

(Low-Resource) NLP Tasks

Graham Neubig
@ CMU Low-resource NLP Bootcamp
5/18/2020



Carnegie Mellon University
Language Technologies Institute



NEULAB

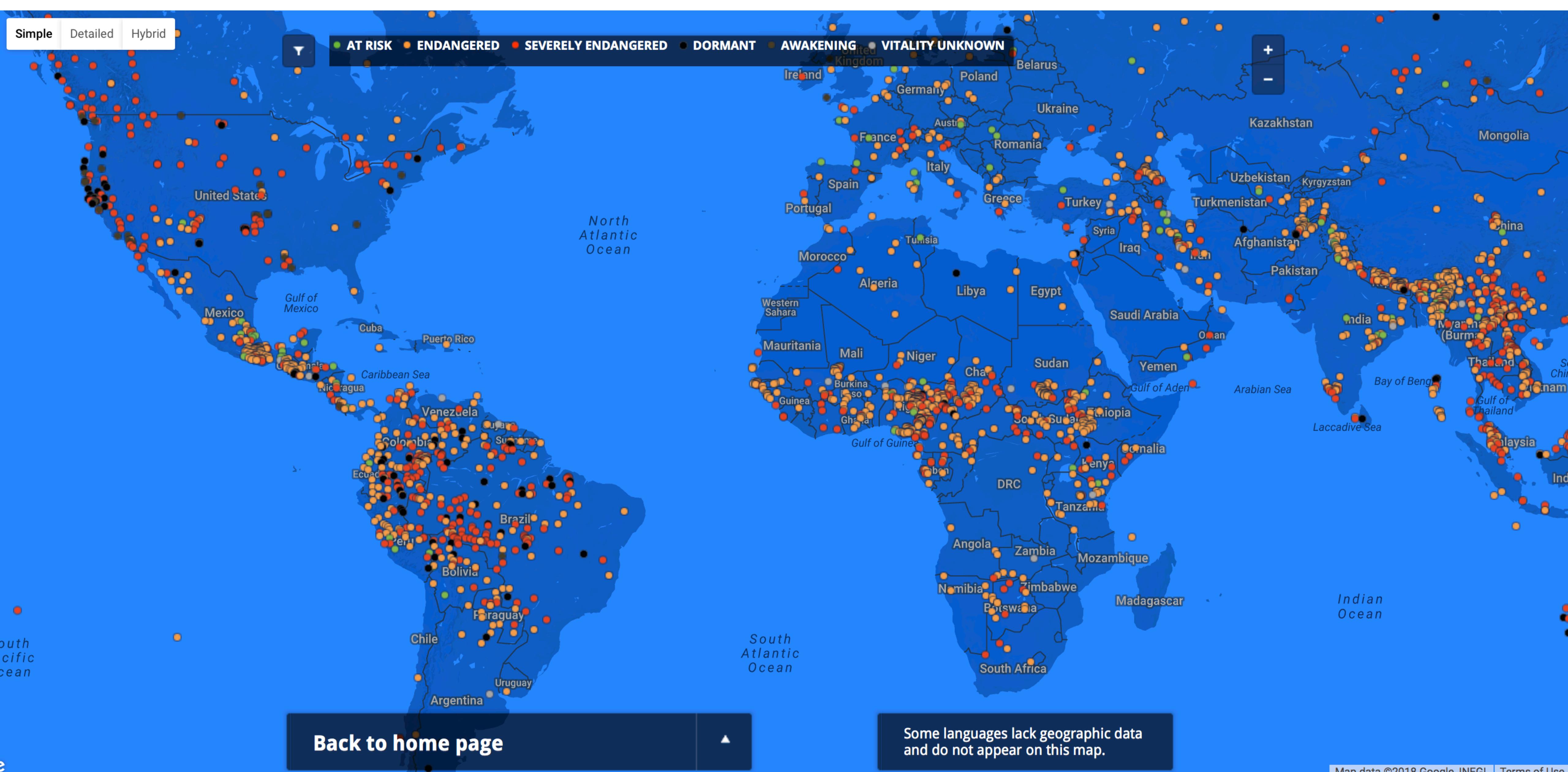
Most Spoken Languages of the World

- | | |
|-----------------------|--------------------------------|
| 1. English (1.132 B) | 6. العربية (273.9 M) |
| 2. 中文(普通话) (1.116 B) | 7. বাংলা (265.0 M) |
| 3. हिन्दी (615.4 M) | 8. Россия (258.2 M) |
| 4. Español (534.4 M) | 9. Português (234.1 M) |
| 5. Français (279.8 M) | 10. Bahasa Indonesia (279.8 M) |

Source: Ethnologue 2019 via Wikipedia

Simple Detailed Hybrid

AT RISK ENDANGERED SEVERELY ENDANGERED DORMANT AWAKENING VITALITY UNKNOWN



Map data ©2018 Google, INEGI | Terms of Use

<http://endangeredlanguages.com/>

Why NLP for All Languages?

- Aid **human-human communication** (e.g. machine translation)
- Aid **human-machine communication** (e.g. speech recognition/synthesis, question answering, dialog)
- **Analyze/understand language** (syntactic analysis, text classification, entity/relation recognition/linking)

Rule-based NLP Systems

- Develop rules, from simple scripts to more complicated rule systems
- Generally must be developed for each language by a linguist
- Appropriate for some simple tasks, e.g. pronunciation prediction in epitran

```
from epitran.backoff import Backoff
>>> backoff = Backoff(['hin-Deva', 'eng-Latn', 'cmn-Hans'])
>>> backoff.transliterate('हिन्दी')
'hindi:'
>>> backoff.transliterate('English')
'ɪŋglɪʃ'
>>> backoff.transliterate('中文')
'tsɔŋwən'
```

<https://github.com/dmort27/epitran>

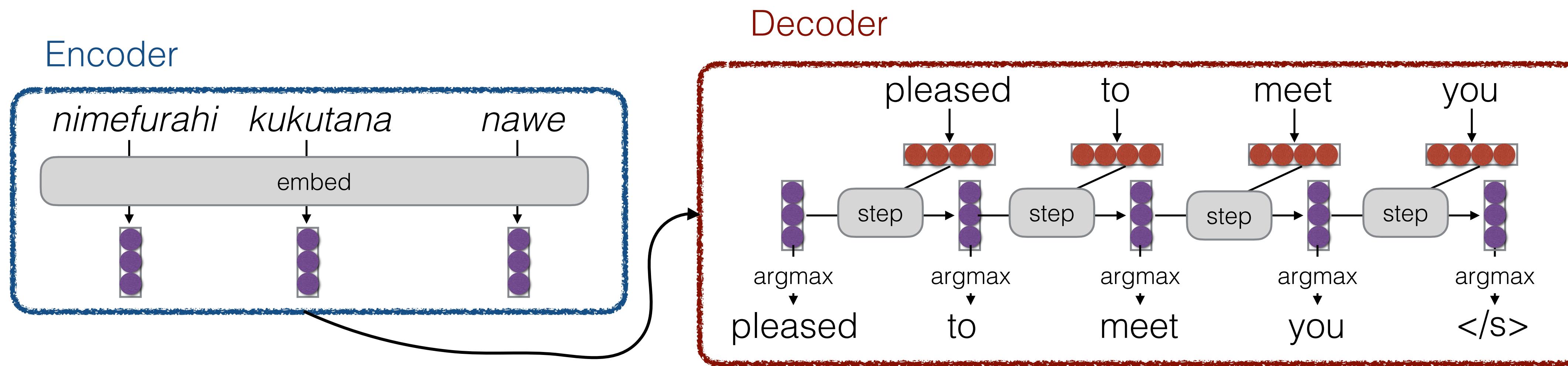
Machine Learning NLP Systems

- Formally, learn a model to map an $\text{input } X$ into an $\text{output } Y$. Examples:

<u>Input X</u>	<u>Output Y</u>	<u>Task</u>
Text	Text in Other Language	Translation
Text	Response	Dialog
Speech	Transcript	Speech Recognition
Text	Linguistic Structure	Language Analysis

- To learn, we can use
 - Paired data $\langle X, Y \rangle$, source data X , target data Y
 - Paired/source/target data in *similar* languages

Example Model: Sequence-to-sequence Model with Attention

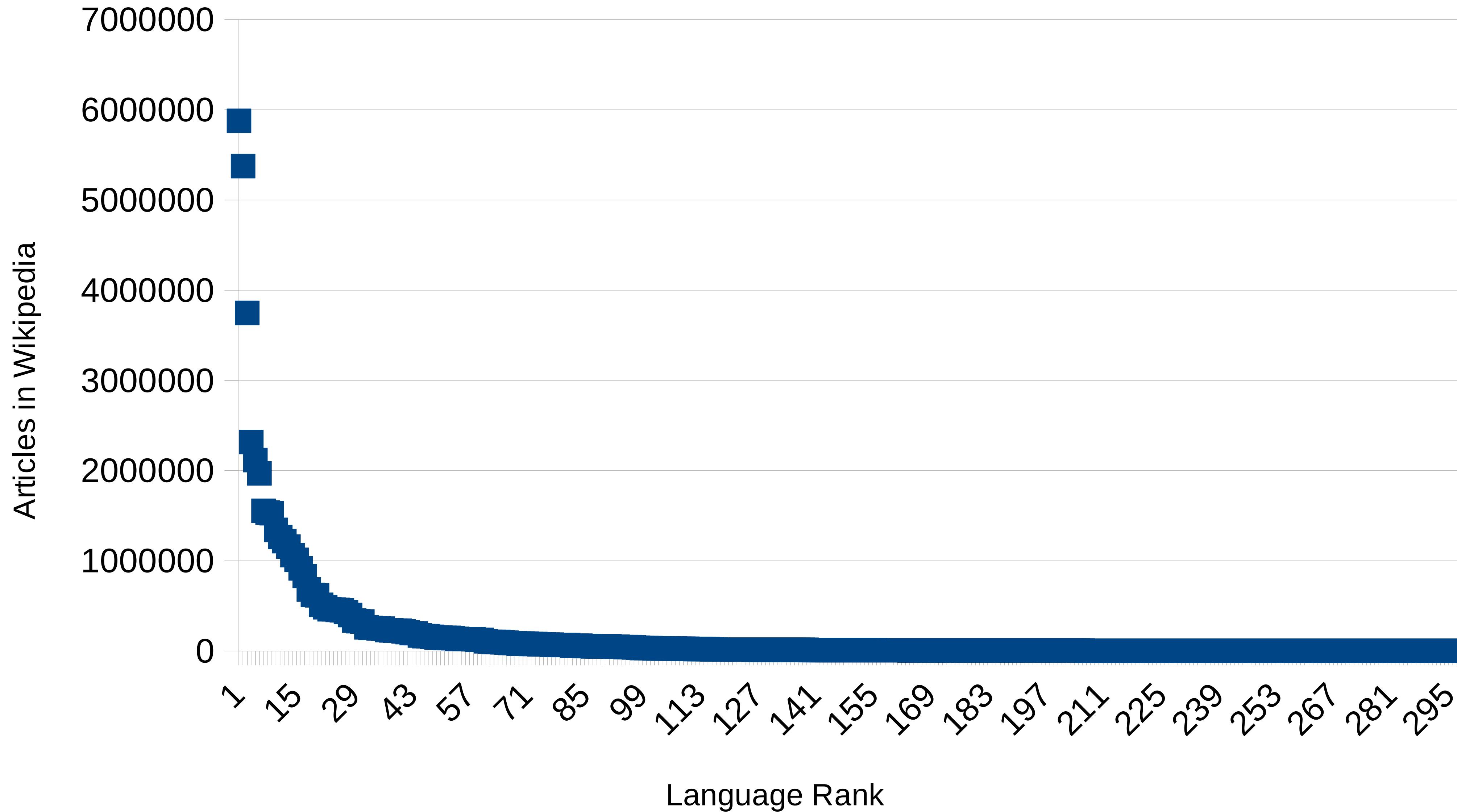


- **Various tasks:** Translation, speech recognition, dialog, summarization, language analysis
- **Various models:** LSTM, transformer
- Generally trained using **supervised learning**: maximize likelihood of $\langle X, Y \rangle$

Evaluating ML-based NLP Systems

- **Train** on training data
- **Validate** the model on "validation" or "development data"
- **Test** the model on unseen data according to a task-specific evaluation metric

The Long Tail of Data



Aiding Human-Human Communication

Machine Translation

Input X

Text

Output Y

Text in Other Language

Task

Translation

Google Translate interface:

- Top bar: Google Translate
- Language pair: JAPANESE ↔ ENGLISH
- Input text: カーネギー・メロン (Carnegie Mellon)
- Output text: Kānegīmeron
- Bottom right: Send feedback button

Microsoft Translator interface:

- Top bar: Microsoft Translator, Text, Conversation, Apps, For business, Help
- Message: We have updated our terms of use. [Learn More](#) [Dismiss](#)
- Input text: 言語技術研究所 (Institute of Language And Technology)
- Output text: Institute of Language And Technology
- Bottom right: 7/5000

Machine Translation Data

Last year I showed these two slides so that demonstrate that the arctic ice cap, which for most of the last three million years has been the size of the lower 48 states, has shrunk by 40 percent.

