Machine Translation

Philipp Koehn

13 June 2024



Agenda



- Understanding the translation problem
- Statistical methods
- Neural method (including Transformer)



some history

An Old Idea



Warren Weaver on translation as code breaking (1947):

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".



Georgetown experiment Machine could translate

1954

• Excited research in 1950s and 1960s

250 words and 6 grammar rules

• 1966 ALPAC report:

- only \$20 million spent on translation in the US per year
- no point in machine translation

Early Efforts and Disappointment





Rule-Based Systems



- Rule-based systems
 - build dictionaries
 - write transformation rules
 - refine, refine, refine
- Météo system for weather forecasts (1976)
- Systran (1968), Logos and Metal (1980s)

```
"have" :=
```

if

```
subject (animate)
and object (owned-by-subject)
then
translate to "kade... aahe"
if
subject (animate)
and object (kinship-with-subject)
then
translate to "laa... aahe"
if
subject (inanimate)
then
translate to "madhye...
aahe"
```

Statistical Machine Translation



- 1980s: IBM
- 1990s: increased research
- Mid 2000s: Phrase-Based MT (Moses, Google)
- Around 2010: commercial viability



Neural Machine Translation



- Late 2000s: neural models for computer vision
- Since mid 2010s: neural models for machine translation
- 2016: Neural machine translation the new state of the art



Hype 8 **Neural** MT Georgetown Hype **Statistica** experiment MT **Expert systems / 5th generation Al** Reality 1950 1960 1970 1980 2000 2010 1990 2020



why is it hard?

Word Translation Problems



• Words are ambiguous

He deposited money in a bank account with a high interest rate.

Sitting on the bank of the Mississippi, a passing ship piqued his interest.

- How do we find the right meaning, and thus translation?
- Context should be helpful

Syntactic Translation Problems



• Languages have different sentence structure

das	behaupten	sie	wenigstens
this	claim	they	at least
the		she	

- Convert from object-verb-subject (OVS) to subject-verb-object (SVO)
- Ambiguities can be resolved through syntactic analysis
 - the meaning the of das not possible (not a noun phrase)
 - the meaning she of sie not possible (subject-verb agreement)

Semantic Translation Problems



• Pronominal anaphora

I saw the movie and it is good.

- How to translate it into German (or French)?
 - it refers to movie
 - movie translates to Film
 - Film has masculine gender
 - ergo: it must be translated into masculine pronoun er
- We are not handling this very well [Le Nagard and Koehn, 2010]

Semantic Translation Problems



• Coreference

Whenever I visit my uncle and his daughters, I can't decide who is my favorite cousin.

- How to translate cousin into German? Male or female?
- Complex inference required

No Single Right Answer



这个 机场 的 安全 工作 由 以色列 方面 负责.

Israeli officials are responsible for airport security. Israel is in charge of the security at this airport. The security work for this airport is the responsibility of the Israel government. Israeli side was in charge of the security of this airport. Israel is responsible for the airport's security. Israel is responsible for safety work at this airport. Israel presides over the security of the airport. Israel took charge of the airport security. The safety of this airport is taken charge of by Israel. This airport's security is the responsibility of the Israeli security officials.



data-driven models

Word Alignment





Phrase-Based Model





- Foreign input is segmented in phrases
- Each phrase is translated into English
- Phrases are reordered
- Workhorse of today's statistical machine translation

Syntax-Based Translation





Semantic Translation



• Abstract meaning representation [Knight et al., ongoing]

- Generalizes over equivalent syntactic constructs (e.g., active and passive)
- Defines semantic relationships
 - semantic roles
 - co-reference
 - discourse relations

Neural Model







what is it good for?



what is it good *enough* for?

Why Machine Translation?



Assimilation — reader initiates translation, wants to know content

- user is tolerant of inferior quality
- focus of majority of research (GALE program, etc.)

