# On the Inter-Linguistic Disparity of Knowledge Graphs: Bridging the Gap between English and non-English Languages

Simone Conia · *UniDive – 2nd General Meeting | WG-{1,3,4}* · Apple & Sapienza University of Rome

UniDive · COST · Funded by the European Union · FAIR Future Artificial Intelligence Research · SAPIENZA Università di Roma · Apple
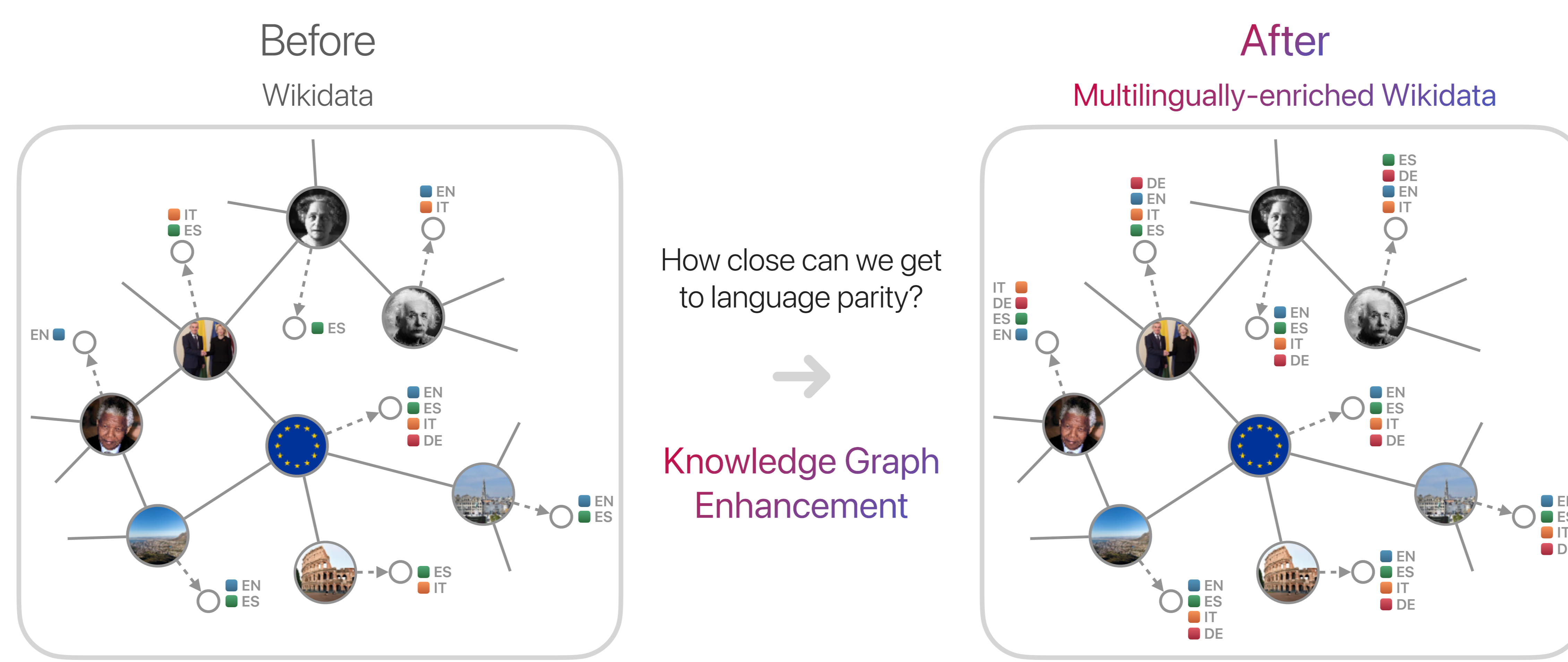
## Abstract

Recent work in Natural Language Processing and Computer Vision has been using textual information – e.g., entity names and descriptions – available in **knowledge graphs to ground neural models to high-quality structured data**. However, when it comes to **non-English languages**, the quantity and quality of textual information are comparatively scarce.

To address this issue, we introduce the novel task of automatic **Knowledge Graph Completion (KGE)** and perform a thorough investigation on bridging the gap in both the quantity and quality of textual information between English and non-English languages. More specifically, we: i) bring to light the problem of increasing multilingual coverage and precision of entity names and descriptions in Wikidata; ii) demonstrate that state-of-the-art methods, namely, Machine Translation (MT), Web Search (WS), and Large Language Models (LLMs), struggle with this task; iii) present **M-NTA, a novel unsupervised approach that combines MT, WS, and LLMs to generate high-quality textual information**; and, iv) study the impact of increasing multilingual coverage and precision of non-English textual information in Entity Linking, Knowledge Graph Completion, and Question Answering.

