

Written by: Edo Liberty  
Presented by: Ryan Rogers  
With Some Slides from: Edo Liberty

# SIMPLE AND DETERMINISTIC MATRIX SKETCHING

# Set up

- $A$  is an  $n \times m$  matrix
- We want to compute the  $m \times m$  matrix:  $A^T A$
- Problem:  $n >$  machine memory.
- Goal: Find 'good' approximate  $d \times m$  matrix  $B$  for any  $\|x\| = 1$

$$\|A^T A - B^T B\| \leq \text{Small}$$

# Set up

- $A$  is an  $n \times m$  matrix
- We want to compute the  $m \times m$  matrix:  $A^T A$
- Problem:  $n >$  machine memory.
- Goal: Find 'good' approximate  $d \times m$  matrix  $B$  for any  $\|x\| = 1$

$$\|A^T A - B^T B\| \leq \varepsilon \|A\|_f^2$$

# Sketches



# Sketches

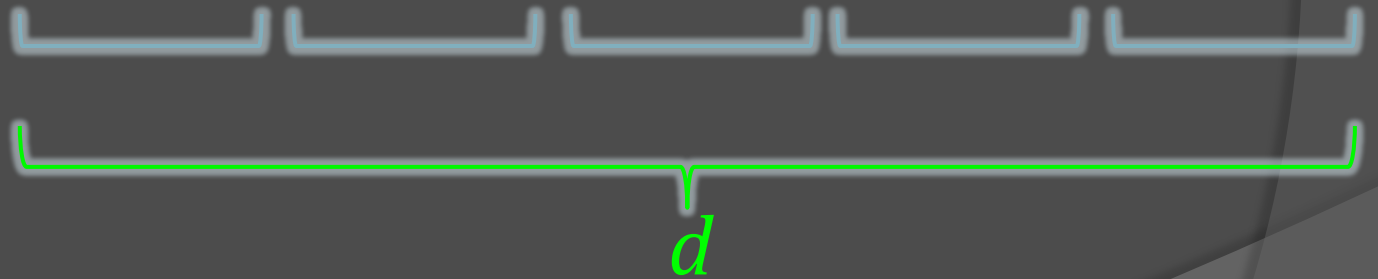
- ⦿ A **sketch** of a matrix  $A$  is another matrix  $B$ , that is significantly smaller than  $A$  but still approximates  $A$  well.
- ⦿ We need this if:
  - Rows of matrix can be processed only once
  - Storage is limited

# Frequent Items

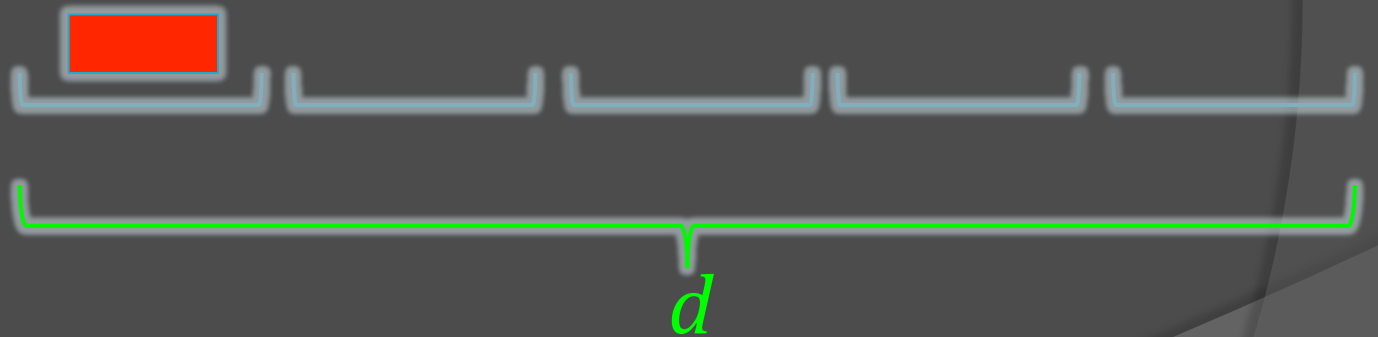
- ⦿ Universe  $U = \{a_1, \dots, a_m\}$  and a stream  $A_1, A_2, \dots, A_n$
- ⦿ Frequency  $f_i$  of item  $a_i$  in the stream
- ⦿ Use only  $O(d)$  space to produce approximate counts  $g_i$ , such that

