

Find us on github!

The A4 library

liba4.net



UNIVERSITY OF LIVERPOOL

Physics data processing with Google Protocol Buffers

Johannes Ebke^{LMU} & Peter Waller^{UL}

Why Google Protocol Buffers?

- **fast serialization** and deserialization of **structured data**
- separation of data structure definition and code
- simple, extendable message definition language
- code generators for many programming languages
- thread safe
- forward and backward binary compatibility
- definition available at runtime
- very **well tested** and widely used

Protocol Buffer definition

```
message Lepton {
  optional double pt = 1;
  optional double eta = 2;
  optional double phi = 3;
  optional int32 charge = 4;
}

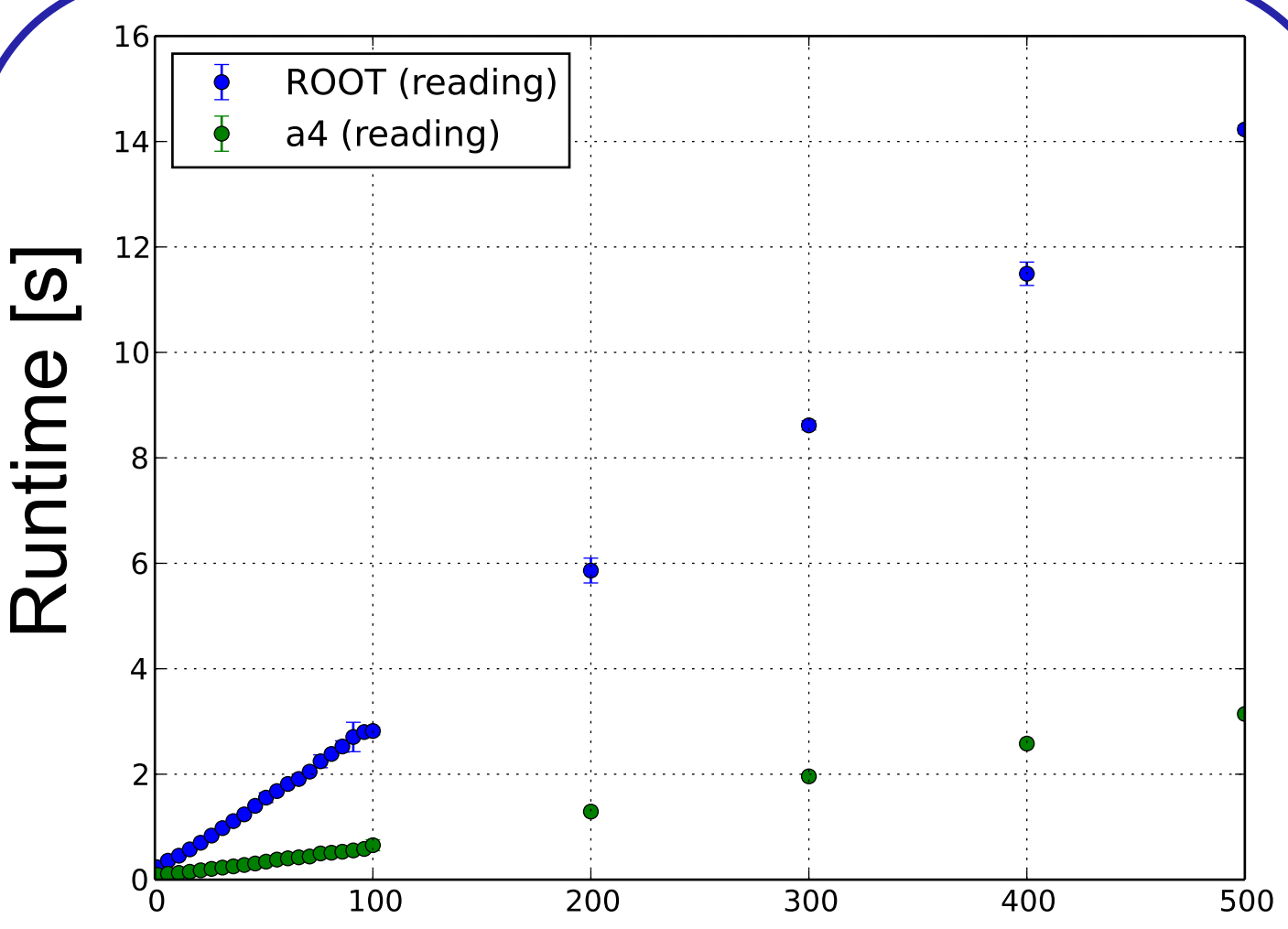
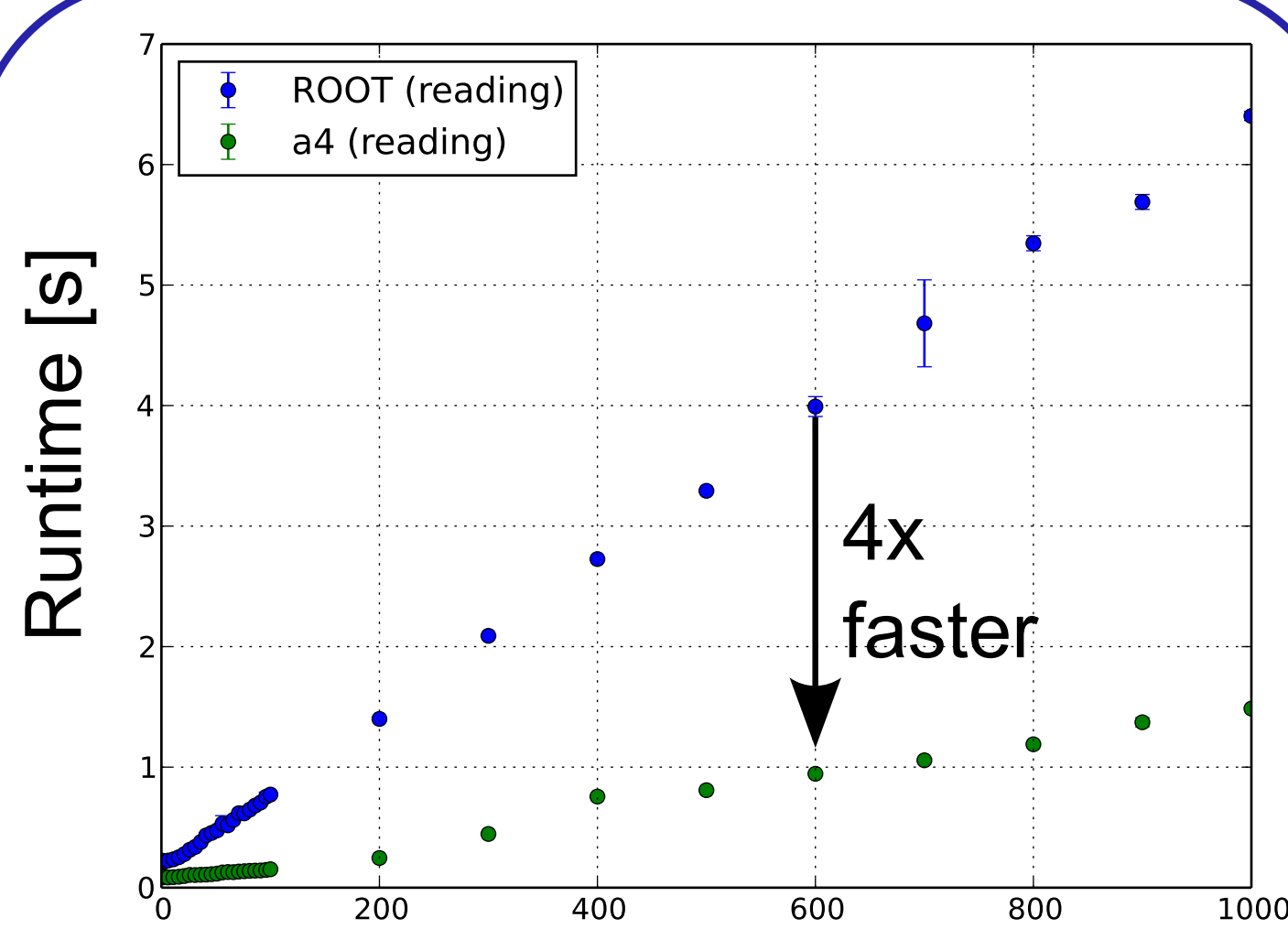
message PhysicsEvent {
  optional int32 run_number = 1;
  optional int32 event_number = 2;
  repeated Lepton electrons = 5;
  repeated Lepton muons = 6;
}
```

