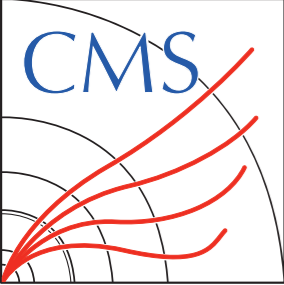


MG4GPU CMS INTEGRATION STATUS

FOR CMS-MG JOINT MEETING

24.07.26

Jin Choi, Sapta Bhattacharya, Stefan Roiser, Andrea Valassi,
Olivier Mattelaer, Zenny Wettersten

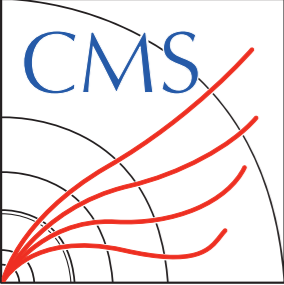


CONTENTS



❖ **Gridpack Production**

❖ **Event Generation**



GRIDPACK PRODUCTION

❖ Test Process

✓ Using Drell-Yan process with jet-binned configuration

⇒ Most common bkg. in CMS Analyses

⇒ No Madspin dep. (Not implemented in GPU yet)

⇒ Growing complexity - Possible to check improvements & bottlenecks w.r.t. final state multiplicity

	DY+0j	DY+1j	DY+2j	DY+3j	DY+4j
diagrams	30	180	3120	27600	412560
processes	15	45	285	435	1455

```

1  import model sm-no_b_mass
2  set nb_core 10
3
4  define p = u d c s b u~ d~ c~ s~ b~ g
5  define j = p
6  define ell+ = e+ mu+ ta+
7  define ell- = e- mu- ta-
8  define nu = ve vm vt
9  define nubar = ve~ vm~ vt~
10
11 generate p p > ell+ ell- j j j j @0
12
13 output madevent_gpu DY4j_LO_5f_CUDA -nojpeg
    
```

DY4j_LO_5f_CUDA_proc_card.dat

```

1  set sde_strategy 1
2  set vector_size 8192
3  set cudacpp_backend CUDA
    
```

DY4j_LO_5f_CUDA_customizecards.dat

❖ Test Process

- ✓ Using Drell-Yan jet-binned samples [\[run cards\]](#)
 - ⇒ Most common bkg. in CMS Analyses
 - ⇒ No Madspin dep. (Not implemented to GPU yet)
 - ⇒ Growing complexity - Possible to check improvements & bottlenecks w.r.t. final state multiplicity

	DY+0j	DY+1j	DY+2j	DY+3j	DY+4j
diagrams	30	180	3120	27600	412560
processes	15	45	285	435	1455

❖ Comparison

- ✓ Cross-section: for least validation
- ✓ Timing: Using time command, compared full production time & ME integration time

- ❖ **Integrating MG4GPU to CMS-genproduction** [[genproduction/mg4gpu](#)]
- ✓ Based on the master branch(for RUN3 production) - updated patches for MG352/mg4gpu
- ✓ Workflow: Environment setup(e.g. CMSSW / CUDA) - download mg4gpu - apply patches
- run ./gridpack_generation.sh
- ✓ Previously used git clone for downloading mg4gpu: large repo, takes ~ 10 min. to clone
⇒ Compressed the repo in EOS area, untar the repo rather than downloading: ~ 4 min.
- ✓ Two major patches for mg4gpu side (for version [mg4gpu_2024-03-14.tar.gz](#))

```

1  diff --git a/madgraph/interface/madevent_interface.py b/madgraph/interface/madevent_interface.py
2  index 8c509e83f..e6e7bd0dc 100755
3  --- a/madgraph/interface/madevent_interface.py
4  +++ b/madgraph/interface/madevent_interface.py
5  @@ -3966,7 +3966,8 @@ Beware that this can be dangerous for local multicore runs. """
6      Pdir = set([os.path.dirname(G) for G in Gdir])
7      for P in Pdir:
8          allG = misc.glob('G*', path=P)
9  -         for G in allG:
10 +         filG = [f for f in allG if not os.path.basename(f).startswith('Gpu')]
11 +         for G in filG:
12             if pjoin(P, G) not in Gdir:
13                 logger.debug('removing %s', pjoin(P,G))
14                 shutil.rmtree(pjoin(P,G))

```

some files start with "Gpu*" and erased when clearing some directories like G3*...

