

PROBLEM

The problem of similarity queries on vector space models is a general problem compared to commonly used functions such as Jaccard and cosine. It is a difficult problem on large-scale data sets due to three aspects.

1. Arbitrary token weights integrated in pair-wise similarity.
2. Variant degrees of each token in different records.
3. Given a threshold, exactly answer a query from millions of records.

CONTRIBUTIONS

We formally define the problem of similarity queries on collections of records using a vector space model, and solve it efficiently based on a revised inverted index. Our main contributions are

1. Pruning methods with various bounds
2. Solid proofs of their correctness
3. An efficient elastic scheme to iteratively tighten the bounds
4. A family of algorithms based on an inverted index for large data sets
5. Evaluating intensive experiments on dozens million of records

FORMAL DEFINITION

Given a collection of records \mathcal{S} , a p-norm selection includes a record Q and a threshold τ . Suppose a token has a weight $w(t)$ and a degree $f(t, R)$ (in record R), it is to find all the records $R \in \mathcal{S}$ such that

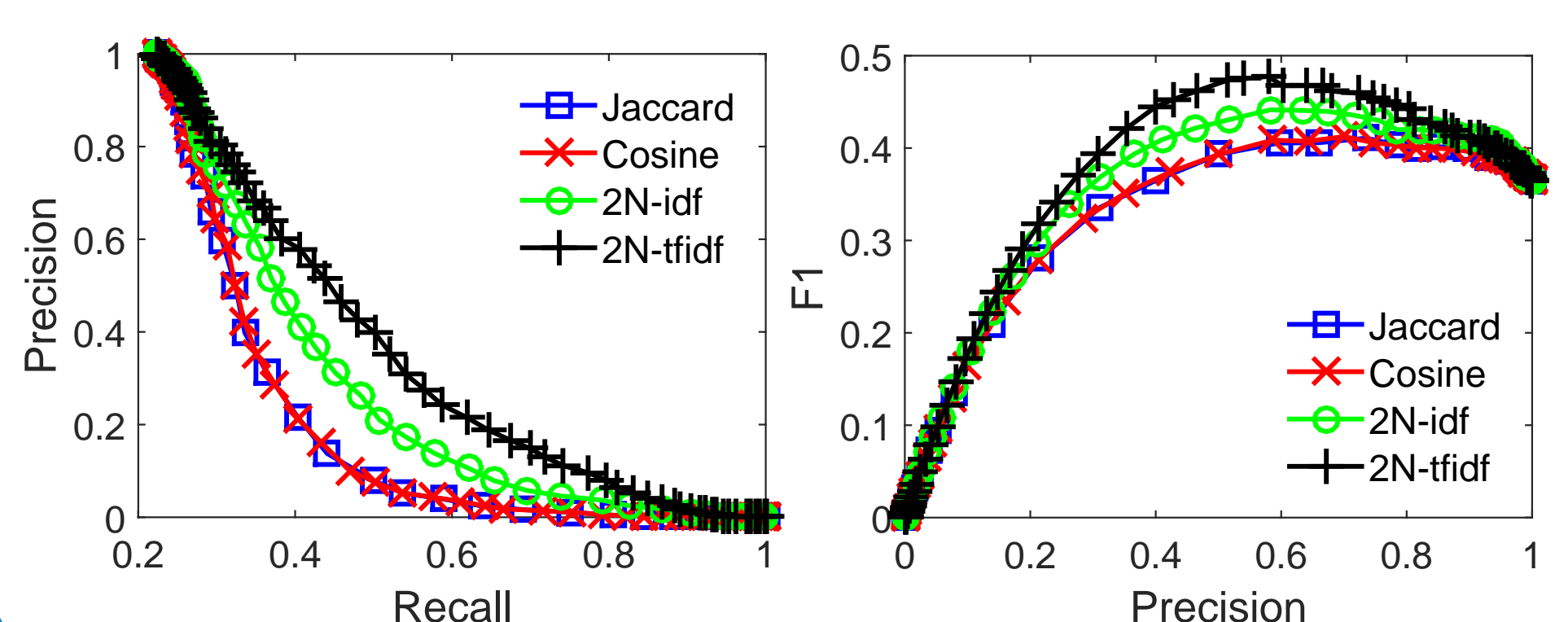
$$\frac{\left(\sum_{t \in R \cap Q} (f(t, R)w(t))^{\frac{p}{2}} (f(t, Q)w(t))^{\frac{p}{2}} \right)^{\frac{2}{p}}}{\left(\sum_{t \in R} (f(t, R)w(t))^p \sum_{t \in Q} (f(t, Q)w(t))^p \right)^{\frac{1}{p}}} \geq \tau$$

