

Section 2: Definition & Preliminaries

A Retrieval-based LM: Definition

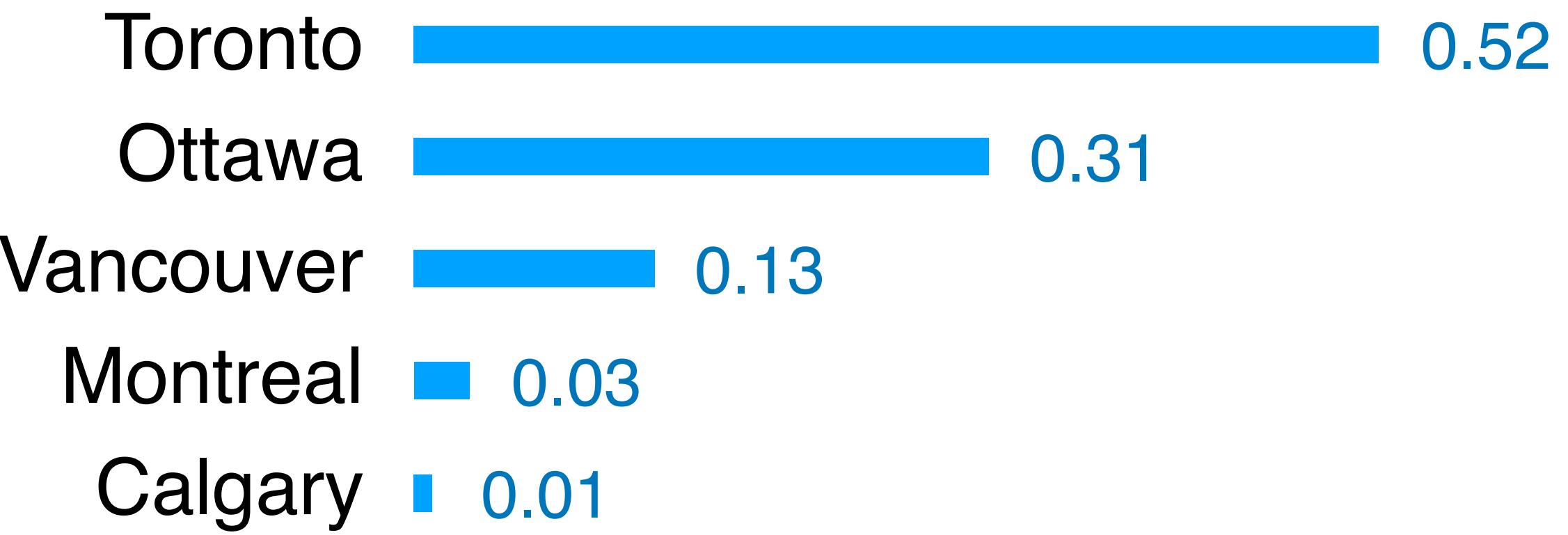
A language model (LM) that uses
an external datastore at test time

A Retrieval-based LM: Definition

**A language model (LM) that uses
an external datastore at test time**

A language model (LM)

$$P(x_n | x_1, x_2, \dots, x_{n-1})$$



Language model (Transformers)

