

Is NLP Ready for Standardization?

Lauriane Aufrant

Findings of EMNLP 2022

What is the common point between...



What is the common point between...



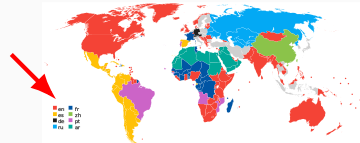
THE
C
PROGRAMMING
LANGUAGE



mp3



```
~~~~~  
!"#$%&'()*+,-./0123456789:;<=>?  
@ABCDEFGHIJKLMNPOQRSTUVWXYZ[\]^_  
`abcdefghijklmnopqrstuvwxyz{|}~`~`  
~~~~~  
¡¢£¥¦§¨ª«¬®¯°±²³´µ¶·¸¹º»¼½¾¿  
ÀÁÂÃÄÅÆÇÈÉÊËÌÍÎÏÐÑÒÓÔÕÖ×ØÙÚÛÜÝÞß  
àáâãäåæçèéêëìíîïðñóôõö÷øùúûüýþ
```



What is the common point between...



ISO/IEC 9899



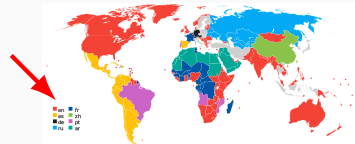
ISO/IEC 11172-3



ISO/IEC 8859-1



```
.....  
!"#$%&'()*+,-./0123456789:;<=>?  
@ABCDEFGHIJKLMNPOQRSTUVWXYZ[\]^_  
`abcdefghijklmnopqrstuvwxyz{|}~  
.....  
¡¢£¥¦§¨ª«¬®¯°±²³´µ¶·¸¹º»¼½¾¿  
ÀÁÂÃÄÅÆÇÈÉÊËÌÍÎÏÐÑÒÓÔÕÖ×ØÙÚÛÜÝÞß  
àáâãäåæçèéêëìíîïðñóôõö÷øùúûüýþ
```

