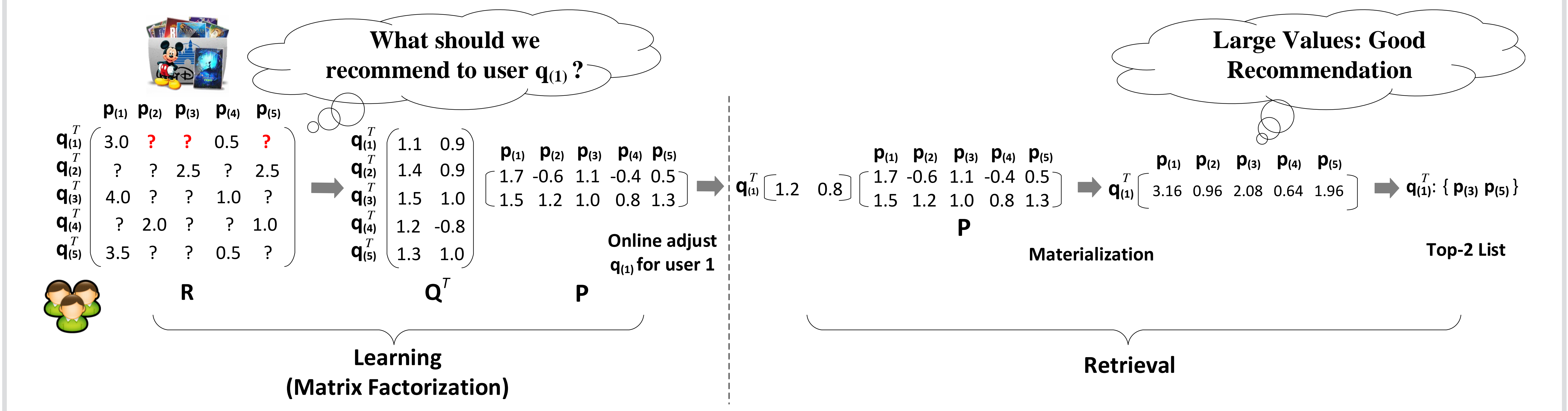


# FEXIPRO: Fast and Exact Inner Product Retrieval in Recommender Systems

Hui Li<sup>†</sup>, Tsz Nam Chan<sup>§</sup>, Man Lung Yiu<sup>§</sup>, Nikos Mamoulis<sup>†</sup>

<sup>†</sup>The University of Hong Kong    <sup>§</sup>Hong Kong Polytechnic University  
<sup>†</sup>{hli2, nikos}@cs.hku.hk    <sup>§</sup>{cstnchan, csmlyiu}@comp.polyu.edu.hk

