# Statistical Machine Translation without Parallel Data

Maryam Siahbani

Supervisor: Dr. Anoop Sarkar

Natural Language Laboratory

Simon Fraser University

**SFU** NatLangLab

# Statistical Machine Translation (SMT)

$$e = \arg \max_{e} \{ \Pr(e \mid f) \}$$

**e:  target sentence**

**f:  source sentence**



f → Translation System → e

# Statistical Machine Translation (SMT)

$$e = \arg \max_{e} \{ \Pr(e \,|\, f) \}$$
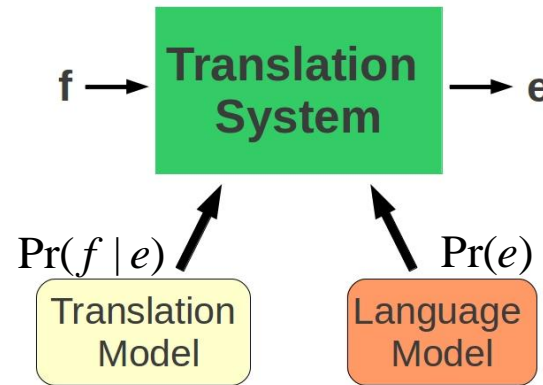
**e:  target sentence**
**f:  source sentence**

**Noisy channel**

$$\Pr(e \,|\, f) \propto \Pr(e).\Pr(f \,|\, e)$$

**Log-linear model**

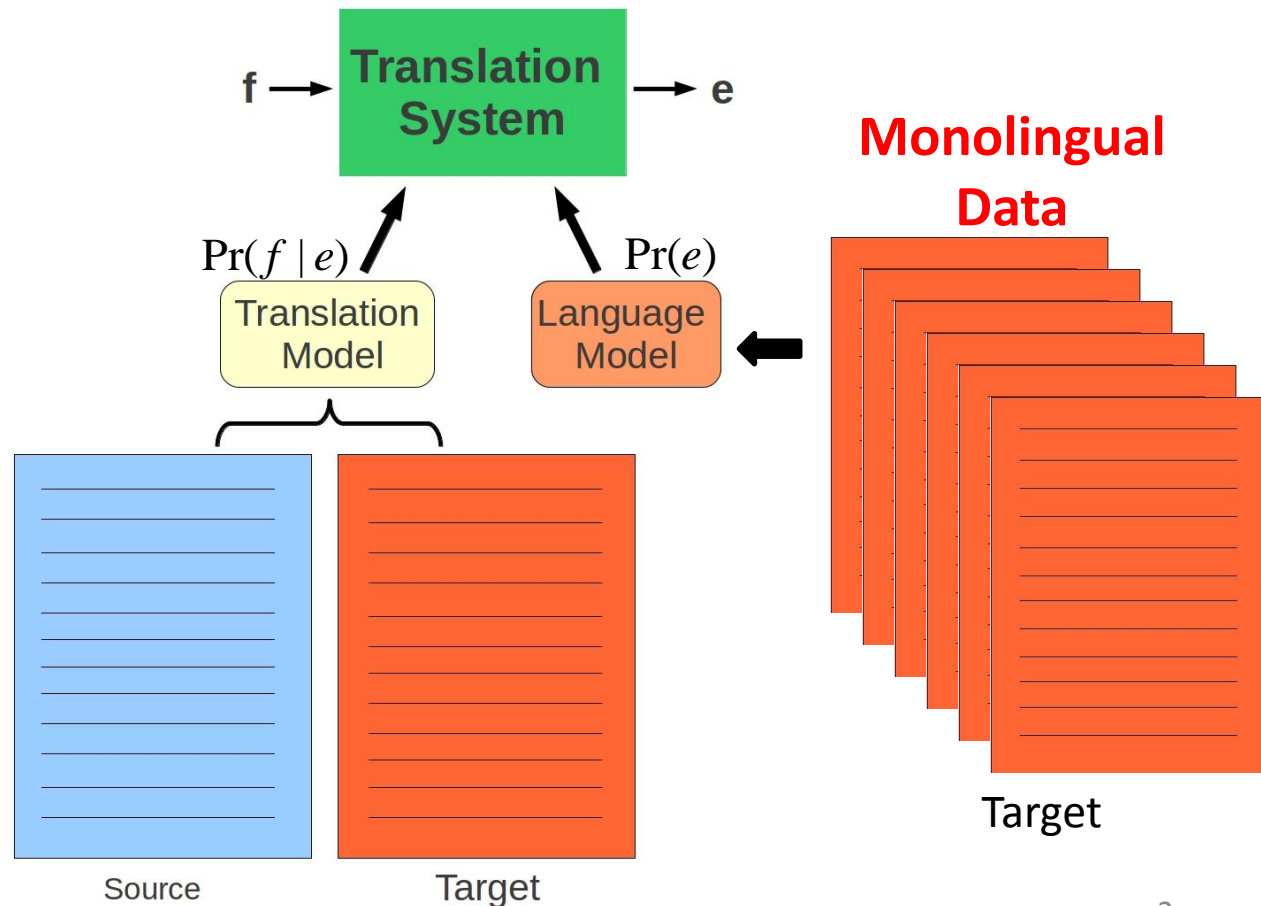$$\Pr(e \,|\, f) \propto \exp \sum_{i} \lambda_{i} h_{i}(e, f)$$

f → **Translation System** → e

$\Pr(f \,|\, e)$    Translation Model

$\Pr(e)$    Language Model

# Statistical Machine Translation (SMT)

$$e = \arg\max_{e} \{ \Pr(e \mid f) \}$$

**e: target sentence**
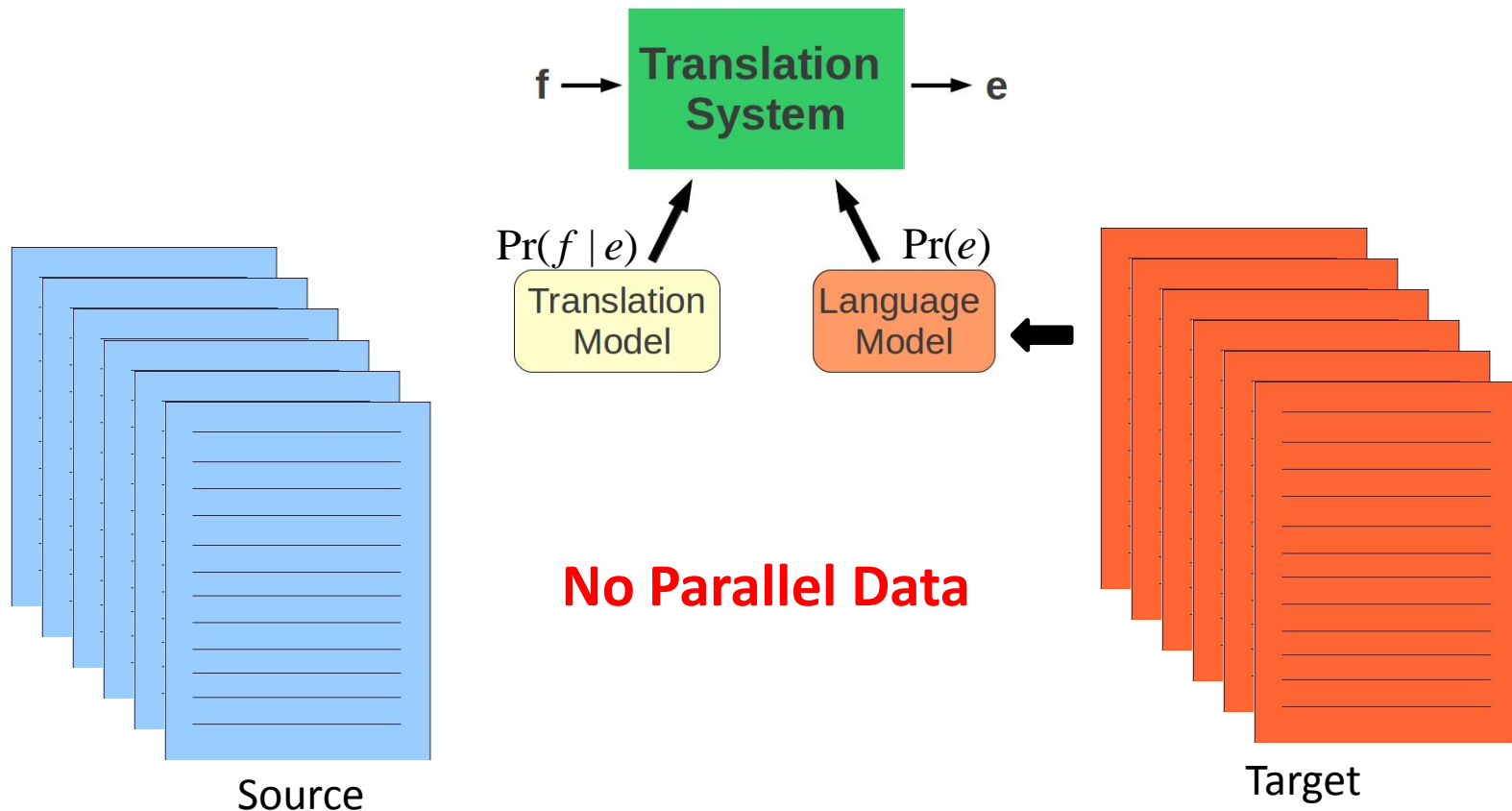**f: source sentence**



**Monolingual Data**

**Parallel Data
(or bilingual data)**

Target

# Statistical Machine Translation (SMT)

$$e = \arg \max_{e} \{ \Pr(e \,|\, f) \}$$

**e: target sentence**
**f: source sentence**



f → **Translation System** → e

$\Pr(f \,|\, e)$   $\Pr(e)$

Translation Model        Language Model
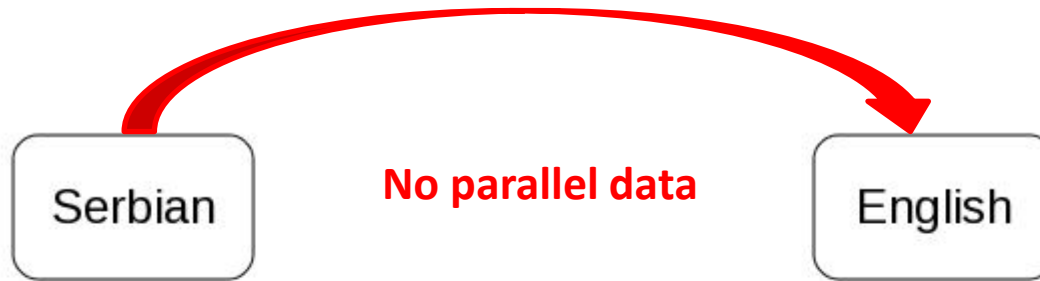
**No Parallel Data**

