

# M<sup>3</sup>: Taking Microkernels to the Next Level

Nils Asmussen

FOSDEM, 02/02/2020, Brussels



- Nils Asmussen
- PhD last year at the OS chair of the TU Dresden
- Low-level system programming and microkernels
- Worked on several microkernel-based OSes in the past
  - Escape, own hobby OS (presented here in 2013): <https://github.com/Nils-TUD/Escape>
  - NRE, userland for NOVA: <https://github.com/TUD-OS/NRE>
  - M<sup>3</sup>, presented today: <https://github.com/TUD-OS/M3>
- Since 2019 at the Barkhausen Institut



- Research institute in Dresden, founded end of 2017
- Currently about 30 people
- Low-latency and secure IoT systems
- Focus on research and demonstrators



- Research institute in Dresden, founded end of 2017
- Currently about 30 people
- Low-latency and secure IoT systems
- Focus on research and demonstrators

Wireless

RF Design

Privacy

Lab

MPSoC

OS

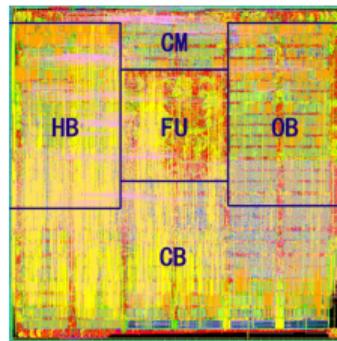


- Microkernel-based systems have proven valuable for several objectives
  - Security
  - Robustness
  - Real time
  - Flexibility
- Recently, new challenges are coming from the hardware side
  - Heterogeneous systems
  - Third-party components
  - Security issues of complex general-purpose cores

# Heterogeneous Systems

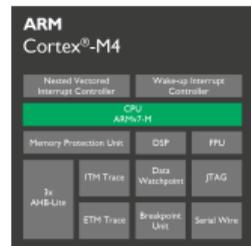
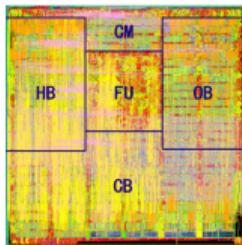


Snapdragon X20 LTE modem	Adreno 630 Visual Processing Subsystem
Wi-Fi	
Hexagon 685 DSP	Qualcomm Spectra 280 ISP
Qualcomm Aqstic Audio	Kryo 385 CPU
System Memory	Qualcomm Mobile Security



- Demanded by performance and energy requirements
- Big challenge for OSeS: single shared kernel on all cores does no longer work
- OSeS need to be prepared for processing elements with different feature sets

# Third-party Components



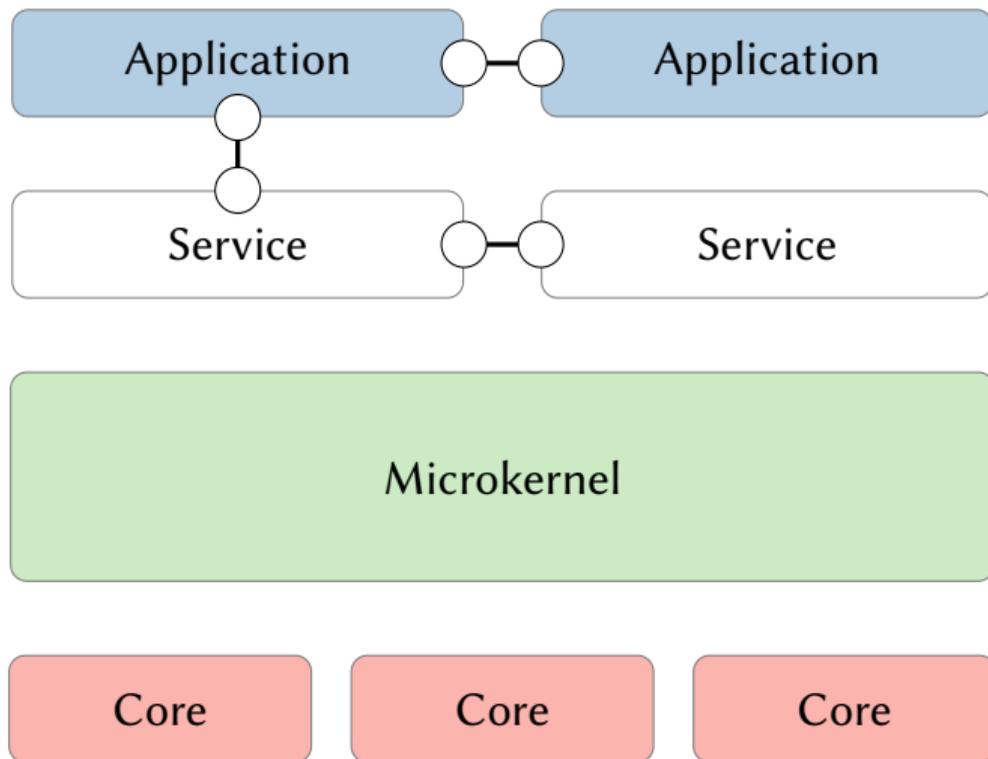
- Market pressure forces us to integrate third-party components
- We should not trust these components
- Currently, often no isolation between them
- Bug in such a component can compromise whole system (see Broadcom incident)

# Security Issues of Complex General-purpose Cores

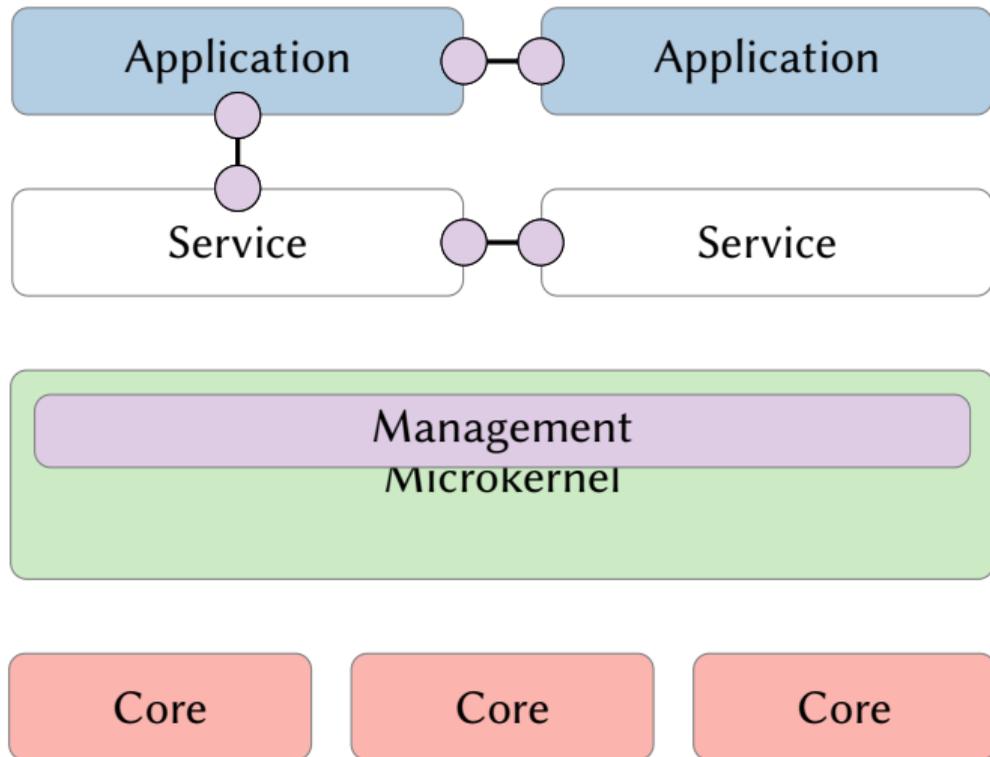


- 20 known attacks (and counting ...)
- Allow to leak private data, sometimes bypassing all security measures of the core
- Mitigations exist, but these are complex and costly
- These security holes have been lurking in CPUs for many years
- Should we still trust these complex cores to properly enforce the isolation between different software components?

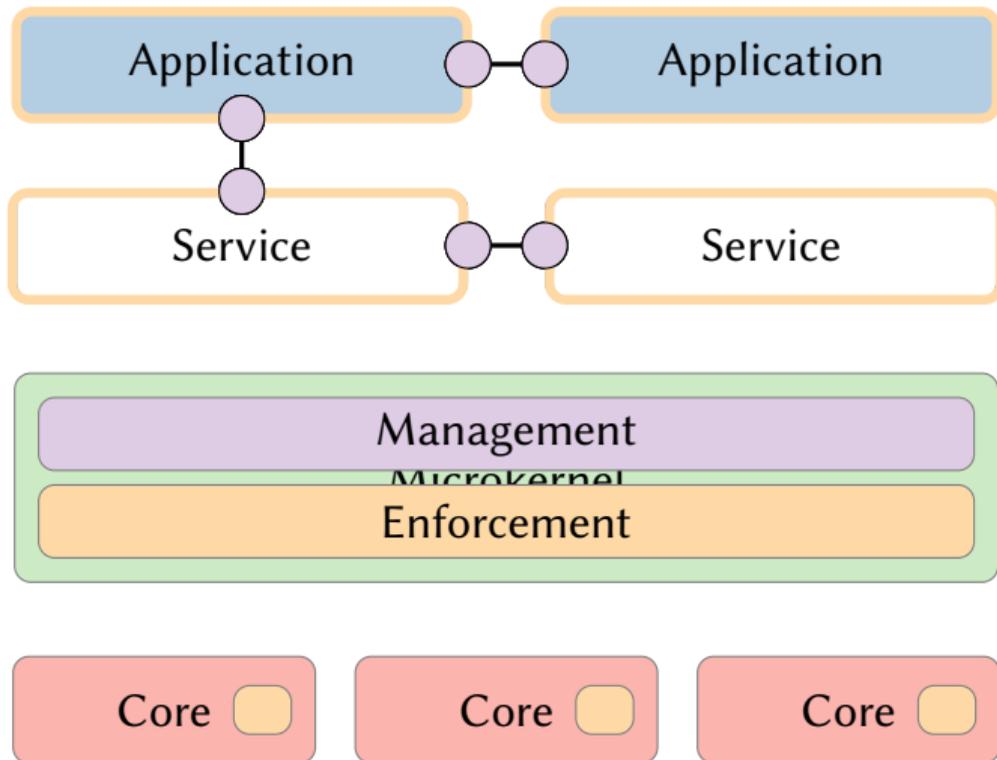
# Microkernel-based System as Foundation



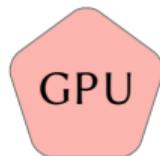
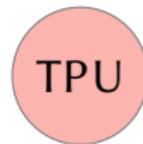
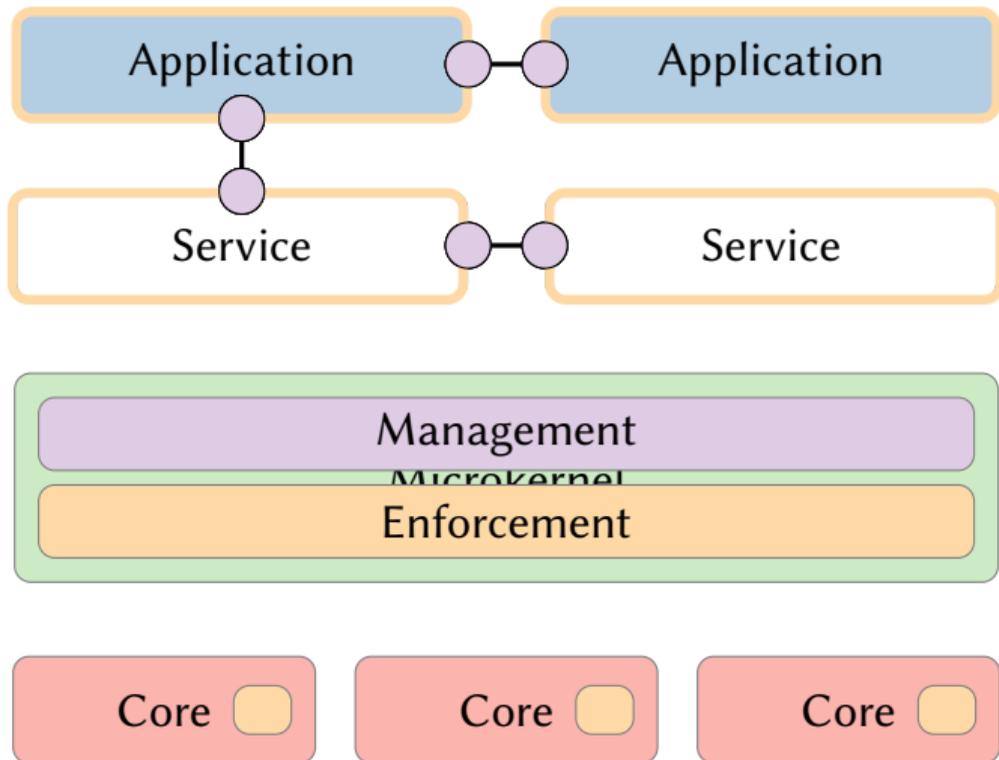
# Microkernel-based System as Foundation



