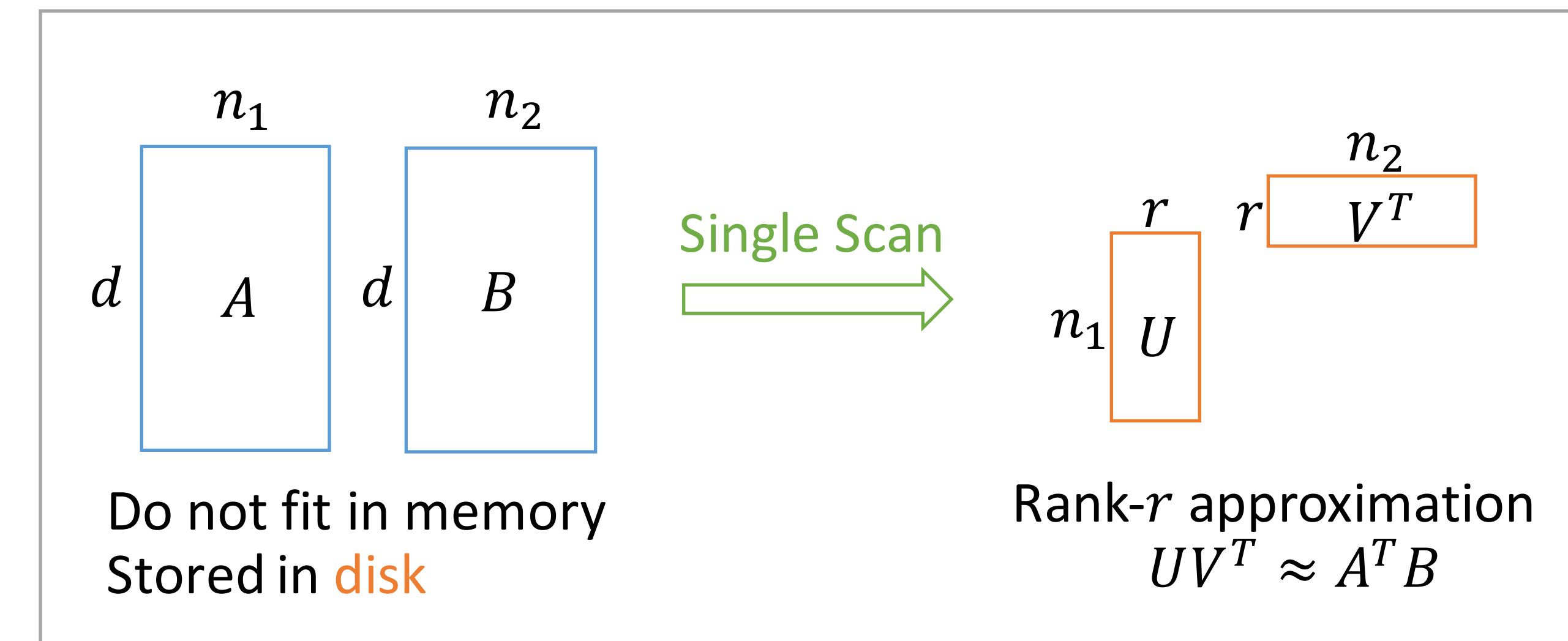


# Single Pass PCA of Matrix Products

Shanshan Wu, Srinadh Bhojanapalli, Sujay Sanghavi, Alexandros G. Dimakis

## [ Problem ]

Given two large matrices  $A$  and  $B$  (assumed too large to fit in memory), find the rank- $r$  approximation of their product  $A^T B$ , using a single pass over the matrices.



Applications

- PCA is a special case when  $A = B$
- Capture co-occurrence (e.g.,  $A$  user-by-query,  $B$  user-by-ad) or cross-covariance relation (e.g.,  $A$  genotype,  $B$  phenotype)

Why care about **single-pass**?

