

# Towards Equitable and Culturally Adapted Multilingual Dialog Systems

Ivan Vulić

*LTL, University of Cambridge*



UNIVERSITY OF  
CAMBRIDGE

UNLP Workshop  
(Online)  
May 25 2024

# Why Multilingual NLP?



*“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 لطفاً

*“¿Puedo tomar un taxi hasta el aeropuerto?”*

*“Molim Vas jedno vozilo do Autobusnog”*

هل يمكنني 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...**  
...and reaching **more users and customers...**

# Why Multilingual NLP and Multilingual Dialogue?

...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

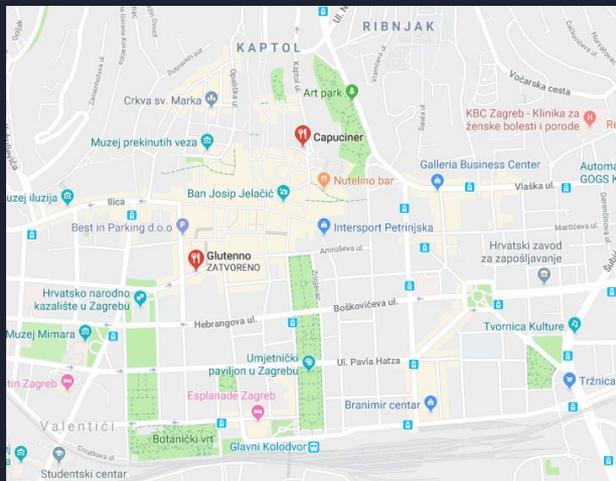
“95% of all languages in use today will never gain traction online” (Andras Kornai)

“The limits of my language *online* mean the limits of my world?”

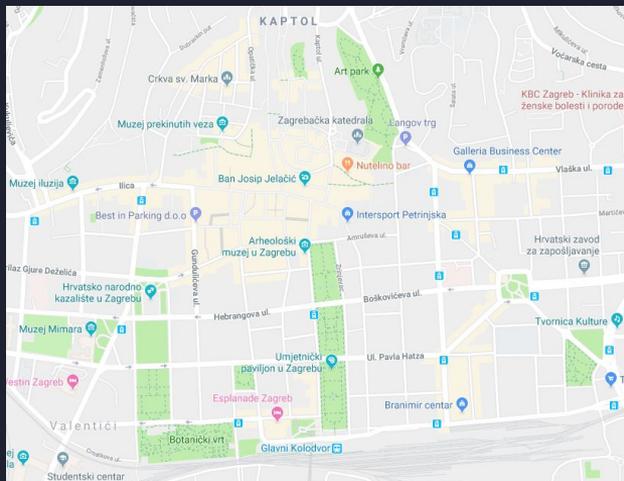
# Why Multilingual NLP?

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

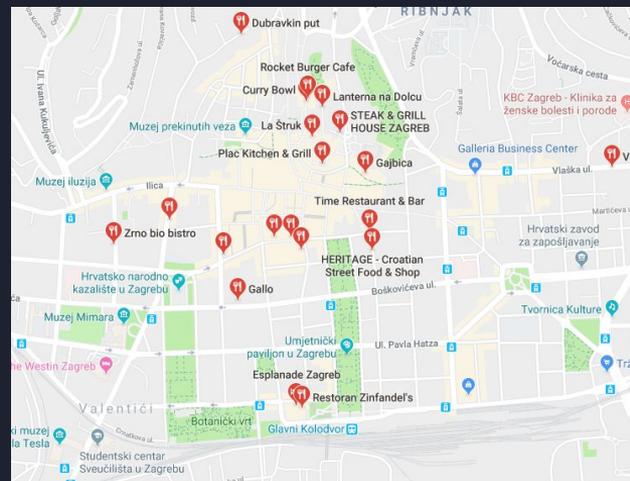
We're in Zagreb searching for...



...éttermek (HU)



...jate (EU)



...restaurants (EN)

# English Dialogue Systems

A successful conversational agent must (at least implicitly) 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 Dialogue Systems?

According to Ethnologue there are **7,000+** living languages

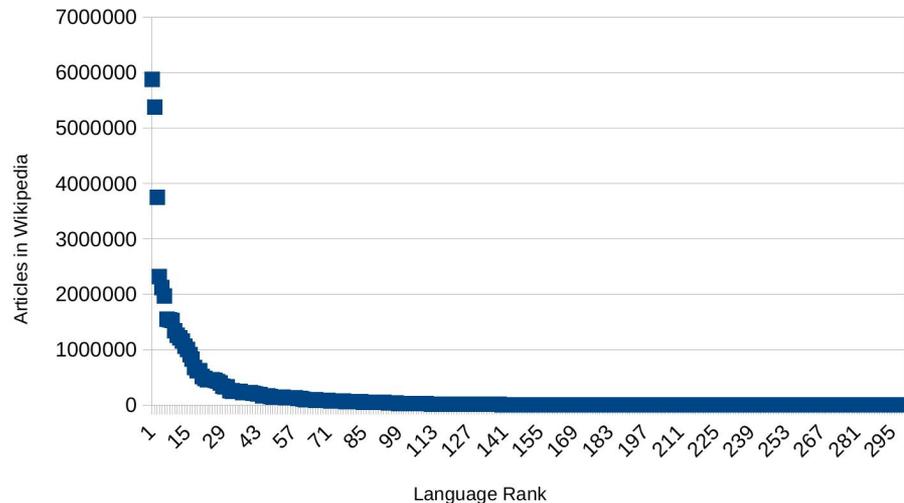
What about **language varieties** and **dialects**?

What about “**social media**” languages and **slang**?

What about “**all those domains**”?

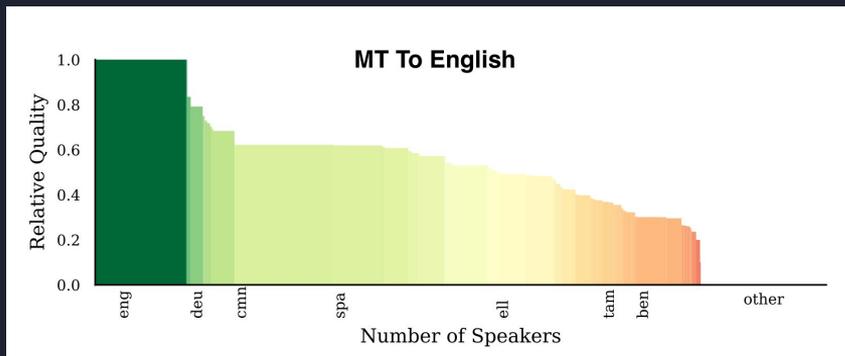


## The Long Tail of Data

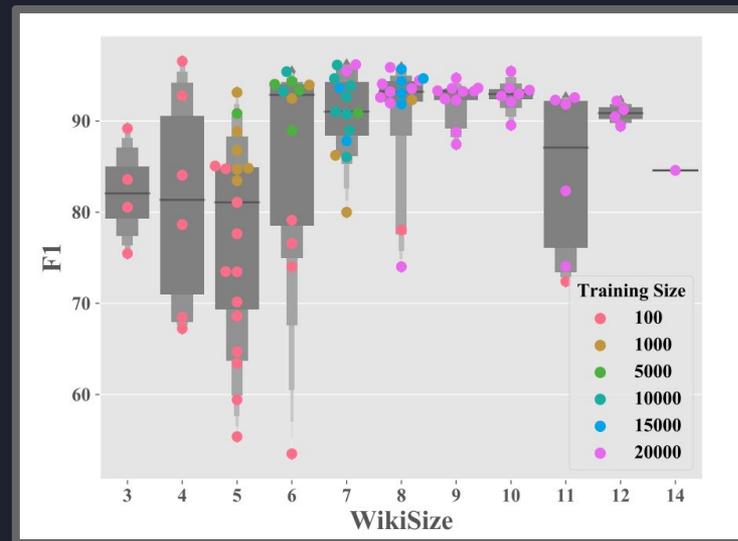


Even getting “raw” unannotated data is problematic for many languages...

