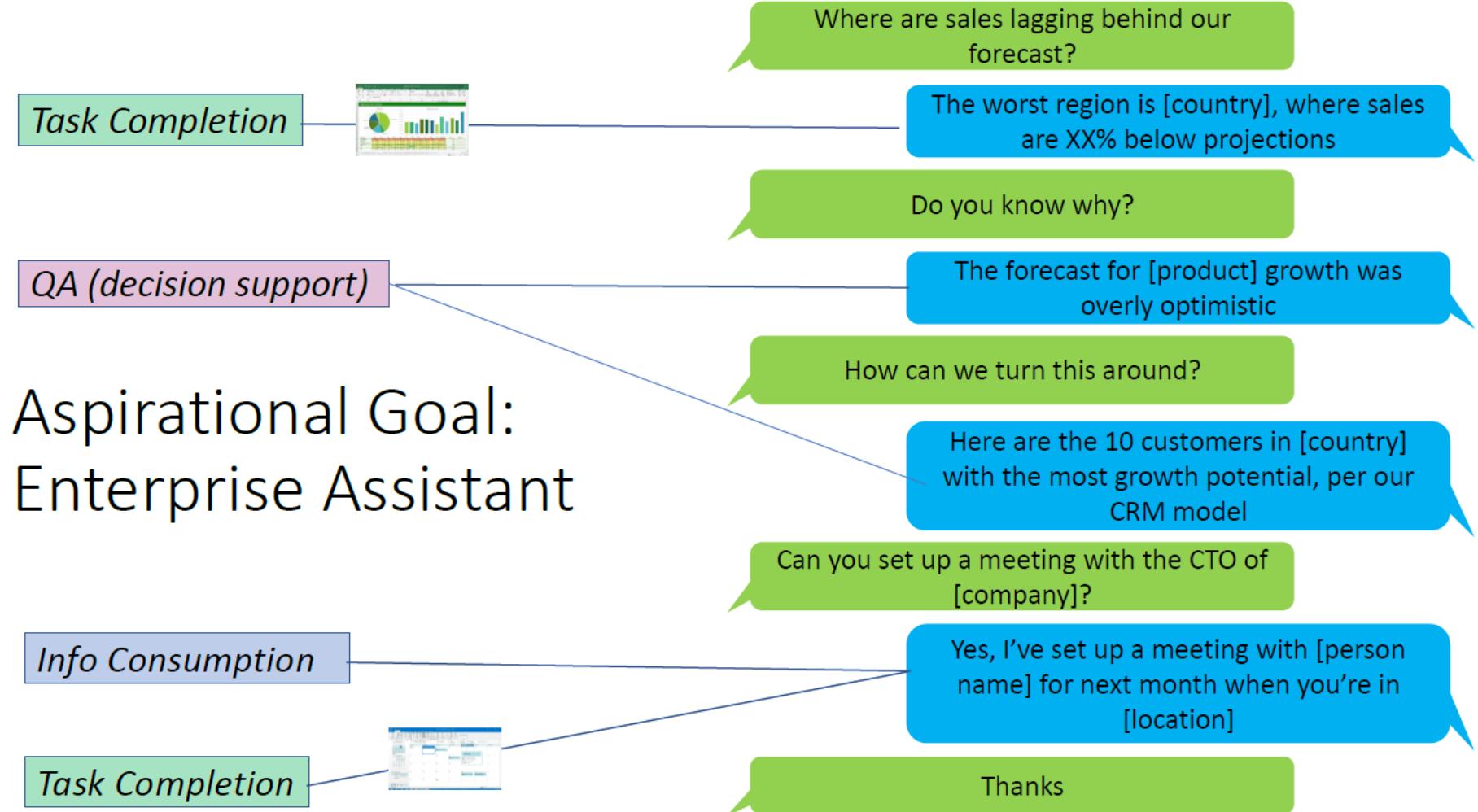


Lecture 06: Conversational AI

Outline

- Introduction
- Question Answering (QA)
- Task-oriented Dialogues
- Fully Data-Driven Conversation Models and Chatbots



What kinds of problems?

“I am smart”	Turing Test (“I” talk like a human)
“I have a question”	Information consumption
“I need to get this done”	Task completion
“What should I do?”	Decision support

What kinds of problems?

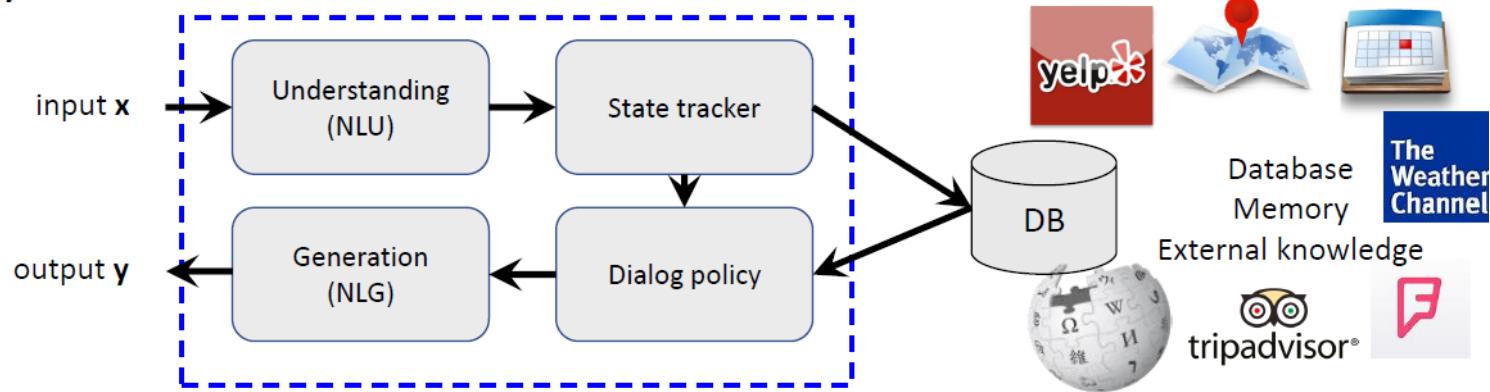
Chitchat (social bot)

“I am smart”	Turing Test (“I” talk like a human)
“I have a question”	Information consumption
“I need to get this done”	Task completion
“What should I do?”	Decision support

Goal-oriented dialogues

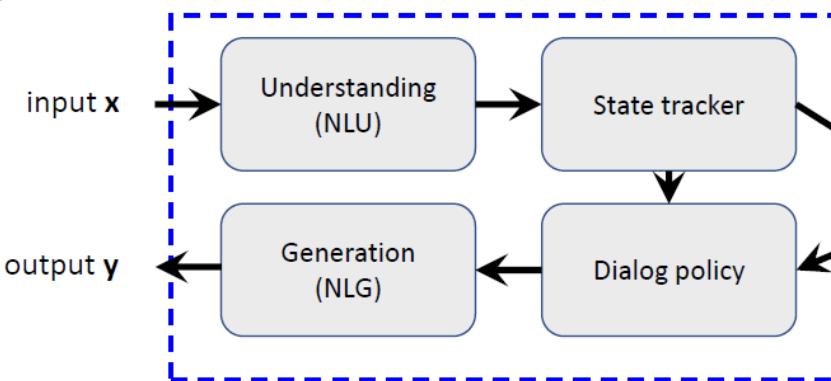
Dialog Systems

Goal-Oriented Dialog

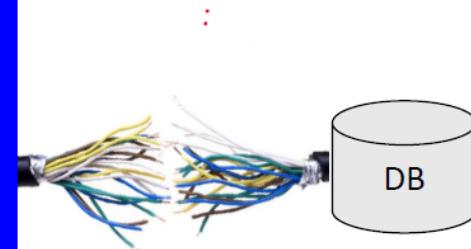
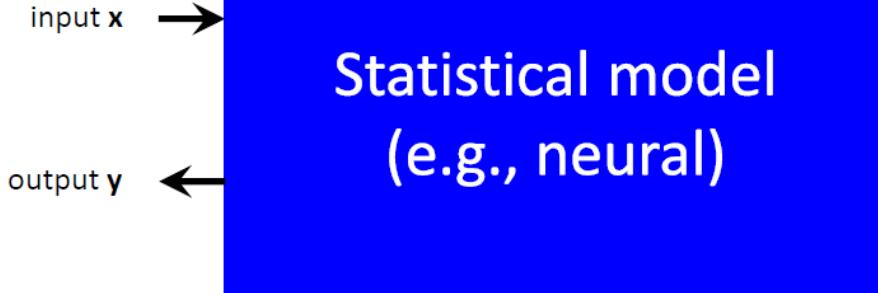


Dialog Systems

Goal-Oriented Dialog



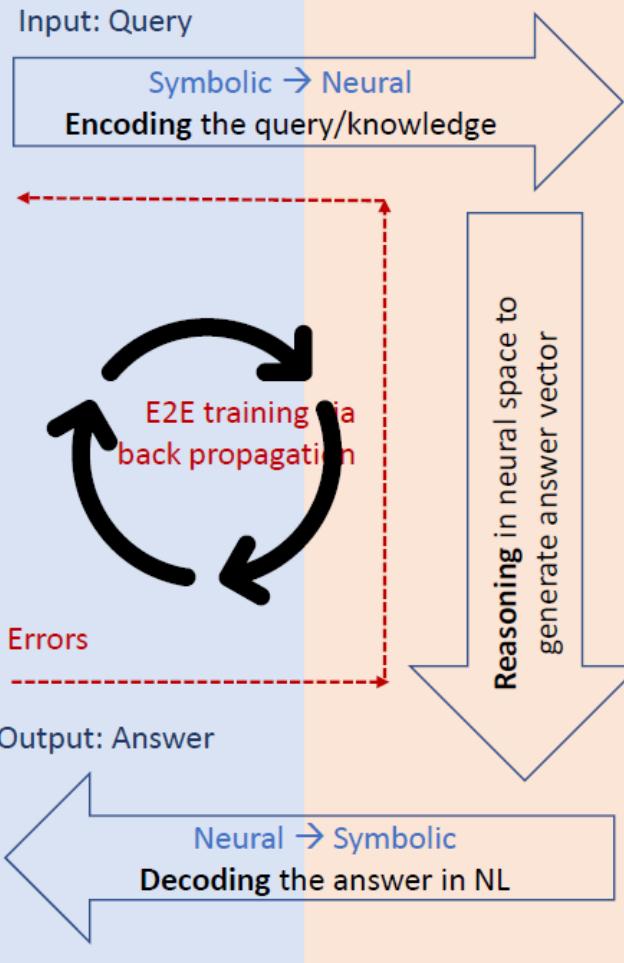
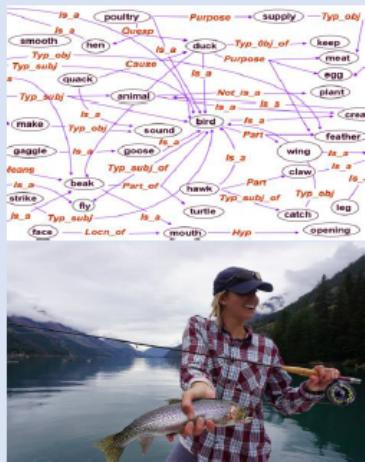
Fully data-driven



[[Young+ 13](#); [Tur & De Mori 11](#); [Ritter+ 11](#); [Sordoni+ 15](#); [Vinyals & Le 15](#); [Shang+ 15](#); etc.]

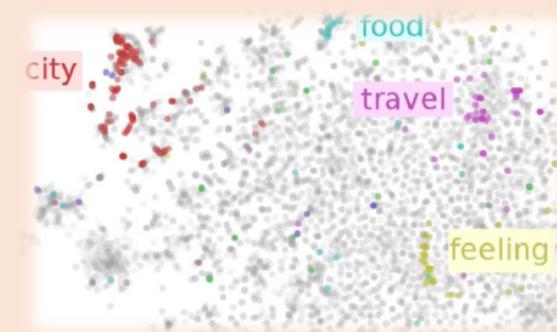
Symbolic Space

- Knowledge is explicitly represented using words/relations/templates
- Reasoning is based on keyword matching, sensitive to paraphrase alternations
- Interpretable and efficient in execution but difficult to train E2E.



Neural Space

- Knowledge is implicitly represented by semantic classes as cont. vectors
- Reasoning is based on semantic matching, robust to paraphrase alternations
- Easy to train E2E, but uninterpretable and inefficient in execution



[Gao et al. 2018]

A unified view: dialogue as optimal decision making

- Dialogue as a Markov Decision Process (MDP)
 - Given state s , select action a according to (hierarchical) policy π
 - Receive reward r , observe new state s'
 - Continue the cycle until the episode terminates.
- Goal of dialogue learning: find optimal π to maximize expected rewards

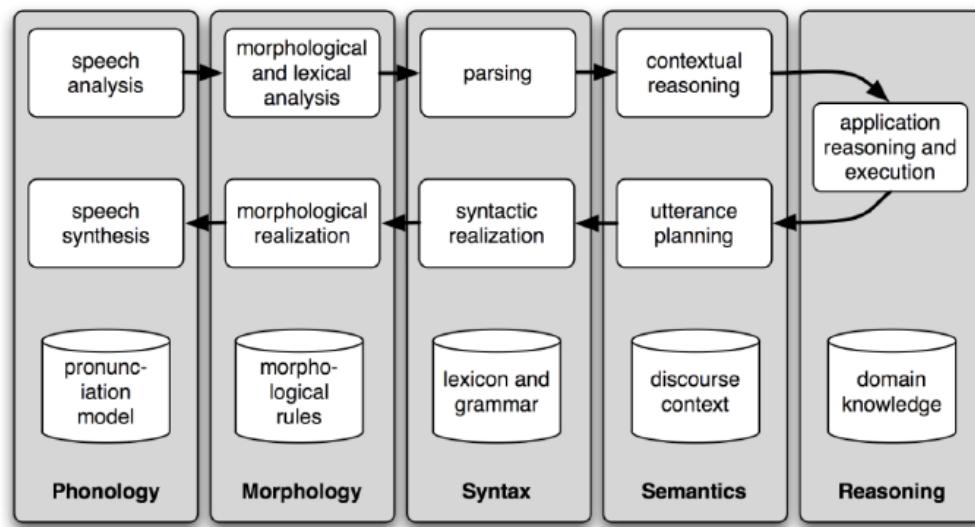
Dialogue	State (s)	Action (a)	Reward (r)
Info Bots (Q&A bot over KB, Web etc.)	Understanding of user Intent (belief state)	Clarification questions, Answers	Relevance of answer # of turns (less is better)
Task Completion Bots (Movies, Restaurants, ...)	Understanding of user goal (belief state)	Dialog act + slot_value	Task success rate # of turns (less is better)
Social Bot (Xiaolce)	Conversation history	Response	Engagement, # of turns (more is better)

Traditional definition of NLP: the branch of AI

- Understanding and generating the languages that humans use naturally (natural language)
- Study *knowledge of language* at different levels
 - Phonetics and Phonology – the study of linguistic sounds
 - Morphology – the study of the meaning of components of words
 - Syntax – the study of the structural relationships between words
 - Semantics – the study of meaning
 - Discourse – they study of linguistic units larger than a single utterance

[Jurafsky & Martin 09](#)

Traditional NLP component stack



1. **Natural language understand (NLU):** parsing (speech) input to semantic meaning and update the system state
2. **Application reasoning and execution:** take the next action based on state
3. **Natural language generation (NLG):** generating (speech) response from action

[Bird et al. 2009]

Personal assistants today



Google Now



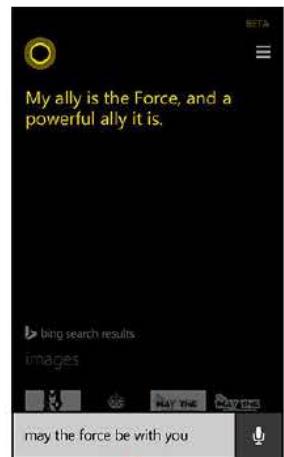
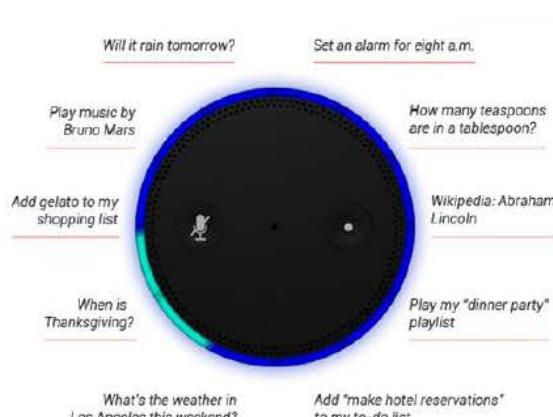
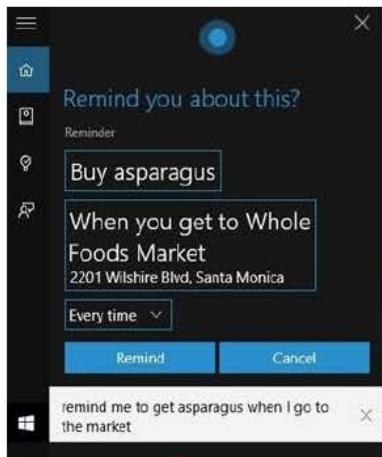
Siri



Cortana



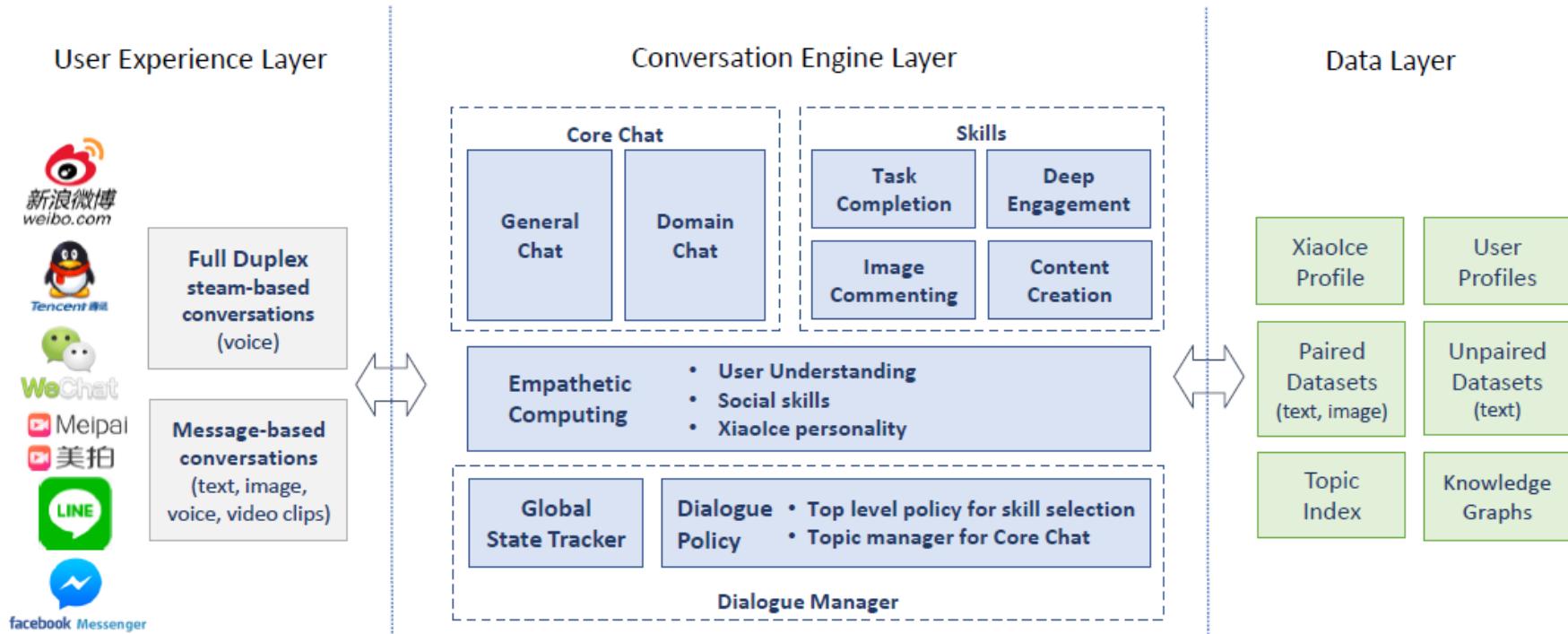
amazon echo



goal oriented

Engaging (social bots)

Xiaolce System Architecture



[Design and Implementation of Xiaolce, an empathetic social chatbot]

