

A MPI/OpenCL hybrid implementation of the Matrix Element Method in the context of the Higgs boson property analyses

G. Grasseau, T. Strebler, A. Chiron,
P. Paganini and F. Beaudette,

Laboratoire Leprince-Ringuet
CNRS/IN2P3, École Polytechnique

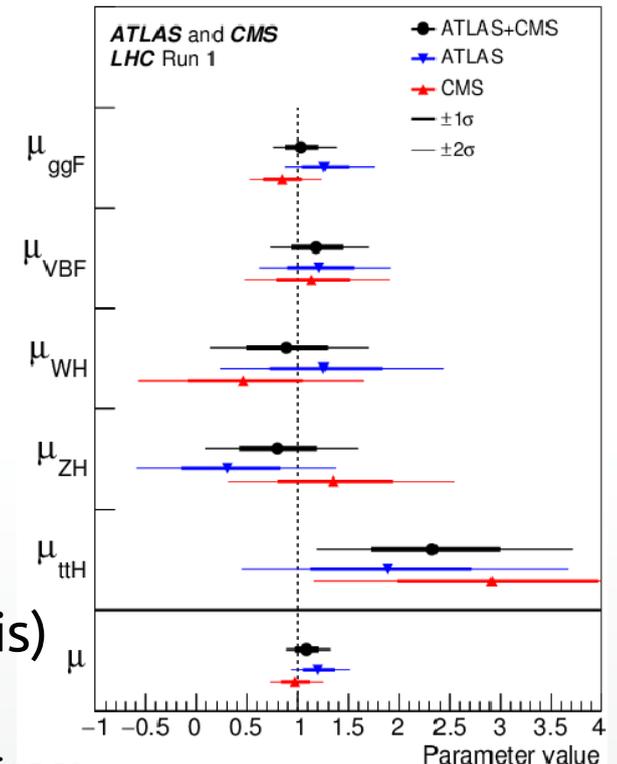
Introduction

Elementary particles of the Standard Model (SM)

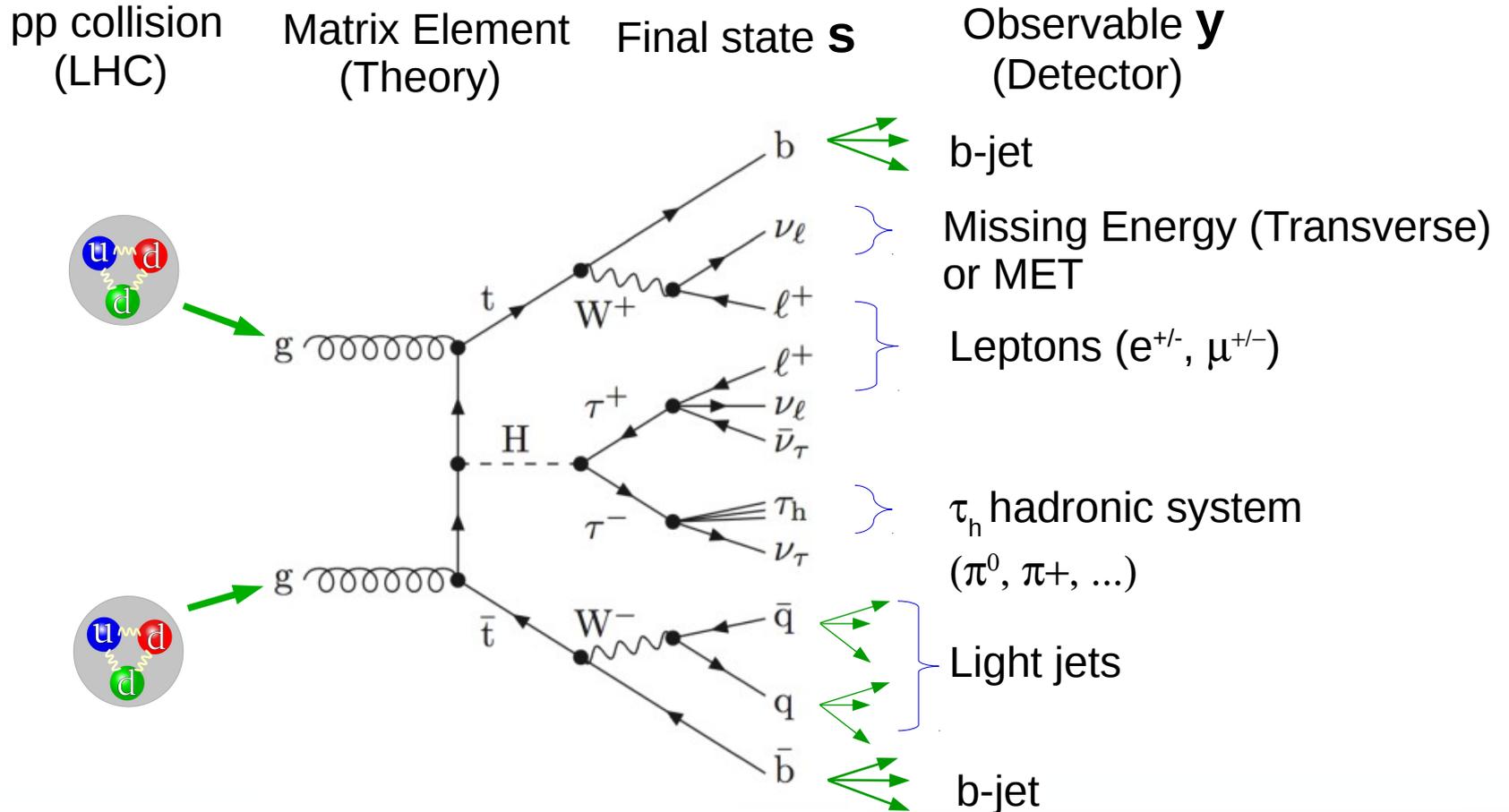
mass → charge → spin →	$\approx 2.3 \text{ MeV}/c^2$ 2/3 1/2 u up	$\approx 1.275 \text{ GeV}/c^2$ 2/3 1/2 c charm	$\approx 173.07 \text{ GeV}/c^2$ 2/3 1/2 t top	0 1 g gluon	$\approx 126 \text{ GeV}/c^2$ 0 0 H Higgs boson
	$\approx 4.8 \text{ MeV}/c^2$ -1/3 1/2 d down	$\approx 95 \text{ MeV}/c^2$ -1/3 1/2 s strange	$\approx 4.18 \text{ GeV}/c^2$ -1/3 1/2 b bottom	0 0 1 γ photon	
	$0.511 \text{ MeV}/c^2$ -1 1/2 e electron	$105.7 \text{ MeV}/c^2$ -1 1/2 μ muon	$1.777 \text{ GeV}/c^2$ -1 1/2 τ tau	0 1 1 Z Z boson	
	$< 2.2 \text{ eV}/c^2$ 0 1/2 ν_e electron neutrino	$< 0.17 \text{ MeV}/c^2$ 0 1/2 ν_μ muon neutrino	$< 15.5 \text{ MeV}/c^2$ 0 1/2 ν_τ tau neutrino	$80.4 \text{ GeV}/c^2$ ± 1 1 W W boson	GAUGE BOSONS

