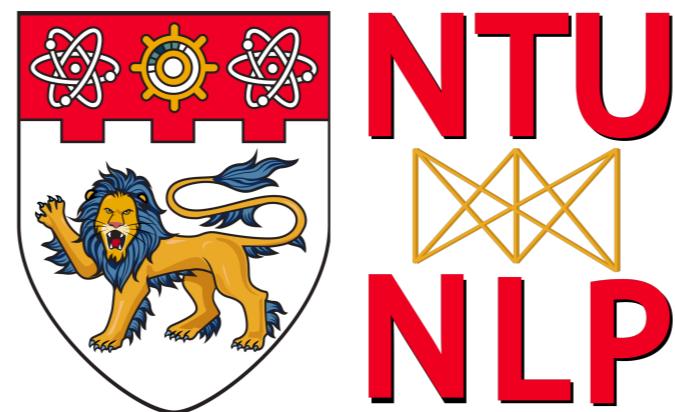


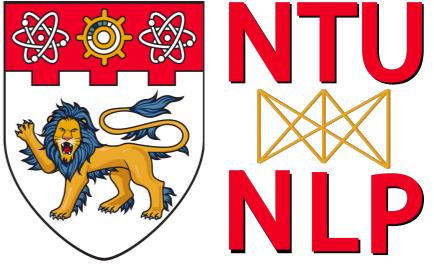
# Deep Learning for Natural Language Processing

Shafiq Joty



Lecture 13: Multilingual NLP & Recap

# Overview



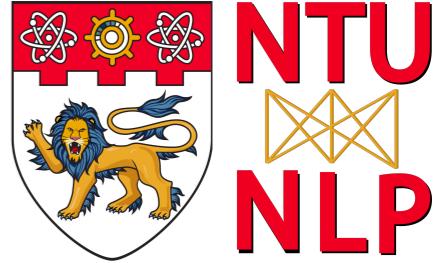
## Models/Algorithms

- Linear models
- Feed-forward Neural Nets (FNN)
- Window-based methods
- Convolutional Nets
- Recurrent Neural Nets
- Recursive Neural Nets

## NLP tasks/applications

- Word meaning
- Language modelling
- Sequence tagging
- Sequence encoding
- Parsing
- Hierarchical encoding

# Overview



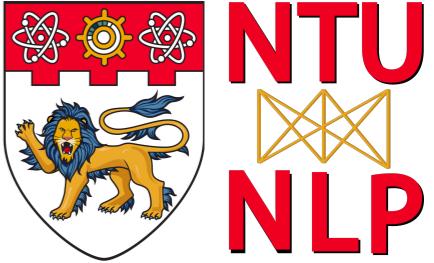
## Models/Algorithms

- Seq2Seq
  - + Attention
  - + Subword
- Seq2Seq Variants
- Transformer Seq2Seq

## NLP tasks/applications

- Machine Translation
- Summarization
- Parsing
- Dialogue generation

# Overview

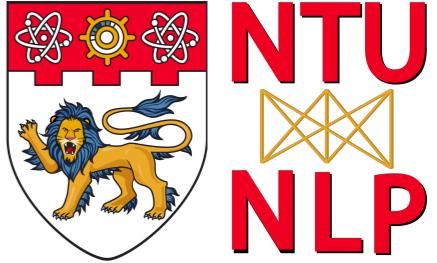


Models/Algorithms

NLP tasks/applications

- Multilingual NLP
- Self-supervised pretraining
- Robust & adversarial NLP

# Today: Multilingual NLP



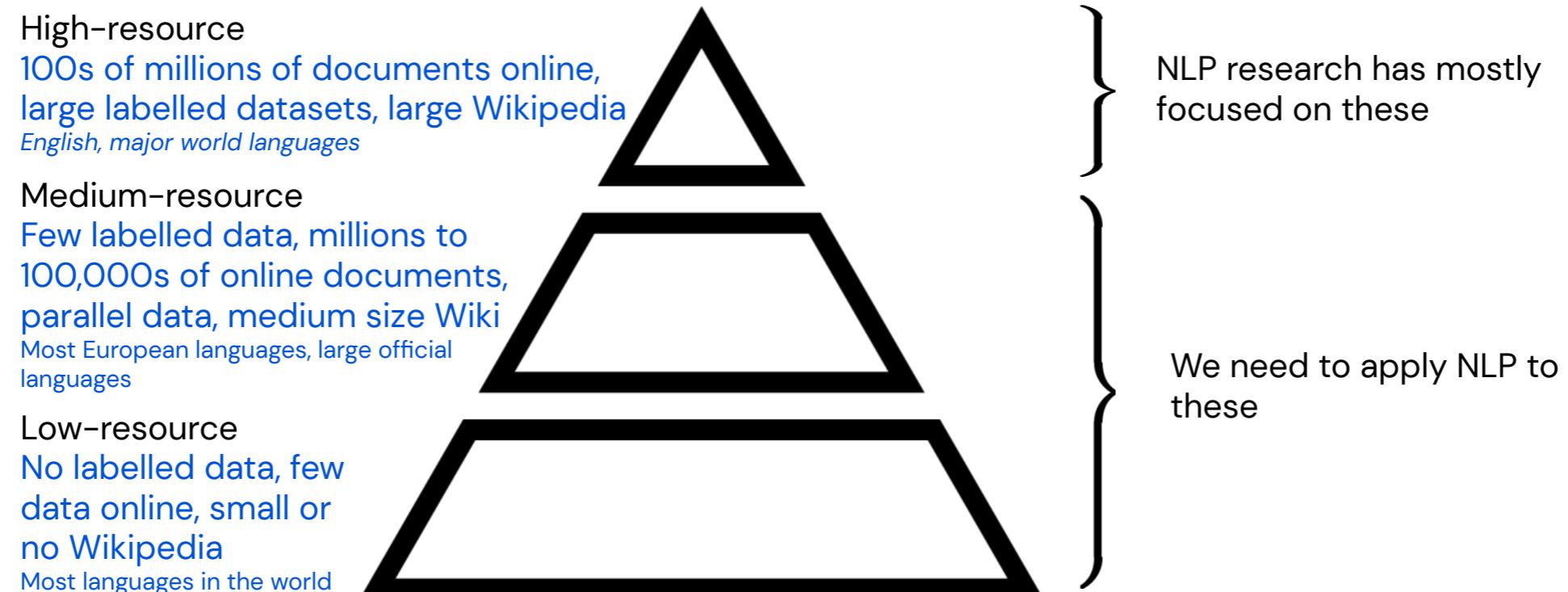
- Why Multilingual NLP?
- Cross-lingual Models

# Why Multilingual NLP?

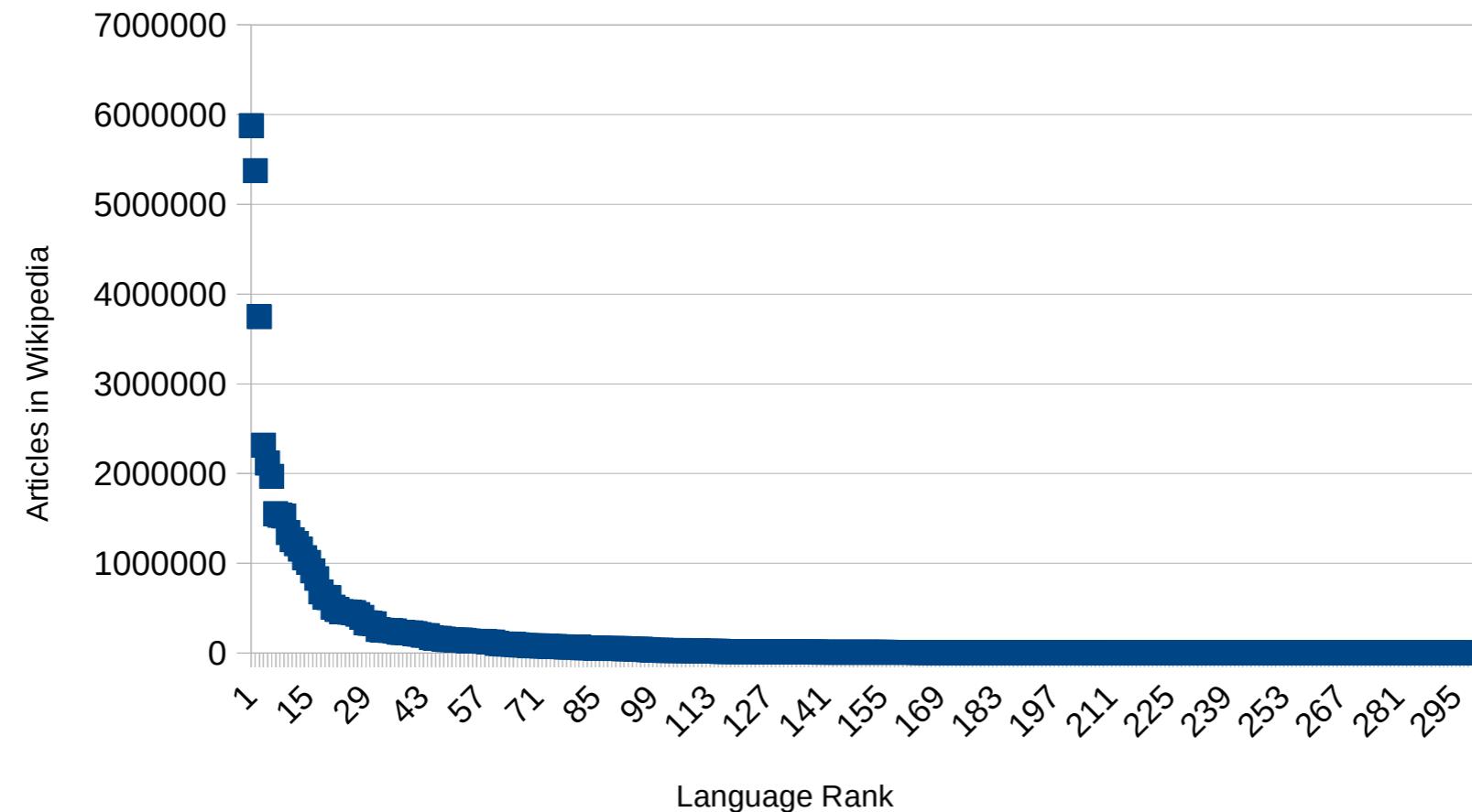
NLP not just for English but for the rest of the world's 7000 languages



# Why Multilingual NLP?



# Why Multilingual NLP?



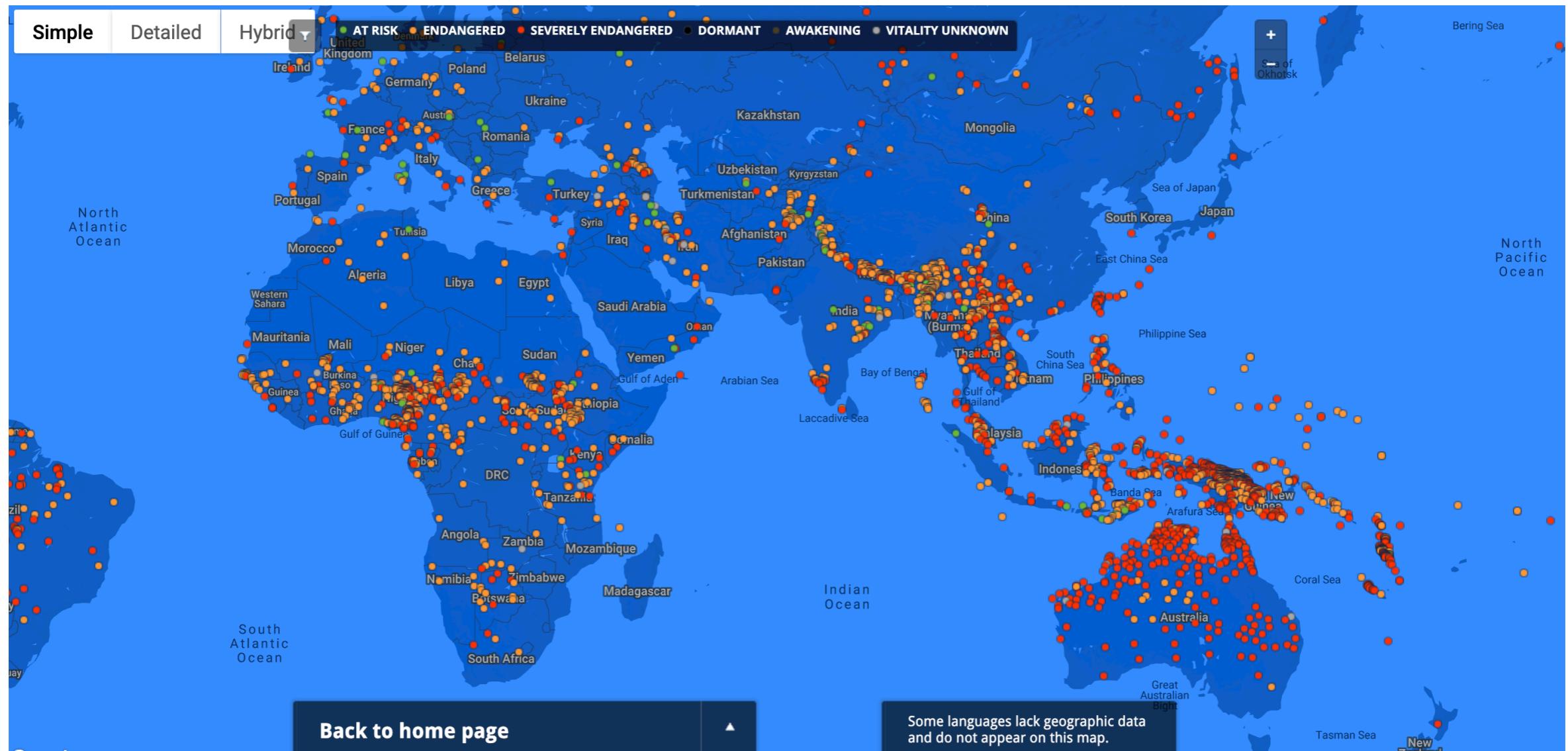
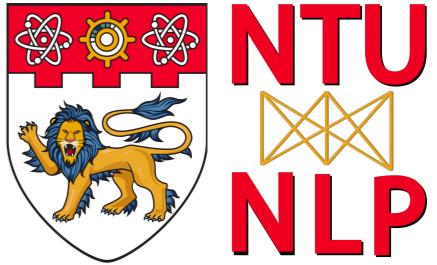
- There is not enough monolingual data for many languages
- Even less annotated data for tasks (MT, sequence labels)

# Old Paradigm

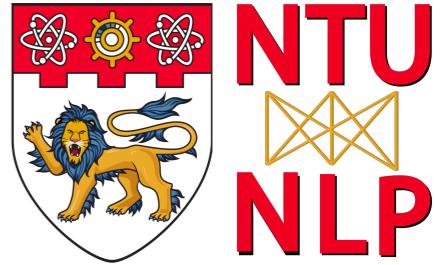
Old paradigm:

- Language-specific NLP
  - Language-specific feature computation and preprocessing
  - Manual curation and annotation of large-scale resources for thousands of languages is infeasible or prohibitively expensive
- 
- |   |  |
|---|--|
| - | <i>p.currentSense + p.lemma</i>        |
| - | <i>p.currentSense + p.pos</i>          |
| - | <i>p.currentSense + a.pos</i>          |
| - | <i>p<sub>-1</sub>.FEAT1</i>            |
| - | <i>p.FEAT2</i>                         |
| - | <i>p<sub>1</sub>.FEAT3</i>             |
| - | <i>p.semrn.semdpred</i>                |
| - | <i>p.lm.dpred</i>                      |
| - | <i>p.form + p.children.dpred.bag</i>   |
| - | <i>p.lemma<sub>n</sub> (n = -1, 0)</i> |
| - | <i>p.lemma + p.lemma<sub>1</sub></i>   |
| - | <i>p.pos<sub>-1</sub> + p.pos</i>      |
| - | <i>p.pos<sub>1</sub></i>               |
| - | <i>p.pos + p.children.dpred.bag</i>    |

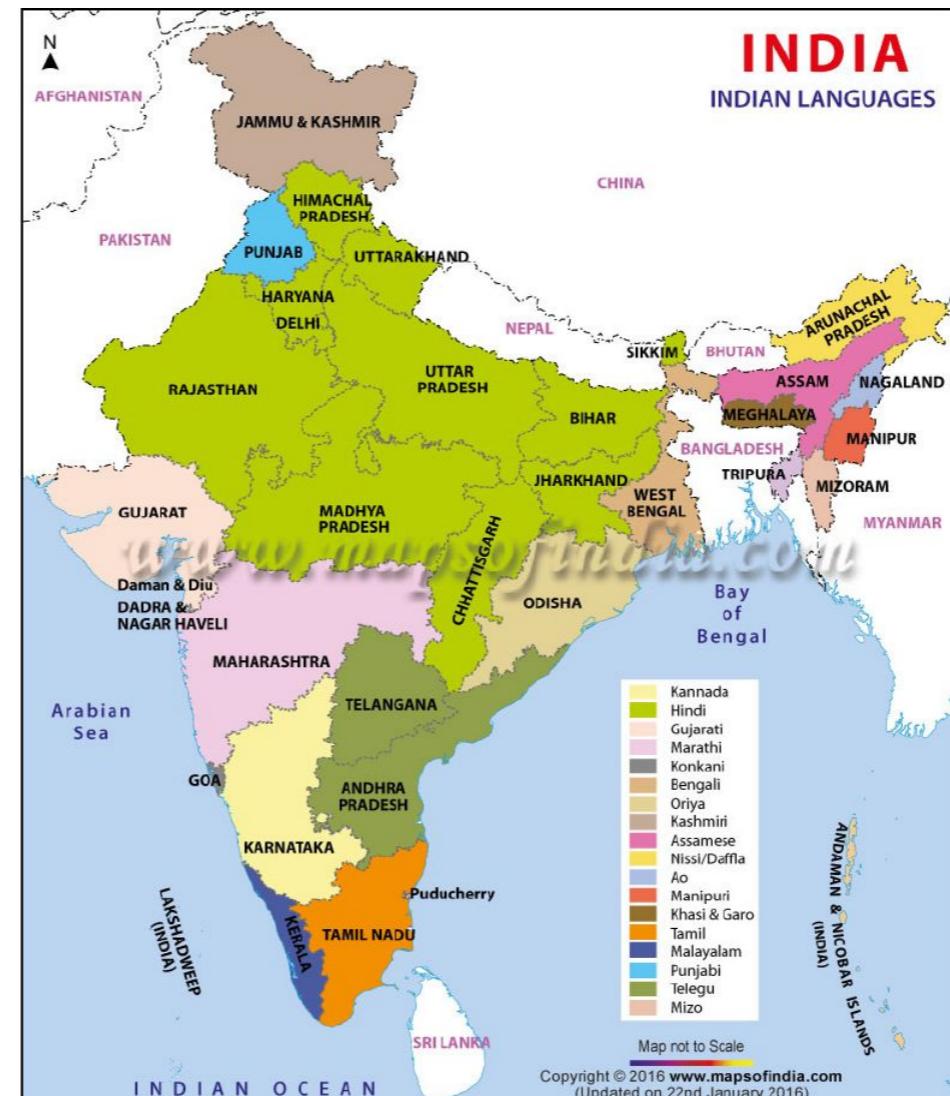
# Why Multilingual NLP?



