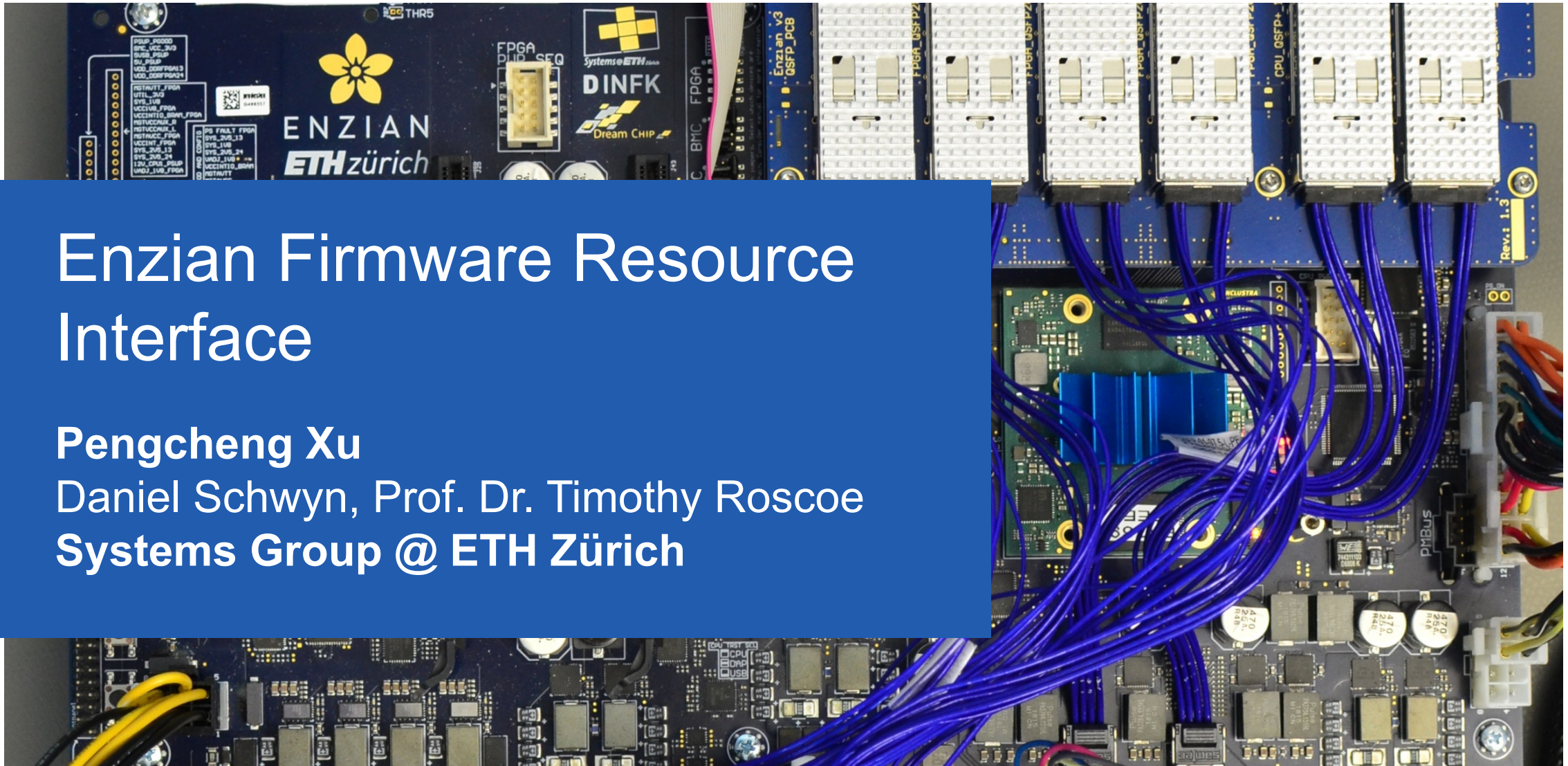


Enzian Firmware Resource Interface

Pengcheng Xu

Daniel Schwyn, Prof. Dr. Timothy Roscoe

Systems Group @ ETH Zürich



- The speaker would like to thank the Enzian team, including (but not limited to):

David Cock, Timothy Roscoe, Daniel Schwyn, Michael Giardino, Ben Fiedler, and many more

- For their great help:
 - Discussions and feedback on the core design
 - Practical help with implementation, development and debugging
 - Constructive comments on the project report
 - Access to previous slide decks for building this presentation
 - Feedback on the presentation dry-run
 - ...