## Two-phase Recommender Systems



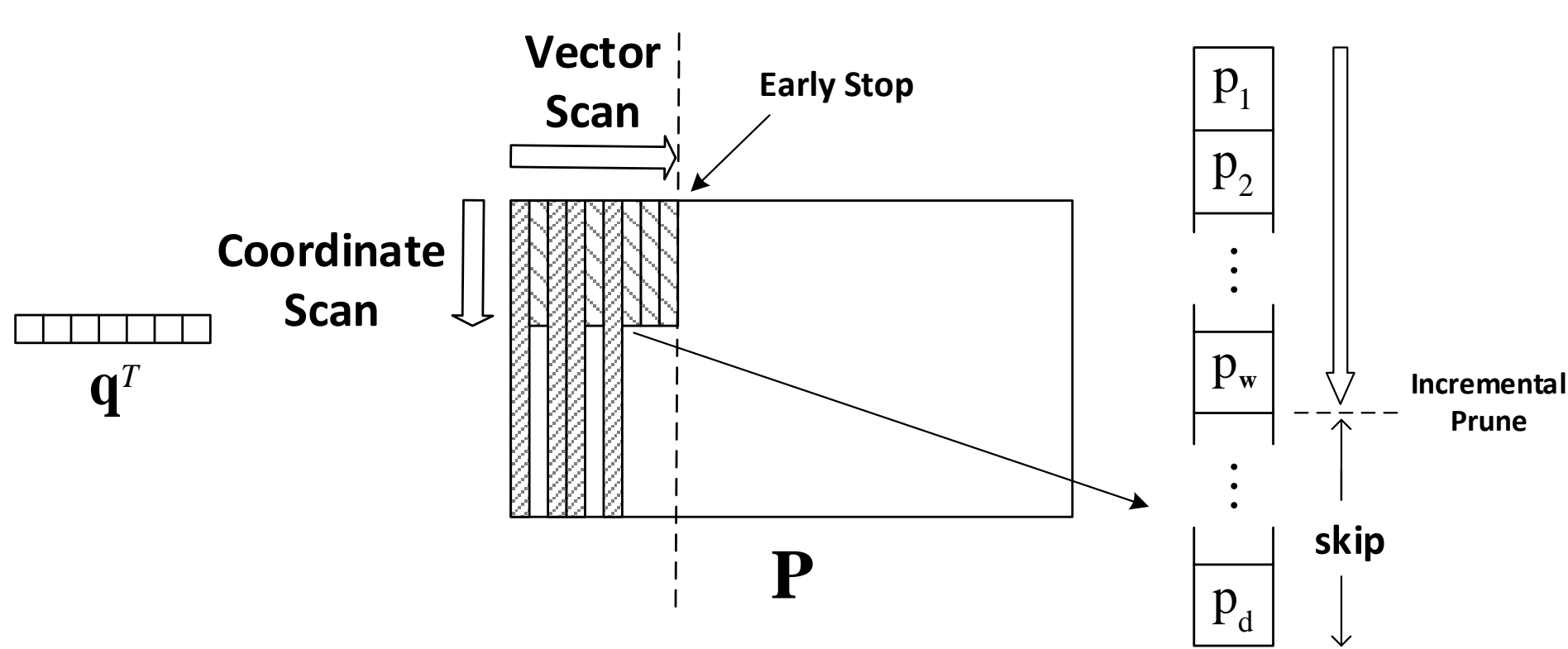
## Inner Product Retrieval

- Given a query vector  $q$  from matrix  $Q$ , find the  $k$  vectors  $p$  from matrix  $P$ , for which the  $q^T p$  values are the largest in  $q^T P$ . Ties are broken arbitrarily.
- Naive method is expensive. Materializing the entire ratings prediction list for all users is practically infeasible.

## Sequential Scan

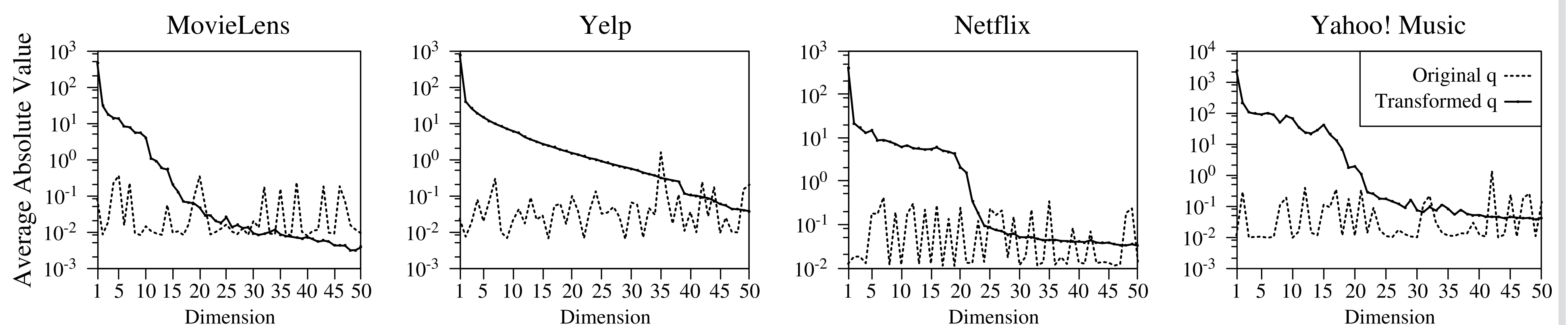
Due to **dimensionality curse**, sequential scan is the most efficient way for **exact** top- $k$  inner product retrieval:

- Vector Scan:  $q^T p \leq \|q\| \|p\|$ .
- Coordinate Scan:  $q^T p \leq q^{\ell T} p^{\ell} + \|q^h\| \|p^h\|$ .

