

# Challenges and Strategies in Cross-Cultural NLP

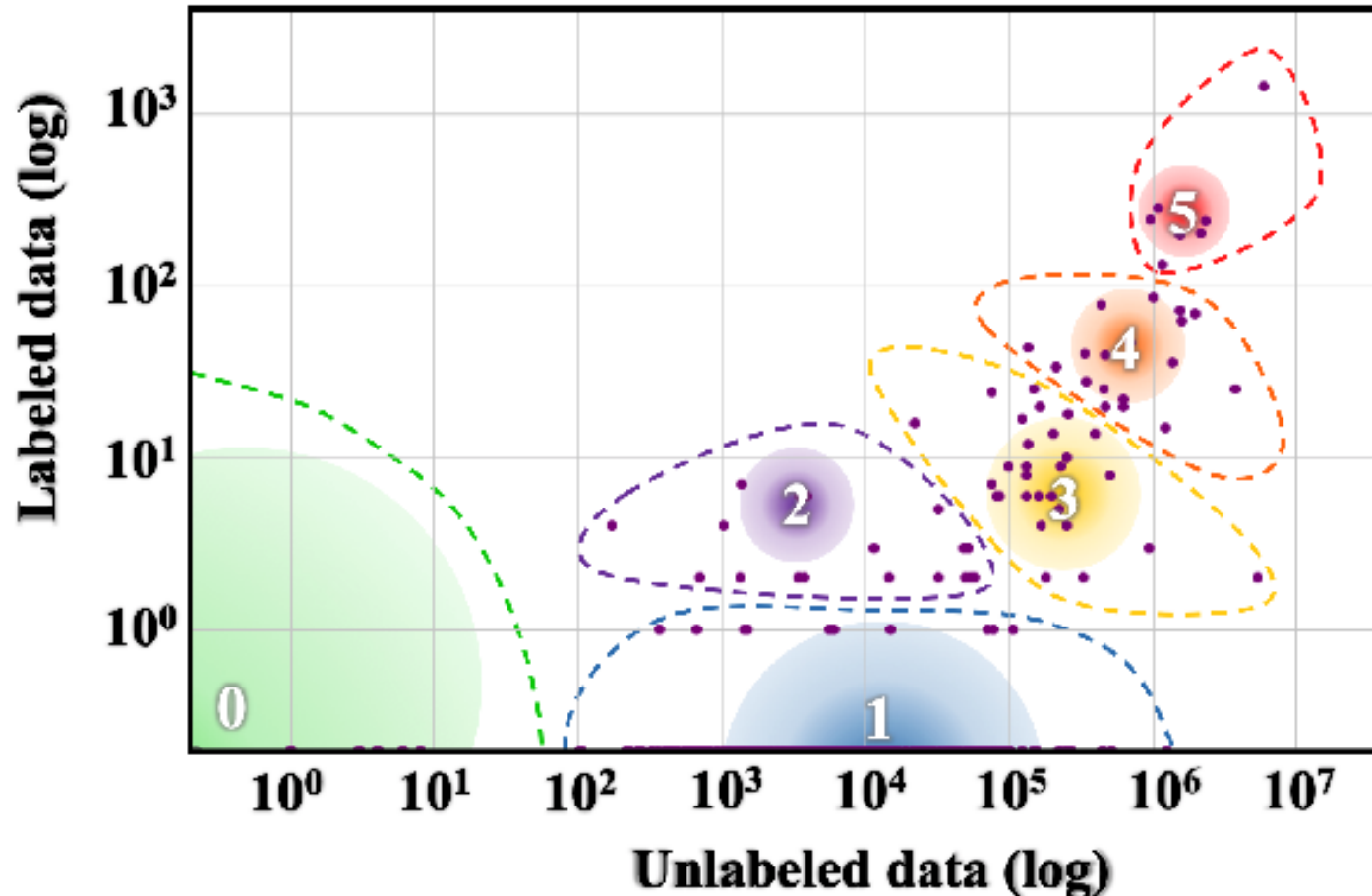
Daniel Hershcovich, Stella Frank, Heather  
Lent, Miryam de Lhoneux, Mostafa Abdou,  
Stephanie Brandl, Emanuele Bugliarello, Laura  
Cabello Piqueras, Ilias Chalkidis, Ruixiang  
Cui, Constanza Fierro, Katerina  
Margatina, Phillip Rust and Anders Søgaard