# The Long Tail of Data Means Inequality

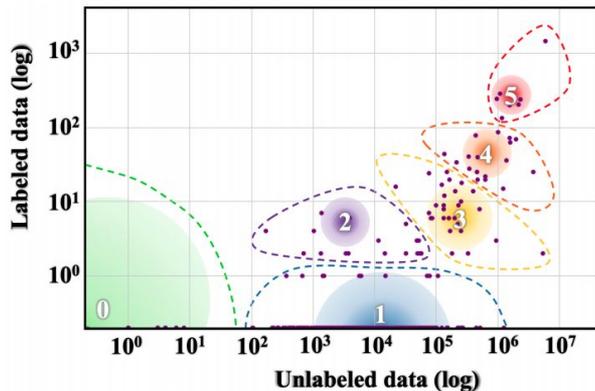


**MT for major versus minor languages**  
(Blasi et al, 2022)



**NER with mBERT on 99 languages**  
(Wu and Dredze, 2020)

# Are All Languages Created Equal?



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?**

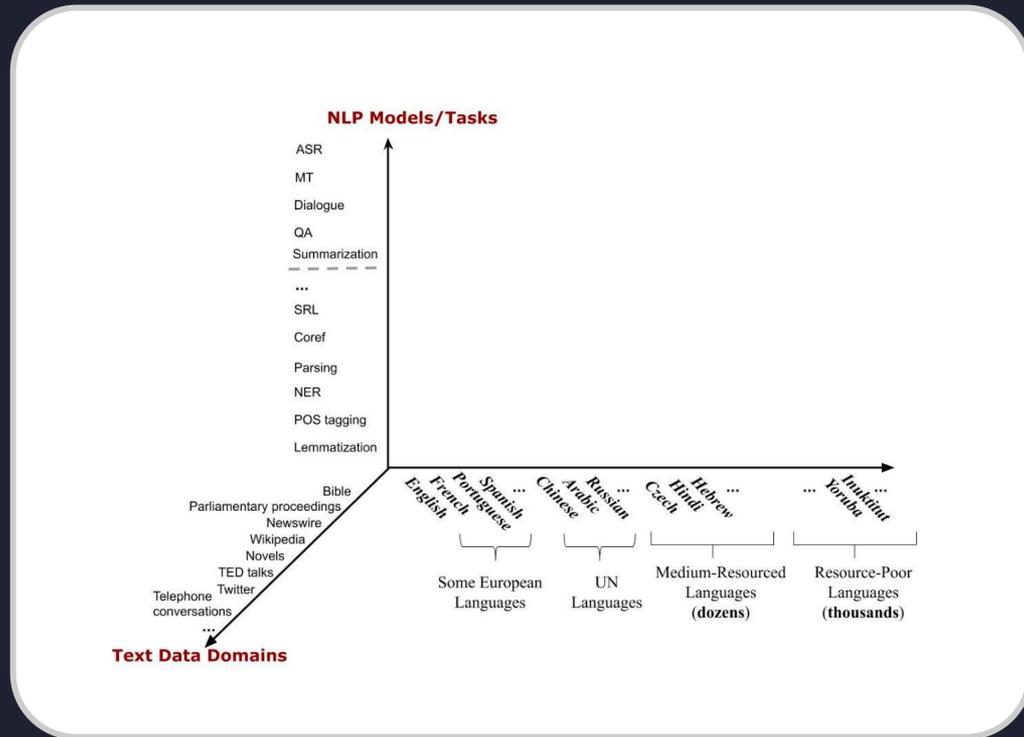
**Can we at least try to ‘approximate’ equality?**

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%

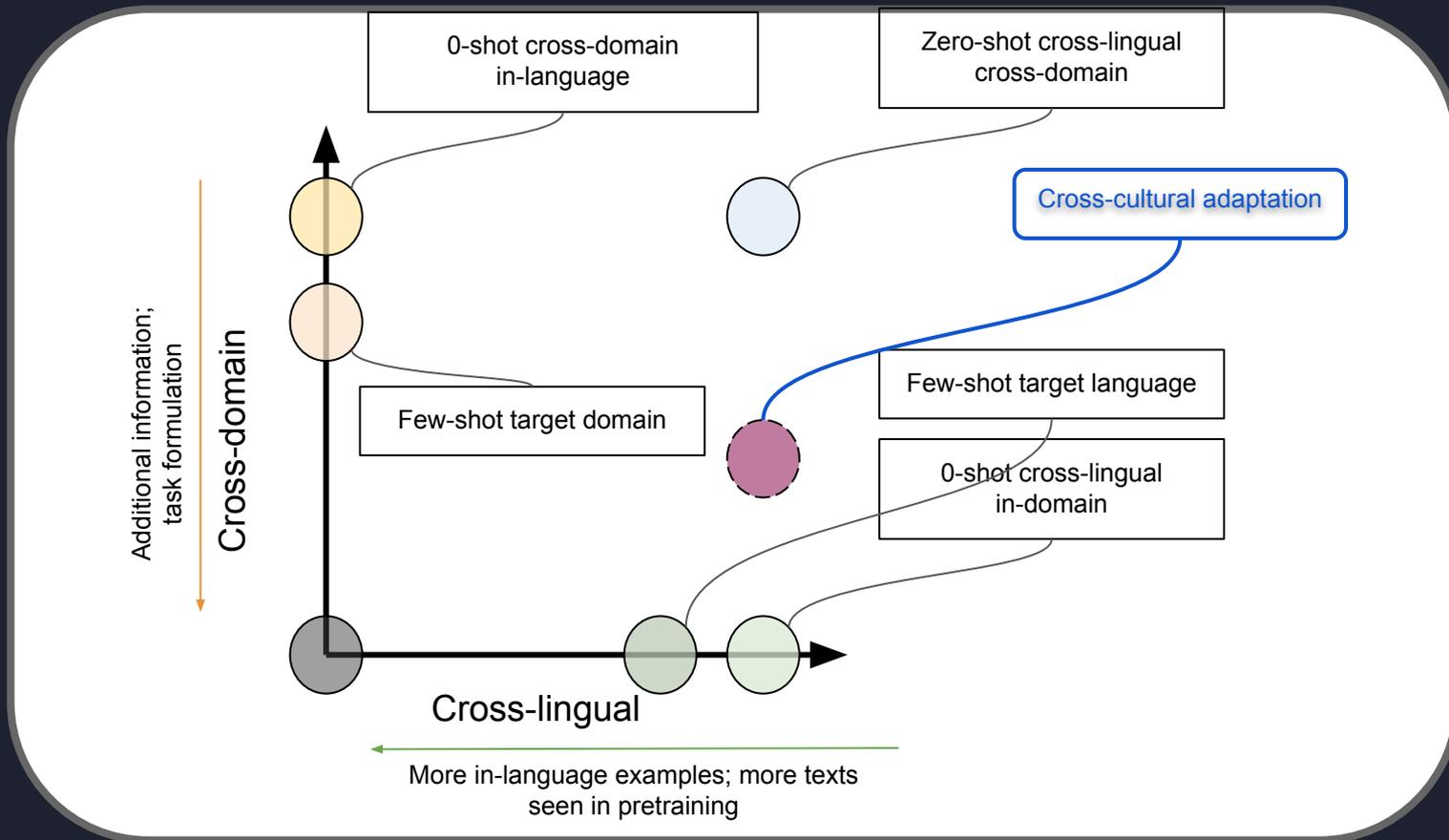
# What Can We Do in Multilingual and Multi-Domain Setups?

- Many NLP tasks and domains **share common knowledge about language** (e.g. linguistic representations, structural similarities)
- Languages and domains **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)

