

# A Hybrid Morpheme-Word Representation for Machine Translation of Morphologically Rich Languages

Minh-Thang Luong, Preslav Nakov & Min-Yen Kan  
EMNLP 2010, Massachusetts, USA





**epäjärjestelmällisyydellistytämättömyydellänsäkäänköhän**



**This is a Finnish word!!!**

epä+ järjestelmä+ llisyy+ dellisty+ ttä+ mättö+ myy+ dellä+ nsä+ kään+ kö+ hän



**epäjärjestelmällisydellistytämättömydellänsäkäänköhän**



**This is a Finnish word!!!**

epä+ järjestelmä+ llisy+ dellisty+ ttä+ mättö+ myy+ dellä+ nsä+ kään+ kö+ hän



**system**



**epäjärjestelmällisyydellistytämättömyydellänsäkäänköhän**

**This is a Finnish word!!!**

epä+ järjestelmä+ llisyy+ dellisty+ ttä+ mättö+ myy+ dellä+ nsä+ kään+ kö+ hän

**unsystem**



**epäjärjestelmällisyydellistytämättömyydellänsäkäänköhän**

**This is a Finnish word!!!**

epä+ järjestelmä+ llisy+ dellisty+ ttä+ mättö+ myy+ dellä+ nsä+ kään+ kö+ hän

**unsystematic**



**epäjärjestelmällisyystyttämättömyydellänsäkäänköhän**

**This is a Finnish word!!!**

epä+ järjestelmä+ llisyy+ dellisty+ ttä+ mättö+ myy+ dellä+ nsä+ kään+ kö+ hän

I wonder if it's not with his act of not having made something be seen as unsystematic

- **Morphologically rich languages (Arabic, Basque, Turkish, Russian, Hungarian, etc. )**
- Extensive use of affixes

Morphologically rich languages are hard for MT  
 → Analysis at the morpheme level is needed.



## Morphological Analysis Helps ?

– Translation into morphologically *poor* languages

➤ **Morpheme** representation alleviates data sparseness

*Arabic* → *English* (Lee, 2004)    *Czech* → *English* (Goldwater & McClosky, 2005)  
*Finnish* → *English* (Yang & Kirchhoff, 2006)

➤ For large corpora, **word** representation captures context better

*Arabic* → *English* (Sadat & Habash, 2006)  
*Finnish* → *English* (de Gispert et al., 2009)

*Our approach*: the basic unit of translation is the *morpheme*, but *word* boundaries are respected at all MT stages.



## Morphological Analysis Helps ?

– Translation into morphologically *rich* languages

---

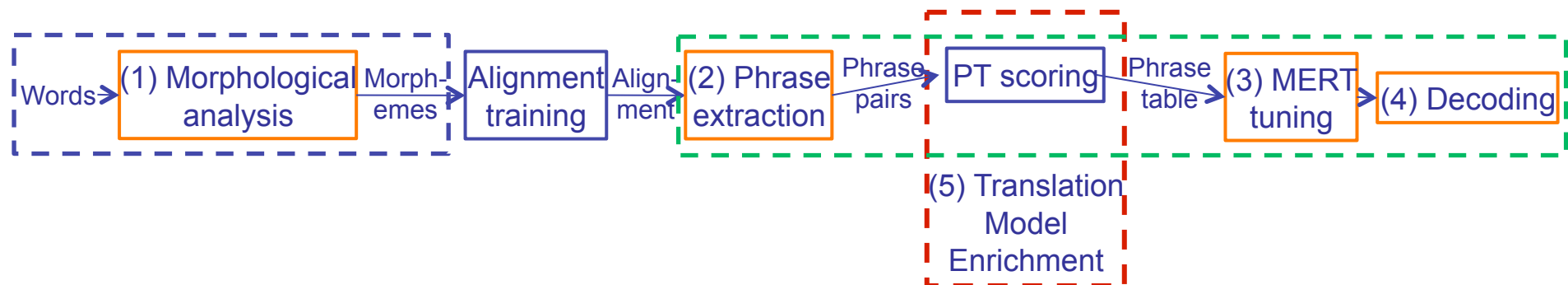
- **A challenging translation direction**
- **Recent interest in:**
  - English → Arabic (Badr et al., 2008)
  - English → Turkish (Oflazer and El-Kahlout, 2007)  
→ enhance the performance for small bi-texts only
  
  - English → Greek (Avramidis and Koehn, 2008),
  - English → Russian (Toutanova et al., 2008)  
→ rely heavily on language-specific tools

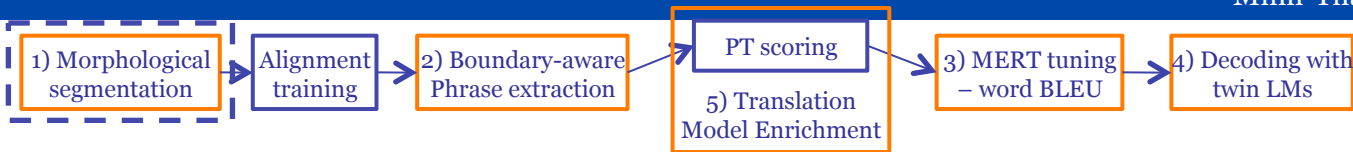
We want an unsupervised approach that works for large training bi-texts.



