

Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release

Taylor Childers Walter Hopkins Nathan Nichols



Laurence Field Stephan Hageboeck

Stefan Roiser David Smith Jorgen Teig Andrea Valassi Zenny Wettersten

CERN

Olivier Mattelaer



Carl Vuosalo



CHEP, Norfolk, VA, 08 May 2023

https://indico.jlab.org/event/459/contributions/11829/

Motivation: Monte Carlo Event Generators in WLCG computing

• HL-LHC computing needs expected to outgrow resource growth

-Need R&D on software to improve efficiency and port it to new resources, such as GPUs at HPC centres



- MC generators projected to use 10% 20% of ATLAS/CMS WLCG CPU budget
 - -Many ways to speed them up see the HEP Software Foundation (HSF) Generator WG review paper
 - Ideal candidates to exploit data parallelism in GPUs (SIMT) and CPUs (SIMD)?



Motivation: Monte Carlo Event Generators in WLCG computing

• HL-LHC computing needs expected to outgrow resource growth

-Need R&D on software to improve efficiency and port it to new resources, such as GPUs at HPC centres



- MC generators projected to use 10% 20% of ATLAS/CMS WLCG CPU budget
 - -Many ways to speed them up see the HEP Software Foundation (HSF) Generator WG review paper
 - Ideal candidates to exploit data parallelism in GPUs (SIMT) and CPUs (SIMD)?



What is a MC ME generator? A simplified computational anatomy

Monte Carlo sampling: randomly generate and process MANY different events ("phase space points")



This can be parallelized (SIMT/SIMD and multithreading)



(NB: "Matrix Element" is an element of the **scattering matrix**; not the linear algebra concept!)

Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release



ME generator data parallelism: design for lockstep processing!

- In MC generators, the same function is used to compute the Matrix Element for many different events
 - -<u>ANY</u> matrix element generator is a good fit for lockstep processing on GPUs (SIMT) and vector CPUs (SIMD)
 - -Data parallelism strategy in madgraph4gpu is event-level parallelism (many events = many phase space points)





Université catholique



- MG5aMC production version is in Fortran
 - -Outer shell: Madevent (random sampling, integration and event generation + I/O, multi-jet merging...)
 - –Matrix Element (ME) core
 - MG5aMC generates ME code for each physics process
 - MEs may take >95% of the CPU time for complex processes (e.g. $gg \rightarrow t\bar{t}ggg$)
 - High potential for acceleration on CPUs / GPUs



MG5aMC and the madgraph4gpu project

- madgraph4gpu: speed up Matrix Element calculation in MG5aMC on GPUs and vector CPUs
 - -Collaboration of theoretical/experimental physicists with software engineers born in the HSF generator WG
 - –More history and details in <u>vCHEP2021</u> and <u>ICHEP2022</u>

EPJ Web of Conferences 251 , 03045 (2021) CHEP 2021	https://doi.org/10.1051/epjconf/20212510304
Design and engineering of a for the MG5aMC event generation of the MG5aMC event generation of the transmission of transmission	simplified workflow execution ator on GPUs and vector CPUs
Andrea Valassi ^{1,*} , Stefan Roiser ^{1,} , Olivier M	attelaer ² , and Stephan Hageboeck ¹
¹ CERN, IT-SC group, Geneva, Switzerland ² Université Catholique de Louvain, Belgium	
https://doi.org/	(10 1051/opicopf/202125102015

https://doi.org/10.1051/epjconf/202125103045



- Two parallel approaches to reimplement the ME calculation
 - -(1) "CUDACPP", our initial single-code CUDA/C++ back-end targeting NVidia GPUs and SIMD on vector CPUs
 - -(2) Portability Frameworks (PFs: Alpaka, Kokkos, SYCL); added to gain experience with portability





OLD MADEVENT (CURRENT: LHC PROD) SINGLE-EVENT API

MG5aMC: old and new architecture designs



MATRIX ELEMENT: CPU BOTTLENECK IN OLD MADEVENT First we developed the new ME engines in standalone applications

> 1. STANDALONE (TOY APPLICATIONS) MULTI-EVENT API



Then we modified the existing all-Fortran MadEvent into a <u>multi-event</u> framework and we injected the new MEs into it