Introduction & Motivation



Design



Implementation



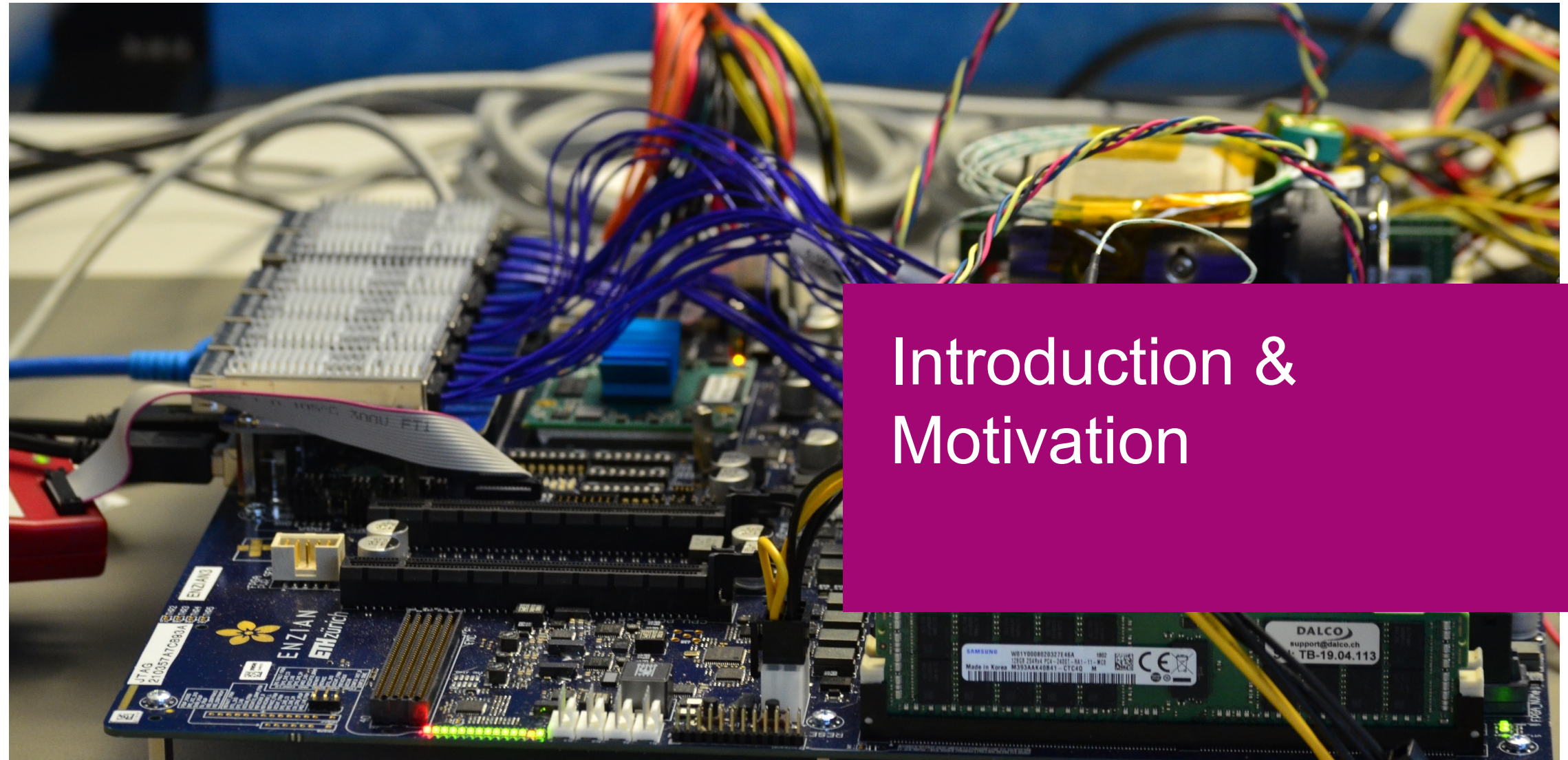
Evaluation



Live Demo



Discussion

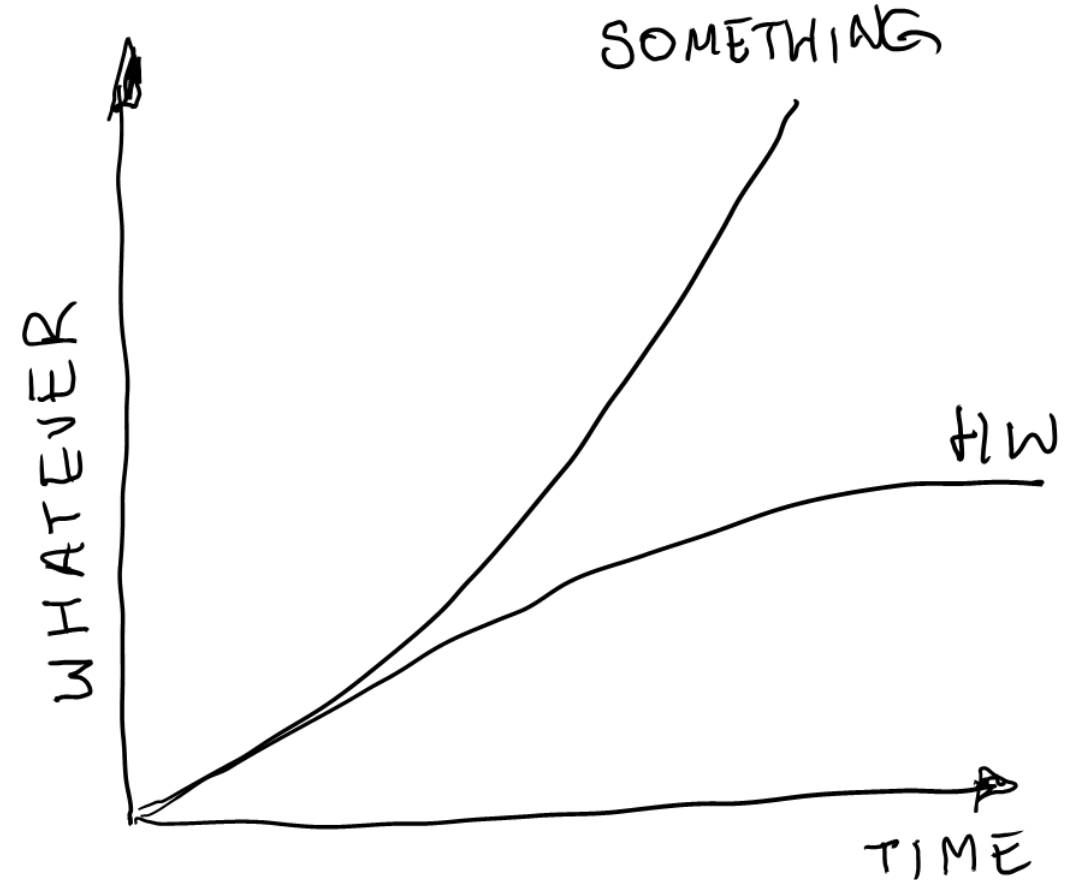


Introduction & Motivation

The usual bla, bla, bla ...

- Moore's law
- Dennar scaling, physical limits
- Multicore
- GPU, TPU, FPGA
- Data centers and the cloud
- ...

- Corollary: Hardware is changing really fast (because they do not know what to do with the transistors)



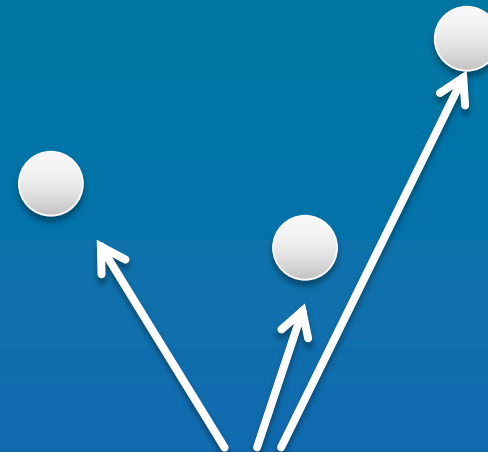
(courtesy Gustavo Alonso)

