

# Thread-scalable evaluation of multi-jet observables

Walter T. Giele<sup>a</sup>, Gerben C. Stavenga<sup>b</sup>, Jan Winter<sup>c</sup>

Theoretical Physics Department, Fermi National Accelerator Laboratory, P.O. Box 500, Batavia, IL 60510, USA

Received: 10 June 2010 / Revised: 3 June 2011 / Published online: 12 July 2011  
© The Author(s) 2011. This article is published with open access at [Springerlink.com](http://Springerlink.com)

**Abstract** We have implemented the leading-color  $n$ -gluon amplitudes using the Berends–Giele recursion relations on a multi-threaded GPU. Speed-up factors between 150 and 300 are obtained compared to the CPU-based implementation of the same event generator. In this first paper, we study the feasibility of a GPU-based event generator with an emphasis on the constraints imposed by the hardware. Some studies of Monte Carlo convergence and accuracy are presented for  $PP \rightarrow 2, \dots, 10$  jet observables using of the order of  $10^{11}$  events.

## 1 Introduction

Leading order (LO) parton-level Monte Carlos (MCs) play a prominent role in collider phenomenology [1–6]. As one needs to average the calculation of the observable over many events, the evaluation time for the event generation is a crucial issue in the development of LO parton level MCs. Furthermore, to make full use of the recent progress in the calculation of virtual corrections [7–10], fast tree-level event generators are needed for the calculation of the radiative contributions in a next-to-leading order MC.

One can use large-scale grids for the generation of the tree-level events. Such grids are expensive and need a large infrastructure. A more preferable solution would be to run the MC on a single, affordable workstation. As we will show this is possible using a massively parallel GPU. The NVIDIA Tesla computing processor is designed for numerical applications [11] and the CUDA C compiler [12] provides a familiar development environment. We will use the NVIDIA<sup>®</sup> Tesla<sup>™</sup> C1060 GPU throughout the paper.<sup>1</sup>

<sup>1</sup>We thank the LQCD Collaboration for giving us access to the Tesla GPU processors.

<sup>a</sup>e-mail: [giele@fnal.gov](mailto:giele@fnal.gov)

<sup>b</sup>e-mail: [stavenga@fnal.gov](mailto:stavenga@fnal.gov)

<sup>c</sup>e-mail: [jwinter@fnal.gov](mailto:jwinter@fnal.gov)

In this paper we will execute all steps that are needed for event generation on the GPU. These steps include the implementation of the unit-weight phase-space generator RAMBO [13], the evaluation of the strong coupling and parton density function using LHAPDF [14], the evaluation of the leading-color  $gg \rightarrow 2, \dots, 10$  gluon matrix elements at LO and the calculation of some observables. The CPU is tasked with calculating the distributions using the event weight and observables provided by the GPU. By utilizing memory with a fast access time only, considerable speed-ups are obtained in the event generation time. This memory is limited in size, requiring some coding effort. As the GPU chips are developing fast, we can enhance the capabilities of our parton-level generator in accordance.

In Refs. [15, 16] methods have been developed to evaluate multi-jet cross sections on GPUs within the framework of the HELAS matrix-element evaluator [17], which forms the basis of the MADGRAPH event generator [1]. The method is based on individual Feynman diagram evaluations. As such the scaling with the number of external particles of the scattering process is faster than factorial. Such an algorithm will have limited scalability properties, which cannot be compensated by deploying a large number of threads. Instead, an algorithm of polynomial complexity will have excellent scaling properties; its only limitation is the available fast-access memory size. Polynomial algorithms for the evaluation of ordered LO multi-parton matrix elements have been formulated in the form of Berends–Giele (BG) recursion relations [18]. For a leading-color generator, any Standard Model matrix element can be evaluated with an algorithm of polynomial complexity of degree 4 [19] or, by using more memory storage, of degree 3 [20]. For any fixed color expansion, the complexity remains polynomial. Therefore, we will use ordered recursive evaluations of the matrix elements instead of Feynman diagram evaluations.

In this paper we present a GPU-based implementation of all basic tools needed for a LO generator. In Sect. 2 we discuss the GPU and its hardware limitations. Accord-

**Table 1** The number of  $n$ -gluon events, which can be simultaneously executed on one MP (and is equal to  $2048/[n \times (n + 1)]$ ) and the number of available threads per event (equal to  $n \times (n + 1)/2$ ). The total number of events evaluated in parallel on the Tesla chip is  $30 \times (\text{events}/\text{MP})$ 

$n$	4	5	6	7	8	9	10	11	12
Events per MP	102	68	48	36	28	22	18	15	13
Threads per event	10	15	21	28	36	45	55	66	78

ing to these limitations, we will determine the optimal running configuration as a function of the number of gluons. The algorithmic implementation of the recursion algorithm and other tools such as phase-space generation, experimental cuts and parton density functions are discussed in Sect. 3. Finally, in Sect. 4 we put all pieces together and construct the leading-color LO parton-level generator capable of generating up to  $PP \rightarrow 10$  jets with sufficient statistics for serious phenomenology. The conclusions and outlook are given in Sect. 5.

## 2 Thread-scalable algorithms for event generators

Monte Carlo algorithms belong to a class of algorithms, which can be trivially parallelized, by dividing the events over the threads. Optimized for graphics processing, the GPU works by having many threads executing essentially the same instructions over different data. For a given class of events, e.g.  $n$ -gluon scattering, the only difference between the events is due to the external sources, i.e. the momentum and polarization four-vectors of each gluon defining the state of the external gluon. The recursive algorithm acts on these input sources in an identical manner. That is, each thread can execute the *same* processor instructions to calculate the matrix-element weight.

However, because of the hardware constraints such a straightforward approach is limited by the amount of available fast access memory. The GPU memory is independent from the CPU memory and divided into the off-chip global memory and the on-chip memory. This distinction is important as the off-chip memory is large (of the order of gigabytes) but slow to access by the threads. Therefore, we want to limit the access to the global memory by using it only for the transfer of results to the CPU memory. The on-chip memory is fast to access, but limited in size (of the order of tens of kilobytes). The first on-chip memory structures are the registers. Each thread has its own registers, which cannot be accessed by other threads. These registers are used within the algorithm for variable storage, function evaluations, etc. The other on-chip memory structure is shared memory, which is accessible to all the threads on a multi-processor (MP). The current GPUs are not yet optimize-able to one event per thread due to these shared memory and register constraints. With the next generation of GPUs the

shared memory will increase significantly, and we will reach the point at which we can evaluate one event per thread up to large multiplicities of gluons.

From this discussion the limitations are clear as each event requires a certain amount of the limited register and shared memory. For the optimal solution, we put the maximum number of events on one MP, such that the evaluation does not exceed the available on-chip memory. The resulting multiple threads per event can be used to *unroll* do-loops etc., thereby help speed up the evaluation. This optimal solution is dependent on the rapidly evolving hardware structure of the GPU chips.

By lowering the number of events per MP below the optimal solution, the number of available threads per event increases. However, this will not lead to an effective speed-up of the overall event generation as the total number of threads per GPU is fixed. Once the number of events to be used per MP has been determined, the GPU evaluation becomes scalable. The MC generator now simply scales with the number of available MPs on the GPU.

We have used the NVIDIA<sup>®</sup> Tesla<sup>™</sup> C1060 chip for the numerical evaluations in this paper. This chip consists of 30 MPs each capable of running up to 1024 threads. In the [Appendix](#) we have described the architecture of this computing processor in more detail; here it is sufficient to know that each MP has 16,384 32-bit registers and an internal shared memory of 16,384 bytes. Each thread is assigned its own registers from the pool. The compilation of the current MC implementation indicates that 35 registers per thread are needed. This gives us an upper maximum based solely on the use of registers of  $16,384/35 = 468$  threads per MP (each of which could potentially be used to evaluate one event). The momenta and current storage is of more concern. As we will see in the next section, for the evaluation of the  $n$ -gluon matrix element, we need to store  $n \times (n + 1)/2$  four-vectors in single (float) precision. This requires  $8 \times n \times (n + 1)$  bytes of shared memory per event on the MP.<sup>2</sup> The resulting maximum number of events per MP as a function of the number of gluons is given in Table 1. Note that up to 44-gluon scattering can be evaluated on the MP (albeit with only one event per MP). Beyond 44 gluons the shared memory is too small to store all the required four-vectors.

<sup>2</sup>A bit of calculus shows that if we have to store  $n \times (n + 1)/2$  real-valued four-vectors in single precision we need  $4 \times 4 \times n \times (n + 1)/2$  bytes of shared memory.

