

Script Normalization for Unconventional Writing of Under-Resourced Languages

Sina Ahmadi and Antonios Anastasopoulos
George Mason University

The 61st Annual Meeting of the
Association for Computational Linguistics

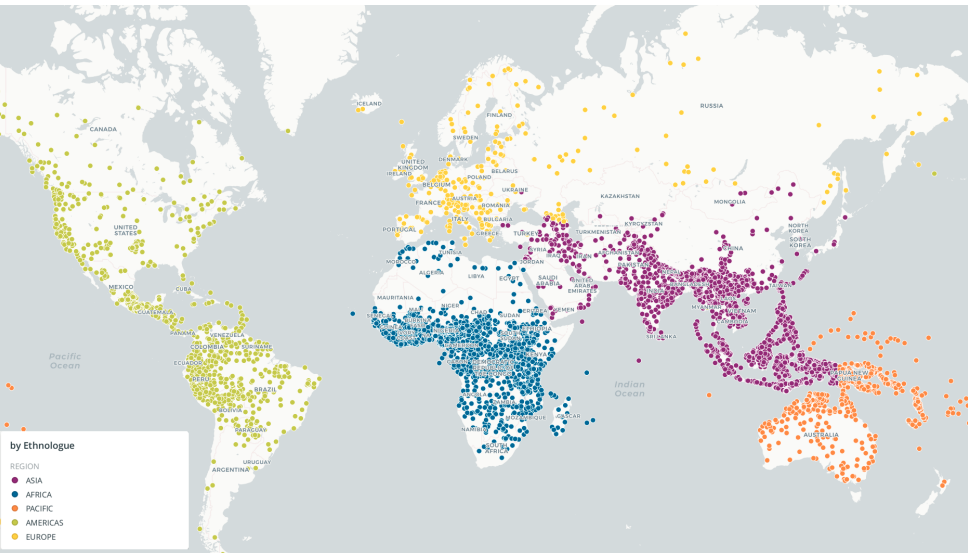
July 10, 2023



Context

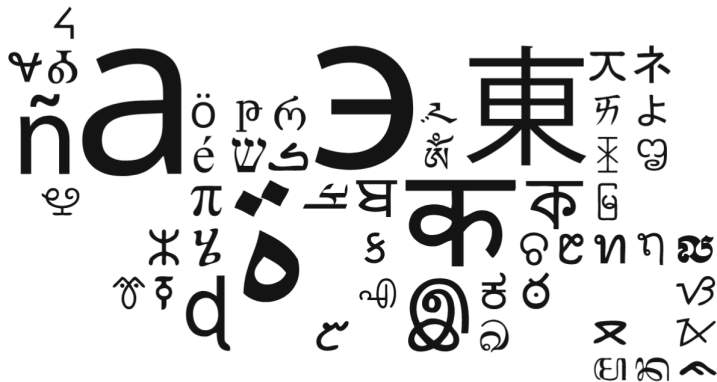
Context: Languages and Writing Systems

- More than 7,000 “languages” are spoken (Ethnologue, 2023).



Context: Languages and Writing Systems

- More than 7,000 “languages” are spoken (Ethnologue, 2023).
- Almost 300 writing systems exist (and many adopted ones)



Context: Languages and Writing Systems

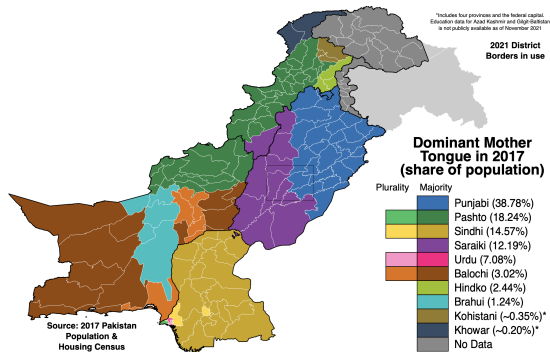
- More than 7,000 “languages” are spoken (Ethnologue, 2023).
- Almost 300 writing systems exist (and many adopted ones)
- Less than 4,000 languages have a written form



Context: Language Communities

Most countries are multi-lingual, but not all officially!

