DeepMind

# RETRO: Improving language models by retrieving from trillions of tokens
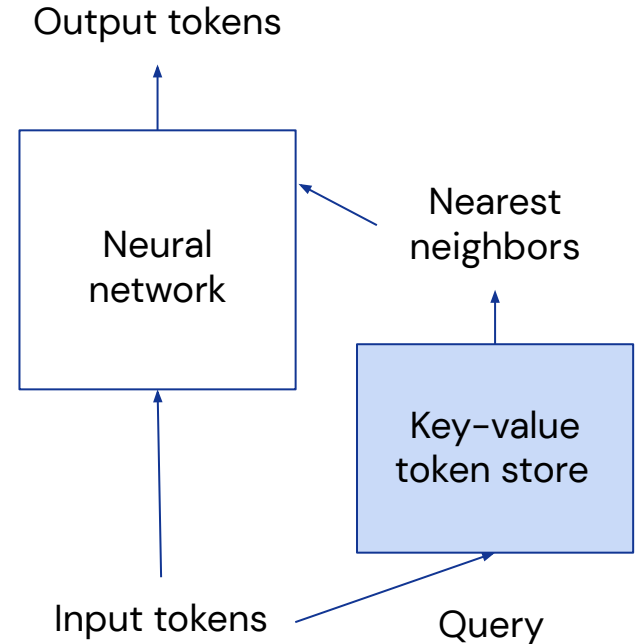
**Sebastian Borgeaud\*, Arthur Mensch\*, Jordan Hoffmann\***, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean–Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae[†], Erich Elsen[†], **Laurent Sifre**[\*,†]

\*Equal contributions, [†]Equal senior authorship

# Adding an explicit memory to language models

- Increasing the model size has two effects:
  - Increase **memorisation** of training data
  - Increase **generalisation** performance

- Can we increase memorisation without the extra parameters ?
  - External memory mechanism
  - Associated to a trainable neural network

- Combine a **parametric model** (neural network) with a **non-parametric model** (data store providing nearest neighbours)

Output tokens

Neural network

Nearest neighbors

Key-value token store

Input tokens

Query

# RETRO: adding a very large database of tokens

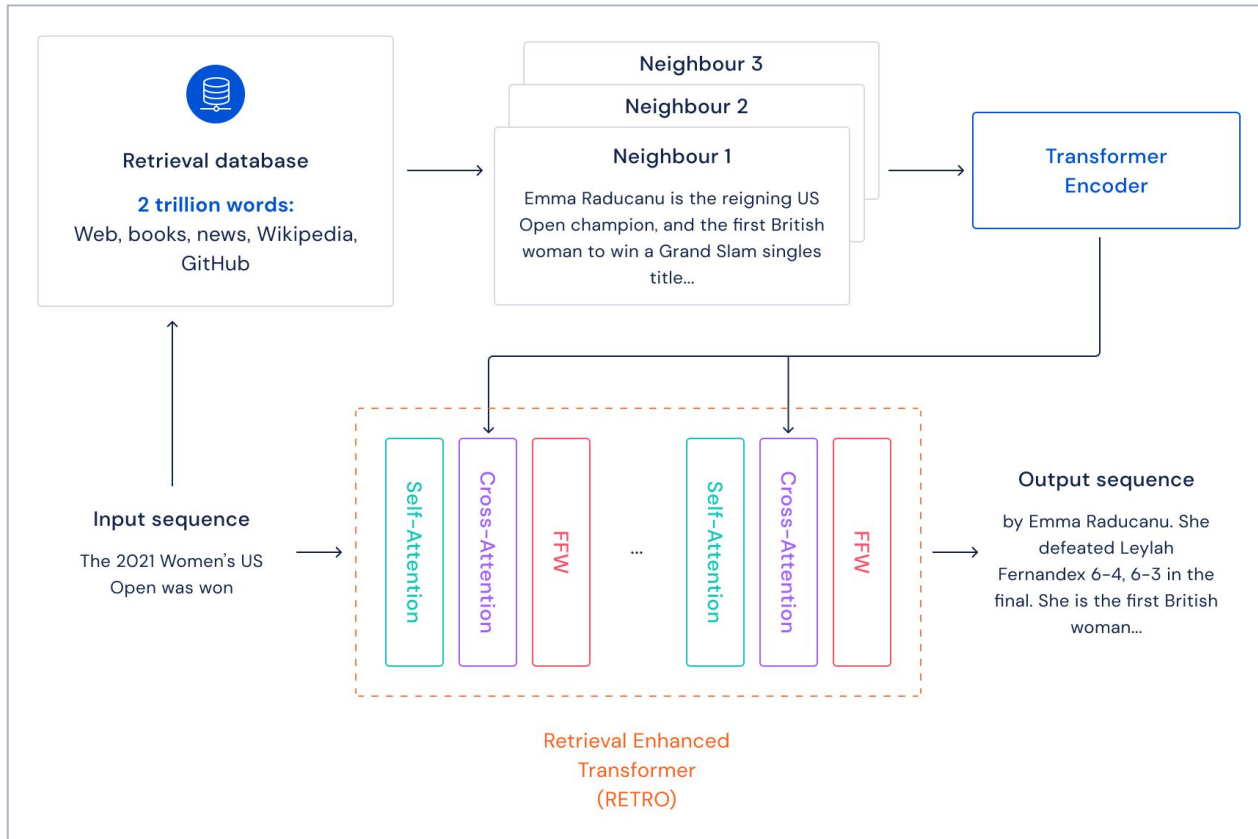Previous work on retrieval have focused on retrieving from Wikipedia (**3B tokens**)

