

On the Scalability of Data Reduction Techniques in Current and Upcoming HPC Systems from an Application Perspective

**A. Huebl^{1,2}, R. Widera¹, F. Schmitt^{2,3}, A. Matthes^{1,2},
N. Podhorszki⁴, J.Y. Choi⁴, S. Klasky⁴ and M. Bussmann¹**

¹ Helmholtz-Zentrum Dresden - Rossendorf

² Technische Universität Dresden

³ NVIDIA ARC GmbH

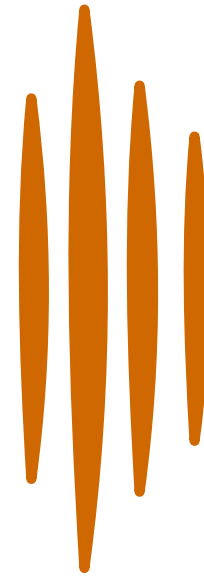
⁴ Oak Ridge National Laboratory



DRBSD-1 Workshop at ISC, Frankfurt *June 22th, 2017*



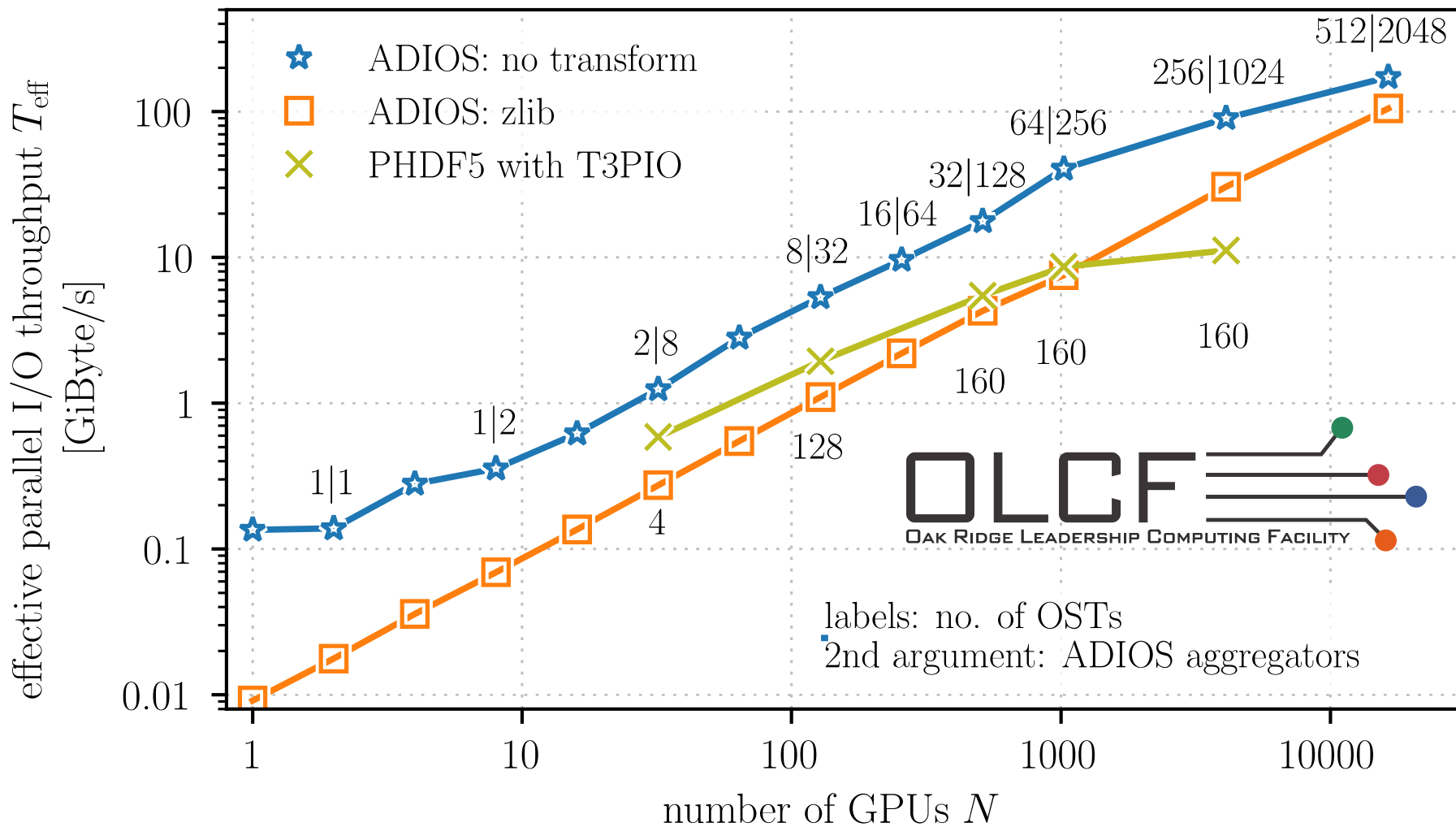
PICon GPU



- fully GPU-parallelized
- scales to full size Titan
- up to 60 GByte/s / node

Titan I/O Weak Scaling

$$T_{\text{eff}} \equiv \frac{N \times S}{t_{\text{I/O}}}$$



- Isn't data reduction always for free?
- Modeling of Parallel I/O
- Trading Compute Resources for I/O