Communication — participants don't speak same language, rely on translation

- users can ask questions, when something is unclear
- chat room translations, hand-held devices
- often combined with speech recognition, IWSLT campaign

Dissemination — publisher wants to make content available in other languages

- high demands for quality
- currently almost exclusively done by human translators

Quality



HTER assessment

0%	publishable
10%	oditable
20%	euitable
30%	gistable
40%	triagable
50%	

(scale developed in preparation of DARPA GALE programme)

Applications



HTER	assessment	application examples
0%	publishable	Seamless bridging of language divide
10%	publishable	Automatic publication of official announcements
	editable	Increased productivity of human translators
20%		Access to official publications Multi-lingual communication (chat, social networks)
30%	gistable	Information gathering Trend spotting
40%	triagable	Identifying relevant documents

50%



ibm model 1

and the

expectation maximization algorithm

Lexical Translation



• How to translate a word \rightarrow look up in dictionary

Haus — house, building, home, household, shell.

- Multiple translations
 - some more frequent than others
 - for instance: house, and building most common
 - special cases: Haus of a snail is its shell
- Note: In all lectures, we translate from a foreign language into English

Collect Statistics



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

Translation of <i>Haus</i>	Count
house	8,000
building	1,600
home	200
household	150
shell	50

Estimate Translation Probabilities



Maximum likelihood estimation

$$p_f(e) = \begin{cases} 0.8 & \text{if } e = \text{house}, \\ 0.16 & \text{if } e = \text{building}, \\ 0.02 & \text{if } e = \text{home}, \\ 0.015 & \text{if } e = \text{household}, \\ 0.005 & \text{if } e = \text{shell}. \end{cases}$$

Alignment



• In a parallel text (or when we translate), we align words in one language with the words in the other



• Word positions are numbered 1–4

Alignment Function



- Formalizing alignment with an alignment function
- Mapping an English target word at position i to a German source word at position j with a function $a : i \to j$
- Example

$$a: \{1 \rightarrow 1, 2 \rightarrow 2, 3 \rightarrow 3, 4 \rightarrow 4\}$$

Reordering



Words may be reordered during translation



 $a: \{1 \rightarrow 3, 2 \rightarrow 4, 3 \rightarrow 2, 4 \rightarrow 1\}$

One-to-Many Translation



A source word may translate into multiple target words



 $a: \{1 \rightarrow 1, 2 \rightarrow 2, 3 \rightarrow 3, 4 \rightarrow 4, 5 \rightarrow 4\}$

Dropping Words



Words may be dropped when translated (German article das is dropped)



 $a: \{1 \rightarrow 2, 2 \rightarrow 3, 3 \rightarrow 4\}$

Inserting Words



- Words may be added during translation
 - The English just does not have an equivalent in German
 - We still need to map it to something: special NULL token


IBM Model 1



- Generative model: break up translation process into smaller steps
 - IBM Model 1 only uses lexical translation
- Translation probability
 - for a foreign sentence $\mathbf{f} = (f_1, ..., f_{l_f})$ of length l_f
 - to an English sentence $\mathbf{e} = (e_1, ..., e_{l_e})$ of length l_e
 - with an alignment of each English word e_j to a foreign word f_i according to the alignment function $a : j \rightarrow i$

$$p(\mathbf{e}, a | \mathbf{f}) = \frac{\epsilon}{(l_f + 1)^{l_e}} \prod_{j=1}^{l_e} t(e_j | f_{a(j)})$$

– parameter
$$\epsilon$$
 is a normalization constant

Example



das		Haus			ist		klein	
e	t(e f)	e	t(e f)		e	t(e f)	e	t(e f)
the	0.7	house	0.8		is	0.8	small	0.4
that	0.15	building	0.16		'S	0.16	little	0.4
which	0.075	home	0.02		exists	0.02	short	0.1
who	0.05	household	0.015		has	0.015	minor	0.06
this	0.025	shell	0.005		are	0.005	petty	0.04

 $p(e, a|f) = \frac{\epsilon}{5^4} \times t(\text{the}|\text{das}) \times t(\text{house}|\text{Haus}) \times t(\text{is}|\text{ist}) \times t(\text{small}|\text{klein})$ $= \frac{\epsilon}{5^4} \times 0.7 \times 0.8 \times 0.8 \times 0.4$ $= 0.00029\epsilon$



finding translations

Centauri-Arcturan Parallel Text



1a. ok-voon ororok sprok .1b. at-voon bichat dat .

2a. ok-drubel ok-voon anok plok sprok .2b. at-drubel at-voon pippat rrat dat .

3a. erok sprok izok hihok ghirok .3b. totat dat arrat vat hilat .

4a. ok-voon anok drok brok jok .4b. at-voon krat pippat sat lat .

5a. wiwok farok izok stok . 5b. totat jjat quat cat .

