# Some recent trends




USE BERT FOR EVERYTHING!!1!

# Some recent trends



Who cares

~~Chinchila~~

~~PaLM~~

~~GPT-2~~

~~ELECTRA~~

~~XLM-R~~

~~RoBERTa~~

USE ~~BERT~~ FOR EVERYTHING!!1!



2

NLP GEORGE MASON

# Make it multilingual!



Good recap of the current state of multilingual AI:
https://ruder.io/state-of-multilingual-ai/

# Lang Tech utility is unequally distributed!

# Lang Tech utility is unequally distributed!

Compare:
- American English speaker
- Arabic speaker
  - Tunisian vs Egyptian vs …
- Bemba speaker

# Global Utility Metrics

**Systematic Inequalities in Language Technology Performance across the World's Languages**

**Damián Blasi**
Harvard University
dblasi@fas.harvard.edu

**Antonios Anastasopoulos**
George Mason University
antonis@gmu.edu

**Graham Neubig**
Carnegie Mellon University
gneubig@cs.cmu.edu

NLP
GEORGE
MASON

5

# Global Utility Metrics

A language technology should be measured by the **utility** it provides to **every person in the world**

# Global Utility Metrics

A language technology should be measured by the **utility** it provides to **every person in the world**

$$M = \sum_i u_i$$

5

# Global Utility Metrics

**Systematic Inequalities in Language Technology Performance across the World's Languages**

**Damián Blasi**
Harvard University
dblasi@fas.harvard.edu

**Antonios Anastasopoulos**
George Mason University
antonis@gmu.edu

**Graham Neubig**
Carnegie Mellon University
gneubig@cs.cmu.edu

A language technology should be measured by the **utility** it provides to **every person in the world**

$$M = \sum_i u_i$$

Two problems:

NLP
GEORGE
MASON

# Global Utility Metrics

A language technology should be measured by the **utility** it provides to **every person in the world**

$$M = \sum_i u_i$$

Two problems:

- **Problem 1:** how to measure utility of an NLP system?
  → Very hard, use standard accuracy metrics as a proxy now (happy to discuss more!)

5

# Global Utility Metrics

A language technology should be measured by the **utility** it provides to **every person in the world**

$$M = \sum_i u_i$$

Two problems:

- **Problem 1:** how to measure utility of an NLP system?
  → Very hard, use standard accuracy metrics as a proxy now (happy to discuss more!)

- **Problem 2:** how to consider different utility provided to every person in the world?
  → Measure over subgroups (here, languages), weighted by demand + coefficient τ.

# Global Utility Metrics

A language technology should be measured by the **utility** it provides to **every person in the world**

$$M = \sum_i u_i$$

Two problems:

- **Problem 1:** how to measure utility of an NLP system?
  → Very hard, use standard accuracy metrics as a proxy now (happy to discuss more!)

- **Problem 2:** how to consider different utility provided to every person in the world?
  → Measure over subgroups (here, languages), weighted by demand + coefficient τ.

$$M_\tau = \sum_{l \in \mathcal{L}} d_l^{(\tau)} \cdot u_l$$

NLP
GEORGE
MASON

# Global Utility Metrics

**Systematic Inequalities in Language Technology Performance across the World's Languages**

Damián Blasi
Harvard University
dblasi@fas.harvard.edu

Antonios Anastasopoulos
George Mason University
antonis@gmu.edu

Graham Neubig
Carnegie Mellon University
gneubig@cs.cmu.edu

A language technology should be measured by the **utility** it provides to **every person in the world**

$$M = \sum_i u_i$$

Two problems:

- **Problem 1:** how to measure utility of an NLP system?
  → Very hard, use standard accuracy metrics as a proxy now (happy to discuss more!)

- **Problem 2:** how to consider different utility provided to every person in the world?
  → Measure over subgroups (here, languages), weighted by demand + coefficient τ.

$$M_\tau = \sum_{l \in \mathcal{L}} \boxed{d_l^{(\tau)}} \cdot u_l$$

"normalized demand"

5

# Global Utility Metrics

A language technology should be measured by the **utility** it provides to **every person in the world**

$$M = \sum_i u_i$$

Two problems:

- **Problem 1:** how to measure utility of an NLP system?
  → Very hard, use standard accuracy metrics as a proxy now (happy to discuss more!)

- **Problem 2:** how to consider different utility provided to every person in the world?
  → Measure over subgroups (here, languages), weighted by demand + coefficient τ.

$$M_\tau = \sum_{l \in \mathcal{L}} \boxed{d_l^{(\tau)}} \cdot \boxed{u_l}$$

"normalized demand"  "utility"

5

# Global Utility Metrics

**Systematic Inequalities in Language Technology Performance across the World's Languages**

Damián Blasi
Harvard University
dblasi@fas.harvard.edu

Antonios Anastasopoulos
George Mason University
antonis@gmu.edu

Graham Neubig
Carnegie Mellon University
gneubig@cs.cmu.edu

A language technology should be measured by the **utility** it provides to **every person in the world**

$$M = \sum_i u_i$$

Two problems:

- **Problem 1:** how to measure utility of an NLP system?
  → Very hard, use standard accuracy metrics as a proxy now (happy to discuss more!)

- **Problem 2:** how to consider different utility provided to every person in the world?
  → Measure over subgroups (here, languages), weighted by demand + coefficient τ.

$$M_\tau = \sum_{l \in \mathcal{L}} \boxed{d_l^{(\tau)}} \cdot \boxed{u_l} \quad d_l^{(\tau)} = \frac{n_l^\tau}{\sum_{l' \in \mathcal{L}} n_{l'}^\tau}$$

"normalized demand"  "utility"

5

# Global Utility Metrics

A language technology should be measured by the **utility** it provides to **every person in the world**

$$M = \sum_i u_i$$

Two problems:

- **Problem 1:** how to measure utility of an NLP system?
  → Very hard, use standard accuracy metrics as a proxy now (happy to discuss more!)

- **Problem 2:** how to consider different utility provided to every person in the world?
  → Measure over subgroups (here, languages), weighted by demand + coefficient τ.

τ=1 : every person equal
("demographic-average utility")

$$M_\tau = \sum_{l \in \mathcal{L}} \boxed{d_l^{(\tau)}} \cdot \boxed{u_l} \quad d_l^{(\tau)} = \frac{n_l^\tau}{\sum_{l' \in \mathcal{L}} n_{l'}^\tau}$$

"normalized demand"  "utility"

5

NLP
GEORGE
MASON

# Global Utility Metrics

**Systematic Inequalities in Language Technology Performance across the World's Languages**

Damián Blasi
Harvard University
dblasi@fas.harvard.edu

Antonios Anastasopoulos
George Mason University
antonis@gmu.edu

Graham Neubig
Carnegie Mellon University
gneubig@cs.cmu.edu

A language technology should be measured by the **utility** it provides to **every person in the world**

$$M = \sum_i u_i$$

Two problems:

- **Problem 1:** how to measure utility of an NLP system?
  → Very hard, use standard accuracy metrics as a proxy now (happy to discuss more!)

- **Problem 2:** how to consider different utility provided to every person in the world?
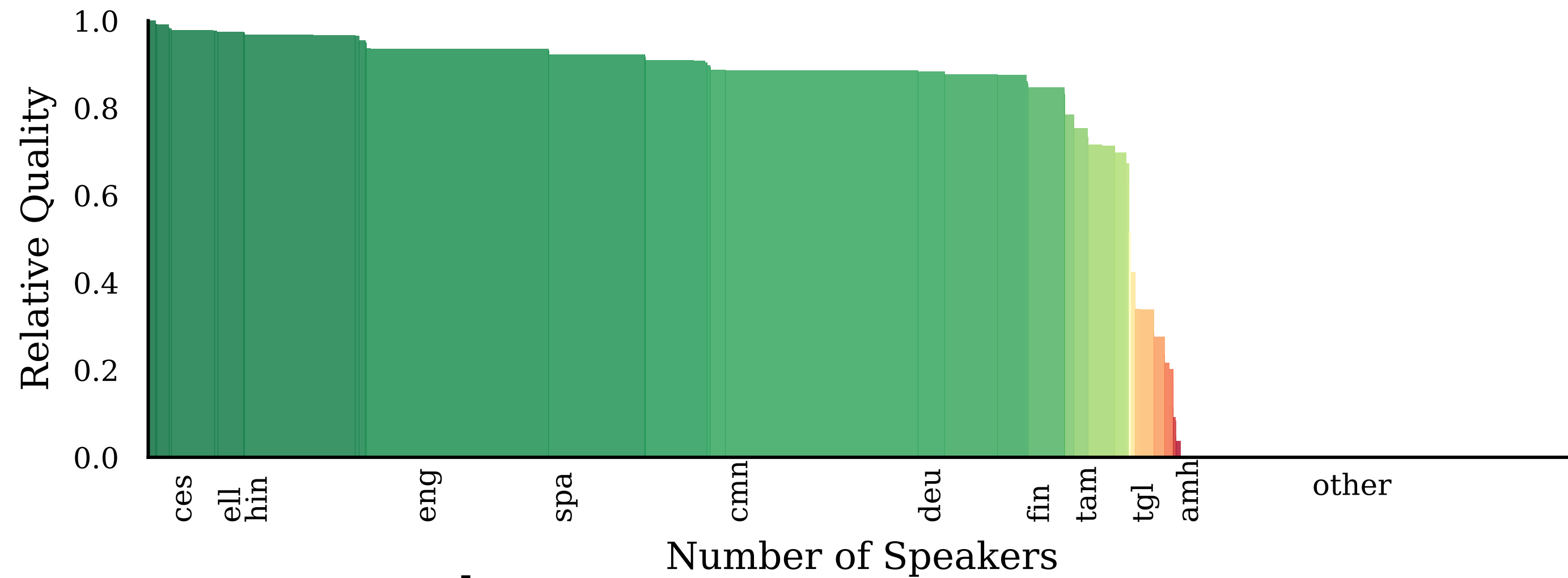  → Measure over subgroups (here, languages), weighted by demand + coefficient τ.

$$M_\tau = \sum_{l \in \mathcal{L}} \boxed{d_l^{(\tau)}} \cdot \boxed{u_l} \quad d_l^{(\tau)} = \frac{n_l^\tau}{\sum_{l' \in \mathcal{L}} n_{l'}^\tau}$$

τ=1 : every person equal
("demographic-average utility")
τ=0 : every subgroup equal
("linguistic-average utility")

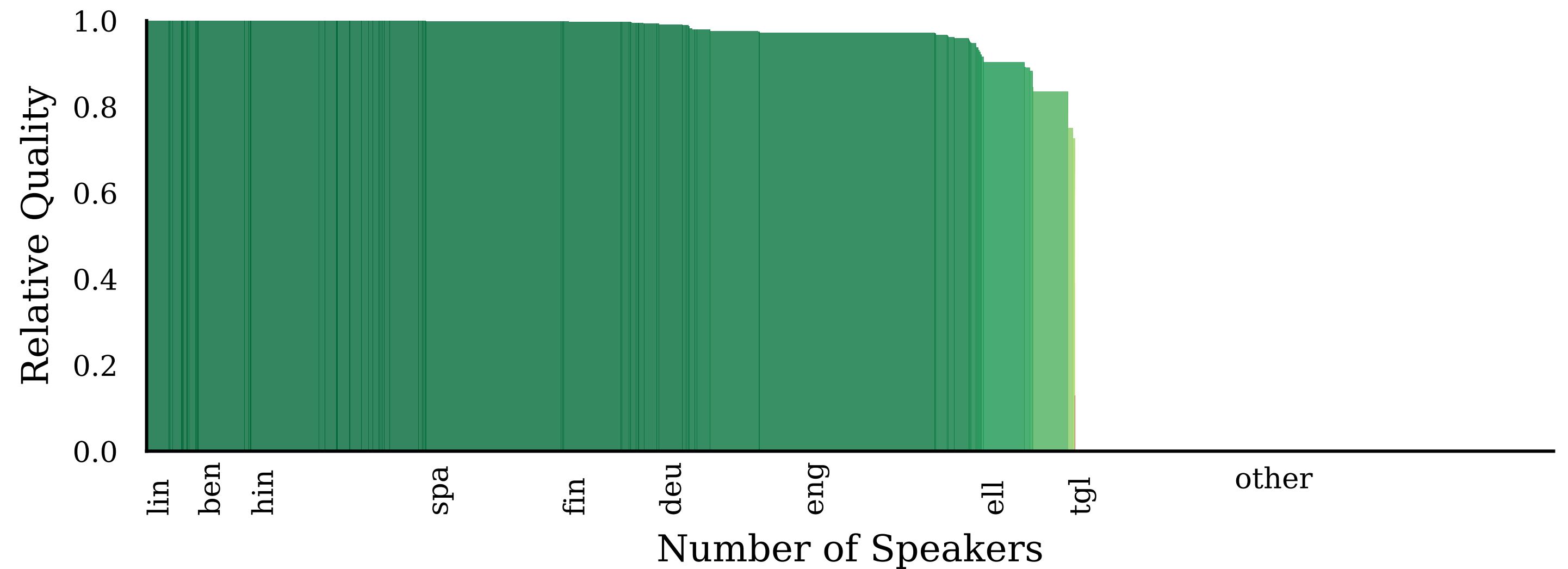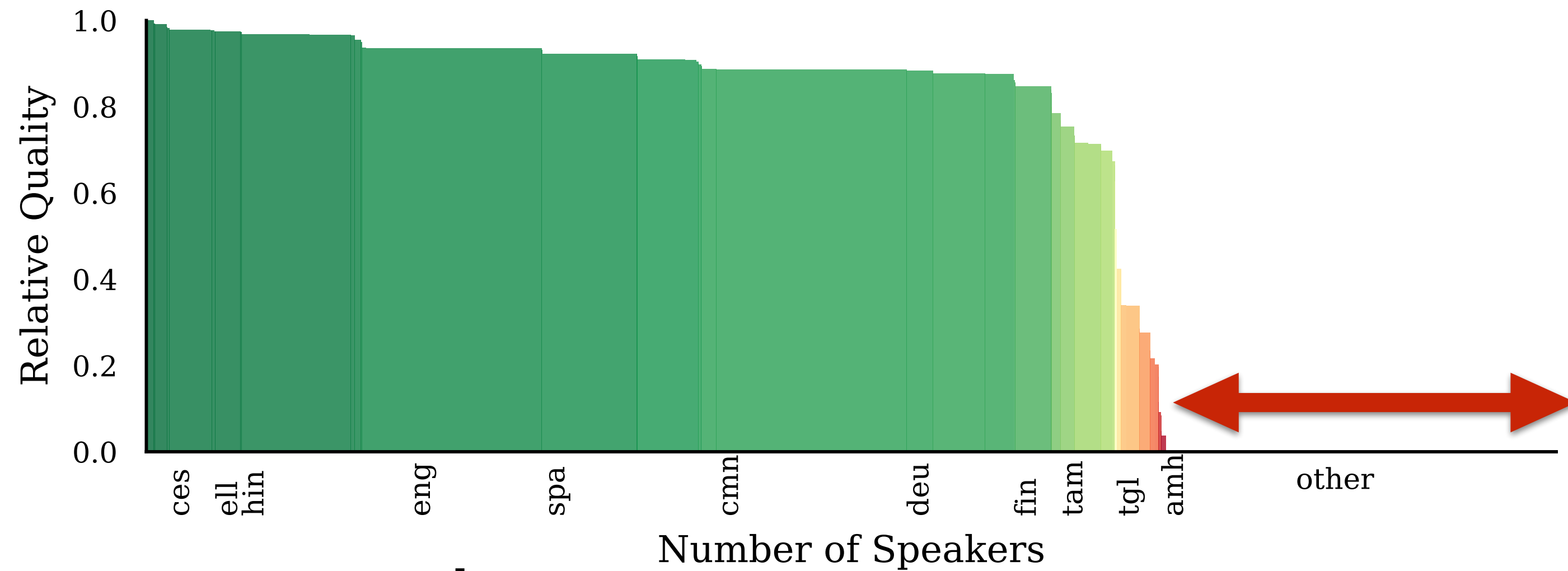"normalized demand"   "utility"

