

Simple, Parallel, High-Performance Virtual Machines for Extreme Computations

Bijan Chokoufe Nejad^{a,b}, Thorsten Ohl^b, Jürgen Reuter^a

^aDESY Theory Group, Notkestr. 85, D-22607 Hamburg

^bUniversity of Würzburg, Emil-Hilb-Weg 22, 97074 Würzburg, Germany

Abstract

We introduce a high-performance virtual machine (VM) written in a numerically fast language like Fortran or C to evaluate very large expressions. We discuss the general concept of how to perform computations in terms of a VM and present specifically a VM that is able to compute tree-level cross sections for any number of external legs, given the corresponding byte code from the optimal matrix element generator, O'MEGA. Furthermore, this approach allows to formulate the parallel computation of a single phase space point in a simple and obvious way. We analyze hereby the scaling behaviour with multiple threads as well as the benefits and drawbacks that are introduced with this method. Our implementation of a VM can run faster than the corresponding native, compiled code for certain processes and compilers, especially for very high multiplicities, and has in general runtimes in the same order of magnitude. By avoiding the tedious compile and link steps, which may fail for source code files of gigabyte sizes, new processes or complex higher order corrections that are currently out of reach could be evaluated with a VM given enough computing power.

Keywords: Virtual Machine, High-Performance Computing, Automation of perturbative calculations, Higher Orders, Parallel Computation

1. Introduction

Computations in high energy physics tend to hit the limits of what is computationally feasible. Setting demanding grid computations aside, one encounters even in perturbative calculations expressions of cross sections of enormous size. Such computations for the Large Hadron Collider (LHC), its upgrade the High Luminosity LHC (HL-LHC) or the planned International Linear Collider (ILC) are and keep getting more challenging as cross sections are needed for a high number of external particles and to increasing precision to match the experimental efforts. When facing such problems, a compromise has to be made, in order to have a maintainable and extendible solution for the developer and at the same time fast execution of the code. The latter

Email addresses: bijan.chokoufe@desy.de (Bijan Chokoufe Nejad), ohl@physik.uni-wuerzburg.de (Thorsten Ohl), juergen.reuter@desy.de (Jürgen Reuter)

Preprint submitted to Computer Physics Communications

November 17, 2014

cannot be overrated as the same code, typically representing a certain process, has to be evaluated billions of times with different input data for the Monte Carlo integration and/or parameter scans.

A popular approach to solve this problem is a meta-programming ansatz, i.e. to determine the expression of a cross section itself in a higher level programming language like `Mathematica`, `OCaml`, `FORM` or `Python` while the numerical evaluation is performed in high-performance languages like `Fortran` or `C`. Hereby, the expression is vastly reduced with computer algebra and tailored algorithms on the higher level to make the execution later on as fast as possible. Examples for this are the tree-level and one-loop matrix element generators `MADGRAPH` [1], `FORMCALC` [2, 3] or `O'MEGA` [4]. A problem, however, arises when the expression becomes so large that it is impossible to compile and link, and hence to evaluate numerically, due to the sheer size. In `Fortran`, which is known for its excellent numerical performance, we typically encounter this problem for source code of gigabyte sizes irrespective of the available memory. This problem is also being addressed by the project `HEPGAME` [5] that is based on `Form` [6, 7] and aims to reduce the code size before compilation by using new concepts from game theory like Monte Carlo tree searches. Furthermore, we should mention `HAGGIES` [8], written in `Java`, which also generates optimised programs for efficient numerical evaluation of mathematical expressions using multivariate Horner-schemes and common subexpression elimination (CSE) to reduce the source code size.

In this paper, we show how to circumvent the tedious compile step in between completely by using a VM. To avoid confusions, we have to define what we mean with the term VM. A VM is in our context a compiled program, an interpreter, that is able to read instructions, in the form of *byte code*, from disk and perform an arbitrary number of operations out of a finite *instruction set*. We do not refer to any sort of operating system emulation that is commonly encountered under the term VM. Also the parallel virtual machine (PVM) [9] is a completely different idea, combining a network of multiple computers to one VM. Far closer to our VM is the VM used in the open-source project `NUMEXPR` [10]. Their VM is written in `C` and specializes on the fast numerical expression evaluation of very large arrays in `Python` by dividing array operands in chunks that easily fit in the cache of the central processing unit (CPU) and avoiding the creation of temporary arrays. Though the idea is related, in our application we have comparably small arrays per instruction and can hence not benefit from this project. We want to stress that a VM allows the complexity of the computation to be only set by the available hardware and not limited by software design or intermediate steps. Furthermore, we will show that a VM is easy to implement and makes parallel evaluation obvious.

An important concern is of course whether the VM can still compete with compiled code in terms of speed. The instructions have to be translated by the VM to actual machine code, which is a potential overhead. However, in the computation of matrix elements there are typically a lot of complex scalar products involved in a single instruction implying that this overhead plays a minor role. What we explicitly give up are the optimizations that the compiler can perform in the context of multiple instructions like CSE and data and instruction prefetching. Of these, at least the CSE can be done beforehand by constructing the byte code with the lowest number of common subexpressions on the higher level. In fact, we will show that a VM can even be faster than compiled code for certain processes and compilers since the VM also benefits from the fact that instruction cache misses are less likely in this formulation.

We will apply the concept of a VM to the tree-level Optimizing Matrix Element Generator, `O'MEGA`, to allow the computation of higher multiplicities of colored particles given the same hardware. This does obviously not imply that the presented computational method is restricted

to tree-level computations. When trying to obtain higher order cross-sections the same problem can arise even for less external particles, due to the inherent complexity of the computation. We expect that VM implementations in such environments are a possible way to go beyond what is nowadays considered as still feasible.

Apart from cross sections, we believe that the problem of evaluating huge expressions numerically is a more general one, just as algebraic tools like FORM [6, 7] or integration tools like CUBA [11] are useful beyond their original field of study. Therefore, we will tackle this problem at first in a rather general way in Section 2, before we turn to the implementation of the O’Mega virtual machine (OVM) in Section 3. Then we benchmark this proposal in Section 4 and conclude with a summary of our findings and a small technical outlook in Section 5.

2. General Virtual Machines

We will describe in this section the necessary components to perform a computation with a VM. The byte code plays a central role as it embodies all nontrivial information about how to compute the object of desire. One might imagine the VM as a machine, which has a number of registers, and is given instructions how to act on them. This picture is quite similar to a CPU, except that we are doing this on a higher level, i.e. our registers are arrays of e.g. wave functions or momenta and the instructions can encode scalar products or more complicated expressions. In Appendix A, we also give a purely mathematical example implementation of a VM. This code is well suited for adaption to other problems as it has no dependency on external libraries and still includes all of the necessary infrastructure.

2.1. Byte Code

For the dynamic construction of the VM, it is necessary to include a *header* in the byte code, which contains the number of objects that have to be allocated. For convenience, it is also useful to have some version numbers that document which physical or mathematical constants should be used together with this byte code or comments to indicate how it was produced. Optionally, one can add after the header tables of precomputed parameters, like information about the involved helicities, color or flavor. After this the body of instructions follows, whereby each line corresponds to a certain operation that the VM should perform on its registers.

