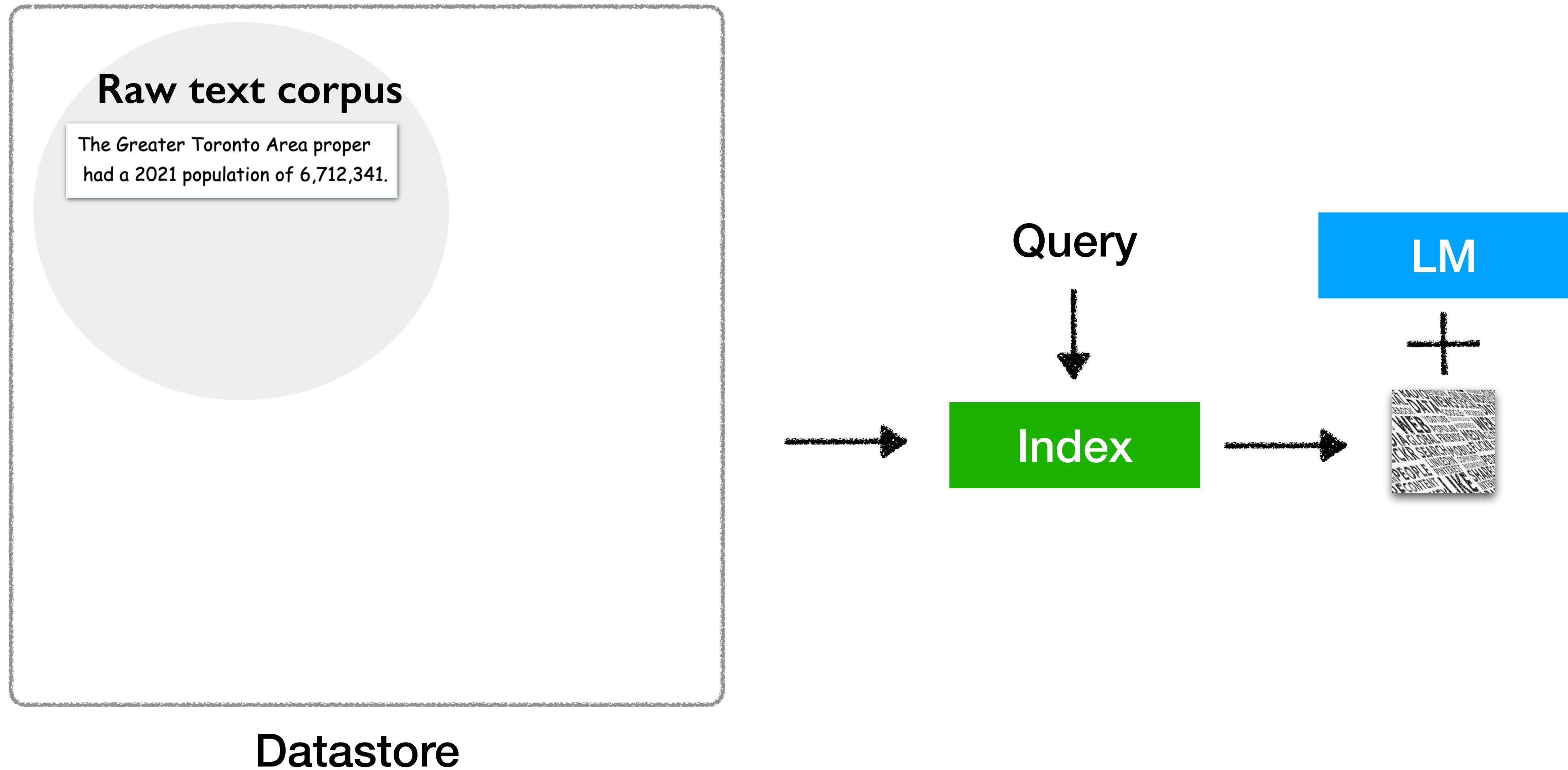
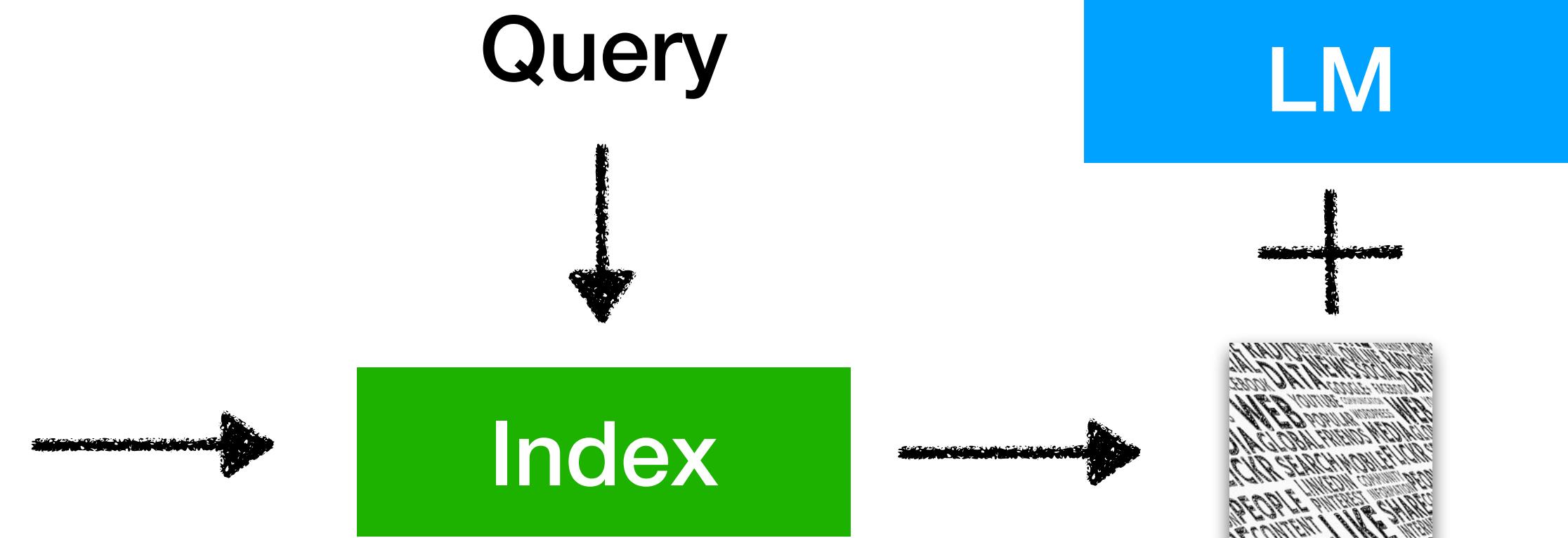
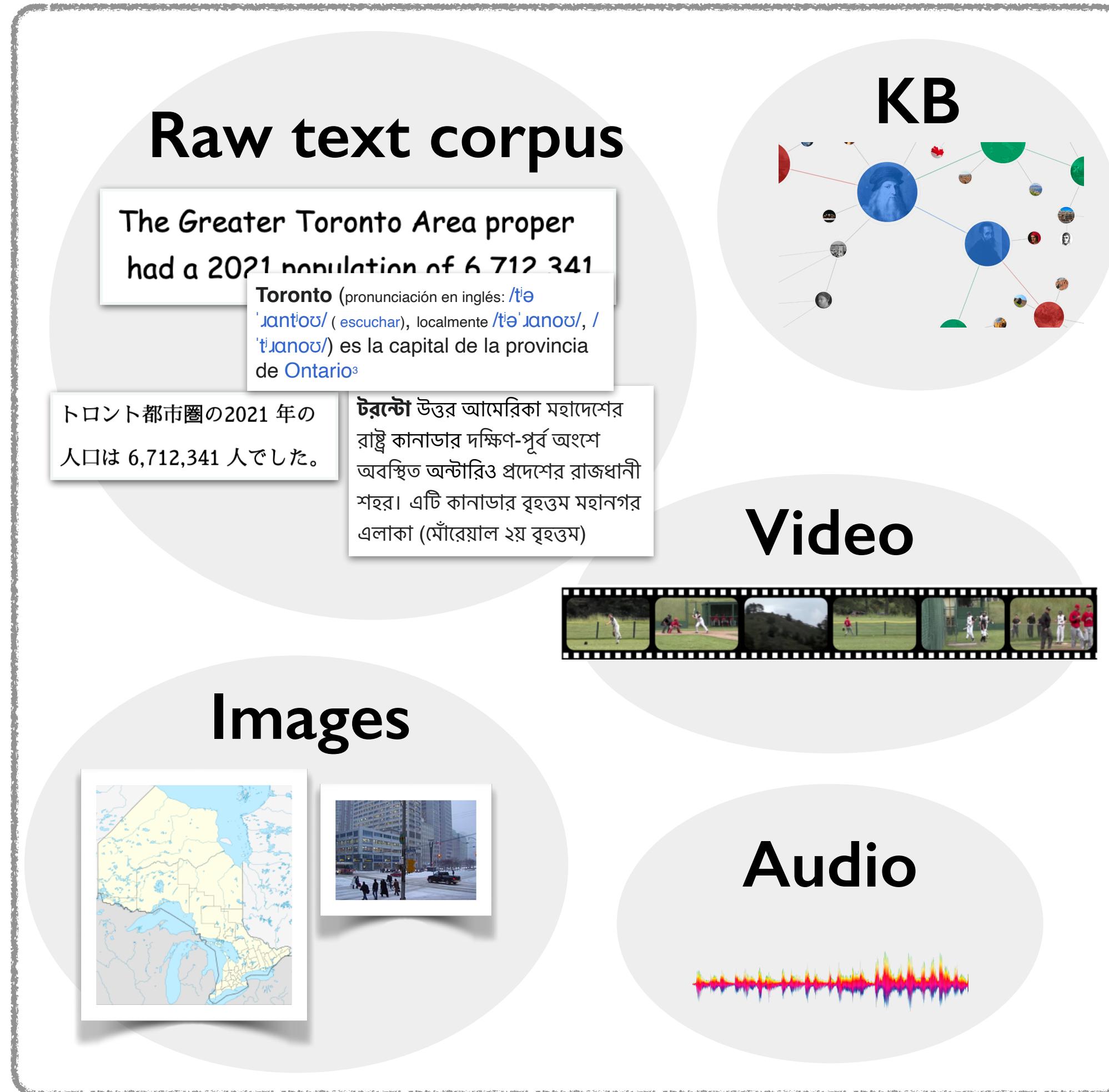


