



A quantitative review of data formats for HEP

Jakob Blomer, CERN

Many thanks to Philippe Canal and Axel Naumann!

ACAT 2017, Seattle

August 24th, 2017

Focus on data formats for the “last mile” of publication creation

- Data sets are relatively small
 - $\mathcal{O}(10^7)$ events
 - Tens to hundreds of properties per event
 - Volume: gigabytes to terabytes
- Little data lock-in: final data sets can be reproduced
e. g. starting from CMS MINIAOD, ATLAS Derived xAOD
- Processing tends to escape central computing workflows and resources
- Number of reads \gg number of writes
- I/O bound: little calculation

What we want from the data format

- 1 Fast turn-around, e. g.
"Let me quickly plot. . ."
 - Tuning cuts (partial reading of the data set)
 - Feeding into machine-learning framework (full reading of the data set)
- 2 Integration, unleashing the data
 - Libraries for C++, Python
 - Machine-learning frameworks (e. g. TMVA, TensorFlow)
 - Big Data schedulers (e. g. Spark)
 - Analytic tool kits (e. g. ROOT, R, SciPy)

- ROOT:**
 - Stream of serialized C++ objects in columnar layout
- Protobuf:**
 - Not a file format per se but (de)-serialization of small records
 - Possible file format: `schema || [size record-blob]*`
- SQLite:**
 - SQL database in a file
- HDF5:**
 - Hierarchical “data sets” that each contain a “data space” (n -dimensional array)
 - Popular in the HPC universe, integrated with MPI-I/O
 - There are HEP machine-learning data sets in HDF5
- Parquet:**
 - From the Apache Hadoop/Spark universe
 - Column-wise binary storage
- Avro:**
 - From the Apache Hadoop/Spark universe
 - JSON schema, row-wise binary storage
- Kudu:**
 - C++ native implementation of a format similar to Parquet
 - Not (yet?) available as library, only client-server APIs

	ROOT	PB	SQLite	HDF5	Parquet	Avro
Well-defined encoding	✓	✓	✓	✓	✓	✓
C/C++ Library	✓	✓	✓	✓	✓	✓
Self-describing	✓	⚡	✓	✓	✓	✓
Nested types	✓	✓	?	?	✓	✓
Columnar layout	✓	⚡	⚡	?	✓	⚡
Compression	✓	✓	⚡	?	✓	✓
Schema evolution	✓	⚡	✓	⚡	?	?

✓ = supported

⚡ = unsupported

? = difficult / unclear



Starting point: “What if I had my data set in format X?”

Example Analysis

- 8.5 million LHC run 1 events $B \rightarrow KKK$ [▶ Link](#)
- Flat n -tuple, 26 branches, mostly floating point numbers
- 21 branches needed for the toy analysis
- 2.4 million events can be skipped because one of the kaon candidates is flagged as a muon
- Toy analysis:
sum over all branches from non-cut Kaons

```
struct BDecay {  
    double h1_px;  
    double h2_px;  
    double h3_px;  
    double h1_py;  
    ...  
};
```

On the simple end of the spectrum,
helps to understand performance base case

Results

Outside managed storage systems (e. g. EOS), files are vulnerable to “bit rot”.
Without checksums, silently flipped bits can change physics results!

Test

- 500 k events, 60 MB to 100 MB
- Run toy analysis
- 100 trials with one random bit flip

Notes

- Usually: checksumming as a byproduct of compression
- ROOT reports errors but does not crash
- ROOT occasionally runs into infinite loop (**under investigation**)

Format	Undetected Errors
ROOT (uncompressed)	32 %
ROOT (zlib)	0 %
ROOT (lzma)	0 %
ROOT (lz4)	32 %
sqlite	32 %
HDF5 (column)	39 %
HDF5 (row)	37 %
protobuf (uncompressed)	34 %
protobuf (zlib)	0 %
parquet (uncompressed)	33 %
parquet (zlib)	0 %
avro (uncompressed)	34 %
avro (zlib)	23 %

Outside managed storage systems (e. g. EOS), files are vulnerable to “bit rot”.
Without checksums, silently flipped bits can change physics results!

