

Wonders and Mysteries of Multilingual Language Models

François Yvon

ISIR — CNRS and Sorbonne Université



Workshop DataIA - ILLS @ CentraleSupélec

Multilingual Natural Language Processing

Computational Language Documentation

(Bayesian) Probabilistic Models

*S. Okabe, P. Godard,
L. Besacier*

Retrieval Augmented MT

Non Autoregressive MT (LevT), imitation learning

J. Grego, J. Xu, M. Bouthors

Multidomain MT

Adapters, curriculum / sampling policy learning

P.M. Quang, J. Grego

Interacting with bitexts

Alignment & edit-based models

*P. Cubaud, J. Grego,
A.K Ngo-Ho, J. Xu*

MT for Open Science

Integrating Terminology & Phraseology in NMT

*S. Abdul-Rauf, R. Bawden,
N. Kübler, etc*

Multilingual LLMs

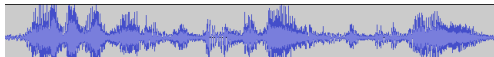
X-lingual Transfer, Probing & Metrics

*G. Wisniewski, H. Schütze,
A. Imani, M. Sabet, etc.*

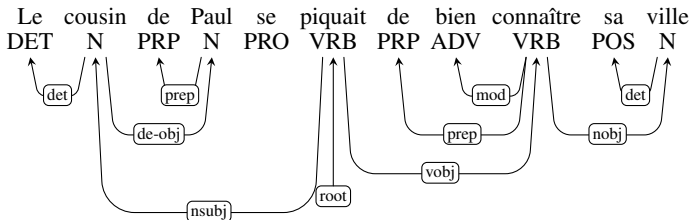
📖 details fyvo.github.io

Natural Language Processing : the Data-Based Revolution

Issues in Linguistic Analysis



lə || kuzɛ̃ || dø || pol || sə || pikɛ̃ || də || bjɛ̃ || ø :də ø :də || bjɛ̃ || konɛtrə || sa || vil



(sp / se-piquer-de-1

: ARG0 (c / cousin)

: ARG1 (k / connaître-2

: ARG0 (c)

: ARG1 (v / ville))

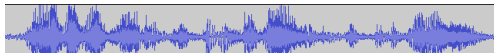
)



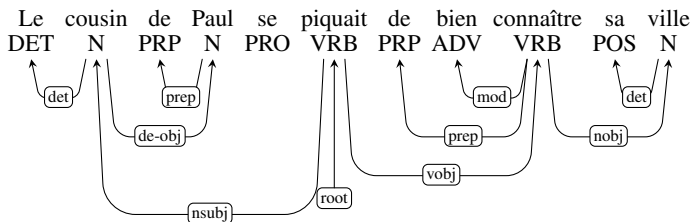
Cousin ?

Natural Language Processing : the Data-Based Revolution

Issues in Linguistic Analysis



lə || kuzɛ̃ || dø || pol || sə || pikɛ̃ || də || bjɛ̃ || ø :də ø :də || bjɛ̃ || konɛtrə || sa || vil



(sp / se-piquer-de-1

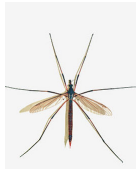
: ARG0 (c / cousin)

: ARG1 (k / connaître-2

: ARG0 (c)

: ARG1 (v / ville))

)



Cousin ?

Natural Language Processing : the Data-Based Revolution

Downstream Applications



Nicolas Sarkozy arbitra la semaine prochaine sur les amendements un bon cru pour le tourisme . Dès cet été des milices et les étrangers sont venus à Paris , bonsoir de fréquentation des hôtels et du tourisme d'affaires sur les six premiers mois de la vie ..

Et direction tout de suite . La Rochelle , où les socialistes font l'heure rentrée politique du PS , réuni au grand complet à l'exception de Dominique Strauss-Kahn , Amandine a -l-elle Hayat . Vous êtes sur place . Ségolène Royal est apparemment à très présente un depuis hier ..

Oui . Marianne et c'est clairement aujourd'hui la journée de Ségolène Royal , elle entend bien marquer des points . L'an dernier , elle avait survolé ses universités d'été avec des sauter , c'est bien cette fois , elle sera présente pendant les trois jours et cela commencé soir avec un dîner de presse . Cela s'est poursuivie ce matin avec de nouveau une foule de médias pour une conférence de presse sur la crise laitière et puis cet après- midi , ce sera bien sûr le traditionnel discours d'ouverture devant 4000 la ..

- ▶ **text analysis** : information retrieval and classification, question-answering, spell checking, information extraction, NLU, etc
- ▶ **text generation** : machine translation, summarization, dialog, TTS, etc

Natural Language Processing : the Data-Based Revolution

Processing Sequences of Symbols with Hidden Structure

Input	$\mathbf{w} = w_1 \dots w_T \in \Sigma_1^T$ w_t : <i>a phoneme, a letter, a morpheme, a word, a sentence</i>
Output	$\mathbf{y} \in \Sigma_2$: <i>sequence classification</i> y : <i>filter, topic, polarity, language, author</i>
Output	$\mathbf{y} \in \Sigma_2^T$: <i>sequence labelling</i> y_t : <i>PoS, feature, bracketing, etc</i>
Output	$\mathbf{y} \in \Sigma_2^*$: <i>sequence transduction</i> y : <i>transcription, translation, answer, correction, summary, etc</i>
Output	$\mathbf{y} \in \mathcal{G}$: <i>sequence parsing</i> y : <i>morphological or syntactic tree, semantic graph, discourse tree</i>

- ▶ Ambiguous or unobserved boundaries
- ▶ Ubiquitous ambiguity of basic units w_t
- ▶ Structural constraints over possible \mathbf{x}, \mathbf{y} (unobserved)
- ▶ Internal structure within Σ_1 and Σ_2 (unobserved)
- ▶ Masive variability of data distribution (domain, genre, style, register)