As part of our effort towards better multilingual knowledge graphs, we also introduce **WikiKGE-10, the first human-curated benchmark to evaluate KGE approaches in 10 languages** across 7 language families.

## Multilingual Knowledge Graphs | Overview

Knowledge graphs (KGs) encode our collective understanding of the world in a structured representation.



**Before** — Wikidata

How close can we get to language parity?

**After** — Multilingually-enriched Wikidata

**Knowledge Graph Enhancement**

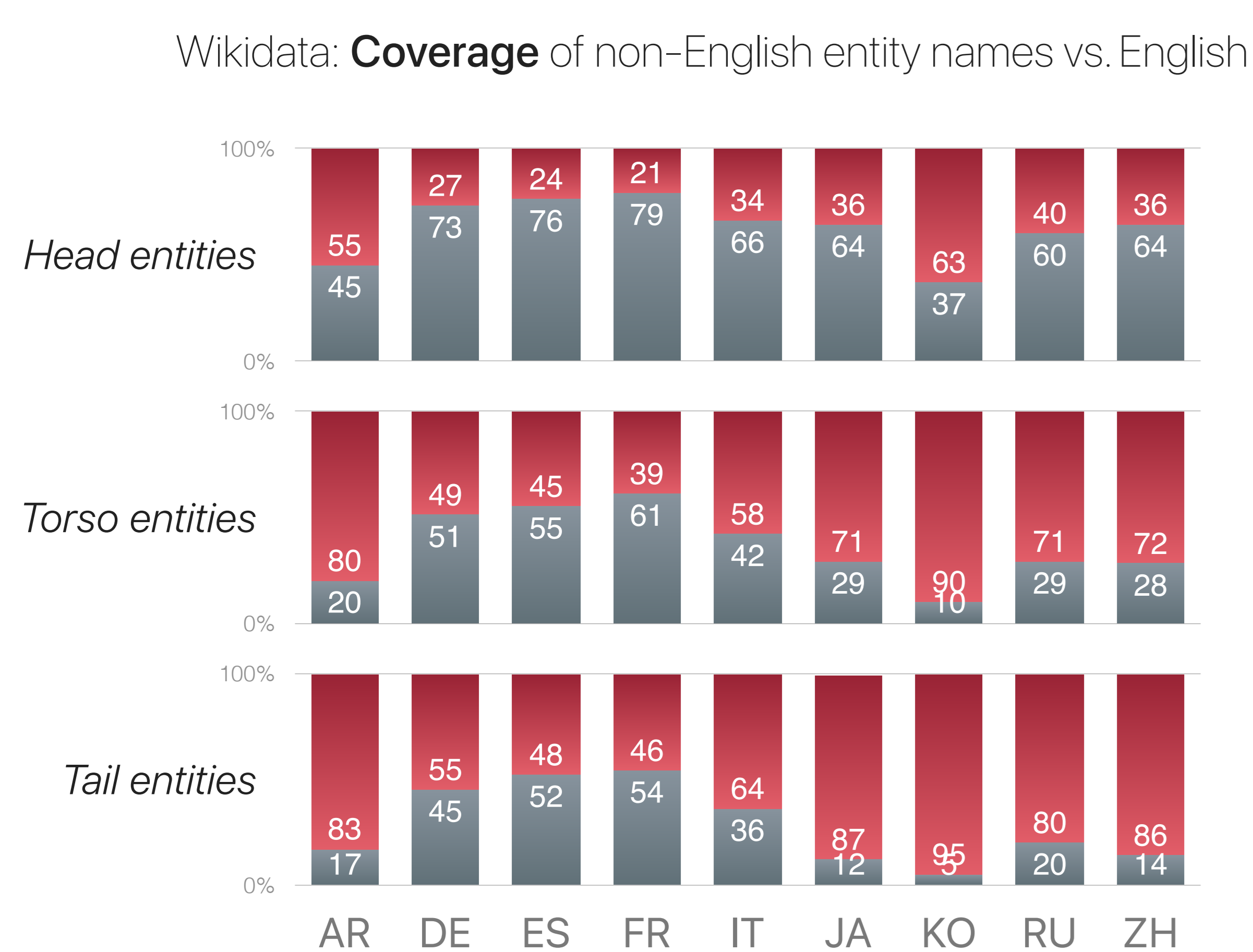Textual information in **multilingual** KGs lacks

# Coverage | Precision

↓

Limited **multilingual** applications

- Language Modeling
- Recommenders
- Question Answering
- Entity Linking
- Image Classification
- Information Retrieval
- Sense Disambiguation
- Other tasks

## Textual Information in Multilingual Knowledge Graphs | Current Limitations

Even for entity names, multilingual coverage is far from 100%.

Wikidata: **Coverage** of non-English entity names vs. English



Head entities / Torso entities / Tail entities — AR DE ES FR IT JA KO RU ZH

Many models use KGs out of the box but textual information in KGs is not always precise.