## Methodology

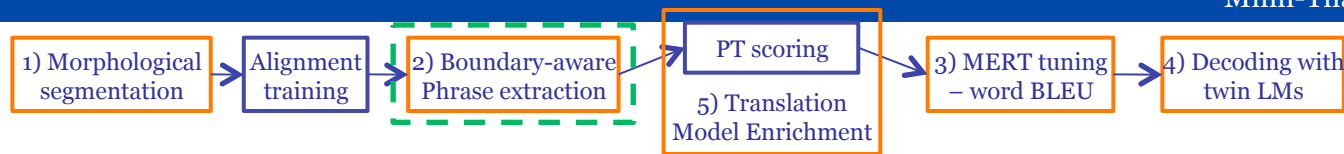
- **Morphological Analysis** – *Unsupervised*
- **Morphological Enhancements** – *Respect word boundaries*
- **Translation Model Enrichment** – *Merge phrase tables (PT)*



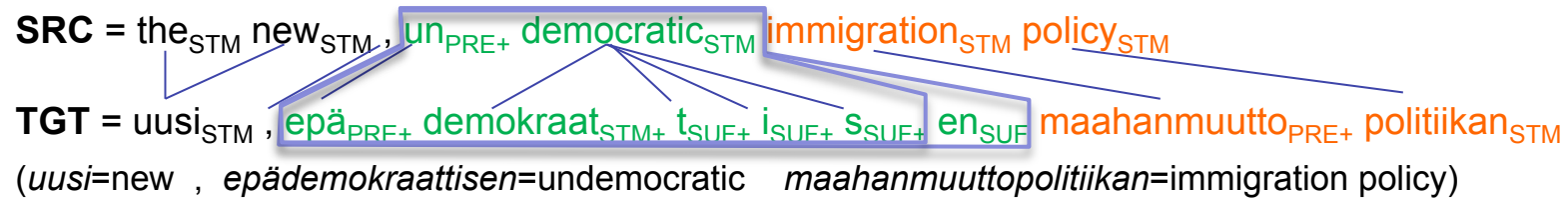


## Morphological Analysis

- Use *Morfessor* (Creutz and Lagus, 2007) - **unsupervised** morphological analyzer
- Segments words → morphemes (PRE, STM, SUF)  
un/PRE+ care/STM+ fu/SUF+ ly/SUF
- “+” sign used to enforce **word boundary constraints** later



## Word Boundary-aware Phrase Extraction



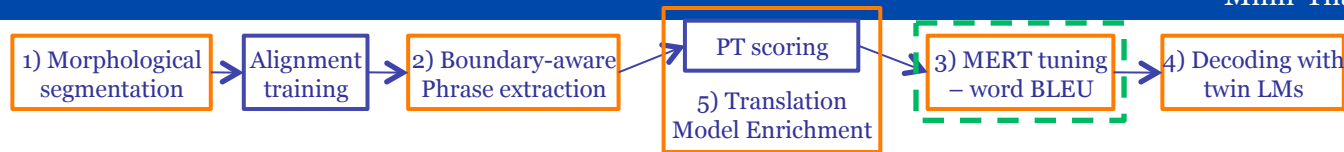
- **Typical SMT:** maximum phrase length  $n=7$  words
- **Problem:** morpheme phrases of length  $n$ 
  - can span less than  $n$  words
  - may only partially span words

This problem is severe for morphologically rich languages.

- **Solution:** morpheme phrases
  - span up to  $n$  words
  - fully span words



TRAINING



## Morpheme MERT Optimizing Word BLEU

### ➤ Why ?

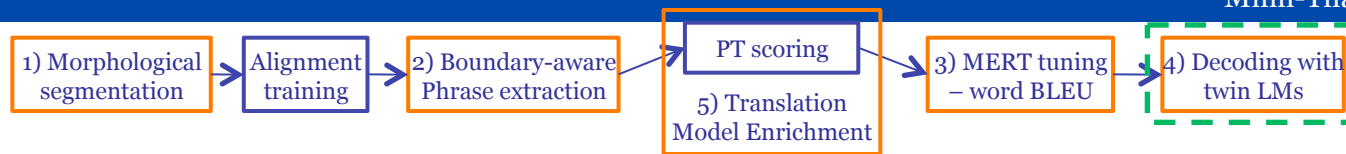
- BLEU's brevity penalty is influenced by sentence length
- The same # of words span different # of morphemes
  - a 6-morpheme Finnish word:  $ep\ddot{a}_{PRE+} demokraat_{STM+} t_{SUF+} i_{SUF+} s_{SUF+} en_{SUF}$
- Suboptimal weight for the SMT word penalty feature