$$|f_i - g_i| < n / d$$

# Frequent Items

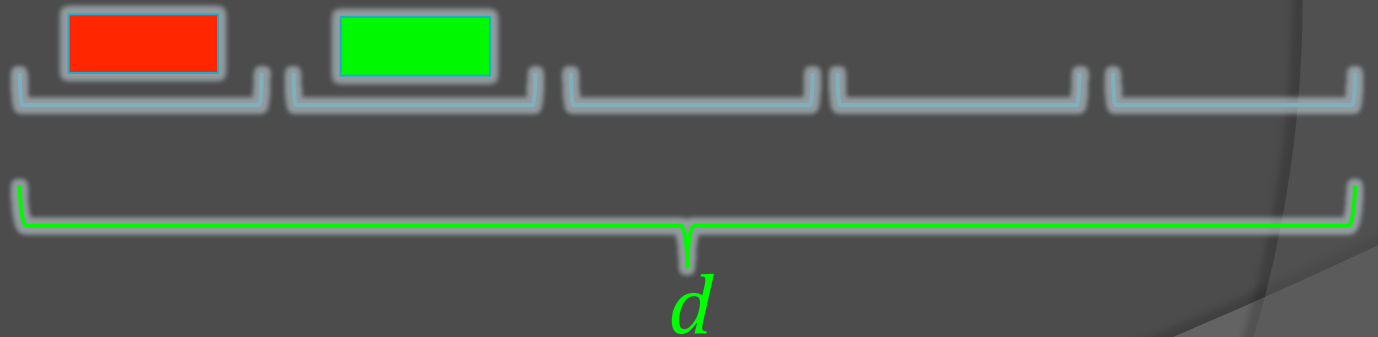


# Frequent Items

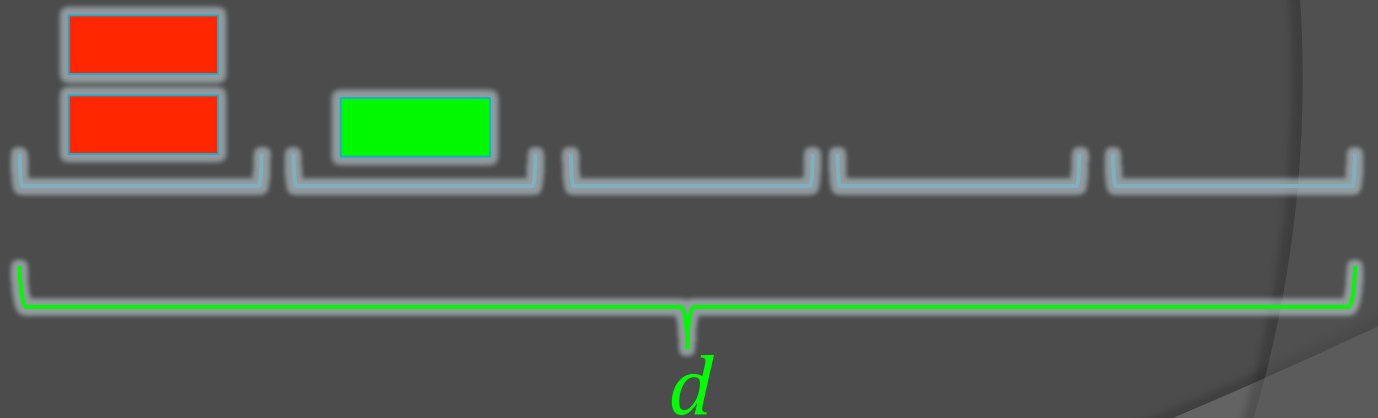




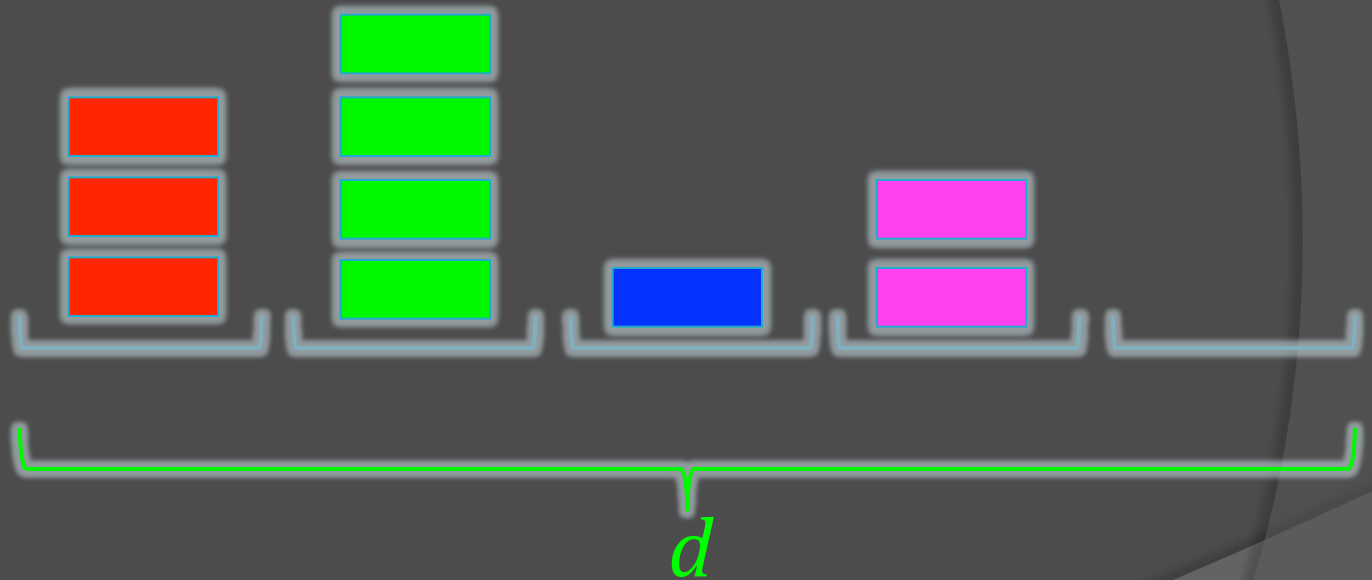
# Frequent Items



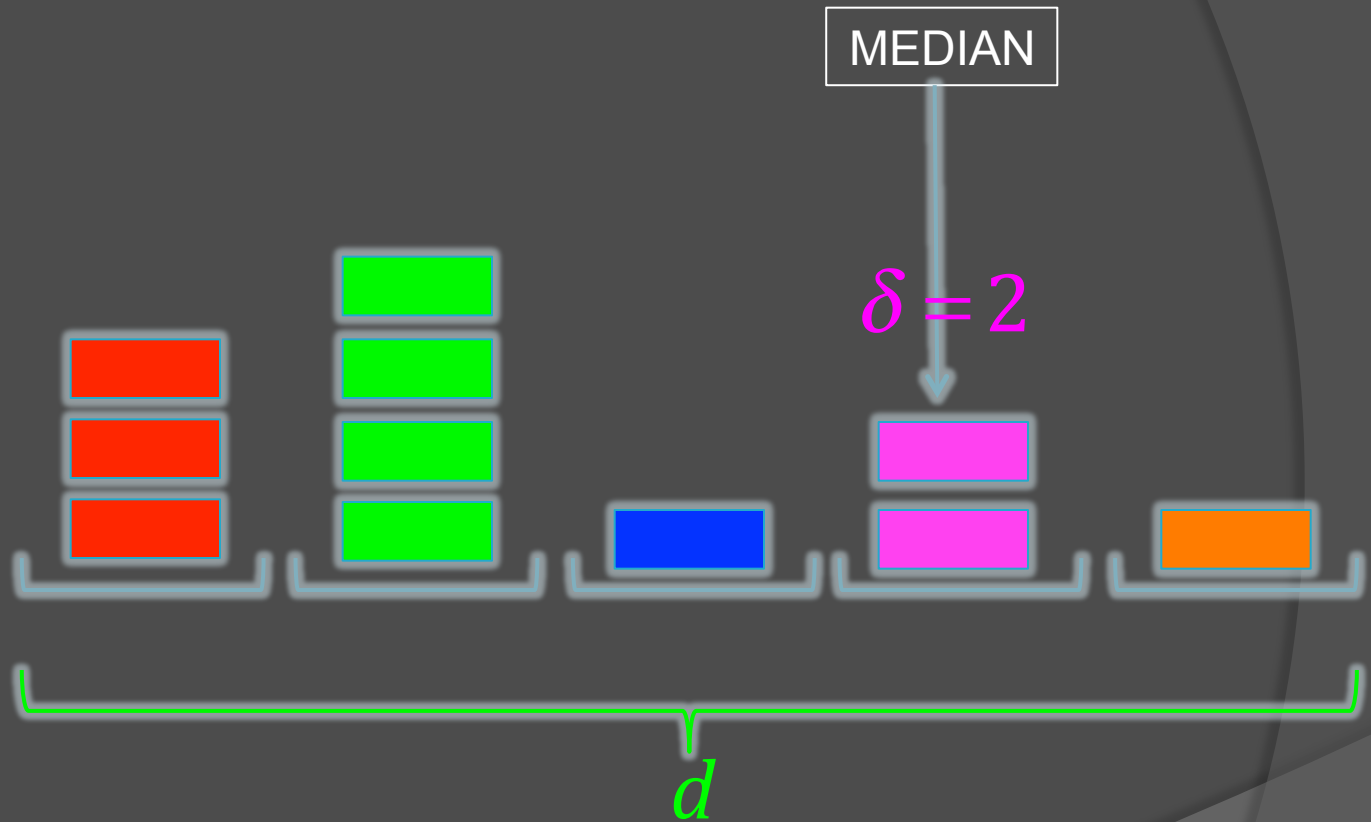
# Frequent Items



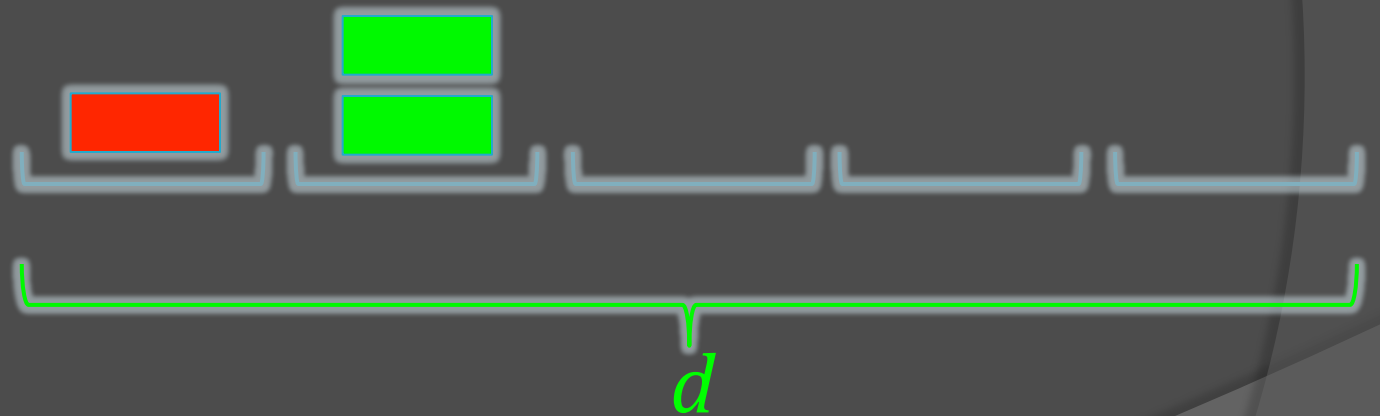
# Frequent Items



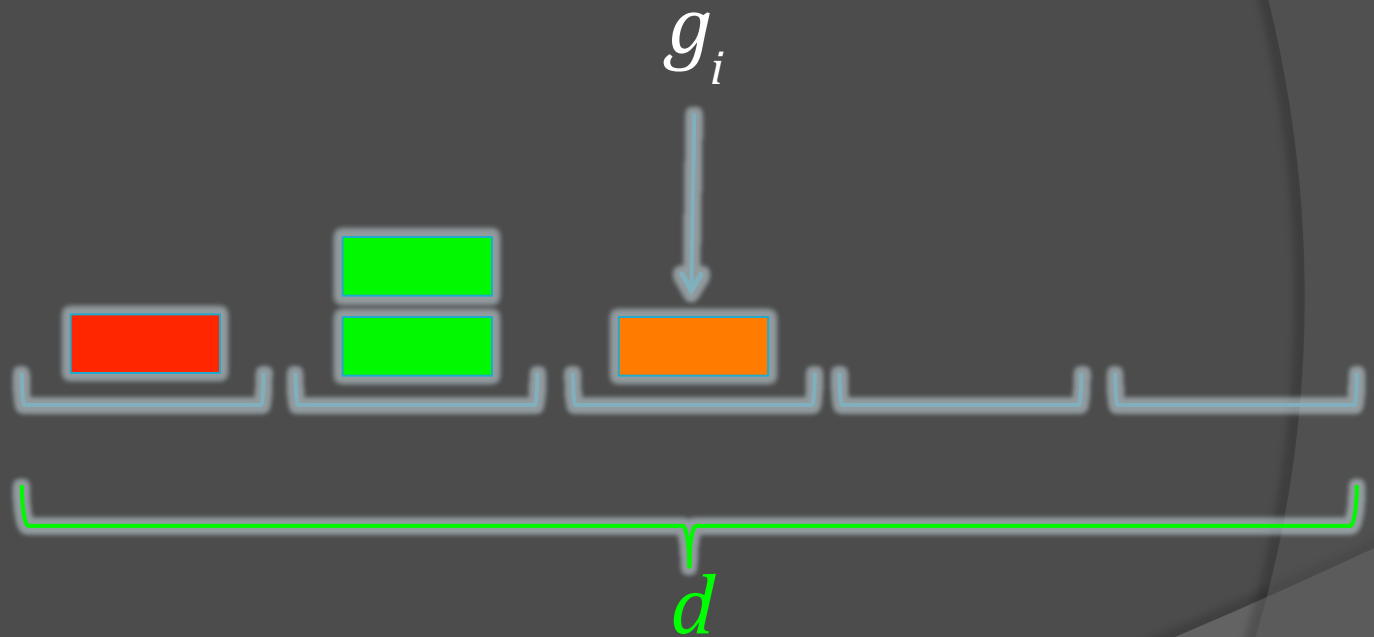
# Frequent Items



# Frequent Items



# Frequent Items



# Frequent Items – Observations

- We always get an undercount  $g_i \leq f_i$
- If we let  $\delta_t$  be the amount we decrease counter at time  $t$  then

$$g_i \geq f_i - \sum_t \delta_t$$

- Sum up the undercounts

$$0 \leq \sum_{i=1}^d g_i \leq \sum_{t=1}^n \left( 1 - \frac{d}{2} \delta_t \right)$$