去年 この2つのスライドをお見せして 過去3百万年 アラスカとハワイを除く米国と一緒に面積があった極域の氷河が 約40%も縮小したことが おわかりいただけたでしょう

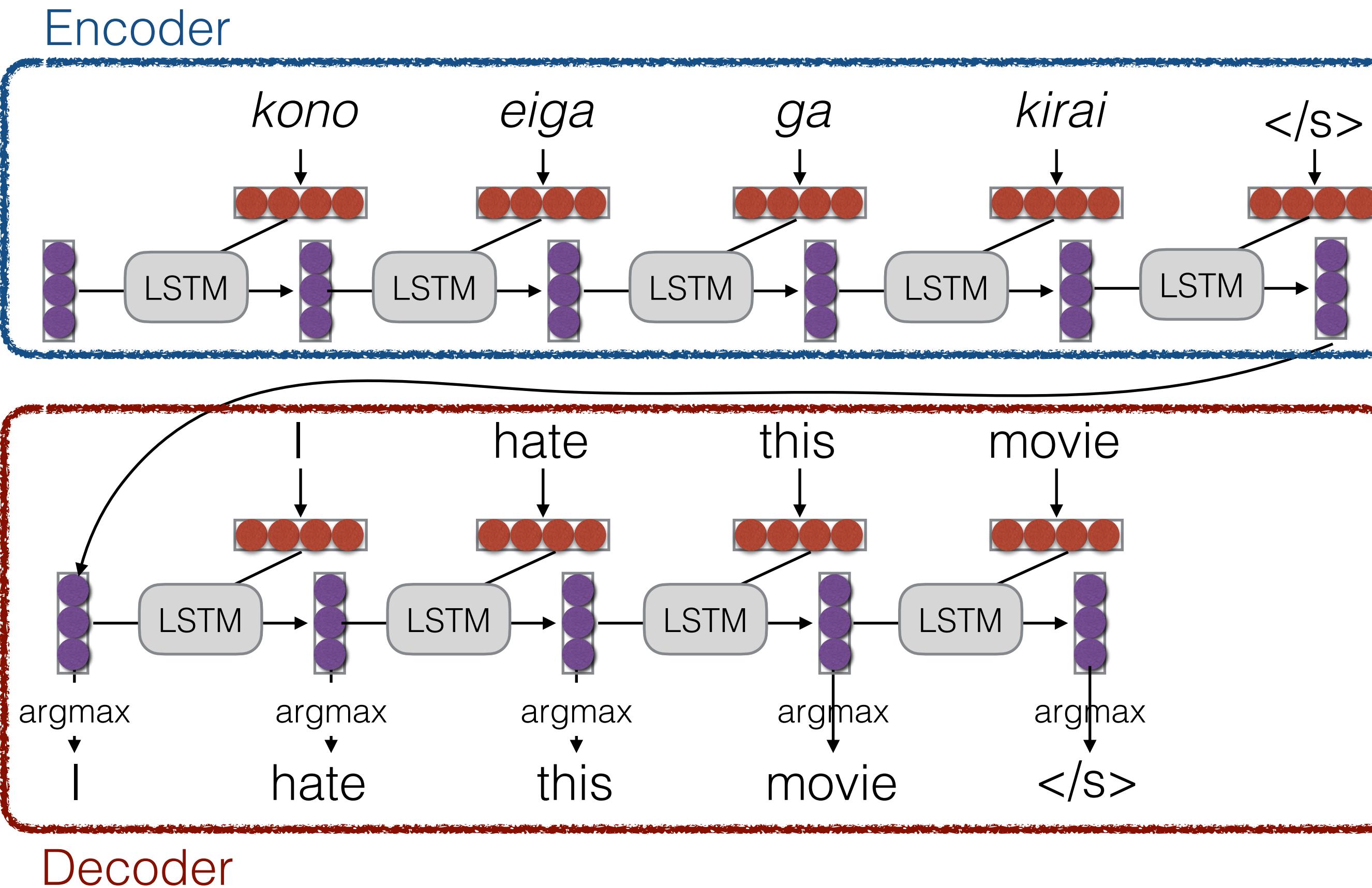
But this understates the seriousness of this particular problem because it doesn't show the thickness of the ice.
しかし もっと深刻な問題というのは 実は氷河の厚さなのです

The arctic ice cap is, in a sense, the beating heart of the global climate system.
極域の氷河は 言うなれば 世界の気候システムの鼓動する心臓で

It expands in winter and contracts in summer.
冬は膨張し夏は縮小します

The next slide I show you will be a rapid fast-forward of what's happened over the last 25 years.
では 次のスライドで 過去25年の動きを早送りにして見てみましょう

MT Modeling Pipeline



<https://github.com/pytorch/fairseq>



Joey NMT
<https://github.com/joeynmt/joeynmt>

Naturally Occurring Sources of MT Data

- Compared to other NLP tasks, data relatively easy to find!
 - **News:** Local news, BBC world service, Voice of America
 - **Government Documents:** Governments often mandate translation
 - **Wikipedia:** Some Wikipedia articles are translated into many languages, identify and
 - **Subtitles:** Subtitles of movies and TED talks
 - **Religious Documents:** Bible, Jehovah's Witness publications



<http://opus.nlpl.eu/>

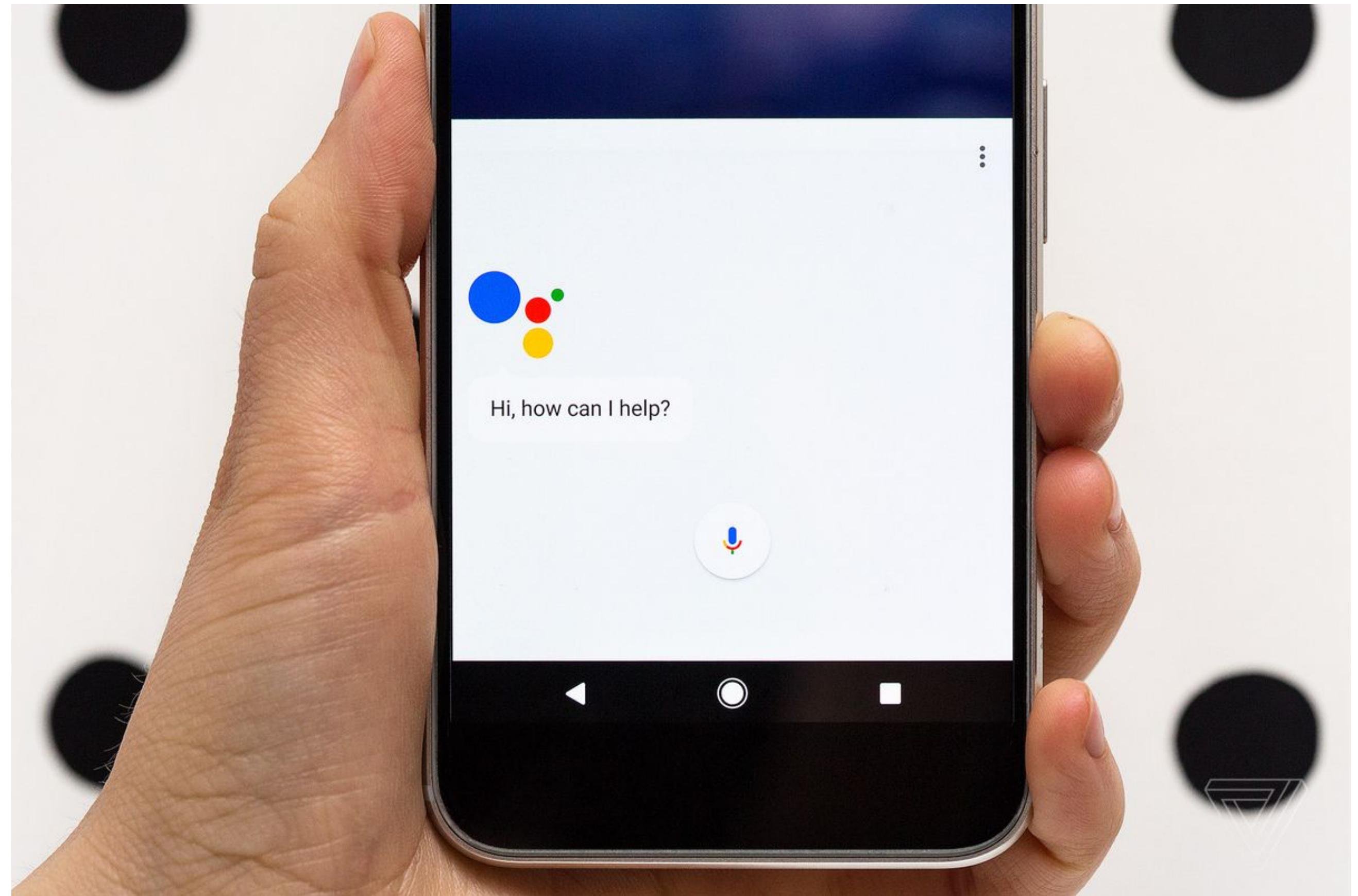
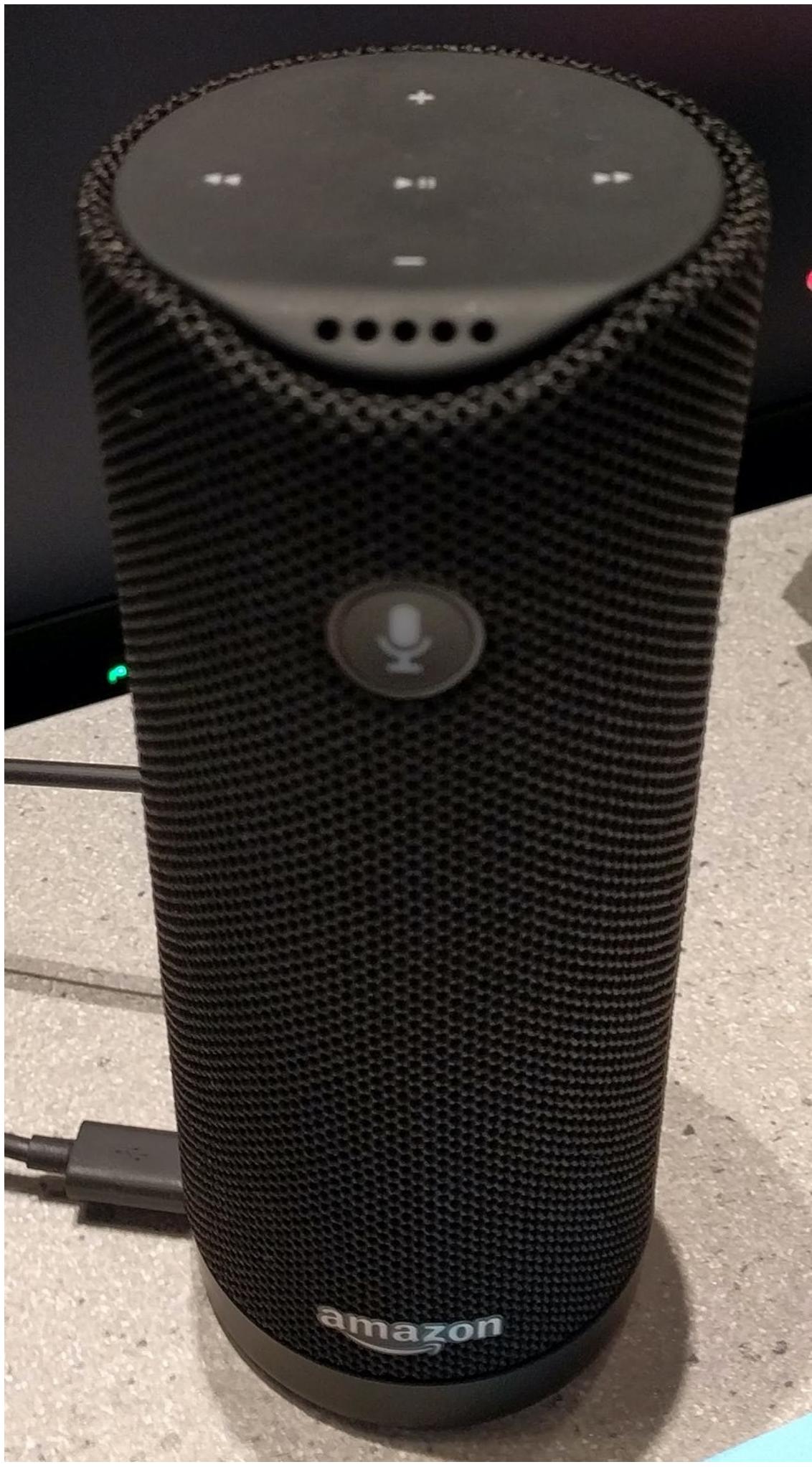
MT Evaluation Metrics

