# Natural Language Learning with Formal Semantics

#### Wenhu Chen

# Outline

- 1. Introduction
- 2. Formal Semantics
- 3. Its Application in NLP
- 4. My Work & Future Direction



Meaning





Linguists can barely reach an agreement on the definition of "meaning"

Raw-String

- 1. Simple
- 2. Straightforward

 $\checkmark$ 

Bag of words Unigram/Bigram





What aspects do you need the meaning representation to capture?

What aspects do you need the meaning representation to capture? Can it lean meaning or just similarity between meaning?

## 1. Raw-String Representation

- Cannot directly link to the external world knowledge.
- Cannot combine pieces of knowledge together to perform inference.

Utterance1: John has an apple Utterance2: John has a banana ?

Utterance: John has an apple and a banana.

# 1. Raw-String Representation

- Assuming we want to build an NLP system to understand math questions?
  - a. Such system cannot leverage the mathematical knowledge.
  - b. The sample complexity becomes extremely high.







Utterance: What is the largest prime less than 10?

- We construct the meaning of natural language in terms of variables, predicates, arguments.

- We construct the meaning of natural language in terms of variables, predicates, arguments.
- Formal Semantics aim to use explicit and grounded meaning representation to convey information.

- We construct the meaning of natural language in terms of variables, predicates, arguments.
- Formal Semantics aim to use explicit and grounded meaning representation to convey information.
- The formal semantics can help us understand human language from a mathematical logic perspective.



Utterance: What is the largest prime less than 10?

- Data-driven learning process
- Suitable for deep learning framework



- Data-driven learning process
- Suitable for deep learning framework





We don't need to worry about the structure of computation in the middle



Tight coupling between machine learning and representation means there's always risk that some new semantic phenomenon arises and suddenly our model is useless

#### **Evaluation Chart**

Туре	Meaning	Similarity Meaning	Grounded Meaning	Lexical Meaning
String	×	×	×	$\checkmark$
Vector				
Formal				

#### **Evaluation Chart**

Туре	Meaning	Similarity Meaning	Grounded Meaning	Lexical Meaning
String	×	×	×	$\checkmark$
Vector	?	$\checkmark$	?	$\checkmark$
Formal				

#### **Evaluation Chart**

Туре	Meaning	Similarity Meaning	Grounded Meaning	Lexical Meaning
String	×	×	×	$\checkmark$
Vector	?	$\checkmark$	?	$\checkmark$
Formal	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$

#### **Formal Semantics**

- For a long time, such an explicit meaning representation has dominated the research of NLP community
  - Reasoning
  - Perception
  - Action

#### **Formal Semantics**

- For a long time, such an explicit meaning representation has dominated the research of NLP community
  - Reasoning
  - Perception
  - Action



#### **Formal Semantics**

- For a long time, such an explicit meaning representation has dominated the research of NLP community
  - Reasoning
  - Perception

\_



## List of Semantic Forms

- Introduction of formal semantics
  - Lambda Calculus (Turing 1937)
  - Combinatory categorial grammar (CCG) (Steedman et al. 1996, 2000)
  - Weighted linear CCG (Clark & Currant et. al 2007)
  - Dependency-based compositional semantics (DCS) (Liang et al. 2011)
  - Lambda DCS (Liang et al. 2013)
  - Abstract Meaning Representation (Banarescu et al. 2013)
  - Domain Specific Languages (DSL)
    - Functional program semantics (Liang et al. 2017, Dawn et al. 2018)
    - SQL, etc (Zhong et al. 2017, Yu et al. 2018)

## Semantic Forms

- Combinatory Categorial Grammar
  - Definition
  - Rules
  - Learning
- Dependency-based compositional semantics
  - Definition
  - Rules
  - Learning
- Neural Symbolic Machine (Functional-Program)
  - Definition
  - Learning

