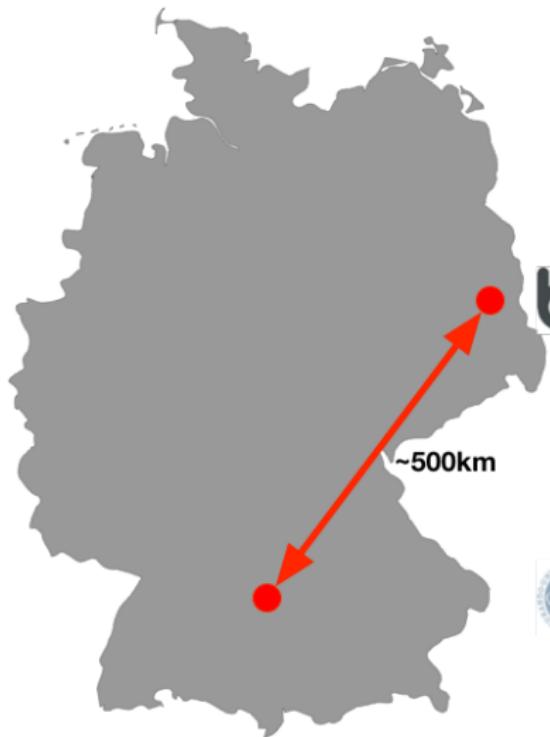March 3, 2017

# SCALABLE OPERATING SYSTEMS
FOR MANY THREADS ON MANY-CORES

Stefan Wesner, Lutz Schubert, Vladimir Nikolov,
Stefan Bonfert, Jörg Nolte, Robert Kuban, Randolf Rotta, . . .

~500km

## Cottbus: Verteilte Systeme & Betriebssysteme



- Jörg Nolte
- Randolf Rotta
- Robert Kuban
- Philipp Gypser
- ...

## Ulm: Organisation und Management von Informationssystemen



- Stefan Wesner
- Lutz Schubert
- Vladimir Nikolov
- Stefan Bonfert
- ...

# RESEARCH FIELDS

## Cottbus

- operating systems, distributed & parallel middleware
- ... for multi- and many-core processors:
  memory consistency, NVRAM, 100Gb/s End2End wireless
- ... and for embedded wireless ad-hoc networks:
  self-stabilization, energy management, event driven processing

## Ulm

- programming models for heterogeneous systems (embedded and HPC)
- scheduling in dynamic systems
- adaptive resource management
- real-time scheduling

# MYTHOS: MANY THREADS OPERATING SYSTEM

- BMBF HPC3 call
  „Anwendungsorientierte HPC-Software für skalierbare Parallelrechner"
- Oct. 2013 – Sep. 2016
- `https://github.com/ManyThreads/mythos`

## Application Partners

- HLRS Stuttgart: Michael Resch - molecular dynamics
- University Siegen: Sabine Roller - fluid dynamics
- Alcatel Lucent: Peter Domschitz - distributed media processing

# APPLICATION SCENARIOS & REQUIREMENTS

## Scientific Computing

- **-** single application owns all cores & memory
- **-** user-level scheduling of application's tasks
- $\Rightarrow$ fast thread suspend and wake-up instead of polling
- $\Rightarrow$ huge number of threads manipulate shared address space

## Media Processing Cloud

- **-** isolated 3rd-party compute tasks
- **-** supervisor configures tasks, communicates data,
  minimizes number of actives cores & processors
- $\Rightarrow$ dynamic creation, migration and teardown
  of simple protection domains & threads

## PROJECT GOALS/CONSTRAINTS

ThOS

### Target Hardware

- Intel XeonPhi Knights Corner:
  PCIe card, 60 cores, 240 threads, 8GiB memory
- use host CPU (Linux) for communication instead of virtualization
- exploit shared memory and caches

