

CS11-747 Neural Networks for NLP

Models of Dialog and Conversation

Graham Neubig



Carnegie Mellon University

Language Technologies Institute

Site

<https://phontron.com/class/anlp2022/>

What is Dialog?

- Understanding utterances, in the context of a conversation
- Generating responses
 - That are **consistent** and **coherent** with the dialog history
 - That are **interesting** and **engaging**
 - That **meaningfully** progress the dialog **towards a goal**

Types of Dialog

- Who is talking?
 - Human-human
 - Human-computer
- Why are they talking?
 - Task driven
 - Open-domain, chat

Open-Domain (Chat)

Two Paradigms

- Generation-based models
 - Take input, generate output
 - Good if you want to be creative
- Retrieval-based models
 - Take input, find most appropriate output
 - Good if you want to be safe

Problem 1: Dialog More Dependent on Global Coherence

- Considering only a single previous utterance will lead to locally coherent but globally incoherent output
- Necessary to consider more context! (Sordoni et al. 2015)



- Contrast to MT, where context sometimes is (Matsuzaki et al. 2015) and sometimes isn't (Jean et al. 2017) helpful

Problem 2: Dialog allows Much More Varied Responses

- For translation, there is lexical variation but content remains the same
- For dialog, content will also be different! (e.g. Li et al. 2016)

Input: What are you doing?	
-0.86 I don't know.	-1.09 Get out of here.
-1.03 I don't know!	-1.09 I'm going home.
-1.06 Nothing.	-1.09 Oh my god!
-1.09 Get out of the way.	-1.10 I'm talking to you.
Input: what is your name?	
-0.91 I don't know.	...
-0.92 I don't know!	-1.55 My name is Robert.
-0.92 I don't know, sir.	-1.58 My name is John.
-0.97 Oh, my god!	-1.59 My name's John.
Input: How old are you?	
-0.79 I don't know.	...
-1.06 I'm fine.	-1.64 Twenty-five.
-1.17 I'm all right.	-1.66 Five.
-1.17 I'm not sure.	-1.71 Eight.

Problem 3: Dialog Agents should have Personality

- If we train on all of our data, our agent will be a mish-mash of personalities (e.g. Li et al. 2016)

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?

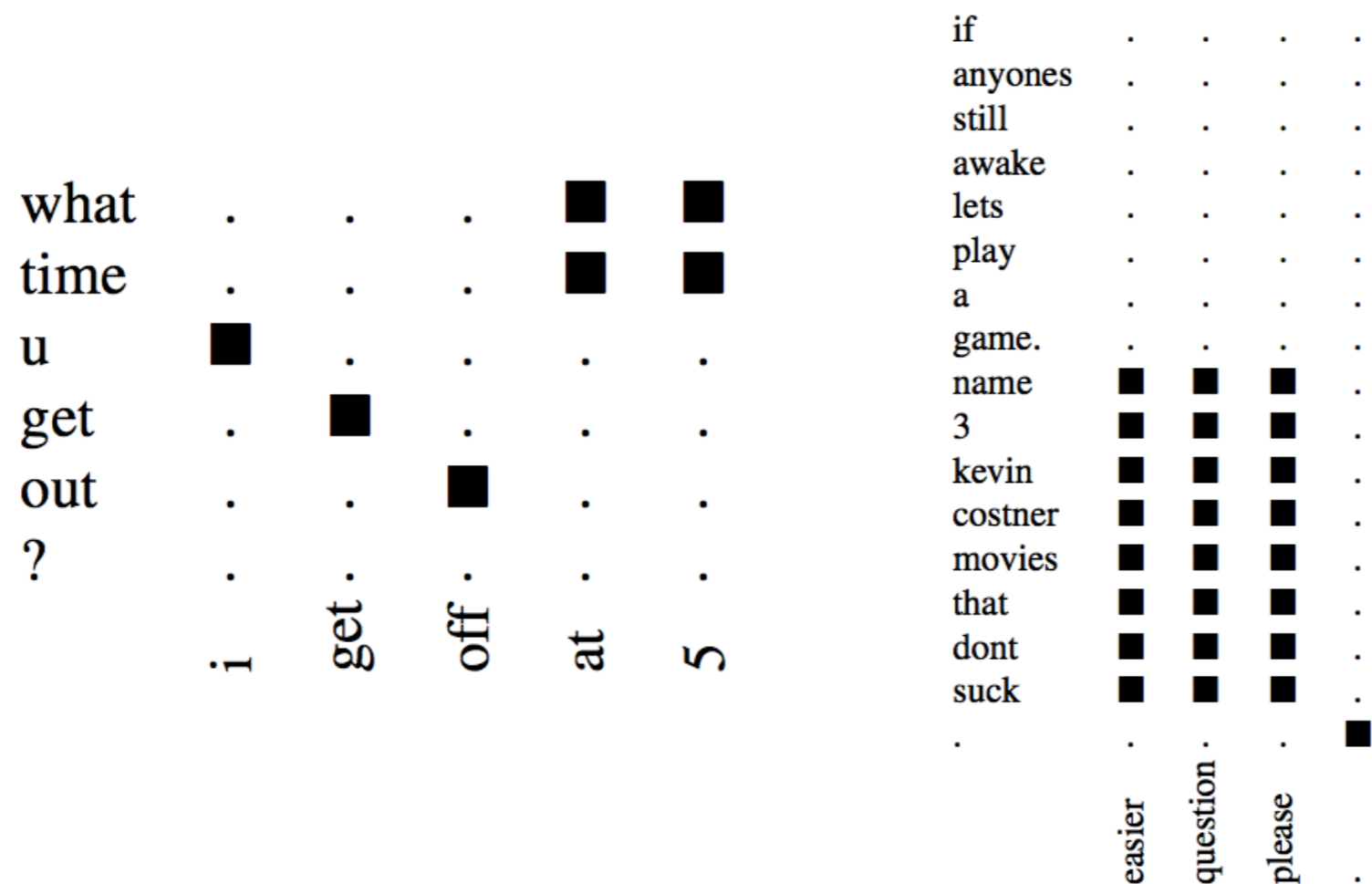
- We would like our agents to be consistent!

