

CS447: Natural Language Processing

<http://courses.engr.illinois.edu/cs447>

Lecture 23:

Statistical

Machine Translation (II)

Julia Hockenmaier

juliahmr@illinois.edu

3324 Siebel Center

The IBM alignment models

The IBM models

Use the noisy channel (Bayes rule) to get the best (most likely) target translation e for source sentence f :

$$\arg \max_e P(e|f) = \arg \max_e P(f|e)P(e)$$

noisy channel

The translation model $P(f|e)$ requires alignments a

$$P(f|e) = \sum_{a \in \mathcal{A}(e,f)} P(f, a|e)$$

marginalize (=sum)
over all alignments a

Generate f and the alignment a with $P(f, a|e)$:

$$P(f, a|e) = \underbrace{P(m|e)}_{\text{Length: } |f|=m} \prod_{j=1}^m \underbrace{P(a_j | a_{1..j-1}, f_{1..j-1}, m, e)}_{\text{Word alignment } a_j} \underbrace{P(f_j | a_{1..j} f_{1..j-1}, e, m)}_{\text{Translation } f_j}$$

$m = \# \text{ words}$
in f_j

probability of
alignment a_j

probability
of word f_j

Model parameters

Length probability $P(m \mid n)$:

What's the probability of generating a source sentence of length m given a target sentence of length n ?

Count in training data

Alignment probability: $P(\mathbf{a} \mid m, n)$:

Model 1 assumes all alignments have the same probability:

For each position $a_1 \dots a_m$, pick one of the $n+1$ target positions uniformly at random

Translation probability: $P(f_j = lac \mid a_j = i, e_i = lake)$:

In Model 1, these are the only parameters we have to learn.

Representing word alignments

		1	2	3	4	5	6	7	8
		Marie	a	traversé	le	lac	à	la	nage
0	NULL								
1	Mary								
2	swam								
3	across								
4	the								
5	lake								



Position	1	2	3	4	5	6	7	8
Foreign	Marie	a	traversé	le	lac	à	la	nage
Alignment	1	3	3	4	5	0	0	2

Every source word $f[i]$ is aligned to one target word $e[j]$ (incl. NULL). We represent alignments as a vector \mathbf{a} (of the same length as the source) with $\mathbf{a}[i] = j$

IBM model 1: Generative process

For each target sentence $e = e_1..e_n$ of length n :

0	1	2	3	4	5
NULL	Mary	swam	across	the	lake

1. **Choose a length m** for the source sentence (e.g $m = 8$)

Position	1	2	3	4	5	6	7	8
----------	---	---	---	---	---	---	---	---

2. **Choose an alignment $a = a_1...a_m$** for the source sentence

Each a_j corresponds to a word e_i in e : $0 \leq a_j \leq n$

Position	1	2	3	4	5	6	7	8
Alignment	1	3	3	4	5	0	0	2

3. **Translate each target word e_{a_j}** into the source language

Position	1	2	3	4	5	6	7	8
Alignment	1	3	3	4	5	0	0	2
Translation	Marie	a	traversé	le	lac	à	la	nage

IBM model 1: details

The **length probability** is constant: $P(m | e) = \epsilon$

The **alignment probability** is uniform

(n = length of target string): $P(a_i | e) = 1/(n+1)$

The **translation probability** depends only on e_{a_i}

(the corresponding target word): $P(f_i | e_{a_i})$

$$\begin{aligned} P(\mathbf{f}, \mathbf{a} | \mathbf{e}) &= \underbrace{P(m | \mathbf{e})}_{\text{Length: } |\mathbf{f}|=m} \prod_{j=1}^m \underbrace{P(a_j | a_{1..j-1}, f_{1..j-1}, m, \mathbf{e})}_{\text{Word alignment } a_j} \underbrace{P(f_j | a_{1..j} f_{1..j-1}, \mathbf{e}, m)}_{\text{Translation } f_j} \\ &= \epsilon \prod_{j=1}^m \frac{1}{n+1} P(f_j | e_{a_j}) \\ &= \frac{\epsilon}{(n+1)^m} \prod_{j=1}^m P(f_j | e_{a_j}) \end{aligned}$$