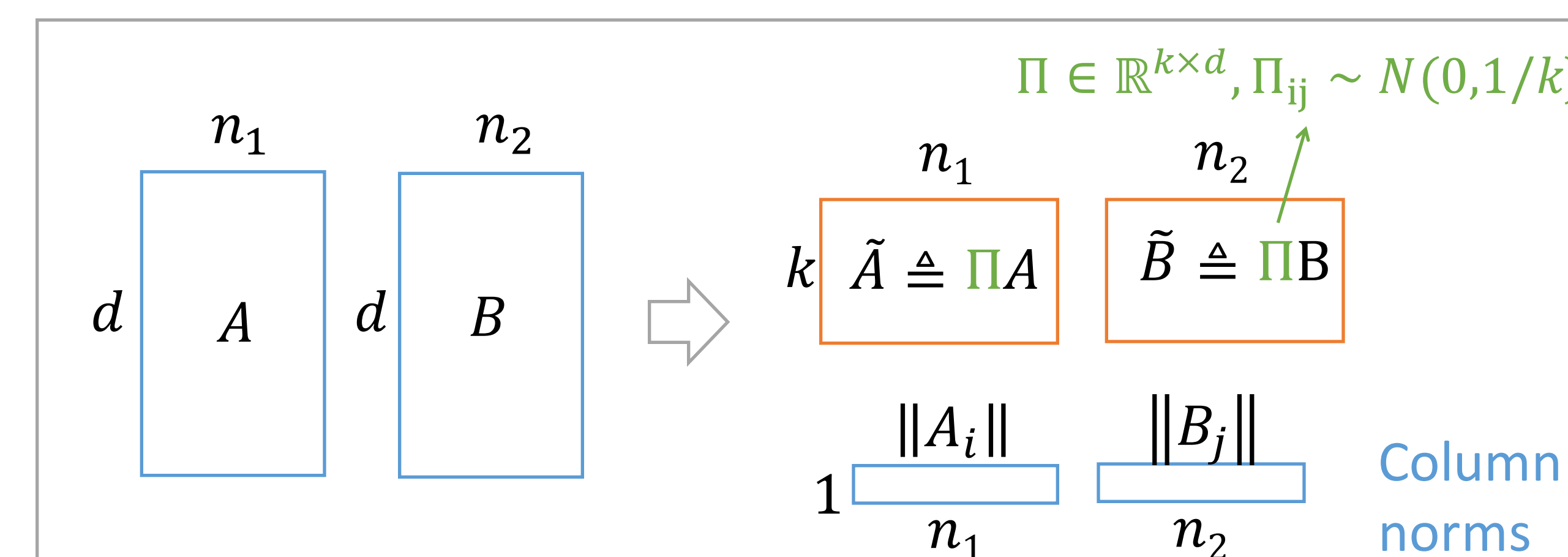
- Reduce disk I/O overhead
- Applicable when the data is streaming

Algorithm	No. of Passes	Memory
$A^T B$ + SVD	1 or more	$O(n_1 n_2)$
Direct power method on $A, B$	$O(r)$	$O(\max\{n_1, n_2\}r)$
LELA [BJS15]	2	$O(\max\{n_1, n_2\}r^3/\epsilon^2)$
SMP-PCA (Our Algorithm)	1	$O(\max\{n_1, n_2\}r^3/\epsilon^2)$

