

# Machine Translation

**CS 490A, Fall 2020**

Applications of Natural Language Processing

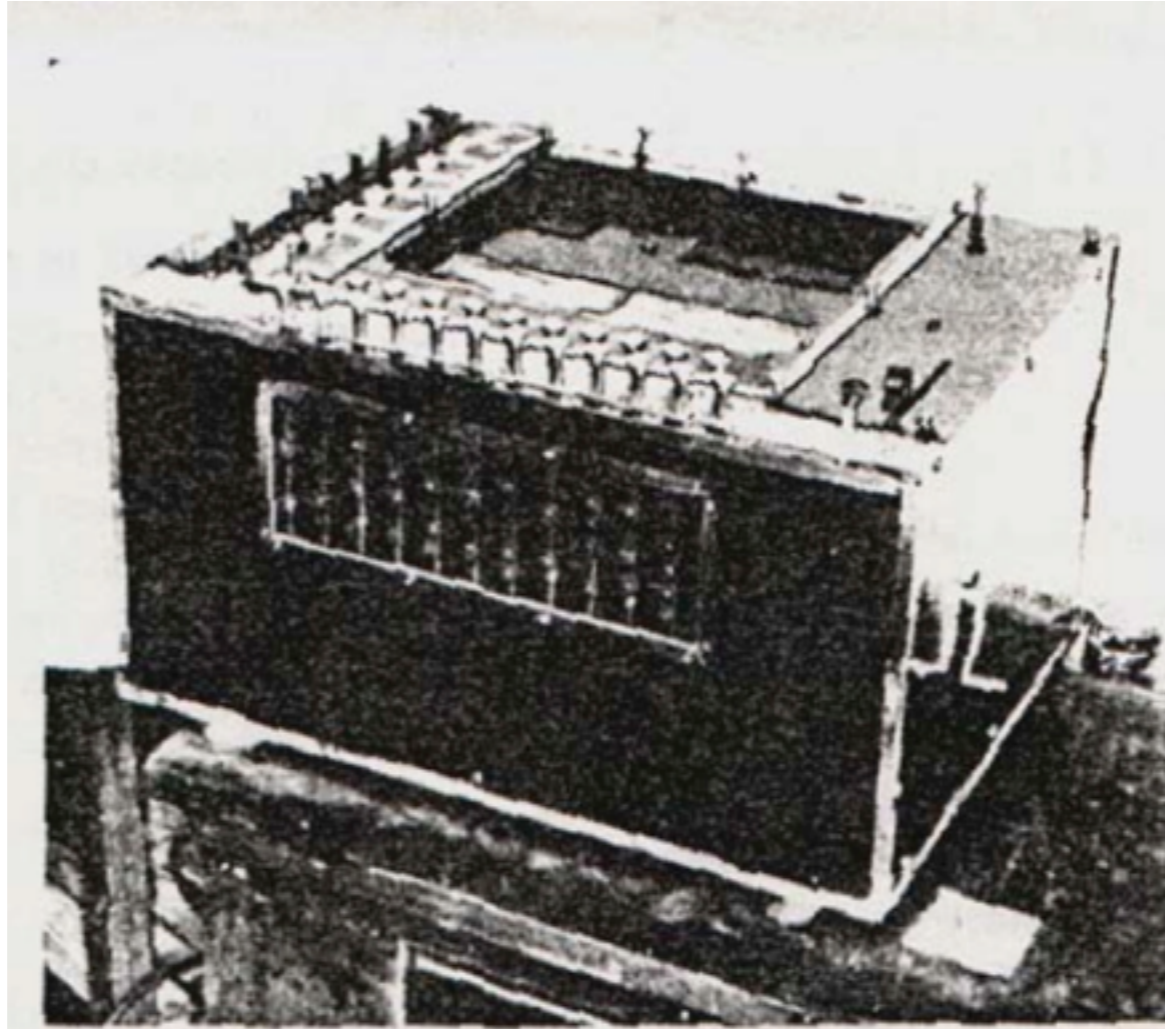
[https://people.cs.umass.edu/~brenocon/cs490a\\_f20/](https://people.cs.umass.edu/~brenocon/cs490a_f20/)

**Brendan O'Connor**

College of Information and Computer Sciences

University of Massachusetts Amherst

# MT is long-sought

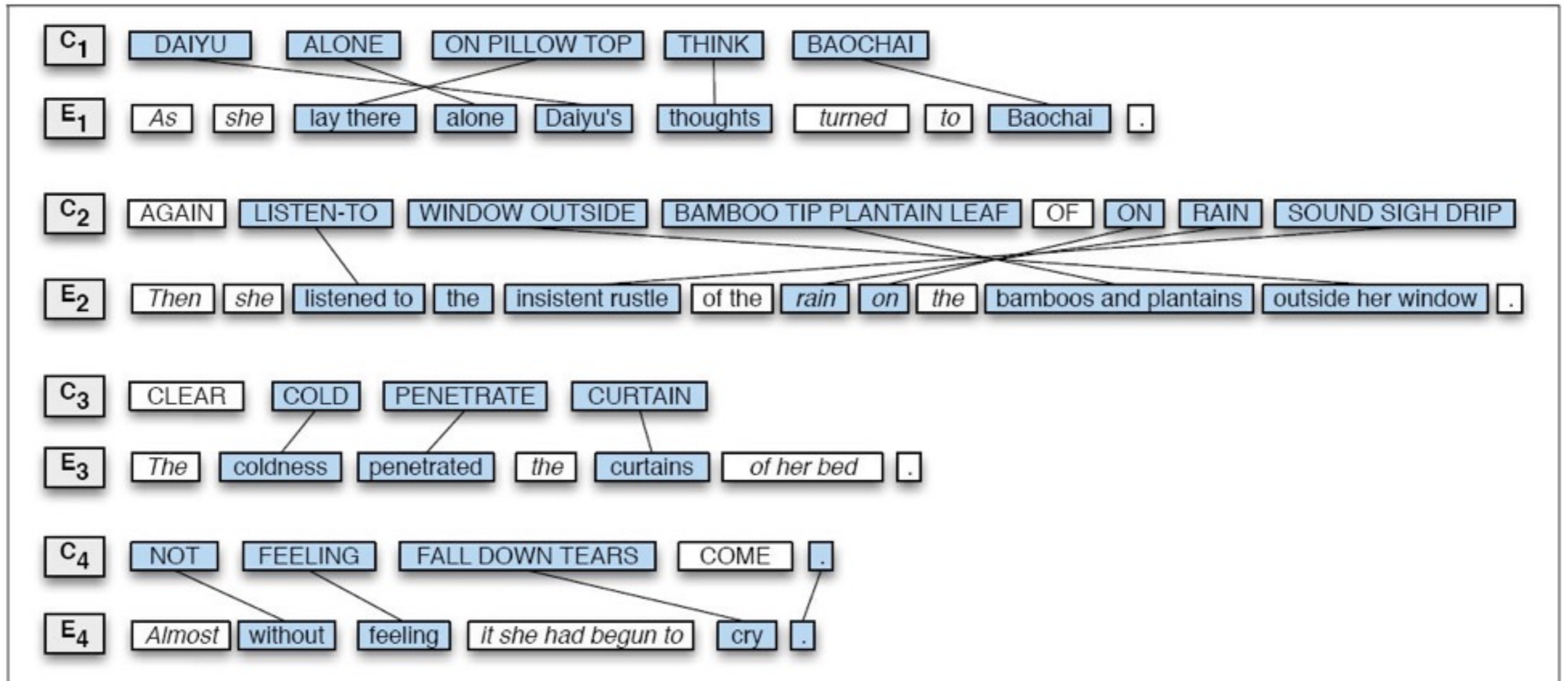


*Georges Artrouni's  
“mechanical brain”,  
a translation device  
patented in France  
in 1933. (Image from  
Corbé by way of  
[John Hutchins](#))*