5

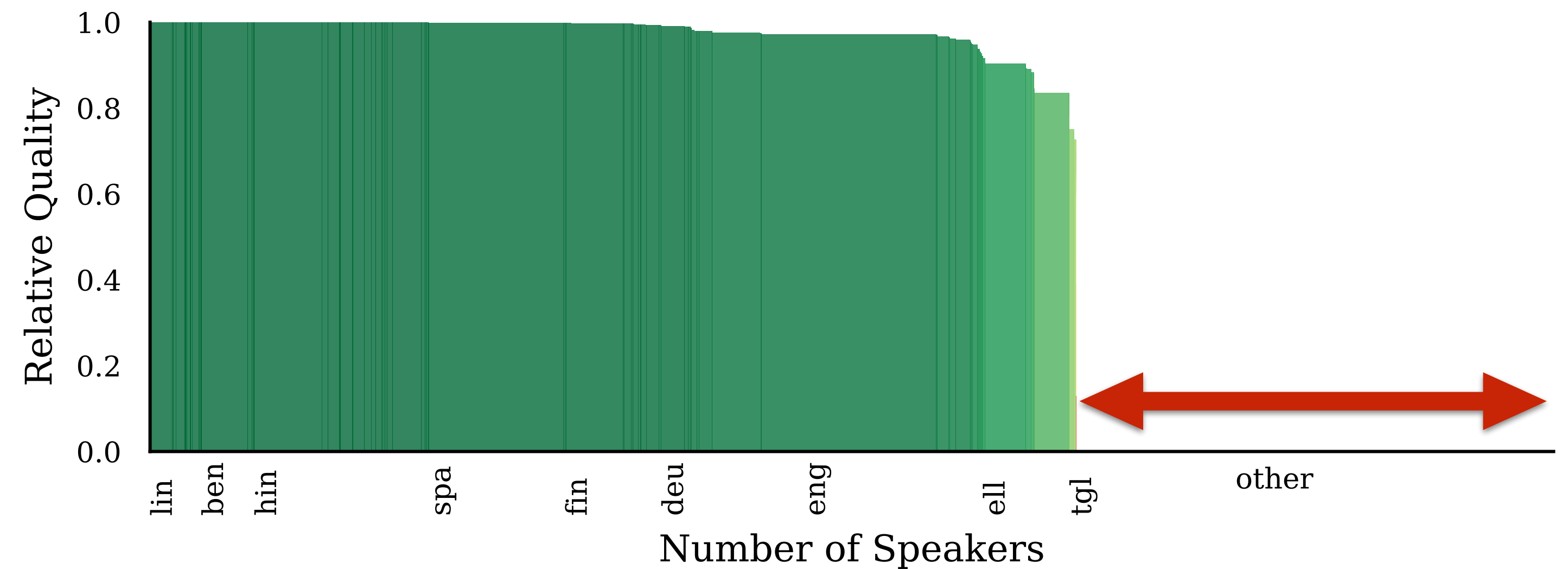# Zooming In (Analysis Tasks)



Dependency Parsing

Inflection

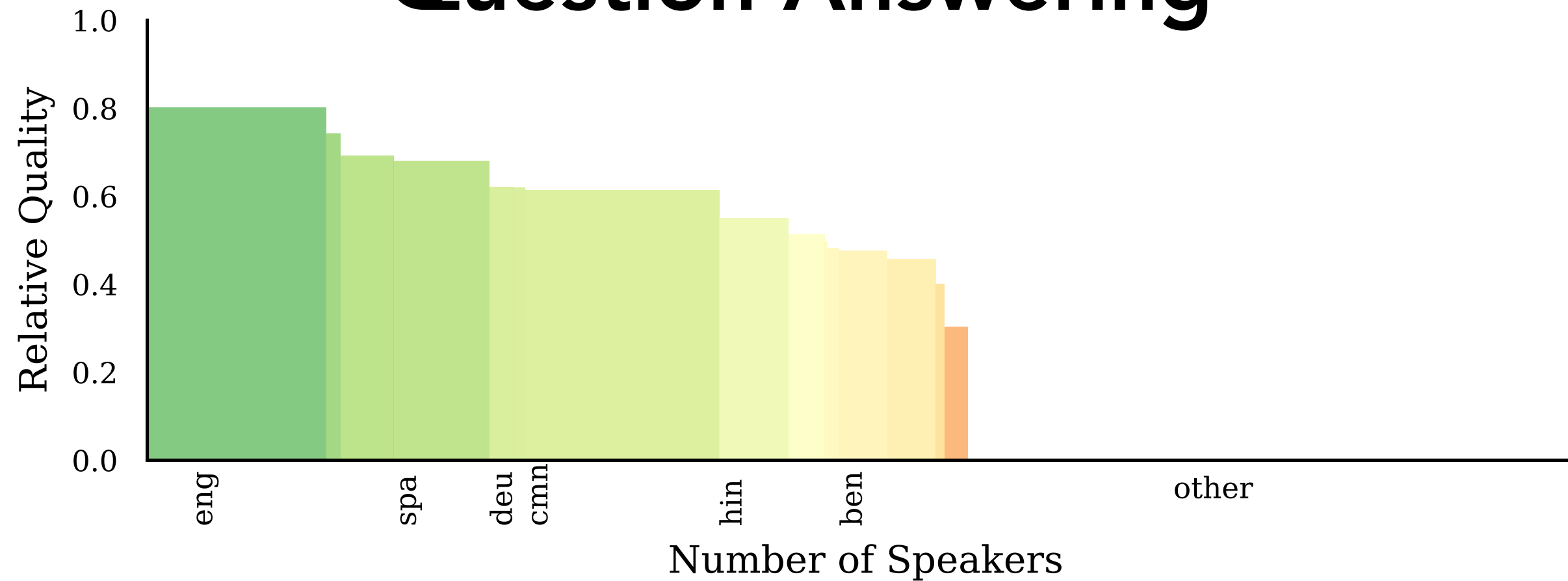# Zooming In (Analysis Tasks)

