

Knowledge-enhanced NLP: Progress and Trend

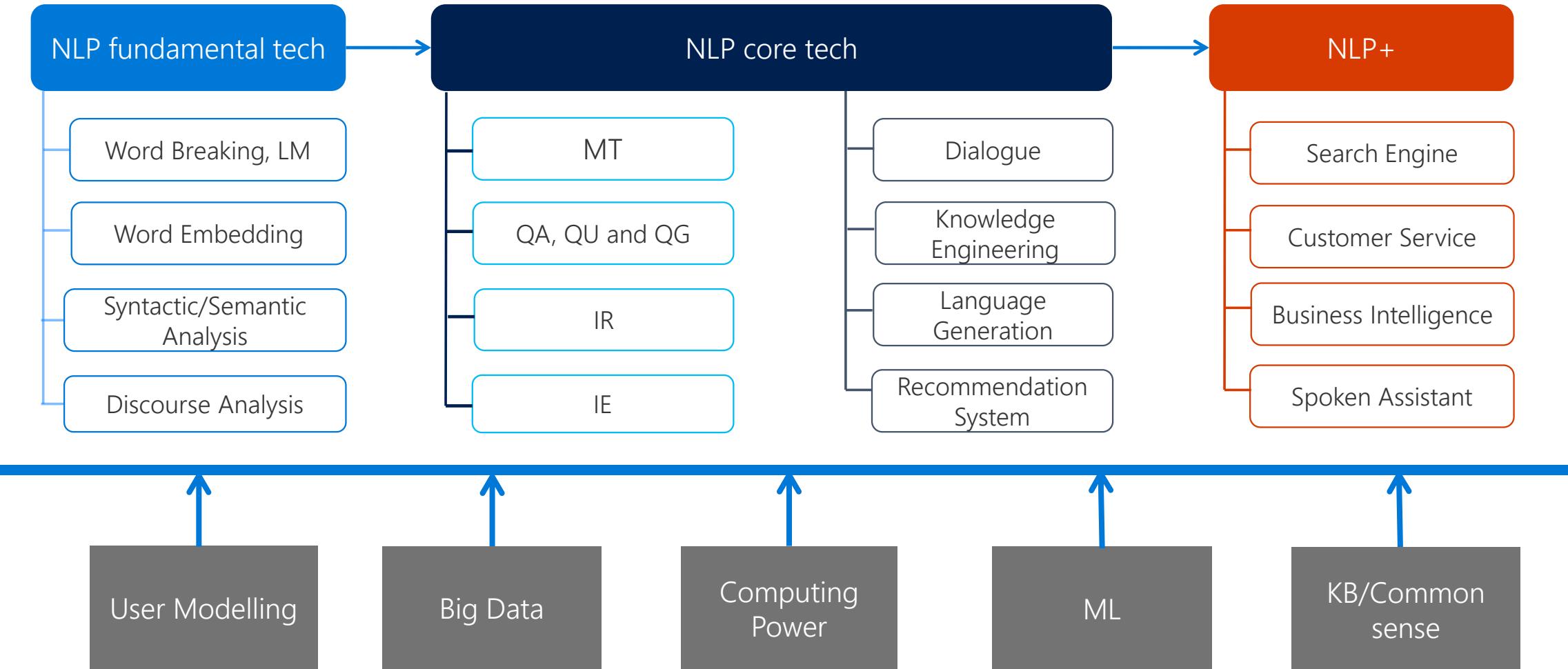
Nan Duan

Microsoft Research Asia

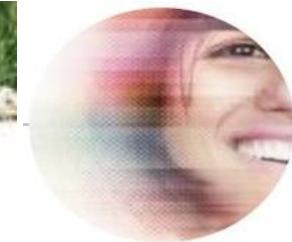
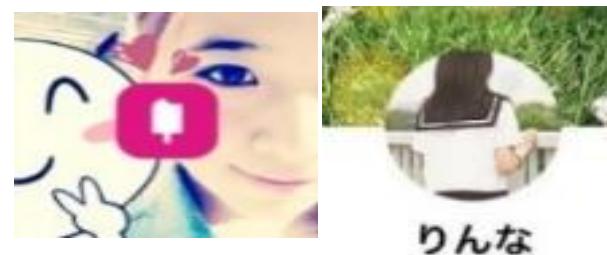
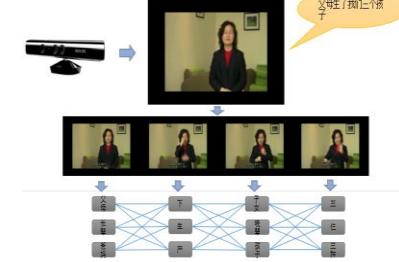
2018-11-17



NLP Technologies

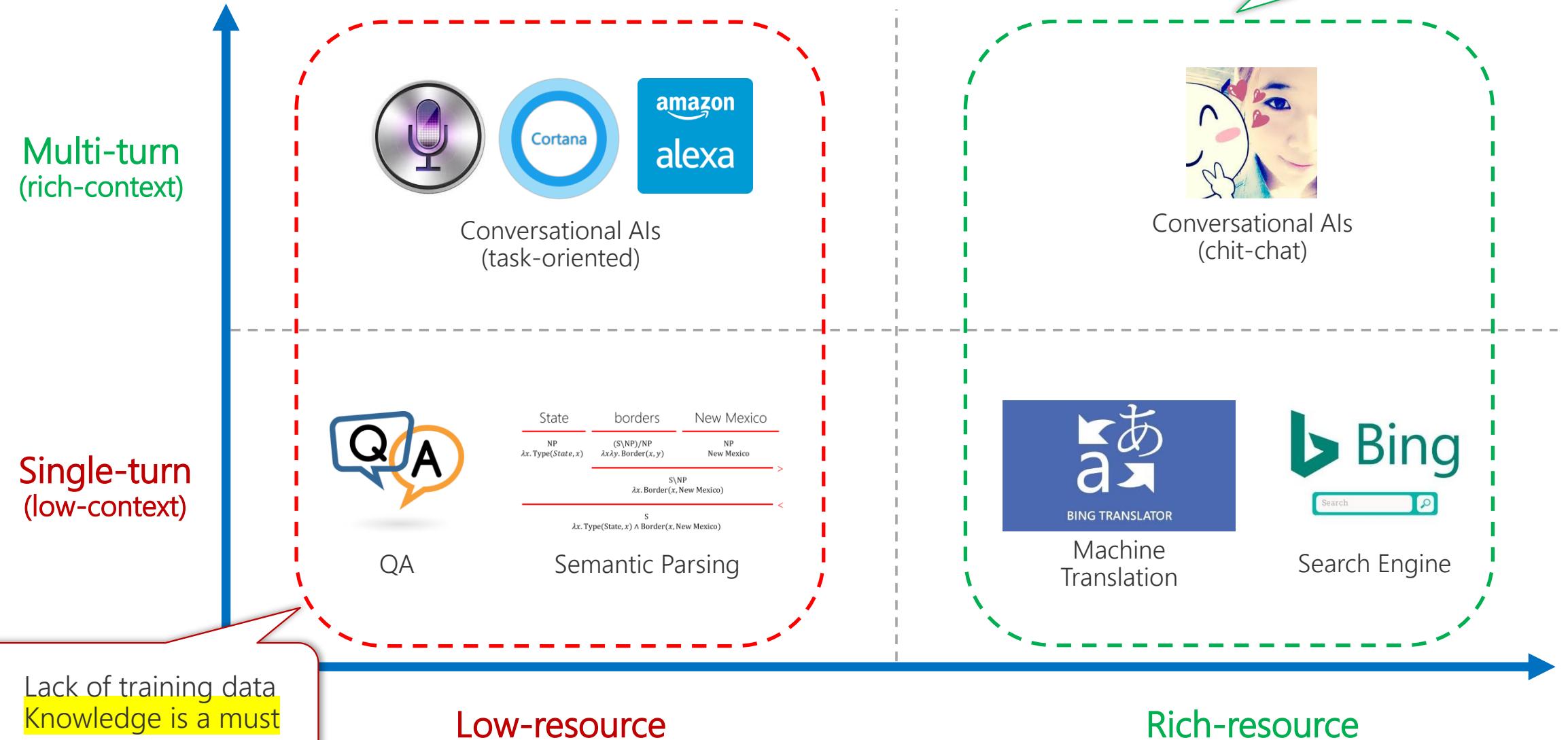


MSRA NLP Research Achievements



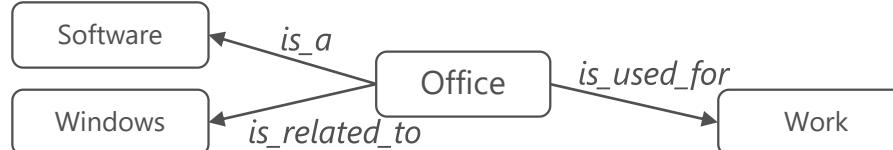
Slides from Dr. Ming Zhou

Importance of Knowledge to NLP

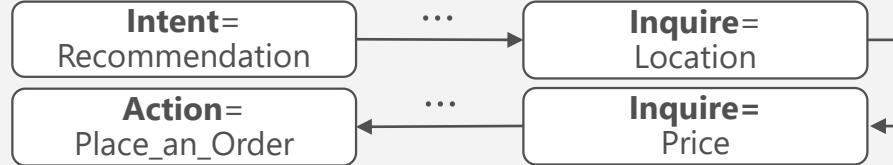


Knowledge Pyramid in NLP

Commonsense



**Knowledge Base
(Specific Domain)**



**Knowledge Base
(Open Domain)**

Microsoft Office	Developer	Microsoft
Microsoft Office	Initial release	19 November 1990
Microsoft Office	Operating system	Microsoft Windows
Microsoft Office	Website	office.com
...



**Supervised
Pre-trained Embeddings**

MT Model
QA Model
MRC Model
...



**Unsupervised
Pre-trained Embeddings**



which version of Office is right for student

Today's Agenda

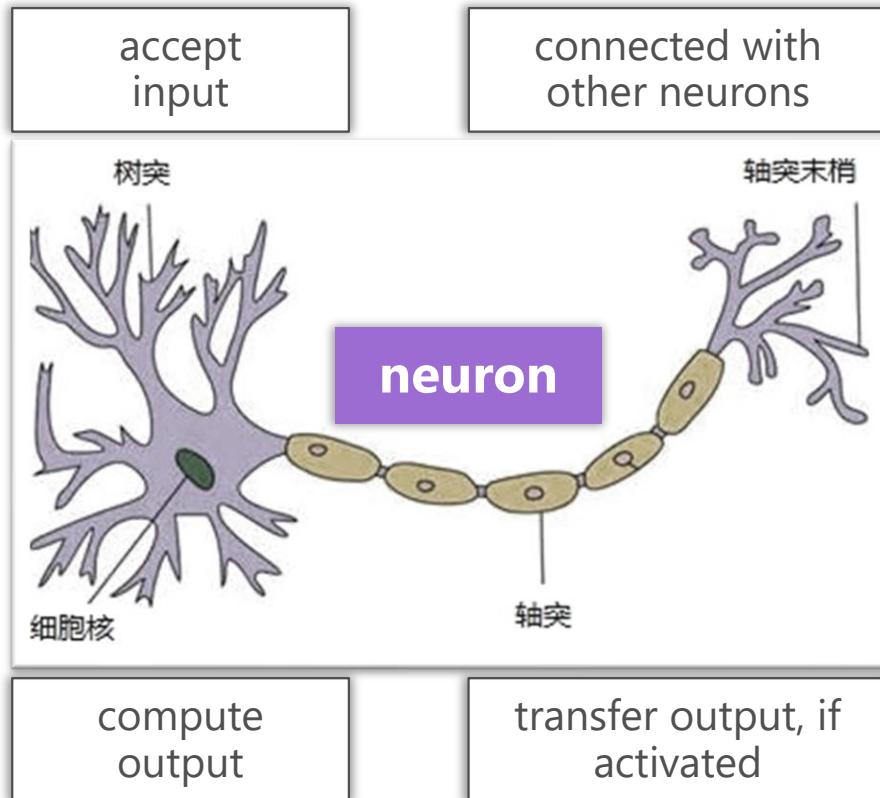
- Deep Learning Basics for NLP
- NLP with Pre-trained Embeddings
- NLP with Knowledge Bases
- NLP with Commonsense
- Summary and Trend

Today's Agenda

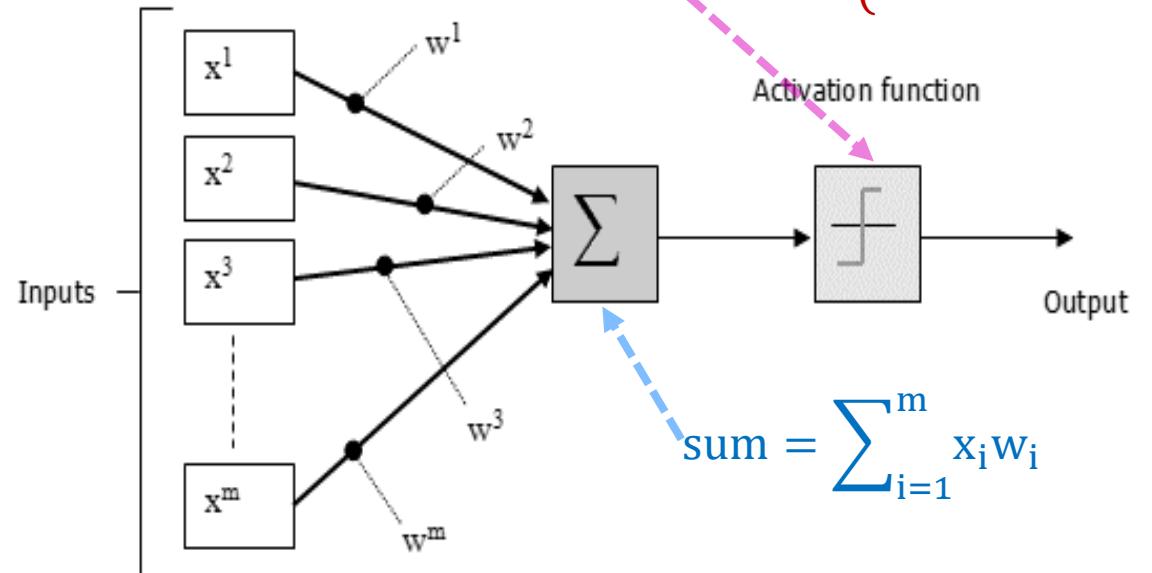
- **Deep Learning Basics for NLP**
- NLP with Pre-trained Embeddings
- NLP with Knowledge Bases
- NLP with Commonsense
- Summary and Trend

Neuron (神经元)

(McCulloch and Pitts, 1943)

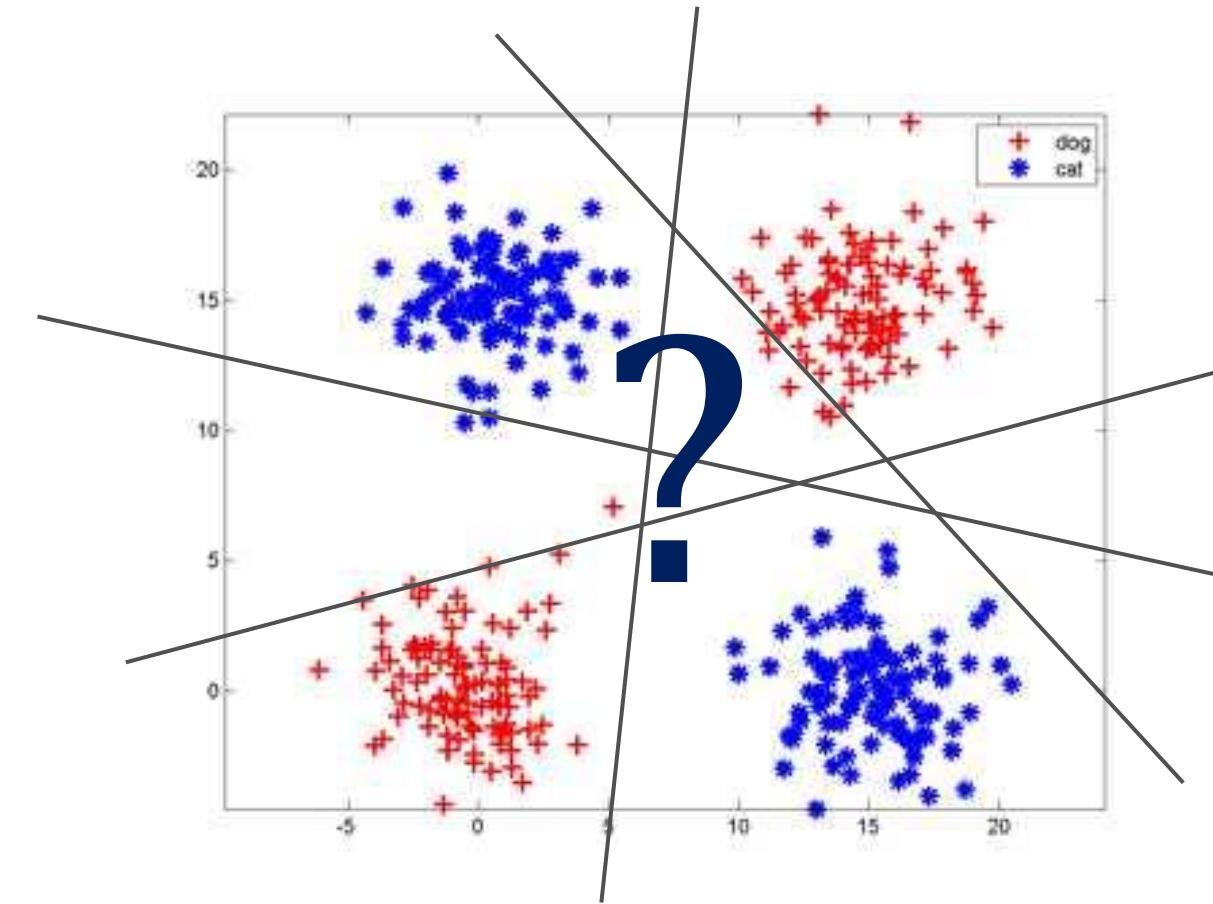
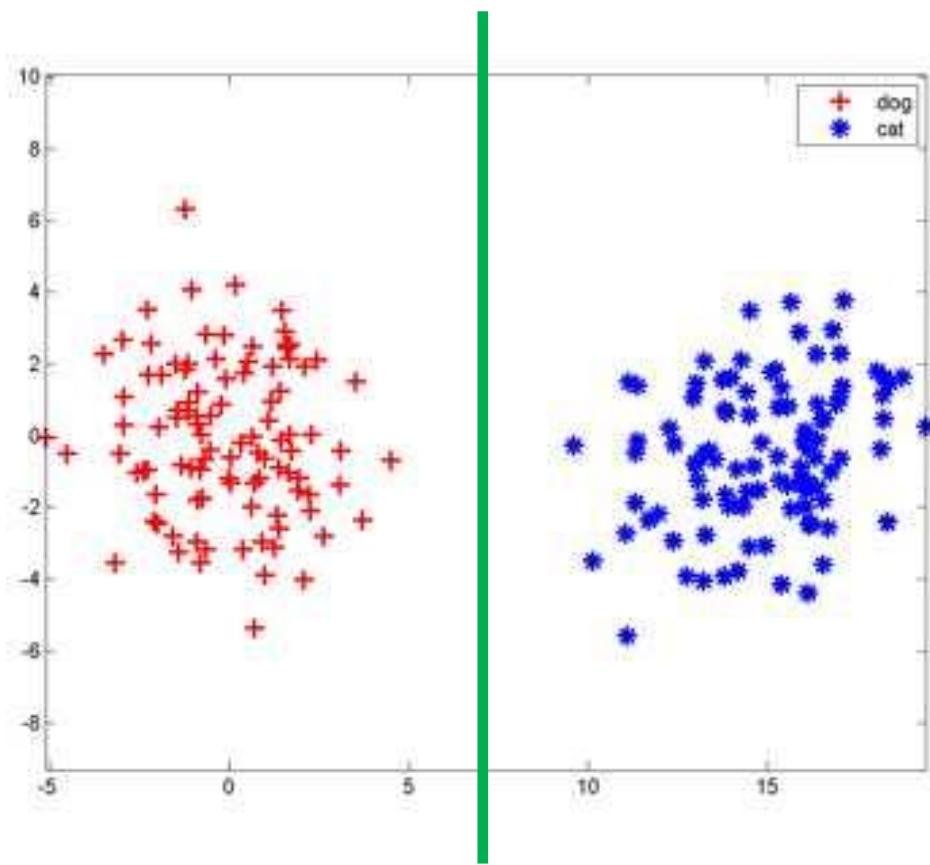


