



## Original paper

# Monte Carlo Processing on a Chip (MCoaC)-preliminary experiments toward the realization of optimal-hardware for TOPAS/Geant4 to drive discovery

Yogindra S. Abhyankar<sup>a</sup>, Sachin Dev<sup>b</sup>, O.S. Sarun<sup>a</sup>, Amit Saxena<sup>a</sup>, Rajendra Joshi<sup>a</sup>, Hemant Darbari<sup>a</sup>, C. Sajish<sup>a</sup>, U.B. Sonavane<sup>a</sup>, Vivek Gavane<sup>a</sup>, Abhay Deshpande<sup>c</sup>, Tanuja Dixit<sup>c</sup>, Rajesh Harsh<sup>c</sup>, Rajendra Badwe<sup>d</sup>, G.K. Rath<sup>e</sup>, Siddhartha Laskar<sup>d</sup>, Bruce Faddegon<sup>f</sup>, Joseph Perl<sup>g</sup>, Harald Paganetti<sup>h</sup>, Jan Schuemann<sup>h</sup>, Anil Srivastava<sup>i</sup>, Ceferino Obcemea<sup>i</sup>, Asheet K. Nath<sup>a</sup>, Ashok Sharma<sup>e</sup>, Jeffrey Buchsbaum<sup>i,\*</sup>

<sup>a</sup> Centre for Development of Advanced Computing (C-DAC), Pune, India

<sup>b</sup> Open Health Systems Laboratory (OHSL), USA

<sup>c</sup> Society for Applied Microwave Electronics Engineering & Research (SAMEER), Mumbai, India

<sup>d</sup> Tata Memorial Centre (TMC), Mumbai, India

<sup>e</sup> All India Institute of Medical Sciences (AIIMS), Delhi, India

<sup>f</sup> University of California San Francisco, USA

<sup>g</sup> SLAC National Accelerator Laboratory, Menlo Park, USA

<sup>h</sup> Massachusetts General Hospital and Harvard Medical School, Boston, USA

<sup>i</sup> National Cancer Institute (NCI), Bethesda, USA

## ARTICLE INFO

**Keywords:**  
Monte Carlo  
Simulation  
Geant4  
TOPAS  
FPGA

## ABSTRACT

Amongst the scientific frameworks powered by the Monte Carlo (MC) toolkit Geant4 (Agostinelli et al., 2003), the TOPAS (Tool for Particle Simulation) (Perl et al., 2012) is one. TOPAS focuses on providing ease of use, and has significant implementation in the radiation oncology space at present. TOPAS functionality extends across the full capacity of Geant4, is freely available to non-profit users, and is being extended into radiobiology via TOPAS-nBIO (Ramos-Mendez et al., 2018). A current “grand problem” in cancer therapy is to convert the dose of treatment from physical dose to biological dose, optimized ultimately to the individual context of administration of treatment. Biology MC calculations are some of the most complex and require significant computational resources. In order to enhance TOPAS’s ability to become a critical tool to explore the definition and application of biological dose in radiation therapy, we chose to explore the use of Field Programmable Gate Array (FPGA) chips to speedup the Geant4 calculations at the heart of TOPAS, because this approach called “Reconfigurable Computing” (RC), has proven able to produce significant (around 90x) (Sajish et al., 2012) speed increases in scientific computing. Here, we describe initial steps to port Geant4 and TOPAS to be used on FPGA.

We provide performance analysis of the current TOPAS/Geant4 code from an RC implementation perspective. Baseline benchmarks are presented. Achievable performance figures of the subsections of the code on optimal hardware are presented; Aspects of practical implementation of “Monte Carlo on a chip” are also discussed.

## 1. Introduction

This technical report presents a novel, global project’s initiation and initial steps, to use new computer hardware technology and novel thinking to increase the speed of computations critical for biological dose calculations in radiation oncology. If successful, the project goal is to enhance the field, by fusing three broad scientific thoughts to potentially create a new, critical capacity crossing into multiple emerging

areas to ultimately help enable precision, adaptive, biologic treatment planning. The long term goal of this project is to create tools to help implement the capacity to biologically and contextually calculate the dose of all forms of radiation therapy including “blends” of radiation, the immunological status of the patient, the genetics and epigenetics of the patient’s tumor and normal tissue, the biologic response to therapy by the patient to that point in treatment, and to be granular enough to address the complex spatial distribution of normal tissue and tumor

\* Corresponding author.

E-mail address: [jeff.buchsbaum@nih.gov](mailto:jeff.buchsbaum@nih.gov) (J. Buchsbaum).

<https://doi.org/10.1016/j.ejmp.2019.06.016>

Received 29 January 2019; Received in revised form 21 May 2019; Accepted 29 June 2019

Available online 16 July 2019

1120-1797/ Published by Elsevier Ltd on behalf of Associazione Italiana di Fisica Medica.