Test

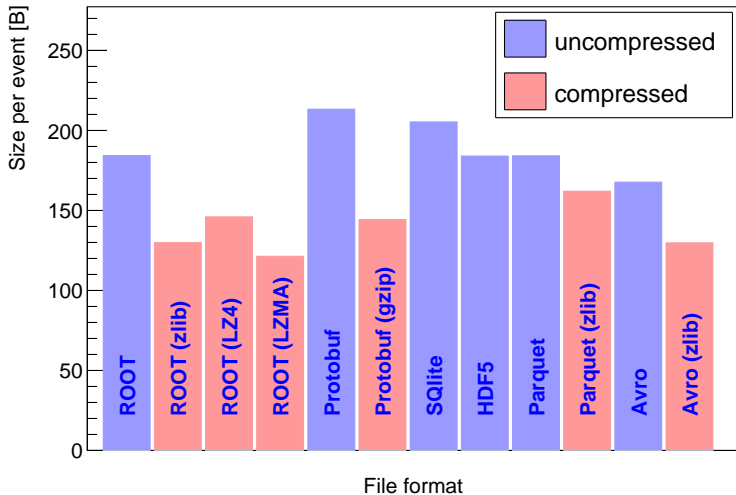
- 500 k events, 60 MB to 100 MB
- Run toy analysis
- 100 trials with one random bit flip

Notes

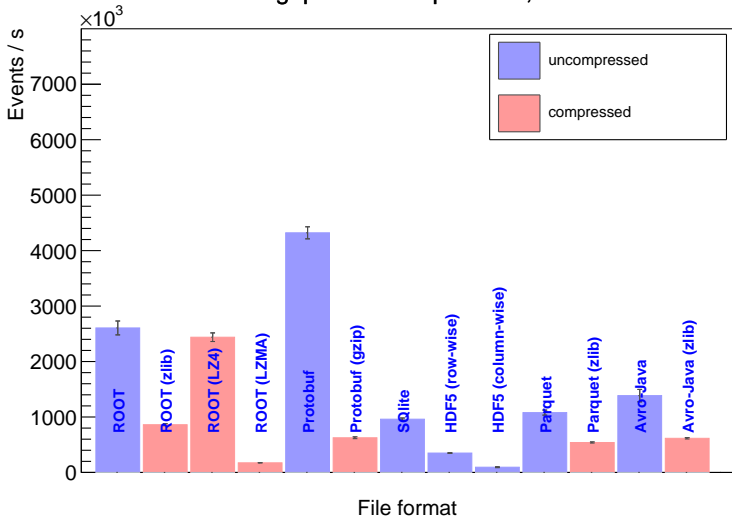
- Usually: checksumming as a byproduct of compression
- ROOT reports errors but does not crash
- ROOT occasionally runs into infinite loop (**under investigation**)

Format	Undetected Errors
ROOT (uncompressed)	32 %
ROOT (zlib)	0 %
ROOT (lzma)	0 %
ROOT (lz4) under investigation	32 %
sqlite	32 %
HDF5 (column)	39 %
HDF5 (row)	37 %
protobuf (uncompressed)	34 %
protobuf (zlib)	0 %
parquet (uncompressed)	33 %
parquet (zlib)	0 %
avro (uncompressed)	34 %
avro (zlib) checksums turned off	23 %

