# BMI: Bottleneck-Minimal Indexing for Generative Document Retrieval

Xin Du* and Kumiko Tanaka-Ishii
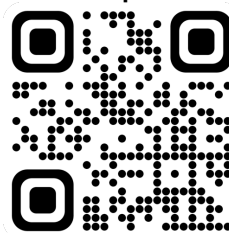
Waseda University
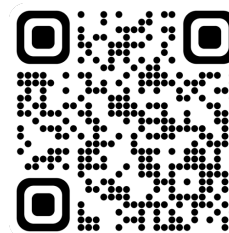
https://ml-waseda.jp

Lixin Xiu*
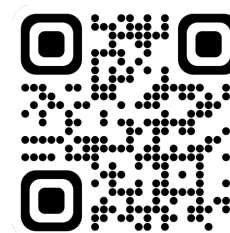
The University of Tokyo

Paper      GitHub      Our lab

* Equal contribution
Corresponding author

# Document Retrieval

**Documents** ➡️ **Index** ⬅️ **Query**

Stage 1. indexing          Stage 2. retrieval

Wikipedia pages

*"lemonade ..."*

*"coffee ..."*

*"cocktail ..."*

*"Calories in lemonade"*
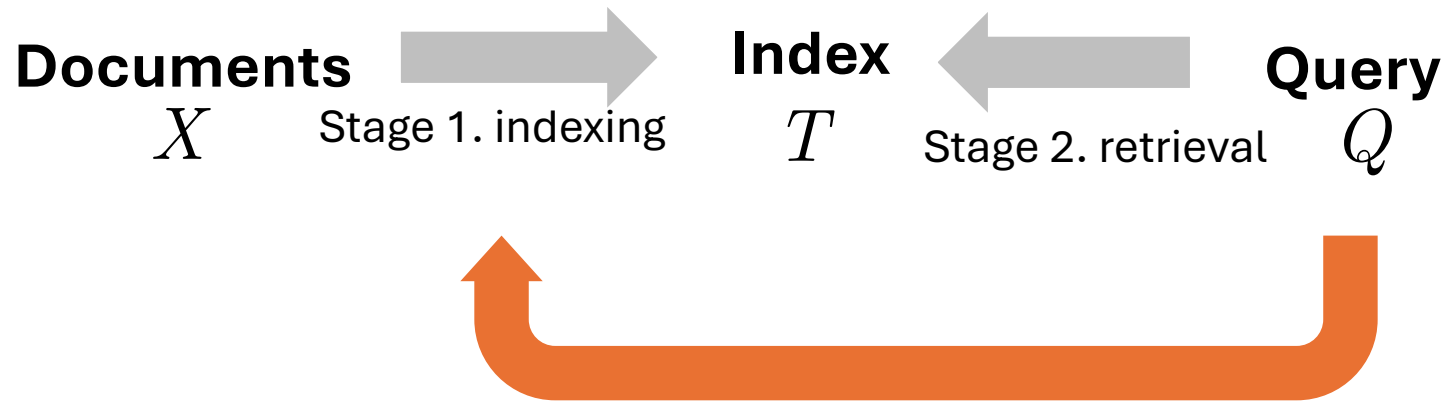
- Boolean vectors
- TF-IDF, BM25, ...
- Semantic hashing (2009)
- BERT embedding (2018)

# **Indexing** for Document Retrieval

**Documents** $\longrightarrow$ **Index** $\longleftarrow$ **Query**
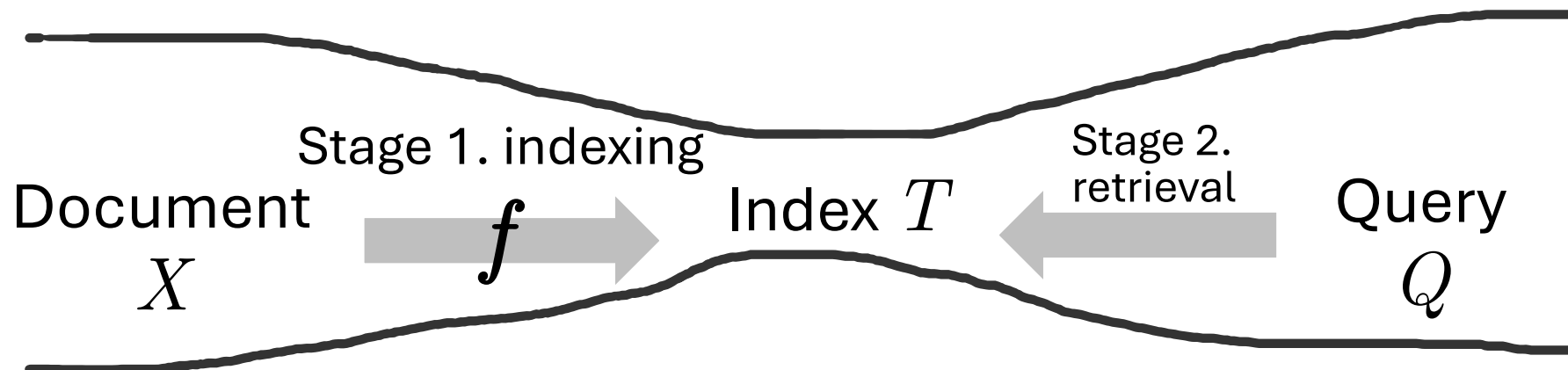$X$     Stage 1. indexing     $T$     Stage 2. retrieval     $Q$