- The recently discovered Higgs boson (2012) can be produced in different ways in pp collisions @LHC
- The combination of LHC experiments shows an excess of events when the H is produced in association with 2 top quarks (ttH channel)

- ttH is an interesting channel to look at Run 2: it allows probing the top-Higgs Yukawa coupling
- In addition, it has several decay channels
- Among all the ttH channels looked by CMS, the H decays in 2 τ (H-> $\tau\tau$) is one of the most challenging
- LLR team is deeply involved in Matrix Element Method : VBF, ttH channel (T. Strebler PHD thesis)



Theory and observables



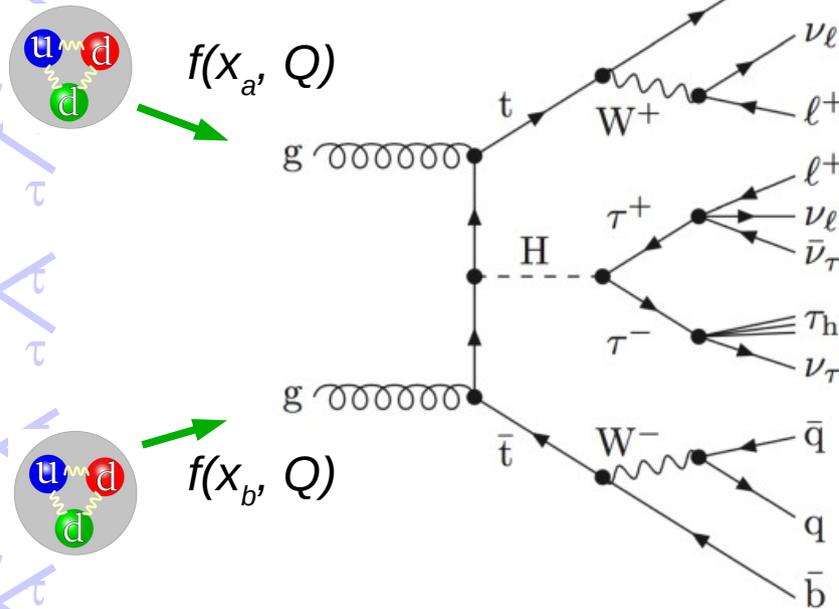
- Leptons $\ell^{+/-}$, hadronic system τ_h , are precisely reconstructed
- Jet energy reconstructed with a finite resolution
- ν 's are unobserved but their global (transverse) momentum can be inferred from the MET

Matrix Element Method (MEM)

pp collision
(LHC)

Matrix Element
(Theory)