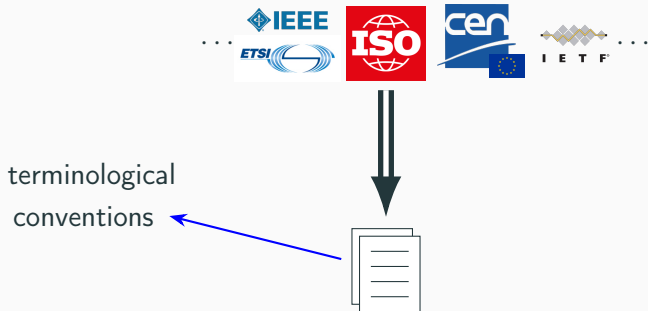


ISO 639-1

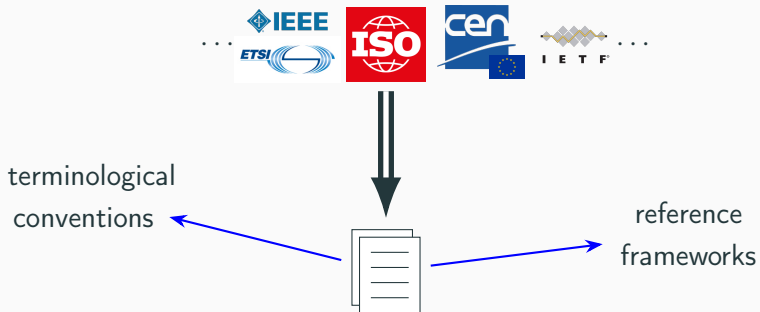
Standards



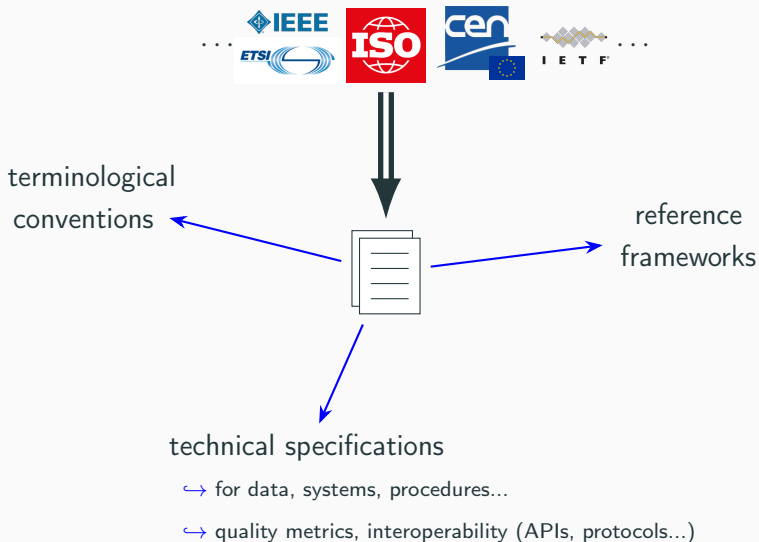
Standards



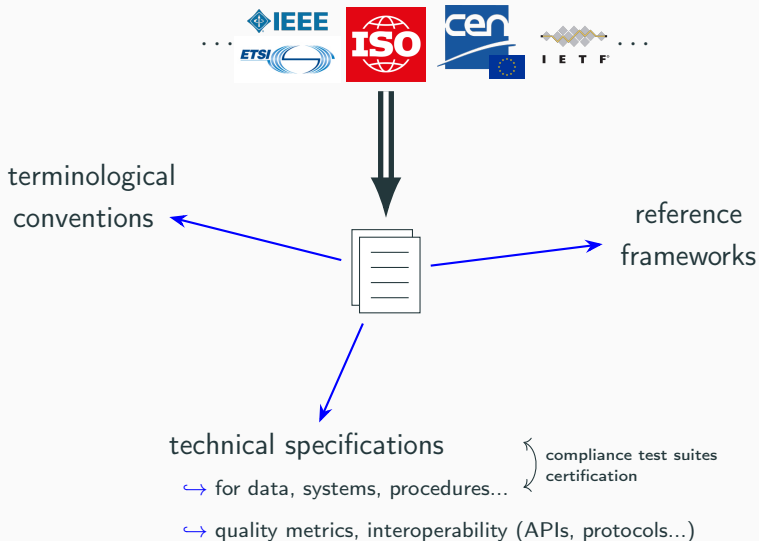
Standards

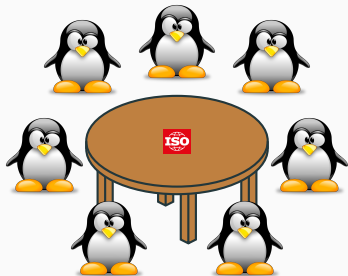


Standards



Standards









And... NLP?

Towards fair and reproducible evaluation: Machine translation & BLEU

$$\text{BLEU} = \exp^{-\max(0, \frac{L}{L_{\text{closest-ref}}} - 1)} \cdot \sqrt[n]{\prod_{k=1}^n \frac{\text{TP}_{n\text{-gram}}(k)}{\text{TP}_{n\text{-gram}}(k) + \text{FP}_{n\text{-gram}}(k)}}$$

Towards fair and reproducible evaluation: Machine translation & BLEU

$$\text{BLEU} = \exp^{-\max\left(0, \frac{L}{L_{\text{closest-ref}}} - 1\right)} \cdot \sqrt[n]{\prod_{k=1}^n \frac{\text{TP}_{n\text{-gram}}(k)}{\text{TP}_{n\text{-gram}}(k) + \text{FP}_{n\text{-gram}}(k)}}$$

A Call for Clarity in Reporting BLEU Scores

Matt Post
Amazon Research
Berlin, Germany

Abstract

The field of machine translation faces an under-recognized problem because of inconsistency in the reporting of scores from its dominant metric. Although people refer to “the” BLEU score, BLEU is in fact a parameterized metric whose values can vary wildly with changes to these parameters. These parameters are often not reported or are hard to find, and consequently, BLEU scores between papers cannot be directly compared. I quantify this variation, finding differences as high as 1.8 between commonly used configura-

hered the field through a decade and a half of quality improvements (Graham et al., 2014).

This is of course not to claim there are no problems with BLEU! Its weaknesses abound, and much has been written about them (cf. Callison-Burch et al. (2006); Reiter (2018)). This paper is not, however, concerned with the shortcomings of BLEU as a proxy for human evaluation of quality; instead, our goal is to bring attention to the relatively narrower problem of the *reporting* of BLEU scores. This problem can be summarized as follows:

- ▶ Tokenization (user’s vs metric-internal)?
 - ▶ N-gram max size?
 - ▶ Number of references?
 - ▶ ...
- ↪ Up to Δ 1.8 BLEU (> gains from BPE!)

