

# Bilingual Word Embeddings and Unsupervised SMT

Viktor Hangya  
*hangyav@cis.lmu.de*

CIS, LMU Munich

2021

# Outline

## Word Embeddings

- Vector representation
- Learning word embeddings

## Bilingual Word Embeddings

- Motivation
- Overview of approaches
- Mapping
- Orthogonal mapping
- Unsupervised training

## Unsupervised Statistical MT

- MT with bilingual word embeddings
- Integrating embeddings into Moses

# Word representation

How do we represent words in neural networks?

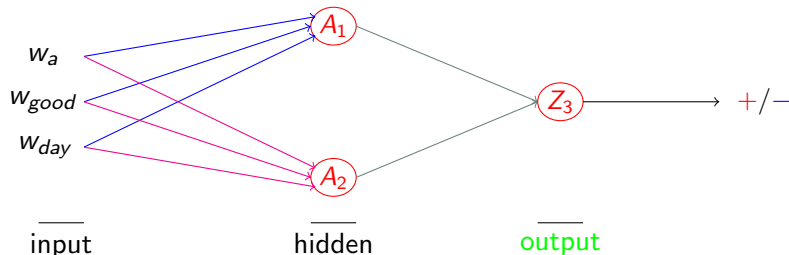
- One-hot vector:

- $w_{good} = [1, 0, 0, 0, \dots, 0]$

- $w_{great} = [0, 1, 0, 0, \dots, 0]$

- $w_{day} = [0, 0, 1, 0, \dots, 0]$

- $w_a = [0, 0, 0, 1, \dots, 0]$



# Word representation

How do we represent words in neural networks?

- One-hot vector:

- ▶  $w_{good} = [1, 0, 0, 0, \dots, 0]$
- ▶  $w_{great} = [0, 1, 0, 0, \dots, 0]$
- ▶  $w_{day} = [0, 0, 1, 0, \dots, 0]$
- ▶  $w_a = [0, 0, 0, 1, \dots, 0]$

- Does not carry word similarity information

- This would be better

- ▶  $w_{good} = [\mathbf{1}, \mathbf{0}, 0, 0, \dots, 0]$
- ▶  $w_{great} = [\mathbf{1}, \mathbf{0}, 0, 0, \dots, 0]$
- ▶  $w_{day} = [0, 0, 1, 0, \dots, 0]$
- ▶  $w_a = [0, 0, 0, 1, \dots, 0]$

# Word representation

How do we represent words in neural networks?

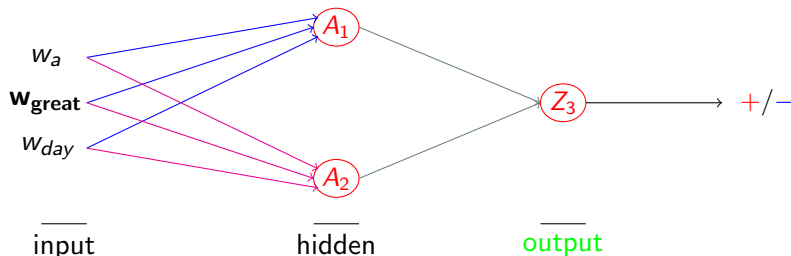
- One-hot vector:

▶  $w_{good} = [1, 0, 0, 0, \dots, 0]$

▶  $w_{great} = [1, 0, 0, 0, \dots, 0]$

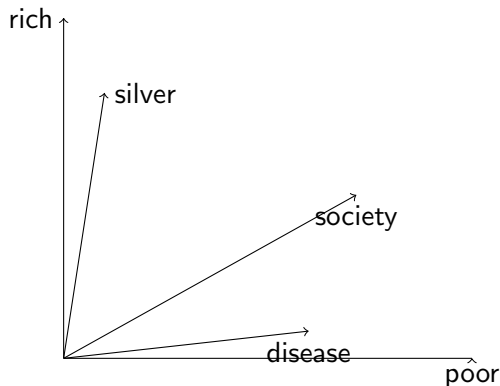
▶  $w_{day} = [0, 0, 1, 0, \dots, 0]$

▶  $w_a = [0, 0, 0, 1, \dots, 0]$



# Word Embeddings

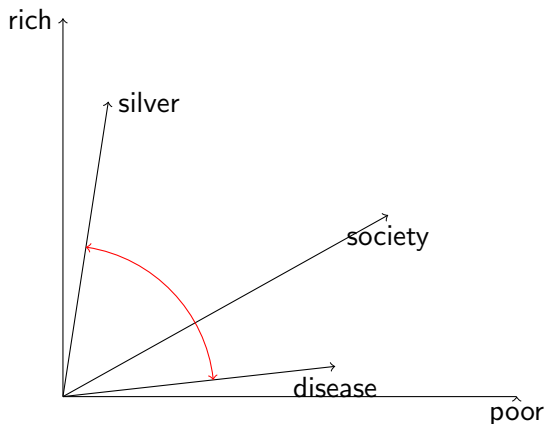
- Representation of words in vector space
- Learn a low dimensional vector representations
  - ▶ typically *dimensions* = 300
  - ▶  $w_{good} = [0.234, 0.001, -0.456, 0.000, \dots, -0.938]$



# Word Embeddings

- Similar words are close to each other  
→ Similarity is the cosine of the angle between two word vectors

$$\text{cosine}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$



## Word Embeddings

- Vector representation
- Learning word embeddings

## Bilingual Word Embeddings

- Motivation
- Overview of approaches
- Mapping
- Orthogonal mapping
- Unsupervised training

## Unsupervised Statistical MT

- MT with bilingual word embeddings
- Integrating embeddings into Moses



# Learning word embeddings

**Distributional similarity:** the meaning of words can be understood from their context

*I drink **water** before bed.*                      or                      *I drink **milk** before bed.*

Count-based methods:

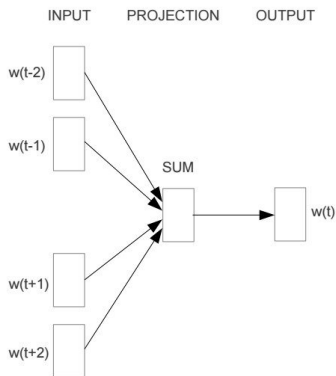
- Compute co-occurrence statistics
- Map sparse high-dimensional vectors to small dense representation
- Matrix factorization approaches: SVD

Neural networks:

- Predict a word from its neighbors
  - ▶ **Word2Vec: CBOW and skipgram** Mikolov et al. (2013a)
- Language Modeling Task
  - ▶ ELMo, BERT Peters et al. (2018); Devlin et al. (2018)

# Learning word embeddings with Continuous Bag-Of-Words

Training example:... *CEO of BMW was fired ...*

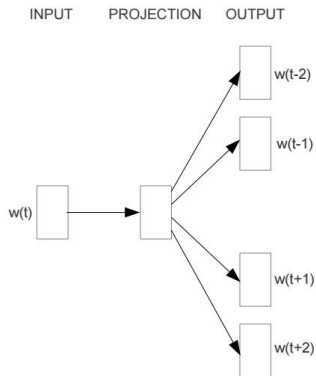


**CBOW**

Mikolov et al. (2013a)

# Learning word embeddings with skip-gram

Training example: ... CEO of *BMW* was fired ...



