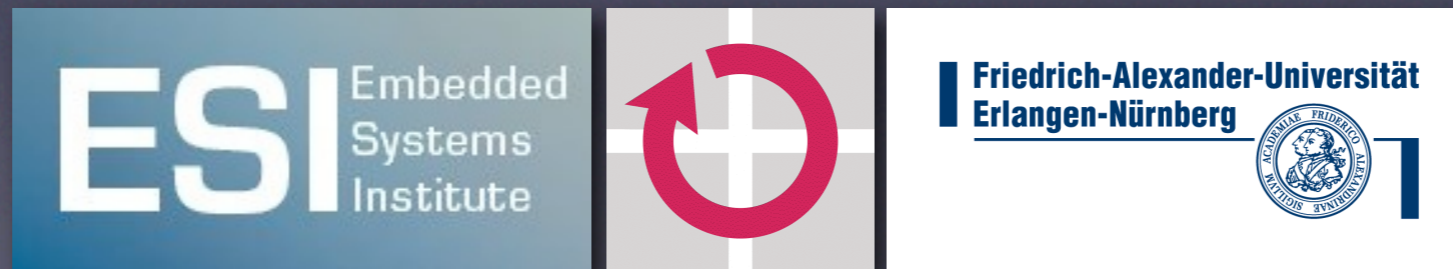
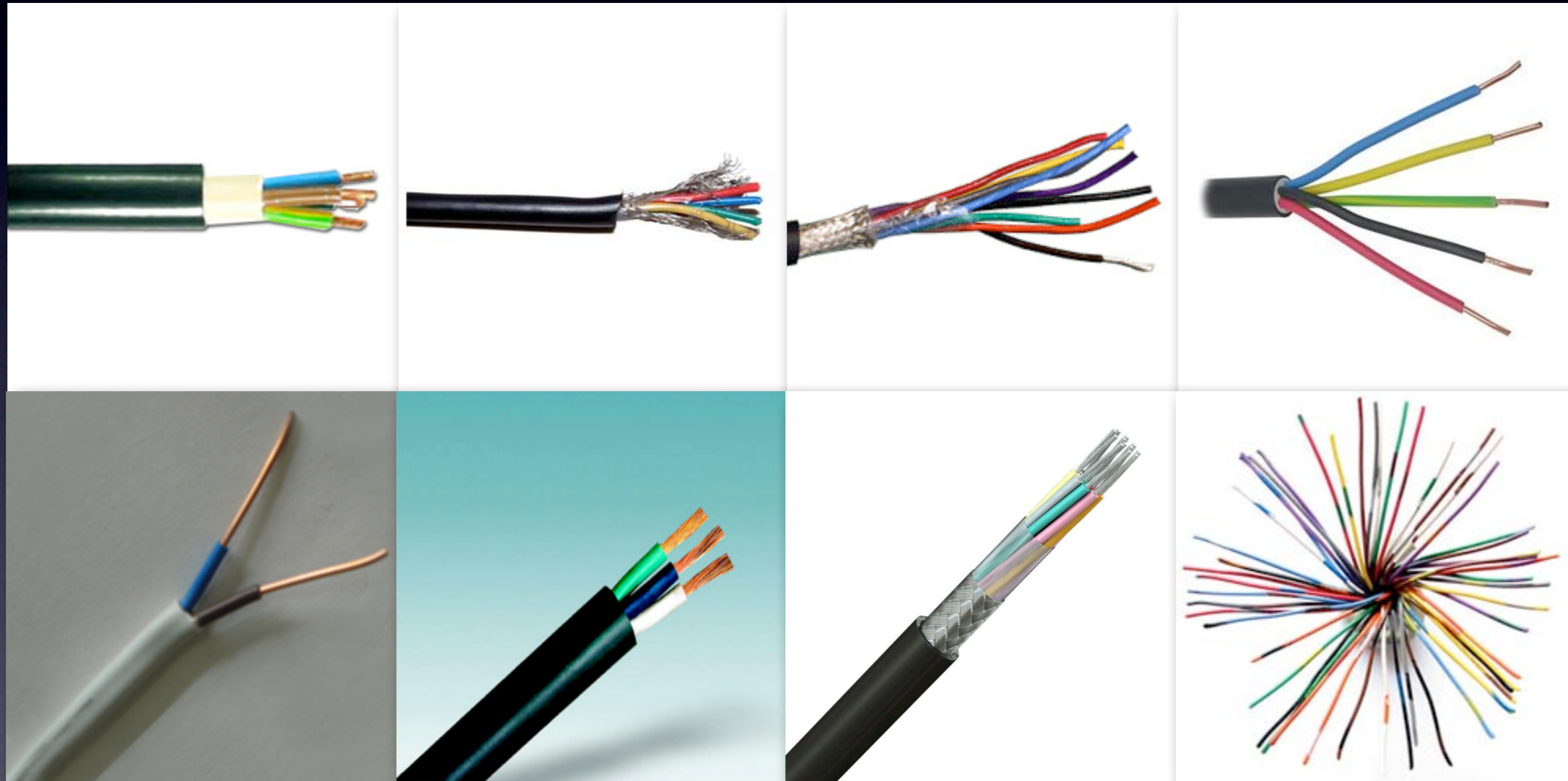