# Microkernel-based System as Foundation



# Microkernel-based System as Foundation





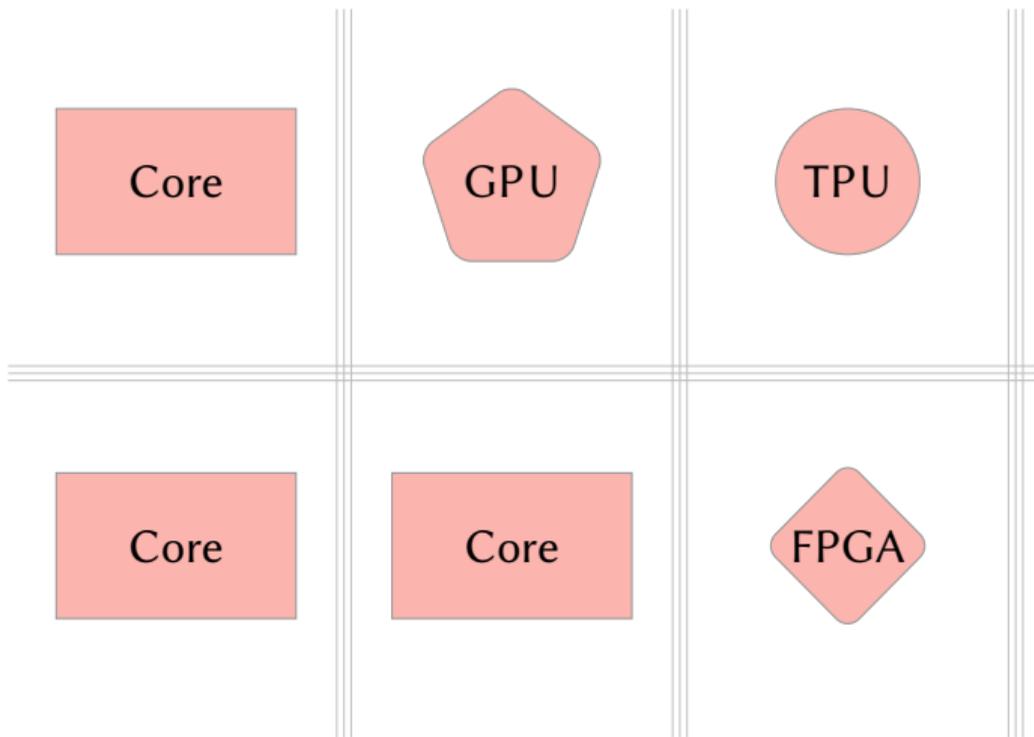
- 1 The New System Architecture
- 2 M<sup>3</sup>: The Operating System
- 3 What are the Benefits?



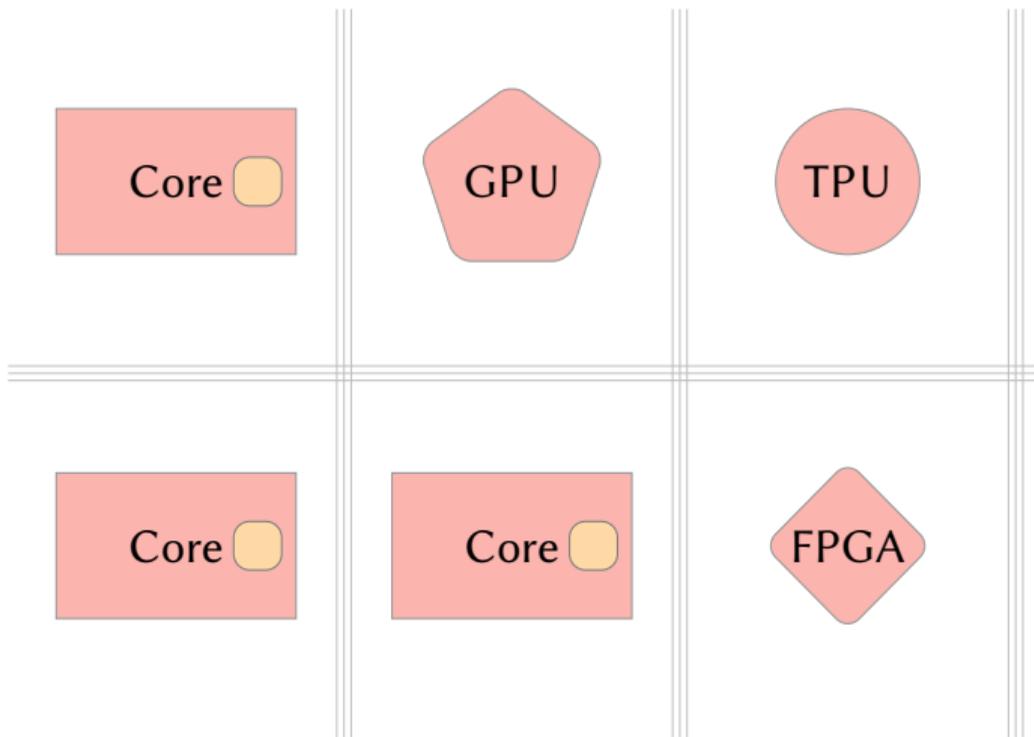
- 1 The New System Architecture
- 2 M<sup>3</sup>: The Operating System
- 3 What are the Benefits?



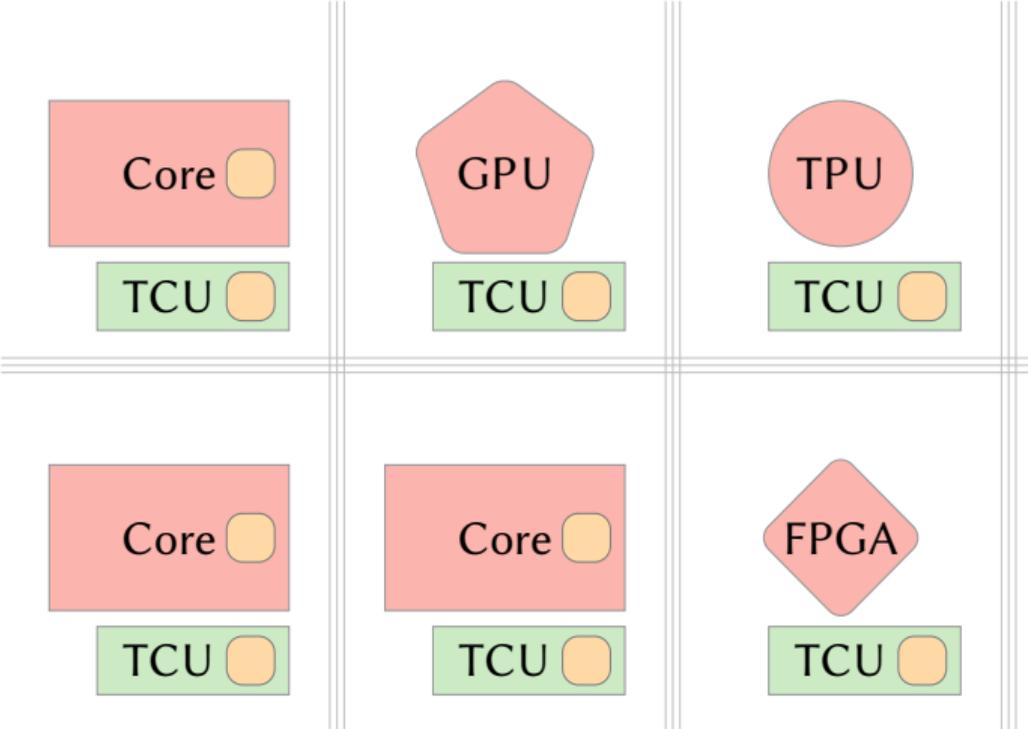
# Hardware/Operating System Co-Design



# Hardware/Operating System Co-Design



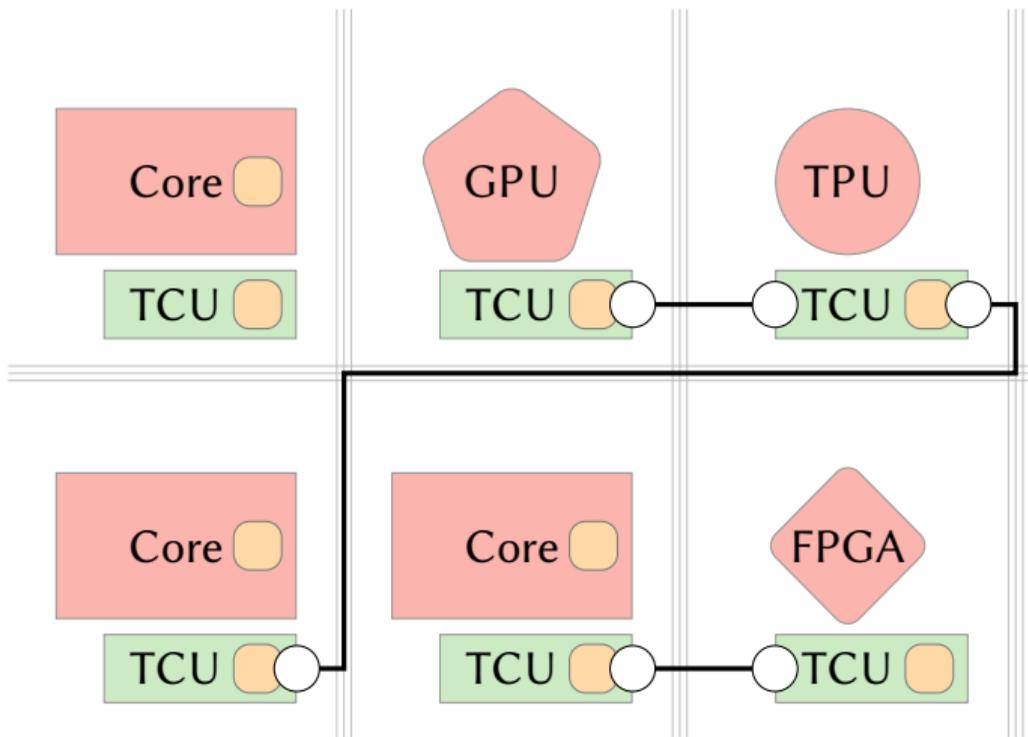
# Hardware/Operating System Co-Design



Key ideas:

- TCU as new hardware component

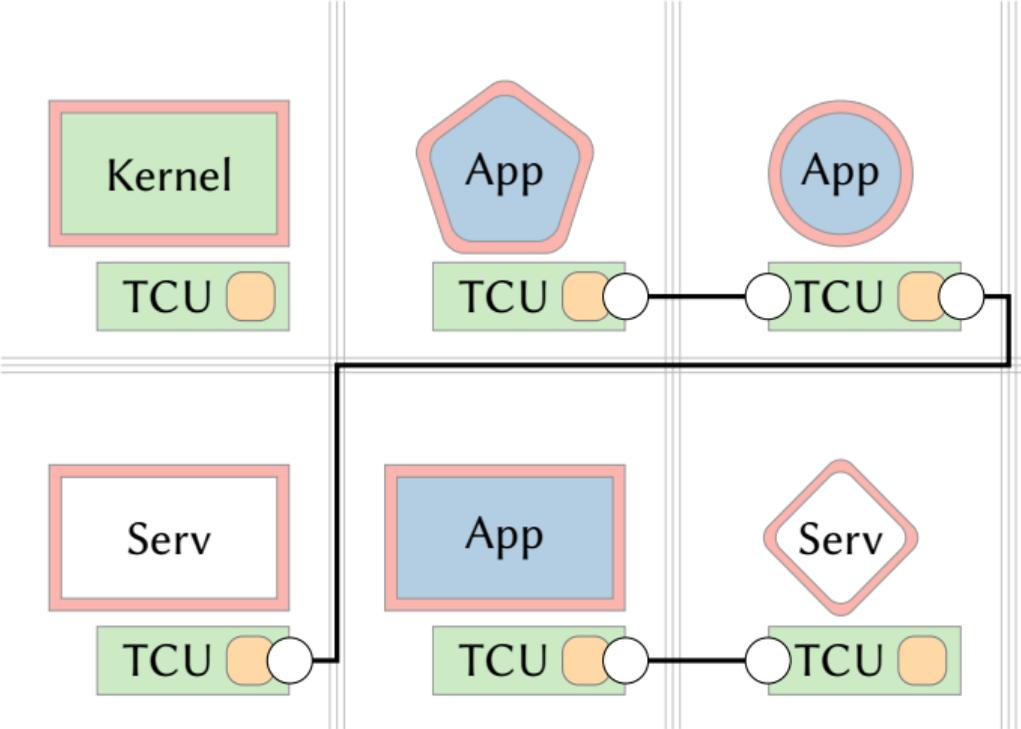
# Hardware/Operating System Co-Design



Key ideas:

- TCU as new hardware component

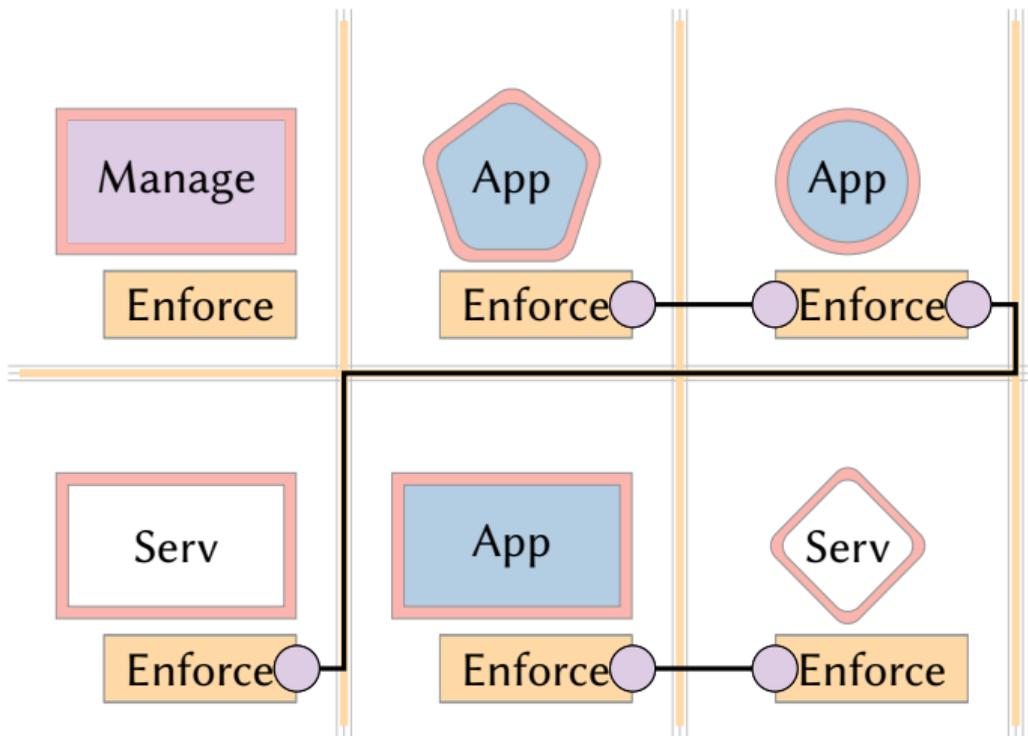
# Hardware/Operating System Co-Design



Key ideas:

- TCU as new hardware component
- Kernel on dedicated PE

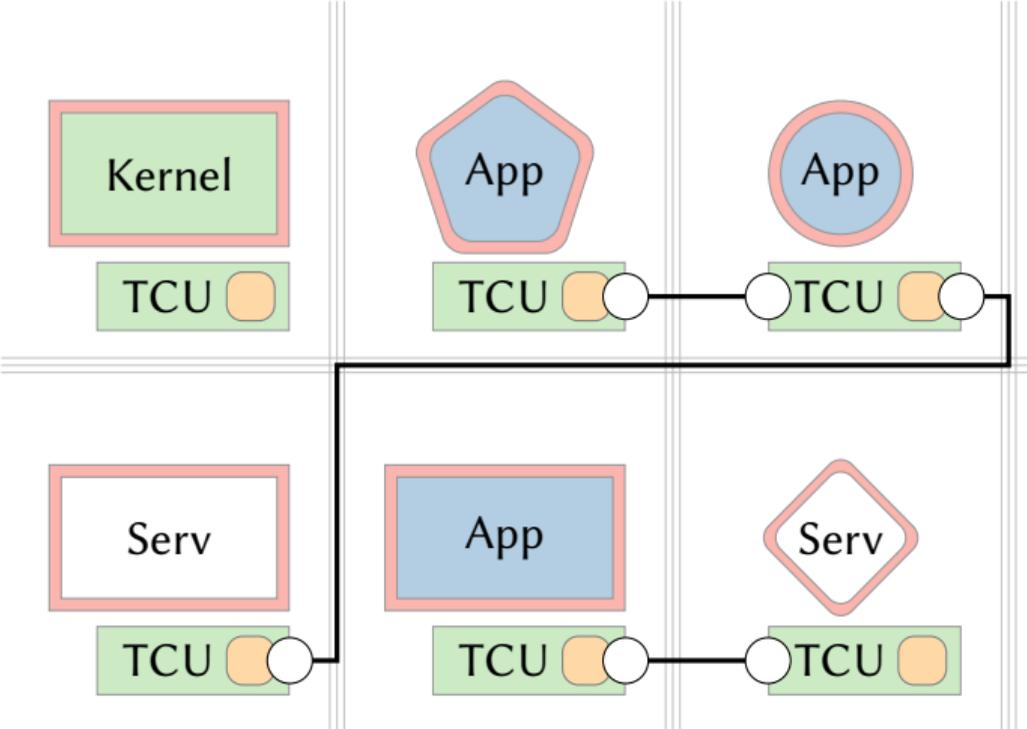
# Hardware/Operating System Co-Design



Key ideas:

- TCU as new hardware component
- Kernel on dedicated PE
- Kernel manages, TCU enforces

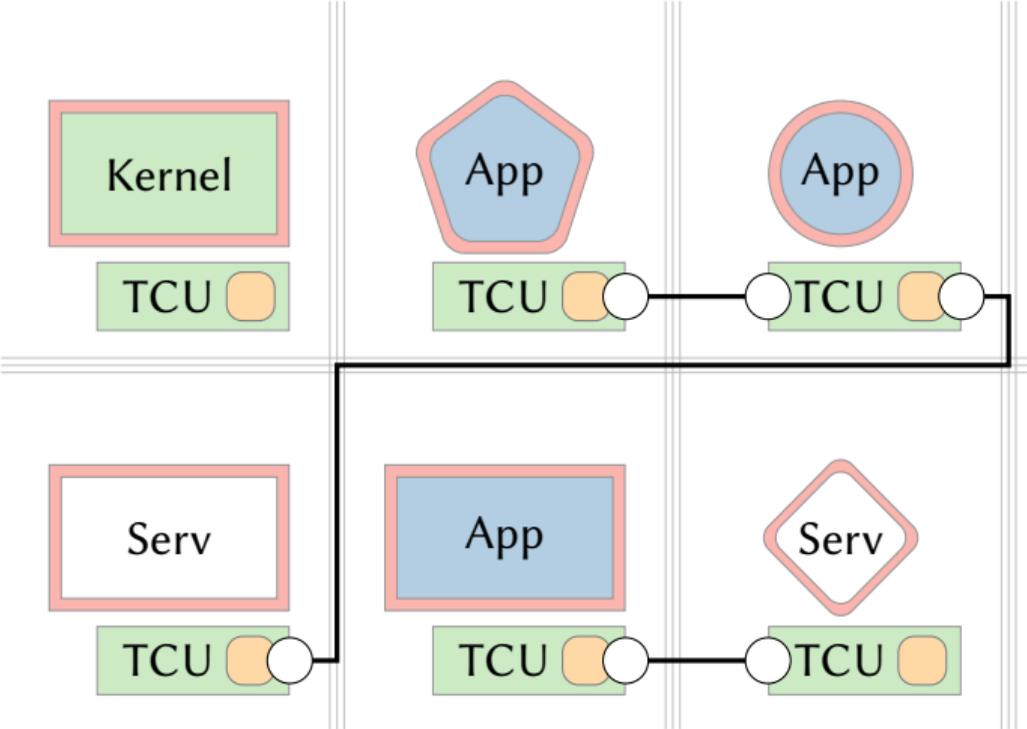
# Hardware/Operating System Co-Design



Takes  $\mu$ -kernels to the next level:

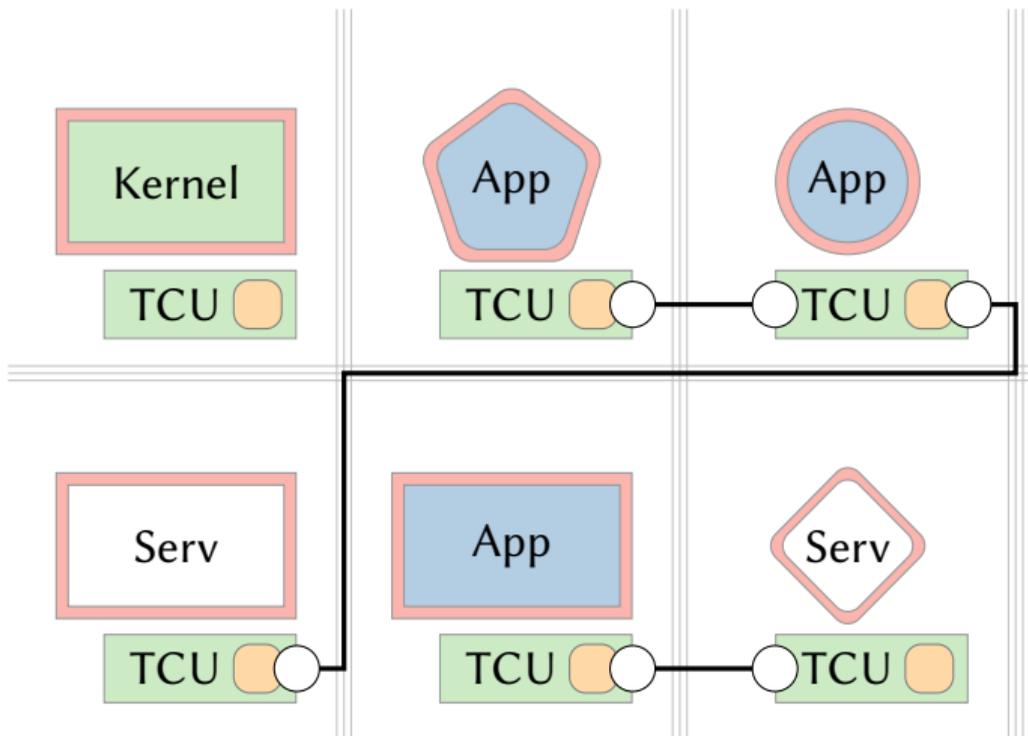
- TCU as secure foundation

# Hardware/Operating System Co-Design



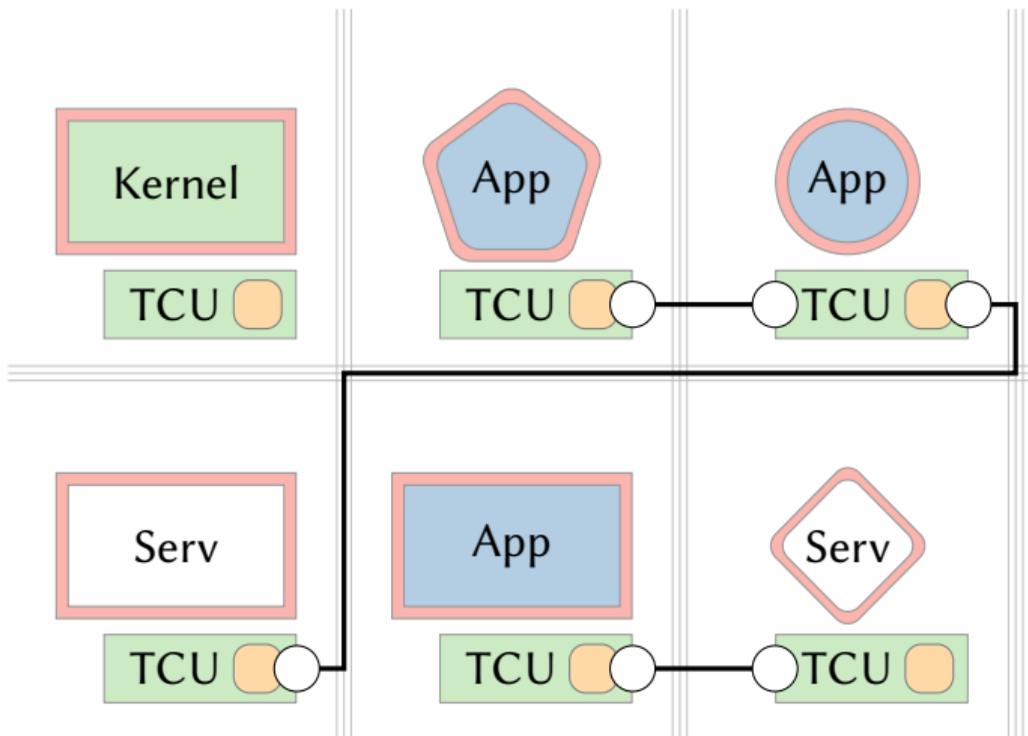
Takes  $\mu$ -kernels to the next level:

- TCU as secure foundation
- Heterogeneity:  
Uniform interface



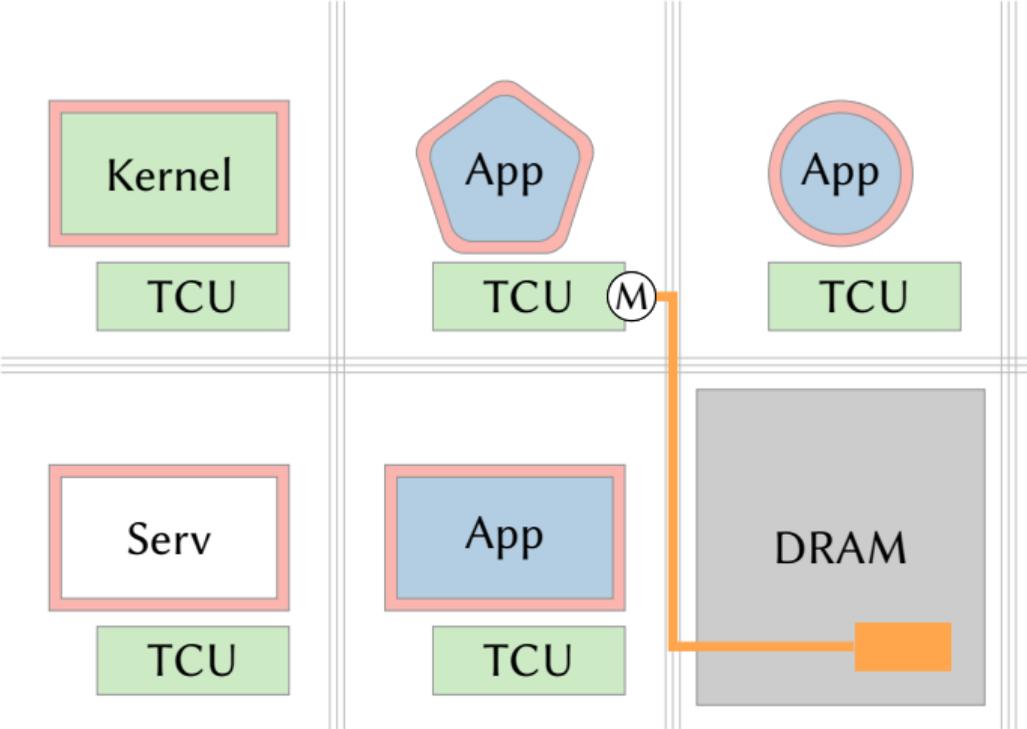
Takes  $\mu$ -kernels to the next level:

- TCU as secure foundation
- Heterogeneity:  
Uniform interface
- Untrusted HW comp.:  
Protected by TCU



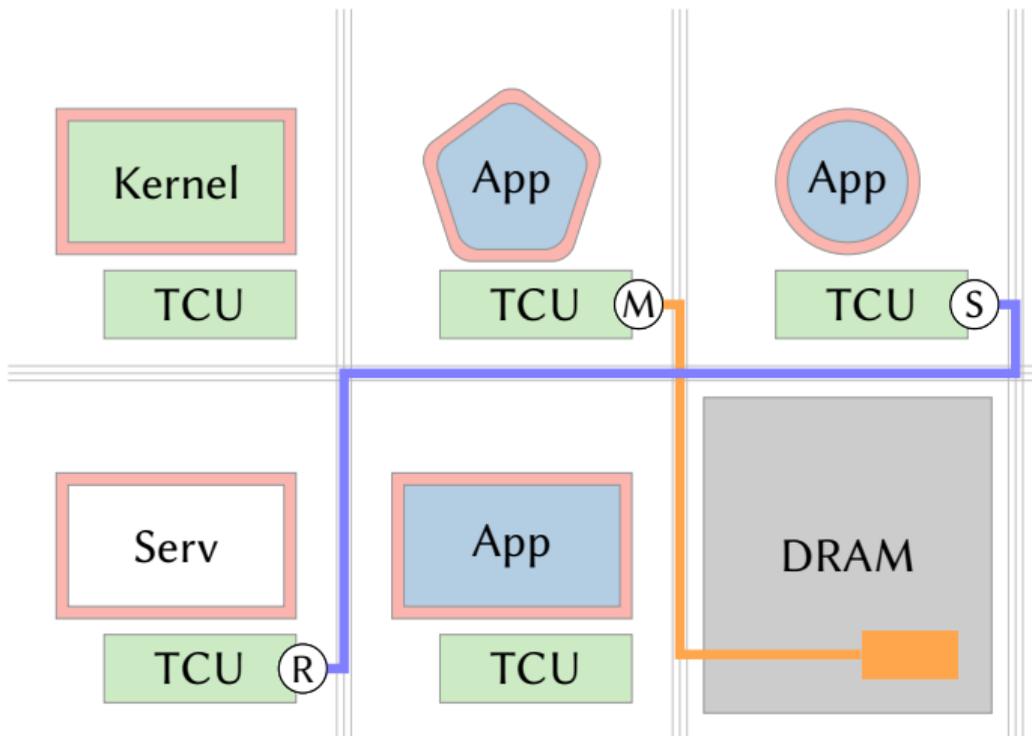
Takes  $\mu$ -kernels to the next level:

- TCU as secure foundation
- Heterogeneity:  
Uniform interface
- Untrusted HW comp.:  
Protected by TCU
- Side channels:  
Physical isolation



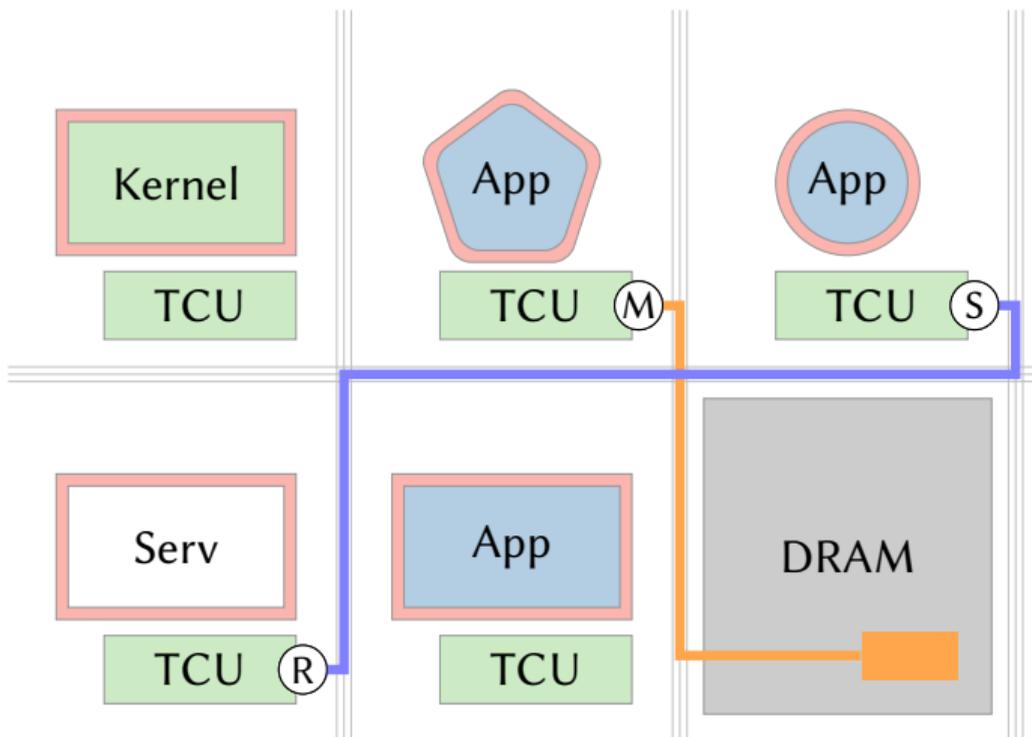
TCU provides *endpoints* to:

- Access memory (contiguous range, byte granular)



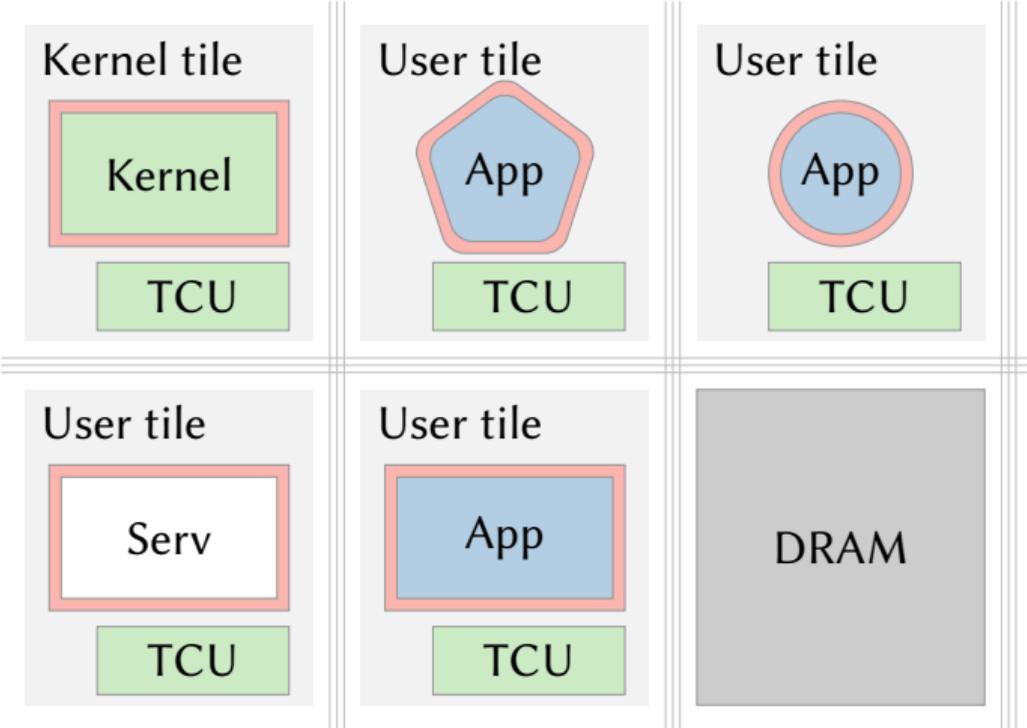
TCU provides *endpoints* to:

- Access memory (contiguous range, byte granular)
- Receive messages into a receive buffer
- Send messages to a receiving endpoint



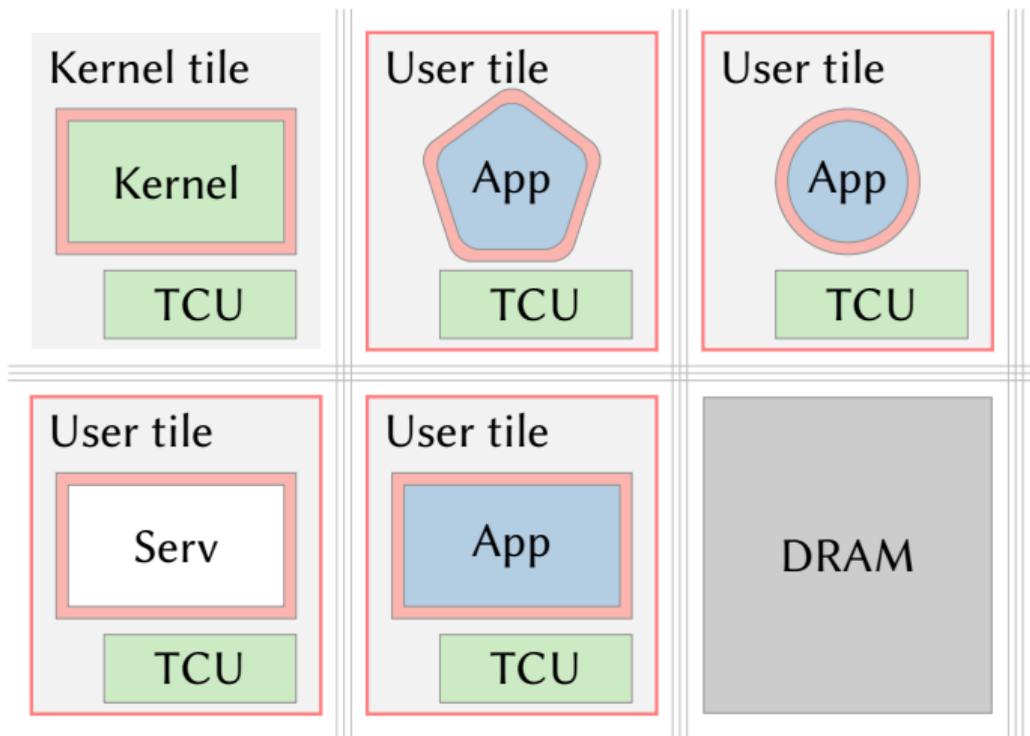
TCU provides *endpoints* to:

- Access memory (contiguous range, byte granular)
- Receive messages into a receive buffer
- Send messages to a receiving endpoint
- Replies for RPC



