# LZ4, BulkIO, and offset removal performance

Jim Pivarski

Princeton University – DIANA

October 11, 2017

Three updates to ROOT I/O are aimed at speeding up or reducing file size for end-user analysis:

- ▶ new compression algorithm: LZ4 (speed)
- ▶ reading TBasket data directly into arrays: BulkIO (speed)
- ▶ removing offset data from TBranches that have a counter (size)

Three updates to ROOT I/O are aimed at speeding up or reducing file size for end-user analysis:

- ▶ new compression algorithm: LZ4 (speed)
- ▶ reading TBasket data directly into arrays: BulkIO (speed)
- ▶ removing offset data from TBranches that have a counter (size)

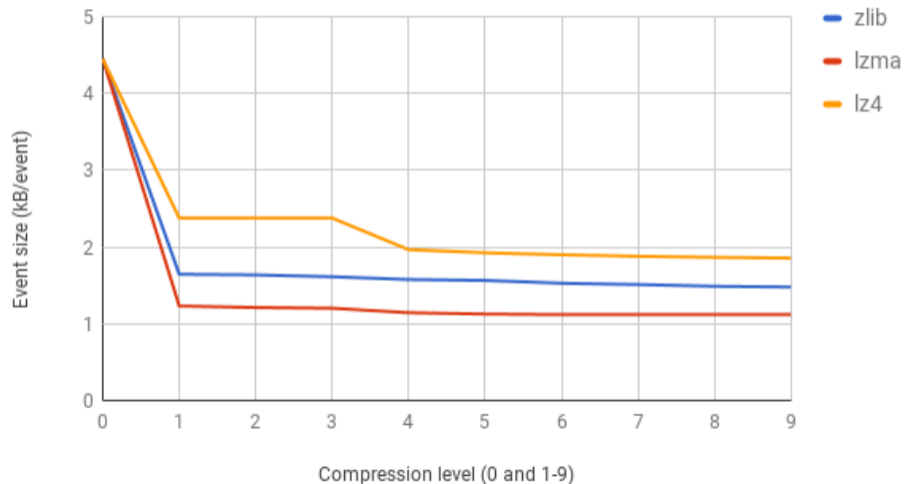Focus on CMS NanoAOD in particular because

- ▶ it is aimed at end-users (1–2 kB/event)
- ▶ it is broadly intended for 30–50% of analyses (not an individual user's ntuple)

Three updates to ROOT I/O are aimed at speeding up or reducing file size for end-user analysis:

- ▶ new compression algorithm: LZ4 (speed)
- ▶ reading TBasket data directly into arrays: BulkIO (speed)
- ▶ removing offset data from TBranches that have a counter (size)

Focus on CMS NanoAOD in particular because

- ▶ it is aimed at end-users (1–2 kB/event)
- ▶ it is broadly intended for 30–50% of analyses (not an individual user's ntuple)
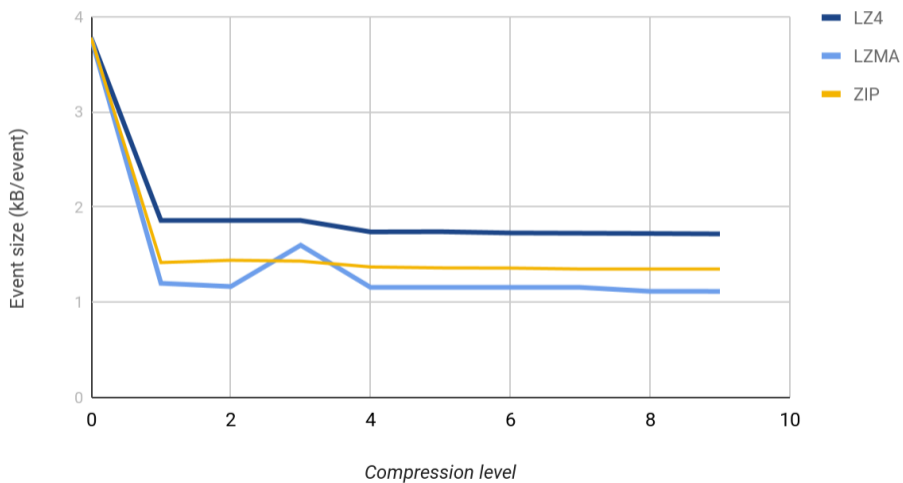
Also including studies of LHCb (thanks, Oksana!).

No ATLAS files because I can't generate new ones or `TTree::CopyTree` old ones.