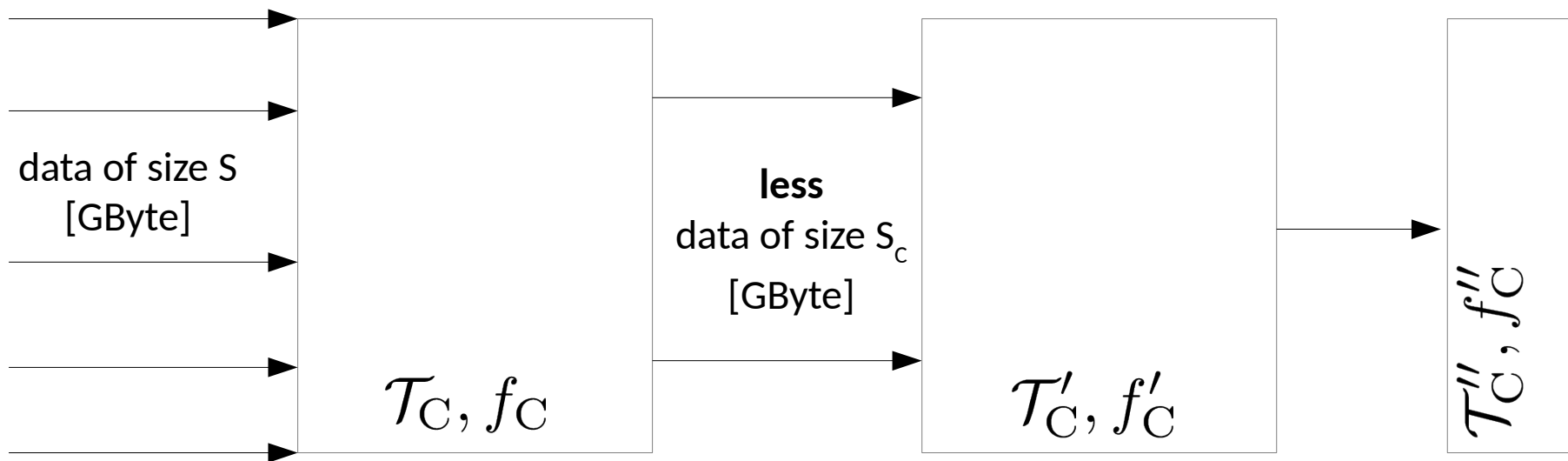


Characterization Of Data Reduction

Data Reduction

f_C reduction (compression) ratio S_c / S

\mathcal{T}_C reduction (compression) throughput T_C / T_{memcpy}

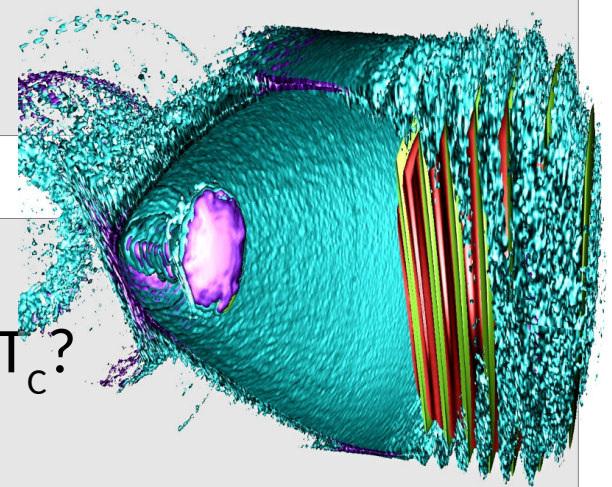


Data Reduction Examples

Binning of a spectrogram $f_c \rightarrow 0, T_c \sim T_{\text{memcpy}}$

cp to NVRAM $f_c = 1, T_c < T_{\text{memcpy}}$

Ray-cast or photo-realistic ray-trace $f_c \rightarrow 0, T_c ?$



Modeling Real-World Data Flows

Application View

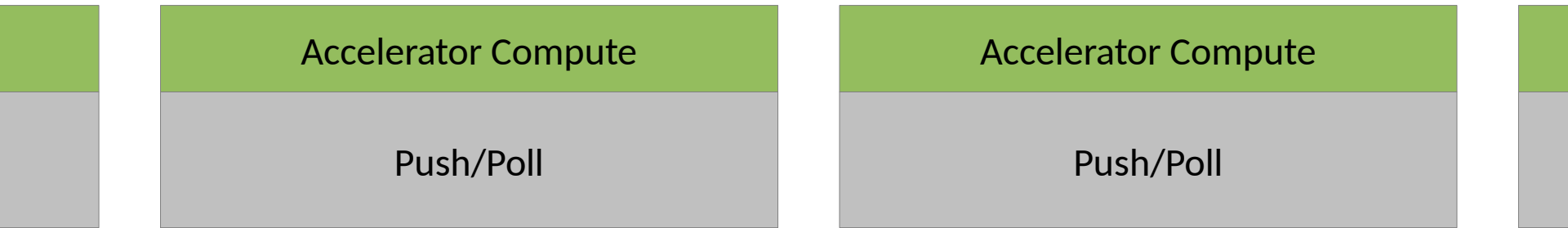
Compute

Compute

time



Application View

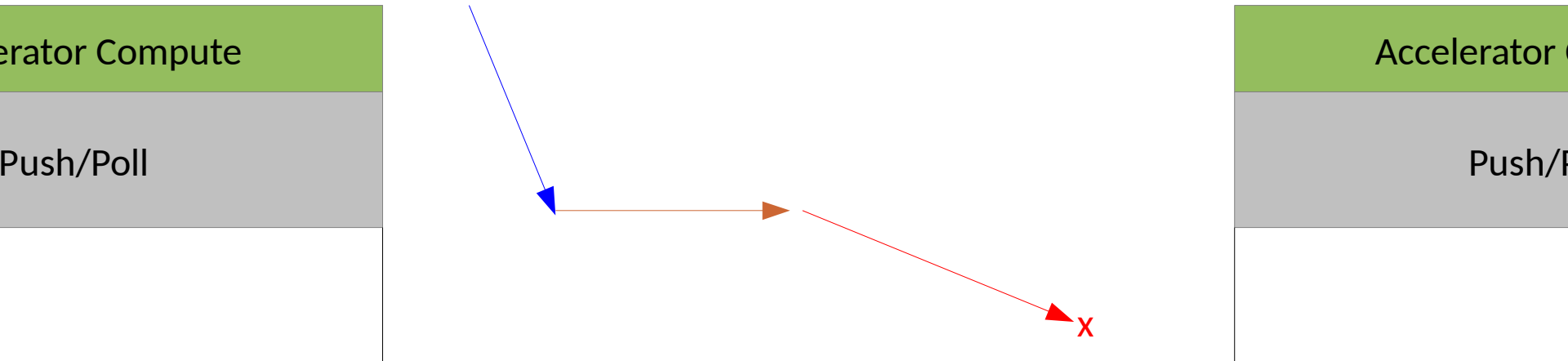


time



Application View

Unreduced, Synchronous I/O

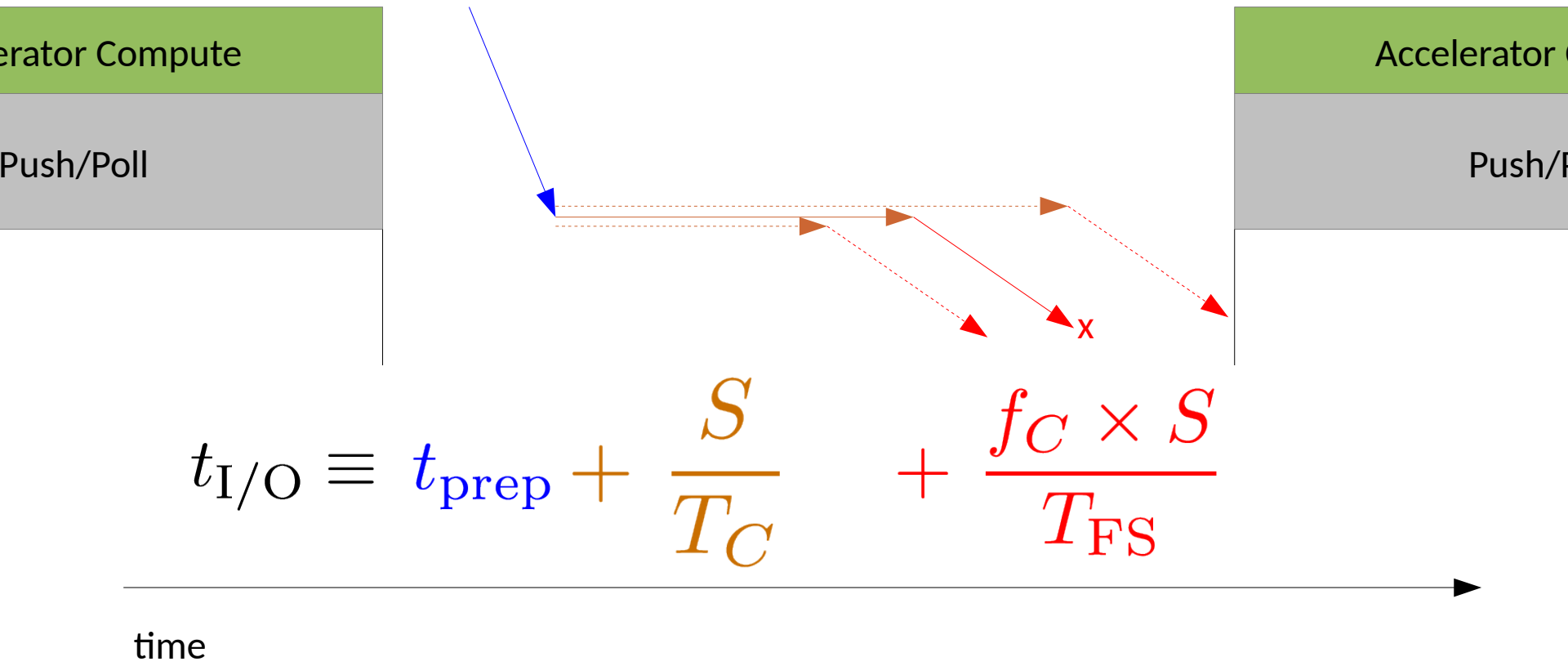