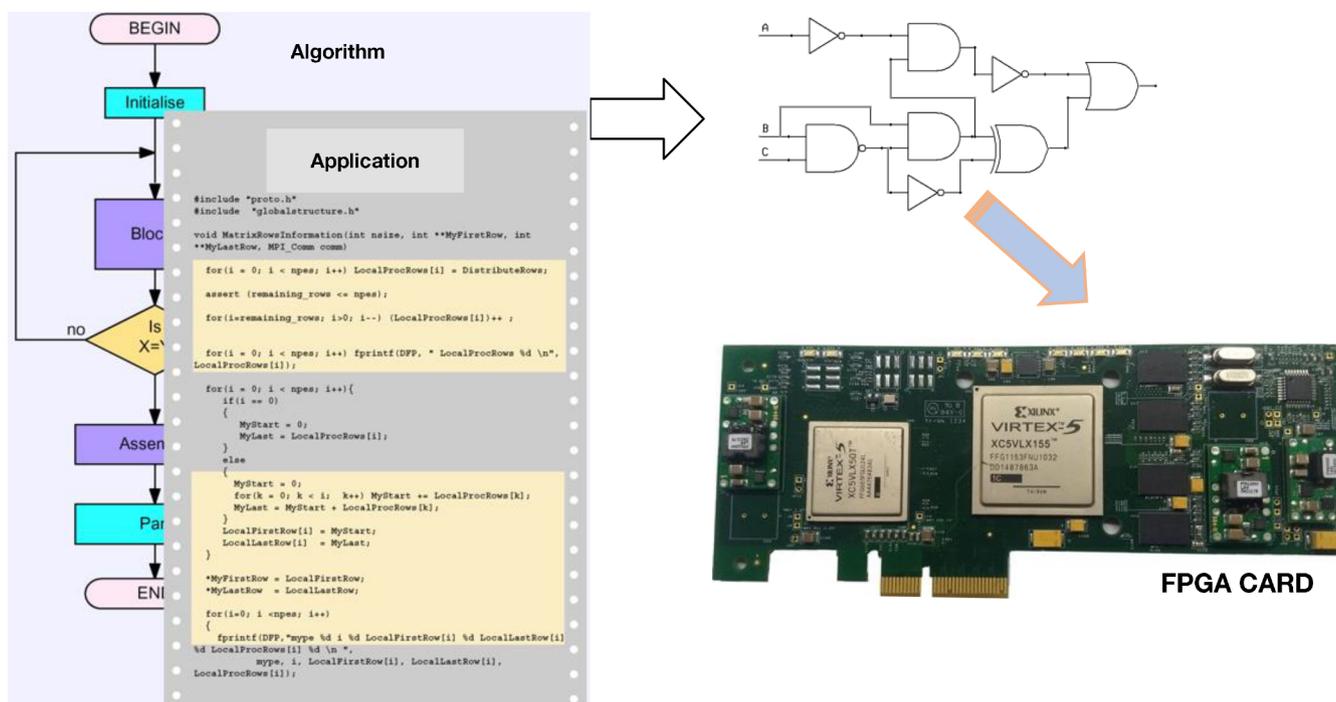


Fig. 1. The algorithm is manually converted to equivalent circuit to extract performance.

tissue. Clearly, this is not the only tool that will be needed to enable this goal but the coalescence of expertise, need, and technology has made its deployment newly possible. Biologic treatment planning is a grand problem of radiation oncology and is driven by the need to properly employ radiation oncology to best treat our patients [5]. Tools to address this problem have broad, worldwide applicability.

The second critical component of this project is the hardware and expertise to implement what has been described as the “reconfigurable supercomputer” via the use of a class of processing unit called the Field Programmable Gate Array (FPGA) [6,7]. The key concept in the use of the FPGA is to achieve a calculation speed increase by converting the series of general CPU instructions, each taking time, into an efficient equivalent-circuit in FPGA (Fig. 1) performing required operations in parallel and thereby saving time. The strength of this approach is the generation of an optimized equivalent circuit; however this approach lacks the inherent flexibility of the generic CPU to adapt to any code. At different times, multiple circuits can be put on FPGA making it reconfigurable, although at any point in time, it has essentially only one configuration. FPGAs are established as one of the best alternatives for solving complex scientific and engineering problems that are energy efficient compared to other available compute accelerators. This is mainly due to the fact that the FPGAs have inherently parallel hardware structure and are customizable as required. Use of this technique achieved a 90-fold improvement in search speed for the open source FASTA [4,8,9] search software. Multiple FPGAs, often on PCIe cards, can be employed at the same time within one system to further enhance this capacity. Beyond radiation oncology, the need for many areas of big science to achieve increased speed makes this approach of universal high impact. The focus on Geant4 [1] specifically will significantly impact the global physics community in this dimension, as many projects in particle physics (such as the large High Energy Physics projects at CERN), particle-astrophysics (such as the Fermi Space Telescope) and materials science (such as studies of radiation effects in electronics) also use Geant4. Thus, the value of MC speedup being implemented in this project is significantly broader than that previously attempted [10].

The third aspect of this project is the transformational ease of using TOPAS in the context of traditional MC programming and package use. It is not an understatement to say that to achieve expertise in Geant4 is

something that can take years. TOPAS provides the end user this power in only a few hours. This will allow biologists and physicians to use MC, in addition to physicists. To achieve the goal of biological treatment planning, familiarizing more people with this tool to assist in research will be critical. In the paper we describe work that is focused on optimizing TOPAS [2] using new approaches in hardware and programming. The biological extension of TOPAS, the TOPAS-nBIO [3], will also benefit directly from this project given the high calculation costs it demands and for the purpose of this paper should be considered part of TOPAS overall. Fast and easy to use MC for biological treatment planning is the desired goal of this project.

## 2. Materials and methods

**Computer:** The computer environment utilized for this work was a multicore Xeon workstation- (E5-2650v2 @2.6 GHz with 128 GB memory), running the Centos version 7 LINUX operating system software with standard programming tools.

**FPGA programming tools:** Xilinx Vivado Design Suite version 2017.3 [11] was used for generating the FPGA design.

**FPGA Card:** C-DAC developed accelerator card [12] with Xilinx FPGA and Xilinx Alveo U200 [13] was used.

**Software Tools:** TOPAS version 3.1.3 [2] including Geant4 [1,14,15] version 10.3.1 was used. Profiler tools OpenSpeedShop v2.3.1 [16], Perf v3.10, Valgrind v3.13 [17], kcachegrind v0.7.2 [18] and IgProf v5.9.16 [19] were used for profiling.