### Minimal Base Kernel

- focus on threads and address spaces
- avoid in-kernel policies, heuristics, magic, . . .
- spatial scheduling: thread placement instead of time sharing
- no POSIX interfaces (yet)

ThOS

## Scalability Bottlenecks from Hidden Sharing

- **–** e.g. memory pools, translation tables, queues, services. . .
- **–** automatic partitioning is wasteful, complex, inflexible
- ⇒ enable explicit sharing, reward localized usage

## Active Waiting wastes Time * Power

- **–** sharing implies synchronization
- **–** assume high latency, high contention
- ⇒ use waiting time for other tasks or sleep,
  reward asynchronous / event-driven applications

ThOS

## Advantages of Sharing

- **-** efficient use of memory (no replication)
- **-** fast HW support for atomic updates
- **-** flexible reconfiguration, migration

## Disadvantages

- **-** HW bottlenecks (directory contention, RFO latency)
- **-** races, inconsistent concurrent operations
- **-** danger of false sharing

$\Rightarrow$ trade-off based on application
  MyThOS: mitigate disadvantages, reward localized usage

ThOS

Tasklets between hardware threads: fire&forget

- small messages: next ptr, function ptr, arguments
- each core has 2 queues: incoming and local tasklets
- if empty return to user mode or sleep; wakeup IPI on first push

Location-Based Serialization

- all objects assigned to specific HW threads
  send method calls as tasklet to *correct* queue
- ⇒ manual placement of responsibilities
  fragile usage (e.g. concurrent deletion)

## Delegating Object Monitors

- method calls encoded in tasklets, supplied by caller
- object monitor serializes them
- delegation chooses responsible HW thread
- tasklet is returned via response method to caller's object

## Deferred Synchronous

- tasklets are statically allocated during object creation
- flow control: block next requests until internal tasklets have returned
- # of tasklets = # of logical control flows

## EXAMPLE CALL

```
Adder* addServer = ...;
Tasklet myTask;
addServer->add(&myTask, this, 47, 11);

class Adder {
 Monitor monitor;
public:
 void add(Tasklet* t, AddResult* res, int a, int b) {
   monitor.request(t,
           [=](Tasklet* t){this->addImpl(t,res,a,b);} );
 }

private:
 void addImpl(Tasklet* t, AddResult* res, int a, int b) {
   res->result(t, a+b);
   monitor.requestDone();
 }
};
```

## EXAMPLE CALL

```
Adder* addServer = ...;
Tasklet myTask;
addServer->add(&myTask, this, 47, 11);

class Adder {
 Monitor monitor;
public:
 void add(Tasklet* t, AddResult* res, int a, int b) {
   monitor.request(t,
           [=](Tasklet* t){this->addImpl(t,res,a,b);} );
 }

private:
 void addImpl(Tasklet* t, AddResult* res, int a, int b) {
   res->result(t, a+b);
   monitor.requestDone();
 }
};
```

## EXAMPLE CALL

```
Adder* addServer = ...;
Tasklet myTask;
addServer->add(&myTask, this, 47, 11);

class Adder {
 Monitor monitor;
public:
 void add(Tasklet* t, AddResult* res, int a, int b) {
   monitor.request(t,
           [=](Tasklet* t){this->addImpl(t,res,a,b);} );
 }

private:
 void addImpl(Tasklet* t, AddResult* res, int a, int b) {
   res->result(t, a+b);
   monitor.requestDone();
 }
};
```

## EXAMPLE CALL

```
Adder* addServer = ...;
Tasklet myTask;
addServer->add(&myTask, this, 47, 11);

class Adder {
 Monitor monitor;
public:
 void add(Tasklet* t, AddResult* res, int a, int b) {
   monitor.request(t,
            [=](Tasklet* t){this->addImpl(t,res,a,b);} );
 }

private:
 void addImpl(Tasklet* t, AddResult* res, int a, int b) {
   res->result(t, a+b);
   monitor.requestDone();
 }
};
```

