

15-494/694: Cognitive Robotics

Dave Touretzky

Lecture 11:

Speech Generation and
Recognition

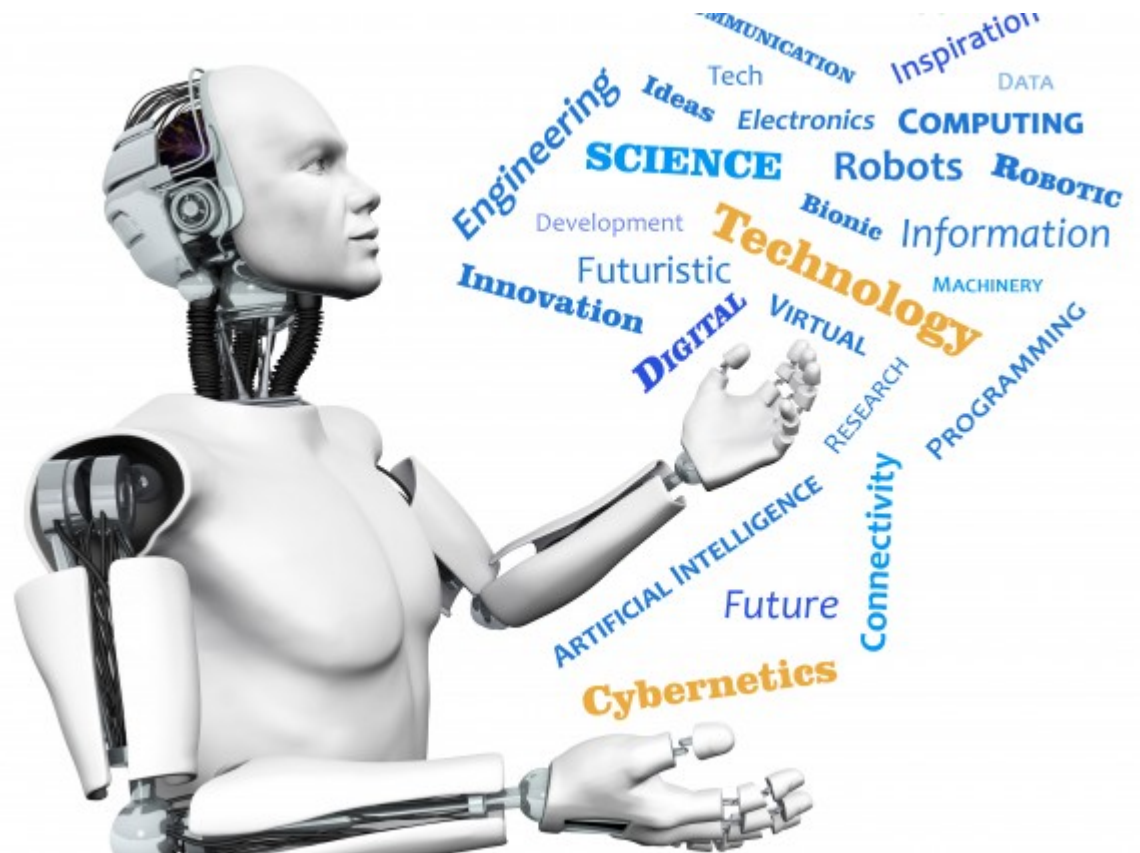


Image from <http://www.futuristgerd.com/2015/09/10>

Speech Generation

- We use the Google Cloud Text to Speech API to generate speech for VEX AIM.
- Parameters are defined in `actuators.py`

```
tts_voice = texttospeech.VoiceSelectionParams(  
    language_code="en-US",  
    name="en-US-Journey-F",  
    ssml_gender=texttospeech.SsmlVoiceGender.FEMALE  
)
```

- Requires Google Cloud credentials. If not available, revert to `gTTS` package which uses Google Translate's speech facility.

Google Cloud Text to Speech

- Info at <https://cloud.google.com/text-to-speech>
- Multiple voice models
 - Basic
 - Studio
 - WaveNet
 - Neural2
 - Journey (now “Chirp HD”)
 - etc.
- Some models allow control of speech rate and pitch.

SSML

- Speech Synthesis Markup Language
- Can be used to control pronunciation of things like acronyms or numbers
- Can be used to insert pauses or emphasis where needed
- Not currently used in vex-aim-tools but might be in the future.