### KGs contains human errors

| Language | Entity name |
|---|---|
| English | Olivier Giroud |
| French | Olivier Giroud |
| Spanish | **Oliver** Giroud |
| Japanese | オリヴィエ・ジルー |
| Chinese | 奥利维耶·吉鲁 |
| ... | |

Image of Olivier Giroud: © Anton Stalius 2018

A spelling error in the primary name of a popular entity in Wikidata

### KGs contains stale entries

| Language | Entity name |
|---|---|
| English | Elliot Page |
| Trad. Chinese | 艾略特·佩吉 (Elliot Page) |
| Simp. Chinese | 艾莲·佩奇 (Ellen Page) |
| ... | |

Image of Elliot Page: © Fosse 2023

This entity name has not been updated to reflect changes in the real world.

### KGs contains under-specific information

| Lang. | Entity description |
|---|---|
| EN | Japanese composer (1952–2023) |
| ES | músico japonés |
| FR | musicien, compositeur, producteur et acteur japonais |

Image of Ryuichi Sakamoto: © Joi Ito 2007

For this entity, different languages have descriptions with different information

## WikiKGE-10 | A underline{human-graded} benchmark for evaluating KGE systems in underline{10 languages}

How did we create WikiKGE-10?

**1.** We select 10 languages across 7 linguistic families

*Arabic  German  English  Spanish  French  Italian  Japanese  Korean  Russian  Chinese*

**2.** We sample 1000 entities from the top-10%


Popularity / Head entities / Entity count

**3.** Humans validate, identify errors, add names

**4.** Agreement?

**5.** Done! WikiKGE-10

WikiKGE-10 contains over 35k manually-graded entity names across 10 languages.

| | AR | DE | EN | ES | FR | IT | JA | KO | RU | ZH | **All** |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Entities | 1,000 | 1,000 | 1,000 | 1,000 | 1,000 | 1,000 | 1,000 | 1,000 | 1,000 | 1,000 | **10,000** |
| Entity names in WikiKGE-10 | 4,213 | 3,498 | 2,837 | 4,320 | 3,548 | 3,156 | 2,999 | 3,874 | 3,901 | 4,088 | **36,434** |
| - **Valid names** in Wikidata | 2,521 | 2,336 | 2,090 | 2,732 | 2,330 | 1,840 | 2,235 | 2,136 | 2,706 | 2,569 | **23,495** |
| - **Invalid names** in Wikidata | 320 | 491 | 219 | 571 | 530 | 236 | 486 | 329 | 507 | 830 | **4,663** |

**Strong agreement** — Krippendorf's alpha = **0.94**

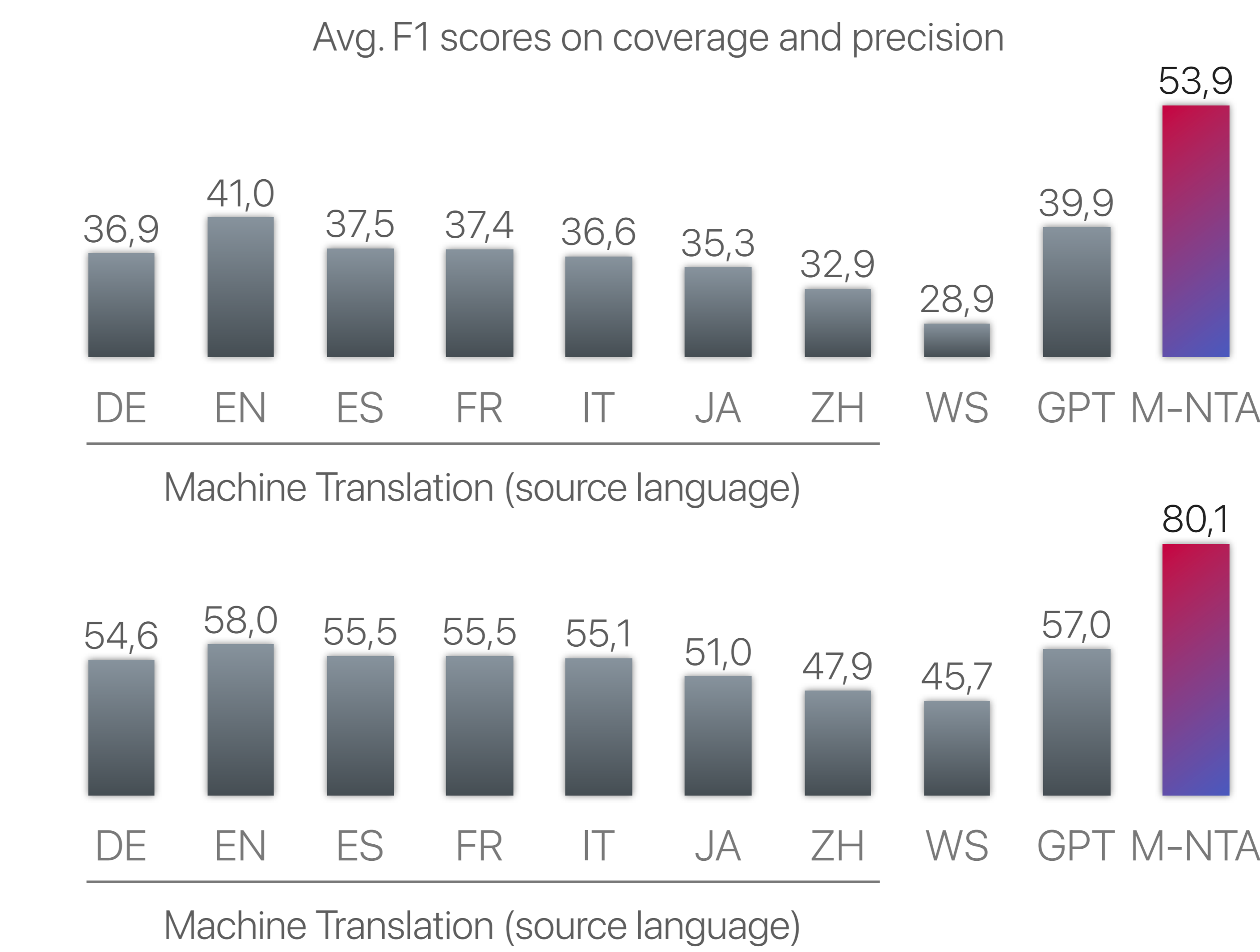**Wikidata is incomplete** — **+35-40%** names in WikiKGE-10

