



# Single-Pass PCA of Large High-Dimensional Data

**Wenjian Yu<sup>1</sup>, Yu Gu<sup>1</sup>, Jian Li<sup>1</sup>, Shenghua Liu<sup>2</sup>, Yaohang Li<sup>3</sup>**

<sup>1</sup>Department of Computer Science and Technology,  
*Tsinghua University*, Beijing 100084, China

<sup>2</sup>Institute of Computing Technology, *Chinese Academy of Science*

<sup>3</sup>Department of Computer Science, *Old Dominion University*, USA  
[yu-wj@tsinghua.edu.cn](mailto:yu-wj@tsinghua.edu.cn)

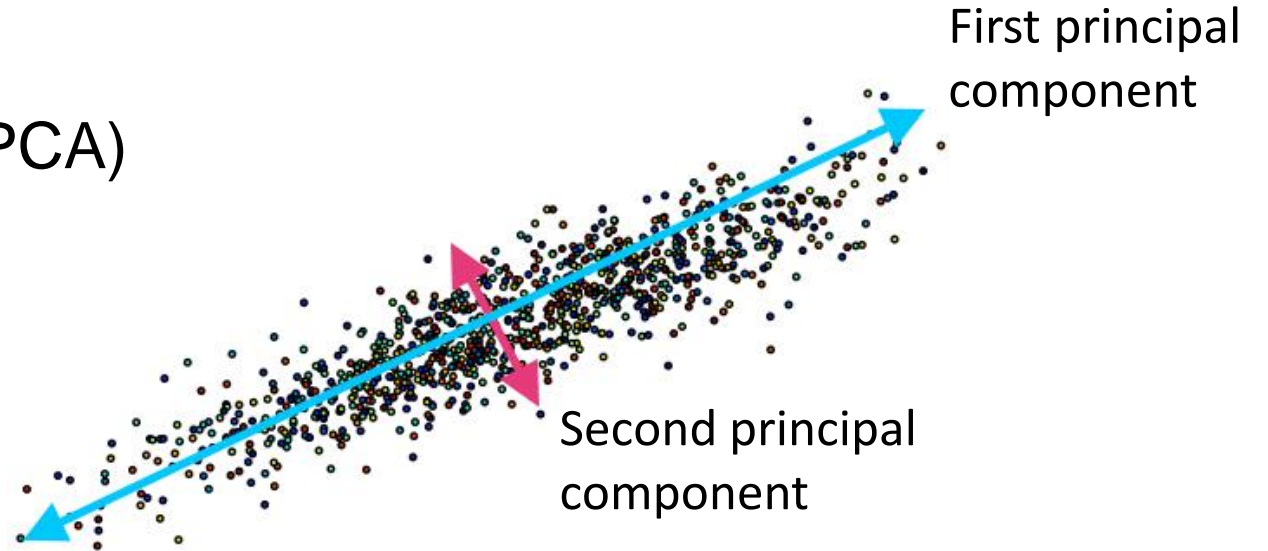
# Outline

- **Introduction**
- **Technical Background**
- **The Single-Pass Algorithm for PCA**
- **Experimental Results**
- **Conclusion**

# Introduction

## ■ Background

- Principal component analysis (PCA)
- **An open problem**: calculate PCA of large-size and high-dimensional dense data in a limited-memory computer
- **A single-pass algorithm**: particularly useful / efficient, for data stored in slow memory or streaming data
- There are single-pass PCA algorithms for **SPSD matrix** or **low-dimensional** data, but the study for the algorithm for more general matrices is not sufficient.