The capital city of Ontario is

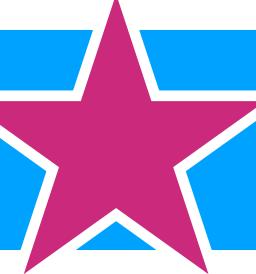
x_1

x_2

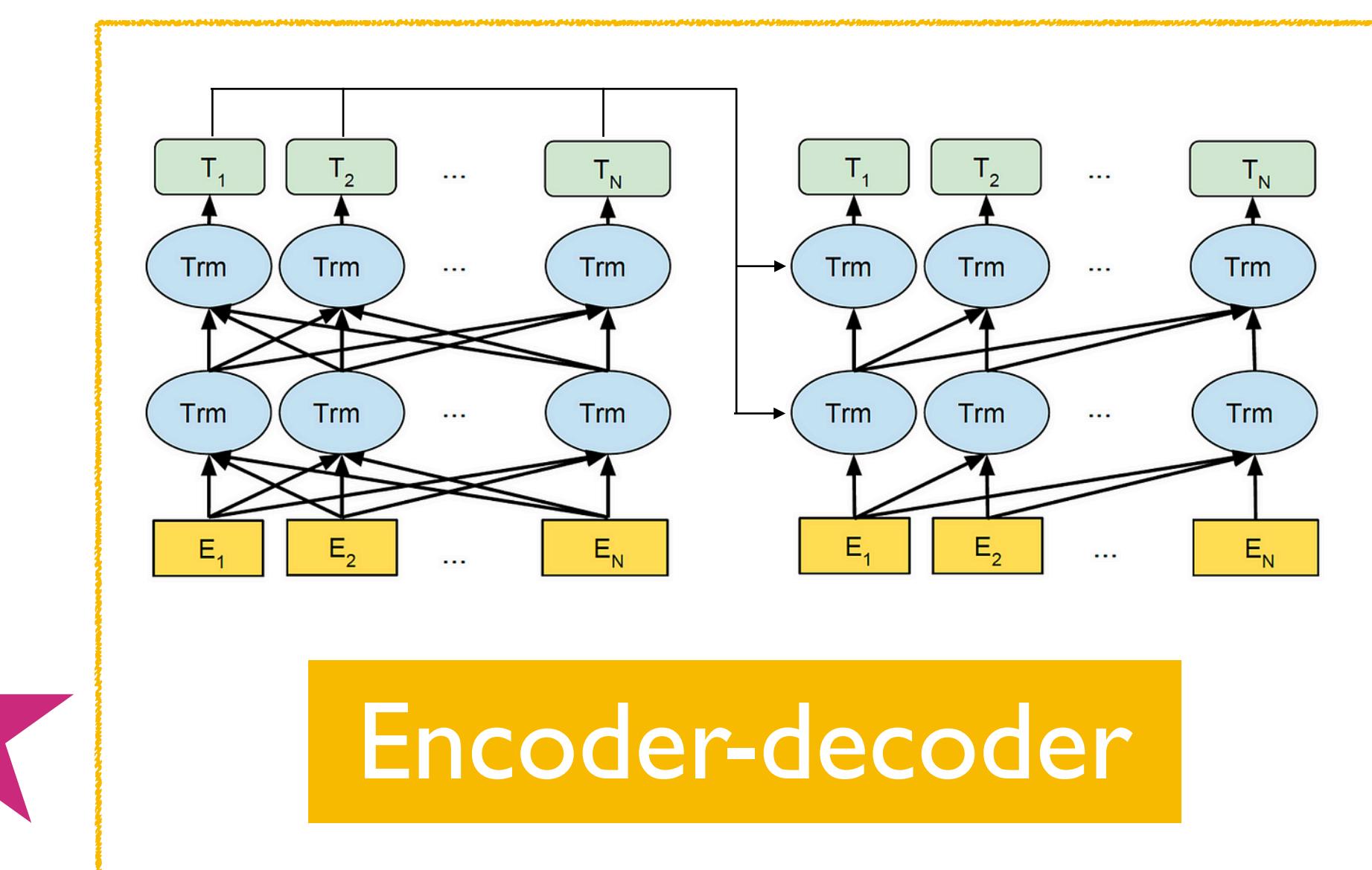