Generative Dialog Response Generation

Generation-based Models

(Ritter et al. 2011)

- Train model to "translate" from context to response
- Like other seq2seq tasks, dialog response generation can be done with encoder-decoders (Sordoni et al. 2015, Sheng et al. 2015, Vinyals and Le 2015)



DialoGPT [Zhang et al. 2019]

- Continue pre-training GPT-2 on conversations from Reddit
- Filter long utterances
- Filter non-English utterances
- Filter URLs
- Filter toxic comments
- Train on 147M dialog instances (1.8B words)
- “Human-level” response generation ability

Meena [Adiwardana et al. 2020]

Conversation C

Human: *Haha, I think mathematics has a presence beyond human existence.*

Meena: *So does philosophy.*

Human: *Can you state with examples how this is true.*

Meena: *Can you state how it isn't?*

Human: *Without human existence there will be no one left to comprehend what's right what's wrong, and thus philosophy will [sic]*

Meena: *That's not the definition of philosophy.*

Human: *Without human existence there will be no one left to comprehend what's right what's wrong, and thus philosophy will not be needed. How do you define philosophy?*

Meena: *The study of general knowledge or the fundamental nature of things.*

Conversation D

Human: *I like how cute cows are*

Meena: *Ah, interesting. You like how they are cute, right?*

Human: *yeah, and they're really smart*

Meena: *I heard they go to college*

Human: *cows go to college?*

Meena: *I heard that a cow went to Harvard.*

Human: *what did the cow study?*

Meena: *Bovine sciences.*

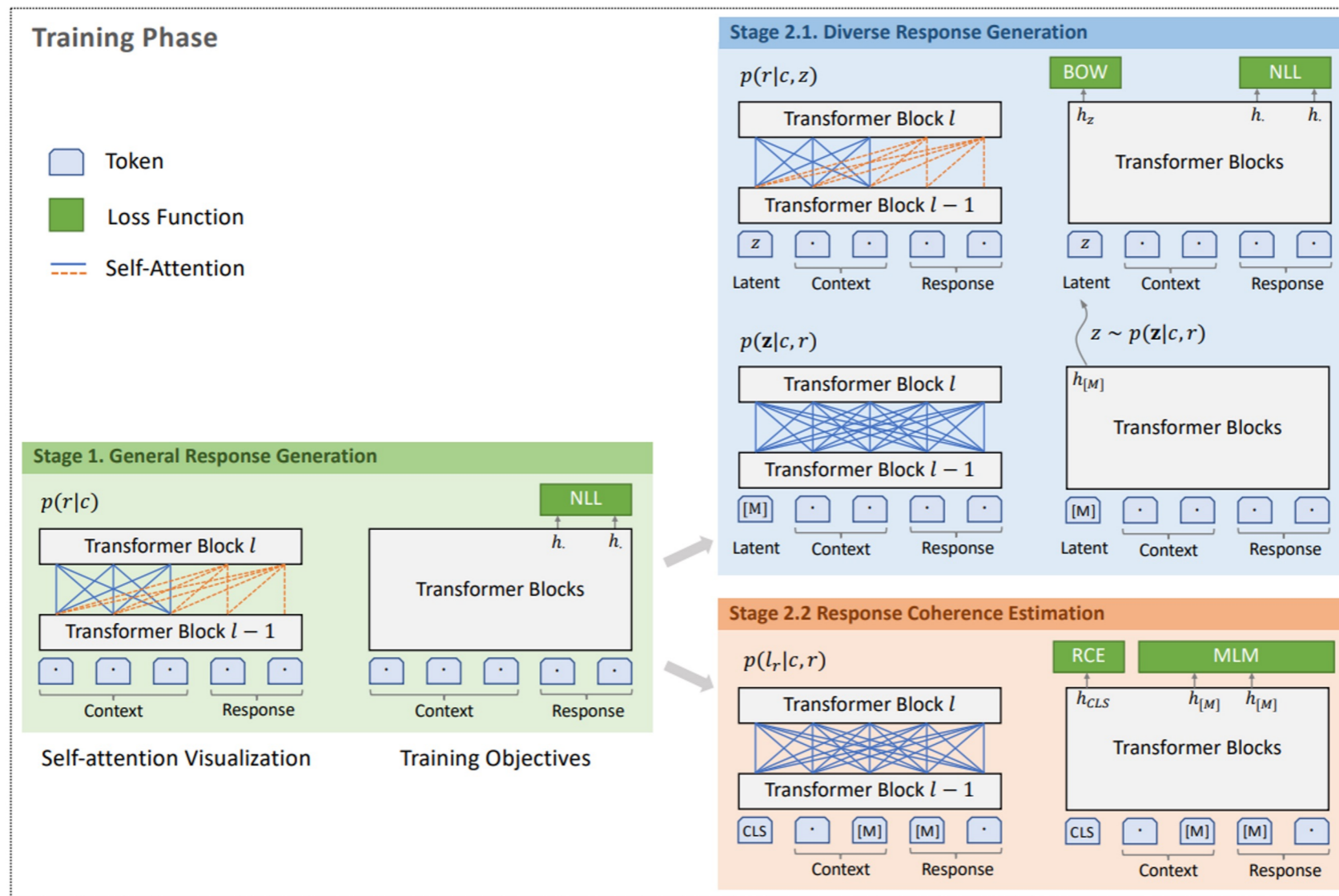
Human: *do horses go to Harvard?*

Meena: *Horses go to Hayvard.*

Human: *that's a pretty good joke, I feel like you led me into it*

Meena: *You were trying to steer it elsewhere, I can see it.*

PLATO-2 [Bao et al. 2021]



PLATO-2 [Bao et al. 2021]

User: Hello
System: hello! hows your day going so far?
User: What are you up to?
System: just finished making coffee! what about you? what do you do?
User: i'm playing video games
System: oh, nice! what kind? and i'm just relaxing here and listening to some classical music.
User: i'm playing league of legends
System: nice! what rank are you? i haven't played in forever.
User: who is your favorite classical musician?
System: probably beethoven. i like classical music.
User: i'm plat 2
System: nice! i'd imagine you are better than that!
User: can you play an instrument?
System: i can play the piano. do you play any instruments?
User: thanks
System: i appreciate it! so what is it like being plat 2?
User: i can play the clarinet, but not very well haha
System: oh nice! i played clarinet in high school.
User: the games are difficult but rewarding
System: yes they are! i played the trumpet in college.

...

Improving Variety and Personalization

Diversity Promoting Objective for Conversation (Li et al. 2016)