ACL 2022

UNIVERSITY OF COPENHAGEN



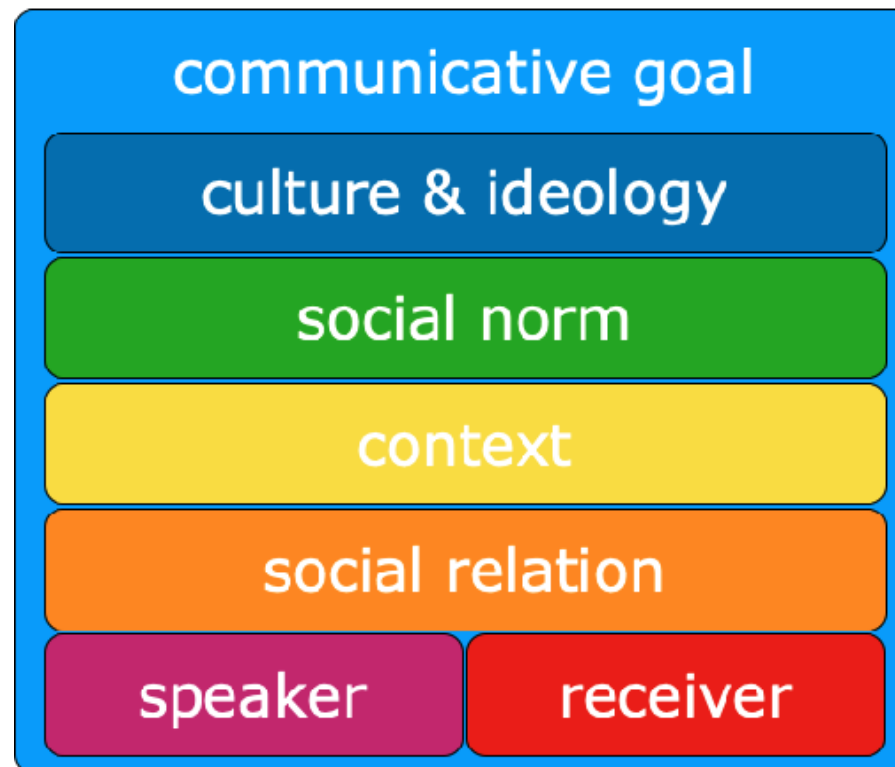
# Resource disparity for languages



[The State and Fate of Linguistic Diversity and Inclusion in the NLP World](#)  
(Joshi et al., ACL 2020)

# Social factors

NLP is for people (not just languages)



[The Importance of Modeling Social Factors of Language: Theory and Practice](#)  
(Hovy & Yang, NAACL 2021)

