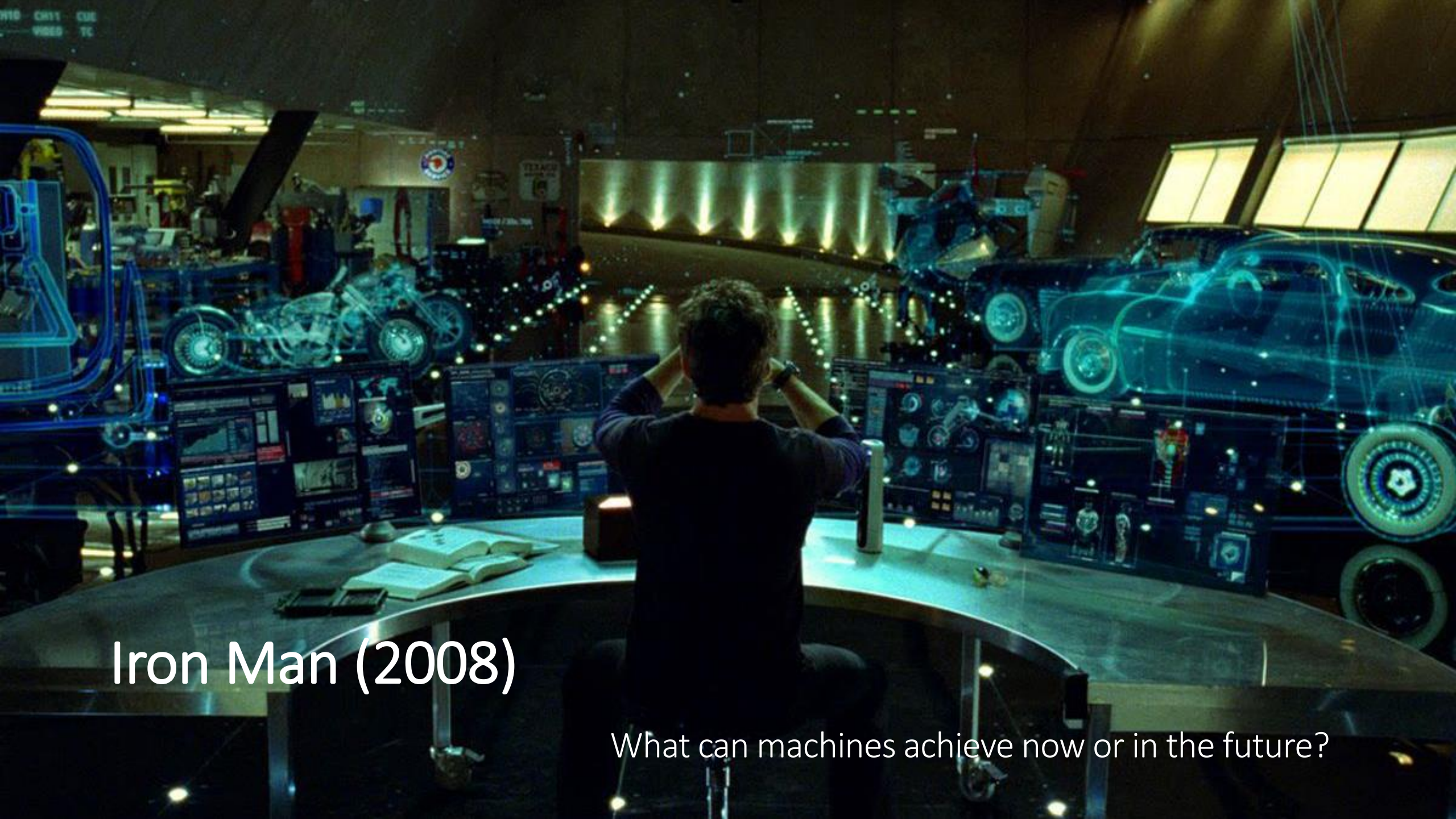


NLP APPLICATIONS III: DIALOGUE SYSTEMS



Yun-Nung Vivian Chen
<http://vivianchen.idv.tw>





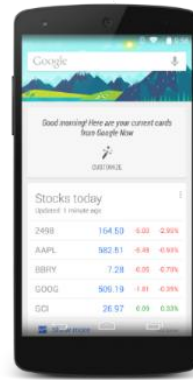
Iron Man (2008)

What can machines achieve now or in the future?

Language Empowering Intelligent Assistants



Apple Siri (2011)



Google Now (2012)
Google Assistant (2016)



Microsoft Cortana (2014)



Amazon Alexa/Echo (2014)



Google Home (2016)



Apple HomePod (2017)



Facebook Portal (2019)

Why Natural Language?

- Global Digital Statistics (2018 January)



Total Population
7.59B



Internet Users
4.02B



Active Social
Media Users
3.20B



Unique Mobile
Users
5.14B



Active Mobile
Social Users
2.96B

The more **natural** and **convenient** input of devices evolves towards **speech**.



Why and When We Need?

“I want to chat”

“I have a question”

“I need to get this done”

“What should I do?”

Turing Test (talk like a human) Social Chit-Chat

Information consumption

Task completion

Decision support

Task-Oriented Dialogues

- *What is today's agenda?*
- *What does NLP stand for?*

- *Book me the train ticket from Kaohsiung to Taipei*
- *Reserve a table at Din Tai Fung for 5 people, 7PM tonight*
- *Schedule a meeting with Vivian at 10:00 tomorrow*

- *Is this summer school good to attend?*

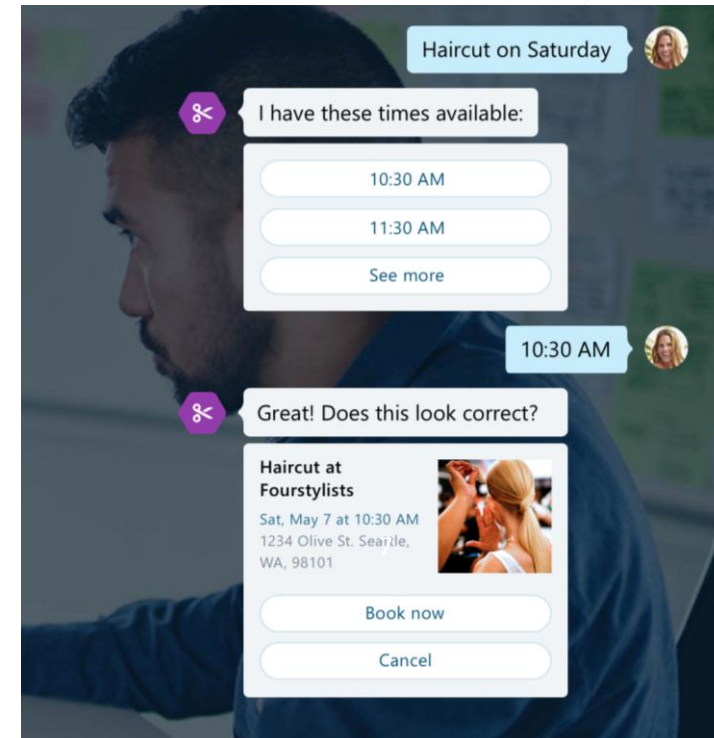
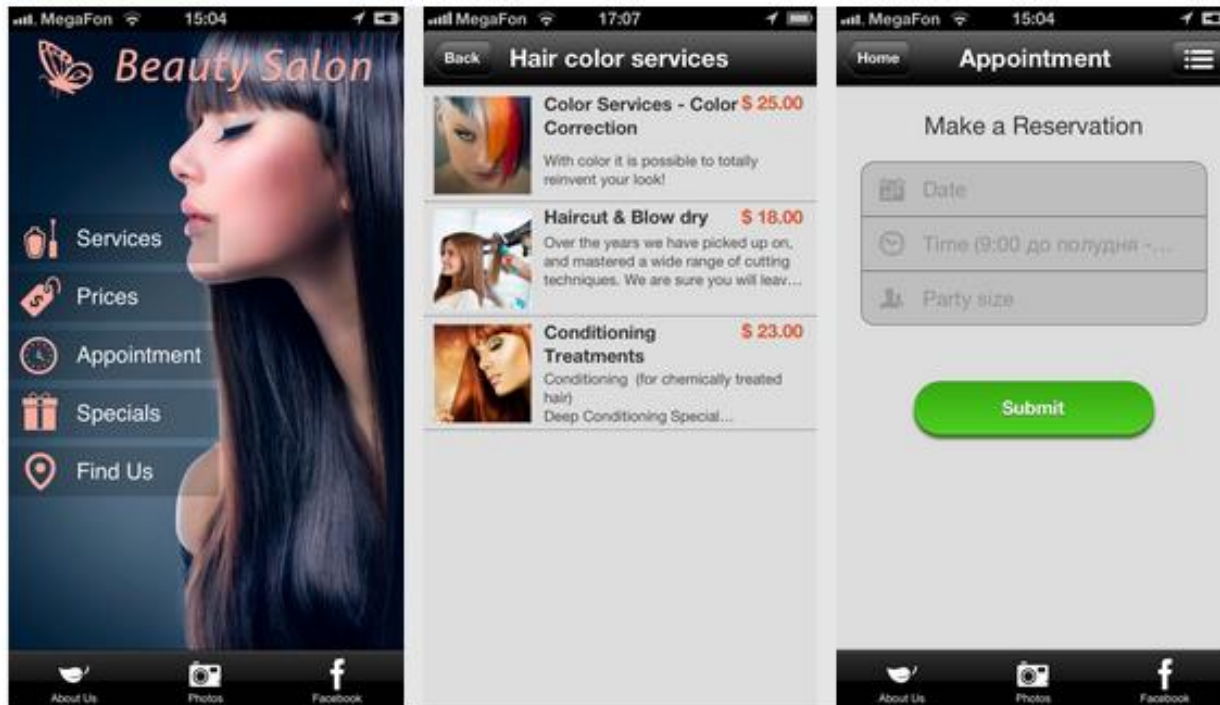
Intelligent Assistants



Task-Oriented

App → Bot

- A **bot** is responsible for a “single” domain, similar to an app



Users can initiate dialogues instead of following the GUI design



Two Branches of Conversational AI



MIULAB

NTU

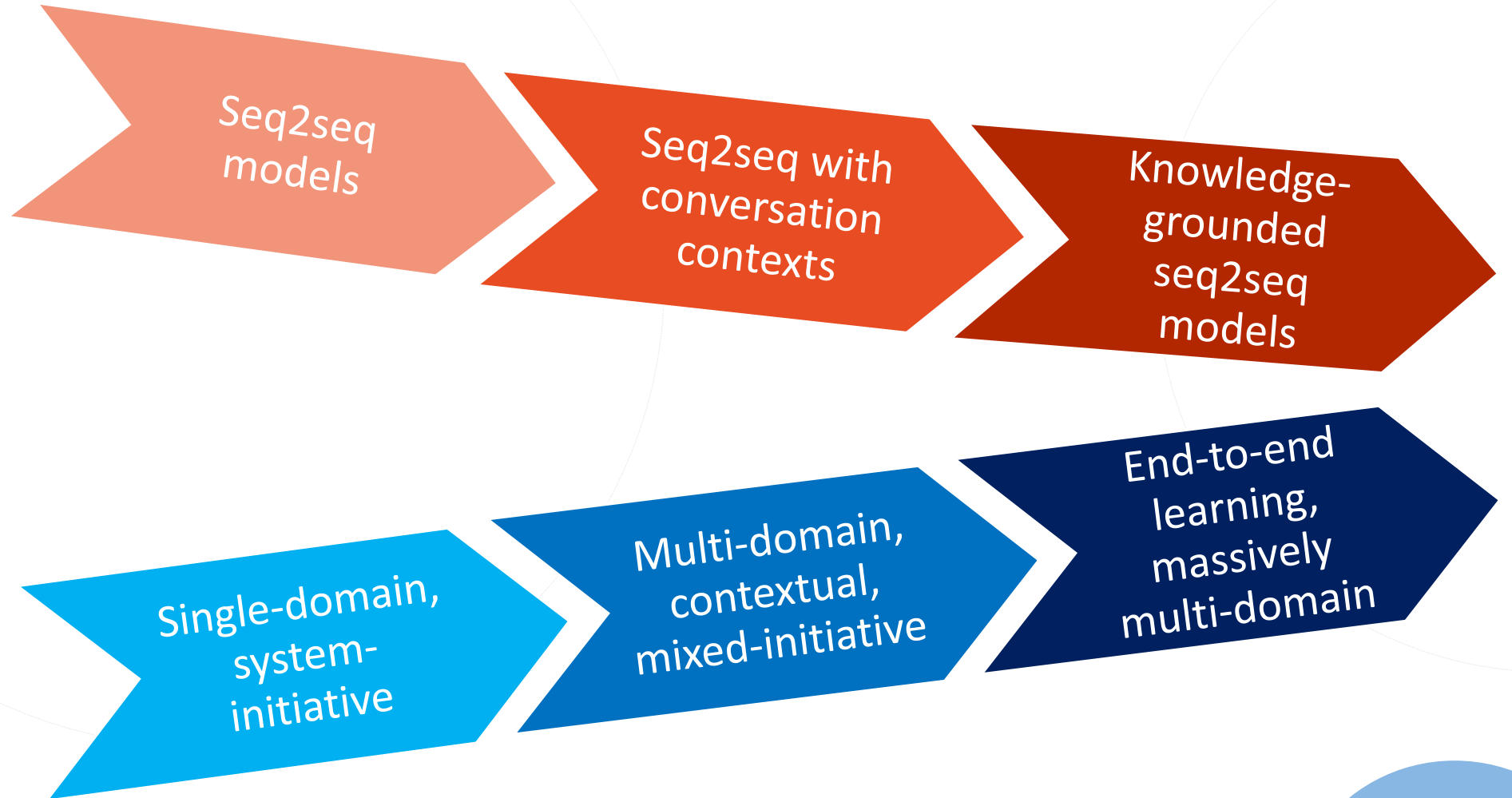
8



Chit-Chat



Task-Oriented





Task-Oriented Dialogues

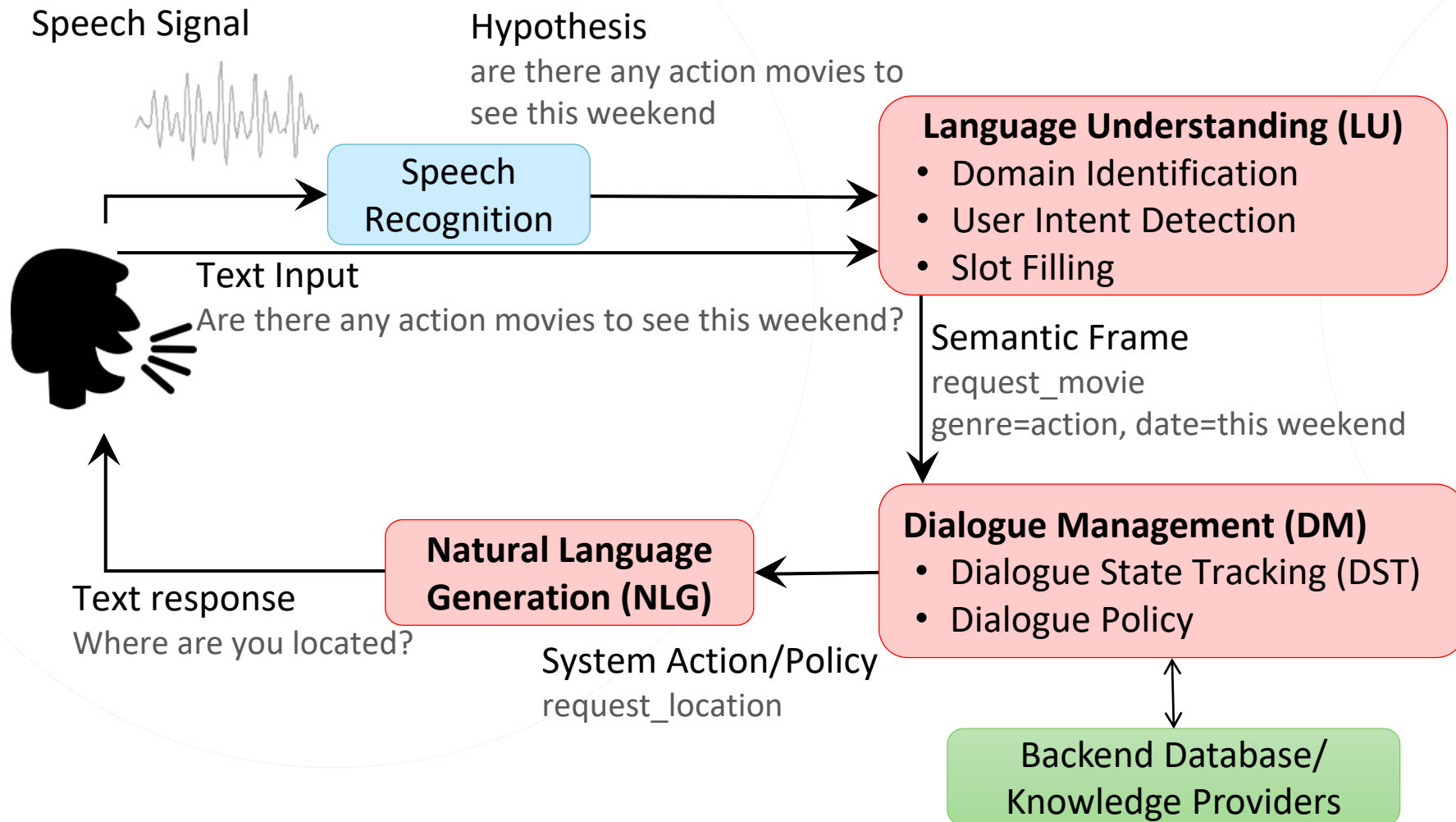


JARVIS – Iron Man's Personal Assistant

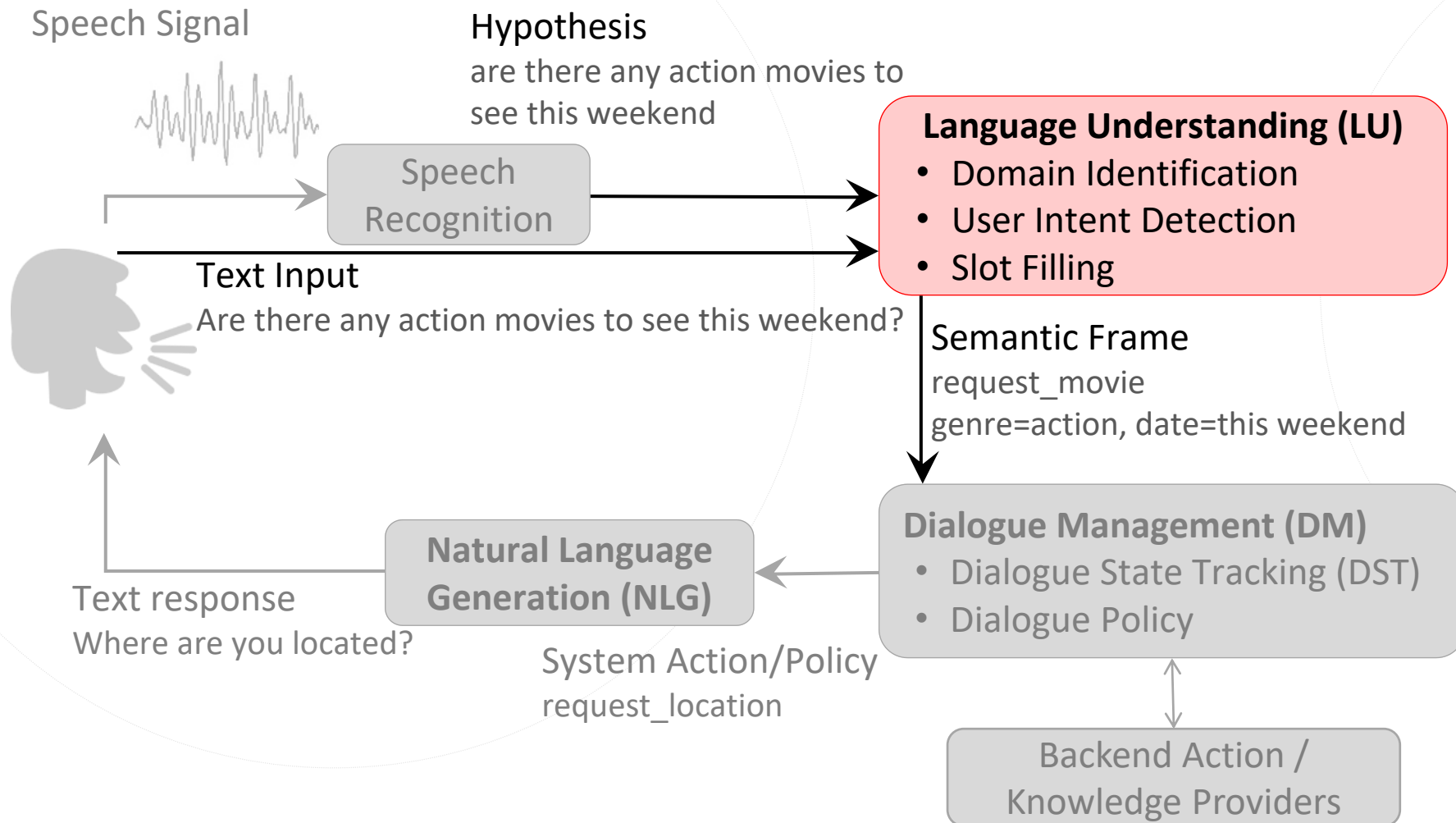


Baymax – Personal Healthcare Companion

Task-Oriented Dialogue Systems (Young, 2000)