Source

Target

# Methods

- Translation lexicon
    - Limited parallel resource
    - Monolingual corpora

- Translation system

# Translation Lexicon Induction

# Limited Parallel Data

# Inducing Translation Lexicon: Bridge language



Serbian → English

**No parallel data**

Prazan

# Inducing Translation Lexicon: Bridge language

**Available translation lexicon**



Serbian        Czech        English

Prazan

# Inducing Translation Lexicon: Bridge language



**Language family**  **Available translation lexicon**

Serbian → Czech → English

Prazan → Prazdny

# Inducing Translation Lexicon: Bridge language
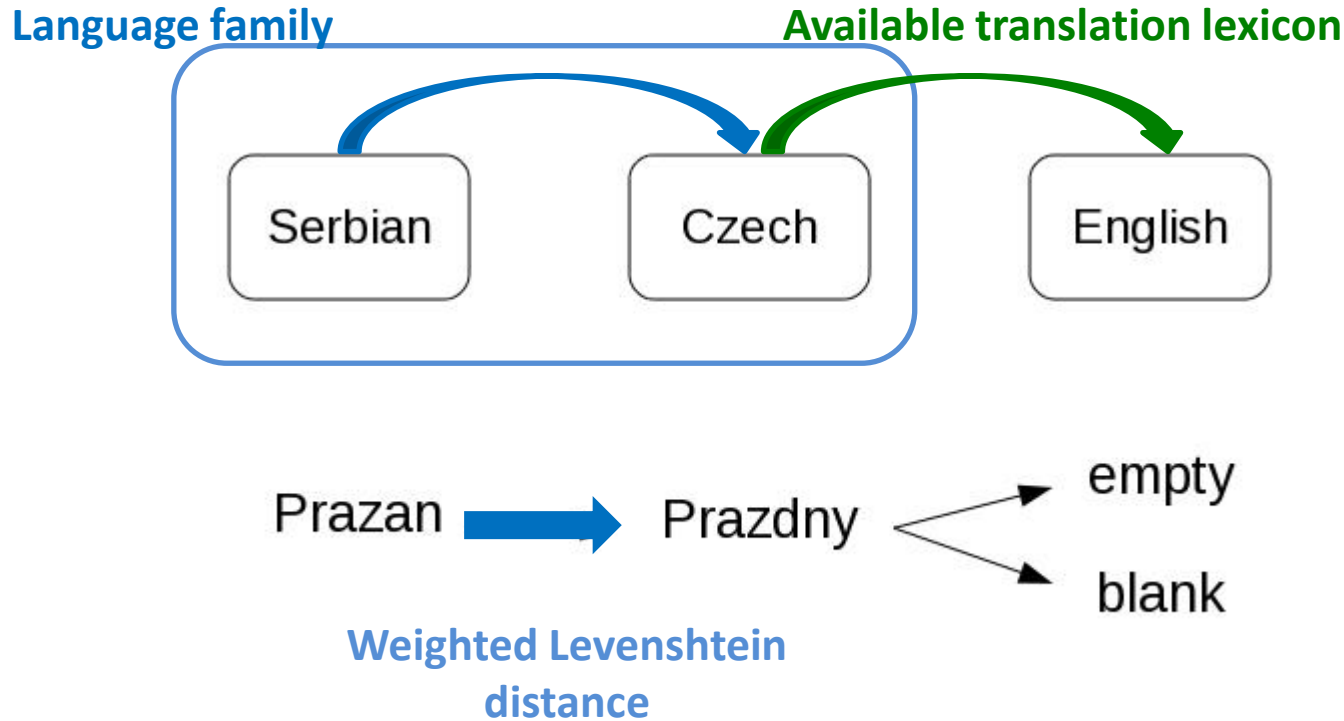
# Inducing Translation Lexicon:
# Bridge language (Mann and Yarowsky, 2001)

**Language family**          **Available translation lexicon**

Serbian → Czech → English

Prazan ⟶ Prazdny < empty / blank

**Weighted Levenshtein distance**

**Operations in Levenshtein:**
Substitution: O**t**por -> O**d**por, cost (t,d)
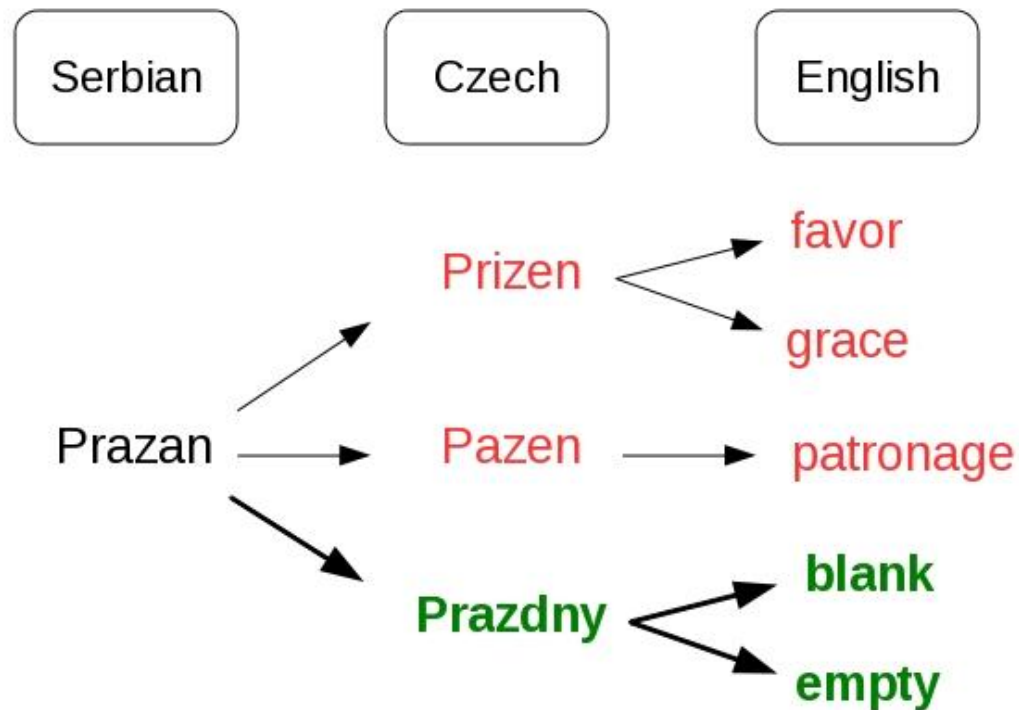Insertion: Pravo -> **V**provo , cost(ε, v)
Deletion: **C**hvala -> Hvaliti , cost (c, ε)

5

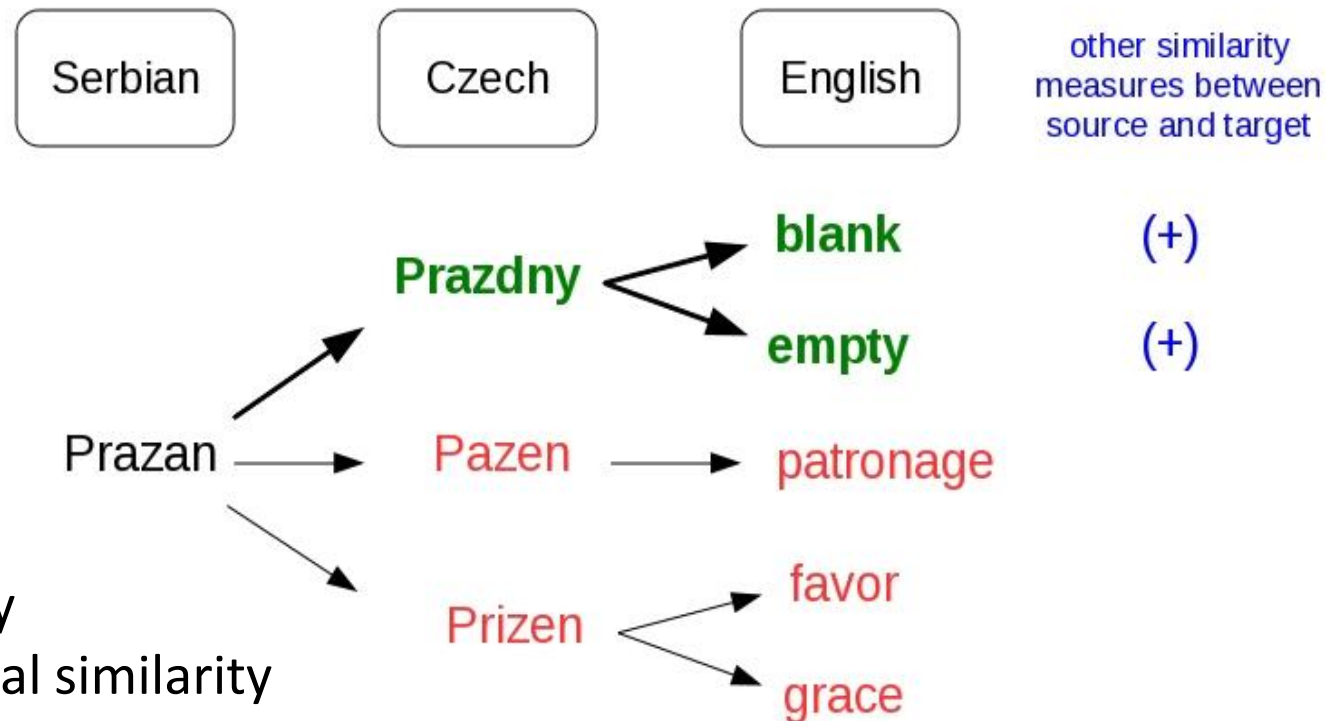# Inducing Translation Lexicon: Bridge language