Natural Language Processing : the Data-Based Revolution

The “Unreasonable Effectiveness” of Supervised Learning in NLP (1993-)

There is nothing but tabular data

sent_id = fr_partut-ud-184

text = *La sécurité des transports a trop fait l'actualité ces derniers temps :*

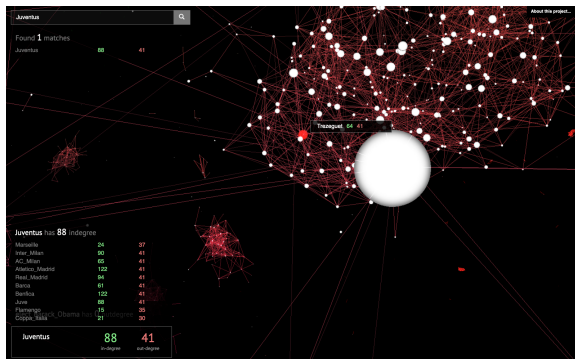
1	La	le	DET	RD	Def=Def Gdr=Fem Num=Sing PronType=Art
2	sécurité	sécurité	NOUN	S	Gdr=Fem Num=Sing
3-4	des	_	_	_	_
3	de	de	ADP	E	_
4	les	le	DET	RD	Def=Def Num=Plur PronType=Art
5	transports	transport	NOUN	S	Gdr=Masc Num=Plur
6	a	avoir	AUX	VA	Mood=Ind Num=Sing Person=3 Tense=Pres
7	trop	trop	ADV	B	
8	fait	faire	VERB	V	Gdr=Masc Num=Sing Tense=Past VbForm=Part
9	l'	le	DET	RD	Def=Def Num=Sing PronType=Art
10	actualité	actualité	NOUN	S	Gdr=Fem Num=Sing
11	ces	ce	DET	DD	Num=Plur PronType=Dem
12	derniers	dernier	ADJ	NO	Gdr=Masc Num=Plur NumType=Ord
13	temps	temps	NOUN	S	Gdr=Masc
14	:	:	PUNCT	FC	_

Picking “low-hanging fruits” with simple learners (classifiers, sequence learning)

Monolingual Pretrained Language Models are Powerful

Modern Lexical Representations : Word Embeddings

Map words w_i in context c onto vector-space representations x as $x = f_{\theta}(w_i; c)$. $f_{\theta}()$ is a neural net eg. Transformer [Vaswani et al., 2017] with parameters θ .



The structure of word spaces : related words are close

Monolingual Pretrained Language Models are Powerful

Learning θ

Gigantic collections of texts, simple **auxiliary tasks** with **natural annotations**

1. Predict next word given prefix : **pure decoder**

Longtemps je me suis couché [mask] ... \Rightarrow train $P_{\theta}(w_t|w_{<t})$

GPT-2 [Radford et al., 2019], OPT [Zhang et al., 2022], GPT-J [Wang and Komatsuzaki, 2021], etc

2. Predict missing word given bidirectional context : **pure encoder**

Longtemps je me suis couché de [mask] bonne heure. \Rightarrow train $P_{\theta}(w_t|w_{-t})$

BERT [Devlin et al., 2019], Roberta [Liu et al., 2019], CamemBERT [Martin et al., 2020] etc.

3. Denoising sequence to sequence : **encoder-decoder**

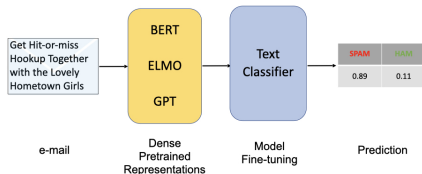
Longtemps je couché suis de bnone || Longtemps je me suis couché de bonne heure.
 \Rightarrow train $P_{\theta}(w|\tilde{w})$

BART [Lewis et al., 2020], T5 [Raffel et al., 2020], etc.

Monolingual Pretrained Language Models are Powerful

Using θ

1. Learn + fine-tune task-adapted model $h_{\phi, \theta}(w; c) = h_{\phi}(f_{\theta}(w; c))$



2. With causal LMs, multi-purpose text generation via **prompting**

GEN Of course. In Chorukur, Monday is ilopagar, Tuesday ilopager, ...
Wednesday ilopagar, Thursday ilopagir ...

Q&A Answer this : What are the birth date and place of Ludvík Vaculík? ...
23 July 1926, in Brumov, Moravia

SA "This Czech writer has written some the most wonderful French novels."
is a positive comment? ... [Yes | No]

Monolingual Pretrained Language Models are Powerful

Using θ

1. Learn + fine-tune task-adapted model $h_{\phi, \theta}(w; c) = h_{\phi}(f_{\theta}(w; c))$
2. With causal LMs, multi-purpose text generation via prompting

- ☺ LM pre-training is mostly unsupervised
- ☺ Instructions can be learned or fine-tuned
- ☺ Extremely successful for “high resource” languages
- ☺ Annotation, corpus and model size matter

Monolingual Pretrained Language Models are Powerful

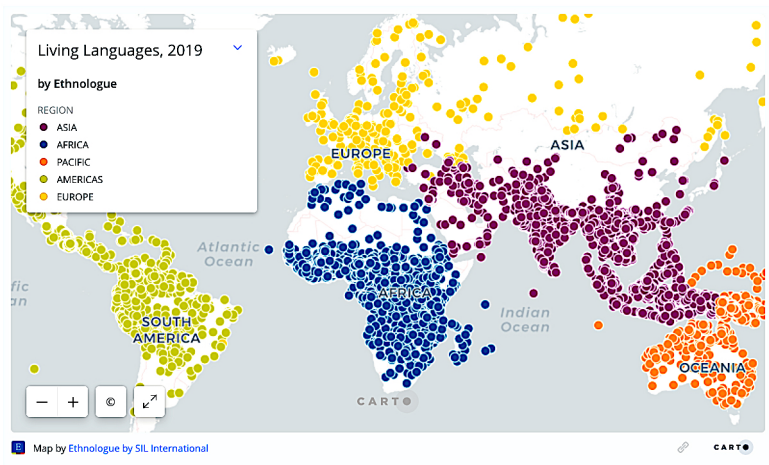
Using θ

1. Learn + fine-tune task-adapted model $h_{\phi, \theta}(w; c) = h_{\phi}(f_{\theta}(w; c))$
2. With causal LMs, multi-purpose text generation via prompting