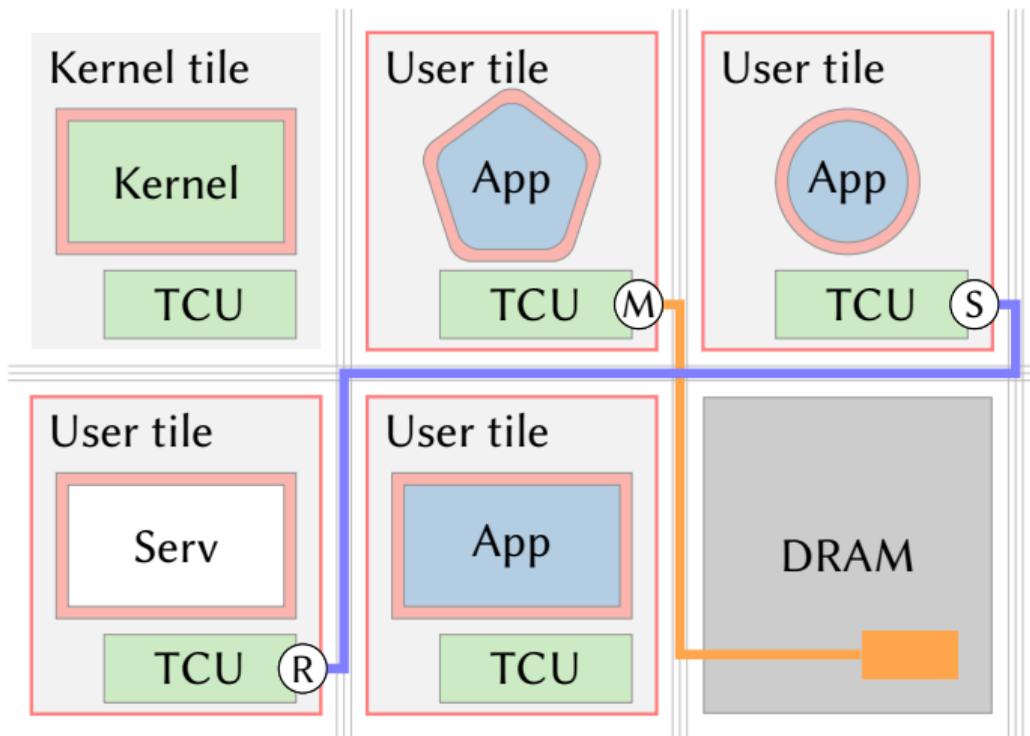
TCU-based isolation:

- Additional protection layer



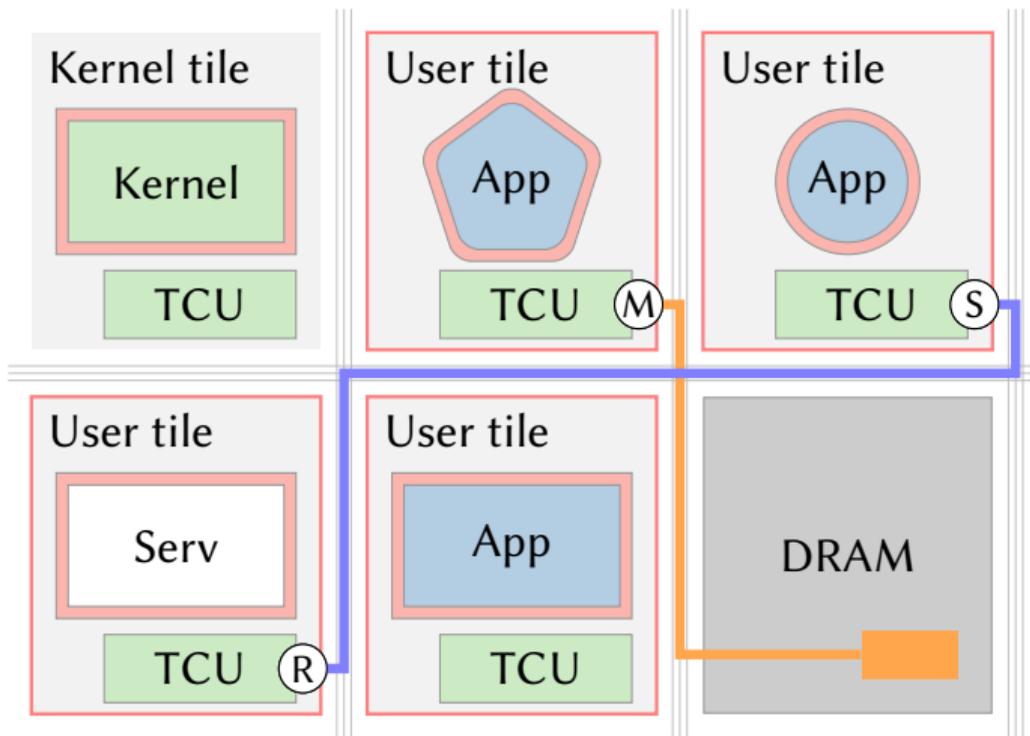
TCU-based isolation:

- Additional protection layer



TCU-based isolation:

- Additional protection layer
- Only kernel tile can establish communication channels



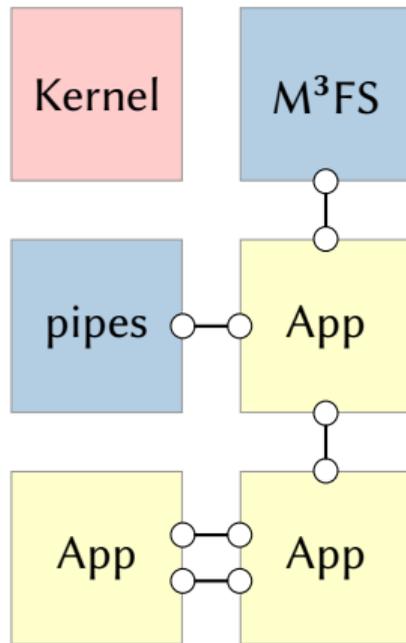
TCU-based isolation:

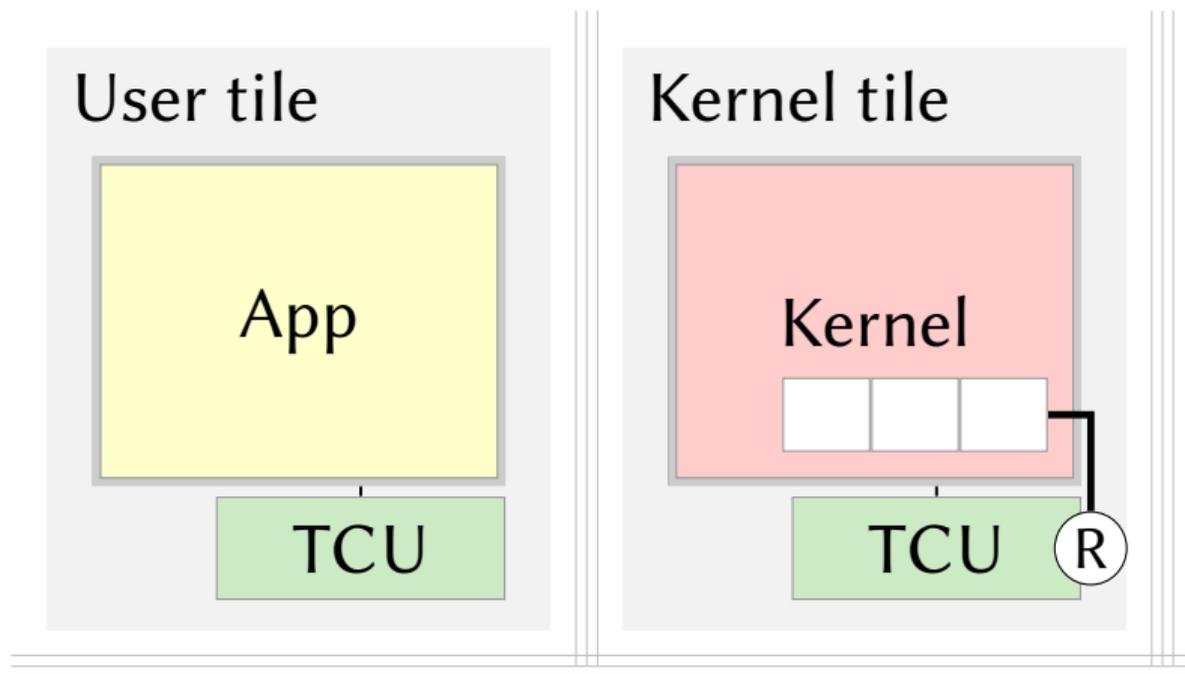
- Additional protection layer
- Only kernel tile can establish communication channels
- User tiles can only use established channels

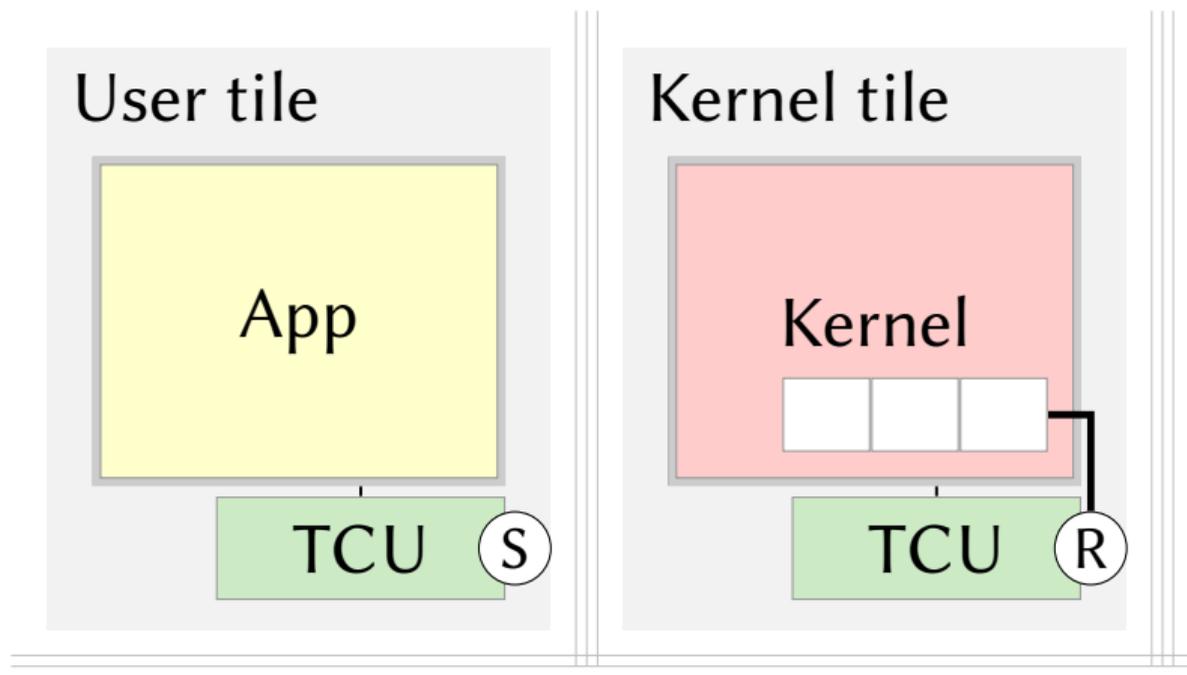


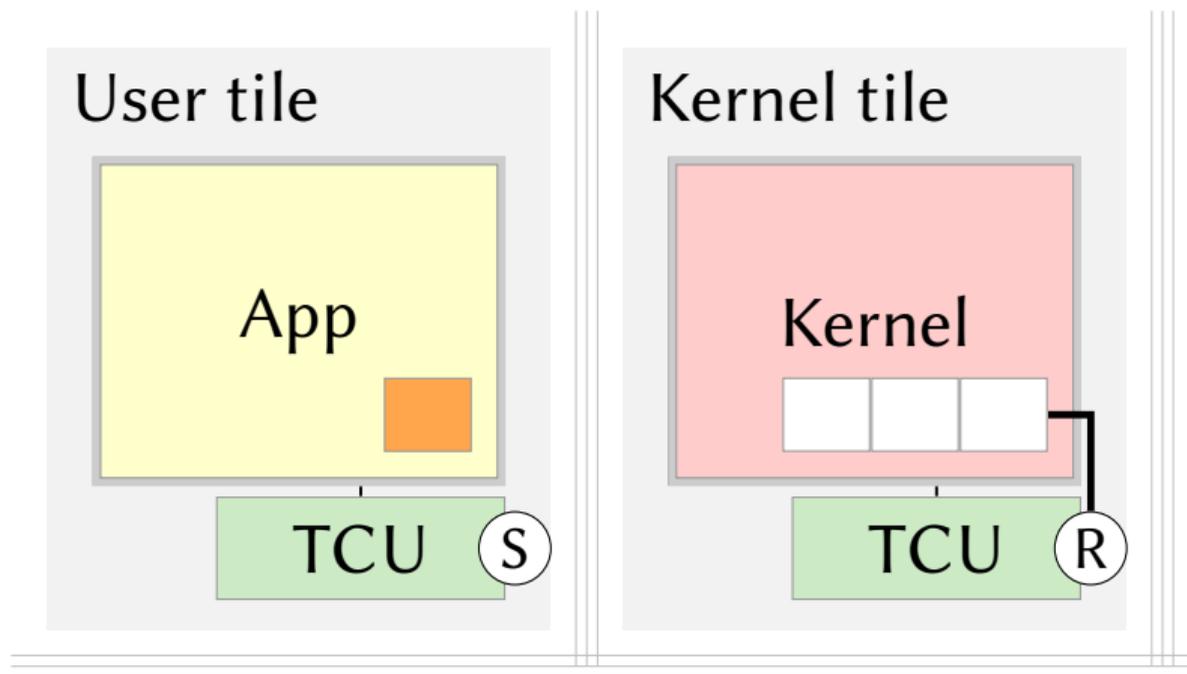
- 1 The New System Architecture
- 2 M<sup>3</sup>: The Operating System**
- 3 What are the Benefits?