- Two varieties of evaluation:
- **Manual Evaluation:** Ask a human annotator how good they think the translation is, including fluency (how natural is the grammar), adequacy (how well does it convey meaning)
- **Automatic Evaluation:** Compare the output to a reference output for lexical overlap (BLEU, METEOR), or attempt to match semantics (MEANT, BERTscore)

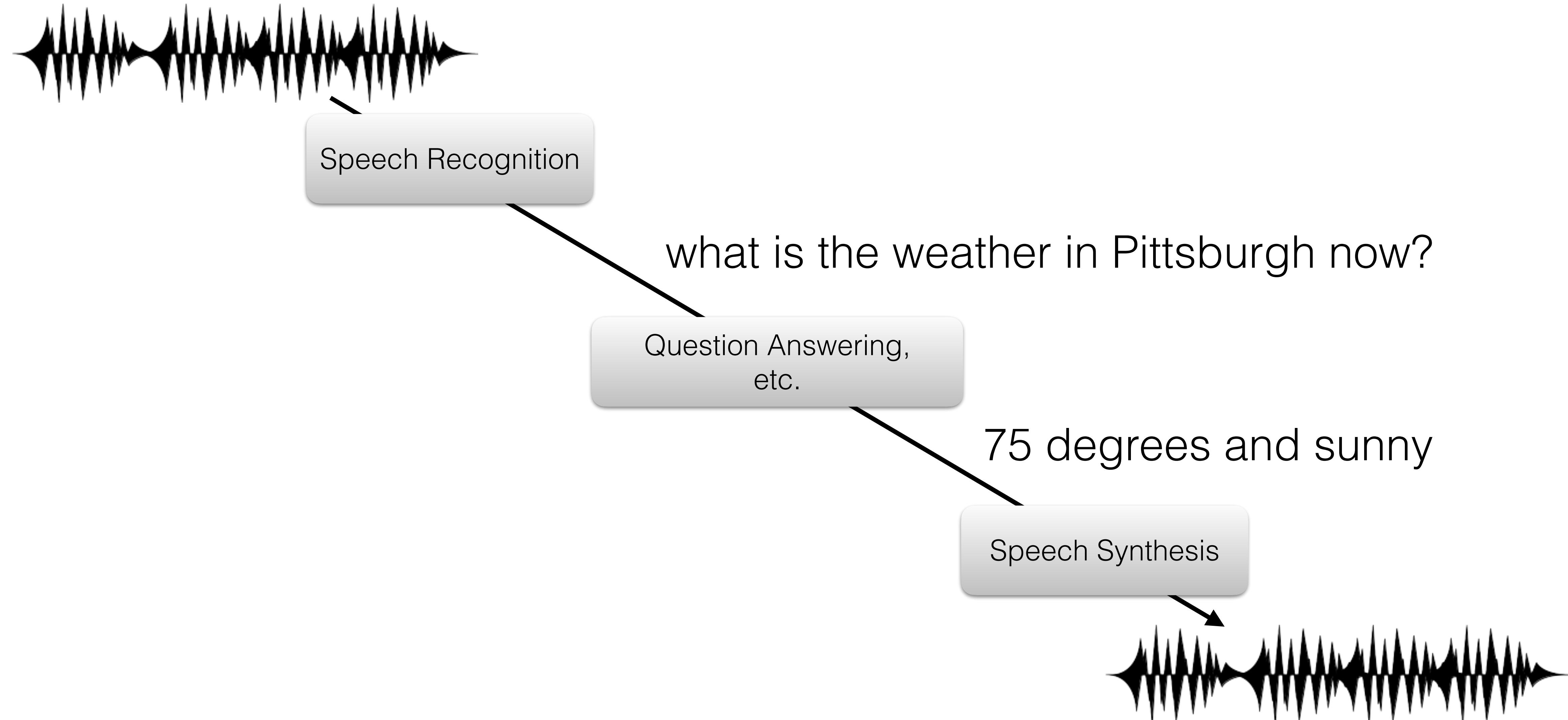
Translation	Fluency	Adequacy	Overlap
please send this package to Pittsburgh	high	high	perfect
send my box, Pittsburgh	low	medium	low
please send this package to Tokyo	high	low	high
I'd like to deliver this parcel, destination Pittsburgh	high	high	low

Aiding Human-Machine Communication

Personal Assistants

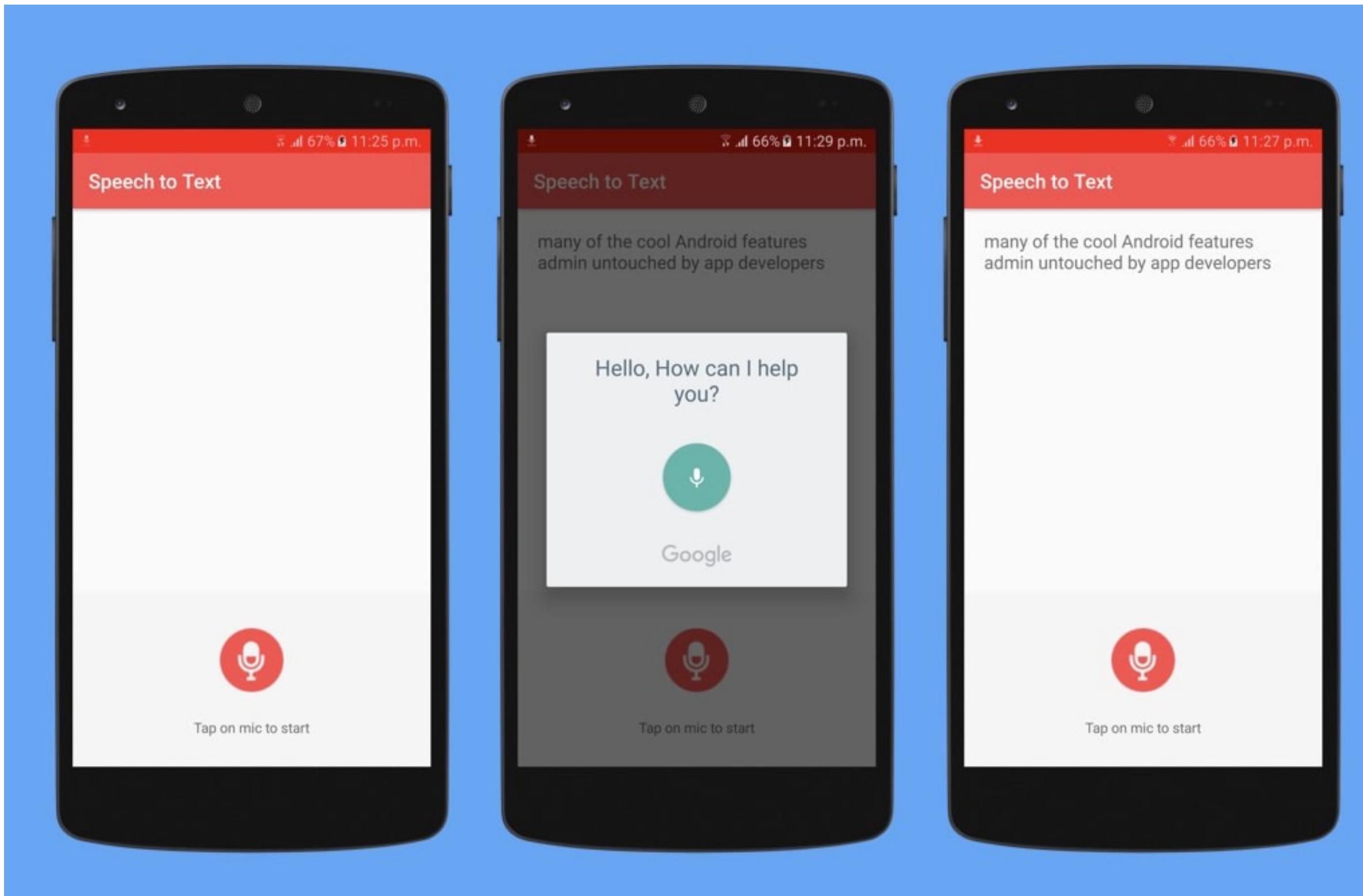


Personal Assistant Pipeline



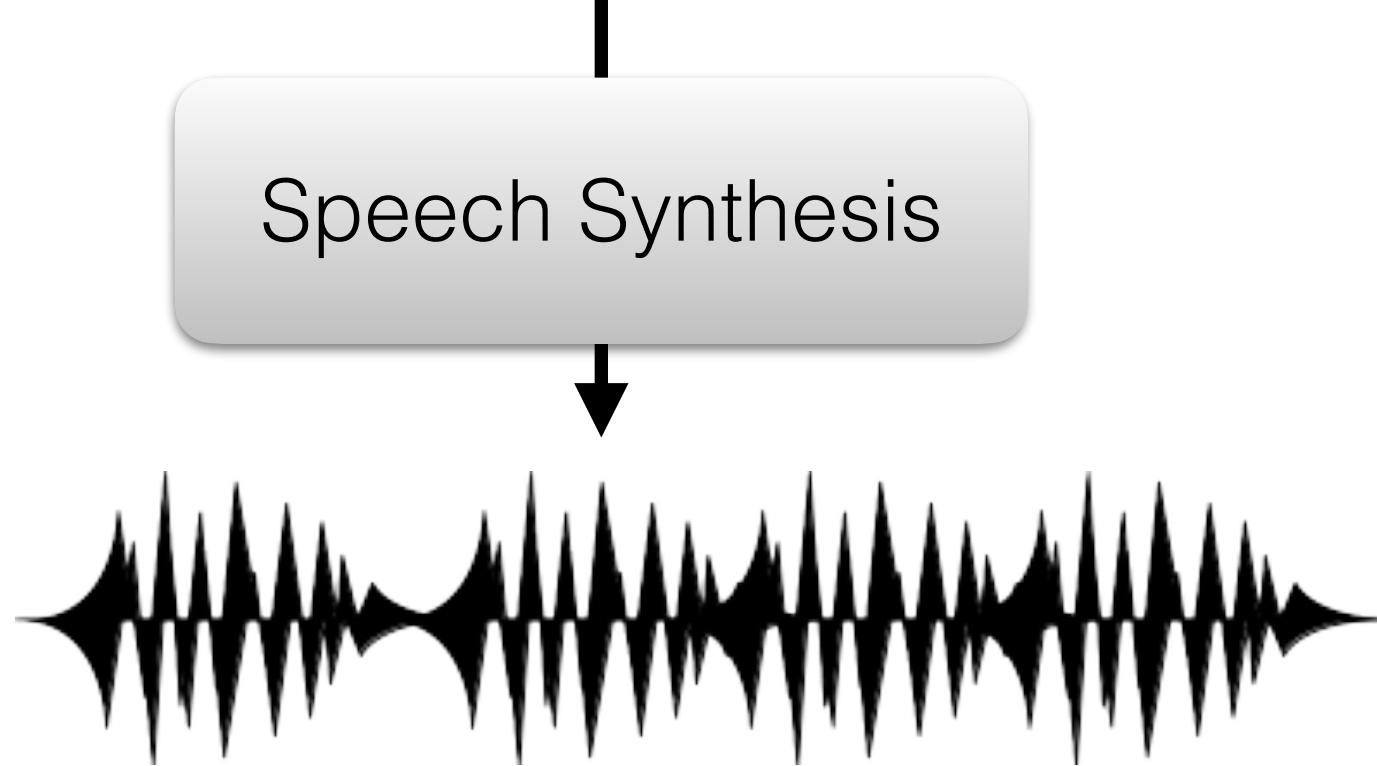
Speech

Input X Output Y Task
Speech Text Speech Recognition



Input X Output Y Task
Text Speech Speech Synthesis

75 degrees and sunny



Speech Data



Last year I showed these two slides so that demonstrate that the arctic ice cap, which for most of the last three million years has been the size of the lower 48 states, has shrunk by 40 percent.



But this understates the seriousness of this particular problem because it doesn't show the thickness of the ice.



