

**TRANSLATION IS LIKE CHOPPING  
AN ONION -  
FIRST, YOU THINK YOU'LL  
MANAGE IT.**

**AND THEN YOU END UP  
CRYING IN THE KITCHEN.**



# CS 2731 Introduction to Natural Language Processing

## Session 22: Machine translation part 1

---

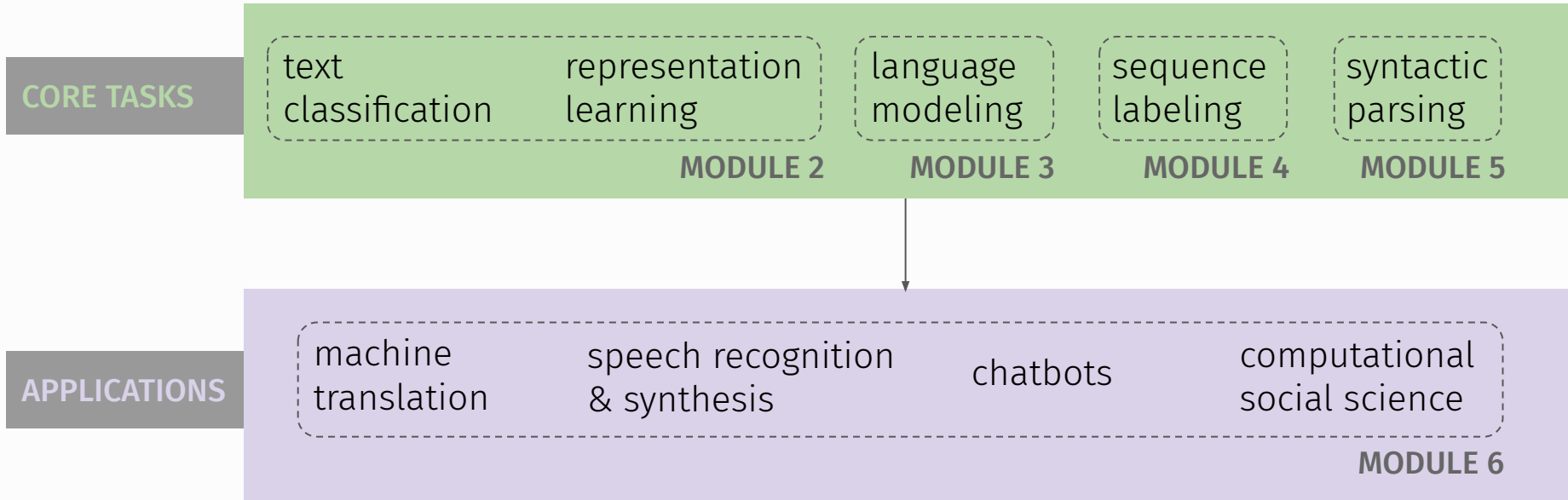
Michael Miller Yoder

November 13, 2023

# Course logistics

- Basic working project system **due this Thu 11-16**
  - 1-2 pages, in ACL LaTeX format that final report will be in
- Office hours the same times, but switching instructor/TA
  - Michael's office hours this week:  
2:45-3:45pm Tue 11-14, Sennott Square 6505
  - Pantho's office hours this week:  
1:30-2:30pm Wed 11-15, Sennott Square 5106
- Pantho will be giving the lecture on Wed

# Core tasks and applications of NLP



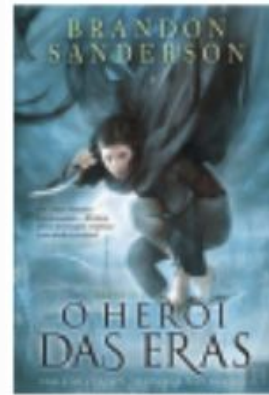
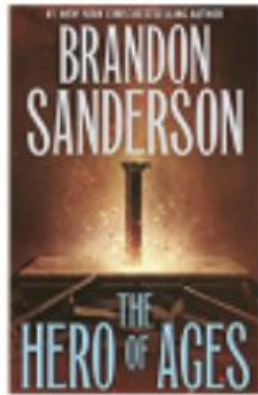
# Overview: Machine translation part 1

- History of machine translation (MT)
- Translation in practice
- Why is translation difficult?
- Parallel corpora
  - Sentence alignment

