

# Dialogue Datasets

CS294s: BUILDING THE BEST VIRTUAL ASSISTANT

Ryan Kearns & Lucas Sato

Mentor: Giovanni Campagna

May 14, 2020

# Outline

## 1. Introduction: Why Datasets?

## 2. MultiWOZ in the Almond/ThingTalk/Genie Context

## 3. What's In a Dataset

- a. Dialogue Generation
- b. Annotation Generation
- c. Annotation Styles

## 4. MultiWOZ Revisited

# 1. Why Datasets?

“Perhaps the most important news of our day is that datasets—not algorithms—might be the key limiting factor to development of human-level artificial intelligence.”

- Alexander Wissner-Gross, 2016  
Harvard University Institute for Applied Computational Science

# Outline

1. Introduction: Why Datasets?
- 2. MultiWOZ in the Almond/ThingTalk/Genie Context**
3. What's In a Dataset
  - a. Dialogue Generation
  - b. Annotation Generation
  - c. Annotation Styles
4. MultiWOZ Revisited

## 2. MultiWOZ in the Almond/ThingTalk/Genie Context

**User:** I need to book a hotel in the east  
that has four stars

Hotel-Rating: 4, Hotel-Location: East DST

**Agent:** What is your price range?

**User:** Price doesn't matter as long  
as it has free wifi.

Hotel-Rating: 4, Hotel-Location: East, DST  
Hotel-Wifi: True, Hotel-Price: dontcare

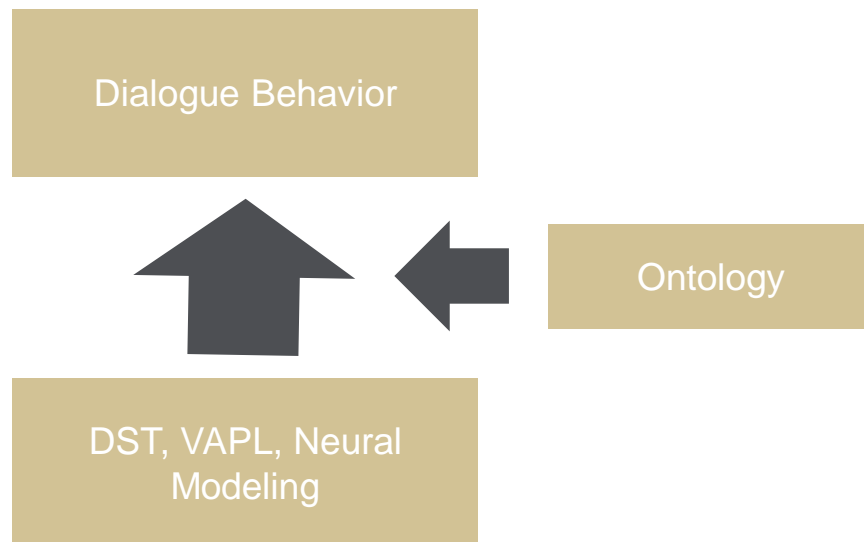
**Agent:** In that case, I would

**User:** Thanks. Please get me a taxi recommend Allenbell.  
from here to the hotel.

Hotel-Rating: 4, Hotel-Location: East, DST  
Hotel-Wifi: True, Hotel-Price: dontcare  
Hotel-Name: Allenbell  
Taxi-Departure: Home, Taxi Destination: Allenbell

## 2. MultiWOZ in the Almond/ThingTalk/Genie Context

- MultiWOZ (and most datasets) has a corpus and annotations.
- We personally only use the former. We don't train on MultiWOZ.