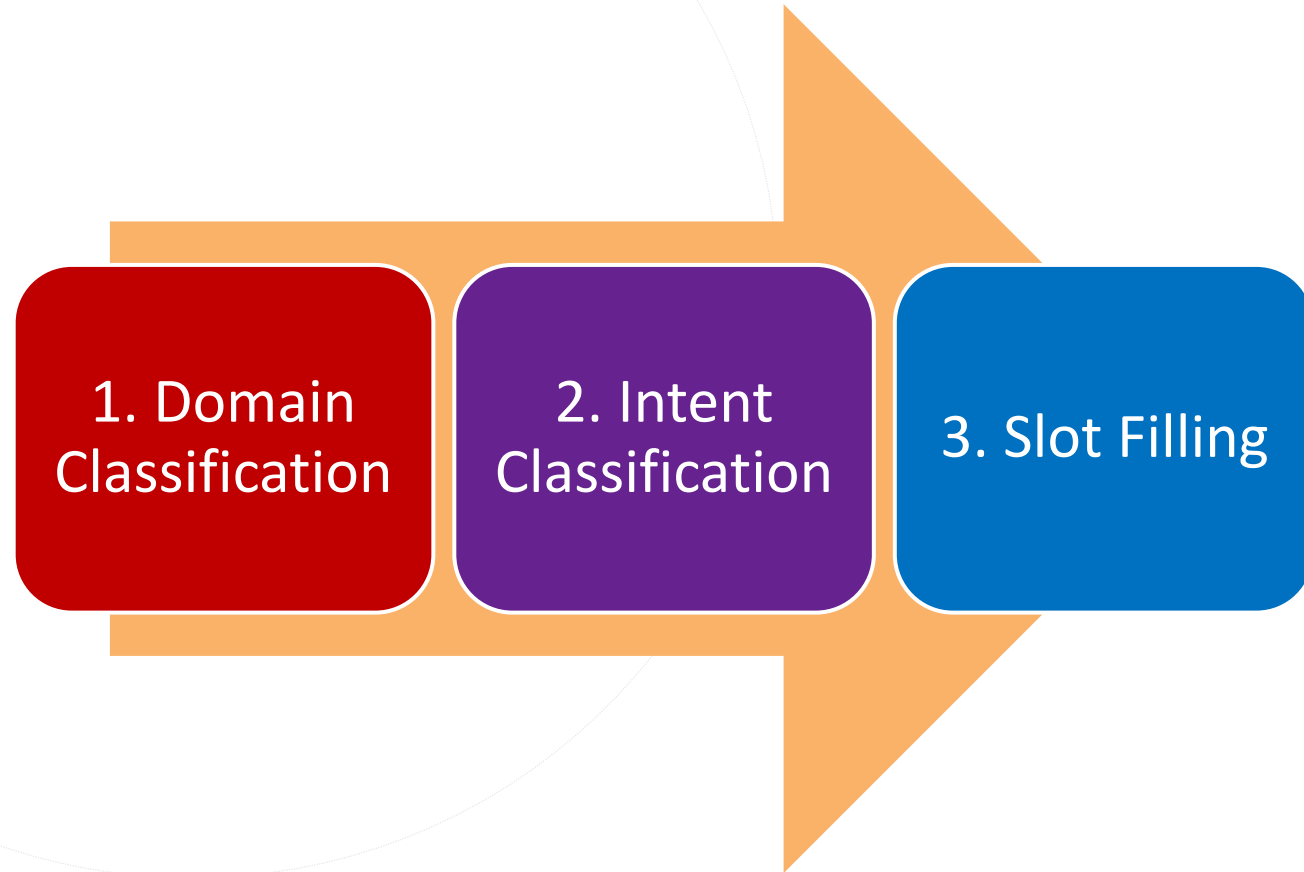


Task-Oriented Dialogue Systems (Young, 2000)



Language Understanding (LU)

- Pipelined



1. Domain Identification

Requires Predefined Domain Ontology

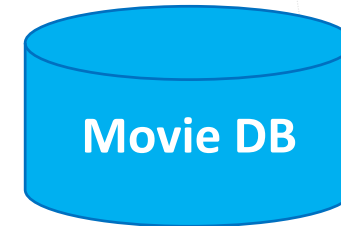
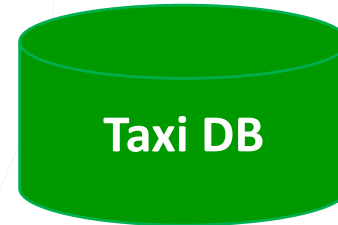
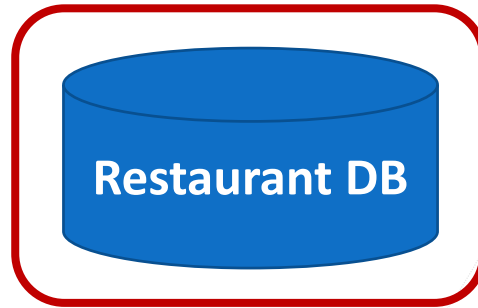
User



find a good eating place for taiwanese food



Intelligent Agent



Organized Domain Knowledge (Database)

Classification!



2. Intent Detection

Requires Predefined Schema

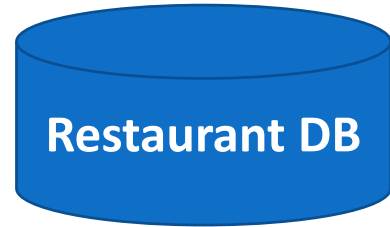
User



find a good eating place for taiwanese food



Intelligent Agent

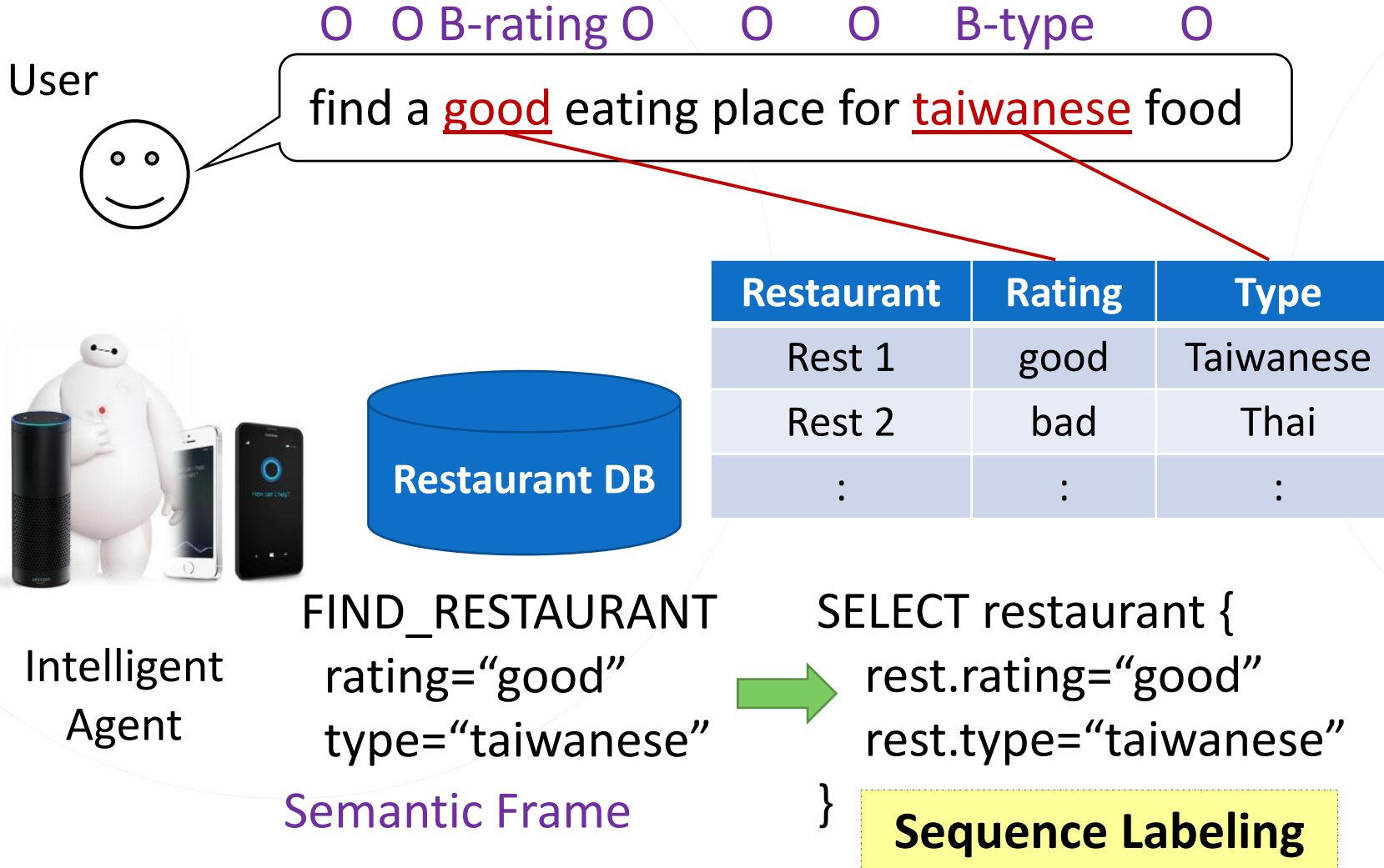


- FIND_RESTAURANT
- FIND_PRICE
- FIND_TYPE
- :

Classification!

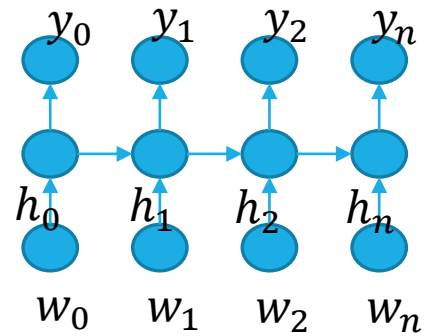
3. Slot Filling

Requires Predefined Schema

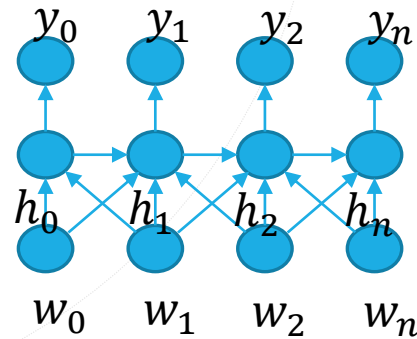


Slot Tagging (Yao et al, 2013; Mesnil et al, 2015)

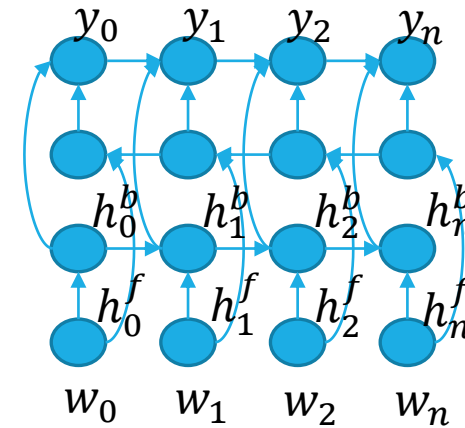
- Variations:
 - RNNs with LSTM cells
 - Input, sliding window of n-grams
 - Bi-directional LSTMs



(a) LSTM



(b) LSTM-LA



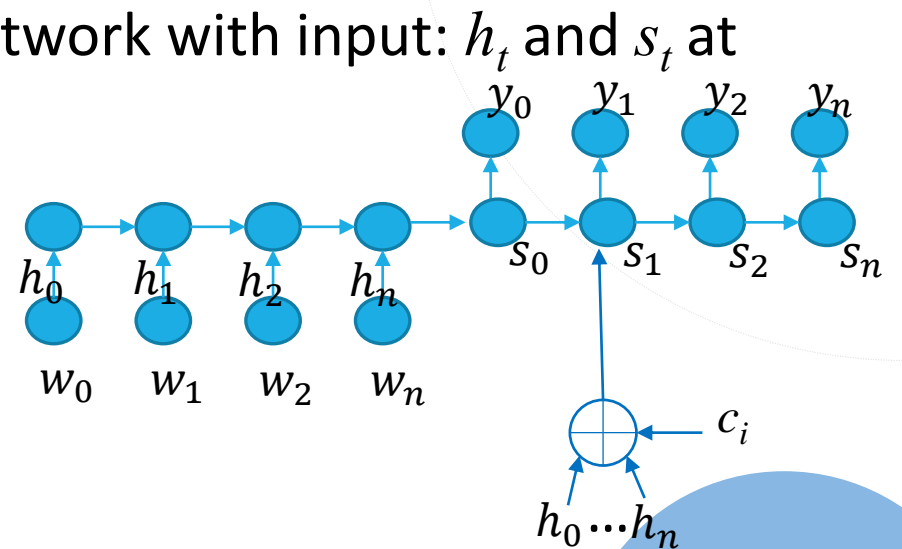
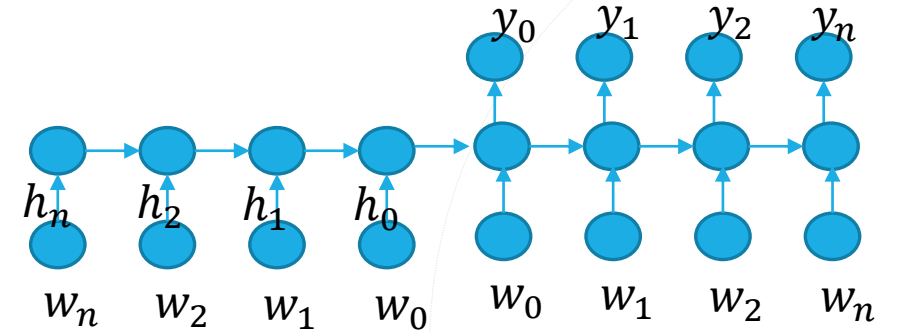
(c) bLSTM

Slot Tagging (Kurata et al., 2016; Simonnet et al., 2015)

- Encoder-decoder networks
 - Leverages sentence level information

- Attention-based encoder-decoder

- Use of attention (as in MT) in the encoder-decoder network
- Attention is estimated using a feed-forward network with input: h_t and s_t at time t



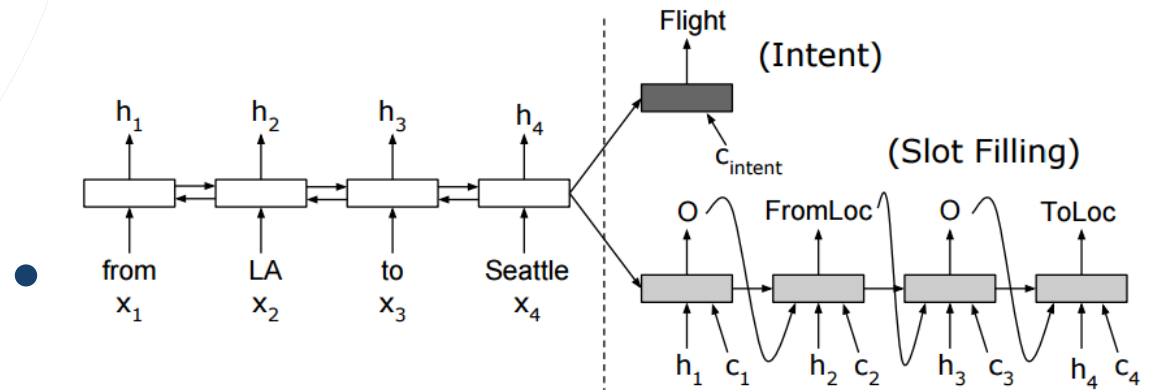
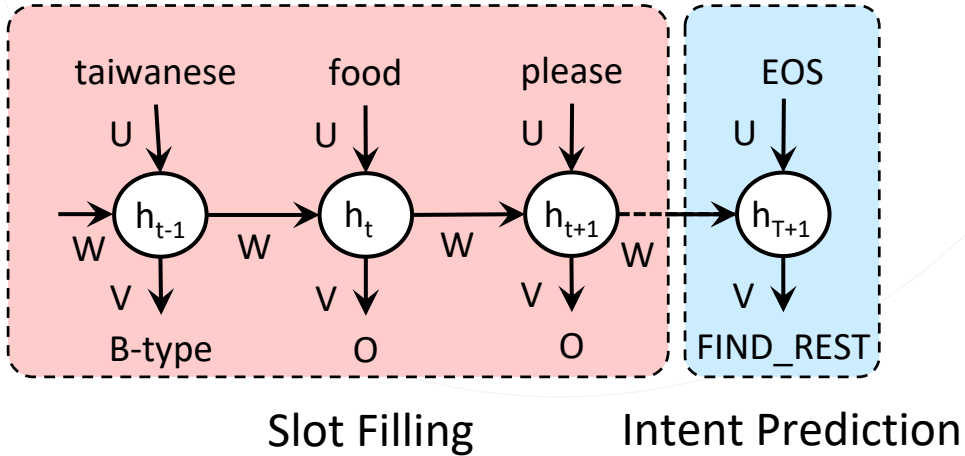
Joint Semantic Frame Parsing

Sequence-based (Hakkani-Tur et al., 2016)

- Slot filling and intent prediction in the same output sequence

Parallel (Liu and Lane, 2016)

- Intent prediction and slot filling are performed in two branches

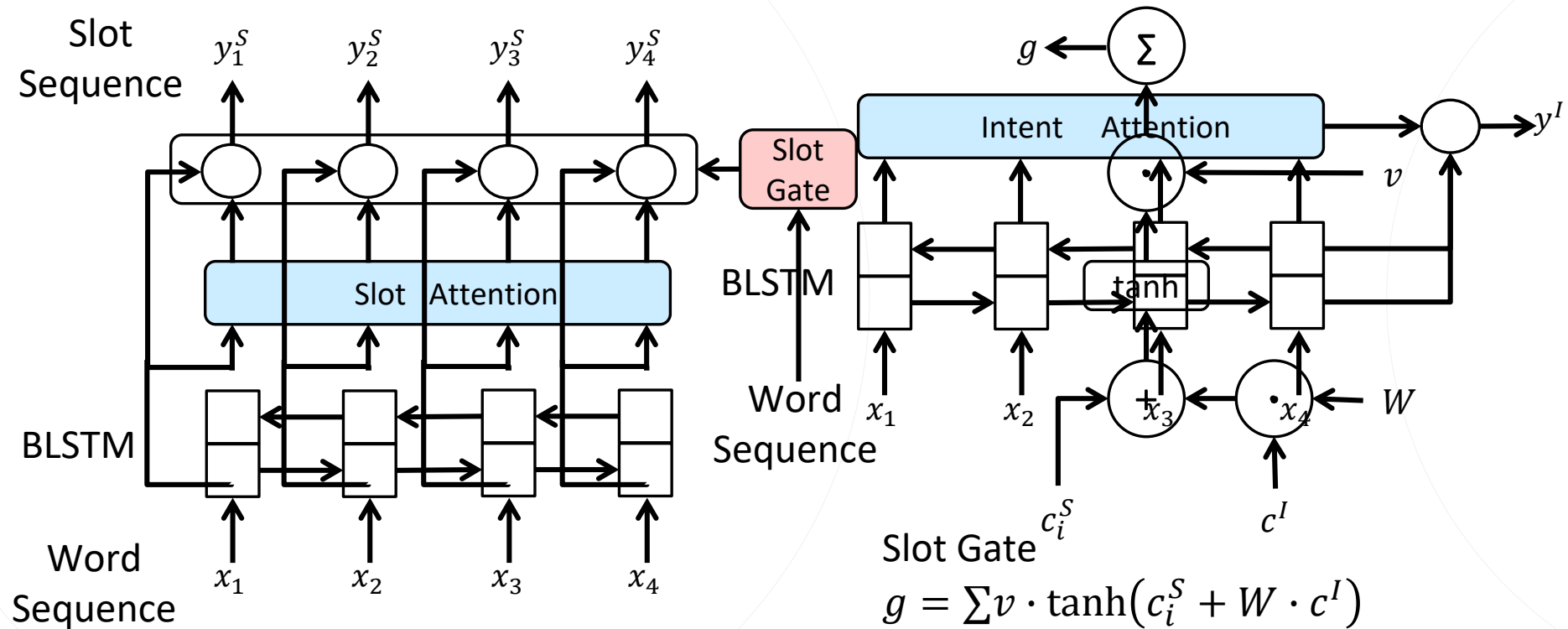


Joint Model Comparison

	Attention Mechanism	Intent-Slot Relationship
Joint bi-LSTM	X	Δ (Implicit)
Attentional Encoder-Decoder	\checkmark	Δ (Implicit)
Slot Gate Joint Model	\checkmark	\checkmark (Explicit)



Slot-Gated Joint SLU (Goo+, 2018)



Slot Gate

$$g = \sum v \cdot \tanh(c_i^S + W \cdot c^I)$$

Slot Prediction

$$y_i^S = \text{softmax}(W^S(h_i + g \cdot c_i^S) + b^S)$$

g will be larger if slot and intent are better related

Contextual Language Understanding

- User utterances are highly ambiguous in isolation

Restaurant
Booking

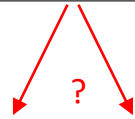


Book a table for 10 people tonight.

Which restaurant would you like to book a table for?



Cascal, for 6.



#people time



End-to-End Memory Networks (Sukhbaatar et al, 2015)

U: "i d like to purchase tickets to see deepwater horizon"

S: "for which theatre"

U: "angelika"

S: "you want them for angelika theatre?"

U: "yes angelika"

S: "how many tickets would you like ?"

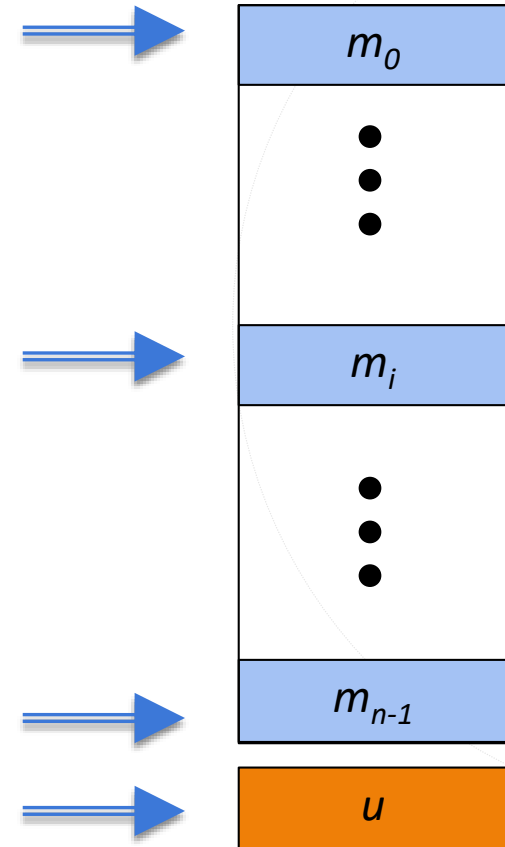
U: "3 tickets for saturday"

S: "What time would you like ?"