Systemsoftware im Zeitalter mehrkerniger Prozessoren

Wolfgang Schröder-Preikschat

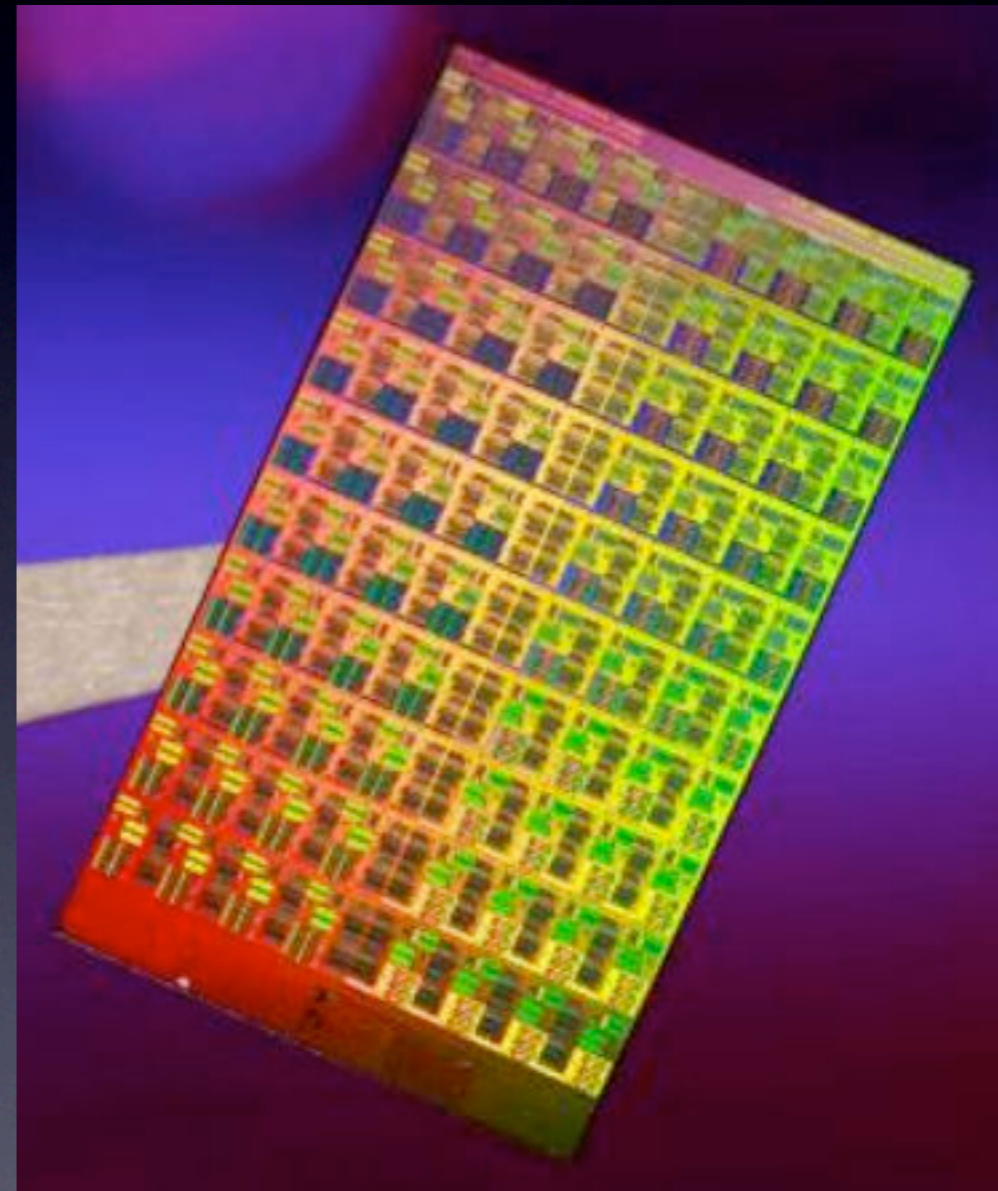
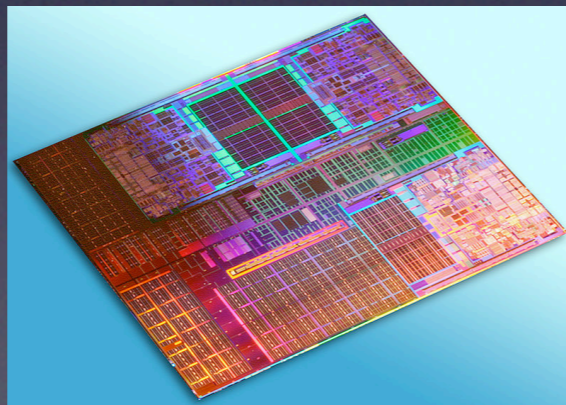
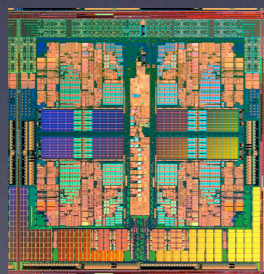
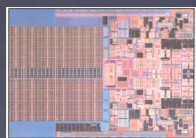