**Skip-gram**

Mikolov et al. (2013a)

# Word embedding quality

- Semantic similarity
  - ▶ Correlation of cosine similarities with given word pair similarity scores
  - ▶ Example: *SimLex-999*
    - ★ coast – shore: 9.00
    - ★ clothes – closet 1.96
  
- Downstream tasks
  - ▶ Embeddings as features in neural networks
  - ▶ Example:
    - ★ Machine translation (BLEU)
    - ★ sentiment analysis (accuracy)
    - ★ etc.

## Word Embeddings

- Vector representation
- Learning word embeddings

## Bilingual Word Embeddings

- **Motivation**
- Overview of approaches
- Mapping
- Orthogonal mapping
- Unsupervised training

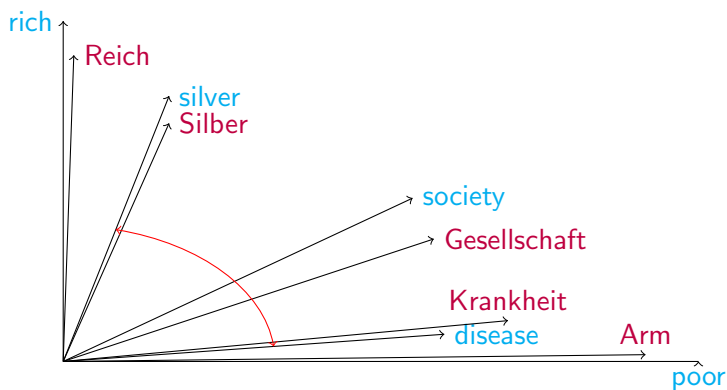
## Unsupervised Statistical MT

- MT with bilingual word embeddings
- Integrating embeddings into Moses

# Bilingual Word Spaces

Representation of words in two languages in same semantic space:

- Similar words are close to each other
- Given by cosine



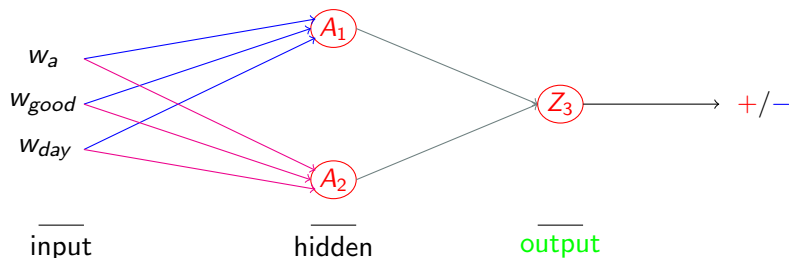
# Cross-lingual transfer learning

- Train on English

- ▶  $w_{good} = [0.23, 0.01, -0.45, 0.00, \dots, 0.93]$

- ▶  $w_{day} = [-0.76, 0.98, 0.23, 0.74, \dots, 0.01]$

- ▶  $w_a = [0.54, -0.39, 0.28, 0.79, \dots, 0.42]$



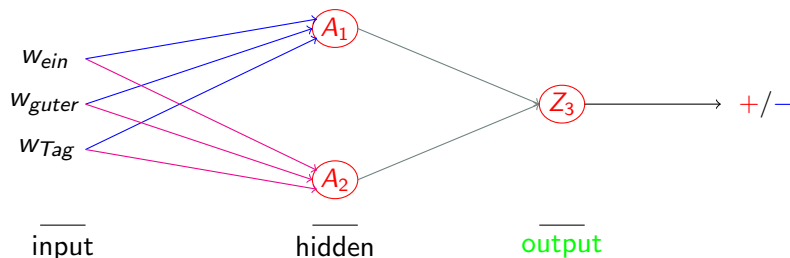
# Cross-lingual transfer learning

- Classify German

- $w_{good} = [0.23, 0.01, -0.45, 0.00, \dots, 0.93] \approx w_{guter}$

- $w_{day} = [-0.76, 0.98, 0.23, 0.74, \dots, 0.01] \approx w_{Tag}$

- $w_a = [0.54, -0.39, 0.28, 0.79, \dots, 0.42] \approx w_{ein}$

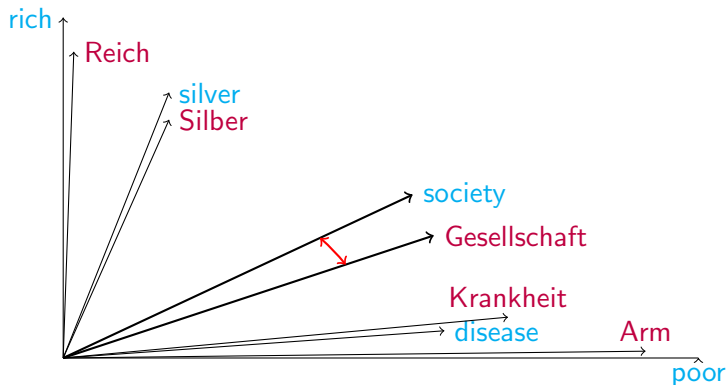




# Translating words

Translating word using cosine similarity

- **society** → **Gesellschaft**



## Word Embeddings

- Vector representation
- Learning word embeddings

## Bilingual Word Embeddings

- Motivation
- **Overview of approaches**
- Mapping
- Orthogonal mapping
- Unsupervised training

