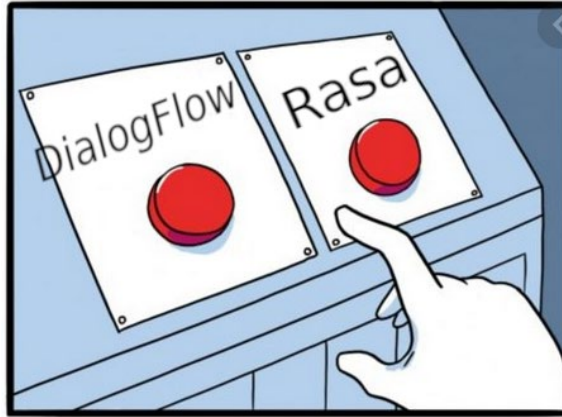# BLS, Integrations, DialogFlow, and Rasa (part 2)



Lecture 14

# Review: Tools and Web Services

- A **Web service** is a **program** that runs on a **server** that interacts with other computers and software over the **Internet**
  - For your prototype front- and back-ends, consider `screen` or `tmux`
    - Run services that persist without constant ssh'ing
  - Web services interact using serialized data (**JSON** payloads)
    - Basically, giant "standardized" strings that contain relevant information
    - Historically, XML or per-application custom protocols have been used
      - e.g., HTTP (still used today) is simply customized strings
        - GET google.com/
        - Response 200, 404, 403, 500, etc.

- Front end → JSON → CAI platform → JSON →  BLS

# One-Slide Summary: DF/Rasa/Integrations

- **Conversational AI platforms** generally are **web-based** tools that provide a way to design:
  - **Intents** for classifying text
  - **Slot extraction** (slot-value pairing, entity extraction) for pulling semantics
  - **Dialogue tracking** (states, contexts) for handling multiple "turns"
- **DialogFlow** is a Google platform
  - Design intents with contexts, use WebHook for "fulfillment" of business logic
- **Rasa** is an open-source set of command-line tools
  - Design intents and "stories", use Actions for business logic

# Comparing Platforms

- Large software systems come with various tradeoffs

- For Conversational AI, platforms balance
  - Speed of development
  - Ease of experience design
  - Testing
  - Training time
  - End-user experience

# Clinc

- You design a **state graph** that models flow through a conversation
  - A state encompasses an **intent**
  - Intent classification model is trained via examples per transition between states

- The platform is reached remotely via the /v1/query API endpoint
  - Your front-end uses this API (e.g., website, Alex, etc.)

- The platform engages your BLS on a per-utterance basis
  - Move through the state graph on each utterance

# Clinc: State Graph and Competencies

- The magic sauce to Clinc is the **State Graph** and notion of **competencies**

- These abstractions are meant to ease **development**
  - HITR: A human that helps you with something is *competent* at a few skills
    - A bank teller *knows* how to open an account
    - A travel agent *knows* how to get good airline tickets
    - A real estate agent (hopefully) *knows* how to not sell you a broken garbage house