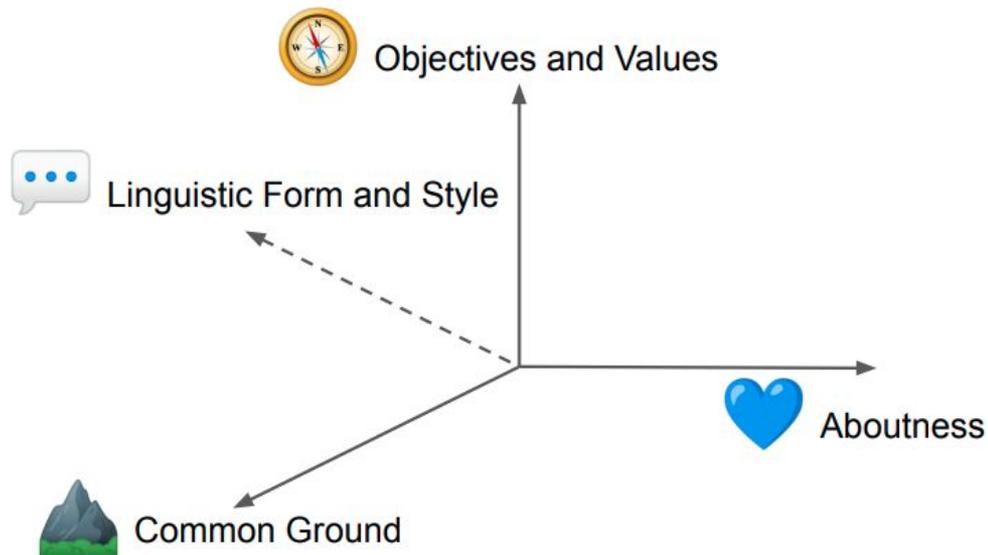
Image courtesy of Yulia Tsvetkov



# Towards Multilingual and Multi-Domain Systems?



# Why Cultural Adaptation?



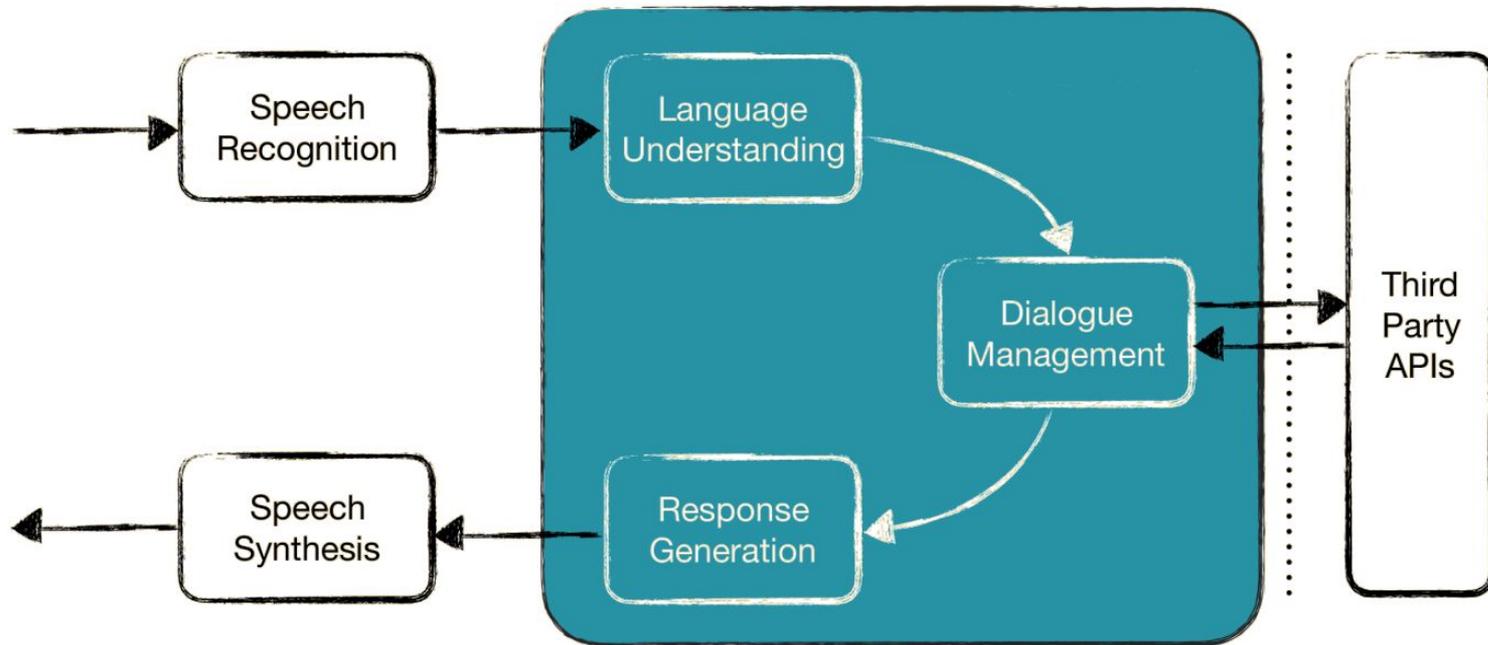
[Hershcovich+, ACL-22]

**Common Ground:** shared / common sense knowledge  
**Aboutness:** what people care to convey and talk about

# Why Cultural Adaptation?

- Mitigating conversational **bias towards source-culture concepts and contexts**
  - *Tailgating in Germany?*
  - *Baseball discussions in Croatia?*
  - *“March Madness” in Turkey?*
- Taking into account **specificities of the target culture**
  - *Postcode patterns vs no postcodes at all?*
  - *Buses vs trains in public transport?*
- **Avoiding ‘atypical’ culturally ungrounded dialogues?**
  - *The concept of “gastropub” in Arabic-speaking countries?*

# ("Old School") Task-Oriented Dialogue Systems



# Preliminaries: Three Pillars of Dialogue (or any ML-Driven Tech)

## Good NLU Performance

### Model



- + More efficient data usage
- + Less annotation effort
- + More accurate predictions
- Research-intensive and experimentation-intensive

### Data



- + Always improves performance
- Expensive to create
- Requires multiple cycles

### Design



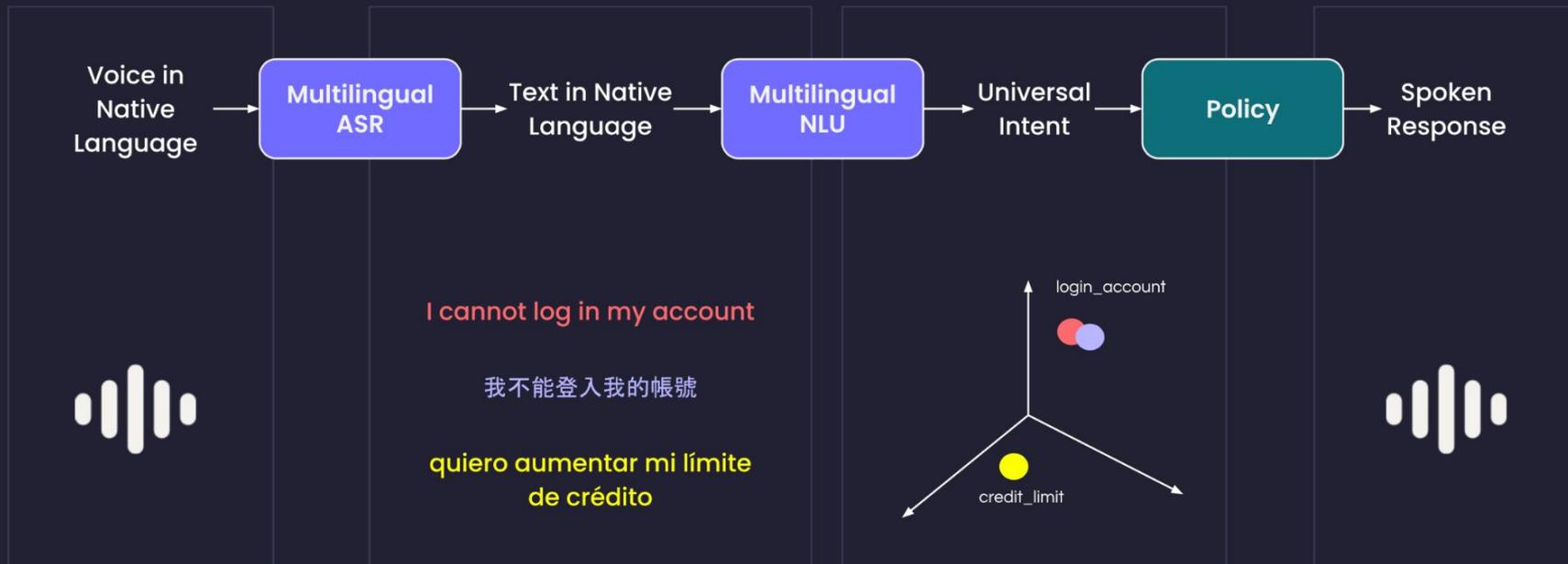
- + Improves performance
- + Improves efficiency and 'semantic support'
- Requires some expertise, intuition and practice

# **MULTI<sup>3</sup>NLU<sup>++</sup>: A Multilingual, Multi-Intent, Multi-Domain Dataset for Natural Language Understanding in Task-Oriented Dialogue**

**Nikita Moghe<sup>\*1\*</sup>, Evgeniia Razumovskaia<sup>\*2</sup>, Liane Guillou<sup>1</sup>,  
Ivan Vulić<sup>2</sup>, Anna Korhonen<sup>2</sup>, Alexandra Birch<sup>1</sup>**  
School of Informatics, University of Edinburgh<sup>1</sup>  
Language Technology Lab, University of Cambridge<sup>2</sup>

[Findings of ACL-23]

# (Multilingual) Intent Detection



# Multi-Label Intent Detection and Modular “Subintents”



Intents: affirm, card, arrival, less\_lower\_before

Yes, I need this card to arrive before 3pm on Jan 14  
time date

Intents: greet, change, spa, booking

Hi, can I change my spa reservation for Friday?  
date

Intents: booking, make, accesibility

One accessible room for two adults from the 24th to the 4th  
rooms adults date\_from date\_to

- Reusability and composability (**across domains**)
- **“Semantic sharing”** and data-efficient generalisation
- Handling more complex scenarios (and with smaller intent sets)

(English-only) NLU++ [Casanueva et al., NAACL-HLT 2022]

# Why Multi<sup>3</sup>NLU++?

## Challenges:

- Enabling training and evaluation of multi-domain NLU models in multiple languages
- When data is scarce, it should be high quality
- Cross-lingual approaches should ideally make improvements **across the 'resourceness' spectrum**
- No datasets for effective **cross-domain** and **cross-lingual** evaluation in **realistic dialogue problems**

# What is Multi<sup>3</sup>NLU++?

Multi<sup>3</sup>NLU++, a dataset for dialogue NLU which is:

- **Multi**-lingual: English, Spanish, Turkish, Marathi, Amharic
- **Multi**-domain: BANKING and HOTELS (and COMBINED)
- **Multi**-intent: each example is labelled with multiple intents  
(and also **Multi**-parallel, it should have been called **Multi<sup>4</sup>NLU++**)
- Realistic, conversational language

As a benchmark, Multi<sup>3</sup>NLU++ allows for systematic, controlled comparison:

- across the languages with different levels of resources
- across domains – on seen and unseen intents
- across dialogue NLU tasks

# What is in Multi<sup>3</sup>NLU++?

Intents: balance, overdraft, how\_much

en: I spent \$58 in overdraft. What is my current balance?

am: ከሂሳቤ ላይ ማውጣት ከምቅለው በላይ 58 ዶላር አውጥቻለሁ አሁን ያለኝ ቀሪ ሂሳብ ስንት ነው?

mr: मी ओव्हरड्राफ्टमध्ये ५८ डॉलर्स खर्च केले. माझी सध्याची शिल्लक किती आहे?

tr: ek hesabımdan 58 dolar harcadım. Mevcut bakiyem ne?

es: Gasté 58 dólares en descubierto. ¿Cuál es mi saldo actual?

