

# HermitCore: A Library Operating System for Cloud and HPC

Jens Breitbart<sup>0</sup>, Stefan Lankes<sup>1</sup>, Simon Pickartz<sup>1</sup>

<sup>0</sup>Bosch Chassis Systems Control, Stuttgart, Germany

<sup>1</sup>RWTH Aachen University, Aachen, Germany

# Agenda

---

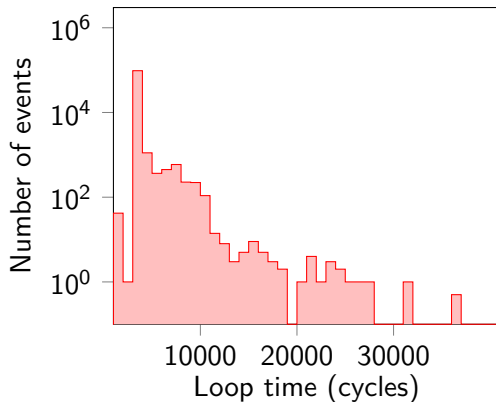
- Motivation
- Challenges for the Future Systems
- OS Architectures
- HermitCore Design
- Performance Evaluation
- Conclusion and Outlook

- Complexity of high-end HPC systems keeps growing
  - ≡ Extreme degree of parallelism
  - ≡ Heterogeneous core architectures
  - ≡ Deep memory hierarchy
  - ≡ Power constrains
    - ⇒ Need for scalable, reliable performance and capability to rapidly adapt to new HW
- Applications have also become complex
  - ≡ In-situ analysis, workflows
  - ≡ Sophisticated monitoring and tools support, etc. . .
  - ≡ Isolated, consistent simulation performance
    - ⇒ Dependence on POSIX, MPI and OpenMP
- Seemingly contradictory requirements. . .

## Operating System = Overhead

---

- Every administration layer has its overhead  $\Rightarrow$  e. g. Hourglass benchmark



- OS noise reduce the scalability / increases latency

# How does the HPC community reduce the overhead?

---

## Light-weight Kernels

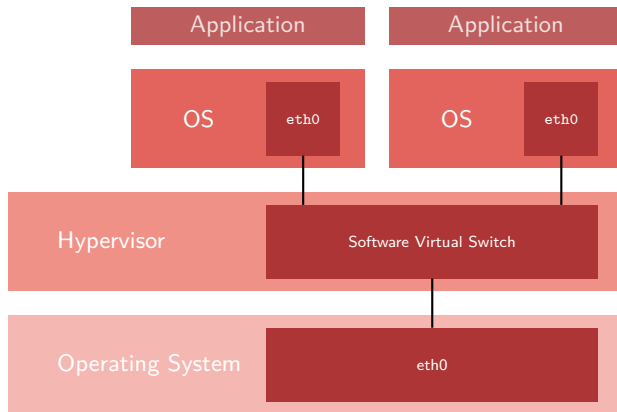
- Typically taken of an existing fat kernel (e. g. Linux)
- Removal of unneeded features to improve the scalability
- e. g. ZeptoOS

## Multi-Kernels

- A specialized kernel runs side-by-side to full-weight kernel (e. g. Linux)
- Applications of the full-weight kernels are able to run on the specialized kernel.
- The specialized kernel catch every system call and delegate them to the full-weight kernel
  - ≡ Binary compatible to the full-weight kernel
- Examples: mOS, McKernel

## OS Designs for Cloud Computing – Usage of Common OS

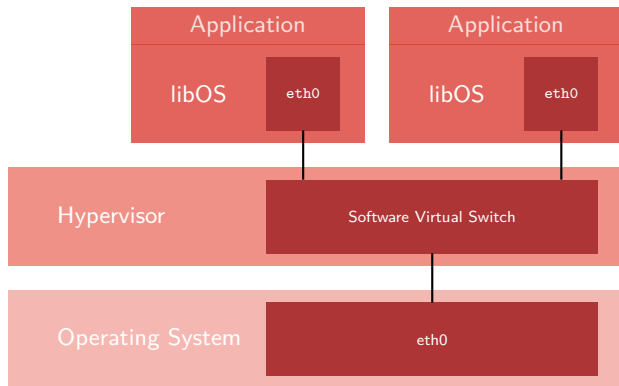
---



- Two operating systems to maintain a single computer? Double Management!
- Why should one run a Multi-User-/Multi-Tasking OS within a hypervisor in the Cloud for a simple webserver?