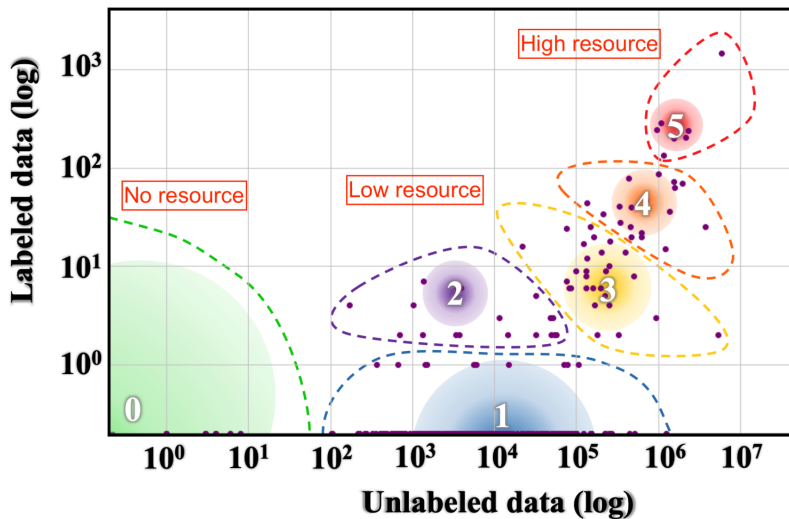
- ☺ LM pre-training is mostly unsupervised
- ☺ Instructions can be learned or fine-tuned
- ☺ Extremely successful for “high resource” languages

- ☺ Annotation, corpus and model size matter

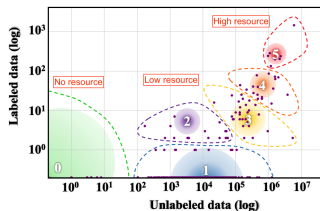
So Many Languages, most of them “Low-Resource”



So Many Languages, most of them “Low-Resource”



So Many Languages, most of them “Low-Resource”



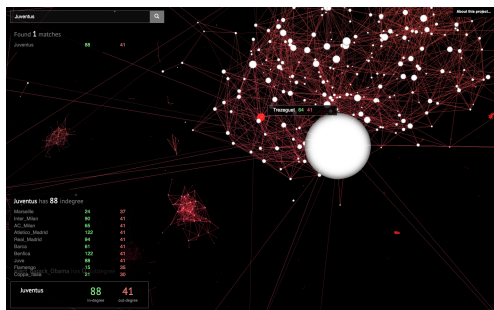
Cls	Example languages	#Lang	#Spkr	% Lang
0	Dahalo, Warlpiri, Popoloca, Wallisian, Bora	2191	1.2B	88.38%
1	Cherokee, Fijian, Greenlandic, Bhojpuri, Navajo	222	30M	5.49%
2	Zulu, Konkani, Lao, Maltese, Irish	19	5.7M	0.36%
3	Indonesian, Ukrainian, Cebuano, Afrikaans, Hebrew	28	1.8B	4.42%
4	Russian, Hungarian,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

Data and analysis from [Joshi et al., 2020]

Multilingual Pretrained Language Models - mPLMs

Multilingual Lexical Representations : Word Embeddings

Map words w onto vector-space representations x in context c as $x = f_{\theta}(w; c)$
 $f_{\theta}()$ a complex neural net with parameters θ .



Neighbours should be semantically similar **cross-linguistically**

Multilingual Pretrained Language Models - mPLMs

Multilingual representations of texts

_Tous _les _être s _humain s _na issent _libre s _et _ég aux _en _digni té _et _en
_droits . _Ils _sont _do u és _de _raison _et _de _conscience _et _doivent _agir
_les _uns _en vers _les _autres _dans _un _esprit _de _frater n ité .

_Všichni _lidé _rod í _se _svobod ní _a _sobě _rov ní _co _do _d ů stoj nosti _a
_práv . _Jsou _na dán í _rozum em _a _s vědomí m _a _mají _spolu _jedna t _v
_du chu _brat r ství .

_Tutti _gli _esse ri _umani _na sconno _liberi _ed _e gu ali _in _digni tà _e _diritti
_ . _Es si _sono _do tati _di _ragione _e _di _coscienza _e _devono _agir e _gli _uni
_verso _gli _altri _in _spirito _di _fra tella nza .

Low-level, language independent segmentations

Multilingual Pretrained Language Models - mPLMs

Multilingual representations of texts

_Tous _les _être s _humain s _na issent _libre s _et _ég aux _en _digni té _et _en
_droits . _Ils _sont **_do** u és _de _raison _et _de _conscience _et _doivent **_agir**
_les _uns _en vers _les _autres _dans _un _esprit _de _frater n ité .

_Všichni _lidé _rod í _se _svobod ní _a _sobě _rov ní _co **_do** _d ů stoj nosti _a
_práv . _Jsou _na dán í _rozum em _a _s vědomí m _a _mají _spolu _jedna t _v
_du chu _brat r ství .

