

Leistungsanalyse

Hochleistungs-Ein-/Ausgabe

Michael Kuhn

Wissenschaftliches Rechnen
Fachbereich Informatik
Universität Hamburg

2018-06-08



Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

1 Leistungsanalyse

- Orientierung
- Einleitung
- Leistungsmessung
- Leistungsbewertung
- Zusammenfassung

2 Quellen

fio [2]

- Flexibler E/A-Tester
 - Autor ist Jens Axboe, Maintainer des Block-Layers
 - Unter anderem verantwortlich für die CFQ-, NOOP- und Deadline-Scheduler sowie blktrace und `splice`
- Beliebige Workloads
 - Statt vieler spezieller Tests
- Unterstützung durch Job-Dateien
 - Gemeinsame und job-spezifische Parameter
 - Alles auch über Kommandozeile steuerbar

fio...

- Ausrichtung
 - Beispielsweise zur Ausrichtung der E/A an Seitengrenzen
- Komprimierbar- und Deduplizierbarkeit
 - Aktuelle SSDs komprimieren die Daten transparent
- Durchsatzlimit
 - Nützlich um Hintergrundlast zu erzeugen
- Verifikation
 - Überprüfung ob die gelieferten Daten den vorher geschriebenen entsprechen

IOR [1]

- Für parallele verteilte Dateisysteme Zugriff von mehreren Knoten notwendig
 - fio erlaubt nur mehrere Prozesse auf einem Knoten
 - IOR nutzt MPI für parallelen Zugriff
- Unterstützung für mehrere Backends
 - POSIX, MPI-IO, HDF5 und Parallel-NetCDF
- Einige unterschiedliche E/A-Modi
 - Gemeinsame oder prozess-lokale Dateien
 - Prozessumordnung zur Umgehung von Cachingproblemen
 - Client X schreibt, Client X+n liest

IOR...

```
1 IOR START
2   api=MPIIO
3   testFile=/path/to/file
4   repetitions=3
5   readfile=0
6   writefile=1
7   filePerProc=0
8   keepfile=0
9   segmentCount=10
10  blockSize=1g
11  transferSize=1m
12 RUN
13 IOR STOP
```

- Schreiben in gemeinsame Datei mit drei Wiederholungen
- Auch komplett über Kommandozeile steuerbar

mdtest...

- Unterstützt ähnliche Funktionalitäten wie IOR
 - Prozessumordnung für Cache-Umgehung
 - Mehrere Wiederholungen
- Optional können Daten in Dateien geschrieben werden
 - Bei Bedarf auch synchronisiert werden
- Unterstützung für mehrere Wurzelverzeichnisse
 - Beispielsweise um mehrere Metadatenserver zu nutzen

Werkzeuge

- Benchmarks erlauben Erfassung der momentanen Leistung
 - Kein Aufschluss über Ursachen etc.
- Analyse und Optimierung erfordern Werkzeuge
 - Einblick in innere Vorgänge notwendig
 - Üblicherweise mit Hilfe von Tracing (Score-P)
- Manchmal abstrakte Betrachtung ausreichend
 - Grobe Übersicht über E/A-Verhalten (Darshan)

Darshan

- Werkzeug zur E/A-Charakterisierung
 - Sanskrit für „Sicht“ oder „Vision“
- Möglichst genaues Bild der Anwendungs-E/A
 - Einschließlich Eigenschaften wie der E/A-Muster
 - Dabei aber minimaler Overhead
- Geeignet für dauerhaften Einsatz
 - Wurde mit über 750.000 Kernen getestet
- Sehr gute Unterstützung für MPICH
 - Argonne National Laboratory
 - Gruppe um OrangeFS, MPICH und ROMIO

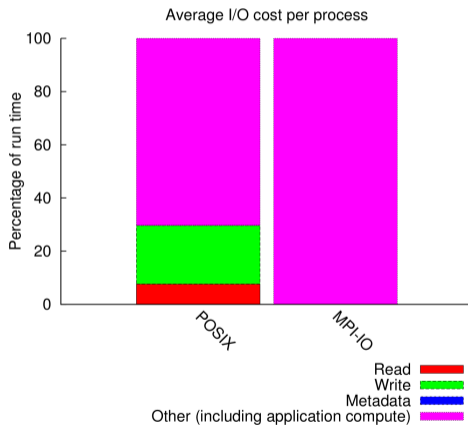
Darshan...

- 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
 - `darshan-job-summary.pl`, `darshan-parser` etc.

Darshan...