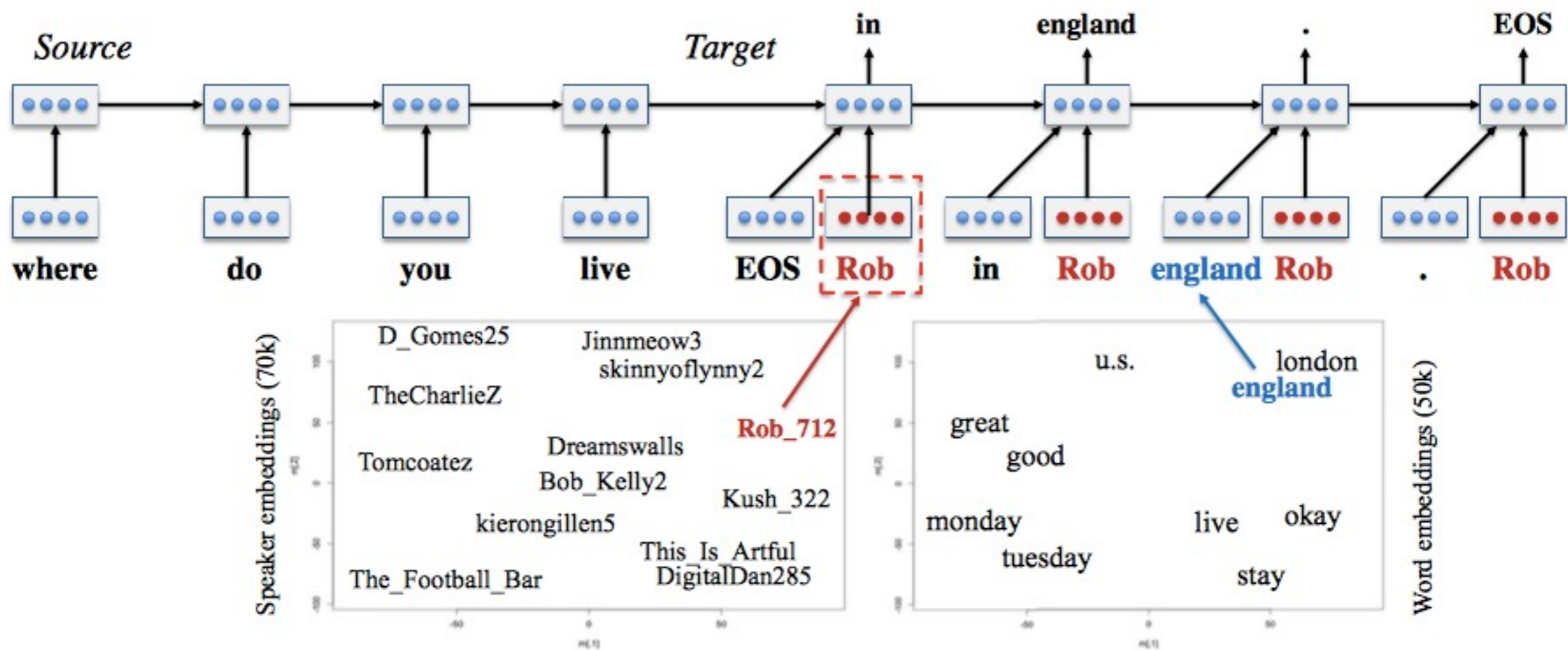
- Basic idea: we want responses that are likely given the context, unlikely otherwise
- Method: subtract weighted unconditioned log probability from conditioned probability (calculated only on first few words)

$$\log p(T|S) - \lambda \log U(T)$$

message	SEQ2SEQ	MMI $P(T S) - \lambda P(T)$
I did not get the report from an mi-6 agent.	I don't think that's a good idea.	You did the right thing, did you?
You haven't been given an assignment in this case.	I don't know what you are talking about.	I've been looking all over for you.
I'm losing my grip.	I don't know what you are talking about.	I'm the only one in the world.
I am ready to help.	Come on, come on	I have something we need to talk about.
You programmed me to gather intelligence. That's all I've ever done.	You know that, don't you?	You do have fun, don't you?
I mean, we'd have to talk to him.	I mean, I don't know.	I mean, he's a good guy.
This has nothing to do with him.	You don't know anything about him, do you?	I know him as much.
How come you never say it?	I don't know	Because I don't want to hurt you
He is the one behind this scheme to crash x-net.	I don't know what you are talking about	Yes, he is, isn't he?

Persona-based Neural Dialog Models (Li et al. 2017)

- Model each speaker in embedding space



- Also model who the speaker is speaking to in speaker-addressee model

Prompting-based Methods to Personalization (Zhang et al. 2018)

- Express information about the speaker in text

Persona 1	Persona 2
I like to ski My wife does not like me anymore I have went to Mexico 4 times this year I hate Mexican food I like to eat cheetos	I am an artist I have four children I recently got a cat I enjoy walking for exercise I love watching Game of Thrones

[PERSON 1:] Hi

[PERSON 2:] Hello ! How are you today ?

[PERSON 1:] I am good thank you, how are you.

[PERSON 2:] Great, thanks ! My children and I were just about to watch Game of Thrones.

[PERSON 1:] Nice ! How old are your children?

[PERSON 2:] I have four that range in age from 10 to 21. You?

[PERSON 1:] I do not have children at the moment.

[PERSON 2:] That just means you get to keep all the popcorn for yourself.

[PERSON 1:] And Cheetos at the moment!

[PERSON 2:] Good choice. Do you watch Game of Thrones?

[PERSON 1:] No, I do not have much time for TV.

[PERSON 2:] I usually spend my time painting: but, I love the show.

Personality Infused Dialog

(Mairesse et al. 2007)

- Train a generation system with controllable “knobs” based on personality traits
- e.g. Extraversion:
- Non-neural, but well done and perhaps applicable

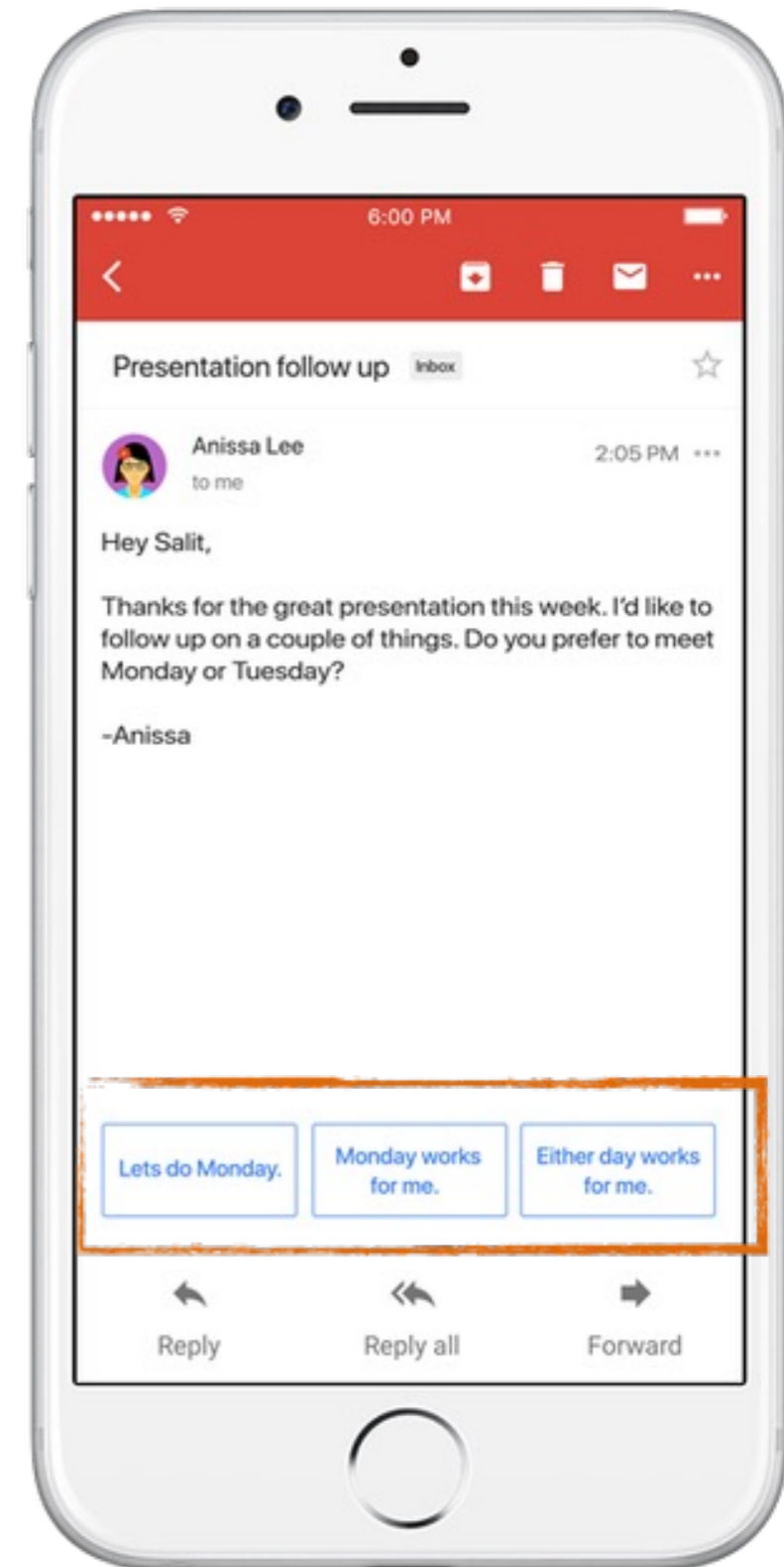
Alt	Realization	Extra
5	Err... it seems to me that Le Marais isn't as bad as the others.	1.83
4	Right, I mean, Le Marais is the only restaurant that is any good.	2.83
8	Ok, I mean, Le Marais is a quite french, kosher and steak house place, you know and the atmosphere isn't nasty, it has nice atmosphere. It has friendly service. It seems to me that the service is nice. It isn't as bad as the others, is it?	5.17
9	Well, it seems to me that I am sure you would like Le Marais. It has good food, the food is sort of rather tasty, the ambience is nice, the atmosphere isn't sort of nasty, it features rather friendly servers and its price is around 44 dollars.	5.83
3	I am sure you would like Le Marais, you know. The atmosphere is acceptable, the servers are nice and it's a french, kosher and steak house place. Actually, the food is good, even if its price is 44 dollars.	6.00
10	It seems to me that Le Marais isn't as bad as the others. It's a french, kosher and steak house place. It has friendly servers, you know but it's somewhat expensive, you know!	6.17
2	Basically, actually, I am sure you would like Le Marais. It features friendly service and acceptable atmosphere and it's a french, kosher and steak house place. Even if its price is 44 dollars, it just has really good food, nice food.	6.17