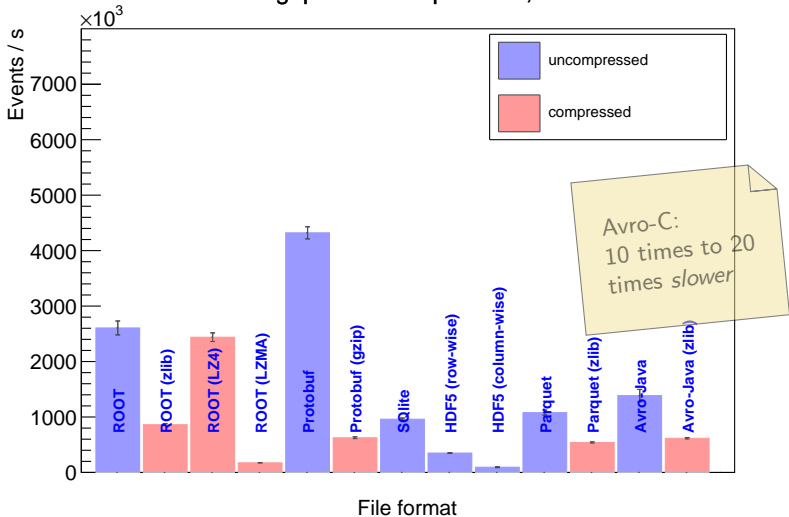
Data size LHCb OpenData



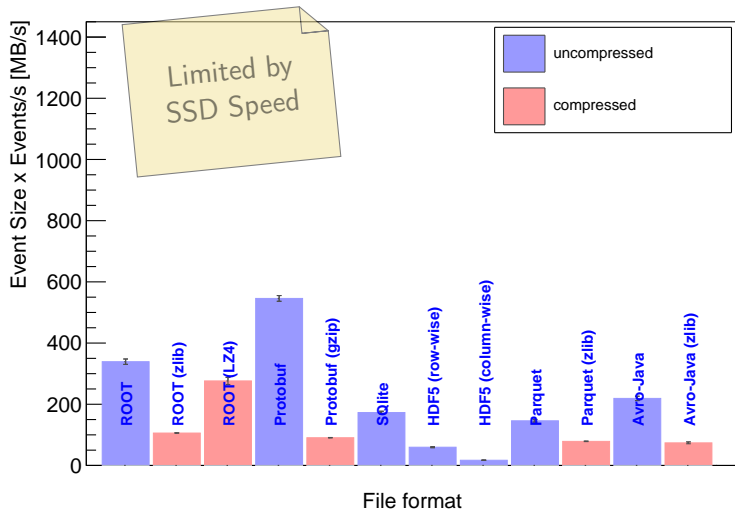
READ throughput LHCb OpenData, warm cache



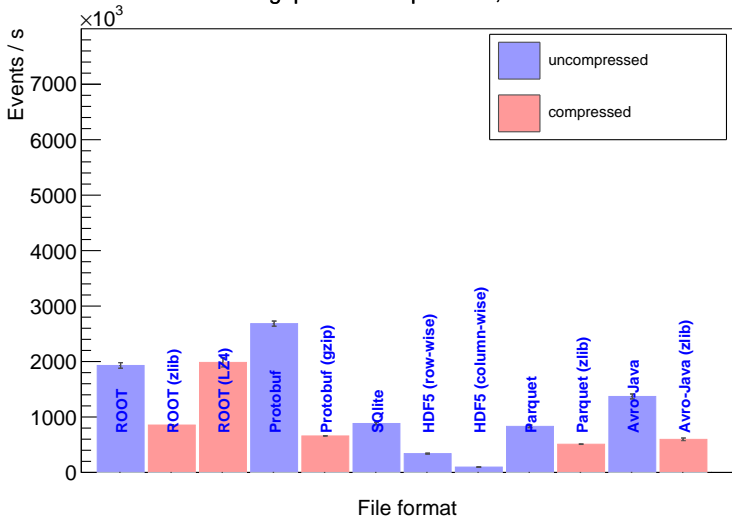
READ throughput LHCb OpenData, warm cache



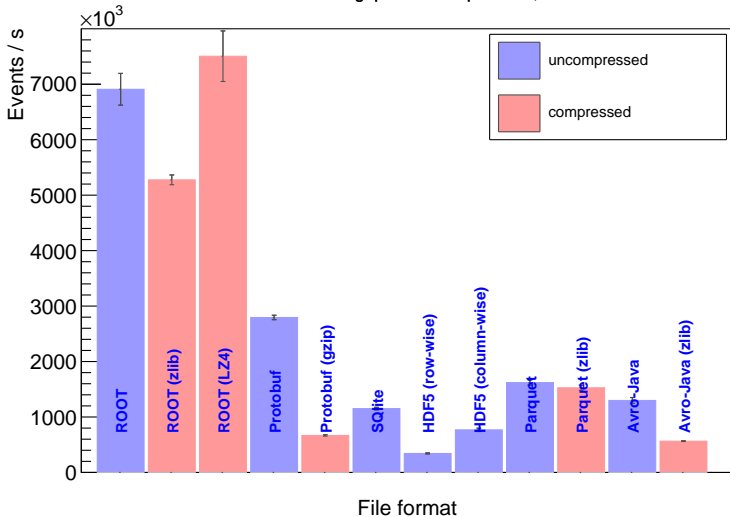
READ throughput LHCb OpenData, SSD cold cache



