

# Machine Translation

A large, dense word cloud centered around the word "hello" in various languages. The words are in different colors and sizes, creating a visual representation of global communication. The most prominent words include "hello" in English, "bonjour" in French, "hallo" in German, "dobar" in Portuguese, "selamat" in Indonesian, "dobre" in Polish, and "namaste" in Sanskrit.



# Tower of Babel

Text

Documents

DETECT LANGUAGE

ENGLISH

SPANISH

FRENCH

^



ENGLISH

SPANISH

ARABIC

▼

← Search languages

 Detect language 

Danish

Hmong

Lithuanian

Romanian

Telugu

Afrikaans

Dutch

Hungarian

Luxembourgish

Russian

Thai

Albanian

English

Icelandic

Macedonian

Samoan

Turkish

Amharic

Esperanto

Igbo

Malagasy

Scots Gaelic

Turkmen

Arabic

Estonian

Indonesian

Malay

Serbian

Ukrainian

Armenian

Filipino

Irish

Malayalam

Sesotho

Urdu

Azerbaijani

Finnish

Italian

Maltese

Shona

Uyghur

Basque

French

Japanese

Maori

Sindhi

Uzbek

Belarusian

Frissian

Javanese

Marathi

Sinhala

Vietnamese

Bengali

Galician

Kannada

Mongolian

Slovak

Welsh

Bosnian

Georgian

Kazakh

Myanmar (Burmese)

Slovenian

Xhosa

Bulgarian

German

Khmer

Nepali

Somali

Yiddish

Catalan

Greek

Kinyarwanda

Norwegian

Spanish

Yoruba

Cebuano

Gujarati

Korean

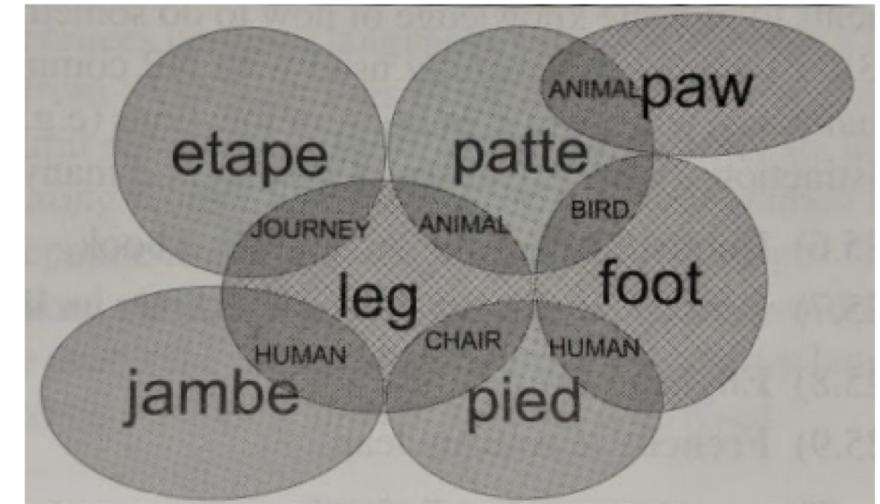
Odia (Oriya)

Sundanese

Zulu

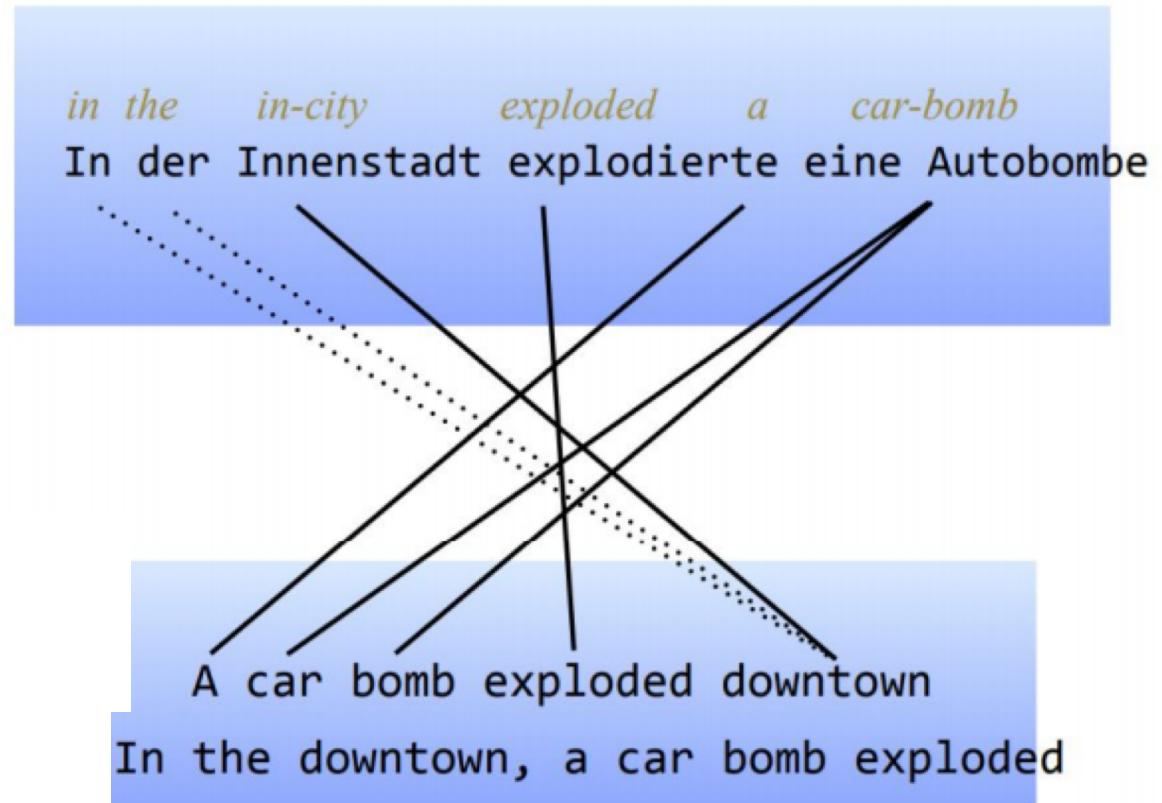
# Dictionaries

- English: leg, foot, paw
- French: jambe, pied, patte, etape



# Challenges

- Ambiguities
  - Words
  - Morphology
  - Syntax
  - Semantics
  - Pragmatics
- Gaps in data
  - Availability of corpus
  - Commonsense knowledge
- Understanding of context, connotation, social norms, etc



# Research Problems

- How can we formalize the process of learning to translate from examples?
- How can we formalize the process of finding translations for new inputs?
- If our model produces many outputs, how do we find the best one?
- If we have a gold standard translation, how can we tell if our output is good or bad?