# Introduction

- Randomized **matrix** algorithm

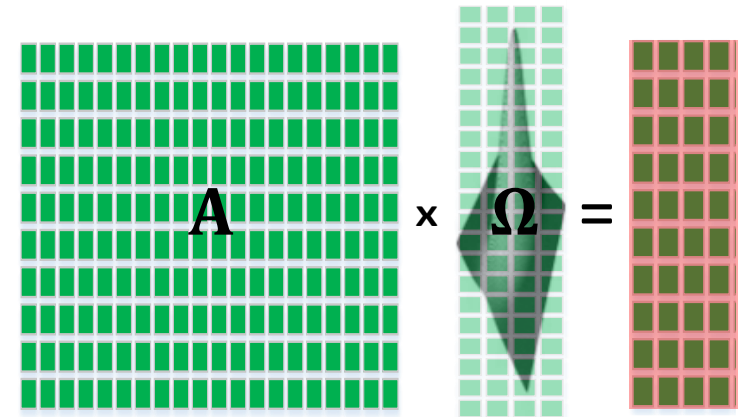
- Has advantages over traditional algorithms (like SVD)  
(faster runtime, better parallelism, pass-efficient; suitable for large data)

- **randQB** based on *random projection* <sup>[1]</sup>

$$\mathbf{A} \approx \mathbf{QB}$$

$\mathbf{Q}$  captures the dominant actions of  $\mathbf{A}$   
Small sketch  $\mathbf{B}$  facilitates computation

- Applied to computing PCA <sup>[2]</sup>
- Useful for distributed PCA; excellent performance on parallel computers
- A blocked version for rank-revealing matrix factorization <sup>[3]</sup>



[1] **N Halko, P-G Martinsson, J A Tropp**, “Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix Decompositions,” **SIAM review**, 2011

[2] **N Halko, et al.**, “An algorithm for the principal component analysis of large data sets,” **SISC**, 2011

[3] **P-G Martinsson and S. Voronin**, “A randomized blocked algorithm for efficiently computing rank-revealing factorizations of matrices,” **SISC**, 2016

# Introduction

- Our contribution

- We reconstruct the blocked randQB algorithm <sup>[3]</sup> to obtain a single-pass PCA algorithm
- Single-pass
  - involves only one pass over specified large high-dimensional data
- Efficiency
  - $O(mnk)$  time complexity and  $O(k(m+n))$  space complexity for computing k-PCA, and well adapts to parallel computing
- Accuracy
  - same theoretic error bounds as the randomized blocked algorithm; much less error than its counterpart

[3] **P-G Martinsson** and **S. Voronin**, “A randomized blocked algorithm for efficiently computing rank-revealing factorizations of matrices,” *SISC*, 2016

# Technical Background – SVD and PCA

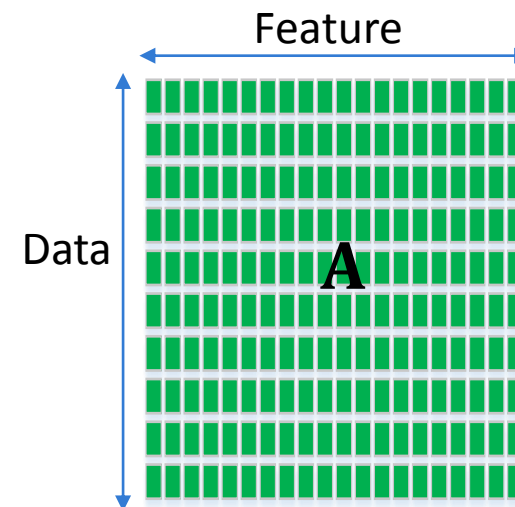
- Truncated singular value decomposition (SVD)

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \quad \longrightarrow \quad \mathbf{A} \approx \mathbf{A}_k = \mathbf{U}_k\mathbf{\Sigma}_k\mathbf{V}_k^T$$

- $\mathbf{A}_k$ : rank-k approximation of  $\mathbf{A}$  (optimal in  $l_2$ -norm and F-norm)

- SVD and PCA are closely related

- Suppose each row of matrix  $\mathbf{A}$  is an observed data
- PCA is realized through truncated SVD
- The leading right singular vectors ( $\mathbf{v}_i$ ) of  $\mathbf{A}$  are the **principal components**.  
Particularly,  $\mathbf{v}_1$  is the first principal component



