

Study of the feasibility of accelerating MC generation with FPGAs

Jornadas ICTEA 2025

Santiago Folgueras¹, <u>Carlos Vico Villalba</u>¹, Pelayo Leguina López¹ Luca Fiorini², Héctor Gutiérrez Arance², Alberto Valero Biot², Francisco Hervás Álvarez², Arantza De Oyanguren Campos², Jiahui Zhuo², Volodymyr Svintozelskyi², Valerii Kholoimov²,

> ¹ Universidad de Oviedo, Asturias, Spain ² Instituto de Física Corpuscular (CSIC-UV), Valencia, Spain

The Standard Model of Particle Physics

- The SM is a renormalizable Quantum Field Theory.
 - It predicts the nature of high energy interactions between the particles of the model.
- The **goal** of the experiments at the **LHC** is to
 - **Measure** the **properties** of **these particles**.
 - Explore further extensions in search of new physics.



Montecarlo modeling

- The **SM predictive power** is **tested against** the recorded **data**.
 - Montecarlo (MC) generation approaches allow to resembles the nature of the data (counting experiments) by generating events.
- Based on a **factorized approach** that **divides** the whole **calculation into steps**.
 - Hard interactions: Leading Order (LO) + Next-to-LO (NLO) α_s expansion terms in a perturbative series at scale μ .
 - Soft interactions: approximation of collinear/soft QCD/weak radiation effects.



Montecarlo modeling

- **MC modeling** (how well the simulations describe the data) is **essential** for the **LHC physics programme**.
 - For instance: ttW (measured in Oviedo) production has shown slight tensions with the theory.
 - (<u>Also in Oviedo</u>) the modeling is thoroughly reviewed.





- Left: momentum in the transverse plane for the radiation emitted when producing a pair of quarks and a W boson.
- **Problem in 2021**: discontinuity at 150 GeV.
 - → Caused a 4 standard deviations disagreement between experiment and theory!

Montecarlo modeling

- MC modeling (how well the simulations describe the data) is essential for the LHC physics programme.
 - For instance: ttW (measured in Oviedo) production has shown slight tensions with the theory.
 - (<u>Also in Oviedo</u>) the modeling is thoroughly reviewed.





- Left: momentum in the transverse plane for the radiation emitted when producing a pair of quarks and a W boson.
- Solution in 2022 (from MC builders):
 - Proper modeling of weak radiation off of W boson.
 - A slight tension remains at ~2 standard deviations.
 - Much more controlled.

Current issues

- **MC modeling** is a challenge from multiple perspectives:
 - In the case of ttW \rightarrow > 2 years of work just to get this MC simulation working in CMS!!
- Timetable for this work:
 - Few months of **compiling all the diagrams** that make up the **hard interaction**.
 - Each time we do this: ~1 day of running (Intel(R) Core(TM) i7 CPU 970 @ 3.20GHz)
 - Validation has to be made in a large enough sample of events (statistically meaningful predictions):
 - We can generate about 1M events in ~1 day
 - A reasonably large sample is made out of 10M events
 - This, in any case, runs relatively fast and efficient.
- <u>The bottleneck is in the compilation time</u>

Current issues

- **This issue** is actually more relevant for much more common processes \rightarrow example: top-antitop pair • production (measured in Oviedo).
- The problem grows as we add more radiation diagrams •



Current issues

- "Does the trick" → assume the radiation is soft/collinear, take the modeling at degraded accuracy from the parton shower.
 - Works very well for smaller datasets \rightarrow not a problem for the Run 2 + Run 3 results.
 - Let's not forget we are aiming at measuring properties of the Higgs in phase spaces with up to 4 jets (HH→4b channel) in the final state → this is extremely sensitive to these effects.
 - One needs to solve the bottleneck by better improving / pipelining the computation through improved running architectures such as GPU



Where the community is going

- The developers of one of the most commonly used MC software package (Madgraph5_aMC@NLO) are already working towards a solution to this problem.
 - As of now Madgraph5_aMC@NLO works in Fortran. Runs on CPU.
 - They ported the code to consider CUDA/C++ which is significantly faster.