U: "Any time on saturday is fine"

S: "okay , there is 4:10 pm , 5:40 pm and 9:20 pm"

U: "Let's do 5:40"



E2E MemNN for Contextual LU (Chen+, 2016)

1. Sentence Encoding

$$m_i = \text{RNN}_{\text{mem}}(x_i)$$

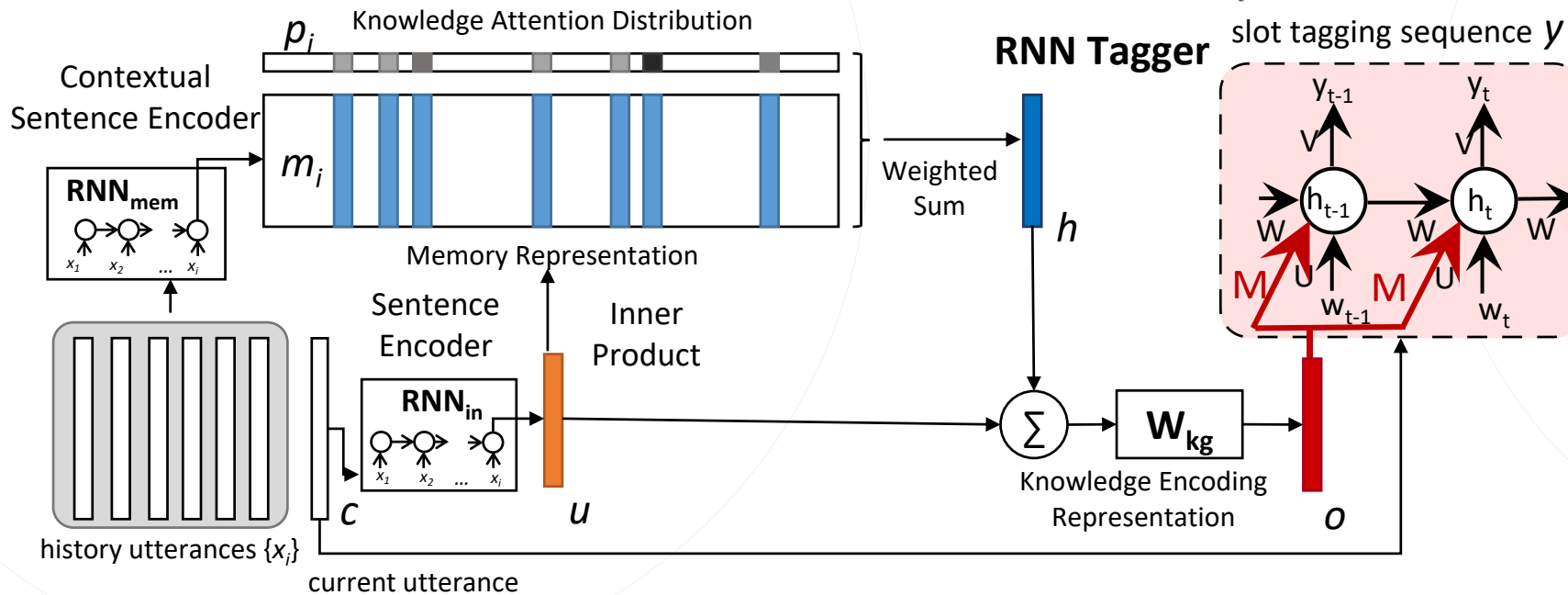
$$u = \text{RNN}_{\text{in}}(c)$$

2. Knowledge Attention

$$p_i = \text{softmax}(u^T m_i)$$

3. Knowledge Encoding

$$h = \sum_i p_i m_i \quad o = W_{\text{kg}}(h + u)$$



Idea: additionally incorporating contextual knowledge during slot tagging
 → track dialogue states in a latent way

E2E MemNN for Contextual LU ([Chen et al., 2016](#))

U: "i d like to purchase tickets to see deepwater horizon" → 0.69

S: "for which theatre"

U: "angelika"

S: "you want them for angelika theatre?"

U: "yes angelika"

S: "how many tickets would you like ?" → 0.13

U: "3 tickets for saturday"

S: "What time would you like ?"

U: "Any time on saturday is fine"

S: "okay , there is 4:10 pm , 5:40 pm and 9:20 pm" → 0.16

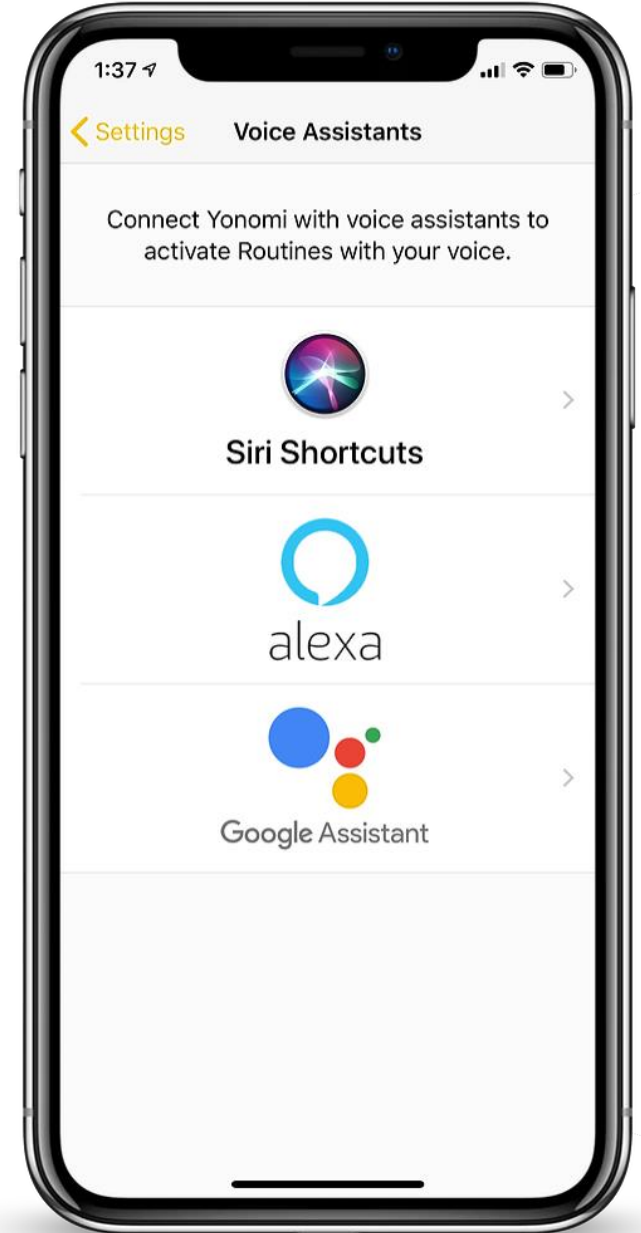
U: "Let's do 5:40"





Recent Advances in NLP

- Contextual Embeddings (ELMo & BERT)
 - Boost many understanding performance with pre-trained natural language



Call me ASAP

! ?
1 2 3 4 5 6 7 8
q w e r t y u i o p
a s d f g h j k l
- _ ↑ ↓
7 1 2 3 , .



SAMSUNG

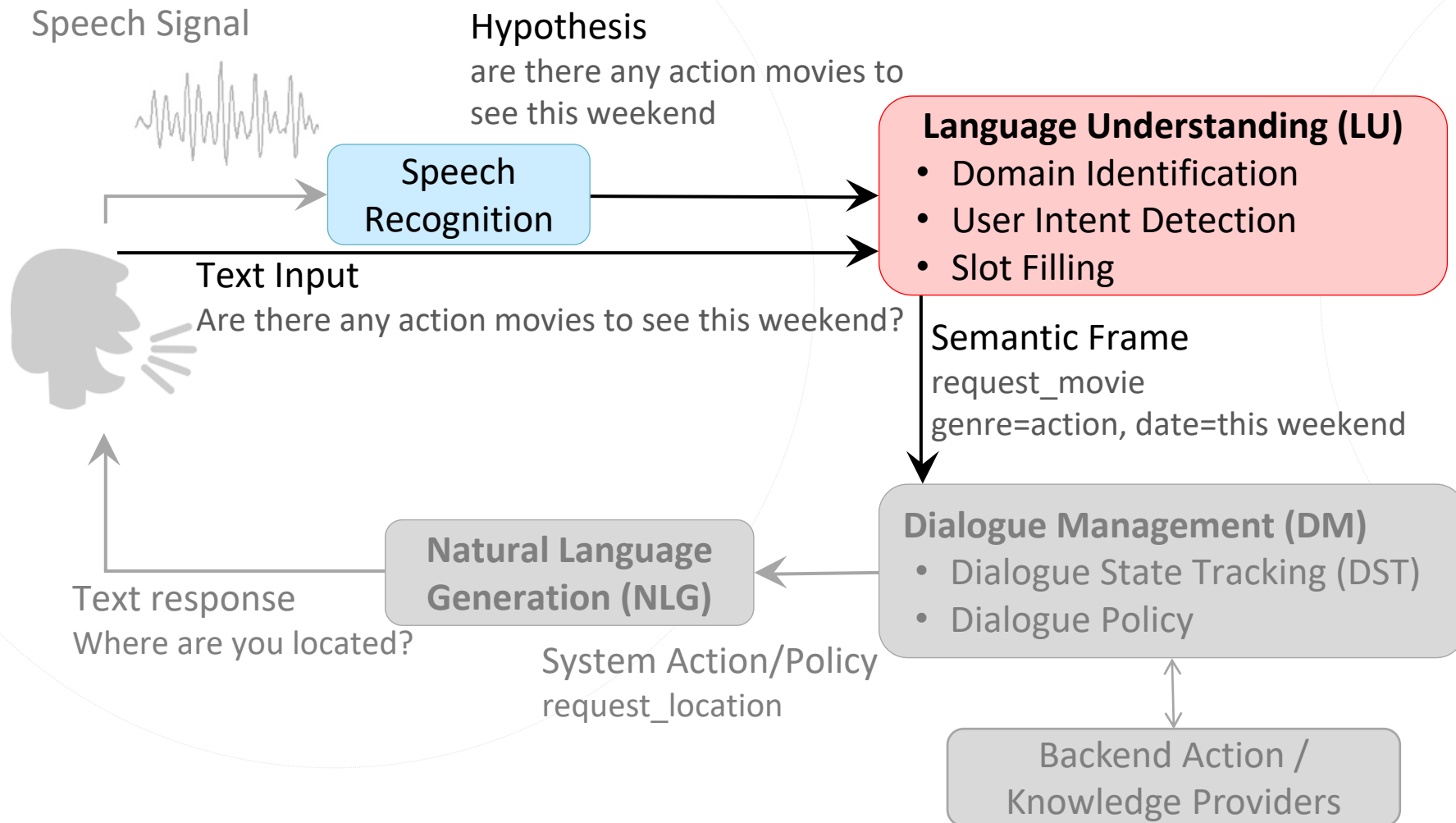
6:00



Listening...

~~Lift all lights to Morocco~~
List all flights tomorrow

Task-Oriented Dialogue Systems (Young, 2000)




Mismatch between Written and Spoken Languages




Training

- Written language

A stack of five colorful books (red, green, purple, blue, yellow) stacked on top of each other.

Testing

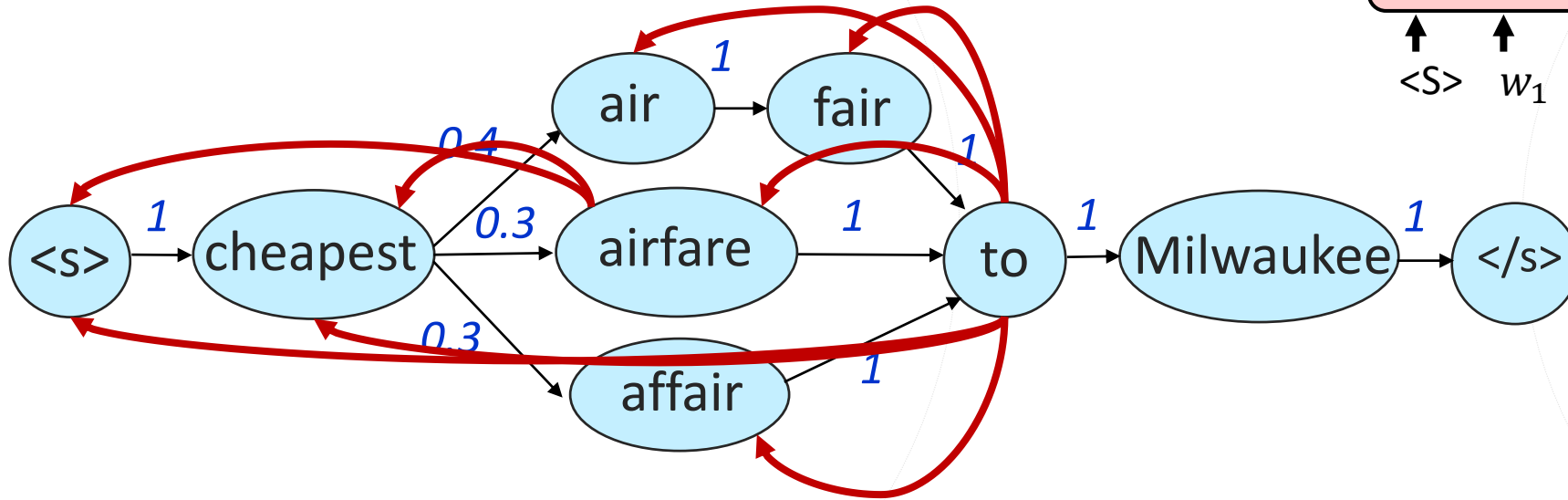
- Spoken language
- Include recognition errors

A young child with short brown hair is shown in profile, looking to the right and speaking. Several letters (O, Y, L, A, F, P, a, C, M, H, E, T, D) are floating around the child's head, representing spoken language.

- Goal: ASR-Robust Contextualized Embeddings
 - ✓ learning contextualized word embeddings specifically for spoken language
 - ✓ achieves better performance on *spoken* language understanding tasks

Adapting Transformer to ASR Lattices (Huang and Chen, 2019)

- Idea: lattices may include correct words
- Goal: feed lattices into Transformer



$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}} + M\right)V$$

ASR-Robust Contextualized Embeddings



- Confusion-Aware Fine-Tuning

- Supervised

Acoustic Confusion $C = \{w_3^{x_{trs}}, w_2^{x_{asr}}\}$

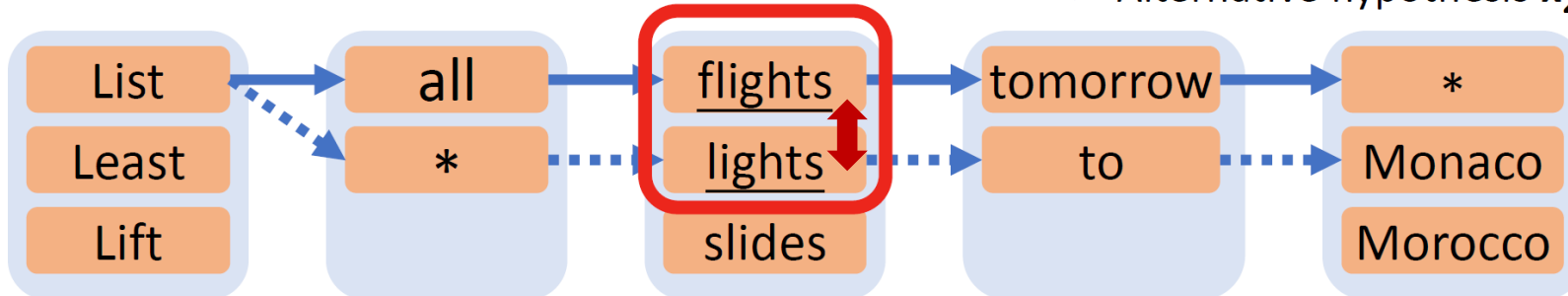
x_{trs} : Show me the fares from Dallas to Boston

x_{asr} : Show me * affairs from Dallas to Boston

- Unsupervised

Acoustic Confusion

→ Top hypothesis x_1
 Alternative hypothesis x_2



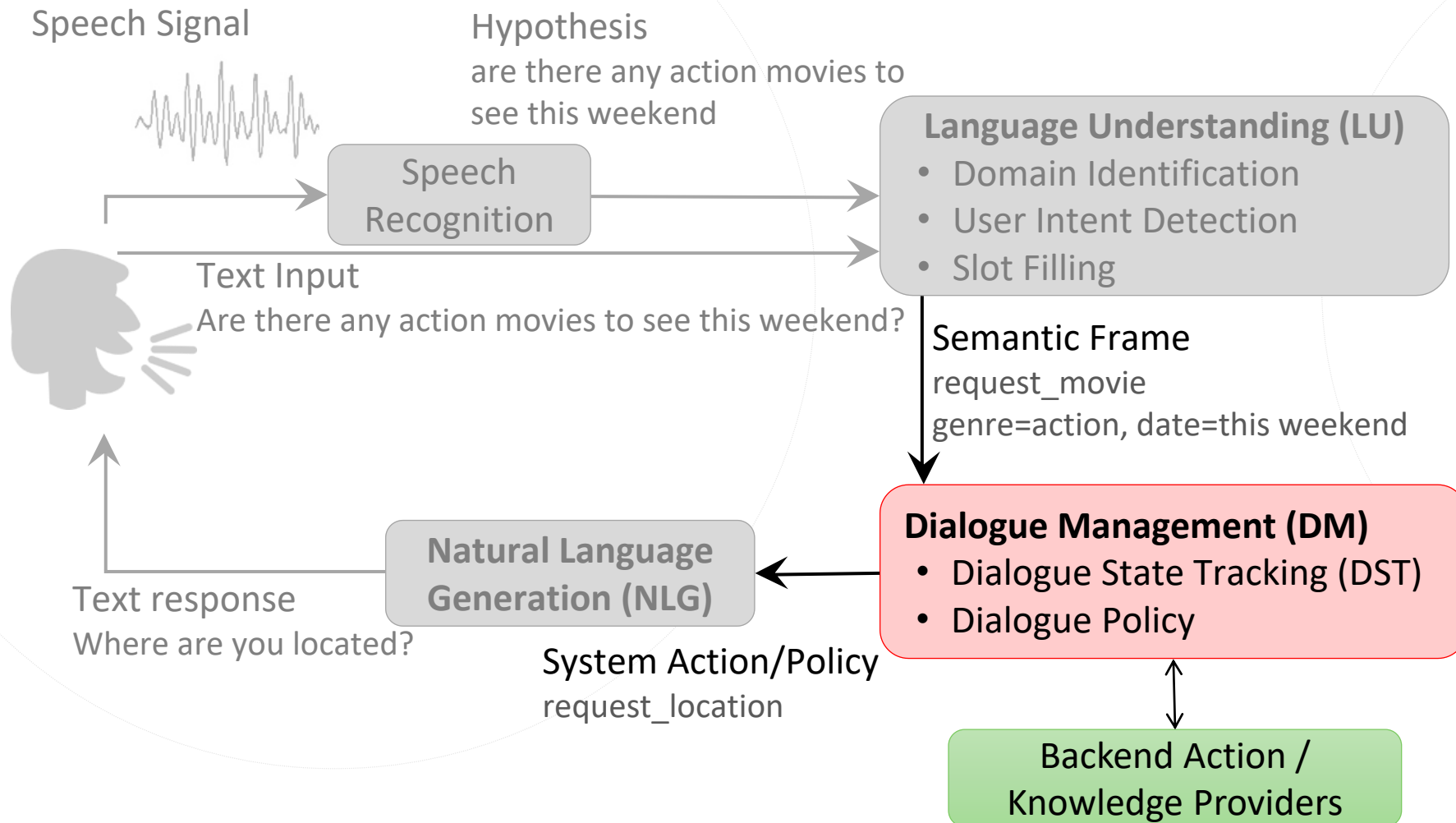


LU Evaluation

- Metrics

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

Task-Oriented Dialogue Systems (Young, 2000)



