

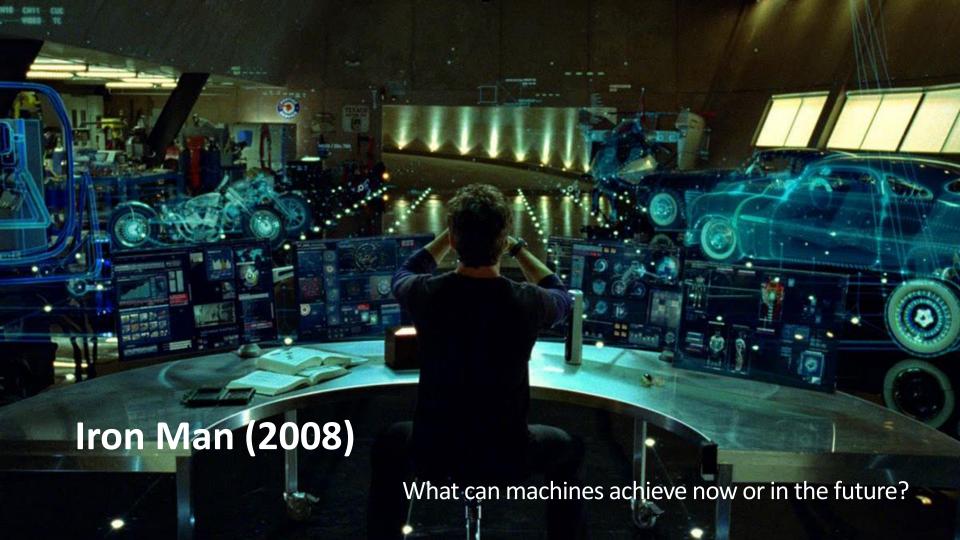
### **Towards Open-Domain Conversational AI**



YUN-NUNG (VIVIAN) CHEN 陳縕儂
HTTP://VIVIANCHEN.IDV.TW









## Language Empowering Intelligent Assistants







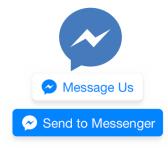


Apple Siri (2011)

Google Now (2012) Google Assistant (2016)

Microsoft Cortana (2014)









Amazon Alexa/Echo (2014)

Facebook M & Bot (2015)

Google Home (2016)

Apple HomePod (2017)



"I want to chat"

"I have a question"

"I need to get this done"

"What should I do?"

Turing Test (talk like a human)

Information consumption

Task completion



"I want to chat"

"I have a question"

"I need to get this done"

"What should I do?"

Turing Test (talk like a human)

Information consumption

Task completion



"I want to chat"

Turing Test (talk like a human)

"I have a question"

Information consumption

"I need to get this done"

Task completion

"What should I do?"

- What is today's agenda?
- Which room is SCAI workshop in?
- What does SCAI stand for?

"I want to chat"

"I have a question"

"I need to get this done"

"What should I do?"

Turing Test (talk like a human)

Information consumption

Task completion

**Decision support** 

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

7



"I want to chat"

"I have a question"

"I need to get this done"

"What should I do?"

Turing Test (talk like a human)

Information consumption

Task completion

- Is this brussels waffle worth to try?
- Is the SCAI workshop good to attend?



"I want to chat"

"I have a question"

"I need to get this done"

"What should I do?"

Turing Test (talk like a human)

