



# A Dense Subset Index for Collective Query Coverage

Kartik Nair<sup>1,2</sup>, Prithish Chakraborty<sup>1</sup>, Atharva Tambat<sup>1</sup>, Indradyumna Roy<sup>1</sup>, Soumen Chakrabarti<sup>1</sup>, Anirban Dasgupta<sup>3</sup>, and Abir De<sup>1</sup>

<sup>1</sup>IIT Bombay <sup>2</sup>Carnegie Mellon <sup>3</sup>IIT Gandhinagar



## Background: Traditional Vector-Bag Retrieval

Consider bag-of-words scoring (MaxSim) popularized by ColBERT [Khattab et al. 2020] — given a query and corpus as sets of **token embeddings**

$$Q = \{q_1, \dots, q_M\} \quad X = \{x_1, \dots, x_L\}$$

$$F(X_c, q) = \max_{x \in X_c} q^T x$$

Traditional retrieval is *competitive*: documents are scored *against each other* to rank higher.

$$S \leftarrow \text{top-}k_{X \in C} F(X, Q)$$

## Collaborate, not compete

Collective scoring is a natural extension to MaxSim. Scoring a set of documents  $S = \{X_1, X_2, \dots, X_k\}$  against the query.

$$F(S, Q) = \sum_{q \in Q} \max_{x \in \cup_{s \in S} X_s} q^T x = \sum_{q \in Q} \max_{s \in S} \max_{x \in X_s} q^T x$$

**Coverage Objective:**

$$\max_{S \subseteq C} F(S, Q) \quad \text{s.t. } |S| \leq K$$

This objective is sub modular in S, hence admits a  $(1 - 1/e)$  greedy approximation with the marginal gain function  $F(c | S, Q)$

- 1: Initialize  $S_0 \leftarrow \emptyset$
- 2: **for**  $k = 1, \dots, K$  **do**
- 3:  $c_k = \arg \max_{c \in C \setminus S_{k-1}} F(c | S_{k-1}, Q)$
- 4:  $S_k \leftarrow S_{k-1} \cup \{c_k\}$
- 5: **return**  $S_K$

The greedy solution requires **exhaustive scoring**, which **grows linearly** in complexity with the size of corpus.

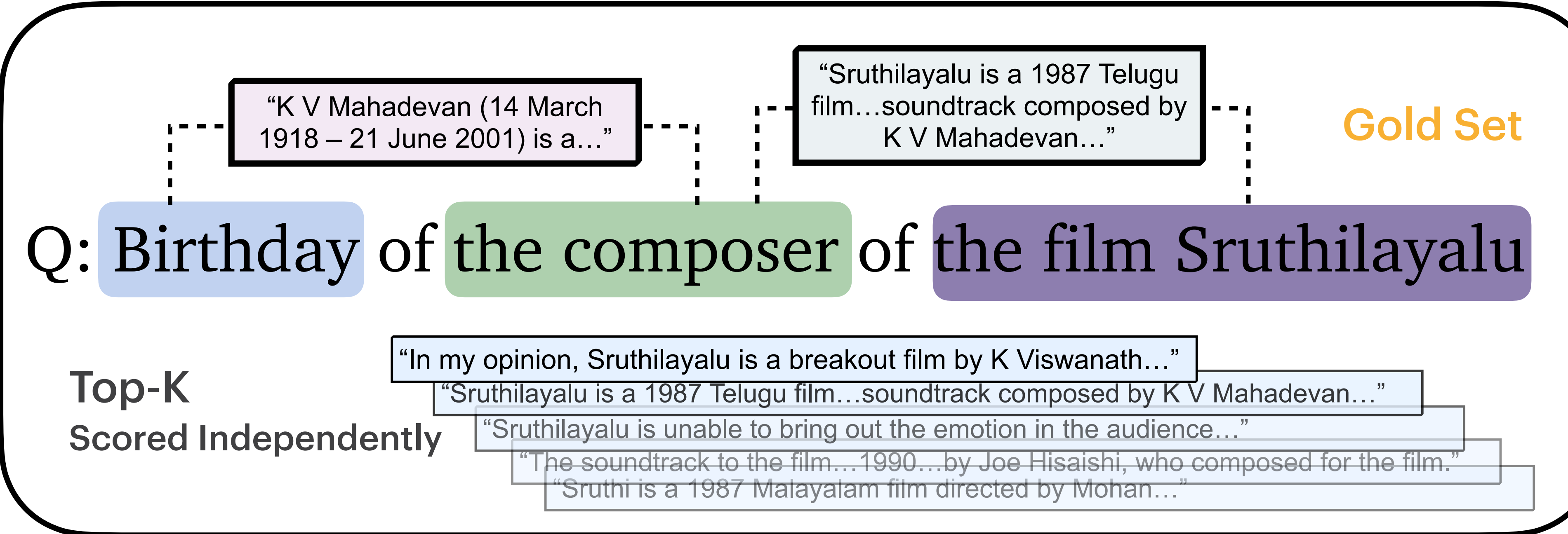
We can instead **formulate this as a first stage retrieval problem!**

