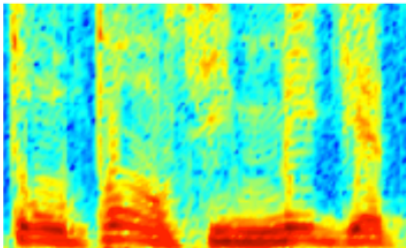




# Introduction

## AUDIO



Audio Spectrogram

DENSE

## IMAGES

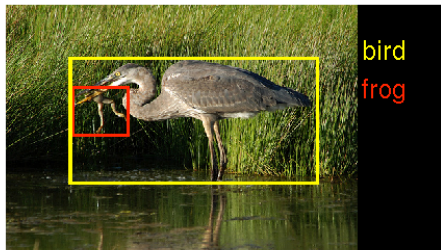


Image pixels

DENSE

## TEXT

0	0	0	0.2	0	0.7	0	0	0	...	...
---	---	---	-----	---	-----	---	---	---	-----	-----

Word, context, or  
document vectors

SPARSE

## Motivation

Monolingual word embeddings

The NLP community has developed good features for several tasks, but finding

- **task-invariant** (POS tagging, NER, SRL); **AND**
- **language-invariant** (English, Danish, Afrikaans,...)

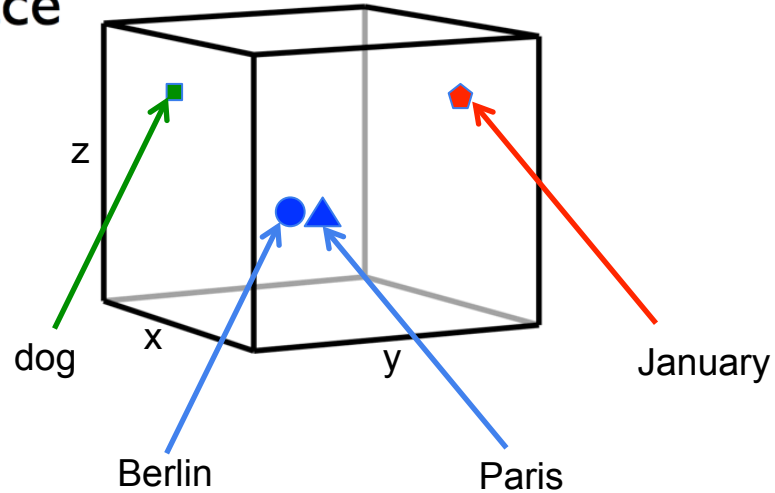
Cross-lingual word embeddings (this talk)

features is non-trivial and time-consuming (been trying for 20+ years).

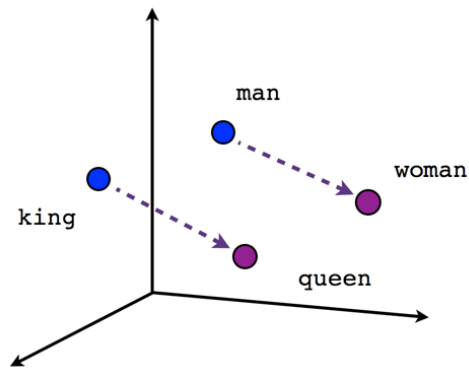
Learn word-level features which generalize across tasks *and* languages.

# Word Embeddings

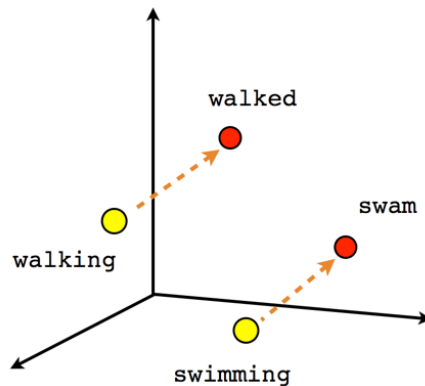
3D embedding  
space



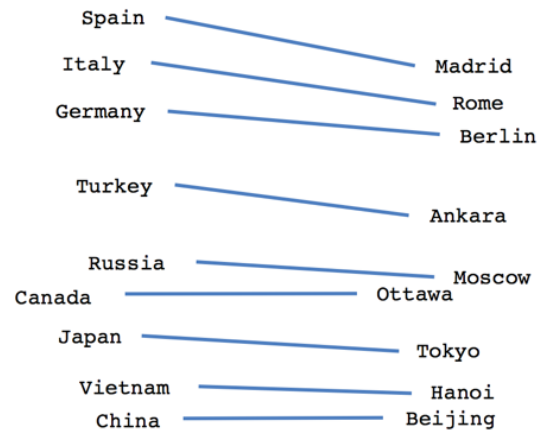
# Word Embeddings Capture Interesting, Universal Features



Male-Female



Verb tense

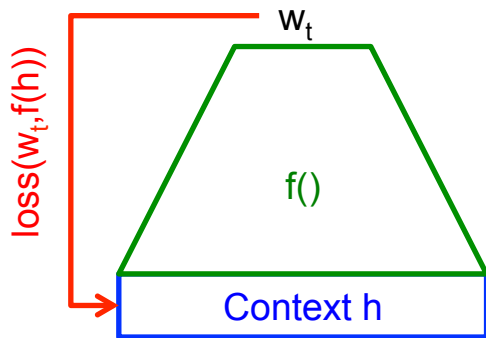


Country-Capital

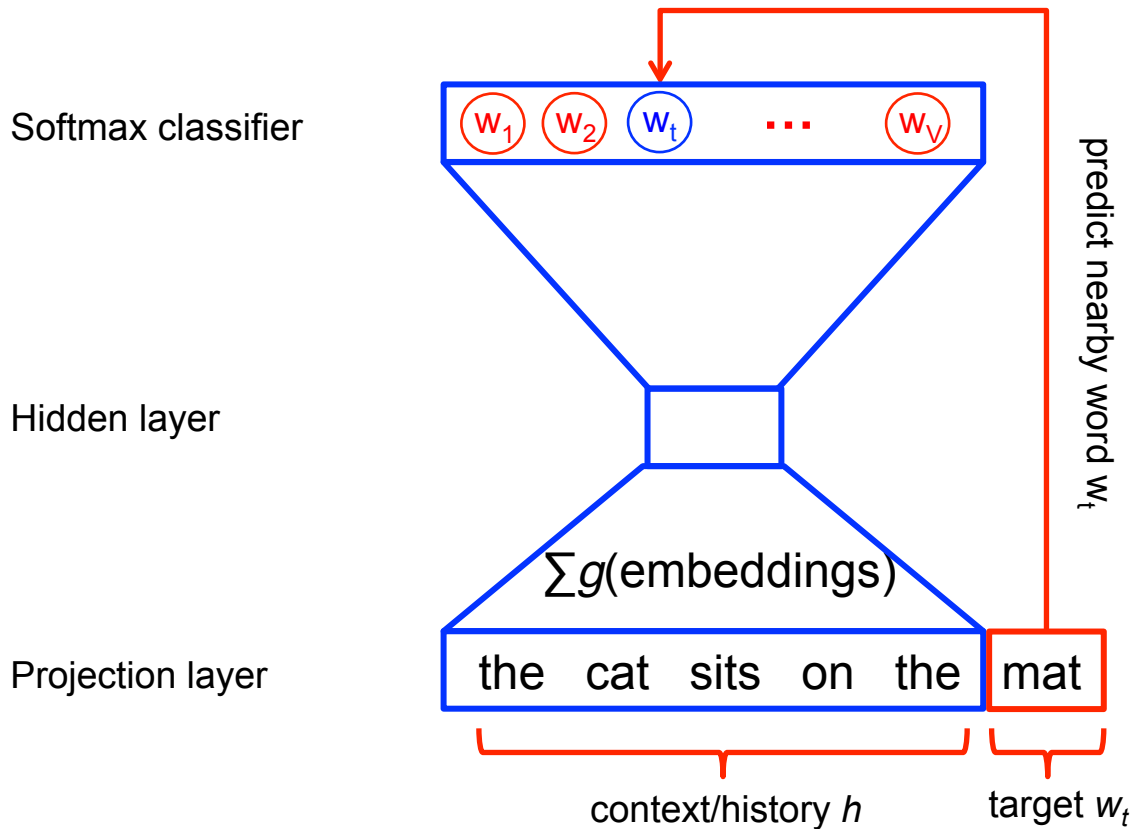
# Word Embedding Models

Models differ based on:

1. How they compute the **context  $h$** 
  - positional / convolutional / bag-of-words
2. How they map context to target  $w_t = f(h)$ 
  - linear, bilinear, non-linear
3. How they measure  $\text{loss}_R(w_t, f(h))$  and how they're trained
  - **language models (NLL, ...)**
  - **word embeddings:** negative sampling (CBOW/ Skipgram), sampled rank loss (Collobert+Weston), squared-error (LSA)



# Learning Word Embeddings: Pre-2008



$$Pr(w_t) = \frac{\text{score}(w_t; \theta)}{\sum_{j=1}^V \text{score}(w_j; \theta)}$$

(Bengio *et al.*, 2003)