### 3 The implementation of the thread-scalable algorithm

Now that we have determined the optimal running configuration, i.e. the number of events per MP, we can implement the algorithm. We will describe the implementation of the `THREADED EVENTS SIMULATOR MC`, which we name `TESS MC`, for the `NVIDIA Tesla` chip.<sup>3</sup> As we have many threads available per event, we will use these threads to speed up the MC. In Fig. 1 we show the thread usage during different stages of the event generator for a  $2 \rightarrow n$  gluon process. The fraction of the evaluation time spent in each stage depends on the gluon multiplicity. For 4 gluons, we get 20%, 20%, 50% and 0% for the `RAMBO`, `PS-weight`, `ME-weight` and the `epilogue` phases, respectively. For 12 gluons, the time consumption divides up into 2%, 18%, 75% and 4% for the four phases.

The initialization phase (not shown in Fig. 1) consists of starting up the kernel on the GPU. This is taken care of by the `CUDA` run-time code and does essentially not depend on the number of threads it has to spawn. However it is a significant part when the total kernel time is small, as is for the 4-gluon case.

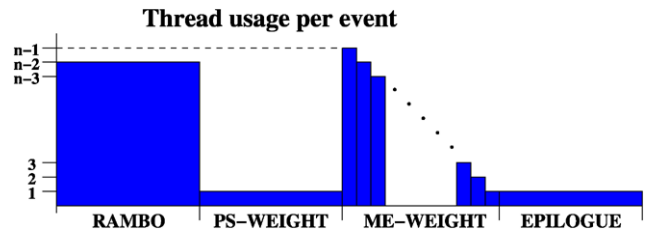
The kernel starts initiating the unit-weight phase-space generator `RAMBO`. On the CPU this algorithm grows linearly with  $n$  as we have to construct the  $n - 2$  outgoing momenta. On the GPU we can employ  $n - 2$  threads to simultaneously generate the outgoing momenta, making the `RAMBO` code in practice independent of  $n$ .<sup>4</sup>

After the momenta are generated, we have to calculate the strong coupling constant, the parton density functions and the observables. We also determine, if the event passes the canonical cuts. If the event fails the cuts, it is only flagged as such; the matrix-element weight will still be evaluated as this has no effect on the overall evaluation time. This means one can deviate from the chosen canonical cuts on the CPU during the histogramming phase if so desired. Note that we could in principle generate more events, which pass the cut before starting the calculation of the matrix-element weights. This should increase the performance of the Monte Carlo, at the cost of additional bookkeeping.

The evaluation of the strong coupling constant and parton density functions requires special attention. As we have used all shared memory for the four-vector storage of the gluon currents and momenta, we have to use the off-chip global memory to store the parton density and strong coupling constant information in the form of grids. Furthermore, interpolation is required between the grid points. To facilitate this,

<sup>3</sup>The `TESS MC` code can be downloaded from the website: <http://vircol.fnal.gov/TESS.html>.

<sup>4</sup>The `RAMBO` algorithm has some summation operations, which grow linearly with  $n$ , but this time scaling is very small compared to the overall evaluation time of the `RAMBO` algorithm.



**Fig. 1** The thread usage for an  $n$ -gluon event as the algorithm progresses through the stages of the event generation: flat phase-space generation, phase-space weight evaluation (including parton density functions and  $\alpha_S$ ), matrix-element evaluation and finalization phase

we use a special type of memory, the so-called texture memory. This off-chip memory was designed for graphics applications and performs hardware interpolations of the grid. Specifically, we set up a 1-dimensional grid for the strong coupling constant. The value of the strong coupling constant is stored as a function of the renormalization scale at integer values of the grid. For the 2-dimensional grid used by the parton density functions, the two dimensions are given by the factorization scale and the parton fractions. This parton density grid is directly obtained from `LHAPDF` [14]. After the grid initialization, the texture memory can be accessed by the GPU and its hardware will perform the appropriate linear interpolation between the grid points when accessing the grid using non-integer values. This way we have a very fast evaluation of the strong coupling and parton density functions taking only about 6% and 0.6% of the total GPU time for 4-gluon and 12-gluon processes, respectively.

The four-momenta are generated and the phase-space weight is determined, hence we have to evaluate the matrix-element weight next. This happens at the core of the event generator where we use recursion relations to compute these weights.

For this proof-of-concept program, we decided to use the recursion relation of Ref. [18] and restricted ourselves to the case of pure gluonic cross sections; quarks can be easily added at a later stage without changing the event generator in a fundamental way. The recursion relations we employ are given by

$$\begin{aligned}
 J_\mu[m, \dots, n] &= \frac{1}{K[m, \dots, n]^2} \left( \sum_{i=m}^{n-1} [J[m, \dots, i], J[i+1, \dots, n]]_\mu \right. \\
 &\quad \left. + \sum_{i=m}^{n-2} \sum_{j=i+1}^{n-1} \{ J[m, \dots, i], J[i+1, \dots, j], \right. \\
 &\quad \left. J[j+1, \dots, n] \}_\mu \right), \tag{1}
 \end{aligned}$$

where  $J_\mu[m, \dots, n]$  is a conserved four-vector current depending on the external gluons  $\{m, \dots, n\}$ . Furthermore, we have used the shorthand notations

$$\begin{aligned}
 K^\mu[m, \dots, n] &= \sum_{i=m}^n k_i^\mu, \\
 [J[\{a\}], J[\{b\}]]_\mu &= 2(J[\{a\}] \cdot K[\{b\}])J_\mu[\{b\}] \\
 &\quad - 2(K[\{a\}] \cdot J[\{b\}])J_\mu[\{a\}] \\
 &\quad + (J[\{a\}] \cdot J[\{b\}]) (K_\mu[\{a\}] - K_\mu[\{b\}]), \\
 \{J[\{a\}], J[\{b\}], J[\{c\}]\}_\mu &= 2(J[\{a\}] \cdot J[\{c\}])J_\mu[\{b\}] \\
 &\quad - (J[\{a\}] \cdot J[\{b\}])J_\mu[\{c\}] \\
 &\quad - (J[\{c\}] \cdot J[\{b\}])J_\mu[\{a\}],
 \end{aligned} \tag{2}$$

where the external gluon labeled  $i$  has momentum  $k_i^\mu$  and polarization state  $J^\mu[i]$ . These four-vectors form the initial conditions for the recursion relation. In addition to the  $n$  momenta, the recursion relation requires  $n \times (n - 1)/2$  four-vector currents to be stored giving a total storage of  $n \times (n + 1)/2$  four-currents per event.

The recursion relations have a polynomial complexity of order  $n^4$  for calculating the currents [19]. By exploiting the available threads for each event, we can reduce the algorithmic complexity of the BG recursion relation. The relation is easily thread-able, which enables us to lower the polynomial scaling of the evaluation time of the recursion relation to  $n^3$ . A full recursion for an  $n$ -gluon process is completed in  $n - 1$  steps. In the first step, we use  $n - 1$  threads (represented by the first of the columns in the ME-weight part of Fig. 1) to calculate the polarization vectors  $\{J[2], \dots, J[n]\}$  needed as a starting point in the recursion relation. We choose each polarization vector as a random unit vector orthogonal to the respective gluon momentum. By doing this, instead of employing the conventional helicity vectors, we obtain real-valued currents. This avoids complex multiplications and reduces the shared-memory usage, resulting in a significant time gain.

After the 1-currents have been determined, we use  $n - 2$  threads (represented by the second of the columns in the ME-weight part of Fig. 1) to calculate the 2-currents  $\{J[2, 3], J[3, 4], \dots, J[n - 1, n]\}$ . We continue with the  $n - 1$  steps until we have determined  $J[2, 3, \dots, n]$  at which point we can calculate the ordered amplitude and, hence, the matrix-element weight. Note that because we make use of the multiple threads we have reduced the computational effort from  $\mathcal{O}(n^4)$  to  $\mathcal{O}(n^3)$  complexity.

