Presenter: Hannes Mohr - KIT, IPE

Contributors: L. Ardilla Perez, M. Balzer,
M. Caselle, S. Chilingaryan, T. Dritschler,
A. Kopmann, L. Rota, T. Schuh, M. Weber

**KIT**

Karlsruhe Institute of Technology



# EVALUATION OF GPUS FOR HIGH–LEVEL TRIGGERS IN HIGH ENERGY PHYSICS

▸ Implement track trigger using GPUs

▸ Use established methods for seeding

▸ Present our own version of the Hough transformation

▸ Compare different GPUs/vendors

▸ Investigate data transfer/latencies

▸ Estimate impact of technological advances

**OUR GOAL IS TO ACHIEVE COMPETITIVE RESULTS, WHILE GAINING FLEXIBILITY**

Key:
- Muon
- Electron
- Charged Hadron (e.g. Pion)
- Neutral Hadron (e.g. Neutron)
- Photon

4T

2T

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoid

Iron return yoke interspersed with Muon chambers

Transverse slice through CMS

D.Barney, CERN, February 2004

# CMS DETECTOR

*Image CMS Collaboration

# BASELINE GEOMETRY – 6 LAYERS 5 DISKS

# SILICON TRACKER

Peak luminosity ● — Integrated luminosity

*Image CMS Collaboration

▸ **Current CMS trigger won't be able to handle:**

  ▸ Increased data rates

  ▸ Increased pile-up

▸ **Currently proposed solution:**

  ▸ Data reduction on detector

  ▸ Raise latency of trigger from 3.4 to 12.5 µs

  ▸ L1 track trigger

**Front End**

**Back End**

| Detector Hit | ≈$10^4$Tb/s → | Stub Building | 50Tb/s → | L1 Tracking | 1Tb/s → | Global Trigger |

▸ Readout at 40 MHz, BX every 25 ns

▸ 6 μs each for L1 Trigger and Global Trigger

▸ L1 Tracking to combine Track seeding and Fitting

hit matching windows (discrete)

pass — fail

~1.6-4mm

~0.1mm

stub

▸ Applies momentum cut to hits

▸ Delivers estimate on track bend

▸ Drastically decreases number of hits by a factor of 100

▸ Associative Memory approach (ASICs)

▸ Time-multiplexed FPGA Hough transformation

▸ …

## CURRENT APPROACHES USE SPECIALIZED HARDWARE

# COMPARISON GPU VS. FPGA

Nvidia Tesla K40c vs. XILINX VIRTEX-7 XC7VX1140T (both 28nm)



**GPU** **FPGA**

**I/O BW [GB/S]**
GPU: 16
FPGA: 348

**MEMORY BANDWIDTH [GB/S]**
GPU: 288
FPGA: 1,8

**GFLOPS***
GPU: 5040
FPGA: 850

**GFLOPS/WATT***
GPU: 21,44
FPGA: 20,8

**GFLOPS/$***
GPU: 1
FPGA: 0,24

0    0,25    0,5    0,75    1

**\*SINGLE PRECISION, FLOATING POINT**

| GPU | FPGA |
| --- | --- |
| Rapid development cycles and high flexibility | Huge I/O Bandwidth |
| Large bandwidth to external memory | Deterministic timings/ runtimes |
| High Floating-point performance | High bit-level performance |

Equivalent to FPGA approach, used by collaboration of KIT and UK Track Trigger Group*

▸ Uncompress data

▸ Perform Hough transformation

➤ Uses module bend information

▸ Apply layer condition

▸ Reject or return track candidates

Calculate possible

$(\varphi_0, q/pt)$ pairs for each hit

Make histogram in

Hough space

Calculate possible

($\varphi_0$, q/pt) pairs for each hit

Make histogram in

Hough space

Calculate possible

$(\varphi_0, q/pt)$ pairs for each hit

Make histogram in

Hough space

TTBar event - PU 140

TTBar event - PU 140

GPU  implementation specifics:

▸ Optimized for minimum latency

▸ Computes $q/p_t$-bins in parallel

▸ Almost no dependence on number of stubs

▸ Kernel scheduling

▸ Kernel launch time

▸ Allocation of shared memory

▸ …?

Invocation and setup of kernels is too costly,

we need to keep it running continuously

➡ **NEEDS SPINNING KERNEL**

Stubs

CURRENTLY NOT POSSIBLE IN OPENCL:
CACHE CAN'T BE FLUSHED FROM KERNEL

# WHAT ABOUT DATA TRANSFER?

(Red) Conventional transfer   (Green) RDMA transfer

KERNEL RUNS
CONTINUOUSLY

PRELOAD DATA

1 GPU:request data

FPGA

GPU

KERNEL RUNS
CONTINUOUSLY

PRELOAD DATA

1 GPU:request data

FPGA

GPU

2 FPGA:Data Transfer

SIMULTANEOUSLY

2 GPU:Poll for data

KERNEL RUNS
CONTINUOUSLY

PRELOAD DATA

1 GPU:request data

FPGA

GPU

2 FPGA:Data Transfer

3 GPU:Computes

SIMULTANEOUSLY

2 GPU:Poll for data

At the moment we don't write back into the FPGA,

read and write 160 stubs (64 bits each)

➤ Start transfer ➤ Poll for data ➤ Write back result



* ···· estimated response time

➤ Read/Uncompress data ➤ Compute ➤ Poll



160 stubs, 1 sector

▸ Computation time is higher than data transfer

▸ We can hide the transfer behind the computation

## INTERLEAVED APPROACH

▸ Start data transfer for current data set

▸ Do calculations on previous dataset (lies in register memory)

▸ Poll new data (should take less time)

## INCREASES THROUGHPUT AT COST OF LATENCY

Data older, throughput higher

➤ (poll) Read/Uncompress data ➤ Ask for data ➤ Compute



* ···· estimated response time

▶ Hexagonal bins in hough space

▶ Suppresses fake candidates

▶ Runtime comparable

▶ only 1 possible bin per row

▶ less algorithmic branching

CAVEAT: NEEDS MORE BINS
(FACTOR OF 2)

Results for TTBar Dataset PU140, whole detector, 1 event

➤ (poll) Read/Uncompress data ➤ Ask for data ➤ Compute



latency limit

* ···· estimated response time

Performance:

▸ Computational time of around 4 μs

▸ Transfer time of around 2 μs

Surpassed our expectations

Development is faster

More complex algorithms are possible:

▸ Example: hexagonal approach

Data transfer using standard interfaces is challenging

▸ Need to process multiple sectors per card in future

▸ Look at performance of newer cards

    ▸ High Bandwidth Memory,

      already in consumer model cards,

      promises 2-4x better throughput

▸ Investigate new transfer technologies

    ▸ PCIe 4.0 (2x faster)

    ▸ nv-link (5-10x faster)

## MOORS LAW IS OUR FRIEND!

# THANK YOU!

# QUESTIONS?

TTBar PU140 CUDA 7.5, spinning kernel

TTBar PU140 CUDA 8, spinning kernel