# Two Views Of MT

---

# MT as Code Breaking

One naturally wonders if the problem of translation could conceivably be treated as a problem in cryptography. When I look at an article in Russian, I say: '*This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.*'



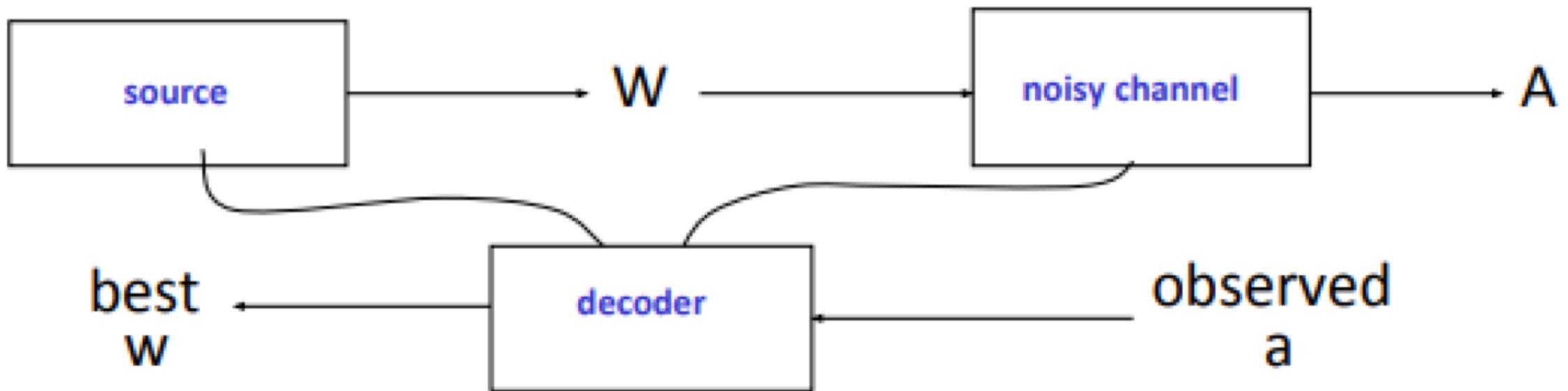
Warren Weaver to Norbert Wiener, March, 1947

# The Noisy-Channel Model

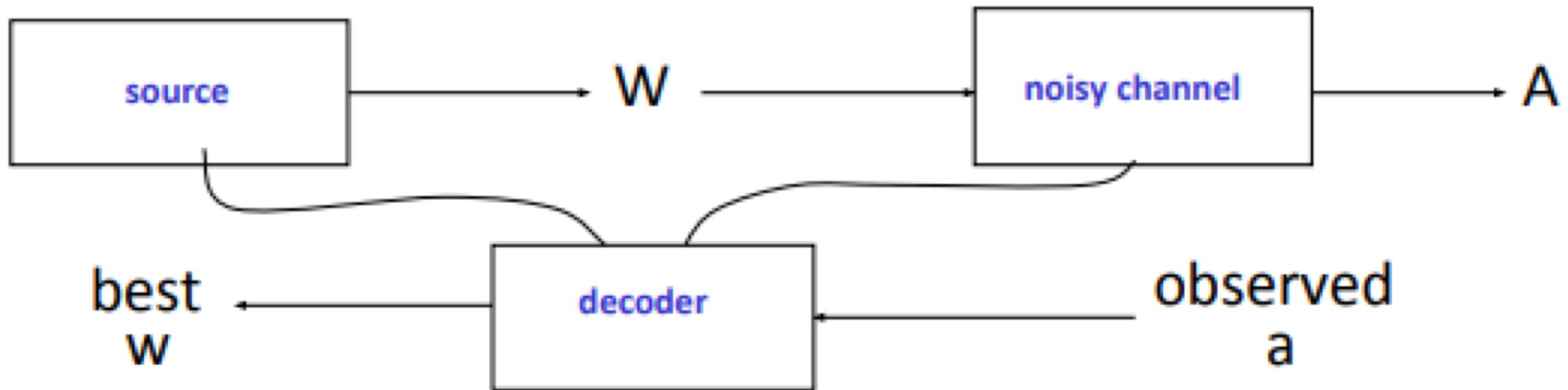


Claude Shannon. "A Mathematical Theory of Communication" 1948.

# The Noisy-Channel Model



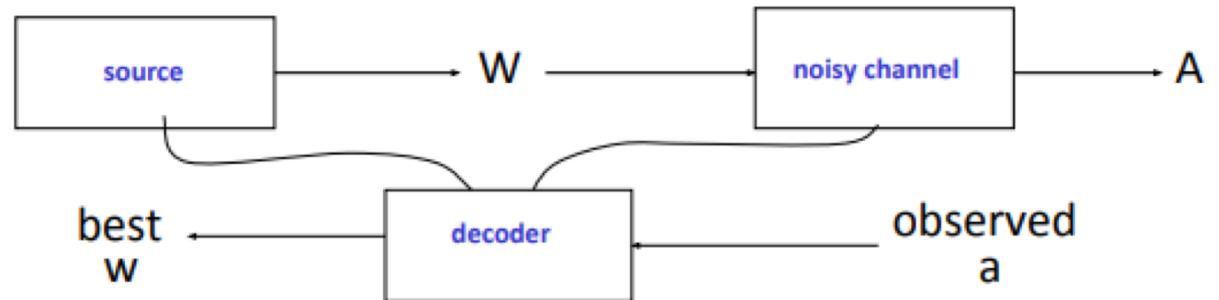
# The Noisy-Channel Model



We want to predict a sentence given acoustics:

$$w^* = \arg \max_w P(w|a)$$

# The Noisy-Channel Model



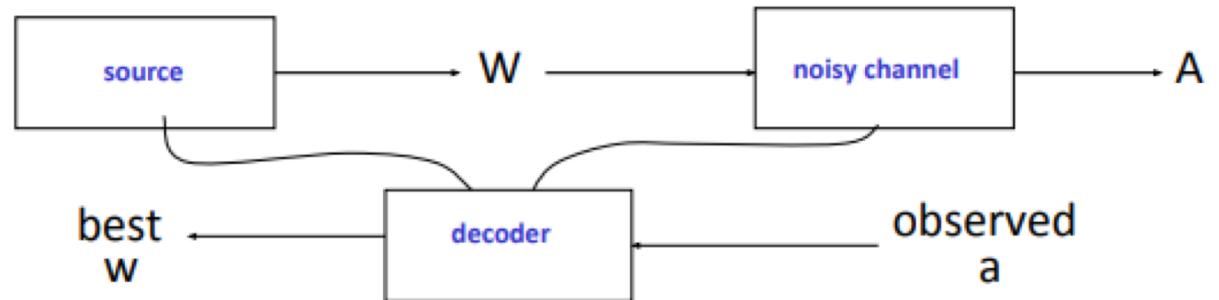
$$\begin{aligned}w^* &= \arg \max_w P(w|a) \\&= \arg \max_w \mathbf{P}(a|w) \mathbf{P}(w) / P(a)\end{aligned}$$

$$= \arg \max_w \mathbf{P}(a|w) \mathbf{P}(w)$$

Channel model

Source model

# The Noisy-Channel Model



$$\begin{aligned}w^* &= \arg \max_w P(w|a) \\&= \arg \max_w \mathbf{P}(a|w) \mathbf{P}(w) / P(a) \\&= \arg \max_w \mathbf{P}(a|w) \mathbf{P}(w)\end{aligned}$$

Likelihood

Acoustic model (HMMs)  
Translation model

Prior

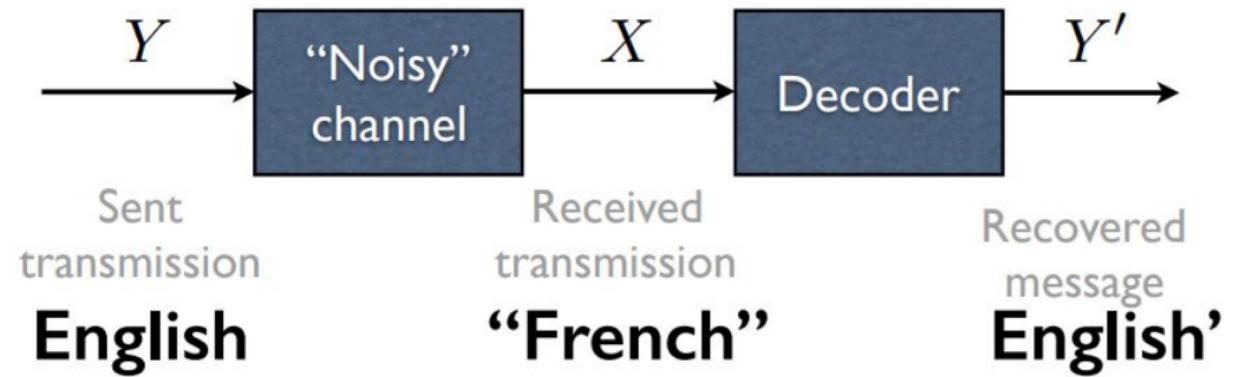
Language model: Distributions  
over sequence of words

# The Noisy-Channel Model

$$\hat{e} = \arg \max_e p_\varphi(e) \times p_\theta(f | e)$$

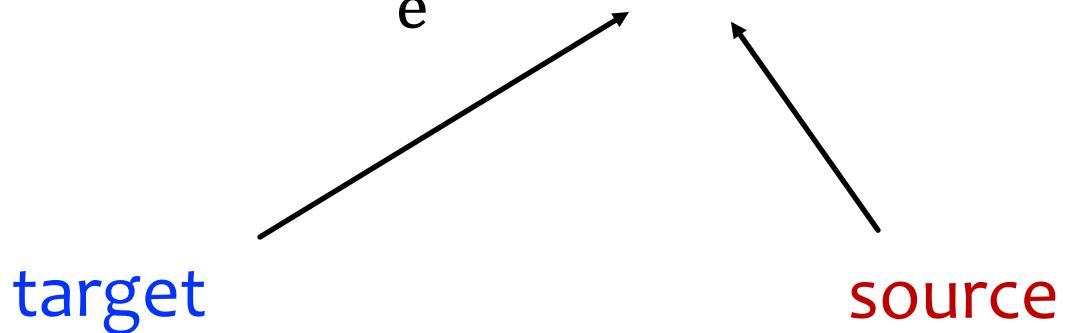
Language model

Translation model



# MT as Direct Modeling

$$\hat{e} = \arg \max_e p_\lambda(e | f)$$



- One model does everything
- Trained to reproduce a corpus of translations

# Two Views of MT

- **Code breaking** (aka the noisy channel, Bayes rule)
  - I know the **target language**
  - I have example **translations texts** (example enciphered data)
- **Direct modeling** (aka pattern matching)
  - I have **really good learning algorithms** and a bunch of **example inputs** (source language sentences) and **outputs** (target language translations)

# Which is Better?

- **Noisy channel** -  $p_\phi(e) \times p_\theta(f | e)$ 
  - Easy to use monolingual target language data
  - Search happens under a product of two models (individual models can be simple, product can be powerful)
- **Direct Model** -  $p_\lambda(e | f)$ 
  - Directly model the process you care about
  - Model must be very powerful

