

Performance of Random Sampling for Low-rank Approximation on GPUs

Théo Mary,

October 1, 2014

Context

- Application with high computational and storage costs.
 - e.g., solve $Ax = b$
- We can use matrix compression techniques to save computations and storage.
 - Because of redundant information;
 - Because we don't need ε_{mach} accuracy;

Context

- Application with high computational and storage costs.
 - e.g., solve $Ax = b$
- We can use matrix compression techniques to save computations and storage.
 - Because of redundant information;
 - Because we don't need ε_{mach} accuracy;

Mathematical Formulation

$$\begin{array}{ccc} A & \approx & BC \\ m \times n & & m \times k \quad k \times n \end{array}$$

- $\|E\| = \|A - BC\| \leq \sigma_{k+1} \leq \varepsilon$
- The ε threshold is to be set by the user, according to their needs.

- Pivoted QR decomposition

$$AP = (Q_1 \quad Q_2) \begin{pmatrix} R_{11} & R_{12} \\ & R_{22} \end{pmatrix}$$

with

- $Q = (Q_1 \quad Q_2)$ an $m \times n$ orthogonal matrix;
 - $R = \begin{pmatrix} R_{11} & R_{12} \\ & R_{22} \end{pmatrix}$ an $n \times n$ upper triangular matrix;
 - P a $n \times n$ pivot matrix.
- Truncated Pivoted QR Decomposition

$$\begin{matrix} AP & \approx & Q_1 & \begin{pmatrix} R_{11} & R_{12} \end{pmatrix} \\ m \times n & & m \times k & \begin{matrix} k \times n \end{matrix} \end{matrix}$$

- Computes QR with column pivoting using BLAS-3.
- To compute truncated version → one line to change in the code.
- Limitations:
 - Not only BLAS-3: also BLAS-2;
 - Synchronization at every step to pick pivot;
 - Limited parallelism and data locality;
 - Communications.

- **Stage A:** generate Q , orthogonal subspace spanning the range of A , i.e.:

$$A \approx AQ^T Q$$

- **Stage B:** use Q to compute low-rank approximations of A (QR, SVD, ...) with standard deterministic methods.

for $i = 1, 2, \dots, k$ **do**

 Draw random vector ω_i .

 Form the product $b_i = \omega_i A$.

end for

- $B = \{b_1, \dots, b_k\}$ is **probably** linearly independent
 $\Rightarrow Q = \text{orth}(B)$.

```
for  $i = 1, 2, \dots, k + p$  do  
    Draw random vector  $\omega_i$ .  
    Form the product  $b_i = \omega_i A$ .  
end for
```

- $B = \{b_1, \dots, b_{k+p}\}$ is **probably** linearly independent
 $\Rightarrow Q = \text{orth}(B)$.
- p is the oversampling.

- $\ell = k + p$

$$\begin{array}{c} B \\ \ell \times n \end{array} = \begin{array}{c} \Omega \\ \ell \times m \end{array} \begin{array}{c} A \\ m \times n \end{array}$$

- $\ell = k + p$

$$\begin{matrix} B & = & \Omega & A \\ \ell \times n & & \ell \times m & m \times n \end{matrix}$$

- How do we generate Ω ?
 - Gaussian
 - FFT

- $\|A - AQ^T Q\| \leq C(\Omega, \rho)\sigma_{k+1}$
 - If $\{\sigma_i\}_{i=1,n}$ decay slowly, $\|A - AQ^T Q\|$ can be big.
- $B = \Omega A (A^T A)^q$
- $\|A - AQ^T Q\| \leq C(\Omega, \rho)^{1/(2q+1)}\sigma_{k+1}$

- $\|A - AQ^T Q\| \leq C(\Omega, \rho)\sigma_{k+1}$
 - If $\{\sigma_i\}_{i=1,n}$ decay slowly, $\|A - AQ^T Q\|$ can be big.
- $B = \Omega A (A^T A)^q$
- $\|A - AQ^T Q\| \leq C(\Omega, \rho)^{1/(2q+1)}\sigma_{k+1}$
- Round-off errors \Rightarrow need to reorthogonalize.

$$B_0 = \Omega A$$

repeat q times:

$$\begin{aligned} \acute{Q}_0 &= \text{orth}(B_0) & ; & \acute{B}_1 = \acute{Q}_0 A^T \\ \acute{Q}_1 &= \text{orth}(\acute{B}_1) & ; & B_1 = \acute{Q}_1 A \end{aligned}$$

- Truncated PQR step:

- Compute pivot P using sampled matrix B :

$$BP \approx \widehat{Q}_k \begin{pmatrix} \widehat{R}_{1:k} & \widehat{R}_{k+1:n} \end{pmatrix} = BP_{1:k} \begin{pmatrix} I_k \widehat{R}_{1:k}^{-1} \widehat{R}_{k+1:n} \end{pmatrix}$$

- Use pivot P on full matrix A :

$$\Rightarrow AP \approx AP_{1:k} \begin{pmatrix} I_k & \widehat{R}_{1:k}^{-1} \widehat{R}_{k+1:n} \end{pmatrix}$$

- QR step:

$$AP_{1:k} = Q\bar{R}$$

- Final approximation:











$$\begin{array}{ccc}
 AP & \approx & Q \quad \bar{R} \begin{pmatrix} I_k & \widehat{R}_{1:k}^{-1} \widehat{R}_{k+1:n} \end{pmatrix} \\
 m \times n & & m \times k \quad \quad k \times n \\
 & & \mathbf{Q} \quad \quad \mathbf{R}
 \end{array}$$

- Two memory levels hierarchy: fast/slow (M = size of fast memory).

	#flops	#words
Sampling	$\mathcal{O}(mnl)$	$\mathcal{O}(mnl/M^{1/2})$
Iter. (mult.)	$\mathcal{O}(mnlq)$	$\mathcal{O}(mnlq/M^{1/2})$
Iter. (orth.)	$\mathcal{O}((m+n)\ell^2)$	$\mathcal{O}((m+n)\ell^2/M^{1/2})$
TruncPQR	$\mathcal{O}(nl^2)$	$\mathcal{O}(nl^2)$
QR	$\mathcal{O}(ml^2)$	$\mathcal{O}(ml^2/M^{1/2})$
Total	$\mathcal{O}(mnl(1+q))$	$\mathcal{O}(mnl(1+q)/M^{1/2})$
TruncQP3	$\mathcal{O}(mnk)$	$\mathcal{O}(mnk)$

- Under assumption that p, q are constant, Randomized PQR converges towards communications lower bound.

- Needed for power method and for QR($AP_{1:k}$) step.

	Stability		Performance
Householder QR		ε	
Cholesky QR		$\kappa(A)^2 \varepsilon$	
CA QR		ε	
CGS		$\kappa(A)^2 \varepsilon$	
MGS		$\kappa(A) \varepsilon$	

- Cholesky QR Algorithm:

1. Form $S = X^T X$.
2. Compute Cholesky factorization $R = \text{chol}(S)$.
3. Solve $Q = XR^{-1}$.

- Hardware: Bunsen (16 threads running on 16 cores Genuine Intel(R) CPU @ 2.60GHz + 3 GPUs Tesla K40c)
- Software: MAGMA.
- Test matrices:
 - power: $A = X\Sigma Y$, with $\sigma_i = -i^\alpha$, ($\alpha = 3$).
 - exponent: $A = X\Sigma Y$, with $\sigma_i = 10^{-i\gamma}$ ($\gamma = 0.1$).
 - hapmap.

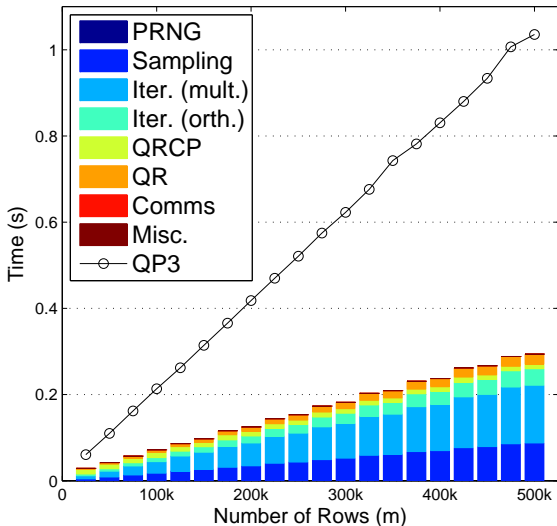
	power	exponent	hapmap
m	500,000	500,000	503,783
n	500	500	506
k	50	50	50
p	10	10	10
ℓ	60	60	60
σ_1	1	1	9.9e+03
σ_{k+1}	8e-06	1.3e-05	5e+02
$\kappa(A)$	1.3e+05	7.9e+04	2e+01