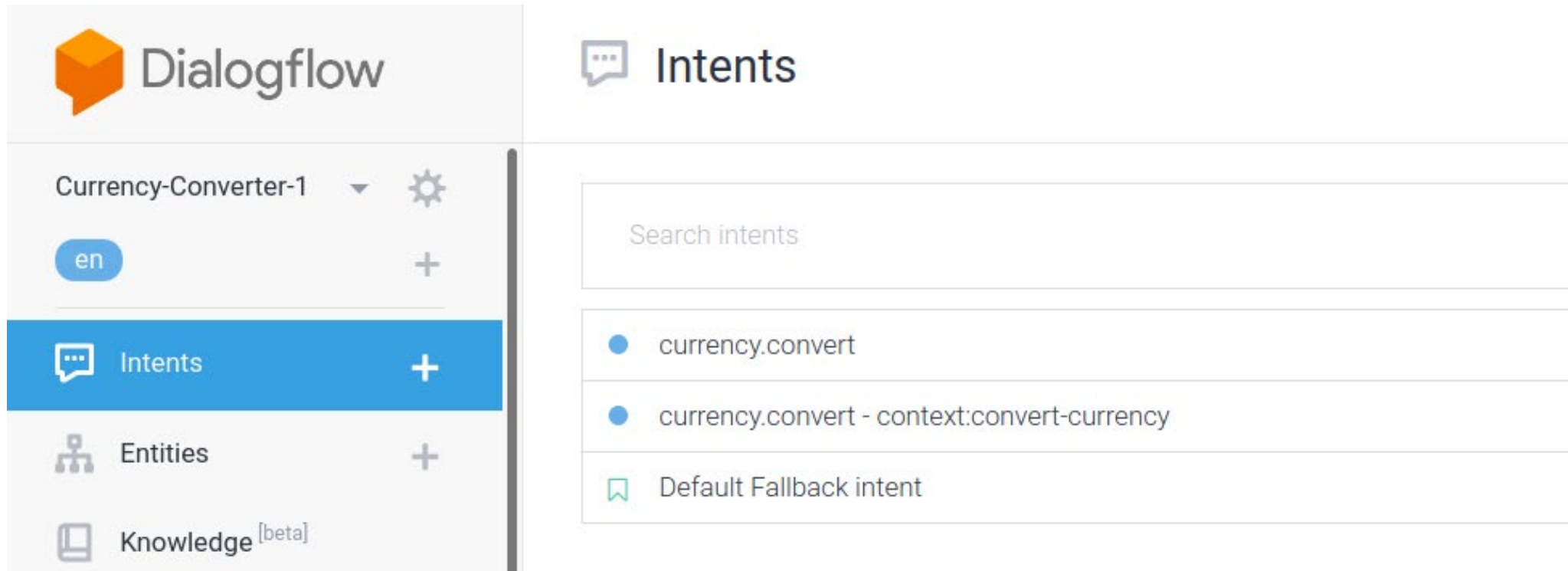# Clinc: Benefits and Drawbacks

- **Benefit:**
  - You design an experience around a state graph (complexity management)
  - Integration of crowdsourcing (simpler data collection and curation)

- **Drawback**:
  - Training is complicated
    - Example overlap
    - Single intent in multiple competencies
    - Stateful classifiers increase training time