# Frequent Items – Observations

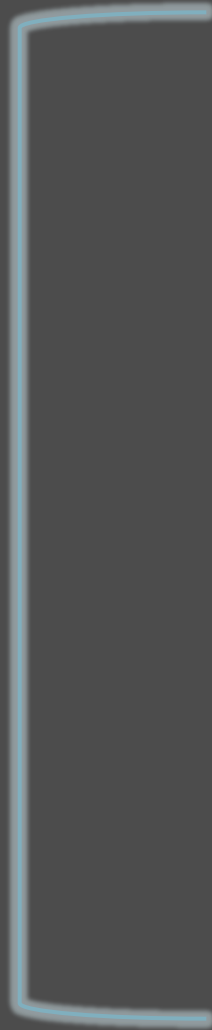
⊙ Thus, we get  $\sum_t \delta_t \leq 2n / d$

⊙ Set  $d = 2 / \varepsilon$ :

$$|f_i - g_i| \leq \varepsilon n$$



# Frequent Directions



d x m

# Frequent Directions

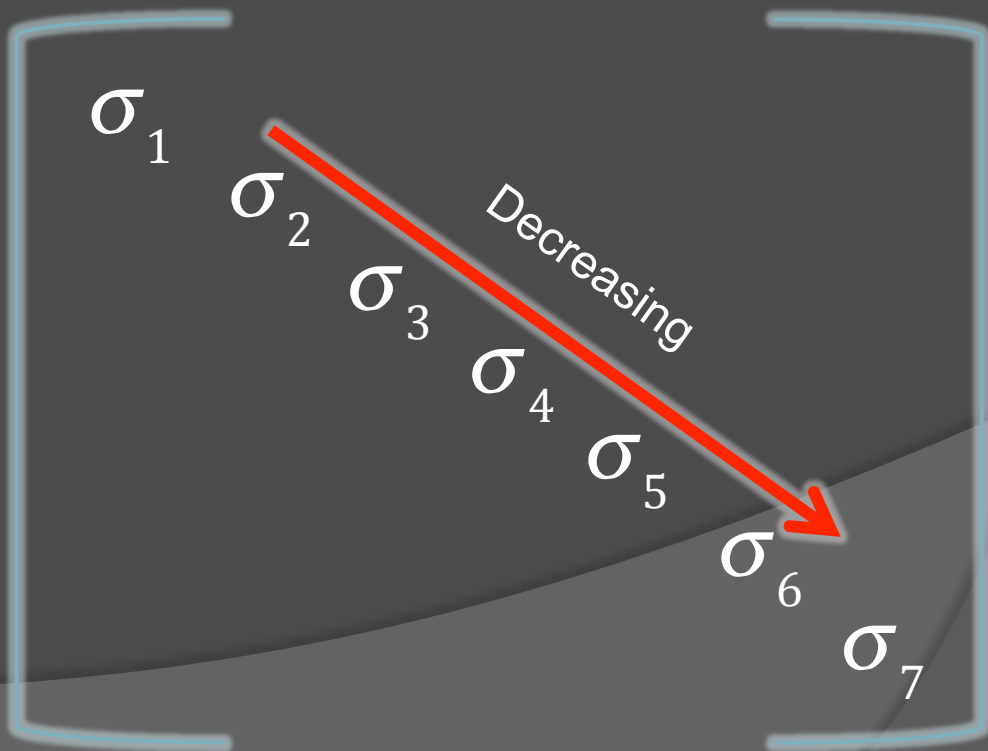
We now need  
to zero out  
some rows to  
make room for  
more!



