

Efficient Algorithms for Evaluating Matrix Polynomials

Sivan Toledo and Niv Hoffman
Blavatnik School of Computer Science
Tel-Aviv University

Oded Schwartz
Hebrew University

Efficiency in Evaluation of

$$q(\mathbf{A}) = c_0\mathbf{I} + c_1\mathbf{A} + c_2\mathbf{A}^2 + \dots + c_d\mathbf{A}^d$$

\mathbf{A} is an n -by- n real or complex matrix

Given $\mathbf{A}^\ell, \mathbf{A}^k,$

- $c_k\mathbf{A}^k$ requires n^2 scalar multiplications, but
- $\mathbf{A}^{\ell+k} = \mathbf{A}^\ell\mathbf{A}^k$ requires n^3 (or $n^{2.81}$ with Strassen)

Minimize the number of matrix-matrix multiplications required to evaluate $q(\mathbf{A})$; these are called in the literature *nonscalar* multiplications

Maybe other ways to reduce arithmetic

Reduce communication (cache misses in this talk)

A Side Note on Functions of Matrices

When A is normal or close to normal, we can compute $q(A)$ or any other function $f(A)$ using an eigendecomposition,

$$\begin{aligned}q(A) &= q(V\Lambda V^*) \\ &= Vq(\Lambda)V^* \\ &= V \begin{pmatrix} q(\lambda_1) & & \\ & \ddots & \\ & & q(\lambda_n) \end{pmatrix} V^* .\end{aligned}$$

It is mostly when A is far from normal when these problems become interesting

Building Blocks

$q(A)$ via Explicit Powers

Allocate three matrices, A , A^k , and Q

Initialize $Q = 0$, and $A^k = I$

for $k \leftarrow 0, \dots, d$, add $Q \leftarrow Q + c_k A^k$ and multiply
 $A^{k+1} \leftarrow A A^k$

In step d we produce $q(A) = Q$

A total of $d - 1$ matrix-matrix multiplications

$q(\mathbf{A})$ via Horner's Rule

Allocate two matrices, \mathbf{A} and \mathbf{Q}

Initialize $\mathbf{Q} = c_{d-1}\mathbf{I} + c_d\mathbf{A}$

In step $k \leftarrow d - 2, \dots, 1, 0$, multiply and add $\mathbf{Q} \leftarrow c_k\mathbf{I} + \mathbf{Q}\mathbf{A}$

In step 0 we produce $q(\mathbf{A}) = \mathbf{Q}$

Again a total of $d - 1$ matrix-matrix multiplications

$q(\mathbf{A})$ via Its Roots

If the roots ξ_1, \dots, ξ_d are given, not the coefficients c , then

Allocate two matrices, \mathbf{A} and \mathbf{Q}

Initialize $\mathbf{Q} = \mathbf{A} - \xi_1 \mathbf{I}$

For $k \leftarrow 2, \dots, d$ add and multiply $\mathbf{Q} \leftarrow \mathbf{Q}(\mathbf{A} - \xi_k \mathbf{I})$

In step d we produce $q(\mathbf{A}) = \mathbf{Q}$

Again a total of $d - 1$ matrix-matrix multiplications

(Computing the roots from the coefficients is not advised;
often ill conditioned)

Clearly, Matrix Multiplication is Important Here

n^3 multiplications ($\approx 2n^3$ arithmetic operations) using native algorithms

$O(n^\omega)$ using so-called *fast methods*; e.g., $O(n^{2.81})$ for Strassen-Winograd

$O(n^2 k^{\omega-2})$ to multiply an n -by- n matrix by an n -by- k (block mat-vec); no gain for $k = 1$!

$\times 6$ arithmetic speedup for triangular mat mult (plus can apply Strassen to square blocks within the recursion)

Does it Really Take $d - 1$ Matrix Multiplications?

Classical/Naive methods (e.g. Horner) require $d - 1$.