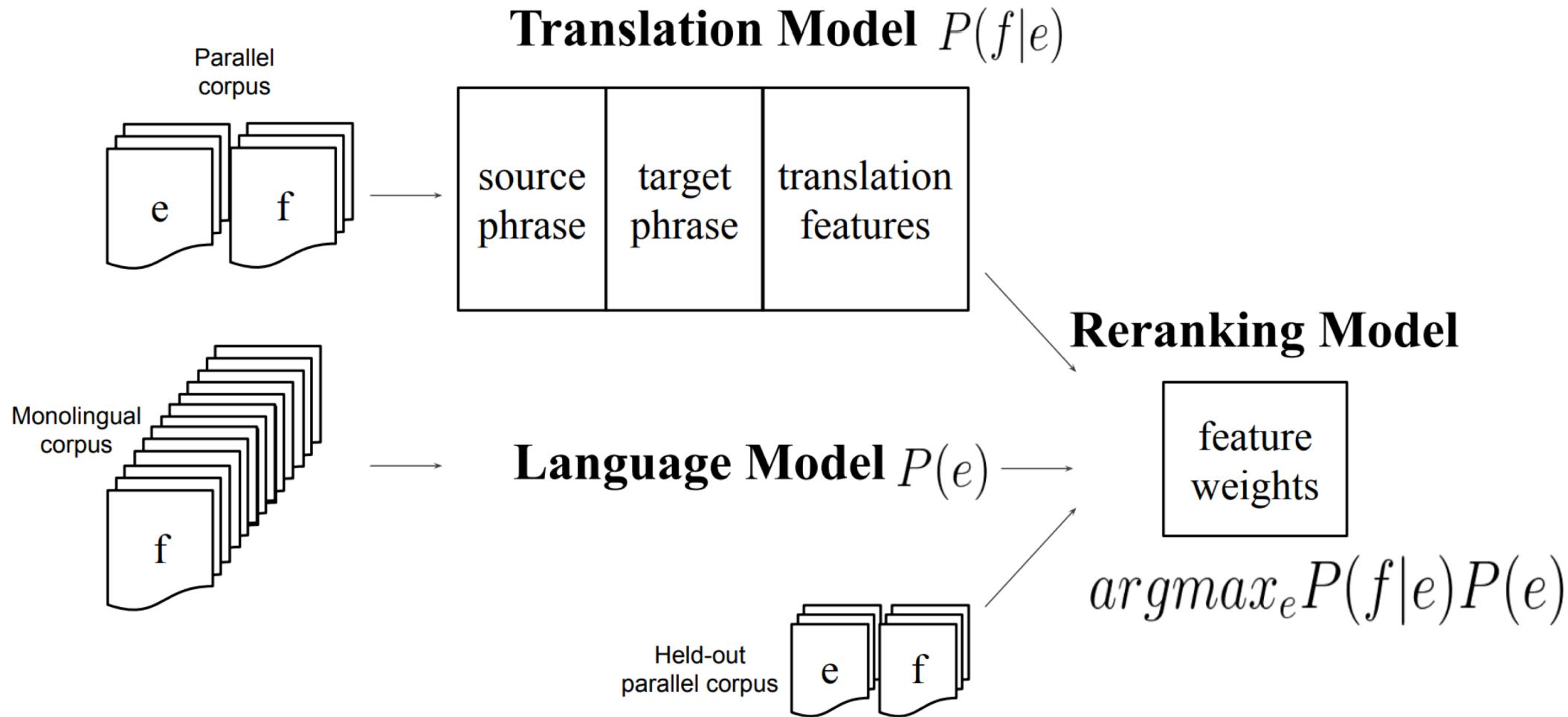
# Where are we in 2020?

- Direct modeling is where most of the action is
  - Neural networks are very good at generalizing and conceptually very simple
  - Inference in “product of two models” is hard
- Noisy channel ideas are incredibly important and still play a big role in how we think about translation

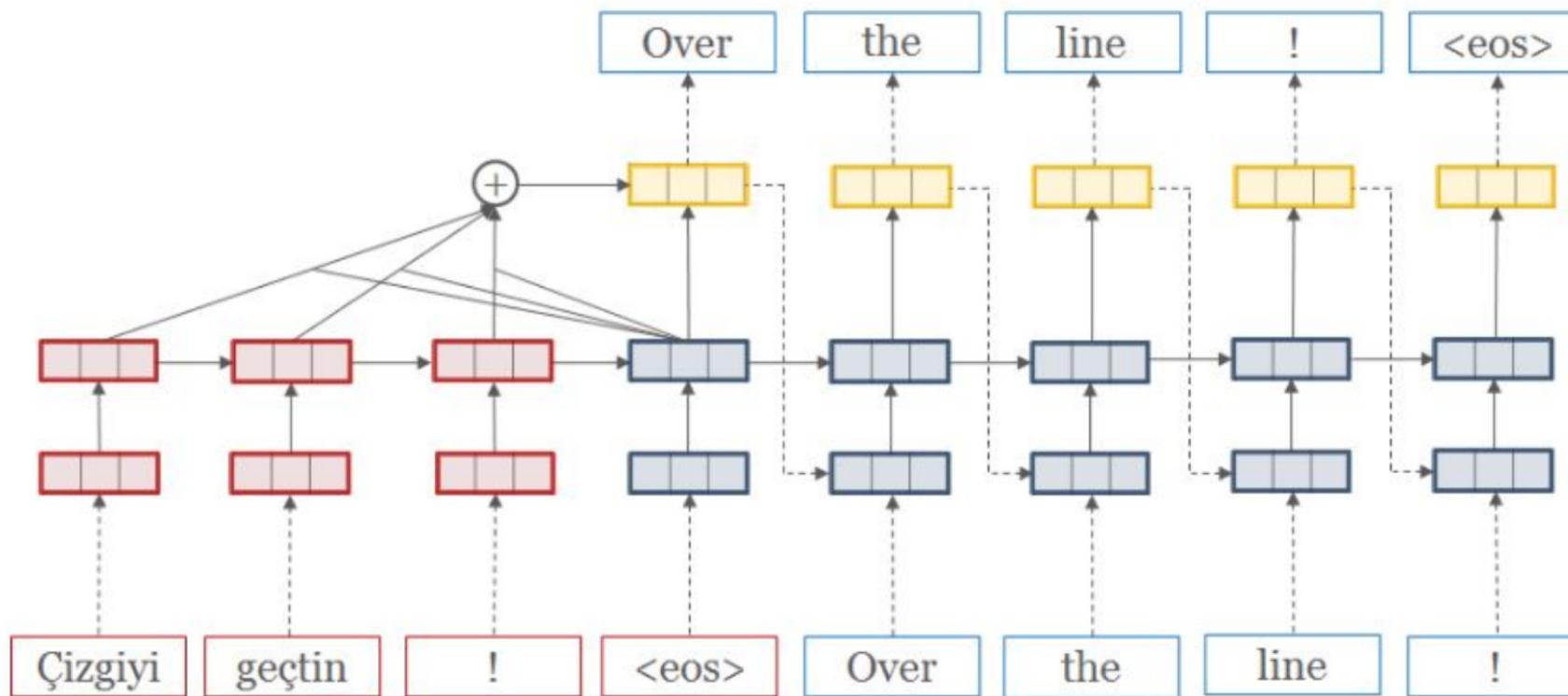
## Two Views of MT

- Noisy channel  $\hat{e} = \arg \max_e p_\varphi(e) \times p_\theta(f | e)$
- Direct  $\hat{e} = \arg \max_e p_\lambda(e | f)$

