# 50.040
# Natural Language Processing

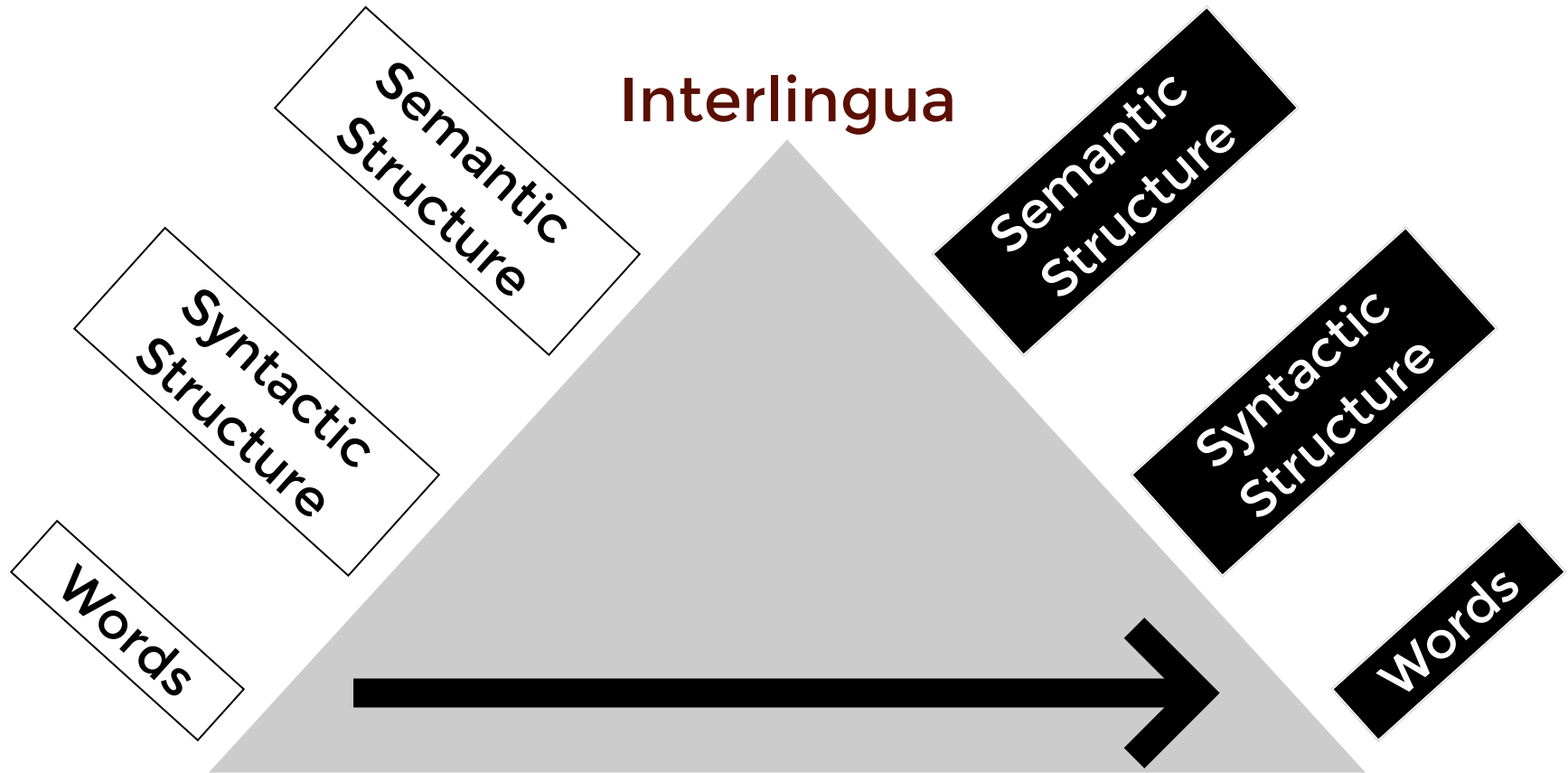Lu, Wei

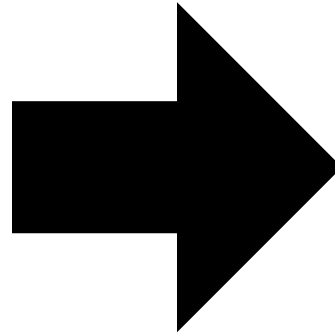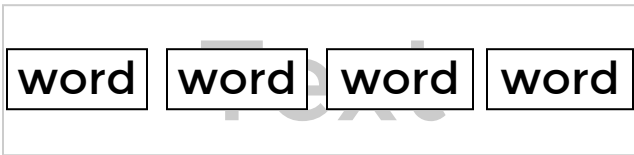SINGAPORE UNIVERSITY OF
TECHNOLOGY AND DESIGN