amount\_of\_money

# How was Multi<sup>3</sup>NLU++ Created?

**Language selection:** diverse level of ‘resourceness’, different language families, scripts and geographical spread

- **Spanish** - high-resource, Romance, Latin script
- **Marathi** - mid-resource, Indo-Aryan, Devanagari script
- **Turkish** - mid-resource, Turkic, Latin script
- **Amharic** - low-resource, Semitic, Ge’ez script

*(p.s. Turkish is agglutinative vs synthetic/fusional languages)*

## **Manual translation:**

Professional, aimed at preserving the colloquial nature of the English utterances

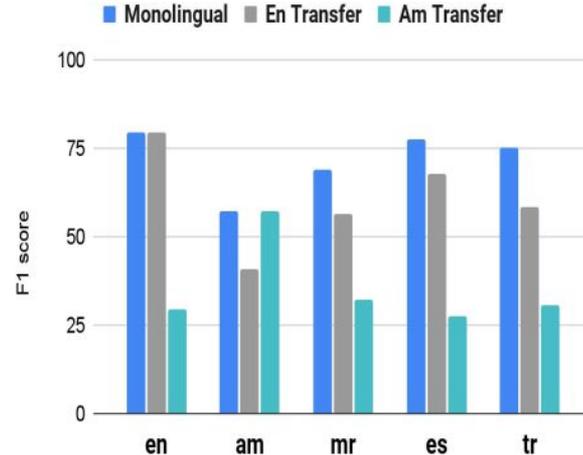
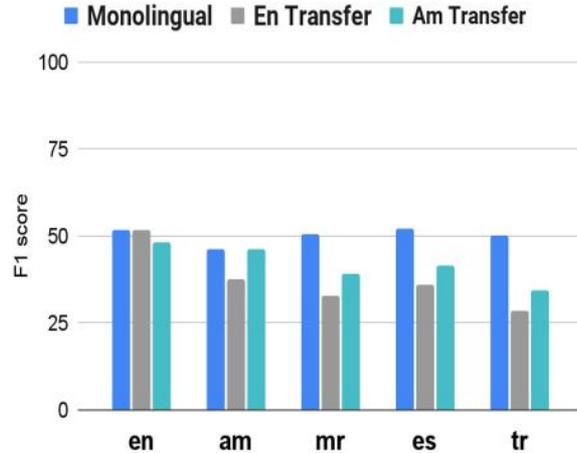
- *Three translators per language*
- *Slot span verification: IAA of ~90%*

# Some Modeling Paradigms in Comparison



- + (Then) state-of-the-art sentence encoders (e.g., LaBSE)
- + (Still) state-of-the-art multilingual encoders (e.g., XLM-R, mDeBERTa)
- + We also test **standard (full) fine-tuning**

# A (Tiny) Sample of Results...



- **QA-based models** are better than MLP-based models and full fine-tuning
- **Performance drop for all languages beyond English** (although the training and test data are multi-parallel!)
- **Amharic (MLP)** and **English (QA)** are best sources for cross-lingual transfer

## ...with Some (More) General Empirical Findings

*(there are many more in the paper)*

- **Overall trend across languages:** with more training data we gain better performance both in-domain and cross-domain
  - **Language-specific:** absolute numbers are indicative of resources available for pretraining of the model
  - **The gap between low- and high-resource languages** is rooted in (i) the amount of in-task training data; (ii) representational power of multilingual models
- 
- In the **cross-domain setup** high-resource languages benefit more from the increase in training data size than lower-resource languages
  - **Cross-domain cross-lingual generalisations:** The lower-resource the language, the lower the performance

# Cross-Lingual Dialogue Dataset Creation via Outline-Based Generation

**Olga Majewska<sup>◇</sup> Evgeniia Razumovskaia<sup>◇</sup> Edoardo M. Ponti<sup>†◇</sup>**

**Ivan Vulić<sup>◇</sup> Anna Korhonen<sup>◇</sup>**

<sup>◇</sup>Language Technology Lab, University of Cambridge

<sup>†</sup>Institute for Language, Cognition and Computation, University of Edinburgh

## Two Main Goals

- 1. Get rid of translation-based and test “translation-free” data creation**

...and get rid of negative effects of “translationese”...

- 2. Verify the ability to do (preliminary) cultural adaptation**

...and its importance in dialogue

*(Proof-of-concept work: the total cost of the whole dataset was £800)*

# Bottom-Up (*Outline-Based*) Dialogue Creation

**Stage 1:** Source Dialogue Sampling

**Stage 2:** Outline Generation

**Stage 3:** Dialogue Writing

**Stage 4:** Slot Span Validation

Language	ISO	Family	Branch	Macro-area	L1 [M]	Total [M]
Russian	RU	Indo-European	Balto-Slavic	Eurasia	153.7	258
Standard Arabic	AR	Afro-Asiatic	Semitic	Eurasia / Africa	0 <sup>†</sup>	274
Indonesian	ID	Austronesian	Malayo-Polynesian	Papunesia	43.6	199
Kiswahili	SW	Niger–Congo	Bantu	Africa	16.3	69

# Stage 1: Source Dialogue Sampling

## Starting point: English Schema-Guided Dialogue (SGD) Dataset

- Readily available (abstracted) dialogue schemata  
[\[Rastogi et al., arXiv-19\]](#)
- We randomly sample dialogues for 11 domains, 10 examples per intent

	Alarm (◇)	Flights	Homes	Movies	Music	Media	Banks	Payment (◇)	RideSharing	Travel	Weather	#turns
<b>Dev</b>	13	12	12	16	14	-	14	-	-	12	18	<b>1138</b>
<b>Test</b>	21	23	13	19	16	17	-	8	11	-	-	<b>1352</b>

service_name: "Payment" description: "Digital wallet to make and request payments"	<b>Service</b>
name: "account_type" description: "Source of money to make payment" possible_values: ["in-app balance", "debit card", "bank"]	category: True <b>Slots</b>
name: "amount" description: "Amount of money to transfer or request"	category: False
name: "contact_name" description: "Name of contact for transaction"	category: False
name: "MakePayment" description: "Send money to your contact" required_slots: ["amount", "contact_name"] optional_slots: ["account_type" = "in-app balance"]	<b>Intents</b>
name: "RequestPayment" description: "Request money from a contact" required_slots: ["amount", "contact_name"]	

## Stage 2: Outline Generation

Act	Slot/Intent	Description	Value	Outline
INFORM_INTENT	SearchOnewayFlight	Search for one-way flights to the destination of choice	-	<i>Express the desire to search for one-way flights</i>
REQUEST	number_checked_bags	Number of bags to check in	2	<i>Ask if the number of bags to check in is 2</i>

(Some) Cultural Adaptation happens here:

- *New York -> Jakarta*
- *American Airlines -> Garuda Indonesia*
- *\$ -> Rp (Rupiah)*
- ...

Split	Localised Slot Values				
	AR	ID	RU	SW	AVG
Dev	42.98	59.68	61.60	76.51	60.19
Test	13.50	57.00	53.81	78.34	50.66

## Stage 3: Dialogue Writing from Outlines

---

### Outlines

---

**USER:** *Express the desire to search for roundtrip flights for a trip*

the name of the airport or city to arrive at: Seattle  
the company that provides air transport services: American Airlines

---

**ASSISTANT/SYSTEM:** *Inform the user that you found 1 such option(s). Offer the following option(s):*

the company that provides air transport services: American Airlines  
departure time of the flight flying to the destination: 7:35am  
departure time of the flight coming back from the trip: 4:15pm  
the total cost of the flight tickets: \$343

