



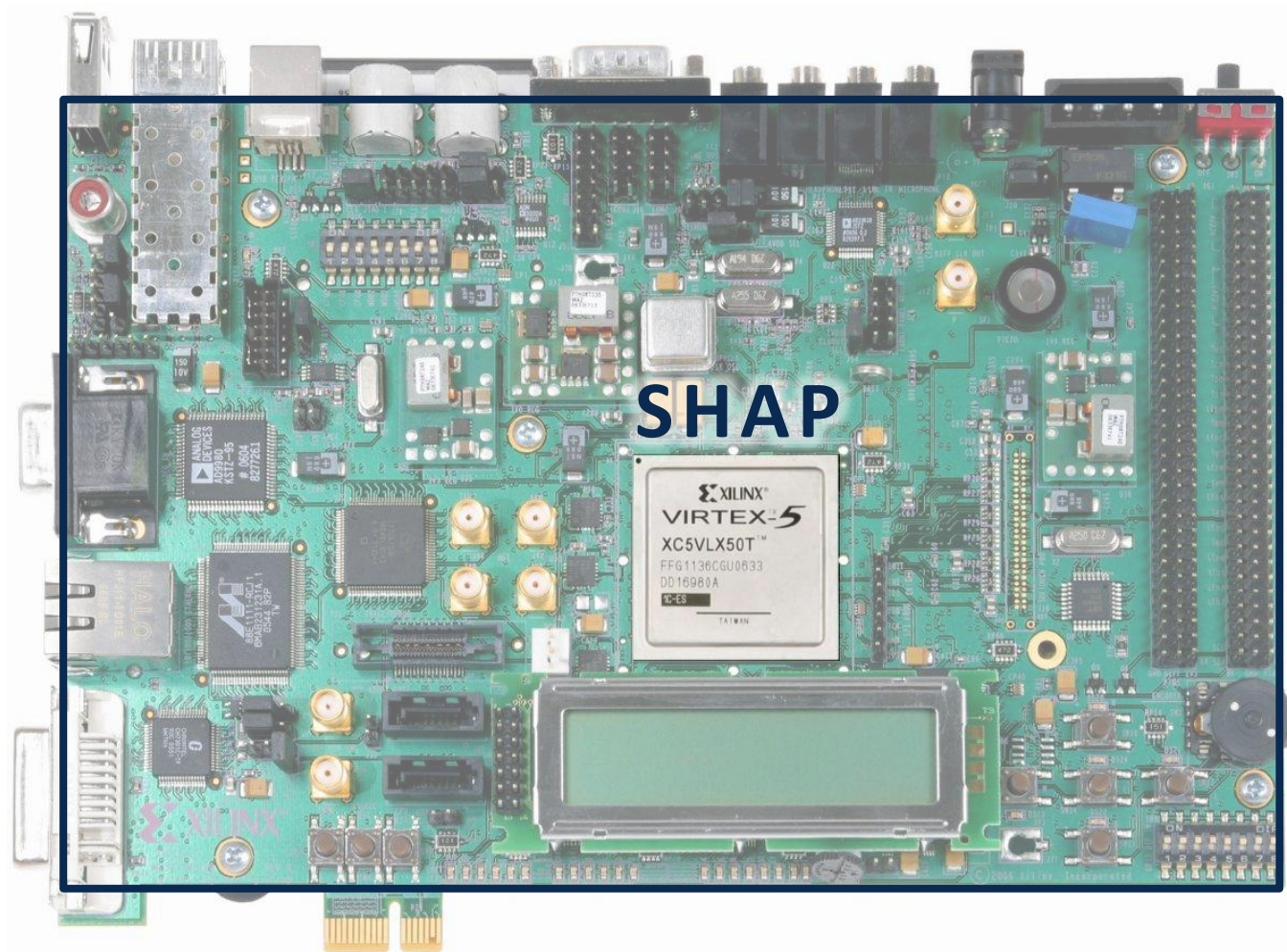
Diplomverteidigung

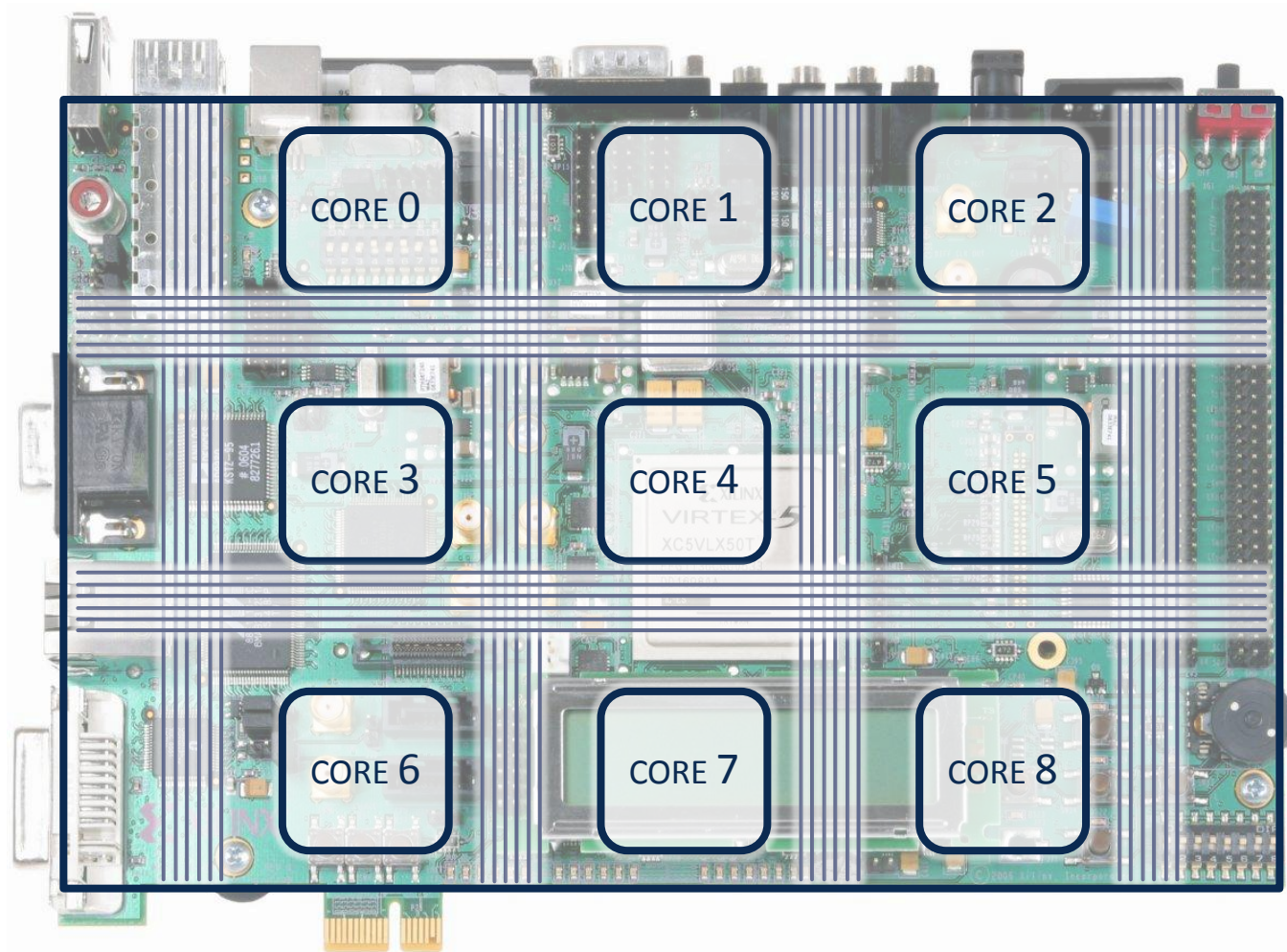
Implementierung eines Scheduling
mit dynamischer Lastverteilung für die
SHAP-Mehrkernarchitektur

Peter Ebert

Dresden, 12.02.2013







1. Grundlagen

- Scheduling-Algorithmen 5-16
- SHAP-Architektur 17-23

2. Implementierung

- Ambitionen 24
- Scheduler 25-29
- Locks 30-33
- Konzept 34-41
- Dialog der Kerne 42-53

3. Auswertung

- Ressourcen 54-60
- Leistung 61-77

4. Zusammenfassung 78-82

Klassische Strategien

- **Simpel** (First In First Out)
- **Vorhersage** (Shortest Job First, Shortest Remaining Time First)
- **Vorhersage ohne Verhungern** (Highest Response Ratio Next)
- **Rundenbasierend** (Round-Robin, VRR, WRR, DPRR)
- **Prioritäten statisch** (Rate Monotonic Sch., Deadline Monotonic Sch.)
- **Prioritäten dynamisch** (Earliest Deadline First, Least Laxity First, Multilevel Queue Sch., Multilevel Feedback Queue Sch.)

Java

- 10 Prioritäten
- Abbildung auf OS nicht vorgeschrieben

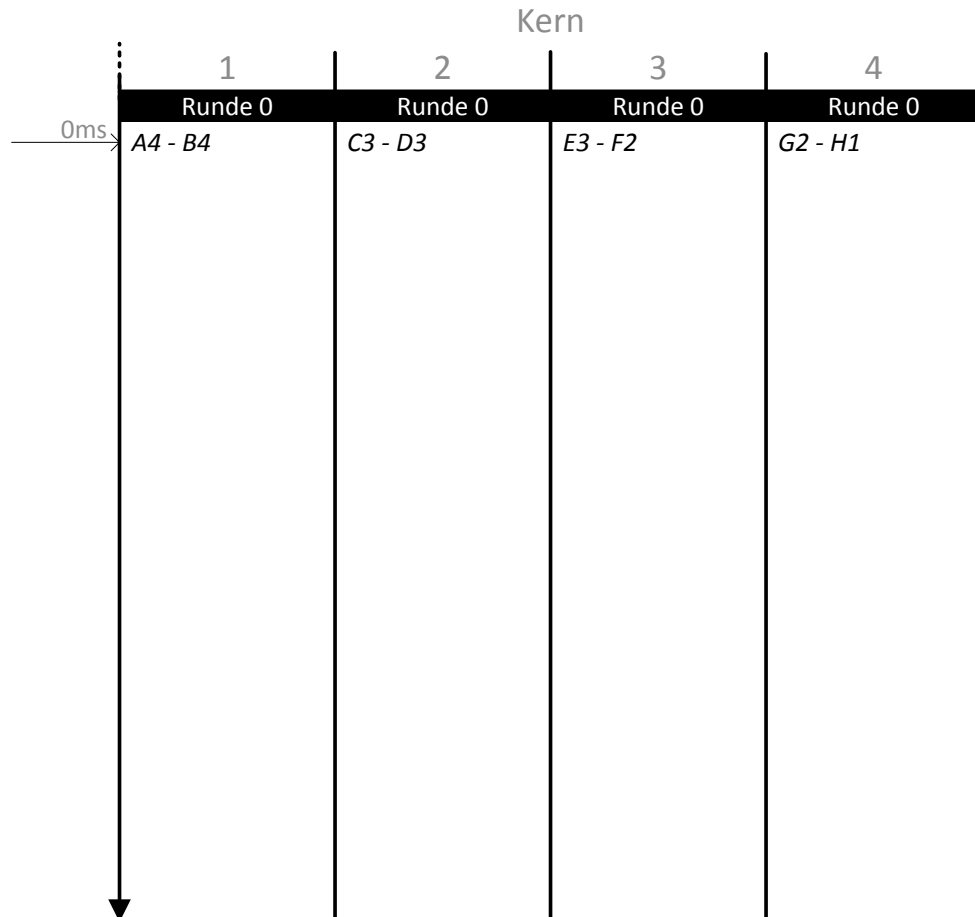
Mehrkern-Scheduler

- **Gang Scheduling** (Nutzung der Inter-Thread-Kommunikation)
- **Distributed Weighted Round-Robin** (Threadmigration zur Gleichverteilung der Rundenlängen)

Distributed Weighted Round-Robin

Wichtungen (Threads):

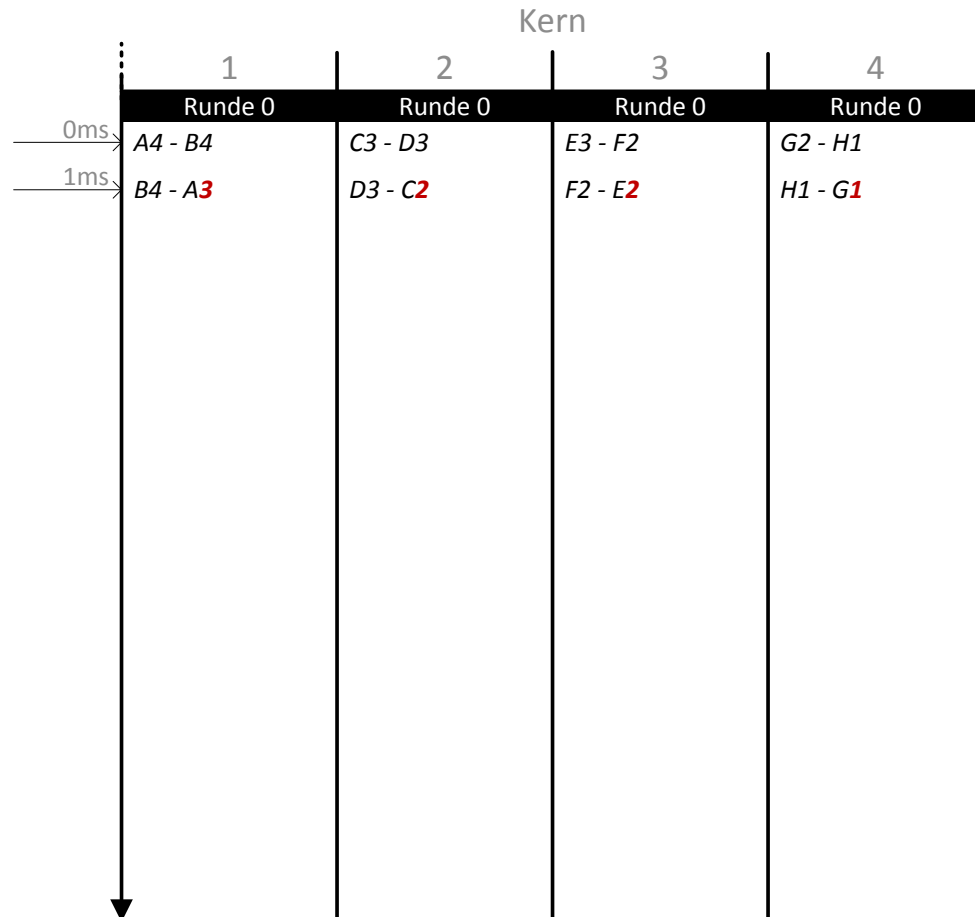
- 4 (A,B)
- 3 (C,D,E)
- 2 (F,G)
- 1 (H)



Distributed Weighted Round-Robin

Wichtungen (Threads):

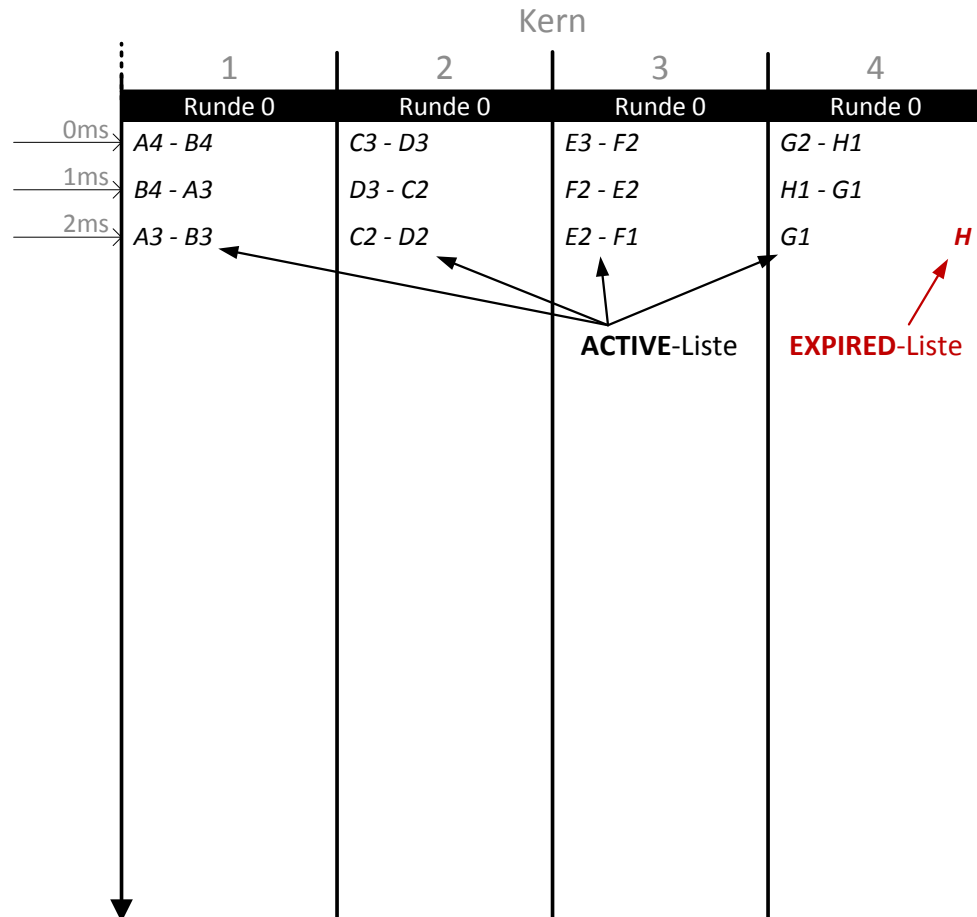
- 4 (A,B)
- 3 (C,D,E)
- 2 (F,G)
- 1 (H)



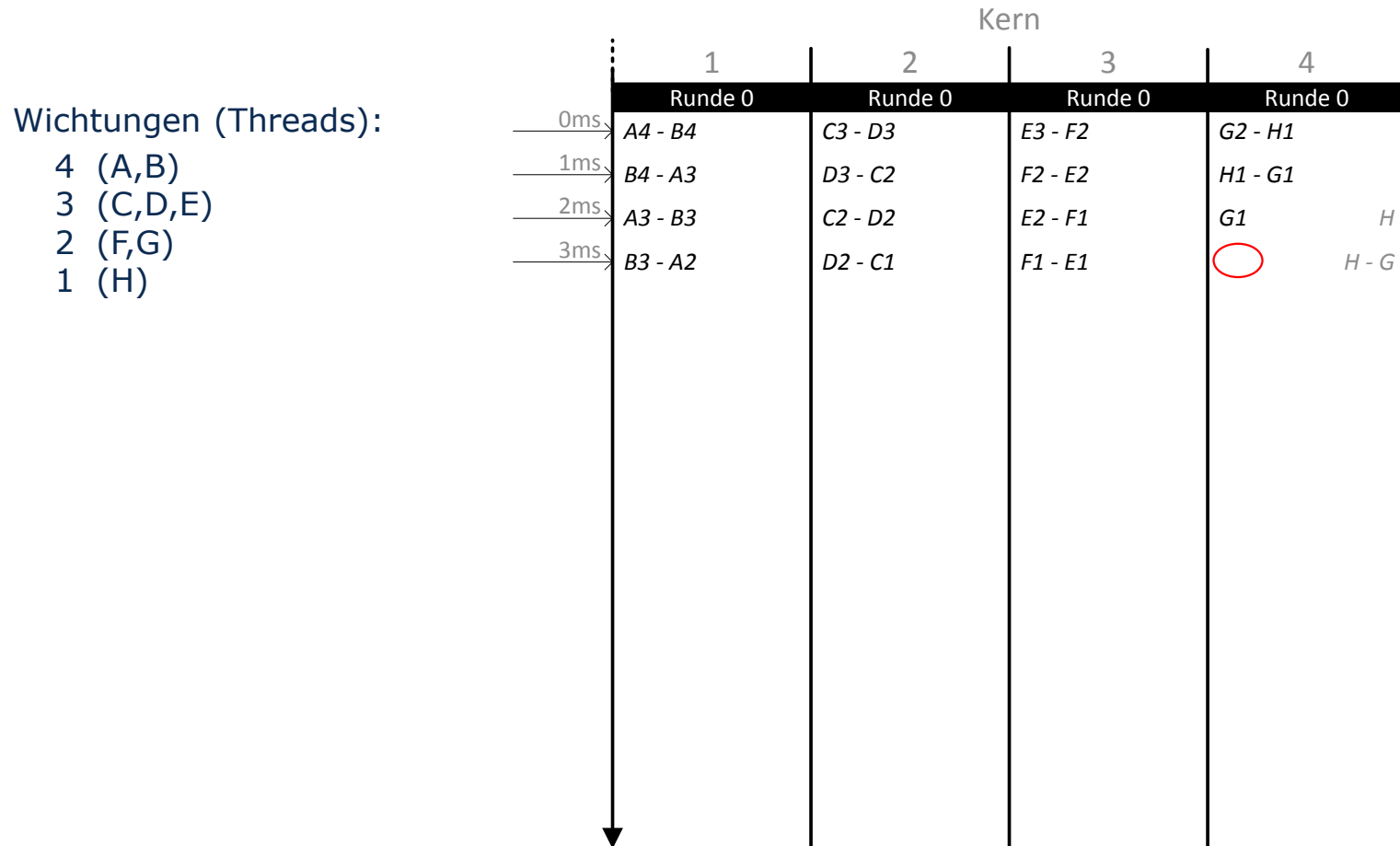
Distributed Weighted Round-Robin

Wichtungen (Threads):

- 4 (A,B)
- 3 (C,D,E)
- 2 (F,G)
- 1 (H)



Distributed Weighted Round-Robin



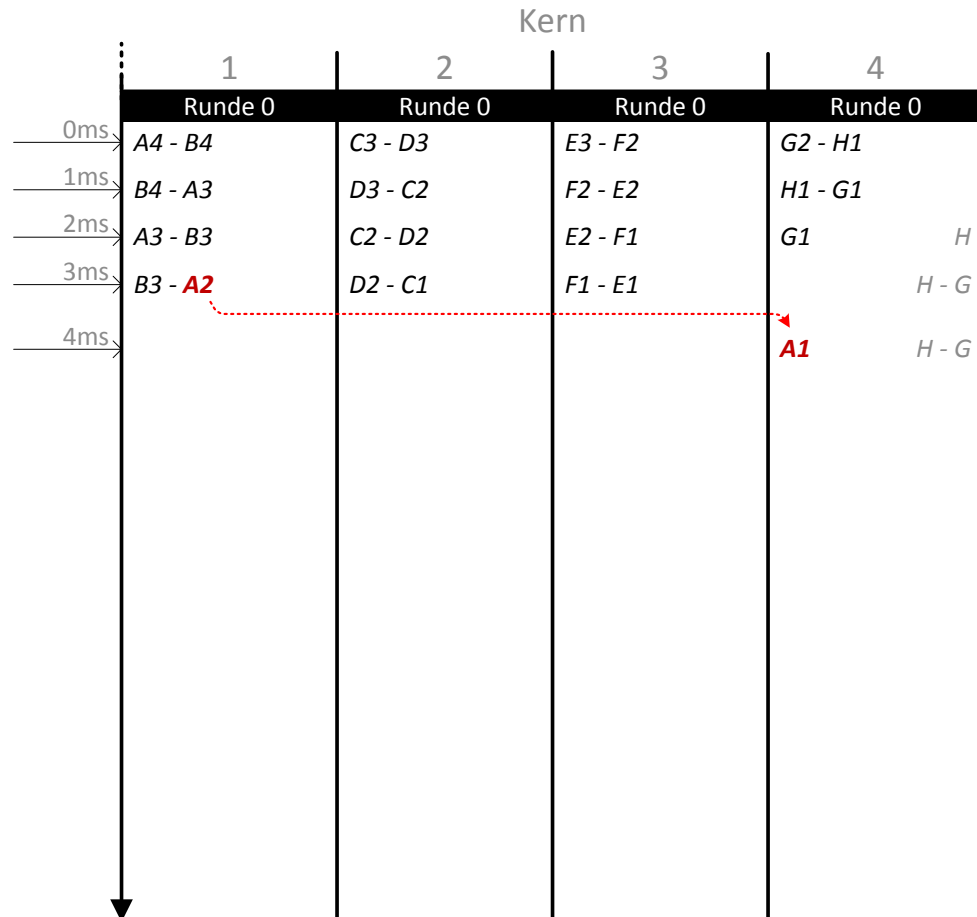
Distributed Weighted Round-Robin

Wichtungen (Threads):

- 4 (A,B)
- 3 (C,D,E)
- 2 (F,G)
- 1 (H)

Vom Kern mit den meisten Threads wird migriert.

Erster Thread aus ACTIVE wird migriert.



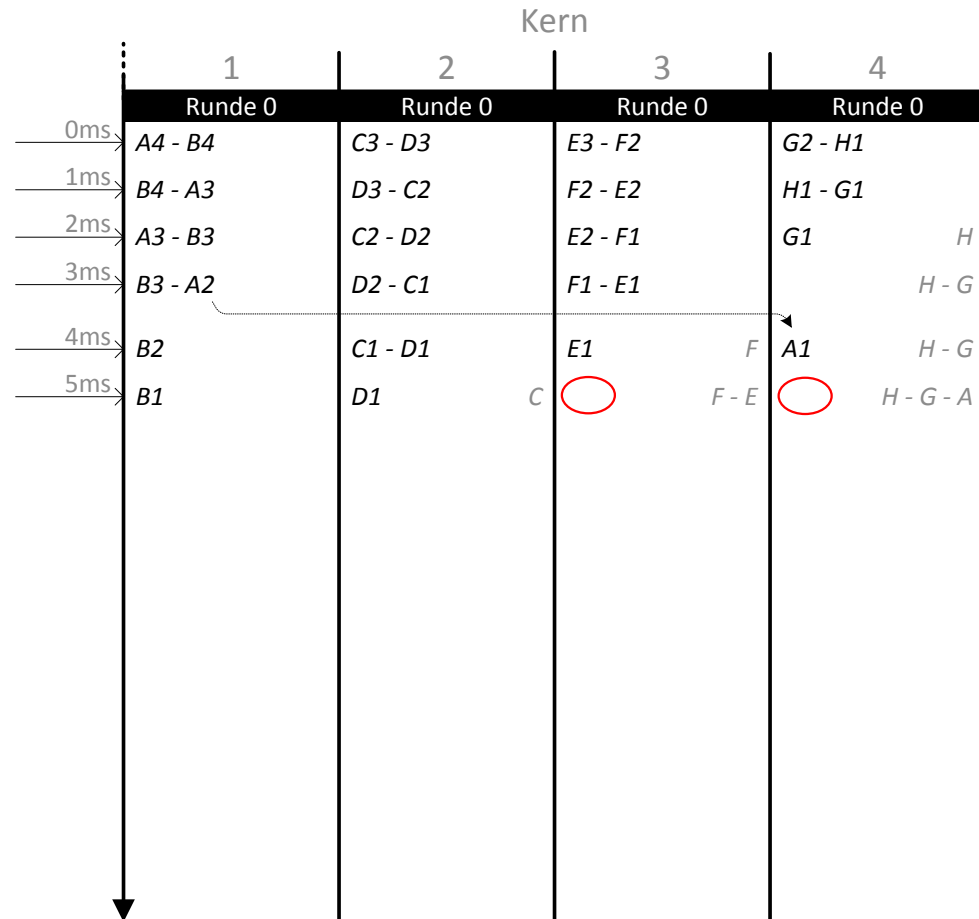
Distributed Weighted Round-Robin

Wichtungen (Threads):

- 4 (A,B)
- 3 (C,D,E)
- 2 (F,G)
- 1 (H)

Vom Kern mit den meisten Threads wird migriert.

Erster Thread aus ACTIVE wird migriert.



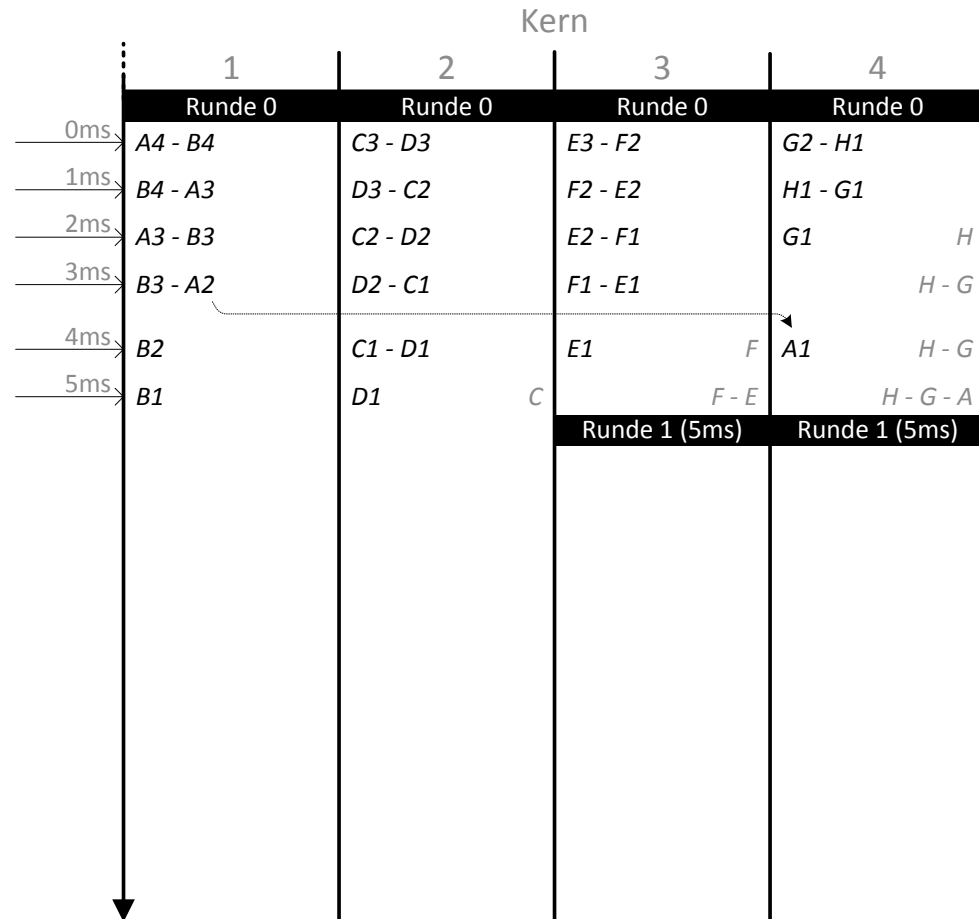
Distributed Weighted Round-Robin

Wichtungen (Threads):

- 4 (A,B)
- 3 (C,D,E)
- 2 (F,G)
- 1 (H)

Vom Kern mit den meisten Threads wird migriert.

Erster Thread aus ACTIVE wird migriert.