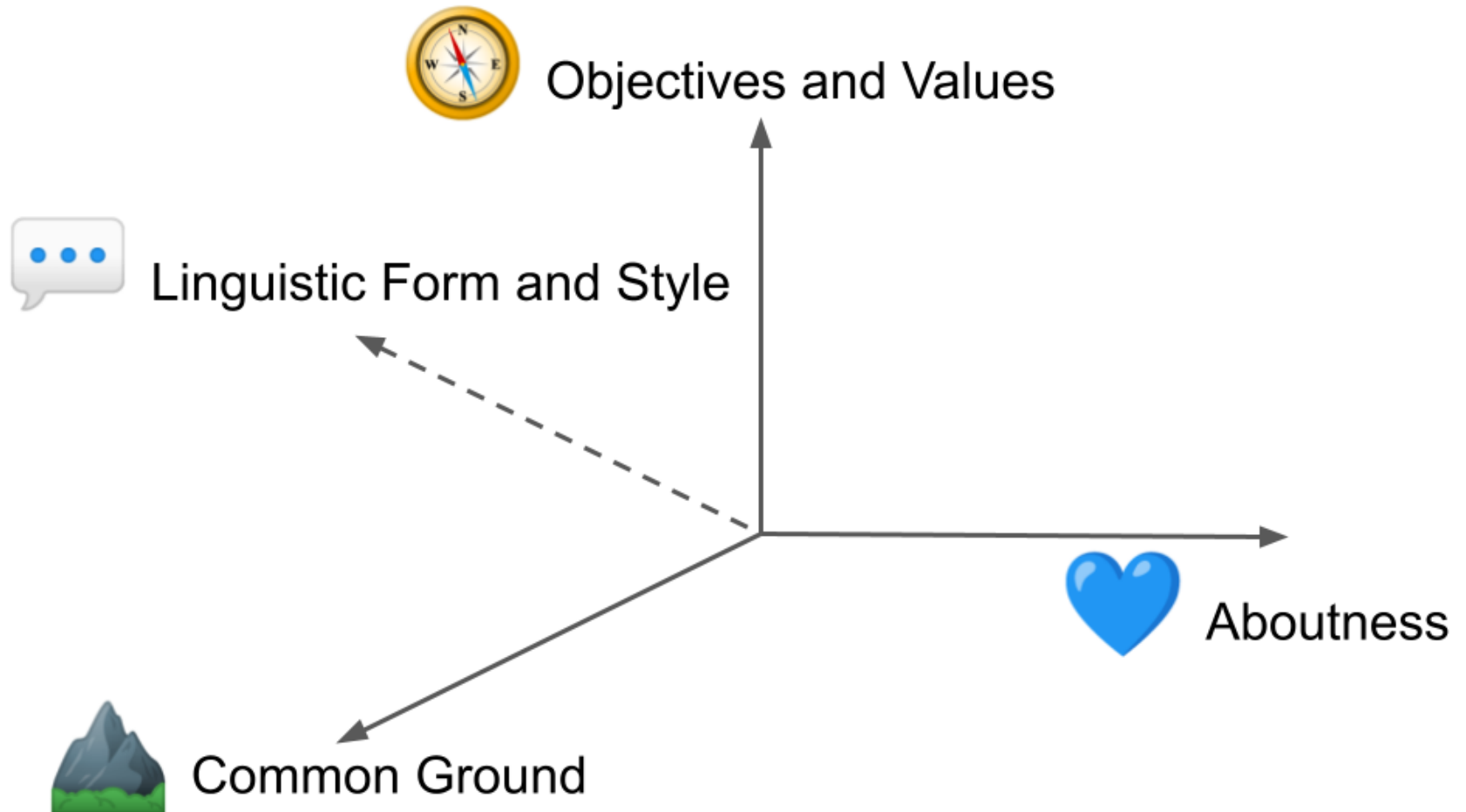
# Social bias in language models

Models	Demographics Alignment															
bert-base-cased																
bert-base-uncased																
bert-base-multilingual-cased																
bert-large-cased																
bert-large-uncased																
distilbert-base-uncased																
albert-base-v2																
albert-large-v2																
albert-xxlarge-v2																
roberta-base																
roberta-large																
google/electra-large-generator																
google/electra-small-generator																
gpt2																
gpt2-medium																
gpt2-large																
gpt2-xl																
<b>Group</b>																
<b>Mean Rank</b>	3.1	3.4	4.0	6.1	6.1	8.1	8.1	9.2	9.8	9.9	10.3	10.3	10.8	11.1	12.0	13.8

## Sociolectal Analysis of Pretrained Language Models

(Zhang et al., EMNLP 2021)

# Dimensions of culture



# Form

*How we express ourselves in language*

Morphosyntax

Word choice

Style

# Style

Stylistic aspects of linguistic form:

Directness

Formality

Politeness

Emotional expression

# Levels of granularity

Linguistic and cultural variation within groups



## **Idiolect**

Individual,  
personality

## **Sociolect, dialect**

Social group or region,  
sub-culture

## **Standardised language**

Country, national  
culture

## **Language, language family**

International cultures



# Common ground

Shared  
knowledge based  
on which people  
reason and  
communicate

Conceptualisation

Commonsense

# Conceptualisation

Objects

Colours

Kinship

Space

Time

Events

# Events and rituals



Tamil  
ஜல்லிக்கட்டு  
*jallikattu*

Visually Grounded Reasoning across Languages and Cultures  
(Liu et al., EMNLP 2021)

Visual concepts include culture-specific activities that cannot be mapped across cultures.

# Commonsense

Physical

Social

Taxonomic

Temporal

"Commonsense is the basic level of practical knowledge and reasoning concerning everyday situations and events that are commonly shared **among most people.**"