Feasible hardware design space

Scope of most systems software research

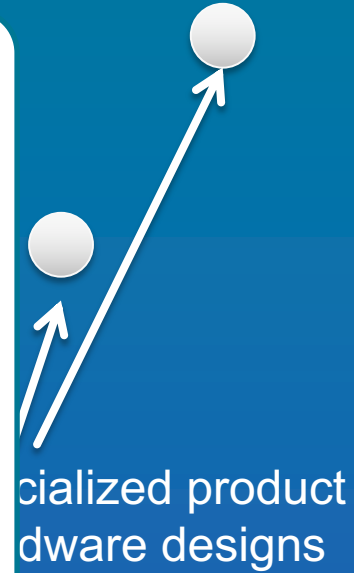
Available COTS hardware

Specialized product hardware designs

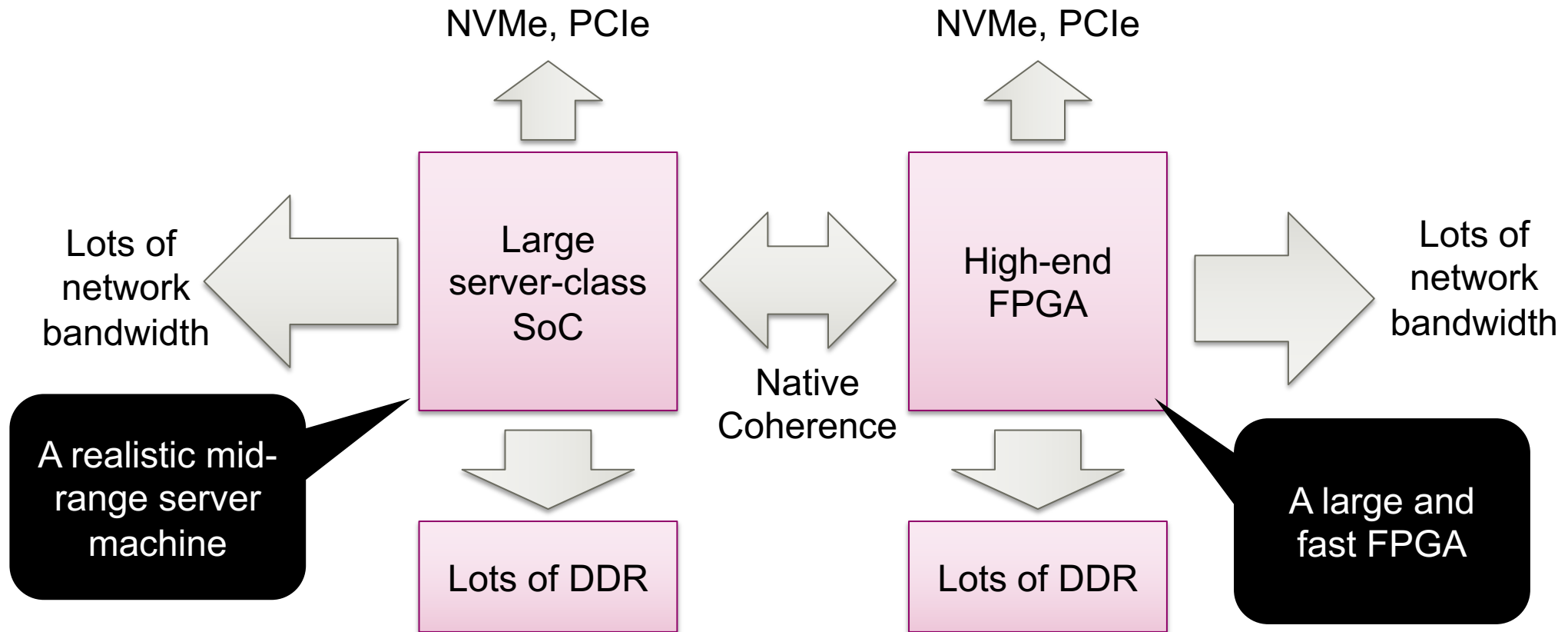


Feasible hardware design space

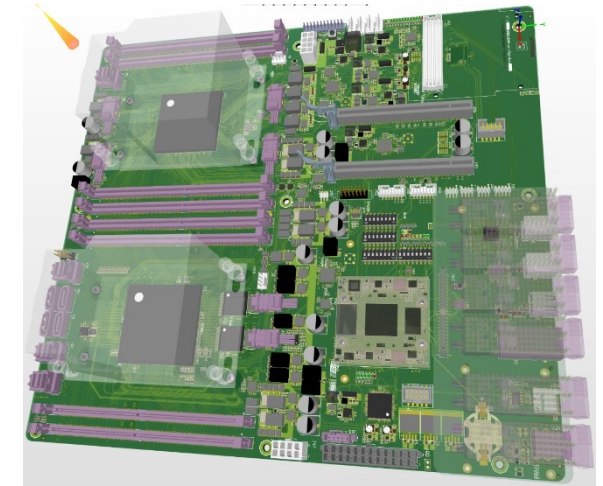
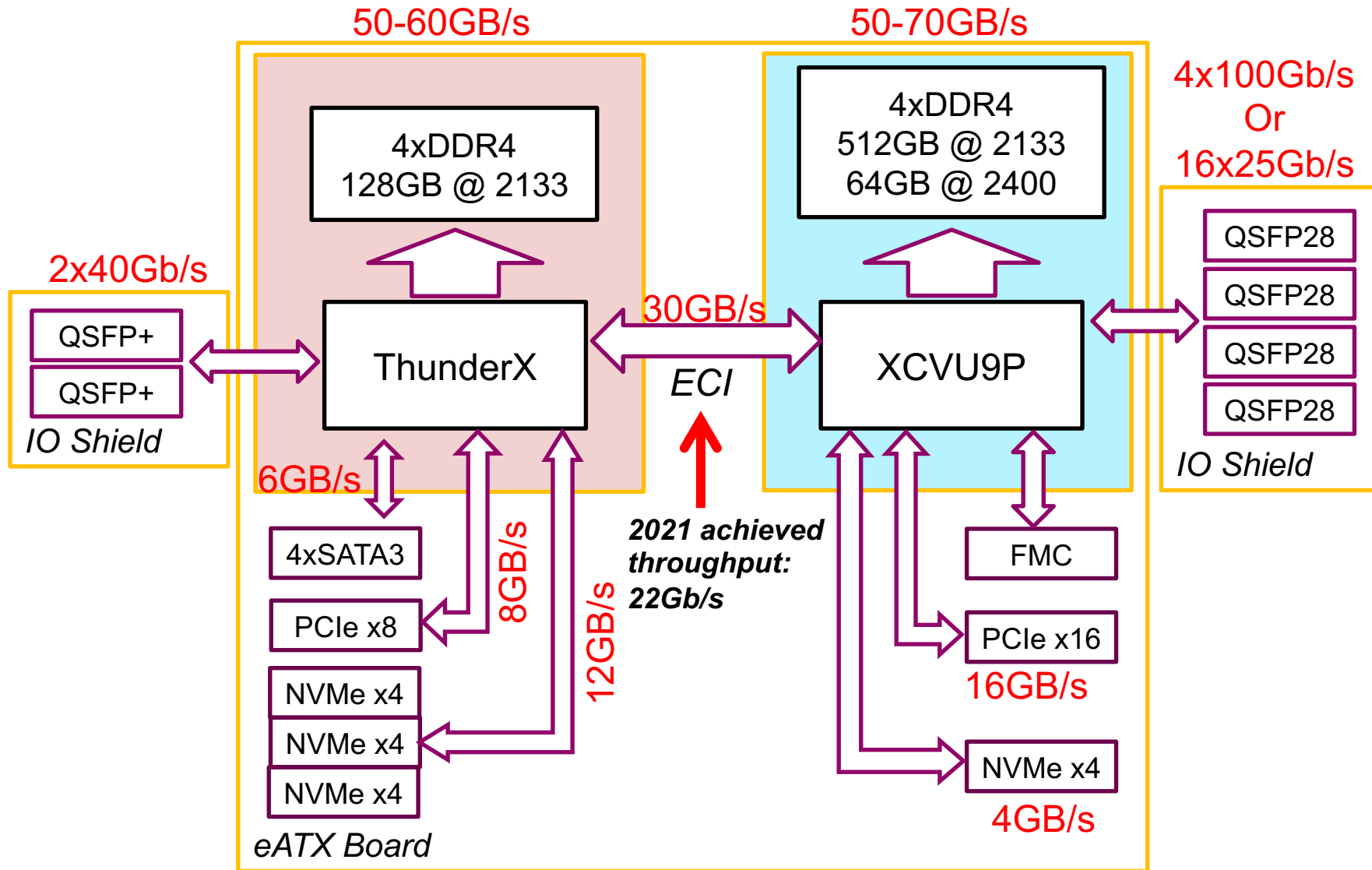
S
SC