$$t_{I/O} \equiv t_{\text{prep}} + t_{\text{memcpy}} + t_{\text{off RAM}}$$



Application View

Reduced, Synchronous I/O



Comparison Synchronous I/O

Unreduced

$$t_{I/O} \equiv t_{\text{prep}} + \frac{S}{T_{\text{memcpy}}} + \frac{S}{T_{\text{FS}}}$$

> when?

Reduced

$$t_{I/O} \equiv t_{\text{prep}} + \frac{S}{T_C} + \frac{f_C \times S}{T_{\text{FS}}}$$



Break-Even Threshold

$$\frac{T_C \times (1 - f_C)}{1 - T_C} > T_{FS}$$

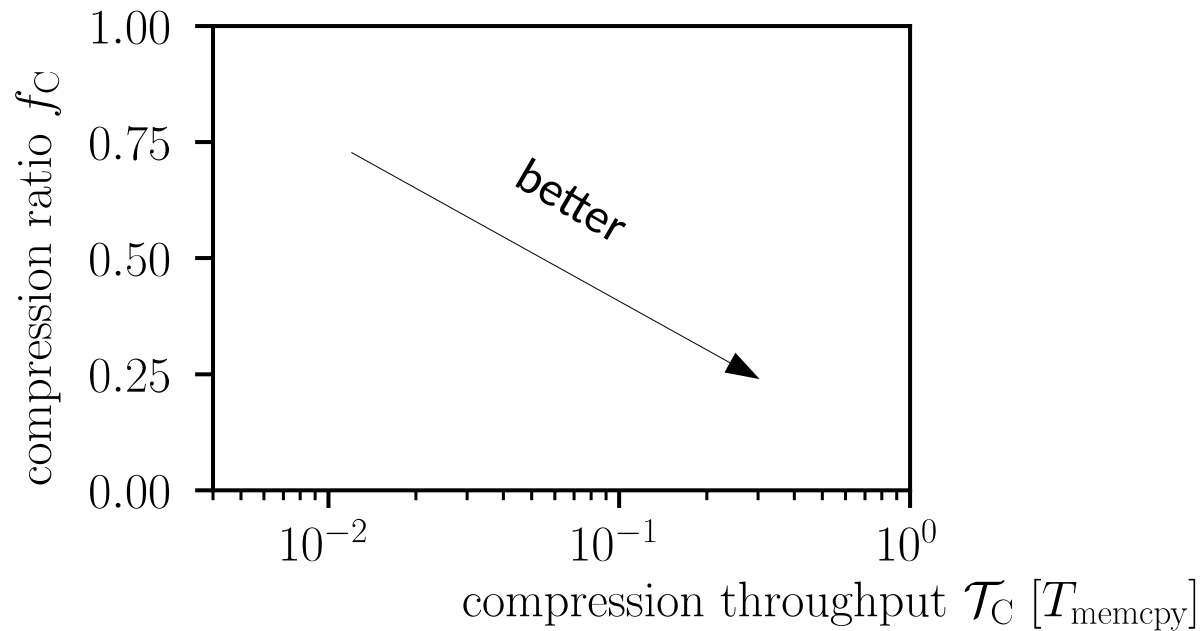
$$f_C \equiv \frac{S_C}{S} \quad T_C \equiv \frac{T_C}{T_{\text{memcpy}}} \quad T_{FS} \equiv \frac{T_{FS}}{T_{\text{memcpy}}}$$



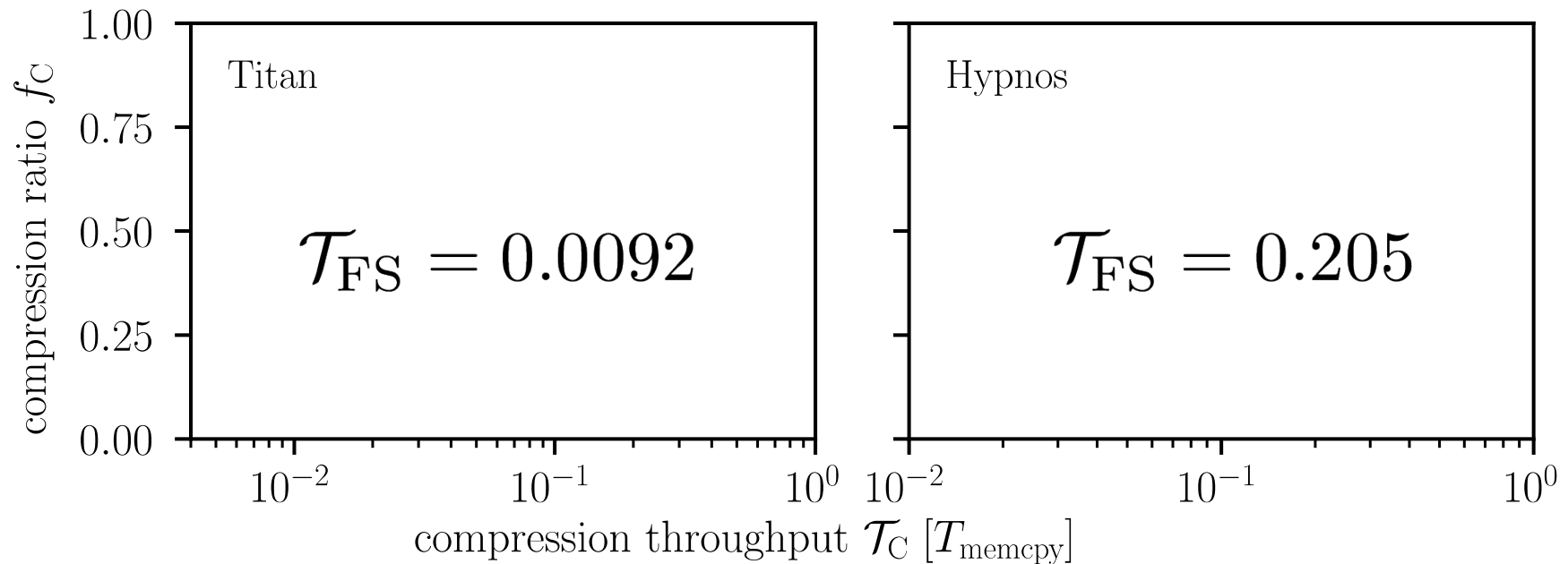
Compression

(as a data reduction example)