Retrieval-based Models

Dialog Response Retrieval

- Idea: many things can be answered with template
- Simply find most relevant response out of existing ones in corpus

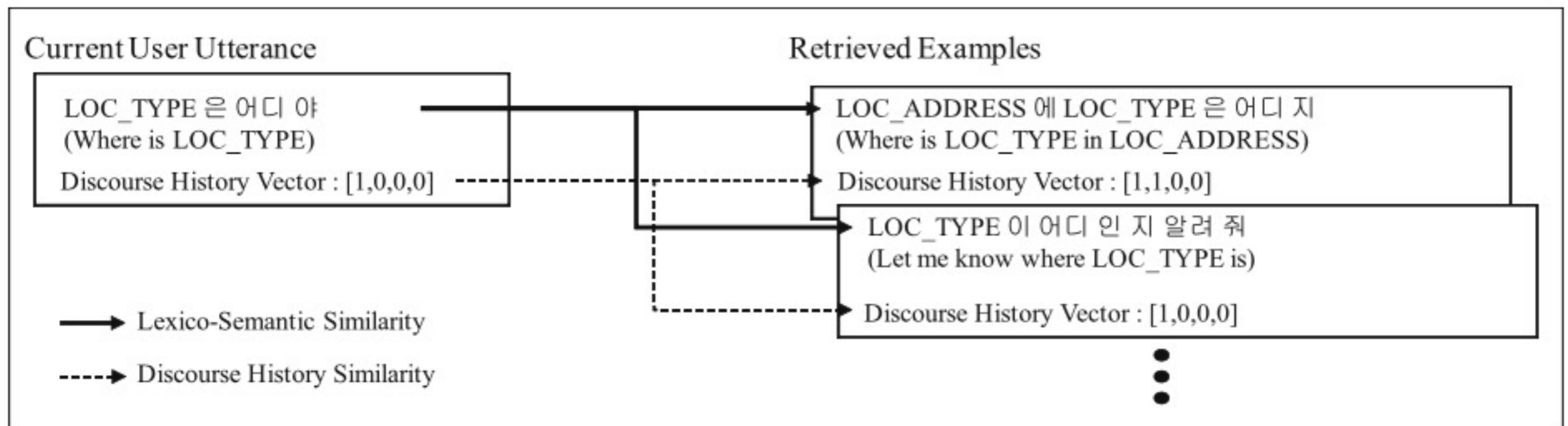
Template responses



Retrieval-based Chat

(Lee et al. 2009)

- Basic idea: given an utterance, find the most similar in the database and return it

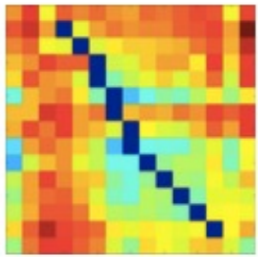
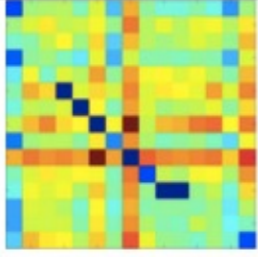
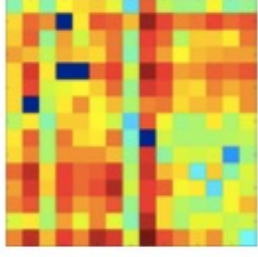


- Similarity based on exact word match, plus extracted features regarding discourse

Neural Response Retrieval

(Nio et al. 2014)

- Idea: use neural models to soften the connection between input and output and do more flexible matching

<i>sim</i>	Sentences	Matrix
0.94	S_1) Captain, we can not keep going fast on these icy roads. S_2) We can not keep going fast on these icy roads!	 A heatmap matrix showing a strong diagonal correlation between the two sentences, indicating high similarity. The diagonal elements are dark blue, while the off-diagonal elements are yellow and red.
0.60	S_1) Hold your fire! He's got a girl. S_2) Looks like he's got a hostage.	 A heatmap matrix showing a moderate correlation between the two sentences. The diagonal elements are dark blue, and there are some yellow and red areas, indicating a more complex relationship.
0.38	S_1) Yes, I can see that too and I don't think it's so terrible. S_2) That's why I do all the thinking.	 A heatmap matrix showing a weak correlation between the two sentences. The diagonal elements are dark blue, and there are many yellow and red areas, indicating a low similarity.