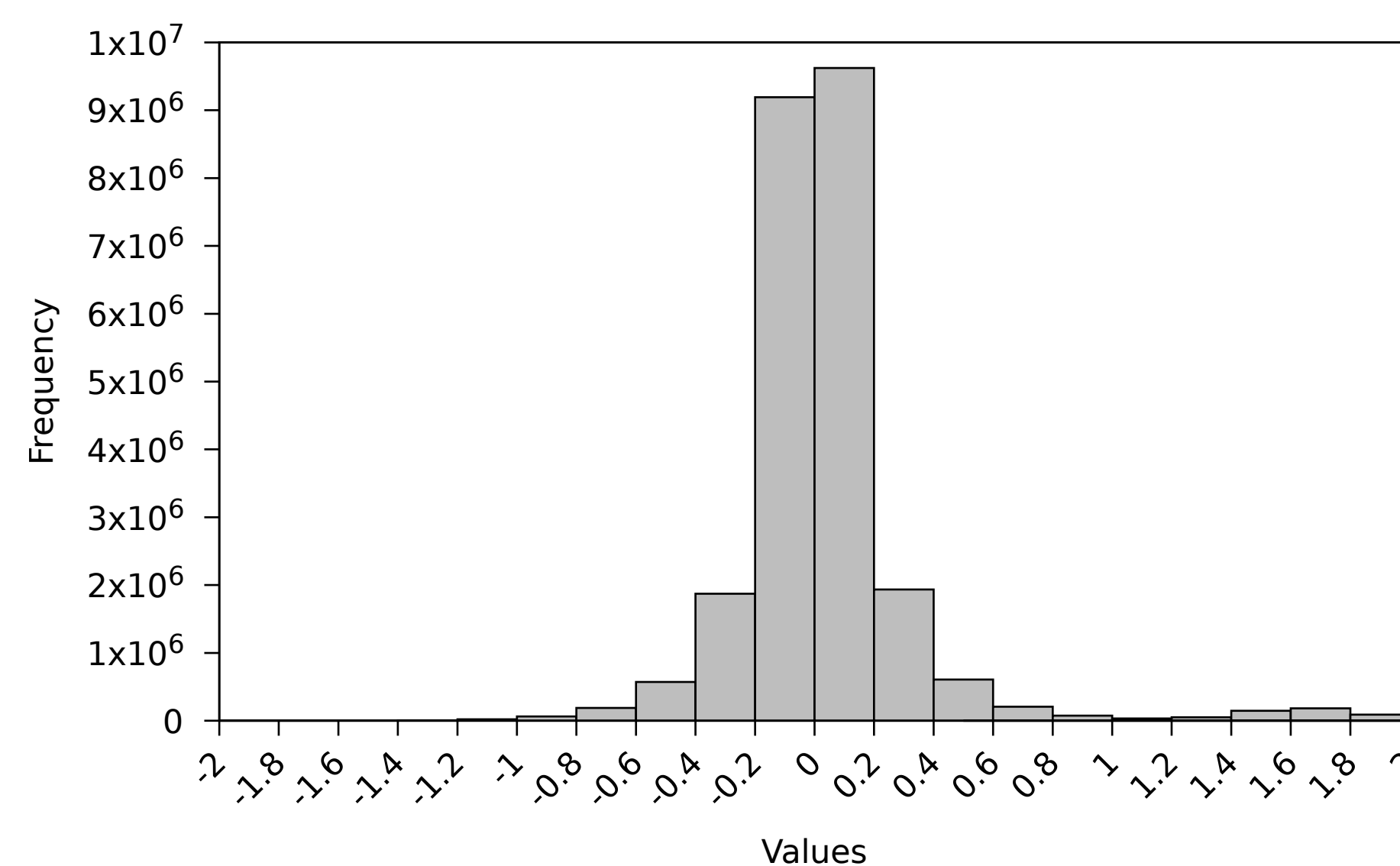


## FEXIPRO Framework

- SVD transformation:**  $q^T P = \bar{q}^T \bar{P}$ , where  $\bar{q} = \Sigma_d U^T q$  and  $\bar{P} = V_1^T$ .