- ❖ **Integrating MG4GPU to CMS-genproduction** [[genproduction/mg4gpu](#)]
- ✓ Based on the master branch(for RUN3 production) - updated patches for MG352/mg4gpu
- ✓ Workflow: Environment setup(e.g. CMSSW / CUDA) - download mg4gpu - apply patches
- run ./gridpack_generation.sh
- ✓ Previously used git clone for downloading mg4gpu: large repo, takes ~ 10 min. to clone
⇒ Compressed the repo in EOS area, untar the repo rather than downloading: ~ 4 min.
- ✓ Two major patches for mg4gpu side (for version [mg4gpu_2024-03-14.tar.gz](#))

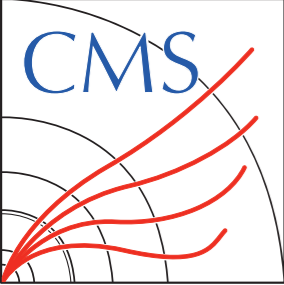
```

1  diff --git a/madgraph/various/systematics.py b/madgraph/various/systematics.py
2  index 28eaed0..5f787de 100644
3  --- a/madgraph/various/systematics.py
4  +++ b/madgraph/various/systematics.py
5  @@ -169,7 +169,7 @@ class Systematics(object):
6      self.orig_ion_pdf = False
7      self.ion_scaling = ion_scaling
8      self.only_beam = only_beam
9  -   if isinstance(self.banner.run_card, banner_mod.RunCardLO):
10 +   if self.banner.run_card.LO:
11         self.is_lo = True
12         if not self.banner.run_card['use_syst']:
13             raise SystematicsError('The events have not been generated with use_syst=True. Cannot evaluate systematics error on thes

```

self.banner.run_card does not work with use_syst option





ENVIRONMENTS

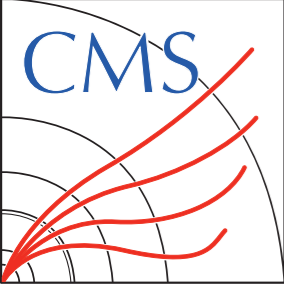


❖ HPCs

- ✓ Ixplus800(GPU): AMD EPYC 7313 16-core processor (AVX2 support), A100 GPU → repeatedly halted
- ✓ SNU-server: Intel(R) Xeon(R) CPU E5-2699 v3 (72 cores, AVX2 support), no GPU
→ tested FORTRAN DY+4j gridpacks
- ✓ **Ixplus condor**: Randomly matched to
 - w/o GPU: Intel(R) Xeon(R) CPU E5-2650 v4 @ 2.20GHz (48 CPUs)
Intel(R) Xeon(R) Silver 4216 CPU @ 2.10GHz (64 CPUs)
 - w/ GPU: AMD EPYC 7313 16-core processor (AVX2 support), A100 GPU

❖ Sidenotes

- ✓ For testing CPU usage in Ixplus condor,
randomly matches to the nodes with 48/64 cores + AVX2 supports
- ✓ There is 4 A100 GPU node but the gridpack production failed if there is multiple GPUs



ENVIRONMENTS

❖ HPCs



lx

OpenStack project with GPU flavors in pass-through mode



SP

This option is identical to the one described in the [Projects](#) section, except that GPU flavors will be assigned to your project. You can then launch instances with GPUs. The available flavors are:



lx

Flavor Name	GPU	RAM	vCPUs	Disk	Ephemeral	Comments
g1.xlarge	V100	16 GB	4	56 GB	96 GB	[^1], deprecated
g1.4xlarge	V100 (4x)	64 GB	16	80 GB	528 GB	[^1]
g2.xlarge	T4	16 GB	4	64 GB	192 GB	[^1], deprecated
g2.5xlarge	T4	168 GB	28	160 GB	1200 GB	[^1]
g3.xlarge	V100S	16 GB	4	64 GB	192 GB	[^1]
g3.4xlarge	V100S (4x)	64 GB	16	128 GB	896 GB	[^1]
g4.p1.40g	A100 (1x)	120 GB	16	600 GB	-	[^1], AMD CPUs
g4.p2.40g	A100 (2x)	240 GB	32	1200 GB	-	[^1], AMD CPUs
g4.p4.40g	A100 (4x)	480 GB	64	2400 GB	-	[^1], AMD CPUs



