

# Leistungsanalyse

## Hochleistungs-Ein-/Ausgabe

Michael Kuhn

Wissenschaftliches Rechnen  
Fachbereich Informatik  
Universität Hamburg

2015-06-22

- 1 Leistungsanalyse
  - Orientierung
  - Einleitung
  - Leistungsmessung
  - Leistungsbewertung
  - Zusammenfassung

- 2 Quellen





























# mdtest

- Datendurchsatz durch Benchmarks gut abgedeckt
  - Metadaten aber auch ein wichtiger Faktor
- mdtest nutzt MPI für parallelen Metadatenzugriff
  - Operationen werden mit POSIX-Schnittstelle durchgeführt
  - MPI-IO bietet nicht ausreichend Funktionalität
- Aufgeteilt in drei Phasen
  - Erstellen, Status abrufen und Löschen
- Arbeitet mit hierarchischer Struktur







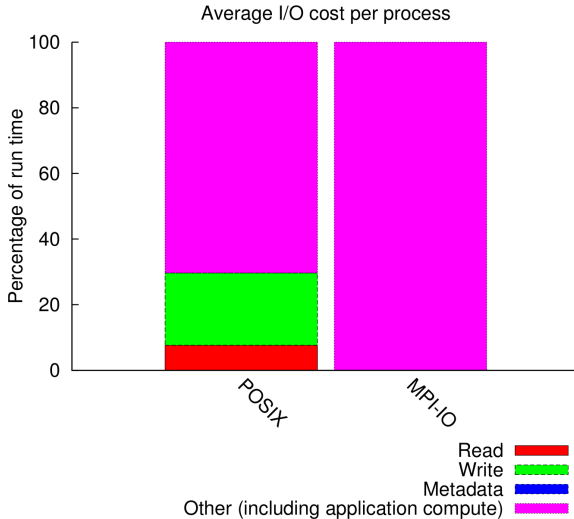




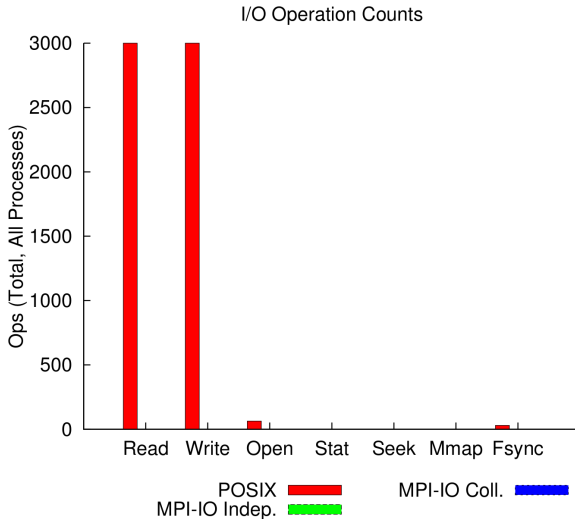
# Darshan...

- Besteht aus zwei Teilen
  - Runtime und Werkzeuge
- Runtime erfasst Anwendungs-E/A
  - Spezifisch für eine MPI-Implementierung
  - Außerdem Optionen für Batch-Scheduler und gemeinsames Protokollverzeichnis
  - Compiler-Wrapper und Preload-Bibliothek `libdarshan.so`
- Werkzeuge analysieren aufgezeichnete Protokolle
  - Befehle: `darshan-job-summary.pl`, `darshan-parser` etc.

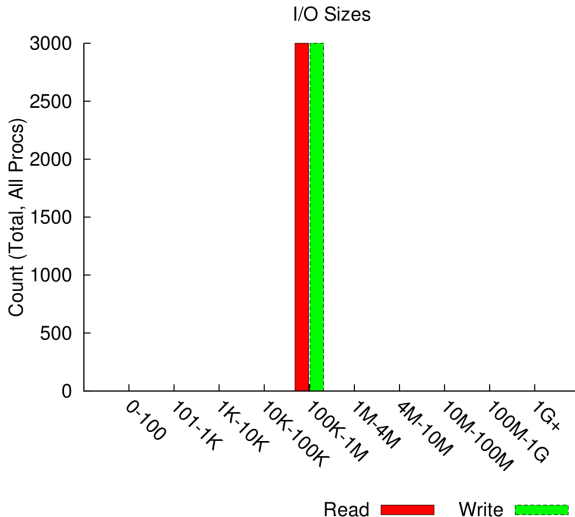
## Darshan...