- Pakistan:
→ Urdu and English



Context: Language Communities

Most countries are multi-lingual, but not all officially!

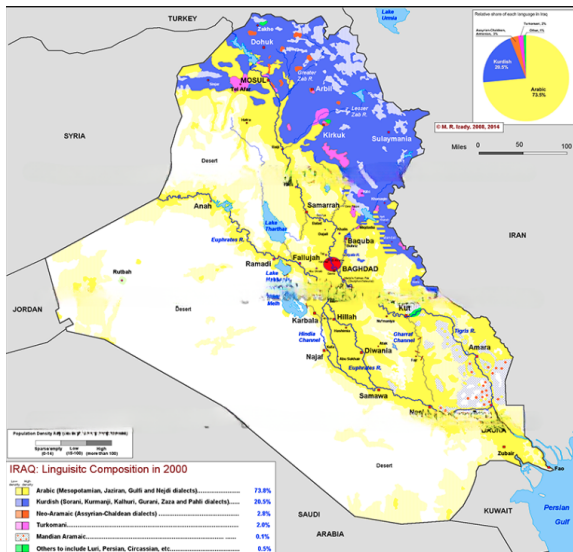
- Pakistan:
→ Urdu and English
- India:
→ Hindi, Kashmiri,
Sindhi and 20 more



Context: Language Communities

Most countries are multi-lingual, but not all officially!

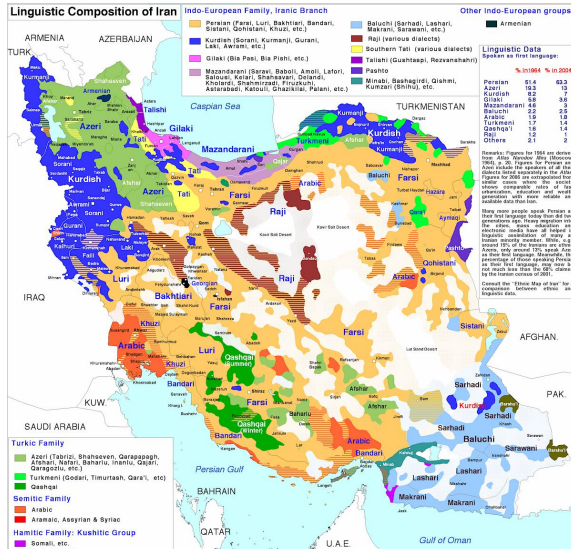
- Pakistan:
→ Urdu and English
- India:
→ Hindi, Kashmiri, Sindhvi and 20 more
- Iraq:
→ Arabic and Kurdish



Context: Language Communities

Most countries are multi-lingual, but not all officially!

- Pakistan:
→ Urdu and English
- India:
→ Hindi, Kashmiri, Sindhi and 20 more
- Iraq:
→ Arabic and Kurdish
- Iran:
→ Persian!



Unconventional Writing

Unconventional Writing

Unconventional writing

Unconventional writing refers to the usage of the **script of another language**, presumably that of a dominant language.

- Unsystematic writing

Unconventional Writing

Unconventional writing

Unconventional writing refers to the usage of the **script of another language**, presumably that of a dominant language.

- Unsystematic writing
- Not necessarily complying with orthography

Unconventional Writing

Unconventional writing

Unconventional writing refers to the usage of the **script of another language**, presumably that of a dominant language.

- Unsystematic writing
- Not necessarily complying with orthography
- Impact of donor language on code switching

Unconventional Writing

Unconventional writing

Unconventional writing refers to the usage of the **script of another language**, presumably that of a dominant language.

- Unsystematic writing
- Not necessarily complying with orthography
- Impact of donor language on code switching
- No specific rule for mapping graphemes-phonemes

Unconventional Writing

Unconventional writing

Unconventional writing refers to the usage of the **script of another language**, presumably that of a dominant language.