S



F

ra



T

(
f
tedly halted
(
(
F
(
VM
E
F
(
1

(
(
A

 **Least validation**

Compatible

	FORTRAN [pb]	CPP [pb]	CUDA [pb]
DY+0j	5704 \pm 10.11	5711 \pm 1.053	5710 \pm 1.484
DY+1j	3539 \pm 8.096	3535 \pm 1.263	3536 \pm 1.442
DY+2j	2228 \pm 3.143	2236 \pm 0.503	2237 \pm 0.4618
DY+3j	1375 \pm 1.265	1387 \pm 0.3515	1385 \pm 0.3288
DY+4j	883.4 \pm 0.3813	845.8 \pm 0.21	843.8 \pm 0.2022



A bit large errors / different xsecs for FORTRAN?

FORTRAN: Original MG
CPP: Vectorized CPU
CUDA: GPU

Results (full production)

	48 Intel FORTRAN	48 Intel, avx2 CPP	16 AMD + 1 A100 GPU CUDA
DY+0j	7m 59s	8m 38s	8m 1s
DY+1j	9m 27s	21m 3s	17m 20s
DY+2j	21m 24s	85m 6s	71m 25s
DY+3j	293m 38s	698m 41s	274m 44s
DY+4j	72 cms2 19362m 13s 13.5 days...	64 Intel, avx2 18509m 11s 12.8 days...	7682m 17s 5.3 days!

- ✓ FORTRAN / CPP run x3-4 processes in parallel
 - production time not differ much with current HPC configuration
- ✓ Can see x2-3 improvement in DY+4j: ME Integration is much faster for high multiplicity processes



PRODUCTION TIME

Results (full time)

	72 Intel cores	72 Intel cores	batch job 16 AMD cores + 1 A100 GPU
	FORTRAN	CPP	CUDA
DY+0j	11m 31s	6m 32s	8m 1s

```
INFO: Building madevent in madevent_interface.py with 'FORTRAN' matrix elements
INFO: P0_dxsx_taptamdxx
INFO: Building madevent in madevent_interface.py with 'FORTRAN' matrix elements
INFO: P0_uux_epemgg
INFO: Building madevent in madevent_interface.py with 'FORTRAN' matrix elements
INFO: P0_ddx_epemgg
INFO: Building madevent in madevent_interface.py with 'FORTRAN' matrix elements
INFO: P0_uux_taptamgg
INFO: Building madevent in madevent_interface.py with 'FORTRAN' matrix elements
INFO: P0_ddx_taptamgg
INFO: Building madevent in madevent_interface.py with 'FORTRAN' matrix elements
```

Compilation (ME)

```
INFO: Idle: 1, Running: 8, Completed: 281 [ current time: 16h42 ]
INFO: Idle: 0, Running: 9, Completed: 281 [ 0.02s ]
INFO: Idle: 0, Running: 6, Completed: 284 [ 3.3s ]
INFO: Idle: 0, Running: 3, Completed: 287 [ 10.5s ]
INFO: Idle: 0, Running: 0, Completed: 290 [ 17.7s ]
INFO: Idle: 0, Running: 0, Completed: 290 [ 17.7s ]
```

Execution (ME)

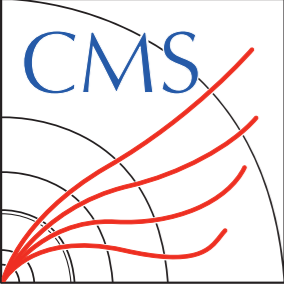
```
sum of cpu time of last step: 58m04s
improvement can only be seen in D114j...
```

ys?

