

XOR QA: Cross-lingual Open-Retrieval Question Answering across many languages

Akari Asai¹, Jungo Kasai¹, Jonathan H. Clark², Kenton Lee², Eunsol Choi³, Hannaneh Hajishirzi¹⁴

¹University of Washington, ²Google Research, ³The University of Texas at Austin, ⁴Allen Institute for AI



Links ▶

Questions?
@AkariAsai



Website
& leaderboard



Repository

Summary

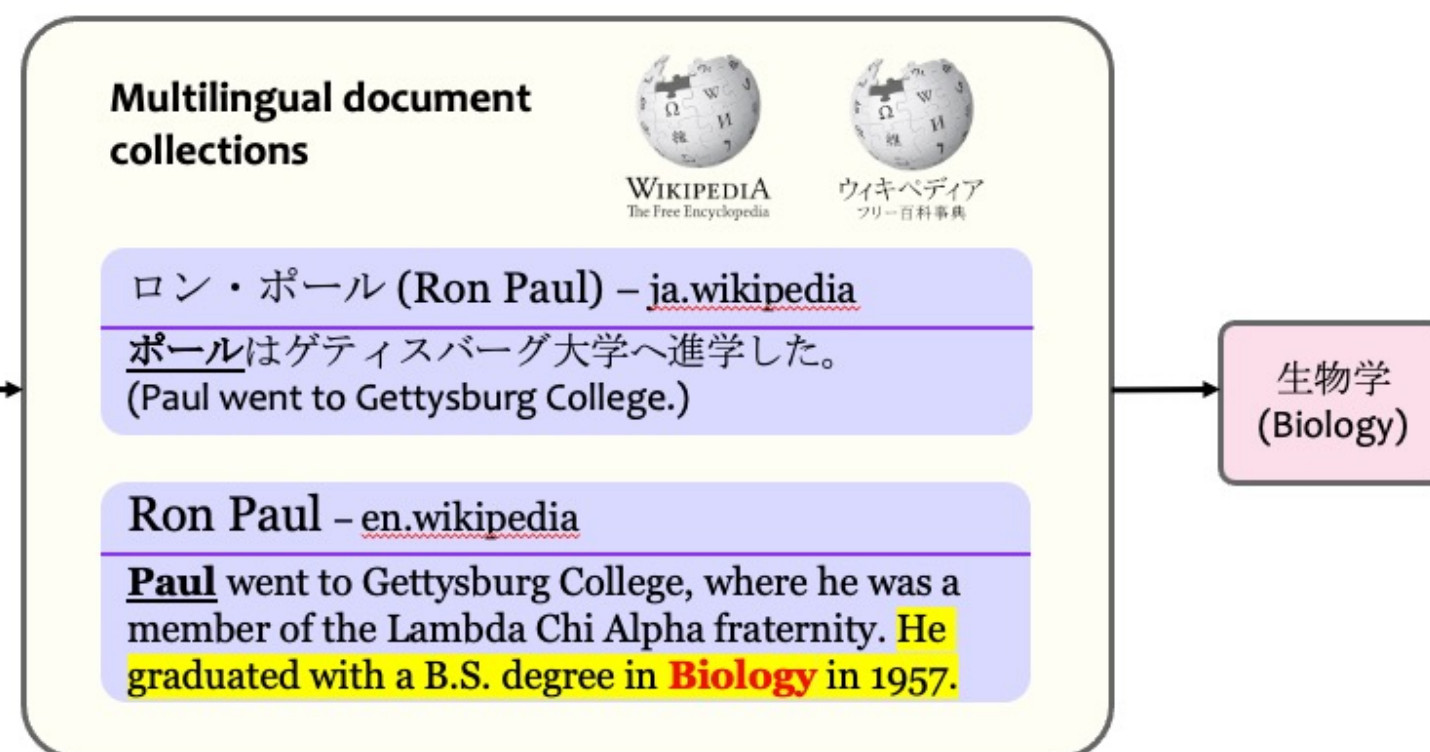
- A large-scale multilingual open-retrieval QA, **XOR-TyDi QA**
 - 40 k** newly annotated data across **7 languages**
 - Fine-grained **paragraph & span level annotations** & **human translated questions**
- A new task, **XOR QA** involving **cross-lingual IR and QA**.
- Experimental results show large spaces for improvements.