- Model uses Socher et al. (2011) recursive auto-encoder + dynamic pooling

Smart Reply for Email Retrieval (Kannan et al. 2016)

- Implemented in GMail smart reply
- Response model with seq2seq scoring, but many improvements
 - **Beam search** over response space for scalability
 - **Canonicalization** of syntactic variants and **clustering** of similar responses
 - **Human curation** of responses
 - Enforcement of diversity through **omission of redundant responses** and **enforcing positive/negative**

Open-domain Dialog Evaluation

Dialog Evaluation

Goal: Construct automatic evaluation metrics for response generation/interactive dialog

27

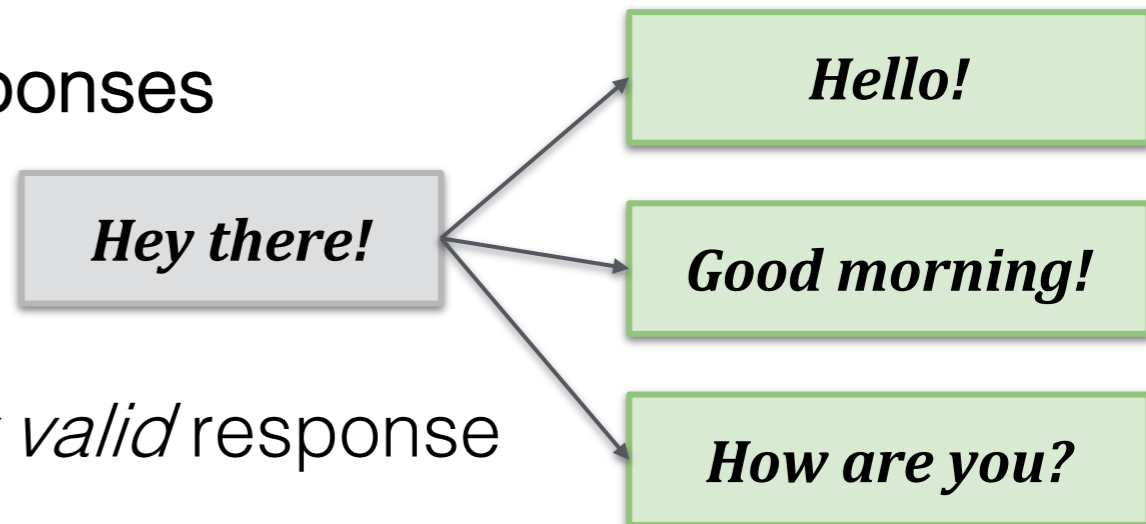
Given: dialog history, generated response, reference response (optional)

Output: a score for the response

Why is evaluating dialog hard? (1/3)

1. One-to-many nature of dialog

- For each dialog there are many valid responses
- Cannot compare to a reference response
 - The reference response isn't the *only valid* response
- Existing metrics won't work
 - BLEU, F-1, etc.



Why is evaluating dialog hard? (2/3)

2. Dialog quality is multi-faceted

- A response isn't just good or bad
- For interpretability, should measure ²⁹multiple qualities
 - Relevance
 - Interestingness
 - Fluency

Why is evaluating dialog hard? (2/3)

3. Dialog is inherently interactive

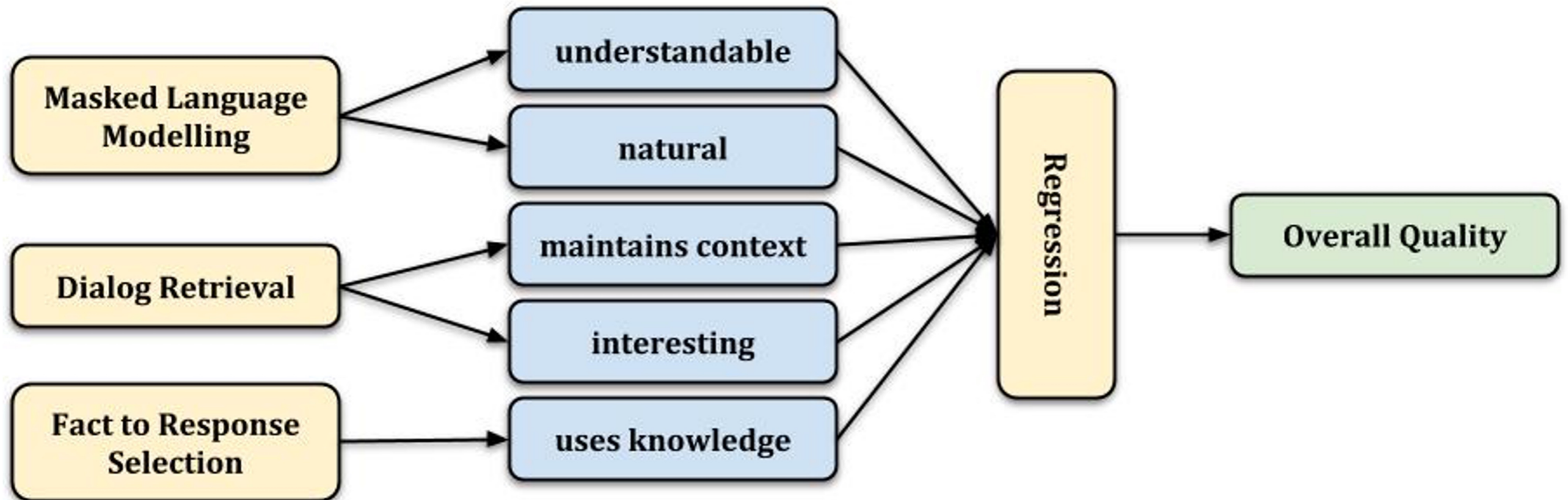
- Dialog systems are designed to have a **back-and-forth interaction** with a user
 - Research largely focuses on **static corpora**³⁰ → Reduces the problem of dialog to **response generation**
- Some properties of a system can't be assessed outside an **interactive environment**
 - Long-term planning, error recovery, coherence.

Dialog Evaluation Metrics

- USR [Mehri and Eskenazi. 2020]
- GRADE [Huang et al. 2020]
- HolisticEval [Pang et al. 2020]
- DSTC6 [Hori and Hori. 2017]
- FED [Mehri and Eskenazi. 2020]
- DSTC9 [Gunasekara et al. 2021]

31

USR [Mehri and Eskenazi. 2020]



Using Multiple References with Human Evaluation Scores (Galley et al. 2015)

- Retrieve good-looking responses, perform human evaluation, up-weight good ones, down-weight bad ones

Context c	Message m	Response r	Score
i was about to text you and my two cousins got excited cause they thought you were "rihanna"	aww, i can imagine their disappointment	they were very disappointed!!!	0.6
<i>yes. my ex-boyfriend, killed my cat. like i say, it was the start of a bad time...</i>	<i>i can imagine!</i>	<i>yes. luckily, the whole thing feels very much of the past now.</i>	0.8
<i>its good.. for some reason i can't name stand out tracks but i've been playing it since it dropped</i>	<i>i can imagine, banks doesn't disappoint</i>	<i>na this is anything but a disappointment..</i>	0.6
<i>at my lil cousins dancing to "dance for you". these kids are a mess.</i>	<i>lmaoo i can imagine.</i>	<i>they were belly rolling, filarious.</i>	0.4
<i>what's sick about it?? do you know how long it is?? no so how is it sick?</i>	<i>i can imagine it</i>	<i>your imagination is wrong, very wrong at that.</i>	-0.1
<i>it's saad oh yeah the snow is very beautiful</i>	<i>yeah i can imagine</i>	<i>the weather in russia is very cool.</i>	-0.7