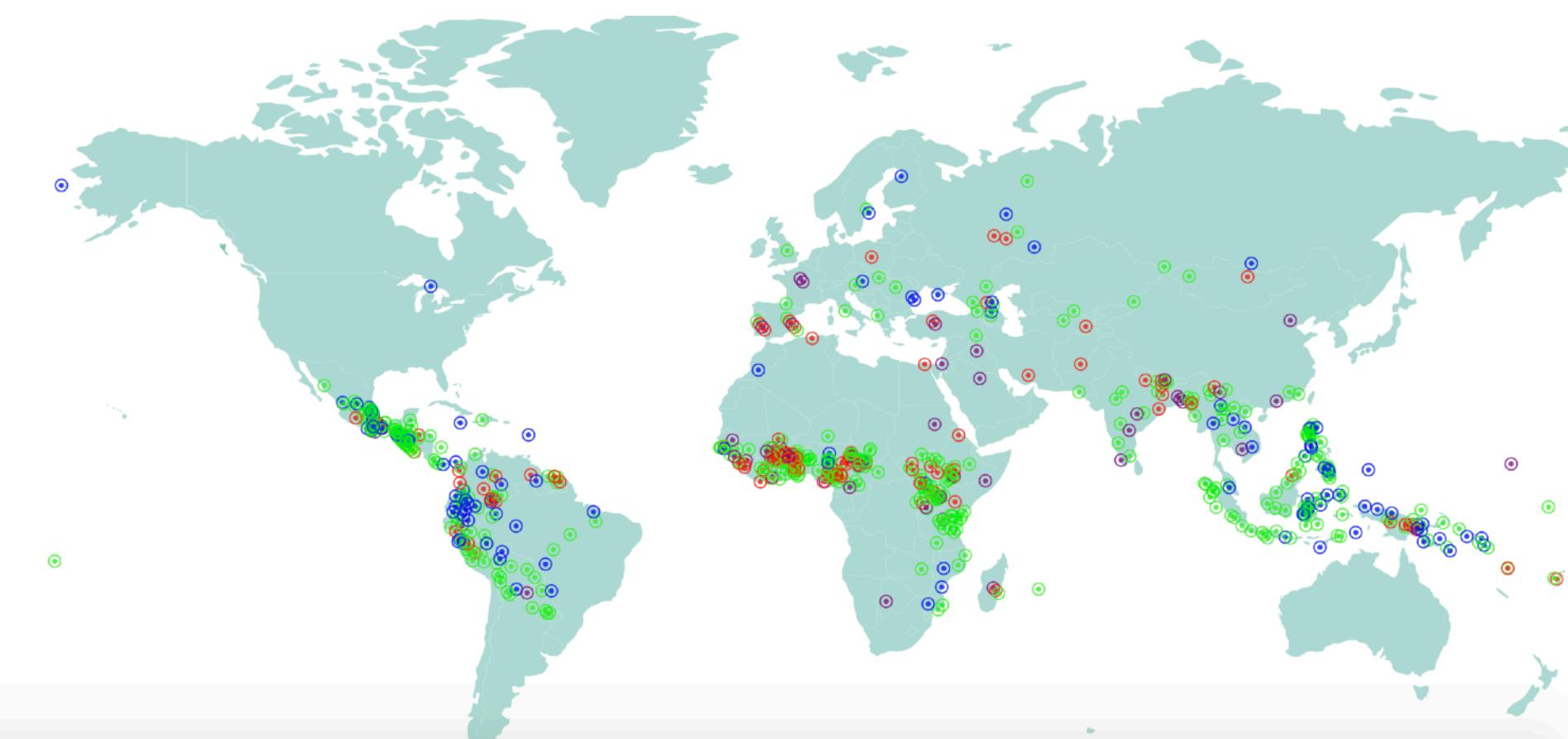
The arctic ice cap is, in a sense, the beating heart of the global climate system.

- **Speech Recognition:** Multi-speaker, noisy, conversational best for robustness
- **Speech Synthesis:** Single-speaker, clean, clearly spoken best for clarity

Naturally Occurring Sources of Speech Data

- **Transcribed News:** Sometimes spoken radio news also has transcriptions
- **Audio Books:** Regular audio books or religious books
- **Subtitled Talks/Videos:** TED(x) talks or YouTube videos often have transcriptions
- **Manually Transcribed Datasets:** Record speech you want and manually transcribe yourself (e.g. CallHome)

CMU Wilderness Multilingual Speech Dataset



https://github.com/festvox/datasets-CMU_Wilderness

Common Voice
moz://a

CONTRIBUTE DATASETS LANGUAGES ABOUT

Today's Progress
501 / 1200
Clips recorded

Speak

Donate your voice

Recording voice clips is an integral part of building our open dataset; some would say it's the fun part too.

Help us get to 1,200

<https://voice.mozilla.org/en>

Speech Recognition Modeling Pipeline

- **Feature Extraction:** Convert raw wave forms to features, such as frequency features
- **Speech Encoder:** Run through an encoder (often reduce the number of frames)
- **Text Decoder:** Decode using sequence-to-sequence model or special-purpose decoder such as CTC



<https://github.com/espnet/espnet>

ASR Evaluation Metrics

- **Automatic evaluation:** word error rate

C=correct, S=substitution, D=deletion, I=insertion

correct:	this	is	some	recognized		speech	
recognized:	this		some	wreck	a	nice speech	
type:	C	D	C	S	I	I	C

$$\text{WER} = (S + D + I) / \text{reference length}$$

$$(2+1+1)/5 = 80\%$$

Speech Synthesis Modeling Pipeline

- **Text Encoder:** Encode text into representations for downstream use
- **Speech Decoder:** Predicts features of speech, such as frequency
- **Vocoder:** Turns spoken features into a waveform

ASR Evaluation Metrics

- **Automatic evaluation:** word error rate

C=correct, S=substitution, D=deletion, I=insertion

correct:	this	is	some	recognized		speech	
recognized:	this		some	wreck	a	nice speech	
type:	C	D	C	S	I	I	C

$$\text{WER} = (S + D + I) / \text{reference length}$$

$$(2+1+1)/5 = 80\%$$

Question Answering

Input X

Textual Question

Output Y

Answer

Task

Question Answering

what languages are spoken in switzerland

X |

All Images Maps News Videos More Settings Tools

Switzerland > Official languages

German	French
Romansh	Italian

what is the difference between weather and climate

X |

All Books News Videos Images More Settings Tools

About 191,000,000 results (0.46 seconds)

Whereas **weather** refers to short-term changes **in the atmosphere**, **climate** describes what the **weather** is like over a long period of time **in a specific area**. **Different regions** can have **different climates**. ... **Weather** tells you what to wear each day. **Climate** tells you what types of clothes to have in your closet. Mar 23, 2018

[www.ncei.noaa.gov > news > weather-vs-climate](http://www.ncei.noaa.gov/news/weather-vs-climate)

[What's the Difference Between Weather and Climate? | News ...](http://www.ncei.noaa.gov/news/whats-difference-between-weather-and-climate-news...)

QA over Knowledge Bases

QA over Text

Example Knowledge Base: WikiData

Richard Feynman

Discuss "Richard Feynman" Hide Empty Fields

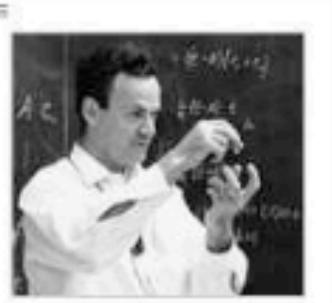


Image 1 of 1

Types: Person (People), Author (Publishing), Physicist (Science), Deceased Person (People), Film writer (Film), Influence Node (mikolov's types), Person Or Being In Fiction (Fictional Universes), Book Subject (Publishing)

Also known as: Richard Phillips Feynman

Gender: Male

Date of Birth: May 11, 1918

Place of Birth: Far Rockaway, Queens

Country Of Nationality: United States

Profession: Physicist, Scientist

Religion: Atheism

Parents: double-click to add

Children: Michelle Louise Feynman, Carl Feynman

Siblings:

Sibling
Joan Fey

Joan Feynman
Person

Richard Feynman ... (Richard Phillips Person, Author, Physicist, Deceased Person, Film ...)

Bell
Ana Gasteyer
H
Gervase of Tilbury
W
Alec Baldwin ... (Alexander Rae Person, Film actor, Film director, Film producer, TV ...)

P
Ernest Thesiger
C
Mean Girls
D
Riverside Drive
P
Portrait of Jennie
T
Televisions Personalities ... (The Television ...)

Description
Create New Person

Page History
Created by Metaweb Oct 22, 2006
Last edited by robert Oct 29, 2007

Web Link(s)
double-click to add

Employment history
Cornell University
California Institute of Technology
Thinking Machines

Education
Princeton University • 1942 • Ph.D.
Massachusetts Institute of Technology • 1939 • Bachelor's degree

Quotations
; like sex: sure, it may give some results, but that's not why we do not create, I do not understand.

lived
California
Los Angeles, New Mexico

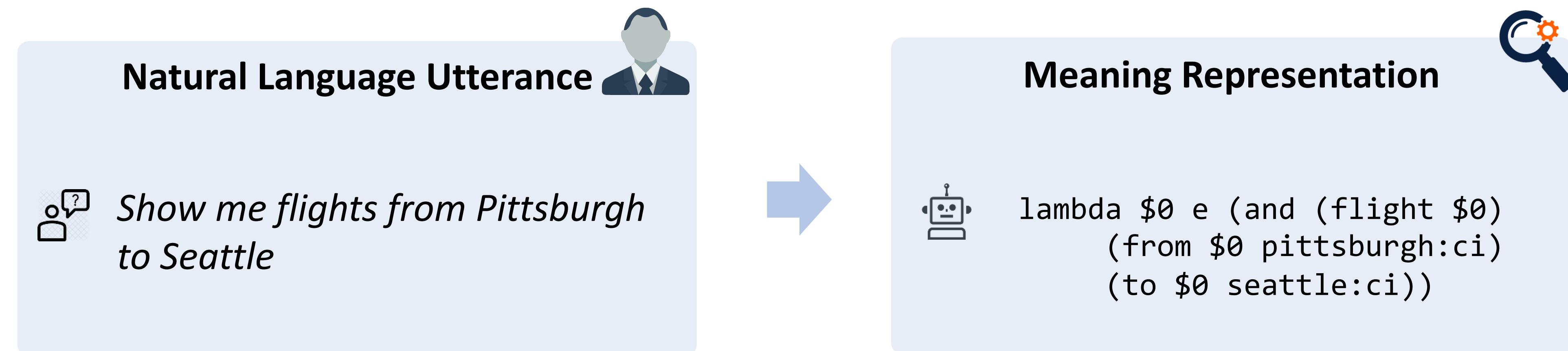
Books Written
What Do You Care What Other People Think?
The Pleasure of Finding Things Out
The Feynman Lectures on Physics
Surely You're Joking, Mr. Feynman!

Language	Label	Description	Also known as
English	Richard Feynman	American theoretical physicist	Richard Phillips Feynman Richard P. Feynman Ofey
Japanese	リチャード・P・ファインマン	アメリカ合衆国の物理学者	Feynman ファインマン リチャード・ファインマン
Spanish	Richard Feynman	físico estadounidense y premio Nobel	Richard Phillips Feynman Feynman Richard P. Feynman Richard P Feynman
Traditional Chinese	No label defined	No description defined	
Afrikaans	Richard Feynman	No description defined	
Amharic	ፋይናንም	No description defined	ፋይናንም
Aragonese	Richard Feynman	No description defined	
Arabic	ريتشارد فاينمان	فيزيائي أمريكي	فیزیائی امریکی
Egyptian Arabic	ريتشارد فاينمان	No description defined	
Assamese	ରିଚାର୍ଡ ଫାଇନମନ	No description defined	
Asturian	Richard Feynman	No description defined	
Azerbaijani	Riçard Feynman	No description defined	
South Azerbaijani	ريچارد فاينمان	No description defined	
Belarusian	Рычард Філіп Фейнман	No description defined	Рычард Фейнман Рычард Файнман
Bulgarian	Ричард Файнман	Американски физик	Ричард Филип Фейнман Ричард Фейнман Фейнмън Фейнман Ричард Филип Файнмън Ричард Файнмън Ричард Фейнмън Файнмън

<https://www.wikidata.org/>

Semantic Parsing

- The process of converting natural language to a more abstract, and often operational semantic representation