Yoav et. al, Learning Compact Lexicons for CCG Semantic Parsing, 2014 Yoav et. al, Weakly Supervised Learning of Semantic Parsers for Mapping Instructions to Actions, 2013 Luke et. al, Learning to Map Sentences to Logical Form: Structured Classification with Probabilistic Categorial Grammars, 2005 Steedman et. al, Surface Structure and Interpretation. 1996 Steedman et. al, The Syntactic Process. 2000

28

- Categorial Formalism
  - Compositional
  - Puts information on the words

- Categorial Formalism
  - Compositional
  - Puts information on the words
- Transparent interface between syntax and semantics

- Categorial Formalism
  - Compositional
  - Puts information on the words
- Transparent interface between syntax and semantics
- Including Predicate-argument structure, quantification and information structure.

- Categorial Formalism
  - Compositional
  - Puts information on the words
- Transparent interface between syntax and semantics
- Including Predicate-argument structure, quantification and information structure.
- Same expressive power as lambda calculus.

# **CCG** Categories

- Basic Building Block
- Capture Syntactic and Semantic Information Jointly

$$\fbox{ADJ} \quad : \quad \fbox{\lambda x. func(x)}$$

# **CCG** Categories

- Basic Building Block
- Capture Syntactic and Semantic Information Jointly



# **CCG** Categories

- Primitive symbols: N, S, NP, ADJ and PP
- Syntactic combination operation (/,\)
- Slashes specify argument order



# **CCG** Lexical Entries

- Pair words and phrases with meaning
- Meaning captured by a CCG category


## **CCG** Lexical Entries

- Pair words and phrases with meaning
- Meaning captured by a CCG category



# **CCG** Operations

- Equivalent to function application
- Two direction of application

$$B:g \quad A \backslash B: f \Rightarrow A: f(g) \quad (<)$$
$$A/B:f \quad B:g \Rightarrow A: f(g) \quad (>)$$

# **CCG** Operations

- Equivalent to function application
- Two direction of application



## **CCG** Parsing

- Use Lexicon to match words and phrases with their CCG categories
- Combine categories using operators



## **CCG** Parsing

- Use Lexicon to match words and phrases with their CCG categories
- Combine categories using operators



## **CCG** Parsing

- Use Lexicon to match words and phrases with their CCG categories
- Combine categories using operators



## Lexicon Problem

- Key component of CCG
- Same words often paired with many different categories
- Difficult to learn from limited data

## Lexicon Problem

- Key component of CCG
- Same words often paired with many different categories
- Difficult to learn from limited data







(garden, {*garden*}) (house, {house})

Templates	
$\lambda(\omega,v_i).$	
$\Big  \; [\omega  o ADJ : \lambda x.  of(x, ly.  v_1(y))] \Big $	
$\lambda(\omega, v_i).$	
$[\omega  o N: \lambda x.  v_1(x)]$	

Lexicon

Original Lexicon	$\begin{array}{l} \text{flight} \ \vdash S NP:\lambda x.flight(x) \\ \text{flight} \ \vdash S NP/(S NP):\lambda f.\lambda x.flight(x) \land f(x) \\ \text{flight} \ \vdash S NP \backslash (S NP):\lambda f.\lambda x.flight(x) \land f(x) \\ \text{ground transport} \ \vdash S NP:\lambda x.trans(x) \\ \text{ground transport} \ \vdash S NP/(S NP):\lambda f.\lambda x.trans(x) \land f(x) \\ \text{ground transport} \ \vdash S NP \backslash (S NP):\lambda f.\lambda x.trans(x) \land f(x) \\ \end{array}$	Compress
Factored	$( ext{flight})$ (ground transport, { $trans$ }) $\lambda(\omega, \{v_i\}_1^n).[\omega \vdash S NP : \lambda x.v_1(x)]$	Compress

 $\lambda(\omega, \{v_i\}_1^n).[\omega \vdash S|NP/(S|NP) : \lambda f.\lambda x.v_1(x) \land f(x)]$  $\lambda(\omega, \{v_i\}_1^n) [\omega \vdash S | NP \setminus (S | NP) : \lambda f . \lambda x . v_1(x) \land f(x) ]$ 