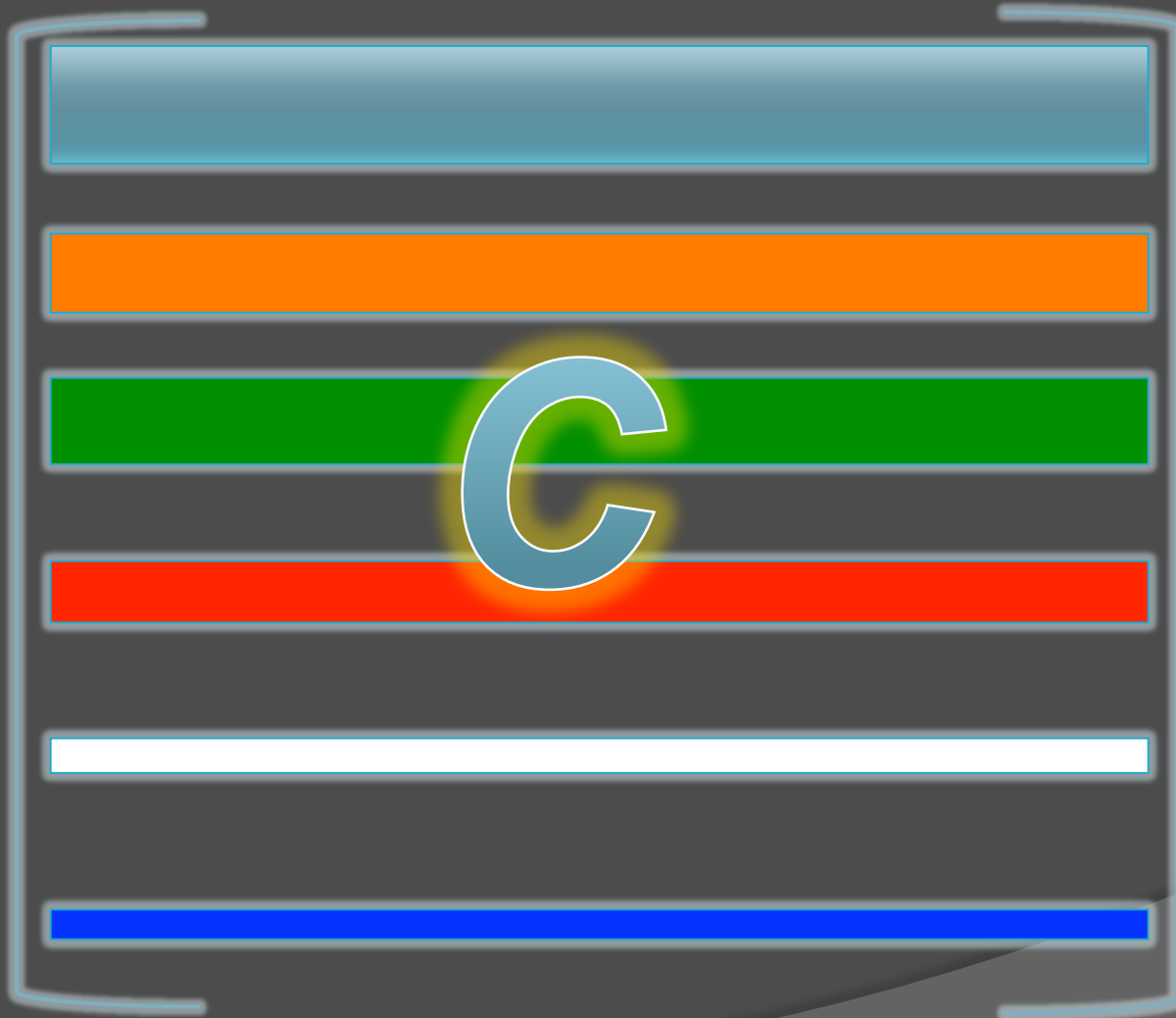
# Frequent Directions

Find SVD of  $B$ :  $U_{d \times d} \Sigma_{d \times m} V_{m \times m}^T = B$

$$\Sigma =$$



# Frequent Directions



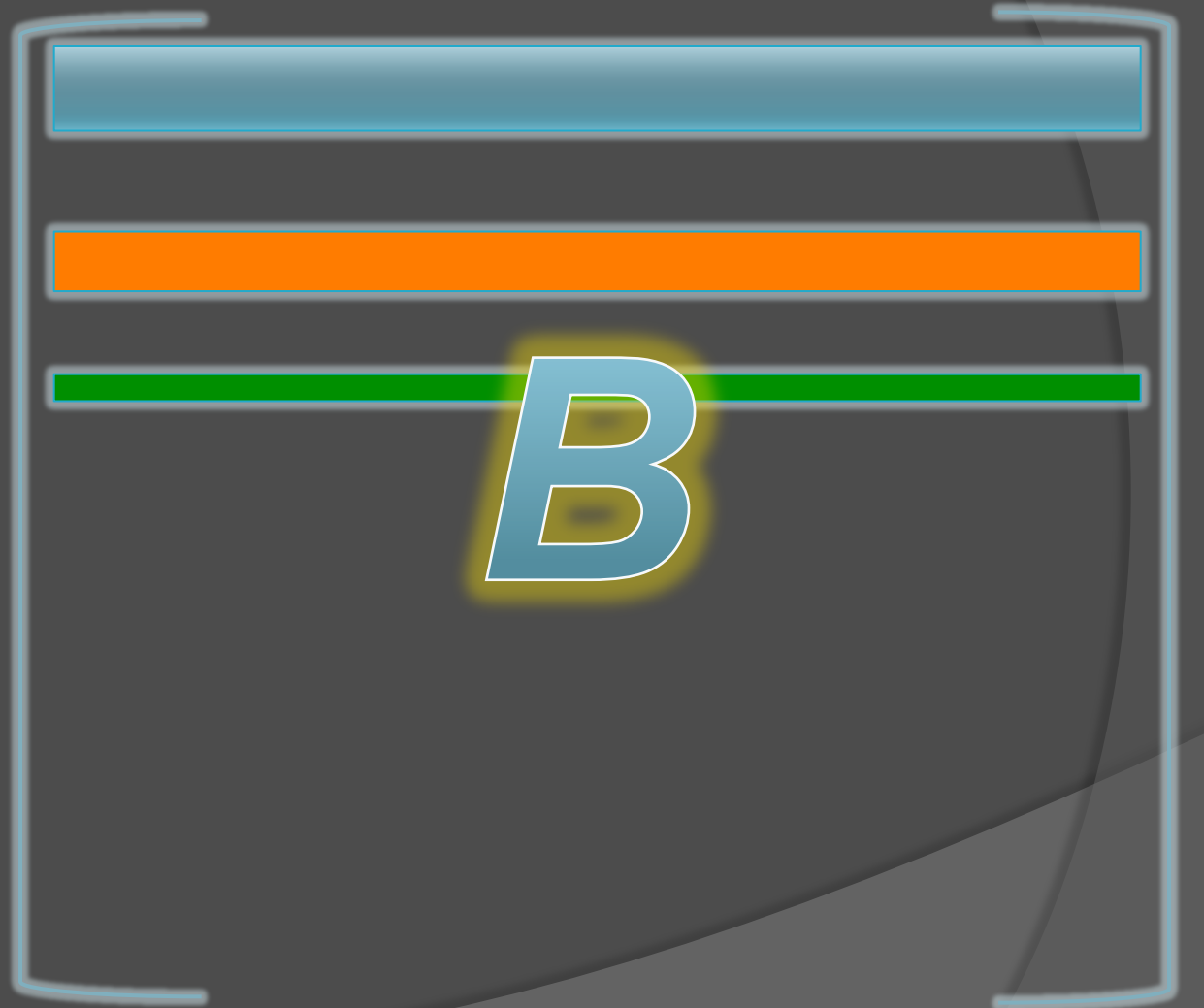
# Frequent Directions

$$\Sigma = \begin{bmatrix} \sigma_1 & & & & & & \\ & \sigma_2 & & & & & \\ & & \sigma_3 & & & & \\ & & & \sigma_4 & & & \\ & & & \uparrow & & & \\ & & & \sigma_{d/2} = \sqrt{\delta} & & & \\ & & & & \sigma_5 & & \\ & & & & & \sigma_6 & \\ & & & & & & \sigma_7 \end{bmatrix}$$