Background

Multilingual Open-retrieval QA is challenging:

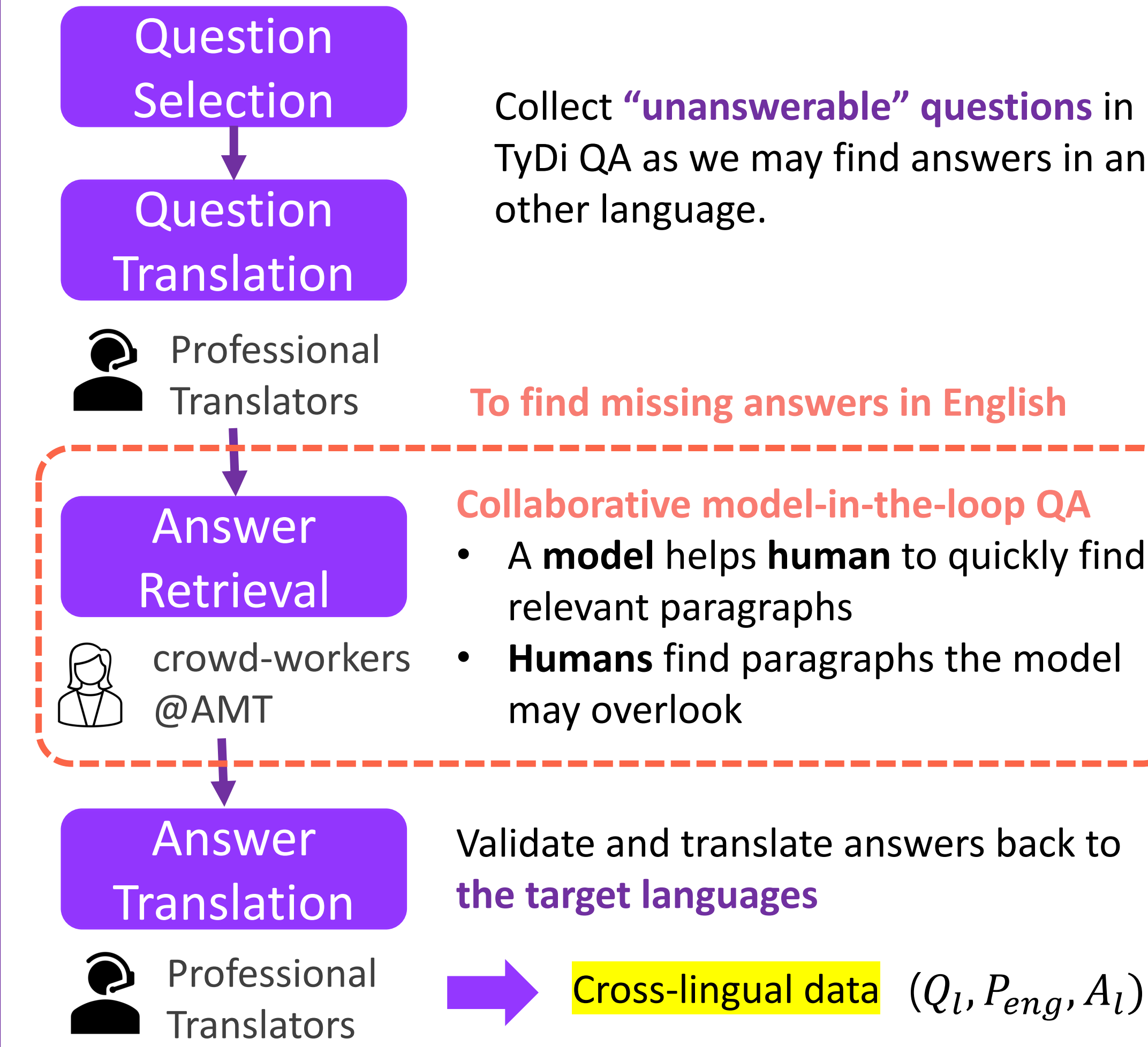
- Information scarcity:** limited web documents available in low resource languages (e.g., Telugu).
- Information asymmetry:** bias towards their own culture or general interests

→ To build a QA system that works well across many languages, **need to retrieve and read across multiple languages.**