The memory was the core of the device. It consisted of a paper band 40 cm wide, which could be up to 40 meters in length, moving over two rolling drums and held in position by perforations on the edges. The dictionary entries were recorded in normal orthographic form (i.e. not coded) line by line in five columns. The first column was for the source language word (or term), the other columns for equivalents in other languages and for other useful information.

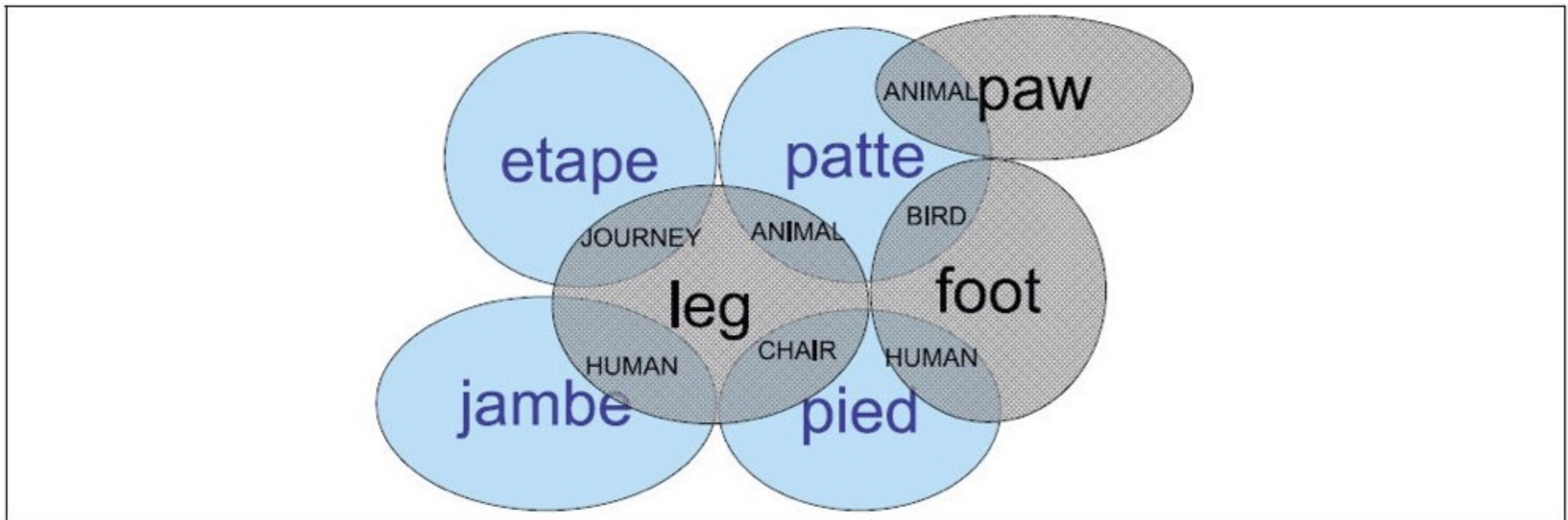
# MT is hard

- Word order, word meanings



# MT is hard

- Word meaning:  
many-to-many and context dependent



- *Translation* itself is hard: metaphors, cultural references, etc.

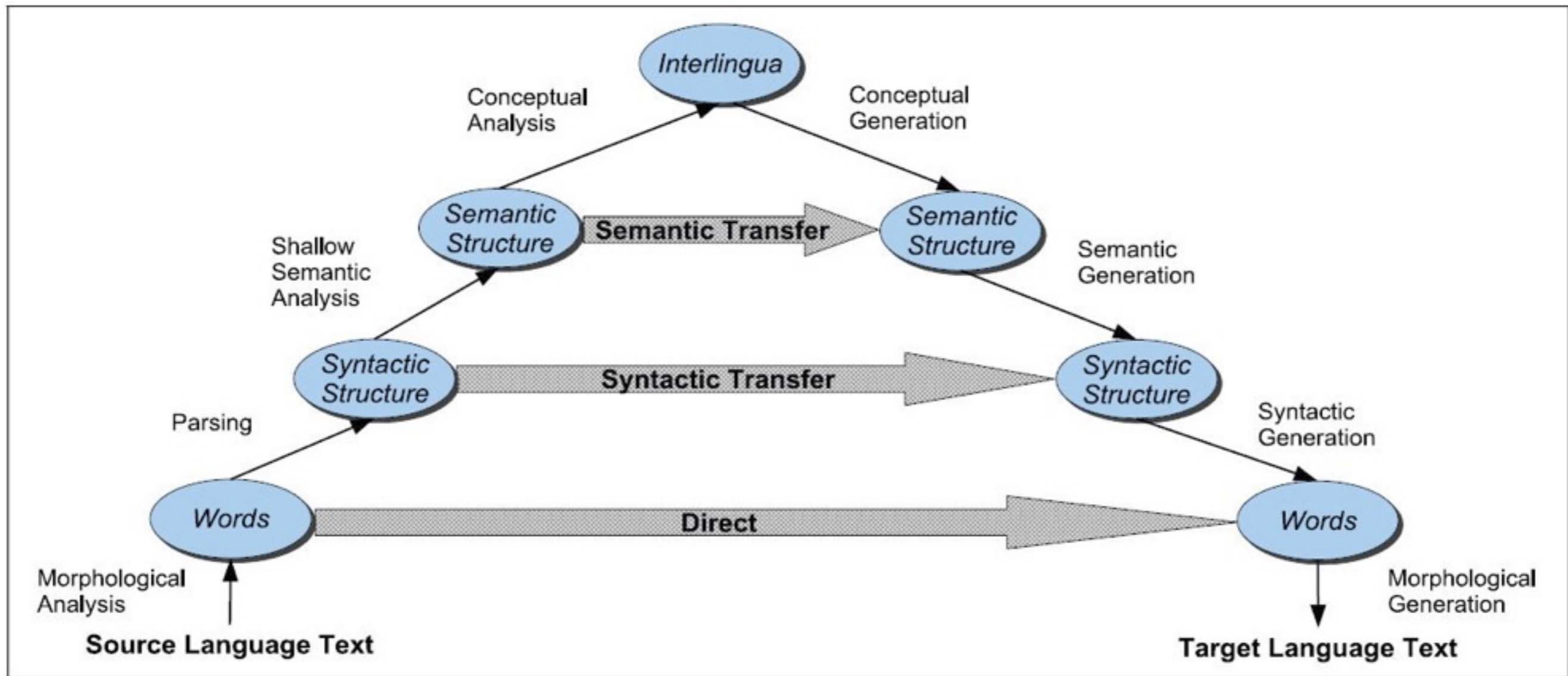
# MT goals

- Motivation: Human translation is expensive
- Rough translation vs. none
- Interactive assistance for human translators
  - e.g. Lilt
    - <https://www.youtube.com/watch?v=YZ7G3gQgpfl>
    - <https://lilt.com/app/projects/details/1887/edit-document/2306>
  - [compare to bilingual dictionary]