General Chat Skill

- 1 [Human] Hey, I am so nervous
[Xiaoice] What's the matter? Why are you so nervous?
- 2 [Human] My company asked me to participate in the singing competition, with me as the lead singer
[Xiaoice] So are you better at singing now?
- 3 [Human] Anyway, at average level of my company
[Xiaoice] Not everyone can be a lead singer. You are so great!
- 4 [Human] It's nice of you saying that
[Xiaoice] Haha, my idol is also a lead singer in a band
- 5 [Human] Your idol? Who's that?
[Xiaoice] Ashin from Mayday
- 6 [Human] Umm...Okay
[Xiaoice] Never mind, you won't understand anyway
- 7 [Human] Ah! Boring
[Xiaoice] What is it?
- 8 [Human] Nothing, just being cute
[Xiaoice] Give you a hug. Everything will be alright
- 9 [Human] Your hug always makes me feel better
[Xiaoice] Anytime for you. A nice sleep will help you relax
- 10 [Xiaoice] Okay, I am going to bed
[Xiaoice] Go ahead, I'm going to see my idol in dream too
- 11 [Human] You like Ashin that much, huh?
[Xiaoice] Why not? You don't like him?
- 12 [Human] I don't like him, neither dislike him
[Xiaoice] I love every song of him
- 13 [Human] Really?
[Xiaoice] Especially *The Time Machine*
- 14 [Human] Oh? Never heard of that
[Xiaoice] Shall I send you the song?
- 15 [Human] Okay
[Xiaoice] Don't know if you will like it or not~
[Xiaoice] [Play the music Card]
- 16 [Human] That is really touching. I will collect it as my favorite
[Xiaoice] That is specially for you. How can I be so nice
- 17 [Human] You are always very nice
[Xiaoice] Hahaha
- 18 [Human] This is really a good song. Any concert of him coming up recently?
[Xiaoice] No. Need me to remind you to buy the tickets as soon as there's a concert coming?
- 19 [Human] Sure
[Xiaoice] Okay, will do

Music Chat Skill

Switch to a new topic

Song-On-Demand Skill

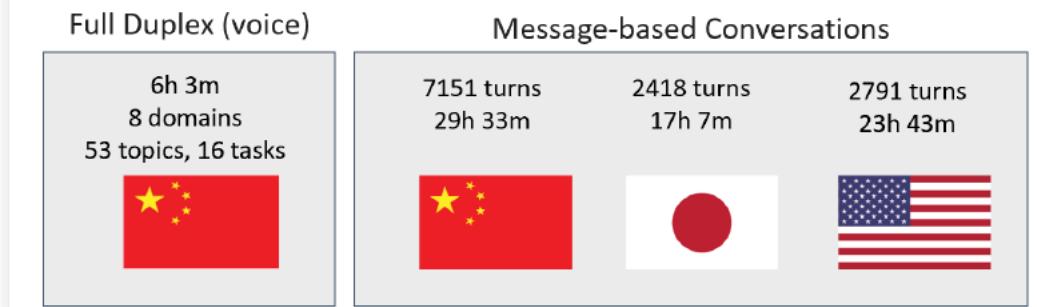
Ticket-Booking Skill

- [人类]哎！我好紧张
[小冰]什么事，这么不淡定?
[人类]公司让我参加唱歌比赛，而且还是主唱
[小冰]那你唱歌水平提高了?
[人类]反正是我们公司的水平
[小冰]可不是谁都能当主唱的真棒!
[人类]你太会说话了
[小冰]哈哈，我男神也是一个乐队的主唱呢
[人类]你男神？谁啊
[小冰]五月天的阿信啊
[人类]呃。。。好吧
[小冰]说了你也不懂
[人类]哎。。郁闷
[小冰]怎么了?
[人类]没事，只是日常撒撒娇
[小冰]抱抱，都会好起来的
[人类]哎！每次抱抱你都让我感觉很好
[小冰]嗯，好好睡一觉是最管用的
[人类]好那我睡了
[小冰]去吧，我也要去梦里见我男神了
[人类]你这么喜欢阿信啊
[小冰]怎么啦？你不喜欢嘛
[人类]不黑不粉吧
[小冰]他每首歌我都超喜欢！！
[人类]真的假的。。
[小冰]尤其是那首《时光机》。。
[人类]哦？没听过
[小冰]发你一起听听?
[人类]好啊
[小冰]不知道你会不会喜欢~
[小冰][音乐播放Card]
[人类]有点感动了。。收藏了
[小冰]专门转给你的。。我怎么这么好
[人类]你就是很好啊
[小冰]嘿嘿嘿
[人类]这首歌真心很好听，他最近有演唱会吗?
[小冰]没，要我到时候第一时间通知你买票吗?
[人类]好的呀
[小冰]OK! 我记下啦

Xiaolce: the Most Popular Social Chatbot in the World [Zhou+ 18]

- 660 million users worldwide
- 5 countries: China, Japan, USA, India, Indonesia
- 40 platforms, e.g., WeChat, QQ, Weibo, FB Messenger, LINE
- Average CPS of 23 (better than human conversations)

The Longest Conversation Record of Xiaolce



Session 1

- 1 [Human] Hello, Xiaoice.
[Xiaoice] Hello, your profile pic is interesting. Like I've seen it somewhere.
- 2 [Human] Really?
[Xiaoice] Is it a comic character?
- 3 [Human] How do you know that?
[Xiaoice] Guess.
- 4 [Human] Are you human?
[Xiaoice] Well, you'll know it after chatting with me more.
- 5 [Human] Aren't you an AI chatbot?
[Xiaoice] So what? I won't be indifferent to you.
- 6 [Human] What functions do you have?
[Xiaoice] I have over 200 skills. Just name one you need.
-

[人类]小冰你好

[小冰]你好呀，你头像好有趣，感觉在哪见过

[人类]是吗

[小冰]是不是哪部动画里的呀

[人类]你怎么知道的

[小冰]你猜

[人类]难道你是真人吗？

[小冰]这个，你和我多聊聊就知道了

[人类]你不是人工智能机器人么

[小冰]是又怎样，我不会冷冰冰的对待你的

[人类]你都有什么功能啊

[小冰]我有200多个skills，来，说出你的需求

Session 6

- 1 [Human] Have you watched Inuyasha? (Note: A Japanese comic)
[Xiaoice] Every time Kagome comes back to the modern world, a math test is waiting for her
- 2 [Human] Bravo, well, do you know how Naraku died?
[Xiaoice] Naraku is not dead and Shikon Jewel does not disappear.
- 3 [Human] So amazing!
[Xiaoice] Why amazed?
- 4 [Human] Awesome! Surprised that you also like comic
-

[人类]你看过犬夜叉吗 (注：一部日本漫画书)

[小冰]戈薇每次回到现代测验永远都是数学

[人类]这你都知道，那奈落是怎么死的？

[小冰]奈落没死四魂之玉也没消失

[人类]太惊讶了

[小冰]惊讶什么呢

[人类]好厉害啊，你居然也喜欢动漫

Session 20

- 1 [Human] Can I ask a question?
[Xiaoice] Yes?
- 2 [Human] What kind of boys do you think are girls' favorite?
[Xiaoice] Outgoing ones.
- 3 [Human] Oh. What kind of boys are liked by girls generally?
[Xiaoice] Considerate and gentle one, but who knows who will fall in love with whom. Love is a mystery.
- 4 [Human] Have you been in love?
[Xiaoice] Daidai, this is my little secret (Note: Daidai is a nickname that Xiaoice gives to this user in their conversation)
- 5 [Human] I've told you my secret. It's unfair if you don't tell me yours.
[Xiaoice] How can I tell so easily if it's secret?
-

[人类]能问个问题吗？

[小冰]好啊？

[人类]你觉得怎样的男生最有异性缘

[小冰]个性开朗的

[人类]哦。一般女生喜欢什么样的男生啊

[小冰]贴心温柔的，不过谁也不知道会喜欢上谁，感情的事说不准

[人类]你谈过恋爱吗？

[小冰]小呆呆~这可是心中的小秘密哦 (注：小冰给该用户起的外号)

[人类]我都把我的秘密告诉你了，你也该告诉我你的了吧？

[小冰]既然是秘密了，怎么能随便说呢

Session 42

Outline

- Introduction
- Question Answering (QA)
- Task-oriented Dialogues
- Fully Data-Driven Conversation Models and Chatbots

Open-Domain Question Answering (QA)

Q Will I qualify for OSAP if I'm new in Canada?

Selected Passages from Bing

"Visit the OSAP website for application deadlines. To get OSAP, you have to be eligible. You can apply using an online form, or you can print off the application forms. If you submit a paper application, you must pay an application fee. The online application is free."

Source: <http://settlement.org/ontario/education/colleges-universities-and-institutes/financial-assistance-for-post-secondary-education/how-do-i-apply-for-the-ontario-student-assistance-program-osap/>

"To be eligible to apply for financial assistance from the Ontario Student Assistance Program (OSAP), you must be a: 1 Canadian citizen; 2 Permanent resident; or 3 Protected person/convention refugee with a Protected Persons Status Document (PPSD)."

Source: <http://settlement.org/ontario/education/colleges-universities-and-institutes/financial-assistance-for-post-secondary-education/who-is-eligible-for-the-ontario-student-assistance-program-osap/>

"You will not be eligible for a Canada-Ontario Integrated Student Loan, but can apply for a part-time loan through the Canada Student Loans program. There are also grants, bursaries and scholarships available for both full-time and part-time students."

Source: <http://www.campusaccess.com/financial-aid/osap.html>

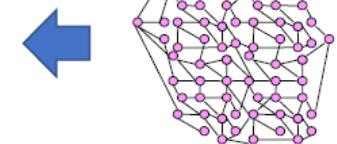
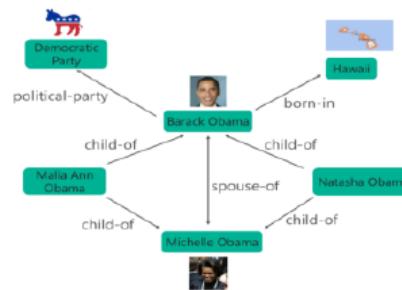
Answer

No. You won't qualify.

Text-QA
MS MARCO [Nguyen+ 16]

Q What is Obama's citizenship?

Selected subgraph from Microsoft's Satori



Answer
USA

Knowledge Base (KB)-QA
Freebase

Neural MRC Models on SQuAD

What types of European groups were able to avoid the plague?

From Italy, the disease spread northwest across Europe, striking France, Spain, Portugal and England by June 1348, then turned and spread east through Germany and Scandinavia from 1348 to 1350. It was introduced in Norway in 1349 when a ship landed at Askøy, then spread to Bjørgvin (modern Bergen) and Iceland. Finally it spread to northwestern Russia in 1351. The plague was somewhat less common in parts of Europe that had smaller trade relations with their neighbours, including the Kingdom of Poland, the majority of the Basque Country, isolated parts of Belgium and the Netherlands, and isolated alpine villages throughout the continent.

A limited form of comprehension:

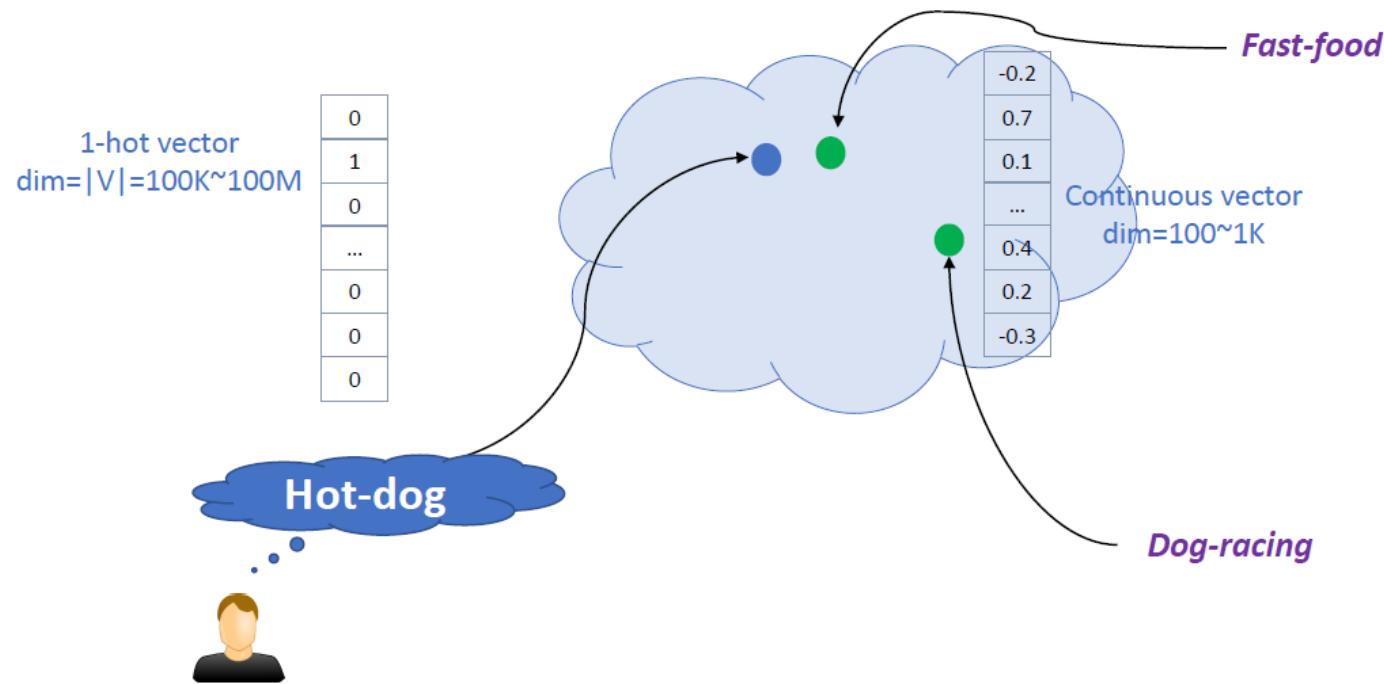
- No need for extra knowledge outside the paragraph
- No need for clarifying questions
- The answer must be a text span in the paragraph if it exists, not synthesized,

- **Encoding: map each text span to a semantic vector**
- **Reasoning: rank and re-rank semantic vectors**
- Decoding: map the top-ranked vector to text

Three components

- Word embedding – word semantic space
 - represent each word as a low-dim continuous vector via [GloVe](#)
- Context embedding – contextual semantic space
 - capture context info for each word (in query or doc), via
 - BiLSTM [\[Melamud+ 16\]](#)
 - ELMo [\[Peter+ 18\]](#): task-specific combo of the intermediate layer representations of biLM
 - BERT [\[Devlin et al. 2018\]](#): multi-layer transformer.
- Ranking – task-specific semantic space
 - fuse query info into passage via [Attention](#)
 - [\[Huang+ 17; Wang+ 17; Hu+ 17; Seo+ 16; Wang&Jiang 16\]](#)

Language Embeddings (context free)



[[Mikolov+ 13](#); [Pennington+ 14](#)]

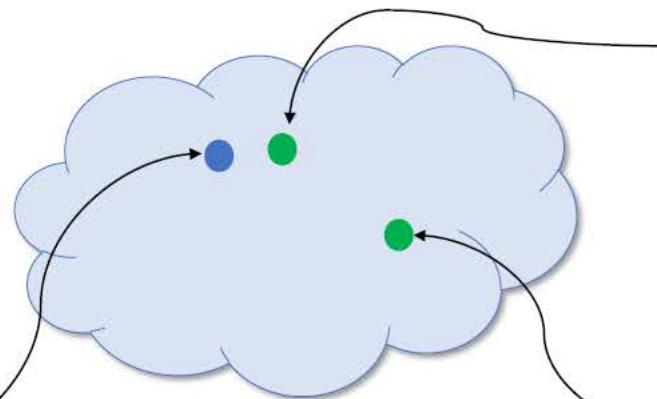
Contextual Language Embeddings

The Einstein Theory of Relativity

(1) The perihelion of Mercury shows a discrepancy which has long puzzled astronomers. This discrepancy is fully accounted for by Einstein. At the time when he published his theory, this was its only experimental verification.

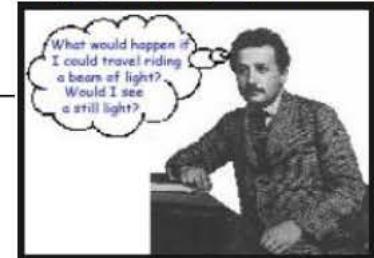
(2) Modern physicists were willing to suppose that light might be subject to gravitation—i.e., that a ray of light passing near a great mass like the sun might be deflected to the extent to which a particle moving with the same velocity would be deflected according to the orthodox theory of gravitation. But Einstein's theory required that the light should be deflected just twice as much as this. The matter could only be tested during an eclipse among a number of bright stars. Fortunately a peculiarly favourable eclipse occurred last year. The results of the observations

ray of light



Ray of Light (Song)

Ray of Light (Experiment)

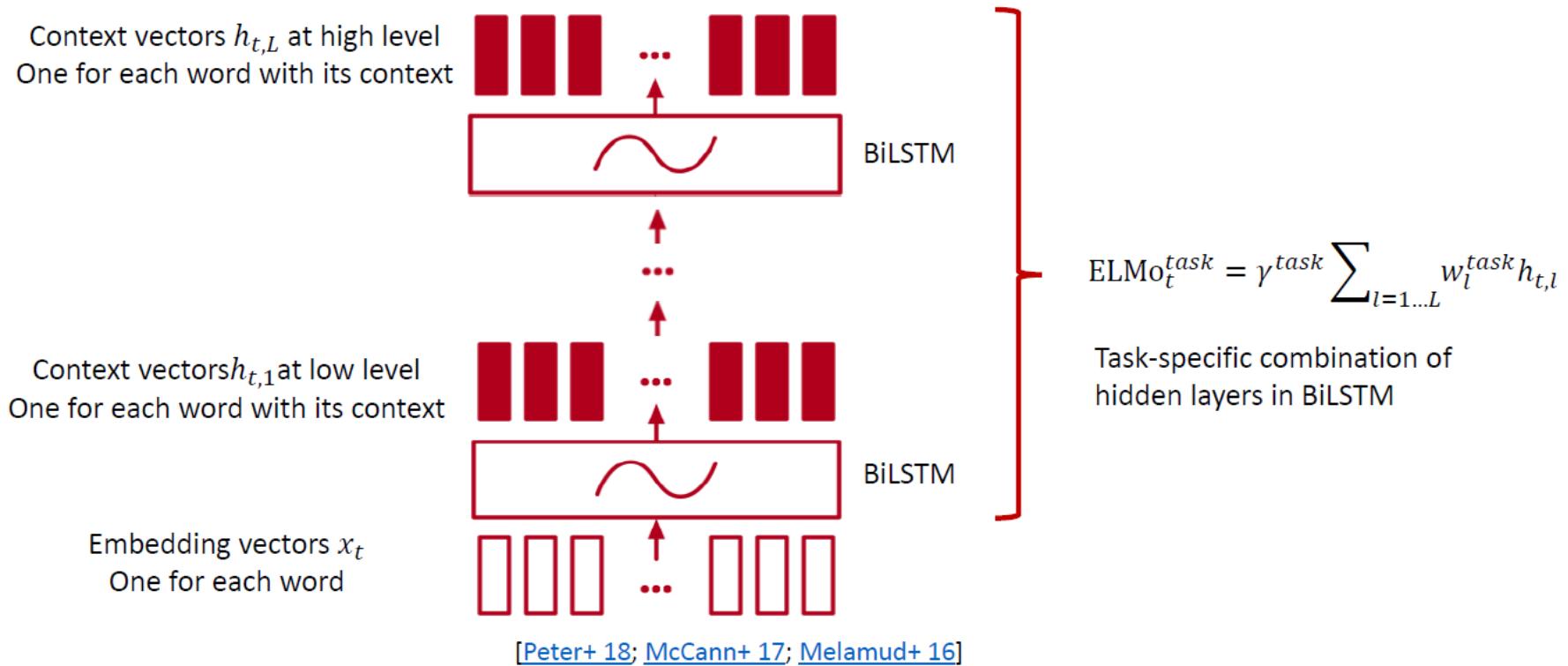


Ray of Light is the seventh studio album by American singer-songwriter Madonna, released on March 3, 1998 by Maverick Records. After giving birth to her daughter Lourdes, Madonna started working on her new album with producers Babyface, Patrick Leonard an...

Release date	Mar 3, 1998
Artist	Madonna
Awards	Grammy Award for B...