6a. lalok sprok izok jok stok .6b. wat dat krat quat cat .

7a. lalok farok ororok lalok sprok izok enemok .7b. wat jjat bichat wat dat vat eneat .

8a. lalok brok anok plok nok .8b. iat lat pippat rrat nnat .

9a. wiwok nok izok kantok ok-yurp .9b. totat nnat quat oloat at-yurp .

10a. lalok mok nok yorok ghirok clok .10b. wat nnat gat mat bat hilat .

11a. lalok nok crrrok hihok yorok zanzanok .11b. wat nnat arrat mat zanzanat .

12a. lalok rarok nok izok hihok mok .12b. wat nnat forat arrat vat gat .

Translation challenge: farok crrrok hihok yorok clok kantok ok-yurp

(from Knight (1997): Automating Knowledge Acquisition for Machine Translation)



em algorithm

Learning Lexical Translation Models



- We would like to estimate the lexical translation probabilities t(e|f) from a parallel corpus
- ... but we do not have the alignments
- Chicken and egg problem
 - − if we had the *alignments*,
 → we could estimate the *parameters* of our generative model
 - if we had the *parameters*,
 - \rightarrow we could estimate the *alignments*



- Incomplete data
 - if we had *complete data*, would could estimate *model*
 - if we had *model*, we could fill in the *gaps in the data*
- Expectation Maximization (EM) in a nutshell
 - 1. initialize model parameters (e.g. uniform)
 - 2. assign probabilities to the missing data
 - 3. estimate model parameters from completed data
 - 4. iterate steps 2–3 until convergence





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





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





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





- Convergence
- Inherent hidden structure revealed by EM





• Parameter estimation from the aligned corpus

IBM Model 1 and EM



- EM Algorithm consists of two steps
- Expectation-Step: Apply model to the data
 - parts of the model are hidden (here: alignments)
 - using the model, assign probabilities to possible values
- Maximization-Step: Estimate model from data
 - take assign values as fact
 - collect counts (weighted by probabilities)
 - estimate model from counts
- Iterate these steps until convergence

IBM Model 1 and EM



- We need to be able to compute:
 - Expectation-Step: probability of alignments
 - Maximization-Step: count collection

IBM Model 1 and EM



- Probabilities p(the|la) = 0.7 p(house|la) = 0.05p(the|maison) = 0.1 p(house|maison) = 0.8
- Alignments

 $\begin{aligned} \mathbf{a} & \bullet & \bullet & \mathsf{the} \\ \mathsf{maison} & \bullet & \mathsf{house} \\ p(\mathbf{e}, a | \mathbf{f}) = 0.56 \\ p(\mathbf{e}, a | \mathbf{f}) = 0.035 \\ p(\mathbf{e}, a | \mathbf{f}) = 0.035 \\ p(\mathbf{e}, a | \mathbf{f}) = 0.08 \\ p(\mathbf{e}, a | \mathbf{f}) = 0.005 \\ p(\mathbf{e}, a | \mathbf{f}) = 0.005 \\ p(\mathbf{e}, a | \mathbf{f}) = 0.005 \\ p(\mathbf{e}, a | \mathbf{f}) = 0.007 \end{aligned}$

• Counts c(the|la) = 0.824 + 0.052 c(house|la) = 0.052 + 0.007c(the|maison) = 0.118 + 0.007 c(house|maison) = 0.824 + 0.118



- We need to compute $p(a|\mathbf{e}, \mathbf{f})$
- Applying the chain rule:

$$p(a|\mathbf{e}, \mathbf{f}) = \frac{p(\mathbf{e}, a|\mathbf{f})}{p(\mathbf{e}|\mathbf{f})}$$

• We already have the formula for $p(\mathbf{e}, \mathbf{a}|\mathbf{f})$ (definition of Model 1)