# Technical Background – Randomized SVD

## ■ Basic randQB scheme

- Produce near-optimal low-rank appr.
- Accuracy can be improved with power iteration scheme
- Well suit to parallel computing
- Result has small random variance
- Better than the column-pivoted QR

## ■ A single-pass variant

- Reduce to 1 visit of  $A$

$$A \approx QQ^T A \tilde{Q} \tilde{Q}^T = QB \tilde{Q}^T$$

- More approximation is included

---

### Algorithm 1 Basic randomized scheme for truncated SVD

---

**Require:**  $A \in \mathbb{R}^{m \times n}$ , rank  $k$ , over-sampling parameter  $s$ .

- 1:  $l = k + s$ ;
- 2:  $\Omega = \text{randn}(n, l)$ ;
- 3:  $Q = \text{orth}(A\Omega)$ ;
- 4:  $B = Q^T A$ ;
- 5:  $[\tilde{U}, S, V] = \text{svd}(B)$ ;
- 6:  $U = Q\tilde{U}$ ;
- 7:  $U = U(:, 1:k)$ ;  $V = V(:, 1:k)$ ;  $S = S(1:k, 1:k)$ ;
- 8: **return**  $U, S, V$ .

}  $A$  is visited twice  
←

---

### Algorithm 2 An existing single-pass algorithm

---

**Require:**  $A \in \mathbb{R}^{m \times n}$ , rank parameter  $k$ .

- 1: Generate random  $n \times k$  matrix  $\Omega$  and  $m \times k$  matrix  $\tilde{\Omega}$ ;
- 2: Compute  $Y = A\Omega$  and  $\tilde{Y} = A^T \tilde{\Omega}$  in a single pass over  $A$ ;
- 3:  $Q = \text{orth}(Y)$ ;  $\tilde{Q} = \text{orth}(\tilde{Y})$ ;
- 4: Solve linear equation  $\tilde{\Omega}^T QB = \tilde{Y}^T \tilde{Q}$  for  $B$ ;
- 5:  $[\tilde{U}, S, \tilde{V}] = \text{svd}(B)$ ;
- 6:  $U = Q\tilde{U}$ ;  $V = \tilde{Q}\tilde{V}$ ;
- 7: **return**  $U, S, V$ .

# The Single-Pass PCA Algorithm

## The blocked randQB algorithm <sup>[3]</sup>

```
function [Q, B] = randQB_b (A, ε, b)
(1) for i = 1, 2, 3, ...
(2)   Ωi = randn(n, b);
(3)   Qi = orth(AΩi);
(4)   Qi = orth(Qi - ∑j=1i-1 QjQjTQi);
(5)   Bi = QiTA;
(6)   A = A - QiBi;
(7)   if ||A|| < ε then stop
(8) end for
(9) Q = [Q1 ... Qi]; B = [B1T ... BiT]T.
```

- Mathematically equivalent to the basic randQB algorithm (Gram-Schmidt procedure)
- Iterative blocked procedure for monitoring approximation error while keeping high efficiency
- Mainly aimed at the problem of adaptive rank determination

**Convert it to a pass-efficient procedure  
(multiplications with **A** moved out of loop)**

**Theorem 1:** The **Q** and **B** obtained with Alg. 3 satisfy: **Q** is orthonormal and **B** = **Q**<sup>T</sup>**A**

## Algorithm 3 A pass-efficient blocked algorithm

**Require:**  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , rank parameter  $k$ , block size  $b$ .

```
1: Q = []; B = [];
2: Ω = randn(n, k);
3: G = AΩ;
4: H = ATG;
5: for i = 1, 2, ..., k/b do
6:   Ωi = Ω(:, (i-1)b + 1 : ib);
7:   Yi = G(:, (i-1)b + 1 : ib) - Q(BΩi);
8:   [Qi, Ri] = qr(Yi);
9:   Bi = Ri-T(H(:, (i-1)b + 1 : ib)T - ΩiTBTB);
10:  Q = [Q, Qi]; B = [BT, BiT]T;
11: end for
```