We encode the instructions in pure integers inside a simple ASCII byte code such that it is in principle human-readable if the meanings of the numbers are known. Though the use of mere numbers does not exploit the full dictionary of ASCII, it allows a very fast integer line-by-line scan of the byte code in the Fortran VM and avoids a translation step because integers are already addresses in arrays. Since the initialization is, however, very fast compared to the runtime, this could be optimized with a binary format, if the size of the byte code becomes a problem. As the byte code size is about a factor of ten smaller compared to the native source code, as shown in Section 4.3 this is not yet a concern for our application. The fact that our byte code is portable and platform independent is a positive surplus when calculations are performed on clusters.

The first number of an instruction is the *operation code (opcode)* that specifies which operation will be performed. For illustration, consider the example

```
1 5 4 3
```

which could be translated into $\text{momentum}(5) = \text{momentum}(4) + \text{momentum}(3)$, a typical operation to compute the s -channel momentum in a $2 \rightarrow 2$ scattering process. Depending on the context, set by the opcode, the following numbers have different meanings but are typically addresses, i.e. indices, of objects, or specify how exactly the function should act on the operands, by what numbers the result should be multiplied, etc.

When designing the byte code of a high-performance VM, the line length should be chosen such that the most frequent operations fit within a line. Complex operations that would increase the line length significantly above the average requirement, can also be split in multiple lines by using *sub-instructions*, which are introduced in Section 2.3.

2.2. Interpreter

The interpreter is a very simple program that reads the byte code into memory and then loops over the instruction block with a `decode` function, which is basically a `select/case` statement depending on the opcode. The instructions can be instantly, compared to the execution time of the relevant instructions, translated to physical machine code, since the different types of operations are already compiled and only the memory locations of the objects have to be inserted. The compilation of the VM itself is very fast and has only to be done once which is handy for the use of many byte codes and necessary for extreme computations as motivated above.

Two things have to be adapted in the interpreter of the VM, when one wants to tackle a new type of problem, e.g. when going from tree-level to one-loop. At first, one has to specify, where and with which types to expect header, comments, tables and instructions¹. Furthermore, the `decode` function needs to be able to translate any instruction line into operations on registers, i.e. all opcodes have to be implemented. The functions can be arbitrarily complex and are also allowed to call external libraries, though most likely better performance is achieved by keeping things as simple as possible. Especially, with parallelization in mind, it is desirable to have roughly the same amount of computation time for different instructions, to ensure an even workload and hereby minimizing idle times at synchronization points.

Given this environment, the byte code file that is given to the VM completely dictates the specific problem, or process in the cross section context, that should be computed. Input data or external parameters are given as arguments to the function call of the VM. The calling application has of course to make sure that these parameters match the corresponding byte code, which can be ensured with version numbers.

2.3. Parallelization

The generation of events for collider physics usually parallelize trivially. Since an integral is in most cases needed, the same code is just evaluated multiple times with different input data. The situation can change, however, for an extreme computation that already uses all caches. Depending on the size of the caches and the scheduler, evaluating such code with multiple data at the same time, can run even slower than the single-threaded execution. Obviously, the computation is then so large, containing numerous objects, that it is worth trying to parallelize the execution with a single set of input data with shared memory.

Developing truly parallel code for a complicated calculation, however, is a non-trivial task since race conditions have to be kept in mind at all times. Furthermore, physicists have to delve into the frameworks like OpenMP or MPI to find the best parallelization method for each piece,

¹One could, in principle, also determine this dynamically by using a certain markup, if one has the desire to do so.

which is time consuming and likely to introduce bugs that are hard to find. The byte code gives us the opportunity to write the parallel computation in an obvious fashion that is both easy to generate and to implement in the VM. The idea is to split the byte code into recursion *levels*, whereby in each level all *building blocks* are non-nested and may be computed in parallel. Different levels are separated by necessary synchronization points at which threads have to wait until intermediate results are communicated and which can be represented in the byte code with a zero opcode. It is clear that one should aim to keep the number of synchronization points to the inherent minimum of the computation for optimal performance.

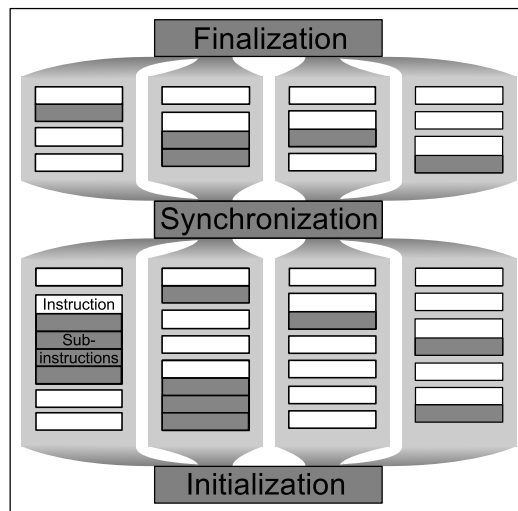


Figure 1: Sketch of the parallelization scheme for byte code of two levels. Instructions and sub-instructions are in white and gray, respectively. Certain instructions imply that all following sub-instructions have to be executed before the next instruction is addressed. This grouping of instructions allows multiple sequential writes while minimizing synchronization points.

The fact that we demand commutativity within a level implies that every virtual register is changed by at most one thread. A potential problem would hence be that the same address might be written to successively multiple times in a computation though still being fully disconnected to other parts. To maintain the parallel nature with respect to the other parts and at the same time the sequential nature of such a subcomputation, we can group instructions addressing the same register to a building block. A building block consists of one instruction and zero or more sub-instructions. Sub-instructions are conveniently represented in the byte code with negative opcodes that are skipped over by the main loop. Normal instructions can imply that all following sub-instructions have to be executed sequentially before the thread computes the next instruction. This is sketched in Figure 1. In OpenMP this translates to a `parallel` region with a `parallel do` loop over the instructions in a level. More specifically, the complete parallelization of the VM can be written as

```
!$omp parallel
do level = 1, vm%N_levels - 1
  !$omp do schedule (static)
  do instruction = vm%levels (level) + 1, vm%levels (level + 1)
```

```

        call decode (vm, instruction)
    end do
    !$omp end do
end do
!$omp end parallel

```

hereby we have an outer loop over levels as `vm%levels` contains the indices in the vector of instructions at which the level changes. Due to our organization of the byte code, we can then perform the inner loop in each level in parallel, whereby the `static` schedule just means that every thread gets the same number of instructions. The `decode` function is given the line number instruction in the `instructions` block that should be translated.

As a side note, we want to mention that the sketched parallelization should be very well suited for an implementation on a graphics processing unit (GPU). A common problem, encountered when trying to do scientific computing on a GPU, is the finite kernel size problem. As noted e.g. in Ref. [12], large source code cannot be processed by the CUDA compiler, which is related to the fact that the numerous cores on a GPU are designed to execute simple operations very fast. Dividing an amplitude into smaller pieces, which are computed one by one, introduces more communication overhead and is no ultimate solution since the compilation can still fail for complex amplitudes [12]. The VM on the other hand is a fixed small kernel, no matter how complex the specific computation is. A potential bottleneck might be the availability of the instruction block to all threads, but this question has to be settled by an implementation and might have a quite hardware dependent answer.

Finally, we note that the phase space parallelization mentioned in the beginning of this subsection can still be applied. When considering heterogeneous cluster or grid environments, where each node is equipped with multi-core processors, a combination of distributed memory parallelization for the combination of different phase space points and shared memory parallelization of a single point seems to be a quite natural and extremely potent combination.