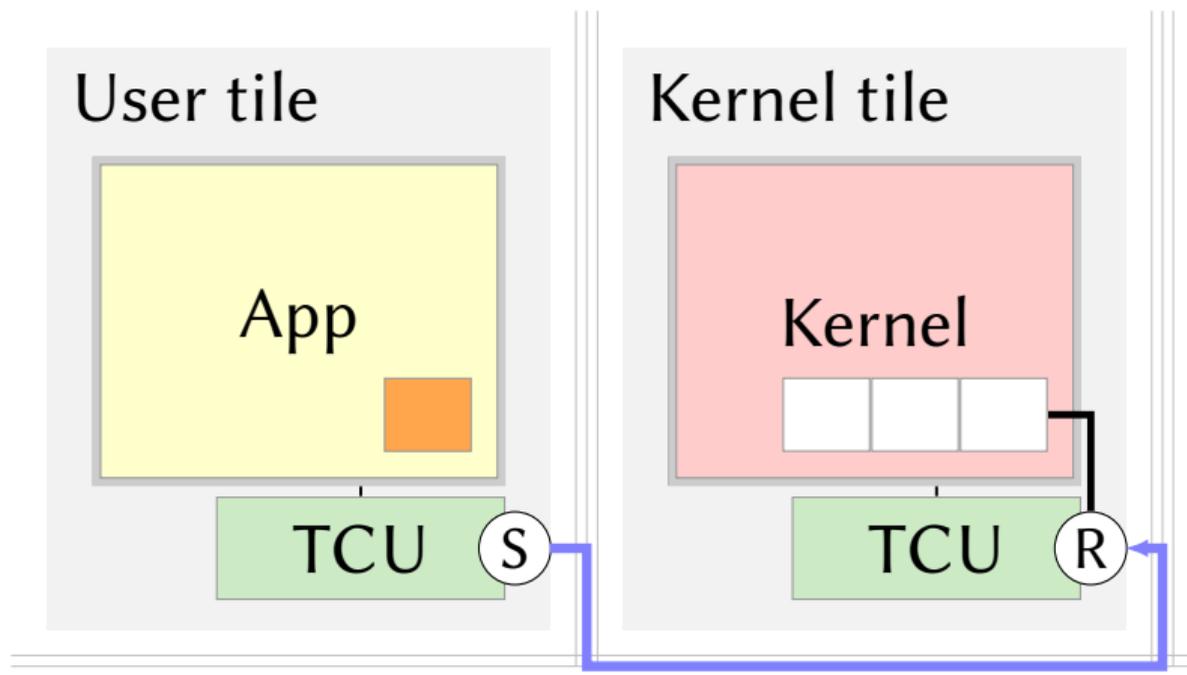
- **M<sup>3</sup>: Microkernel-based system** for het. **m**anycores (or  $L4 \pm 1$ )
- Implemented from scratch in Rust and C++
- Drivers, filesystems, etc. implemented on user tiles
- Kernel manages permissions, using capabilities
- TCU enforces permissions (communication, memory access)
- Kernel is independent of other tiles

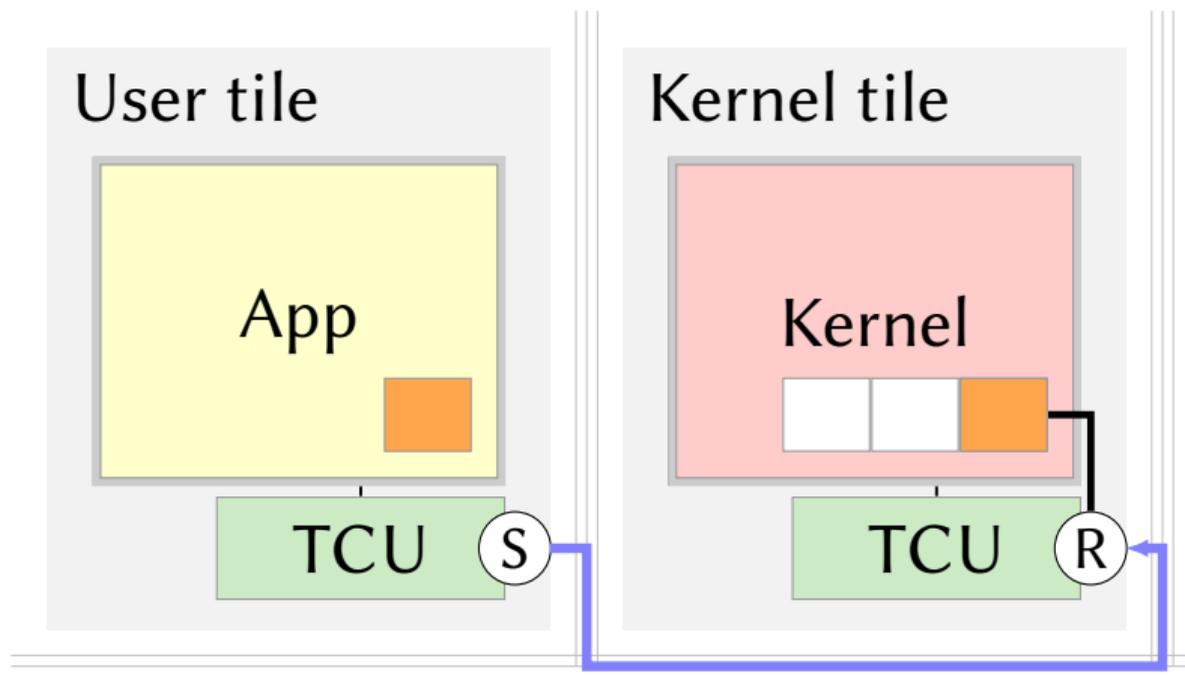




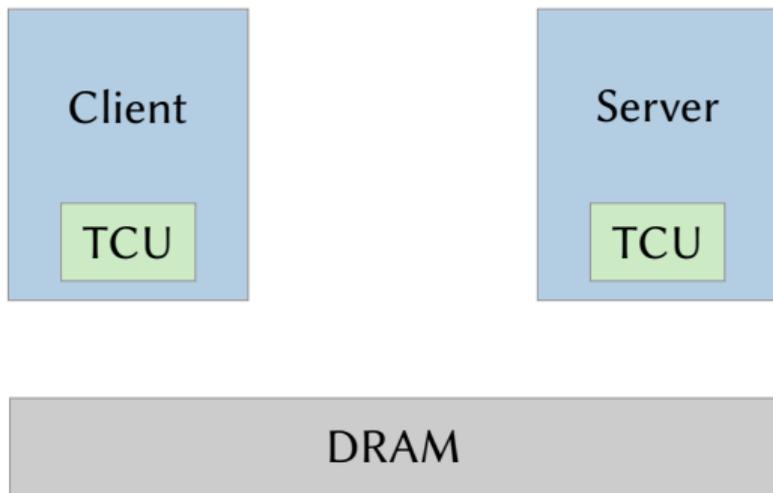


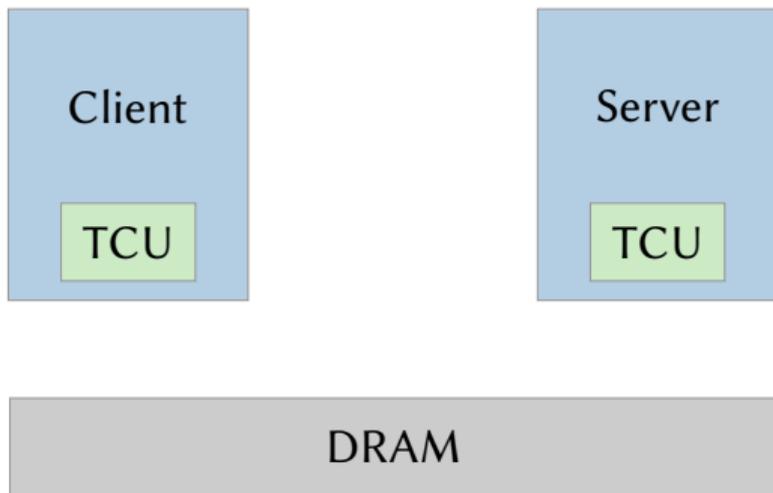






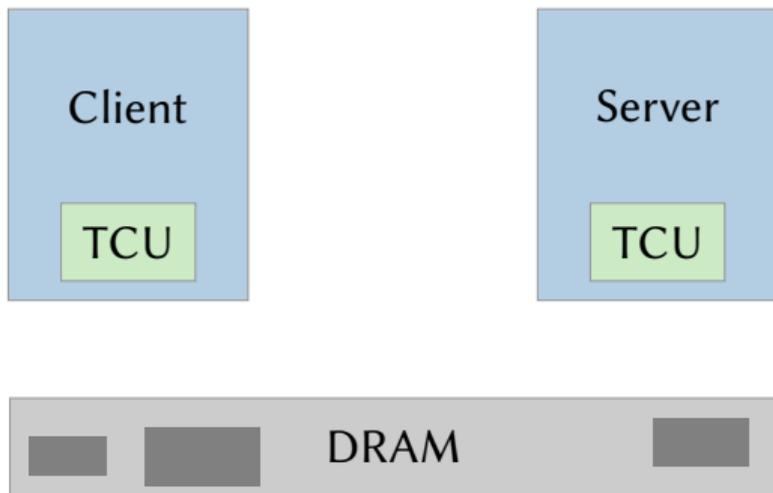
# OS Service Access





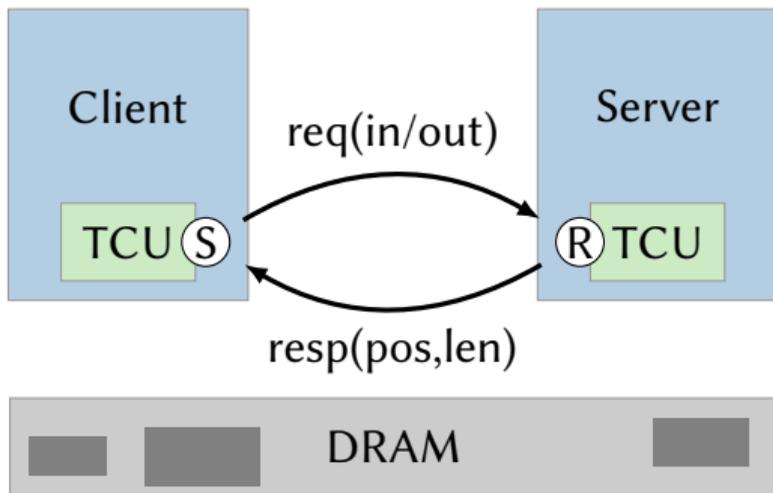
## File Protocol:

- Used for: files, pipes, ...



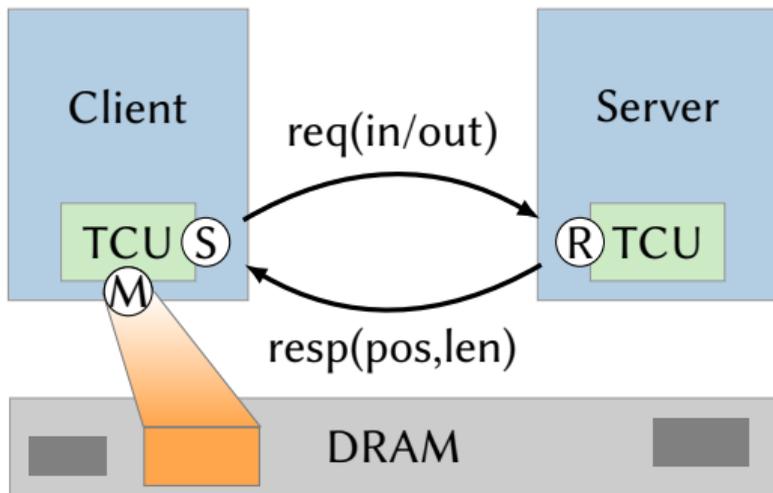
## File Protocol:

- Used for: files, pipes, ...
- Data in memory



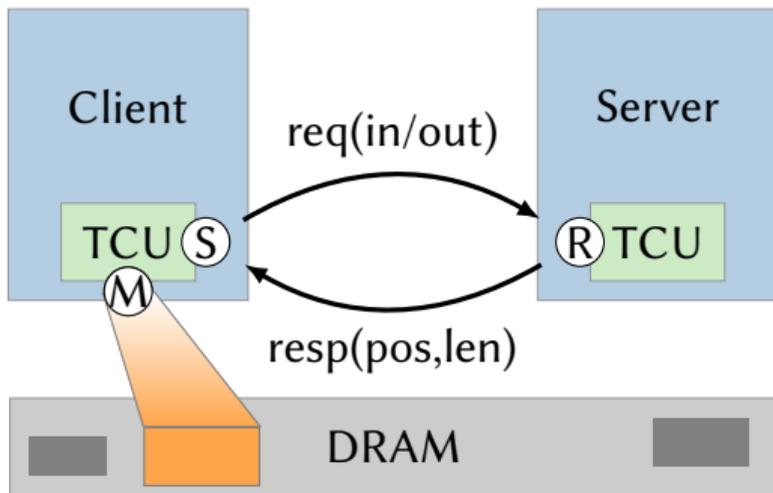
## File Protocol:

- Used for: files, pipes, ...
- Data in memory
- Msg channel between client and server
  - req(in) for next input piece
  - req(out) for next output piece



## File Protocol:

- Used for: files, pipes, ...
- Data in memory
- Msg channel between client and server
  - req(in) for next input piece
  - req(out) for next output piece
- Server configures client's memory EP



