Towards Language Technology for a Truly Multilingual World?

Ivan Vulić LTL, University of Cambridge & PolyAI



ECIR 2022, April 13 2022





Why Multilingual NLP?



...and reaching more users and customers...

"I'd like a ride to Russell Square" אני רוצה מונית לתחנה המרכזית בתל אביב "Posso fare un giro per sei persone a Roma Termini?" "Један ауто до главне железничке молим Вас" "یک کابین در ایستگاه اصلی اتوبوس لطفا" "¿Puedo tomar un taxi hasta el aeropuerto?" "Molim Vas jedno vozilo do Autobusnog" هل يمكننى الحصول على سيارة أجرة من ميدان التحرير؟ "可以載我去故宮博物館嗎?"

"私は銀座にタクシーを手に入れることはできますか?"

Speaking more languages means communicating with more people...

Why Multilingual NLP?

...but there are more profound and democratic reasons to work in this area:

- decreasing the digital divide
- dealing with inequality of information (access)
- mitigating cross-cultural biases
- deploying language technology for **underrepresented** languages, dialects, minorities; societal impact
- understanding cross-linguistic differences

"95% of all languages in use today will never gain traction online" (Andras Kornai)

"The limits of my language online mean the limits of my world?"

Source: http://labs.theguardian.com/digital-language-divide/



Why Multilingual NLP?

We're in Zagreb searching for...



...éttermek (HU)

...jatetxea (EU)

Inequality of information and representation can also affect how we understand places, events, processes...

...restaurants (EN)



English Conversational Al

A successful conversational agent must perform:

Automatic speech recognition (ASR)

• Language analysis:

- Language modeling, spelling correction
- Syntactic analysis: POS tagging, parsing
- Semantic analysis: named entity recognition, event detection, semantic role labeling, WSD
- Coreference resolution, entity linking, commonsense reasoning, world knowledge

Dialog modeling:

- Natural language understanding, intent detection, language generation, dialog state tracking
- Information Search and QA
- Text-to-Speech



Multilingual Conversational AI?

According to Ethnologue there are 7,000+ living languages

What about language varieties and dialects?

What about "social media" languages and slang?

IN AN IDEAL WORLD WE WOULD HAVE ALL THE DATA WE NEED!

THIS WORLD IS THE ONLY ONE WE'VE GOT

imatip.com

The Long Tail of Data



Even getting "raw" unannotated data is problematic for many languages...



Are All Languages Created Equal?



Class	5 Example Languages	#Langs	#Speakers	% of Total Langs
0	Dahalo, Warlpiri, Popoloca, Wallisian, Bora	2191	1.2B	88.38%
1	Cherokee, Fijian, Greenlandic, Bhojpuri, Navajo	222	30M	5.49%
2	Zulu, Konkani, Lao, Maltese, Irish	19	5.7M	0.36%
3	Indonesian, Ukranian, Cebuano, Afrikaans, Hebrew	28	1.8B	4.42%
4	Russian, Hungarian, Vietnamese, Dutch, Korean	18	2.2B	1.07%
5	English, Spanish, German, Japanese, French	7	2.5B	0.28%

Most languages are "Left-Behinds" [Joshi et al., ACL-20; Blasi et al., ACL-22]

Is creating equitable language technology across different languages then even possible?

Can we at least try to 'approximate' equality?





Language Variety and Variability

Afro-Asiatic Nilo-Sa			aharan?	Nig	jer–Co	ngo	Khois	Khoisan (areal)		
Indo- European	Cauca (areal	asian I)	Uralic		Dra	avidian	Altai (are	ic al)	Paleosit (areal)	
Sino-Tibet	an		Hmong-	Mie	en	Kra-D)ai		Austroas	
Austronesian			Papuan / (areal)			Austra	alian (areal)	Andama (areal)	
Eskimo– Aleut	Algic	U	Ito-Azteca	an	Na Yei	-Dené niseiar	(and ?)	Dené-	- America (areal)	
Creole/Pic	lgin/M	ixed La	anguage solate	Sig Ian	n gua	Co age lar	onstru nguag	cted L	Inclassifie	

[Image from: Wikipedia]





Figure 2: Density of WALS typological features of the world's languages reduced to 2 dimensions via PCA. Red dots are languages covered by UD. Darkness corresponds to more probable regions.