## Unsupervised Statistical MT

- MT with bilingual word embeddings
- Integrating embeddings into Moses

# Learning Bilingual Word Embeddings

- Required bilingual training signal:
  - ▶ Parallel sentences
    - ★ Hermann and Blunsom (2014), Gouws et al. (2015), Gouws and Søgaard (2015), Duong et al. (2016)

das	Haus	ist	klein
the	house	is	small

# Learning Bilingual Word Embeddings

- Required bilingual training signal:
  - ▶ Document-aligned data
    - ★ Vulic and Moens (2015); Vulic and Korhonen (2016)

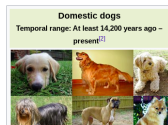
## Dog

From Wikipedia, the free encyclopedia

*This article is about the domestic dog. For related species known as "dogs", see [Canidae](#). For other uses, see [Dog \(disambiguation\)](#).*

The **dog** (*Canis familiaris* when considered a distinct species or *Canis lupus familiaris* when considered a subspecies of the wolf)<sup>[5]</sup> is a member of the genus *Canis* (canines), which forms part of the wolf-like *canids*,<sup>[6]</sup> and is the most widely abundant terrestrial *carnivore*.<sup>[7][8][9][10][11]</sup> The dog and the extant *gray wolf* are *sister taxa*<sup>[12][13][14]</sup> as modern wolves are not closely related to the wolves that were first *domesticated*.<sup>[13][14]</sup> which implies that the direct ancestor of the dog is *extinct*.<sup>[15]</sup> The dog was the first species to be domesticated,<sup>[14][16]</sup> and has been *selectively bred* over millennia for various behaviors, sensory capabilities, and physical attributes.<sup>[17]</sup>

Their *long association with humans* has led dogs to be uniquely attuned to human behavior<sup>[18]</sup> and they are able to thrive on a starch-rich diet that would be inadequate for other canids.<sup>[19]</sup> Dogs vary widely in shape, size and colors.<sup>[20]</sup> They perform many roles for humans, such as *hunting*, *herding*, *pulling loads*, *protection*, *assisting police and military*, *companionship* and, more recently, *aiding disabled people* and *therapeutic* roles. This influence on human society has given them the sobriquet of "man's best friend".



## Haushund

↪ Hund ist eine Weiterleitung auf diesen Artikel. Weitere Bedeutungen sind unter [Hund \(Begriffsklärung\)](#) aufgeführt.

Der **Haushund** (*Canis lupus familiaris*) ist ein *Haustier* und wird als *Halm-* und *Nutztier* gehalten. Seine wilde Stammform ist der *Wolf*, dem er als *Unterart* zugeordnet wird. Wann die *Domestizierung* stattfand, ist umstritten; wissenschaftliche Schätzungen variieren zwischen 15.000 und 100.000 Jahren v. u. Z.

Im engeren Sinn bezeichnet man als Haushund die Hunde, die überwiegend im Haus gehalten werden, und kennzeichnet damit also eine *Halbungsform*. Historisch wurde ein Hund, der zur Bewachung des Hauses gehalten wird, als Haushund bezeichnet.<sup>[1]</sup> Eine weitere Verwendung des Begriffs ist die Einschränkung auf *sozialisierte* (Haus-)Hunde, also Hunde, die an das Zusammenleben mit Menschen in der menschlichen Gesellschaft gewöhnt und an dieses angepasst sind. Damit wird der *Haushund* abgegrenzt gegen wild lebende, verwilderte oder streunende Hunde, die zwar auch domestiziert, aber nicht sozialisiert sind.<sup>[2]</sup>

Der *Dingo* ist ebenfalls ein Haushund, wird jedoch provisorisch als eigenständige Unterart des Wolfes geführt.<sup>[3]</sup>



# Learning Bilingual Word Embeddings

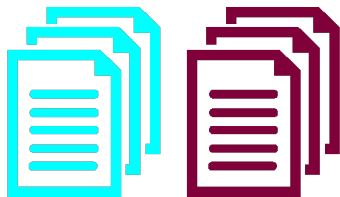
- Required bilingual training signal:
  - ▶ Monolingual data and a seed dictionary
    - ★ Mikolov et al. (2013b); Faruqui and Dyer (2014); Lazaridou et al. (2015)



```
dog    -- Hund
apple  -- Apfel
100    -- 100
...    -- ...
```

# Learning Bilingual Word Embeddings

- Required bilingual training signal:
  - ▶ Monolingual data only
    - ★ Conneau et al. (2017); Artetxe et al. (2018a)



# Learning Bilingual Word Embeddings

- Approaches:
  - ▶ Mapping
    - ★ Step 1: build source language embeddings
    - ★ Step 2: build target embeddings
    - ★ Step 3: map them to a shared space
  - ▶ Joint training
    - ★ Step 1: build both source and target embeddings in one step

## Word Embeddings

- Vector representation
- Learning word embeddings

## Bilingual Word Embeddings

- Motivation
- Overview of approaches
- **Mapping**
- Orthogonal mapping
- Unsupervised training