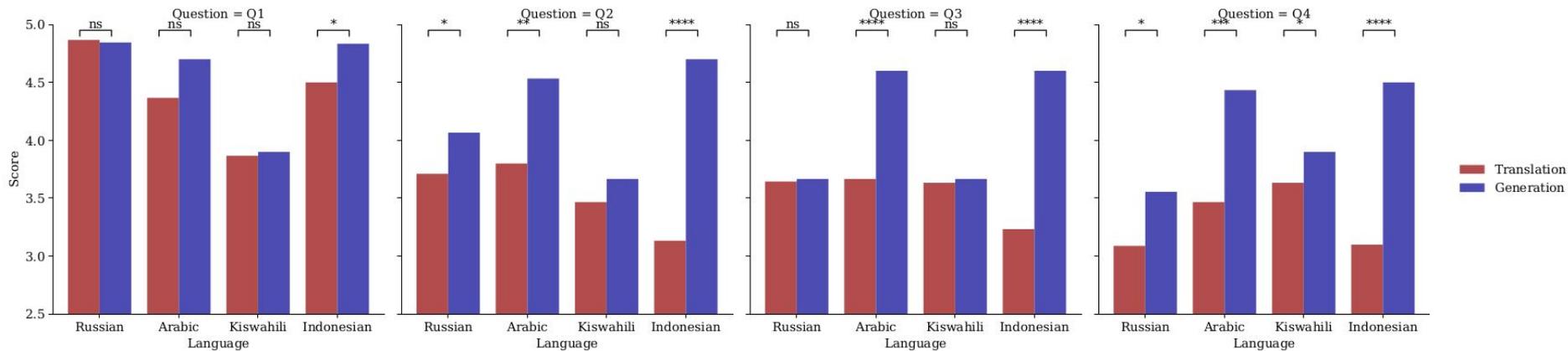
---

## Stage 4: Slot span verification (~99% agreement)

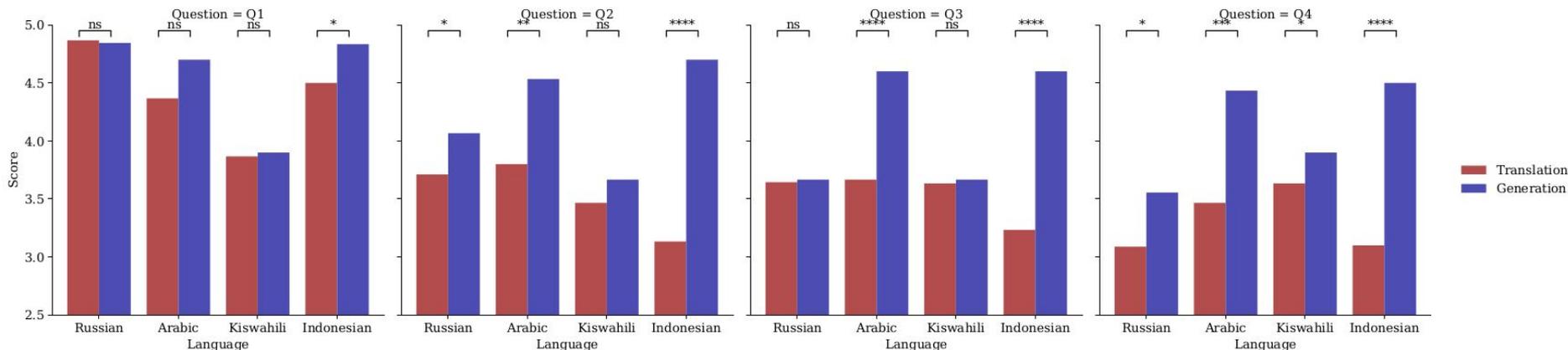
# Impact of Outline-Based Creation and Cultural Adaptation?

## Questions

- Q1. The ASSISTANT helps satisfy the USER's requests.  
Q2. The USER speaks naturally and sounds like an Arabic native speaker.  
Q3. The ASSISTANT speaks naturally and sounds like an Arabic 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 Arabic friends or family.

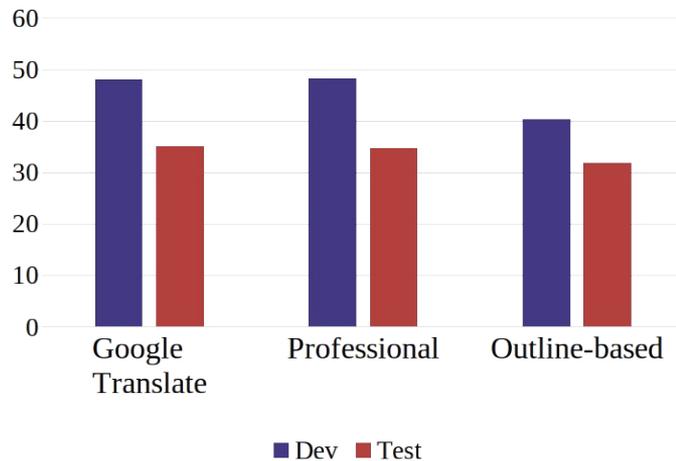
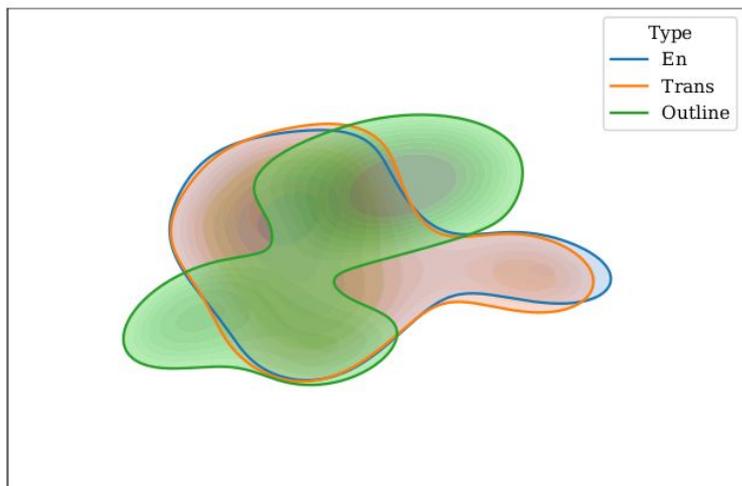


# Impact of Outline-Based Creation and Cultural Adaptation?



- Improved **naturalness, target-language fluency** (Q2, Q3) and **cultural familiarity** of entities (Q4)
- Effects of “translationese” in direct translation output:
  - **syntactic calques** (“*The meeting has been scheduled*”), **source lexical bias**
- **A/B Test** with 15 human participants per language: COD-based dialogues are more **natural-sounding** (80%+ in all 4 languages)

# MT-Based Data Creation Inflates Performance



...and this is the consequence...

“Non-natural” alignment of data samples?

# **MULTI<sup>3</sup>WOZ: A Multilingual, Multi-Domain, Multi-Parallel Dataset for Training and Evaluating Culturally Adapted Task-Oriented Dialog Systems**

**Songbo Hu<sup>1\*</sup> Han Zhou<sup>1\*</sup> Mete Hergul<sup>1</sup>  
Milan Gritta<sup>2</sup> Guchun Zhang<sup>2</sup> Ignacio Iacobacci<sup>2</sup>  
Ivan Vulić<sup>1†</sup> Anna Korhonen<sup>1†</sup>**

<sup>1</sup>Language Technology Lab, University of Cambridge, UK

<sup>2</sup>Huawei Noah's Ark Lab, London, UK

# Summary of Multilingual ToD Datasets

Dataset	# Langs	# Domains	# Train	# Test	No Translation?	Culturally Adapted?	Multi-P?
WOZ 2.0	3	1	600	400	✗	✗	✓
BiToD	2	5	2,894	451	✓	✓	✗
AllWOZ	8	5	40	50	✗	✗	✓
GlobalWOZ	21	7	0 (8,437)	500 (1,000)	✗	✓	✗
Multi <sup>2</sup> WOZ	5	7	0	1,000	✗	✗	✓
<b>Multi3WOZ</b>	4	7	7,440	860	✓	✓	✓

**The need has been recognised**

The solutions have been (too) quick or inadequate

# Solutions Have Been (Too) Quick and Inadequate

**GlobalWOZ** is full of design-triggered issues and inconsistencies:

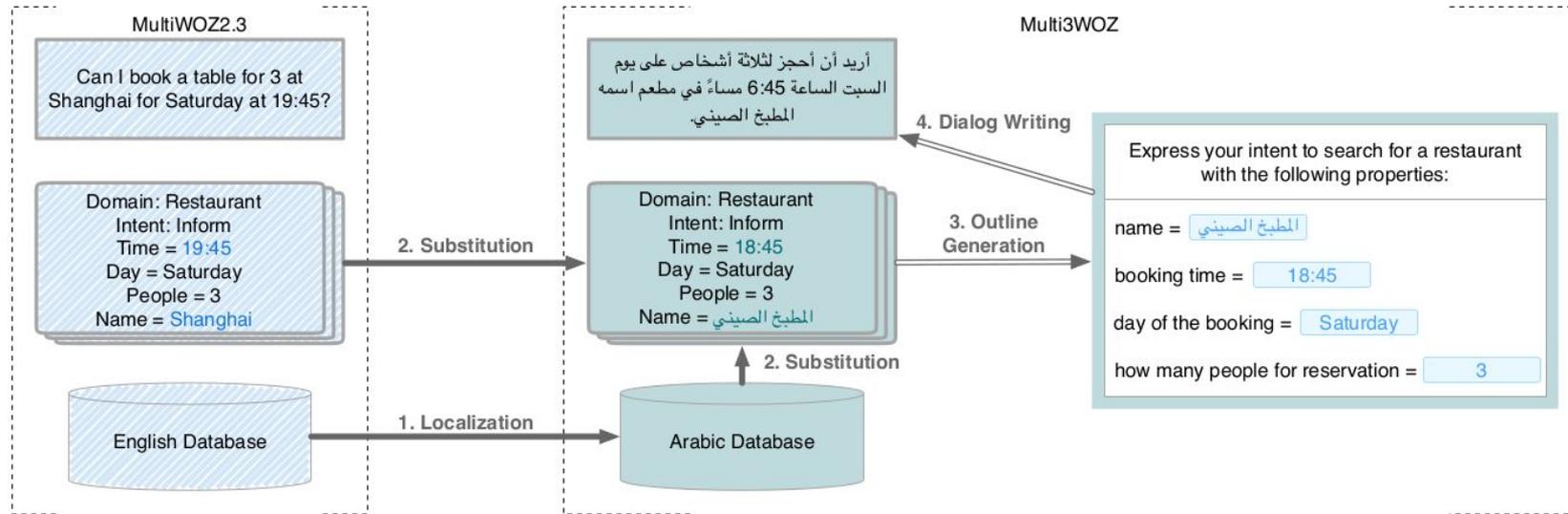
- **Inconsistent script-switched** and **code-switched** dialogues
- Erroneous **slot-value annotations inconsistent with dialogue ontology** and **database**
- **Contextual inconsistencies**
- **Translation-based:** automatic NMT for training data + PEMT for test data
- **Test sets** in different languages are not parallel
  - + *A heuristic for selection of dialogues penalises lexical variation*

# Main Goals

1. Get rid of translation-based and collect “translation-free” data multilingually **on a large scale**
  - *High-quality training, dev, and test data (e.g., getting rid of GlobalWOZ-style issues)*
2. Training, evaluation, and analysis of multilingual and cross-lingual ToD systems
  - ***Multi-parallel, multi-lingual, multi-domain***

*(Not a proof-of-concept any more: the total cost of the whole dataset was ~£55,000)*

# Multi<sup>3</sup>WOZ Dataset Construction: Bottom-Up (again)



- **"COD"-verified outline-based design!**
- **Cultural adaptation:** 1. localization, 2. substitution

# Example (Parallel) Dialogues

## English

I would like to eat Indian food.