Image: courtesy of Edoardo Maria Ponti



Multilingual Representation Learning Cross-lingual Transfer Learning

Why Cross-Lingual Transfer?



Image: courtesy of Sebastian Ruder

Why Cross-Lingual Transfer and Multilingual Learning?

- Many NLP tasks share common knowledge about language (e.g. linguistic representations, structural similarities)
- Languages share common structure (on the lexical, syntactic, and semantic level)
- Annotated data is rare, make use of as much supervision as available
- Empirically, transfer learning has resulted in SOTA for many supervised NLP tasks (e.g. classification, information extraction, QA, etc)

Image courtesy of Yulia Tsvetkov





Joint Multilingual Learning in a Nutshell



Joint multilingual learning – train a single model on a mix of datasets in all languages, to enable data and parameter sharing where possible

Cross-Lingual Transfer in a Nutshell



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

generalizes out-of-the-box in a low-resource language/domain

Few-shot learning: train a model in one language/domain and use only few examples from a low-resource language/domain to adapt it

Zero-shot learning: train a model in one language/domain and assume it

Cross-Lingual Representations Enable Transfer



Crossing the Chasm: Representation Learning

Multilingual / cross-lingual representation of meaning

Word-Level

Cross-lingual word embeddings

 Words with similar meaning across
 languages have similar representations

Text Encoding

Multilingual unsupervised pretraining

- o mbert
- \circ XLM(-R)
- \circ mT5
- 0 ...





Old School: Cross-Lingual Word Embeddings

A range of different methods (with different data requirements), but the same goal: Induce a semantic vector space in which words with similar meaning end up with similar vectors, regardless of whether they come from the same language or from different languages.



Inspired by: Slides from Eneko Agirre, Mikel Artetxe



Old School*: Cross-Lingual Word Embeddings

*Old (NLP/ML/IR dialectal): dominantly used in 2019 and even in 2020...

How to use CLWEs for cross-lingual transfer for supervised tasks?

(Assumption: zero-shot transfer)

Step 1. Induce the cross-lingual (bilingual) word embedding space

Step 2. Train the (neural) model using the task-annotated data in the source language • e.g., for NER train a BiLSTM+classifier using embeddings of the source language words as the input

Step 3. At prediction time, for texts in the target language, use embeddings of target language words as the input to the trained classifer



New School: Multilingual Language Models

Deep Transformer networks pretrained on large multilingual corpora via (masked) language modeling objectives

They work with an automatically induced shared subword vocabulary across all languages they represent

They are unsupervised from the perspective of not using any explicit cross-lingual learning signal.

At first, praised for their effective (zero-shot) cross-lingual performance • "Surprising cross-lingual effectiveness of mBERT"

• "mBERT surprisingly good at zero-shot cross-lingual transfer"



New School: Zero-Shot Transfer to (Low-Resource) Languages

Step 1: Train a multilingual model. **Step 2:** Step 3: Why?

Training **data** is **expensive** and not available for many languages, especially ones that are considered "low-resource".

Fine-tune model on a task in a high resource source language.

Transfer and evaluate the model on a low resource target language.

So... We Have Solved Zero-Shot Cross-Lingual Transfer?

No! Settings in which they were evaluated were too simple and too favourable...

Study 1. **Tasks:** NER, POS tagging (Pires et al., ACL-19)

Study 2. **Tasks:** NER, NLI (Karthikeyan et al., ICLR-20)

In most studies the selected target languages were:

from the same language family as the source (English)
 with large corpora in pretraining

Target Languages: DE, NL, ES

Target Languages: ES, HI, RU

So... We Have Solved Zero-Shot Cross-Lingual Transfer?