Information consumption

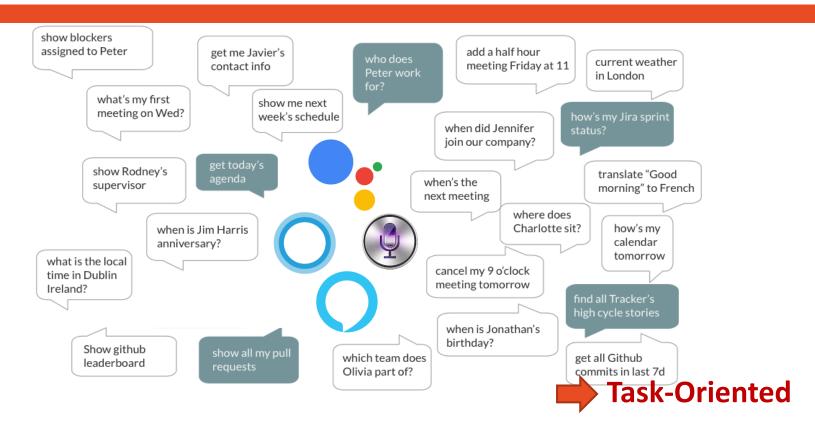
Task completion

Decision support

Task-Oriented
Dialogues

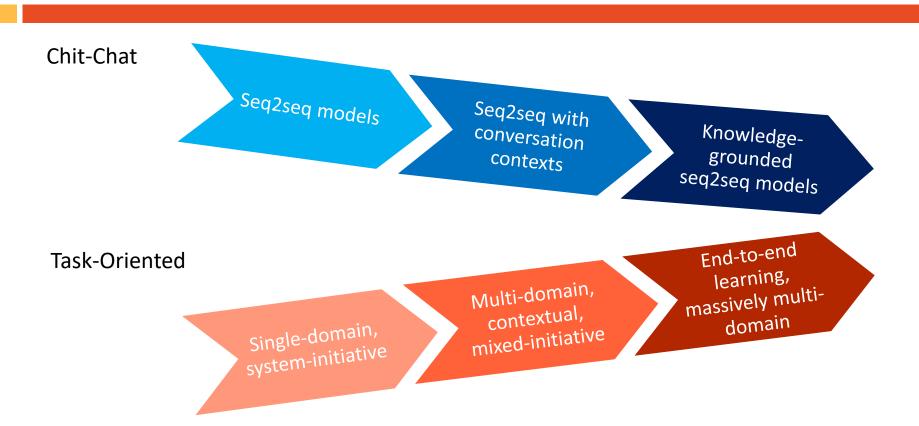


## Intelligent Assistants





# **Conversational Agents**



## Task-Oriented Dialogue Systems



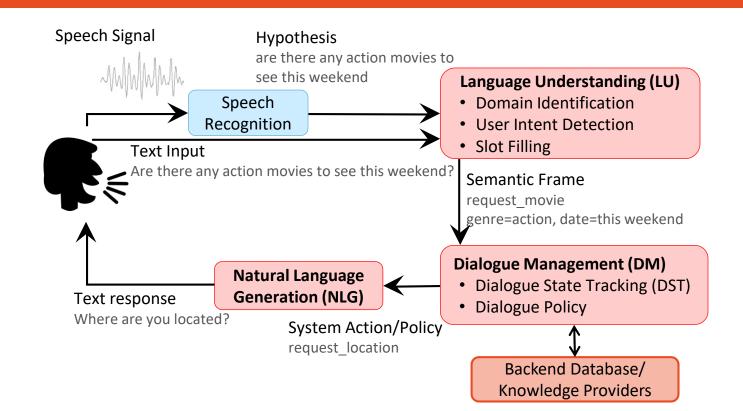
JARVIS – Iron Man's Personal Assistant



Baymax – Personal Healthcare Companion

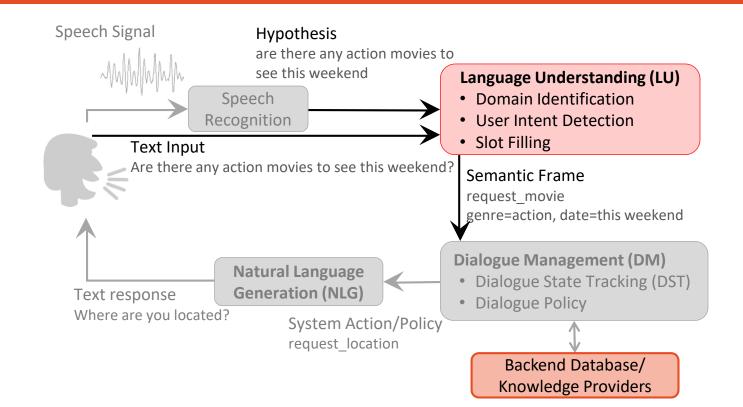


## Task-Oriented Dialogue Systems (Young, 2000)





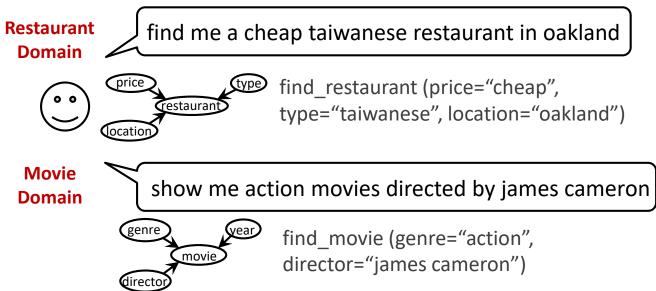
## Task-Oriented Dialogue Systems (Young, 2000)





## Semantic Frame Representation

- Requires a domain ontology: early connection to backend
- Contains core concept (intent, a set of slots with fillers)

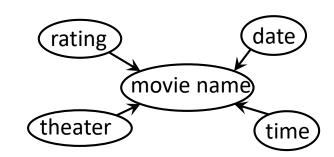




## Backend Database / Ontology

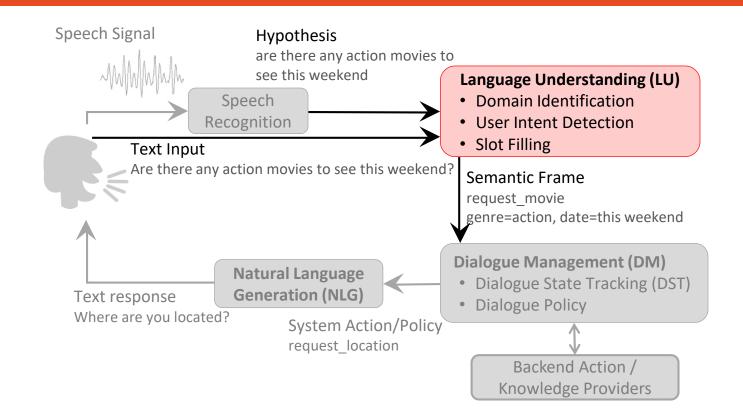
- Domain-specific table
  - □ Target and attributes
- Functionality
  - □ Information access: find specific entries
  - □ Task completion: find the <u>row</u> that satisfies the constraints

Movie Name	Theater	Rating	Date	Time
Iron Man Last	Taipei A1	8.5	2018/10/31	09:00
Iron Man Last	Taipei A1	8.5	2018/10/31	09:25
Iron Man Last	Taipei A1	8.5	2018/10/31	10:15
Iron Man Last	Taipei A1	8.5	2018/10/31	10:40





## Task-Oriented Dialogue Systems (Young, 2000)





## Language Understanding (LU)

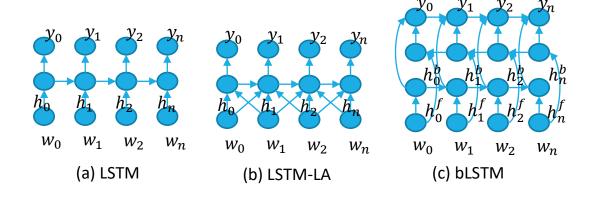
Pipelined 1. Domain 2. Intent 3. Slot Filling Classification Classification