Finding the best alignment

How do we find the **best alignment** between **e** and **f**?

$$\begin{aligned}\hat{\mathbf{a}} &= \arg \max_{\mathbf{a}} P(\mathbf{f}, \mathbf{a} | \mathbf{e}) \\ &= \arg \max_{\mathbf{a}} \frac{\epsilon}{(n+1)^m} \prod_{j=1}^m P(f_j | e_{a_j}) \\ &= \arg \max_{\mathbf{a}} \prod_{j=1}^m P(f_j | e_{a_j})\end{aligned}$$

$$\hat{a}_j = \arg \max_{a_j} P(f_j | e_{a_j})$$

Learning translation probabilities

The only parameters that need to be learned are the **translation probabilities** $P(f | e)$

$$P(f_j = lac \mid e_i = lake)$$

If the training corpus had word alignments, we could simply count how often ‘lake’ is aligned to ‘lac’:

$$P(lac \mid lake) = \text{count}(lac, lake) / \sum_w \text{count}(w, lake)$$

But we don’t have word alignments.

So, instead of relative frequencies, we have to use *expected* relative frequencies:

$$P(lac \mid lake) = \langle \text{count}(lac, lake) \rangle / \langle \sum_w \text{count}(w, lake) \rangle$$

Training Model 1 with EM

The only parameters that need to be learned are the **translation probabilities** $P(f | e)$

We use the **EM algorithm** to estimate these parameters from a corpus with S sentence pairs $s = \langle f^{(s)}, e^{(s)} \rangle$ with alignments $\mathcal{A}(f^{(s)}, e^{(s)})$

- **Initialization:** guess $P(f | e)$
- **Expectation step:** compute expected counts

$$\langle c(f, e) \rangle = \sum_{s \in S} \langle c(f, e | e^{(s)}, f^{(s)}) \rangle$$

- **Maximization step:** recompute probabilities $P(f | e)$

$$\hat{P}(f | e) = \frac{\langle c(f, e) \rangle}{\sum_{f'} \langle c(f', e) \rangle}$$

Expectation-Maximization (EM)

1. Initialize a first model, M_0

2. **Expectation (E) step:**

Go through training data to gather expected counts
 $\langle \text{count}(lac, lake) \rangle$

3. **Maximization (M) step:**

Use expected counts to compute a new model M_{i+1}
 $P_{i+1}(lac | lake) = \langle \text{count}(lac, lake) \rangle / \langle \sum w \text{count}(w, lake) \rangle$

4. **Check for convergence:**

Compute log-likelihood of training data with M_{i+1}
If the difference between new and old log-likelihood
smaller than a threshold, stop. Else go to 2.

The E-step

Compute the expected count $\langle c(f, e | \mathbf{f}, \mathbf{e}) \rangle$:

$$\langle c(f, e | \mathbf{f}, \mathbf{e}) \rangle = \sum_{\mathbf{a} \in \mathcal{A}(\mathbf{f}, \mathbf{e})} P(\mathbf{a} | \mathbf{f}, \mathbf{e}) \cdot \underbrace{c(f, e | \mathbf{a}, \mathbf{e}, \mathbf{f})}_{\text{How often are } f, e \text{ aligned in } \mathbf{a}?$$

$$P(\mathbf{a} | \mathbf{f}, \mathbf{e}) = \frac{P(\mathbf{a}, \mathbf{f} | \mathbf{e})}{P(\mathbf{f} | \mathbf{e})} = \frac{P(\mathbf{a}, \mathbf{f} | \mathbf{e})}{\sum_{\mathbf{a}'} P(\mathbf{a}', \mathbf{f} | \mathbf{e})}$$

$$P(\mathbf{a}, \mathbf{f} | \mathbf{e}) = \prod_j P(f_j | e_{a_j})$$

$$\langle c(f, e | \mathbf{f}, \mathbf{e}) \rangle = \sum_{\mathbf{a} \in \mathcal{A}(\mathbf{f}, \mathbf{e})} \frac{\prod_j P(f_j | e_{a_j})}{\sum_{\mathbf{a}'} \prod_j P(f_j | e_{a'_j})} \cdot c(f, e | \mathbf{a}, \mathbf{e}, \mathbf{f})$$

Other translation models

Model 1 is a very simple (and not very good) translation model.