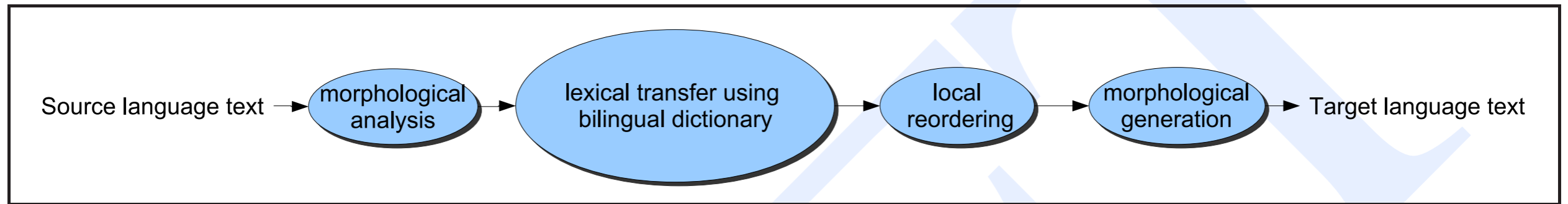
# MT paradigms

- Rule-based *transfer rules*
  - Manually program lexicons/rules
  - SYSTRAN (AltaVista BabelFish; originally from 70s)
- Statistical MT
  - Word-to-word, phrase-to-phrase probs
  - Learn translation rules from data, search for high-scoring translation outputs
    - Phrase or syntactic transformations
  - Key research in the early 90s
  - Google Translate (mid 00s)
  - Open-source: Moses
- Neural MT
  - Research in early 10s; very recently deployed
  - Latent representations of words/phrases

# Vauquois Triangle



# Direct (word-based) transfer



Input:

After 1: Morphology

After 2: Lexical Transfer

After 3: Local reordering

After 4: Morphology

Mary didn't slap the green witch

Mary DO-PAST not slap the green witch

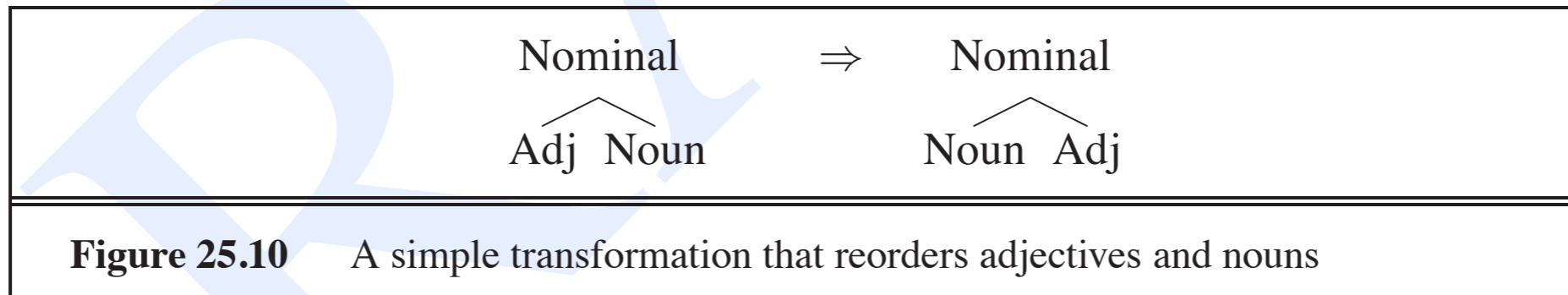
Maria PAST no dar una bofetada a la verde bruja

Maria no dar PAST una bofetada a la bruja verde

Maria no dió una bofetada a la bruja verde



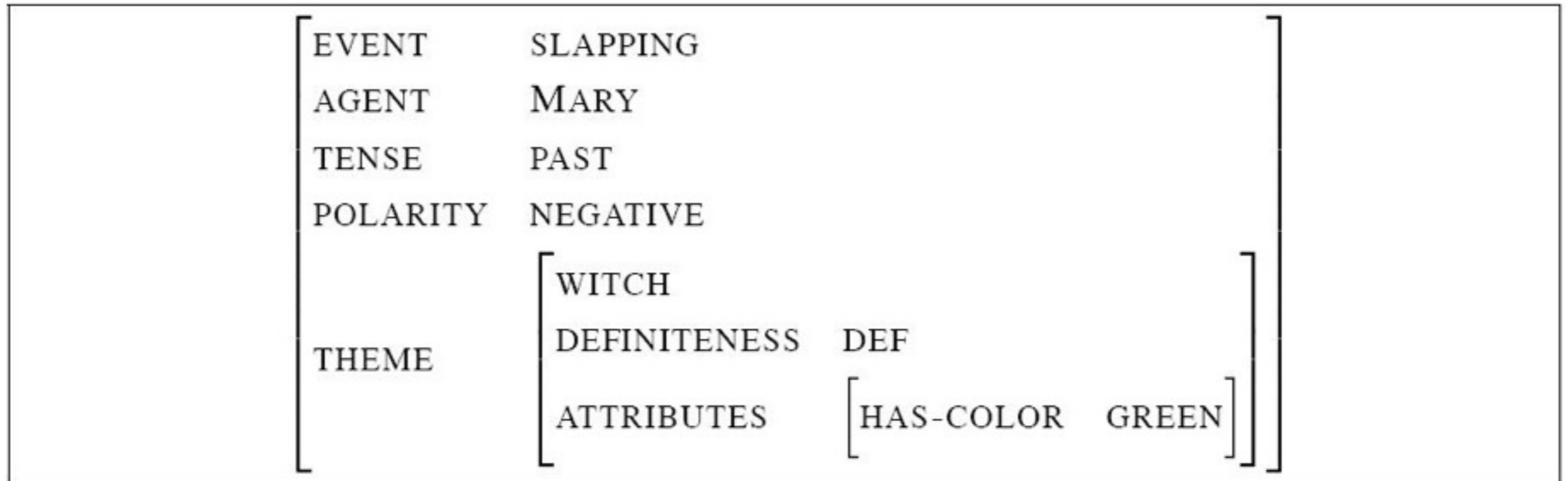
# Syntactic transfer



<b>English to Spanish:</b>			
1.	$\text{NP} \rightarrow \text{Adjective}_1 \text{ Noun}_2$	$\Rightarrow$	$\text{NP} \rightarrow \text{Noun}_2 \text{ Adjective}_1$
<b>Chinese to English:</b>			
2.	$\text{VP} \rightarrow \text{PP}[+\text{Goal}] \text{ V}$	$\Rightarrow$	$\text{VP} \rightarrow \text{V PP}[+\text{Goal}]$
<b>English to Japanese:</b>			
3.	$\text{VP} \rightarrow \text{V NP}$	$\Rightarrow$	$\text{VP} \rightarrow \text{NP V}$
4.	$\text{PP} \rightarrow \text{P NP}$	$\Rightarrow$	$\text{PP} \rightarrow \text{NP P}$
5.	$\text{NP} \rightarrow \text{NP}_1 \text{ Rel. Clause}_2$	$\Rightarrow$	$\text{NP} \rightarrow \text{Rel. Clause}_2 \text{ NP}_1$

# Interlingua

*“Mary did not slap the green witch”*



- More like classic logic-based AI
- Works in narrow domains
- Broad domain currently fails
  - Coverage: Knowledge representation for all possible semantics?
  - Can you parse to it?
  - Can you generate from it?