## RNN for Slot Tagging — I (Yao et al, 2013; Mesnil et al, 2015)

#### Variations:

- a. RNNs with LSTM cells
- b. Input, sliding window of n-grams
- c. Bi-directional LSTMs

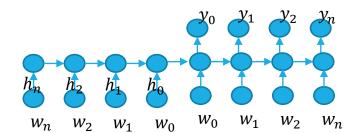


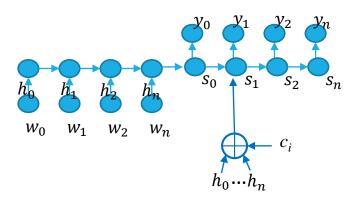


## RNN for Slot Tagging — II (Kurata et al., 2016; Simonnet et al., 2015)

- Encoder-decoder networks
  - Leverages sentence level information

- Attention-based encoder-decoder
  - Use attention (as in MT) in the encoderdecoder network







### Joint Semantic Frame Parsing

 Intent prediction Slot filling and and slot filling intent prediction Sequence-Parallel are performed in the same based (Liu and in two branches output sequence (Hakkani-Tur Lane, 2016) et al., 2016) (Intent) taiwanese food please **EOS** (Slot Filling) FromLoc ^ ToLoc from Seattle FIND\_REST, B-type Intent **Slot Filling** Prediction

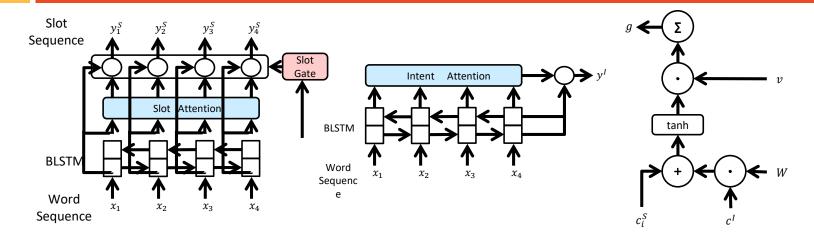


## Joint Model Comparison

	Attention Mechanism	Intent-Slot Relationship
Joint bi-LSTM	X	Δ (Implicit)
Attentional Encoder-Decoder	٧	Δ (Implicit)
Slot Gate Joint Model	٧	√ (Explicit)



## Slot-Gated Joint SLU (Goo et al., 2018)



Slot Gate

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

 $W^{S}$ : matrix for output layer  $b^{S}$ : bias for output layer

Slot Prediction 
$$y_i^S = softmax(W^S(h_i + c_i^S) + b^S)$$
  $\longrightarrow y_i^S = softmax(W^S(h_i + \boldsymbol{g} \cdot c_i^S) + b^S)$ 

g will be larger if slot and intent are better related

 $c_i^S$ : slot context vector  $c^{I}$ : intent context vector

W: trainable matrix v: trainable vector g: scalar gate value



## Contextual Language Understanding



```
U_1 send email to bob

S_1 B-contact_name

\Rightarrow send_email(contact_name="bob")

U_2 are we going to fish this weekend

S_2 B-message I-message I-message
I-message I-message
\Rightarrow send_email(message="are we going to fish this weekend")
```



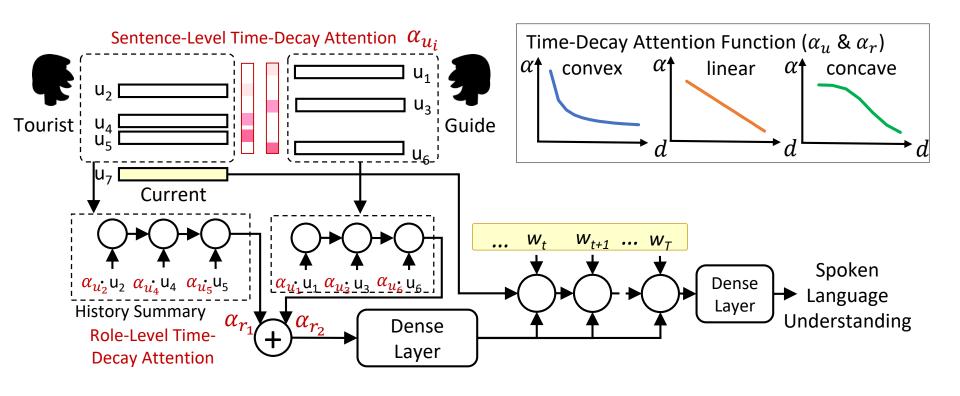
#### E2E MemNN for Contextual LU (Chen et al., 2016)

U: "Let's do 5:40"

