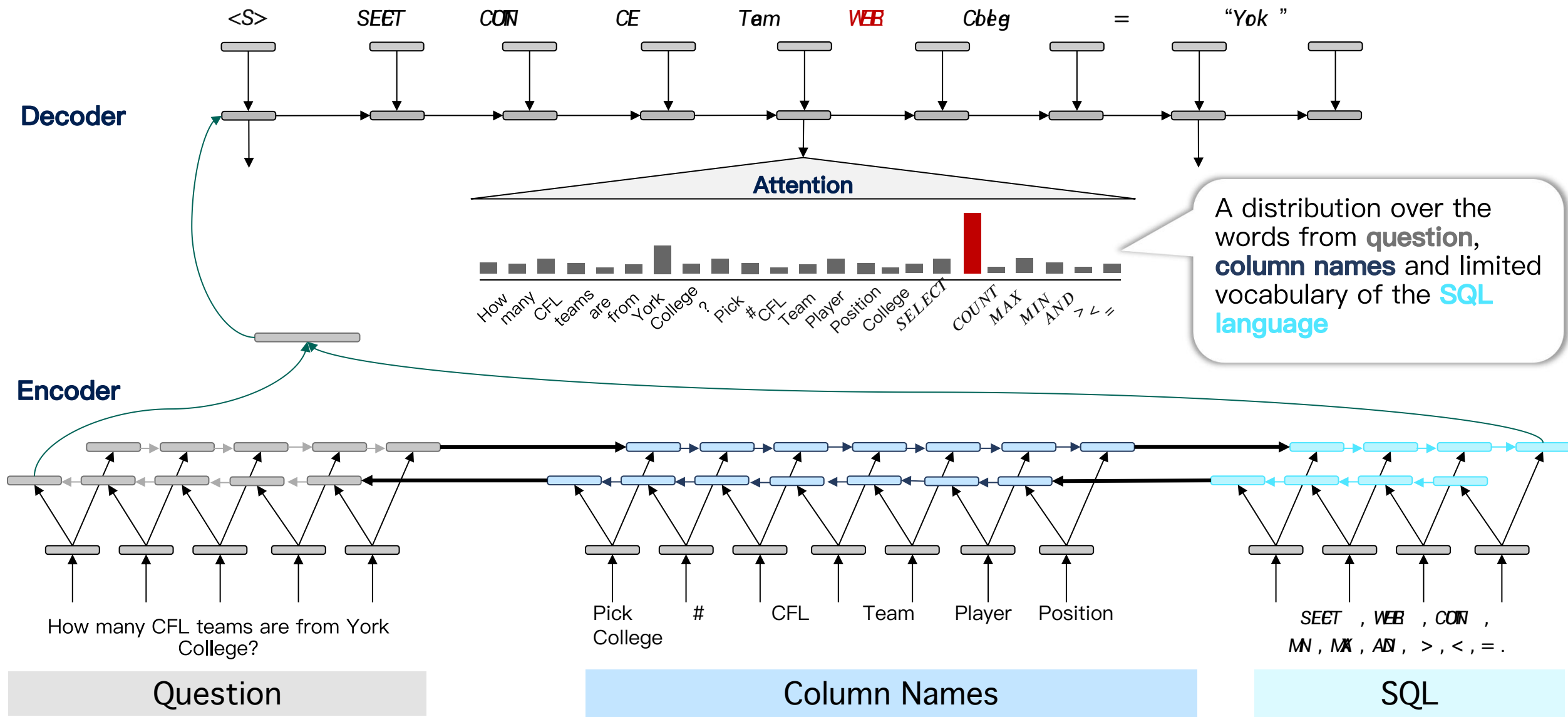


Seq2Seq with 2-layer RNN

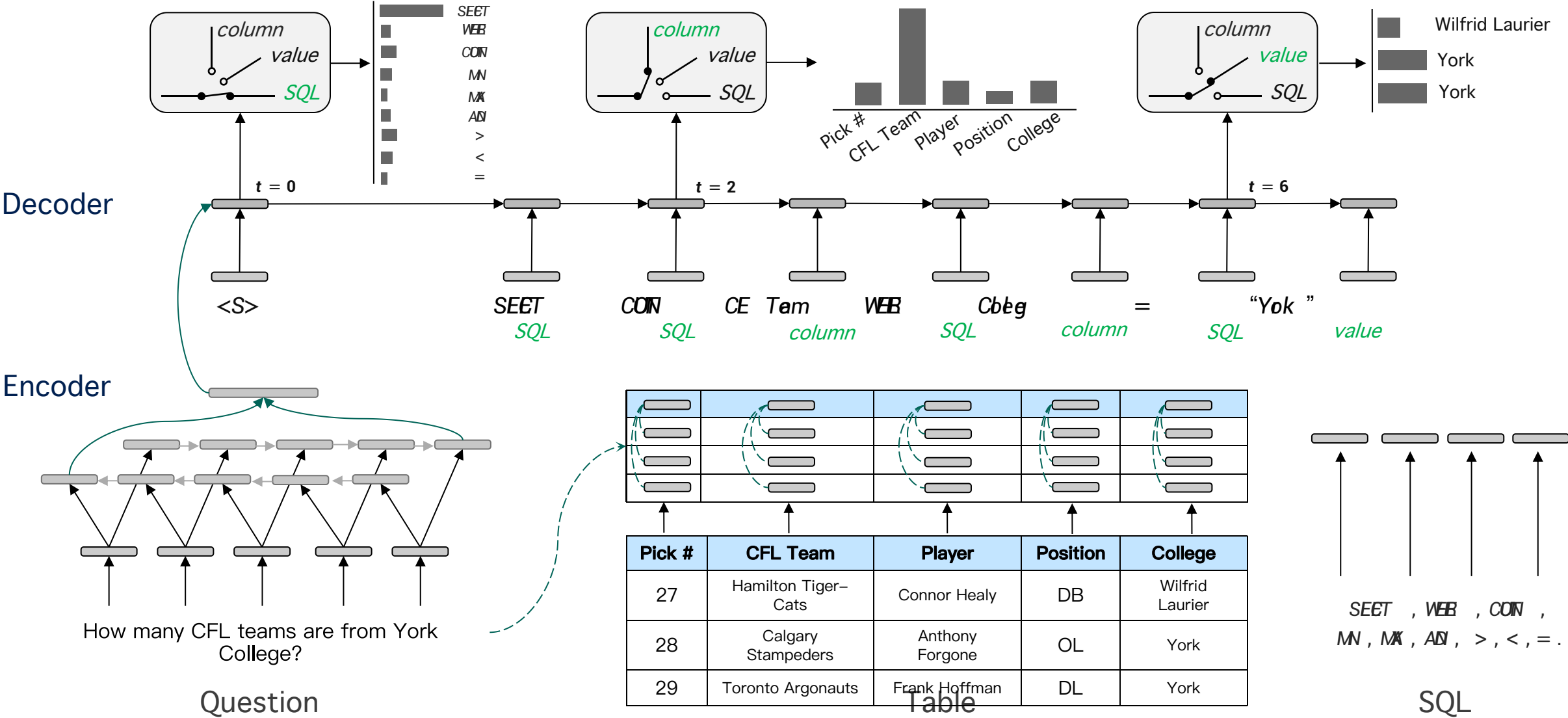


Seq2Tree with a hierarchical tree decode

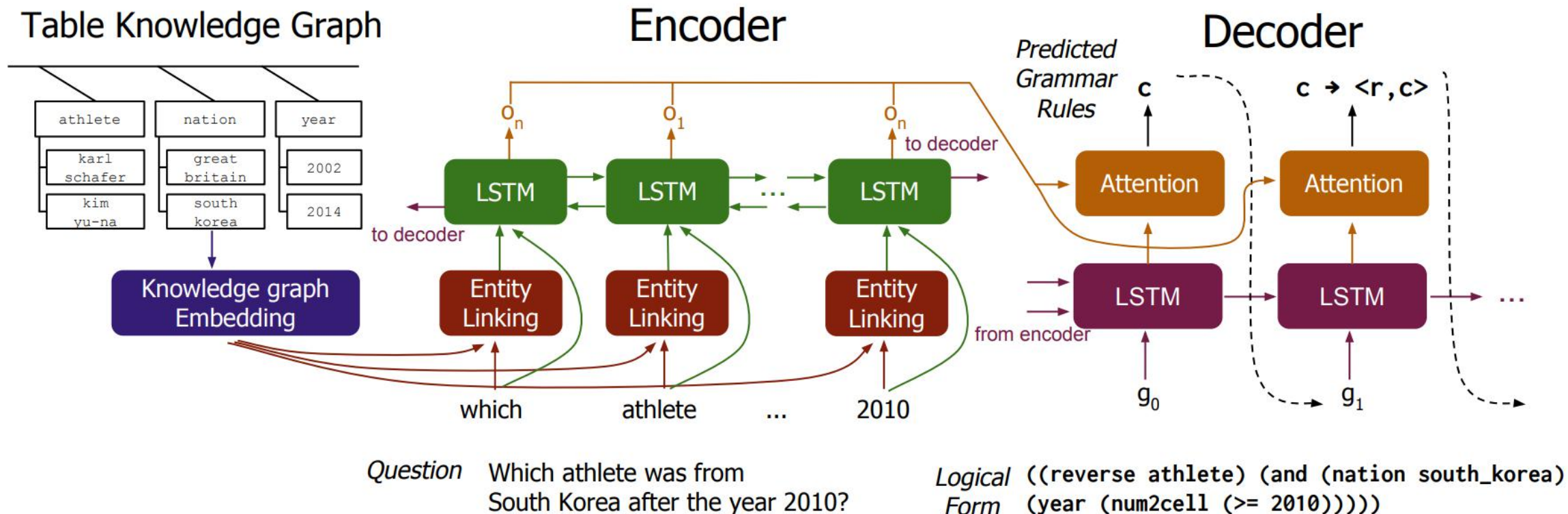
# Seq-to-Seq with Pointer Network



# Seq-to-Seq with Structural Decoding



# Seq-to-Seq with Typed Constrained Decoding



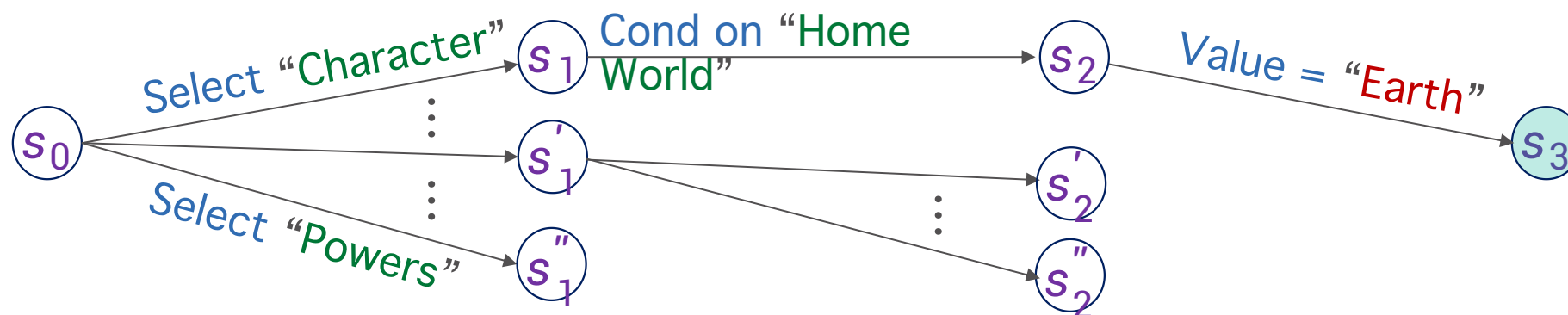
# Action and Module

➤ *Which super heroes came from Earth and first appeared after 2009?*

```
SELECT      Character
WHERE {Home World = Earth} ∧
      {First Appeared > 2009}
```

Legion of Super Heroes Post-Infinite Crisis			
Character	First Appeared	Home World	Powers