Multi-Core

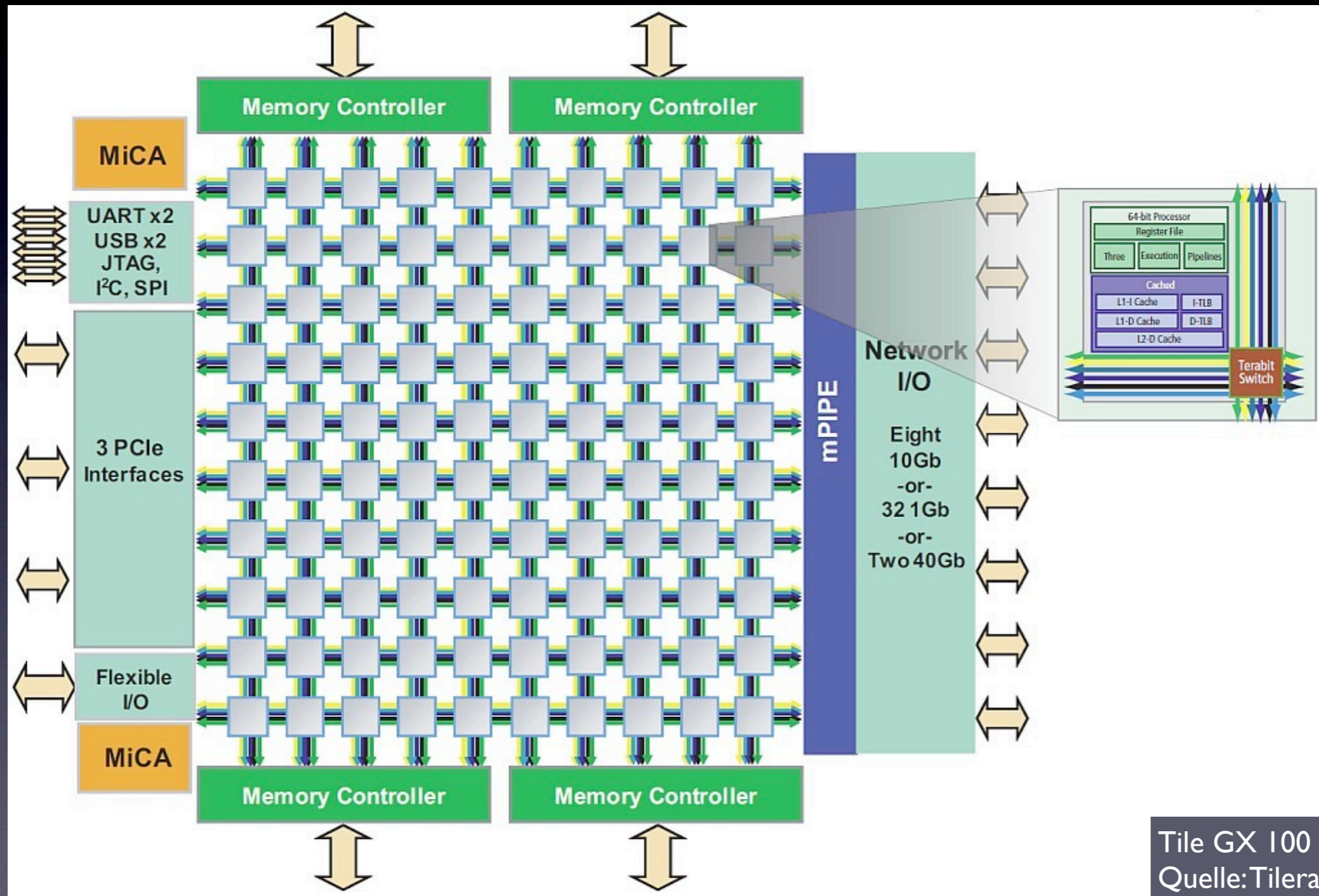


Multi-Core

Mehr- bzw. vielkernige Prozessoren, homogener oder heterogener Art.