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

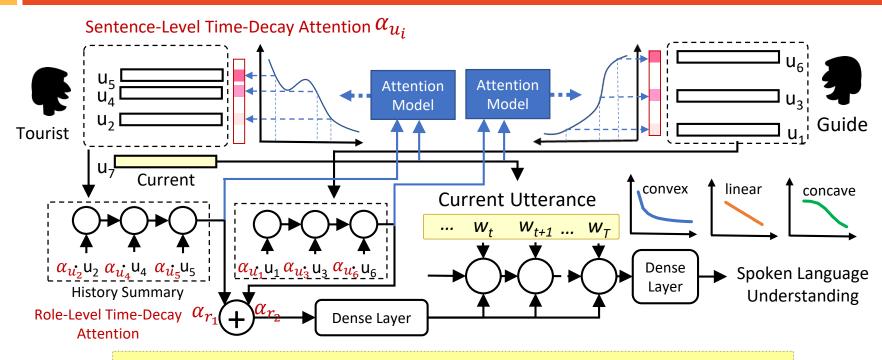


#### Role-Based & Time-Aware Attention (Su et al., 2018)





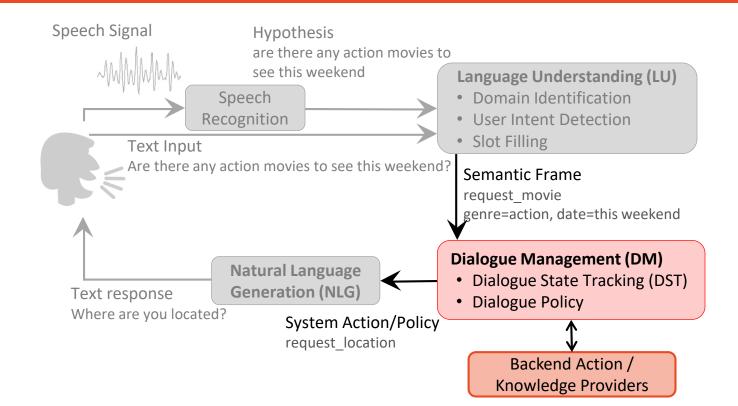
### Context-Sensitive Time-Decay Attention (Su et al., 2018)



Time-decay attention significantly improves the understanding results



## Task-Oriented Dialogue Systems (Young, 2000)



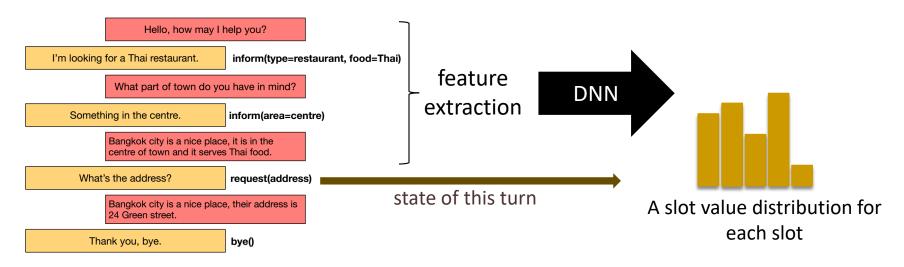


## Dialogue State Tracking

Hello, how may I help you? request (restaurant; foodtype=Thai) I'm looking for a Thai restaurant. What part of town do you have in mind? 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. What's the address? request (address) Bangkok city is a nice place, their address is 24 Green street. Thank you, bye. bye ()



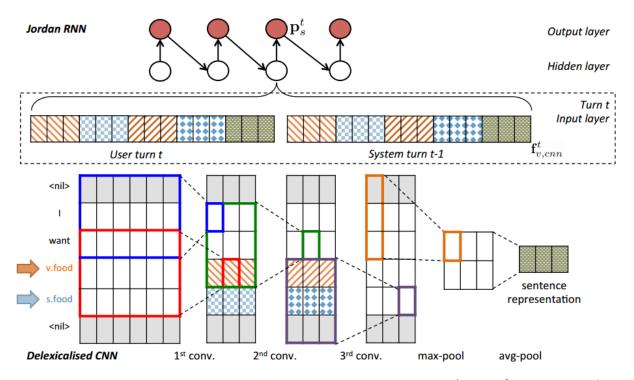
#### **DNN** for DST



multi-turn conversation



### RNN-CNN DST (Wen et al., 2016)



(Figure from Wen et al, 2016)



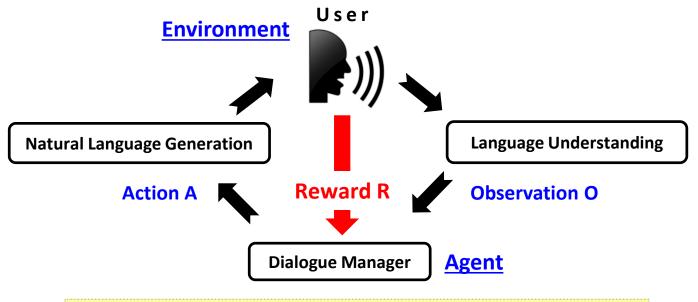
## Dialogue Policy Optimization





## Dialogue Policy Optimization

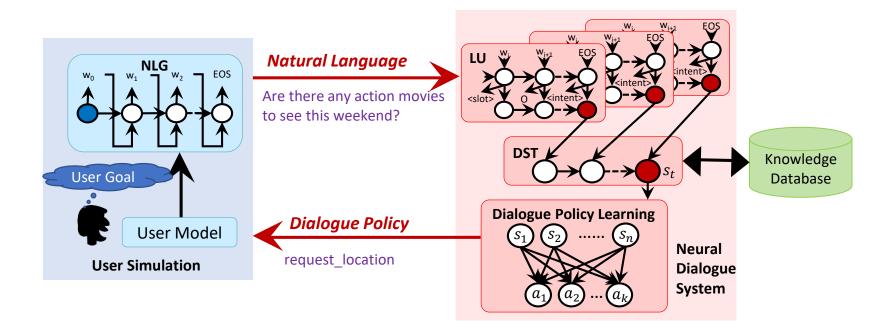
Dialogue management in a RL framework



Select the best action that maximizes the future reward



## 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'.

Issue 1: no notion about

what requests can be skipped

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!

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



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

User goal – Two tickets for "Hail Caesar" tomorrow 7:15PM at AMC Lowes Oak tree 6, Seattle.