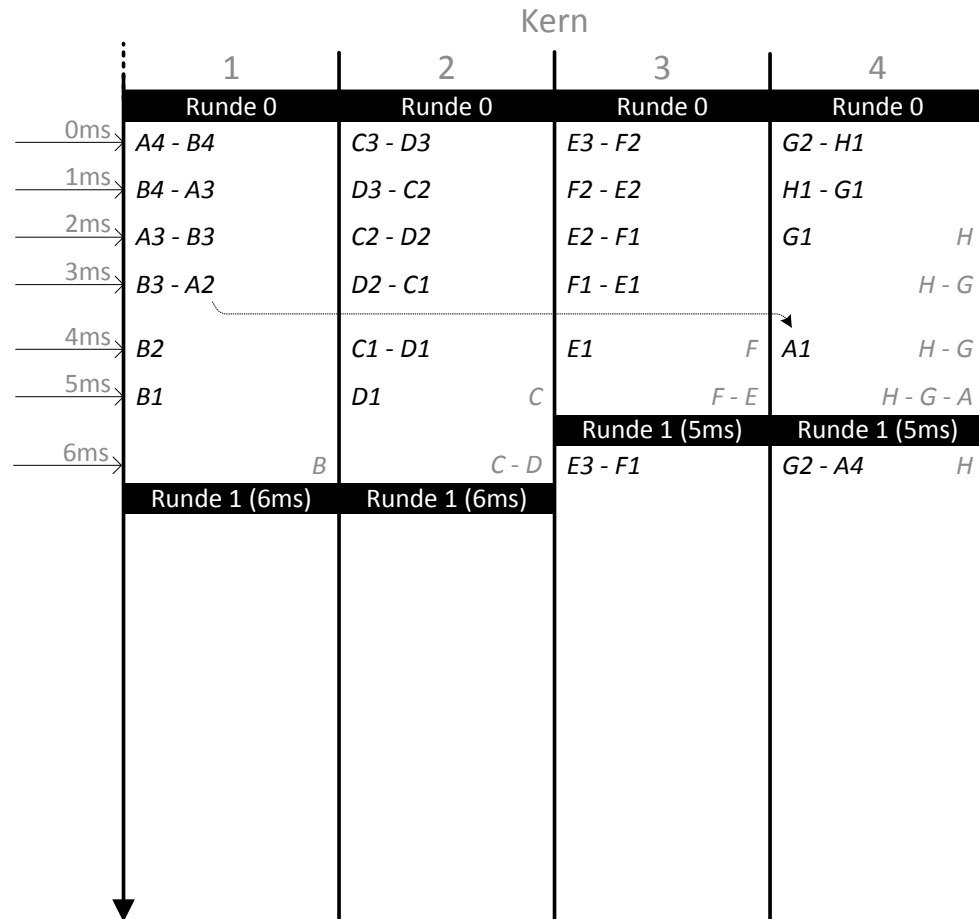
Distributed Weighted Round-Robin

Wichtungen (Threads):

- 4 (A,B)
- 3 (C,D,E)
- 2 (F,G)
- 1 (H)

Vom Kern mit den meisten Threads wird migriert.

Erster Thread aus ACTIVE wird migriert.



1.1 Scheduling-Algorithmen

Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur

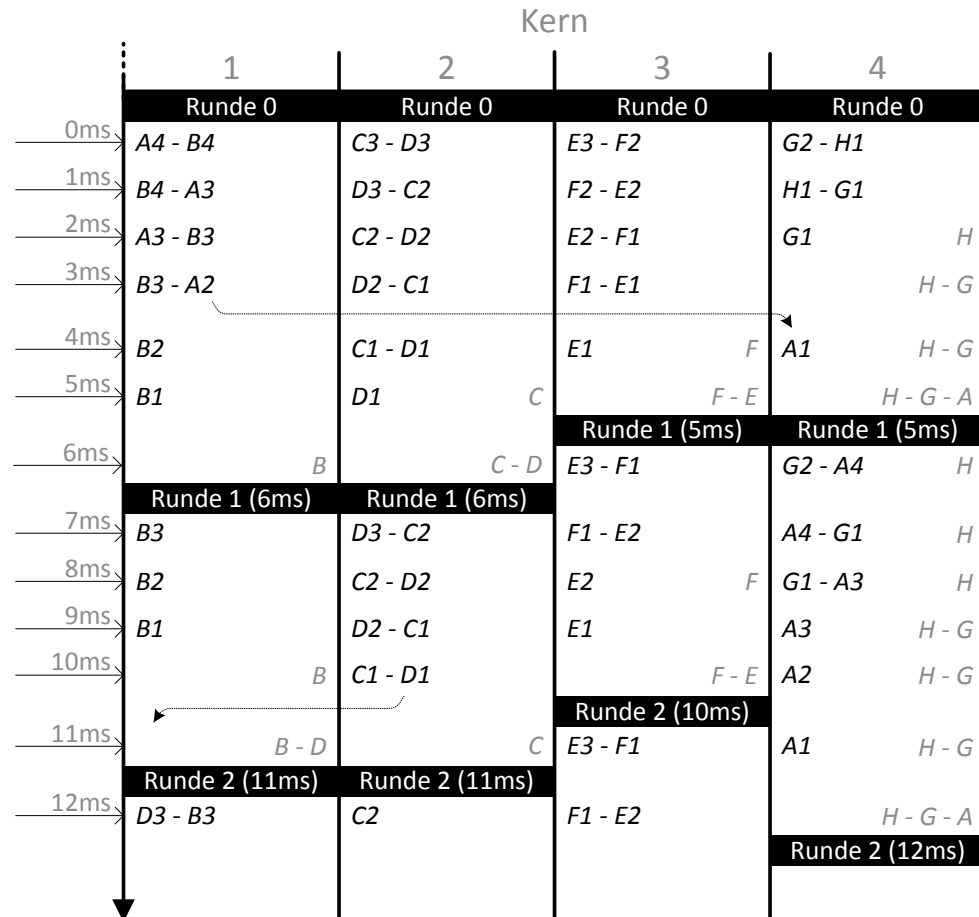
Distributed Weighted Round-Robin

Wichtungen (Threads):

- 4 (A,B)
- 3 (C,D,E)
- 2 (F,G)
- 1 (H)

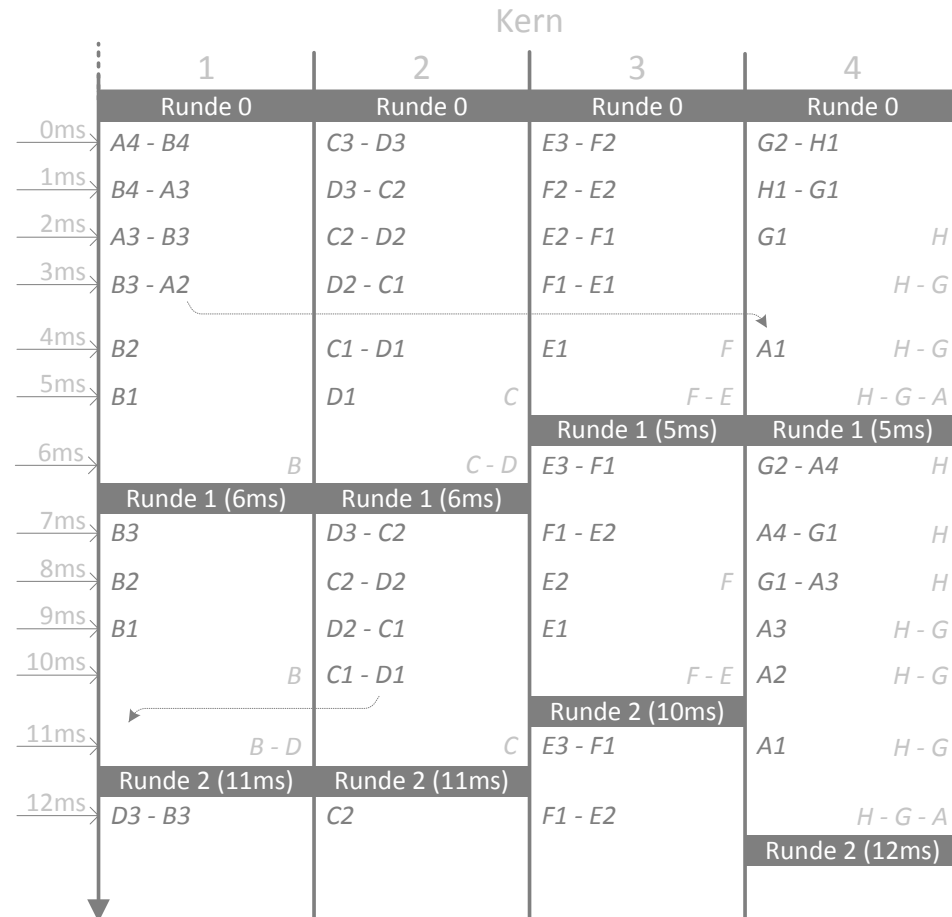
Vom Kern mit den meisten Threads wird migriert.

Erster Thread aus ACTIVE wird migriert.



Distributed Weighted Round-Robin

- + Ressourcen-schonend
- + Freiheit in Detailfragen
- + kein Leerlaufen
- + keine Synchronisierung d. Kerne (keine Wartezeiten)
- + Unterstützung von Prioritäten (Wichtungen)
- + Fairness
- + zusätzl. inneres Scheduling
- + gute Skalierbarkeit
- + kleine max. Antwortzeit
- keine Inter-Thread-Kommunikation
- Übersteuern möglich



1.1 Scheduling-Algorithmen

Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur

Distributed Weighted Round-Robin

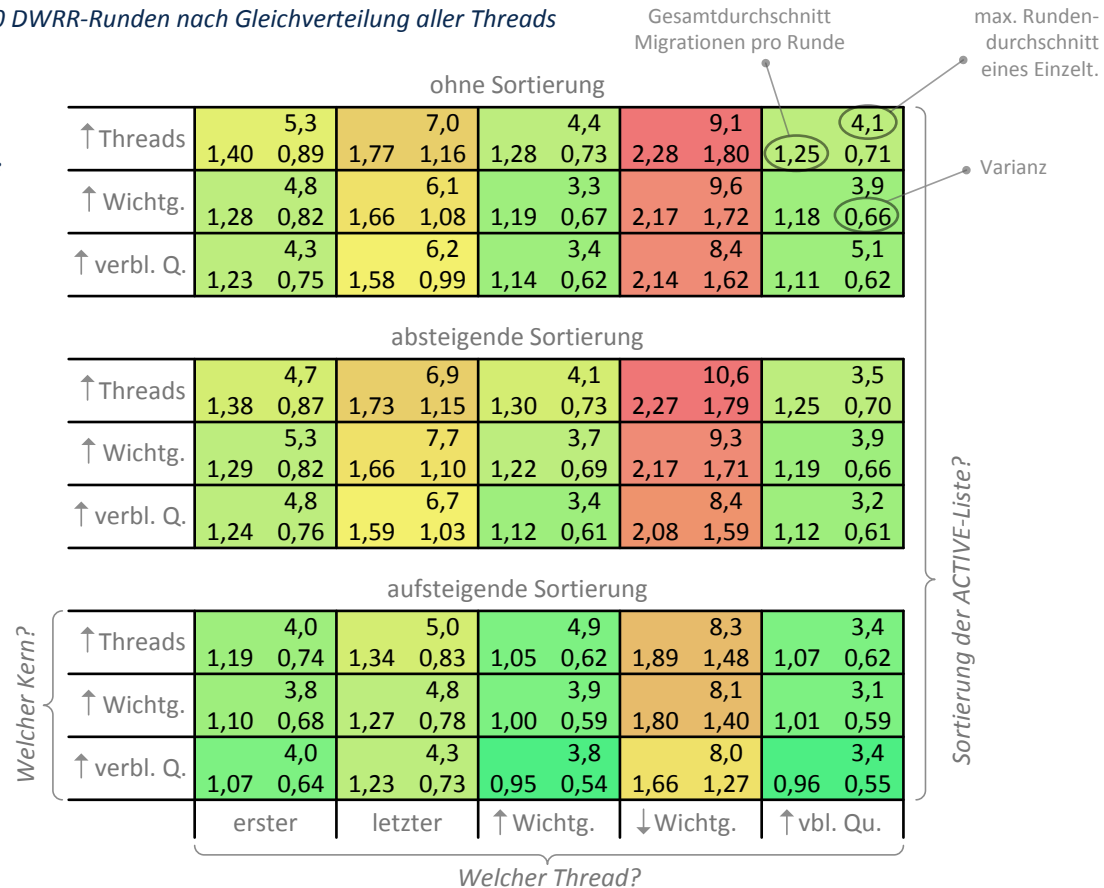
5000 Einzeltests mit je 30 DWRR-Runden nach Gleichverteilung aller Threads

Kerne: [2,18]

Threads: [3, 198]*

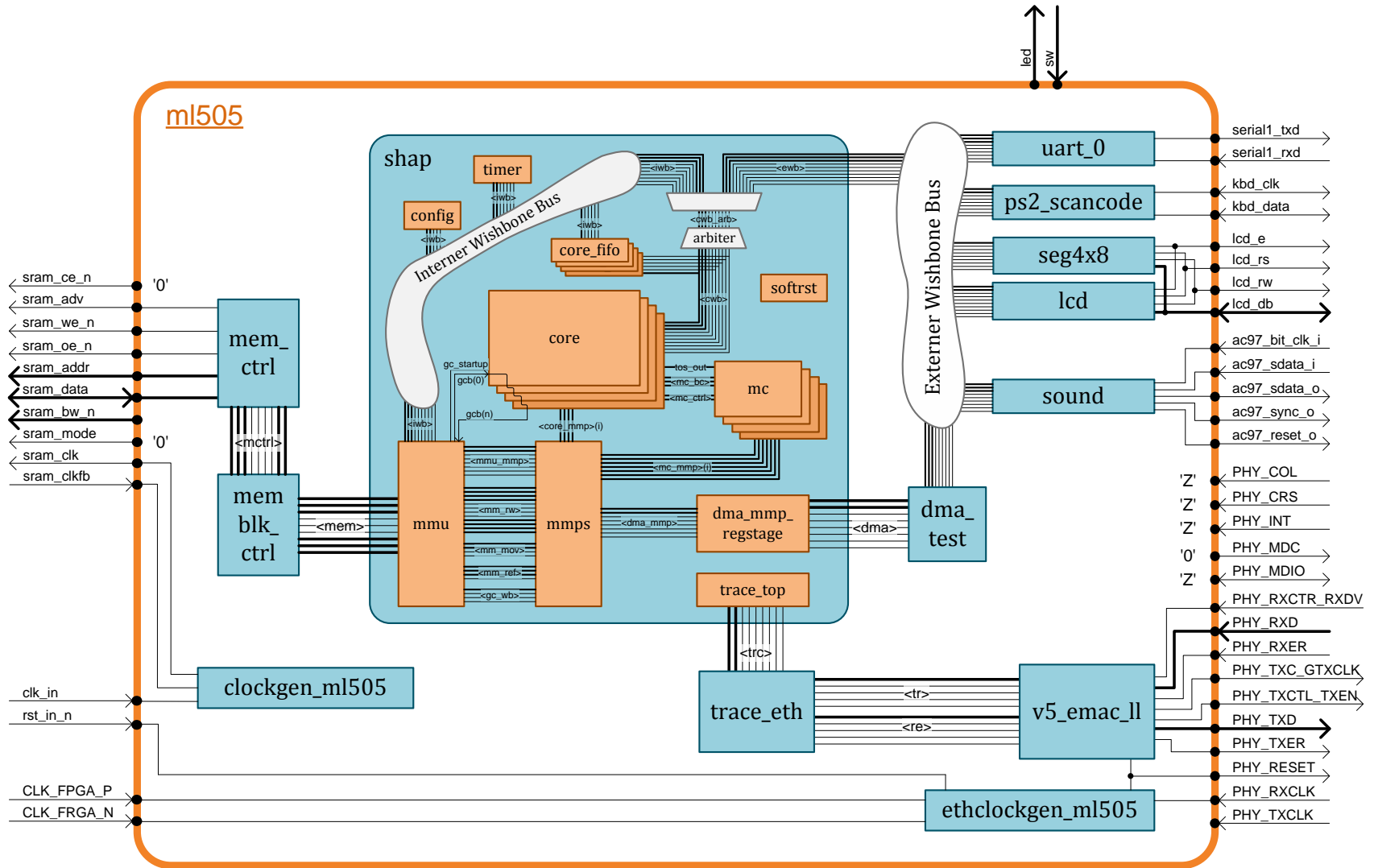
Wichtigungen: [1, 10]

* abhängig von der Anzahl der Kerne



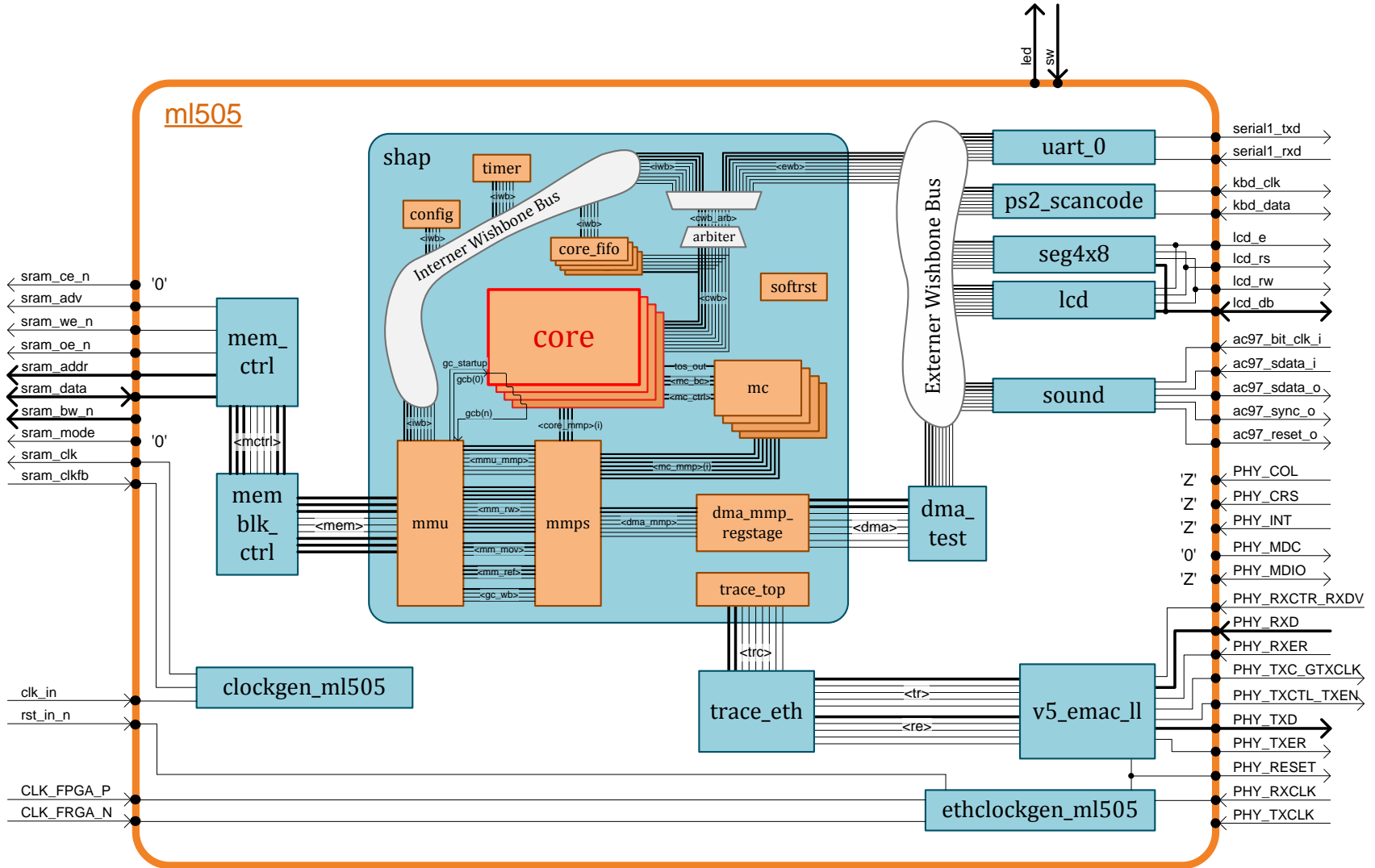
1.2 SHAP-Architektur

Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur



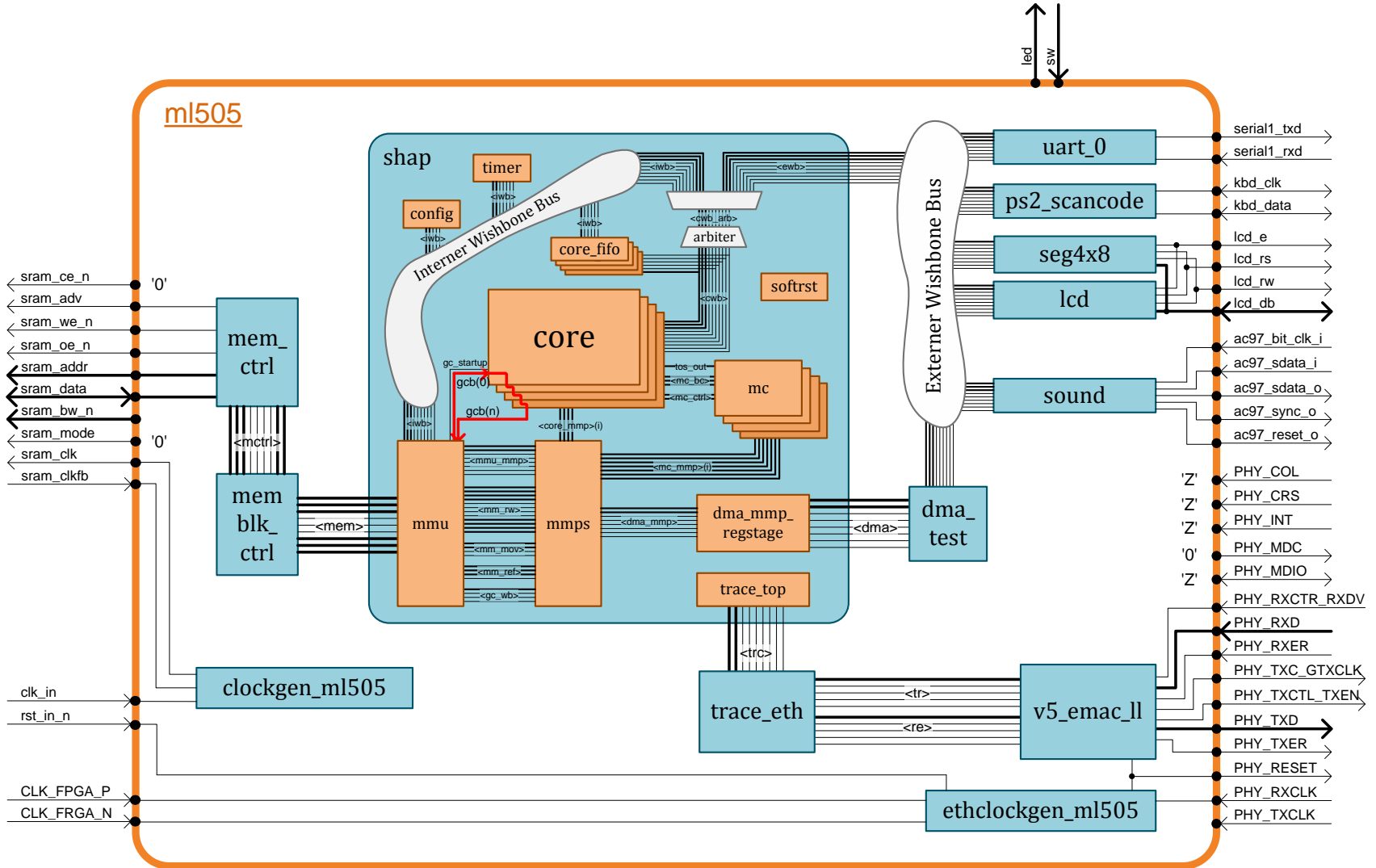
1.2 SHAP-Architektur

Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur

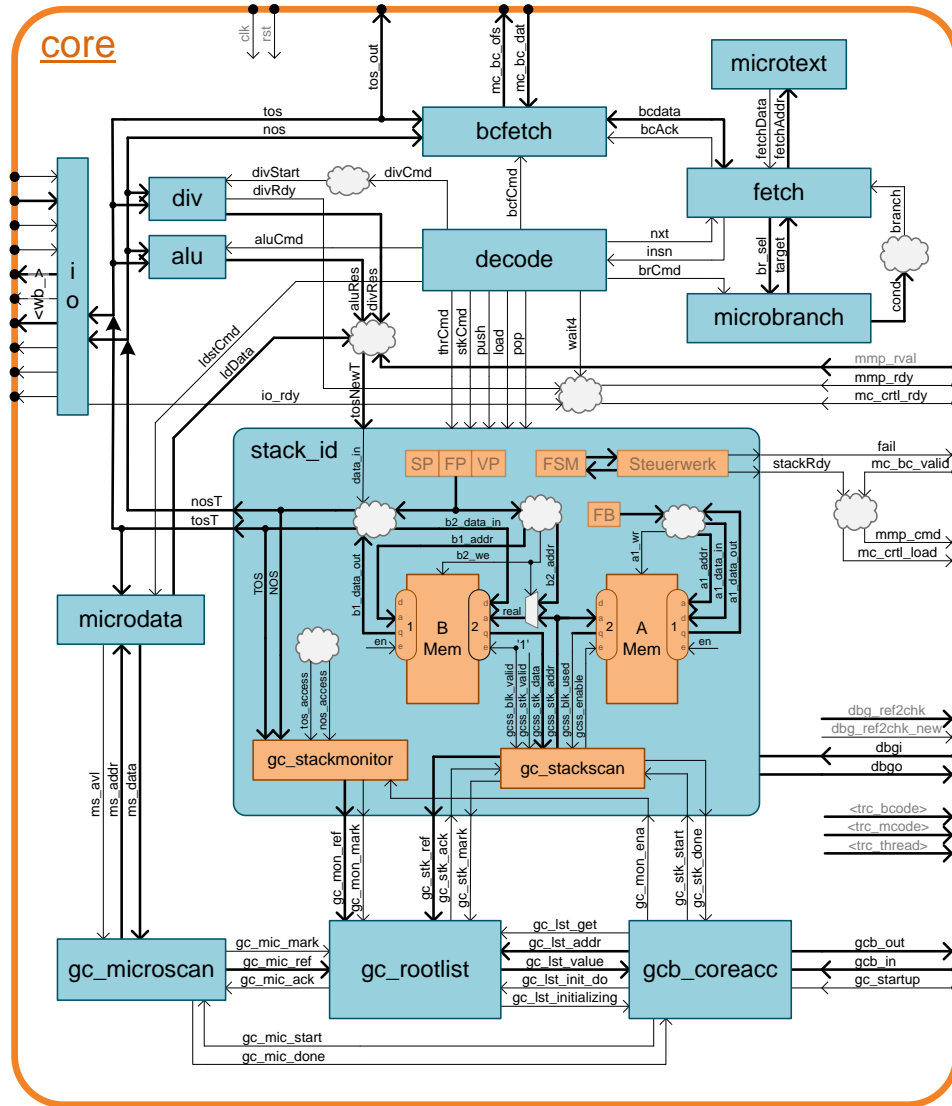


1.2 SHAP-Architektur

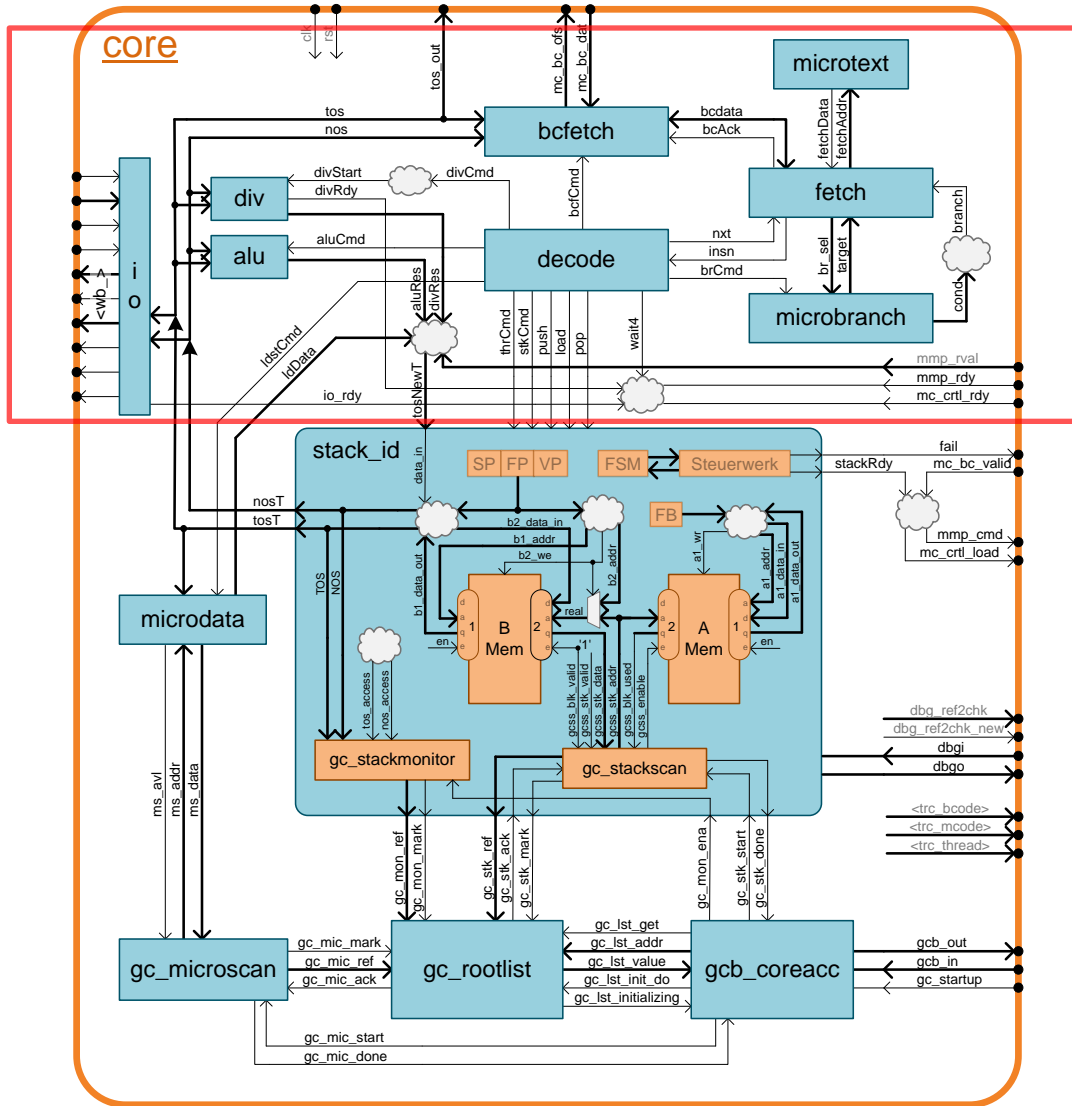
Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur



Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur



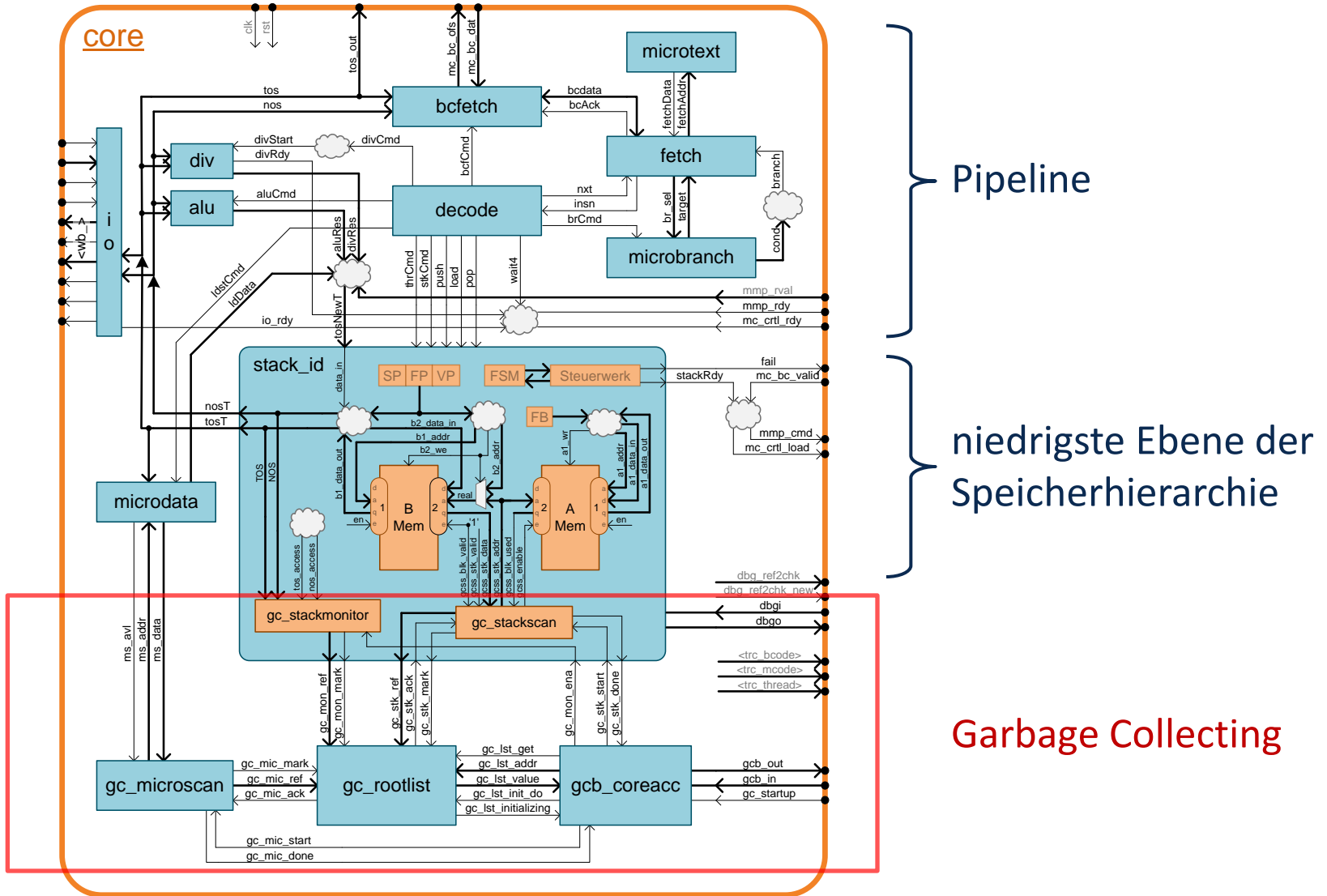
Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur



Pipeline

1.2 SHAP-Architektur

Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur



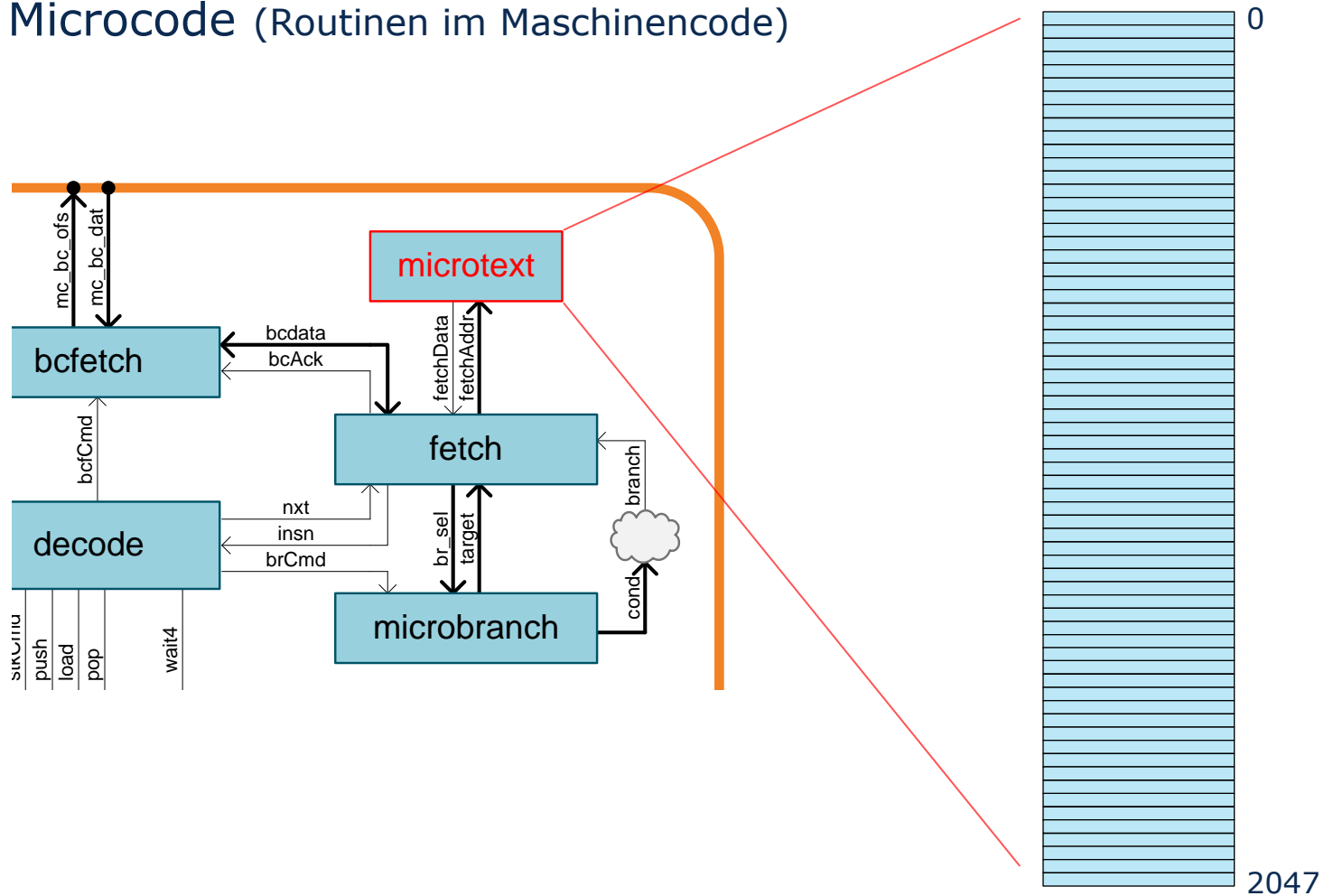
Sollkriterien

- geringe Umbauten an SHAP
 - geringer Scheduling-Overhead
 - geringer Migrations-Overhead
 - geringe Zusatzhardware
- } einfache Berechnungen

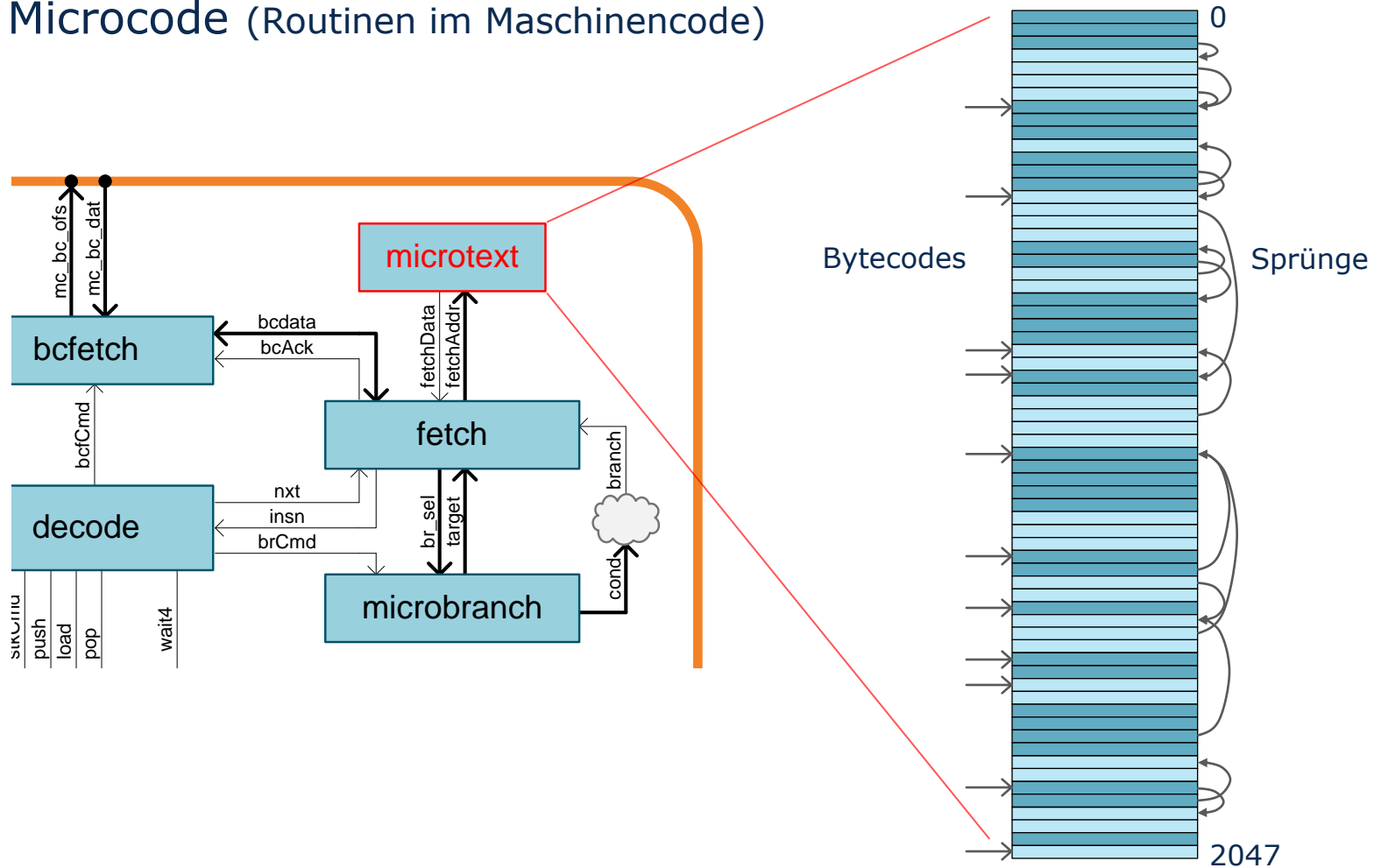
Kannkriterien

- ↳ Unterstützung der 10 Java-Prioritäten

Microcode (Routinen im Maschinencode)



Microcode (Routinen im Maschinencode)

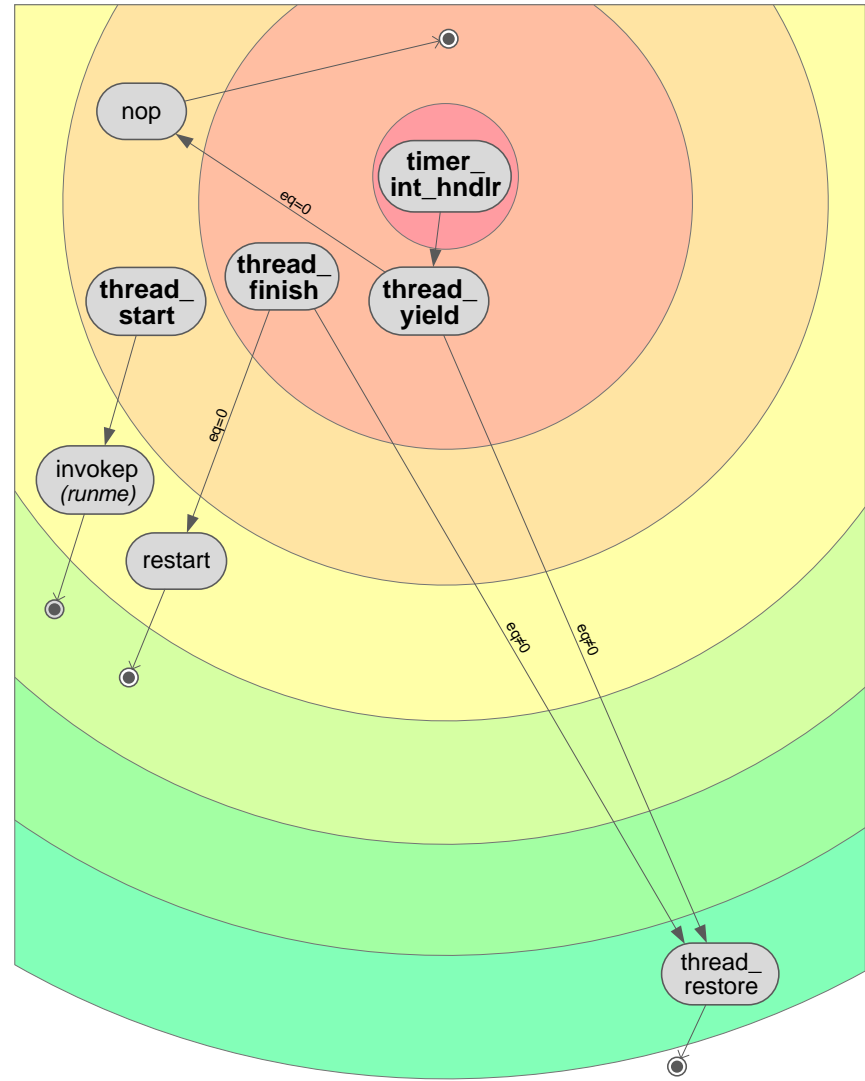


2.2 Scheduler

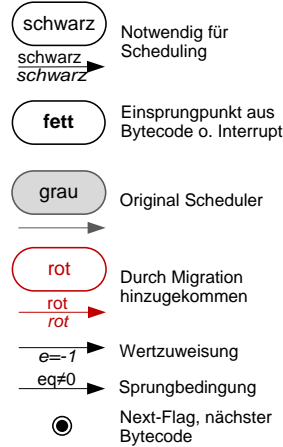
Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur

Original

- schwarz
Notwendig für Scheduling
- schwarz
→
schwarz
→
schwarz
- fett
Einsprungpunkt aus Bytecode o. Interrupt
- grau
Original Scheduler
-
- eq≠0
Sprungbedingung
- Next-Flag, nächster Bytecode
- ^
Bitw. XOR
- eq ... ct^ct.next

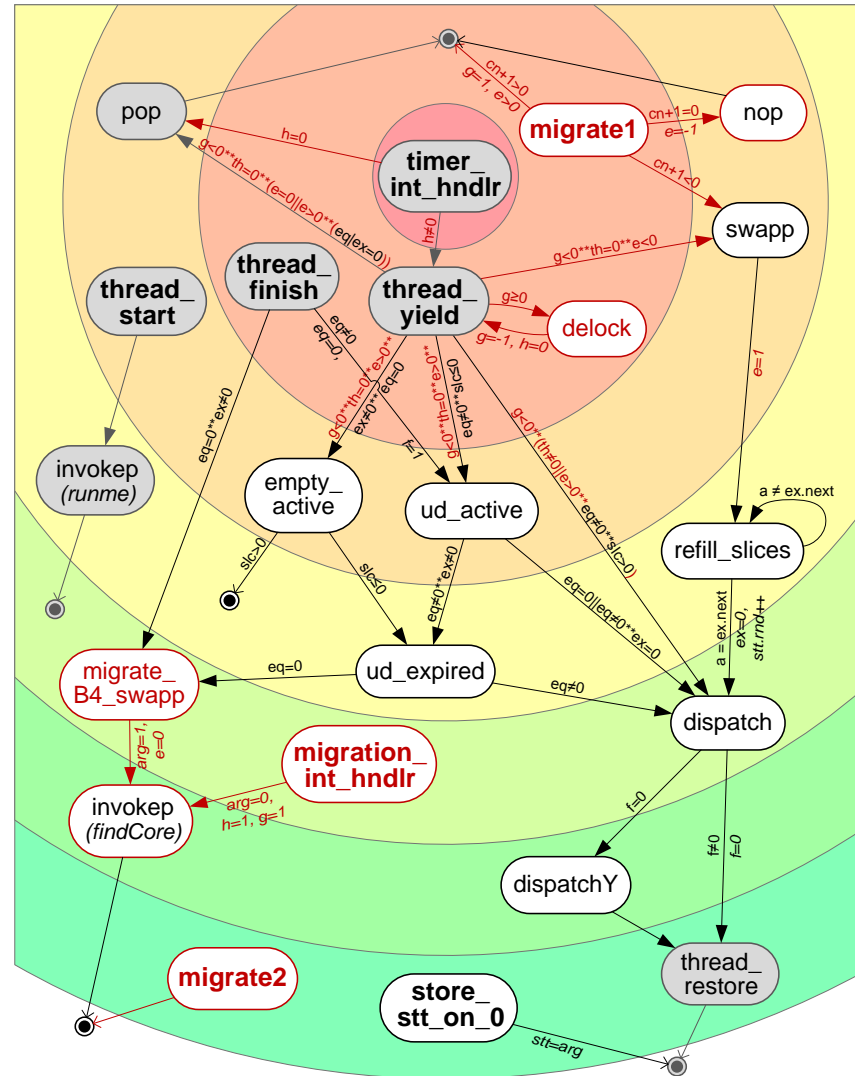


DWRR



** || Log. AND u. OR
 ^ | Bitw. XOR u. OR

slc ... verbleib. Slices
 stt ... StarterThread
 cn ... Kernnummer
 arg ... Argument
 eq ... $ct^{ct.next}$
 th ... Thread
 rnd ... Runde



Locks

GC-Lock (Garbage Collector nicht zeitgleich mit Migration)

Migration-Lock (max. eine Threadmigration gleichzeitig)

Yield-Lock (kein Threadwechsel während der Migration)

Locks

GC-Lock (Garbage Collector nicht zeitgleich mit Migration)

- keine Sicherungsmaßnahmen für Referenzen im Stack des Zielthreads
- lückenlose Stackübertragung ohne Adressierung möglich
 - **Lösung:** Deaktivierungsbefehl für ZPU über IWB

Migration-Lock (max. eine Threadmigration gleichzeitig)

Yield-Lock (kein Threadwechsel während der Migration)

Locks

GC-Lock (Garbage Collector nicht zeitgleich mit Migration)

- keine Sicherungsmaßnahmen für Referenzen im Stack des Zielthreads
- lückenlose Stackübertragung ohne Adressierung möglich
 - **Lösung:** Deaktivierungsbefehl für ZPU über IWB

Migration-Lock (max. eine Threadmigration gleichzeitig)

- keine Synchronisierung der Kern- und Threadauswahl
- keine Sicherungsmaßnahmen bei GC-Deaktivierung und GC-Bus-Kommunikation
- Vermeidung von Übersteuern
- lückenlose Stackübertragung ohne Adressierung möglich
 - **Lösung:** `static-synchronized`-Methode zum Setzen der Lock-Variable

Yield-Lock (kein Threadwechsel während der Migration)

Locks

GC-Lock (Garbage Collector nicht zeitgleich mit Migration)

- keine Sicherungsmaßnahmen für Referenzen im Stack des Zielthreads
- lückenlose Stackübertragung ohne Adressierung möglich
- **Lösung:** Deaktivierungsbefehl für ZPU über IWB

Migration-Lock (max. eine Threadmigration gleichzeitig)

- keine Synchronisierung der Kern- und Threadauswahl
- keine Sicherungsmaßnahmen bei GC-Deaktivierung und GC-Bus-Kommunikation
- Vermeidung von Übersteuern
- lückenlose Stackübertragung ohne Adressierung möglich
- **Lösung:** `static-synchronized`-Methode zum Setzen der Lock-Variable

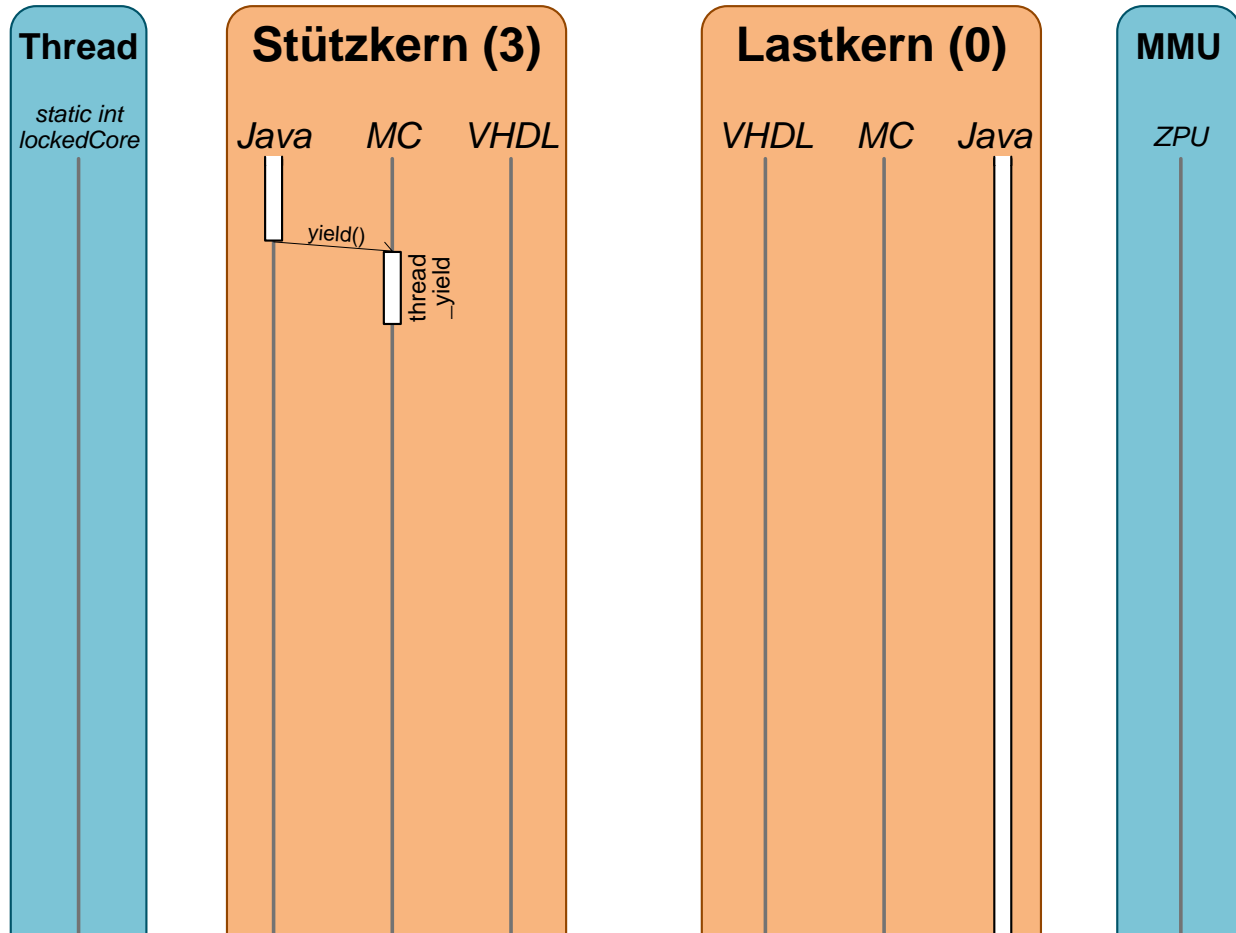
Yield-Lock (kein Threadwechsel während der Migration)

- um den noch unfertigen Threadwechsel nicht abubrechen
- benötigt bei: direkten und indirektem Aufruf der MC-Routine `thread_yield`
- gilt auch für Timer-Interrupt
- **Lösung:** Vermeidung von `yield()`, von `io_yield` im MC und Sprünge im MC anhand neuer Variablen

2.4 Konzept

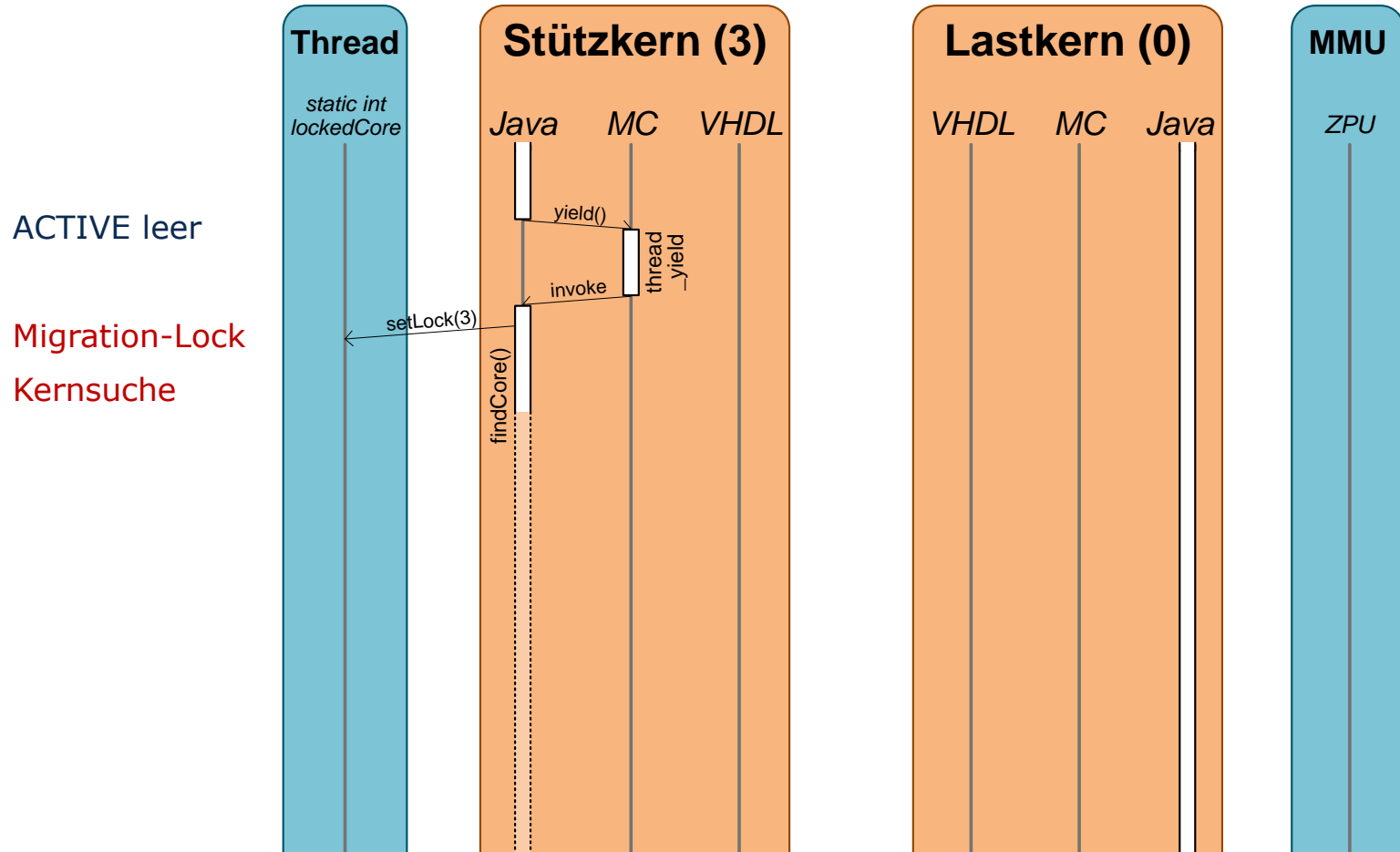
Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur

ACTIVE leer



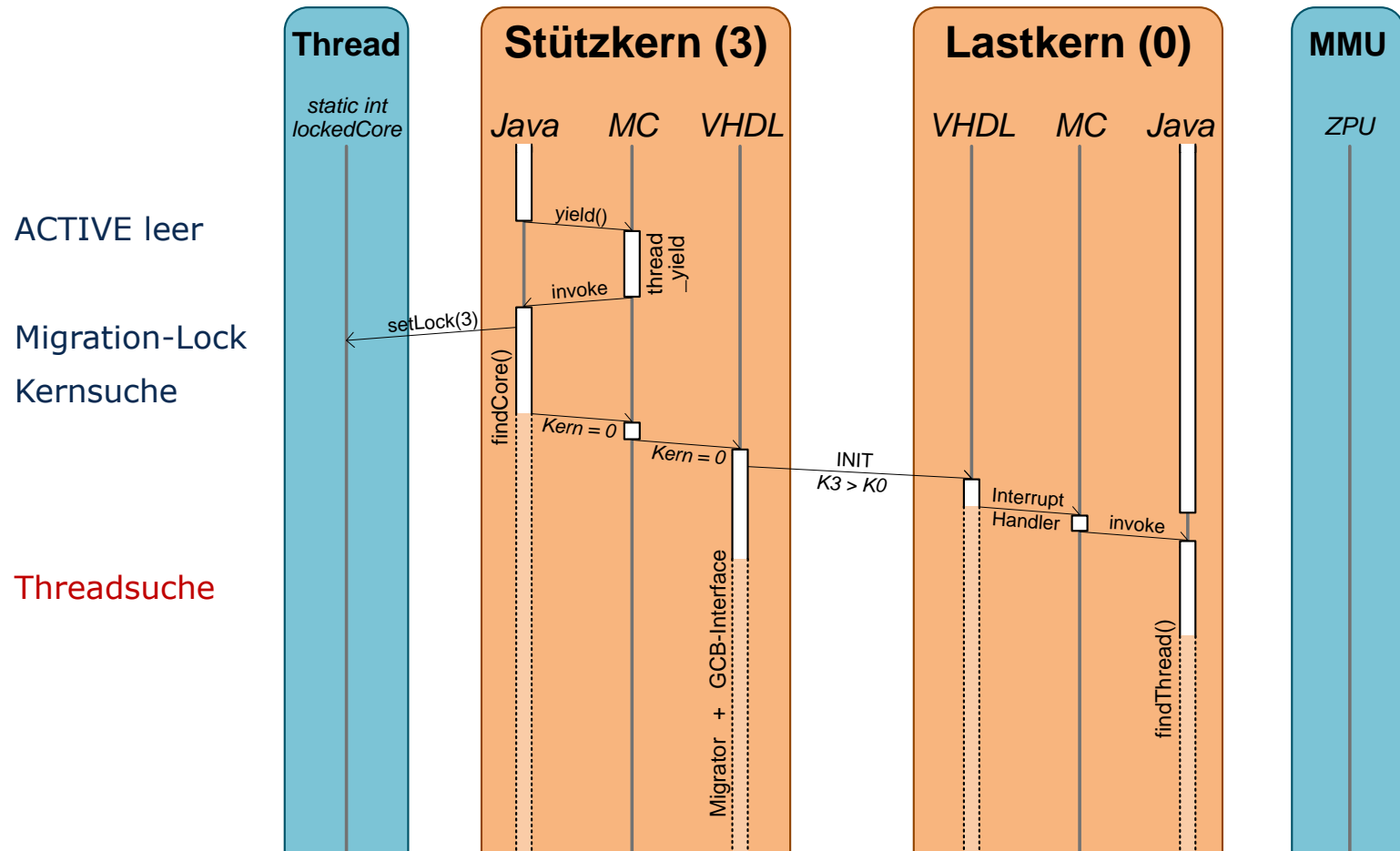
2.4 Konzept

Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur



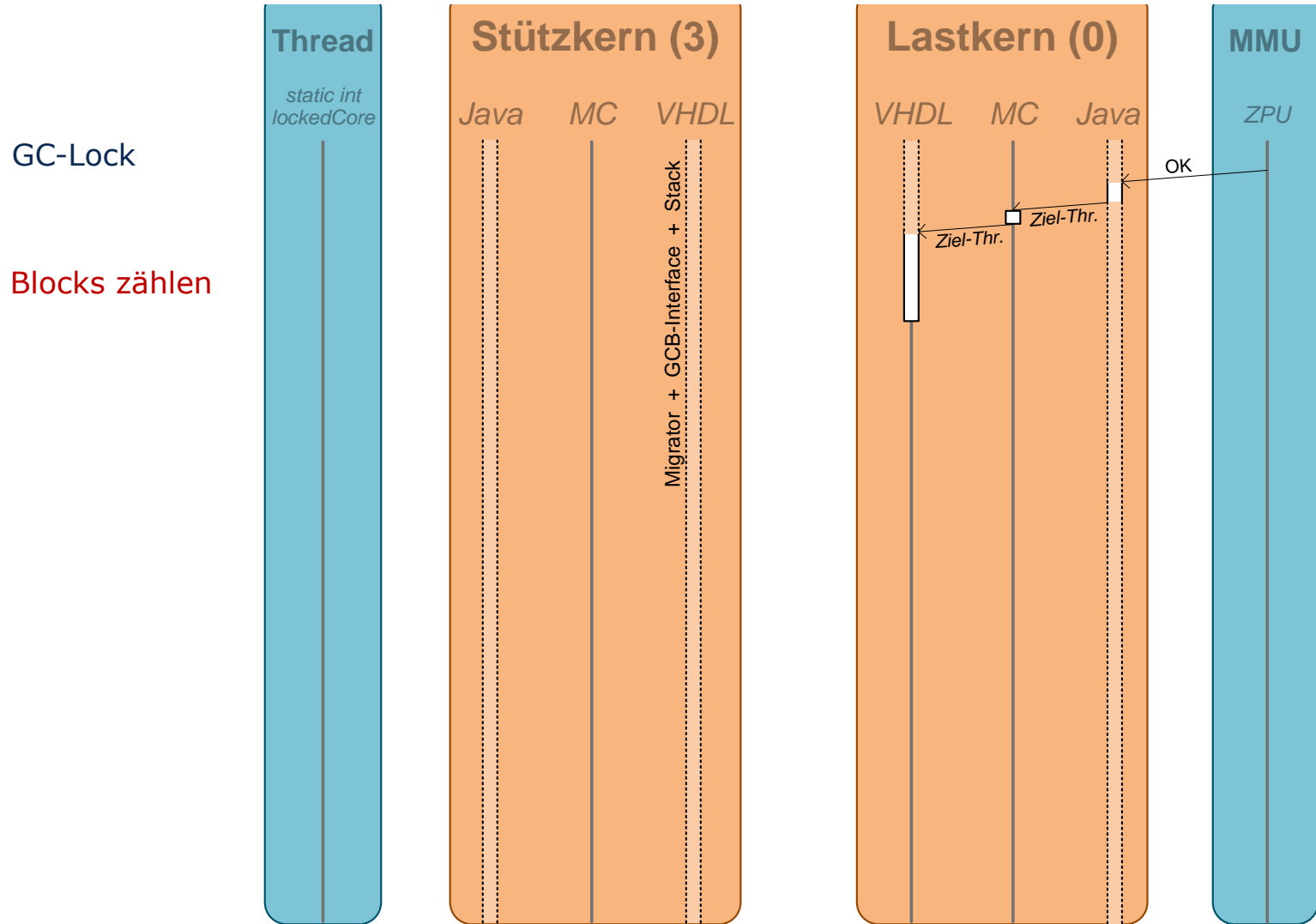
2.4 Konzept

Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur



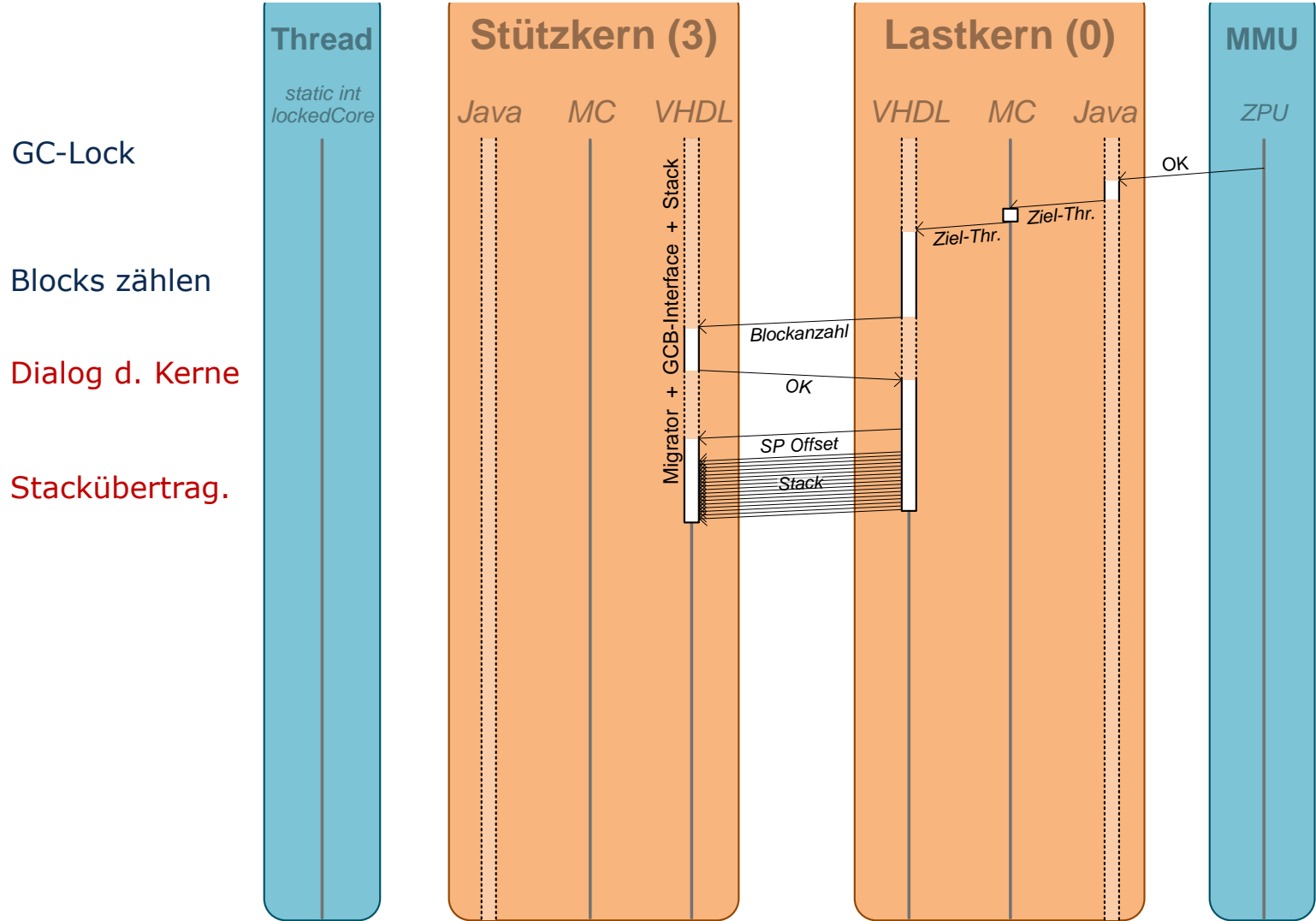
2.4 Konzept

Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur



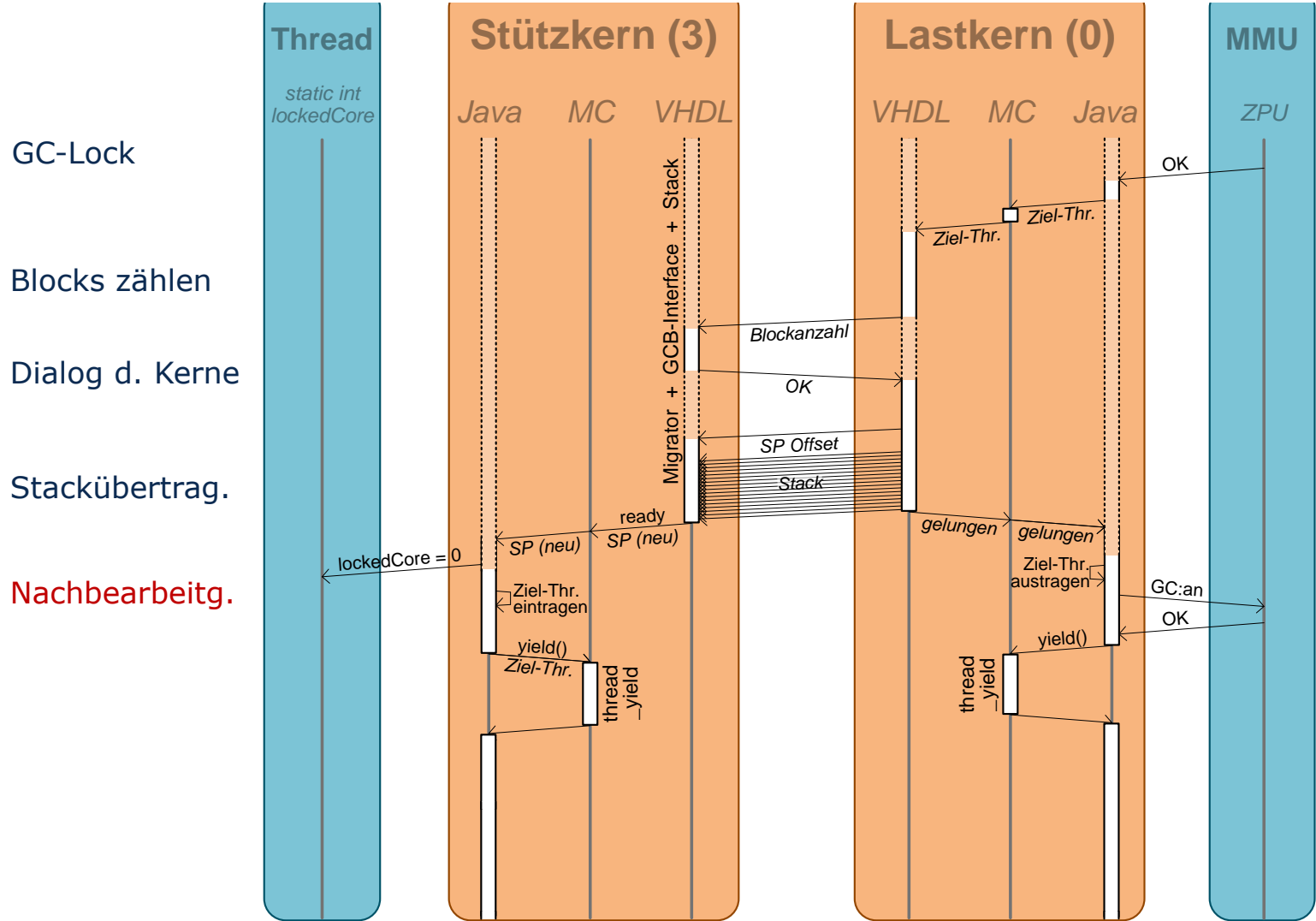
2.4 Konzept

Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur



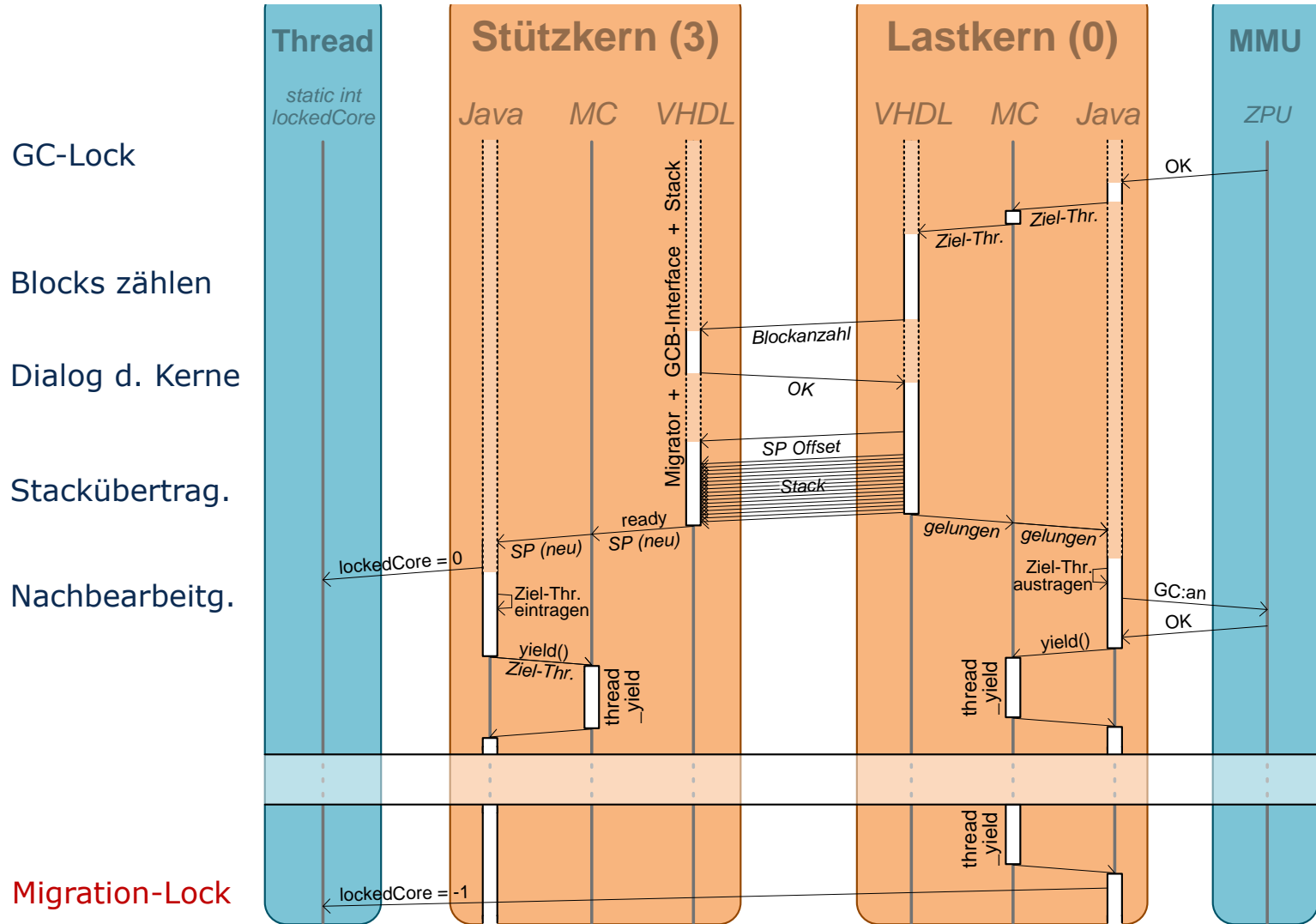
2.4 Konzept

Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur

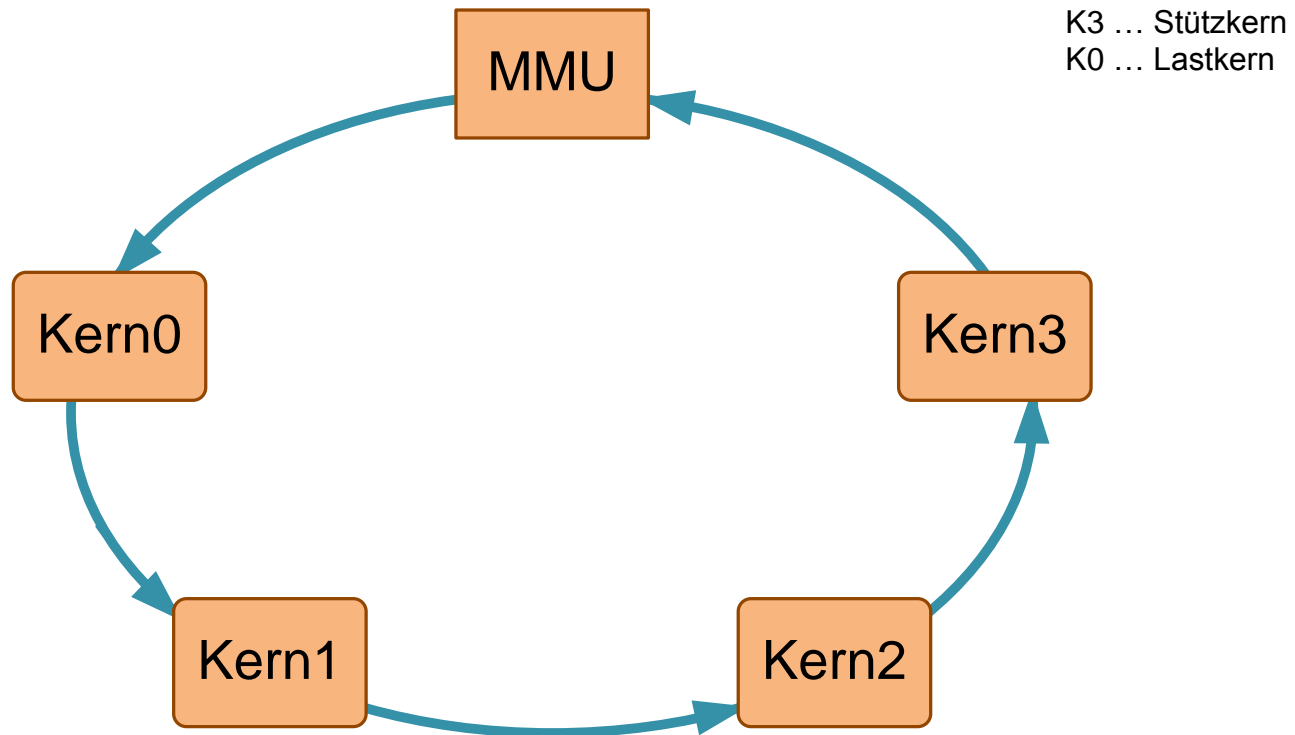


2.4 Konzept

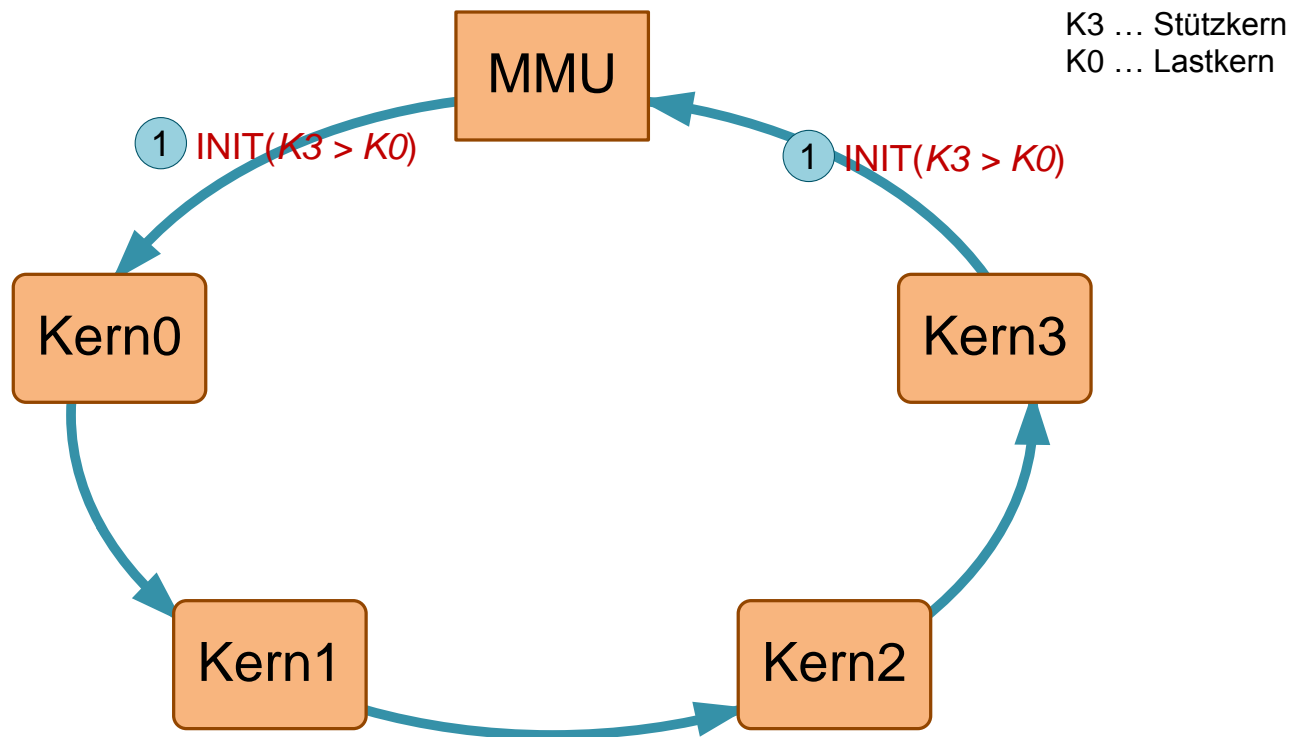
Implementierung eines Scheduling mit dynamischer Lastverteilung für die SHAP-Mehrkernarchitektur



