

# Datenreduktion

## Hochleistungs-Ein-/Ausgabe

Michael Kuhn

Wissenschaftliches Rechnen  
Fachbereich Informatik  
Universität Hamburg

2018-06-15



Universität Hamburg  
DER FORSCHUNG | DER LEHRE | DER BILDUNG

## 1 Datenreduktion

- Orientierung
- Motivation
- Speicherkostenmodell
- Wiederberechnung
- Deduplikation
- Kompression
- Erweiterte Kompression
- Zusammenfassung

## 2 Quellen



# Rechen- und Speicherleistung

- Kapazität und Leistung erhöhen sich weiterhin exponentiell
  - Komponenten verbessern sich unterschiedlich schnell
- E/A wird zunehmend zu einem Problem
  - Daten können immer schneller produziert werden
  - Speicherung ist nicht immer problemlos möglich
- Konsequenz: Höhere Ausgaben für Speicherhardware
  - In der Folge weniger verfügbare Mittel für Rechenleistung
  - Oder insgesamt teurere Systeme
- Speicherhardware ist ein maßgeblicher Teil der TCO
  - DKRZ: Ca. 20 % der Gesamtkosten
  - Entspricht 6.000.000 € Anschaffungskosten



## Beispiel: DKRZ

	2009	2015	Faktor
Leistung	150 TF/s	3 PF/s	20x
Knotenzahl	264	2.500	9,5x
Knotenleistung	0,6 TF/s	1,2 TF/s	2x
Arbeitsspeicher	20 TB	170 TB	8,5x
Speicherkapazität	5,6 PB	45 PB	<b>8x</b>
Speicherdurchsatz	30 GB/s	400 GB/s	13,3x
Festplatten	7.200	8.500	1,2x
Archivkapazität	53 PB	335 PB	<b>6,3x</b>
Archivdurchsatz	9,6 GB/s	21 GB/s	<b>2,2x</b>
Stromverbrauch	1,6 MW	1,4 MW	0,9x
Beschaffungskosten	30 M€	30 M€	1x

## Beispiel: DKRZ...

	2020	2025	Exascale (2020)
Leistung	60 PF/s	1,2 EF/s	1 EF/s
Knotenanzahl	12.500	31.250	100k–1M
Knotenleistung	4,8 TF/s	38,4 TF/s	1–15 TF/s
Arbeitsspeicher	1,5 PB	12,8 PB	3,6–300 PB
Speicherkapazität	270 PB	1,6 EB	0,15–18 EB
Speicherdurchsatz	2,5 TB/s	15 TB/s	20–300 TB/s
Festplatten	10.000	12.000	100k–1M
Archivkapazität	1,3 EB	5,4 EB	7,2–600 EB
Archivdurchsatz	57 GB/s	128 GB/s	—
Stromverbrauch	1,4 MW	1,4 MW	20–70 MW
Beschaffungskosten	30 M€	30 M€	200 M\$



























# Überblick

- Daten werden in Blöcke aufgeteilt (4–16 KB)
  - Unterschiedliche Aufteilungsmethoden
- Jeder einzigartige Datenblock wird nur einmal gespeichert
  - Mehrfach vorkommende Blöcke werden referenziert
- Bisherige Studien zeigen 20–30 % Einsparung für HPC-Daten
  - Insgesamt mehr als 1 PB untersucht
  - Deduplikation ganzer Dateien: 5–10 %
- Deduplikation hat auch Nachteile
  - Deduplikationstabellen müssen im RAM gehalten werden
    - Je 1 TB Daten ca. 5–20 GB für Tabellen

# Overhead

- Deduplikationstabellen speichern Referenzen zwischen Hashes und Daten
  - SHA256-Hash (256 Bit = 32 Byte)
  - 8 KB große Dateisystemblöcke (mit Offsets der Größe 8 Byte)
  - Zusätzlicher Overhead von 8 Byte pro Hash
- Für effiziente Online-Deduplikation müssen sie im Arbeitsspeicher gehalten werden
  - Duplikate müssen bei jeder Schreiboperation gesucht werden
  - Schnelle Speichergeräte (SSDs) sind immer noch um Größenordnungen langsamer

$$1 \text{ TB} \div 8 \text{ KB} = 125.000.000$$

$$125.000.000 \cdot (32 \text{ B} + 8 \text{ B} + 8 \text{ B}) = 6 \text{ GB} \quad (0,6 \%)$$