- Capture systematic variations in word usage
- Each variation can then be applied to compact units in lexical matching
- Abstracts the compositional nature of the word

## Learning (Data)

- (Natural Language, Lambda Form)

Show me flights to Boston  $\lambda x.flight(x) \wedge to(x, BOSTON)$ 

I need a flight from baltimore to seattle  $\lambda x.flight(x) \wedge from(x, BALTIMORE) \wedge to(x, SEATTLE)$ 

what ground transportation is available in san francisco  $\lambda x.ground\_transport(x) \land to\_city(x, SF)$ 

#### Learning (Structure Prediction)



#### Learning (Log-Linear Scorer)



#### Learning (Log-Linear Scorer)



- Given a weighted linear model:
- CCG lexicon  $\Lambda$
- Feature function  $f: X \times Y \to \mathbb{R}^m$
- Weights  $w \in \mathbb{R}^m$
- The best parse is:

$$y^* = rg\max_{y} w \cdot f(x, y)$$

• We consider all possible parses y for sentence x given the lexicon  $\Lambda$ 



Lexical Generation Procedure

 $GENLEX(x,\mathcal{V};\Lambda, heta)$ 

Given:

Sentence xValidation Function  $\mathcal{V}$ Log-Linear Model  $\theta$ Lexicon  $\Lambda_0$ 





### **Unified Learning Algorithm**

Initialize  $\theta$  using  $\Lambda_0\;$  ,  $\Lambda \leftarrow \Lambda_0$ 

For t = 1 ... T, i = 1 ... n:

Step 1: (Lexical generation)
Step 2: (Update parameters)

**Output:** Parameters  $\theta$  and lexicon  $\Lambda$ 

Online

• Input:

$$\{(x_i, \mathcal{V}_i) : i = 1 \dots n\}$$

- 2 steps:
  - Lexical generation
  - Parameter update

## Unified Learning Algorithm

Initialize  $\theta$  using  $\Lambda_0\;$  ,  $\Lambda \leftarrow \Lambda_0$ 

For t = 1 ... T, i = 1 ... n:

Step 1: (Lexical generation)
Step 2: (Update parameters)



Initialize  $\theta$  using  $\Lambda_0\;$  ,  $\Lambda \leftarrow \Lambda_0$ 

For t = 1 ... T, i = 1 ... n:

#### Step 1: (Lexical generation)

- a. Set  $\lambda_G \leftarrow GENLEX(x_i, \mathcal{V}_i; \Lambda, \theta)$ ,  $\lambda \leftarrow \Lambda \cup \lambda_G$
- b. Let Y be the k highest scoring parses from  $GEN(x_i; \lambda)$
- c. Select lexical entries from the highest scoring valid parses:

 $\lambda_i \leftarrow \bigcup_{y \in MAXV_i(Y;\theta)} LEX(y)$ 

d. Update lexicon:  $\Lambda \leftarrow \Lambda \cup \lambda_i$ 

Step 2: (Update parameters)

**Output:** Parameters  $\theta$  and lexicon  $\Lambda$ 

Generate a large set of potential lexical entries

eta weights  $x_i$  sentence  $\mathcal{V}_i$  validation function  $GENLEX(x_i, \mathcal{V}_i; \Lambda, \theta)$ lexical generation function

- Initialize  $\theta$  using  $\Lambda_0\;$  ,  $\Lambda \leftarrow \Lambda_0$
- For t = 1 ... T, i = 1 ... n:

#### Step 1: (Lexical generation)

- a. Set  $\lambda_G \leftarrow GENLEX(x_i, \mathcal{V}_i; \Lambda, \theta)$ ,  $\lambda \leftarrow \Lambda \cup \lambda_G$
- b. Let Y be the k highest scoring parses from  $GEN(x_i; \lambda)$
- c. Select lexical entries from the highest scoring valid parses:
  - $\lambda_i \leftarrow \bigcup_{y \in MAXV_i(Y;\theta)} LEX(y)$
- d. Update lexicon:  $\Lambda \leftarrow \Lambda \cup \lambda_i$