Section 7: Multilingual & Multimodal

Retrieval-based LM for diverse knowledge sources

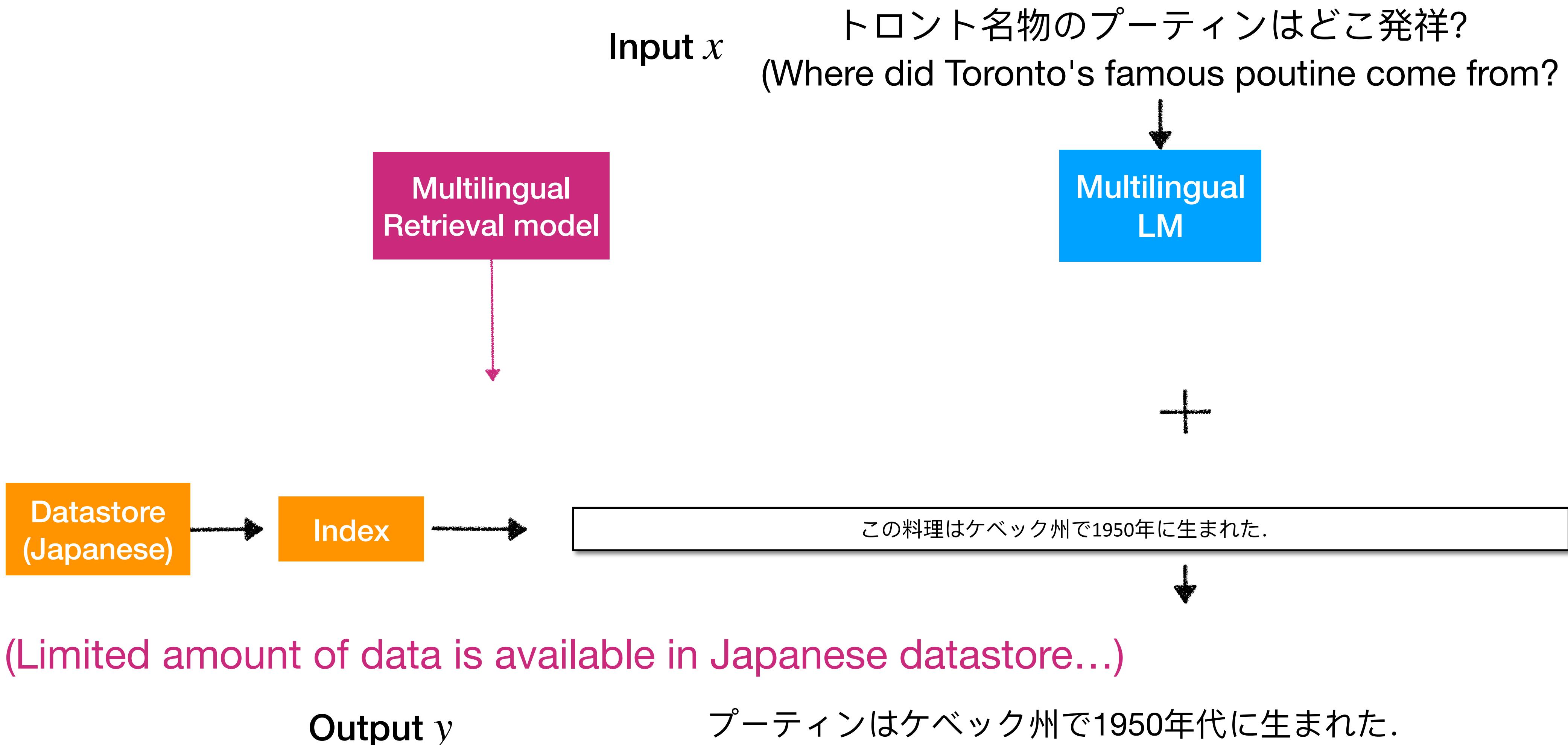


Retrieval-based LM for diverse knowledge sources

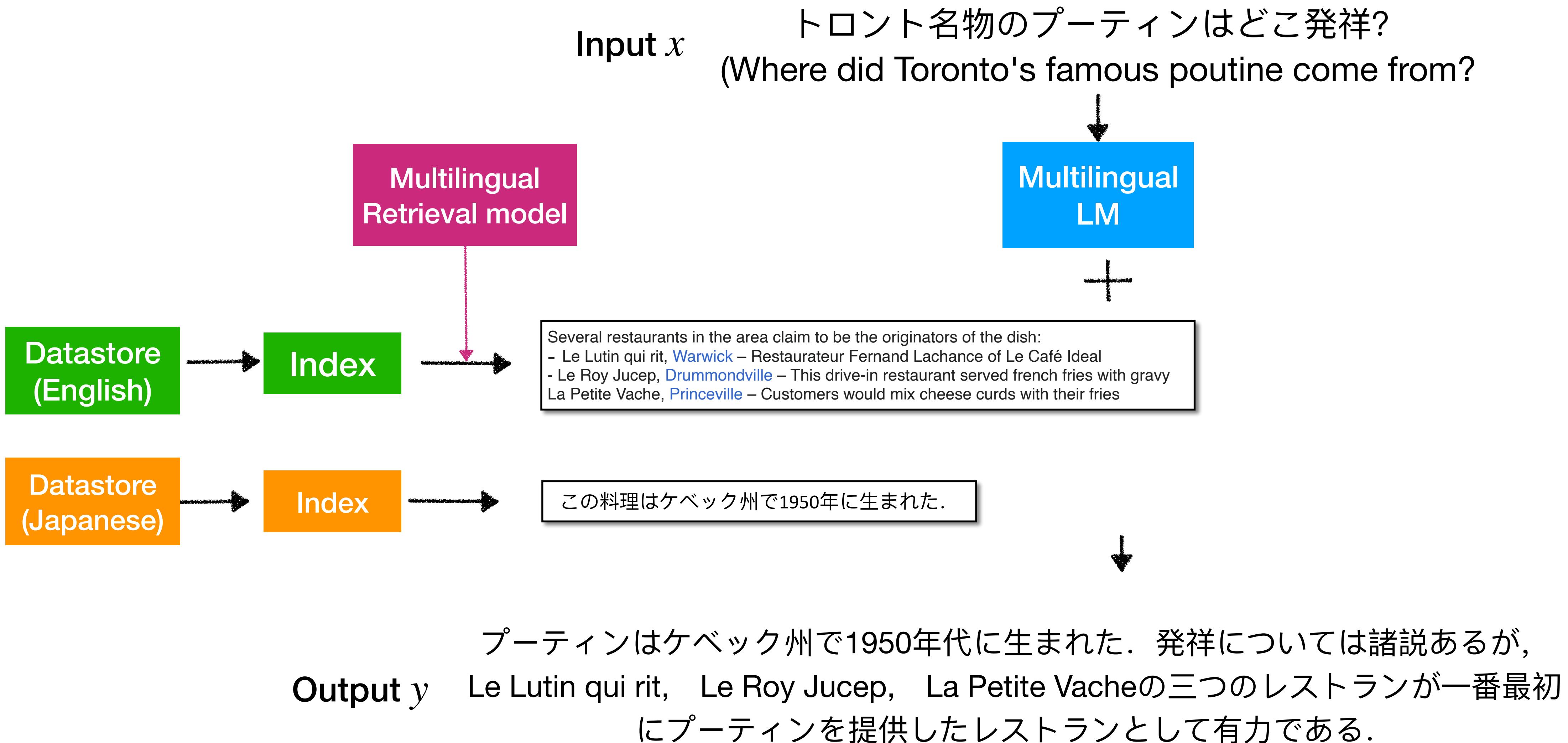


Datastore