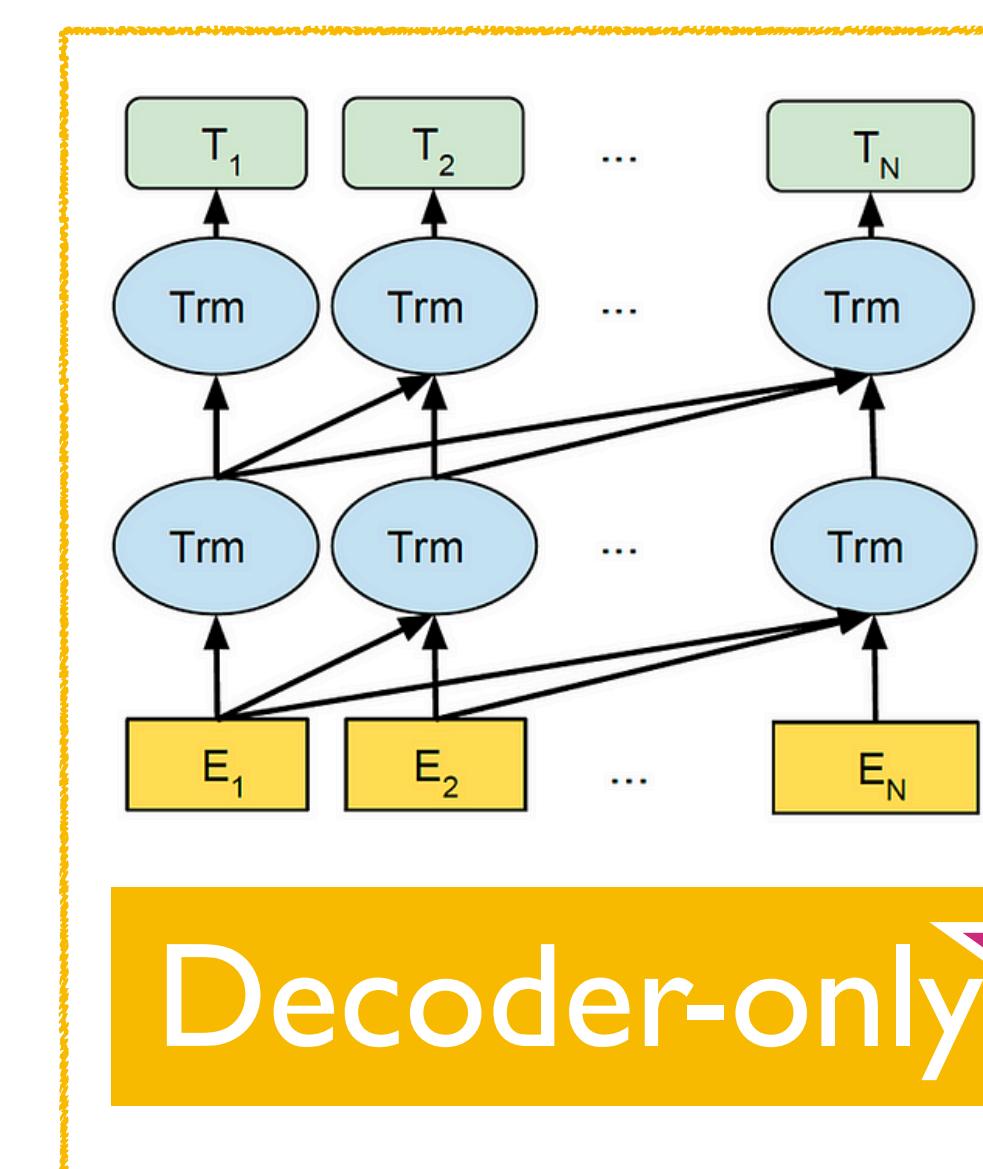
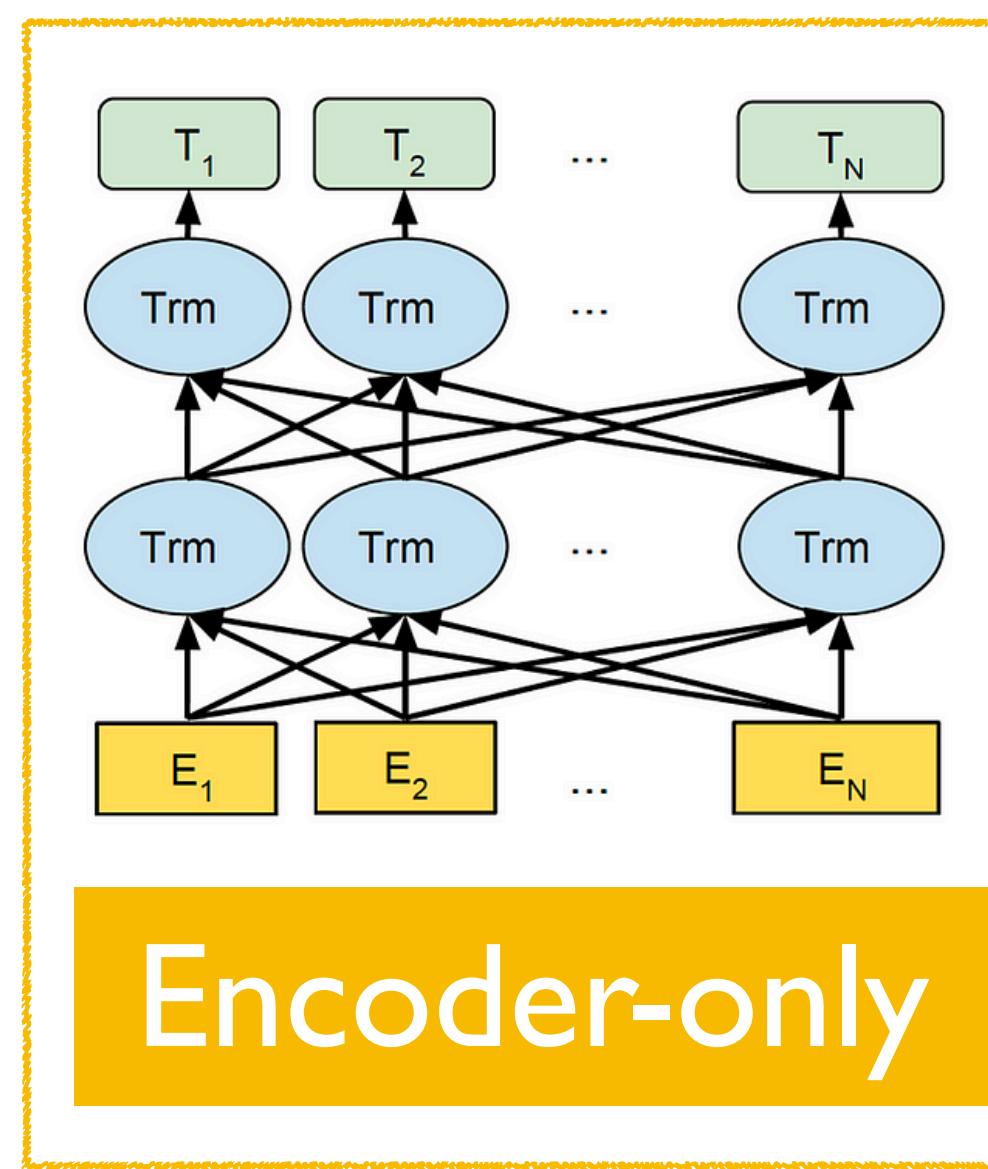
...