# Analyse

	2009	2015	2020	2025
Speicher	5,6+ <b>1,68</b> PB	45+ <b>13,5</b> PB	270+ <b>81</b> PB	1,6+ <b>0,48</b> EB
RAM	20+ <b>33,6</b> TB	170+ <b>270</b> TB	1,5+ <b>1,62</b> PB	12,8+ <b>9,6</b> PB
Strom	1,6+ <b>0,24</b> MW	1,4+ <b>0,20</b> MW	1,4+ <b>0,14</b> MW	1,4+ <b>0,09</b> MW
Kosten	30+ <b>2,52</b> M€	30+ <b>2,38</b> M€	30+ <b>1,62</b> M€	30+ <b>1,13</b> M€

- Optimistische Annahme von 30 % Einsparung
- Benötigt mehr zusätzlichen RAM als für Berechnung vorhanden ist (außer 2025)
- Benötigt deutlich mehr Strom (5–15 %)
- Erhöht Gesamtkosten (3–8 %)

# Analyse...

2009	2015	2020	2025
4,3+ <b>1,3</b> PB	34,6+ <b>10,4</b> PB	207,7+ <b>62,3</b> PB	1,2+ <b>0,4</b> EB
20+ <b>25,8</b> TB	170+ <b>207,7</b> TB	1,5+ <b>1,2</b> PB	12,8+ <b>7,4</b> PB
1,54+ <b>0,19</b> MW	1,34+ <b>0,15</b> MW	1,34+ <b>0,1</b> MW	1,34+ <b>0,07</b> MW
28,27+ <b>1,94</b> M€	28,27+ <b>1,83</b> M€	28,27+ <b>1,25</b> M€	28,27+ <b>0,87</b> M€

- Deduplikation wird genutzt, um selbe Kapazität zu erreichen
- Benötigt immer noch deutlich mehr Arbeitsspeicher
- Stromverbrauch erhöht sich (bis zu 8 %)
- Gesamtkosten sinken ab 2020

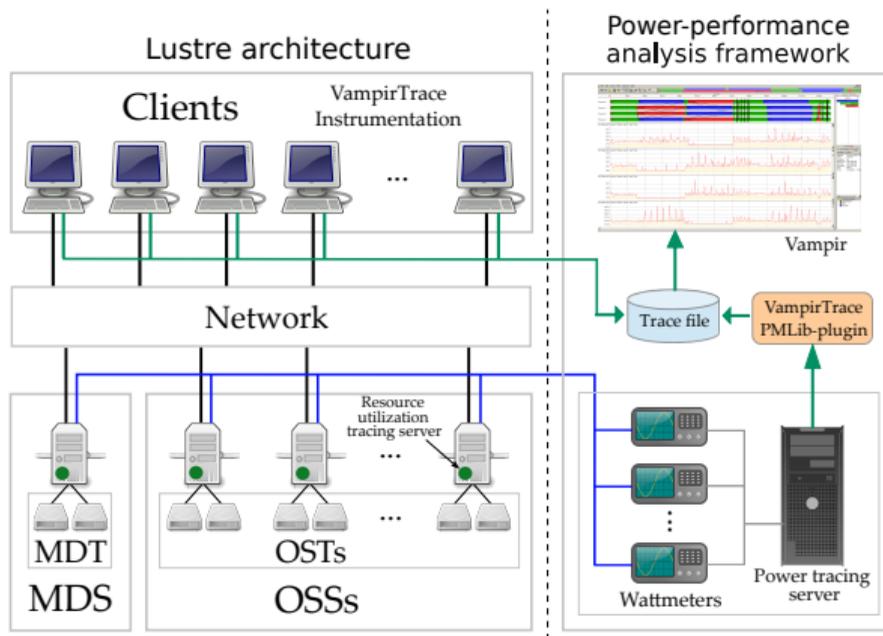
# Zusammenfassung

- Größere Blöcke senken den Overhead
  - 8 KB  $\rightarrow$  0,6 %, 16 KB  $\rightarrow$  0,3 %, 32 KB  $\rightarrow$  0,15 %
  - Einfluss auf Deduplikationsrate muss beachtet werden
- Deduplikation ganzer Dateien
  - E/A-Durchsatz bleibt gleich
  - Dateien müssen immer erst komplett geschrieben werden
- Offline-Deduplikation
  - Einfacher möglich mit modernen Copy-on-Write-Dateisystemen
  - Nützlich für Deduplikation ganzer Dateien
  - Nicht ganz so leistungskritisch
    - Tabellen müssen nicht immer im RAM gehalten werden