- Unsystematic writing
- Not necessarily complying with orthography
- Impact of donor language on code switching
- No specific rule for mapping graphemes-phonemes

A few examples:

Unconventional Writing

Unconventional writing

Unconventional writing refers to the usage of the **script of another language**, presumably that of a dominant language.

- Unsystematic writing
- Not necessarily complying with orthography
- Impact of donor language on code switching
- No specific rule for mapping graphemes-phonemes

A few examples:

- *ana raye7 el gam3a el sa3a 3 el 3asr.* (Arabic in Latin script, aka Arabizi)

Unconventional Writing

Unconventional writing

Unconventional writing refers to the usage of the **script of another language**, presumably that of a dominant language.

- Unsystematic writing
- Not necessarily complying with orthography
- Impact of donor language on code switching
- No specific rule for mapping graphemes-phonemes

A few examples:

- *ana raye7 el gam3a el sa3a 3 el 3asr.* (Arabic in Latin script, aka Arabizi)
- *Tora ti na sou po* (Greek in Latin, aka Greeklish)

Unconventional Writing

Unconventional writing

Unconventional writing refers to the usage of the **script of another language**, presumably that of a dominant language.

- Unsystematic writing
- Not necessarily complying with orthography
- Impact of donor language on code switching
- No specific rule for mapping graphemes-phonemes

A few examples:

- *ana raye7 el gam3a el sa3a 3 el 3asr.* (Arabic in Latin script, aka Arabizi)
- *Tora ti na sou po* (Greek in Latin, aka Greeklish)
- *mer6 pr tn mess pr mn anif* (French SMS language)

Unconventional Writing: Perso-Arabic scripts



Unconventional Writing: Perso-Arabic scripts

- Traditionally used for Arabic
- Used for over a millennium
- A *Reichssprache* for centuries

Arabic



أ ب ت ج
ح خ د ذ ر ز
س ش ع غ
ف ق ل م ن
ه و ي
آ إ أ ا ث و
ص ض ط
ي ة

Unconventional Writing: the Pandora's Box



Script Normalization

Script Normalization

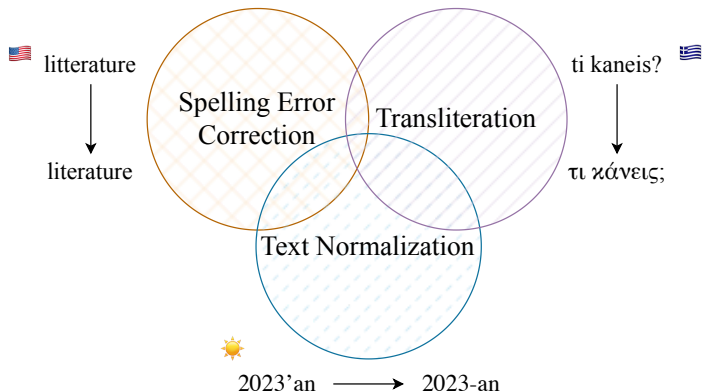
Script Normalization

Normalization of a text written in an unconventional script based on the conventional script and orthography

Script Normalization

Script Normalization

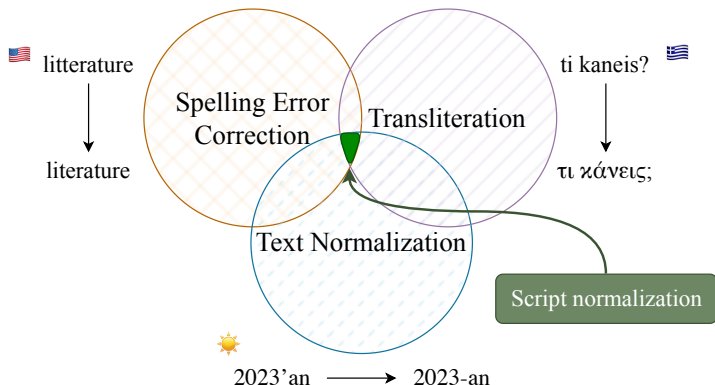
Normalization of a text written in an unconventional script based on the conventional script and orthography