We retrieve from our entire training set: **2T tokens (700x more)**
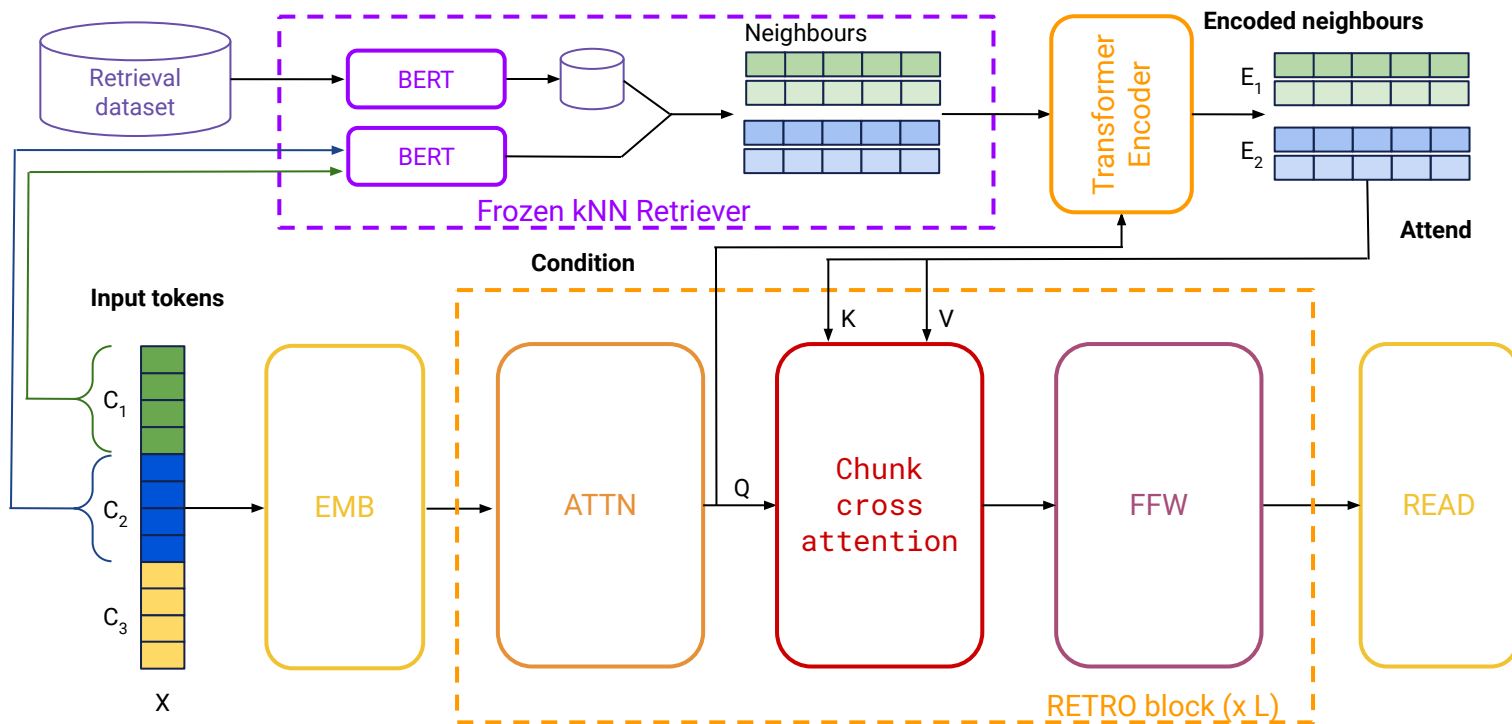
Three major ingredients:

a. **Frozen dense retriever** (pooled BERT on small sequences)
b. **Approx k-nearest neighbors** run as a preprocessing step (quantification and tree search)
c. **DB element is a chunk and not a token:** 30B keys

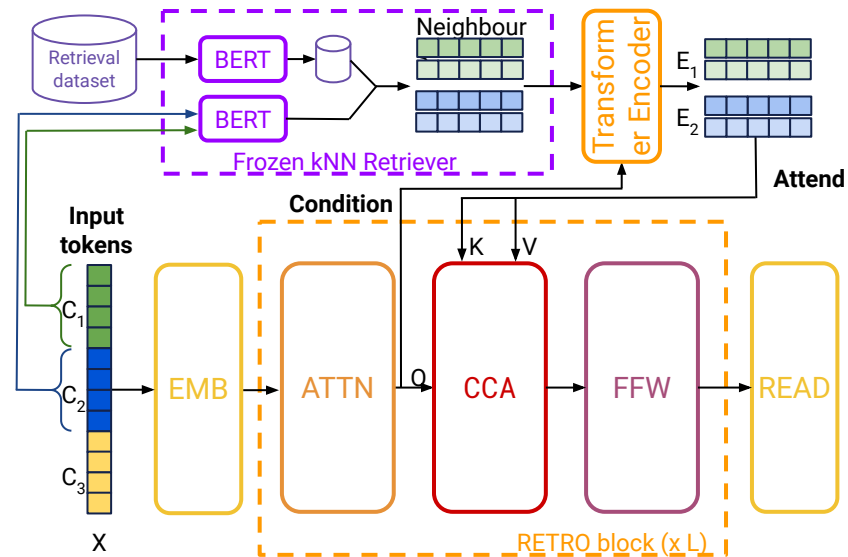# Retrieving at chunk level, predicting at sequence level

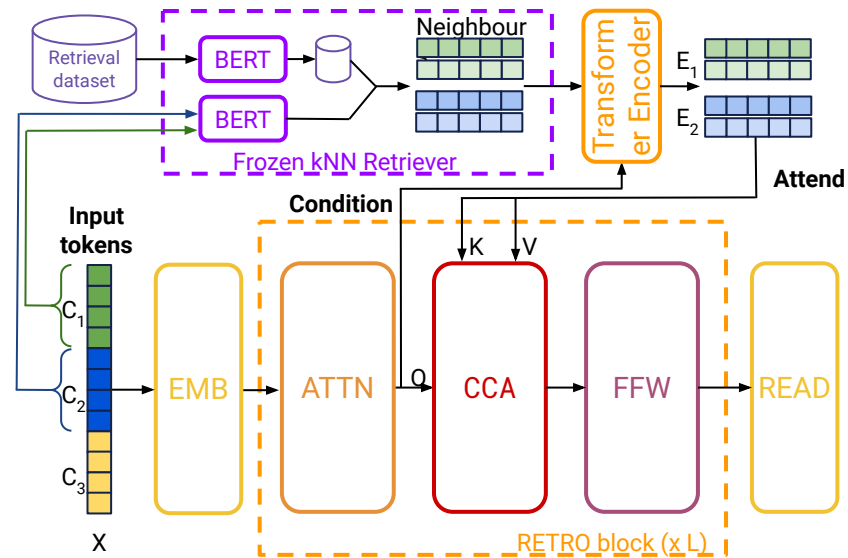# Condition new chunk generation on previous chunk neighbours