- False cognates

# Inducing Translation Lexicon: Bridge language (Schafer and Yarowsky 2002)

- Additional similarity measures



- Context similarity
- Date distributional similarity
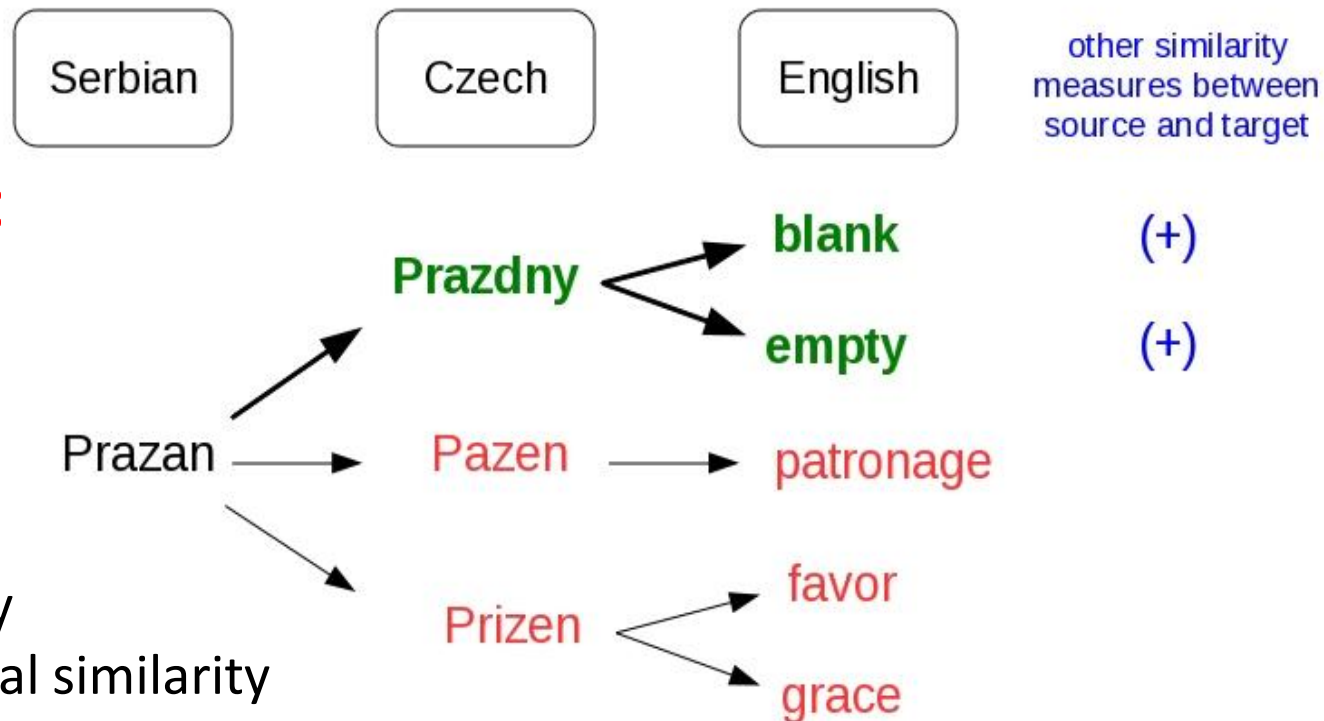- Word frequency similarity

# Inducing Translation Lexicon: Bridge language (Schafer and Yarowsky 2002)

- Additional similarity measures

**9% Improvement**

Accuracy measured against a gold translation lexicon

- Context similarity
- Date distributional similarity
- Word frequency similarity



other similarity measures between source and target

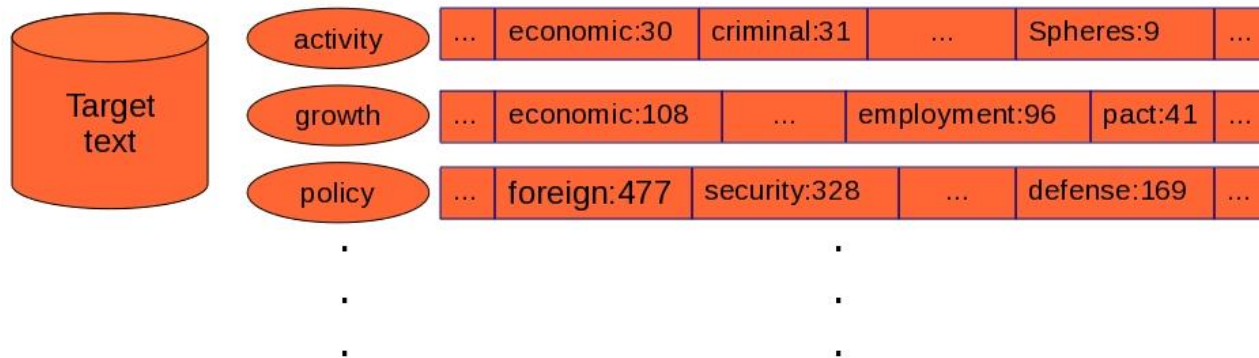# Inducing Translation Lexicon: Using Seed Lexicon

- Given: an initial small lexicon between source and target

- Goal: a translation lexicon between S and T

- Methods: extending the seed lexicon
  - Context similarity

# Inducing Translation Lexicon:
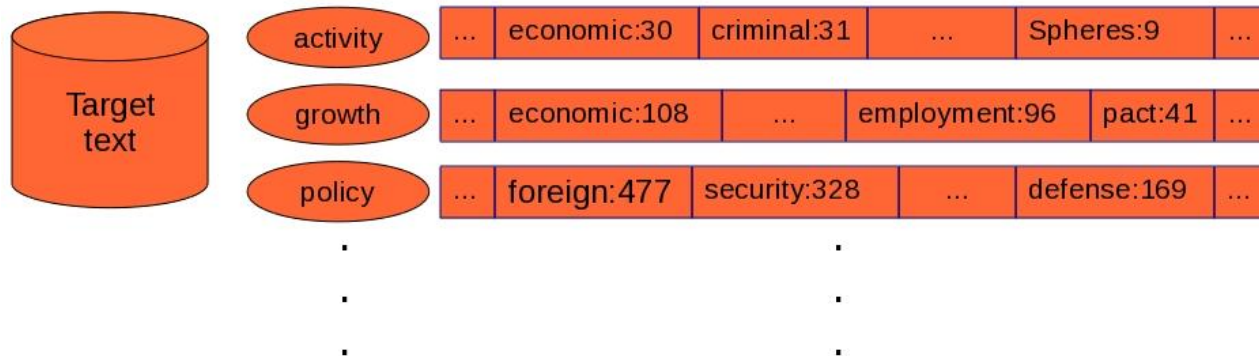# Using Seed Lexicon (Context similarity)

- (Rapp 95),(Rapp 99),(Fung & Yee 98)

# Inducing Translation Lexicon:
# Using Seed Lexicon (Context similarity)

- (Rapp 95),(Rapp 99),(Fung & Yee 98)

# Inducing Translation Lexicon: Using Seed Lexicon (Context similarity)

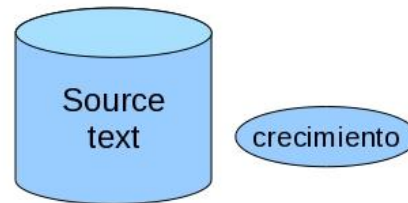- (Rapp 95),(Rapp 99),(Fung & Yee 98)