Night Girl	2007	Kathoon	Super strength
Dragonwing	2010	Earth	Fire breath
Gates	2009	Vyrge	Teleporting
XS	2009	Aarok	Super speed
Harmonia	2011	Earth	Elemental

$Q = \text{"Which super heroes came from Earth?"}$ ,  $A^* = \{\text{Dragonwing, Harmonia}\}$



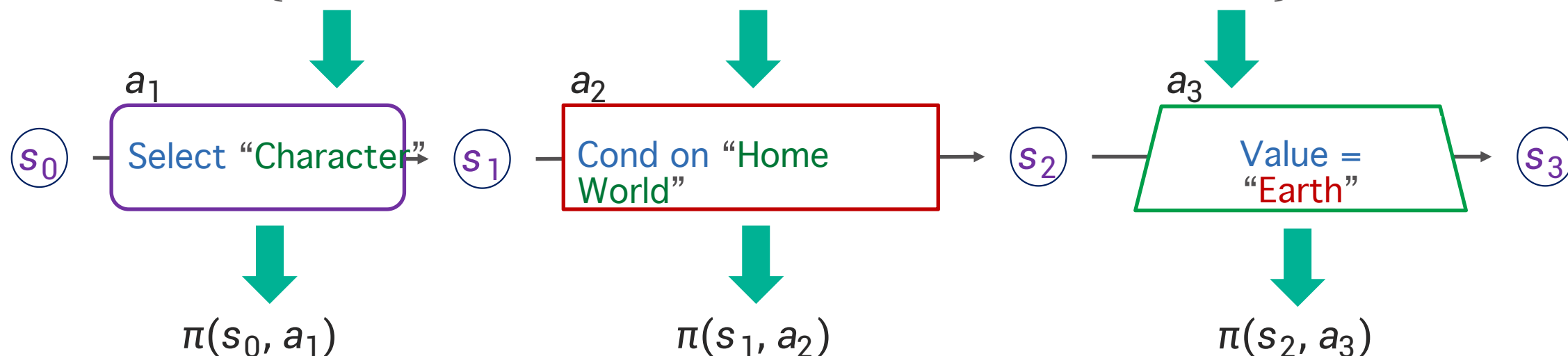
# Action and Module

- The goodness of a state:  $V(s_t) = V(s_{t-1}) + \pi(s_{t-1}, a_t)$ ,  $V(s_0) = 0$
- Value of  $\pi(s, a)$  is determined by a neural-network model
- Actions of the same type (e.g., `select-column`) share the same neural-network module

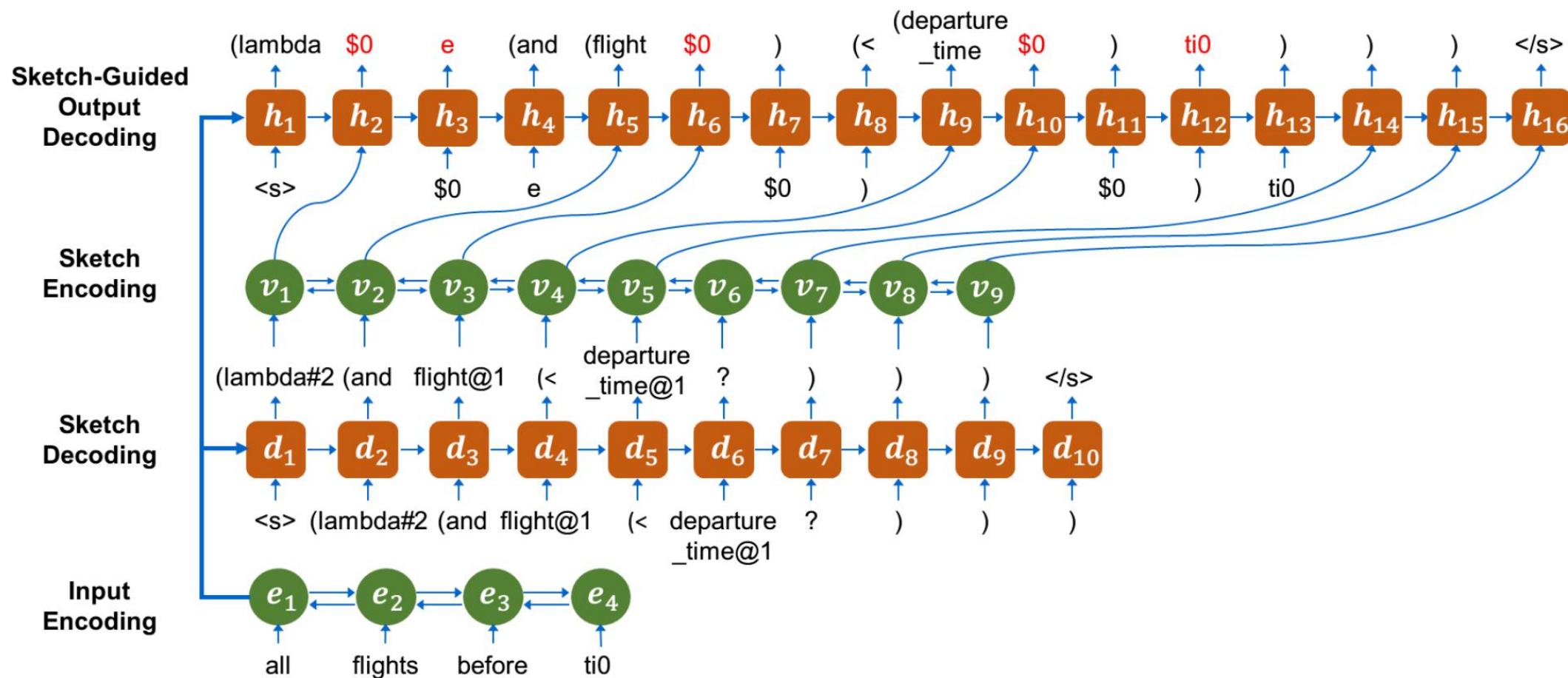
➤ *Which super heroes came from Earth?* ,