Final states s



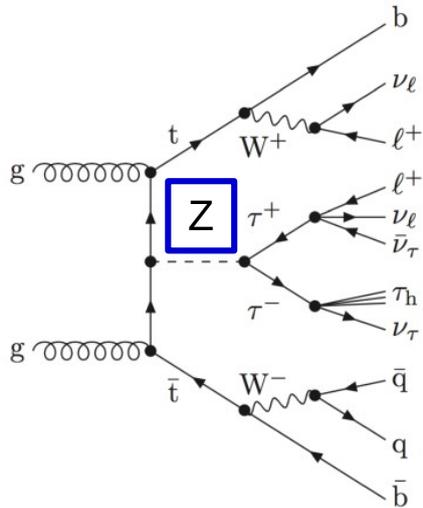
- Studying ttH with the final state s : $2b, 2q, \tau_h, 2 leptons$
- Observable \mathbf{y} lead to several possible states s which are of interest (signal) or not (background)
- MEM: explore (give a weight) to all the possible value of \mathbf{x} which lead to the observable \mathbf{y}

$$w_i(\mathbf{y}) = \frac{1}{\sigma_i} \sum_p \int dx dx_a dx_b \frac{f(x_a, Q) f(x_b, Q)}{x_a x_b s} \delta^2(x_a P_a + x_b P_b - \sum p_k) |\mathcal{M}_i(\mathbf{x})|^2 W(\mathbf{y}|\mathbf{x})$$

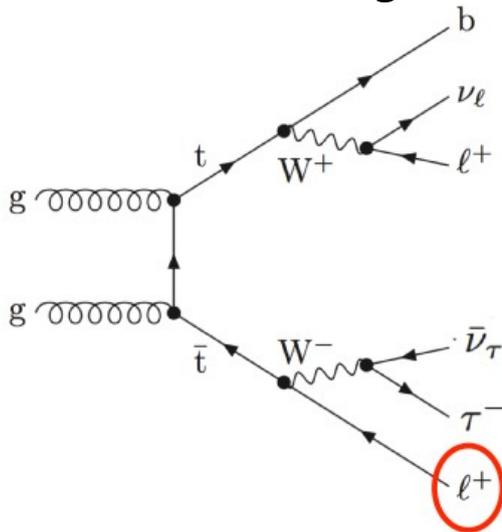
x_a, x_b parton impulsion fraction f Parton Density Fraction (PDF) Momentum conservation ME Transfer Function (TF)

TF: detector response

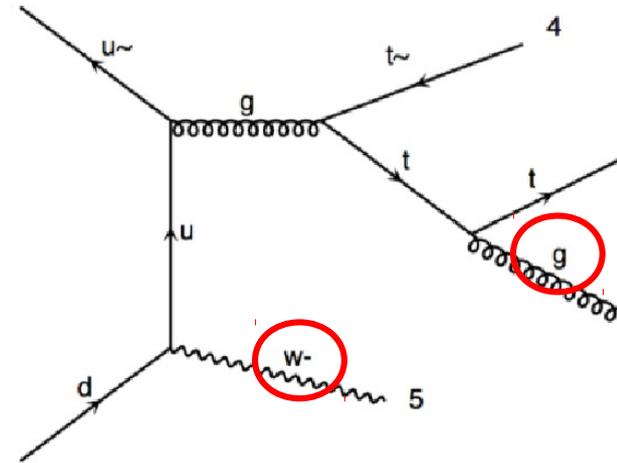
Final states: backgrounds



ttZ: Z production (decays in $\tau\tau_h$)
irreducible background



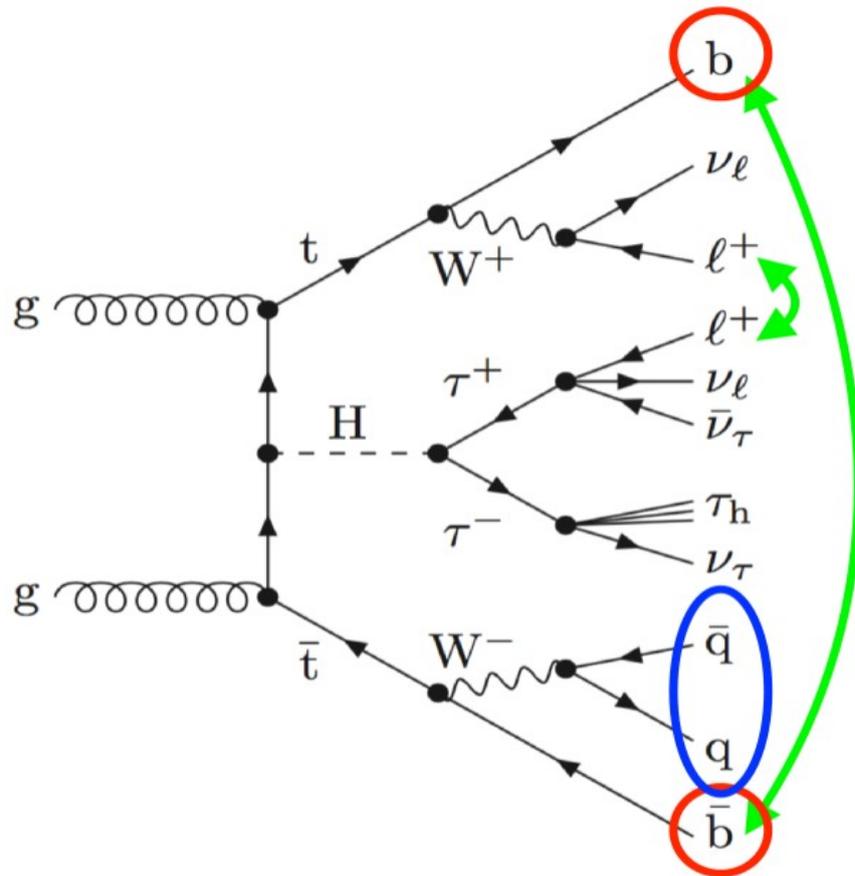
tt+jets: others jets (gluons) can be present in the event (Initial State Radiation)



ttW: g misidentified as τ_h
with $W \rightarrow$ lepton ...
and others possibilities ...