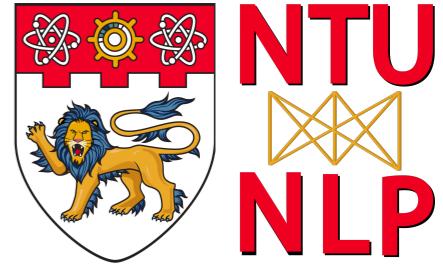
# Low Resource Languages



There are about 460 languages in India.  
1.38 billion people



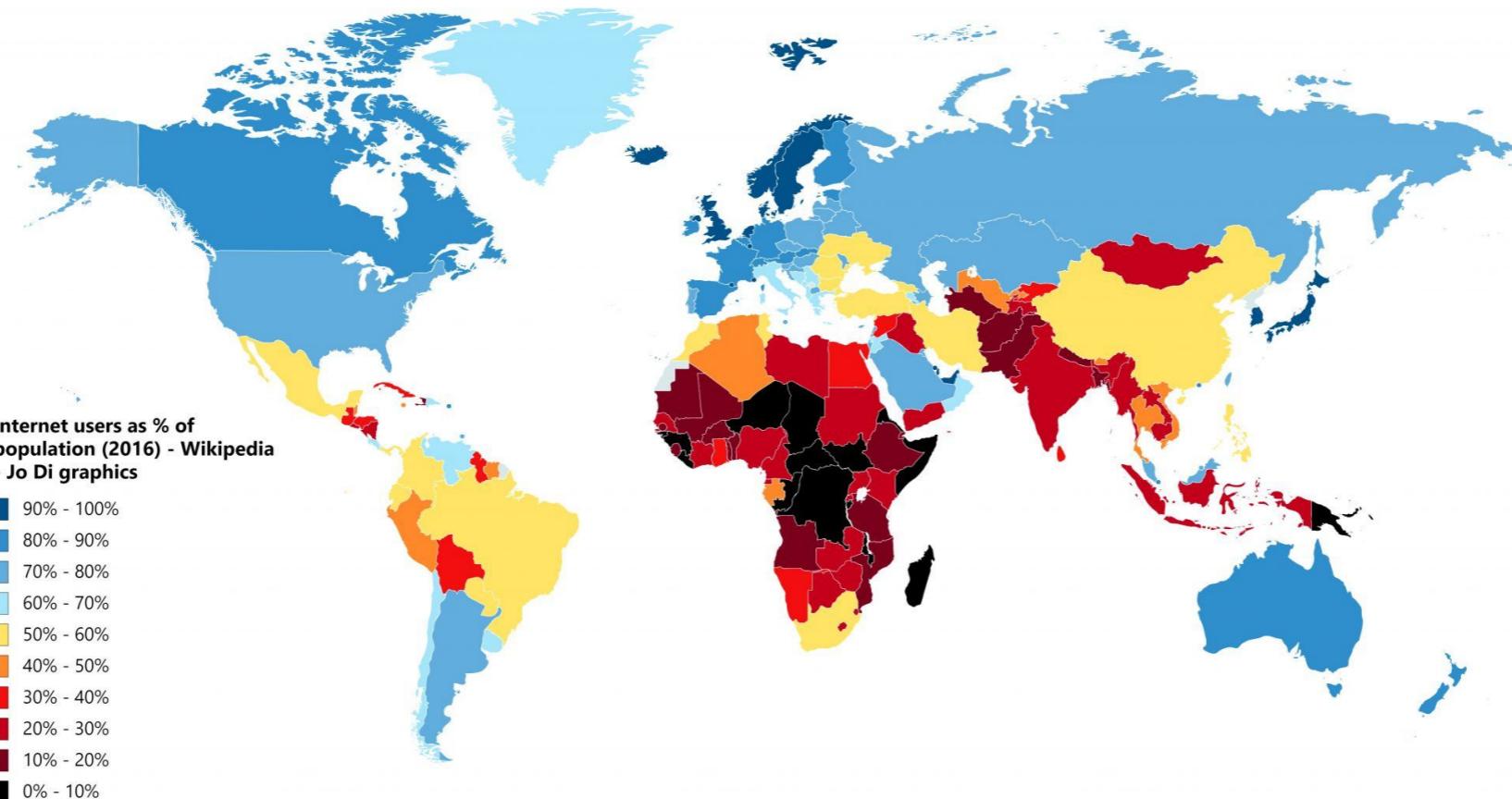
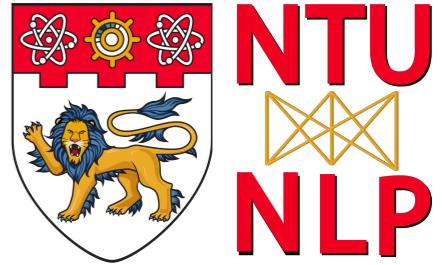
# Low Resource Languages



Africa is a continent with a very high linguistic diversity:  
there are an estimated 1.5-2K African languages from 6 language families.  
**1.33 billion people**

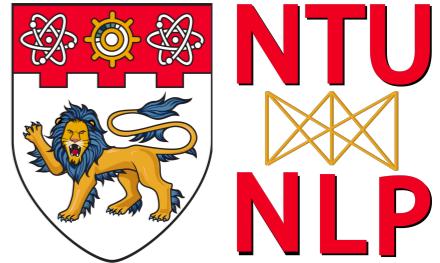


# Low Resource Languages



40% of world's population: South Asia - 1.75 billion, Africa - 1.3 billion, etc.

# How to define similarity across Languages



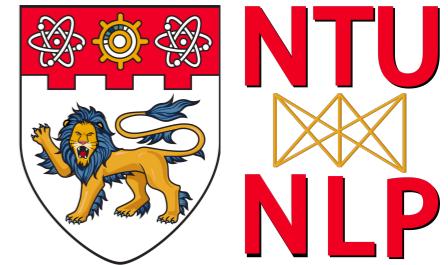
## Word overlap and sub-word overlap

- |              |              |            |            |
|--------------|--------------|------------|------------|
| ○ Russian    | – Русский    | ○ Japanese | – 日本人      |
| ○ Ukrainian  | – Українська | ○ Turkish  | – Türk     |
| ○ Chinese    | – 中文         | ○ Hebrew   | – עברית    |
| ○ Korean     | – 한국어        | ○ Arabic   | – عربی     |
| ○ Vietnamese | – Tiếng Việt | ○ Hindi    | – हिन्दी   |
| ○ Georgian   | – ქართული    | ○ Xhosa    | – isiXhosa |

## Areal similarity

## Demographic similarity

# How to define similarity across Languages

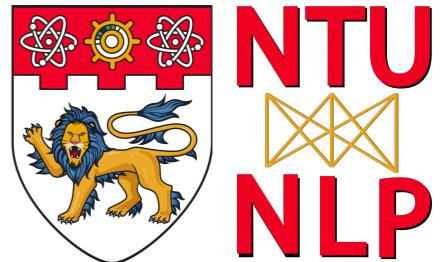


## Genealogical similarity

1. Niger–Congo (1,542 languages) (21.7%)
2. Austronesian (1,257 languages) (17.7%)
3. Trans–New Guinea (482 languages) (6.8%)
4. Sino-Tibetan (455 languages) (6.4%)
5. Indo-European (448 languages) (6.3%)
6. Australian [dubious] (381 languages) (5.4%)
7. Afro-Asiatic (377 languages) (5.3%)
8. Nilo-Saharan [dubious] (206 languages) (2.9%)
9. Oto-Manguean (178 languages) (2.5%)
10. Austroasiatic (167 languages) (2.3%)
11. Tai–Kadai (91 languages) (1.3%)
12. Dravidian (86 languages) (1.2%)
13. Tupian (76 languages) (1.1%)

[www.ethnologue.com](http://www.ethnologue.com)

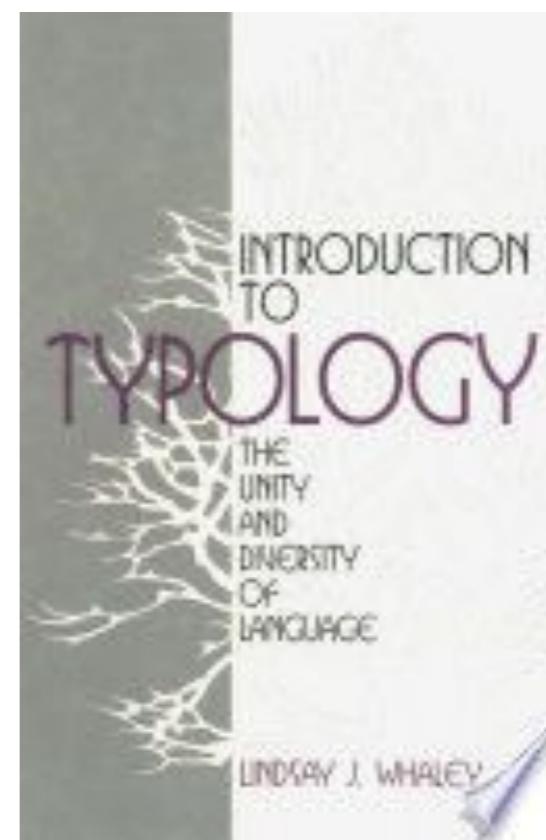
# How to define similarity across Languages



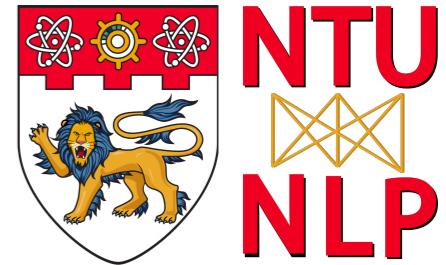
## Typological similarity

- Linguistic typology: classification of languages according to their functional and structural properties
  - explains common properties across languages
  - explains structural diversity across languages

“The classification of languages or components of languages based on shared formal characteristics.”



# How to define similarity across Languages



## Typological similarity

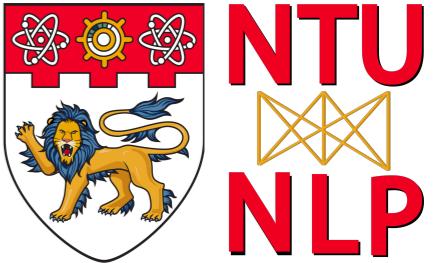
- 2,676 languages, 192 attributes

ID#	Feature Name	Category	Feature Values
1	Consonant Inventories	Phonology (19)	{1:Large, 2:Small, 3:Moderately Small, 4:Moderately Large, 5:Average}
23	Locus of Marking in the Clause	Morphology (10)	{1:Head, 2:None, 3:Dependent, 4:Double, 5:Other}
30	Number of Genders	Nominal Categories (28)	{1:Three, 2:None, 3:Two, 4:Four, 5:Five or More}
58	Obligatory Possessive Inflection	Nominal Syntax (7)	{1:Absent, 2:Exists}
66	The Perfect	Verbal Categories (16)	{1:None, 2:Other, 3:From 'finish' or 'already', 4:From Possessive}
81	Order of Subject, Object and Verb	Word Order (17)	{1:SVO, 2:SOV, 3:No Dominant Order, 4:VSO, 5:VOS, 6:OVS, 7:OSV}
121	Comparative Constructions	Simple Clauses (24)	{1:Conjoined, 2:Locational, 3:Particle, 4:Exceed}
125	Purpose Clauses	Complex Sentences (7)	{1:Balanced/deranked, 2:Deranked, 3:Balanced}
138	Tea	Lexicon (10)	{1:Other, 2:Derived from Sinitic 'cha', 3:Derived from Chinese 'te'}
140	Question Particles in Sign Languages	Sign Languages (2)	{1:None, 2:One, 3:More than one}
142	Para-Linguistic Usages of Clicks	Other (2)	{1:Logical meanings, 2:Affective meanings, 3:Other or none}

Example from Georgi, Xia and Lewis (2010)

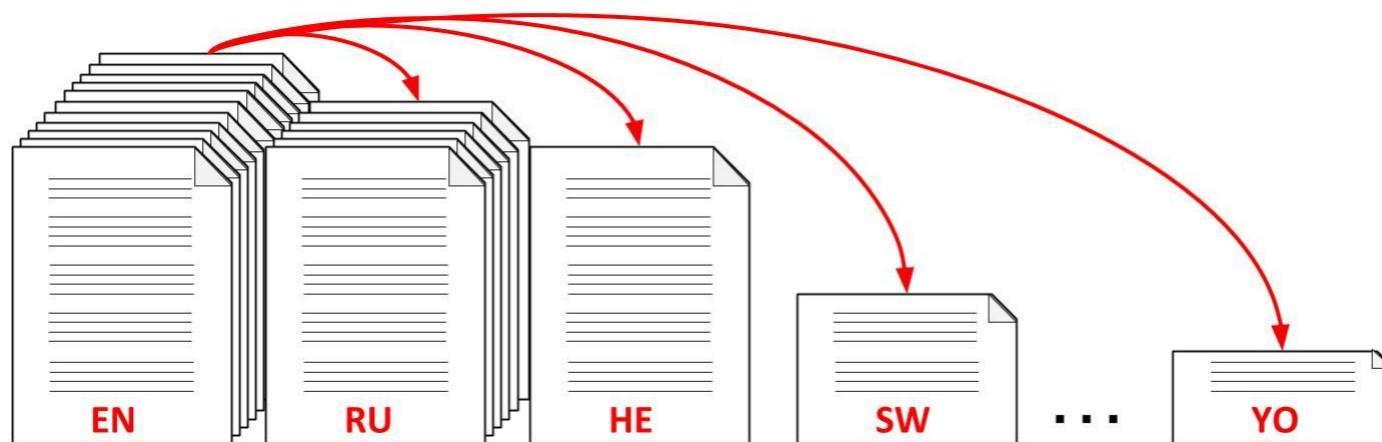
Dryer, Matthew S. & Haspelmath, Martin (eds.) 2013.  
[The World Atlas of Language Structures Online](#).  
Leipzig: Max Planck Institute for Evolutionary Anthropology.

# Approaches to Low-resource NLP

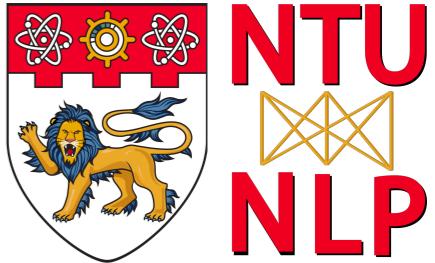


## 1. Cross-lingual transfer learning

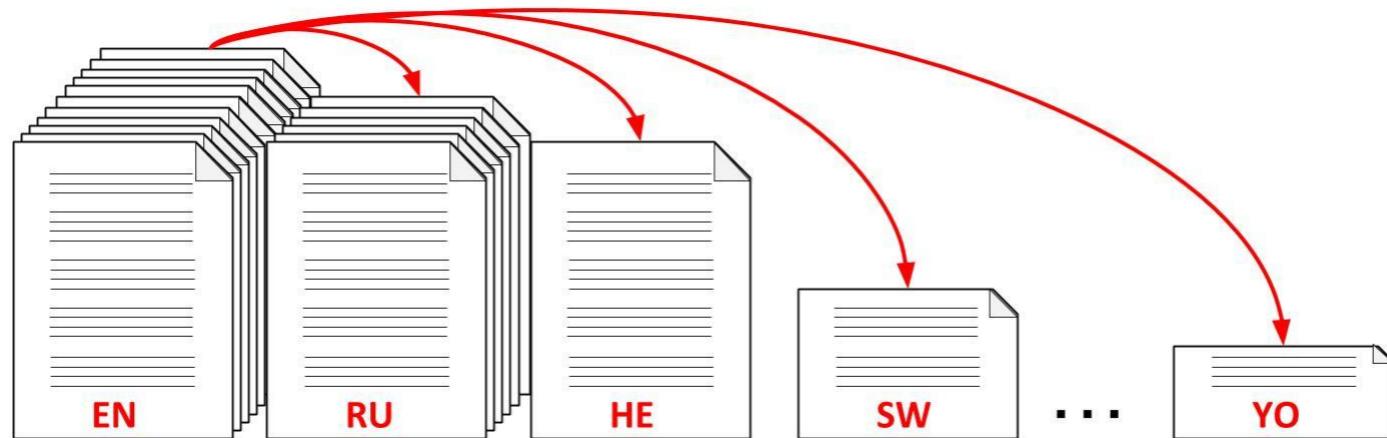
- Transfer of resources and models from resource-rich source to resource-poor target languages