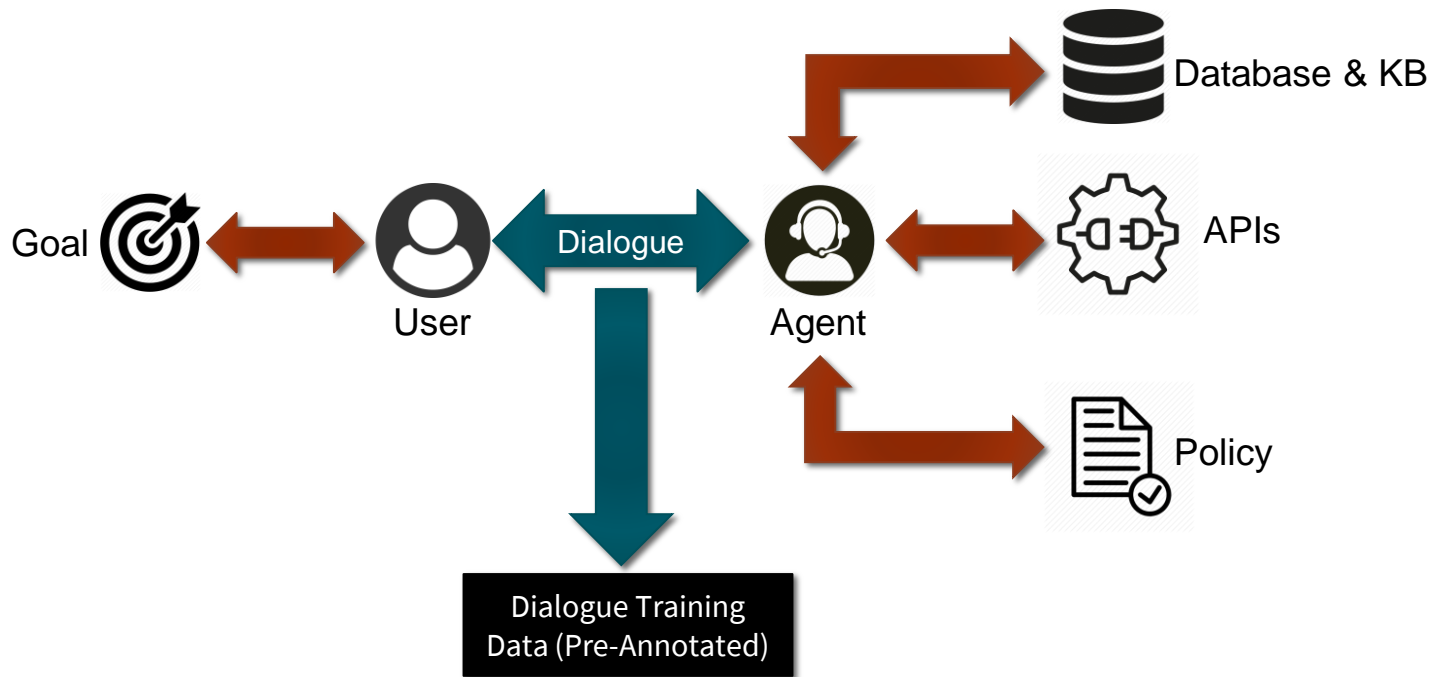


# Outline

1. Introduction: Why Datasets?
2. MultiWOZ in the Almond/ThingTalk/Genie Context
- 3. What's In a Dataset**
  - a. Dialogue Generation
  - b. Annotation Generation
  - c. Annotation Styles
4. MultiWOZ Revisited

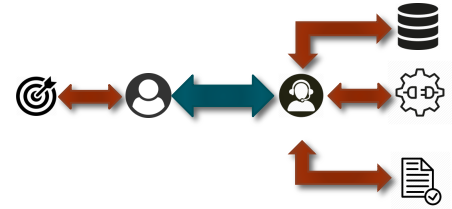
## 3a. Dialogue Generation

Our General Paradigm:











# Human-to-Machine



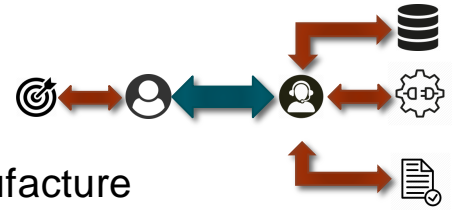
Bootstrap from an existing dialogue system to build a new task-oriented dialogue corpora.