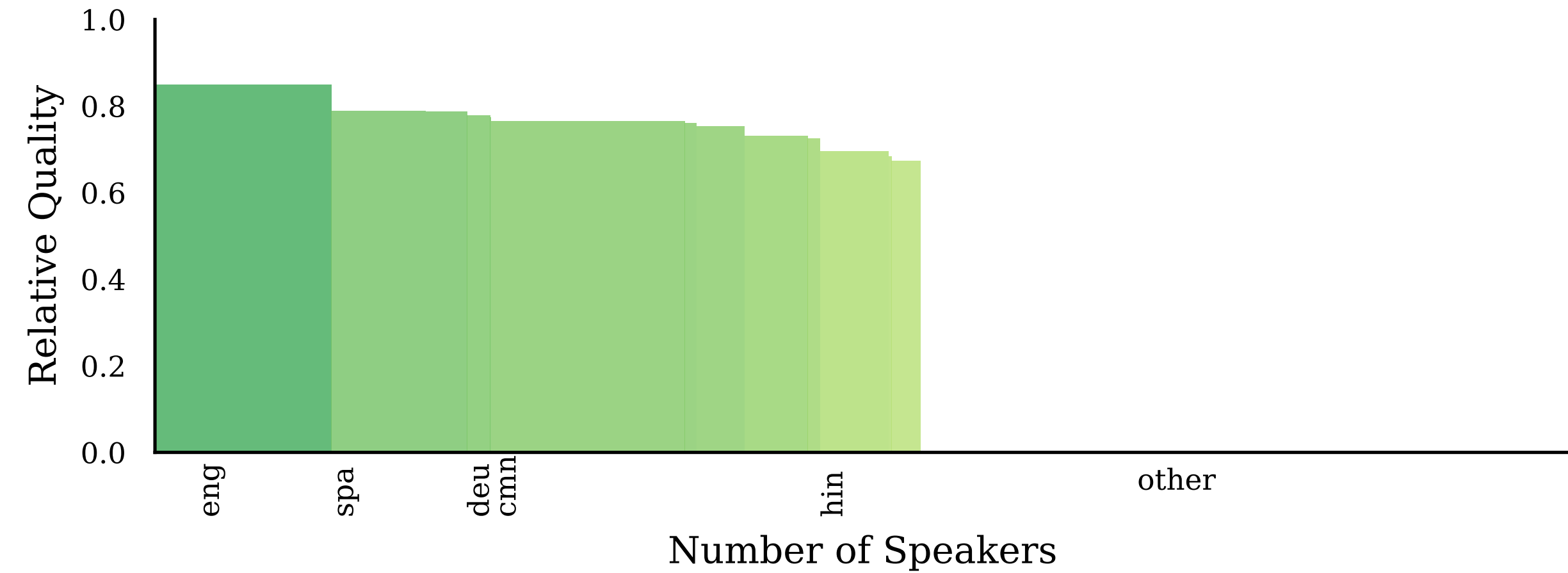

**Dependency Parsing**

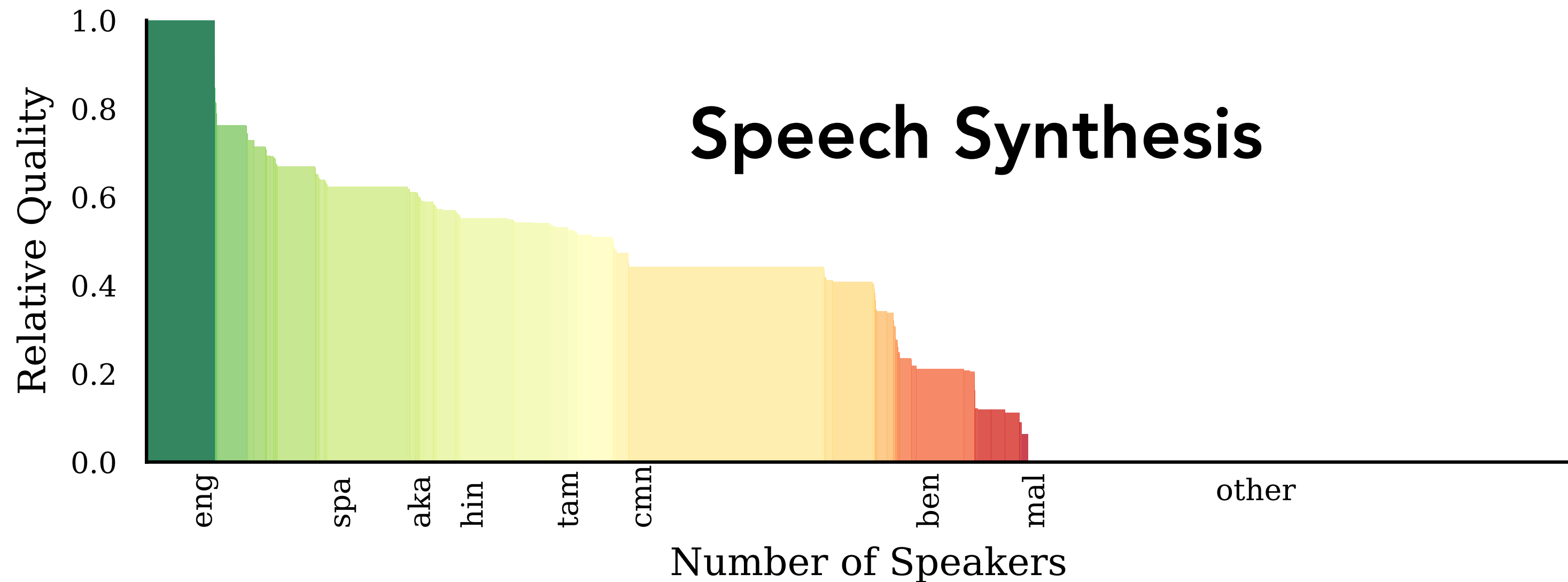**Inflection**

# Zooming In (User-facing Tasks)



**Question Answering**

**Natural Language Inference**

**Speech Synthesis**
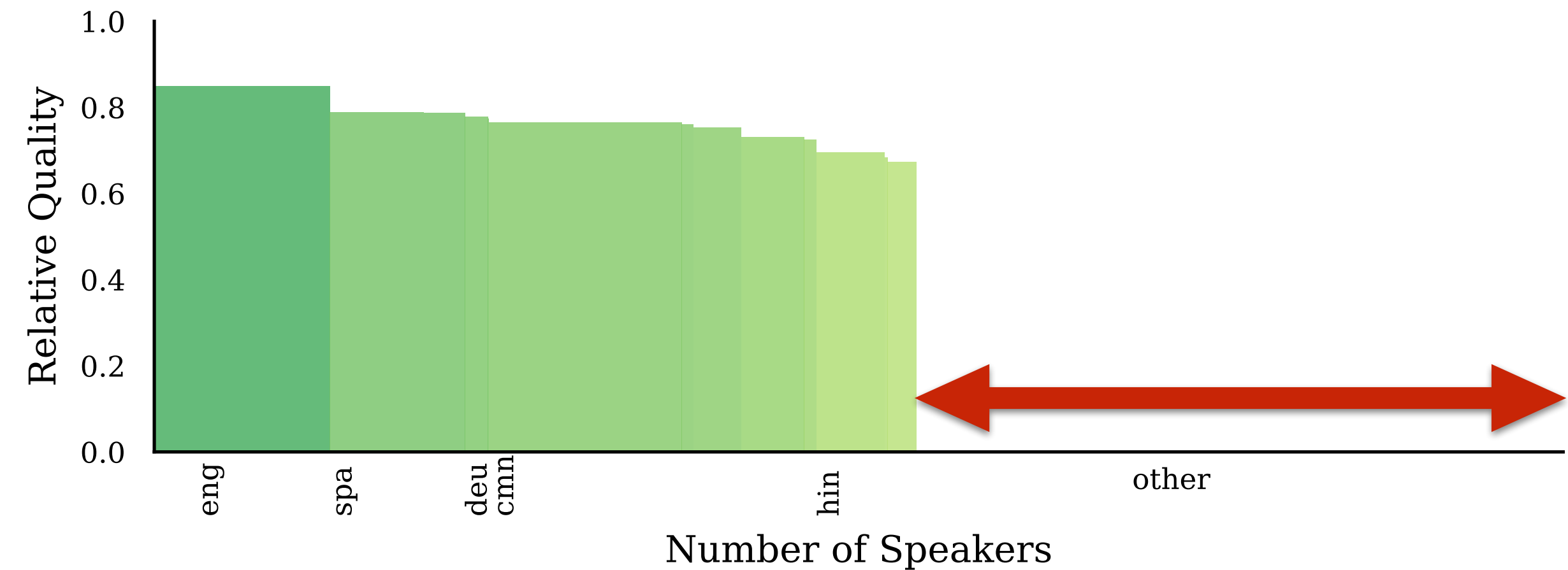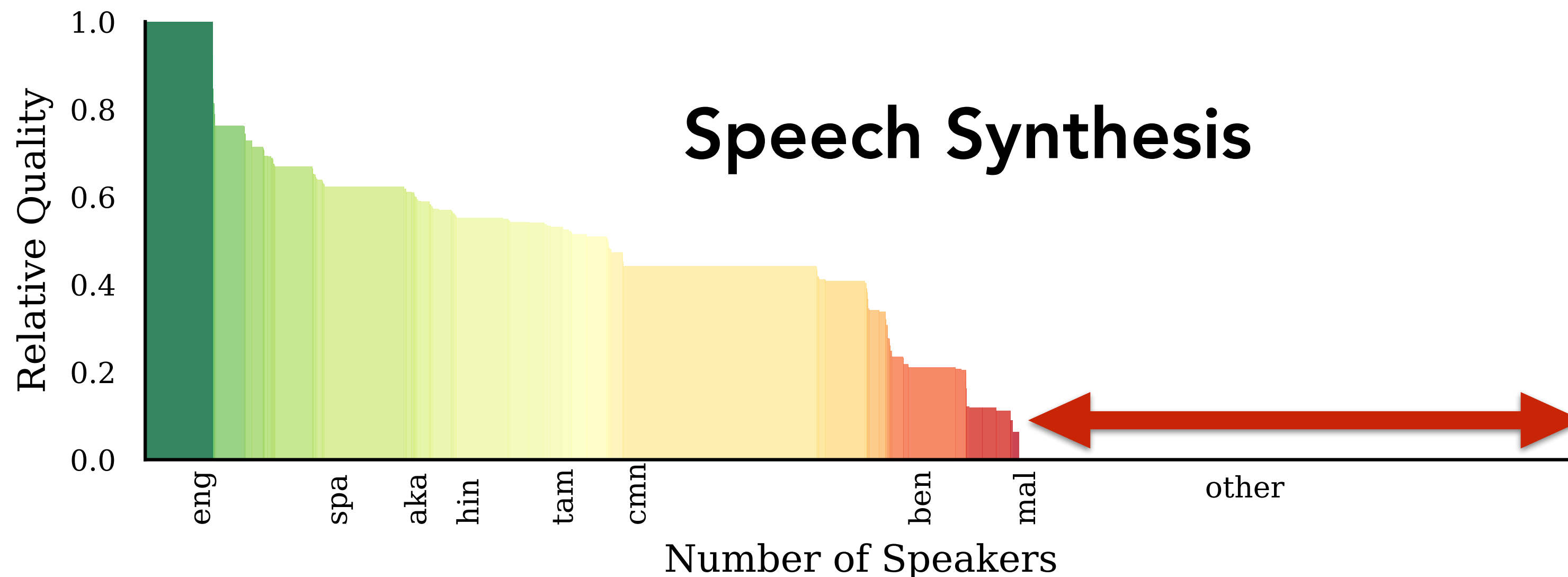
# Zooming In (User-facing Tasks)

# Zooming In (Machine Translation)



MT To English

MT To Spanish

MT To Bengali

# Going Deeper: Dialects

# Going Deeper: Dialects



Very few languages are monoliths!

# Going Deeper: Dialects



Very few languages are monoliths!

Need to model dialectal/regional/user variations.

# Going Deeper: Dialects



Very few languages are monoliths!

Need to model dialectal/regional/user variations.

Problem: most are *spoken (like 45% of all languages)*

# Going Deeper: **Dialects**



Mappa delle Lingue e Gruppi dialettali di'Italia

Very few languages are monoliths!

Need to model dialectal/regional/user variations.

Problem: most are *spoken (like 45% of all languages)*

## SD-QA: Spoken Dialectal Question Answering for the Real World

Fahim Faisal, Sharlina Keshava, Md Mahfuz ibn Alam, Antonios Anastasopoulos
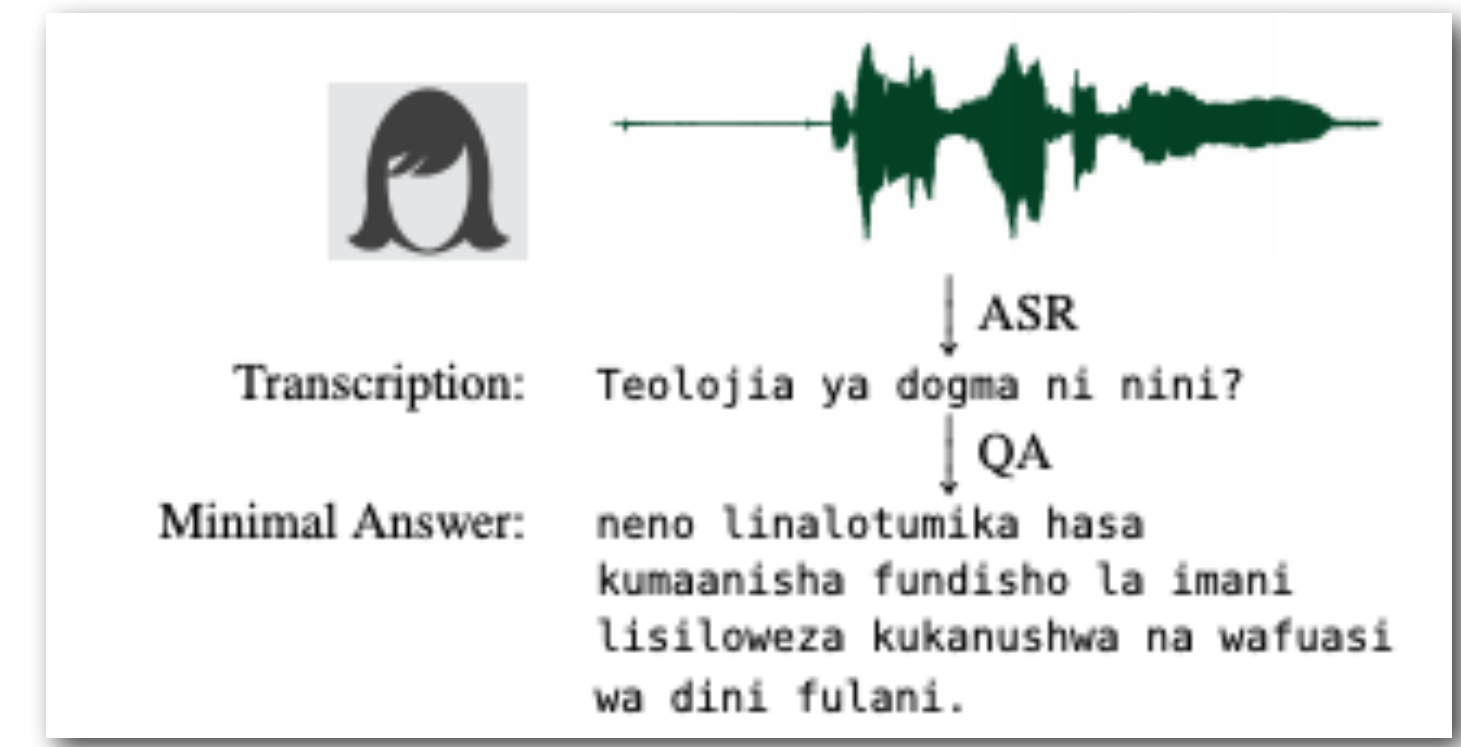Department of Computer Science, George Mason University
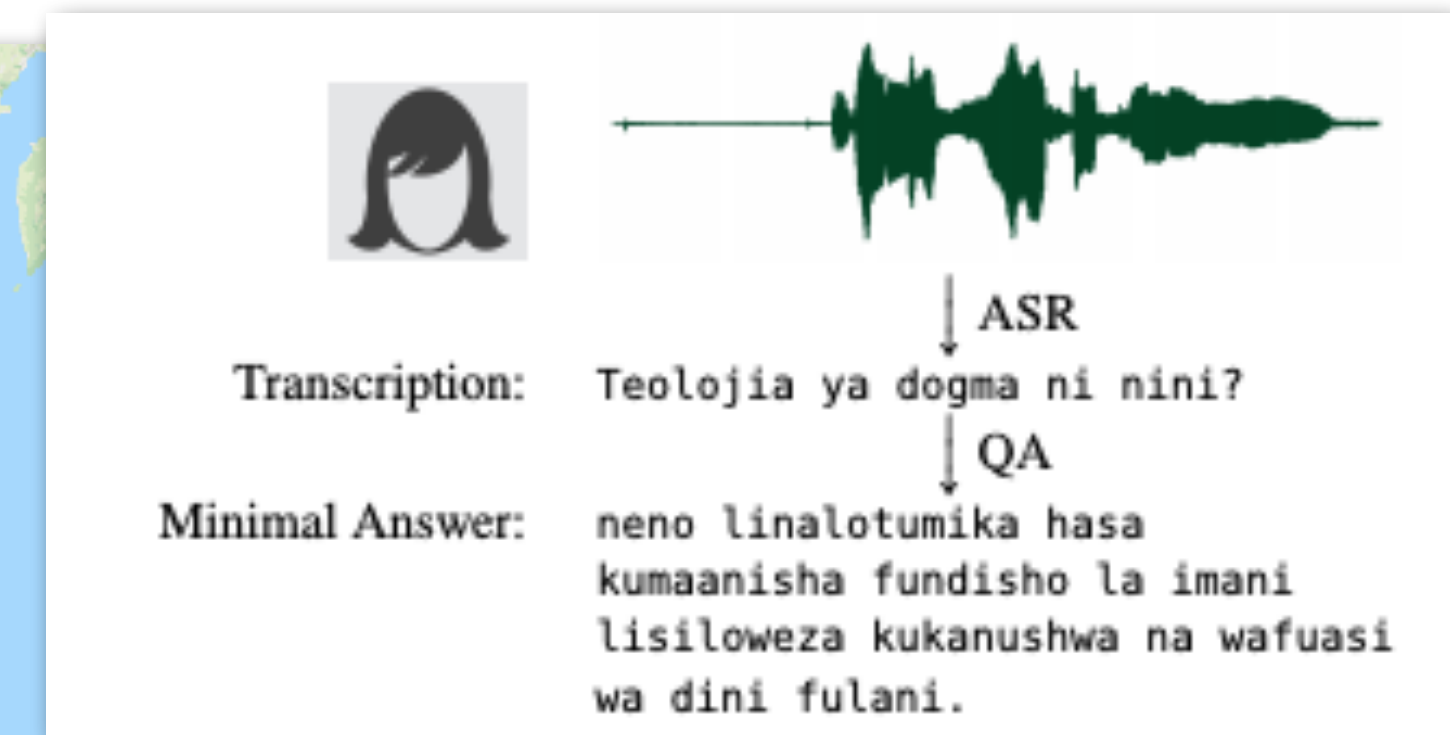{ffaisal,skeshav,malam21,antonis}@gmu.edu

(EMNLP Findings 2021)

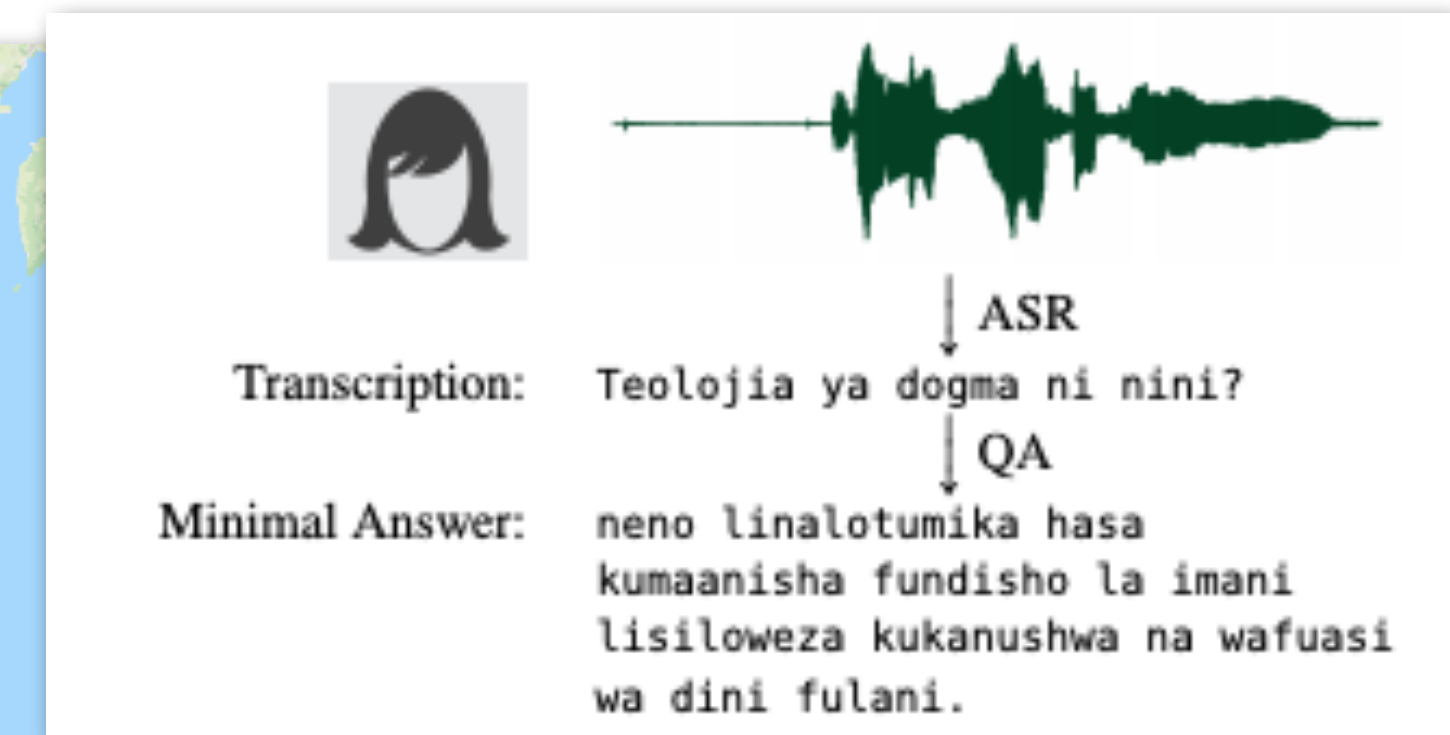# SD-QA: Spoken, Dialectal, Multilingual QA

# SD-QA: Spoken, Dialectal, Multilingual QA

# SD-QA: Spoken, Dialectal, Multilingual QA

# SD-QA: Spoken, Dialectal, Multilingual QA



| Transc. | Arabic Variety | | | | | | | Avg. |
|---------|-----|------|------|------|------|------|------|------|
| | DZA | BHR | EGY | JOR | MAR | SAU | TUN | |
| Gold | 51.3 | | | | | | | |
| ara-VAR | 51.3 | 46.2 | 44.8 | 45.6 | 47.6 | 46.3 | 46.7 | 46.9 |

Transcription: Teolojia ya dogma ni nini?