### ➤ **Solution:** optimize on word BLEU.

### ➤ Each MERT iteration:

- Decode at the morpheme level
- Convert morpheme translation → word sequence
- Compute word BLEU
- Convert back to morphemes

## TUNING



## Decoding with Twin Language Models

- **Morpheme language model (LM)**
  - *Pros*: alleviates data sparseness
  - *Cons*: phrases span fewer words
- **Introduce a second LM at the word level**
  - *Log-linear model*: add a separate feature
  - *Moses decoder*: add *word-level* “view” on the *morpheme-level* hypotheses

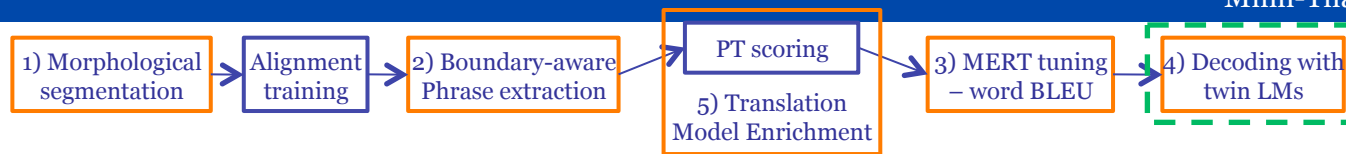
*Previous hypotheses*
*Current hypothesis*

---

uusi<sub>STM</sub> , epä<sub>PRE+</sub> demokraat<sub>STM+</sub> t<sub>SUF+</sub> i<sub>SUF+</sub> s<sub>SUF+</sub> en<sub>SUF</sub> maahanmuutto<sub>PRE+</sub> politiikan<sub>STM</sub>

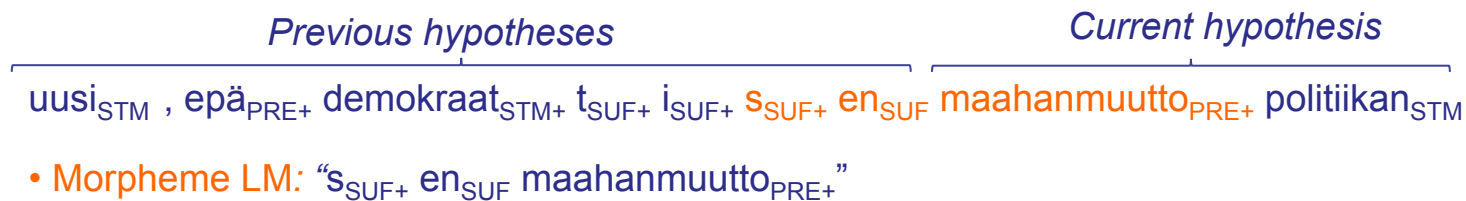
- **Morpheme LM:**

### D E C O D I N G

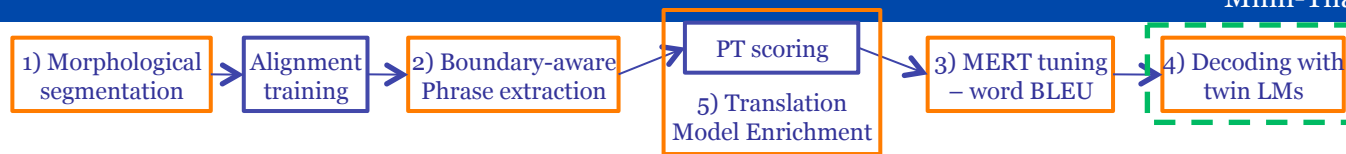


## Decoding with Twin Language Models

- **Morpheme language model (LM)**
  - *Pros*: alleviates data sparseness
  - *Cons*: phrases span fewer words
- **Introduce a second LM at the word level**
  - *Log-linear model*: add a separate feature
  - *Moses decoder*: add *word-level* “view” on the *morpheme-level* hypotheses