# Noisy Channel: Phrase-Based MT



# Neural MT: Conditional Language Modeling



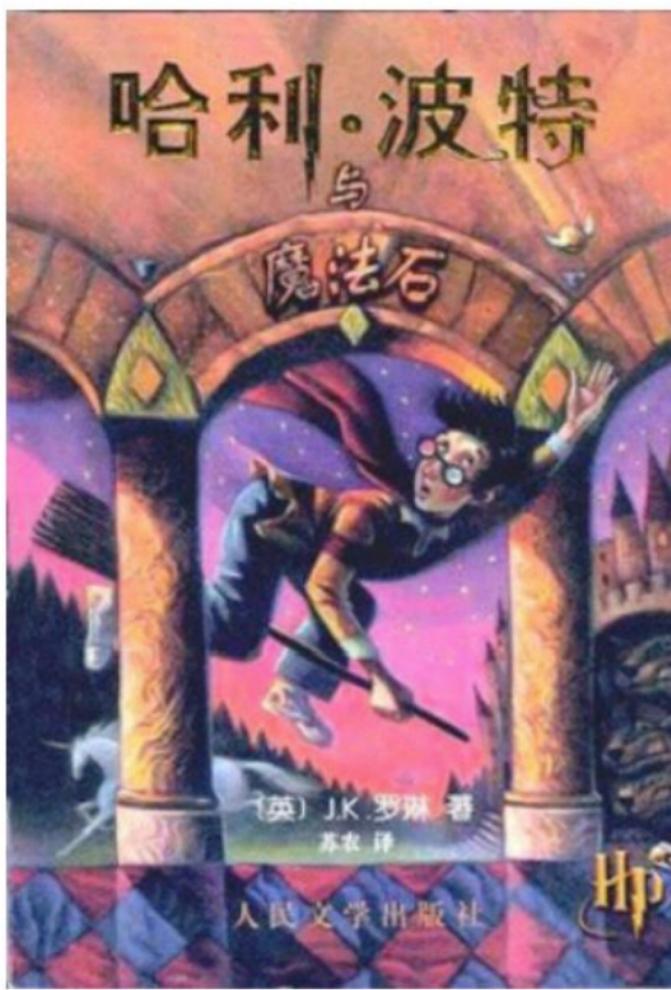
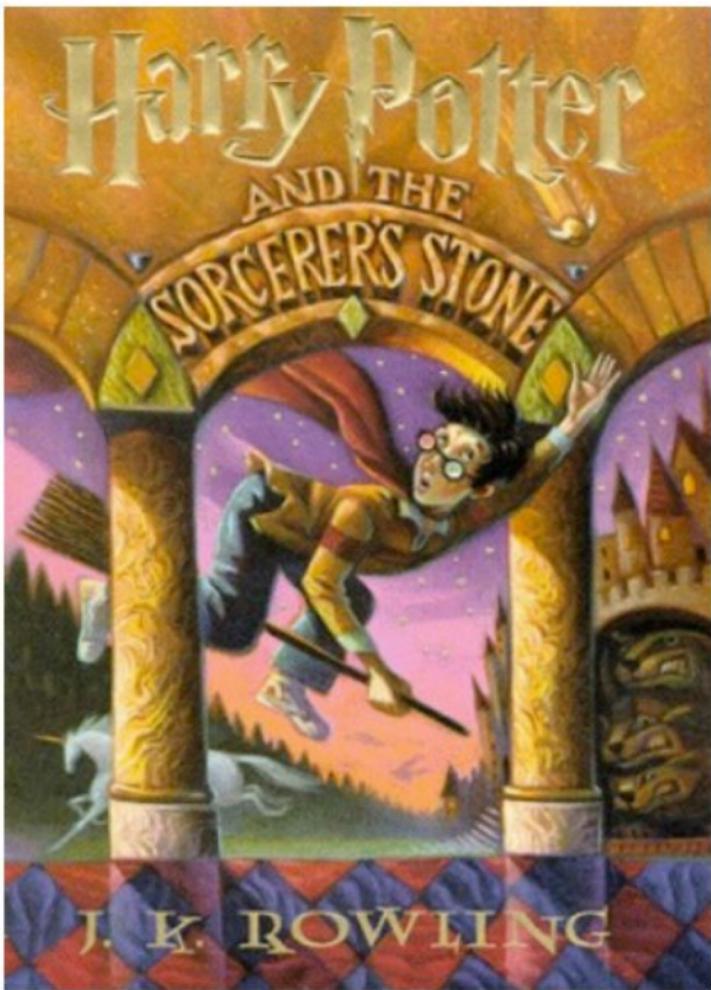
# A Common Problem

- Noisy channel  $\hat{e} = \arg \max_e p_\varphi(e) \times \textcolor{red}{p_\theta(f | e)}$
- Direct  $\hat{e} = \arg \max_e \textcolor{red}{p_\lambda(e | f)}$
- Both models must assign probabilities to how a sentence in one language translates into a sentence in another language

# Learning From Data

---

# Parallel Corpora



# Parallel Corpora

## CLASSIC SOUPS

Sm. Lg.