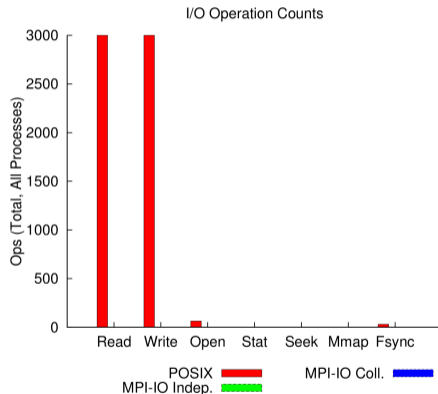
- MPI-parallelierter POSIX-Benchmark mit zehn Prozessen
 - Erst Schreibphase, dann Lesephase
 - Durch Barriers getrennt
 - Blockgröße von 1 MiB
 - 100 Blöcke
 - Leeren des Caches zwischen Schreib- und Lesephase
 - `echo 3 > /proc/sys/vm/drop_caches`
 - `fsync` beim Schließen
 - Drei Wiederholungen

Darshan...



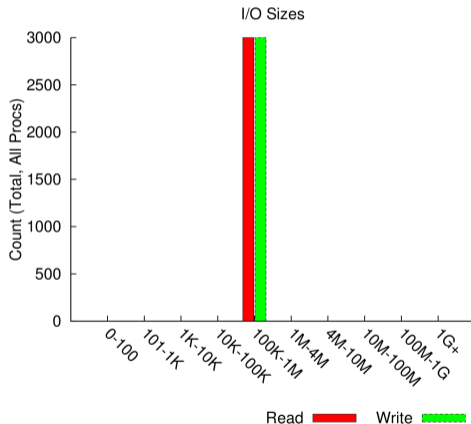
- Keine Berechnungen, trotzdem sehr hoher Other-Anteil

Darshan...



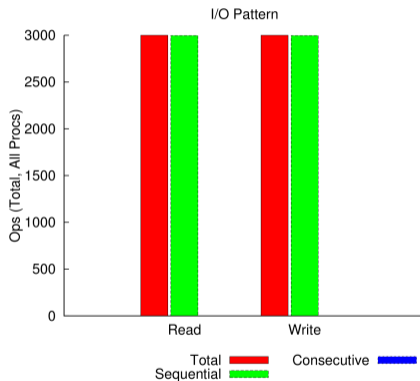
- Wie erwartet 3.000 Schreib- und Leseoperationen
 - 10 Prozesse × 100 Operationen × 3 Wiederholungen

Darshan...



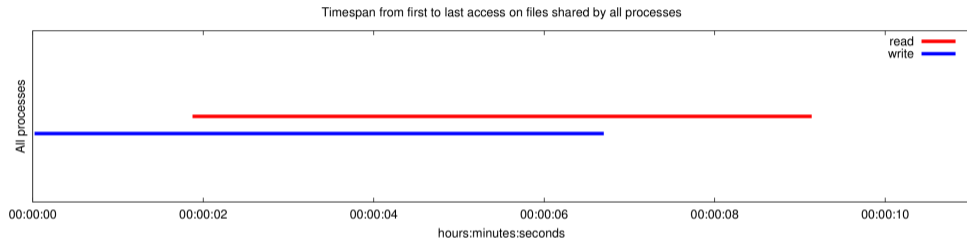
- Alle Operationen haben 1 MiB Größe

Darshan...



- Sequential: Zugriffe mit aufsteigendem Offset
- Consecutive: Zugriffe direkt hintereinander

Darshan...



- Anzeige durch drei Wiederholungen trügerisch
- Grobe Übersichten der Kosten von E/A
 - Bezüglich Aufrufanzahl, Zugriffsgröße und -muster
- Erlaubt Abschätzung ob Optimierung notwendig ist
 - Häufig genauere Analyse notwendig
 - Neuer Modus: Darshan eXtended Tracing (DxT)

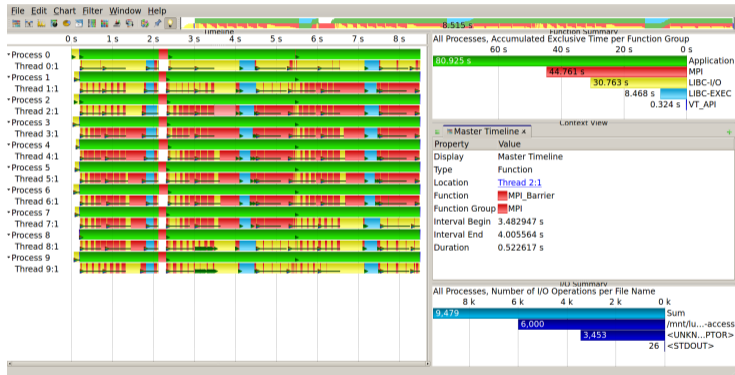
Score-P

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

Score-P...

