

# Track Reconstruction on GPUs for the Mu3e Experiment

Dorothea vom Bruch  
for the Mu3e Collaboration

DPG Frühjahrstagung, T41: Detektoren und DAQ 1

March 10, 2015

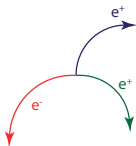


Physikalisches Institut Heidelberg



Mu3e searches for the charged lepton-flavour violating decay  $\mu^+ \rightarrow e^+ e^+ e^-$  with a sensitivity better than  $10^{-16}$

$\mu$  decays at rest  $\rightarrow E_e < 53 \text{ MeV}$



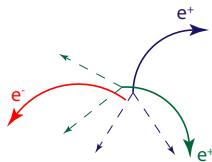
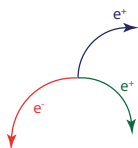
## Signal

- Coincident in time
- Single vertex
- $\Sigma \vec{p}_i = 0$
- $E = m_\mu$



Mu3e searches for the charged lepton-flavour violating decay  $\mu^+ \rightarrow e^+ e^+ e^-$  with a sensitivity better than  $10^{-16}$

$\mu$  decays at rest  $\rightarrow E_e < 53 \text{ MeV}$



## Signal

- Coincident in time
- Single vertex
- $\Sigma \vec{p}_i = 0$
- $E = m_\mu$

## Random Combinations

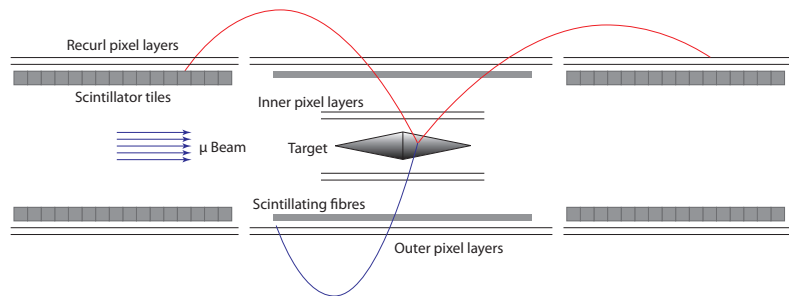
- Not coincident in time
- No single vertex
- $\Sigma \vec{p}_i \neq 0$
- $E \neq m_\mu$

# The Mu3e Detector

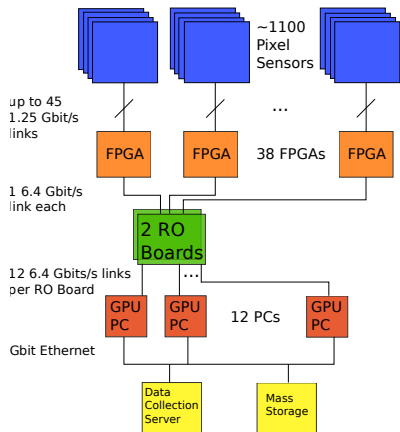


## Requirements:

- Excellent momentum resolution:  $< 0.5 \text{ MeV}/c$
- Good timing resolution: 100 ps for tiles, 1 ns for fibres,  $< 20 \text{ ns}$  for pixels
- Good vertex resolution:  $300 \mu\text{m}$
- High rates  $\mathcal{O}(10^8 - 10^9 \mu/s)$



# Readout Scheme



- Triggerless readout → 1 Tbit/s data rate
- Online data reduction required
- Track reconstruction and vertex fitting on GPUs
- ⇒ Reduction factor of ~ 1000



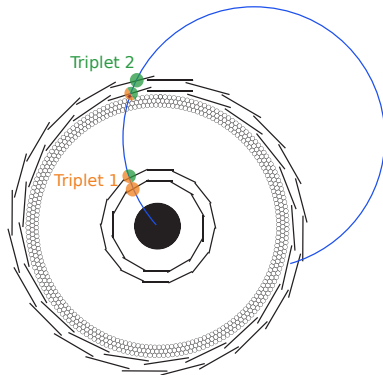
## Low Momentum

Electrons: 15 - 53 MeV

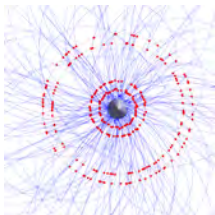
Resolution dominated by multiple Coulomb scattering

- Ignore hit uncertainty
- Describe track as sequence of hit triplets
- Multiple scattering at middle hit of triplet
- Minimize multiple scattering

$$\chi^2 = \frac{\phi_{MS}^2}{\sigma_{MS}^2} + \frac{\theta_{MS}^2}{\sigma_{MS}^2}$$



