


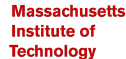
Acts-as-a-Service hackathon

Haoran Zhao, Yuan-Tang Chou, Xiangyang Ju, Elham E Khoda, Andrew Naylor, Yao Yao, Paolo Calafiura, Steven Farrell, Shih-Chieh Hsu, William Patrick McCormack , Philip Coleman Harris , Dylan Sheldon Rankin, Yongbin Feng

ACTS Workshop 2023, Orsay 

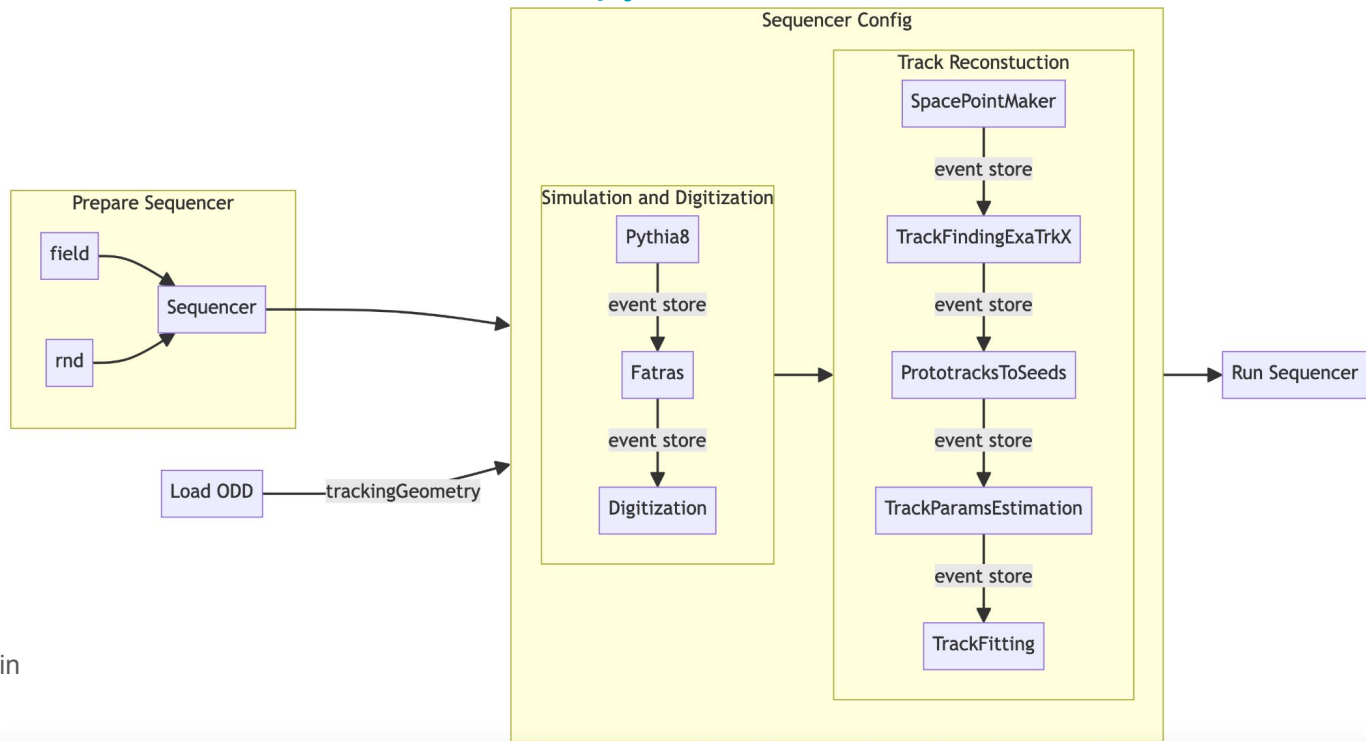


OAC-2117997



Baseline workflow

exatrkx-acts-demonstrator: [inference.py](https://github.com/exatrkx/acts-demonstrator)