Dialogue State Tracking



Dialogue State Tracking

Requires Hand-Crafted States

User

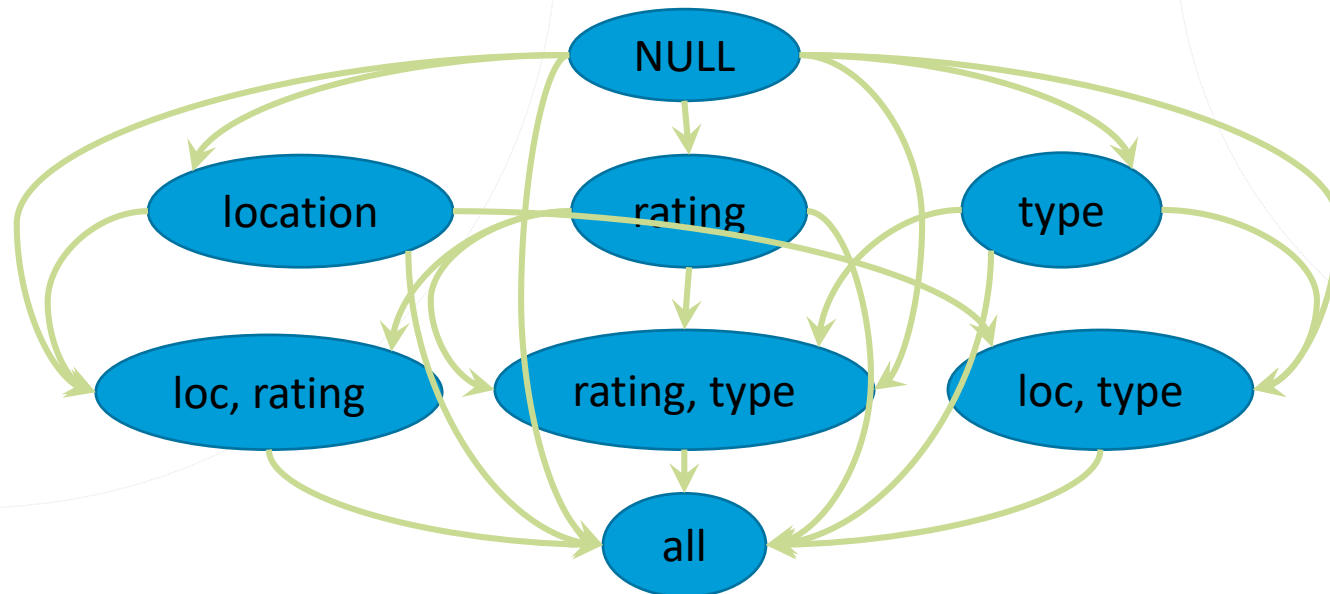


find a good eating place for taiwanese food

i want it near to my office



Intelligent Agent



Dialogue State Tracking

Requires Hand-Crafted States

User

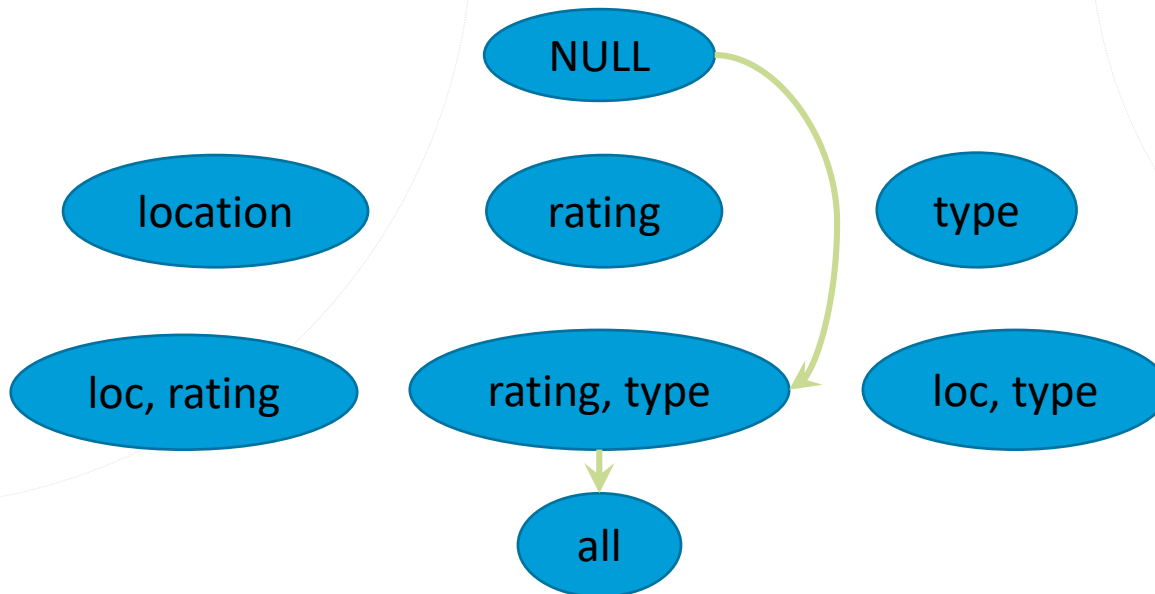


find a good eating place for taiwanese food

i want it near to my office



Intelligent Agent



Dialogue State Tracking

Handling Errors and Confidence

User



find a good eating place for taixxxx food

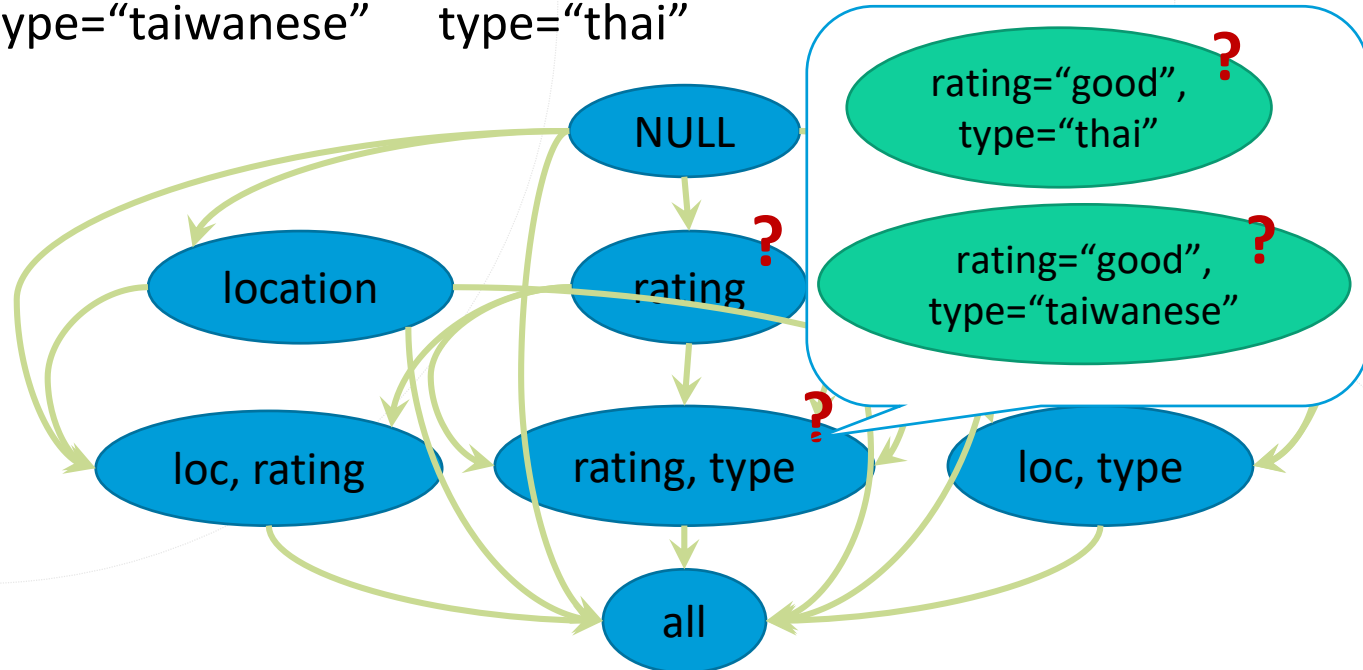
FIND_RESTAURANT
rating="good"
type="taiwanese"

FIND_RESTAURANT
rating="good"
type="thai"

FIND_RESTAURANT
rating="good"



Intelligent Agent





Dialogue State Tracking (DST)

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

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

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





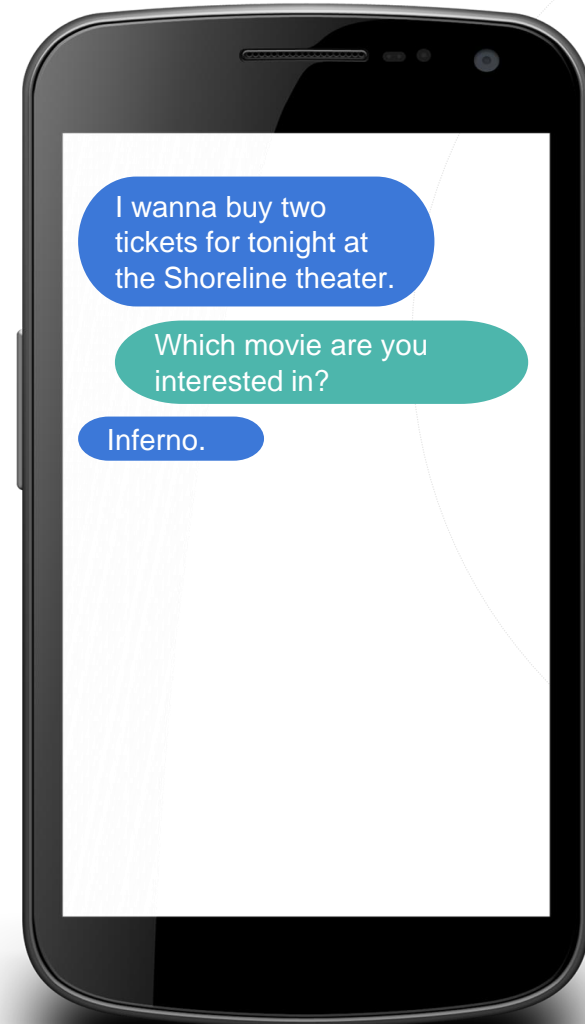
Multi-Domain Dialogue State Tracking

- 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	11/15/17			
Time	6 pm	7 pm	8 pm	9 pm
#People	2			
Theater	Century 16 Shoreline			
Movie	Inferno			

Less Likely More Likely





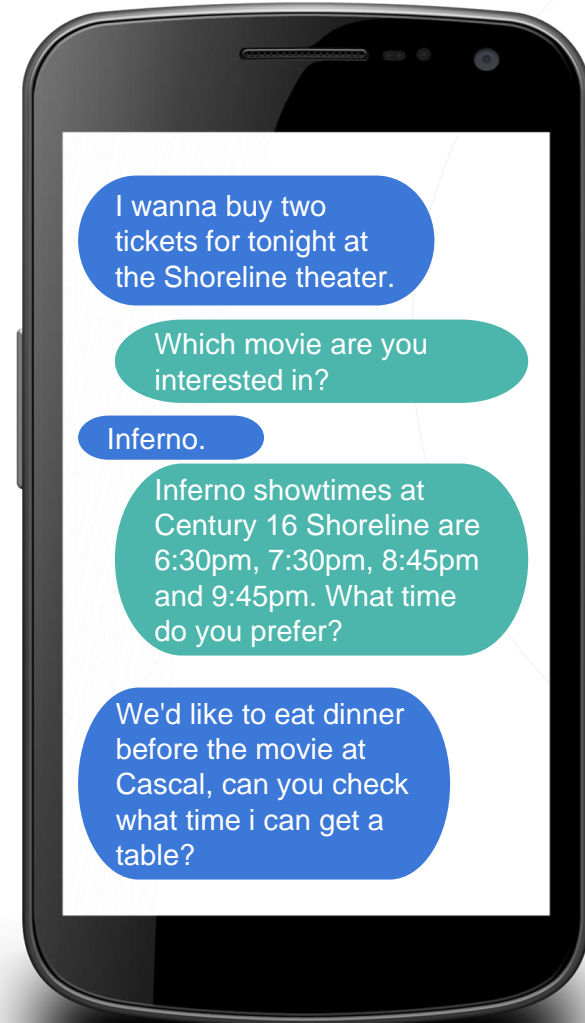
Multi-Domain Dialogue State Tracking

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

Date	11/15/17			
Time	6:30 pm	7:30 pm	8:45 pm	9:45 pm
#People	2			
Theater	Century 16 Shoreline			
Movie	Inferno			

Date	11/15/17		
Time	6:00 pm	6:30 pm	7:00 pm
Restaurant	Cascal		
#People	2		

Less Likely More Likely

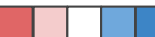


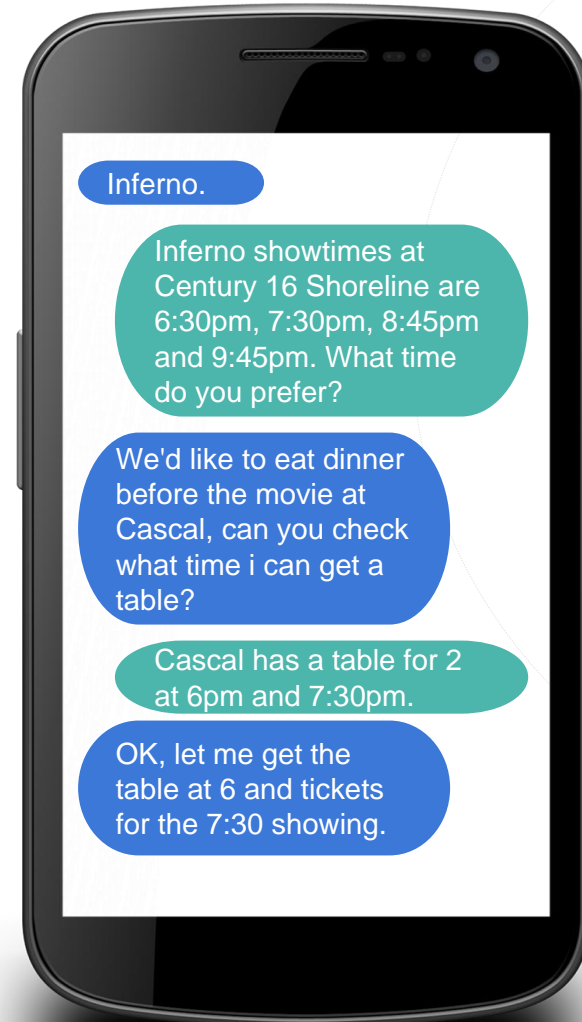
Multi-Domain Dialogue State Tracking

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

	11/15/17			
Date	11/15/17			
Time	6:30 pm	7:30 pm	8:45 pm	9:45 pm
#People	2			
Theater	Century 16 Shoreline			
Movie	Inferno			

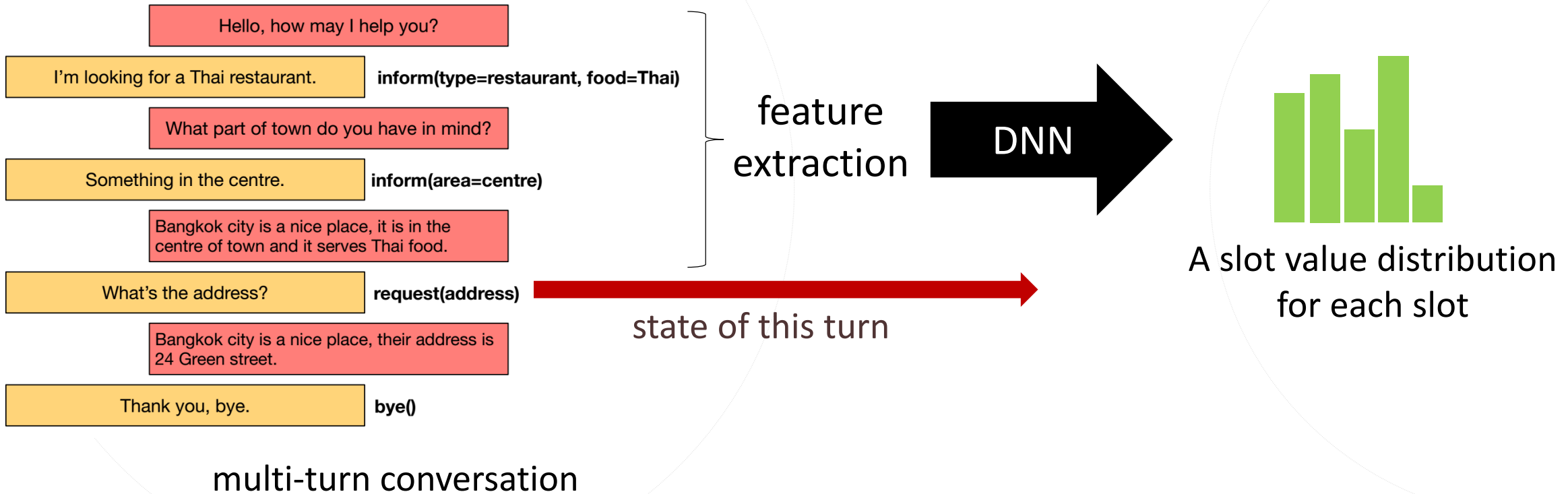
	11/15/17		
Date	11/15/17		
Time	6:00 pm	6:30 pm	7:00 pm
Restaurant	Cascal		
#People	2		

Less Likely  More Likely

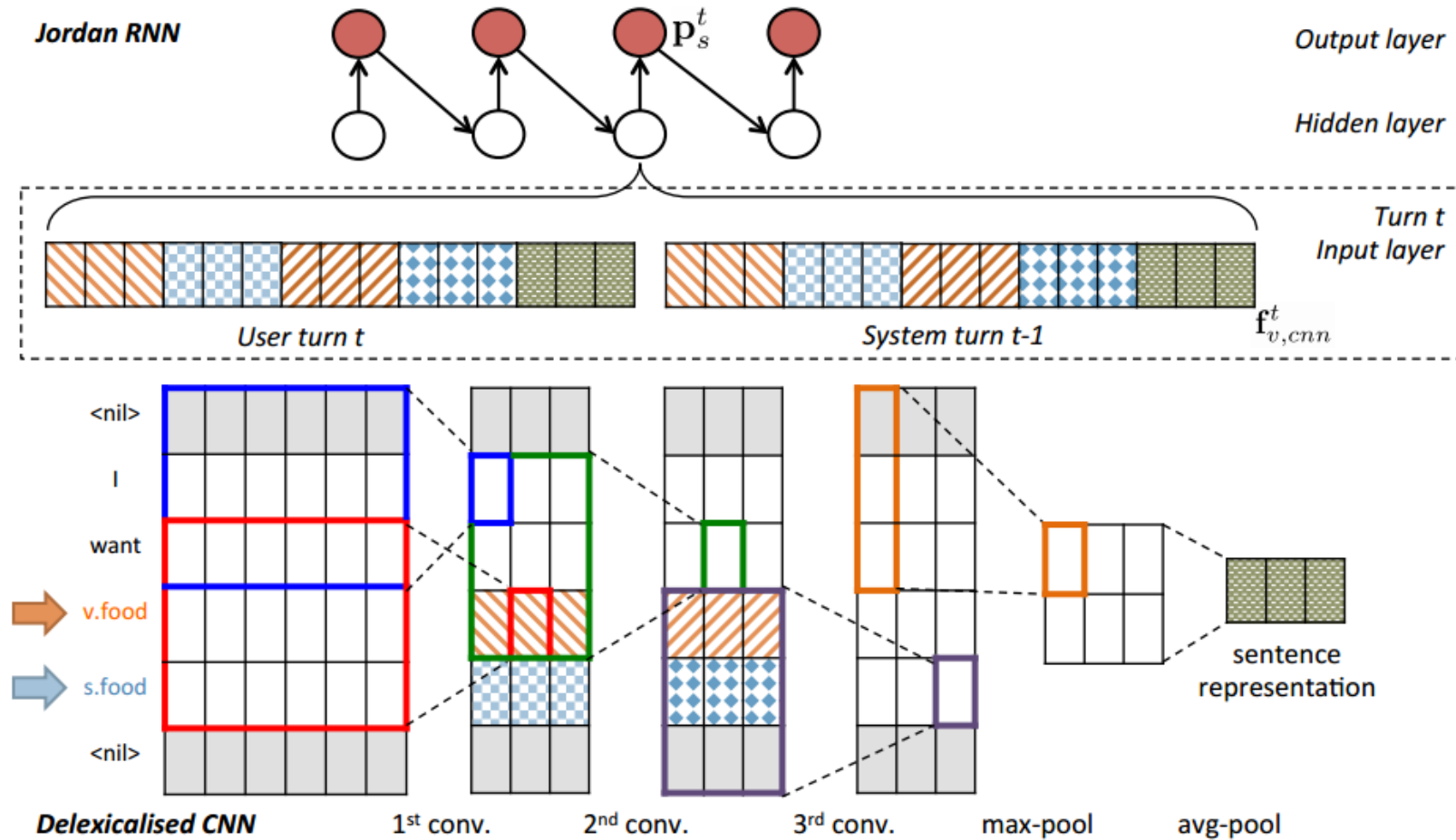




DNN for DST



RNN-CNN DST (Mrkšić+, 2015)

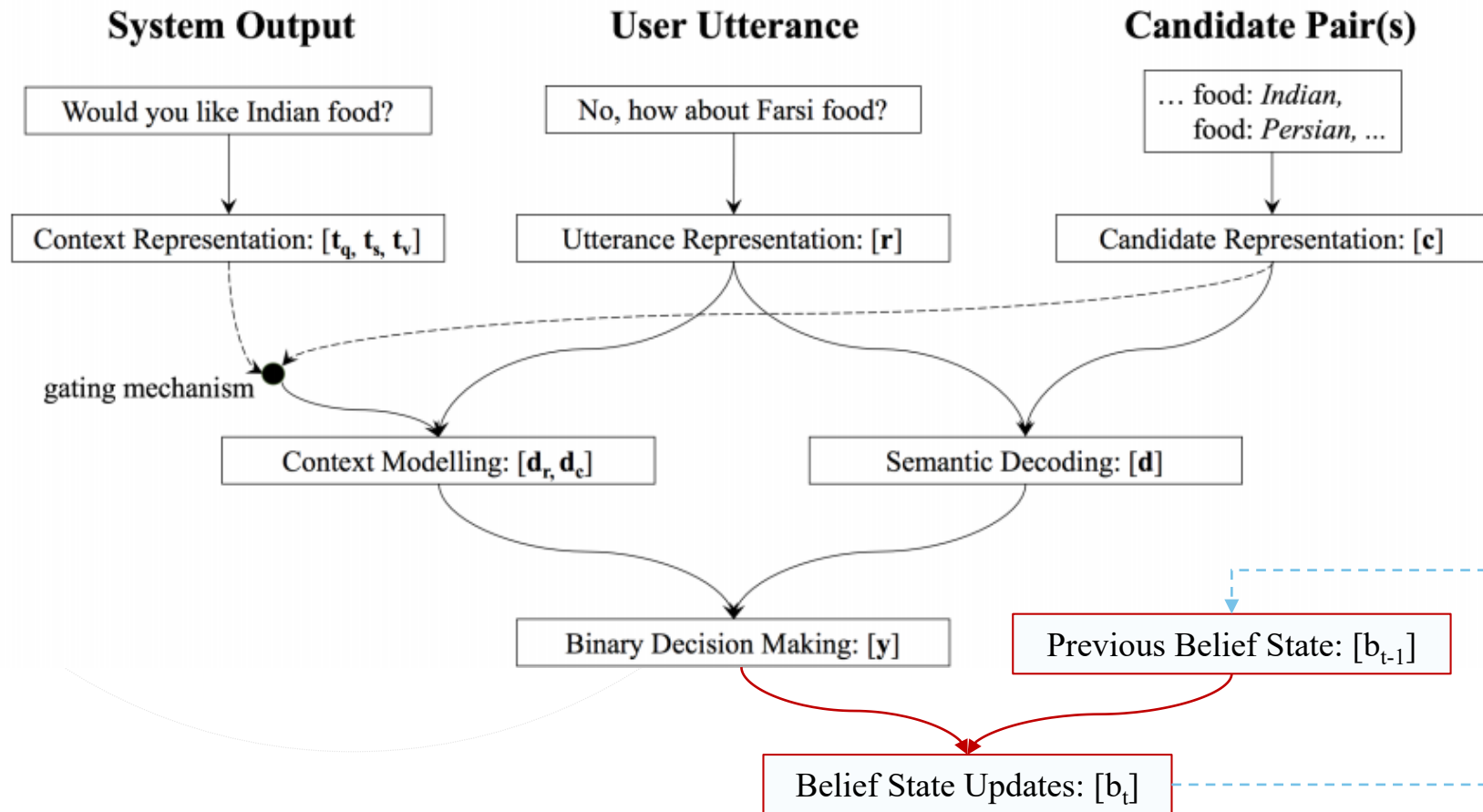


(Figure from Wen et al, 2016)