	A	B	C
1			
2			
3			
4			
5			
6			



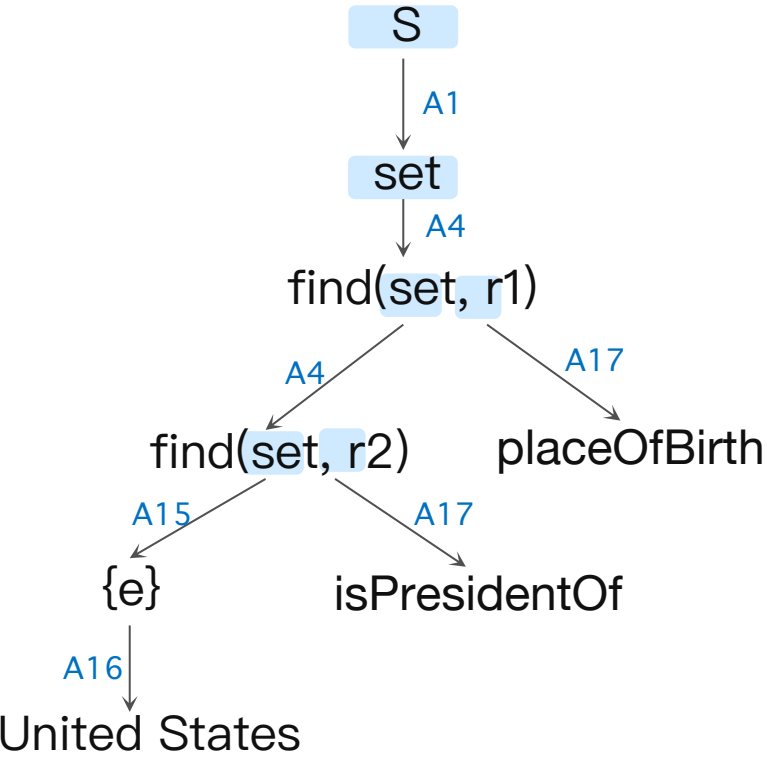
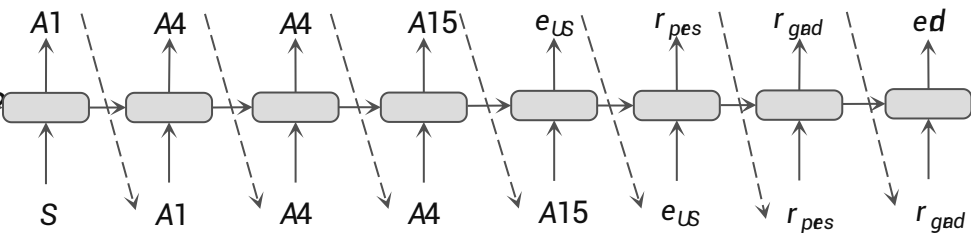
# Coarse-to-Fine Decoding





# KBQA with Semantic Parsing (single-turn)

Where was the president of the United States born?

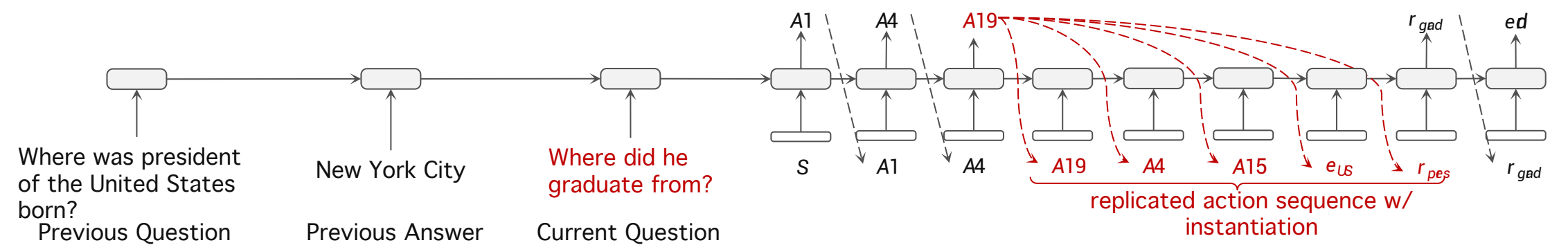


- A1:  $S \rightarrow s\mathbf{e}$
- A4:  $s\mathbf{e} \rightarrow find(s\mathbf{e}, r1)$
- A4:  $s\mathbf{e} \rightarrow find(s\mathbf{e}, r2)$
- A15:  $s\mathbf{e} \rightarrow \{e\}$
- A16:  $e \rightarrow United States$
- A17:  $r2 \rightarrow isPresidentOf$
- A17:  $r1 \rightarrow placeOfBirth$

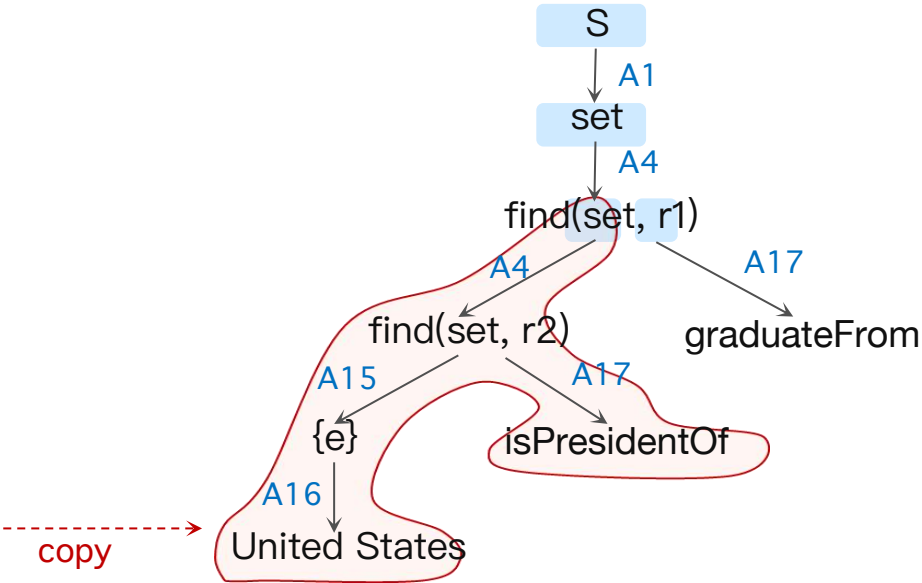
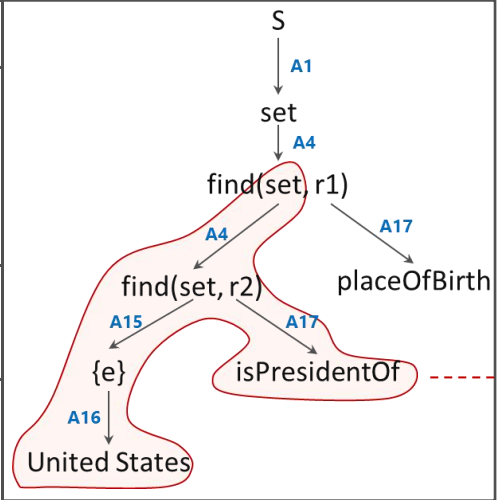
Action	Operation	Description
A1-A3	$S \rightarrow Set \mid Num \mid Bool$	S is start symbol
A4	$Set \rightarrow Find(R, E)$	Set of entities with a r edge to e
A5	$Num \rightarrow Count(Set)$	Total number of set
A6	$Bool \rightarrow (\in, E, Set)$	Whether $E \in Set$
A7	$Set \rightarrow Set \cup Set$	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	$Set \rightarrow arxmax(set, r, num)$	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