Results (ME integration)

	48 Intel	48 Intel, avx2	16 AMD + 1 A100 GPU
	FORTRAN	CPP	CUDA
DY+0j	2.4s	1m 5s	17.7s
DY+1j	12.1s	2m 59s	31.6s
DY+2j	38.5s	8m 1s	2m 29s
DY+3j	3h 33m	6h 30m	33m 34s
DY+4j	72 cms2 315h 38m	64 Intel, avx2 244h 24m	108h 46m

- ✓ FORTRAN / CPP run x3-4 processes in parallel
 - production time not differ much with current HPC configuration
- ✓ Improvement starts with DY+3j (FORTRAN vs CUDA), ~x3 for DY+4j

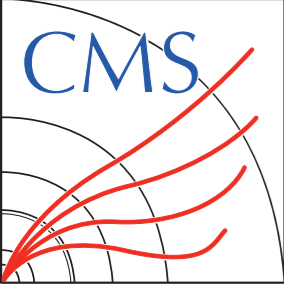


SUMMARY

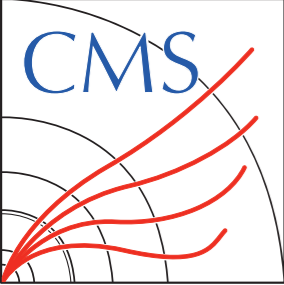


❖ **Comparing timing estimations for FORTRAN/CPP/CUDA**

- ✓ Compared HPC based improvements for gridpack production, for DY+0/1/2/3/4j LO
- ✓ Major bottleneck is **compilation time** for CUDA
- ✓ Both compilation and execution slow in CPP?
- ✓ With current usage, expecting highest gain in processes with **small no. of procces / ≥ 6 final states**



EVENT GENERATION



EVENT GENERATION



❖ From CMS gridpacks

✓ Basic command for evt generation would be:

```
./runcmsgrid.sh $NEVT $RANDOMSEED $NB_CORE
```

Shell

In this scripts, it does:

1. Set up production environment (i.e. CMSSW)
2. Modify madevent/Cards/me5_configuration.txt (e.g. nb_core...)
3. Determine the no. of evts to be generated in each iteration.
4. Generate events. i.e.

```
./bin/gridrun 5000 $SEED $GRAN
```

Shell

5. Combine events / Check the no. of evts / Add scale and PDF weights to LHE files.

❖ From CMS gridpacks

✓ Basic command for evt generation would be

```
./runcmsgrid.sh $NEVT $RANDOMSEED $NB_CORE
```

Shell

In this scripts, it does:

1. Set up production environment (i.e. CMSSW)
2. Modify madevent/Cards/me5_configuration.txt (e.g. nb_core...)
3. Determine the no. of evts to be generated in each iteration.
4. Generate events. i.e.

```
./bin/gridrun 5000 $SEED $GRAN
```

Shell

5. Combine events / Check the no. of evts / Add scale and PDF weights to LHE files.

Parallelize Event Generation?

- Current version of MG does not support nb_core option
- Results in the slides are based on nb_core=1

❖ From CMS gridpacks

✓ Basic command for evt generation would be

```
./runcmsgrid.sh $NEVT $RANDOMSEED $NB_CORE
```

Shell

In this scripts, it does:

1. Set up production environment (i.e. CMSSW)
2. Modify madevent/Cards/me5_configuration.txt (e.g. nb_core...)
3. **Determine the no. of evts** to be generated in each iteration.
4. Generate events. i.e.

```
./bin/gridrun 5000 $SEED $GRAN
```

Shell

5. Combine events / Check the no. of evts / Add scale and PDF weights to LHE files.

- Basic blocks are generating 5000 events in each iteration.
- For each iteration, need to prepare running directories again, e.g. copy & pasting ./madevent executable, ajobs, etc.
- Done by *process/madevent/bin/internal/restore_data*

❖ **restore_data: copy & pasting, untarring each subprocess directories**

✓ It takes a bit long time to prepare running directories before actual execution of madevent...

e.g. DY+4j has 1455(412560) processes(diagrams)

```

39 for i in `cat subproc.mg` ; do
40   cd $i
41   echo $i
42   rm -f ftn25 ftn26 >& /dev/null
43   if [[ -e $1_results.dat ]]; then
44     cp $1_results.dat results.dat >& /dev/null
45   else
46     cp results.dat $1_results.dat >& /dev/null
47   fi
48   for k in G* ; do
49     if [[ ! -d $k ]]; then
50       continue
51     fi
52     cd $k
53     for j in $1_results.dat ; do
54       if [[ -e $j ]]; then
55         cp $j results.dat
56       else
57         cp results.dat $j
58       fi
59     done
60     for j in $1_ftn26.gz ; do
61       if [[ -e $j ]]; then
62         rm -f ftn26 >& /dev/null
63         rm -f $1_ftn26 >& /dev/null
64         gunzip $j
65         cp $1_ftn26 ftn26
66         gzip $1_ftn26
67       fi
68     done
69     cd ../
70   done
71   cd ../
72 done