$$\hat{\Sigma} = \sqrt{\max(\Sigma^2 - \delta I, 0)}$$

# Frequent Directions

$$B = \hat{\Sigma} V^T$$



$d \times m$

# Frequent Directions

---

**Algorithm 1** *Frequent-directions*

---

**Input:**  $\ell, A \in \mathbb{R}^{n \times m}$

$B \leftarrow$  all zeros matrix  $\in \mathbb{R}^{\ell \times m}$

**for**  $i \in [n]$  **do**

    Insert  $A_i$  into a zero valued row of  $B$

**if**  $B$  has no zero valued rows **then**

$[U, \Sigma, V] \leftarrow \text{SVD}(B)$

$C \leftarrow \Sigma V^T$       # Only needed for proof notation

$\delta \leftarrow \sigma_{\ell/2}^2$

$\tilde{\Sigma} \leftarrow \sqrt{\max(\Sigma^2 - I_{\ell}\delta, 0)}$

$B \leftarrow \tilde{\Sigma} V^T$  # At least half the rows of  $B$  are all zero

**end if**

**end for**

**Return:**  $B$

---

# Analysis – Claim 1

- $B^TB$ ,  $A^TA$ ,  $A^TA - B^TB$  are all P.S.D.
- Proof: Check

$$\|Ax\|_2^2 - \|Bx\|_2^2 \geq 0$$



# Analysis – Claim 2

- With sketch  $B$  of size  $d$  from Frequent Directions we have

$$||A^T A - B^T B|| \leq 2 ||A||_f^2 / d$$

- Proof: First prove that for any unit vector  $x$

$$||Ax||^2 - ||Bx||^2 \leq 2/d \left( ||A||_f^2 - ||B||_f^2 \right)$$

# Analysis – Proof Continued

- Now we must show that for the largest e-vector  $x$  that

$$||A^T A - B^T B|| = ||Ax||^2 - ||Bx||^2$$

# Run Time

- SVD of an  $d \times m$  matrix of rank  $r$  takes

$$O(dmr) = O(d^2m)$$

- SVD is done once every  $d/2$  rows
- When SVD is not done, it takes time

$$O(m)$$

- Total run time:

$$O(dnm)$$

# Parallelization

If we have  $A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$  and  $B_i = FD(A_i)$   
then

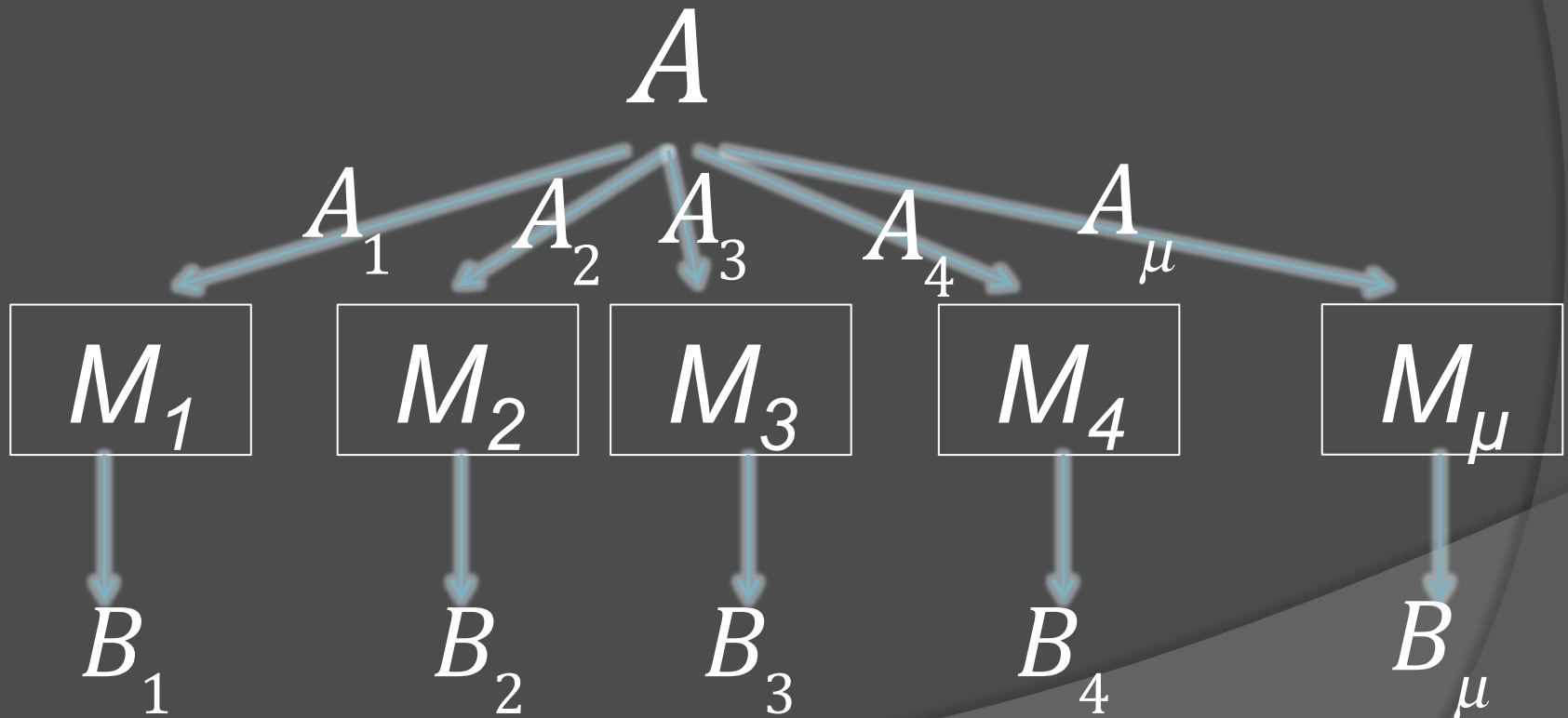
$$\|A^T A - D^T D\| \leq 2 \|A\|_f^2 / d$$