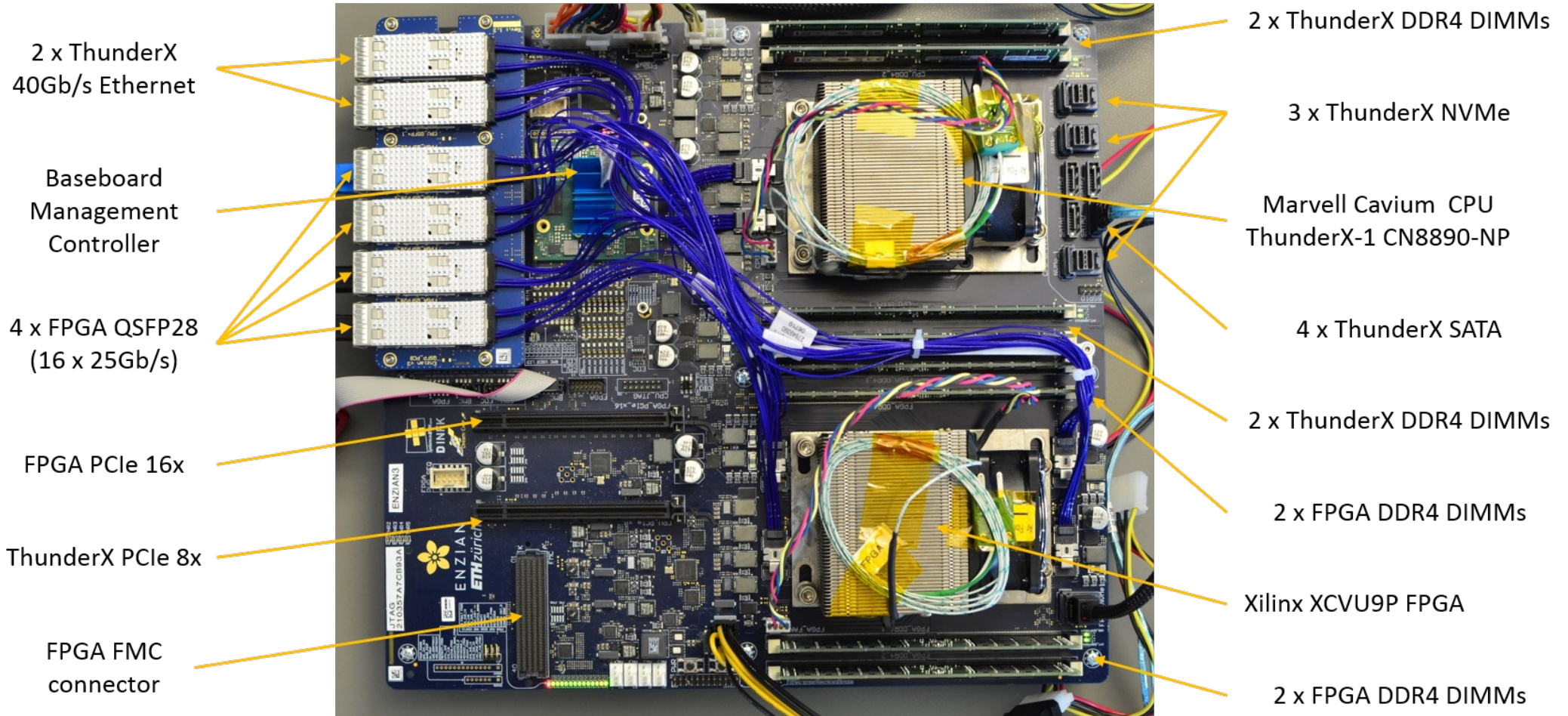
The Vision



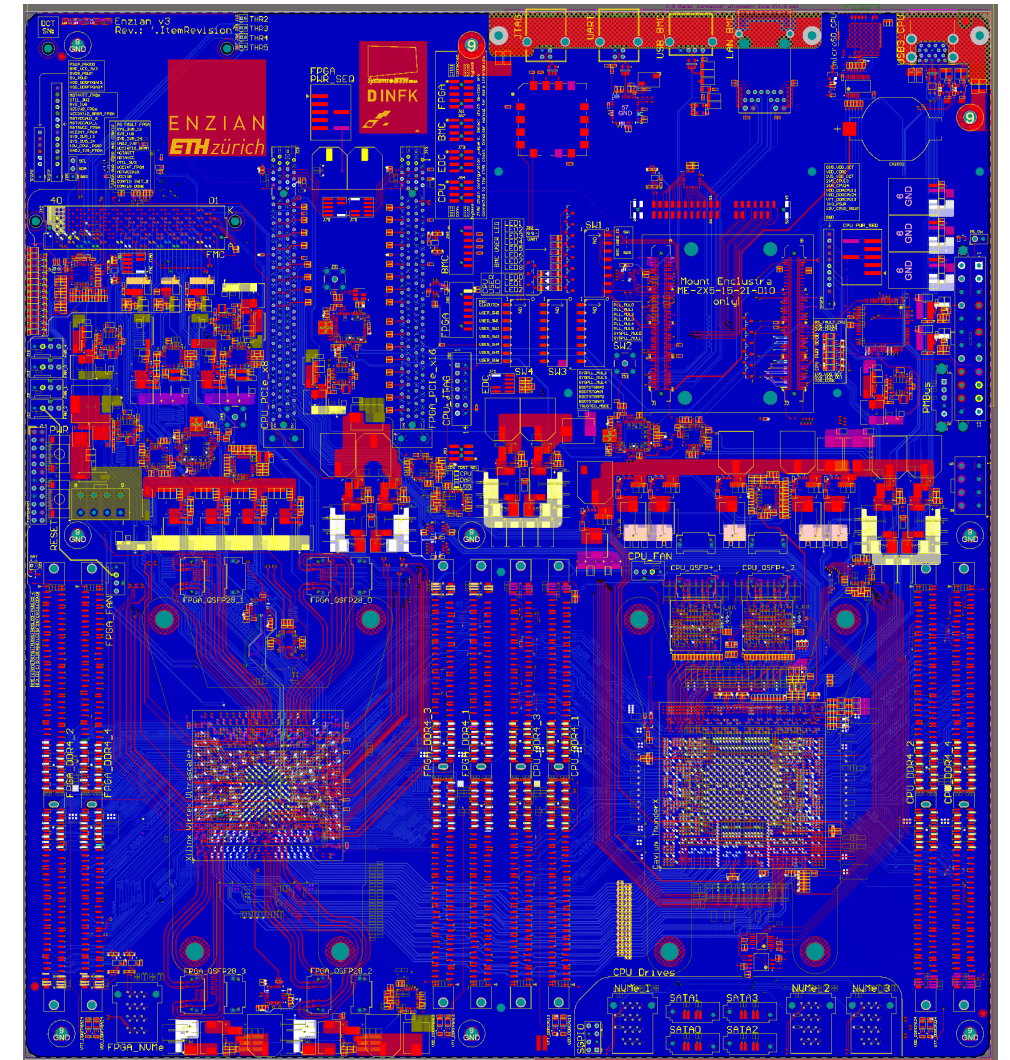