READ throughput LHCb OpenData, SSD cold cache

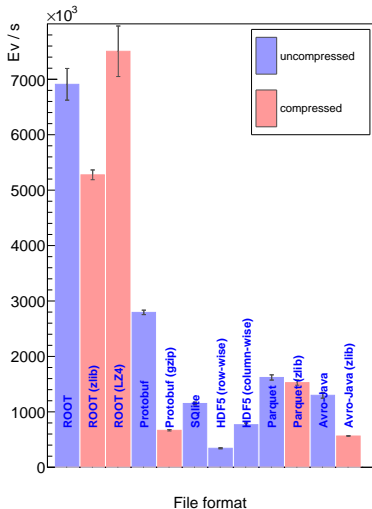


PLOT 2 VARIABLES throughput LHCb OpenData, SSD cold cache

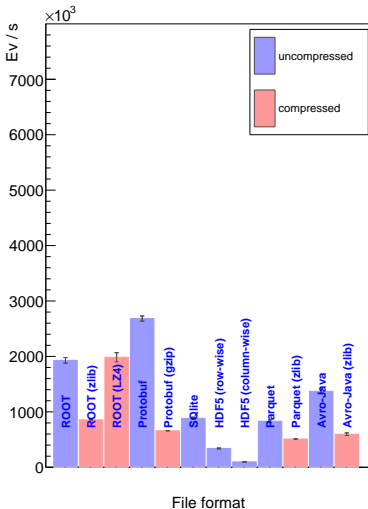


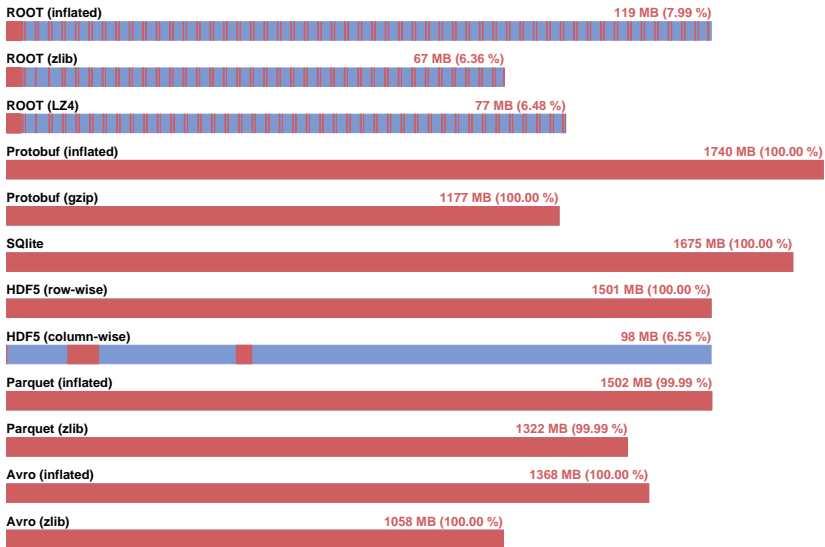
Comparison Reading vs. Plotting (SSD)

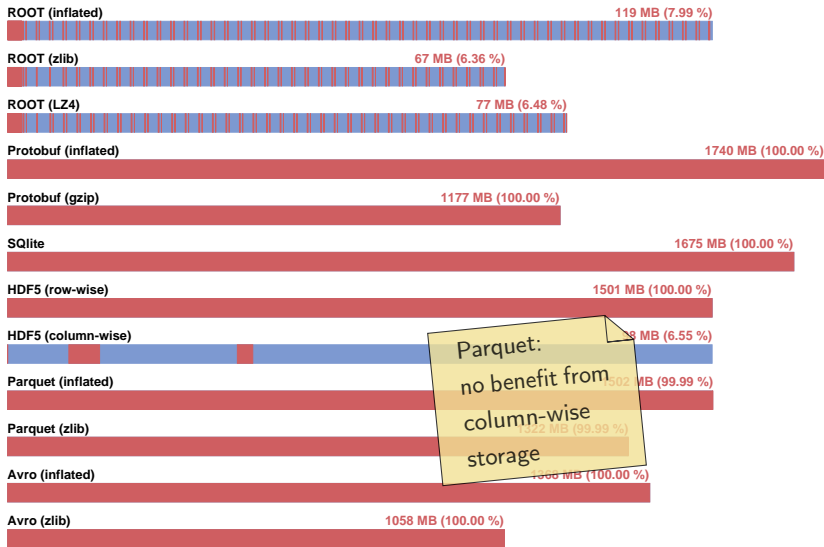
PLOT 2 VARIABLES throughput LHCb OpenData, SSD cold cache