USABILITY

Based on the ground truth about similar pairs of records on "AskUbuntu", the two figures compare the precision and recall of this search using different functions. The results given above say that:

- The tf-idf 2-norm function gave the best accuracy results.
- The two 2-norm weighting functions were better than Jaccard and Cosine.

It was not to show that the p-norm function is the best for all applications. Instead, but a better similarity in certain applications.



THEORETIC BACKGROUNDS

R	t ₁	t ₃	t ₅	t ₆	t ₇	t ₈	t ₉	t ₁₀	Q	t ₂	t ₄	t ₆	t ₇	t ₁₀
c	0.8	0.5	0.6	1.6	0.8	0.4	0.1	0.2	c	0.2	1.2	0.8	1.6	1.2
c ²	0.64	0.25	0.36	2.56	0.64	0.16	0.01	0.04	c ²	0.04	1.44	0.64	2.56	1.44
R	[Visual representation of record R with token weights]								Q	[Visual representation of query Q with token weights]				
i	1	2	3	4	5	6	7	8	j	1	2	3	4	5

⊙	t ₁	t ₂	t ₃	t ₄	t ₅	t ₆	t ₇	t ₈	t ₉	t ₁₀
w(t)	0.1	0.2	0.1	0.4	0.5	0.8	0.8	0.4	0.1	0.2
f(t,R)	8	0	5	0	3	2	1	1	1	1
f(t,Q)	0	1	0	3	0	1	2	0	0	6

Processed and remaining tokens of records.

We first define the processed and remaining tokens in a record when we consider all of its tokens based on a total order.

- The bound on remaining tokens is proposed to reduce the candidate size.
- The bounds on length and prefix can be used to save sequential IOs.
- Using the concepts in the figure we give solid proofs for all bounds.

