

Power Efficiency for Software Algorithms running on Graphics Processors

Björn Johnsonson

Per Ganestam

Michael Doggett

Tomas Akenine-Möller



Overview

-
- Motivation
 - Goal
 - Project
 - Applications
 - Methodology
 - Results
 - Observations



Motivation

-
- Energy efficiency increasingly important
 - Phones, tablets, laptops, desktops...
 - Hard area
 - Harder to analyze than regular performance
 - Need hardware and/or hardware support
 - Less intuitive / harder to predict



Motivation

-
- Lots of papers looking at algorithms for power efficient hardware
 - What can be done from the software side on current hardware?
 - Learn more about energy and power



Goal

-
- Say we want to optimize to lower energy usage
 - Do we have to measure power?
 - Or can we make an estimate based on rendering times alone?



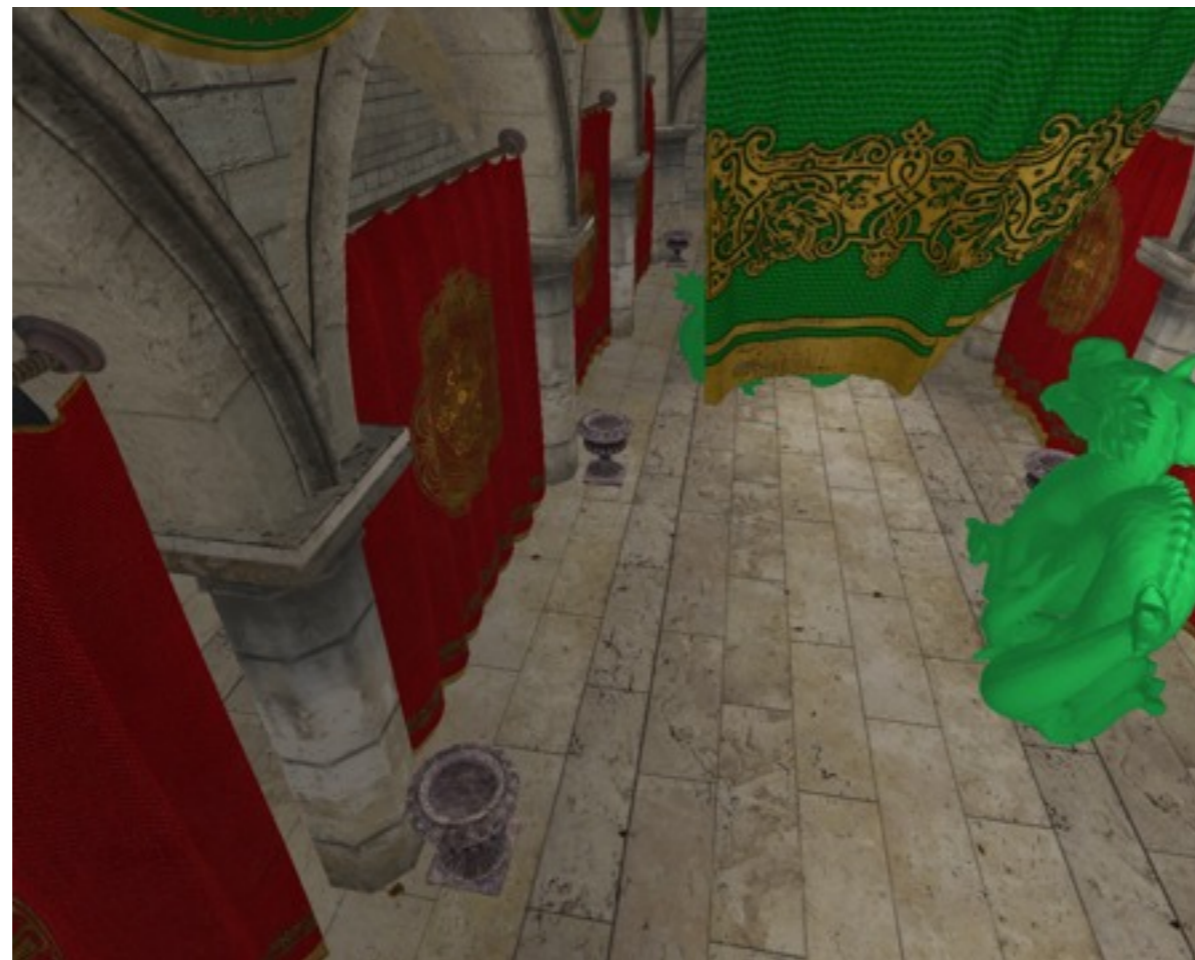
Project

-
- Two common graphics problems
 - Implemented several different solutions solving those problems
 - Measure their power usage and rendering time



Applications

- Primary visibility and shading
 - Forward rendering
 - Forward rendering with pre-Z pass
 - Deferred shading



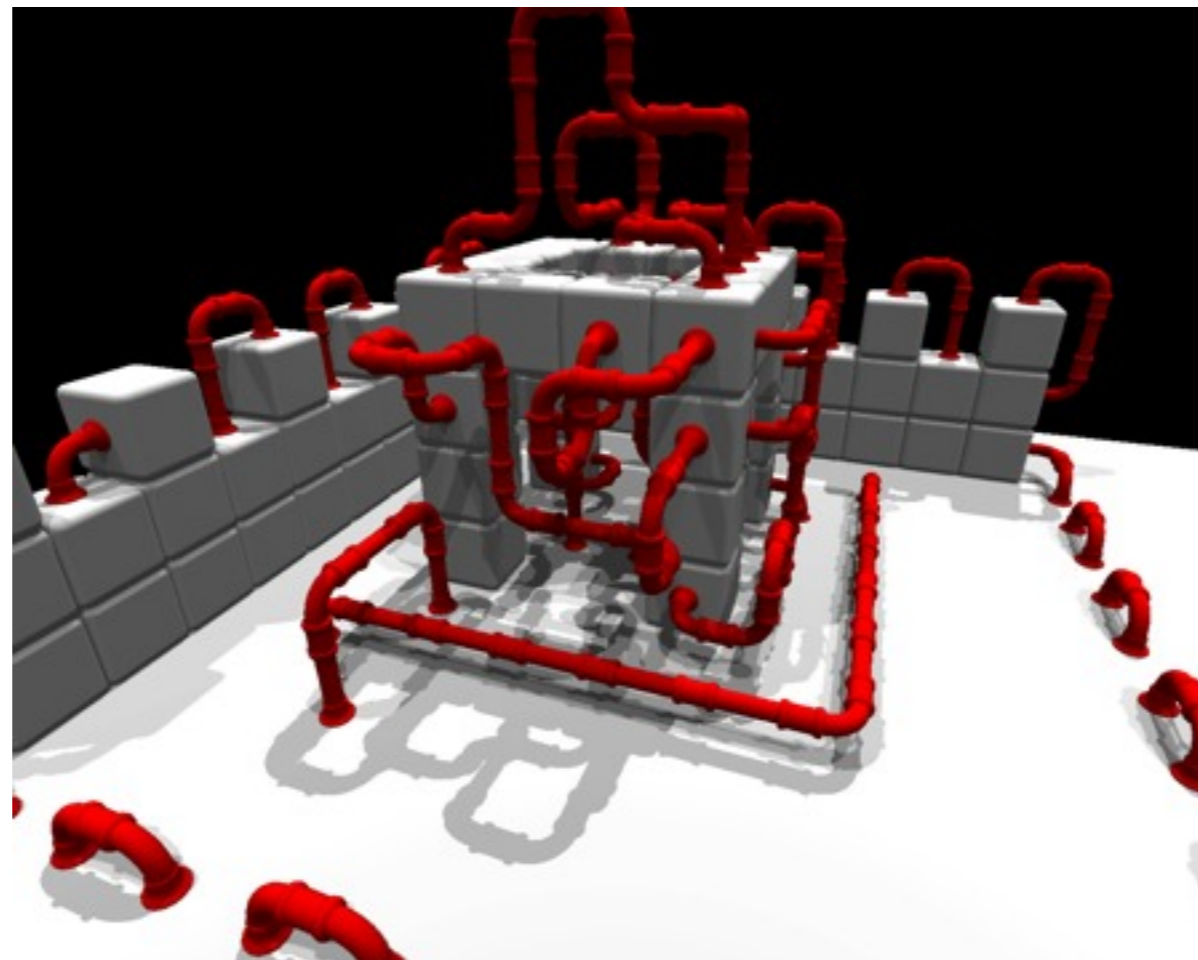
Applications

- Primary visibility and shading
 - Forward rendering
 - Forward rendering with pre-Z pass
 - Deferred shading