# Approaches to Low-resource NLP

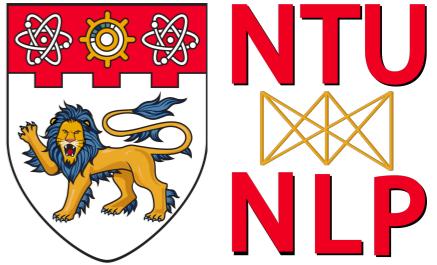


## 1. Cross-lingual transfer learning

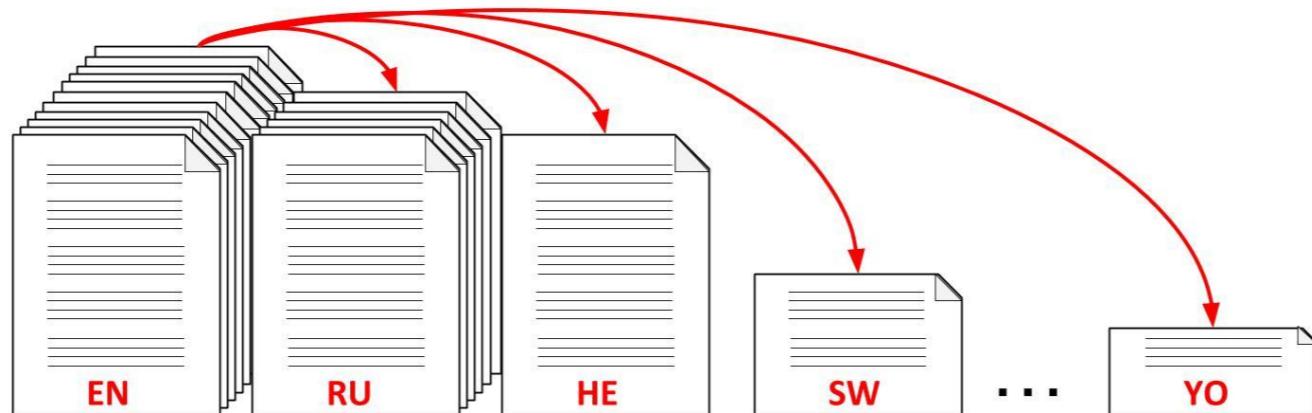


- **Instance** level transfer: transfer annotations (e.g., NER, POS tags) via cross-lingual bridges (e.g., word or phrase alignments)
- **Model** level transfer: train a model in a resource-rich language and adapt (e.g. fine-tune) it in a resource-poor language

# Approaches to Low-resource NLP



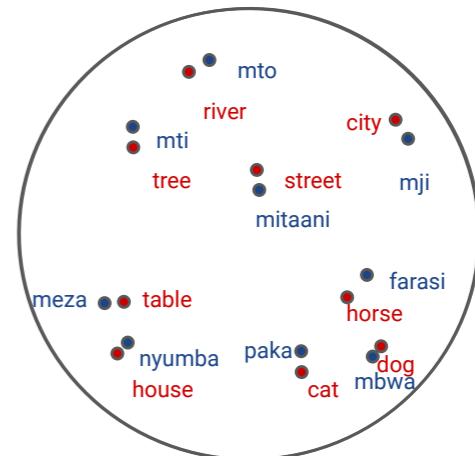
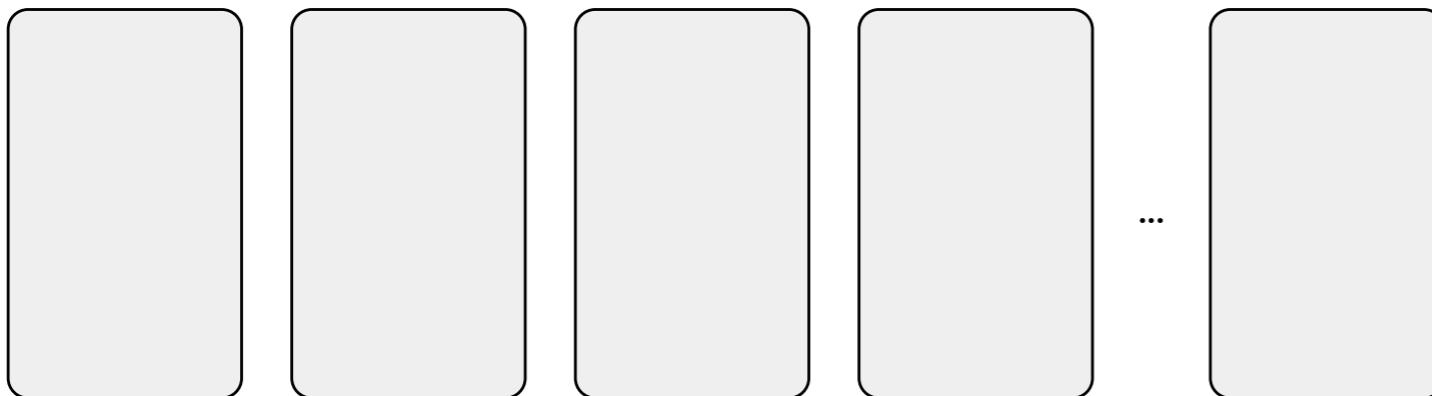
## 1. Cross-lingual transfer learning



- **Model level transfer:** train a model in a resource-rich language and adapt (e.g. fine-tune) it in a resource-poor language
  - **Zero-shot** learning – train a model in one domains and assume it generalizes more or less out-of-the-box in a low-resource domain
  - **Few shot** learning – train a model in one domain and use only few examples from a low-resource domain to adapt it

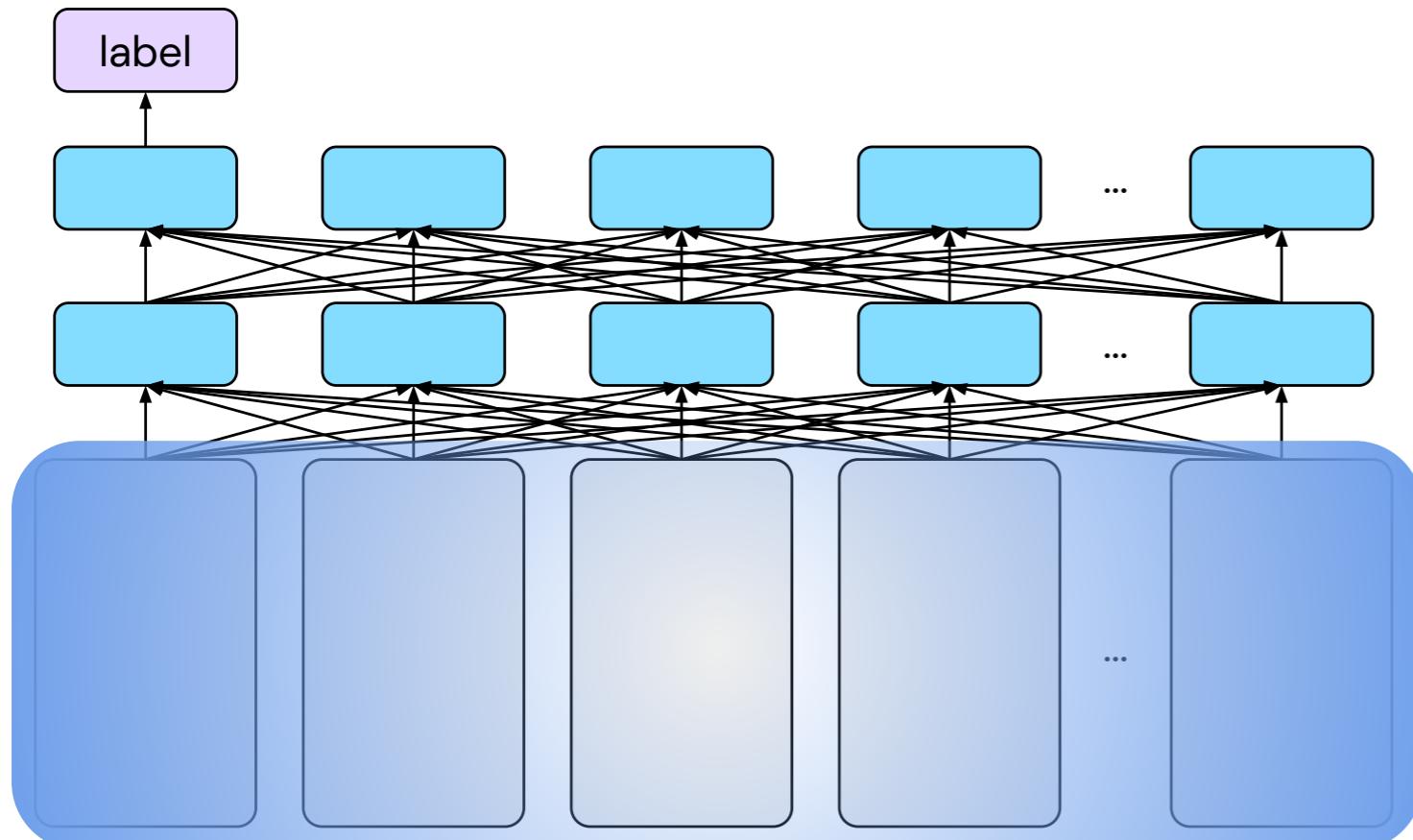
# Cross-lingual Model Learning Paradigm

- Step 1: Learn cross-lingual representations



# Cross-Lingual Model Learning Paradigm

entailment

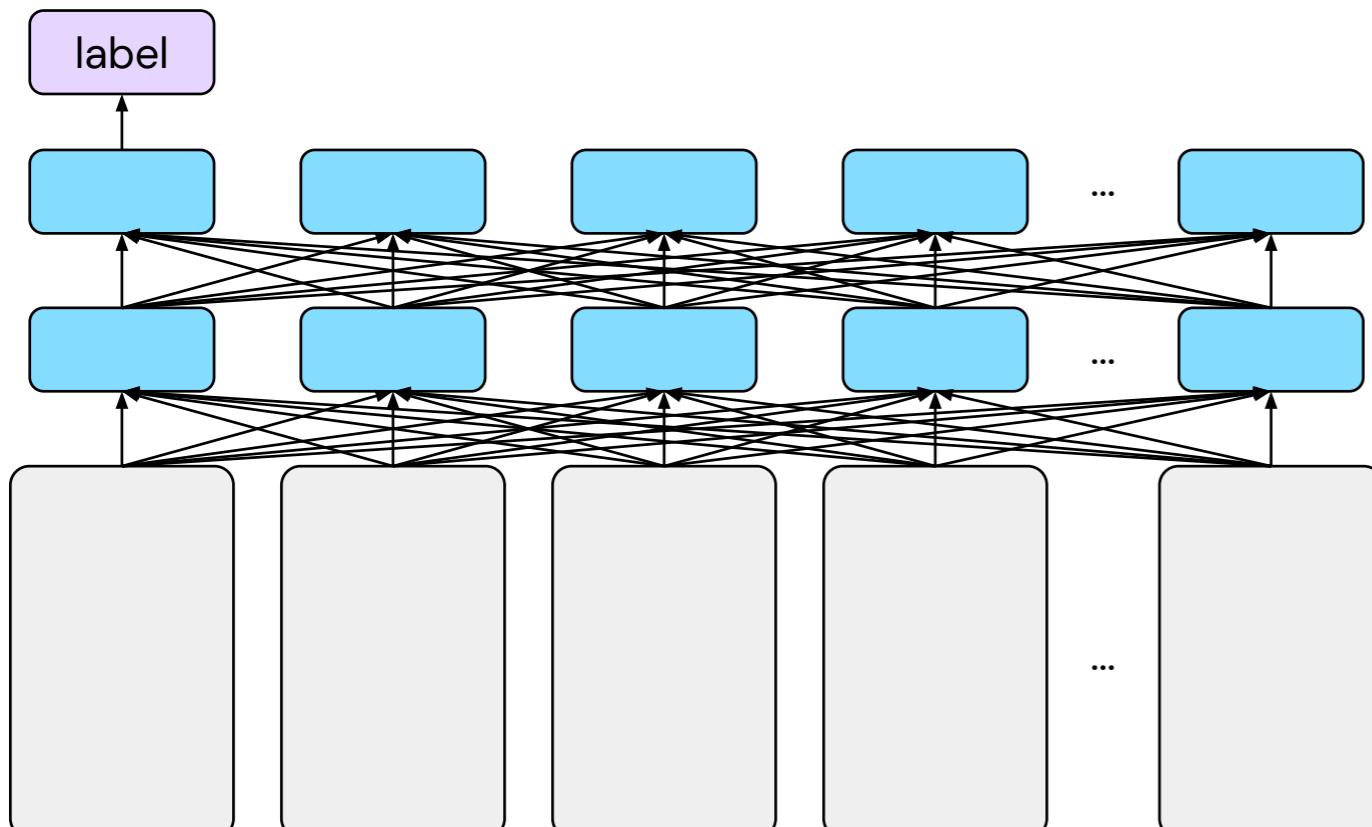


a soccer game with multiple males playing [SEP] some men are playing a sport

- Step 1: Learn cross-lingual representations
- Step 2: Fine-tune model on labelled data of a high-resource language (mostly English) *cross-lingual parameters are often kept fixed*

# Cross-Lingual Model Learning Paradigm

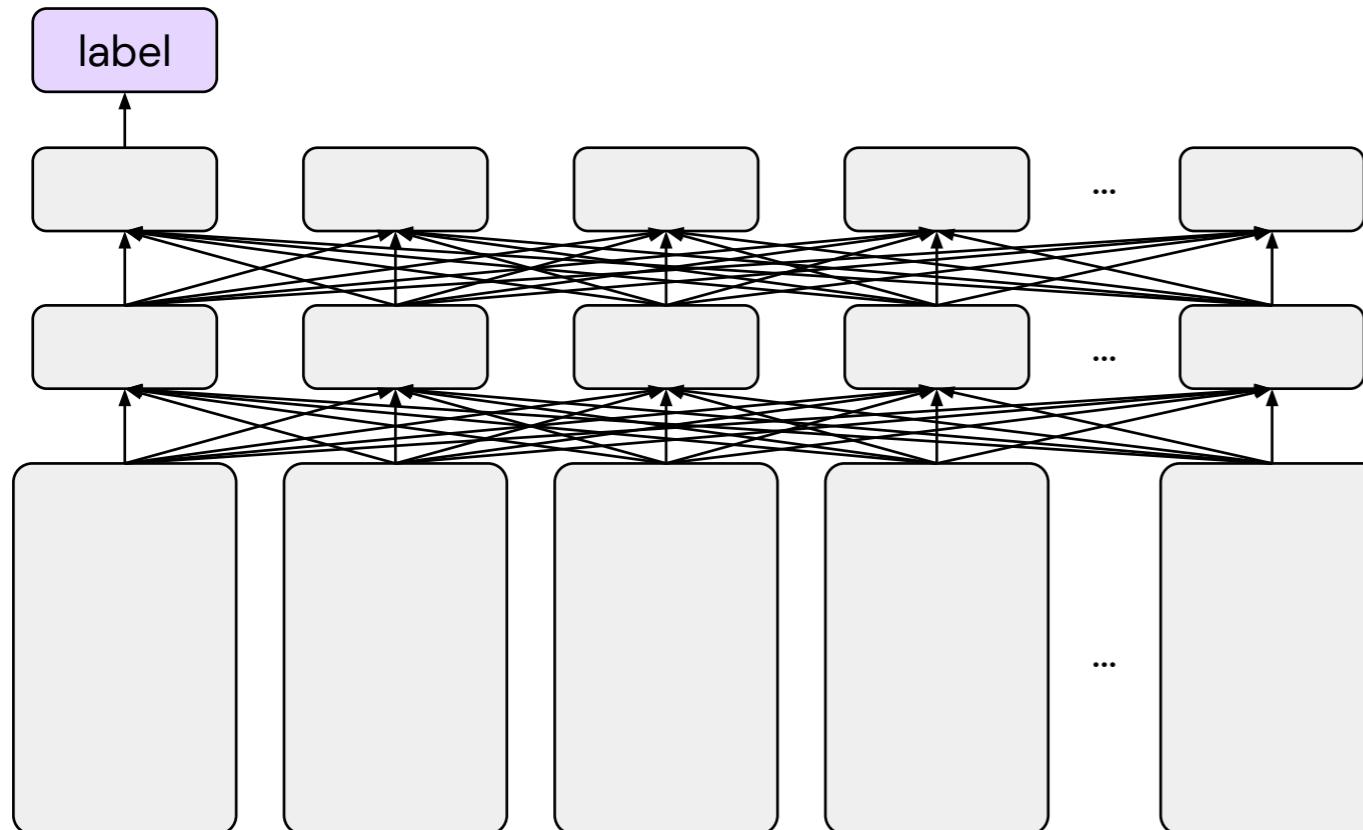
contradiction



el público se partió de risa [SEP] a nadie le hizo gracia

- Step 1: Learn cross-lingual representations
- Step 2: Fine-tune model on labelled data of a high-resource language
- Step 3: Zero-shot transfer to a low-resource language