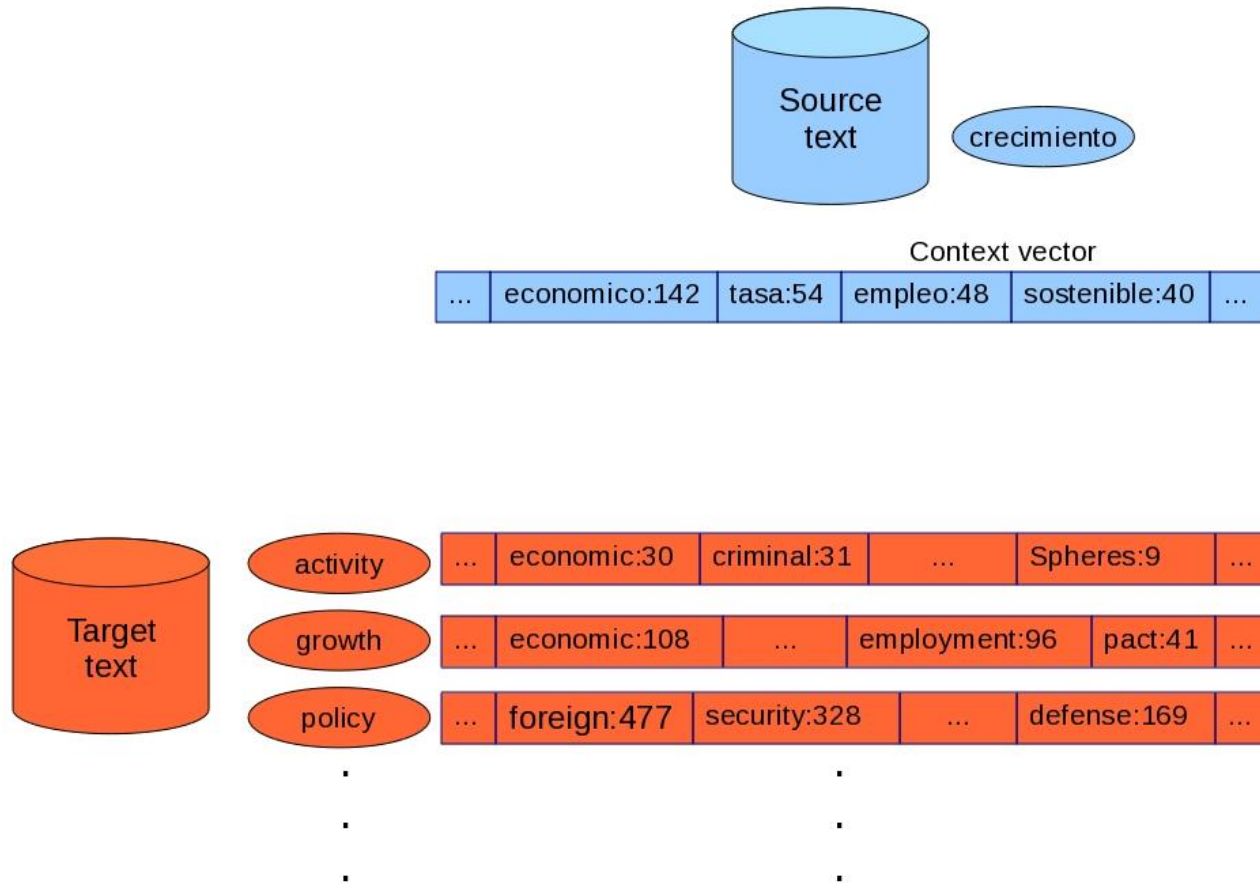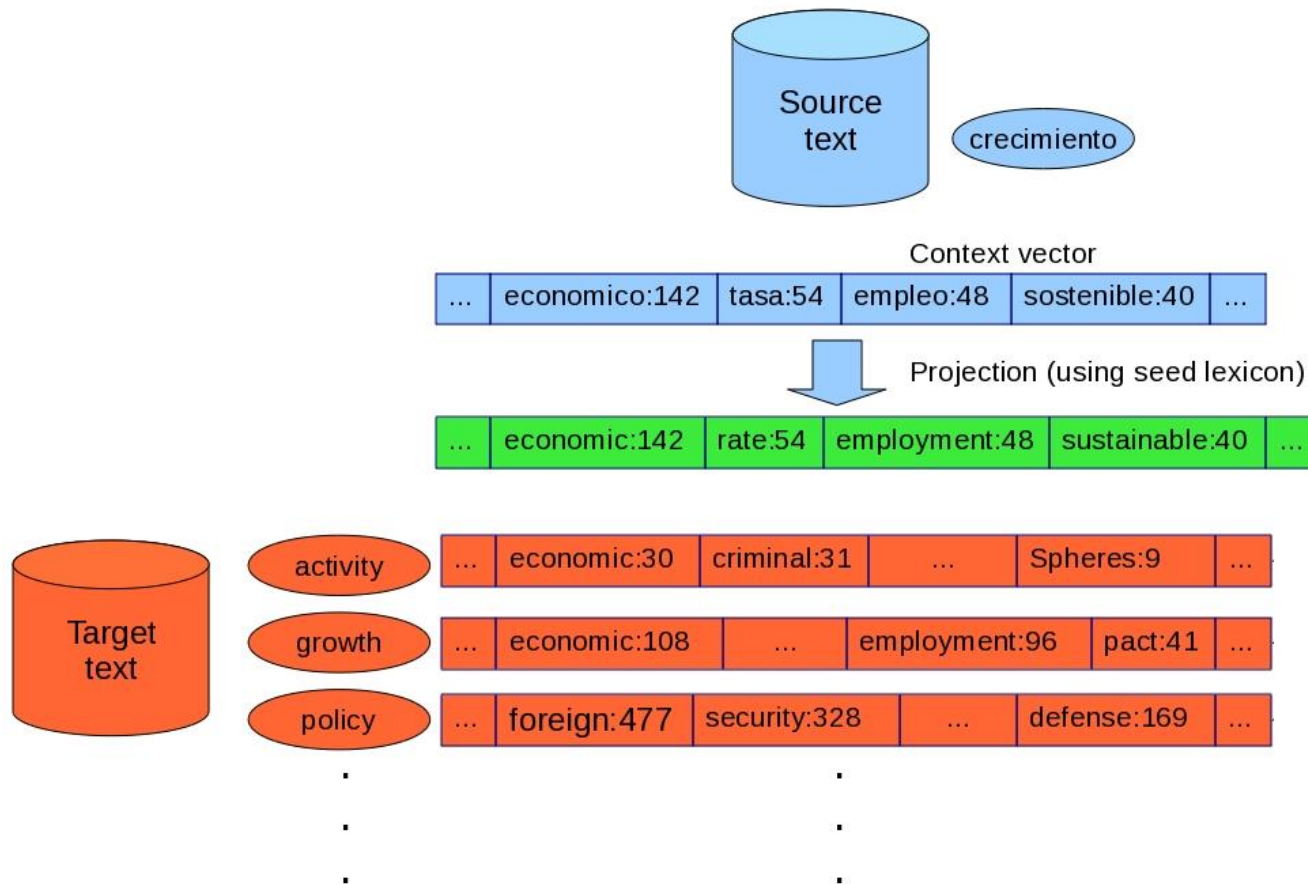# Inducing Translation Lexicon: Using Seed Lexicon (Context similarity)

- (Rapp 95),(Rapp 99),(Fung & Yee 98)

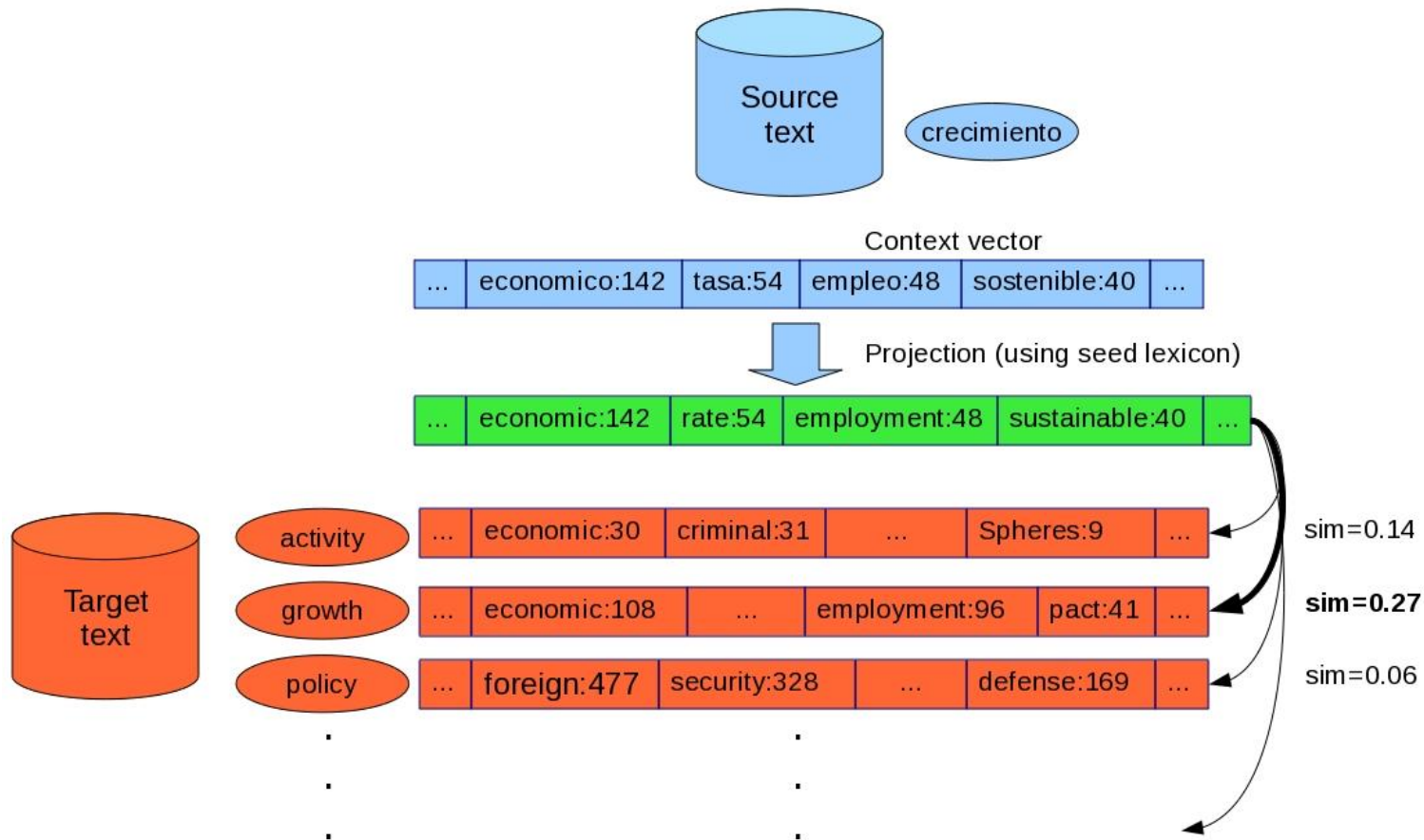# Inducing Translation Lexicon: Using Seed Lexicon (Context similarity)

- (Rapp 95),(Rapp 99),(Fung & Yee 98)

# Inducing Translation Lexicon: Using Seed Lexicon (Context similarity)

- ## Modeling the context
  - – Window of fixed size (Rapp 95),(Rapp 99),(Fung & Yee 98)