_Tutti _gli _esse ri _umani _na sconno _liberi _ed _e gu ali _in _digni tà _e _diritti
_Es si _sono **_do** tati _di _ragione _e _di _coscienza _e _devono **_agir** e _gli _uni
_verso _gli _altri _in _spirito _di _fra tella nza .

Low-level, language independent segmentations

Multilingual Pretrained Language Models - mPLMs

Learning multilingual θ

Multilingual Mixtures of corpora, simple **auxiliary tasks** with **natural annotations** :

1. Predict next word given prefix : **pure decoder**.
mGPT [Shliachko et al., 2022], XGLM [Lin et al., 2021], etc
2. Predict missing word given full context : **pure encoder**.
mBERT [Devlin et al., 2019], XLM-R [Conneau et al., 2020]
3. Denoising sequence to sequence : **encoder-decoder**
mBART50 [Liu et al., 2020b], mT5 [Xue et al., 2020], etc

+ **Complementary objectives** to bridge between languages :
parallel corpora, bilingual dictionaries, synthetic code-switched data [Conneau et al., 2020, Chi et al., 2021], etc.

Multilingual Pretrained Language Models - mPLMs

Using multilingual θ

1. Learn / fine-tune task-adapted model $h_{\phi, \theta}(w; c) = h_{\phi}(f_{\theta}(w; c))$ on L_1 ,
perform zero-shot X-lingual transfer to L_2
2. Multilingual text generation with prompting

Translate into English

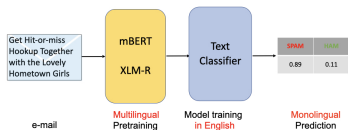
“ By the end of the year, we will have seven new pharmacists. ” :

D’ici la fin de l’année, nous aurons sept nouveaux pharmaciens.

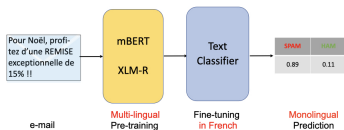
- ☺ Hardly more difficult than monolingual
- ☺ Bring low resource languages up to steam

The effectiveness of multilingual PLMs

One model handling multiple languages



English spam filter

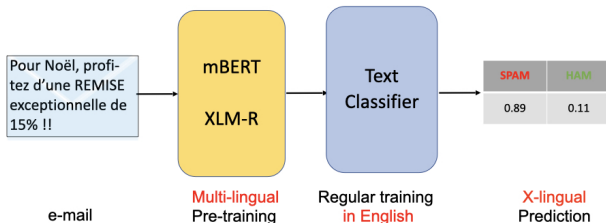


French spam filter

Effective for many languages / tasks [Pires et al., 2019, Wu and Dredze, 2020]

The effectiveness of multilingual PLMs

Knowledge transfer between languages

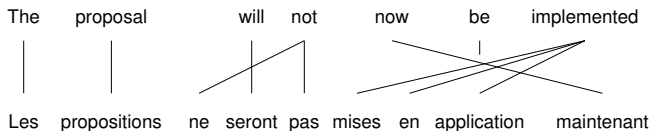


Multilingual zero-shot spam filter

- ☺ No foreign training data required
- ☺ Effective for many languages and tasks : X-GLUE [Liang et al., 2020]
X-TREME+ [Hu et al., 2020, Ruder et al., 2023]

The effectiveness of multilingual PLMs

Multilingual Embeddings Align Lexical Representations



- ▶ Word Alignment = Bipartite Graph
- ▶ Multilingual embeddings + graph algorithms = **unsupervised alignments**
- ▶ Aligns code-switching and multilingual texts
- ▶ Also : **unsupervised sentence alignments**

Method	ENG-CES		ENG-DEU		ENG-FAS	
	F_1	AER	F_1	AER	F_1	AER
IBM2	.76	.25	.71	.29	.57	.43
IBM4	.75	.26	.77	.23	.51	.49
eflomal	.85	.15	.77	.23	.61	.39
fastText	.70	.30	.60	.40	.50	.50
mBERT[8]	.87	.13	.79	.21	.67	.33
XLM-R[8]	.87	.13	.79	.21	.70	.30

Word alignment, Argmax method

Jalili Sabet et al. [2020]

The effectiveness of multilingual PLMs

Unsupervised Alignments in Large Scale Evaluation

Glott-500 [ImaniGooghari et al., 2023] : a multilingual
(XML-type) model for 500 languages