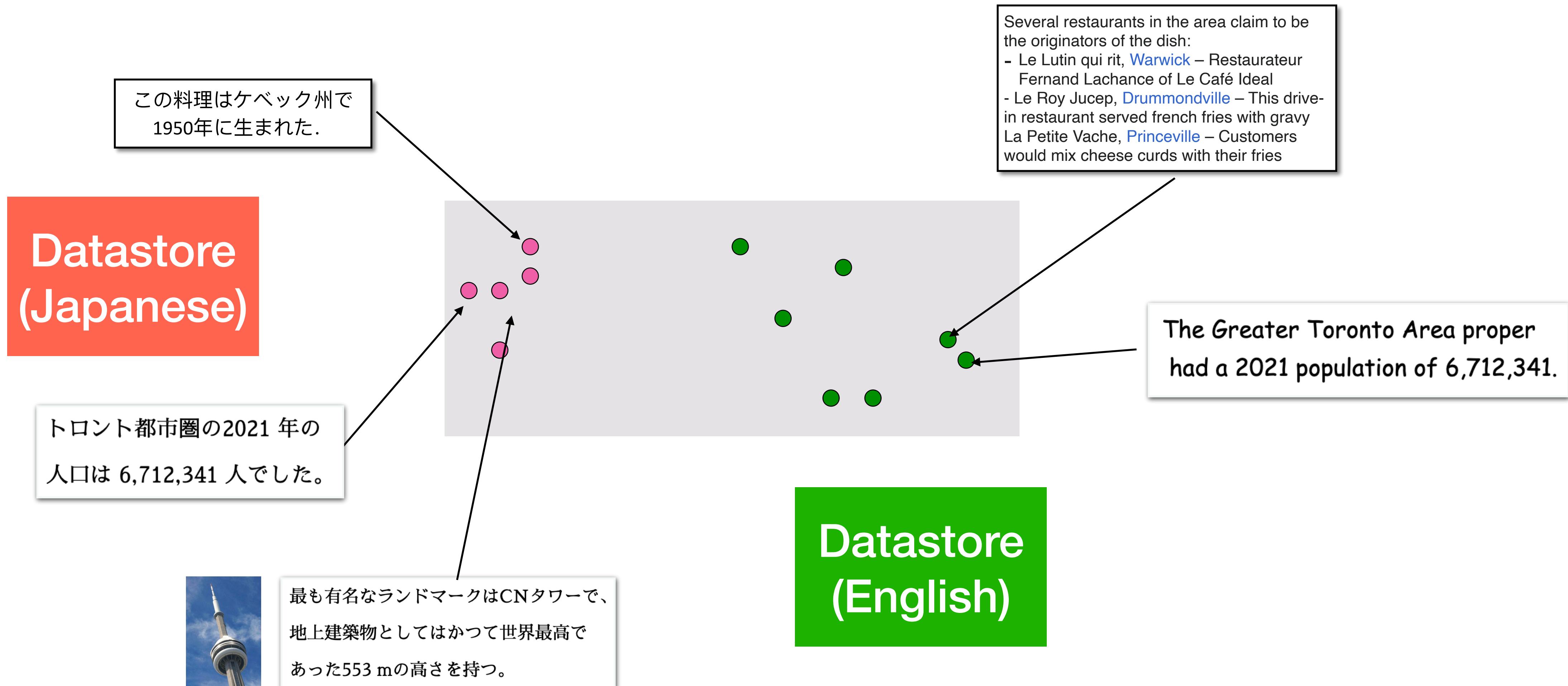
Multilingual Retrieval-based LM



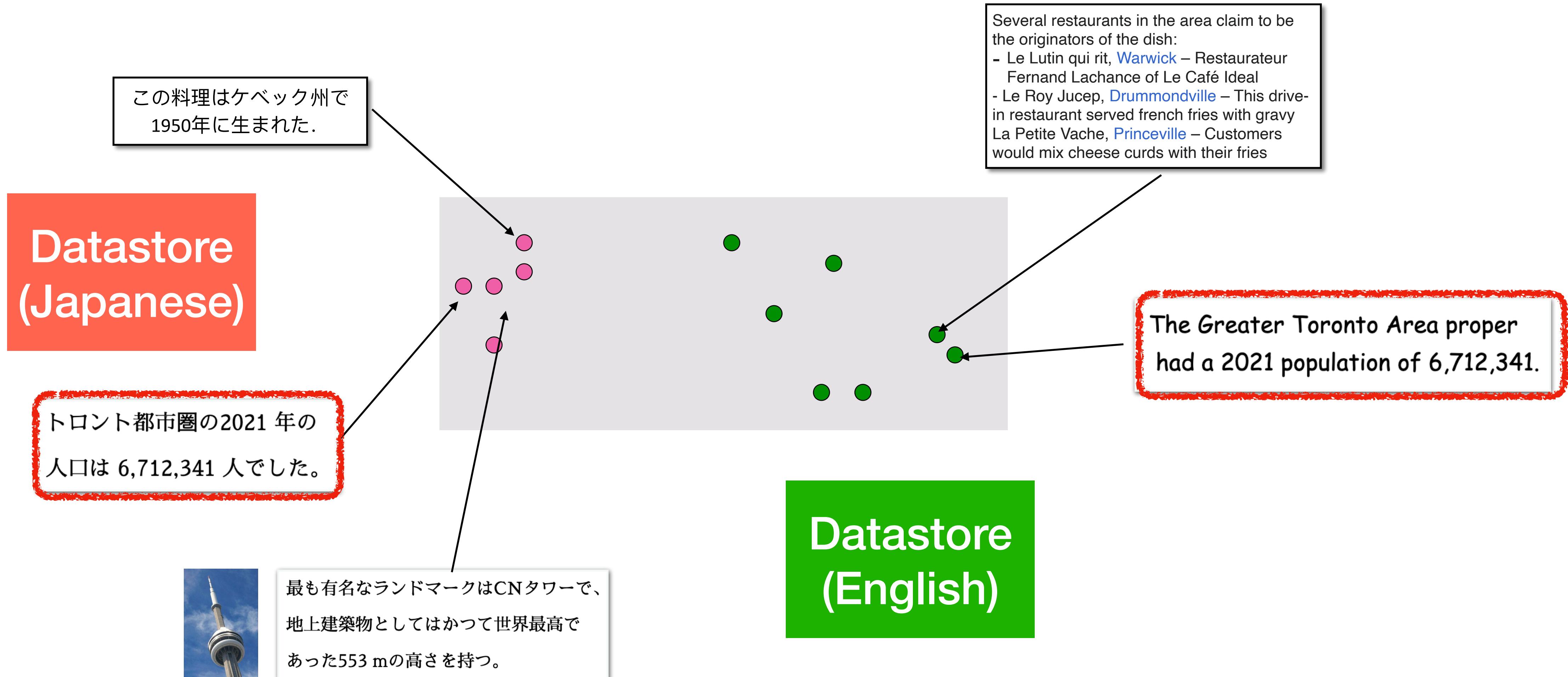
Multilingual Retrieval-based LM



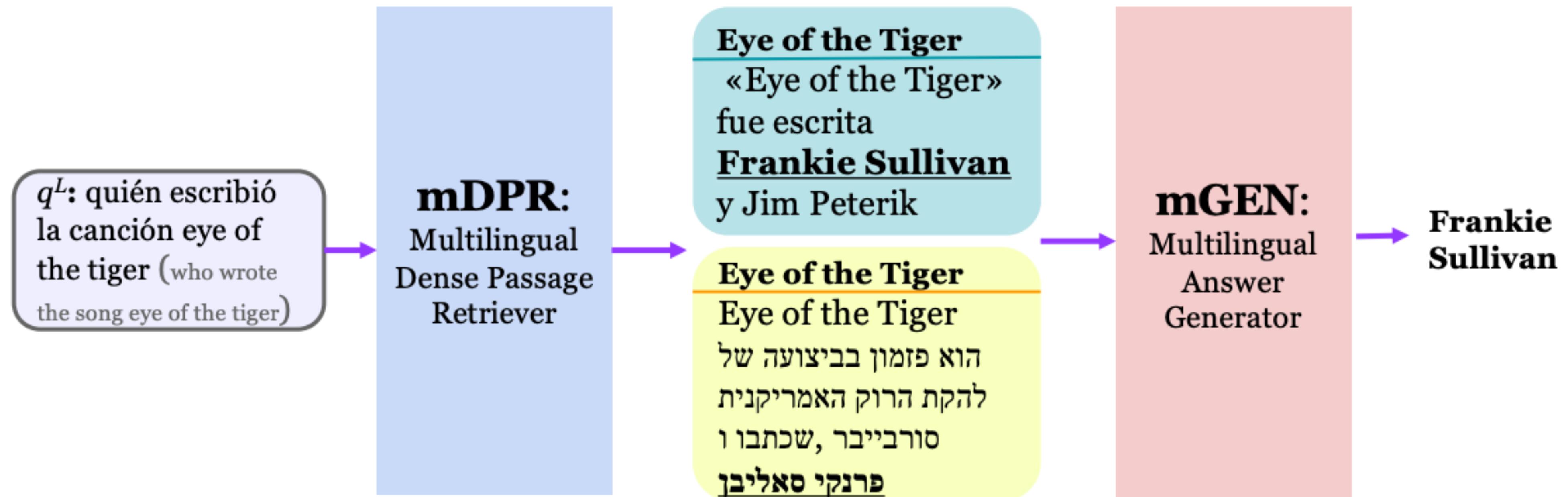
Language biases in representation spaces