... el camino para el **crecimiento** y la prosperidad económica ...
(the) (path) (to) (the) (growth) (and) (the) (prosperity) (economic)

# Inducing Translation Lexicon:
# Using Seed Lexicon (Context similarity)

• Modeling the context

  – Window of fixed size (Rapp 95),(Rapp 99),(Fung & Yee 98)

... el camino para el **crecimiento** y la prosperidad económica ...
(the) (path)    (to) (the)    (growth)    (and) (the) (prosperity)    (economic)

| Position | Adjacent context |
|----------|------------------|
| -2       | para             |
| -1       | el               |
| +1       | y                |
| +2       | la               |

# Inducing Translation Lexicon: Using Seed Lexicon (Context similarity)

- Modeling the context
  - Window of fixed size (Rapp 95),(Rapp 99),(Fung & Yee 98)
  - Using dependency trees (Garera et al., 2009)

... el camino para el **crecimiento** y la prosperidad económica ...
(the) (path)   (to) (the)    (growth)   (and) (the) (prosperity)   (economic)

| Position | Adjacent context | Dependency context |
|---|---|---|
| -2 | para | |
| -1 | el | |
| +1 | y | |
| +2 | la | |

# Inducing Translation Lexicon:
# Using Seed Lexicon (Context similarity)

- ## Modeling the context
  - Window of fixed size (Rapp 95),(Rapp 99),(Fung & Yee 98)
  - Using dependency trees (Garera et al., 2009)

... el camino para el **crecimiento** y la prosperidad económica ...
(the) (path)    (to) (the)    (growth)    (and) (the) (prosperity)    (economic)

| Position | Adjacent context | Dependency context |
|----------|------------------|--------------------|
| -2       | para             | camino             |
| -1       | el               | para               |
| +1       | y                |                    |
| +2       | la               |                    |

# Inducing Translation Lexicon:
# Using Seed Lexicon (Context similarity)

- ## Modeling the context
  - Window of fixed size (Rapp 95),(Rapp 99),(Fung & Yee 98)
  - Using dependency trees (Garera et al., 2009)

... el camino para el **crecimiento** y la prosperidad económica ...
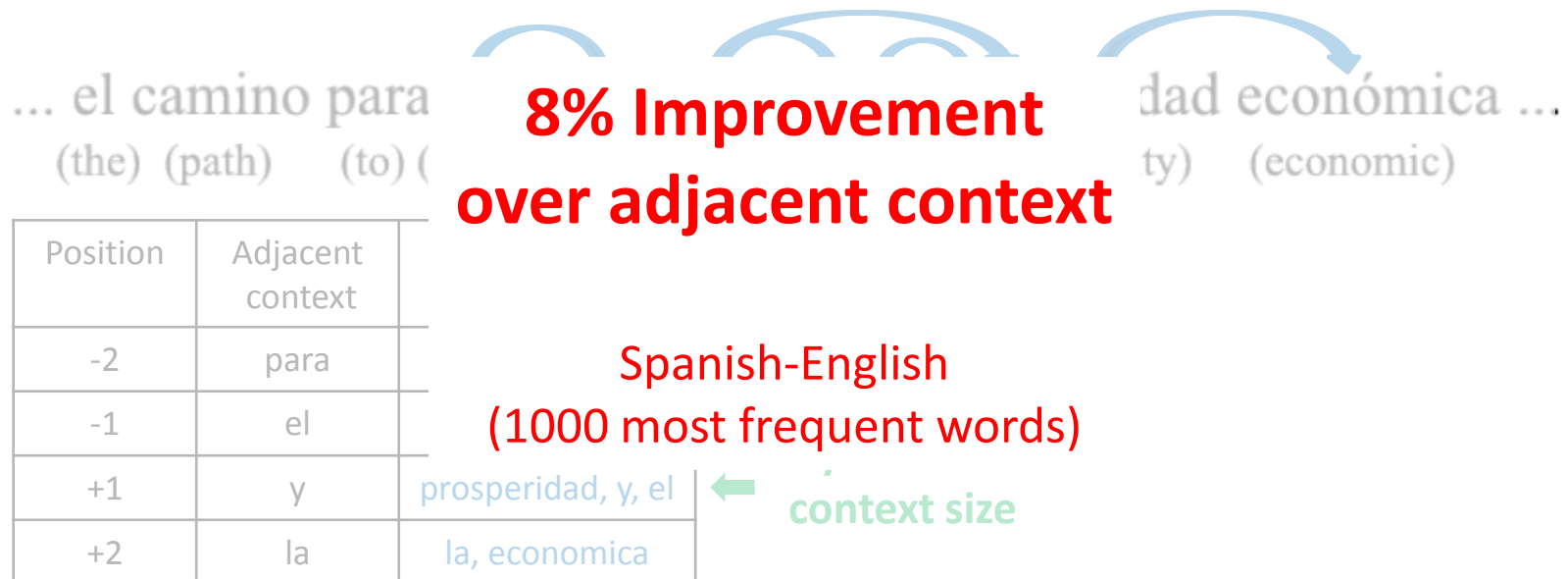(the) (path)    (to) (the)    (growth)    (and) (the) (prosperity)    (economic)

| Position | Adjacent context | Dependency context |
|----------|------------------|--------------------|
| -2 | para | camino |
| -1 | el | para |
| +1 | y | prosperidad, y, el |
| +2 | la | la, economica |

**Dynamic context size**

# Inducing Translation Lexicon: Using Seed Lexicon (Context similarity)

- Modeling the context
  - Window of fixed size (Rapp 95),(Rapp 99),(Fung & Yee 98)
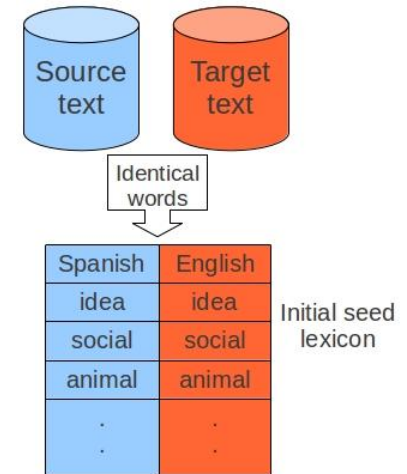  - Using dependency trees (Garera et al., 2009)

… el camino para

(the) (path)   (to) (

**8% Improvement
over adjacent context**

dad económica …

ty)   (economic)

| Position | Adjacent context | |
|----------|------------------|--|
| -2 | para | |
| -1 | el | |
| +1 | y | prosperidad, y, el |
| +2 | la | la, economica |

Spanish-English
(1000 most frequent words)

← **context size**

# No Parallel Data

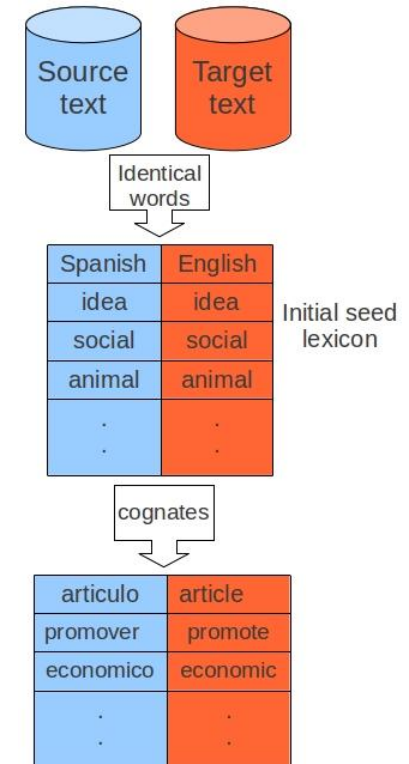# Inducing Translation Lexicon: Monolingual Data (Similarity measures)

- Similarities in monolingual corpora (Koehn & Knight 2002)
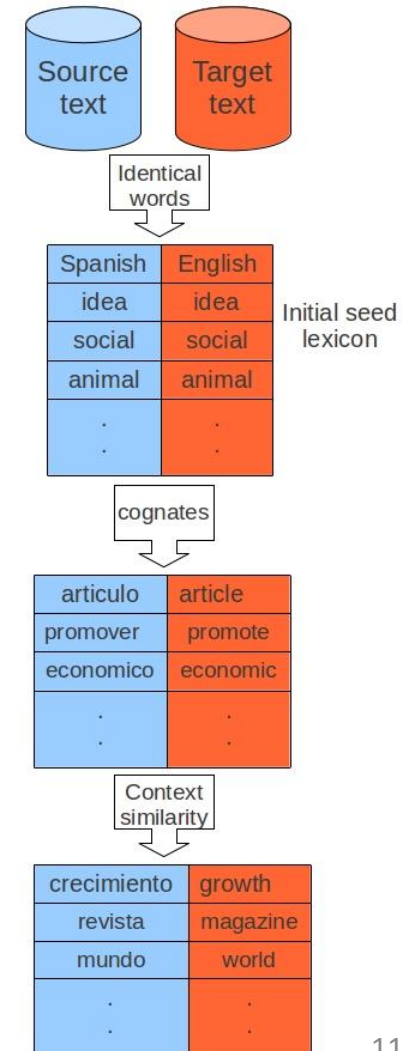  - Identical words

# Inducing Translation Lexicon: Monolingual Data (Similarity measures)

- Similarities in monolingual corpora (Koehn & Knight 2002)
  - Identical words
  - Orthographic similarity (cognates)