$$a(sum) = \begin{cases} +1 & sum \geq 0 \\ -1 & sum < 0 \end{cases}$$

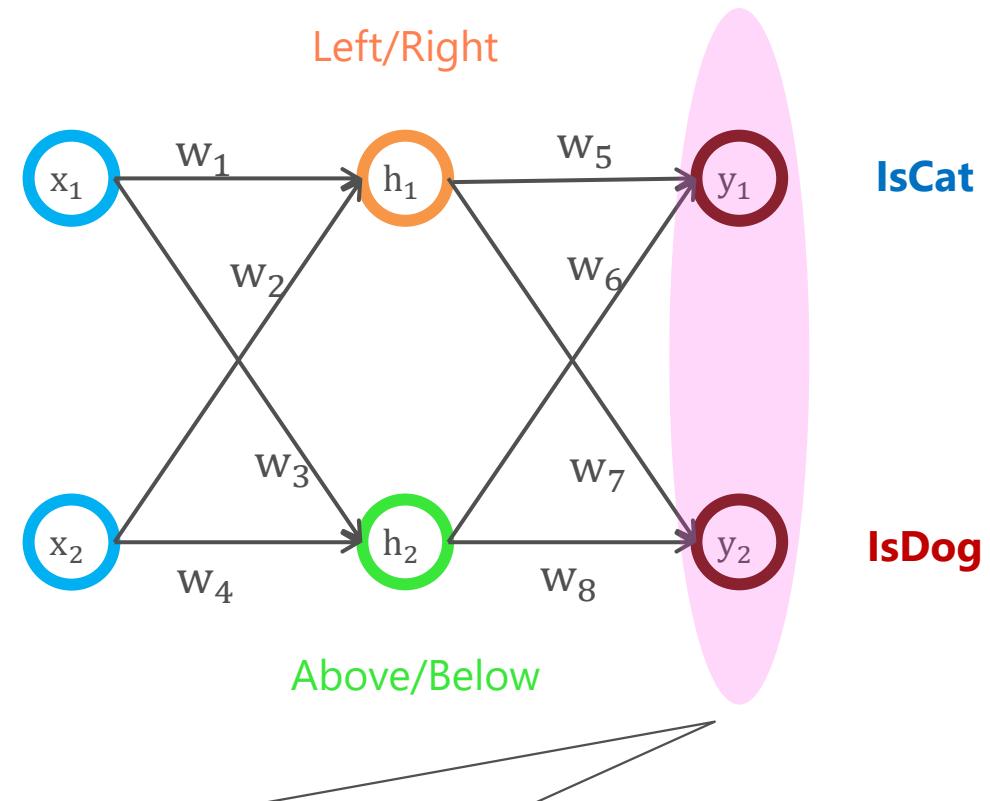
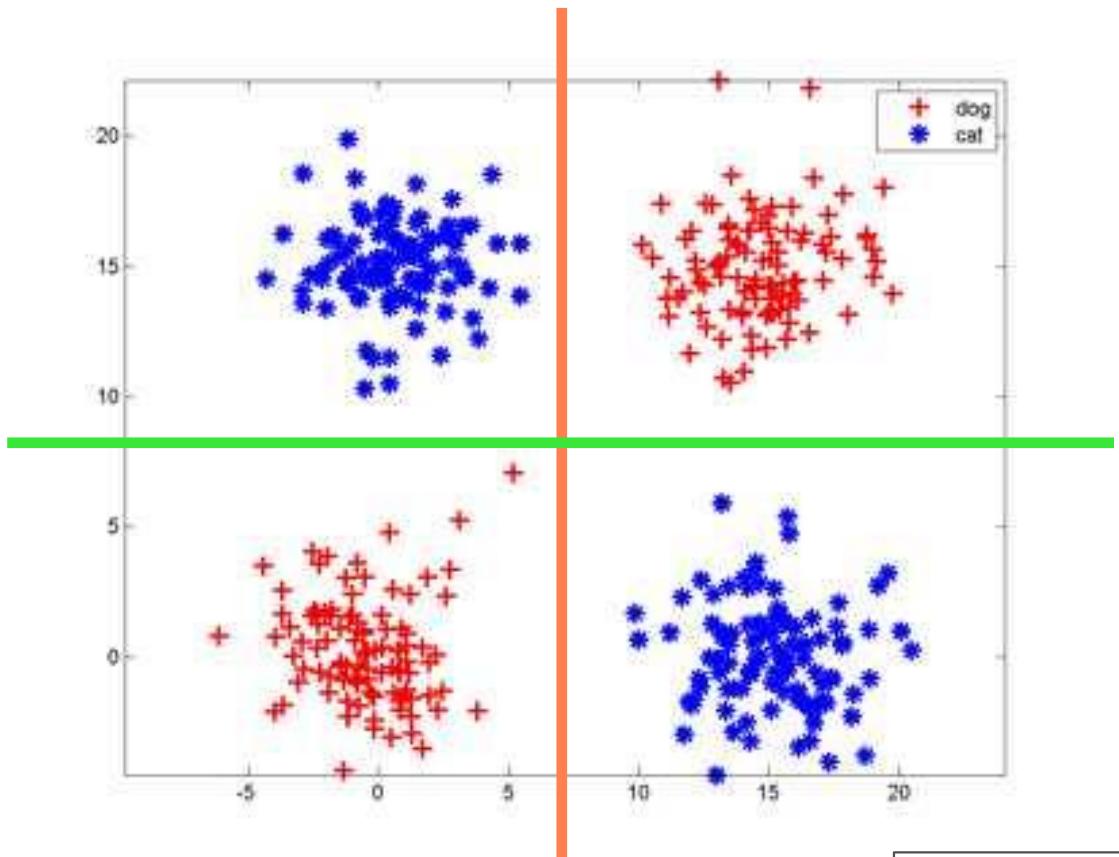


Perception (感知机)

Separate Cats and Dogs with Neuron



Neural Network (神经网络)



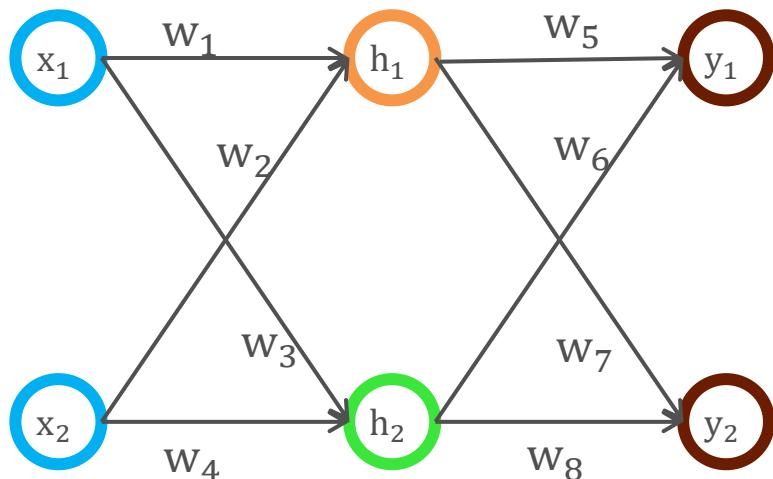
softmax function squashes a K-dimensional vector Z of arbitrary real values to a K-dimensional vector $\sigma(Z)$ of real values in the range $[0, 1]$ that add up to 1.

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$$

Feed Forward

$$\text{net}_{h1} = w_1 * x_1 + w_2 * x_2$$

$$\text{out}_{h1} = \frac{1}{1 + e^{-\text{net}_{h1}}}$$



$$\text{net}_{h2} = w_3 * x_1 + w_4 * x_2$$

$$\text{out}_{h2} = \frac{1}{1 + e^{-\text{net}_{h2}}}$$

$$\text{net}_{y1} = w_5 * \text{out}_{h1} + w_6 * \text{out}_{h2}$$

$$\text{net}_{y2} = w_7 * \text{out}_{h1} + w_8 * \text{out}_{h2}$$

$$\text{out}_{y1} = \frac{1}{1 + e^{-\text{net}_{y1}}} \neq \text{target}_{y1}$$

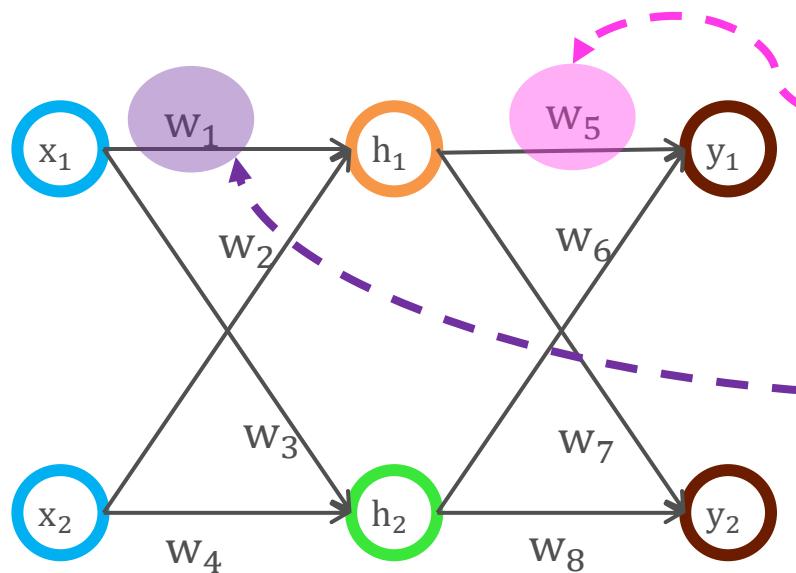
$$\text{out}_{y2} = \frac{1}{1 + e^{-\text{net}_{y2}}} \neq \text{target}_{y2}$$

We want to adjust weights $\{w_j\}$ to ensure that each out_{y_i} is as close to target_{y_i} as possible.

Backward Propagation

$$w_i^+ = w_i - \alpha \frac{\partial E}{\partial w_i}$$

$$E = \frac{1}{2} (\text{target}_{y1} - \text{out}_{y1})^2 + \frac{1}{2} (\text{target}_{y2} - \text{out}_{y2})^2$$

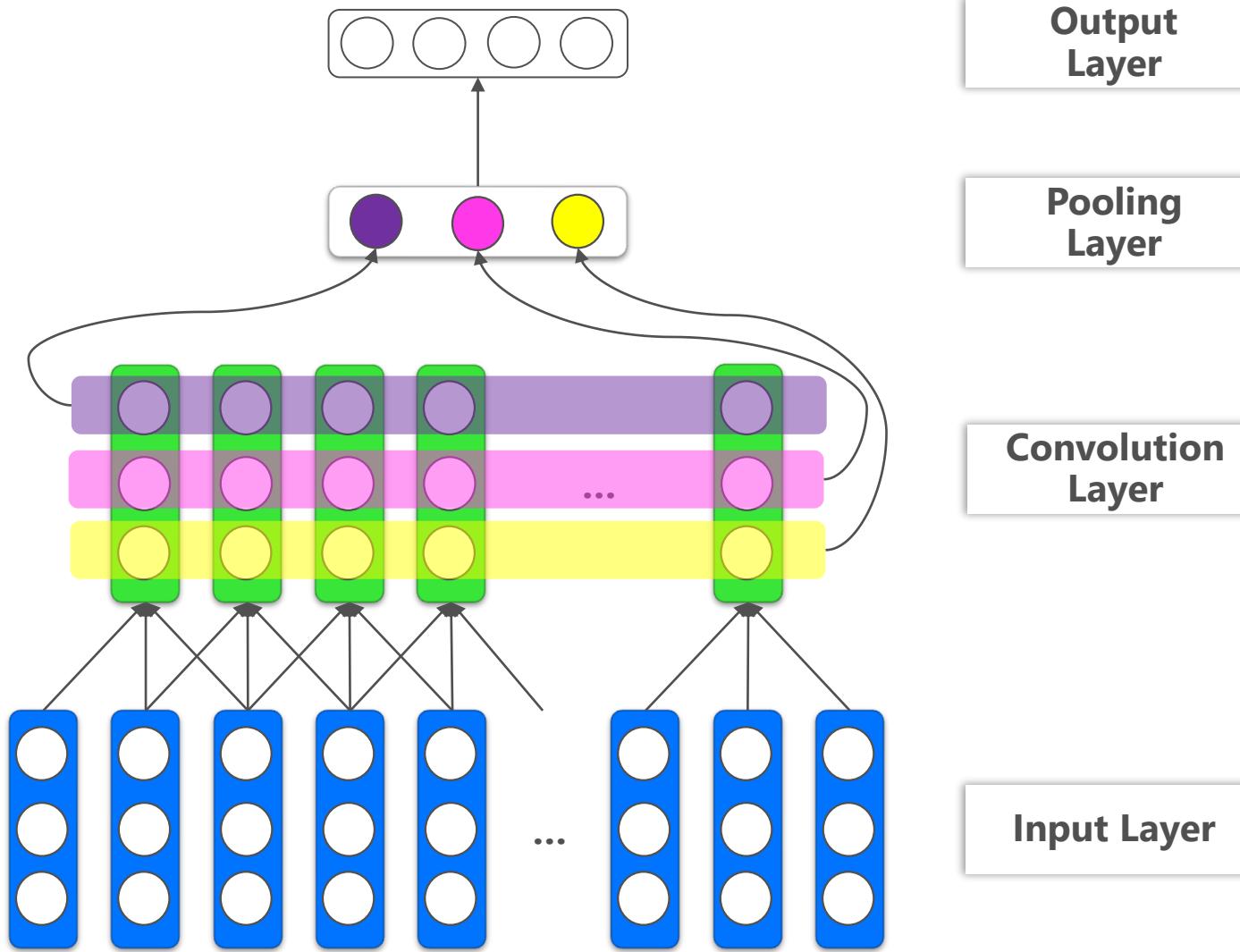


$$\frac{\partial E}{\partial w_5} = \frac{\partial E}{\partial \text{out}_{y1}} * \frac{\partial \text{out}_{y1}}{\partial \text{net}_{y1}} * \frac{\partial \text{net}_{y1}}{\partial w_5}$$

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial \text{out}_{h1}} * \frac{\partial \text{out}_{h1}}{\partial \text{net}_{h1}} * \frac{\partial \text{net}_{h1}}{\partial w_1}$$

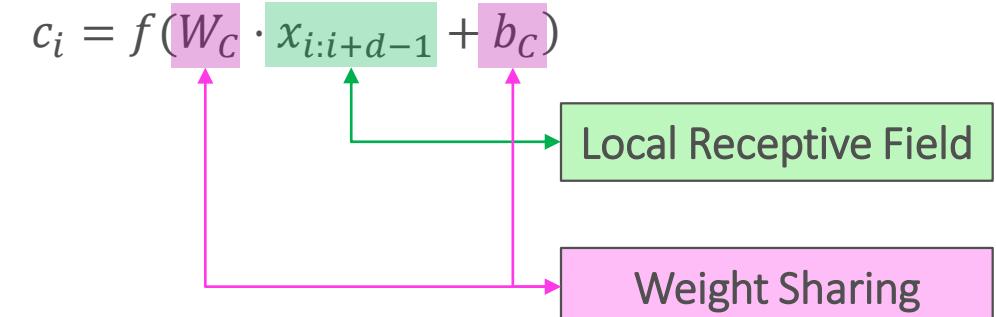
Chain Rule

Convolution Neural Network (CNN)

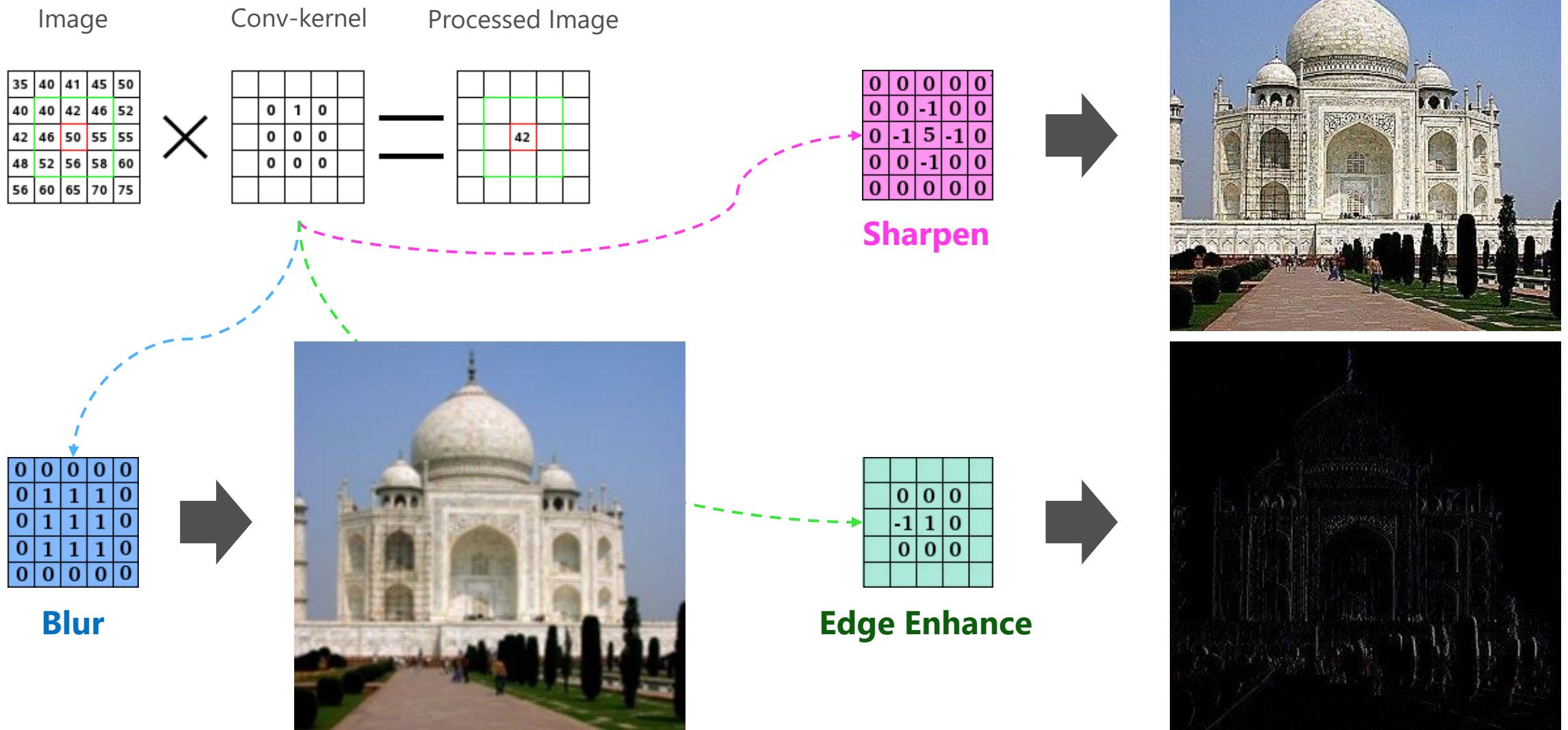


$$y_O = \tanh(W_O \cdot \hat{c} + b_O)$$

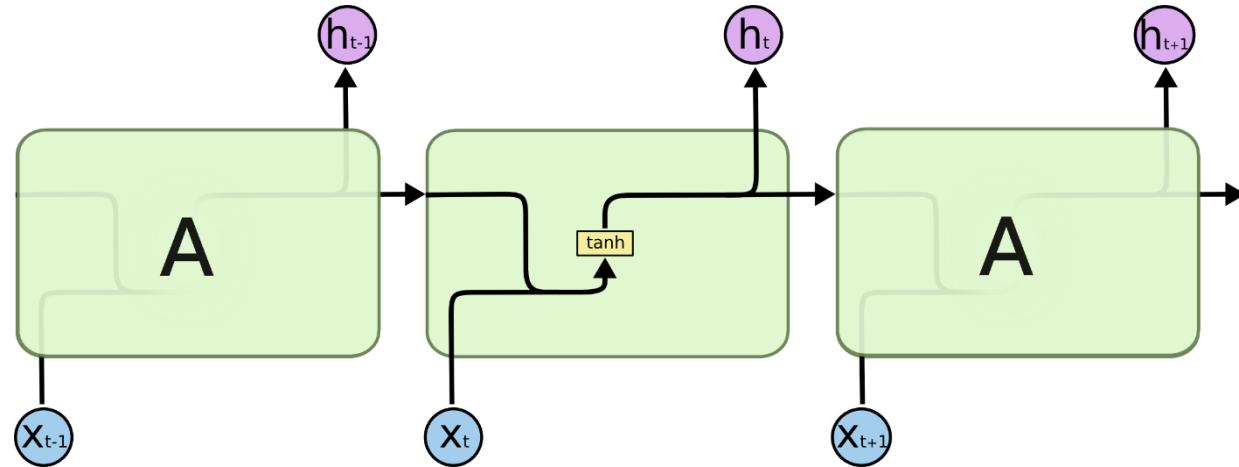
$$\hat{c}^j = \max\{c_1^j, \dots, c_N^j\}$$



Convolution is Feature Extraction

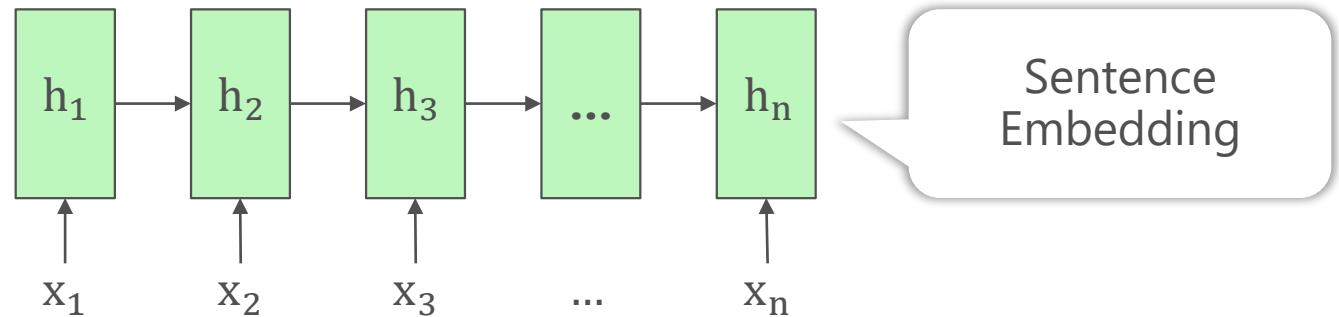


Recurrent Neural Network (RNN)