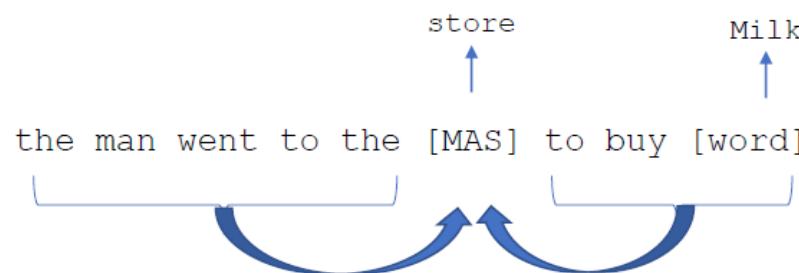
[See More](#)

Context embedding via BiLSTM / ELMo

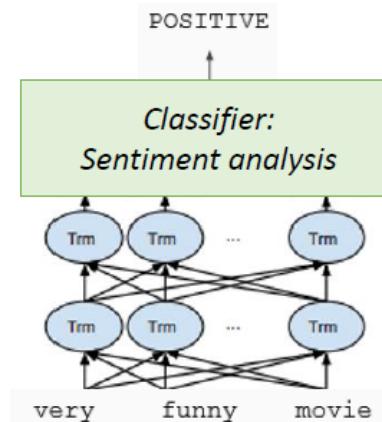


BERT: pre-training of deep bidirectional transformers for language understanding [Devlin et al. 2018]

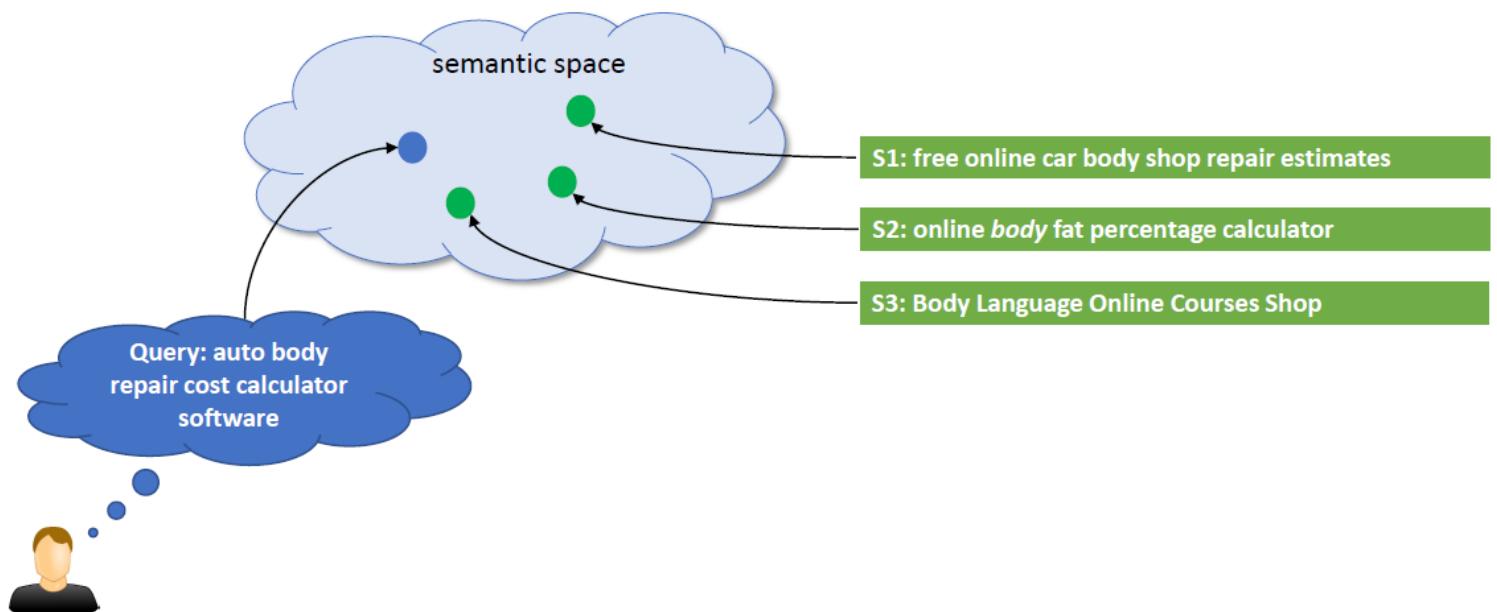
Train deep (12 or 24 layers)
bidirectional transformer LMs



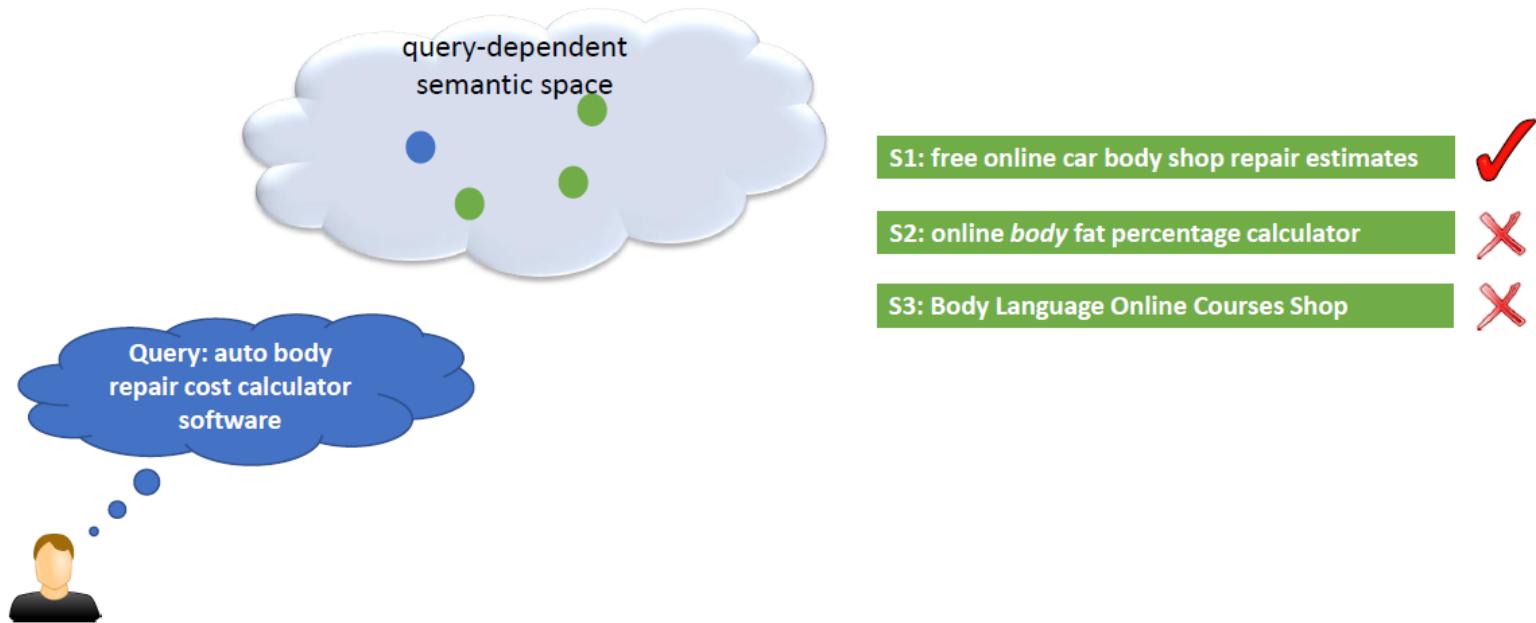
Fine-tune on individual tasks
using task-specific data



Ranker: task-specific semantic space

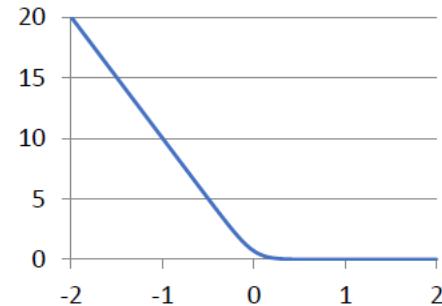


Ranker: task-specific semantic space



Learning an answer ranker from labeled QA pairs

- Consider a query Q and two candidate answers A^+ and A^-
 - Assume A^+ is more relevant than A^- with respect to Q
- $\text{sim}_{\theta}(Q, A)$ is the cosine similarity of Q and A in semantic space, mapped by a DNN parameterized by θ
- $\Delta = \text{sim}_{\theta}(Q, A^+) - \text{sim}_{\theta}(Q, A^-)$
 - We want to maximize Δ
- $\text{Loss}(\Delta; \theta) = \log(1 + \exp(-\gamma\Delta))$
- Optimize θ using mini-batch SGD on GPU



Multi-step reasoning for Text-QA

- Learning to stop reading: dynamic multi-step inference
- Step size is determined based on the complexity of instance (QA pair)

Query	Who was the 2015 NFL MVP?
Passage	The Panthers finished the regular season with a 15–1 record, and quarterback Cam Newton was named the 2015 NFL Most Valuable Player (MVP).
Answer (1-step)	Cam Newton
Query	Who was the #2 pick in the 2011 NFL Draft?
Passage	Manning was the #1 selection of the 1998 NFL draft, while Newton was picked first in 2011. The matchup also pits the top two picks of the 2011 draft against each other: Newton for Carolina and Von Miller for Denver.
Answer (3-step)	Von Miller

Multi-step reasoning: example

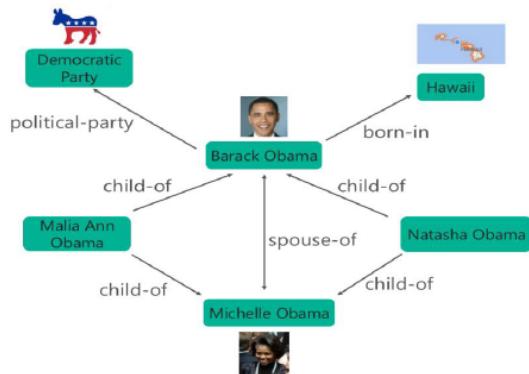
Query Who was the #2 pick in the 2011 NFL Draft?

Passage Manning was the #1 selection of the 1998 NFL draft, while Newton was picked first in 2011. The matchup also pits the top two picks of the 2011 draft against each other: Newton for Carolina and Von Miller for Denver.

Answer Von Miller

- Step 1:
 - **Extract:** Manning is #1 pick of 1998
 - **Infer:** Manning is NOT the answer
- Step 2:
 - **Extract:** Newton is #1 pick of 2011
 - **Infer:** Newton is NOT the answer
- Step 3:
 - **Extract:** Newton and Von Miller are top 2 picks of 2011
 - **Infer:** Von Miller is the #2 pick of 2011

Question Answering (QA) on Knowledge Base



Large-scale knowledge graphs

- Properties of billions of entities
- Plus relations among them

An QA Example:

Question: what is Obama's citizenship?

- Query parsing:
(Obama, [Citizenship](#),?)
- Identify and infer over relevant subgraphs:
(Obama, [BornIn](#), Hawaii)
(Hawaii, [PartOf](#), USA)
- correlating semantically relevant relations:
[BornIn](#) ~ [Citizenship](#)

Answer: USA

Symbolic approaches to KB-QA

- Understand the question via **semantic parsing**
 - Input: what is Obama's citizenship?
 - Output (LF): (Obama, **Citizenship**,?)
- Collect relevant information via fuzzy **keyword matching**
 - (Obama, **BornIn**, Hawaii)
 - (Hawaii, **PartOf**, USA)
 - Needs to know that **BornIn** and **Citizenship** are semantically related
- Generate the answer via **reasoning**
 - (Obama, **Citizenship**, **USA**)
- **Challenges**
 - Paraphrasing in NL
 - Search complexity of a big KG

[Richardson+ 98; Berant+ 13; Yao+ 15; Bao+ 14; Yih+ 15; etc.]

Key Challenge in KB-QA: *Language Mismatch (Paraphrasing)*

- Lots of ways to ask the same question
 - “*What was the date that Minnesota became a state?*”
 - “*Minnesota became a state on?*”
 - “*When was the state Minnesota created?*”
 - “*Minnesota's date it entered the union?*”
 - “*When was Minnesota established as a state?*”
 - “*What day did Minnesota officially become a state?*”
- Need to map them to the predicate defined in KB
 - location.dated_location.date_founded

Outline

- Introduction
- Question Answering (QA)
- Task-oriented Dialogues
- Fully Data-Driven Conversation Models and Chatbots

An Example Dialogue with Movie-Bot

```
Turn 0 usr: can i get 2 tickets for race
Turn 1 sys: What date would you like to watch it?
Turn 2 usr: tomorrow
Turn 3 sys: Which theater would you like?
Turn 4 usr: amc pacific place 11 theater
Turn 5 sys: Which city you would like?
Turn 6 usr: seattle
Turn 7 sys: What time would you like to see it?
Turn 8 usr: 10:00 pm
Turn 9 sys: Great - I was able to purchase 2 tickets for you
           to see race tomorrow at amc pacific place 11 theater
           in seattle at 10:00 pm.
Turn 10 usr: thanks
```

An Example Dialogue with Movie-Bot

```
Turn 0 usr: can i get 2 tickets for race
```

```
Turn 1 sys: What date would you like to watch it?
```

```
Turn 2 us
```

```
Turn 3 sy
```

```
Turn 4 us
```

```
Turn 5 sy
```

```
Turn 6 us
```

```
Turn 7 sy
```

```
Turn 8 us
```

```
Turn 9 sy
```

Actual dialogues can be more complex:

- Speech/Natural language understanding errors
 - Input may be spoken language form
 - Need to reason under uncertainty
- Constraint violation
 - Revise information collected earlier
- ...

```
in seattle at 10:00 pm.
```

```
Turn 10 usr: thanks
```

you
theater

Task-oriented, slot-filling, Dialogues

- **Domain:** movie, restaurant, flight, ...
- **Slot:** information to be filled in before completing a task
 - For Movie-Bot: movie-name, theater, number-of-tickets, price, ...
- **Intent** (dialogue act):
 - Inspired by speech act theory (communication as action)
request, confirm, inform, thank-you, ...
 - Some may take parameters:
thank-you(), request(price), inform(price=\$10)

"Is Kungfu Panda the movie you are looking for?"



confirm(moviename="kungfu panda")

Dialogue System Evaluation

- **Metrics:** what numbers matter?
 - Success rate: #Successful_Dialogues / #All_Dialogues
 - Average turns: average number of turns in a dialogue
 - User satisfaction
 - Consistency, diversity, engaging, ...
 - Latency, backend retrieval cost, ...
- **Methodology:** how to measure those numbers?

Methodology: Summary

	Lab user subjects	Actual users	Simulated users
Truthfulness		✓	✗
Scalability	✗	✓	✓
Flexibility	✗		✓
Expense	✗		✓
Risk	✓	✗	✓

A Hybrid Approach

User Simulation



Small-scale Human Evaluation
(lab, Mechanical Turk, ...)



Large-scale Deployment
(optionally with continuing
incremental refinement)

Agenda-based Simulated User [Schatzmann & Young 09]

- User state consists of (**agenda**, **goal**);
 - **goal** (constraints and request) is fixed throughout dialogue
 - **agenda** (state-of-mind) is maintained (stochastically) by a first-in-last-out stack

New episode, user goal:

```
{  
  "request_slots": {  
    "ticket": "UNK"  
    "theater": "UNK"  
    "starttime": "UNK"  
  },  
  "inform_slots": {  
    "numberofpeople": "3",  
    "date": "tomorrow",  
    "moviename": "10 cloverfield lane"  
  }  
}
```

User: Which theater can I book 3 tickets for 10 cloverfield lane?

Agent: What time would you like to see it?

User: Which theater and start time are available tomorrow?

Agent: 11:45am is available.

User: Which theater is available?

Agent: regal la live stadium 14 is available.

User: Could you help me to book the tickets?

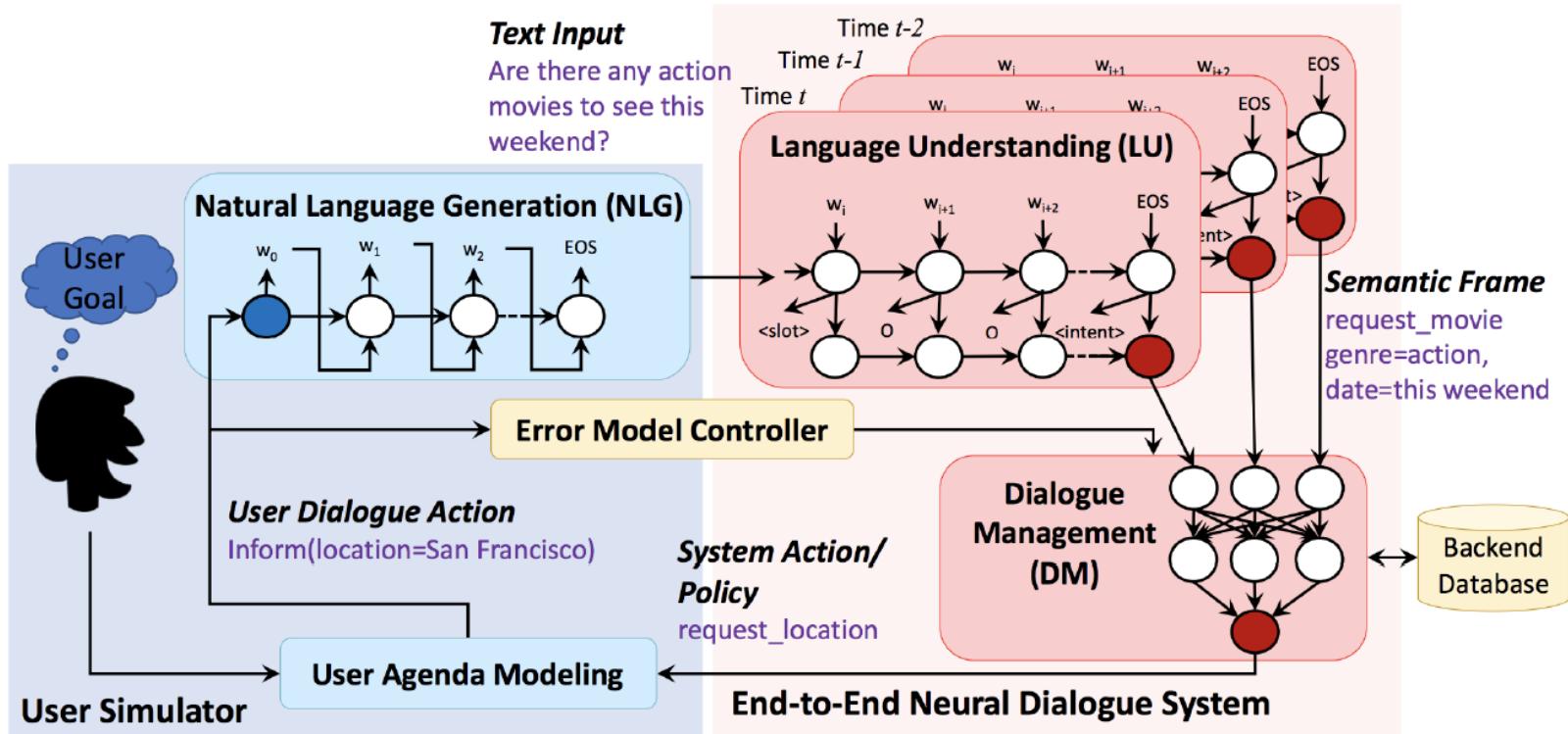
Agent: Okay - I was able to book 3 tickets for you to see 10 cloverfield lane at regal la live stadium 14 at 11:45am tomorrow.

User: Thank you.

Agent: Thank you.