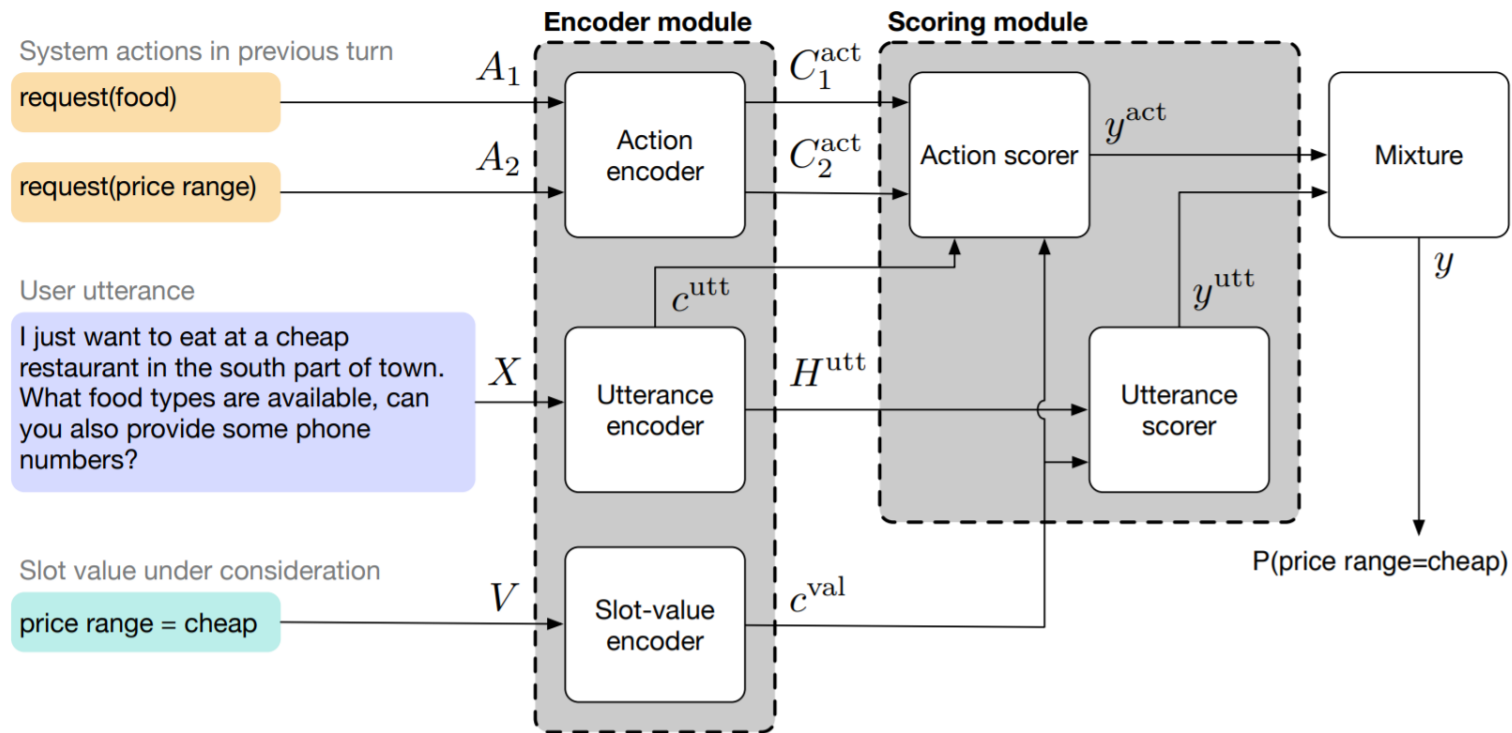
Neural Belief Tracker (Mrkšić+, 2016)

- Candidate pairs are considered



Global-Locally Self-Attentive DST (Zhong+, 2018)

- More advanced encoder
 - Global modules share parameters for all slots
 - Local modules learn slot-specific feature representations



Dialog State Tracking Challenge (DSTC)

(Williams et al. 2013, Henderson et al. 2014, Henderson et al. 2014, Kim et al. 2016, Kim et al. 2016)



Challenge	Type	Domain	Data Provider	Main Theme
<u>DSTC1</u>	Human-Machine	Bus Route	CMU	Evaluation Metrics
<u>DSTC2</u>	Human-Machine	Restaurant	U. Cambridge	User Goal Changes
<u>DSTC3</u>	Human-Machine	Tourist Information	U. Cambridge	Domain Adaptation
<u>DSTC4</u>	Human-Human	Tourist Information	I2R	Human Conversation
<u>DSTC5</u>	Human-Human	Tourist Information	I2R	Language Adaptation



DSTC4-5

- Type: Human-Human
- Domain: Tourist Information

{Topic: Accommodation; NAME: InnCrowd Backpackers Hostel; GuideAct: REC; TouristAct: ACK}

Guide: Let's try this one, okay?

Tourist: Okay.

Guide: It's InnCrowd Backpackers Hostel in Singapore. If you take a dorm bed per person only twenty dollars. If you take a room, it's two single beds at fifty nine dollars.

Tourist: Um. Wow, that's good.

Guide: Yah, the prices are based on per person per bed or dorm. But this one is room. So it should be fifty nine for the two room. So you're actually paying about ten dollars more per person only.

Tourist: Oh okay. That's- the price is reasonable actually. It's good.

{Topic: Accommodation; Type: Hostel; Pricerange: Cheap; GuideAct: ACK; TouristAct: REQ}

Tourist: Can you give me some uh- tell me some cheap rate hotels, because I'm planning just to leave my bags there and go somewhere take some pictures.

Guide: Okay. I'm going to recommend firstly you want to have a backpack type of hotel, right?

Tourist: Yes. I'm just gonna bring my backpack and my buddy with me. So I'm kinda looking for a hotel that is not that expensive. Just gonna leave our things there and, you know, stay out the whole day.

Guide: Okay. Let me get you hm hm. So you don't mind if it's a bit uh not so roomy like hotel because you just back to sleep.

Tourist: Yes. Yes. As we just gonna put our things there and then go out to take some pictures.

Guide: Okay, um-

Tourist: Hm.

DST Evaluation

- Metric
 - Tracked state accuracy with respect to user goal
 - Recall/Precision/F-measure individual slots



Dialogue Policy Optimization



Hello, how may I help you?

greeting ()

I'm looking for a Thai restaurant.

request (restaurant; foodtype=Thai)

What part of town do you have in mind?

request (area)

Something in the centre.

inform (area=centre)

Bangkok city is a nice place, it is in the centre of town and it serves Thai food.

inform (restaurant=Bangkok city, area=centre of town, foodtype=Thai)

What's the address?

request (address)

Bangkok city is a nice place, their address is 24 Green street.

inform (address=24 Green street)

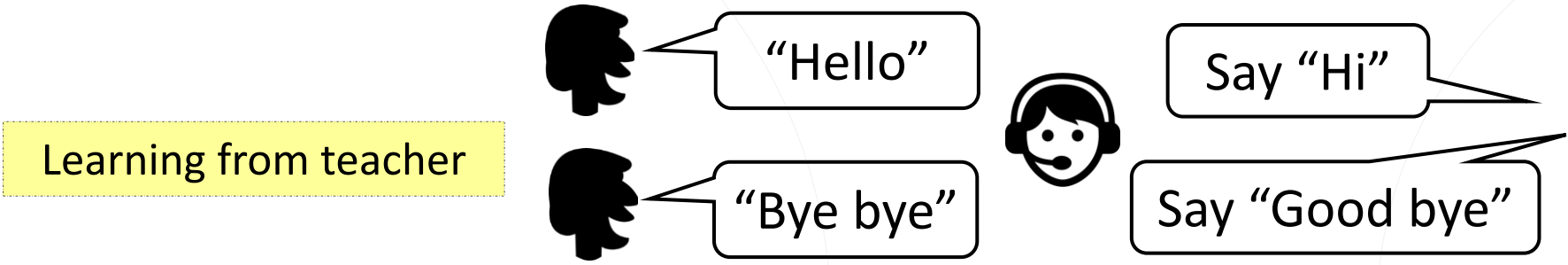
Thank you, bye.

bye ()

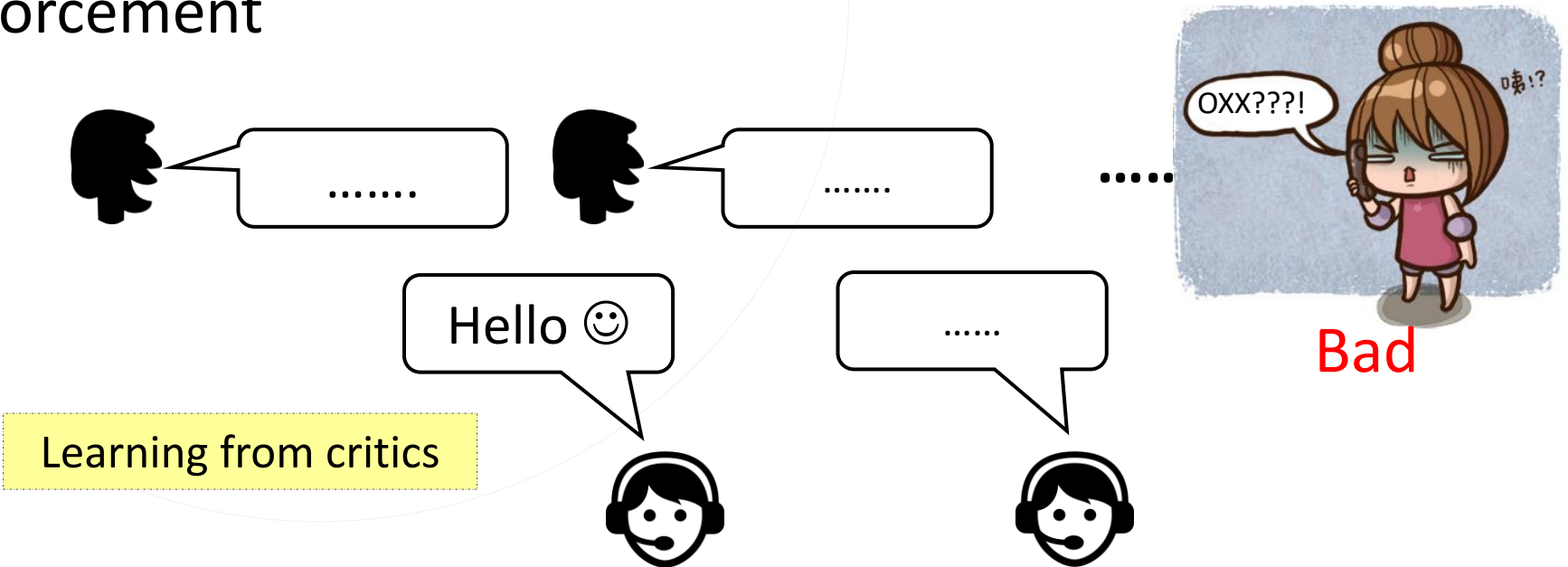


Supervised v.s. Reinforcement

- Supervised



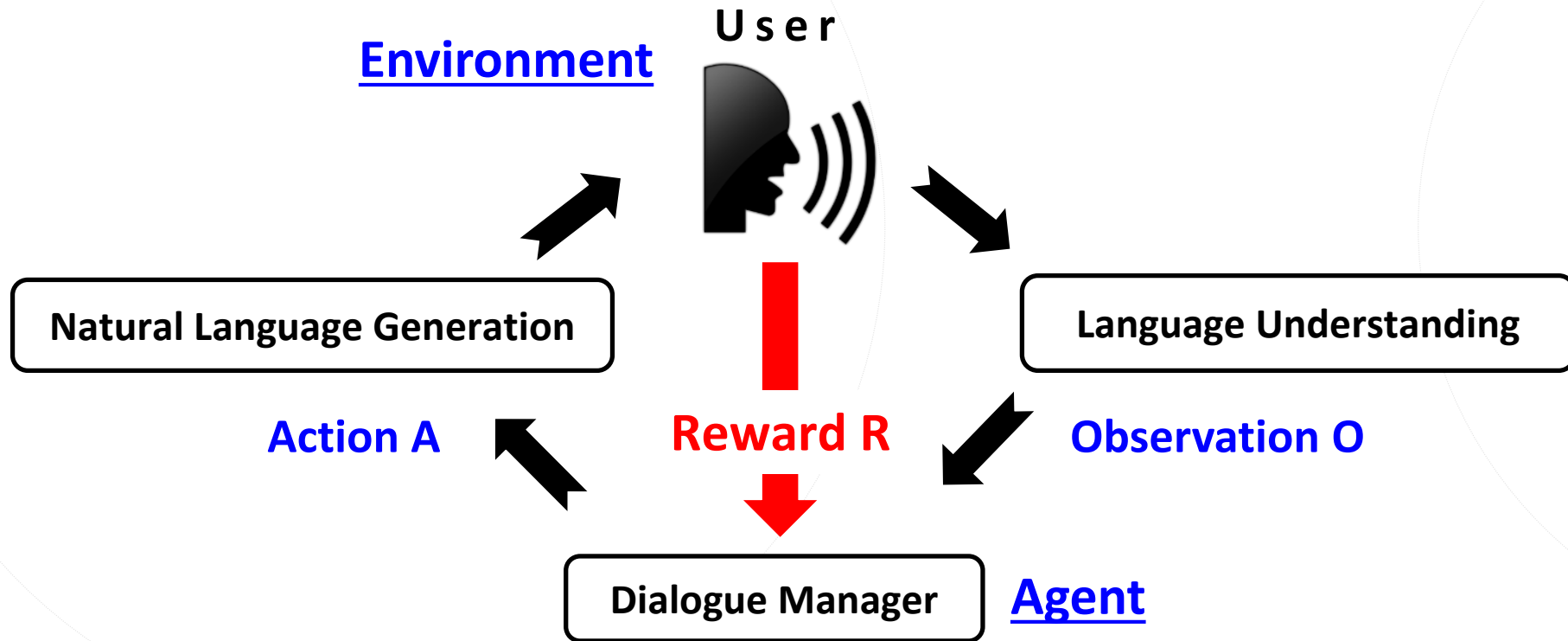
- Reinforcement





Dialogue Policy Optimization

- Dialogue management in a RL framework



Select the best action that maximizes the future reward

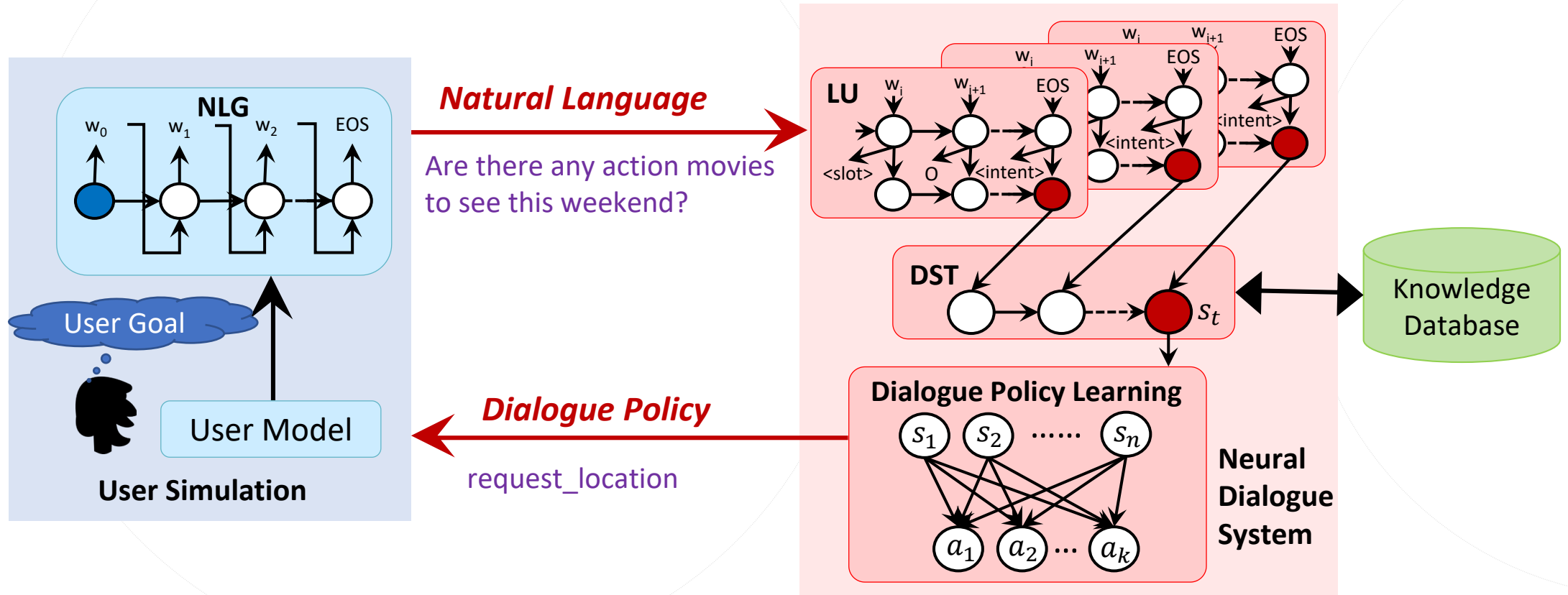
Reward for RL \cong Evaluation for System

- Dialogue is a special RL task
 - Human involves in interaction and rating (evaluation) of a dialogue
 - Fully human-in-the-loop framework
- Rating: correctness, appropriateness, and adequacy

- Expert rating	high quality, high cost
- User rating	unreliable quality, medium cost
- Objective rating	Check desired aspects, low cost



E2E Task-Completion Bot (TC-Bot) (Li et al., 2017)



Idea: SL for each component and RL for end-to-end training the neural dialogue system

Reinforcement Learning Agent [\(Li et al., 2017\)](#)

User goal – *Two tickets for “Deadpool” tomorrow 9PM at AMC Pacific Place 11 theater, Seattle.*

RULE BASED SYSTEM

User: Do you have any ticket for movie ‘Deadpool’ at 9:00 pm?

Agent: Which theater do you want?

User: AMC pacific place 11 theater, please!

Agent: How many tickets do you want?

User: Two!

Agent: What date?

User: Tomorrow!

Agent: Which city do you want?

User: Seattle!

Agent: Do you have any specification for zip code?

User: No, I do not care.

Agent: Do you have any distance constraint?

User: No, I do not care.

Agent: How about video format? 2d? Imax 3d?

User: Any kind is okay.

Agent: Do you have any constraint to price?

User: No!

Agent: Okay, I will issue 2 tickets for you, tomorrow 9:00 pm at AMC pacific place 11 theater, Seattle, movie ‘Deadpool’.

User: Okay, thanks.

Agent: Thanks!

REINFORCEMENT LEARNING SYSTEM

User: I want two tickets for movie ‘Deadpool’.

Agent: What is the starttime you prefer?

User: 9:00 PM

Agent: Which day?

User: Tomorrow

Agent: Which city?

User: Seattle

Agent: Which theater do you prefer?

User: AMC Pacific Place 11 theater.

Agent: Okay, I will issue 2 tickets for you, tomorrow 9:00 pm at AMC pacific place 11 theater, Seattle, movie ‘Deadpool’.

User: Okay, thanks.

Agent: Thanks!