Language-Script	XML-R-B	XML-L	Glott500	Language-Script	XML-R-B	XML-L	Glott500	Language-Script	XML-R-B	XML-L	Glott500
af_Latn	2.50	2.25	6.23	nb_Latn	7.61	4.09	6.26	nl_Latn	3.19	4.18	4.33
af_Grek	2.79	3.38	4.77	nl_Latn	2.78	3.28	3.48	nl_Latn	1.43	1.84	2.25
ar_Latn	4.02	4.49	6.39	np_Latn	3.32	4.06	6.39	no_Latn	2.01	2.64	3.63
ar_Arab	3.77	4.60	7.19	or_Oyja	2.75	3.28	3.22	or_Latn	1.56	2.16	2.19
az_Latn	4.01	4.83	5.88	os_Latn	2.62	3.05	3.78	os_Latn	2.44	2.71	5.93
az_Cyrl	4.54	5.24	8.21	pa_Latn	2.27	3.15	4.13	pa_Latn	2.79	3.59	4.67
ba_Latn	3.12	3.80	4.19	pa_Cyrl	3.99	4.56	7.37	pa_Latn	2.82	3.12	4.63
ba_Cyrl	2.22	2.67	4.74	pl_Latn	3.54	3.81	6.83	pl_Latn	2.81	3.08	6.88
bn_Latn	3.85	4.62	5.75	pl_Latn	2.94	3.37	4.87	pl_Latn	2.88	3.14	4.28
bn_Arab	4.54	4.48	7.66	ps_Latn	2.31	3.24	4.65	ps_Latn	1.43	1.83	2.36
bi_Latn	2.81	3.17	4.94	ps_Latn	1.28	1.65	3.85	ps_Latn	2.47	2.83	4.23
bi_Arab	3.26	3.92	4.88	pt_Latn	2.81	3.46	6.49	pt_Cyrl	2.74	3.68	4.13
bo_Latn	4.06	4.19	6.03	pt_Latn	2.69	3.24	4.58	pt_Latn	2.43	3.23	4.74
bo_Arab	1.63	1.89	3.79	ro_Latn	3.17	3.90	5.38	ro_Latn	3.41	4.13	6.13
br_Latn	3.19	3.57	5.88	ro_Latn	1.74	3.99	3.20	ro_Latn	3.18	4.06	7.45
br_Arab	3.36	3.98	4.92	ru_Latn	1.75	2.11	2.31	ru_Latn	3.05	4.06	6.78
ca_Latn	2.74	3.28	6.01	ru_Cyrl	1.54	1.35	2.46	ru_Cyrl	2.71	2.83	3.33
ca_Arab	2.76	3.28	4.58	ru_Cyrl	2.90	3.42	4.46	ru_Latn	2.16	2.56	3.68
ca_Latn	3.00	3.43	5.48	ru_Latn	2.62	3.02	4.10	ru_Latn	2.01	2.29	3.77
ca_Arab	1.83	2.07	3.51	sv_Latn	1.99	2.51	2.56	sv_Cyrl	2.89	3.48	4.72
ch_Latn	1.88	2.23	3.48	sq_Latn	2.42	3.15	4.41	sq_Arab	2.58	3.11	3.61
ch_Arab	2.80	3.48	3.61	sr_Latn	1.75	2.03	2.71	sq_Latn	2.26	2.78	3.79
ch_Latn	3.52	4.24	4.49	sr_Latn	3.09	3.87	4.85	sr_Cyrl	3.71	5.96	7.47
ch_Arab	2.76	3.38	4.45	sv_Latn	2.18	2.74	3.32	sr_Arab	1.88	2.88	3.96
cm_Latn	3.03	3.98	4.88	sv_Latn	2.64	3.49	3.69	sv_Latn	2.29	3.07	3.63
cm_Arab	2.26	2.73	3.71	th_Latn	3.32	3.85	6.87	th_Cyrl	2.73	3.28	7.24
cn_Latn	1.11	1.06	2.81	th_Latn	4.00	4.60	6.84	th_Latn	3.12	3.96	4.91
cn_Arab	1.46	1.80	2.23	tl_Latn	3.21	3.85	5.41	tl_Cyrl	2.41	3.06	3.66
en_Latn	2.60	3.23	4.79	tl_Latn	2.58	3.09	4.79	tl_Latn	2.96	3.64	5.14
en_Cyrl	3.18	4.15	4.28	tr_Deva	3.02	3.97	6.21	tl_Arab	3.99	4.68	6.09
en_Latn	2.14	2.39	3.62	tr_Latn	1.88	2.34	3.39	tl_Latn	2.87	3.59	4.24
en_Arab	2.18	2.54	4.26	tr_Latn	2.75	3.47	4.24	tr_Latn	3.84	3.74	4.33
en_Latn	1.94	2.23	4.60	uk_Latn	2.85	3.60	4.90	uk_Latn	2.44	2.86	4.53
en_Arab	1.44	1.78	2.63	uk_Cyrl	3.39	4.27	5.19	uk_Latn	4.77	4.19	4.19
es_Latn	3.21	3.64	6.39	uz_Latn	2.46	3.14	4.88	uz_Latn	3.87	1.03	1.12
es_Arab	3.85	4.24	5.90	uz_Latn	3.60	4.41	7.41	uz_Cyrl	4.23	4.23	6.09
es_Latn	2.23	2.63	3.68	uz_Arab	3.88	4.81	5.83	uz_Latn	2.93	3.35	4.19
es_Arab	1.85	2.41	3.92	uz_Latn	3.21	4.14	5.82	uz_Latn	2.01	2.66	3.87
fr_Latn	2.80	3.43	3.68	uz_Deva	3.28	3.80	5.81	uz_Latn	3.41	3.41	3.48
fr_Latn	2.92	3.48	4.88	uz_Latn	3.29	4.06	5.74	uz_Latn	3.25	4.00	5.17
fr_Deva	3.59	3.80	4.13	uz_Latn	3.06	3.82	5.41	uz_Latn	2.44	2.68	3.68
fr_Latn	2.94	3.28	4.77	uz_Latn	2.76	3.09	5.96	uz_Latn	2.94	2.26	2.96
fr_Deva	3.61	3.90	5.19	uz_Latn	2.37	3.28	3.59	uz_Latn	2.19	2.57	3.67
fr_Latn	2.48	2.93	4.05	uz_Latn	2.21	2.74	2.85	uz_Latn	3.22	3.74	5.31
fr_Latn	3.14	3.33	4.28	uz_Latn	2.69	2.70	4.20	uz_Latn	3.77	4.18	5.63
fr_Latn	3.52	3.86	5.19	uz_Deva	2.71	2.77	2.91	uz_Latn	4.29	5.13	5.19
fr_Latn	4.14	4.24	7.62	uz_Deva	3.27	3.20	4.39	uz_Latn	3.65	4.65	5.36
fr_Latn	4.4	4.18	5.47	uz_Cyrl	3.20	3.32	5.85	uz_Cyrl	3.52	4.60	5.63
fr_Latn	4.54	4.18	5.62	uz_Latn	1.89	3.23	2.46	uz_Latn	3.67	4.39	6.44
fr_Latn	1.70	2.00	2.42	uz_Latn	2.93	3.44	4.56				

Table 22: Accuracy of XLM-R-B, XLM-R-L, and Glott500-m on Round Trip Alignment (Part II).