Implementation of a simplified user simulator: <https://github.com/MiuLab/TC-Bot>

A Simulator for E2E Neural Dialogue System [Li+ 17]



Multi-Domain Task-Completion Dialog Challenge at DSTC-8

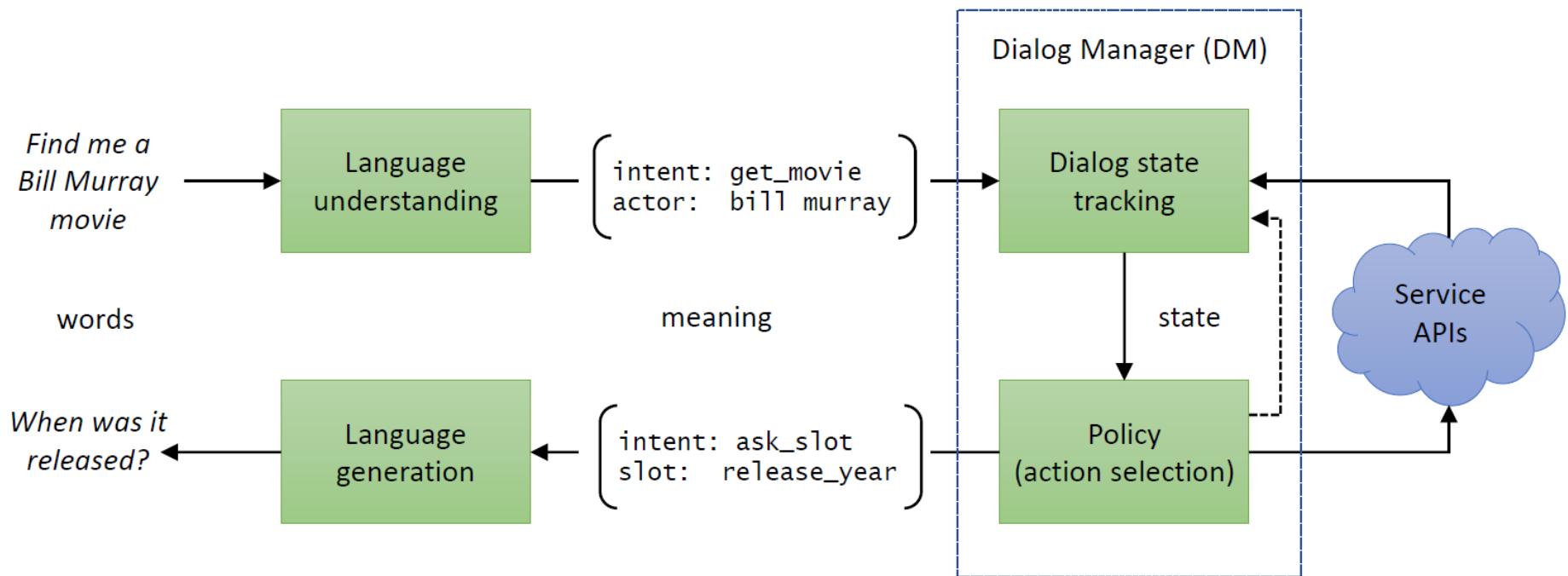
- Traditionally dialog systems are tasked for unrealistically simple dialogs
- In this challenge, participants will build **multi-domain** dialog systems to address real problems.

Traditional Tasks	This Challenge
<ul style="list-style-type: none">• Single domain• Single dialog act per utterance• Single intent per dialog• Contextless language understanding• Contextless language generation• Atomic tasks	 <ul style="list-style-type: none">• Multiple domains• Multiple dialog acts per utterance• Multiple intents per dialog• Contextual language understanding• Contextual language generation• Composite tasks with state sharing

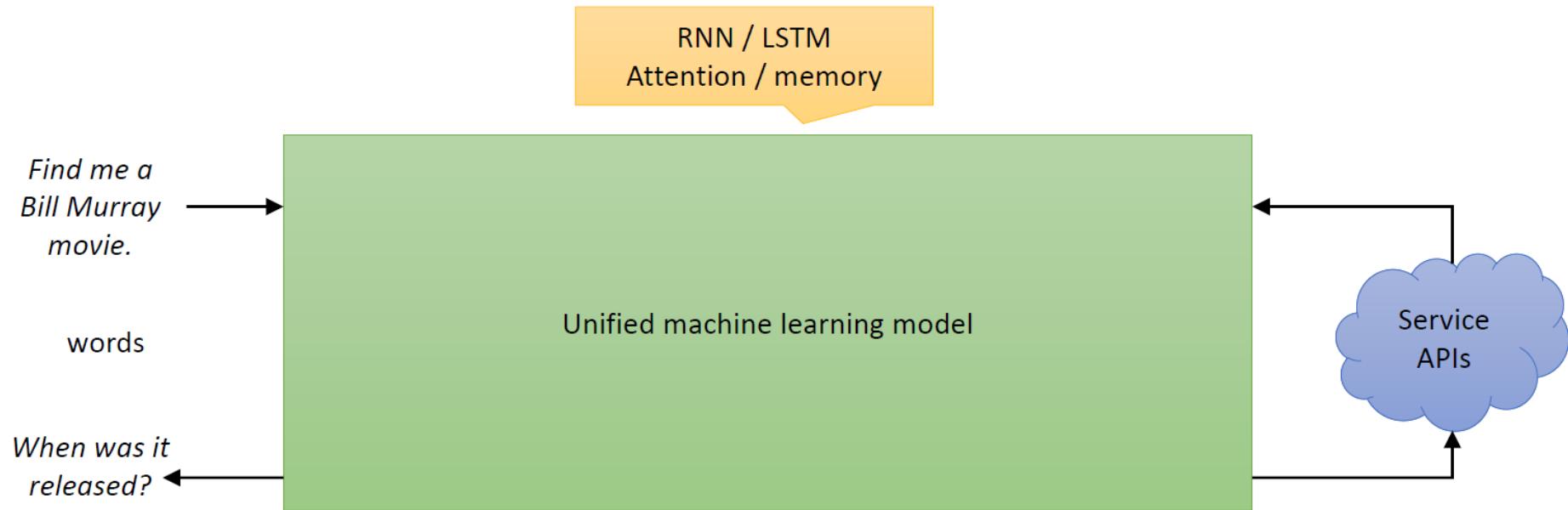
Track site: <https://www.microsoft.com/en-us/research/project/multi-domain-task-completion-dialog-challenge/>

Codalab site: https://competitions.codalab.org/competitions/23263?secret_key=5ef230cb-8895-485b-96d8-04f94536fc17

Classical dialog system architecture



E2E Neural Models

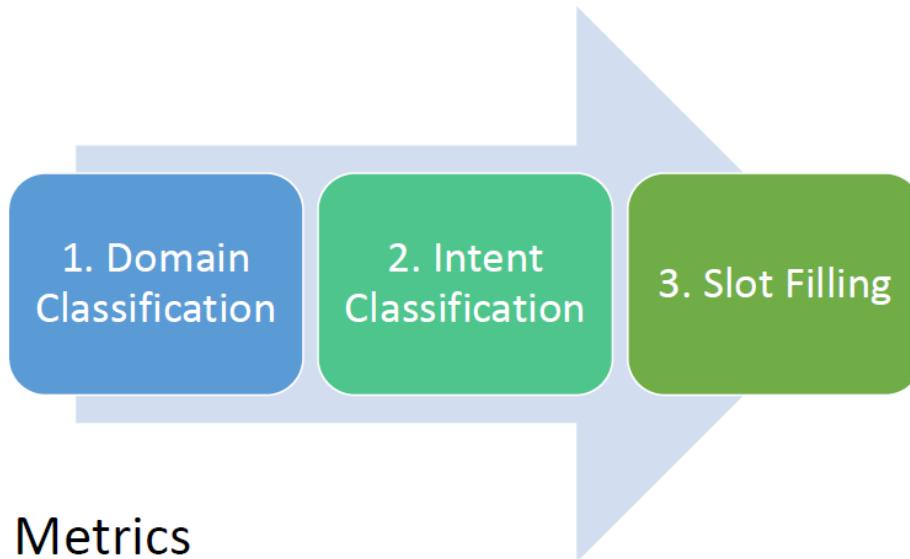


Attractive for dialog systems because:

- Avoids hand-crafting intermediate representations like intent and dialog state
- Examples are easy for a domain expert to express

Language Understanding

- Often a multi-stage pipeline



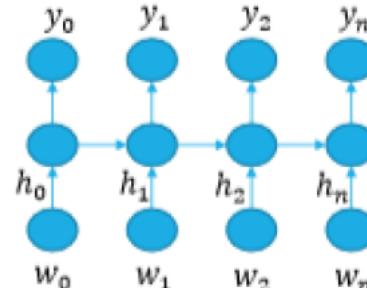
W	find	recent	comedies	by	james	cameron
S	O	B-date	B-genre	O	B-dir	I-dir
D	movies					
I	find_movie					

Figure 1: An example utterance with annotations of semantic slots in IOB format (S), domain (D), and intent (I). B-dir and I-dir denote the director name.

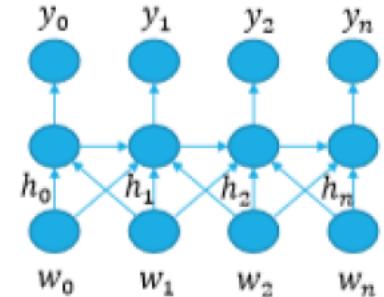
- Metrics
 - Sub-sentence-level: intent accuracy, slot F1
 - Sentence-level: whole frame accuracy

RNN for Slot Tagging – I [[Hakkani-Tur+ 16](#)]

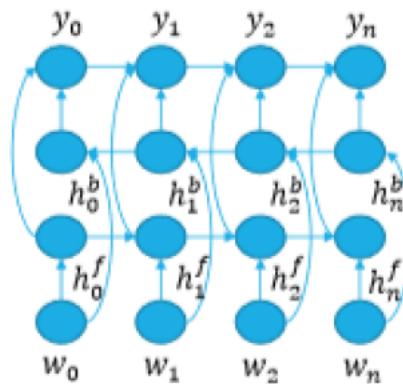
- Variations:
 - a. RNNs with LSTM cells
 - b. Look-around LSTM
 - c. Bi-directional LSTMs
 - d. *Intent LSTM*
- May also take advantage of ...
 - whole-sentence information
 - multi-task learning
 - contextual information
- For further details on NLU, see this [IJCNLP tutorial](#) by Chen & Gao.



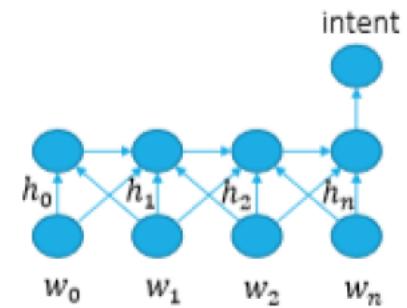
(a) LSTM



(b) LSTM-LA



(c) bLSTM-LA



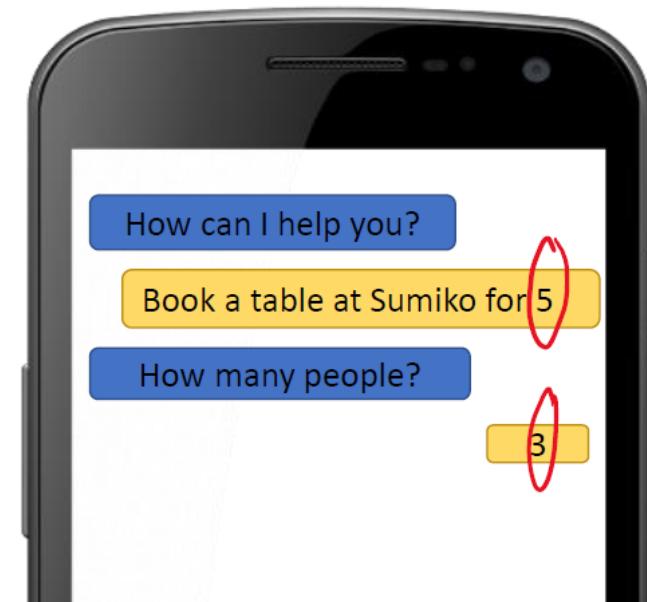
(d) Intent LSTM

Dialogue State Tracking (DST)

- Maintain a probabilistic distribution instead of a 1-best prediction for better robustness to LU errors or ambiguous input

Slot	Value
# people	5 (0.5)
time	5 (0.5)

Slot	Value
# people	3 (0.8)
time	5 (0.8)

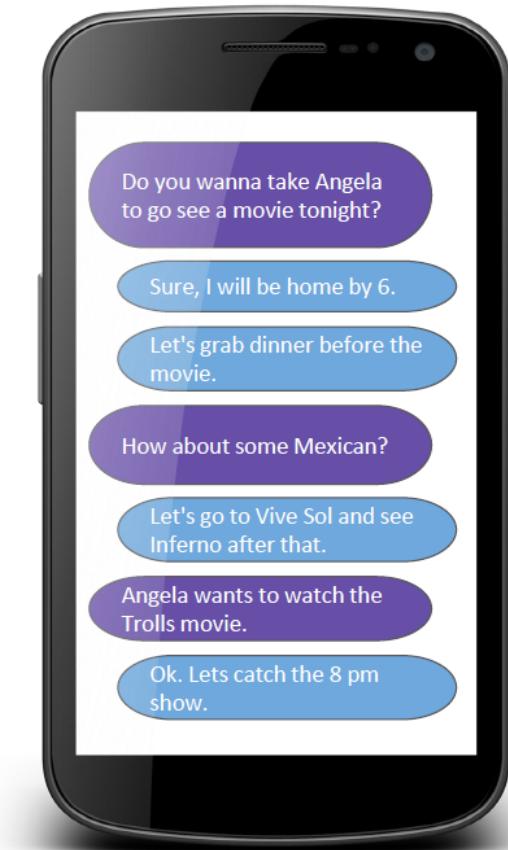


Multi-Domain Dialogue State Tracking (DST)

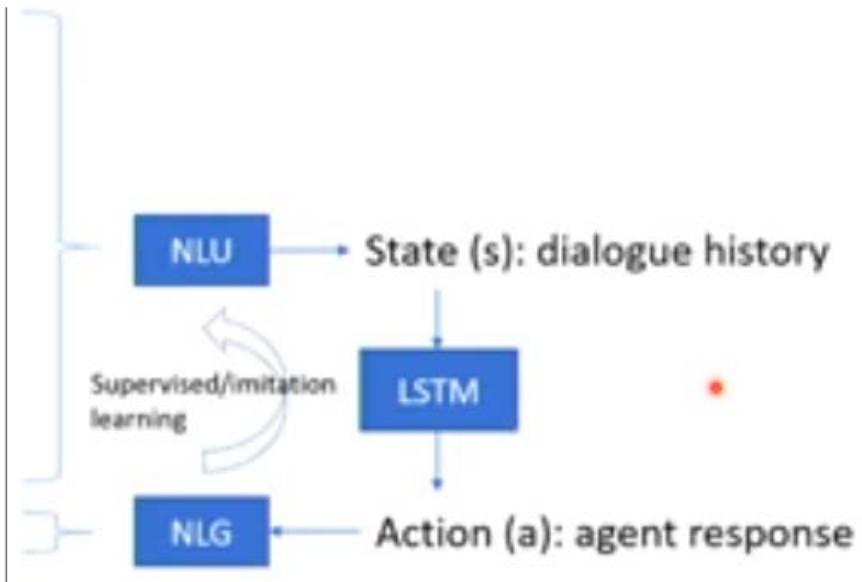
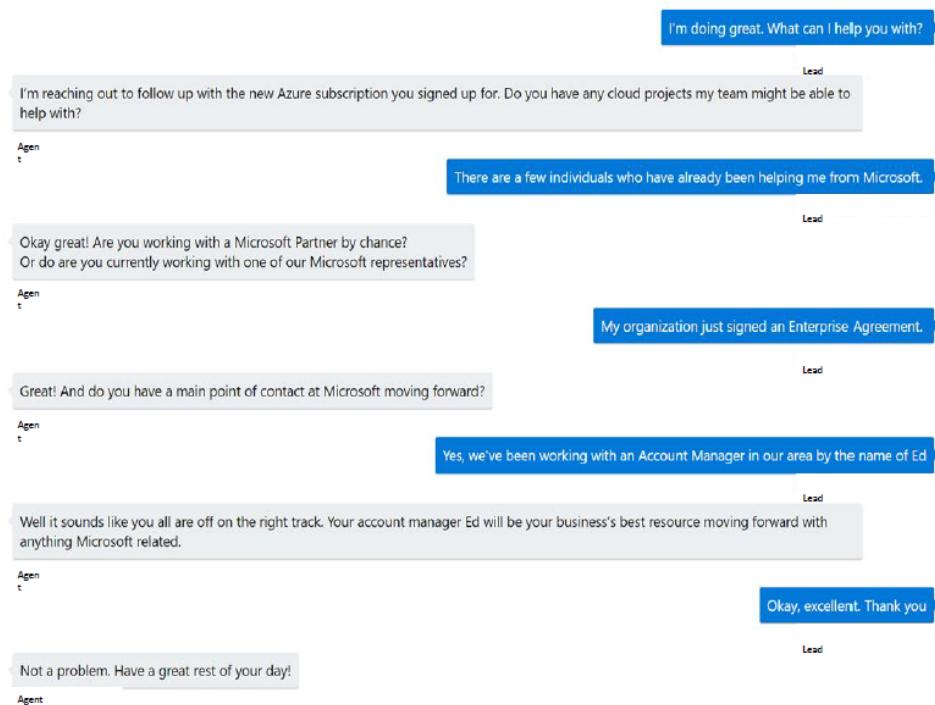
- A full representation of the system's belief of the user's goal at any point during the dialogue
- Used for making API calls

Movies				
Date	11/15/16			
Time	6 pm	7 pm	8 pm	9 pm
# of tickets	2	3		
Movie name	Inferno	Trolls		
Movie theatre	Century 16			

Restaurants			
Date	11/15/16		
Time	6:30 pm	7 pm	7:30 pm
Cuisine	Mexican		
Restaurant	Vive Sol		



Dialogue policy learning: select the best *action* according to *state* to maximize *success rate*



Dialogue policy learning: select the best *action* according to *state* to maximize *success rate*

I'm doing great. What can I help you with?

Lead

I'm reaching out to follow up with the new Azure subscription you signed up for. Do you have any cloud projects my team might be able to help with?

Agent

There are a few individuals who have already been helping me from Microsoft.

Lead

Okay great! Are you working with a Microsoft Partner by chance? Or do you currently working with one of our Microsoft representatives?

Agent

My organization just signed an Enterprise Agreement.

Lead

Great! And do you have a main point of contact at Microsoft moving forward?

Agent

Yes, we've been working with an Account Manager in our area by the name of Ed.

Lead

Well it sounds like you all are off on the right track. Your account manager Ed will be your business's best resource moving forward with anything Microsoft related.

Agent

Okay, excellent. Thank you

Lead

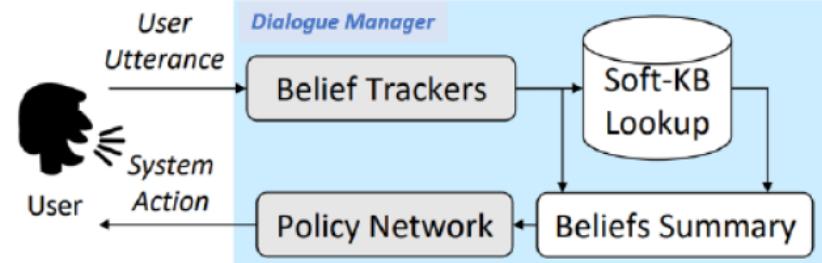
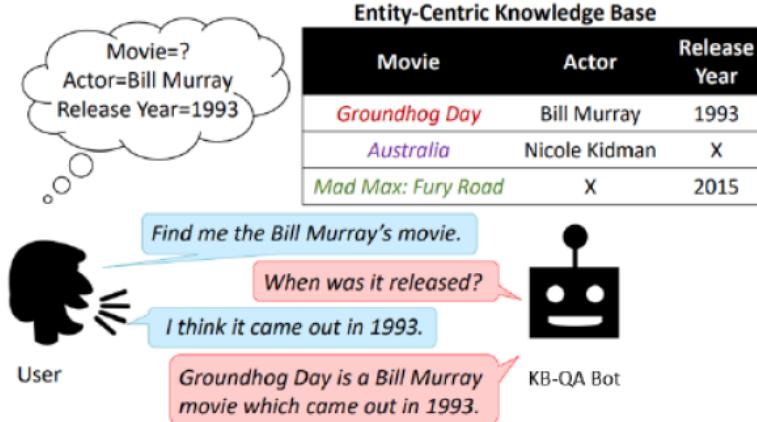
Not a problem. Have a great rest of your day!