# Google DialogFlow
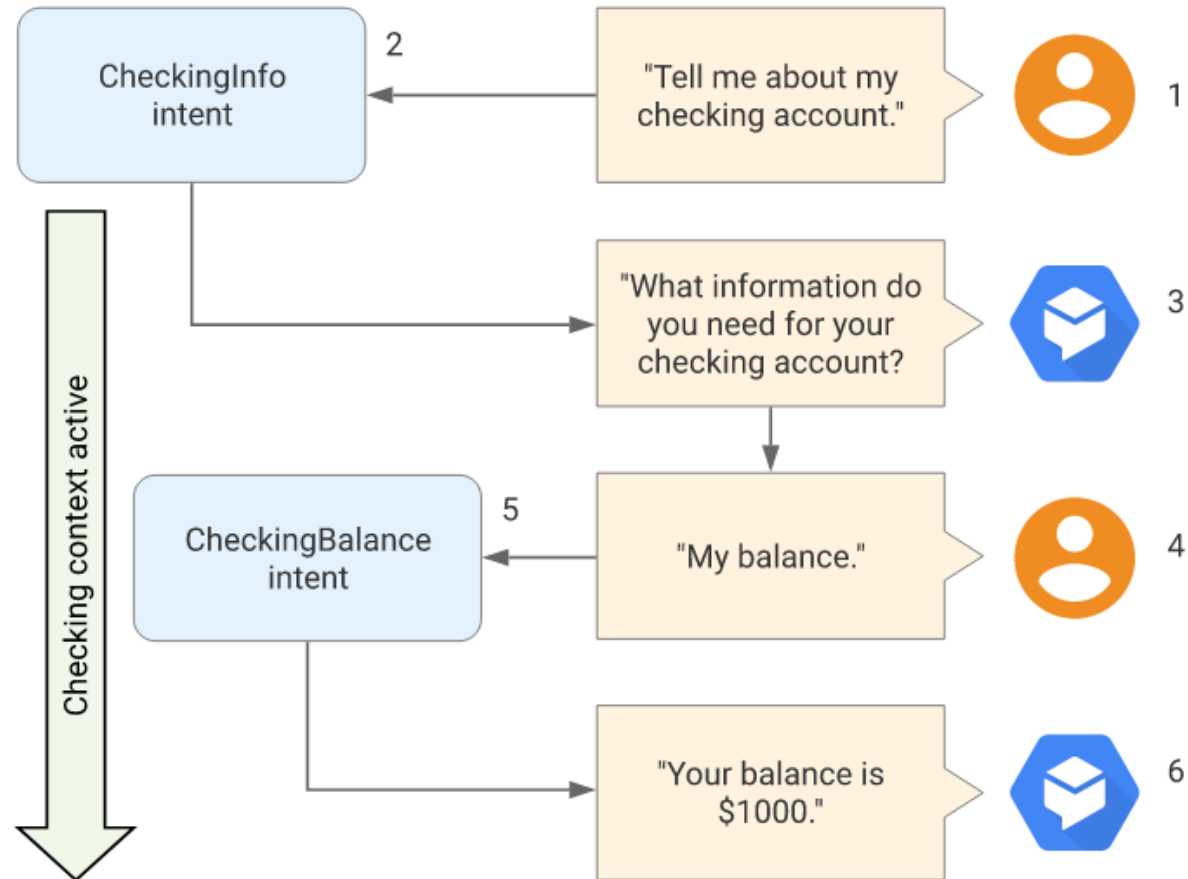
- DialogFlow allows specifying a list of **intents**

# Google DialogFlow

- No formal notion of *state*; instead, DialogFlow uses **contexts**
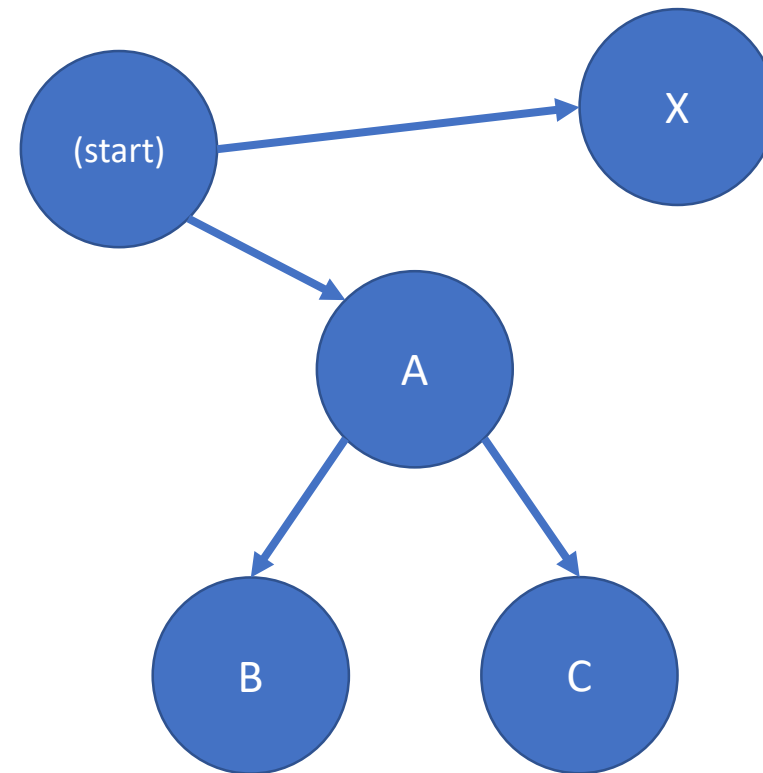
# DialogFlow: Context

- Similar to a Clinc **competency**, a DialogFlow **context** is a set of intents that share **slots (entities/parameters)**

- You can design responses and intents knowing that you can extract and fill slots using contexts

- Like with an intent or state, multiple **contexts** can become *active* when an intent is matched

# DialogFlow Contexts

- Contexts bias intent classification
  - Input contexts: helps to decide between potentially-ambiguous intents
    - "car" vs. "truck" in a vehicle-shopping bot
  - Output contexts: When an intent is classified, assign set of output contexts
    - "select_vehicle" in a vehicle-shopping bot

- Contexts also support timing
  - Default 20 minutes
  - 5 conversational turns

# DialogFlow: Follow-up intents

- Within an intent, a **follow-up** allows intent classification having already classified an intent
  - Helps maintain statefulness

