# Integration of CUDA GPU code in Garfield++

2nd DRD1 Collaboration Meeting

Tom Neep

University of Birmingham

Kostas Nikolopoulos

University of Birmingham

University of Hamburg

Mark Slater

University of Birmingham

2024-06-24

UNIVERSITY OF
BIRMINGHAM

# Introduction

- We are interested in simulating Townsend avalanches using `Garfield++`, specifically using `AvalancheMicroscopic`

- In the past we have found that simulating some high-gain detectors can take a very long time (tracking $10^5, 10^6, \ldots$ electrons)

- As the tracking of each electron in the avalanche is performed independently, this is a good (not perfect) scenario for performing in parallel

- GPUs are excellent at this and are becoming more common in academic environments (your institute likely has GPUs you can access if not then CERN does)

- In the past few years we have been working on adding CUDA (NVIDIA GPU) support to `Garfield++`

- Our code has now been added to the master branch of `Garfield++` (!398)

Integration of CUDA GPU code in Garfield++

UNIVERSITY OF BIRMINGHAM

# Updates

- Some of you may have been in my similar talk from RD51 meeting this time last year

- I will not cover some of the "work in progress" details that I discussed then, but please look back at the slides or ask questions if you are interested

- A lot of discussion in that talk about how we ensured that we get the same results on the CPU and GPU
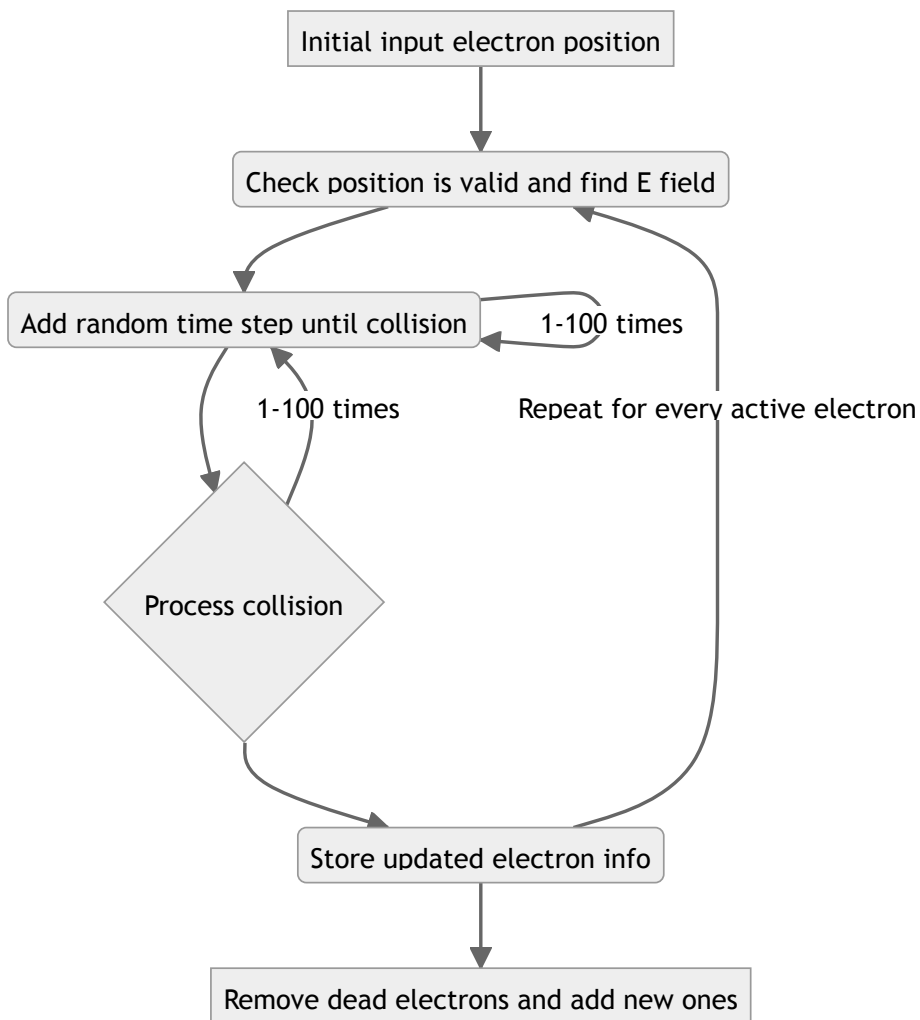
Integration of CUDA GPU code in Garfield++