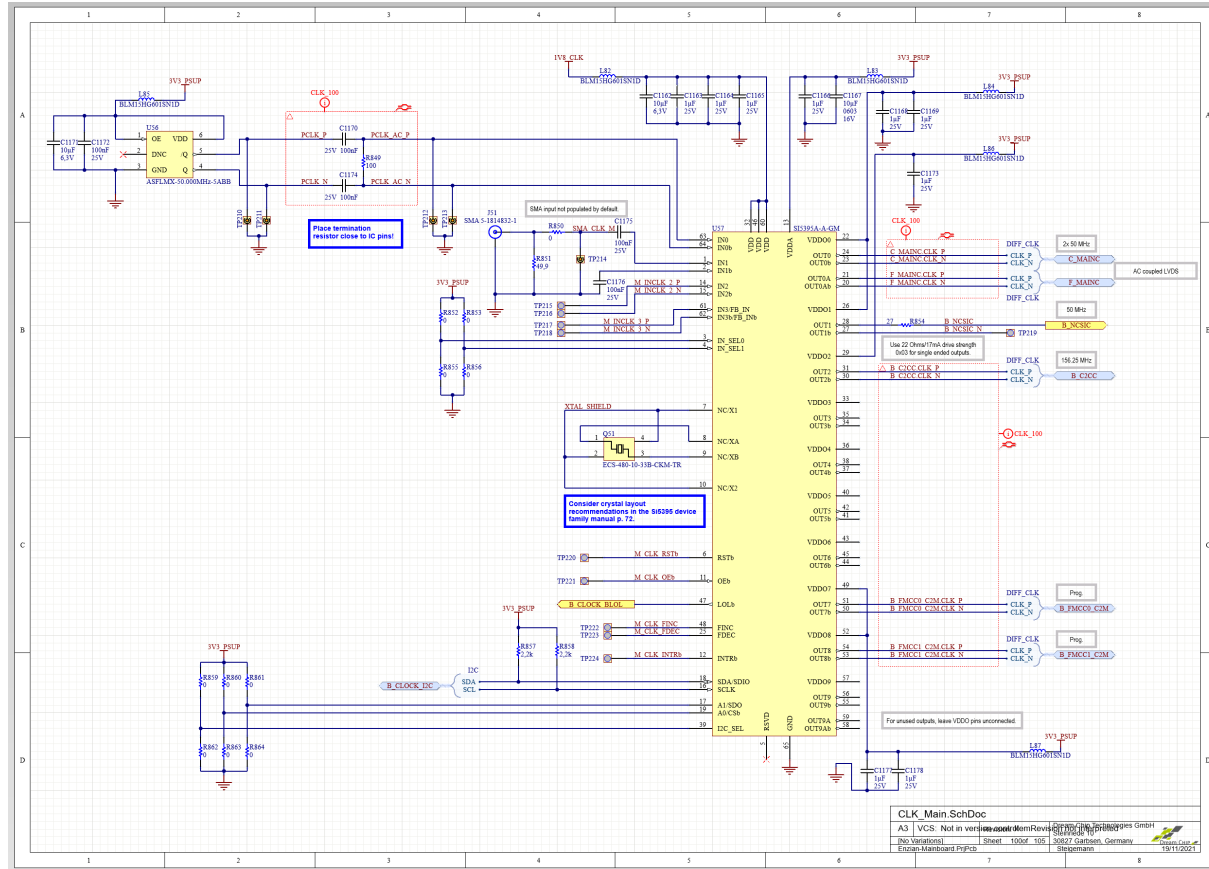
The Reality