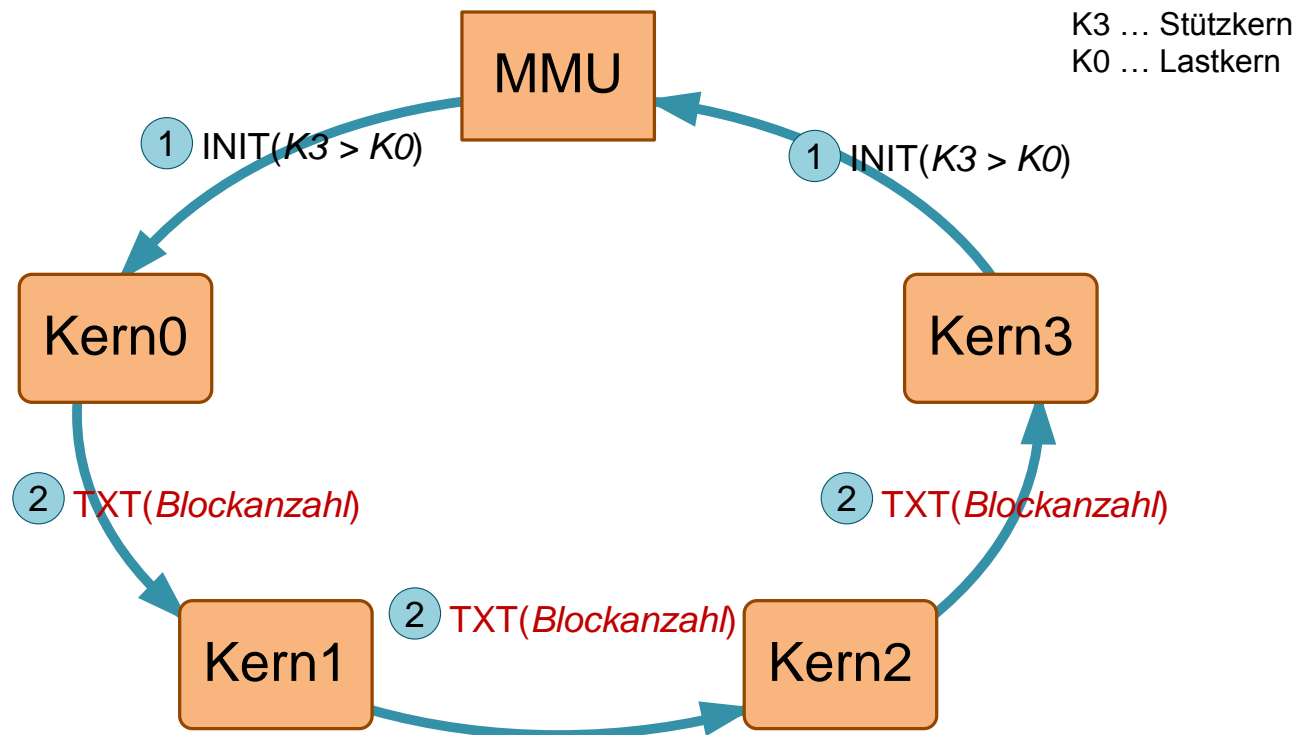
Kommunikation auf GC-Bus



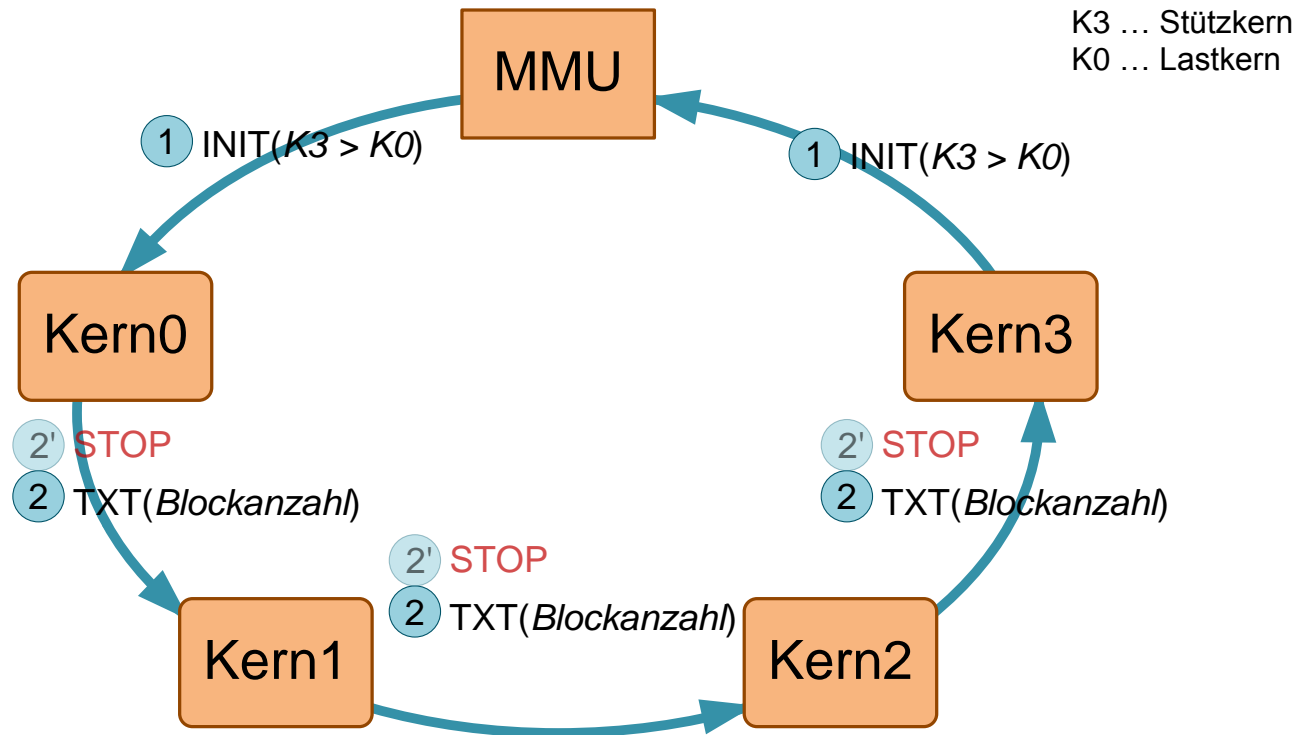
Kommunikation auf GC-Bus



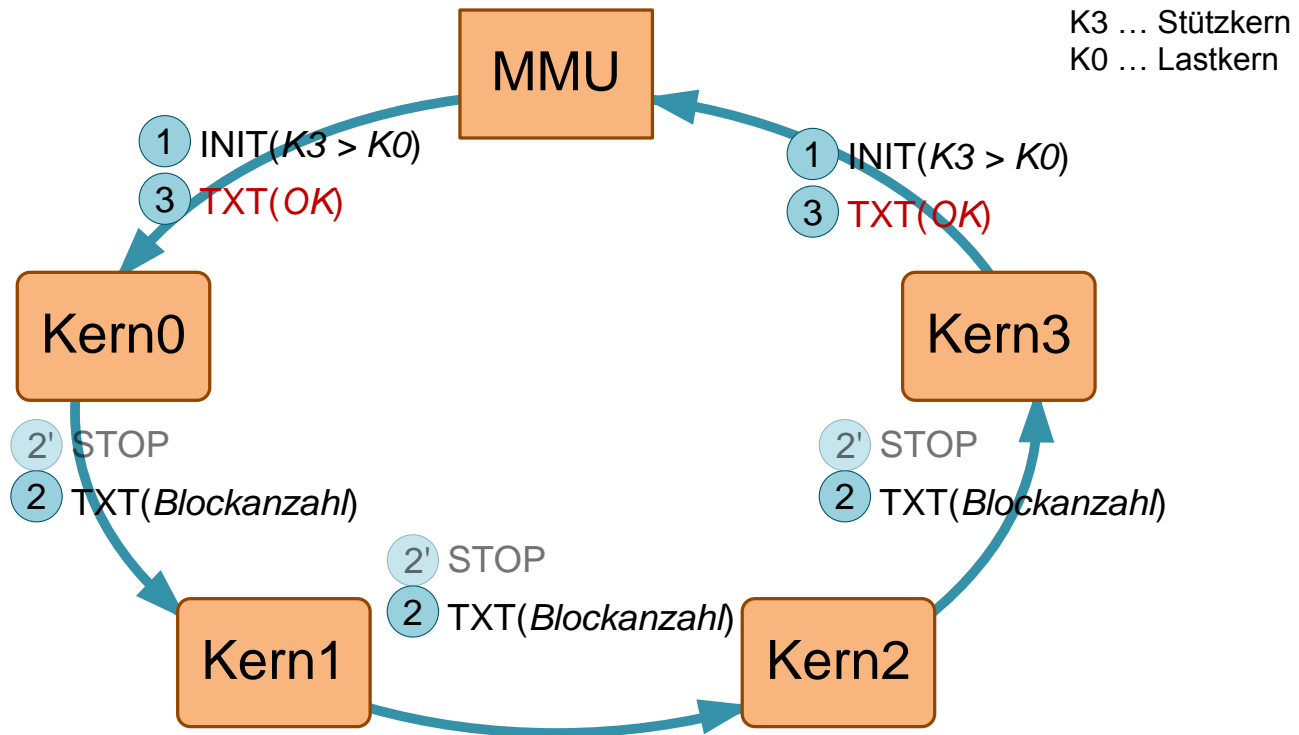
Kommunikation auf GC-Bus



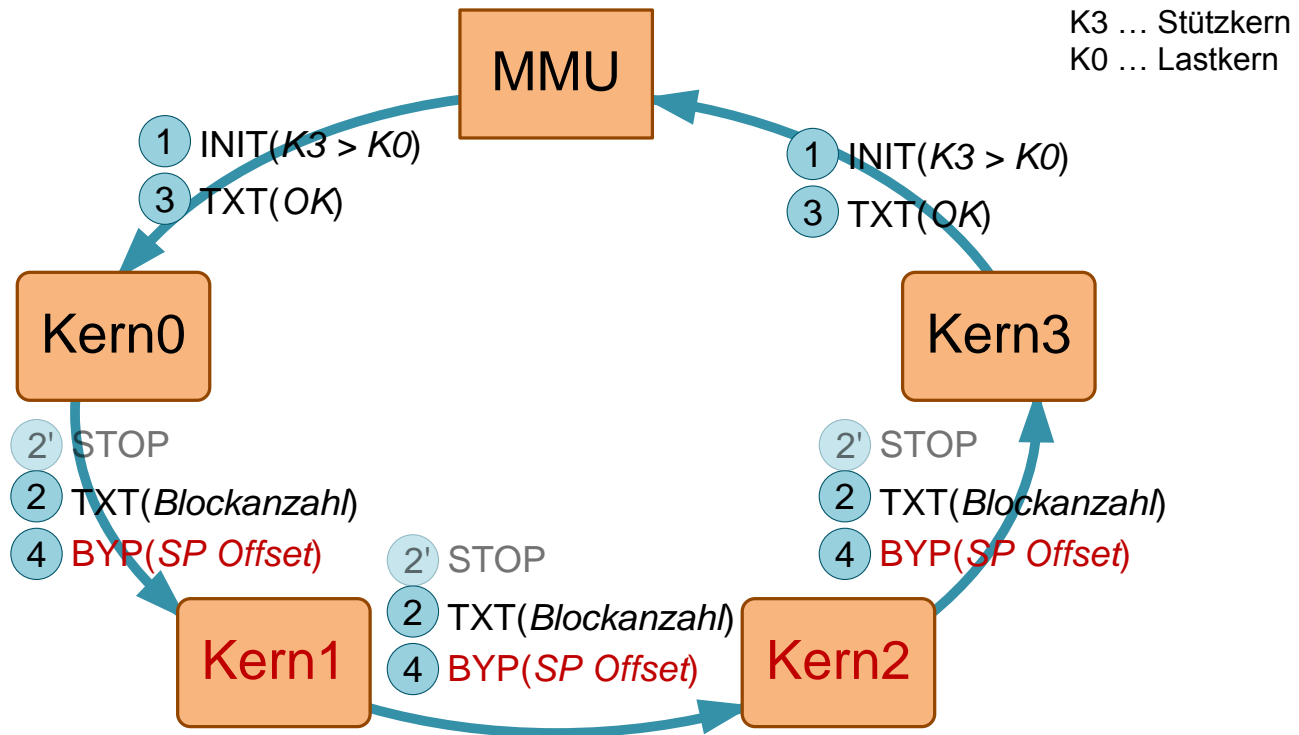
Kommunikation auf GC-Bus



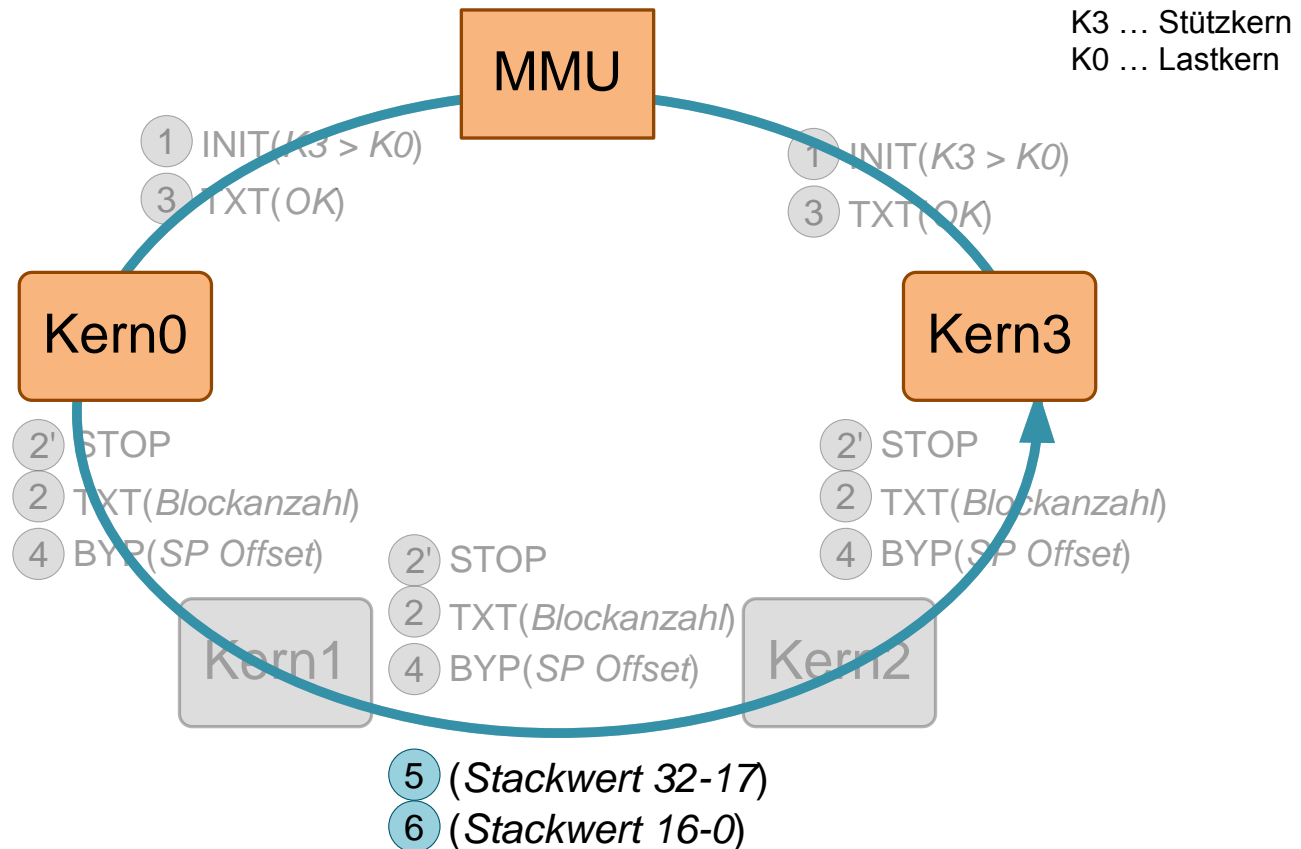
Kommunikation auf GC-Bus



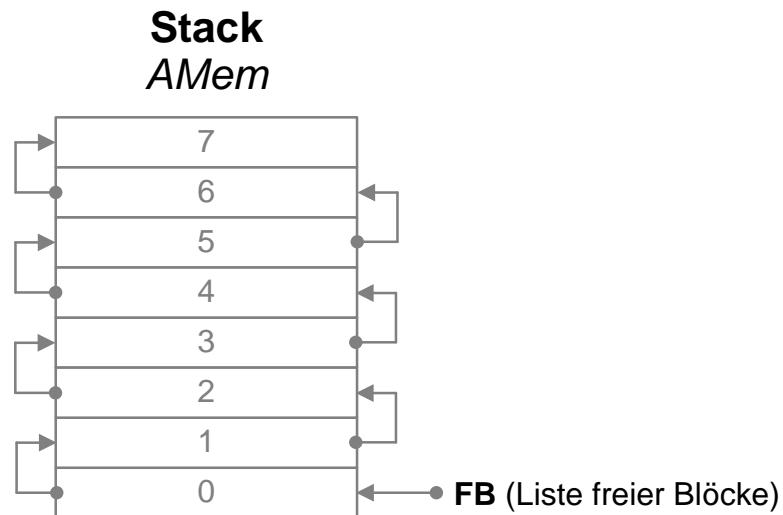
Kommunikation auf GC-Bus



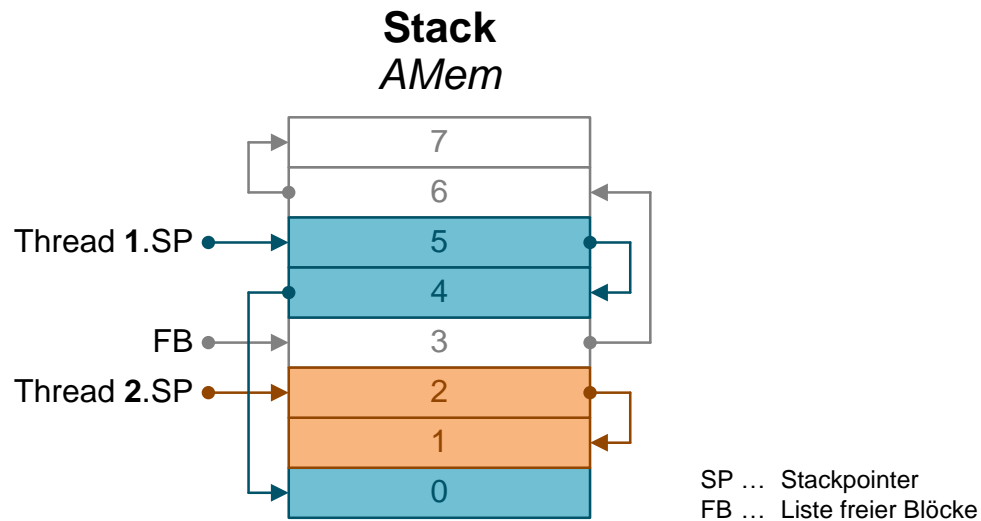
Kommunikation auf GC-Bus



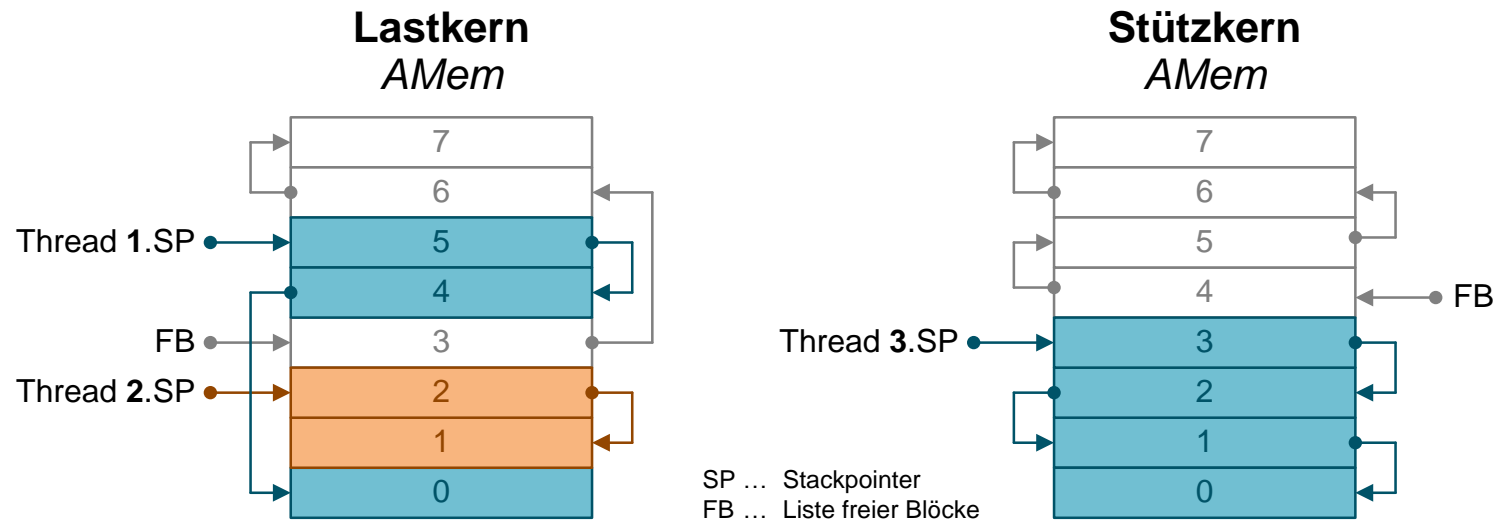
Stackübertragung



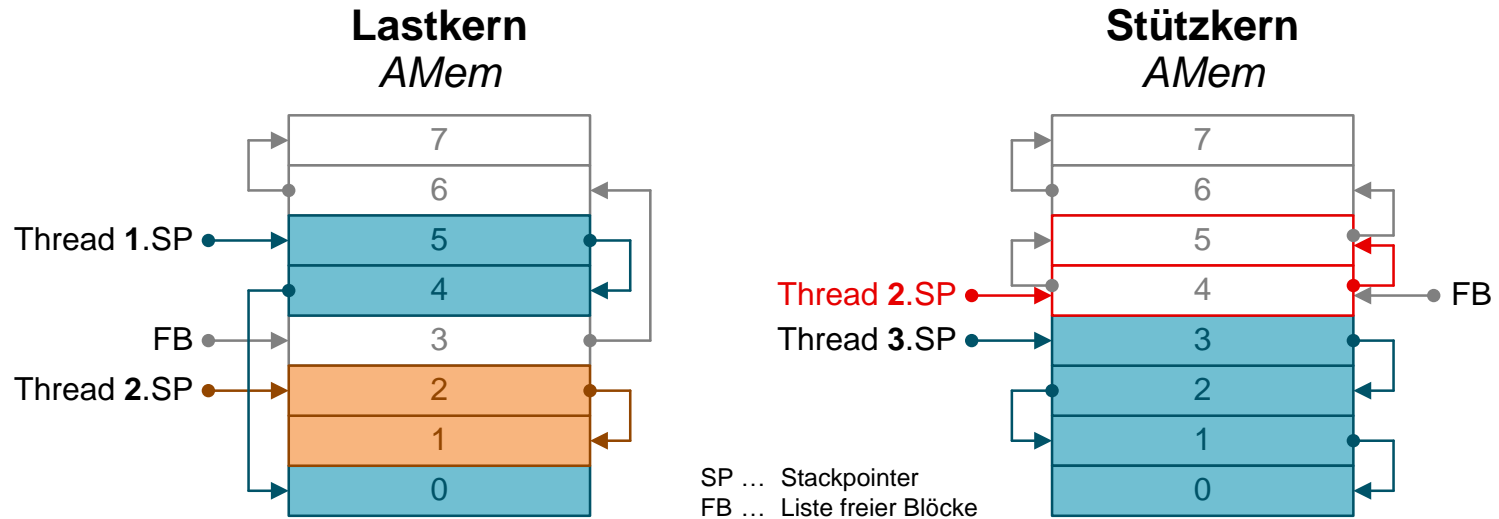
Stackübertragung



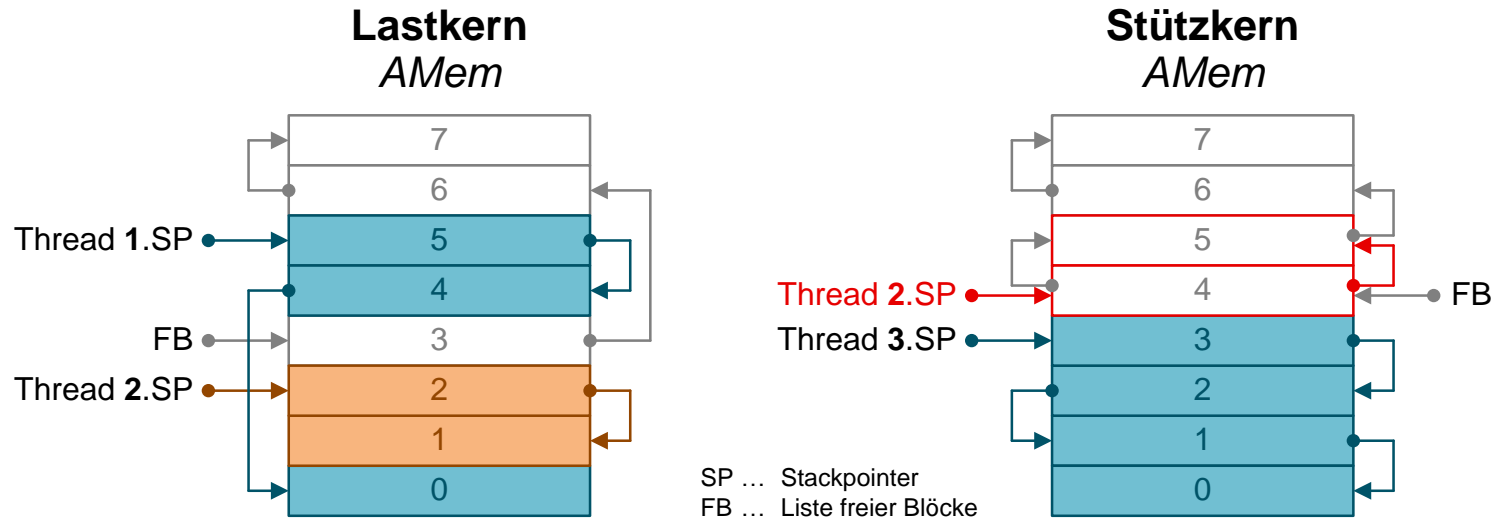
Stackübertragung



Stackübertragung

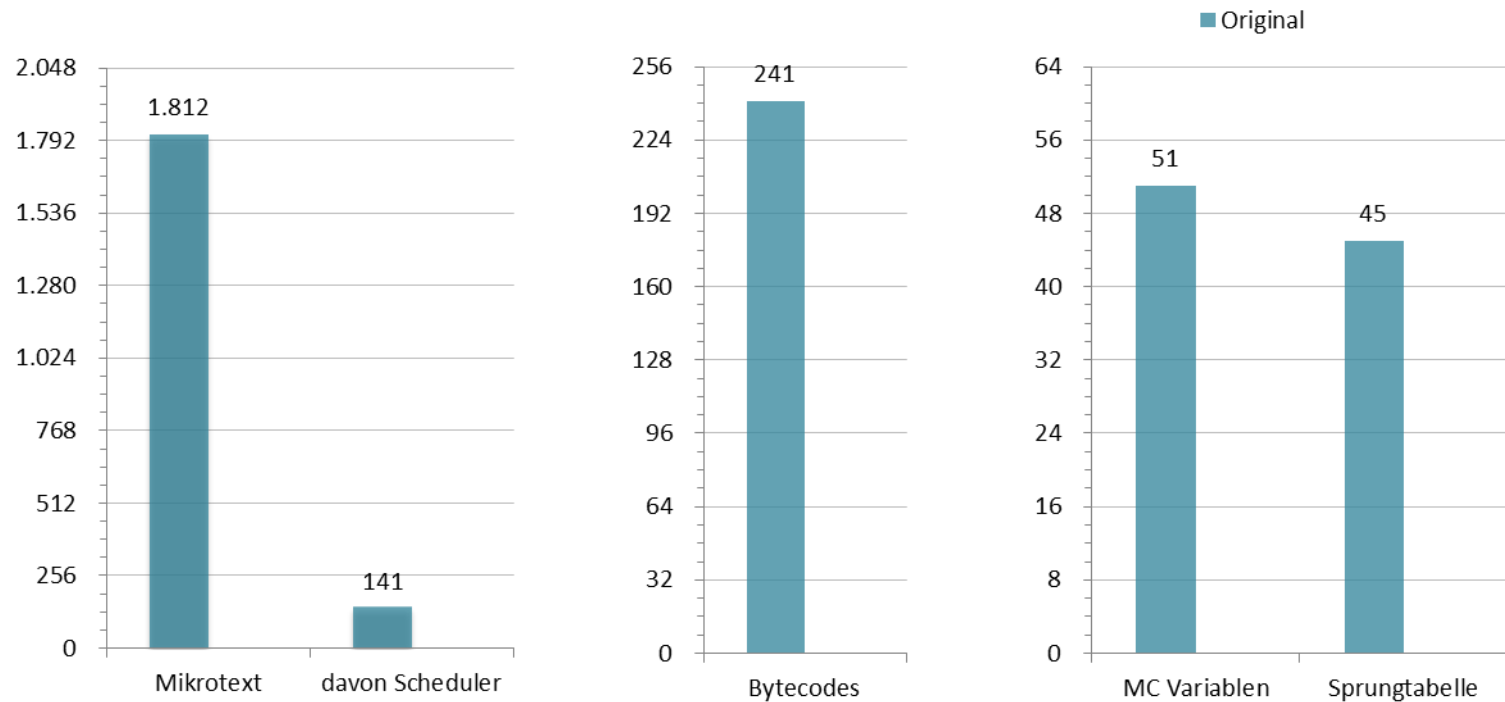


Stackübertragung



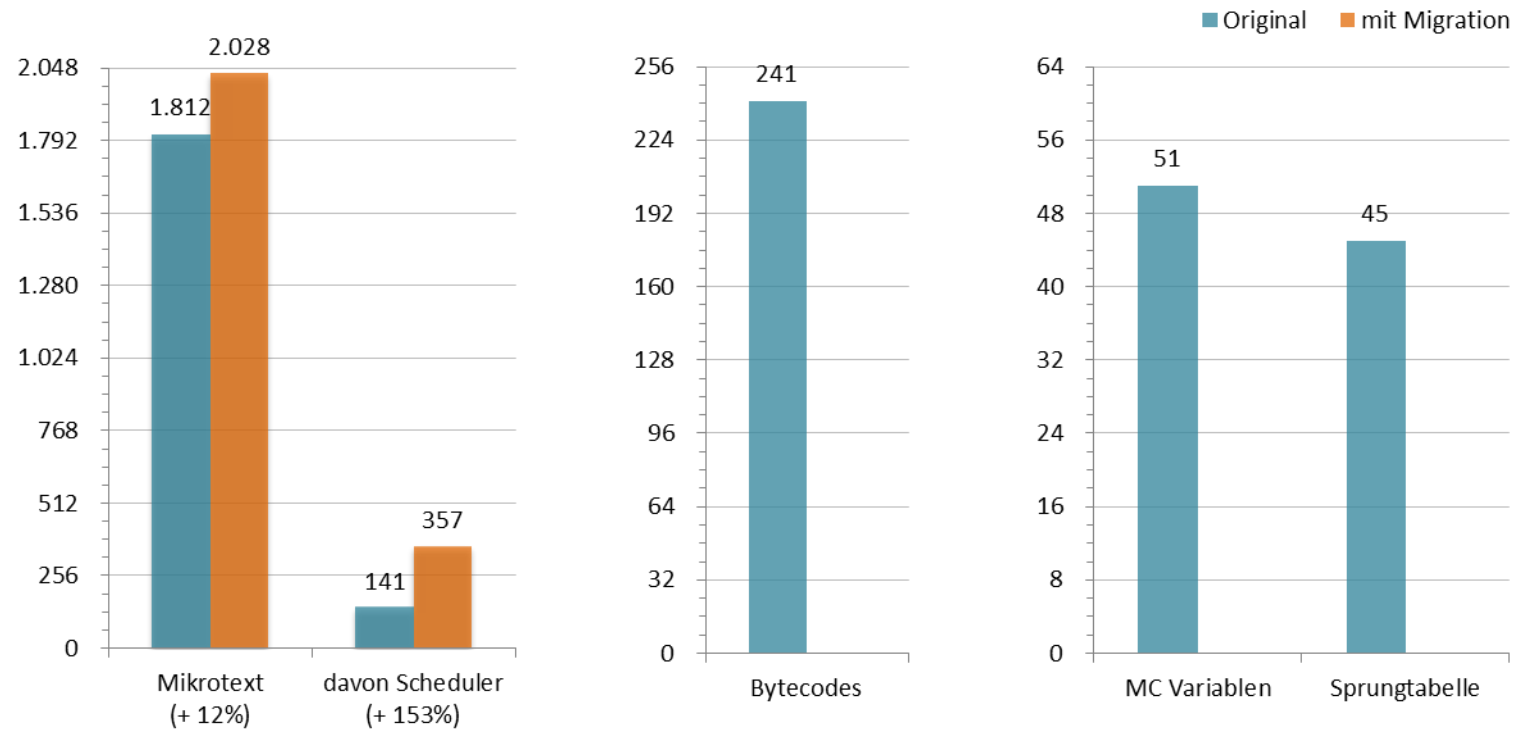
- Stackpointer anpassen
- Datenregister anpassen

Code



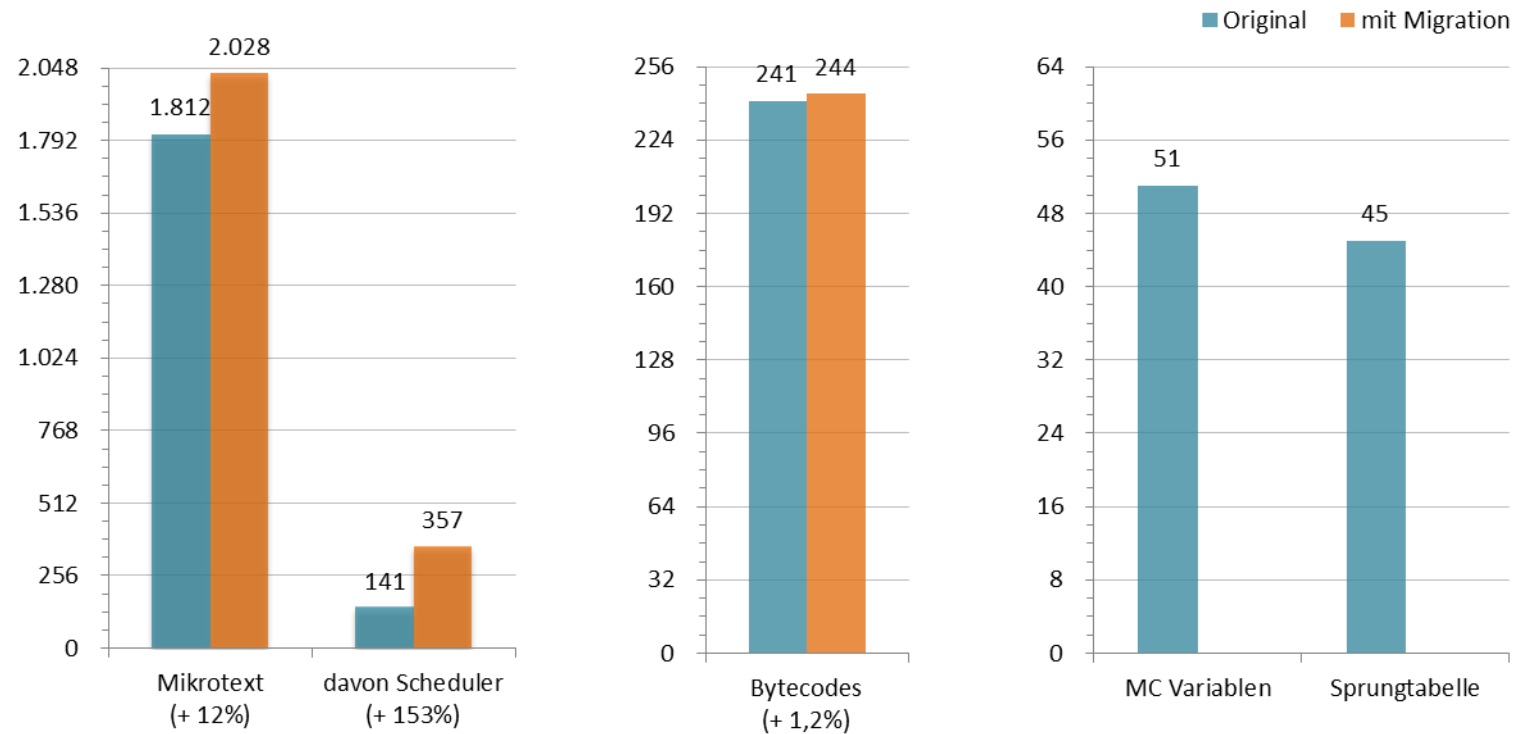
Mikrocode für Floating Point ist nicht enthalten