The definition of the final state drives the S/B

Permutations



- Problem to associate the b -jets measures to the (b, \bar{b}) of final state, idem for the 2 leptons
- 4 permutations (green arrows)
- 1 missing q or \bar{q} in the reconstruction:
 - 2 more integration variables (direction)
 - (4 x) permutations on all possible “light jets”

Integration space dimension	ttH, $H \rightarrow \tau\tau$	ttZ, $Z \rightarrow \tau\tau$	ttW, $W \rightarrow l\nu$	tt+jets
no missing jet	5	5	6	4
with missing jets	7	7	8	6

MEM MPI implementation

- PDF: LHAPDF library
- ME computation: MadGraph5 2.2.1 *code generator* (C++)
- ROOT: I/O, Lorentz/geometric arithmetics
- Integration: VEGAS algorithm (GSL)
- Mean CPU time per event 13 min.
- Parallel version (MPI) to tune the analysis method (T. Strebler)
- One run takes several days on 200-400 physical cores

“Daily used” of MEM-MPI on 400-core platform

OpenCL Implementation

Requirements:

Aggregate all the computing powers of the \neq nodes (MPI + OCL)

Benefit of all device computing power, including CPUs

→ several OCL queues in a node

VEGAS: keep the computation of the chi-square (GSL)

Features:

- Minimize host/devices communications:
 - 1 event is assigned to a queue/device
 - All the integration part (VEGAS) must be done inside the device (including reductions)
- No blocking calls (kernels, communications) → OCL events
- Minimize the synchronization points (reductions)

OCL Kernels

Main kernel (one Vegas iteration) :

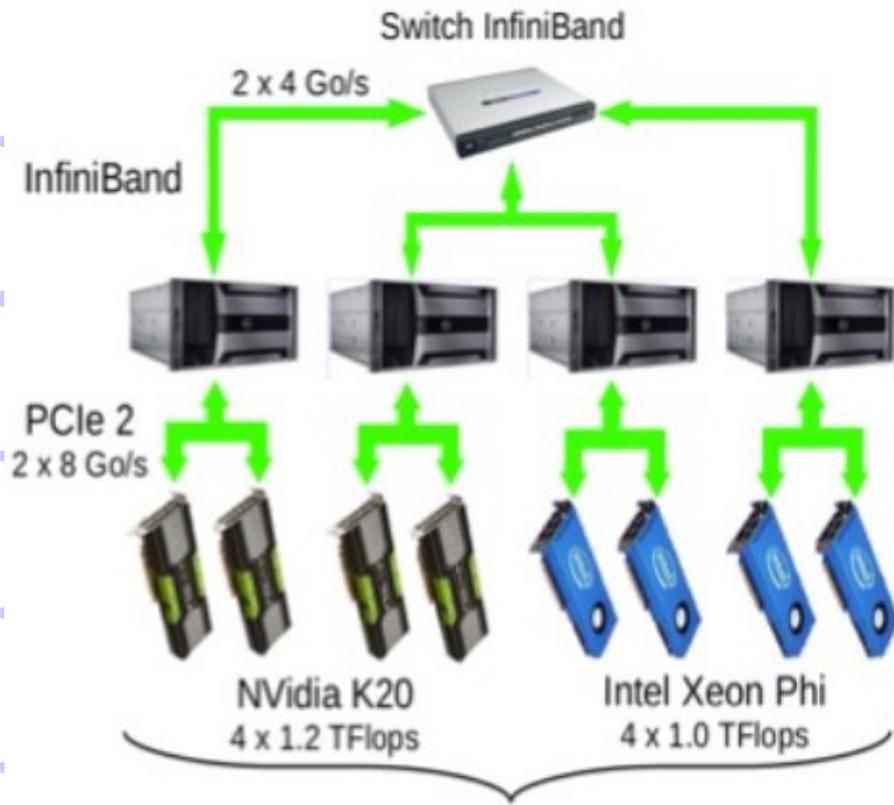
- We developed MadGraph extension to generate the OCL kernel codes
- LHAPDF lib.: Fortran to C-kernel translation
- ROOT tools: Lorentz/geometric arithmetics

→ big kernels (10-20 x 10³ lines)

Data/Work flow

```
host→devices( Config )
Loop on Events
host→devices( Event )
Loop on Permutations
Loop on IntegrationTypes
host→devices( VegasState )
kVegasSetUp()
Loop on  $\chi^2$  // for Vegas
kVegasCompute()
kVegasFinalize()
Devices→host( VegasState )
```