In principle we may be able to improve even further. The initial  $\mathcal{O}(n^4)$  growth of the one-threaded recursion relation

to calculate the  $J[2, 3, \dots, n]$  current can be reduced by rewriting the recursion relation as

$$\begin{aligned}
 J_\mu[m, \dots, n] &= \sum_{i=m}^{n-1} [(W[i + 1, \dots, n] \cdot J[m, \dots, i])_\mu \\
 &\quad - (W[m, \dots, i] \cdot J[i + 1, \dots, n])_\mu],
 \end{aligned} \tag{3}$$

where the tensor  $W_{\mu\nu}$  is defined as

$$\begin{aligned}
 W_{\mu\nu}[m, \dots, n] &= 2J_\mu[m, \dots, n]K_\nu[m, \dots, n] \\
 &\quad - K_\mu[m, \dots, n]J_\nu[m, \dots, n] \\
 &\quad + \sum_{i=m}^{n-1} (J_\mu[m, \dots, i]J_\nu[i + 1, \dots, n] \\
 &\quad - J_\mu[i + 1, \dots, n]J_\nu[m, \dots, i]).
 \end{aligned} \tag{4}$$

By undoing the nested summations in the second term of (1) we have lowered the complexity of the algorithm to  $\mathcal{O}(n^3)$ . However, this is only achieved at the cost of using significantly more storage. For each event, one would have to store  $n \times (n - 1)/2$  tensors of dimension  $4 \times 4$  in addition to the  $n \times (n + 1)/2$  momenta and current four-vectors. Up to  $n \approx 10$  the extra work of doing matrix multiplications together with the fact that the relative prefactor of the  $n^4$ -algorithm is small,  $1/4$ , compared to the  $n^3$ -algorithm actually make the  $n^3$ -algorithm slower than the  $n^4$ -algorithm. Moreover, the extra storage demand does not make the  $n^3$ -algorithm attractive for our GPU implementation.

From the current for  $n - 1$  gluons we then obtain the amplitude for the  $n$ -gluon matrix element by putting the off-shell leg on-shell, contracting in with the final polarization vector and symmetrizing over the gluons in the current. Specifically,

$$\begin{aligned}
 \mathcal{A}(1, \dots, n) &= \sum_{\pi}^{(n-1)!} \text{Tr}[T^{a_{\pi_1}} \dots T^{a_{\pi_{n-1}}} T^{a_n}] m(\pi_1, \dots, \pi_{n-1}, n),
 \end{aligned} \tag{5}$$

with

$$\begin{aligned}
 m(\pi_1, \dots, \pi_{n-1}, n) &= (J[\pi_1, \dots, \pi_{n-1}] \cdot J[n]) \\
 &\quad \times K^2[1, \dots, n - 1]_{K[1, \dots, n]=0}.
 \end{aligned} \tag{6}$$

Notice that for a given phase-space point, we have to perform the permutation sum requiring  $(n - 1)!$  steps to arrive at the full amplitude. This would immediately lead to a factorial growth in the computer time. We can circumvent the super-exponential sum over permutations in (5). In the

leading-color approximation this is easily accomplished and the color-summed squared amplitude is given by

$$\begin{aligned}
 & |\mathcal{A}(1, \dots, n)|^2 \\
 & \sim N_C^{n-2} (N_C^2 - 1) \\
 & \times \left( \sum_{\pi}^{(n-1)!} |m(\pi_1, \dots, \pi_{n-1}, n)|^2 + \mathcal{O}\left(\frac{1}{N_C^2}\right) \right). \quad (7)
 \end{aligned}$$

As we will use this matrix element in a  $2 \rightarrow n - 2$  gluon-scattering phase-space integration, we can use the symmetry of the final state to remove the permutation sum over the ordered amplitudes. In detail,

$$\begin{aligned}
 & d\sigma (PP \rightarrow n - 2 \text{ jets}) \\
 & = \int dx_1 dx_2 \frac{F_g(x_1) F_g(x_2)}{4p_1 \cdot p_2} \\
 & \times \frac{1}{(n - 2)!} \int d\Phi (p_1 p_2 \rightarrow p_3 \cdots p_n) \\
 & \times \sum_{\pi}^{(n-1)!} |m(\pi_1, \dots, \pi_{n-1}, n)|^2 \\
 & = \int dx_1 dx_2 \frac{F_g(x_1) F_g(x_2)}{4p_1 \cdot p_2} (n - 1) \\
 & \times \int d\Phi (p_1 p_2 \rightarrow p_3 \cdots p_n) |m(1, \sigma_2, \dots, \sigma_n)|^2, \quad (8)
 \end{aligned}$$

where  $p_1 = x_1 P_1$ ,  $p_2 = x_2 P_2$ , the parton density function is given by  $F_g(x)$ ,  $d\Phi$  is the phase-space integration measure and  $\{\sigma_2, \dots, \sigma_n\}$  is a permutation of the list  $\{2, \dots, n\}$  assigned randomly for each MC phase-space point evaluation.

Eventually, in the very last step of our threaded event simulation all the results are put together and returned to the CPU for processing.

By using the TESS MC, we can evaluate the differential  $n$ -jet cross sections in the leading-color approximation. The algorithm is of polynomial complexity and scales as  $n^3$  with the number  $n$  of gluons.

#### 4 A numerical study of the threaded events simulator

The first issue to study is the timing behavior of the TESS MC. We show our results for several gluon multiplicities in Fig. 2 where we plot the GPU timing as a function of events per MP (up to the respective maximum number of events per MP as determined earlier in Sect. 2 and given in Table 1). In a sweep each MP will evaluate a number of events in parallel using the threads. In principle the sweep time should be independent of the number of events evaluated by each MP as long as the shared-memory constraints are not exceeded, cf. Table 1. However, we have to execute a substantial amount of transcendental function calls per event, which

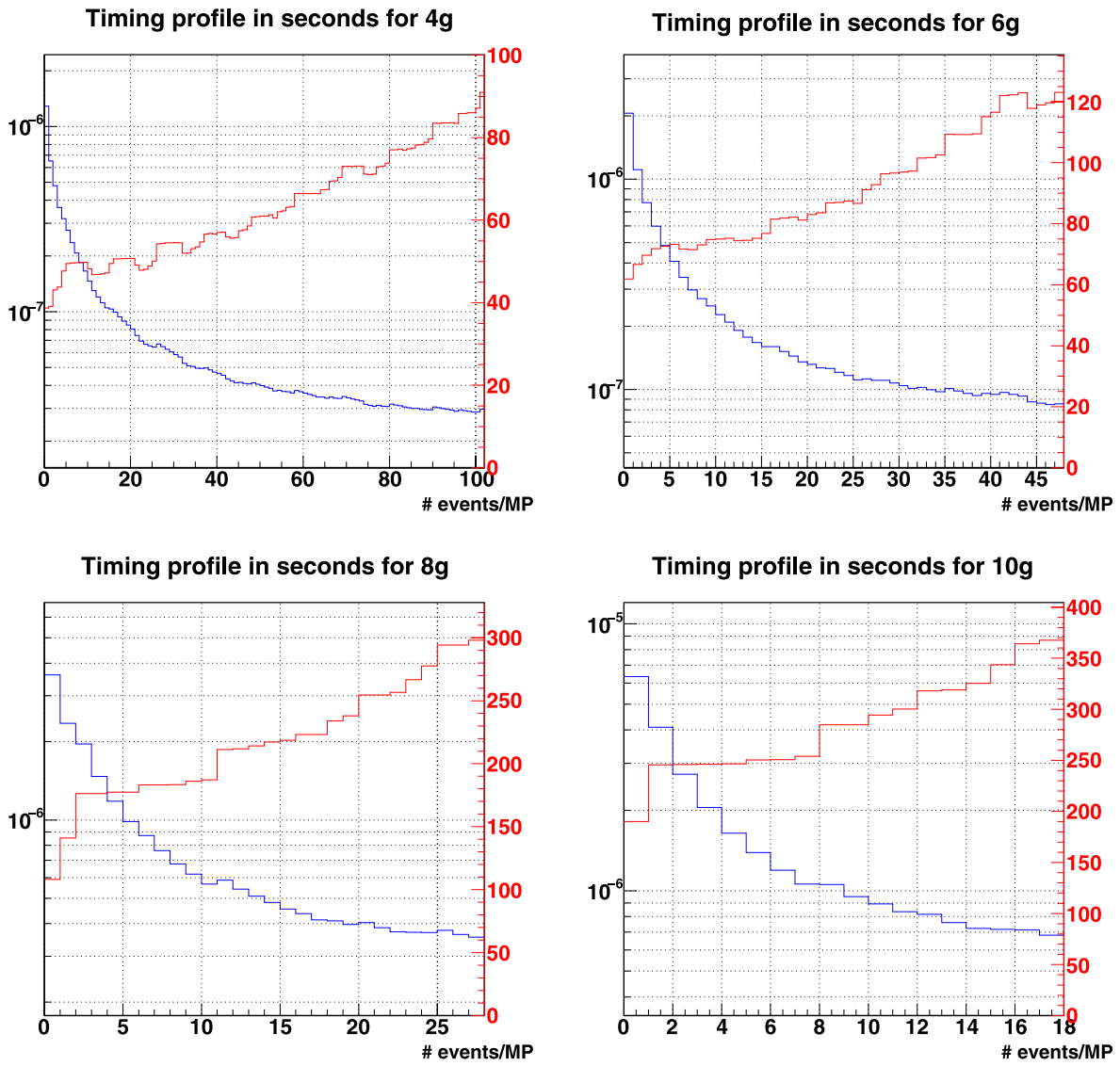
induces some queuing at the special-function units each MP uses for evaluating these functions. This queuing effect will increase as the number of events per MP rises and, hence, lead to a slower execution of the sweep. In Fig. 2, one can see this complicated timing behavior, which is controlled by the GPU hardware. As discussed the overall evaluation time increases with the number of events per MP, see the red curves in the plots. In fact, the increase of the overall evaluation time is overcome by the gain we achieve in evaluating more events per MP. The more relevant quantity therefore is the evaluation time per event, defined as the GPU evaluation time divided by the total number of generated events. As clearly indicated by the blue curves in Fig. 2, the time consumption per event steadily decreases as the number of events per MP increases. The best performance will be achieved by using the maximal number of events available per MP.