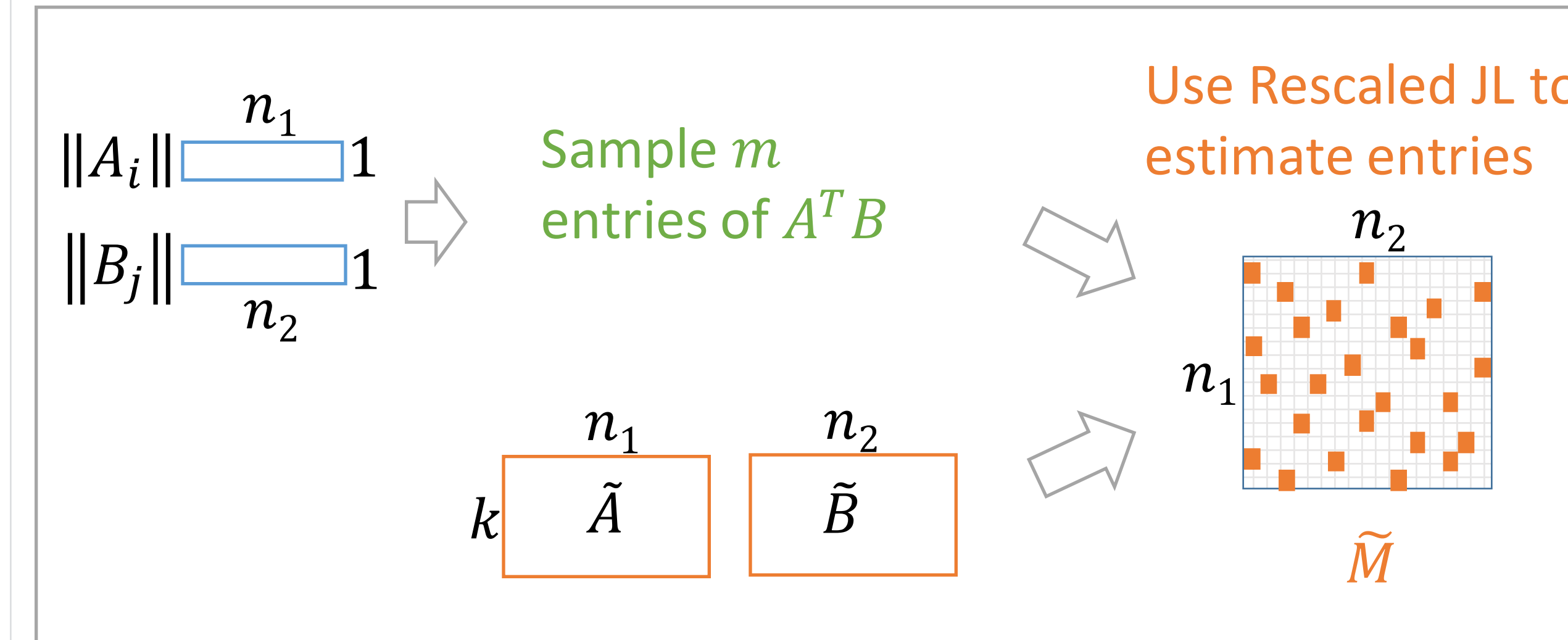
[BJS15] S. Bhojanapalli, P. Jain, and S. Sanghavi. Tighter low-rank approximation via sampling the leveraged element, SODA 2015.

## [ Our Algorithm SMP-PCA (3 steps) ]

### Step 1 Sketching



### Step 2 Sampling & estimating entries of $A^T B$



### Key Idea I Entrywise Sampling

$$\text{Sample } (A^T B)_{ij} \propto q_{ij} = m \left( \frac{\|A_i\|^2}{2n_2 \|A\|_F^2} + \frac{\|B_j\|^2}{2n_1 \|B\|_F^2} \right)$$