**Design Flow:** FPGAs are semiconductor devices that can be reprogrammed to desired application or functionality requirements after manufacturing. They implement an actual circuit, corresponding to the desired functions. The design flow for porting applications on FPGA accelerators is different than that of a standard software-design-flow. The FPGA design flow consists of various steps such as application profiling, converting algorithm into Hardware Description Language (HDL) design or customization of algorithm to FPGA architecture, functional simulation of the obtained design and implementing the design into FPGA. To extract maximum performance, benchmarking and optimization of the design is performed. It is very difficult to put the whole application on the FPGA; only the compute intensive portion

of the application is run on the FPGA and the rest of the application is run on the CPU. Application profiling is used to identify the most compute intense portions of the application. After the identification of the compute intense portions of the application, they are converted to FPGA design. There are two approaches for doing this. The first one is traditional, that of manually converting the code fragment to HDL design. The second approach that is more recent is to directly convert the high-level C or OpenCL code to HDL using high level synthesis tools. With this approach, the user is not required to have hardware programming knowledge. After the conversion, functional simulation is performed to verify the functionality of the design. Further, the design is synthesized and converted to a programming file for the FPGA. After porting of the compute intense portion of the application onto the FPGA, optimization and benchmarking is done to extract further performance from the FPGA.

**Use case:** SAMEER, India, has developed a low energy oncology system, named Siddhartha, capable of delivering 6 MeV energy photons with flattened dose of 240 cGy/min at 1 m distance [20]. The radiation field generated by this machine is square in shape due to symmetric movement of X-Y jaws and has a maximum field size of 35x35 cm<sup>2</sup>. To test applicability of Geant4/TOPAS on the FPGA card, Geant4 code which describes the basic geometry of SAMEER 6 MeV electron linear accelerator as shown in Fig. 2 was developed. A pencil beam of 6 MeV energy was taken as input to study the bremsstrahlung pattern generated after falling on a high Z target like tungsten or tantalum. The photon output is collimated in the forward direction using primary and secondary collimators. The dose profiles were obtained in a water phantom of dimensions 50 cm × 50 cm with a voxel size of 1 cm<sup>3</sup> for 20 million histories. The dose profile was also compared with the experimentally obtained data to verify the Geant4/TOPAS code. This code was profiled to identify the compute intense functions in Geant4.

Although the main goal is to improve the computational efficiency of TOPAS, the work so far has focused on developments affecting the underlying Geant4 code. A Geant4 MC simulation toolkit-based application described as the use-case above was profiled using multiple open-source profilers such as OpenSpeedShop, pref, Valgrind and

IgProf. Multiple profilers were used to validate the profiled output. For profiling the application, the application source code as well as the Geant4 toolkit was compiled by gcc version 4.8.5 using the `-pg` and `-g` options with no optimization. This helps the profiler in creating a function call-list and file references by extracting information from the application compiled using the above said options.

### 3. Results

The function “G4PhysicsVector::Value” that takes most of the compute time as seen in the data obtained using Valgrind profiler (Fig. 3). The same function surfaced as the most compute function when the application was profiled using other three profiler tools, confirming the selection of the function. As can be seen from the Fig. 3, the function “G4PhysicsVector::Value” is shown as taking the most compute time, but its share in the full execution time is only 4%. Since we were performing the preliminary investigation, we did not plan to do any changes in the Geant4 software, we went ahead with the selection of this function, knowing fully that getting overall speedup will be difficult. The “G4PhysicsVector::Value” function as shown in Fig. 4, calls two inline functions called “FindBin” and “Interpolation”. The code analysis shows that the interpolation function calls “SplineInterpolation” which has many mathematical expressions as compared to “FindBin”. The inline functions do not appear in our profile call-list due to the fact that for the inline functions the compiler at compile time places a copy of the code of that function at each point where the function is called. The call graph for the G4PhysicsVector (Fig. 5) indicates the path that the application code follows, while calling the SplineInterpolation function.

The Spline-interpolation as shown in Fig. 6, is based on a piecewise polynomial called cubic spline, known to reduce interpolation errors. This is a widely used numerical method for interpolation in the scientific domain, compared to other interpolation methods. In order to measure the software execution time of this function independently, we isolated the function code from the Geant4 toolkit. When executed independently, each call to the spline-interpolation function took around