**Adaptive Probing**: updated marginal gain at each step

$$S_0 = \emptyset \xrightarrow[F(\bullet | S_0, Q)]{\text{probe with}} S_1 \xrightarrow[F(\bullet | S_1, Q)]{\text{probe with}} S_2 \rightarrow \dots \rightarrow S_K$$

**Challenge:** adapt and/or approximate score for retrieval

- (1) Separate into corpus-dependent term for indexing
- (2) Convert to inner-product form



## Making the Marginal Gain Index-able

1. Difference to dot product: separating state dependence

$$F(c | S, Q) = \sum_{q \in Q} [F(X_c, q) - F(S, q)]_+ = \sum_{q \in Q} \max_{x \in X_c} \left[ \underbrace{q^T x}_{=: \hat{q}_S} \right]_+ \quad \text{state dependence absorbed by the query}$$

2. Dealing with the hinge score with sign-hashing

$$w \sim \mathcal{N}(0, \mathbf{I}_{d+1})$$

$$\text{Randomised Feature Map } \phi_w(x) = \frac{1}{\sqrt{2}} \begin{bmatrix} x \\ \text{sign}(w^T x) \end{bmatrix}$$

$$\Pr(\phi_w(x)^T \phi_w(y) = [x^T y]_+) = \frac{\text{green} + \text{red}}{\text{green} + \text{red} + \text{blue}} \geq \frac{1}{2}$$

We can improve our approximation by drawing multiple hyperplanes  $\{w_1, w_2, \dots, w_R\} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_{d+1})$

This yields a new scoring function  $G_{1:R}$

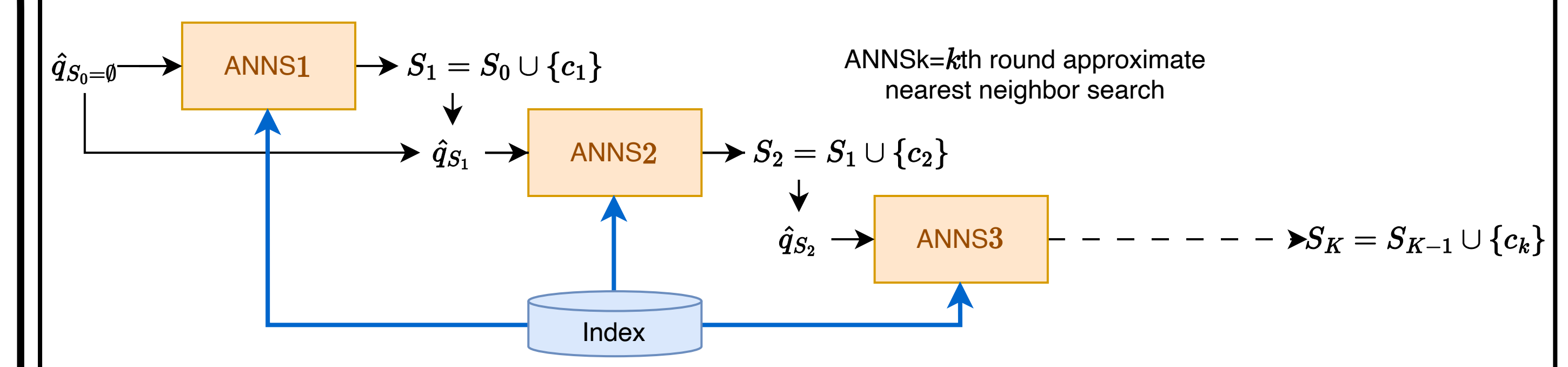
$$G_{1:R}(c | S, Q) = \sum_{q \in Q} \max_{x \in X_c} \left[ \max_{r \in [R]} \phi_{w_r}(\hat{q}_S)^T \phi_{w_r}(\hat{x}) \right]$$

Replacing the marginal with  $G_{1:R}$  gives us a  $(1 - 1/e - \delta)$ -approximation where  $\delta \leq |Q|/2^R$