3. O'MEGA Virtual Machine

The concept of a VM can be easily applied to evaluate tree-level matrix elements of arbitrary multiplicity. The Optimizing Matrix Element Generator, O'MEGA [4], avoids the redundant representation of amplitudes in the form of Feynman-diagrams by using one-particle off-shell wave functions (1POWs) recursively. Just like the first two numerical codes ALPHA [13] and HELAC [14], which focussed on the standard model (SM), O'MEGA tames herewith the computational growth with the number of external particles from a factorial to an exponential one but is completely general with respect to the used Lagrangian. Other programs with a very similar approach are COMIX [15], based on the color-dressed Berends-Giele recursion formulation, and RECOLA [16], which follows more closely the Dyson-Schwinger formulation incorporating an important generalization [17] that allows to compute one-loop amplitudes in the SM.

The model-independence is achieved in O'MEGA with the meta-programming ansatz mentioned earlier whereby the symbolic representation is determined in OCaml. This abstract expression is then translated to valid Fortran code that is automatically compiled and used in WHIZARD [18] for event generation. As O'MEGA has been designed in a modular way, it has been rather straightforward to add an additional output module that produces byte code instead of Fortran code. Some additional technical details about the implementation can be found in Ref. [19] and more completely in the documented source code [20].

Table 1: Byte code cheat sheet. Each instruction line consists of eight numbers having a different meaning depending on the first one, the opcode. In general, the objects on the left hand side (lhs) are constructed from the right hand side (rhs). X, Y and Z are placeholders for the different Lorentz types of wave functions like fermions, scalars, etc. The value for width indicates which width scheme is used while its value and the one of the mass is inferred from the PDG code. outer_ind denotes spin and momentum index of the wave function. sym is the symmetry factor computed from the number of identical particles involved.

code	coupl	coeff	lhs	rhs ₁	rhs ₂	rhs ₃	rhs ₄
ADD_MOMENTA	0	0	p_lhs	p_rhs ₁	p_rhs ₂	p_rhs ₃	0
LOAD_X	PDG	0	wf	outer_ind	0	0	amp
PROPAGATE_Y	PDG	width	wf	p	0	0	amp
FUSE_Z	coupl	coeff	lhs	rhs ₁	rhs ₂	rhs ₃	rhs ₄
CALC_BRAKET	sign	0	amp	sym	0	0	0

The number of distinct operations that have to be performed in the computation of a cross section is related to the Feynman rules and therefore quite limited. As such, these operations are very good candidates for the translation to byte code. In fact, this results in only about 80 different opcodes for the complete SM, which have been implemented in the OVM. In order to support completely general Lagrangians with arbitrary tensor structures as in [21, 22], the subroutines implementing the vertices can be mapped to opcodes dynamically. They can be classified as described in Table 1 as ADD_MOMENTA, LOAD_X, PROPAGATE_Y, FUSE_Z and CALC_BRAKET, i.e. the addition of momenta, the construction of external wave functions, the propagation of wave functions, the fusion of wave functions according to the Feynman rules and the computation of the final bracket, which yields the amplitude with appropriate prefactors, respectively. This limited set of instructions as well as the objects in a calculation can each be identified unambiguously with an integer. To obtain this integer in O'MEGA, we apply a map from a given set of objects, e.g. wave functions, to the numbers from 1 to N , where N is the cardinality of the set, by using an ordering that ensures that distinct objects are not assigned the same number. To discriminate between particle flavors, there is of course already a well-known ordering that we can use, namely the Particle Data Group (PDG) integers [23].

For the parallel execution, we identify the different levels by the number of external momenta a wave function is connected to or equivalently the number of summands in the momentum of the wave function. This is depicted in Figure 2. Furthermore, we have to group the FUSE_Z instructions to building blocks together with either PROPAGATE_Y or CALC_BRAKET instructions. In this sense, all FUSE_Z instructions are sub-instructions that can belong to either of these two building blocks and the OVM will either form a 1POW $\phi(p+q) = \phi(p)\phi(q)$ or the amplitude $\mathcal{A} = \phi(p)\phi(q)$ depending on the main instruction.

The OVM is initialized with a call that specifies where to find the byte code file, what versions of the OVM and physics model are used, as well as input arrays for masses, widths and couplings, which hold the numeric values for the different types of particles and interactions. In the header of the byte code file the OVM finds the number of momenta, amplitudes (due to multiple color flows and flavor combinations) and wave functions that should be allocated. This is followed by fixed tables for spin, flavor, color flows and color ghosts, for details concerning the color flow formulation cf. Ref. [24], as well as whether a certain flavor-color combination is allowed. Finally, the body of instructions completes the necessary information to compute the cross section.

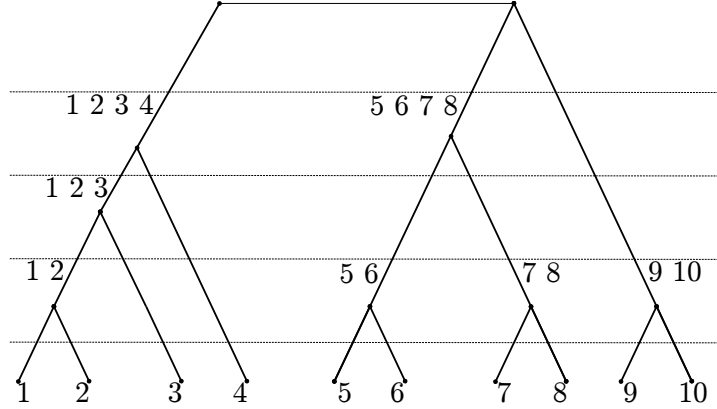


Figure 2: The classification of levels by the number of summands in the momenta yields an unambiguous organization of the calculation whereby each level can be calculated in parallel. We emphasize that this illustration is only one of thousands of possible partitions, whereby each IPOW is heavily reused.

4. Speed Benchmarks

In this section, we benchmark the OVM against the compiled code in Section 4.1, analyze the scaling behavior with multiple cores, which indicates to which degree we are computing in parallel and how much speed up we can expect for more cores, in Section 4.2 and end with a remark on the byte code generation performance in Section 4.3. All processes shown here, and various others, have been validated against the compiled versions where possible for random massless momenta, generated by RAMBO [25], with the help of an automated test suite that is run when `make check` is started in the build folder of O'MEGA. Further tests or benchmarks can be added by appending a single line to the two steering files. We stress that every process is computed in its respective model (QED, QCD or SM) to *full* tree-level order including all interferences and we have not restricted e.g. the Drell-Yan amplitudes to only one electroweak propagator. For simplicity of the test and benchmark suite, we use massless momenta but are in no way restricted to massless theories and do not use simplifications that would render the massless code faster.

To avoid complete compiler dependence of the results, we use two different compilers that are commonly used in scientific projects. These are the GNU and Intel compilers, `gfortran 4.7.1` and `ifort 14.0.3`, respectively. We do not claim that our results are necessarily representative for all Fortran compilers or even compiler versions, but they should still give a good impression of the expected variance in performance. For multi-threading, we use the OPENMP library of the compilers as we are only interested in shared memory parallelization as discussed in Section 2.3. The evaluation time measurements are performed on a computer with two Intel(R) Xeon(R) E5-2440 @ 2.40GHz CPUs, having 16 MiB L3 cache on each socket, and 2x 32 GiB RAM running under Scientific Linux 6.5. The machine has been locked down exclusively for these runs to minimize context switches as far as possible.