gū→τ⁺τ⁻ggū (subprocess of DY+3j) 100 Feynman diagrams 6x6 color matrix



A

		madevent											
$g\bar{u} \rightarrow \tau^+ \tau^- gg\bar{u}$ 81920 weighted events)	MEs precision	$t_{\text{TOT}} = t_{\text{Mad}} + t_{\text{MEs}}$ [sec]	$N_{\text{events}}/t_{\text{TOT}}$ [events/sec]	N _{events} /t _{MEs} [MEs/sec]									
Portran(scalar)	double	52.2 = 17.0 + 35.2	1.57E3 (=1.0)	2.32E3 (=1.0									
C++/none(scalar)	double	50.9 = 17.0 + 33.9	1.61E3 (x1.0)	2.42E3 (x1.0									
C++/sse4(128-bit)	double	35.5 = 17.0 + 18.5	2.31E3 (x1.5)	4.44E3 (x1.9									
C++/avx2(256-bit)	double	24.5 = 16.9 + 7.6	3.34E3 (x2.1)	1.08E4 (x4.7									
C++/512y(256-bit)	double	23.9 = 16.9 + 7.0	3.43E3 (x2.2)	1.17E4 (x5.0									
C++/512z(512-bit)	double	26.7 = 17.0 + 9.6	3.09E3 (x2.0)	8.57E3 (x3.7									
CUDA/GPU	double	17.6 = 17.4 + 0.3	4.65E3 (x3.0)	3.26E5 (x140									
C++/none(scalar)	mixed	50.9 = 16.9 + 33.9	1.61E3 (x1.0)	2.41E3 (x1.0									
C++/sse4(128-bit)	mixed	33.9 = 16.9 + 17.0	2.41E3 (x1.5)	4.82E3 (x2.1									
C++/avx2(256-bit)	mixed	24.8 = 17.2 + 7.6	3.31E3 (x2.1)	1.08E4 (x4.7									
C++/512y(256-bit)	mixed	24.1 = 17.1 + 7.0	3.40E3 (x2.2)	1.18E4 (x5.0									
C++/512z(512-bit)	mixed	26.5 = 17.0 + 9.6	3.09E3 (x2.0)	8.57E3 (x3.7									
CUDA/GPU	mixed	17.7 = 17.4 + 0.3	4.64E3 (x3.0)	3.23E5 (x138									
C++/none(scalar)	float	50.1 = 16.9 + 33.1	1.64E3 (x1.0)	2.47E3 (x1.1									
C++/sse4(128-bit)	float	26.3 = 16.9 + 9.4	3.11E3 (x2.0)	8.70E3 (x3.7									
C++/avx2(256-bit)	float	20.8 = 16.9 + 3.9	3.94E3 (x2.5)	2.12E4 (x9.1									
C++/512y(256-bit)	float	20.6 = 16.9 + 3.6	3.99E3 (x2.5)	2.27E4 (x9.7									
C++/512z(512-bit)	float	21.7 = 16.9 + 4.8	3.78E3 (x2.4)	1.71E4 (x7.3									
CUDA/GPU	float	17.6 = 17.4 + 0.2	4.66E3 (x3.0)	4.46E5 (x191									

By A. Valassi (CHEP 2024)

For one typical subprocess of DY+3jets: Fortran MEs ~ 67% of the total time => Max overall speedup is x3 (Amdahl)

Achieved speedups (mixed FP precision): - x3.0 on GPU (NVidia V100) - x2.2 on SIMD ("512y" on Intel Silver)

In Oviedo we went for an alternative way :)

- One liner: to perform viability studies on whether MC generation workflows can be accelerated using Field Programmable Gate Arrays (FPGAs).
- The project is a proof of concept
 - Can we use FPGAs for MC generation?
 - At this stage, we have identified a task to be optimized in FPGA, and started making comparisons against CPU performance.
 - The ultimate goal is to test against GPU performance to see if FPGA can actually be competitive on this task.



1

- **Technical description:** A Field Programmable Gate Array (FPGA) is an integrated circuit with a programmable hardware fabric that allows it to be reconfigured to behave like another circuit.
- Non-technical description: an electronic chip built inside a board that can be reprogrammed at any time (unlike GPUs/CPUs or any kind of ASIC).
- The advantages of an FPGA are:
 - **Fine-grained customization**
 - **Spatial compute** 0
 - Hardware flexibility Ο
 - **Diverse IOs** \bigcirc
- The disadvantages are:
 - Little access to non-expert users. 0
 - "Limited" resources for complex operations 0 (e.g divisions, exponentiations).