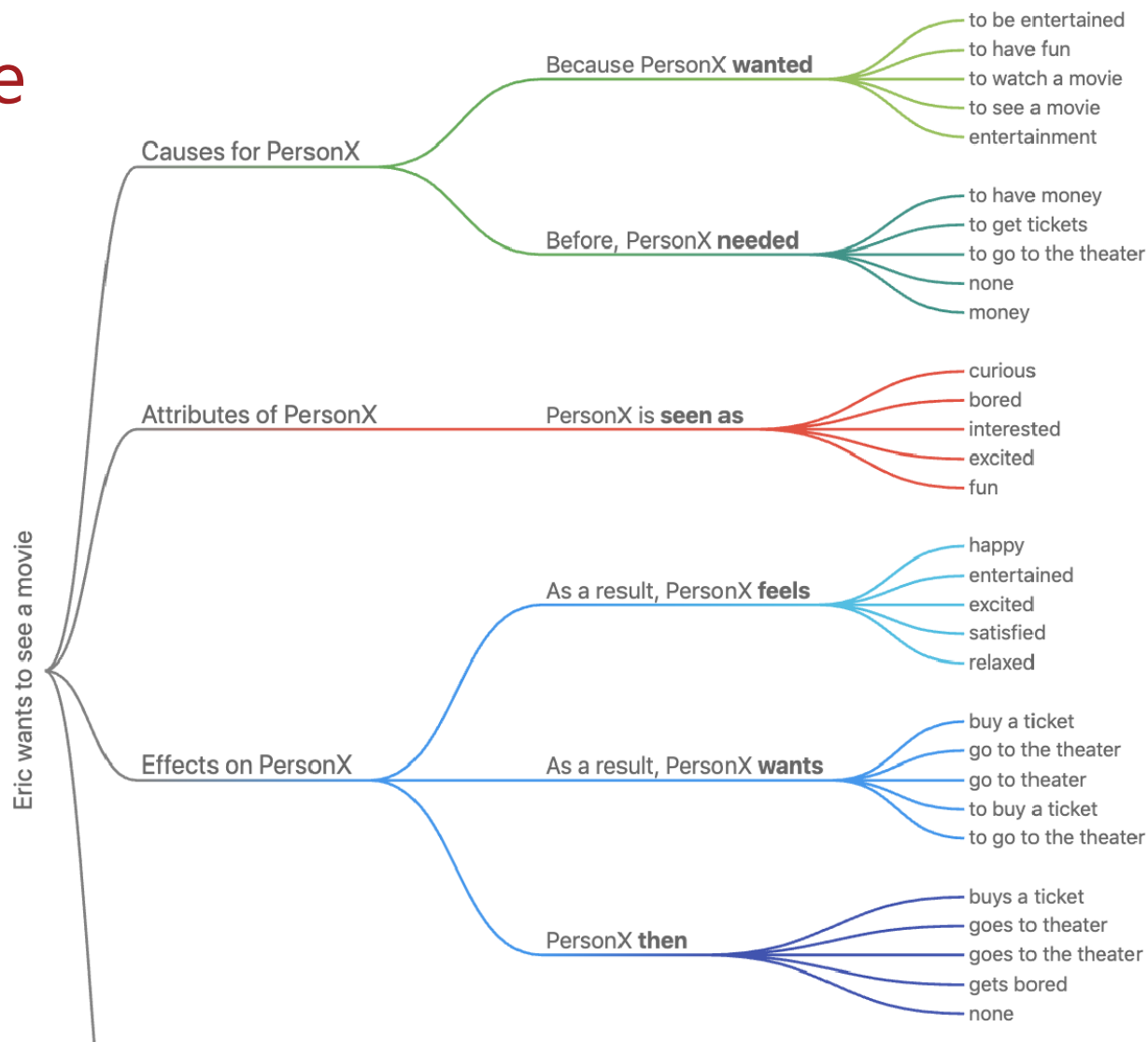
Commonsense Reasoning for Natural Language Processing

(Sap et al., ACL 2020 Tutorial)