清 燉 雞 湯	57.	House Chicken Soup (Chicken, Celery, Potato, Onion, Carrot) .....	1.50	2.75
雞 飯 湯	58.	Chicken Rice Soup .....	1.85	3.25
雞 麵 湯	59.	Chicken Noodle Soup .....	1.85	3.25
廣 東 雲 吞	60.	Cantonese Wonton Soup.....	1.50	2.75
蕃 茄 蛋 湯	61.	Tomato Clear Egg Drop Soup .....	1.65	2.95
雲 吞 湯	62.	Regular Wonton Soup .....	1.10	2.10
酸 辣 湯	63.	Hot & Sour Soup .....	1.10	2.10
蛋 花 湯	64.	Egg Drop Soup.....	1.10	2.10
雲 蛋 湯	65.	Egg Drop Wonton Mix.....	1.10	2.10
豆 腐 菜 湯	66.	Tofu Vegetable Soup .....	NA	3.50
雞 玉 米 湯	67.	Chicken Corn Cream Soup .....	NA	3.50
蟹 肉 玉 米 湯	68.	Crab Meat Corn Cream Soup.....	NA	3.50
海 鮮 湯	69.	Seafood Soup.....	NA	3.50

# Parallel Corpora (mining parallel data from microblogs Ling et al., 2013)

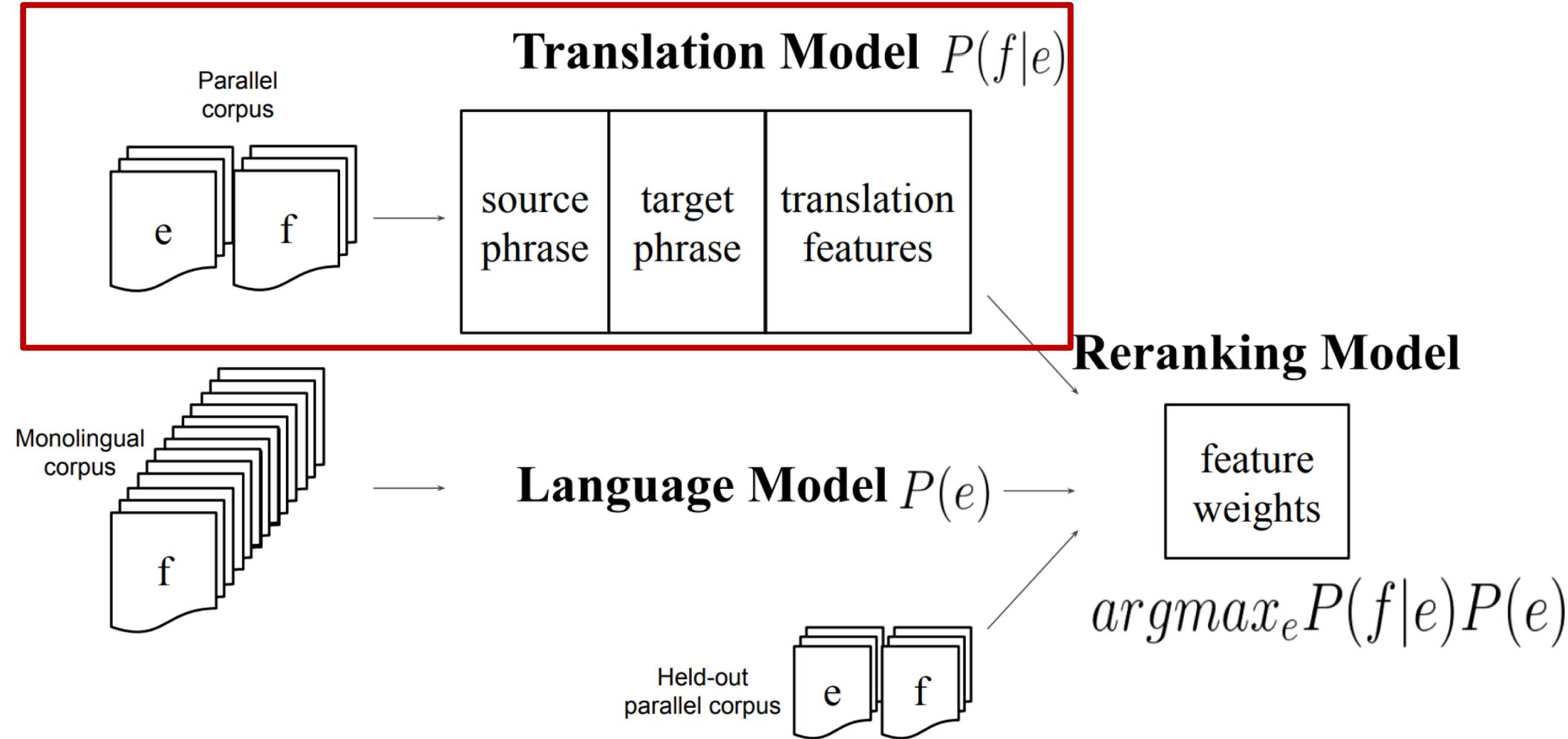
	ENGLISH	MANDARIN
1	i <b>wanna</b> live in a wes anderson world	我想要生活在Wes Anderson的世界里
2	Chicken soup, corn never truly digests. <b>TMI</b> .	鸡汤吧，玉米神马的从来没有真正消化过.恶心
3	To DanielVeuleman <b>yea iknw imma</b> work on that	对DanielVeuleman说，是的我知道，我正在向那方面努力
4	<b>msg 4</b> Warren G his <b>eday</b> is today 1 <b>yr</b> 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 <b>اليابي</b> حر
7	onni this gift only <b>4 u</b>	أوني هذه 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 لظاهرات في عدة مناطق غدا

Table 2: Examples of English-Mandarin and English-Arabic sentence pairs. The English-Mandarin sentences were extracted from Sina Weibo and the English-Arabic sentences were extracted from Twitter. Some messages have been shorted to fit into the table. Some interesting aspects of these sentence pairs are marked in bold.