Step 2: (Update parameters)



- Initialize  $\theta$  using  $\Lambda_0\;$  ,  $\Lambda \leftarrow \Lambda_0$
- For t = 1 ... T, i = 1 ... n:
  - Step 1: (Lexical generation)
    - a. Set  $\lambda_G \leftarrow GENLEX(x_i, \mathcal{V}_i; \Lambda, \theta)$ ,  $\lambda \leftarrow \Lambda \cup \lambda_G$
    - b. Let Y be the k highest scoring parses from  $GEN(x_i; \lambda)$
    - c. Select lexical entries from the highest scoring valid parses:
      - $\lambda_i \leftarrow \bigcup_{y \in MAXV_i(Y;\theta)} LEX(y)$
    - d. Update lexicon:  $\Lambda \leftarrow \Lambda \cup \lambda_i$

Step 2: (Update parameters)



- Initialize  $\theta$  using  $\Lambda_0\;$  ,  $\Lambda \leftarrow \Lambda_0$
- For t = 1 ... T, i = 1 ... n:
  - Step 1: (Lexical generation)
    - a. Set  $\lambda_G \leftarrow GENLEX(x_i, \mathcal{V}_i; \Lambda, \theta)$ ,  $\lambda \leftarrow \Lambda \cup \lambda_G$
    - b. Let Y be the k highest scoring parses from  $GEN(x_i; \lambda)$
    - c. Select lexical entries from the highest scoring valid parses:
      - $\lambda_i \leftarrow \bigcup_{y \in MAXV_i(Y;\theta)} LEX(y)$
    - d. Update lexicon:  $\Lambda \leftarrow \Lambda \cup \lambda_i$

**Step 2:** (Update parameters)



Initialize  $\theta$  using  $\Lambda_0$  ,  $\Lambda \leftarrow \Lambda_0$ 

For t = 1 ... T, i = 1 ... n:

Step 1: (Lexical generation)

Step 2: (Update parameters)

- a. Set  $G_i \leftarrow MAXV_i(GEN(x_i; \Lambda); \theta)$ and  $B_i \leftarrow \{e | e \in GEN(x_i; \Lambda) \land \neg \mathcal{V}_i(y)\}$
- b. Construct sets of margin violating good and bad parses:

$$R_{i} \leftarrow \{g | g \in G_{i} \land \exists b \in B_{i} \\ s.t. \langle \theta, \Phi_{i}(g) - \Phi_{i}(b) \rangle < \gamma \Delta_{i}(g, b) \} \\ E_{i} \leftarrow \{b | b \in B_{i} \land \exists g \in G_{i} \\ s.t. \langle \theta, \Phi_{i}(g) - \Phi_{i}(b) \rangle < \gamma \Delta_{i}(g, b) \}$$

c. Apply the additive update:

$$egin{aligned} & heta \leftarrow heta + rac{1}{|R_i|} \sum_{r \in R_i} \Phi_i(r) \ & -rac{1}{|E_i|} \sum_{e \in E_i} \Phi_i(e) \end{aligned}$$



Initialize  $\theta$  using  $\Lambda_0$  ,  $\Lambda \leftarrow \Lambda_0$ 

For t = 1 ... T, i = 1 ... n:

- **Step 1:** (Lexical generation)
- Step 2: (Update parameters)
  - a. Set  $G_i \leftarrow MAXV_i(GEN(x_i; \Lambda); \theta)$ and  $B_i \leftarrow \{e | e \in GEN(x_i; \Lambda) \land \neg \mathcal{V}_i(y)\}$
  - b. Construct sets of margin violating good and bad parses:

$$R_{i} \leftarrow \{g | g \in G_{i} \land \exists b \in B_{i} \\ s.t. \langle \theta, \Phi_{i}(g) - \Phi_{i}(b) \rangle < \gamma \Delta_{i}(g, b) \} \\ E_{i} \leftarrow \{b | b \in B_{i} \land \exists g \in G_{i} \\ s.t. \langle \theta, \Phi_{i}(g) - \Phi_{i}(b) \rangle < \gamma \Delta_{i}(g, b) \}$$



c. Apply the additive update:

$$\theta \leftarrow \theta + \frac{1}{|R_i|} \sum_{r \in R_i} \Phi_i(r) \\ - \frac{1}{|E_i|} \sum_{e \in E_i} \Phi_i(e)$$