Round Trip alignments, Part II

From 100^2 to 500^2 unsupervised word alignments

- ▶ New challenges
 - ▶ data selection and filtering
 - ▶ language identification at scale
- ▶ New evaluation methods
 - ▶ Round trip alignment
 - ▶ Sentence retrieval

The effectiveness of multilingual PLMs

Multilingual Text Generation

If the French says : Ça n'arrête pas, nous sommes bien placés pour le savoir...

then the English should say : It doesn't stop, we know that well...

Model	En-Fr	Fr-En
GPT-2 [1,5b] [Radford et al., 2019]	5	11.5
GPT-3 [175b] [Brown et al., 2020]	21.2	25.2
PALM [540B] [Chowdhery et al., 2022]	38.5	41.1
BLOOM [176B] [BigScience et al., 2022]	??	??
SOTA [Liu et al., 2020a]	44.1	-

Closing the gap with well trained bilingual MT, really ?

The effectiveness of multilingual PLMs

Understanding MT performance of BLOOM

- ▶ Methods :
 - ▶ Multiple prompts, in multiple languages
 - ▶ 3 datasets and tasks
 - ▶ Dozens of language pairs (H-H, L-H, L-L)
 - ▶ Multiple metrics
- ▶ Main takeaways
 - ▶ a call for clarity when reporting results
 - ▶ the failures of 0-shot learning, “language hallucinations”
 - ▶ “Free-rider effects” and language similarities
 - ▶ poor MT quality for LR languages

Src ↓	Trg →	ar	en	es	fr	zh
ar	BLOOM	–	40.3	23.3	33.1	17.7
	M2M	–	25.5	16.7	25.7	13.1
en	BLOOM	28.2	–	29.4	45.0	26.7
	M2M	17.9	–	25.6	42.0	19.3
es	BLOOM	18.8	32.7	–	24.8	20.9
	M2M	12.1	25.1	–	29.3	14.9
fr	BLOOM	23.4	45.6	27.5	–	23.2
	M2M	15.4	37.2	25.6	–	17.6
zh	BLOOM	15.0	30.5	20.5	26.0	–
	M2M	11.6	20.9	16.9	24.3	–

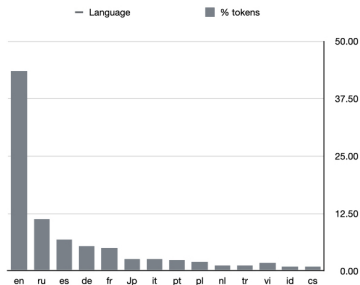
(a) High-resource language pairs.

Src ↓	Trg →	en	fr	hi	id	vi
en	BLOOM	–	45.0	27.2	39.0	28.5
	M2M	–	42.0	28.1	37.3	35.1
fr	BLOOM	45.6	–	18.5	31.4	32.8
	M2M	37.2	–	22.9	29.1	30.3
hi	BLOOM	35.1	27.6	–	–	–
	M2M	27.9	25.9	–	–	–
id	BLOOM	43.2	30.4	–	–	–
	M2M	33.7	30.8	–	–	–
vi	BLOOM	38.7	26.8	–	–	–
	M2M	29.5	25.8	–	–	–

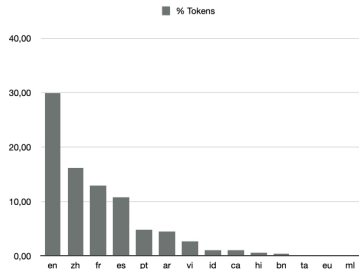
(b) High→mid-resource language pairs.

The effectiveness of multilingual PLMs

Caveats



Languages in mT5
(104 languages)



Languages in Bloom
(59 languages)

- ☹ English centric
- ☹ Poorly documented, hard to reproduce – what is a language anyway?

What we would like to know

Why can this work ?

- ▶ Lexical overlap not necessary

L1 : Longtemps je me suis couché de bonne heure .

L2 : Mpohufnqt kf nf tvjt dpvdiê ef cpoof ifvsf !

- ▶ Transfer breaks with reversal or words shuffle

L1 : Longtemps je me suis couché de bonne heure .

L2 : heure bonne de douché suis me je longtemps

- ▶ Vocabulary alignment helps
- ▶ Language typological similarity helps

Analyses in [K et al., 2020, Dufter and Schütze, 2020, Deshpande et al., 2022, ImaniGooghari et al., 2023], etc

What we would like to know

Why can this work ?

- ▶ Lexical overlap not necessary

L1 : Longtemps je me suis couché de bonne heure .

L2 : Mpozufnqt kf nf tvjt dpvdiê ef cpoof ifvsf !

- ▶ Transfer breaks with reversal or words shuffle

L1 : Longtemps je me suis couché de bonne heure .

L2 : . heure bonne de douché suis me je longtemps

- ▶ Vocabulary alignment helps
- ▶ Language typological similarity helps

Analyses in [K et al., 2020, Dufter and Schütze, 2020, Deshpande et al., 2022, ImaniGooghari et al., 2023], etc

? Impact of language distributions

? Impact of model size

? Impact of number of languages - curse of multilinguality

? Which linguistic properties help / break X-lingual transfer

? Metrics for “coverage”

What we would like to know

Can mPLMs handle truly multilingual texts?

- ? Generate **code-switched** output

Et le premier ministre nous répond que *a farmer is a farmer a Canadian is a Canadian* d'un bout à l'autre du Canada.

Autrement dit *they are getting out of the closet* parce que cela leur donne le droit d avoir deux enfants.

(Exemples from [Carpuat, 2014])

- ? Generate text with **multilingual prompts**

Přeložit do angličtiny : Ça n'arrête pas, nous sommes bien placés pour le savoir... ; :

It doesn't stop, we know that well...