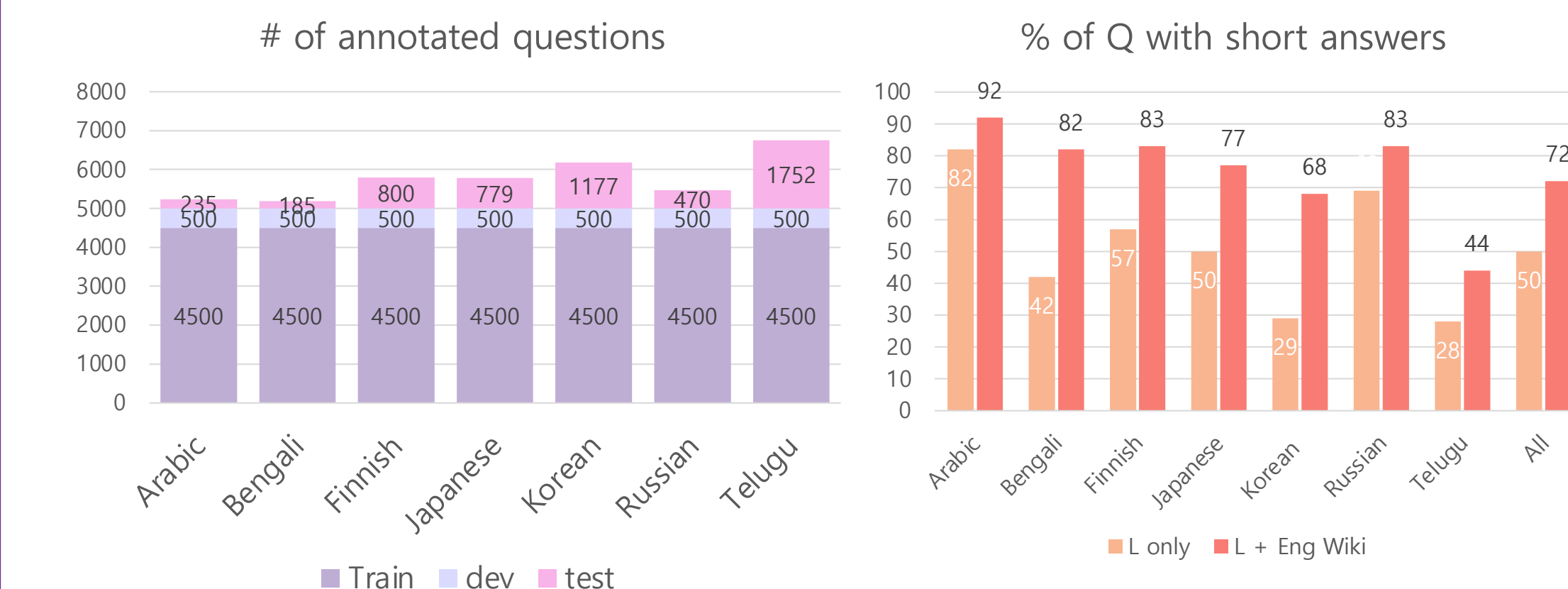
Dataset Collections of XOR-TyDi

- For scalability and quality, **decompose tasks into subtasks**
- Introduce a new annotation framework, **collaborative model-in-the-loop QA**



Dataset Statistics of XOR-TyDi

- Train & evaluation data for 7 languages**
- Increase the answer coverage up to **40%**



XOR QA: 3 increasingly complex and useful tasks

XOR-Retrieve

- Input:** Query in l and English Wikipedia
- Output:** top k documents
- Metric:** Top- k tokens recall ($k = 2000, 5000$)

XOR-EngSpan

- Input:** Query in l and English Wikipedia
- Output:** an answer in English, A_{eng}
- Metric:** Token F1, EM