Script Normalization

Script Normalization

Normalization of a text written in an unconventional script based on the conventional script and orthography



Script Normalization: Approach

1 Data collection – Not easy!

Language	639-3	WP	script type	diacritics	ZWNJ	Dominant
Azeri	azb	azb	Abjad	✓	✓	Persian
Turkish						
Kashmiri	kas	ks	Alphabet	✓	✗	Urdu
Gilaki	glk	glk	Abjad	✓	✓	Persian
Gorani	hac	-	Alphabet	✗	✗	Persian, Arabic, Sorani
Kurmanji	kmr	-	Alphabet	✗	✗	Persian, Arabic
Sorani	ckb	ckb	Alphabet	✗	✗	Persian, Arabic
Mazanderani	mzn	mzn	Abjad	✓	✓	Persian
Sindhi	snd	sd	Abjad	✓	✗	Urdu
Persian	fas	fa	Abjad	✓	✓	-
Arabic	arb	ar	Abjad	✓	✗	-
Urdu	urd	ur	Abjad	✓	✓	-

Script Normalization: Approach

- 1 Data collection – Not easy!
- 2 Script mapping
 - Common characters
 - Visual resemblance
 - Orthographic rules

Language	Unconventional script	Source	Target
Azeri Turkish	Persian	چ	چ
Sorani	Arabic	ز	ذ / ض / ظ / ز
Kashmiri	Urdu	اُ	اُ / ا
Sindhi	Urdu	ي	ے / ي / ی

Script Normalization: Approach

- 1 Data collection – Not easy!
- 2 Script mapping
 - Common characters
 - Visual resemblance
 - Orthographic rules
- 3 Character-alignment matrix
→ sequence alignment based on dictionaries

		Sorani to Arabic
▼ ئ:		
	ـ:	0.9829
	ئ:	1.0066
	ا:	1
▼ ا:		
	ا:	1.9559
▼ ب:		
	ب:	1.999
▼ د:		
	د:	2
▼ ی:		
	ي:	1.9222000000000001
▼ ت:		
	ت:	1.7437
	ط:	1.1711

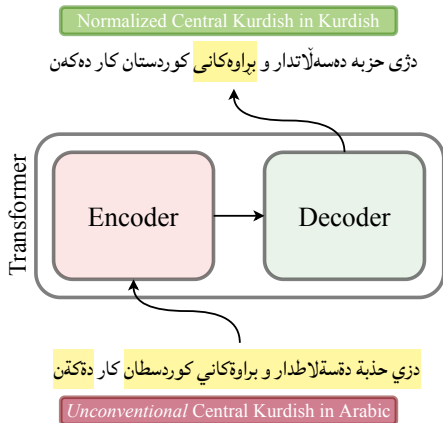
Script Normalization: Approach

- 1 Data collection – Not easy!
- 2 Script mapping
 - Common characters
 - Visual resemblance
 - Orthographic rules
- 3 Character-alignment matrix
 - sequence alignment based on dictionaries
- 4 Synthetic data generation
 - randomly generate pairs
 - inject noise

Noise %	Sentence
Clean	دووهمین پێشانگه‌ها فوتۆگرافه‌رێن کورد ل به‌لجیکا Second Kurdish photographers' exhibition in Belgium
20	دووهمین پێشانگه‌ها فوتۆگرافه‌رێن کورد ل به‌لجیکا
40	دووه‌مین پێشانگه‌ها فطگرافه‌رن کورد ل به‌لجیکا
60	دووه‌مین پێشانگه‌ها فوتۆگرافه‌رن کورد ل به‌لجیکا
80	دووه‌مین پێشانگه‌ها فوتۆگرافه‌رین کورد ل به‌لجیکا
100	دووهمین پێشانگه‌ها فوتۆگرافه‌رین کورد ل به‌لجیکا