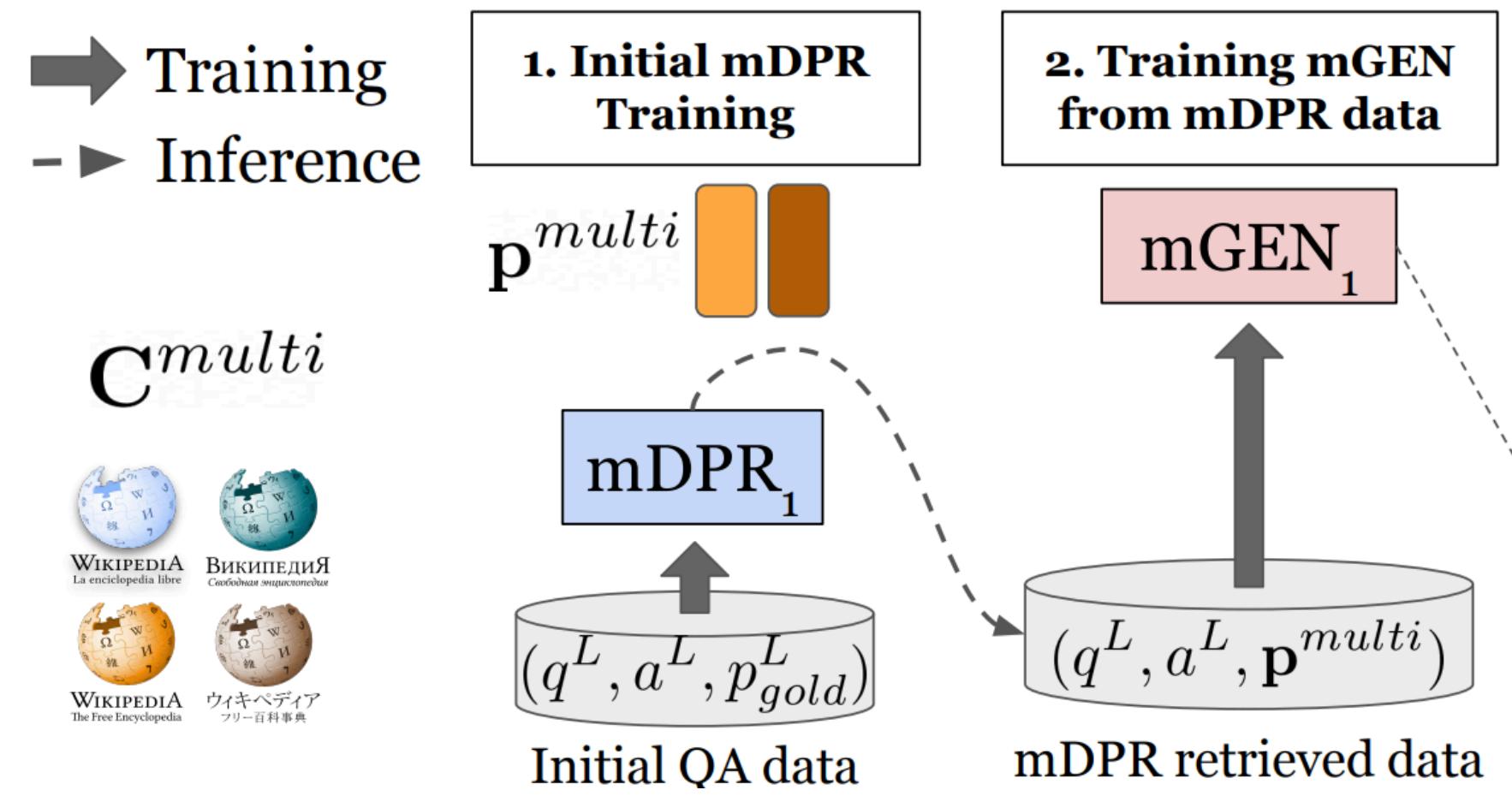
Language biases in representation spaces



CORA (Asai et al., 2021)

