

# Zero-Copy Parallelized Empirical Gramians

Christian Himpe ([christian.himpe@uni-muenster.de](mailto:christian.himpe@uni-muenster.de))  
Mario Ohlberger ([mario.ohlberger@uni-muenster.de](mailto:mario.ohlberger@uni-muenster.de))

WWU Münster  
Institute for Computational and Applied Mathematics

PACO 2015  
2015-07-06

Now Tell Me:

# What are empirical gramians?

→ Empirical gramians are a system-theoretic tool for model reduction and system identification of input-output systems!

## Why **Model Order Reduction**<sup>1</sup> (MOR)?

- Enable or accelerate large-scale system evaluation,
- repeatedly in many-query settings,
- such as model-constrained optimization or control problems.

## Why **Power-Aware** Reduced Order Model (ROM) Computation?

- Reduce waste of computationally resources.
- Realistically, offline time is not infinite.
- Save the Planet!

---

<sup>1</sup> Model Reduction is in itself PACO by using ROMs, i.e. on low-power devices.

# Outline<sup>2</sup>

- 1 Mathematical Background
- 2 Computational Means
- 3 Power Aware COmputing

---

<sup>2</sup>Disclaimer: There are no experimental results yet due to unreleased hardware.

# Mathematical Background

# State-Space Systems

(Nonlinear) State-Space System:

$$\dot{x}(t) = f(x(t), u(t), \theta),$$

$$y(t) = g(x(t), u(t), \theta),$$

$$x(0) = x_0$$

$$\text{State: } x(t), \dim(x(t)) \gg 1$$

$$\text{Input: } u(t), \dim(u(t)) \ll \dim(x(t))$$

$$\text{Output: } y(t), \dim(y(t)) \ll \dim(x(t))$$

$$\text{Parameters: } \theta, \dim(\theta) \gg 1$$

Linear State-Space System:

$$\dot{x}(t) = Ax(t) + Bu(t),$$

$$y(t) = Cx(t),$$

$$x(0) = x_0$$

$$\text{System Matrix: } A \in \mathbb{R}^{N \times N}$$

$$\text{Input Matrix: } B \in \mathbb{R}^{N \times J}$$

$$\text{Output Matrix: } C \in \mathbb{R}^{O \times N}$$

# Projection-Based Model Reduction

Reduced Order Model (ROM):

$$\dot{x}_r(t) = f_r(x_r(t), u(t), \theta),$$

$$y_r(t) = g_r(x_r(t), u(t), \theta),$$

$$x_r(0) = x_{r,0}$$

$$\dim(x_r(t)) \ll \dim(x(t))$$

$$\|y(\theta) - y_r(\theta)\| \ll 1$$

Projection-Based ROM:

$$\dot{x}_r(t) = Vf(Ux_r(t), u(t), \theta),$$

$$y_r(t) = g(Ux_r(t), u(t), \theta),$$

$$x_r(0) = Vx_0$$

$$U \in \mathbb{R}^{\dim(x(t)) \times \dim(x_r(t))}$$

$$V \in \mathbb{R}^{\dim(x_r(t)) \times \dim(x(t))}$$

$$VU = \mathbb{1}$$

# Gramian-Based Model Reduction [Moore'81, Fernando & Nicholson'83]

**Controllability** Gramian:

$$W_C = \int_0^{\infty} e^{At} B B^T e^{A^T t} dt$$

$$\Rightarrow A W_C + W_C A^T = -B B^T$$

**Observability** Gramian:

$$W_O = \int_0^{\infty} e^{A^T t} C^T C e^{At} dt$$

$$\Rightarrow A^T W_O + W_O A = -C^T C$$

→ Balanced Truncation

**Cross Gramian:**

$$W_X = \int_0^{\infty} e^{At} B C e^{At} dt$$

$$\Rightarrow A W_X + W_X A = -B C$$

→ Direct Truncation



## Empirical Gramians [Lall et al'99]

System Gramians,  
from primal and adjoint impulse response trajectories:

$$W_C = \int_0^{\infty} (e^{At}B)(e^{At}B)^T dt = \int_0^{\infty} x_{\delta}(t)(x_{\delta}(t))^T dt$$
$$W_O = \int_0^{\infty} (e^{At}C)^T (Ce^{At}) dt = \int_0^{\infty} x_{\delta}^*(t)(x_{\delta}^*(t))^T dt$$

Extend to nonlinear system<sup>3</sup>, since only based on trajectories!

$$W_X = \sum_{\{\Delta u, \Delta x_0\}} \int_0^{\infty} x_{\delta}(t)(x_{\delta}^*(t))^T dt$$

Empirical Gramians Computation:

- 1 Simulate Trajectories
- 2 Gramian Assembly

---

<sup>3</sup>This is better than linearization!

