# Patatrack - day 1

Andrea, Felice, Marco, Vincenzo, Vishal

# CUDA Graph

- Our RAW to Tracks is running many small kernels

- Got a 4% speed up: baseline 1065Hz → 1100Hz

```cpp
428    bool graphFlag = true;
429    cudaGraph_t graph;
430    cudaGraphExec_t graphExec;
431    void CAHitQuadrupletGeneratorKernels::classifyTuples(HitsOnCPU const & hh, TuplesOnGPU & tuples, cudaStream_t cudaStream) {
432
433        if(graphFlag){
434            cudaStreamBeginCapture(cudaStream);
435            auto blockSize = 64;
436            auto numberOfBlocks = (CAConstants::maxNumberOfQuadruplets() + blockSize - 1)/blockSize;
437
438            kernel_VerifyFit<<<numberOfBlocks, blockSize, 0, cudaStream>>>(tuples.tuples_d, tuples.helix_fit_results_d, tuples.qual
439
440            numberOfBlocks = (CAConstants::maxNumberOfDoublets() + blockSize - 1)/blockSize;
441            kernel_fishboneCleaner<<<numberOfBlocks, blockSize, 0, cudaStream>>>(device_theCells_, device_nCells_,tuples.quality_d)
442
443            numberOfBlocks = (CAConstants::maxNumberOfDoublets() + blockSize - 1)/blockSize;
444            kernel_fastDuplicateRemover<<<numberOfBlocks, blockSize, 0, cudaStream>>>(device_theCells_, device_nCells_,tuples.tuple
445
446            kernel_countTracks<<<numberOfBlocks, blockSize, 0, cudaStream>>>(tuples.tuples_d,tuples.quality_d,counters_);
447            cudaStreamEndCapture(cudaStream, &graph);
448            cudaGraphInstantiate(&graphExec, graph, NULL, NULL, 0 );
449            graphFlag=false;
450        }
451        cudaGraphLaunch(graphExec, cudaStream);
452    }
```

# Tomorrow

- Port more kernels to CUDA Graph and try to launch different streams in the same graph

- Understand how to port CUB kernel calls to Graph

- Profile more kernels