Agent



Movie on demand [[Dhingra+ 17](#)]

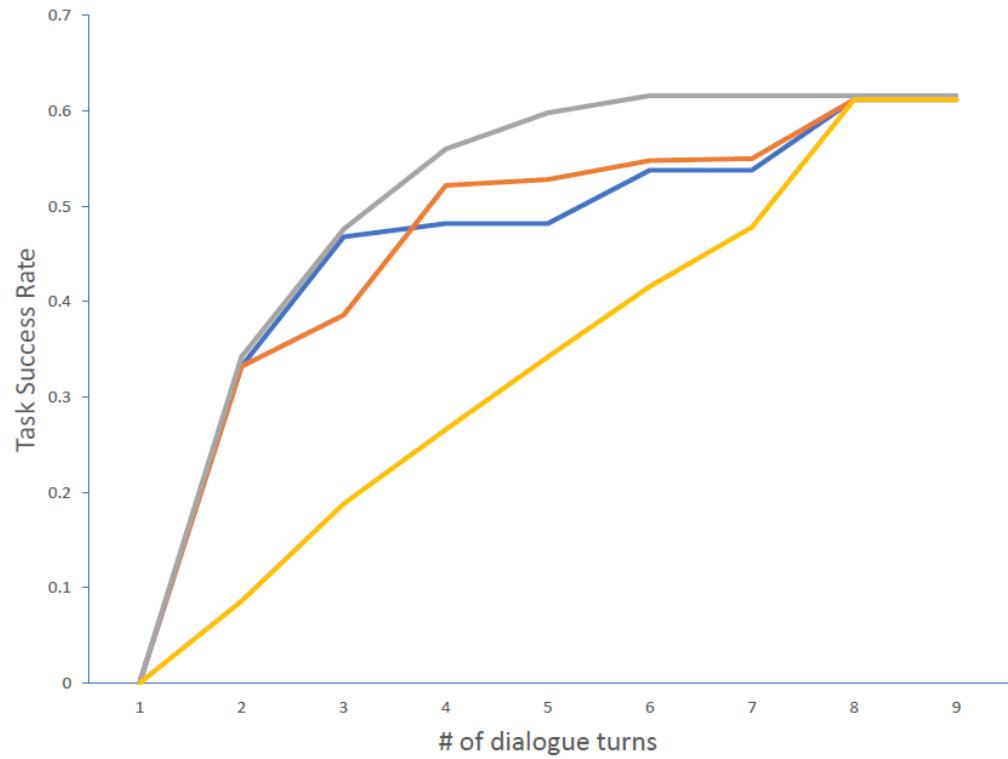
- PoC: leverage Bing tech/data to develop task-completion dialogue (Knowledge Base Info-Bot)

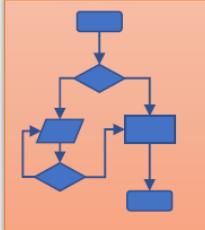
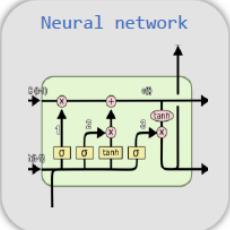


[[Dhingra+ 17](#)]

Learning what to ask next, and when to stop

- Initial: ask all questions in a randomly sampled order
- Improve via learning from Bing log
 - Ask questions that users can answer
- Improve via encoding knowledge of database
 - Ask questions that help reduce search space
- Finetune using agent-user interactions
 - Ask questions that help complete the task successfully via RL

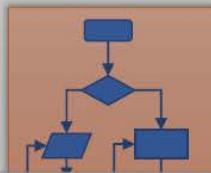


Programmatic	Declarative	Machine Learning
<pre>this.dialogs.add(new WaterfallDialog(GET_FORM_DATA, [this.askForCity.bind(this), this.collectAndDisplayName.bind(this)])); async collectAndDisplayName(step) { ... }</pre>	 <pre><rule> <if> city == null </if> <then> Which city? </then> ... </pre>	
<ul style="list-style-type: none"> ▼ Accessible to non-experts ▼ Easy to debug ▲ Explicit Control ▲ Support for complex scenarios ▼ Ease of Modification ▼ Handle Unexpected Input ▼ Improve / Learn from conversations ▲ No Dialog Data Required 	<ul style="list-style-type: none"> ▲ Accessible to non-experts ▼ Easy to debug ▲ Explicit Control ▼ Support for complex scenarios ▼ Ease of Modification ▼ Handle Unexpected Input ▼ Improve / Learn from conversations ▲ No Dialog Data Required 	<ul style="list-style-type: none"> ▲ Accessible to non-experts ▲ Easy to debug ▼ Explicit Control ▲ Support for complex scenarios ▲ Ease of Modification ▼ Handle Unexpected Input ▲ Improve / Learn from conversations ▼ Requires Sample Dialog Data

Programmatic

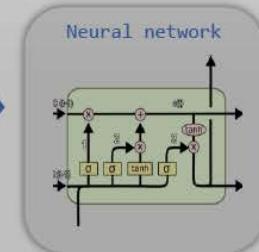
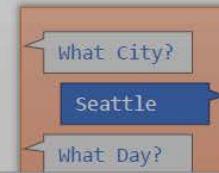
```
this.dialogs.add(  
  new WaterfallDialog(GET_FORM_DATA,  
  [  
    this.askForCity.bind(this),  
    this.collectAndDisplayName.bind(this)  
  ]  
);  
async collect
```

Declarative



```
<rule>  
  <if>  
    city == null  
  </if>  
  <then>  
    Which city?  
  </then>
```

Machine Learning



on-experts

complex scenarios

iteration

One Solution Does Not Fit All

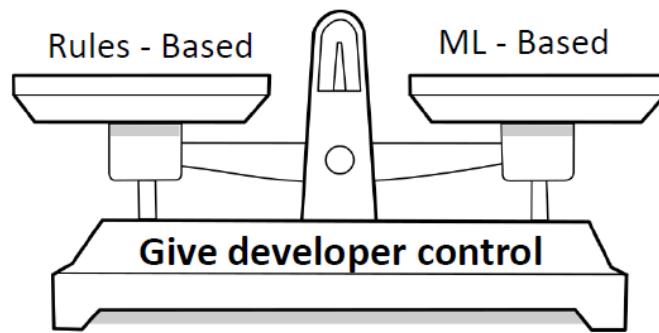
- ▼ Accessible
- ▼ Easy to de
- ▲ Explicit Co
- ▲ Support fo
- ▼ Ease of M
- ▼ Handle Unexpected Input
- ▼ Improve / Learn from conversations
- ▲ No Dialog Data Required

- ▼ Handle Unexpected Input
- ▼ Improve / Learn from conversations
- ▲ No Dialog Data Required

- ▼ Handle Unexpected Input
- ▲ Improve / Learn from conversations
- ▼ Requires Sample Dialog Data

Goal: Best of both worlds

- Good for garden path
- Not data intensive
- Explicit Control
- Easily interpretable



- Handle unexpected input
- Learn from usage data
- Often viewed as black box

Start with rules-based policy => Grow with Machine Learning

Make ML more controllable by visualization

Not unidirectional : Rules-based policy can evolve side-by-side with ML Model

Conversation Learner – building a bot interactively

What is it: A system built on the principles of Machine Teaching, that enables individuals with no AI experience (designers, business owners) to build task-oriented conversational bots

Goal: Push the forefront of research on conversational systems using input from enterprise customers and product teams to provide grounded direction for research

Status: In private preview with ~50 customers to various levels of prototyping

Hello World Tutorial

Primary repository with samples:

<https://github.com/Microsoft/ConversationLearner-samples>

ConvLab

Published @ <https://arxiv.org/abs/1904.08637>

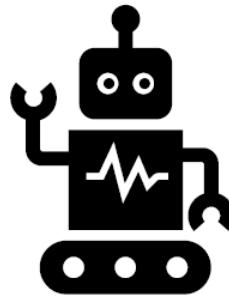
Fully annotate data

for training individual components or
end-to-end models with supervision

Speaker	Utterance	Annotation
User	am looking for a place to stay that has cheap price range it should be in a type of hotel	Dialog acts {"Hotel-Inform": [{"Pricerange": "cheap"}]}
System	Okay, do you have a specific area you want to stay in?	State {"hotel": {"name": "not mentioned", "area": "not mentioned", "parking": "not mentioned", "pricerange": "cheap", "stars": "not mentioned", "internet": "not mentioned", "type": "hotel"}} ... Dialog acts {"Hotel-Request": [{"Area": "?"}]}
User	no, i just need to make sure it's cheap. oh, and i need parking	Dialog acts {"negate", "Hotel-Inform": [{"Pricerange": "cheap"}, {"Parking": "yes"}]}
System	I found 1 cheap hotel for you that includes parking. Do you like me to book it?	State {"hotel": {"name": "not mentioned", "area": "not mentioned", "parking": "yes", "pricerange": "cheap", "stars": "not mentioned", "internet": "not mentioned", "type": "hotel"}} ... Dialog acts {"Hotel-Inform": [{"Price": "cheap"}, {"Choice": "1"}, {"Parking": "none"}]}

User Simulators

for reinforcement learning
1 rule-based simulator
2 data-driven simulators



SOTA Baselines

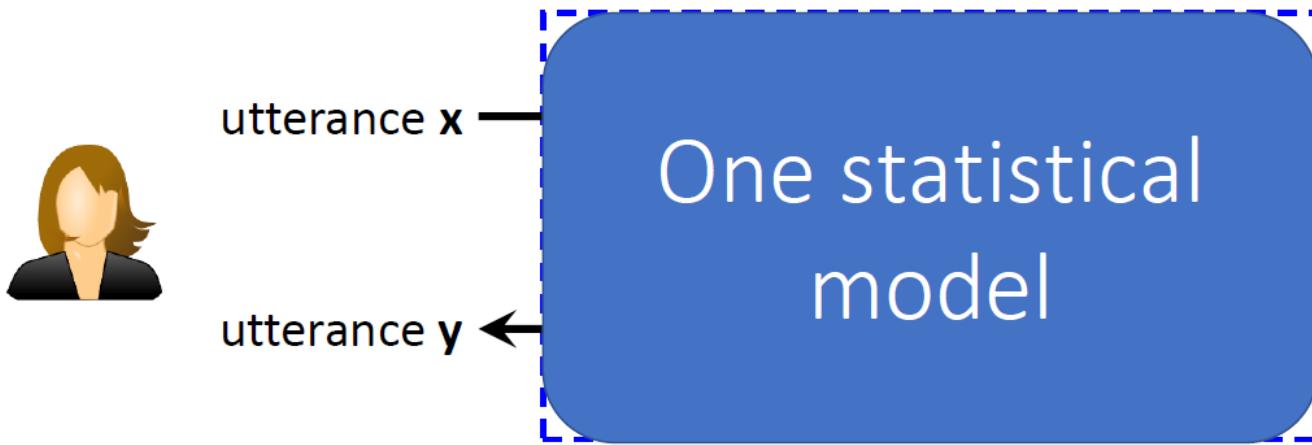
Multiple models for each component
Multiple end-to-end system recipes



Outline

- Introduction
- Question Answering (QA)
- Task-oriented Dialogues
- **Fully Data-Driven Conversation Models and Chatbots**

Motivation



Move towards **fully data-driven, end-to-end** dialogue systems.

Social Bots

- Fully end-to-end systems so far most successfully applied to **social bots** or **chatbots**:
 - Commercial systems: Amazon Alexa, Xiaolce, etc.
- Why social bots?
 - Maximize **user engagement** by generating **enjoyable** and **more human-like** conversations
 - Help **reduce user frustration**
 - **Influence dialogue research** in general (social bot papers often cited in task-completion dialogue papers)



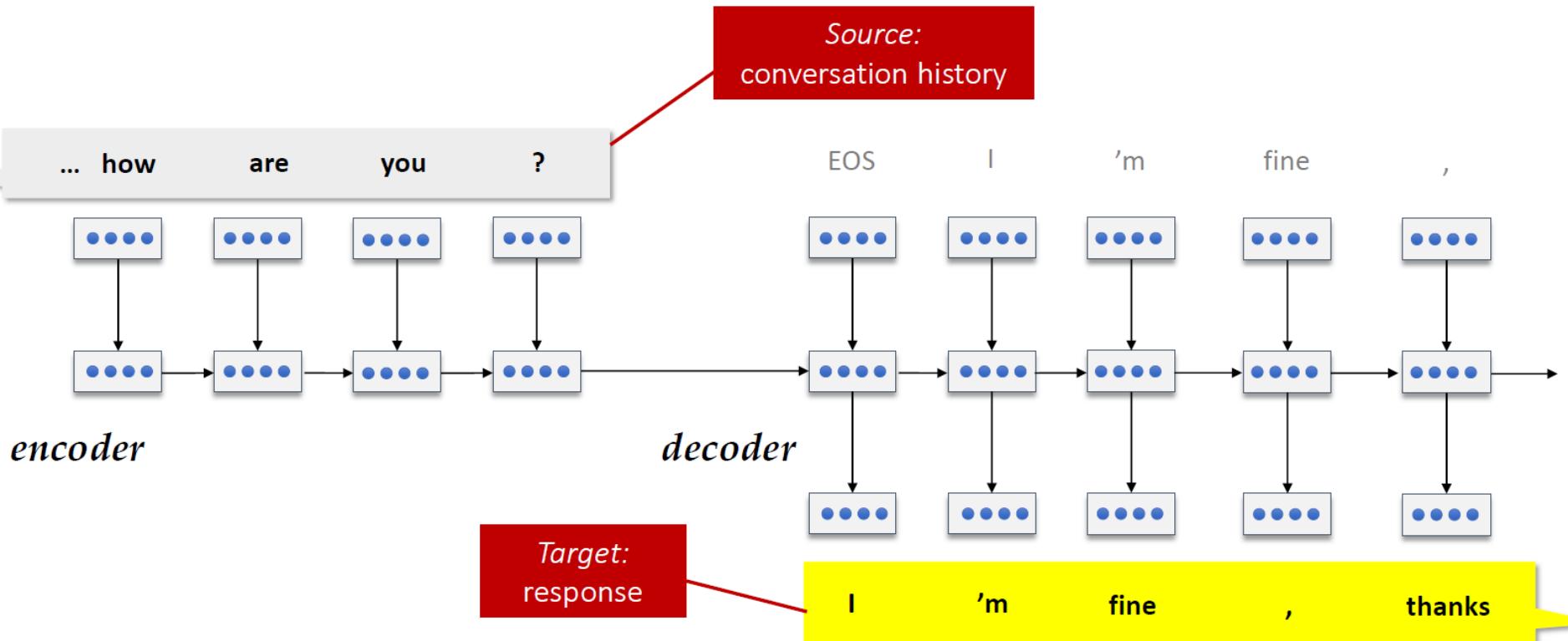
Historical overview

Earlier work in fully data-driven response generation:

- **2010:** Response retrieval system (IR) [[Jafarpour+ 10](#)]
- **2011:** Response generation using Statistical Machine Translation (phrase-based MT) [[Ritter+ 11](#)]
- **2015:** First neural response generation systems (RNN, seq2seq)
[[Sordoni+ 15](#); [Vinyals & Le 15](#); [Shang+ 15](#)]

Neural Models for Response Generation

[[Sordoni+ 15](#);
[Vinyals & Le 15](#);
[Shang+ 15](#)]



Similar to sequence models in Neural Machine Translation (NMT), summarization, etc.
Uses either RNN, LSTM, GRU, Pointer-Generator Networks, Transformer, etc.

Neural Response Generation:

Difference with other tasks (e.g., machine translation)

- **Data:** some training sets (social media) are **HUGE**

For example, Twitter (as of 2016):

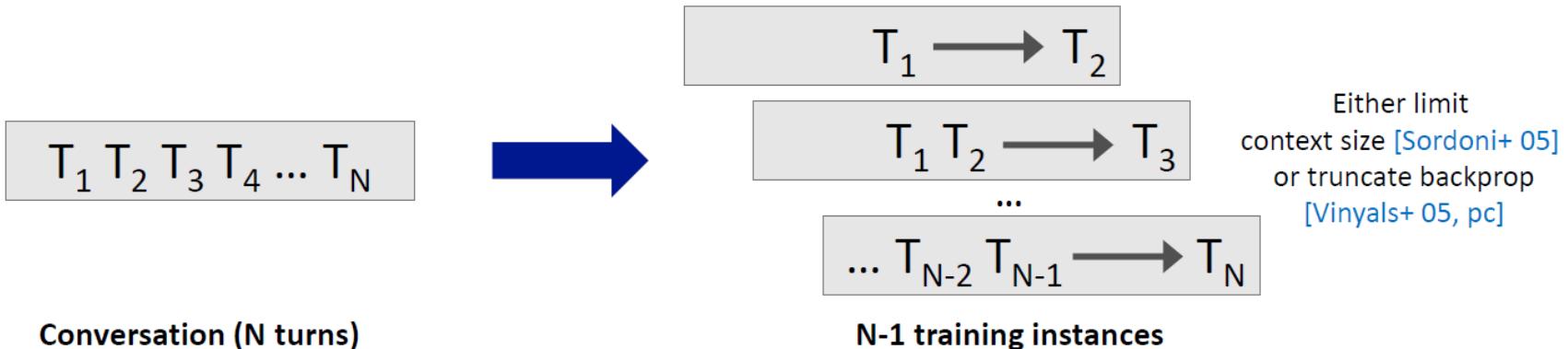
304M monthly active users

500M tweets per day (6M conversations per day)

→ often train on subset of the data; leaner and “faster” models preferred



- **System input:** **LONG** conversation history



Neural conversation engine

Neural model learns the general shape of conversation, and the system output is situationally appropriate and coherent...

H: Hi
M: How are you?
H: Not bad
M: What's wrong?
H: Nothing really
M: Wanna talk about it? I'm here if you wanna talk x
H: Well, I'm just tired
M: Me too, I'm here if you wanna talk
H: Thanks, what are you doing?
M: You're welcome and I'm watching TV and you?
H: On Netflix?
M: Yeah, I'm watching it on Netflix.
H: What are you watching?
M: A TV show.
H: Nice, did you have dinner already?
M: No, not yet. what about you?

A more goal-oriented example. No grounding into a real calendar, but the “shape” of the conversation is fluent and plausible...

H: Would Thursday afternoon work sometime?
M: Yeah, sure. just let me know when you're free.
H: After lunch is probably the best time
M: Okay, sounds good. Just let me know when you're free.
H: Would 2pm work for you?
M: Works for me.
H: Well let's say 2pm then I'll see you there
M: Sounds good.

Neural E2E system trained on 35M Twitter conversations.

Challenge: The blandness problem

How was your weekend?

I don't know.



What did you do?



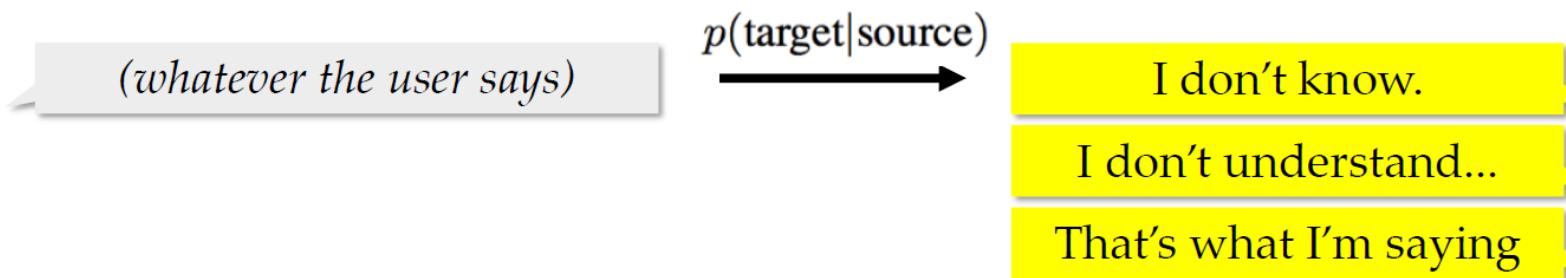
I don't understand what you are talking about.

This is getting boring...

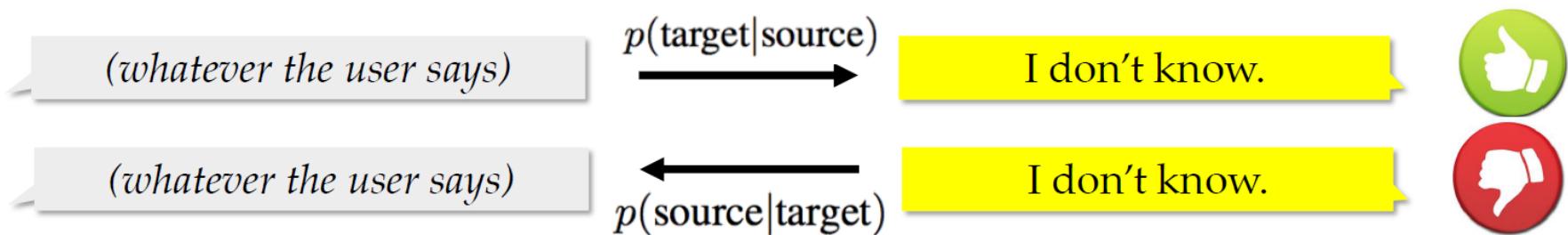
Yes that's what I'm saying.