# OS Designs for Cloud Computing – LibraryOS

---



- Now, every system call is a function call  $\Rightarrow$  Low overhead
- Whole optimization of an image possible (including the library OS)
  - ≡ Link Time Optimization (LTO)
- Removing unneeded code  $\Rightarrow$  reduces the attack vector

- Combination of the Unikernel and Multi-Kernel to reduce the overhead
  - ≡ The same binary is able to run
    - = in a VM (classical unikernel setup)
    - = or bare-metal side-by-side to Linux (multi-kernel setup)
- Support for dominant programming models (MPI, OpenMP)
- Single-address space operating system
  - ≡ No TLB Shutdown



## Classical Unikernel Setup

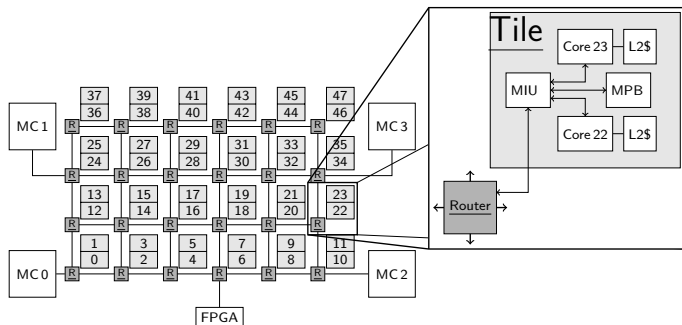
- Applications are able to boot directly within a VM
  - ≡ Tested with Qemu / KVM
  - ≡ Tested with uhyve (experimental KVM-based Hypervisor)
    - = Qemu emulates more than HermitCore needs ⇒ large setup time
    - = uhyve reduce the boot time (from 2s to ~30ms)
  - ≡ Google Compute Engine

## Multi-Kernel Setup

- One kernel per NUMA node
  - ≡ Only local memory accesses (UMA)
  - ≡ Message passing between NUMA nodes
- One FWK (Linux) in the system to get access to a broader driver support
  - ≡ Only a backup for pre- / post-processing
  - ≡ Critical path should be handled by HermitCore

# Runtime Support

- SSE, AVX2, AVX512, FMA,...
- Full C-library support (newlib)
- HBM support similar to memkind
- IP interface & BSD sockets (LwIP)
  - ≡ IP packets are forwarded to Linux
  - ≡ Shared memory interface
- Pthreads
  - ≡ Thread binding at start time
  - ≡ No load balancing ⇒ less housekeeping
- OpenMP via Intel's Runtime
- iRCCE- & MPI (via SCC-MPICH)
- Full support for the Go runtime