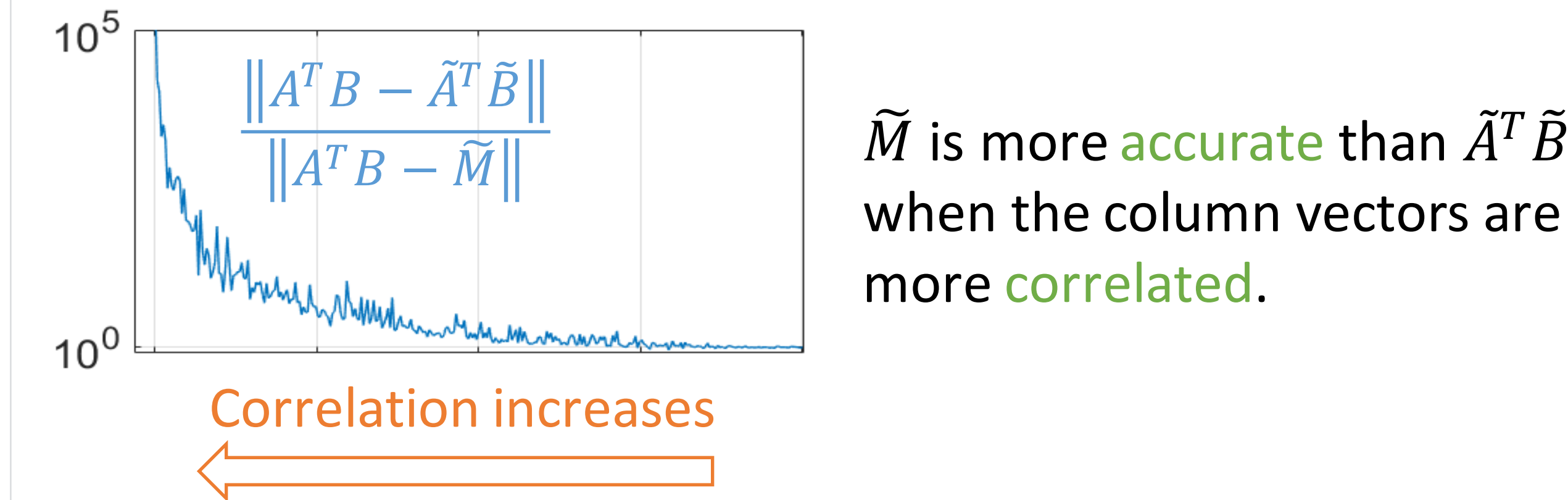
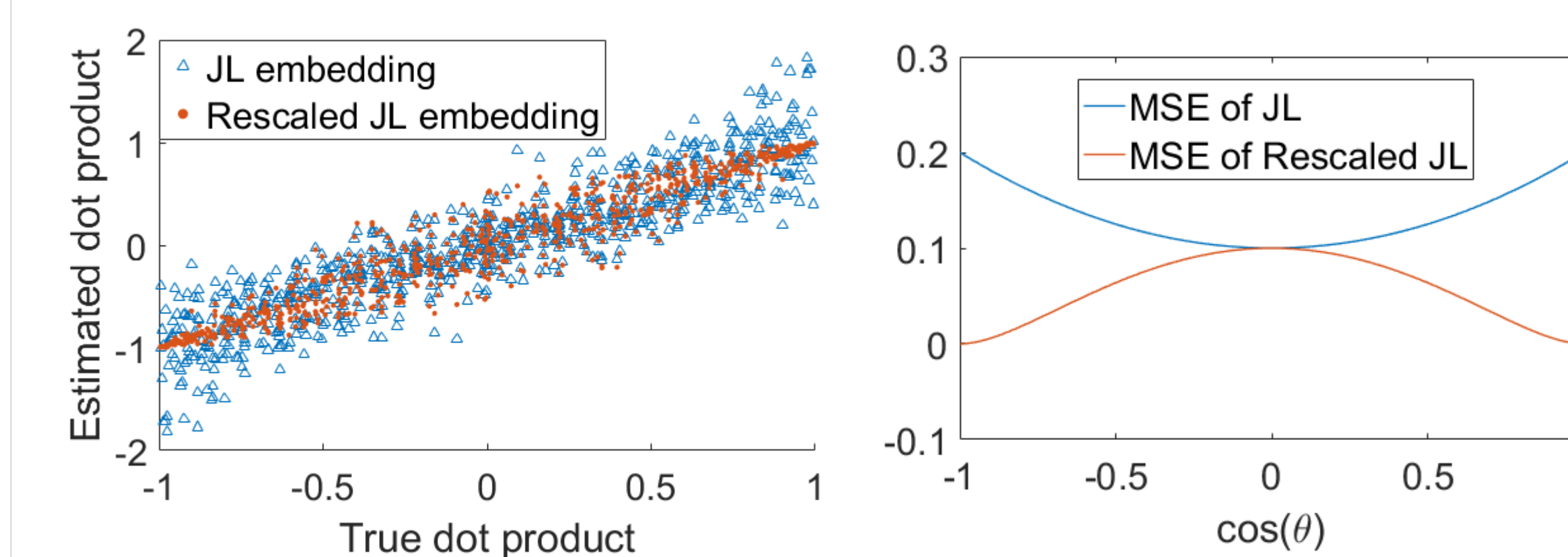
Higher weights are given to heavy rows and columns