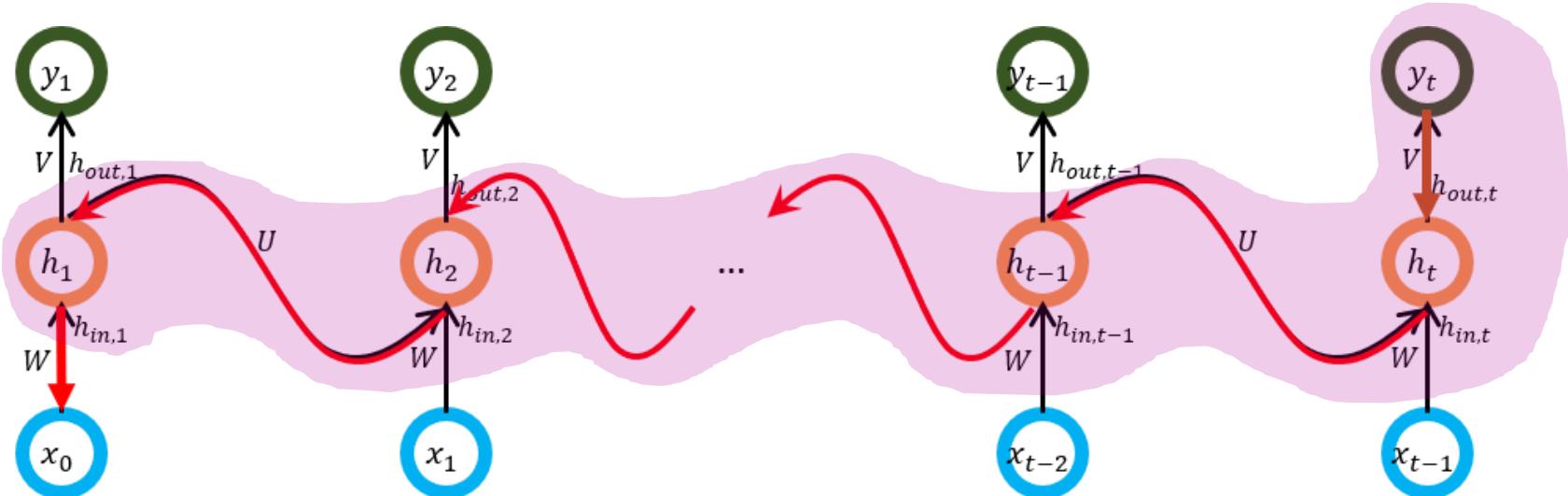


$$h_t = \tanh(W_h \cdot h_{t-1} + W_x \cdot x_t + b) = \tanh(W \cdot [h_{t-1}; x_t] + b)$$

- Goal
 - Represent an ordered sequence of words
- Input
 - A sequence of word embeddings
- Output
 - A sequence of hidden states

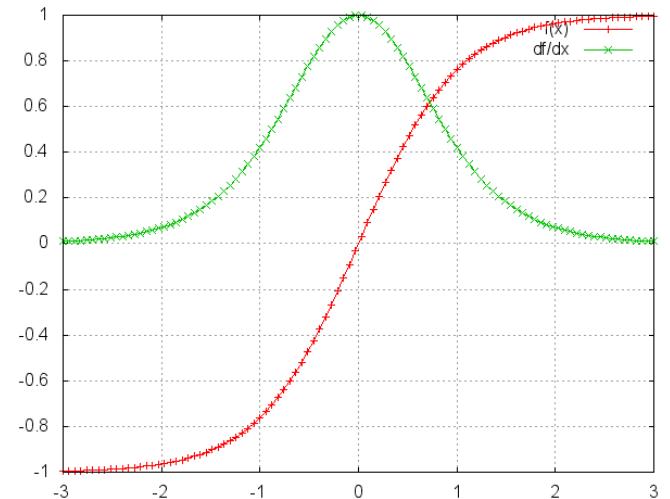


Vanishing Gradients

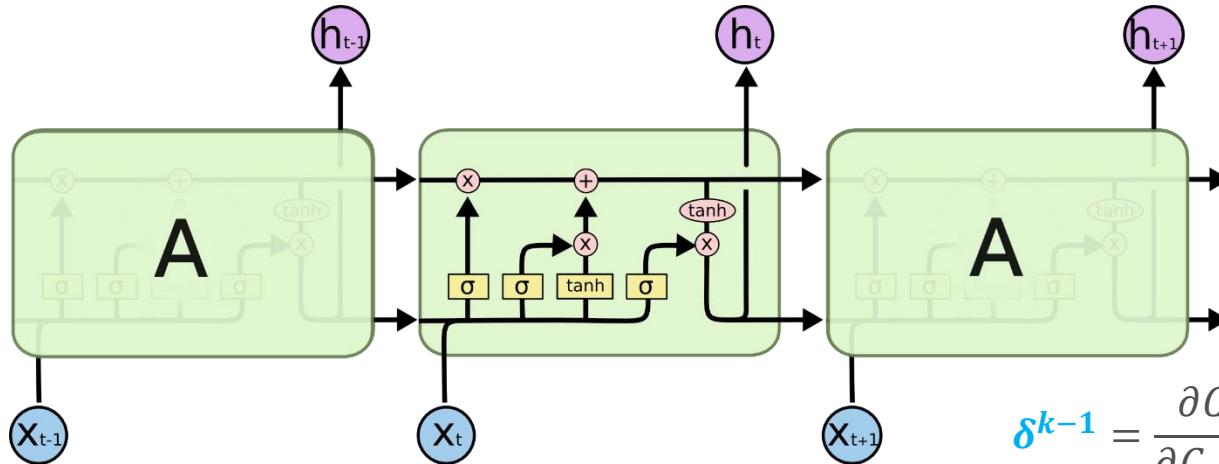


$$\delta_{in,1} = \delta_{out,t} * \frac{\partial h_{out,t}}{\partial h_{in,t}} * \frac{\partial h_{in,t}}{\partial h_{out,t-1}} * \dots * \frac{\partial h_{in,2}}{\partial h_{out,1}} * \frac{\partial h_{out,1}}{\partial h_{in,1}}$$

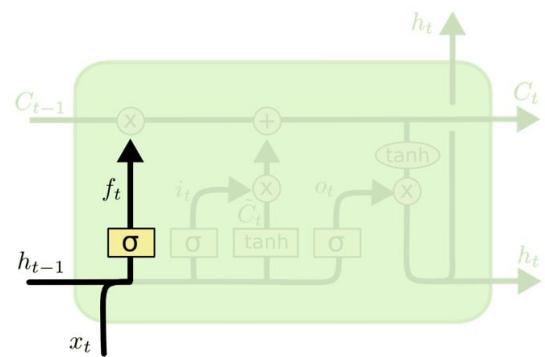
Value range of $\tanh(x)$ and $\tanh'(x)$



Long Short-Term Memory (LSTM)



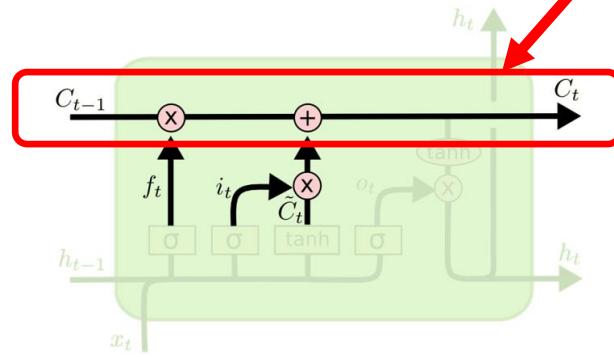
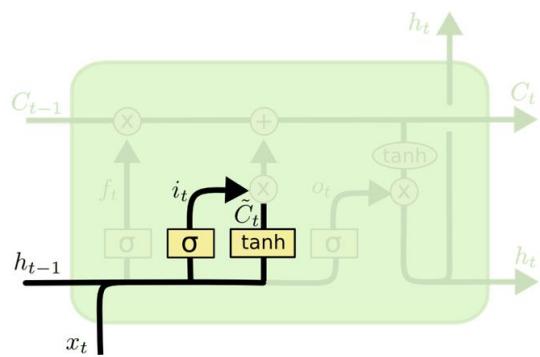
$$\begin{aligned}\delta^{k-1} &= \frac{\partial C_t}{\partial C_{k-1}} = \frac{\partial C_t}{\partial C_k} * \frac{\partial C_k}{\partial C_{k-1}} = \delta^k * \frac{\partial C_k}{\partial C_{k-1}} \\ &= \delta^k * (f_k + \dots)\end{aligned}$$



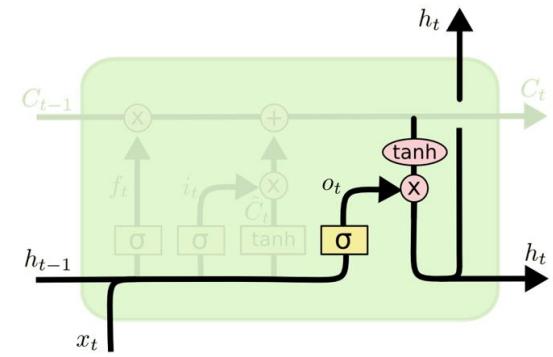
$$\mathbf{f}_t = \sigma(W_f \cdot [h_{t-1}; x_t] + b_f)$$

$$\mathbf{i}_t = \sigma(W_i \cdot [h_{t-1}; x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}; x_t] + b_C)$$



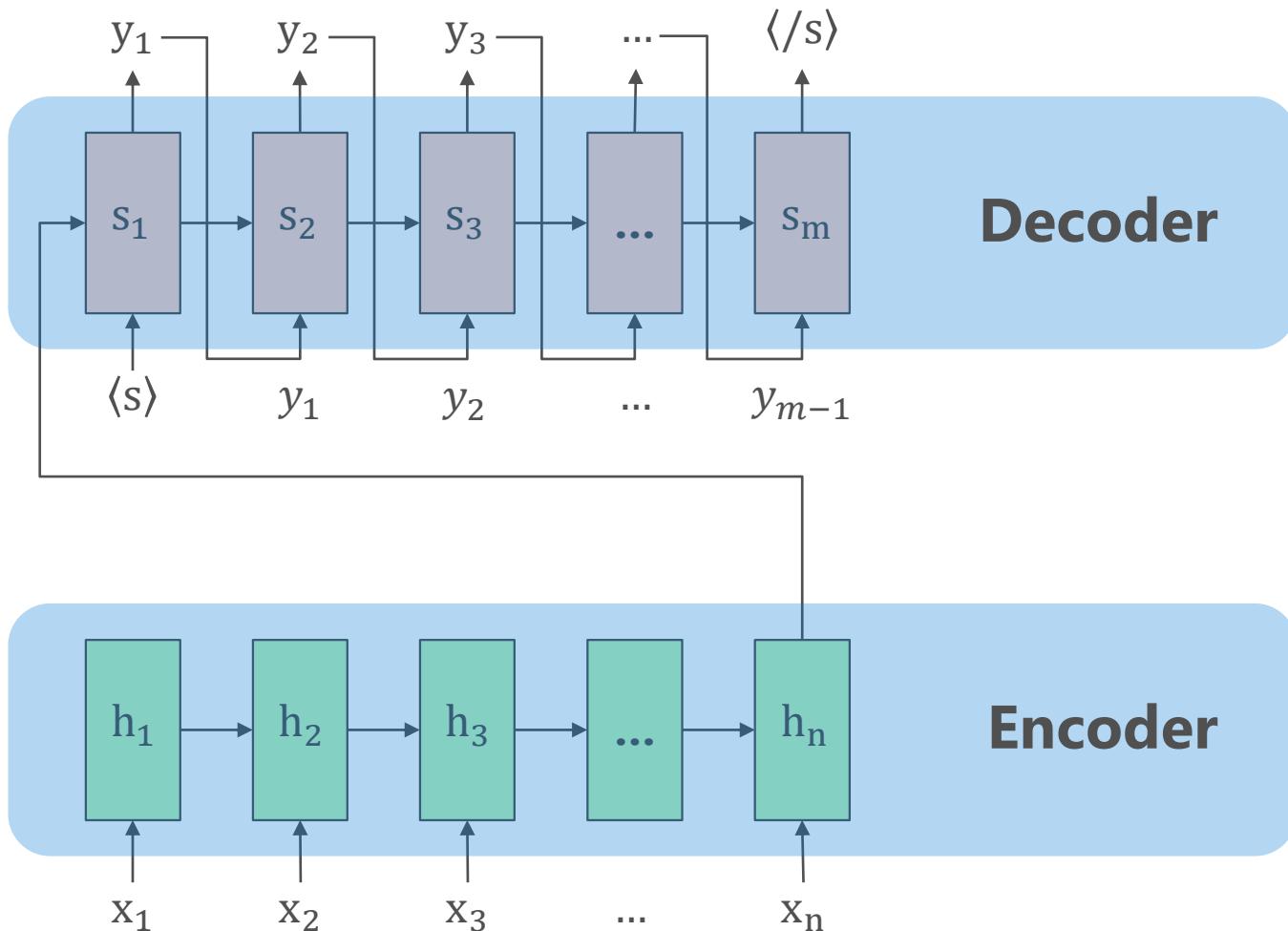
$$C_t = \mathbf{f}_t \cdot C_{t-1} + \mathbf{i}_t \cdot \tilde{C}_t$$



$$\mathbf{o}_t = \sigma(W_o \cdot [h_{t-1}; x_t] + b_o)$$

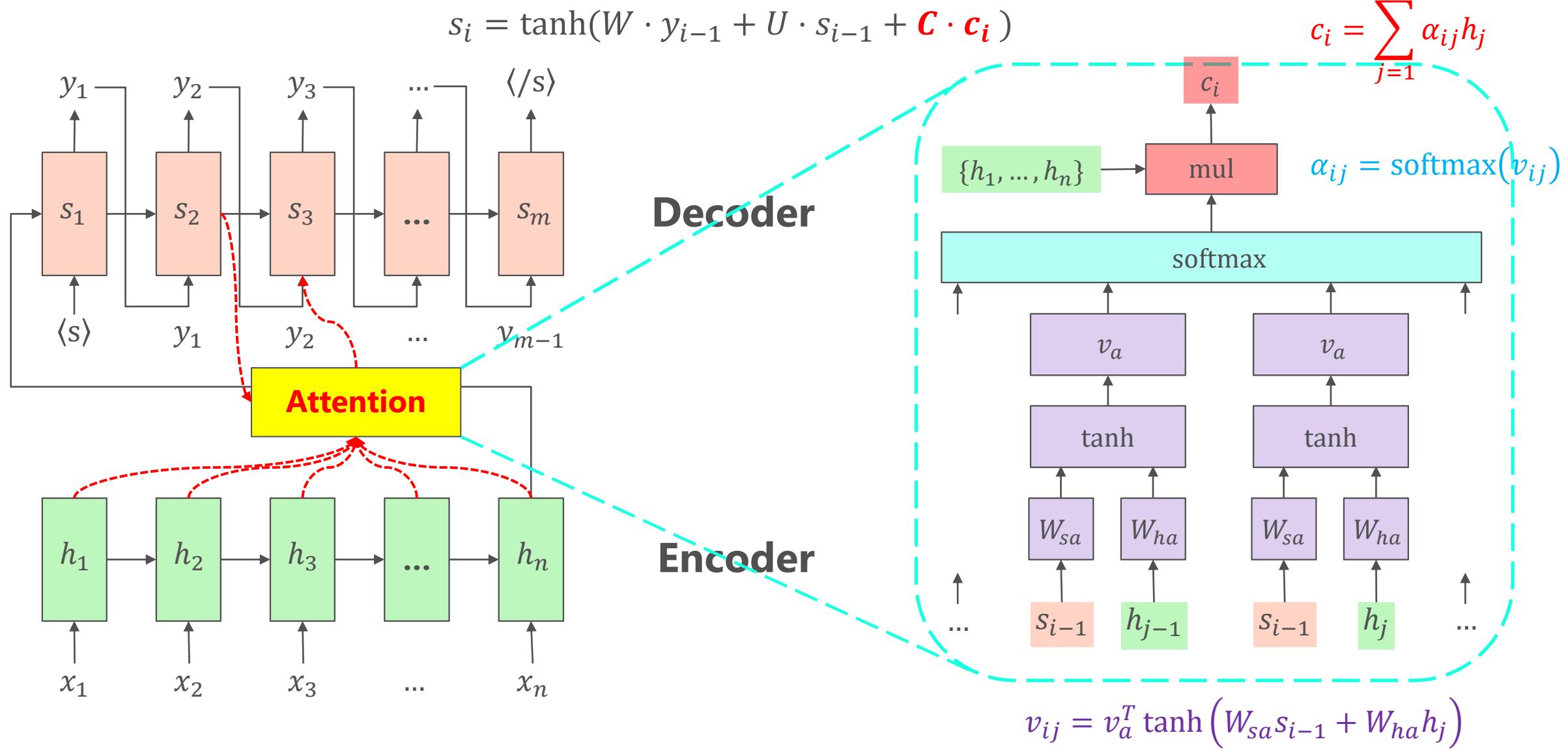
$$h_t = \mathbf{o}_t \cdot \tanh(C_t)$$

Encoder-Decoder Framework for Sequence Generation

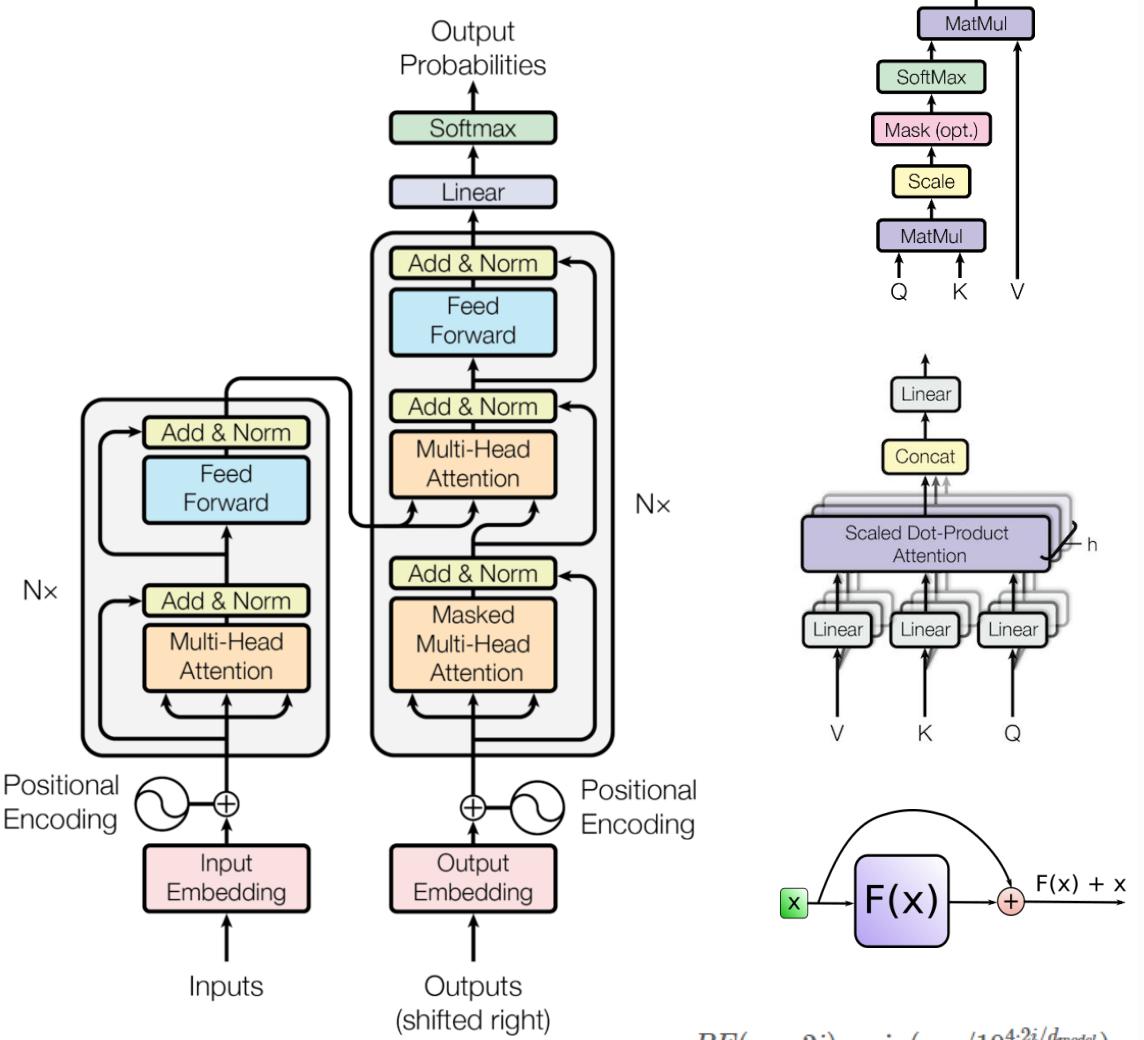


$$s_i = \tanh(W \cdot y_{i-1} + U \cdot s_{i-1})$$

Encoder-Decoder Framework with Attention



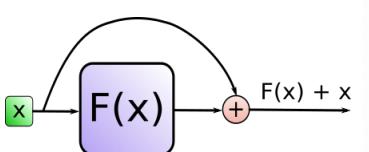
Transformer



$$PE(pos, 2i) = \sin(pos/10^{4 \cdot 2i/d_{model}})$$

$$PE(pos, 2i + 1) = \cos(pos/10^{4 \cdot 2i/d_{model}})$$

Transformer



Today's Agenda

- Deep Learning Basics for NLP
- **NLP with Pre-trained Embeddings**
- NLP with Knowledge Bases
- NLP with Commonsense
- Summary and Trend

Bag-of-words

- A text (such as a sentence or a document) is represented as the bag (multiset) of its words, disregarding grammar and even word order but keeping multiplicity.