## Combined Reduction [H. & Ohlberger'14]

Combined State and Parameter Reduction:

$$\begin{aligned}\dot{x}_r(t) &= Vf(Ux_r(t), u(t), P\theta_r), & P &\in \mathbb{R}^{\dim(\theta_r) \times \dim(\theta)}, P^T P = \mathbb{1} \\ y_r(t) &= g(Ux_r(t), u(t), P\theta_r), & \dim(\theta_r) &\ll \dim(\theta) \\ x_r(0) &= Vx_0 & \|y(\theta) - y_r(\theta_r)\| &\ll 1 \\ \theta_r &= P^T \theta\end{aligned}$$

Augmented System:

$$\begin{aligned}\begin{pmatrix} \dot{x}(t) \\ \dot{\theta}(t) \end{pmatrix} &= \begin{pmatrix} f(x(t), u(t), \theta(t)) \\ 0 \end{pmatrix} \\ y(t) &= g(x(t), u(t), \theta(t)) \\ \begin{pmatrix} x(0) \\ \theta(0) \end{pmatrix} &= \begin{pmatrix} x_0 \\ \theta_0 \end{pmatrix}\end{aligned}$$

# Empirical Gramian Flavors<sup>4</sup> [H. & Ohlberger'13]

- Empirical Controllability Gramian
- Empirical Observability Gramian
- Empirical Cross Gramian
- Empirical Linear Cross Gramian
- Empirical Non-Symmetric Cross Gramian
- Empirical Sensitivity Gramian (Parameter Controllability)
- Empirical Identifiability Gramian (Parameter Observability)
- Empirical Joint Gramian (Cross-Gramian-Based Combined Reduction)

---

<sup>4</sup> Compute with: emgr - Empirical Gramian Framework (<http://gramian.de>)

# Computational Means

## CPU vs GPU

	CPU	GPU
Cores	$10^1$	$10^3$
Clock	~ 4 Ghz	~ 1 Ghz
Flow-Control	☺	☹
GEMM	☹	☺
Parallelism	SMT, MT, SIMD	SIMT, SIMD

# CPU + iGPU

Previously, on *Game of CPUs*:

**1989** 80486DX: CPU + iFPU

**2004** GMA900: CPU + iGPU (Somewhat of a GPU)

**2011** Llano: CPU + iGPU (=:APU)

**2012** Intel HD: CPU + iGPU

**2014** Kaveri: HSA-Ready APU (+hUMA)

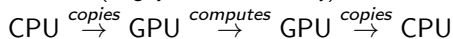
Generally:

- CPU and GPU in one package
- CPU and GPU use shared memory

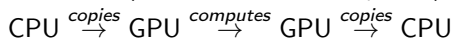
# Zero-Copy

GPGPU with:

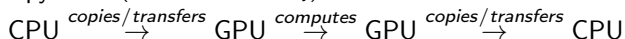
- discrete GPU (roughly Distributed Memory):



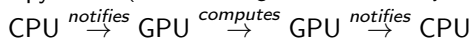
- integrated GPU (nUMA - Non-Uniform Memory Access):



- zero-copy GPU (Shared Virtual Memory):



- zero-copy iGPU (hUMA - Heterogenous Unified Memory Access):





# Easy As BLAS

## Requirements:

- capable of zero-copy
- should not affect high-level code (i.e. Matlab)
- on BLAS layer

## Weapon of Choice: **Offloading** via ACML<sup>5</sup>

**Pro** No changes in HLL implementation required

**Con** LUA configuration requires knowledge of the problem

**Pro** Effectively uses OpenCL

**Con** Works only on supported systems (i.e. HSA/huma)

**Pro** ACML backend is the open-source cBLAS

**Con** ACML is free (of charge), but not open-source

---

<sup>5</sup>The MKL can offload to Xeon Phi, too.