4.1. Runtime Performance

In Figure 3, 4 and 5, we show the measured CPU times for QCD, SM and QED processes with two different optimization levels for the compiled code and the OVM using the GNU and

Intel compiler. Since the evaluation times are highly reproducible, we use only three runs to obtain mean and standard deviation. In most cases this results in vanishing error bars. We stress that we show here the relative times normalized for each process to `gfortran-03`, which is why the times are not growing with the number of particles. Absolute times for fully color and helicity summed amplitudes are increasing at least like 2^n due to helicity and like $(n-1)!$ (for the gluon amplitude) due to the number of color flows if no Monte Carlo methods are employed to include these sums in the integration. Lower optimization levels than `-02` are not competitive in terms of run time. For `gfortran`, we observe for most processes the fastest performance with `-03` and for `ifort` with `-02`, which is an effect commonly encountered. The fastest performance is given by the source code compiled with `ifort-02` being roughly 0.75 times the time needed by `gfortran-03`.

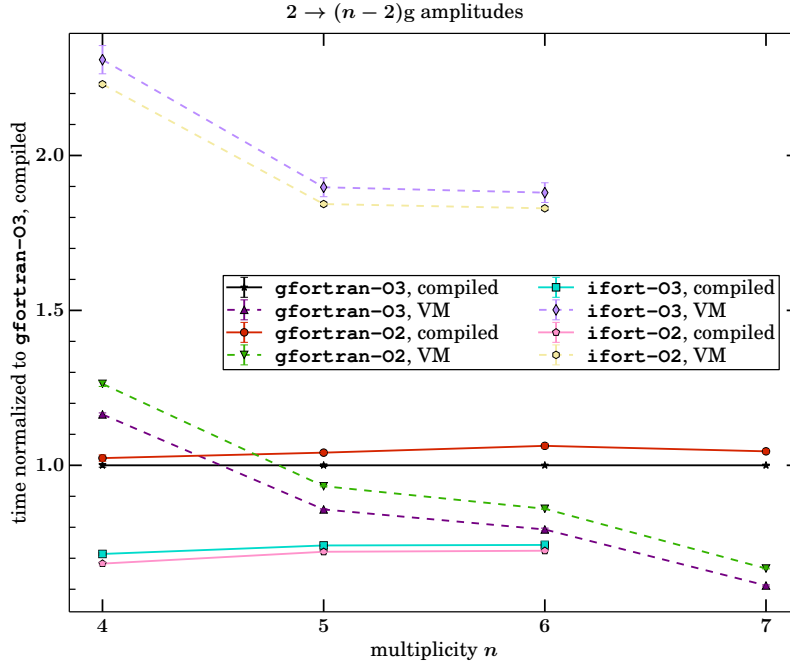


Figure 3: CPU times measured with the Fortran intrinsic `cpu_time` and normalized for each process to the compiled source code using `gfortran -03`. Dashed (solid) lines represent the OVM (compiled source code). The error bars correspond to the standard deviation of three runs.

The crucial point, however, is that `ifort` fails to compile the $n = 7$ gluon and the $u\bar{u} \rightarrow e^+e^-6j$ Drell-Yan process while the OVM immediately starts computing. The GNU compiler is usually able to compile one multiplicity higher compared to the Intel before breaking down. This fits together with the better performance of the compilable processes and longer compile times as `ifort` seems to apply more sophisticated optimization methods to the source code. Disabling the optimizations with `-00` still does not allow to compute the aforementioned processes with both compilers.

Another interesting observation is that the OVM gets faster compared to the compiled code with increasing multiplicity of external particles though this feature is more pronounced in SM

and QCD processes. This is no initialization effect since we allocate the arrays in the beginning and only measure the generation time of matrix elements for M different phase space points. M has been set beforehand for each process with the known approximate scaling for higher multiplicities such that it takes a couple of minutes to complete the computation to have a reliable measurement. The absolute costs for translating an instruction line to actual machine code, i.e. the virtualization costs, are proportional to the number of instructions resulting hence in a constant factor in the relative, normalized time and can not account for this scaling behavior. The most important difference between the compiled source code and the VM is then the explicit double loop in the VM, which goes over the instructions in a level and over all levels as shown in the code excerpt in Section 2.3. The advantages and disadvantages of the double loop are basically the same as general loop unrolling considerations. The native source code represents hereby the unrolled loop that does not have to check for the loop variables, can use latency hiding to start the next instruction while waiting for memory, potentially use CSE² and optimize the prefetching of the processor. The double loop of the VM on the other hand has the advantage of having a higher probability to keep the decode function in the instruction cache. This can potentially explain the scaling behavior with growing complexity compared to the compiled code. We observe roughly the same effect for both compilers, but the OVM compiled with `ifort` is about a factor of two slower than the version with `gfortran`, rendering it not really useful for production runs. This could eventually be solved with a profile-guided optimization, but this is beyond the scope of this work.

Finally, we want to understand the performance difference of the OVM between the QED and QCD amplitudes. To get an impression of the computational complexity, consider that the $e^+e^- \rightarrow 9\gamma$ amplitude is represented by 125 KiB and the $gg \rightarrow 4g$ by 269 KiB of byte code consisting of 3373 and 6780 instructions, respectively. Since the complexity grows exponentially with the number of external particles, those processes can be considered approximatively equally expensive. The difference is, however, that the QED amplitude consists, due to the very high number of external particles, of 8 levels while the QCD amplitude has only 4 levels. This results in about 422 and 1695 instructions per level on average, which is why we can expect worse parallel performance for the QED amplitude due to higher synchronization costs compared to the work to be done per level. Furthermore, this is accompanied with higher memory needs: For the QED amplitude, we need 549 momenta, 256 spinors, 256 conjugated spinors and 9 vector wave functions. The QCD amplitude, on the other hand, requires only 31 momenta and 330 vector wave functions. Returning with this information to the argument made in the last paragraph, the compiled code can gain more from data prefetching in the case of the QED amplitude while the VM improves for more instructions on less data as it is the case for QCD due to the lower chance for instruction cache misses. Overall, we can expect QED to be the worst case scenario for the OVM as it has the lowest number of flavors and the simplest gauge structure one can think of³. Considering all results, we find that the runtimes are in the same order of magnitude and that a VM can be competitive in terms of speed with the compiled version, especially for extreme computations with a high amount of operations per memory.

4.2. Parallelization

Amdahl's idealized law [26] simply divides an algorithm into parallelizable parts p and strictly serial parts $1 - p$. Therefore, the possible speedup s for a computation with n processors

²Although all common subexpressions have in our case already been avoided by O'MEGA.

³Excluding toy models like ϕ^4 theory.