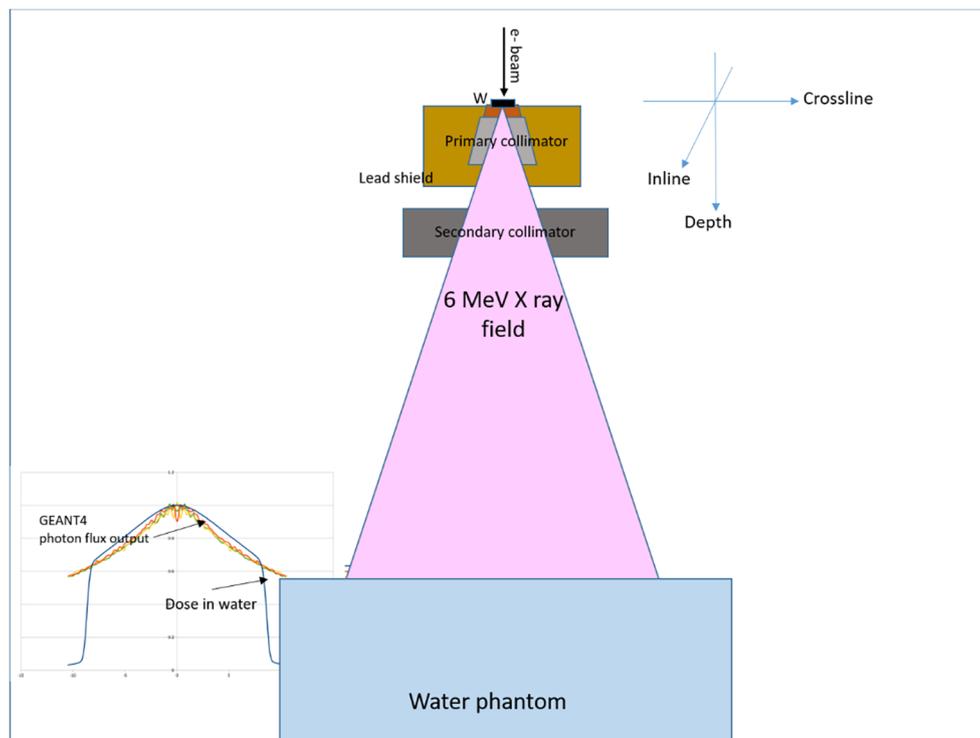


Fig. 2. Experimental setup developed using Geant4 toolkit capable of delivering 6 MeV energy photons.

Self	Called	Function
4.49	196 537 758	G4PhysicsVector::Value(double, unsigned long&) const
3.58	38 425 349	G4SteppingManager::DefinePhysicalStepLength()
3.03	51 271 826	G4NavigationLevel::G4NavigationLevel(G4VPhysicalVolume*, G4A...
2.71	31 434 707	G4Navigator::LocateGlobalPointAndSetup(CLHEP::Hep3Vector cons...
2.62	15 617 666	G4UniversalFluctuation::SampleFluctuations(G4MaterialCutsCouple...
2.36	38 425 350	G4Navigator::ComputeStep(CLHEP::Hep3Vector const&, CLHEP::He...
2.30	38 486 639	G4SteppingManager::Stepping()
2.29	19 114 848	CLHEP::RanecuEngine::flatArray(int, double*)
2.23	163 960 164	CLHEP::RanecuEngine::flat()
1.90	9 917 544	G4UrbanMscModel::SampleCosineTheta(double, double)
1.89	52 660 384	_ieee754_atan2_avx
1.86	38 425 350	G4Transportation::AlongStepGetPhysicalInteractionLength(G4Track ...
1.81	33 299 397	G4Sphere::DistanceToIn(CLHEP::Hep3Vector const&, CLHEP::Hep3V...
1.80	75 445 260	G4VEMProcess::PostStepGetPhysicalInteractionLength(G4Track con...
1.70	3 639 905	G4ParameterisedNavigation::ComputeStep(CLHEP::Hep3Vector con...
1.67	67 811 615	_sin_avx
1.65	47 113 529	G4SteppingManager::InvokePSDIP(unsigned long)
1.62	38 425 349	G4SteppingManager::InvokeAlongStepDoItProcs()
1.53	10 939 070	G4VoxelNavigation::ComputeStep(CLHEP::Hep3Vector const&, CLH...
1.40	40 058 946	G4VEnergyLossProcess::PostStepGetPhysicalInteractionLength(G4T...
1.30	40 594 995	_ieee754_acos_sse2
1.22	38 425 349	G4SteppingManager::InvokePostStepDoItProcs()
1.19	8 564 789	__strtod_l_internal
1.18	8 555 250	std::num_get<>::_M_extract_float(std::istreambuf_iterator<>, std::...
1.17	38 425 350	G4Transportation::PostStepDoIt(G4Track const&, G4Step const&)
1.05	38 425 530	G4ParticleChange::CheckIt(G4Track const&)
1.01	20 029 415	G4UrbanMscModel::ComputeTruePathLengthLimit(G4Track const&, ...
0.94	54 365 620	G4Sphere::DistanceToIn(CLHEP::Hep3Vector const&) const
0.93	27 054 710	_cos_avx
0.93	38 425 350	G4Transportation::AlongStepDoIt(G4Track const&, G4Step const&)
0.91	45 001 252	_int_free
0.90	46 709 703	G4Navigator::LocateGlobalPointWithinVolume(CLHEP::Hep3Vector c...
0.87	104 957 287	G4ProcessManager::GetAttribute(int) const
0.84	48 043 204	G4StepPoint::operator=(G4StepPoint const&)
0.82	8 748 455	G4ParticleHPVector::GetXsec(double)
0.78	3 640 088	G4ParameterisedNavigation::LevelLocate(G4NavigationHistory&, G...
0.78	20 029 563	G4VEnergyLossProcess::AlongStepDoIt(G4Track const&, G4Step co...
0.77	20 029 415	G4UrbanMscModel::ComputeGeomPathLength(double)
0.76	4 815 930	G4SeltzerBergerModel::SampleSecondaries(std::vector<>*, G4Mat...
0.76	45 510 029	_int_malloc
0.73	26 707 245	G4Navigator::ComputeSafety(CLHEP::Hep3Vector const&, double, ...
0.73	101 139 283	std::string::push_back(char)
0.68	307 953 683	G4TouchableHistory::GetVolume(int) const
0.68	8 472 792	G4UrbanMscModel::SampleDisplacement(double, double)
0.68	38 425 350	G4ParticleChangeForTransport::UpdateStepForAlongStep(G4Step*)
0.68	35 958 440	G4Sphere::Inside(CLHEP::Hep3Vector const&) const
0.65	39	G4ParticleHPVector::ThinOut(double)
0.65	20 029 563	G4VMultipleScattering::AlongStepDoIt(G4Track const&, G4Step con...
0.65	23 846 372	G4NormalNavigation::ComputeStep(CLHEP::Hep3Vector const&, CL...
0.64	9 556 564	G4SteppingManager::SetInitialStep(G4Track*)
0.64	38 426 476	G4CrossSectionDataStore::GetCrossSection(G4DynamicParticle con...
0.63	90 171 260	std::basic_string<>::basic_string(std::string const&)
0.59	38 426 295	G4VDiscreteProcess::PostStepGetPhysicalInteractionLength(G4Trac...
0.51	45 509 725	malloc