## Is the software stack difficult to maintain?

---

- Changes to the common software stack determined with cloc

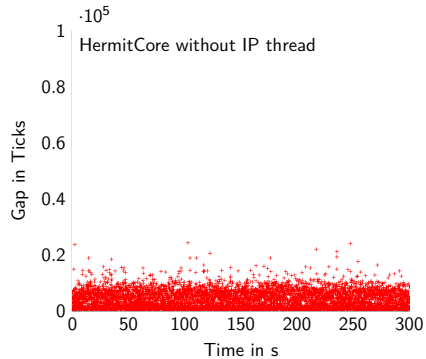
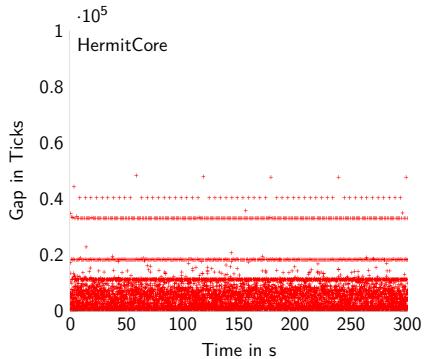
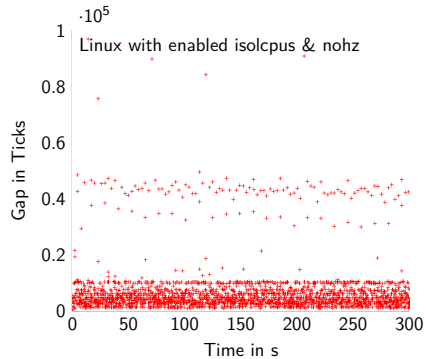
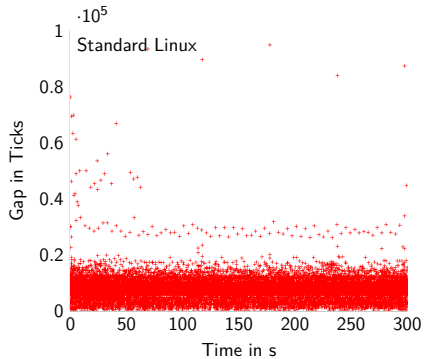
<b>Software Stack</b>	<b>LoC</b>	<b>Changes</b>
binutils	5 121 217	226
gcc	6 850 382	4 821
Linux	15 276 013	1 296
Newlib	1 040 826	5 472
LwIP	38 883	832
Pthread	13 768	466
OpenMP RT	61 594	324
HermitCore	–	10 597

- Test systems

- ≡ Intel Haswell CPUs (E5-2650 v3) clocked at 2.3 GHz
- ≡ Intel KNL (Phi 7210) clocked at 1.3 GHz, SNC mode with four NUMA nodes

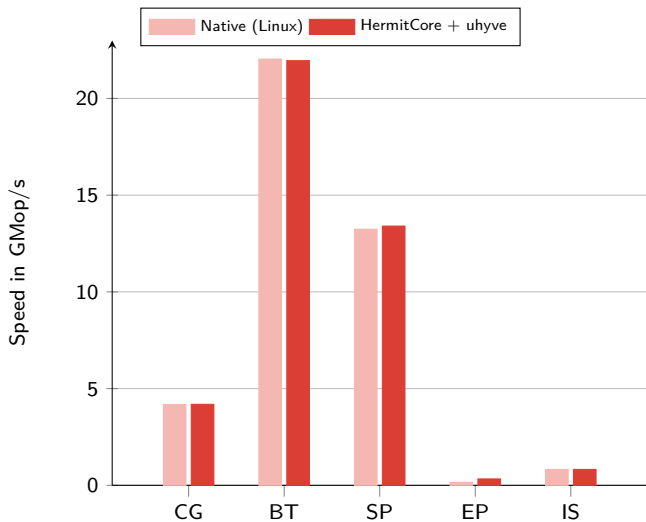
- Results in CPU cycles

System activity	KNL		Haswell	
	HermitCore	Linux	HermitCore	Linux
getpid()	15	486	14	143
sched_yield()	197	983	97	370
malloc()	3 051	12 806	3715	6575
first write access to a page	2 078	3 967	2018	4007



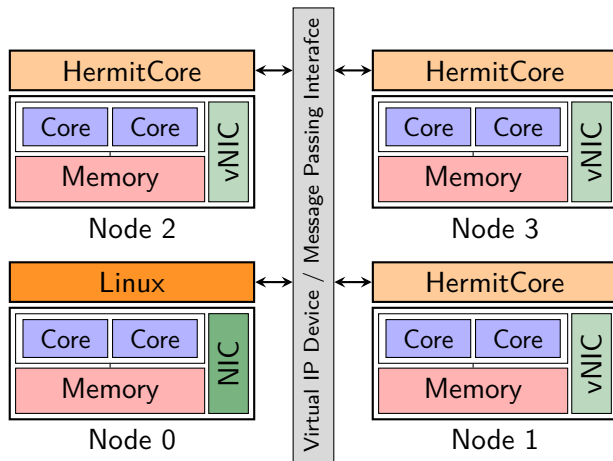
## Overhead of VMs – Determined via NAS Parallel Benchmarks (Class B)

---



# Outlook

- Arm v8 support
- SR-IOV simplifies the coordination between Linux & HermitCore



## Conclusions

---

- It works! ⇒ <https://youtu.be/gDYCJ1D0TKw>
- Binary packages are available
- Reduce the OS noise significantly
- Try it out!

