Using GPUs to Accelerate the Solution of Large-Scale Model Reduction Problems

Peter Benner¹ Pablo Ezzatti² Francisco Igual³ Daniel Kressner⁴ Enrique S. Quintana-Ortí³ Alfredo Remón³

¹Fakultät für Mathematik, Chemnitz University of Technology, Chemnitz (Germany)

²Centro de Cálculo-Instituto de la Computación, Universidad de la República, Montevideo (Uruguay)

³Seminar für angewandte Mathematik, ETH Zürich, Zurich (Switzerland).

⁴Depto. de Ingeniería y Ciencia de Computadores, Universidad Jaume I, Castellón (Spain).



Dynamical Linear Systems

j'i

Linear time-invariant systems:

$$\dot{x}(t) = Ax(t) + Bu(t), \quad t > 0, \quad x(0) = x^0,$$

 $y(t) = Cx(t) + Du(t), \quad t \ge 0,$

- *n* state-space variables, i.e., *n* is the order of the system;
- *m* inputs,
- p outputs.

Corresponding TFM:

$$G(s) = C(sI_n - A)^{-1}B + D.$$

Goal for Model Reduction



Find a reduced-order model

$$\begin{split} \dot{\hat{x}}(t) &= \hat{A}\hat{x}(t) + \hat{B}u(t), & t > 0, & \hat{x}(0) = \hat{x}^0, \\ \hat{y}(t) &= \hat{C}\hat{x}(t) + \hat{D}u(t), & t \ge 0, \end{split}$$

of order $r \ll n$ such that the output error

$$y - \hat{y} = Gu - \hat{G}u = (G - \hat{G})u$$

is "small".

Motivation



Ϊī

 μ -mechanical Gyroscope [The Imego Institute (Sweden) + SAAB BOFORS DYNAMICS AB]

- Commercial rate sensor with applications in inertial navigation systems.
- Simulation problem: Improve the design with respect to a number of parameters.

● *n* = 17,361 states.

Can we obtain a reduced-order model with similar behavior?





- Truncation methods for model reduction
- 2 Solution of Lyapunov equations
- GPU implementationMatrix inversion
- Iterative refinement

5 Conclusions

Outline



Truncation methods for model reduction

- Solution of Lyapunov equations
- GPU implementationMatrix inversion
- Iterative refinement
- 5 Conclusions

Balanced Truncation (I)

j,

Balanced Truncation is an absolute error method, which aims at

 $\min \|G - \hat{G}\|_{\infty}$

Composed of the following three steps:

Step 1. Solve the coupled Lyapunov matrix equations

$$AW_c + W_c A^T + BB^T = 0,$$

$$A^T W_o + W_o A + C^T C = 0,$$

for the observability and controllability Gramians, W_c and W_o . Actually, we need the Cholesky factors *S* and *R* such that

$$W_c = S^T S, \quad W_o = R^T R.$$

S and R are dense.

Balanced Truncation (I)



Balanced Truncation is an absolute error method, which aims at

$$\min \|G - \hat{G}\|_{\infty}$$

Composed of the following three steps:

Step 1. Solve the coupled Lyapunov matrix equations

$$AW_c + W_c A^T + BB^T = 0,$$

$$A^T W_o + W_o A + C^T C = 0,$$

for the observability and controllability Gramians, W_c and W_o . Actually, we need the Cholesky factors *S* and *R* such that

$$W_c = S^T S, \quad W_o = R^T R.$$

S and R are dense.

Balanced Truncation (II)



Step 2. Compute the Hankel singular values (HSV) from

$$SR^{T} = U\Sigma V^{T} = \begin{bmatrix} U_{1} & U_{2} \end{bmatrix} \begin{bmatrix} \Sigma_{1} & \\ & \Sigma_{2} \end{bmatrix} \begin{bmatrix} V_{1}^{T} \\ V_{2}^{T} \end{bmatrix},$$

with U, V, and Σ partitioned at a certain order r.

The HSV in $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$, measure how much a state is involved in energy transfer from a given input to a certain output!