Fig. 3. Profiles of Geant4 based code generated by Valgrind profiler. The function “G4PhysicsVector::Value” appears at the top of the profile list. The first column ‘Self’ show the time taken by a function as a percentage of the total time. The second column ‘called’ show the number of times the function was called and the third column ‘Function’ show the name of the function.

```
G4double G4PhysicsVector::Value(G4double theEnergy, size_t& lastIdx) const
{
    G4double y;
    if(theEnergy <= edgeMin) {
        lastIdx = 0;
        y = dataVector[0];
    } else if(theEnergy >= edgeMax) {
        lastIdx = numberOfNodes-1;
        y = dataVector[lastIdx];
    } else {
        lastIdx = FindBin(theEnergy, lastIdx);
        y = Interpolation(lastIdx, theEnergy);
    }
    return y;
}
```

Fig. 4. G4PhysicsVector::Value is defined here and the function in-turn calls a function “interpolation”.

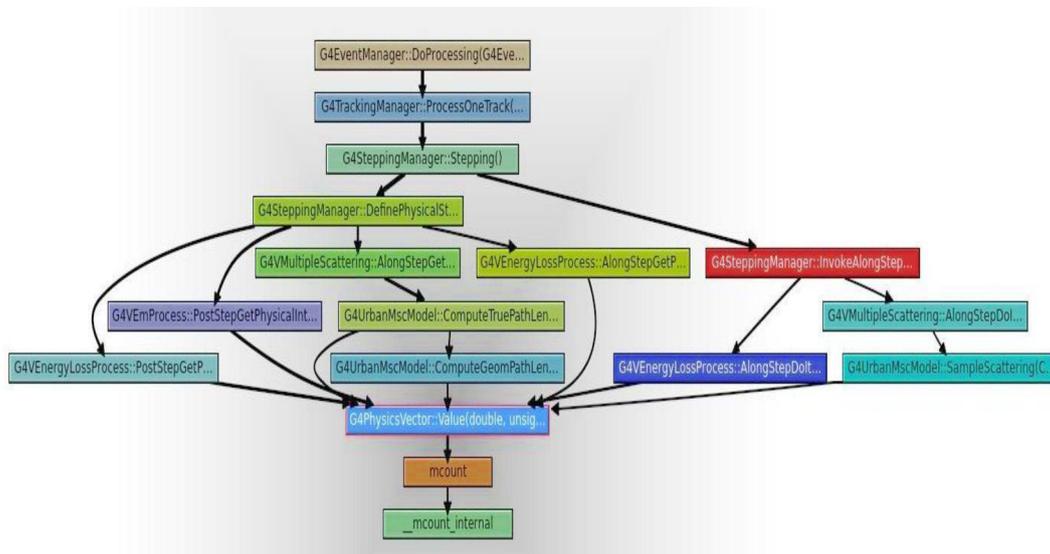


Fig. 5. The call graph for the G4PhysicsVector indicates the path that the Geant4 code follows.

```

inline
G4double G4PhysicsVector::LinearInterpolation(size_t idx, G4double e) const
{
    // Linear interpolation is used to get the value. Before this method
    // is called it is ensured that the energy is inside the bin
    // 0 < idx < numberOfNodes-1

    return dataVector[idx] +
        ( dataVector[idx + 1]-dataVector[idx] ) * ( e - binVector[idx] )
        / ( binVector[idx + 1]-binVector[idx] );
}

//-----

inline
G4double G4PhysicsVector::SplineInterpolation(size_t idx, G4double e) const
{
    // Spline interpolation is used to get the value. Before this method
    // is called it is ensured that the energy is inside the bin
    // 0 < idx < numberOfNodes-1

    static const G4double onesixth = 1.0/6.0;

    // check bin value
    G4double x1 = binVector[idx];
    G4double x2 = binVector[idx + 1];
    G4double delta = x2 - x1;

    G4double a = (x2 - e)/delta;
    G4double b = (e - x1)/delta;

    // Final evaluation of cubic spline polynomial for return
    G4double y1 = dataVector[idx];
    G4double y2 = dataVector[idx + 1];

    G4double res = a*y1 + b*y2 +
        ( (a*a*a - a)*secDerivative[idx] +
          (b*b*b - b)*secDerivative[idx + 1] ) * delta * delta * onesixth;

    return res;
}

//-----

inline
G4double G4PhysicsVector::Interpolation(size_t idx, G4double e) const
{
    return useSpline ? SplineInterpolation(idx, e) : LinearInterpolation(idx, e);
}
    
```

Fig. 6. The spline interpolation is an inline function called by G4PhysicsVector::value and this doesn't show up in our profile list.