Is this really necessary?

d/2 Mat-Mults: Rabin-Winograd/Patterson-Stockmeyer

No (Rabin-Winograd 1971, Patterson-Stockmeyer 1973).

$$\begin{aligned}q(\mathbf{A}) &= c_0\mathbf{I} + c_1\mathbf{A} + c_2\mathbf{A}^2 + \cdots + c_{2r-1}\mathbf{A}^{2r-1} \\ &= c_0\mathbf{I} + \cdots + c_{r-1}\mathbf{A}^{r-1} + \mathbf{A}^r(\tilde{c}_0\mathbf{I} + \cdots + \tilde{c}_{r-1}\mathbf{A}^{r-1})\end{aligned}$$

use recursion and denote # mat-mults by $N(d)$ (repeated squaring for \mathbf{A}^r). For powers of two,

$$N(2r - 1) = 2N(r - 1) + 1$$

$$\text{total (inc. repeated squaring)} \approx d/2 + \log(d)$$

for non powers of two, split into powers of two so

$$N(d) \leq d/2 + 2\log(d)$$

The Algorithm of Paterson & Stockmeyer

Consider an example,

$$\begin{aligned}q(A) &= 2I + 3A + 4A^2 \\ &\quad + 5A^3 + 6A^4 + 7A^5 \\ &\quad + 8A^6 + 9A^7 + 2A^8 \\ &= (2I + 3A + 4A^2) \\ &\quad + (5A^3 + 6A^4 + 7A^5) \\ &\quad + (8A^6 + 9A^7 + 2A^8) \\ &= (2I + 3A + 4A^2) I \\ &\quad + (5I + 6A + 7A^2) A^3 \\ &\quad + (8I + 9A + 2A^2) (A^3)^2\end{aligned}$$

Analysis of the Example

Consider an example,

$$\begin{aligned}q(A) &= (2I + 3A + 4A^2) I \\ &\quad + (5I + 6A + 7A^2) A^3 \\ &\quad + (8I + 9A + 2A^2) (A^3)^2\end{aligned}$$

We perform 1 matmult to produce I, A, A^2 , 1 to produce A^3 , 1 square it, and 2 more to multiply degree-2 polynomials by powers of A^3 .

A polynomial in A^3 whose coefficients are quadratics in A .

The General Case of the PS Algorithm

Let $ps = d$ (remainder is easy to handle),

$$\begin{aligned}q(\mathbf{A}) &= c_0\mathbf{I} + c_1\mathbf{A} + c_2\mathbf{A}^2 + \cdots + c_d\mathbf{A}^d \\ &= c_0\mathbf{I} + c_1\mathbf{A} + \cdots + c_{p-1}\mathbf{A}^{p-1} \\ &\quad + (c_p\mathbf{I} + c_{p+1}\mathbf{A} + \cdots + c_{2p-1}\mathbf{A}^{p-1}) \mathbf{A}^p \\ &\quad + \cdots \\ &\quad + (c_{d-p+1}\mathbf{I} + c_{d-p+2}\mathbf{A} + \cdots + c_d\mathbf{A}^{p-1}) (\mathbf{A}^p)^{s-1} \\ &= c_0\mathbf{I} + c_1\mathbf{A} + \cdots + c_{p-1}\mathbf{A}^{p-1} \\ &\quad + (c_p\mathbf{I} + c_{p+1}\mathbf{A} + \cdots + c_{2p-1}\mathbf{A}^{p-1}) \mathbf{A}^p \\ &\quad + \cdots \\ &\quad + (c_{(s-1)p}\mathbf{I} + c_{(s-1)p+1}\mathbf{A} + \cdots + c_{(s-1)p+p-1}\mathbf{A}^{p-1}) (\mathbf{A}^p)^{s-1}\end{aligned}$$

Arithmetic and Memory Complexity of PS

Let $ps = d$ (remainder is easy to handle),

$$\begin{aligned}q(\mathbf{A}) &= (\mathbf{c}_0\mathbf{I} + \mathbf{c}_1\mathbf{A} + \cdots + \mathbf{c}_{p-1}\mathbf{A}^{p-1}) \\ &\quad + (\mathbf{c}_p\mathbf{I} + \mathbf{c}_{p+1}\mathbf{A} + \cdots + \mathbf{c}_{2p-1}\mathbf{A}^{p-1}) \mathbf{A}^p \\ &\quad + \cdots \\ &\quad + (\mathbf{c}_{(s-1)p}\mathbf{I} + \mathbf{c}_{(s-1)p+1}\mathbf{A} + \cdots + \mathbf{c}_{(s-1)p+p-1}\mathbf{A}^{p-1}) (\mathbf{A}^p)^{s-1}\end{aligned}$$

Form and store $\mathbf{A}^2, \dots, \mathbf{A}^{p-1}, \mathbf{A}^p$ explicitly ($p-1$ MMs, $p+1$ matrices to store)