**Wikidata is imprecise** — **20%** names are incorrect

## M-NTA | Combining MT, WS, and LLMs for KGE

**M-NTA** leverages the complementary knowledge across locales and tools to provide accurate predictions.

**M-NTA** — For entity names & entity descriptions



KG / Ensembling / Entity selection / Alignment (text-to-triple) / Naturalization (triple-to-text) / Machine Translation / Web Search / LLMs

**M-NTA** can successfully combine information from multiple tools, sources, and languages.

Avg. F1 scores on coverage and precision



Machine Translation (source language) — DE 36,9 | EN 41,0 | ES 37,5 | FR 37,4 | IT 36,6 | JA 35,3 | ZH 32,9 | WS 28,9 | GPT 39,9 | M-NTA 53,9

Machine Translation (source language) — DE 54,6 | EN 58,0 | ES 55,5 | FR 55,5 | IT 55,1 | JA 51,0 | ZH 47,9 | WS 45,7 | GPT 57,0 | M-NTA 80,1

**Main takeaway**

Generating a fact from **multiple sources and languages may offer complementary pieces of information**, which provide varying views on our world knowledge.

## KGE | Impact on downstream tasks

### Multilingual **Entity Linking**

"El comandante **Armstrong** fue el primer ser humano que pisó la superficie del satélite terrestre el 21 de julio de 1969 a las 2:56 (hora internacional UTC) [...]."

Neil Armstrong — *American astronaut*

Edwin Howard Armstrong — *American electrical engineer*

mGENRE on Wikinews-7



| FR | IT | RU | ZH | Avg |
FR 73,4 / 74,1 · IT 56,8 / 58,2 · RU 65,8 / 66,2 · ZH 52,8 / 55,0 · Avg 62,2 / 63,3

w/out M-NTA · w/ M-NTA

### Multilingual **KG Completion**


EN IT FR ZH

AlignKGC on DBP-5L

| FR | IT | RU | ZH | Avg |
FR 47,4 / 47,5 · IT 64,6 / 66,3 · RU 64,4 / 66,0 · ZH 62,8 / 64,2 · Avg 59,8 / 61,1

w/out M-NTA · w/ M-NTA

### Multilingual **KGQA**

Q: *Quanto è alto Barack Obama?*
A: *1,87 metri*

Q: 영화 식스티 세컨즈의 배우는
A: 안젤리나 졸리, 윌리엄 리 스콧, ...

Unanswerable queries in MKQA



DE -52,0% / -15,4% · EN · ES -41,9% / 39,9% · FR · IT -53,8% · JA -32,1% · KO · ZH -22,4% · Avg -36,8%

w/ M-NTA

## Conclusion

WikiKGE-10 is available at:

https://github.com/apple/ml-kge



Simone Conia
simone.conia@uniroma1.it

**Special thanks** to

Min Li
min_li6@apple.com

Daniel Lee
daniel_lee1@ucalgary.ca

Umar Farooq Minhas
ufminhas@apple.com

Ihab Ilyas
ilyas@apple.com

Yunyao Li
yunyaol@adobe.com