Now that we have determined the optimal running conditions, we give in Table 2 the evaluation time per event on the GPU compared to the evaluation time of the same algorithm when executed on the CPU.<sup>5</sup> As can be seen the speed-up in evaluation time is substantial, ranging from almost a factor of 300 for 4-gluon processes to a factor of around 150 for 12-gluon processes. Note that the speed-up is completely due to the fact that we evaluate in parallel 3060 and 390 events for the 4- and 12-gluon case, respectively. Because of the substantial time gains, a single GPU can replace a large grid of hundreds of CPUs.

As one example for rather unoptimized GPU code running, we have tested the performance of executing the events on the GPU sequentially: using only one event per sweep results in an event evaluation time, which is slower than the corresponding CPU evaluation time as given in Table 2. In particular, we found factors of 10 and 2 for the 4- and 12-gluon computations, respectively. Speaking of code optimization there are many factors affecting the performance of GPU computing. We tried to integrate the capabilities of the Tesla chip and CUDA framework into the program design. We however did not go as far as to optimize the code to exploit—instead of  $n - 1$  threads per event (cf. Fig. 1)—all available threads per event as shown in Table 2.

Also of interest is the scaling behavior of the algorithm. As expected, on the CPU it is simply polynomial scaling with a factor of 4 in the limit of a large number of gluons. We see from the table that this scaling is setting in quickly. The GPU algorithm scales with a factor of 3 as discussed in Sect. 3. However, as the number of gluons increases, the number of events per MP decreases. This makes the timing

<sup>5</sup>Beside holding the NVIDIA GPU Tesla chip, the workstation, which we used for our studies, comes with a quadruple core 3 GHz processor of the type AMD Phenom™ II X4 940. Here and elsewhere in this section we refer to this specific model when we mention the CPU.



**Fig. 2** (Color online) The horizontal axis is the number of events per MP in a sweep, giving a total number of  $30 \times$  (events per MP) evaluated events per sweep. The red curves used together with the vertical

axes on the right indicate the total GPU time in seconds for 1,000,000 sweeps. The blue curves depict the evaluation time of one event in seconds as labeled by the vertical axes on the left

**Table 2** The GPU and CPU evaluation times per event,  $T_n^{\text{GPU}}$  and  $T_n^{\text{CPU}}$ , given as a function of the number  $n$  of gluons for  $gg \rightarrow (n-2)g$  processes. The polynomial scaling measures are also shown, for the

GPU,  $P_n(3)$ , and for the CPU,  $P_n(4)$ . The  $P_n(m)$  are defined as  $P_n(m) = [(n-1)/n] \times \sqrt[m]{T_n/T_{n-1}}$ . The rightmost column finally displays the gain  $G_n = T_n^{\text{CPU}}/T_n^{\text{GPU}}$

$n$	$T_n^{\text{GPU}}$ (s)	$P_n(3)$	$T_n^{\text{CPU}}$ (s)	$P_n(4)$	$G_n$
4	$2.975 \times 10^{-8}$		$8.753 \times 10^{-6}$		294
5	$4.438 \times 10^{-8}$	0.91	$1.247 \times 10^{-5}$	0.87	281
6	$8.551 \times 10^{-8}$	1.03	$1.966 \times 10^{-5}$	0.93	230
7	$2.304 \times 10^{-7}$	1.19	$3.047 \times 10^{-5}$	0.96	132
8	$3.546 \times 10^{-7}$	1.01	$4.736 \times 10^{-5}$	0.98	133
9	$4.274 \times 10^{-7}$	0.94	$7.263 \times 10^{-5}$	0.99	170
10	$6.817 \times 10^{-7}$	1.05	$1.044 \times 10^{-4}$	0.99	153
11	$9.750 \times 10^{-7}$	1.02	$1.529 \times 10^{-4}$	1.00	157
12	$1.356 \times 10^{-6}$	1.02	$2.129 \times 10^{-4}$	1.00	158

more dependent on specific hardware issues. As can be seen from Table 2 the polynomial scaling is trending towards a factor of 3.

Given the fast evaluation of events, we can easily generate  $\mathcal{O}(10^{11})$  events for the calculation of the LO cross sections. With these large numbers of generated events, one has to carefully consider the performance of the random number generators. In our case this should cause no issues, since the number of generated random numbers is of the order of the square root of the generator's cycle length. However, as we average over  $\mathcal{O}(10^{11})$  numbers, care has to be taken concerning the loss of precision, which would result in a systematic underestimation of the cross section. This is demonstrated in the first graph of Fig. 3 where we have used a single-precision summation to calculate the 4-gluon cross section. As can be seen the effect becomes dramatic as the number of sweeps is rising and we end up with a totally wrong determination of the cross section. We avoid this problem by using the Kahan summation algorithm [21]. All other graphs of the figure are produced by following this procedure. These additional graphs display the convergence of the cross section estimates including their respective mean standard deviation as a function of the number of GPU sweeps. The vertical axis has been normalized to the respective best estimate of the cross section; all of which are listed in Table 3. For this study, we have used RAMBO as the momenta generator, therefore, a severe under-sampling of small phase-space regions with large weights may occur especially for larger gluon multiplicities. Because the RAMBO phase-space generation is flat and does not reflect the scattering amplitudes' strong dipole structure, such under-sampling effects are expected and cause the peaking behavior of our cross section estimates. Even with  $\mathcal{O}(10^{10})$  phase-space points an estimate of the 12-gluon cross section using the RAMBO event generator is quite unreliable and the mean standard deviation error estimate does not fully reflect the true uncertainty. In a further development step, one may implement a phase-space generator like SARGE [22], which is capable of adapting to the QCD antenna structures as occurring in the matrix elements. As pointed out in Ref. [6], this would resolve the phase-space integration issues we have seen here.

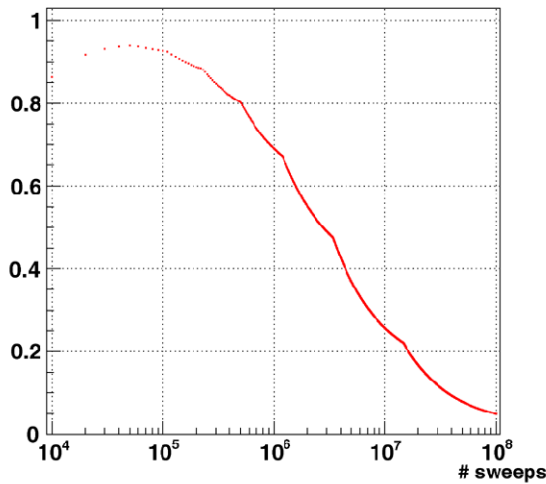
The convergence issues reflected in Fig. 3 should be taken into account when interpreting the uncertainties of our best cross section estimates, which are listed in Table 3. For these cross section calculations of  $gg \rightarrow (n-2)g$  scattering processes at a center-of-mass energy of 14 TeV, we have used the CTEQ6L1 parton density function set [23] as implemented in LHAPDF [14] with a fixed renormalization and factorization scale taken at  $M_Z = 91.188$  GeV. For the jet cuts, we have chosen  $p_T^{\text{jet}} > 20$  GeV,  $|\eta^{\text{jet}}| < 2.5$  and  $\Delta R_{\text{jet-jet}} > 0.4$ . The cut efficiencies for different numbers  $n$  of gluons can be read off Table 3. Employing this set of

cuts we were also able to verify the jet production cross sections that we have produced using COMIX with the results reported in Ref. [6]. To have a stringent comparison, we ran COMIX [6] for pure gluon scatterings as provided within SHERPA [24] version 1.2.3 yielding LO cross sections that take the full color dependence into account. These results are also listed in the table; for the 4-gluon and 5-gluon processes, they can be directly compared to the cross section estimates obtained with the TESS MC on the GPU, since the leading-color approximation already gives the exact result. The agreement is found to be satisfactory.