# Rules are hard

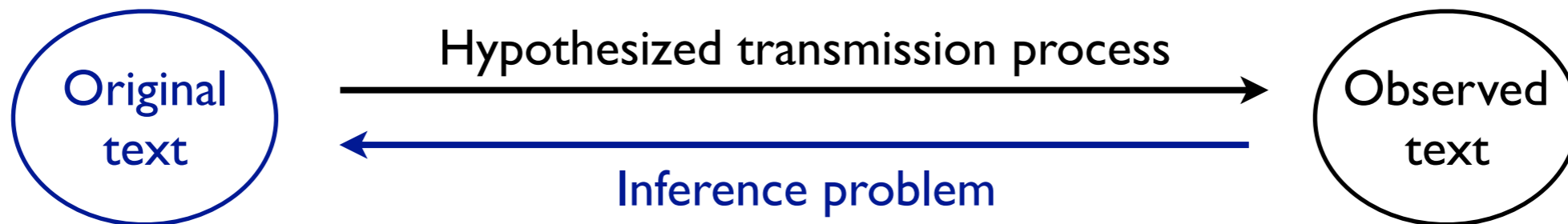
- Coverage
- Complexity (context dependence)
- Maintenance

```
function DIRECT_TRANSLATE_MUCH/MANY(word) returns Russian translation  
  
if preceding word is how return skol'ko  
else if preceding word is as return stol'ko zhe  
else if word is much  
    if preceding word is very return nil  
    else if following word is a noun return mnogo  
else /* word is many */  
    if preceding word is a preposition and following word is a noun return mnogii  
    else return mnogo
```

# Machine learning for MT

- MT as ML: Translation is something people do naturally. Learn rules from data?
- Parallel data: (source, target) text pairs
  - E.g. 20 million words of European Parliament proceedings  
<http://www.statmt.org/europarl/>
- Training: learn parameters to predict {source => target}
  - Sequence-to-sequence problem
- Test time: given source sentence, search for high-scoring target (e.g. beam search)

# Noisy channel model



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 (1955)

## Machine translation

$$P(\text{target text} \mid \text{source text}) \propto P(\text{source text} \mid \text{target text}) P(\text{target text})$$

# Statistical MT

- Pioneered at IBM, early 1990s  
(Forerunner of 90s-era statistical revolution in NLP)

## The COLING Paper Review

The validity of statistical (information theoretic) approach to MT has indeed been recognized, as the authors mention, by Weaver as early as 1949. And was universally recognized as mistaken by 1950. (cf. Hutchins, MT: Past, Present, Future, Ellis Horwood, 1986, pp. 30ff. and references therein) The crude force of computers is not science. The paper is simply beyond the scope of COLING.

Historical notes: <http://cs.jhu.edu/~post/bitext/>

# Statistical MT

- Pioneered at IBM, early 1990s  
(Forerunner of 90s-era statistical revolution in NLP)
- Noisy channel model borrowed from  
speech recognition processing

"Every time I fire a linguist,  
the performance of the speech recognizer goes up"  
[Fred Jelinek]

# IBM Models

- [Brown et al. 1993, “The Mathematics of Statistical Machine Translation: Parameter Estimation”]
- *Lexical translations*: each source word has word-level translations to target language
- *Alignments*: hypothesizes that individual input words get translated to outputs (potentially in different order)
- Training
  - Problem: don't know which came from which!
  - Solution: use the *EM algorithm*



# Example: learning with parallel data

Sentence pair 3 is much more challenging.  
So far, we have

erok	srok	izok	hihok	ghirok
totat	dat	arrat	vat	hilat

The Centauri word *izok* would be translated as either *totat*, *arrat*, or *vat*, yet when you look at *izok* in sentence pair 6, none of those three words appear in the Arcturan. Therefore, *izok* appears to be ambiguous. The word *hihok*, however, is fixed in sentence pair 11 as *arrat*. Both sentence pairs 3 and 12 have *izok hihok* sitting directly on top of *arrat vat*; so, in all possibility, *vat* seems a reasonable translation for (ambiguous) *izok*. Sentence pairs 5, 6, and 9 suggest that *quat* is its other translation. Through process of elimination, you connect the words *erok* and *totat*, finishing off the analysis:

erok	srok	izok	hihok	ghirok
totat	dat	arrat	vat	hilat

AI Magazine Volume 18 Number 4 (1997) (© AAI)

## Automating Knowledge Acquisition for Machine Translation

Kevin Knight

Notice that aligning the sentence pairs helps you to build the translation dictionary and that building the translation dictionary also helps you decide on correct alignments. You might call this the *decipherment method*.

# Example: learning with parallel data

1a. Garcia and associates.

1b. Garcia y asociados.

2a. Carlos Garcia has three associates.

2b. Carlos Garcia tiene tres asociados.

3a. his associates are not strong.

3b. sus asociados no son fuertes.

4a. Garcia has a company also.

4b. Garcia tambien tiene una empresa.

5a. its clients are angry.

5b. sus clientes están enfadados.

6a. the associates are also angry.

6b. los asociados tambien están enfadados.

7a. the clients and the associates are enemies.

7b. los clientes y los asociados son enemigos.

8a. the company has three groups.

8b. la empresa tiene tres grupos.

9a. its groups are in Europe.

9b. sus grupos están en Europa.

10a. the modern groups sell strong pharmaceuticals.

10b. los grupos modernos venden medicinas fuertes.

11a. the groups do not sell zanzanine.

11b. los grupos no venden zanzanina.

12a. the small groups are not modern.

12b. los grupos pequeños no son modernos.