Initialize heta using  $\Lambda_0$  ,  $\Lambda \leftarrow \Lambda_0$ 

For t = 1 ... T, i = 1 ... n:

Step 1: (Lexical generation)

Step 2: (Update parameters)

- a. Set  $G_i \leftarrow MAXV_i(GEN(x_i; \Lambda); \theta)$ and  $B_i \leftarrow \{e | e \in GEN(x_i; \Lambda) \land \neg \mathcal{V}_i(y)\}$
- b. Construct sets of margin violating good and bad parses:

$$\begin{aligned} R_i \leftarrow \{g | g \in G_i \land \exists b \in B_i \\ s.t. & \langle \theta, \Phi_i(g) - \Phi_i(b) \rangle < \gamma \Delta_i(g, b) \} \\ E_i \leftarrow \{b | b \in B_i \land \exists g \in G_i \\ s.t. & \langle \theta, \Phi_i(g) - \Phi_i(b) \rangle < \gamma \Delta_i(g, b) \} \end{aligned}$$

c. Apply the additive update:

$$eta \leftarrow heta + rac{1}{|R_i|} \sum_{r \in R_i} \Phi_i(r) \ - rac{1}{|E_i|} \sum_{e \in E_i} \Phi_i(e)$$



Initialize heta using  $\Lambda_0$  ,  $\Lambda \leftarrow \Lambda_0$ 

For t = 1 ... T, i = 1 ... n:

Step 1: (Lexical generation)

Step 2: (Update parameters)

- a. Set  $G_i \leftarrow MAXV_i(GEN(x_i; \Lambda); \theta)$ and  $B_i \leftarrow \{e | e \in GEN(x_i; \Lambda) \land \neg \mathcal{V}_i(y)\}$
- b. Construct sets of margin violating good and bad parses:

$$R_{i} \leftarrow \{g | g \in G_{i} \land \exists b \in B_{i} \\ s.t. \langle \theta, \Phi_{i}(g) - \Phi_{i}(b) \rangle < \gamma \Delta_{i}(g, b) \} \\ E_{i} \leftarrow \{b | b \in B_{i} \land \exists g \in G_{i} \\ s.t. \langle \theta, \Phi_{i}(g) - \Phi_{i}(b) \rangle < \gamma \Delta_{i}(g, b) \}$$

c. Apply the additive update:

$$\theta \leftarrow \theta + \frac{1}{|R_i|} \sum_{r \in R_i} \Phi_i(r) \\ - \frac{1}{|E_i|} \sum_{e \in E_i} \Phi_i(e)$$



### **Semantic Forms**

- Combinatory Categorial Grammar
  - Definition
  - Rules
  - Learning
- Dependency-based compositional semantics
  - Definition
  - Rules
  - Learning
- Neural Symbolic Machine (Functional-Program)
  - Definition
  - Learning

Liang et. al, Learning dependency-based compositional semantics, 2011 Liang et. al, Lambda dependency-based compositional semantics, 2013 Jonathan et. al, Semantic parsing on Freebase from question-answer pairs, 2013 Jonathan et. al, Semantic parsing via paraphrasing, 2014

- Weakness of CCG Parse
  - Too much hand-coded rules for lexicon construction.
  - Learning algorithm is complicated.
  - The grammar rules are too strict.
  - Learning requires annotating logic form.