- In CPU programming:
 - Write code in a given language (C, C++, fortran).
 - \circ Compile the code.
 - Code is translated into instructions that can be understood by the CPU.
 - We execute the compiled code.



- In CPU programming:
 - Write code in a given language (C, C++, fortran).
 - \circ Compile the code.
 - Code is translated into instructions that can be understood by the CPU.
 - We execute the compiled code.
- In FPGA programming:
 - We compile code in a given language (VHDL, HLS).
 - The code is translated into Register Transfer Level language (RTL).
- The RTL tells the FPGA what circuit do we have to build to solve the task.

```
D ≤ not Q;
process(clk) VHDL
begin
    if rising_edge(clk) then
        Q ≤ D;
    end if;
end process;
```



• In CPU programming:

- Write code in a given language (C, C++, fortran).
- \circ Compile the code.
- Code is translated into instructions that can be understood by the CPU.
- \circ We execute the compiled code.

• In FPGA programming:

- We compile code in a given language (VHDL, HLS).
- The code is translated into Register Transfer Level language (RTL).
- The RTL tells the FPGA what circuit do we have to build to solve the task.
 - Which resources from the FPGA are needed to solve the task at hand.

Signals	-Waves	-					_		_	_				_				
Time	3		100 ns	s		20	90 ns					300	ns				40	90 ns
ncolor_mult																		
valid_in=1																		
valid_counter=3	θ		(1	2)(з	θ												
continious_read=0																		
first_valid=1																		
bram_enable=1																		
cf_addr_counter=1	θ				,1	2	3)4)5)6)	7 8	9)+) E)F	DÐ)	+ 1	5			
bram address mult[5:0]=1	Θ					1	2	3 4)5)	6 7	8 9			+ +	15			
bram data in[31:0]=00000000	0000000				+ 00	3+ (+	000+	+ (+)+)•)(e)e	++)+)+)Ð	+)+	+ 0	0100	010	
ncolor_counter=3	0		(1	2)(3	4	5 6	7 8	9	+)+	Ð)+)+	0	1 2	3 4	5	6 7 8	9 +
i idx=0	0						1		2		3			0		1		2
j idx=3	θ				3		0 1	2 3	0	1 2	3 0	1 2	3	9 1	2 3	0	1 2 3	0 1
latency counter[3:0] =8	0		(1	2 3 4	5 6 7	8 9	(F)F)+)1	3									
sumM x																		
r[15:0]=24	0		6	15	24	33	6)+)+		6 +	++	6 +		+ 6)+)+		6 + +	+ 6
i[15:0]=24	θ		6	15	24	33	6) -)+) .)	6 +	+++	6+)÷)	+ 6)+)+	(F)	5 + +	+ 6
sumM y conj																		
r[15:0]=6	θ		6)(1	5		24		33		6		-y	15	24
i[15:0]=-6	0		-6)(-	15	Y	-24		-33		-	6		- 15	-2
mult1																		
i[15:0]=24	θ		6	15	24	33)6	`) +)+)(+)(6 +	++	6+)F)	+)6)+)+)(=)(6 + +	+ 6
r[15:0]=24	θ		6	15	24	33	6)+)+		6 +	-	6 +		+ 6)+)+		5 + +	+ 6
i[15:0]=-6	θ		-6				γ.	15	Ťγ	-24		-33)-(6		- 15	
r[15:0]=6	θ		6)1	5		24		33		6			15	24
p1 zr sig[15:0]=72	(DOOX)0			72		180	¥28	38	39	6 +) T	γ . γ	-)-) T)F)	+ + + +	+++
p1 zi sig[15:0]=0	DXXX	θ																
p1 out																		
p1 pipe																		
mult2																		
result valid mult2=0																		
p2 zr sig[15:0] =0	()xxx)0				[+] e	<u>_</u>	10			Ð			-)-)-	+)+
p2 zi sig[15:0]=0	(DXXX			6				[+	6	[+	0	Ē	N		ŶŦŢŦ	ή-γ	+ + +	+++
partial regl												_^						
cf reg																		
i[15:0]=0	Θ				1	θ	5 0	9		2 6	+)+	3)7	ne)	+ 4	8 +	16		
r[15:0]=0	0				(1	0	5 0	9).	2 6	++	3 7		+ 4	8+	16		
p2																		
p2 reg																		
matrix1 reg																		
i[15:0]=0	θ								72		972)(+		+)+)+)+		+ 3276	7
r[15:0]=0	0								72		972	7		+ +)+)+		+ 3276	7
												_						

• We focus all of our examples in one of the main processes used in experimental High Energy Physics: ttbar multijets (considering up to 3 jets).