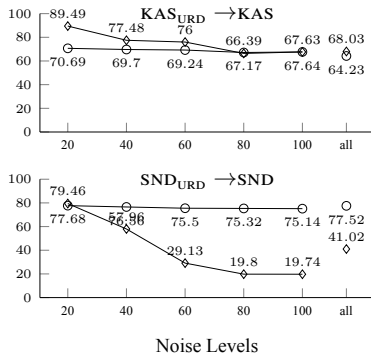
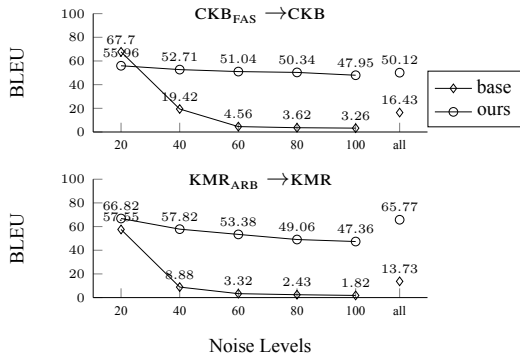
Script Normalization: Approach

- 1 Data collection – Not easy!
- 2 Script mapping
 - Common characters
 - Visual resemblance
 - Orthographic rules
- 3 Character-alignment matrix
 - sequence alignment based on dictionaries
- 4 Synthetic data generation
 - randomly generate pairs
 - inject noise
- 5 Model
 - encoder-decoder with self-attention



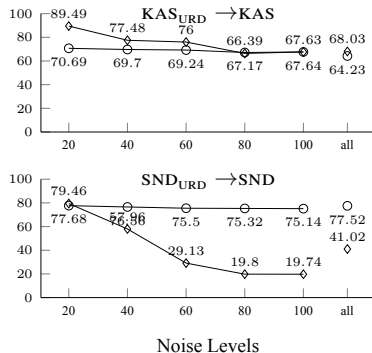
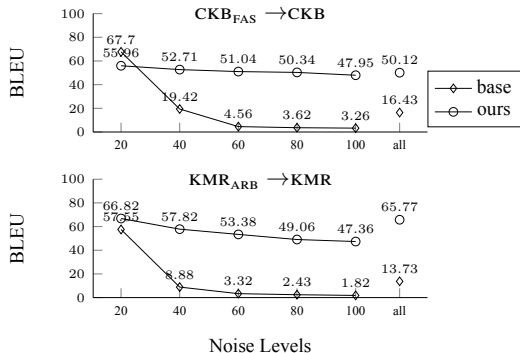
Script Normalization: Intrinsic Experiments

- Baseline: a naive “copy” system



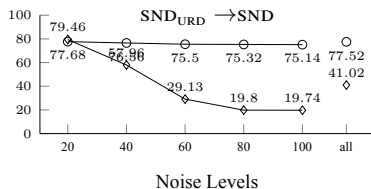
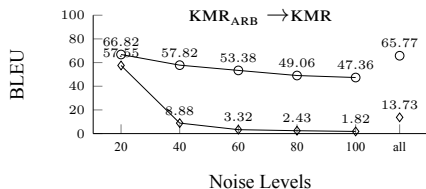
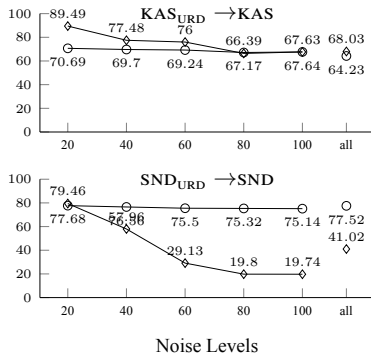
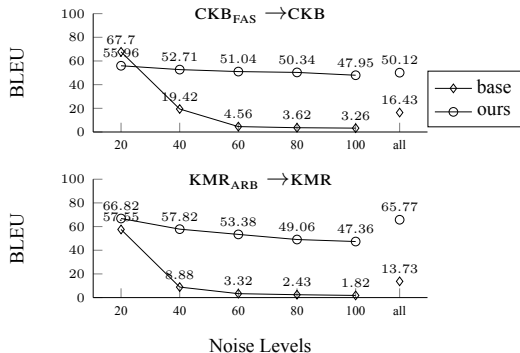
Script Normalization: Intrinsic Experiments