## EXAMPLE CALL

```
Adder* addServer = ...;
Tasklet myTask;
addServer->add(&myTask, this, 47, 11);

class Adder {
 Monitor monitor;
public:
 void add(Tasklet* t, AddResult* res, int a, int b) {
   monitor.request(t,
           [=](Tasklet* t){this->addImpl(t,res,a,b);} );
 }

private:
 void addImpl(Tasklet* t, AddResult* res, int a, int b) {
   res->result(t, a+b);
   monitor.requestDone();
 }
};
```

ThOS

### Problem

- pending tasklets while user requests deletion
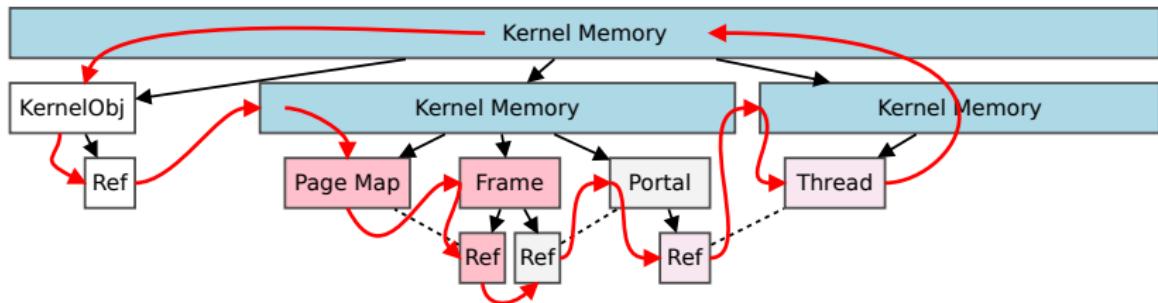
### Solution