Minimal Answer: neno linalotumika hasa kumaanisha fundisho la imani lisiloweza kukanushwa na wafuasi wa dini fulani.

# SD-QA: Spoken, Dialectal, Multilingual QA

Let's make a plan

NLP beyond the top-100 languages

GEORGE MASON UNIVERSITY

# Going beyond the top-100 languages

# Going beyond the top-100 languages

# Going beyond the top-100 languages

# Are all unseen languages equally hard?

When Being Unseen from mBERT is just the Beginning:
Handling New Languages With Multilingual Language Models

Benjamin Muller[†]    Antonis Anastasopoulos[‡]    Benoît Sagot[†]    Djamé Seddah[†]
[†]Inria, Paris, France
[‡]Department of Computer Science, George Mason University, USA
firstname.lastname@inria.fr    antonis@gmu.edu

(NAACL 2021)

https://github.com/benjamin-mlr/mbert-unseen-languages.git

# Are all unseen languages hard?

# Are all unseen languages hard?

Some are "easy"

# Are all unseen languages hard?

Some are "easy"

> Similar languages in pre-training + same script

> e.g. Faroese, Swiss German

# Are all unseen languages hard?

Some are "easy"

    Similar languages in pre-training +
same script

    e.g. Faroese, Swiss German

Some are "intermediate"

# Are all unseen languages hard?

Some are "easy"

    Similar languages in pre-training + same script

    e.g. Faroese, Swiss German

Some are "intermediate"

    Simple approach (continued fine-tuning) leads to good results

    e.g. Maltese, Bambara, Wolof

# Are all unseen languages hard?

Some are "easy"

 Similar languages in pre-training + same script

 e.g. Faroese, Swiss German

Some are "intermediate"

 Simple approach (continued fine-tuning) leads to good results

 e.g. Maltese, Bambara, Wolof

Some seem "hard"

# Are all unseen languages hard?

Some are "easy"

Similar languages in pre-training + same script

e.g. Faroese, Swiss German

Some are "intermediate"

Simple approach (continued fine-tuning) leads to good results

e.g. Maltese, Bambara, Wolof

Some seem "hard"

Hard == Static monolingual embeddings >> mBERT adaptation

# Are all unseen languages hard?

Some are "easy"

    Similar languages in pre-training + same script

    e.g. Faroese, Swiss German

Some are "intermediate"

    Simple approach (continued fine-tuning) leads to good results

    e.g. Maltese, Bambara, Wolof

Some seem "hard"

    Hard == Static monolingual embeddings >> mBERT adaptation

    Similar languages in pre-training BUT different script

NLP
GEORGE
MASON

# Are all unseen languages hard?

Some are "easy"

    Similar languages in pre-training + same script

    e.g. Faroese, Swiss German

Some are "intermediate"

    Simple approach (continued fine-tuning) leads to good results

    e.g. Maltese, Bambara, Wolof

Some seem "hard"

    Hard == Static monolingual embeddings >> mBERT adaptation

    Similar languages in pre-training BUT different script

    e.g. Uyghur, Sorani, Sindhi

# Are all unseen languages hard?

Some are "easy"

Similar languages in pre-training + same script

e.g. Faroese, Swiss German

Some are "intermediate"

Simple approach (continued fine-tuning) leads to good results

e.g. Maltese, Bambara, Wolof

Some seem "hard"

Hard == Static monolingual embeddings >> mBERT adaptation

Similar languages in pre-training BUT different script

e.g. Uyghur, Sorani, Sindhi

Transliteration helps

# Doing better by hard-coding linguistic information

## Phylogeny-Inspired Adaptation of Multilingual Models to New Languages

**Fahim Faisal, Antonios Anastasopoulos**

Department of Computer Science, George Mason University

{ffaisal,antonis}@gmu.edu

(AACL 2022)