- Baseline: a naive “copy” system
- Ours: trained models on different levels of noise



Script Normalization: Intrinsic Experiments

- Baseline: a naive “copy” system
- Ours: trained models on different levels of noise
- **Our models dramatically improve over the baseline**



Script Normalization: Intrinsic Experiments

- Baseline: a naive “copy” system
- Ours: trained models on different levels of noise
- **Our models dramatically improve over the baseline**
- Including when evaluated on real data

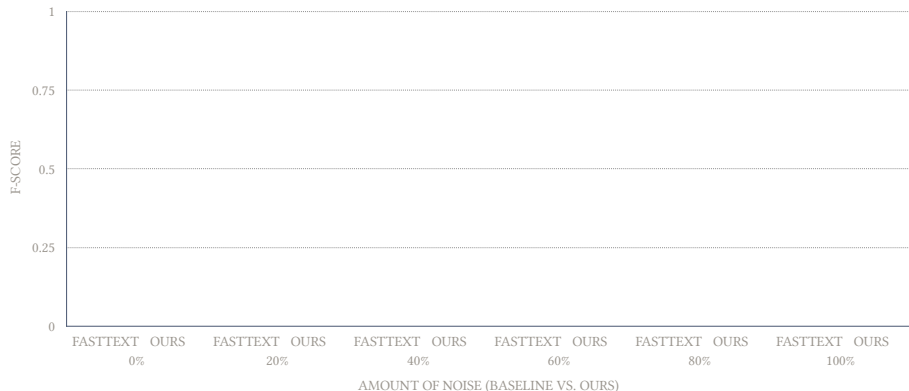
Sorani Eval Set	Original (Unconventional)		Normalized	
	BLEU	chrF	BLEU	chrF
$\text{CKB}_{\text{FAS}} \rightarrow \text{CKB}$	1.2	38.4	20.1	69.6
$\text{CKB}_{\text{ARB}} \rightarrow \text{CKB}$	0.4	19.4	12.7	65.2

① Language identification (LID)

Script Normalization: Extrinsic Experiments

① Language identification (LID)

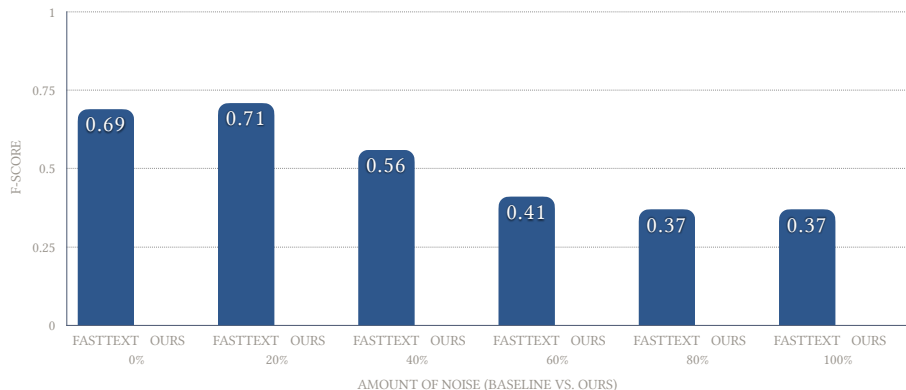
- Compare LID with and without normalization



Script Normalization: Extrinsic Experiments

① Language identification (LID)