Ex-Situ Measurements for Compression Algorithms

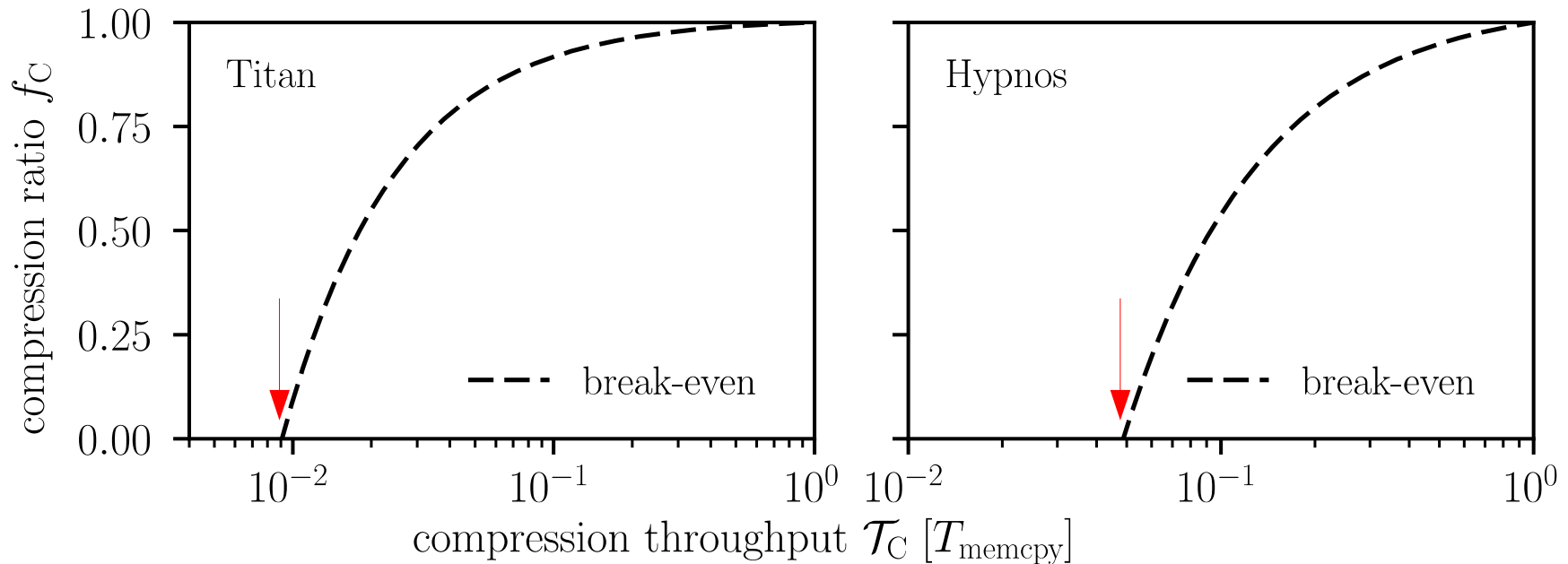


Ex-Situ Measurements for Compression Algorithms

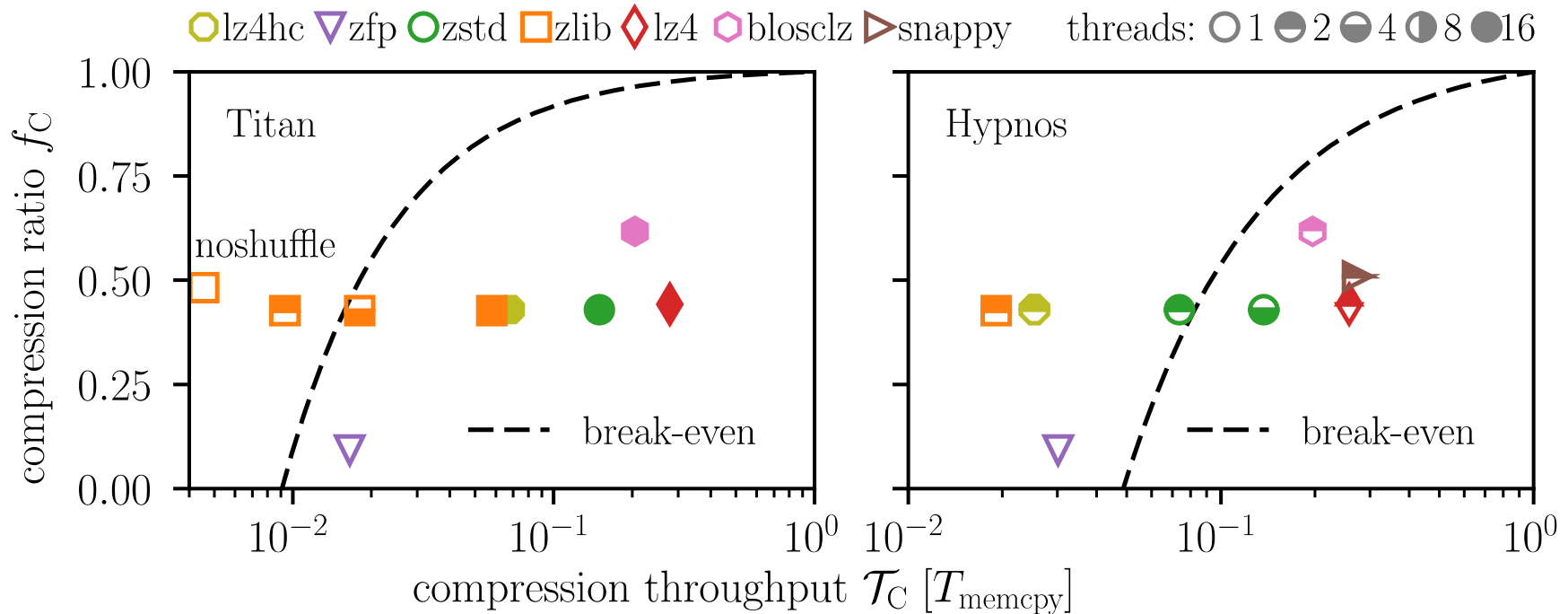


Ex-Situ Measurements for Compression Algorithms

$$\mathcal{T}_C > \frac{\mathcal{T}_{FS}}{\mathcal{T}_{FS} + 1}$$



Ex-Situ Measurements for Compression Algorithms



Zfp 0.5.1: three uncompressed bits / scalar; on *particle data*

Blosc 1.11.4-dev: add bitshuffle pre-conditioner



Impact on Applications

Predicting Feasible Reduction Algorithms to Reduce I/O Time