Set $\mathbf{Q} \leftarrow$ highest coefficient polynomial $= \sum_{\ell=0}^{p-1} \mathbf{c}_{\dots} \mathbf{A}^{\ell}$ ($p-1$ scale-add, 0 MMs)

For $k \leftarrow s-1, \dots, 1$, 0 multiply and add

$\mathbf{Q} \leftarrow \mathbf{Q}\mathbf{A}^p + \sum_{\ell=0}^{p-1} \mathbf{c}_{\dots} \mathbf{A}^{\ell}$ (s m-multiplications)

Parameter Optimization for PS

Total number of matrix multiplications is $p - 1 + s$ where $s = \lceil d/p \rceil - 1$

Therefore $ps \approx d$ so $p + s - 1$ is minimized near $p \approx \sqrt{d}$ at about $2\sqrt{d}$ MMs

Can we go even lower?

No. Paterson & Stockmeyer proved that there are qs for which # of MMs is at least \sqrt{d}

($\sqrt{d} - 1/2$ MMs if the coefficients are integers), under some reasonable assumptions

A Comparison

	Work (# MMs)	Memory (# words)
Explicit Powers	$d - 1$	$3n^2$
Horner's Rule	$d - 1$	$2n^2$
MM PS	$p - 1 + s$	$(p + 1) n^2$
MM PS for $ps \approx d$	$\approx 2\sqrt{d}$	$\approx \sqrt{d}n^2$

A work-storage tradeoff

Memory-Efficient Matrix-Vector PS (Van Loan 1979)

$$q(\mathbf{A}) = (2\mathbf{I} + 3\mathbf{A} + 4\mathbf{A}^2) \mathbf{I} + \cdots + (8\mathbf{I} + 9\mathbf{A} + 2\mathbf{A}^2) (\mathbf{A}^3)^2$$

So the j th column is

$$\begin{aligned} q(\mathbf{A})_{:,j} &= (\mathbf{I} (2\mathbf{I} + 3\mathbf{A} + 4\mathbf{A}^2))_{:,j} \\ &\quad + (\mathbf{A}^3 (5\mathbf{I} + 6\mathbf{A} + 7\mathbf{A}^2))_{:,j} \\ &\quad + \left((\mathbf{A}^3)^2 (8\mathbf{I} + 9\mathbf{A} + 2\mathbf{A}^2) \right)_{:,j} \\ &= \mathbf{I} (2\mathbf{I}e_j + 3\mathbf{A}e_j + 4\mathbf{A}^2e_j) \\ &\quad + \mathbf{A}^3 (5e_j + 6\mathbf{A}e_j + 7\mathbf{A}^2e_j) \\ &\quad + (\mathbf{A}^3)^2 (8e_j + 9\mathbf{A}e_j + 2\mathbf{A}^2e_j) \end{aligned}$$

Matrix-Vector SPH (Van Loan 1979)

$$\begin{aligned}q(A)_{:,j} &= I(2Ie_j + 3Ae_j + 4A^2e_j) \\ &\quad + A^3(5e_j + 6Ae_j + 7A^2e_j) \\ &\quad + (A^3)^2(8e_j + 9Ae_j + 2A^2e_j)\end{aligned}$$

Compute A^p ($\log_2 p$ m-multiplications, can use Strassen)

For $j \leftarrow 1, \dots, n$

Compute $Ae_j, \dots, A^{p-1}e_j$ ($p-1$ matvecs)

Set $Q_j \leftarrow \sum_{\ell=0}^{p-1} c_{\ell} A^{\ell} e_j$ (vec ops)

For $k \leftarrow s-1, \dots, 1, 0$ multiply and add

$Q_j \leftarrow A^p Q_j + \sum_{\ell=0}^{p-1} c_{\ell} A^{\ell} e_j$ (s MVs, vec ops)

$n(p-1) + ns$ matvecs $\equiv s + p - 1$ (conventional, no Strassen) MMs

A Comparison (now with Van Loan's schedule)

	Work (# MMs)	Memory (# words)
Explicit Powers	$d - 1$	$3n^2$
Horner's Rule	$d - 1$	$2n^2$
MM PS	$p - 1 + s$	$(p + 1) n^2$
MV PS	$\log_2 p + (s + p)_{\text{conv}}$	$3n^2 + pn$

Communication (Cache-Miss) Analysis

Communications Analysis in Two-Level Memories

We assume that memory consists of a large slow memory and a fast memory (cache) of size M , and count the number of accesses to slow memory (cache misses) under an ideal replacement policy.