x_{n-1}

A language model (LM): Categories

Toronto
Autoregressive LM 
The capital city of Ontario is _____

capital Ontario
Masked LM
The _____ city of _____ is Toronto



A language model (LM): Prompting

The capital city of Ontario is

LM

Toronto

Fact probing

Cheaper than an iPod. It was

LM

great
terrible

Sentiment analysis

“Hello” in French is

LM

Bonjourno

Translation

I'm good at math. $5 + 8 \times 12 =$

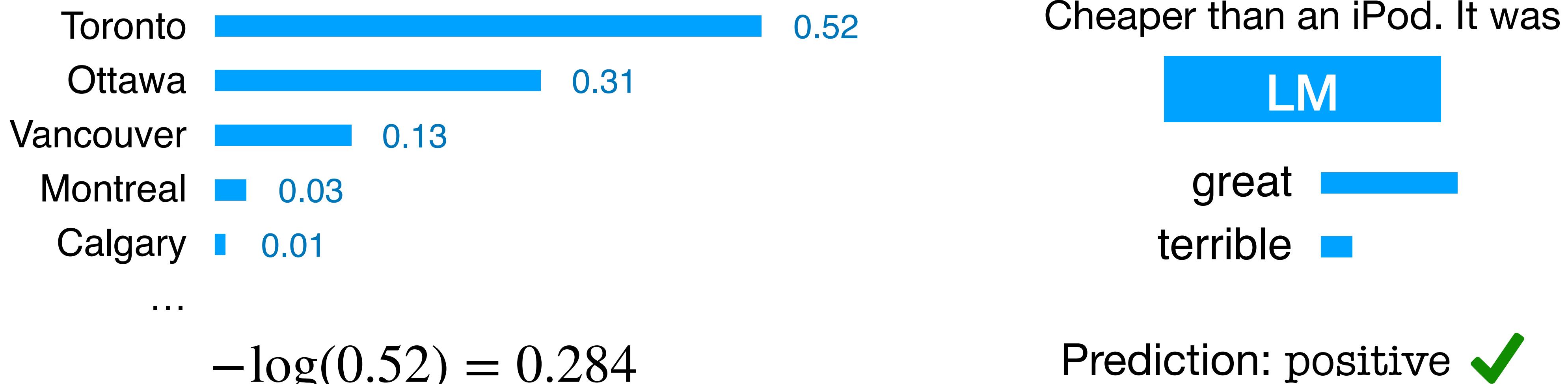
LM

101

Arithmetic

A language model (LM)