Code



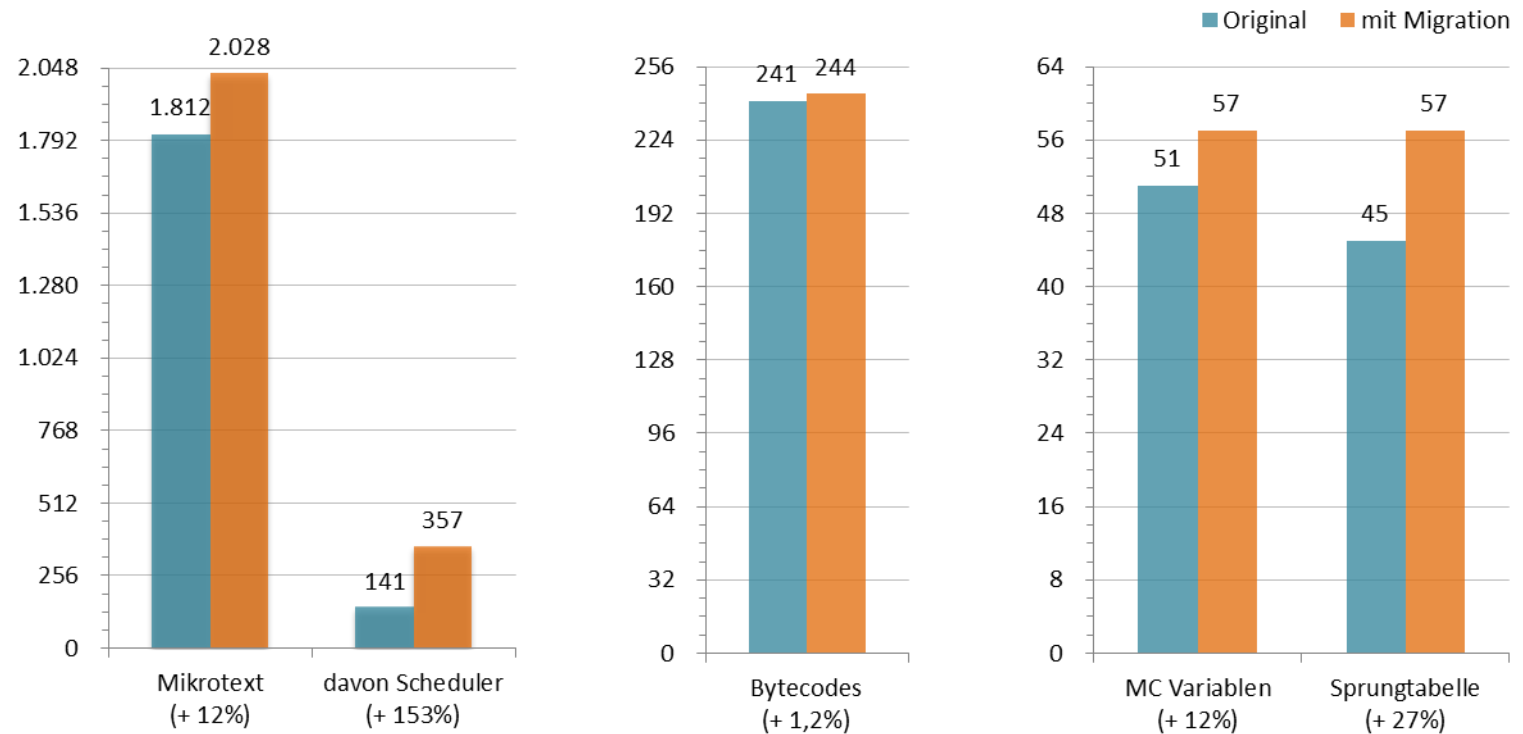
Mikrocode für Floating Point ist nicht enthalten

Code



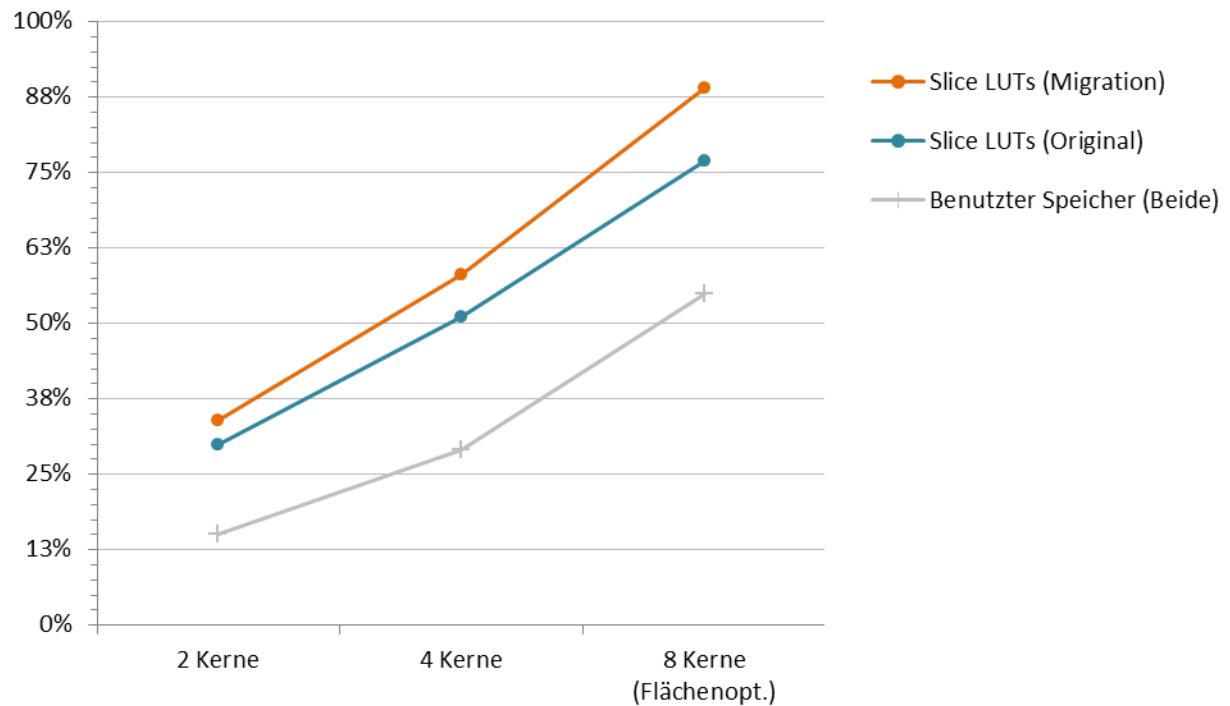
Mikrocode für Floating Point ist nicht enthalten

Code

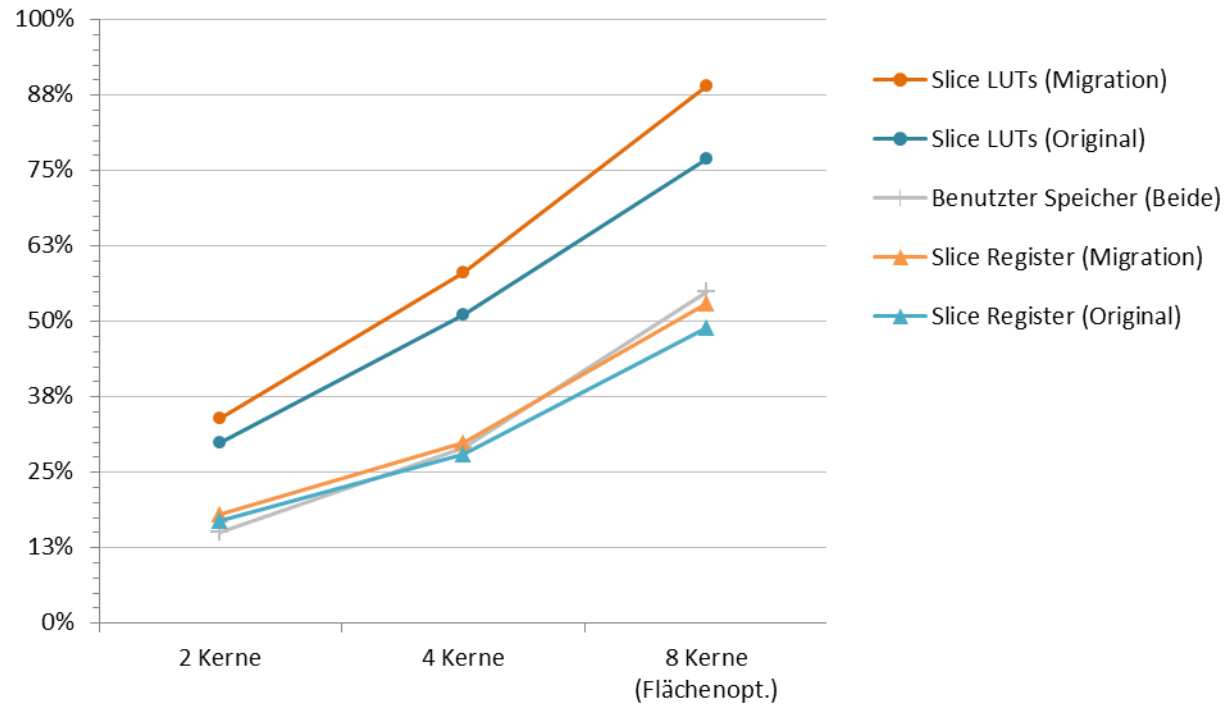


Mikrocode für Floating Point ist nicht enthalten

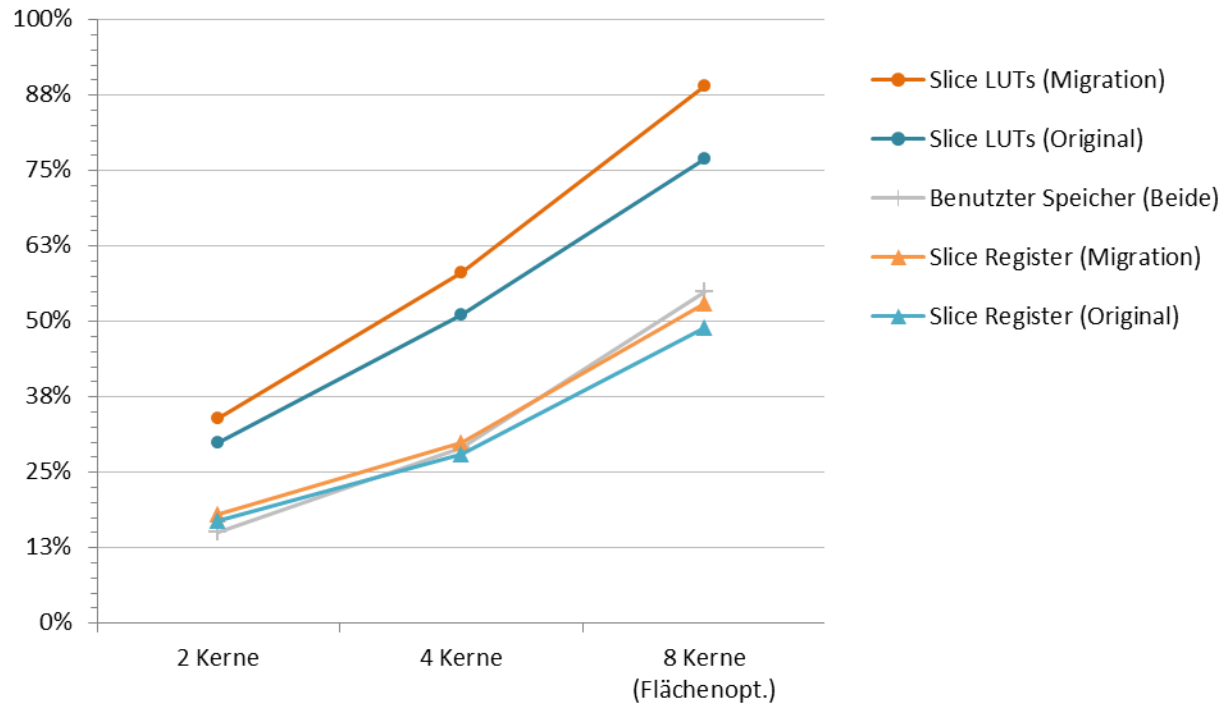
Hardware



Hardware



Hardware



$\text{Verbrauch}_{original} = 2848 \text{ LUTs} + \text{Kernanzahl} * 2989 \text{ LUTs}$
 $\text{Verbrauch}_{mit Migration} = 2992 \text{ LUTs} + \text{Kernanzahl} * 3471 \text{ LUTs}$

Keine umfassende Analyse und Bewertung sinnvoll

Abschätzung



Keine umfassende Analyse und Bewertung sinnvoll

Abschätzung

└─● Benchmarks

└─● JGF-Crypt (angepasst)

└─● JGF-SparseMatmult (angepasst)

Keine umfassende Analyse und Bewertung sinnvoll

Abschätzung

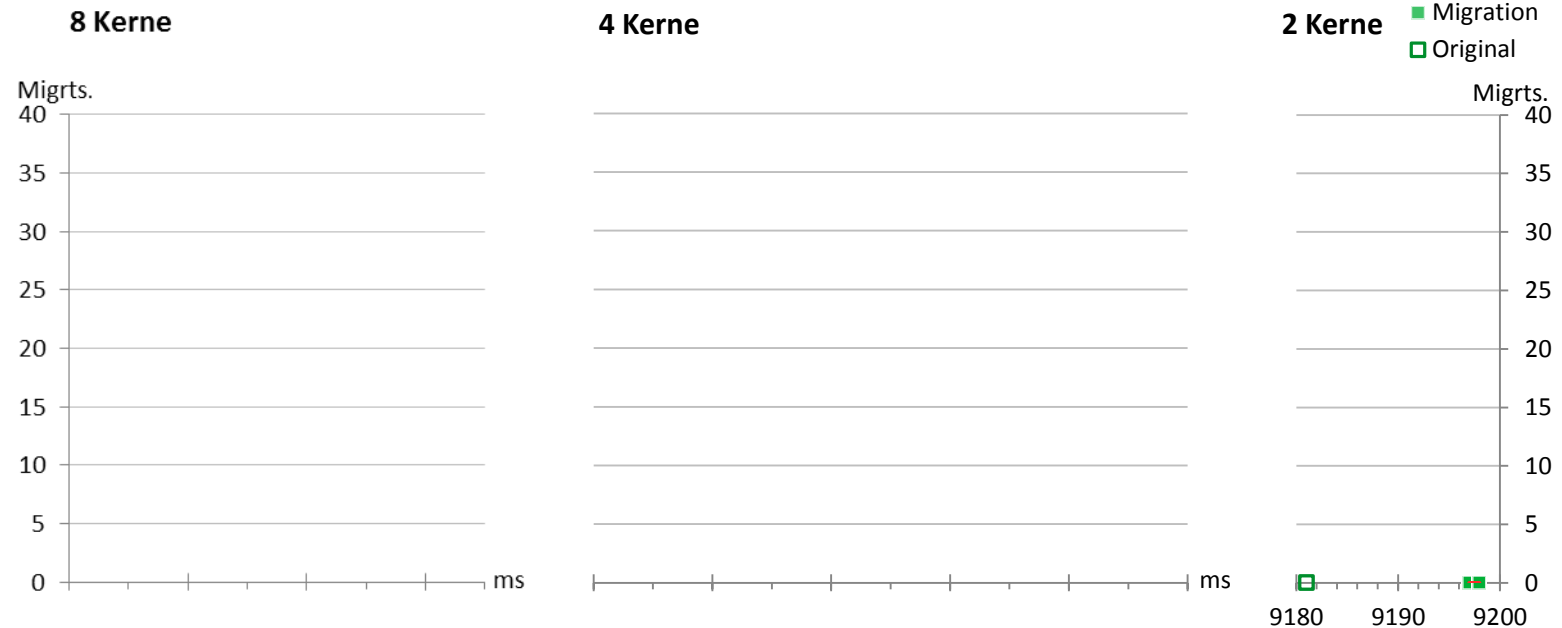
- Benchmarks
 - JGF-Crypt (angepasst)
 - JGF-SparseMatmult (angepasst)
- Auswahl-Strategien
 - Kern mit den meisten Threads
 - Thread mit den meisten verbleibenden Zeitscheiben
- Messung
 - Migrationsanzahl
 - mittlerer und maximaler Speed-Up

$$S_{p, mittel} = \frac{T_{Original, mittel}}{T_{Migration, mittel}}$$

$$S_{p, max} = \frac{T_{Original, max}}{T_{Migration, max}}$$

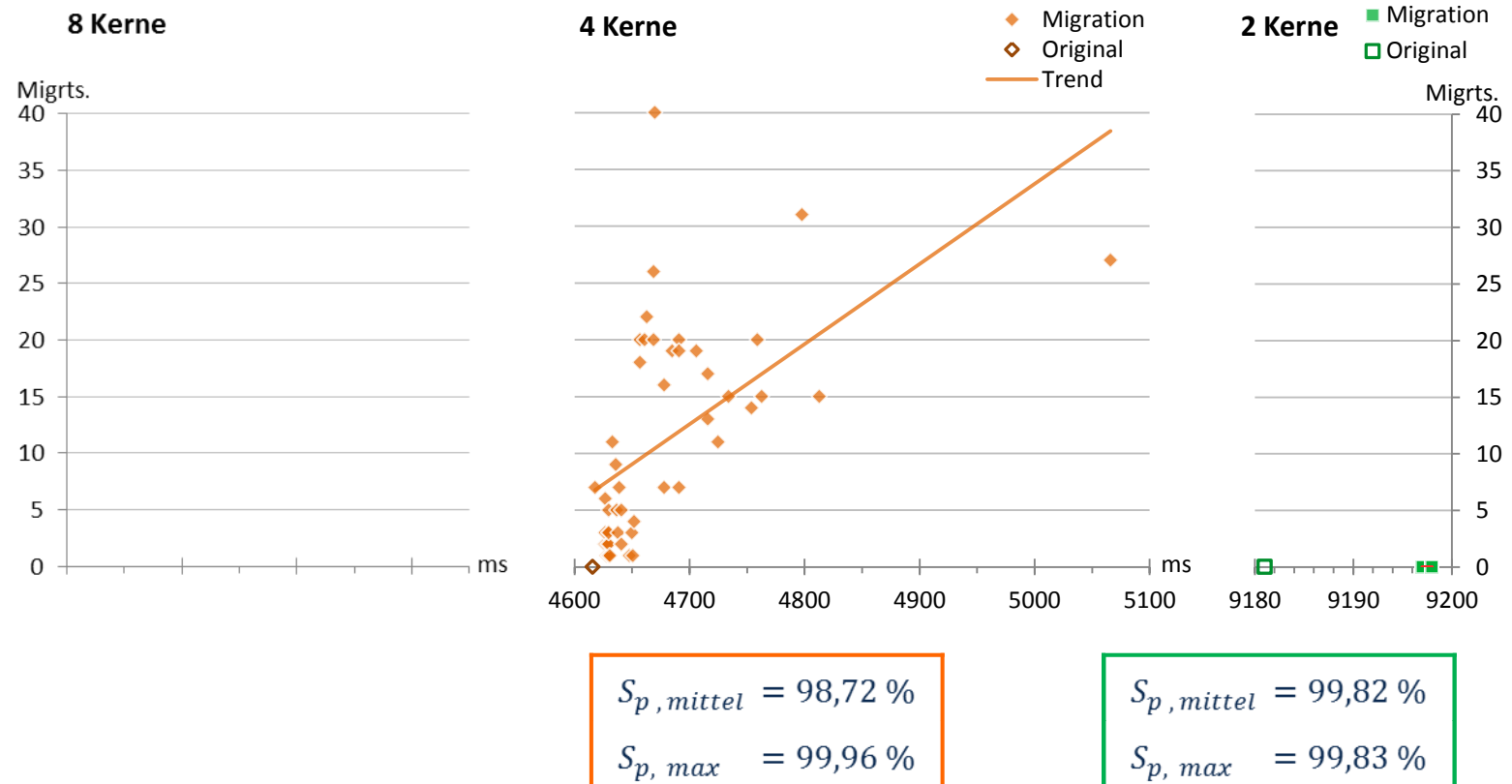
Ungünstigster Vergleich (ausgeglichenes System)

JGF-Crypt-Benchmark (Threadanzahl = 5 * Kernanzahl)



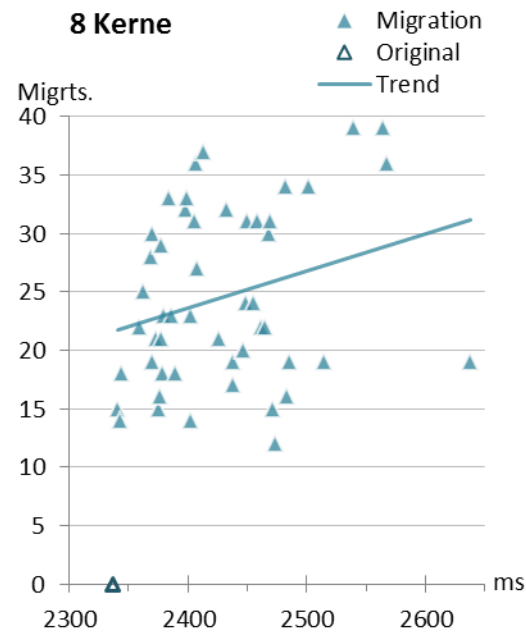
Ungünstigster Vergleich (ausgeglichenes System)

JGF-Crypt-Benchmark (Threadanzahl = 5 * Kernanzahl)



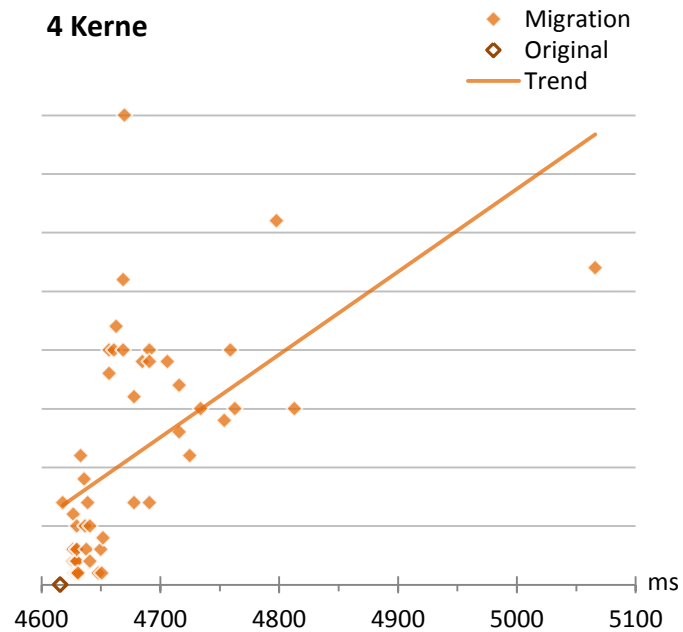
Ungünstigster Vergleich (ausgeglichenes System)

JGF-Crypt-Benchmark (Threadanzahl = 5 * Kernanzahl)



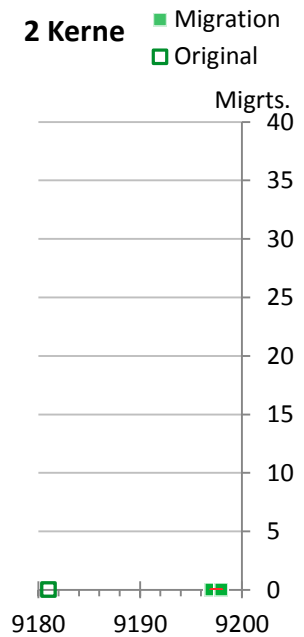
$$S_{p, mittel} = 96,20 \%$$

$$S_{p, max} = 99,83 \%$$



$$S_{p, mittel} = 98,72 \%$$

$$S_{p, max} = 99,96 \%$$



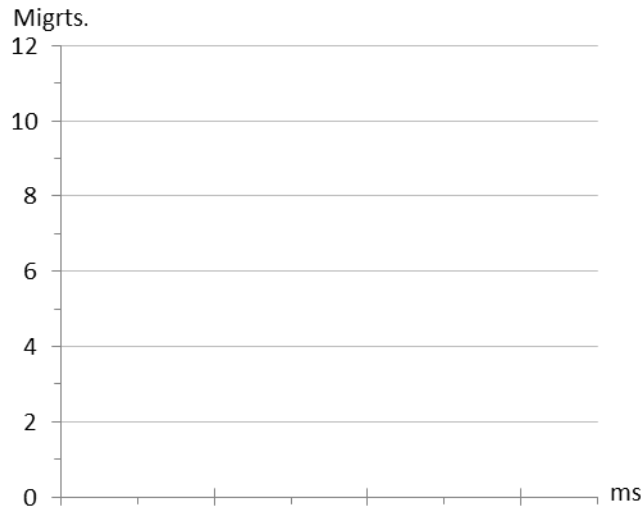
$$S_{p, mittel} = 99,82 \%$$

$$S_{p, max} = 99,83 \%$$

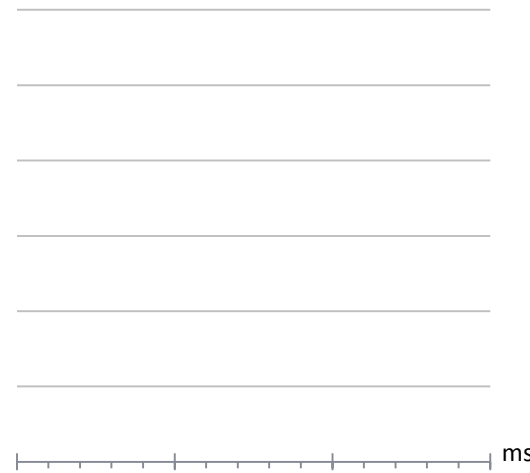
Ungünstigster Vergleich (ausgeglichenes System)

JGF-SparseMatmult-Benchmark (Threadanzahl = 10 * Kernanzahl)

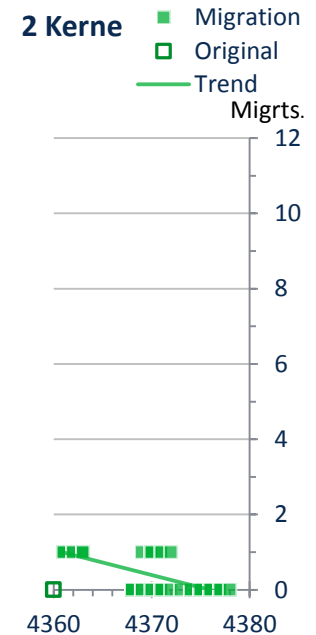
8 Kerne



4 Kerne



2 Kerne



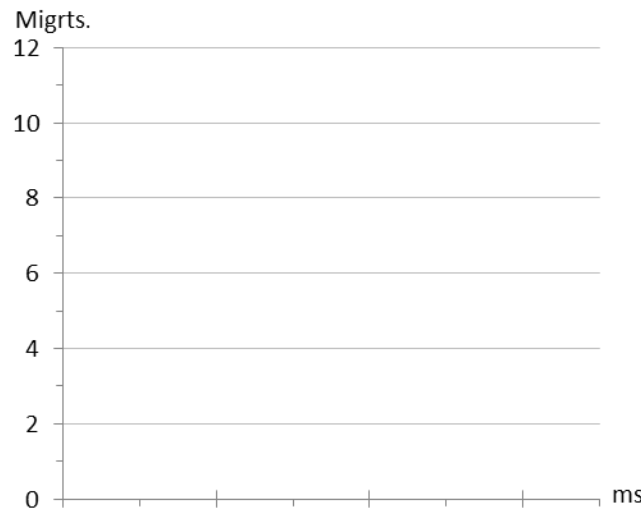
$$S_{p, mittel} = 99,74 \%$$

$$S_{p, max} = 99,98 \%$$

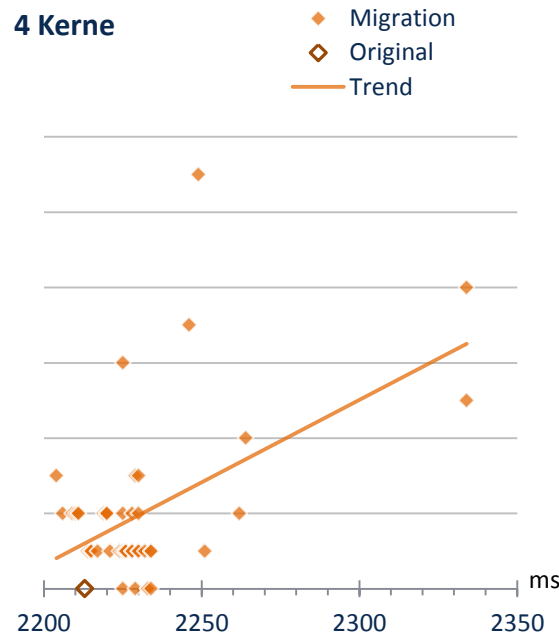
Ungünstigster Vergleich (ausgeglichenes System)

JGF-SparseMatmult-Benchmark (Threadanzahl = 10 * Kernanzahl)

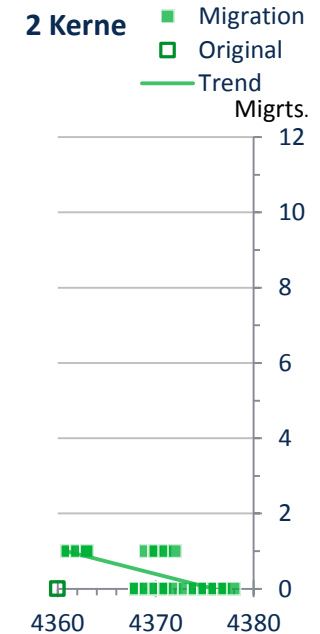
8 Kerne



4 Kerne



2 Kerne



$$S_{p, mittel} = 99,20 \%$$

$$S_{p, max} = 100,41 \%$$

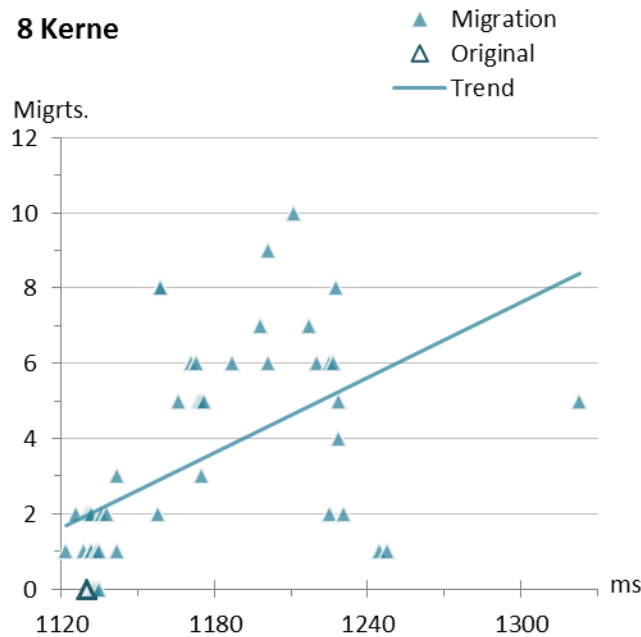
$$S_{p, mittel} = 99,74 \%$$

$$S_{p, max} = 99,98 \%$$

Ungünstigster Vergleich (ausgeglichenes System)

JGF-SparseMatmult-Benchmark (Threadanzahl = 10 * Kernanzahl)