# Inducing Translation Lexicon: Monolingual Data (Similarity measures)
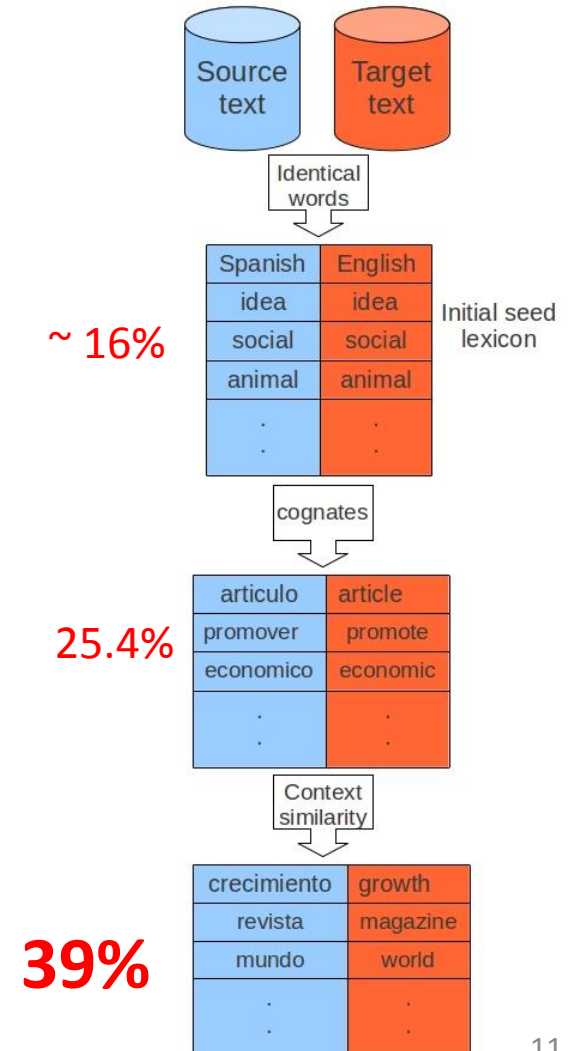
- Similarities in monolingual corpora (Koehn & Knight 2002)
  - Identical words
  - Orthographic similarity (cognates)
  - Context similarity



11

# Inducing Translation Lexicon: Monolingual Data (Similarity measures)

- Similarities in monolingual corpora (Koehn & Knight 2002)
  - Identical words
  - Orthographic similarity (cognates)
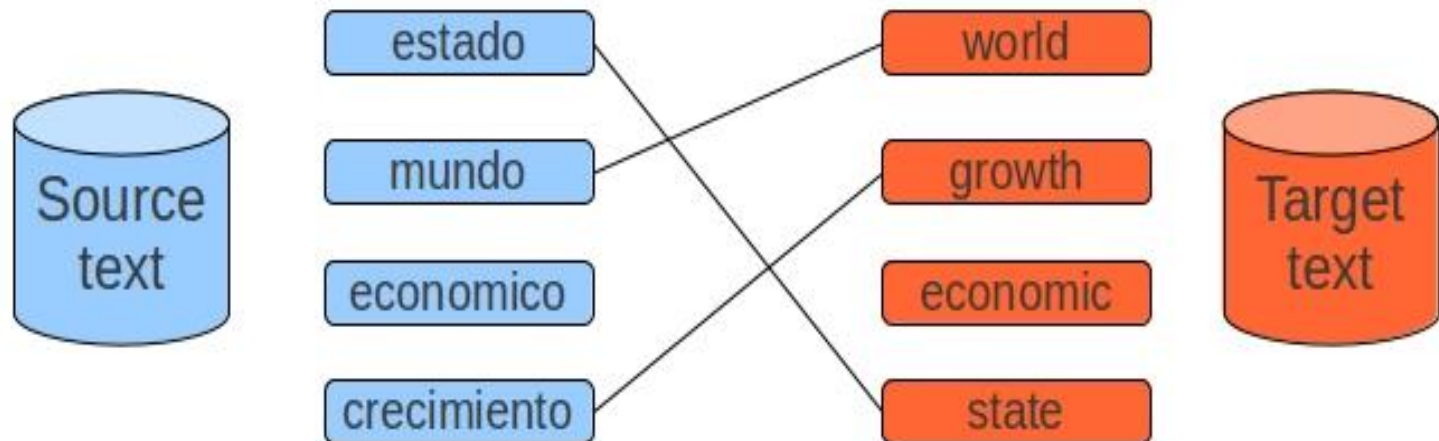  - Context similarity

German-English
1000 most frequent nouns



~ 16%

25.4%

**39%**

11

# Inducing Translation Lexicon:
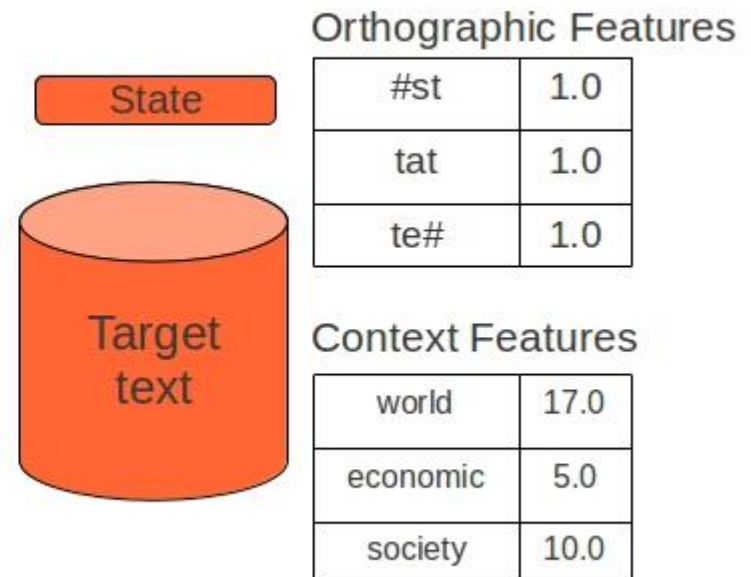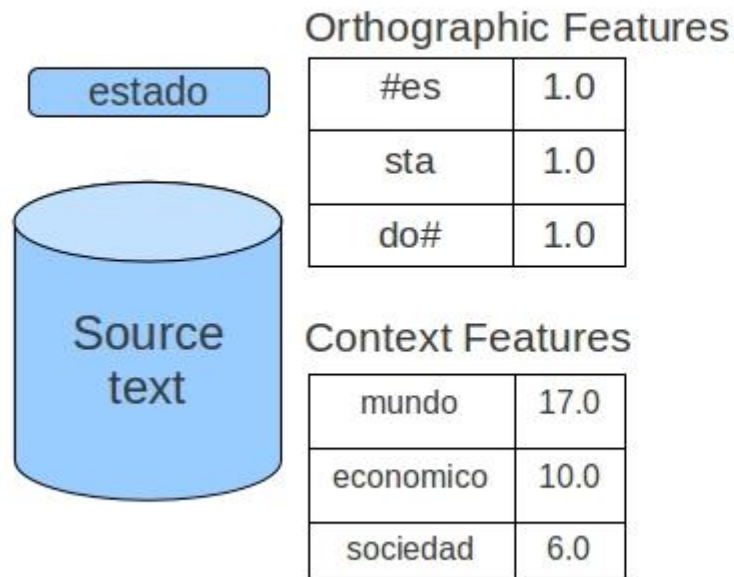No Parallel Data (Model Translation Lexicon as a Mapping)

- A generative model for matching between S and T (Haghighi et al., 2008)