```



```

97 ./run.sh 5000 10
98 Now generating 5000 events with random seed 10 and granularity 1
99 *****
100 *
101 *           W E L C O M E  t o
102 *           M A D G R A P H 5 _ a M C @ N L O
103 *           M A D E V E N T
104 *
105 *
106 *           *
107 *         * * * * 5 * * * *
108 *           *
109 *
110 *
111 *           V E R S I O N  3 . 5 . 3 _ l o _ v e c t
112 *
113 *           The MadGraph5_aMC@NLO Development Team - Find us at
114 *           https://server06.fynu.ucl.ac.be/projects/madgraph
115 *
116 *           Type 'help' for in-line help.
117 *
118 * *****
119 INFO: load configuration from /srv/work/process/madevent/Cards/m
120 INFO: load configuration from /srv/work/mgbasedir/input/mg5_conf
121 INFO: load configuration from /srv/work/process/madevent/Cards/m
122 Using default text editor "vi". Set another one in ./input/mg5_c
123 No valid eps viewer found. Please set in ./input/mg5_configurati
124 No valid web browser found. Please set in ./input/mg5_configurat
125 WRITE GRIDCARD /srv/work/process/madevent
126 generate 5000 events
127 P0_gg_epemuux
128 P0_gg_epemddx
129 P0_gg_taptamuux
130 P0_gg_taptamddx
131 P0_gu_epemgu
132 P0_gd_epemgd
133 P0_gux_epemgux
134 P0_gdx_epemgdx
135 P0_gu_taptamgu
136 P0_gd_taptamgd
137 P0_gux_taptamgux
138 P0_gdx_taptamgdx
139 P0_uu_epemuu
140 P0_uux_epemuux
141 P0_dd_epemdd
142 P0_ddx_epemddx
143 P0_uux_epemuxux
144 P0_dxdx_epemdxdx
145 P0_ud_epemud
146 P0_uc_epemuc
147 P0_uux_epemddx
148 P0_uux_epemccx
149 P0_udx_epemudx

```

✓ Has been parallelized by MG team [[restore_data](#)]
 ~ 30m reduced for DY+4j after parallelization with 16 CPUs



EVENT GENERATION



❖ From CMS gridpacks

✓ Basic command for evt generation would be

```
./runcmsgrid.sh $NEVT $RANDOMSEED $NB_CORE
```

Shell

In this scripts, it does:

1. Set up production environment (i.e. CMSSW)
2. Modify madevent/Cards/me5_configuration.txt (e.g. nb_core...)
3. Determine the no. of evts to be generated in each iteration.
4. Generate events. i.e.

```
./bin/gridrun 5000 $SEED $GRAN
```

Shell

5. **Combine events / Check the no. of evts / Add scale and PDF weights to LHE files.**
This part does not take much time

❖ Patches

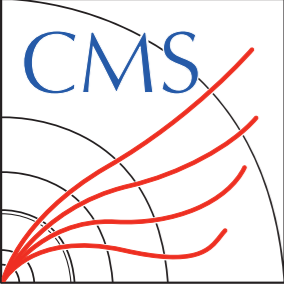
✓ After un-tarring gridpacks, two patches applied:

1. restore_data: To parallelize copy & past subprocesses directory - Utilize full CPU cores in machine
2. runcmsgrid.sh: Max Evts per iteration changed to 500k - avoid repeatedly call restore_data

```

75 max_events_per_iteration=500000
76
77 # define max event per iteration as 5000 if n_evt<45000 or n_evt/9 otherwise
78 #max_events_per_iteration=$(( $nevt > 5000*9 ? ($nevt / 9) + ($nevt % 9 > 0) : 5000 ))
79 # set starting variables
80 produced_lhe=0
81 run_counter=0
82 # if rnum allows, multiply by 10 to avoid multiple runs
83 # with the same seed across the workflow
84 run_random_start=$(( $rnum*10 ))
85 # otherwise don't change the seed and increase number of events as 10000 if n_evt<50000 or n_evt/9 otherwise
86 #if [ $run_random_start -gt "89999990" ]; then
87 #   run_random_start=$rnum
88 #   max_events_per_iteration=$(( $nevt > 10000*9 ? ($nevt / 9) + ($nevt % 9 > 0) : 10000 ))
89 #fi
90

