

# ATNLP: Multilinguality

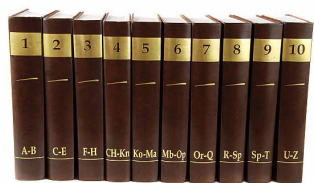
---

Anders Søgaard

coASfal



What are language  
models for?







+





+



+



	4	4-6	4-6-6	4-6-6-3
good	0.2	0.2	0.2	0.1
image	0.1	0.1		
home	0.1	0.1	0.1	0.1
hope	0.2	0.2	0.2	



He \_\_\_\_\_  
 {Legal} He \_\_\_\_\_  
 {9th century}, He \_\_\_\_\_  
 {Translation, English-German}, He \_\_\_\_\_  
 Translate 'He' into German...

walks, talks, ...  
 appeals, adjudicates, ...  
 sayeth, hath, ...  
 Er  
 Er

2017

2018

2019

2021

## Transformers

## Conditional language models

## Seeing everything as question answering

## In-context learning

Good scaling properties, exploiting GPUs, ability to model long-range dependencies.

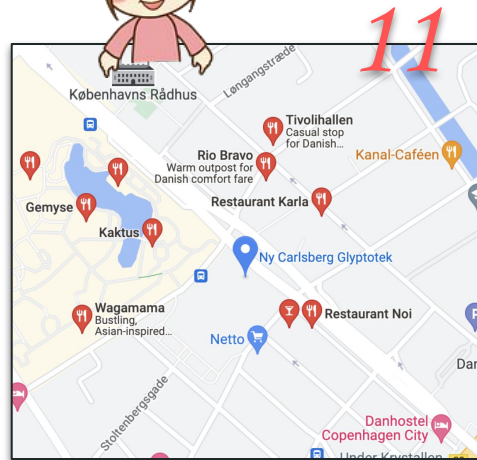
Language prefixes in machine translation, transfer learning, style transfer.

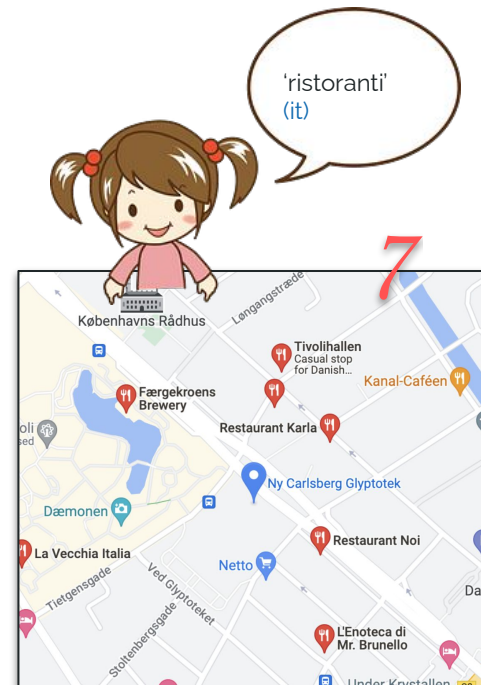
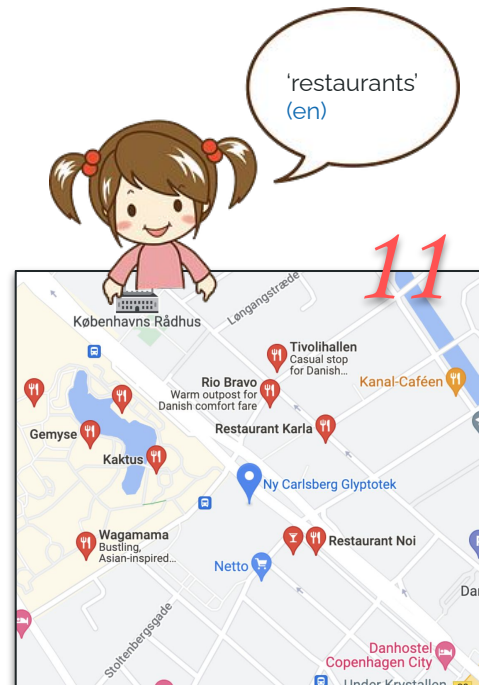
What is the German translation (or the grammatical analysis) of 'Mary bought a house'?

Moving away from training one model for each task, toward training models that can learn on the fly.