## Advances in Learning Word Embeddings

Not interested in language modelling (for which we need normalized probabilities), so we don't need the expensive softmax. Can use much faster

- **hierarchical softmax** (Morin + Bengio, 2005),
- **sampled rank loss** (Collobert + Weston, 2008),
- **noise-contrastive estimation** (Mnih + Teh, 2012)
- **negative sampling** (Mikolov et al., 2013)

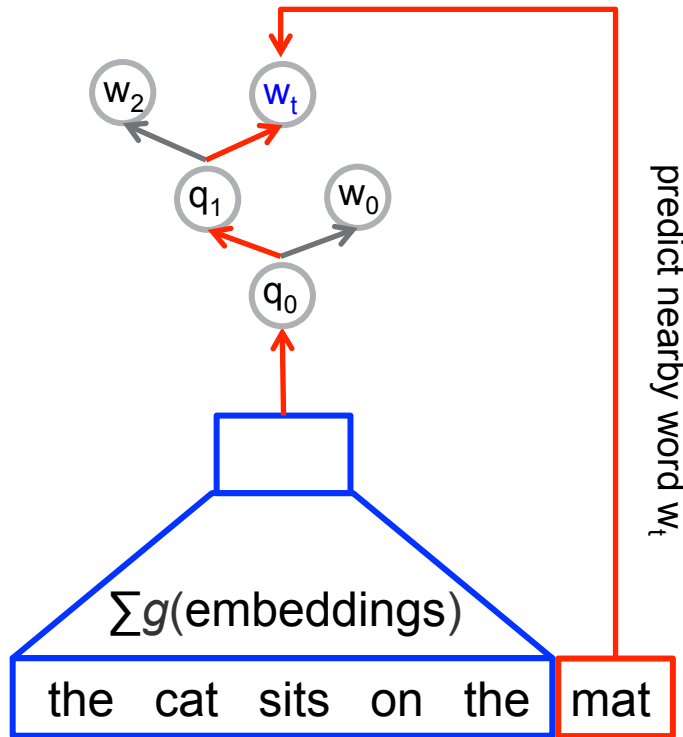


# Learning Word Embeddings: Hierarchical Softmax

Hierarchical  
Softmax classifier

Hidden layer

Projection layer

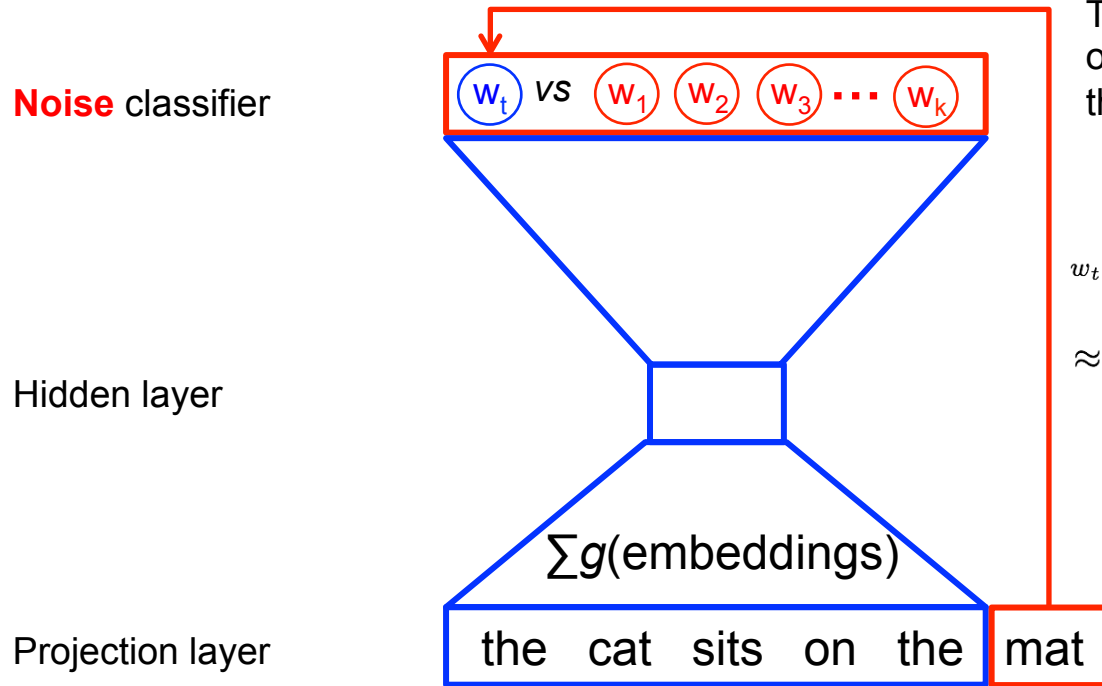


$$P(w_t|h) = \prod_{i \in L(w_t)} P(q_i|h)$$

Significant savings since  $|L(w)| \ll V$

(Morin + Bengio, 2005;  
Mikolov *et al.*, 2013)

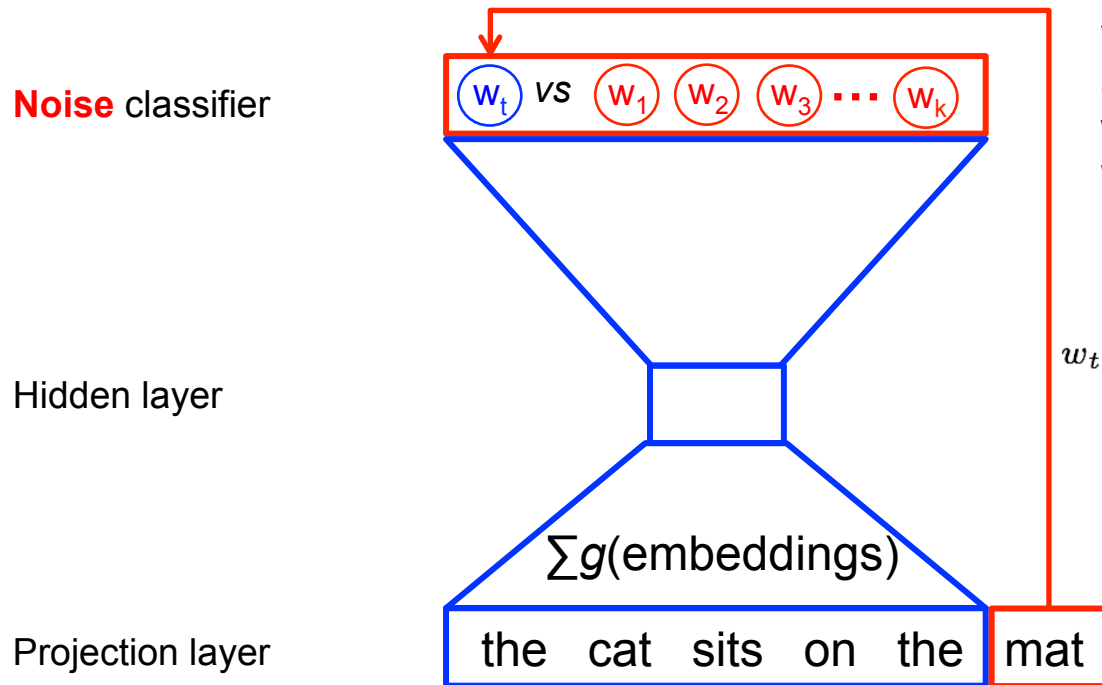
# Learning Word Embeddings: **Sampled rank loss**



Train a **non-probabilistic** model to **rank** an observed word  $w_t \sim P_{\text{data}}$  some margin higher than  $k$  ( $\ll V$ ) sampled noise words  $w_{\text{noise}} \sim P_{\text{noise}}$

$$\sum_{w_t, h \sim \text{data}} \sum_{j=1}^V \max\{0, 1 - (s_{\theta}(w_t; h) - s_{\theta}(w_j, h))\}$$
$$\approx \sum_{w_t, h \sim \text{data}} \sum_{j \sim P(w)}^k \max\{0, 1 - (\underbrace{s_{\theta}(w_t; h)}_{\text{increase}} - \underbrace{s_{\theta}(w_j, h)}_{\text{decrease}})\}$$

# Learning Word Embeddings: Noise-contrastive Estimation



Train a **probabilistic** model  $P(w|h)$  to be able to discriminate an observed nearby word  $w_t \sim P_{\text{data}}$  from sampled noise words  $w_{\text{noise}} \sim P_{\text{noise}}$

$$\sum_{w_t, h \sim P_{\text{data}}} P(w_t|h) + \sum_{j \sim P_{\text{noise}}} (1 - P(w_j|h))$$

# Neural Word Embeddings

Why do “similar” words have similar embeddings?