# Überblick

- Erfassung der wichtigsten Leistungsmetriken unterschiedlicher Kompressionsalgorithmen
  - Kompressionsrate, Prozessorlast, Stromverbrauch und Laufzeit
- $\approx$  500 GB Klimadaten (MPI-OM)
  - Vorabtests mit sich wiederholenden und zufälligen Daten
  - Serielle Tests für Grundleistung
  - Parallele Tests für echte Anwendungen

# Tracing



**Abbildung:** Framework zur Messung der Leistung und des Stromverbrauchs

# Tracing...

- Normale Lustre-Installation
  - Clients und Server auf unterschiedlichen Maschinen
- Zusätzliche Instrumentierung
  - VampirTrace für Client-Anwendungen
  - pmserver für Dateisystemserver
  - Server zur Erfassung des Stromverbrauchs
    - Verbunden mit Strommessgeräten
- pmlib-Plugin erlaubt Korrelation von Client- und Serveraktivität

# Algorithmen

Kompressionsalgorithmus	Kompressionsrate	Prozessorauslastung	Laufzeitverhältnis
none	1,00	23,7	1,00
zle	1,13	23,8	1,04
lzjb	1,57	24,8	1,09
lz4	1,52	22,8	1,09
gzip-1	2,04	56,6	1,06
gzip-9	2,08	83,1	13,66

**Tabelle:** Leistungsdaten für Klimadaten

- Laufzeit erhöht sich nur leicht (außer für höhere gzip-Level)
- gzip erhöht Prozessorauslastung deutlich
- lz4 (und gzip-1) am interessantesten

# Algorithmen...

Kompressionsalgorithmus	Kompressionsrate	Prozessorauslastung	Laufzeitverhältnis
none	1,00	23,7	1,00
lz4	126,96	15,8	1,28
gzip-1	126,96	23,3	1,24

**Tabelle:** Leistungsdaten für sich wiederholende Daten

- Mit Hilfe des `yes`-Kommandos erzeugt
- lz4 erzeugt niedrigere Prozessorlast als keine Kompression
- Beide Algorithmen erhöhen die Laufzeit um ca. 25 %

# Algorithmen...

Kompressionsalgorithmus	Kompressionsrate	Prozessorauslastung	Laufzeitverhältnis
none	1,00	23,5	1,00
lz4	1,00	24,1	0,97
gzip-1	1,00	66,1	1,03

Tabelle: Leistungsdaten für zufällige Daten

- Mit Hilfe des `frandom`-Kernelmoduls erzeugt
- `gzip-1` benötigt deutlich mehr Prozessorzeit
- Beide Algorithmen haben kaum Einfluss auf die Laufzeit
  - Erinnerung: serieller Test mit einer Festplatte



# Analyse

	2009	2015	2020	2025
Speicher	5,6+ <b>2,8</b> PB	45+ <b>22,5</b> PB	270+ <b>135</b> PB	1,6+ <b>0,8</b> EB
Strom	1,6+ <b>0,025</b> MW	1,4+ <b>0,025</b> MW	1,4+ <b>0,025</b> MW	1,4+ <b>0,025</b> MW

**Tabelle:** Vor- und Nachteile von Kompression

- Angenommene Kompressionsrate von 1,5 für LZ4
- Pessimistische Annahme von 10 % höherem Stromverbrauch
- Laufzeitverhältnis von 1,0
  - Nicht notwendig zusätzliche Prozessoren zu beschaffen

# Zusammenfassung

- Kompression kann Speicherkapazität deutlich erhöhen
  - Passende Algorithmen haben vernachlässigbaren Overhead
- Oft keine zusätzliche Hardware notwendig
- Geringer zusätzlicher Stromverbrauch
  - Insgesamt trotzdem lohnenswert
- Anwendungsspezifische Kompression kann Kompressionsraten deutlich erhöhen
  - Dadurch verlustbehaftete Kompression möglich
  - Kompressionsraten  $> 10$  sind möglich

# Überblick

- Kompression im Dateisystem kann bereits genutzt werden
  - Lustre unterstützt ein ZFS-Backend
  - Kompression kann einfach in ZFS aktiviert werden
- Momentan allerdings nur statische Ansätze für Kompression
  - Ein Kompressionsalgorithmus pro Dateisystem
  - Dynamischere Ansätze könnten effizienter komprimieren
- Semantische Informationen zur Verbesserung der Kompression
  - Adaptive Kompression muss auch raten
  - Effizientere anwendungsspezifische Kompression