  - A: I like candy
    - Yes:
      - B: Do you want chocolate?
    - No:
      - C: Do you want asparagus?
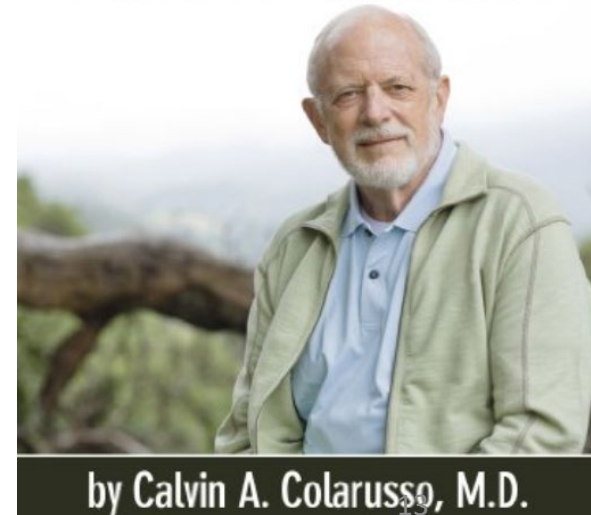  - X: I like turtles

B and C are follow-ups to A

# DialogFlow: Fulfillment

- (same as Business Logic)

- Basically, you run a web service that "fulfills" a request

- Set contexts, set intents, set responses, etc.

- You can also create fulfillment "Functions"
  that are integrated into DialogFlow



IN PURSUIT OF
**HAPPINESS AND FULFILLMENT**

by Calvin A. Colarusso, M.D.

# DialogFlow: Benefits and Drawbacks
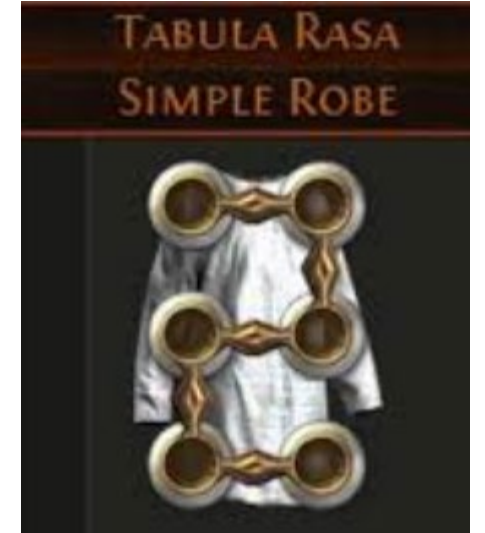
- **Benefits:**
  - Massive Google resources and availability
  - List of intents may be "simpler" than a state graph
  - Built-in support for BLS (i.e., integrated "Fulfillment")

- **Drawbacks:**
  - Intents, contexts, and follow-ups may introduce design complexity
  - Limitations of contexts and follow-ups promote serial dialogues
  - Cloud deployment?
    - Where is the data stored?

# Rasa

- Open source Conversational AI platform

- Lots of plaintext, command-line interactions

- Intents and slots as expected
  - Define a list of intents and corresponding examples
  - Label a few example slots within each example

- Rasa runs as a background service that you run locally

# Rasa: Domains

- A Conversational AI built in Rasa is defined by a **domain**
  - Provides a notion of "scope" for what the virtual assistant does

- The domain lists out all the intents, slots, responses

```
 1  intents:
 2    - greet
 3    - goodbye
 4    - affirm
 5    - deny
 6    - mood_great
 7    - mood_unhappy
 8    - bot_challenge
 9
10  actions:
11  - utter_greet
12  - utter_cheer_up
13  - utter_did_that_help
14  - utter_happy
15  - utter_goodbye
16  - utter_iamabot
17  - action_myendpoint
18
19  responses:
20    utter_greet:
21    - text: "Hey! How are you?"
22
23    utter_cheer_up:
24    - text: "Here is something to cheer you up:"
25      image: "https://i.imgur.com/nGF1K8f.jpg"
26
```

# Rasa: Intents

- Defined in MarkDown files

- Just a big list of intent names and examples

```
 1 ## intent:greet
 2 - hey
 3 - hello
 4 - hi
 5 - good morning
 6 - good evening
 7 - hey there
 8
 9 ## intent:goodbye
10 - bye
11 - goodbye
12 - see you around
13 - see you later
14
15 ## intent:affirm
16 - yes
17 - indeed
18 - of course
19 - that sounds good
20 - correct
```