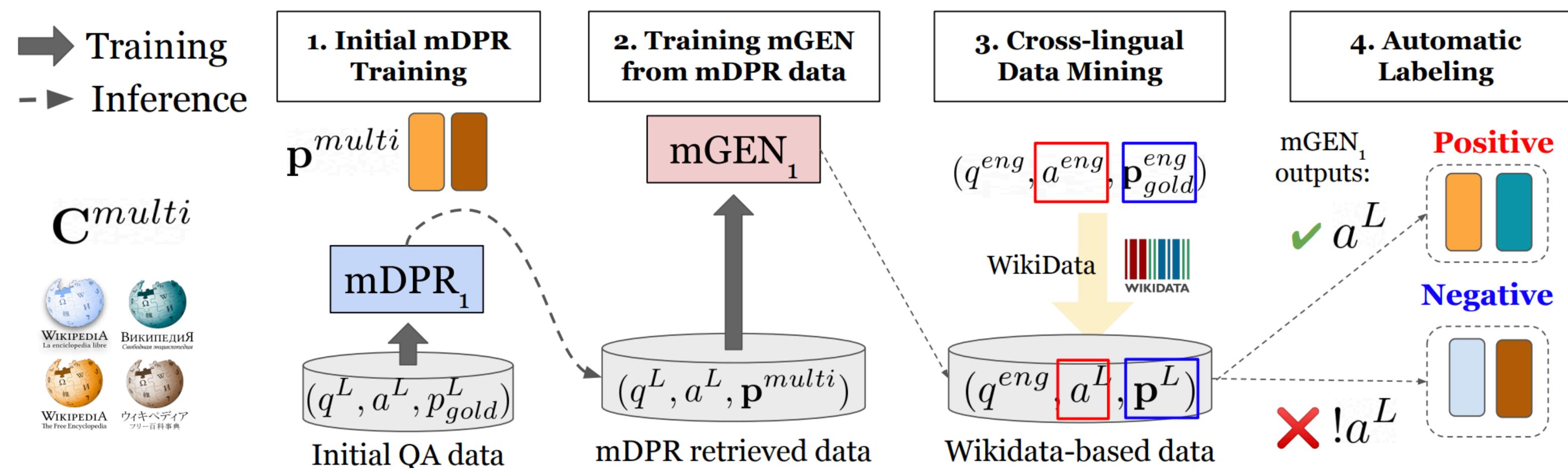


CORA: iteratively training multilingual LM & retriever



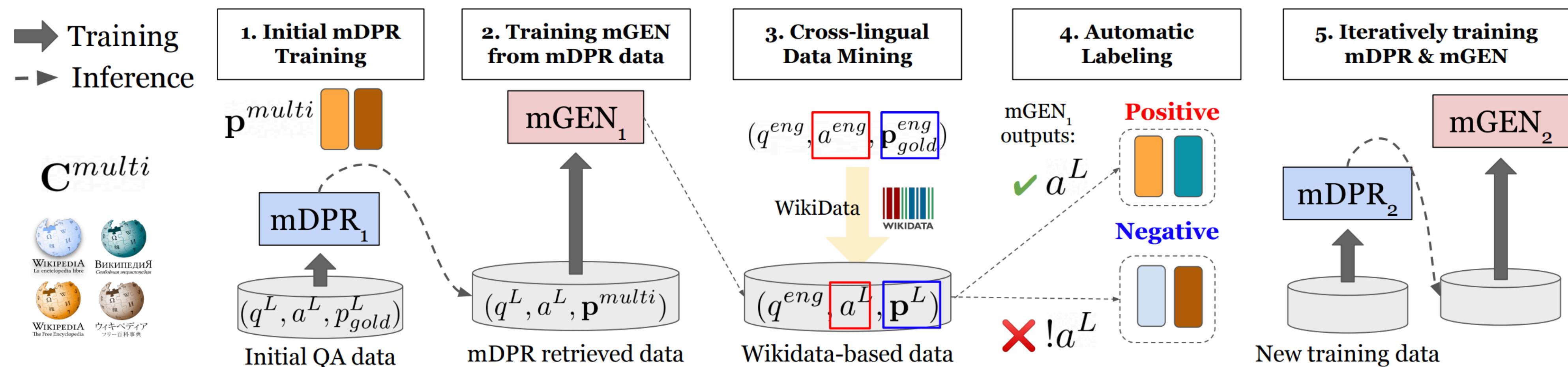
Initial training of retriever and LM

CORA: iteratively training multilingual LM & retriever

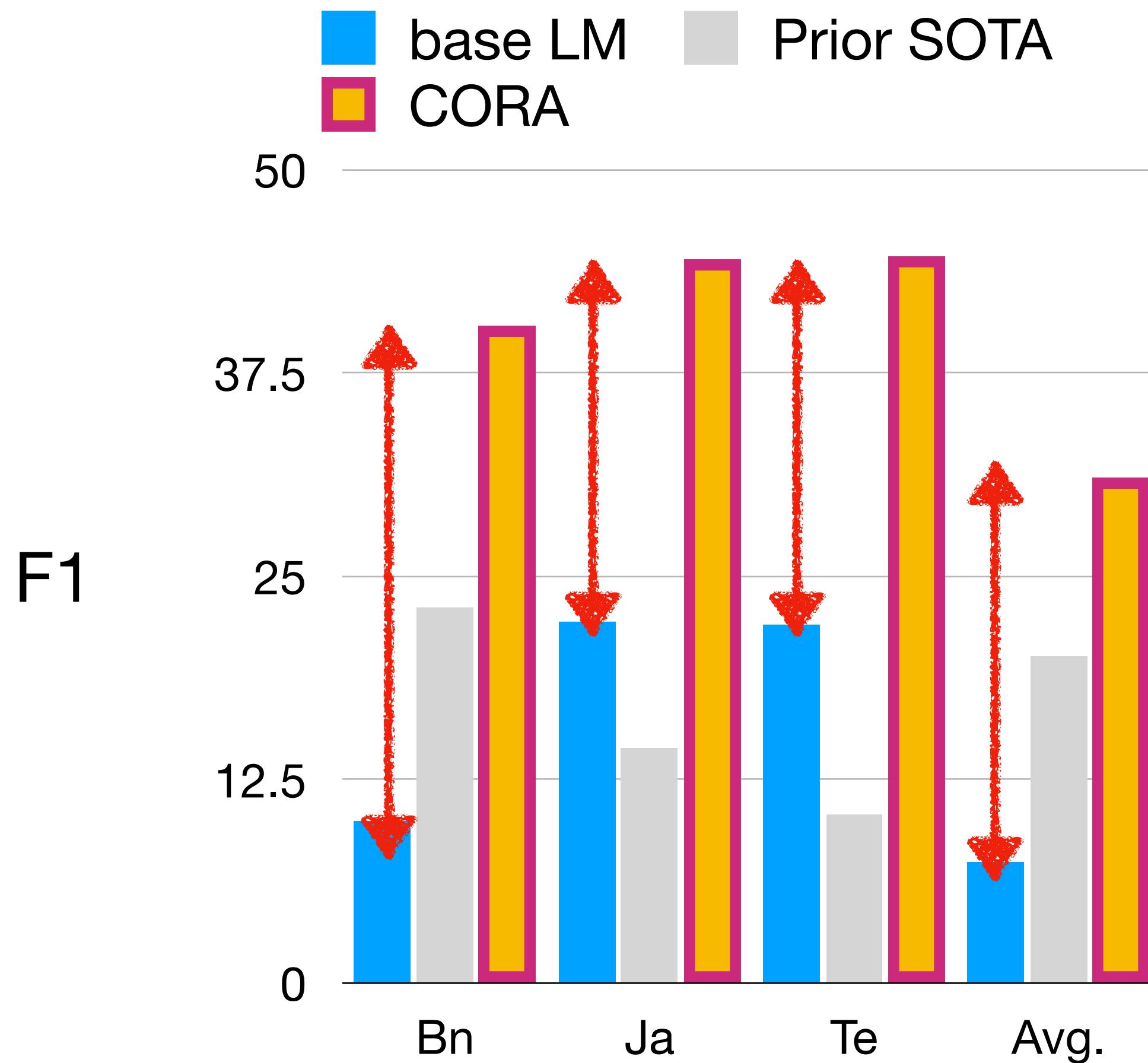


Expand training data using trained models
as well as structured cross-lingual data.