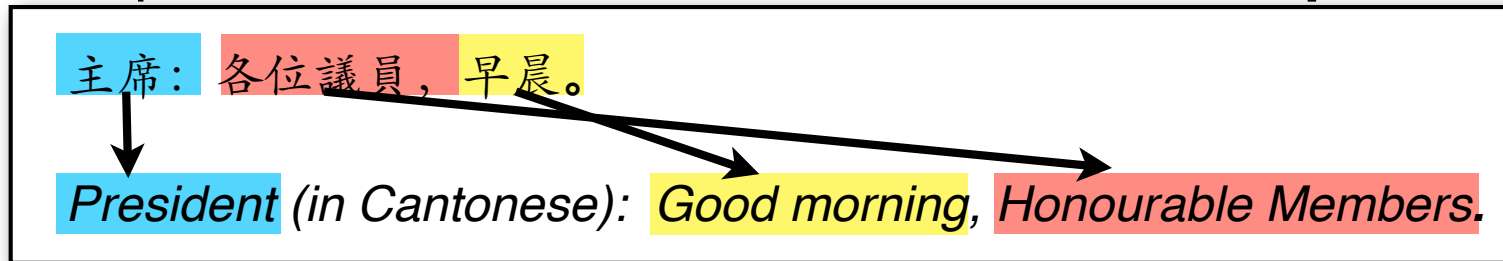
IBM models 2-5 are more complex. They take into account:

- “**fertility**”: the number of foreign words generated by each target word
- the **word order** and **string position** of the aligned words

Phrase-based translation models

Phrase-based translation models

Assumption: fundamental units of translation are **phrases**:



Phrase-based model of $P(F | E)$:

1. Split target sentence deterministically into phrases $ep_1 \dots ep_n$
2. Translate each target phrase ep_i into source phrase fp_i with translation probability $\phi(fp_i | ep_i)$
3. Reorder foreign phrases with distortion probability

$$d(a_i - b_{i-1}) = c^{|a_i - b_{i-1} - 1|}$$

a_i = start position of the source phrase generated by the current (i-th) target phrase, e_i

b_{i-1} = end position of the source phrase generated by the preceding (i-1-th) target phrase, e_{i-1}

$|a_i - b_{i-1} - 1|$ = #words in source sentence between translations of e_i , e_{i-1}

Phrase-based models of $P(f | e)$

Split target sentence $e = e_{1..n}$ into phrases $\mathbf{ep}_1.. \mathbf{ep}_N$:

[The green witch] [is] [at home] [this week]

Translate each target phrase \mathbf{ep}_i into source phrase \mathbf{fp}_i with **phrase translation probability** $\varphi(\mathbf{fp}_i | \mathbf{ep}_i)$:

[The green witch] = [die grüne Hexe], ...

Arrange the set of source phrases $\{ \mathbf{fp}_i \}$ to get f with **distortion probability** $d(a_i - b_{i-1})$

[Diese Woche] [ist] [die grüne Hexe] [zuhause]

NB: we can redefine $P_{\text{trans}}(\mathbf{fp}_i | \mathbf{ep}_i)$ as

$$P_{\text{trans}}(\mathbf{fp}_i | \mathbf{ep}_i) = \varphi(\mathbf{fp}_i | \mathbf{ep}_i) d(a_i - b_{i-1})$$

Translation probability $\varphi(fp_i | ep_i)$

Phrase translation probabilities can be obtained from a **phrase table**:

EP	FP	count
green witch	grüne Hexe	...
at home	zu Hause	10534
at home	daheim	9890
is	ist	598012
this week	diese Woche

This requires **phrase alignment**

Word alignment

	Diese	Woche	ist	die	grüne	Hexe	zuhaus
The							
green							
witch							
is							
at							
home							
this							
week							

Phrase alignment

	Diese	Woche	ist	die	grüne	Hexe	zuhaus
The							
green							
witch							
is							
at							
home							
this							
week							

Decoding (for phrase-based MT)

Translating

How do we translate a foreign sentence (e.g. “*Diese Woche ist die grüne Hexe zuhause*”) into English?

- We need to find $\hat{e} = \operatorname{argmax}_e P(f | e)P(e)$
- There is an exponential number of candidate translations e
- But we can look up phrase translations ep and $\varphi(fp | ep)$ in the phrase table

diese	Woche	ist	die	grüne	Hexe	zuhaus
this 0.2	week 0.7	is 0.8	the 0.3	green 0.3	witch 0.5	home 1.00
these 0.5			the green 0.4		sorceress 0.6	
this week 0.6				green witch 0.7		
is this week 0.4			the green witch 0.7			

Generating a (random) translation

1. Pick the first Target phrase ep_1 from the candidate list.

$$P := P_{LM}(<S> ep_1) P_{Trans}(fp_1 | ep_1)$$

$$E = \text{the}, F = \langle \dots \text{die} \dots \rangle$$

2. Pick the next target phrase ep_2 from the candidate list

$$P := P \times P_{LM}(ep_2 | ep_1) P_{Trans}(fp_2 | ep_2)$$

$$E = \text{the green witch}, F = \langle \dots \text{die grüne Hexe} \dots \rangle$$

3. Keep going: pick target phrases ep_i until the entire source sentence is translated

$$P := P \times P_{LM}(ep_i | ep_{1 \dots i-1}) P_{Trans}(fp_i | ep_i)$$

$$E = \text{the green witch is}, F = \langle \dots \text{ist die grüne Hexe} \dots \rangle$$

diese	Woche	ist	1 die	grüne	Hexe	zuhaus
this 0.2	week 0.7	3 is 0.8	1 the 0.3	green 0.3	witch 0.5	5 at home 0.5
these 0.5			the green 0.4		sorceress 0.6	
4 this week 0.6			2 green witch 0.7			
is this week 0.4			the green witch 0.7			

Finding the best translation

How can we find the *best* translation efficiently?

There is an exponential number of possible translations.

We will use a *heuristic search algorithm*

We cannot guarantee to find the best (= highest-scoring) translation, but we're likely to get close.

We will use a *“stack-based” decoder*

(If you've taken Intro to AI: this is A* (“A-star”) search)

We will score partial translations based on how good we expect the corresponding completed translation to be.

Or, rather: we will score partial translations on how **bad** we expect the corresponding complete translation to be.

That is, our scores will be **costs (high=bad, low=good)**

Scoring partial translations

Assign **expected costs** to *partial* translations (E, F) :

$$\begin{aligned} \text{expected_cost}(E, F) = & \text{current_cost}(E, F) \\ & + \text{future_cost}(E, F) \end{aligned}$$

The **current cost** is based on the score of the partial translation (E, F)

The **(estimated) future cost** is an **upper bound** on the actual cost of completing the partial translation (E, F) :

$$\begin{aligned} \text{true_cost}(E, F) & (= \text{current_cost}(E, F) + \text{actual_future_cost}(E, F)) \\ & \leq \text{expected_cost}(E, F) (= \text{current_cost}(E, F) + \text{future_cost}(E, F)) \end{aligned}$$

because $\text{actual_future_cost}(E, F) \leq \text{future_cost}(E, F)$

Stack-based decoding

Maintain a **priority queue** (=‘stack’) of **partial translations** (hypotheses) with their **expected costs**.