Task-driven Dialog

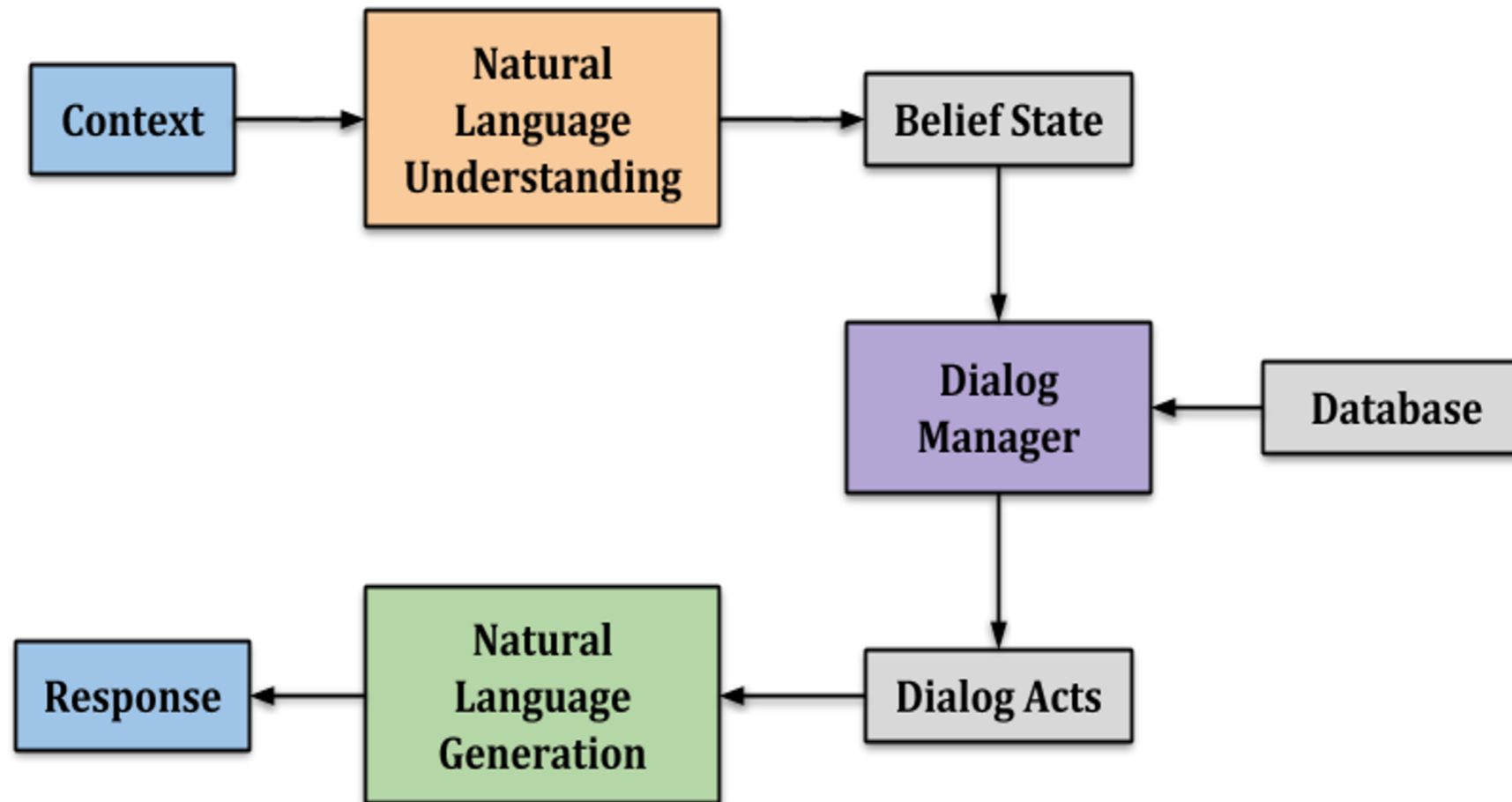
Chat vs. Task Oriented

- Chat is basically to keep the user entertained
- What if we want to do an actual task?
 - Book a flight
 - Access information from a database

Task-oriented Dialog Framework

- In semantic frame based dialog:
- **Natural language understanding** to fill the slots in the frame based on the user utterance
- **Dialog state tracking** to keep track of the overall dialog state over multiple turns
- **Dialog control** to decide the next action based on state
- **Natural language generation** to generate utterances based on current state

Pipeline Dialog System



Natural Language Understanding

Natural language understanding in dialog involves several key tasks:

- DialoGLUE [Mehri et al. 2020]
- Intent prediction: ATIS, SNIPS, Banking77, CLINC150, HWU64
- Slot filling: ATIS, SNIPS, DSTC8-SGD, Restaurant8k
- State tracking: MultiWOZ (2.X)

NLU (for Slot Filling) w/ Neural Nets (Mesnil et al. 2015)

- Slot filling expressed as BIO scheme

Sentence	<i>show</i>	<i>flights</i>	<i>from</i>	<i>Boston</i>	<i>To</i>	<i>New</i>	<i>York</i>	<i>today</i>
Slots/Concepts	O	O	O	B-dept	O	B-arr	I-arr	B-date
Named Entity	O	O	O	B-city	O	B-city	I-city	O
Intent	<i>Find Flight</i>							
Domain	<i>Airline Travel</i>							

- RNN-CRF based model for tags

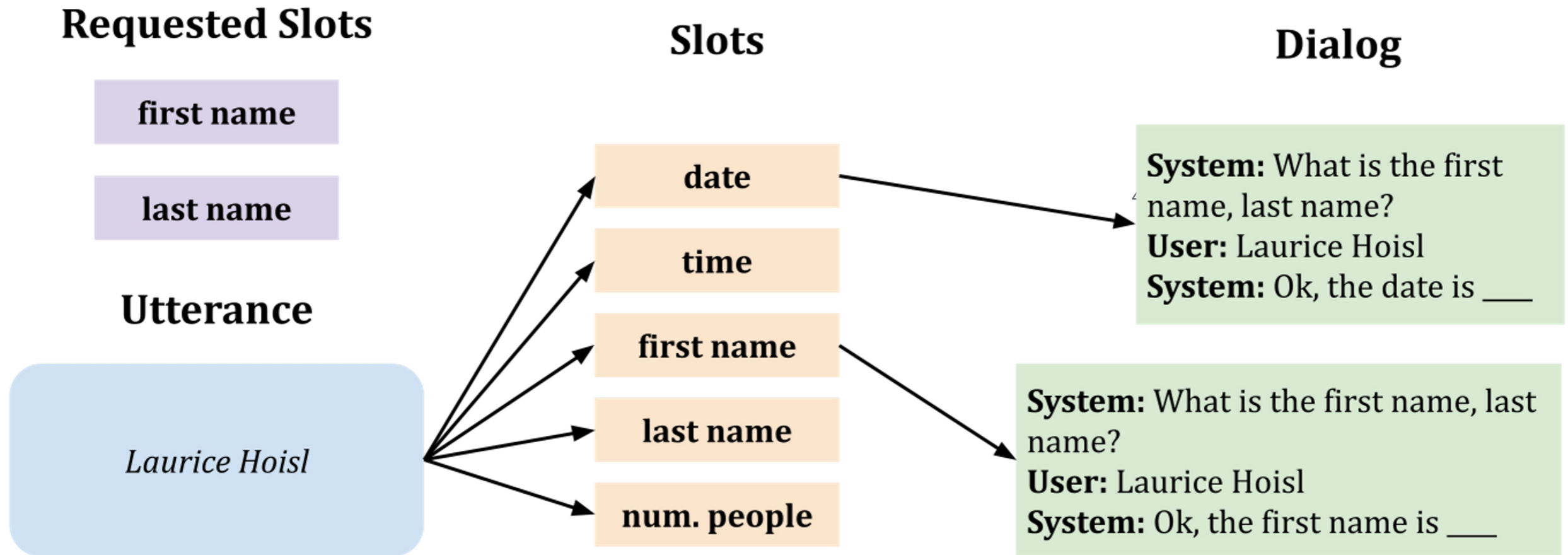
ConVEx [Henderson and Vulic. 2020]