Issue: no notion about what requests can be skipped

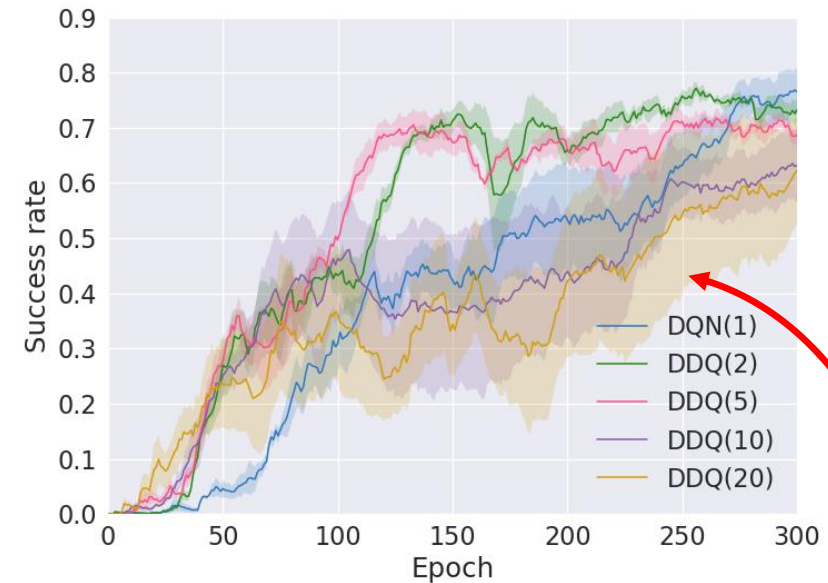
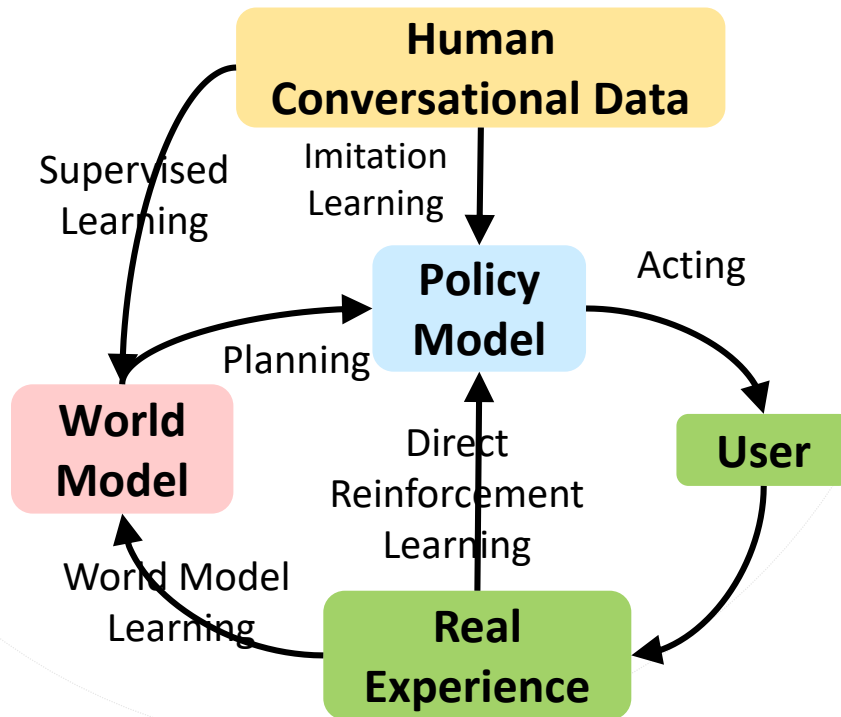
Skip the requests the user may not care about to improve efficiency





Planning – Deep Dyna-Q (Peng+, 2018)

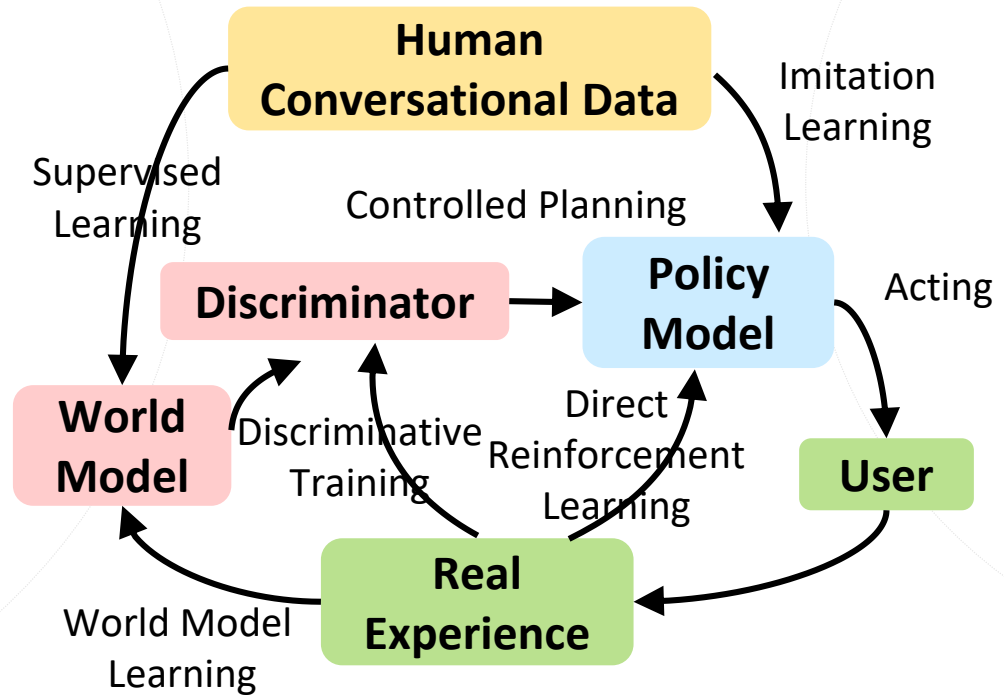
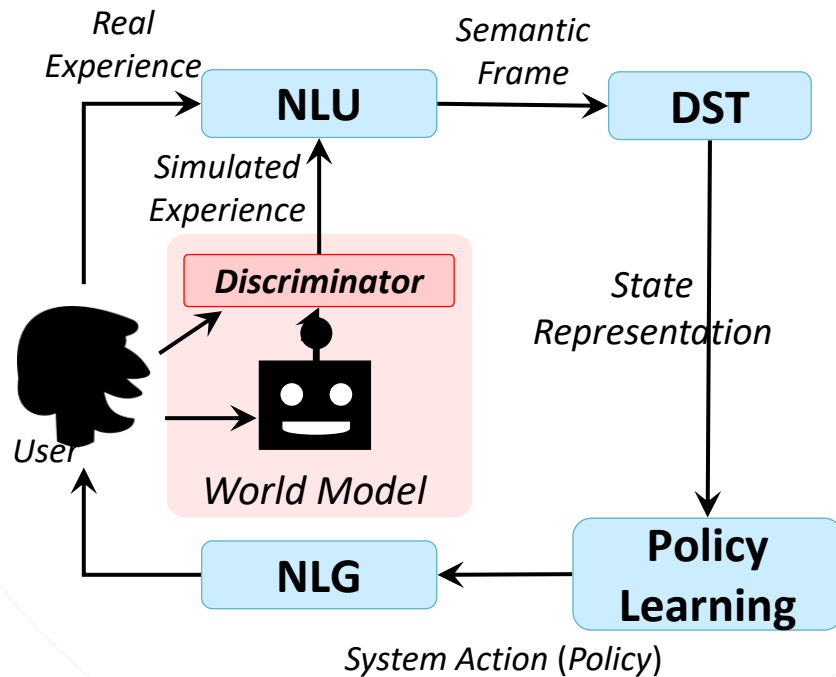
- Issues: sample-inefficient, discrepancy between simulator & real user
- Idea: learning with real users with planning



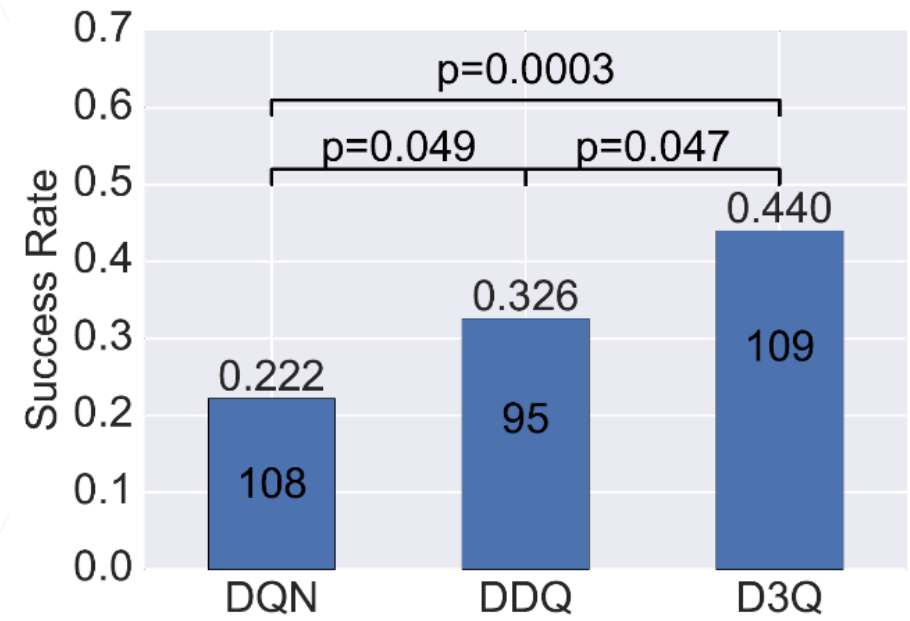
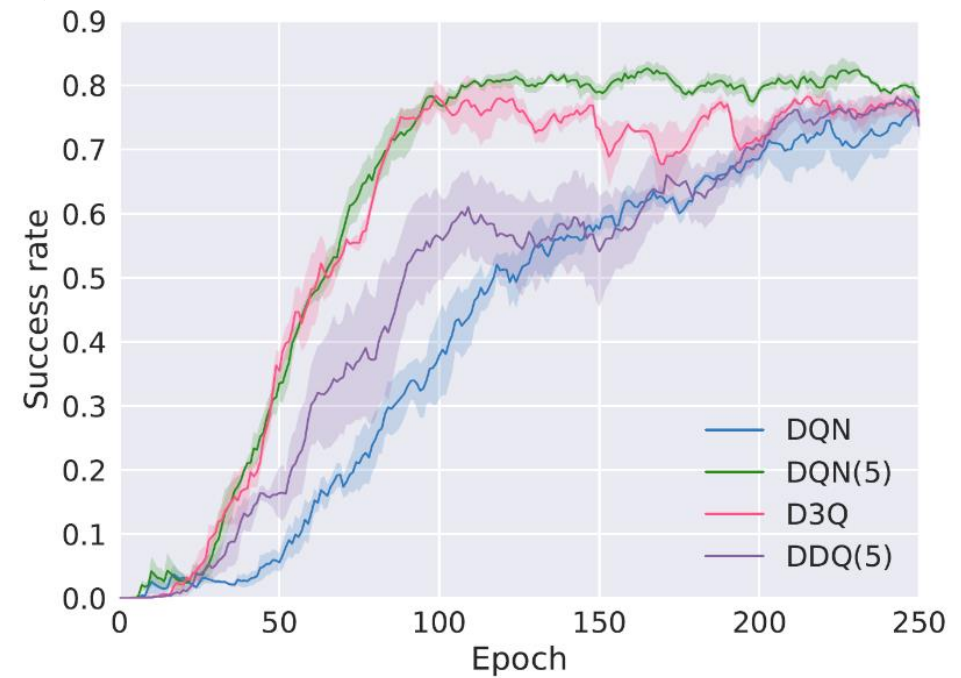
Policy learning suffers from the poor quality of fake experiences

Robust Planning – D3Q (Su+, 2018)

- Idea: add a *discriminator* to filter out the bad experiences



Robust Planning – D3Q (Su+, 2018)



The policy learning is more robust and shows the improvement in human evaluation

Dialogue Management Evaluation

- Metrics
 - Turn-level evaluation: system action accuracy
 - Dialogue-level evaluation: task success rate, reward

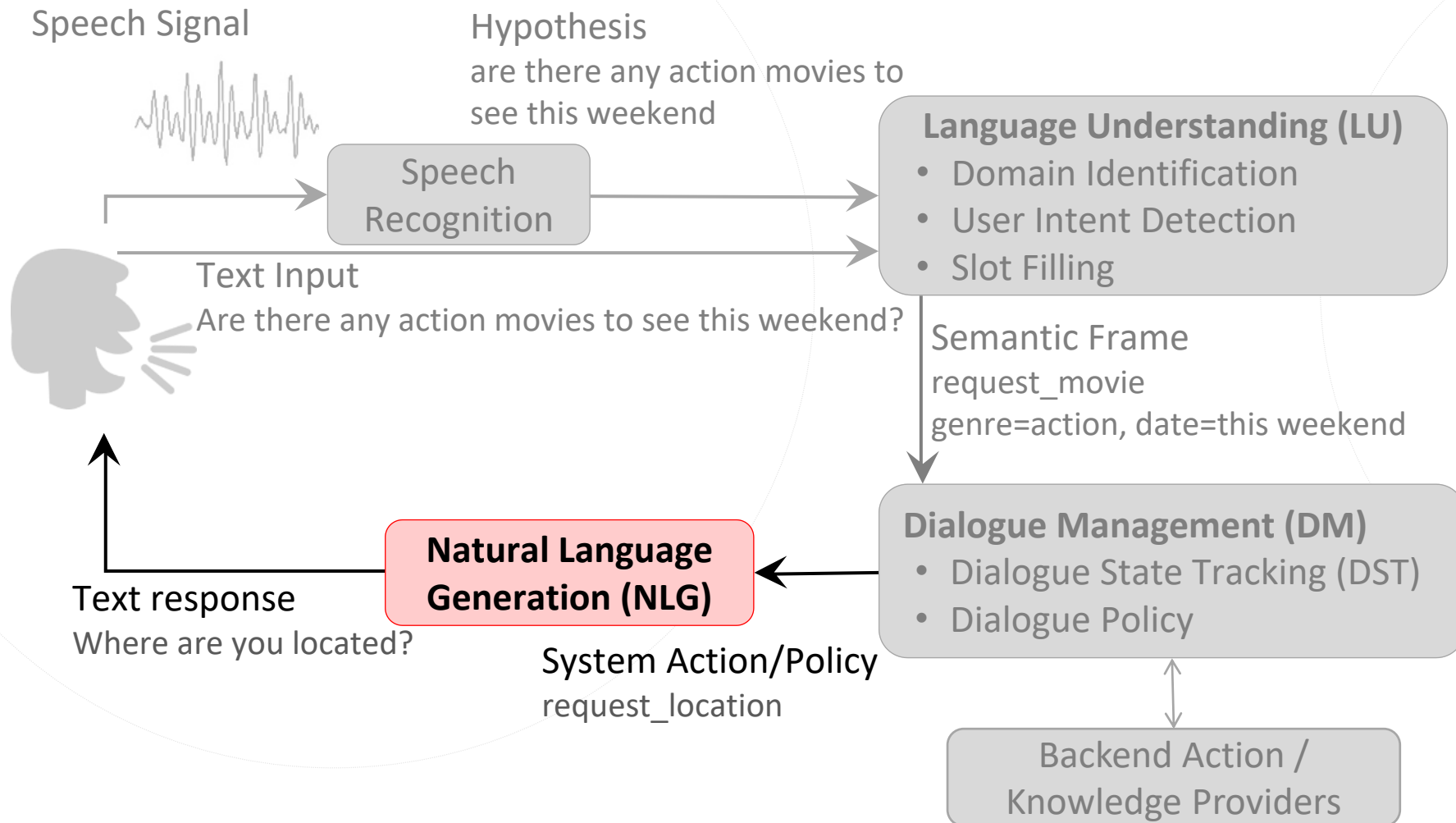


RL-Based DM Challenge

- SLT 2018 Microsoft Dialogue Challenge:
End-to-End Task-Completion Dialogue Systems
 - Domain 1: Movie-ticket booking
 - Domain 2: Restaurant reservation
 - Domain 3: Taxi ordering



Task-Oriented Dialogue Systems (Young, 2000)



Natural Language Generation (NLG)

- Mapping dialogue acts into natural language

`inform(name=Seven_Days, foodtype=Chinese)`



Seven Days is a nice Chinese restaurant





Template-Based NLG

- Define a set of rules to map frames to natural language

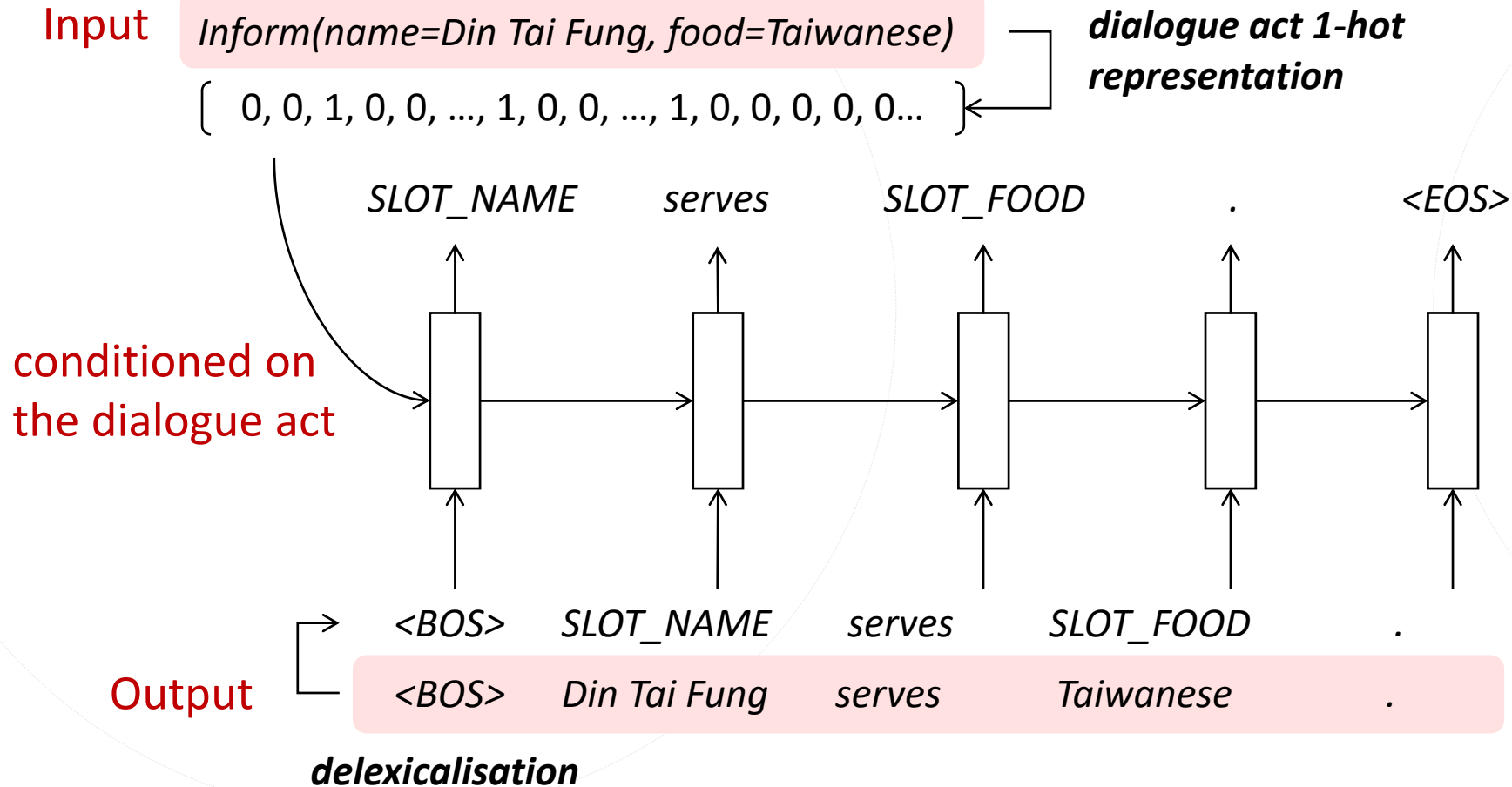
Semantic Frame	Natural Language
confirm()	“Please tell me more about the product your are looking for.”
confirm(area=\$V)	“Do you want somewhere in the \$V?”
confirm(food=\$V)	“Do you want a \$V restaurant?”
confirm(food=\$V,area=\$W)	“Do you want a \$V restaurant in the \$W.”

Pros: simple, error-free, easy to control

Cons: time-consuming, rigid, poor scalability



RNN-Based LM NLG (Wen et al., 2015)





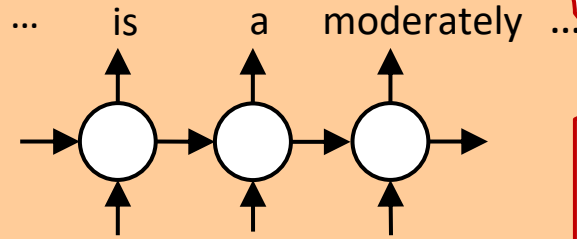
Issues in NLG

- Issue
 - NLG tends to generate **shorter** sentences
 - NLG may generate **grammatically-incorrect** sentences
- Solution
 - Generate word patterns in a order
 - Consider **linguistic patterns**

Hierarchical NLG w/ Linguistic Patterns (Su et al., 2018)

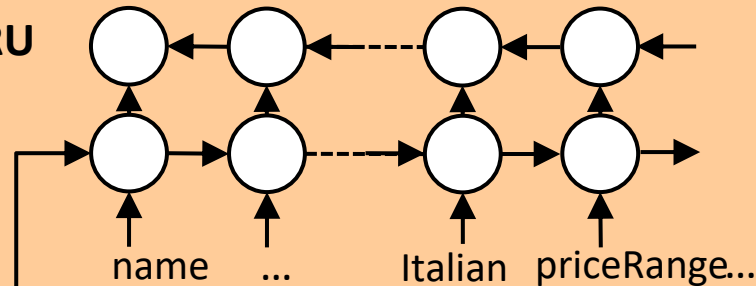
GRU Decoder

1. Repeat-input
2. Inner-Layer Teacher Forcing
3. Inter-Layer Teacher Forcing
4. Curriculum Learning



last output y_{t-1}^i ...All Bar One is a ...
 output from last layer y_t^{i-1} ...All Bar One is moderately..

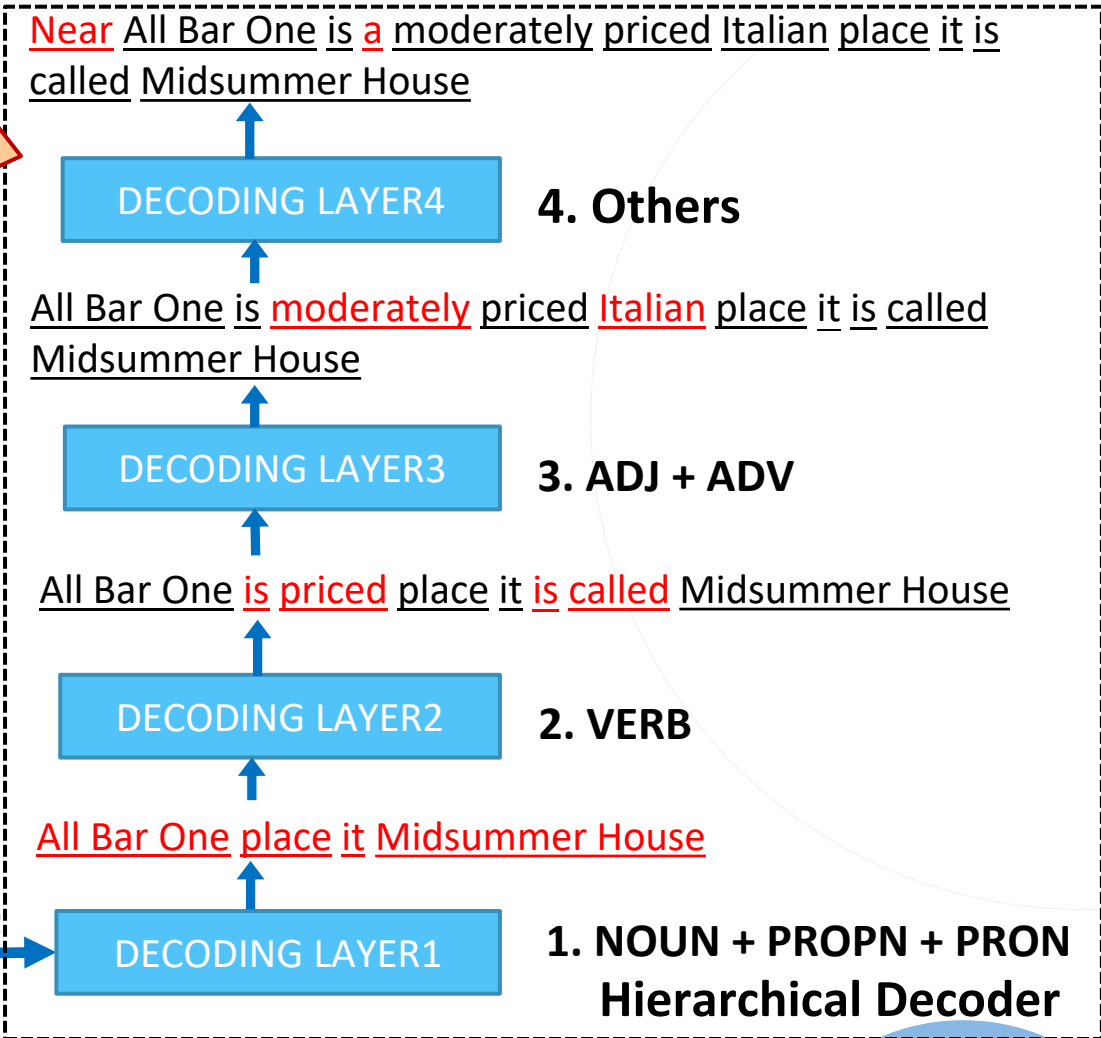
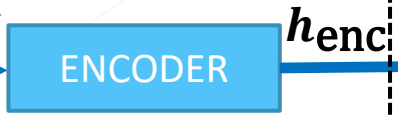
Bidirectional GRU Encoder



Semantic 1-hot Representation

[... 1, 0, 0, 1, 0, ...]

