

89688: Statistical Machine Translation **Phrase-Based Models and Decoding**

Roee Aharoni Computer Science Department Bar Ilan University

Based in part on slides by Phillip Koehn

May 2020



• So far, we discussed IBM Model1 - ignores word positions

2020

 $p(f_1 \ldots f_m, a_1 \ldots a_m | e_1 \ldots e_l, m) =$





- So far, we discussed IBM Model1 ignores word positions
- The additional IBM models are:

2020

 $p(f_1 \ldots f_m, a_1 \ldots a_m | e_1 \ldots e_l, m) =$





- So far, we discussed IBM Model1 ignores word positions
- The additional IBM models are:
 - Model 2 Word positions matter! (and Model1 is a special case)

2020

$$p(f_1 \dots f_m, a_1 \dots a_m | e_1 \dots e_l, m) =$$



 $p(f_1 \dots f_m, a_1 \dots a_m | e_1 \dots e_l, m) =$

$$=\prod_{i=1}^{m} q(a_i|i,l,m)t(f_i|e_{a_i})$$



- So far, we discussed IBM Model1 ignores word positions
- The additional IBM models are:
 - Model 2 Word positions matter! (and Model1 is a special case)
 - Model 3 Introduce *Fertility*



Unabhaengigkeitserklaerung

0	0.00008
I	0.1
2	0.0002
3	0.8
4	0.009
5	0

zum = (zu + dem)

0	0.01
I	0
2	0.9
3	0.0009
4	0.0001
5	0

Haus

0	0.0
I	0.9
2	0.0
3	0
4	0
5	0



I
2
7

- So far, we discussed IBM Model1 ignores word positions
- The additional IBM models are:
 - Model 2 Word positions matter! (and Model1 is a special case)
 - Model 3 Introduce *Fertility*
 - Model 4 *Relative* word positions





- So far, we discussed IBM Model1 ignores word positions
- The additional IBM models are:
 - Model 2 Word positions matter! (and Model1 is a special case)
 - Model 3 Introduce *Fertility*
 - Model 4 *Relative* word positions
 - Model 5 Avoids aligning target words to the same source word







• IBM models create a many-to-one mapping







- IBM models create a many-to-one mapping
- One source word can generate several target words







- IBM models create a many-to-one mapping
- One source word can generate several target words
- But we need many-to-many mappings...







- IBM models create a many-to-one mapping
- One source word can generate several target words
- But we need many-to-many mappings...
 - "Prime Minister"







ROEE AHARONI







• Koehn, Och & Marcu (2003)





- Koehn, Och & Marcu (2003)
 - Input is segmented into phrases





- Koehn, Och & Marcu (2003)
 - Input is segmented into phrases
 - Each phrase is translated into another phrase





- Koehn, Och & Marcu (2003)
 - Input is segmented into phrases
 - Each phrase is translated into another phrase
 - Phrases are reordered





- Koehn, Och & Marcu (2003)
 - Input is segmented into phrases
 - Each phrase is translated into another phrase
 - Phrases are reordered



 Many-to-many translations can handle non-compositional cases





- Koehn, Och & Marcu (2003)
 - Input is segmented into phrases
 - Each phrase is translated into another phrase
 - Phrases are reordered



2020

- Many-to-many translations can handle non-compositional cases
- Models local context





- Koehn, Och & Marcu (2003)
 - Input is segmented into phrases
 - Each phrase is translated into another phrase
 - Phrases are reordered



2020

- Many-to-many translations can handle non-compositional cases
- Models local context
- First Google Translate models were phrase-based (2006)







Start with word alignments





- Start with word alignments
- Preprocess using a heuristic "Grow-Diag-Final" (Och & Ney, 2003)

2020





- Start with word alignments
- Preprocess using a heuristic "Grow-Diag-Final" (Och & Ney, 2003)
- Collect **all** phrase pairs that are **consistent** with the word alignment







2020

- Start with word alignments
- Preprocess using a heuristic "Grow-Diag-Final" (Och & Ney, 2003)
- Collect **all** phrase pairs that are **consistent** with the word alignment
 - Phrase alignment has to contain all alignment points for **all** covered words









• We can learn phrase pairs with varying lengths:



(bruja verde, green witch), (Maria no daba una bofetada, Mary did not slap), (no daba una bofetada a la, did not slap the), (a la bruja verde, the green witch)

(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch), (verde, green), (Maria no, Mary did not), (no daba una bofetada, did not slap), (daba una bofetada a la, slap the),



- We can learn phrase pairs with varying lengths:
- Too sparse with little data!



(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch), (verde, green), (Maria no, Mary did not), (no daba una bofetada, did not slap), (daba una bofetada a la, slap the), (bruja verde, green witch), (Maria no daba una bofetada, Mary did not slap), (no daba una bofetada a la, did not slap the), (a la bruja verde, the green witch)



- We can learn phrase pairs with varying lengths:
- Too sparse with little data!
- Should be tuned



(Maria, Mary), (no, did not), (slap, daba una bofetada), (a la, the), (bruja, witch), (verde, green), (Maria no, Mary did not), (no daba una bofetada, did not slap), (daba una bofetada a la, slap the), (bruja verde, green witch), (Maria no daba una bofetada, Mary did not slap), (no daba una bofetada a la, did not slap the), (a la bruja verde, the green witch)



ROEE AHARONI

2020



• We want to solve the following problem:

2020



- We want to solve the following problem:
- This task is called **decoding**

2020



- We want to solve the following problem:
- This task is called **decoding**
- We need to **search** for the most probable translation