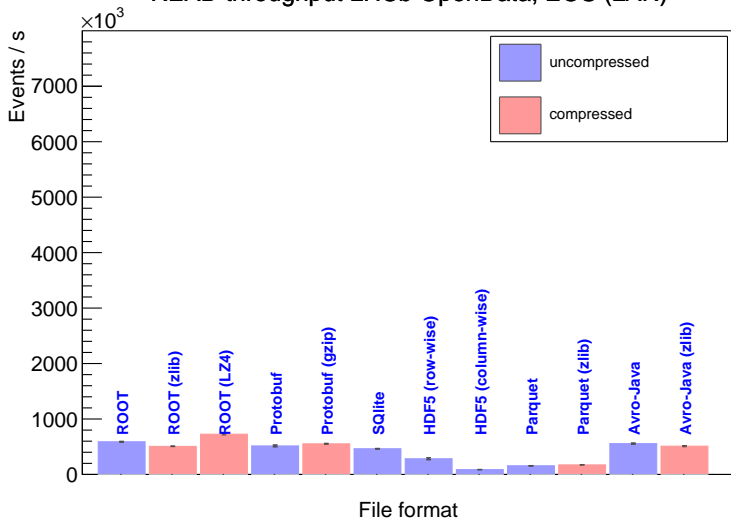
READ throughput LHCb OpenData, SSD cold cache





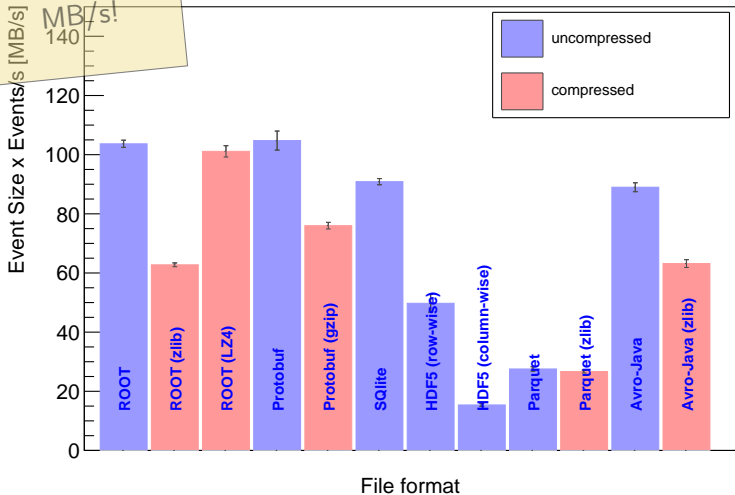


READ throughput LHCb OpenData, EOS (LAN)



Zoomed, MB/s!

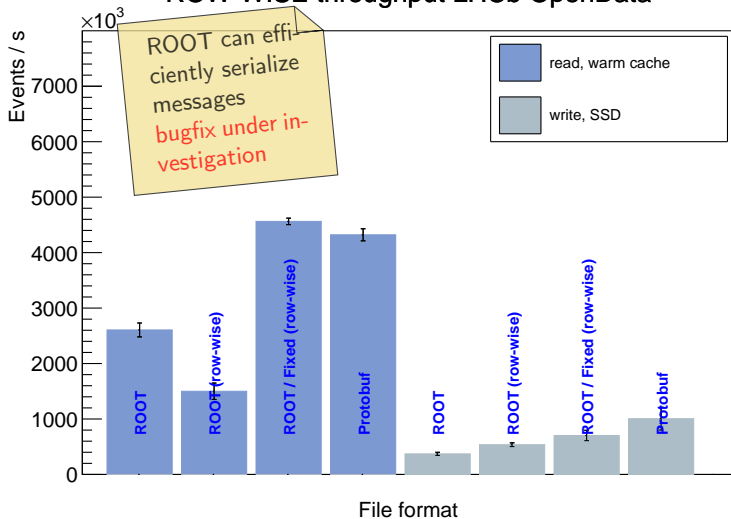
READ throughput LHCb OpenData, EOS (LAN)