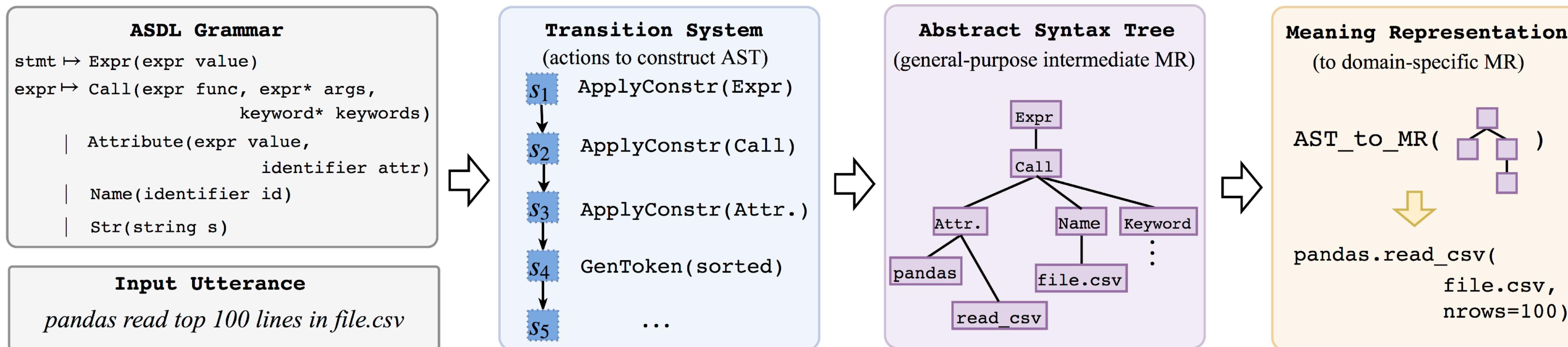


- These can be used to query **databases** (SQL), **knowledge bases** (SPARQL), or even generate **programming code** (Python)

Semantic Parsing Modeling Pipeline

- **Text Encoder:** Encode text into representations for downstream use
- **Tree/Graph Decoder:** Predict a tree or graph structured

TranX

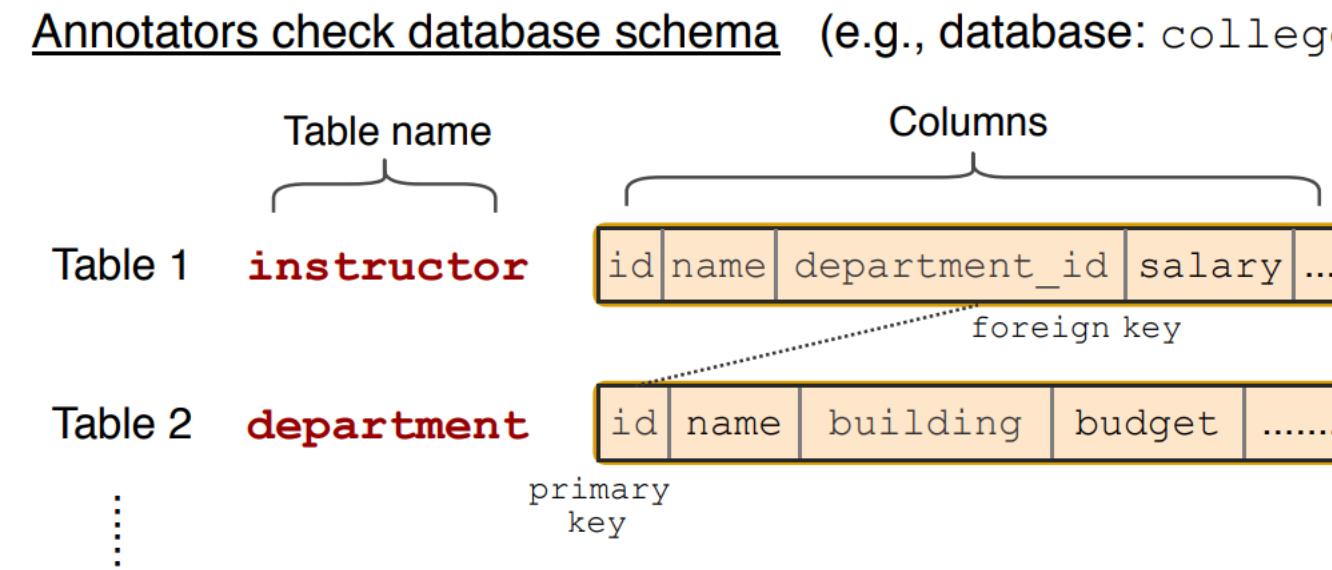


<https://github.com/pcyin/tranX>

Semantic Parsing Datasets

- **Text-to-SQL:** WikiSQL, Spider datasets
- **Text-to-knowledge graph:** WebQuestions, ComplexWebQuestions
- **Text-to-program:** CoNaLa, CONCODE

Spider



Annotators create:

Complex question What are the name and budget of the departments with average instructor salary greater than the overall average?

Complex SQL

```
SELECT T2.name, T2.budget
FROM instructor as T1 JOIN department as T2
ON T1.department_id = T2.id
GROUP BY T1.department_id
HAVING avg(T1.salary) >
(SELECT avg(salary) FROM instructor)
```

CoNaLa

```
{
  "question_id": 36875258,
  "intent": "copying one file's contents to another in python",
  "rewritten_intent": "copy the content of file 'file.txt' to file 'file2.txt'",
  "snippet": "shutil.copy('file.txt', 'file2.txt')",
}

{
  "intent": "How do I check if all elements in a list are the same?",
  "rewritten_intent": "check if all elements in list `mylist` are the same",
  "snippet": "len(set(mylist)) == 1",
  "question_id": 22240602
}
```

<https://yale-lily.github.io/spider>

<https://conala-corpus.github.io/>

Example Tasks/Datasets for QA over Text

Span Selection (SQuAD)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **graupel** and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain in scattered locations are called "showers".

What causes precipitation to fall?
gravity

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?
graupel

Where do water droplets collide with ice crystals to form precipitation?
within a cloud

Multiple Choice (MCTest)

James the Turtle was always getting in trouble. Sometimes he'd reach into the freezer and empty out all the food. Other times he'd sled on the deck and get a splinter. His aunt Jane tried as hard as she could to keep him out of trouble, but he was sneaky and got into lots of trouble behind her back.

One day, James thought he would go into town and see what kind of trouble he could get into. He went to the grocery store and pulled all the pudding off the shelves and ate two jars. Then he walked to the fast food restaurant and ordered 15 bags of fries. He didn't pay, and instead headed home.

His aunt was waiting for him in his room. She told James that she loved him, but he would have to start acting like a well-behaved turtle.

After about a month, and after getting into lots of trouble, James finally made up his mind to be a better turtle.

1) What is the name of the trouble making turtle?

- A) Fries
- B) Pudding
- C) James
- D) Jane

2) What did James pull off of the shelves in the grocery store?

- A) pudding
- B) fries
- C) food
- D) splinters

Cloze (CNN Daily Mail)

Original Version

Context

The BBC producer allegedly struck by Jeremy Clarkson will not press charges against the “Top Gear” host, his lawyer said Friday. Clarkson, who hosted one of the most-watched television shows in the world, was dropped by the BBC Wednesday after an internal investigation by the British broadcaster found he had subjected producer Oisin Tymon “to an unprovoked physical and verbal attack.” ...

Query

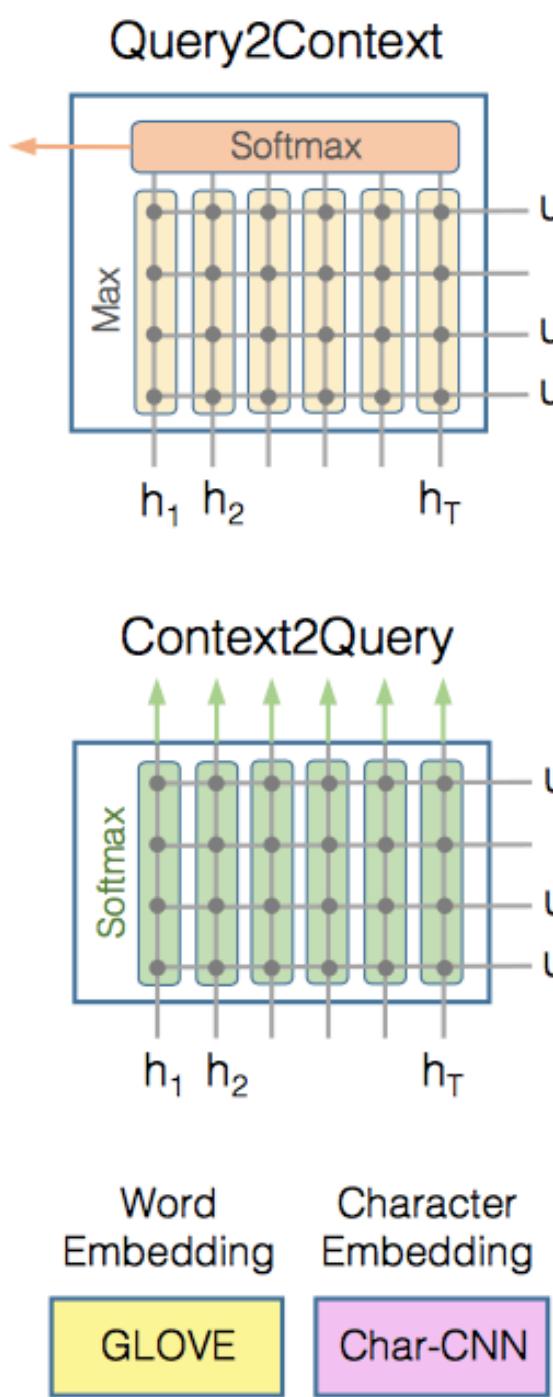
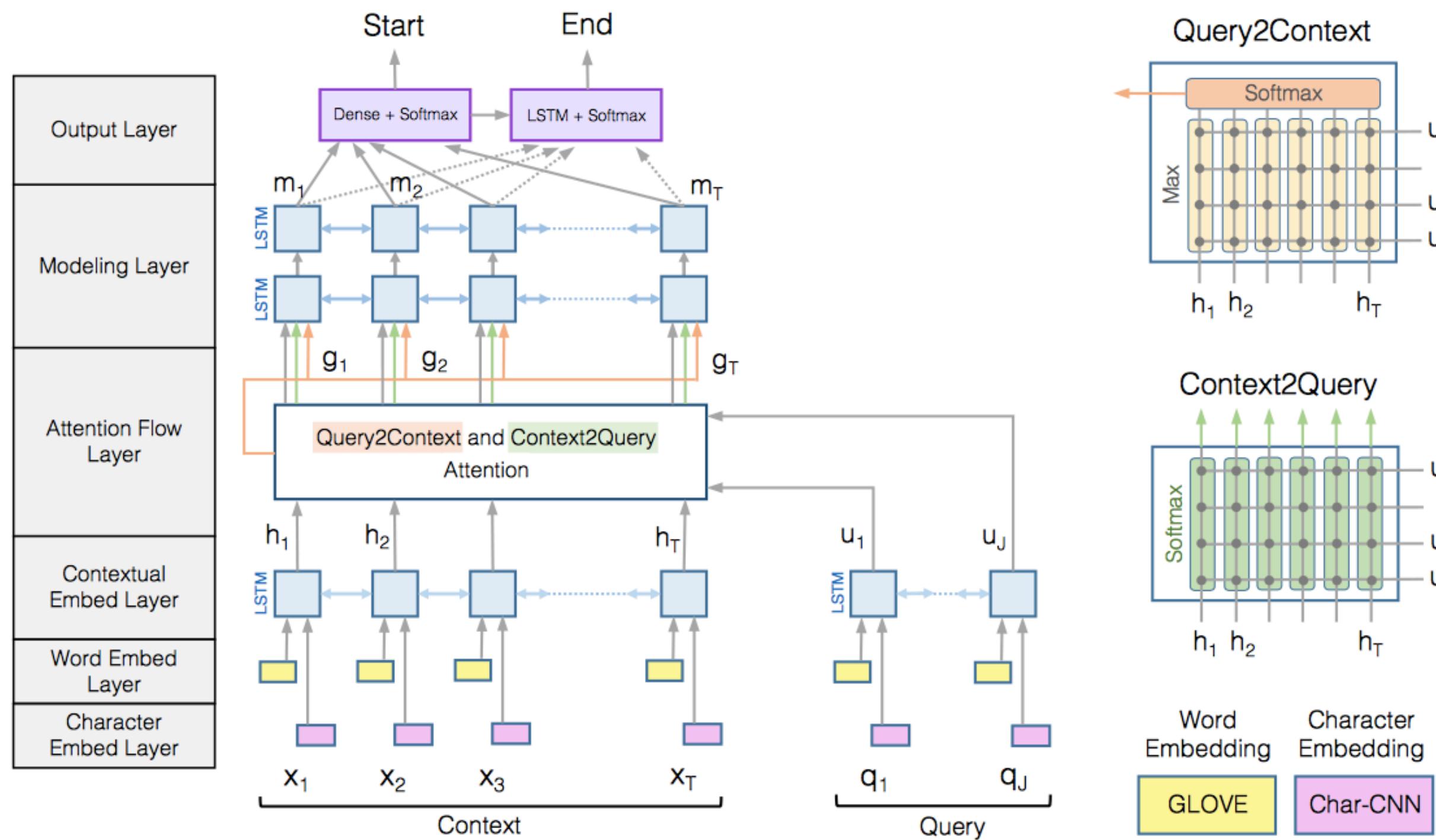
Producer X will not press charges against Jeremy Clarkson, his lawyer says.