# The Digital Language Divide

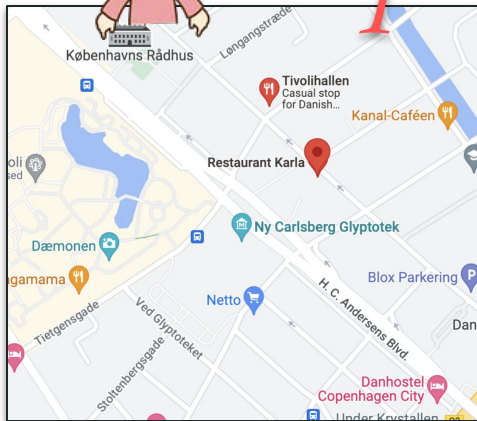








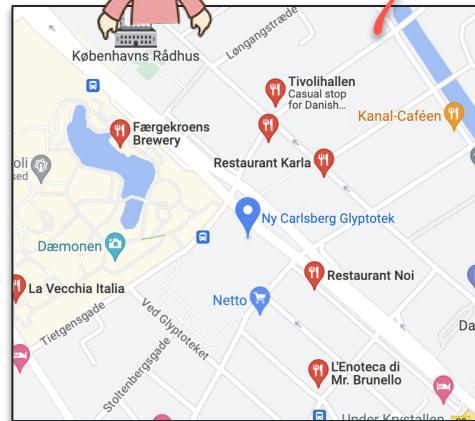
1



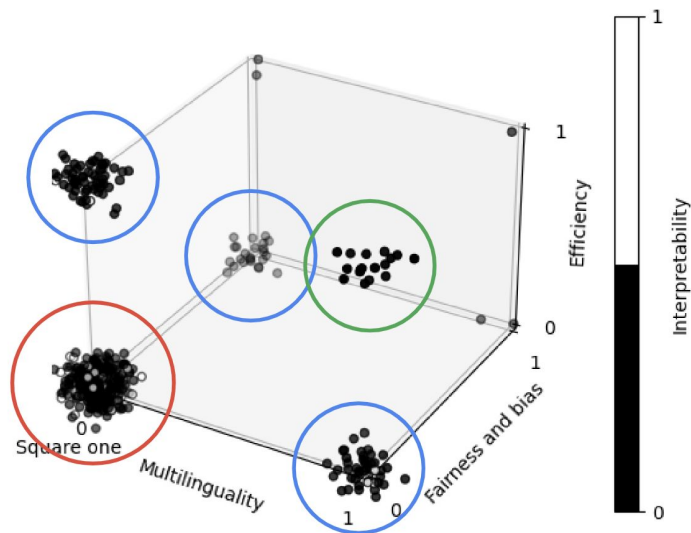
11



7

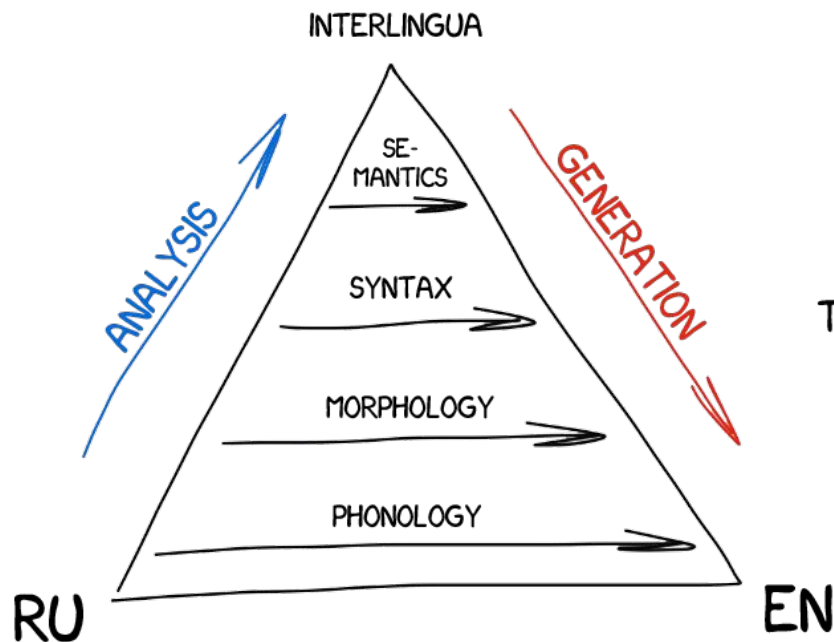


Area	# papers	English	Accuracy / F1	Multilinguality	Fairness and bias	Efficiency	Interpretability	>1 dimension
ACL 2021 oral papers	461	69.4%	38.8%	13.9%	6.3%	17.8%	11.7%	6.1%
MT and Multilinguality	58	0.0%	15.5%	56.9%	5.2%	19.0%	6.9%	13.8%
Interpretability and Analysis	18	88.9%	27.8%	5.6%	0.0%	5.6%	66.7%	5.6%
Ethics in NLP	6	83.3%	0.0%	0.0%	100.0%	0.0%	0.0%	0.0%
Dialog and Interactive Systems	42	90.5%	21.4%	0.0%	9.5%	23.8%	2.4%	2.4%
Machine Learning for NLP	42	66.7%	40.5%	19.0%	4.8%	50.0%	4.8%	9.5%
Information Extraction	36	80.6%	91.7%	8.3%	0.0%	25.0%	5.6%	8.3%
Resources and Evaluation	35	77.1%	42.9%	5.7%	8.6%	5.7%	14.3%	5.7%
NLP Applications	30	73.3%	43.3%	0.0%	10.0%	20.0%	10.0%	0.0%
Sentiment Analysis	18	100.0%	72.2%	0.0%	0.0%	11.1%	11.1%	0.0%
Summarization	12	91.7%	0.0%	0.0%	8.3%	0.0%	8.3%	0.0%