# Translation

- Mapping a “text” in a source language to a target language

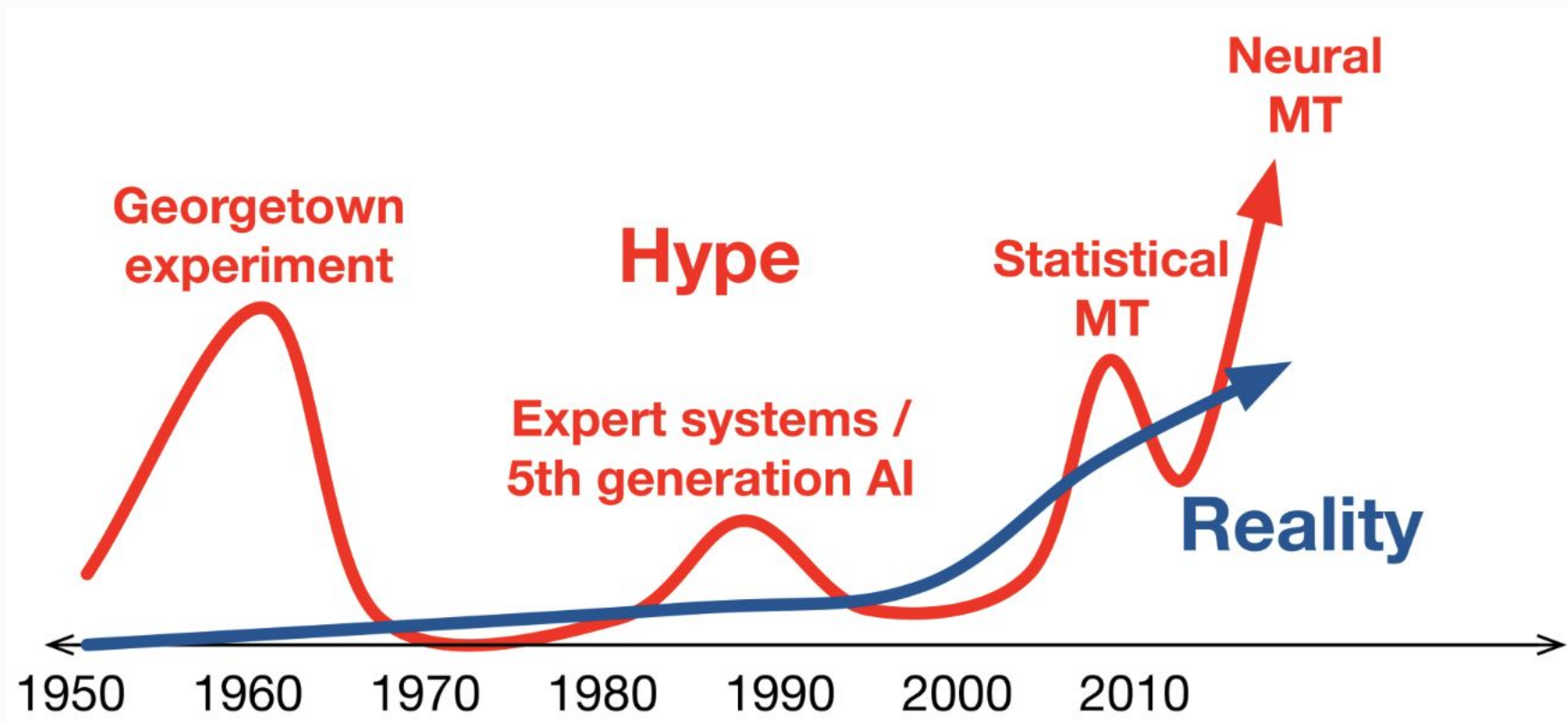
“I went to the store to buy eggs” → “Eu fui à loja comprar ovos”



# History of machine translation

---

# MT history: hype vs reality





# When did people start using computers to translate?



- Roughly 1950s
- Research stopped in the US for about 15-20 years after a 1967 report deemed it impossible
- Research resumed in the US in the early 1980s