- reference counting delays removal until safe [Magenta]
- smart pointers enable forced removal of references, hierarchy of dependent pointers [like OC tree]
- fence multicast over all active HW threads resolves temporary references [like RCU gc]
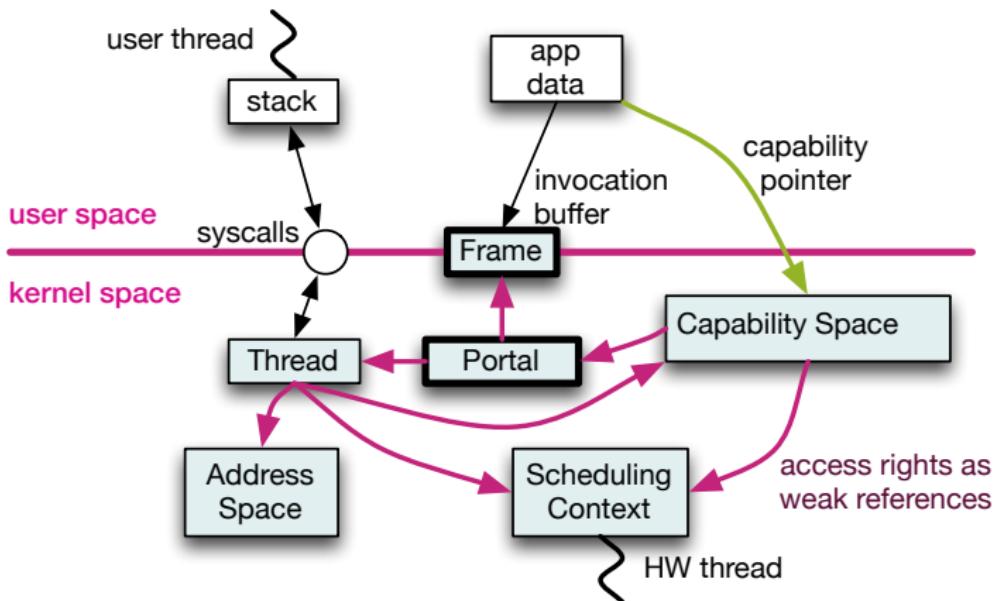
- range of unused memory for new kernel objects,
  can be split $\Rightarrow$ user can distribute load over cores
- MyThOS: atomic ops and locking just on neighbor nodes,
  local synchronization on partitioned memory,
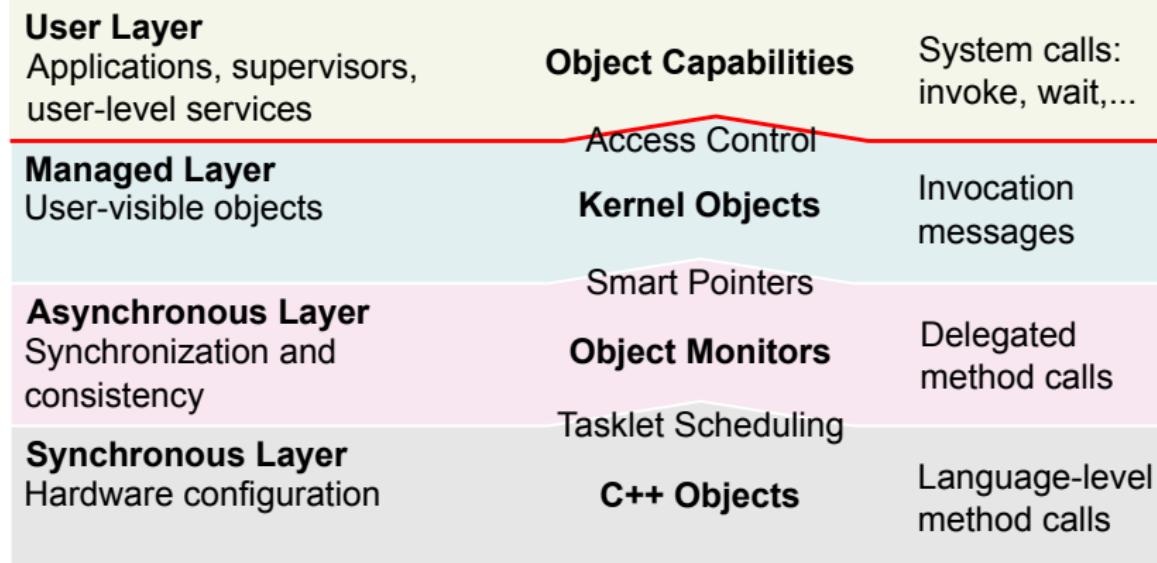  detect user's race conditions

**ASYNCHRONOUS CAPABILITY INVOCATIONS**
**Object Capabilities [seL4, Fiasco.OC], Portals separated from Threads**

## LAYERED ARCHITECTURE

myThOS

| | | |
|---|---|---|
| **User Layer**<br>Applications, supervisors,<br>user-level services | **Object Capabilities** | System calls:<br>invoke, wait,... |
| | Access Control | |
| **Managed Layer**<br>User-visible objects | **Kernel Objects** | Invocation<br>messages |
| | Smart Pointers | |
| **Asynchronous Layer**<br>Synchronization and<br>consistency | **Object Monitors** | Delegated<br>method calls |
| | Tasklet Scheduling | |
| **Synchronous Layer**<br>Hardware configuration | **C++ Objects** | Language-level<br>method calls |

ThOS

Challenge: eliminate hidden sharing, reduce active waiting

Solution: Delegation, Object Monitors, Explicit Kernel Memory

|                        | on-demand allocation | explicit sharing |
|------------------------|----------------------|------------------|
| big kernel lock        | ?                    | seL4             |
| fine-grained locking   | Fiasco.OC, Linux     | MyThOS           |
| message passing        | ?                    | Barrelfish       |

Results

- chaotic recursive fork/join: doesn't crash (too often)
- core sleep too eager + high IPI latency $\Rightarrow$ needs tuning