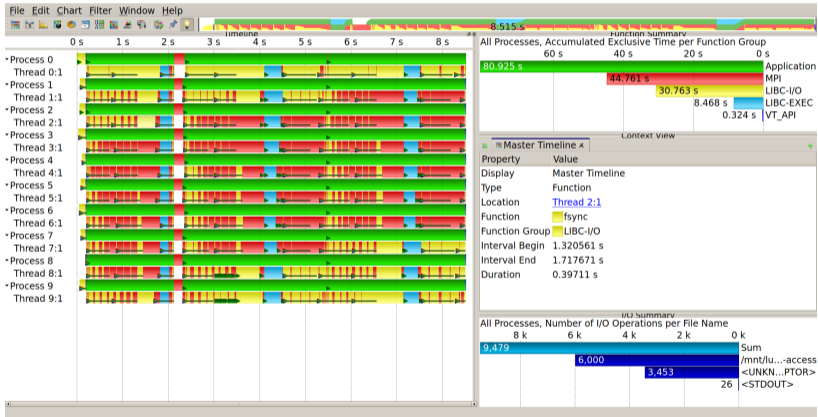
- Tracing erzeugt signifikanten Overhead
 - Laufzeit deutlich höher
 - Eventuell abweichendes Laufzeitverhalten
- Momentan noch keine E/A-Unterstützung
 - Dafür noch VampirTrace notwendig

Score-P...



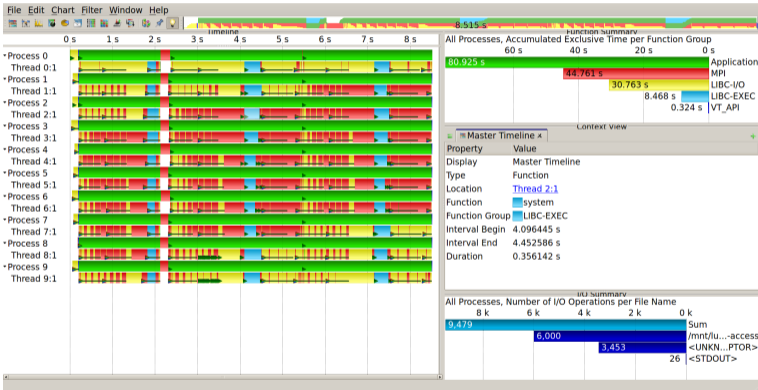
- Hauptthreads schlafen durchgängig
- Sehr ungleiche E/A-Zeiten, dadurch lange Barrieren

Score-P...



- Synchronisieren dauert teilweise sehr lange

Score-P...



- Leeren des Caches auch sichtbar
 - Blockiert bis Blöcke geschrieben sind

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

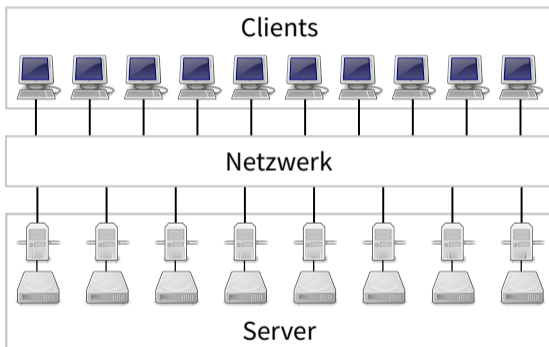
Beispiel

- Beispiel: Echtzeit-Twitter-Analyse
 - System kann in Echtzeit Stimmung in Tweets bestimmen
 - Alternativ Tweets speichern und später analysieren
- Frage: Wie ist die Leistung zu bewerten?

Beispiel

- Beispiel: Echtzeit-Twitter-Analyse
 - System kann in Echtzeit Stimmung in Tweets bestimmen
 - Alternativ Tweets speichern und später analysieren
- Frage: Wie ist die Leistung zu bewerten?
 - Ca. 6.000 Tweets pro Sekunde
 - 140 Bytes pro Tweet entspricht 840 KB/s
 - 26,5 TB pro Jahr

Übersicht



- 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`
 - Idee: Viele kleine E/A-Operationen ausführen

```

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-Knoten
- ≈ 1.000.000 IOPS
 - Blockgröße von 1 ist wichtig, da 0 eventuell durch libc behandelt
- ≈ 300.000 Kontextwechsel
- ≈ 3 GiB/s Durchsatz
 - Hauptspeicher üblicherweise höher, da Overhead durch tmpfs

Netzwerk

- Deutliche Unterschiede je nach Netzwerktechnologie
 - InfiniBand vs. Ethernet
- Durchsatz des Netzwerks
 - Sollte höher sein als Leistung der Speichergeräte
- Anzahl der Pakete pro Sekunde
 - Beschränkung bei vielen kleinen Nachrichten
 - Wichtig bei Metadatenoperationen
- Messung mit Hilfe von ping und iperf

Netzwerk...