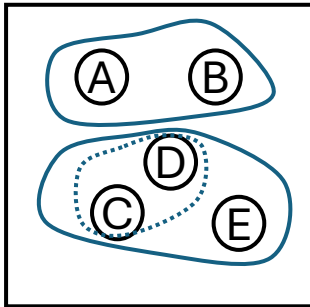A long-existing idea:     **Queries $Q$ should be incorporated** in indexing

**But how?**
It lacks a theoretical guide.

# This Work: An Information Bottleneck Model [1] for Generative Document Retrieval

(GDR explained in the next page)

Stage 1. indexing

Document $X$ → $f$ → Index $T$

Stage 2. retrieval ← Query $Q$

- Retrieval as "data transmission" from query $Q$ to document $X$.
- Index $T$ as the *"bottleneck"*.

## Optimal indexing $f^*$: $X \mapsto T$ should be ***bottleneck-minimal.***

[1] Tishby, Pereira, and Bialek. The information bottleneck method (2000)

# Generative Document Retrieval (GDR)[2,3]
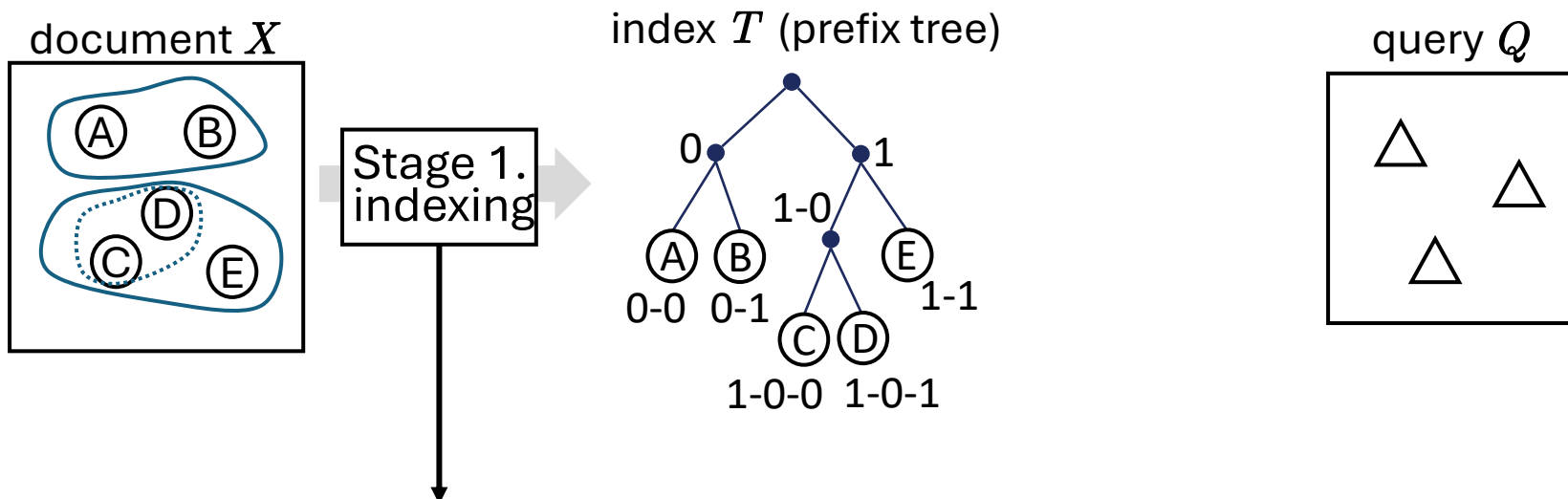
document $X$

query $Q$

[2] Cao et al. Auto-regressive entity retrieval. ICLR 2021.
[3] Tay et al. Transformer memory as a differentiable search index. NeurIPS 2022.
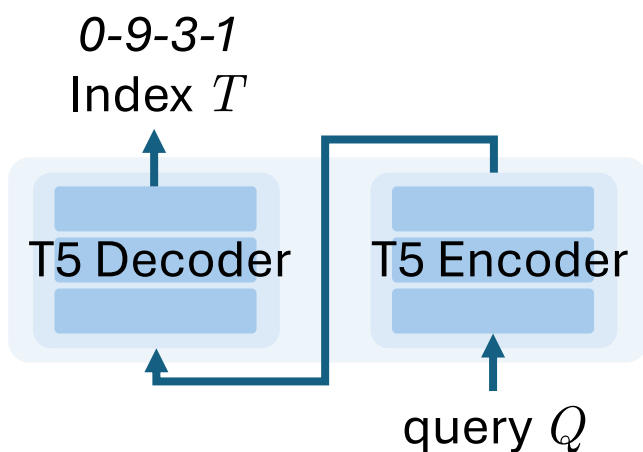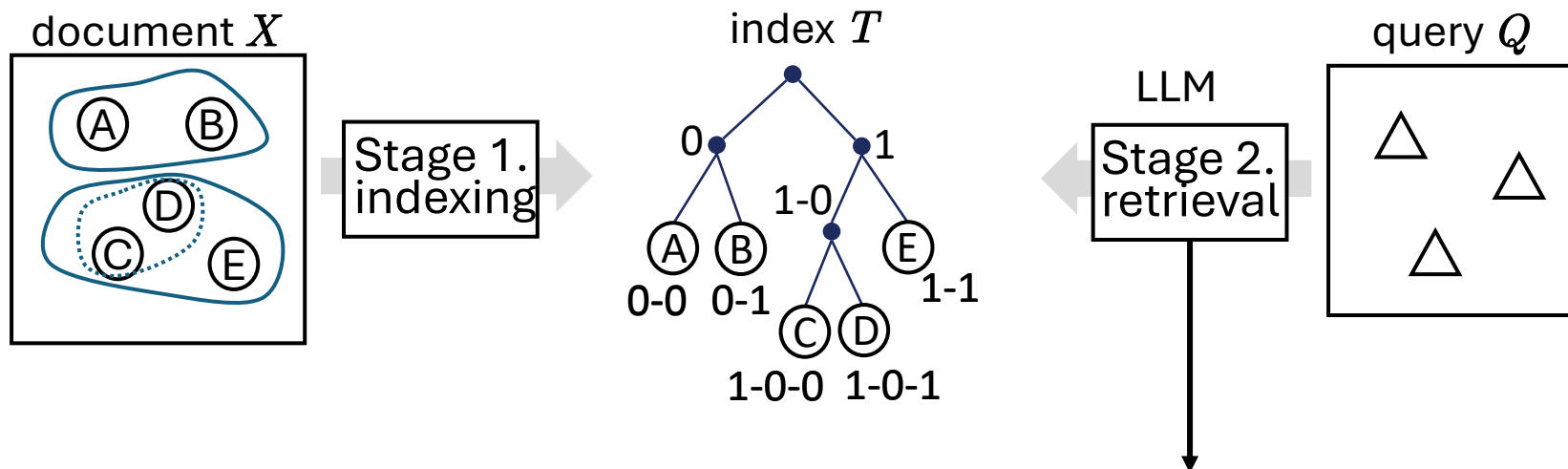
# Generative Document Retrieval (GDR)

document $X$



Stage 1. indexing

index $T$ (prefix tree)



query $Q$



## Stage 1.
Hierarchical clustering applied to document vectors. [2-3]

• Good Intuition but lacks theoretical support.