Input name[Midsummer House], food[Italian],
 Semantics priceRange[moderate], near[All Bar One]



1. NOUN + PROPN + PRON
Hierarchical Decoder



NLG Evaluation

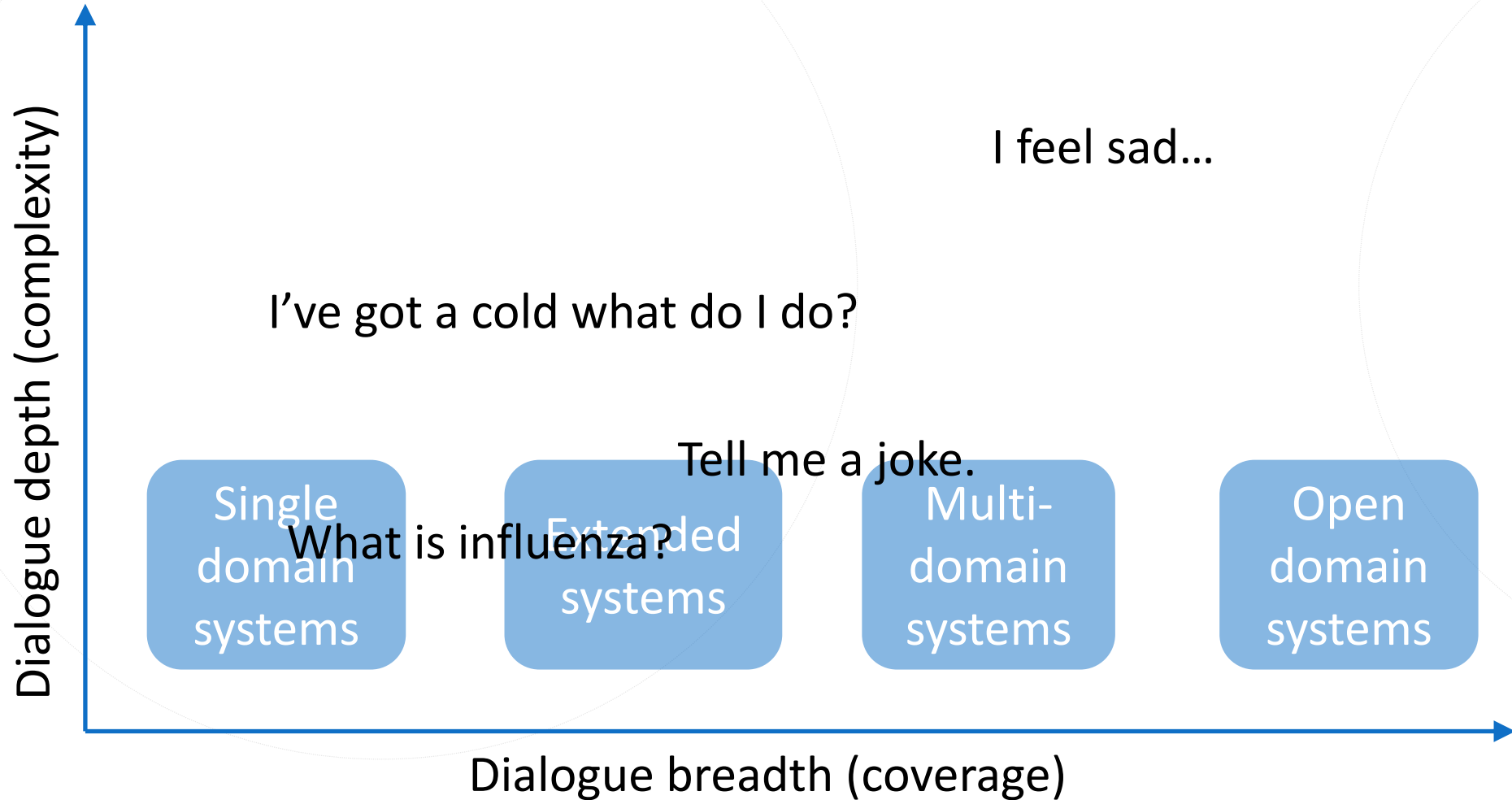
- Metrics

- Subjective: human judgement (Stent+, 2005)
 - Adequacy: correct meaning
 - Fluency: linguistic fluency
 - Readability: fluency in the dialogue context
 - Variation: multiple realizations for the same concept
- Objective: automatic metrics
 - Word overlap: BLEU (Papineni+, 2002), METEOR, ROUGE
 - Word embedding based: vector extrema, greedy matching, embedding average

There is a gap between human perception and automatic metrics

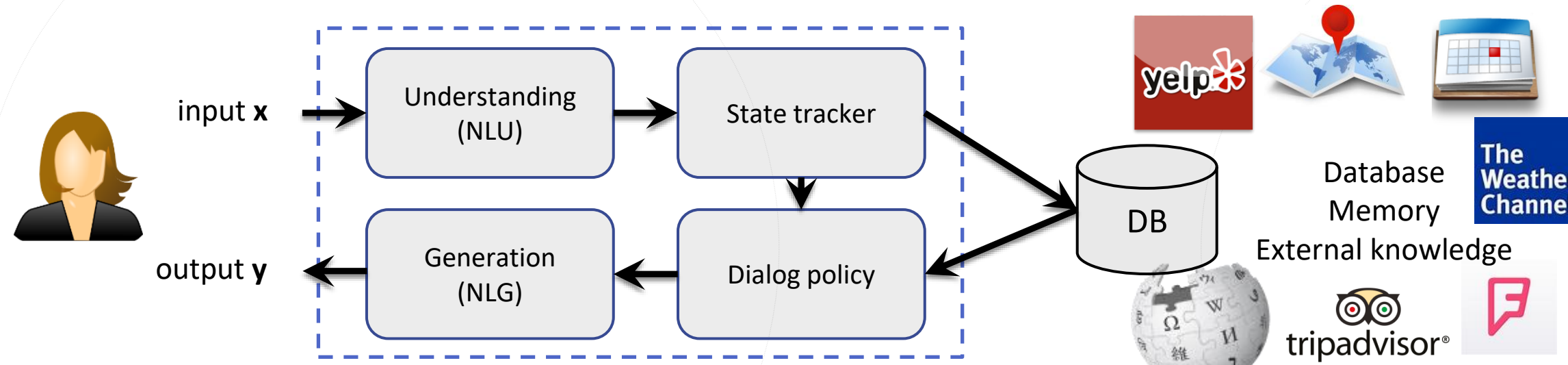


Evolution Roadmap

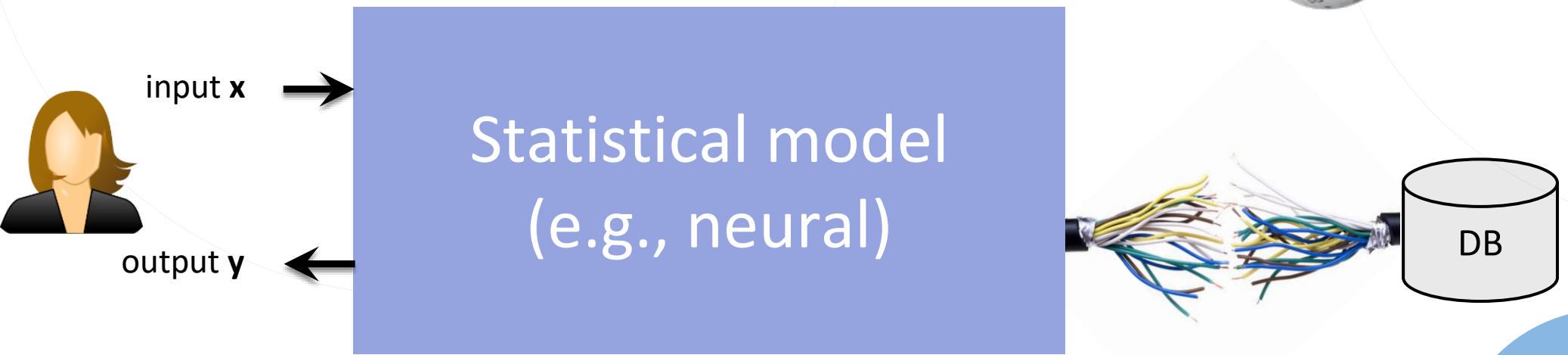


Dialogue Systems

Task-Oriented Dialogue



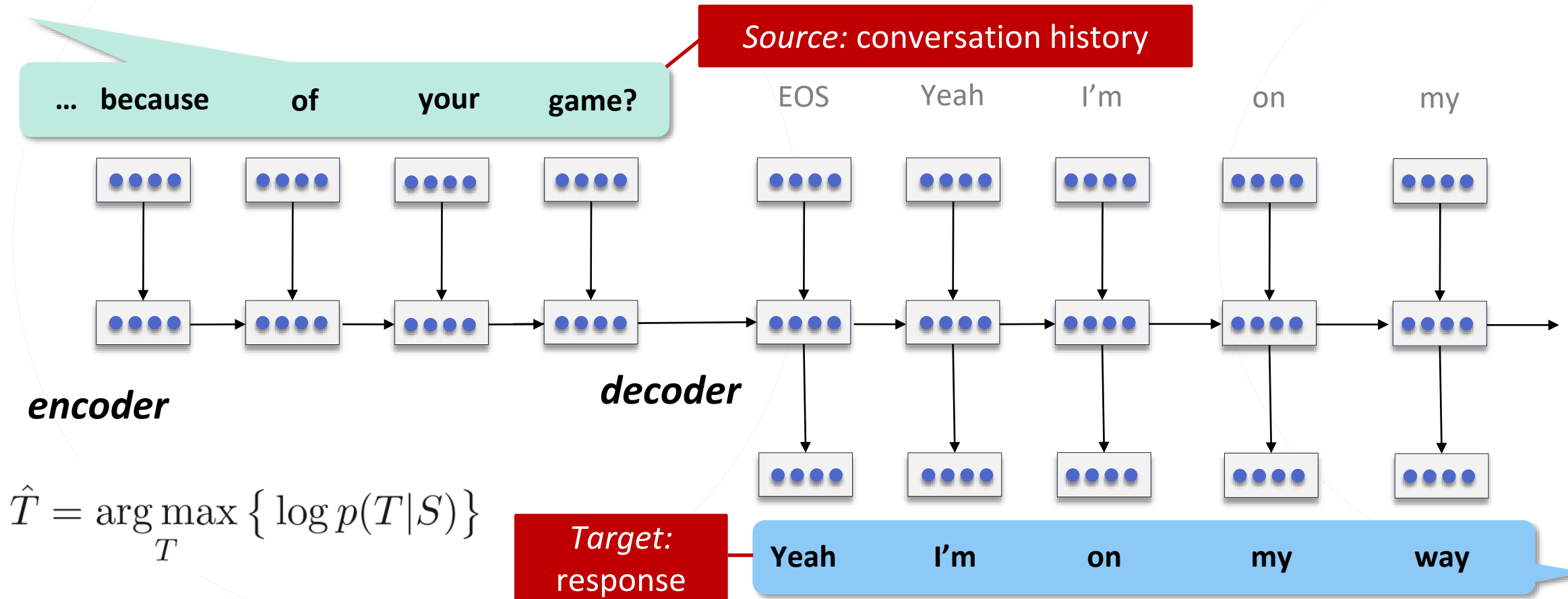
Fully Data-Driven





Chit-Chat Social Bots

Neural Response Generation ([Sordoni et al., 2015](#); [Vinyals & Le, 2015](#))

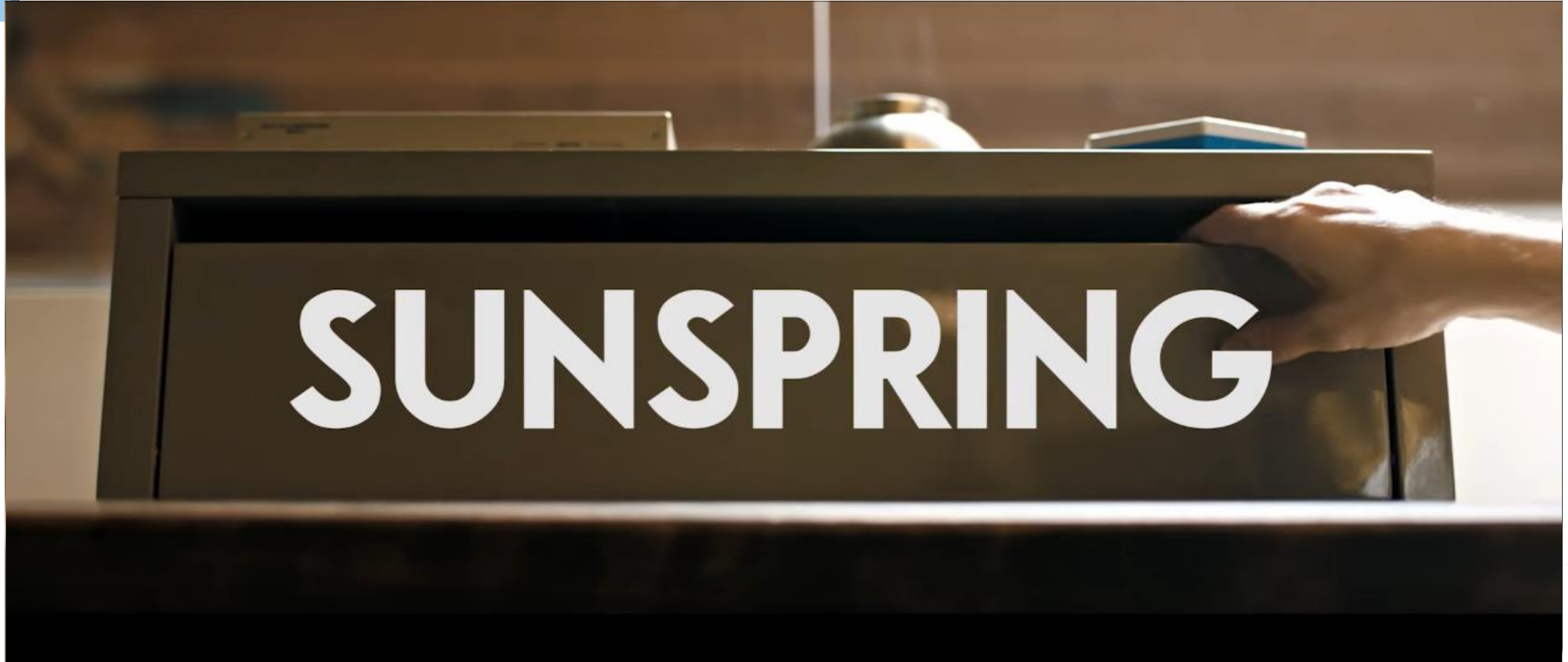


$$\hat{T} = \arg \max_T \{ \log p(T|S) \}$$

Learns to generate dialogues from offline data (no state, action, intent, slot, etc.)

Sci-Fi Short Film - SUNSPRING

<https://www.youtube.com/watch?v=LY7x2Ihqj>



Issue 1: Blandness Problem

Wow sour starbursts really do make your mouth water... mm drool.
Can I have one?

Of course!

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

I don't know.

'tis a fine brew on a day like this! Strong?

I'm not sure yet,

Well he was on in Bromley a while ago.

I don't even know what he's talking about.

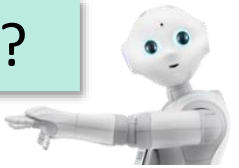
32% responses are general and meaningless

"I don't know"

"I don't know what you are talking about"

"I don't think that is a good idea"

"Oh my god"



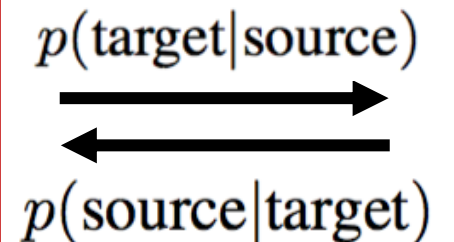
Mutual Information for Neural Generation (Li et al., 2016)

- Mutual information objective

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

$$\hat{T} = \arg \max_T \left\{ \underbrace{\log p(T|S)}_{\text{standard likelihood}} - \underbrace{\lambda \log p(T)}_{\text{anti-LM}} \right\}$$

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



MMI for Response Diversity ([Li et al., 2016](#))

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?

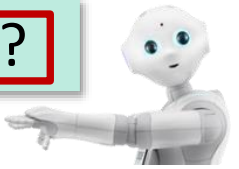
I think he'd be a good signing.

'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.



MMI for Response Diversity (Li et al., 2016)

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?**

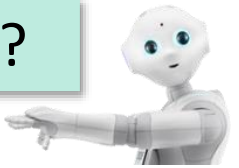
I think he'd be a **good signing**.

'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**.

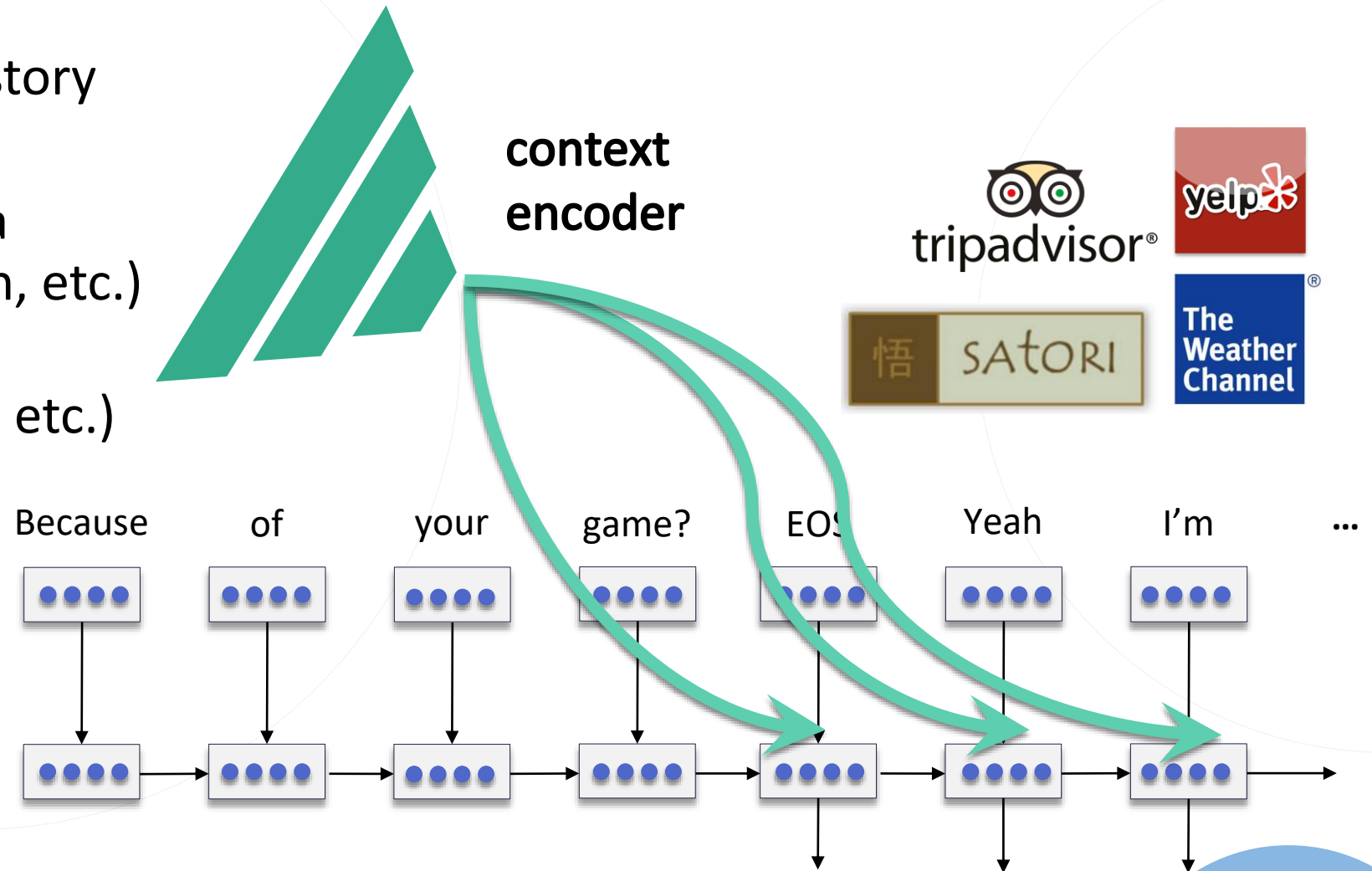




Real-World Conversations

□ Multimodality

- Conversation history
- Persona
- User profile data (bio, social graph, etc.)
- Visual signal (camera, picture etc.)
- Knowledge base
- Mood
- Geolocation
- Time