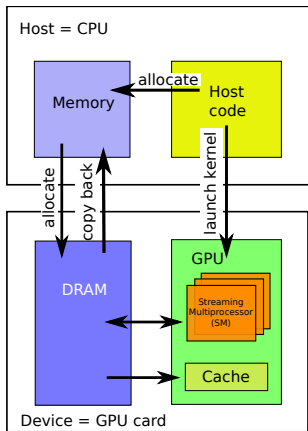
# Fit on the GPU



- Number of possible track candidates  $\sim n^3$
- With  $n$ : # hits per layer
- On GPU: Loop over all possible combinations
  - Geometrical selection cuts
  - Triplet fit
  - Vertex fit

⇒ Reduction factor of  $\sim 1000$

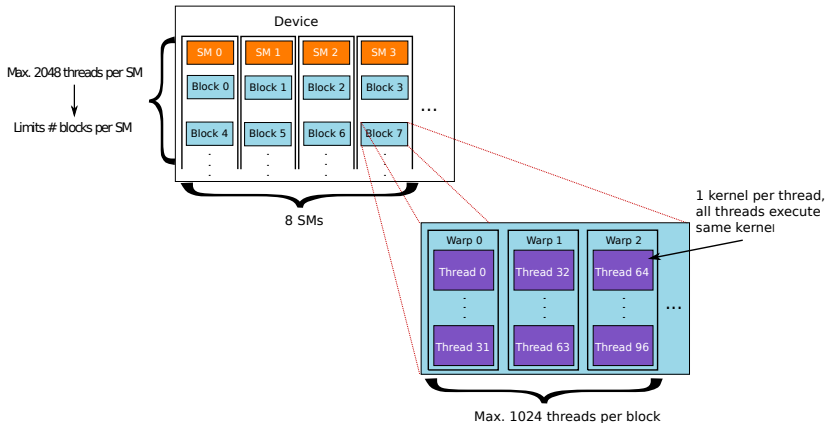
# GPU Properties



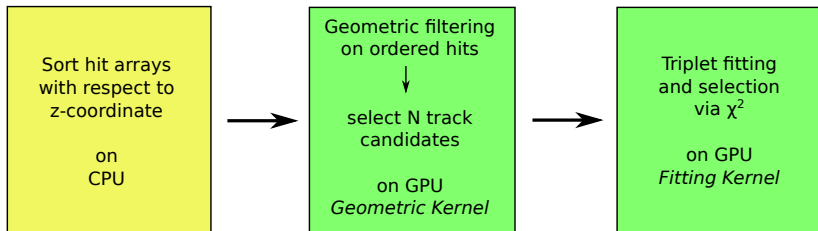
- Highly parallel structure
- Process large blocks of data
- Nvidia: API extension to C: CUDA (Compute Unified Device Architecture)



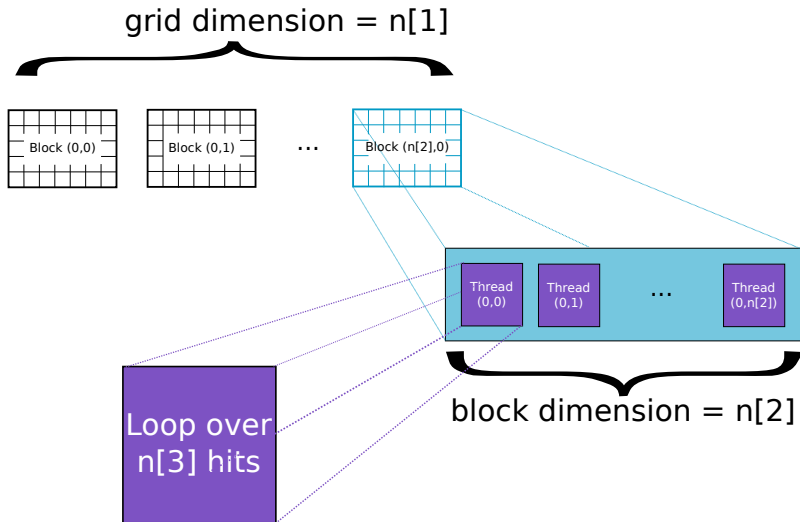
# GPU Architecture (GTX 680 as example)