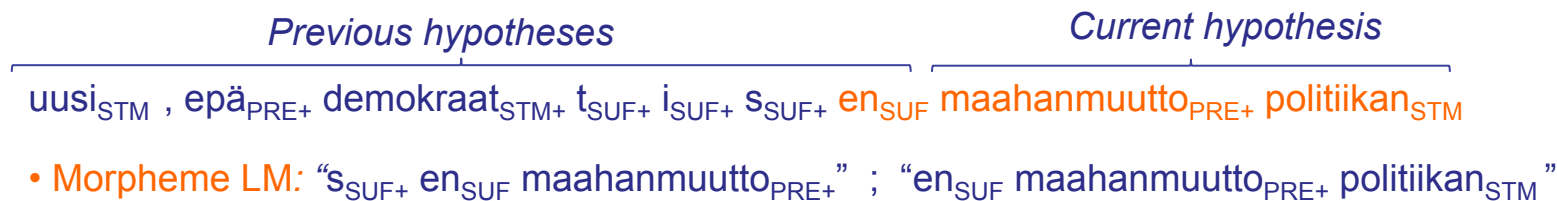


### DECODING

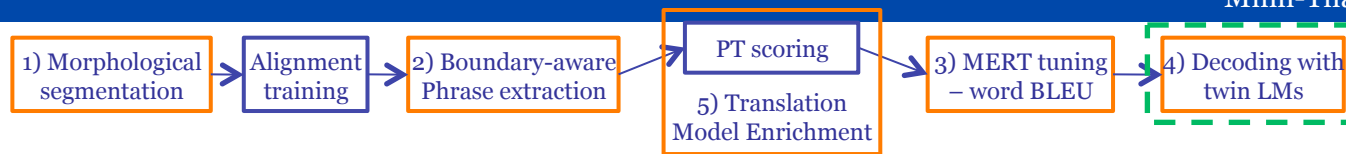


## Decoding with Twin Language Models

- **Morpheme language model (LM)**
  - *Pros*: alleviates data sparseness
  - *Cons*: phrases span fewer words
- **Introduce a second LM at the word level**
  - *Log-linear model*: add a separate feature
  - *Moses decoder*: add *word-level* “view” on the *morpheme-level* hypotheses

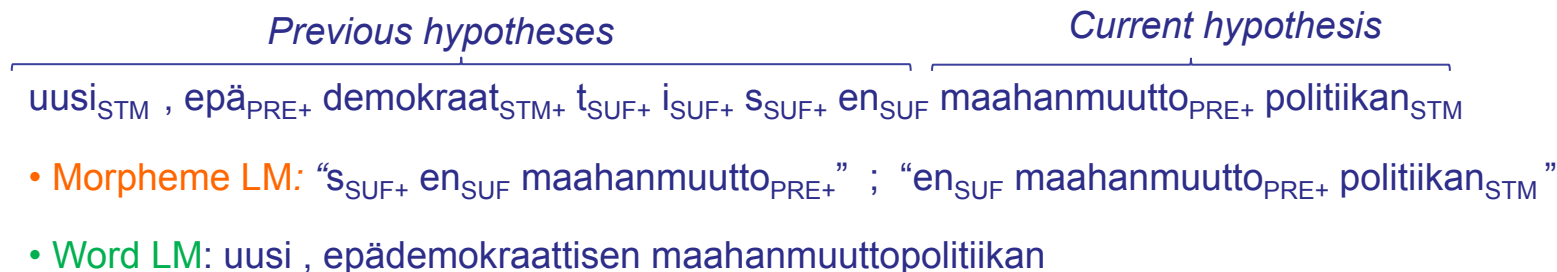


### DECODING



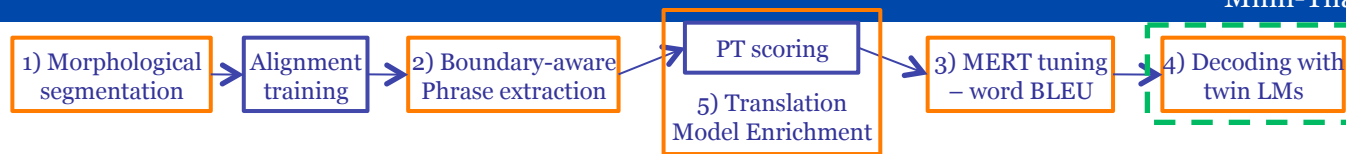
## Decoding with Twin Language Models

- **Morpheme language model (LM)**
  - *Pros*: alleviates data sparseness
  - *Cons*: phrases span fewer words
- **Introduce a second LM at the word level**
  - *Log-linear model*: add a separate feature
  - *Moses decoder*: add *word-level* “view” on the *morpheme-level* hypotheses



### DECODING





## Decoding with Twin Language Models

- **Morpheme language model (LM)**
  - *Pros*: alleviates data sparseness
  - *Cons*: phrases span fewer words
- **Introduce a second LM at the **word level****
  - *Log-linear model*: add a separate feature
  - *Moses decoder*: add *word-level* “view” on the *morpheme-level* hypotheses

*Previous hypotheses*
*Current hypothesis*

$uusi_{STM}, epä_{PRE+}, demokraat_{STM+}, t_{SUF+}, i_{SUF+}, s_{SUF+}, en_{SUF}, maahanmuutto_{PRE+}, politiikan_{STM}$