• We need to compute $p(\mathbf{e}|\mathbf{f})$

$$p(\mathbf{e}|\mathbf{f}) = \sum_{a} p(\mathbf{e}, a|\mathbf{f})$$

= $\sum_{a(1)=0}^{l_f} \dots \sum_{a(l_e)=0}^{l_f} p(\mathbf{e}, a|\mathbf{f})$
= $\sum_{a(1)=0}^{l_f} \dots \sum_{a(l_e)=0}^{l_f} \frac{\epsilon}{(l_f+1)^{l_e}} \prod_{j=1}^{l_e} t(e_j|f_{a(j)})$



$$p(\mathbf{e}|\mathbf{f}) = \sum_{a(1)=0}^{l_f} \dots \sum_{a(l_e)=0}^{l_f} \frac{\epsilon}{(l_f+1)^{l_e}} \prod_{j=1}^{l_e} t(e_j|f_{a(j)})$$
$$= \frac{\epsilon}{(l_f+1)^{l_e}} \sum_{a(1)=0}^{l_f} \dots \sum_{a(l_e)=0}^{l_f} \prod_{j=1}^{l_e} t(e_j|f_{a(j)})$$
$$= \frac{\epsilon}{(l_f+1)^{l_e}} \prod_{j=1}^{l_e} \sum_{i=0}^{l_f} t(e_j|f_i)$$

- Note the trick in the last line
 - removes the need for an exponential number of products
 - \rightarrow this makes IBM Model 1 estimation tractable

The Trick



(case
$$l_e = l_f = 2$$
)

$$\begin{split} \sum_{a(1)=0}^{2} \sum_{a(2)=0}^{2} &= \frac{\epsilon}{3^{2}} \prod_{j=1}^{2} t(e_{j}|f_{a(j)}) = \\ &= t(e_{1}|f_{0}) \ t(e_{2}|f_{0}) + t(e_{1}|f_{0}) \ t(e_{2}|f_{1}) + t(e_{1}|f_{0}) \ t(e_{2}|f_{2}) + \\ &+ t(e_{1}|f_{1}) \ t(e_{2}|f_{0}) + t(e_{1}|f_{1}) \ t(e_{2}|f_{1}) + t(e_{1}|f_{1}) \ t(e_{2}|f_{2}) + \\ &+ t(e_{1}|f_{2}) \ t(e_{2}|f_{0}) + t(e_{1}|f_{2}) \ t(e_{2}|f_{1}) + t(e_{1}|f_{2}) \ t(e_{2}|f_{2}) = \\ &= t(e_{1}|f_{0}) \ (t(e_{2}|f_{0}) + t(e_{2}|f_{1}) + t(e_{2}|f_{2})) + \\ &+ t(e_{1}|f_{1}) \ (t(e_{2}|f_{0}) + t(e_{2}|f_{1}) + t(e_{2}|f_{2})) + \\ &+ t(e_{1}|f_{2}) \ (t(e_{2}|f_{0}) + t(e_{2}|f_{1}) + t(e_{2}|f_{2})) = \\ &= (t(e_{1}|f_{0}) + t(e_{1}|f_{1}) + t(e_{1}|f_{2})) \ (t(e_{2}|f_{0}) + t(e_{2}|f_{1}) + t(e_{2}|f_{2})) \end{split}$$



• Combine what we have:

 $p(\mathbf{a}|\mathbf{e}, \mathbf{f}) = p(\mathbf{e}, \mathbf{a}|\mathbf{f}) / p(\mathbf{e}|\mathbf{f})$ $= \frac{\frac{\epsilon}{(l_f + 1)^{l_e}} \prod_{j=1}^{l_e} t(e_j|f_{a(j)})}{\frac{\epsilon}{(l_f + 1)^{l_e}} \prod_{j=1}^{l_e} \sum_{i=0}^{l_f} t(e_j|f_i)}$ $= \prod_{j=1}^{l_e} \frac{t(e_j|f_{a(j)})}{\sum_{i=0}^{l_f} t(e_j|f_i)}$

IBM Model 1 and EM: Maximization Step 56

- Now we have to collect counts
- Evidence from a sentence pair **e**,**f** that word *e* is a translation of word *f*:

$$c(e|f; \mathbf{e}, \mathbf{f}) = \sum_{a} p(a|\mathbf{e}, \mathbf{f}) \sum_{j=1}^{l_e} \delta(e, e_j) \delta(f, f_{a(j)})$$

• With the same simplication as before:

$$c(e|f; \mathbf{e}, \mathbf{f}) = \frac{t(e|f)}{\sum_{i=0}^{l_f} t(e|f_i)} \sum_{j=1}^{l_e} \delta(e, e_j) \sum_{i=0}^{l_f} \delta(f, f_i)$$