- Config: LHAPDF, MagGraph (sub) processes, etc.
- Event: coming from MPI msg → device/queue
- IntegrationTypes loop: asynchronous mode non-blocking calls (cl::Events)



Each node

2 x Intel E5-2650: 2 GHz, **16** physical cores, with AVX (4 doubles), 64 Go memory

Interconnection

switch InfiniBand

Devices

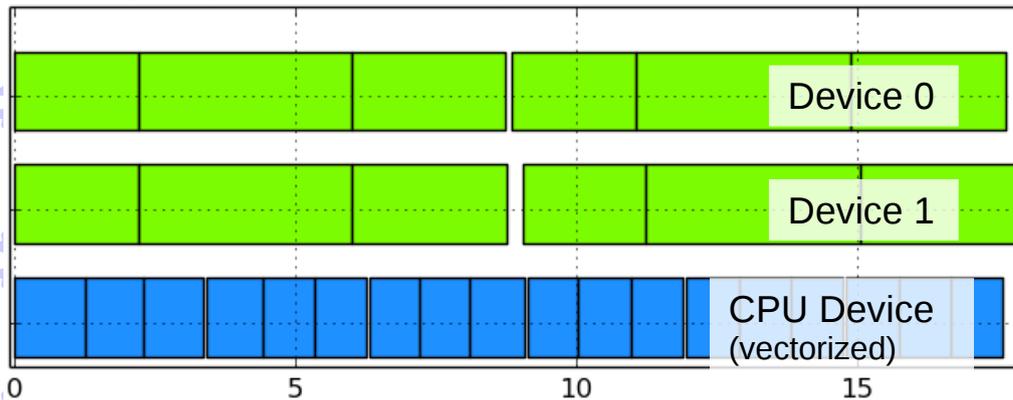
Nvidia K20, Titans

Xeon Phi

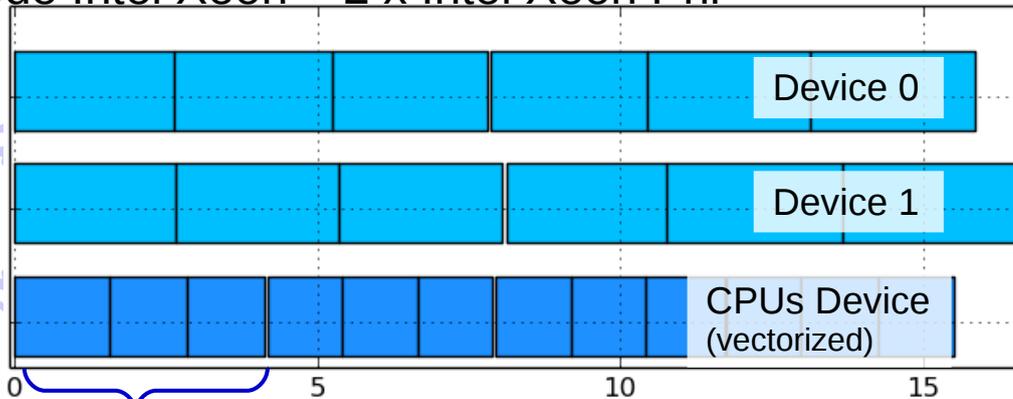
AMD FirePro S9170

Load-balancing in a single node

Node Intel Xeon + 2 x NVidia K20



Node Intel Xeon + 2 x Intel Xeon Phi



1 event: 3 possible permutations → 3 x ttH integrations

The NodeScheduler feeds all the Devices/CPUs inside one node

Same event, 3 permutations, ttH (signal hypothesis)

- NVidia specificities:
 - Buffer must be “pinned” in the memory **not to block** the copy call
- In OCL for NVidia GPUs:
created with `CL_MEM_ALLOC_HOST_PTR`
allocated `enqueueMapBuffer()`