J.R. Firth

Citizens of  $\left\{ \begin{array}{c} \text{France} \\ \text{Denmark} \\ \dots \\ \text{Sweden} \end{array} \right\}$  protested today

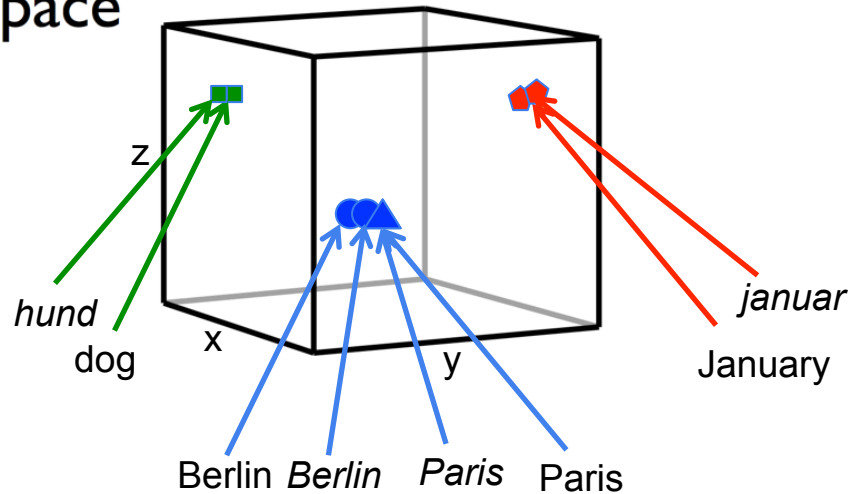
All training objectives have the form:

$$\begin{aligned} \min_{\mathbf{R}} \sum_i J(w_t^{(i)}, h^{(i)}) \\ = \min_{\mathbf{R}} \sum_i \text{distance}(w_t^{(i)}, h^{(i)}) \end{aligned}$$

I.e. for a **fixed context**, all **distributionally similar words** will get updated towards a common point.

# Cross-lingual Word Embeddings

3D embedding  
space



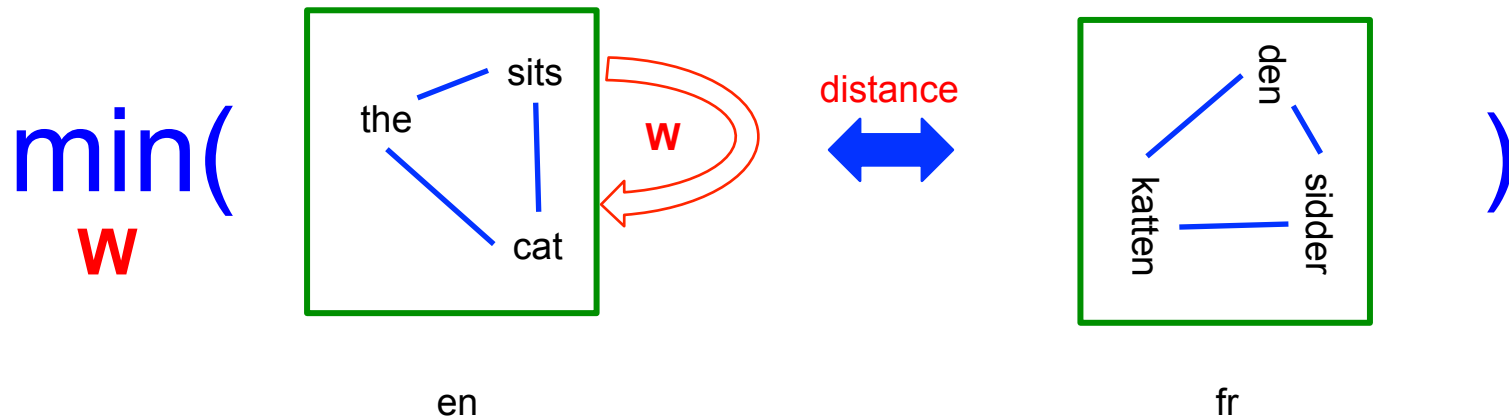
We want to learn an alignment between the two embedding spaces s.t. translation pairs are close.

# Learning **Cross-lingual** Word Embeddings: Approaches

1. Align pre-trained embeddings (**offline**)
2. Jointly learn and align embeddings (**online**) using parallel-only data
3. Jointly learn and align embeddings (**online**) using monolingual **and** parallel data

# Learning Cross-lingual Word Embeddings I

Offline methods: “Translation Matrix”



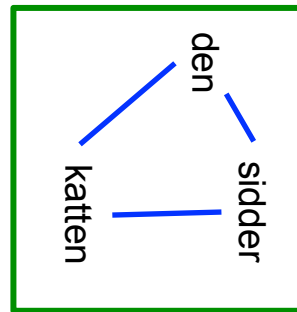
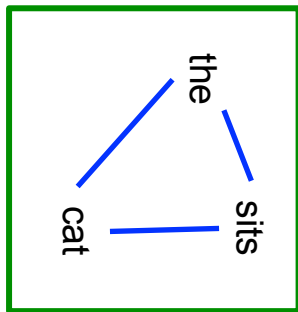
Learn  $\mathbf{W}$  to transform the pre-trained English embeddings into a space where the distance between a word and its translation-pair is minimized

$$\min_{\mathbf{W}} ||\mathbf{R}^{en}\mathbf{W} - \mathbf{R}^{fr}||^2$$

(Mikolov et al., 2013)

# Learning Cross-lingual Word Embeddings I

## Offline methods



en x **W**

fr

Transformed embeddings

Learn **W** to transform the pre-trained English embeddings into a space where the distance between a word and its translation-pair is minimized

Can also learn a *separate* **W** for *en* and *fr* using **Multilingual CCA**  
(Faruqui et al, 2014)

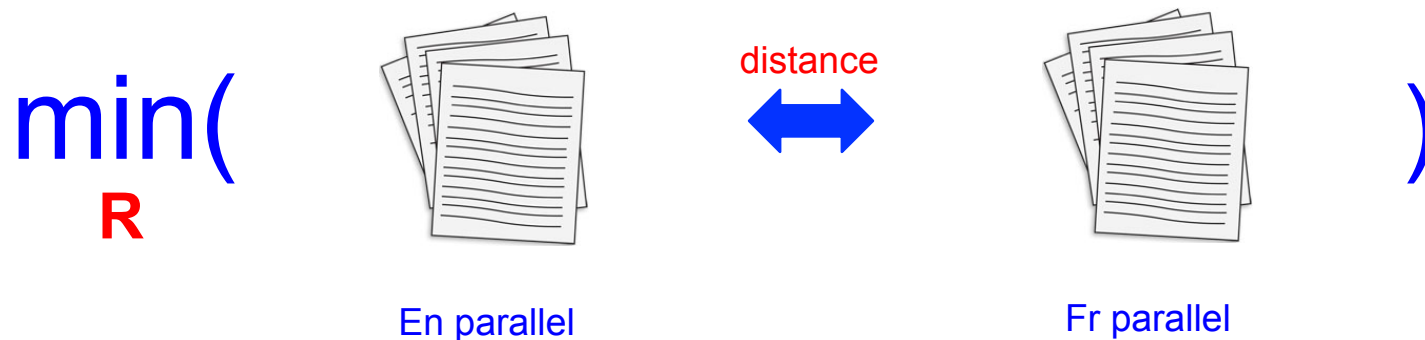
$$\min_{\mathbf{W}} ||\mathbf{R}^{en} \mathbf{W} - \mathbf{R}^{fr}||^2$$

(Mikolov et al., 2013)



# Learning **Cross-lingual** Word Embeddings II

## Parallel-only methods



**Bilingual Auto-encoders**

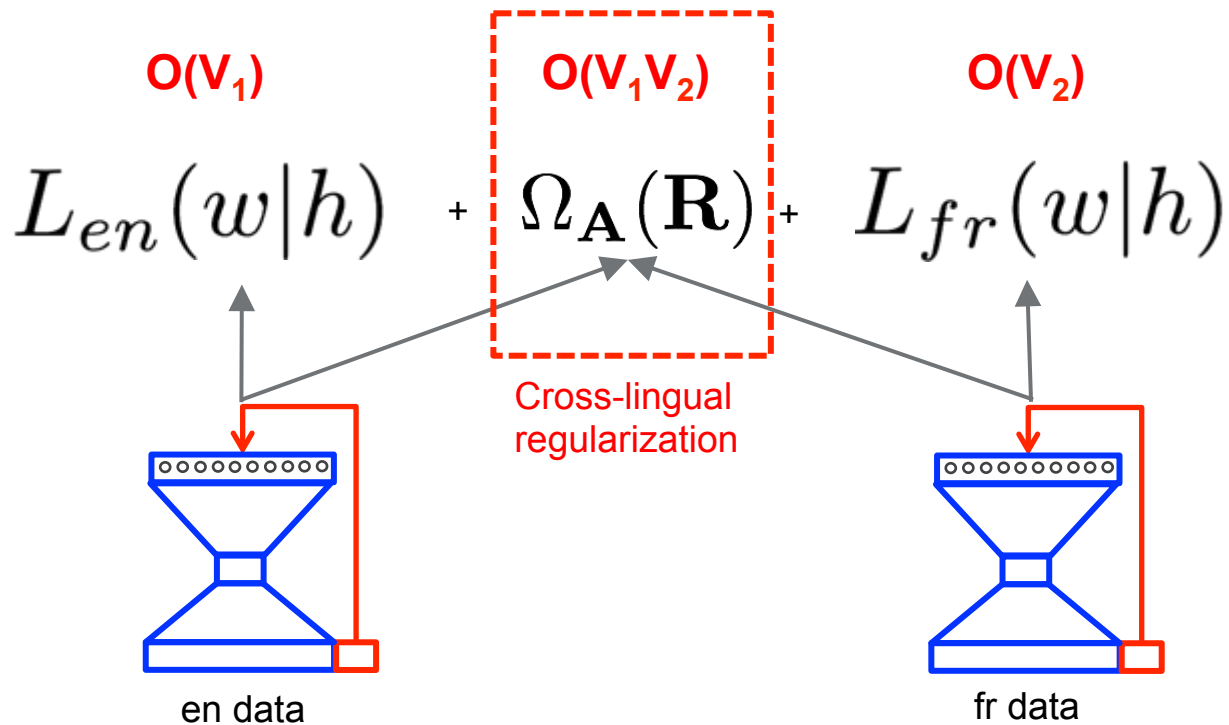
(Chandar *et al.*, 2013)