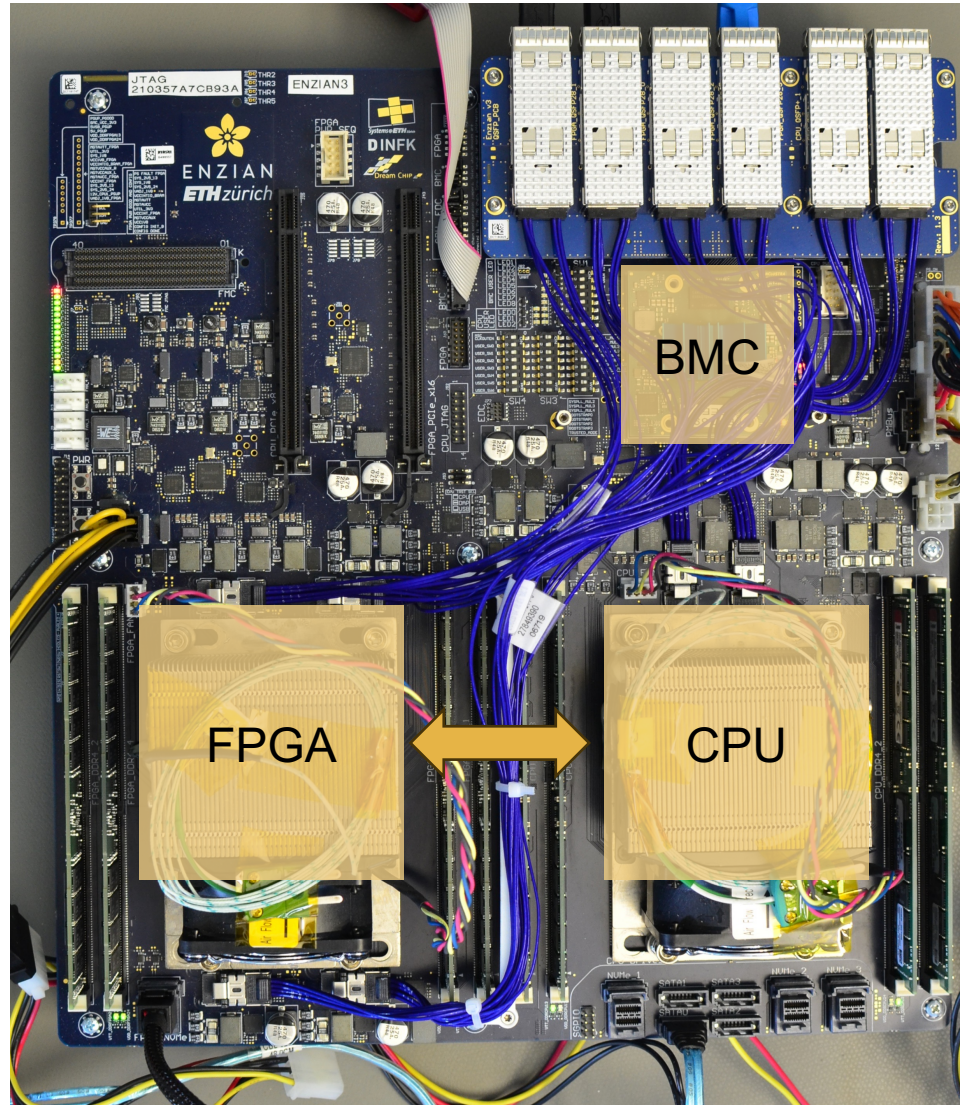
So far we have 9 working machines (5 more in testing)