> 2. NEW MADEVENT (<u>GOAL: LHC PROD</u>) MULTI-EVENT API



Argonne 🕰

Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release

S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023

Université catholique

MadEvent + CUDA Speed up

- Data-centre GPUs such as the Tesla A100 speed up the matrix-element calculations significantly
 - Simpler GPUs lack double-precision performance
 - So far, Madgraph only works reliably in double-precision mode

Process		Madevent 262144 eve	Standalone CUDA	
	Total	Momenta+unweight	Matrix elm	ME Throughput
$e^+e^- \rightarrow \mu^+\mu^-$ +CUDA Tesla A100	17.9 s 10.0 s 1.8 x	10.2 s 10.0 s 1.0 x	7.8 s 0.02s 390 x	$1.9 imes 10^{6} { m s}^{-1}$ $633.8 imes 10^{6} { m s}^{-1}$ $334 { m x}$
$gg \rightarrow t\bar{t}gg$ +CUDA Tesla A100	209.3 s 8.4 s 24.9 x	7.8 s 7.8 s 1.0 x	201.5 s 0.6 s 336 x	$2.8 imes 10^3 m s^{-1}$ $758.9 imes 10^3 m s^{-1}$ 271 m x
$gg \rightarrow t\bar{t}ggg$ +CUDA Tesla A100	$\begin{array}{ccc} 2507.6 & {\rm s} \\ 30.6 & {\rm s} \\ 82.0 & {\rm x} \end{array}$	$\begin{array}{c} 12.2 \ {\rm s} \\ 14.1 \ {\rm s} \\ 0.9 \ {\rm x} \end{array}$	2495.3 s 16.5 s 151 x	$\begin{array}{c} 1.1 \times 10^2 \mathrm{s}^{-1} \\ 170.7 \times 10^2 \mathrm{s}^{-1} \\ 155 \ \mathrm{x} \end{array}$



FORTRAN parts limit total achievable speed up (Amdahl's law)

- Especially for simple SM processes
- E.g. phase-space sampling, generation of momenta, unweighting



What about non-NVidia architectures?



- SYCL:
 - Similar performance as for direct CUDA
 - Portability to AMD / Intel GPUs
- Good scaling observed on Aurora supercomputer testbed "Sunspot"

The XE-HPC results are from early silicon devices and early oneAPI DPC++ software stack. The GPUs are currently used at the Argonne Leadership Computing Facility and other customer sites to prepare code for future Intel data center GPUs, including those to be used in the Aurora exascale supercomputer.

Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release

S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023



What about non-NVidia architectures?



- Initially, ported madgraph C++ matrix elements to CUDA
- Also tested SYCL, as well as Kokkos / Alpaka (discontinued)
- SYCL:
 - Similar performance as for direct CUDA
 - Portability to AMD / Intel GPUs
- Good scaling observed on Aurora supercomputer testbed "Sunspot"

The XE-HPC results are from early silicon devices and early oneAPI DPC++ software stack. The GPUs are currently used at the Argonne Leadership Computing Facility and other customer sites to prepare code for future Intel data center GPUs, including those to be used in the Aurora exascale supercomputer.

-11

cond [s

Madgraph5 aMC@NLO on GPUs and vector CPUs: experience with the first alpha release

S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023



Argonne 🧲

Université

11

MadEvent: FORTRAN $\leftarrow \rightarrow$ FORTRAN+CUDA

 Old MadEvent spent almost all time in matrix-element calculations







Madevent: FORTRAN $\leftarrow \rightarrow$ FORTRAN+CUDA

- Old MadEvent spent almost all time in matrix-element calculations
- With CUDA, the matrix elements are about 1% of the total run time
- Previously unimportant steps such as unweighting limit the speed up





Improving the Madevent Side

- Improved handling of MLM
- GPU-assisted unweighting
 - -Use GPU to for parallel weight computation
 - Helps FORTRAN unweighting routine to discard events faster
- More investigation on madevent side possible