**BiCVM**

(Hermann *et al.*, 2014)

# Learning **Cross-lingual** Word Embeddings III

Online methods

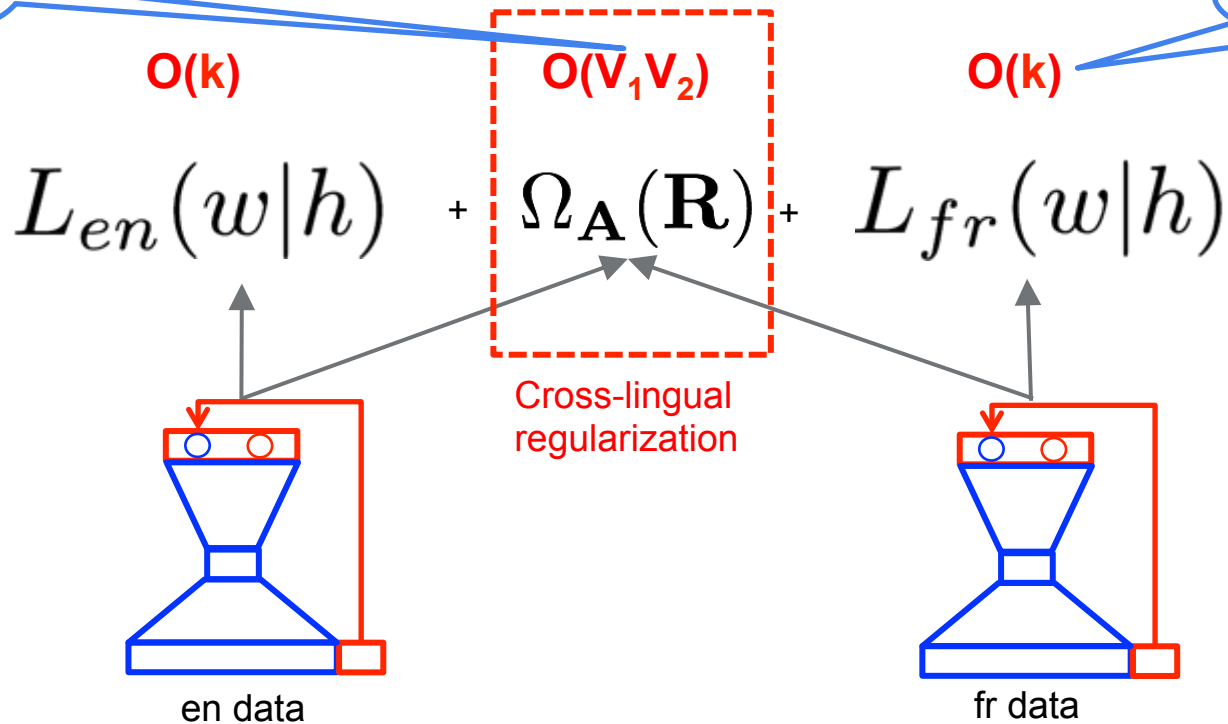


# Learning Cross-lingual Word Embeddings III

Online methods

Slow!

Fast sampled rank loss



# Multilingual distributed feature learning: Trade-offs

OFFLINE

METHOD	PROS	CONS
Translation Matrix ( <i>Mikolov et al. 2013</i> )	<ul style="list-style-type: none"><li>• FAST</li><li>• Simple to implement</li></ul>	<ul style="list-style-type: none"><li>• Assumes a global, linear, <i>one-to-one mapping</i> exists between words in 2 languages.</li><li>• Requires accurate dictionaries</li></ul>
Multilingual CCA ( <i>Faruqui et al. 2014</i> )		

# Multilingual distributed feature learning: Trade-offs

OFFLINE

PARALLEL

METHOD	PROS	CONS
Translation Matrix ( <i>Mikolov et al. 2013</i> )	<ul style="list-style-type: none"><li>• FAST</li><li>• Simple to implement</li></ul>	<ul style="list-style-type: none"><li>• Assumes a global, linear, <i>one-to-one mapping</i> exists between words in 2 languages.</li><li>• Requires accurate dictionaries</li></ul>
Multilingual CCA ( <i>Faruqui et al. 2014</i> )		
Bilingual Auto-encoders ( <i>Chandar et al. 2014</i> )	Simple to implement (?)	<ul style="list-style-type: none"><li>• Bag-of-words models</li><li>• Learns more semantic than syntactic features</li><li>• Reduced training data</li><li>• Big <i>domain bias</i></li></ul>
BiCVM ( <i>Hermann et al., 2014</i> )	Allows arbitrary differentiable sentence composition function	

# Multilingual distributed feature learning: Trade-offs

OFFLINE

PARALLEL

ONLINE

METHOD	PROS	CONS
Translation Matrix ( <i>Mikolov et al. 2013</i> )	<ul style="list-style-type: none"> <li>FAST</li> <li>Simple to implement</li> </ul>	<ul style="list-style-type: none"> <li>Assumes a global, linear, <i>one-to-one mapping</i> exists between words in 2 languages.</li> <li>Requires accurate dictionaries</li> </ul>
Multilingual CCA ( <i>Faruqui et al. 2014</i> )		
Bilingual Auto-encoders ( <i>Chandar et al. 2014</i> )	Simple to implement (?)	<ul style="list-style-type: none"> <li>Bag-of-words models</li> <li>Learns more semantic than syntactic features</li> <li>Reduced training data</li> <li>Big <i>domain bias</i></li> </ul>
BiCVM ( <i>Hermann et al., 2014</i> )	Allows arbitrary differentiable sentence composition function	
<i>Klementiev et al., 2012</i>	Can learn fine-grained, cross-lingual syntactic/semantic features (depends on window-length)	<ul style="list-style-type: none"> <li>SLOW</li> <li>Requires word-alignments (GIZA++/Fastalign)</li> </ul>
<i>Zhu et al., 2013</i>		

# Multilingual distributed feature learning: Trade-offs

OFFLINE

PARALLEL

ONLINE

METHOD	PROS	CONS
Translation Matrix ( <i>Mikolov et al. 2013</i> )	<ul style="list-style-type: none"> <li>FAST</li> <li>Simple to implement</li> </ul>	<ul style="list-style-type: none"> <li>Assumes a global, linear, <i>one-to-one mapping</i> exists between words in 2 languages.</li> <li>Requires accurate dictionaries</li> </ul>
Multilingual CCA ( <i>Faruqui et al. 2014</i> )		
Bilingual Auto-encoders ( <i>Chandar et al. 2014</i> )	Simple to implement (?)	<ul style="list-style-type: none"> <li>Bag-of-words models</li> <li>Learns more semantic than syntactic features</li> <li>Reduced training data</li> <li>Big <i>domain bias</i></li> </ul>
BiCVM ( <i>Hermann et al., 2014</i> )	Allows arbitrary differentiable sentence composition function	
<i>Klementiev et al., 2012</i>	Can learn fine-grained, cross-lingual syntactic/semantic features (depends on window-length)	<ul style="list-style-type: none"> <li>SLOW</li> <li>Requires word-alignments (GIZA++/FastAlign)</li> </ul>
<i>Zhu et al., 2013</i>		

*This work makes multilingual distributed feature learning more efficient for transfer learning and translation.*