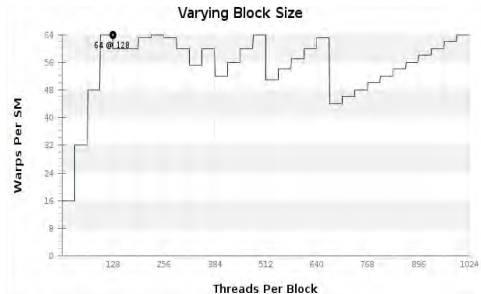
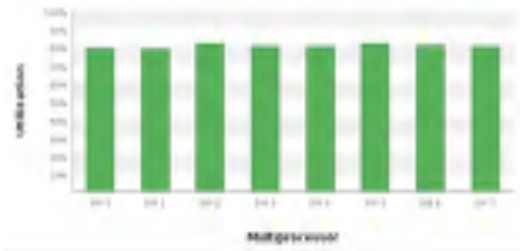
# Main Steps



# Grid for Geometric Kernel



# Compute Utilization for Fitting Kernel





- Process  $1.4 \times 10^{10}$  triplets / s
- Most time spent on geometric kernel  $\rightarrow$  outsource to FPGA
- Ratio of copying data from CPU to GPU to compute time: 40 %, will improve when selection cuts are applied on FPGA
- For  $10^8 \mu/s$ :  $10^{12}$  hit combinations  $\rightarrow$  48 GPU computers
- For  $10^9 \mu/s$ : More improvements needed



- Study data transmission from FPGA to GPU:

- $\text{FPGA} \xrightarrow{\text{PCIe, DMA}} \text{CPU} \xrightarrow{\text{PCIe}} \text{GPU}$

- $\text{FPGA} \xrightarrow{\text{PCIe, DMA}} \text{GPU}$

- Outsource selection to FPGA
- Include vertex fit

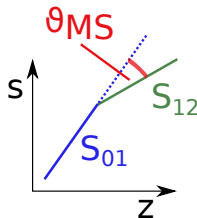
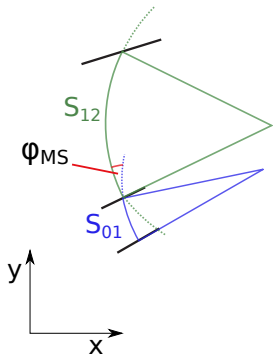


## Backup Slides

# Multiple Scattering at middle hit of triplet



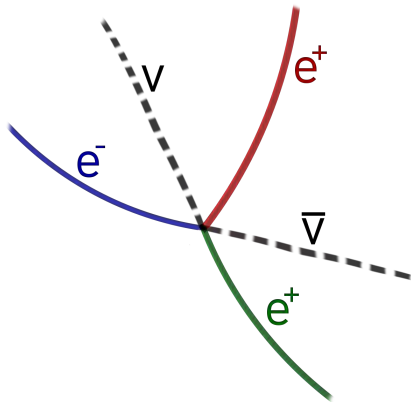
Multiple scattering at middle hit of triplet



$$\text{Minimize } \chi^2 = \frac{\phi_{MS}^2}{\sigma_{MS}^2} + \frac{\theta_{MS}^2}{\sigma_{MS}^2}$$



# Irreducible Background



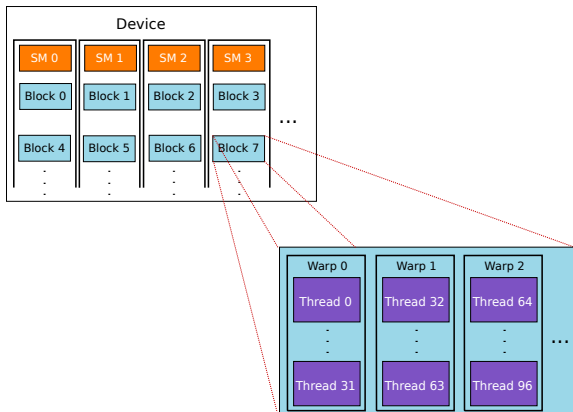


| Phase 1                               | Phase 2                                 |
|---------------------------------------|---|
| $\mathcal{O}10^8 \mu/s$               | $\mathcal{O}10^9 \mu/s$                 |
| Central part<br>+<br>1 recurl station | Central part<br>+<br>2 recurls stations |

# Hardware Implementation: Warps



- After block is assigned to SM  
→ Division into units called warps
- On GTX 680: 1 warp = 32 threads