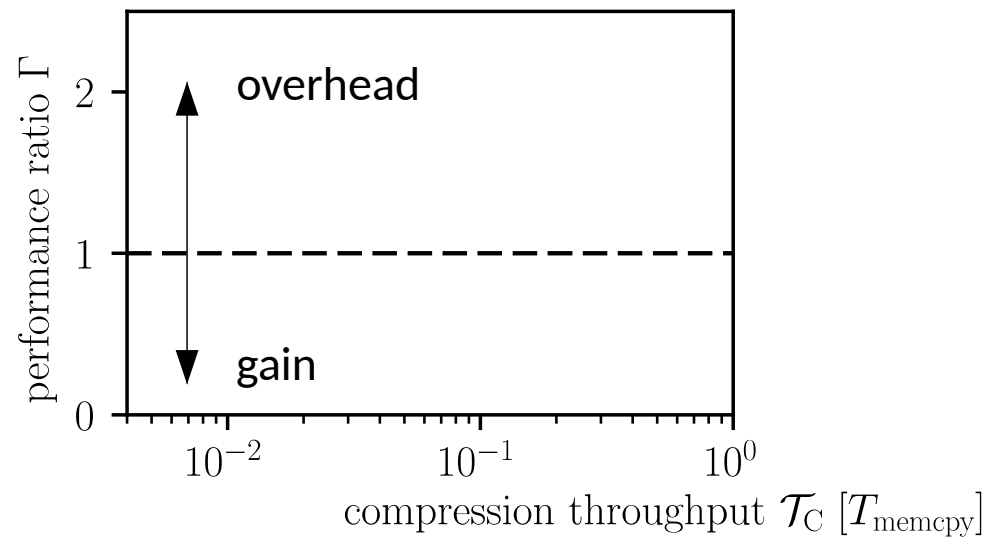
Summary of System Characteristics

f_C	reduction (compression) ratio	[0:1]
\mathcal{T}_C	reduction (compression) throughput	[T_{memcpy}]
\mathcal{T}_{FS}	write to next stage (e.g. FS)	[T_{memcpy}]
t_{prep}	preparation for the reduction stage	(seconds)



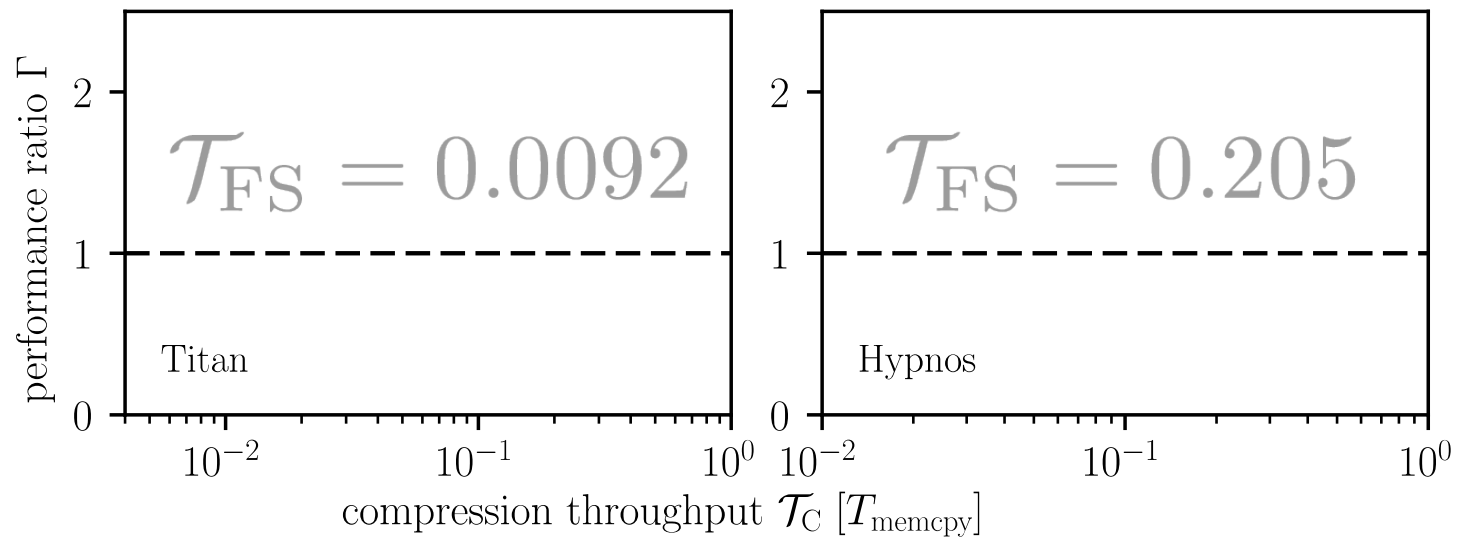
Application Performance: Gain or Loss?

$$\Gamma \equiv \frac{t_{I/O}^{\text{reduced}}}{t_{I/O}^{\text{unreduced}}}$$



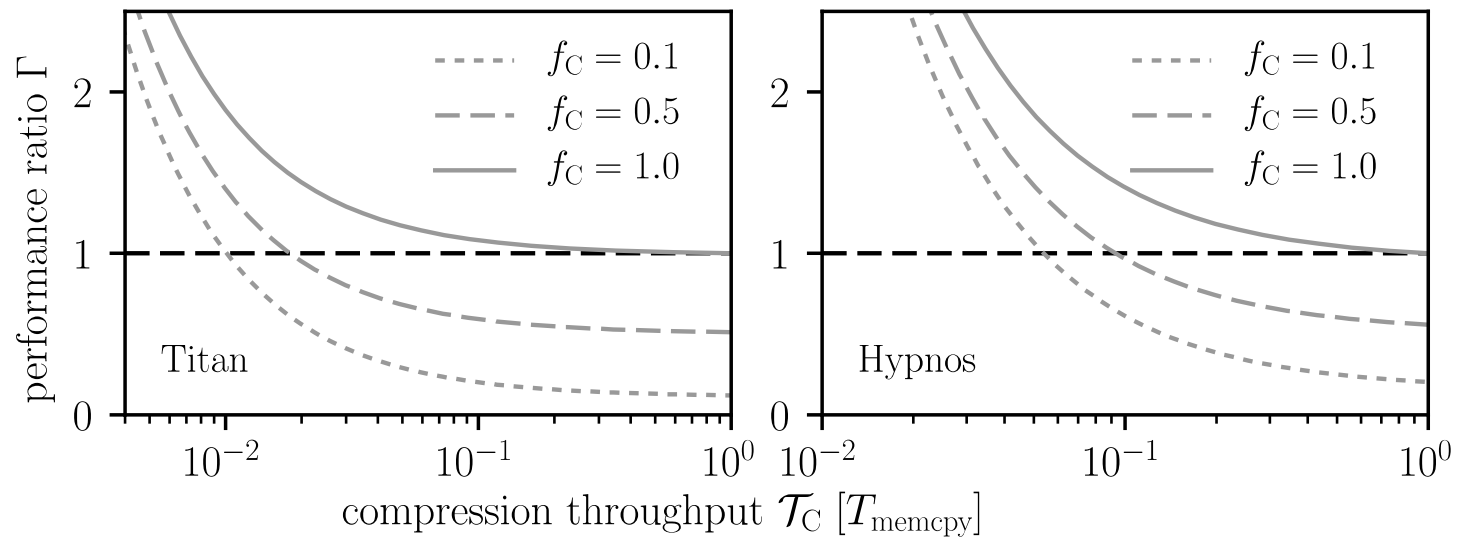
Application Performance: Gain or Loss?