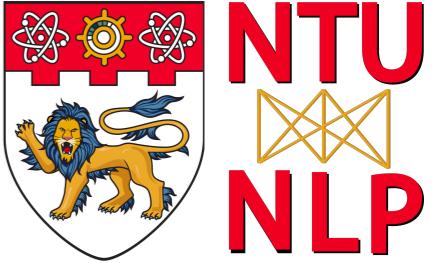
# Cross-lingual Learning Paradigm



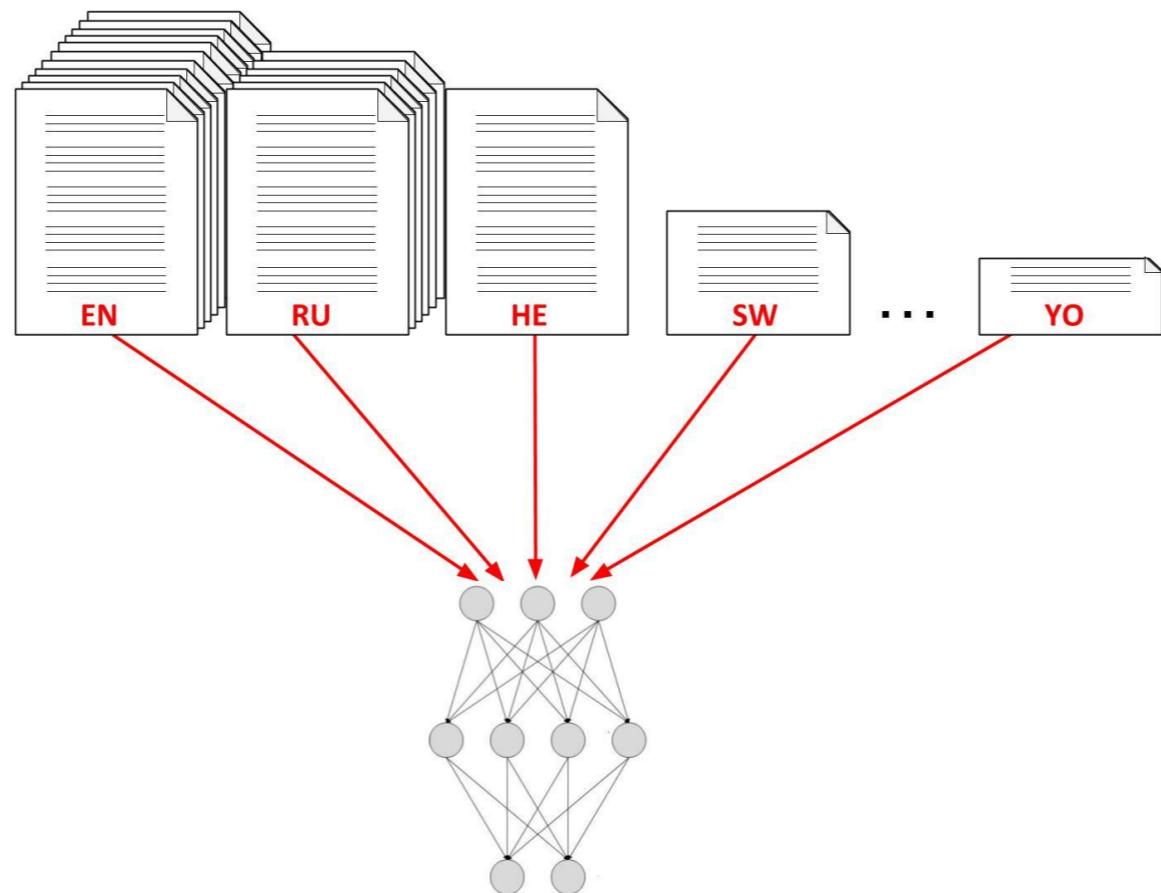
- Cross-lingual representations can be learned at different levels: from the word level to the sentence level

E.g., multi-lingual BERT

# Approaches to Low-resource NLP



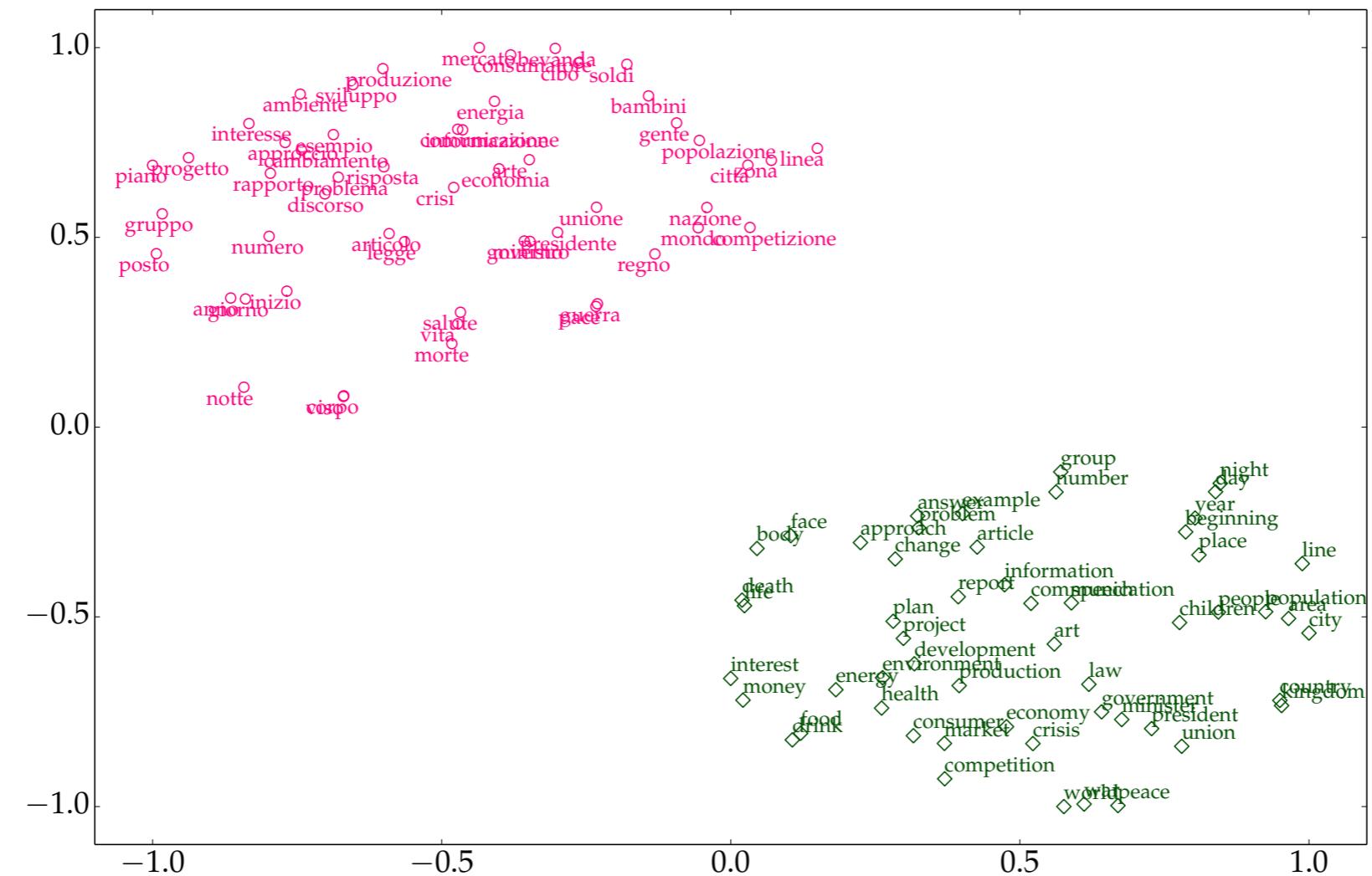
## 2. Joint Multi-lingual learning



Train a single model on a mix of datasets in all languages, to enable data and parameter sharing where possible

# Cross-lingual Word Embeddings

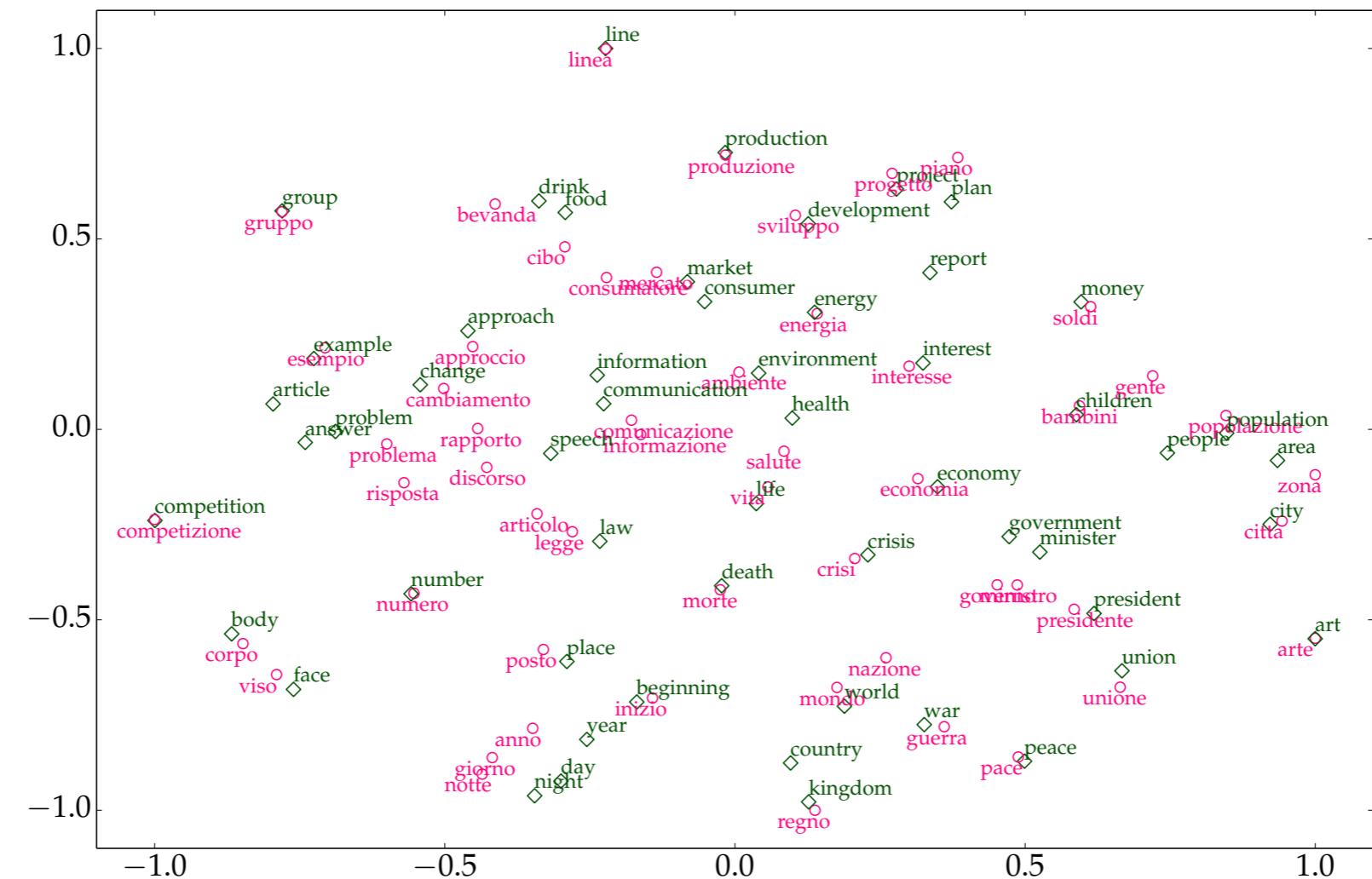
Cross-lingual representations of words in a joint embedding space



Two monolingual word embeddings (Ruder, 2019)

# Cross-lingual Word Embeddings

Cross-lingual representations of words in a joint embedding space



Cross-lingual word embeddings (Ruder et al, 2019)

# Cross-lingual Word Embeddings

## Why do we need cross-lingual embeddings?

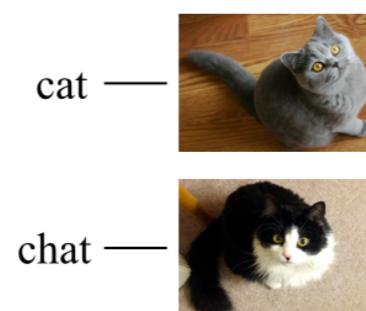
- Compare word vectors (meaning of words) across languages.
  - Key to bilingual lexicon induction and machine translation
- Transfer of knowledge (semantic & syntactic) between languages.
  - Cross-lingual NER, QA, Sentiment Analysis

# Cross-lingual Word Embeddings

What kind of data is used to learn cross-lingual embeddings?

	Parallel	Comparable
Word	Dictionaries	Images
Sentence	Translations	Captions
Document	-	Wikipedia

cat — chat  
dog — chien



The dog chases  
the cat.  
|  
Le chien poursuit  
le chat.

The dog chases the  
cat in the grass.  
|  
  
|  
Le chat s'envole  
du chien.

There are a lot of  
dogs in the park. They  
like to chase cats.  
|  
Les chats se relaxent.  
Ils fuient les chiens  
dès qu'ils les voient.

(a) Word, par.

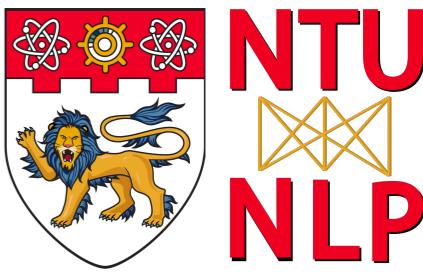
(b) Word, comp.

(c) Sentence, par.

(d) Sentence, comp.

(e) Doc., comp.

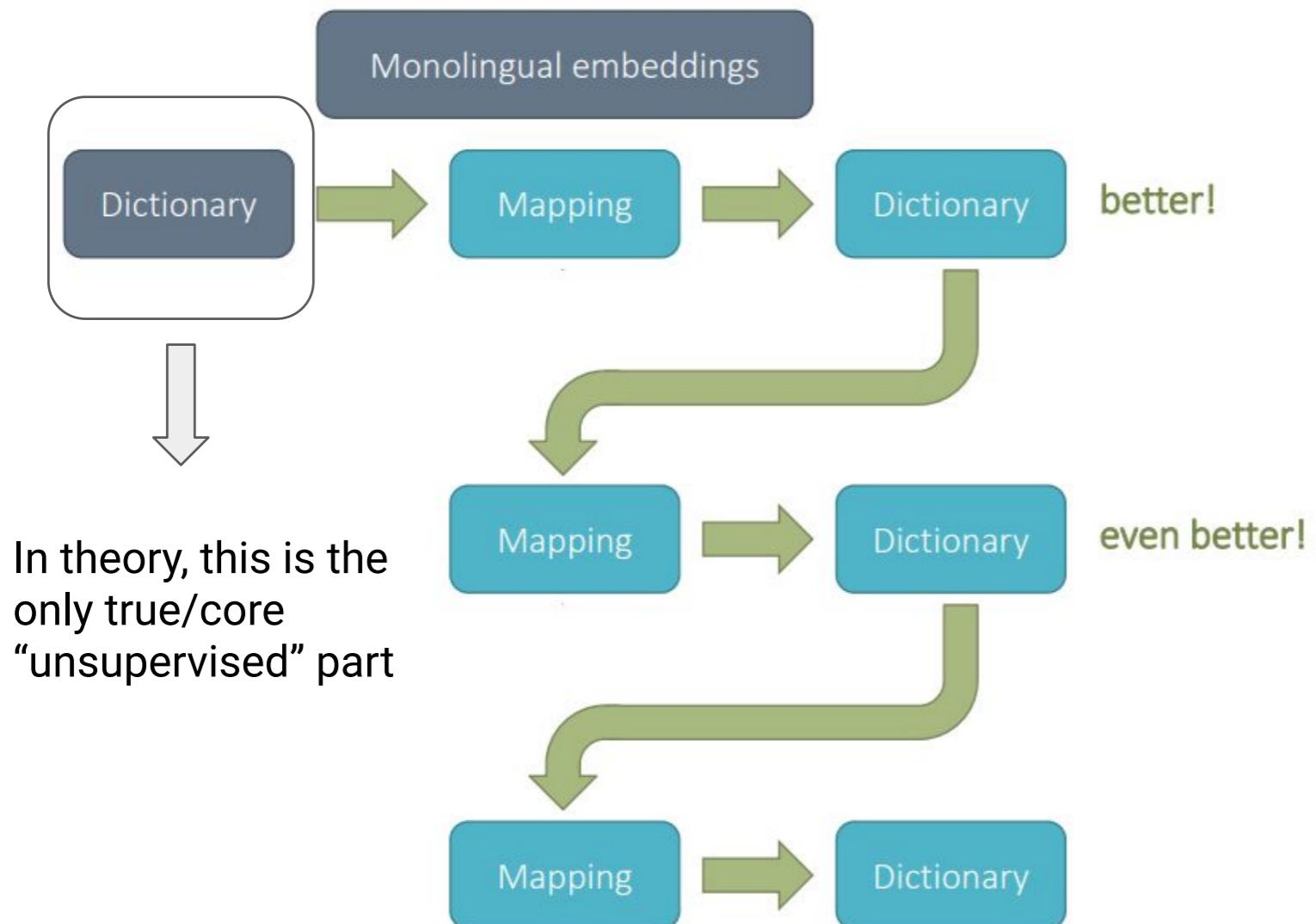
# Cross-lingual Embedding Learning Methods



	Parallel	Comparable
Word —Mapping	Mikolov et al. (2013b) Faruqui and Dyer (2014) Lazaridou et al. (2015) Dinu et al. (2015) Xing et al. (2015) Lu et al. (2015) Vulić and Korhonen (2016) Ammar et al. (2016b) Zhang et al. (2016b, 2017ab) Artexte et al. (2016, 2017, 2018ab) Smith et al. (2017) Hauer et al. (2017) Mrkšić et al. (2017b) Conneau et al. (2018a) Joulin et al. (2018) Alvarez-Melis and Jaakkola (2018) Ruder et al. (2018) Glavaš et al. (2019)	Bergsma and Van Durme (2011) Kiela et al. (2015) Vulić et al. (2016)
Word —Pseudo-bilingual	Xiao and Guo (2014) Duong et al. (2015) Gouws and Søgaard (2015) Duong et al. (2016) Adams et al. (2017)	
Word —Joint	Klementiev et al. (2012) Kočiský et al. (2014)	
Sentence —Matrix factorization	Zou et al. (2013) Shi et al. (2015) Gardner et al. (2015) Guo et al. (2015) Vyas and Carpuat (2016)	
Sentence —Compositional	Hermann and Blunsom (2013, 2014) Soyer et al. (2015)	
Sentence —Autoencoder	Lauly et al. (2013) Chandar et al. (2014)	
Sentence —Skip-gram	Gouws et al. (2015) Luong et al. (2015) Coulmance et al. (2015) Pham et al. (2015)	
Sentence —Other	Levy et al. (2017) Rajendran et al. (2016)	Calixto et al. (2017) Gella et al. (2017)
Document		Vulić and Moens (2013a, 2014, 2016) Søgaard et al. (2015) Mogadala and Rettinger (2016)