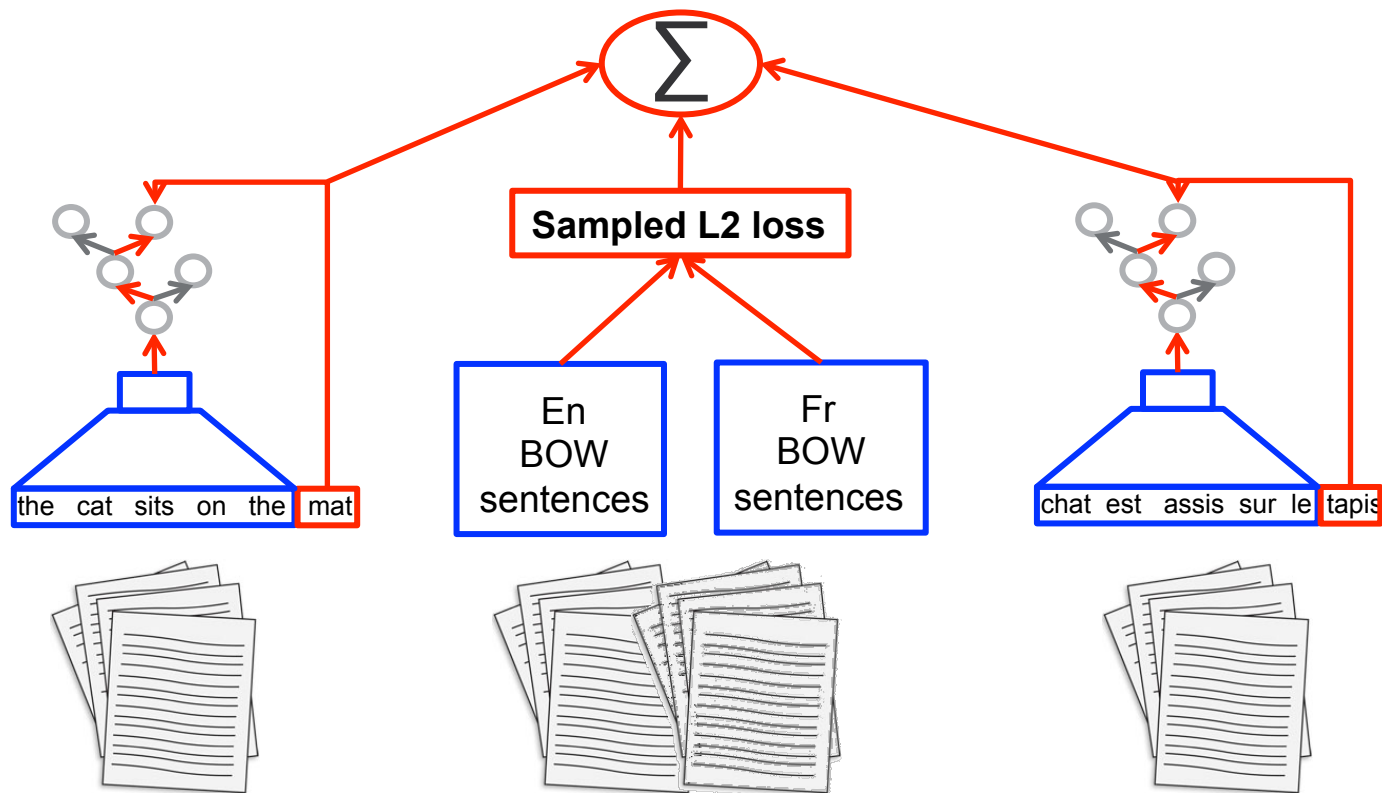


# BILBOWA:

Fast **B**ilingual **B**ag-Of-**W**ords  
Embeddings without **A**lignments



# BiBOWA Architecture

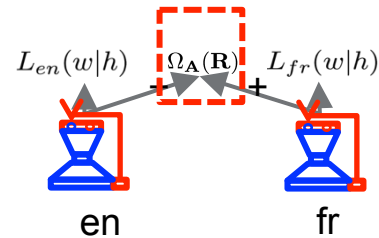


En monolingual

En-Fr parallel

Fr monolingual

# The BiBOWA Cross-lingual Objective I



We want to learn similar embeddings for translation pairs. The **exact** cross-lingual objective to minimize is the weighted sum over all distances of word-pairs:

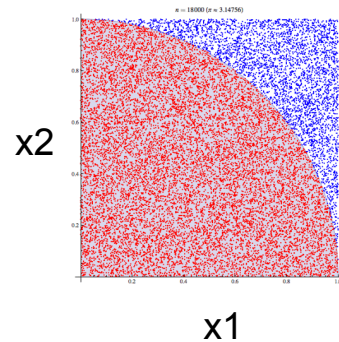
$$\Omega_{\mathbf{A}}(\mathbf{R}) = \sum_{w_i \in V_{\text{en}}} \sum_{w_j \in V_{\text{fr}}} a_{ij} ||\mathbf{R}_{[i,:]}^e - \mathbf{R}_{[j,:]}^f||^2$$

*Alignment  
"score"*

**Main contribution:** We **approximate** this by sampling parallel sentences.

## Quick Aside: Monte Carlo integration

$$\begin{aligned} & \mathbb{E}_{p(x)} f(x) \\ &= \int_x p(x) f(x) \text{ (Continuous)} \\ &= \sum_{x_i} Pr(x_i) f(x_i) \text{ (Discrete)} \\ &\approx 1/N \sum_{x \sim p(x)} f(x) \end{aligned}$$



## The BiBOWA Cross-lingual Objective II

$$\begin{aligned}
 \Omega_{\mathbf{A}}(\mathbf{R}) &= \sum_{i,j} a_{ij} \|\mathbf{R}_{[i,:]}^e - \mathbf{R}_{[j,:]}^f\|^2 \\
 &= \mathbb{E}_{(i,j) \sim P(w^e, w^f)} \left[ \|\mathbf{R}_{[i,:]}^e - \mathbf{R}_{[j,:]}^f\|^2 \right] \\
 &\approx \frac{1}{S} \sum_{s \in S} \frac{1}{mn} \sum_{(i,j) \in s} \|\mathbf{R}_{[i,:]}^e - \mathbf{R}_{[j,:]}^f\|^2
 \end{aligned}$$

$P(w^e, w^f)$  is the distribution of en-fr **word alignments**

Assume  $P(w^e, w^f)$  is uniform.  $m/n$  length of en/fr sentence.

Now we set **S = 1** at each time step  $t$ :

$$\Omega_{\mathbf{A}}^{(t)}(\mathbf{R}) = \left\| \underbrace{\frac{1}{m} \sum_{i \in s_e} \mathbf{R}_{[i,:]}^e}_{\text{mean en sentence-vector}} - \underbrace{\frac{1}{n} \sum_{j \in s_f} \mathbf{R}_{[j,:]}^f}_{\text{mean fr sentence-vector}} \right\|^2$$

## Implementation Details: Open-source

Implemented in C (part of word2vec, soon). Multi-threaded: One thread per language (monolingual), and one additional thread per language-pair (cross-lingual) (i.e. **asynchronous SGD**)

Runs at ~10-50K words per second on MBP.

Can process 500M words (monolingual) and 45M words (parallel, recycles) in about 2.5h on my MBP.

## Cross-lingual subsampling for better results

At training step  $t$ , draw a random number  $u \sim U[0,1]$ . Then:

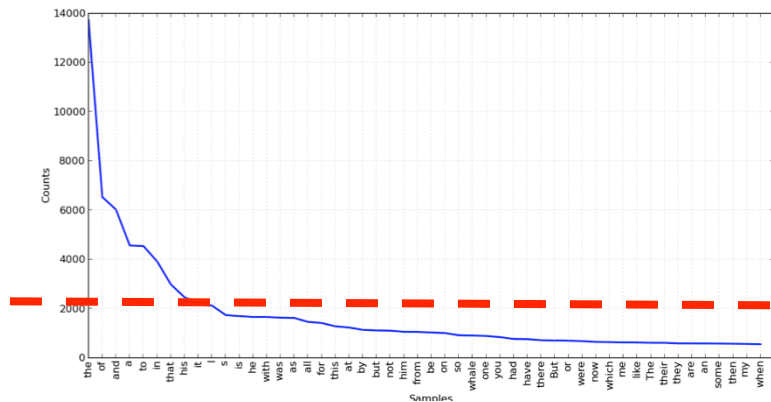
$$\Omega_{\mathbf{A}}^{(t)}(\mathbf{R}) = \left\| \frac{1}{m} \sum_{i \in s_e} \mathbb{1}_{u < f(w_i)} \mathbf{R}_{[i,:]}^e - \frac{1}{n} \sum_{j \in s_f} \mathbb{1}_{u < f(w_j)} \mathbf{R}_{[j,:]}^f \right\|^2$$

Want to estimate **alignment** statistics  $P(e,f)$ .

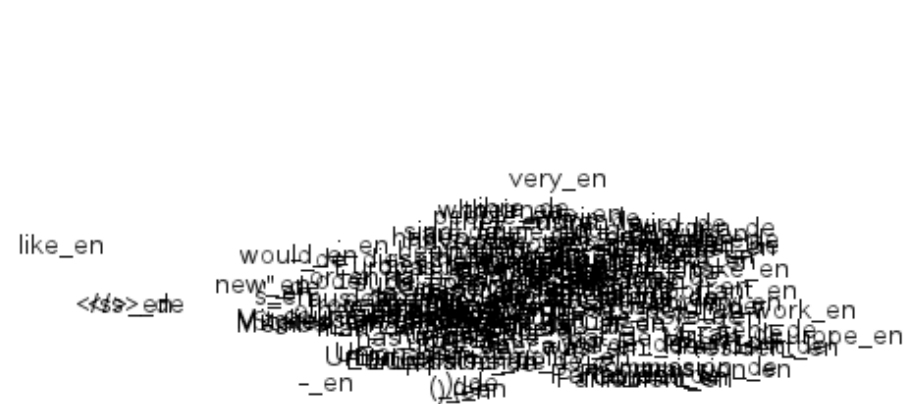
Skewed at the sentence-level by  
(unconditional) unigram word frequencies.