CORA: iteratively training multilingual LM & retriever

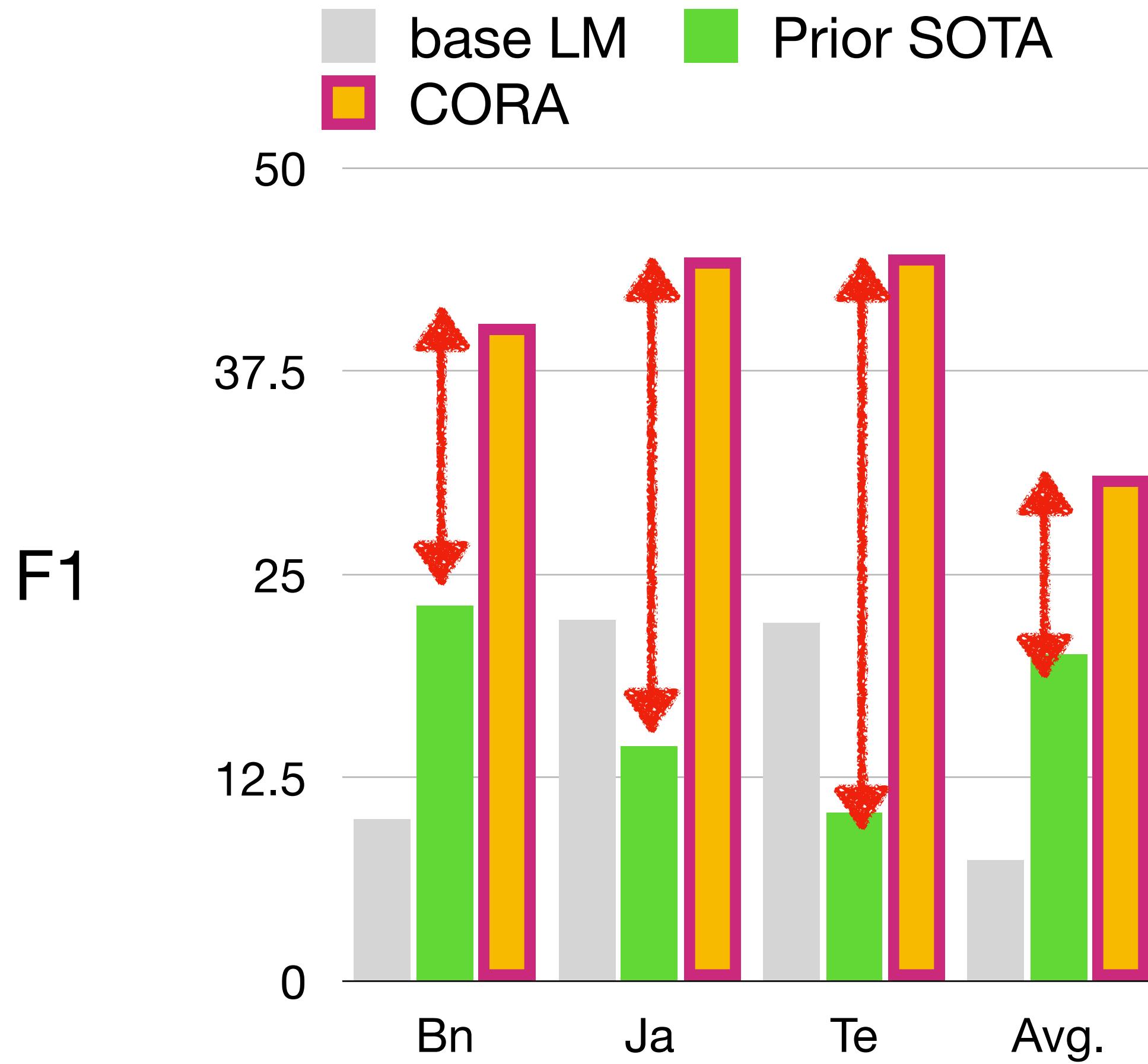


Results



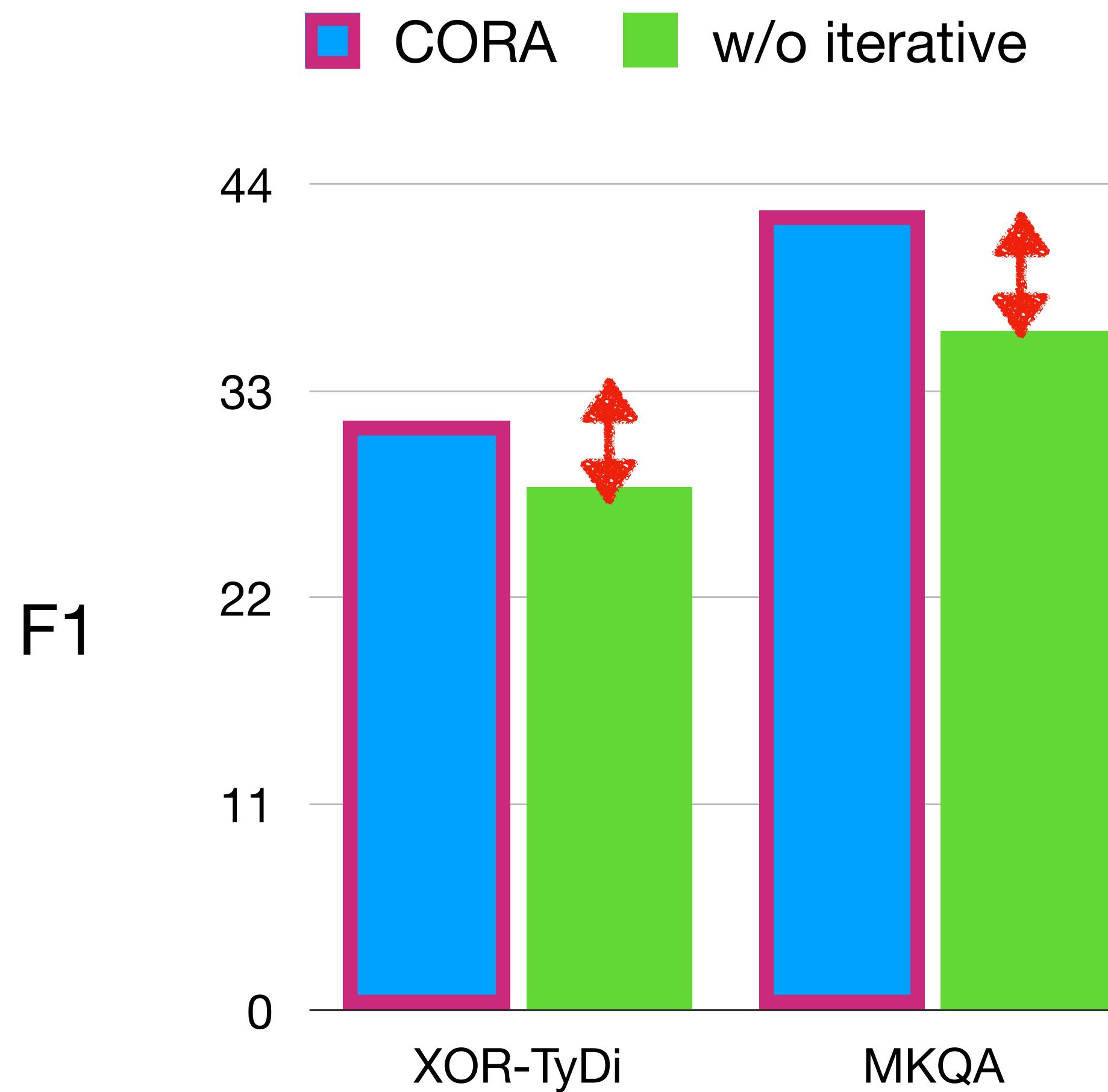
Large gains from fine-tuned LM
without retrieval

Results



Significantly Outperforms prior SOTA

Results



Iterative training of retriever and LM
gives large performance improvements

Multilingual retrieval-based LMs for diverse tasks

Question Answering

- * CL-ReLKT (Limkonchotiwat et al., 2022): knowledge transfer for better cross-lingual retrieval training
- * Gen-TyDi QA (Muller et al., 2023): generate full sentence answers for cross-lingual QA.

Fact Verification

- * CONCRETE (Hung et al., 2022): Improving cross-lingual fact-checking with cross-lingual retrieval

Dialogue

- * Cross-lingual Knowledge-grounded Dialogue (Kim et al 2021): a Korean knowledge-grounded dialogue system that learns to generate Korean response given English & Korean knowledge

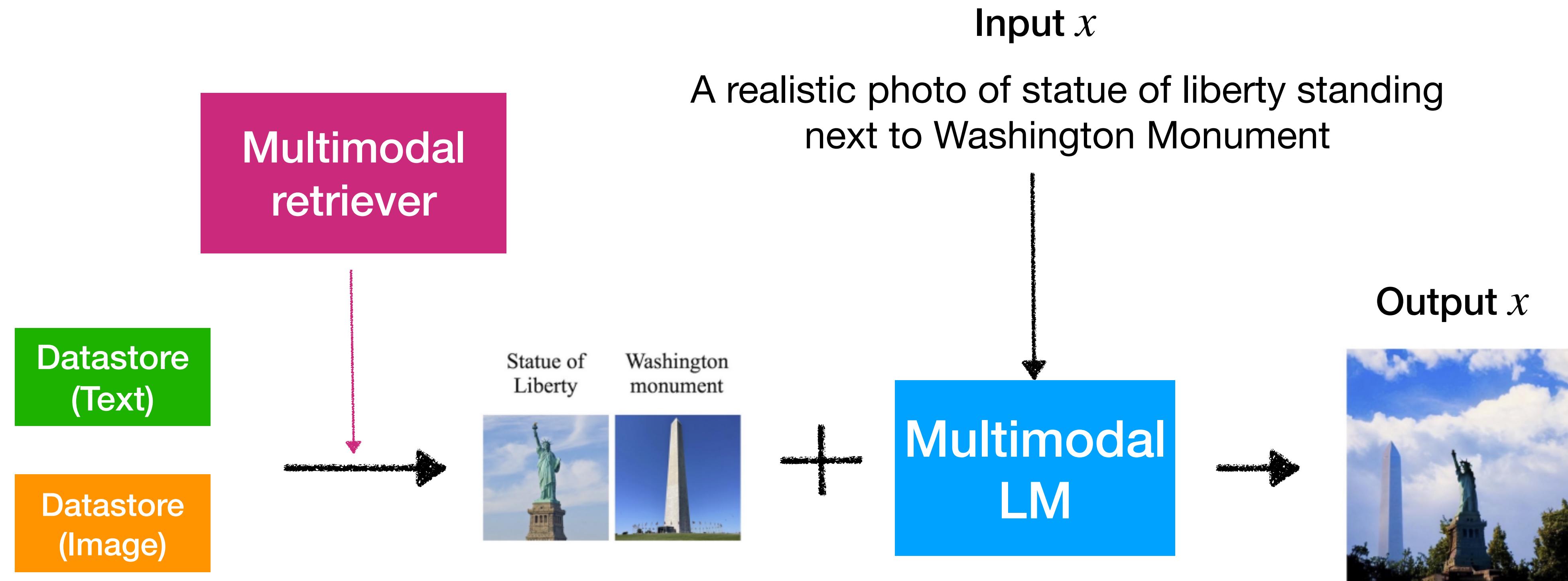
Event Extraction

- * R-GQA (Du and Ji, 2022): retrieve similar QA pairs for event argument extraction.

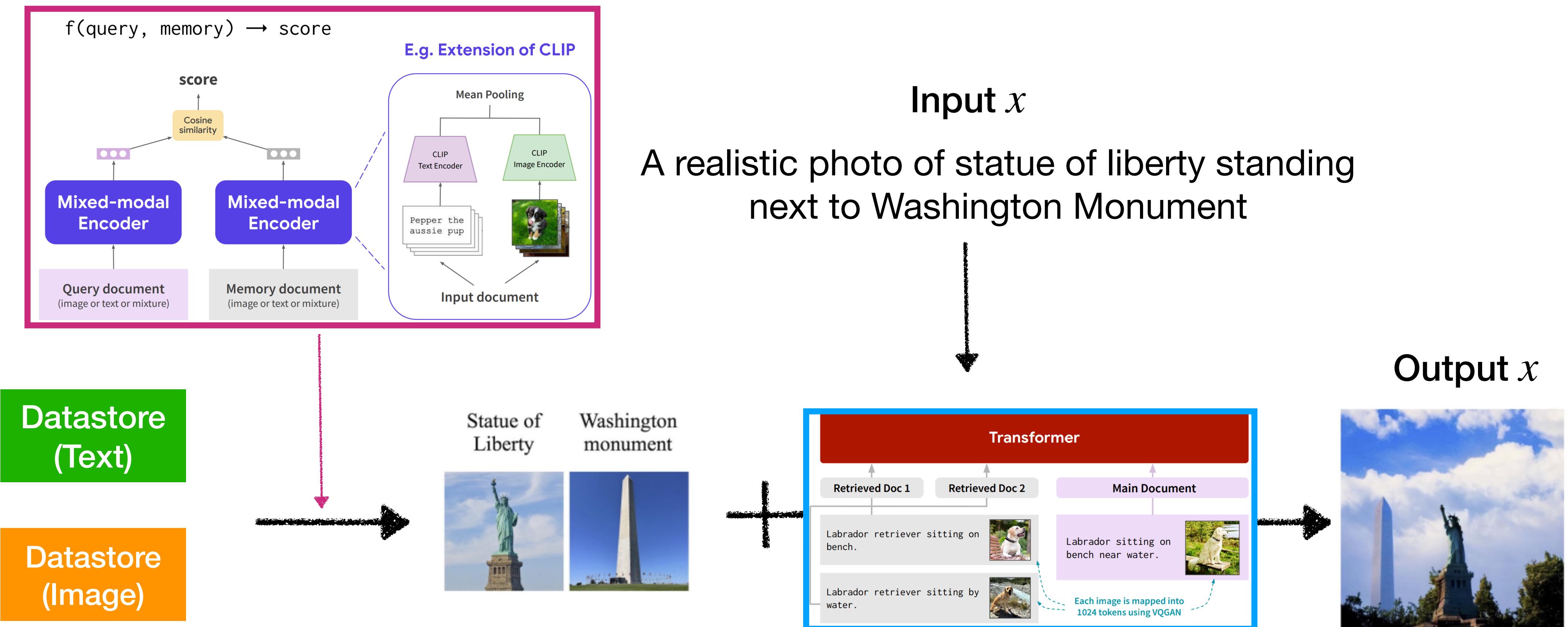
Key-phrase generations

- * Retrieval-augmented Multilingual Key phrase Generation (Gao et al 2022): Using iterative training to improve retrieval & LM for key phrase generations