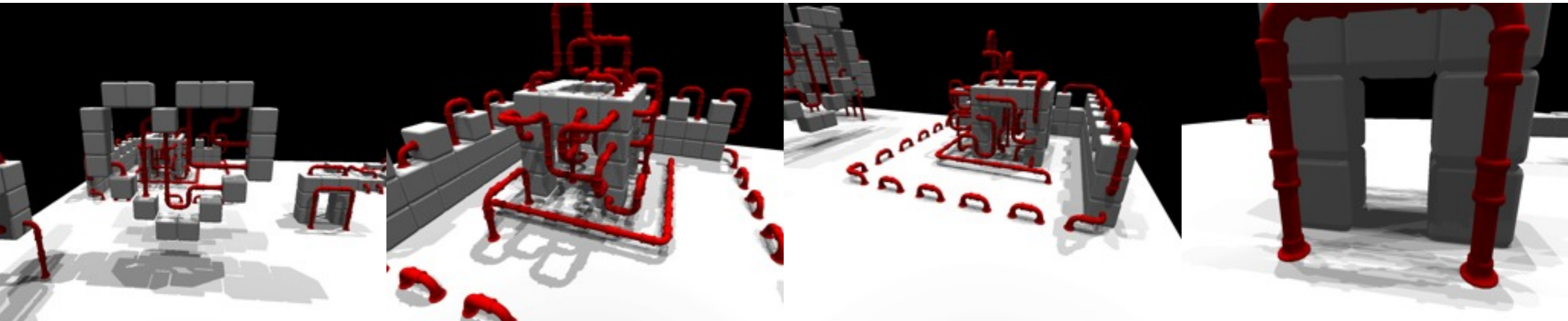
Applications

- Shadow algorithms
 - Stencil shadow volumes
 - Shadow mapping
 - Variance shadow mapping



Applications

- Shadow algorithms
 - Stencil shadow volumes
 - Shadow mapping
 - Variance shadow mapping



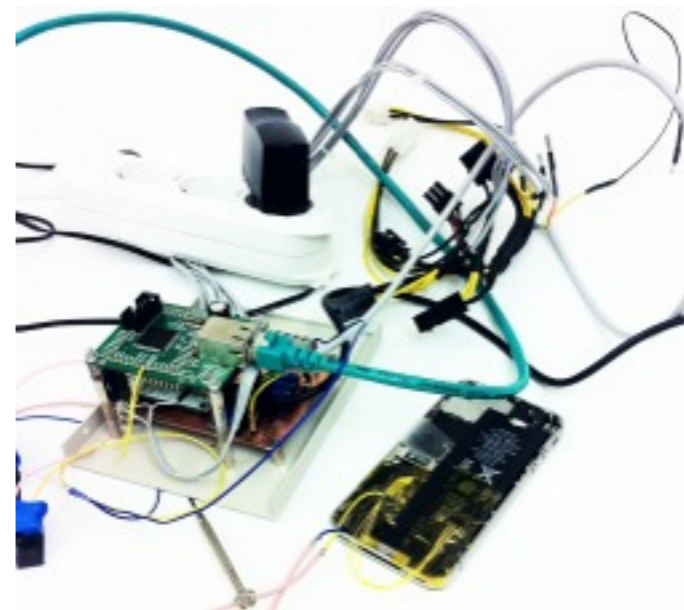
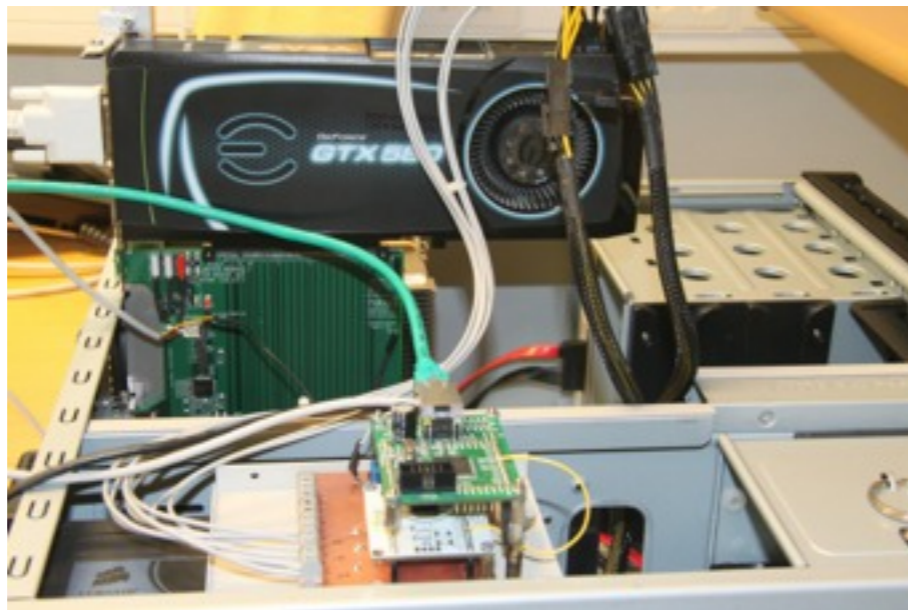
Applications

-
- OpenGL / OpenGL ES
 - Full speed (no capping)
 - Widely different platforms gives different frame times (5ms to 2s)
 - Timestamp at beginning and end of frame
 - Only integrate over actual rendering time
 - Animated camera path (~2000 frames)



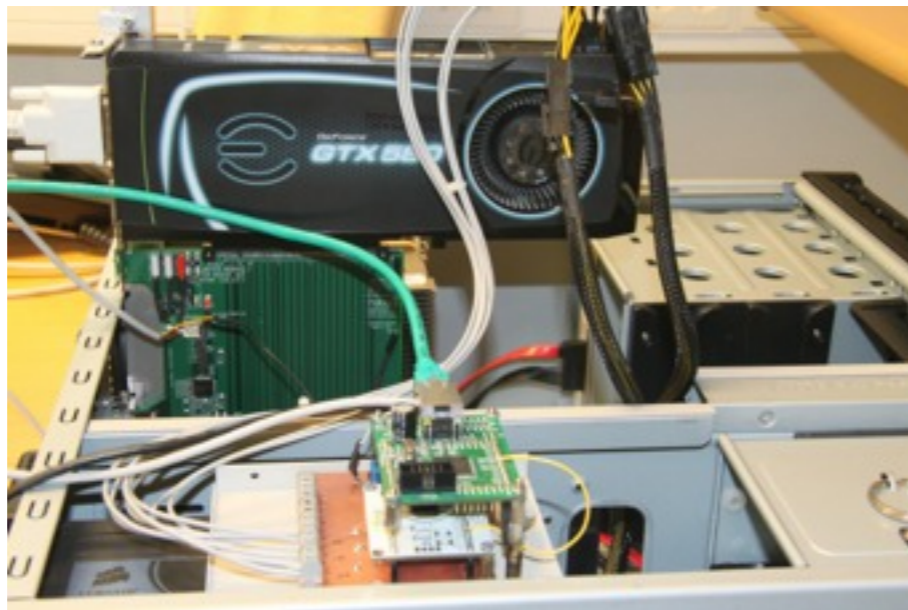
Methodology

- We have built a measurement station
 - Measure at 40kHz
 - 4 ACS710 Hall effect current sensors (<12A)
 - 2 shunt current sensors (<1A)
- Connected between GPU and power source
 - Different places on different platforms