- AWS instance with a fast SSD disk (i2.xlarge).
- No resource contention because I paid for exclusive access.
- "Writing" means a `TTree::CopyTree` with new TFile compression.
- "Reading" means filling a class made by MakeClass.
- "BulkIO" means filling arrays through `GetEntriesSerialized`.
- Always *reading* from warmed cache.
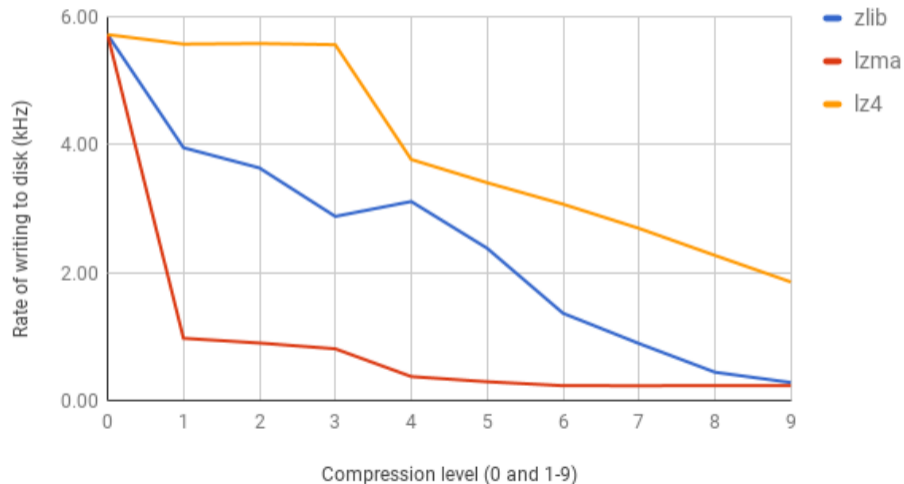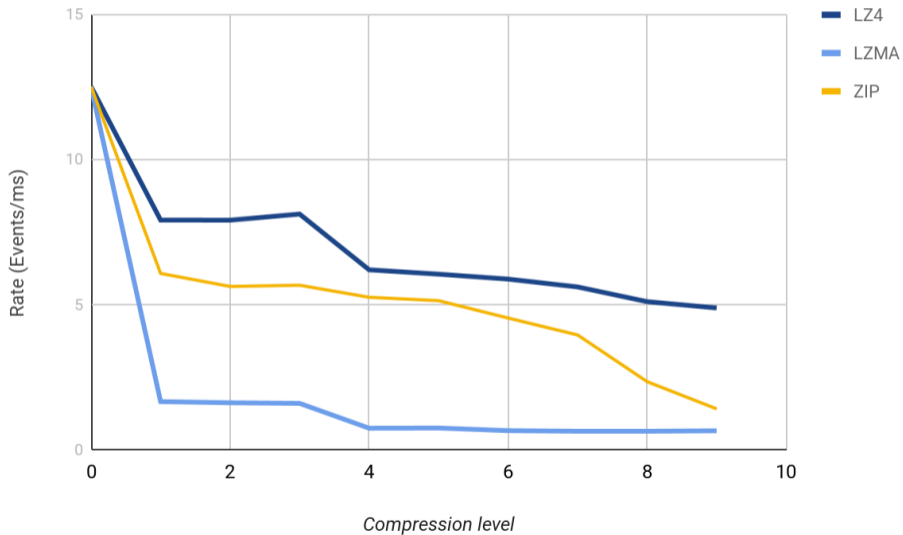- Five repeated trials; standard deviations are small compared to trends.

CMS NanoAOD

LHCB B2ppKK2011_md_noPIDstrip.root (22920 entries)

CMS NanoAOD

CMS NanoAOD

CMS NanoAOD

# write speed vs size
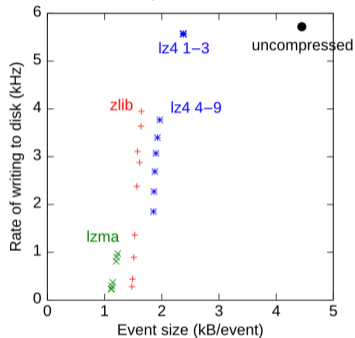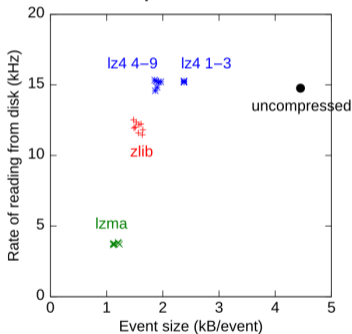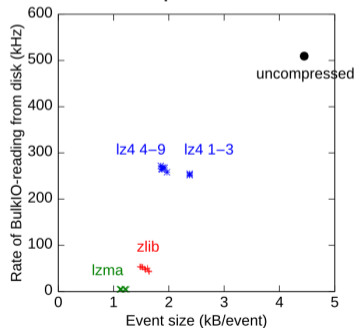