- Any mat-mult, $M > 2n^2$, only $\Theta(n^2)$ compulsory misses
- Conventional mat-mult, $M < n^2/2$, $\Theta(n^3/\sqrt{M})$ misses (blocking, recursion)
- Fast mat-mult, small cache, $O(n^\omega/M^{\omega/2-1})$ misses (less work, smaller data-reuse ratio)

Cache Misses in MM PS

If $M > (p + 1)n^2$, only $\Theta(n^2)$ compulsory misses

If $M < n^2/2$,

$$\Theta\left(\frac{(p + s)n^3}{\sqrt{M}} + spn^2\right) = \Theta\left(\frac{(p + d/p)n^3}{\sqrt{M}} + dn^2\right)$$

This is minimized for $p \approx \sqrt{d}$ so $\Theta(\sqrt{d}n^3/\sqrt{M} + dn^2)$

- $\Theta(\sqrt{d}n^3/\sqrt{M})$ for $M \leq n^2/d$ (tiny cache, MMs dominate)
- $\Theta(dn^2)$ for $n^2/d \leq M < n^2/2$ (small cache, matrix scale-adds dominate)

Mind the gap

Cache Misses in MM PS: Filling the Gap

for $3n^2 \leq M \leq (p+1)n^2$ MMs are comm-cheap but MSAs are not

Number of cache misses $\Theta(sp n^2) = \Theta(dn^2)$

Cache misses minimized by setting $p \approx M/n^2$, to get out of this regime

But arithmetic is not; yuck

Cache Misses in the MV PS

If $M > 3n^2 + pn$, only $\Theta(n^2)$ compulsory misses

If $pn \leq M < n^2/2$, $\Theta(\log p \cdot n^3/\sqrt{M})$ to generate A^p , plus n iterations in which we pay $\Theta(pn^2)$ to construct $Ae_j, \dots, A^{p-1}e_j$ and $\Theta(sn^2)$ for s MVs with A^p

In total,

$$\begin{aligned}\Theta\left(\frac{\log p \cdot n^3}{\sqrt{M}} + n(p+s)n^2\right) &= \Theta\left(\frac{\log p \cdot n^3}{\sqrt{M}} + \left(p + \frac{d}{p}\right)n^3\right) \\ &= \Theta\left(\left(p + \frac{d}{p}\right)n^3\right)\end{aligned}$$

Picking a (Theoretical) Winner

If $M > (p + 1)n^2$, use MM-PS (Strassen okay)

If $M > 3n^2 + pn$, use MV-PS (optimal arithmetic sans Strassen)

Otherwise (tiny cache), use MM-PS

But we can do better with new variants

New Algorithmic Ideas

1: Block-Column PS

Compute A^p ($\log_2 p$ m -multiplications)

Pick a block size b

For n/b blocks $E_j \leftarrow I_{:,j:j+b}$

 Compute $AE_j, \dots, A^{p-1}E_j$ ($p-1$ b -mat-vecs)

 Set $Q_j \leftarrow$ (highest coefficient polynomial) E_j (b -vec scale-add)

 For $k \leftarrow s-1, \dots, 1, 0$ multiply and add

$Q_j \leftarrow A^p Q_j +$ (next lower coefficient polynomial) E_j

$\Theta\left(\frac{\log p \cdot n^3}{\sqrt{M}} + \frac{n}{b}(p+s)n^2\right)$ cache misses for $pnb \leq M \leq n^2/2$

2: Use Fast Matrix Multiplication in MM-PS, BMV-PS

Great opportunity for fast mat-mult

Benefit in block-MV PS drops with block size (both arithmetic and data-reuse ratio)

3: Transform to (Real) Schur Form

If $A = VTV^*$, then $q(A) = Vq(T)V^*$, so the expensive part ($> O(n^3)$) is performed on triangular matrices, $\times 6$ arithmetic benefit

\sqrt{d} needs to be high enough to offset the cost of the Schur decomposition

Unfortunately, level-3 BLAS do not have triangular-triangular matrix multiplication, but can implement fairly easily using recursion

4: Schur Form, Reorder, Parlett+Davies-Higham

If $A = VTV^*$, then $q(A) = Vq(T)V^*$

Apply block Schur-Parlett substitution to $q(T)$ (evaluate q on diagonal blocks, solve Sylvester equations for off-diagonal blocks)