# Inducing Translation Lexicon:
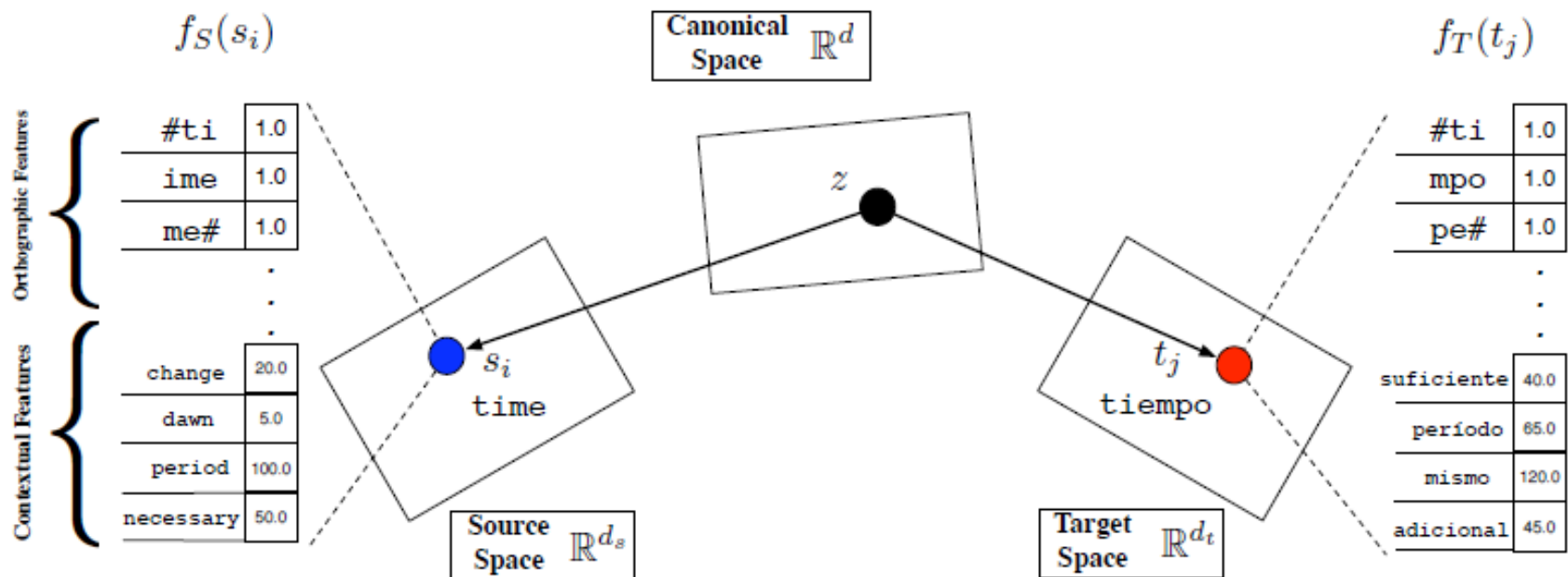## No Parallel Data (Haghighi et al., 2008)

- Orthographic (n-gram characters)
- Contextual features



Orthographic Features

estado

| #es | 1.0 |
|-----|-----|
| sta | 1.0 |
| do# | 1.0 |

Source text

Context Features

| mundo | 17.0 |
|-------|------|
| economico | 10.0 |
| sociedad | 6.0 |

Orthographic Features

State

| #st | 1.0 |
|-----|-----|
| tat | 1.0 |
| te# | 1.0 |

Target text

Context Features

| world | 17.0 |
|-------|------|
| economic | 5.0 |
| society | 10.0 |

# Inducing Translation Lexicon:
# No Parallel Data (Haghighi et al., 2008)

**Observed words in source and target spaces are projected to a common latent space**

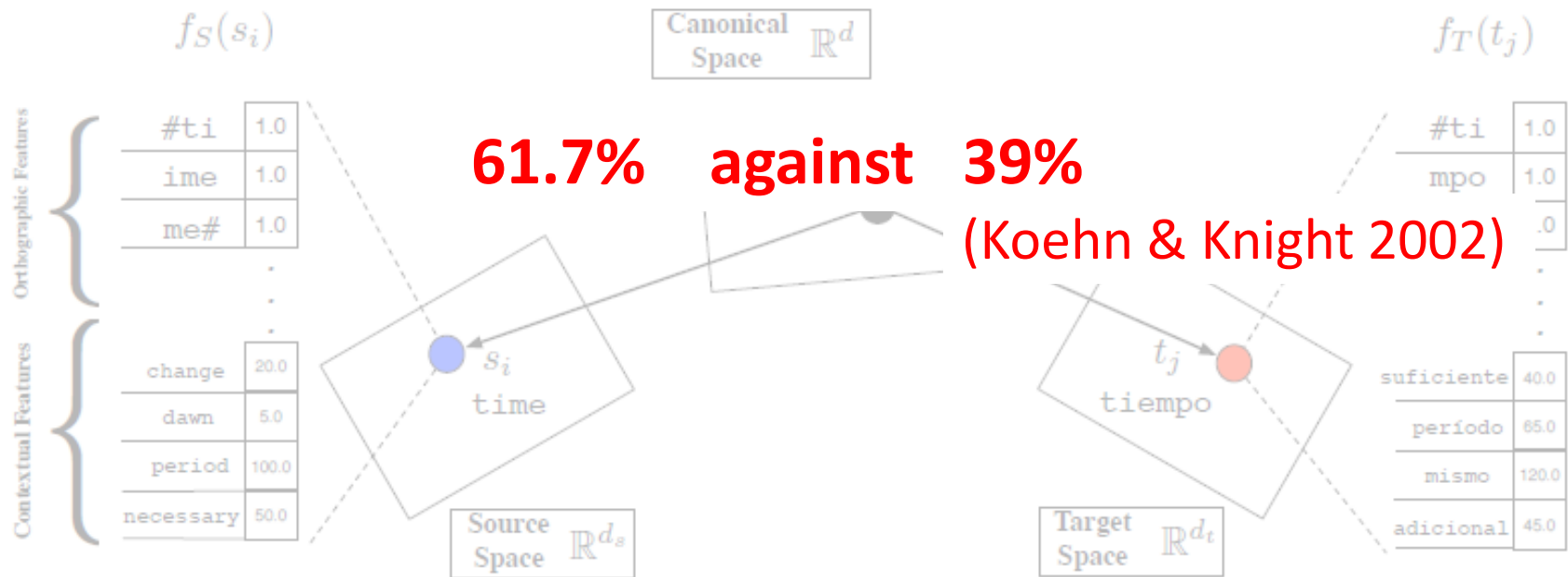**Canonical Correlation Analysis (CCA)**



(Haghighi et al., 2008)

# Inducing Translation Lexicon:
# No Parallel Data (Haghighi et al., 2008)

**Observed words in source and target spaces are projected to a common latent space**

**Canonical Correlation Analysis (CCA)**



**61.7%   against   39%**
(Koehn & Knight 2002)

(Haghighi et al., 2008)

# Methods

- Inducing Translation Lexicon
  - Limited Parallel Data
  - Monolingual Data

- Translation System
  - Word-based
  - Phrase-based

# Translation System:
## Word-based: MT as decipherment (Ravi and Knight 2011)

**f: source sentence**        **e: target sentence**

Train parameters to maximize probability of **observed sentence pairs (e,f)**:

$$\underset{\theta}{argmax}\, P_\theta(e,f) \simeq \underset{\theta}{argmax} \prod_{e,f} P_\theta(f|e)$$

$P(f|e)$
English-Spanish
Translation Model

Spanish        English

**Training with Parallel data**

# Translation System:
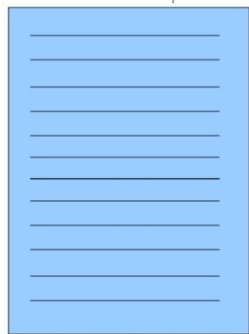## Word-based: MT as decipherment (Ravi and Knight 2011)

**f: source sentence**      **e: target sentence**

Train parameters to maximize probability of **observed sentence pairs (e,f)**:

$$argmax_{\theta} P_{\theta}(e,f) \simeq argmax_{\theta} \prod_{e,f} P_{\theta}(f|e)$$

$P(f|e)$
English-Spanish
Translation Model

Spanish          English

**Training with Parallel data**

English

Train parameters to maximize probability of **observed source text f**:

$$argmax_{\theta} P_{\theta}(f) \simeq argmax_{\theta} \prod_{f} P_{\theta}(f)$$
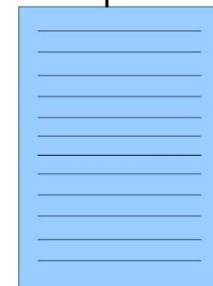
$$argmax_{\theta} \prod_{f} \sum_{e} P(e) P_{\theta}(f|e)$$

$P(e)$
Language Model

$P(f|e)$
English-Spanish
Translation Model

Spanish

**Training without Parallel data**

16

# Translation System:
## Word-based: MT as decipherment (Ravi and Knight 2011)

**f: source sentence**     **e: target sentence**

Train parameters to maximize probability of **observed sentence pairs (e,f)**:

$$argmax_\theta P_\theta(e,f) \simeq argmax_\theta \prod_{e,f} P_\theta(f|e)$$

$P(f|e)$
English-Spanish Translation Model

Spanish    English

**Training with Parallel data**

English

Train parameters to maximize probability of **observed source text f**:
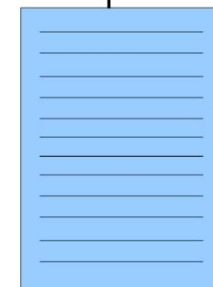
$$argmax_\theta P_\theta(f) \simeq argmax_\theta \prod_f P_\theta(f)$$

$$argmax_\theta \prod_f \sum_e P(e)P_\theta(f|e)$$

$P(e)$
Language Model

$P(f|e)$
English-Spanish Translation Model

Spanish

**Training without Parallel data**

# Translation System:
## Word-based: MT as decipherment (Ravi and Knight 2011)
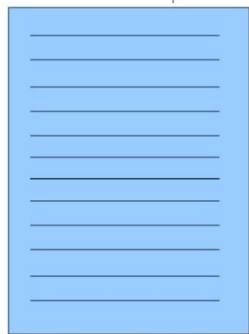
**f: source sentence**      **e: target sentence**

Train parameters to maximize probability of **observed sentence pairs (e,f)**:

$$argmax_{\theta} P_{\theta}(e,f) \simeq argmax_{\theta} \prod_{e,f} P_{\theta}(f|e)$$

$P(f|e)$
English-Spanish Translation Model

Spanish      English

**Training with Parallel data**

English

Train parameters to maximize probability of **observed source text f**:

$$argmax_{\theta} P_{\theta}(f) \simeq argmax_{\theta} \prod_{f} P_{\theta}(f)$$

$$argmax_{\theta} \prod_{f} \sum_{e} P(e) P_{\theta}(f|e)$$

$P(e)$
Language Model

$P(f|e)$
English-Spanish Translation Model

Scalability Challenges

Spanish

**Training without Parallel data**
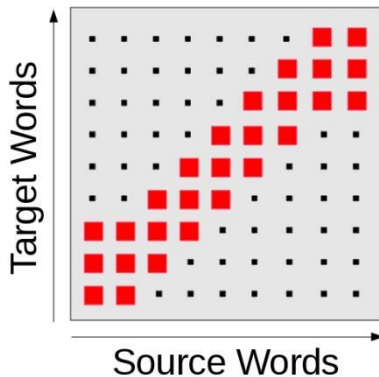
16

# Translation System: Word-based

- Restricting the translation model (Nuhn et al., 2012)
  - Determining a set of active translations
  - Estimating the probabilities of active translations

# Translation System: Word-based

- Restricting the translation model (Nuhn et al., 2012)
  - Determining a set of active translations
  - Estimating the probabilities of active translations



Initialization using
word frequency ranks

# Translation System: Word-based

- Restricting the translation model (Nuhn et al., 2012)
  - Determining a set of active translations
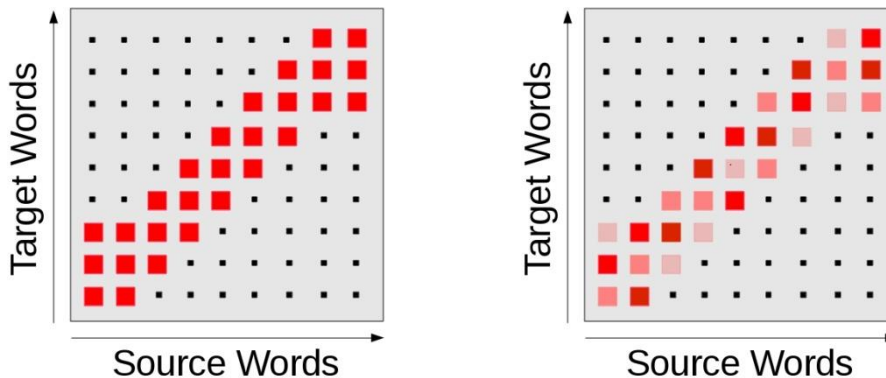  - Estimating the probabilities of active translations