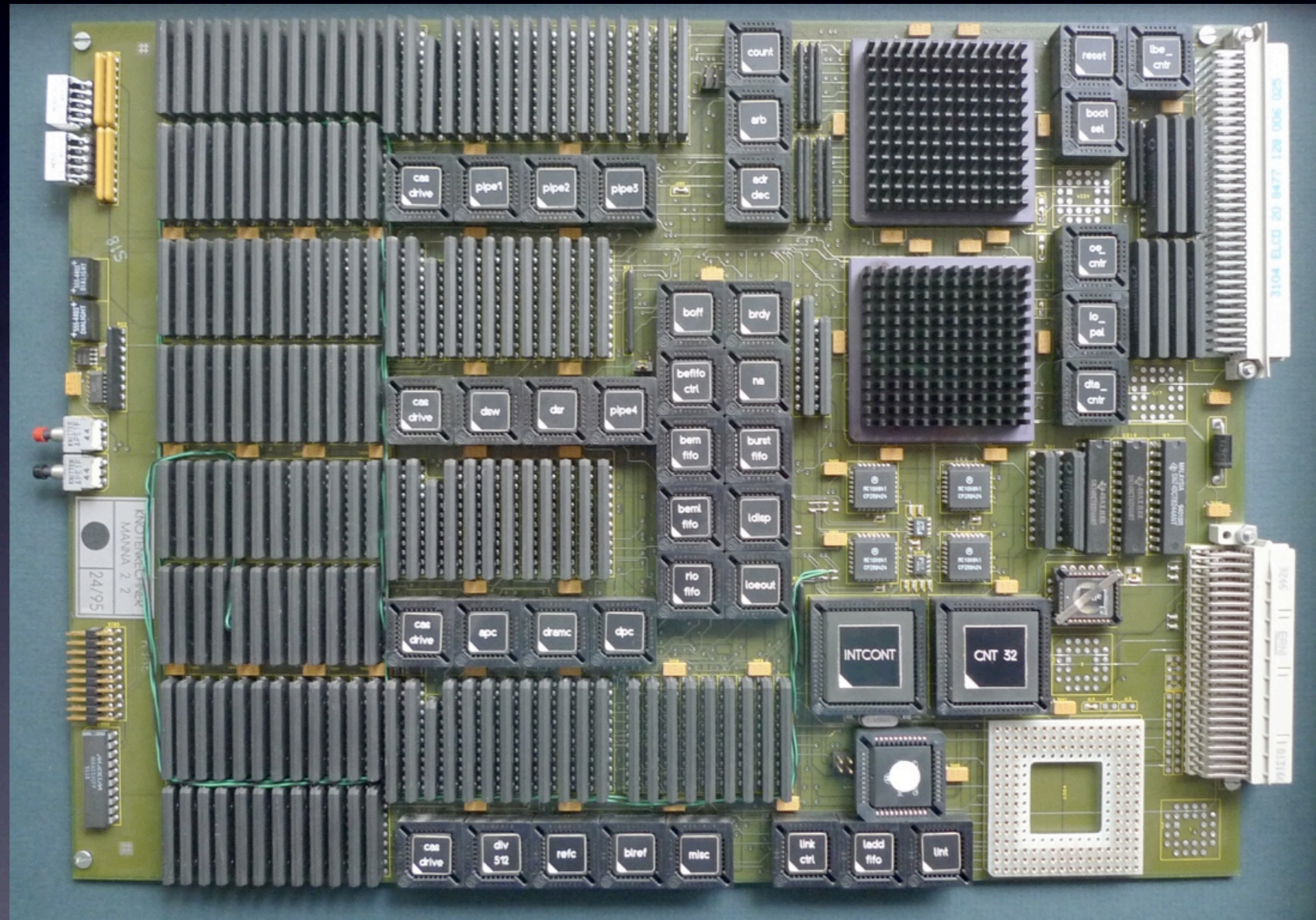
Multi-Core Embedded



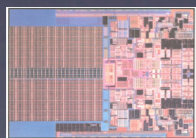
Tile GX 100
Quelle: Tiler

Déjà vu

Dualprocessor



Dualcore



Problemfelder

- Wettstreitigkeit

- Interferenz gleichzeitiger Prozesse

- Synchronisation

- Kooperation & Konkurrenz koordinieren

- Latenzen

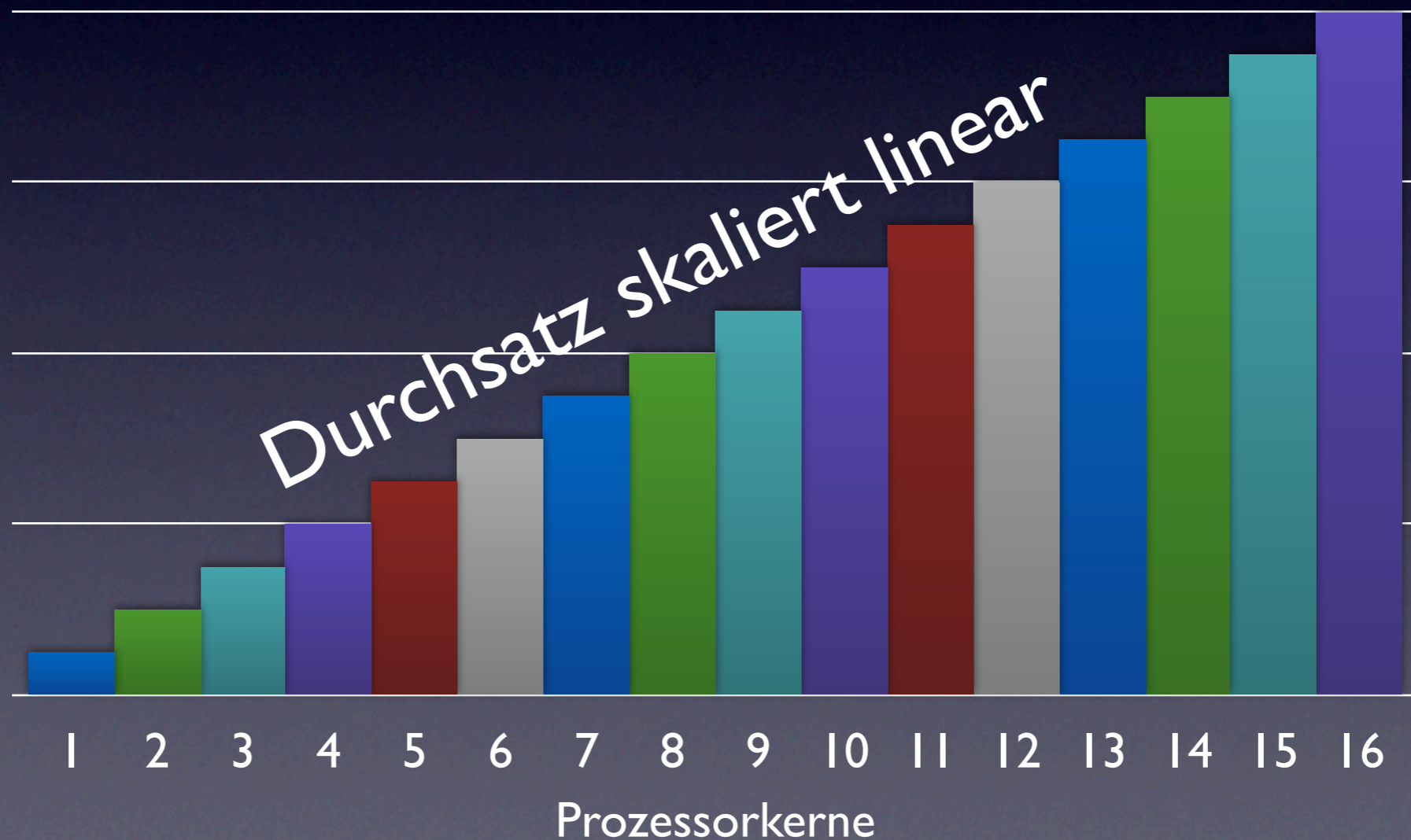
- Kaschierung von Systemaktivitäten

Wettstreitigkeit

- Fallstudie: Dateideskriptortabelle
 - gemeinsam benutzte Datenstruktur
- Testlauf: Kosten bei gleichzeitigen Zugriffen
 - `dup/close` fadenlokaler Deskriptoren
 - pro Kern genau ein `dup/close`-Faden
 - 16-Kern AMD Opteron, Linux 2.6.27