Pre-training paradigm specifically for slot filling → strong few-shot/zero-shot performance

<i>Template Sentence</i>	<i>Input Sentence</i>
I get frustrated everytime I browse /r/all. I stick to my <i>BLANK</i> most of the time.	/r/misleadingpuddles Saw it on the frontpage , plenty of content if you like the premise.
Why Puerto Rico? It's Memphis at Dallas, which is in Texas where <i>BLANK</i> hit	Hurricane Harvey . Just a weird coincidence.
<i>BLANK</i> is my 3rd favorite animated Movie	Toy Story 3 ended perfectly, but Disney just wants to keep milking it.
It really sucks, as the V30 only has <i>BLANK</i> . Maybe the Oreo update will add this.	Thanks for the input, but 64GB is plenty for me :)
I took <i>BLANK</i> , cut it to about 2 feet long and duct taped Vive controllers on each end. Works perfect	Yeah, I just duct taped mine to a broom stick . You can only play no arrows mode but it's really fun.
I had <i>BLANK</i> and won the last game and ended up with 23/20 and still didn't get it.	I know how you feel my friend and I got 19/20 on the tournament today

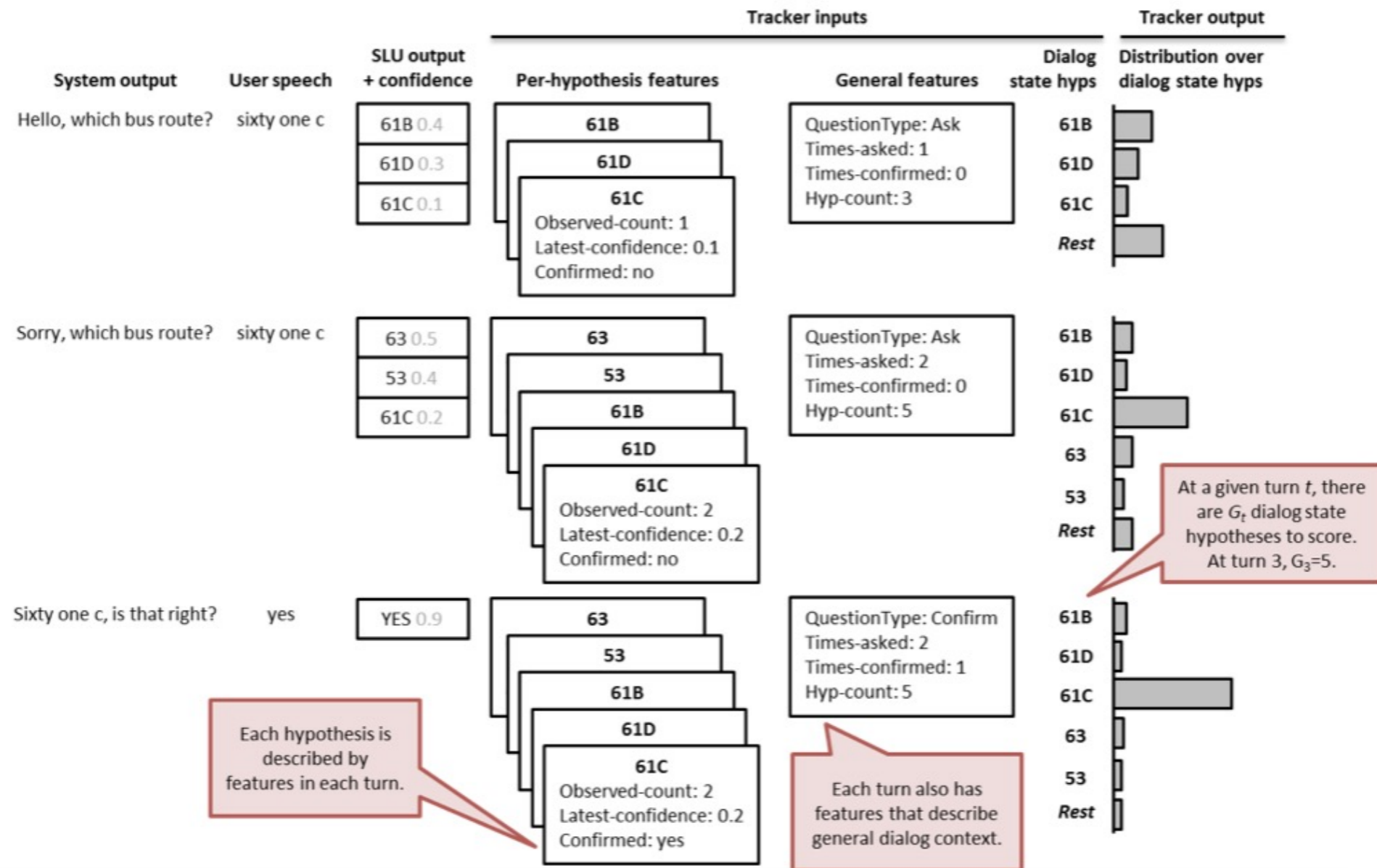
Table 1: Sample data from Reddit converted to sentence pairs for the ConVEx pretraining via the pairwise cloze task. Target spans in the input sentence are denoted with bold, and are “*BLANKed*” in the template sentence.