SSML Example

```
<speaK xmlns="http://www.w3.org/2001/10/synthesis" xml:lang="en-US">
  <p>
    <s>This is a normal sentence with a short pause.</s>
    <s>I can make this word sound a little
      <emphasis level="moderate">important</emphasis>!</s>
    <s><say-as interpret-as="ordinal">1st</say-as> in line.</s>
  </p>
</speaK>
```

“Say” Node

- Constant case:

`Say('hello there') =C=> next`

- List (will choose at random):

`Say(['hello', 'hi', 'howdy']) =C=> next`

- Event-driven case:

`Compute() =SayData=> Say() =C=> next`

- Subclassing “Say”:

`class SpeakBattery(Say)`

Make A SpeakBattery Node

```
class SpeakBattery(Say):  
    def start(self, event=None):  
        cap = self.robot.battery_capacity  
        self.text = f'battery capacity {cap} %'  
        super().start(event)
```

Speech Recognition

- VEX AIM has no microphone
- Use the laptop's mic or a USB mic
- Recognition via the browser's speech API
 - Must have network access to function.
 - Biased towards conversational English, not arbitrary robot commands
 - Accuracy in 2026 is quite good.
- Sometimes uses special characters we don't want, e.g., "15 degrees" is transcribed as "15°".

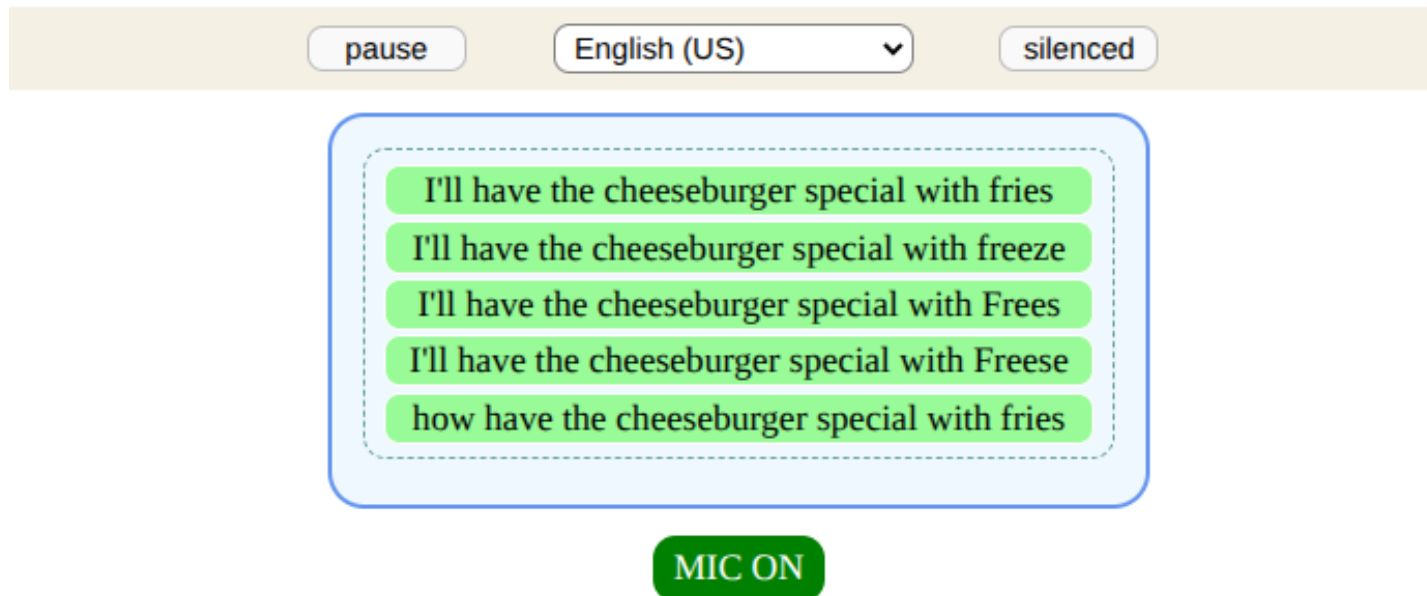
Demo: Google Speech API

<https://www.cs.cmu.edu/~dst/SpeechDemo>

Speech Recognition Demo

Speak into your microphone; see the results below.

Click [here](#) for experiments to try.



The screenshot shows the Google Speech API demo interface. At the top, there is a control bar with three buttons: "pause", "English (US)" (with a dropdown arrow), and "silenced". Below this bar is a large light blue rounded rectangle containing a dashed border box. Inside this box are five green rounded rectangles, each containing a different transcription of the same sentence: "I'll have the cheeseburger special with fries", "I'll have the cheeseburger special with freeze", "I'll have the cheeseburger special with Frees", "I'll have the cheeseburger special with Freese", and "how have the cheeseburger special with fries". Below the dashed box is a green button with the text "MIC ON".