There are 5 Indian restaurants available in the centre. What day are you looking to dine?

I need a reservation for Thursday. There will be 1 person at 11:15. Any of the 5 restaurants would be fine.

I was able to get you a table at [The Golden Curry](#). Your reference number is 7DD4DK3N and they 'll hold your table for 15 minutes. Can I help with anything else?

Thank you and goodbye.

You 're welcome. Goodbye.

## Arabic

أريد أن أجد مطعم هندي

هناك خمسة اختيارات لي مطاعم هندية في وسط المدينة. ما هو يوم الحجز الذي تريده؟

أريد أن أحجز لشخص واحد اليوم الخميس في الساعة ١١:١٥.

إن اسم المطعم هو [سانجام](#) ورقم حجزك هو SVRJYBEW. هل تحتاج أي مساعدة أخرى؟

شكرا للمساعدة

عفوا وحياك الله يا عزيزي.

## French

Je souhaiterais avoir celui qui propose les plats indiens.

Je souhaiterais avoir celui qui propose les plats indiens.

Je réserve pour 1 personne à 10 h 45 le jeudi.

Bien réservé! Numéro de référence UXQXMAA8. Nom de l'emplacement: [Indien 6](#). Autres choses?

Non, je vous remercie.

Bonne journée à vous.

## Turkish

Türk yemekleri sunan bir restoran anyorum.

merkez bölgesinde [türk](#) yemekleri sunan beş farklı seçenek bulunuyor. Hangi güne rezervasyon yaptırmak istersiniz?

perşembe günü saat 11:15 için 1 kişilik rezervasyon yaptırmak istiyorum.

[Göksu Lokantalan](#) restoranına yapılan A1IA291R referans kodlu rezervasyonunuz onaylanmıştır.

Teşekkürler.

Rica ederim. İyi günler.

- **Cultural adaptation:** slot-value redistribution, slot-value randomization, controlled entity replacement

# **A Systematic Study of Performance Disparities in Multilingual Task-Oriented Dialogue Systems**

**Songbo Hu<sup>1</sup> Han Zhou<sup>1</sup> Zhangdie Yuan<sup>2</sup> Milan Gritta<sup>3</sup>  
Guchun Zhang<sup>3</sup> Ignacio Iacobacci<sup>3</sup> Anna Korhonen<sup>1</sup> Ivan Vulić<sup>1</sup>**

<sup>1</sup>Language Technology Lab, University of Cambridge, UK

<sup>2</sup>Department of Computer Science and Technology, University of Cambridge, UK

<sup>3</sup>Huawei Noah's Ark Lab, London, UK

**Multi<sup>3</sup>WOZ** is multi-parallel and contains abundant (high-quality) training data:

**Analyses across different:**

- Source and target languages (cross-lingual transfer)
- Domains (cross-domain transfer)
- Learning setups (“many”-shot vs few-shot vs zero-shot)

# Preliminaries and Notation

- $P(\cdot)$ : a dialogue model
- $D$ : a task-specific dialogue dataset
- $D^{src}$ : a typically high-resource source language dataset
- $D^{tgt}$ : a low-resource target language dataset with equal size and quality as  $D^{src}$
- $D_{few}^{tgt}$ : a realistic low-resource target language dataset, which is considerably smaller compared to  $D^{src}$  and  $D^{tgt}$
- $M(\cdot)$ : an automatic evaluation metric

# Notions of Equivalence in Performance

- **Absolute  $\theta$ -Equivalence:** we define that two systems achieve absolute  $\theta$ -equivalence iff  $M(P^{tgt}(\cdot)) \geq \theta \cdot M(P^{src}(\cdot))$ , where  $\theta \in [0, 1]$ .
- **Relative  $\theta$ -Equivalence:** We define that the two systems achieve relative  $\theta$ -equivalence iff the metric  $M(P_{few}^{tgt}(\cdot)) \geq \theta \cdot M(P^{tgt}(\cdot))$ , where  $\theta \in [0, 1]$ .

## **(RQ1) Supervised, Translation-Based, Zero-Shot**

*RQ1) Given recent progress in multilingual LMs, machine translation, and cross-lingual transfer, is language-specific data still necessary for the development of a T O D system for a new language?*

# (RQ1) Supervised, Translation-Based, Zero-Shot

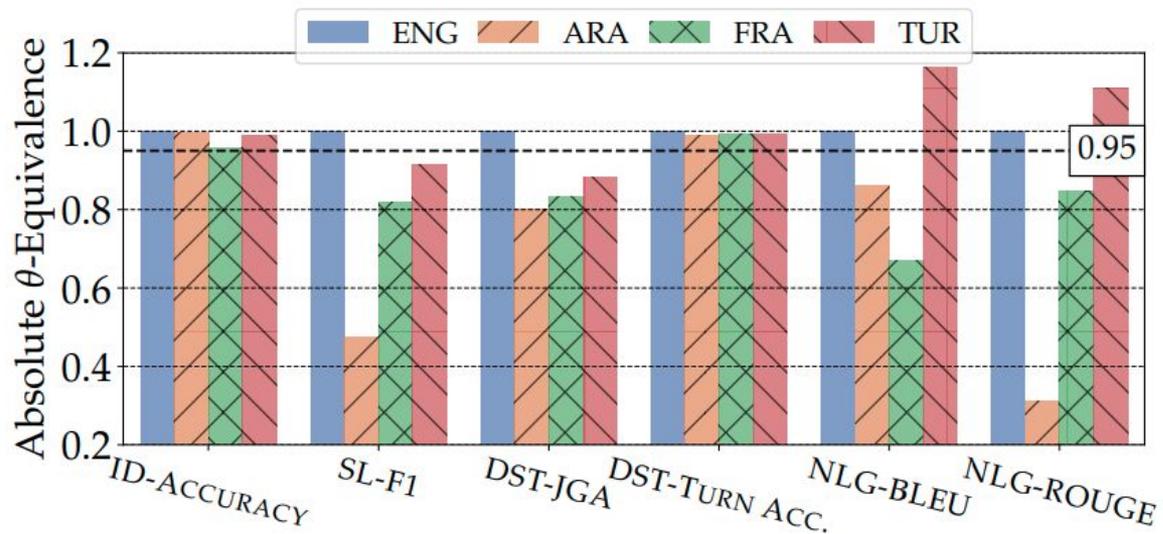
Language	Intent Detection		Slot Labelling			Dialogue State Tracking			Natural Language Generation		
	Accuracy	F1	Precision	Recall	F1	JGA	Turn Acc.	F1	BLEU	ROUGE	METEOR
<b>Fully Supervised</b>											
ENG	93.2 <sub>92.0</sub>	96.1 <sub>95.3</sub>	94.6 <sub>93.6</sub>	95.7 <sub>96.0</sub>	95.1 <sub>94.8</sub>	57.2 <sub>59.8</sub>	97.7 <sub>97.9</sub>	92.5 <sub>93.5</sub>	20.1 <sub>20.8</sub>	47.3 <sub>48.4</sub>	42.9 <sub>44.1</sub>
ARA	92.7 <sub>92.1</sub>	95.0 <sub>94.6</sub>	42.4 <sub>42.2</sub>	48.5 <sub>48.1</sub>	45.2 <sub>45.0</sub>	42.0 <sub>47.9</sub>	96.4 <sub>96.9</sub>	88.0 <sub>89.4</sub>	6.8 <sub>17.9</sub>	0.8 <sub>15.0</sub>	19.4 <sub>36.0</sub>
FRA	89.2 <sub>88.6</sub>	93.0 <sub>92.6</sub>	76.9 <sub>77.1</sub>	79.2 <sub>79.1</sub>	78.0 <sub>78.1</sub>	47.6 <sub>49.7</sub>	96.8 <sub>97.0</sub>	89.4 <sub>90.1</sub>	12.9 <sub>13.9</sub>	39.6 <sub>40.9</sub>	33.8 <sub>35.2</sub>
TUR	92.2 <sub>91.5</sub>	95.0 <sub>94.4</sub>	76.9 <sub>77.1</sub>	87.6 <sub>87.3</sub>	87.1 <sub>86.9</sub>	50.5 <sub>52.9</sub>	97.1 <sub>97.3</sub>	90.5 <sub>91.2</sub>	5.5 <sub>24.2</sub>	24.7 <sub>53.7</sub>	22.5 <sub>48.6</sub>
<b>Zero-shot Cross-lingual Transfer</b>											
ARA	82.1 <sub>65.7</sub>	88.2 <sub>74.8</sub>	27.4 <sub>17.2</sub>	31.2 <sub>27.7</sub>	29.2 <sub>21.2</sub>	1.9 <sub>1.5</sub>	82.5 <sub>80.7</sub>	17.0 <sub>5.8</sub>	0.2 <sub>0.2</sub>	2.5 <sub>2.1</sub>	2.4 <sub>2.0</sub>
FRA	83.9 <sub>77.0</sub>	89.8 <sub>85.0</sub>	58.5 <sub>49.1</sub>	61.2 <sub>62.4</sub>	59.8 <sub>54.9</sub>	5.5 <sub>3.7</sub>	86.6 <sub>85.1</sub>	40.1 <sub>32.8</sub>	0.5 <sub>0.4</sub>	4.2 <sub>4.7</sub>	6.1 <sub>5.9</sub>
TUR	87.0 <sub>74.9</sub>	91.4 <sub>81.7</sub>	68.1 <sub>48.5</sub>	74.7 <sub>66.6</sub>	71.2 <sub>56.2</sub>	3.5 <sub>1.3</sub>	85.2 <sub>82.1</sub>	34.4 <sub>15.2</sub>	0.3 <sub>0.4</sub>	3.7 <sub>4.4</sub>	6.1 <sub>5.8</sub>
<b>Translate Train</b>											
ARA	72.0 <sub>67.3</sub>	81.9 <sub>78.9</sub>	0 <sub>0</sub>	0 <sub>0</sub>	0 <sub>0</sub>	9.2 <sub>32.4</sub>	89.1 <sub>94.2</sub>	52.7 <sub>79.9</sub>	1.1 <sub>1.2</sub>	6.3 <sub>6.7</sub>	7.4 <sub>7.6</sub>
FRA	66.2 <sub>63.4</sub>	77.4 <sub>74.9</sub>	0 <sub>0</sub>	0 <sub>0</sub>	0 <sub>0</sub>	10.4 <sub>9.8</sub>	90.6 <sub>90.6</sub>	60.0 <sub>58.7</sub>	2.6 <sub>3.2</sub>	20.4 <sub>23.2</sub>	15.1 <sub>17.8</sub>
TUR	71.2 <sub>66.5</sub>	82.2 <sub>78.6</sub>	0 <sub>0</sub>	0 <sub>0</sub>	0 <sub>0</sub>	10.5 <sub>32.9</sub>	90.5 <sub>94.3</sub>	60.4 <sub>79.7</sub>	1.0 <sub>1.0</sub>	16.9 <sub>17.4</sub>	12.7 <sub>13.0</sub>