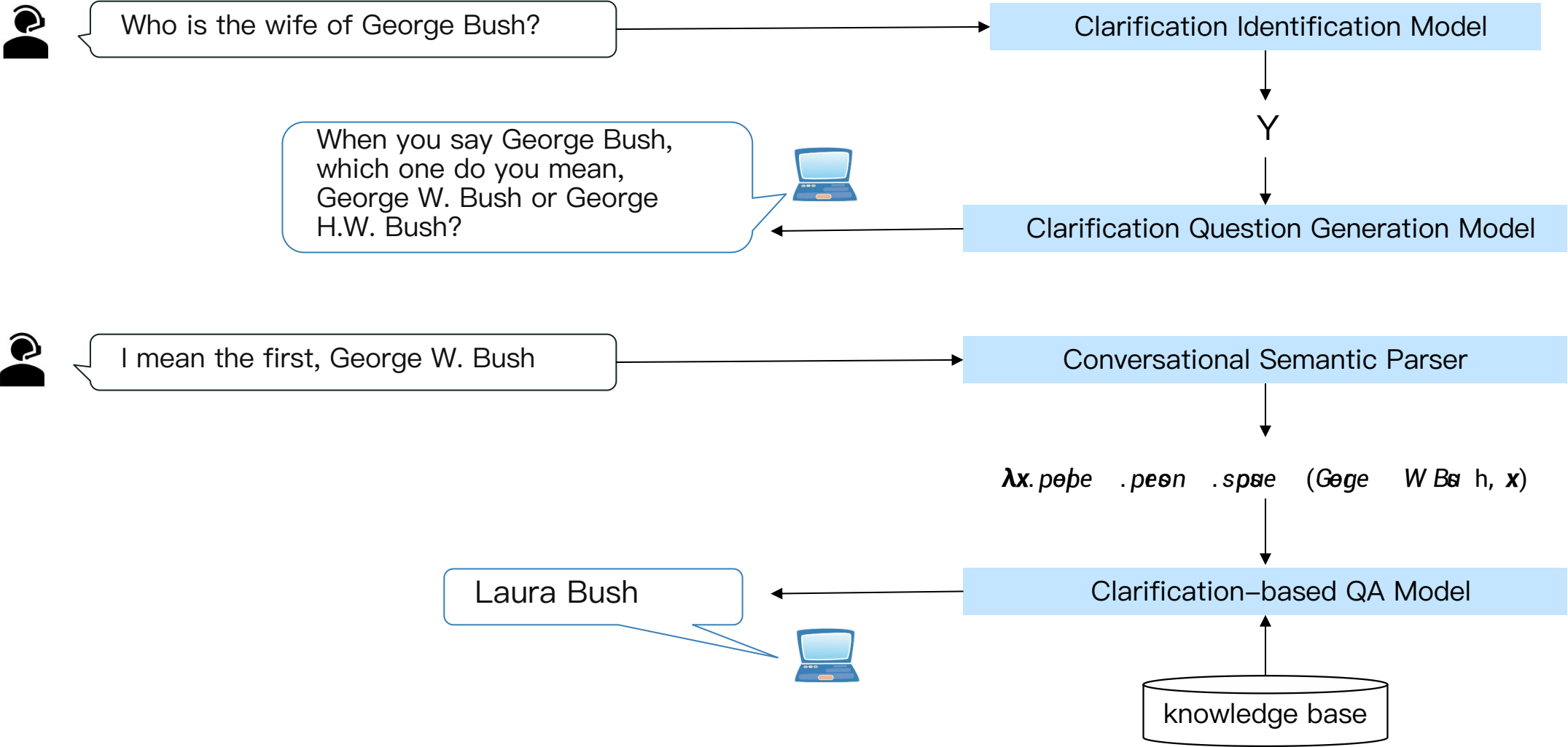
# KBQA with Semantic Parsing (multi-turn)



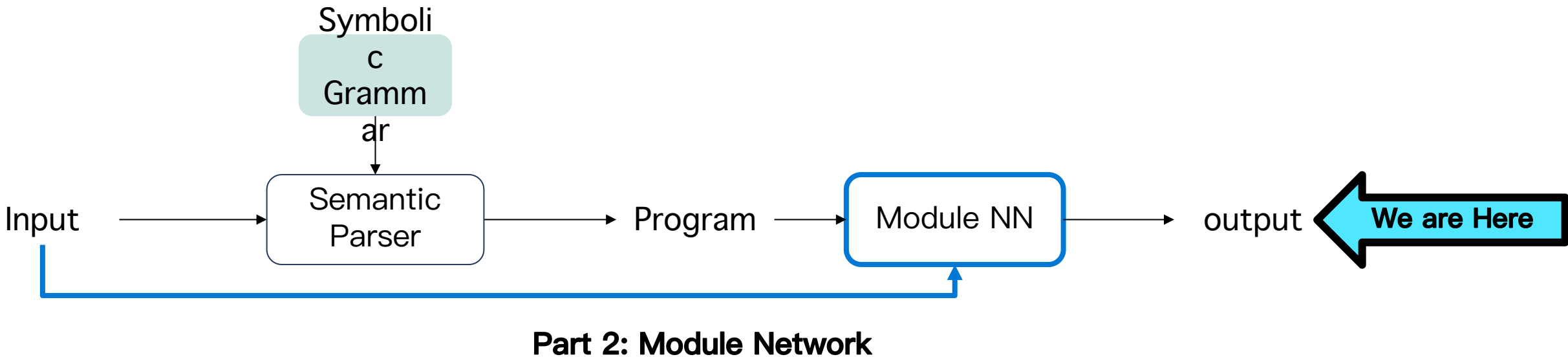
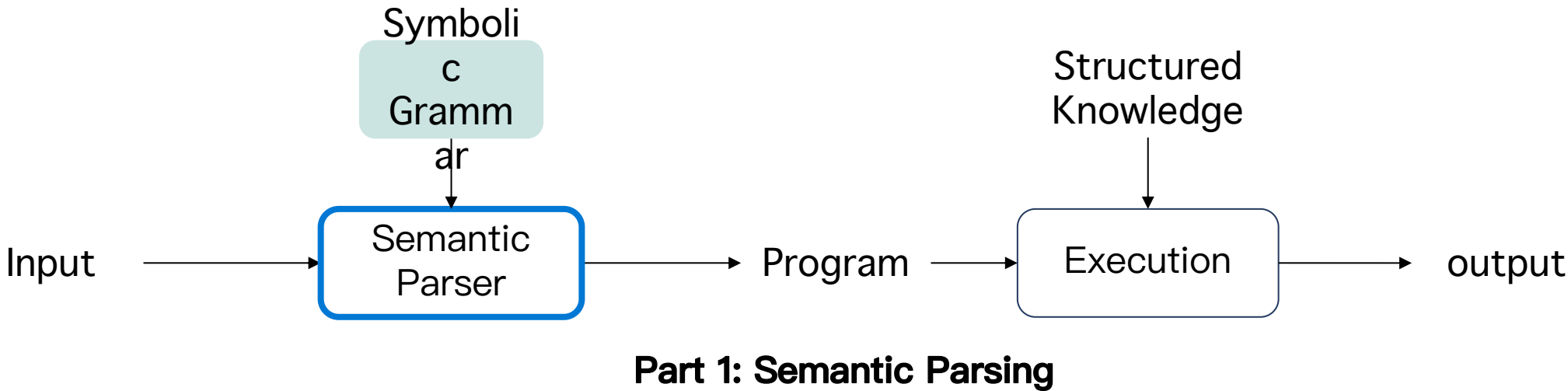
Dialog Memory	
Entity	{United States, tag=utterance} {New York City, tag=answer}
Predicate	{isPresidentOf} {placeOfBirth}
Action Subsequence	$se \rightarrow A4 A15 e_{US} r_{pes}$ $se \rightarrow A4 A15$ $se \rightarrow A4 A4 A15 e_{US} r_{pes} r_{bth}$ $se \rightarrow A4 A4 A15$