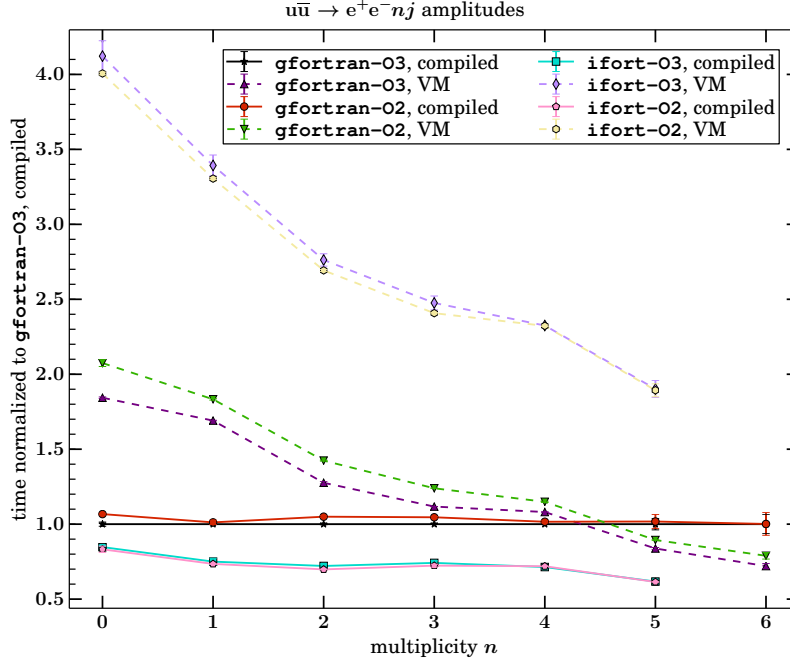


Figure 4: Same as Figure 3 but for the SM Drell-Yan process $u\bar{u} \rightarrow e^+e^-nj$ where $j = u, \bar{u}, g$.

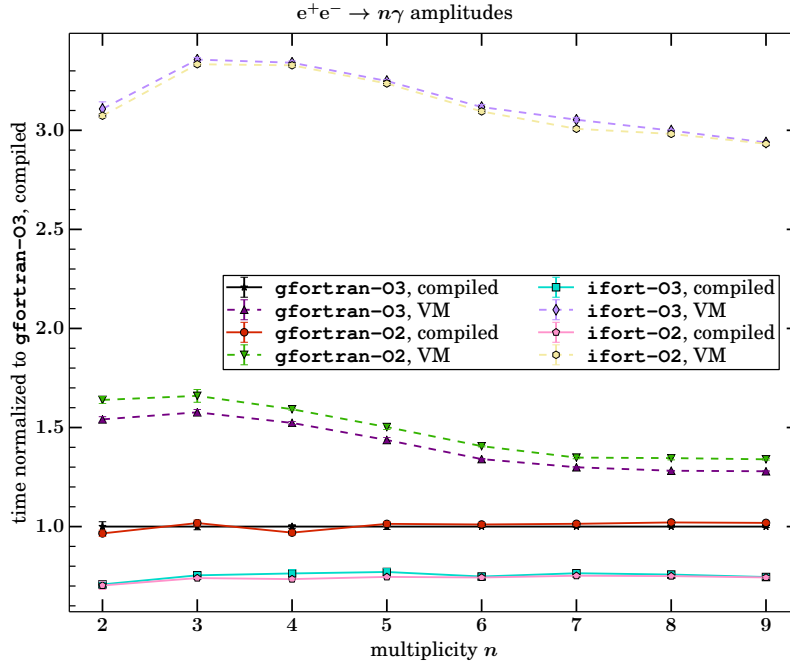


Figure 5: Same as Figure 3 but for quantum electrodynamics (QED) photon production $e^+e^- \rightarrow n\gamma$.

is

$$s(n) \equiv \frac{t(1)}{t(n)} = \frac{1}{(1-p) + \frac{p}{n}}. \quad (1)$$

Communication costs between processors $O(n)$ have been neglected hereby in the denomina-

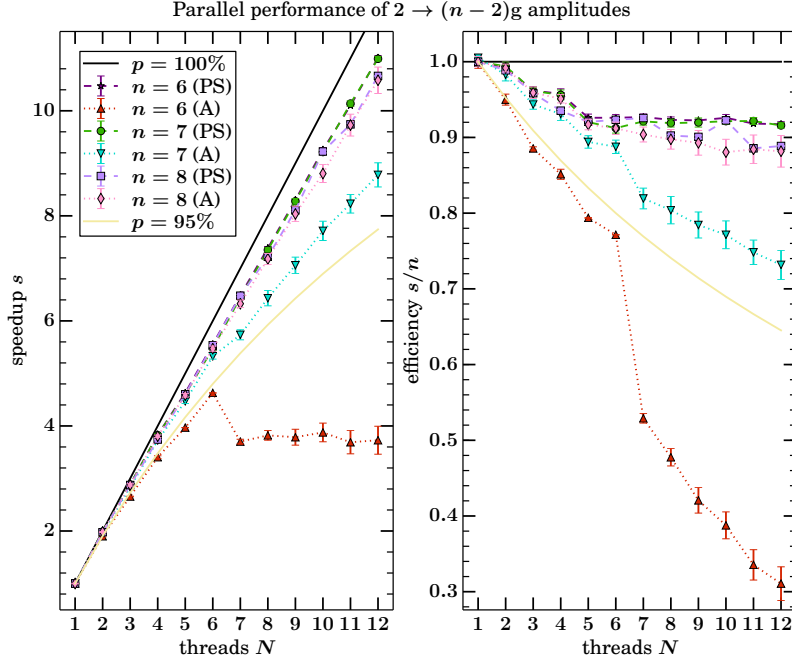


Figure 6: Speedup and efficiency to compute a fixed number of phase space points for the parallel evaluation of multiple phase space points (PS) and the parallel evaluation of the amplitude itself (A) are shown as dashed and dotted lines. The error bars correspond to the standard deviation of three runs. The solid lines represent Amdahl's law for a fixed value of the parallelizable part p .

tor of eq. (1). This means that we have $\lim_{n \rightarrow \infty} s(n) = 1/(1-p)$ in the idealized case and $\lim_{n \rightarrow \infty} s(n) = 0$ including communication costs. In reality, we are interested in high speedups for finite n and also have to care about efficient cache usage. The picture becomes more complicated in modern Non-Uniform Memory Access (NUMA) environments with multiple CPUs on the same board where each socket has its own memory that the others can access as distributed shared memory. For our machine, the two sockets have even and odd numbers for the CPUs on them. To improve the thread scheduling, we have pinned the threads to the cores via

```
GOMP_CPU_AFFINITY='0 2 4 6 8 10 1 3 5 7 9 11'
```

corresponding to using the first CPU for the threads 1–6 and then the second for 7–12. Hyper-threading is disabled as it is not expected to speedup such a calculation. Sadly, we could not achieve any $s > 1$ for the parallelization of the OVM with the Intel compiler neither by using multiple phase space points at once nor by computing the amplitude in parallel. The reason for this is quite unclear, as the exact same code shows the expected speedup with the GNU compiler,

and seems to be correlated with the bad single-core performance of the OVM compiled with `ifort`.