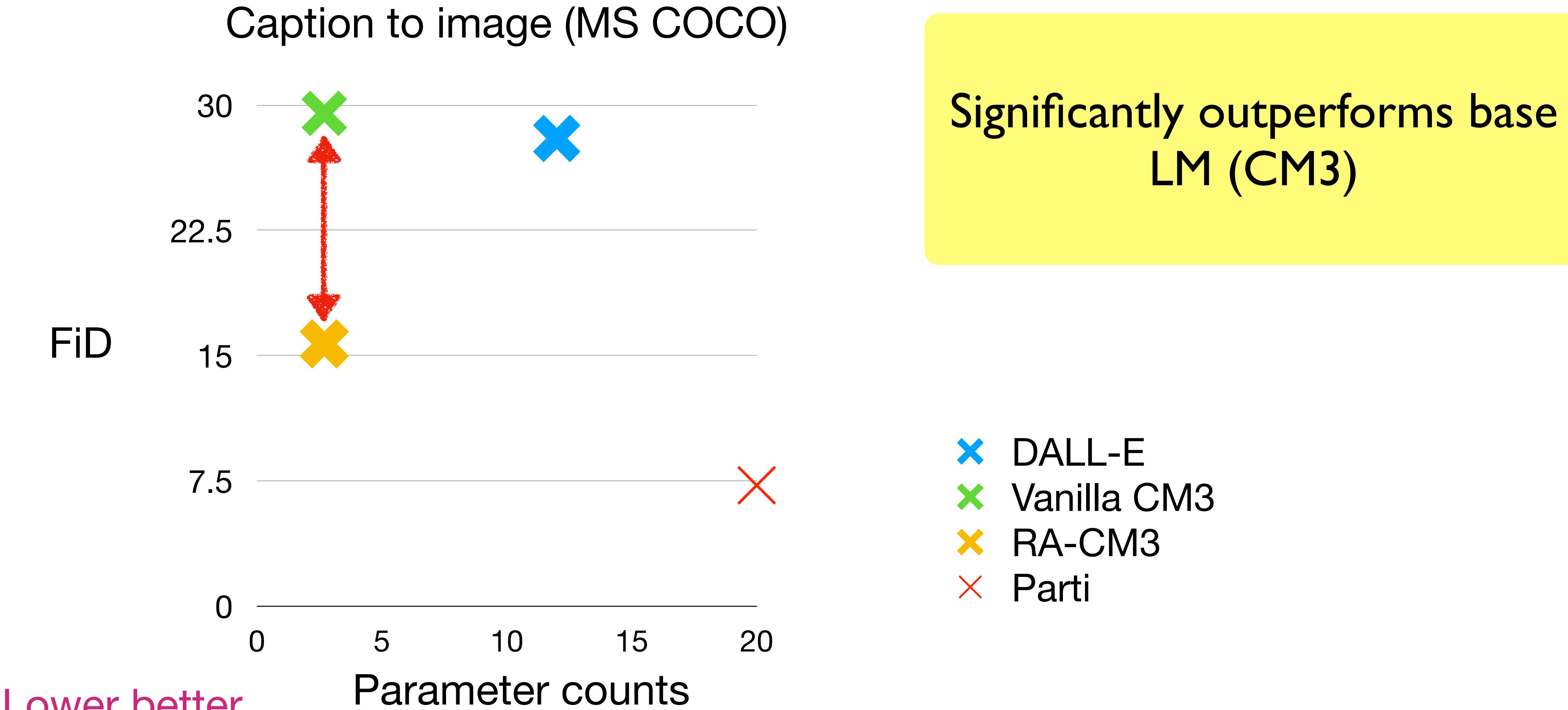
Multi-modal retrieval-augmented generations



RA-CM3 (Yasunaga et al., 2023)

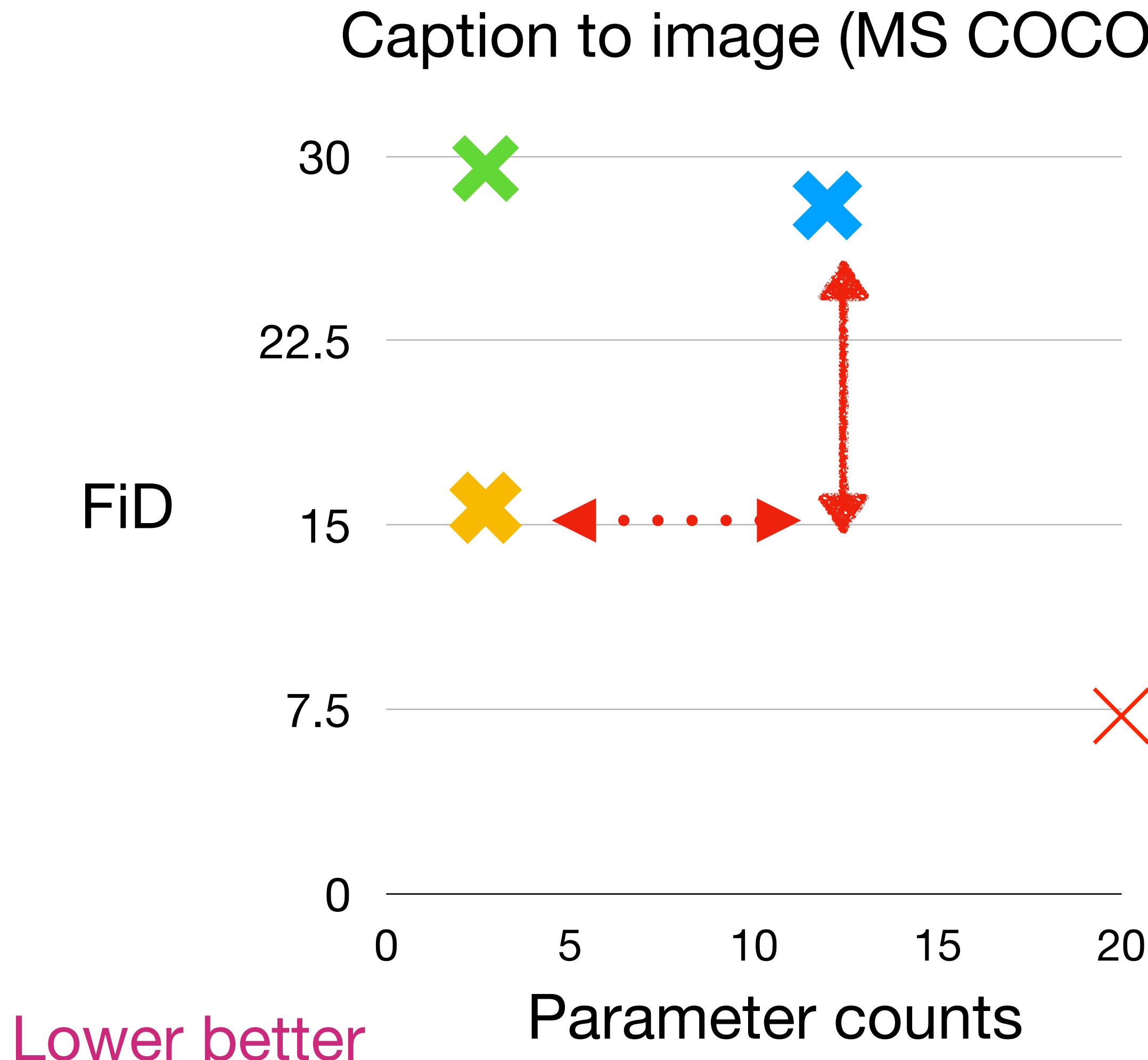


Results



Yasunaga et al. 2023. "Retrieval-Augmented Multimodal Language Modeling"