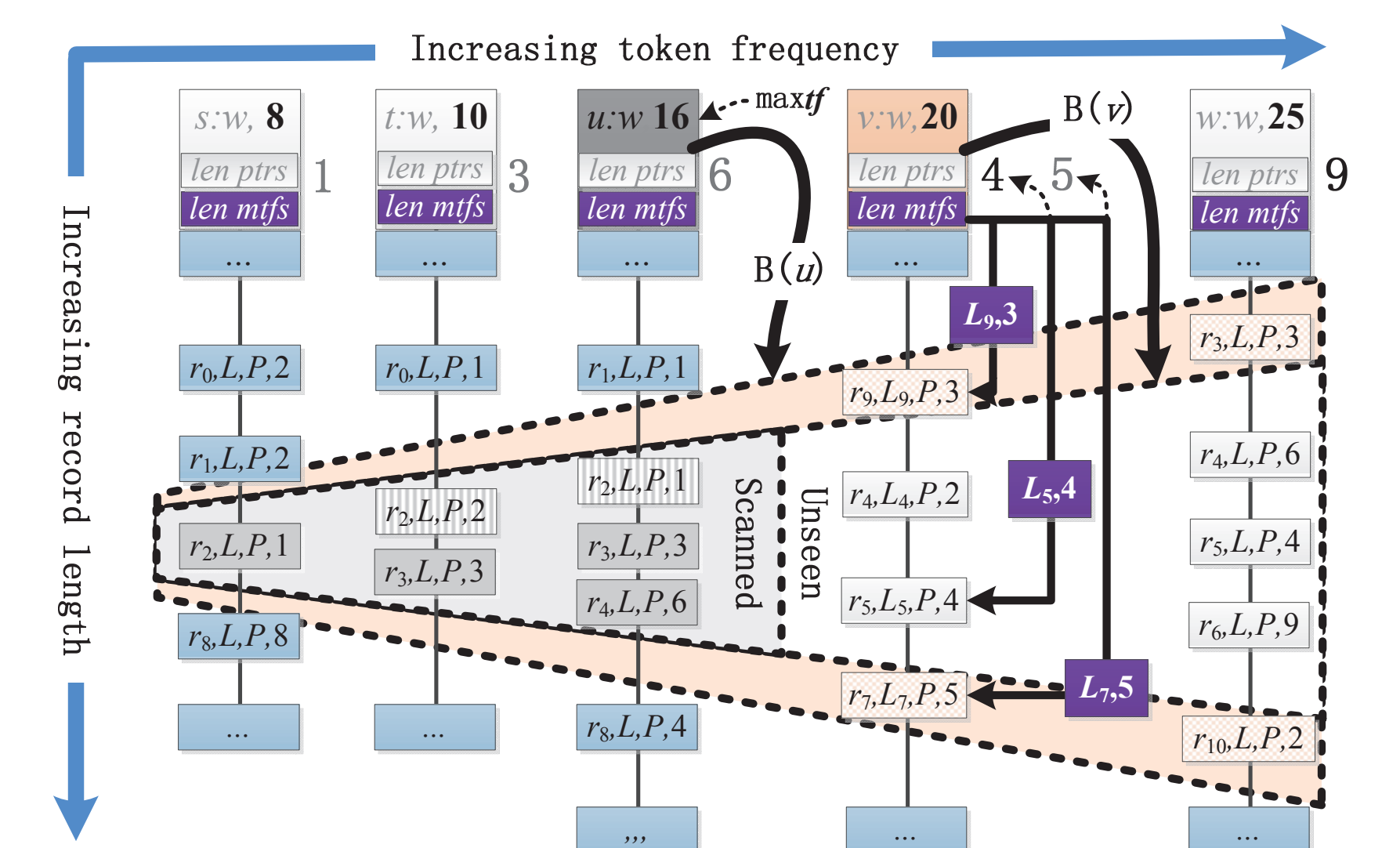
The inverted index is revised to support efficient pruning based on external storages.

- Associate token weights and degrees in the token and record entries.
- Support a skipping scheme and prefix in scanning the record lists of a token.

ELASTIC BOUNDS

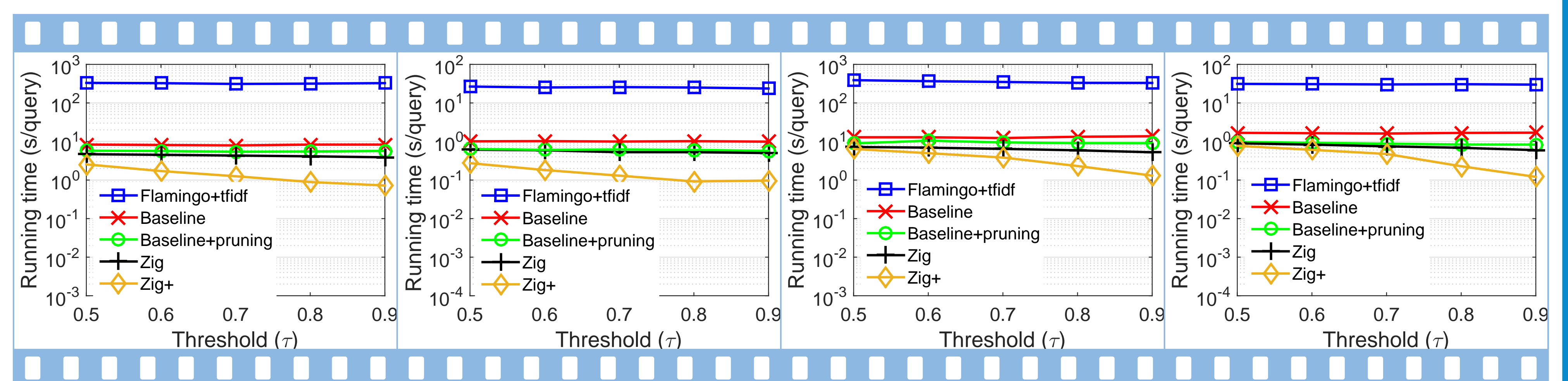
Significant improvement can be achieved as long as we iteratively tighten the bounds based on decreasing maximal degrees and reversely shrink the length bounds:

- After each step, the minimal and maximal degrees will be recomputed by the decreasing length bounded by the candidate set.
- The updated length bounds are used to shrink the range of the minimal and maximal degrees of unseen tokens.
- The maximal token degrees are associated in the index to iteratively tighten length bounds and maximal degrees until no more changes are observed.

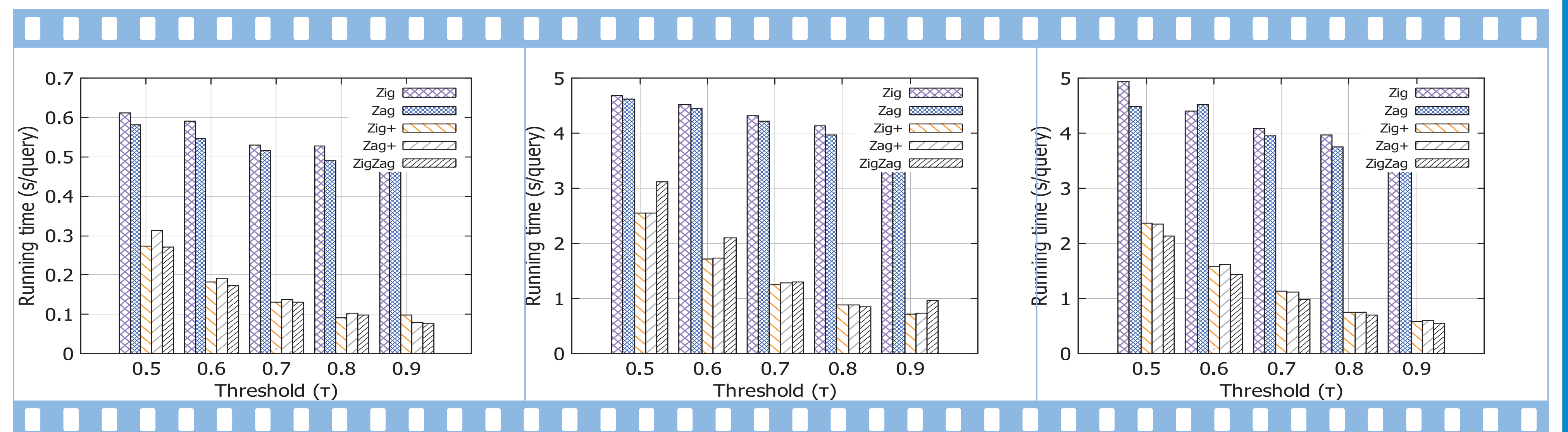


Scanning record entries using Elastic Bounds.

RESULTS



MT on HDD MT in Memory Twitter on HDD Twitter in Memory



MT in-memory MT on HDD MT on SSD


The comparisons with exiting methods demonstrated the superiority of the proposed schemes. The naïve Baseline is at least 10X faster than Flamingo, and Baseline+pruning is on average 2X better than Baseline. Zig+ is

1.2~10X better than Zig. ZigZag is the best one in the cases of using in-memory and SSD platforms, and falls short of the slow random IOs of HDD. Also, Zig+(Zag+) is always better than Zig(Zag).

A FUTURE DIRECTION

Can we associate part of the "remaining tokens" in the inverted index, allowing us to do "in-place" pruning. How to decide the associated length needs in-deep researches.

SOURCE CODE

Please contact by the following email or from Wechart by  Email: lwh@whu.edu.cn