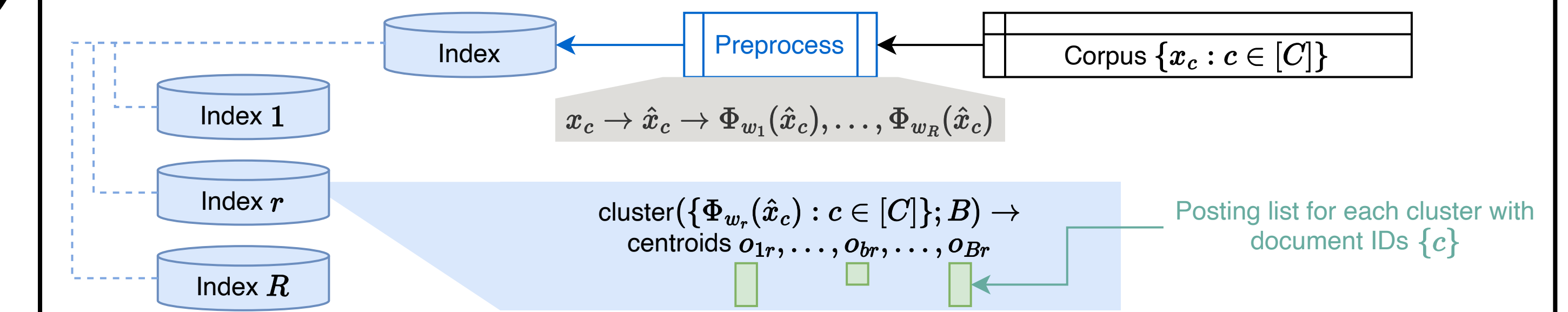
However, indexing this now requires  $R$  Replica indexes: with each replica storing the corresponding feature-mapped document tokens  $\{\phi_{w_r}(\hat{x}) | x \in X\}$

## The Retrieval Pipeline

Each greedy step is replaced by an approximate nearest neighbour search



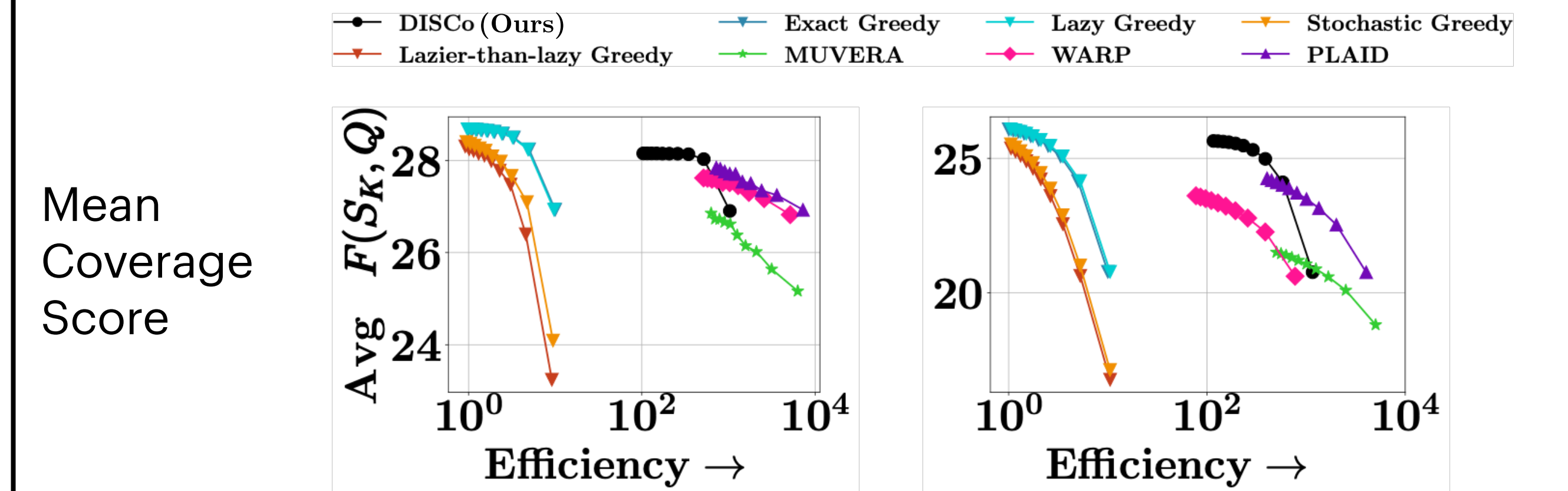
Each replica index has a token-level IVF, with each token stored as a centroid and quantized residual.



Retrieval consists of a token level lookup followed by progressive levels of candidate pruning

## Experiments

We compare our method (DISCo) against variants of greedy sub modular solvers, and traditional multi-vector retrieval engines.



Dataset ↓ Method →	DISCo (Ours)	Exact Greedy	Lazy Greedy	PLAID	WARP
2WikiMultiHopQA	0.90	0.91	0.91	0.89	0.82
HotpotQA	0.84	0.83	0.83	0.81	0.77

MAP on Multi-Hop QA