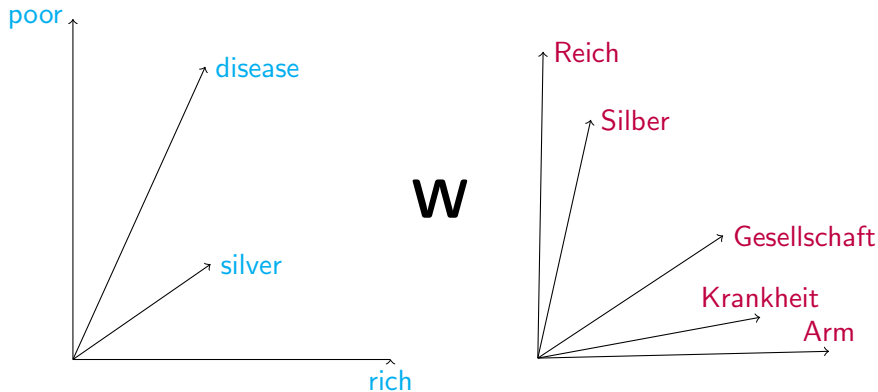
## Unsupervised Statistical MT

- MT with bilingual word embeddings
- Integrating embeddings into Moses



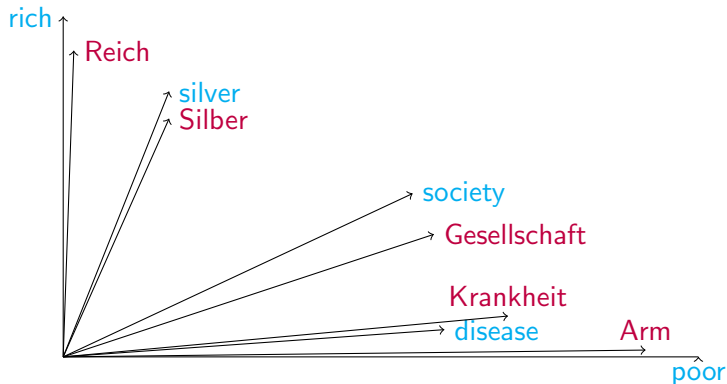
## Mapping with seed dictionary

- Learn monolingual word embeddings
- Learn a linear mapping  $W$



## Mapping with seed dictionary

- Project source words into target space



# Mapping with seed dictionary

1. Train **monolingual** word embeddings (Word2vec) in **English**
  - ▶ Need **English** monolingual data
2. Train **monolingual** word embeddings (Word2vec) in **German**
  - ▶ Need **German** monolingual data
3. Learn mapping **W** using a seed dictionary
  - ▶ Need a list of **5000 English words and their translation**

# Learning $W$ by minimizing Euclidean distance

Regression (Mikolov et al. (2013b))

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \sum_i^n \| \mathbf{W} \mathbf{x}_i - \mathbf{y}_i \|^2$$

$\mathbf{x}_i$  : **embedding** of i-th **source** (English) word in the seed dictionary.

$\mathbf{y}_i$  : **embedding** of i-th **target** (German) word in the seed dictionary.

# Learning $W$ by minimizing Euclidean distance

Regression (Mikolov et al. (2013b))

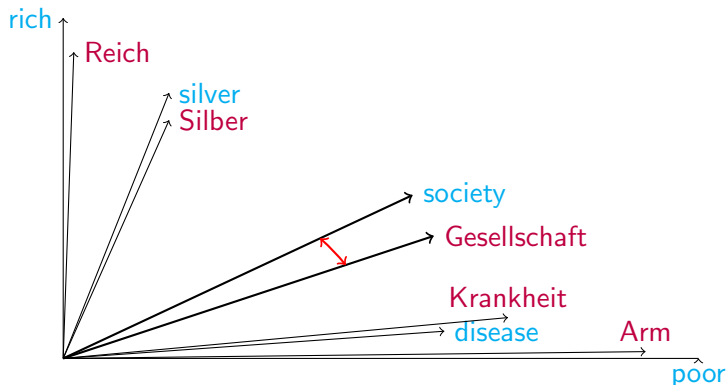
$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \sum_i^n \| \mathbf{W}\mathbf{x}_i - \mathbf{y}_i \|^2$$

- For each pair  $(x_i, y_i)$  in the seed dictionary:
  - ▶ Predict projection  $\mathbf{y}^*$  by computing  $\mathbf{W}\mathbf{x}_i$
  - ▶ Compute **squared error** between  $\mathbf{y}^*$  and  $\mathbf{y}_i$
- Find  $\mathbf{W}$  such that squared error over training set is minimal

# Bilingual dictionary induction

- Task to evaluate bilingual word embeddings intrinsically
- Given a set of source words, find the corresponding translations:
  - ▶ Given **society**, find its vector in the BWE space
  - ▶ Retrieve the **German** word whose vector is the most similar (cosine)

$$\text{cosine}(\mathbf{W}_x, y)$$



# Bilingual dictionary induction

- Evaluation: *precision@n*
  - ▶ is the correct translation in the  $n$  most similar translations?
  - ▶ dog → Katze, Hunde, **Hund**, Giraffe, Maus

	English to Italian			Italian to English		
	P@1	P@5	P@10	P@1	P@5	P@10
<hr/> <i>Methods with cross-lingual supervision (WaCky)</i> <hr/>						
Mikolov et al. (2013b) †	33.8	48.3	53.9	24.9	41.0	47.4

Conneau et al. (2017)

## Word Embeddings

- Vector representation
- Learning word embeddings

## Bilingual Word Embeddings

- Motivation
- Overview of approaches
- Mapping
- **Orthogonal mapping**
- Unsupervised training

## Unsupervised Statistical MT

- MT with bilingual word embeddings
- Integrating embeddings into Moses



# Isomorphism

- Can we improve the approach?
- Isomorphism (approximate)
  - ▶ If languages convey similar information in similar contexts their monolingual embeddings should be isomorphic.
  - ▶ it is only true to some extent for each language pair
- restrictions on the mapping ( $\mathbf{W}$ )
  - ▶ e.g. allow rotation only
  - ▶ preserves word similarities in the monolingual spaces



(Conneau et al., 2017)



# Vector normalization

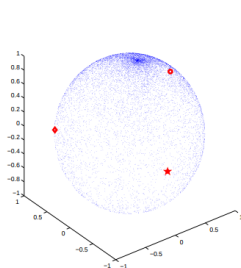
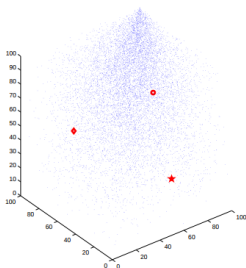
- Mismatch:

- ▶ Training (squared-error):  $\mathbf{W}^* = \arg \min_{\mathbf{W}} \sum_i^n \| \mathbf{W}\mathbf{x}_i - \mathbf{y}_i \|^2$

- ▶ Test:  $\text{cosine}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$

- Normalize vectors to length 1:  $x = \frac{x}{\|x\|}$

- ▶  $\mathbf{W}^* = \arg \max_{\mathbf{W}} \sum_i^n (\mathbf{W}\mathbf{x}_i)^T \mathbf{y}_i$



# Procrustes problem

- Orthogonality constraint on  $\mathbf{W}$
- Xing et al. (2015)

$$\mathbf{W}^* = UV^T \quad U\Sigma V^T = SVD(\mathbf{YX}^T)$$

$\mathbf{X}$  : contains the embedding of the  $i^{th}$  source word ( $x_i$ ) in the seed dictionary in row  $i$ .

$\mathbf{Y}$  : contains the embedding of the  $i^{th}$  target word ( $y_i$ ) in the seed dictionary in row  $i$ .