#### **RULE BASED SYSTEM**

User: Do you any ticket for movie 'Hail Caesar' at theater

AMC Lowes Oak Tree 6?

Agent: What start time do you prefer?

User: 7:15 pm Tomorrow

Agent: Which city you are looking for?

User: Seattle

Agent: How many tickets do you want?

User: Two

Agent: Okay, I have issued two tickets for movie 'Hail Caeser'

tomorrow 7:50 pm at theater AMC Lowes Oak Tree 6,

Seattle.

User: Oh, no. Issue 2: no principle way to deal Agent: Thanks! with LU errors/uncertainty

#### REINFORCEMENT LEARNING SYSTEM

User: Do you any ticket for movie 'Hail Caesar' at theater AMC

Lowes Oak Tree 6?

Agent: What start time do you prefer?

User: 7:15 pm Tomorrow

Agent: Which city you are looking for?

User: Seattle

Agent: Okay, I will issue you tickets for movie 'Hail Caeser' tomorrow 7:50 pm at theater AMC Lowes Oak Tree 6, Seattle.

User: Oh, no.

Agent: It is 7:50 pm, right?

User: Two tickets for 7:15 pm, please.

Agent: Okay, I will issue you two tickets for movie 'Hail Caeser' tomorrow **7:15 pm** at theater AMC Lowes Oak Tree 6, Seattle.

User: Thanks! Agent: Thanks!

Automatically learn when and what to confirm for better robustness



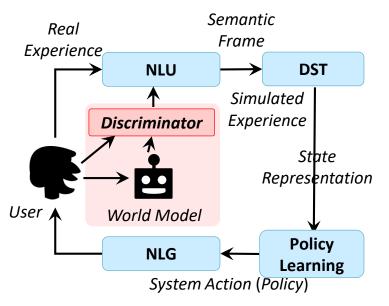
### RL in Dialogue Systems

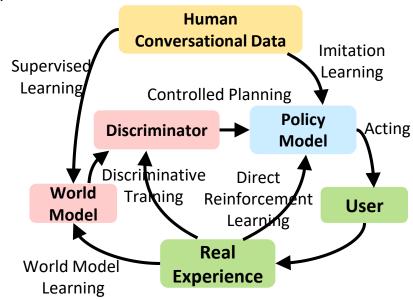
- Sample inefficient, hard to design reward function, local optima...
- 2. Real users are expensive
- Discrepancy between real users and simulators



### D3Q: Discriminative Deep Dyna-Q (Su et al., 2018)

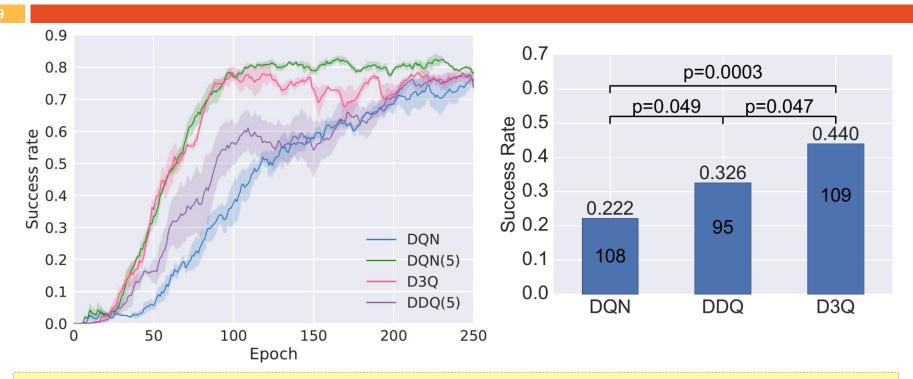
- Idea
  - learning with real users with planning
  - **a** add a *discriminator* to filter out the bad experiences







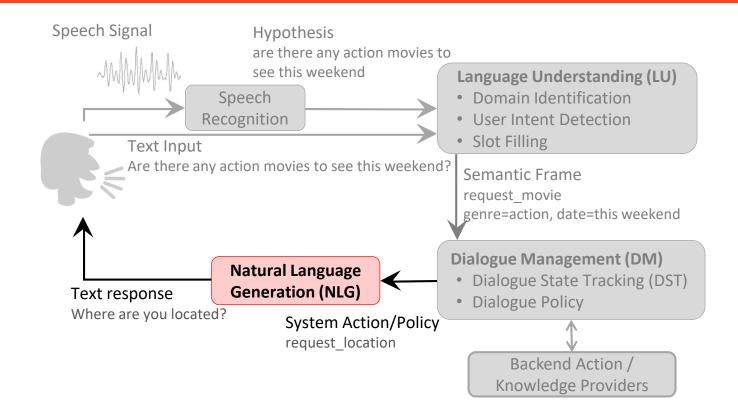
### D3Q: Discriminative Deep Dyna-Q (Su et al., 2018)



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



# 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

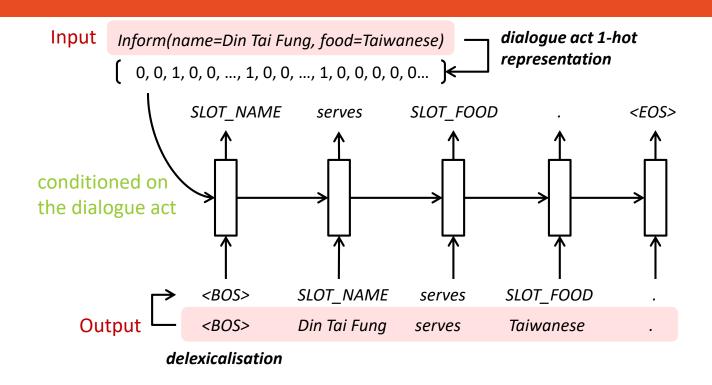
Semantic Frame	Natural Language	
confirm()	"Please tell me more about the product your are looking for."	
confirm(area=\$V)	"Do you want somewhere in the \$V?"	
confirm(food=\$V)	"Do you want a \$V restaurant?"	
confirm(food=\$V,area=\$W)	"Do you want a \$V restaurant in the \$W."	