Often evaluated with



Perplexity

Downstream accuracy

(Zero-shot or few-shot in-context learning,
or fine-tuning)

(More in Section 5)

A Retrieval-based LM: Definition

**A language model (LM) that uses
an external datastore at test time**

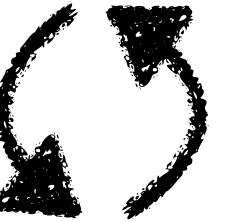
A Retrieval-based LM: Definition

A language model (LM) that uses
an external datastore at test time

External datastore



The capital city of Ontario is **Toronto**



LM

Training time

The capital city of Ontario is _____



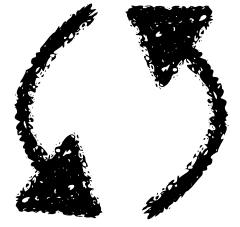
LM

Test time

External datastore



The capital city of Ontario is **Toronto**



LM

Training time



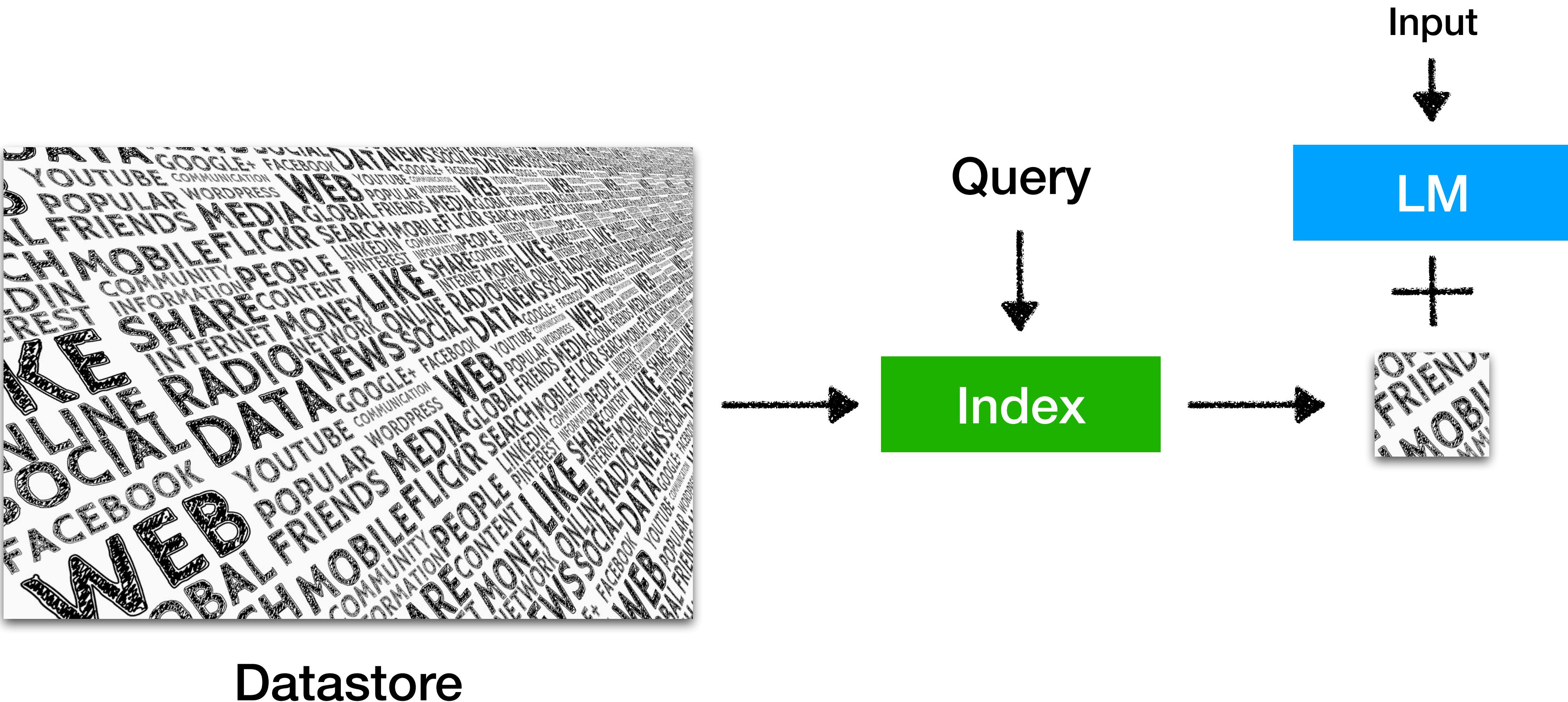
The capital city of Ontario is _____



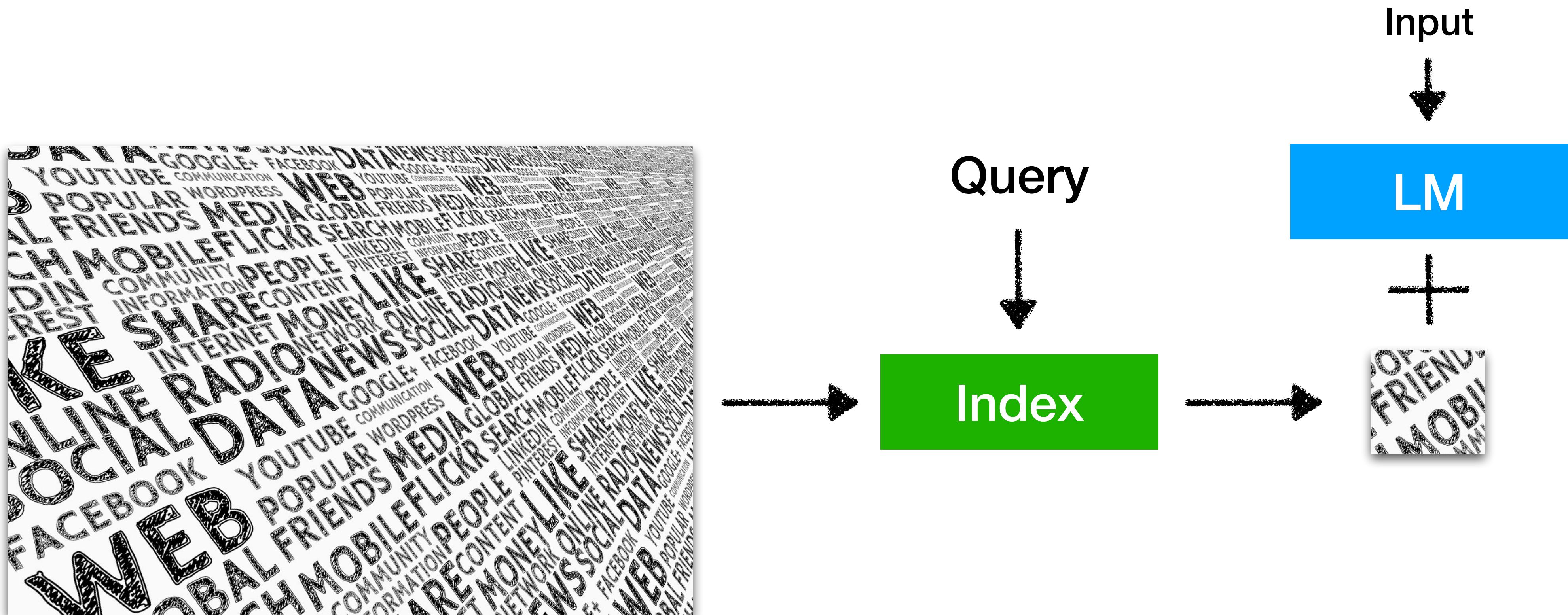
LM

Test time

Inference



Inference: Datastore



Datastore

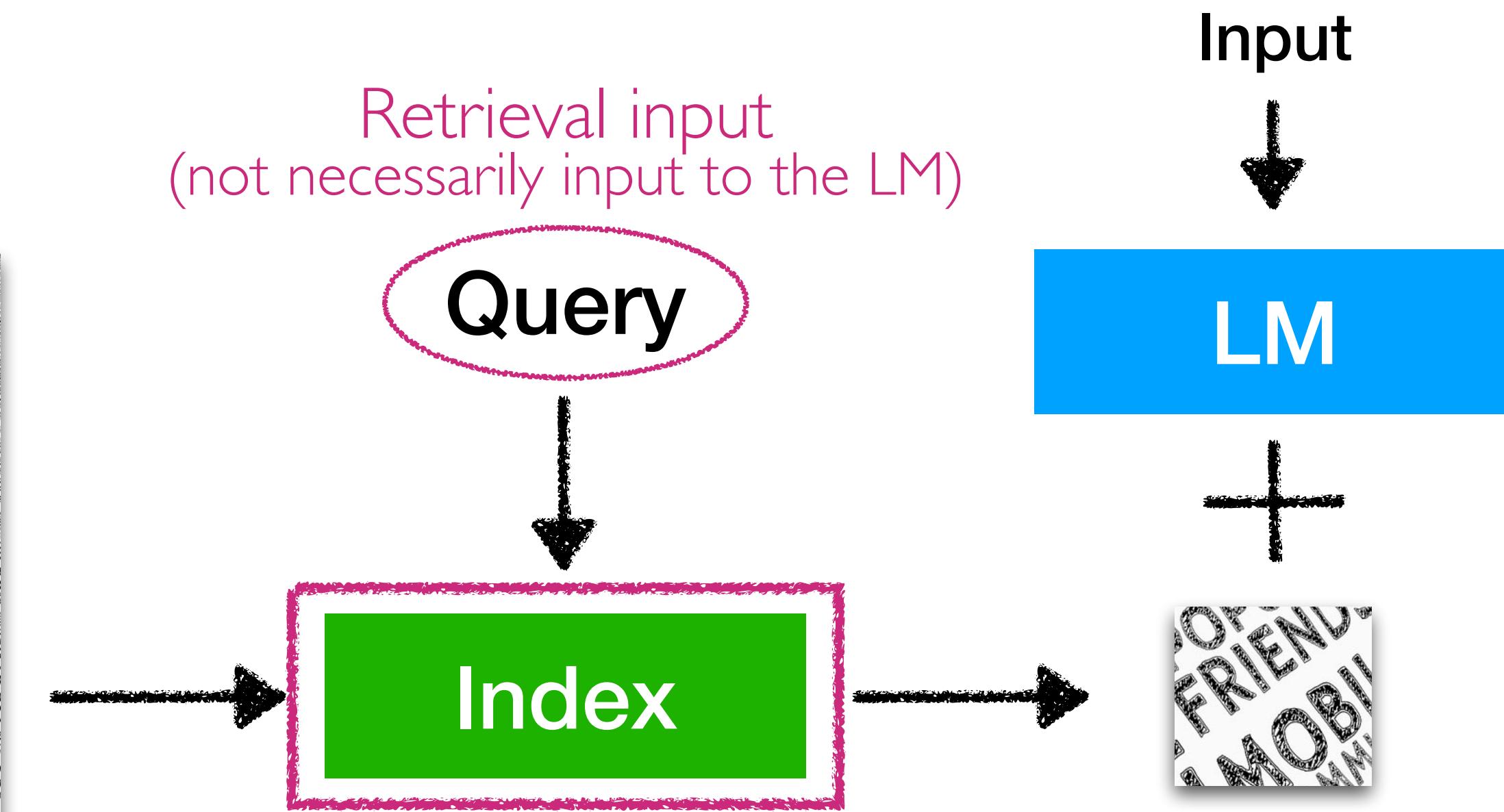
Raw text corpus

At least billions~trillions of tokens
Not labeled datasets
Not structured data (knowledge bases)
(We'll talk about other cases later)

Inference: Index



Datastore



Find a small subset of elements in a datastore
that are the most similar to the query

Inference: Index

Goal: find a small subset of elements in a datastore
that are the most similar to the query

sim: a similarity score between two pieces of text

Example

$$\text{sim}(i, j) = \frac{\text{tf}_{i,j} \times \log \frac{N}{\text{df}_i}}{\# \text{ of occurrences of } i \text{ in } j}$$

of total docs
of docs containing i

Example

$$\text{sim}(i, j) = \underline{\text{Encoder}(i)} \cdot \underline{\text{Encoder}(j)}$$

Maps the text into an h -dimensional vector

An entire field of study on how to get (or learn) the similarity function better
(We'll see some in Section 4)

Inference: Index

Goal: find a small subset of elements in a datastore
that are the most similar to the query

sim: a similarity score between two pieces of text

Can be a totally separate research area on
how to do this fast & accurate

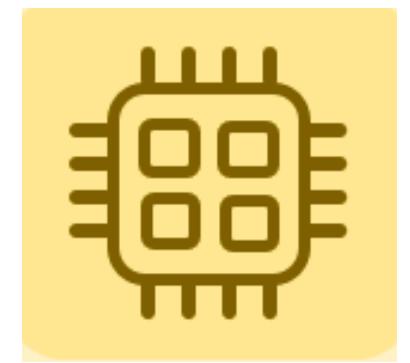
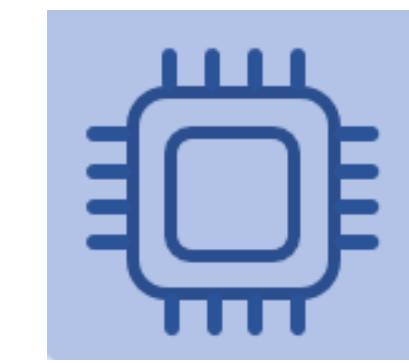