# Tasks in NLP

# Machine Translation

Semantic Structure

Interlingua

Semantic Structure

Syntactic Structure
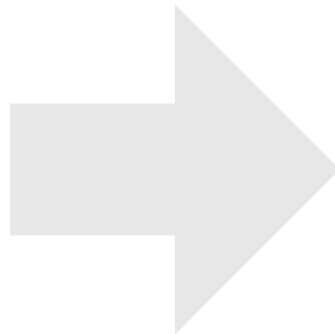
Syntactic Structure

Words

Words

Text-to-text Problem
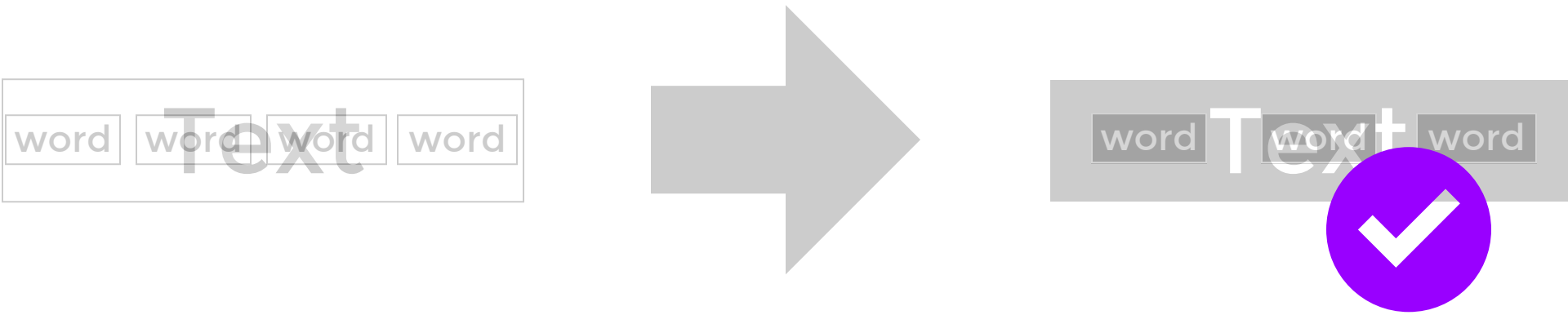
# Machine Translation