- ? Answer questions X-linguistically
- ? Generate summaries from **multilingual sources**

Take Aways

Large mPLMs serve practical purposes

They learn linguistic features

They display amazing emerging properties

Multilingual models also useful in MT, parsing, etc

They improve technological support for many languages [Ruder et al., 2023]

They remain poorly understood

Bibliography I

Workshop BigScience, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, Dragomir Radev, Eduardo González Ponferrada, Efrat Levkovizh, Ethan Kim, Eyal Bar Natan, Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady Elsahar, Hamza Benyamina, Hieu Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jörg Frohberg, Joseph Tobing, Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro Von Werra, Leon Weber, Long Phan, Loubna Ben allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, María Grandury, Mario Šaško, Max Huang, Maximin Coavoux, Mayank Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad A. Jauhar, Mustafa Ghaleb, Nishant Subramani, Nora Kassner, Nurulaqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rheza Harliman, Rishi Bommasani, Roberto Luis López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, Shayne Longpre, Somaieh Nikpoor,

Bibliography II

Stanislav Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, Davut Emre Taşar, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaffin, Arnaud Stiegler, Debajyoti Datta, Eliza Szczechla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M Saiful Bari, Maged S. Al-shaibani, Matteo Manica, Nihal Nayak, Ryan Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Fevry, Trishala Neeraj, Urmish Thakker, Vikas Raunak, Xiangru Tang, Zheng-Xin Yong, Zhiqing Sun, Shaked Brody, Yallow Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre François Lavallée, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, Suraj Patil, Tim Dettmers, Ahmed Baruwa, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, Aurélie Névéol, Charles Lovering, Dan Garrette, Deepak Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, Jessica Zosa Forde, Jordan Clive, Jungo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, Miruna Clinciu, Najoung Kim, Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, Rui Zhang, Ruochen Zhang, Sebastian Gehrmann, Shachar Mirkin, Shani Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, Tomasz Limisiewicz, Verena Rieser, Vitaly

Bibliography III

Protasov, Vladislav Mikhailov, Yada Pruksachatkun, Yonatan Belinkov, Zachary Bamberger, Zdeněk Kasner, Alice Rueda, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ana Santos, Anthony Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh HajiHosseini, Bahareh Behroozi, Benjamin Ajibade, Bharat Saxena, Carlos Muñoz Ferrandis, Danish Contractor, David Lansky, Davis David, Douwe Kiela, Duong A. Nguyen, Edward Tan, Emi Baylor, Ezinwanne Ozoani, Fatima Mirza, Frankline Ononiwu, Habib Rezanejad, Hessie Jones, Indrani Bhattacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jesse Passmore, Josh Seltzer, Julio Bonis Sanz, Livia Dutra, Mairon Samagaio, Maraim Elbadri, Margot Mieskes, Marissa Gerchick, Martha Akinlolu, Michael McKenna, Mike Qiu, Muhammed Ghauri, Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nour Fahmy, Olanrewaju Samuel, Ran An, Rasmus Kromann, Ryan Hao, Samira Alizadeh, Sarmad Shubber, Silas Wang, Sourav Roy, Sylvain Viguier, Thanh Le, Tobi Oyebade, Trieu Le, Yoyo Yang, Zach Nguyen, Abhinav Ramesh Kashyap, Alfredo Palasciano, Alison Callahan, Anima Shukla, Antonio Miranda-Escalada, Ayush Singh, Benjamin Beilharz, Bo Wang, Caio Brito, Chenxi Zhou, Chirag Jain, Chuxin Xu, Clémentine Fourrier, Daniel León Perrián, Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrmann, Gabriel Altay, Giyaseddin Bayrak, Gully Burns, Helena U. Vrabc, Imane Bello, Ishani Dash, Jihyun Kang, John Giorgi, Jonas Golde, Jose David Posada, Karthik Rangasai Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, Maria A Castillo, Marianna Nezhurina, Mario Sängner, Matthias Samwald, Michael Cullan, Michael Weinberg, Michiel De Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, Nicholas Michio Broad, Nikolaus Muellner, Pascale Fung, Patrick Haller,

Bibliography IV

Ramya Chandrasekhar, Renata Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda, Shlok S Deshmukh, Shubhanshu Mishra, Sid Kiblawi, Simon Ott, Sinee Sang-aaronsiri, Srishti Kumar, Stefan Schweter, Sushil Bharati, Tanmay Laud, Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yash Shailesh Bajaj, Yash Venkatraman, Yifan Xu, Yingxin Xu, Yu Xu, Zhe Tan, Zhongli Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and Thomas Wolf. BLOOM : A 176b-parameter open-access multilingual language model. CoRR, abs/2211.05100, 2022. doi : 10.48550/arXiv.2211.05100. URL <https://doi.org/10.48550/arXiv.2211.05100>.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.

Marine Carpuat. Mixed language and code-switching in the Canadian Hansard. In Proceedings of the First Workshop on Computational Approaches to Code Switching, pages 107–115, Doha, Qatar, October 2014. Association for Computational Linguistics. doi : 10.3115/v1/W14-3913. URL <https://aclanthology.org/W14-3913>.

Bibliography V

- Zewen Chi, Li Dong, Furu Wei, Nan Yang, Saksham Singhal, Wenhui Wang, Xia Song, Xian-Ling Mao, Heyan Huang, and Ming Zhou. InfoXLM : An information-theoretic framework for cross-lingual language model pre-training. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics : Human Language Technologies, pages 3576–3588, Online, June 2021. Association for Computational Linguistics. doi : 10.18653/v1/2021.naacl-main.280. URL <https://aclanthology.org/2021.naacl-main.280>.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. PaLM : scaling language modeling with pathways, 2022. URL <https://arxiv.org/abs/2204.02311>.