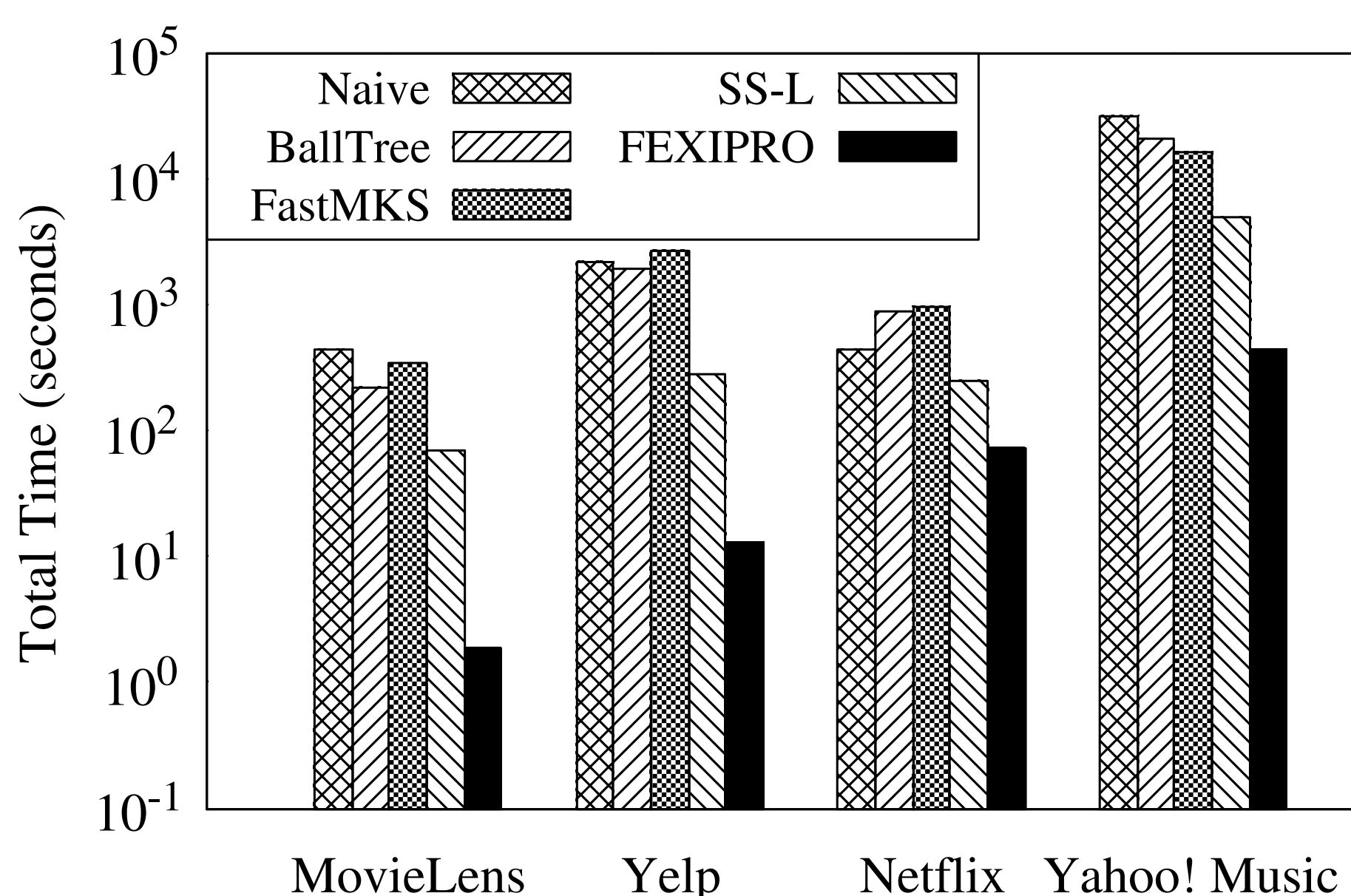
- Integer-based pruning:**  $q^T p < \sum_{s=1}^d (|q_s| \cdot |p_s| + |q_s| + |p_s| + 1)$ . It is a fast-to-compute upper bound. In order for it to be effective in recommender systems, we should scale the vector values before application.



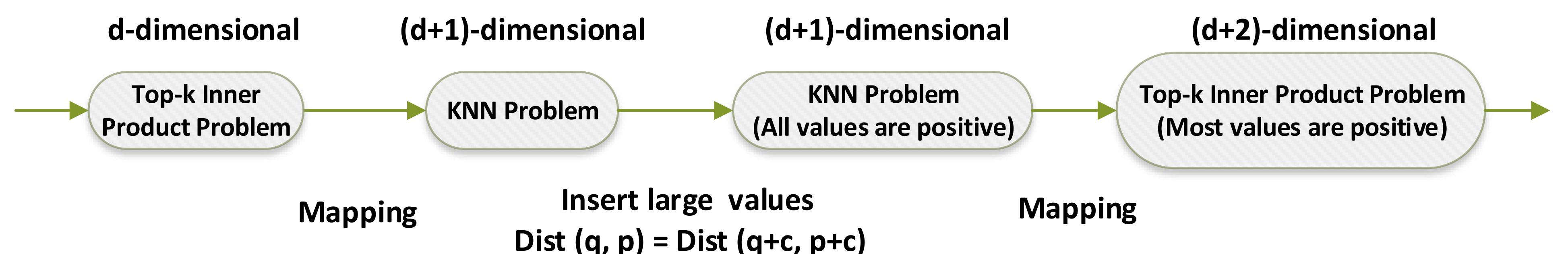
$$\begin{array}{c}
 q^T \\
 p^T \\
 \hline
 \hat{q}^T \\
 \hat{p}^T \\
 \hline
 [\hat{q}] + \Delta \hat{q} \\
 [\hat{p}] + \Delta \hat{p} \\
 \hline
 \hat{q}^T \hat{p} \leq IU(\hat{q}, \hat{p}) = 5726
 \end{array}$$

## Results for Top-1 Retrieval

FEXIPRO outperforms current exact inner product retrieval approaches typically by an order of magnitude. **Single-threaded** FEXIPRO is even faster than Intel MKL, a high-performance **multithreaded** matrix kernel library.



- Monotonicity reduction** transforms most values to positive ones, hence rendering inner products to monotonically increase as more dimensions are processed. As a result,  $q^{\ell T} p^{\ell}$  occupies the most of  $q^T p$  in incremental pruning.



## Results for Top-k Retrieval