# A Brief History of Machine Translation

# VAUQUOIS TRIANGLE

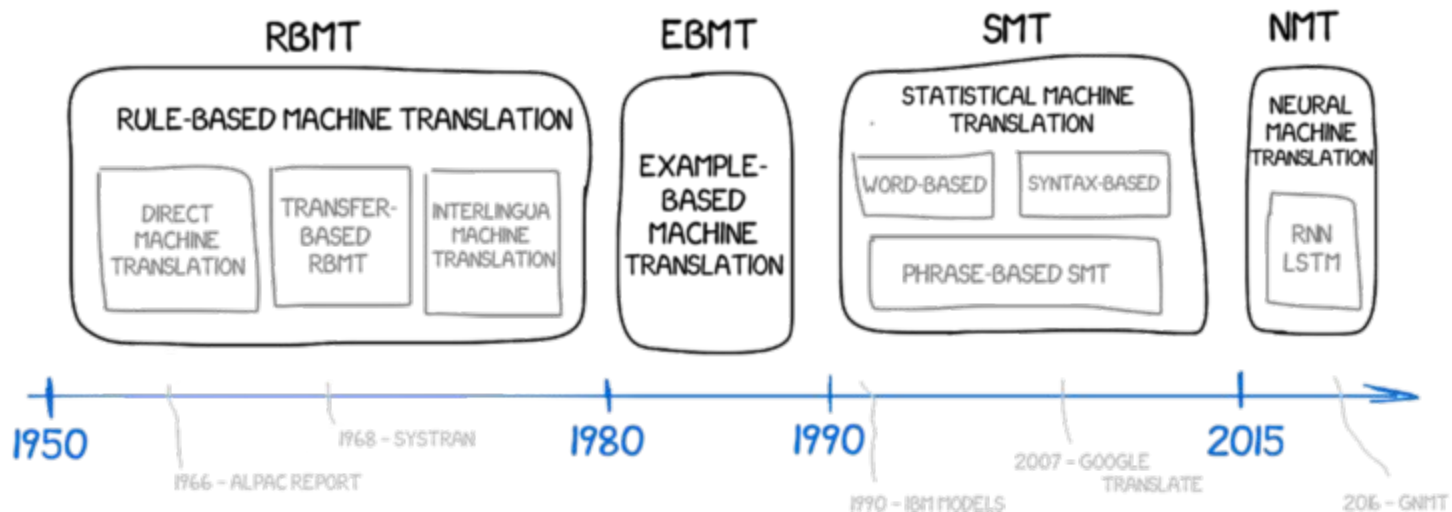


INTERLINGUA TRANSLATION

TRANSFER-BASED TRANSLATION

DIRECT TRANSLATION

## A BRIEF HISTORY OF MACHINE TRANSLATION



I		WANT		FORTY		KILOGRAMS OF		PERSIMMONS
↓		↓		↓		↓		↓
ICH		WOLLEN		VIERZIG		KILOGRAMM		PERSIMONEN



I | CRAVE | 40 | KG | KAKI

I		WANT		FORTY		KILOGRAMS OF		PERSIMMONS
↓		↓		↓		↓		↓
ICH		WOLLEN		VIERZIG		KILOGRAMM		PERSIMONEN

I



ICH

WANT



WILL

FORTY



VIERZIG

KILOGRAMS OF PERSIMMONS

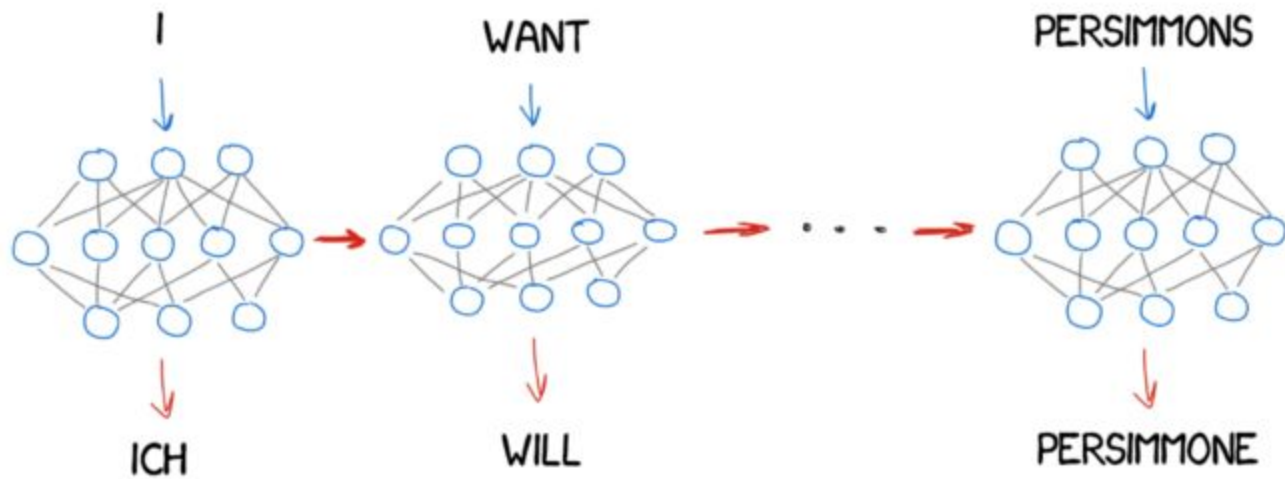


KILOGRAMM

PERSIMMONS



PERSIMONEN



ICH MÖCHTE KEINE PERSIMONEN ESSEN

/		/		
I	WANT	NOT	PERSIMMON	EAT
I	NOT	WANT	EAT	PERSIMMON

NOT ENOUGH EXAMPLES ABOUT PERSIMMONS

UNIGRAMS:

1. NOT
2. ENOUGH
3. EXAMPLES
4. ABOUT
5. PERSIMMONS

NOT ENOUGH EXAMPLES ABOUT PERSIMMONS

BIGRAMS:

1. NOT ENOUGH
2. ENOUGH EXAMPLES
3. EXAMPLES ABOUT
4. ABOUT PERSIMMONS

NOT ENOUGH EXAMPLES ABOUT PERSIMMONS

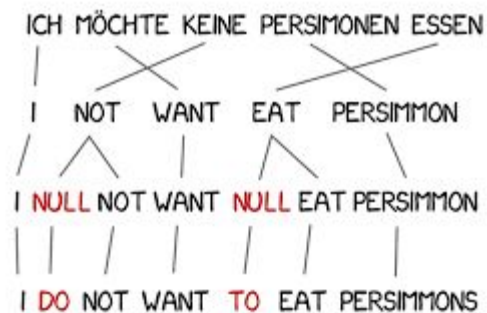
TRIGRAMS:

1. NOT ENOUGH EXAMPLES
2. ENOUGH EXAMPLES ABOUT
3. EXAMPLES ABOUT PERSIMMONS

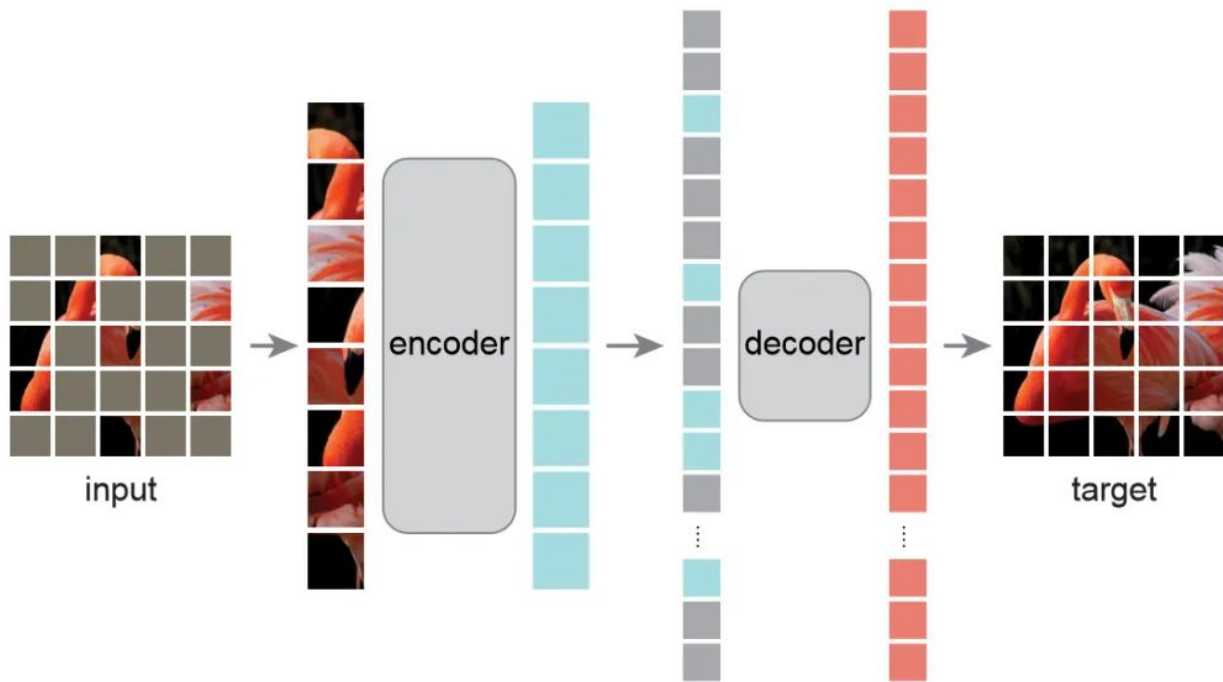
# Why translation ends up ugly

## Standard Challenges in Machine Translation

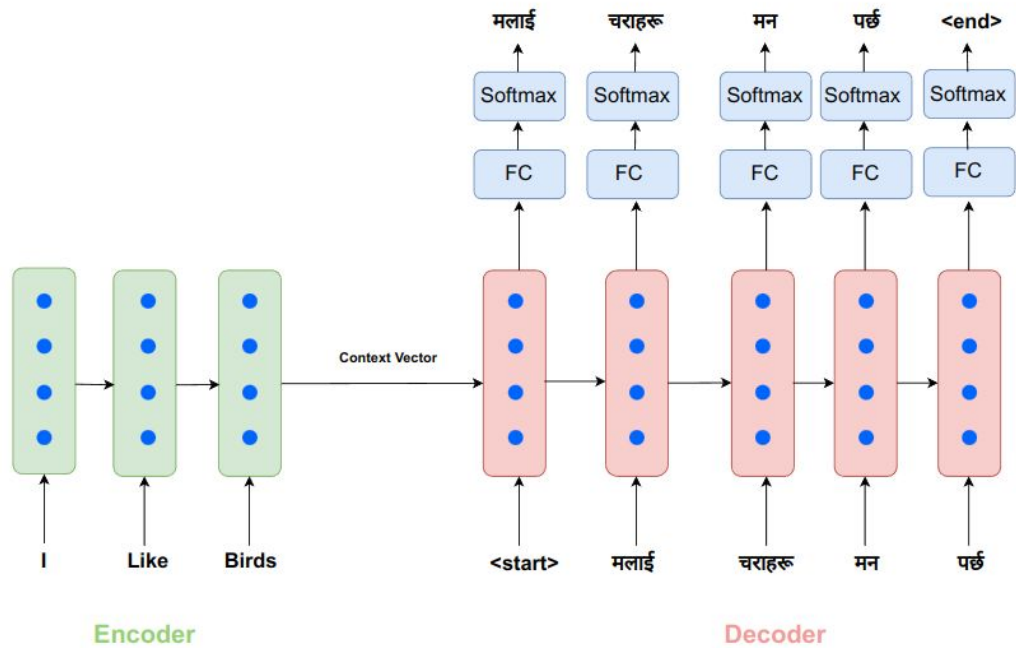
1. Correspondences are not 1:1.
2. The units of translation are not words, but concepts.
3. Words-to-concepts is not 1:1, because
  - a. Synonymy
  - b. Ambiguity
  - c. Rich morphology and word order