# Conversational Question Answering

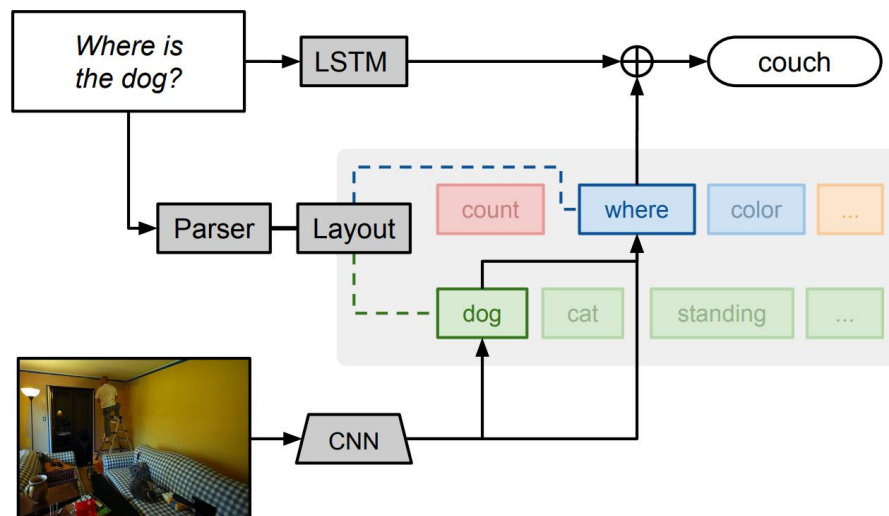


# Outline

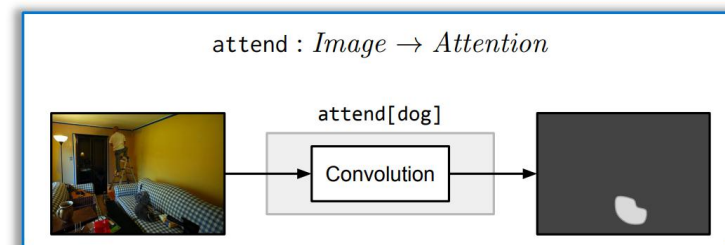


# Module Network

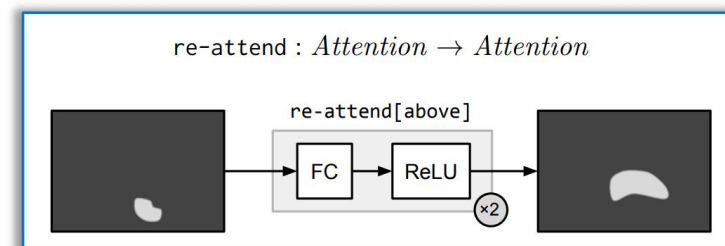
## Reusable neural modules with different architectures



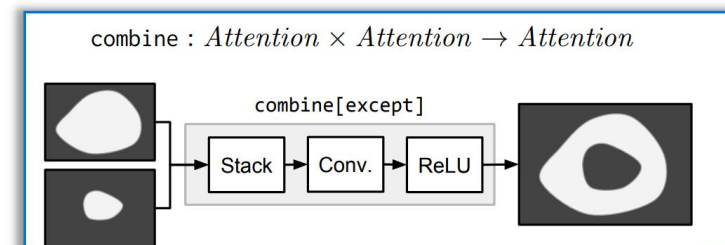
### Attention



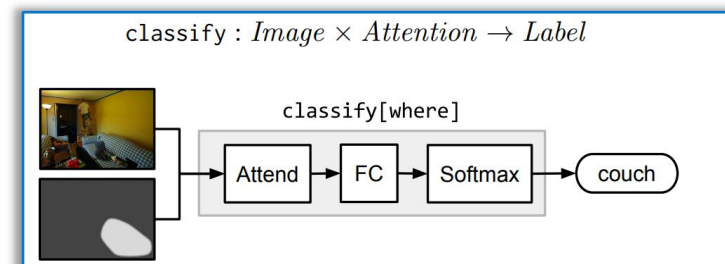
### Re-Attention



### Combination



### Classification



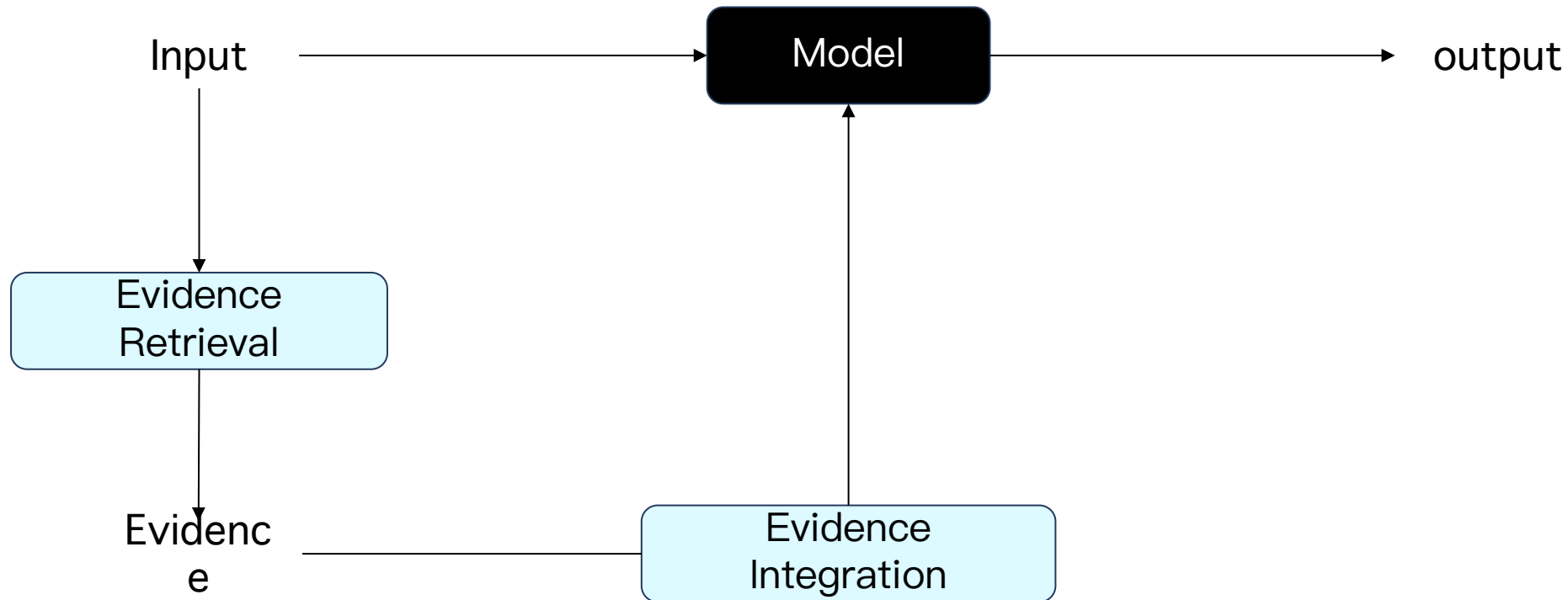
# Performance V.S. Interpretability

- Semantic Parser
  - Good interpretability, good performance on limited applications
  - The extension of grammar to open domain is challenging
- Module Network
  - Good performance with neural models as backbone, limited interpretability
  - Moderate interpretability, compared with semantic parser.
  - The definition of grammar is typically task-specific, and manually designed by experts