$$\begin{split} \lambda c \, \exists m \, \exists \ell \, \exists s \, . \\ \texttt{city}(c) \wedge \texttt{major}(m) \wedge \\ \texttt{loc}(\ell) \wedge \texttt{CA}(s) \wedge \\ c_1 = m_1 \wedge c_1 = \ell_1 \wedge \ell_2 = s_1 \end{split}$$

major cities in CA



major cities in CA



$$\begin{split} \lambda c \, \exists m \, \exists \ell \, \exists s \, . \\ \texttt{city}(c) \wedge \texttt{major}(m) \wedge \\ \texttt{loc}(\ell) \wedge \texttt{CA}(s) \wedge \\ c_1 = m_1 \wedge c_1 = \ell_1 \wedge \ell_2 = s_1 \end{split}$$

major cities in CA


# Strong vs. Weak Supervision

**Detailed Supervision (current)** 

- doesn't scale up
- representation-dependent

What is the largest city in California?

expert

 $\operatorname{argmax}(\{c: \operatorname{city}(c) \land \operatorname{loc}(c, \operatorname{CA})\}, \operatorname{population})$ 

# Strong vs. Weak Supervision

#### **Detailed Supervision (current)**

- doesn't scale up
- representation-dependent

What is the largest city in California? expert  $argmax(\{c: city(c) \land loc(c, CA)\}, population)$ 

#### Natural Supervision (new)

- scales up
- representation-independent



### Lexicon vs. Triggers



requires strict grammar rule during parsing and harsh type constraint

### Lexicon vs. Triggers





CCG

{(city, city),(city, state),(city, river), . . . }

requires strict grammar rule during parsing and harsh type constraint requires only trigger words for predicates without grammar rule

# Syntactic & Semantic Alignment



Utterance: major cities in CA

# Syntactic & Semantic Alignment



Utterance: major cities in CA

# Parsing: Lexical Triggers

city city state state river river argmax population population CA What is the most populous city in CA ?

# Parsing: Lexical Triggers



#### Lexical Triggers:

- 1. String match  $CA \Rightarrow CA$
- 2. Function words (20 words)  $most \Rightarrow \operatorname{argmax}$
- 3. Nouns/adjectives  $city \Rightarrow city \text{ state river population}$

# Parsing: Predicates to DCS Trees (DP)

 $C_{i,j} = \text{set of DCS trees for span } [i,j]$ 



# Parsing: Predicates to DCS Trees (DP)

 $C_{i,j} = \mathsf{set} \text{ of DCS trees for span } [i,j]$ 



## Parsing: Predicates to DCS Trees (DP)

 $C_{i,j} = \text{set of DCS trees for span } [i,j]$ 





#### Database

city
San Francisco
Chicago
Boston

#### loc

Mount Shasta	California
San Francisco	California
Boston	Massachusetts

CA California

Assuming we have our world knowledge containing 1) unary predicates 2) binary predicates in a database.



#### Database

city	
San Francisco	0
Chicago	
Boston	

#### loc

Mount Shasta	California
San Francisco	California
Boston	Massachusetts

CA California

How can we assign values to different nodes?



#### Database

city
San Francisco
Chicago
Boston

#### loc

Mount Shasta	California
San Francisco	California
Boston	Massachusetts

CA California

Constraint 1: the value in the each node should come from its corresponding predicate table.

DCS tree	Constraints
city	$c \in \texttt{city}$
1	$c_1 = \ell_1$
loc	$\ell \in \texttt{loc}$



Mount Shasta	California
San Francisco	California
Boston	Massachusetts

CA California

Constraint 2: the assignment of adjacent nodes need to meet the requirement in the edge.