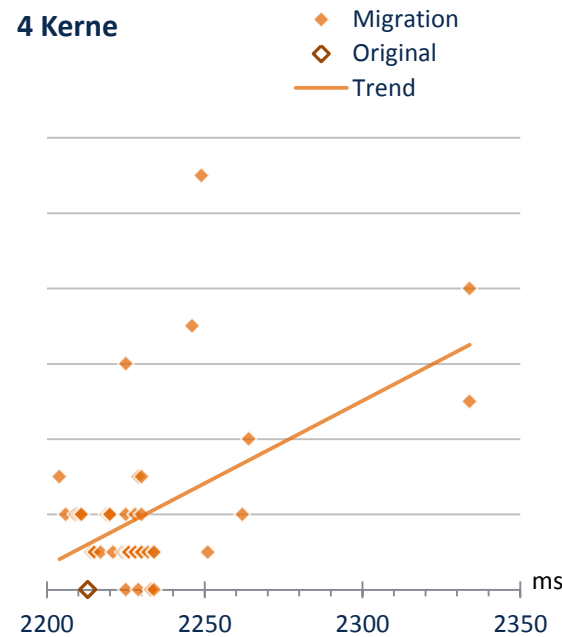
8 Kerne



$$S_{p, mittel} = 96,31 \%$$

$$S_{p, max} = 100,71 \%$$

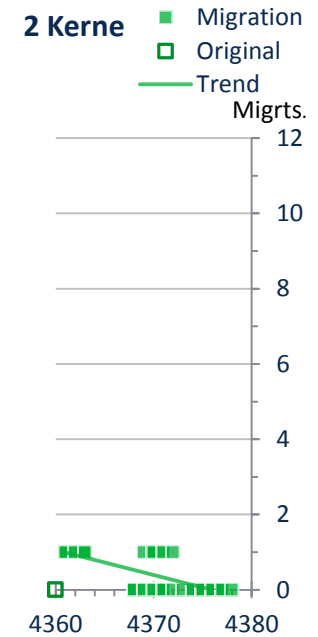
4 Kerne



$$S_{p, mittel} = 99,20 \%$$

$$S_{p, max} = 100,41 \%$$

2 Kerne

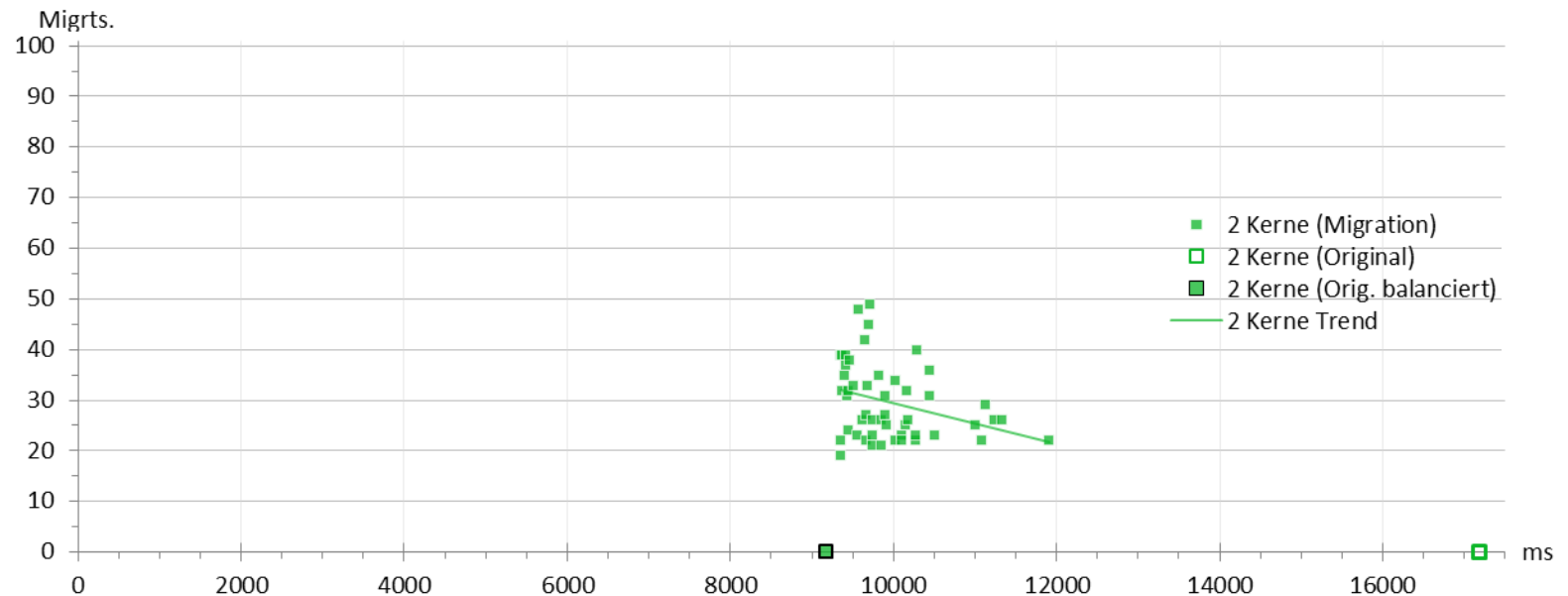


$$S_{p, mittel} = 99,74 \%$$

$$S_{p, max} = 99,98 \%$$

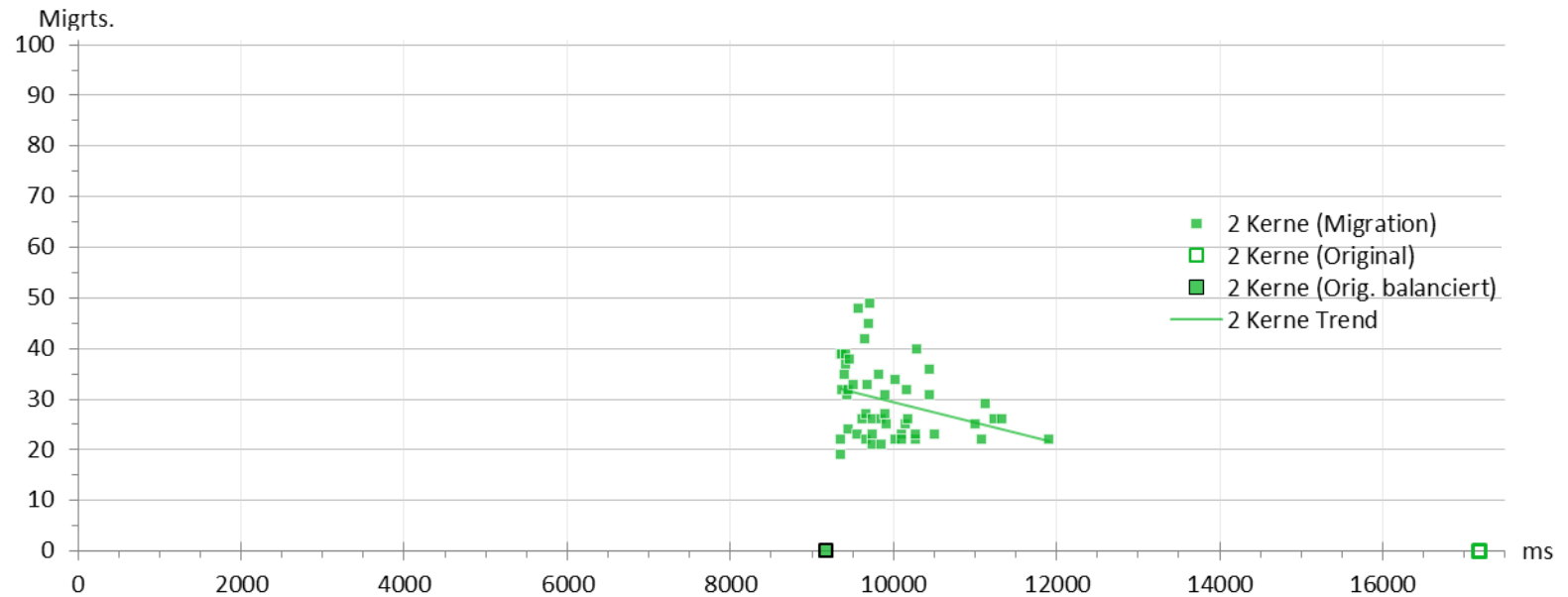
Günstigster Vergleich (unausgeglichenes System)

JGF-Crypt-Benchmark (16 Threads)



Günstigster Vergleich (unausgeglichenes System)

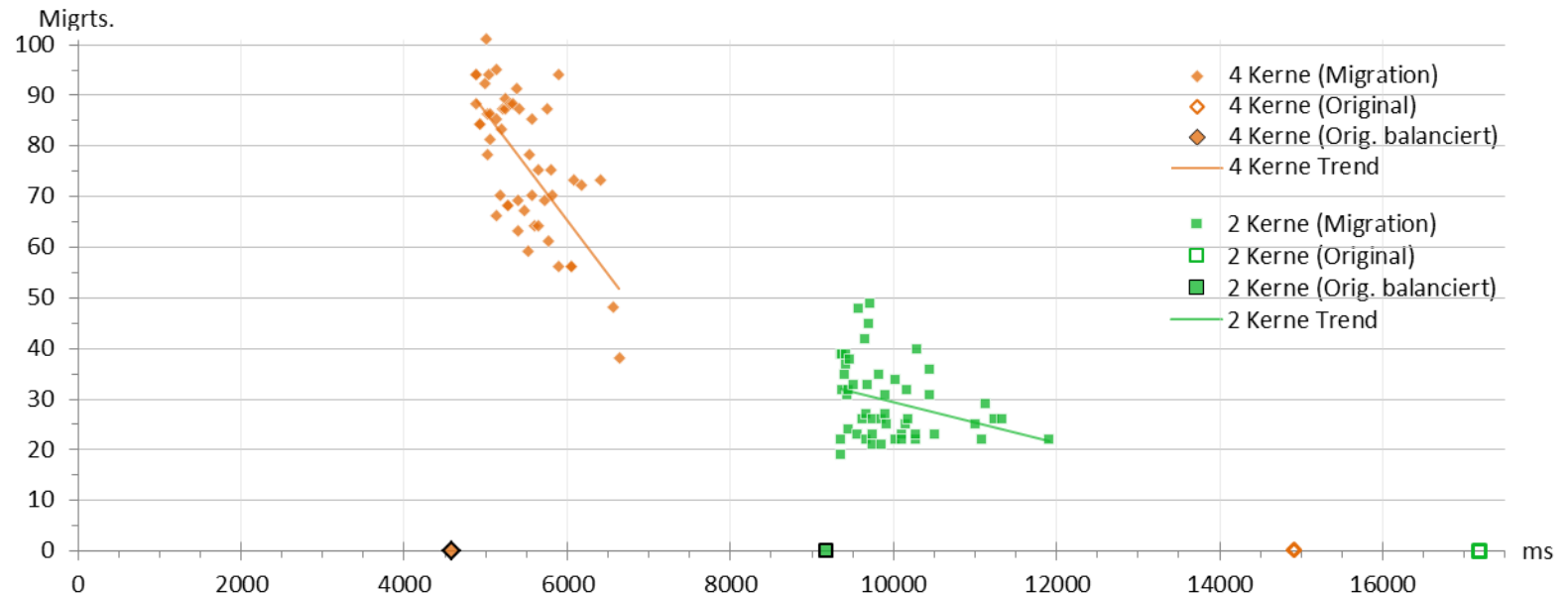
JGF-Crypt-Benchmark (16 Threads)



$S_{p, mittel}$	=	1,725
$S_{p, max}$	=	1,837
$S_{p, theor}$	=	1,875

Günstigster Vergleich (unausgeglichenes System)

JGF-Crypt-Benchmark (16 Threads)



$$S_{p, \text{mittel}} = 2,73$$

$$S_{p, \text{max}} = 3,05$$

$$S_{p, \text{theor}} = 3,25$$

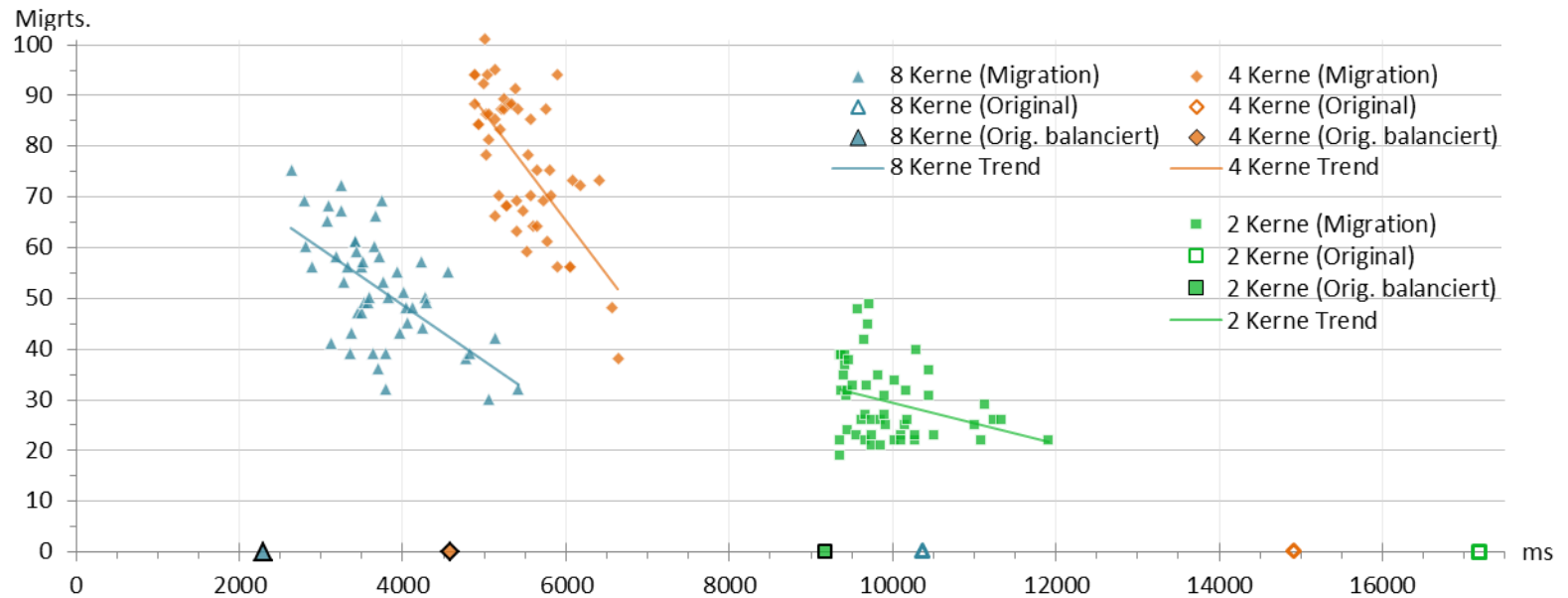
$$S_{p, \text{mittel}} = 1,725$$

$$S_{p, \text{max}} = 1,837$$

$$S_{p, \text{theor}} = 1,875$$

Günstigster Vergleich (unausgeglichenes System)

JGF-Crypt-Benchmark (16 Threads)



$$S_{p, \text{mittel}} = 2,78$$

$$S_{p, \text{max}} = 3,94$$

$$S_{p, \text{theor}} = 4,50$$

$$S_{p, \text{mittel}} = 2,73$$

$$S_{p, \text{max}} = 3,05$$

$$S_{p, \text{theor}} = 3,25$$

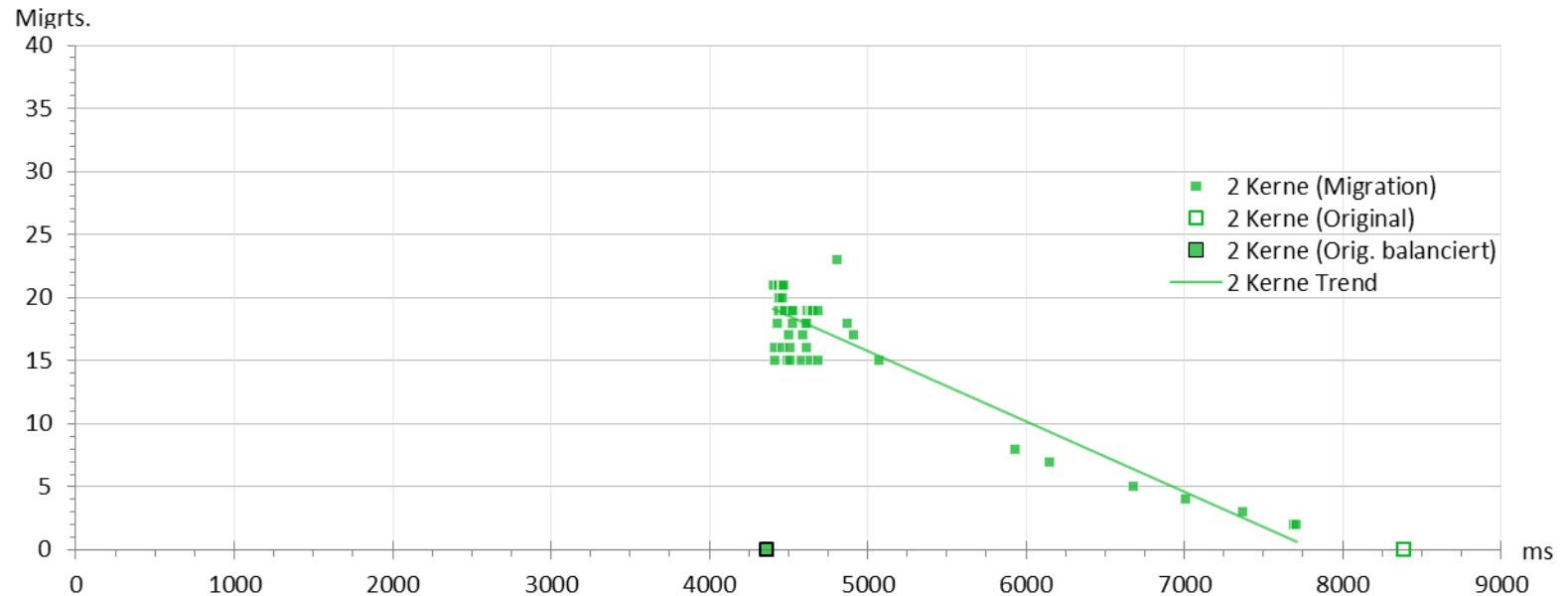
$$S_{p, \text{mittel}} = 1,725$$

$$S_{p, \text{max}} = 1,837$$

$$S_{p, \text{theor}} = 1,875$$

Günstigster Vergleich (unausgeglichenes System)

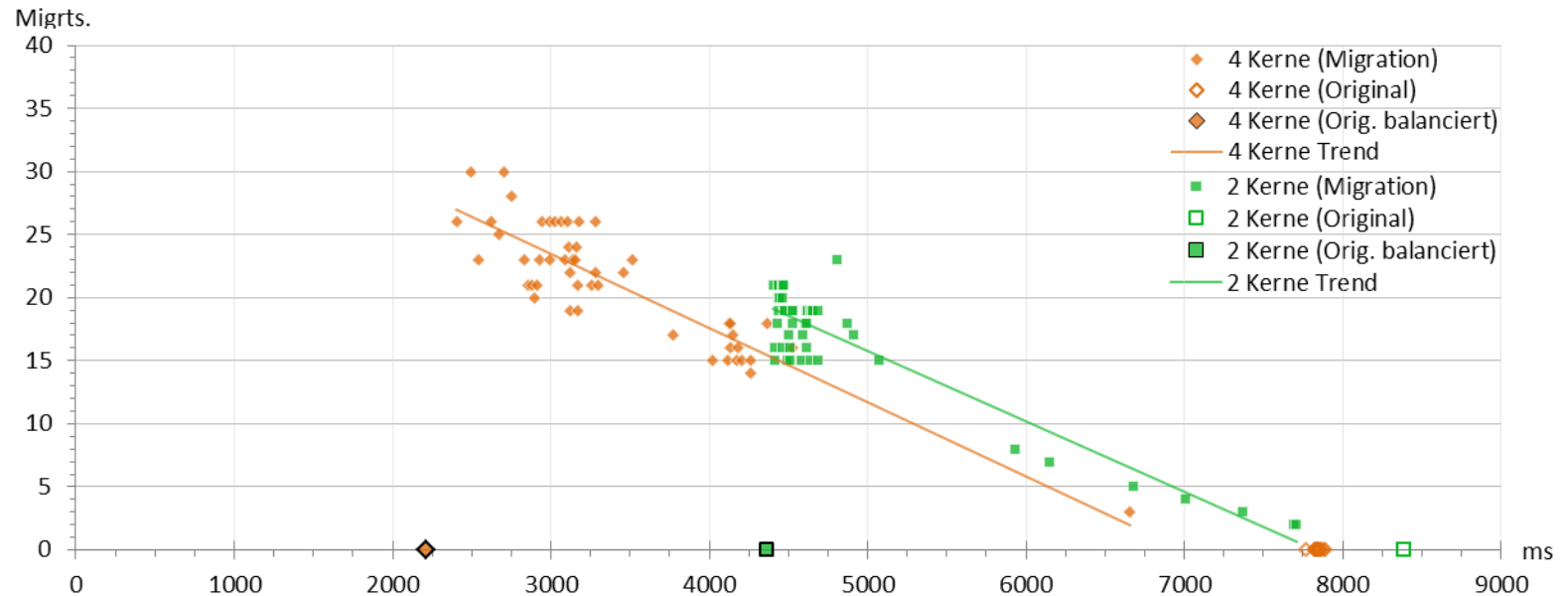
JGF-SparseMatmult-Benchmark (24 Threads)



$S_{p, mittel}$	= 1,715
$S_{p, max}$	= 1,904
$S_{p, theor}$	= 1,917

Günstigster Vergleich (unausgeglichenes System)

JGF-SparseMatmult-Benchmark (24 Threads)



$$S_{p, mittel} = 2,30$$

$$S_{p, max} = 3,23$$

$$S_{p, theor} = 3,50$$

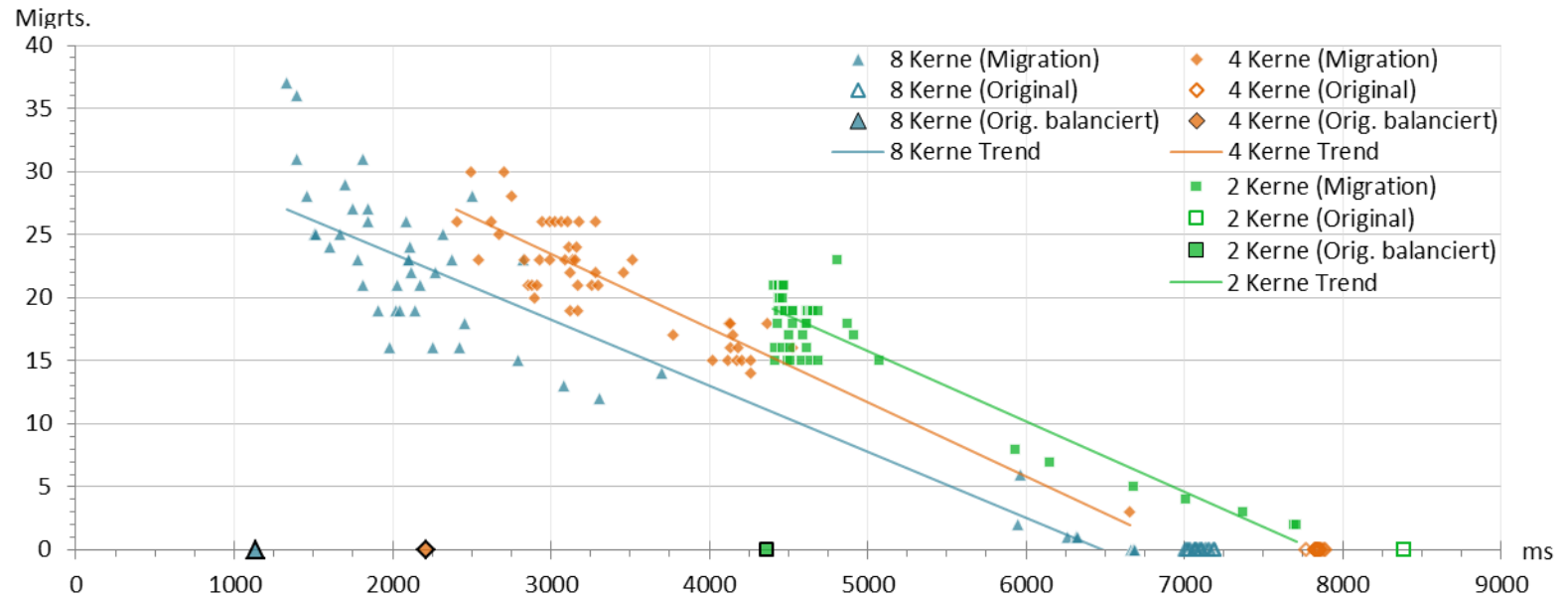
$$S_{p, mittel} = 1,715$$

$$S_{p, max} = 1,904$$

$$S_{p, theor} = 1,917$$

Günstigster Vergleich (unausgeglichenes System)

JGF-SparseMatmult-Benchmark (24 Threads)



$$S_{p, \text{mittel}} = 2,42$$

$$S_{p, \text{max}} = 5,26$$

$$S_{p, \text{theor}} = 5,67$$

$$S_{p, \text{mittel}} = 2,30$$

$$S_{p, \text{max}} = 3,23$$

$$S_{p, \text{theor}} = 3,50$$

$$S_{p, \text{mittel}} = 1,715$$

$$S_{p, \text{max}} = 1,904$$

$$S_{p, \text{theor}} = 1,917$$

Entwurf

- Distributed Weighted Round-Robin als Strategie
- GC-Bus als Verbindungs-Hardware

Entwurf

- Distributed Weighted Round-Robin als Strategie
- GC-Bus als Verbindungs-Hardware

Anpassungen

- Mikrocode erweitert
- 3 Locks vor Migration
- Zielthreadsuche und Nachbearbeitung überwiegend in Java
- Migrator als neues Hardware-Modul
- Erweiterung des GC-Bus
- Erweiterung anderer Module (vor allem Stack)

Funktionalität

- Threadmigration funktioniert
- fehlendes Feature bei Leerlauf eines Kerns

Funktionalität

- Threadmigration funktioniert
- fehlendes Feature bei Leerlauf eines Kerns

Ergebnisse

- 16 % mehr LUTs
- mittlerer Speed-Up: 42 - 92 % des theor. Maximums
- beste Speed-Up: 87 - 99 % des theor. Maximums
- schlechtester Speed-Up: 96 % zum Original

Funktionalität

- Threadmigration funktioniert
- fehlendes Feature bei Leerlauf eines Kerns

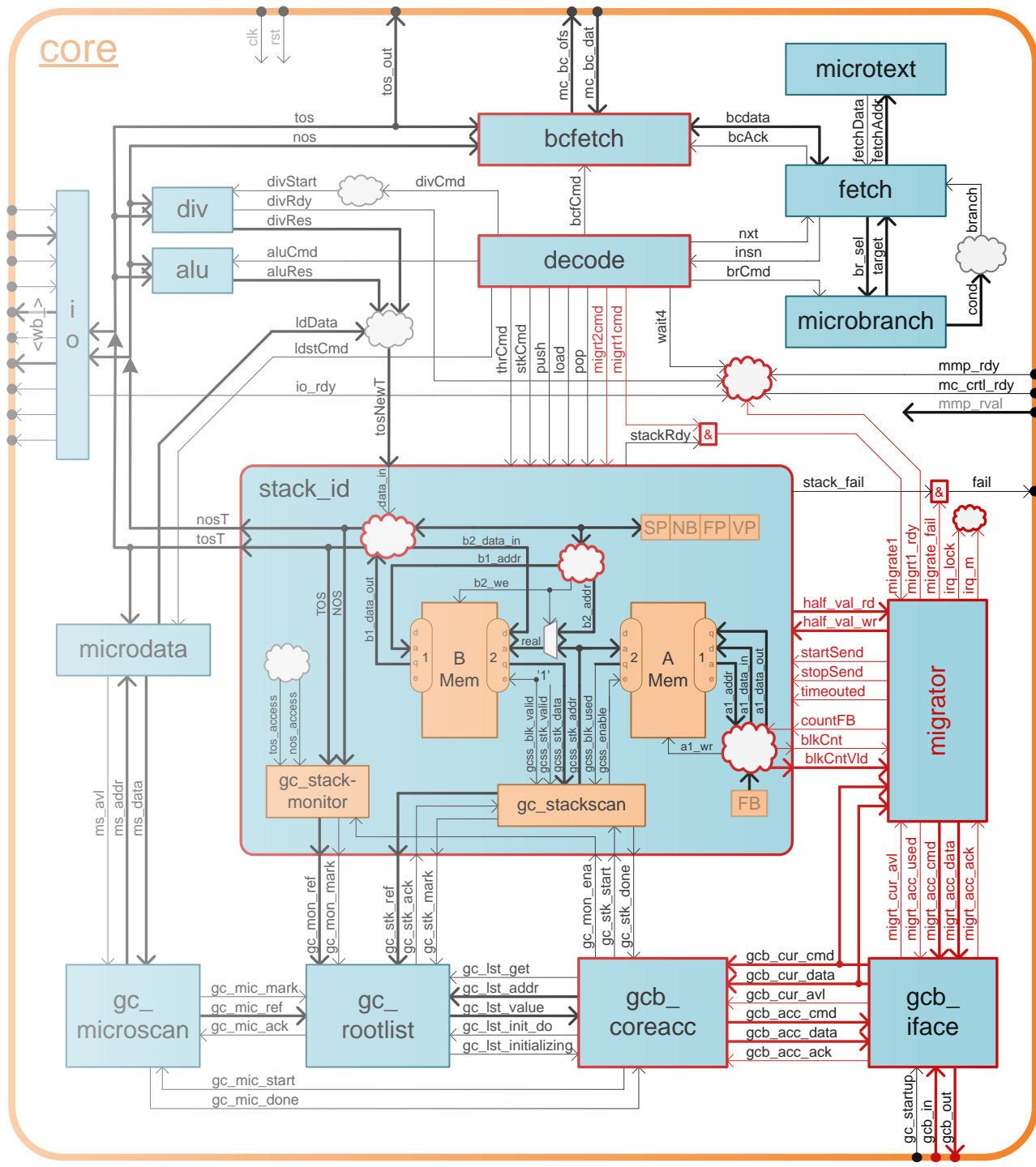
Ergebnisse

- 16 % mehr LUTs
- mittlerer Speed-Up: 42 - 92 % des theor. Maximums
- beste Speed-Up: 87 - 99 % des theor. Maximums
- schlechtester Speed-Up: 96 % zum Original

Aussichten

- leerlaufende Kerne nicht isolieren
- Vergrößerung des Mikrocode-Speichers
- Kompatibilität mit optimiertem ShapLinker

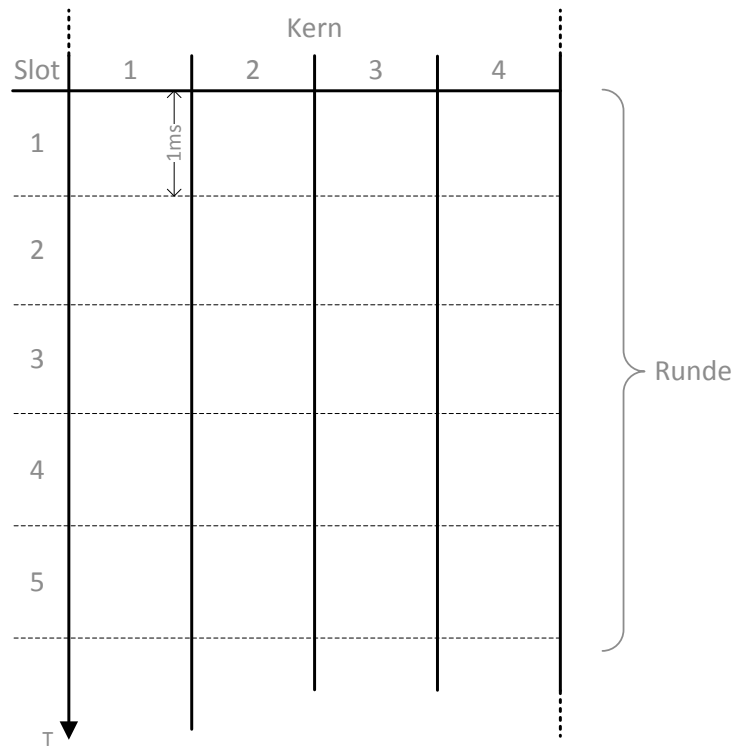
Vielen
Dank!