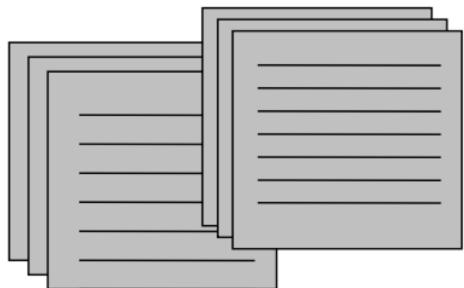
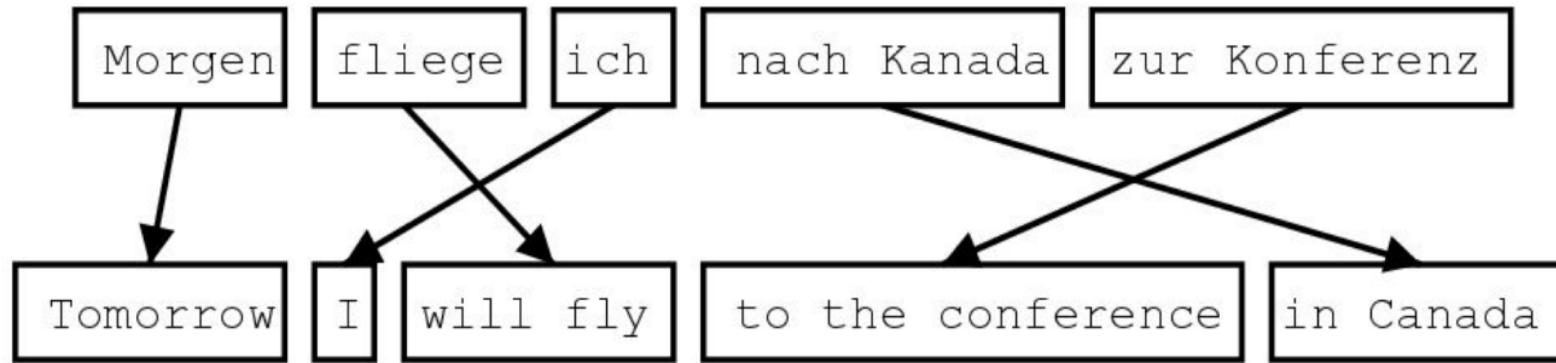
# Discussions

- There is a lot more monolingual data in the world than translated data
- Easy to get about 1 trillion words of English by crawling the web
- With some work, you can get 1 billion translated words of English-French
  - What about Japanese-Turkish?

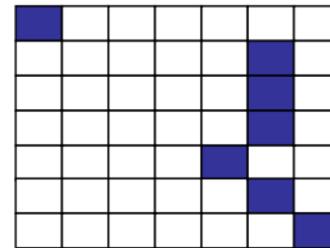
# Phrase-Based MT



# Construction of t-table



Sentence-aligned  
corpus



Word alignments



cat     chat     0.9
the cat     le chat     0.8
dog     chien     0.8
house     maison     0.6
my house     ma maison     0.9
language     langue     0.9
...

Phrase table  
(translation model)

# Word Alignment Models

---

# Lexical Translation

- How do we translate a word? Look it up in the dictionary  
*Haus – house, building, home, household, shell*
- Multiple translations
  - Some more frequent than others
  - Different word senses, different registers, different functions
  - *House, home* are common
- *Shell* is specialized (the Haus of a snail is a shell)

# How Common is Each?

- Look at a parallel corpus (German text along with English translation)

Translation of Haus	Count
house	8000
building	1600
home	200
household	150
shell	50

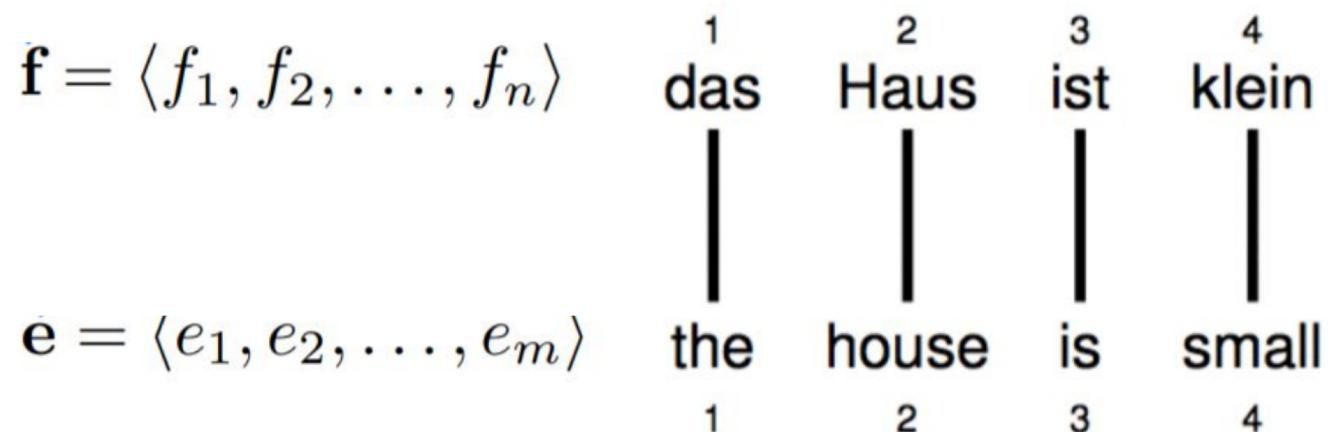
# Estimate Translation Probabilities

- Maximum likelihood estimation

$$\hat{p}_{\text{MLE}}(e \mid \text{Haus}) = \begin{cases} 0.8 & \text{if } e = \text{house}, \\ 0.16 & \text{if } e = \text{building}, \\ 0.02 & \text{if } e = \text{home}, \\ 0.015 & \text{if } e = \text{household}, \\ 0.005 & \text{if } e = \text{shell}. \end{cases}$$

# Word Alignment

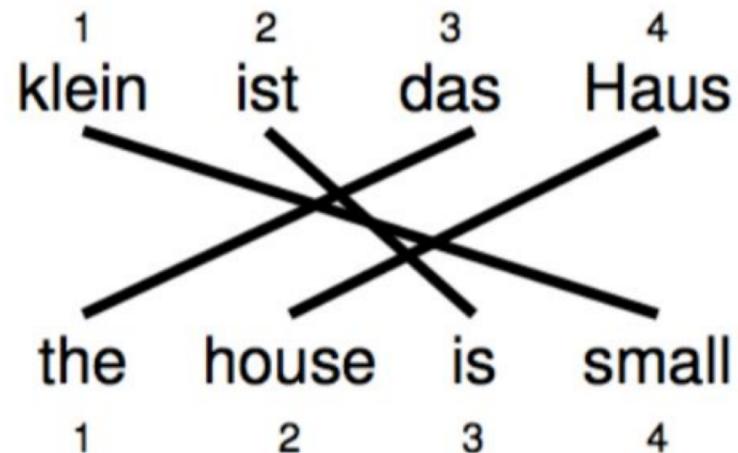
- Alignment can be visualized by drawing links between two sentences, and they are represented as vectors of positions