```



GENERATION TIME

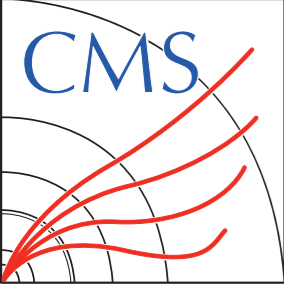


Results

✓ Producing 100K events w/ single thread

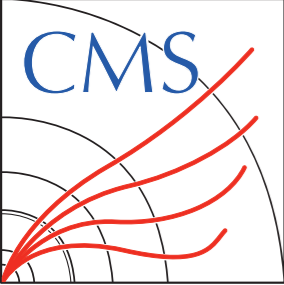
	FORTRAN	CPP	CUDA
DY+2j	80m 10s	59s 2s	40m 2s
DY+3j	130m 51s	153m 46s	101m 25s
DY+4j	never ends (>4000m)	1366m 49s	426m 54s

✓ Improvement starts with DY+2j, ~x10 faster for DY+4j



INTEGRATION IN CMS RELVAL WORKFLOW

- ❖ **Integrated the workflow updated with GPU gridpack in the standard CMS ReVal workflow**
- ✓ Used the existing DY+4Jet workflow and updated the path to the gridpack
- ✓ Workflow directed toward GPU node [\[ReVal\]](#)
- ✓ Successfully produced events [\[JIRA\]](#)
- ✓ The full chain with pythia works and the gridpacks are now virtually production ready

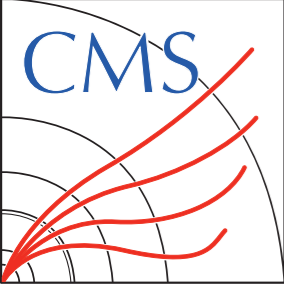


SUMMARY



❖ **Event Generation**

- ✓ Compared thread-based event generation time
- ✓ Repeatedly copy & pasting subprocess directories matters
- ✓ Can see x10 improvement comparing DY+4j FORTRAN vs CUDA
- ✓ Implementation with ReVal chain has been tested