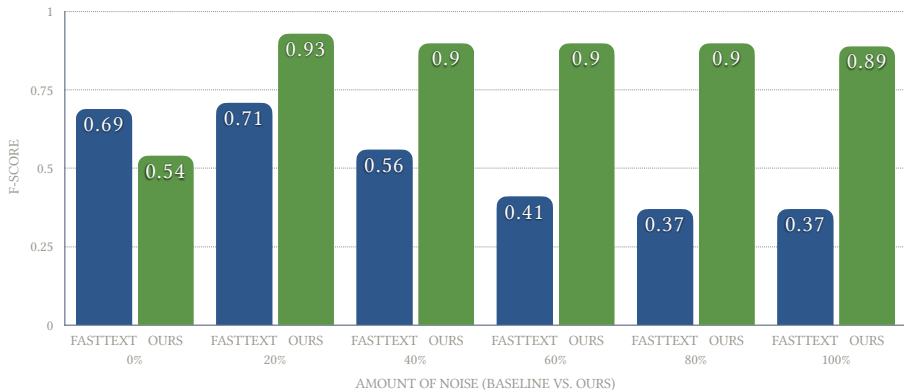
- Compare LID with and without normalization
- Terrible performance by any existing model



Script Normalization: Extrinsic Experiments

① Language identification (LID)

- Compare LID with and without normalization
- Terrible performance by any existing model
- Models trained on normalized datasets improve the F-scores

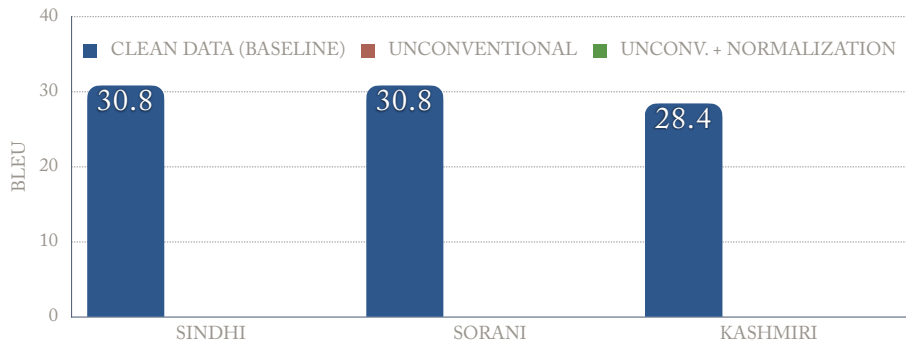


① Language identification (LID)

- Compare LID with and without normalization
- Terrible performance by any existing model
- Models trained on normalized datasets improve the F-scores
- Closely-related languages (scripts) are confused!

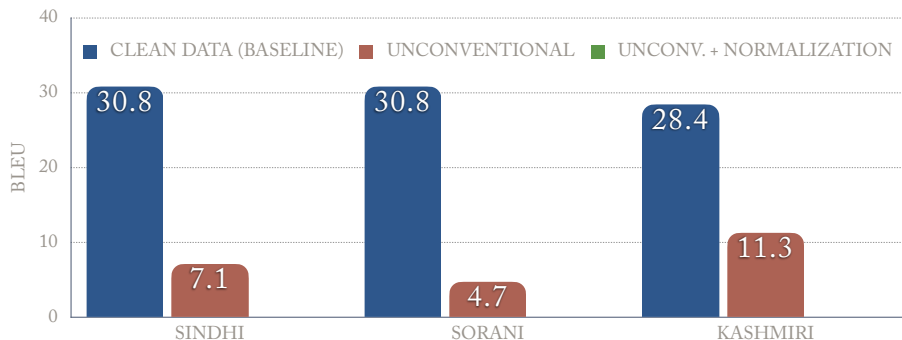
Script Normalization: Extrinsic Experiments

- 1 Language identification (LID)
- 2 **Machine Translation (MT)**
 - Evaluate MT with and without normalization



Script Normalization: Extrinsic Experiments

- 1 Language identification (LID)
- 2 **Machine Translation (MT)**
 - Evaluate MT with and without normalization
 - Terrible performance on noisy data (NLLB as baseline)

