

15. März 2016



**MYTHOS**  
VIELKERN-BETRIEBSSYSTEME FÜR HPC

Stefan Wesner, Lutz Schubert, Vladimir Nikolov, Jörg Nolte,  
**Randolf Rotta**, Mai Krüger, Robert Kuban, uvm. . .



**Ziel** dynamische HPC Anwendungen auf **Vielkern-Prozessoren**

**Fokus** schnelle Thread-Verwaltung

- Durchsatz statt Echtzeit, Overheads reduzieren
- BMBF Projekt<sup>1</sup> von Uni Ulm, mit Uni Siegen, BTU Cottbus, Alcatel-Lucent Deutschland AG, HLRS Stuttgart

---

<sup>1</sup>Fördernummer BMBF 01IH13003

## 1. Wozu Vielkern-Betriebssysteme?

- Herausforderungen - Vielkern-Prozessoren
- Anforderungen - Zwei Szenarien
- Konsequenzen - Entwurfsziele

## 2. Zwischenergebnisse

- Die Kunst der Weglassens
- Wartefreie Aufgabenverwaltung
- Pragmatische Kompromisse

## 3. Zusammenfassung

## 1. Wozu Vielkern-Betriebssysteme?

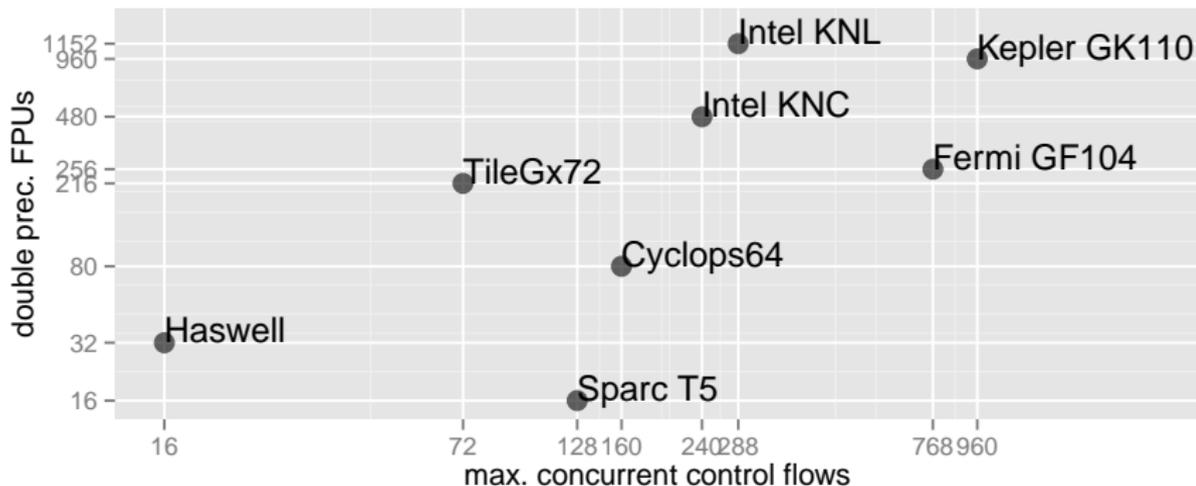
- Herausforderungen - Vielkern-Prozessoren
- Anforderungen - Zwei Szenarien
- Konsequenzen - Entwurfsziele

## 2. Zwischenergebnisse

- Die Kunst der Weglassens
- Wartefreie Aufgabenverwaltung
- Pragmatische Kompromisse

## 3. Zusammenfassung

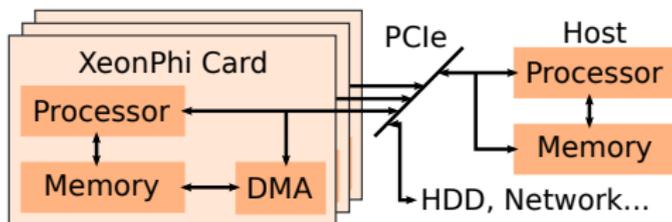
## MULTI- VERSUS MANY-CORE



- max. paralleler Durchsatz in geg. Power- & Platz-Budget
  - Energie- & Platzeffizienz statt Single-Thread Performance
- ⇒ **Durchsatz nur durch Aufgaben- und Daten-Parallelität**