Task	Model	EN	<mark>ZH</mark> Д	TR Δ	RU Δ	AR	HI Δ	EU Δ	FI Δ	HE Δ	IT A	$\frac{JA}{\Delta}$	ко Д	sv A	VI Δ	$\frac{TH}{\Delta}$	ES Δ	EL A	DE A	FR	BG ∆	sw A	UR A
DEP	B X	91.2 92.0	-43.9 - 85.4	-46.0 -44.2	-28.1 -29.7	-56.4 -54.6	-36.1 -39	-50.2 -49.5	-30.7 -26.7	-36.1 -39	-17.1 -23.5	-60.1 -80.5	-56.1 -56.0	-14.3 -16.3	÷	-	÷	-	1	-	-	-	-
POS	B X	95.8 96.3	-38.0 -69.2	-35.9 -27.7	-16.0 -14.3	-40.1 -37.1	-33.4 -27.3	-34.6 -31.9	-21.9 -17.9	-33.4 -27.3	-19.8 -19.0	-46.1 -77.0	-42.0 -37.3	-9.6 -10.7	-	-	-	-	-	-	÷	-	-
NER	B X	92.4 91.6	-23.3 - 34.8	-11.6 -6.2	-10.7 -13.7	-31.7 -24.6	-11.1 -16.5	-12.8 -8.0	-3.8 -0.9	-11.1 -16.5	-2.6 -2.4	-25.7 -30.1	-13.8 -15.6	-6.7 -2.2	-	-	-	2	-	-	-	-	2
XNLI	B X	82.8 84.3	-13.6 -11.0	-20.6 -11.3	-13.5 -9.0	-17.3 -13.0	-21.3 -14.2	Ī	-	-	-	÷	-	-	-11.9 -9.7	-28.1 -12.3	-8.1 -5.8	-14.1 -8.9	-10.5 -7.8	-7.8 -6.1	-13.3 -6.6	-33.0 -20.2	-23.4 -17.3
XQuAD	B X	71.1 72.5	-22.9 -26.2	-34.2 -18.7	-19.2 -15.4	-24.7 -24.1	-28.6 -22.8	Ę	-	-	-	-	-	-	-22.1 -19.7	-43.2 -14.8	-16.6 -14.5	-28.2 -15.7	-14.8 -16.2	-	-	-	-

(Lauscher et al., EMNLP-20)

B=mBERT; X=XLM-R

Huge drops for: 1. Distant target languages 2. Target languages with small pretraining corpora



So... We Have Solved Zero-Shot Cross-Lingual Transfer?

More problems...

"The Curse of Multilinguality"

(Conneau et al., ACL-20)



And what about low-resource languages not covered at all in the pretraining data?



Some Recap and Basics...

(Almost) everything I am about to cover in more detail involves:

- Transfer learning in NLP
- We **only** look at deep neural networks, specifically the Transformer architecture.
- We leverage pre-trained transformer-based models such as BERT/ROBERTa/XLM-R/mBERT.
- These have been trained on **massive** amount of text data using Masked Language Modeling (MLM).
- Transfer learning with these pretrained models usually involves stacking a prediction head on top of the model.
- Usually all parameters are the fine-tuned on the downstream task (e.g. using cross-entropy loss).



Problems of Multi-Task and Transfer Learning

Multi-Task Learning:



Sequential Fine-Tuning:



Catastrophic Interference: Sharing all parameters **\Theta** between tasks results in deterioration of performance for a subset of tasks.

Catastrophic Forgetting: Sequential fine-tuning on tasks results in forgetting information learned in earlier stages of transfer learning.



Modular and Parameter-Efficient? $\Theta - \operatorname{argmin} L(D_{NL}; \Theta)$

A single Transformer (encoder) layer





= Parameters are fine-tuned

D_{NLI} = NLI Dataset

L = Loss function, e.g. cross entropy loss

 Θ = Parameters of the model



Modular and Parameter-Efficient: Adapters $\Theta - \operatorname{argmin} L(D_{NL}; \Theta)$



= Parameters are frozen

= Parameters are fine-tuned

Houlsby, Neil, et al. "Parameter-Efficient Transfer Learning for NLP." International Conference on Machine Learning. 2019.



Modular and Parameter-Efficient: Adapters $\Theta - \operatorname{argmin} L(D_{NI}; \Theta)\phi)$ A single Transformer

Φ

A



= Parameters are frozen

Houlsby, Neil, et al. "Parameter-Efficient Transfer = Parameters are fine-tuned Learning for NLP." International Conference on Machine Learning. 2019.