<http://www.hermitcore.org>

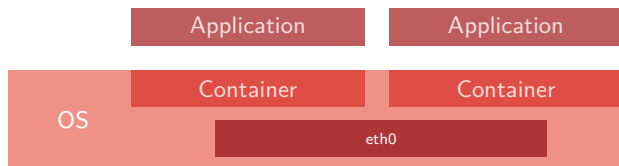
Thank you for your kind attention!



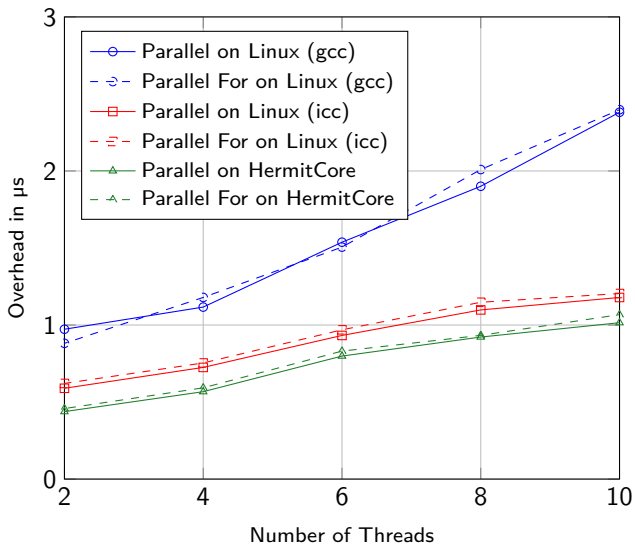
**Backup slides**

# OS Designs for Cloud Computing – Container

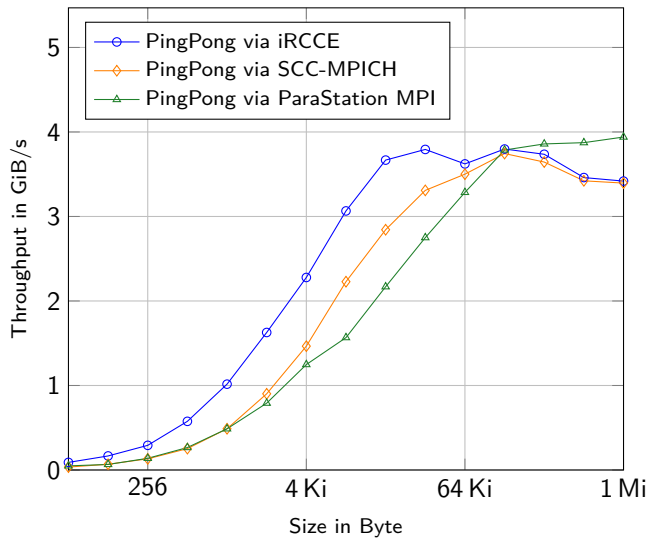
---



- Building of virtual borders (namespaces)
- Containers and their processes doesn't see each other
- Fast access to OS services
- Less secure because an exploit for the container attacks also the host OS
- Doesn't reduce the OS noise of the host system