Higham-Davies: ensure Sylvester equations are well conditioned by partitioning $\Lambda(A)$ into well-separated clusters & reordering the Schur form (Bai-Demmel-Kressner)

Benefit: $\Theta(\sqrt{d}n^3)$ super-cubic algorithm applied to diagonal blocks of T , not all of it

Cost: $O(n^3)$ Schur decomposition and reordering

[Not implemented, but shows how novel discoveries can improve old algorithms]

5: Remainder Evaluation (Really an Open Problem)

If $d > n$, let χ be the Characteristic polynomial of A and let

$$q(A) = \chi(A)\delta(A) + \rho(A) = \rho(A) .$$

Clearly, evaluating $\rho(A)$ is cheaper than applying $q(A)$

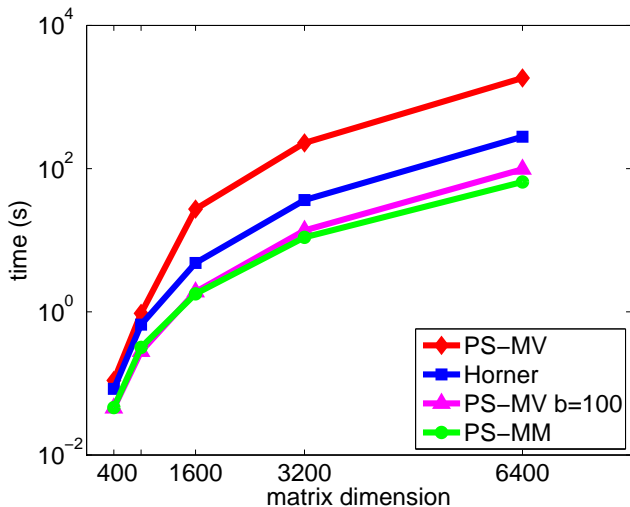
But can we determine the coefficients of ρ in a stable way?

Open problem (AFAIK)

I'll mention another interesting and related open problem at the end of the talk

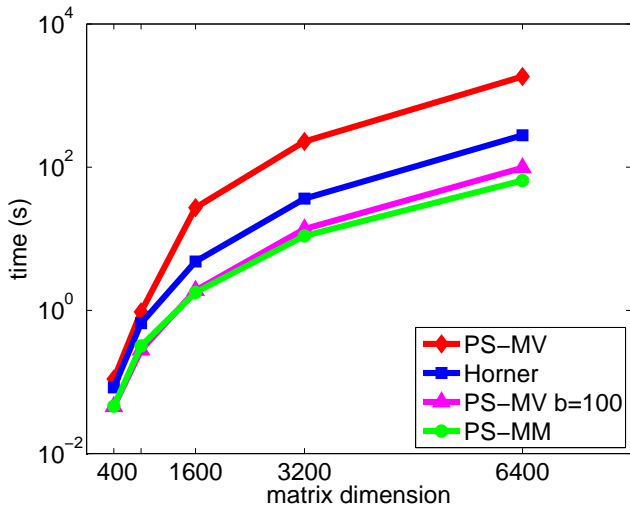
Experimental Results

$d = 100$; MM-PS and Block-MM-PS Better than Horner

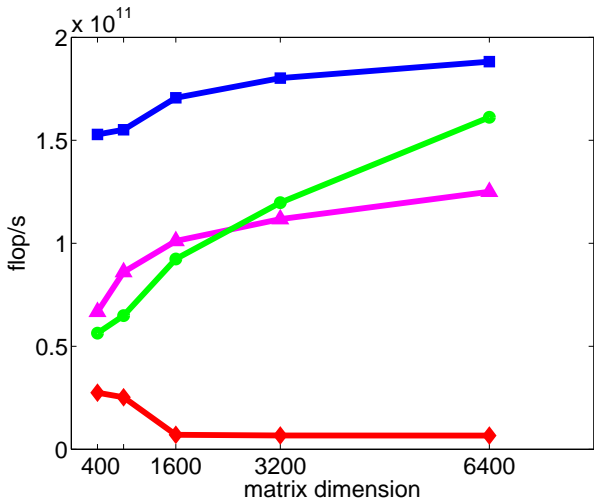


Quad-core
i7
C+MKL

MV PS is Terrible! Cache Misses!