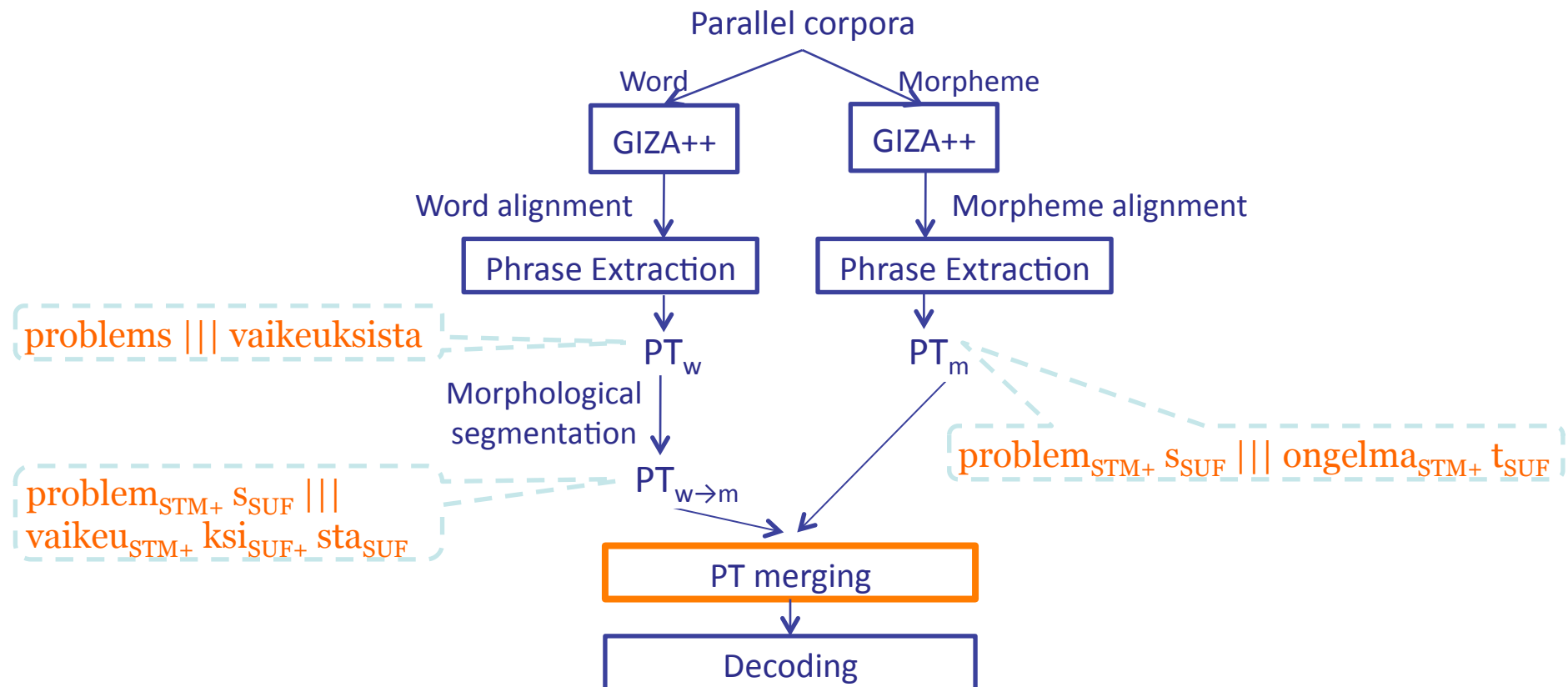
- **Morpheme LM**: “ $s_{SUF+} en_{SUF} maahanmuutto_{PRE+}$ ” ; “ $en_{SUF} maahanmuutto_{PRE+} politiikan_{STM}$ ”

This is (1) different from scoring with two word-level LMs &  
 (2) superior to n-best rescoring.

D E C O D I N G

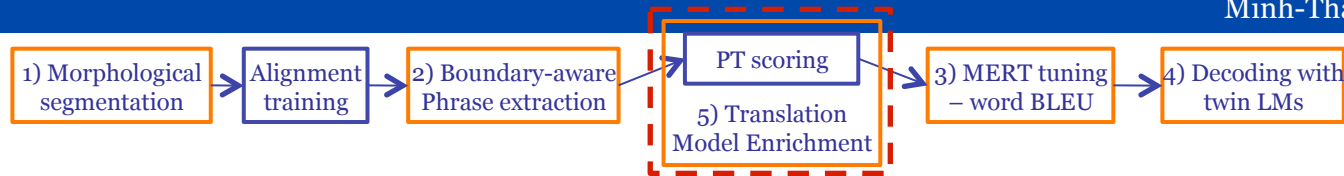


# Building Twin Translation Models



From the same source, we generate two translation models.