After collecting these counts over a corpus, we can estimate the model:

$$t(e|f;\mathbf{e},\mathbf{f}) = \frac{\sum_{(\mathbf{e},\mathbf{f})} c(e|f;\mathbf{e},\mathbf{f}))}{\sum_{e} \sum_{(\mathbf{e},\mathbf{f})} c(e|f;\mathbf{e},\mathbf{f}))}$$

IBM Model 1 and EM: Pseudocode



Inp	Put: set of sentence pairs (\mathbf{e}, \mathbf{f}) trut: translation prob $t(e f)$
1:	initialize $t(e f)$ uniformly
2:	while not converged do
3:	// initialize
4:	count(e f) = 0 for all e, f
5:	total(f) = 0 for all f
6:	for all sentence pairs (e , f) do
7:	// compute normalization
8:	for all words e in e do
9:	s-total $(e) = 0$
10:	for all words <i>f</i> in f do
11:	s-total(e) += $t(e f)$
12:	end for
13:	end for

14:	// collect counts
15:	for all words <i>e</i> in e do
16:	for all words <i>f</i> in f do
17:	$\operatorname{count}(e f) += \frac{t(e f)}{s - total(e)}$
18:	$total(f) += \frac{t(e f)}{s-total(e)}$
19:	end for
20:	end for
21:	end for
22:	// estimate probabilities
23:	for all foreign words <i>f</i> do
24:	for all English words <i>e</i> do
25:	$t(e f) = \frac{\operatorname{count}(e f)}{\operatorname{total}(f)}$
26:	end for
27:	end for
28:	end while

Convergence





Perplexity



- How well does the model fit the data?
- Perplexity: derived from probability of the training data according to the model

$$\log_2 PP = -\sum_s \log_2 p(\mathbf{e}_s | \mathbf{f}_s)$$

• Example (ϵ =1)

	initial	1st it.	2nd it.	3rd it.	•••	final
p(the haus das haus)	0.0625	0.1875	0.1905	0.1913	•••	0.1875
p(the book das buch)	0.0625	0.1406	0.1790	0.2075	•••	0.25
p(a book ein buch)	0.0625	0.1875	0.1907	0.1913	•••	0.1875
perplexity	4095	202.3	153.6	131.6	•••	113.8

Higher IBM Models



IBM Model 1	lexical translation
IBM Model 2	adds absolute reordering model
IBM Model 3	adds fertility model
IBM Model 4	relative reordering model
IBM Model 5	fixes deficiency

- Only IBM Model 1 has global maximum
 - training of a higher IBM model builds on previous model
- Computionally biggest change in Model 3
 - trick to simplify estimation does not work anymore
 - \rightarrow exhaustive count collection becomes computationally too expensive
 - sampling over high probability alignments is used instead



phrase-based models

Motivation



- Word-Based Models translate *words* as atomic units
- Phrase-Based Models translate *phrases* as atomic units
- Advantages:
 - many-to-many translation can handle non-compositional phrases
 - use of local context in translation
 - the more data, the longer phrases can be learned
- "Standard Model", used by Google Translate and others until about 2017

Phrase-Based Model





- Foreign input is segmented in phrases
- Each phrase is translated into English
- Phrases are reordered

Real Example



• Phrase translations for den Vorschlag learned from the Europarl corpus:

English	$\phi(ar{e} ar{f})$	English	$\phi(ar{e} ar{f})$
the proposal	0.6227	the suggestions	0.0114
's proposal	0.1068	the proposed	0.0114
a proposal	0.0341	the motion	0.0091
the idea	0.0250	the idea of	0.0091
this proposal	0.0227	the proposal,	0.0068
proposal	0.0205	its proposal	0.0068
of the proposal	0.0159	it	0.0068
the proposals	0.0159		•••

- lexical variation (proposal vs suggestions)
- morphological variation (proposal vs proposals)
- included function words (the, a, ...)
- noise (it)

Noisy Channel Model



- We would like to integrate a language model
- Bayes rule

$$\begin{aligned} \operatorname{argmax}_{\mathbf{e}} p(\mathbf{e}|\mathbf{f}) &= \operatorname{argmax}_{\mathbf{e}} \frac{p(\mathbf{f}|\mathbf{e}) \ p(\mathbf{e})}{p(\mathbf{f})} \\ &= \operatorname{argmax}_{\mathbf{e}} p(\mathbf{f}|\mathbf{e}) \ p(\mathbf{e}) \end{aligned}$$



decoding

Decoding



• We have a mathematical model for translation

$p(\mathbf{e}|\mathbf{f})$

• Task of decoding: find the translation \mathbf{e}_{best} with highest probability