Each element on the stack is **open** (we haven’t yet pursued this hypothesis) or **closed** (we have already pursued this hypothesis)

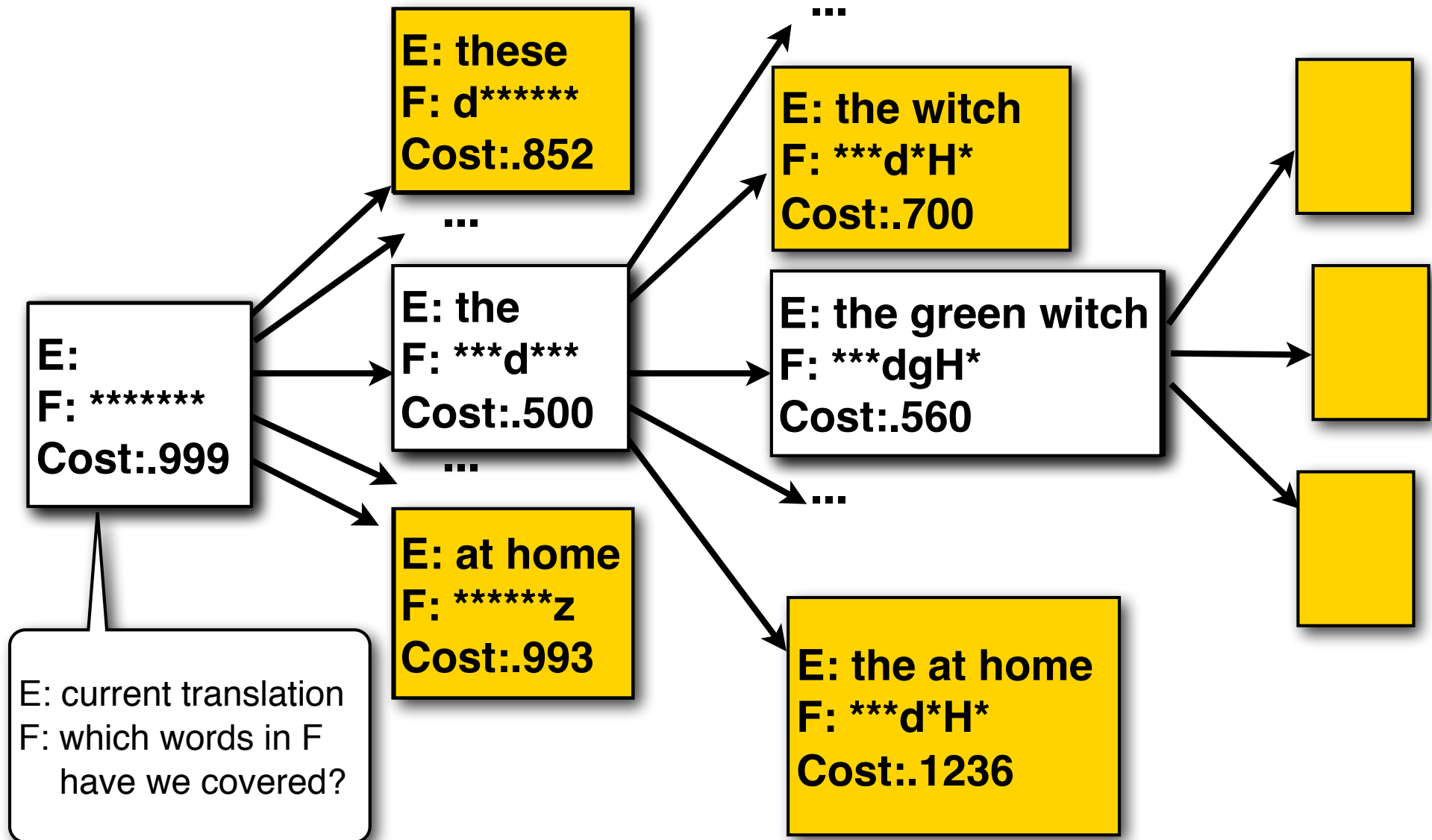
At each step:

- **Expand** the best open hypothesis (the open translation with the lowest expected cost) in all possible ways.
- These new translations become **new open elements** on the stack.
- **Close** the best open hypothesis.