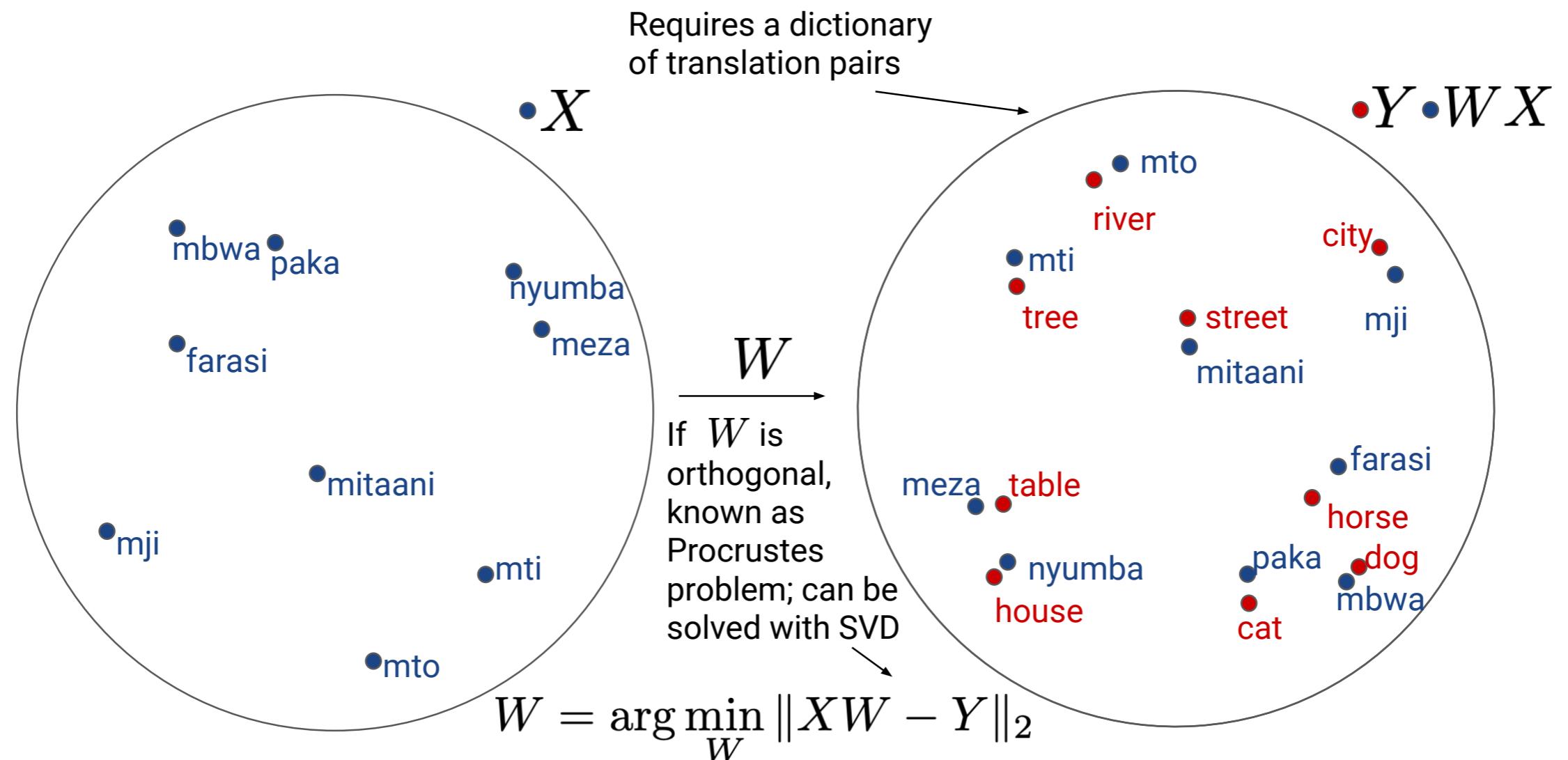
Check also more recent work including ours (see suggested readings)

# Supervised/Unsupervised Mapping based methods

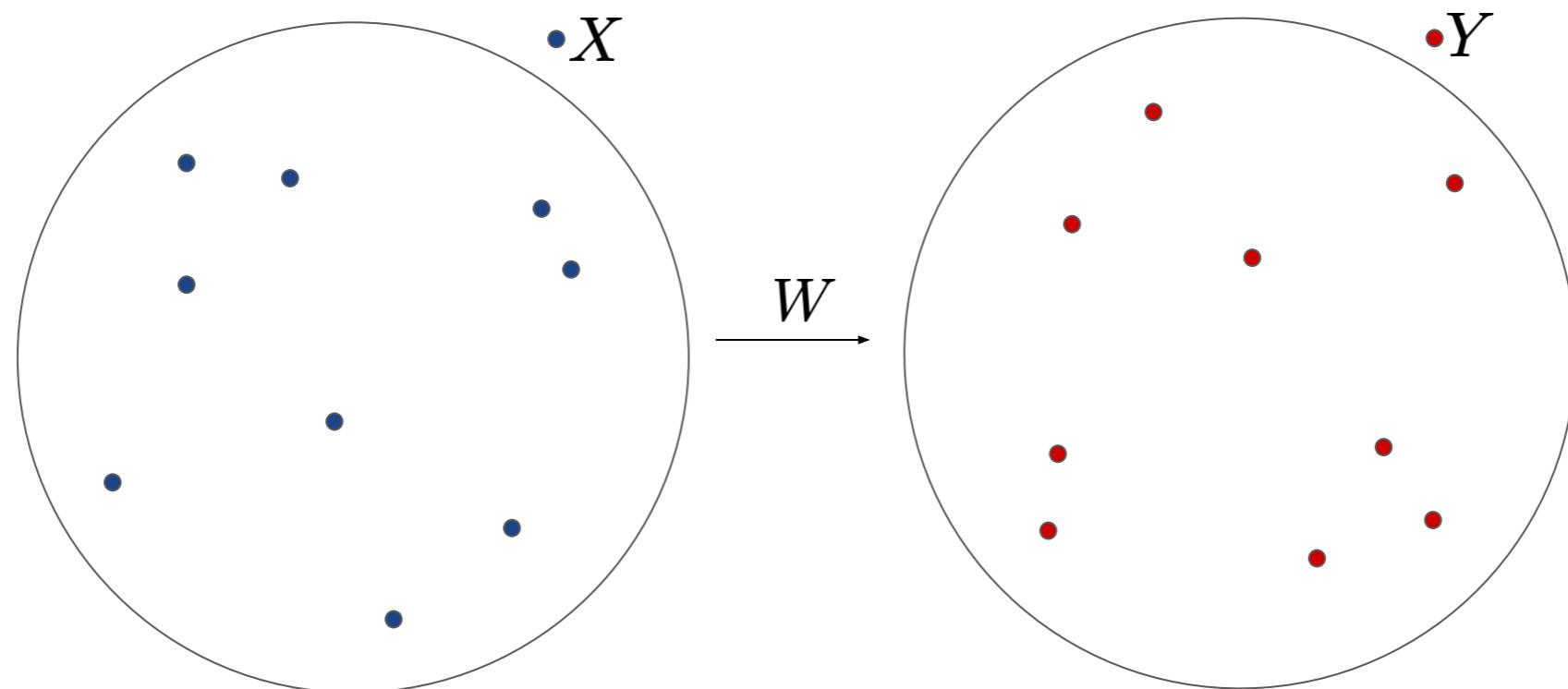


- Most current approaches use self-learning
- Particularly useful with limited supervision
- The seed dictionary improves over time
- Initial seed dictionary is either given (supervised) or learned (unsupervised)
- Many ways to learn seed dictionary
- One common approach: adversarial mapping