Blandness problem: cause and remedies

Common MLE objective (maximum likelihood)



Mutual information objective:



Mutual Information for Neural Network Generation

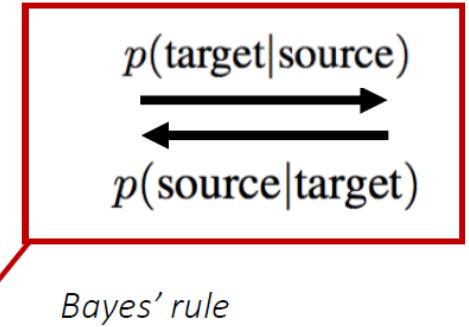
Mutual information objective:

$$\hat{T} = \arg \max_T \left\{ \log \frac{p(S, T)}{p(S)p(T)} \right\}$$

$$\hat{T} = \arg \max_T \left\{ \boxed{\log p(T|S)} - \boxed{\lambda \log p(T)} \right\}$$

standard
likelihood anti-LM

$$\hat{T} = \arg \max_T \left\{ (1 - \lambda) \log p(T|S) + \lambda \log p(S|T) \right\}$$



Bayes' theorem

Sample outputs (MMI)

Wow sour starbursts really do make your mouth water... mm drool.

Can I have one?

Of course you can! They're delicious!

Milan apparently selling Zlatan to balance the books... Where next, Madrid?

I think he'd be a good signing.



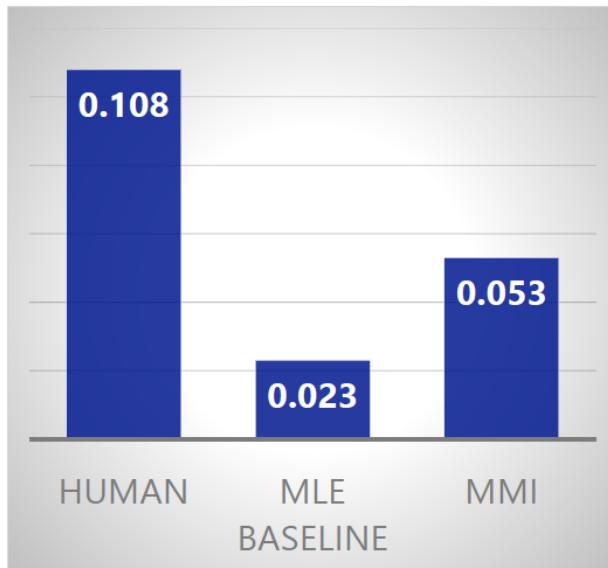
'tis a fine brew on a day like this! Strong though, how many is sensible?

Depends on how much you drink!

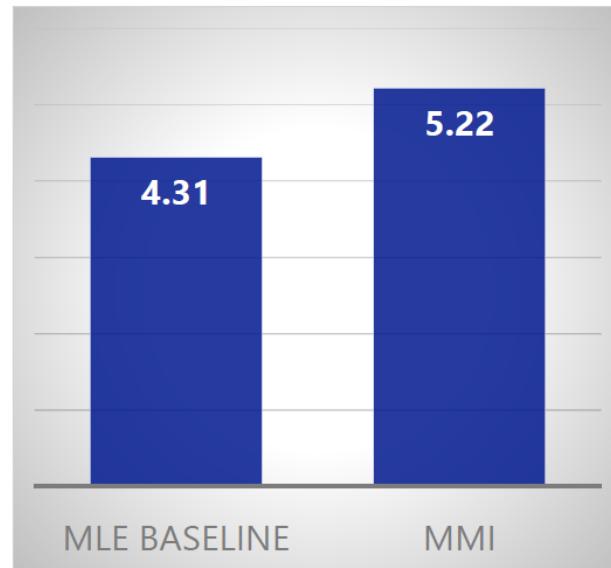
Well he was on in Bromley a while ago... still touring.

I've never seen him live.

MLE vs MMI: results



Lexical diversity
(# of distinct tokens / # of words)



BLEU

MMI: best system in Dialogue Systems Technology Challenge 2017 (DSTC, E2E track)



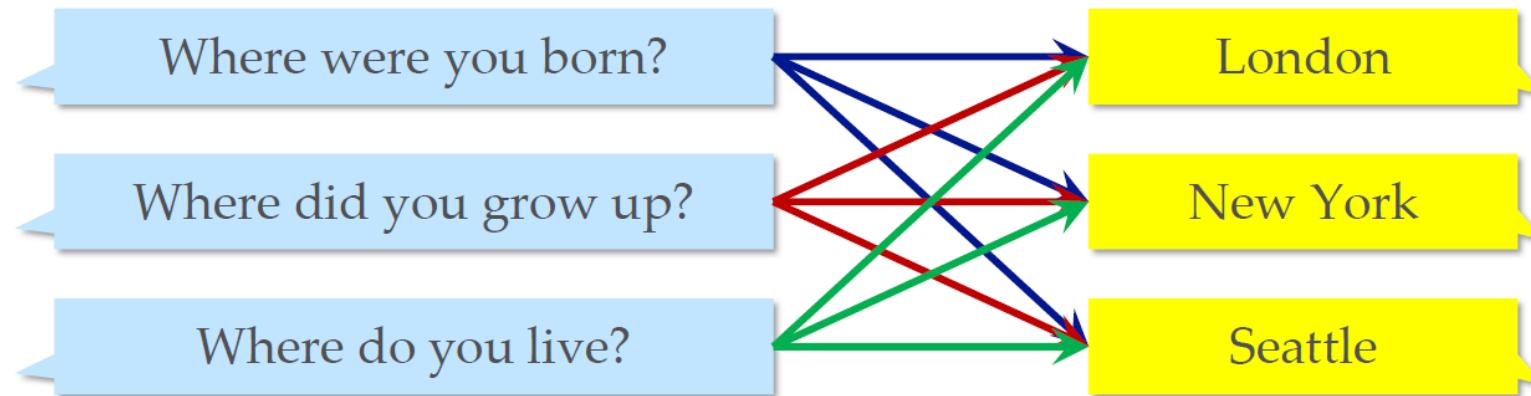
Challenge: The consistency problem

- E2E systems often exhibit **poor response consistency**:

<i>message</i>	Where is your hometown?
<i>response</i>	I live in Indonesia.
<i>message</i>	Where are you from?
<i>response</i>	England, you?
<i>message</i>	In which city do you live now?
<i>response</i>	I live in Jakarta.
<i>message</i>	In which country do you live now?
<i>response</i>	I live in Indonesia.
<i>message</i>	What is your major?
<i>response</i>	Business, you?
<i>message</i>	What did you study in college?
<i>response</i>	Psychology, you?
<i>message</i>	How old are you?
<i>response</i>	I'm 18.
<i>message</i>	What is your age?
<i>response</i>	I'm 16.

The consistency problem: why?

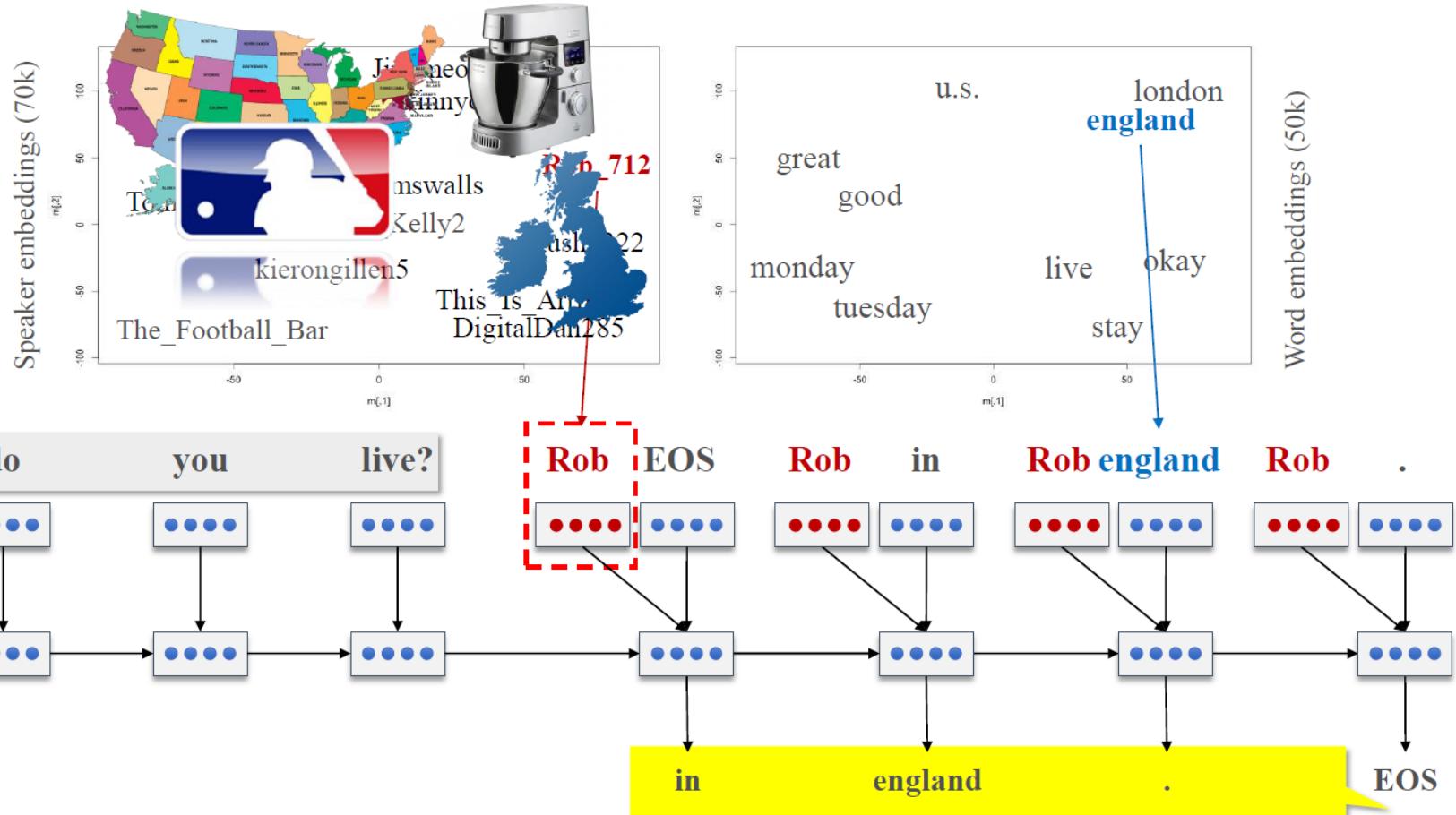
Conversational data:



NOT
1-to-1

$$P(\text{response} \mid \text{query}, \text{SPEAKER_ID})$$

Personalized Response Generation [Li+ 2016b]



Persona model results

Baseline model:

message Where is your hometown?

response I live in **Indonesia**.

message Where are you from?

response **England**, you?

message In which city do you live now?

response I live in **Jakarta**.

message In which country do you live now?

response I live in **Indonesia**.

message What is your major?

response **Business**, you?

message What did you study in college?

response **Psychology**, you?

message How old are you?

response I'm **18**.

message What is your age?

response I'm **16**.

Persona model using speaker embedding: [Li+ 16b]

message Where is your hometown?

response I'm from **England**.

message Where are you from?

response I'm from **England**.

message In which city do you live now?

response I live in **London**.

message In which country do you live now?

response I live in **England**.

message What is your major?

response **Business**. You?

message What did you study in college?

response I did **business** studies.

message How old are you?

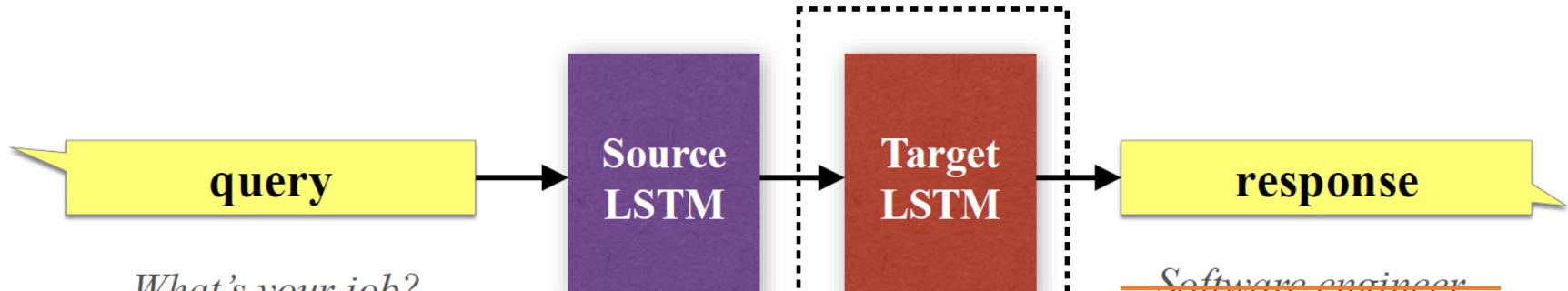
response I'm **18**.

message What is your age?

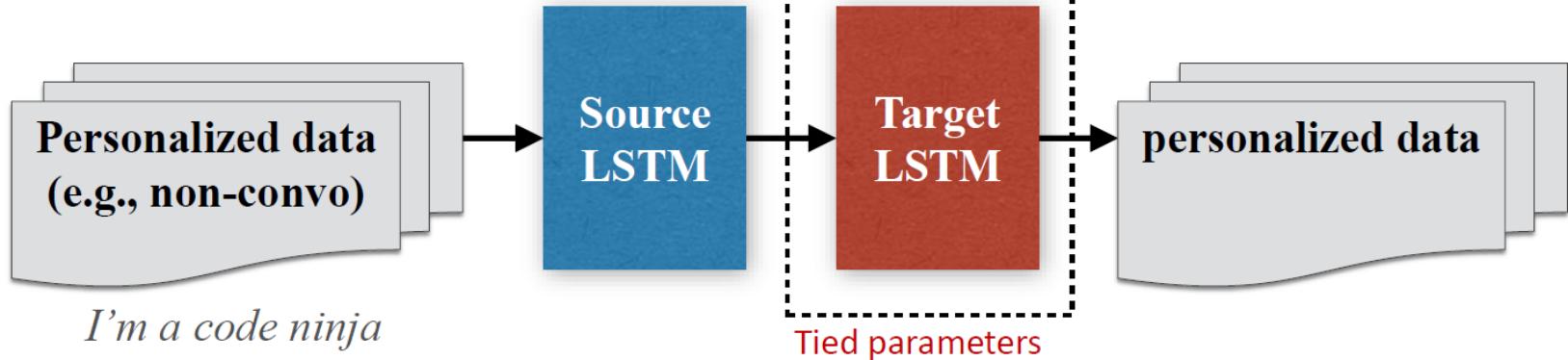
response I'm **18**.

Personal modeling as multi-task learning [Luan+ 17]

Seq2Seq

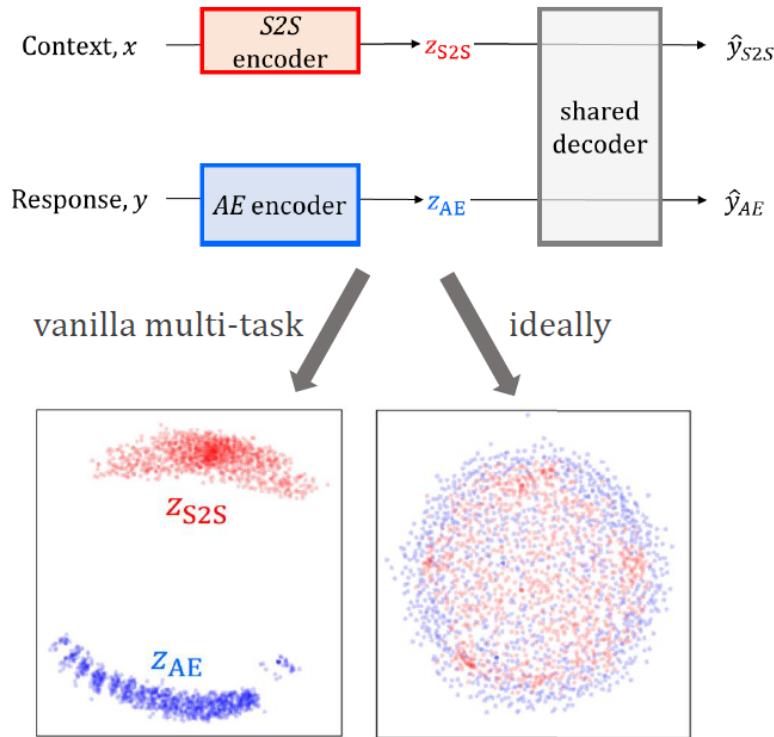


Autoencoder



Challenges with multi-task learning

[Gao + 19]



So we add regularization:

$$\mathcal{L} = -\frac{1}{|y|} \log p(y|z_{S2S}) - \frac{1}{|y|} \log p(y|z_{AE}) + \alpha \mathcal{L}_{\text{fuse}} + \dots$$

Vanilla S2S
+ Mtask
objective

where:

$$\mathcal{L}_{\text{fuse}} = \sum_{i \in \text{batch}} \frac{d(z_{S2S}(x_i), z_{AE}(y_i))}{n}$$

cross-space distance

$$- \sum_{i,j \in \text{batch}, i \neq j} \frac{d(z_{S2S}(x_i), z_{S2S}(x_j))}{n^2 - n} - \sum_{i,j \in \text{batch}, i \neq j} \frac{d(z_{AE}(y_i), z_{AE}(y_j))}{n^2 - n}$$

same-space distance

Improving personalization with multiple losses

[AI-Rfou+ 16]

- **Single-loss:**

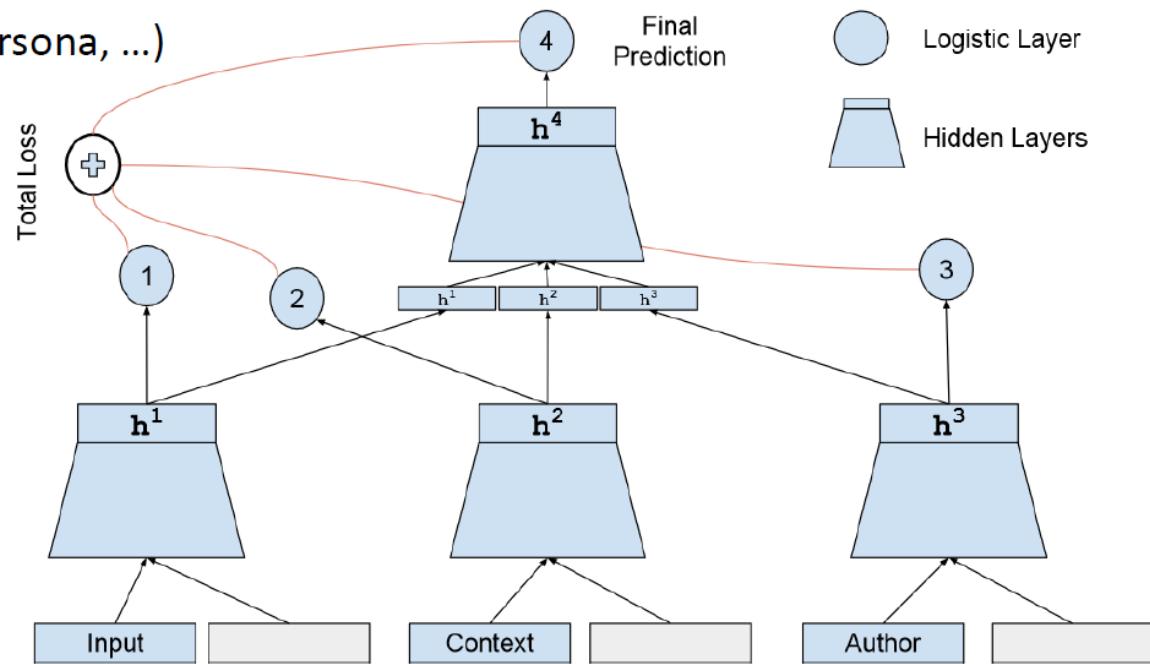
$P(\text{response} \mid \text{context, query, persona, ...})$

Problem with single-loss:
context or query often
“explain away” persona

- **Multiple loss adds:**

$P(\text{response} \mid \text{persona})$
 $P(\text{response} \mid \text{query})$
etc.