<u>Adapter</u> parameters \$ are encapsulated between transformer layers with parameters **\Theta** which are frozen







Parameter Efficiency of Adapters in Transformers

Training adapters instead of full model fine-tuning achieves similar results.

Adapters are smaller in size than training the full model.

R' Μ S7 C SS Q Μ 0

Houlsby, Neil, et al. "Parameter-Efficient Transfer Learning for NLP." International Conference on Machine Learning. 2019. Pfeiffer, Jonas, et al. "AdapterFusion: Non-destructive task composition for transfer learning." arXiv preprint (2020a). Pfeiffer, Jonas, et al. "Adapterhub: A framework for adapting transformers." Proceedings of EMNLP: Systems Demonstrations (2020b)

Performance on GLUE tasks

	Full	Pfeif.	Houl.
TE (Wang et al., 2018)	66.2	70.8	69.8
RPC (Dolan and Brockett, 2005)	90.5	89.7	91.5
FS-B (Cer et al., 2017)	88.8	89.0	89.2
oLA (Warstadt et al., 2019)	59.5	58.9	59.1
ST-2 (Socher et al., 2013)	92.6	92.2	92.8
NLI (Rajpurkar et al., 2016)	91.3	91.3	91.2
NLI (Williams et al., 2018)	84.1	84.1	84.1
QP (Iyer et al., 2017)	91.4	90.5	90.8



Houlsby et al., 2019

Number of newly introduced Parameters

Rate	Base #Params	Size	Large #Params	Size
64	0.2M	0.9Mb	0.8M	3.2Mb
16	0.9M	3.5Mb	3.1M	13Mb
2	7.1M	28Mb	25.2M	97Mb



Pfeiffer et al., 2020a











Encapsulated Adapters?

MLM (English) MLM (Quechuan)



- Adapters learn transformations that make the underlying model more suited to a task or language.
- Using masked language modelling (MLM), we can learn language-specific transformations for e.g.
 English and Quechua.
- As long as the underlying model is kept fixed, these transformations are **roughly interchangeable**.





MAD-X: An Adapter-Based Framework for Transfer

Step 1: Train Language Adapters

We train **language adapters** for the **source language** and the **target language** with masked language modelling on Wikipedia.



MLM (English)



MLM (Quechuan)



MAD-X

Step 2: Train a Task Adapter

We train task adapters in the source language stacked on top of the source language adapter.

The language adapter $\phi_{\rm l}$ as well as the transformer weights **O** are frozen while only the task adapter parameter ϕ_{+} are trained.



MAD-X

Step 3: Zero-Shot transfer to unseen language

We **replace** the **source** language adapter with the **target** language adapter, while **keeping** the "language agnostic" **task adapter**.





Datasets: Inclusion of Diverse and Low-Resource

NER: WikiAnn Dataset We chose a diverse set of languages from different language families.

XQuAD (Cross-lingual Question Answering Dataset)

XCOPA (Ponti et al. 2020b)

1	Language	ISO code	Language family	# of Wiki articles	Covered by SOTA
	English	en	Indo-European	6.0M	\checkmark
	Japanese	ja	Japonic	1.2M	\checkmark
	Chinese	zh	Sino-Tibetan	1.1 M	\checkmark
	Arabic	ar	Afro-Asiatic	1.0M	\checkmark
	Javanese	jv	Austronesian	$\overline{57k}$	
	Swahili	SW	Niger-Congo	56k	\checkmark
	Icelandic	is	Indo-European	49k	\checkmark
	Burmese	my	Sino-Tibetan	45k	\checkmark
	Quechua	qu	Quechua	$\overline{22k}$	
	Min Dong	cdo	Sino-Tibetan	15k	
	Ilokano	ilo	Austronesian	14k	
	Mingrelian	xmf	Kartvelian	13k	
	Meadow Mari	mhr	- Ūralic	<u>10k</u>	
	Maori	mi	Austronesian	7k	
	Turkmen	tk	Turkic	6k	
	Guarani	gn	Tupian	4k	
		25			



Relative **F1 improvement** of **MAD-X^{Large} over XLM-R^{Large}** in cross-lingual NER transfer.