## Überblick...

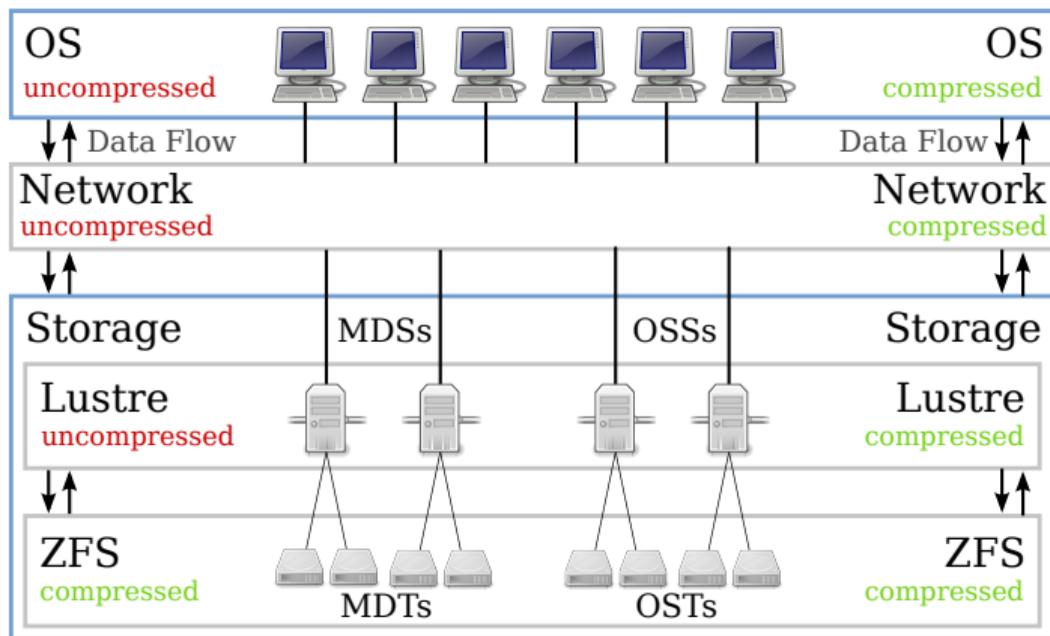


Abbildung: Lustre-Architektur mit Kompression

# Unterstützung

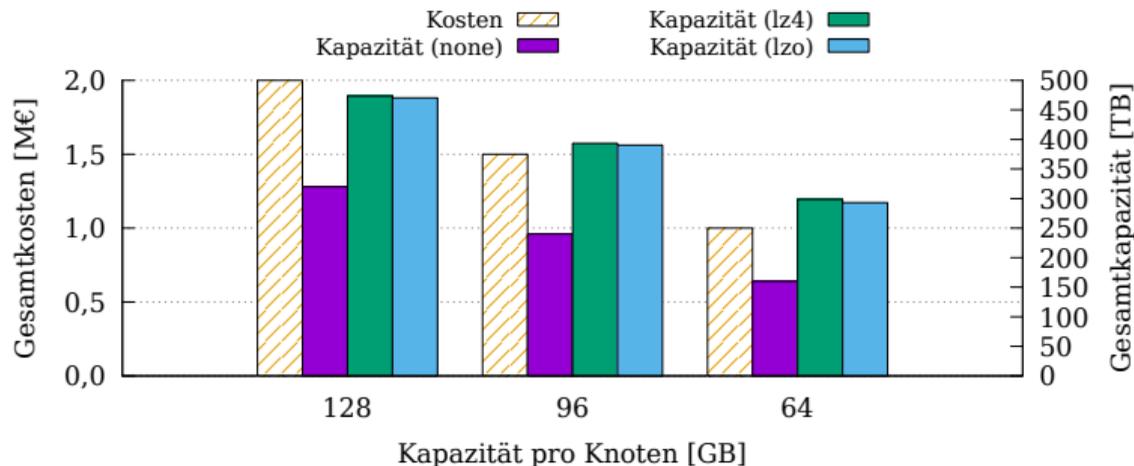
- Unterstützung auf mehreren Ebenen denkbar
  - Jeder Ansatz hat Vor- und Nachteile
  - Kompression auf dem Client beeinflusst die Berechnung, kann aber auch Netzwerkdurchsatz erhöhen
- Bisher keine Unterstützung für clientseitige Kompression
  - Komplet트 transparent für Anwendungen
  - Konfigurierbar mittels `ladvise`
- Kompression ist statisch
  - Nutzung von Informationen über die Daten, die aktuelle Last etc.
  - Sowohl auf Clients als auch Servern nützlich