## ENRICHING



## Phrase Table (PT) Merging

Phrase translation probabilities

Lexicalized translation probabilities

problem <sub>STM+</sub> s <sub>SUF</sub>     vaikeu <sub>STM+</sub> ksi <sub>SUF+</sub> sta <sub>SUF</sub>	0.07 0.11 0.01 0.01 2.7	} Phrase penalty
problem <sub>STM+</sub> s <sub>SUF</sub>     ongelma <sub>STM+</sub> t <sub>SUF</sub>	0.37 0.60 0.11 0.14 2.7	

- **Add-feature methods** e.g., (Chen et al., 2009)

problem <sub>STM+</sub> s <sub>SUF</sub>     vaikeu <sub>STM+</sub> ksi <sub>SUF+</sub> sta <sub>SUF</sub>	0.07 0.11 0.01 0.01 2.7	2.71
problem <sub>STM+</sub> s <sub>SUF</sub>     ongelma <sub>STM+</sub> t <sub>SUF</sub>	0.37 0.60 0.11 0.14 2.7	1 2.7

➔ **Heuristic-driven**

in the 1<sup>st</sup> PT

in the 2<sup>nd</sup> PT

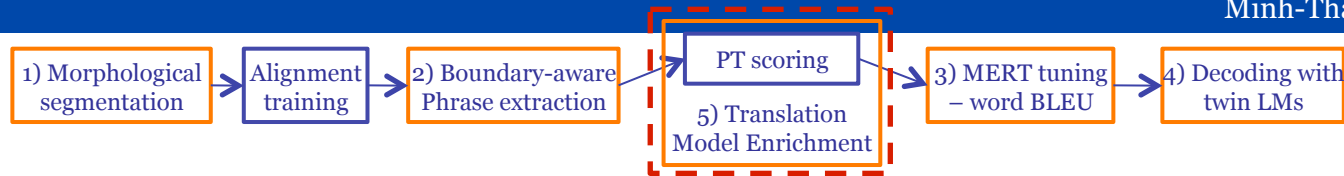
- **Interpolation-based methods** e.g., (Wu & Wang, 2007)

- Linear interpolation of phrase and lexicalized translation probabilities

➔ For two PTs originating from different sources

We take into account the fact that our twin translation models are of equal quality.

ENRICHING



## Our Method: Phrase Translation Probabilities

- Preserve the normalized ML estimations (Koehn et al., 2003)

$$\phi(\bar{f}|\bar{e}) = \frac{\#(\bar{f}, \bar{e})}{\sum_{\bar{f}} \#(\bar{f}, \bar{e})}$$

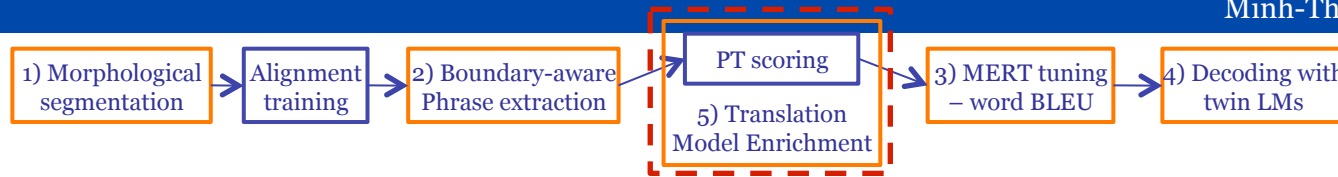
The number of times the pair  $(\bar{f}, \bar{e})$  was extracted from the training dataset

- Use the raw counts of both models to compute

$$\phi(\bar{f}, \bar{e}) = \frac{\#_m(\bar{f}, \bar{e}) + \#_{w \rightarrow m}(\bar{f}, \bar{e})}{\sum_{\bar{f}} \#_m(\bar{f}, \bar{e}) + \sum_{\bar{f}} \#_{w \rightarrow m}(\bar{f}, \bar{e})}$$

$\text{PT}_m$ 
 $\text{PT}_{w \rightarrow m}$

### ENRICHING



# Lexicalized Translation Probabilities

- Use linear interpolation
- What if a phrase pair belongs to one PT only?

$\begin{matrix} \text{problem}_{\text{STM}+} \text{ s}_{\text{SUF}} \parallel\parallel \text{vaikeu}_{\text{STM}+} \text{ ksi}_{\text{SUF}+} \text{ sta} \\ \text{problem}_{\text{STM}+} \text{ s}_{\text{SUF}} \parallel\parallel \text{ongelma}_{\text{STM}+} \text{ t}_{\text{SUF}} \end{matrix}$ <p style="text-align: center;"><math>\text{PT}_m</math></p>	$\begin{matrix} \text{problem}_{\text{STM}+} \text{ s}_{\text{SUF}} \parallel\parallel \text{vaikeu}_{\text{STM}+} \text{ ksi}_{\text{SUF}+} \text{ sta} \\ \text{problem}_{\text{STM}+} \text{ s}_{\text{SUF}} \parallel\parallel \text{ongelma}_{\text{STM}+} \text{ sta}_{\text{SUF}} \end{matrix}$ <p style="text-align: center;"><math>\text{PT}_{w \rightarrow m}</math></p>
--	--