- We want to solve the following problem:
- This task is called **decoding**
- We need to **search** for the most probable translation
- For example:


- We want to solve the following problem:
- This task is called **decoding**
- We need to **search** for the most probable translation
- For example:

ja geht er

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

nicht nach hause



- We want to solve the following problem:
- This task is called **decoding**
- We need to **search** for the most probable translation
- For example:



2020

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

nicht nach hause



- We want to solve the following problem:
- This task is called **decoding**
- We need to **search** for the most probable translation
- For example:



2020

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



- We want to solve the following problem:
- This task is called **decoding**
- We need to **search** for the most probable translation
- For example:



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



- We want to solve the following problem:
- This task is called **decoding**
- We need to **search** for the most probable translation
- For example:



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



ROEE AHARONI



In phrase-based translation, score partial hypotheses by:





 $\mathbf{e}_{\mathsf{best}} = \mathsf{argmax}_{\mathbf{e}} \prod \phi(\bar{f}_i | \bar{e}_i) \ d(start_i - end_{i-1} - 1) \ p_{\text{LM}}(\mathbf{e})$

- In phrase-based translation, score **partial hypotheses** by: $\mathbf{e}_{\mathsf{best}} = \mathsf{argmax}_{\mathbf{e}} \prod \phi($ i=1
- Translation model, reordering model, language model

$$\bar{f}_i|\bar{e}_i) d(start_i - end_{i-1} - 1) p_{\text{LM}}(\mathbf{e})$$

- In phrase-based translation, score **partial hypotheses** by: $\mathbf{e}_{\mathsf{best}} = \mathsf{argmax}_{\mathbf{e}} \prod \phi($ i=1
- Translation model, reordering model, language model
- A small problem...

$$\bar{f}_i|\bar{e}_i) d(start_i - end_{i-1} - 1) p_{\text{LM}}(\mathbf{e})$$

ROEE AHARONI

• Large search space many options to choose from

\supset
\supset
\supset
\supset
it is
he will
it goes
he goe

- Large search space many options to choose from
- Use words? Phrases? Larger phrases?

- Large search space many options to choose from
- Use words? Phrases? Larger phrases?
- 2727 relevant phrase pairs for this sentence in Europarl alone

- Large search space many options to choose from
- Use words? Phrases? Larger phrases?

- 2727 relevant phrase pairs for this sentence in Europarl alone
- Exponential explosion

2020

• Start with empty hypothesis

- Start with empty hypothesis
- Fetch all initial hypotheses and score them

- Start with empty hypothesis
- Fetch all initial hypotheses and score them
- Maintain a coverage vector for each hypothesis

- Start with empty hypothesis
- Fetch all initial hypotheses and score them
- Maintain a coverage vector for each hypothesis
- Iterate: expand each hypothesis until full cover

- Start with empty hypothesis
- Fetch all initial hypotheses and score them
- Maintain a coverage vector for each hypothesis
- Iterate: expand each hypothesis until full cover
- Backtrack to get best candidate

2020

• The proposed approach is still exponential.

2020

- The proposed approach is still exponential.
- How can we improve?

2020

- The proposed approach is still exponential.
- How can we improve?
 - Recombination merge hypotheses with similar traits

- The proposed approach is still exponential.
- How can we improve?
 - Recombination merge hypotheses with similar traits
 - Depends on the modelled context

- The proposed approach is still exponential.
- How can we improve?
 - Recombination merge hypotheses with similar traits
 - Depends on the modelled context
 - Still not enough, not controllable

 Idea: throw away bad hypothesis early

- Idea: throw away bad hypothesis early
- Put comparable hypotheses into stacks (same coverage)

no word translated

one word translated

two words translated

three words translated

- Idea: throw away bad hypothesis early
- Put comparable hypotheses into stacks (same coverage)
- Limit the number of hypotheses in each stack throw away the worst hypothesis when stack is full

1	П	Т	Т	Т

no word translated

one word translated

two words translated

three words translated

- Idea: throw away bad hypothesis early
- Put comparable hypotheses into stacks (same coverage)
- Limit the number of hypotheses in each stack throw away the worst hypothesis when stack is full
- Beam search is similar, but without stacks

_	_	_	_
μ	-	-	-
			_

 \Box

no word translated

one word translated

two words translated

three words translated

- for all hypotheses in stack do for all translation options do if applicable then create new hypothesis place in stack recombine with existing hypothesis **if** possible prune stack **if** too big

- 3: 5:
- 1: place empty hypothesis into stack 0 2: for all stacks 0...n - 1 do 4: 6: 7: 8: 9:
- end if 10:
- end for 11:
- end for 12:
- 13: **end for**

Stack Decoding Algorithm

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

- end if 10:
- end for 11:
- end for 12:
- 13: **end for**

 $O(\max \text{ stack size} \times \text{ translation options} \times \text{ sentence length})$

Stack Decoding Algorithm

Summary

Summary

• Additional lexical models
- Additional lexical models
 - Position based (absolute Model 2, relative Model 4)

- Additional lexical models
 - Position based (absolute Model 2, relative Model 4)
 - Fertility (Model 3)

- Additional lexical models
 - Position based (absolute Model 2, relative Model 4)
 - Fertility (Model 3)
- Phrase based models

- Additional lexical models
 - Position based (absolute Model 2, relative Model 4)
 - Fertility (Model 3)
- Phrase based models
- Decoding

- Additional lexical models
 - Position based (absolute Model 2, relative Model 4)
 - Fertility (Model 3)
- Phrase based models
- Decoding
 - Recombination, Pruning (Stack decoding/Beam search)