In Figure 6, 7 and 8, we show the speedup with multiple cores N by either using the parallelization procedure, discussed in Section 2.3, to compute one amplitude in parallel or by computing multiple amplitudes for multiple phase space points in parallel again for processes with different multiplicities n in QCD, SM and QED. In a real application the phase space parallelization can not be as efficient as the naive version here, where we can just parallelize the do loop over N_{points} , since usually VEGAS [27] grids are used to approximate the matrix-element and these have to be adjusted iteratively. These book-keeping tasks reduce the parallel parts and the phase-space parallelization shown here (PS) can therefore be regarded as upper bounds. For the parallelization, we chose to only compare the runtime of a single helicity combination to reduce the overall time needed to perform the tests since numerical off-shell recursion algorithms have the same runtime for every helicity, opposed to the closed analytical formulae [28]. To measure the speedup we have used wall clock times as given by the OpenMP function `omp_get_wtime`. In Figure 6, we can see that the $n = 7$ and $n = 8$ gluon amplitudes parallelize very well with both methods with parallelizable parts above 95%. In the shared memory parallel evaluation of the amplitude (A), the impact of the architecture is quite obvious. For $N = 7$, i.e. when the second CPU of the NUMA environment is activated, we see a sharp drop in efficiency, which can be expected since there will be synchronization costs at the end of each level and costs to maintain cache coherency after each instruction inside the amplitude. For the $n = 6$ gluon and the $n = 4$ Drell-Yan process, this effect even leads to a saturation in the speedup indicating that the performance is bound by the memory transfer rates between both CPUs. We observe that this effect becomes less important for more complex amplitudes. To understand this, note that Sandy Bridge with its Intel Quick Path Interconnect (QPI) is actually a cache coherent NUMA architecture, meaning that the cache controllers are required to maintain a consistent memory image when more than one local cache stores the same memory location. Such cache coherency effects have been studied e.g. for the related Nehalem microarchitecture in Ref. [29]. They have shown that the bandwidth to other cores strongly depend on the coherency state of the accessed data. If the latest copy is in the local caches of the remote core, which is more likely to occur for smaller processes, read bandwidths decrease significantly. As expected by the discussion in Section 4.1, the QED amplitudes do only parallelize well if phase space parallelization is used.

The very good performance of the phase space parallelization can be explained by the available cache. The size of the L3 cache per core, 2.7 MiB, is more than enough to host $N = 12$ independent versions of the OVM even for the $n = 8$ gluon amplitude, where momenta, amplitudes and wave functions account to 464.77 KiB. This will break down for this architecture, however, for one multiplicity higher if we extrapolate the given scaling for the number of objects involved in the calculation. The current version of O'MEGA will produce code for *all* color flows of a given process simultaneously. Therefore we have not included the $n = 9$ gluon amplitude in the tests, because the $(9 - 1)! = 40\,320$ different color flow amplitudes do not fit into memory. For real world applications the summation of all color amplitudes will have to be replaced by a sampling of color space.

Either way, it is important to also have the possibility to compute one amplitude in parallel since architectures change and e.g. the Intel Xeon Phi has only 512 KiB cache per core, rendering already the $n = 8$ case close to inappropriate for phase space parallelization.

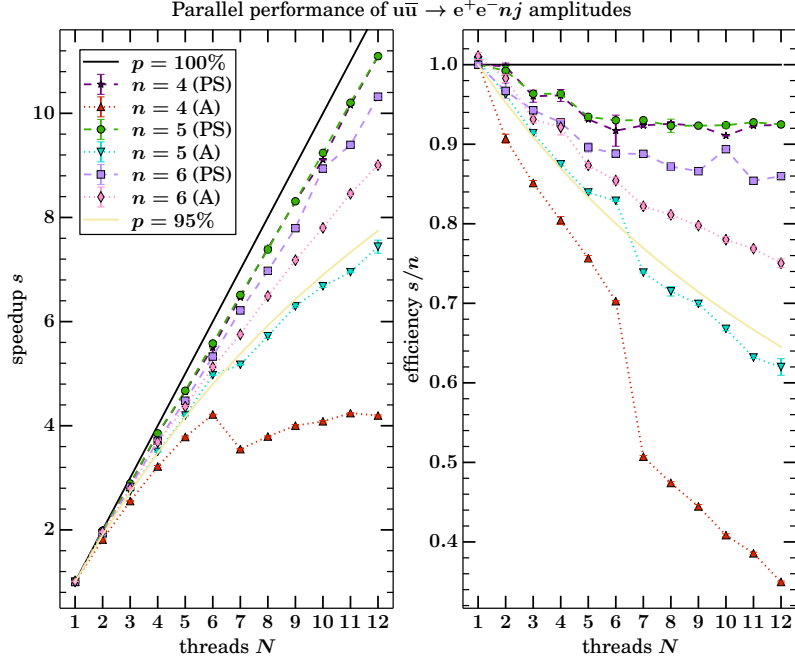


Figure 7: Same as Figure 3 but for the SM Drell-Yan process $u\bar{u} \rightarrow e^+e^-nj$ where $j = u, \bar{u}, g$.

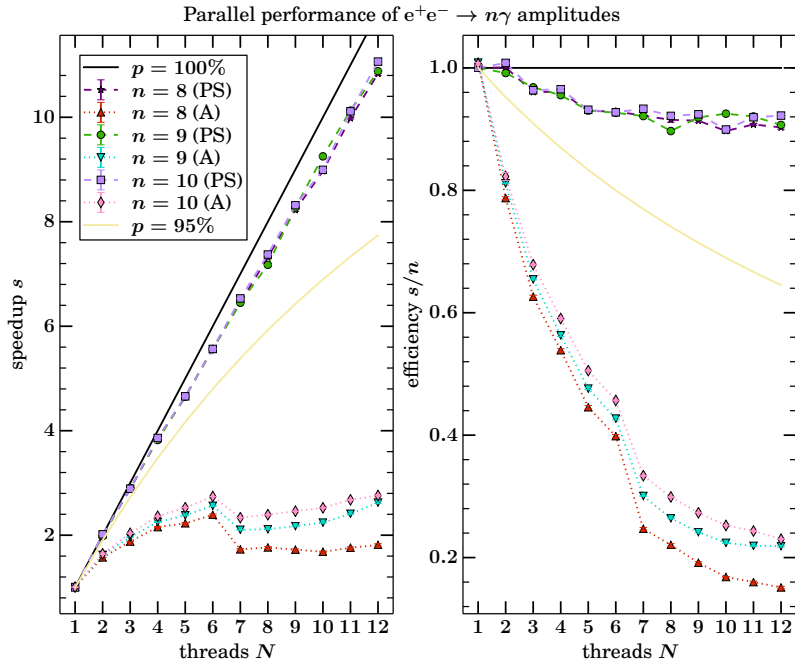


Figure 8: Same as Figure 3 but for QED photon production $e^+e^- \rightarrow n\gamma$.

4.3. Bytecode Generation

It is intuitively clear that integer byte code is smaller than syntactically correct Fortran source code. Furthermore, we use long strings in the source code for debugging purposes, i.e. to directly see to which color flow and momentum combination a IPOW belongs. To be specific, we note that the byte code for the OVM is about one order of magnitude smaller. For convenience, some values together with their old compile times are shown in Table 2. The bytecode size has been furthermore almost halved for very colorful amplitudes in a later version, by using the symmetry of the color factor table, but this could have been achieved with the Fortran output as well and is not shown here. The smaller output format leads to less required RAM and time to produce it. Especially for many color flows, where the generation time of O'MEGA is dominated by the output procedure, we observe e.g. for $gg \rightarrow 6g$ a reduction in memory from 2.17 GiB to 1.34 GiB and in generation time from 11 min 52 s to 3 min 35 s, while staying roughly the same for small processes.