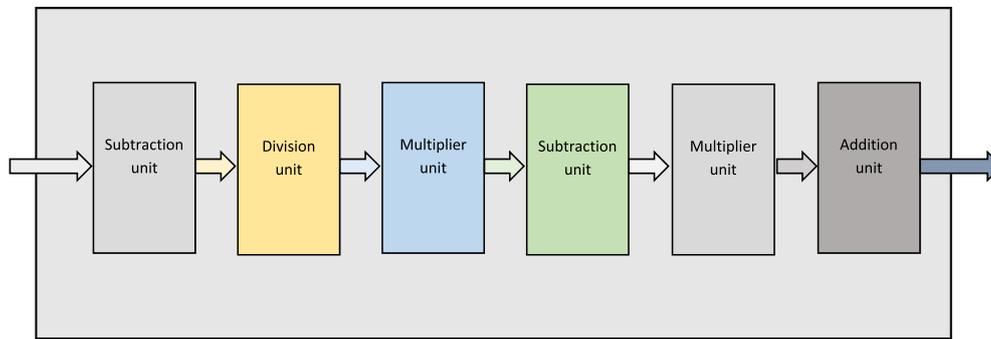


Fig. 7. The simplified block diagram of a spline interpolation calculation.

```

module spline_interpolation #
(
  parameter DOUBLE_DATA_WIDTH = 64,
  parameter ONE_SIXTH = 64'h3FC5555555555555
)
(
  input logic          clk,
  input logic          rst,
  input logic          start,
  input logic [DOUBLE_DATA_WIDTH-1:0] x1, // binVector[idx]
  input logic [DOUBLE_DATA_WIDTH-1:0] x2, // binVector[idx+1]
  input logic [DOUBLE_DATA_WIDTH-1:0] e, // energy
  input logic [DOUBLE_DATA_WIDTH-1:0] y1, // dataVector[idx]
  input logic [DOUBLE_DATA_WIDTH-1:0] y2, // dataVector[idx+1]
  input logic [DOUBLE_DATA_WIDTH-1:0] z1, // secDerivative[idx]
  input logic [DOUBLE_DATA_WIDTH-1:0] z2, // secDerivative[idx+1]
  output logic [DOUBLE_DATA_WIDTH-1:0] res // a*y1 + b*y2 + [(a*a*a - a) * z1 + (b*b*b - b) * z2] * delta *delta * onesixth
);
// delta = x2 - x1 ; a = (x2 - e) / delta; b= (e - x1) / delta ; onesixth= 0x3FC5555555555555

logic [DOUBLE_DATA_WIDTH-1:0] addsub1_a_tdata;
logic          addsub1_a_tvalid;
logic [DOUBLE_DATA_WIDTH-1:0] addsub1_b_tdata;
logic          addsub1_b_tvalid;
logic [DOUBLE_DATA_WIDTH-1:0] addsub1_result_tdata;
logic          addsub1_result_tvalid;
logic          addsub1_operation_tvalid;
logic [7 : 0]          addsub1_operation_tdata;

logic [DOUBLE_DATA_WIDTH-1:0] addsub2_a_tdata;
logic          addsub2_a_tvalid;
logic [DOUBLE_DATA_WIDTH-1:0] addsub2_b_tdata;
logic          addsub2_b_tvalid;
logic [DOUBLE_DATA_WIDTH-1:0] addsub2_result_tdata;
logic          addsub2_result_tvalid;
logic          addsub2_operation_tvalid;
logic [7 : 0]          addsub2_operation_tdata;

logic [DOUBLE_DATA_WIDTH-1:0] addsub3_a_tdata;
logic          addsub3_a_tvalid;
logic [DOUBLE_DATA_WIDTH-1:0] addsub3_b_tdata;
logic          addsub3_b_tvalid;
logic [DOUBLE_DATA_WIDTH-1:0] addsub3_result_tdata;
logic          addsub3_result_tvalid;
logic          addsub3_operation_tvalid;
logic [7 : 0]          addsub3_operation_tdata;

logic [DOUBLE_DATA_WIDTH-1:0] mult1_a_tdata;
logic          mult1_a_tvalid;

```

Fig. 8. System Verilog code for spline interpolation to generate the necessary hardware.

23 ns. This execution time was obtained by using the clock\_gettime (CLOCK\_MONOTONIC, tick) function, a standard method prescribed to measure execution time. To get credible measurement we took the average of multiple runs and the code was compiled with the -O3

optimization option. Since the cubic spline interpolation is a piecewise polynomial equation, we expect the number of calls to this particular function to be extremely high, when the full application is executed. This was also observed in the profile output of our use case where this