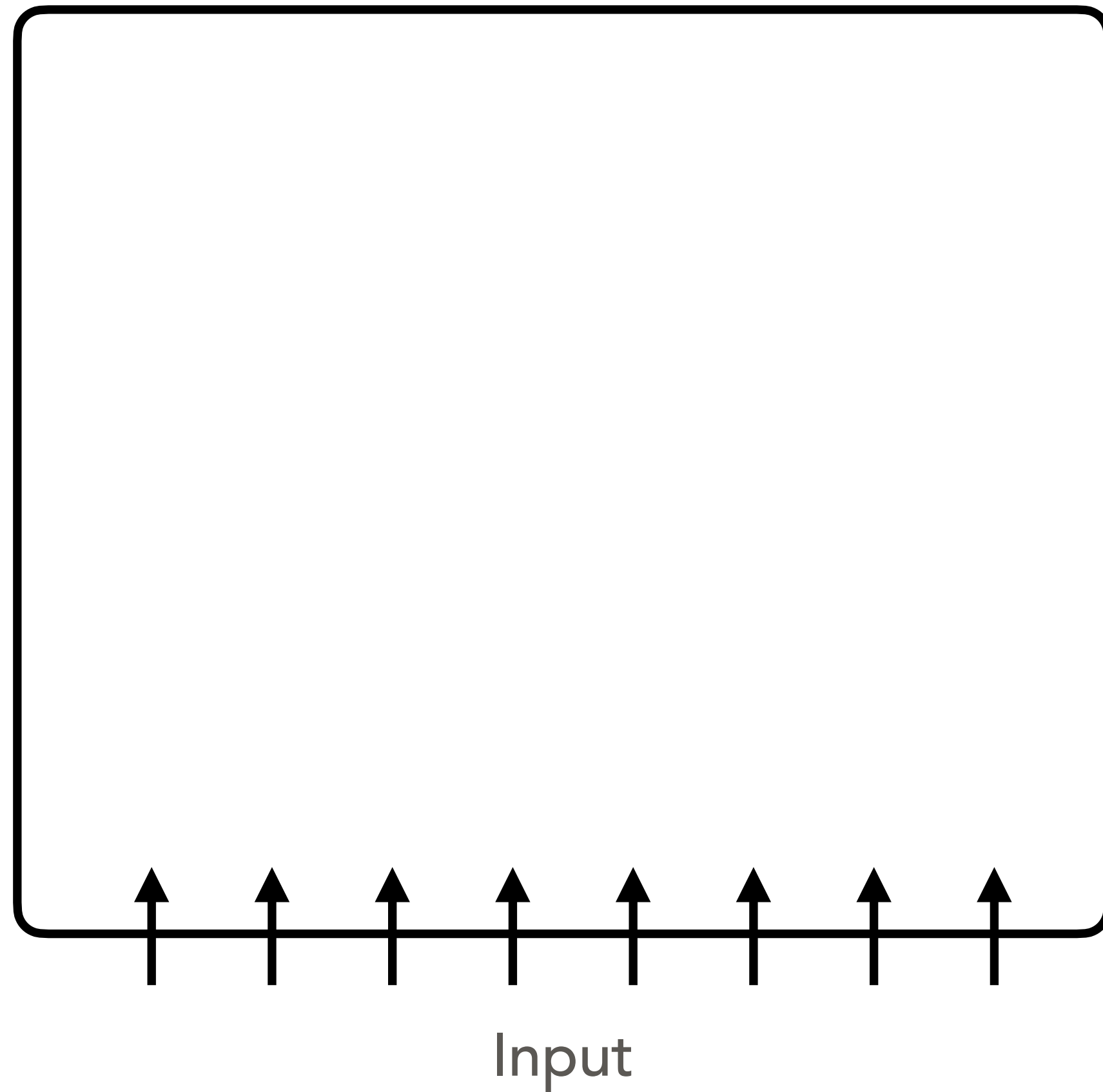https://github.com/ffaisal93/adapt_lang_phylogeny
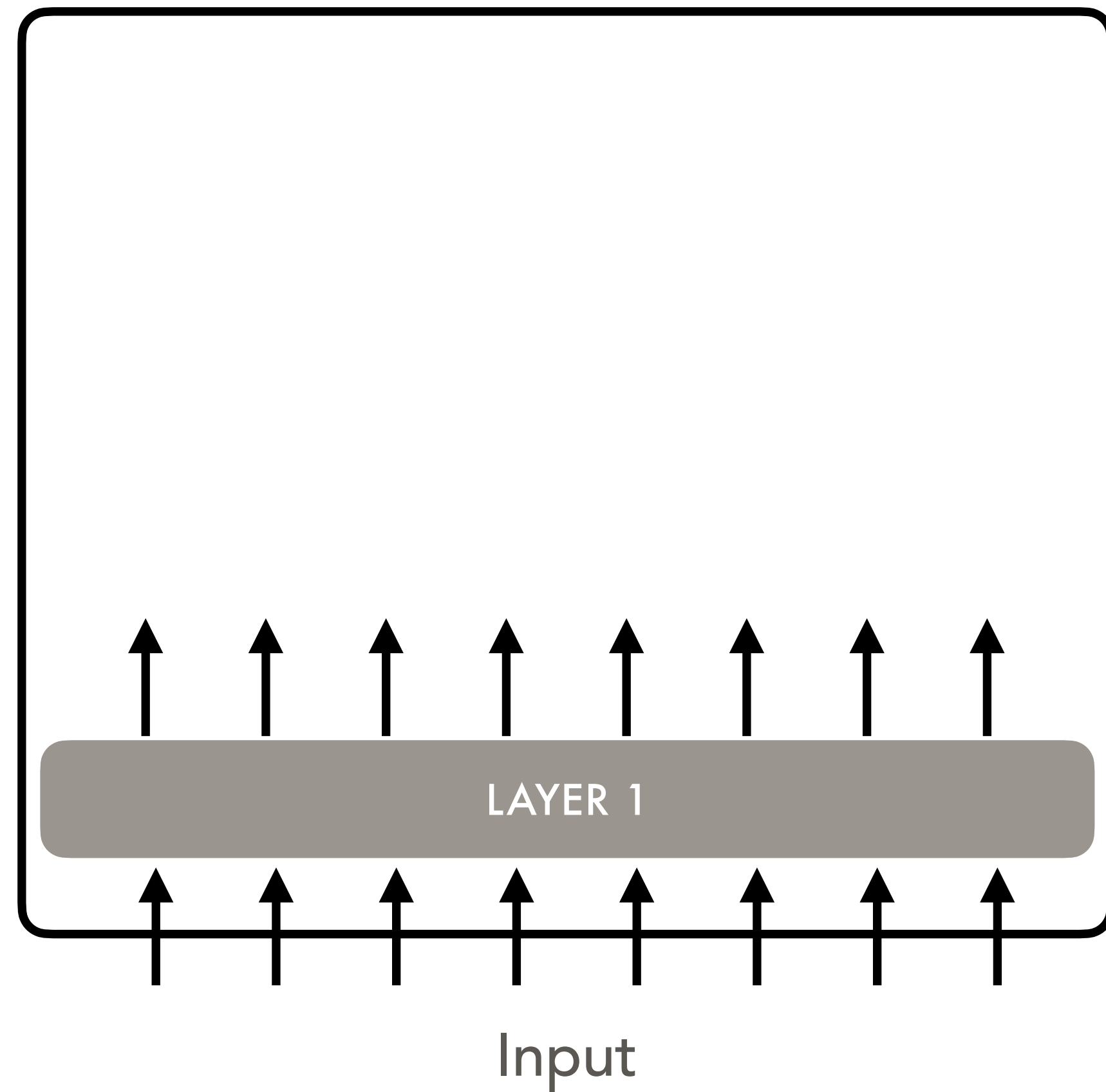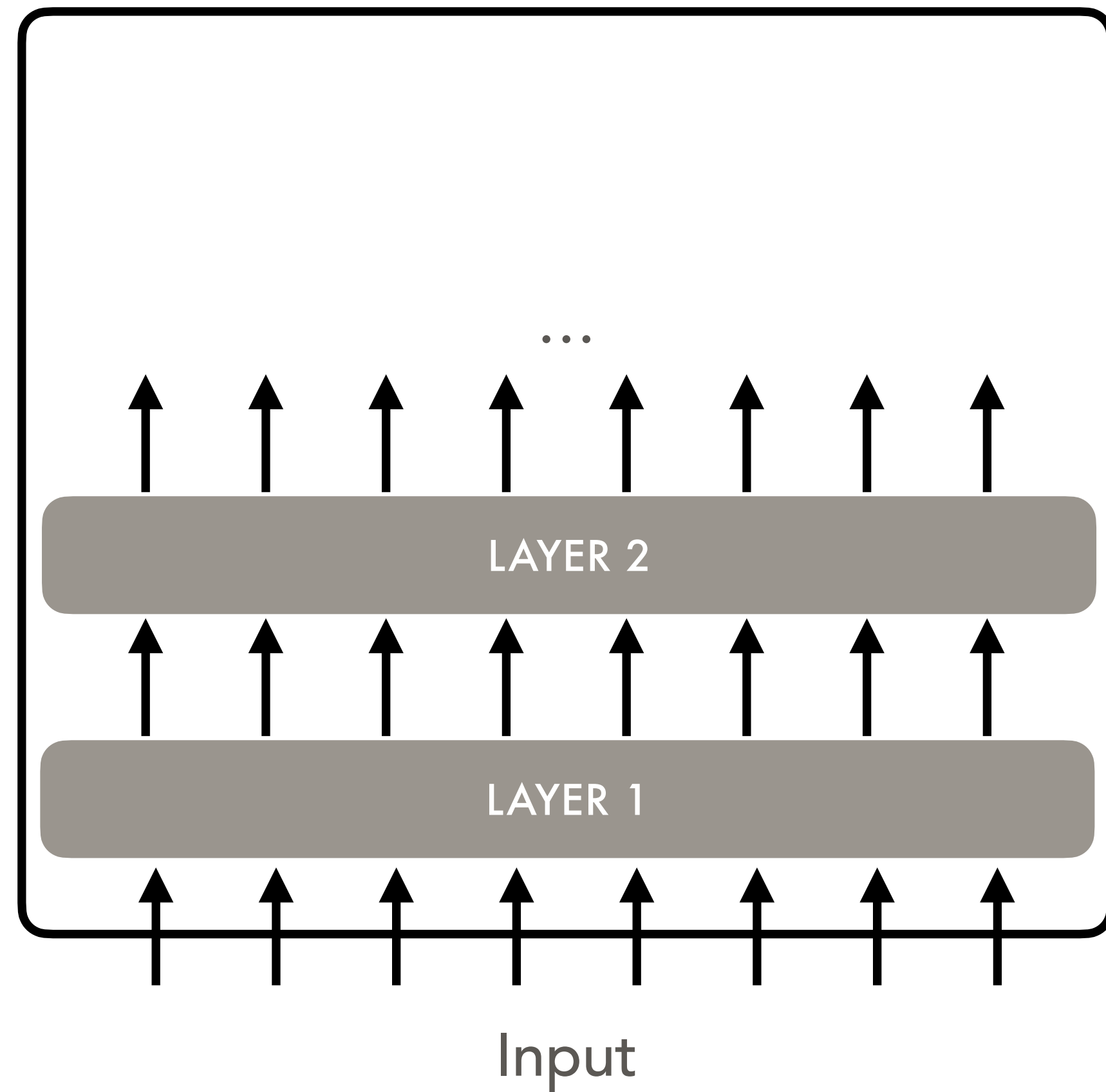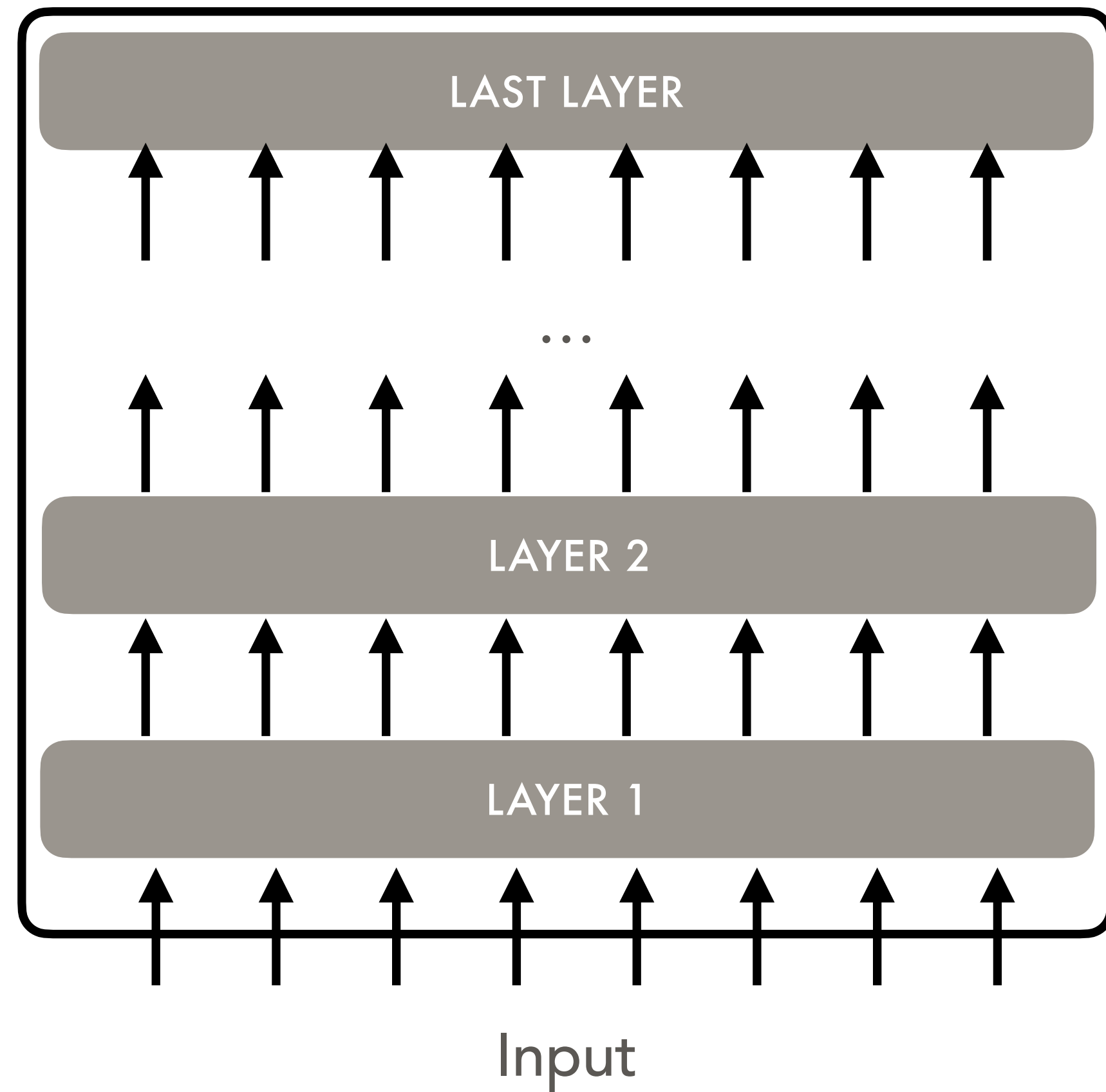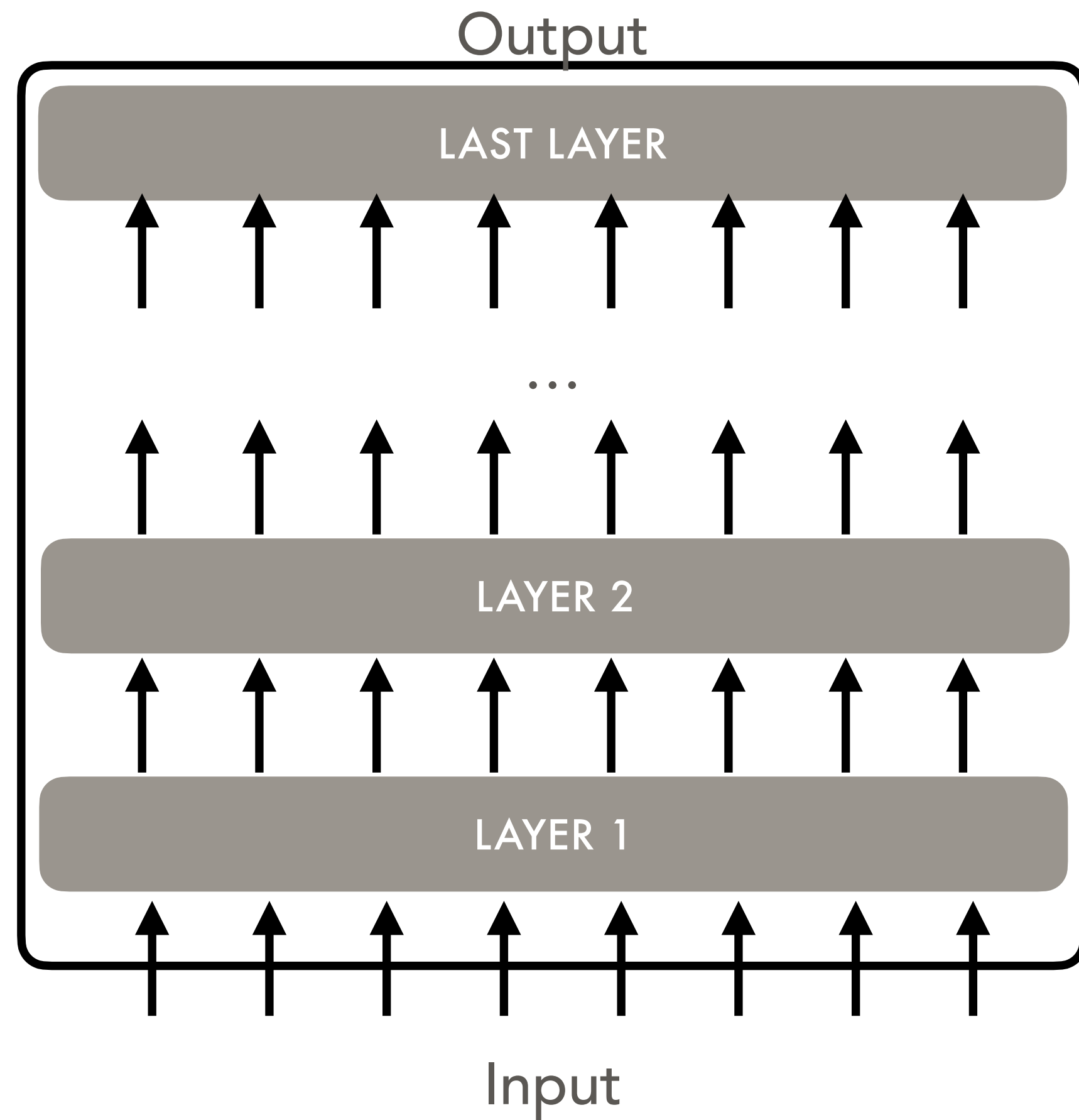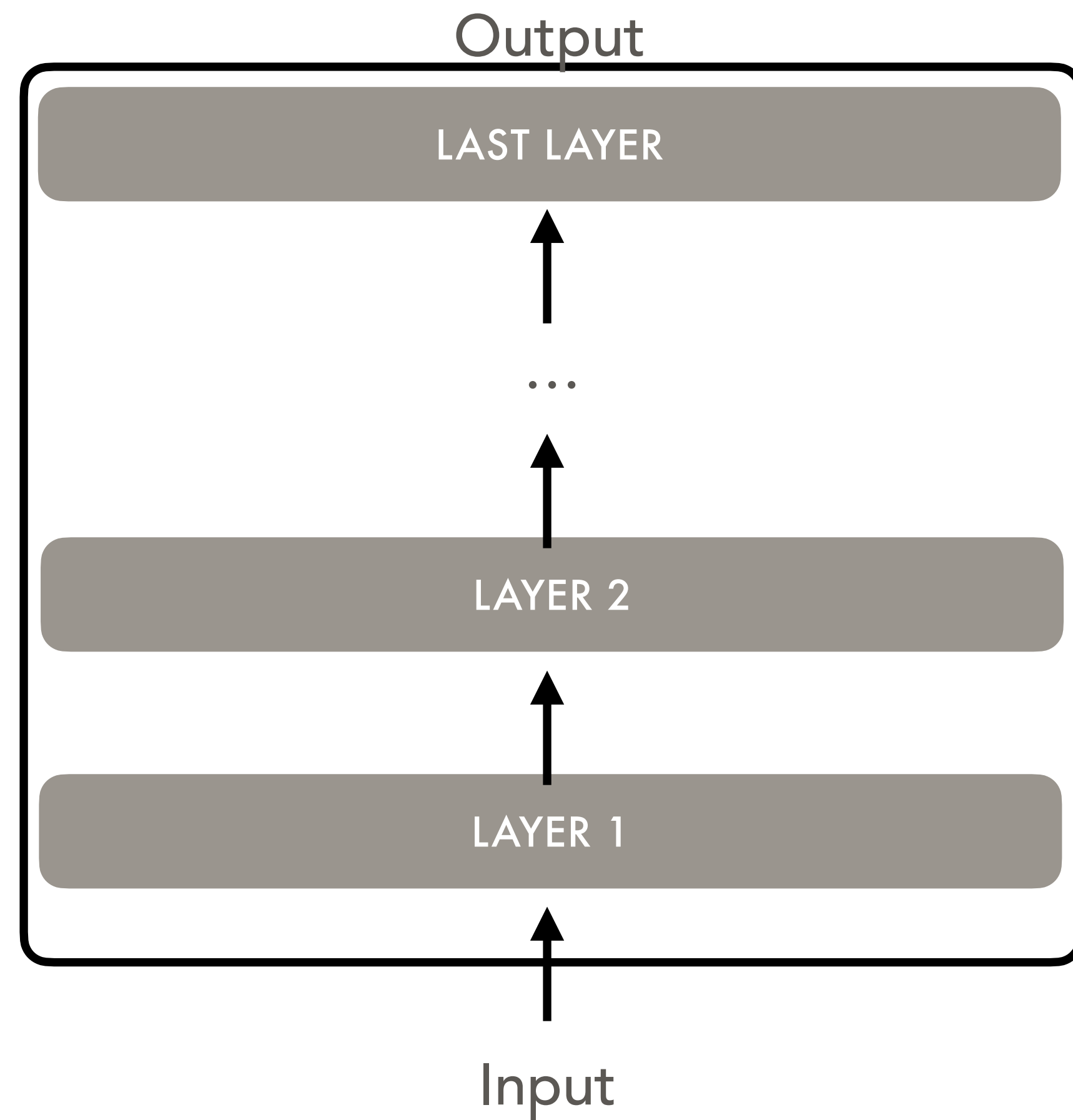
# Revisiting Adapters