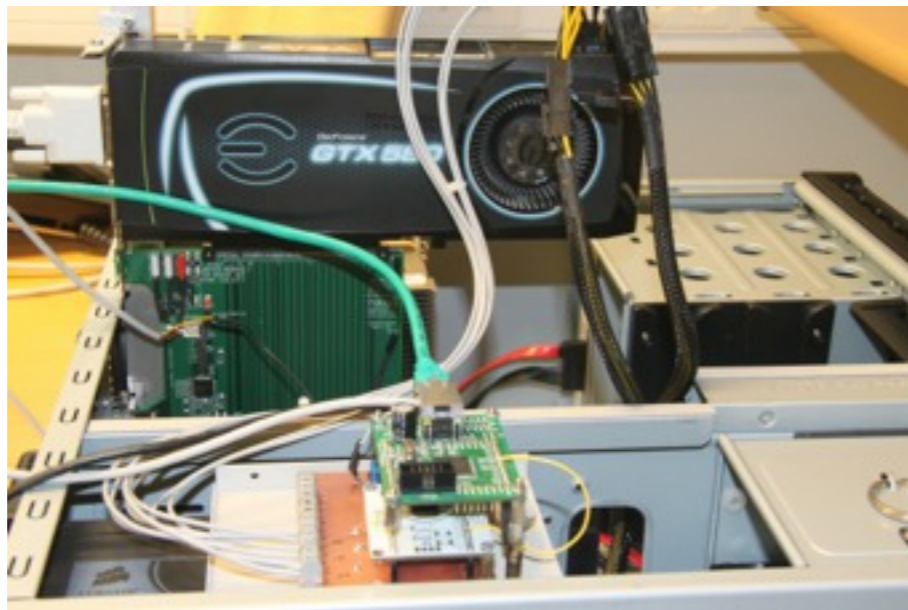
Discrete graphics cards

- Connected on the PCI-Express bus



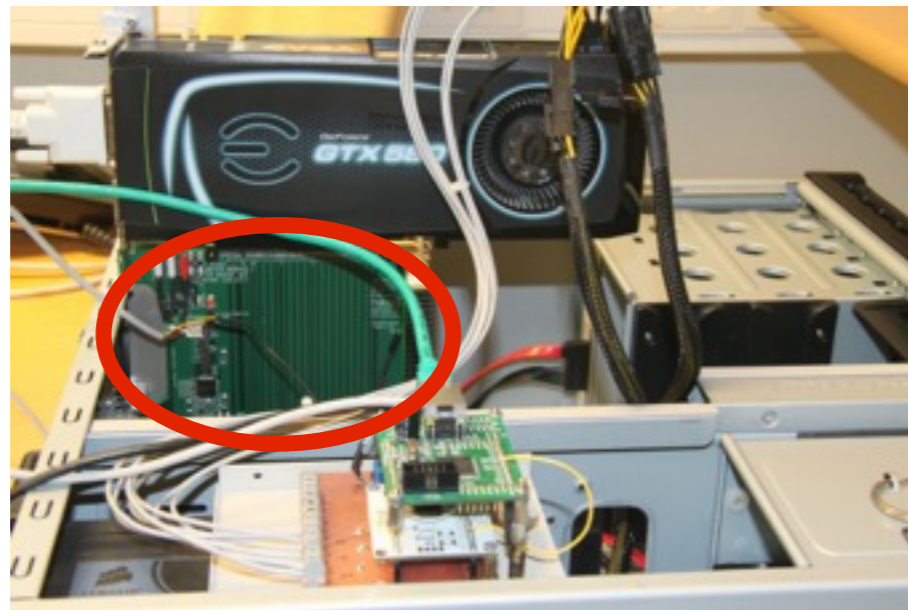
Discrete graphics cards

- Connected on the PCI-Express bus
 - PCI-Express bus provides $\leq 75W$



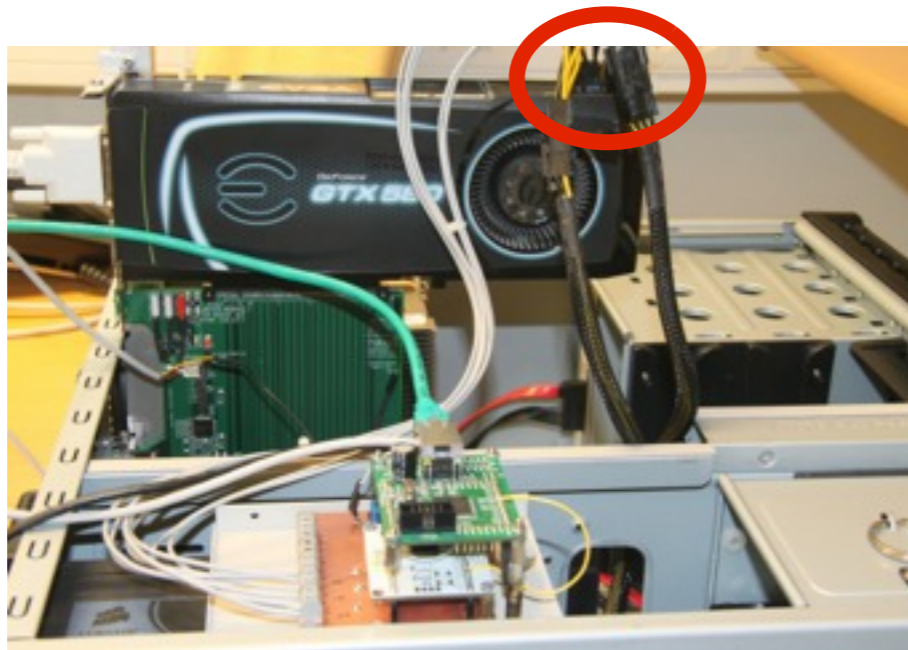
Discrete graphics cards

- Connected on the PCI-Express bus
 - PCI-Express bus provides $\leq 75W$
 - Measured through an Ultraview PCIeEXT-16HOT expander card



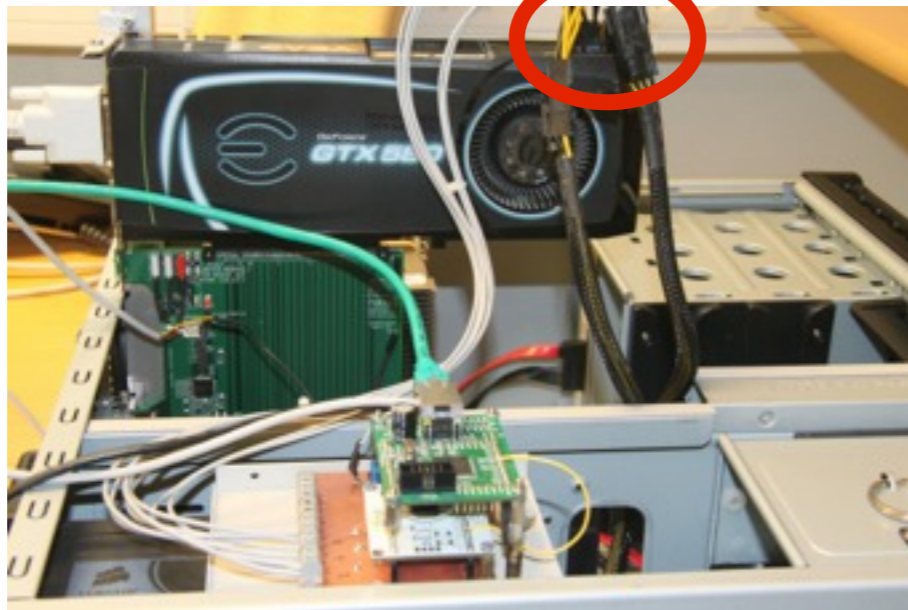
Discrete graphics cards

- Connected on the PCI-Express bus
 - PCI-Express bus provides $\leq 75W$
 - Measured through an Ultraview PCIeEXT-16HOT expander card
 - PCI-Express power connectors