# Generative Document Retrieval (GDR)



document $X$

index $T$

query $Q$

Stage 1. indexing

LLM

Stage 2. retrieval

0

1

1-0

A   B                    E

0-0   0-1                1-1

C   D

1-0-0   1-0-1

*0-9-3-1*
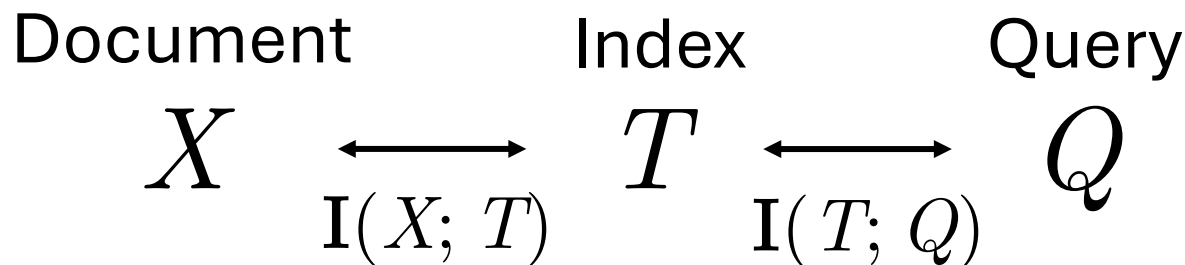Index $T$

T5 Decoder      T5 Encoder

query $Q$

*"Calories in lemonade."*

Stage 2. Train a language model to "translate" a query into index.
(i.e., underline{generate} the index digit by digit)

Promising, but limited by the model's *finite* and often *insufficient* size.

# Information Bottleneck (IB) Model for GDR

Document      Index      Query

$$X \quad \longleftrightarrow \quad T \quad \longleftrightarrow \quad Q$$

$$\mathbf{I}(X;\,T) \qquad\qquad \mathbf{I}(T;\,Q)$$

IB studies the *tradeoff* by the Lagrangian:

$$L[p(T|X),\,\beta] = \mathbf{I}(X;\,T)\ -\ \beta\,\mathbf{I}(T;\,Q) \quad \beta \geq 0$$

assuming Markov chain $T \leftrightarrow X \leftrightarrow Q$

## Why is this *tradeoff* essential for GDR ?

# Why is this *tradeoff* essential for GDR ?
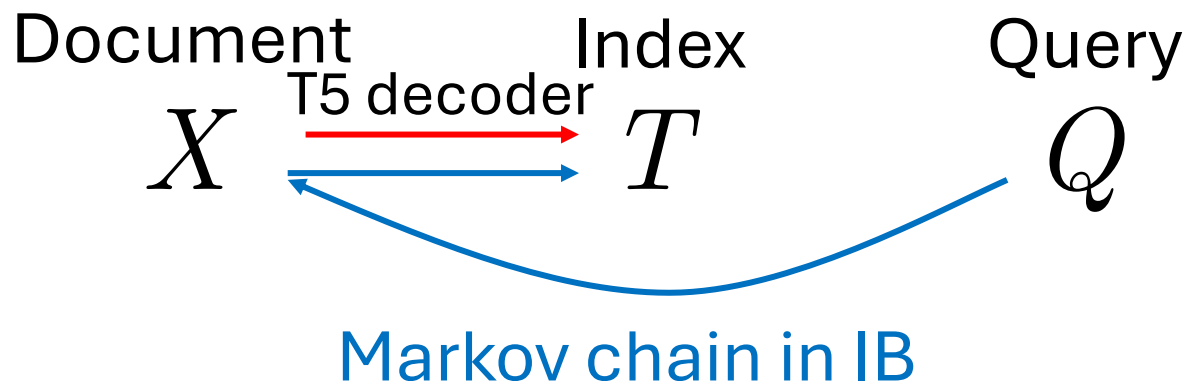
Document

$X$

Index

$T$

Query

$\longleftarrow Q$

Retrieve by
generation

**Performance**

$$\mathbf{I}(T; Q)$$

Maximal
retrieval accuracy

# Why is this *tradeoff* essential for GDR ?

Document $\qquad$ Index $\qquad$ Query

$$X \xrightarrow{\text{T5 decoder}} T \qquad Q$$

Markov chain in IB

**Cost of Memory**

$$\mathbf{I}(X; T)$$

How many bits the language model must *"memorize"* per document

**Performance**

$$\mathbf{I}(T; Q)$$

Maximal
retrieval accuracy

# Tradeoff: Model Size v.s. Retrieval Performance

**minimize**

$$L[p(T|X), \beta] = \mathbf{I}(X; T) - \beta \mathbf{I}(T; Q)$$

**Model size**    **Performance**

# Tradeoff: Model Size v.s. Retrieval Performance

**minimize**

$$L[p(T|X), \beta] = \mathbf{I}(X; T) - \beta \, \mathbf{I}(T; Q)$$

**Model size (T5)** **Performance**

Optimal indexing achieves:

 highest accuracy for a specific model size, or

# Tradeoff: Model Size v.s. Retrieval Performance

**minimize**

$$L[p(T|X), \beta] = \mathbf{I}(X; T) \ - \ \beta \, \mathbf{I}(T; Q)$$

**Model size (T5)** **Performance**

Optimal indexing achieves:

smallest model size for a given accuracy

# Tradeoff: Model Size v.s. Retrieval Performance

**minimize**

$$L[p(T|X), \beta] = \mathbf{I}(X; T) - \beta \, \mathbf{I}(T; Q)$$

**Model size (T5)** **Performance**

Such optimal indexing under a **limited model-size budgets** is called the ***bottleneck-minimal indexing*** (next page),

*NOT* the case when a model can be infinitely-large (**unlimited model size**), e.g., a maximum inner-product search model.

# This Work: Bottleneck-Minimal Indexing

**Our proposed definition of BMI:**
An indexing function $f: X \mapsto T$ is called an BMI if it maximizes the likelihood function $p(\text{dataset}| f)$

$$f^* = \underset{f}{\text{argmax}} \prod_{\text{doc } x} p^*(X{=}x \mid T{=}f(x))$$

$$= \underset{f}{\text{argmax}} \prod_{\text{doc } x} \underbrace{\frac{p^*(T{=}f(x)|X{=}x)}{p^*(f(x))}}_{\text{Optimal Solution to IB}} \underbrace{p(x)}_{\text{constant}}$$

$$= \underset{f}{\text{argmin}} \sum_{\text{doc } x} \text{KL}[p(Q|x) \parallel p(Q|f(x))]$$

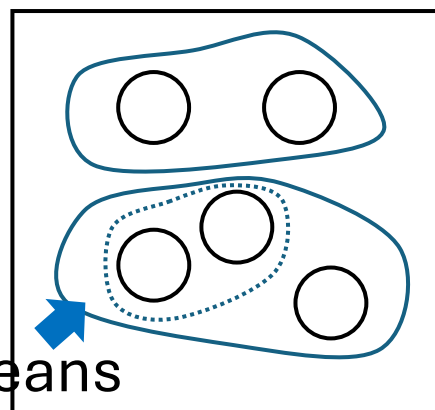# How we acquire $f^*$

Assume $p(Q|x)$ and $p(Q|f(x))$ to be Gaussian

center vector of Gaussian

$$f^* = \underset{f}{\mathrm{argmin}} \sum_{\mathrm{doc}\ x} \| \mathrm{E}[p(Q|x)] - \mathrm{E}[Q|f(x)] \|^2$$

Estimated with a doc2query[4] model

Essentially, we applying k-means clustering to query center vectors $\{\mathrm{E}[p(Q|x)]\}$ instead of document vectors $\{x\}$.
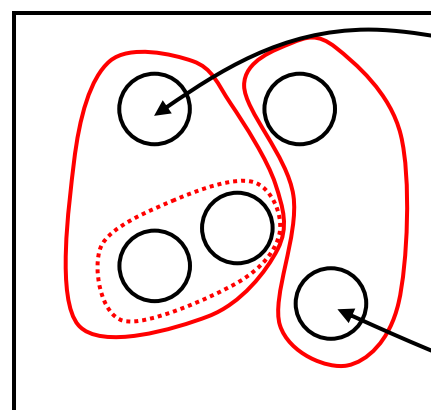
[4] Nogueira and Jimmy Lin. From doc2query to docTTTTTquery (2019).

# How we acquire $f^*$

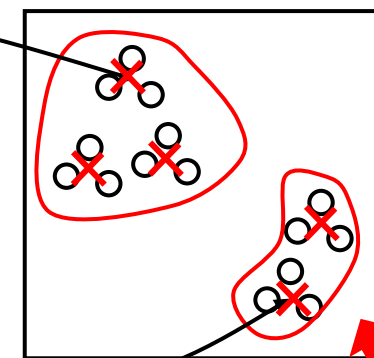Previous method [2-3]