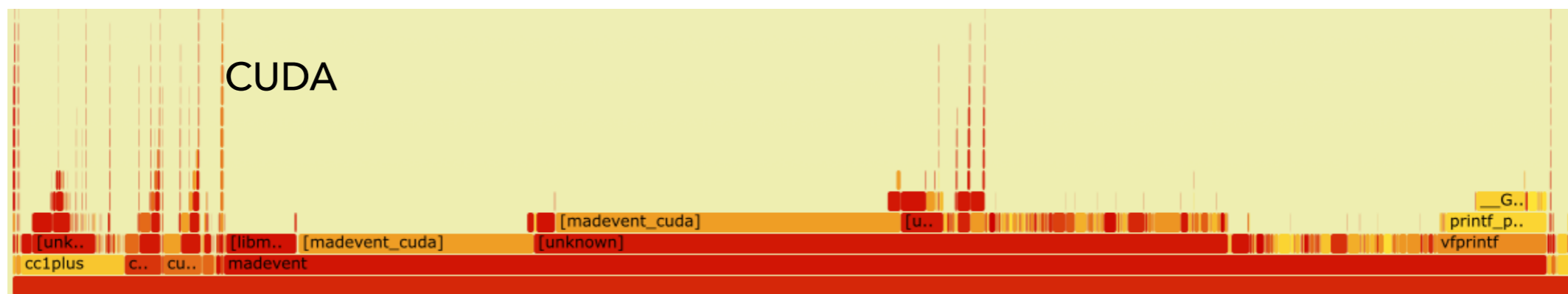
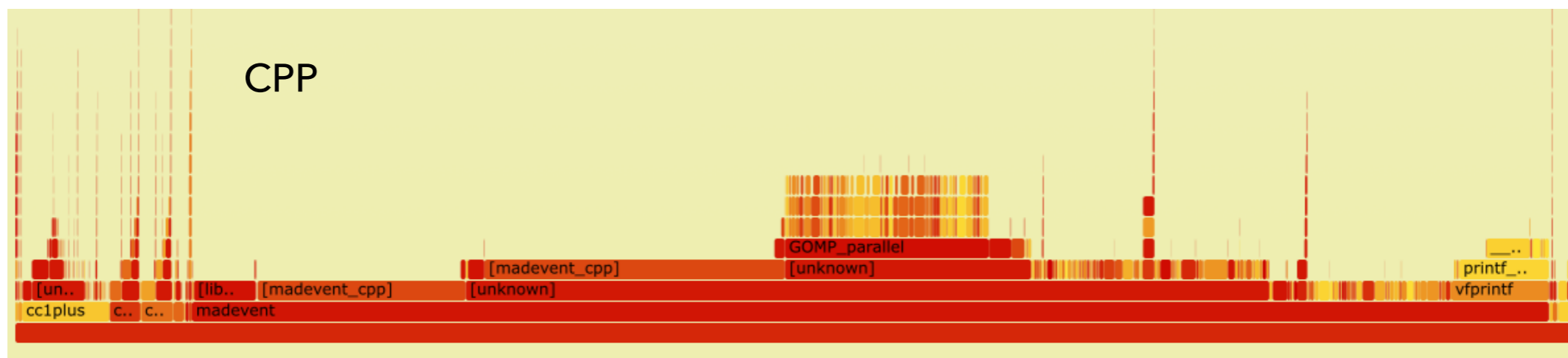
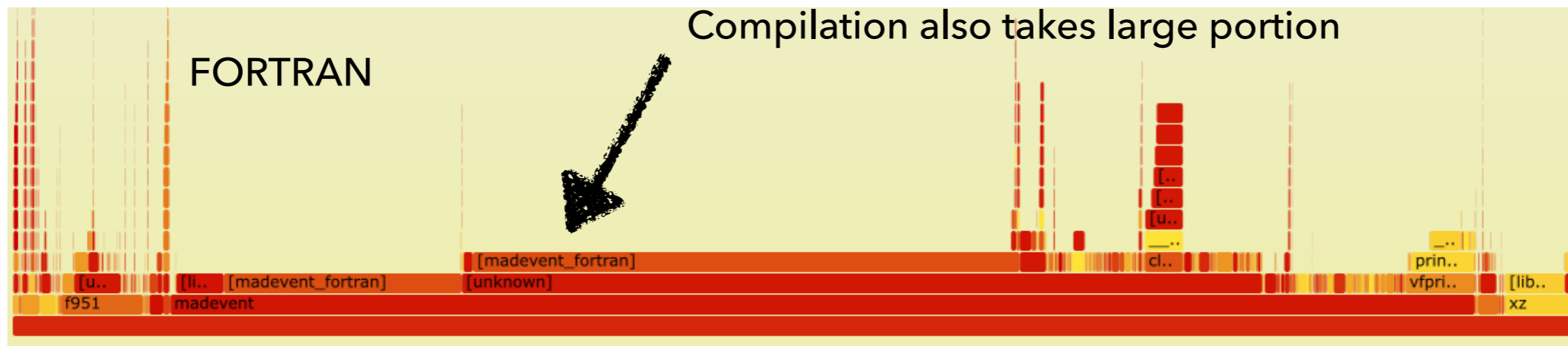


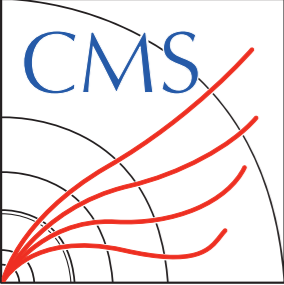
BACK UP

FLAMEGRAPHS - GRIDPACK PRODUCTION

DY+2j

Most of the time consuming part is still madevent...
Compilation also takes large portion