# The Single-Pass PCA Algorithm

## ■ Algorithm 3

- Equivalent to the randQB alg.
  - Steps 3 and 4 can be executed with only one pass over  $A$
- Add re-orthogonalization steps to alleviate round-off error
  - Memory cost is about  $(m + 2n)l$  floating numbers
  - Time complexity (flop count) close to the basic randQB (Alg.1)
  - With the power scheme, accuracy can be largely improved at the cost of one more visit of  $A$

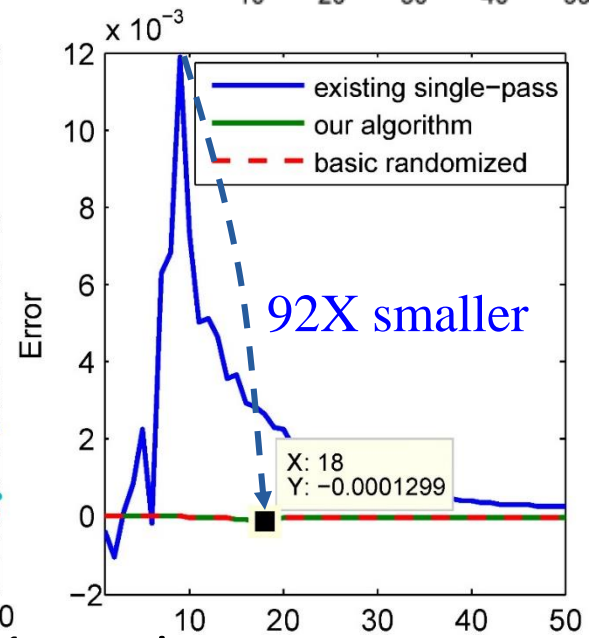
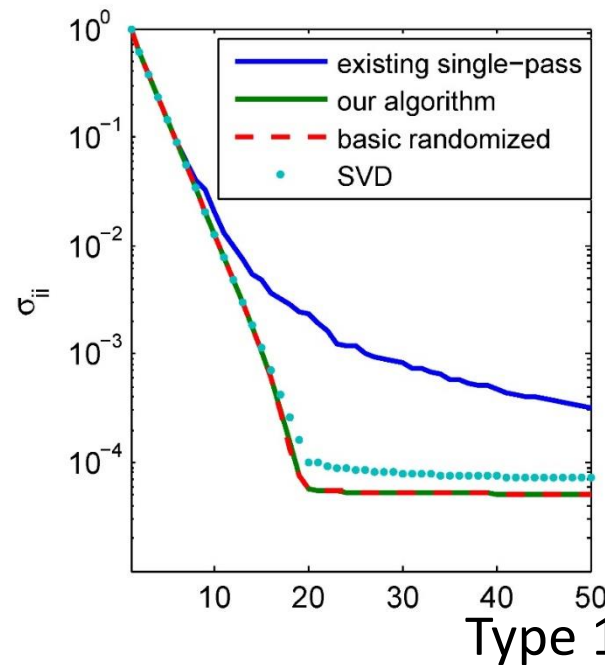
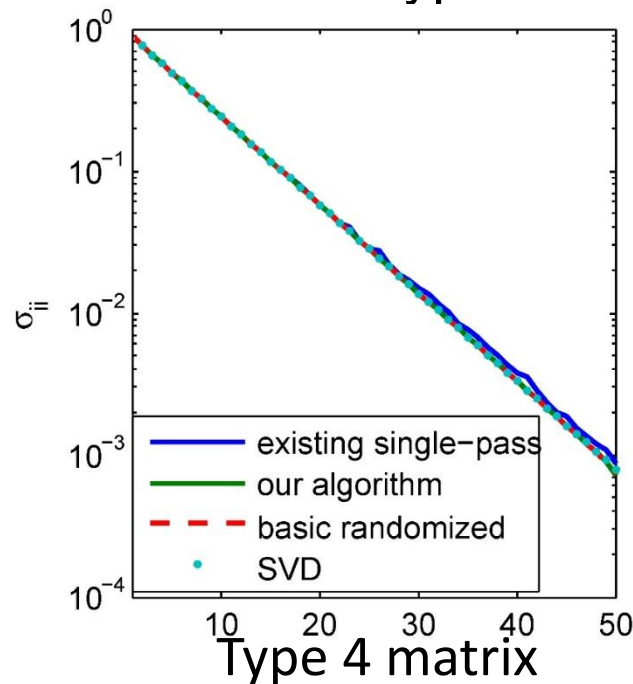
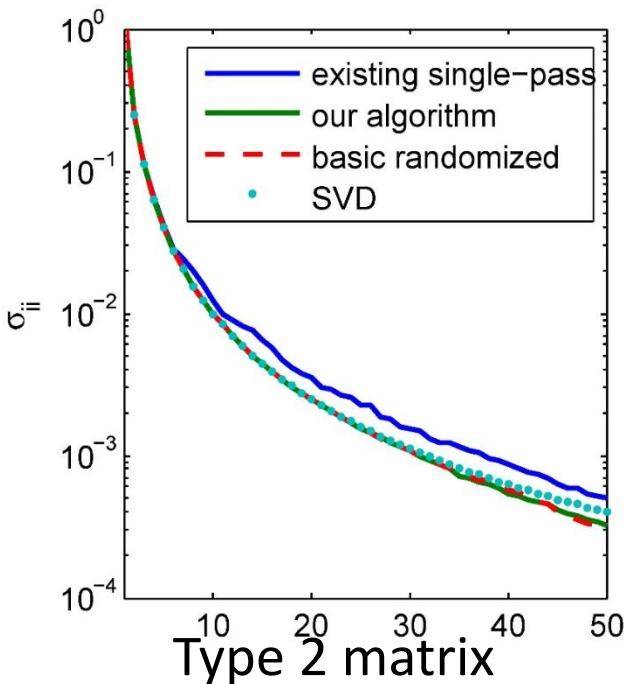
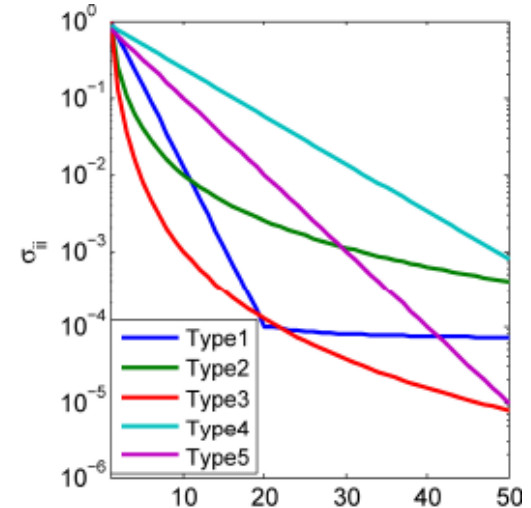
## Algorithm 4 A single-pass algorithm for computing PCA