Index: given q , return $\arg\text{Topk}_{d \in \mathcal{D}} \text{sim}(q, d)$ through fast nearest neighbor search

A subset of a datastore

Software: FAISS, Distributed FAISS, SCaNN, etc...

Method	Class name	<code>index_factory</code>	Main parameters	Bytes/vector	Exhaustive	Comments
Exact Search for L2	<code>IndexFlatL2</code>	"Flat"	<code>d</code>	<code>4*d</code>	yes	brute-force
Exact Search for Inner Product	<code>IndexFlatIP</code>	"Flat"	<code>d</code>	<code>4*d</code>	yes	also for cosine (normalize vectors beforehand)
Hierarchical Navigable Small World graph exploration	<code>IndexHNSWFlat</code>	"HNSW,Flat"	<code>d, M</code>	<code>4*d + x * M * 2 * 4</code>	no	
Inverted file with exact post-verification	<code>IndexIVFFlat</code>	"IVFx,Flat"	<code>quantizer, d, nlists, metric</code>	<code>4*d + 8</code>	no	Takes another index to assign vectors to inverted lists. The 8 additional bytes are the vector id that needs to be stored.
Locality-Sensitive Hashing (binary flat index)	<code>IndexLSH</code>	-	<code>d, nbits</code>	<code>ceil(nbites/8)</code>	yes	optimized by using random rotation instead of random projections
Scalar quantizer (SQ) in flat mode	<code>IndexScalarQuantizer</code>	"SQ8"	<code>d</code>	<code>d</code>	yes	4 and 6 bits per component are also implemented.
Product quantizer (PQ) in flat mode	<code>IndexPQ</code>	"PQx", "PQM"x"nbits"	<code>d, M, nbits</code>	<code>ceil(M * nbites / 8)</code>	yes	
IVF and scalar quantizer	<code>IndexIVFScalarQuantizer</code>	"IVFx,SQ4" "IVFx,SQ8"	<code>quantizer, d, nlists, qtype</code>	<code>SQfp16: 2 * d + 8, SQ8: d + 8 or SQ4: d/2 + 8</code>	no	Same as the <code>IndexScalarQuantizer</code>
IVFADC (coarse quantizer+PQ on residuals)	<code>IndexIVFPQ</code>	"IVFx,PQ"y"x"nbits"	<code>quantizer, d, nlists, M, nbits</code>	<code>ceil(M * nbites/8)+8</code>	no	
IVFADC+R (same as IVFADC with re-ranking based on codes)	<code>IndexIVFPQR</code>	"IVFx,PQy+z"	<code>quantizer, d, nlists, M, nbits, M_refine, nbits_refine</code>	<code>M+M_refine+8</code>	no	