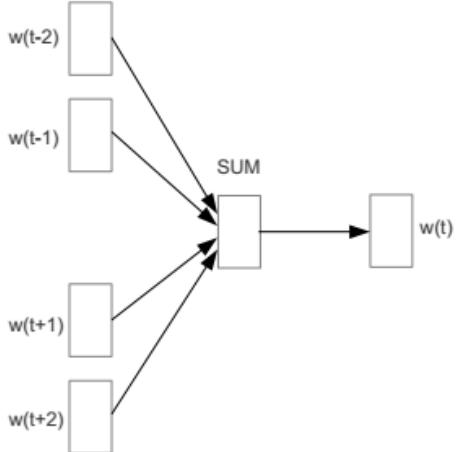
the dog is on the table

0	0	1	1	0	1	1	1
are	cat	dog	is	now	on	table	the

Word2Vec

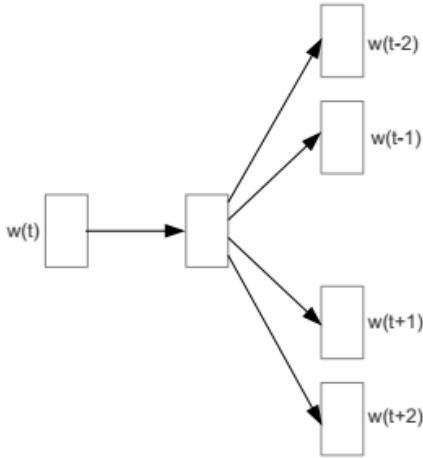
(Mikolov et al., 2013)

INPUT PROJECTION OUTPUT



CBOW

INPUT PROJECTION OUTPUT



Skip-gram

$$p^{CBOW}(w_t | w_{t-n}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+n}) = \frac{\exp(h^T v'_t)}{\sum_{w_i \in V} \exp(h^T v'_{w_i})}$$

$$p^{skip-gram}(w_{t+j} | w_t) = \frac{\exp(v_{w_t}^T v'_{w_{t+j}})}{\sum_{w_i \in V} \exp(v_{w_t}^T v'_{w_i})}$$

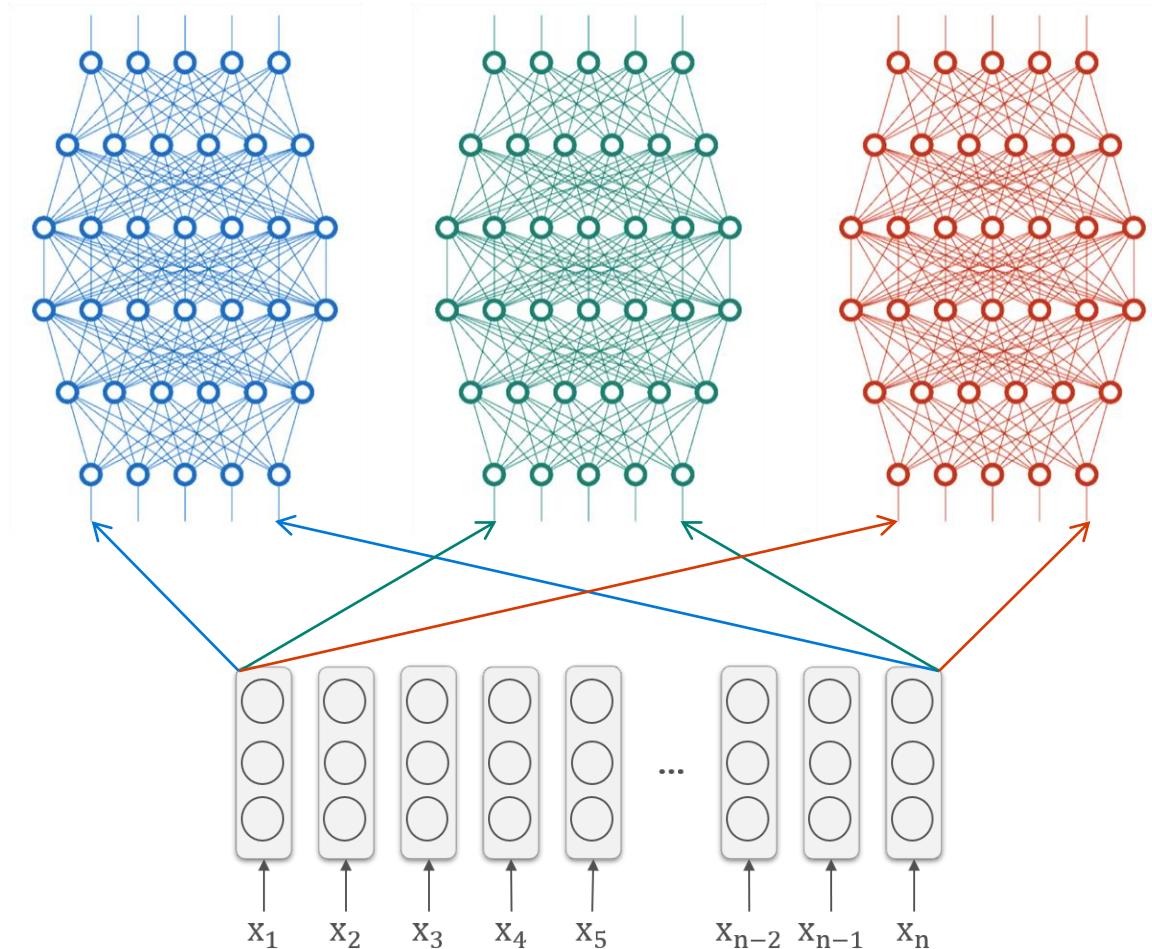
Objective function of CBOW: minimize $-\frac{1}{T} \sum_{t=1}^T \log p^{CBOW}(w_t | w_{t-n}, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+n})$

Objective function of Skip-gram: minimize $-\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p^{skip-gram}(w_{t+j} | w_t)$

Word2Vec in NLP

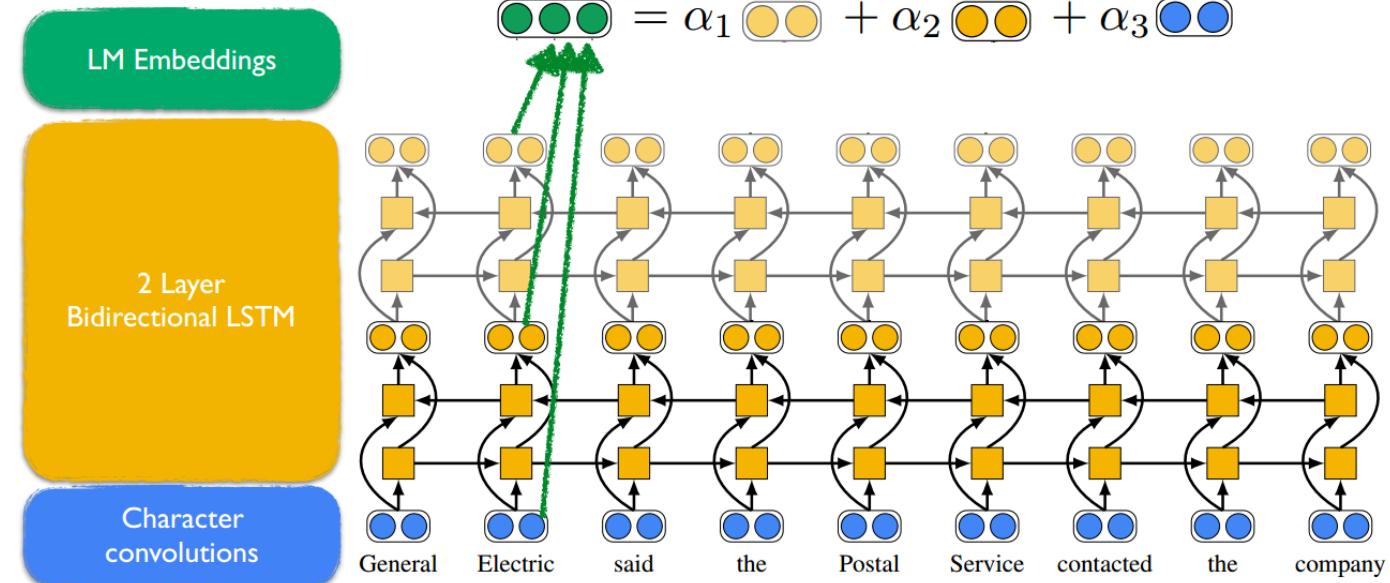
(Mikolov et al., 2013)

- Initialize the first layer of a neural network



ELMo

(Peters et al., 2018)



Forward LM: $p(t_1, \dots, t_N) = \prod_{k=1}^N p(t_k | t_1, \dots, t_{k-1})$

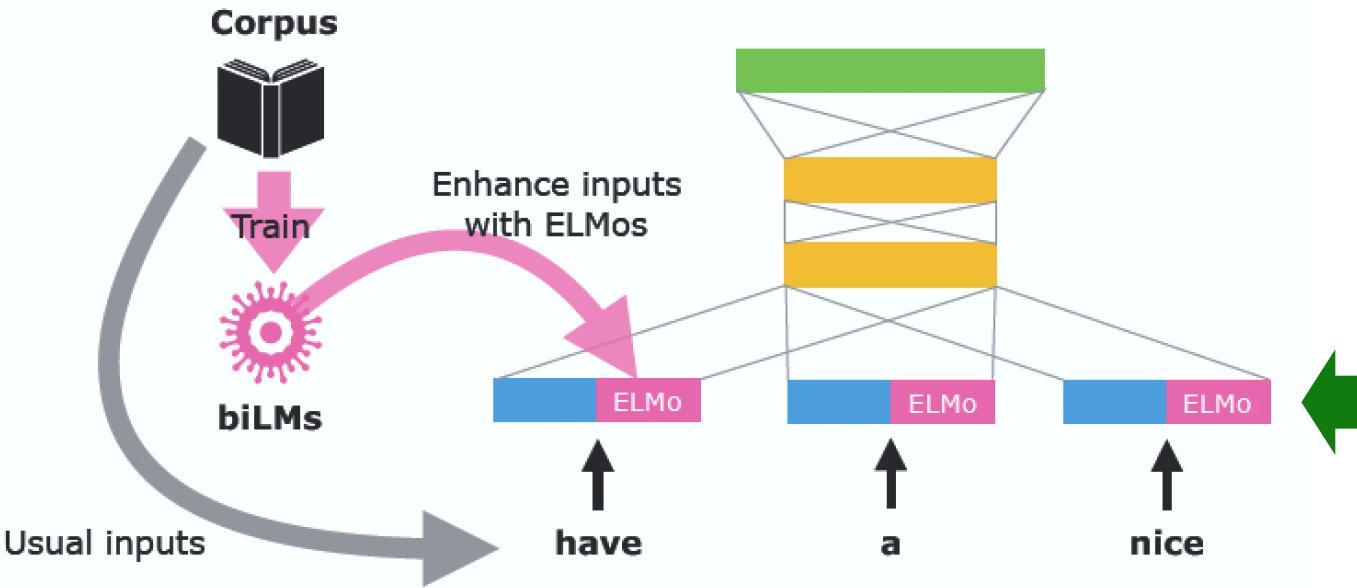
Backward LM: $p(t_1, \dots, t_N) = \prod_{k=1}^N p(t_k | t_{k+1}, \dots, t_N)$

Objective function: maximize $\sum_{t=1}^T (\log p(t_k | t_1, \dots, t_{k-1}; \Theta_x, \overrightarrow{\Theta_{LSTM}}, \Theta_{softmax}) + \log p(t_k | t_{k+1}, \dots, t_N; \Theta_x, \overleftarrow{\Theta_{LSTM}}, \Theta_{softmax}))$

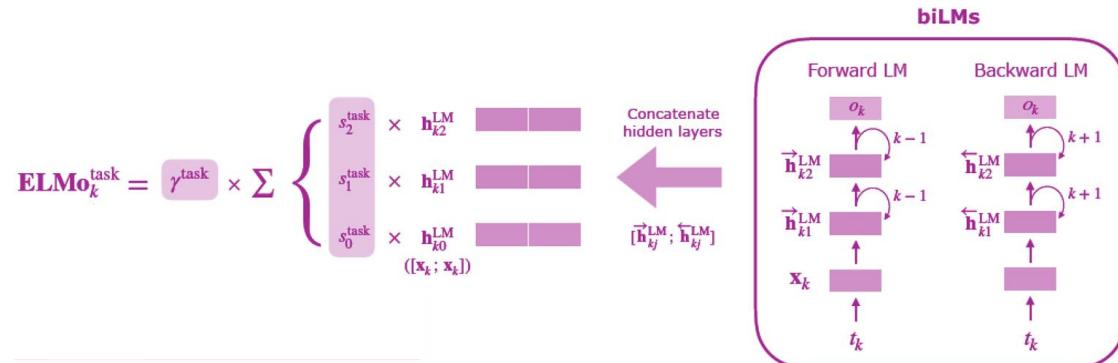
ELMo in NLP

(Peters et al., 2018)

ELMo can be integrated to almost all neural NLP tasks with simple concatenation to the embedding layer

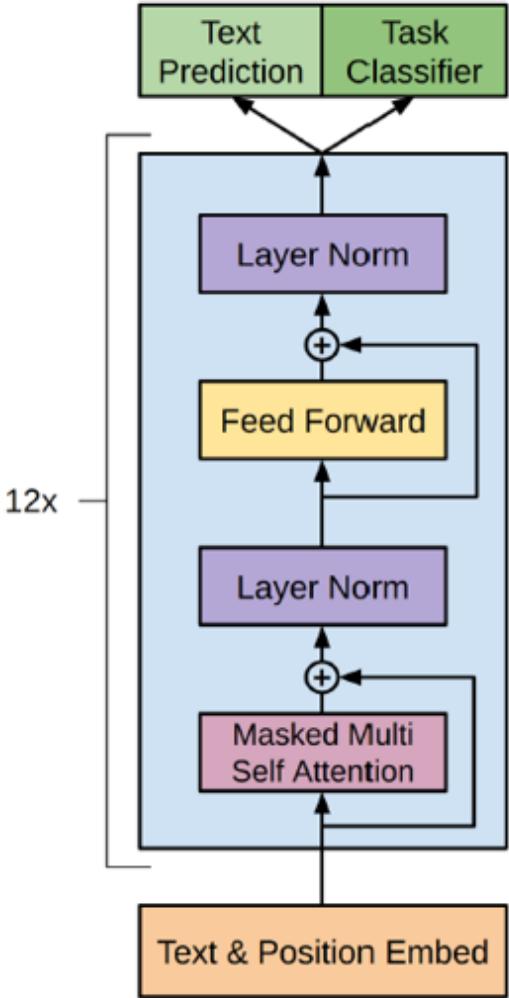


TASK	PREVIOUS SOTA	OUR BASELINE	ELMO +	INCREASE (ABSOLUTE/ RELATIVE)
			BASELINE	
SQuAD	Liu et al. (2017)	84.4	81.1	8.7 / 24.9%
SNLI	Chen et al. (2017)	88.6	88.0	0.7 / 5.8%
SRL	He et al. (2017)	81.7	81.4	3.2 / 17.2%
Coref	Lee et al. (2017)	67.2	67.2	3.2 / 9.8%
NER	Peters et al. (2017)	91.93 ± 0.19	90.15	92.22 ± 0.10 / 2.06 / 21%
SST-5	McCann et al. (2017)	53.7	51.4	54.7 ± 0.5 / 3.3 / 6.8%



OpenAI GPT

(Radford et al., 2018)



$$h_0 = UW_e + W_p$$

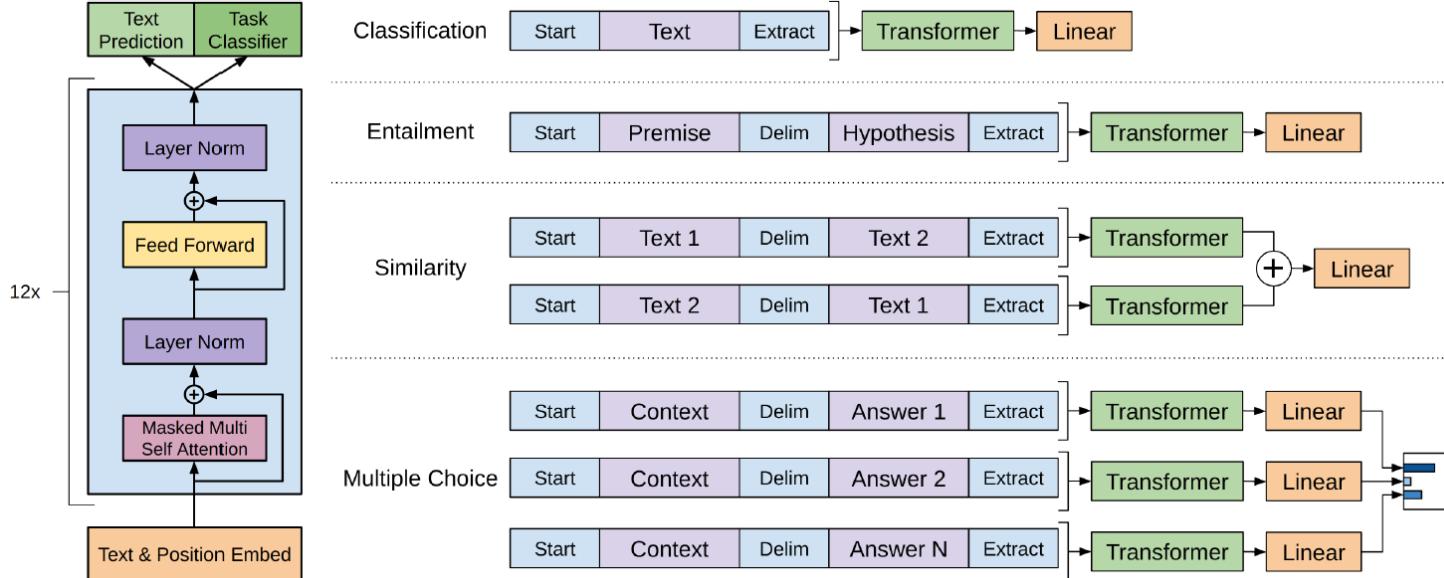
$$h_l = \text{transformer}_{\text{block}(h_{l-1})} \forall i \in [1, n]$$

$$p(u) = \text{softmax}(h_n W_e^T)$$

Objective function: maximize $\sum_i \log p(u_i | u_{i-k}, \dots, t_{i-1}; \Theta)$

OpenAI GPT in NLP

(Radford et al., 2018)



Method	MNLI-m	MNLI-mm	SNLI	SciTail	QNLI	RTE
ESIM + ELMo [44] (5x)	-	-	89.3	-	-	-
CAFE [58] (5x)	80.2	79.0	89.3	-	-	-
Stochastic Answer Network [35] (3x)	80.6	80.1	-	-	-	-
CAFE [58]	78.7	77.9	88.5	83.3	-	-
GenSen [64]	71.4	71.3	-	-	82.3	59.2
Multi-task BiLSTM + Attn [64]	72.2	72.1	-	-	82.1	61.7
Finetuned Transformer LM (ours)	82.1	81.4	89.9	88.3	88.1	56.0

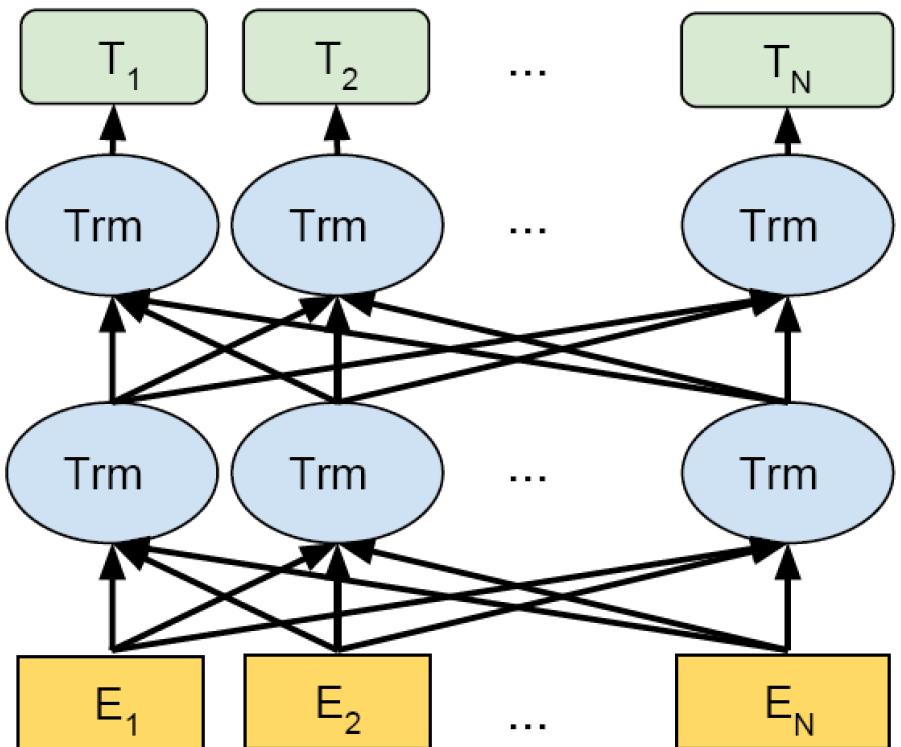
Method	Story Cloze	RACE-m	RACE-h	RACE
val-LS-skip [55]	76.5	-	-	-
Hidden Coherence Model [7]	77.6	-	-	-
Dynamic Fusion Net [67] (9x)	-	55.6	49.4	51.2
BiAttention MRU [59] (9x)	-	60.2	50.3	53.3
Finetuned Transformer LM (ours)	86.5	62.9	57.4	59.0