# A model made for sampling

- Fully autoregressive model

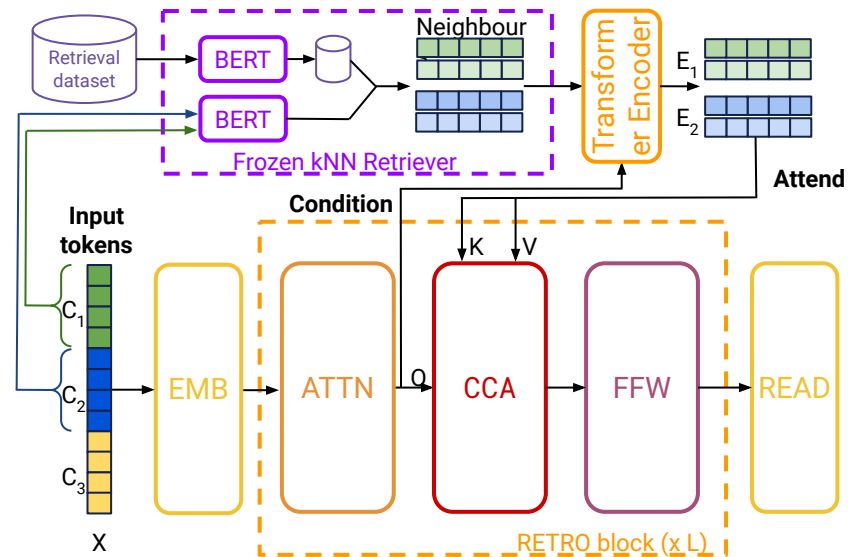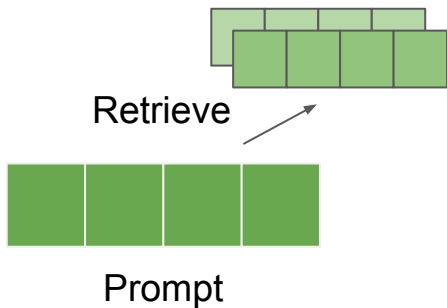- Sampling with retrieval queries in the loop