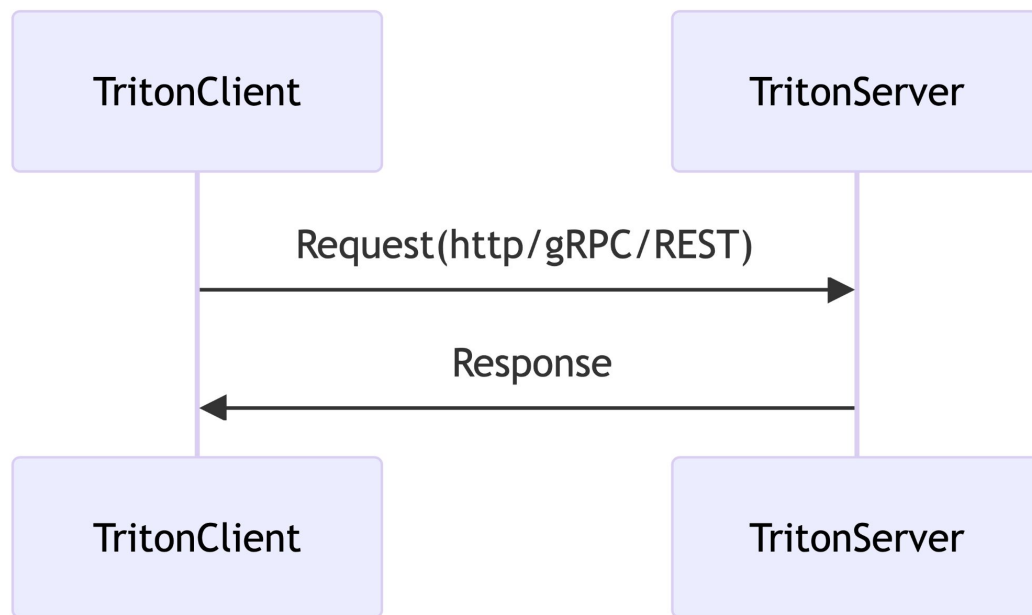


Thanks to Benjamin Huth for the nice demonstrator

Triton Inference Server

Request: inputs for the ML model

Response: inference results



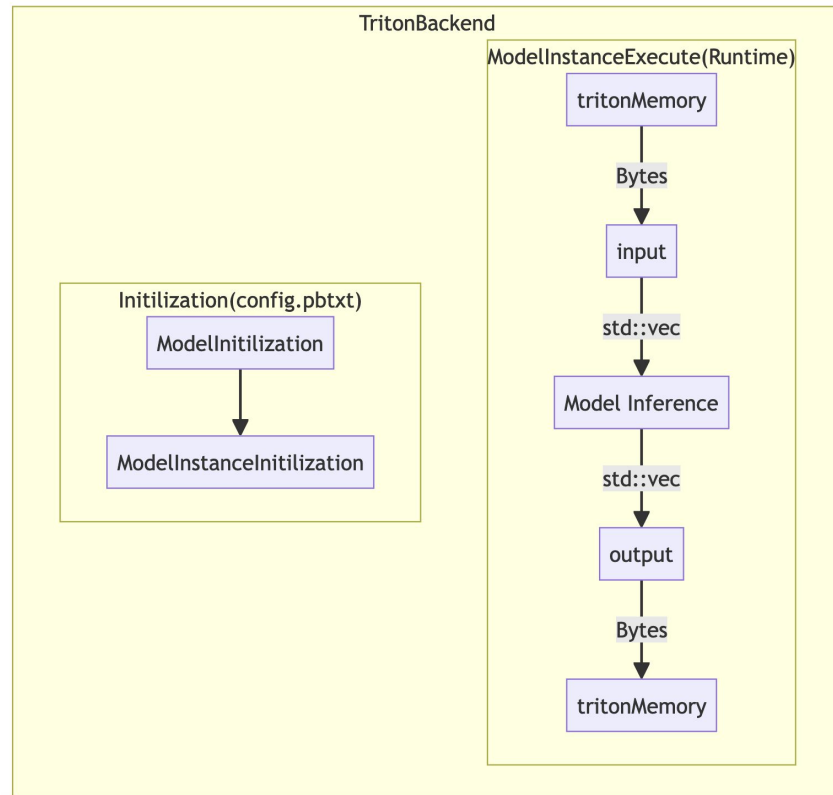
Triton Inference Server Backend

Main components:

- Initialization
- Execution

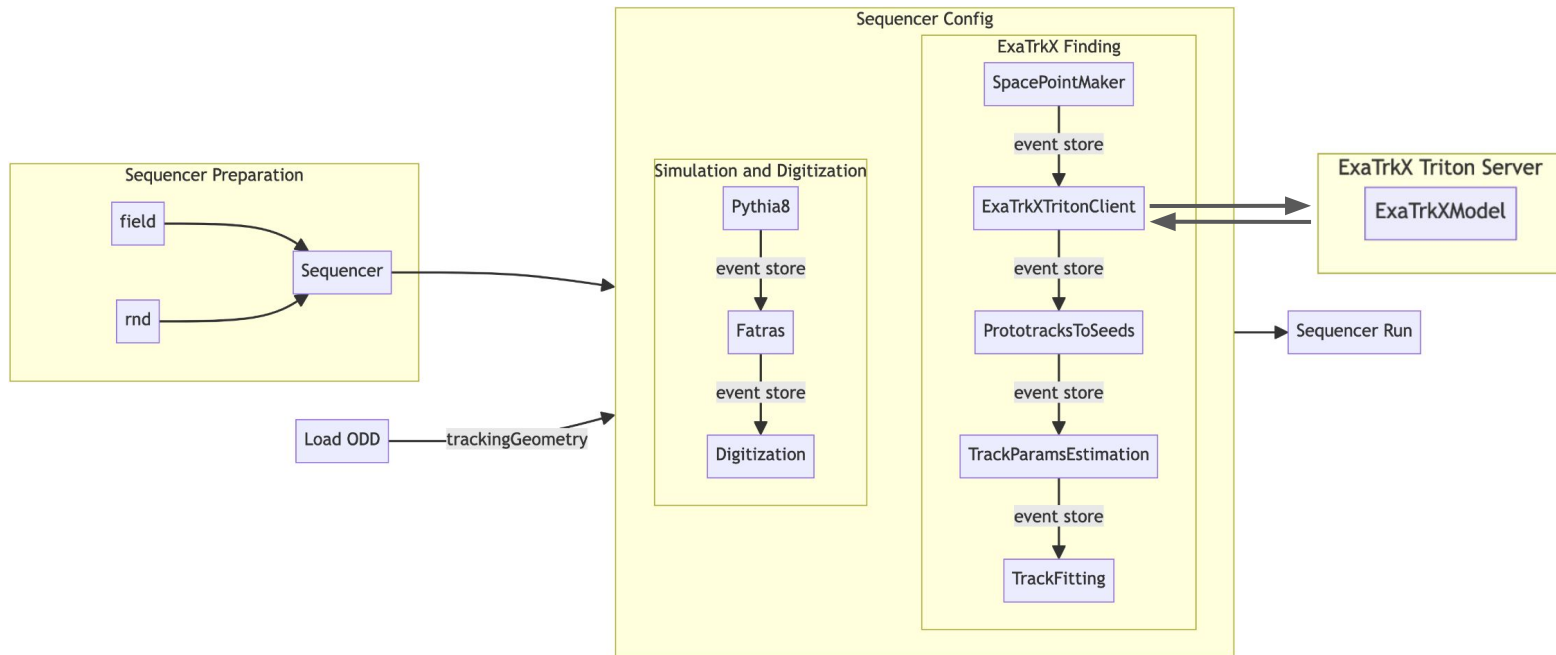
Example:

[exatrkx_gpu backend](#)



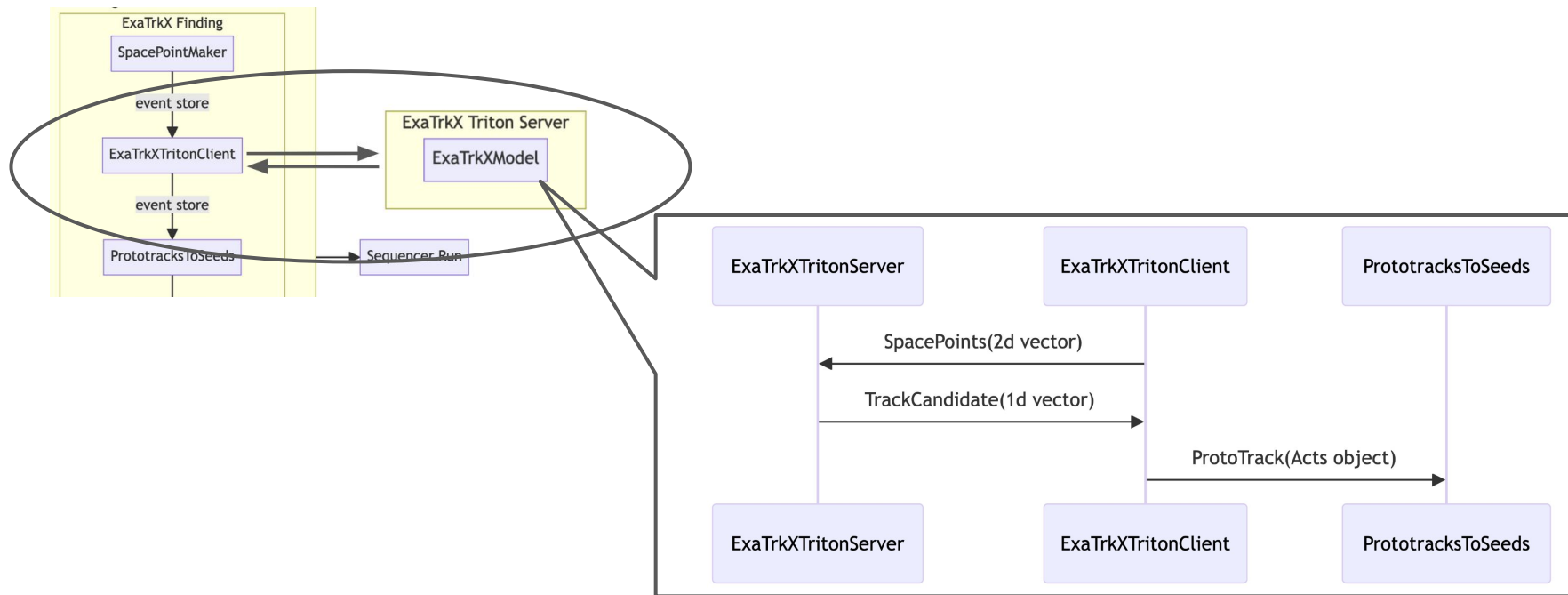
Integration of the ExaTrkX-as-a-Service to ACTS