Gang Scheduling (Inter-Thread-Kommunikation)

Programme (Threads):

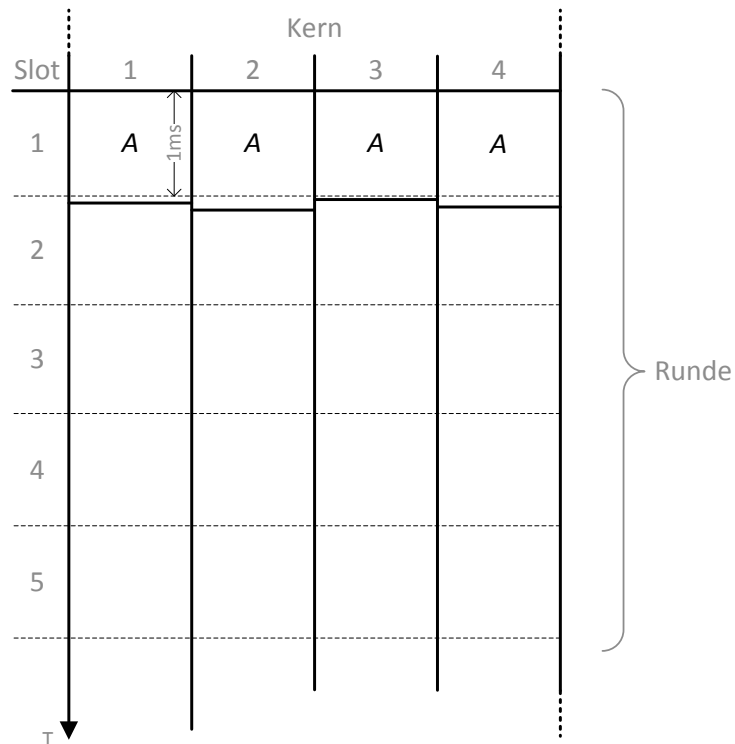
$A(6), B(5), C(3), D(3), E(2), F(1)$



Gang Scheduling (Inter-Thread-Kommunikation)

Programme (Threads):

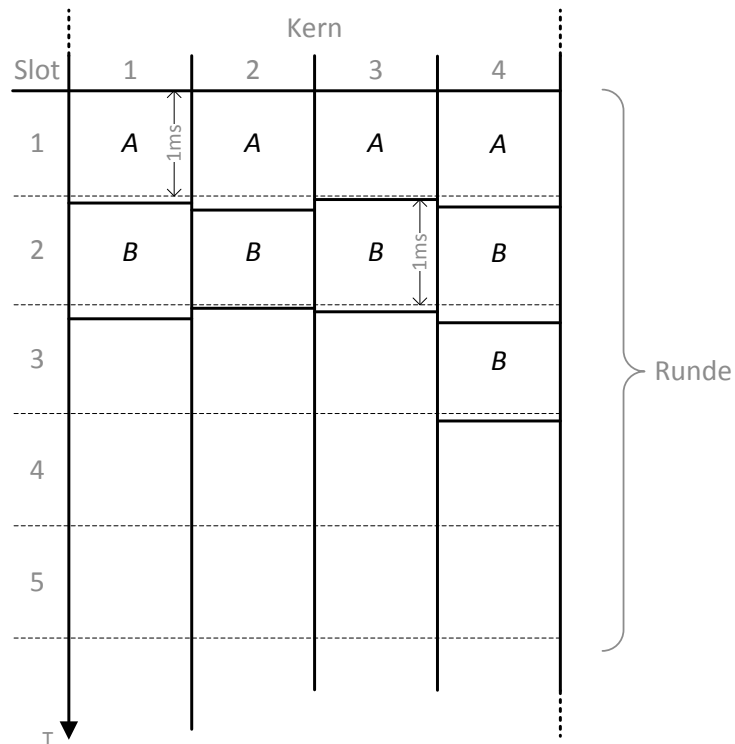
A(6), B(5), C(3), D(3), E(2), F(1)



Gang Scheduling (Inter-Thread-Kommunikation)

Programme (Threads):

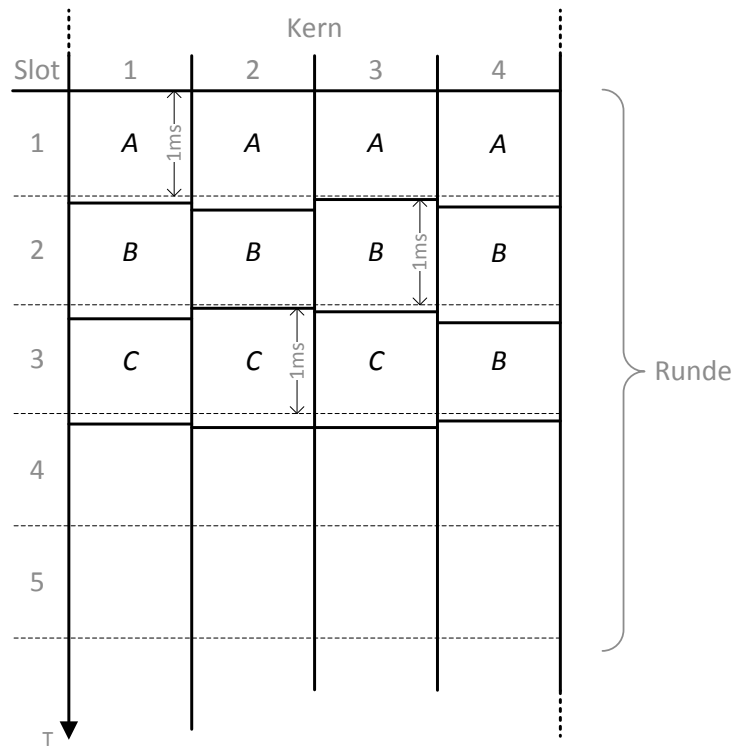
$A(6), B(5), C(3), D(3), E(2), F(1)$



Gang Scheduling (Inter-Thread-Kommunikation)

Programme (Threads):

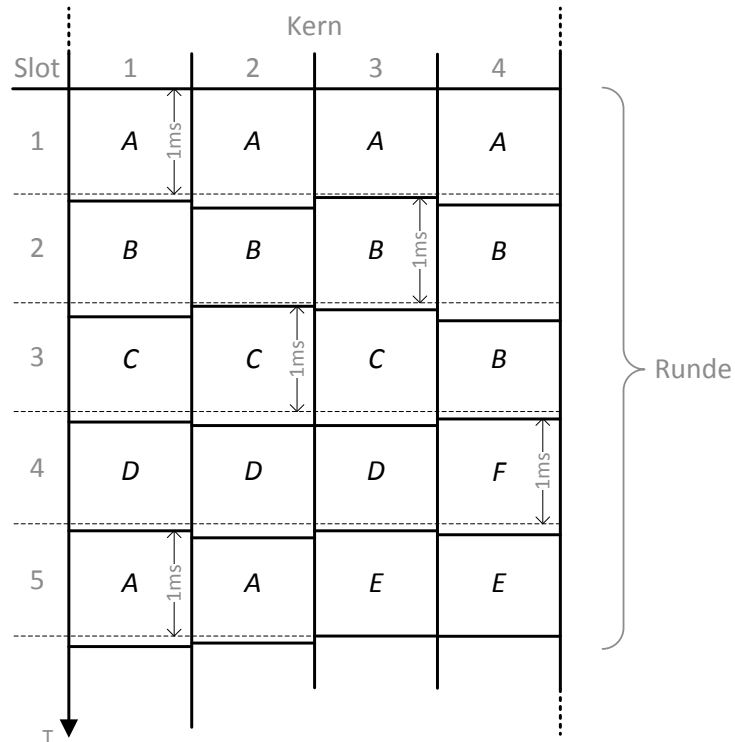
$A(6), B(5), C(3), D(3), E(2), F(1)$



Gang Scheduling (Inter-Thread-Kommunikation)

Programme (Threads):

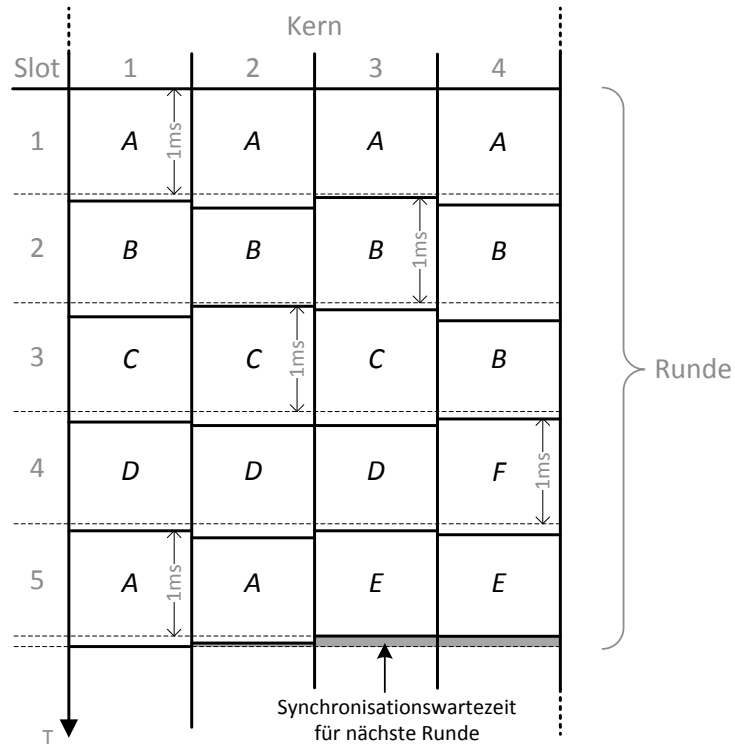
$A(6), B(5), C(3), D(3), E(2), F(1)$



Gang Scheduling (Inter-Thread-Kommunikation)

Programme (Threads):

A(6), B(5), C(3), D(3), E(2), F(1)

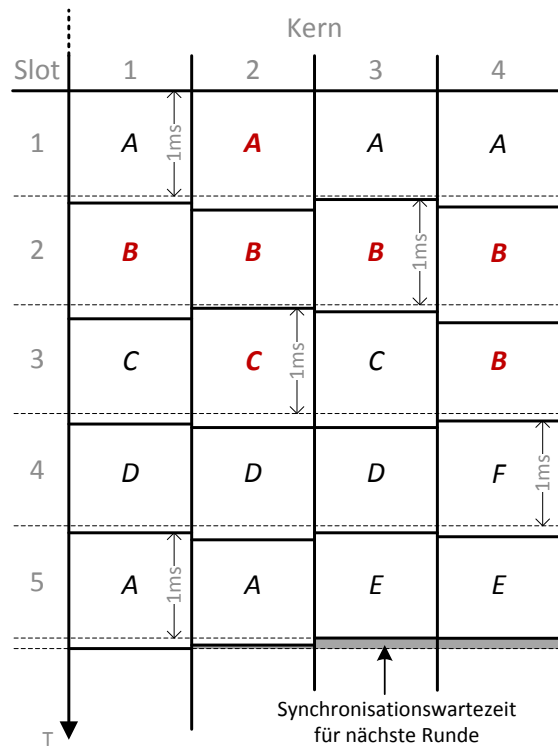


Gang Scheduling (Inter-Thread-Kommunikation)

Programme (Threads):

a) A(6), B(5), C(3), D(3), E(2), F(1)

b) **A(5), B(0), C(2), D(3), E(2), F(1)**

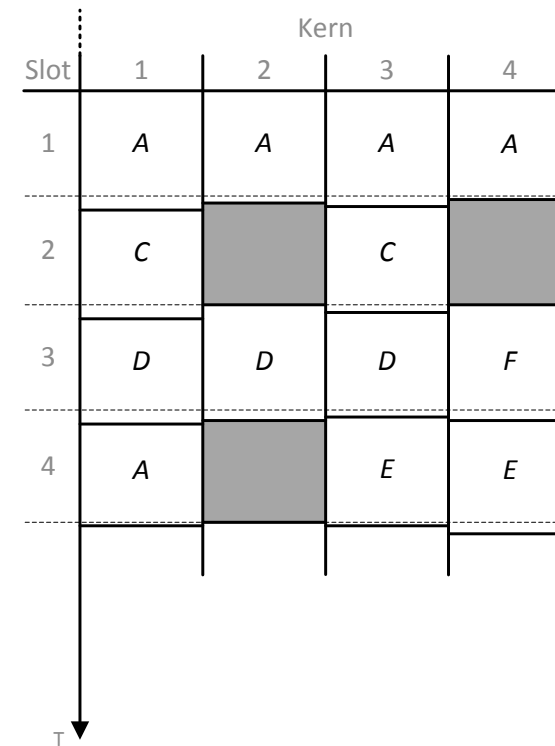
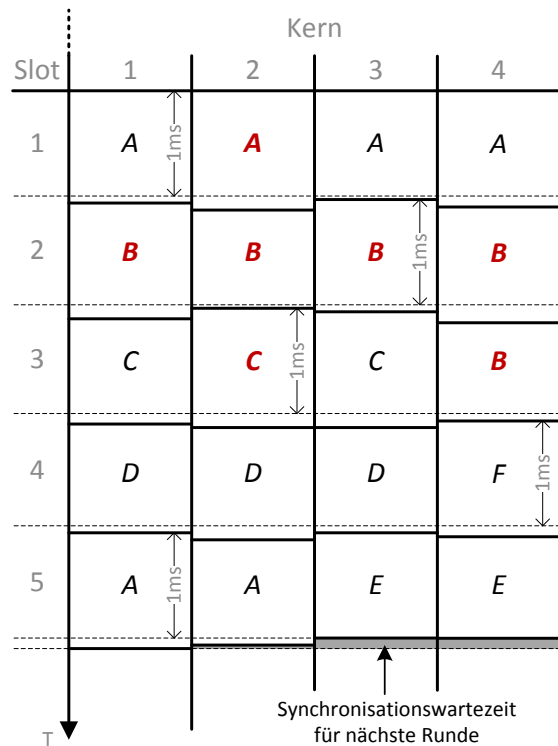


Gang Scheduling (Inter-Thread-Kommunikation)

Programme (Threads):

a) A(6), B(5), C(3), D(3), E(2), F(1)

b) A(5), B(0), C(2), D(3), E(2), F(1)

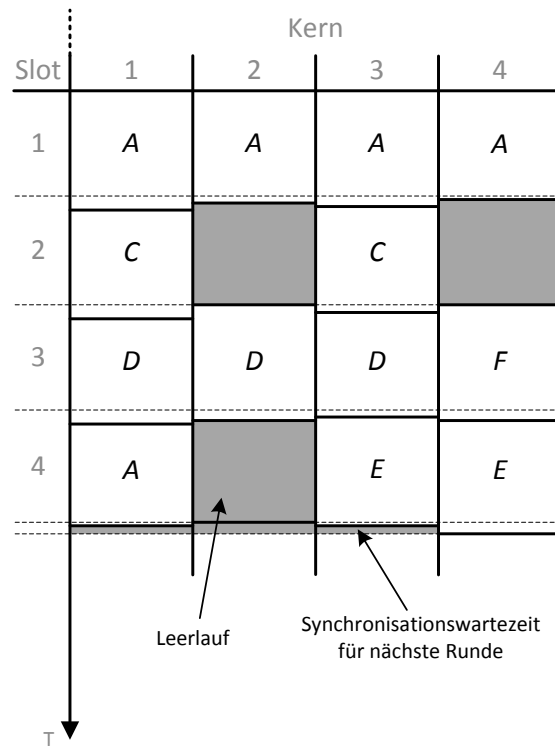
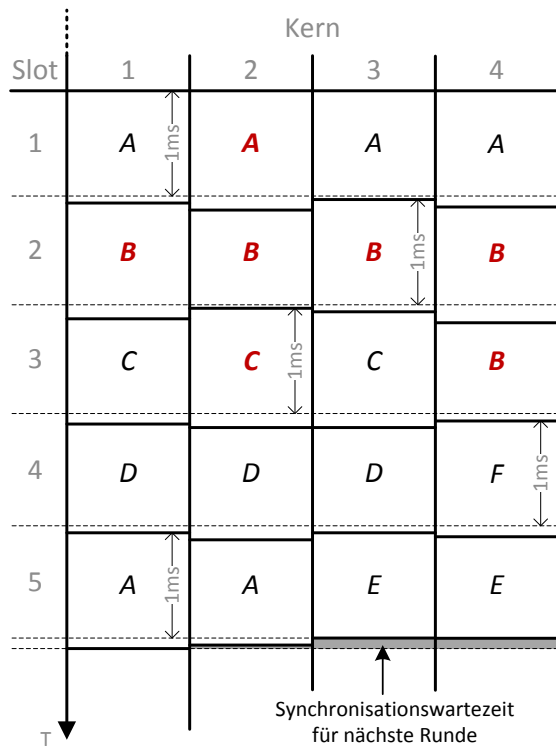


Gang Scheduling (Inter-Thread-Kommunikation)

Programme (Threads):

a) A(6), B(5), C(3), D(3), E(2), F(1)

b) A(5), B(0), C(2), D(3), E(2), F(1)

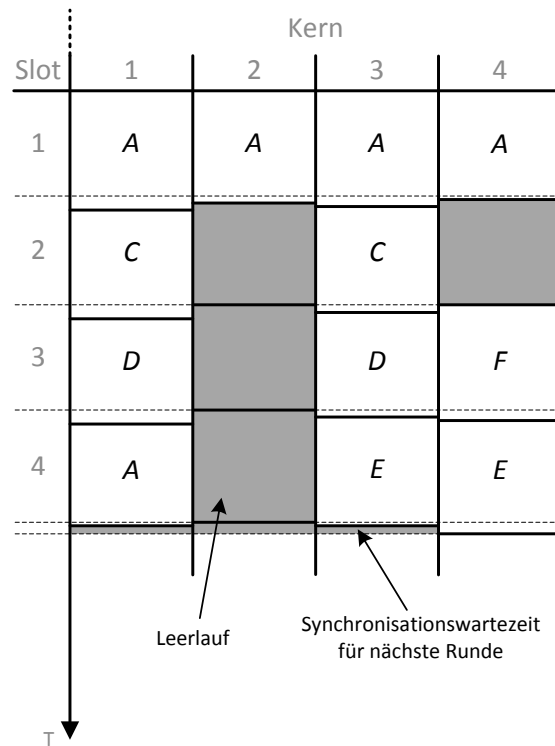
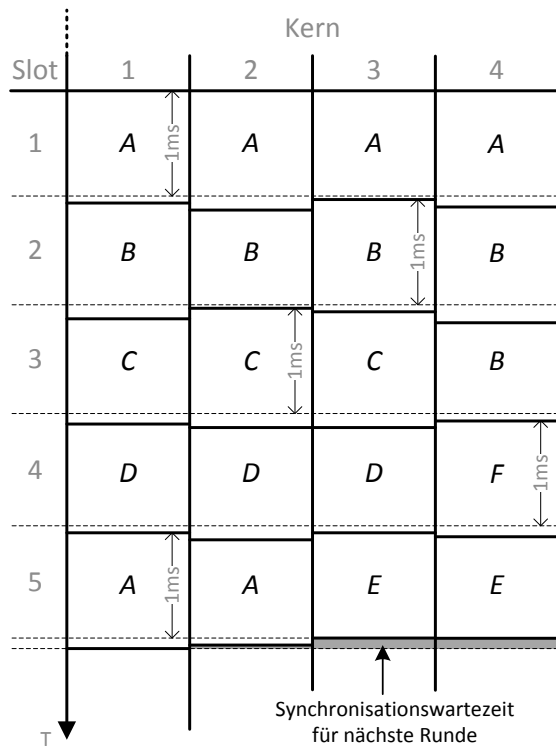


Gang Scheduling (Inter-Thread-Kommunikation)

Programme (Threads):

a) A(6), B(5), C(3), D(3), E(2), F(1)

b) A(5), C(2), D(2), E(2), F(1)

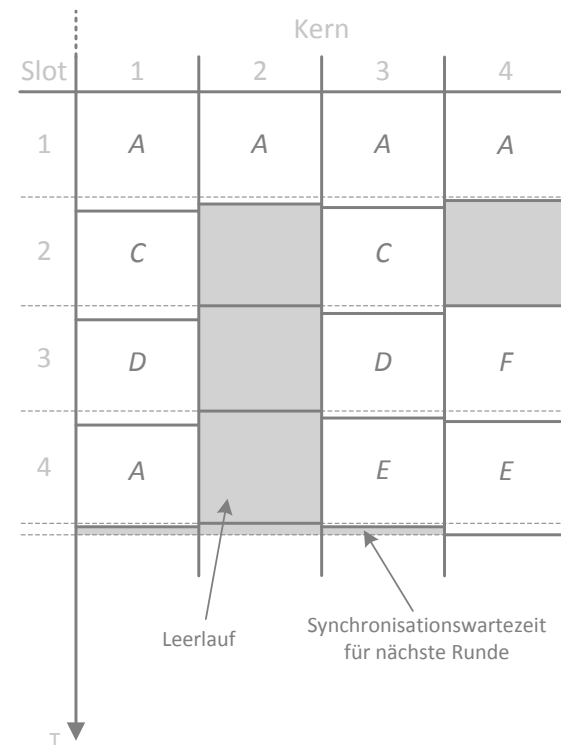


Gang Scheduling (Inter-Thread-Kommunikation)

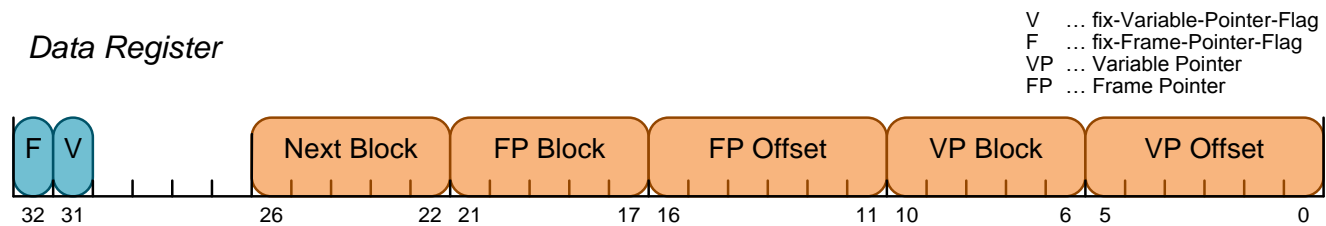
- + Inter-Thread-Kommunikation
- + Fairness
- + Konzept der Threadgruppen passt zu SHAP
- + viele Pack-Algorithmen

- ineffizient ohne synchr. Kerne
- Vorausplanung der Runde (Wartezeit)
- Lücken im Plan
- yield() wenig sinnvoll
- keine Prioritäten
- großer Migrations-Overhead (Wartezeit)

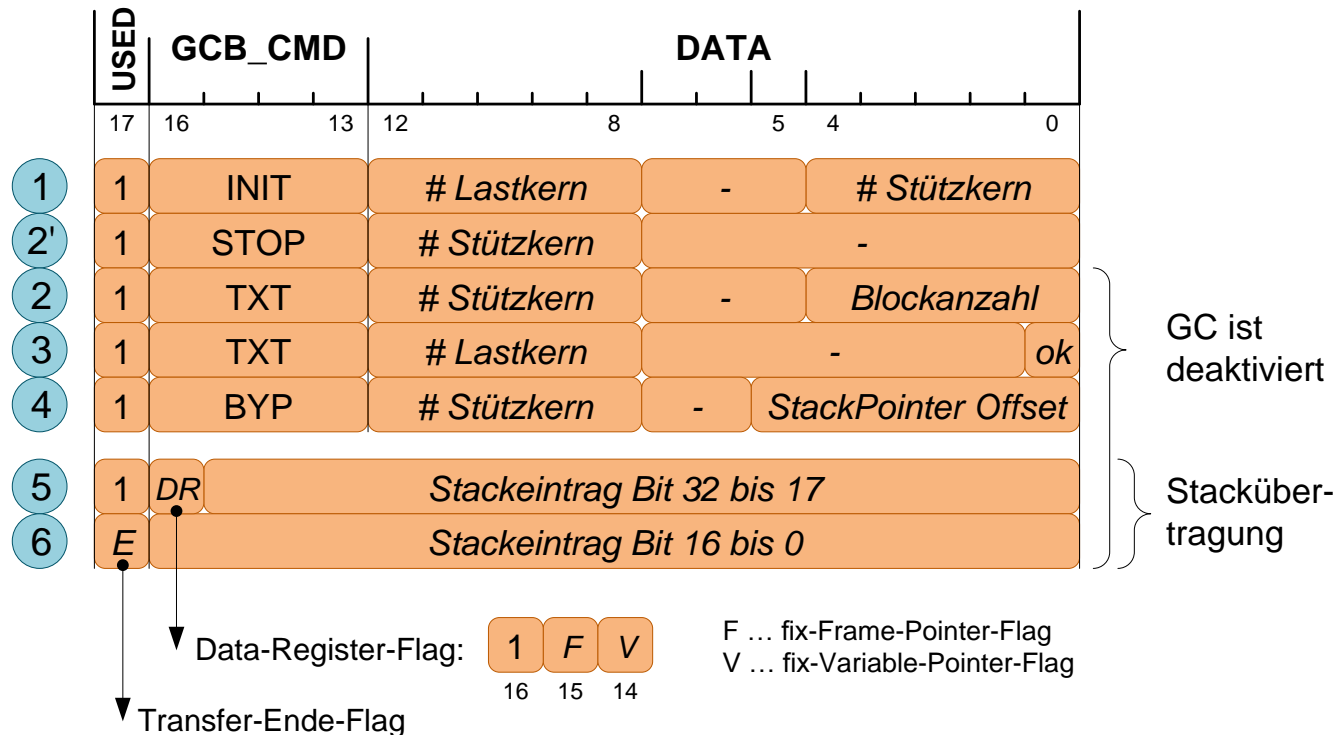
b) A(5), C(2), D(2), E(2), F(1)



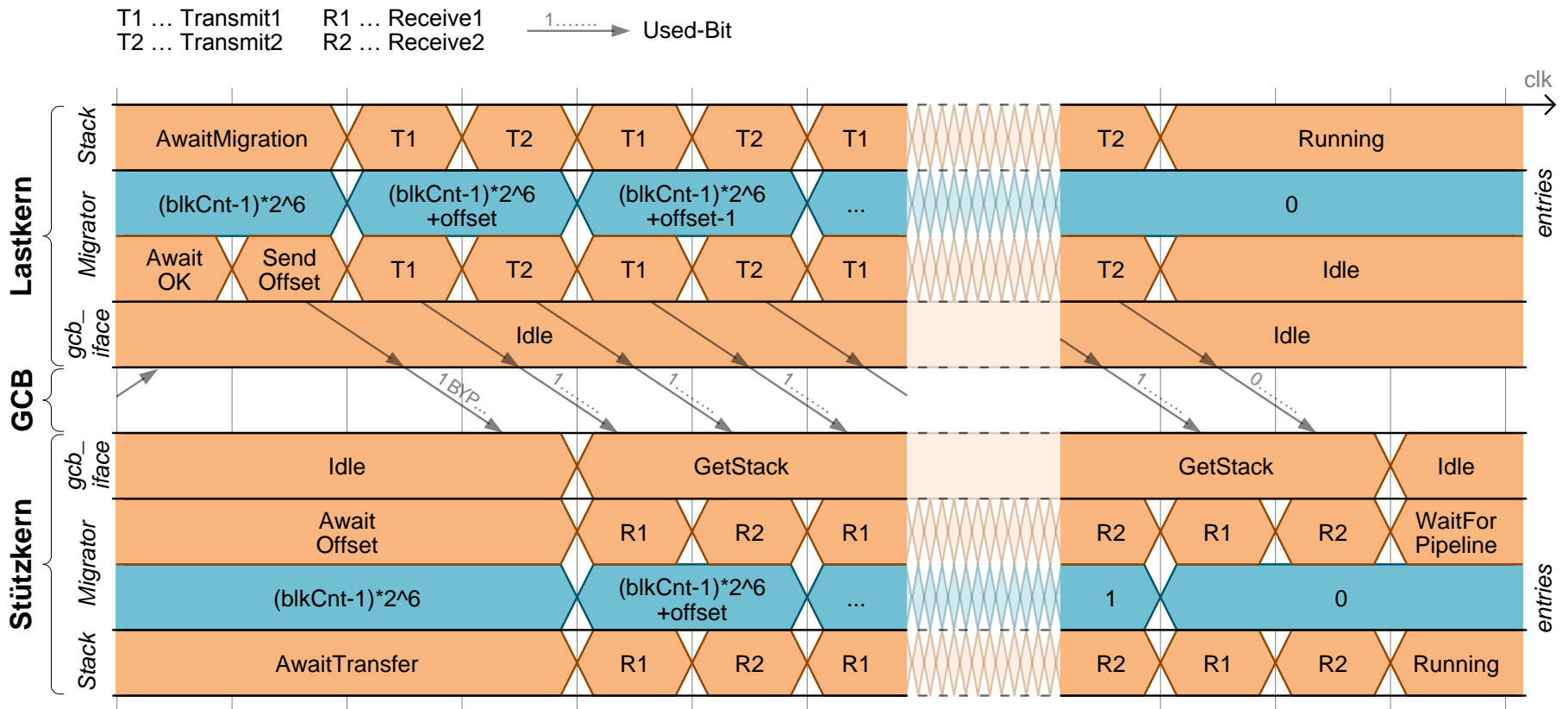
Datenregister Patching



GC-Bus Protokoll



Stackübertragung Zustände



Intuitiver Ansatz

a) *Informationsbasis: Threadanzahl*

Thread	Kern 1		2		3		4	
	μs	%						
1	1000	25	1000	40	1000	49	300	75
2	1000	25	900	36	1000	49	50	12
3	1000	25	400	16	50	2	50	12
4	1000	25	200	8				
Σ	4000		2500		2050		400	

b) *Informationsbasis: Rundenlänge*

Thread	Kern 1		2		3		4	
	μs	%						
1	1000	43	1000	45	1000	46	1000	43
2	1000	43	1000	45	1000	46	900	39
3	300	13	200	9	50	2	400	17
4					50	2		
5					50	2		
Σ	2300		2200		2150		2300	

Cassavant & Kuhl