$$\Gamma \equiv \frac{t_{I/O}^{\text{reduced}}}{t_{I/O}^{\text{unreduced}}}$$



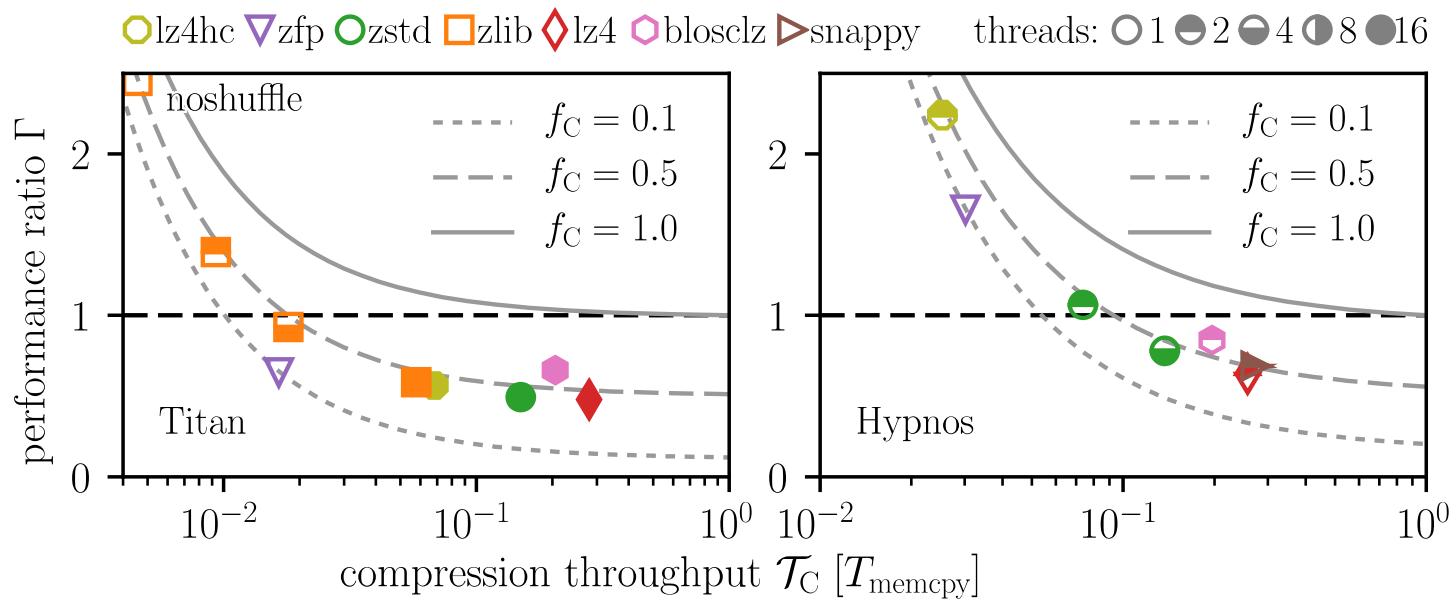
Application Performance: Gain or Loss?

$$\Gamma \equiv \frac{t_{I/O}^{\text{reduced}}}{t_{I/O}^{\text{unreduced}}}$$



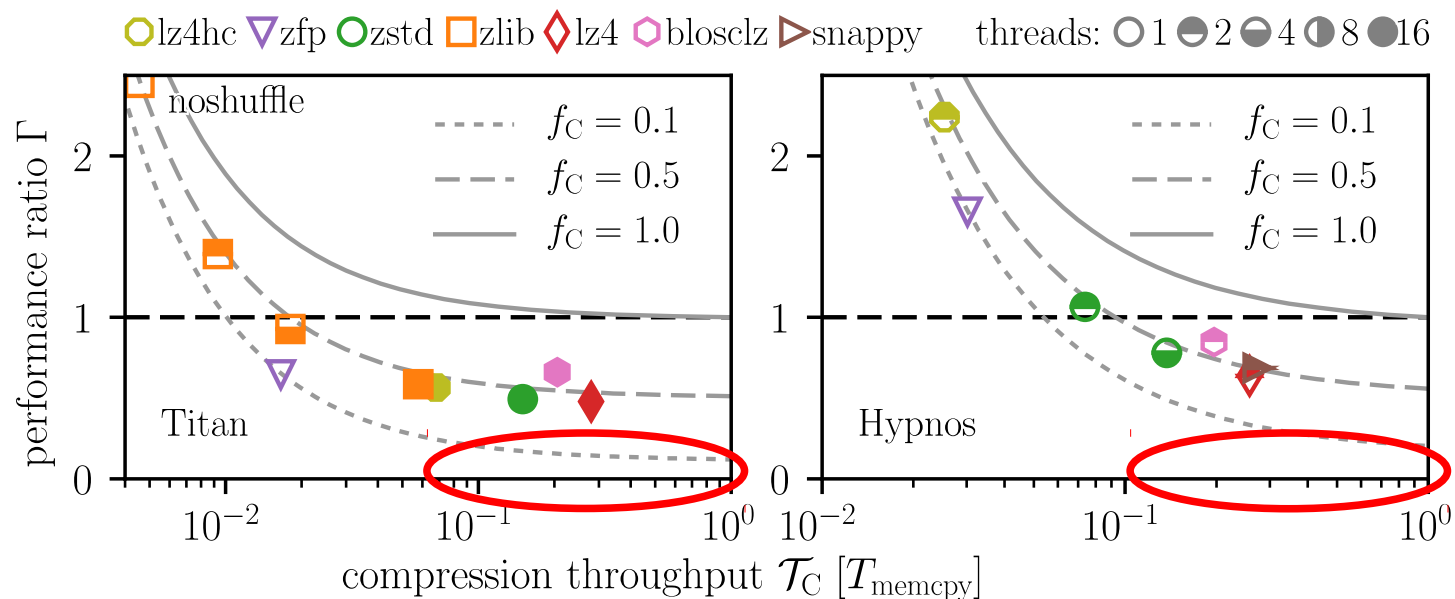
Application Performance: Gain or Loss?