# Evidence-based Models in NLP

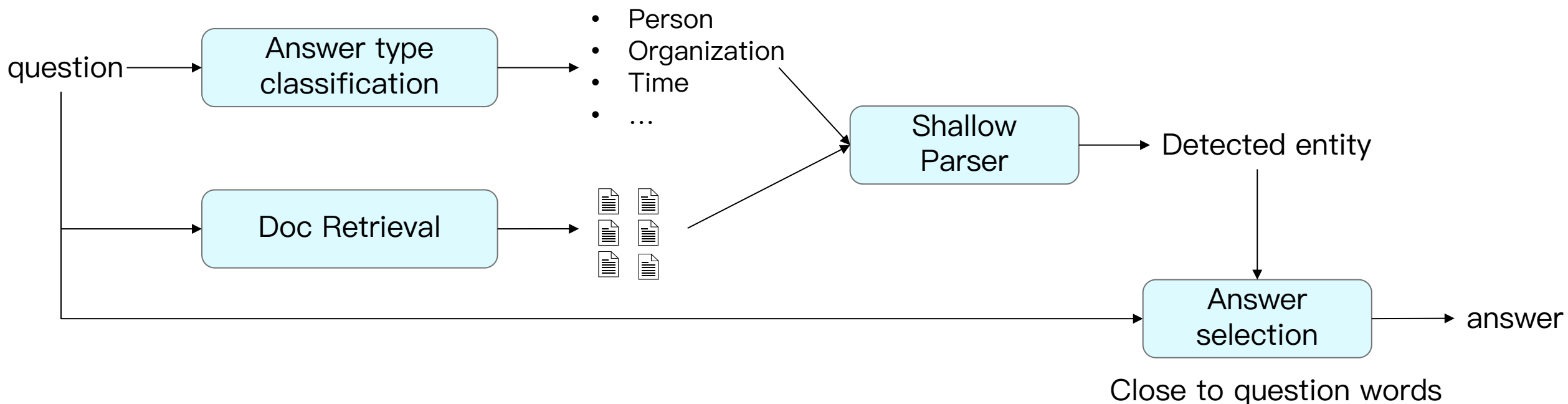
# General Framework

- Consider evidence as an additional input of the model



# Open Question Answering

- 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.

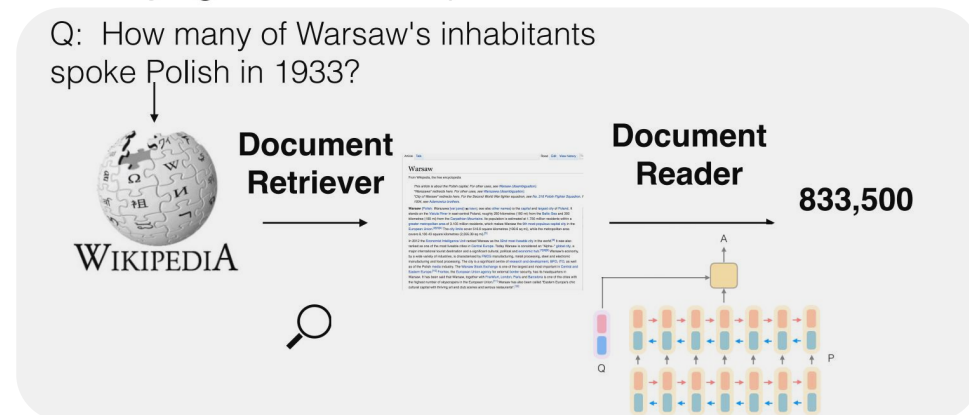




# Sparse Retrieval Model (DrQA)

## 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  Freebase (Berant et al, 2013)
  - WikiMovies (Miller et al, 2016)

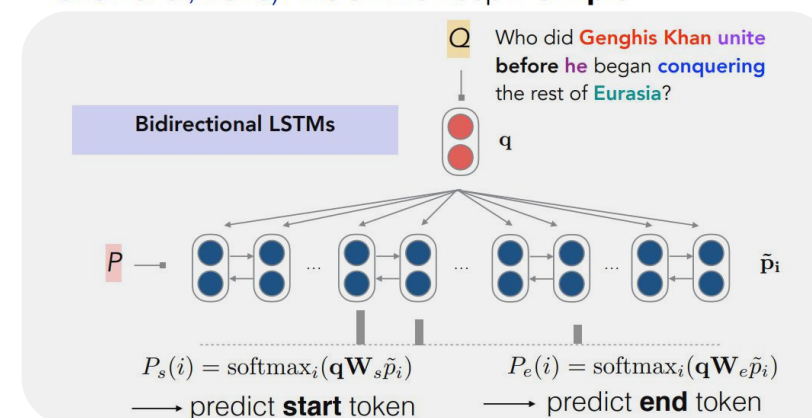
## Document Retriever

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**!



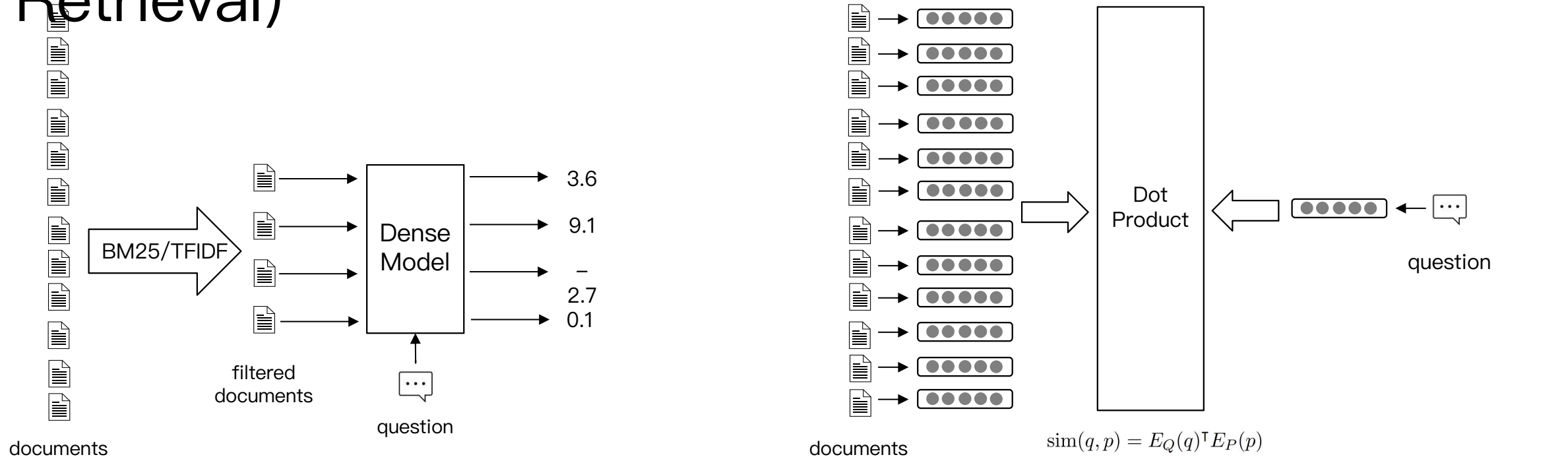
The input vectors consist of:

- 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..