# A model made for sampling

- Fully autoregressive model

- Sampling with retrieval queries in the loop

# A model made for sampling

- Fully autoregressive model
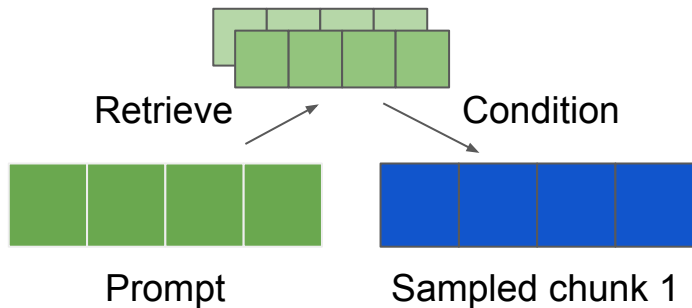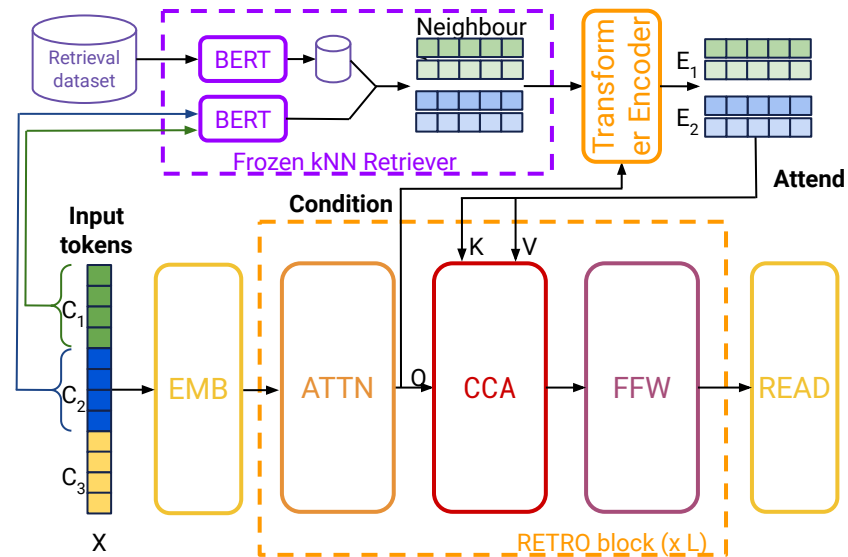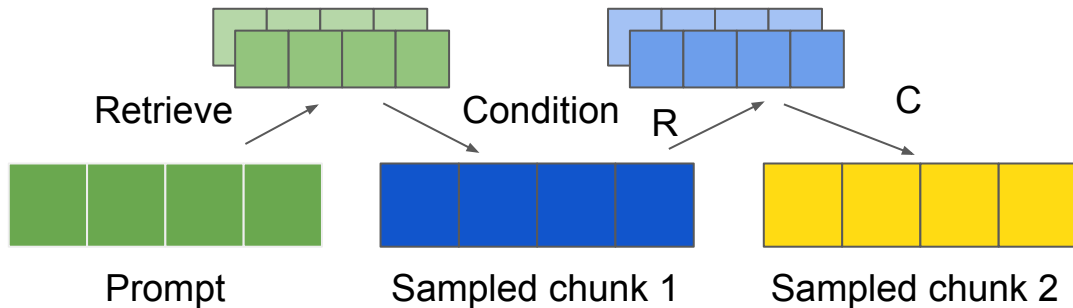
- Sampling with retrieval queries in the loop