# What did early MT systems look like?

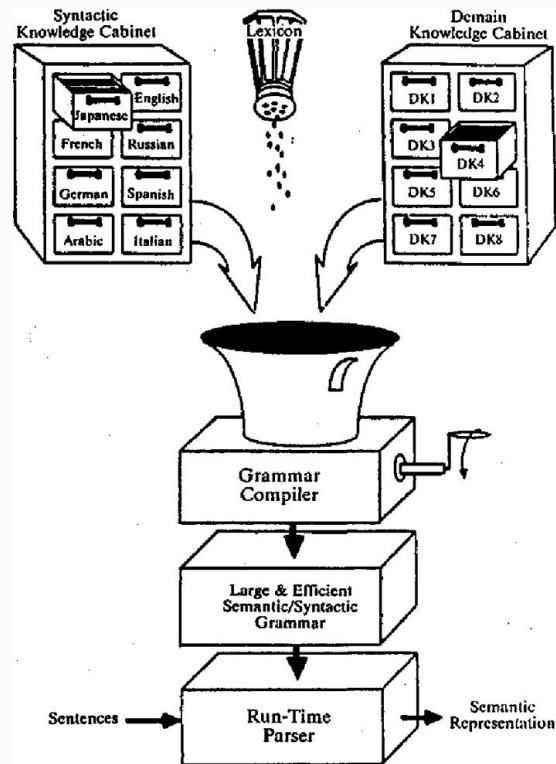
Human linguists wrote elaborate rules involving syntax, semantics, etc

```
(<S> <--> (<V>)  
  ((x0 = x1)))
```

```
(<S> <--> (<NP> <S>)  
  ((x2 subj-case) = *defined*)  
  ((x2 subj-case) = (x1 case))  
  (x0 = x2)  
  ((x0 subj) = x1)))
```

```
(<S> <--> (<NP> <S>)  
  ((x2 obj-case) = *defined*)  
  ((x2 obj-case) = (x1 case))  
  (x0 = x2)  
  ((x0 obj) = x1)))
```

```
(emap *insert  
  <=l=> insert ((CAT v) (SUBCAT trans))  
  (role =sem (*physical-action))  
  (:agent =syn (SUBJECT))  
  (:theme =syn (DOBJECT))  
  (:goal =syn (PPADJUNCT  
    ((PREP into) (CAT n))))))
```



# Learning to translate from data

Since the late 1980s, Machine Translation researchers have been using parallel corpora to train Machine Translation systems.

	ENGLISH	MANDARIN
1	i <b>wanna</b> live in a wes anderson world	我想要生活在Wes Anderson的世界里
2	Chicken soup, corn never truly digests. <b>TMI.</b>	鸡汤吧，玉米神马的从来没有真正消化过.恶心
3	To DanielVeuleman <b>yea iknw imma</b> work on that	对DanielVeuleman说，是的我知道，我正在向那方面努力
4	<b>msg 4</b> Warren G his <b>cday</b> is today 1 yr older.	发信息给Warren G，今天是他的生日，又老了一岁了。
5	Where <b>the hell</b> have you been all these years?	这些年你 <b>TMD</b> 到哪去了
	ENGLISH	ARABIC
6	It's <b>gonna</b> be a warm week!	الاسبوع الياي حر
7	onni this gift only <b>4 u</b>	أوني هذة الهدية فقط لك
8	sunset in aqaba :)	غروب الشمس في العقبة:)
9	RT @MARYAMALKHAWAJA: there is a call for widespread protests in #bahrain <b>tmrw</b>	هناك نداء لمظاهرات في عدة مناطق غدا

# Statistical machine translation (1990s-2010s)

- Core idea: Learn a probabilistic model from data
- For French  $\rightarrow$  English, we want to find best English sentence  $y$ , given French sentence  $x$
- Use Bayes Rule to break this down into two components to be learned separately:

$$\operatorname{argmax}_y P(y|x) = \operatorname{argmax}_y \underbrace{P(x|y)} \underbrace{P(y)}$$

## Translation Model

Models how words and phrases should be translated (*fidelity*).

Learned from parallel data.

## Language Model

Models how to write good English (*fluency*).