These 10 lines of code actually consume ~60% of the total computation time!!!



P1_gg_ttxggg

- 119 integration channels after gensym.
- Focused on one (G1).
- We measured the time that is spent in integrating this channel.



- t_{it} is the time per iteration of the color matrix operation.
- C_{t} is the cumulated time spent in doing this operation
- $C_t = t_{it} \times N_{calls}$, where N_{calls} is the number of calls to functions that integrate the diagram)

Using an optimization for these kind of computations

• Translation to FPGA language

from Sign Process Syst 95, 543-550 (2023) TS(K) DO M = 1, NAMPSO nampso acc DO I = 1, NCOLOR ZTEMP = (0.D0, 0.D0)ncolor mult DO J = 1, NCOLOR ZTEMP = ZTEMP CF(J,I) JAMP(J,M) OFFFO DO N = 1, NAMPSO ncolor mult SEEEO TS(K) = TS(K) + ZTEMP*DCONJG(JAMP(I,N)) COEFFS CF BRAM N binary tree adder nampso acc JAMP sumM_c_0 conj JAMP D D sumM 0 BRAM 0 Text sumM c 1 conj JAMP D D sumM 1 BRAM 1 Output sumM_c_N coni JAMP D sumM_0_N D Matrix1 BRAM N

Car

• We have computed a first preliminary comparison of the performance between CPU/ FPGA DSPs / FPGA AI Cores.

Time spent in one call to the color matrix calculation for different architectures, expressed in microseconds (μs). The numbers for the fortran column have been obtained using the default output from madgraph. For the CPU we use an Intel(R) Core(TM) i7 CPU 970 @ 3.20GHz. For the FPGA (DSP) we use a VCU13P FPGA running at 645 MHz. For the FPGA (AI CORES) we are using a VERSAL board.

N _{color}	t _{it} CPU(fortran)	t _{it} FPGA(DSP)	t _{it} FPGA(AI Cores)
6	1.00	0.043	2.7
27	2.00	0.14	Work in progress
120	29.00	0.45	Work in progress

Conclusions

- In this talk we have presented a **first temptative implementation** of the **color calculation in a FPGA** device.
 - How does an FPGA work.
 - A priori limitations.
- We have started to compare the performance between CPU and two different implementations of the color matrix calculation in FPGA.
- Our metric so far: How much time does it get an output per call.
 - Good improvement in timings between CPU and FPGA (DSP).
 - Further optimization needs to be done in the FPGA (AI CORE) implementation to get competitive result.
- To be done in following iterations:
 - Test another implementation in **High Level Synthesis** language (a C++ interface to VHDL).
 - Comparisons with GPU performance!
 - Test compatibility of results between CPU/FPGA.
 - Test against cudacpp CPU implementation (possibly faster?)
 - Test the overall time taken to perform the whole integration.

carlos.vico.villalba@cern.ch

backup

- There are two main considerations to be taken into account.
 - Area constraints: do we have enough resources to perform the task?



- There are two main considerations to be taken into account.
 - **Area constraints:** do we have enough resources (LUTs, FFs, BRAMS, DSPs, etc...) to perform the task?
 - The more operations we have to do (in parallel) \rightarrow the more resources we consume.



- There are two main considerations to be taken into account.
 - **Area constraints:** do we have enough resources (LUTs, FFs, BRAMS, DSPs, etc...) to perform the task?
 - The more operations we have to do (in parallel) \rightarrow the more resources we consume.
- Things that significantly affect AREA: Divisions, exponentiations, ... Ο Floating point operations. 0 Ż TUC Image used for illustration purposes of a more complex operation. Occupancy >> 1%

- There are two main considerations to be taken into account.
 - Area constraints: do we have enough resources (LUTs, FFs, BRAMS, DSPs, etc...) to perform the task?
 - The more operations we have to do (in parallel) \rightarrow the more resources we consume.
- Things that significantly affect AREA:
 - Divisions, exponentiations, ...
 - Floating point operations.