# Lexical Translation

- How do we translate a word? Look it up in the dictionary

*Haus : house, home, shell, household*

- Multiple translations
  - Different word senses, different registers, different inflections (?)
  - *house, home* are common
  - *shell* is specialized (the Haus of a snail is a shell)

# How common is each?

Translation	Count
house	5000
home	2000
shell	100
household	80

# Lexical Translation

- Goal: a model  $p(\mathbf{e} \mid \mathbf{f}, m)$
- where  $\mathbf{e}$  and  $\mathbf{f}$  are complete English and Foreign sentences


$$\mathbf{e} = \langle e_1, e_2, \dots, e_m \rangle \quad \mathbf{f} = \langle f_1, f_2, \dots, f_n \rangle$$

# Lexical Translation

- Goal: a model  $p(\mathbf{e} \mid \mathbf{f}, m)$
- where  $\mathbf{e}$  and  $\mathbf{f}$  are complete English and Foreign sentences
- Lexical translation makes the following **assumptions**:
  - Each word in  $e_i$  in  $\mathbf{e}$  is generated from exactly one word in  $\mathbf{f}$
  - Thus, we have an *alignment*  $a_i$  that indicates which word  $e_i$  “came from”, specifically it came from  $f_{a_i}$ .
  - Given the alignments  $\mathbf{a}$ , translation decisions are conditionally independent of each other and depend *only* on the aligned source word  $f_{a_i}$ .

# Lexical Translation

$$\mathbf{e} = \langle e_1, e_2, \dots, e_m \rangle \quad \mathbf{f} = \langle f_1, f_2, \dots, f_n \rangle$$
$$\mathbf{a} = \langle a_1, a_2, \dots, a_m \rangle \quad \text{each } a_i \in \{0, 1, \dots, n\}$$

Modeling assumptions

$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in \{0, 1, \dots, n\}^m} p(\mathbf{a} \mid \mathbf{f}, m) \times \prod_{i=1}^m p(e_i \mid f_{a_i})$$

[Alignment] × [Translation | Alignment]

Goal: a model  $p(\mathbf{e} \mid \mathbf{f}, m)$

# Alignment

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

	1	2	3	4
<b>f</b> =	das	Haus	ist	klein
<b>e</b> =	the	house	is	small
	1	2	3	4

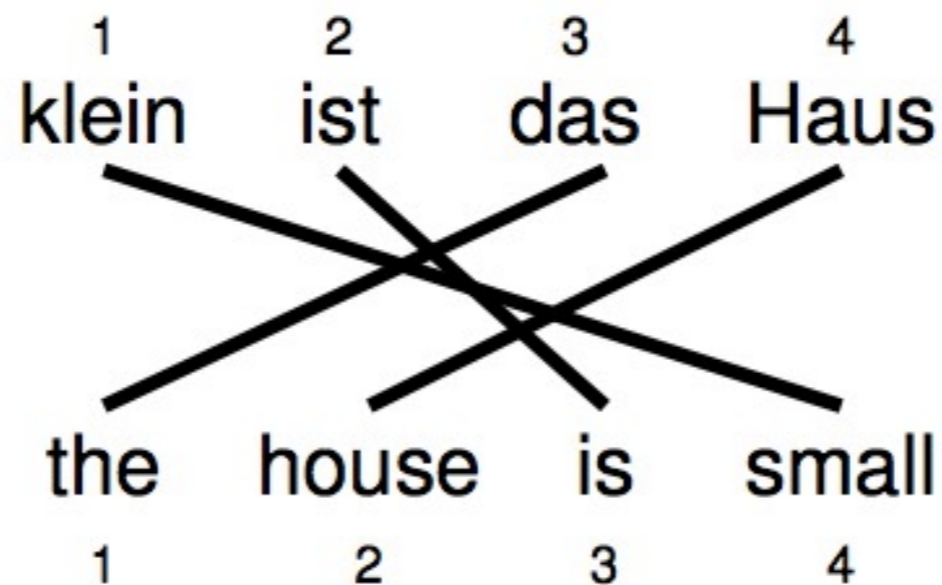
$$\mathbf{a} = (1, 2, 3, 4)^\top$$



Goal: a model  $p(\mathbf{e} \mid \mathbf{f}, m)$

# Reordering

- Words may be reordered during translation.

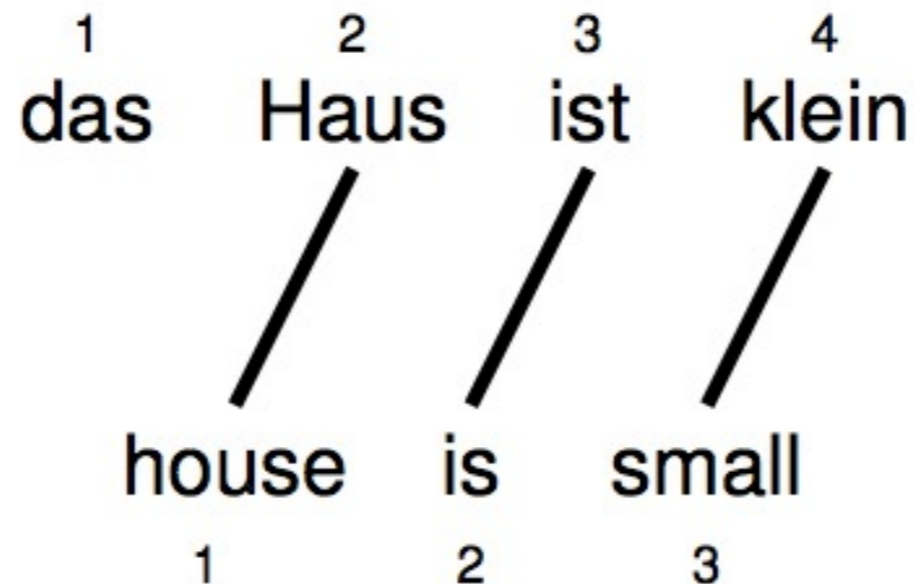


$$\mathbf{a} = (3, 4, 2, 1)^\top$$

Goal: a model  $p(\mathbf{e} \mid \mathbf{f}, m)$

# Word Dropping

- A source word may not be translated at all



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

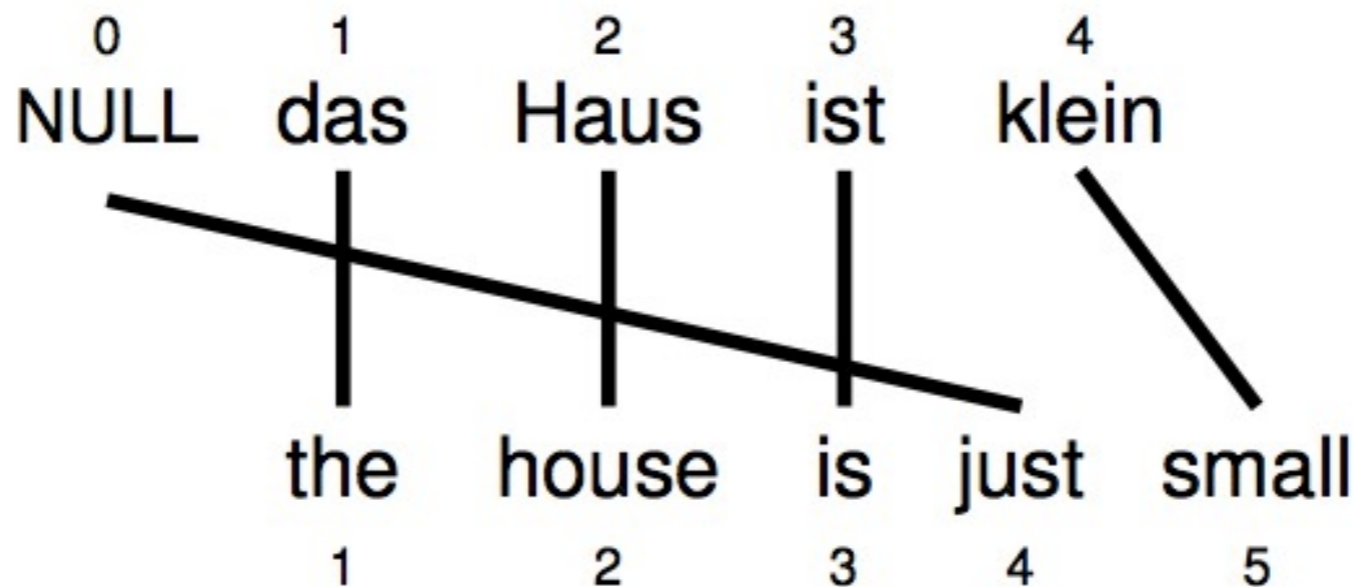
Goal: a model  $p(\mathbf{e} \mid \mathbf{f}, m)$