**Simple solution:**

*Subsample* frequent words to  
flatten the distribution!



## Cross-lingual subsampling for better results

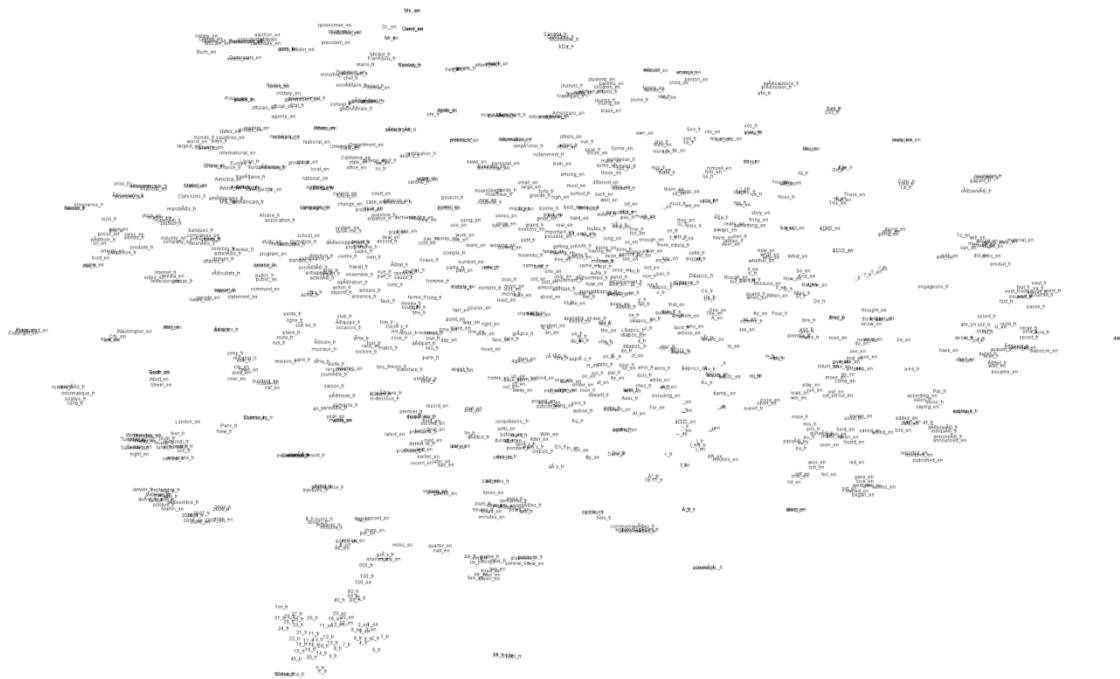


## Without SS



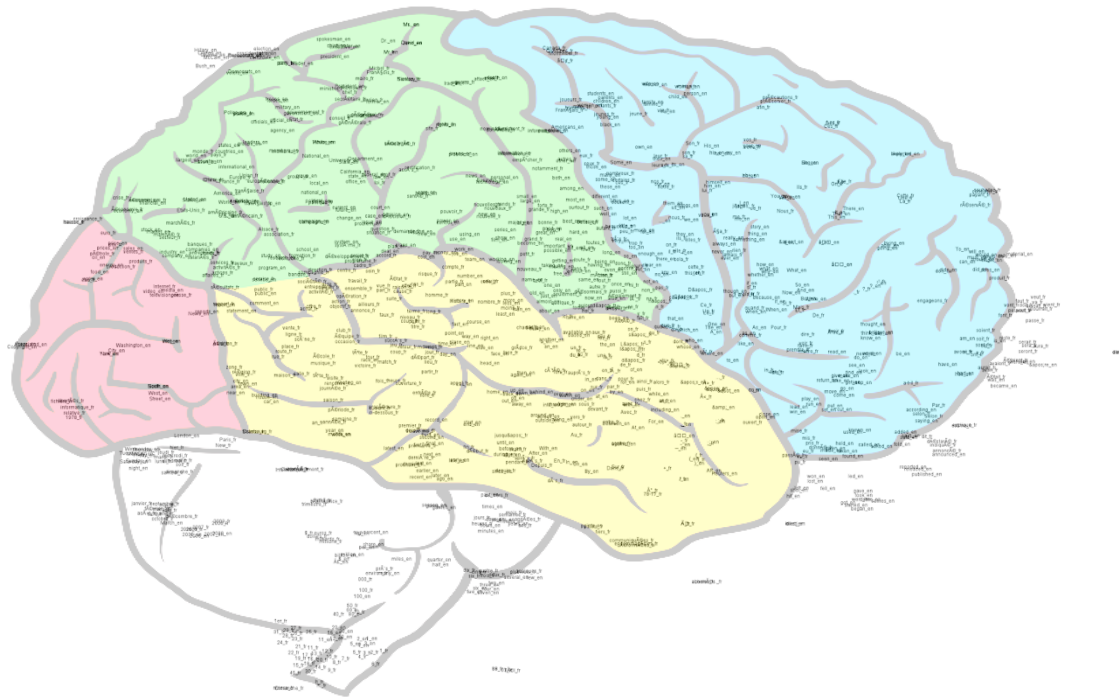
## With SS

# Qualitative Analysis: *en-fr* t-SNEs I

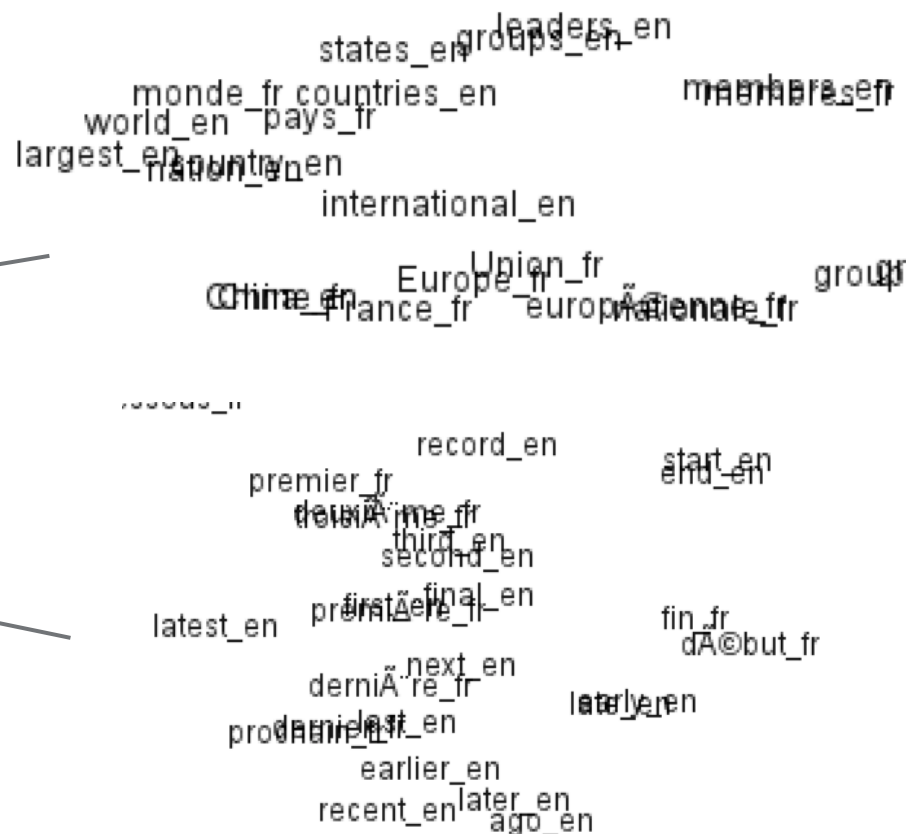
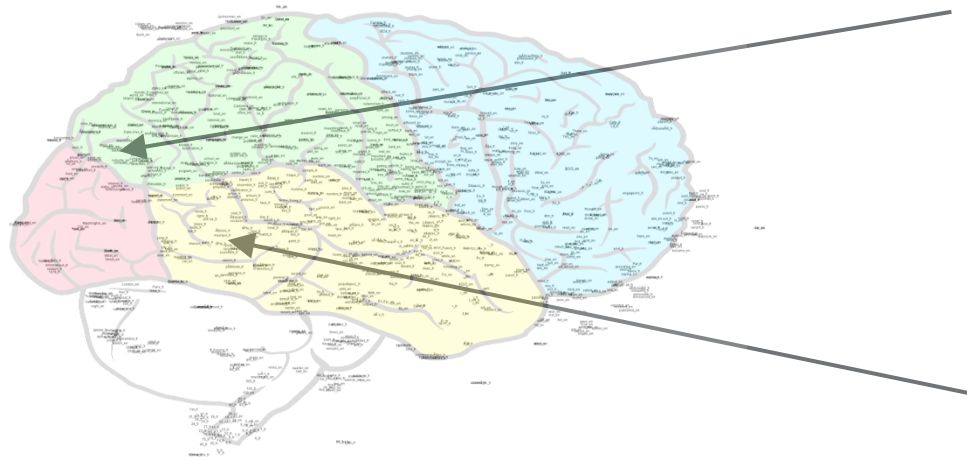