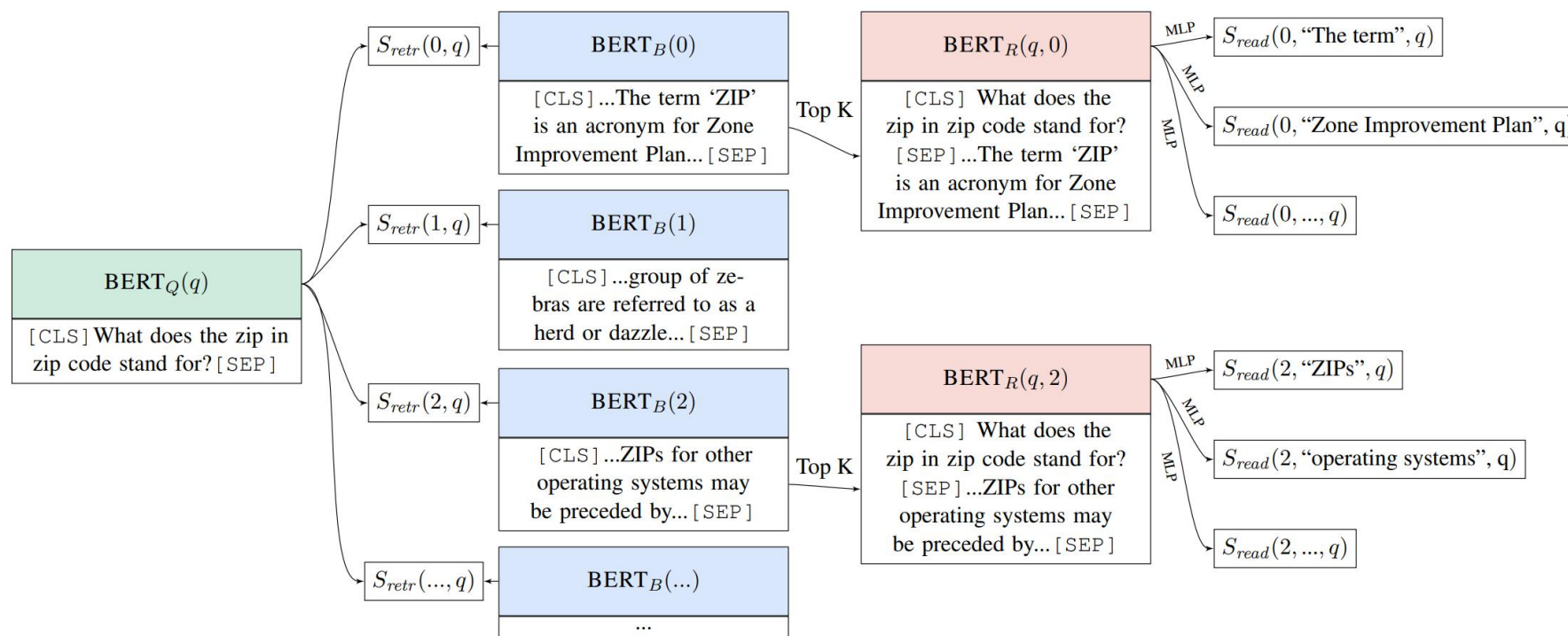
# DPR Results

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	-	-	-	<b>56.5</b>
Single	REALM <sub>Wiki</sub> (Guu et al., 2020)	39.2	-	40.2	46.8	-
Single	REALM <sub>News</sub> (Guu et al., 2020)	40.4	-	40.7	42.9	-
Single	BM25	32.6	52.4	29.9	24.9	38.1
	DPR	<b>41.5</b>	56.8	34.6	25.9	29.8
	BM25+DPR	39.0	57.0	35.2	28.0	36.7
Multi	DPR	<b>41.5</b>	56.8	<b>42.4</b>	49.4	24.1
	BM25+DPR	38.8	<b>57.9</b>	41.1	<b>50.6</b>	35.8

Table 4: End-to-end QA (Exact Match) Accuracy. The first block of results are copied from their cited papers. REALM<sub>Wiki</sub> and REALM<sub>News</sub> 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.

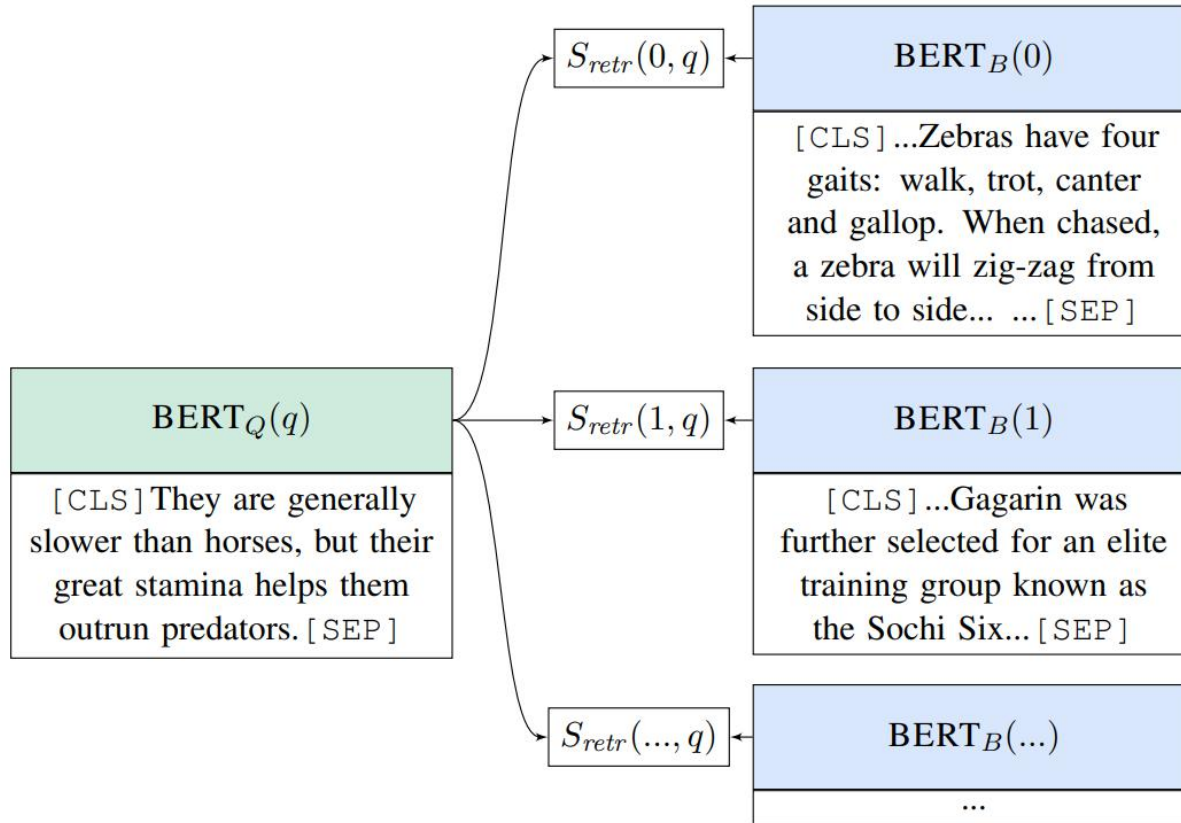
# Joint Retrieval and Reader

- 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.





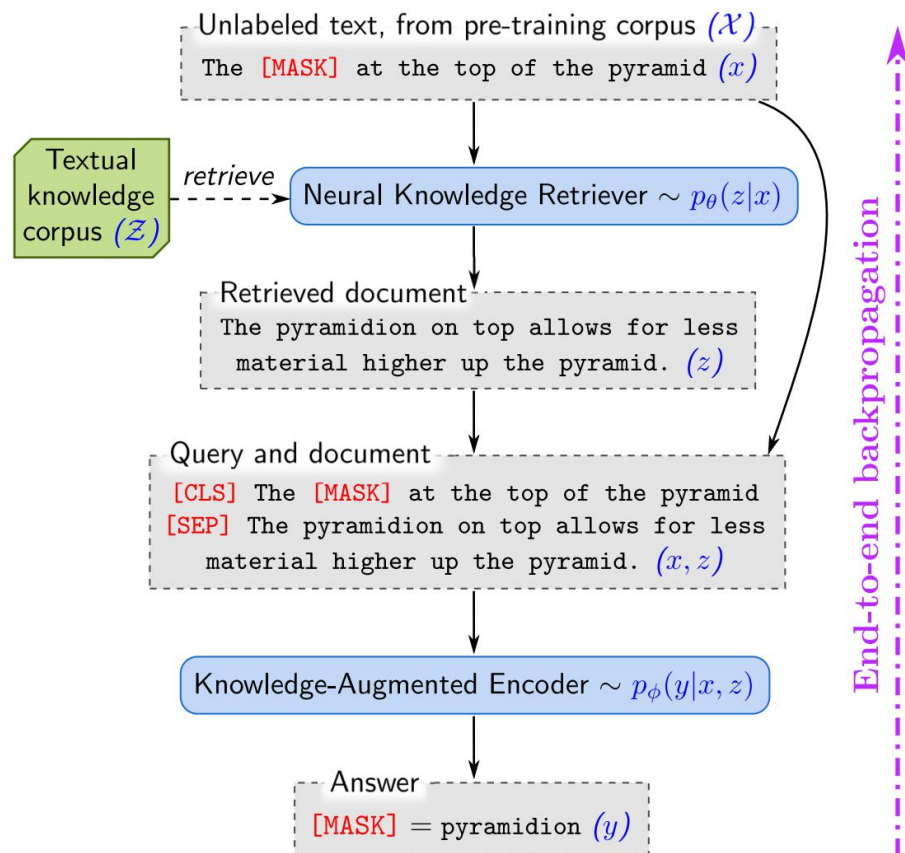
# Pre-train Retrieval Model with Inverse Cloze Task