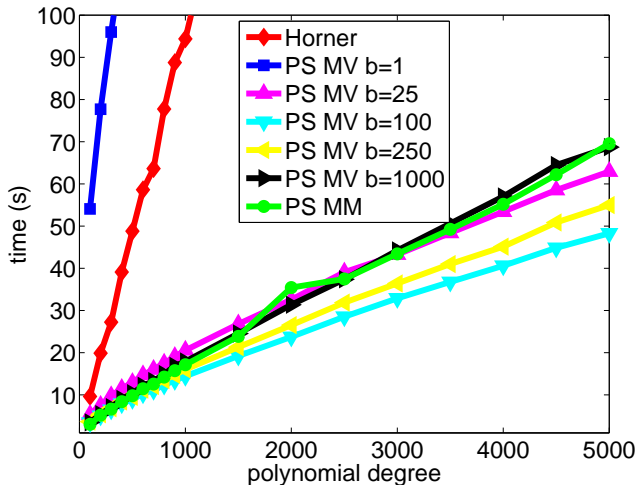


Flop/s Rates



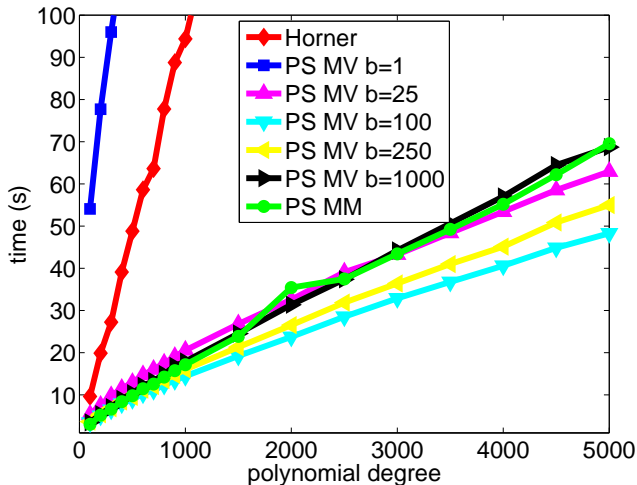
Block-Size Sensitivity (Not High but Not Monotone)

$n = 2000$



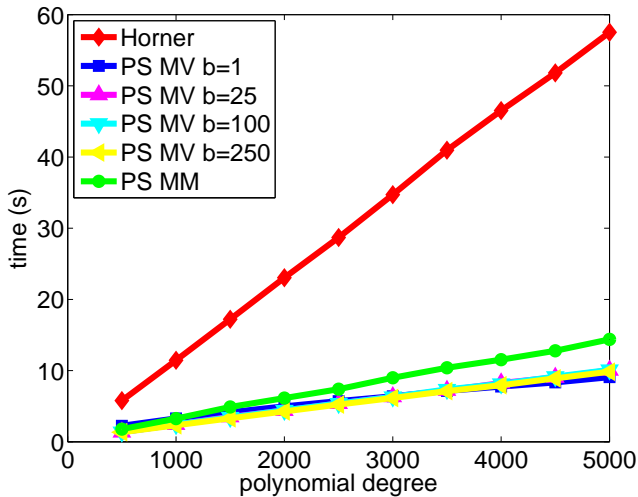
Also, For High Degrees Horner is Terrible