# Kompression: Schichten

Algorithmus	Kompression	Dekompression	Rate
lz4fast	2.945 MB/s	6.460 MB/s	1.825
lz4	1.796 MB/s	5.178 MB/s	1.923
lz4hc	258 MB/s	4.333 MB/s	2.000
lzo	380 MB/s	1.938 MB/s	1.887
xz	26 MB/s	97 MB/s	2.632
zlib	95 MB/s	610 MB/s	2.326
zstd	658 MB/s	2.019 MB/s	2.326

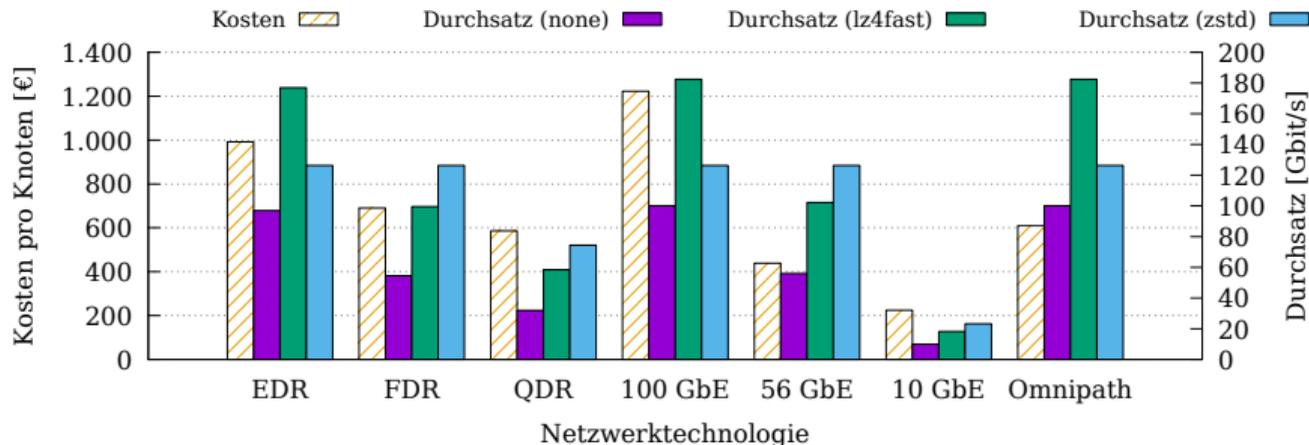
- lz4 und lz4fast allgemein sehr gut, zstd auch interessant
  - lz4fast und zstd können über Parameter angepasst werden
- Mehrere gute Kandidaten für Archivierung

# Kompression: Arbeitsspeicher



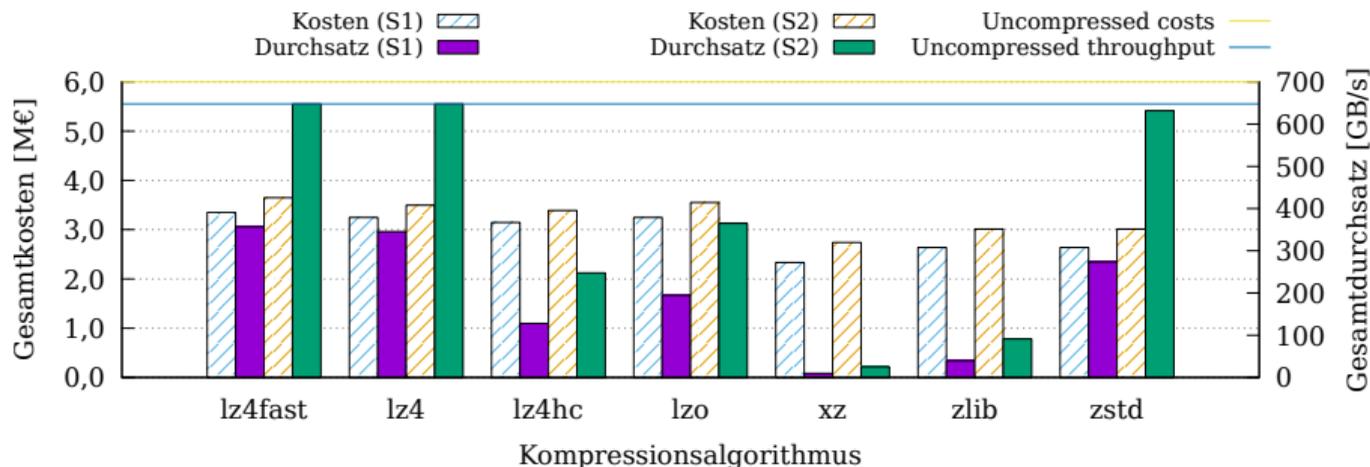
- Schätzungen mit zram
- Ziel: Kapazität pro Knoten von 128 GB
  - Nicht möglich mit 64 GB
  - 60 GB komprimiert, 4 GB unkomprimiert