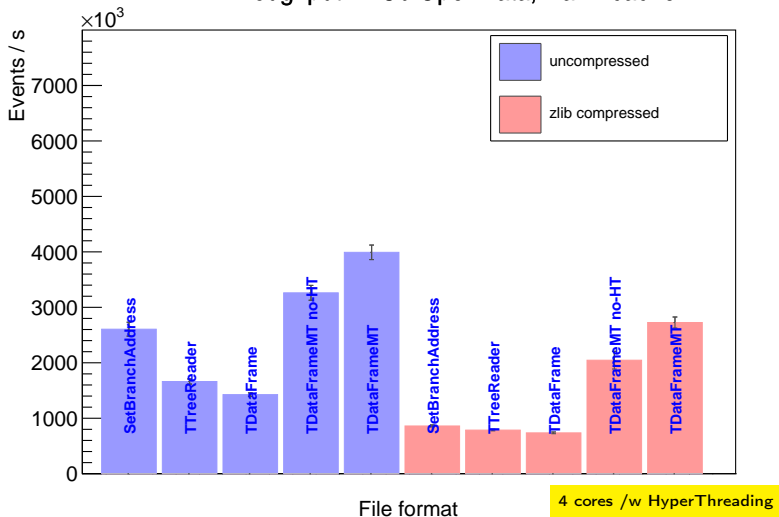
Looking into ROOT I/O

The following results are mostly relevant
on very fast storage or memory

ROW-WISE throughput LHCb OpenData

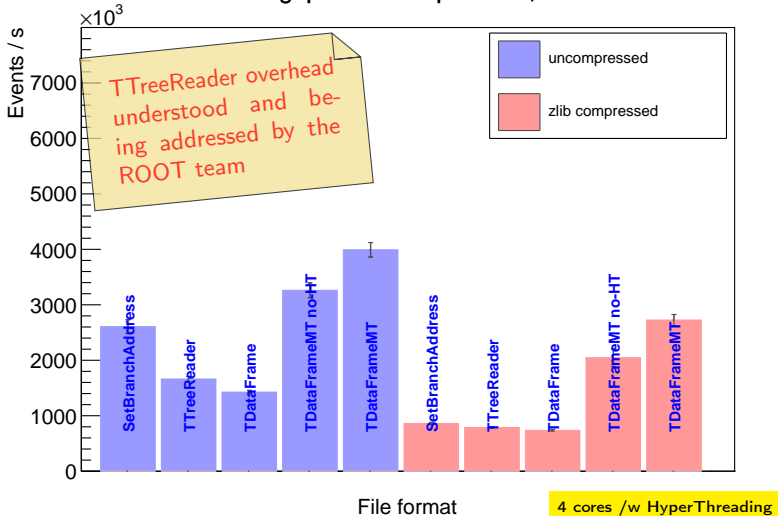


READ Throughput LHCb OpenData, warm cache



TDataFrame: see talk by G. Amadio in track 2

READ Throughput LHCb OpenData, warm cache



TDataFrame: see talk by G. Amadio in track 2

“Flat” Event Structure

```
struct FlatEvent {  
    double h1_px;  
    double h2_px;  
    double h3_px;  
    double h1_py;  
    ...  
};
```

“Deep” Event Structure

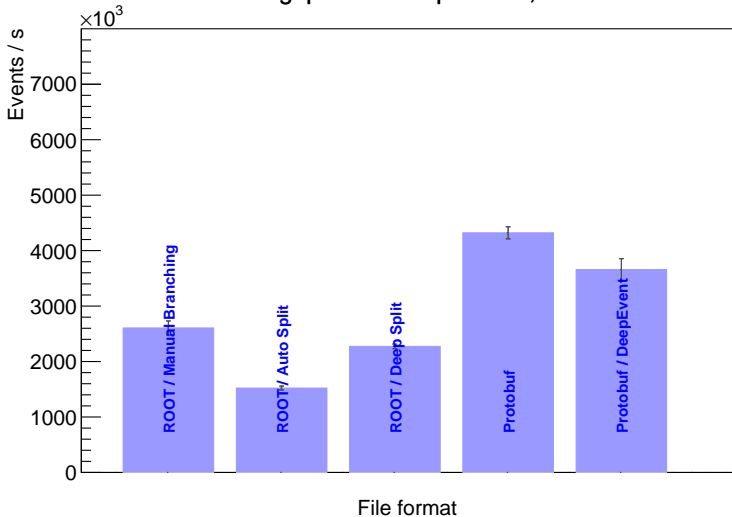
```
struct Kaon {  
    double h_px;  
    double h_py;  
    ...  
};  
struct DeepEvent {  
    std::vector<Kaon> kaons;  
};
```