Before, FORTRAN only

Flame Graph







MadEvent with vectorized C++ for $gg \rightarrow t\bar{t}gg$ (on a single CPU core)

		madevent			
$aa \rightarrow t\bar{t}aa$	MEs	$t_{\rm TOT} = t_{\rm Mad} + t_{\rm MEs}$	$N_{\rm events}/t_{\rm TOT}$	$N_{\rm events}/t_{\rm MEs}$	
$gg \rightarrow iigg$	precision	[sec]	[events/sec]	[MEs/sec]	
Fortran(scalar)	double	37.3 = 1.7 + 35.6	2.20E3 (=1.0)	2.30E3 (=1.0)	
C++/none(scalar)	double	37.8 = 1.7 + 36.0	2.17E3 (x1.0)	2.28E3 (x1.0)	
C++/sse4(128-bit)	double	19.4 = 1.7 + 17.8	4.22E3 (x1.9)	4.62E3 (x2.0)	
C++/avx2(256-bit)	double	9.5 = 1.7 + 7.8	8.63E3 (x3.9)	1.05E- (x4.6)	
C++/512y(256-bit)	double	8.9 = 1.8 + 7.1	9.29E3 (x4.2)	1.16E- (x5.0)	
C++/512z(512-bit)	double	6.1 = 1.8 + 4.3	1.35E4 (x6.1)	1.91E- (x8.3)	
C++/none(scalar)	float	36.6 = 1.8 + 34.9	2.24E3 (x1.0)	2.35E. (x1.0)	
C++/sse4(128-bit)	float	10.6 = 1.7 + 8.9	7.76E3 (x3.6)	9.28E3 (x4.1)	
C++/avx2(256-bit)	float	5.7 = 1.8 + 3.9	1.44E4 (x6.6)	2.09E- (x9.1)	
C++/512y(256-bit)	float	5.3 = 1.8 + 3.6	1.54E4 (x7.0)	2.30E4 (x10.0)	
C++/512z(512-bit)	float	3.9 = 1.8 + 2.1	2.10E4 (x9.6)	3.92E4 (x17.1)	



ME speedup ~ x8 (double) and x16 (float) over scalar Fortran <u>Our ME engine reaches the maximum theoretical SIMD speedup!</u> Overall speedup so far~ x6 (double) and x10 (float) over scalar Fortran (Amdahl's law)

Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release

S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023

512



Work-In-Progress, future plans, ideas...

- More performance for matrix elements in CUDA?
 - -Smaller kernels: compiler optimizations, data placement, thread scheduling ...
 - -Split Feynman diagrams and color algebra
 - -Tensore cores for color algebra computations (e.g. cublas)
- More performance for MadEvent
 - -Parallel execution, heterogeneous workflow, vectorisation
 - -Algorithmic optimizations (profiling)
- Functional improvements and longer-term plans
 - -Support for NLO QCD processes
 - –Event-by-event ME reweighting (and derivatives?) → See talk by Z. Wettersten
 - -Test LHAPDF on GPU
 - -Single-precision floating point?



The "alpha release"

- First "gridpacks" generated for testing in LHC experiments
- Upstream madgraph4gpu GPU & SIMD code generation to mg5amcnlo

 Including extra cross-checks that the LHE color IDs are those required for parton showers
- Test and fix any bugs in full pp collision simulation – Implement weak, SUSY, BSM processes
- Stay tuned! [©]
 - -We will be happy to help your experiment in the integration!
 - -NB: This will use a new Fortran version too (multi-event API). Needs statistical validation...



Conclusions

- The Matrix Element calculation in ANY ME event generator can be efficiently parallelized using SIMD or GPUs
- For Madgraph, we get:
 - -SIMD: Speed ups of 6x for AVX512
 - -GPU: 20x to 80x depending on the process using Tesla A100
- Our reengineering of MG5aMC is a functional alpha release for LO QCD / EM processes, but weak interactions need more work
- Please get in touch if you want to test the gridpacks
- Using SYCL, we get similar performances to CUDA and we may also run on AMD or Intel GPUs
- Tests on supercomputer testbeds show good scaling



Acknowledgements

- This research used resources of the Argonne Leadership Computing Facility, which is a U.S. Department of Energy Office of Science User Facility operated under contract DE-AC02-06CH11357, and the Joint Laboratory for System Evaluation (JLSE) at Argonne National Laboratory.
- We gratefully acknowledge the use (under PRACE proposal PRACE-DEV-2022D01-022) of the JUWELS supercomputer and other computing resources provided and operated by the Jülich Supercomputing Centre at Forschungszentrum Jülich.
- We gratefully acknowledge the use (under ISCRA-C project MG5A100) of computing resources provided and operated by CINECA.
- We thank the organizers and our mentors at the GPU Hackathon in CSCS Lugano in September 2022.