Hearing Our Own Speech

- To avoid the robot hearing its own speech, the Say node temporarily disables speech recognition before speaking.
- It re-enables recognition when the speech is complete.
- This process is imperfect. Mistakes will be made.

Declining Speech Recognition

Speech recognition is turned on by default.

To turn it off: use `speech=False` in `StateMachineProgram`.

```
class VEXCommand(StateMachineProgram):  
    def __init__(self):  
        super().__init__(speech=False)
```

Disabling Speech Recognition

- On Linux or macOS, disable speech recognition for simple_cli with

```
export ROBOT_NO_SPEECH=1
```

- You can simulate speech input by typing in the simple_cli REPL, beginning the line with a period and a space:

```
C> . this is speech input
```

When To Listen

- Microphone is always on
- We could use a wake word to indicate we're addressing the robot.
 - “Celeste, please grab a barrel”
- You've seen this trick before:
 - “Alexa, ...”
 - “Hey Siri, ...”
 - “OK Google, ...”
- Celeste doesn't do this because it's awkward for longer dialogs.

The =Hear()=> Transition

```
dispatch: Say('What now?')
```

```
dispatch =Hear('celeste turn left')=>  
  Turn(90) =C=> dispatch
```

```
dispatch =Hear('celeste drive forward')=>  
  Forward(50) =C=> dispatch
```

String Matching for =Hear=>

- Remove all punctuation
- Convert everything to lowercase
- Normalize homophones

Homophones

- “Thesaurus” data structure defined in `aim_fsm/thesaurus.py`
- Words:
 - `cozmo` ← `cosmo`, `cosmos`, `cosimo`, ...
 - `right` ← `write`, `wright`
 - `cube1` ← `q1`, `coupon`, `cuban`
- Phrases:
 - `cube1` ← `cube 1`
 - `paperclip` ← `paper clip`

Regular Expression Matching

- Uses the Python re package

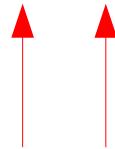
- Example: optional words

'celeste ?(please|) drive forward'

- Be careful about spaces!

- Example: scanning for keywords:

'celeste .* grab.*'



spaces on both sides of .* will be a problem
if the .* matches the null string

Regular Expressions

.	match any char.	*	zero or more, e.g.
^	start of string	.	* matches any
\$	end of string	.	*? is non-greedy
\d	digit	+	+? one or more
\w	word char. (alpha-numeric or underscore)	?	?? zero or one
\s	whitespace	(...)	grouping
\b	word boundary	(x y)	alternation
		r'regexp \d syntax'	
		'string \d syntax'	

Checking the Match Results

- When a =Hear=> transition fires because of a match, it offers a SpeechEvent to the target node(s).
- The SpeechEvent contains three items:
 - **string:** the input string
 - **words:** list of words in the string
 - **result:** the match result from re.match
 - contains the groups defined by ()

Extracting Groups (1)

```
from aim_fsm import *

class Speech1(StateMachineProgram):

    class Heard(Say):
        def start(self, event):
            obj = event.result.groups()[1]
            self.text = 'I will grab %s' % obj
            super().start(event)
```

Extracting Groups (2)

```
$setup{
  loop: Say( 'what now' )
  loop =Hear( 'celeste ?(please|) grab a
    (barrel|ball)' )=>
    self.Heard() =C=> loop
  loop =Hear=> Say( 'Pardon me?' )
    =C=> loop
}
```

Dialoging with GPT-4o

- Instead of parsing user utterances with regular expressions in Hear transitions, we can let GPT-4o do that work.
 - Much better strategy!
- Now the problem is how to get GPT-4o's understanding back into our program logic.
 - Use #hashtag tokens for actions to take
 - Ask GPT to parse user input and return info in a specific form (e.g., dominoes)

Dialog

- Dialog requires access to a knowledge base, state, and conversational history.
 - The current pose and world map
 - The robot's recent actions (and plans?)
 - Recent object references
 - Necessary to resolve “it”, e.g.:
 - “Do you see an orange barrel?”
 - “Pick it up.”
 - Timestamps
- These are automatically provided by the AskGPT node.

What Else Is Needed?

- A simple Query-Response structure is somewhat inflexible.
 - GPT must wait for input from the human or a request from the robot.
 - What if we want GPT to initiate action on its own?
- Could poll GPT with the current state and see if it wants to take any action. We haven't tried this yet.