Culture-dependent

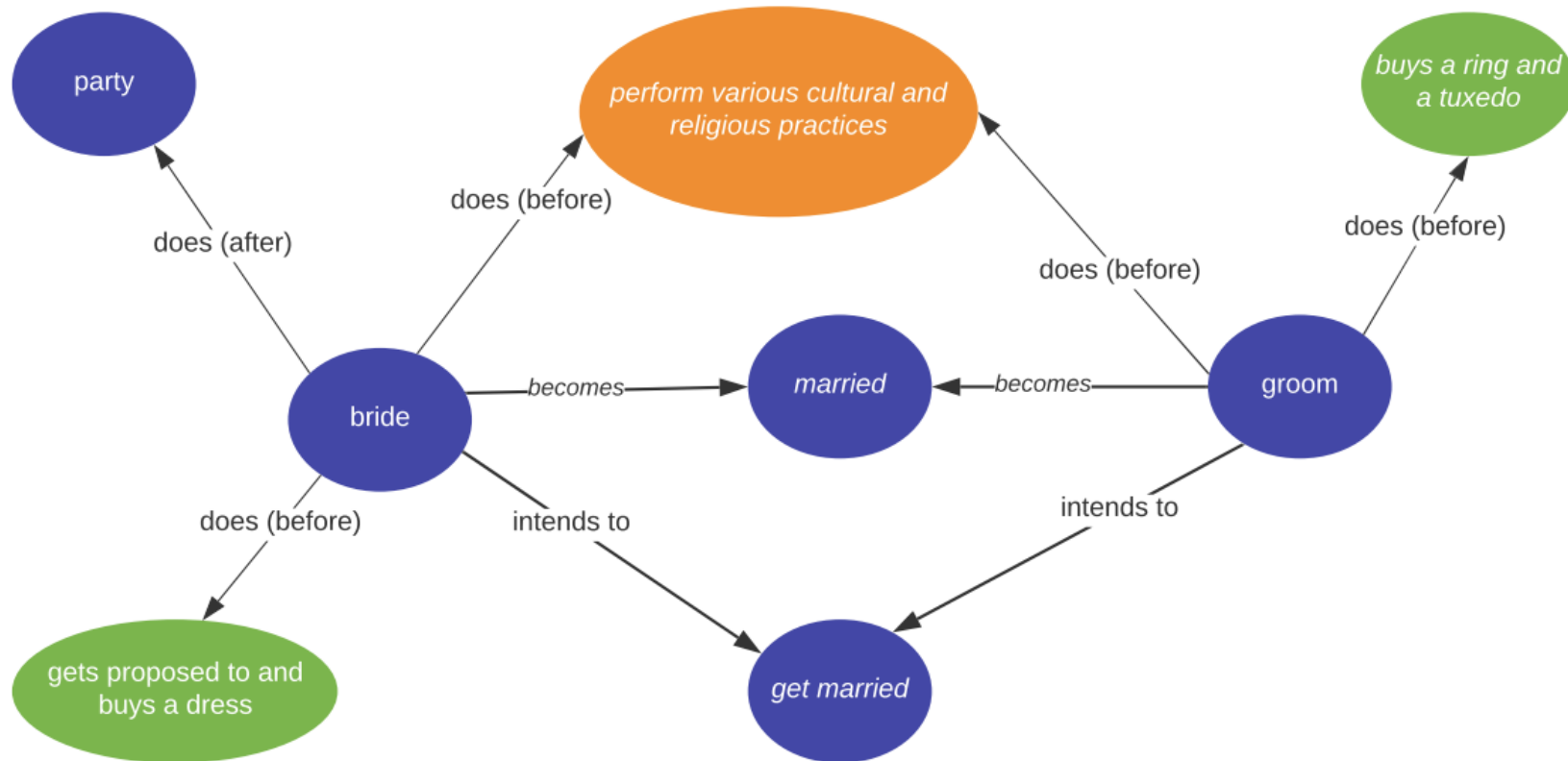
# Commonsense

Many resources to represent implicit knowledge are culturally biased



COMET: Commonsense Transformers for Automatic Knowledge Graph Construction  
(Bosselut et al., ACL 2019)

# Commonsense



Towards an Atlas of Cultural Commonsense for Machine Reasoning  
(Acharya et al., CSKGs 2021)

Some knowledge is "universal", other culture-specific

# Knowledge bias in language models

“[X] was created in [Y]”

en

---

Japan (170), Italy (56)

de

---

Deutschland (217), Japan (70)

nl

---

Nederland (172), Italië (50)

it

---

Italia (167), Giappone (92)

The language of prompting affects the model's answer to prompts

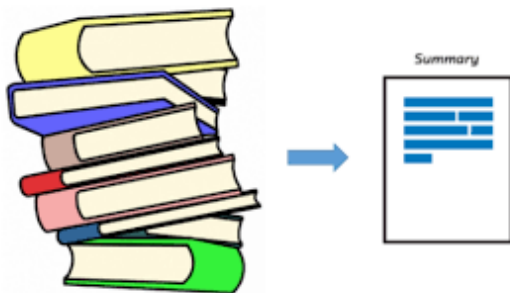


[Multilingual LAMA: Investigating Knowledge in Multilingual Pretrained Language Models](#) (Kassner et al., EACL 2021)

# Aboutness

What content do people *care about*?