Our method

Document $X$

Document $X$
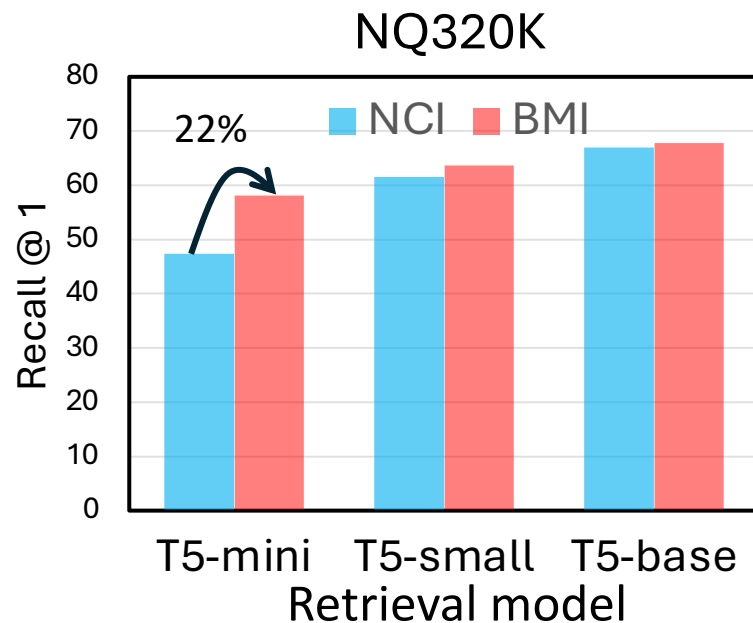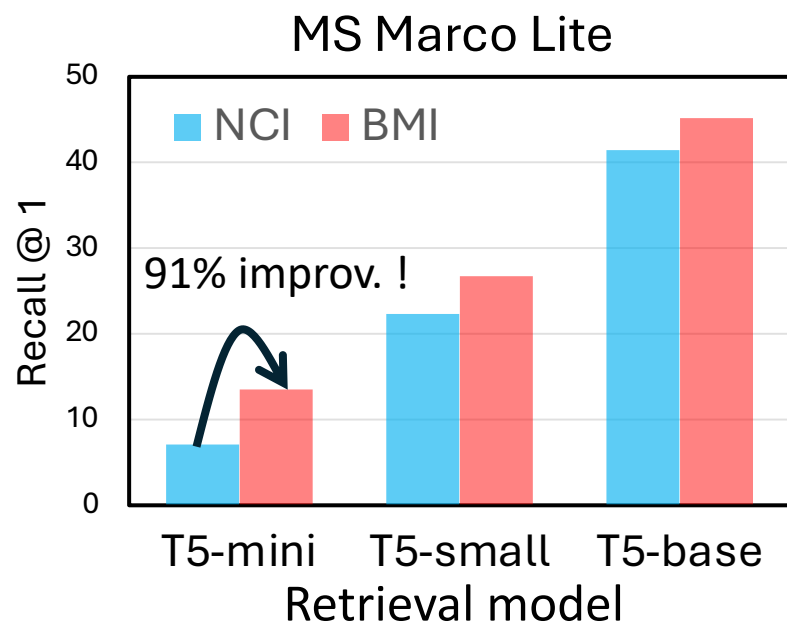
Query $Q$



K-means

Doc2query [4]

K-means

Essentially, we applying k-means clustering to query center vectors $\{\mathrm{E}[p(Q|x)]\}$ instead of document vectors $\{x\}$.

[2] Cao et al. Auto-regressive entity retrieval. ICLR 2021.
[3] Tay et al. Transformer memory as a differentiable search index. NeurIPS 2022.
[4] Nogueira and Jimmy Lin. From doc2query to docTTTTTquery (2019).

# Retrieval Accuracy



BMI organize indexing "knowledge" in a much more efficient way than previous GDR methods [2-3] when model size is insufficient, which is common in real-world applications.
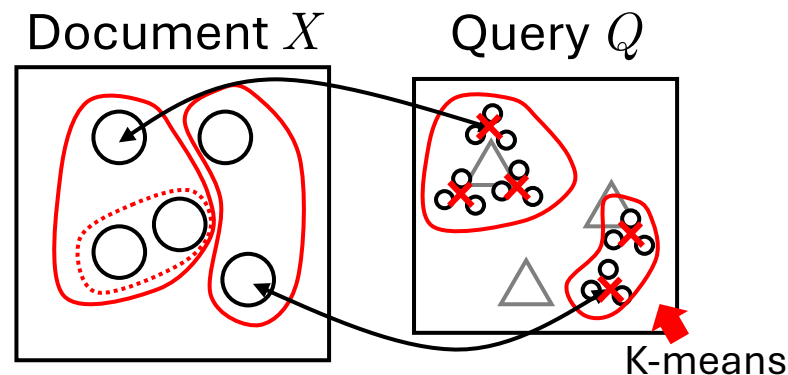
# Takeaway