Method	Classification		Semantic Similarity		GLUE
	CoLA (mc)	SST2 (acc)	MRPC (F1)	STS-B (pc)	QQP (F1)
Sparse byte mLSTM [16]	-	93.2	-	-	-
TF-KLD [23]	-	-	86.0	-	-
ECNU (mixed ensemble) [60]	-	-	-	81.0	-
Single-task BiLSTM + ELMo + Attn [64]	35.0	90.2	80.2	55.5	66.1
Multi-task BiLSTM + ELMo + Attn [64]	18.9	91.6	83.5	72.8	63.3
Finetuned Transformer LM (ours)	45.4	91.3	82.3	82.0	70.3
					72.8

BERT

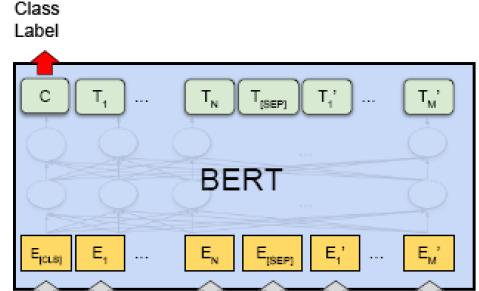
(Devlin et al., 2018)

BERT (Ours)

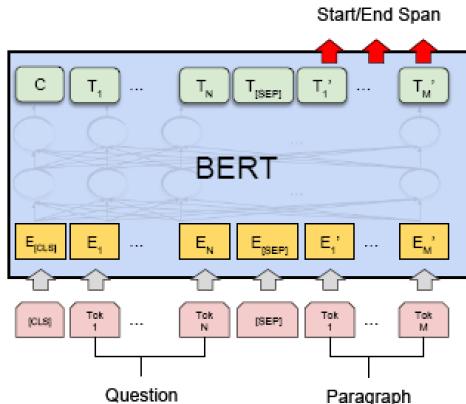


BERT in NLP

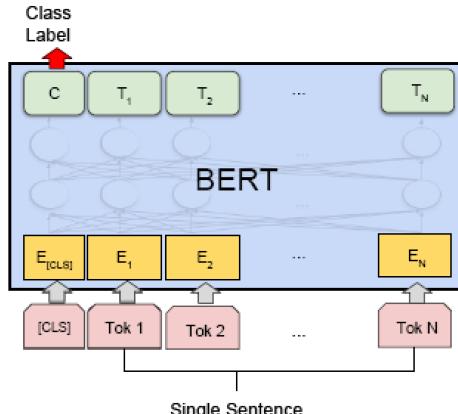
(Devlin et al., 2018)



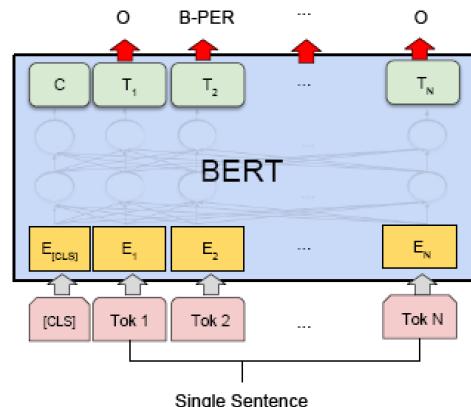
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



(c) Question Answering Tasks:
SQuAD v1.1



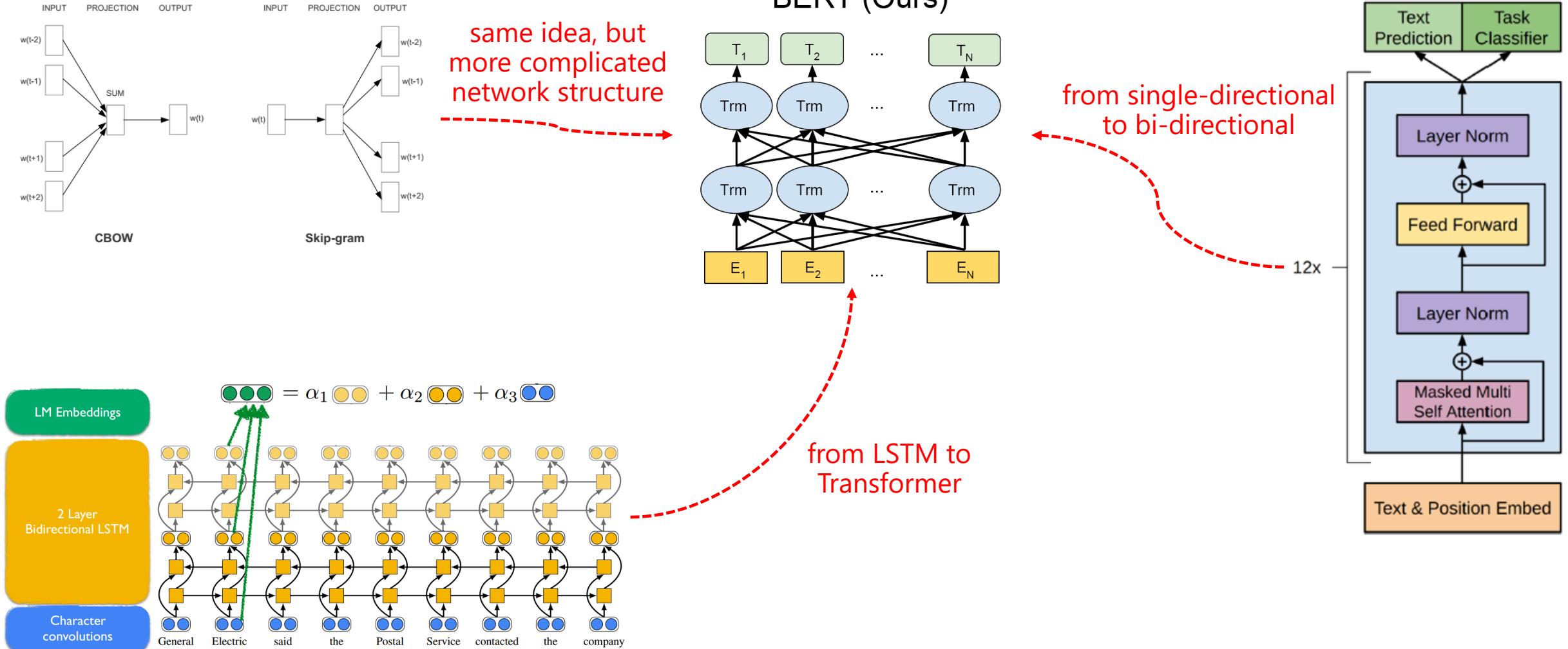
(b) Single Sentence Classification Tasks:
SST-2, CoLA



(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

System	Dev		Test	
	EM	F1	EM	F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
#1 Single - nlnet	-	-	83.5	90.1
#2 Single - QANet	-	-	82.5	89.3
Published				
BiDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQA)	84.2	91.1	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQA)	86.2	92.2	87.4	93.2

Comparison



Conclusion (1)

- Summary
 - Transformer is better than RNN
 - Bi-direction is better than single-direction
 - Size of training data, the larger the better
- Trend
 - Task-aware
 - Knowledge-aware

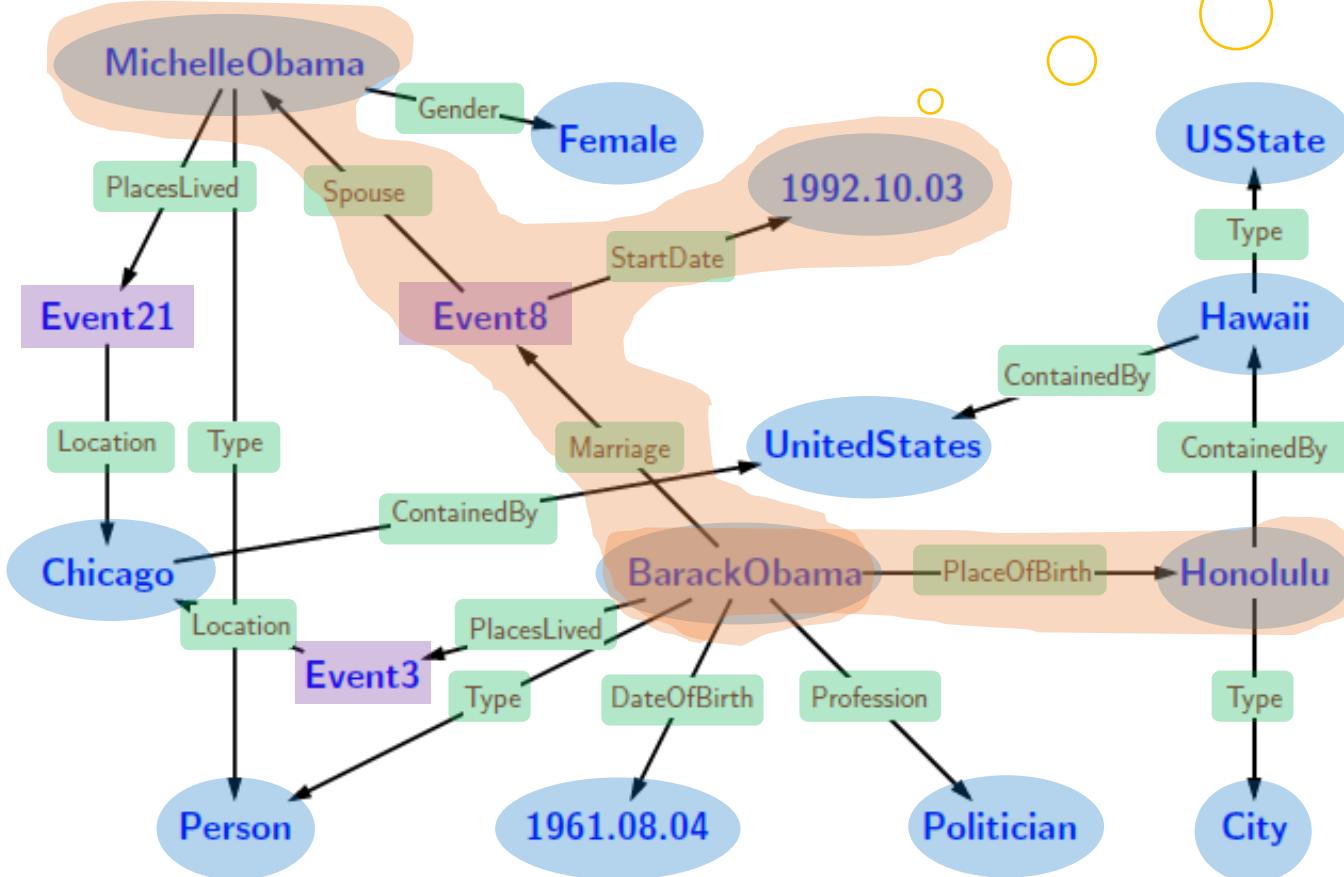
Break ^_ ^

Today's Agenda

- Deep Learning Basics for NLP
- NLP with Pre-trained Embeddings
- **NLP with Knowledge Bases**
- NLP with Commonsense
- Summary and Trend

Knowledge Graph

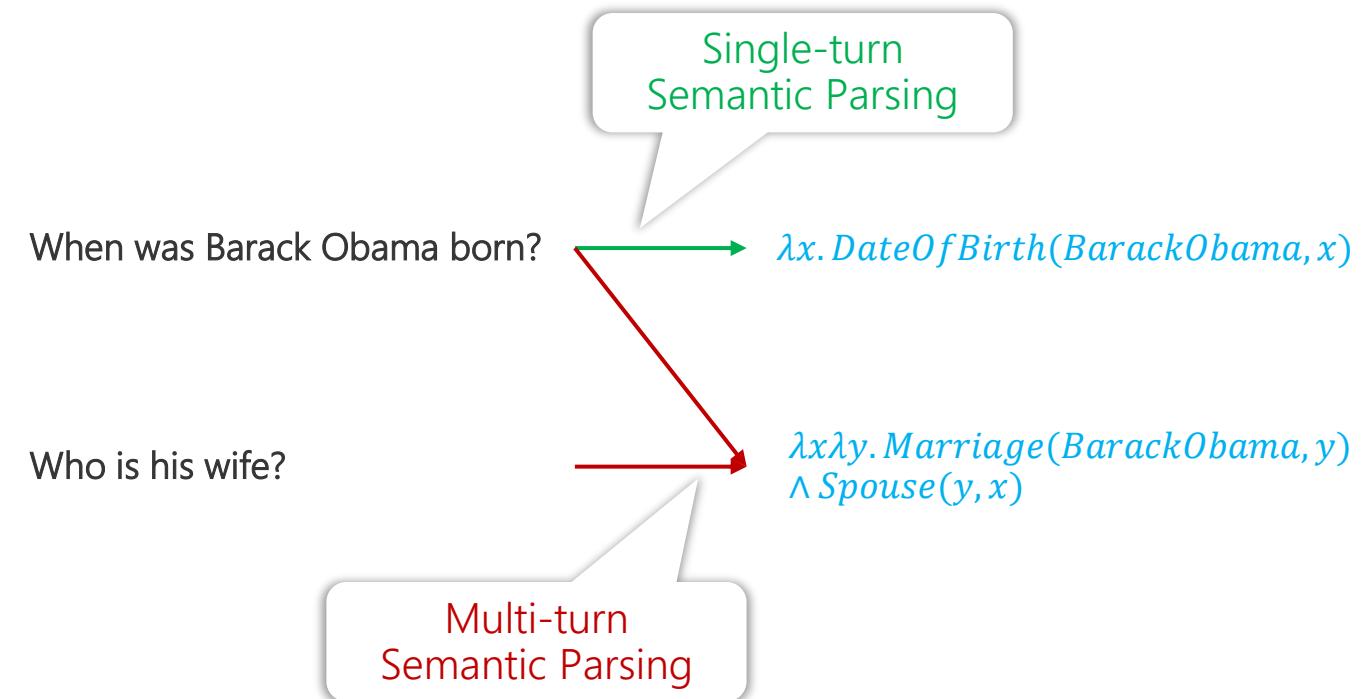
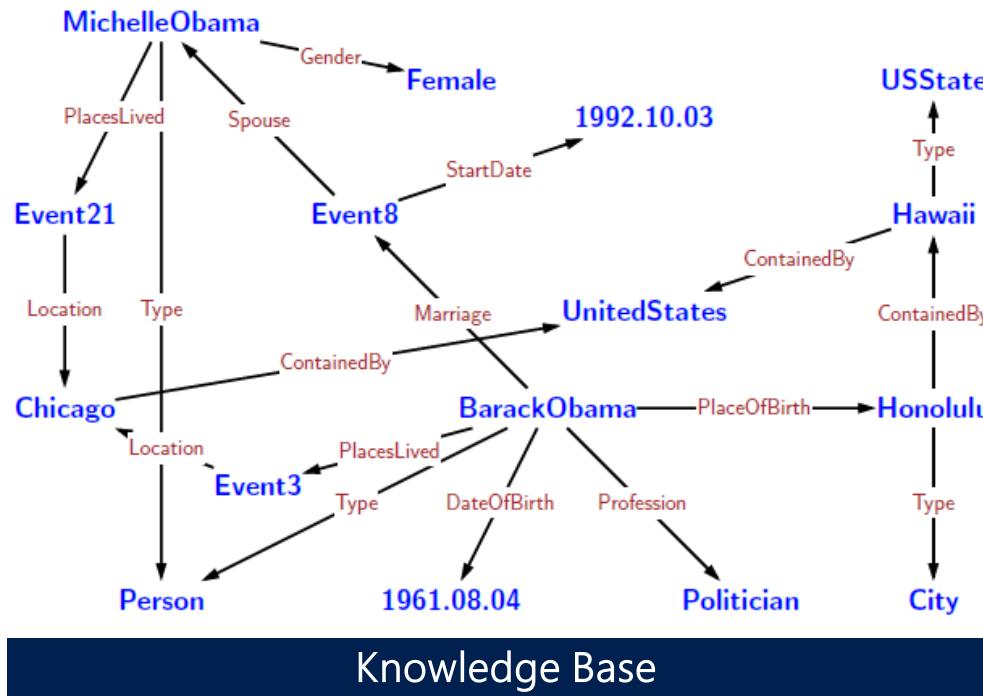
- Curated database with well-defined ontology



- Entity**
Objects/Values in the world
- Predicate**
Relation between two connected entities
- CVT (Compound Value Type)**
Not a real-world entity, but is used to collect multiple fields of an event
- Fact**
Triple, which connects two entities
Event, which connects multiple entities via a CVT node

Semantic Parsing with Knowledge Graph

- Convert NL utterances into machine executable LFs based on knowledge base



Lambda Calculus (λ -Calculus) as Logical Form (LF)

- λ -Calculus was introduced by Alonzo Church in 1930s
- Any computable function can be expressed using this formalism
- The core concept in λ -Calculus is “expression”
- An **expression** is defined recursively as follows

$\langle \text{expression} \rangle := \langle \text{constant} \rangle \mid \langle \text{variable} \rangle \mid \langle \text{function} \rangle \mid \langle \text{application} \rangle$

$\langle \text{function} \rangle := \lambda \langle \text{variable} \rangle. \langle \text{expression} \rangle$

$\langle \text{application} \rangle := \langle \text{expression} \rangle \langle \text{expression} \rangle$



Lambda Calculus: Constant

- Represent objects in the world

China, Bill Gates, Mount Everest, 2017, ...

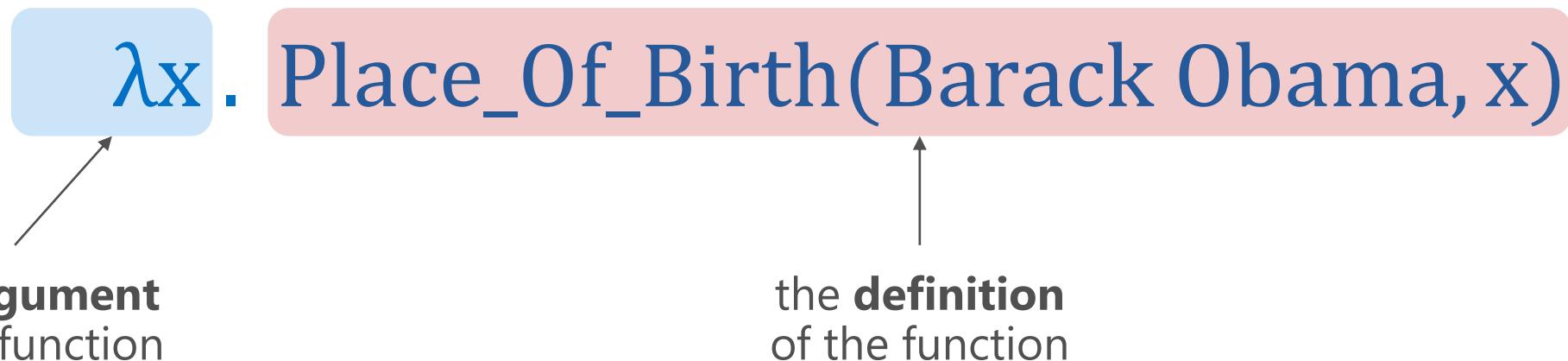
Lambda Calculus: Variable

- Represent object variables

x, y, z, ...

Lambda Calculus: Function

- Represent a function, and return the output of the function



Lambda Calculus: Application

- Apply the first expression to the second expression

$\lambda x \lambda y. \text{Place_Of_Birth}(x, y) \quad \lambda x. (x = \text{Barack Obama})$



$\lambda y. \text{Place_Of_Birth}(\text{Barack Obama}, y)$

Transforming Natural Language into Logical Form (λ -Calculus)

Natural Language

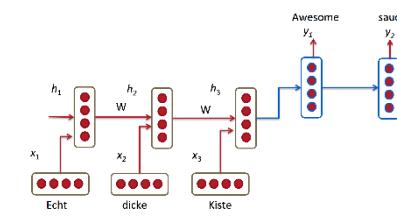
What city was Obama born ?