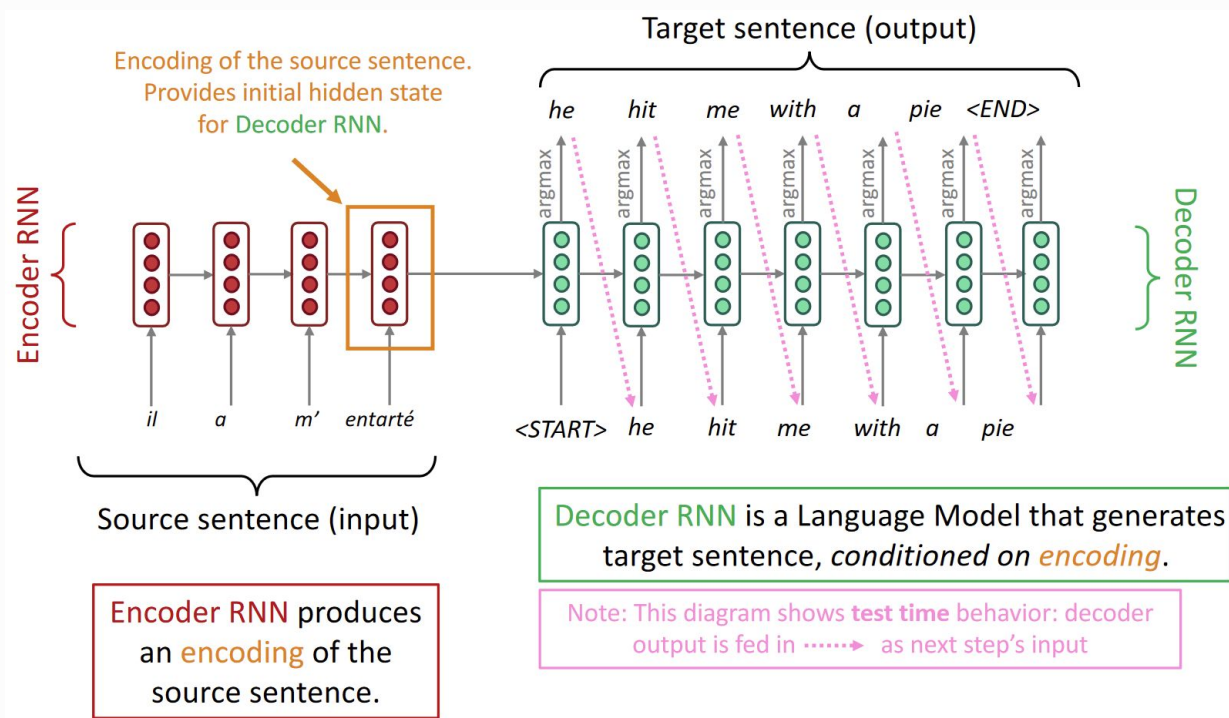
Learned from monolingual data.

# Statistical machine translation (1990s-2010s)

- The best SMT systems were extremely complex
  - Hundreds of important details
- Systems had many separately-designed subcomponents
  - Lots of feature engineering
  - Need to design features to capture particular language phenomena
- Required compiling and maintaining extra resources, like tables of equivalent phrases
  - Lots of human effort to maintain
- Repeated effort for each language pair

# Neural machine translation (2010s on)

- Single end-to-end neural network
- Encoder-decoder (sequence-to-sequence, seq2seq) framework



# Translation in practice

---

# Machine translation is a \$3 billion market

## Translation of text

The screenshot shows the Google Translate interface. At the top left is the Google Translate logo. Below it are three tabs: 'Text' (selected), 'Documents', and 'Websites'. The language selection bar shows 'ENGLISH' selected on the left and 'JAPANESE' selected on the right. The input text is 'Machine translation is a \$3 billion market.' and the output is '機械翻訳は 30 億ドルの市場です。'. Below the output is the phonetic transcription 'Kikai hon'yaku wa 30 oku-doru no ichibadesu.' and a speaker icon for audio playback. The character count '43 / 5,000' is visible at the bottom of the input area.



# Machine translation is a \$3 billion market

## Translation of speech

Person: Alexa, how do you say, “I hate this movie” in Japanese.

Alexa: “I hate this movie” in Japanese is “Kono eiga wa kirai da.”

Person : Alexa, how do you say, “I hate this movie in Japanese” in Japanese.

Alexa: “I hate this movie in Japanese” in Japanese is “Kono eiga wa nihongo de wa kirai da.”