Preliminary Performance on a single device

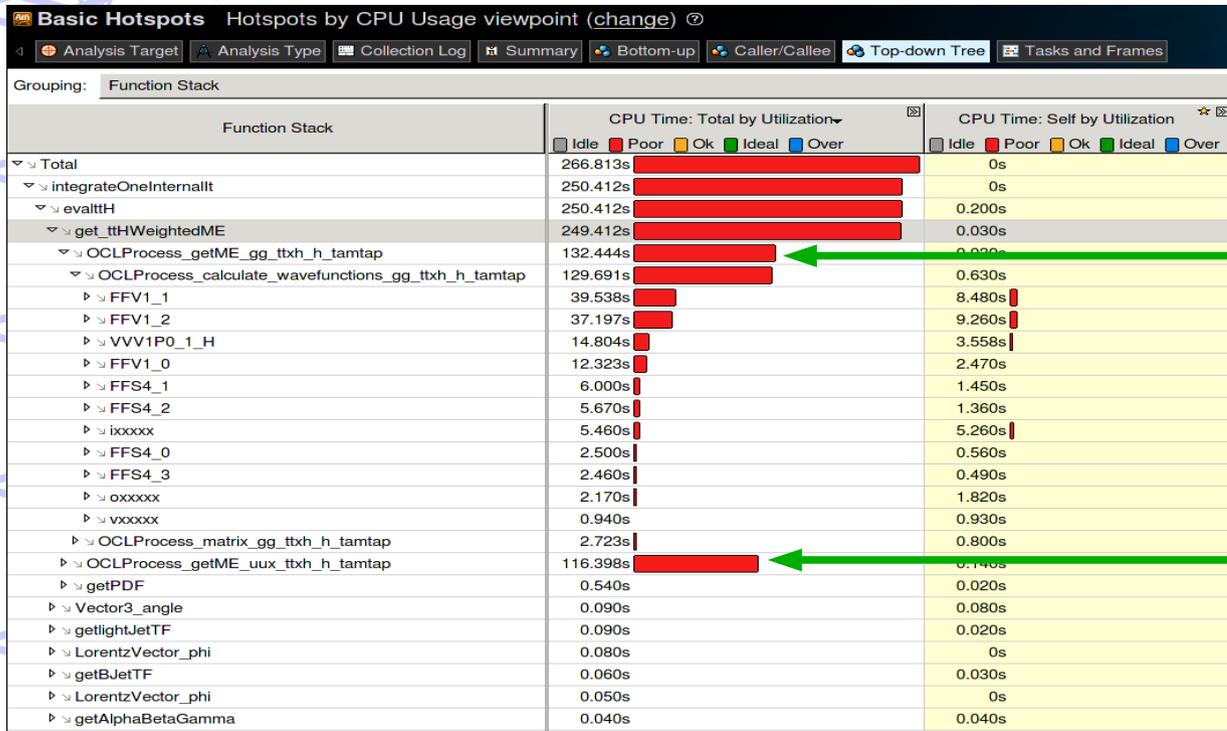
One event, 3 permutations, ttH with 15552 integration points

	C++ -O3	OCL K20	OCL X.Phi	OCL CPUs	OCL AMD
Time (s)	91.6	8.74	6.90	3.16	-
Speedup	1.00	10.74	13.3	29.0	-
Speedup with 16 MPI proc.	1.00	0.66	0.83	1.81	-

We obtained better performance on smaller kernels (simplified ME, speedup > 50 on K20)

How to get performance analysis of kernels with OCL ?

Performance analysis tools

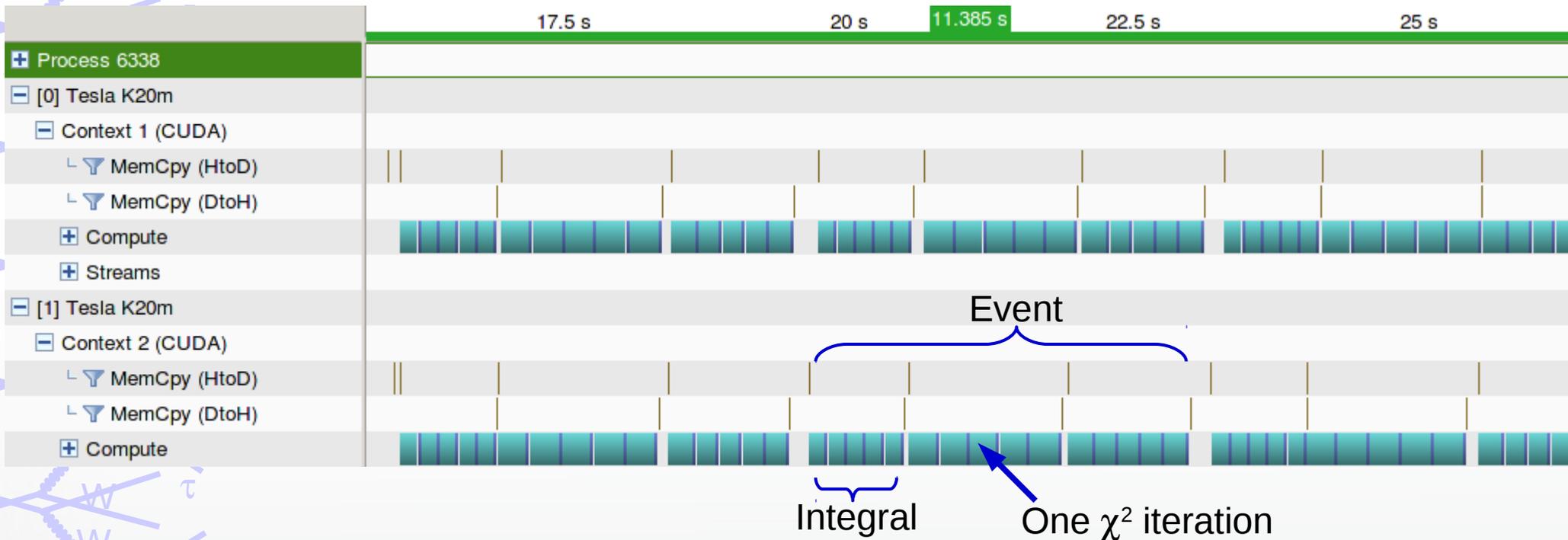


VTune analysis on kernels (w/o OCL): workload dominated by the ME computation (green arrows)