# Revisiting Adapters

# Revisiting Adapters



Input

# Revisiting Adapters



Input

# Revisiting Adapters

# Revisiting Adapters

# Revisiting Adapters

# Revisiting Adapters



Output

| LAST LAYER |
| :---: |

…

| LAYER 2 |
| :---: |

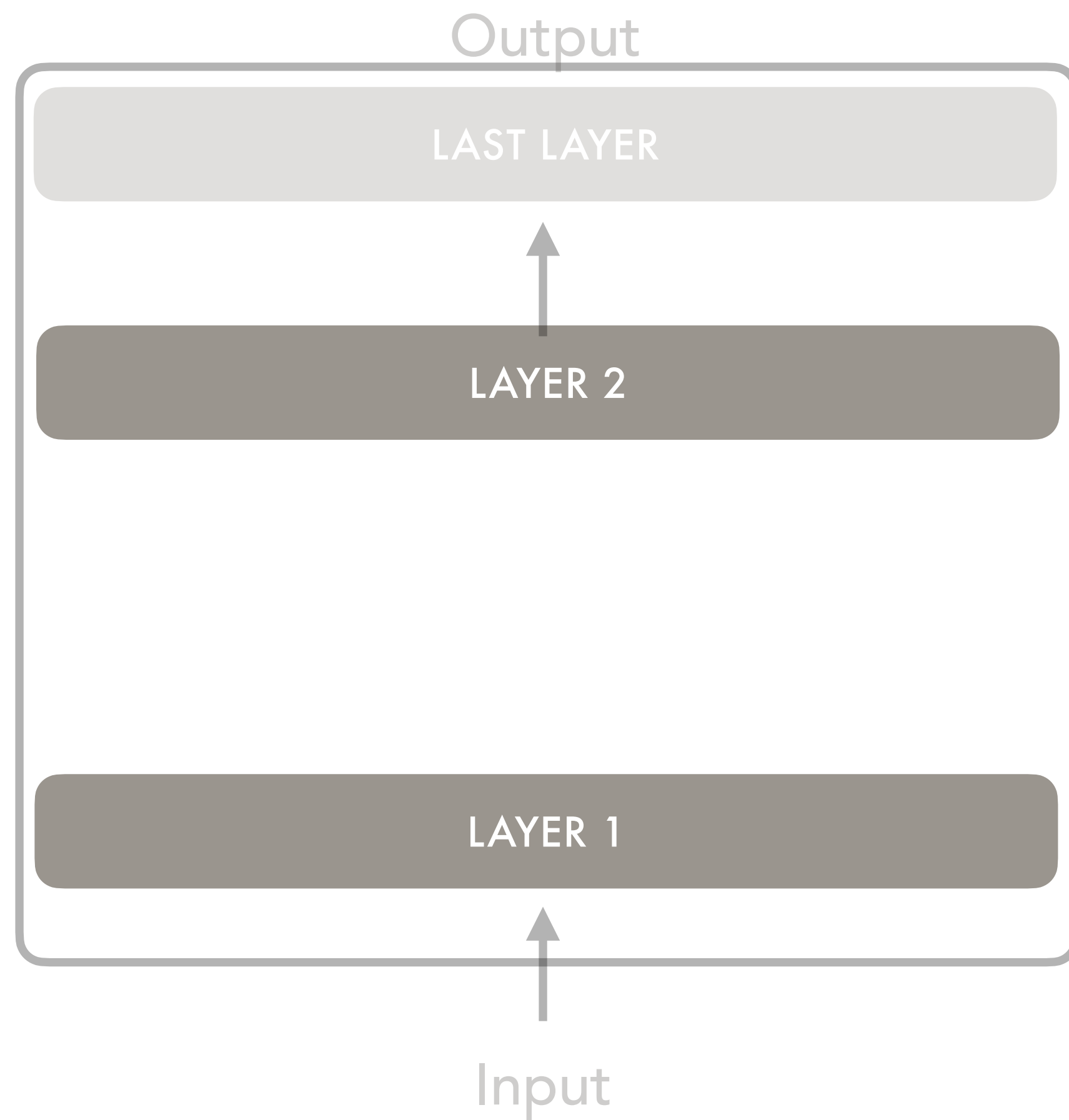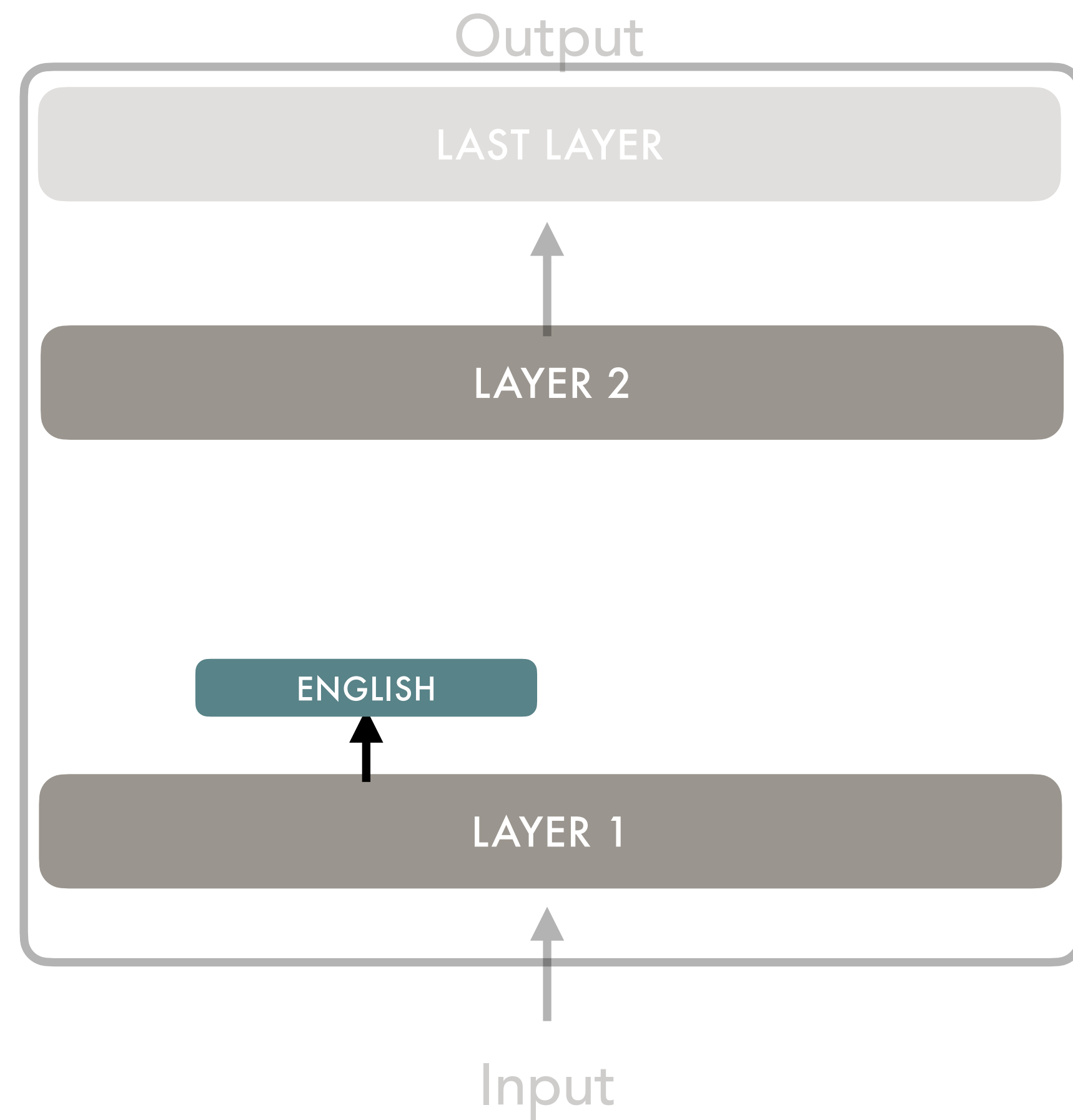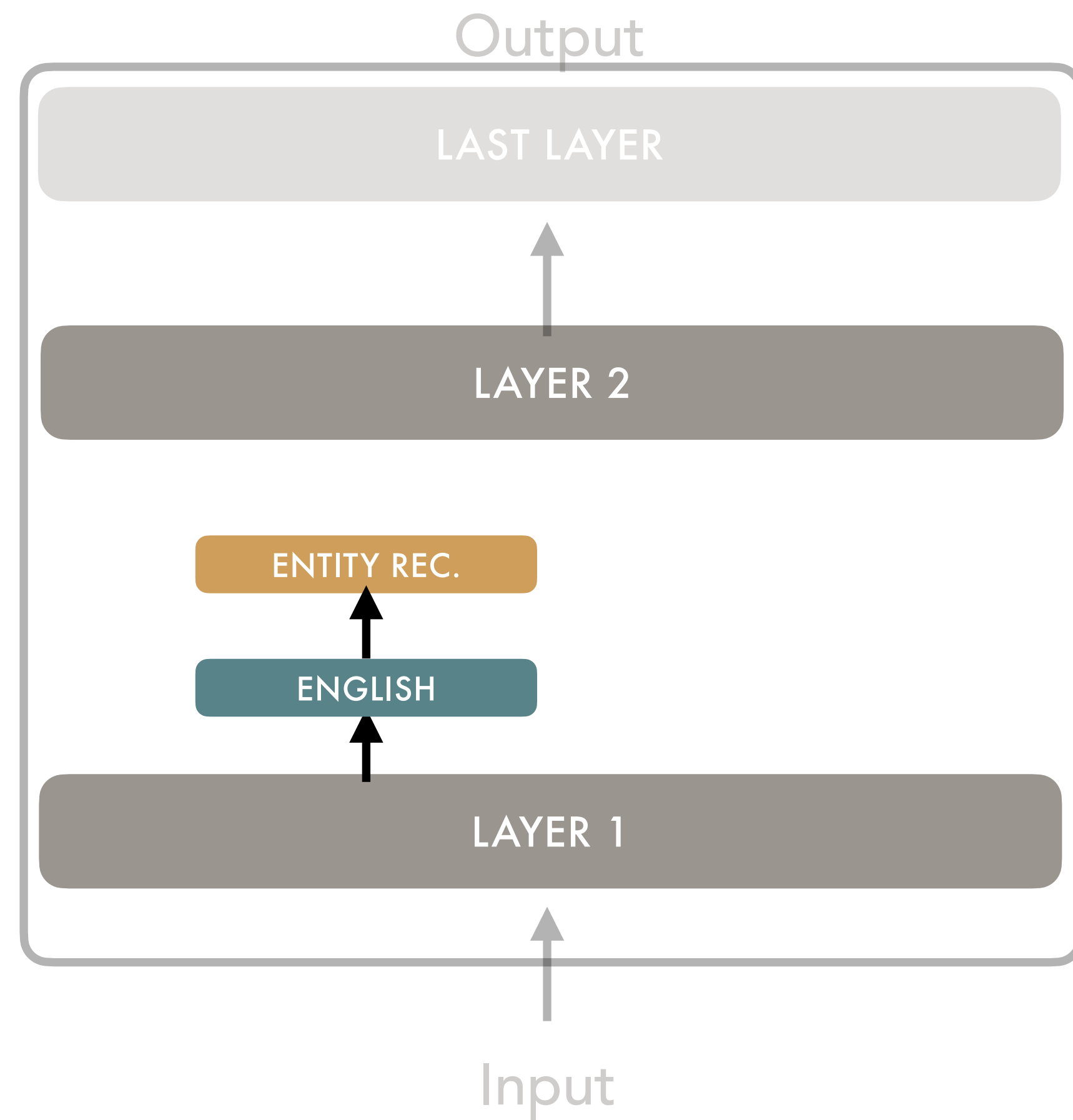| LAYER 1 |
| :---: |

Input

NLP GEORGE MASON
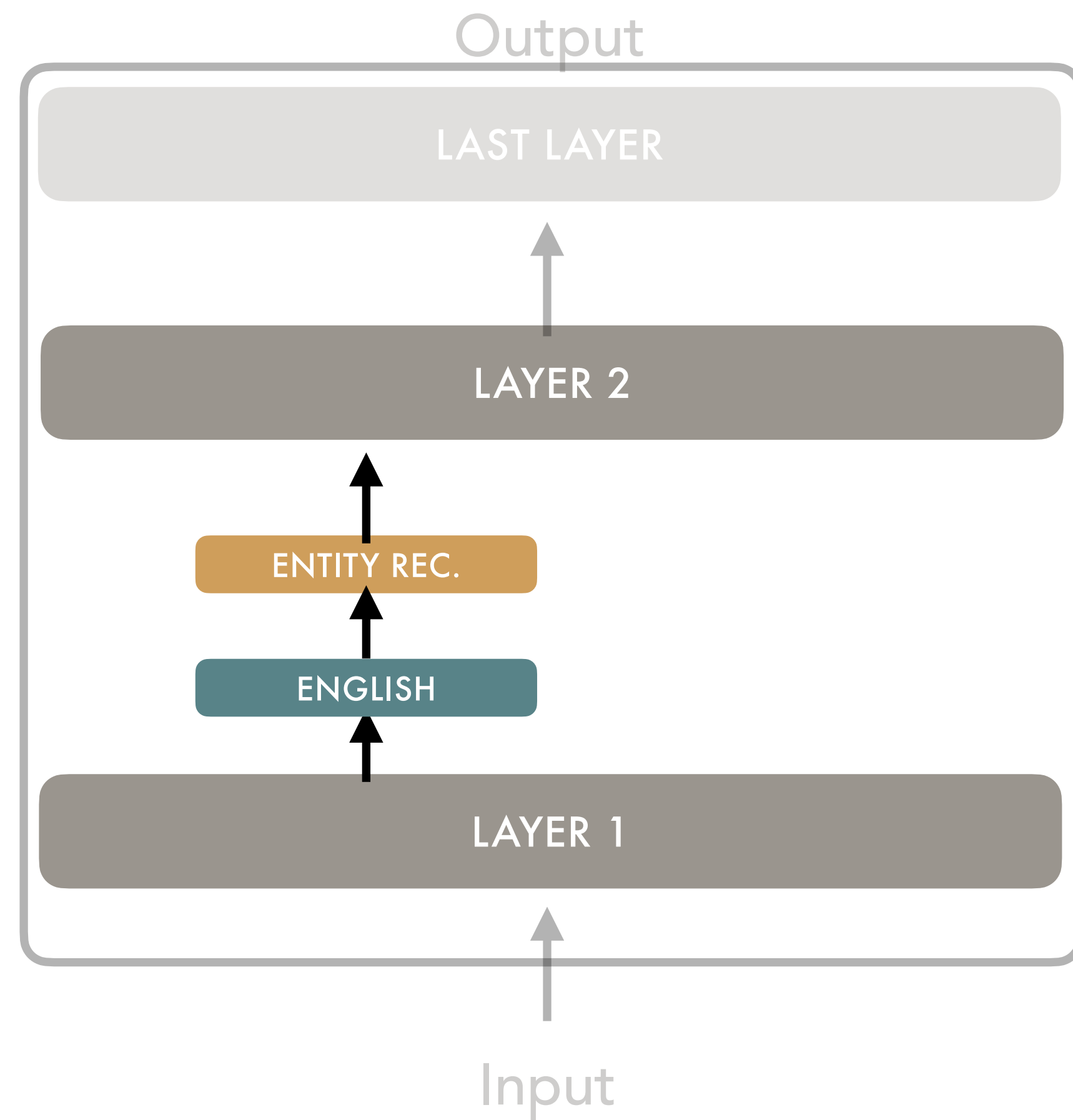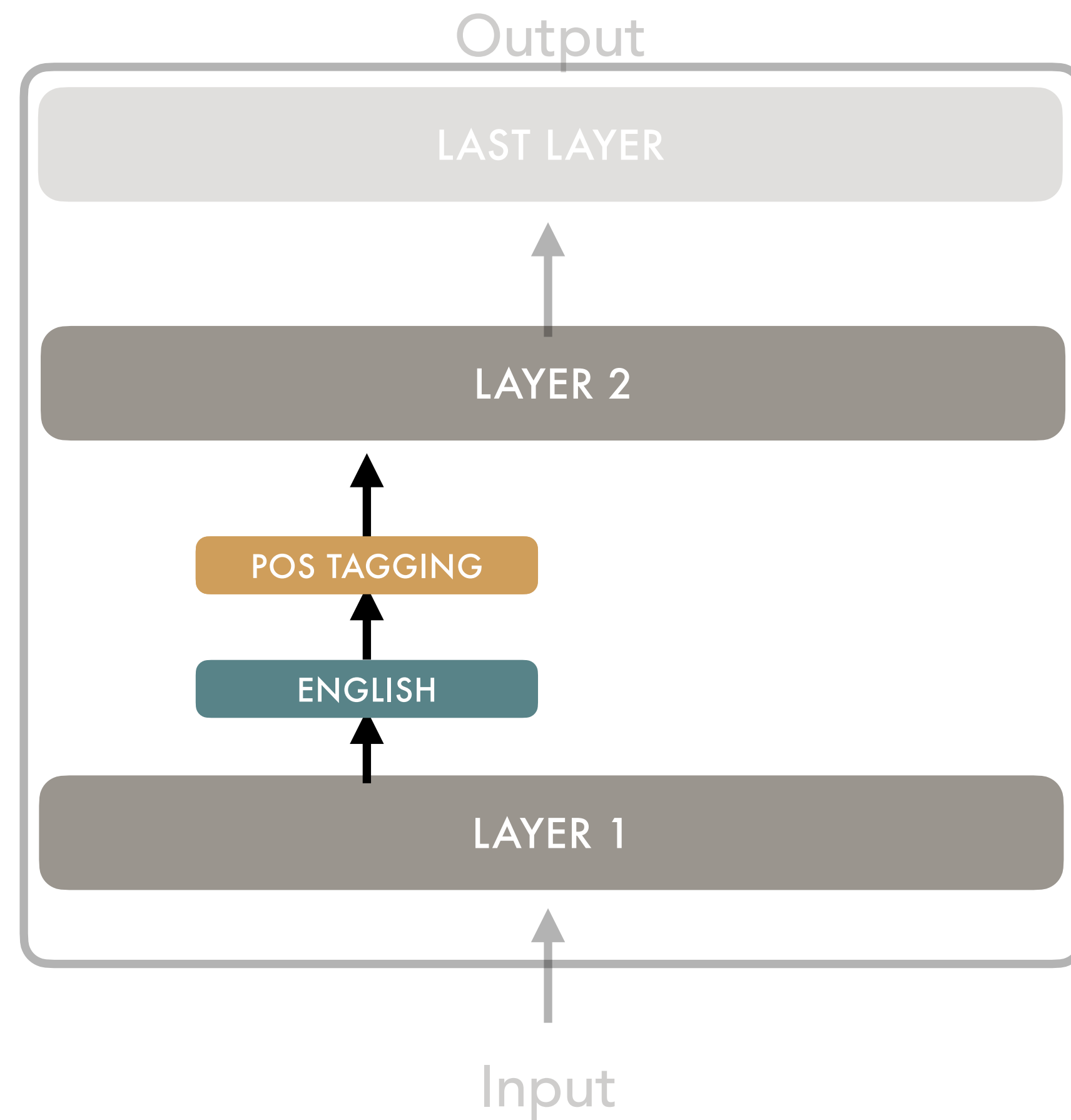
# Revisiting Adapters