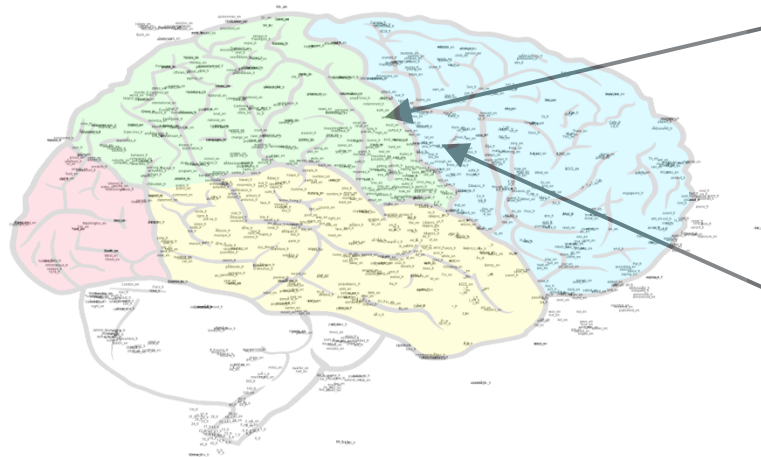
## Qualitative Analysis: *en-fr* t-SNEs I



# Qualitative Analysis: *en-fr*

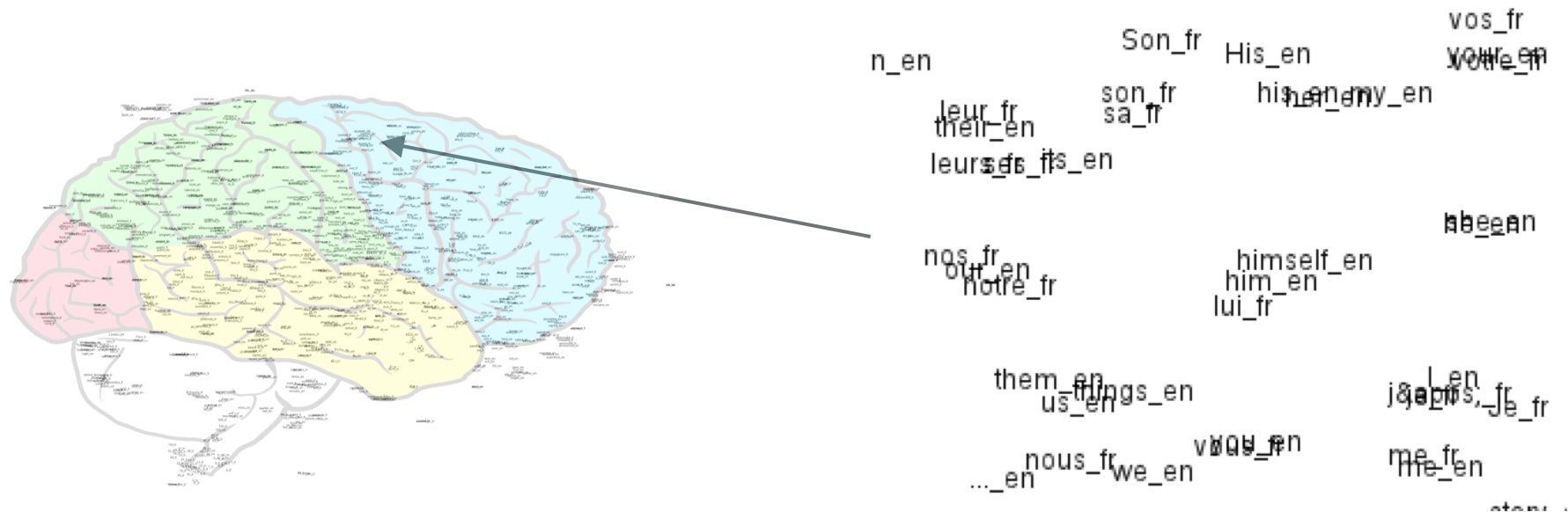


# Qualitative Analysis: *en-fr*

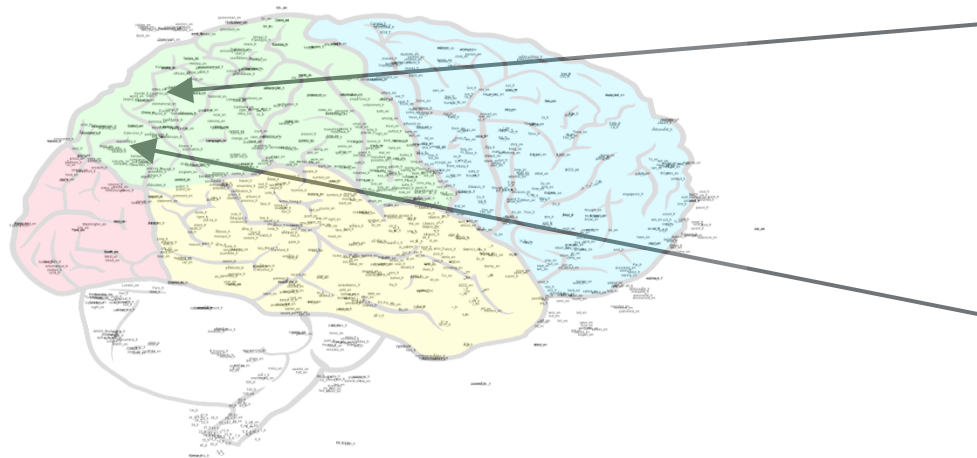


jusqu'&\_fr Au\_fr  
until\_en  
before\_en\_1\_fr With\_en  
durant\_en\_1\_fr After\_en  
après\_fr  
pendant\_fr s\_fr En\_fr In\_en Oh\_en  
depuis\_fr  
side\_en\_1\_fr By\_  
d'&\_fr  
parteniers\_fr  
times\_en  
mois\_fr  
jours\_fr semaines\_fr  
heures\_fr jours\_en années\_fr  
minutes\_en

# Qualitative Analysis: *en-fr*



# Qualitative Analysis: *en-fr*



troops\_en  
military\_en  
gouvernement\_fr  
Police\_en  
police\_en  
officials\_en  
official\_en  
Etat\_fr  
agency\_en  
America\_en  
franÃ§aise\_fr  
United\_en  
States\_en  
World\_en  
American\_en  
Etats-Unis\_fr  
amÃ©ricaine\_fr  
U.S.\_en  
amÃ©ricain\_fr

# Experiments: *En-De* Cross-lingual Document Classification

Exact replication (obtained from the authors) of *Klementiev et al.*'s cross-language document classification setup:

**Goal:** Classify documents in target language using only labelled documents in source language.

**Data:** English-German RCV1 data (5K test, 100 - 10K training, 1K validation)

## 4 Labels:

- CCAT (Corporate/Industrial),
- ECAT (Economics),
- GCAT (Government/Social), and
- MCAT (Markets)

## Results: English-German Document Classification

	<i>en2de</i>	<i>de2en</i>	<i>Training Size</i>	<i>Training Time (min)</i>
<i>Majority class</i>	46.8	46.8	-	-
<i>Glossed</i>	65.1	68.6	-	-
<i>MT</i>	68.1	67.4	-	-
Klementiev et al.	77.6	71.1	50M words	14,400 (10 days)

## Results: English-German Document Classification

	<i>en2de</i>	<i>de2en</i>	<i>Training Size</i>	<i>Training Time (min)</i>
<i>Majority class</i>	46.8	46.8	-	-
<i>Glossed</i>	65.1	68.6	-	-
<i>MT</i>	68.1	67.4	-	-
Klementiev et al.	77.6	71.1	50M words	14,400 (10 days)
Bilingual Autoencoders	<b>91.8</b>	72.8	50M words	4,800 (3.5 days)
BiCVM	83.7	71.4	50M words	15
<b>BiBOWA</b> (this work)	86.5	<b>75</b>	50M words	<b>6</b>



## Experiments: WMT11 *English-Spanish* Translation

- Trained BilBOWA model on En-Es Wikipedia/Europarl data.
  - Vocabulary = 200K
  - Embedding dimension = 40,
  - Crosslingual  $\lambda$ -weight in {0.1, 1.0, 10.0, 100.0}
- Exact replica of (*Mikolov, Le, Sutskever, 2013*):
  - Evaluated on WMT11 lexicon, translated using GTranslate
  - Top 5K-6K words as test set

## Experiments: WMT11 *English*→*Spanish* Translation

	<i>Dimension</i>	<i>Prec@1</i>	<i>Prec@5</i>	<i>% Coverage</i>
<i>Edit distance</i>	-	13	24	92.9
<i>Word co-occurrence</i>	-	19	30	92.9
Translation Matrix	<i>300-800</i>	33	51	92.9
<b>BilBOWA (This work)</b>	40	<b>39 (+6%)</b>	<b>51</b>	92.7