- CodeXL (AMD) works well for simple kernels (compiler) ...
- NVVP (Nvidia) not allowed with OCL ...
- VTune (Intel) with OCL (CPU) ... difficult

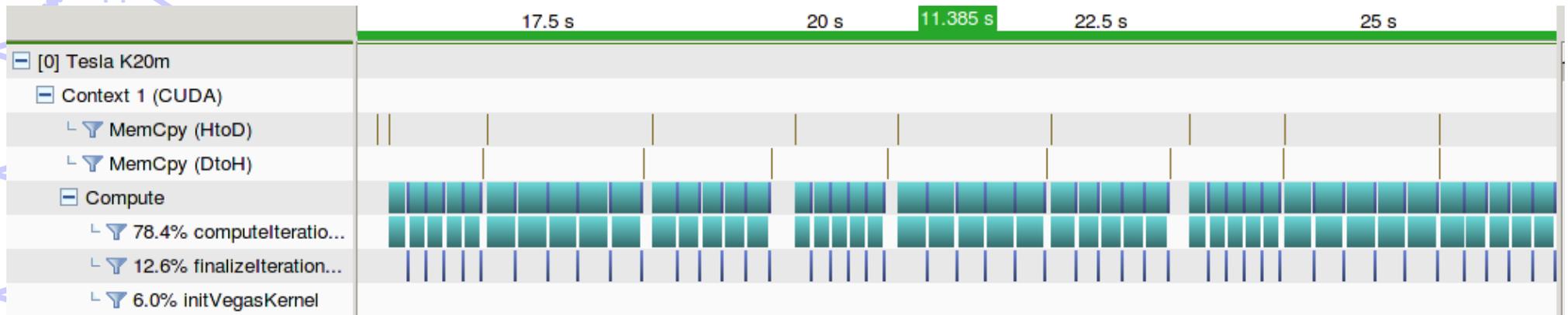
CL-CUDA/cl.hpp

- LLR development, motivations: for debugging, to preserve our OCL developments, ...
- Principles : routes `<cl.hpp>` calls/methods to CUDA calls. Handle heterogeneous devices
- Change: `#include <CL/cl.hpp>` by `#include <CL-CUDA/cl.hpp>` and `-lcuda`



Kernel performance

- Good device occupancy (asynch. mechanisms)
- Host ↔ Devices copies are negligible
- Kernel performance: ~2 x faster with CL-CUDA

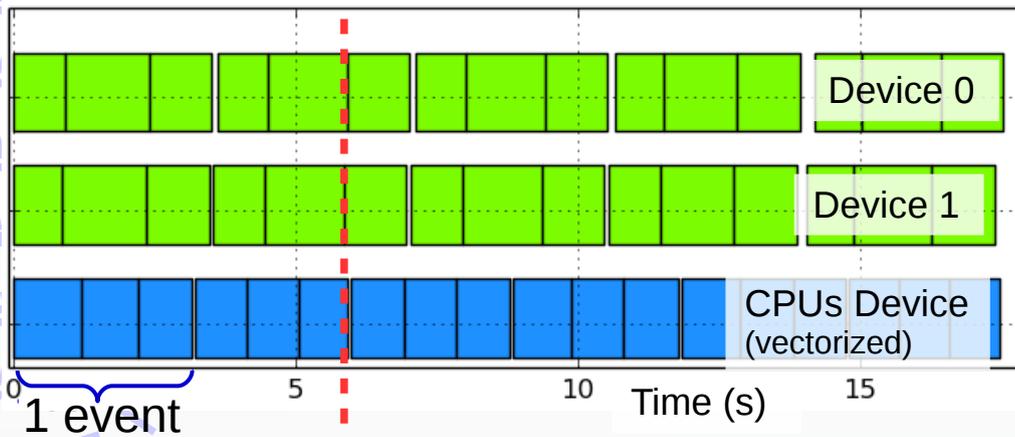


- Kernel performance is limited by the use of 255 registers per threads
- `--maxrregcount` doesn't improve performance
- VTune targets MadGraph expressions
- Better use of `__local` (`__shared__`) space memory to avoid register spilling and/or to reduce register use by thread

Conclusion/Perspectives

- CL-CUDA takes advantage of both CUDA/Intel-OCL compilers, speedup ~ 5.5 for one K20's node (speedup ~ 90 compared with a single MPI process)
- Optimization: better use of data locality (`__local`)