Exact Search

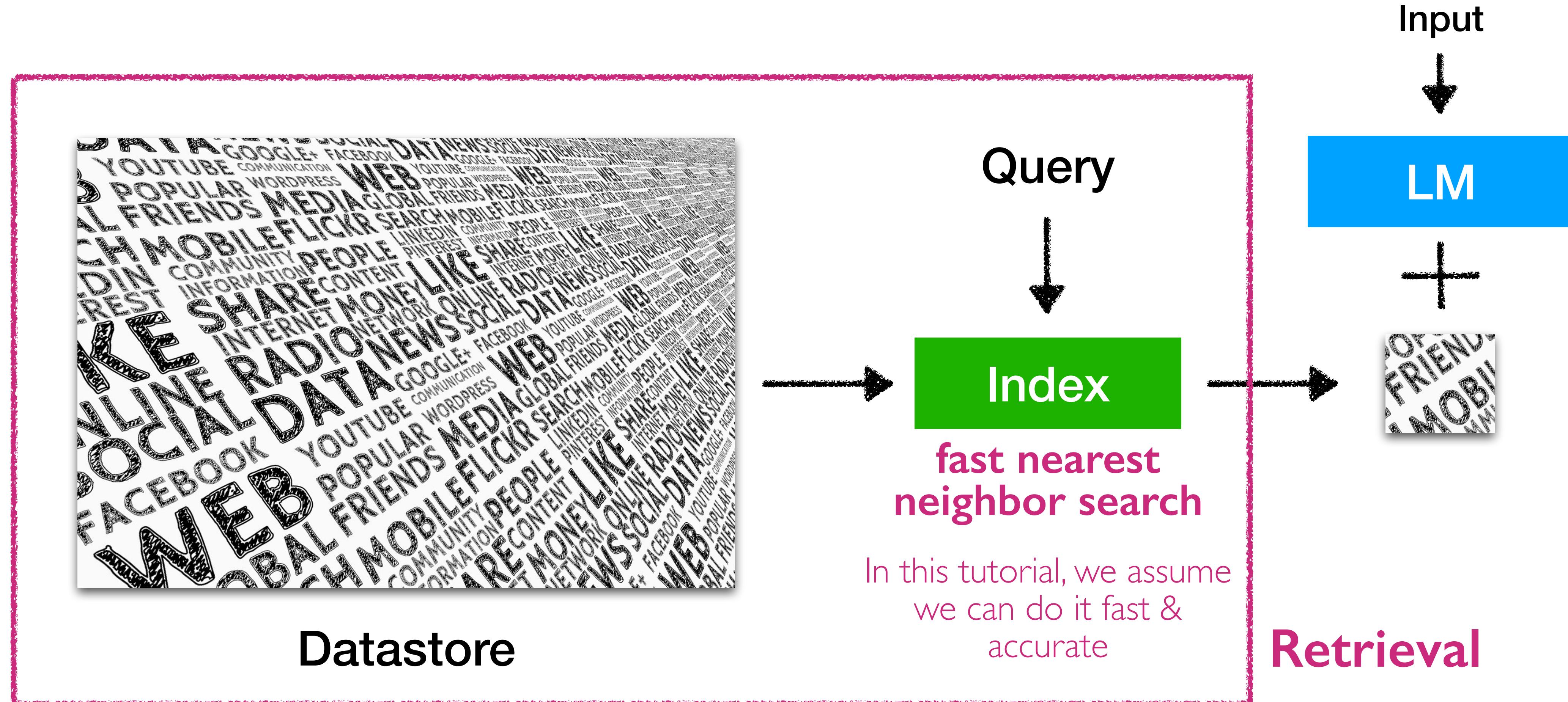


CPU vs. GPU

Approximate Search
(Relatively easy to scale to ~1B elements)