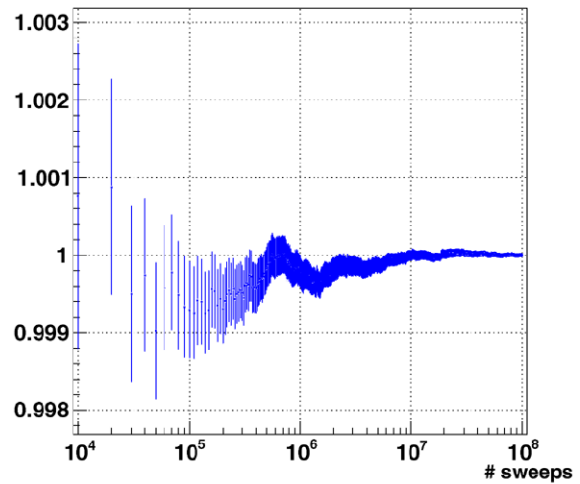
While the GPU code will evolve over the coming years, partly due to new hardware developments, it is of some interest to do a snapshot performance comparison of the GPU code with a highly optimized, state-of-the-art CPU code. Again our choice is to use the parton-level event generator COMIX [6] run in single-thread mode on a quadruple core AMD Phenom™ II X4 940 (3 GHz) processor. We also have utilized the option that in the SHERPA framework the matrix-element generation of COMIX can be combined with a RAMBO-like phase-space integration. As the quantity we base the comparison on, we choose the total computation time needed to reach a certain precision in evaluating the cross sections as listed in Table 3. Our results are shown in Table 4 with the benchmark precision taken from the respective GPU calculation. The values marked by a “\*” were extracted from runs terminated before completion; they are rough estimates of how long the integration would have taken if forced to reach the target precision. Some numbers only represent lower bounds taken from Monte Carlo integrations with uncertainties larger than the benchmarks. From Table 4 we see faster evaluations being accomplished by the CPU codes for the processes lowest in multiplicity. Because of their simplicity, obviously in our case it is impedimental to deal with the extra overhead of steering the GPU calculation from the CPU. For medium multiplicities, the GPU computations are faster by factors of a few or comparable in speed with the COMIX calculations while downgrading from the COMIX-specific to a RAMBO-like phase-space generation of the gluon momenta results in manifestly slower evaluation times for  $n \geq 6$  gluons. For large gluon multiplicities ( $n \geq 11$ ), the values for the CPU run times turn hopeless and one in fact clearly benefits from performing the cross section calculations on the GPU. However, for two reasons, the numbers for the computation time ratios need to be interpreted carefully: (i) CPU codes can be designed to make use of CPU multi-threading and (ii) for the same reason as already discussed above, the problematic convergence behavior of the RAMBO-like phase-space generation leads to larger than ordinarily expected uncertainties on the ratios of Table 4.

We show differential distributions in Fig. 4. To obtain them we again used  $10^9$  sweeps where, for a certain gluon

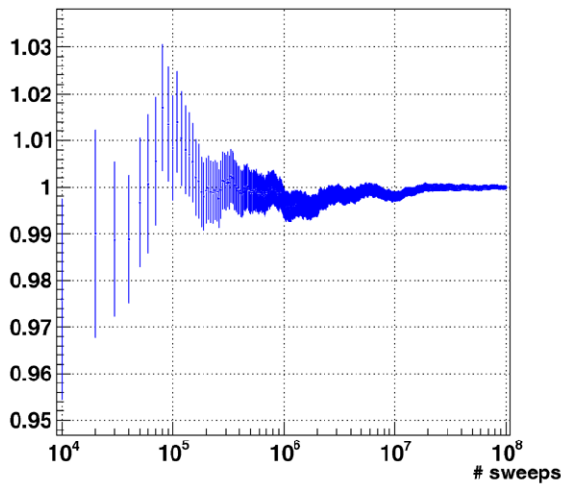
The normalized 4g cross section without Kahan summation



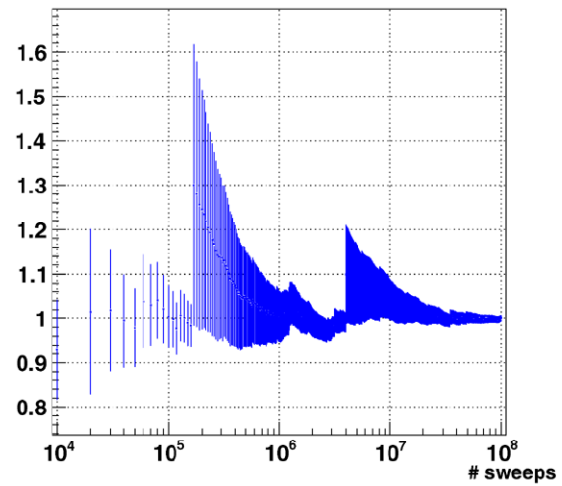
The normalized 4g cross section



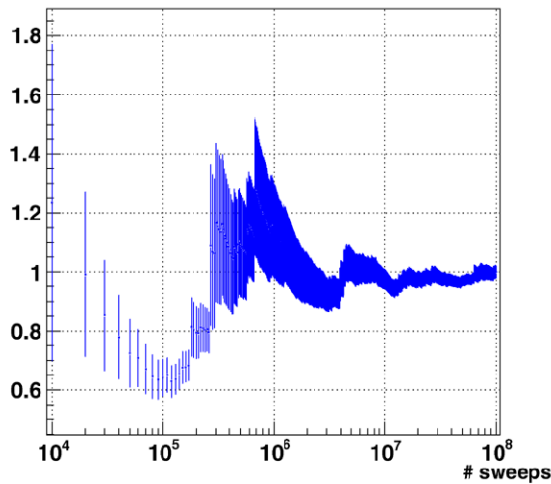
The normalized 6g cross section



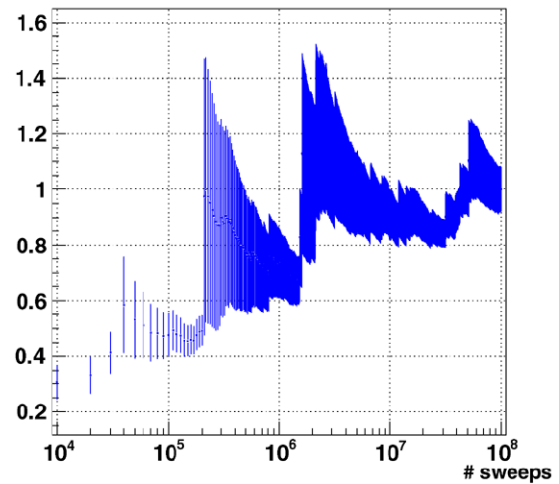
The normalized 8g cross section



The normalized 10g cross section



The normalized 12g cross section



**Fig. 3** The number of sweeps versus several  $gg \rightarrow (n - 2)g$  cross sections normalized to their respective best cross section estimates as given in Table 3. The error is the mean standard deviation. The plot in

the *top left pane* is an example of false cross section determination, if one does not rely on Kahan summation



**Table 3** The cross sections  $\sigma_n$  for  $gg \rightarrow (n - 2)g$  and their mean standard deviations in pb as calculated by the TESS MC using  $10^9$  sweeps. The two columns to the right show the total number of generated events and the number of events passing the jet cuts. For comparison, the cross sections  $\sigma_n^{\text{COMIX}}$  in pb that were computed by COMIX [6] as implemented in SHERPA [24] version 1.2.3 are also given. Note that the tree-level matrix elements generated by COMIX encode the full color dependence

$n$	$\sigma_n$ (pb) [ $\sigma_n^{\text{COMIX}}$ (pb)]	$N_{\text{generated}}/10^{10}$	$N_{\text{accepted}}/10^{10}$
4	$(2.32421 \pm 0.00047) \times 10^8$	30.6	19.6848
4	$[(2.32584 \pm 0.00047) \times 10^8]$		
5	$(1.4353 \pm 0.0011) \times 10^7$	20.4	11.2939
5	$[(1.4348 \pm 0.0011) \times 10^7]$		
6	$(2.84780 \pm 0.00096) \times 10^6$	14.4	6.98918
6	$[(2.85714 \pm 0.00097) \times 10^6]$		
7	$(6.356 \pm 0.012) \times 10^5$	10.8	4.49985
7	$[(6.422 \pm 0.012) \times 10^5]$		
8	$(1.608 \pm 0.011) \times 10^5$	8.40	2.93316
8	$[(1.670 \pm 0.011) \times 10^5]$		
9	$(4.38 \pm 0.11) \times 10^4$	6.60	1.88182
9	$[(4.97 \pm 0.13) \times 10^4]$		
10	$(1.193 \pm 0.024) \times 10^4$	5.40	1.22356
10	$[(1.489 \pm 0.027) \times 10^4]$		
11	$(3.550 \pm 0.020) \times 10^3$	4.50	0.788017
11	$[(4.80 \pm 0.13) \times 10^3]$		
12	$(9.64 \pm 0.74) \times 10^2$	3.90	0.513041
12	$[(17.7 \pm 3.1) \times 10^2]$		