## MYTHOS ZIELARCHITEKTUR

### Intel Xeon Phi Knights Corner



- 60 Kerne mit je 512bit VPU, 4 HW Threads, 1.05GHz  
In-Order, keine Sprungvorhersage, 2 Takte Instruction Decoder
  - 8GB kohärenter gemeinsamer Speicher, Streuung über Verzeichnisse  
150ns–1000+ns Latenz, Flaschenhals für Synchronisation
  - PCIe2.0 Karten: eigene BS Kernel einfach startbar,  
Kommunikation mit Host über nicht-kohärenten gem. Speicher
- ⇒ Ideale Worst-Case Plattform!?

## SZENARIO I: HPC SIMULATIONEN

Moleküldynamik (HLRS), Strömungen und Elektrodynamik (Uni Siegen)



Brandenburg  
University of Technology  
Cottbus - Senftenberg

- eine Anwendung besitzt “alle” Kerne und Speicher
- Aufgabenverwaltung und Kommunikation auf Anwendungsebene
- dynamische Speicherverwaltung

### Anforderungen

- effizientes Warten statt Polling (min Energie, min Contention)  
Semaphore/Futex `wait, signal`
- viele Threads manipulieren gemeinsamen Adressraum (IO, PGAS, DSM)  
`mmap, mremap, mprotect`

## SZENARIO II: COMPUTE CLOUD

### Verteilte Medienverarbeitung (Alcatel/CenoLabs)

- dynamisch konfigurierter Datenfluss-Graph
- Knoten aktiviert durch Input Daten, laufen in Isolation+Zeitschranke
- Cloud Manager bewegt Daten, migriert Knoten für min. Energie

### Anforderungen

- viele Threads mit eigenem Adressraum erzeugen und aktivieren  
`create`, `destroy`, `activate` mit Zeitschranke
- Manager manipuliert fremde Adressräume, z.B. `zero-copy comm`  
`mmap`, `mremap`, `mprotect` auf Kind-Threads
- Weiterleiten der Systemaufrufe, Exceptions etc. an Manager  
`event queue`

1. **Overheads reduzieren durch Weglassen:**  
Policies vereinfachen und auslagern  
z.B. kein Demand-Paging, kein Swapping, nur FIFO Scheduling
2. **BS-Aufgaben parallelisieren:**  
Aufräumaufgaben über freie Kerne verteilen  
z.B. Rahmen freigeben und nullen
3. **Wartezeiten sinnvoll füllen:**  
Asynchrone BS-Aufgaben und Kommunikation  
z.B. Active Messages, Delegation Queues, Continuation Passing
4. **Kommunikationskosten reduzieren:**  
kleine Aufgaben verschicken statt viele Daten bewegen  
lokal gemeinsame Caches ausnutzen

## 1. Wozu Vielkern-Betriebssysteme?

Herausforderungen - Vielkern-Prozessoren  
Anforderungen - Zwei Szenarien  
Konsequenzen - Entwurfsziele

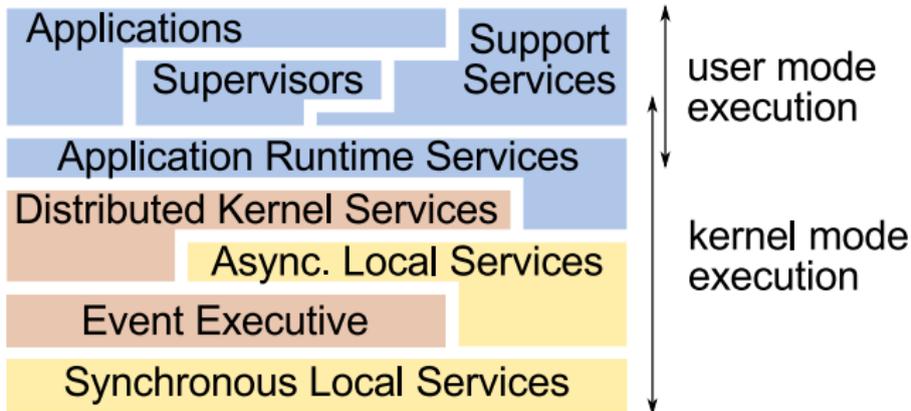
## 2. Zwischenergebnisse

Die Kunst der Weglassens  
Wartefreie Aufgabenverwaltung  
Pragmatische Kompromisse