Answer

Oisin Tymon

Machine Reading Modeling Pipeline

- **Document Encoder:** Encode text into representations for downstream use
- **Question Encoder:** Encode the question into some usable representation
- **Matcher:** Match between the input and output



<https://github.com/allenai/bi-att-flow>

Multilingual QA

- e.g. TyDiQA

Q: من هو موزارت ؟
 mn hw mwzArt ?
Who is Mozart?

A: أماديوس موتسارت
 >mAdyws mwtsArt
...Amadeus Mozart ...

Q: Kuka keksi viiko-n-päivä-t ?
 who invented week-GEN-day-PL ?
Who invented the days of the week?

A: Seitsem-päivä-inen viikko on
 seven-NOM-day-PL.ADJ week-NOM is
 todennäköisesti lähtöisin Babylonia-sta...
 likely origin Babylonia-ELA
The seven-day week is likely from Babylonia.

Q: Как далеко Уран от
 how far Uranus-SG.NOM from
 Земл-и?
 Earth-SG.GEN?
How far is Uranus from Earth?

A: Расстояние между Уран-ом
 distance between Uranus-SG.INSTR
 и Земл-ёй меняется от 2,6
 and Earth-SG.INSTR varies from 2,6
 до 3,15 млрд км...
 to 3,15 bln km...
The distance between Uranus and Earth fluctuates from 2.6 to 3.15 bln km...

LANGUAGE	LATIN SCRIPT ^a	WHITE SPACE TOKENS	SENTENCE BOUNDARIES	WORD FORMATION ^b	GENDER ^c	PRODROP
ENGLISH	+	+	+	+	+	—
ARABIC	—	+	+	++	+	+
BENGALI	—	+	+	+	+	+
FINNISH	+	+	+	+++	—	—
INDONESIAN	+	+	+	+	—	+
JAPANESE	—	—	+	+	—	+
KISWAHILI	+	+	+	+++	— ^e	+
KOREAN	—	+	+	+++	—	+
RUSSIAN	+	+	+	++	+	+
TELUGU	—	+	+	+++	+	+
THAI	—	—	—	+	+	+

<https://github.com/google-research-datasets/tydiqa>

Dialog Systems

Input X
Utterance

Output Y
Response

Task
Response Generation

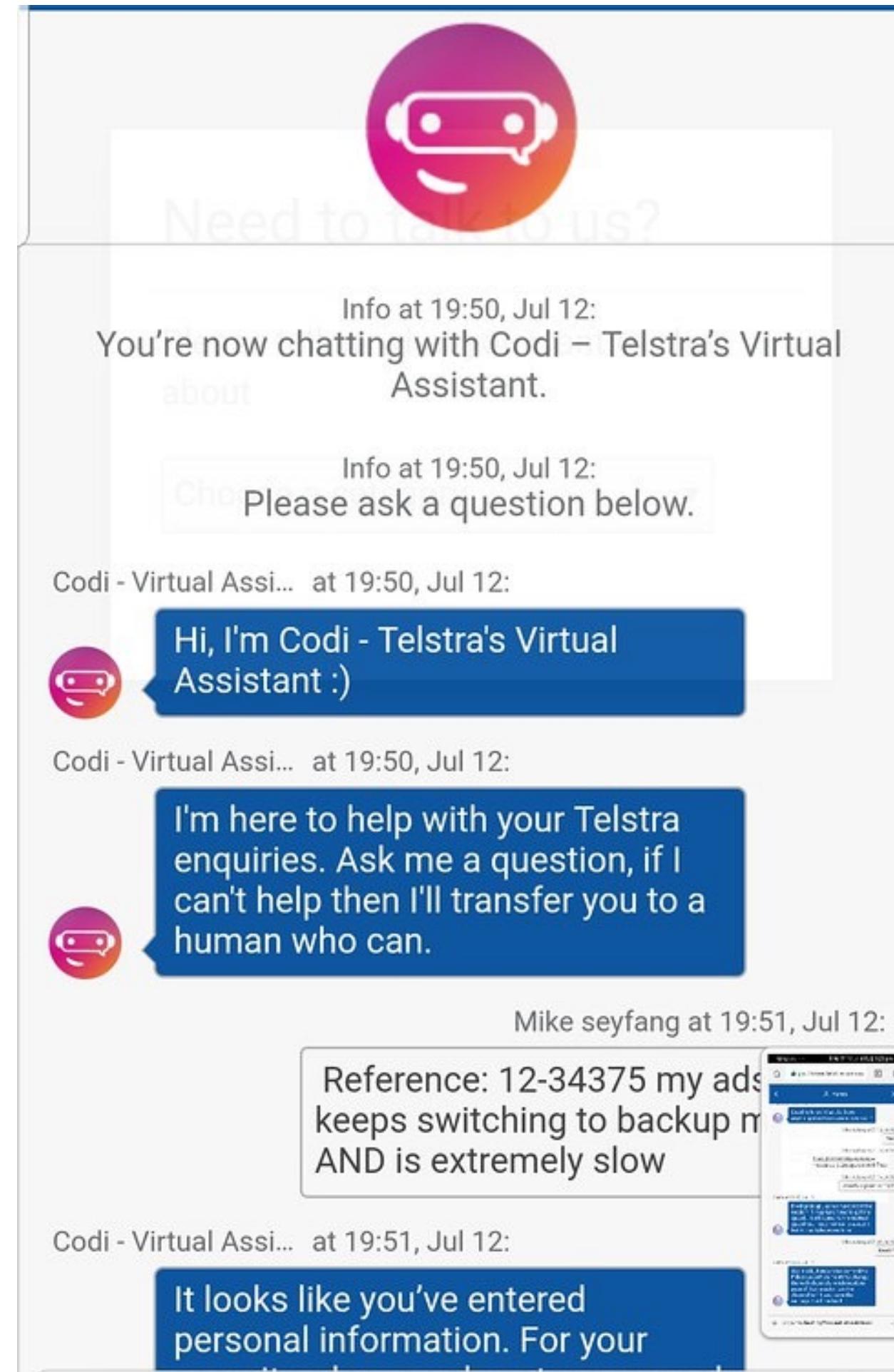


Image: Mike Seyfang @ flickr

Dialog System Data

Task-oriented Dialog

ASSISTANT: How can I help you?
USER: Hi, could you help me with booking movie tickets for tonight?
ASSISTANT: What movie are you interested in?
USER: The Upside.
ASSISTANT: Did you have a theater in mind?
USER: Could you check if the Regal Neshaminy...
No, AMC Neshaminy in Neshaminy, PA is playing it?
ASSISTANT: Could you spell that?
USER: Sure, n e s h a m i n y.
ASSISTANT: I have a showtime at 7:30 and at 10:30, is that okay?
USER: Yes, could you get two tickets for the 7:30?
ASSISTANT: One moment. Okay so that's 2 tickets for 7:30 at the AMC Neshaminy 24?
USER: Yes.
ASSISTANT: It'll be twenty-four ninety-nine for your tickets.
USER: That sounds great.
ASSISTANT: I've confirmed your tickets, they'll arrive via text shortly. Did you need any other information?
USER: No, that was it. Thank you so much for your help.
ASSISTANT: Great, no problem. I hope you have fun.
USER: I hope so, too. Thank you so much.

Chat

[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.

Naturally Occurring Sources of Conversation Data

- **Human-Machine Dialog:** Let's Go!, CMU Communicator
- **Human-Human Constrained Dialog:** Map Task, Debates
- **Human-Human Spontaneous Dialog:** Switchboard, AMI Meeting Corpus
- **Human-Human Scripted Dialog:** Movie Dialog
- **Human-Human Written Dialog:** Twitter, Reddit, Ubuntu Chat

<https://breakend.github.io/DialogDatasets/>

Dialogue Modeling Pipeline

- **Context Encoder:** Encode the entire previous context
- (optionally) **Explicit Belief Tracking, Database Access**
- **Utterance Decoder:** Generate the output utterance

Dialogue Evaluation

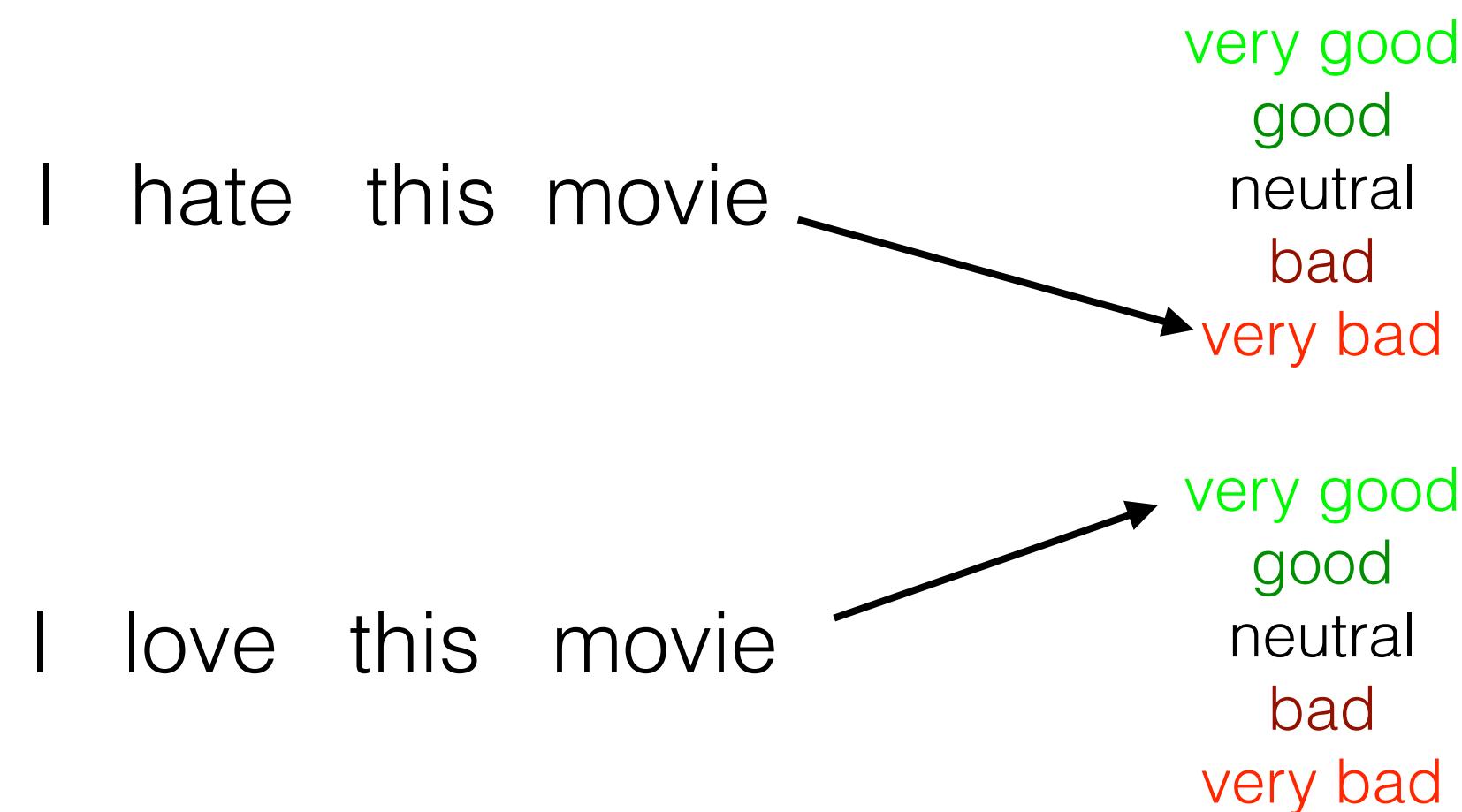
- **Task Oriented Dialog:**
 - Task completion rate
 - Time to task completion
 - User satisfaction
- **Free-form Dialog**
 - Attempts to use overlap with reference utterances not successful
 - Largely resort to human evaluation