- Approximation error $\|AP - QR\|/\|A\|$

	QP3	Rand $q = 0$	Rand $q = 1$	Rand $q = 2$
Power	4.4679e-05	9.0784e-05	4.5948e-05	4.4489e-05
Exponent	2.6892e-05	5.1795e-05	2.6899e-05	2.6898e-05
Hapmap	5.9887e-01	9.8563e-01	8.7363e-01	8.1816e-01

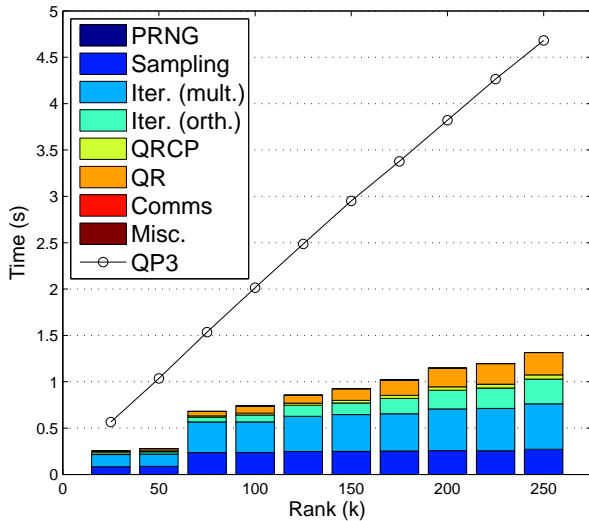
- Usefulness of oversampling: roughly an order of magnitude.

Time with increasing number of rows



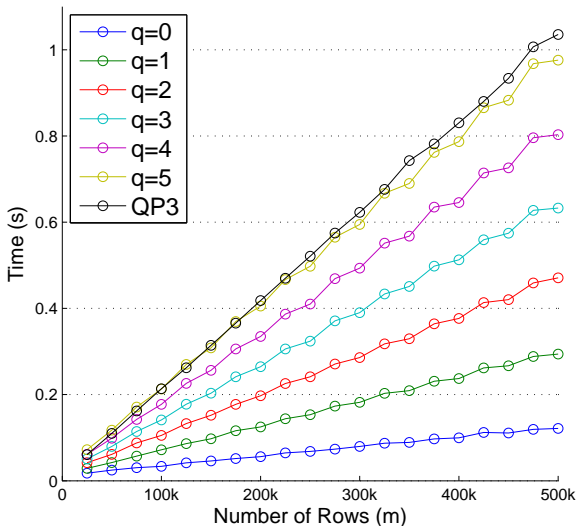
$(m = 25,000 \rightarrow 500,000; n = 500; \ell = k + p = 50 + 10; q = 1)$

Time with increasing rank



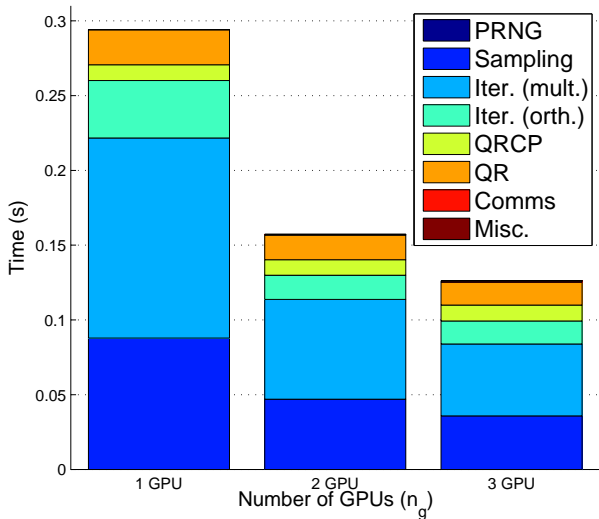
$(m = 500,000; n = 500; \ell = k + p; k = 25 \rightarrow 250; p = 10; q = 1)$

Time with increasing number of iterations



$(m = 25,000 \rightarrow 500,000; n = 500; \ell = k + p = 50 + 10; q = 0 \rightarrow 5)$

Time on 1,2,3 GPUs



($m = 500,000$; $n = 500$; $\ell = k + p = 50 + 10$; $q = 1$)

Summary

- Substitute QP3 with Randomized, which is dominated by matrix-matrix multiplications \Rightarrow very nice properties: data locality, higher parallelism, no synchronizations, minimized communications.
- Comparable accuracy on the test matrices used.
- Performance speedups above 8 and scaling on multiple GPUs.
- For the future:
 - Use other error measurements (applications e.g. clustering).
 - Orthogonalization studies.

PhD

- Using low-rank approximations to improve multifrontal solvers.

Thank you!
Questions?