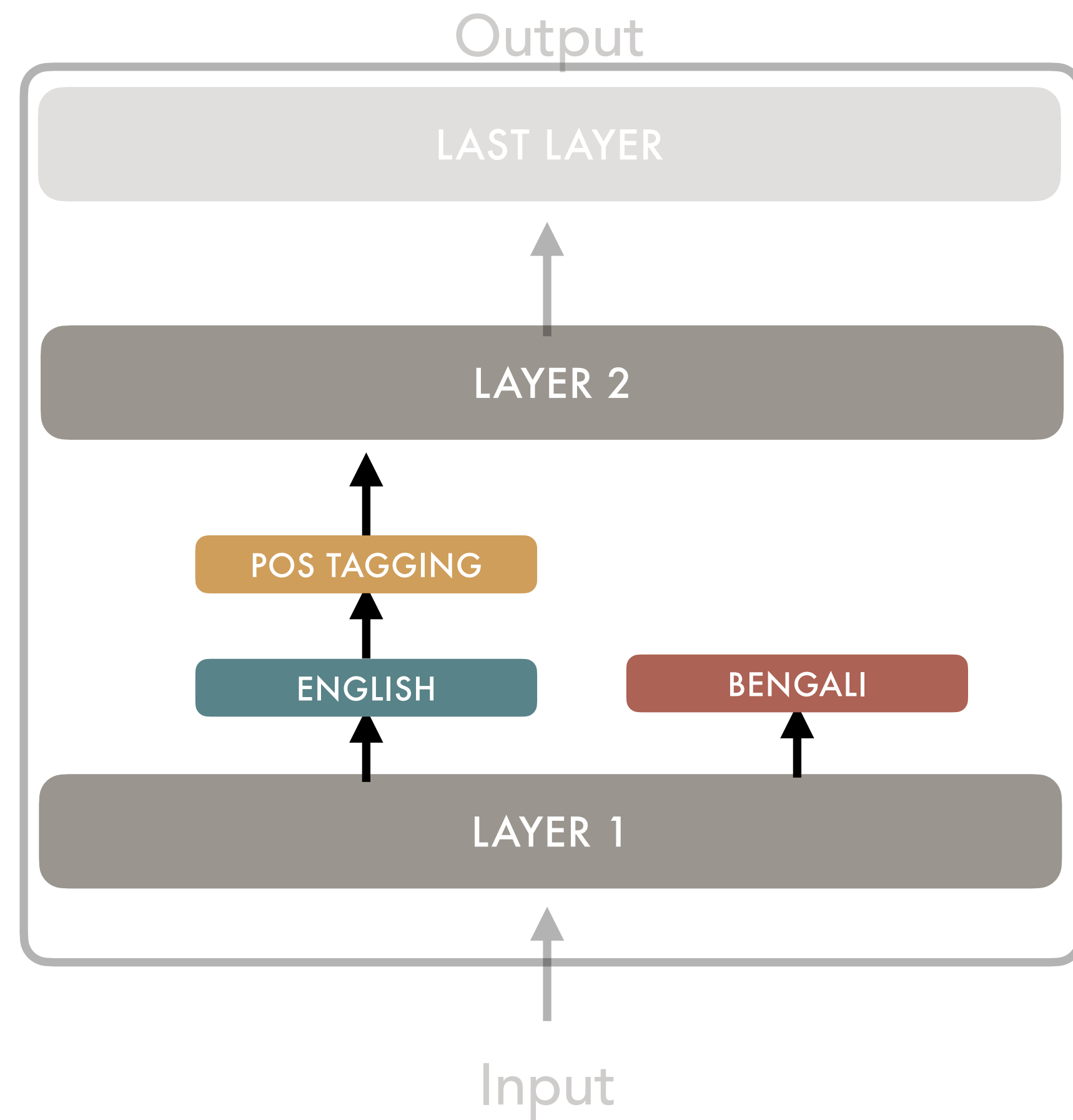# Revisiting Adapters

# Revisiting Adapters

# Revisiting Adapters

# Revisiting Adapters

# Revisiting Adapters

# Revisiting Adapters

# Revisiting Adapters

# Revisiting Adapters



Easy zero-shot adaptation to new languages at a low cost (additional parameters)

# Revisiting Adapters



Easy zero-shot adaptation to new languages at a low cost (additional parameters)

Avoids catastrophic forgetting

# Revisiting Adapters



Easy zero-shot adaptation to new languages at a low cost (additional parameters)

Avoids catastrophic forgetting

Performance comparable to full-model fine-tuning

# Revisiting Adapters



Easy zero-shot adaptation to new languages at a low cost (additional parameters)

Avoids catastrophic forgetting

Performance comparable to full-model fine-tuning

**Can we do better?**

# Follow Phylogeny for Parameter Sharing

# Follow Phylogeny for Parameter Sharing

For Dutch input

# Follow Phylogeny for Parameter Sharing

For Bengali input

# Results

# Results

## DEPENDENCY PARSING



Chart with y-axis labeled "UAS" ranging from 0 to 80. X-axis shows [T], [LT], [FGLT] for GERMANIC (12) and [T], [LT], [FGLT] for URALIC (11).

# Results on unseen languages

# Results on unseen languages
## DEPEDENCY PARSING

# Results on unseen languages
## DEPEDENCY PARSING



Much larger improvements for *new, unseen* languages

# Ablations

# Ablations

## DEPENDENCY PARSING ON URALIC LANGS

# Ablations

## DEPENDENCY PARSING ON URALIC LANGS



Even constraining to the same number of parameters, still improvements!
Is it language sharing or network depth?

# Ablations

Even constraining to the same number of parameters, still improvements!
Is it language sharing or network depth?

Same idea applied to Translation:

   - 2nd best constrained system at WMT Shared Task on Large-Scale Multilingual Systems for African Languages!

# Going forward and beyond

No matter what, we need data in these languages.

What data do we need, though?

# Few-Shot is the way

Let's leave script issues aside for a minute

    —since we can find solutions, e.g.

**MAD-X: An Adapter-Based Framework for Multi-Task Cross-Lingual Transfer**

Jonas Pfeiffer[1], Ivan Vulić[2], Iryna Gurevych[1], Sebastian Ruder[3]
[1]Ubiquitous Knowledge Processing Lab, Technical University of Darmstadt
[2]Language Technology Lab, University of Cambridge
[3]DeepMind

**CANINE: Pre-training an Efficient Tokenization-Free Encoder for Language Representation**

Jonathan H. Clark, Dan Garrette, Iulia Turc, John Wieting
Google Research

**Parsing with Multilingual BERT, a Small Corpus, and a Small Treebank**

Ethan C. Chau[†◇]    Lucy H. Lin[†]    Noah A. Smith[†*]
[†]Paul G. Allen School of Computer Science & Engineering, University of Washington
[◇]Department of Linguistics, University of Washington
[*]Allen Institute for Artificial Intelligence

**When Being Unseen from mBERT is just the Beginning: Handling New Languages With Multilingual Language Models**

Benjamin Muller[†]    Antonis Anastasopoulos[‡]    Benoît Sagot[†]    Djamé Seddah[†]
[†]Inria, Paris, France
[‡]Department of Computer Science, George Mason University, USA
firstname.lastname@inria.fr   antonis@gmu.edu

It seems that we can do ~~great~~ well with just a few in-domain in-language task data + data augmentation!

28

# Going forward and beyond

# Going forward and beyond

**Towards More Equitable Question Answering Systems:
How Much More Data Do you Need?**

Arnab Debnath, Navid Rajabi, Fardina Fathmiul Alam, Antonios Anastasopoulos
Department of Computer Science, George Mason University
{adebnath, nrajabi, falam5, antonis}@gmu.edu

(ACL 2021)

**How much data?**

# Going forward and beyond

**Towards More Equitable Question Answering Systems:**
**How Much More Data Do you Need?**

Arnab Debnath, Navid Rajabi, Fardina Fathmiul Alam, Antonios Anastasopoulos
Department of Computer Science, George Mason University
{adebnath,nrajabi,falam5,antonis}@gmu.edu

(ACL 2021)

**Dataset Geography: Mapping Language Data to Language Users**

Fahim Faisal, Yinkai Wang, Antonios Anastasopoulos
Department of Computer Science, George Mason University, USA
{ffaisal, ywang88, antonis}@gmu.edu

(ACL 2022)

**How much data?**

**What data?**

29

# Example: (Extractive) Question Answering

We have a lot of English-only datasets for QA (e.g. SQuAD)

Can we leverage them, and investigate few-shot approaches in new languages?

Study on TyDi-QA dataset (7 languages)

**Towards More Equitable Question Answering Systems:
How Much More Data Do you Need?**

Arnab Debnath, Navid Rajabi, Fardina Fathmiul Alam, Antonios Anastasopoulos
Department of Computer Science, George Mason University
{adebnath, nrajabi, falam5, antonis}@gmu.edu

(ACL 2021)

# Few-Shot adaptation

# Few-Shot adaptation



MACRO-AVG F1-SCORE

90

67.5

45

22.5

0

ZERO-SHOT (SQUAD)     XL +10/LANG     XL +50/LANG     XL +100/LANG     XL +500/LANG     XL +500/LANG +DATAAUG     SKYLINE

NLP
GEORGE
MASON

# Few-Shot adaptation

# Few-Shot adaptation



MACRO-AVG F1-SCORE

90 — 67.5 — 45 — 22.5 — 0

ZERO-SHOT (SQUAD): 58
XL +10/LANG: 64.2
XL +50/LANG
XL +100/LANG
XL +500/LANG
XL +500/LANG +DATAAUG
SKYLINE

31

# Few-Shot adaptation

# Few-Shot adaptation



31

# Few-Shot adaptation



Within 98% of skyline with less than 10% in-language training data!

# So how should we spend our annotation budget?

# So how should we spend our annotation budget?

**My view:**

Focus on building high-quality evaluation sets

Spend only a fraction of your budget on training data — combine with stronger baselines

# So how should we spend our annotation budget?

**My view:**

Focus on building high-quality evaluation sets

Spend only a fraction of your budget on training data — combine with stronger baselines

*\*\*Terms and Conditions apply*

# So how should we spend our annotation budget?

**My view:**

Focus on building high-quality evaluation sets

Spend only a fraction of your budget on training data — combine with stronger baselines

*\*\*Terms and Conditions apply*

4500 training examples in 1 language

# So how should we spend our annotation budget?

**My view:**

Focus on building high-quality evaluation sets

Spend only a fraction of your budget on training data — combine with stronger baselines

Avg F-score on
6 other languages:

*\*\*Terms and Conditions apply*

| | |
|---|---|
| 4500 training examples in 1 language | 72.3 |

NLP GEORGE MASON

# So how should we spend our annotation budget?

**My view:**

Focus on building high-quality evaluation sets

Spend only a fraction of your budget on training data — combine with stronger baselines

Avg F-score on
6 other languages:

*\*\*Terms and Conditions apply*

4500 training examples in 1 language          72.3
<
1500 training examples in 3 languages

NLP
GEORGE
MASON

# So how should we spend our annotation budget?

**My view:**

Focus on building high-quality evaluation sets

Spend only a fraction of your budget on training data — combine with stronger baselines

Avg F-score on
6 other languages:

*\*\*Terms and Conditions apply*

| | |
|---|---|
| 4500 training examples in 1 language | 72.3 |
| < | |
| 1500 training examples in 3 languages | 74.5 |

NLP
GEORGE
MASON

# So how should we spend our annotation budget?

**My view:**

Focus on building high-quality evaluation sets

Spend only a fraction of your budget on training data — combine with stronger baselines

Avg F-score on
6 other languages:

*\*\*Terms and Conditions apply*

| | |
|---|---|
| 4500 training examples in 1 language | 72.3 |
| < | |
| 1500 training examples in 3 languages | 74.5 |
| < | |
| 500 training examples in 6 languages | |

N L P
GEORGE
MASON

# So how should we spend our annotation budget?

**My view:**

Focus on building high-quality evaluation sets

Spend only a fraction of your budget on training data — combine with stronger baselines

Avg F-score on
6 other languages:

*\*\*Terms and Conditions apply*

4500 training examples in 1 language    72.3
<
1500 training examples in 3 languages    74.5
<
500 training examples in 6 languages    78.7

# So how should we spend our annotation budget?

**My view:**

Focus on building high-quality evaluation sets

Spend only a fraction of your budget on training data — combine with stronger baselines

Avg F-score on
6 other languages:

*\*\*Terms and Conditions apply*

| | |
|---|---|
| 4500 training examples in 1 language | 72.3 |
| < | |
| 1500 training examples in 3 languages | 74.5 |
| < | |
| 500 training examples in 6 languages | 78.7 |
| < | |
| 250 training examples in 12 languages | … |

# How Representative Are your Data?

# How Representative Are your Data?

Where does your data come from?

NLP GEORGE MASON

# How Representative Are your Data?

Where does your data come from?

Which speakers are modeled?

# How Representative Are your Data?

Where does your data come from?

Which speakers are modeled?

Study for country-level representation

**Dataset Geography: Mapping Language Data to Language Users**

Fahim Faisal, Yinkai Wang, Antonios Anastasopoulos
Department of Computer Science, George Mason University, USA
{ffaisal, ywang88, antonis}@gmu.edu

(ACL 2022)

33

# Idea

# Idea

**Named Entities can reveal the information we need!**

# Idea

**Named Entities can reveal the information we need!**

# Idea

Named Entities can reveal the information we need!