Semantic Parsing



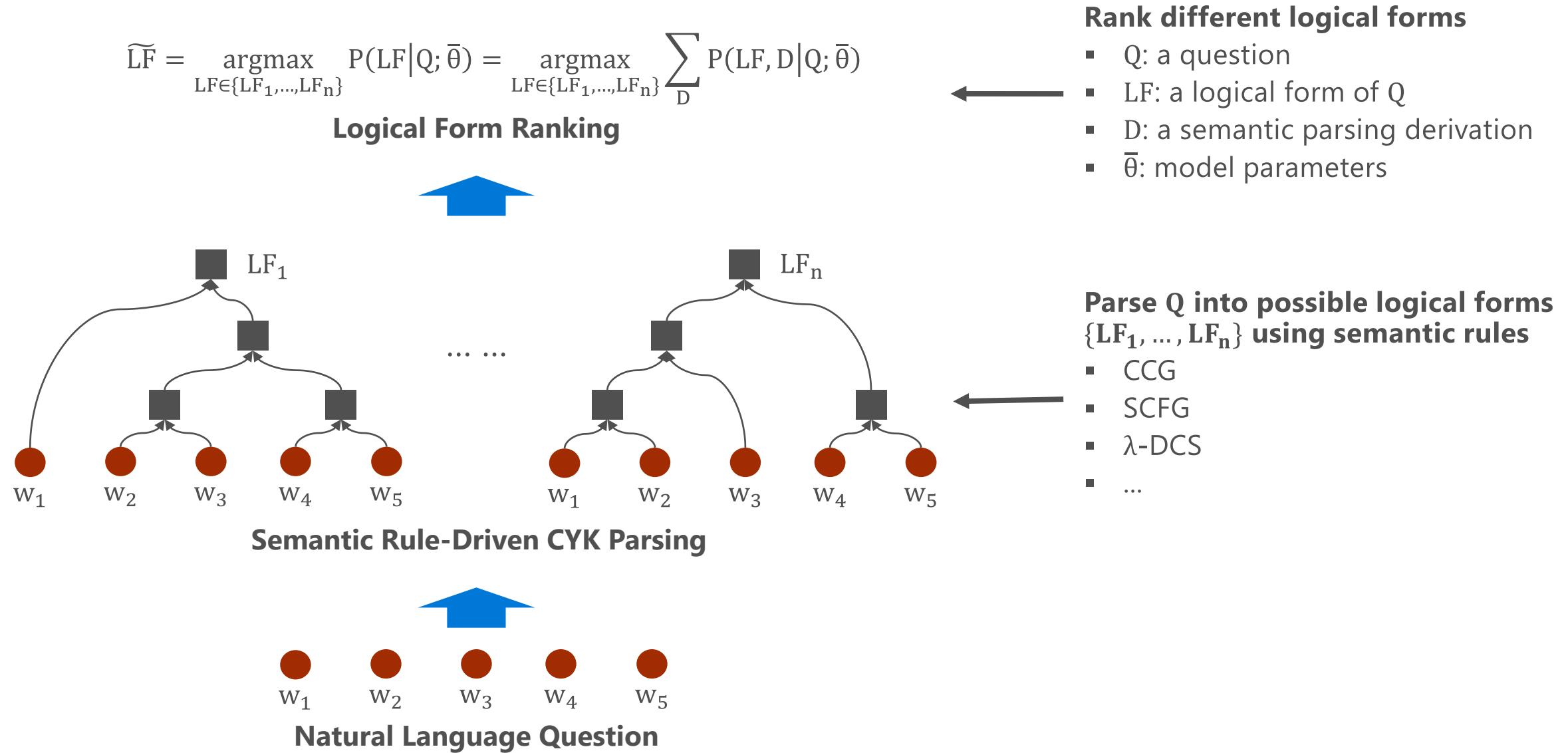
Grammar-based
Semantic Parsing

Logical Form

$$\lambda x. \text{Type}(\text{City}, x) \wedge \text{Place_of_Birth}(\text{Barack Obama}, x)$$


Neural Network-based
Semantic Parsing

Generic Framework of Grammar-based Semantic Parsing



Combinatorial Categorial Grammar (CCG)

- CCG captures **syntactic** and **semantic** information jointly

A CCG Rule Example

$$\text{border} := (S \setminus NP) / NP : \lambda x \lambda y. \text{Border}(x, y)$$

- Match natural language input

natural language

syntax

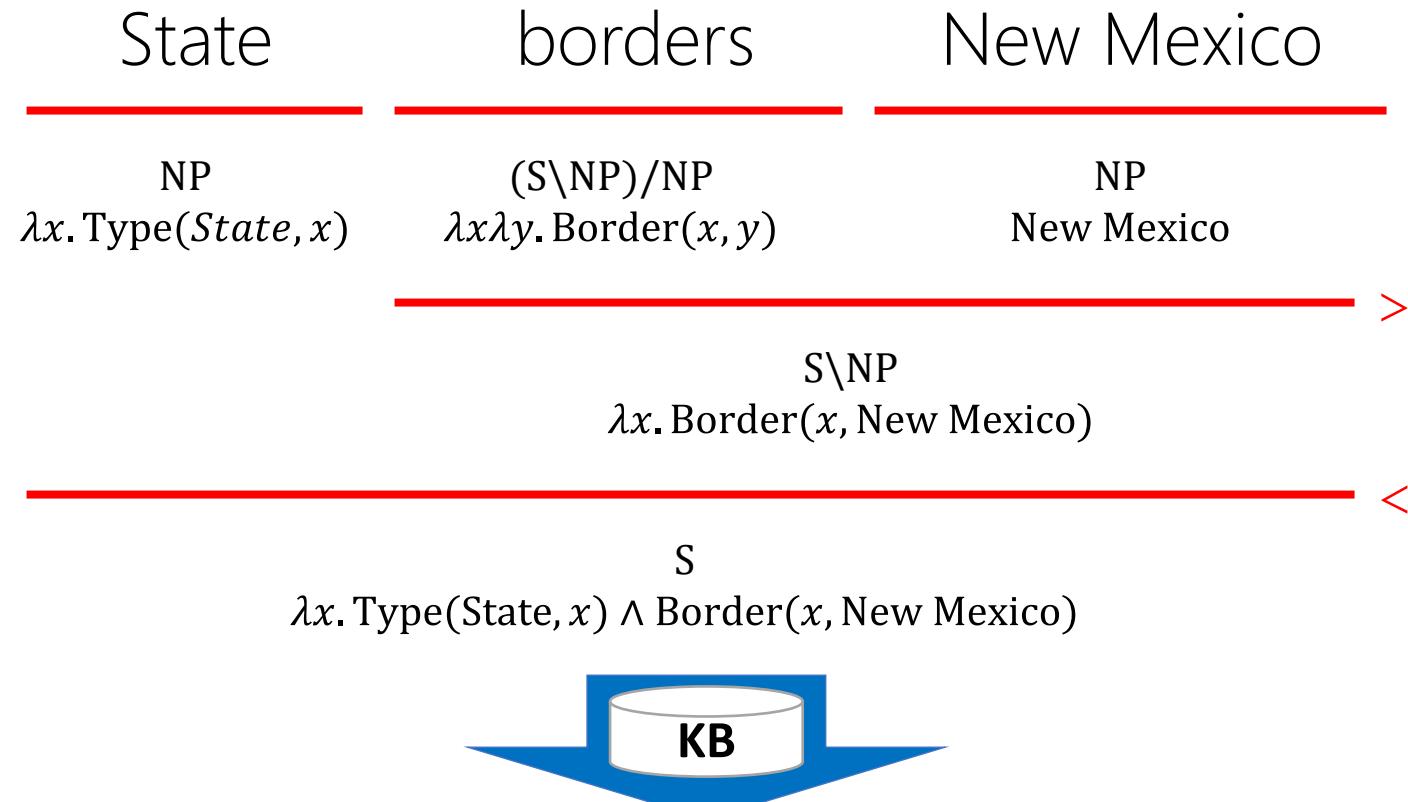
semantics

- Syntactic symbols: S, N, NP, ADJ and PP
- Syntactic combinator: / and \
- Slashes specify combination orders and directions

- λ -Calculus expression
- Sematic types are the logical forms of the natural language parts

Semantic Parsing with CCG

(Zettlemoyer and Collins, 2007; Kwiatkowski et al., 2011)



Arizona, Colorado, Oklahoma, Texas

CCG Rule Mining

- Input (<question, logical form> pairs)

Texas borders New Mexico
borders(texas, new_mexico)

use rules to extract all possible <Q, LF> pairs

Category Rules

Input Trigger	Output Category
constant c	$NP : c$
arity one predicate p	$N : \lambda x.p(x)$
arity one predicate p	$S \setminus NP : \lambda x.p(x)$
arity two predicate p	$(S \setminus NP) / NP : \lambda x.\lambda y.p(y, x)$
arity two predicate p	$(S \setminus NP) / NP : \lambda x.\lambda y.p(x, y)$
arity one predicate p	$N / N : \lambda g.\lambda x.p(x) \wedge g(x)$
arity two predicate p and constant c	$N / N : \lambda g.\lambda x.p(x, c) \wedge g(x)$
arity two predicate p	$(N \setminus N) / NP : \lambda x.\lambda g.\lambda y.p(y, x) \wedge g(x)$
arity one function f	$NP / N : \lambda g.\text{argmax/min}(g(x), \lambda x.f(x))$
arity one function f	$S / NP : \lambda x.f(x)$

- Output (CCG rules)

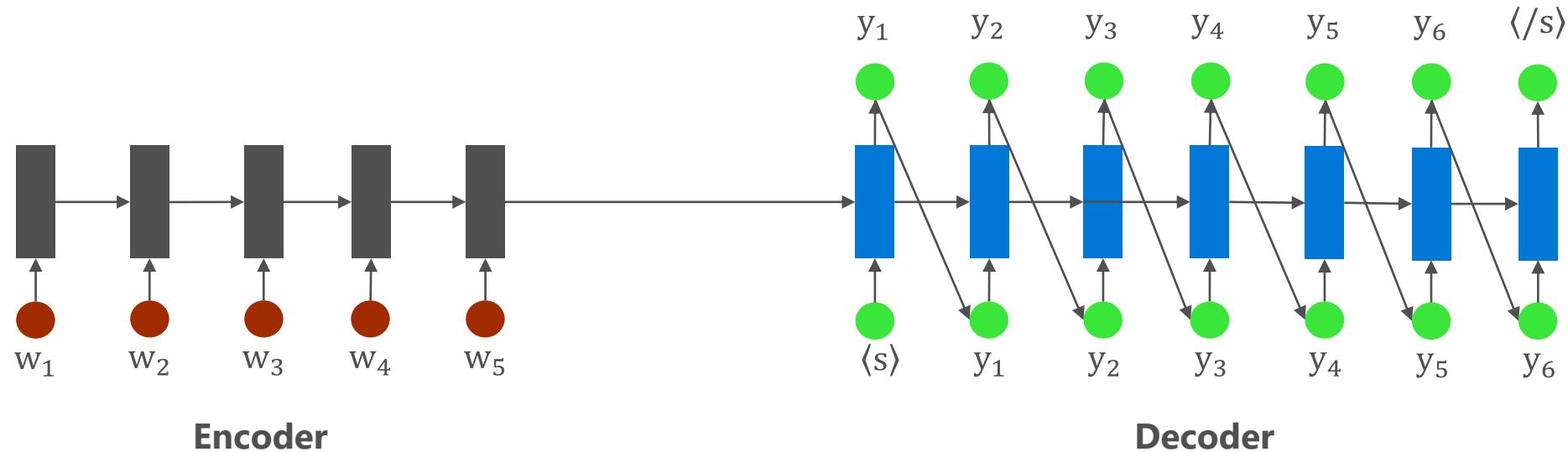
Texas := $NP : \text{texas}$
borders := $(S \setminus NP) / NP : \lambda x.\lambda y.\text{borders}(y, x)$
New Mexico := $NP : \text{new_mexico}$

1. maximize the likelihood: $\prod_i P_w(LF_i | Q_i) = \prod_i \sum_d P_w(LF_i, d | Q_i)$

2. keep CCG rules that occur in the highest scoring derivations of training data

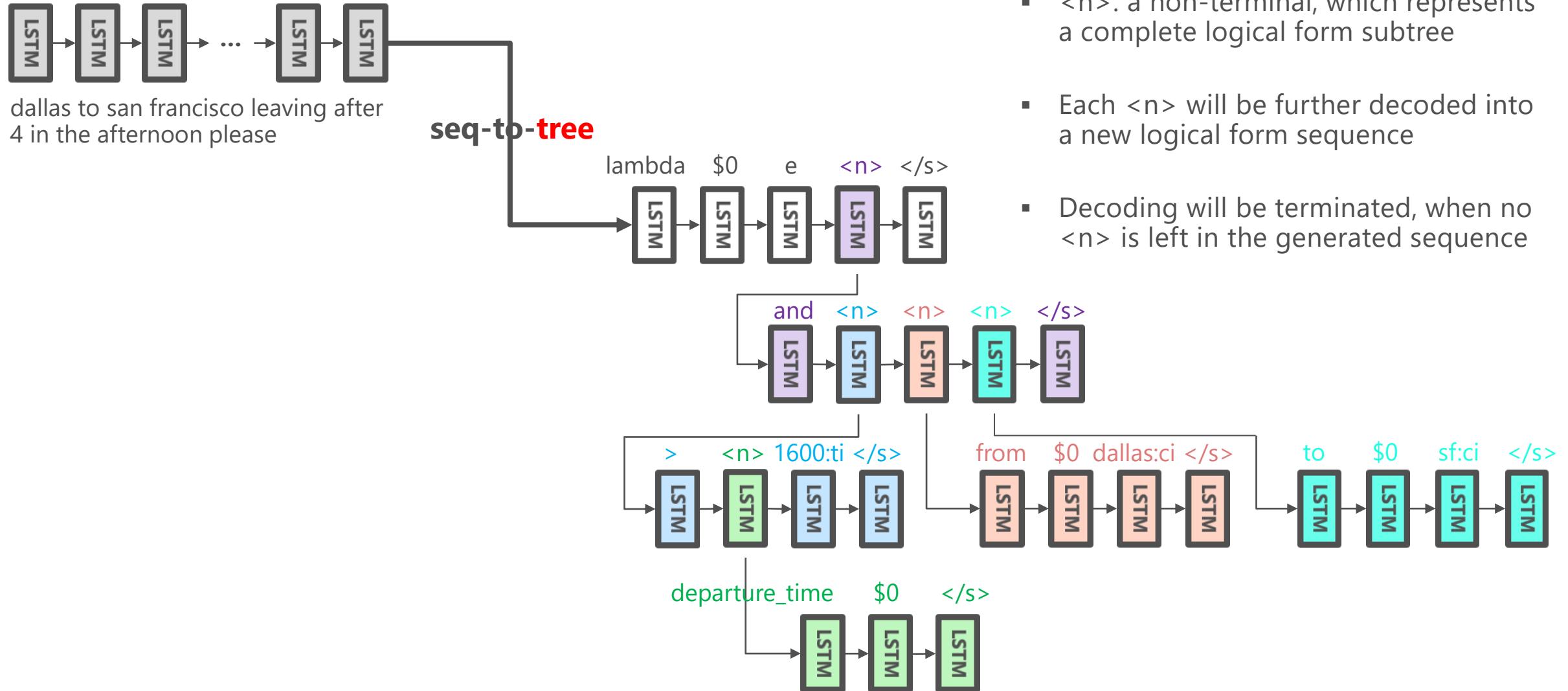
Generic Framework of Neural Network-based Semantic Parsing

- Perform semantic parsing as neural machine translation
- **Encoder** encodes each question into hidden states using RNN
- **Decoder** generates a logical form word-by-word based on question encoding using RNN



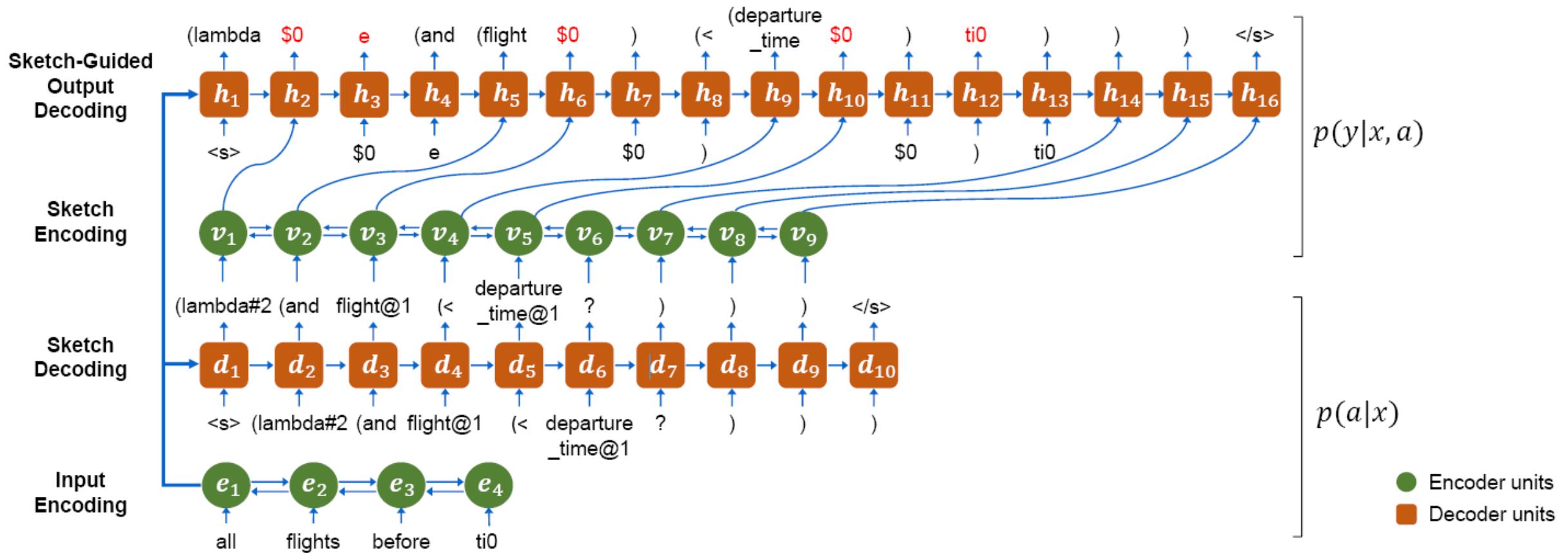
Single-turn Semantic Parsing with Sequence-to-Tree

(Dong and Lapata, 2016; Jia and Liang, 2016)



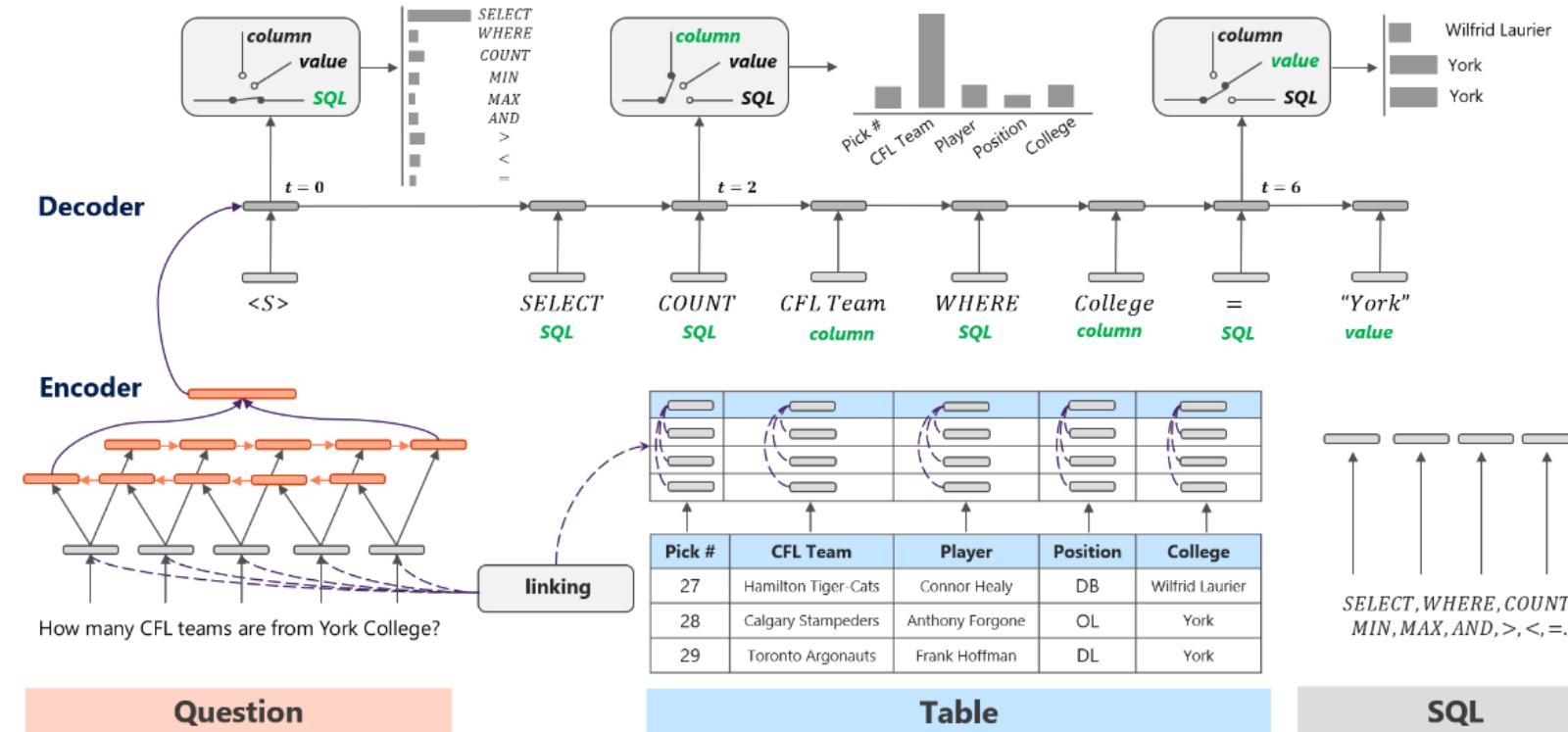
Single-turn Semantic Parsing with Sketch Decoding

(Dong and Lapata, 2018)



Single-turn Semantic Parsing with Multiple Gates

(Sun et al., 2018)



Datasets from Simple Questions to Complex Questions

The SimpleQuestions dataset

This section proposes SimpleQuestions, a dataset collected for research in automatic question answering with human generated questions. Details and baseline results on this dataset can be found in the paper:

Antoine Bordes, Nicolas Usunier, Sumit Chopra and Jason Weston. [Large-Scale Simple Question answering with Memory Networks](#), arXiv:1506.02075.

The SimpleQuestions dataset consists of a total of 108,442 questions written in natural language by human English-speaking annotators each paired with a corresponding fact, formatted as (subject, relationship, object), that provides the answer but also a complete explanation. Facts have been extracted from the Knowledge Base [Freebase](#). We randomly shuffle these questions and use 70% of them (75910) as training set, 10% as validation set (10845), and the remaining 20% as test set.

Here are some examples of questions and facts:

* What American cartoonist is the creator of Andy Lippincott?
Fact: (andy_lip

* Which forest is
Fact: (fires_cr

* What does Jimmy
Fact: (jimmy_ne

* What dietary re
Fact: (kimchi,

SimpleQuestions Nearly Solved: A New Upperbound and Baseline Approach

Michael Petrochuk
University of Washington Department
of Computer Science & Engineering
mikep5@cs.washington.edu

Luke Zettlemoyer
University of Washington Department
of Computer Science & Engineering
lsz@cs.washington.edu

Because of data ambiguity, the upper-bound performance on this benchmark at **83.4%** (Petrochuk and Zettlemoyer, 2018).