UNIVERSITY OF BIRMINGHAM

# GPUs

| | **CPU** | **GPU** |
|---|---|---|
| Function | Generalized component that handles main processing functions of a server | Specialized component that excels at parallel computing |
| Processing | Designed for serial instruction processing | Designed for parallel instruction processing |
| Design | Fewer, more powerful cores | More cores than CPUs, but less powerful than CPU cores |
| Best suited for | General purpose computing applications | High-performance computing applications |

Integration of CUDA GPU code in Garfield++

UNIVERSITY OF BIRMINGHAM

# AvalancheMicroscopic



- A rough representation of a single iteration of `AvalancheMicroscopic`

- This is essentially the loop that we want to run on the GPU

Integration of CUDA GPU code in Garfield++

UNIVERSITY OF BIRMINGHAM

# Changes to your code

To install `Garfield++` with GPU support, follow the instructions adding the `USEGPU` option when calling `cmake`

```
cmake -DUSEGPU=ON <src dir>
make -j 8
make install
```

and then change your code to run the avalanche on your GPU

```
1  AvalancheMicroscopic aval;
2  aval.SetSensor(&sensor);
3  aval.SetRunModeOptions(MPRunMode::GPUExclusive, 0);
```

You can run your code as normal and then retrieve the results with the `GetNumberOfElectronEndpointsGPU` and `GetElectronEndpointGPU` methods

```
 1    unsigned int endpoints = aval.GetNumberOfElectronEndpointsGPU();
 2    double xe1, ye1, ze1, te1, e1, xe2, ye2, ze2, te2, e2;
 3    int status;
 4    for (unsigned int i=0; i<endpoints; ++i) {
 5      aval.GetElectronEndpointGPU(
 6            i,
 7            xe1, ye1, ze1, te1, e1,
 8            xe2, ye2, ze2, te2, e2,
 9            status);
10    }
```

Integration of CUDA GPU code in Garfield++

UNIVERSITY OF BIRMINGHAM

# Changes to Garfield++

- We had to make several changes to the core `Garfield++` code

- Changes largely involve making GPU versions of classes

- Preprocessor macros used to "mark" code for GPU or CPU

```
❯ git diff --numstat master -- Source/ Include/Garfield/ | sort -nr
256 24   Source/ComponentFieldMap.cc
255 59   Source/AvalancheMicroscopic.cc
159 13   Include/Garfield/ComponentFieldMap.hh
138 14   Source/MediumMagboltz.cc
97  14   Include/Garfield/Medium.hh
92  0    Include/Garfield/AvalancheMicroscopic.hh
88  28   Include/Garfield/TetrahedralTree.hh
87  8    Include/Garfield/MediumMagboltz.hh
86  8    Include/Garfield/Component.hh
63  6    Include/Garfield/Sensor.hh
47  6    Source/TetrahedralTree.cc
34  6    Source/Sensor.cc
18  1    Include/Garfield/MediumGas.hh
15  0    Include/Garfield/MultiProcessInterface.hh
12  1    Include/Garfield/ComponentAnsys123.hh
8   0    Source/ComponentAnsys123.cc
3   0    Source/MediumGas.cc
3   0    Source/Medium.cc
3   0    Source/Component.cc
3   0    Include/Garfield/RandomEngineRoot.hh
3   0    Include/Garfield/RandomEngine.hh
```

Integration of CUDA GPU code in Garfield++

UNIVERSITY OF BIRMINGHAM

# Example changes required: arrays

```
ComponentFieldMap.cc
 1  #ifdef __GPUCOMPILE__
 2  __device__ void ComponentGPU::Jacobian13(
 3    const double xn[10],
 4    const double yn[10],
 5    const double zn[10],
 6    const double fourt0, const double fourt1,
 7    const double fourt2, const double fourt3,
 8    double& det, double jac[4][4])
 9  #else
10  void ComponentFieldMap::Jacobian13(
11      const std::array<double, 10>& xn,
12      const std::array<double, 10>& yn,
13      const std::array<double, 10>& zn,
14      const double fourt0, const double fourt1,
15      const double fourt2, const double fourt3,
16      double& det, double jac[4][4])
17  #endif
```

- As an example, here is the signature of the `Jacobian13` method

- In this case, we simply need to replace `std::array` with C-style arrays

> The CUDA standard library has advanced since we started this work, and it may now be possible to use `cuda::std::array` to simplify these kind of changes in the future