# Kompression: Netzwerk



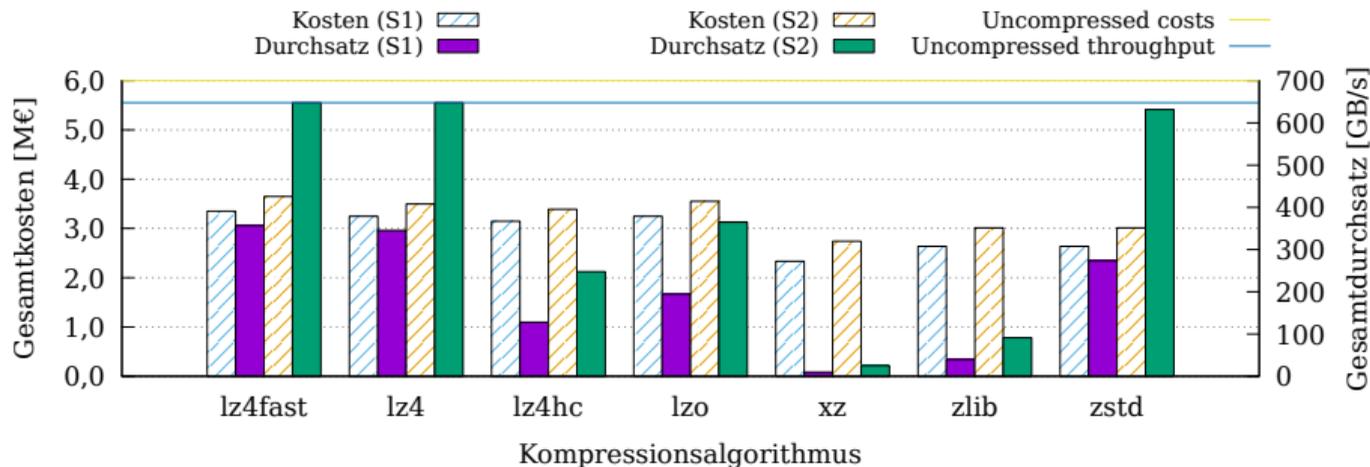
- zstd reduziert den Durchsatz bei Netzwerken mit hohem Durchsatz ( $> 54$  Gbit/s)
- FDR-InfiniBand kann mit QDR-InfiniBand ersetzt werden, wenn lz4fast genutzt wird (Kostenreduktion 15%)
- Durchsatz kann mit lz4fast bzw. zstd auf 100 Gbit/s bzw. 125 Gbit/s erhöht werden

# Kompression: Speicher



- S1: So viele SSU/ESU-Paare wie für 50 PB notwendig (geringere Kosten/Durchsatz)
- S2: 50 SSU/ESU-Paare und so viele HDDs wie für 50 PB notwendig (höhere Kosten/Durchsatz als S1)

# Kompression: Speicher...



- lz4 und lz4fast beeinflussen Leistung nicht
  - Kosten sinken auf 3.500.000 €
- zstd senkt Durchsatz um 20 GB/s
  - Kosten sinken um 50 % auf 3.000.000 €