Balanced Truncation (III)



Step 3. In the square-root balance truncation (SRBT) method

$$T_l = \Sigma_1^{-1/2} V_1^T R$$
 and $T_r = S^T U_1 \Sigma_1^{-1/2}$,

and $(\hat{A}, \hat{B}, \hat{C}, \hat{D}) = (T_l A T_r, T_l B, C T_r, D)$ for the TFM:

$$\hat{G}(s) = CT_r(sI_n - T_lAT_r)^{-1}T_lB + D.$$

• Computable error bound: $\|G - \hat{G}\|_{\infty} \leq 2\sum_{k=r+1}^{n} \sigma_{k}$.

Balanced Truncation (IV)



Given (A, B, C, D, x^0) with A large, and $m, p \ll n \dots$

How do we solve the previous numerical problems?

Coupled Lyapunov equations.

- SVD of matrix product.
- Application of the SRBT formulae to obtain the reduced-order model.

Outline



Truncation methods for model reduction

Solution of Lyapunov equations

GPU implementationMatrix inversion

Iterative refinement

5 Conclusions

Sign Function Method



Given $\alpha \in \mathbb{R}$,

sign (
$$\alpha$$
) =

$$\begin{cases}
1 \text{ if } \alpha > 0, \\
-1 \text{ if } \alpha < 0, \\
\text{undefined otherwise.}
\end{cases}$$

For a matrix $A \in \mathbb{R}^{n \times n}$, sign (A) is a function of the signs of its eigenvalues.

Given

$$H = \begin{bmatrix} A & 0 \\ C^T C & -A^T \end{bmatrix}, \quad \operatorname{sign}(H) = \begin{bmatrix} -I_n & 0 \\ 2W_o & I_n \end{bmatrix},$$

where W_o is the observability Gramian.

So, how do we compute the sign function?

Sign Function Methods (I)

For
$$H = \begin{bmatrix} A & 0 \\ C^T C & -A^T \end{bmatrix}$$
 the classical Newton iteration boils down to
 $A_{j+1} = \frac{1}{2}(A_j + A_j^{-1})/2, \quad A_0 = A,$
 $R_{j+1} = \frac{1}{\sqrt{2}}\begin{bmatrix} R_j \\ R_j A_j^{-1} \end{bmatrix}, \quad R_0 = C,$

which converges to R, the Cholesky factor of W_o .

At each iteration R_i is increased in p rows ($p \Rightarrow$ number of outputs).

The computation of the inverse represents the main part of the computation ($O(2n^3)$ flops).

Sign Function Methods (II)

As in model reduction *R* (and *S*) is usually rank-deficient the cost of the iteration and subsequent steps can be greatly reduced:

At the *j*th iteration, compute the rank-revealing QR (RRQR) factorization

$$\frac{1}{\sqrt{2}} \left[\begin{array}{c} R_j \\ R_j A_j^{-1} \end{array} \right] = \bar{Q}\bar{R}\Pi$$

and then set

$$R_{j+1} = (\bar{R}\Pi)^T.$$

On convergence the iteration produces dense, full-rank \hat{R} with $l \ll n$ columns, such that

$$\hat{R}^T \hat{R} \approx R^T R = W_o.$$

Outline



- Truncation methods for model reduction
- 2 Solution of Lyapunov equations
- GPU implementationMatrix inversion
- Iterative refinement

5 Conclusions

Hybrid approach for the sign function

Hybrid approach

Each step is performed in the most suitable device:

•
$$A_{j+1} = \frac{1}{2}(A_j + A_j^{-1})/2, \quad A_0 = A \Rightarrow \text{Matrix inverse on GPU}$$

• $R_{j+1} = \frac{1}{\sqrt{2}} \begin{bmatrix} R_j \\ R_j A_j^{-1} \end{bmatrix}, \quad R_0 = C \Rightarrow \text{GEMM on CPU or GPU}$
• RRQR $\Rightarrow \text{Executed on CPU}$

Via LU factorization

- $\bigcirc PA = LU$
- $\bigcirc U \to U^{-1}$
- Solve the system $XL = U^{-1}$ for X
- Undo the permutations $A^{-1} := XP$

Implementation

- The algorithm sweeps through the matrix four times
- Presents a mild load imbalance, due to the work with triangular factors