Node Intel Xeon + 2 x NVidia K20



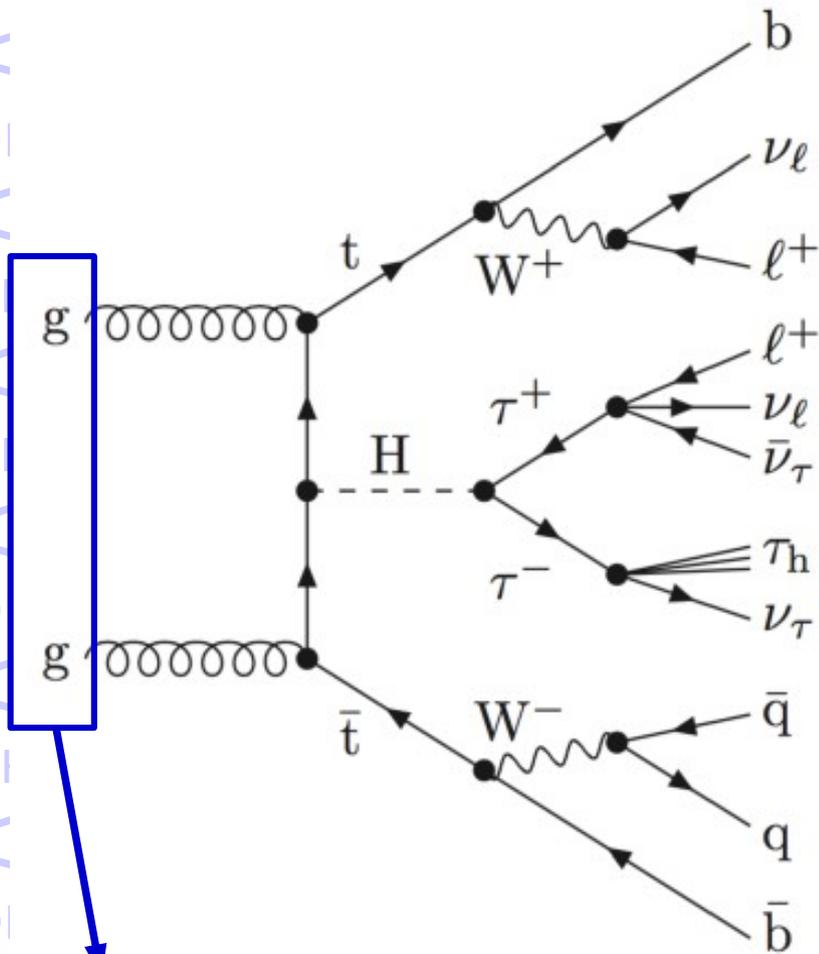
Time per ev. with 16 MPI processes

Next steps:

- Physics: include ttW, tt+jets in the next weeks
- Production on 10 nodes x 2 K80s (CC-IN2P3)
- Allows to compute more accurately integrals (dim. > 5 , 15k points)
- Evaluate on recent platforms: NVidia Pascal, Intel KNL (GENCI)
- Evaluate OpenMP 4.x

Backup

Final States: signal

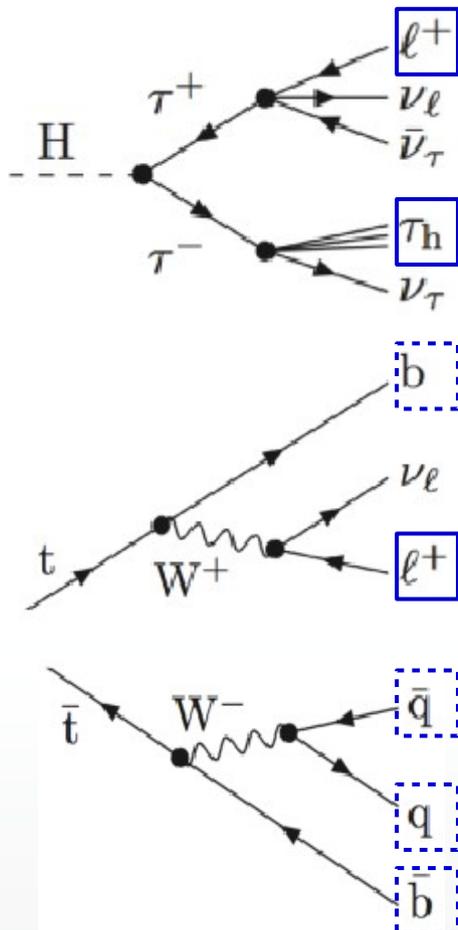


Others initial states ($u\bar{u}$, $d\bar{d}$, $c\bar{c}$, $s\bar{s}$) leading to the same final state ($t\bar{t}H$, $H\tau\tau$)

- Final state chosen to optimize the “S/B” ratio:
 - 2 tops production (studied channel)
 - Higgs boson decaying in 2 τ 's:
 - one τ decay into hadrons (hadronic system),
 - other τ decays in a lepton
 - Top quarks decays:
 - One decay in a single lepton+b(+neutrino)
 - One decays in quarks $q\bar{q} + b$
- And the 2 leptons with same sign

Integration variables for ttH/Z

ttH, H→ττ: 3 x 11 variables with the measure constrains, the mass invariant constrains and the momentum conservation → 5 integration variables



- **Higgs/Z decay to ττ**
2 integration var.: $|\vec{\tau}^+|, \cos(\theta_{\tau\tau})$
- **Leptonic top decay**
Setting ν direction (θ_{lv}, η_{lv}) → E_v → E_b
2 integration var.: neutrino's direction (θ_{lv}, η_{lv})
- **Hadronic top decay**
Setting E_q → E_{qbar} → E_{bbar}
1 integration var.: E_q variable

Example: mass invariant for (W, q, q_{bar}):

$$m_W^2 = E_W^2 - \vec{P}_W^2 = (E_q + E_{qbar})^2 - (\vec{P}_q + \vec{P}_{qbar})^2$$