XOR-Full

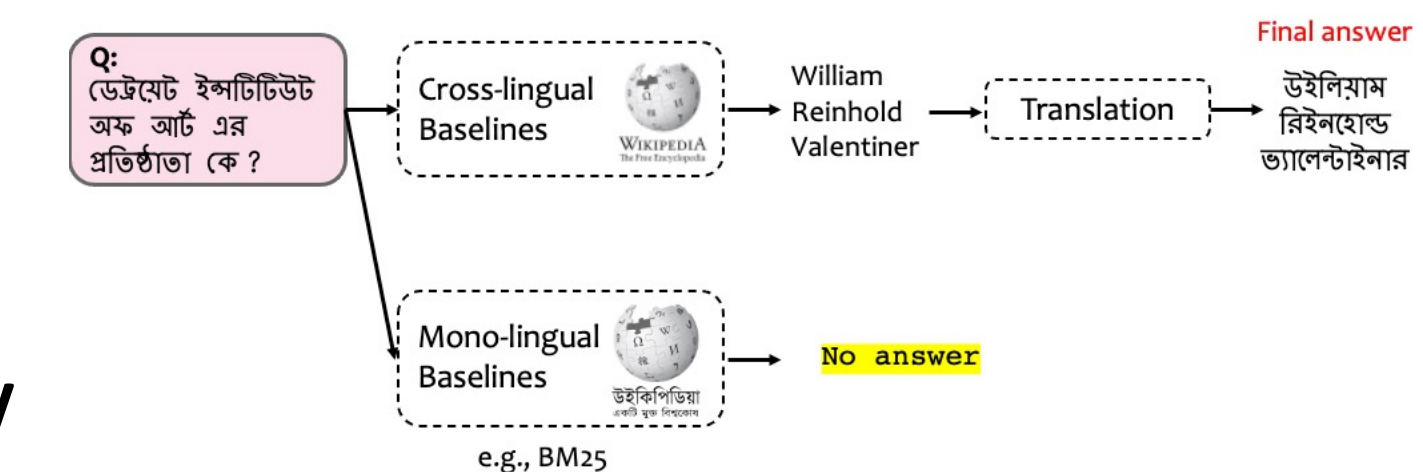
- Input:** Query in l , English Wikipedia and the Wikipedia in l
- Output:** an answer in l
- Metric:** Token F1, EM, BLEU

Cross-lingual baselines (for: XOR-Retrieve, XOR-English Span)

- Translate-test:** translate query into English, run English DPR
- Multilingual:** use multilingual BERT & XLM-R for encoders

Full baselines (for: XOR-full)

- Full:** Cross-lingual Baselines + Monolingual baseline
- English only**
- Monolingual-only**