Note modern day FPGAs have a ton of resources, so area is in some cases not the most limiting factor.



Image used for illustration purposes of a more complex operation.



Occupancy >> 1%

- There are two main considerations to be taken into account.
 - Area constraints: do we have enough resources (LUTs, FFs, BRAMS, DSPs, etc...) to perform the task?
 - The more operations we have to do (in parallel) \rightarrow the more resources we consume.
 - **Timing constraints:** algorithm meets the required clock frequency while maintaining correct functionality?
 - Pipelining: Pipelining is a technique that improves data processing speed by breaking operations into smaller sequential stages, allowing partial results to be processed concurrently.
 - After a certain number of cycles, the pipeline reaches a steady state where one output is produced per clock cycle

- There are two main considerations to be taken into account.
 - Area constraints: do we have enough resources (LUTs, FFs, BRAMS, DSPs, etc...) to perform the task?
 - The more operations we have to do (in parallel) \rightarrow the more resources we consume.
 - **Timing constraints:** algorithm meets the required clock frequency while maintaining correct functionality?
 - Pipelining: Pipelining is a technique that improves data processing speed by breaking operations into smaller sequential stages, allowing partial results to be processed concurrently.
 - After a certain number of cycles, the pipeline reaches a steady state where one output is produced per clock cycle



Example of a **non-pipelined** designed. This design requires 3 clocks to finish, but only 1 is used, leading to time slack.



Example of a **pipelined** designed. This design uses more registers to keep processing inputs while the final result of the first set of inputs is not finished yet. **Consequence: more area.**

- There are two main considerations to be taken into account.
 - Area constraints: do we have enough resources (LUTs, FFs, BRAMS, DSPs, etc...) to perform the task?
 - The more operations we have to do (in parallel) \rightarrow the more resources we consume.
 - **Timing constraints:** algorithm meets the required clock frequency while maintaining correct functionality?
 - Pipelining: Pipelining is a technique that improves data processing speed by breaking operations into smaller sequential stages, allowing partial results to be processed concurrently.
 - After a certain number of cycles, the pipeline reaches a steady state where one output is produced per clock cycle



Example of a **non-pipelined** designed. This design requires 3 clocks to finish, but only 1 is used, leading to time slack.



Example of a **pipelined** designed. This design uses more registers to keep processing inputs while the final result of the first set of inputs is not finished yet. **Consequence: more area.**

Implementation of the color decomposition in FPGA AI Cores

- The previous implementation is done using **Digital Signal Processors** (DSP), **Blocks of ram** (BRAM), **Lookup Tables** (LUTs), and other kids of hardware components.
 - **Problem:** DSP do not handle well floating point precission.
- Modern day FPGA boards (see <u>Versal FPGA</u>) include what is called Adaptive Intelligence (AI) engines.
 - Specialized processing blocks optimized for high-throughput, low-latency workloads

• Al engines are:

- Arrays of cores for parallel processing.
- Each core can process vectors in parallel at higher speeds (1.25 GHz).
- Each core can handle:
 - 32 KB data memory
 - 16 KB program memory
 - Efficient data flow.



Implementation of the color decomposition in FPGA AI Cores



• We have computed a first preliminary comparison of the performance between CPU/ FPGA DSPs / FPGA AI Cores.

Time spent in one call to the color matrix calculation for different architectures, expressed in microseconds (μs). The numbers for the fortran column have been obtained using the default output from madgraph. For the CPU we use an Intel(R) Core(TM) i7 CPU 970 @ 3.20GHz. For the FPGA (DSP) we use a VCU13P FPGA running at 645 MHz. For the FPGA (AI CORES) we are using a VERSAL board.

N _{color}	t _{it} CPU (fortran)	t _{it} FPGA(DSP)	t _{it} FPGA(AI Cores)
6	1.00	0.043	2.7
27	2.00	0.14	Work in progress
120	29.00	0.45	Work in progress

Implementation of the color decomposition in FPGA (VHDL)

bram addr bram

latency

Time

Second step: validation in simulation

Schematic view of the signals throughout one call to the matrix color code.