Theorie

Dateideskriptortabelle: # dup/close pro Sekunde

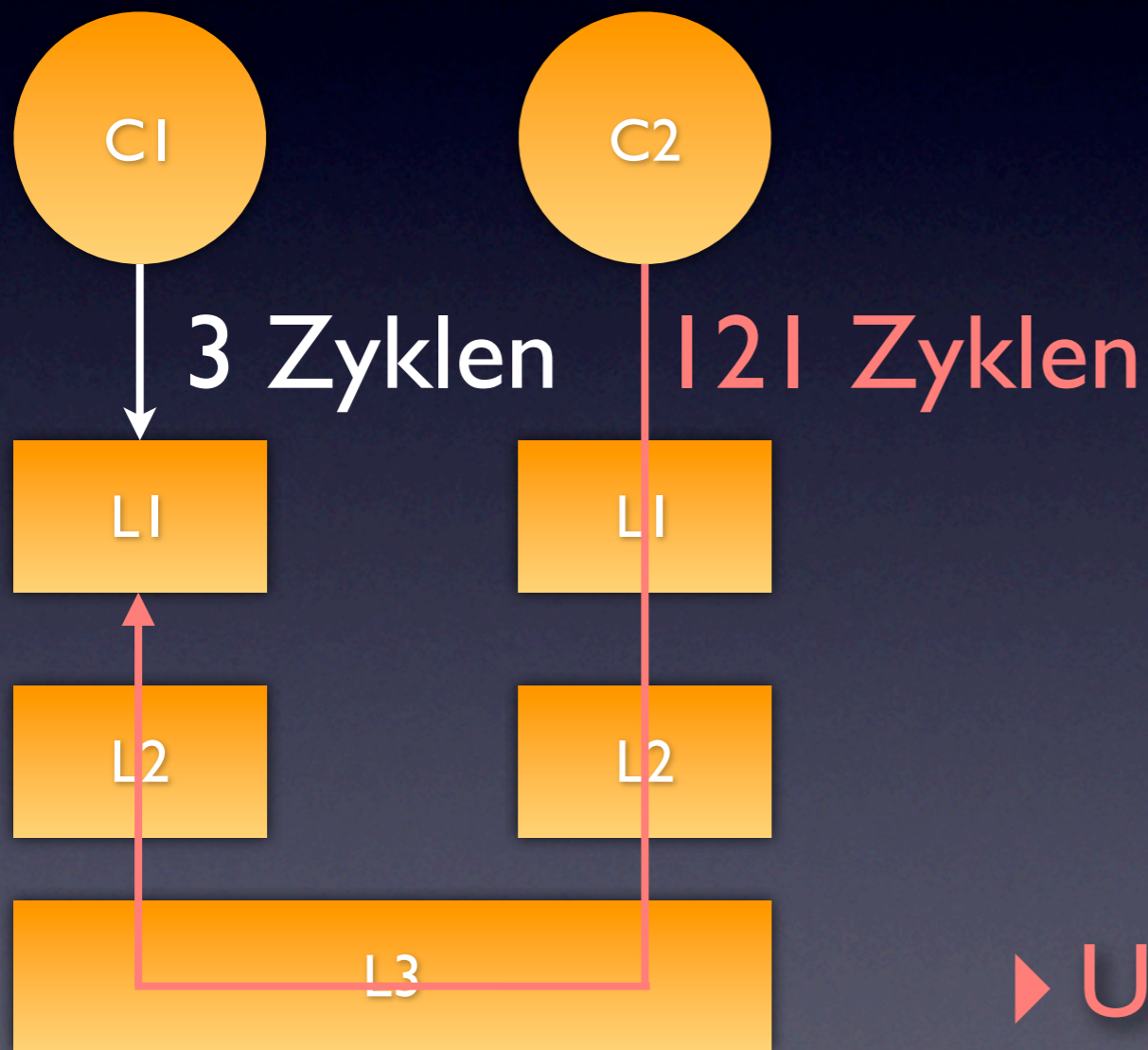


Praxis

Dateideskriptortabelle: # dup/close pro Sekunde



Zugriffsmuster



```
fd_alloc () {  
    lock(fd_table);  
    fd = get_free_fd();  
    set_fd_used(fd);  
    fix_smallest_fd();  
    unlock(fd_table);  
}
```

► Ursache für Durchsatzabfall

Konkurrenz

Die **Vorgangssperre** sequentialisiert alle Zugriffe auf `fd_table`, obwohl nicht alle Fäden auf alle Tabelleneinträge zugleich zugreifen!

```
fd_alloc () {  
    lock(fd_table);  
    fd = get_free_fd();  
    set_fd_used(fd);  
    fix_smallest_fd();  
    unlock(fd_table);  
}
```

false sharing

► Ursache keiner Verbesserung

Wo der Schuh drückt...

- Implementierung des Betriebssystems
 - mangelnde Skalierbarkeit
- Implementierung der Anwendung
 - mangelnde Parallelität
- Nutzung von Betriebssystemfunktionen
 - mangelndes Bewusstsein

Skalierungsanleitungen für den Linux-Kern

- mehrkernige Paketverarbeitung
- Elimination falscher Mitbenutzung
- „schludrige“ (*sloppy*) Zähler
- prozessorkernlokale Datenstrukturen
- sperrfreie Vergleiche
- Vermeidung unnötiger Blockierungen

Empfehlungen

- Mitbenutzung kontrollierbar machen
 - Anwendungen entscheiden lassen können
- geschickte Einplanung gleichzeitiger Fäden
 - Anwendungswissen einfließen lassen
- Softwaretechnik (wirklich) praktizieren
 - Entwurfsentscheidungen zurückstellen

Synchronisation



Synchronisation

- altbekannte Techniken dominieren
 - Schlossvariable, Semaphor, Monitor
- naiver Einsatz blockierender Verfahren
 - kritischer Abschnitt \neq unteilbar
 - CPU \neq unteilbares Betriebsmittel
- funktionale vs. nichtfunktionale Eigenschaft

Umbau von Altsoftware
ist nicht trivial...