## Darshan...

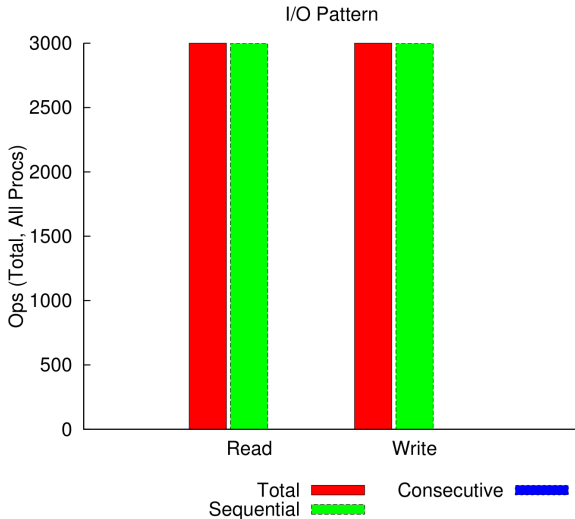


## Darshan...

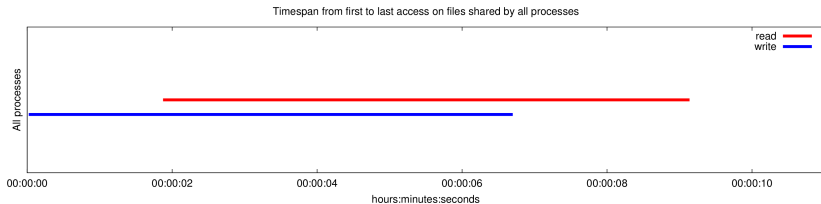




## Darshan...



# Darshan...

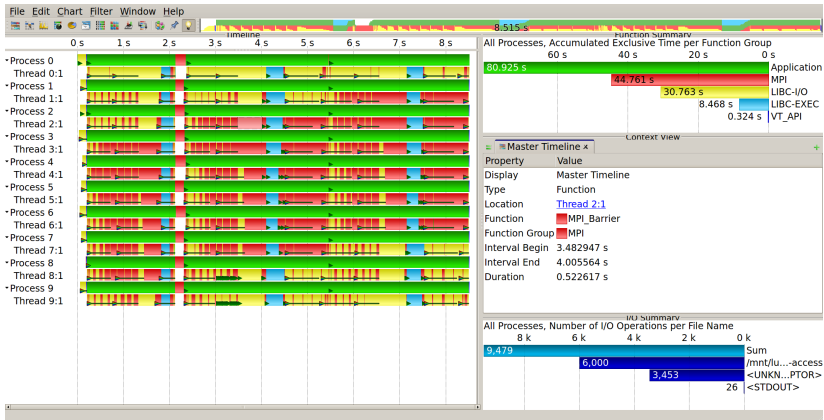


- Grobe Übersichten der Kosten von E/A
  - Bezüglich Aufrufanzahl, Zugriffsgröße und -muster
- Erlaubt Abschätzung ob Optimierung notwendig ist

# VampirTrace

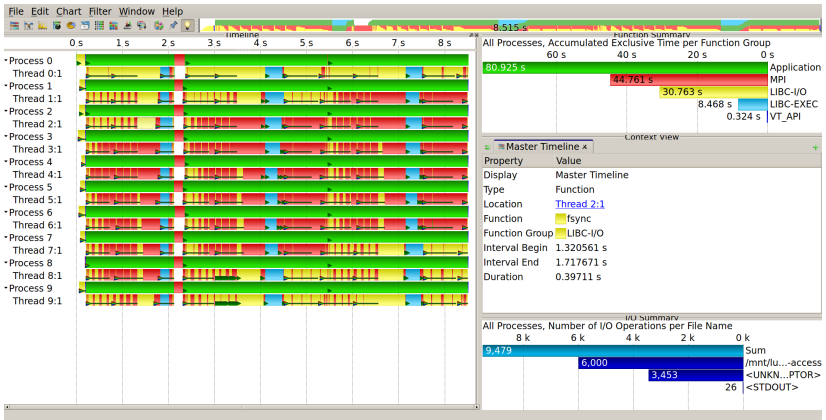
- VampirTrace zeichnet Spurdaten auf
  - VampirTrace ist Open Source
  - Nachfolger: Score-P
  - Spezifisch für eine MPI-Implementierung
  - Compiler-Wrapper vtcc
- Vampir zeigt Spurdaten an
  - Vampir ist kommerziell
  - Evaluationslizenzen verfügbar
- Spurdaten sind deutlich größer als Darshan-Protokolle
  - Im getesteten Fall mehr als Faktor 100