# 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

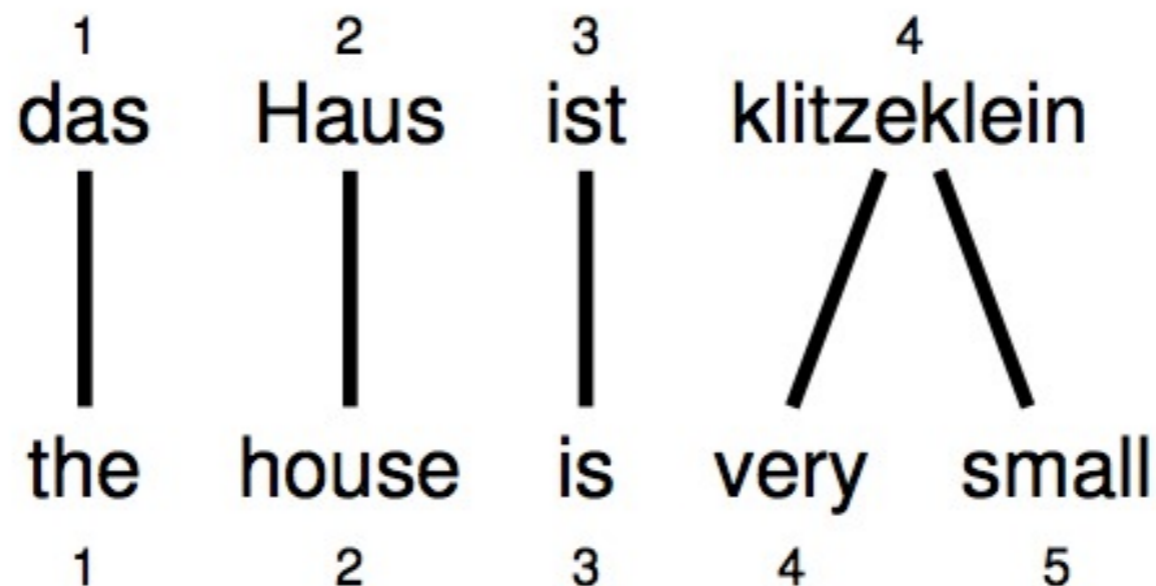


$$\mathbf{a} = (1, 2, 3, 0, 4)^\top$$

Goal: a model  $p(\mathbf{e} \mid \mathbf{f}, m)$

# One-to-many Translation

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

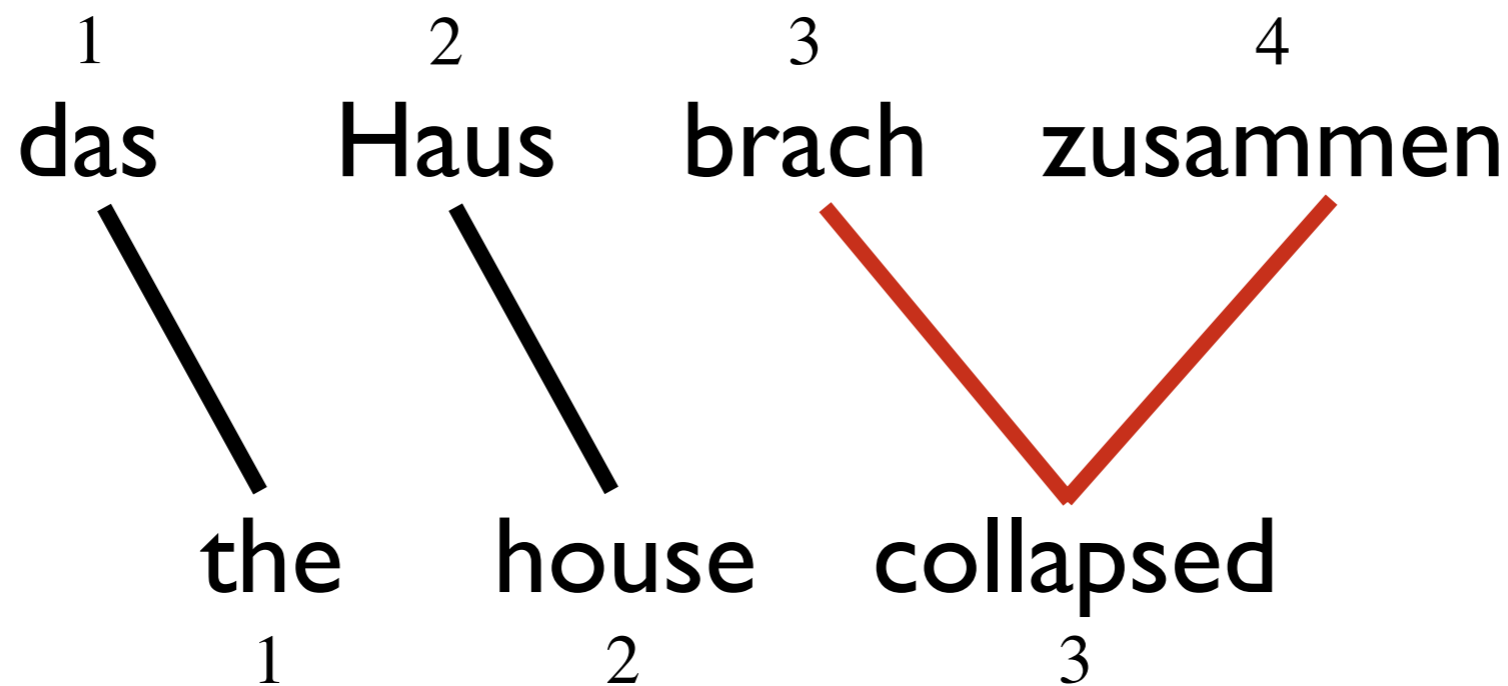


$$\mathbf{a} = (1, 2, 3, 4, 4)^\top$$

Goal: a model  $p(\mathbf{e} \mid \mathbf{f}, m)$

# Many-to-one Translation

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



$\mathbf{a} = ???$

[IBM Model 1 can't do this]