**Pros:** simple, error-free, easy to control

Cons: time-consuming, rigid, poor scalability



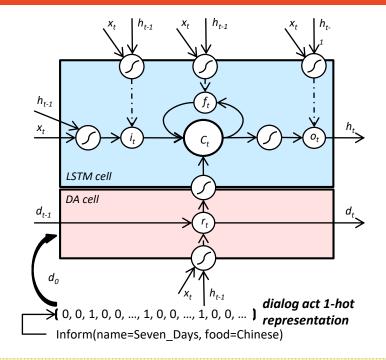
### RNN-Based LM NLG (Wen et al., 2015)





### Semantic Conditioned LSTM (Wen et al., 2015)

- Issue: semantic repetition
  - Din Tai Fung is a great Taiwanese restaurant that serves Taiwanese.
  - Din Tai Fung is a child friendly restaurant, and also allows kids.



Idea: using gate mechanism to control the generated semantics (dialogue act/slots)

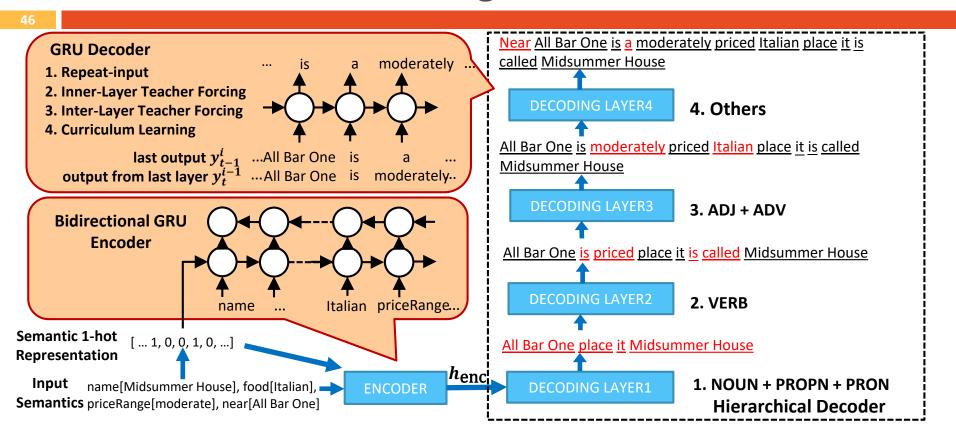


### Issues in NLG

- 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)

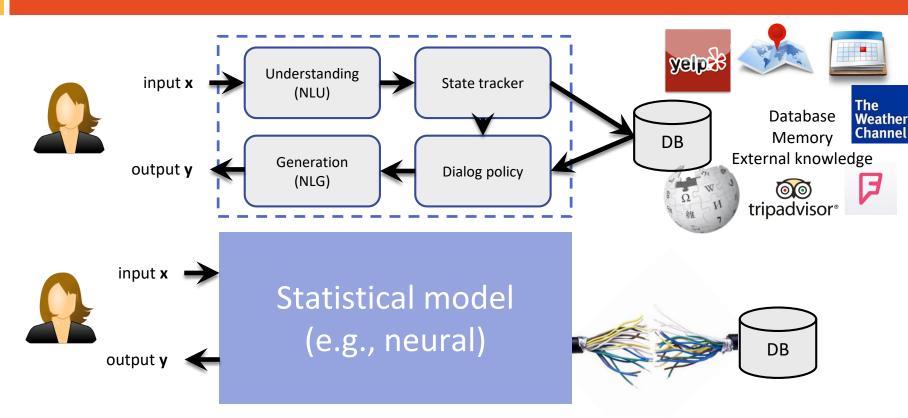




### **Evolution Roadmap**

depth (complexity) I feel sad... I've got a cold what do I do? Tell me a joke. Single domain is influenza?led Open Multi-Dialogue domain domain systems systems systems systems Dialogue breadth (coverage)

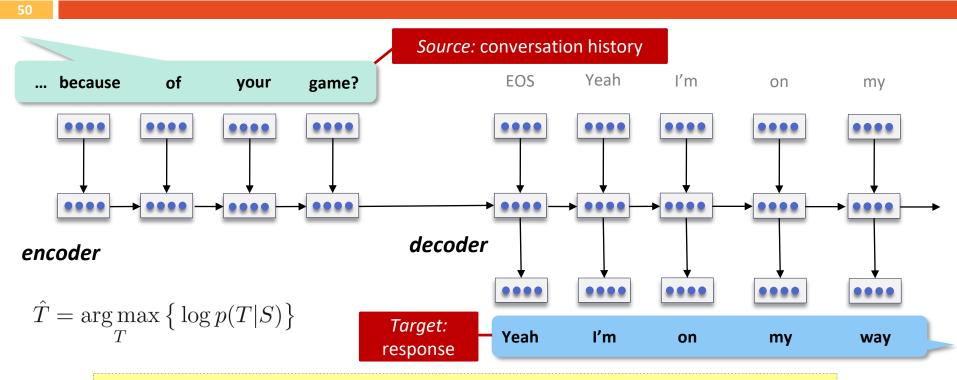
# **Dialogue Systems**



# <sup>49</sup> Chit-Chat Social Bots



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



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



### 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! Stro

I'm not sure yet,

Well he was on in Bromley a wh

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"

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



### 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.

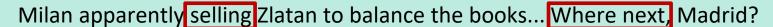


### MMI for Response Diversity (Li et al., 2016)

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

Can I have one?