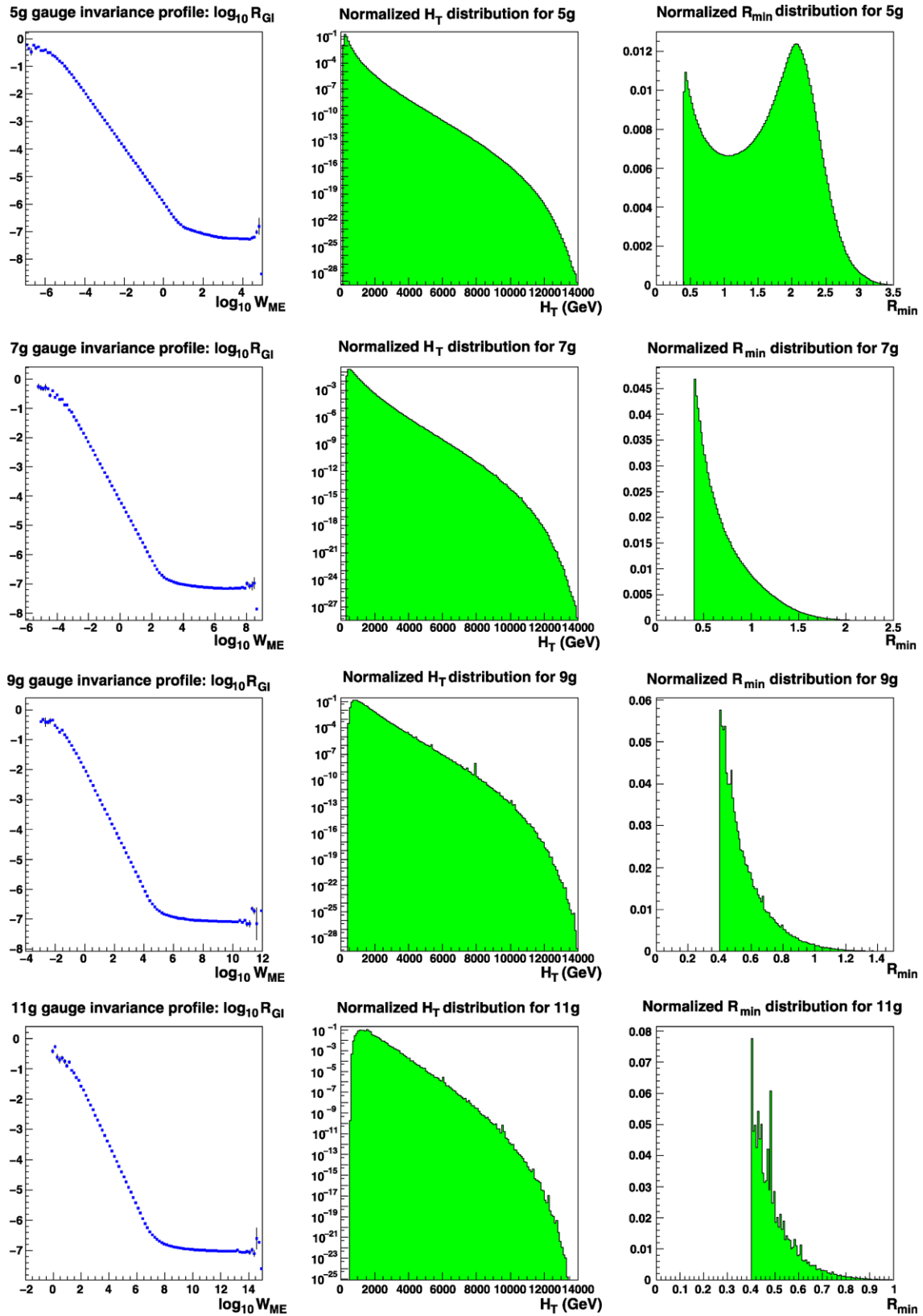
**Table 4** The total computation times  $\tau_n$  in minutes needed to obtain the LO  $gg \rightarrow (n - 2)g$  cross sections with accuracies as given in Table 3 when using the TESS and COMIX [6] Monte Carlo programs. Note that the integration (CPU) times marked by \* are taken from projections in runs terminated before the target precision was reached. Only lower bounds determined with  $\sim 3.1$ ,  $\sim 4.8$  and  $\sim 2.3$  times target precision are shown for the  $n = 10$ ,  $n = 11$  and  $n = 12$  cases, respectively. The SHERPA [24] event generator version 1.2.3 was deployed to run COMIX and get the LO cross sections, which include the full color information. The times given in square brackets refer to using COMIX matrix elements together with a RAMBO-like phase-space integration as opposed to the default treatment in COMIX where the phase-space integration has been optimized to be conform to the matrix-element generation. The rightmost column shows the ratios with respect to the computation time of the TESS Monte Carlo. COMIX was run on a single AMD Phenom™ II X4 940 (3 GHz) processor

$n$	$\tau_n^{\text{TESS}}$ (min.)	$\tau_{n, [\text{RAMBO}]}^{\text{COMIX}}$ (min.)	$\tau_{n, [\text{RAMBO}]}^{\text{COMIX}}/\tau_n^{\text{TESS}}$
4	151.7	40.83 [103.5]	0.27 [0.68]
5	150.9	29.87 [113.1]	0.20 [0.75]
6	205.2	1540 [13500*]	7.50 [66*]
7	414.7	434.8 [11500*]	1.05 [28*]
8	496.4	278.3 [2580*]	0.56 [5.2*]
9	470.1	640.4 [8711]	1.36 [18.5]
10	613.5	7538 [>62950]	12.3 [>102]
11	731.3	>61850	>84
12	881.4	>52350	>59

multiplicity, the total number of generated events can be read off Table 3. We kept most of the input parameters unaltered except for the jet cuts, which we changed to  $p_T^{\text{jet}} > 60$  GeV,  $|\eta^{\text{jet}}| < 2.0$  and  $\Delta R_{\text{jet-jet}} > 0.4$ , and the choice of the renormalization and factorization scales, which we decided to set dynamically using  $H_T$  as a scale. On the right hand side of Fig. 4 we show for 3, 5, 7, 9 gluon jets in the final state the normalized distributions for the  $H_T$  observable and the minimum  $R$ -separation,  $R_{\text{min}}$ , which we define through the jet-jet pair being closest in  $R$ -space,  $R_{\text{min}} = \min\{\Delta R_{ij}\}$ . As can be seen smooth distributions are easily obtained using the RAMBO phase-space generator. They are normalized to the total cross sections, which have been calculated by TESS as

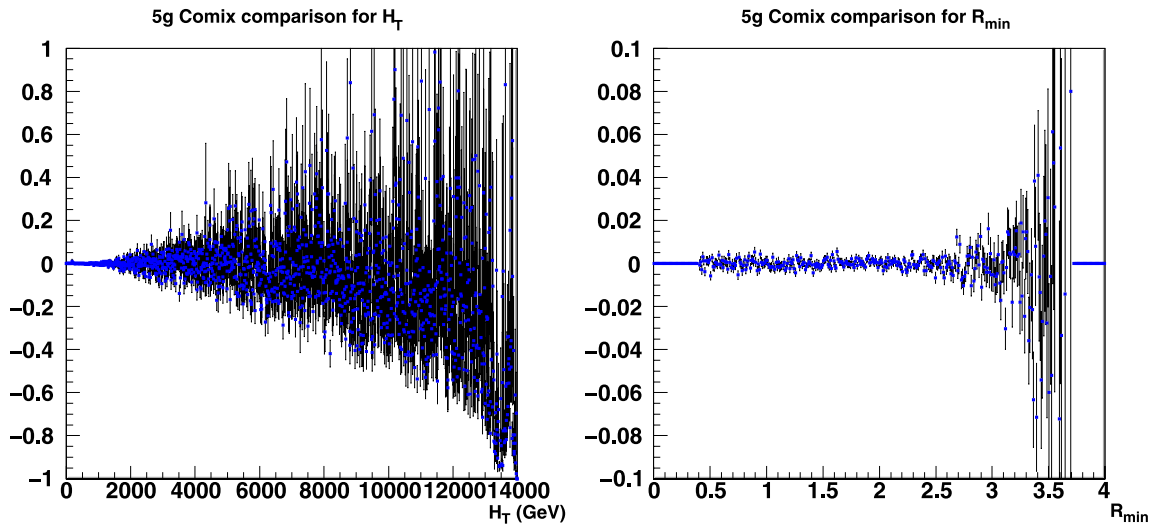
$$\begin{aligned} \sigma_5 &= (6.97838 \pm 0.00044) \times 10^4 \text{ pb}, \\ \sigma_7 &= (4.9761 \pm 0.0043) \times 10^2 \text{ pb}, \\ \sigma_9 &= (4.532 \pm 0.044) \text{ pb}, \\ \sigma_{11} &= (4.51 \pm 0.19) \times 10^{-2} \text{ pb}. \end{aligned} \tag{9}$$