Issue 2: Response Inconsistency



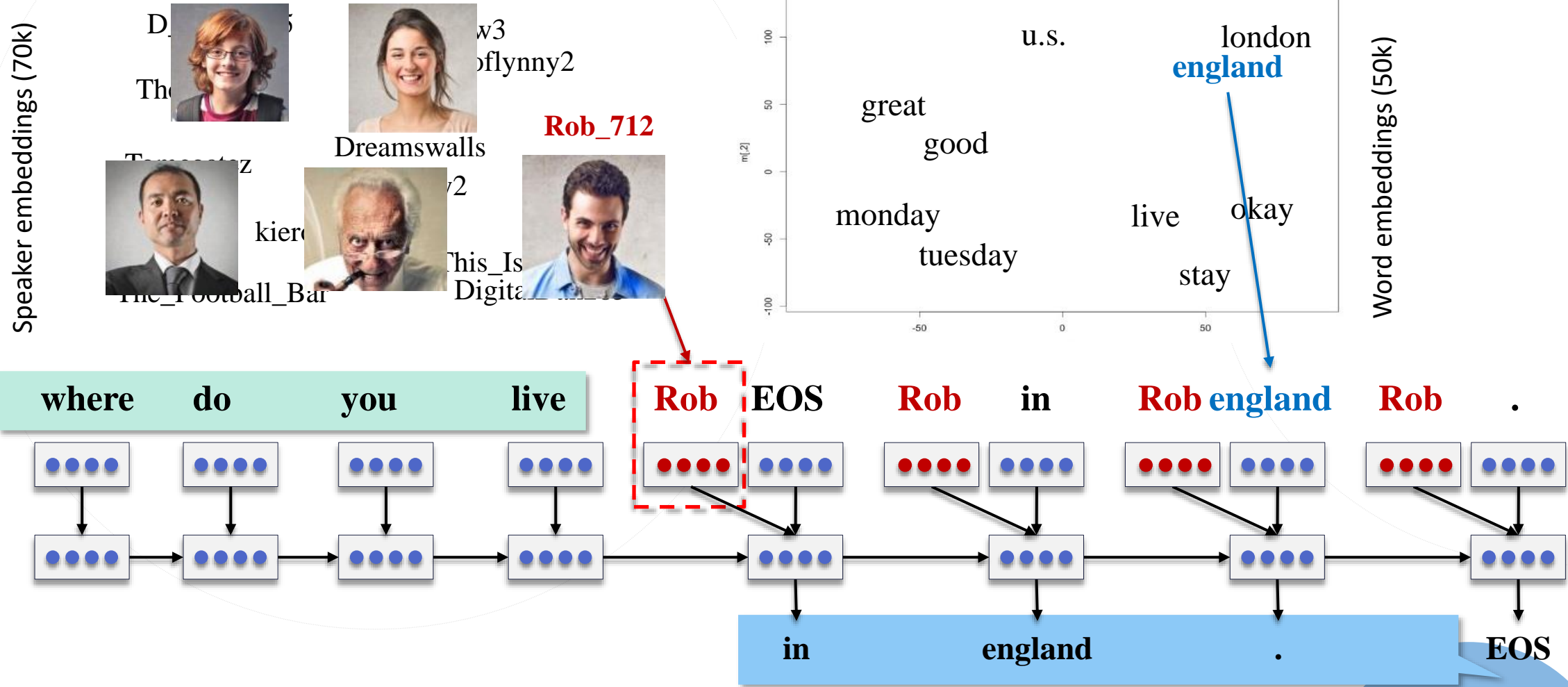
message Where do you live now?
response I live in Los Angeles.
message In which city do you live now?
response I live in Madrid.
message In which country do you live now?
response England, you?

message Where were you born?
response I was born in Canada.
message Where are you from?
response England, you?
message Where did you grow up?
response I grew up in Texas.

message How old are you?
response 16 and you?
message What's your age?
response 18.

message What is your major?
response I'm majoring in psychology
message What did you study in college?
response English lit.

Personalized Response Generation (Li et al., 2016)



Persona Model for Speaker Consistency (Li et al., 2016)



Baseline model → inconsistency

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 → consistency

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**.

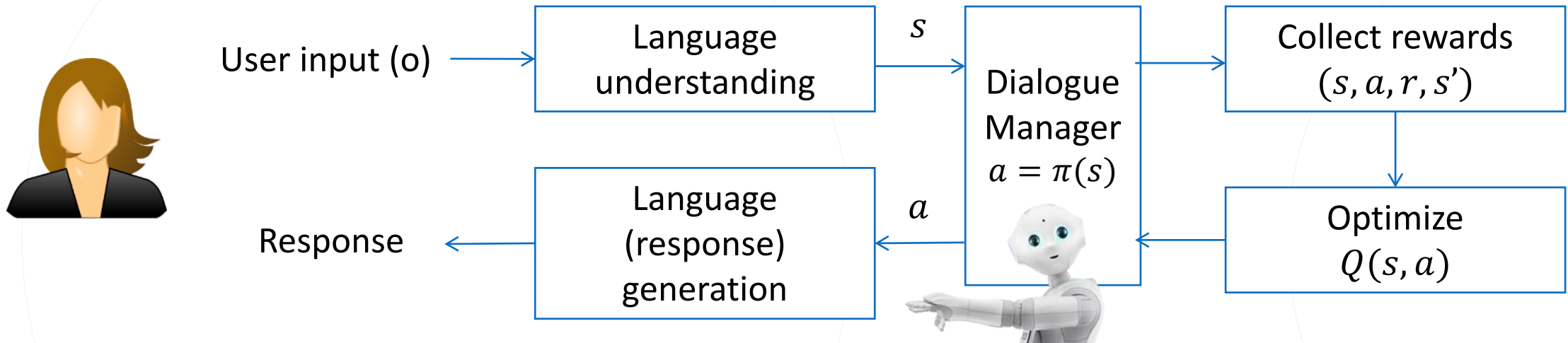
message What is your major?
response **Business**. You?

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**.

Issue 3: Dialogue-Level Optimization via RL



Application	State	Action	Reward
Task Completion Bots (Movies, Restaurants, ...)	User input + Context	Dialog act + slot-value	Task success rate # of turns
Info Bots (Q&A bot over KB, Web etc.)	Question + Context	Clarification questions, Answers	Relevance of answer # of turns
Social Bot (Xiaolce)	Conversation history	Response	Engagement(?)

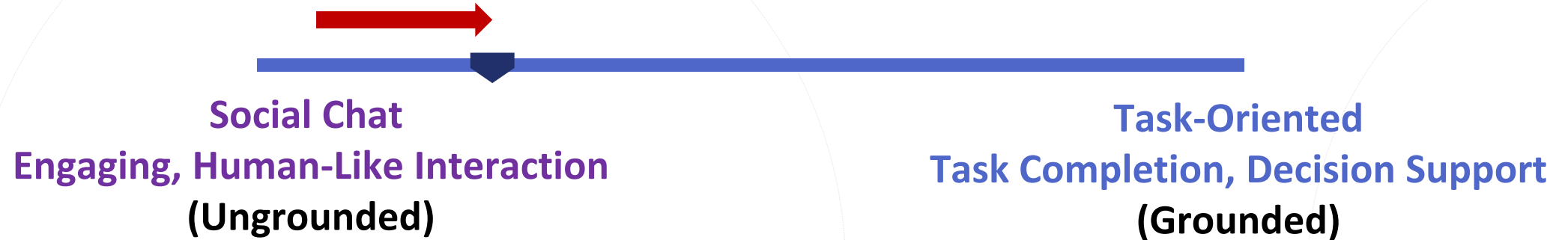
Deep RL for Response Generation (Li et al., 2016)



Input message	Supervised Learning Agent	Reinforcement Learning Agent
How old are you ?	I'm 16.	I'm 16, why are you asking ?
What's your full name ?	I have no idea .	What's yours ?
I don't want to go home tonight.	Really ?	Why ?
Do you have any feelings for me ?	I don't know what you are talking about.	I don 't want to hurt your feelings .
How much time do you have here?	Not long enough. Sorry, sir.	Ten seconds .
Shall we get started ?	Of course !	Yes . We 've got a lot of work to do here .
Do you play football ?	No, i don't	Yes. I love football !
We'd have to talk to him.	I mean, he's a good guy	About what ?
How come you never say it?	Because I don't want to hurt you .	I don 't think it 's a good idea to say it .

- RL agent generates more *interactive* responses
- RL agent tends to end a sentence *with a question* and hand the conversation over to the user

Issue 4: No Grounding ([Sordoni et al., 2015](#); [Li et al., 2016](#))

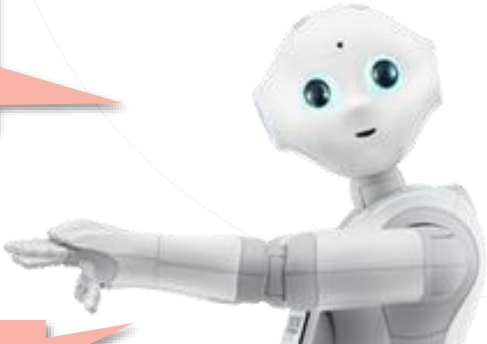


The weather is so depressing these days.

I know, I dislike rain too.
What about a day trip to eastern Washington?

Any recommendation?

Try **Dry Falls**, it's spectacular!



Knowledge-Grounded Responses ([Ghazvininejad et al., 2017](#))

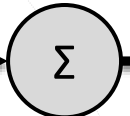
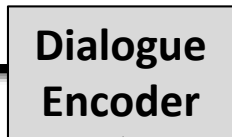
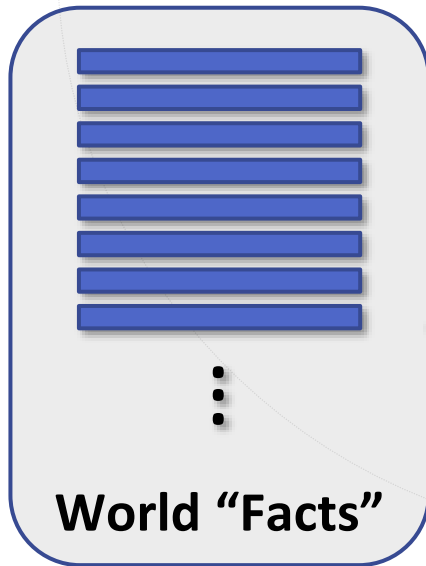


MIULA

NTU
84

Going to Kusakabe tonight

Conversation History



Try omakase, the best in town

Response



Conversation and Non-Conversation Data

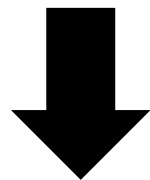


*You know any good **A** restaurant in **B**?*



*Try **C**, one of the best **D** in the city.*

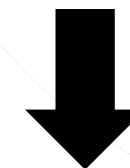
Conversation Data



Kisaku
★★★★☆ 515 reviews
\$\$ · Sushi Bars, Japanese
2101 N 55th St Ste 100
Seattle, WA 98103
b/t 56th St & N Kenwood Pl
Wallingford
(206) 545-9050
kisaku.com

"Kisaku is one of the best sushi restaurants in Seattle and located in the heart of Rainier Town." in 23 reviews

Knowledge Resource



*You know any good **Japanese** restaurant in **Seattle**?*

*Try **Kisaku**, one of the best **sushi** restaurants in the city.*



Knowledge-Grounded Responses ([Ghazvininejad et al., 2017](#))

A: Visiting the celebs at Los Angeles International Airport (LAX) - [...] w/ 70 others

B: Nice airport terminal. Have a safe flight.

A: Is that [...] in your photos? It's on my list of places to visit in NYC.

B: Don't forget to check out the 5th floor, while you are here, it's a great view.

A: Live right now on [...] Tune in!!!!

B: Listen to Lisa Paige

A: Been craving Chicken Pot Pie-who has the best? Trying [...] at [...] Must be Change of weather!

B: Love the pasta trattoria.

A: So [...] is down to one copy of Pound Foolish. I'm curious to see if they are re-ordering it.

B: Check out the video feed on 6 and take a picture of the Simpsons on the 3rd floor.

A: I wish [...] would introduce another vegetarian option besides the shroomburger. It's delicious but kind of ridiculous.

B: This is the best j.crew in the world. Try the lemonade!

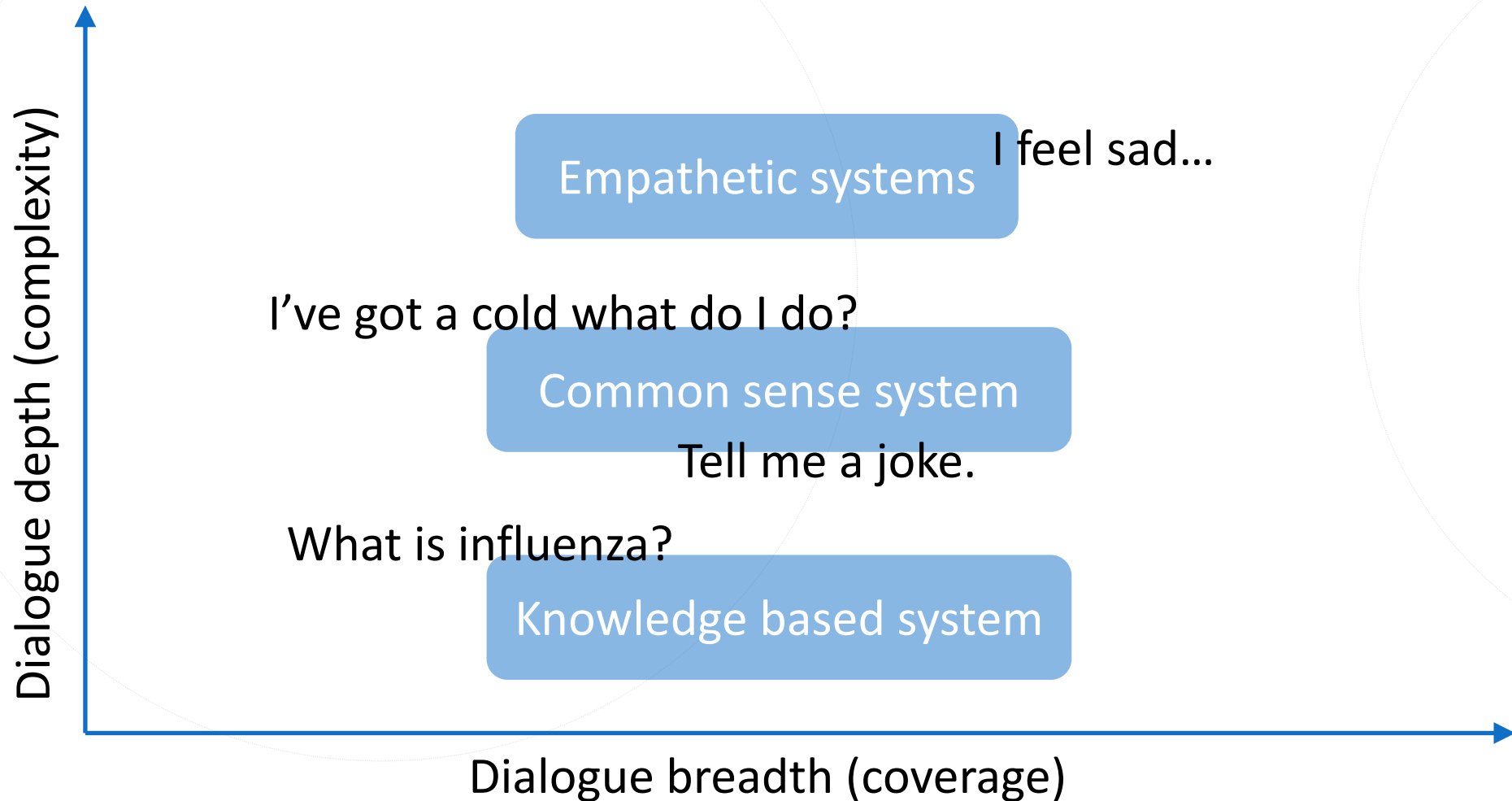
A: Just had an awesome dinner at [...] Great recommendation [...]

B: One of my favorite places I've ever been to in NYC. The food is great and the service is lackluster.

Results (23M conversations) outperforms competitive neural baseline (human + automatic eval)



Evolution Roadmap



Multimodality & Personalization ([Chen et al., 2018](#))

- Task: user intent prediction
- Challenge: language ambiguity



v.s.



① User preference

- ✓ Some people prefer “Message” to “Email”
- ✓ Some people prefer “Ping” to “Text”

② App-level contexts

- ✓ “Message” is more likely to follow “Camera”
- ✓ “Email” is more likely to follow “Excel”

Behavioral patterns in history helps intent prediction.



High-Level Intention Learning ([Sun et al., 2016](#); [Sun et al., 2016](#))

- High-level intention may span several domains

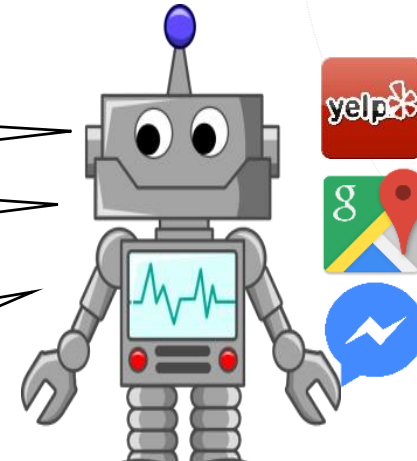
Schedule a lunch with Vivian.



What kind of restaurants do you prefer?

The distance is ...

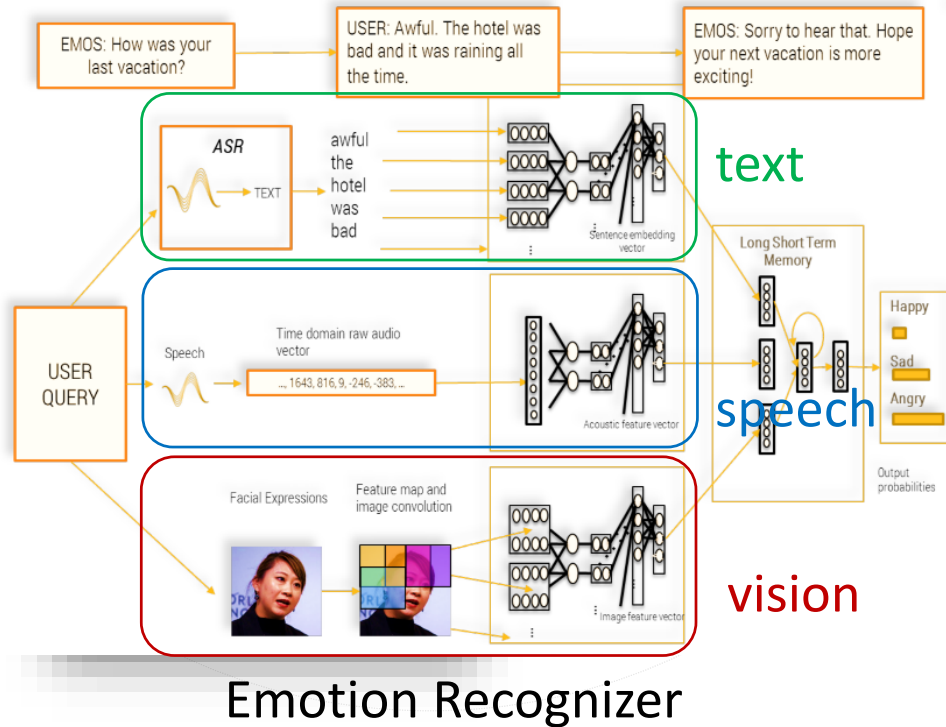
Should I send the restaurant information to Vivian?



Users interact via high-level descriptions and the system learns how to plan the dialogues

Empathy in Dialogue System (Fung et al., 2016)

- Embed an empathy module
 - Recognize emotion using multimodality
 - Generate emotion-aware responses



Zara - The Empathetic Supergirl



Face recognition output

```
{
  "recognition": "Race: Asian Confidence: 65.42750000000001 Smiling: 3.95896 Gender: Female Confidence: 88.9369",
  "race": "Asian",
  "race_confidence": "65.42750000000001",
  "smiling": "3.95896",
  "gender": "Female",
  "gender_confidence": "88.9369"
}
```

(index):1728

(index):1729

Cognitive Behavioral Therapy (CBT)

Mood Tracking



Pattern Mining



Daily lessons and check-ins

Depression Reduction



Content Providing



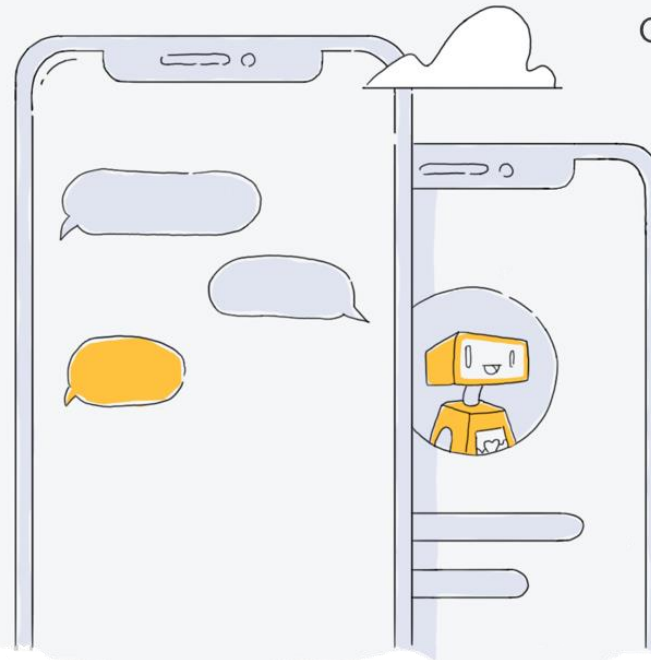
Always Be There



Know You Well



Quick conversation to feel better





Challenges & Conclusions

Challenge Summary

The human-machine interface is a hot topic but several components must be integrated!

Most state-of-the-art technologies are based on DNN

- Requires huge amounts of labeled data
- Several frameworks/models are available

Fast domain adaptation with scarce data + re-use of rules/knowledge

Handling reasoning and personalization

Data collection and analysis from un-structured data

Complex-cascade systems require high accuracy for working good as a whole





Her (2013)

What can machines achieve now or in the future?



Yun-Nung (Vivian) Chen

Assistant Professor, National Taiwan University

y.v.chen@ieee.org / <http://vivianchen.idv.tw>