# Idea

**Named Entities can reveal the information we need!**

For a given dataset

RUI COSTA  FROM AMADORA PLAYED FOR FIORENTINA

# Idea

**Named Entities can reveal the information we need!**

For a given dataset

☑ Identify named entities

RUI COSTA FROM AMADORA PLAYED FOR FIORENTINA

# Idea

**Named Entities can reveal the information we need!**

For a given dataset

☑ Identify named entities

☑ Link entities to countries
through wikidata

RUI COSTA FROM AMADORA PLAYED FOR FIORENTINA

WIKIDATA ID: Q11571
COUNTRY: PORTUGAL

WIKIDATA ID: Q1422
COUNTRY: ITALY

# Idea

## Named Entities can reveal the information we need!

For a given dataset

- ☑ Identify named entities
- ☑ Link entities to countries through wikidata
- ☑ Aggregate through dataset
  - ☑ Representativeness measures
  - ☑ Fairness measures
  - ☑ Visualizations

RUI COSTA FROM AMADORA PLAYED FOR FIORENTINA

WIKIDATA ID: Q11571
COUNTRY: PORTUGAL

WIKIDATA ID: Q1422
COUNTRY: ITALY

country
Portugal
Italy

# Dataset Geography

Code
&
Dataset

https://github.com/ffaisal93/dataset_geography

Project Webpage
&
Additional Visualizations

https://nlp.cs.gmu.edu/project/datasetmaps

# Dataset Geography

# Dataset Geography



Dataset Map: Masakhaner kinyarwanda

Dataset Entities Map

Top-10 Represented Countries
Countries Missing: 171 of 243 (70.37%)

Main Countries where language is spoken.
Percentage in-country: 49.74%
(Green indicates allocation proportional to the population)

# Dataset Geography



Dataset Map: Masakhaner wolof

Dataset Entities Map

Top-10 Represented Countries
Countries Missing: 177 of 243 (72.84%)

Main Countries where language is spoken.
Percentage in-country: 24.20%
(Green indicates allocation proportional to the population)

# Dataset Geography

# What do communities need/want?

# What do communities need/want?

Work *with* the communities *for* the communities

# What do communities need/want?

Work *with* the communities *for* the communities

NLP
GEORGE
MASON

# What do communities need/want?

Work *with* the communities *for* the communities

# What do communities need/want?
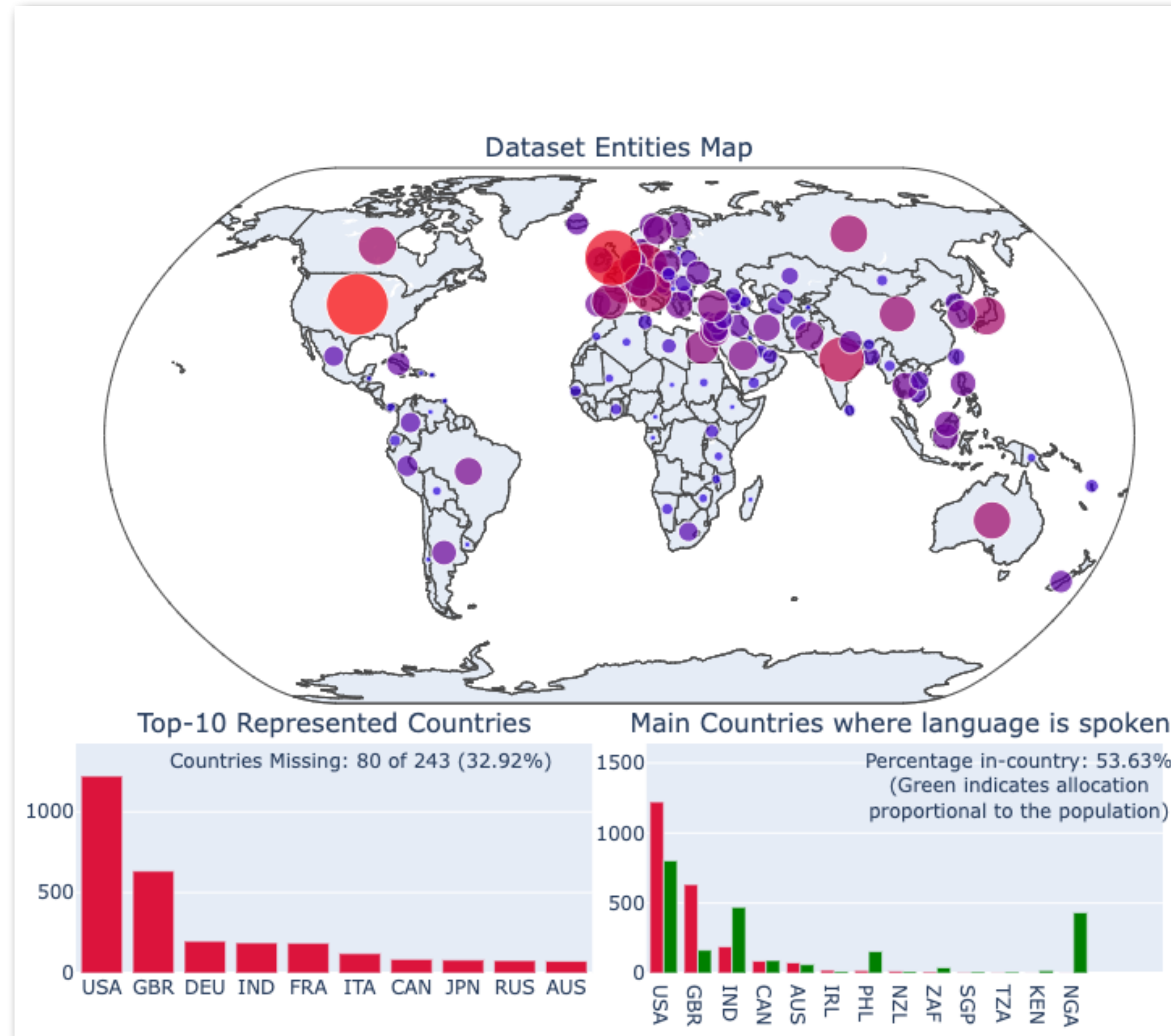
Work *with* the communities *for* the communities



**Educational Tools for Mapuzugun**

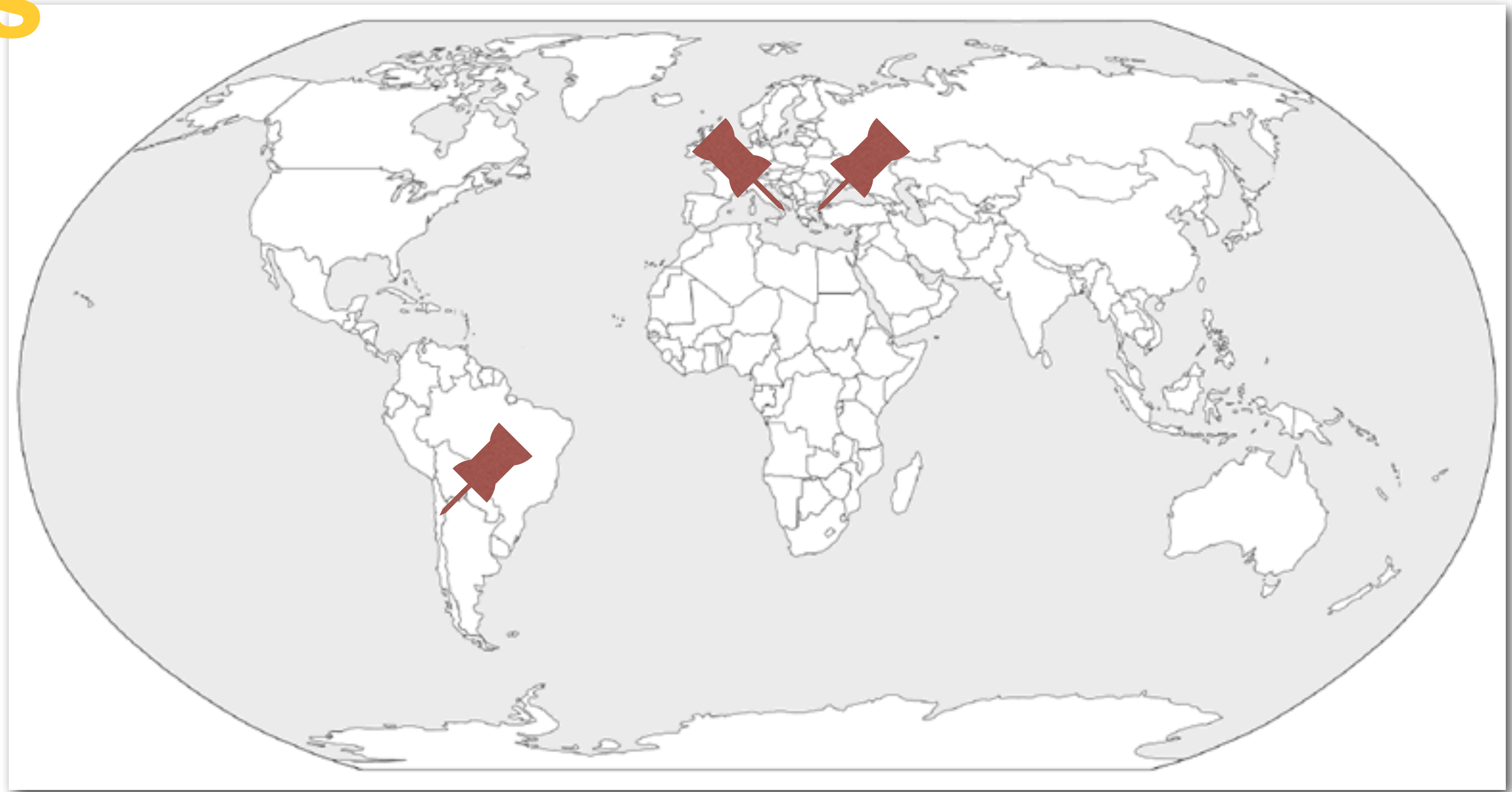Cristian Ahumada[1]    Claudio Gutierrez[1]    Antonios Anastasopoulos[2]
[1]Department of Computer Science, Universidad de Chile
[2]Computer Science Department, George Mason University
ahumada.860@gmail.com    cgutierr@dcc.uchile.cl    antonis@gmu.edu

# What do communities need/want?

Work *with* the communities *for* the communities

# What do communities need/want?

Work *with* the communities *for* the communities



BembaSpeech: A Speech Recognition Corpus for the Bemba Language

Claytone Sikasote*
Department of Computer Science
University of Zambia
Zambia
claytone.sikasote@cs.unza.zm

Antonios Anastasopoulos
Department of Computer Science
George Mason University
USA
antonis@gmu.edu

BIG-C: Multimodal Dataset for the Bemba Language

Claytone Sikasote[1], Eunice Mukonde[1], and Antonios Anastasopoulos[2]
[1]Department of Computer Science, University of Zambia, Zambia
[2]Department of Computer Science, George Mason University, USA
claytone.sikasote@cs.unza.zm, antonis@gmu.edu

36

# What do communities need/want?

Work *with* the communities *for* the communities



BembaSpeech: A Speech Recognition Corpus for the Bemba Language

Claytone Sikasote*
Department of Computer Science
University of Zambia
Zambia
claytone.sikasote@cs.unza.zm

Antonios Anastasopoulos
Department of Computer Science
George Mason University
USA
antonis@gmu.edu

BIG-C: Multimodal Dataset for the Bemba Language

Claytone Sikasote[1], Eunice Mukonde[1], and Antonios Anastasopoulos[2]
[1]Department of Computer Science, University of Zambia, Zambia
[2]Department of Computer Science, George Mason University, USA
claytone.sikasote@cs.unza.zm, antonis@gmu.edu

NLP
GEORGE
MASON

36

# Thank you!

Shoutout to collaborators:
Graham Neubig, Damian Blasi,
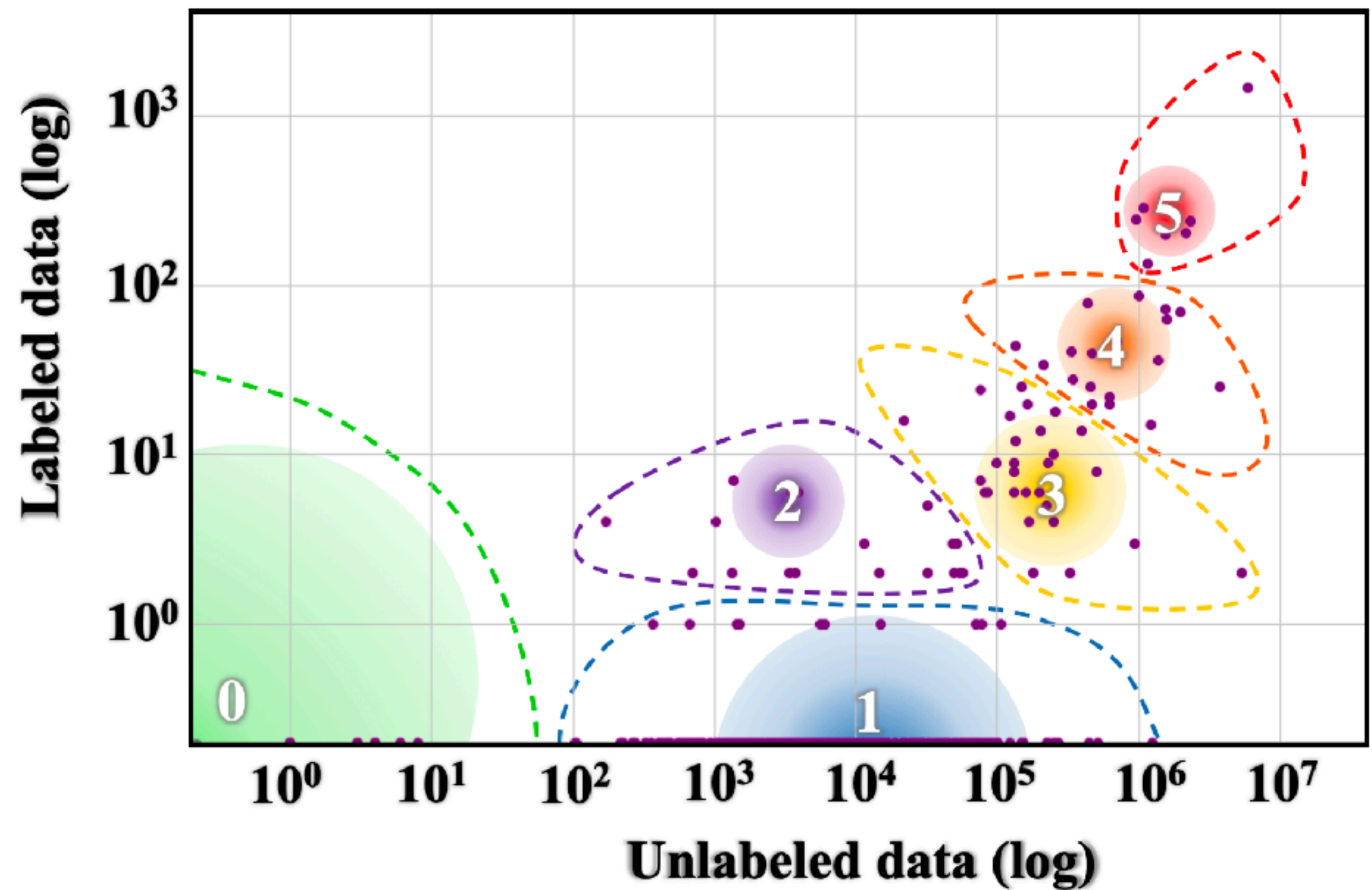Benjamin Muller, Benoît Sagot,
Djamé Seddah

And students:
Fahim Faisal, Sharlina Keshava,
Mahfuz ibn Alam, Yinkai Wang

Other things I'm working on:

- NLP for endangered languages (e.g. OCR for scanned documents from Latin America, building basic tools for Griko, Mapudungun, Pomak)

- NLP for linguists (Machine-aided annotation)

- Machine Translation from/into dialects

- Cross-Lingual and Cross-Cultural Fairness

- Geospatial Language Understanding and Navigation

- SLT for Crisis Response

- …

## GMU and GMNLP is hiring!
## Faculty/postdocs/PhD students

N L P
GEORGE
MASON

The amount of data, labeled or unlabeled, varies wildly across languages!

Image from Joshi et al 2020