## 3. Zusammenfassung

# MEHRSCHICHTIGE ARCHITEKTUR

Lokale und Verteilte Dienste separieren, Steuerung von Außen

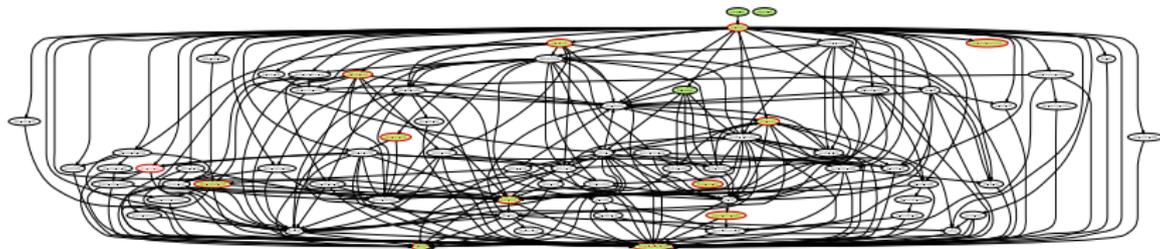


# DIE KUNST DER WEGLASSENS

mcconf: Kernel maßschneidern für Plattform und Anwendung

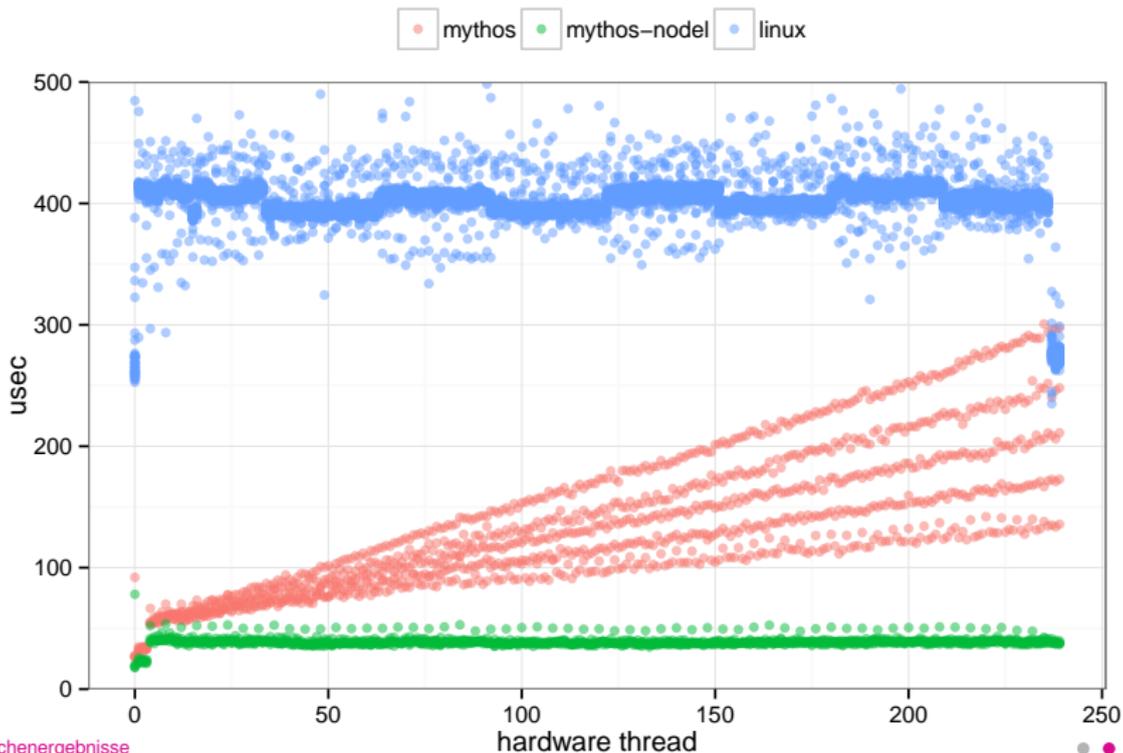
**Komponenten** Vielfalt durch Laufzeit-Komposition  
einfaches Replizieren & Platzieren auf Kerne