Untiefen...

FreeBSD

cpu_spinwait turnstile_try

lock_mtx
mtx_lock
mtx_lock_flags
_get_sleep_lock
mtx_lock_sleep
turnstile_wait
mtx_lock_spin
lock_spin_flags
lock_spin_flags
_get_spin_lock
_mtx_lock_spin
_obtain_lock
atomic_cmpset_acp_ptr
atomic_cmpset_acp_long
atomic_cmpset_long
cmpxchgq

thread_lock
thread_lock_flags
_thread_lock_flags
spinlock_enter
intr_disable

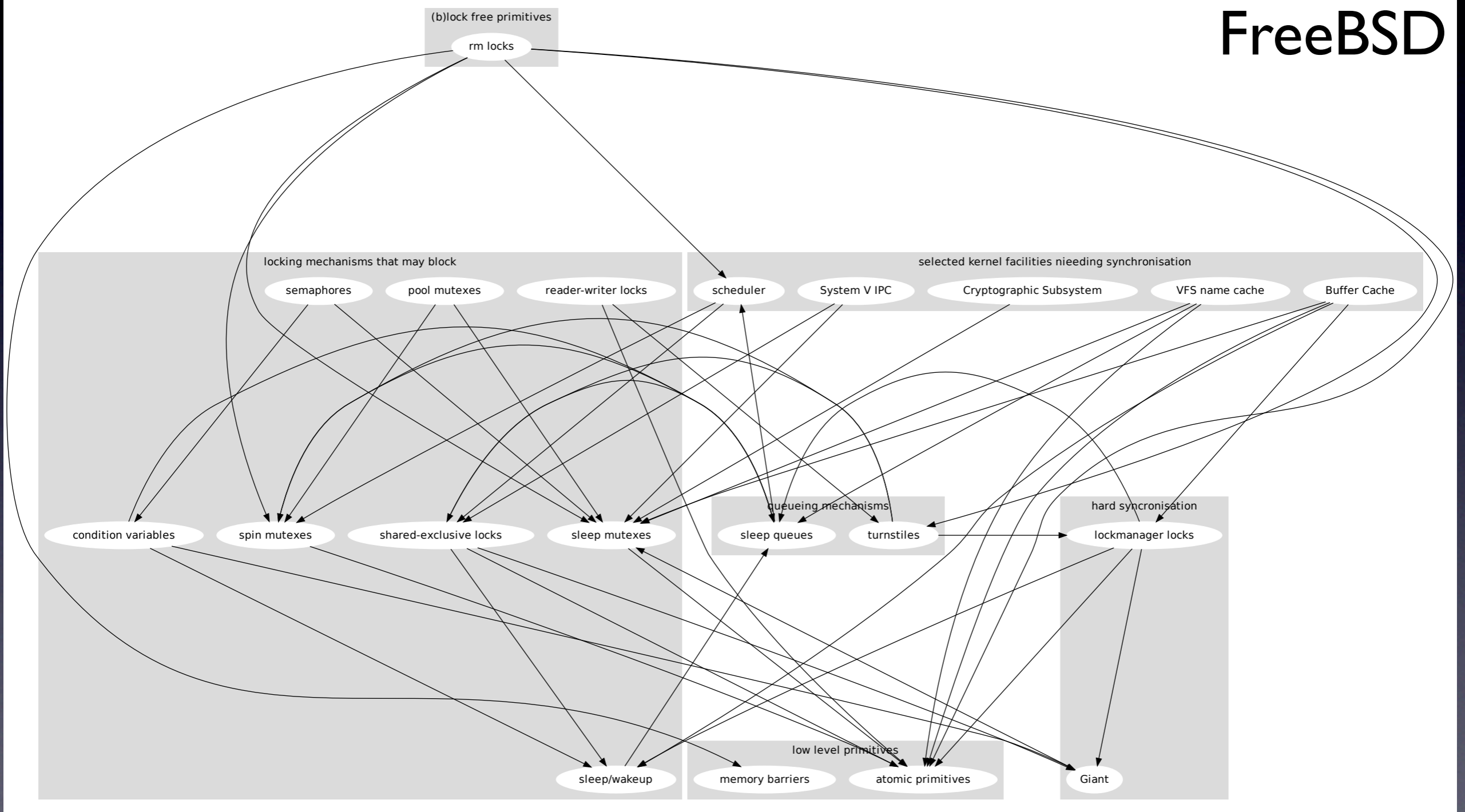


```
void enable_mmioTRACE (void) {  
→ mutex_lock(&mmioTRACE_mutex);  
  if (is_enabled())  
    goto out;  
  if (nommioTRACE)  
    pr_info("MMIO trace disabled\n");  
  kmmio_init();  
→ enter_uniprocessor();  
→ spin_lock_irq(&trace_lock);  
→ atomic_inc(&mmioTRACE_enabled);  
→ spin_unlock_irq(&trace_lock);  
  pr_info("enabled.\n");  
  out:  
→ mutex_unlock(&mmioTRACE_mutex);  
}
```