# Bilingual dictionary induction

	English to Italian			Italian to English		
	P@1	P@5	P@10	P@1	P@5	P@10
<hr/> <i>Methods with cross-lingual supervision (WaCky)</i> <hr/>						
Mikolov et al. (2013b) †	33.8	48.3	53.9	24.9	41.0	47.4
Procrustes - CSLS	44.9	61.8	66.6	38.5	57.2	63.0

Conneau et al. (2017)

\* Cross-Domain Similarity Local Scaling (CSLS): cosine alternative Conneau et al. (2017)

## Word Embeddings

- Vector representation
- Learning word embeddings

## Bilingual Word Embeddings

- Motivation
- Overview of approaches
- Mapping
- Orthogonal mapping
- **Unsupervised training**

## Unsupervised Statistical MT

- MT with bilingual word embeddings
- Integrating embeddings into Moses

# Unsupervised mapping

- Low-resource languages
  - ▶ No seed dictionary is available for training
  - ▶ only monolingual corpora for both languages

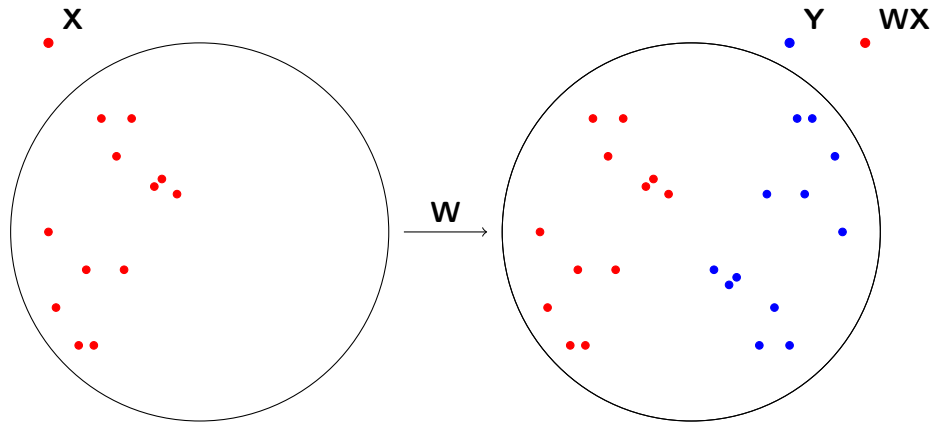
0. Generate an initial seed dictionary automatically
1. Learn mapping
2. Induce a better dictionary
3. Goto 1. until convergence



# Adversarial training

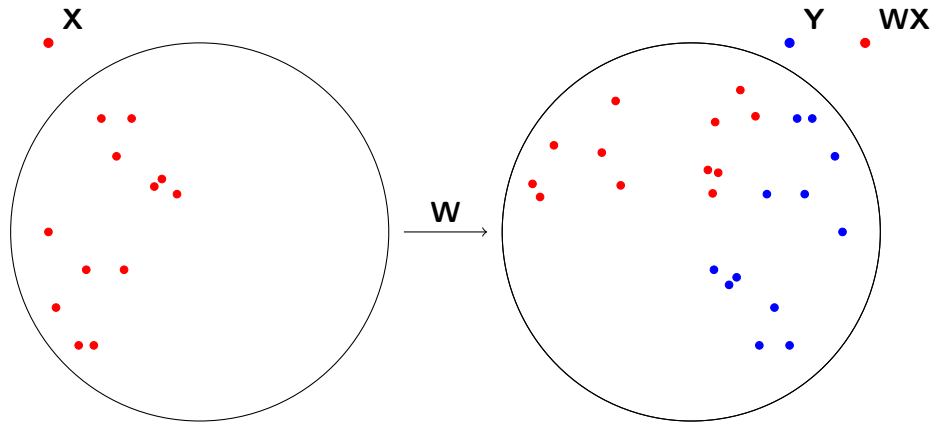
- Conneau et al. (2017)
- Two player game
  - ▶ **Discriminator**: discriminate mapped source language vectors  $\mathbf{W}x_i$  from target  $y_i$ 
    - ★  $x_i$  and  $y_i$  are the embeddings of any word in the vocabulary (we have no seed dictionary)
  - ▶ **Generator**: generate the mapping  $\mathbf{W}$  such that the discriminator fails
- Iterative process: both make a step after each other