**Real time translation of meetings is also now viable.**

# Most translation is still done by human translators

## Translation and Localization Industry Grows 11.8% in 2021 to USD 26.6bn



# Post-editing and computer-assisted translation

- Checking and correcting of machine translation by humans is called **post-editing**



Evacuation Ladder



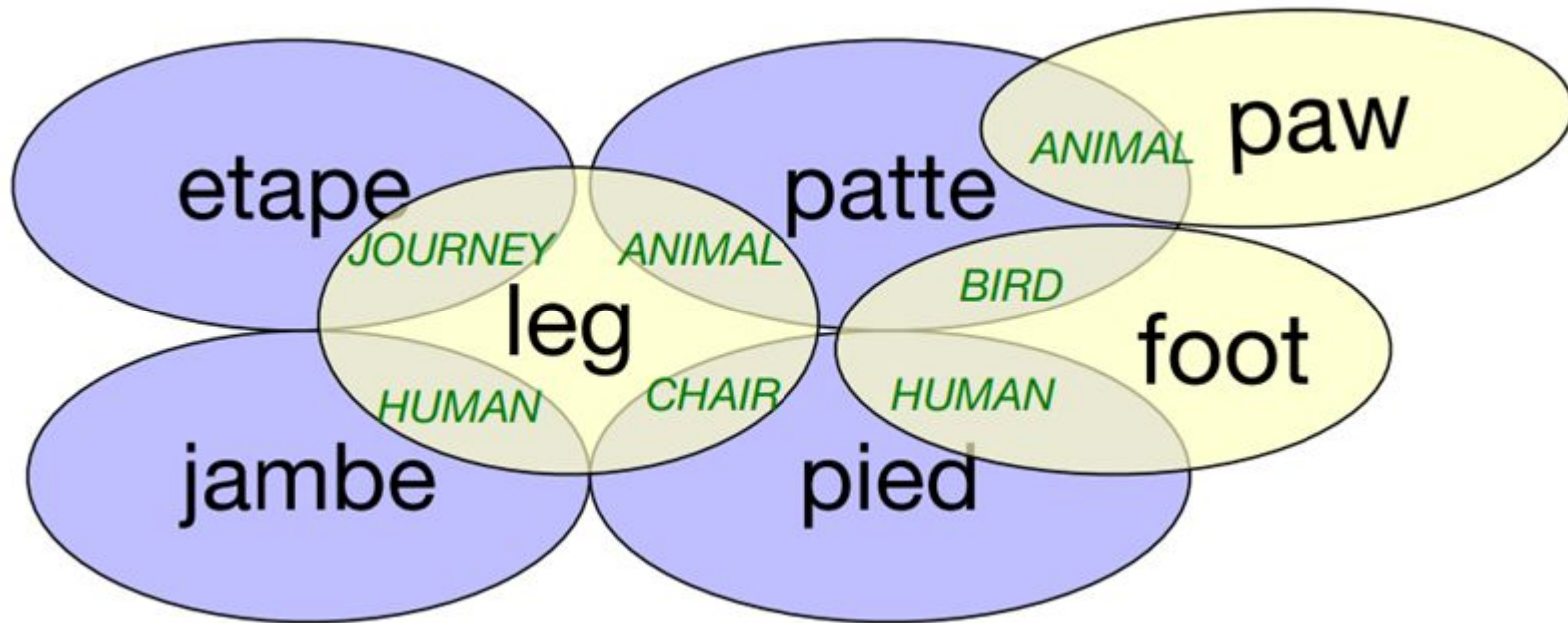
Do not yell

# Why is translation difficult?

---

# Why not just look up each word in a dictionary and translate word-for-word?

Many-to-many mappings of words



# Why not translate word-for-word: grammar distinctions

The grammars of some languages make distinctions that other languages don't make:

- Russian *kniga* translates to English as *the book* or *a book*.
  - English grammar makes a distinction in definiteness
  - Russian grammar does not.
- English *it* translates to French *il/le* (masculine) or *elle/la* (feminine).
- English *a* translates to French as *un* (masculine) or *une* (feminine).
  - *Une chaise* (a chair) vs *un livre* (a book)
  - French grammar makes a distinction in gender
  - English grammar does not.

# Why not translate word-for-word: Different numbers of words to say the same thing

uygarlaştıramadıklarımızdanmışsınızcasına

“(behaving) as if you are among those whom we were not able to civilize”

<u>uygar</u>	“civilized”
<u>+laş</u>	“become”
<u>+tır</u>	“cause to”
<u>+ama</u>	“not able”
<u>+dık</u>	past participle
<u>+lar</u>	plural
<u>+ımız</u>	first person plural possessive (“our”)
<u>+dan</u>	ablative case (“from/among”)
<u>+mış</u>	past
<u>+sınız</u>	second person plural (“y’all”)
<u>+casına</u>	finite verb → adverb (“as if”)

# Why not translate word-by-word: word order

English: *He wrote a letter to a friend* ← SVO (verb-medial)

Japanese: *tomodachi ni tegami-o kaita* ← SOV (verb-final)  
friend to letter wrote

Arabic: *katabt risāla li šadq* ← VSO (verb-initial)  
wrote letter to friend



# Exercise: Tajik

There are 3,344,720 speakers of *Tajik* in Tajikistan (one of the Central Asian republics of the former Soviet Union) and another million speakers in surrounding countries.