$$\Gamma \equiv \frac{t_{I/O}^{\text{reduced}}}{t_{I/O}^{\text{unreduced}}}$$



Application Performance: Gain or Loss?

$$\Gamma \equiv \frac{t_{I/O}^{\text{reduced}}}{t_{I/O}^{\text{unreduced}}}$$



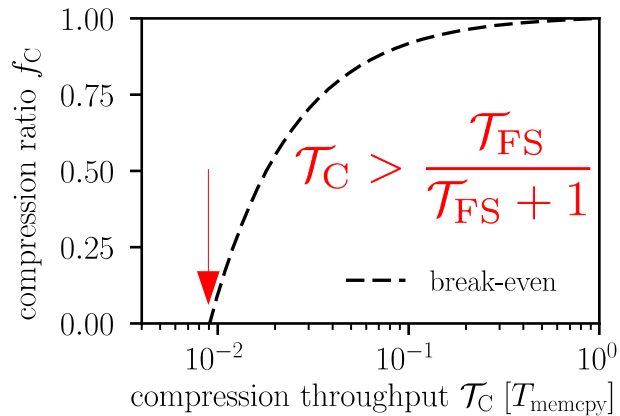
Enable application I/O on upcoming systems

via $\Gamma \ll 1$ ($\Gamma=1$ will be expensive)

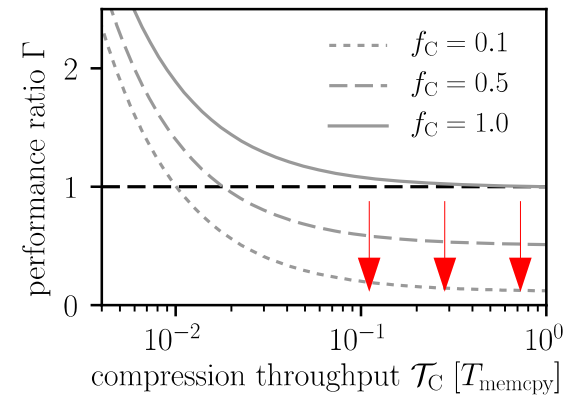


Consequences for Reduction Techniques

Throughput first



Aggressive Reduction next



Copy to burst buffers immediately “asynchronous I/O”
interleave computation with: t_{reduce} , $t_{\text{off RAM}}$
still: backlog on subsequent I/O

Users need easily programmable I/O stages

... & each stage bound by the above (else: backlog)



Even with data reduction,
there is no free lunch...

... unless you are *fast*

\mathcal{T}_C

... only then take a big bite

f_C

Thank you for your attention!

This research used resources of the Oak Ridge Leadership Computing Facility located in the Oak Ridge National Laboratory, which is supported by the Office of Science of the Department of Energy under Contract DE-AC05-00OR22725.

This project has received funding from the European Unions Horizon 2020 research and innovation programme under grant agreement No 654220.



Backup Slides

(Scenarios of Interest)

On-GPU Compression via Copy then → NIC

$$t_{I/O} \equiv t_{\text{prep}} + t_{\text{reduce}} + f_c \times t_{\text{off RAM}}$$

- - - = =

needs free global GPU RAM
(GPU apps are greedy, better have $f_c \rightarrow 0$)

Who compresses? User or I/O library?



Persistently In-GPU Compressed → direct to NIC

$$t_{I/O} \equiv \cancel{t_{\text{prep}}} + \cancel{t_{\text{reduce}}} + f_C \times t_{\text{off RAM}}$$

~ 0 ~ 0 $=$

needs algorithmic & data structure support

Every access costs:

- Flop/s
- registers

Examples exists: e.g. zfp



Titan vs. Summit Prediction

$$N_{\text{nodes}} T_{\text{FS}} : \times 2.5^{\text{faster}} \times 4^{\text{less nodes}}$$

$$S/\text{sec} : \times 10^{\text{faster}} \times 4^{\text{fatter}}$$

$$T_{\text{FS}} \approx \times 0.25$$

I/O : 4.0 more costly

$$t_{\text{I/O}} \equiv t_{\text{prep}} + \frac{S}{T_C} + \frac{f_C \times S}{T_{\text{FS}}}$$

== or +
==
x 4

- slow reductions hit break-even earlier
- but $t_{\text{I/O}}$ (starts $\Gamma \times 4$!) does *not* decrease linear with f_C alone



Titan vs. Hypnos

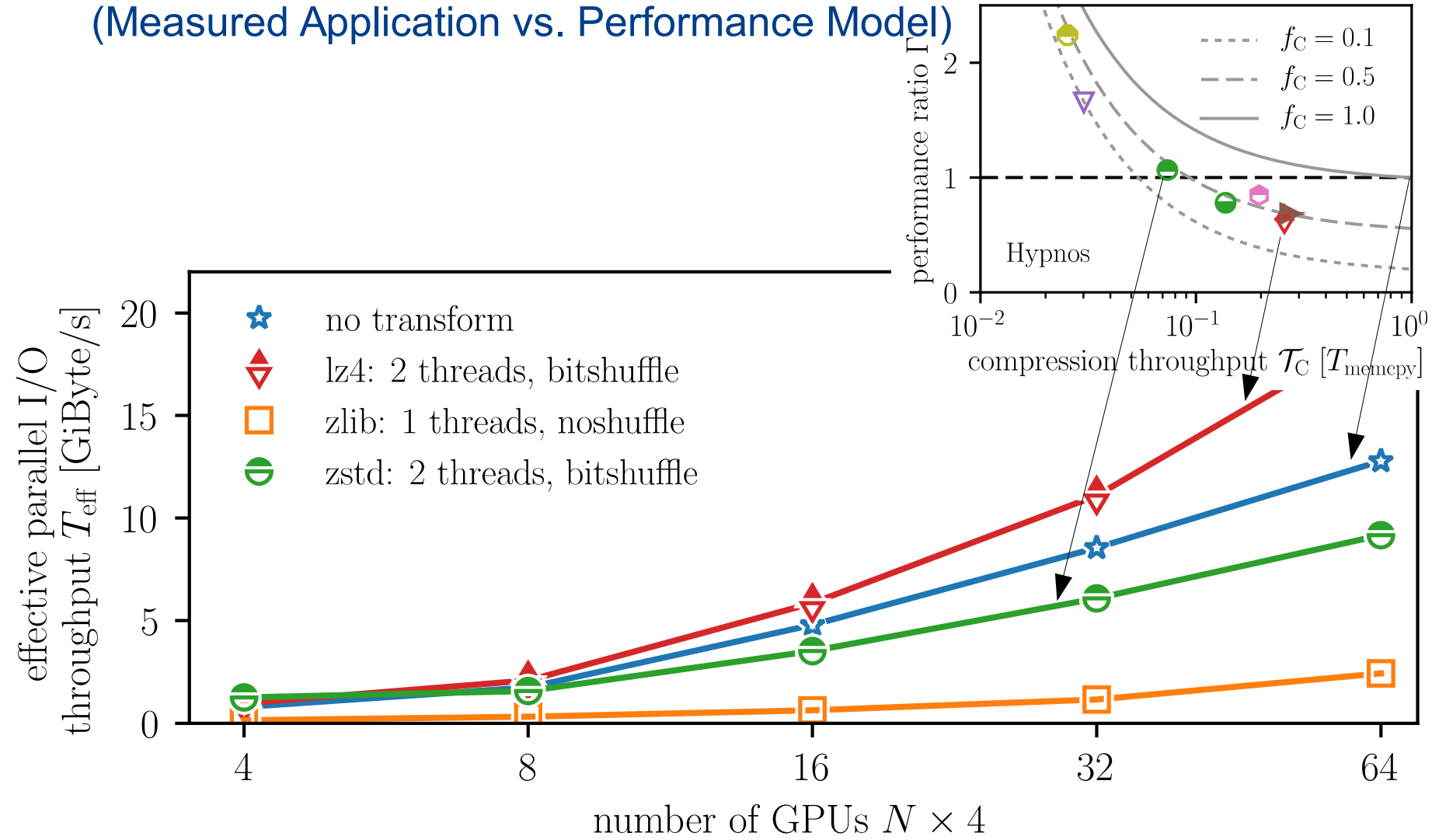
Table 1. PIconGPU I/O benchmark systems, both commissioned in 2012/ 13: relevant system characteristics and single node average file system throughput T_{FS} , defined as the design parallel bandwidth $B_{parallel}$ divided by N nodes