External ExaTrkX server, see [Yuan-Tang's slide](#)



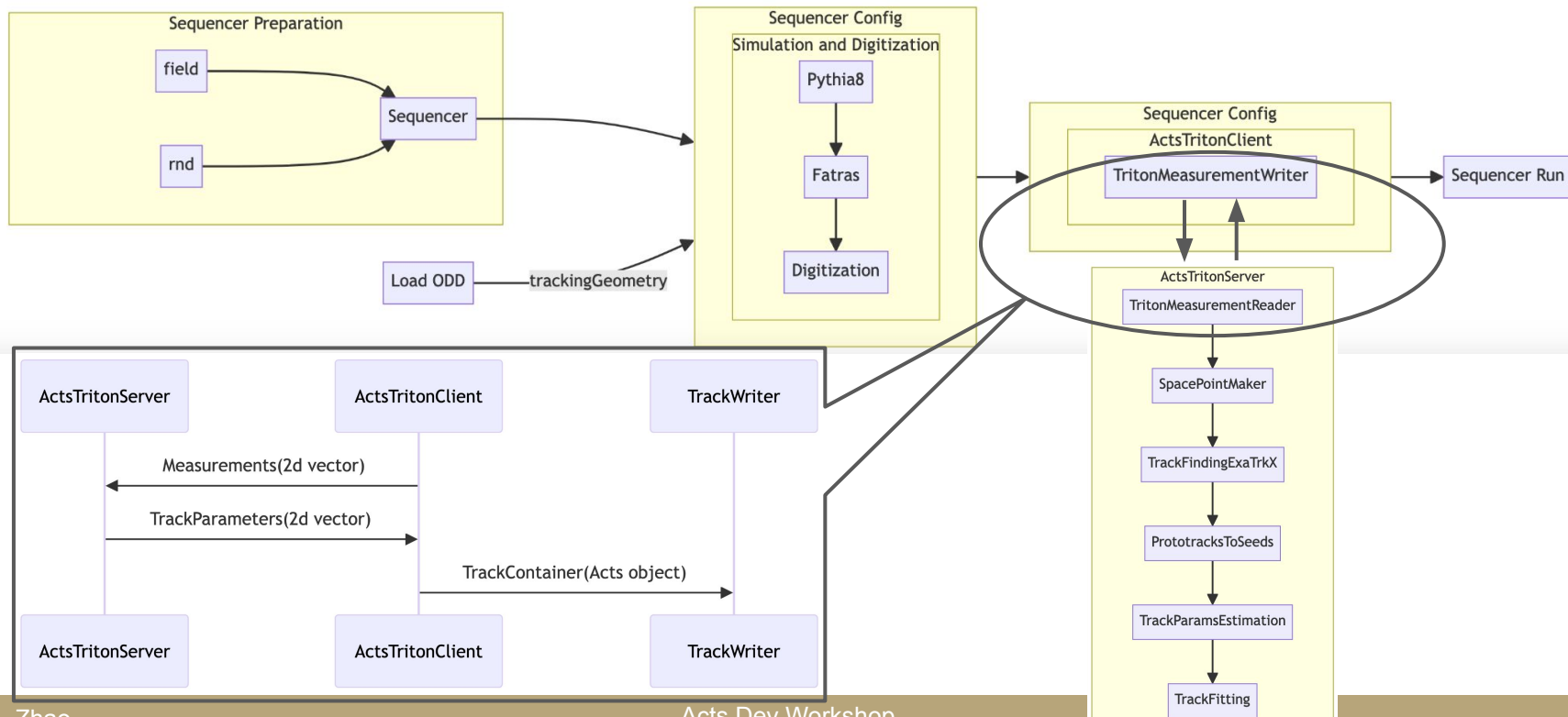
Integration of the ExaTrkX-as-a-Service to ACTS

External ExaTrkX server, see [Yuan-Tang's slide](#)



Acts-aaS

Our proposed workflow, measurements in, track out



Discussions

1. I/O between the client and the server
2. Dataflow on the server side
3. Miscellaneous questions

Discussions

1. I/O between the client and the server

a. What parts does the server do?