дуусти хуби ҳамсоҷай суро  
ҳамсоҷай дуусти хуби суро  
ҳамсоҷай хуби дуусти суро

a good friend of your neighbor  
a neighbor of your good friend  
a good neighbor of your friend

Above are three phrases in Tajik with their English translations. Your task is to give the English translations of all four Tajik words. The possibilities are simply "good," "friend," "neighbor," and "your." The order of the words – which is not the same order as in English! – does the rest.

дуусти \_\_\_\_\_  
ҳамсоҷай \_\_\_\_\_  
хуби \_\_\_\_\_  
суро \_\_\_\_\_

# What is difficult about translation?

- People in NLP and MT have reduced “language divergences” to six major word order features from WALS, or seven lexical features
- But language typology is a system of “morphosyntactic strategies”, of which there are 1000s



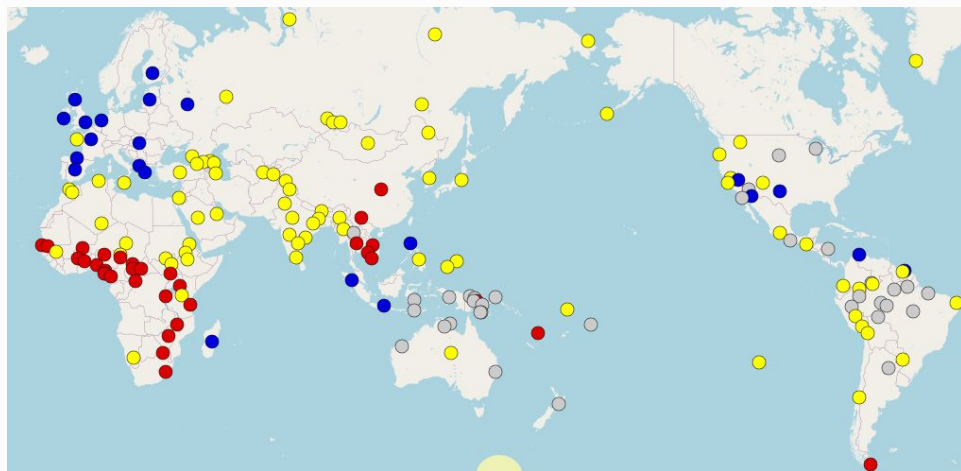
## Feature 121A: Comparative Constructions

Yellow: X is big from Y, or X is big to Y

Red: X is big, exceeds Y

Grey: X is big, Y is small

Blue: X is big than Y



# But the picture is not so gloomy

- MT researchers have made much progress on handling language divergence
- Use data from typologically similar languages
- Use a multilingual model trained on many typologically different languages

# Why is translation difficult? Style and genre

錨玉自在枕上感念寶釵

dai yu zi zai zhen shang gan nian bao chai

From “Dream of the Red Chamber”, Cao Xue Qin (1792)

Chinese: Daiyu alone at bed top think baochai.

English: Daiyu alone on **the** bed thought **about** baochai.

# Why is translation difficult? Style and genre

錨玉自在枕上感念寶釵

dai yu zi zai zhen shang gan nian bao chai

From “Dream of the Red Chamber”, Cao Xue Qin (1792)

Chinese:

DAIYU ALONE ON BED TOP

THINK

BAOCHAI

English:

As she lay there alone Daiyu's thoughts turned to Baochai .