Towards fair and reproducible evaluation: Named entity recognition & F1

$$F1 = \frac{2 \cdot P \cdot R}{P + R}, P = \frac{TP}{TP + FP}, R = \frac{TP}{TP + FN}$$

Towards fair and reproducible evaluation: Named entity recognition & F1

$$F1 = \frac{2 \cdot P \cdot R}{P + R}, P = \frac{TP_?}{TP_? + FP_?}, R = \frac{TP_?}{TP_? + FN_?}$$

Towards fair and reproducible evaluation: Named entity recognition & F1

$$F1 = \frac{2 \cdot P \cdot R}{P + R}, P = \frac{TP_?}{TP_? + FP_?}, R = \frac{TP_?}{TP_? + FN_?}$$

Ref: B-PER I-PER 0 B-LOC 0

Pred: B-PER 0 0 B-LOC I-LOC

Towards fair and reproducible evaluation: Named entity recognition & F1

$$F1 = \frac{2 \cdot P \cdot R}{P + R}, \quad P = \frac{TP_?}{TP_? + FP_?}, \quad R = \frac{TP_?}{TP_? + FN_?}$$

Ref: B-PER I-PER 0 B-LOC 0

Pred: B-PER 0 0 B-LOC I-LOC
 TP FN TN TP FP

↪ 67 F1 on token-level (BIO) labels

↪ 0 F1 on predicted chunks

Towards fair and reproducible evaluation: Named entity recognition & F1

$$F1 = \frac{2 \cdot P \cdot R}{P + R}, \quad P = \frac{TP_?}{TP_? + FP_?}, \quad R = \frac{TP_?}{TP_? + FN_?}$$

Ref: B-PER I-PER 0 B-LOC 0
Pred: B-PER 0 0 B-LOC I-LOC
 TP FN TN TP FP

- ↪ 67 F1 on token-level (BIO) labels
- ↪ 0 F1 on predicted chunks

- ▶ Typed matching vs boundaries only
- ▶ 2 correct boundaries vs 1 correct vs at least 1 correct token
- ▶ Non 1-1 matches (ref B-PER I-PER vs pred B-PER B-PER)

Towards fair and reproducible evaluation: Named entity recognition & F1

$$F1 = \frac{2 \cdot P \cdot R}{P + R}, \quad P = \frac{TP_?}{TP_? + FP_?}, \quad R = \frac{TP_?}{TP_? + FN_?}$$

Ref: B-PER I-PER 0 B-LOC 0
Pred: B-PER 0 0 B-LOC I-LOC
 TP FN TN TP FP

- ▶ Typed matching vs boundaries only
- ▶ 2 correct boundaries vs 1 correct vs at least 1 correct token
- ▶ Non 1-1 matches (ref B-PER I-PER vs pred B-PER B-PER)

↪ 67 F1 on token-level (BIO) labels
↪ 0 F1 on predicted chunks

And what about...

- ? invalid sequences
B-PER I-LOC 0 I-PER 0
- ? micro vs macro-F1

- ? ignoring Other/MISC classes
- ? dependency of evaluation tools to the encoding scheme (BIO, BILUO...)

Consistent data and annotations

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE

Consistent data and annotations

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE

Consistent data and annotations

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE

Let's write a
conversion tool!



Consistent data and annotations

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE
BILUO	B-PER	I-PER	L-PER	O	O	U-LOC	U-DATE

Let's write a
conversion tool!



Consistent data and annotations

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE
BILUO	B-PER	I-PER	L-PER	O	O	U-LOC	U-DATE
BIOES	B-PER	I-PER	E-PER	O	O	S-LOC	S-DATE

Let's write a
conversion tool!



Consistent data and annotations

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE
BILUO	B-PER	I-PER	L-PER	O	O	U-LOC	U-DATE
BIOES	B-PER	I-PER	E-PER	O	O	S-LOC	S-DATE
BMEOW	B-PER	M-PER	E-PER	O	O	W-LOC	W-DATE

Let's write a
conversion tool!



Consistent data and annotations

SGML

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE
BILUO	B-PER	I-PER	L-PER	O	O	U-LOC	U-DATE
BIOES	B-PER	I-PER	E-PER	O	O	S-LOC	S-DATE
BMEOW	B-PER	M-PER	E-PER	O	O	W-LOC	W-DATE

Let's write a
conversion tool!