On the left hand side of Fig. 4 we have added profile plots displaying the relative gauge invariance versus the decimal logarithm of the matrix-element weight. Specifically, we



**Fig. 4** Left panels: the profile plots of the relative gauge invariance as a function of the decimal logarithm of the matrix-element weight,  $\log_{10} W_{ME}$ ; center panels: the normalized  $H_T$  distributions and right

panels: the normalized minimum  $R$ -separations between pairs of jets. All of which is shown for  $gg \rightarrow (n - 2)g$  scatterings at a 14 TeV center-of-mass energy for  $n = 5, 7, 9, 11$



**Fig. 5** The ratio  $(\sigma_5^{\text{TESS}} \times d\sigma_5^{\text{COMIX}}/dX)/(\sigma_5^{\text{COMIX}} \times d\sigma_5^{\text{TESS}}/dX) - 1$  for the 5-gluon  $X = H_T$  (left panel) and  $X = R_{\min}$  distributions. The  $R_{\min}$  variable is defined by the smallest separation in  $R$ -space of any

pair of jets. The mean standard deviation error bars of the COMIX calculation are also shown

show the average

$$\frac{|K[1] \cdot J[2, \dots, n]|^2}{|J[1] \cdot J[2, \dots, n]|^2}$$

and its mean standard deviation as a function of the matrix-element weight

$$W_{\text{ME}} = |m(1, 2, \dots, n)|^2 = |(J[1] \cdot J[2, \dots, n]) \times K^2[2, \dots, n]|^2. \tag{10}$$

The behavior is as expected; for large weights, we see gauge cancellations up to float precision. For small weights, the gauge cancellations are less precise. However, these small-weight events are not important since they do not contribute to the calculation of the observables.

Finally, in Fig. 5 we compare our results to the results obtained with the SHERPA event generator [24] version 1.2.0 where the tree-level matrix elements and phase-space integrators again have been generated by COMIX [6]. For this comparison, we use both the  $H_T$  and  $R_{\min}$  5-gluon distributions and we fix the renormalization and factorization scale through  $M_Z = 91.118$  GeV to avoid any issues resulting from slight differences in the evolution codes for running scales between the two MCs. Furthermore, to have a sole shape comparison, we plot the ratio  $(\sigma_5^{\text{TESS}} \times d\sigma_5^{\text{COMIX}}/dX)/(\sigma_5^{\text{COMIX}} \times d\sigma_5^{\text{TESS}}/dX) - 1$  with the results shown in Fig. 5 and  $X$  being the observable in consideration. Note that for the minimum  $R$ -separation distribution, we find excellent agreement with the SHERPA prediction given by COMIX. For the  $H_T$  distribution, we have to realize that the cross section spans 28 orders of magnitude. As

COMIX relies on importance sampling, it only sparsely populates the tail of the distribution. This leads to large uncertainties at large values of  $H_T$  and, in these regions, COMIX will hence tend to underestimate the value for the cross section.

### 5 Conclusions and outlook

In our first exploration of the potential of using multi-threaded GPU-based workstations for Monte Carlo programs, we obtained very encouraging results. We implemented the entire TESS Monte Carlo on the GPU chip; the only off-chip usage occurs through utilizing the texture memory for the evaluation of the parton density function and the strong coupling constant. The GPU global memory is solely used for transferring the Monte Carlo results to the CPU memory. At this exploratory phase of the project, we limited ourselves to the calculation of leading-color leading-order  $n$ -gluon matrix elements. With respect to the CPU-based implementation of our Monte Carlo we have found impressive speed-ups in the computations reaching from  $\mathcal{O}(300)$  for  $PP \rightarrow 2$  jets to  $\mathcal{O}(150)$  for  $PP \rightarrow 10$  jets. In a comparison with a (CPU-based) modern parton-level generator such as COMIX [6] we notice that the calculation of cross sections can become more efficient for more complex processes like  $gg \rightarrow 4g$  and larger multiplicities in the final state. The full potential of the GPU calculation unfolds itself when one deals with the determination of cross sections for  $gg \rightarrow \geq 9g$  cross sections; here the benefits with respect to the CPU-based evaluations are huge. Nevertheless, the results, in particular, the uncertainties of the phase-space

integration in TESS should be interpreted with some care, since the employed RAMBO algorithm is very sensitive to the occurrence of suddenly largely fluctuating weights.

Given these results we are encouraged to further develop the TESS Monte Carlo by including quarks, vector bosons and subleading color contributions. We are also planning to implement on the GPU a dipole-based phase-space generator like SARGE as an alternative to the unit-weight phase-space generator RAMBO. This will avoid the under-sampling issues in high jet-multiplicity final states and render a performance comparison to efficient CPU-based tree-level Monte Carlos more meaningful. These improvements will result in a full leading-order parton-level event generator, which has the potential to be two orders of magnitude faster than existing leading-order parton-level generators.

More importantly, a GPU-based Monte Carlo may be used as the generator for the real corrections in an automated next-to-leading order parton-level MC generator. The virtual corrections can be calculated by using a generalized-unity based method [25–33].

Finally, GPU chips for numerical evaluations are still evolving rapidly. This will lead to additional significant speed-ups over CPU-based Monte Carlos in the coming years. The next generation GPUs are already on the market: the NVIDIA<sup>®</sup> Fermi<sup>™</sup> chip (released in Fall 2010) comes with improvements over the Tesla chip such as 32 kb instead of 16 kb memory for registers, 48 kb instead of 16 kb shared memory and support for double-precision calculations. Through a unified memory pointer the Fermi chip has full support for the C++ programming language. Running TESS as is in the new environment will lead to a performance increase varying between 1.3 and 2.0 owing to the faster clock rate and more special-function units. In addition and with some effort, the memory layout of the program could be adjusted to make use of the increase in shared memory and hence run more events in parallel. This will give significant additional speed-ups.

**Acknowledgements** We want to thank Jim Simone for suggesting the Tesla GPU for use in event generators. We thank the High-Performance Computing Department at Fermilab for giving us support and access to the LQCD Tesla-based workstations.

We would also like to thank Patrick Fox, who came up with the name “TESS” and suggested it to us as the name for our Monte Carlo program.

Fermilab is operated by Fermi Research Alliance, LLC, under contract DE-AC02-07CH11359 with the United States Department of Energy.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution Noncommercial License which permits any noncommercial use, distribution, and reproduction in any medium, provided the original author(s) and source are credited.

## Appendix: The GPU architecture

Here we give a rough overview of the CUDA framework. We, however, cannot discuss all of the many more subtleties, which can be detrimental to the performance, if not aware of.

The CUDA framework is build to utilize the NVIDIA<sup>®</sup> Tesla<sup>™</sup> GPUs. To program the GPUs efficiently, one needs to have a global understanding of their underlying structure. The programming model is build around three concepts: threads, blocks and grids. Threads are exactly what one expects—a single thread of code execution. These threads are the basis of the parallelization with the same thread being executed multiple times working on different data, but simultaneously. Identical threads are organized in blocks where each thread in a block has a unique thread index to specify what data to work on. Identical blocks are organized in a grid where each block has a unique block index. The reason for this extra hierarchy is that threads within one block are always executed on a single multi-processor (MP). Therefore, threads within a block can share information using the shared memory and can synchronize their execution. Threads in different blocks may be executed on different MPs and are therefore unable to share information or synchronize execution. A single MP is able to execute more than one block if the total number of threads, the total register count and the total amount of shared memory are within the limits of the MP.

One MP consists of 8 streaming processors, or “CUDA cores”, and 2 special-function units (SFUs).<sup>6</sup> It is possible to run up to 1024 threads on one MP. In a single instruction multiple data (SIMD) design, a CUDA core can execute 32 threads simultaneously. The 32 threads executing on a single CUDA core are called a warp. Within a warp every instruction needs to be the same. Branchings, which diverge within a warp are therefore not executed concurrently, but sequentially. So, if part of the threads in a warp branch and another part does not, then the chip will first calculate one of the two code paths and when finished the other. Diverging branchings within warps are therefore expensive as opposed to those in threads belonging to different warps that are free. This is important to keep in mind while designing code. Note that each MP has only 2 SFUs; calculating special functions like logarithms can easily result in a bottleneck, because the threads have to share the SFUs.