Goal in execution: TrackFinding(GPU) + TrackFitting(GPU) + SpacePointMaker?

Sourcelinks needed for fitting

In the initialization: load the ML models, config the B-field, tracking geometry etc..

Discussions

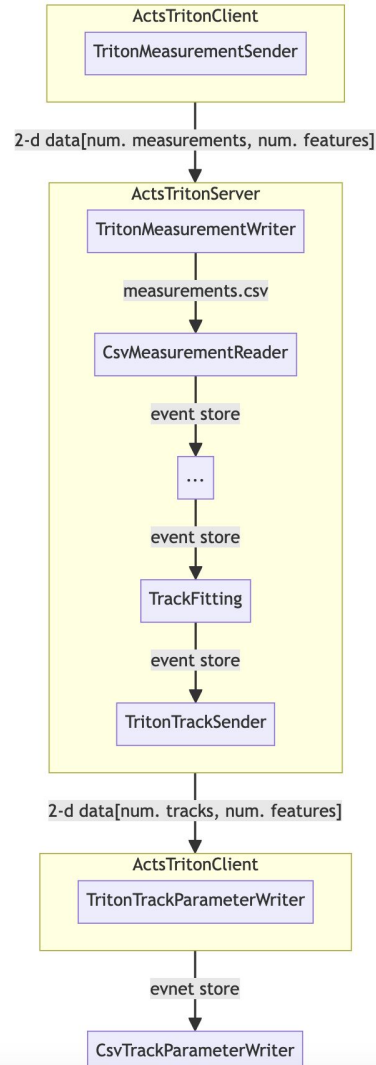
1. I/O between the client and the server
 - a. What parts does the server do?
 - b. Only the primitive data types are supported;
Bool, int, float, string

Proposal: send 2-d data with shape [number of measurements, 13]
Wrap to the composite data type on the server

```
1 measurement_id, geometry_id, local_key, local0, local1, phi, theta, time, var_local0, var_local1, var_phi, var_theta, var_time
2 0, 1152921779484760322, A_k, -0.988179684, 27.5145588, 0, 0, 0, 9.99999975e-05, 9.99999975e-05, 0, 0, 0
3 1, 1152921779484760322, A_k, -1.59969449, 25.6581497, 0, 0, 0, 9.99999975e-05, 9.99999975e-05, 0, 0, 0
4 2, 1152921779484760322, A_k, -7.61218023, 25.5525475, 0, 0, 0, 9.99999975e-05, 9.99999975e-05, 0, 0, 0
5 3, 1152921779484760578, A_k, -8.14472675, 28.2101498, 0, 0, 0, 9.99999975e-05, 9.99999975e-05, 0, 0, 0
```

Discussions

1. I/O between the client and the server
2. Dataflow: Sequencer vs no sequencer option on the server side
 - a. Sequencer + python backend
Port inference.py to the server
 - i. Pro: utilises the event store; least development effort
 - ii. Con: I/O time with storing csv files



Discussions

1. I/O between the client and the server
2. Dataflow: Sequencer vs no sequencer option on the server side
 - a. Sequencer + python backend
 - b. No sequencer + custom backend
 - i. Pro: reduce the I/O
 - ii. Con: code refactoring

```
369 // the Kalman fitter
370 ActsExamples::TrackFittingAlgorithm::Config trackFittingConfig;
371 trackFittingConfig.inputMeasurements = "measurements";
372 trackFittingConfig.inputSourceLinks = "sourcelinks";
373 trackFittingConfig.inputProtoTracks = "extrkrx_prototracks";
374 trackFittingConfig.inputInitialTrackParameters = "trackParameters";
375 trackFittingConfig.outputTracks = "tracks";
376
377 Acts::FreeToBoundCorrection freeToBoundCorrection;
378 freeToBoundCorrection.apply =
379     false; // Explicitly set to false, though it's the default
380
381 std::shared_ptr<ActsExamples::TrackFitterFunction> myKalmanFitter =
382     ActsExamples::makeKalmanFitterFunction(trackingGeometry, field, true,
383     true, 0.0, freeToBoundCorrection);
384
385 Acts::CalibrationContext calibContext;
386 std::shared_ptr<ActsExamples::MeasurementCalibrator> calibrator;
387 ActsExamples::TrackFittingAlgorithm trackFitting(trackFittingConfig,
388     Acts::Logging::INFO);
389
390 ActsExamples::ConstTrackContainer consttracksContainer =
391     trackFitting.executeTrackFitting(measurements, sourceLinks, protoTracks,
392     trackParameters, nullptr, geoContext,
393     magFieldContext, calibContext,
394     calibrator, myKalmanFitter);
395
```