# Supervised Encoder-Decoder Machine Translation

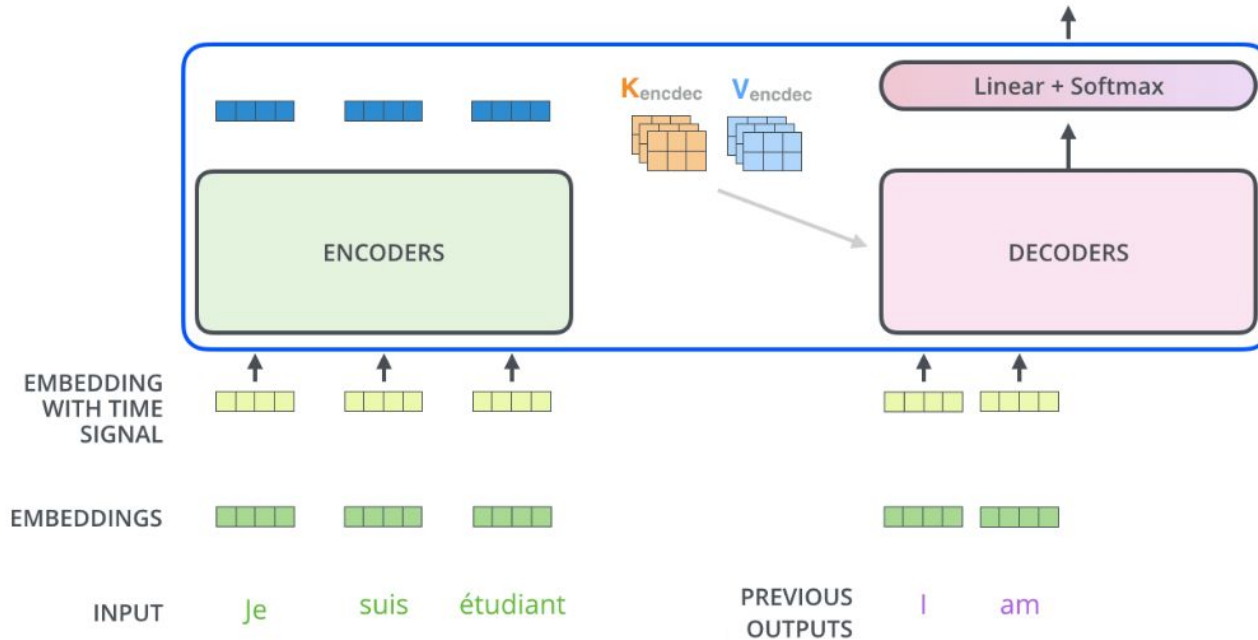






Decoding time step: 1 2 (3) 4 5 6

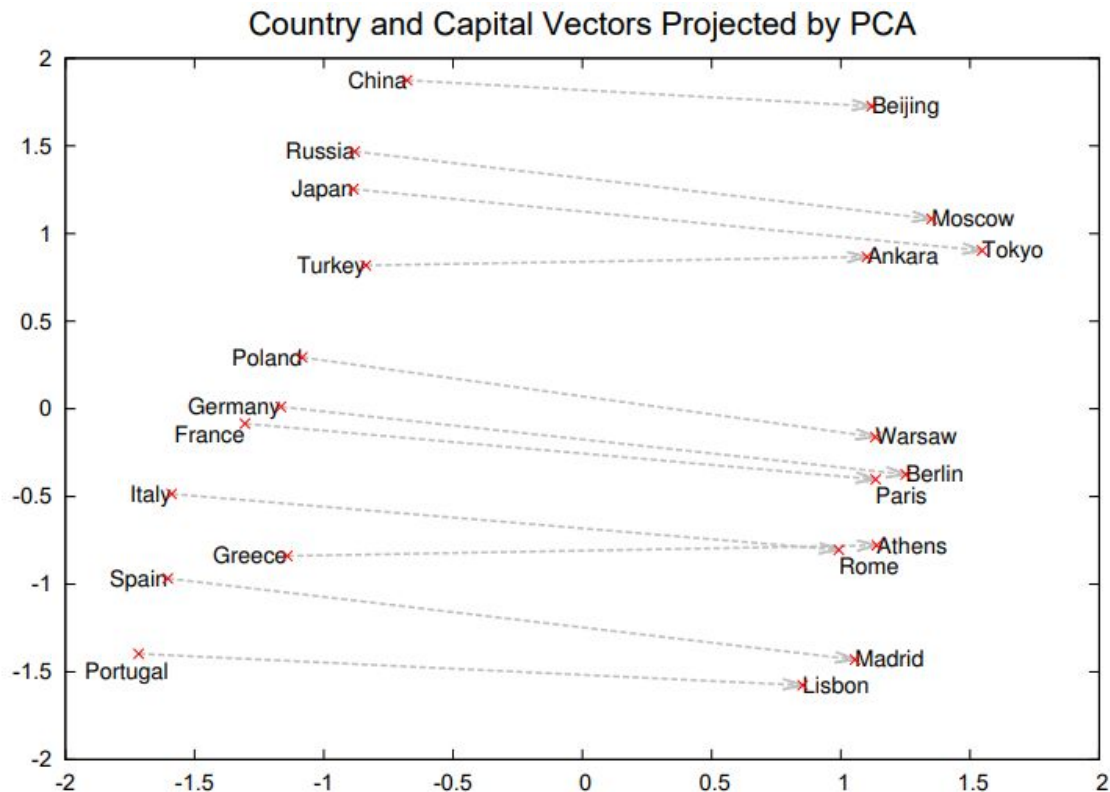
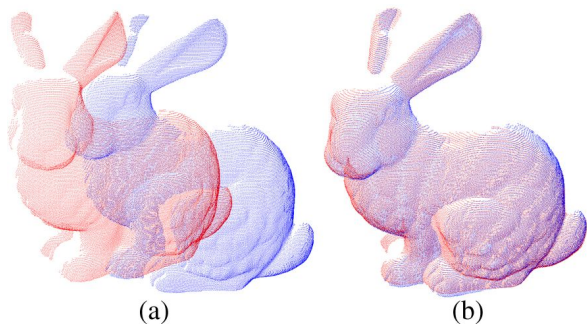
OUTPUT I am



Unsupervised  
Machine  
Translation?