Consistent data and annotations

SGML

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE
BILUO	B-PER	I-PER	L-PER	O	O	U-LOC	U-DATE
BIOES	B-PER	I-PER	E-PER	O	O	S-LOC	S-DATE
BMEOW	B-PER	M-PER	E-PER	O	O	W-LOC	W-DATE

JSON

Let's write a
conversion tool!



Consistent data and annotations

SGML

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE
BILUO	B-PER	I-PER	L-PER	O	O	U-LOC	U-DATE
BIOES	B-PER	I-PER	E-PER	O	O	S-LOC	S-DATE
BMEOW	B-PER	M-PER	E-PER	O	O	W-LOC	W-DATE

tabular

JSON

Let's write a
conversion tool!



Consistent data and annotations

SGML

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE
BILUO	B-PER	I-PER	L-PER	O	O	U-LOC	U-DATE
BIOES	B-PER	I-PER	E-PER	O	O	S-LOC	S-DATE
BMEOW	B-PER	M-PER	E-PER	O	O	W-LOC	W-DATE

tabular

brat standoff

JSON

Let's write a
conversion tool!



Consistent data and annotations

TMX

SGML

brat standoff

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE
BILUO	B-PER	I-PER	L-PER	O	O	U-LOC	U-DATE
BIOES	B-PER	I-PER	E-PER	O	O	S-LOC	S-DATE
BMEOW	B-PER	M-PER	E-PER	O	O	W-LOC	W-DATE

tabular

JSON

Let's write a
conversion tool!



Consistent data and annotations

TMX

SGML

brat standoff

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE
BILUO	B-PER	I-PER	L-PER	O	O	U-LOC	U-DATE
BIOES	B-PER	I-PER	E-PER	O	O	S-LOC	S-DATE
BMEOW	B-PER	M-PER	E-PER	O	O	W-LOC	W-DATE

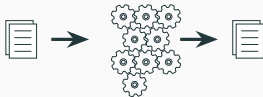
tabular

JSON

Let's write a
conversion tool!



XML-XCES



Consistent data and annotations

TMX

SGML

brat standoff

	John	D.	Smith	arrived	in	Toronto	yesterday
BIO	B-PER	I-PER	I-PER	O	O	B-LOC	B-DATE
IOB	I-PER	I-PER	I-PER	O	O	I-LOC	I-DATE
BILUO	B-PER	I-PER	L-PER	O	O	U-LOC	U-DATE
BIOES	B-PER	I-PER	E-PER	O	O	S-LOC	S-DATE
BMEOW	B-PER	M-PER	E-PER	O	O	W-LOC	W-DATE

tabular

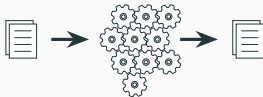
bitext

JSON

Let's write a
conversion tool!



XML-XCES



Consistent data and annotations

Universal Dependencies


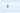




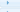
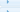
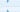
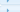
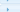

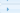

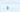
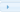
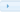
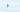
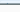
Universal Dependencies (UD) is a framework for consistent annotation of grammar (parts of speech, morphological features, and syntactic dependencies) across different human languages. UD is an open community effort with over 300 contributors producing nearly 200 treebanks in over 100 languages. If you're new to UD, you should start by reading the first part of the Short Introduction and then browsing the annotation guidelines.

- [Short introduction to UD](#)
- [UD annotation guidelines](#)
- More information on UD:
 - [How to contribute to UD](#)
 - [Tools for working with UD](#)
 - [Changes to the UD guidelines](#)
 - [UD-related events](#)
- Query UD treebanks online:
 - [SETS treebank search](#) maintained by the University of Turku
 - [PML Tree Query](#) maintained by the Charles University in Prague
 - [Kortext](#) maintained by the Charles University in Prague
 - [Grew-match](#) maintained by Inria in Nancy
 - [INESS](#) maintained by the University of Bergen
- [Download UD treebanks](#)

If you want to receive news about Universal Dependencies, you can subscribe to the [UD mailing list](#). If you want to discuss individual annotation questions, use the [Github issue tracker](#).

Current UD Languages

Information about language families (and genera for families with multiple branches) is mostly taken from [WALS Online](#) (IE = Indo-European).