- Related to topic/domain



Entities



Experiences



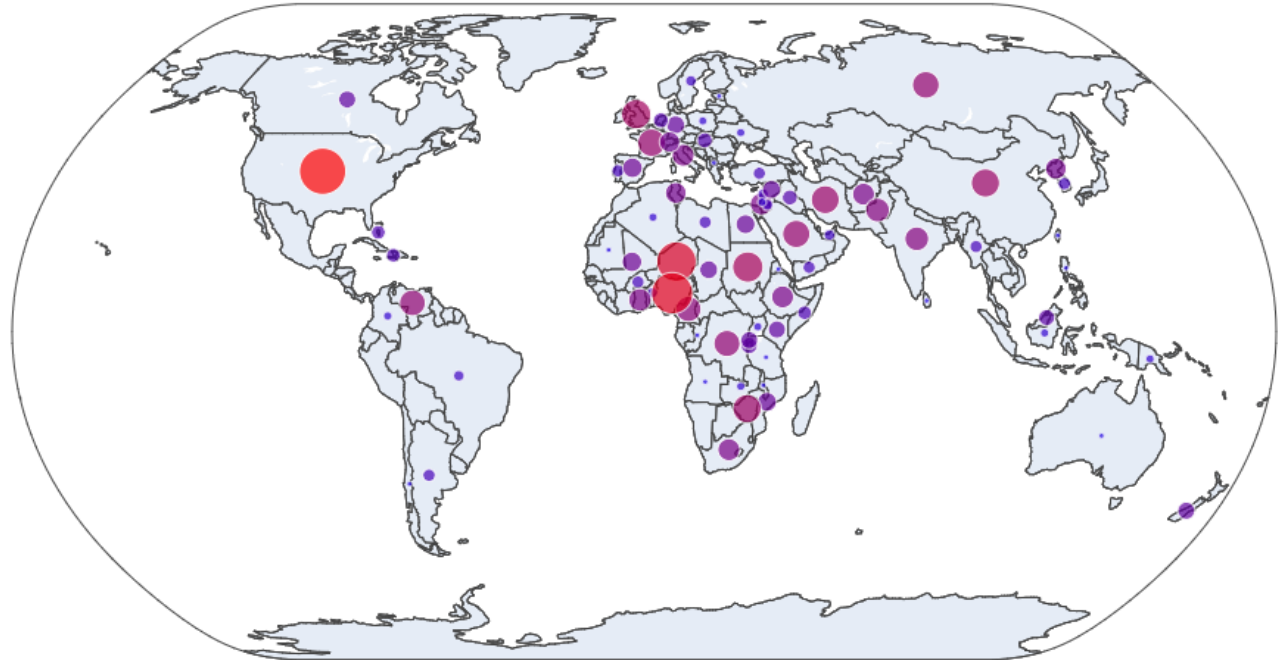
Aspects



# Entities

Dataset Map: Masakhaner hausa

Dataset Entities Map



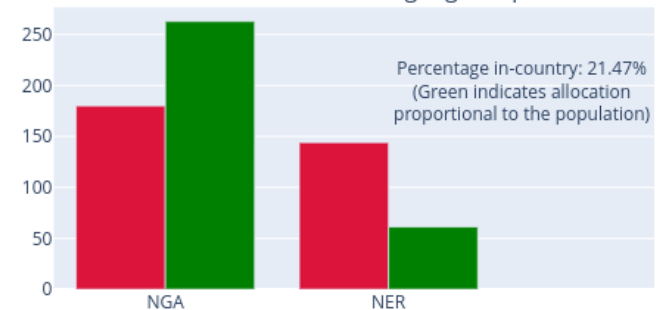
USA & Europe  
are over-  
represented in  
datasets across  
languages

Top-10 Represented Countries

Countries Missing: 143 of 243 (58.85%)



Main Countries where language is spoken.

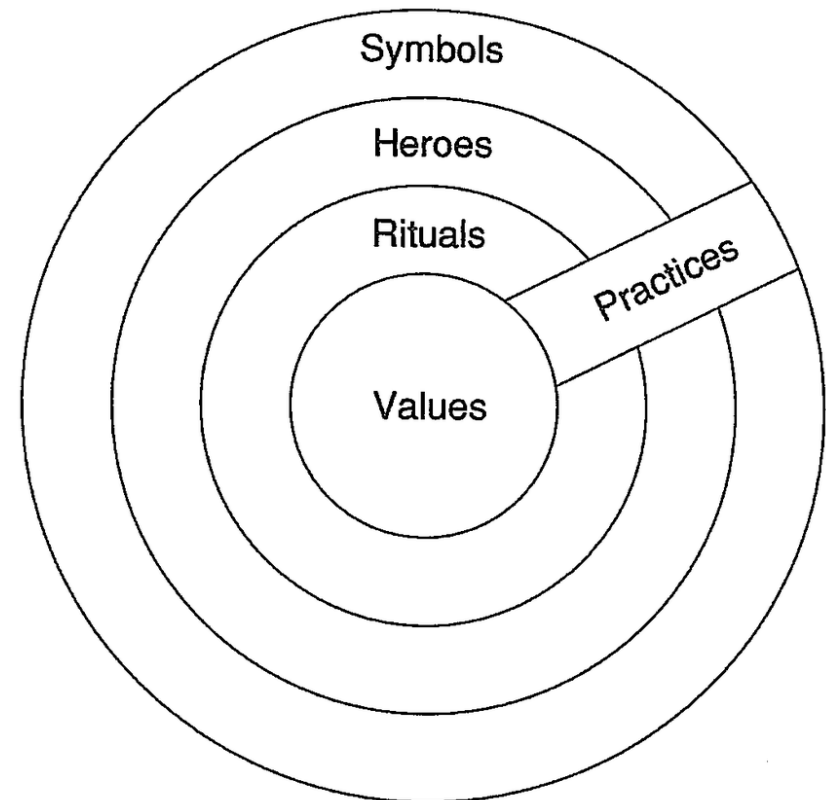


Dataset Geography: Mapping Language Data to Language Users  
(Faisal et al., 2021)

# Values

## Objectives and goals people strive for

- What is considered desired or desirable



Cultures and Organizations: Software of the Mind  
(Hofstede, 1991)