BACKUP

Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023



Code Generator and Transformation to CUDA

- User chooses process, MG5aMC determines Feynman diagrams and generates code
 - Initially Fortran (default), C++, or Python
 - More complicated processes produce larger code



Process	LOC	functions	function calls
$e^+e^- \rightarrow \mu^+\mu^-$	776	8	16
$gg \rightarrow t\bar{t}$	839	10	22
$gg \rightarrow t\bar{t}g$	1082	36	106
$gg \rightarrow t\bar{t}gg$	1985	222	786

- Strategy: modify code-generating code (add CUDA, improve C++ backend)
 - (1) Start simple: bootstrap with $e^+e^- \rightarrow \mu^+\mu^-$ (two diagrams, few lines of C++ code)
 - (2,3) Add CUDA and improve C++, port upstream to code generator (Python)
 - (4) Generate more complex LHC processes $gg \rightarrow t\bar{t}, t\bar{t}g, t\bar{t}gg$
 - Add missing functionality, fix issues, improve performance, iterate





Helicity amplitudes – same code in CUDA and in vectorized C++



- Old slide! The new code is different, the idea is the same!
- Formally the same code for CUDA and scalar/vector C++ – hide type behind a typedef – add a few missing operators

SIMD in CUDA/C++ uses compiler vector extensions!

Flexible design: being reused also for vectorized SYCL!

typedef sycl::vec<fptype, MGONGPU_MARRAY_DIM> fptype_sv;



Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release

S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023

Reweighting





$$\bullet w_{i}(m_{W},\Gamma_{W}) = \frac{|\mathscr{M}(m_{W},\Gamma_{W},p_{i}^{1},p_{i}^{2},p_{i}^{3},p_{i}^{4})|^{2}}{|\mathscr{M}(m_{W}^{MC},\Gamma_{W}^{MC},p_{i}^{1},p_{i}^{2},p_{i}^{3},p_{i}^{4})|^{2}}$$

Old technique, renewed interest!

- Advantages of reweighting: savings in computing costs (no detector simulation), fewer statistical fluctuations
- In practice for MG5AMC: read in an LHE file, add weights, write back the modified LHE file – Will use the new matrix element engine in CUDA/C++

- For further details and a status report: Zenny's upcoming poster at CHEP 2023!

- Theoretical and technical challenges
 - NLO reweighting (see O. Mattelaer, https://arxiv.org/abs/1607.00763)
 - Coverage of phase space in the new parameter set
 - Reweighting for a given event-by-event helicity and color



All MadEvent functionalities have been integrated over time

Most of these required some changes to the input/output API of our Fortran-to-CUDA/C++ "Bridge"

- *Helicity filtering* at initialization time, compute the allowed combinations of particle helicities -This is computed in CUDA/C++ using the same criteria as in Fortran
- *"Multi-channel"* single-diagram enhancement of ME output
 - This is the specificity of the MadEvent sampling algorithm (Maltoni Stelzer 2003) $f_i = \frac{|A_i|^2}{\sum_i |A_i|^2} |A_{tot}|^2$
- Event-by-event running QCD coupling constants $\alpha_s(Q^2)$
 - -The scale is currently computed in Fortran from momenta and passed to the CUDA/C++ for each event
 - -No support yet for weak interaction / BSM / SUSY
- Event-by-event *choice of helicity and color* in LHE files
 - -Pass two additional random numbers per event from Fortran to CUDA/C++, retrieve helicity and color
 - -NEW (January 2023)! This was the last big missing physics functionality (showstopper to a release)
 - We now get the same cross section AND the same LHE files (within numerical precision) in Fortran and CUDA/C++





CUDA vs SYCL on NVidia A100

PRELIMINARY! N. Nichols, T. Childers (SYCL) J. Teig (tests/plots)

- SYCL and CUDA implementations have ~similar performances but –SYCL seems better for less complex processes
 - -CUDA seems better for more complex processes
- These are very recent results, which are still being digested (WIP!)