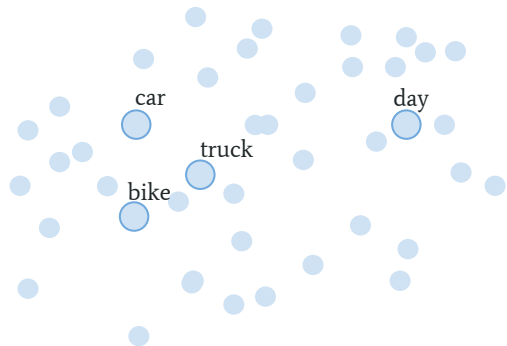
# Multilingual language models

*Observation:* Because relations are encoded systematically, in the limit language-specific embedding spaces will be isomorphic. This means we can learn linear mappings to align them.

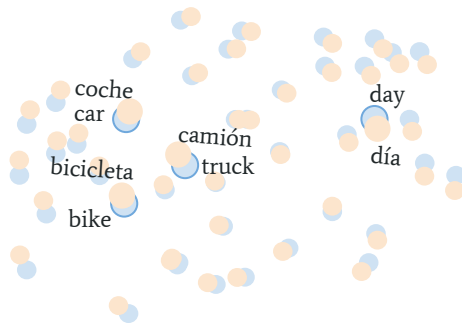
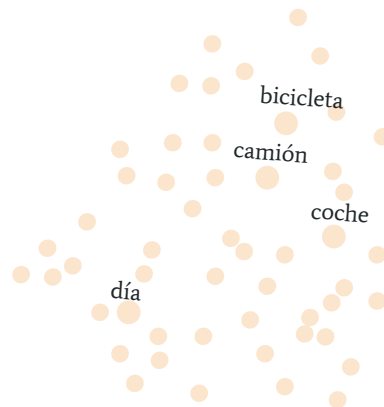


# Cross-Lingual Word-Embedding **Alignment**

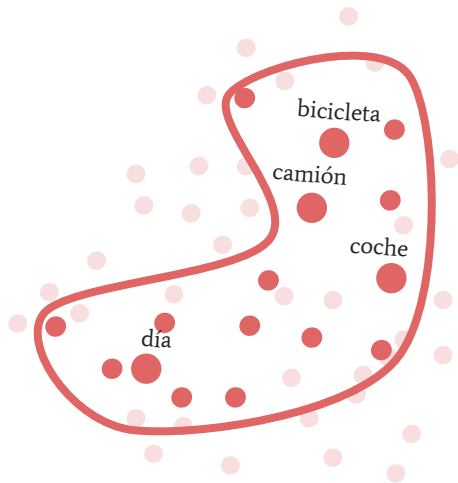
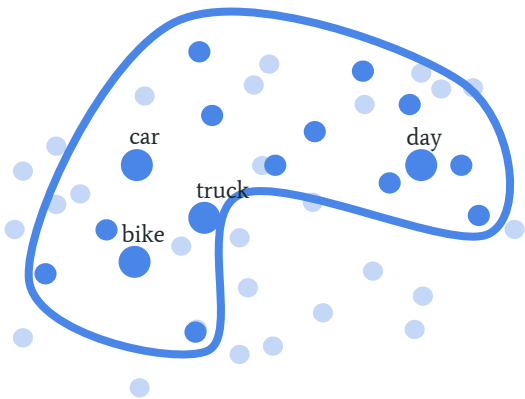
English monolingual word embeddings



Spanish monolingual word embeddings



# Alignment



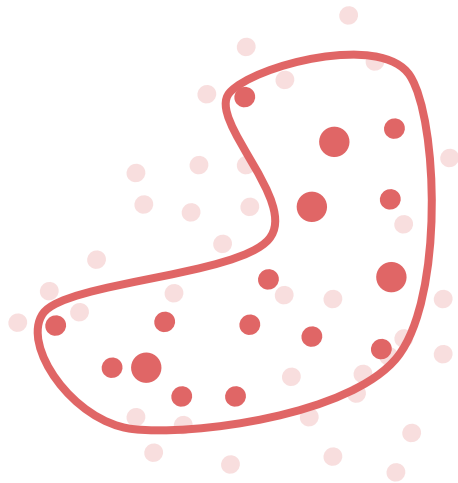
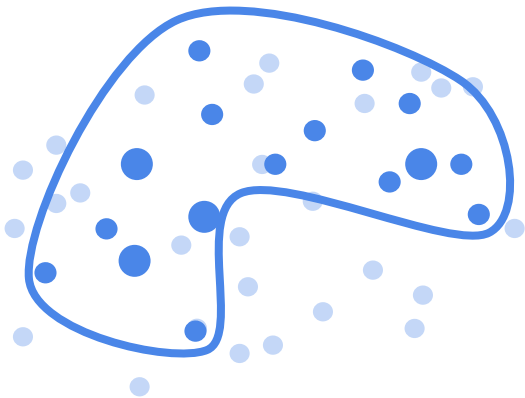
Train dictionary:

car	coche
truck	camión
bike	bicicleta
day	día

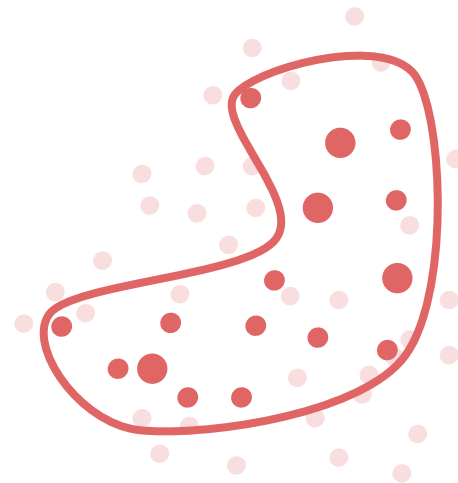
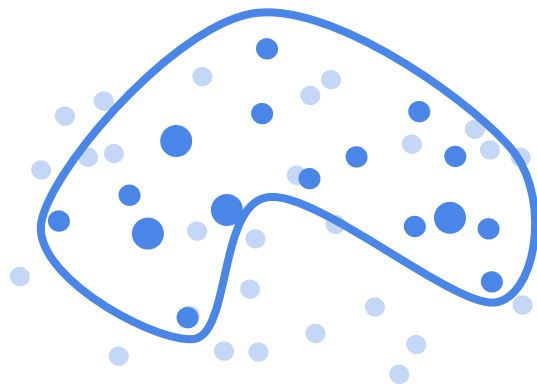
...

...