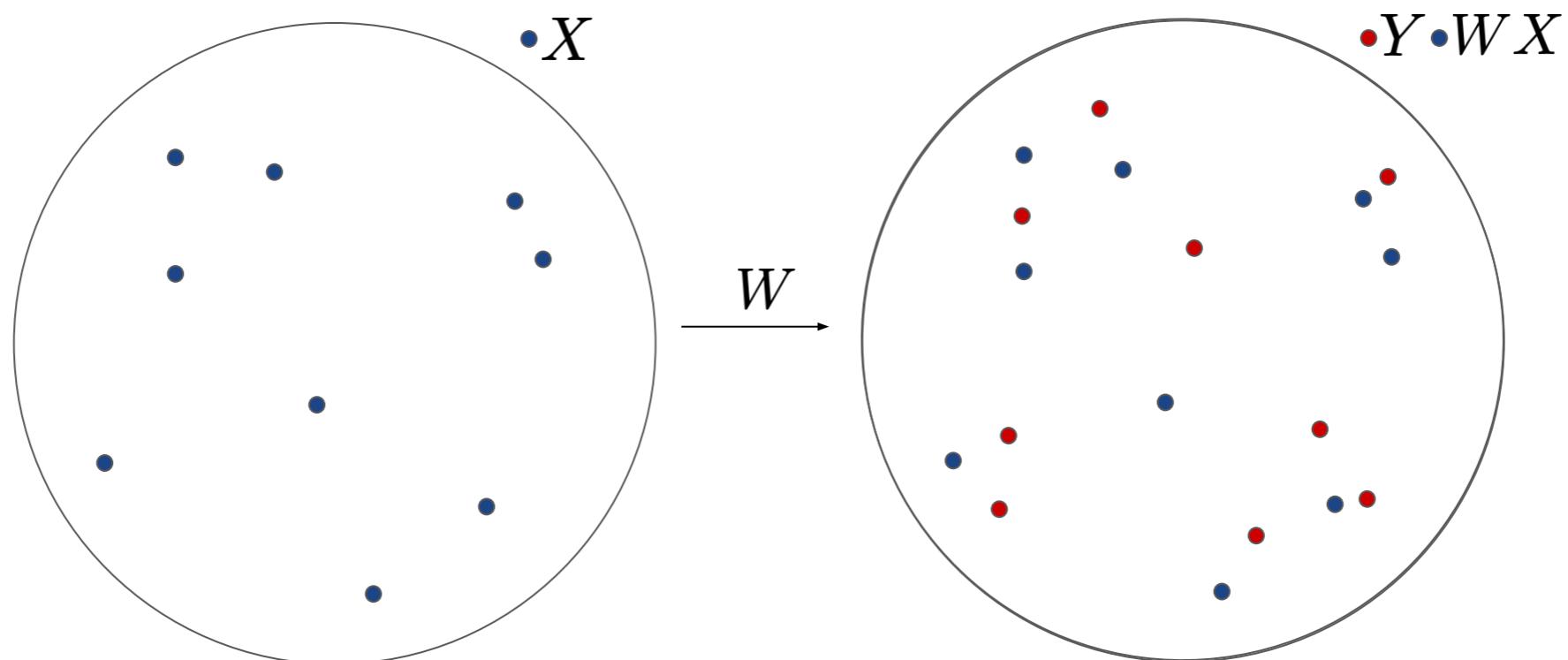
# Supervised Mapping based methods



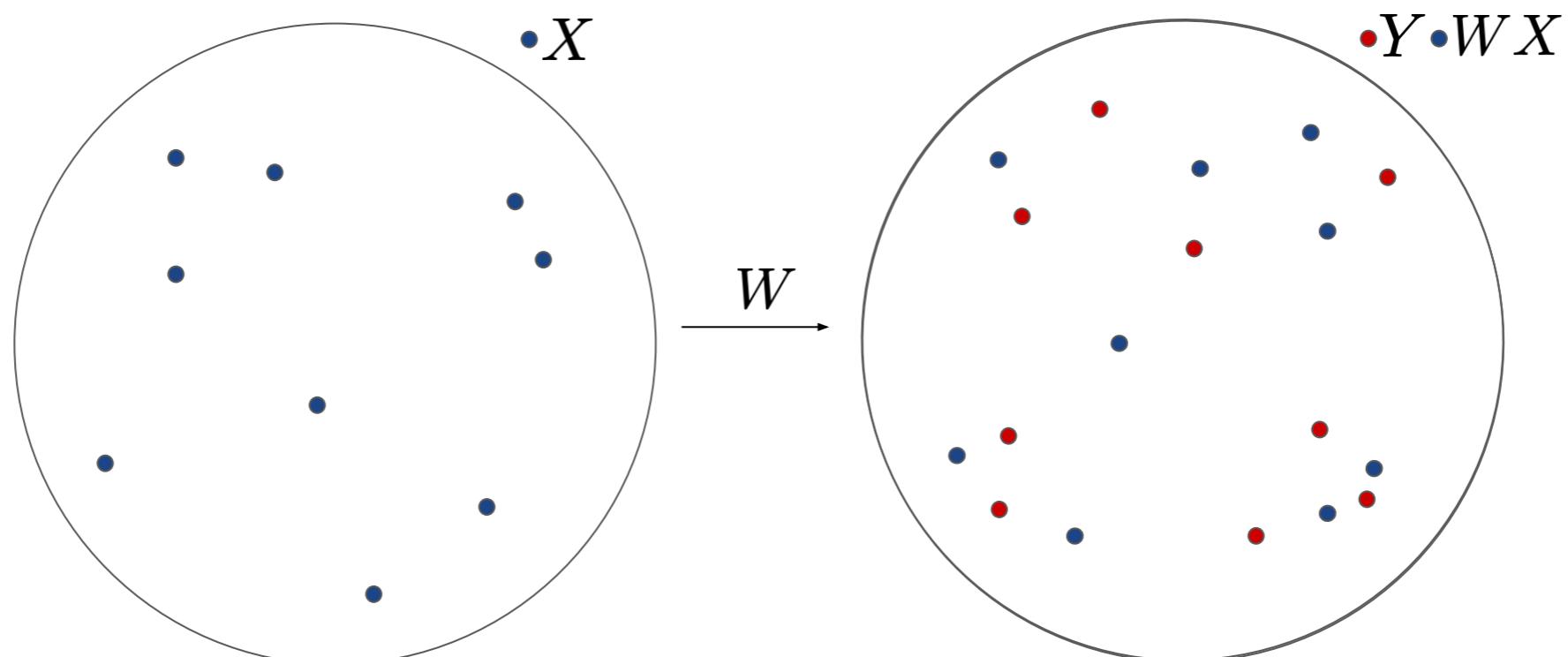
# Unsupervised Mapping



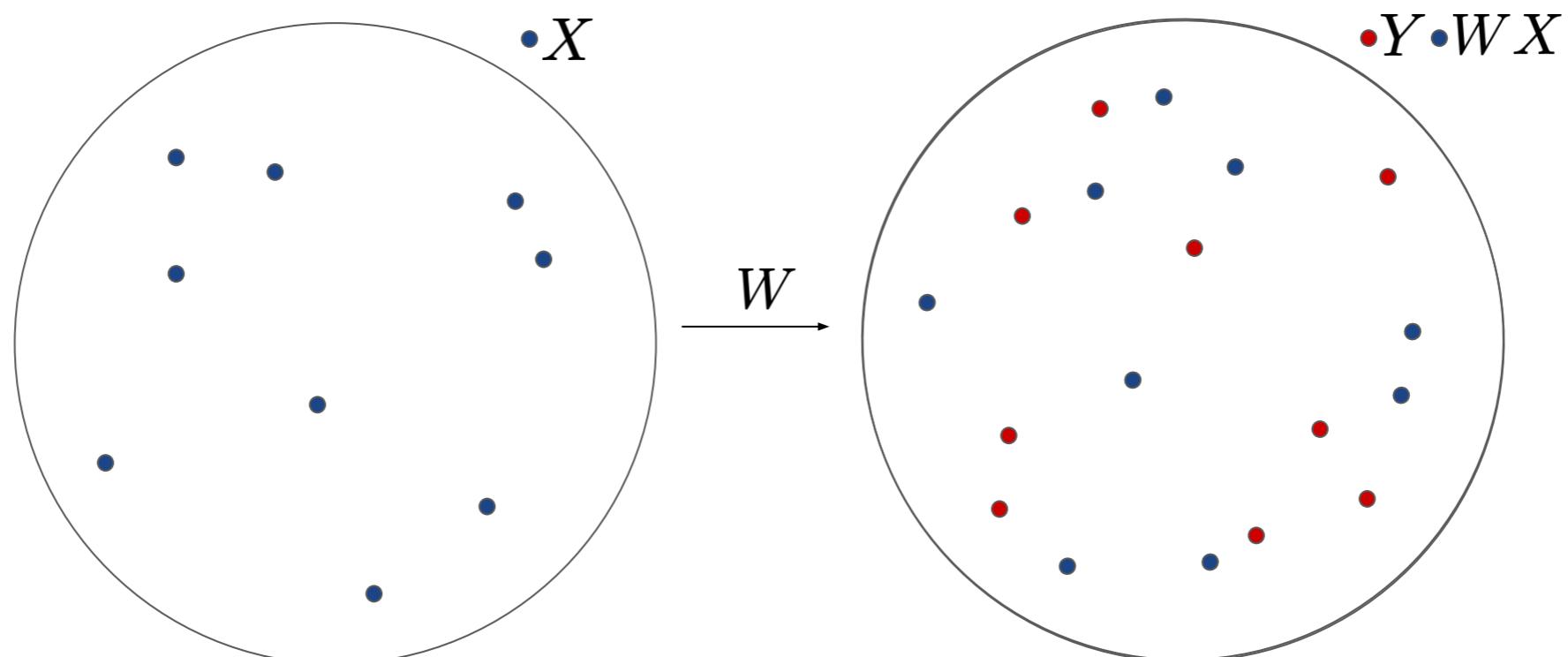
# Unsupervised Mapping



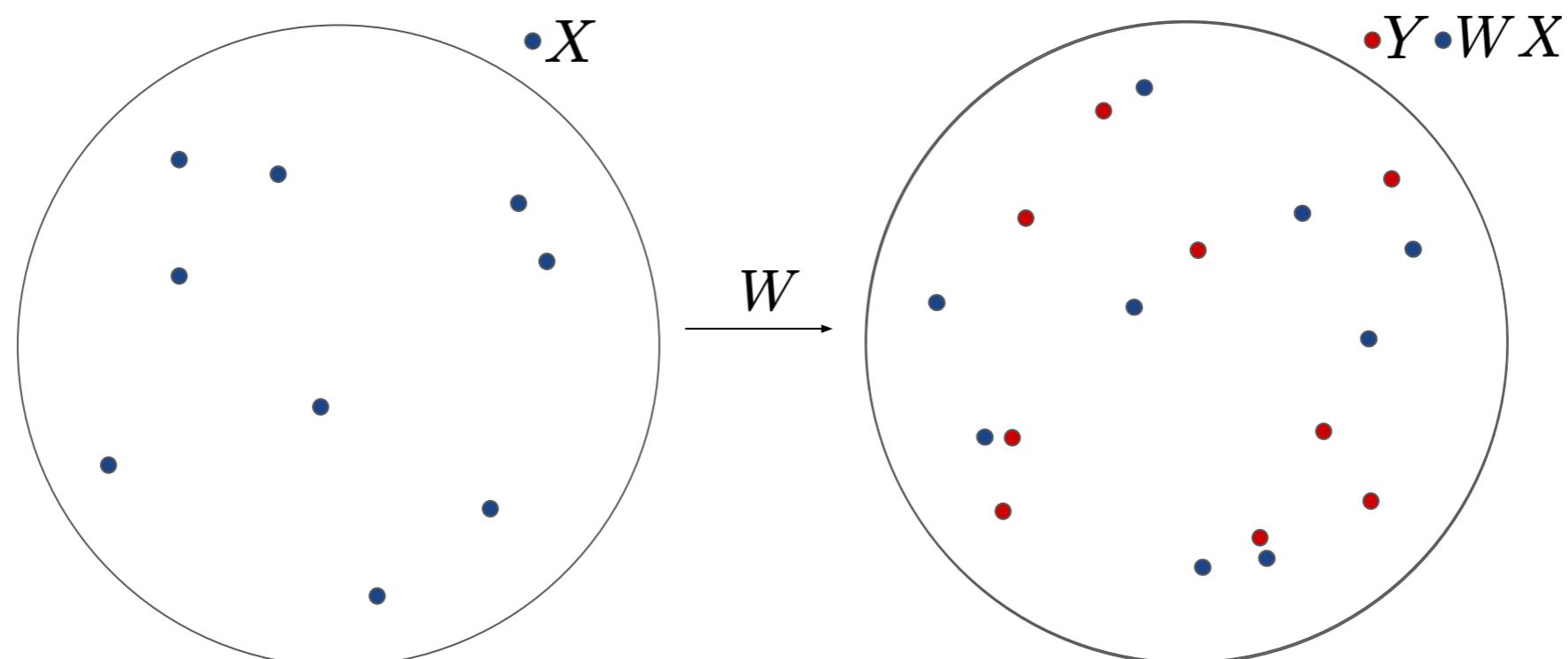
# Unsupervised Mapping



# Unsupervised Mapping

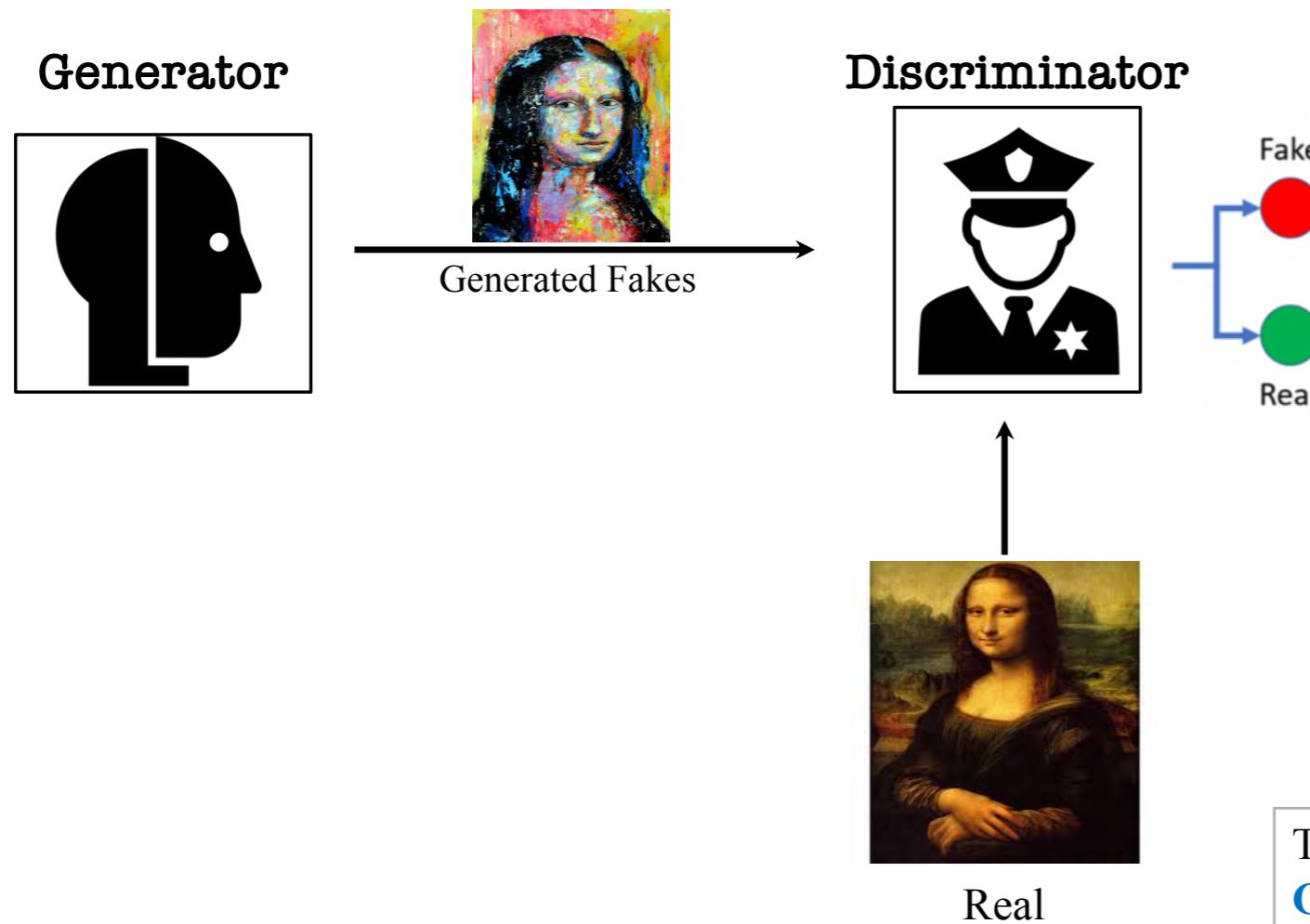


# Unsupervised Mapping



# Generative Adversarial Nets

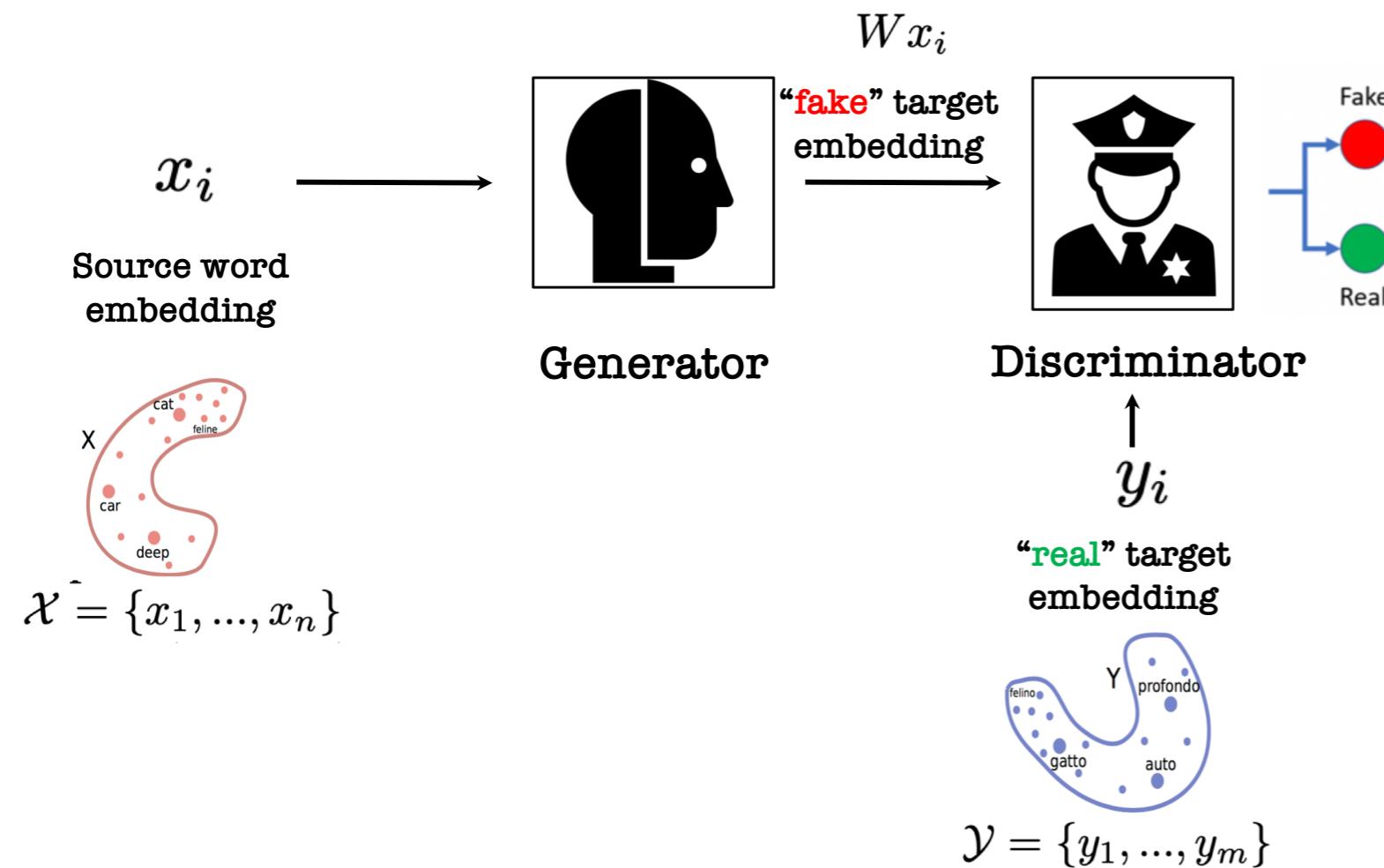
[Goodfellow et al., 2014]



Two player game, where  
**Generator**: generates **fake** examples  
**Discriminator**: differentiate between **real** and **fake**

# Applications of GANs in NLP

## Unsupervised Word Translation



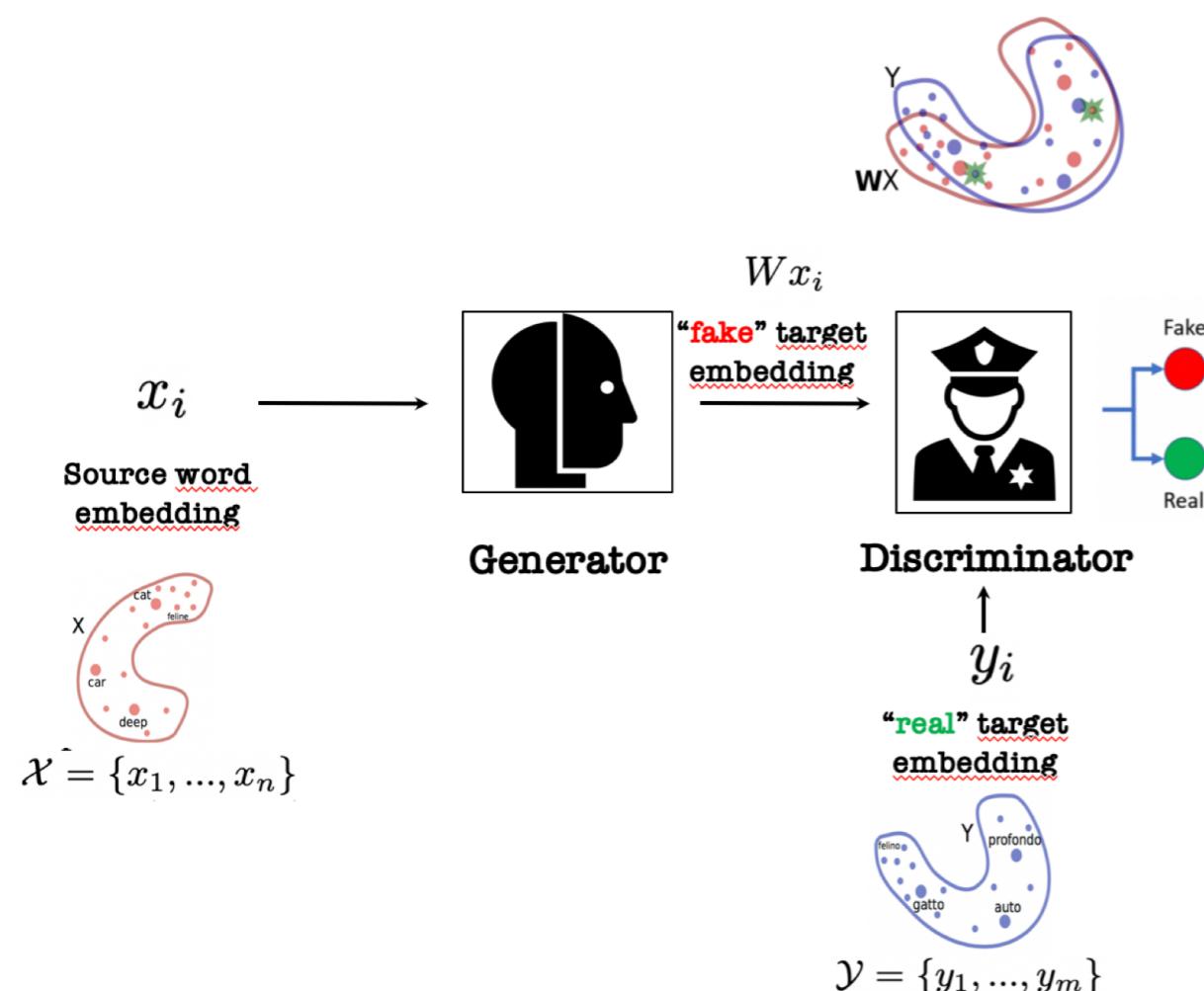
- Generator needs to **match distribution** of target language in order to **fool discriminator** consistently.
- **Hypothesis:** Best way to do this is to align words with their translations.

**Generator:** projects source word embedding  $x_i$  into the target language using linear mapper  $W \Rightarrow Wx_i$

**Discriminator:** differentiate between “fake” projected embeddings and “true” target language embeddings  $y_i$

# Applications of GANs in NLP

## Unsupervised Word Translation



### Objective

$\theta_D \Rightarrow$  Discriminator parameters

$P_{\theta_D}(\text{source} = 1|z) \Rightarrow$  Probability that z is the mapping of a source embedding

### Discriminator Loss

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

An embedding in  $X$

An embedding in  $Y$

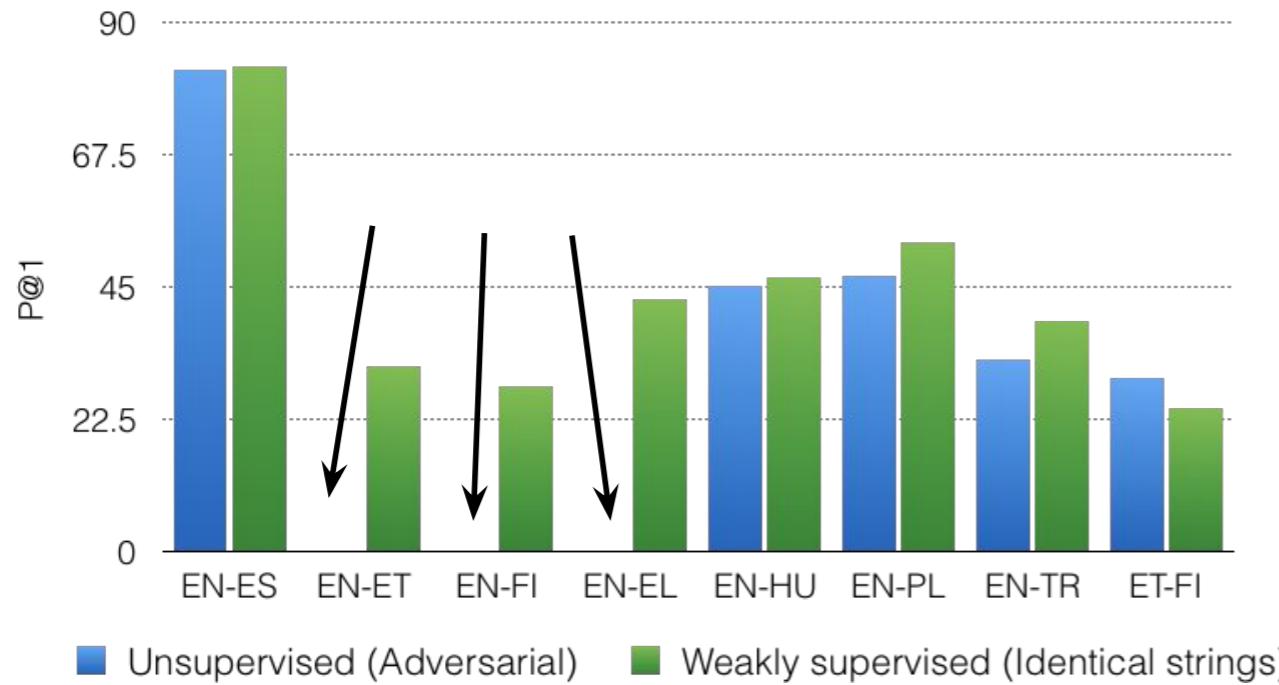
Maximize probability of predicting correct source

### Generator/Mapper Loss

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

Maximize probability of fooling discriminator

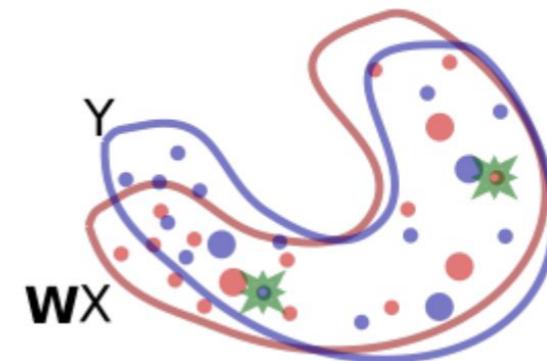
# Adversarial Mapping



[Søgaard et al. \(2018\)](#)