#### **Constraint Satisfaction Problem**



D	a	ta	b	a	s	e

city
San Francisco
Chicago
Boston

	10	C
Mount S	Shasta	California
San Fra	ncisco	California
Boston		Massachusetts

CA California

A DCS tree assignment becomes a constraint satisfaction problem (CSP)

#### **Constraint Satisfaction Problem**



city
San Francisco
Chicago
Boston
•••

Database

loc		
California		
California		
Massachusetts		

CA California

A DCS tree assignment becomes a constraint satisfaction problem (CSP) Dynamic Programming => time complexity = O(#node)

## Learning



## Learning

**Objective Function:** 

 $\max_{\boldsymbol{\theta}} \sum_{\boldsymbol{z}} p(\boldsymbol{y} \mid \boldsymbol{z}, \boldsymbol{w}) \, p(\boldsymbol{z} \mid \boldsymbol{x}, \boldsymbol{\theta})$ 

Interpretation Semantic parsing

## Learning



# **Semantic Forms**

- Combinatory Categorial Grammar
  - Definition
  - Rules
  - Learning
- Dependency-based compositional semantics
  - Definition
  - Rules
  - Learning
- Neural Symbolic Machine (Functional-Program)
  - Definition
  - Learning

Liang, Chen et al. Neural Symbolic Machines: Learning Semantic Parsers on Freebase with Weak Supervision. 2017 Liang, Chen et al. Memory Augmented Policy Optimization for Program Synthesis and Semantic Parsing. 2018 Agarwal et al. Learning to Generalize from Sparse and Underspecified Rewards. 2019

# Neural Symbolic Machine

- Weakness of traditional Semantic Parsing
  - The previous methods heavily rely on the predicate mapping, the mapping procedure is based on rules.
  - The coverage is problematic when dealing with more complex utterance.
  - The previous frameworks are dominated by the composition rules enforced as prior, only few parameters to learn.
  - The symbolic executor and the ranking model are separated without interaction.

# Neural Symbolic Machine



# Neural Symbolic Machine



# **Domain Specific Language**

GO (Hop M1 ParentOf) (Argmin M2 BornIn) RETURN

```
(Argmin M2 BornIn)
```

```
 \begin{split} &= \{e_1 \mid e_1 \in \mathsf{M2}, \\ &\quad \exists e_2: \quad (e_1, \texttt{BornIn}, e_2) \in \mathcal{K}, \\ &\quad \forall e \in \mathsf{M2}, e_3: \quad (e, \texttt{BornIn}, e_3) \in \mathcal{K} \Longrightarrow e_2 \leq e_3 \} \end{split}
```

MEMORY		
M1	{m.Obama}	
M2	{m.Malia,m.Sasha}	
MЗ		

={m.Malia}

${\mathcal K}$ : BornIn			
m.Obama	1961		
m.Malia	1998		
m.Sasha	2001		

# **Domain Specific Language**





# **Function definitions**

- (Hop v p): entities reached from v using relation p
- (ArgHop v1 v2 p): entities in v1 from which v2 can be reached using relation p
- (ArgMin v n): entities in v which have the lowest number in field n
- (ArgMax v n): entities in v which have the highest number in field n
- ...
- (Count v): number of entities in v

# Compositionality

x: Largest city in the US => y: NYC



# Memory-aided Sequence

x: Largest city in the US => y: NYC



The graph structure is flattened as sequence

#### Large Search Space



## Compiler-aided prune











# Memory Enc-Dec Model

- Overall model is conditional distribution over all token sequences.

$$P_{\theta}(\mathbf{z}|\mathbf{x}) = P_{\theta}(z_{1}, \dots, z_{T}|x_{1}, \dots, x_{N}) = P_{\theta}(z_{1:T}|x_{1:N})$$

 $= P_{\theta}(z_1|x_{1:N}) \cdot P_{\theta}(z_2|x_{1:N}, z_1) \cdot P_{\theta}(z_3|x_{1:N}, z_1, z_2) \cdot \ldots \cdot P_{\theta}(z_T|x_{1:N}, z_{1:T-1})$
### Memory Enc-Dec Model

- Overall model is conditional distribution over all token sequences.

$$P_{\theta}(\mathbf{z}|\mathbf{x}) = P_{\theta}(z_1, \dots, z_T | x_1, \dots, x_N) = P_{\theta}(z_{1:T} | x_{1:N})$$