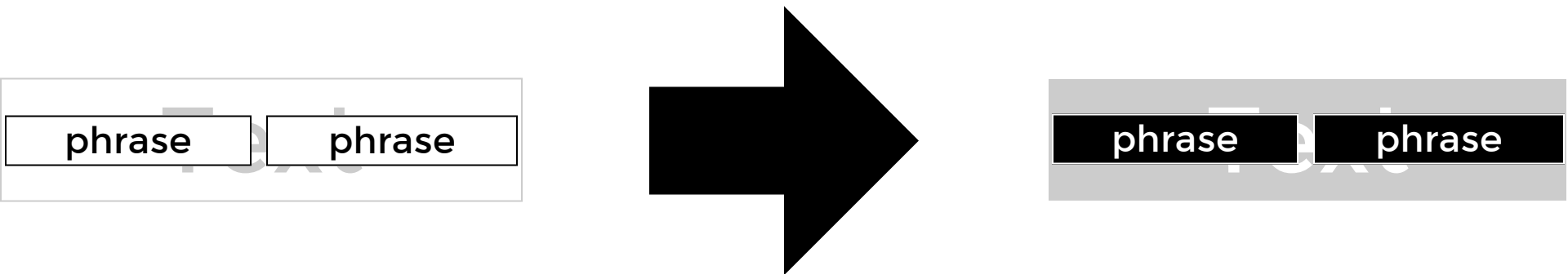

Word-based Translation

Phrase-based Translation

# Machine Translation



Word-based Translation



Phrase-based Translation

# Phrase-based Translation

The Alignment Template Approach to Statistical Machine Translation

Franz Josef Och*
Google

Hermann Ney†
RWTH Aachen

A phrase-based statistical machine translation approach — the alignment template approach — is described. This translation approach allows for general many-to-many relations between words. Thereby, the context of words is taken into account in the translation model, and local changes in word order from source to target language can be learned explicitly. The model is described using a log-linear modeling approach, which is a generalization of the often used source–channel approach. Thereby, the model is easier to extend than classical statistical machine translation systems. We describe in detail the process for learning phrasal translations, the feature functions used, and the search algorithm. The evaluation of this approach is performed on three different tasks. For the German–English VERBMOBIL task, we analyze the effect of various system components. On the French–English Canadian HANSARDS task, the alignment template system obtains significantly better results than a single-word-based translation model. In the Chinese–English 2002 National Institute of Standards and Technology (NIST) machine translation evaluation it yields statistically significantly better NIST scores than all competing research and commercial translation systems.

## 1. Introduction

Machine translation (MT) is a hard problem, because natural languages are highly complex, many words have various meanings and different possible translations, sentences might have various readings, and the relationships between linguistic entities are often vague. In addition, it is sometimes necessary to take world knowledge into account. The number of relevant dependencies is much too large and those dependencies are too complex to take them all into account in a machine translation system. Given these boundary conditions, a machine translation system has to make decisions (produce translations) given incomplete knowledge. In such a case, a principled approach to solving that problem is to use the concepts of statistical decision theory to try to make optimal decisions given incomplete knowledge. This is the goal of statistical machine translation.

The use of statistical techniques in machine translation has led to dramatic improvements in the quality of research systems in recent years. For example, the statistical approaches of the VERBMOBIL evaluations (Wahlster 2000) or the U.S. National

* 1600 Amphitheatre Parkway, Mountain View, CA 94043. E-mail: och@google.com.
† Lehrstuhl für Informatik VI, Computer Science Department, RWTH Aachen–University of Technology, Ahornstr. 55, 52056 Aachen, Germany. E-mail: ney@cs.rwth-aachen.de.

Submission received: 19 November 2002; Revised submission received: ...
publication: 1 June 2004

© 2004 Association for Computational Linguistics

Probably the first successful approach to phrase-based MT, but is complicated!

6

# Phrase-based Translation



Simpler phrase-based SMT, was one of the main approaches to MT in the last decade.

# Word Alignment

SUTD is the only university in the East .

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

# Phrase Alignment

SUTD is the only university in the East .

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

SUTD is the only university in the East .

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

# Phrase Lexicon

SUTD is the only university in the East .

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

SUTD is the only university in the East .

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

新加坡 科技 设计 大学　　⇒　　　　SUTD
唯一 的 一所 大学　　⇒　　the only university

# Word Alignment

$\mathbf{A} : p(\boldsymbol{e}|\boldsymbol{f})$

SUTD  is  the  only  university  in  the  East  .

French-English: one-many

新加坡
科技
设计
大学
是
东部
唯一
的
一所
大学
。

# Word Alignment

**B** : $p(\boldsymbol{f}|\boldsymbol{e})$

| | SUTD | is | the | only | university | in | the | East | . |
|---|---|---|---|---|---|---|---|---|---|
| 新加坡 | ■ | | | | | | | | |
| 科技 | ■ | | | | | | | | |
| 设计 | ■ | | | | | | | | |
| 大学 | ■ | | | | | | | | |
| 是 | | ■ | | | | | | | |
| 东部 | | | | | | | | ■ | |
| 唯一 | | | | ■ | | | | | |
| 的 | | | | | | ■ | | | |
| 一所 | | | | | ■ | | | | |
| 大学 | | | | | ■ | | | | |
| 。 | | | | | | | | | ■ |

French-English: many-one

# Phrase Alignment

$$\mathbf{M} = h(\mathbf{A}, \mathbf{B})$$

some heuristics

SUTD · is · the · only · university · in · the · East · .