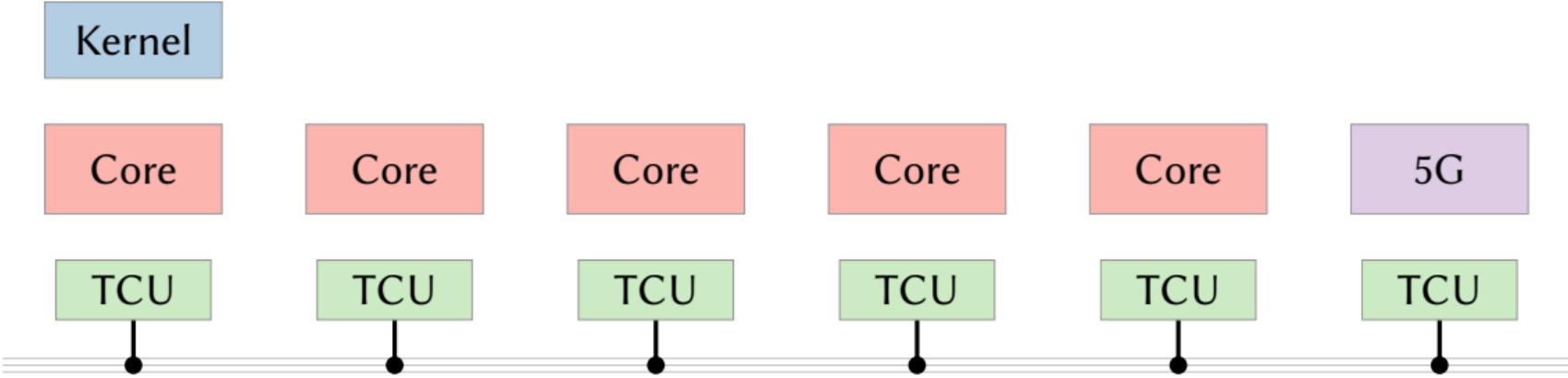
## File Protocol:

- Used for: files, pipes, ...
- Data in memory
- Msg channel between client and server
  - req(in) for next input piece
  - req(out) for next output piece
- Server configures client's memory EP
- Client accesses data via TCU

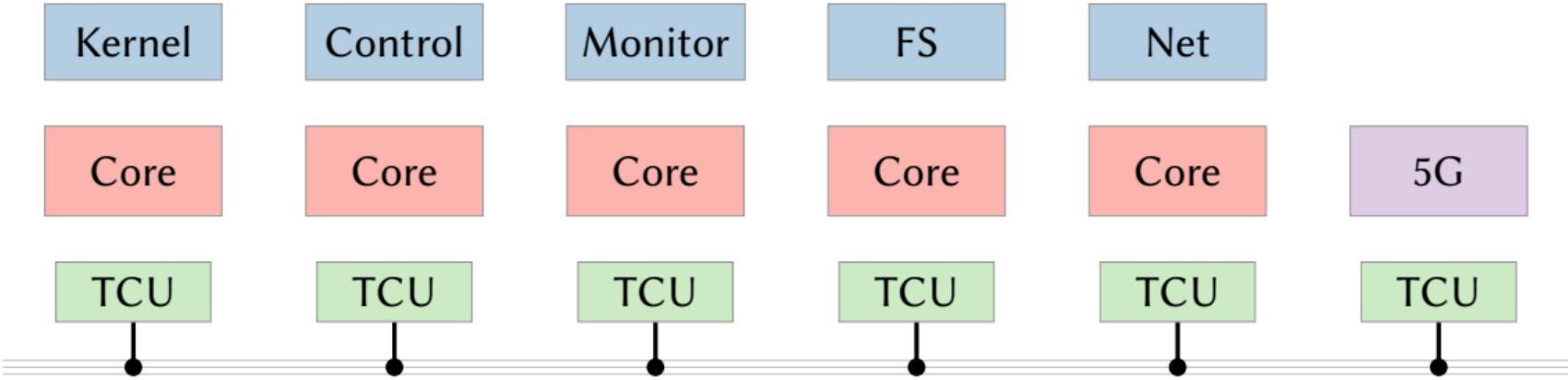


- 1 The New System Architecture
- 2 M<sup>3</sup>: The Operating System
- 3 What are the Benefits?**