+		Afrikaans	1	49K	🇳🇵	IE, Germanic
+		Akkadian	2	25K	🇮🇶	Afro-Asiatic, Semitic
+		Akuntsu	1	1K	🇳🇮	Tupian, Tupari
+		Albanian	1	<1K	🇲🇰	IE, Albanian
+		Amharic	1	10K	🇪🇹	Afro-Asiatic, Semitic
+		Ancient Greek	2	416K	🇬🇷	IE, Greek
+		Ancient Hebrew	1	39K	🇮🇱	Afro-Asiatic, Semitic
+		Apurina	1	<1K	🇬🇹	Arawakan
+		Arabic	3	1,042K	🇸🇦	Afro-Asiatic, Semitic
+		Armenian	2	94K	🇦🇲	IE, Armenian
+		Assyrian	1	<1K	🇮🇶	Afro-Asiatic, Semitic
+		Bambara	1	13K	🇲🇱	Mande
+		Basque	1	121K	🇫🇷	Basque
+		Beja	1	<1K	🇱🇾	Afro-Asiatic, Cushitic
+		Belarusian	1	305K	🇧🇪	IE, Slavic
+		Bengali	1	<1K	🇬🇧	IE, Indic
+		Bhojpuri	1	6K	🇮🇳	IE, Indic
+		Breton	1	10K	🇫🇷	IE, Celtic
+		Bulgarian	1	156K	🇧🇬	IE, Slavic

Working on the same task?