Table 2: Size of the byte code (BC) compared to the Fortran source code together with the corresponding compile time with `gfortran`. The compile times were measured on a computer with an i7-2720QM CPU. The $2g \rightarrow 6g$ process fails to compile.

process	BC size	Fortran size	t_{compile}
$gg \rightarrow gggggg$	428 MiB	4.0 GiB	-
$gg \rightarrow gggggg$	9.4 MiB	85 MiB	483(18) s
$gg \rightarrow q\bar{q}q'\bar{q}'q''\bar{q}''g$	3.2 MiB	27 MiB	166(15) s
$e^+e^- \rightarrow 5(e^+e^-)$	0.7 MiB	1.9 MiB	32.46(13) s

5. Summary and Outlook

A VM circumvents the compile and link problems that are associated with huge source code as it emerges from very complex algebraic expressions. This work is a, to our knowledge first, proof of principle that VMs are indeed a viable option that is maintaining relatively high-performance in the numerical evaluation of these expressions and allows to approach the hardware limits. In practice, a VM saves hours of compile time that would result often enough in internal compiler errors instead of working code. The concept has been successively applied to construct the OVM that is now an alternative method to compute tree-level matrix elements in the publicly available package O'MEGA and will also be integrated in WHIZARD in an upcoming release of the package. Any computation can in principle be performed with a VM though the benefits are clearly in the regime of extreme computations that are not solvable with the conventional method. Here, we have seen that VMs can even perform better than compiled code. Also the parallelization of the amplitude is for very complex processes close to the optimum.

It would be an interesting experiment to reduce the virtualization overhead by using an actual machine to compute matrix elements. The number of instructions corresponding to different wave function fusions and propagators is finite for renormalizable theories (including effective theories up to a fixed mass dimension) and implemented similarly in the various matrix element generators. If the authors can agree on a common set of instructions and conventions this machine could therefore be used by all those programs. The LHC collaborations might actually have a need for this, especially in the light of the HL-LHC, where the number of events for simulation

and reconstruction increases by an order of magnitude and new computing clusters will most likely be needed. Field programmable gate arrays (FPGAs) can serve as such a machine as they have comparable if not superior floating-point performance with respect to current microprocessors and the OVM and its instruction set is the first step to test the feasibility and potential gains of computing matrix elements in this environment. The hardware integration might be quite easy as Intel has recently revealed [30] that Xeon processors can in future be paired with a FPGA in a single socket.

While GPUs and FPGAs are rather unconventional devices that will need large code modifications, similar speedups could be achieved with the Many Integrated Cores (MIC) platform. Various existing scientific applications in Fortran and C++ have been analyzed in Ref. [31] with an early development environment release of the upcoming Intel Xeon Phi. They have shown that it is possible to compile libraries that utilize the Autotools build system for the MIC environment just by setting the proper `./configure` options, at least for static builds. This is a clear advantage as no rewriting is necessary while the speedup can still be in the order of 20 for about 100 threads. Obviously, this strongly depends on having a highly parallel code. We would expect for the OVM speedups in the range of 17 – 50 for processes that exhibit 95 % to 99 % parallel fractions by extrapolating the data of Section 4 and assuming no severe memory problems. In fact, the Xeon Phi possesses no L3 cache at all but a set of coherent L2 caches with less overall cache per core. Thus, we might see a break down in the efficiency of the phase space parallelization, when the objects of one matrix element exceed the L2 cache, while on the other hand high speedups in the parallelization of the amplitude can be maintained.

Acknowledgments

BCN thanks Danny van Dyk, Tomas Jezo and Jos Vermaseren for useful discussions of the idea.

Appendix A. A Trivial Example

To try out the ideas of Section 2, we evaluate a trivial problem here. Consider the identity, for $x \in \mathbb{R}$,

$$\begin{aligned} (-1)^x &= e^{i\pi x} + e^{\ln 2} + \frac{1}{1 - \frac{1}{2}} - \frac{1}{(1 - \frac{1}{2})^2} \\ &\equiv C_1(x) + e^{R_1} + R_2 - R_3, \end{aligned} \quad (\text{A.1})$$

which can be written in terms of the known series C_1 and R_i as

$$\begin{aligned} C_1(x) &= \sum_{n=0}^{\infty} \frac{1}{n!} (i\pi x)^n & R_1 &= \sum_{n=1}^{\infty} (-1)^{n+1} \frac{1}{n} \\ R_2 &= \sum_{n=0}^{\infty} \left(\frac{1}{2}\right)^n & R_3 &= \sum_{n=1}^{\infty} n \left(\frac{1}{2}\right)^{n-1}. \end{aligned} \quad (\text{A.2})$$

To test the VM, we can compute these series explicitly numerically. This means we create byte code for a given N , whereby the above equations follow for $N \rightarrow \infty$, and execute it in the VM. Though these series are not particularly interesting by themselves, it allows us to test the whole

VM infrastructure in a self-contained way, i.e. without dependencies on external libraries⁴. The corresponding template code is freely accessible at <https://github.com/bijancn/basic-vm> and can be used to create a VM for any purpose. It is written in a subset of Fortran2003 that is supported by most modern compilers for mere convenience of the author and due to the environment in which the O'Mega virtual machine (OVM) is used. A translation to C or Fortran95 is straightforward as the code structure is very simple. An earlier version in Fortran95 had the same performance characteristics in the tested cases as the one shown here, indicating that possible performance penalties for the use of some higher-level constructs on the top-level are negligible.

The first step in writing a VM is to identify the set of operations that are needed to perform the computation. Let us associate an instruction line of five integers as operation code (opcode), left-hand side (LHS) and right-hand side (RHS)

```
OPCODE LHS RHS1 RHS2 RHS3.
```

The three fundamental operations, identified by the opcodes 1-3, could then be

```
1 tmp_real(LHS) += RHS1 / RHS2
2 tmp_cplx(LHS) += (const(RHS1) * input(RHS1))RHS2 / table(RHS3)
3 output(LHS) = tmp_cplx(LHS) + exp(tmp_real(RHS1))
  + tmp_real(RHS2) - tmp_real(RHS3).
```

Hereby, we have chosen to put the factorials needed into a separate table. These are given in the byte code as a simple line-by-line array

```
1
1
2
6...
```

The calling application has to supply the constant block

$$\text{const}(1) = i\pi \tag{A.3}$$

and as input data

$$\text{input}(1) = x. \tag{A.4}$$

The byte code for the first elements of $C_1 = 1 + i\pi - \frac{1}{2}\pi^2 + \dots$ now reads, e.g.,

```
2 1 1 0 1
2 1 1 1 2
2 1 1 2 3
2 1 1 3 4
```

and so on. Full example byte codes are part of the repository as well as a Python script to dynamically construct such byte code for any N . Building and running the code is explained in the README.

⁴Except OpenMP, which is needed for parallelization, but the single-threaded execution also works without the library.