Example: Let's Go Bus Information System, used for the first Dialogue State Tracking Challenge (DSTC)

-  **User:** real humans interacting with the dialogue system
-  **Agent:** existing dialogue system, likely following rigid rule-based dialogue policy
-  **Goal:** derived from existing dialogue system
-  **Database / KB:** derived from existing dialogue system
-  **APIs:** derived from existing dialogue system
-  **Policy:** derived from existing dialogue system







*Great for expanding the capabilities of an existing domain, but can we generalize beyond this domain?*

# Machine-to-Machine



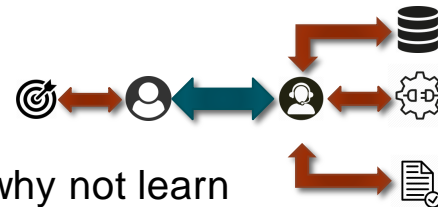
Engineer a simulated user plus a transaction environment to manufacture dialogue templates en masse, then map those dialogue templates to natural language.

Example: Shah et al., 2018, “a framework combining automation and crowdsourcing to rapidly bootstrap end-to-end dialogue agents for goal-oriented dialogues”

-  **User:** engineered, agenda-based simulator
-  **Agent:** engineered, likely from a finite-state machine
-  **Goal:** derived from scenarios produced by Intent+Slots task schema
-  **Database / KB:** domain-specific, wrapped into API client
-  **APIs:** provided by developer
-  **Policy:** engineered specifically for agent

*Great for exhaustively exploring the space of possible dialogues, but will the training data actually match real-world scenarios?*

# Human-to-Human



If we really want our agents mimicking human dialogue behavior, why not learn from real human conversations?

Example: Twitter dataset (*Ritter et al., 2010*), Reddit conversations (*Schradling et al., 2015*), Ubuntu technical support corpus (*Lowe et al., 2015*)

 **User:** real humans on the Internet

 **Agent:** real humans on the Internet

 **Goal:** ???

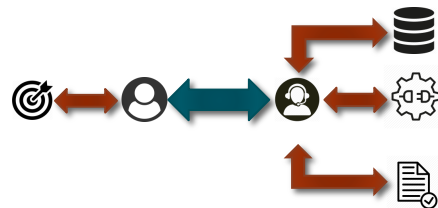
 **Database / KB:** ???

 **APIs:** ???

 **Policy:** real human dialogue policies!

*Great for teaching a system real human dialogue patterns, but how will we ground dialogues to the KB + API required by our dialogue agent?*


# Human-to-Human (WOZ)




Humans produce the best dialogue behavior. Let's use humans to *simulate* a machine dialogue agent, grounding the dialogue in our KB+APIs.

Example: WOZ2.0 (Wen et al., 2017), FRAMES (El Asri et al., 2017), MultiWOZ{1.0, 2.0, 2.1} (Budzianowski et al., 2018)


 **User:** crowdworker

 **Agent:** crowdworker, simulating a human-quality dialogue system

 **Goal:** provided by the task description

 **Database / KB:** domain-specific, provided to the agent by experimenters

 **APIs:** domain-specific, provided to the agent by experimenters

 **Policy:** up to the crowdworker – nuanced, but maybe idiosyncratic

*Great for combining human dialogue policies with grounding in the specific transaction domain, but annotations will be nontrivial – how do we ensure their correctness?*

# Dialogue Generation – Summary

## Human-to-Machine

Bootstrap from an existing dialogue system to build a new task-oriented dialogue corpora.

## Human-to-Human

If we really want our agents mimicking human dialogue behavior, why not learn from real human conversations?

## Machine-to-Machine

Engineer a simulated user plus a transaction environment to manufacture dialogue templates en masse, then map those dialogue templates to natural language.

## Human-to-Human (WOZ)

Humans produce the best dialogue behavior. Let's use humans to *simulate* a machine dialogue agent, grounding the dialogue in our KB+APIs.

# Dialogue Generation – Pros & Cons

## Human-to-Machine

- + Intuitive to use existing dialogue data for dialogue system development
- Only possible to improve existing, working systems. No generalizations to new domains
- Initial system's capacities & biases may encourage behaviors that perform in testing but don't generalize

## Human-to-Human

- + Training data will map directly onto real-world interactions
- No grounding in any existing knowledge base or API limits usability

## Machine-to-Machine

- + Full coverage of all dialogue outcomes in domain
- Naturalness of the dialogue mismatches with real interactions
- Hard to simulate noisy conditions typical of real interactions

## Human-to-Human (WOZ)

- + Ground realistic human dialogue within the capacities of the dialogue system
- High prevalence of misannotation errors

## Question

WHICH DIALOGUE  
GENERATION TECHNIQUE  
SEEMS MOST SUITED FOR  
YOUR OWN PROJECT'S  
DOMAIN?