Bottleneck-minimal indexing: an optimal indexing principle under a limited model-size budget.

$$L[p(T|X), \beta] = \mathbf{I}(X; T) - \beta\, \mathbf{I}(T; Q)$$
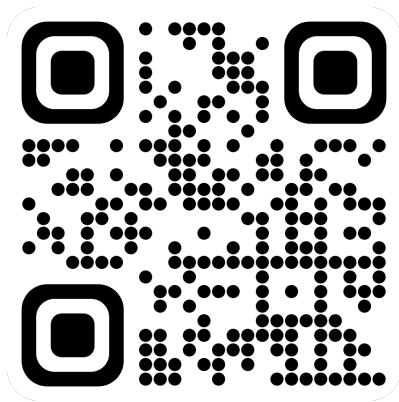
**Model size**   **Performance**

In GDR, better to watch queries rather than documents.
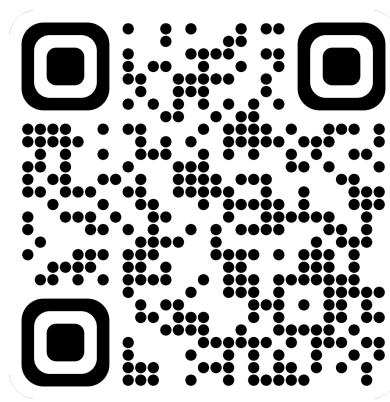


Document $X$     Query $Q$

K-means

Information-bottleneck theory worked well to explain how knowledge like indexing can be organized efficiently.
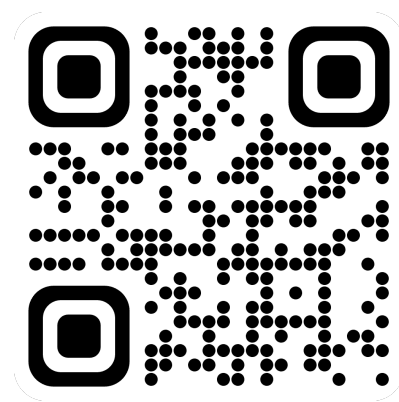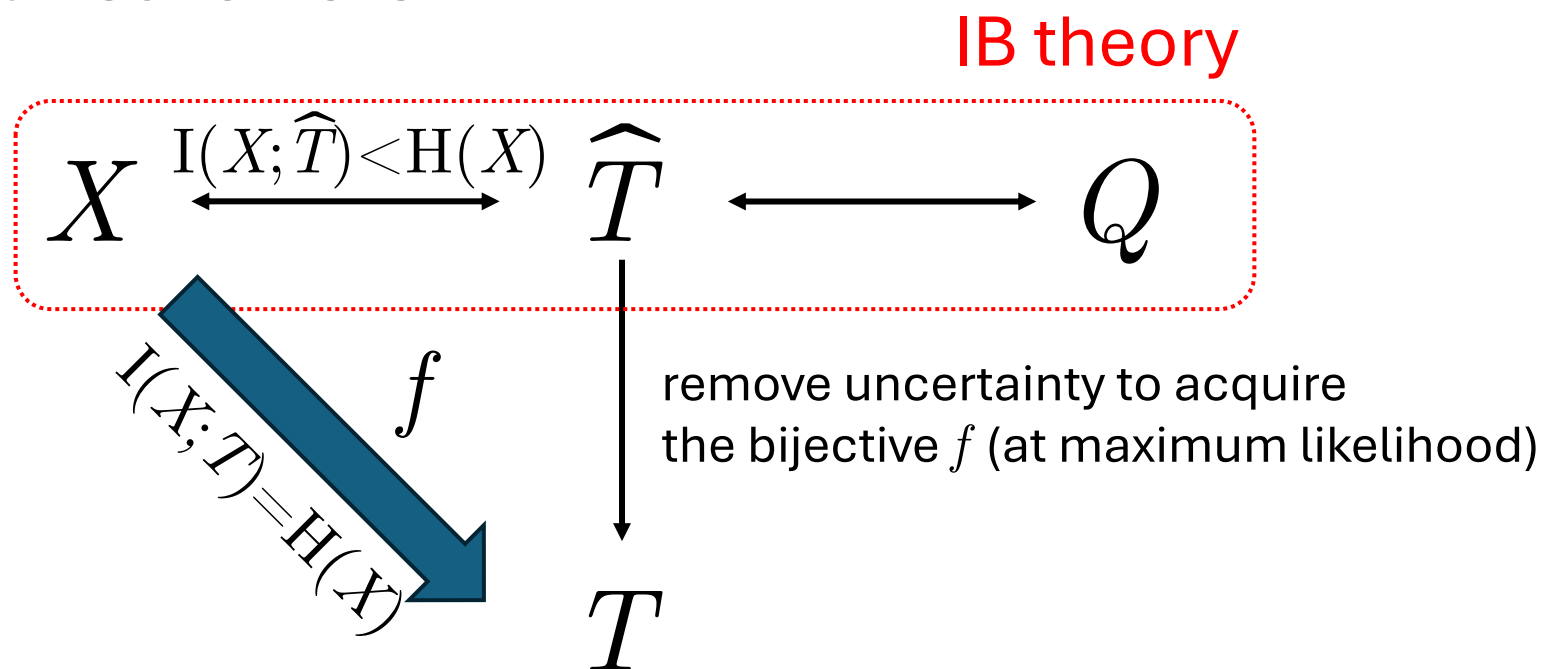
# Thank you for listening!