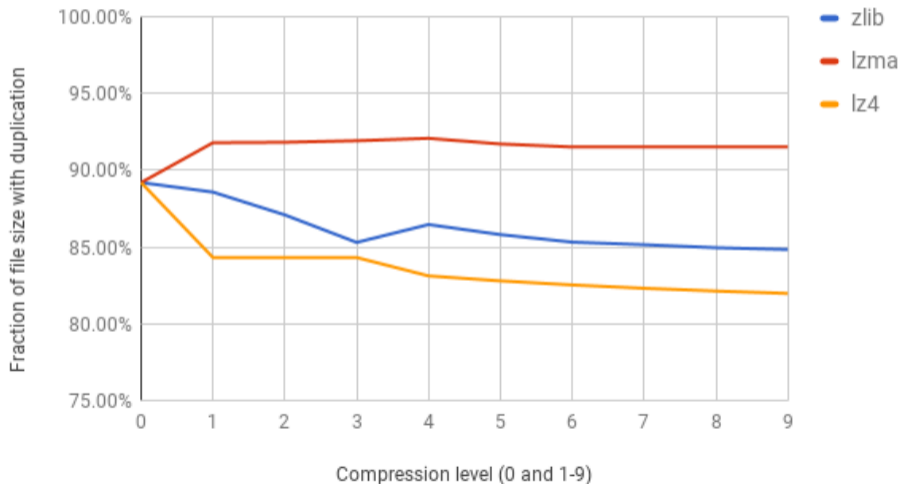# read speed vs size



# BulkIO speed vs size

⊛dianahep

TBranches for variable-sized data contain offsets indicating where each entry starts.

- ▶ This is unnecessary for branches with counters (e.g. `"Muon.pt[nMuons]/F"`).
- ▶ A fix is in progress (PR #1003) to optionally not write these offsets.
- ▶ May also write counts, instead of offsets, since repeated values might be more compressible.

> My study pre-dated (inspired) this PR; I constructed a copy
> of NanoAOD without offsets by putting all muon data into a
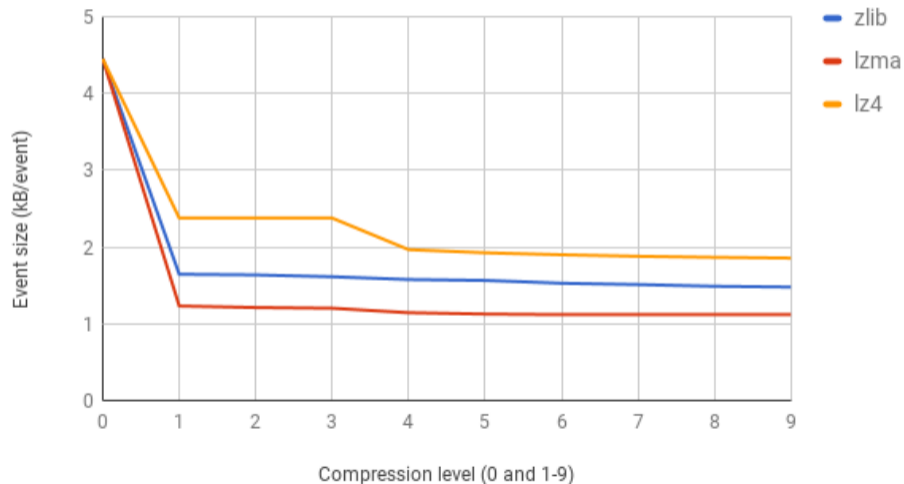> flat TTree, all jet data into a flat TTree, etc.

File size without duplication of particle counts

CMS NanoAOD

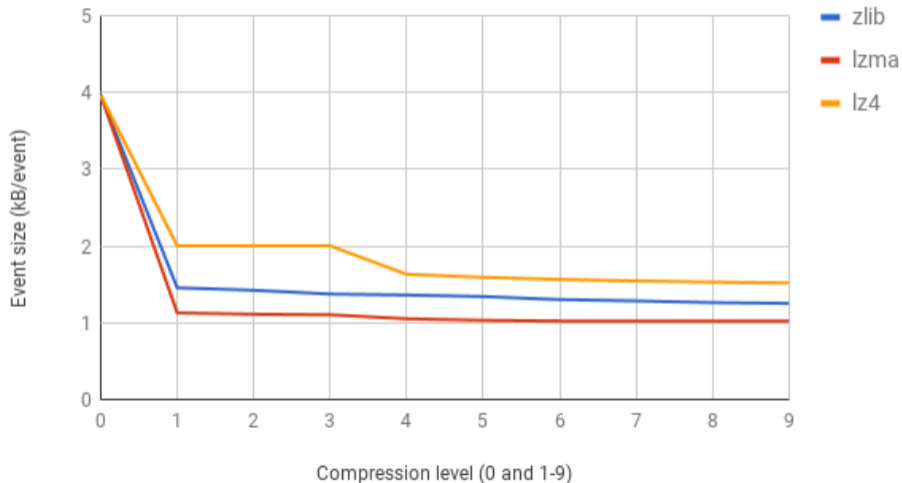CMS NanoAOD without particle count duplication

# Do offsets vs. counts matter? Yes for LZ4.

⬭dianahep

Synthetic test:

I generated
Poisson-random
counts and
integrated them
to make offsets,
then ZLIB and
LZ4 compressed
them.



Mean of Poisson random variable

⊖dianahep

LZ4 is as fast as uncompressed data for traditional GetEntry jobs.

BulkIO is an order of magnitude faster than GetEntry, especially with LZ4.

Unnecessary offsets add ∼10% to file size; may be removed.

Counts compress better than offsets, especially for LZ4.