  - EM training (Ravi and Knight 2011)
    with limited iterations (20-30)

# Translation System: Word-based

- Restricting the translation model (Nuhn et al., 2012)
  - Determining a set of active translations
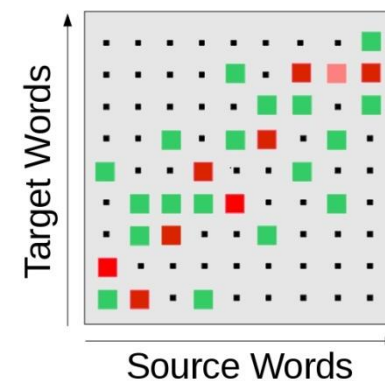  - Estimating the probabilities of active translations



Inducing active translations using results of EM

# Translation System: Word-based

- Restricting the translation model (Nuhn et al., 2012)
  - Determining a set of active translations
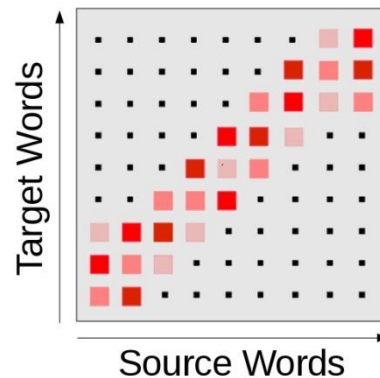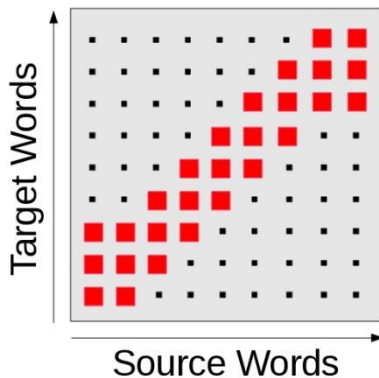  - Estimating the probabilities of active translations

# Translation System: Word-based

- Restricting the translation model (Nuhn et al., 2012)
  - Determining a set of active translations
  - Estimating the probabilities of active translations
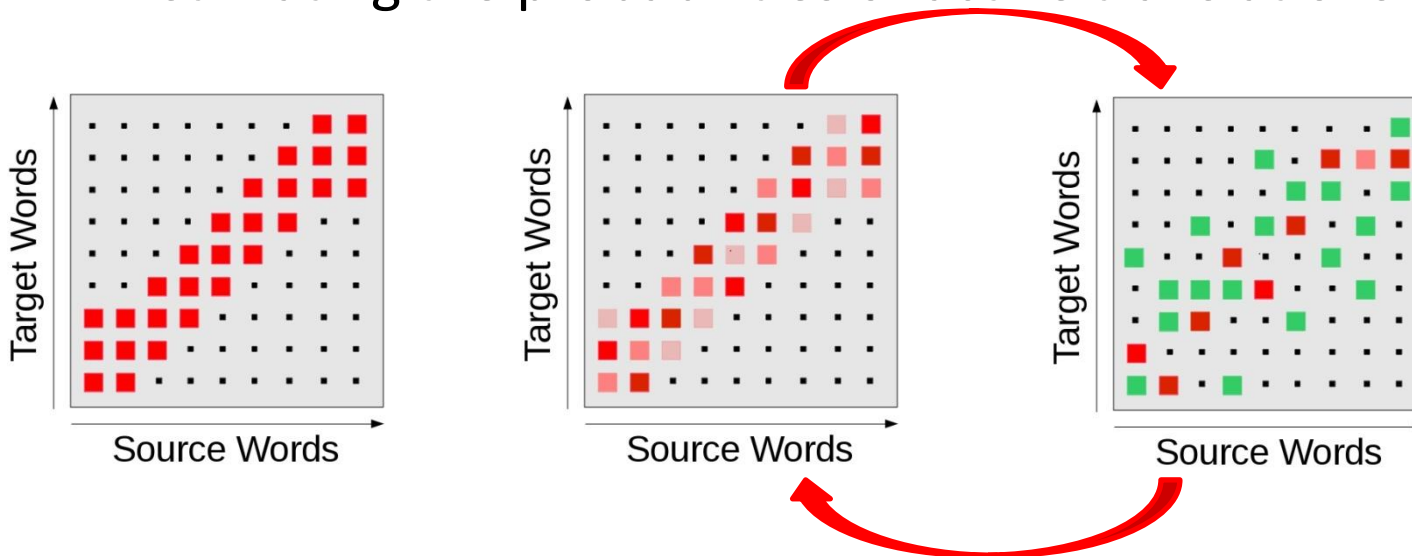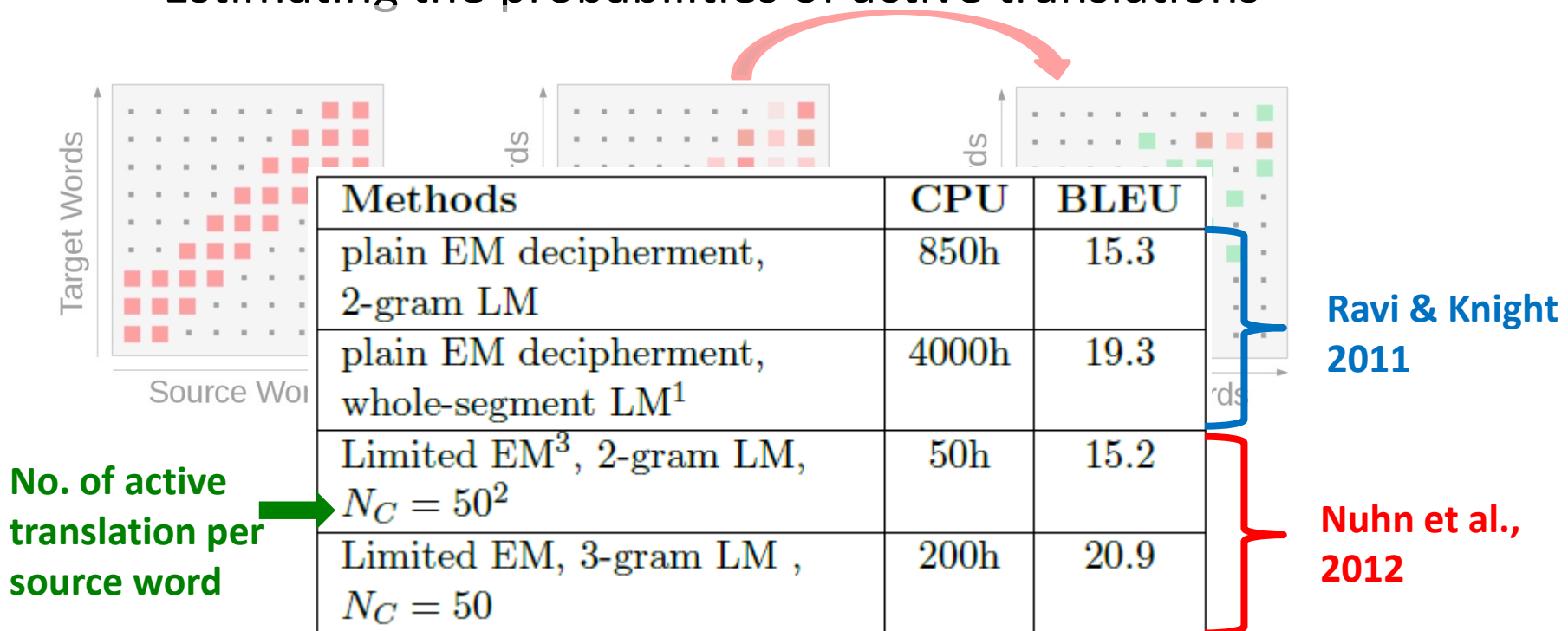
**No. of active translation per source word**

| Methods | CPU | BLEU |
|---|---|---|
| plain EM decipherment, 2-gram LM | 850h | 15.3 |
| plain EM decipherment, whole-segment LM[1] | 4000h | 19.3 |
| Limited EM[3], 2-gram LM, $N_C = 50^2$ | 50h | 15.2 |
| Limited EM, 3-gram LM, $N_C = 50$ | 200h | 20.9 |

**Ravi & Knight 2011**

**Nuhn et al., 2012**

# Translation System

## Limited Parallel Data

# Phrase-based Translation System: Limited Parallel Data(Klementiev et al.,2012)

- Parameters in phrase-based translation systems
  - Large amounts of parallel data for estimating parameters

# Phrase-based Translation System: Limited Parallel Data(Klementiev et al.,2012)

- Parameters in phrase-based translation systems
  - Large amounts of parallel data for  estimating parameters
- Extract features from monolingual data
  - Extend the idea of translation lexicon induction to phrases (using seed lexicon)
  - Reordering model

# Phrase-based Translation System: Limited Parallel Data(Klementiev et al.,2012)

- Parameters in phrase-based translation systems

  – Large amounts of parallel

  **Monolingual** &

  **Bilingual Features**
  **21.87** BLEU score

  Bilingual Features
  **22.92** BLEU score

  from m

  – Extend the idea of translation lexicon induction to phrases (using seed lexicon)

  – Reordering model

# Conclusion

- Translation lexicon induction:
  - Depend on context and orthographic  similarities
    - Does not work on historically unrelated language pairs
  - Scalability: Just applied on small word sets(1000 or so)

- Translation Systems:
  - New research direction
  - Scalability: applicable to limited vocabularies and data sets
  - Low translation quality