# Alignment

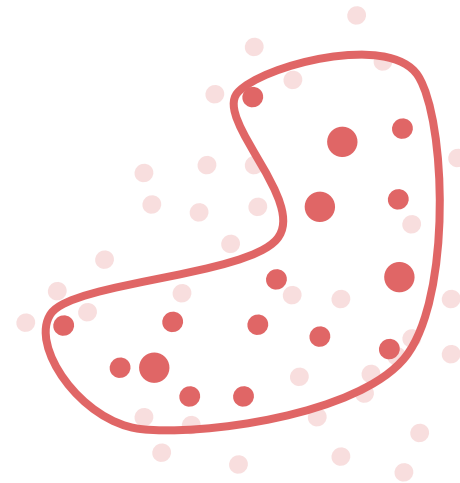
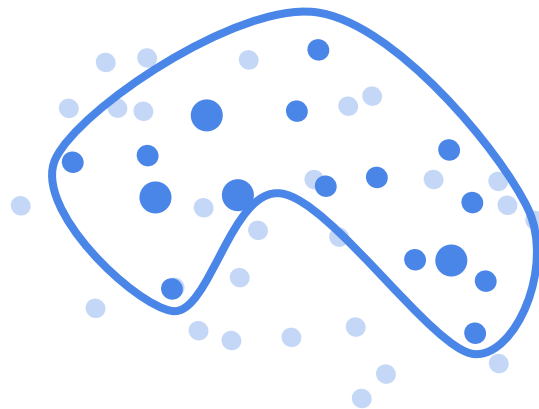


# Alignment

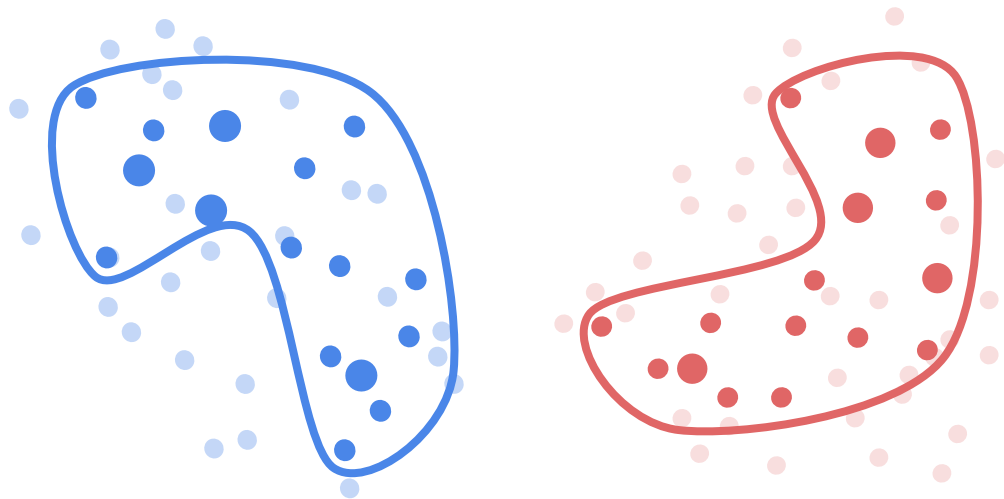




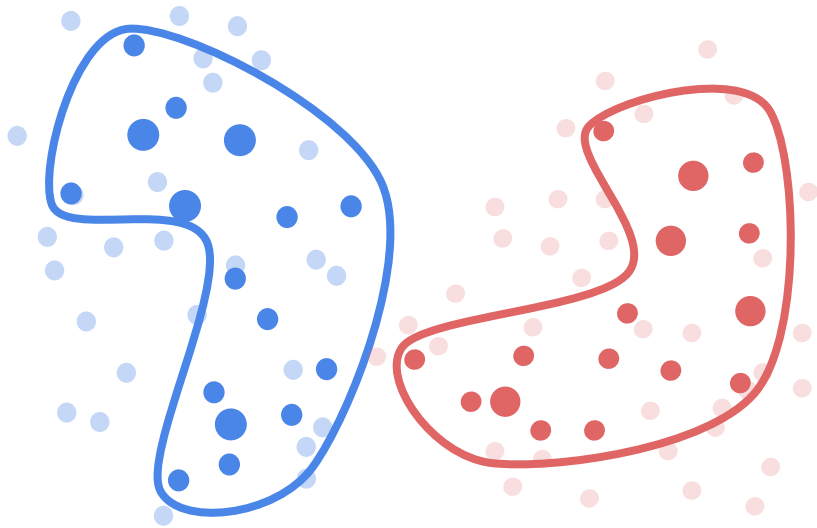
# Alignment



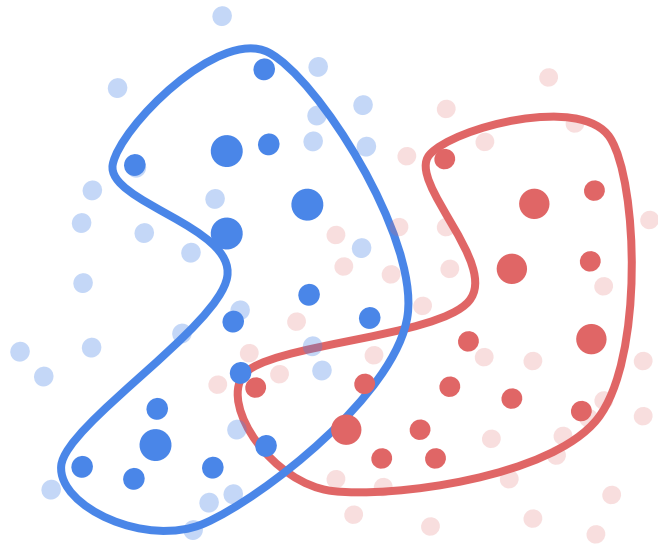
# Alignment



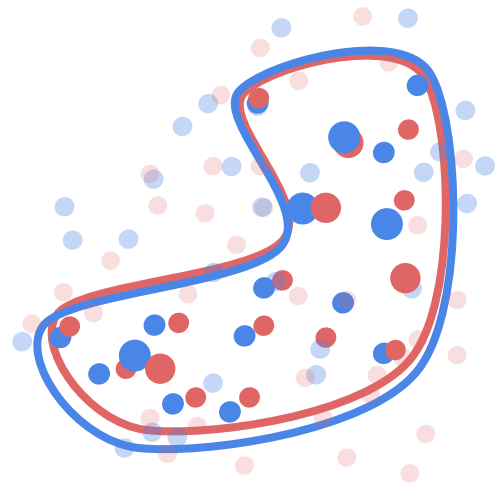
# Alignment



# Alignment



# Alignment



# Alignment

	X'	
car	-3.2	1.2
truck	-3.4	-1.1
bike	-2.9	0.5
day	2.3	1.3
...	...	...

	Z'	
coche	1.2	-2.3
camión	1.7	-1.1
bicicleta	1.4	3.1
día	4.2	2.1
...	...	...