	Titan	Hypnos (queue: 'k20')
GPUs / node	1x K20x	4x K20m
CPUs / node	1x AMD Opteron 6274	2x Intel Xeon E5-2609
CPU-cores / GPU	16 (8FP)	2
GPU / CPU Flop/s (DP)	9.3:1	7.6:1
file system	Spider/Lustre	GPFS
$B_{parallel} = T_{FS} * N$ [GiByte/s]	1000	20
T_{FS} [GiByte/s]	0.055	1.25
CPU T_{memcpy} [GiByte/s]	6.0	6.1
maximum number of nodes N_{max}	18000	16



Hypnos I/O Weak Scaling

(Measured Application vs. Performance Model)

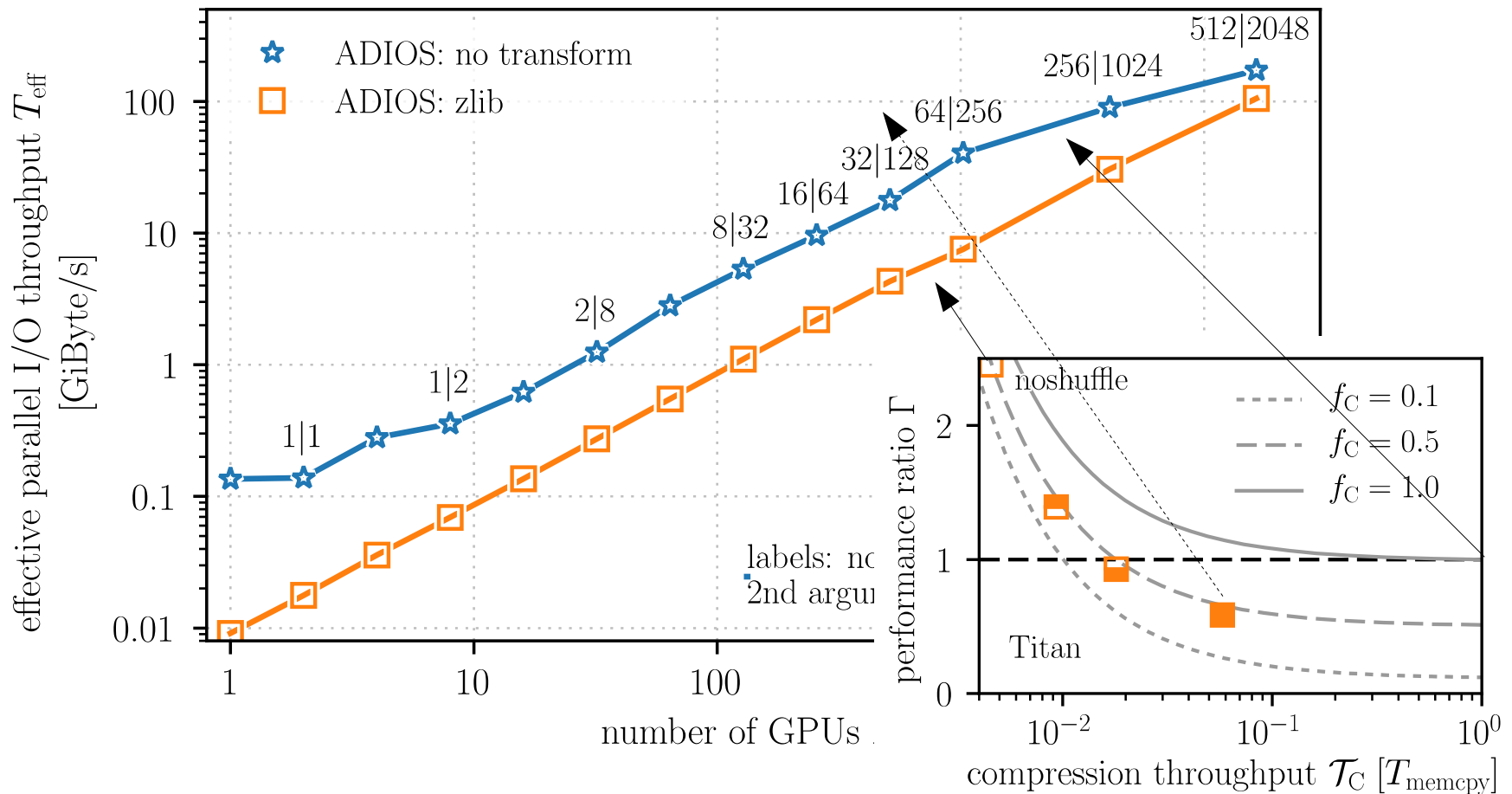


○ lz4hc
 ▽ zfp
 ○ zstd
 □ zlib
 ◇ lz4
 ◊ blosclz
 ▷ snappy
 threads:
 ○ 1
 ◐ 2
 ◑ 4
 ◒ 8
 ◓ 16



Titan I/O Weak Scaling

(Measured Application vs. Performance Model)



○ lz4hc
 ▽ zfp
 ○ zstd
 □ zlib
 ◇ lz4
 ◊ blosclz
 ▷ snappy
 threads: ○ 1 ◐ 2 ◑ 4 ◒ 8 ◓ 16



Titan I/O Weak Scaling

(Atlas Monitoring)

$$T_{\text{eff}} \equiv \frac{N \times S}{t_{\text{I/O}}} < B_{\text{parallel}}$$