**Module** Effizienz durch Compile-Time-Komposition  
Abhängigkeitsverwaltung mit `#include`,  
Varianten filtern anhand Prozessor, Umgebung...



# FORK-JOIN BENCHMARK

Einen Thread erzeugen und aktivieren; Freigeben kostet noch...



## WARTEFREIE AUFGABENVERWALTUNG

### Asynchrones Ausführungsmodell mit lokaler Balanzierung

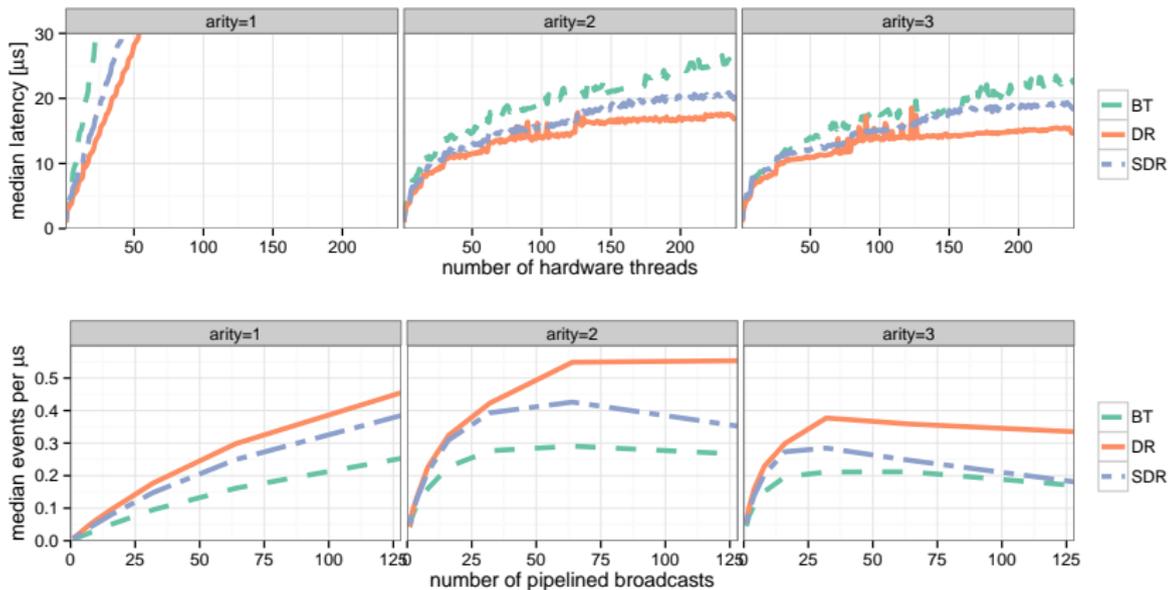
- Ereignisse, Aktive Nachrichten, Kritische Abschnitte, Interrupt-Epilog, Continuations. . .
- ⇒ "Tasklets" basierend auf C++11 Lambdas
- pro HW-Thread: Kernel-, User-, Idle-Warteschlangen
- Thread-Gruppen mit Work Sharing/Stealing, Delegation Locks regeln lokale Nebenläufigkeit

```
threads[40]->send( [] () { mlog::kernel.info("Hello"); } );
```

```
Mutex* m;  
m->protect (  
    [=] () { process_event(args...); ... }  
);
```