# VampirTrace...



- Sehr ungleiche E/A-Zeiten, dadurch lange Barries

# VampirTrace...



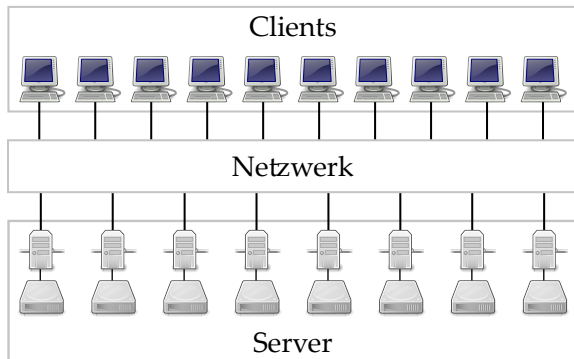
- Synchronisieren dauert teilweise sehr lange



# Einleitung

- Bewertung der Leistung durch Modellierung der theoretisch möglichen Leistung
- Dazu sind einige Informationen notwendig
  - Involvierte Komponenten
  - Leistungscharakteristika der Komponenten
- Zusätzliche Leistungsmessungen der Komponenten
  - Dazu wieder andere Werkzeuge

# Einleitung...



- Clients: IOPS, RAM, Bus zum Netzwerk
- Netzwerk: Durchsatz und Latenz
- Server: Durchsatz und IOPS der Speichergeräte, Bus



# Client

- Anzahl der E/A-Operationen pro Sekunde
  - Kontextwechsel könnten Beschränkung sein
- Durchsatz und Latenz des Hauptspeichers
  - Üblicherweise kein Problem
- Abschätzung mit Hilfe von `tmpfs` und `fio`

```
1 $ mkdir /tmp/fs
2 $ mount -t tmpfs tmpfs /tmp/fs
3 $ ...
4 $ umount /tmp/fs
```

# Client...

```
1 $ fio --name=switch --filename=/tmp/fs/foo --rw=write --bs=1
   ↪ --size=1g --runtime=60 [--numjobs=n]
2 $ vmstat 1
3 $ fio --name=bw --filename=/tmp/fs/foo --rw=write --bs=1m
   ↪ --size=$size --runtime=60
```

- Messung auf west
- $\approx 1.000.000$  IOPS
- $\approx 300.000$  Kontextwechsel
- $\approx 3$  GiB/s Durchsatz
  - Hauptspeicher vermutlich höher, da Overhead durch tmpfs



















# Zusammenfassung...

- Grobe Modellierung ist häufig schon ausreichend
  - Lässt sich bei Bedarf verfeinern
- Unvorhersehbares Verhalten macht Analyse schwierig
  - Siehe beispielsweise Leseleistung bei Lustre
- Analyse der konkreten Implementierung notwendig
  - Optimierung setzt viele Detailkenntnisse voraus

## 1 Leistungsanalyse

- Orientierung
- Einleitung
- Leistungsmessung
- Leistungsbewertung
- Zusammenfassung

## 2 Quellen

# Quellen I

- [1] CHAOS Development Team. IOR - Parallel filesystem I/O benchmark. <https://github.com/chaos/ior>.
- [2] Jens Axboe. fio - Flexible IO Tester. <http://git.kernel.dk/?p=fio.git;a=summary>.
- [3] Hongzhang Shan and John Shalf. Using IOR to Analyze the I/O Performance for HPC Platforms. In *In: Cray User Group Conference (CUG'07, 2007)*.