Language Analysis

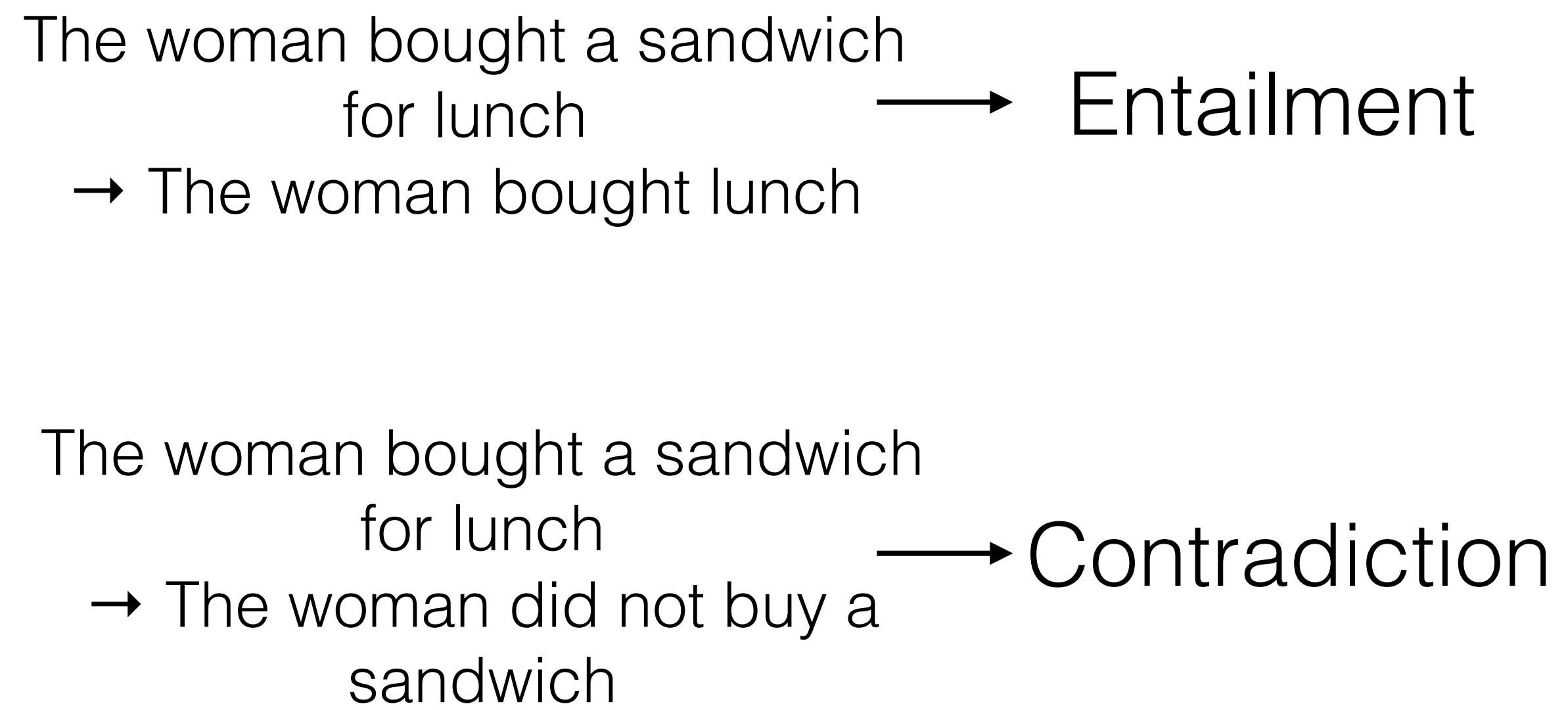
Text Classification

<u>Input X</u>	<u>Output Y</u>	<u>Task</u>	<u>Input X</u>	<u>Output Y</u>	<u>Task</u>
Text	Label	Text Classification	Text Pair	Label	Text Pair Classification

Sentiment Analysis



Textual Entailment



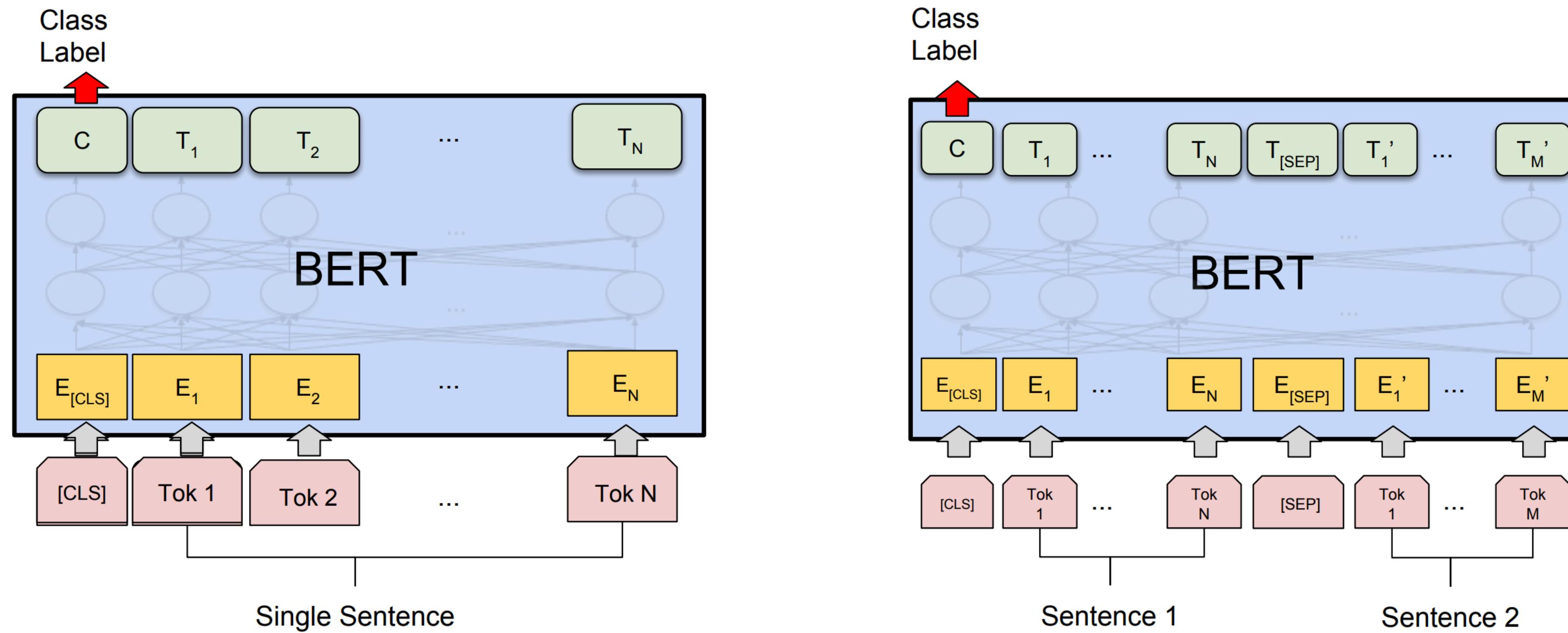
Text Classification Datasets

- **Sentiment Analysis:** Stanford sentiment treebank, Amazon reviews
- **Topic Classification:** 20 newsgroups, Wiki-500k
- **Paraphrase Identification:** Microsoft Paraphrase Corpus
- **Textual Entailment:** Stanford/Multi Natural Language Inference
- many many others!

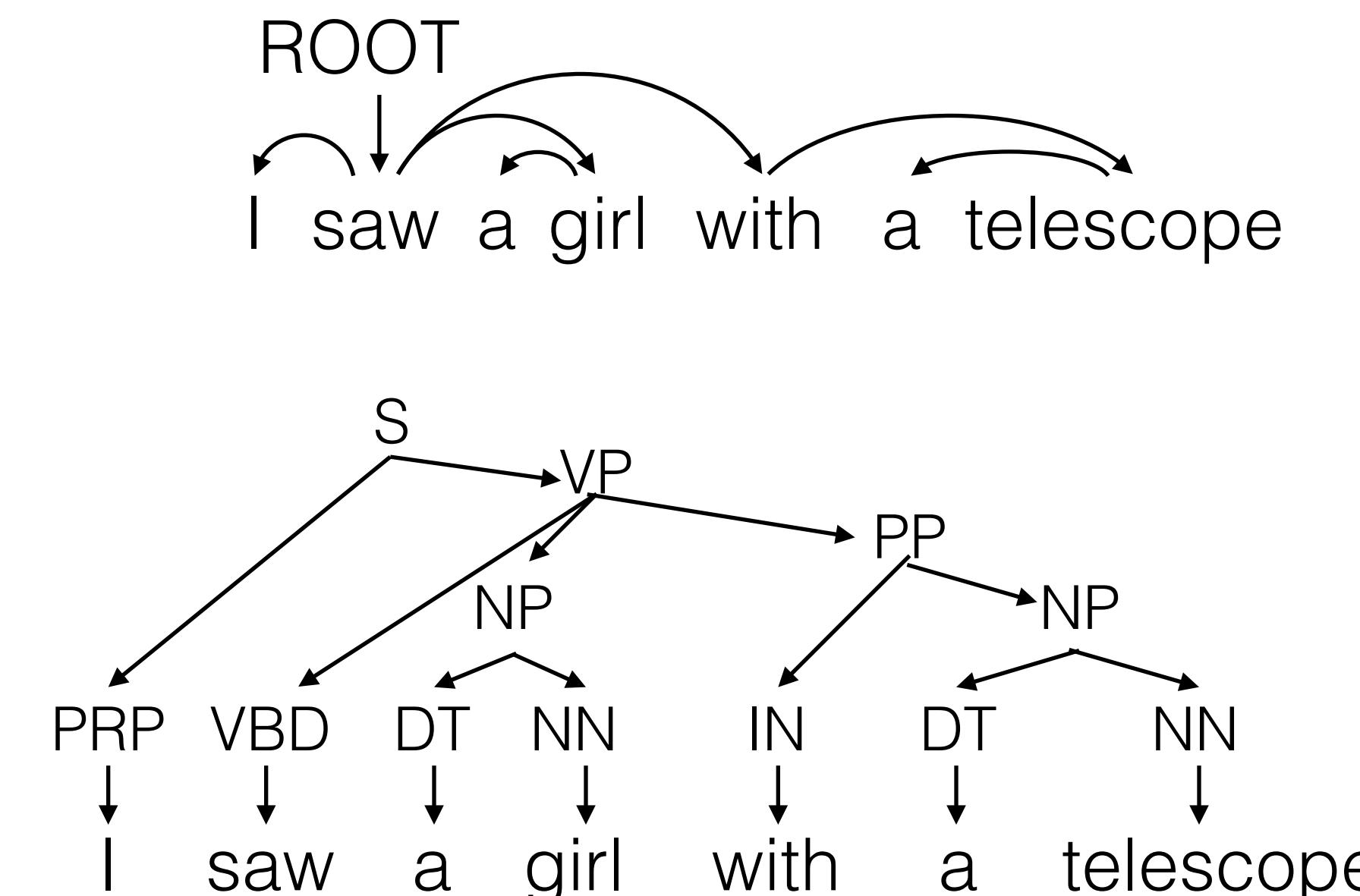
Text Classification Pipeline

- **Text Encoder:** Encode the text you want to analyze
- **Predictor:** Predict using a label

Example: BERT



Text Analysis Tasks

<u>Input X</u>	<u>Output Y</u>	<u>Task</u>	<u>Input X</u>	<u>Output Y</u>	<u>Task</u>
Text	Per-word Tags	Sequence Labeling	Text	Syntax Trees	Syntactic Parsing
I watched the movie	PRN VBD DET NN		I saw a girl with a telescope		
CMU is in Pittsburgh	ORG X X LOC				