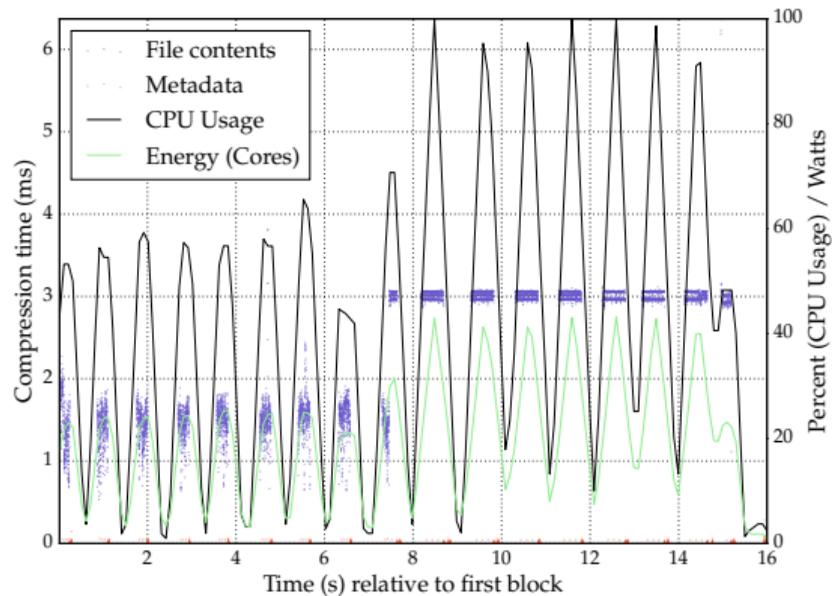
# Feature-Wunschliste

- Unterstützung für Kompression im parallelen Dateisystem
  - Interaktion mit anwendungsspezifischer Kompression
- Entwickler sollen nützliche Informationen spezifizieren können
  - Zusätzliches Wissen über die Daten (Varianz, Muster etc.)
  - Semantische Informationen im ganzen Stack nutzen
- Datenreduktion in einer zentralen Schicht
  - Momentan implementieren alle Schichten eigene Lösungen
  - Redundante Operationen, falsche Reihenfolge etc.

## Adaptive Kompression [2]

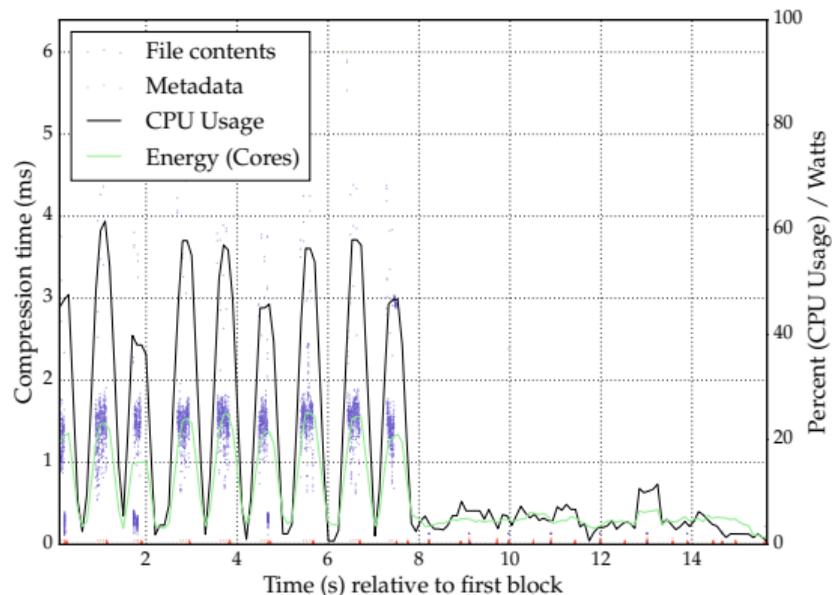
- Adaptive Kompression in ZFS
  - Direkt durch Lustre nutzbar
- Unterstützung für unterschiedliche Modi
  - Leistung, Archivierung, Energieverbrauch
- Unterschiedliche Heuristiken zur Bestimmung der Kompression
  - Basierend auf Dateitypen und Kostenfunktionen
- Für Kostenfunktion werden alle Algorithmen getestet
  - Der beste wird für die nächsten Operationen genutzt