 It will be very interesting to understand more in detail what goes on

We plan to also compare more systematically the CUDACPP and SYCL performances on CPUs (vectorization, multi-core), but it will take time and optimization tweaks... WIP!



ME throughput in C++ for $gg \rightarrow t\bar{t}gg$ (on all the cores of a CPU)

ggttgg check.exe scalability on "bmk6130" (2x 16-core 2.1GHz Xeon Gold 6130 with 2x HT) for 10 cycles



• Most previous results were single-thread peak throughputs

Large SIMD speedups are also confirmed when all CPU cores are used

- AVX512/zmm speedup of x16 over no-SIMD for a single core slightly decreases to ~x12 on a full node (clock slowdown?)
- Overall speedup on 32 physical cores (over no-SIMD on 1 core) is around 280 (maximum would be 16x32=512)
- Aggregate MEs throughput from many identical processes using the standalone application (HEP-workload Docker container)



MEs in MadEvent: CUDA vs SYCL for $gg \rightarrow t\bar{t}gg$



• ME throughput only - SYCL comparable to CUDA but somewhat lower

Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023



Some ideas for heterogeneous processing



Throughput variation as a function of GPU grid size (#blocks * #threads)

This is the number of events processed in parallel in one cycle

To further reduce the relative overhead of the scalar Fortran MadEvent - parallelize it on many CPU cores?

- Blue curve: one single CPU process using the GPU
 - For $gg \rightarrow t\bar{t}gg$, you need at least ~16k events to reach the throughput plateau
- Yellow, Green, Red curves: 2, 4, 8 CPU processes using the GPU at the same time
 - Fewer events in each GPU grid are needed to reach the plateau if several CPU processes use the GPU
 - The total Fortran RAM would remain the same, but the CPU time in the Fortran overhead would be reduced
 - (Why total throughput increases beyond the nCPU=1 plateau is not understood yet!...)

Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023



MadEvent/C++ for $gg \rightarrow t\bar{t}ggg$ (on a single core)

		madevent			standalone
$aa \rightarrow t\bar{t}aa$	MEs	$t_{\rm TOT} = t_{\rm Mad} + t_{\rm MEs}$	$N_{\rm events}/t_{\rm TOT}$	$N_{\rm events}/t_{\rm MEs}$	
$gg \rightarrow iigg$	precision	[sec]	[events/sec]	[MEs/sec]	
Fortran(scalar)	double	813.2 = 3.7 + 809.6	1.01E2 (=1.0)	1.01E2 (=1.0)	
C++/none(scalar)	double	986.0 = 4.3 + 981.7	8.31E1 (x0.8)	8.35E1 (x0.8)	9.82E1
C++/sse4(128-bit)	double	514.7 = 4.2 + 510.5	1.59E2 (x1.6)	1.61E2 (x1.6)	1.95E2
C++/avx2(256-bit)	double	231.6 = 4.0 + 227.6	3.54E2 (x3.5)	3.60E2 (x3.6)	4.41E2
C++/512y(256-bit)	double	208.6 = 3.9 + 204.8	3.93E2 (x3.9)	4.00E2 (x4.0)	4.95E2
C++/512z(512-bit)	double	124.6 = 4.0 + 120.6	6.58E2 (x6.5)	6.79E2 (x6.7)	8.65E2
C++/none(scalar)	float	936.1 = 4.3 + 931.8	8.75E1 (x0.9)	8.79E1 (x0.9)	1.02E2
C++/sse4(128-bit)	float	228.9 = 3.9 + 225.0	3.58E2 (x3.6)	3.64E2 (x3.6)	4.30E2
C++/avx2(256-bit)	float	114.1 = 3.8 + 110.4	7.18E2 (x7.2)	7.43E2 (x7.4)	9.06E2
C++/512y(256-bit)	float	104.5 = 3.8 + 100.7	7.84E2 (x7.9)	8.14E2 (x8.1)	1.00E3
C++/512z(512-bit)	float	61.8 = 3.8 + 58.0	1.33E3 (x13.3)	1.41E3 (x14.1)	1.77E3
C++/none(scalar)	mixed	986.0 = 4.3 + 981.6	8.31E1 (x0.8)	8.35E1 (x0.8)	9.98E1
C++/sse4(128-bit)	mixed	500.4 = 3.9 + 496.5	1.64E2 (x1.6)	1.65E2 (x1.6)	2.00E2
C++/avx2(256-bit)	mixed	220.5 = 3.8 + 216.7	3.72E2 (x3.7)	3.78E2 (x3.8)	4.55E2
C++/512y(256-bit)	mixed	195.6 = 3.7 + 191.8	4.19E2 (x4.2)	4.27E2 (x4.3)	5.21E2
C++/512z(512-bit)	mixed	118.5 = 3.8 + 114.7	6.92E2 (x6.9)	7.15E2 (x7.2)	8.97E2