$$\mathbf{a} = (1, 2, 3, 4)^\top$$

# Reordering

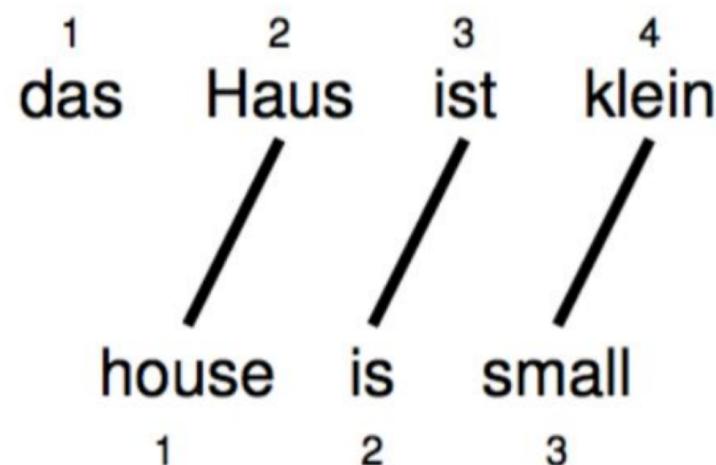
- Words may be reordered during translation



$$\mathbf{a} = (3, 4, 2, 1)^\top$$

# Word Dropping

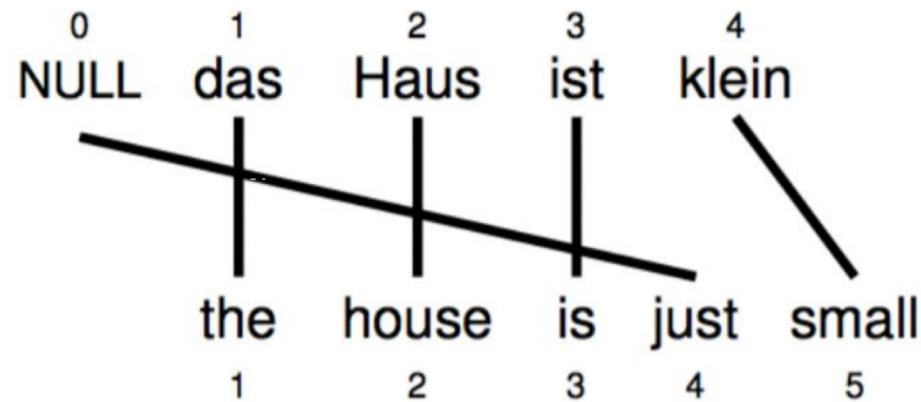
- A source word may not be translated at all



$$\mathbf{a} = (2, 3, 4)^\top$$

# Word Insertion

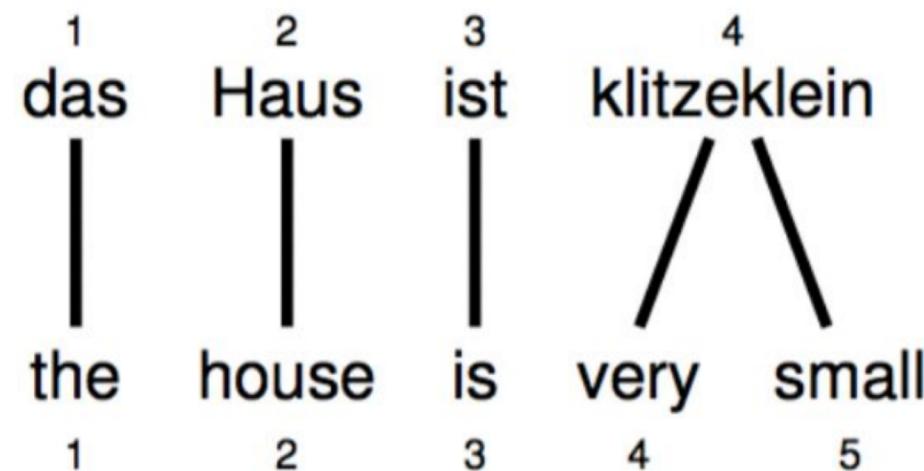
- Words may be inserted during translation
  - English **just** does not have an equivalent
  - But it must be explained – we typically assume every source sentence contains a **NULL** token



$$\mathbf{a} = (1, 2, 3, 0, 4)^\top$$

# One-to-many Translation

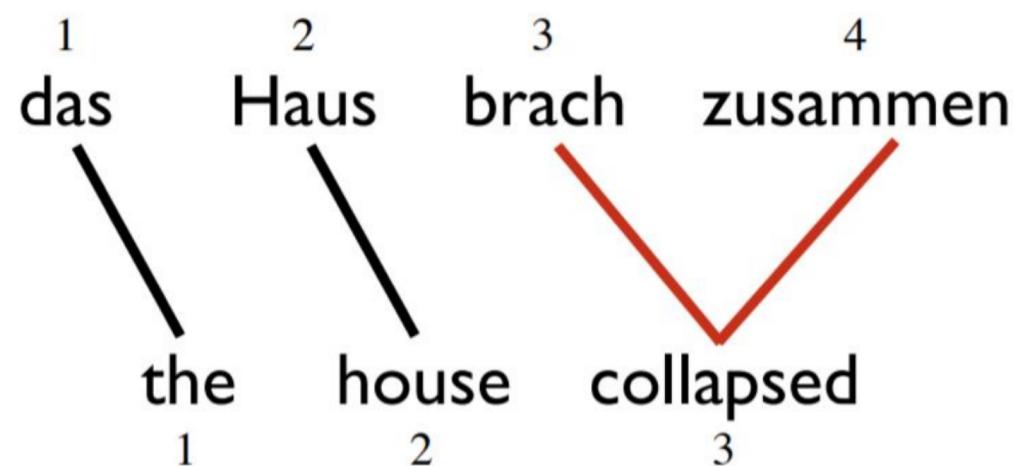
- A source word may translate into **more than one** target word



$$\mathbf{a} = (1, 2, 3, 4, 4)^\top$$

# Many-to-one Translation

- More than one source word may **not** translate as a unit in lexical translation



$$\mathbf{a} = ???$$

$$\mathbf{a} = (1, 2, (3, 4)^\top)^\top ?$$

# Computing Word Alignments

- Word alignments are the basis for most translation algorithms
- Given two sentences F and E, find a good alignment
- But a word-alignment algorithm can also be part of a mini-translation model itself
- One the most basic alignment models is also a simplistic translation model