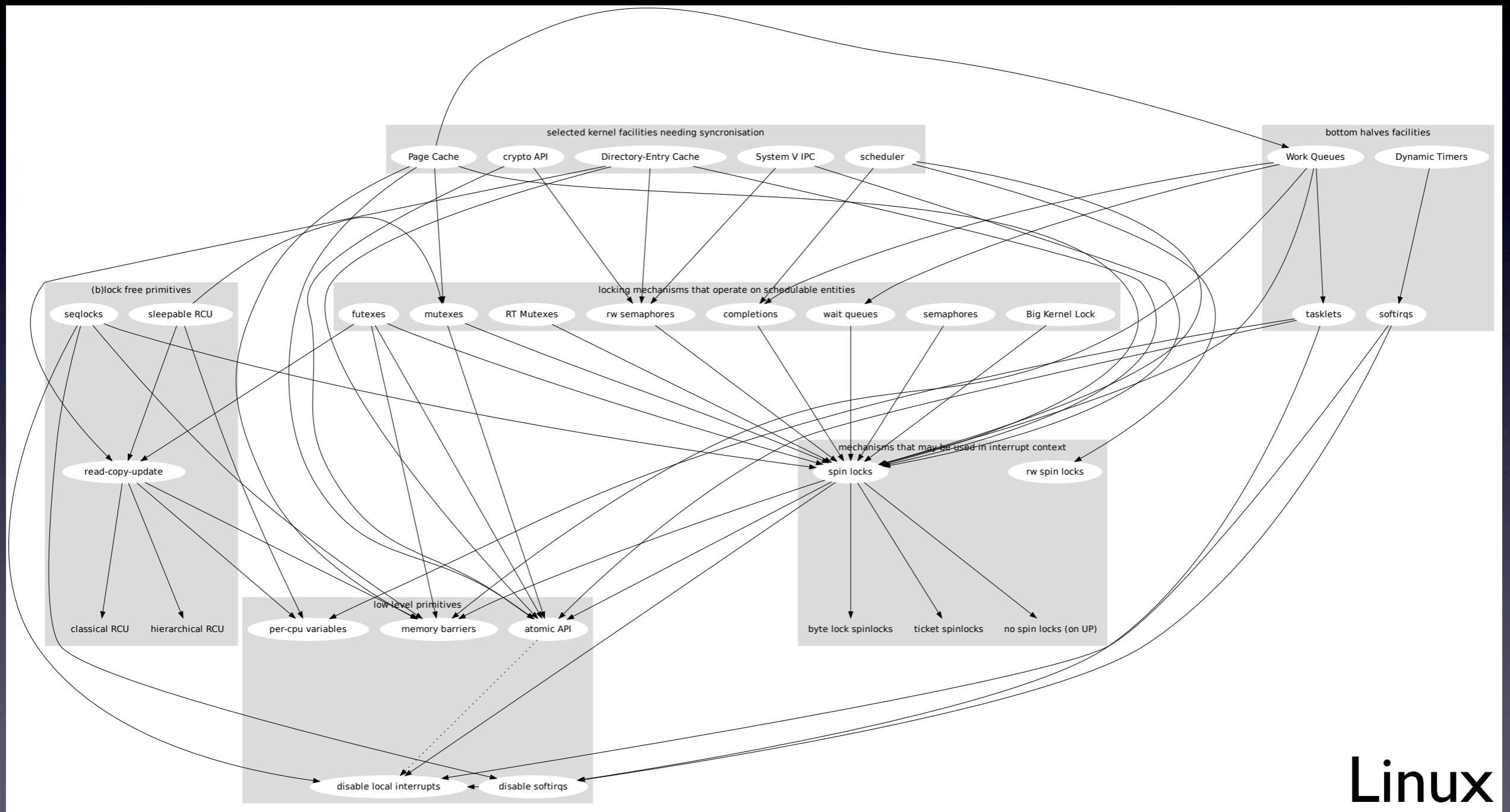
Linux

Untiefen...

FreeBSD



Untiefen...