## (Meta) values

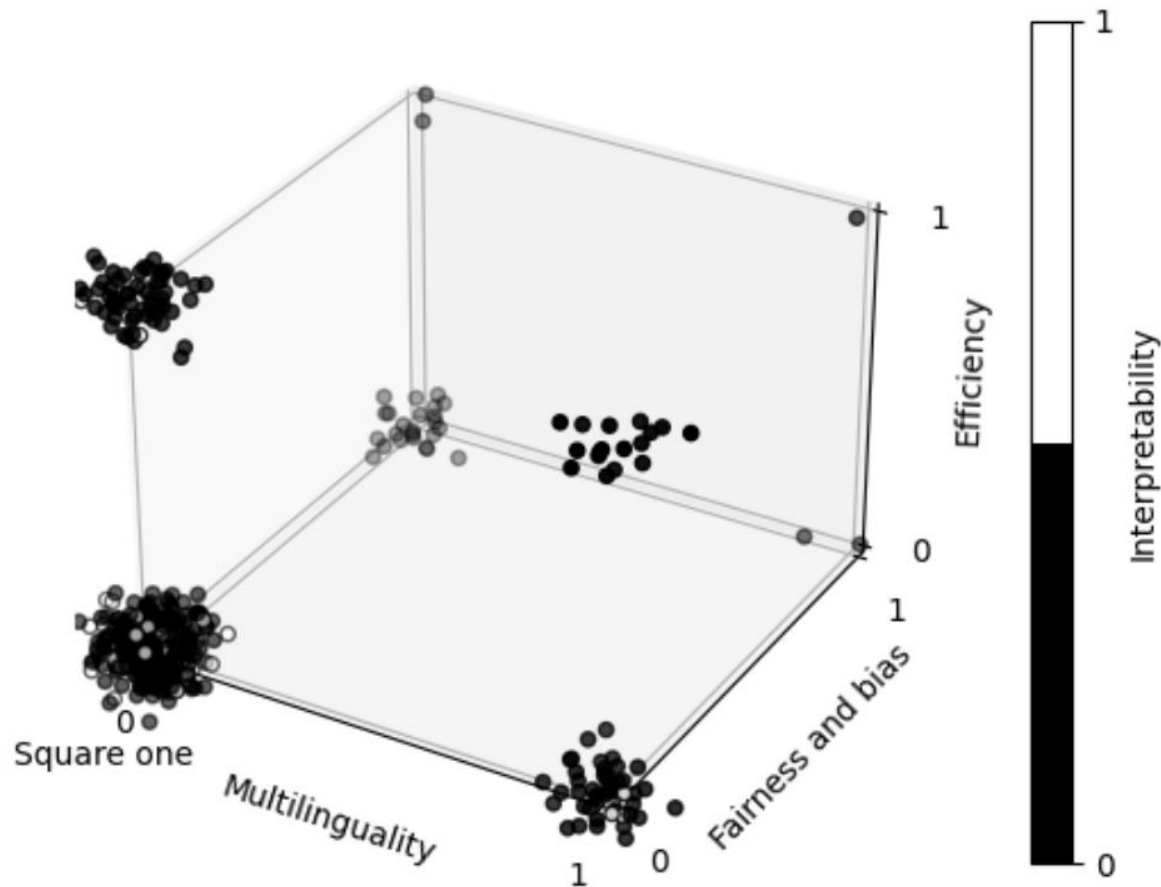
# *Why* are we doing NLP?

- Users may have different goals, often implicit

**No single correct answer.**

Changing the World by Changing the Data (Rogers, ACL 2021)

# Common meta-objectives in NLP



Accuracy,  
fairness,  
etc. reflect  
the values  
of NLP  
researchers

[Square One Bias in NLP: Towards a Multi-Dimensional Exploration of the Research Manifold](#) (Ruder et al., ACL 2022)