# Rasa: Slots

- Defined in domain, provided in examples

- Rasa supports **types** of slots
  - List
  - String
  - Float
  - Range
  - Enumerated

```
## venue search
* search venues
    - action search venues
    - slot{"venues": [{"name": "Big Arena", "reviews": 4.5}]}

## concert search
* search concerts
    - action search concerts
    - slot{"concerts": [{"artist": "Foo Fighters", "reviews": 4.5}]}
```

- Rasa also works with a variety of channels
  - (slots can be strange... e.g., JSON objects)

# Rasa: Dialogue tracking with Stories

- Instead of *states*, Rasa uses **stories**

- You provide examples of sequences of intents and actions to take
  - Rasa *learns* sequences of intents to classify based on those examples

```
## story_email_not provided
*greet
    - utter_greet
*subscribe_newsletter
    - utter_ask_email
*inform{email:'example@example.com'}
    - slot{email:'example@example.com'}
    - action_subscribe_newsletter
```

# Rasa: Actions

- (same as business logic and fulfillment)

- **Actions** are a kind of special "intent" that can get classified
  - The virtual assistant can predict going to perform an action
  - You indicate Actions in Stories

- You can write a web service that is engaged
  - You need an "action" in your user stories

# Rasa: Benefits and Drawbacks

- **Benefits:**
  - Command-line and plaintext
  - Example-based input allows coercing statefulness
  - Built-in handling of "actions" (very Unix-y)

- **Drawbacks:**
  - Open-source project
    - Ever tried OpenOffice?
  - Lack of designed statefulness entails complex examples

# Conversational Platforms: Summary

- These platforms provide coarse abstractions for modeling conversations

- They *assume* conversations take place turn-at-a-time
  - Intent-based
  - Slots to extract
  - Etc.

- The developer still needs to:
  - Model the conversation
  - Collect and curate the training data
  - Write code to do real stuff (i.e., front end and business logic)

# Response Generation

- Conversational AI platforms can figure out what a user is talking about

- However, AI is still far too limited to actually "think"
  - Need human designer to figure out how to **take action** and **form responses**

- This manifests in two ways
  - Business Logic / Webhook / Action: Code you write to actually do stuff
  - **Response Generation:** filling out a response that an end-user will understand

# Response Generation

- This is *not* the same as **Natural Language Generation**

| Context | All of a sudden, [*Emily*]₁ walked towards [*the dragon*]₂. |
|---|---|
| Current Sentence | [*Seth*]₃ yelled at [*her*]₁ to get back but _____ |

Figure 1: An example of entity-labeled story data. The brackets indicate which words are part of entity mentions. Mentions marked with the same number refer to the same entity. The goal is to continue the story in a coherent way. The actual story reads, "*Seth yelled at her to get back but she ignored him.*"

**INSTA-TRUMP**

...merica for our country, thriving. It can happen. They shouldn't be allowed to have a nuclear weapons that are obsolete.

We're a debtor nation. He agreed.

...So I will devote 1,000%. I couldn't answer the question, is Iraq a good time. His brand – " Like I care about that country falls apart.

Well, I got it for two reasons. Number one, he won't want it. Probably a lobbyists, by the donors, and by the way, she was not awake to take that call at 3 o'clock in the morning.

And I want to tell you it's been no crowd like this.

**Document Planning** — content planning, document outline

**Microplanning** — referring expressions, word choice, aggregation

**Realisation** — converting specifications to a real text

# Response Generation

- Usually **template** based with **rules** (e.g., code)

# Response Generation

- The Conversational AI platform you use entails a specific method for generating responses

- Clinc: Jinja based on slots filled and intents
  - The weather in {% slots['city'] %} is {% slots['temperature'] %}.

- DialogFlow: Responses with syntax
  - The weather in $city is $temperature.

- Rasa: Built into stories
  - The weather in {city} is {temperature}.

# Response Generation

- When you design the experience, a few considerations:
  - Is the response being **read aloud**? Or read on a **screen**?
  - Is the response going to be **interpreted** by some front-end?
    - Images?
    - HTML?
  - Does the response make **logical** sense?
    - Does it account for slots received already?
    - Does the response flow from previous responses?
  - Is the response an appropriate **length**?
    - Can the user find the information or do the task faster than it takes to read the response?