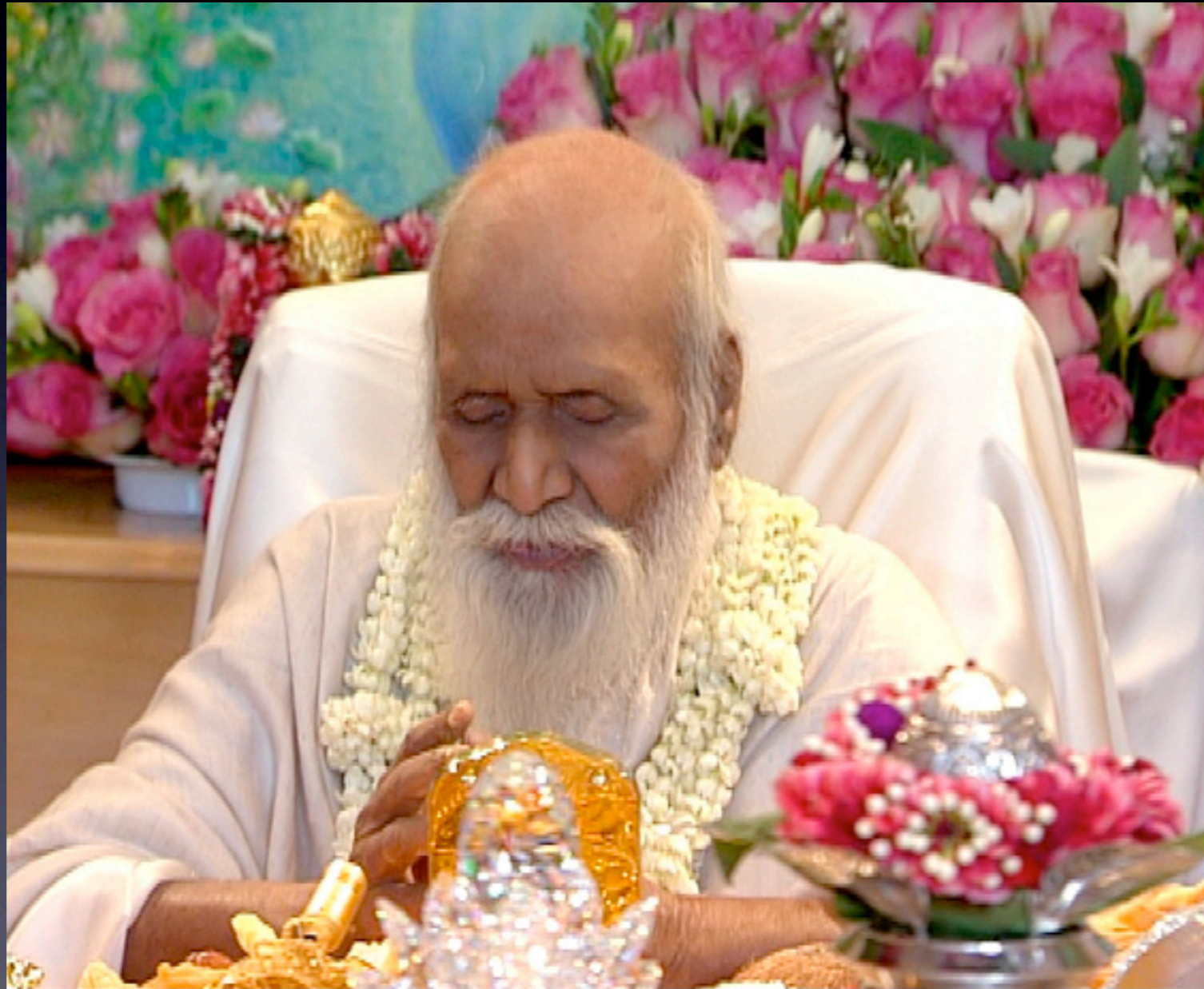


Linux

Fortschrittsgarantien

- behinderungsfrei (*obstruction-free*)
 - außer Konkurrenz, Fortschritt garantiert
- sperrfrei (*lock-free*)
 - Fortschritt des Systems garantiert
- wartefrei (*wait-free*)
 - Fortschritt jedes Fadens garantiert

TM



Transactional Memory

- Hardware (1986)

HTM

- Sun Rock (ASPLOS'09) ?



- Software (1995)

STM

- auf Basis konventioneller Primitiven

- Hybrid (2005)

- HTM, solange Ressourcen es hergeben
- STM, sonst...

(2009)

Transaktionen...

```
void stm_ahead (chain_t *this, chain_t *item) {  
    transactional {  
        item->link = this->link;  
        this->link = item;  
    }  
}
```

```
chain_t *stm_strip (chain_t *this) {  
    chain_t *node;  
    transactional {  
        node = this->link;  
        if (node != 0)  
            this->link = node->link;  
    }  
    return node;  
}
```


...considered harmful!

```
void ewd_prolaag (semaphore_t *this) {  
    transactional {  
        if (this->load-- <= 0)  
            sad_sleep(&this->line);  
    }  
}
```



```
void ewd_verhoog (semaphore_t *this) {  
    transactional {  
        if (this->load++ < 0)  
            sad_awake(&this->line);  
    }  
}
```



Pro und Contra TM

- geeignet für „einfache“ kritische Abschnitte
- komponierbar und verklemmungsfrei: wie CAS bzw. LL/SC, aber semantisch mächtiger
- Nutzung in Betriebssystemen gestaltet sich nicht einfacher als mit CAS bzw. LL/SC
- erfordert Hardwareunterstützung — und fordert einiges an Hardwareressourcen

Latenzverbergung

- altbekannte Technik im Bereich HPC
 - Minimierung von Startzeiten (IPC, E/A)
 - asymmetrische Mehrprozessorsysteme
 - weitgehend ungenutzt für Echtzeitbetrieb
 - deterministische Programmabläufe
- Prozessorkerne für Betriebssystemzwecke

Verwendung der Kerne

Mehrkernprozessor

— Universalkerne

— Spezialkerne

— Wandelkerne

— Fixkerne

— Treiberkerne

— Gerätekerne

— Unterbrecherkerne

— Dienstkerne

— Einplanungskern

— Ein-/Ausgabekern

— Kommunikationskern

— Seitenumlagerungskern

⋮

Balancierung allgemeiner Systemlast

Kaschierung systembedingter Latenzen

„horizontale“ Parallelität (Systemaufrufe)

„vertikale“ Parallelität (Fließbandverfahren)

abgesetzte Geräteprogrammierung

Steuerung der Peripherie

Annahme von (asynchronen) Unterbrechungen

abgesetzte Dienstleistung

Zustellung von Arbeitsaufträgen

Operationen auf Dateien

Abwicklung des Protokollstapels

„external pager“ (Ersetzungsstrategie)

Aktuelle Entwicklungen

- Corey (MIT, Exokern-artig)
- Barrelfish (ETH, Mikrokern-artig)
- Factored OS (MIT, Mikrokern-artig)
- TxLinux (UT, Monolith) — Simulation
- Helios (MSR, Mikrokern-artig)
- TxOS (UT, Monolith)

TM

Resümee

- Multi-Core — des Kaisers neue Kleider!?
- Wettstreitigkeit vorbeugen/vermeiden
- viel stärker optimistische Steuerung von Nebenläufigkeit in Erwägung ziehen
- Systemlatenzen mittels Kernen verbergen
- TM kein Allheilmittel (für Systemsoftware)