Integration of CUDA GPU code in Garfield++

UNIVERSITY OF BIRMINGHAM

# Example changes required: vectors

```
ComponentFieldMap.hh

// Bounding boxes of the elements.
#ifdef __GPUCOMPILE__
GPUFLOAT** m_bbMin = nullptr;
GPUFLOAT** m_bbMax = nullptr;
#else
std::vector<std::array<double, 3> > m_bbMin;
std::vector<std::array<double, 3> > m_bbMax;
#endif
```
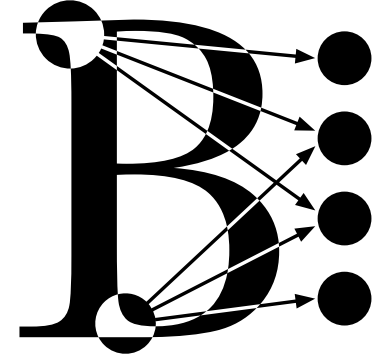
- Likewise we need to replace `std::vector` with dynamically allocated arrays (`cudaMallocManaged`)

- In this case we'd allocate the memory when the GPU version of the class is created and copy from the vector on the CPU

Integration of CUDA GPU code in Garfield++

UNIVERSITY OF BIRMINGHAM

# Benchmarking setup

- To benchmark the code we have adapted a model from `Garfield/Examples/Gem`

- This allows us to use an existing electric field map without needing to create a new one

- This example gives a low, unrealistic gain ($\approx 4$) but serves to demonstrate the performance of the CPU versus the GPU

- To test the performance, we start $N$ electrons at the same position at the entrance of the GEM and see how long it takes to simulate the whole avalanche

Integration of CUDA GPU code in Garfield++

UNIVERSITY OF BIRMINGHAM

# Baskerville

- For this presentation benchmarks were performed on the Baskerville HPC at the University of Birmingham

- 52 liquid-cooled compute trays, each two 36 core Intel Xeon Platinum 8360Y CPUs and four NVIDIA A100 GPUs (meshed with NVIDIA NVLINK)

- Currently 445th on the top500.org list

> **You do not need a supercomputer to benefit from this work!**

# Benchmark results

- Run different number of electrons and fit a straight line to the timings

- GPUs up to 70× faster than CPUs for a large number of initial electrons

- Latest Nvidia GPUs have >2× more cores as A100 - would be interesting to test!



Integration of CUDA GPU code in Garfield++

# Running at CERN

1. `ssh` to `lxplus-gpu.cern.ch` (likely to get a NVIDIA Tesla T4 in my experience)

2. Run an HTCondor job (possibly interactively)

- Run the following command on lxplus to see the GPUs available:

```
❯ condor_status -constraint '!isUndefined(DetectedGPUs)' \
-compact \
-af GPUs_DeviceName TotalGPUs | sort | uniq -c

    31 NVIDIA A100-PCIE-40GB 1
     1 NVIDIA A100-PCIE-40GB 4
     1 Tesla P100-SXM2-16GB 1
     5 Tesla V100-PCIE-32GB 1
     1 Tesla V100-PCIE-32GB 4
    19 Tesla V100S-PCIE-32GB 1
```

Integration of CUDA GPU code in Garfield++

UNIVERSITY OF BIRMINGHAM

# Future work

- Our implementation of `AvalancheMicroscopic` is a starting point and there are many features in `Garfield++` that still don't work on the GPU e.g.

  - Calculating induced signals, which will likely have some *performance impact*

  - Using field maps not created by `Ansys` (should be fairly straight forward to add)

  - Some scope for tidying the user interface

  - Work on generic multithreading

  - **Isabella Oceano** from the University of Hamburg will be investigating some of these issues

- Additionally, there are some technical improvements that could be made:

  - As mentioned, using more of the CUDA standard library might further minimise changes required to CPU code to work on the GPU

  - It should be possible to start an avalanche on the CPU and then switch to the GPU when it reaches a certain size

  - Using multiple GPUs? Using non-Nvidia GPUs?

- We need to identify the features people need to perform their research effectively

UNIVERSITY OF BIRMINGHAM

# Summary

- You can now run on your Nvidia GPU by cloning the `Garfield++` master branch

- `AvalancheMicroscopic` has been implemented, with caveats

- Big speed ups for large avalanches

- Please try running your model on a GPU and see what problems (if any) you run into

- Thank you to **Mark Slater** for the initial GPU implementation and **Heinrich** for reviewing the merge request

> We a preparing a paper which describes the work presented today

Integration of CUDA GPU code in Garfield++

UNIVERSITY OF BIRMINGHAM

# Backup

# Benchmark results

- Run different number of electrons and fit a straight line to the timings
- GPUs up to 70× faster than CPUs for a large number of initial electrons



Integration of CUDA GPU code in Garfield++