# Adaptive Kompression... [2]



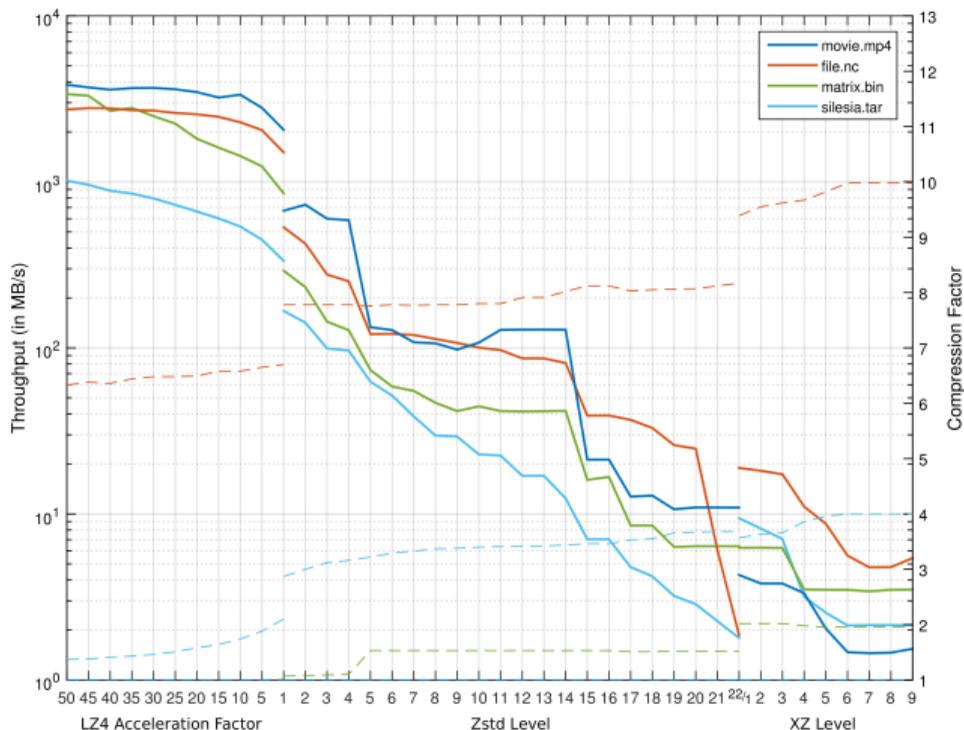
**Abbildung:** Komprimieren einer gemischten Datei mit `gzip-1`

## Adaptive Kompression... [2]



**Abbildung:** Komprimieren einer gemischten Datei mit dem Archivmodus

# Adaptive Kompression... [3]





# Zusammenfassung

- **Wiederberechnung**
  - Nicht alle Ergebnisse werden gespeichert
  - Negative Bilanz bei häufiger Wiederberechnung
- **Deduplikation**
  - Duplikate werden nicht gespeichert, sondern Referenzen auf existierende Blöcke
  - Overhead durch zusätzliche Checks
- **Kompression kann die TCO deutlich senken**
  - Positiver Effekt auf Arbeitsspeicher und Netzwerkdurchsatz
  - Nützlich für Daten, die nicht explizit durch Anwendungen komprimiert werden

## Zusammenfassung...

- Genaue Analyse der Kostenfaktoren und Nutzung notwendig
  - Rechenleistung steigt mit jeder Generation um Faktor 20
  - Speicherkapazität steigt nur um Faktor 8
- Schätzung der Kosten für Berechnung, Speicherung und Archivierung
  - Kostenmodelle für Langzeitarchivierung
  - Speicherung von Daten verursacht die größten Kosten
  - Das Problem wird sich weiter verschlimmern

- 1 Datenreduktion
  - Orientierung
  - Motivation
  - Speicherkostenmodell
  - Wiederberechnung
  - Deduplikation
  - Kompression
  - Erweiterte Kompression
  - Zusammenfassung

- 2 Quellen

# Quellen I

- [1] Konstantinos Chasapis, Manuel Dolz, Michael Kuhn, and Thomas Ludwig. Evaluating Power-Performance Benefits of Data Compression in HPC Storage Servers. In Steffen Fries and Petre Dini, editors, *IARIA Conference*, pages 29–34. IARIA XPS Press, 04 2014.
- [2] Florian Ehmke. Adaptive Compression for the Zettabyte File System. Master's thesis, Universität Hamburg, 02 2015.
- [3] Janosch Hirsch. Dynamic decision-making for efficient compression in parallel distributed file systems. Master's thesis, Universität Hamburg, 08 2017.
- [4] Michael Kuhn, Julian Kunkel, and Thomas Ludwig. Data Compression for Climate Data. *Supercomputing Frontiers and Innovations*, pages 75–94, 06 2016.