### Key Idea II Rescaled JL

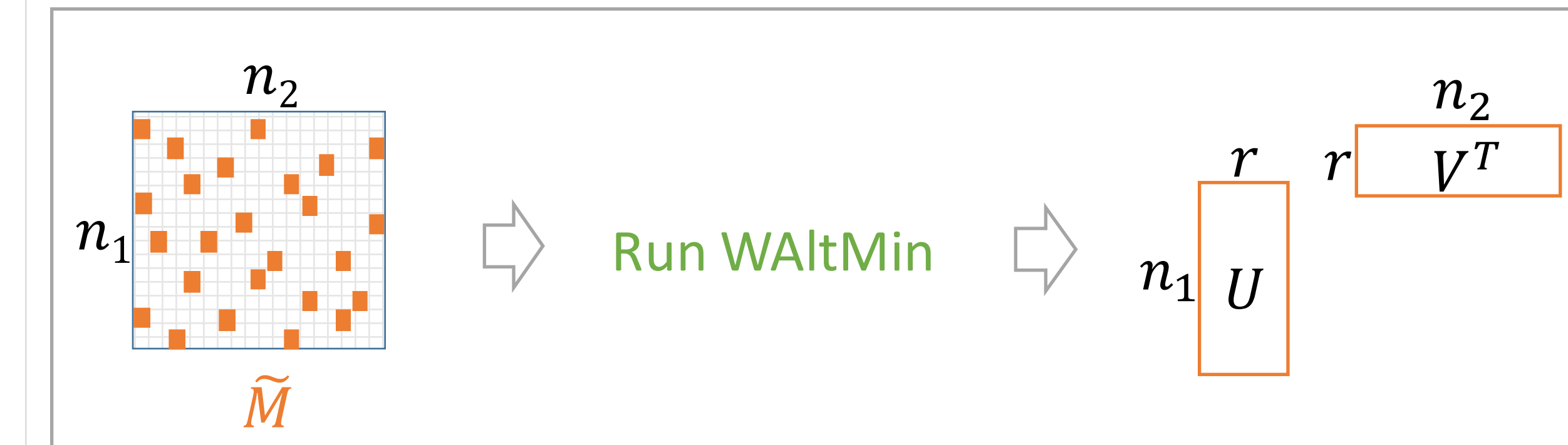
$$\tilde{M}_{ij} = \|A_i\| \cdot \|B_j\| \cdot \frac{\tilde{A}_i^T \tilde{B}_j}{\|\tilde{A}_i\| \cdot \|\tilde{B}_j\|}$$

Rescaling reduces error from distorted norms

**Observation** Rescaled JL ( $\tilde{M}_{ij}$ ) has smaller mean squared error than standard JL ( $\tilde{A}_i^T \tilde{B}_j$ ).



### Step 3 Low rank approx. from the sampled entries



WAltMin: optimize over  $U(V)$  assuming  $V(U)$  is fixed.

$$\min_{U, V} \sum_{(i, j) \in \Omega} \frac{(e_i^T U V^T e_j - \tilde{M}_{ij})^2}{\hat{q}_{ij}}$$

Weights ensure unbiasedness

## [ Theoretical Guarantee ]

**Theorem** (Informal) Let

- $\rho$  be the condition number of  $(A^T B)_r$
- $\tilde{r}$  be the maximum stable rank of  $A$  and  $B$

If the input parameters satisfy

- Sketch size  $k = O\left(\frac{r^3 \rho^2 \tilde{r} \log n}{\epsilon^2}\right)$
- The number of samples  $m = O\left(\frac{nr^3 \rho^2 \tilde{r}^2 T^2 \log n}{\gamma \epsilon^2}\right)$
- The number of WAltMin iterations  $T = O\left(\log \frac{1}{\zeta}\right)$

Then w.p.  $\geq 1 - \gamma$ , we get

$$\|(A^T B)_r - UV^T\| \leq \epsilon \|A^T B - (A^T B)_r\|_F + \zeta + \epsilon \sigma_r^*$$

Vanishes if  $A^T B$  is exactly rank- $r$

Captures the error caused by sketching

Computation Complexity:

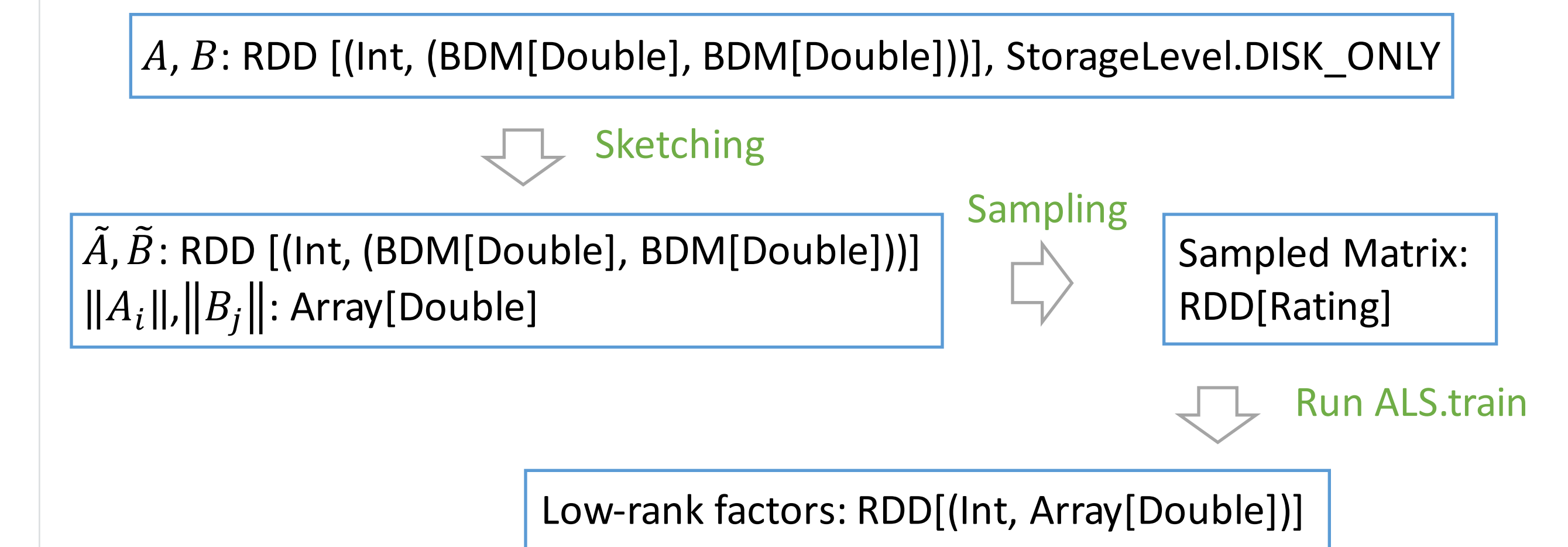
$$O(\underbrace{\text{nnz}(A)k + \text{nnz}(B)k}_{\text{Sketching}} + \underbrace{mk}_{\text{Sampling}} + \underbrace{mr^2 T}_{\text{Alternating least squares}})$$

## [ Numerical Experiments ]

### Spark Implementation

Source code <https://github.com/wushanshan/MatrixProductPCA>

Flow chat of SMC-PCA in Spark-1.6.2:

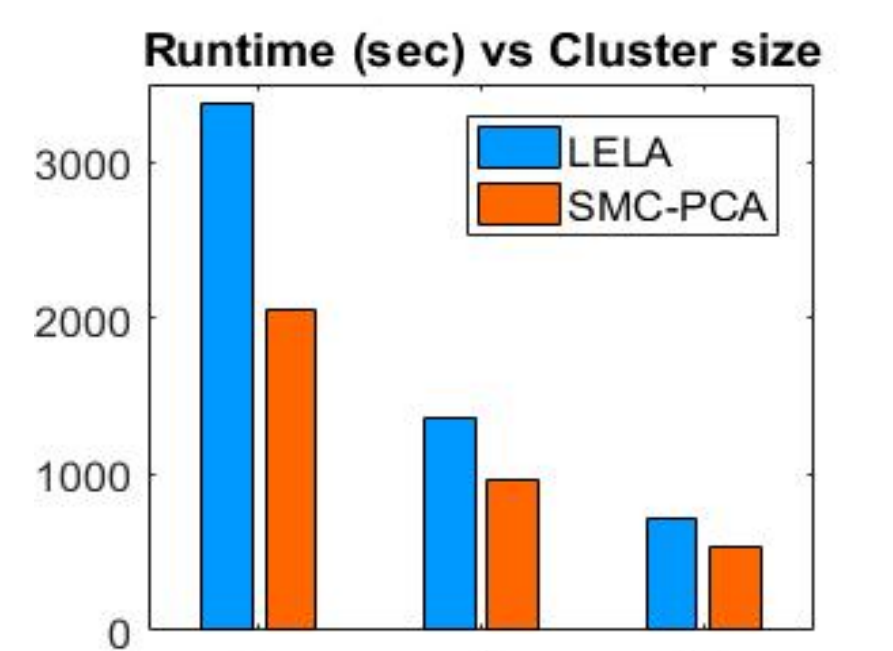


### Synthetic Dataset

- Spark-1.6.2, m3.2xlarge EC2 instances
- $A, B$ : 100k-by-100k dense matrices (150GB)
- Use SRHT as the sketching step