Discussions - code refactoring example

TrackFittingAlgorithm.hpp

```
63 // Framework execute method of the fitting algorithm
64 //
65 // @param ctx is the algorithm context that holds event-wise information
66 // @return a process code to steer the algorithm flow
67 ActsExamples::ProcessCode execute(const AlgorithmContext& ctx) const final;
68
69 // Refactor the execute method without the need of an AlgorithmContext
70 //
71 // @param measurements is the input measurements
72 // @param sourceLinks is the input source links
73 // @param protoTracks is the input proto tracks
74 // @param initialParameters is the input initial track parameters
75 // @param clusters is the input clusters
76 // @param geoContext is the geometry context
77 // @param magFieldContext is the magnetic field context
78 // @param calibContext is the calibration context
79 // @param calibrator is the measurement calibrator
80 // @param fit is the track fitter function
81 // @return the fitted tracks in a container
82 ConstTrackContainer executeTrackFitting( You, 2 days ago • adding TrackFi
83     const MeasurementContainer& measurements,
84     const IndexSourceLinkContainer& sourceLinks,
85     const ProtoTrackContainer& protoTracks,
86     const TrackParametersContainer& initialParameters,
87     const ClusterContainer* clusters,
88     Acts::GeometryContext& geoContext,
89     Acts::MagneticFieldContext& magFieldContext,
90     Acts::CalibrationContext& calibContext,
91     std::shared_ptr<MeasurementCalibrator>& calibrator,
92     std::shared_ptr<TrackFitterFunction>& fit) const;
```

TrackFittingAlgorithm.cpp

```
69 ActsExamples::ProcessCode ActsExamples::TrackFittingAlgorithm::execute(
70     const ActsExamples::AlgorithmContext& ctx) const {
71     // Read input data
72     const auto& measurements = m_inputMeasurements(ctx);
73     const auto& sourceLinks = m_inputSourceLinks(ctx);
74     const auto& protoTracks = m_inputProtoTracks(ctx);
75     const auto& initialParameters = m_inputInitialTrackParameters(ctx);
76
77     const ClusterContainer* clusters =
78         m_inputClusters.isInitialized() ? &m_inputClusters(ctx) : nullptr;
79
80     Acts::GeometryContext geoContext = ctx.geoContext;
81     Acts::MagneticFieldContext magFieldContext = ctx.magFieldContext;
82     Acts::CalibrationContext calibContext = ctx.calibContext;
83     std::shared_ptr<MeasurementCalibrator> calibrator = m_cfg.calibrator;
84     std::shared_ptr<TrackFitterFunction> fit = m_cfg.fit;
85
86     auto constTracks = executeTrackFitting(measurements, sourceLinks, protoTracks,
87                                           clusters, geoContext,
88                                           magFieldContext, calibContext, calibrator, fit);
89
90     m_outputTracks(ctx, std::move(constTracks));
91     return ActsExamples::ProcessCode::SUCCESS;
92 }
```

Discussions - Miscellaneous

Surface creation failed w/o geoIDhook(thought optional..)

A global geoContext for odd.py?

odd.py - getOpenDataDetector()

```
56 def geoid_hook(geoid, surface): Benjamin Huth, 10 months ago • fea
57     if geoid.volume() in volumeRadiusCutsMap:
58         r = sqrt(surface.center()[0] ** 2 + surface.center()[1] ** 2)
59
60         geoid.setExtra(1)
61         for cut in volumeRadiusCutsMap[geoid.volume()]:
62             if r > cut:
63                 geoid.setExtra(geoid.extra() + 1)
64
65         return geoid
66
67 dd4hepConfig = acts.examples.dd4hep.DD4hepGeometryService.Config(
68     xmlFileNames=[str(odd_xml)],
69     logLevel=customLogLevel(),
70     dd4hepLogLevel=customLogLevel(),
71     geometryIdentifierHook=acts.GeometryIdentifierHook(geoid_hook),
72 )
73 detector = acts.examples.dd4hep.DD4hepDetector()
```

```
38 class VolumeRadiusGeometryIdentifierHook : public Acts::GeometryIdentifierHook {
39     public:
40         static const std::unordered_map<unsigned int, std::vector<double>>
41             volumeRadiusCutsMap;
42
43         // overwrite decorateIdentifier method
44         virtual Acts::GeometryIdentifier decorateIdentifier(
45             Acts::GeometryIdentifier identifier,
46             const Acts::Surface& surface) const override {
47             Acts::GeometryContext geoContext;
48             auto volumeIter = volumeRadiusCutsMap.find(identifier.volume());
49             if (volumeIter != volumeRadiusCutsMap.end()) {
50                 // FIXME: surface needs a GeometryContext, not sure how odd.py does it
51                 Acts::Vector3 centerPoint = surface.center(geoContext);
52
53                 double r = std::sqrt(centerPoint.x() * centerPoint.x() +
54                                     centerPoint.y() * centerPoint.y());
55
56                 unsigned int extraValue = 1;
57                 for (double cut : volumeIter->second) {
58                     if (r > cut) {
59                         extraValue++;
60                     }
61                 }
62                 identifier.setExtra(extraValue);
63             }
64             return identifier;
65         }
66     };
```