One has to take into account that in a FPGA, variables become signals.

This is the first working version for the computation with VHDL implementation.

s	-Waves-																								
			100 ns			200	ns			300	15		4	00 ns			00 ns			60	90 ns			700	0 ns
ncolor_mult																									
valid_in=1																									
valid_counter=3	Θ		1	2	3	0																			
<pre>continious_read =0</pre>																									
first_valid=1																									
bram_enable=1																									
cf_addr_counter=1	Ð				1	2	3 4	5 6 7	89	333	++=15														
_address_mult[5:0]=1	0e					1	2 3 4	156	7 8	9.4.4		15													
bram_data_in[31:0] =00000000	0000000				+ 00	0+ + 0	00+ + +					+ 00	100010												
ncolor_counter=3	θ		1	2	3	4 5	678	39+	(H)		0 1 2	3 4	5678	9++	+ + + +	θ 1 2	3 4	5 6 7	8 9	+ + +	+ + +	0 1	2 3 4	56	7 8
i_idx=0	0					1					0					0						0			
j_idx =3	θ				3	Θ	127	3 0 1	2 3	0 1 2	3 0 1	2 3	0123	012	3 0 1 7	301	2 3	012	30	123	0 1 2	30	123	01	23
tency counter[3:0] =8	θ		1	234	567	89+		13																	
sumM_x																									
r[15:0] =24	θ		6	15	24	33)6 + (+ + 6	++	+6+	++6	+ + /	+6++	+6+	++6+	++6	+ +	+ 6 +)+)+)	6 + +	+ 6 +) - -	6 + +	+ 6	+)Z
i[15:0] =24	le le		6	15	24	33	6+	4 + 6		-	++6	+	+ 6 + +	+6+	+ + 6 +	++6)+ (+	+ 6 +	++	6 + +	+ 6 +)+ (+)	6 + +	+ 6	+ 2
sumM y conj																									
r[15:0] =6	Θ		6				15	24		33	6		15	24	33	6		15		24	33		6	15	
i[15:0] =-6	Ø		-6				-15		4	-33	-6		-15	-24	-33		6	-15		-24	-33		-6	a di	5
mult1																									
i[15:0] =24	0		6	15	24	33	6+			+ 6 +	++6	++	+ 6 + +	+ 6 +	+ + 6 +	++6	++	+ 6 +	++	6 + +	+ 6 +)+)+	6++	+ 6	+ 2
r[15:0] =24	Θ		6	15	24	33	6 + -	+ + 6		+ 6 +	+ + 6	+	+ 6 + +	+ 6 +	+ + 6 +	+ + 6	-	+ 6 +	(I)	6 + +	+ 6 +	++	6 + +	- 6	+ 2
i[15:0] =-6	Θ		-6				-15		24	-33	-6		-15	-24	-33		6	-15		-24	-33		-6	0.1	5
r[15:0] =6	0		6				15	24		33	6		15	24	33	6		15	Y	24	33		6	15	
p1 zr sig[15:0] =72	()xxxx]0			72	1	80 (1	288	396	000	D	(+)+Y	E (E) E (E				X-X-	(F)(F)(F	(E)		FIF	DE	(III)	JODY	Ŧ
p1 zi sig[15:0] =0	()poox	0								T															
pl out																									
pl pipe																									
mult2																									
result valid mult2=0																									
p2 zr sig[15:0] =0	()xxx)e			l l	4 0	+ 0	DO	E)EE	++	+++++		- (-)-)-		X-T-				E E	DE	-		EDE
p2 zi sig[15:0] =0	()xxxx			0				+ 0	+ 0			ŌÐ	+++++		+)+)+			E)E)E	NOD			100		100	Đĩ
partial reg1																									
cf reg																									
i[15:0]=0	Ð				1	0 5	0	9 + 2	6 +	+ 3 7	+ + 4	8 +	16												
r[15:0]=0	Ð				1	0 5	0	1 + 2	6 +	+ 3 7	+++4	8 +	16												
p2																									
p2 reg																									
matrix1 reg																									
i[15:0]=0	0							72	977	-	++++	+ +	+ + 3276	7											
r[15:0]=0	Θ							72	977		-	++	+ + 3276	7											