**Require:**  $A \in \mathbb{R}^{m \times n}$ , rank parameter  $k$ , block size  $b$ .

```
1:  $Q = []$ ;  $B = []$ ;
2: Choose  $l = tb$ , which is slightly larger than  $k$ ;
3:  $\Omega = \text{randn}(n, l)$ ;  $G = []$ ; Set  $H$  to an  $n \times l$  zero matrix;
4: while  $A$  is not completely read through do
5:   Read next few rows of  $A$  into RAM, denoted by  $a$ ;
6:    $g = a\Omega$ ;  $G = [G; g]$ ;
7:    $H = H + a^T g$ ;
8: end while
9: for  $i = 1, 2, \dots, t$  do
10:   $\Omega_i = \Omega(:, (i-1)b + 1 : ib)$ ;
11:   $Y_i = G(:, (i-1)b + 1 : ib) - Q(B\Omega_i)$ ;
12:   $[Q_i, R_i] = \text{qr}(Y_i)$ ;
13:   $[Q_i, \tilde{R}_i] = \text{qr}(Q_i - Q(Q^T Q_i))$ ;
14:   $R_i = \tilde{R}_i R_i$ ;
15:   $B_i = R_i^{-T} (H(:, (i-1)b+1 : ib))^T - Y_i^T Q B - \Omega_i^T B^T B$ ;
16:   $Q = [Q, Q_i]$ ;  $B = [B^T, B_i^T]^T$ ;
17: end for
18:  $[\tilde{U}, S, V] = \text{svd}(B)$ ;
19:  $U = Q\tilde{U}$ ;
20:  $U = U(:, 1 : k)$ ;  $V = V(:, 1 : k)$ ;  $S = S(1 : k, 1 : k)$ ;
21: return  $U, S, V$ .
```