新加坡科技设计大学是东部唯一的一所大学。

French-English: many-many

# Phrase Alignment

$$\mathbf{M} = h(\mathbf{A}, \mathbf{B})$$

Some sample heuristics

Start with the intersection of $\mathbf{A}$ and $\mathbf{B}$

Incrementally add points from union of $\mathbf{A}$ and $\mathbf{B}$

First only add points to words which are not aligned

Give priority to points with neighboring points

# Phrase Lexicon

We need an algorithm to extract the phrase pairs

A phrase pair $(\bar{f}, \bar{e})$ is **consistent** if:

At least one word in $\bar{f}$ aligns with a word in $\bar{e}$

No words in $\bar{e}$ align to words outside $\bar{f}$

No words in $\bar{f}$ align to words outside $\bar{e}$

Extract all consistent phrase pairs from the training set

# Phrase-based Translation

**LM : Language Model**

How well the translated sentence reads

**TM : Phrase Translation Model**

How faithful the translation is to the original

**DM : Distortion Model**

How much efforts on "moving the eyes" in translation is required

# Phrase Translation Model

We need to score each extracted phrase pair

A phrase pair $(\bar{f}, \bar{e})$ can be scored as:

$$t(\bar{f}|\bar{e}) = \frac{\text{count}(\bar{f}, \bar{e})}{\text{count}(\bar{e})}$$

This is the "phrase translation probability", which says something about the quality of this translation pair

# Distortion Model

A distortion parameter that says the **significance** of the amount of distortion (usually **negative**).

Position of the **first word** in the French phrase that corresponds to the **current translated** English phrase

$$\eta \times |\text{pl}(p_k) + 1 - \text{pf}(p_{k+1})|$$

Position of the **last** word in the French phrase that corresponds to the **previous translated** English phrase

$$e = \underbrace{p_1 p_2 \ldots p_{L-1} p_L}_{L \text{ phrases}}$$

This quantity measures how much efforts on "moving the eyes" is needed when the translator is doing the translation.

# Distortion Model

Typically, there is a limit to the maximal distortion we can tolerate.

$$|\mathrm{pl}(p_k) + 1 - \mathrm{pf}(p_{k+1})| \leq d$$

Large distortion can lead to poor translation quality in practice

# Phrase-based Translation

$p_1 = (1, 4, \text{SUTD})$

**SUTD**

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$\underbrace{\log q(\text{SUTD}|\langle \textbf{START} \rangle, \langle \textbf{START} \rangle)}_{\text{Language Model}}$$

$$+ \quad \underbrace{\log t(\text{新加坡 科技 设计 大学}|\text{SUTD})}_{\text{Phrase translation model}}$$

$$+ \quad \underbrace{\eta \times 0}_{\text{Distortion model}}$$

# Phrase-based Translation

$p_2 = (5, 5, \text{is})$

**SUTD** **is**

新加坡 科技 设计 大学 | 是 | 东部 唯一 的 一所 大学 。

$$\underbrace{\log q(\text{is}|\langle \textbf{START} \rangle, \text{SUTD})}_{\text{Language model}}$$

$$+ \quad \underbrace{\log t(是|\text{is})}_{\text{Phrase translation model}}$$

$$+ \quad \underbrace{\eta \times 0}_{\text{Distortion model}}$$

# Phrase-based Translation

$$p_3 = (7, 8, \text{the only})$$

SUTD is the only

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$\underbrace{\log q(\text{the}|\text{SUTD}, \text{is}) + \log q(\text{only}|\text{is}, \text{the})}_{\text{Language model}}$$

$$+ \quad \underbrace{\log t(\text{唯一 的}|\text{the only})}_{\text{Phrase translation model}}$$

$$+ \quad \underbrace{\eta \times 1}_{\text{Distortion model}}$$

# Phrase-based Translation

$$p_4 = (9, 10, \text{university})$$

SUTD is the only university

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$\underbrace{\log q(\text{university}|\text{the}, \text{only})}_{\text{Language model}}$$

$$+ \quad \underbrace{\log t(\text{一所 大学}|\text{university})}_{\text{Phrase translation model}}$$

$$+ \quad \underbrace{\eta \times 0}_{\text{Distortion model}}$$

# Phrase-based Translation

$$p_5 = (6, 6, \text{in the East})$$

SUTD is the only university in the East

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$\underbrace{\log q(\text{in}|\text{only}, \text{university}) + \log q(\text{the}|\text{university}, \text{in}) + \log q(\text{East}|\text{in}, \text{the})}_{\text{Language model}}$$

$$+ \quad \underbrace{\log t(\text{东部}|\text{in the East})}_{\text{Phrase translation model}}$$

$$+ \quad \underbrace{\eta \times 5}_{\text{Distortion model}}$$

# Phrase-based Translation

$$p_6 = (11, 11, .)$$

SUTD is the only university in the East .