# Outline

1. Introduction: Why Datasets?
2. MultiWOZ in the Almond/ThingTalk/Genie Context
- 3. What's In a Dataset**
  - a. Dialogue Generation
  - **b. Annotation Generation**
  - c. Annotation Styles
4. MultiWOZ Revisited



## 3b. Annotation generation

### "Built-in" annotations (Machine-generated utterances)

- If the utterance is machine-generated, that it probably already has a formal language annotation
  - Annotation is not really separate from the dialogue generation
  - WikiSQL [Zhong et al. 2017]
- + Only skill needed is paraphrasing
- Still less natural and diverse
  - Requires good utterance synthesis



## 3b. Annotation generation

### Manual annotations (Human-generated utterances)

- Annotation as an explicit step in the process
- Usually done on top of provided data, possibly as a separate process
- Spider [Yu et al. 2019]

- + The dataset and the annotations are probably pretty good
- Potentially very expensive (experts often required)
- Sometimes not actually very good



## 3b. Annotation generation

### Machine-assisted annotations (Human-generated utterances)

- Technology used in making the annotation process seamless or easier for humans
  - Not necessarily a separate step in the process
  - QA-SRL [He et al. 2015]
- + The dataset and the annotations are probably pretty good
- Some upfront cost of developing a good system
  - Not always possible



## Question

HOW DO YOU THINK  
MACHINE-ASSISTED  
ANNOTATION COULD WORK  
IN YOUR PARTICULAR  
PROJECT?



# Outline

1. Introduction: Why Datasets?
2. MultiWOZ in the Almond/ThingTalk/Genie Context
- 3. What's In a Dataset**
  - a. Dialogue Generation
  - b. Annotation Generation
  - **c. Annotation Styles**
4. MultiWOZ Revisited

# A Fundamental Tradeoff

**Expressiveness of  
your representation**

vs.

**Ease of parsing,  
annotation, and execution**



## 3c. Annotation styles

Key Tradeoff: **expressiveness of the representation** vs. **ease of annotation/parsing/execution**

- Logical forms [Zettlemoyer & Collins, 2012; Wang et al. 2015]
- Intent and slot tagging [Goyal et al., 2017; Rastogi et al., 2020; many others...]
- Hierarchical representations [Gupta et al., 2018]
- Executable representations
  - SQL [Zhong et al., 2017; Yu et al., 2019]
  - ThingTalk [Campagna et al., 2019]

# Logical forms

Zettlemoyer & Collins, 2012; Wang et al. 2015

Rigid logical formalisms for queries results in a **precise, machine-learnable, and brittle** representation.

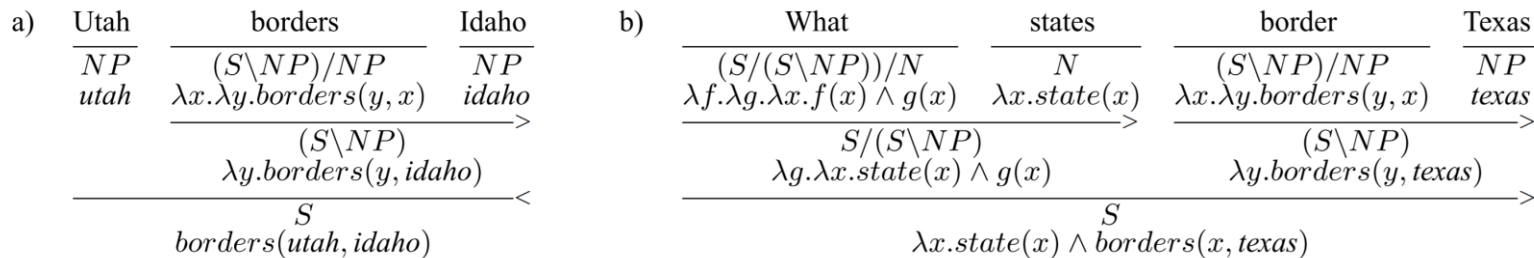


Figure 2: Two examples of CCG parses.