A4 use at the LHC

A4 is currently being used to process **~1 billion** ATLAS events per run on a local cluster. Results from the histograms, cutflows and systematic uncertainties evaluated using A4 have contributed to published ATLAS results.

Benchmarking A4

- flat ntuple-like synthetic events with random entries
- entries can be **numbers** or variable-length **arrays**
- processing = adding up all entries
- comparison: ROOT 5.32 TTree



ROOT conversion

A4 files with histograms or events can be converted to ROOT files, and ROOT TTrees can be converted to A4 files.

A4 File

Histograms are stored in A4 files together with the metadata

Histogram Store

The histogram store enables creating, initializing, filling and saving histograms on one short line of code:

```
S.T<H1>("h")(5,0,5).fill(n_jets);
```

The store object "S" can be a subdirectory, and also stores event weights.

All this is **very fast**, since no strings are involved.

Channels

The call `channel("mu")` in a processor triggers a re-run of the whole analysis on this event, but all histograms will have the prefix `"channel/mu/"`. Only in this run the `channel("mu")` call will return true.

Systematics

Similar to **channels**, the function `systematics("JES")` returns true in a special run, but this time "jes" has to be added on the command line. All histograms will also have the prefix `"systematic/JES"`.

MCViz:

Visualization of HEP Monte Carlo

The background image of this poster has been generated using **MCViz**, another project of the authors. It can convert HepMC, LHE and Pythia events into beautiful graphs in a variety of styles. It is available at mcviz.net

Benchmark Results:

- The A4 library is competitive in terms of speed, in many cases up to 4x faster than ROOT
- No significant nonlinear behavior has been observed
- Handling of large arrays in events could be improved

A4 Processor Example: Skim, slim and thin

```
class SkimSlimThin : public ProcessorOf<Event, MetaData> {
public: void process(const Event& event) {
  // Cut on at least two muons (skim)
  if (event.size_muons() < 2) return;
  // Remove some fields (slim)
  Event e = event; e.clear_tracks();
  foreach(auto& m, *e.mutable_muons()) {
    // Smear muons before writing them
    if (metadata().simulation()) smear(m);
    m.clear_id_hits(); // Remove hits from muons (thin)
  }
  write(e); // Write modified event to output file
}
```

Summary and Outlook

- The A4 library, while still experimental, can be used successfully for fast HEP data analysis
- Almost no manual bookkeeping is necessary due to the use and automatic handling of metadata
- Histograms and cutflows can be added quickly, and powerful tools to structure the analysis are available
- Use and collaboration is invited, visit us at liba4.net