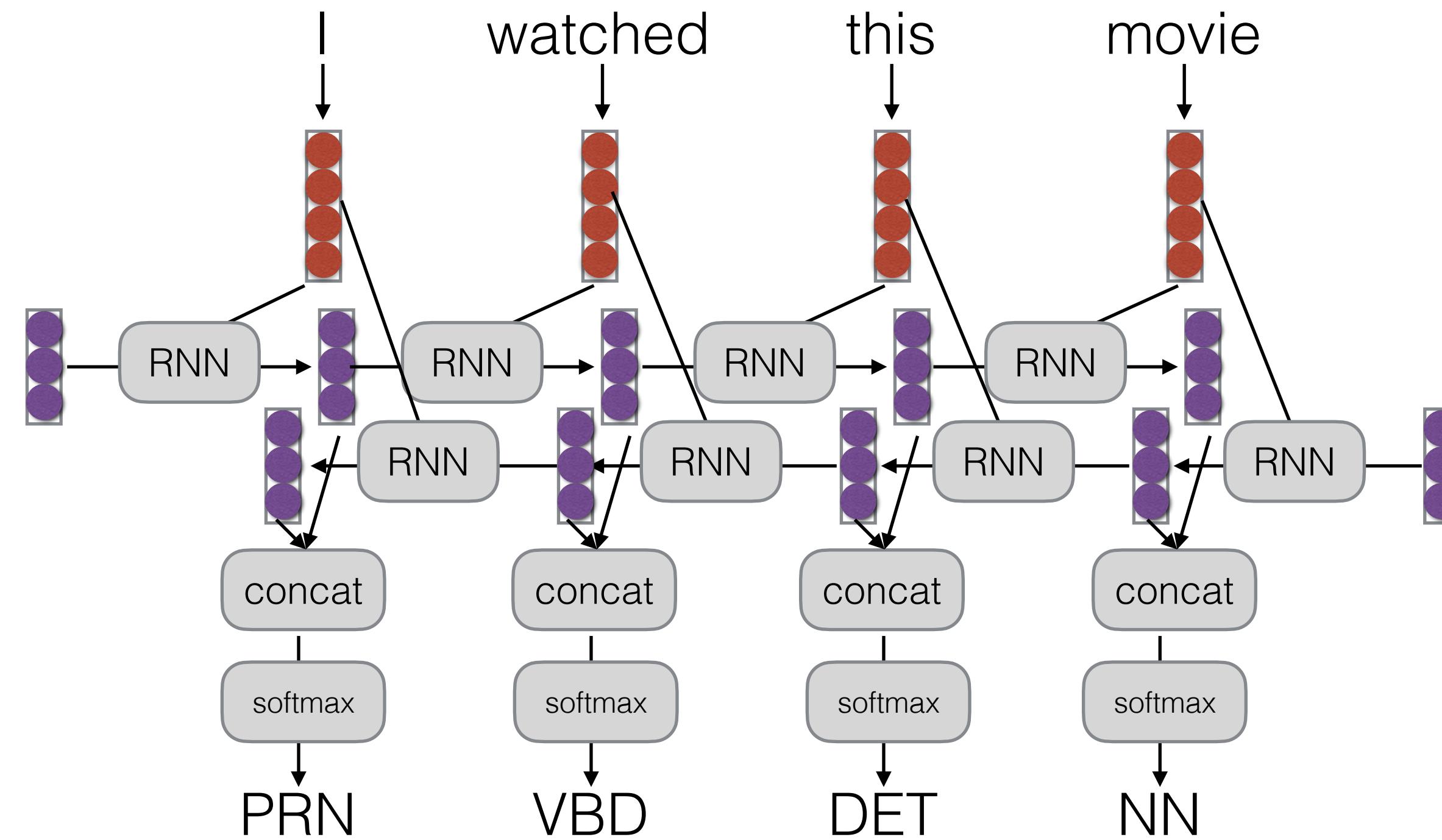
Text Analysis Data

- **OntoNotes:** A large corpus with many different annotations in English.
- **Universal Dependencies Treebank:** Includes parts-of-speech and dependency trees for many languages

▶  Abaza	1	3K		Northwest Caucasian
▶  Afrikaans	1	49K		IE, Germanic
▶  Akkadian	1	1K		Afro-Asiatic, Semitic
▶  Albanian	1	<1K		IE, Albanian
▶  Amharic	1	10K		Afro-Asiatic, Semitic
▶  Ancient Greek	2	416K		IE, Greek
▶  Arabic	3	1,042K		Afro-Asiatic, Semitic
▶  Armenian	1	52K		IE, Armenian
▶  Assyrian	1	<1K		Afro-Asiatic, Semitic
▶  Bambara	1	13K		Mande
▶  Basque	1	121K		Basque
▶  Belarusian	1	13K		IE, Slavic
▶  Bhojpuri	2	6K		IE, Indic
▶  Breton	1	10K		IE, Celtic
▶  Bulgarian	1	156K		IE, Slavic
▶  Buryat	1	10K		Mongolic
▶  Cantonese	1	13K		Sino-Tibetan

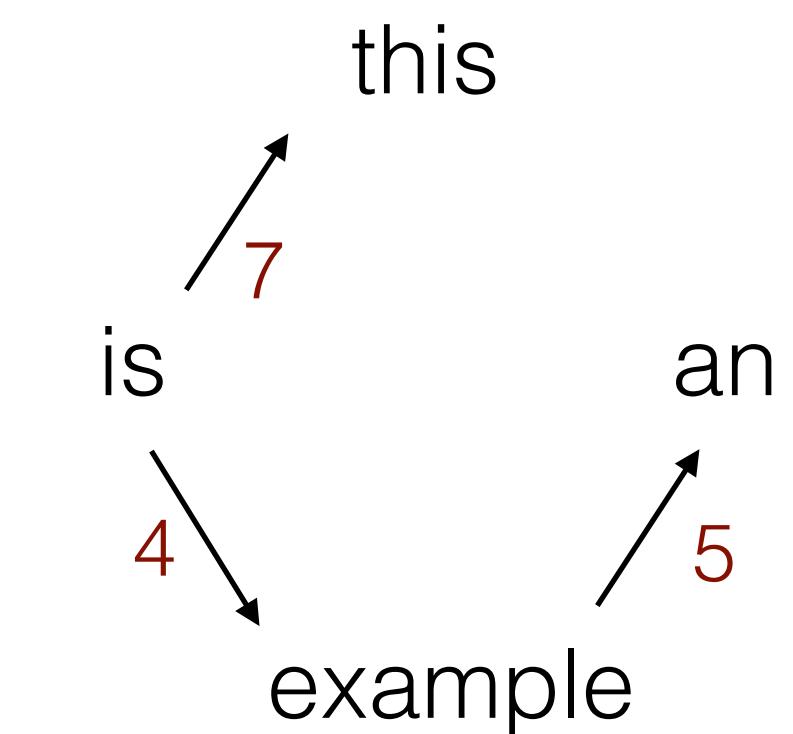
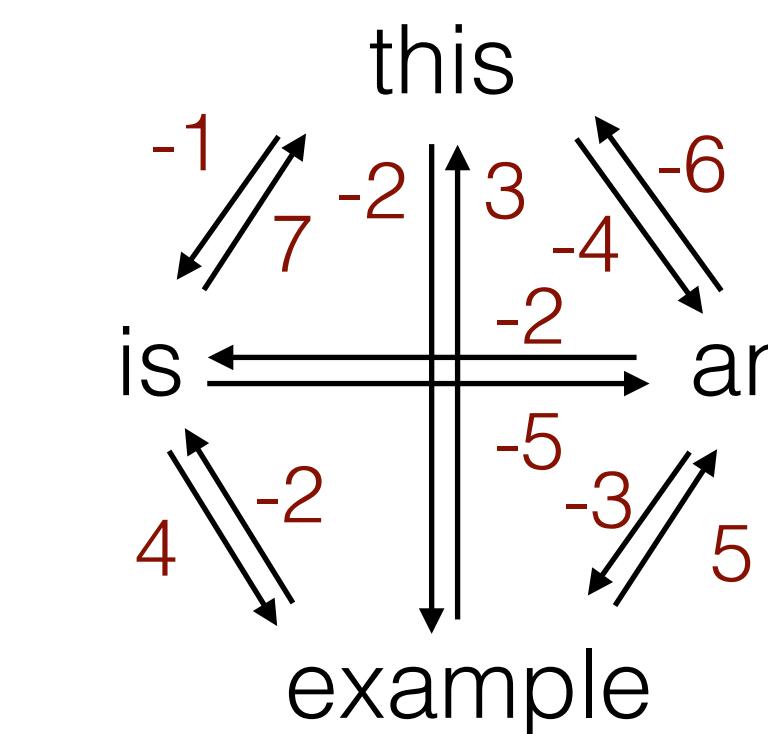
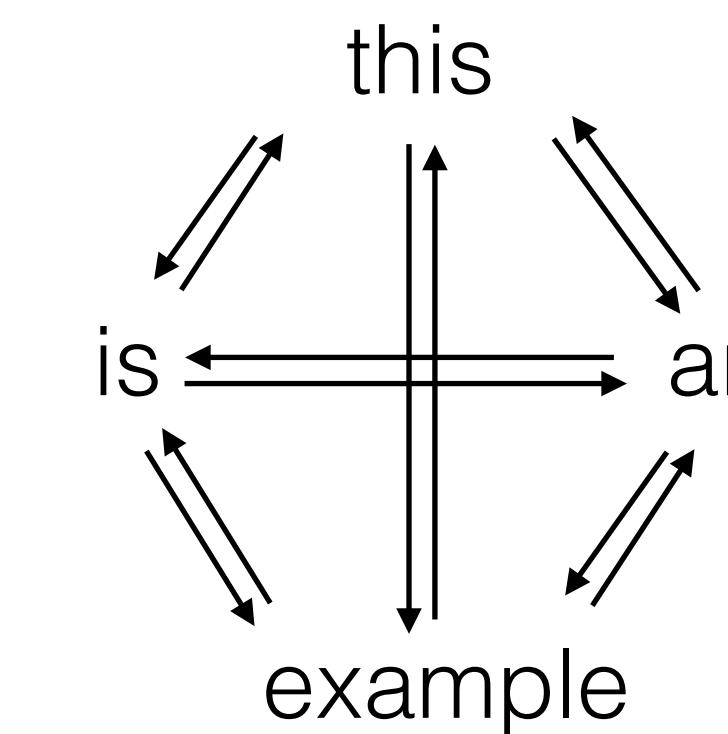
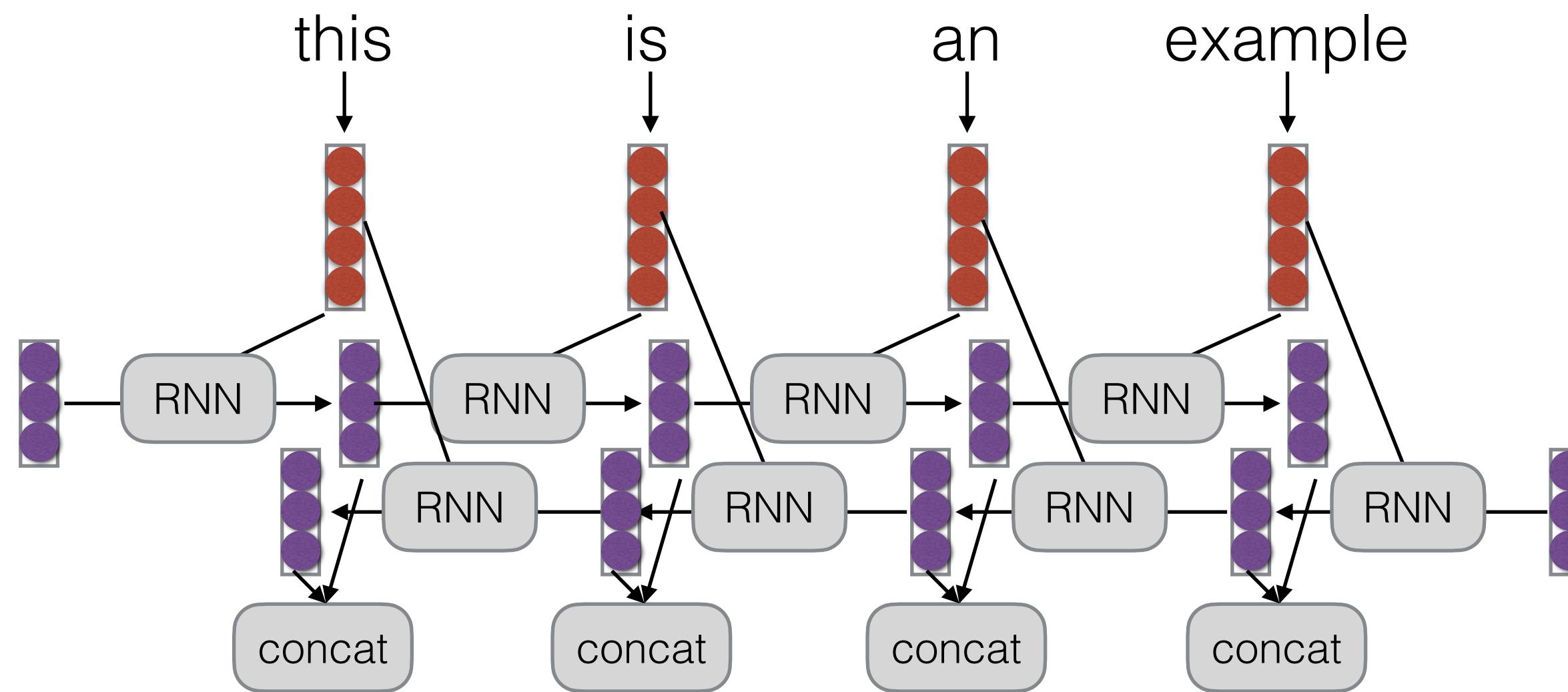
Sequence Labeling Pipeline

- **Text Encoder:** Encode the text you want to analyze
- **Sequence Labeler:** Predict each label independently, or in a joint fashion



Syntactic Parsing Pipeline

- **Text Encoder:** Encode the text you want to analyze
- **Tree Generation:** Generate dependency trees or constituency trees



Evaluation

- **POS Tagging:** Word-by-word accuracy
- **Entity Recognition:** Entity F-measure
- **Dependency Parsing:** Labeled or unlabeled attachment score
- **Constituency Parsing:** Phrase F-measure

Conclusion

Tackling NLP Tasks

- Data collection
- Modeling
- Train/test evaluation

Questions?