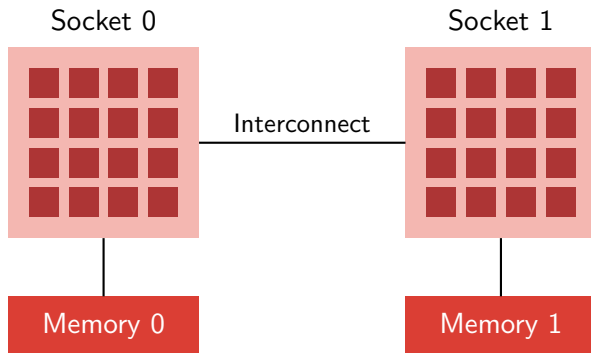


# Throughput Results of the Inter-kernel Communication Layer



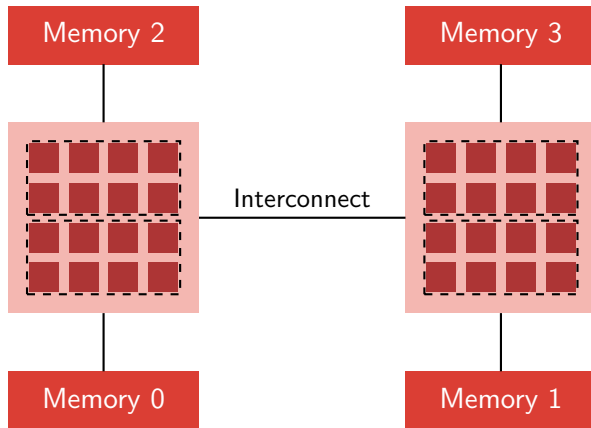
## Non-Uniform Memory Access

- Costs for memory access may vary
- Run processes where memory is allocated
- Allocate memory where the process resides
- Implications for the performance
  - ≡ Where should the applications store the data?
  - ≡ Who should decide the location?
    - = The operating system?
    - = The application developers?

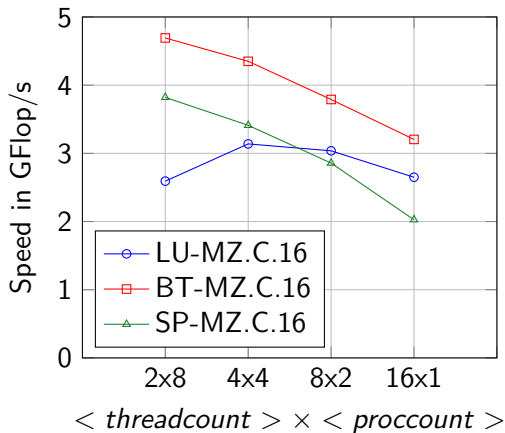


## Non-Uniform Memory Access

- Costs for memory access may vary
- Run processes where memory is allocated
- Allocate memory where the process resides
- Implications for the performance
  - ≡ Where should the applications store the data?
  - ≡ Who should decide the location?
    - = The operating system?
    - = The application developers?



- Parallelization via Shared Memory (OpenMP)
  - ≡ Many side-effects and error-prone
  - ≡ Incremental parallelization
- Parallelization via Message Passing (MPI)
  - ≡ Restructuring of the sequential code
  - ≡ Less side-effects
- Performance Tuning
  - ≡ Bind MPI applications on one NUMA node
  - ⇒ No remote memory access



- GCC includes a OpenMP Runtime (libgomp)
  - ≡ Reuse synchronization primitives of the Pthread library
  - ≡ Other OpenMP runtimes scales better
  - ≡ In addition, our Pthread library was originally not designed for HPC
- Integration of Intel's OpenMP Runtime
  - ≡ Include its own synchronization primitives
  - ≡ Binary compatible to GCC's OpenMP Runtime
  - ≡ Changes for the HermitCore support are small
    - = Mostly deactivation of function to define the thread affinity
  - ≡ Transparent usage
    - = For the end-user, no changes in the build process