 $= P_{\theta}(z_1|x_{1:N}) \cdot P_{\theta}(z_2|x_{1:N}, z_1) \cdot P_{\theta}(z_3|x_{1:N}, z_1, z_2) \cdot \ldots \cdot P_{\theta}(z_T|x_{1:N}, z_{1:T-1})$ 

- NSM models the probability at each time step  $P_{\theta}(z_2|x_{1:N}, z_1)$  conditioned on previous output tokens and utterance

### Memory Enc-Dec Model

- Overall model is conditional distribution over all token sequences.

$$P_{\theta}(\mathbf{z}|\mathbf{x}) = P_{\theta}(z_{1}, \dots, z_{T}|x_{1}, \dots, x_{N}) = P_{\theta}(z_{1:T}|x_{1:N})$$

 $= P_{\theta}(z_1|x_{1:N}) \cdot P_{\theta}(z_2|x_{1:N}, z_1) \cdot P_{\theta}(z_3|x_{1:N}, z_1, z_2) \cdot \dots \cdot P_{\theta}(z_T|x_{1:N}, z_{1:T-1})$ 

- NSM models the probability at each time step  $P_{\theta}(z_2|x_{1:N}, z_1)$  conditioned on previous output tokens and utterance
- NSM uses interpreter to aid the decoder by removing impossible tokens at each step to dramatically shrink search space.

#### Memory Enc-Dec Model



### Syntactic Constraint



#### Semantic Constraint



# Augmented REINFORCE

- REINFORCE (explore)
- Iterative Maximum Likelihood (exploit)

 $\nabla_{\theta} \mathcal{J}_{\boldsymbol{x}}^{RL}(\theta) = \mathbb{E}_{\boldsymbol{z} \sim P_{\theta}}(\cdot | \boldsymbol{x}) [\nabla_{\theta} \log P_{\theta}(\boldsymbol{z} | \boldsymbol{x}) \cdot F_{1}(\boldsymbol{y}^{\boldsymbol{z}}, \boldsymbol{y}^{*})]$ 

$$\approx 0.9 \cdot \frac{1}{\lambda} \sum_{n=1}^{k} P_{\theta}(\boldsymbol{z}_{n} | \boldsymbol{x}) \cdot \nabla_{\theta} \log P_{\theta}(\boldsymbol{z}_{n} | \boldsymbol{x}) \cdot F_{1}(\boldsymbol{y}^{\boldsymbol{z}_{n}}, \boldsymbol{y}^{*}) + 0.1 \cdot P_{\theta}(\boldsymbol{z}_{best} | \boldsymbol{x}) \cdot \nabla_{\theta} \log P_{\theta}(\boldsymbol{z}_{best} | \boldsymbol{x}) \cdot F_{1}(\boldsymbol{y}^{\boldsymbol{z}_{best}}, \boldsymbol{y}^{*})$$

$$\{\mathbf{z}_n\} = k. \operatorname{argmax}_{\mathbf{z}} P_{\theta}(\mathbf{z}|\mathbf{x}), \quad \lambda = \Sigma_i P_{\theta}(\mathbf{z}_n|\mathbf{x})$$

# Wrap up

- Formal Semantics does have many advantages like strong compositionality, explainability.
- However, the annotation effort or rigid structure makes it hard to scale up to large-scale domains or more realistic datasets.
- The vectorized representation in Deep Learning is still the favorable by the community.
- How to neuralize the formal semantics and apply it to modern architecture remains an open problem.

# My Work

- 1. Natural Language Understanding on Structured Data
  - a. Relation Question Answering on Knowledge Graph (Chen et al. NAACL 18)
  - b. Table-based Fact Verification on semi-structured table (Chen et al. Arxiv)
  - c. Multi-model Question Answering on structured scene graph (Chen et al. Arxiv)

- 2. Natural Language Generation on Structured Data
  - a. Semantically conditioned dialog generation on structured semantic form (Chen et al. ACL 19)

# **Open Question**

- Philosophical problem



Meaning serves for the success of downstream tasks

# **Open Question**

- Philosophical problem



If we can already achieve high accuracy, why do we still need to care about meaning?