- More specifically, monolingual embedding spaces should be **approximately isomorphic**, i.e. same number of vertices, connected the same way
- Does not strictly hold** even for related languages

- Fail to map** between distant languages  
→ Why?
- Remember the unsupervised alignment step:



- Embedding spaces need to be **similar** for this to work

Our recent work (ACL, 2020) is independent of this assumption

# Evaluation of Cross Lingual Vectors

Word Translation /Dictionary Induction (standard/intrinsic task):

	En-Ms		En-Fi		En-Et		En-Tr		En-El		En-Fa		En-He		En-Ta		En-Bn		En-Hi	
	→	←	→	←	→	←	→	←	→	←	→	←	→	←	→	←	→	←	→	←
<b>GH Distance</b>	0.49		0.54		0.68		0.41		0.46		0.39		0.45		0.29		0.49		0.56	
<b>Unsupervised Baselines</b>																				
Artetxe et al. (2018b)	49.0	49.7	49.8	63.5	33.7	51.2	52.7	63.5	47.6	63.4	33.4	40.7	43.8	57.5	0.0	0.0	18.4	23.9	39.7	48.0
Conneau et al. (2018)	46.2	0.0	38.4	0.0	19.4	0.0	46.4	0.0	39.5	0.0	30.5	0.0	36.8	53.1	0.0	0.0	0.0	0.0	0.0	0.0
Supervision With “1K Unique” Seed Dictionary																				
<b>Supervised Baselines</b>																				
Artetxe et al. (2017)	36.5	41.0	40.8	56.0	21.3	39.0	39.5	56.5	34.5	56.2	24.1	35.7	30.2	51.7	5.4	12.7	6.2	19.9	22.6	38.8
Artetxe et al. (2018a)	35.3	34.0	30.8	40.8	21.6	32.6	33.7	43.3	32.0	46.4	22.8	27.6	32.27	39.1	7.3	11.9	11.3	15.7	26.2	30.7
Conneau et al. (2018)	46.2	44.7	46.0	58.4	29.3	40.0	44.8	58.5	42.1	56.5	31.6	38.4	38.3	52.4	11.7	16.0	14.3	19.7	32.5	42.3
Joulin et al. (2018)	31.4	30.7	30.4	41.4	20.1	26.0	30.7	36.5	28.8	43.6	18.7	23.1	33.5	34.3	6.0	10.1	7.6	11.3	20.7	25.7
LNMAP	<b>50.1</b>	<b>50.6</b>	<b>54.3</b>	<b>64.3</b>	<b>39.1</b>	<b>51.1</b>	<b>52.9</b>	<b>64.0</b>	<b>48.3</b>	<b>62.1</b>	<b>35.3</b>	<b>41.3</b>	<b>45.3</b>	<b>54.9</b>	<b>18.3</b>	<b>24.3</b>	<b>19.7</b>	<b>30.3</b>	<b>36.9</b>	<b>49.3</b>
Supervision With “5K Unique” Seed Dictionary																				
<b>Supervised Baselines</b>																				
Artetxe et al. (2017)	36.5	42.0	40.8	57.0	22.4	39.6	39.6	56.7	37.2	56.4	26.0	35.3	31.6	51.9	6.2	13.4	8.2	21.3	23.2	38.3
Artetxe et al. (2018a)	54.6	52.5	48.8	65.2	38.2	54.8	52.0	65.1	47.5	64.6	38.4	42.4	47.4	57.4	18.4	25.79	21.9	31.8	40.3	49.5
Conneau et al. (2018)	46.4	45.7	46.0	59.2	31.0	41.7	45.9	60.1	43.1	56.8	31.6	37.7	38.4	53.4	14.27	19.1	15.0	22.6	32.9	42.8
Joulin et al. (2018)	50.0	49.3	53.0	66.1	39.8	52.0	<b>54.0</b>	61.7	47.6	63.4	<b>39.6</b>	42.2	<b>53.0</b>	56.3	16.0	24.25	21.3	27.0	38.3	47.5
LNMAP	<b>50.2</b>	<b>54.6</b>	<b>54.7</b>	<b>69.2</b>	<b>41.8</b>	<b>57.2</b>	52.0	<b>66.6</b>	<b>49.7</b>	<b>65.5</b>	36.8	<b>44.5</b>	47.3	<b>58.7</b>	<b>20.7</b>	<b>30.7</b>	<b>22.2</b>	<b>36.2</b>	<b>39.1</b>	<b>51.6</b>
Supervision With “Whole” (“5K Unique” Source Words) Seed Dictionary																				
<b>Supervised Baselines</b>																				
Artetxe et al. (2017)	37.0	41.6	40.8	57.0	22.7	39.5	38.8	56.9	37.5	57.2	25.4	36.3	32.2	52.1	5.9	14.1	7.7	21.7	22.4	38.3
Artetxe et al. (2018a)	55.2	51.7	48.9	64.6	37.4	54.0	52.2	63.7	48.2	65.0	39.0	42.6	47.6	58.0	19.6	25.2	21.1	30.6	<b>40.4</b>	50.0
Conneau et al. (2018)	46.3	44.8	46.4	59.0	30.9	42.0	45.8	59.0	44.4	57.4	31.8	38.8	39.0	53.4	15.1	18.4	15.5	22.4	32.9	44.4
Joulin et al. (2018)	<b>51.4</b>	49.1	<b>55.6</b>	65.8	40.0	50.2	<b>53.8</b>	61.7	<b>49.1</b>	62.8	<b>40.5</b>	42.4	<b>52.2</b>	57.9	17.7	24.0	20.2	26.9	38.2	47.1
LNMAP	50.4	<b>56.0</b>	55.1	<b>71.5</b>	<b>41.9</b>	<b>57.5</b>	52.6	<b>66.3</b>	<b>49.1</b>	<b>67.1</b>	36.5	<b>43.7</b>	47.4	<b>61.0</b>	<b>19.9</b>	<b>32.0</b>	<b>22.5</b>	<b>36.6</b>	38.7	<b>52.7</b>

Table 1: Translation accuracy (P@1) on **low-resource** languages on **MUSE dataset** using fastText embeddings.

**MUSE:** <https://github.com/facebookresearch/MUSE>

source: (Mohiuddin et al, 2020)

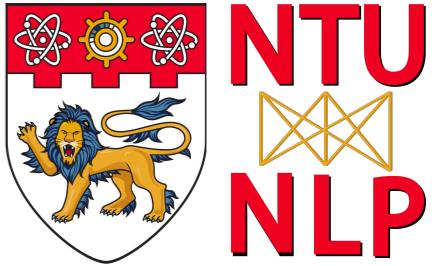
# Evaluation of Cross Lingual Vectors

Extrinsic tasks: Cross-lingual Natural Language Inference (XNLI)

	<i>Supervised</i>	Dict	EN-DE	EN-FR	EN-TR	EN-RU	Avg
PROC	1K	0.561	0.504	0.534	0.544	0.536	
PROC	5K	0.607	0.534	0.568	0.585	0.574	
PROC-B	1K	0.613	0.543	0.568	0.593	0.579	
PROC-B	3K	0.615	0.532	0.573	0.599	0.580	
DLV	5K	0.614	0.556	0.536	0.579	0.571	
RCSLS	1K	0.376	0.357	0.387	0.378	0.374	
RCSLS	5K	0.390	0.363	0.387	0.399	0.385	
	<i>Unsupervised</i>						
VECMAP		0.604	0.613	0.534	0.574	0.581	
MUSE		0.611	0.536	0.359*	0.363*	0.467	
ICP		0.580	0.510	0.400*	0.572	0.516	
GWA		0.427*	0.383*	0.359*	0.376*	0.386	

Other Extrinsic tasks:

- Cross-lingual Document Classification
- Cross-lingual Information Retrieval
- Cross-lingual Paraphrase Identification
- Cross-lingual Question Answering



mBERT ≠ XLM

# XLM (Lample & Conneau)

- Extend approach to multiple languages and show the effectiveness of cross-lingual pretraining
- Introduce a new **unsupervised method** for learning cross-lingual representations using XLM and investigate **two monolingual pretraining objectives**
- Introduce a new **supervised learning objective** that improves cross-lingual pretraining **when parallel data is available**
- Significantly outperform the previous state of the art on cross-lingual classification, unsupervised MT and supervised MT

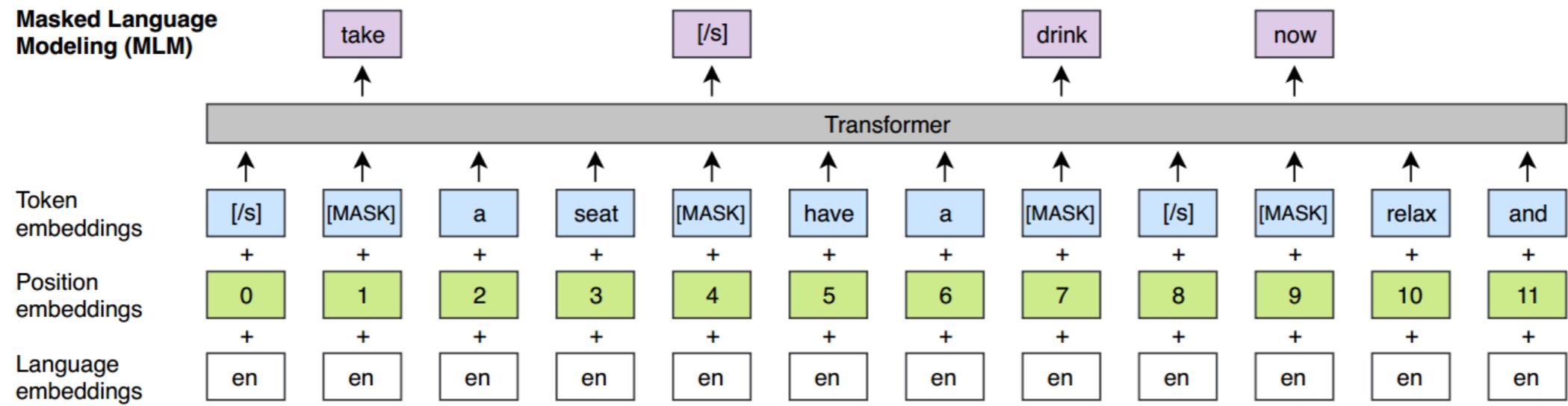
# XLM (Lample & Conneau)

- Proposed three language modeling objectives
- Two of them only require monolingual data (**unsupervised**), while the third one requires parallel sentences (**supervised**)

# XLM (Lample & Conneau)

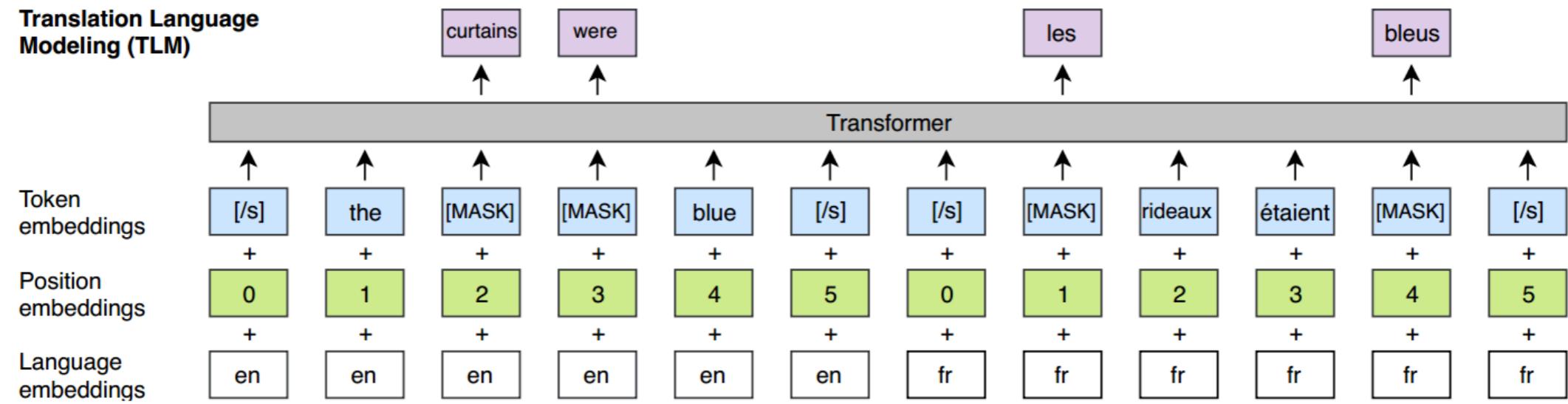
- Process all languages with the same shared vocabulary **created through BPE**
- Learn **BPE splits** on the concatenation of sentences sampled randomly from monolingual corpora.
- Sampling increases the number of tokens associated to low-resource languages and **alleviates bias towards high-resource languages**

# XLM: Pre-training



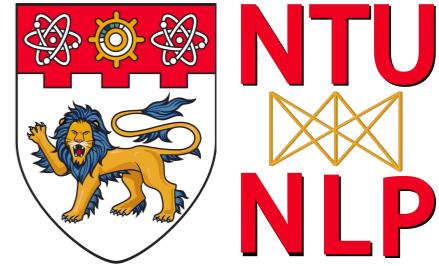
- Similar to **BERT**
- Instead of using pairs of sentences, use **text stream of arbitrary number of sentences** (truncated at 256 tokens)
- To counter the imbalance between rare and frequent tokens, they **subsample the frequent outputs**

# XLM: Pre-training



- Extension of MLM, where instead of considering monolingual text streams, **concatenate parallel sentences**
- Randomly mask words in both the source and target sentences
- For prediction, model can either attend to **surrounding same language words** or to the **other language words**
- This encourages the model to align the English and French representations

# XLM: Finetuning



## Cross-lingual classification

- XLM works as a **better initialization** of sentence encoders for zero-shot cross-lingual classification
- Add a **linear classifier** on top of the first hidden state of the pretrained Transformer, and fine-tune all parameters on the English NLI training dataset
- Evaluate the capacity of the model to make correct NLI predictions in the 15 XNLI languages

# XLM: Finetuning

## Machine Translation

- XLM works as a **better initialization** of supervised and unsupervised NMT systems
- Pretrain entire encoder and decoder with a cross-lingual language model

# XLM: Classification Results

- MultiNLI & XNLI

	en	fr	es	de	el	bg	ru	tr	ar	vi	th	zh	hi	sw	ur	$\Delta$
<i>Machine translation baselines (TRANSLATE-TRAIN)</i>																
Devlin et al. (2018)	81.9	-	77.8	75.9	-	-	-	-	70.7	-	-	76.6	-	-	61.6	-
XLM (MLM+TLM)	<u>85.0</u>	<u>80.2</u>	<u>80.8</u>	<u>80.3</u>	<u>78.1</u>	<u>79.3</u>	<u>78.1</u>	<u>74.7</u>	<u>76.5</u>	<u>76.6</u>	<u>75.5</u>	<u>78.6</u>	<u>72.3</u>	<u>70.9</u>	63.2	<u>76.7</u>
<i>Machine translation baselines (TRANSLATE-TEST)</i>																
Devlin et al. (2018)	81.4	-	74.9	74.4	-	-	-	-	70.4	-	-	70.1	-	-	62.1	-
XLM (MLM+TLM)	<u>85.0</u>	79.0	79.5	78.1	77.8	77.6	75.5	73.7	73.7	70.8	70.4	73.6	69.0	64.7	65.1	74.2
<i>Evaluation of cross-lingual sentence encoders</i>																
Conneau et al. (2018b)	73.7	67.7	68.7	67.7	68.9	67.9	65.4	64.2	64.8	66.4	64.1	65.8	64.1	55.7	58.4	65.6
Devlin et al. (2018)	81.4	-	74.3	70.5	-	-	-	-	62.1	-	-	63.8	-	-	58.3	-
Artetxe and Schwenk (2018)	73.9	71.9	72.9	72.6	73.1	74.2	71.5	69.7	71.4	72.0	69.2	71.4	65.5	62.2	61.0	70.2
XLM (MLM)	83.2	76.5	76.3	74.2	73.1	74.0	73.1	67.8	68.5	71.2	69.2	71.9	65.7	64.6	63.4	71.5
XLM (MLM+TLM)	<u>85.0</u>	<u>78.7</u>	<u>78.9</u>	<u>77.8</u>	<u>76.6</u>	<u>77.4</u>	<u>75.3</u>	<u>72.5</u>	<u>73.1</u>	<u>76.1</u>	<u>73.2</u>	<u>76.5</u>	<u>69.6</u>	<u>68.4</u>	<u>67.3</u>	<u>75.1</u>

# XLM: UNMT Results

		en-fr	fr-en	en-de	de-en	en-ro	ro-en
<i>Previous state-of-the-art - Lample et al. (2018b)</i>							
NMT		25.1	24.2	17.2	21.0	21.2	19.4
PBSMT		28.1	27.2	17.8	22.7	21.3	23.0
PBSMT + NMT		27.6	27.7	20.2	25.2	25.1	23.9
<i>Our results for different encoder and decoder initializations</i>							
EMB	EMB	29.4	29.4	21.3	27.3	27.5	26.6
-	-	13.0	15.8	6.7	15.3	18.9	18.3
-	CLM	25.3	26.4	19.2	26.0	25.7	24.6
-	MLM	29.2	29.1	21.6	28.6	28.2	27.3
CLM	-	28.7	28.2	24.4	30.3	29.2	28.0
CLM	CLM	30.4	30.0	22.7	30.5	29.0	27.8
CLM	MLM	32.3	31.6	24.3	32.5	31.6	29.8
MLM	-	31.6	32.1	<b>27.0</b>	33.2	31.8	30.5
MLM	CLM	<b>33.4</b>	32.3	24.9	32.9	31.7	30.4
MLM	MLM	<b>33.4</b>	<b>33.3</b>	26.4	<b>34.3</b>	<b>33.3</b>	<b>31.8</b>

# XLM: SNMT Results

Pretraining	-	CLM	MLM
Sennrich et al. (2016)	33.9	-	-
ro → en	28.4	31.5	35.3
ro ↔ en	28.5	31.5	35.6
ro ↔ en + BT	34.4	37.0	<b>38.5</b>

# mBART: Multilingual BART

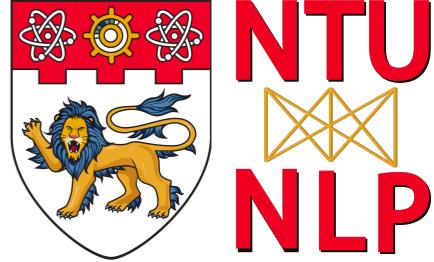
## mBART

- A Fully Trained
- Sequence-to-Sequence
- Denoising Auto-Encoder

pre-trained on large-scale **monolingual corpora** in many languages using the **BART** objective.