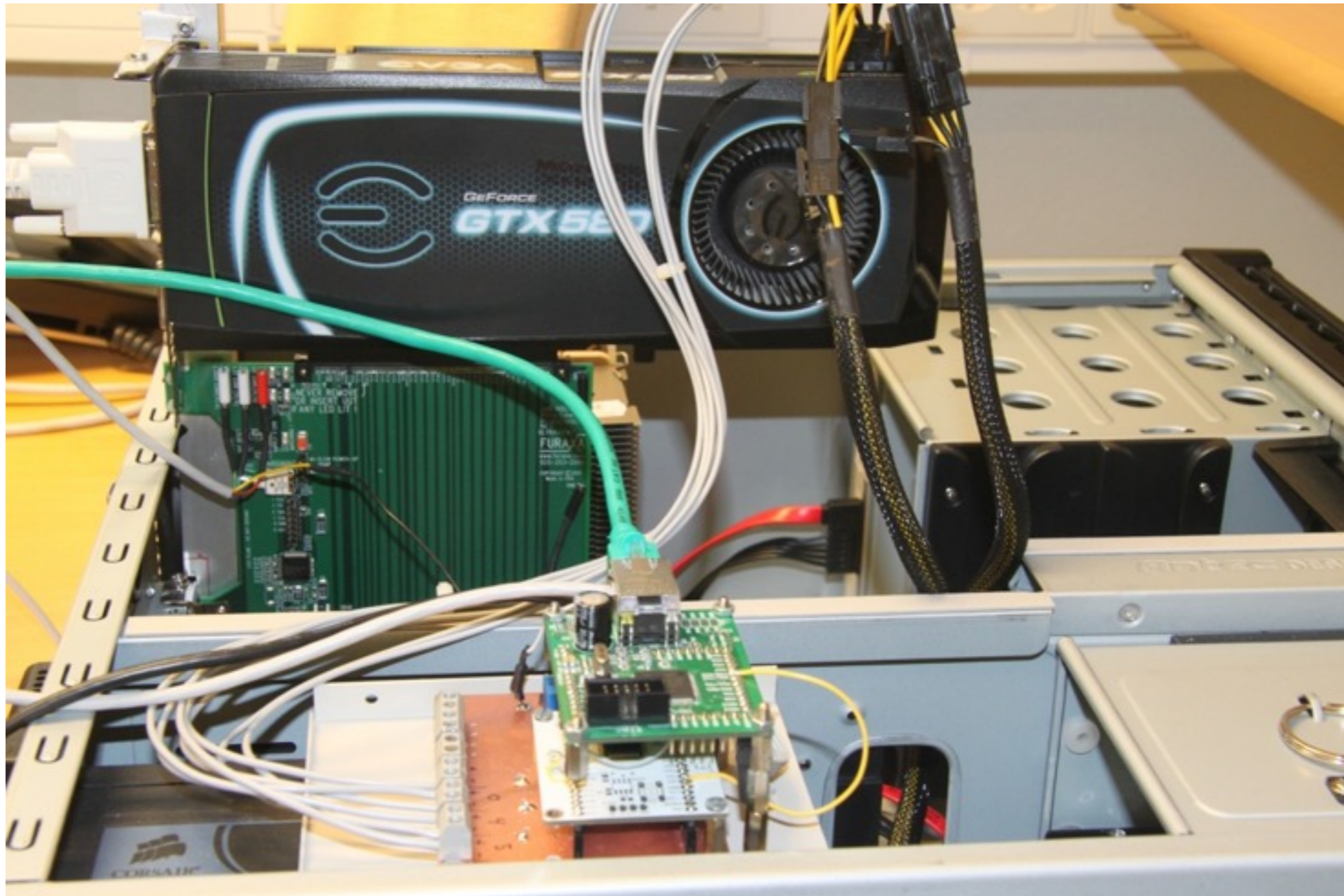


Discrete graphics cards

- Connected on the PCI-Express bus
 - PCI-Express bus provides 75W
 - Measured through an Ultraview PCIeEXT-16HOT expander card
 - PCI-Express power connectors
 - 8-pin provides $\leq 150W$
 - 6-pin provides $\leq 75W$



Discrete graphics cards



Discrete graphics cards

- AMD Radeon 7970 (28nm)
- NVIDIA GeForce GTX 580 (40nm)



-
- Intel Sandy Bridge, HD3000 GPU (32nm)
 - Connected on motherboards 4-pin power connector
 - Provides power to CPU, GPU, and parts of the memory system
 - Two runs
 - Rendering pass
 - Idle pass, all gl* calls removed from code

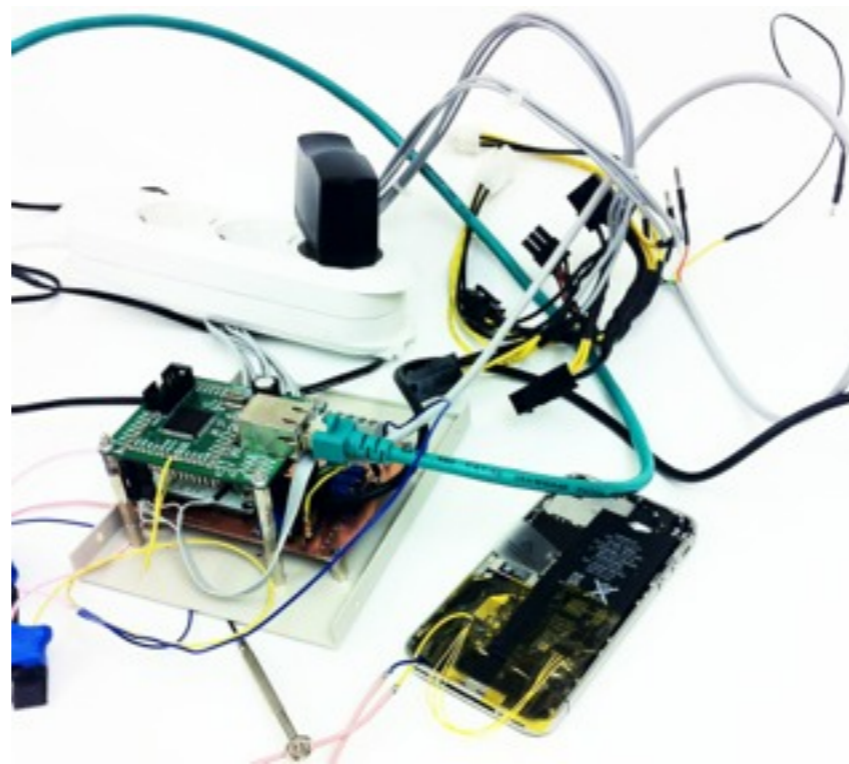


-
- Rendering pass - idle pass = ?
 - GPU power
 - Parts of the memory power
 - Memory bandwidth generated from the graphics workload
 - Not including memory refresh power
 - CPU power for driver execution