Parallel data is more likely to match styles (like literary style) than be an “exact” translation

# Preparing for machine translation

1. Collect a parallel corpus
2. Align sentences
3. Tokenization
  - Split words into sub-word units, e.g., using BPE (Byte Pair Encoding)

# Parallel corpora

---

## Bao - Pitt Campus

### Food

#### Appetizers 头台



**Tea Egg 茶叶蛋**  
\$4.00



**Pork Belly Slider 五花肉刈包**  
\$7.95



**Popcorn Chicken 盐酥鸡**  
\$8.95



**Cantonese Style Chicken Feet 广式凤爪**  
\$8.95



**Rolled Pancakes w/Roast Beef 牛肉卷饼**  
\$12.95



**Pan Fried Radish Cake 萝卜糕**  
\$7.95



**Crab Rangoon 蟹角**  
\$7.95



**Indian Pan Fried Pancake 印度薄煎饼**  
\$6.95



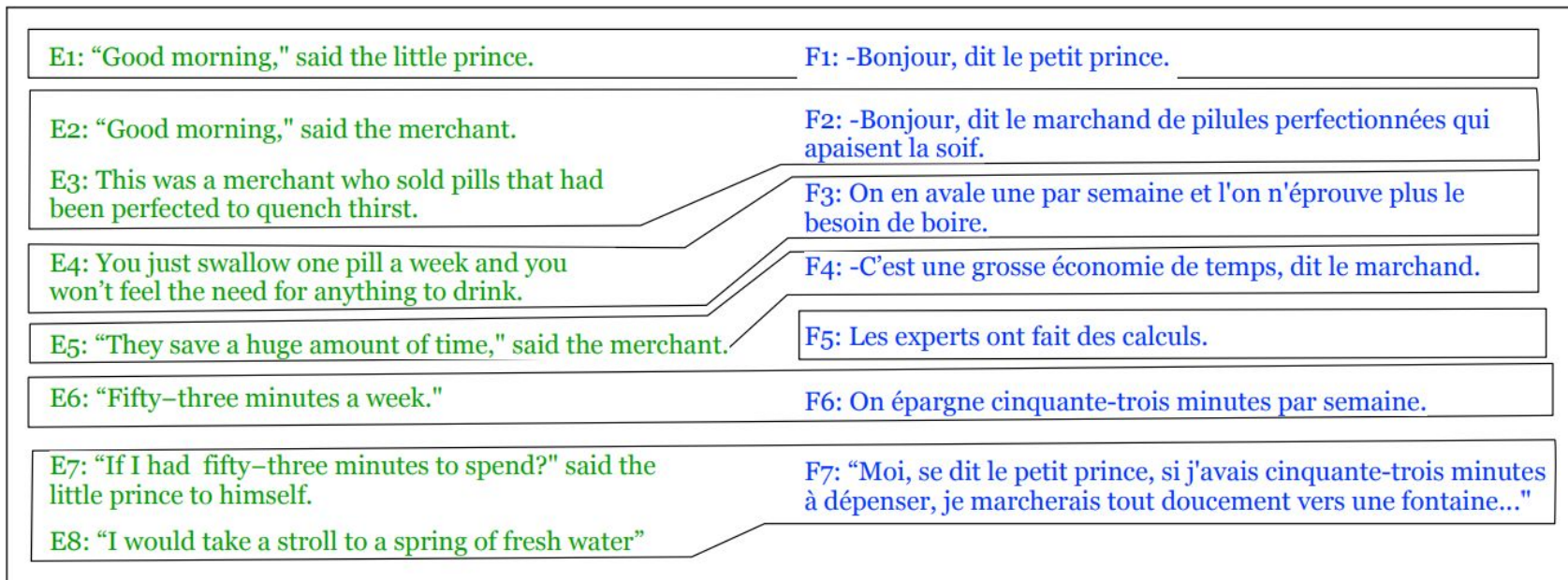
# Parallel corpora examples

- Europarl: Proceedings of the European Parliament; 21 languages; up to 2 million sentences
- United Nations Parallel Corpus: 10 million sentences in Arabic, Chinese, English, French, Russian, Spanish
- OpenSubtitles: movie and TV subtitles
- ParaCrawl: 223 million sentences in 23 EU languages

# What about parallel corpora for the other 7000 languages?

- For many languages, the only parallel text is the Christian Bible.
- Low-resource MT is a large area of research
  - How to leverage monolingual texts (backtranslation)
  - Humans in the loop
  - Leverage multilingual models

# Sentence alignment



**Figure 10.17** A sample alignment between sentences in English and French, with sentences extracted from Antoine de Saint-Exupéry's *Le Petit Prince* and a hypothetical translation. Sentence alignment takes sentences  $e_1, \dots, e_n$ , and  $f_1, \dots, f_n$  and finds minimal sets of sentences that are translations of each other, including single sentence mappings like  $(e_1, f_1)$ ,  $(e_4, f_3)$ ,  $(e_5, f_4)$ ,  $(e_6, f_6)$  as well as 2-1 alignments  $(e_2/e_3, f_2)$ ,  $(e_7/e_8, f_7)$ , and null alignments  $(f_5)$ .