[slide: Chris Dyer]

# IBM Model I: Inference and learning

- Inferring alignments: assume lexical translations are independent conditional on alignments. That implies it's easy to compute

$$p(\mathbf{a} \mid \mathbf{e}, \mathbf{f}, \theta)$$

- How do we learn translation parameters?

$$\arg \max_{\theta} p(\mathbf{e} \mid \mathbf{f}, \theta)$$

- Chicken and egg problem:  
If we knew alignments, translation parameters would be trivial (just counting):

$$\arg \max_{\theta} p(\mathbf{e} \mid \mathbf{a}, \mathbf{f}, \theta)$$



# EM Algorithm

- pick some random (or uniform) parameters
- Repeat until you get bored (~ 5 iterations for lexical translation models)

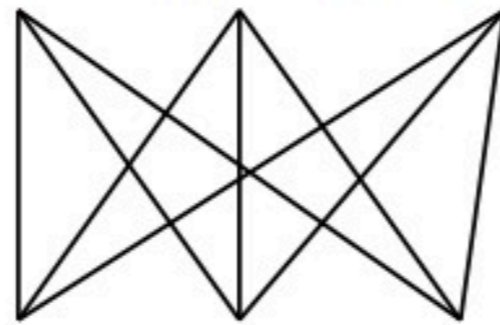
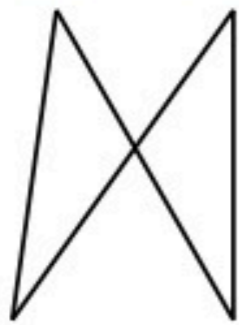
- using your current parameters, compute “expected” alignments for every target word token in the training data

$$p(a_i | \mathbf{e}, \mathbf{f}) \quad (\text{on board})$$

- keep track of the expected number of times  $f$  translates into  $e$  throughout the whole corpus
- keep track of the expected number of times that  $f$  is used as the source of any translation
- use these expected counts as if they were “real” counts in the standard MLE equation

# EM for Model 1

... la maison ... la maison blue ... la fleur ...



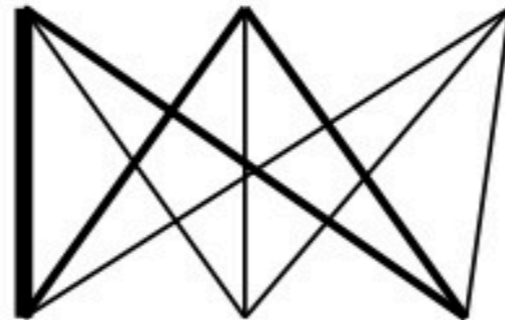
... the house ... the blue house ... the flower ...

- Initial step: all alignments equally likely
- Model learns that, e.g., *la* is often aligned with *the*



# EM for Model 1

... la maison ... la maison blue ... la fleur ...

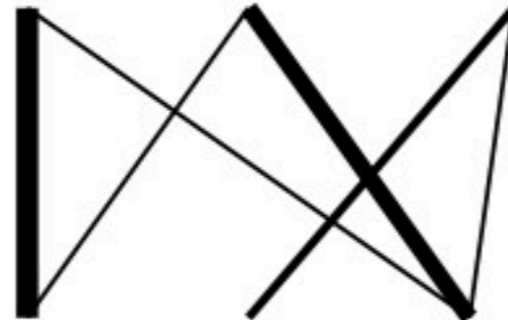


... the house ... the blue house ... the flower ...

- After one iteration
- Alignments, e.g., between *la* and *the* are more likely

# EM for Model 1

... la maison ... la maison bleu ... la fleur ...



... the house ... the blue house ... the flower ...

- After another iteration
- It becomes apparent that alignments, e.g., between **fleur** and **flower** are more likely (pigeon hole principle)

# EM for Model 1

... la maison ... la maison bleu ... la fleur ...



... the house ... the blue house ... the flower ...

- Convergence
- Inherent hidden structure revealed by EM

# EM for Model 1


... la maison ... la maison bleu ... la fleur ...  
/ | | X | |  
... the house ... the blue house ... the flower ...





$p(\text{la}|\text{the}) = 0.453$   
 $p(\text{le}|\text{the}) = 0.334$   
 $p(\text{maison}|\text{house}) = 0.876$   
 $p(\text{bleu}|\text{blue}) = 0.563$   
...

- Parameter estimation from the aligned corpus

# Convergence

das Haus  
  
 the house

das Buch  
  
 the book

ein Buch  
  
 a book

<i>e</i>	<i>f</i>	initial	1st it.	2nd it.	3rd it.	...	final
the	das	0.25	0.5	0.6364	0.7479	...	1
book	das	0.25	0.25	0.1818	0.1208	...	0
house	das	0.25	0.25	0.1818	0.1313	...	0
the	buch	0.25	0.25	0.1818	0.1208	...	0
book	buch	0.25	0.5	0.6364	0.7479	...	1
a	buch	0.25	0.25	0.1818	0.1313	...	0
book	ein	0.25	0.5	0.4286	0.3466	...	0
a	ein	0.25	0.5	0.5714	0.6534	...	1
the	haus	0.25	0.5	0.4286	0.3466	...	0
house	haus	0.25	0.5	0.5714	0.6534	...	1

# EM algorithm

- Very general meta-algorithm when we have
  - observed data  $x$
  - latent (hidden) variables  $z$  (alignments)
  - parameters  $\theta$
- and
  - we want  $\operatorname{argmax}_{\theta} P(x | \theta)$  but it's intractable
  - but this is easy:  $\operatorname{argmax}_{\theta} P(x, z | \theta)$
- EM: iterate
  - E-step: Infer  $P(z | x, \theta)$  [[make your best guess]]
  - M-step: Infer with the usual MLE but with *weighted* counts from the E-step
- Many applications in NLP:
  - Unsupervised training of HMMs, NB, topic models...
  - Semi-supervised learning: see only some of the labels

# MT Evaluation

- Problem: {source => target} has very large output space
- Ideally: bilingual humans judge every (source, translation) pair. Typically Likert scale for:
  - Faithfulness
  - Fluency [[monolingual humans can do this one]]
- *BLEU score*: an automatic metric
  - Given (source, target) gold-standard, score translation
  - ~Precision of translation's ngrams: are they in the gold-standard translation?
  - Brevity penalty (so can't game it with short sentences)
- Problem: there are multiple legitimate ways to say the same thing!
  - => Use multiple alternate gold-standard translations

# Multiple Reference Translations

## Reference translation 1:

The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport.

## Reference translation 2:

Guam International Airport and its offices are maintaining a high state of alert after receiving an e-mail that was from a person claiming to be the wealthy Saudi Arabian businessman Bin Laden and that threatened to launch a biological and chemical attack on the airport and other public places.

## Machine translation:

The American [?] international airport and its the office all receives one calls self the sand Arab rich business [?] and so on electronic mail, which sends out; The threat will be able after public place and so on the airport to start the biochemistry attack, [?] highly alerts after the maintenance.