## Support of compilers beside GCC

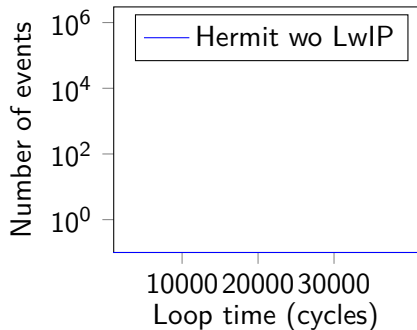
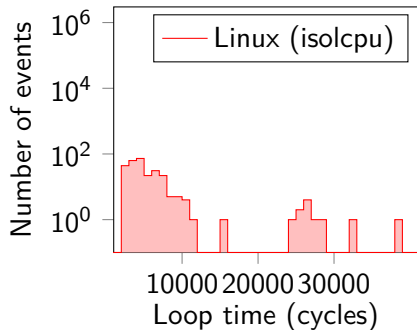
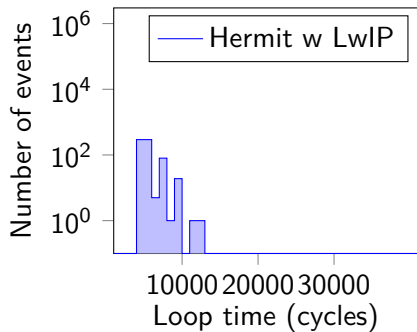
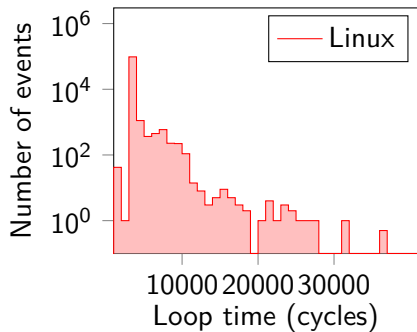
---

- Just avoid the standard environment (`--ffreestanding`)
- Set include path to HermitCore's toolchain
- Be sure that the ELF file use HermitCore's ABI
  - ≡ Patching object files via `elfedit`
- Use the GCC to link the binary

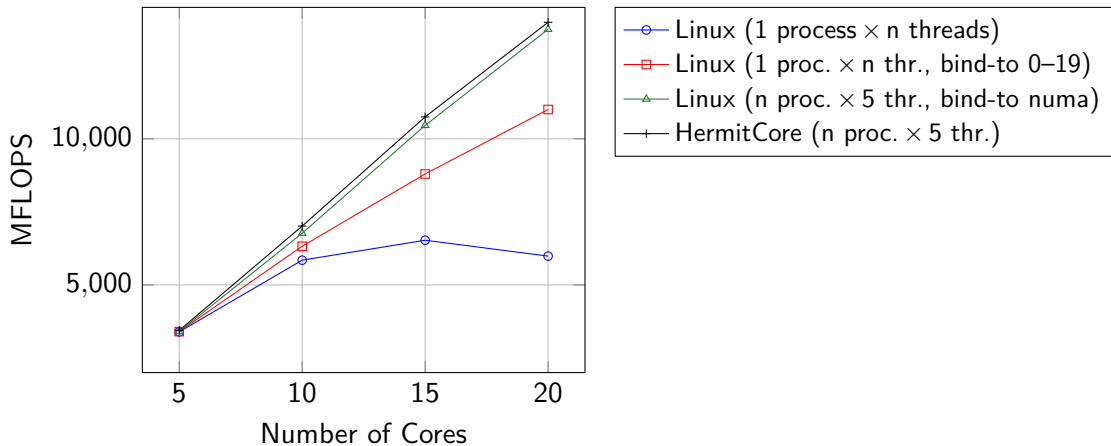
```
LD = x86_64-hermit-gcc
#CC = x86_64-hermit-gcc
#CFLAGS = -O3 -mtune=native -march=native -fopenmp -mno-red-zone
CC = icc -D__hermit__
CFLAGS = -O3 -xHost -mno-red-zone -ffreestanding -I$(HERMIT_DIR) -openmp
ELFEDIT = x86_64-hermit-elfedit
```

```
stream.o: stream.c
    $(CC) $(CFLAGS) -c -o $@ $<
    $(ELFEDIT) --output-osabi HermitCore $@
```

```
stream: stream.o
    $(LD) -o $@ $< $(LDFLAGS) $(CFLAGS)
```



## Hydro (preliminary results)



Thank you for your kind attention!

**Jens Breitbart et al.** – [jens.breitbart@de.bosch.com](mailto:jens.breitbart@de.bosch.com)

[www.jensbreitbart.de](http://www.jensbreitbart.de)

HermitCore logo is provided by EmojiOne.