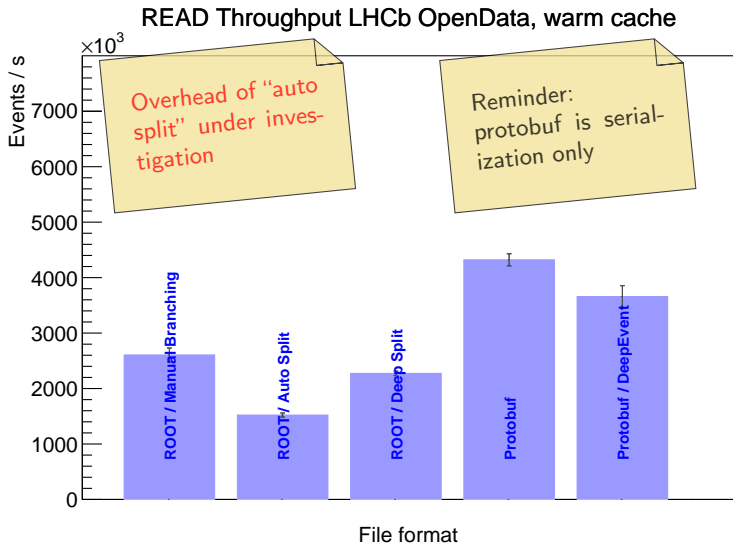
Three options considered

- 1 “Manual”: FlatEvent, branches created by hand
- 2 “Auto”: FlatEvent, branches created by ROOT with split level 99
- 3 “Deep”: DeepEvent, branches created by ROOT with split level 99

(For many branches, manual splitting becomes infeasible)

READ Throughput LHCb OpenData, warm cache





Reminder: reading 21 of 26 columns

ROOT / Manual Branching 1216 MB (80.98 %)



ROOT / Auto Split 1216 MB (80.98 %)



ROOT / Deep Split 1499 MB (81.98 %)



Protobuf 1740 MB (100.00 %)



Protobuf / DeepEvent 1795 MB (100.00 %)



Conclusion

Data Formats

- Large performance differences by different data formats and libraries
- Protobuf is a good benchmark for serialization performance
- HDF5 is most suitable for HPC-style processing of multi-dimensional arrays
- Parquet is an interesting format but C++ libraries are not yet ready

ROOT I/O

- In general best performance in these tests
- Several issues found & being addressed
- LZ4 turns out to be a good trade-off for analysis data sets
- Optionally checksummed ROOT files might be desirable

Many thanks to LHCb for the OpenData data set!

Backup

Hardware	Type
CPU	i7-6820HQ @ 2.7 GHz (4 cores + HT)
Memory	2×16 GB DDR4 2133 MHz
SSD (flash)	1 TB Toshiba XG3 PCIe
HDD (spinning)	Western Digital WD20NMVW

Library	Version
ROOT	6.10/04
protobuf	3.3.2
sqlite	3.19.3
hdf5	1.10.0_patch1
avro-c	1.8.2
avro-java	1.8.2
parquet-cpp	1.2.0

Operating system: Linux 4.12, glibc 2.25, gcc 7.1.1, EOS 4.1.26

Universally efficient data layout and encoding is challenging given large spectrum of storage performance characteristics

- Latency:
10 ns RAM \rightarrow 100 μ s 3D XPoint \rightarrow 1 ms SSD \rightarrow 10 ms HDD, LAN
between fastest and slowest: $\times 100\,000$
- Throughput:
20 GB/s RAM \rightarrow 2 GB/s SSD, 3D XPoint \rightarrow 100 MB/s HDD, LAN
between fastest and slowest: $\times 200$

Additionally: uneven scaling of components, e. g.
memory size vs. network bandwidth in laptops

Links

<https://github.com/lhcb/opendata-project>

<https://github.com/lhcb/opendata-project/blob/master/Background-Information-Notebooks/EventData.ipynb>

<https://github.com/jblomer/iotools/tree/acat17>