 $\mathbf{e}_{\text{best}} = \operatorname{argmax}_{\mathbf{e}} p(\mathbf{e}|\mathbf{f})$

- Two types of error
 - the most probable translation is bad \rightarrow fix the model
 - search does not find the most probably translation \rightarrow fix the search
- Decoding is evaluated by search error, not quality of translations (although these are often correlated)

Translation Process



• Task: translate this sentence from German into English



• Pick phrase in input, translate

Translation Process



• Task: translate this sentence from German into English



- Pick phrase in input, translate
 - it is allowed to pick words out of sequence reordering
 - phrases may have multiple words: many-to-many translation

Translation Process



• Task: translate this sentence from German into English



• Pick phrase in input, translate
Translation Process



• Task: translate this sentence from German into English



• Pick phrase in input, translate



decoding process

Translation Options





- Many translation options to choose from
 - in Europarl phrase table: 2727 matching phrase pairs for this sentence
 - by pruning to the top 20 per phrase, 202 translation options remain

Translation Options





- The machine translation decoder does not know the right answer
 - picking the right translation options
 - arranging them in the right order
- \rightarrow Search problem solved by heuristic beam search

Decoding: Precompute Translation Options 76

er	geht	ja	nicht	nach	hause

consult phrase translation table for all input phrases





er	geht	ja	nicht	nach	hause



initial hypothesis: no input words covered, no output produced

Decoding: Hypothesis Expansion







pick any translation option, create new hypothesis

Decoding: Hypothesis Expansion







create hypotheses for all other translation options

Decoding: Hypothesis Expansion





also create hypotheses from created partial hypothesis

Decoding: Find Best Path





backtrack from highest scoring complete hypothesis



dynamic programming

Computational Complexity



- The suggested process creates exponential number of hypothesis
- Machine translation decoding is NP-complete
- Reduction of search space:
 - recombination (risk-free)
 - pruning (risky)

Recombination



- Two hypothesis paths lead to two matching hypotheses
 - same foreign words translated
 - same English words in the output



• Worse hypothesis is dropped



Recombination



- Two hypothesis paths lead to hypotheses indistinguishable in subsequent search
 - same foreign words translated
 - same last two English words in output (assuming trigram language model)
 - same last foreign word translated



• Worse hypothesis is dropped



Restrictions on Recombination



- Translation model: Phrase translation independent from each other
 → no restriction to hypothesis recombination
- Language model: Last n 1 words used as history in *n*-gram language model \rightarrow recombined hypotheses must match in their last n 1 words
- **Reordering model:** Distance-based reordering model based on distance to end position of previous input phrase
 - \rightarrow recombined hypotheses must have that same end position
- Other feature function may introduce additional restrictions



pruning

Pruning



• Recombination reduces search space, but not enough (we still have a NP complete problem on our hands)

- Pruning: remove bad hypotheses early
 - put comparable hypothesis into stacks
 (hypotheses that have translated same number of input words)
 - limit number of hypotheses in each stack

Stacks





- Hypothesis expansion in a stack decoder
 - translation option is applied to hypothesis
 - new hypothesis is dropped into a stack further down

Stack Decoding Algorithm



- 1: place empty hypothesis into stack 0
- 2: for all stacks 0...n 1 do
- 3: **for all** hypotheses in stack **do**
- 4: **for all** translation options **do**
- 5: **if** applicable **then**
- 6: create new hypothesis
- 7: place in stack
- 8: recombine with existing hypothesis **if** possible
- 9: prune stack **if** too big
- 10: **end if**
- 11: **end for**
- 12: **end for**
- 13: **end for**

Pruning



- Pruning strategies
 - histogram pruning: keep at most k hypotheses in each stack
 - stack pruning: keep hypothesis with score $\alpha \times$ best score ($\alpha < 1$)
- Computational time complexity of decoding with histogram pruning

O(max stack size × translation options × sentence length)

• Number of translation options is linear with sentence length, hence:

 $O(\max \text{ stack size} \times \text{ sentence length}^2)$

• Quadratic complexity



break