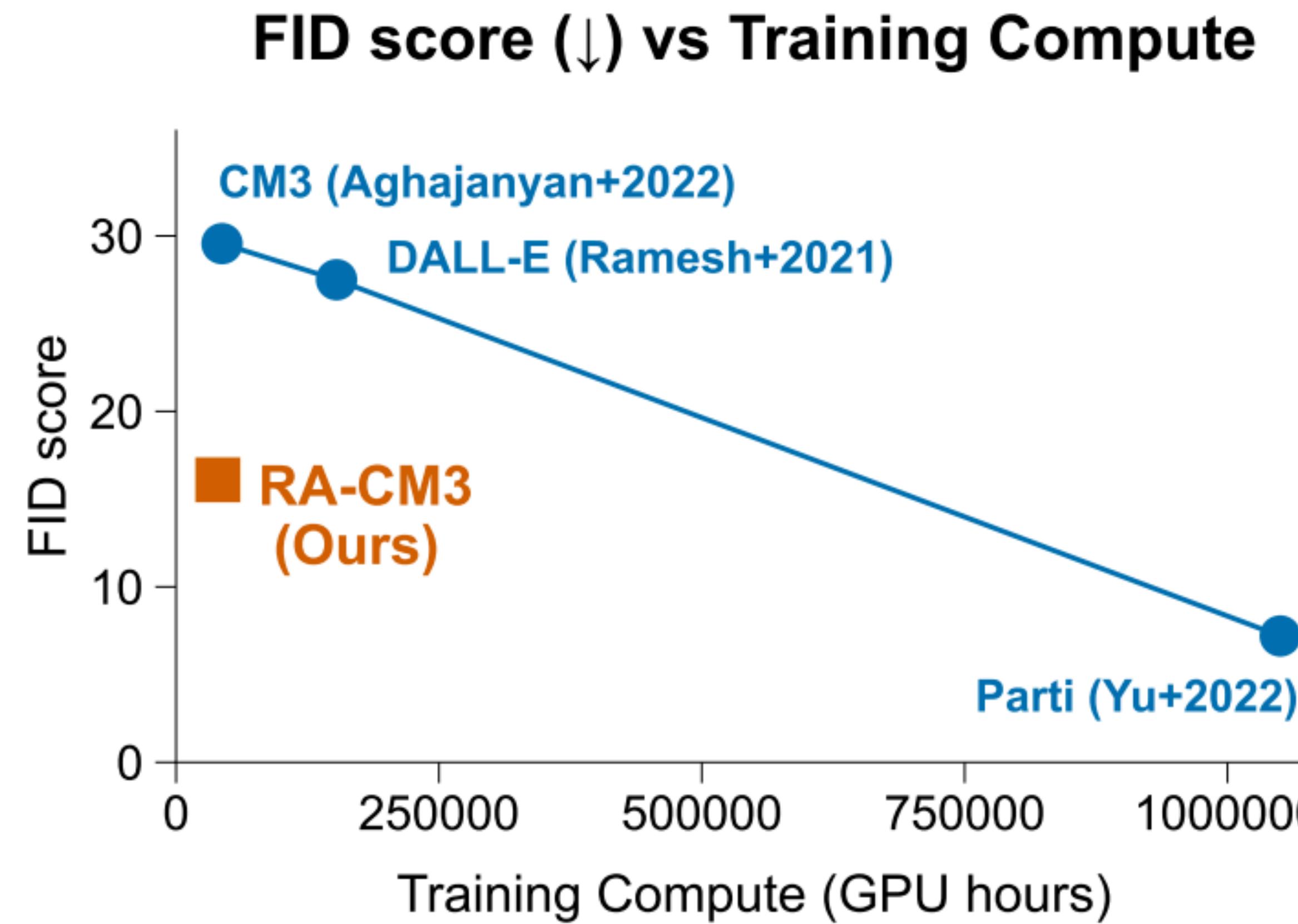
Results



Yasunaga et al. 2023. "Retrieval-Augmented Multimodal Language Modeling"

Results

Caption to image (MS COCO)



Retrieval-augmented generation achieves significantly better training efficiency

More applications beyond text

Multi-modal Retrieval-augmented Pre-training

- * RAVEAL (Hu et al 2023): Pretraining visual-language model using knowledge memory

Multi-modal Question Answering

- * MuRAG (Chen et al., 2022)

Multi-modal Classification

- * ALIGN (Gur et al., 2021)

Multimodal using image and text have been actively studied

More applications beyond text

Multi-modal Retrieval-augmented Pre-training

- * RAVEAL (Hu et al 2023): Pretraining visual-language model using knowledge memory

Multi-modal Question Answering

- * MuRAG (Chen et al., 2022)

Multi-modal Classification

- * ALIGN (Gur et al., 2021)

Retrieval-augmented training for molecules

- * Retrieval-based Molecule Generation (Wang et al., 2023)

Retrieval-augmented 3D motion generations

- * ReMoDiffus (Zhang et al., 2023)

New extensions for new input / output modality!

Wrapping up ...

The diagram illustrates the evolution of data from a raw text corpus to a knowledge base, and finally to multimodal extensions. It consists of several circular nodes connected by lines, forming a network. The top node is labeled 'KB' (Knowledge Base). Below it are four main nodes: 'Raw text corpus', 'Images', 'Video', and 'Audio'. Each of these nodes contains sub-nodes representing specific types of data. The 'Raw text corpus' node contains text snippets in English, Spanish, and Japanese. The 'Images' node contains a map of a region and a street-level photo. The 'Video' node contains a film strip with video frames. The 'Audio' node contains a waveform visualization.

- Raw text corpus**
 - The Greater Toronto Area proper had a 2021 population of 6,712,341
 - Toronto (pronunciación en inglés: /tə'jɑntoʊ/ escuchar), localmente /tə'janoo/, /t̬janoo/) es la capital de la provincia de Ontario³
 - トロント都市圏の2021 年の人口は 6,712,341 人でした。
- Images**
 - A map of a region showing water bodies and land areas.
 - A street-level photograph of a city street with people walking and buildings in the background.
- Video**
 - A film strip showing a sequence of video frames of people playing baseball or softball.
- Audio**
 - An audio waveform visualization.

Extension to multilingual

Cross-lingual retrieval and generation to overcome **datastore scarcity** in many world languages

Extension to multimodal

Key effectivenesses (Section 5; long-tail, efficiency) apply to diverse modality

References (I)

Akari Asai, Xinyan Yu, Jungo Kasai, Hannaneh Hajishirzi. One Question Answering Model for Many Languages with Cross-lingual Dense Passage Retrieval. *NeurIPS* 2021.

Peerat Limkonchotiwat, Wuttikorn Ponwitayarat, Can Udomcharoenchaikit, Ekapol Chuangsawanich, Sarana Nutanong. CL-ReLKT: Cross-lingual Language Knowledge Transfer for Multilingual Retrieval Question Answering. *NAACL Findings* 2022.

Benjamin Muller, Luca Soldaini, Rik Koncel-Kedziorski, Eric Lind, Alessandro Moschitti. Cross-Lingual GenQA: A Language-Agnostic Generative Question Answering Approach for Open-Domain Question Answering. *AACL* 2022.

Kung-Hsiang Huang, ChengXiang Zhai, Heng Ji. CONCRETE: Improving Cross-lingual Fact-checking with Cross-lingual Retrieval. *COLING* 2022.

Yifan Gao, Qingyu Yin, Zheng Li, Rui Meng, Tong Zhao, Bing Yin, Irwin King, Michael R. Lyu. Retrieval-Augmented Multilingual Keyphrase Generation with Retriever-Generator Iterative Training. *NAACL Findings* 2022.

References (2)

- Ziniu Hu, Ahmet Iscen, Chen Sun, Zirui Wang, Kai-Wei Chang, Yizhou Sun, Cordelia Schmid, David A. Ross, Alireza Fathi. REVEAL: Retrieval-Augmented Visual-Language Pre-Training with Multi-Source Multimodal Knowledge Memory. *CVPR* 2023.
- Wenhu Chen, Hexiang Hu, Xi Chen, Pat Verga, William W. Cohen. MuRAG: Multimodal Retrieval-Augmented Generator for Open Question Answering over Images and Text. *EMNLP* 2022.
- Shir Gur, Natalia Neverova, Chris Stauffer, Ser-Nam Lim, Douwe Kiela, Austin Reiter. Cross-Modal Retrieval Augmentation for Multi-Modal Classification. *EMNLP Findings* 2021.
- Zichao Wang, Weili Nie, Zhuoran Qiao, Chaowei Xiao, Richard Baraniuk, Anima Anandkumar. Retrieval-based Controllable Molecule Generation. *ICLR* 2023.
- Mingyuan Zhang, Xinying Guo, Liang Pan, Zhongang Cai, Fangzhou Hong, Huirong Li, Lei Yang, Ziwei Liu. ReMoDiffuse: Retrieval-Augmented Motion Diffusion Model. Arxiv 2023.