In-language data is crucial for performance

## **(RQ2) Intrinsic Bias in Multilingual Language Models**

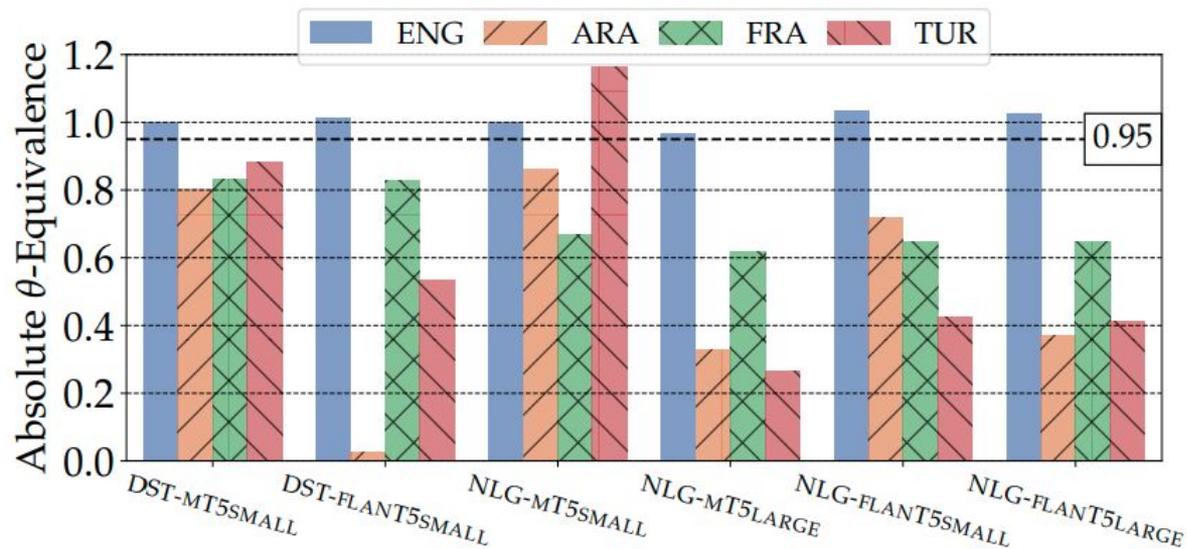
*(RQ2) Given access to the same mPLMs, equivalent amounts of high-quality in-language training data, and a similar development approach as that used to create an English ToD dataset, is it possible to develop a ToD system for a new language that achieves near-English performance?*

## (RQ2) Intrinsic Bias in Multilingual Language Models



Intrinsic bias is prominent, but depends on task complexity as well as evaluation metrics.

## (RQ2) Intrinsic Bias in Multilingual Language Models

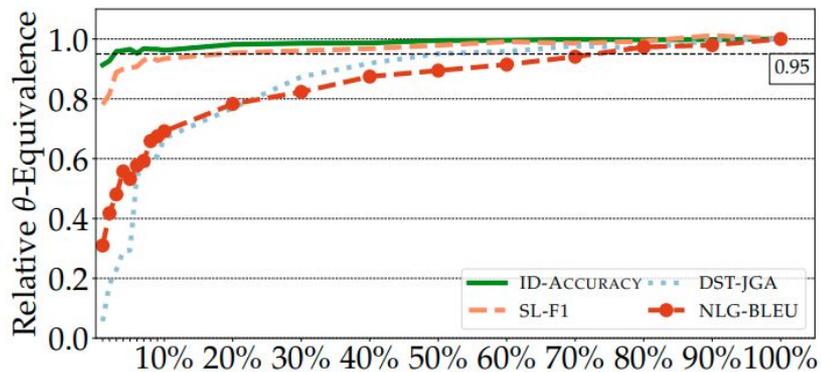


Intrinsic bias is prominent, but also depends on the chosen model, and it also exists with monolingual models

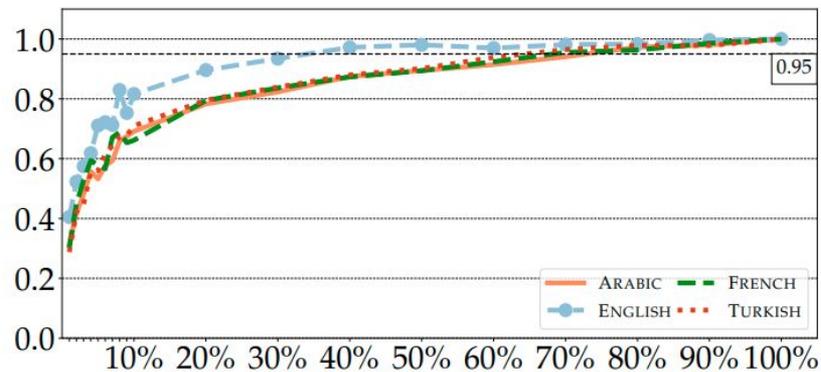
## **(RQ3) Adaptation Bias in Few-Shot Learning**

*(RQ3) How much training data is required in a new language to achieve performance comparable to a ToD system trained with an equivalent amount of in-domain, in-language data as in English?*

## (RQ3) Adaptation Bias in Few-Shot Learning



(a) Different Tasks on Arabic



(b) NLG on Different Languages

English is favoured even when it comes to collecting annotated task data...

## **(RQ4) Cost Efficiency in ToD Data Collection**

*(RQ4) Which data collection strategy maximises system performance across metrics while minimising the amount of annotation required? Such a strategy could optimise the cost-efficiency of annotation for a new language.*

## (RQ4) Cost Efficiency in ToD Data Collection

Strategy	ID	SL	DST	NLG
	Accuracy	F1	JGA	BLEU
Random Sampling	87.9	65.2	20.7	10.4
Max N-gram	<b>88.8</b>	<b>66.5</b>	23.6	<b>12.2</b>
Equal Domain	87.2	65.3	21.1	10.1
Equal Slot	87.9	65.0	26.2	11.3
Max Length	88.3	66.4	<b>26.7</b>	11.5

Averages over the three target languages based on 5% of target language data sampled/created using one of the strategies

**Tip:** Be less random than random sampling

**Future work:** Active learning? More sophisticated heuristics?

# Run Your Own Experiments and Analyses

Use DiaLight [\[Hu et al., NAACL-24: Demos\]](#)

Toolkit	Human Evaluation	Multilinguality	LLM+E2E	Comparative Experiment
PyDial	✓	✗	✗	✗
ConvLab2	✓	✗	✗	✗
ConvLab3	✓	✓	✗	✗
to-llm-bot	✗	✗	✓	✗
other E2E baselines	✗	✗	✗	✗
DIALIGHT(this work)	✓	✓	✓	✓

DiaLight supports:

- fine-tuning and in-context learning for development
- a comprehensive and simple framework for human evaluation
- creation of interactive systems that you can chat with