# iGPU Offloading from ACML

Where the magic happens:

`acmlbasedir/gfortran64_mp/lib/resource`

1 create device context: `context.lua`

see: `createContexts( platforms )`

2 BLAS config: `DEVICE_NAME/BLAS_CALL.lua`

see for example: `Spectre/gemm.lua`

CPU / GPU heuristics in: `tableOfThresholds`

# Zero-Copy Empirical Gramians

$$W_X = \sum_{\{\Delta u, \Delta x_0\}} \int_0^{\infty} x_{\delta}(t)(x_{\delta}^*(t))^T dt$$

## 1 Trajectory Simulation

- General Linear Methods such as (Two-Step) Runge-Kutta
- Rather Serial
- Parallelized on CPU (i.e. multiple trajectories in parallel each SIMD'd)

## 2 Gramian Assembly

- Dense Matrix Multiplication
- Very Parallel
- (to be) Parallelized on GPU

Now, how is this PACO?

# Why APUs?

- Discrete GPU needs (more) power
- Discrete GPU has relatively little memory
- Virtual Shared Memory uses PCIe or copies
  
- Integrated GPU needs less power
  - (i.e. A10-7850 APU (95W) vs Athlon X4 CPU (65W) + HD 7750 GPU (55W))
- Integrated GPU<sup>6</sup> can use full(!) system main memory
- Literally, 0 copies!

---

<sup>6</sup>with hUMA or VSM

# Mathematical Setup

Hyperbolic Network Model:

$$\dot{x}(t) = A \tanh(K(\theta)x(t)) + Bu(t)$$

$$y(t) = Cx(t)$$

System Dimensions:

- $x(t) \in \mathbb{R}^{1024}$
- $u(t) \in \mathbb{R}$
- $y(t) \in \mathbb{R}$
- $\dim(\theta) = \dim(x(t))$

# Numerical Experiment

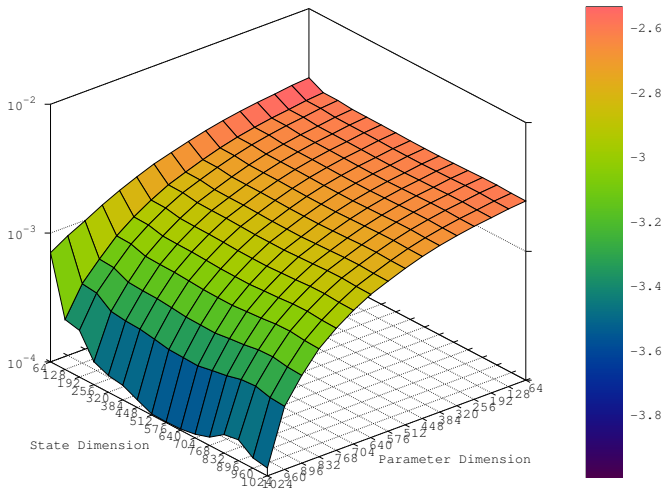


Figure :  $L_2$ -output-error for reduced state and parameter dimensions.

## Offline Time Distribution:

- Trajectory Computation: 99%
- Gramian Assembly Time<sup>7</sup>: 1%

## But:

- With  $N = 1024$ , this is still a small system.
- The trajectory computation parallelizes communication free.
- Can also be used for vector field evaluations<sup>8</sup>.
- It is basically free.

---

<sup>7</sup>This has been reduced by using Generalized Transpositions.

<sup>8</sup>For  $N > 512$  the default thresholds are met.



# Outlook

Next, on *Game of CPUs* (“OpenCL 2.0 is coming!”):

**2015** Carrizo: Full HSA 1.0 (+hUMA)

**2015** Broadwell: OpenCL 2.0 support

hUMA / VSM is on the rise<sup>9</sup>,

- attractive for scientific computing,
- accompanied by “smart” abstractions.

---

<sup>9</sup>This was one of the top developer requests for PS4 and XBone.

# Rough Road Ahead!

So, hardware and opportunities are (about to be) available, yet...

- the software stack:

- 1 Driver (Kernel, Propriety, HSA)

- 2 OpenCL (ICD, Dispatch)

- 3 Library (BLAS, LAPACK)

interaction is not trivial;

- optimization of heuristics is tedious;
- distributed memory systems with APUs are not on the horizon;
- system memory is usually slower than video memory.

- GPGPU becomes more accessible,
- and potentially computes LA more energy-efficiently.
- PACO (planned) by exploiting iGPUs
- with zero-copy capabilities
- for empirical gramian computation<sup>10</sup>.

<http://wwwmath.uni-muenster.de/u/himpe>

Thanks!

---

<sup>10</sup> Get the companion code: <http://j.mp/paco15>