Algorithm implemented by LAPACK

Ĵ۲



Via Gauss-Jordan elimination (GJE)

- Reordering of the computations of LU-based methods
- Requires the same arithmetic cost

Implementation

- The algorithm sweeps through the matrix once
- Most of the computations are highly parallel

Matrix inversion

Matrix inversion (GJE)



Figure: Blocked algorithm for matrix inversion via GJE without pivoting.

j

Matrix inversion (GJE)

GPU implementation

- The matrix is transferred to the GPU
- The inverse is computed completely on the GPU
- Result is transferred back to the CPU

Hybrid implementation

- GPU computes all the matrix-matrix products
- CPU computes the GJEUNB
- Only small (column) panels are transferred

Experimental setup

CPU	Dual Xeon QuadCore E5410
CPU frequency	2.33 Ghz
RAM memory	8 Gbytes
GPU	Tesla C1060
Processor	Nvidia GT200
GPU frequency	1.3 Ghz
Video memory	4 Gbytes DDR3
Interconnection	PCIExpress Gen2
CUDA (CUBLAS) version	2.1
BLAS implementation	GOTOBlas 1.26
Driver version	185.18

Results for matrices with $1000 \le n \le 8000$ and $b \le 200$ Transfer times included in all results

Matrix inversion

Experimental Results

(Matrix inverse - GotoBLAS)





Matrix inversion on Caton2 + Goto 1.26

Matrix inversion

Experimental Results

(Single precision Matrix Sign Function - GotoBLAS)



Sign Function on Caton2 + Goto 1.26

バ

Matrix inversion

Experimental Results

(Double Precision Matrix Sign Function - GotoBLAS)



バ

Outline



- Truncation methods for model reduction
- Solution of Lyapunov equations
- GPU implementationMatrix inversion

Iterative refinement

5 Conclusions

Iterative refinement (I)



Given a Lyapunov equation:

$$AX + XA^T = -BB^T$$

Goals ● Exploit the single-precision capabilities of GPUs. ● Get an approximation of the solution in low-precision arithmetic: L = ApproxLyap(B), X ≈ LL^T ● Refine the result to regain full accuracy.

Iterative refinement (II)



Let

 $L_0 = \operatorname{ApproxLyap}(B)$ SinglePrecision(GPU)

to improve L_0 we construct a correction based on the residual

$$\mathbf{Res} = AL_0L_0^T + L_0L_0^TA^T + BB^T \qquad DoublePrecision$$

and solve

$$L_1 = \text{ApproxLyap}(-\text{Res})$$
 SinglePrecision

to get the correction term.

Problem

Res is usually indefinite

Iterative refinement (III)

Ĵ

Solution

Decompose Res into a positive definite and a negative definite part:

$$\mathsf{Res} = \mathbf{R}_p \mathbf{R}_p^T - \mathbf{R}_n \mathbf{R}_n^T$$

Each term corresponds to a Lyapunov equation (solved in SP):

$$AX_p + X_pA^T = -R_pR_p^T \qquad A(-X_n) + (-X_n)A^T = -R_nR_n^T$$

Then $X_c = X_p + X_n$ solves the correction equation

$$AX_c + X_c A^T = -\text{Res}$$

Corrected solution:

$$X_1 = L_0 L_0^T + L_p L_p^T - L_n L_n^T$$

Iterative refinement

Numerical example (MATLAB)

- $A \rightarrow 900 \times 900$ symmetric negative definite matrix
- $\bullet~$ it \rightarrow number of sign function iterations
- $\bullet \ tol \rightarrow tolerance$ for sign function
- $res_i \rightarrow residual$ after *i* steps of iterative refinement







- Truncation methods for model reduction
- 2 Solution of Lyapunov equations
- GPU implementationMatrix inversion

Iterative refinement



Solution of (large) model reduction problems applying GPUs

- Truncation Methods \rightarrow Lyapunov equations
- Hybrid approach for the Sign Function to solve Lyapunov equations

Iterative refinement approach to combine full-accuracy and high performance

Thank you!

More information...

- HPCA Group at University Jaume I (http://www.hpca.uji.es)
- {figual, quintana, remon}@icc.uji.es