where

$$D = FD \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}$$

# Parallelization

- Let there be  $\mu$  machines and each takes  $n/\mu$  many rows



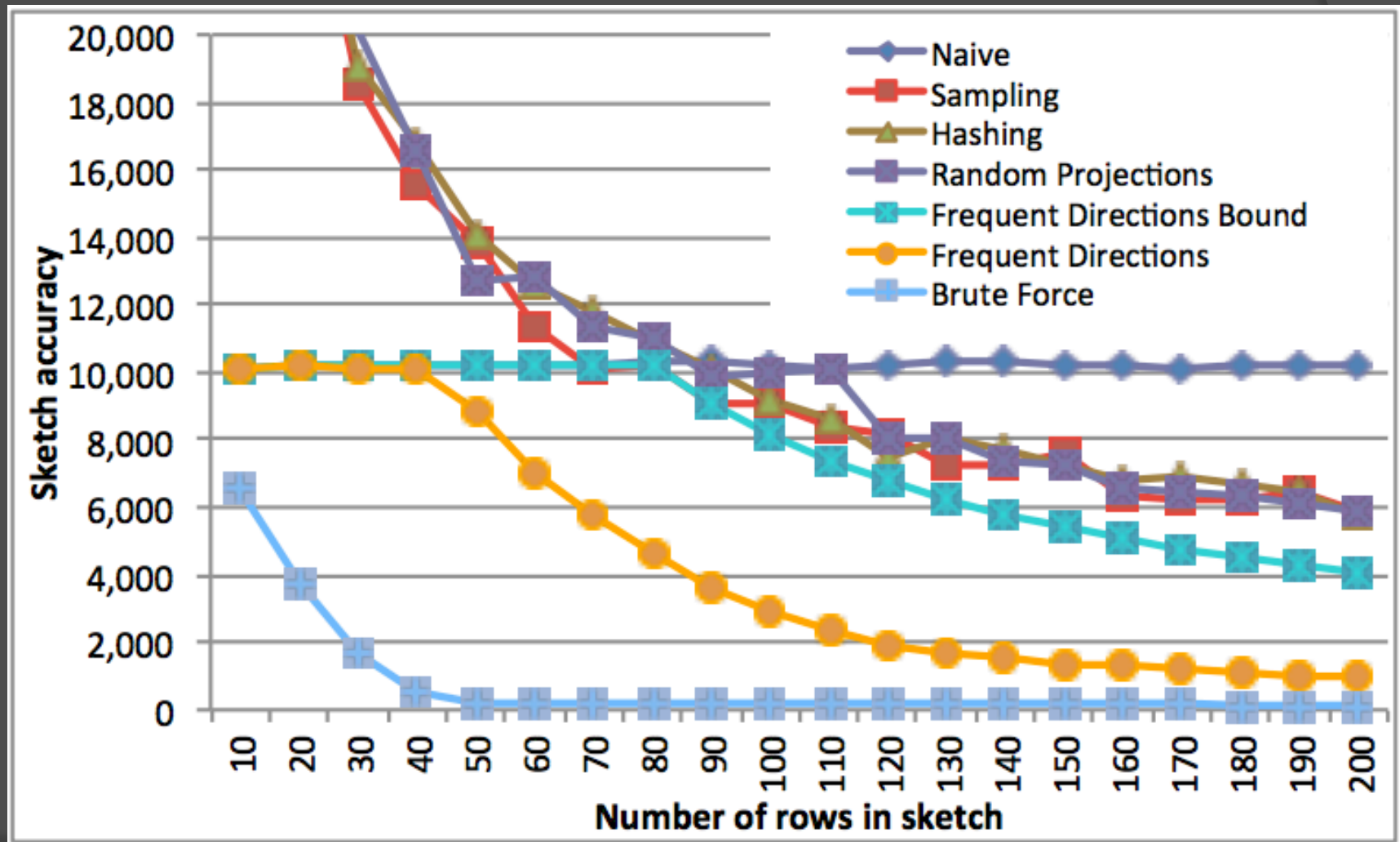
# Parallelization

- Each  $B_i$  has dimension  $d \times m$ .
- Each  $M_i$  took time  $O(dmn / \mu)$
- To then combine the others, can take  $\mu$  more machines, and total run time