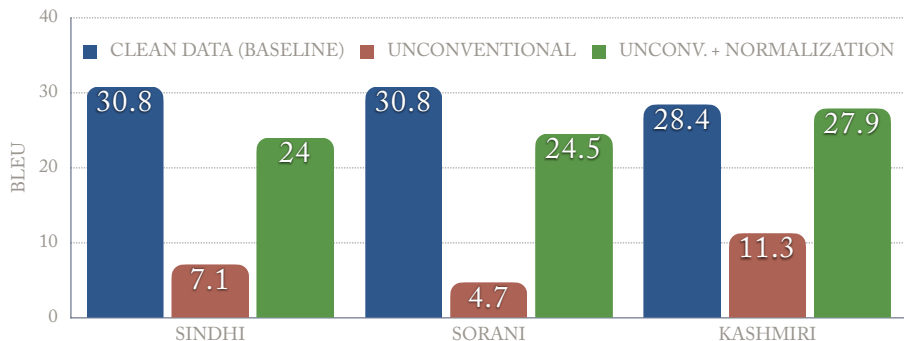


Script Normalization: Extrinsic Experiments

1 Language identification (LID)

2 Machine Translation (MT)

- Evaluate MT with and without normalization
- Terrible performance on noisy data (NLLB as baseline)
- Models trained on normalized datasets improve the F-scores



Conclusion

Key Takeaways

- 1 Unconventional writing is more widespread than thought of.

Key Takeaways

- ① Unconventional writing is more widespread than thought of.
- ② It is an open problem and non-trivial to solve.

Key Takeaways

- 1 Unconventional writing is more widespread than thought of.
- 2 It is an open problem and non-trivial to solve.
- 3 It negatively affects NLP for low-resourced languages.

Key Takeaways

- ① Unconventional writing is more widespread than thought of.
- ② It is an open problem and non-trivial to solve.
- ③ It negatively affects NLP for low-resourced languages.
- ④ We can effectively remediate it, *but only to some extent...*

Key Takeaways

- 1 Unconventional writing is more widespread than thought of.
- 2 It is an open problem and non-trivial to solve.
- 3 It negatively affects NLP for low-resourced languages.
- 4 We can effectively remediate it, *but only to some extent...*
- 5 Do we always need to write a language?
 - Multi-modal NLP
 - Multi-lingual NLP
 - Multi-task NLP
 - Better adaptation in NLP

Key Takeaways

- 1 Unconventional writing is more widespread than thought of.
- 2 It is an open problem and non-trivial to solve.
- 3 It negatively affects NLP for low-resourced languages.
- 4 We can effectively remediate it, *but only to some extent...*
- 5 Do we always need to write a language?
 - Multi-modal NLP
 - Multi-lingual NLP
 - Multi-task NLP
 - Better adaptation in NLP
- 6 Models and codes:
<https://github.com/sinaahmadi/ScriptNormalization>

Any questions?

Kiitos धन्यवाद Köszönöm
Rahmat Tak

شكراً 谢谢 Gracias

Спасибо

Adúpé

Hvala

Mulțumesc

Paldies

Bedankt

Dzięki

ขอบคุณ

Obrigado

koe

Shurkan

Danke Diolch

Danko

Daalu

ơn

Cảm

Grazie

Teşekkürler

감사합니다

Tānan

Merci



ευχαριστώ Mahalo

References I

-  Ahmadi, Sina, Milind Agarwal, and Antonios Anastasopoulos (May 2023). “PALI: A Language Identification Benchmark for Perso-Arabic Scripts”. In: *Proceedings of the 10th Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial)*. Dubrovnik, Croatia: The 17th Conference of the European Chapter of the Association for Computational Linguistics.
-  Ahmadi, Sina and Antonios Anastasopoulos (2023). “Script Normalization for Unconventional Writing of Under-Resourced Languages in Bilingual Communities”. In: Toronto, Canada: The 61st Annual Meeting of the Association for Computational Linguistics.
-  Sheyholislami, Jaffer (2012). “Kurdish in Iran: A case of restricted and controlled tolerance”. In: *International Journal of the Sociology of Language* 2012.217, pp. 19–47.