# A model made for sampling

- Fully autoregressive model

- Sampling with retrieval queries in the loop

# A model made for sampling

- Fully autoregressive model
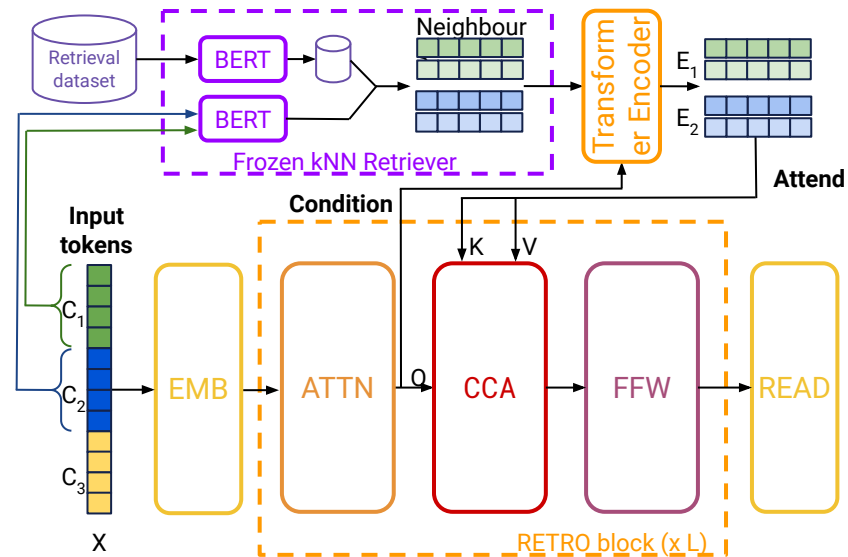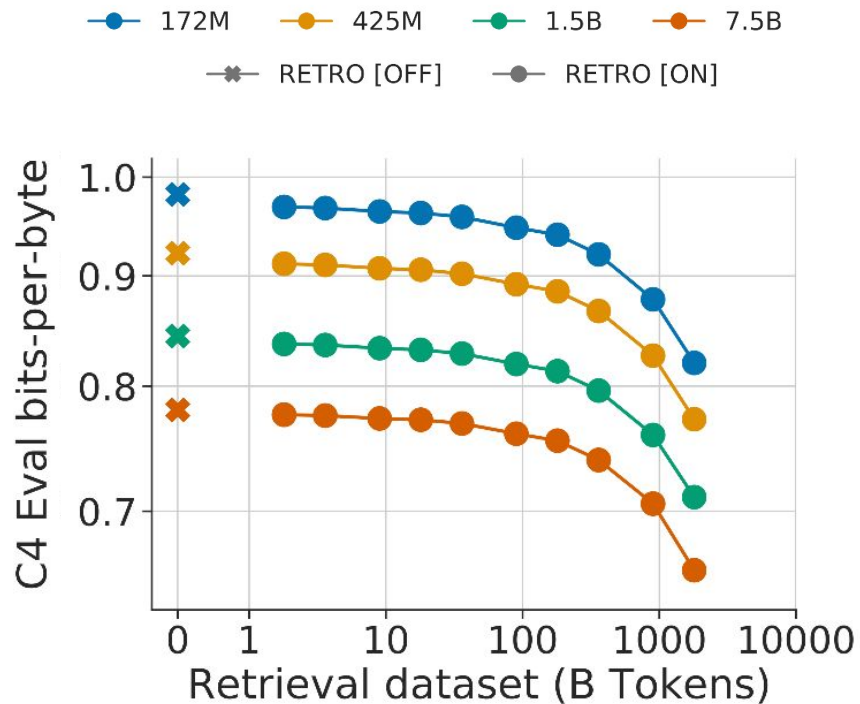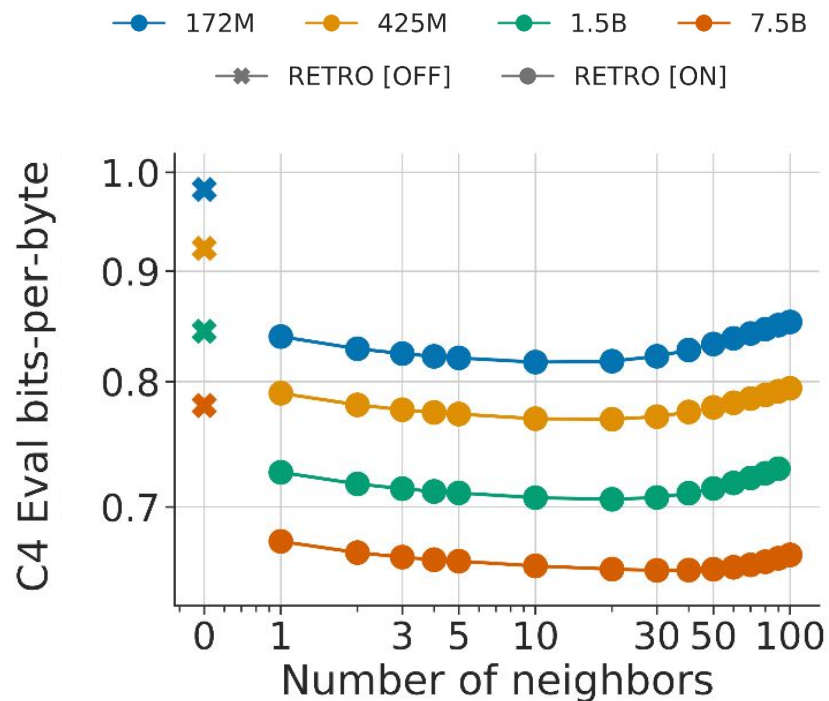
- Sampling with retrieval queries in the loop