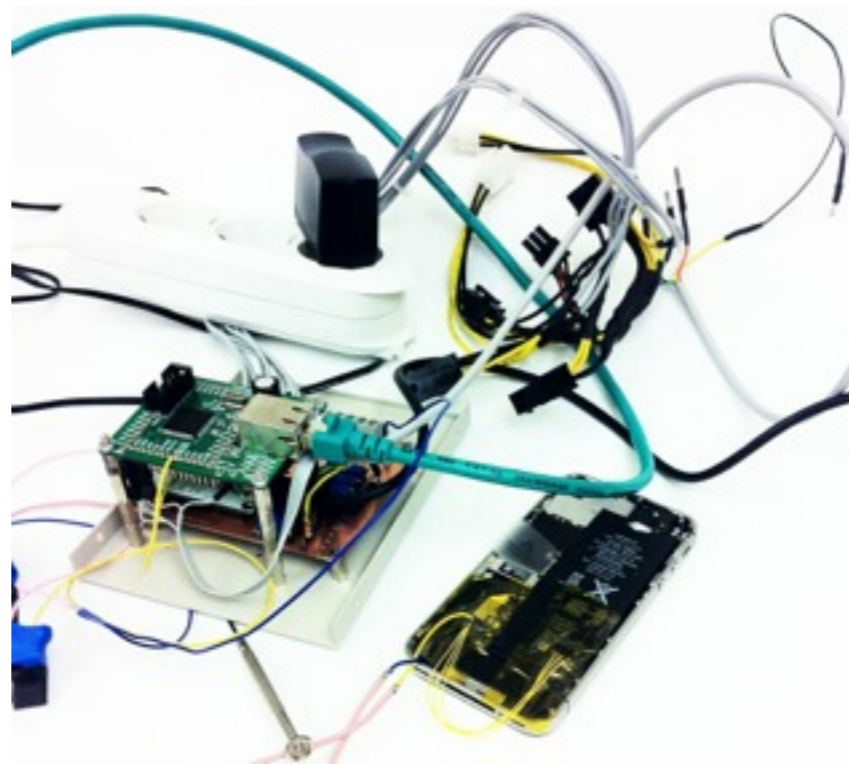


SoC #2

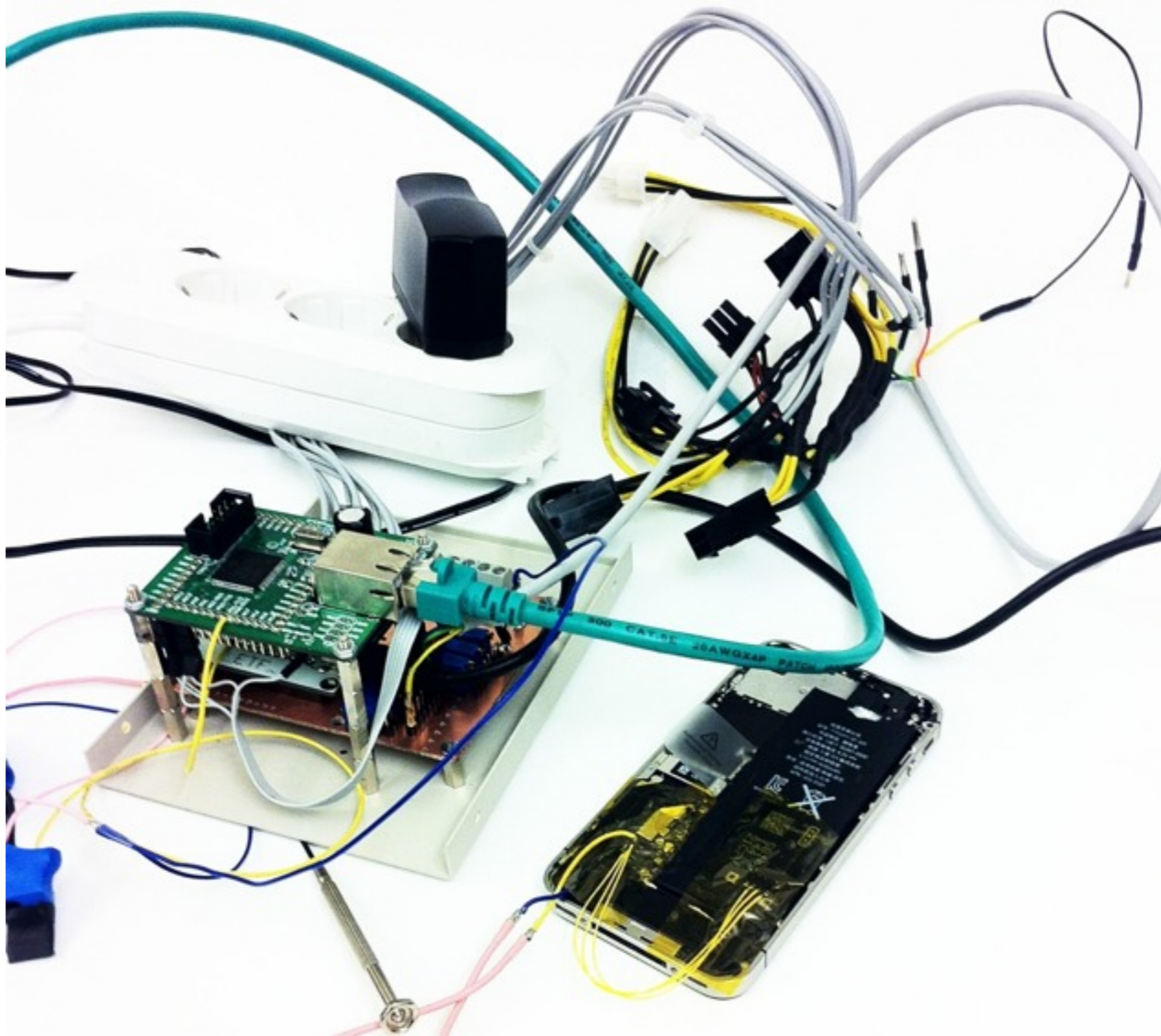
- iPhone 4S, PowerVR SGX543MP2 GPU (45nm)
 - Connected on battery connectors
 - Provides power to everything
 - Two runs
 - Rendering pass
 - Idle pass



- Rendering pass - idle pass = ?
 - GPU power
 - Memory power
 - Only for memory bandwidth generated from the graphics workload
 - Not including memory refresh power
 - CPU power for driver execution