Bibliography VI

- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 8440–8451, Online, July 2020. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/2020.acl-main.747>.
- Ameet Deshpande, Partha Talukdar, and Karthik Narasimhan. When is BERT multilingual? isolating crucial ingredients for cross-lingual transfer. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics : Human Language Technologies, pages 3610–3623, Seattle, United States, July 2022. Association for Computational Linguistics. doi : 10.18653/v1/2022.naacl-main.264. URL <https://aclanthology.org/2022.naacl-main.264>.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT : Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics : Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi : 10.18653/v1/N19-1423. URL <https://www.aclweb.org/anthology/N19-1423>.

Bibliography VII

- Philipp Dufter and Hinrich Schütze. Identifying elements essential for BERT’s multilinguality. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4423–4437, Online, November 2020. Association for Computational Linguistics. doi : 10.18653/v1/2020.emnlp-main.358. URL <https://aclanthology.org/2020.emnlp-main.358>.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. XTREME : A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation. In Hal Daumé III and Aarti Singh, editors, Proceedings of the 37th International Conference on Machine Learning, volume 119 of Proceedings of Machine Learning Research, pages 4411–4421. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/hu20b.html>.
- Ayyoob ImaniGooghari, Peiqin Lin, Amir Hossein Kargaran, Silvia Severini, Masoud Jalili Sabet, Nora Kassner, Chunlan Ma, Helmut Schmid, André F. T. Martins, François Yvon, and Hinrich Schütze. Glot500 : Scaling Multilingual Corpora and Language Models to 500 Languages, 2023.
- Masoud Jalili Sabet, Philipp Dufter, François Yvon, and Hinrich Schütze. SimAlign : High quality word alignments without parallel training data using static and contextualized embeddings. In Findings of the Association for Computational Linguistics : EMNLP 2020, pages 1627–1643, Online, November 2020. Association for Computational Linguistics. doi : 10.18653/v1/2020.findings-emnlp.147. URL <https://www.aclweb.org/anthology/2020.findings-emnlp.147>.

Bibliography VIII

- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. The state and fate of linguistic diversity and inclusion in the NLP world. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 6282–6293, 2020.
- Karthikeyan K, Zihan Wang, Stephen Mayhew, and Dan Roth. Cross-lingual ability of multilingual bert : An empirical study. In Proc. International Conference on Learning Representations, 2020. URL <https://openreview.net/forum?id=HJeT3yrtDr>.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART : Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7871–7880, Online, July 2020. Association for Computational Linguistics. doi : 10.18653/v1/2020.acl-main.703. URL <https://www.aclweb.org/anthology/2020.acl-main.703>.
- Yaobo Liang, Nan Duan, Yeyun Gong, Ning Wu, Fenfei Guo, Weizhen Qi, Ming Gong, Linjun Shou, Daxin Jiang, Guihong Cao, et al. XGLUE : A new benchmark dataset for cross-lingual pre-training, understanding and generation. arXiv preprint arXiv :2004.01401, 2020.
- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona Diab, Veselin Stoyanov, and Xian Li. Few-shot learning with multilingual language models, 2021. URL <https://arxiv.org/abs/2112.10668>.

Bibliography IX

- Xiaodong Liu, Kevin Duh, Liyuan Liu, and Jianfeng Gao. Very deep transformers for neural machine translation, 2020a. URL <https://arxiv.org/abs/2008.07772>.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa : A robustly optimized bert pretraining approach. [arXiv preprint arXiv :1907.11692](https://arxiv.org/abs/1907.11692), 2019.
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine translation. [arXiv preprint arXiv :2001.08210](https://arxiv.org/abs/2001.08210), 2020b.
- Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamé Seddah, and Benoît Sagot. CamemBERT : a tasty French language model. In [Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics](#), pages 7203–7219, Online, July 2020. Association for Computational Linguistics. doi : 10.18653/v1/2020.acl-main.645. URL <https://aclanthology.org/2020.acl-main.645>.
- Telmo Pires, Eva Schlinger, and Dan Garrette. How multilingual is multilingual BERT? In [Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics](#), pages 4996–5001, Florence, Italy, July 2019. Association for Computational Linguistics. doi : 10.18653/v1/P19-1493. URL <https://www.aclweb.org/anthology/P19-1493>.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.

Bibliography X

- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research, 21(140) :1–67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- Sebastian Ruder, Jonathan H. Clark, Alexander Gutkin, Mihir Kale, Min Ma, Massimo Nicosia, Shruti Rijhwani, Parker Riley, Jean-Michel A. Sarr, Xinyi Wang, John Wieting, Nitish Gupta, Anna Katanova, Christo Kirov, Dana L. Dickinson, Brian Roark, Bidisha Samanta, Connie Tao, David I. Adelani, Vera Axelrod, Isaac Caswell, Colin Cherry, Dan Garrette, Reeve Ingle, Melvin Johnson, Dmitry Panteleev, and Partha Talukdar. Xtreme-up : A user-centric scarce-data benchmark for under-represented languages, 2023.
- Oleh Shliazhko, Alena Fenogenova, Maria Tikhonova, Vladislav Mikhailov, Anastasia Kozlova, and Tatiana Shavrina. mGPT : few-shot learners go multilingual, 2022. URL <https://arxiv.org/abs/2204.07580>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems 30, pages 5998–6008. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>.
- Ben Wang and Aran Komatsuzaki. GPT-J-6B : A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>, 2021.

Bibliography XI

- Shijie Wu and Mark Dredze. Do explicit alignments robustly improve multilingual encoders? In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4471–4482, Online, November 2020. Association for Computational Linguistics. doi : 10.18653/v1/2020.emnlp-main.362. URL <https://www.aclweb.org/anthology/2020.emnlp-main.362>.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mT5 : A massively multilingual pre-trained text-to-text transformer. arXiv preprint arXiv :2010.11934, 2020.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. OPT : Open pre-trained transformer language models, 2022. URL <https://arxiv.org/abs/2205.01068>.