More info: <https://github.com/facebookresearch/faiss/wiki>

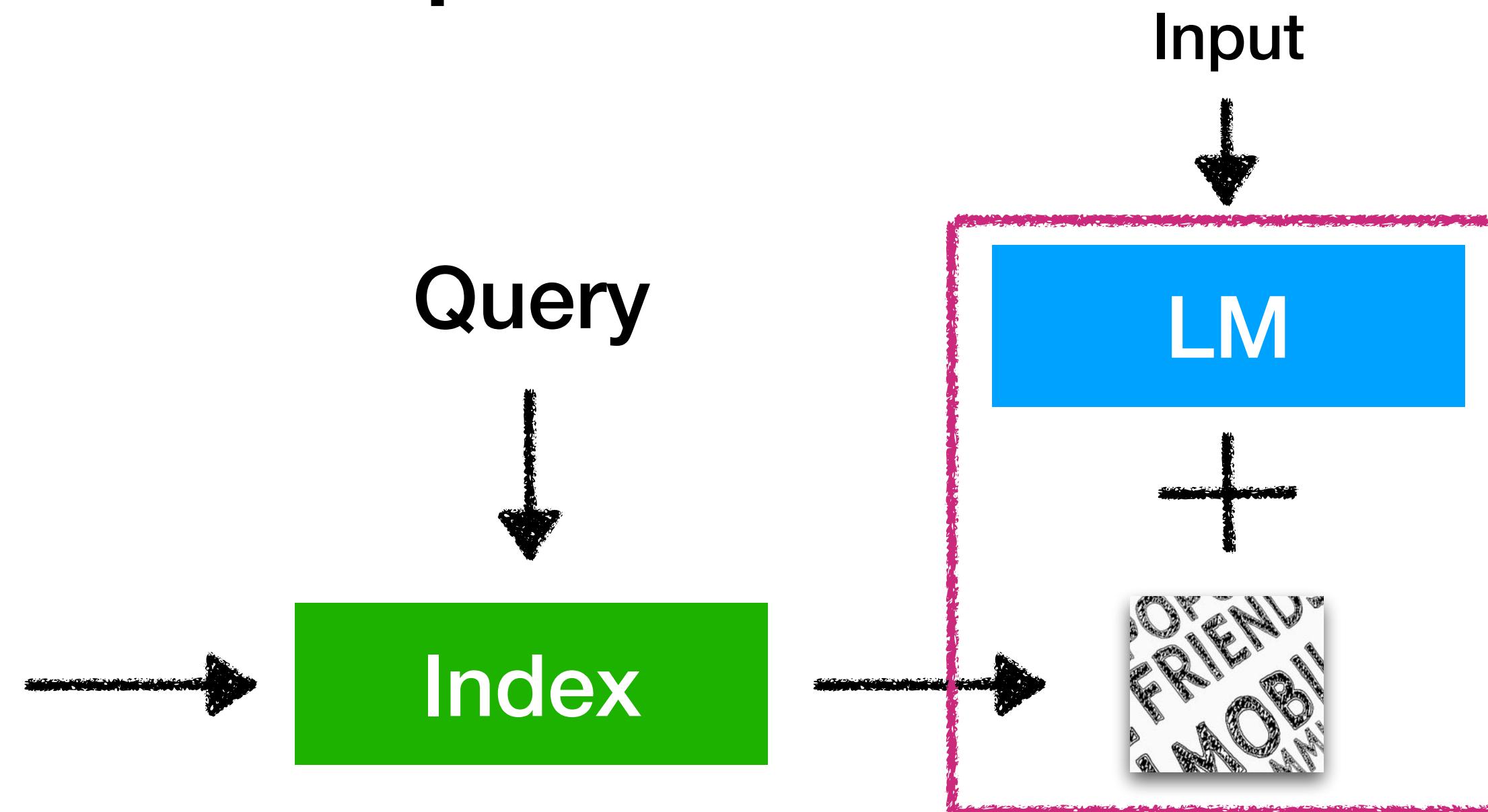
Inference: Search



Inference: Incorporation



Datastore

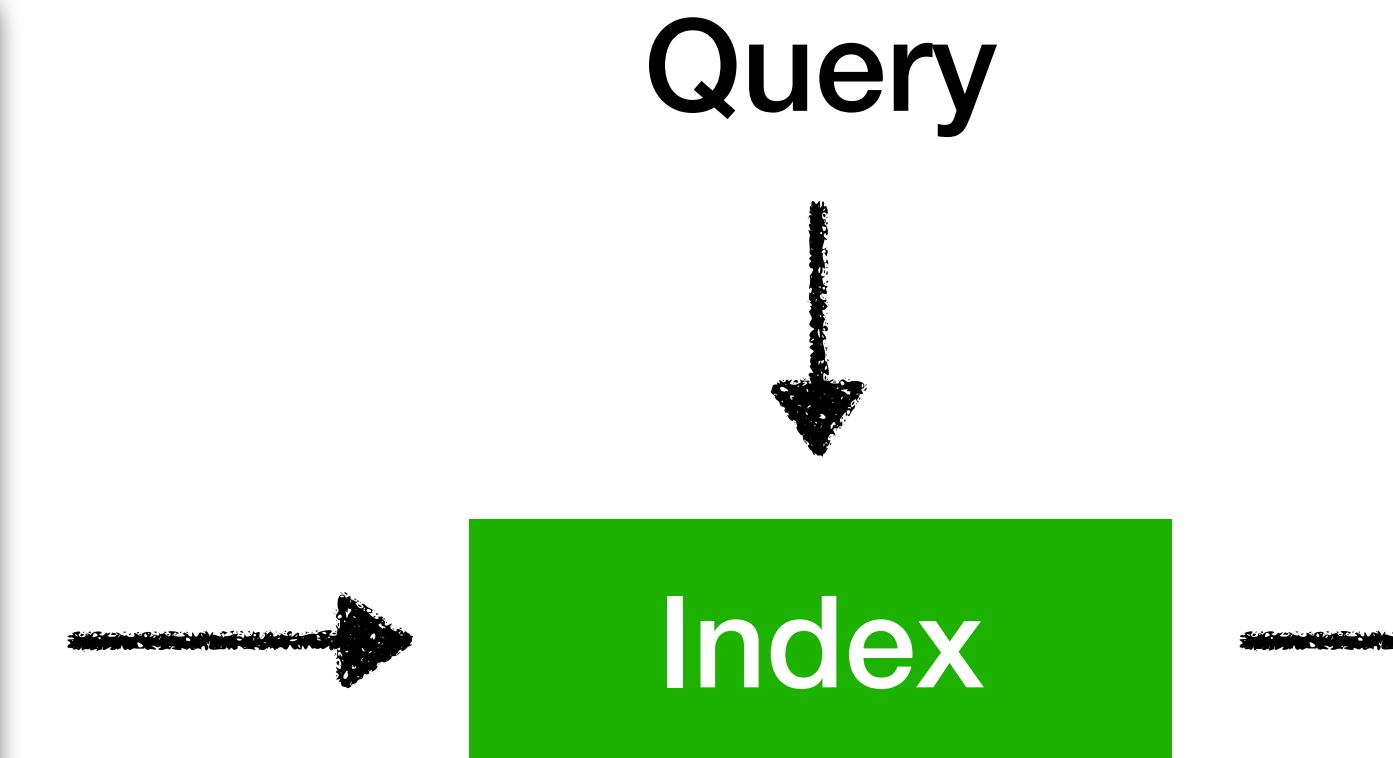


Questions to answer



Datastore

What's the query &
when do we retrieve?



Input

LM

+
+ FRIENDS
+ MOBILE

How do we
use retrieval?

What do we
retrieve?

We'll answer these questions in Section 3!

Notations



Datastore

9