Top right corner represent the realistic scenario of transfering from high resource to low resource

en ja zh ar jv -Language sw is my qu -Source cdo ilo xmf mi mhr tk gn

Languages are more low-resource or unseen during pre-training

0.2	-0.7	0.3	-11.8	5.8	7.7	4.7	5.9	19.7	26.0	8.1	15.9	10.2	12.8	15.5
-0.3	0.8	4.0	-4.4	2.0	4.2	0.5	6.3	14.6	37.4	-3.6	16.3	23.7	2.9	2.2
4.7	3.4	-0.1	-0.1	2.6	4.7	6.7	1.3	21.2	36.6	3.2	13.0	26.8	15.8	-0.7
7.4	-1.7	-1.6	0.7	9.8	12.1	9.8	-4.3	24.5	44.8	26.9	19.2	20.7	21.7	20.1
-8.4	-3.5	-5.0	-5.7	2.9	-2.2	0.5	-1.0	3.7	18.3	4.1	-3.1	8.2	6.5	7.5
-1.4	-4.4	-8.3	-2.9	5.1	2.2	3.8	1.4	17.1	28.6	16.7	11.2	9.3	13.8	14.4
-3.2	-4.4	-7.7	-10.8	9.8	-7.4	1.6	0.6	7.3	27.8	3.4	7.7	12.6	14.2	10.4
-7.5	-1.7	-3.1	-9.4	5.3	-3.7	-2.8	-10.6	-3.3	15.9	-3.9	-0.5	-4.5	2.2	2.6
-4.5	-0.2	0.2	-5.8	4.7	-0.1	12.1	0.4	9.0	25.3	-0.3	18.4	22.1	9.1	16.4
13.7	-0.0	-1.3	7.2	15.4	5.9	22.3	15.7	5.9	16.3	-0.5	22.7	18.0	9.4	16.6
6.0	-2.9	-5.0	5.2	14.9	12.5	14.6	-0.3	11.6	16.6	12.0	12.3	7.9	14.7	6.1
-4.5	-0.2	1.7	-9.2	-0.3	-5.9	-3.0	-5.7	3.0	24.3	-5.8	0.1	-6.0	5.8	4.9
-1.0	-0.1	-0.4	-1.4	4.0	-1.8	8.5	2.3	5.0	10.9	-3.5	19.8	5.2	9.0	8.5
-5.6	-2.5	-4.8	-7.0	-2.0	0.7	2.6	-0.3	3.8	10.0	-10.3	-0.4	9.0	7.6	5.1
-2.7	-1.7	-3.8	-5.5	6.3	-9.5	3.0	-3.5	2.7	14.4	-11.2	6.7	3.2	11.2	12.5
16.1	-1.6	-3.6	-14.5	-13.8	-11.6	-9.6	-6.2	-9.4	8.7	-14.6	-9.9	-5.4	-4.2	-8.5
en	ja	zh	ar	jv	sw	is	my	qu	cdo	ilo	xmf	mi	mhr	tk
					Т	arg	et L	ang	uage	e				



Towards Efficient (and Typologically Driven) Cross-Lingual Transfer

Generating Adapter Parameters

generate adapters on the fly, conditioned on language vectors?

This can be seen as factorising adapters to language (and layer) parameters

More efficient than keeping dedicated language adapters

- The main idea: instead of having dedicated single-language adapters, can we learn to
- It can even work (in theory) in zero-shot (no text data whatsoever!) and few-shot setups...



Generating Adapter Parameters: MAD-G



CPG-Adapters offer better initialisation for further target-specific MLM-ing in low-data scenarios...

(Ansell et al., 2021) learn to generate a monolithic multilingual CPG*-adapter by MLM-ing on 95 languages. A simpler and more efficient MAD-X-style transfer learning.

*CPG=Contextual Parameter Generation (Platanios et al., 2018)

Multi-Source Transfer works much better across different tasks: the model is forced to learn more general language-invariant transfer features?







MAD-G is efficient...