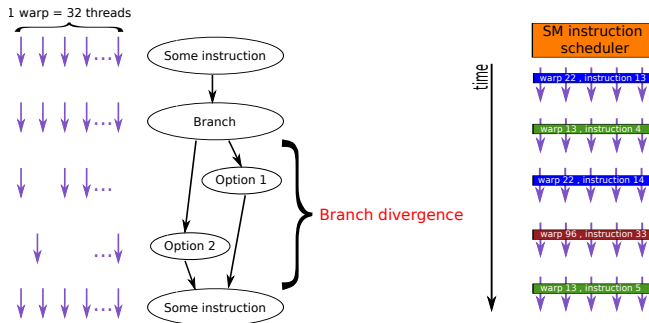


# Warp Scheduling

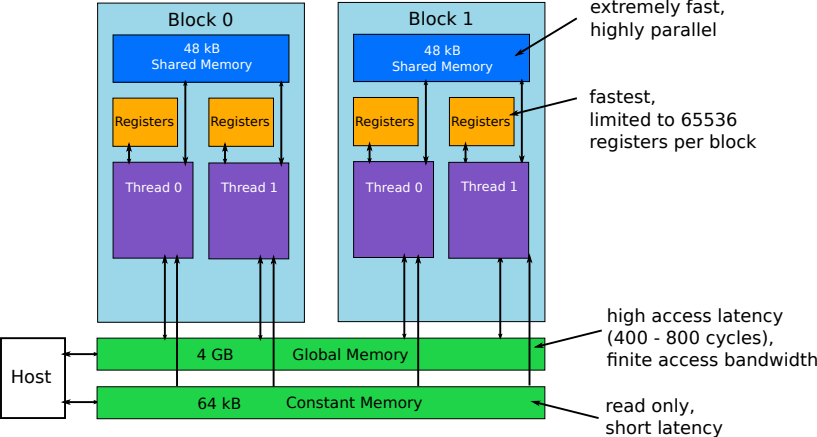


## Warps execute

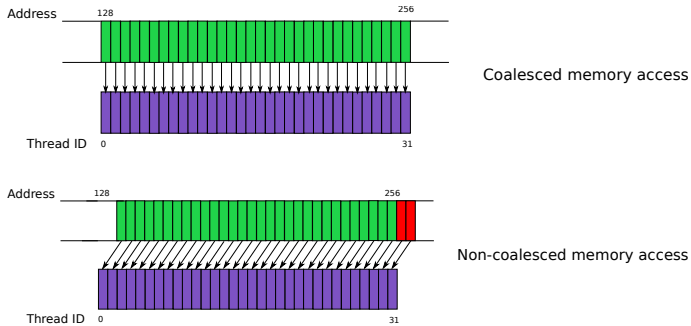
In SIMD fashion (Single Instruction, Multiple Data)  
Not ordered



# GPU Memory



# Memory Access



## Warp Memory Access

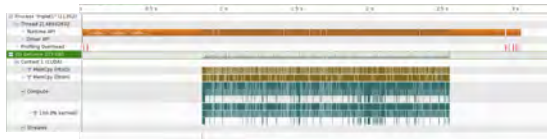
128 bytes in single transaction

# Diagnostic Tools (One Kernel)



Nvidia provides several diagnostic tools:

- Profiler for terminal usage: Time spent by CPU and GPU
- Memory check: Memory allocation errors, misuse, ...
- Visual profiler: Diagnostics for performance of GPU code



| GPU utilization 100% 100% |                           |
|---------------------------|---------------------------|
| GPU utilization 100% 100% |                           |
| Device                    | 1.0000 (1.000000000000)   |
| Attributes                |                           |
| - Compute Capability      | 3.0                       |
| - Memory                  |                           |
| Shared Memory per Block   | 32000                     |
| Registers per Block       | 65536                     |
| Grid Dimension            | 1 (128x128x1 0x0x0 0x0x0) |
| Blocks Dimension          | 1 (128x 128x 88)          |
| Waves per Multiprocessor  | 32                        |
| Blocks per Multiprocessor | 24                        |
| - Multiprocessor          |                           |
| Multiprocessor            | 8                         |
| Clock Rate                | 1.137 GHz                 |
| Execution Model           | SM                        |
| Max PTX                   | 3                         |
| Threads per Warp          | 32                        |
| - Memory                  |                           |
| Global Memory Bandwidth   | 780.0 GB/s                |
| Global Memory Size        | 8 GB                      |
| Global Memory Type        | DDR5                      |
| L2 Cache Size             | 512 KB                    |
| Memory Engines            | 1                         |
| - P2C                     |                           |
| Streams                   | 8                         |
| GPU APIs                  | 5.0000                    |
| L2C APIs                  | 38                        |
| - Environment             |                           |

# Diagnostic Tools (One Kernel)



## Kernel profile: Cuda code and machine code