# Intent and slot tagging

*Goyal et al., 2017; Rastogi et al., 2020; many others...*

More ubiquitous, less expert-reliant representation allows coverage of more possible dialogue states.

Table 2: Full ontology for all domains in our data-set. The upper script indicates which domains it belongs to. \*: universal, 1: restaurant, 2: hotel, 3: attraction, 4: taxi, 5: train, 6: hospital, 7: police.

act type	inform* / request* / select <sup>123</sup> / recommend/ <sup>123</sup> / not found <sup>123</sup> request booking info <sup>123</sup> / offer booking <sup>1235</sup> / inform booked <sup>1235</sup> / decline booking <sup>1235</sup> welcome* / greet* / bye* / reqmore*
slots	address* / postcode* / phone* / name <sup>1234</sup> / no of choices <sup>1235</sup> / area <sup>123</sup> / pricerange <sup>123</sup> / type <sup>123</sup> / internet <sup>2</sup> / parking <sup>2</sup> / stars <sup>2</sup> / open hours <sup>3</sup> / departure <sup>45</sup> destination <sup>45</sup> / leave after <sup>45</sup> / arrive by <sup>45</sup> / no of people <sup>1235</sup> / reference no. <sup>1235</sup> / trainID <sup>5</sup> / ticket price <sup>5</sup> / travel time <sup>5</sup> / department <sup>7</sup> / day <sup>1235</sup> / no of days <sup>123</sup>

Figure from MultiWOZ (Budzianowski et al., 2018)

# Hierarchical Annotations

*Gupta et al., 2018*

Nesting additional intents within slots allows for **function composition & nested API calls**.

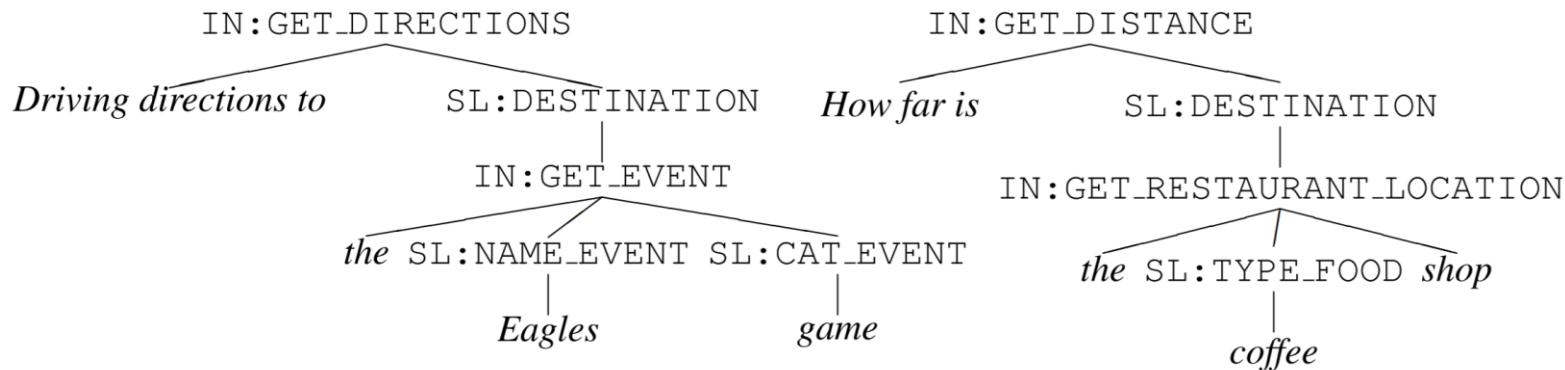


Figure 1: Example TOP annotations of utterances. Intents are prefixed with IN: and slots with SL:. In a traditional intent-slot system, the SL:DESTINATION could not have an intent nested inside it.