Algorithm	Error	Runtime [1]
Exact SVD[2]	0.0271	23 hrs
LELA	0.0274	56 mins
SMP-PCA	0.0280	34 mins

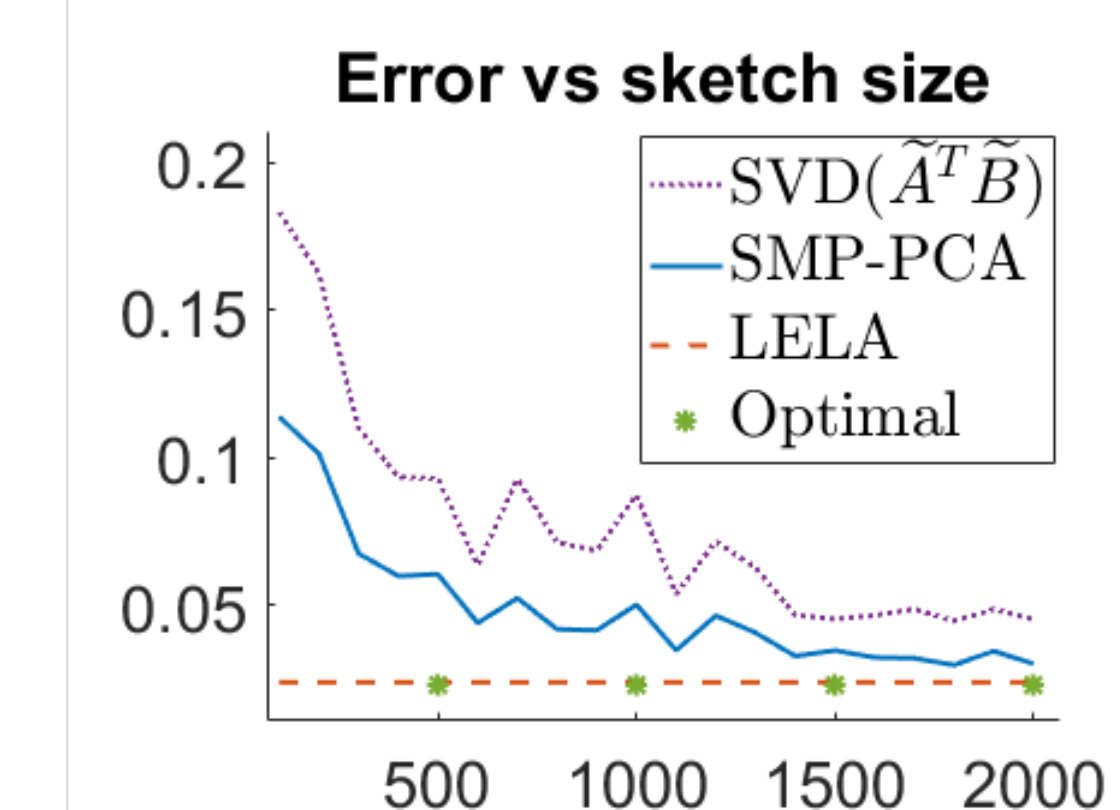
[1] Runtime on a cluster with two m3.2xlarge instances  
[2] Perform power method directly on  $A, B$



Runs faster with smaller loss of accuracy   Scales as the cluster size grows

### Real Dataset

- SMP-PCA outperforms SVD( $\tilde{A}^T \tilde{B}$ ) because of Rescaled JL
- SMP-PCA achieves similar error as the two-pass LELA



Dataset	Optimal	LELA	SMP-PCA
URL-malicious Size: 792k-by-10k	0.0163	0.0182	0.0188
URL-benign Size: 1603k-by-10k	0.0103	0.0105	0.0117