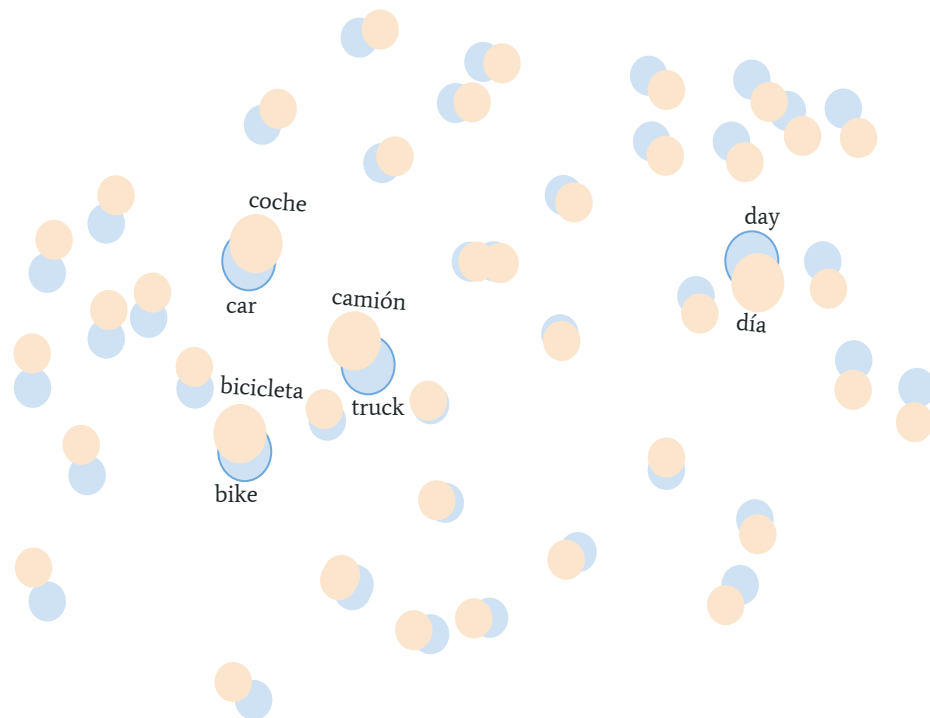
# Alignment

	X'				Z'	
car	-3.2	1.2			coche	1.2   -2.3
truck	-3.4	-1.1			camión	1.7   -1.1
bike	-2.9	0.5			bicicleta	1.4   3.1
day	2.3	1.3			día	4.2   2.1
...	...	...			...	...   ...

W		=
?	?	
?	?	

# Bilingual Dictionary Induction (BDI)



car coche  
bike bicicleta  
truck camión  
day día

table mesa  
toy juguete  
hour hora

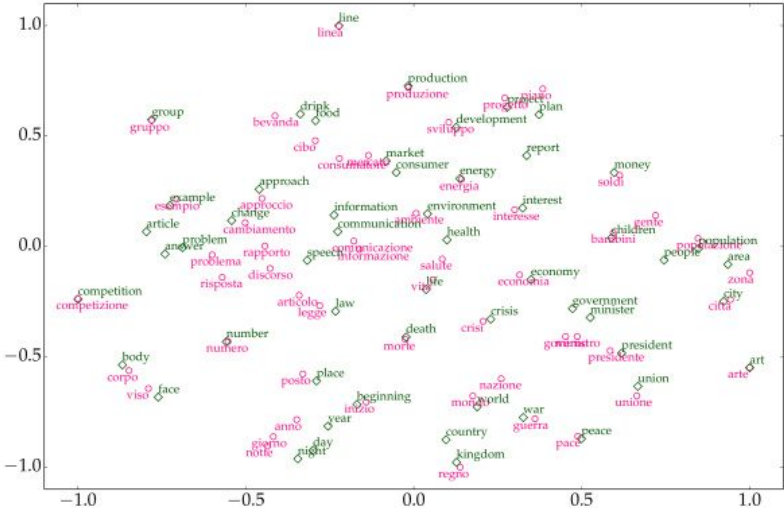
$$NN(\text{blue circle}_{table}) == \text{orange circle}_{mesa}$$



# Unsupervised machine translation

Unsupervised machine translation begins with vocabulary alignment, using point set registration algorithms. Once you know that 'line' and 'linea' are co-referential, we can begin to translate.

**Key idea:** If LM and CV models were aligned in the same way, we could translate and do VQA.



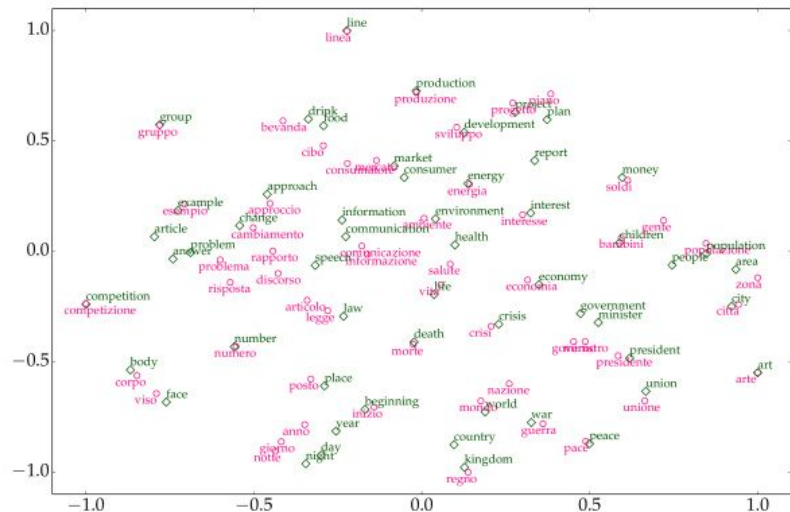
# Unsupervised machine translation

Unsupervised machine translation begins with vocabulary alignment, using point set registration algorithms. Once you know that ‘line’ and ‘linea’ are co-referential, we can begin to translate.

**Key idea:** If LM and CV models were aligned in the same way, we could translate and do VQA.

**Learned lesson:** Unsupervised alignment (e.g., using GANs) only work when spaces are very similar.

	Unsupervised (Adversarial)	Supervised (Identical)
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>
ET-FI	<b>29.62</b>	24.35



### Input image classes



ID: n02834778

### Input words & sentences

bike bicycle cycle wheel

Bike riders should follow the directional signs on ...

Bicycle theft is a crime involving theft of a bicycle.

Cell division occurs as part of a larger cell cycle.

It had a spoked steering wheel and bucket seats.

All had the required height adjustable steering wheel.

The throttle was controlled with a lever on the steer...

Vision  
Encoder

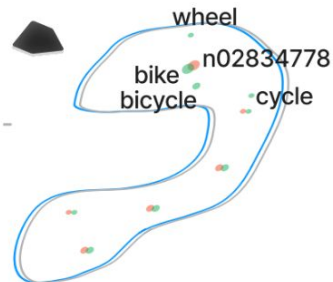
### Image embeddings



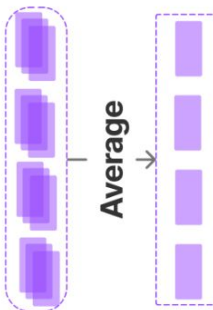
### Source Space



### Aligned Space

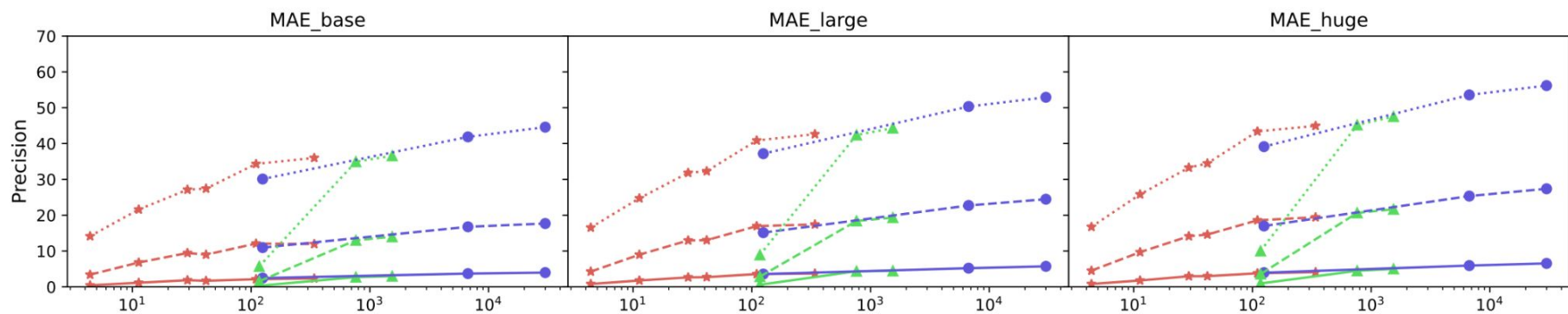


### Word embeddings



### Target Space



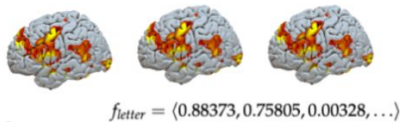


★ BERT models, ▲ GPT2 models, ● OPT models; Dotted line: P@100, dashed line: P@10, solid line: P@1.