- 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

$$p(z | x) = \frac{\exp f(x, z)}{\sum_{z'} \exp f(x, z')},$$

$$f(x, z) = \text{Embed}_{\text{input}}(x)^\top \text{Embed}_{\text{doc}}(z),$$

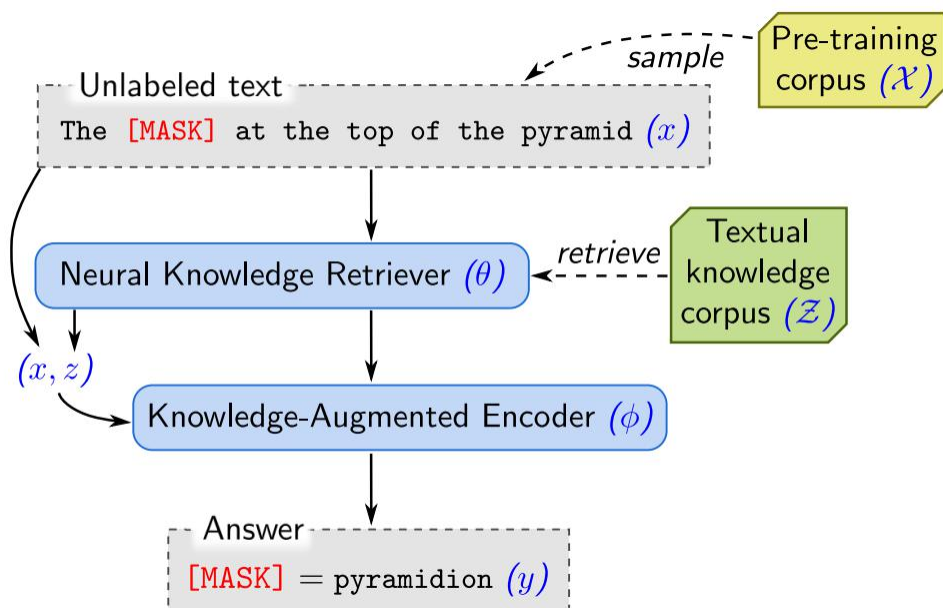
$$\text{Embed}_{\text{input}}(x) = \mathbf{W}_{\text{inputBERTCLS}}(\text{join}_{\text{BERT}}(x))$$

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

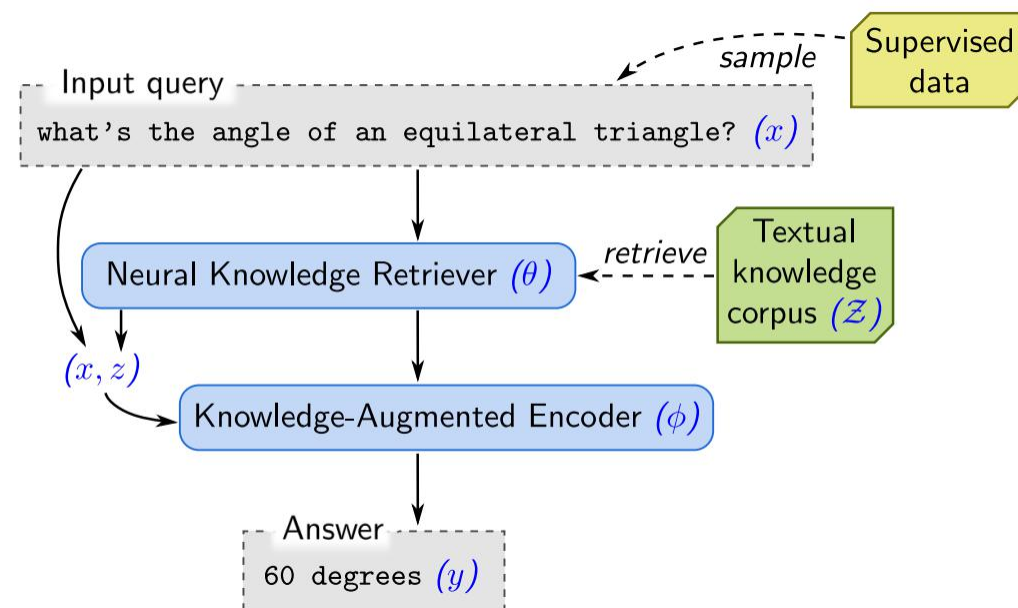
$$\text{join}_{\text{BERT}}(x) = [\text{CLS}] x [\text{SEP}]$$

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

# Pre-train Retrieval Model with REALM



**Unsupervised pre-training**



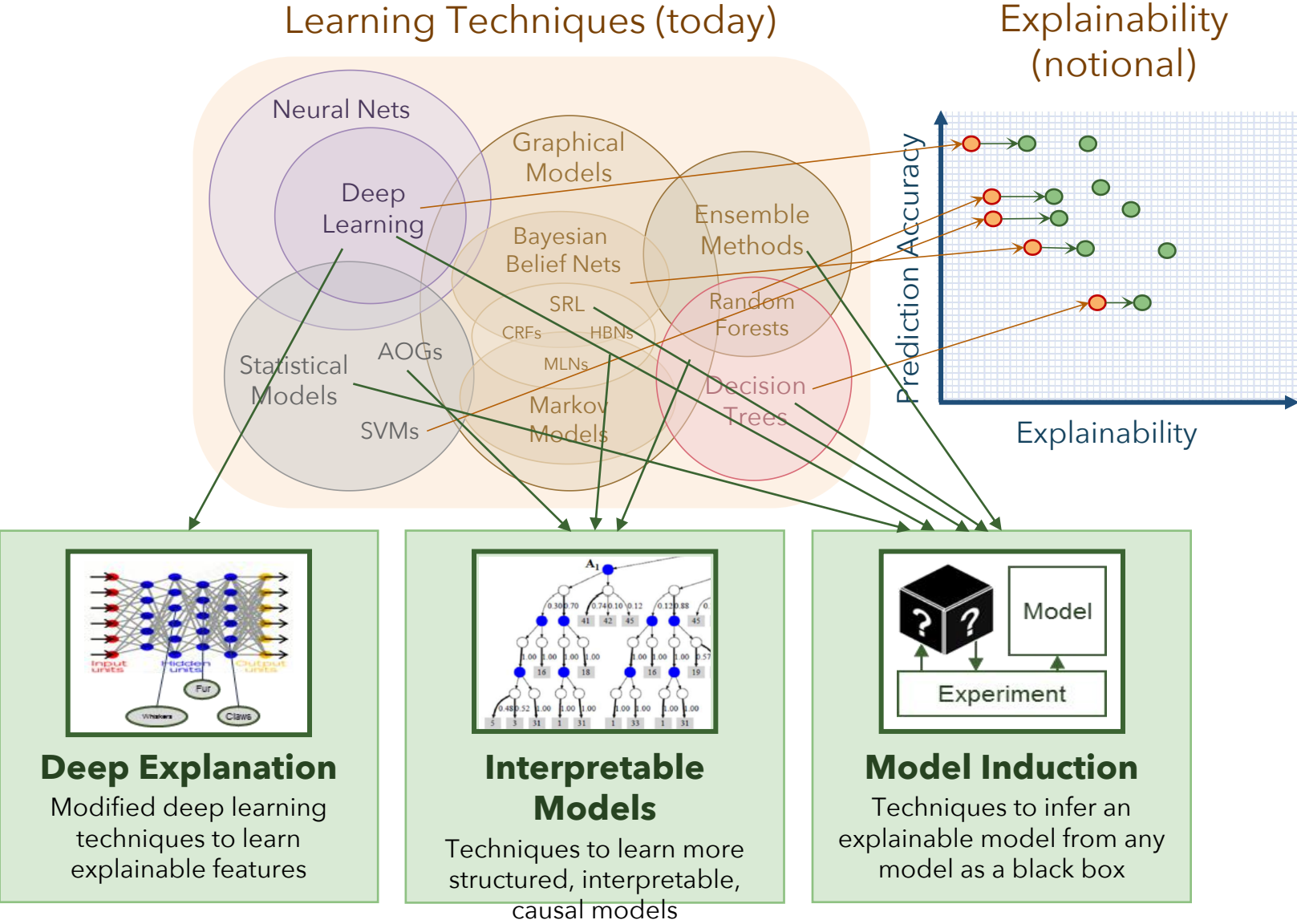
**Supervised fine-tuning**

# Summary

- Topics covered by this talk
  - Logic-based Models in NLP
  - Neural-Symbolic Models in NLP
  - Evidence-based Models in NLP
- Directions worth pursuing
  - Interpretable models and methods
  - Deep understanding with reasoning ability



# Challenge: Performance vs. Explainability



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

# We are hiring!

- Both interns and employees.
- Topics:
  - Semantic Parsing
  - Code Intelligence
  - Machine Reasoning
- Send email to [dutang@microsoft.com](mailto:dutang@microsoft.com)