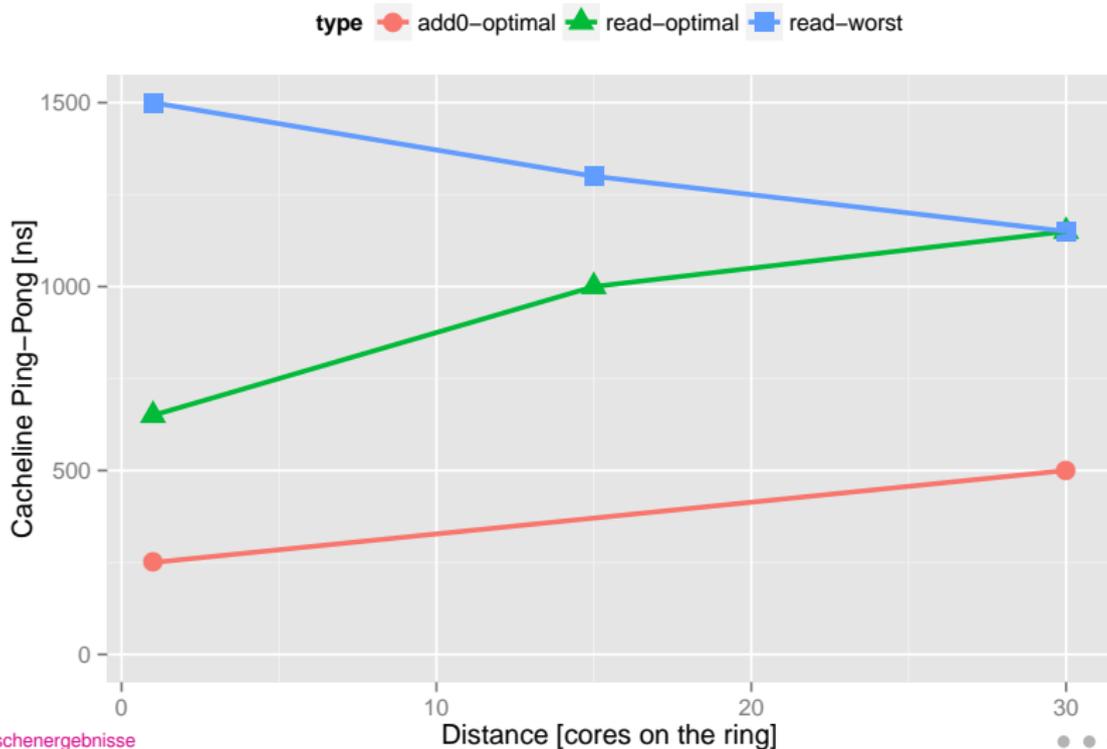
# BROADCAST BENCHMARK

## Latenz und Durchsatz auf dem XeonPhi



# SYNCHRONISATIONSLATENZ

## Geschickte Platzierung der Synchronisationsvariablen



- keine echten Prioritäten, keine Flusskontrolle  
⇒ zentrale Dienste beginnen neue Aufgaben statt alte abzuschließen
- doch viele Dienste und Policies im Kernel implementiert  
⇒ duplizierte Infrastruktur kernel- vs. user mode
- 8GiB/240 = 34MiB pro HW-Thread  
⇒ kein reiner Multikernel, gemeinsamer Speicherpool
- aktives Work-Stealing stört andere Threads

### Einsichten

- statische Ressourcenplanung wie in Echtzeitsystemen
- Microkernel mit optionalen kernel-mode Tasks
- wartefreie Datenstrukturen für Speicherpools

## 1. Wozu Vielkern-Betriebssysteme?

Herausforderungen - Vielkern-Prozessoren  
Anforderungen - Zwei Szenarien  
Konsequenzen - Entwurfsziele

## 2. Zwischenergebnisse

Die Kunst der Weglassens  
Wartefreie Aufgabenverwaltung  
Pragmatische Kompromisse

## 3. Zusammenfassung

- Vielkern-Prozessoren steigern die Herausforderungen
- MyThOS als Experiment-Plattform für parallele BS
- Zeigt Potentiale für dynamische Anwendungen

- Vielkern-Prozessoren steigern die Herausforderungen
- MyThOS als Experiment-Plattform für parallele BS
- Zeigt Potentiale für dynamische Anwendungen

Vielen Dank für Ihre Fragen!