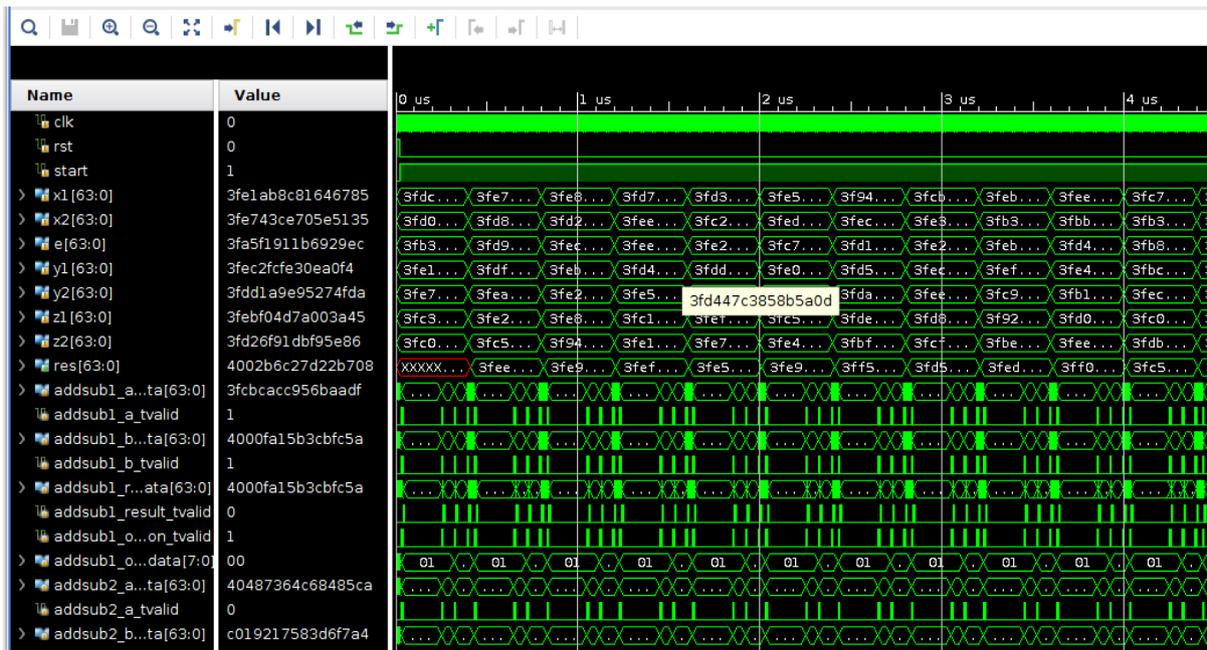


Fig. 9. Simulation waveform of spline interpolation system verilog code.

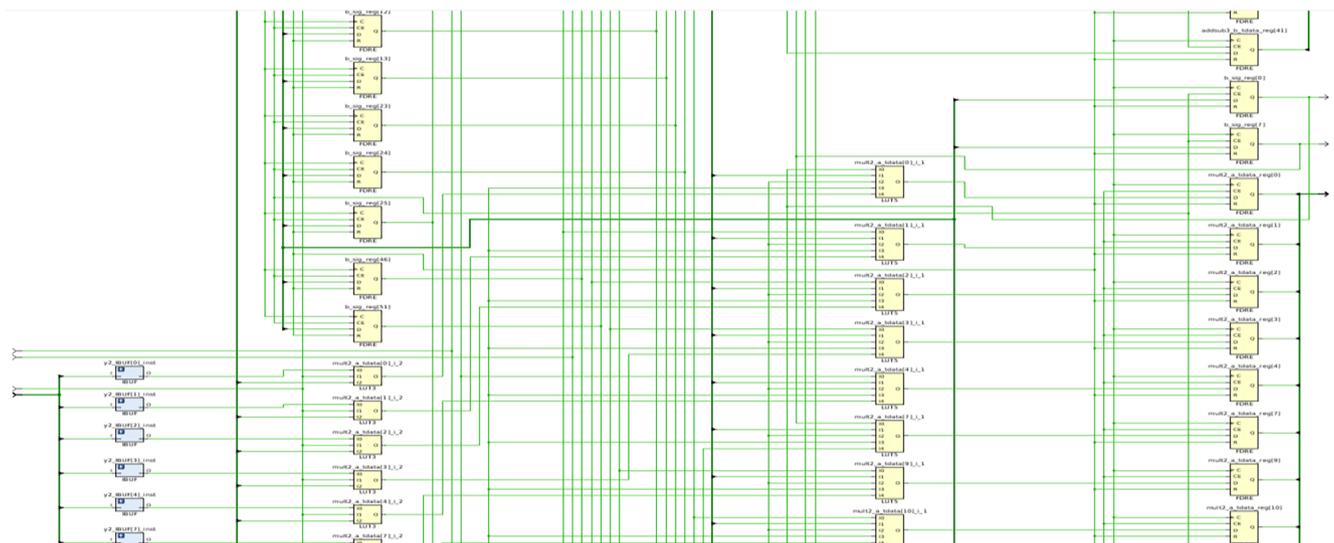


Fig. 10. Synthesis output of spline interpolation system verilog code.

**Table 1**  
Resource and power utilization of single instance of spline interpolation function.

Resource Utilization (Xilinx UltraScale+)	
CLB LUT	0.82%
CLB Registers	0.16%
DSP's	0.37%
Power Utilization	
Total Power	2.76 W
Dynamic Power	0.29 W
Static Power	2.47 W

function was called more than 196 million times. We implemented the spline interpolation function on the FPGA. A simplified block diagram of the spline interpolation function is shown in Fig. 7.

The SplineInterpolation function was coded in System Verilog, a hardware description language to generate hardware. A snapshot of this

code is shown in Fig. 8. The System Verilog code was simulated (Fig. 9) using the Isim simulator that is part of the Xilinx Vivado design suite. To verify the output of our spline-interpolation implementation, we provided the double precision inputs to our implementation and the same was also provided to the spline interpolation software isolated from Geant4. We observed that the double precision output data from our spline-interpolation implementation matched with the output of the isolated code, thus validating our hardware code.

After the simulation, the function was synthesized and implemented into the FPGA using the Xilinx Vivado design suite (Fig. 10). The resource utilization for a single instance of the hardware implementation is shown in Table 1. This preliminary, non-pipelined implementation performs one spline interpolation in 35 clock cycles. This translates to 280 ns when clocked at 125 MHz, the maximum achievable frequency as indicated by the Xilinx tool. The pipelined design can generate output every clock cycle, but requires more resources that bring down the maximum achievable frequency to 100 MHz. By this we can effectively perform spline interpolation in 10 ns. Based on this calculation, a

pipelined design is around 2x faster than the spline interpolation software. As our FPGA is a part of the PCI Express based add-on card, the data transfer time should also be accounted for accurate performance analysis. Currently the time taken to transfer data over PCI Express is significantly higher than the current execution time. One way of mitigating this is by way of porting additional functions of GEANT4 on the FPGA; this in-turn will lead to higher ratio of execution time versus data transfer time.

To validate the energy efficiency of the FPGAs, the power consumed by the spline interpolation hardware-block was measured using the Xilinx tool. A single block consumed around 2.7 W as shown in Table 1. If we populate the whole FPGA, it will consume around 30 W as indicated by the Xilinx power estimator. This value is considerably less than the power consumed by general processors or other accelerators as they are in the order of 90–100 W for CPU and 225–300 W for GPU respectively. Ultimately, these numbers are estimates and the actual power used will vary from case to case. Our preliminary modelling shows, that the FPGAs may be more energy efficient, however further measurements will be required to confirm the same.