$n = 2000$

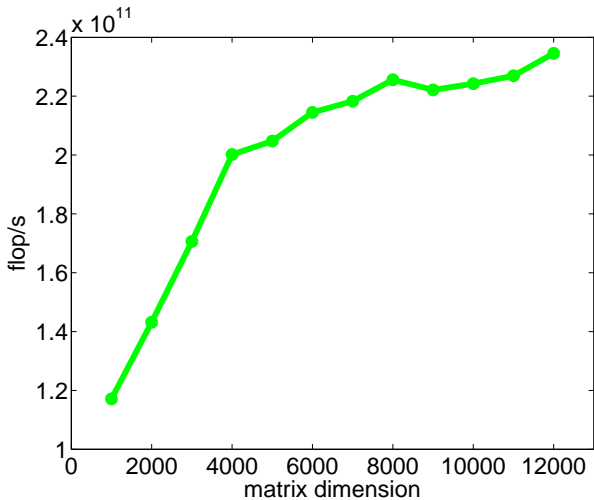


Reduced Sensitivity at Smaller Sizes

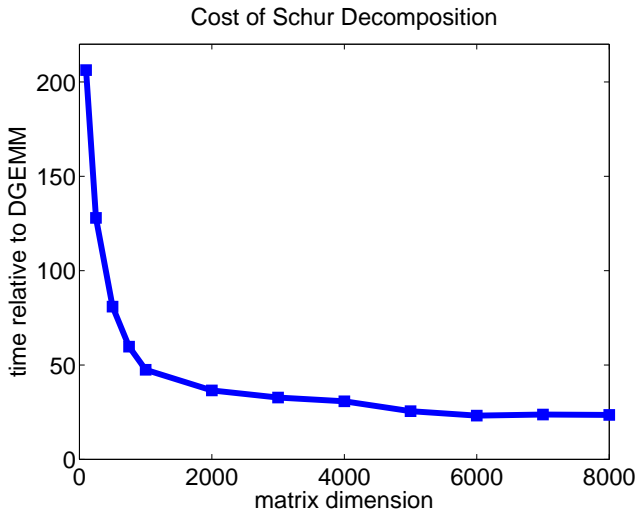
$n = 750$



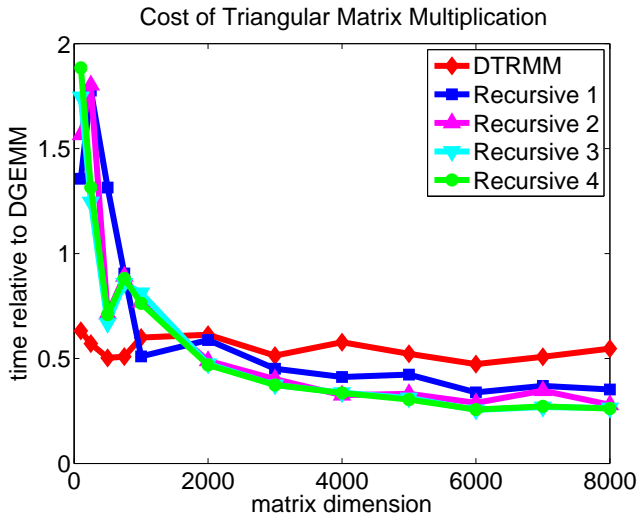
$d = 10$ Emphasizes Cost of Matrix Additions



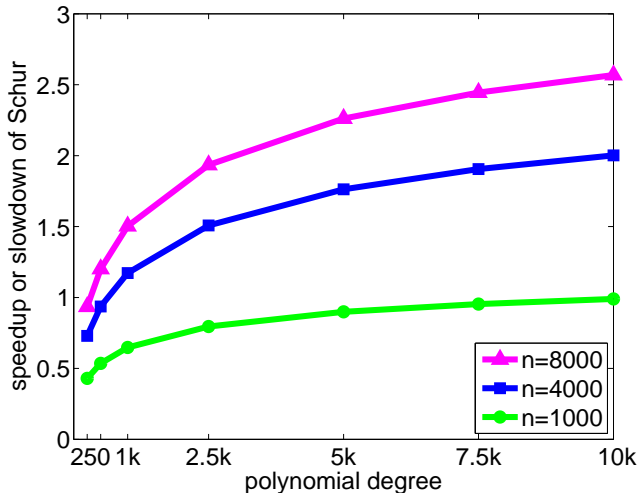
The Cost of the Schur Decomposition (≈ 25 Mat-Mults)



Savings from Triangular Matrix Multiplication ($\approx \times 4$)



Estimating The Cost of Schur + Triangular PS ($> \times 2.5$)



Conclusions and Open Problems

Conclusions

Old algorithms must be revisited (fairly obvious);

MM-PS is memory inefficient, MV-PS super slow due to cache misses

Block MV-PS fixes that

More novel tricks thanks to various innovations, mostly in extreme cases (very high degrees)

Open Problem 1: Applications?

Do these algorithms have interesting applications?

We were attracted to this due to the complex tradeoffs in Patterson-Stockmeyer variants

Seems that rational approximations are used more often in applications

Open Problem 2: Alternate Representations

Numerical analysis of Patterson-Stockmeyer in other (more useful) representations of q

Is there an efficient and numerically-stable version of PS for $d > n$? Think of $\rho(\mathcal{A})$ given in terms of the coefficients of $q(\mathcal{A}) = \chi(\mathcal{A})\delta(\mathcal{A}) + \rho(\mathcal{A})$

Is there an efficient and numerically-stable version of PS for Newton polynomials, or for any other forms of interpolation/least-squares polynomials?

Looking forward to discussions during the rest of the week!