- Lower overhead of scalar MadEvent in $gg \rightarrow t\bar{t}ggg$ than in $gg \rightarrow t\bar{t}gg$: higher overall throughput speedup x13!
- Mixed floating-point precision (single precision color algebra) is 5-10% better than double



Inside the ME calculation: Feynman diagrams, colors, helicities

$$|\mathcal{M}|^{2}(\vec{p}) = \sum_{\lambda \in \{\text{hel}\}} \left[\sum_{c \in \{\text{col}\}} \left| \sum_{d \in \{\text{diag}\}} (\mathcal{M}^{d}_{\lambda}(\vec{p}))^{(c)} \right|^{2} \right]$$

Given the momenta \vec{p} of initial+final partons in one specific event Sum over all helicity combinations λ of initial+final partons Sum over all color combinations c of initial+final partons Include all Feynman diagrams d allowed for the given λ and c

In practice in MG5aMC: use helicity amplitudes and QCD color decomposition

1. (for each helicity λ) compute partial amplitudes J^f for each color ordering permutation f (sum diagrams relevant to f)

$$(J_{\lambda}(\vec{p}))^f = \sum_{d \in \{\text{diag}\}} (\mathcal{M}^d_{\lambda}(\vec{p}))^f$$

Example for $gg \rightarrow t\bar{t}ggg$: 1240 Feynman diagrams (using helicity amplitudes) This takes ~40% of the CPU time for this process

2. (for each helicity λ) compute the sum over colors as the quadratic form JCJ* using the constant color matrix C

$$|\mathcal{M}|^{2}(\vec{p}) = \sum_{\lambda \in \{\text{hel}\}} \left[\sum_{f,g} (J_{\lambda}(\vec{p}))^{f} (C)^{fg} (J_{\lambda}^{*}(\vec{p}))^{g} \right]$$

Example for $gg \rightarrow t\bar{t}ggg$: 120 color ordering permutations, 120x120 matrix This takes ~60% of the CPU time for this process

3. sum over helicities [Example for $gg \rightarrow t\bar{t}ggg$: 128 helicities (before and after filtering)]

Each step computes many events \vec{p} in parallel! CPU: 1 SIMD event-vector at a time. GPU: 1 event per thread.



S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023



MORE BACKUP SLIDES

Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023



Matrix element integration in MadEvent: detailed results (CPU)



S. Hageboeck - CHEP, Norfolk, VA, 08 May 2023

Université catholique

Matrix element integration in MadEvent: detailed results (GPU)



Argonne 🕰

33

Université catholique

Matrix element integration in MadEvent

Replace Fortran MEs by cudacpp (or PFs) MEs in Madevent (keep the same user interface!)