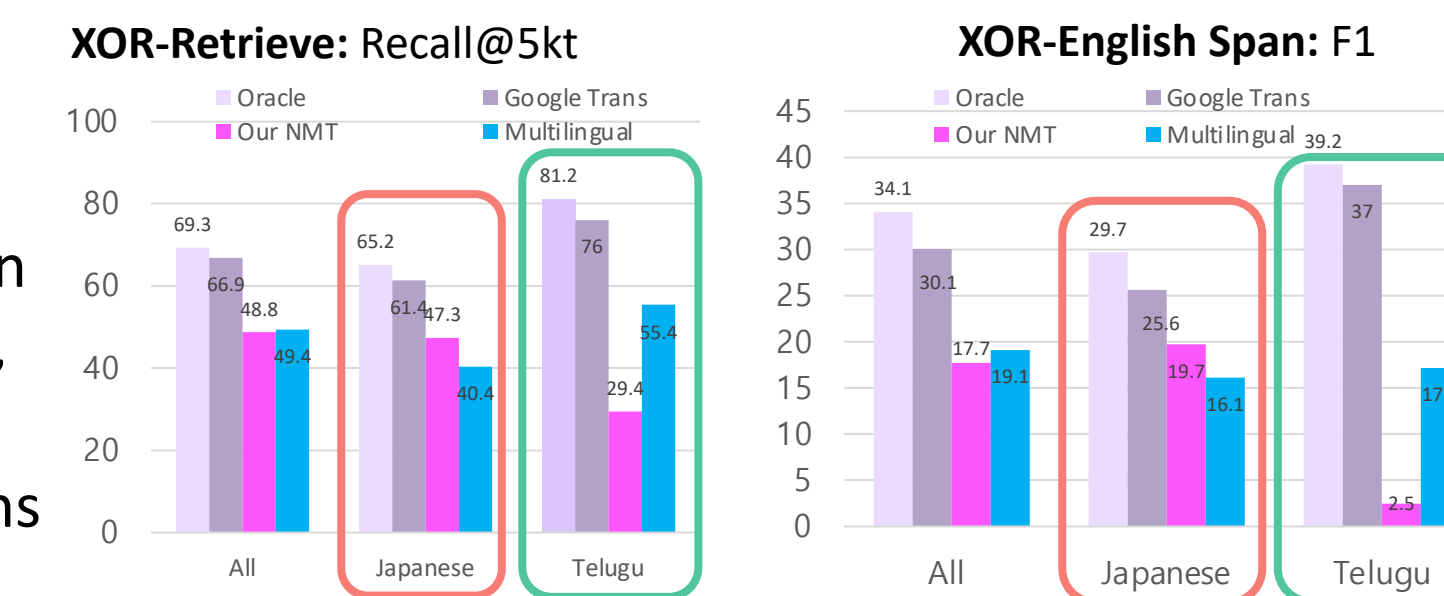


Results

XOR-Retrieve

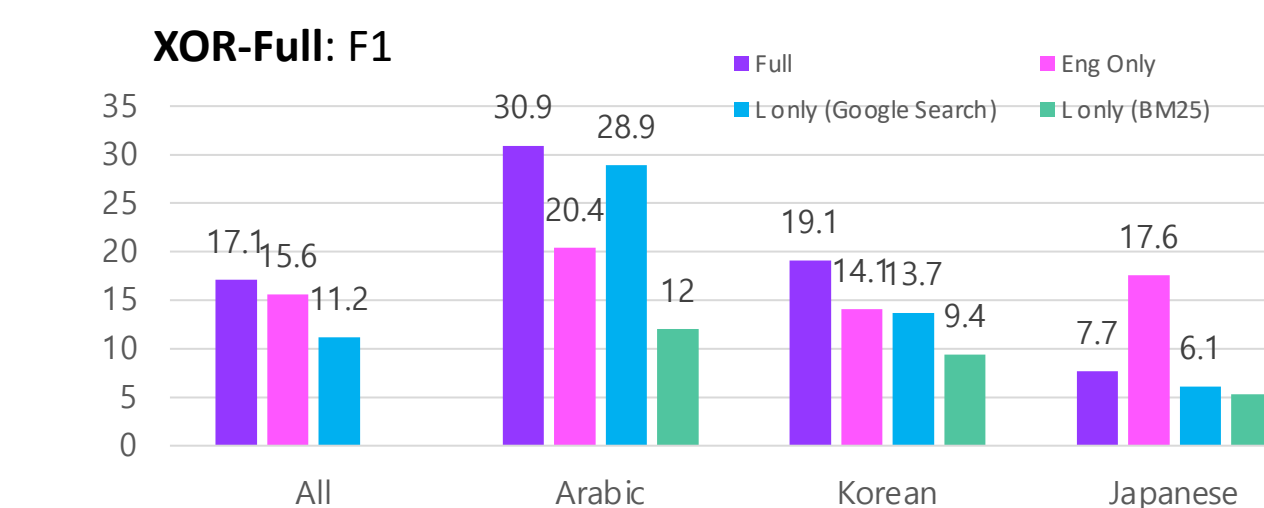
XOR-EngSpan

- High resources L:** Translate-test (w/ human, GMT) > Multilingual
- Low resources L:** Translate-test suffers



XOR-Full

- All models show limited performance**
- Full model performs best



Comparison with other datasets

Dataset	Asked by native speakers	Open-retrieval over millions of documents	Typological Diversity	Cross-lingual	train data
TyDi QA [1]	✓		✓		✓
MLQA [2]			✓	✓	
XQuAD [3]			✓		
MKQA [4]			✓		
MLQA-R [5]				✓	
QA@CLEF 2002-2008 [6]	✓				
XOR-TyDi QA	✓	✓	✓	✓	✓

References

- Clark et al., 2020. “TyDi QA: A Benchmark for Information-Seeking Question Answering in Typologically Diverse Languages”. *TACL*
- Lewis et al., 2020. “MLQA: Evaluating Cross-Lingual Extractive Question Answering”. In *ACL*.
- Artetxe et al., 2020. “On the cross-lingual transferability of monolingual representations”. In *ACL*.
- Longpre et al., 2020. “MKQA: A linguistically diverse benchmark for multilingual open domain question answering”. *arXiv preprint*.
- Roy et al., 2020. “LAREQA: Language-agnostic answer retrieval from a multilingual pool”. In *EMNLP*.
- Former et al., 2019. “Evaluating Multilingual Question Answering System at CLEF”. In *LREC*.