



Preliminary

Performance Profiling  
of madgraph4gpu  
with Nvidia A100 GPU

Carl Vuosalo

Harshit Sharma

University of Wisconsin-Madison



# Profiling Overview--Previous

- Performance profiling to find parameters for highest performance
- Platform at T2\_US\_Wisconsin:
  - CentOS 7.9.2009 (Core)
  - 125.74 GiB RAM
  - 1x AMD EPYC 7402P 24-Core Processor, 48 threads
  - GPU 2x TU104GL [**Tesla T4**] (only 1 GPU used)
- Physics process to profile: g g  $\rightarrow$  t t g g
- Using [Nsight Compute](#) tool from Nvidia to perform profiling
- Majority of time (~90%) spent on SigmaKin kernel that does matrix element (ME) calculation (not counting I/O)



# Profiling Overview--New

- centos:7 Docker container:
  - NVIDIA-Linux-x86\_64-525.116.04
  - cuda\_12.1.1\_530.30.02\_linux
  - HEAD of <https://github.com/madgraph5/madgraph4gpu.git>
  - epochX/cudacpp/gg\_ttgg.sa/SubProcesses/  
P1\_Sigma\_sm\_gg\_ttxgg
  - AVX=avx2 FPTYPE=d HELINL=0
  - Nvidia A100-SXM4 with 40GB global GPU memory
- Using [Nsight Compute](#) tool from Nvidia to perform profiling



# Tesla T4 Performance -- Old

- Check gg\_ttgg GPU performance with varied max registers
- Double precision
- Number of blocks = 65536
- Number of iterations = 12
- Best performance with 160 max registers
- Memory usage:
  - 369 MB for 32 threads/block
  - 873 MB for 128 threads/block

32 threads per block	128	140	160	180	216	240	256
Events/s (ME)	1.13e4	1.12e4	2.30e4	1.11e4	1.11e4	1.11e4	1.11e4

128 threads per block	128	140	160	180	216	240	256
Events/s (ME)	1.91e4	1.90e4	2.31e4	1.89e4	1.90e4	1.89e4	1.90e4



# A100 Performance

- Double precision
- Number of blocks = 65536
- Number of iterations = 1

**Preliminary**

256 threads per block	64 blocks/ grid						
Events/s (ME)	1.02e5						