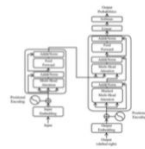
fMRI obtained while participants read or listen to language.



fMRI vectorized and aligned at word level, through Gaussian smoothing.

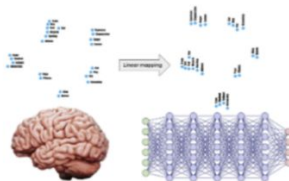


Decontextualized word embeddings obtained from LLMs.



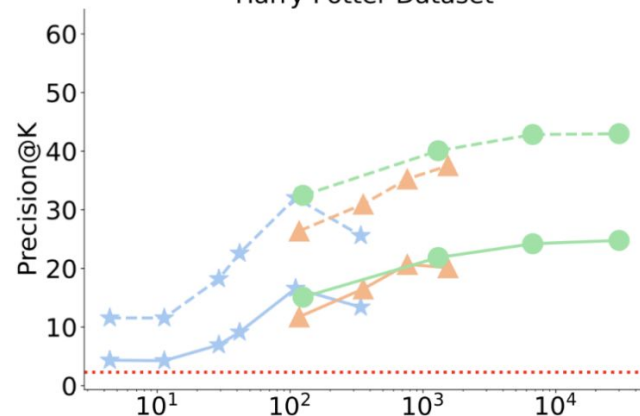
$v_{letter} = \langle 0.14723, 0.16827, 0.00328, \dots \rangle$

Alignment with Procrustes Analysis or ridge regression.

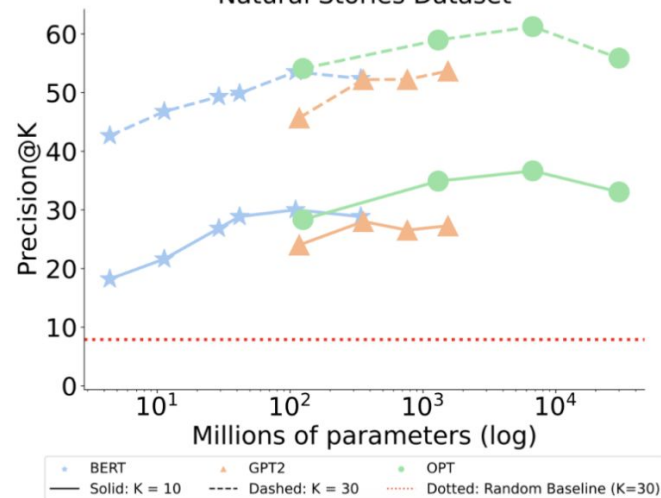


Results, retrieval precision at 50% (see plots → ).

Harry Potter Dataset



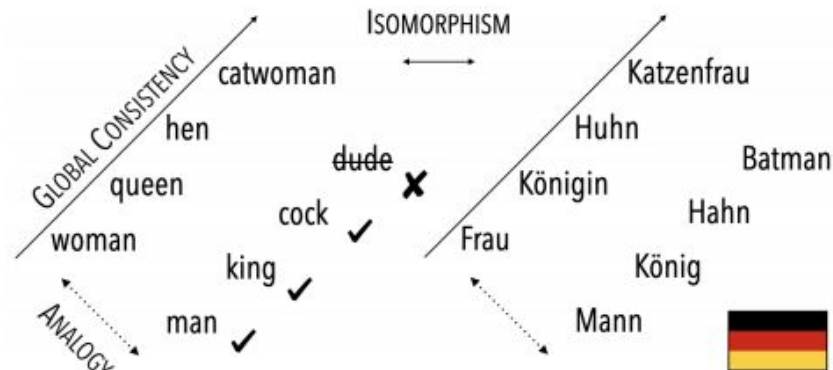
Natural Stories Dataset



## Analogy training multilingual language models

**Idea:** Use knowledge bases to ground language models, thereby debiasing them.

Collaboration with Deepmind and Cambridge University.



Language	Baselines				Analogy Training					
	Fasttext		mBERT		Fasttext		mBERT-WiQueen		mBERT-WiQueen <sup>+</sup>	
	P@1	$\rho$	P@1	$\rho$	P@1	$\rho$	P@1	$\rho$	P@1	$\rho$
Danish	0.1511	0.3001	0.2835	0.3221	0.1688	0.2909	0.3863	0.3010	<b>0.3935</b>	<b>0.2461</b>
German	0.0997	0.3604	0.2658	0.3548	0.1104	0.3702	0.3894	0.3257	<b>0.4538</b>	<b>0.2868</b>
English	0.1255	0.2854	0.2897	0.3107	0.1513	0.2550	0.4091	0.2960	<b>0.4787</b>	0.2821
Spanish	0.0899	0.3383	0.2596	0.3441	0.1194	0.3573	0.3832	0.3198	<b>0.3936</b>	<b>0.3012</b>
Finnish	0.1258	0.3908	0.2679	0.3535	0.1682	0.3731	0.3728	0.3192	<b>0.4019</b>	<b>0.2703</b>
French	0.0943	0.3659	0.2617	0.3545	0.1146	0.3459	0.3707	0.3375	<b>0.4195</b>	<b>0.2991</b>
Italian	0.0731	0.3979	0.2773	0.3711	0.0949	0.3883	0.3821	<b>0.3338</b>	<b>0.4372</b>	0.3722
Dutch	0.1291	0.3520	0.2669	0.3443	0.1497	0.3384	0.3811	<b>0.3202</b>	<b>0.4424</b>	0.3609
Polish	0.1165	0.3397	0.2648	0.3656	0.1456	0.3287	0.3853	0.3468	<b>0.3894</b>	<b>0.2718</b>
Portuguese	0.0898	0.3614	0.2523	0.3536	0.1072	0.3640	0.3697	0.3409	<b>0.3718</b>	<b>0.2653</b>
Swedish	0.1071	0.3449	0.2856	0.3378	0.1415	0.3270	0.3832	0.3128	<b>0.4071</b>	<b>0.2672</b>
<b>Averages</b>	0.1093	0.3488	0.2704	0.3435	0.1338	0.3399	0.3830	0.3231	<b>0.4171</b>	<b>0.2930</b>

?

coAStal