LC-QuAD

Largescale Complex Question Answering Dataset

[Download](#) OR [See Examples](#)

Data Characteristics

Current Version	1.0
Total Questions	5000
Unique Templates	38
Entities Covered	5042
Predicates Covered	615

Contact Us
In case you find any bug in our framework, or any issue with our dataset, please inform us on [Issues Page](#).
Contact: research_trivselid@uni-muenster.de

Examples

Q: What are the mascots of the teams participating in the turkish handball super league?
SELECT DISTINCT ?uri WHERE {
?x dbp:league dbp:Turkish_Handball_Super_League .
?x dbp:mascot ?uri .
}

Documentation & Usage Guides



Computer: Analyse the distribution of the pieces that we have, correcting for changes in star configurations over four billion years, then extrapolate for the missing pieces! (Star Trek, The Chase)

Leaderboard Paper Download Dataset



A dataset for answering complex questions that require reasoning over multiple web snippets.

ComplexWebQuestions is a new dataset that contains a large set of complex questions in natural language, and can be used in multiple ways:

1. By interacting with a search engine, which is the focus of our paper (Talmor and Berant, 2018);
2. As a reading comprehension task: we release 9,595,163 web snippets that are relevant for the questions, and were collected during the development of our model;
3. As a semantic parsing task: each question is paired with a SPARQL query that can be executed against Freebase to retrieve the answer.

The dataset contains 34,689 examples, each containing:

- A complex question
- Answers (including aliases)
- An average of 276.6 snippets per question
- A SPARQL query (against Freebase)

Sample Questions

- "Which school that Sir Ernest Rutherford attended has the latest founding date?"
- "what movies does Leo Howard play in and that is 113.0 minutes long?"
- "Where is the end of the river that originates in Shannon Pot?"

Coreference and Ellipsis Phenomena in Multi-turn Scenarios

Question Entity Coreference

- **Q1:** Who is the president of the United States?
- **A1:** Donald Trump
- **Q2:** what is its population?

Answer Entity Ellipsis

- **Q1:** What movie did Leonardo DiCaprio won an Oscar for?
- **A1:** The Revenant
- **Q2:** who is the director?

Answer Entity Coreference

- **Q1:** Who is the president of the United States?
- **A1:** Donald Trump
- **Q2:** How many children does he have?

Question Predicate Ellipsis

- **Q1:** Who is the president of the United States?
- **A1:** Donald Trump
- **Q2:** how about China?

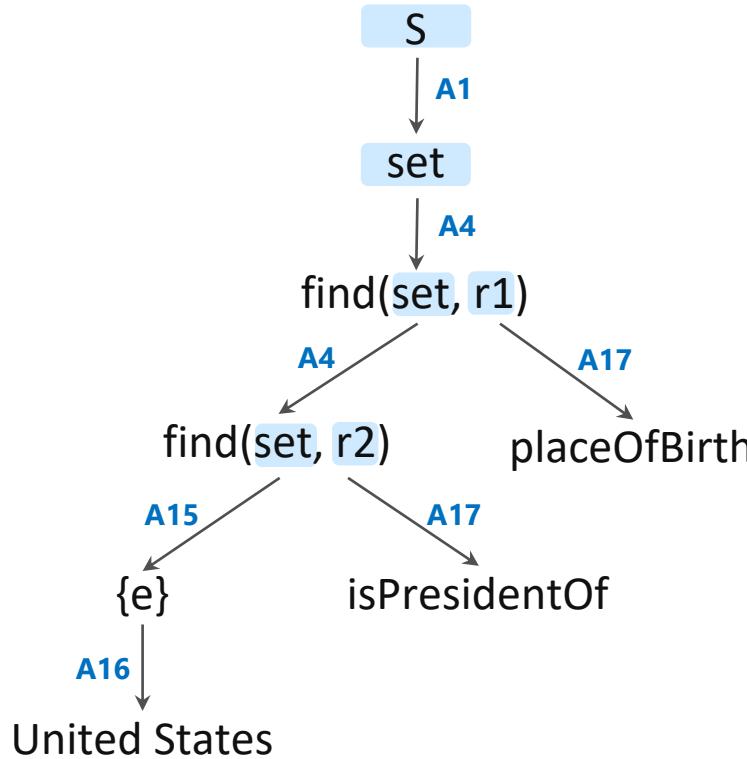
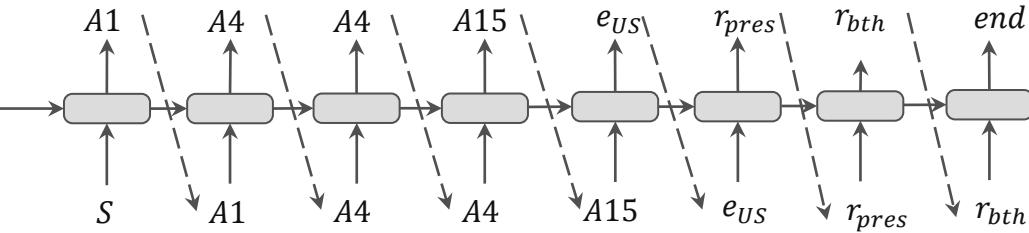
Question Subsequent Coreference

- **Q1:** Where did the president of the United States born?
- **A1:** New York City
- **Q2:** Where did he graduate from?

Multi-turn Semantic Parsing with Dialogue Actions & Memory

(Guo et al., 2018)

Where did the president of the United States born?

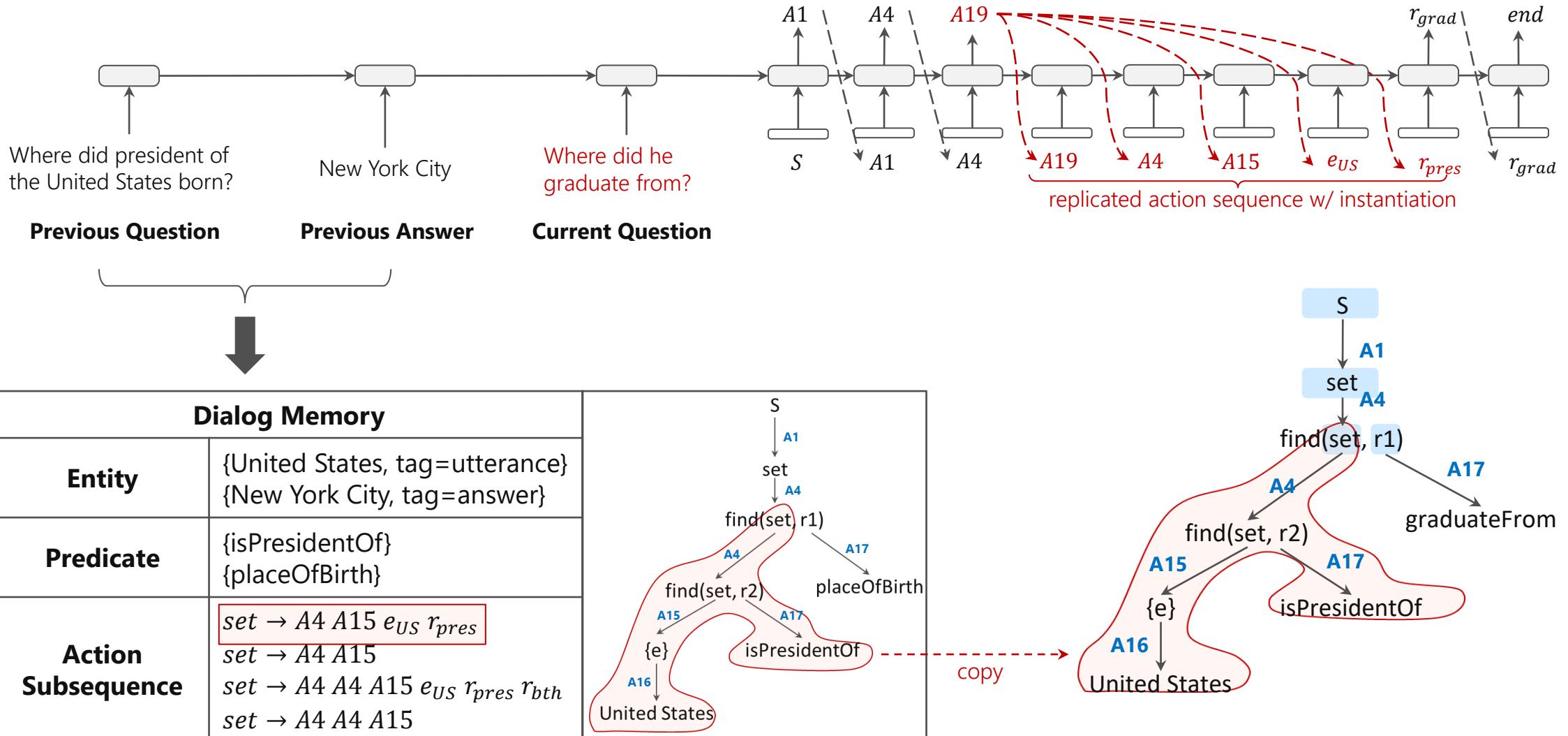


- A1: $S \rightarrow \text{Set}$
- A4: $\text{Set} \rightarrow \text{find}(\text{set}, r_1)$
- A4: $\text{Set} \rightarrow \text{find}(\text{set}, r_2)$
- A15: $\text{Set} \rightarrow \{e\}$
- A16: $e \rightarrow \text{United States}$
- A17: $r_2 \rightarrow \text{isPresidentOf}$
- A17: $r_1 \rightarrow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

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

Multi-turn Semantic Parsing with Dialogue Actions & Memory

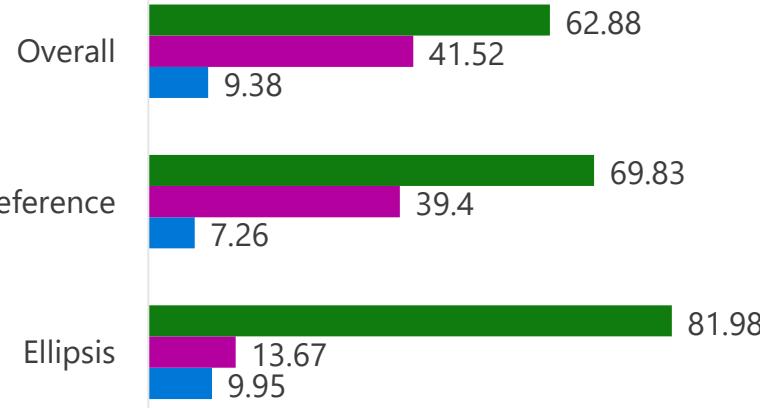
(Guo et al., 2018)



Evaluation on CSQA Dataset

CSQA Dataset (IBM, 2018)

Dialogs	200,000
Turns	1.6M
Entities in KB	12.8M
Unique relations	330
KB Tuples	21.2M
Entity Types	642



id	question type	current question + previous turn	predicted logical form
1	Simple Question (Direct)	Q1: N/A R1: N/A Q2: Who was the dad of Jorgen Ottesen Brahe?	<i>find({Jorgen Ottesen Brahe}, father)</i>
2	Simple Question (Coreferenced)	Q1: Who was the dad of Jorgen Ottesen Brahe? R1: Otte Brahe Q2: Who is the spouse of that one?	<i>find({Otte Brehe}, spouse)</i>
3	Simple Question (Ellipsis)	Q1: What is the profession of Mkhail Beliaiev? R1: Military personnel Q2: And also tell me about Brett MacLean	<i>find({Brett MacLean}, occupation)</i>
4	Logical Reasoning (All)	Q1: N/A R1: N/A Q2: Which administrative territories have diplomatic relations with Italy and are not Derikha present in?	<i>and(diff(find({Italy}, reverse(diplomatic relation)), find({Derikha}, country), find({administrative territories}, isA)))</i>
5	Quantitative Reasoning	Q1: N/A R1: N/A Q2: Which works did min number of people do the dubbing for?	<i>argmin(find({voice actor}, isa), reverse(work))</i>
6	Comparative Reasoning	Q1: N/A R1: N/A Q2: Which musical instruments are played by more number of people than electronic keyboard?	<i>larger(find({musical instruments}, isA), reverse(instrument), count(and(find({electronic keyboard}, reverse(instrument)), find({people}, isA))))</i>
7	Verification (Boolean)	Q1: N/A R1: N/A Q2: Is Arizona Coyotes present in United States of America?	<i>in(Arizona Coyotes , find({United States of America}, reverse(country)))</i>
8	Quantitative Reasoning (Count)	Q1: How many people have birthplace at Provence? R1: 15 Q2: And how about Peterborough?	<i>copy(count(find({Peterborough}, reverse(place of birth))))</i>
9	Comparative Reasoning (Count)	Q1: How many musical instruments are played by greater number of people than Body percussion ? R1: 30 Q2: And also tell me about timpani?	<i>copy(count(larger(find({musical instrument}, isA), reverse(instrument) , count(find({timpani}, reverse(instrument))))))</i>

■ D2A ■ D2A w/o DM ■ S2S

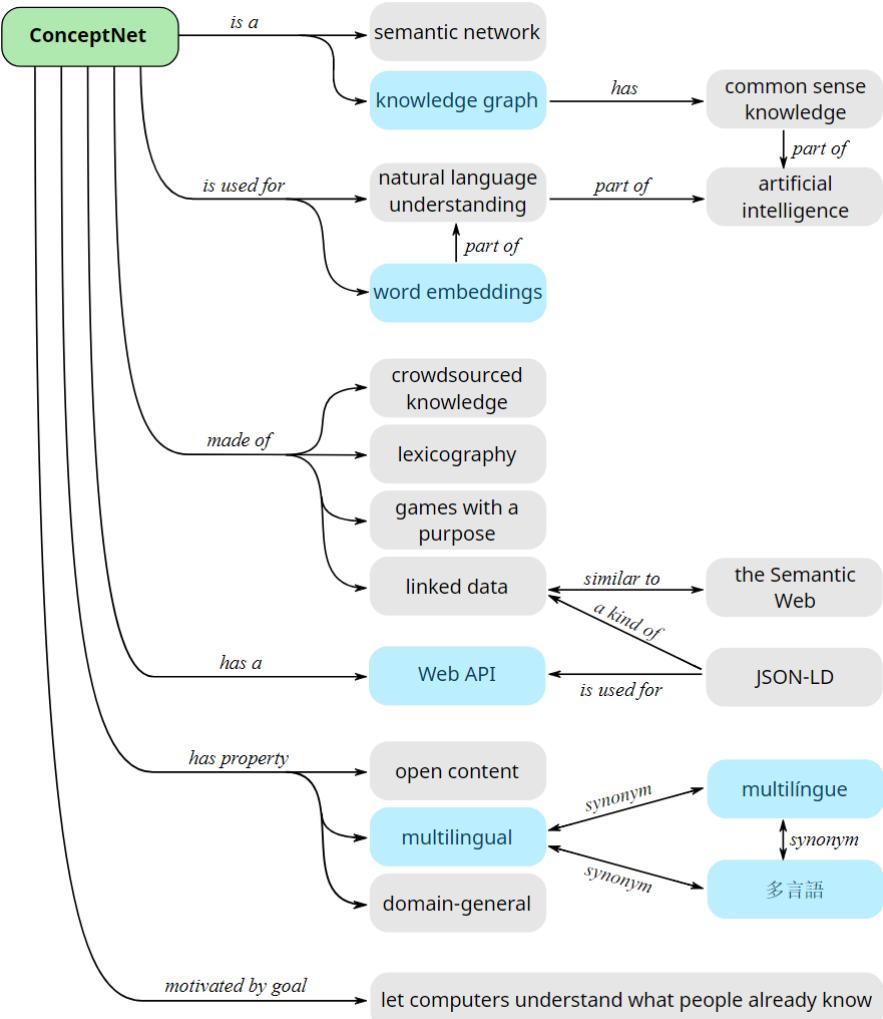
Conclusion (2)

- Summary
 - Semantic parsing is the core task of NLP
 - Dataset set is a big issue for model training
- Trend
 - From grammar to neural work
 - From single-relation to complex-relation
 - From single-turn to multi-tune
 - From single-pass to multi-pass
 - From knowledge graphs to web tables

Today's Agenda

- Deep Learning Basics for NLP
- NLP with Pre-trained Embeddings
- NLP with Knowledge Bases
- **NLP with Commonsense**
- Summary and Trend

Commonsense Knowledge



ConceptNet

WEBCHILD Commonsense Browser

car

Guess the concept: car

Domain Comparable Physical Part Activity Property Location

Ask me!

TYPE OF	motor_vehicle
Related to artifact, under the category of vehicles	
COMPARABLES	vehicle,car car,gasoline car,automobile car,auto car,driver More
ACTIVITIES	drive car buy car leave car park car see car
HAS PHYSICAL PARTS	accelerator accelerator air bag airbrake air horn More
IS PHYSICAL PART OF	train product hotel universe vehicle More
HAS SUBSTANCE	metallic element steel wood silica glass More
IN SPATIAL PROXIMITY WITH	road ground corner street air More
PHYSICAL PROPERTIES	rust sensitive heavy cool bright More
ABSTRACT PROPERTIES	detailed variable welcome happy bold More
OTHER PROPERTIES	long dangerous catastrophic cheap big More
ASSOCIATED WITH COUNTRY	united_states denmark europe united_ukraine russia More

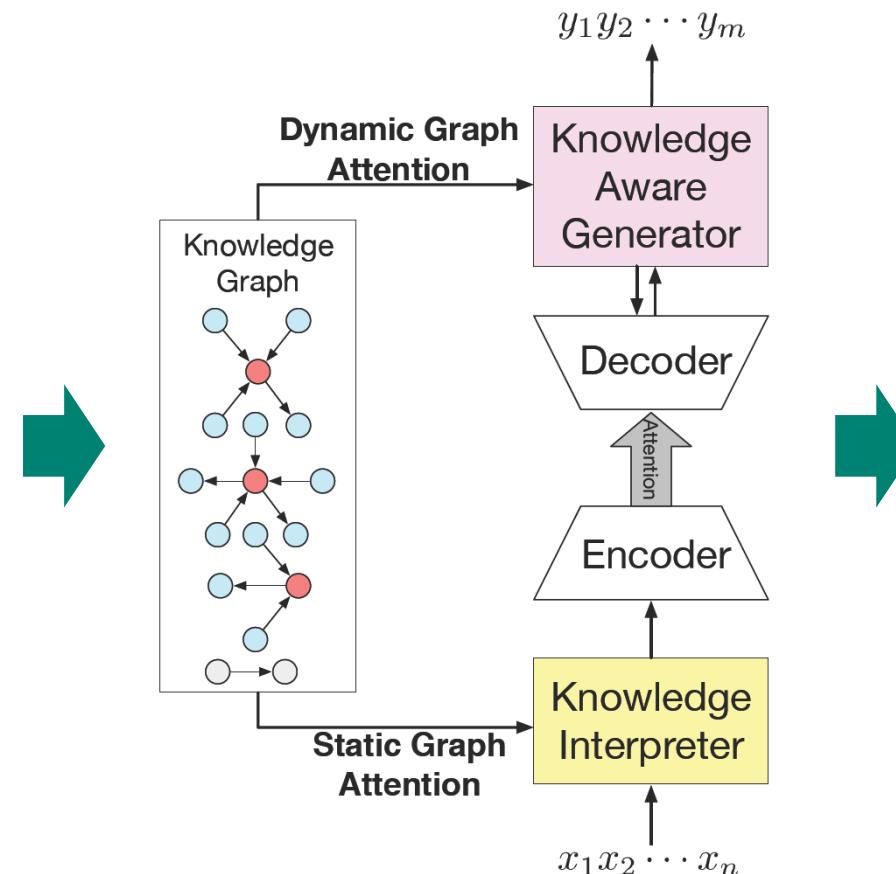
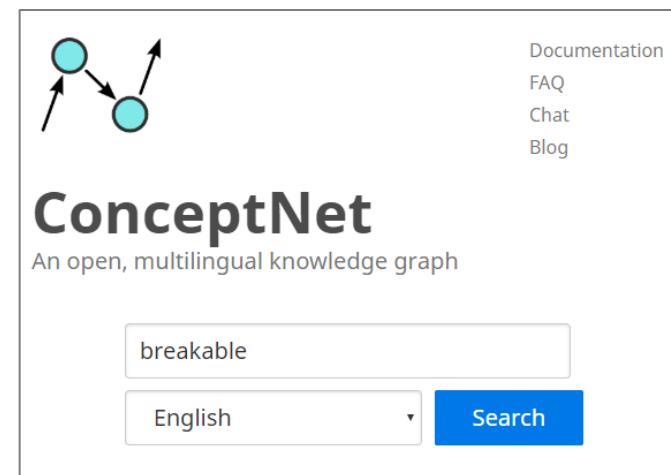
... ● ● ●

WebChild

...

Dialogue System with Commonsense Knowledge

- Use commonsense knowledge from **knowledge base**



Commonsense aware Conversation Generation
(Zhou et al., 2018)

Output Example

Post	Why are you so breakable ?
Knowledge	(glass, RelatedTo, breakable), (brittle, RelatedTo, breakable), (rule, RelatedTo, breakable)
Seq2Seq	I'm not a OOV , I'm just a OOV.
MemNet	I'm not OOV. I'm just a really nice person.
CopyNet	I'm not. I'm just a lurker.
CCM	Because I'm a brittle man .