The other important piece of the architecture is the memory structure. The GPU has its own memory which essentially is large and of the order of gigabytes. This memory

<sup>6</sup>The total number of streaming processors is  $8 \times \#(\text{MP})$ . The number of MPs can vary between cards, but often, as well as in our case, there are 30 MPs on a card—10 TPCs times 3 MPs per TPC (Thread Processing Clusters). From a programmer’s point of view, it is sufficient to know how a single core performs considering the fact that 8 cores have to share 2 SFUs which may result in performance slow-downs.

cannot be accessed directly by the CPU, but must be copied to and from CPU memory using the CUDA API functions. This is the only way to exchange information between CPU and GPU. This copying can easily slow down the program, hence, it is important to minimize copying in designing the program. Furthermore, the bandwidth between the GPU memory and the GPU chip is limited; thus, if all threads ask for global memory, a bottleneck is created easily. The second type of memory is provided by the (32 bit) registers; these are local to the thread and there is a total of 16 Kb registers available to be divided among the threads. These registers are fast, so one wants to store as much information as possible in registers. Design trade-offs are to be made here: more registers per thread means less total threads per MP. The third type of memory is shared memory. This memory with total size of 16 Kb is much like registers, it is fast, but shared among the threads in the same block. Using shared memory of course comes with the usual synchronization problems when different threads access the same element.

## References

1. T. Stelzer, W.F. Long, *Comput. Phys. Commun.* **81**, 357 (1994). [arXiv:hep-ph/9401258](https://arxiv.org/abs/hep-ph/9401258)
2. F. Krauss, R. Kuhn, G. Soff, *J. High Energy Phys.* **0202**, 044 (2002). [arXiv:hep-ph/0109036](https://arxiv.org/abs/hep-ph/0109036)
3. P.D. Draggotis, R.H.P. Kleiss, C.G. Papadopoulos, *Eur. Phys. J. C* **24**, 447 (2002). [arXiv:hep-ph/0202201](https://arxiv.org/abs/hep-ph/0202201)
4. M.L. Mangano, M. Moretti, F. Piccinini, R. Pittau, A.D. Polosa, *J. High Energy Phys.* **0307**, 001 (2003). [arXiv:hep-ph/0206293](https://arxiv.org/abs/hep-ph/0206293)
5. E. Boos et al. (CompHEP Collaboration), *Nucl. Instrum. Methods A* **534**, 250 (2004). [arXiv:hep-ph/0403113](https://arxiv.org/abs/hep-ph/0403113)
6. T. Gleisberg, S. Höche, *J. High Energy Phys.* **0812**, 039 (2008). [arXiv:0808.3674](https://arxiv.org/abs/0808.3674) [hep-ph]
7. W.T. Giele, Z. Kunszt, K. Melnikov, *J. High Energy Phys.* **0804**, 049 (2008). [arXiv:0801.2237](https://arxiv.org/abs/0801.2237) [hep-ph]
8. C.F. Berger et al., *Phys. Rev. D* **78**, 036003 (2008). [arXiv:0803.4180](https://arxiv.org/abs/0803.4180) [hep-ph]
9. A. van Hameren, C.G. Papadopoulos, R. Pittau, *J. High Energy Phys.* **0909**, 106 (2009). [arXiv:0903.4665](https://arxiv.org/abs/0903.4665) [hep-ph]
10. V. Hirschi, R. Frederix, S. Frixione, M.V. Garzelli, F. Maltoni, R. Pittau, [arXiv:1103.0621](https://arxiv.org/abs/1103.0621) [hep-ph]
11. See the websites <http://www.nvidia.com/object/personal-supercomputing.html> and [http://www.nvidia.com/object/product-tesla\\_c1060\\_us.html](http://www.nvidia.com/object/product-tesla_c1060_us.html) for further information
12. See the website [http://www.nvidia.com/object/cuda\\_home.html](http://www.nvidia.com/object/cuda_home.html) for further information
13. R. Kleiss, W.J. Stirling, S.D. Ellis, *Comput. Phys. Commun.* **40**, 359 (1986)
14. M.R. Whalley, D. Bourilkov, R.C. Group, [arXiv:hep-ph/0508110](https://arxiv.org/abs/hep-ph/0508110)
15. K. Hagiwara, J. Kanzaki, N. Okamura, D. Rainwater, T. Stelzer, *Eur. Phys. J. C* **66**, 477–492 (2010). [arXiv:0908.4403](https://arxiv.org/abs/0908.4403) [physics.comp-ph]
16. K. Hagiwara, J. Kanzaki, N. Okamura, D. Rainwater, T. Stelzer, *Eur. Phys. J. C* **70**, 513–524 (2010). [arXiv:0909.5257](https://arxiv.org/abs/0909.5257) [hep-ph]
17. H. Murayama, I. Watanabe, K. Hagiwara, HELAS: HELicity Amplitude Subroutines for Feynman diagram evaluations, KEK-91-11
18. F.A. Berends, W.T. Giele, *Nucl. Phys. B* **306**, 759 (1988)
19. R. Kleiss, H. Kuijff, *Nucl. Phys. B* **312**, 616 (1989)
20. P. Draggotis, R.H.P. Kleiss, C.G. Papadopoulos, *Phys. Lett. B* **439**, 157 (1998). [arXiv:hep-ph/9807207](https://arxiv.org/abs/hep-ph/9807207)
21. W. Kahan, Further remarks on reducing truncation errors. *Commun. ACM* **8**(1), 40 (1965)
22. P.D. Draggotis, A. van Hameren, R. Kleiss, *Phys. Lett. B* **483**, 124 (2000). [arXiv:hep-ph/0004047](https://arxiv.org/abs/hep-ph/0004047)
23. J. Pumplin, D.R. Stump, J. Huston, H.L. Lai, P.M. Nadolsky, W.K. Tung, *J. High Energy Phys.* **0207**, 012 (2002). [arXiv:hep-ph/0201195](https://arxiv.org/abs/hep-ph/0201195)
24. T. Gleisberg, S. Höche, F. Krauss, M. Schönherr, S. Schumann, F. Siegert, J. Winter, *J. High Energy Phys.* **0902**, 007 (2009). [arXiv:0811.4622](https://arxiv.org/abs/0811.4622) [hep-ph]
25. R.K. Ellis, K. Melnikov, G. Zanderighi, *J. High Energy Phys.* **0904**, 077 (2009). [arXiv:0901.4101](https://arxiv.org/abs/0901.4101) [hep-ph]
26. C.F. Berger, Z. Bern, L.J. Dixon, F. Febres Cordero, D. Forde, T. Gleisberg, H. Ita, D.A. Kosower et al., *Phys. Rev. Lett.* **102**, 222001 (2009). [arXiv:0902.2760](https://arxiv.org/abs/0902.2760) [hep-ph]
27. K. Melnikov, M. Schulze, *J. High Energy Phys.* **0908**, 049 (2009). [arXiv:0907.3090](https://arxiv.org/abs/0907.3090) [hep-ph]
28. G. Bevilacqua, M. Czakon, C.G. Papadopoulos, R. Pittau, M. Worek, *J. High Energy Phys.* **0909**, 109 (2009). [arXiv:0907.4723](https://arxiv.org/abs/0907.4723) [hep-ph]
29. G. Bevilacqua, M. Czakon, C.G. Papadopoulos, M. Worek, *Phys. Rev. Lett.* **104**, 162002 (2010). [arXiv:1002.4009](https://arxiv.org/abs/1002.4009) [hep-ph]
30. C.F. Berger, Z. Bern, L.J. Dixon, F.F. Cordero, D. Forde, T. Gleisberg, H. Ita, D.A. Kosower et al., *Phys. Rev. D* **82**, 074002 (2010). [arXiv:1004.1659](https://arxiv.org/abs/1004.1659) [hep-ph]
31. K. Melnikov, M. Schulze, *Nucl. Phys. B* **840**, 129–159 (2010). [arXiv:1004.3284](https://arxiv.org/abs/1004.3284) [hep-ph]
32. R. Frederix, S. Frixione, K. Melnikov, G. Zanderighi, *J. High Energy Phys.* **1011**, 050 (2010). [arXiv:1008.5313](https://arxiv.org/abs/1008.5313) [hep-ph]
33. C.F. Berger, Z. Bern, L.J. Dixon, F.F. Cordero, D. Forde, T. Gleisberg, H. Ita, D.A. Kosower et al., *Phys. Rev. Lett.* **106**, 092001 (2011). [arXiv:1009.2338](https://arxiv.org/abs/1009.2338) [hep-ph]