Paper

GitHub

Our lab

# As $f$ is bijective, isn't $\mathrm{I}(X;T)$ constant ?

We reused $T$ for two different variables $T$ and $\widehat{T}$ for simpler presentation.
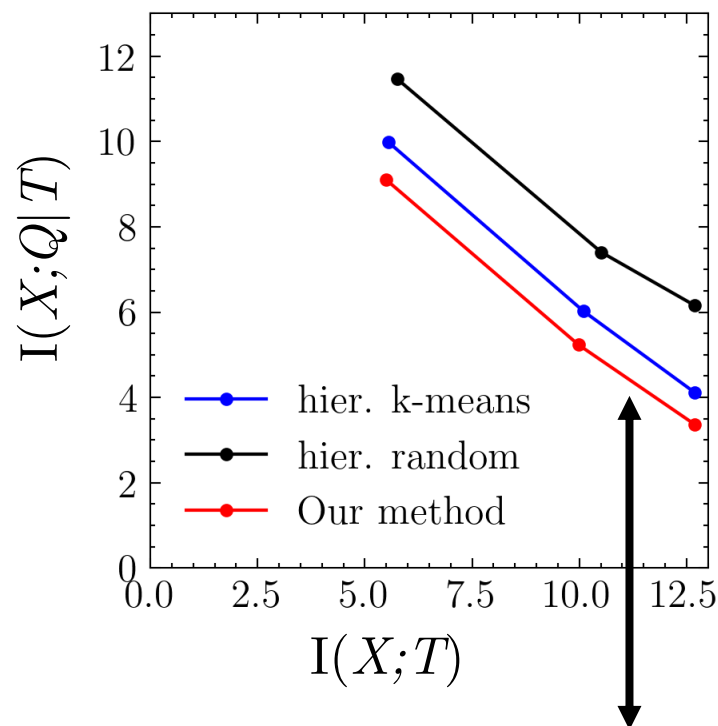
The full scheme is:



IB theory

$X \xleftrightarrow{\mathrm{I}(X;\widehat{T})<\mathrm{H}(X)} \widehat{T} \longleftrightarrow Q$

$\mathrm{I}(X;T)=\mathrm{H}(X)$

$f$

remove uncertainty to acquire the bijective $f$ (at maximum likelihood)

$T$

# Theoretical & Empirical IB Curves

Theoretical IB curve

$\mathrm{I}(X;Q|T)$
$= \mathrm{const} - \mathrm{I}(T;Q)$

theoretical IB
curve (Section 4.1)

empirical-dist.
approxim. (Section 5.1)

neural-network fits
(different model size)

Unfeasible
region

$0$            $\mathrm{I}(X;T)$

Empirical IB curve



IB theory: balance between
$\mathrm{I}(X;T)$ and $\mathrm{I}(T;Q)$

Our method is closer to
the theoretical IB curve.

# Datasets & Settings

|  | NQ320K | MS Marco Lite |
|---|---|---|
| # Documents | 109,739 | 138,457 |
| # Queries (train) | 307,373 | 183,947 |
| # Queries (test | 7,830 | 2,792 |

Used docT5query (base) [5] for query generation.

Retrieval model training: used the implementation of NCI [3]

The change from NCI [3] is only the indexing $f\colon X \mapsto T$