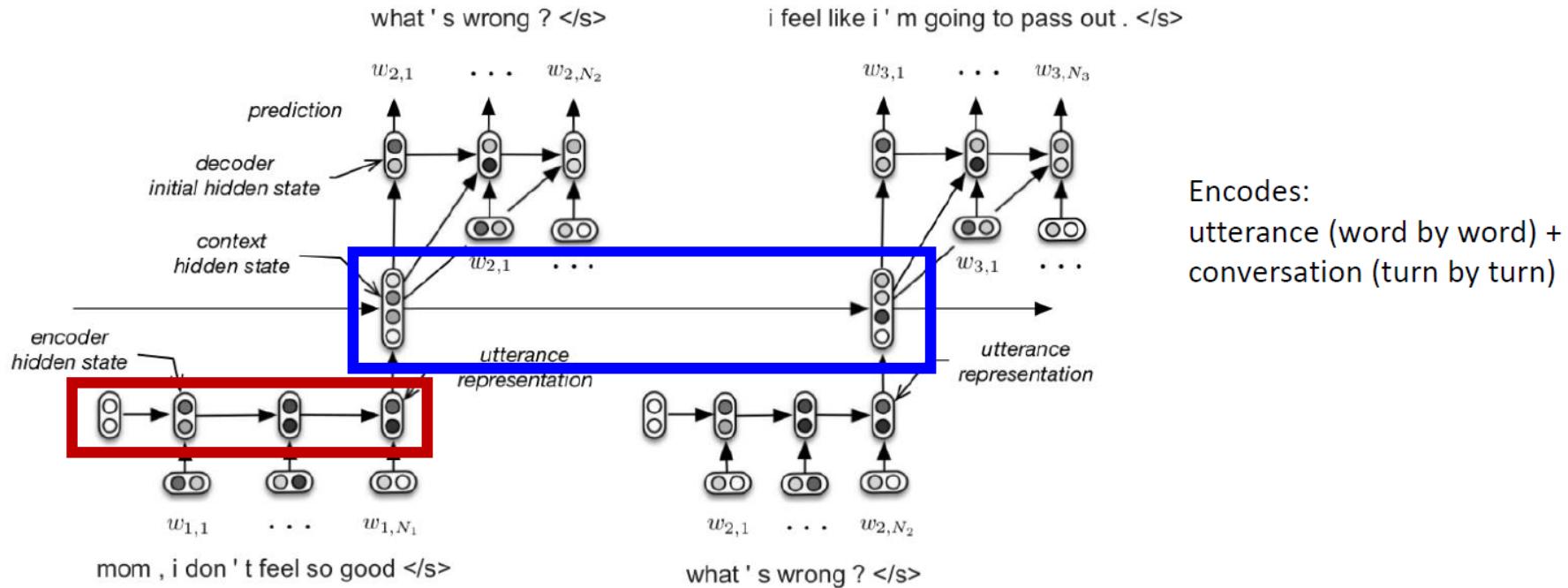
Optimized so that
persona can “predict”
response all by itself
→ more robust
speaker embeddings



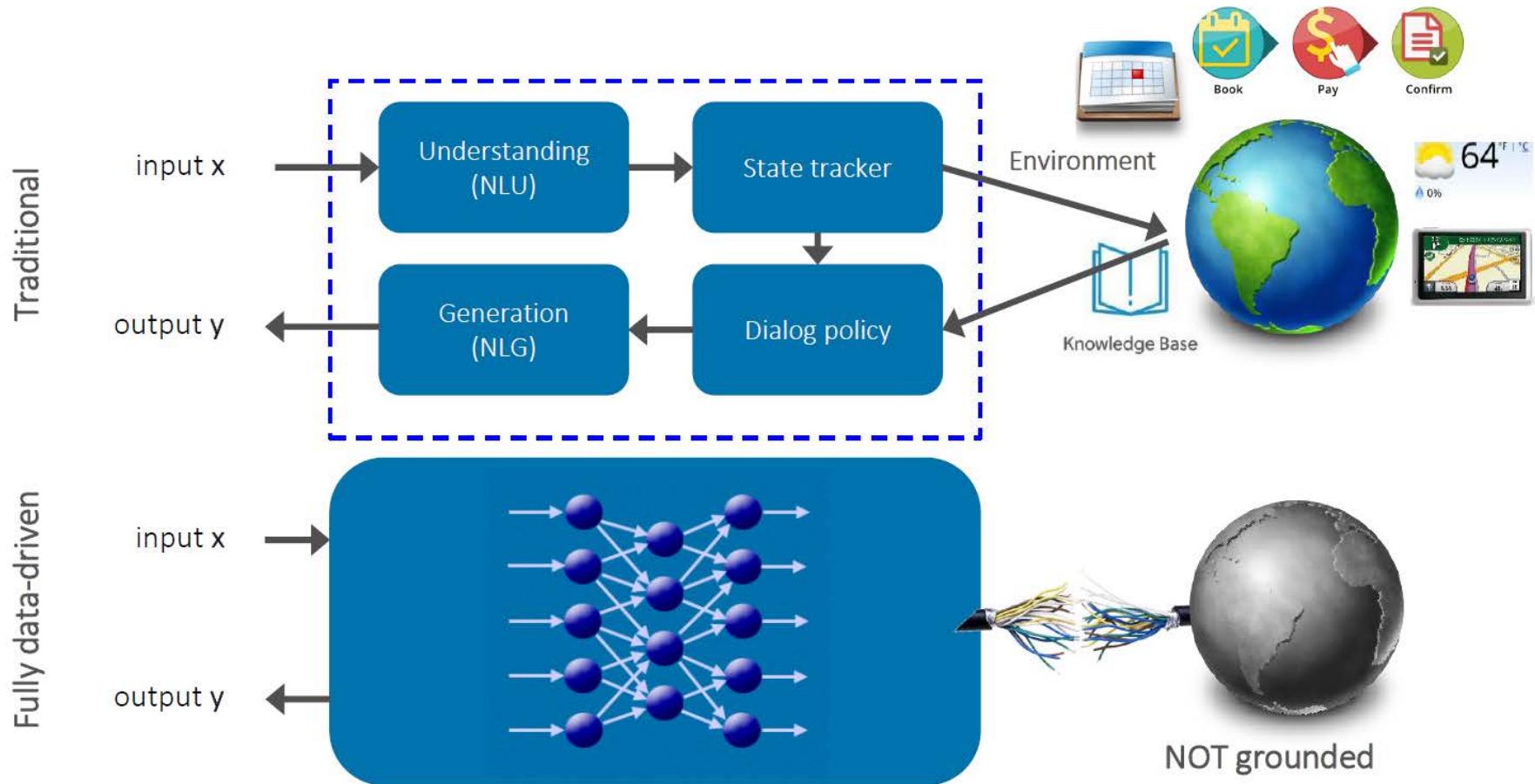
Challenge: Long conversational context

It can be challenging for LSTM/GRU to encode very long context
(i.e. more than 200 words: [[Khandelwal+ 18](#)])

- **Hierarchical Encoder-Decoder (HRED) [[Serban+ 16](#)]**



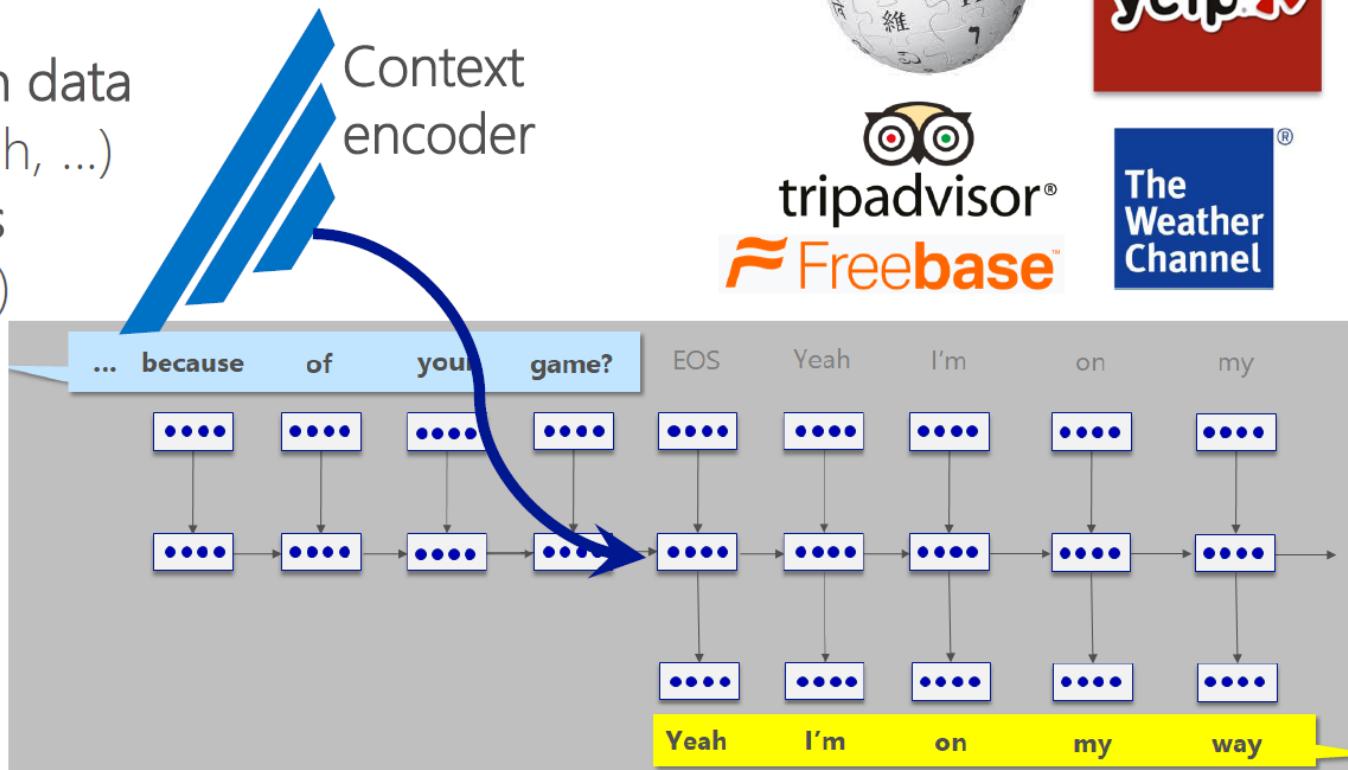
Towards Grounded E2E Conversation Models



E2E Conversation Models in the real world

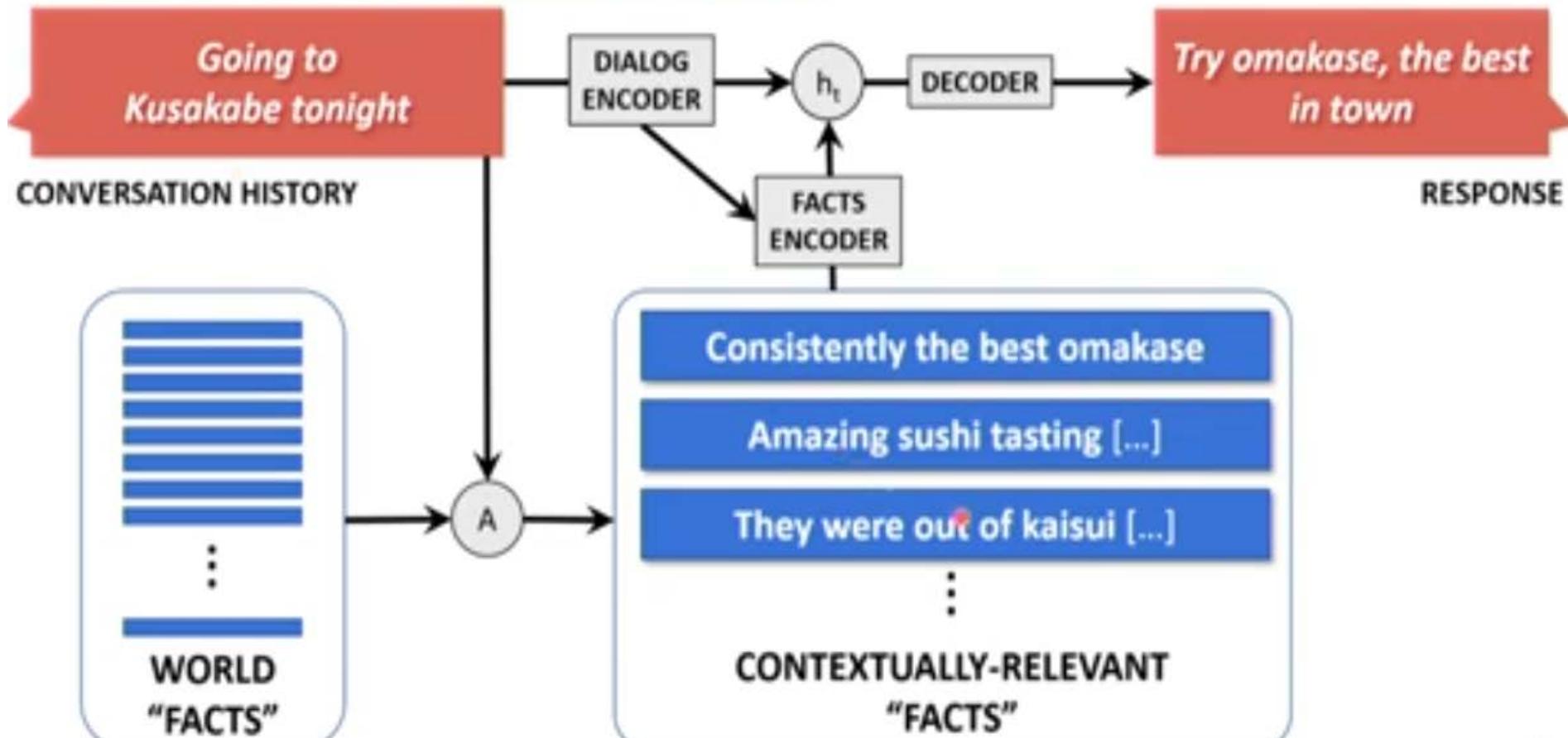


Personalization data
(ID, social graph, ...)
Device sensors
(GPS, vision, ...)
External
"knowledge"



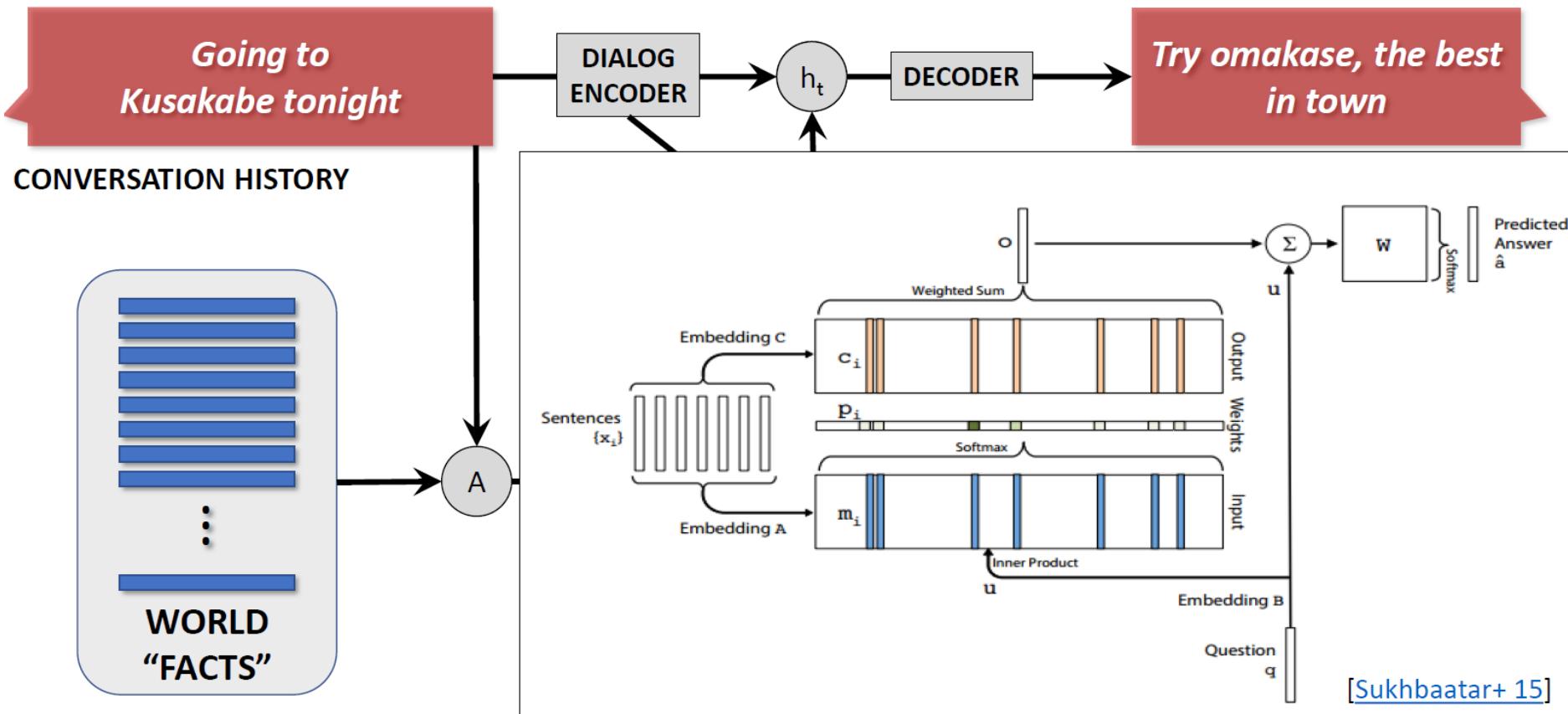
Knowledge-Grounded Neural Conversation Model

[Ghazvininejad+ 17; Agarwal+ 18; Liu+ 18]



Knowledge-Grounded Neural Conversation Model

[Ghazvininejad+ 17; Agarwal+ 18; Liu+ 18]





Sample knowledge-grounded responses

I would give ALMOST anything for some [Mexican restaurant] right now.

Me too [Jalapeno sauce is really good.]

Obsessed with [jewelry company] :-*

oh my gosh [obsessed with their bracelets and the meaning behind them!]

Visiting the celebs at Los Angeles airport - [...] w/ 70 others

Nice airport terminal. [Have a safe flight.]

Experimental results (23M conversations): outperforms competitive neural baseline (human + automatic eval)

Conversations grounded in *Full-Length Documents*

Vesna Vulović

From Wikipedia, the free encyclopedia

Vesna Vulović (Serbian Cyrillic: Весна Вуловић; pronounced [v̞e:sn̞a vu:lovič]; 3 January 1950 – 23 December 2016) was a Serbian flight attendant. She holds the Guinness world record for surviving the highest fall without a parachute: 10,160 metres (33,330 ft). Her fall took place after an explosion tore through the baggage compartment of JAT Flight 367 on 26 January 1972, causing it to crash near Srbská Kamenice, Czechoslovakia. She was the sole survivor of the crash that air safety investigators attributed to a briefcase bomb. The Yugoslav authorities suspected that Croatian nationalists were to blame, but no one was ever arrested. Following the crash, Vulović spent days in a coma and was hospitalized for several months. She suffered a fractured skull, three broken vertebrae, two broken legs, broken ribs and a fractured pelvis. These injuries resulted in her being temporarily paralyzed from the waist down. She made an almost complete recovery but continued to walk with a limp. Vulović maintained that she had no memory of the incident and thus had no qualms about flying in the



Vulović in the early 1970s

Born	3 January 1950 Belgrade, PR Serbia, FPR Yugoslavia
Died	23 December 2016 (aged 66) Belgrade, Serbia

A woman fell 30,000 feet from an airplane and survived.

The page states that a 2009 report found the plane only fell several hundred meters.

Well if she only fell a few hundred meters and survived then I'm not impressed at all.