```
1 $ ping -c 600 -f $host
2 $ iperf --server --port $port
3 $ iperf --client $host --port $port
```

- Messung zwischen west- und sandy-Knoten
- Round Trip Time $\approx 0,110$ ms

- Durchsatz ≈ 112 MiB/s

Netzwerk...

```
1 $ ping -c 600 -f $host
2 $ iperf --server --port $port
3 $ iperf --client $host --port $port
```

- Messung zwischen west- und sandy-Knoten
- Round Trip Time \approx 0,110 ms
 - Entspricht \approx 9.090 Nachrichten pro Sekunde
- Durchsatz \approx 112 MiB/s
 - Entspricht Ethernet mit 1 Gbit/s

Speichergerät

- Deutliche Unterschiede je nach Speichertechnologie
 - HDD vs. SSD
- Durchsatz der Speichergeräte
 - Niedriger als Netzwerkdurchsatz für maximale Ausnutzung
 - Höher als Netzwerkdurchsatz für Leistungsreserven
- IOPS wichtig für Metadatenoperationen
 - Auch für zufällige Datenzugriffe
- Speicherbus kann auch beschränkender Faktor sein
 - Teilweise immer noch SATA 2.0 (300 MB/s)

Beispiel

- Messung des Datendurchsatzes
 - Lustre und OrangeFS
- Unterschiedliche Blockgrößen
 - 1 MiB und 64 KiB
 - Entspricht den Streifenbreiten von Lustre und OrangeFS
- Bei 9.090 Nachrichten pro Sekunde
 - 9,09 GiB/s (1 MiB) bzw. 582 MiB/s (64 KiB) pro Knoten
- Flaschenhals ist Netzwerkdurchsatz
 - Maximal 1.120 MiB/s

Beispiel...

Dateisystem	Blockgröße	1 PPN	6 PPN	12 PPN
Lustre	1 MiB	640 MiB/s	105 MiB/s	110 MiB/s
OrangeFS	1 MiB	160 MiB/s	390 MiB/s	430 MiB/s
OrangeFS	64 KiB	250 MiB/s	115 MiB/s	180 MiB/s

Tabelle: Schreiben

Dateisystem	Blockgröße	1 PPN	6 PPN	12 PPN
Lustre	1 MiB	1.095 MiB/s	735 MiB/s	875 MiB/s
OrangeFS	1 MiB	130 MiB/s	265 MiB/s	430 MiB/s
OrangeFS	64 KiB	505 MiB/s	140 MiB/s	195 MiB/s

Tabelle: Lesen

- Massiver Leistungseinbruch bei Lustre
 - Beim Schreiben kein exklusiver Zugriff auf OST
- Streifenbreite als Blockgröße teilweise vorteilhaft

Beispiel...

Dateisystem	Blockgröße	1 PPN	4 PPN	8 PPN
Lustre	$\frac{1}{PPN}$ MiB	640 MiB/s	620 MiB/s	605 MiB/s
OrangeFS	$\frac{64}{PPN}$ KiB	250 MiB/s	280 MiB/s	210 MiB/s

Tabelle: Schreiben

Dateisystem	Blockgröße	1 PPN	4 PPN	8 PPN
Lustre	$\frac{1}{PPN}$ MiB	1.095 MiB/s	1.800 MiB/s	525 MiB/s
OrangeFS	$\frac{64}{PPN}$ KiB	505 MiB/s	655 MiB/s	455 MiB/s

Tabelle: Lesen

- Bessere Leistung in beiden Fällen
- Anomalie beim Lesen von Lustre

Zusammenfassung

- Benchmarks erlauben die Erfassung der aktuellen Leistung
 - Viele unterschiedliche Benchmarks für viele unterschiedliche Anwendungsfälle
- Werkzeuge sind notwendig zur eingehenden Analyse
 - Sowohl für Grobüberblick als auch zur Detailanalyse
- Messung sagt nichts über mögliche Leistungsfähigkeit aus
 - Dafür ist eine Bewertung der Leistung notwendig
 - Häufig durch Leistungsmodellierung

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] high performance computing. IOR - Parallel filesystem I/O benchmark.
<https://github.com/hpc/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*.