# Adversarial training

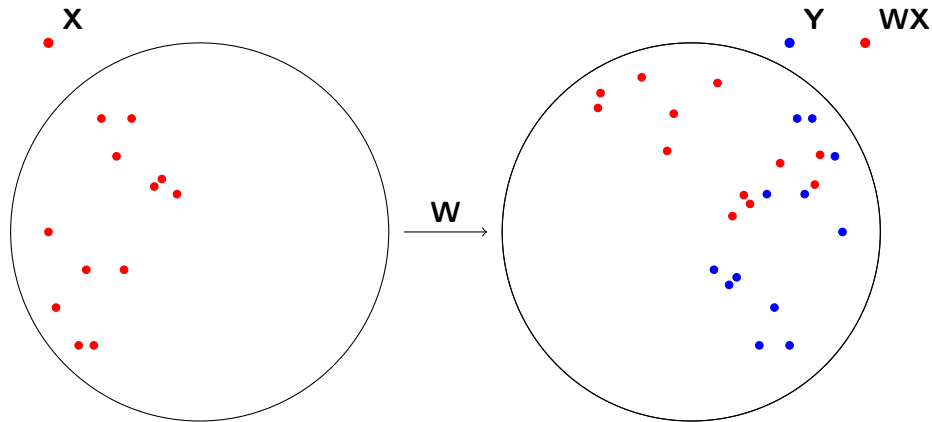




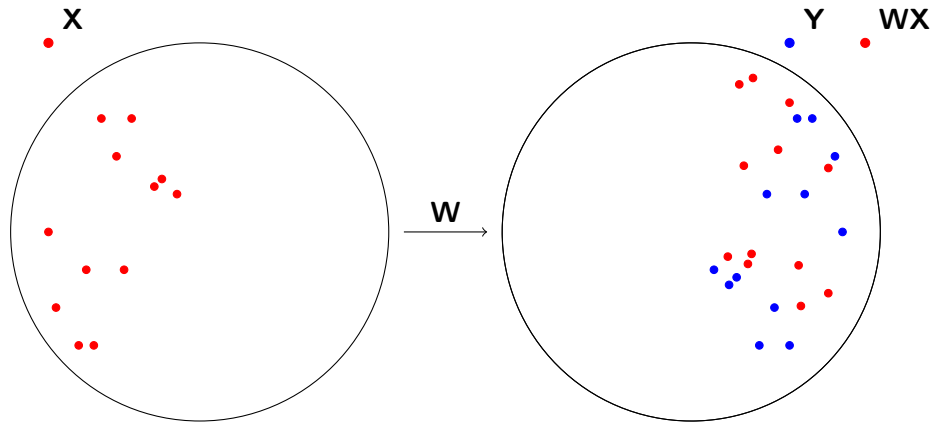
# Adversarial training



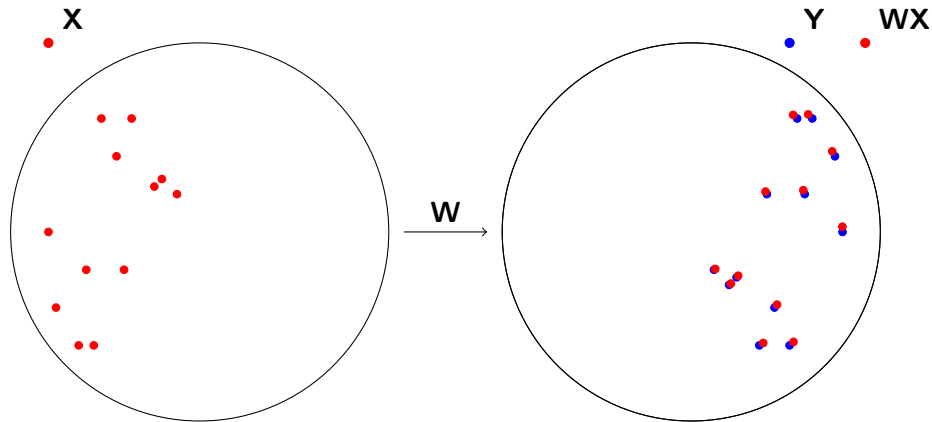
# Adversarial training



# Adversarial training



# Adversarial training



# Adversarial training

- Discriminator

- ▶ Feed forward network: probability of vector  $x$  is the embedding of a source word  $P_{\theta_D}(\text{source} = 1|x)$

$$\mathcal{L}_D(\theta_D|W) = -\frac{1}{n} \sum_{i=1}^n \log P_{\theta_D}(\text{source} = 1|Wx_i) - \frac{1}{m} \sum_{i=1}^m \log P_{\theta_D}(\text{source} = 0|y_i)$$

# Adversarial training

- Discriminator

- ▶ Feed forward network: probability of vector  $x$  is the embedding of a source word  $P_{\theta_D}(\text{source} = 1|x)$

$$\mathcal{L}_D(\theta_D|W) = -\frac{1}{n} \sum_{i=1}^n \log P_{\theta_D}(\text{source} = 1|Wx_i) - \frac{1}{m} \sum_{i=1}^m \log P_{\theta_D}(\text{source} = 0|y_i)$$

- Generator

$$\mathcal{L}_W(W|\theta_D) = -\frac{1}{n} \sum_{i=1}^n \log P_{\theta_D}(\text{source} = 0|Wx_i) - \frac{1}{m} \sum_{i=1}^m \log P_{\theta_D}(\text{source} = 1|y_i)$$

# Adversarial training

- Discriminator

- ▶ Feed forward network: probability of vector  $x$  is the embedding of a source word  $P_{\theta_D}(\text{source} = 1|x)$

$$\mathcal{L}_D(\theta_D|W) = -\frac{1}{n} \sum_{i=1}^n \log P_{\theta_D}(\text{source} = 1|Wx_i) - \frac{1}{m} \sum_{i=1}^m \log P_{\theta_D}(\text{source} = 0|y_i)$$

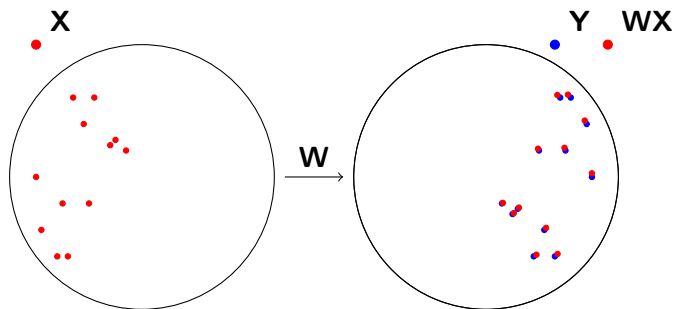
- Generator

$$\mathcal{L}_W(W|\theta_D) = -\frac{1}{n} \sum_{i=1}^n \log P_{\theta_D}(\text{source} = 0|Wx_i) - \frac{1}{m} \sum_{i=1}^m \log P_{\theta_D}(\text{source} = 1|y_i)$$

- Minimize both losses with gradient descent

# Initial dictionary

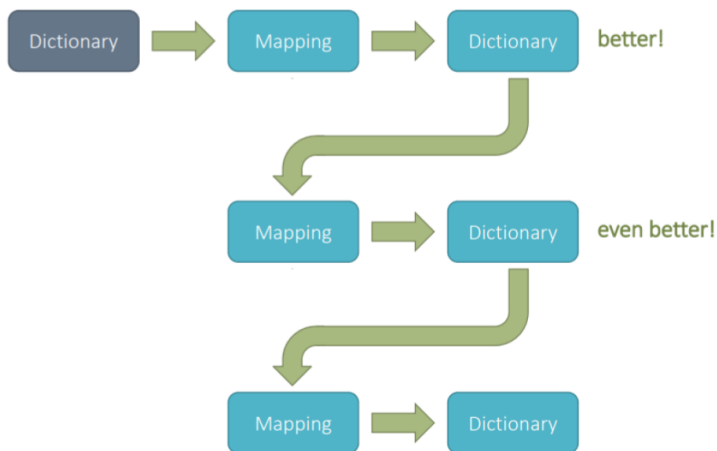
- The iterative process gives us an initial  $W$
- Noisy but works well for frequent words
- Translate frequent words with bilingual dictionary induction
- Use this dictionary for orthogonal mapping





# Iterative refinement

- Increase the quality and size of the dictionary in each step



Ruder et al. (2019)