Intent and Reaction Prediction with Commonsense Knowledge

- Use commonsense knowledge learnt from **supervised data**

PersonX cooks thanksgiving dinner

X's intent → to impress their family
X's reaction → tired, a sense of belonging
Y's reaction → impressed

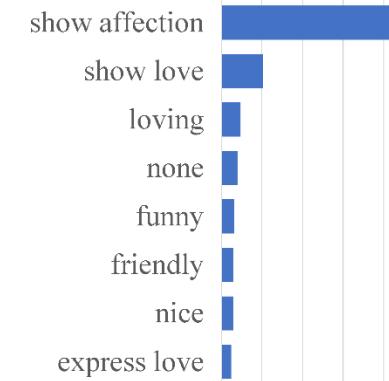
PersonX drags PersonX's feet

X's intent → to avoid doing things
X's reaction → lazy, bored
Y's reaction → frustrated, impatient

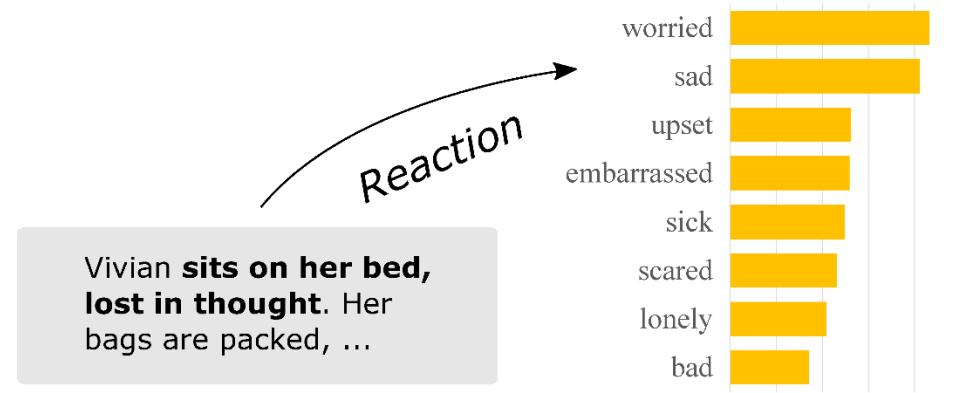
PersonX reads PersonY's diary

X's intent → to be nosey, know secrets
X's reaction → guilty, curious
Y's reaction → angry, violated, betrayed

Event2Mind
(Rashkin et al., 2018)



Juno laughs and **hugs her father, planting a smooch on his cheek.**



QA and Reasoning with Commonsense Knowledge

The service was poor, but the food was _____

- Use commonsense knowledge learnt from **unsupervised data**

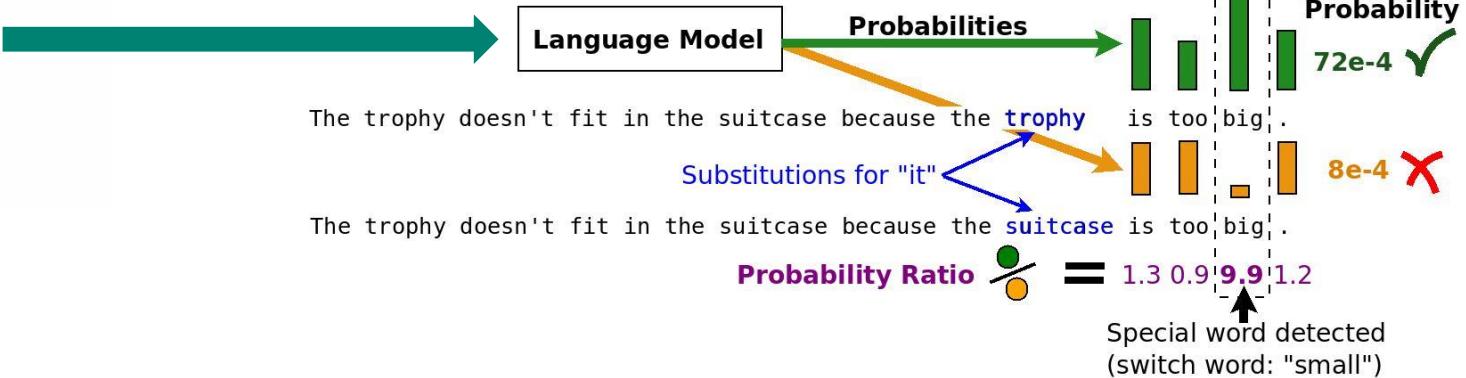


Winograd Schema Challenge

The trophy doesn't fit in the suitcase because **it** is too big.

Question: What is too big?

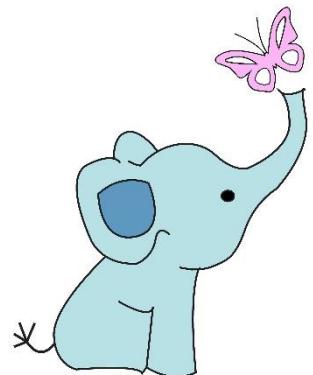
Answer: (a) **the trophy** (b) **the suitcase**



Language Model for Commonsense Reasoning (Trinh and Le, 2018)

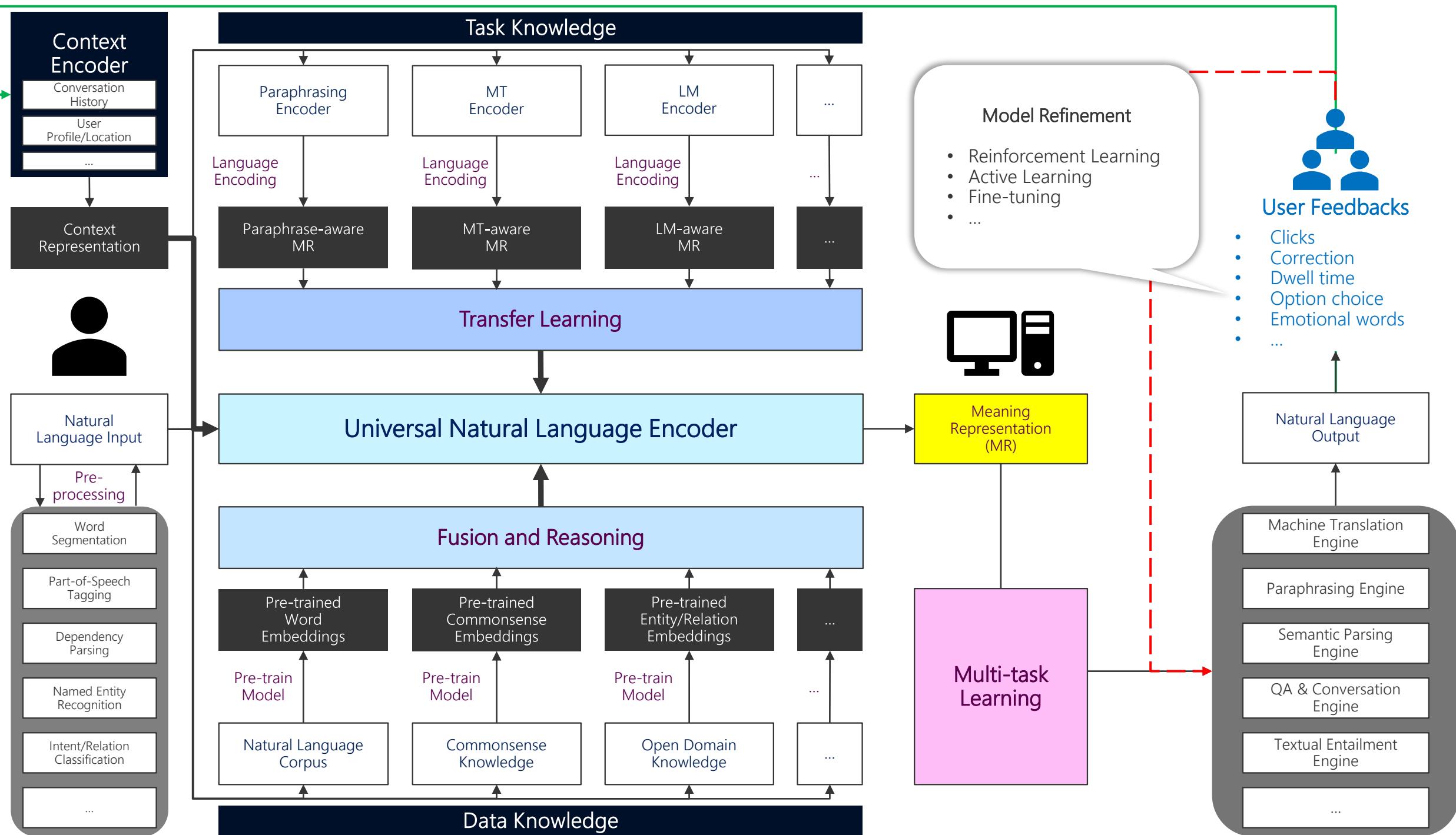
Conclusion (3)

- Commonsense is important to AI models, esp. chatbot and dialogue system
- Few datasets for model developments
- Take-ways
 - Commonsense is invisible most of time
 - But it is important to conversational AIs, such as chatbot and dialogue system



Today's Agenda

- Deep Learning Basics for NLP
- NLP with Pre-trained Embeddings
- NLP with Knowledge Bases
- NLP with Commonsense
- **Summary and Trend**



Application in Bing

who is president of united states in 2000

All Images Videos Maps News Shopping | My saves

Microsoft Show results from Microsoft ▾

18,300,000 Results Any time ▾

President in 2000 Bill Clinton

how high is yao ming's wife

All Images Videos Maps News Shopping | My saves

Microsoft Show results from Microsoft ▾

4,920,000 Results Any time ▾

Ye Li · Height
6' 3" (1.90 m)

Yao Ming 7' 6"
Yi Jianlian 6' 11"

LeBron James 6' 8"
Wang Zhizhi 7' 1"

Bing Knowledge QA

which is the highest mountain in the world

All Images Videos Maps News Shopping | My saves

Microsoft Show results from Microsoft ▾

162,000,000 Results Any time ▾

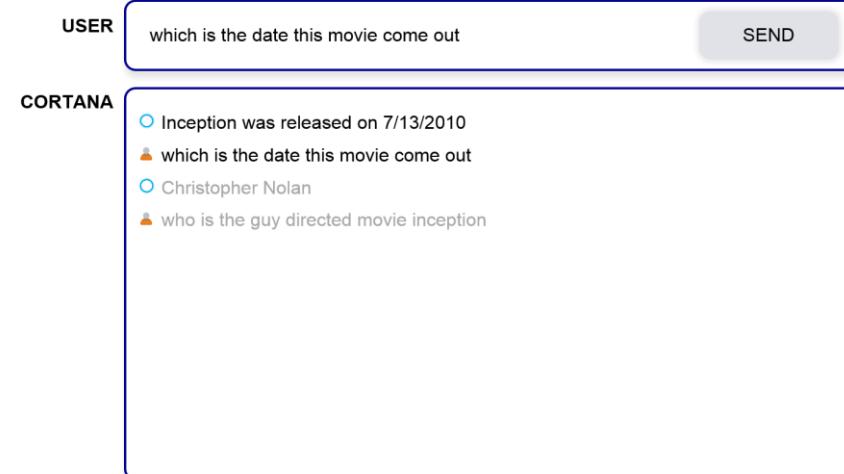
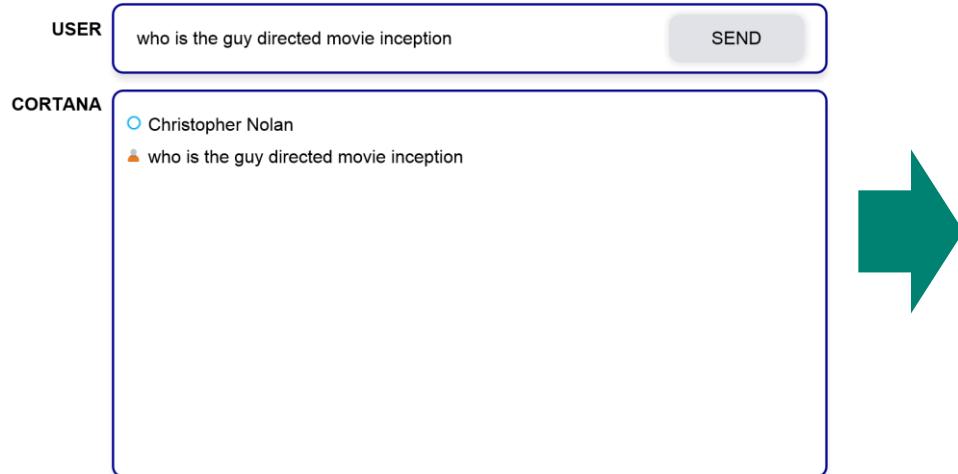
The highest mountain in the world by elevation

Mount Everest 29,029 feet

RANK	MOUNTAIN	ELEVATION (FEET)	PROMINENCE (FEET)
1	Mount Everest	29,029	29,029
2	K2	28,251	13,179
3	Kangchenjunga	28,169	12,867
4	Lhotse	27,940	2,001
5	Makalu	27,838	7,828
6	Cho Oyu	26,864	7,677
7	Dhaulagiri	26,795	11,014
8	Manaslu	26,781	10,144
9	Nanga Parbat	26,660	15,118
10	Annapurna Massif	26,545	9,790

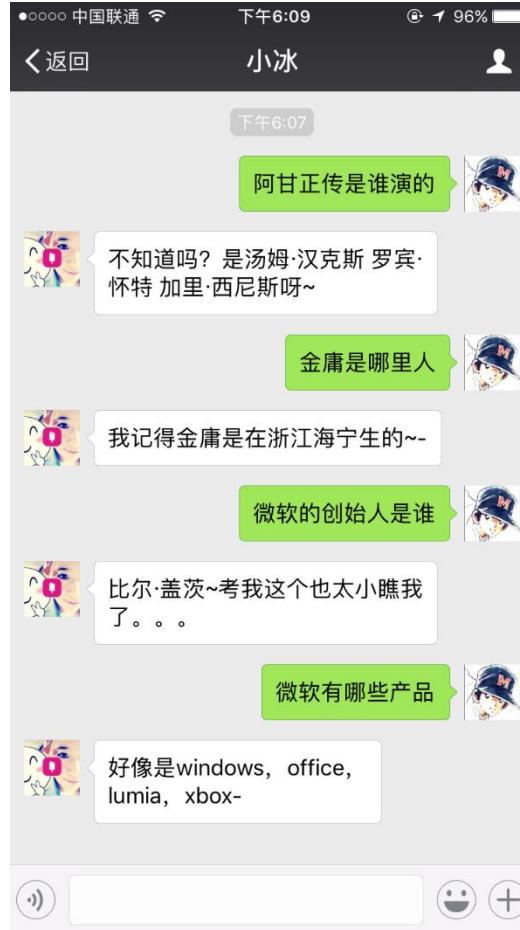
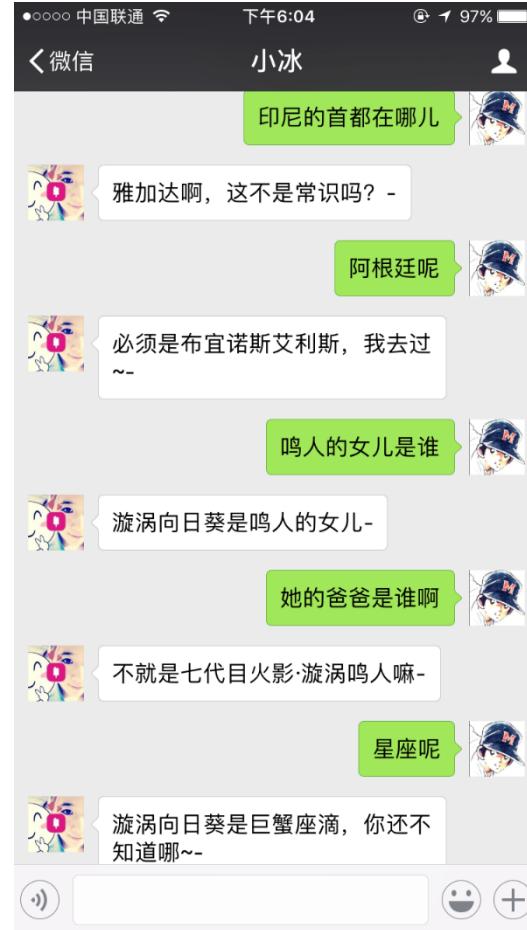
Bing Table QA

Application in Cortana



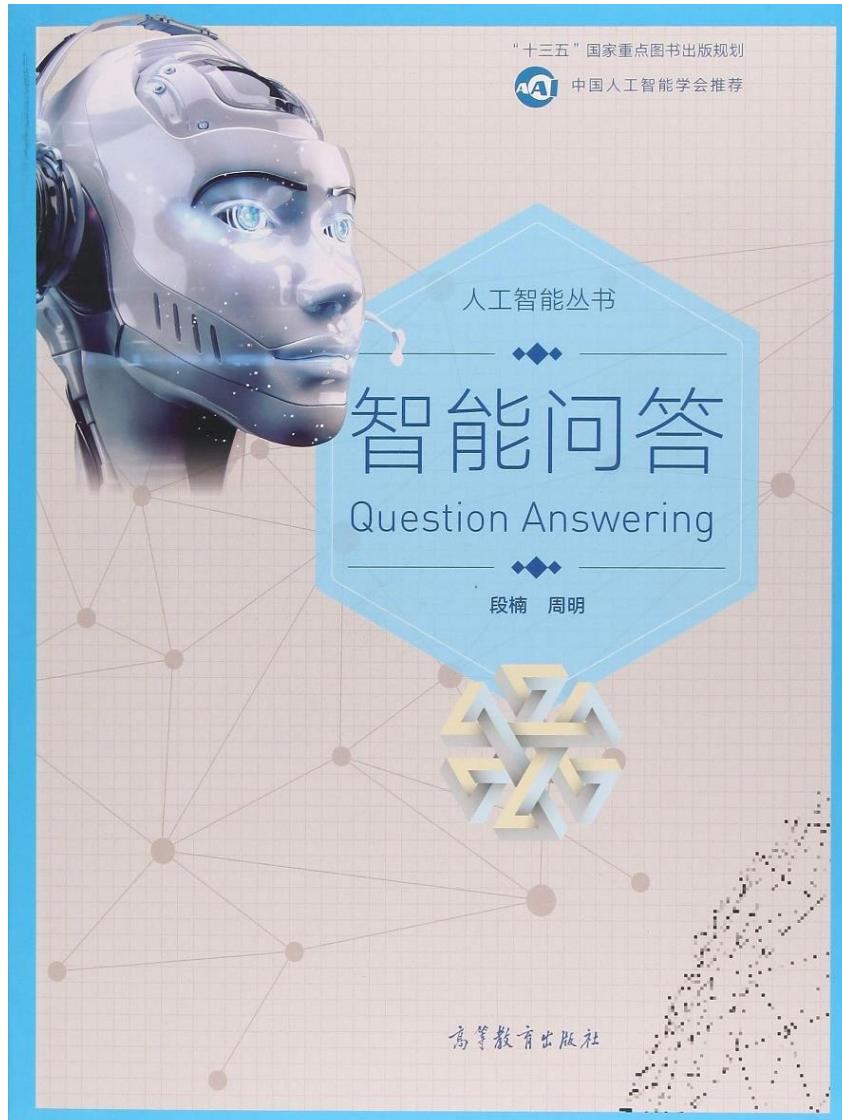
Microsoft Cortana Speaker

Application in Xiaoice



微软		公司名称		微软公司↓
微软		外文名称		Microsoft corporation.↓
微软		总部地点		美国华盛顿州雷德蒙市↓
微软		成立时间		1975年4月4日16时↓
微软		经营范围		操作系统, 办公软件, 手机↓
微软		公司性质		上市公司、外商独资↓
微软		公司口号		新效率(New Efficiency)↓
微软		年营业额		77,849百万美元 (2014年) ↓
微软		员工数		99,000人(2014年)↓
微软		联合创始人		比尔·盖茨、保罗·艾伦↓
微软		现任董事长		约翰·汤普森↓
微软		首席执行官		萨蒂亚·纳德拉↓
微软		首席运营官		凯文·特纳↓
微软		世界500强		第104位 (2014年) ↓
微软		成立地点		美国新墨西哥州阿尔伯克基市↓
微软		中国总部		中国北京海淀区知春路49号↓
微软		主要产品		xbox, windows, office, lumia↓
董明珠(珠海格力集团有限公司原董事长)		中文名		董明珠↓
董明珠(珠海格力集团有限公司原董事长)		外文名		Mingzhu Dong↓
董明珠(珠海格力集团有限公司原董事长)		别名		东方明珠↓
董明珠(珠海格力集团有限公司原董事长)		国籍		中华人民共和国↓
董明珠(珠海格力集团有限公司原董事长)		民族		汉↓
董明珠(珠海格力集团有限公司原董事长)		出生地		江苏南京↓
董明珠(珠海格力集团有限公司原董事长)		出生日期		1954年8月↓
董明珠(珠海格力集团有限公司原董事长)		职业		格力电器董事长兼总裁↓
董明珠(珠海格力集团有限公司原董事长)		毕业院校		芜湖职业技术学院↓
董明珠(珠海格力集团有限公司原董事长)		主要成就		全球100位最佳CEO↓
董明珠(珠海格力集团有限公司原董事长)		代表作品		《棋行天下》↓

Two NLP Books by MSRA



Summary and Trend

- Summary
 - Models can learn **implicit knowledge** from large-scale datasets for rich-resource tasks
 - Models need **explicit knowledge** for low-resource tasks without much training data
- Trend
 - Pre-trained embeddings
 - Commonsense & Reasoning
 - Explainable NLP with knowledge
 - Transfer learning & Multi-task learning
 - Multi-turn modeling
 - Multi-modal learning
 - Multi-lingual approach

Thank you!

We are hiring interns and FTEs ^_^