## Experiments: WMT11 *Spanish*→*English* Translation

	<i>Dimension</i>	<i>Prec@1</i>	<i>Prec@5</i>	<i>% Coverage</i>
<i>Edit distance</i>	-	18	27	92.9
<i>Word co-occurrence</i>	-	20	30	92.9
Translation Matrix	300-800	35	52	92.9
<b>BiBOWA (This work)</b>	40	<b>44 (+11%)</b>	<b>55 (+3%)</b>	92.7



# BARISTA:

Bilingual Adaptive Reshuffling  
*with*  
Individual Stochastic Alternatives

# Motivation

Word embedding models learn to predict targets from contexts by clustering similar words into soft (distributional) equivalence classes

For some tasks, we may easily obtain the desired equivalence classes:

- **POS:** Wiktionary
- **Super-sense** (SuS) tagging: WordNet
- **Translation:** Google Translate / dictionaries

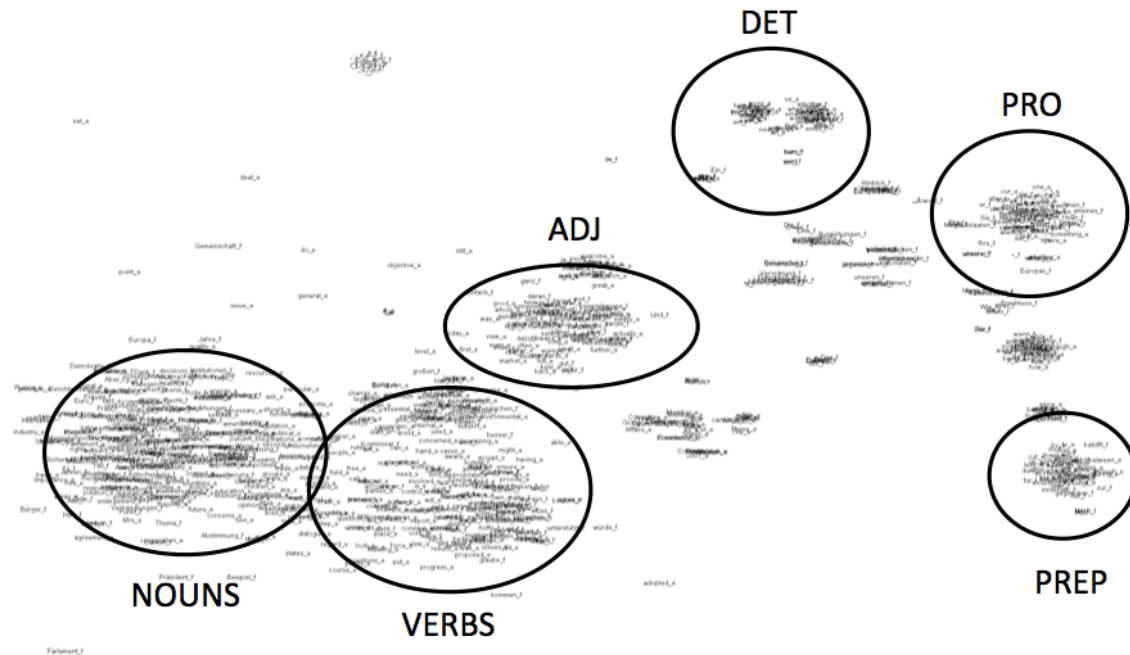
**BARISTA** embeds additional task-specific semantic information by corrupting the training data according to known equivalences  $C(w)$ .

# BARISTA Algorithm

1. Shuffle  $D_{\text{en}}$  &  $D_{\text{fr}}$   $\rightarrow D$
2. For  $w$  in  $D$ :
3.   If  $w$  in  $C$  then  $w' \sim C(w)$  else  $w' = w$
4.    $D' += w'$
4. Train off-the-shelf embedding model on  $D'$

For example: “*we build the house*”:

1. *we build la voiture / they run la house* (**POS**)
2. *we construire the maison / nous build la house* (**Translations**)



kind more e us\_e which has f  
 every thing us\_e what any e  
 the thing e einende many e  
 ein\_e e weany e e f  
 next there e that e e f  
 an\_e e in\_e f

अनंत

„sowohl\_“

vornehmlich\_f  
 informationalen\_f  
 par nationalen\_f  
 seren\_f  
 einigenden\_f

Hauptfunktion\_f  
Kommission\_f  
Was\_Was\_f  
Was\_Was\_f  
Was\_Was\_f  
Erachtens\_f

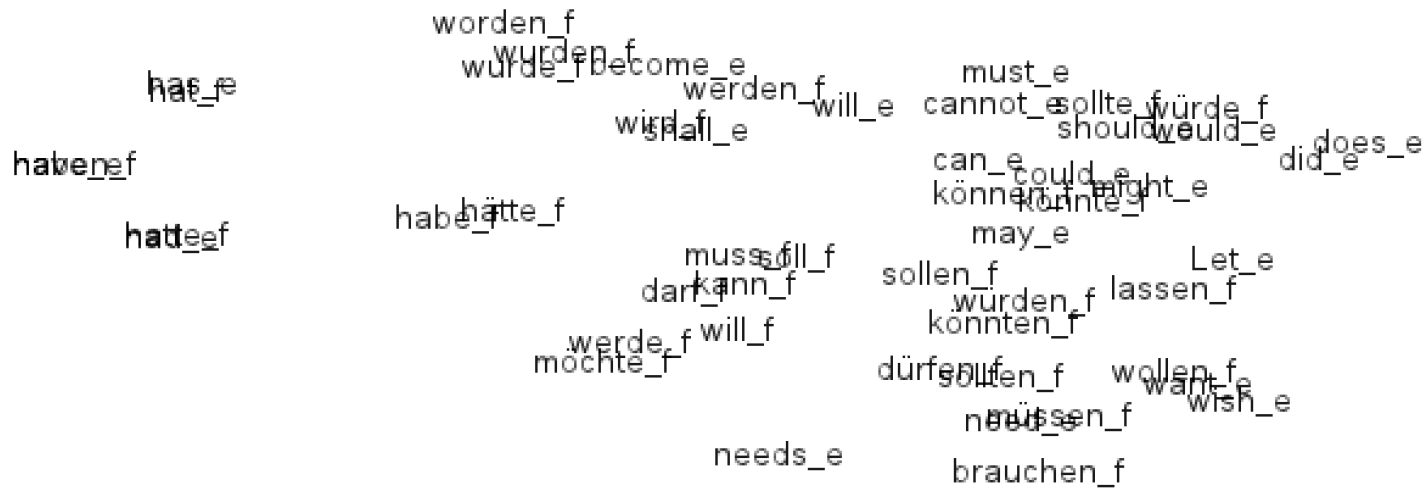
[illegible]



## Qualitative (Translations)

# Prepositions

## Qualitative (Translations)



# Modals

# Cross-lingual POS tagging

Language	Baseline	Random	Klmtv	POS-50	POS-300	Tr-50	Tr-300	DP	B-K
Spanish	80.6	81.8	79.8	82.4	81	81.6	82.6	<b>84.2</b>	80.2
German	80.4	82.7	82.8	81.8	84.1	82.6	<b>84.8</b>	82.8	81.3
Danish	63	68.9	-	68.9	72.4	71.8	78.4	<b>83.2</b>	69.1
Swedish	71.6	73.7	-	75	76	75.4	77.5	<b>80.5</b>	70.1
Italian	80.1	81.3	-	82.1	80.9	82.1	80.7	<b>86.8</b>	68.1
Dutch	74.5	77.2	-	78.3	77.4	78.7	<b>80.3</b>	79.5	65.1
Portuguese	76.9	78.1	-	77.3	76.1	80.6	80.5	<b>87.9</b>	78.4
Avg	75.3	77.7	-	78	78.3	79	80.7	<b>83.6</b>	73.2

## Cross-lingual SuperSense (SuS) tagging

	Baselines		BARISTA	
	MFS	bl	Tr-300	WN-300
blogs	49.1	46.7	50.5	<b>61.9</b>
forum	44.5	41.2	45.0	<b>53.9</b>
magazine	46.5	45.2	50.4	<b>51.5</b>
newswire	48.4	45.4	52.7	<b>60.9</b>
reviews	48.4	44.9	50.4	<b>55.8</b>
speech	51.1	48.4	51.5	<b>58.4</b>
Avg	48.1	45.4	50.5	<b>58.3</b>

## Conclusion

I presented **B<sub>IL</sub>BOWA**, an efficient, bilingual word embedding model with an open-source C implementation (part of word2vec, soon) and **B<sub>AR</sub>ISTA**, a simple technique for embedding additional task-specific cross-lingual information.

*Qualitative* experiments on En-Fr & En-De show that the learned embeddings capture fine-grained cross-lingual linguistic regularities.

*Quantitative* results on Es, De, Da, Sv, It, Nl, Pt for:

- semantic transfer (document classification, xling SuS-tagging)
- lexical transfer (word-level translation, xling POS-tagging)

Thanks!

Questions / Comments?