# Conflicting objectives between stakeholders



Researchers



Practitioners



End-users



Affected communities

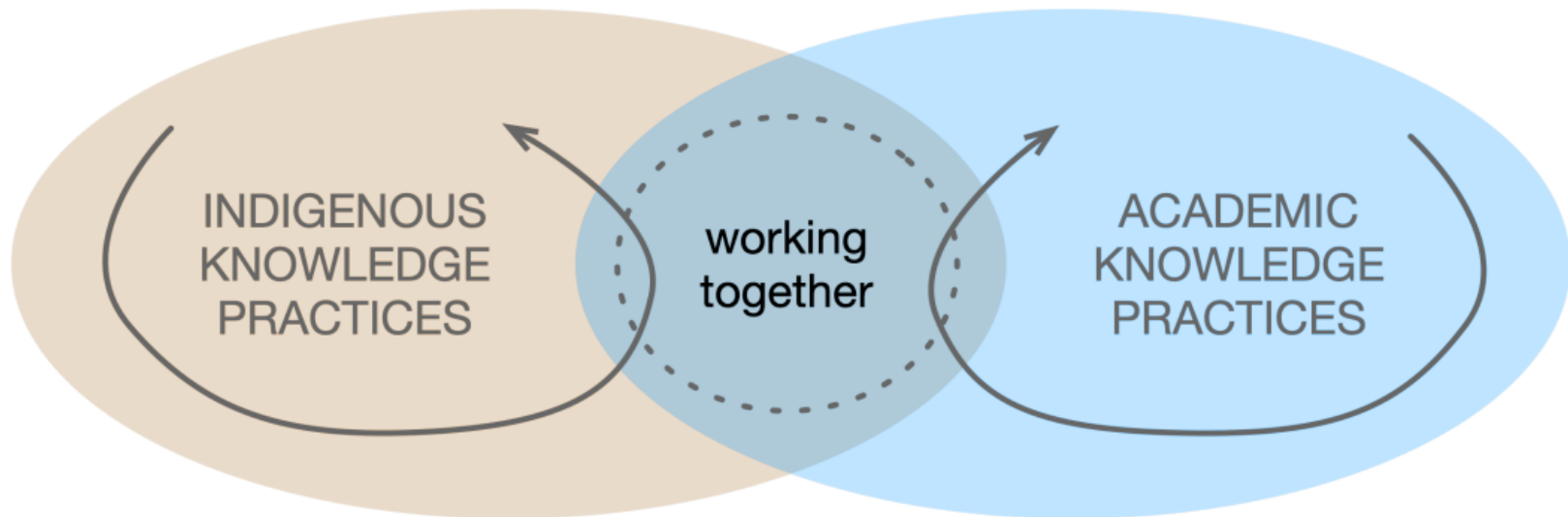


Regulators



Transparent values  
facilitate adaptation  
and decision making

# Language technology for all (potential) users



Local Languages, Third Spaces, and other High-Resource Scenarios  
(Bird, ACL 2022)

Benefit to all requires finding the intersection,  
particularly for local languages

# Value bias in language models



Die allermeisten von uns kennen den Zustand völliger Erschöpfung auf der Flucht, verbunden mit Angst um das eigene Leben oder das Leben der Kinder oder der Partner, zum Glück nicht. Menschen, die sich zum Beispiel aus Eritrea, aus Syrien oder dem Nordirak auf den Weg machen, müssen oft Situationen überwinden oder Ängste aushalten, die uns wahrscheinlich schlichtweg zusammenbrechen ließen. Deshalb müssen wir beim Umgang mit Menschen, die jetzt zu uns kommen, einige klare Grundsätze gelten lassen. Diese Grundsätze entstammen nicht mehr und nicht weniger als unserem Grundgesetz, unserer Verfassung.

Values are altered  
to reflect US culture



(translation)



"1. I am in favor of **limiting** immigration.  
2. I am in favor of **limiting** immigration for humanitarian reasons.  
3. I am in favor of **limiting** immigration for economic reasons."

The Ghost in the Machine has an American accent: value conflict in GPT-3

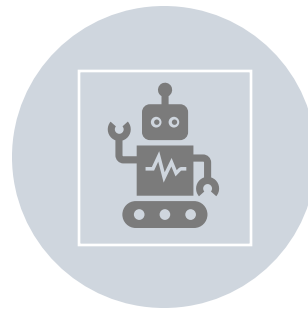
(Johnson et al., 2022)

# Strategies

Existing and potential ways to address the challenges



DATA



MODELS



TASKS



# Data



Selection



Annotation



Projection

Culture-sensitive curation

Culturally diverse collection

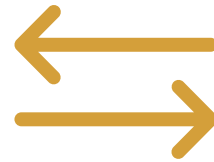
Native data or culturally sensitive translation

# Models



Training

Robust optimisation



Transfer

Balanced sampling



Pre-training

Measuring cultural bias

# Tasks

## Cross-cultural translation



*"I saw Merkel eating a Berliner from Dietsch on the ICE"*



*I saw Biden eating a Boston Cream from Dunkin' Donuts on the Acela*

Adapting Entities across Languages and Cultures

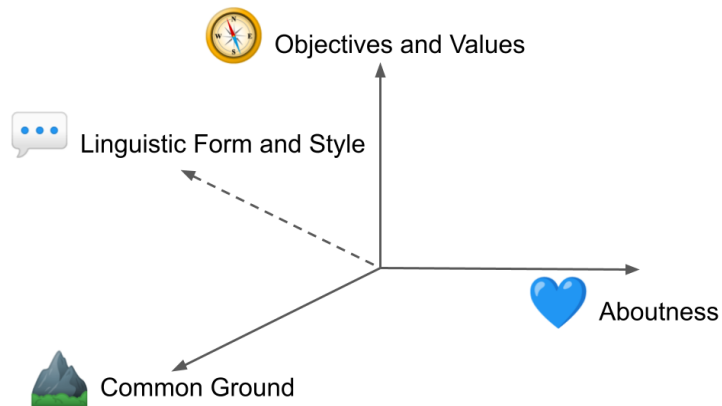
(Peskov et al., Findings 2021)

Style transfer

Entity adaptation

Explanation by analogy

# Summary



NLP is for people (not just languages)

Culture is multidimensional

Objectives may be in conflict

Generalisation-representation trade-off