# Experimental Results

- Five types of test matrices
  - Specified singular value spectrum with various decaying trend
  - For type 1 and 2 matrices, the singular value decays asymptotically slow
- Accuracy of computed singular value
  - A 3000x3000 matrix for each type



# Experimental Results

- Accuracy of the principal components
  - Test on type 1 matrix (with *slowly-decayed*  $\sigma_{ii}$ 's); smaller error for other matrices
  - For  $\mathbf{v}_1$ , only  $2.8 \times 10^{-5}$  difference in  $l_\infty$ -norm
  - For the first 10 principal components, the correlation coefficients are calculated: 0.9993 ~1 even for the 10<sup>th</sup> component
- Runtime comparison (200,000x200,000 matrices)
  - Each matrix stored as a 149 GB hard-disk file
  - Alg. 4 is **2X faster** than Alg. 1; more accurate than Alg. 2

Experiment on a computer  
with two 12-core Xeon  
CPUs and 32GB memory

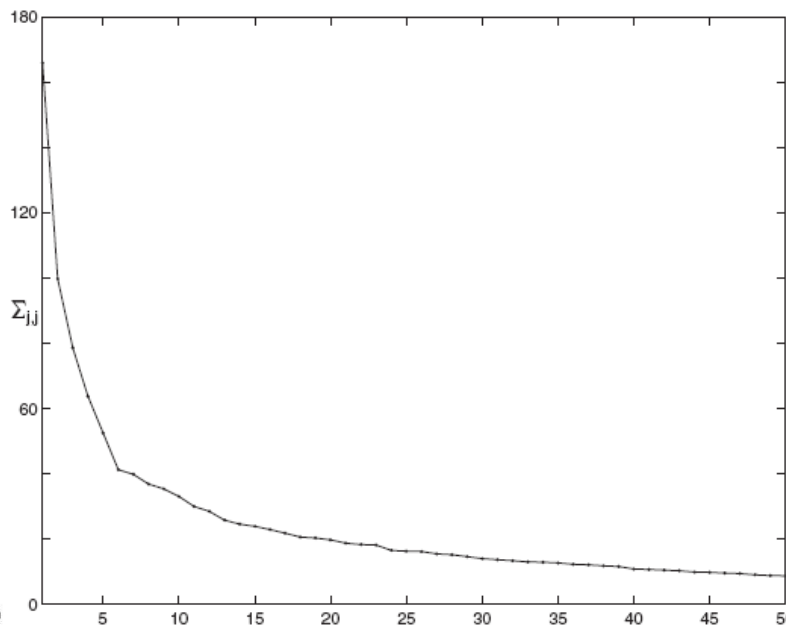
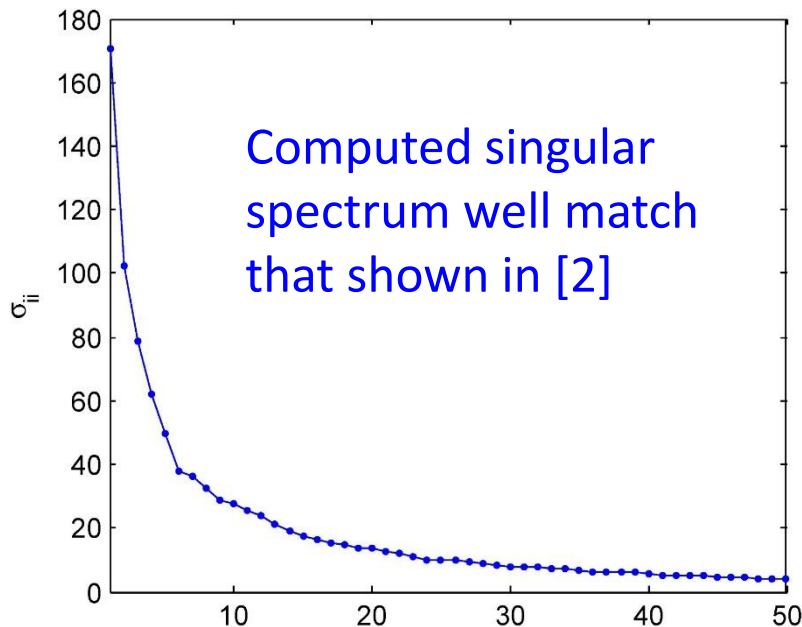
Matrix	$k$	Algorithm 1			Algorithm 2			Algorithm 4		
		$t_{read}$	$t_{PCA}$	max_err	$t_{read}$	$t_{PCA}$	max_err	$t_{read}$	$t_{PCA}$	max_err
Type1	16	2390	2607	1.7e-3	1186	1404	2.2e-2	1206	1426	1.8e-3
Type1	20	2420	2616	9e-4	1198	1380	1.6e-1	1217	1413	1.2e-3
Type1	24	2401	2593	1e-3	1216	1400	1.5e-1	1216	1414	1.2e-3
Type2	12	2553	2764	5e-4	1267	1477	3e-2	1276	1490	5e-4
Type3	24	2587	2777	1e-5	1312	1500	1.7e-3	1310	1502	2e-5