- Full fine-tuning of mBERT for 95 languages requires: 95 * 178M = 16.91B parameters (!!)
- Having MAD-X for 95 languages requires: 728M additional parameters (95 single-language adapters)
- MAD-G for 95 languages conditioned on language vectors only: **228M** parameters
- positions: **38M** additional parameters
- Average per-language training time: x50 shorter than for MAD-X (II)

MAD-G for 95 languages conditioned on language vectors and layer

...but also effective...

method	hau MAD-G-seen	ibo MAD-G-seen	kin MAD-G-seen	lug unseen	luo unseen	pcm unseen	swa mBERT-seen	wol unseen	yor mBERT-seen	8
MAD-G	77.1	69.9	66.1	<u>54.2</u>	32.5	72.6	72.6	32.1	68.8	
MAD-G-LS	<u>72.8</u>	<u>67.5</u>	<u>63.0</u>	55.7	33.3	<u>72.4</u>	71.3	36.7	68.4	
MAD-G-en	44.9	54.5	51.4	50.6	<u>32.9</u>	70.4	69.2	<u>36.4</u>	63.9	
TA-only	43.4	55.7	52.8	47.9	32.8	72.3	68.6	32.1	65.3	
mBERT-ft	43.2	45.5	49.9	49.3	31.6	70.5	65.8	28.1	54.3	
XLM-R-ft	66.4	45.5	36.1	34.8	31.9	68.4	74.5	21.6	33.4	

Table 3: F_1 scores on the MasakhaNER dataset for African languages. Task adapter training/model fine-tuning is conducted on the CoNLL 2003 English NER dataset. XLM-R-ft results are as reported by Adelani et al. (2021).



We can initialise dedicated language adapters via MAD-G for improved transfer





Another Parameter-Efficient Paradigm

Sparse composable masks for cross-lingual transfer



(Ansell et al., ACL-2022)

Performance gains over MAD-X and MAD-G

Parameter-Efficient Cross-Lingual (Re)ranking



(WIP)



Notes on Zero-Shot versus Few-Shot Learning

Should we focus more on few-shot transfer scenarios and quick annotation cycles?



Massively Multilingual Transformer





(Lauscher et al., ACL-20; Zhao et al., ACL-21)