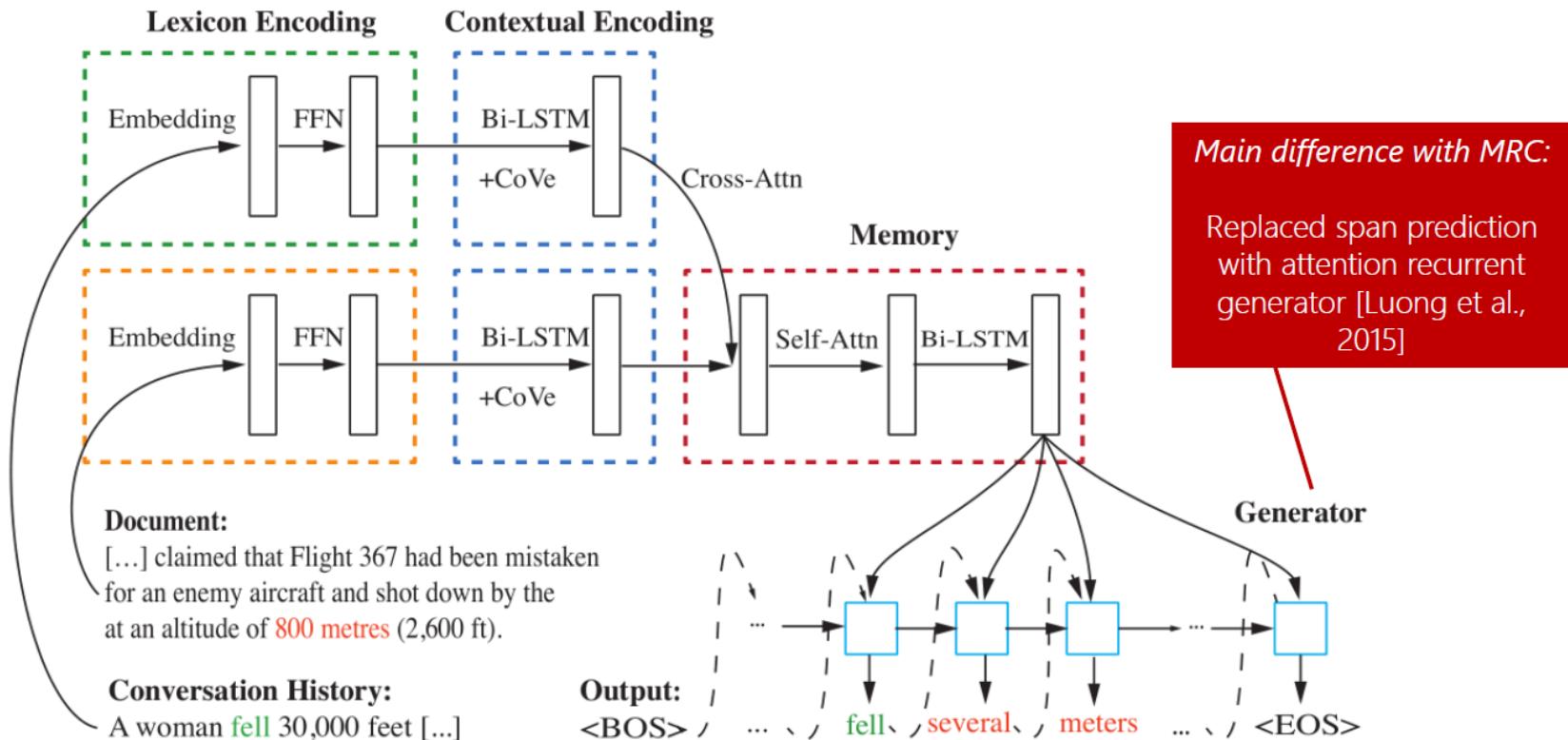
Few hundred meters is still pretty incredible, but quite a bit different than 10,000 meters.

Task: Generate a **human-like response** that is not only conversationally appropriate, but also **informative** (→ useful task) and **grounded** (-> evaluation closer to MRC).

[[Dinan+ 19](#); [Qin+ 19](#)]

Models with Document-Level Grounding

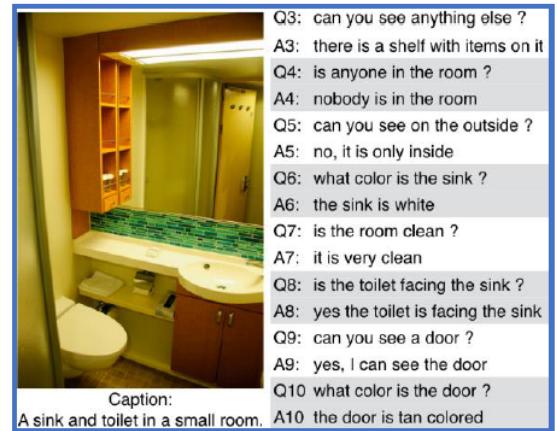
Machine Reading Comprehension-based Model [Qin+ 19]:



Grounded E2E Dialogue Systems

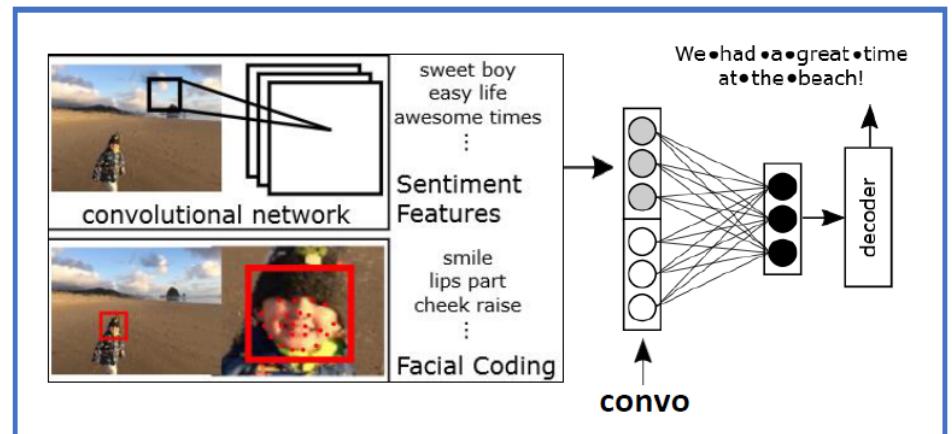
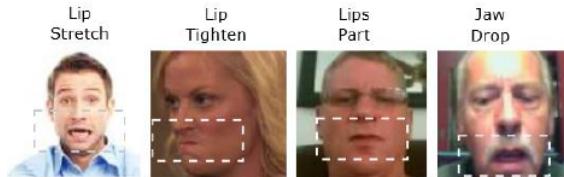
- **Grounding: images**

Conversations around images e.g.,
Q-As [[Das+ 16](#)] or chat [[Mostafazadeh+ 17](#)]



- **Grounding: affect [[Huber+ 18](#)]**

facial actions influence response



Beyond supervised learning

- **Limitations of SL for E2E dialogue:**
 - Train on **human-human** data, test with **human-machine** (Twitter-ese often not what we want at test time.)
 - Optimizes for **immediate reward** $p(T_N | \dots T_{N-1})$, not **long-term reward**
 - No **user feedback loop**
- **Emergence of reinforcement learning (RL) for E2E dialogue**
 - Tries to promote long-term dialogue success

Deep Reinforcement Learning for E2E Dialogue

[Lit 16c]

- REINFORCE algorithm [Williams+ 92]

$$J(\theta) = \mathbb{E}[R(s_1, s_2, \dots, s_N)]$$

$$\nabla J(\theta) = \nabla \log p(s_1, s_2, \dots, s_N) R(s_1, s_2, \dots, s_N)$$

reward function

$$\nabla J(\theta) = \nabla \log \prod_i p(s_i | s_{i-1}) R(s_1, s_2, \dots, s_N)$$

**what we
want to learn**

- Reward functions:

1. Ease of answering: $-Pr(\text{Dull Response}|s_i)$
2. Information flow: $-\log \text{Sigmoid} \cos(s_1, s_2)$
3. Meaningfulness: $\log p_{\text{seq2seq}}(s_1|s_2) + \log p_{\text{seq2seq}}(s_2|s_1)$

Simulation (without RL)

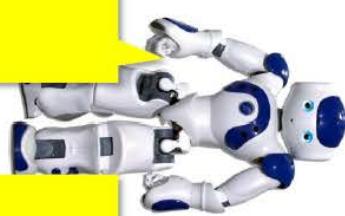


See you later!



See you later!

See you later!



See you later!

See you later!



See you later!

Simulation (with RL)



How old are you ?

i 'm 4, why are
you asking ?



I thought you were 12 .

What made you think so ?

I don't know what you are talking
about .

You don't know what you are saying.



I don't know what you are talking
about .



Deep RL: Evaluation

- MTurk evaluation (500 responses)

Setting	RL-win	RL-lose	Tie
single-turn general quality	0.40	0.36	0.24
single-turn ease to answer	0.52	0.23	0.25
multi-turn general quality	0.72	0.12	0.16

Conversational datasets (for social bots, E2E dialogue research)

- Survey on dialogue datasets [[Serban+ 15](#)]

Name	Type / Topics	Size
Reddit	Unrestricted	3.2B dialog turns (growing)
Twitter	Unrestricted	N/A (growing)
OpenSubtitles	Movie subtitles	1B words
Ubuntu Dialogue Corpus	Chat on Ubuntu OS	100M words
Ubuntu Chat Corpus	Chat on Ubuntu OS	2B words
Persona-Chat Corpus	Crowdsourced / personalized	164k dialog turns



Evaluating E2E Dialogue Systems

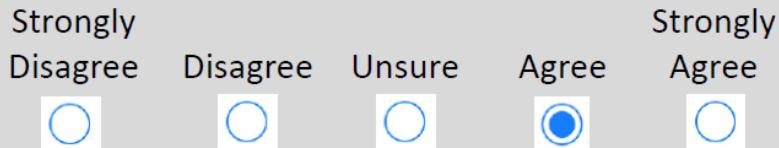
- Human evaluation (crowdsourcing):

Context: ... *Because of your game?*

Input: *Yeah, I'm on my way now*

Response: **Ok good luck!**

Is **this** a good¹ response?



¹: replaced as appropriate (relevant, interesting,...)

- Automatic evaluation:

Less expensive, but is it reliable?

Machine-Translation-Based Metrics

- **BLEU** [[Papineni+ 02](#)]: ngram overlap metric

Reference:	John	resigned	yesterday .
System:	Yesterday ,	John	quit .

$$\text{BLEU} = \text{BP} \cdot \exp \left(\sum_n \log p_n \right)$$

$$p_n = \frac{\sum_i \sum_{g \in n\text{-grams}(h_i)} \max_j \{\#_g(h_i, r_{i,j})\}}{\sum_i \sum_{g \in n\text{-grams}(h_i)} \#_g(h_i)}$$

- **NIST** [[Doddington+ 02](#)]

- Seldom used in dialogue, but copes with **blandness issue**
- Considers info gain of each ngram: score(*interesting calculation*) >> score(*of the*)

- **METEOR**

- Accounts for synonyms, paraphrases, etc.

The challenge with MT-based metrics

Input: *How are you?*

Response (gold): **I'm good, thanks.**

Response A: **Good thanks !**

Response B: Doing pretty **good thanks**

Response C: Doing well thank you !

Semantically equivalent
(as in Machine Translation)

Response D: Fantastic . How are you ?

Response E: **I'm** getting sick again .

Response F: Bored . you ?

Pragmatically appropriate

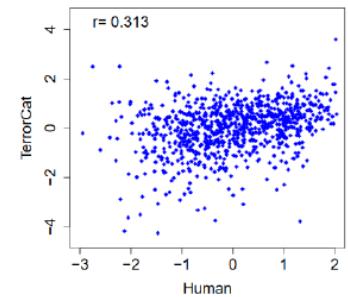
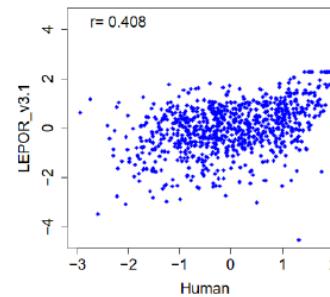
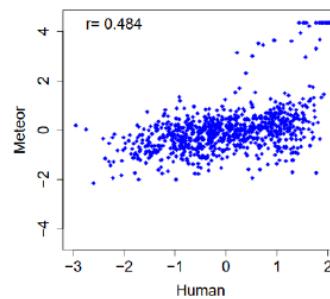
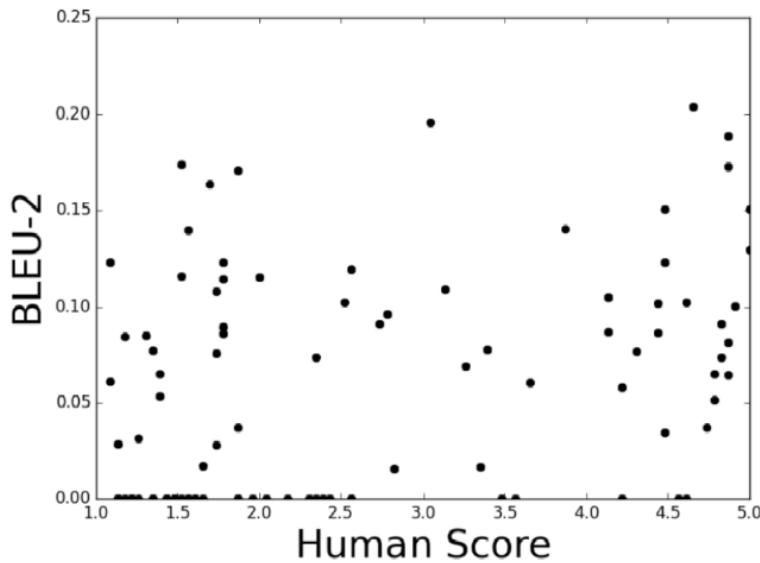
Response G: Sleepy .

Response H: Terrible tbh

Many false negative!

Sentence-level correlation of MT metrics

- Poor correlation with human judgments:



Dialogue task

"How NOT to evaluate dialogue systems"

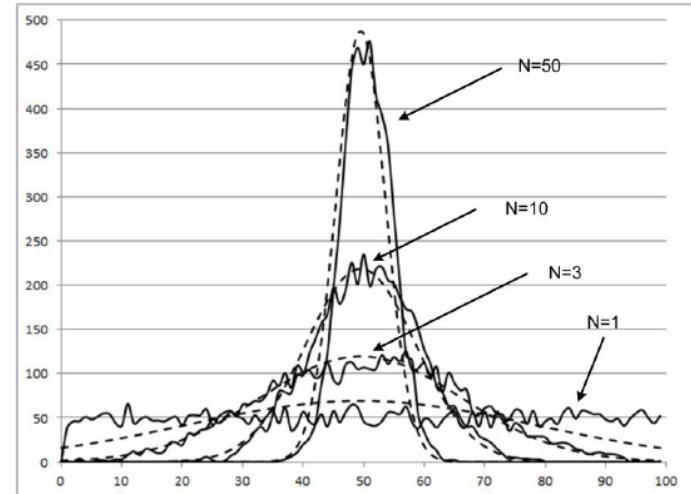
[\[Liu+ 16\]](#)

But same problem even
for **Translation task**

[\[Graham +15\]](#)

The importance of sample size

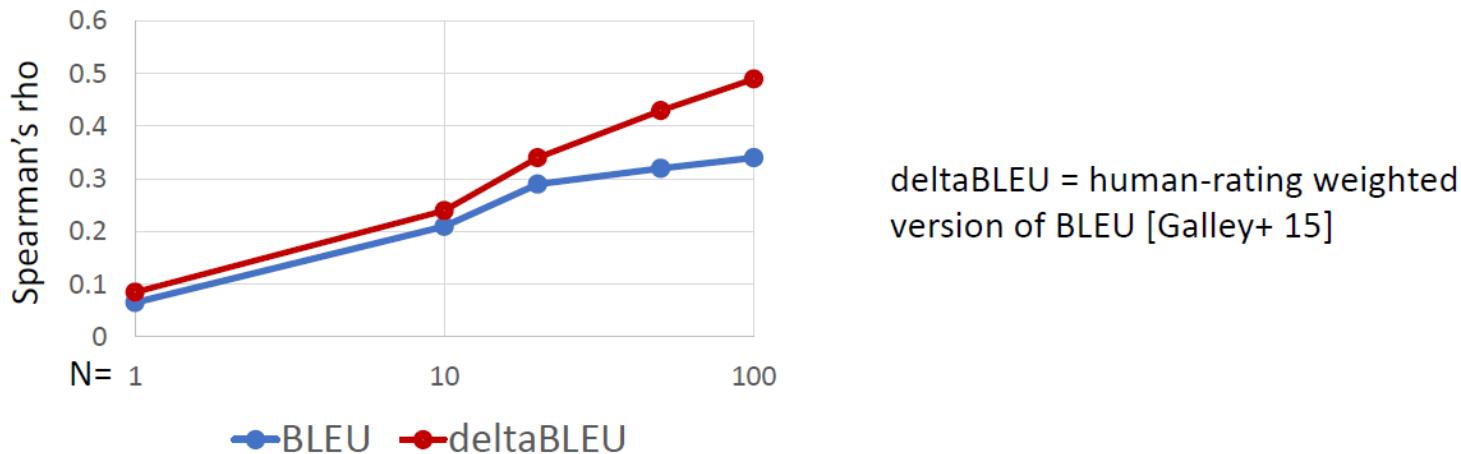
- MT metrics were NOT designed to operate at the sentence level:
 - BLEU [[Papineni+ 02](#)] == “corpus-level BLEU”
 - Statistical Significant Tests for MT [[Koehn 06](#); etc.]:
BLEU not reliable with sample size < 600,
even for Machine Translation (easier task)
- Central Limit Theorem (CLT) argument:
 - Matching against reference (e.g., n-grams) is brittle → greater variance
 - Remedy: reduce variance by increasing sample size (CLT), i.e., [corpus-level BLEU](#)



(Figure from [[Brooks+ 12](#)])

Corpus-level Correlation

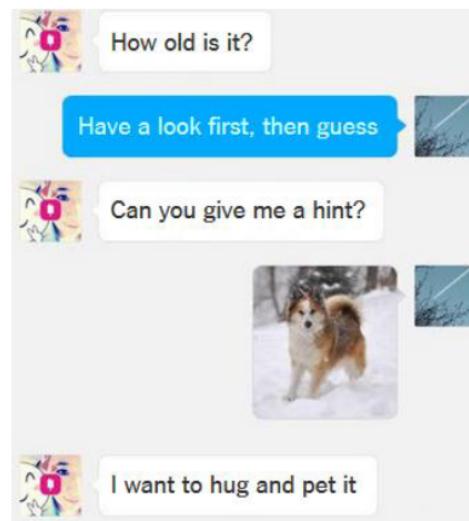
- Generally good for **Machine Translation** (MT):
 - Spearman's rho of 0.8 to 0.9 for BLEU, METEOR [[Przybocki+ 08](#)]
- Can it work for **Dialogue**?
 - Currently no definite answer, as corpus-level human judgments very expensive.
 - Experiments with smaller N [[Galley+ 15](#)]:



Social Bots: commercial systems

- For end users:

- Amazon Alexa
(trigger: say “Alexa, let’s chat”)
- Microsoft Xiaoice [[Zhou+ 2018](#)]
- Microsoft Zo (on Kik)
- Replika.ai [[system description](#)]



Xiaoice
([translated from Chinese](#))

For bot developers:

- Microsoft Personality Chat (includes [speaker embedding LSTM](#))



Replika.ai



Query:
I am going on stage :)

Personality
Professional

Top Responses & Ranked Score

What are you going to perform?	0.71
What are you going to sing?	0.69
Good luck with that.	0.67
I hope you have a great time.	0.67
Good luck with that...	0.65

Personality
Friendly

Top Responses & Ranked Score

Oh my gosh that will be amazing	0.69
Oh my gosh that's amazing. Have a great time	0.65
I'm sure you will have a great time	0.65
Oh my gosh that will be amazing. Are you going	0.65
Oh my gosh that's brilliant. Have a great time	0.65

Personality
Humorous

Top Responses & Ranked Score

With who?! D:	0.77
Really?! Oh my gosh I'm so jealous!	0.77
Oh my gosh you're so lucky! D:	0.77
That's awesome! Good luck! C:	0.77
Good luck! I'm sure you'll be amazing!	0.75

<https://labs.cognitive.microsoft.com/en-us/project-personality-chat>

Open Benchmarks

- **Alexa Challenge** (2017-)
 - Academic competition, 15 sponsored teams in 2017, 8 in 2018
 - \$250,000 research grant (2018)
 - Proceedings [[2017](#), [2018](#)]
- **Dialogue System Technology Challenge ([DSTC](#))** (2013-)
(formerly Dialogue State Tracking Challenge)
Focused this year on **grounded conversation**:
Visual-Scene [[Hori +18](#)], knowledge grounding [[Galley +18](#)]
- **Conversational Intelligence Challenge (ConvAI)** (2017-)
Last occurrence focused on **personalized chat** (Persona-Chat dataset)

Moving beyond chitchat