Test on a  $10^4 \times 10^4$  matrix:  
Our Alg. 4 is **over 300X**  
**faster** than *svd/svds*

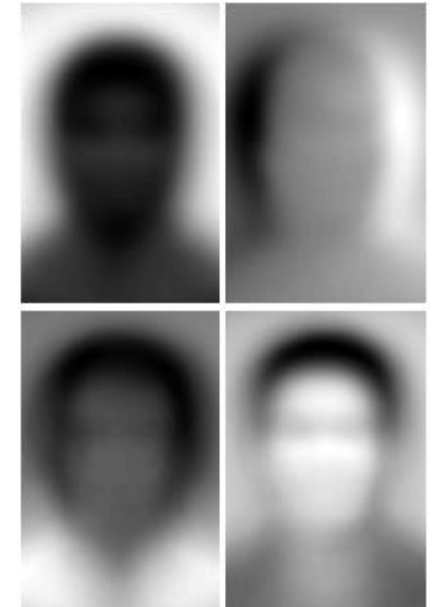
# Experimental Results

## ■ A test of real data

- Face images from the FERET [4]
- As in [2], construct a 102,042x392,216 matrix (*150GB file on hard disk*)
- Compute 50 eigenfaces on the machine with 24 CPU cores
- Runtime of our algorithm: ~ **24 minutes**



Four eigenfaces



[2] **N Halko, et al.**, "An algorithm for the principal component analysis of large data sets," **SISC**, 2011

[4] **P J Phillips, et al.**, "The FERET evaluation methodology for face-recognition algorithms," **T-PAMI**, 2000

# Conclusion

- **A single-pass PCA algorithm for large and high-dimensional data**
- **Only one pass over data matrix, providing that the matrix is stored in a row-major format**
- **Comparable accuracy to existing randomized algorithm; much less error than an existing single-pass algorithm**
- **Experiments demonstrate the algorithm's effectiveness for large-size high-dimensional data (~150 GB disk file), in terms of runtime and memory usage**

*The codes of the proposed algorithm and experimental data are shared on:*  
<https://github.com/WenjianYu/rSVD-single-pass>

*Thank You !*