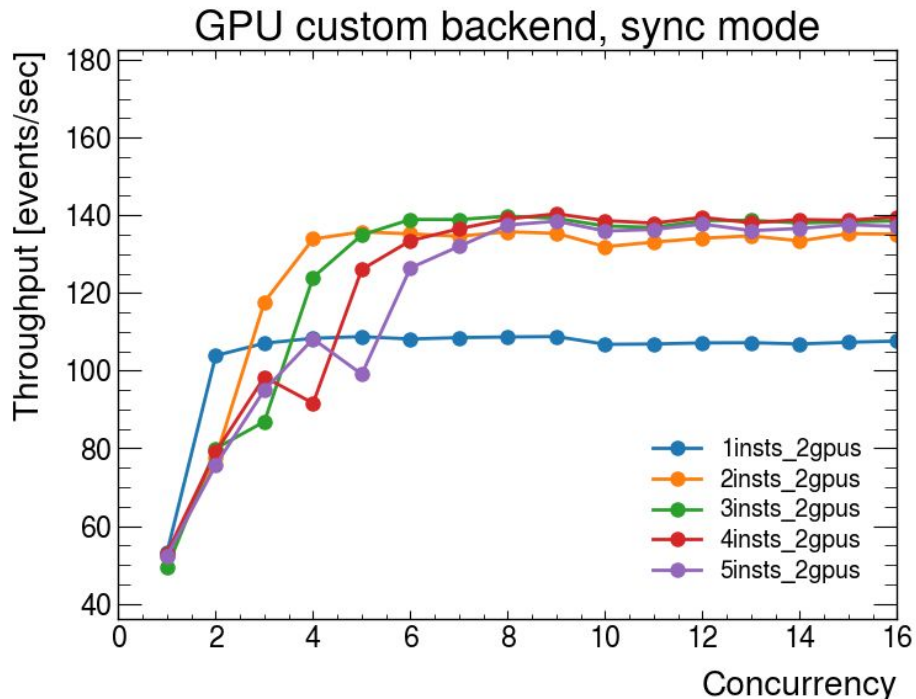
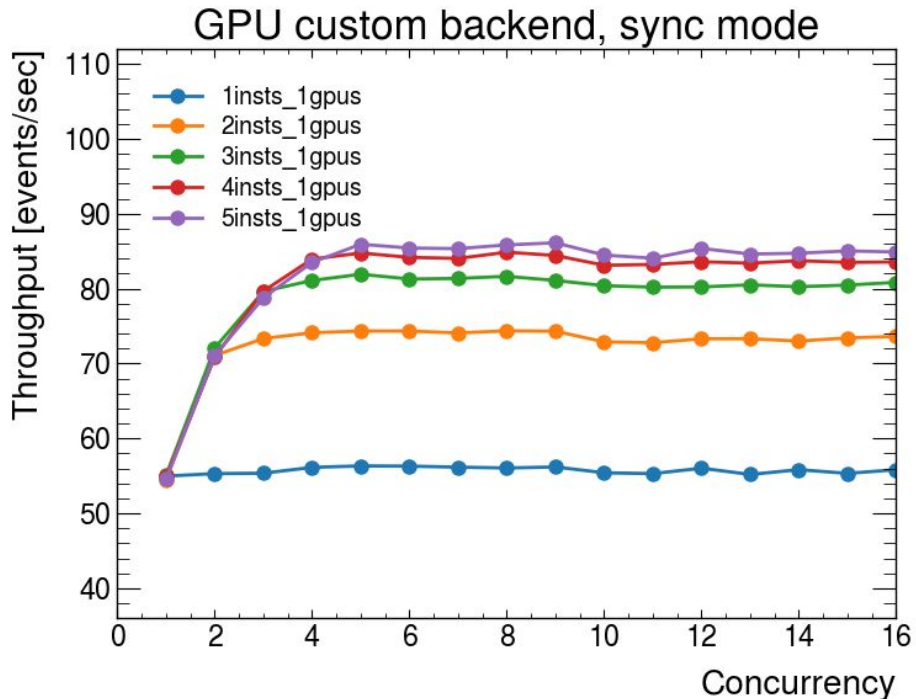
Roadmap

1. Standalone version of cpp
CsvMeasurementReader -> trackFitting
2. Object Oriented version
3. Triton backend implementation
 - a. Develop for I/O between the client and server