Of course you can! They're delicious!



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!

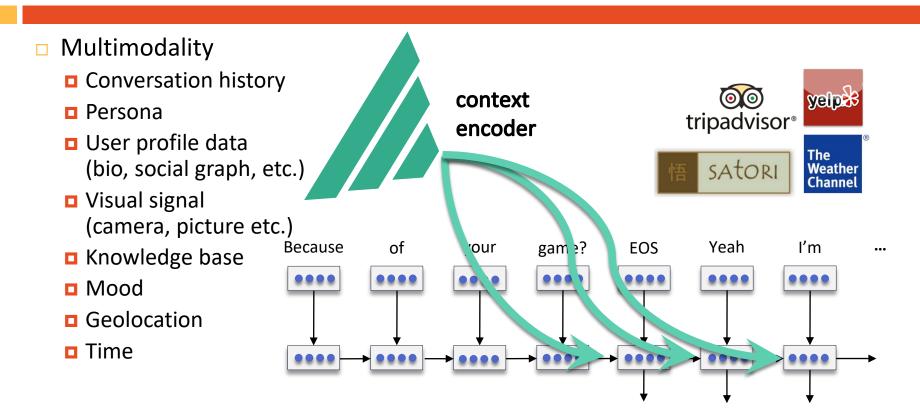
Well he was on in Bromley a while ago... still touring

I've never seen him live.





### **Real-World Conversations**



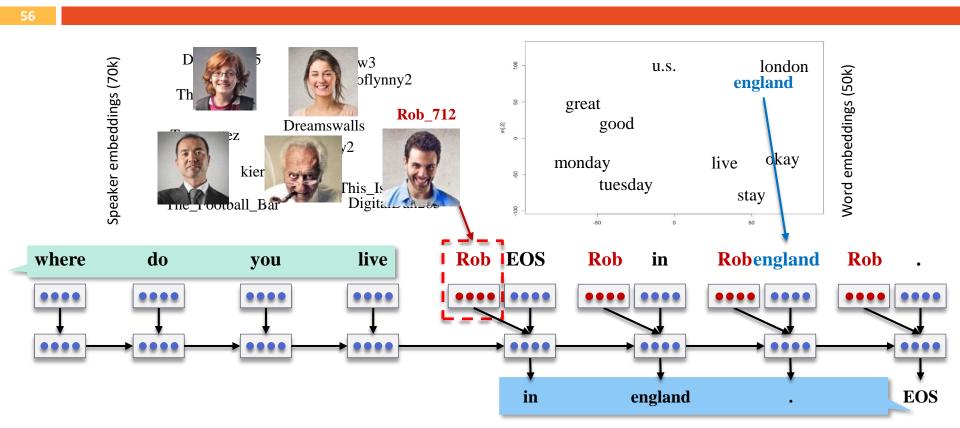


# Issue 2: Response Inconsistency

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

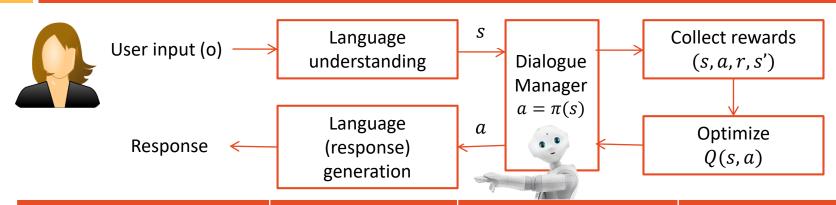
```
Where is your hometown?
message
          I live in Indonesia.
response
          Where are you from?
message
         England, you?
response
          In which city do you live now?
message
          I live in Jakarta.
response
          In which country do you live now?
message
          I live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

#### Persona model using speaker embedding → consistency

```
Where is your hometown?
message
          I'm from England.
response
          Where are you from?
message
          I'm from England.
response
          In which city do you live now?
message
          I live in London.
response
          In which country do you live now?
message
          I live in England.
response
          What is your major?
message
         Business. You?
response
          What did you study in college?
message
          I did business studies.
response
          How old are you?
message
          I'm 18.
response
          What is your age?
message
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)

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

No grounding into a real calendar, but the "shape" of the conversation is fluent and plausible.

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

H: well, i'm just tired

M: me too, i'm here if you wanna talk

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 2 pm work for you?

M: works for me.

H: well let's say 2 pm then i 'll see you there

M: sounds good.



### Chit-Chat v.s. Task-Oriented



Task-Oriented
Task Completion, Decision Support
(Grounded)

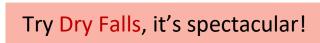
The weather is so depressing these days.



I know, I dislike rain too.

What about a day trip to eastern Washington?

Any recommendation?





### Image-Grounded Dialogues (Mostafazadeh et al., 2017)

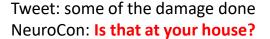
Neural model, trained on 250K Twitter convos:

-Turn 1: tweet + accompanying photo

-Turn 2: tweet is a question

Given data, other signals can be exploited in the same way: program state, database info, geolocation, etc.