GenSF [Mehri and Eskenazi. 2021]



Dialog State Tracking

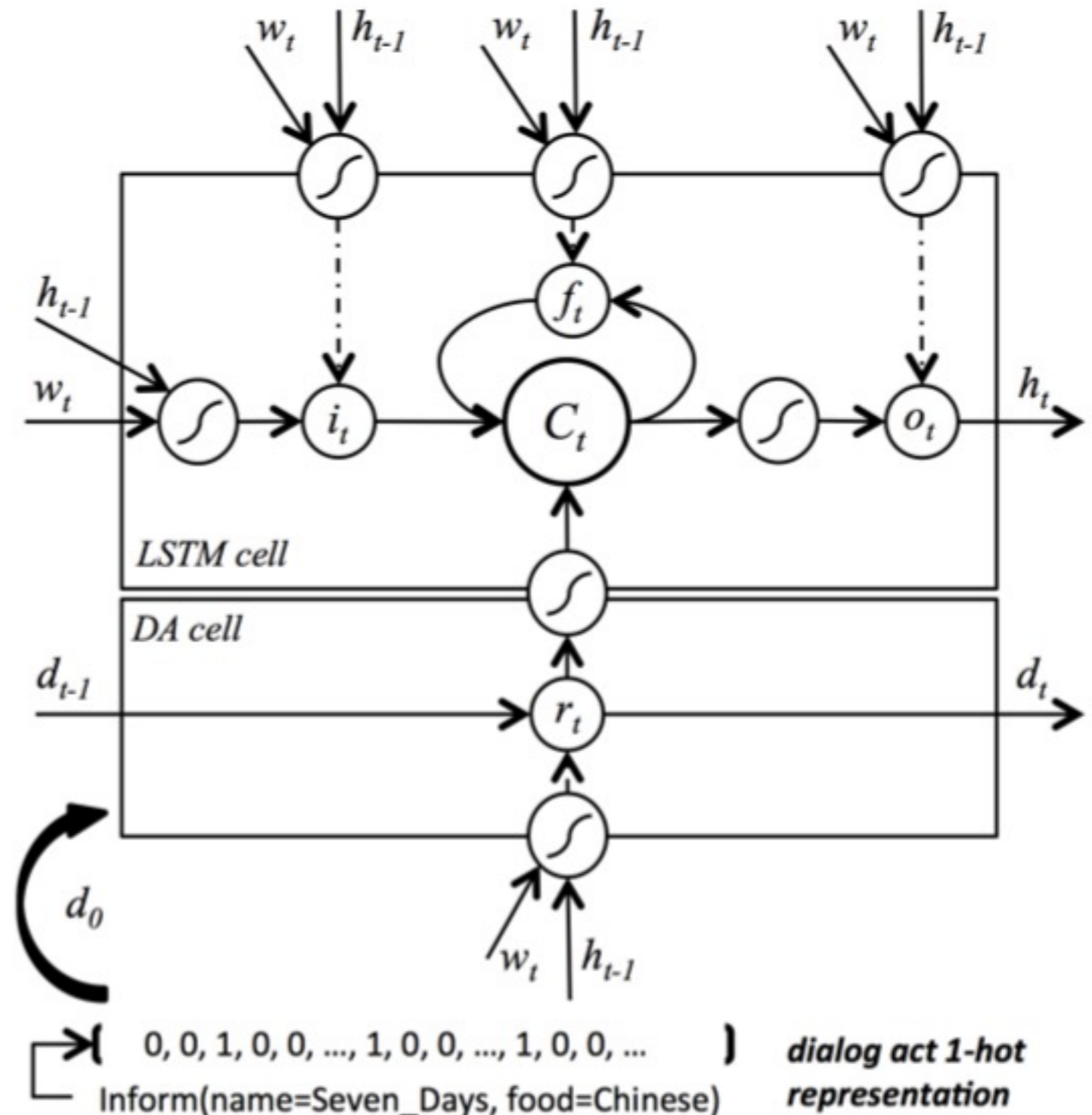
- Track the belief about our current frame-filling state (Williams et al. 2013)



- Henderson et al. (2014) present RNN model that encodes multiple ASR hypotheses and generalizes by abstracting details

Language Generation from Dialog State w/ Neural Nets (Wen et al. 2015)

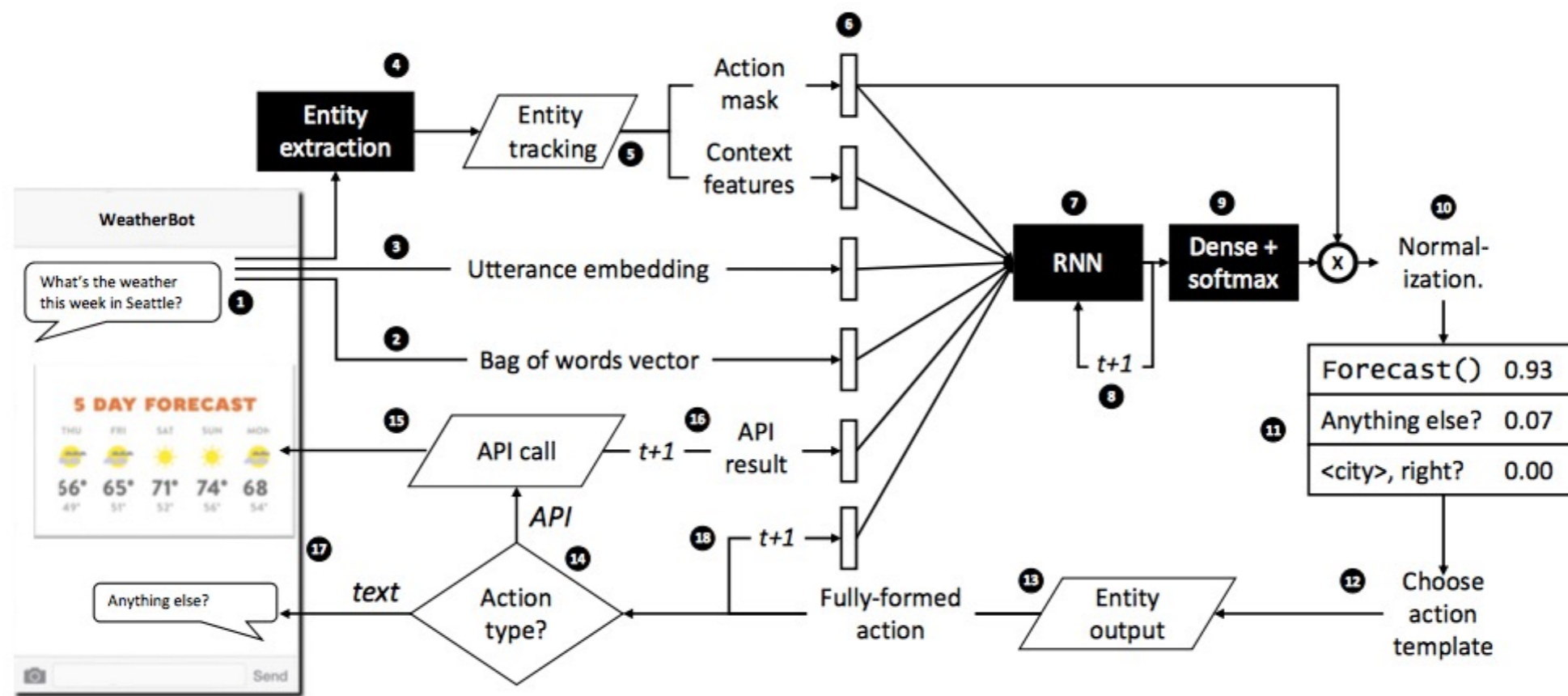
- Condition LSTM units based on the dialog input, output English



End-to-end Dialog Control

(Williams et al. 2017)

- Train an LSTM that takes in text and entities and directly chooses an action to take (reply or API call)



- Trained using combination of supervised and reinforcement learning

Questions?