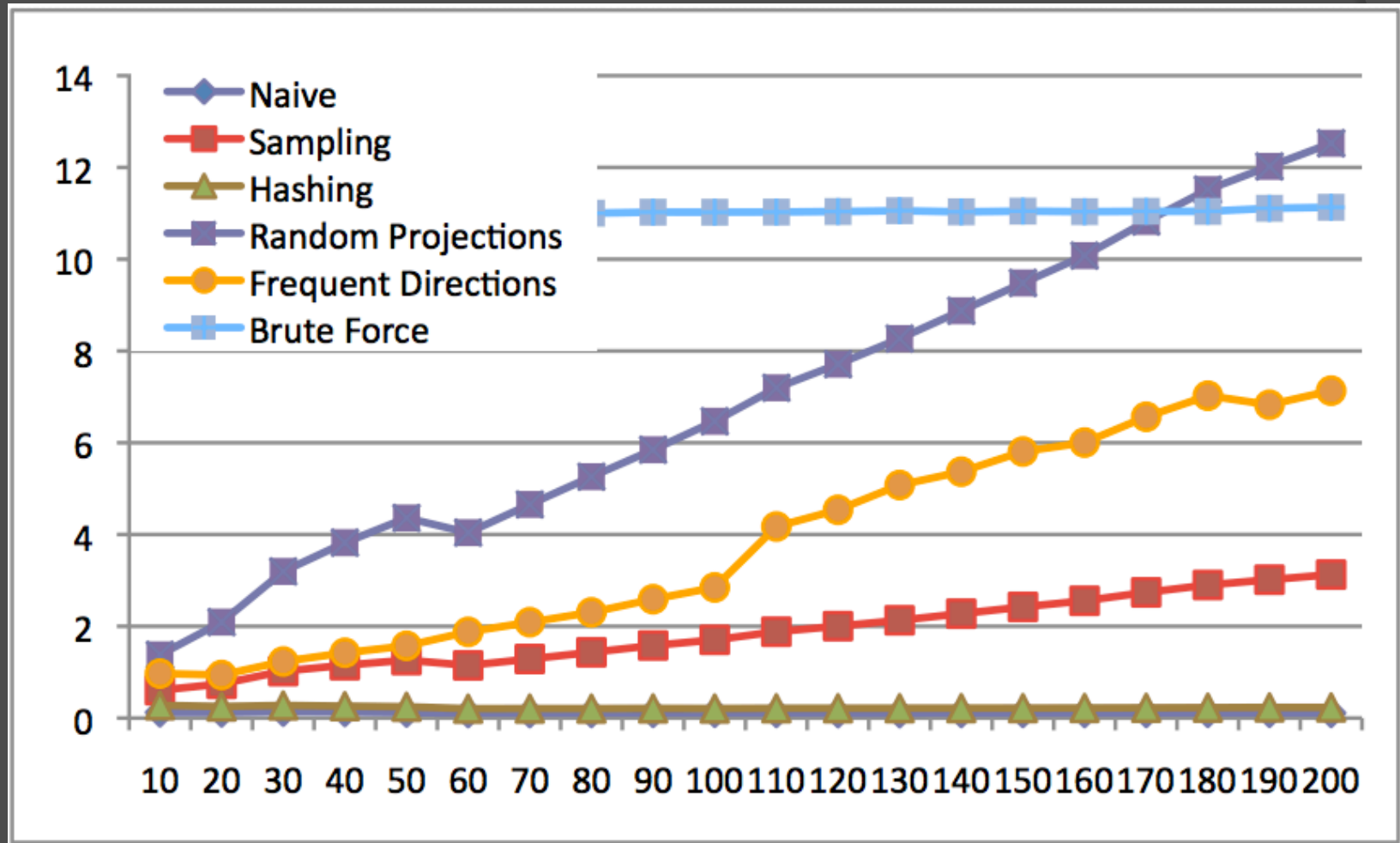
$$O(dmn / \mu + \log(\mu)d^2m)$$

- Set  $\mu = \Theta\left(\frac{n}{d}\right) = \Theta(\varepsilon n) \Rightarrow$  run time  $O\left(\frac{m \log(n)}{\varepsilon^2}\right)$

# Results – Accuracy

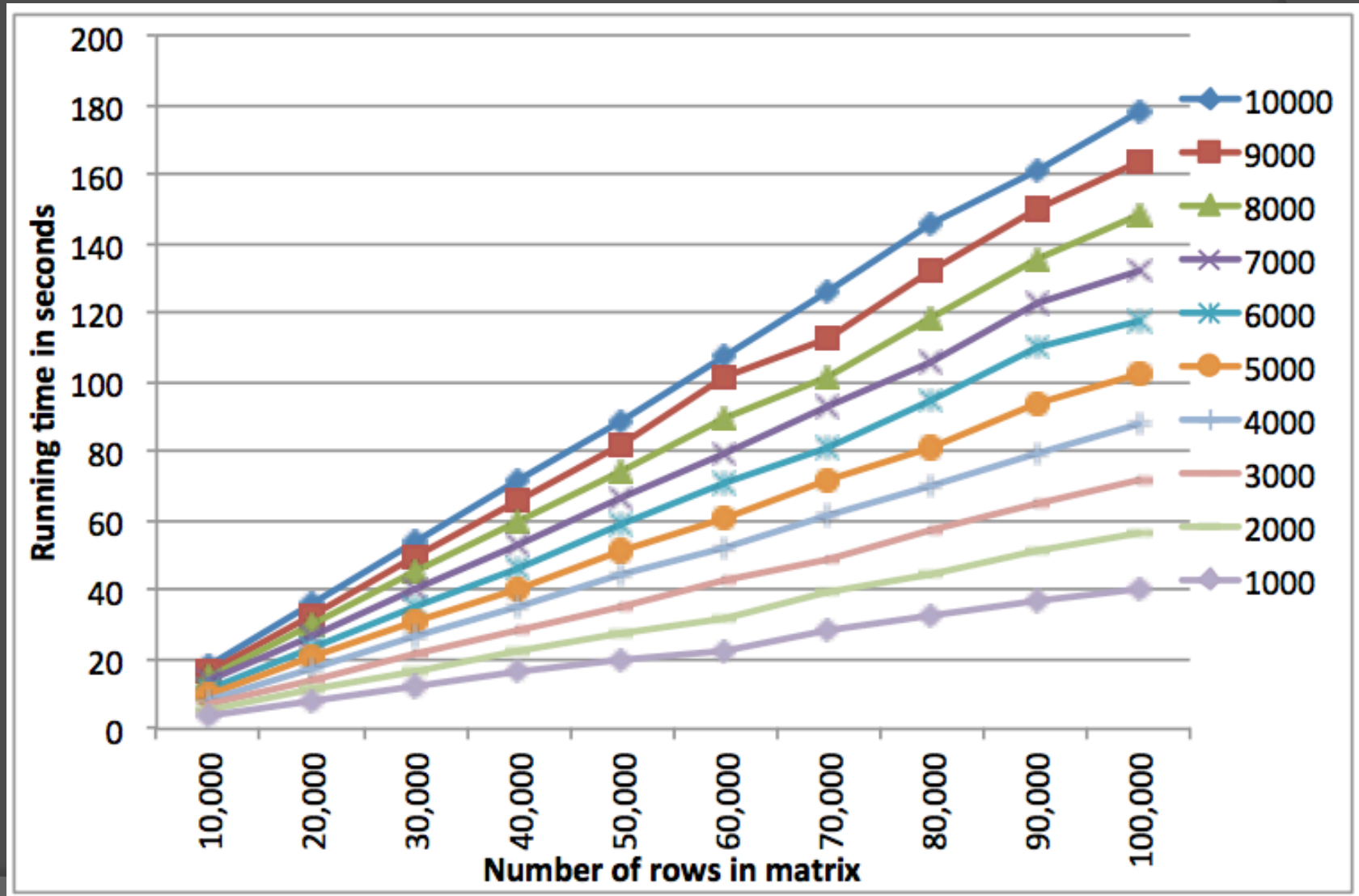


# Results – Run Time vs. Others





# Results – Run Time for FD



The