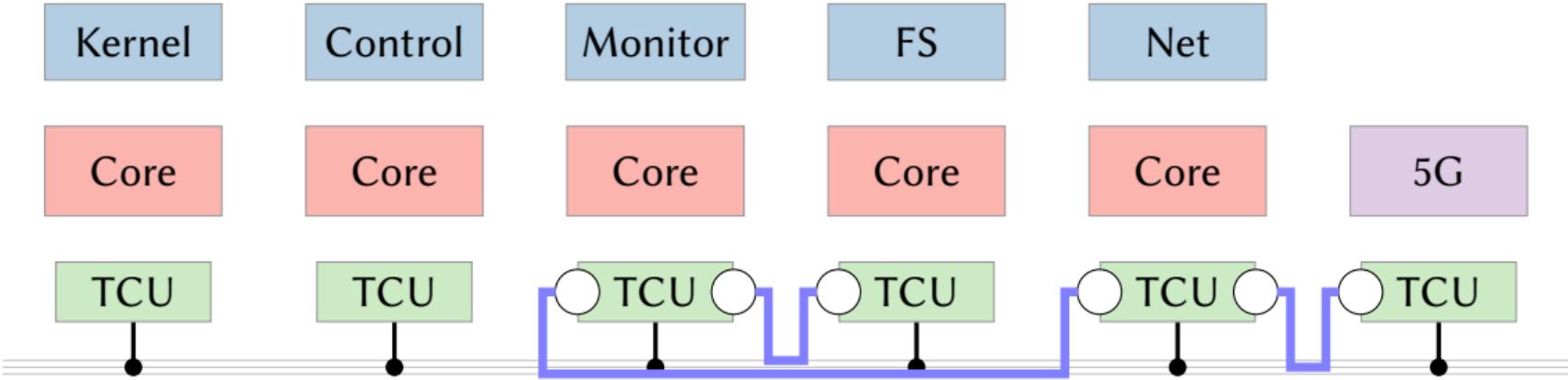
# Example System



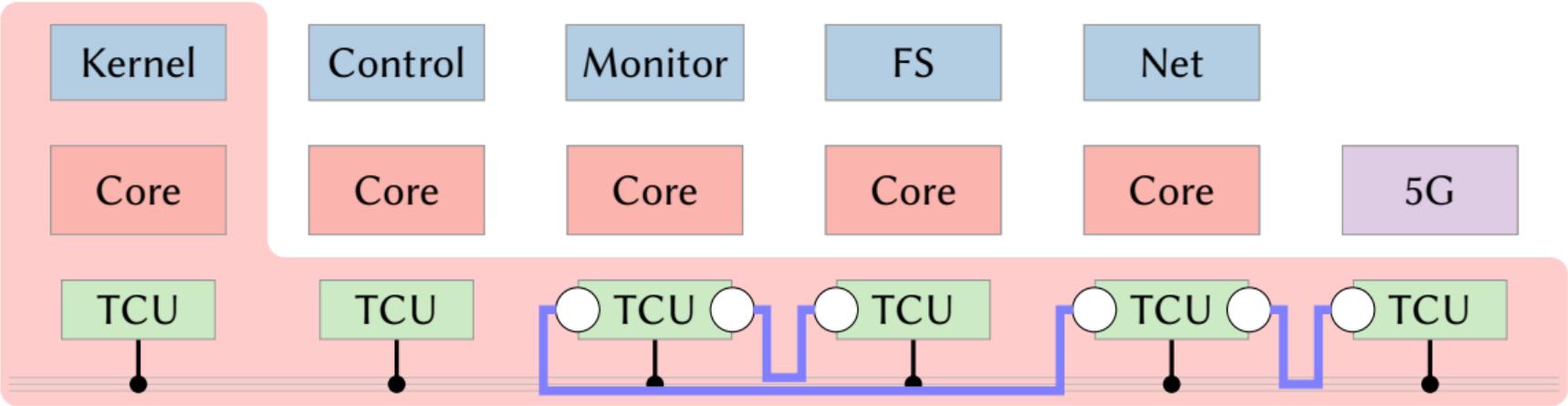
# Example System



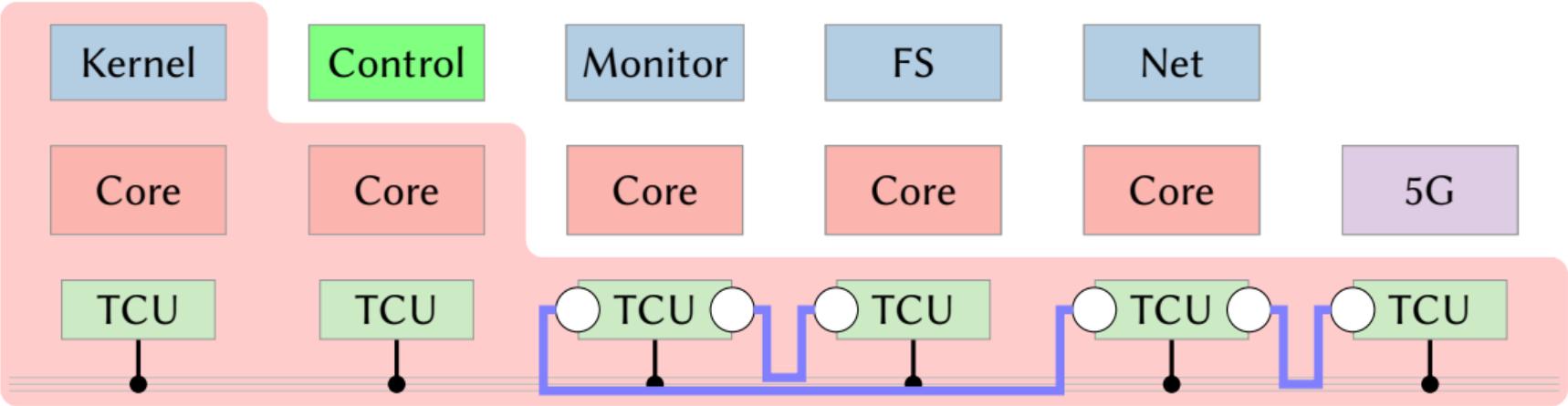
# Example System



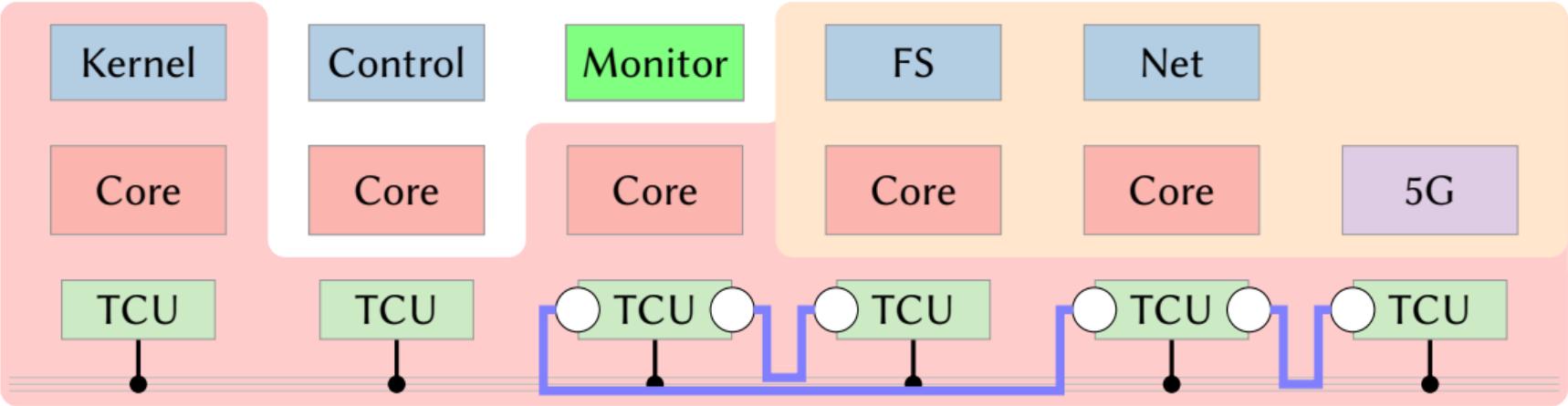
# Example System – TCB



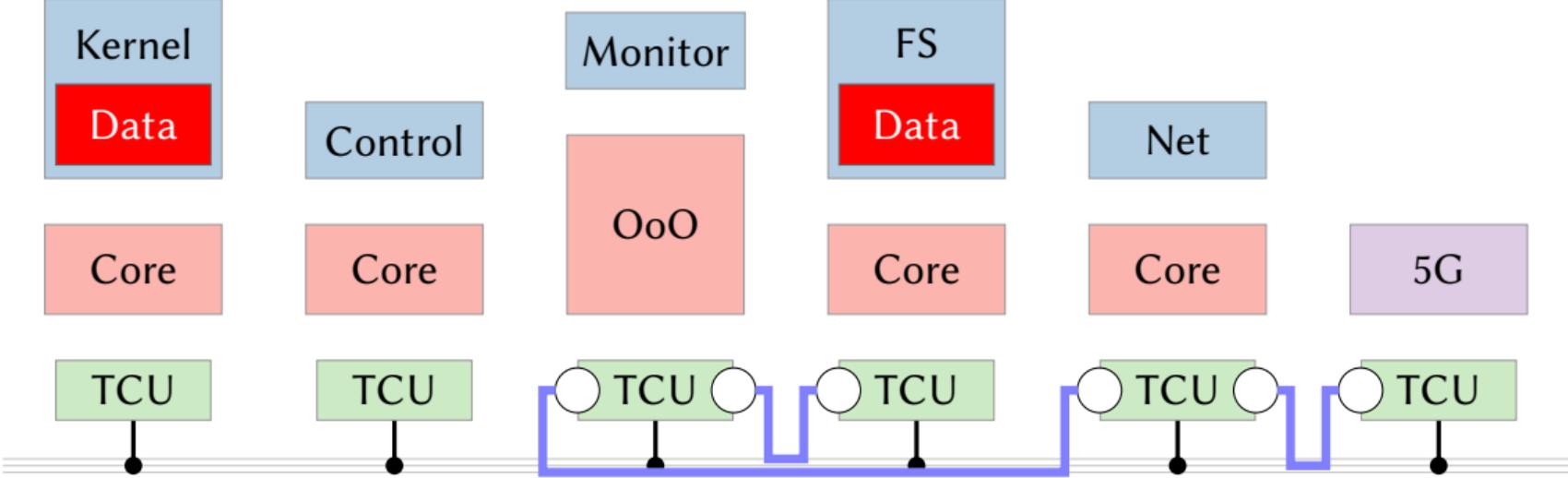
# Example System – TCB



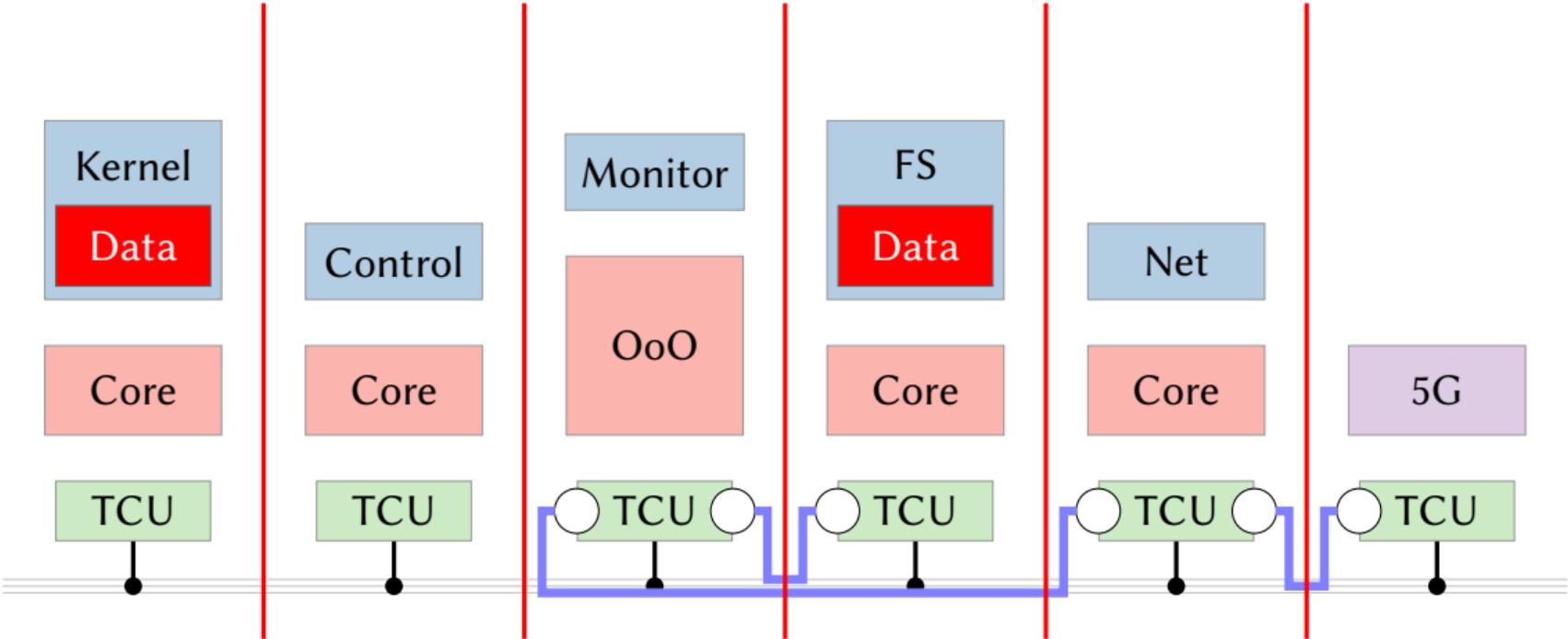
# Example System – TCB



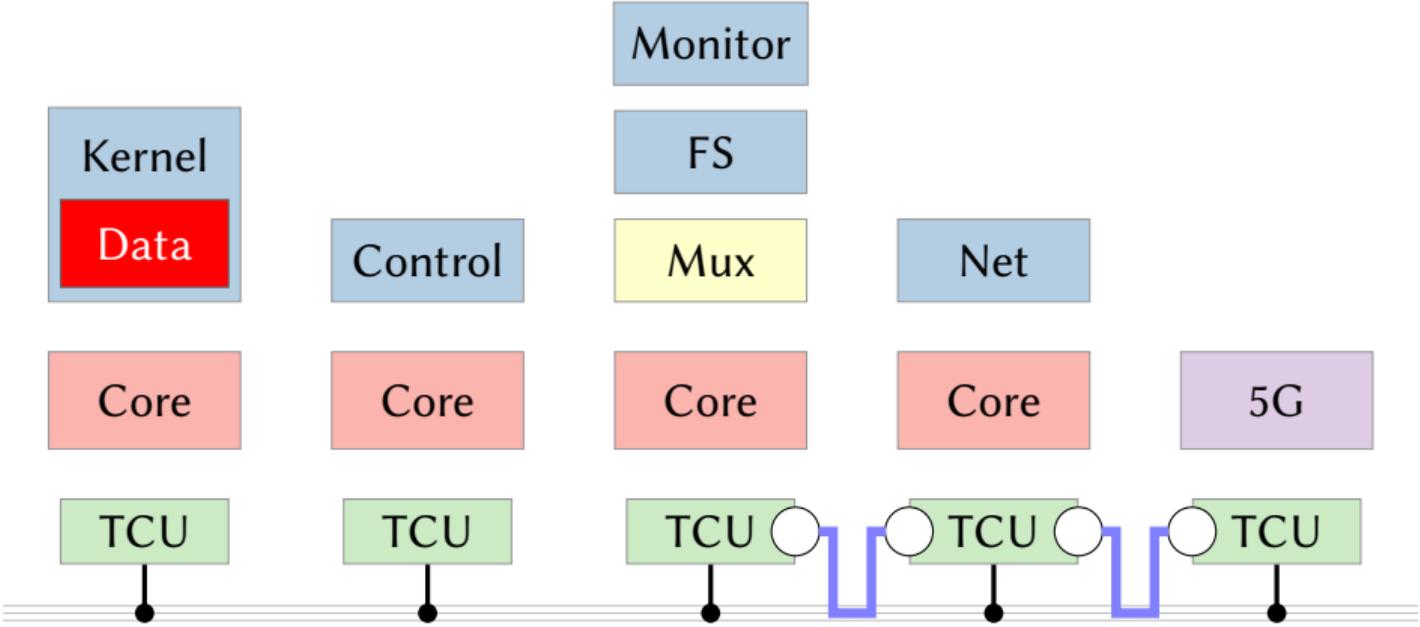
# Example System – Untrusted Core



# Example System – Untrusted Core



# Example System – Sharing (WIP)



```
LL1S =32 KiB (2-way assoc, 4 cycles)
LL1D =32 KiB (2-way assoc, 4 cycles)
L2S  =256 KiB (8-way assoc, 12 cycles)
Comp =Core -> DTU+AT -> L1S -> L2S

PE03: build/gen5-x86_64-release/bin/rctmux
Core =TimingSimpleCPU x86_64 @ 1GHz
DTU  =eps:16, bufsz:1024 B, blocksz:64 B, count:4, tlb:128, walker:1
LL1S =32 KiB (2-way assoc, 4 cycles)
LL1D =32 KiB (2-way assoc, 4 cycles)
L2S  =256 KiB (8-way assoc, 12 cycles)
Comp =Core -> DTU+AT -> L1S -> L2S

PE04: build/gen5-x86_64-release/default.img x 1
DTU  =eps:16, bufsz:1024 B, blocksz:1024 B, count:0, tlb:0, walker:0
Imem =3145728 KiB
Comp =DTU -> DRAM

Global frequency set at 1000000000000 ticks per second
Info: kernel located at: build/gen5-x86_64-release/bin/kernel
Info: kernel located at: build/gen5-x86_64-release/bin/rctmux
Info: kernel located at: build/gen5-x86_64-release/bin/rctmux
Info: kernel located at: build/gen5-x86_64-release/bin/rctmux
warn: DRAM device capacity (49152 Mbytes) does not match the address range assigned (4096 Mbytes)
Info: No kernel set for full system simulation. Assuming you know what you're doing
Info: No kernel set for full system simulation. Assuming you know what you're doing
platform.com_1.device: Listening for connections on port 3456
@: pe00.remote_gdb: listening for remote gdb on port 7000
@: pe01.remote_gdb: listening for remote gdb on port 7001
@: pe02.remote_gdb: listening for remote gdb on port 7002
@: pe03.remote_gdb: listening for remote gdb on port 7003
warn: CoherentXBAR pe04.xbar has no snooping ports attached!
Info: Loaded 'root' to 0x8400000000000000 .. 0x8400000000043800
Info: Loaded 'hello' to 0x8400000000044000 .. 0x840000000016eb00
Info: Loaded 'hello' to 0x840000000016f000 .. 0x8400000000299b00
Info: Loaded 'rctmux' to 0x840000000029a000 .. 0x84000000002aed00
Info: Entering event queue @ 0. Starting simulation...
[kernel_00] Kernel is ready
Hello World
Hello World
[kernel_00] Shutting down
Exiting @ tick 6355915000 because m5 exit instruction encountered
```

gem5 simulator

# Prototype Platforms



```
LL1S =32 KiB (2-way assoc, 4 cycles)
LL1D =32 KiB (2-way assoc, 4 cycles)
L2S =256 KiB (8-way assoc, 12 cycles)
Comp =Core -> DTU+AT -> LL1S -> L2S

PE03: build/gen5-x86_64-release/bin/rctmux
Core =TimingSimpleCPU x86_64 @ 1GHz
DTU =eps:16, bufsz:1024 B, blocksz:64 B, count:4, tlb:128, walker:1
LL1S =32 KiB (2-way assoc, 4 cycles)
LL1D =32 KiB (2-way assoc, 4 cycles)
L2S =256 KiB (8-way assoc, 12 cycles)
Comp =Core -> DTU+AT -> LL1S -> L2S

PE04: build/gen5-x86_64-release/default.img x 1
DTU =eps:16, bufsz:1024 B, blocksz:1024 B, count:0, tlb:0, walker:0
Imem =3145728 KiB
Comp =DTU -> DRAM

Global frequency set at 100000000000 ticks per second
Info: kernel located at: build/gen5-x86_64-release/bin/kernel
Info: kernel located at: build/gen5-x86_64-release/bin/rctmux
Info: kernel located at: build/gen5-x86_64-release/bin/rctmux
Info: kernel located at: build/gen5-x86_64-release/bin/rctmux
warn: DRAM device capacity (49152 Mbytes) does not match the address range assigned (4096 Mbytes)
Info: No kernel set for full system simulation. Assuming you know what you're doing
Info: No kernel set for full system simulation. Assuming you know what you're doing
platform.com_1.device: Listening for connections on port 3456
0: pe00.remote_gdb: listening for remote gdb on port 7000
0: pe01.remote_gdb: listening for remote gdb on port 7001
0: pe02.remote_gdb: listening for remote gdb on port 7002
0: pe03.remote_gdb: listening for remote gdb on port 7003
warn: CoherentXbar pe04_xbar has no snooping ports attached!
Info: Loaded 'root' to 0x8400000000000000 .. 0x8400000000043868
Info: Loaded 'hello' to 0x8400000000044000 .. 0x84000000000816b00
Info: Loaded 'hello' to 0x840000000016f000 .. 0x8400000000299b00
Info: Loaded 'rctmux' to 0x840000000029a000 .. 0x84000000002aed00
Info: Entering event queue @ 0. Starting simulation...
[kernel @0] Kernel is ready
Hello World
Hello World
[kernel @0] Shutting down
Exiting @ tick 6355915000 because m5 exit instruction encountered
```

gem5 simulator



FPGA

# Demo

- Microkernels are great!
- Their ideas can also be applied to hardware:
  - Trusted communication unit per tile
  - Isolated software *and* hardware components on top
- Has several additional benefits:
  - Allows to securely integrate untrusted third-party components
  - Prevents (known) side-channel attacks by physical isolation
  - Simplifies heterogeneous systems by uniform interface
- $M^3$  is available at <https://github.com/TUD-OS/M3>,  
gem5 extensions at <https://github.com/TUD-OS/gem5-dtu>