## Cultural Adaptation in Dialogue is:

- *Necessary*
- *Multi-Faceted / Multi-Layered*
- *Difficult*
- *Contextual*
- *Task-Specific*
- *Underexplored*
  - *Both from data and methodology angle*



## Achieving Multilingual (Performance) Equity in Dialogue is:

- *Necessary*
- *Multi-Faceted / Multi-Layered*
- *Difficult*
- *Contextual*
- *Task-Specific*
- *Underexplored*
  - *Both from data and methodology angle*



**While we haven't even scratched the surface of both. What about:**

- *Low-resource languages?*
- *Non-standard language varieties?*
- *Very complex and specific domains?*
- *Proper end-to-end learning (LLMs with RAG?)*
- *Many other types of equity (and more generally DEI) beyond performance only*

## Bonus: Tackling Data Scarcity with Data-Efficient Methods

# **SQATIN: Supervised Instruction Tuning Meets Question Answering for Improved Dialogue NLU**

**Evgeniia Razumovskaia<sup>1</sup>, Goran Glavaš<sup>2</sup>, Anna Korhonen<sup>1</sup>, Ivan Vulić<sup>1,3</sup>**

<sup>1</sup> Language Technology Lab, University of Cambridge

<sup>2</sup> Center for Artificial Intelligence and Data Science, University of Würzburg

<sup>3</sup> PolyAI Limited

# QA-Based Instruction Tuning of “Small LLMs”

Intent classification

The user says: How much in advance do I have to book a table for 8 people?  
Question: did the user intend to to talk about some booking?

yes

The user says: How much in advance do I have to book a table for 8 people?  
Question: did the user intend to to ask about check in?

no

Slot labelling

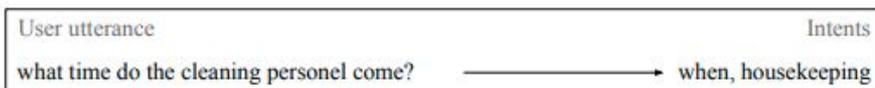
The user says: How much in advance do I have to book a table for 8 people?  
Question: what is the number of people mentioned in this sentence?

8

The user says: How much in advance do I have to book a table for 8 people?  
Question: what is the specific time in the day mentioned in this sentence?

unanswerable

# QA-Based Instruction Tuning of “Small LLMs”



<b>None</b>	Intent: wifi	what time do the cleaning personel come? Did the user intend to ask something related to wifi or wireless?	No
	Intent: housekeeping	what time do the cleaning personel come? Did the user intend to talk about housekeeping issues?	Yes
<b>Descriptive</b>	Intent: wifi	The user says: what time do the cleaning personel come? Question: did the user intend to ask something related to wifi or wireless?	No
	Intent: housekeeping	The user says: what time do the cleaning personel come? Question: did the user intend to talk about housekeeping issues?	Yes

# QA-Based Formulation Wins

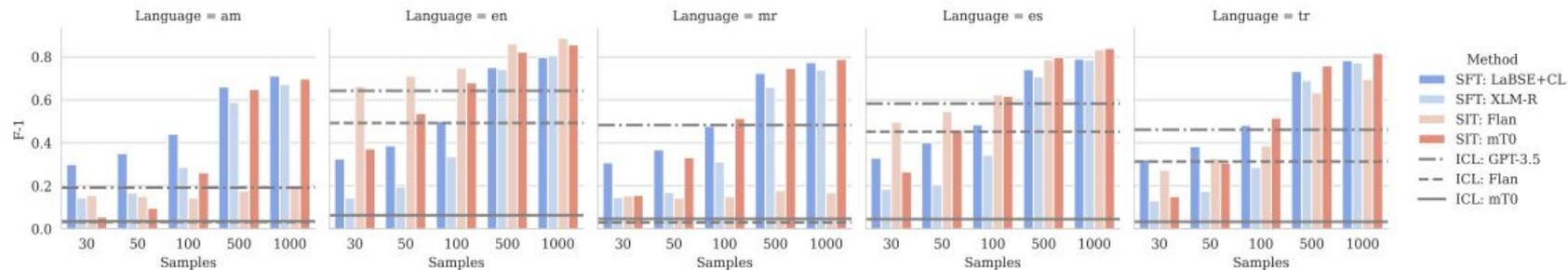
Model	Templ.	ID		VE	
		20-F	10-F	20-F	10-F
<b>BANKING</b>					
CL-SE		58.1	68.8	N/A	N/A
QA-FT: RoBERTa		80.3	85.6	50.5	56.7
QA-FT: mDeBERTa		80.8	85.0	59.7	66.5
QA-FT: T5		82.7	86.8	61.5	73.5
SQATIN	<i>None</i>	85.6	<b>88.5</b>	64.9	75.4
	<i>Desc.</i>	<b>85.8</b>	88.4	<b>66.3</b>	<b>76.3</b>
<b>HOTELS</b>					
CL-SE		51.9	61.8	N/A	N/A
QA-FT: RoBERTa		67.4	73.3	48.1	52.4
QA-FT: mDeBERTa		66.9	73.2	<b>61.6</b>	67.3
QA-FT: T5		69.2	76.5	57.2	<b>67.9</b>
SQATIN	<i>None</i>	73.1	78.0	58.0	<b>67.7</b>
	<i>Desc.</i>	<b>73.4</b>	<b>78.1</b>	58.7	67.0

The results are on English-only NLU++

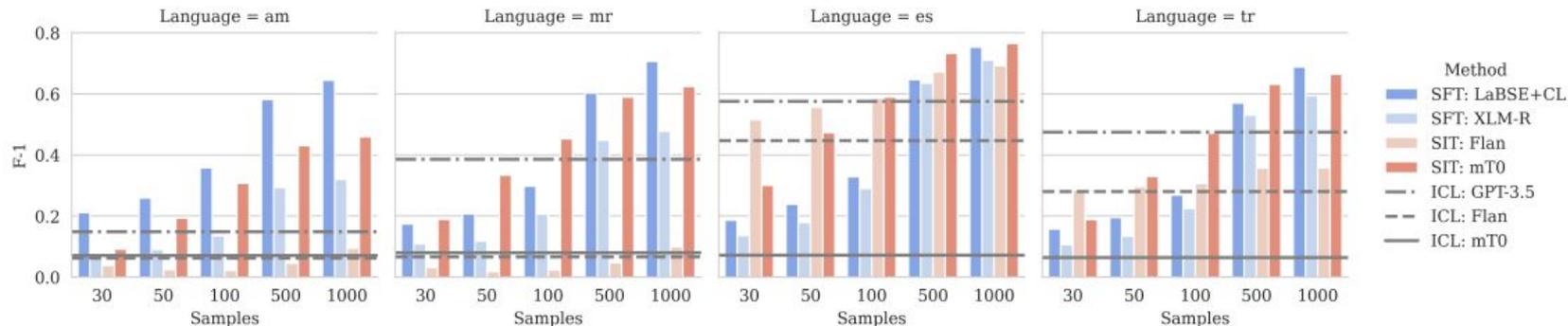
SQATIN is the most robust method in low-resource setups and across the board

Figures taken from [\[Razumovskaia et al., arXiv-24\]](#) (*under review in TACL*)

# What about Multilingual and Cross-Lingual Setups?



(a) ID: In-Language In-Domain

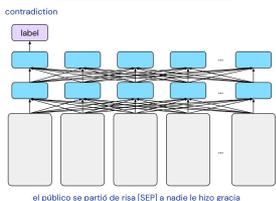


(c) ID: Cross-Lingual In-Domain

# Towards Inclusive, Sustainable, Equitable Multilingual TOD

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

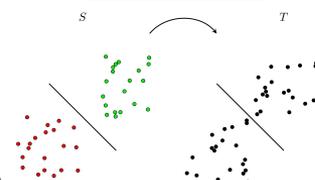
## Deep (machine) learning



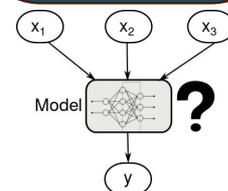
## Multilingual representation learning



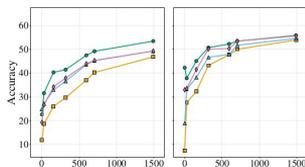
## Cross-lingual knowledge transfer



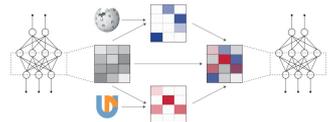
## Transparency



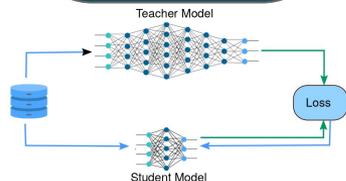
## Sample Efficiency



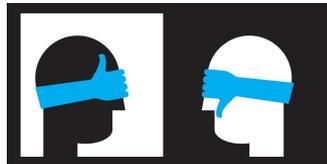
## Modularity



## Model Compactness



## Fairness and Debiasing



Plus Other Crucial Aspects: Cross-Cultural Adaptation, Multi-Modal Learning, Commonsense and World Knowledge, User Experience



UNIVERSITY OF  
CAMBRIDGE

λ ä Language  
宇 ش Technology  
w й Lab

**The talk is largely based on the following papers:**



## Massive thanks to my co-authors!