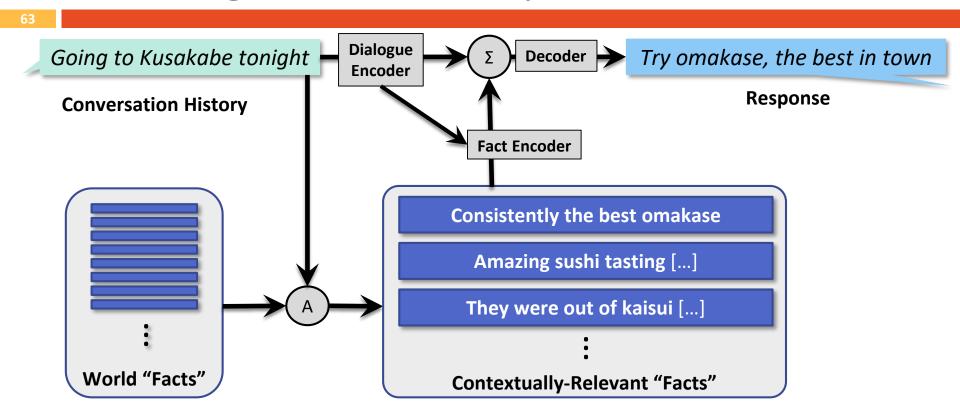
Tweet: my poor baby moms saying it 's probably totaled so rip Bessie NeuroCon: How did you manage to do that?



Tweet: as penitence I offer this NeuroCon: can I have one of those?



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





### Conversation and Non-Conversation Data



**Conversation Data** 



Know



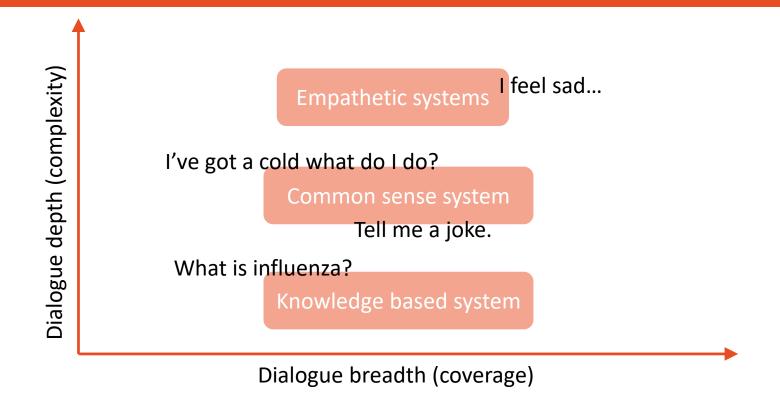




Try **Kisaku**, one of the best **sushi restaurants** in the city.



### **Evolution Roadmap**





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

- Task: user intent prediction
- Challenge: language ambiguity





- ✓ 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)

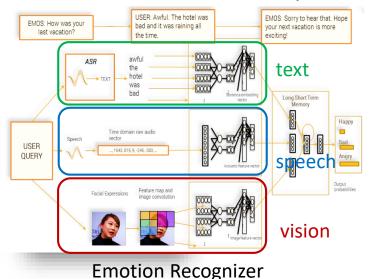
High-level intention may span several domains Schedule a lunch with Vivian. yelpas find restaurant check location contact play music What kind of restaurants do you prefer? The distance is ... Should I send the restaurant information to Vivian?

Users interact via <u>high-level descriptions</u> and the system learns <u>how to plan the dialogues</u>



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

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



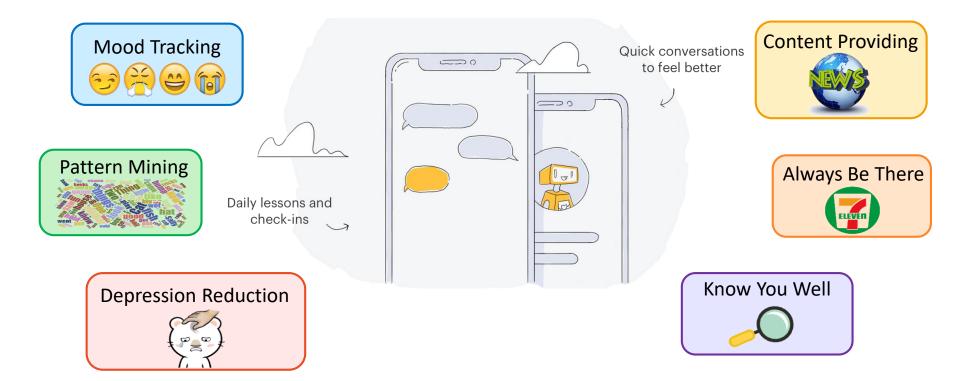
Zara - The Empathetic Supergirl







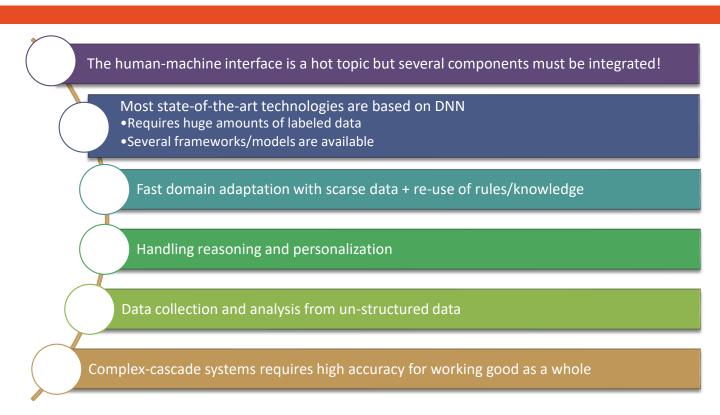
# Cognitive Behavioral Therapy (CBT)



# Challenges and Conclusions



### Challenge Summary





# Thanks for Your Attention!

Q&A



Yun-Nung (Vivian) Chen **Assistant Professor National Taiwan University** y.v.chen@ieee.org / http://vivianchen.idv.tw