# Executable Representations: SQL

Zhong et al., 2017; Yu et al., 2019

Structured nature of the SQL representation helps **prune the space of possibly generated queries**, simplifying the generation problem.

Table: CFLDraft

Pick #	CFL Team	Player	Position	College
27	Hamilton Tiger-Cats	Connor Healy	DB	Wilfrid Laurier
28	Calgary Stampeders	Anthony Forgone	OL	York
29	Ottawa Renegades	L.P. Ladouceur	DT	California
30	Toronto Argonauts	Frank Hoffman	DL	York
...	...	...	...	...

Question:

How many CFL teams are from York College?

SQL:

```
SELECT COUNT CFL Team FROM  
CFLDraft WHERE College = "York"
```

Result:

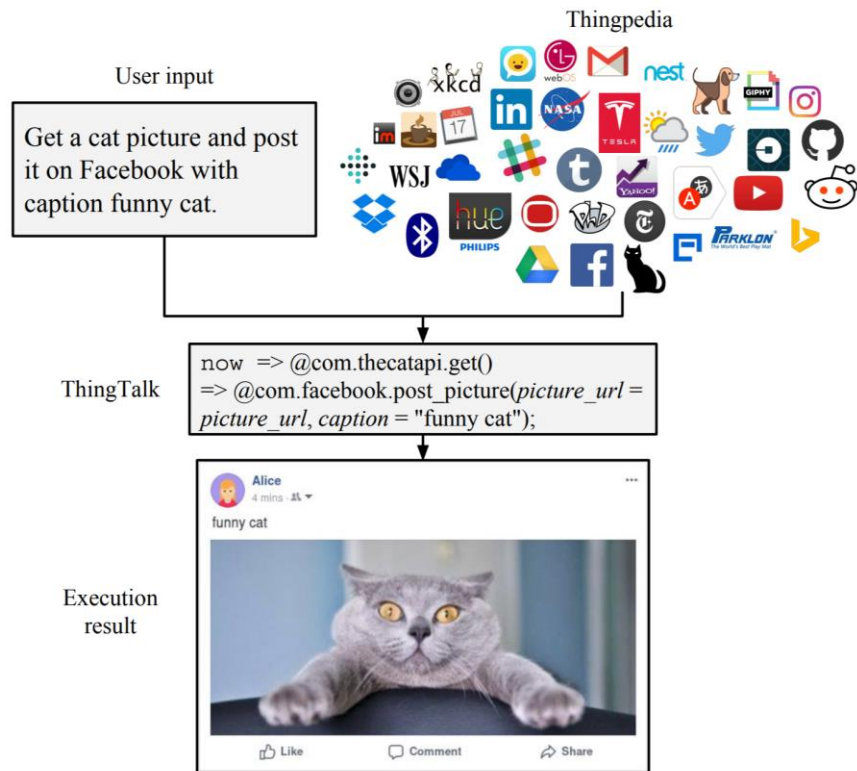
2

Figure 2: An example in WikiSQL. The inputs consist of a table and a question. The outputs consist of a ground truth SQL query and the corresponding result from execution.

# Executable Representations: ThingTalk

Campagna et al., 2019

Semantic-preserving transformation rules mean **canonical examples** for training the neural semantic parser.



# Outline

1. Introduction: Why Datasets?
2. MultiWOZ in the Almond/ThingTalk/Genie Context
3. What's In a Dataset
  - a. Dialogue Generation
  - b. Annotation Generation
  - c. Annotation Styles
- 4. MultiWOZ Revisited**

## 4. MultiWOZ Revisited

- MultiWOZ is a human-human dataset, mostly annotated, with intent and slot tagging.
  - But we don't use it fully, so that ends up being less important.
- MultiWOZ proposes itself as a benchmark dataset for:
  - Dialogue State Tracking
  - Dialogue Context-to-Text Generation
  - Dialogue Act-to-Text Generation

## Question

ARE THERE "BENCHMARKING  
BLIND SPOTS" OR BIASES  
THAT YOUR PROJECT MIGHT  
SUFFER BECAUSE OF THE  
DATASET CHOICE?



Thank you!