Additional Pruning (n -best / beam search):

Only keep the n best open hypotheses around

Stack-based decoding



MT evaluation

Automatic evaluation: BLEU

Evaluate candidate translations against several reference translations.

C1: It is a guide to action which ensures that the military always obeys the commands of the party.

C2: It is to insure the troops forever hearing the activity guidebook that party direct

R1: It is a guide to action that ensures that the military will forever heed Party commands.

R2: It is the guiding principle which guarantees the military forces always being under the command of the Party.

R3: It is the practical guide for the army always to heed the directions of the party.

The **BLEU score** is based on **N-gram precision**:

How many n-grams in the candidate translation occur also in one of the reference translation?

BLEU details

For $n \in \{1, \dots, 4\}$, compute the (modified) **precision of all n -grams**:

$$Prec_n = \frac{\sum_{c \in C} \sum_{n\text{-gram} \in c} \text{MaxFreq}_{\text{ref}}(n\text{-gram})}{\sum_{c \in C} \sum_{n\text{-gram} \in c} \text{Freq}_c(n\text{-gram})}$$

$\text{MaxFreq}_{\text{ref}}(\text{'the party'})$ = max. count of *'the party'* in **one** reference translation.

$\text{Freq}_c(\text{'the party'})$ = count of *'the party'* in candidate translation c .

Penalize short candidate translations by a **brevity penalty BP**

c = length (number of words) of the whole candidate translation corpus

r = Pick for each candidate the reference translation that is closest in length;
sum up these lengths.

Brevity penalty $BP = \exp(1 - c/r)$ for $c \leq r$; $BP = 1$ for $c > r$
(BP ranges from e for $c = 0$ to 1 for $c = r$)

BLEU score

The BLEU score is the geometric mean of the precision of the unigrams, bigrams, trigrams, quadrigrams, weighted by the brevity penalty BP.

$$\mathbf{BLEU} = BP \times \exp \left(\frac{1}{N} \sum_{n=1}^N \log Prec_n \right)$$

Human evaluation

We want to know whether the translation is “**good**” **English**, and whether it is an **accurate translation** of the original.

- Ask human raters to judge the **fluency** and the **adequacy** of the translation (e.g. on a scale of 1 to 5)
- Correlated with **fluency** is accuracy on **cloze task**:
Give rater the sentence with one word replaced by blank.
Ask rater to guess the missing word in the blank.
- Similar to **adequacy** is **informativeness**
Can you use the translation to perform some task
(e.g. answer multiple-choice questions about the text)

Summary:

Machine Translation

Machine translation models

Current MT models all rely on statistics.

Many current models do estimate $P(E \mid F)$ directly, but may use features based on language models (capturing $P(E)$) and IBM-style translation models ($P(F \mid E)$) internally.

There are a number of syntax-based models, e.g. using synchronous context-free grammars, which consist of pairs of rules for the two languages in which each RHS NT in language A corresponds to a RHS NT in language B:

Language A: $XP \rightarrow YP ZP$ Language B: $XP \rightarrow ZP YP$

Very recent developments

Neural network-based approaches:

Recurrent neural networks (RNN) can model sequences (e.g. strings, sentences, etc.)

Use one RNN (the encoder) to process the input in the source language

Pass its output to another RNN (the decoder) to generate the output in the target language

See e.g. http://www.tensorflow.org/tutorials/seq2seq/index.md#sequence-to-sequence_basics

Today's key concepts

Why is machine translation hard?

Linguistic divergences: morphology, syntax, semantics

Different approaches to machine translation:

Vauquois triangle

Statistical MT: Noisy Channel, IBM Model 1 (more on this next time)