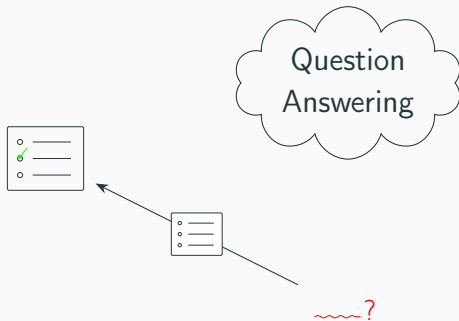


Working on the same task?



~~~~~?

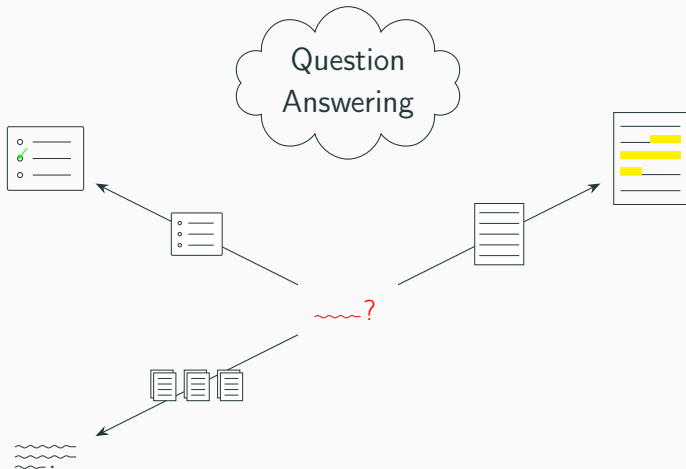
## Working on the same task?



## Working on the same task?

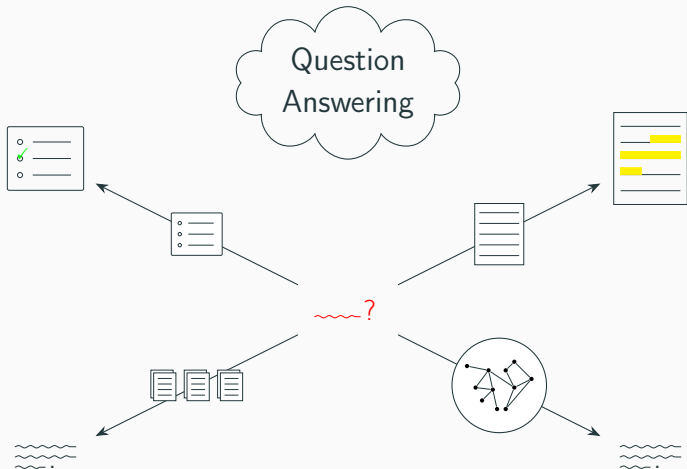


## Working on the same task?



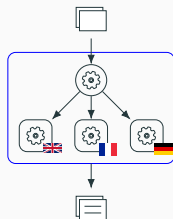


## Working on the same task?

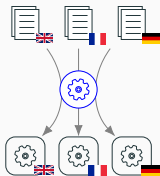
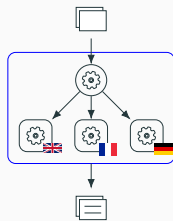




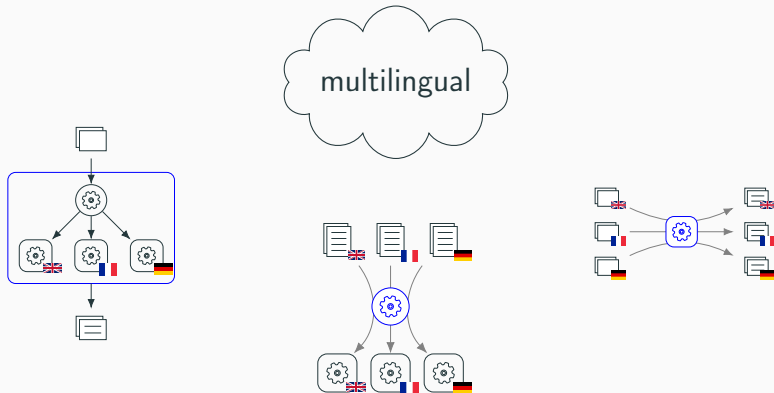
## Reflecting upon our work



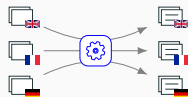
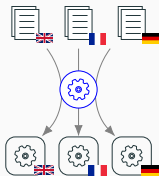
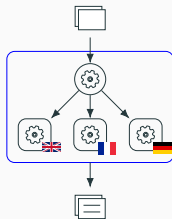
# Reflecting upon our work



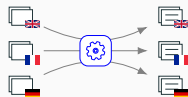
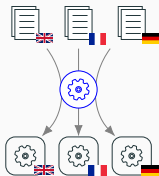
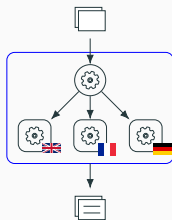
# Reflecting upon our work



# Reflecting upon our work



# Reflecting upon our work



**Standards to the rescue!**



# Benefits and caveats of formal standards

85 F1??  
80 F1??



# Benefits and caveats of formal standards

85 F1??  
80 F1??



# Benefits and caveats of formal standards

85 F1??  
80 F1??



# Benefits and caveats of formal standards

85 F1??  
80 F1??



# Benefits and caveats of formal standards

85 F1??  
80 F1??



## Key to success: finding an appropriate positioning

- ▶ standardize what is consensual, not more, not less
  - ↪ a forum to build further consensus?
- ▶ formalizing/referencing vs prescribing
  - ↪ preserve research freedom in a fast moving field
- ▶ keep up with societal & technological evolutions
  - ↪ standards need to be revised every few years
- ▶ build upon existing work
  - ↪ promoting existing guidelines rather than recrafting

## NLP standardization: current state



TC 37



*“Language and Terminology”*

# NLP standardization: current state



TC 37



*“Language and Terminology”*



JTC 1/SC 42



*“Artificial intelligence”*



# NLP standardization: current state



TC 37



*“Language and Terminology”*



JTC 1/SC 42



*“Artificial intelligence”*



# NLP standardization: current state



TC 37



*“Language and Terminology”*



JTC 1/SC 42



*“Artificial intelligence”*



JTC 21



*“Artificial intelligence”*

– with NLP projects

# NLP standardization: current state

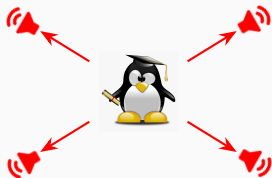
 TC 37

*"Language and Terminology"*



 JTC 1/SC 42

*"Artificial intelligence"*



  JTC 21

*"Artificial intelligence"*

– with NLP projects



## Take-home messages

- Purposes of standards: **formalizing** common knowledge and guidelines, ensuring **consensus** among experts
- Standardization issues exist **all across NLP**
  - ↳ on terminology, concepts, evaluation metrics, data formats...
- A useful contribution to **society at large**
  - ↳ expected benefits for researchers, industry, consumers, regulators
- Some good work initiated: guidelines, workshops, projects...
  - ↳ next step now is to give them **formal** existence
- Existing groups, to be expanded for better coverage of NLP
  - ↳ **researchers needed**: to share your expertise + voice any concern

*Thank you for watching!*

For any question, comment, suggestion: [first.last@inria.fr](mailto:first.last@inria.fr)