## Reference translation 3:

The US International Airport of Guam and its office has received an email from a self-claimed Arabian millionaire named Laden, which threatens to launch a biochemical attack on such public places as airport. Guam authority has been on alert.

## Reference translation 4:

US Guam International Airport and its office received an email from Mr. Bin Laden and other rich businessman from Saudi Arabia. They said there would be biochemistry air raid to Guam Airport and other public places. Guam needs to be in high precaution about this matter.



# MT paradigms

- Statistical MT: phrase-based and fancier variants of Model 1
- Neural MT
  - Use NN representations of words and sentences
  - Condition on source, generate sentence in target language
  - Alignment models reborn as *attention-based neural networks*: choose which source words to look at, when translating next target word (previous lecture)

# Major challenge: low resource settings

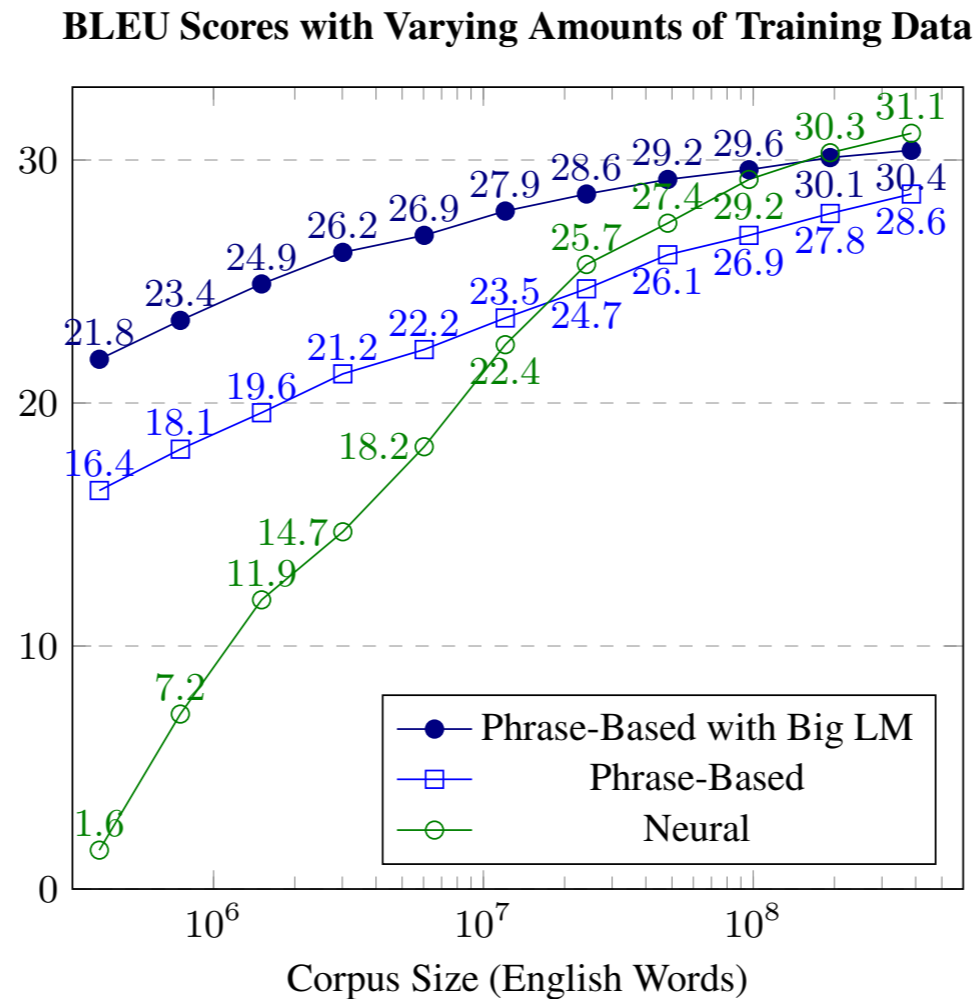


Figure 3: BLEU scores for English-Spanish systems trained on 0.4 million to 385.7 million words of parallel data. Quality for NMT starts much lower, outperforms SMT at about 15 million words, and even beats a SMT system with a big 2 billion word in-domain language model under high-resource conditions.

# Tuning in low data environments

	SMT	NMT				
		chosen	best	75%	50%	25%
Somali-English	15.1	14.4	14.4	12.7	11.7	9.9
Swahili-English	24.4	24.8	24.8	20.5	18.7	15.6

Table 2: **SMT vs. NMT: BLEU score on the Text Testset.** Models trained with 24k baseline dataset. For NMT, we trained approximately 600 systems with different hyperparameters. The “chosen” column shows the BLEU score on the test set based on a model chosen based on the validation set (which is a fair comparison to the SMT score), and the “best” column shows the best possible attainable score (in this case, chosen models happen to be the best models). We also show the 75, 50, 25 percentile of BLEU scores on the test set. **The wide range of scores for NMT indicates the sensitivity of NMT to design choices and the importance of careful tuning in low-resource scenarios.**

# Data availability

Language	Family	CommonCrawl (#documents)	Wikipedia (#documents)	OPUS (#sents)
Afrikaans (afr)	Indo-European	387k	84.0k	1.6m
Akan (aka)	Niger-Congo	3k	0.7k	0.2k
Amharic (amh)	Afroasiatic	66k	14.8k	1m
Arabic (ara)	Afroasiatic	17,772k	945.7k	70m
Berber (ber)	Afroasiatic	0	0	0.1m
Chewa (nya)	Niger-Congo	8k	0.5k	0.9m
Hausa (hau)	Afroasiatic	45k	3.7k	0.4m
Igbo (ibo)	Niger-Congo	8k	1.4k	0.5m
French (fra)	Indo-European	133,401k	2136.3k	180m
Fulani (ful)	Niger-Congo	0	0.2k	0.3k
Kinyarwanda (kin)	Niger-Congo	71k	1.8k	0.8m
Kirundi (run)	Niger-Congo	3k	0.6k	0
Malagasy (mlg)	Austronesian	126k	91.9k	0.9m
Mossi (mos)	Niger-Congo	0	0	0
Oromo (orm)	Afroasiatic	15k	0.8k	0.2m
Portuguese (por)	Indo-European	60,762k	1013.0k	72m
Shona (sna)	Niger-Congo	8k	4.8k	0.8m
Somali (som)	Afroasiatic	117k	5.4k	0.2m
Swahili (swa)	Niger-Congo	234k	53.7k	1.2m
Tigrinya (tir)	Afroasiatic	21k	0.2k	0.4m
Xhosa (xho)	Niger-Congo	12k	1.0k	1.5m
Yoruba (yor)	Niger-Congo	21k	31.9k	0.5m
Zulu (zul)	Niger-Congo	24k	1.3k	1.1m

Table 4: Potential digital resources for an abridged list of languages in Africa. We show the potential monolingual resources (Number of CommonCrawl and Wikipedia documents) and bilingual resources (Number of bilingual sentence pairs via OPUS). One can compare the low-resource condition of these languages, using Somali and Swahili as a reference point. Please refer to Section 5 for details, since these numbers need to be interpreted with care. Languages that are not on this list might have even fewer resources.