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$\underbrace{\log q(.|\text{the}, \text{East})}_{\text{Language model}}$$

$$+ \quad \underbrace{\log t(。 \;|.)}_{\text{Phrase translation model}}$$

$$+ \quad \underbrace{\eta \times 4}_{\text{Distortion model}}$$

# Phrase-based Translation

SUTD is the only university in the East .

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$\text{score}(\boldsymbol{e}) = \underbrace{\log q(\boldsymbol{e})}_{\text{Language model}}$$

$$+ \underbrace{\sum_{k=1}^{L} \log t(p_k)}_{\text{Phrase translation model}}$$

$$+ \underbrace{\sum_{k=1}^{L-1} \eta \times |\text{pl}(p_k) + 1 - \text{pf}(p_{k+1})|}_{\text{Distortion model}}$$

# Decoding

We know how to score a translation derivation, but how do we search for the **most optimal** derivation?

The position of the last French word in the previous French phrase translated

The score of the partial derivation so far

A state $s = (e_1, e_2, b, r, \alpha)$

The last two English words in the previous translated English phrase

A bit string indicating which words in French are (not yet) translated.

# States

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$(\langle \mathbf{START} \rangle, \langle \mathbf{START} \rangle, 00000000000, 0, 0)$$

Initial State

# States

$p_1 = (1, 4, \mathrm{SUTD})$

SUTD

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$(\langle \mathbf{START} \rangle, \mathrm{SUTD}, 11110000000, 4, 3.7)$$

# States

$p_2 = (5, 5, \text{is})$

SUTD is

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$(\text{SUTD}, \text{is}, 11111000000, 5, 8.9)$$

# States

$$p_3 = (7, 8, \text{the only})$$

SUTD is the only

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$(\text{the}, \text{only}, 11111011000, 8, -0.9)$$

# States

$$p_4 = (9, 10, \text{university})$$

SUTD is the only university

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$(\text{only}, \text{university}, 11111011110, 10, 2.2)$$

# States

$$p_5 = (6, 6, \text{in the East})$$

SUTD is the only university in the East

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$(\text{the}, \text{East}, 11111111110, 6, 7.1)$$

# States

$$p_6 = (11, 11, .)$$

SUTD is the only university in the East .

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

$$(\text{East}, ., 11111111111, 11, 5.0)$$

One Final State

# States

SUTD is the only university in the East .

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

**It is similar to a transition-based parser!**

$$(East, ., 11111111111, 11, 5.0)$$

One Final State

# State Transition

$p_2 = (5, 5, \text{is})$

SUTD is

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

Each $p_k$ is essentially an action!

$$(\langle \mathbf{START} \rangle, \text{SUTD}, 11111000000, 5, 8.9)$$

$$(\ldots, 11111000010, \ldots) \quad \ldots \quad (\ldots, 11111000111, \ldots)$$

# State Transition

$p_2 = (5, 5, \text{is})$

SUTD is

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

---

## Allowable actions

1. The action $p_k$ must be compatible with $b$.

2. The distortion limit $d$ must be respected.

$(\dots, 11111000010, \dots)$    $\dots$    $(\dots, 11111000111, \dots)$

# State Transition

$p_2 = (5, 5, \text{is})$

SUTD is

新加坡 科技 设计 大学 是 东部 唯一 的 一所 大学 。

The greedy procedure may not yield the best translation.
Beam search is typically used in practice,
which can also be used for "top-$k$" decoding.

$(\backslash START\backslash, SUTD, 1111100000, 5, 8.9)$

$(\ldots,$

Beam search: instead of committing to a single next action,
we explore a few options at each point in the search process.

$\ldots)$

# Weighted Score

$$\text{score}(\boldsymbol{e}) = \lambda_{\text{LM}} \times \underbrace{\log q(\boldsymbol{e})}_{\text{Language model}}$$

$$+ \ \lambda_{\text{TM}} \times \underbrace{\sum_{k=1}^{L} \log t(p_k)}_{\text{Phrase translation model}}$$

$$+ \ \lambda_{\text{DM}} \times \underbrace{\sum_{k=1}^{L-1} \eta \times |\text{pl}(p_k) + 1 - \text{pf}(p_{k+1})|}_{\text{Distortion model}}$$

Tunable hyper-parameters

# <u>Question</u>

## How to tune the hyper-parameters?

We shall tune the hyper-parameters to optimize some evaluation metric!