References

- [1] J. Alwall, M. Herquet, F. Maltoni, O. Mattelaer, T. Stelzer, MadGraph 5: going beyond, *Journal of High Energy Physics* 2011 (6) (2011) 128. [arXiv:1106.0522v1](#), [doi:10.1007/JHEP06\(2011\)128](#).
- [2] T. Hahn, M. Pérez-Victoria, Automated one-loop calculations in four and D dimensions, *Computer Physics Communications* 118 (2-3) (1999) 153–165. [arXiv:hep-ph/9807565](#), [doi:10.1016/S0010-4655\(98\)00173-8](#).
- [3] B. Chokoufe Nejad, J.-N. Lang, T. Hahn, E. Mirabella, FormCalc 8: Better Algebra and Vectorization, *Acta Physica Polonica B* 44 (11) (2013) 2231. [arXiv:1310.0274](#), [doi:10.5506/APhysPolB.44.2231](#).
- [4] M. Moretti, T. Ohl, J. Reuter, O’Mega: An Optimizing Matrix Element Generator, *arXiv hep-ph (0102195)*. [arXiv:hep-ph/0102195](#).
- [5] B. Ruijl, J. Vermaseren, A. Plaata, J. V. D. Herik, HEPGAME and the Simplification of Expressions, *arXiv* 1405 (6369). [arXiv:1405.6369](#).
- [6] J. Kuipers, T. Ueda, J. Vermaseren, J. Vollinga, FORM version 4.0, *Computer Physics Communications* 184 (5) (2013) 1453–1467. [arXiv:1203.6543](#), [doi:10.1016/j.cpc.2012.12.028](#).
- [7] T. Ueda, J. Vermaseren, Recent developments on FORM, *Journal of Physics: Conference Series* 523 (2014) 012047. [doi:10.1088/1742-6596/523/1/012047](#).
- [8] T. Reiter, Optimising code generation with haggies, *Computer Physics Communications* 181 (7) (2010) 1301–1331. [arXiv:1404.4328](#), [doi:10.1016/j.cpc.2010.01.012](#).
- [9] V. S. Sunderam, PVM: A Framework for Parallel Distributed Computing, *Concurrency: Practice and Experience* 2 (4) (1990) 315–339.
URL http://www.csm.ornl.gov/pvm/pvm_home.html
- [10] D. Cooke, T. Hochberg, F. Alted, I. Vilata, G. Thalhammer, M. Wiebe, G. de Menten, A. Valentino, numexpr - Fast numerical array expression evaluator for Python, NumPy, PyTables, pandas, bcolz and more.
URL <https://github.com/pydata/numexpr>
- [11] T. Hahn, Cuba - a library for multidimensional numerical integration, *Computer Physics Communications* 168 (2) (2005) 78–95. [doi:10.1016/j.cpc.2005.01.010](#).
- [12] K. Hagiwara, J. Kanzaki, N. Okamura, D. Rainwater, T. Stelzer, Fast calculation of HELAS amplitudes using graphics processing unit (GPU), *The European Physical Journal C* 66 (3-4) (2010) 477–492. [arXiv:0908.4403](#), [doi:10.1140/epjc/s10052-010-1276-8](#).
- [13] F. Caravaglios, M. Moretti, An Algorithm to Compute Born Scattering Amplitudes without Feynman Graphs, *Physics Letters B* 358 (3-4) (1995) 332–338. [arXiv:9507237](#), [doi:10.1016/0370-2693\(95\)00971-M](#).
- [14] A. Kanaki, C. G. Papadopoulos, HELAC: A package to compute electroweak helicity amplitudes, *Computer Physics Communications* 132 (3) (2000) 306–315. [arXiv:0002082](#), [doi:10.1016/S0010-4655\(00\)00151-X](#).
- [15] T. Gleisberg, S. Höche, Comix, a new matrix element generator, *Journal of High Energy Physics* 2008 (12) (2008) 039–039. [arXiv:0808.3674v2](#), [doi:10.1088/1126-6708/2008/12/039](#).
- [16] S. Actis, A. Denner, L. Hofer, A. Scharf, S. Uccirati, Recursive generation of one-loop amplitudes in the Standard Model, *arXiv* 1211 (6316). [arXiv:1211.6316](#).
- [17] A. van Hameren, Multi-gluon one-loop amplitudes using tensor integrals, *Journal of High Energy Physics* 2009 (07) (2009) 088–088. [arXiv:0905.1005](#), [doi:10.1088/1126-6708/2009/07/088](#).
- [18] W. Kilian, T. Ohl, J. Reuter, WHIZARD - simulating multi-particle processes at LHC and ILC, *The European Physical Journal C* 71 (9) (2011) 1742. [doi:10.1140/epjc/s10052-011-1742-y](#).
- [19] B. Chokoufe Nejad, *Numerical Calculations of Multi-Jet Cross Sections* (2014).
URL http://www.physik.uni-wuerzburg.de/fileadmin/11030200/Master_Arbeiten/chokoufe_nejad-bijan_master.pdf
- [20] T. Ohl, J. Reuter, W. Kilian, O’Mega: Manual and Commented Source Code (2014).
URL <http://projects.hepforge.org/whizard/omega.pdf>
- [21] C. Degrande, C. Duhr, B. Fuks, D. Grellscheid, O. Mattelaer, T. Reiter, UFO The Universal FeynRules Output, *Computer Physics Communications* 183 (6) (2012) 1201–1214. [doi:10.1016/j.cpc.2012.01.022](#).
- [22] P. de Aquino, W. Link, F. Maltoni, O. Mattelaer, T. Stelzer, ALOHA: Automatic libraries of helicity amplitudes for Feynman diagram computations, *Computer Physics Communications* 183 (10) (2012) 2254–2263. [arXiv:1108.2041v2](#), [doi:10.1016/j.cpc.2012.05.004](#).
- [23] K.A. Olive et al. (Particle Data Group), Review of Particle Physics, *Chinese Physics C* 38 (9) (2014) 090001. [doi:10.1088/1674-1137/38/9/090001](#).
- [24] W. Kilian, T. Ohl, J. Reuter, C. Speckner, QCD in the color-flow representation, *Journal of High Energy Physics* 2012 (10) (2012) 22. [arXiv:1206.3700v2](#), [doi:10.1007/JHEP10\(2012\)022](#).
- [25] R. Kleiss, W. Stirling, S. Ellis, A new Monte Carlo treatment of multiparticle phase space at high energies, *Computer Physics Communications* 40 (1986) 359–373.
- [26] G. Amdahl, *Validity of the Single Processor Approach to Achieving Large-Scale Computing Capabilities*, AFIPS Conference Proceedings 30 (1967) 483–485.
URL <http://www-inst.eecs.berkeley.edu/~n252/paper/Amdahl.pdf>

- [27] G. Lepage, A new algorithm for adaptive multidimensional integration, *Journal of Computational Physics* 27 (2) (1978) 192–203.
- [28] S. Badger, B. Biedermann, L. Hackl, J. Plefka, T. Schuster, P. Uwer, Comparing efficient computation methods for massless QCD tree amplitudes: Closed analytic formulas versus Berends-Giele recursion, *Physical Review D* 87 (3) (2013) 034011. doi:10.1103/PhysRevD.87.034011.
- [29] D. Molka, D. Hackenberg, R. Schone, M. S. Muller, Memory Performance and Cache Coherency Effects on an Intel Nehalem Multiprocessor System, in: 2009 18th International Conference on Parallel Architectures and Compilation Techniques, IEEE, 2009, pp. 261–270. doi:10.1109/PACT.2009.22.
- [30] Intel, *Disrupting the Data Center to Create the Digital Services Economy* (2014).
URL <https://communities.intel.com/community/itpeernetwork/datastack/blog/2014/06/18/disrupting-the-data-center-to-create-the-digital-services-economy>
- [31] K. W. Schulz, R. Ulerich, N. Malaya, P. T. Bauman, R. Stogner, C. Simmons, *Early Experiences Porting Scientific Applications to the Many Integrated Core (MIC) Platform*, TACC - Intel Highly Parallel Computing Symposium MIC.
URL <http://users.ices.utexas.edu/~rhys/papers/SchulzPS2cTIHPCS12.pdf>