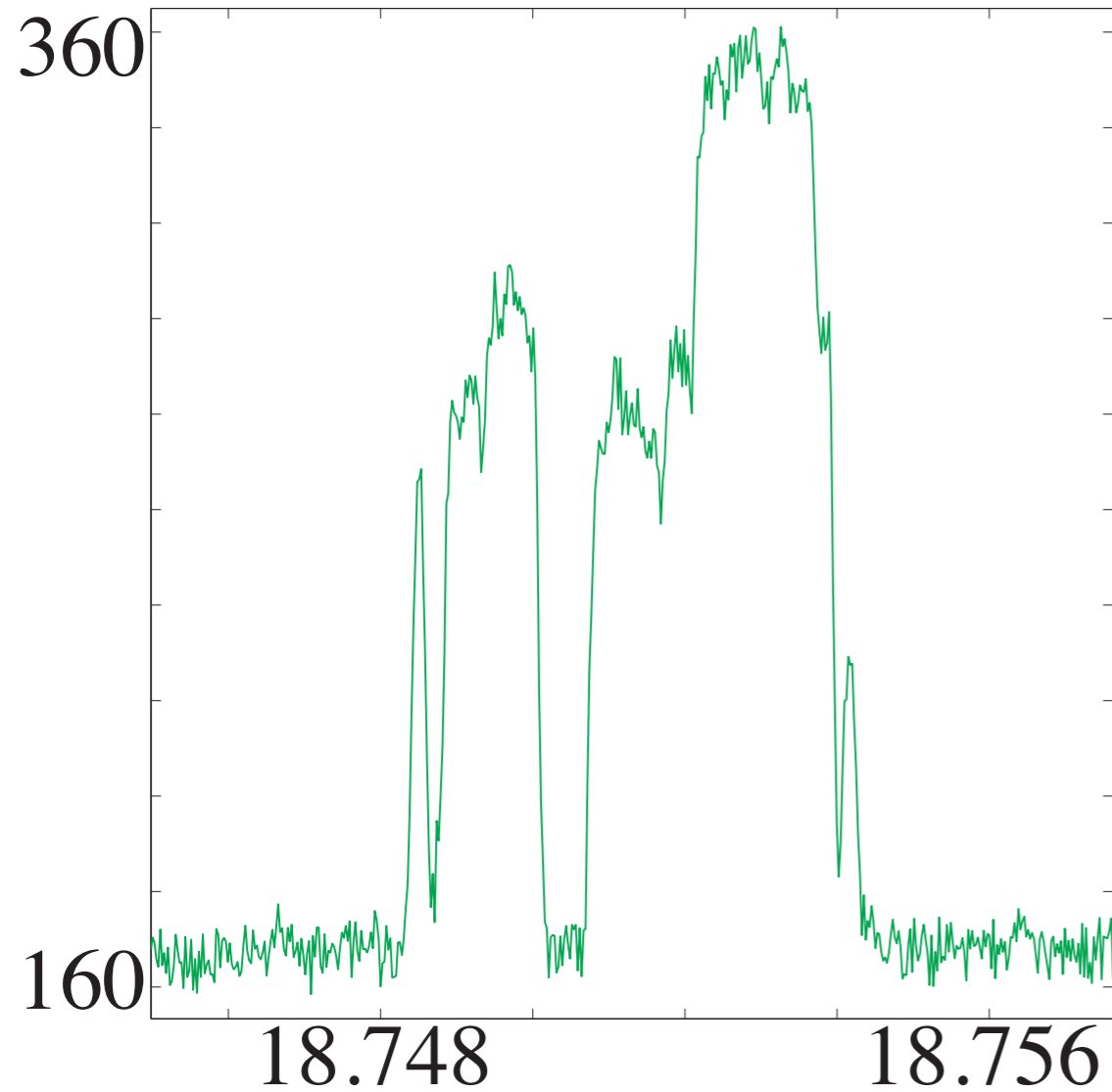
SoC #2



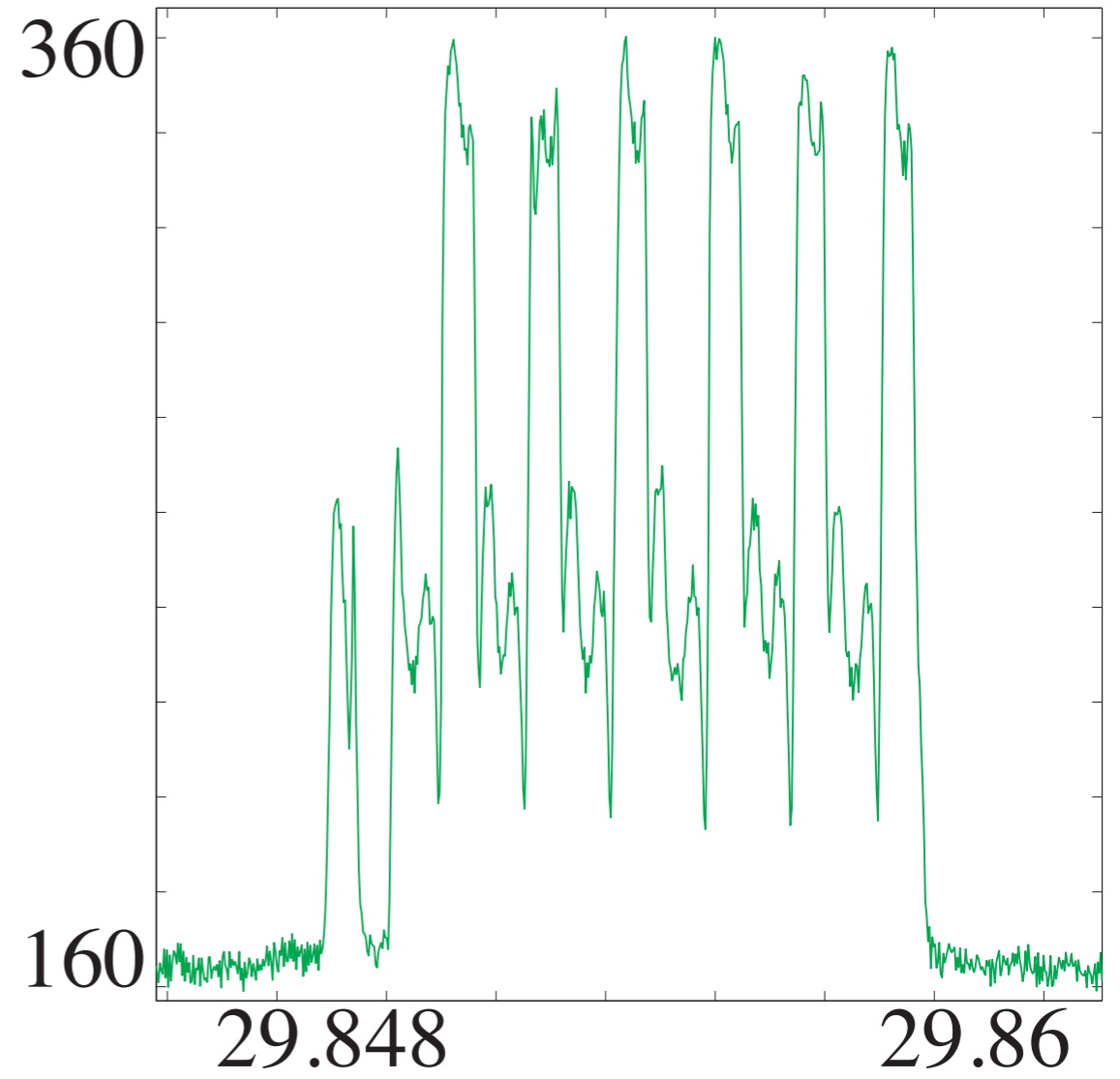
-
- What we measure:
 - High-frequency power data (40kHz)
 - Rendering time per frame