# BLEU



The most widely adopted evaluation metric for measuring MT quality.

41

# MERT



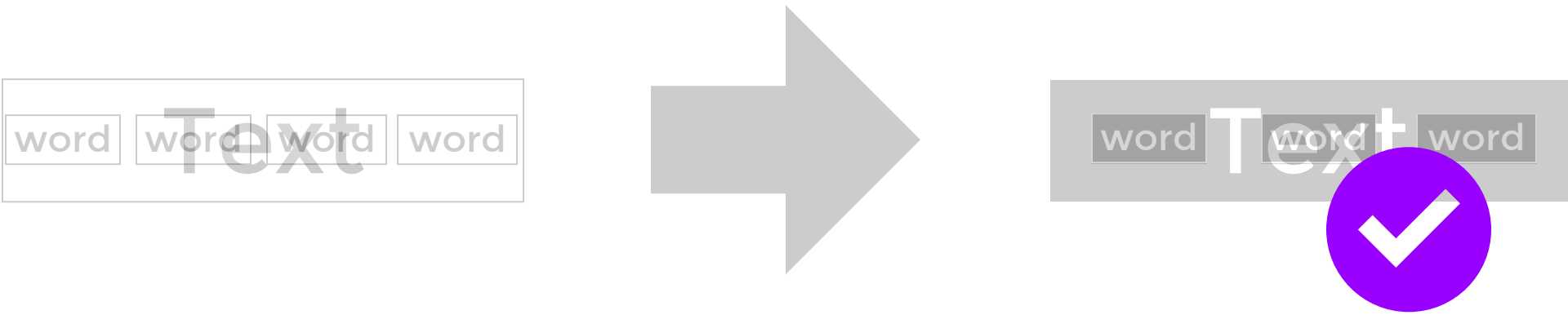Minimum Error Rate Training, which can be used to directly optimize the BLEU score by tuning the hyper-parameters on the development set.

42

# Phrase-based Translation

Training set
(parallel text)

Translation Model $\quad \lambda_{\text{TM}}$

Distortion Model $\quad \lambda_{\text{DM}}$

English Data
(monolingual text)

Language Model $\quad \lambda_{\text{LM}}$

Development set
(parallel text)

"Top-$k$" decoding

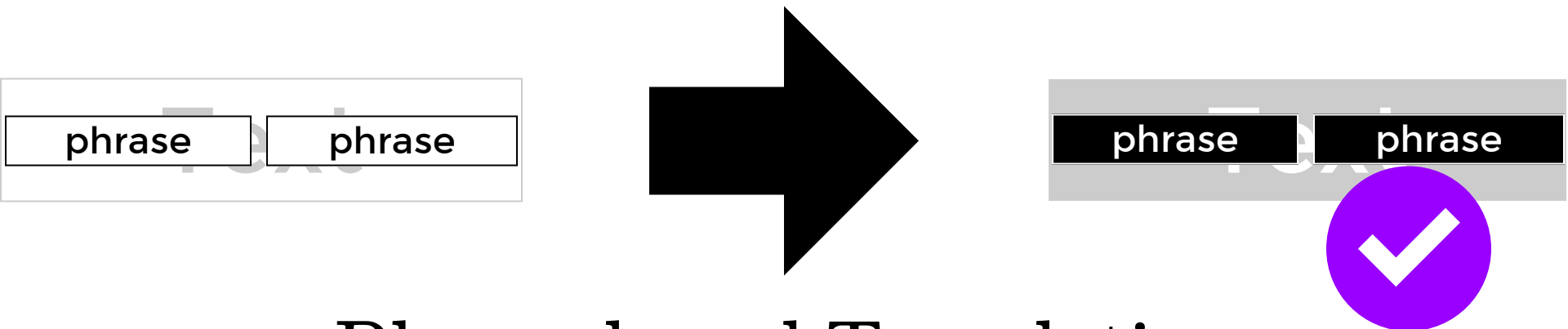MERT

# Machine Translation



Word-based Translation

Phrase-based Translation

# Phrase-based Translation

The process involves a transition-based procedure, which was introduced when discussing parsing.

Text

Parser??

Text

Is it possible to involve a parser in the translation process in some way?

# Machine Translation



Interlingua

Semantic Structure

Syntactic Structure

Words

Upcoming...

Semantic Structure

Syntactic Structure

Words

Syntactic Parsing

Syntactic Transfer

Language Generation