Backup

Multiple instances on 1 GPU and 2 GPUs

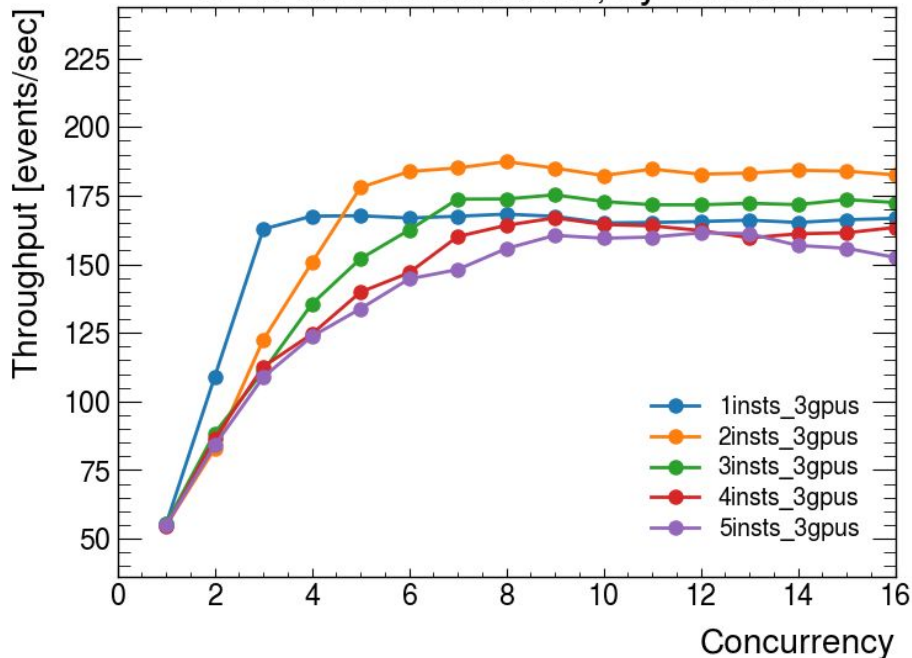
Left plot: 1 GPU max throughput is ~ 80; right plot: scale good with 1 instance



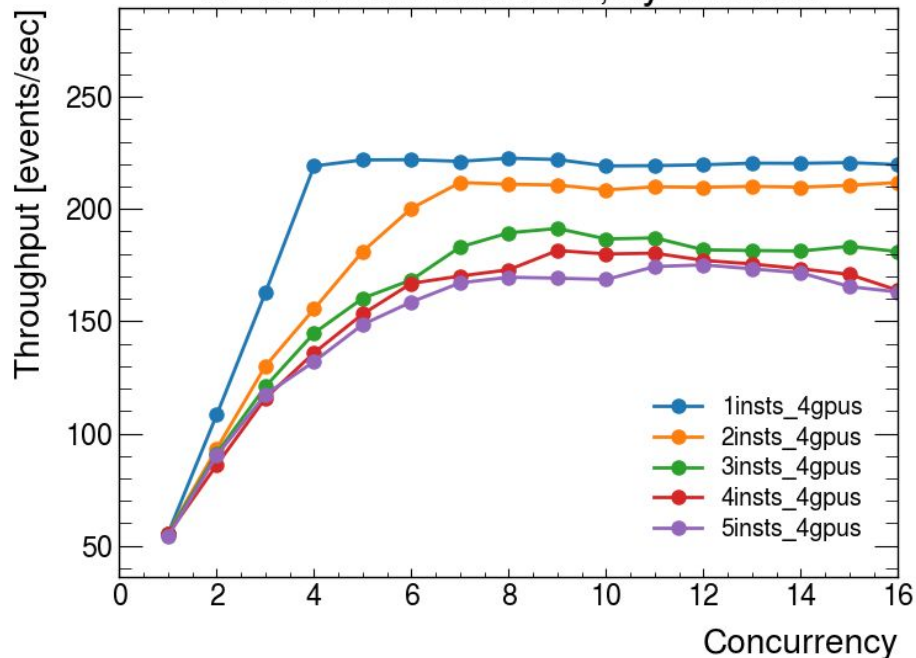
Multiple instances on 3 GPUs and 4 GPUs

Multi instances on multi-GPUs are not scaled properly

GPU custom backend, sync mode



GPU custom backend, sync mode



Discussions

1. I/O between the client and the server
2. Dataflow: Sequencer vs no sequencer option on the server side
 - a. Sequencer + python backend
 - i. Pro: least development effort
 - ii. Con: I/O time with storing csv files
 - b. No sequencer + custom backend
 - i. Pro: reduce the I/O
 - ii. Con: code refactoring

Acts-aaS

Our proposed workflow, measurements in, track out