- Previous methods: interpolate with 0
  - Might cause some good phrases to be penalized
- Our method: induce all scores before interpolation
  - Use the lexical model of one PT to score phrase pairs for the other one

$$\text{lex}(\bar{f}|\bar{e}) = \alpha \times \boxed{\text{lex}_m(\bar{f}_m|\bar{e}_m)} + (1 - \alpha) \times \boxed{\text{lex}_w(\bar{f}_w|\bar{e}_w)}$$

$(\text{problem}_{\text{STM}+} \text{ s}_{\text{SUF}} | \text{ongelma}_{\text{STM}+} \text{ t}_{\text{SUF}})$ 

Phrase Lexical Model (PT<sub>m</sub>)

$(\text{problems} | \text{ongelmat})$ 

Word Lexical Model (PT<sub>w</sub>)

ENRICHING



## Dataset & Settings

---

- **Dataset**
  - Past shared task WPT05 (en/fi)
  - 714K sentence pairs
  - Split into T1, T2, T3, and T4 of sizes 40K, 80K, 160K, and 320K
  
- **Standard phrase-based SMT settings:**
  - Moses
  - IBM Model 4
  - Case insensitive BLEU

---

## EXPERIMENTS

---

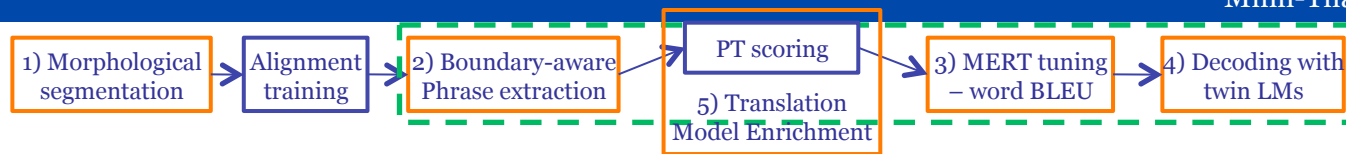
## SMT Baseline Systems

	w-system		m-system	
	BLEU	m-BLEU	BLEU	m-BLEU
<b>T1</b>	11.56	45.57	11.07	49.15
<b>T2</b>	12.95	48.63	12.68	53.78
<b>T3</b>	13.64	50.30	13.32	54.40
<b>T4</b>	14.20	50.85	13.57	54.70
<b>Full</b>	14.58	53.05	14.08	55.26

- **w-system**: word level
- **m-system**: morpheme level
- **m-BLEU**: morpheme version of BLEU

Either the *m-system* does not perform as well as the *w-system* or BLEU is not capable of measuring morpheme improvements.

### EXPERIMENTS



## Morphological Enhancements: Individual

System	T1 (40K)	Full (714K)
w-system	11.56	14.58
m-system	11.07	14.08
m+phr	11.44 <sup>+0.37</sup>	14.43 <sup>+0.35</sup>
m+tune	11.73 <sup>+0.66</sup>	14.55 <sup>+0.47</sup>
m+lm	11.58 <sup>+0.51</sup>	14.53 <sup>+0.45</sup>

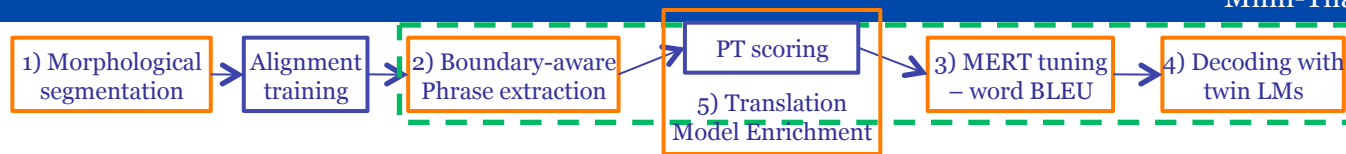
- **phr**: boundary-aware phrase extraction
- **tune**: MERT tuning for word BLEU
- **lm**: decoding with twin LMs



The individual enhancements yield improvements for both small and large corpora.

### EXPERIMENTS





## Morphological Enhancements: Combined

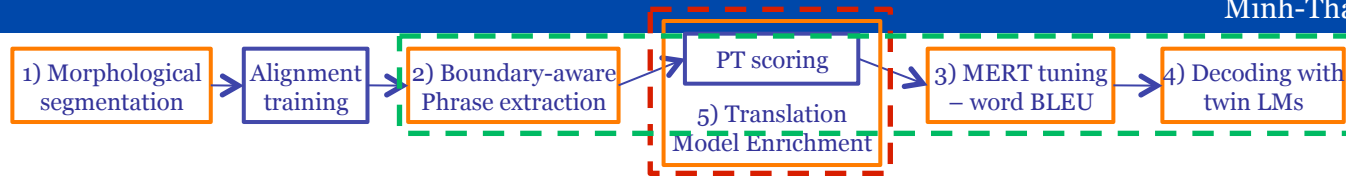
System	T1 (40K)	Full (714K)
w-system	11.56	14.58
m-system	11.07	14.08
m+phr+lm	11.77 <sup>+0.70</sup>	<b>14.58<sup>+0.50</sup></b>
m+phr+lm+tune	<b>11.90<sup>+0.83</sup></b>	14.39 <sup>+0.31</sup>

- **phr**: boundary-aware phrase extraction
- **tune**: MERT tuning for word BLEU
- **lm**: decoding with twin LMs



The morphological enhancements are on par with the *w-system* and yield sizable improvements over the *m-system*.