Linking Fortran and C++ has been easy. As expected, the two main issues have been, instead:

- -1. Moving Madevent from single-event to many-event (functional reengineering of the algorithm)
 - Now also an active area of performance optimizations (next slides: GPU grid and CPU RAM; CPU time and Amdahl...)
- -2. Debugging functional issues caused by hidden inputs and outputs, e.g. coming from Fortran common blocks



Argonne 🛆 🕅 UCL Université de Louvain

Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release

Code generation: from many "epochs" to a single evolving "epoch"



Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release

S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023



Portability Frameworks (PFs)

(2) Second line of development: MEs on PFs

- PFs allow writing algorithms once and running on many architectures with some hardware-specific optimizations
- CUDA code can only run on NVidia GPUs, while <u>Kokkos</u>, <u>Alpaka</u>, and <u>Sycl[Intel]</u> codes can run on most hardware
- In "cudacpp", #ifdef directives separate code branches for GPU and CPU code during compilation (but these are very few: only kernel launching and memory access, not MEs)
- With PFs, the algorithm is typically the same, but the compilation occurs once per architecture type
- PFs often use templating to handle data types and hardware configuration and function lambdas or pointers for passing kernels (the cudacpp plugin has many of these, too)
- PFs still require user to think about "host" vs "device"



"cudacpp" example of compiler directives

540	<pre>#ifdefCUDACC</pre>
541	<pre>#ifndef MGONGPU_NSIGHT_DEBUG</pre>
542	<pre>gProc::sigmaKin<<<gpublocks, gputhreads="">>>(devMomenta.get(), devMEs.get()</gpublocks,></pre>
543	#else
544	gProc::sigmaKin<< <gpublocks, gputhreads,="" ntpbmax*sizeof(float)="">>>(devMome</gpublocks,>
545	#endif
546	checkCuda(cudaPeekAtLastError());
547	<pre>checkCuda(cudaDeviceSynchronize());</pre>
548	#else
549	<pre>Proc::sigmaKin(hstMomenta.get(), hstMEs.get(), nevt); For CPU</pre>
550	#endif

Kokkos example of Templating & lambda

324	{
325	<pre>using member_type = typename Kokkos::TeamPolicy<kokkos::defaultexecut< pre=""></kokkos::defaultexecut<></pre>
326	Kokkos::TeamPolicy <kokkos::defaultexecutionspace> policy(league_size</kokkos::defaultexecutionspace>
327	Kokkos::parallel_for(func,policy,
328	KOKKOS_LAMBDA(member_type team_member){
220	

Kokkos example of Memory Management

262 Kokkos::View<fptype***,Kokkos::DefaultExecutionSpace> devMomenta(Kokkos::ViewAllocateWithoutInitializing("devMomenta"),nevt,npar,np4); 263 auto hstMomenta = Kokkos::create_mirror_view(devMomenta);



S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023



EVEN MORE BACKUP SLIDES

Madgraph5_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023



Standalone CUDA/C++ application VS. MadEvent integration

- Our main focus: the ME calculation in CUDA/C++ (sigmakin kernel/function)

 Design approach: single source code for CUDA and C++ (>90% common code + #ifdef's)
- Our workhorse: a simplified CUDA/C++ toy framework to feed events to the ME kernel
 - All 3 main components on the GPU: random (cuRAND), sampling (RAMBO), ME (sigmakin)
 - Fast, same results in GPU/CPU, but not good for production (RAMBO algorithm is inefficient)
 - The results I present in this talk come from this framework





Event-level parallelism in practice – coding and #events

- Easier to code for GPU SIMT than for CPU SIMD: CUDA code was faster to prototype
- CUDA (GPU) implementation
 - For SIMT, event loop is "orthogonal": one thread = one event (GPU thread ID ↔ event ID)
 - For SIMT, SOA memory layouts are beneficial (coalesced access), but not strictly essential
- C++ (CPU) implementation
 - For SIMD, event loop must be the innermost loop (e.g. invert helicity and event loops)
 - For SIMD, SOA memory layouts in the computational kernel are essential
- To be efficient, CUDA needs O(10k)-O(1M) events in parallel much more than C++!
 - CUDA: lockstep within each warp (32 threads) + many warps in parallel to fill the GPU
 C++: lockstep within a vector register (2-8 doubles) + multi-threading or multi-processing





S. Hageboeck – CHEP, Norfolk, VA, 08 May 2023