Fine-tuning on target task in English

POS

Few-shot fine-tuning on target task in target language

Prediction on target task in target language



Few-Shot > (or >>) Zero-Shot

		K=0	K=1	K=2	K=4	K=8
	EN	96.88	24		÷	23
	DE	88.30	90.36 ± 1.48	90.77 ± 0.87	91.85 ± 0.83	91.98 ± 0.82
S	FR	83.05	88.94 ± 2.46	89.71 ± 1.68	90.80 ± 0.88	91.01 ± 0.94
õ	ES	81.90	83.99 ± 2.35	85.65 ± 1.60	86.30 ± 1.85	88.46 ± 1.90
Π	IT	74.13	74.97 ± 2.04	75.29 ± 1.57	76.43 ± 1.41	78.12 ± 1.25
N	RU	72.33	77.40 ± 4.27	80.57 ± 1.37	81.33 ± 1.33	81.91 ± 1.21
	ZH	84.38	87.18 ± 1.45	87.31 ± 1.53	88.33 ± 1.11	88.72 ± 1.05
	JA	74.58	76.23 ± 1.59	76.71 ± 2.12	78.60 ± 2.43	81.17 ± 1.72
	EN	64.52		-	-	6 - 3
7)	DE	49.62	51.50 ± 1.58	52.76 ± 0.87	52.78 ± 1.00	53.32 ± 0.59
R	FR	47.30	49.32 ± 1.34	49.70 ± 1.43	50.64 ± 0.94	51.23 ± 0.76
IA	ES	48.44	49.72 ± 1.24	49.96 ± 1.12	50.45 ± 1.22	51.25 ± 0.93
2	ZH	40.40	43.19 ± 1.76	44.45 ± 1.36	45.40 ± 1.26	46.40 ± 0.93
	JA	38.84	41.95 ± 2.09	43.63 ± 1.30	43.98 ± 0.89	44.44 ± 0.69
<u></u>	EN	82.67	-	-		
	DE	70.32	70.58 ± 0.36	70.60 ± 0.34	70.61 ± 0.39	70.70 ± 0.50
	FR	73.57	73.41 ± 0.48	73.74 ± 0.46	73.57 ± 0.49	73.77 ± 0.44
	ES	73.71	73.84 ± 0.40	73.87 ± 0.44	73.74 ± 0.48	73.87 ± 0.46
	RU	68.70	68.81 ± 0.52	68.76 ± 0.54	68.87 ± 0.55	68.81 ± 0.77
	ZH	69.32	69.73 ± 0.94	69.75 ± 0.94	70.56 ± 0.76	70.62 ± 0.86
Γ	AR	64.97	64.75 ± 0.36	64.82 ± 0.23	64.82 ± 0.23	64.94 ± 0.37
Z	BG	67.58	68.15 ± 0.69	68.19 ± 0.75	68.55 ± 0.67	68.32 ± 0.70
X	EL	65.67	65.64 ± 0.40	65.73 ± 0.36	65.80 ± 0.41	66.00 ± 0.53
	HI	56.57	56.94 ± 0.82	57.07 ± 0.82	57.21 ± 1.14	57.82 ± 1.18
	SW	48.08	50.33 ± 1.08	50.28 ± 1.24	51.08 ± 0.62	51.01 ± 0.79
	TH	46.17	49.43 ± 2.60	50.08 ± 2.42	51.32 ± 2.07	52.16 ± 2.43
	TR	60.40	61.02 ± 0.68	61.20 ± 0.61	61.35 ± 0.49	61.31 ± 0.56
	UR	57.05	57.56 ± 0.85	57.83 ± 0.91	58.20 ± 0.93	58.67 ± 1.03
	VI	69.82	70.04 ± 0.59	70.14 ± 0.75	70.23 ± 0.63	70.41 ± 0.70
	EN	93.90	-		-	-
	DE	83.80	84.14 ± 0.40	84.08 ± 0.42	84.04 ± 0.47	84.23 ± 0.66
SX	FR	86.90	87.07 ± 0.27	87.06 ± 0.37	87.03 ± 0.31	86.94 ± 0.41
M	ES	88.25	87.90 ± 0.54	87.80 ± 0.56	87.84 ± 0.53	87.85 ± 0.75
A	ZH	77.75	77.71 ± 0.37	77.63 ± 0.47	77.68 ± 0.51	77.82 ± 0.64
	JA	73.30	73.78 ± 0.75	73.71 ± 1.04	73.48 ± 0.69	73.79 ± 1.28
	ко	72.05	73.75 ± 1.30	73.11 ± 1.05	73.79 ± 0.92	73.31 ± 0.61

(...not only for token-level tasks)

Source Fine-Tuning Helps

	ML	Doc	PAW	VSX	l i	PC	S	NER		
	K=1	K=8	K=1	K=8		K=1	K=4	K=1	K=4	
DE	-37.73	-7.67	-31.11	-30.82	RU	-15.89	-3.20	-48.19	-35.77	
FR	-38.14	-13.21	-33.02	-32.34	ES	-9.51	-0.93	-63.98	-41.53	
ES	-33.69	-14.38	-33.76	-33.97	VI	-7.82	-0.36	-54.41	-41.45	
IT	-33.63	-12.62	-	-	TR	-15.05	-8.08	-54.35	-34.52	
RU	-30.66	-11.08	-	2	TA	-13.72	-4.40	-34.70	-24.81	
ZH	-37.31	-12.57	-23.74	-23.65	MR	-11.34	-3.63	-40.10	-25.68	
JA	-29.82	-14.32	-20.97	-20.82	-	-	-	-	-	
ко	-	-	-19.83	-19.68	-	-	-	-	-	

Huge drops without source-language fine-tuning, using only target-language shots



Random Fundamental Problems with Data

We still do not have good (even) evaluation data for many tasks

- If we are working with low-resource languages, we can evaluate only on "lower-level" tasks such as POS tagging, NER, parsing, sentence matching...
- New benchmarking initiatives such as XTREME(-R) and XGLUE help...
- New datasets (e.g., NER data for 10 African languages; NLI dataset for 10 languages of the Americas; a renewed TREC interest in CLIR) help...