GeForce GTX 580



Deferred shading



Variance shadow maps
6 light sources



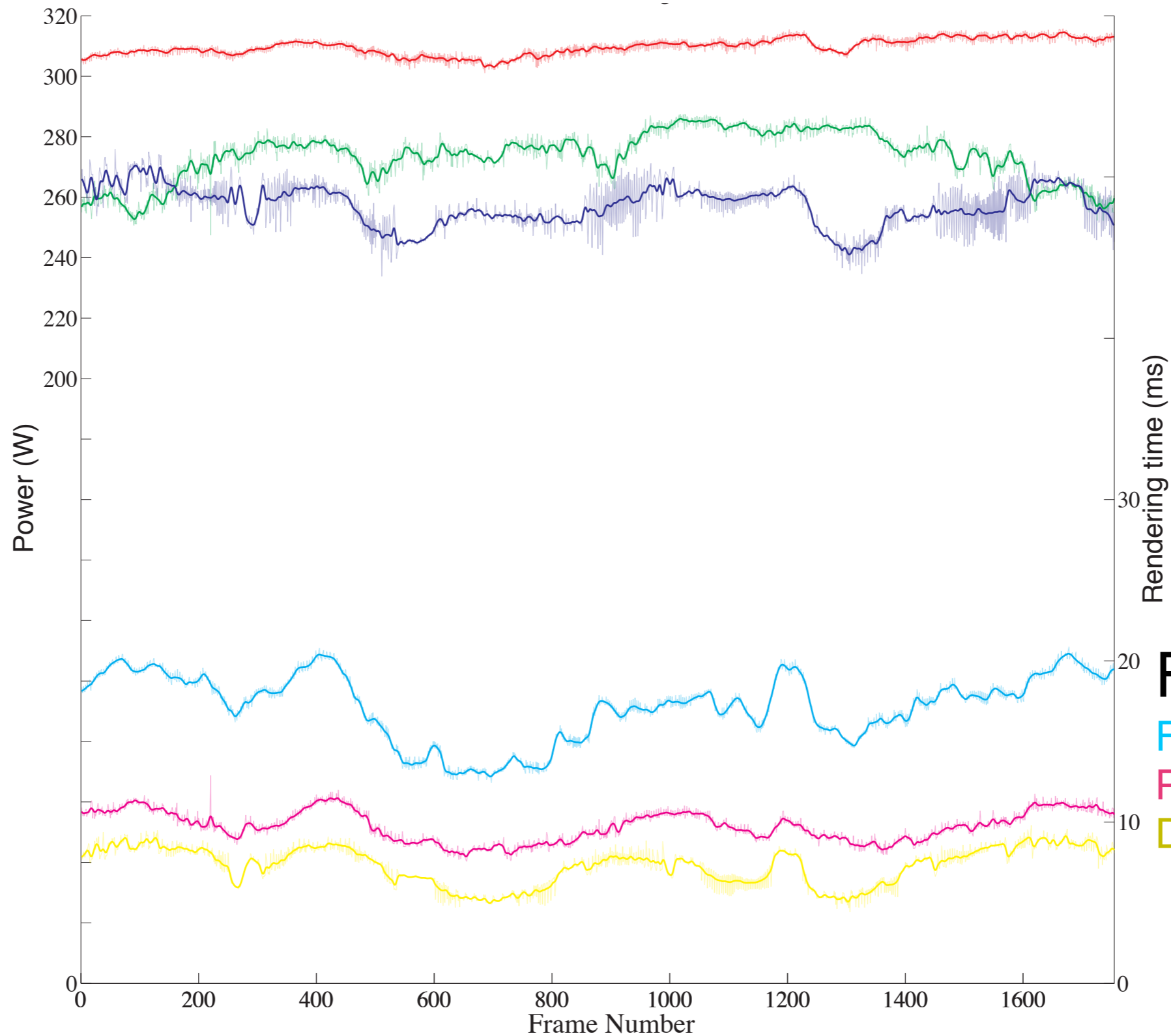
GeForce GTX 580: Primary rendering

Power

Forward

Pre-Z

Deferred



Frametime

Forward

Pre-Z

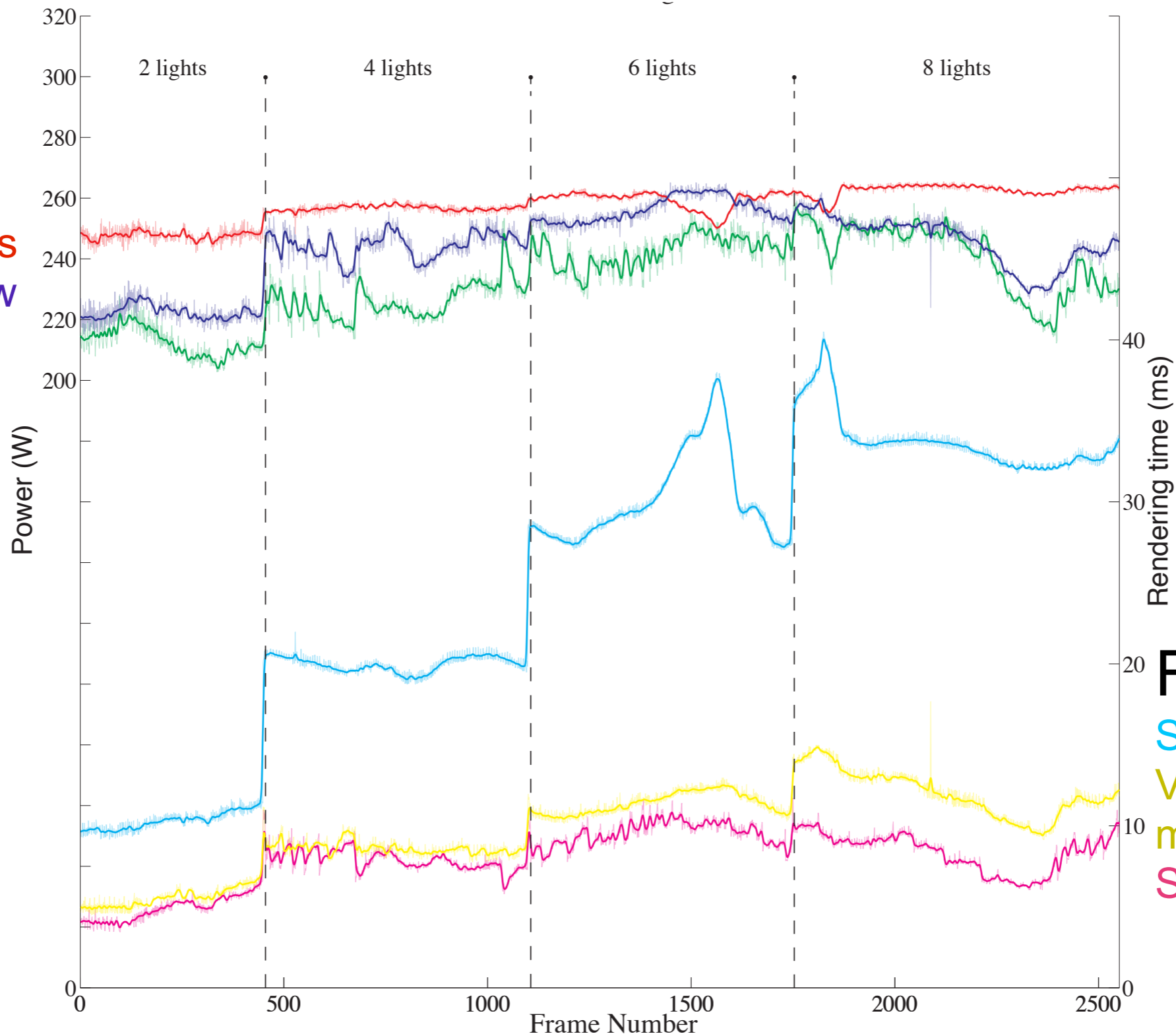
Deferred



GeForce GTX 580: Shadows

Power

- Shadow volumes
- Variance shadow maps
- Shadow maps

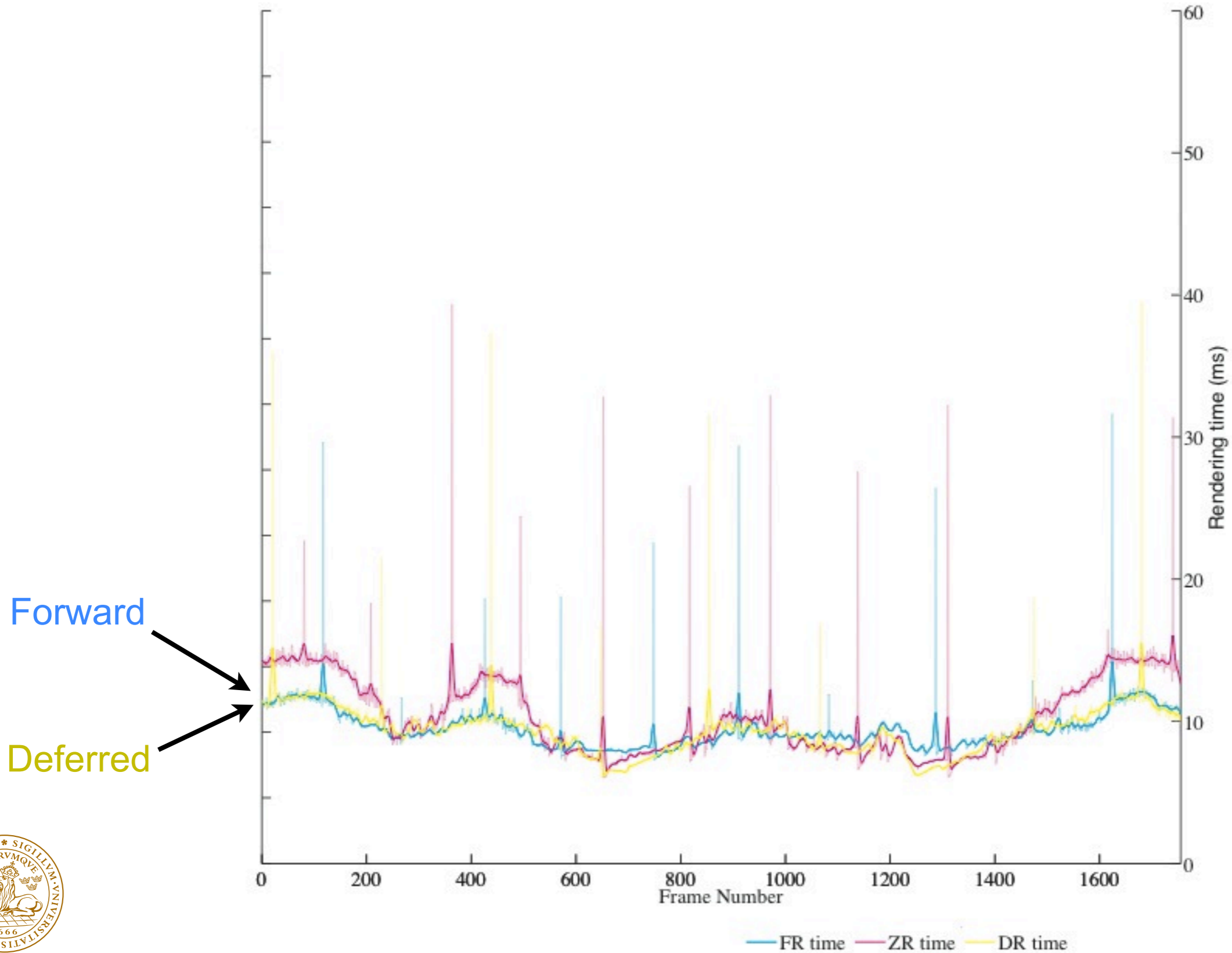


Frametime

- Shadow volumes
- Variance shadow maps
- Shadow maps



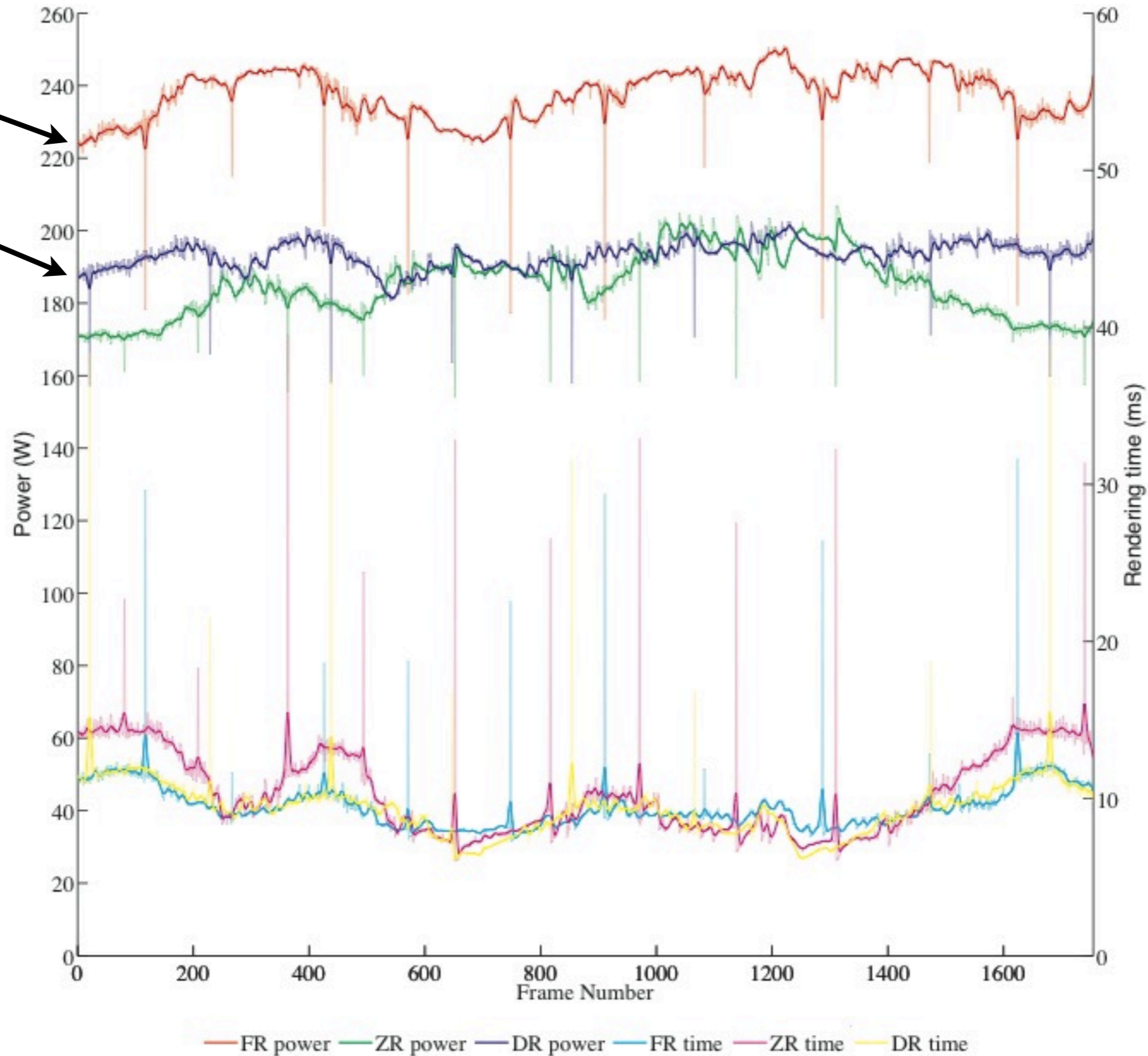
AMD Radeon 7970: Primary rendering



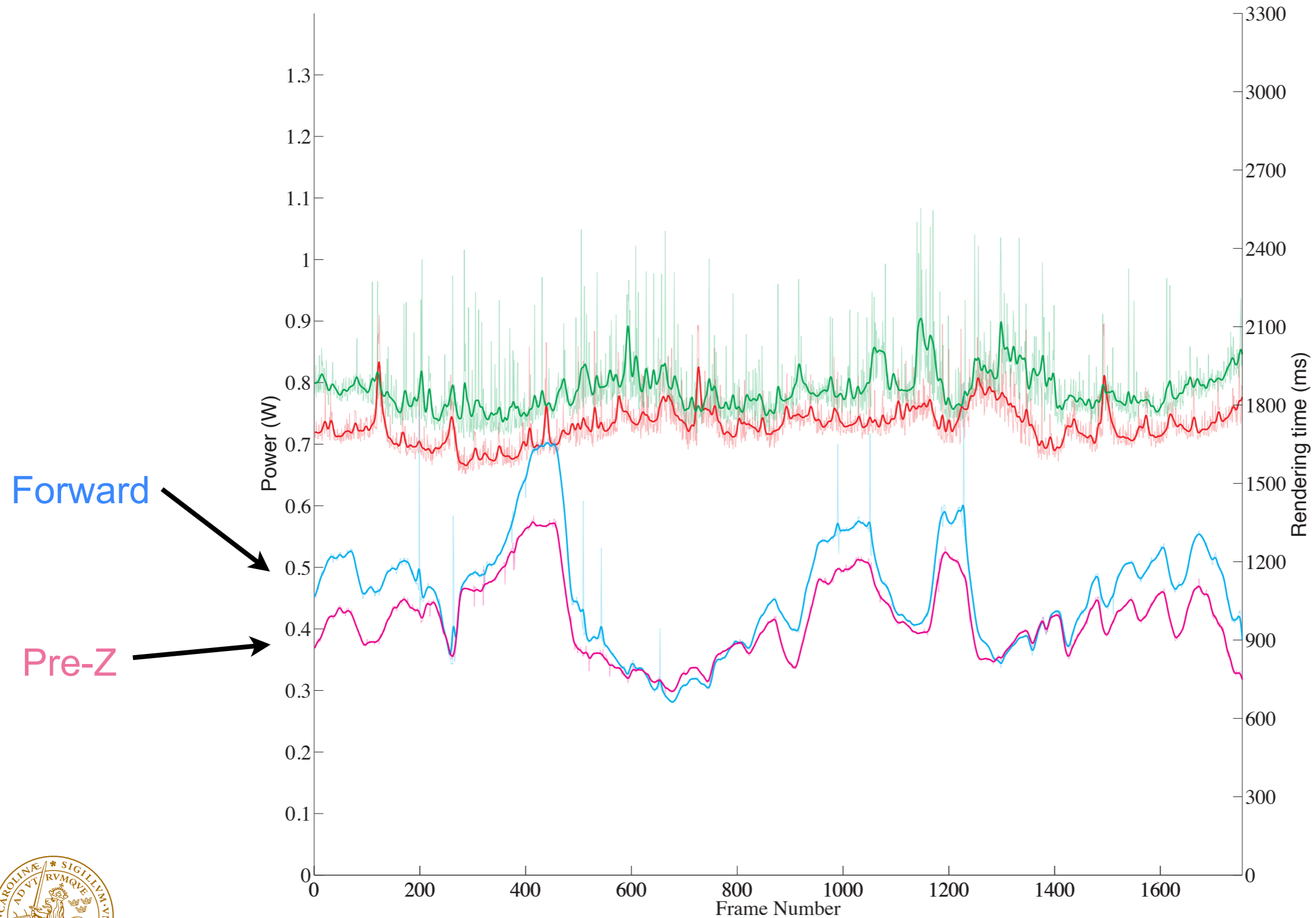
AMD Radeon 7970: Primary rendering

Forward

Deferred



iPhone 4S



iPhone 4S

- Higher energy than expected
 - Mostly due to long rendering times
 - Probably pushing sort-middle to flush buffers

Forward

Pre-Z



Results

-
- What numbers are we interested in?



Results

-
- What other numbers are we interested in?
 - Power



Results

-
- What other numbers are we interested in?
 - Power
 - Energy



-
- What other numbers are we interested in?
 - Power
 - Energy
 - More interesting for battery lifetime



-
- What other numbers are we interested in?
 - Power
 - Energy
 - More interesting for battery lifetime
 - But what metric should we use?



Metric

nJ/pixel



nJ/pixel

- Largely resolution independent



nJ/pixel

- Largely resolution independent
- Easy to divide frame into segments



nJ/pixel

- Largely resolution independent
- Easy to divide frame into segments
- Easy to calculate fillrate for a given TDP



Energy, nJ/pixel

	Primary visibility		
	FR	ZR	DR
GTX 580	1,443	722	511
Radeon 7970	607	512	489
Sandy Bridge	872	314	280
iPhone 4S	2,234	2,015	-



Energy, nJ/pixel

	Shadows		
	SV	SM	VSM
GTX 580	1,325	447	532
Radeon 7970	953	469	804
Sandy Bridge	1,317	311	511
iPhone 4S	-	461	-



-
- Not possible to estimate power / energy only from frame times
 - Surprisingly, pre-Z proved useful on a tiled, deferred, sort-middle architecture for our application
 - Pushing too much geometry?
 - Still similar energy per pixel, despite rendering times that differ by an order of magnitude or more
 - nJ/pixel



Thanks for listening

Any questions?