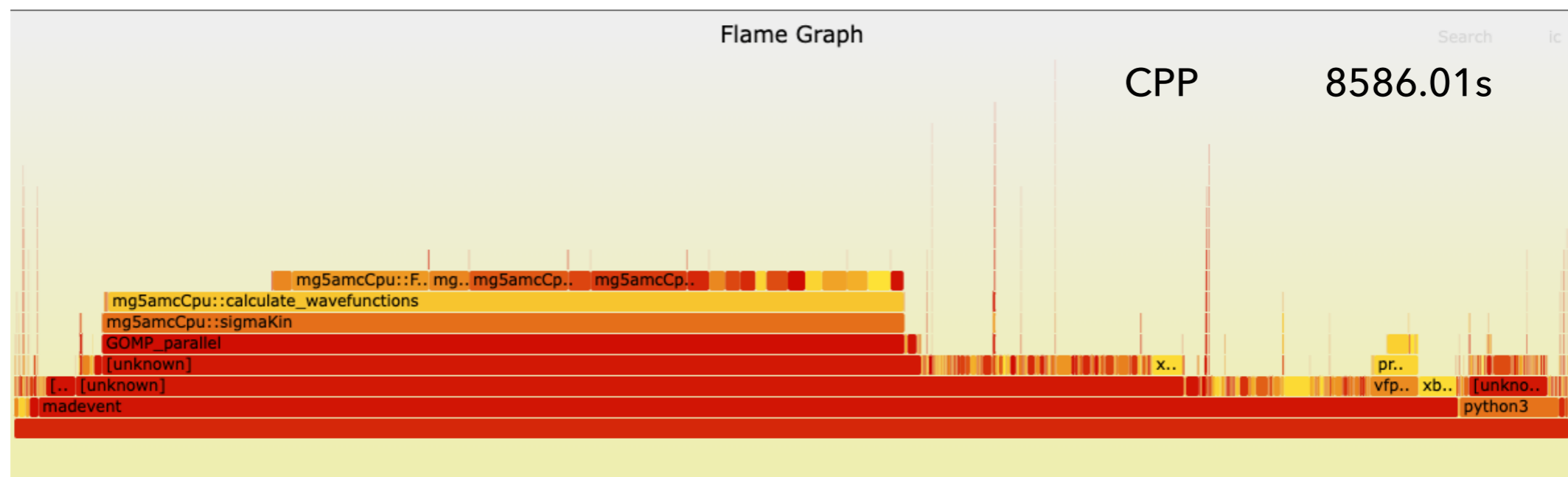
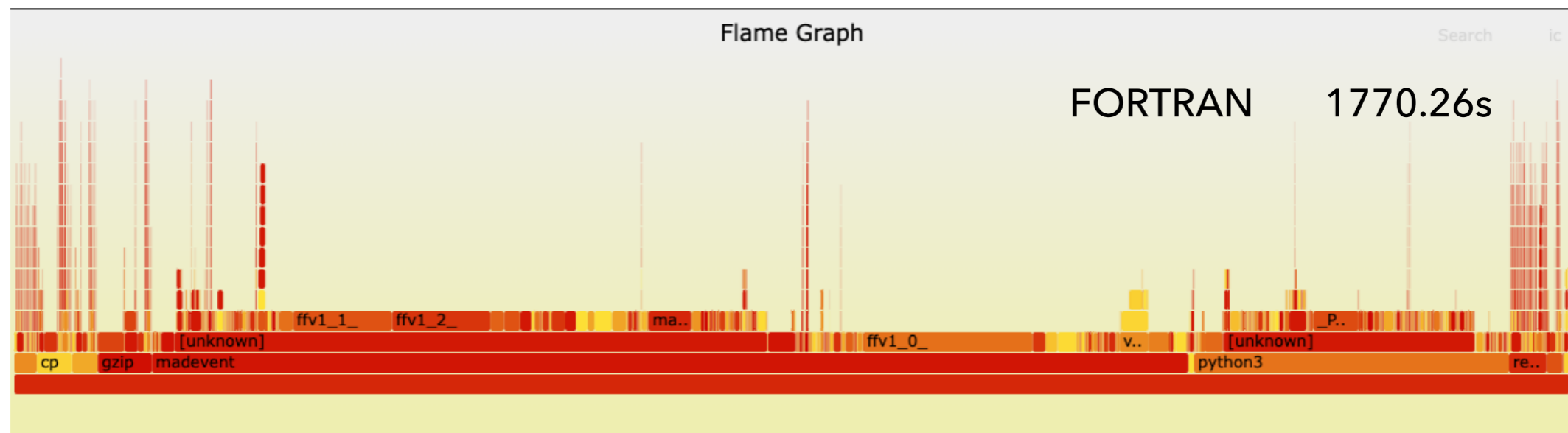




FLAMEGRAPHS - EVENT GENERATION

DY+3j (generating 20000 events)

svg files in [\[lxplus\]](#)





FLAMEGRAPHS - EVENT GENERATION

DY+3j (generating 20000 events)

svg files in [\[lxplus\]](#)