Discussion Point: Can We Make Better Training and Evaluation Data for Low-Resource Languages? (A collective bottom-up effort?)



Random Fundamental Problems with Data

Translation-based data creation? Data "imperialism"? (Bird, 2020)

Instructions

Please state to what extent you agree/disagree with each statement on the scale of 1-5 (1-Strongly disagree, 5-Strongly agree)

Questions

Q1. The ASSISTANT helps satisfy the USER's requests. Q2. The USER speaks naturally and sounds like a Russian native speaker.

Q3. The ASSISTANT speaks naturally and sounds like a Russian native speaker.

Q4. I can easily imagine myself mentioning or hearing the proper names referred to in the dialogue (e.g., titles of films or songs, people, places) in a conversation with my Russian friends or family.

A curious study of dataset creation for task-oriented dialogue (Majewska et al., 2022)



Take-Home Messages: Episode I (We only scratched the surface in this talk...)

- mission of democratising language technology
- adaptable solutions, we need to learn from whatever we've got...
- - cross-lingual word embeddings
 - massively multilingual language models -
 - multilingual representation learning; cross-lingual transfer methods -
 - learning
- with)

Multilingual and cross-lingual NLP and IR are vibrant research fields in the

Too many domains, too many languages, dialects -> we need general and

We have covered (at a very shallow level) high-level approaches as well as lower-level cutting-edge approaches to multilingual and cross-lingual NLP

some more advanced topics: modular and parameter-efficient transfer, few-shot

Despite positive trends, many languages are still left behind (and difficult to work



Advanced Topics (We only scratched the surface in this talk...)

- Active learning
- Meta-learning and few-shot adaptation strategies
- Model adaptation to languages with unseen scripts
- Semantic specialisation of general-purpose models

- Injection of linguistic and world knowledge into multilingual text-based models
- Multi-modal multilingual modeling
- Creative applications of multilingual models
- Multilingual speech recognition Speech translation

Data annotation and resource creation in low-resource languages Induction of linguistic structure from pretrained multilingual LMs Learning multilingual word, sentence, and document encoders Unsupervised and weakly supervised Neural Machine Translation

Multilingual and Cross-Lingual NLP and IR: How to Cope? **Better Models and Algorithms**:

- sophisticated modeling/training methods know NLP/ML
- linguistically informed methods know linguistics
- task knowledge know your task

Better Data and Evaluation:

- every piece of relevant data can help be resourceful
- make data if necessary be connected
- track progress with challenging (and natural!) evaluation data

Better Adaptation:

- leverage similarity between languages
- adapt quickly to low-data regimes and new domains





(Multilingual) Natural Language Processing ...and its many applications

Conversational Systems Virtual Assistants Information Search Question Answering



Digital Education Language Learning Assisted Translation



The grand challenge of multilinguality

Digital language divide versus equal opportunities

7,000+ languages; a wide spectrum of tasks and domains

Fact Checking Verification

- High-resource
 languages
- Medium-resource
- Low-resource
- Endangered



Towards <u>Inclusive, Sustainable, Equitable</u> Multilingual NLP

Widening the global reach of NLP: Far-reaching technological and socioeconomic consequences



Are we too obsessed with task performance only and leaderboards?

Towards a more holistic bottom-up approach to multilingual language technology

An analysis of ACL 2021 papers in a "research meta-space"

Søgaard, Vulić, Ruder: Square One Bias in NLP: Towards a Multi-Dimensional Exploration of the Research Manifold (ACL 2022)



Towards <u>Inclusive, Sustainable, Equitable</u> Multilingual NLP

Widening the global reach of NLP: Far-reaching technological and socioeconomic consequences



Other Aspects of Inclusivity and Equity: *Fairness, Cross-Cultural Adaptation, Multi-Modal Learning*





Featuring:











































Once again, big big thankyous and credits to my collaborators...

ObrigadoСпасибо_ • Dank U Danke Obrigado Ngiyabonga Diolch Kasih rima lac)

iv250@cam.ac.uk