It can finetune both supervised and unsupervised settings.

# mBART: Multilingual BART

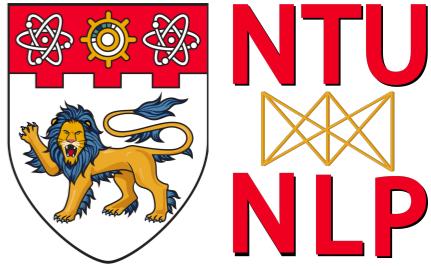


Large-Scale common crawl (CC) corpus.

- Short Infos

- Pre-train on a subset of 25 languages **CC25**
- Balance the corpus by  $\lambda_i = \frac{1}{p_i} \cdot \frac{p_i^\alpha}{\sum_i p_i^\alpha}$ ,  $\alpha = .7$
- 250000 sub-word tokens.
- token covers more lang. than the **CC25** for inference.
- No additional **pre-processing** (true- casing or normalizing punctuation/characters)

# mBART: Multilingual BART

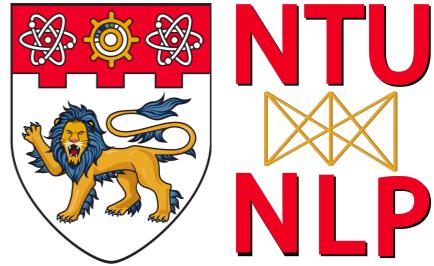


- Regular Transformer (Vaswani et al. 2017)
- 12 layer enc-dec, 1024 dim representation, 16 heads.
- 680M Params.

Noise function  $g$ ,

- Span masking.
- Mask 35% of the words in each instance.
- Random sampling a span length according to poison ( $\lambda = 3.5$ ) dist.
- permute the order of sentences within each instance.
- Decoder input starts with a <LID>

# mBART: Multilingual BART



- Sample a Language ID  $\textcolor{red}{<\text{LID}>}$
- Mask 35% of the words in each instance.
- Pack as many consecutive sentences as possible.
- Ends with either doc boundary or 512 token.

# mBART: Multilingual BART

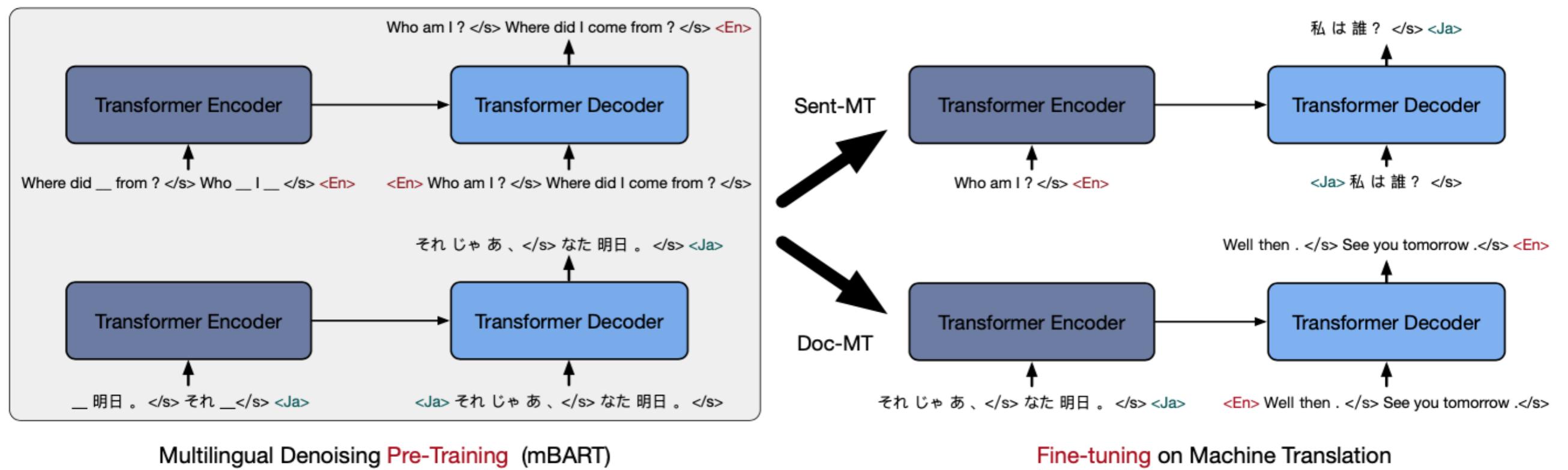
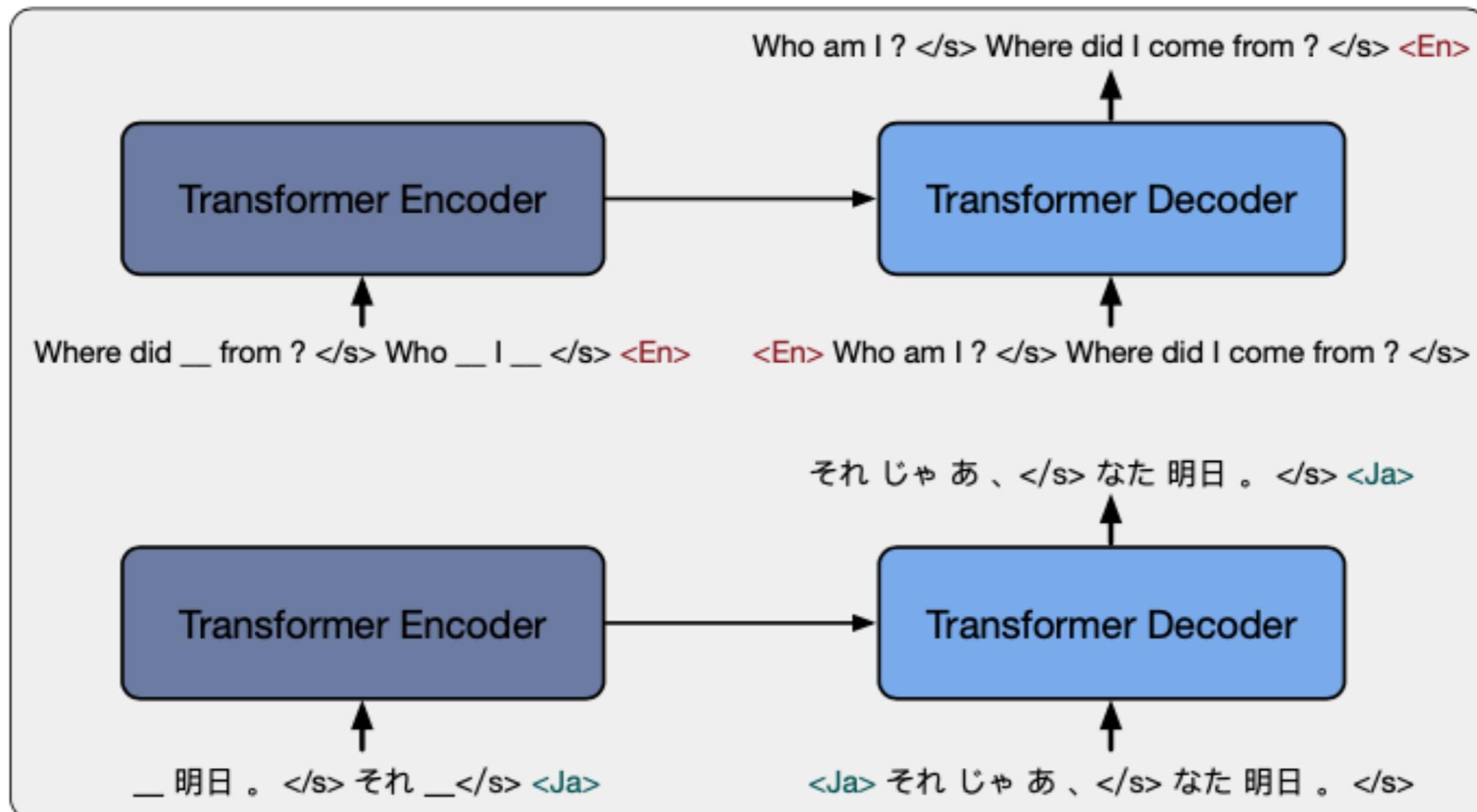


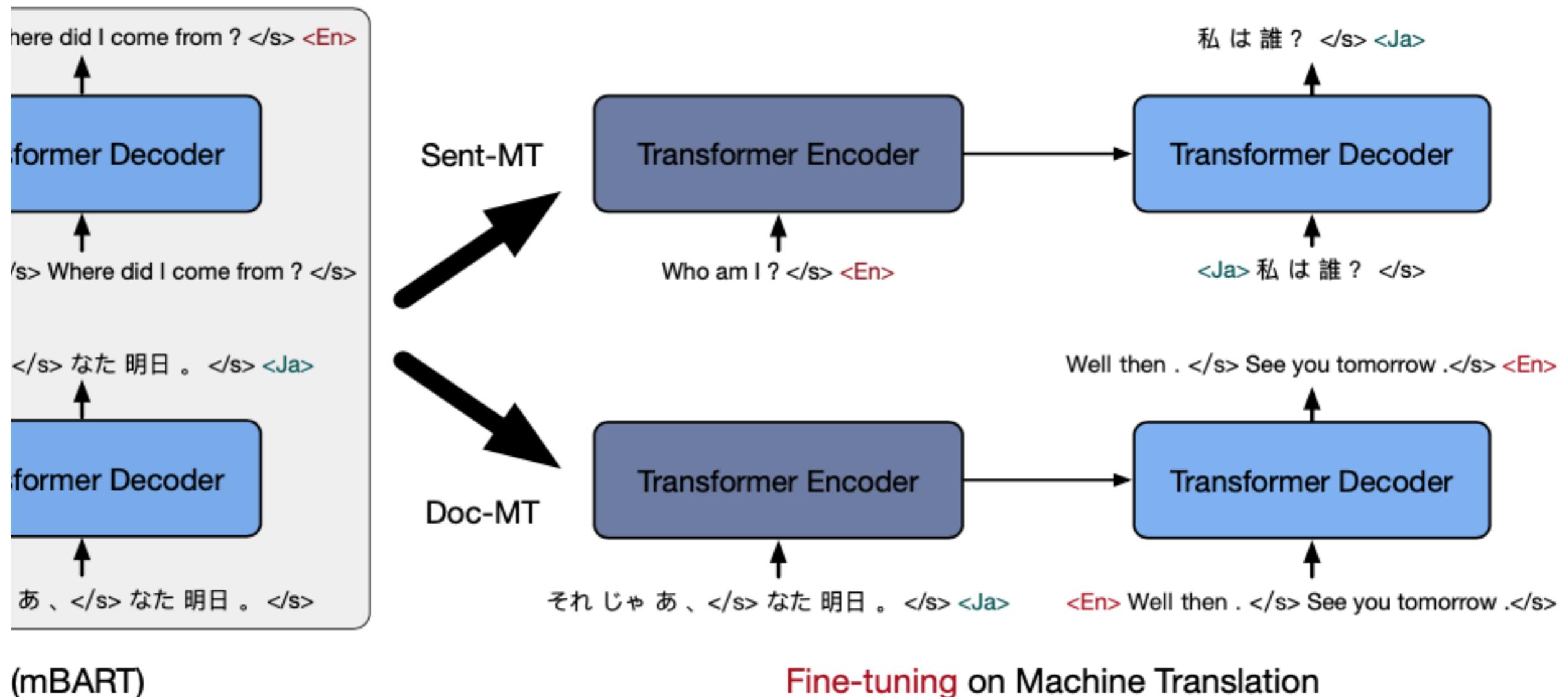
Figure 1: Framework for our Multilingual Denoising Pre-training (left) and fine-tuning on downstream MT tasks (right), where we use (1) sentence permutation (2) word-span masking as the injected noise. A special language id token is added at both the encoder and decoder. One multilingual pre-trained model is used for all tasks.

# mBART: Multilingual BART



Multilingual Denoising Pre-Training (mBART)

# mBART: Multilingual BART



# mBART: Results (Translation)

Languages	En-Gu		En-Kk		En-Vi		En-Tr		En-Ja		En-Ko	
Data Source	WMT19		WMT19		IWSLT15		WMT17		IWSLT17		IWSLT17	
Size	10K		91K		133K		207K		223K		230K	
Direction	←	→	←	→	←	→	←	→	←	→	←	→
<b>Random</b>	0.0	0.0	0.8	0.2	23.6	24.8	12.2	9.5	10.4	12.3	15.3	16.3
<b>mBART25</b>	<b>0.3</b>	<b>0.1</b>	<b>7.4</b>	<b>2.5</b>	<b>36.1</b>	<b>35.4</b>	<b>22.5</b>	<b>17.8</b>	<b>19.1</b>	<b>19.4</b>	<b>24.6</b>	<b>22.6</b>
Languages	En-Nl		En-Ar		En-It		En-My		En-Ne		En-Ro	
Data Source	IWSLT17		IWSLT17		IWSLT17		WAT19		FLoRes		WMT16	
Size	237K		250K		250K		259K		564K		608K	
Direction	←	→	←	→	←	→	←	→	←	→	←	→
<b>Random</b>	34.6	29.3	27.5	16.9	31.7	28.0	23.3	34.9	7.6	4.3	34.0	34.3
<b>mBART25</b>	<b>43.3</b>	<b>34.8</b>	<b>37.6</b>	<b>21.6</b>	<b>39.8</b>	<b>34.0</b>	<b>28.3</b>	<b>36.9</b>	<b>14.5</b>	<b>7.4</b>	<b>37.8</b>	<b>37.7</b>
Languages	En-Si		En-Hi		En-Et		En-Lt		En-Fi		En-Lv	
Data Source	FLoRes		ITTB		WMT18		WMT19		WMT17		WMT17	
Size	647K		1.56M		1.94M		2.11M		2.66M		4.50M	
Direction	←	→	←	→	←	→	←	→	←	→	←	→
<b>Random</b>	7.2	1.2	10.9	14.2	22.6	17.9	18.1	12.1	21.8	20.2	15.6	12.9
<b>mBART25</b>	<b>13.7</b>	<b>3.3</b>	<b>23.5</b>	<b>20.8</b>	<b>27.8</b>	<b>21.4</b>	<b>22.4</b>	<b>15.3</b>	<b>28.5</b>	<b>22.4</b>	<b>19.3</b>	<b>15.9</b>

# mBART: Results (Translation)

Languages	Cs	Es	Zh	De	Ru	Fr
Size	11M	15M	25M	28M	29M	41M
<b>Random</b>	16.5	33.2	<b>35.0</b>	<b>30.9</b>	<b>31.5</b>	<b>41.4</b>
<b>mBART25</b>	<b>18.0</b>	<b>34.0</b>	33.3	30.5	31.3	41.0

Table 3: **High Resource Machine Translation** where all the datasets are from their latest WMT competitions. We only evaluate our models on En-X translation.

# mBART: Results (Translation)

Languages	Cs	Es	Zh	De	Ru	Fr
Size	11M	15M	25M	28M	29M	41M
<b>Random</b>	16.5	33.2	<b>35.0</b>	<b>30.9</b>	<b>31.5</b>	<b>41.4</b>
<b>mBART25</b>	<b>18.0</b>	<b>34.0</b>	33.3	30.5	31.3	41.0

Table 3: **High Resource Machine Translation** where all the datasets are from their latest WMT competitions. We only evaluate our models on En-X translation.

# mBART: Results (Translation)

Model	Pre-training Data	Fine-tuning		
		En→Ro	Ro→En	+BT
<b>Random</b>	None	34.3	34.0	36.8
<b>XLM (2019)</b>	En Ro	-	35.6	38.5
<b>MASS (2019)</b>	En Ro	-	-	39.1
<b>BART (2019)</b>	En	-	-	38.0
<b>XLM-R (2019)</b>	CC100	35.6	35.8	-
<b>BART-En</b>	En	36.0	35.8	37.4
<b>BART-Ro</b>	Ro	37.6	36.8	38.1
<b>mBART02</b>	En Ro	<b>38.5</b>	<b>38.5</b>	<b>39.9</b>
<b>mBART25</b>	CC25	37.7	37.8	38.8

Table 4: Comparison with Other Pre-training Approaches on WMT16 Ro-En.

# mBART: Results (Translation)

	<b>Monolingual</b>	<b>Nl-En</b>	<b>En-Nl</b>	<b>Ar-En</b>	<b>En-Ar</b>	<b>Nl-De</b>	<b>De-Nl</b>
<b>Random</b>	None	34.6 (-8.7)	29.3 (-5.5)	27.5 (-10.1)	16.9 (-4.7)	21.3 (-6.4)	20.9 (-5.2)
<b>mBART02</b>	En Ro	41.4 (-2.9)	34.5 (-0.3)	34.9 (-2.7)	21.2 (-0.4)	26.1 (-1.6)	25.4 (-0.7)
<b>mBART06</b>	En Ro Cs It Fr Es	43.1 (-0.2)	34.6 (-0.2)	37.3 (-0.3)	21.1 (-0.5)	26.4 (-1.3)	25.3 (-0.8)
<b>mBART25</b>	All	<b>43.3</b>	<b>34.8</b>	<b>37.6</b>	<b>21.6</b>	<b>27.7</b>	<b>26.1</b>

Table 7: **Generalization to Unseen Languages** Language transfer results, fine-tuning on language-pairs without pre-training on them. mBART25 uses all languages during pre-training, while other settings contain at least one unseen language pair. For each model, we also show the gap to mBART25 results.