#### 4. Discussion

The analysis of Geant4 code has shown it to be a highly optimized code, where no single piece of code occupies over 10% of the total time spent by the CPU. This was not surprising to our group considering the extensive history of the work done to Geant4 for optimizing it across traditional hardware. Despite this fact, our initial work of porting a single piece of Geant4 code on FPGA has shown encouraging results. Using the FPGAs full potential to handle parallel code we expect to see significant reduction in execution time. The results shown here are preliminary and do not reflect full FPGA utilization; nor do they reflect possible code optimization of Geant4 for FPGA implementation. But from our analysis, it becomes clear that a direct one to one implementation of the Geant4 code will not yield enormous speedup. To achieve significant speedup, the Geant4 code reorganization is required in which different functions of Geant4 will be combined and replaced with equivalent FPGA friendly algorithms.

It may be possible using the latest 16 nm and upcoming 7 nm FPGA products to make very large portions of Geant4 fit in one FPGA. Finally, we show that power usage is relatively low for the FPGA solution employed. The high-level synthesis tools, that directly convert high level C or OpenCL code to HDL, look very promising. We plan to use these high-level tools for targeting more functions of Geant4, as this approach will reduce the development time, enabling evaluation of a number of functions in a short span of time.

#### 5. Conclusions

FPGAs represent a possible way to significantly increase the speed of MC calculations and to decrease the power consumption needed to perform them. Our preliminary work shows that the MC code in Geant4 can be ported to an FPGA. To get considerable FPGA speedup, further work is needed to reorganize the Geant4 code. Given the central

importance of treatment planning in radiation oncology and the need to calculate biological dose, FPGA use for MC calculations will have an enormous impact on the field, if successful. As a bonus, FPGA acceleration offers new computational efficiency to all users of the Geant4 toolkit, from radiobiologists to particle physicists, particle-astro-physicists and materials scientists.

#### Acknowledgments

We wish to thank Xilinx (<http://www.xilinx.com>) for providing an FPGA card to the C-DAC (<http://www.cdac.in>) team in Pune, India for benchmarking and code analysis usage.

#### Declaration of Competing Interest

The authors have none to report.

#### References

- [1] Agostinelli S, et al. GEANT4-a simulation toolkit. *Nucl Instrum Methods Phys Sect A* 2003;506(3):250–303.
- [2] Perl J, et al. TOPAS: an innovative proton Monte Carlo platform for research and clinical applications. *Med Phys* 2012;39(11):6818–37.
- [3] Ramos-Mendez J, et al. Monte Carlo simulation of chemistry following radiolysis with TOPAS-nBio. *Phys Med Biol* 2018;63(10):105014.
- [4] Sajish C, Abhyankar Y, Joshi R. Sequence similarity search on reconfigurable computing system. *Int J Computer Electr Eng* 2012;4(5):771–4.
- [5] Ahmed MM, et al. Workshop report for cancer research: defining the shades of gy: utilizing the biological consequences of radiotherapy in the development of new treatment approaches-meeting viewpoint. *Cancer Res* 2018;78(9):2166–70.
- [6] Brown SD, et al. Field programmable gate arrays. Norwell, MA: Kluwer Academic Publishers; 1992.
- [7] Abhyankar Y. Reconfigurable computing (RC). Available from: [https://www.cdac.in/index.aspx?id=hpc\\_cc\\_reconfigurable\\_computing](https://www.cdac.in/index.aspx?id=hpc_cc_reconfigurable_computing).
- [8] Pearson WR. Rapid and sensitive sequence comparison with FASTP and FASTA. *Methods Enzymol* 1990;183:63–98.
- [9] Lipman DJ, Pearson WR. Rapid and sensitive protein similarity searches. *Science* 1985;227(4693):1435–41.
- [10] Fanti V et al. Dose calculation for radiotherapy treatment planning using Monte Carlo methods on FPGA based hardware. In: Real time conference, 2009. RT '09. 16th IEEE-NPSS. 2009; Beijing, China.
- [11] Xilinx. Vivado. 2019; Available from: <https://www.xilinx.com/products/design-tools/vivado.html>.
- [12] C-DAC. Reconfigurable computing. Available from: [https://www.cdac.in/index.aspx?id=hpc\\_cc\\_RCS\\_V4\\_Brochure\\_Nov13](https://www.cdac.in/index.aspx?id=hpc_cc_RCS_V4_Brochure_Nov13).
- [13] Xilinx. u200. 2019; Available from: <https://www.xilinx.com/products/boards-and-kits/alveo/u200.html>.
- [14] Allison J, et al. Geant4 developments and applications. *IEEE Trans Nucl Sci* 2006;53(1):270–8.
- [15] Allison J, et al. Recent developments in GEANT4. *Nucl Instrum Methods Phys Sect A* 2016;835:186–225.
- [16] Open|SpeedShop. 2019; Available from: <https://openspeedshop.org/>.
- [17] Nethercote N, Walsh R, Fitzhardinge J. Building workload characterization tools with Valgrind (invited tutorial). 2006 IEEE international symposium on workload characterization. San Jose, California, USA: IEEE; 2006.
- [18] Weidendorfer J. kcachegrind. 2019; Available from: <https://kcachegrind.github.io/html/Home.html>.
- [19] Eulisse G, Tuura L. IgProf profiling tool, in Computing in High Energy and Nuclear Physics (CHEP 2004). 2004: Interlaken, Switzerland.
- [20] Krishnan R et al. 'S' Band LInac Tube Development Work in Sameer, in PAC2009. 2009.