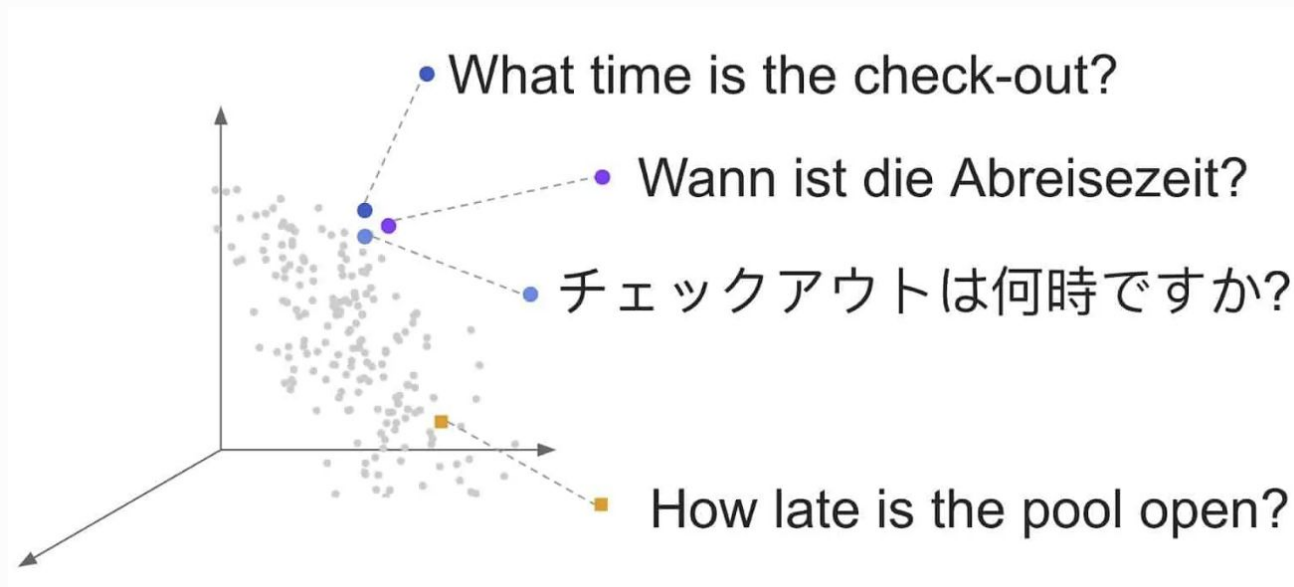
# How to align sentences

Need:

1. Cost function: how likely are a source language span and a target language span to be translations?
2. Alignment algorithm: uses scores between spans to find a good alignment between documents

# Multilingual embedding space

- Cost function: score similarity of sentences across languages with cosine similarity of embeddings in **multilingual embedding space**



# Sentence alignment: cost function and alignment alg

- Cost function using cosine similarity of embeddings in multilingual embedding space [Thompson + Koehn 2019]

$$c(x, y) = \frac{(1 - \cos(x, y))nSents(x) nSents(y)}{\sum_{s=1}^S 1 - \cos(x, y_s) + \sum_{s=1}^S 1 - \cos(x_s, y)}$$

- Dynamic programming algorithm [Gale + Church 1993] as the alignment algorithm
  - Minimize cost over the entire sequence of spans

# Subword tokenization review

- Create a shared vocabulary between source and target language with **subword tokenization**
- Example: Byte-pair encoding (BPE, Sennrich et al. 2016)
  - Merges frequently seen sequences of characters together into tokens
- More powerful alternatives
  - Wordpiece
    - Merge tokens based on what increases language model probability of a training corpus
  - SentencePiece/unigram
    - Start with huge vocabulary of all frequent sequences of characters, remove sequences that don't have a high probability in the training corpus iteratively

# Wrapping up

- Modern machine translation methods use the neural encoder-decoder framework
- MT is often used in conjunction with human translators
- Language divergences (in word meaning, syntax structure, etc) make MT difficult
- Parallel corpora are used for training MT systems
- Sentences must be aligned in parallel corpora
- Subword tokenization is used for a shared vocabulary between languages



*Questions?*