# RETRO improves strongly with database size



Continuous improvement with database size

Improvements for all models

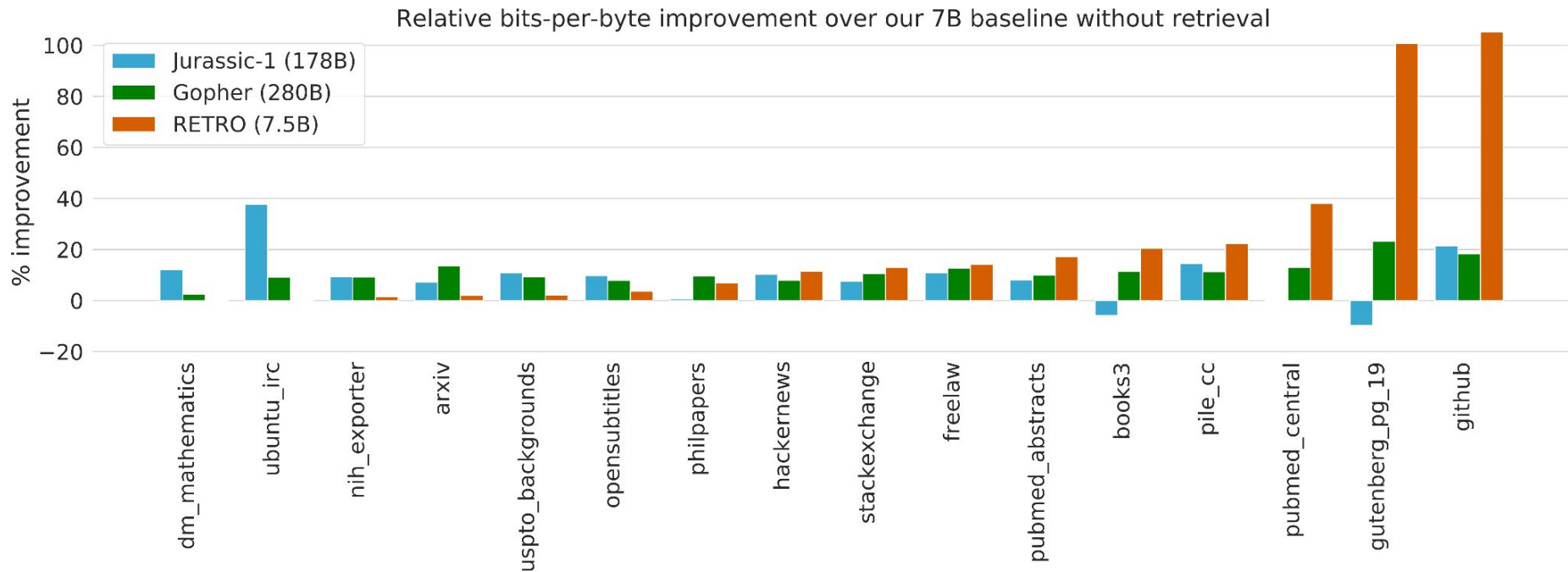# RETRO improves with additional retrieved neighbours



→ Trained with 2 neighbours

→ Improvements up to 40 neighbours at evaluation with 7.5B

→ Larger models can better utilise extra neighbours

# Language modelling: The Pile



Relative bits-per-byte improvement over our 7B baseline without retrieval

Legend:
- Jurassic-1 (178B)
- Gopher (280B)
- RETRO (7.5B)

Y-axis: % improvement

X-axis categories: dm_mathematics, ubuntu_irc, nih_exporter, arxiv, uspto_backgrounds, opensubtitles, philpapers, hackernews, stackexchange, freelaw, pubmed_abstracts, books3, pile_cc, pubmed_central, gutenberg_pg_19, github

# Conclusion

- **RETRO** is a **general** architecture, that is fully **autoregressive** and enables **large scale retrieval**

- Adding a **2T token database** yields a performance improvement that's constant with model size:
  - Similar performance to models with 10x more parameters on the Pile

- Consistent performance across benchmarks
  - Retrieval does exploit train–test leakage more than standard language models
  - But performance also improves on held–out tokens

- Future work on few-shot evaluation