# Bilingual dictionary induction

	English to Italian			Italian to English		
	P@1	P@5	P@10	P@1	P@5	P@10
<hr/> <i>Methods with cross-lingual supervision (WaCky)</i> <hr/>						
Mikolov et al. (2013b) <sup>†</sup>	33.8	48.3	53.9	24.9	41.0	47.4
Procrustes - CSLS	44.9	61.8	66.6	38.5	57.2	63.0
<hr/> <i>Methods without cross-lingual supervision (WaCky)</i> <hr/>						
Adv - Refine - CSLS	45.1	60.7	65.1	38.3	57.8	62.8

Conneau et al. (2017)

# Bilingual dictionary induction

- Results are lower for distant language pairs
- Isomorphism is weaker

	<b>Unsupervised (Adversarial)</b>	<b>Supervised (Identical)</b>
EN-ES	81.89	<b>82.62</b>
EN-ET	00.00	<b>31.45</b>
EN-FI	00.09	<b>28.01</b>
EN-EL	00.07	<b>42.96</b>
EN-HU	45.06	<b>46.56</b>
EN-PL	46.83	<b>52.63</b>
EN-TR	32.71	<b>39.22</b>

Søgaard et al. (2018)

## Word Embeddings

- Vector representation
- Learning word embeddings

## Bilingual Word Embeddings

- Motivation
- Overview of approaches
- Mapping
- Orthogonal mapping
- Unsupervised training

## Unsupervised Statistical MT

- **MT with bilingual word embeddings**
- Integrating embeddings into Moses

# Unsupervised MT

- Most MT systems need a large amount of parallel sentences (few millions)
- Low-resource languages lack this resource
- We use only monolingual data to build an MT system
- Idea:
  - ▶ Start simple: word translation
  - ▶ Extend the initial system to sentences
- We don't expect to beat supervised systems but for many languages this is the best we can do

# Word-by-word translation

- Translate each word in a sentence independently of the others
- Bilingual dictionary induction
  - ▶ unsupervised bilingual word embeddings
  - ▶ cosine similarity of words

# Word-by-word translation

- Translate each word in a sentence independently of the others
- Bilingual dictionary induction
  - ▶ unsupervised bilingual word embeddings
  - ▶ cosine similarity of words

Der Himmel ist blau .

The sky is blue .

# Word-by-word translation

- Translate each word in a sentence independently of the others
- Bilingual dictionary induction
  - ▶ unsupervised bilingual word embeddings
  - ▶ cosine similarity of words

**Der** Himmel ist blau .

**The** sky is blue .



# Word-by-word translation

- Translate each word in a sentence independently of the others
- Bilingual dictionary induction
  - ▶ unsupervised bilingual word embeddings
  - ▶ cosine similarity of words

Der **Himmel** ist blau .

The **sky** is blue .

# Word-by-word translation

- Translate each word in a sentence independently of the others
- Bilingual dictionary induction
  - ▶ unsupervised bilingual word embeddings
  - ▶ cosine similarity of words

Der Himmel **ist** blau .

The sky **is** blue .

# Word-by-word translation

- Translate each word in a sentence independently of the others
- Bilingual dictionary induction
  - ▶ unsupervised bilingual word embeddings
  - ▶ cosine similarity of words

Der Himmel ist **blau** .

The sky is **blue** .

# Word-by-word translation

- Translate each word in a sentence independently of the others
- Bilingual dictionary induction
  - ▶ unsupervised bilingual word embeddings
  - ▶ cosine similarity of words

Der Himmel ist blau .

The sky is blue .

# Results

	en-fr	fr-en	de-en	en-de
Supervised	27.97	26.13	25.61	21.33
word-by-word	6.28	10.09	10.77	7.06

Lample et al. (2018)

# Problems

- Translating compound words

Der Himmel ist **dunkelblau** .

The sky is ~~**dark**~~ .

The sky is ~~**blue**~~ .

The sky is **\*dark blue\*** .

# Problems

- Translating multi word expressions

Fix            und            fertig  
~~Fixed~~      ~~and~~      ~~ready~~  
**Exhausted**

# Phrase embeddings

- One embedding for n-grams

- ▶  $w_{dark\_blue} = [0.23, 0.01, -0.45, 0.00, \dots, 0.93]$

1. look for frequently co-occurring n-grams in the monolingual corpus

(Mikolov et al. (2013c))

$$score(w_i, w_j) = \frac{count(w_i, w_j)}{count(w_i) \times count(w_j)}$$

2. Concatenate n-grams in the corpus

- ▶ The sky is **dark\_blue** .

3. Train monolingual embeddings as before

4. Mapping of monolingual spaces



# Phrase translation

- Translate tokens with cosine
- Tokens can be phrases as well

Der Himmel ist **dunkelblau** .

- ▶  $\text{cosine}(\text{dunkelblau}, \text{dark}) = 0.7$
- ▶  $\text{cosine}(\text{dunkelblau}, \text{blue}) = 0.75$
- ▶  $\text{cosine}(\text{dunkelblau}, \text{dark\_blue}) = 0.83$

fix und fertig

- ▶  $\text{cosine}(\text{fix}, \text{fixed}) = 0.68$
- ▶  $\text{cosine}(\text{und}, \text{and}) = 0.8$
- ▶  $\text{cosine}(\text{fertig}, \text{ready}) = 0.7$
- ▶  $\text{cosine}(\text{fix\_und\_fertig}, \text{exhausted}) = 0.95$

## Word Embeddings

- Vector representation
- Learning word embeddings

## Bilingual Word Embeddings

- Motivation
- Overview of approaches
- Mapping
- Orthogonal mapping
- Unsupervised training

## Unsupervised Statistical MT

- MT with bilingual word embeddings
- Integrating embeddings into Moses

# Problems

- Translating words independently leads to problems:
  - ▶ Fluency/word order

Ich	denke	dass	der	Himmel	blau	ist
I	think	that	the	sky	<del>blue</del>	<del>is</del>
I	think	that	the	sky	<b>is</b>	<b>blue</b>

# Problems

- Translating words independently leads to problems:
  - ▶ Multi sense words and morphology

Ich	sitze	auf	der	Bank
<del>I</del>	<del>sit</del>	on	the	<del>bank</del>
<b>I'm</b>	<b>sitting</b>	on	the	<b>bench</b>

# Log-linear model

- Feature functions:

- ▶ Phrase-table → adequacy

- ★ Generate  $n$  (100) most probable translations for each source word/phrase with bilingual dictionary induction
- different morphological variations and senses

$$\phi(f|e) = \frac{e^{\cos(e,f)}}{\sum_{\hat{f}} e^{\cos(e,\hat{f})}}$$

- ▶ Language model → fluency

- right sense and morphology of words given the source sentence
- right order of words
- ★ Train the same way as for supervised SMT
- ★ using the same monolingual corpus as for the embeddings

# Tuning weights

- Finding the right feature weights needs parallel data
- Back-translation:
  0. Build systems for both directions:  $S_{L_{trg} \rightarrow L_{src}}$  and  $S_{L_{src} \rightarrow L_{trg}}$ 
    - ★ uniform feature weights
  1. Generate synthetic parallel data
    - ★  $D_{L_{src} \rightarrow L_{trg}}$ : (back-)translate  $L_{trg}$  monolingual data with  $S_{L_{trg} \rightarrow L_{src}}$
    - ★  $D_{L_{trg} \rightarrow L_{src}}$ : (back-)translate  $L_{src}$  monolingual data with  $S_{L_{src} \rightarrow L_{trg}}$
    - ★ Source language is noisy but target is not
  2. Optimize weight with MERT
    - ★  $D_{L_{src} \rightarrow L_{trg}}$  for  $S_{L_{src} \rightarrow L_{trg}}$
    - ★  $D_{L_{trg} \rightarrow L_{src}}$  for  $S_{L_{trg} \rightarrow L_{src}}$
  3. Goto 1. until convergence

# Iterative back-translation

- Only the weights are tuned in the previous step
- Build SMT systems from scratch using synthetic parallel data
  - ▶ word align sentences
  - ▶ build phrase tables
  - ▶ etc.
- Weights of  $S_{L_{src} \rightarrow L_{trg}}$  and  $S_{L_{trg} \rightarrow L_{src}}$  are now tuned
- Back-translation:
  1. Generate synthetic parallel data
    - ★  $D_{L_{trg} \rightarrow L_{src}}$ : (back-)translate  $L_{src}$  monolingual data with  $S_{L_{src} \rightarrow L_{trg}}$
    - ★  $D_{L_{src} \rightarrow L_{trg}}$ : (back-)translate  $L_{trg}$  monolingual data with  $S_{L_{trg} \rightarrow L_{src}}$
  2. Build **supervised** SMT systems from scratch
    - ★  $D_{L_{src} \rightarrow L_{trg}}$  for  $S_{L_{src} \rightarrow L_{trg}}$
    - ★  $D_{L_{trg} \rightarrow L_{src}}$  for  $S_{L_{trg} \rightarrow L_{src}}$
  3. Goto 1. until convergence

# Results

	en-fr	fr-en	de-en	en-de
Supervised	27.97	26.13	25.61	21.33
word-by-word	6.28	10.09	10.77	7.06

Lample et al. (2018)

	<b>FR-EN</b>	<b>EN-FR</b>	<b>DE-EN</b>	<b>EN-DE</b>
Proposed system	<b>25.87</b>	<b>26.22</b>	<b>17.43</b>	<b>14.08</b>

Artetxe et al. (2018b)



# Summary

- Word embeddings
  - ▶ Similar words have similar word vectors
  - ▶ Training them with *word2vec* using monolingual data
- Bilingual word embeddings
  - ▶ Mapping monolingual embeddings to shared space using seed dictionary
  - ▶ Unsupervised mapping with adversarial initialization
- Unsupervised SMT
  - ▶ Word-by-word translation with unsupervised bilingual embeddings
  - ▶ Integrating bilingual dictionary induction to SMT

Thank you !

## References I

- Artetxe, M., Labaka, G., and Agirre, E. (2018a). A robust self-learning method for fully unsupervised cross-lingual mappings of word embeddings. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 789—798.
- Artetxe, M., Labaka, G., and Agirre, E. (2018b). Unsupervised Statistical Machine Translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3632—3642.
- Conneau, A., Lample, G., Ranzato, M., Denoyer, L., and Jégou, H. (2017). Word translation without parallel data. *arXiv preprint arXiv:1710.04087*.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

## References II

- Duong, L., Kanayama, H., Ma, T., Bird, S., and Cohn, T. (2016). Learning crosslingual word embeddings without bilingual corpora. In *Proc. EMNLP*.
- Faruqui, M. and Dyer, C. (2014). Improving vector space word representations using multilingual correlation. In *Proc. EACL*.
- Gouws, S., Bengio, Y., and Corrado, G. (2015). Bilbowa: Fast bilingual distributed representations without word alignments. In *Proc. ICML*.
- Gouws, S. and Søgaard, A. (2015). Simple task-specific bilingual word embeddings. In *Proc. NAACL*.
- Hermann, K. M. and Blunsom, P. (2014). Multilingual models for compositional distributed semantics. In *Proc. ACL*, pages 58–68, Baltimore, Maryland. Association for Computational Linguistics.
- Lample, G., Denoyer, L., and Ranzato, M. (2018). Unsupervised Machine Translation Using Monolingual Corpora Only. In *Proceedings of the International Conference on Learning Representations*, pages 1–12.

## References III

- Lazaridou, A., Dinu, G., and Baroni, M. (2015). Hubness and pollution: Delving into cross-space mapping for zero-shot learning. In *Proc. ACL*.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. In *Proceedings of Workshop at ICLR*.
- Mikolov, T., Le, Q. V., and Sutskever, I. (2013b). Exploiting similarities among languages for machine translation. *CoRR*, abs/1309.4168.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013c). Distributed Representations of Words and Phrases and their Compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119.
- Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. (2018). Deep contextualized word representations. In *Proc. of NAACL*.

## References IV

- Ruder, S., Søgaard, A., and Vulić, I. (2019). Unsupervised cross-lingual representation learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, pages 31–38.
- Søgaard, A., Ruder, S., and Vulić, I. (2018). On the Limitations of Unsupervised Bilingual Dictionary Induction. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 778–788.
- Vulic, I. and Korhonen, A. (2016). On the Role of Seed Lexicons in Learning Bilingual Word Embeddings. In *Proc. ACL*, pages 247–257.
- Vulic, I. and Moens, M. (2015). Bilingual word embeddings from non-parallel document-aligned data applied to bilingual lexicon induction. In *Proc. ACL*.
- Xing, C., Wang, D., Liu, C., and Lin, Y. (2015). Normalized word embedding and orthogonal transform for bilingual word translation. In *Proc. NAACL*.