### EXPERIMENTS



## Translation Model Enrichment

Merging methods	Full (714K)
m-system	14.08
w-system	14.58
add-1	14.25 <sup>+0.17</sup>
add-2	13.89 <sup>-0.19</sup>
interpolation	14.63 <sup>+0.55</sup>
<b>ourMethod</b>	<b>14.82<sup>+0.74</sup></b>

- **add-1:** *one* extra feature
- **add-2:** *two* extra features
- **Interpolation:** linear interpolation
- **ourMethod**

Our method outperforms the *w-system* baseline.

### EXPERIMENTS



## Results Significance

---

	BLEU
m-system	14.08
w-system	14.58
<b>ourSystem</b>	<b>14.82<sup>+0.74</sup></b>

- **Absolute improvement of 0.74 BLEU over the *m-system*,**
  - non-trivial relative improvement of 5.6%
- **Outperformed the *w-system* by 0.24 points (1.56% relative)**
- **Statistically significant with  $p < 0.01$  (Collins' sign test)**

---

### ANALYSIS

---



## Translation Proximity Match

### ➤ Automatically extract phrase triples

(economic and social, taloudellis**ia** ja sosiaalis**ia**, taloudellisten ja sosiaalisten)

src    out    ref

- high character-level similarity between *out* & *ref*
- longest common subsequence ratio
- **16,262 triples**: 6,758 match the references exactly
- the remaining triples were close wordforms

*Hypothesis*: our approach yields translations close to the reference wordforms, but it is unjustly penalized by BLEU.

### ANALYSIS

## Human Evaluation

	<i>our</i> vs. <i>m</i>		<i>our</i> vs. <i>w</i>		<i>w</i> vs. <i>m</i>	
<b>Judge 1</b>	25	18	19	12	21	19
<b>Judge 2</b>	24	16	19	15	25	14
<b>Judge 3</b>	27 <sup>†</sup>	12	17	11	27 <sup>†</sup>	15
<b>Judge 4</b>	25	20	26 <sup>†</sup>	12	22	22
<b>Total</b>	101 <sup>‡</sup>	66	81 <sup>‡</sup>	50	95 <sup>†</sup>	70

### ➤ 4 native Finnish speakers

- 50 random test sentences
- follow WMT'09 evaluation:
  - provided judges with the source sentence, its reference translation & outputs of (*m-system*, *w-system*, *ourSystem*) shown in random order
  - asked for three pairwise judgments

The judges consistently preferred: (1) *ourSystem* to the *m-system*, (2) *ourSystem* to the *w-system*, (3) *w-system* to the *m-system*.

## ANALYSIS

## Sample Translations

src: we were very constructive and we negotiated until the last minute of these talks in the hague .

ref: olimme erittäin **rakentavia** ja neuvottelimme haagissa **viime hetkeen saakka** .

Match reference

our: olemme olleet hyvin **rakentavia** ja olemme neuvotelleet **viime hetkeen saakka** naiden neuvottelujen haagissa .

w : ol **Wrong case** in **rakentavia** ja olemme neuvotelleet viime *tippaan niin* naiden neuv **Confusing meaning**

m : olimme erittäin **rakentavan** ja neuvottelimme **viime hetkeen saakka** naiden neuvotteluiden haagissa .

Rank: **our** > **m** ≥ **w**

src: it would be a very dangerous situation if the europeans were to become logistically reliant on russia .

ref: olisi **erittäin** vaarallinen tilanne , jos **eurooppalaiset** tulisivat **logistisesti** riippuvaisia venäjältä .

Match reference

our: olisi **erittäin** vaarallinen tilanne , jos **eurooppalaiset** tulee **logistisesti** riippuvaisia venäjän .

w : se olisi **erittäin** vaaral **Wrong case** s **eurooppalaisten** tulisi *logistically* riippuvaisia **OOV** .

m : se olisi *hyvin* vaarallinen tilanne , jos **eurooppalaiset** *haluavat* tulla **logistisesti** riippuvaisia venäjän .

Rank: **our** > **w** ≥ **m**

Change the meaning

*ourSystem* consistently outperforms the *m-system* & *w-system*, and seems to blend well translations from both baselines.



## Conclusion

---

- **Our approach:**
  - The basic unit of translation is the *morpheme*
  - But *word* boundaries are respected at all MT stages
  - Unsupervised method that works for large training bi-texts
- **Future work:**
  - Extend the morpheme-level framework
  - Incorporate morphological analysis directly into the translation process

**Thank you!**