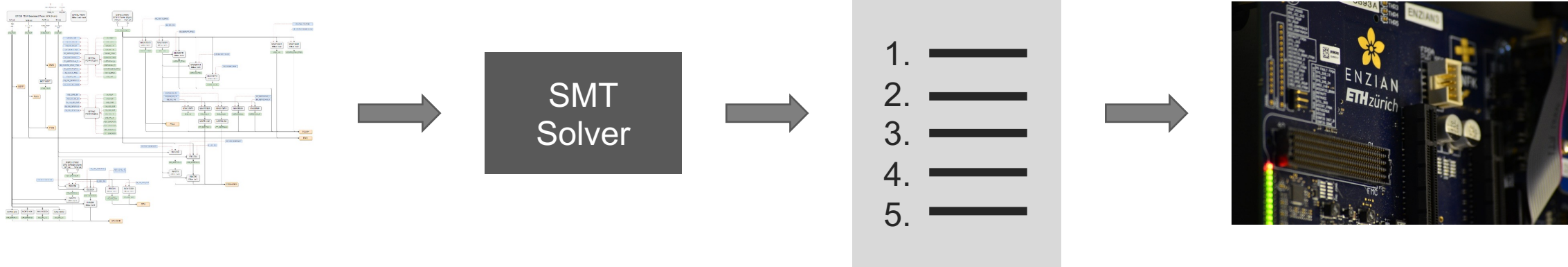
Fully Open-Source HW Design



The Computer in the Computer



... and we need a more systematic approach!

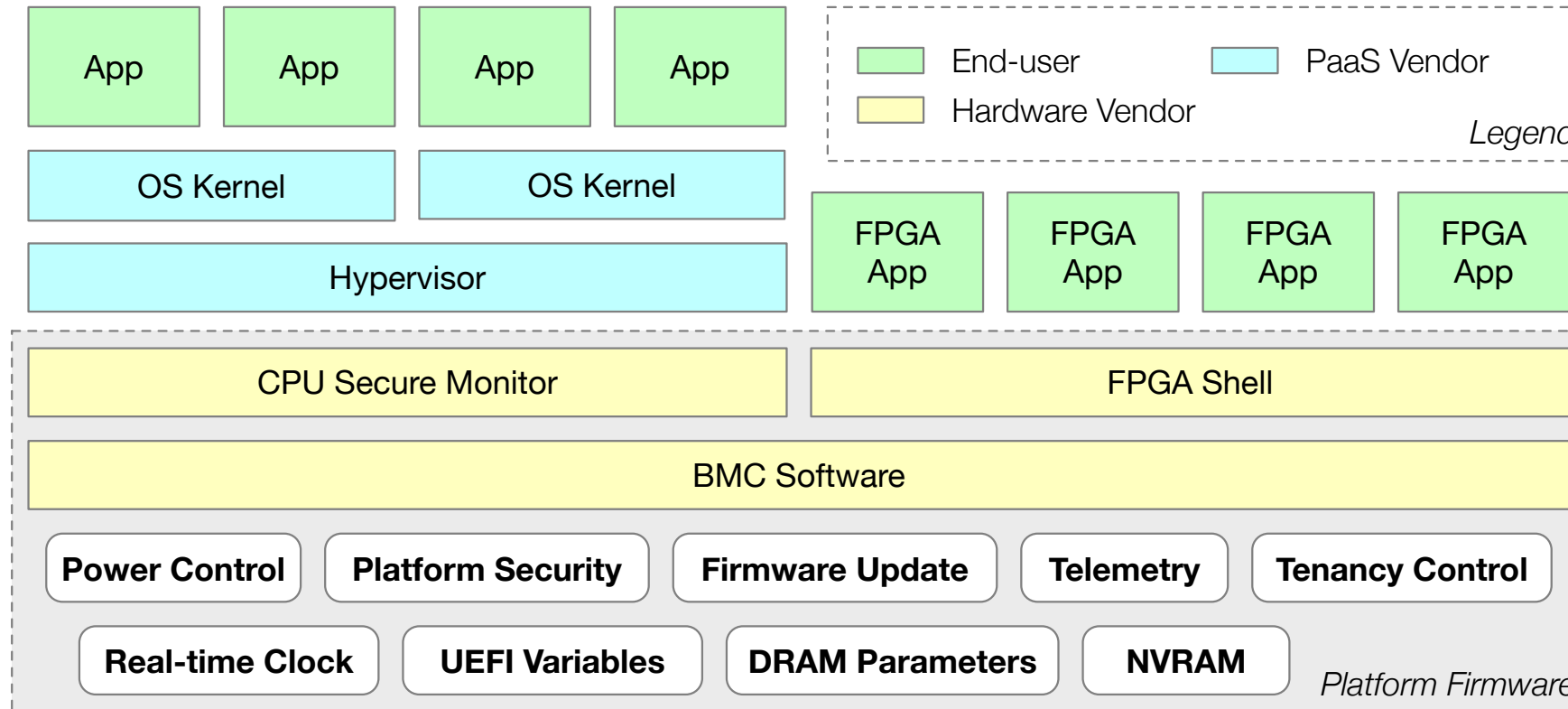


- We can derive sequences efficiently (< 10s)
- Works on real hardware
- More confidence in power-up sequences
- Basis for rigorous specification of correct behavior

Jasmin Schult, Daniel Schwyn, Michael Giardino, David Cock, Reto Achermann, and Timothy Roscoe. 2021. Declarative Power Sequencing. ACM Trans. Embed. Comput. Syst. 20, 5s, Article 84

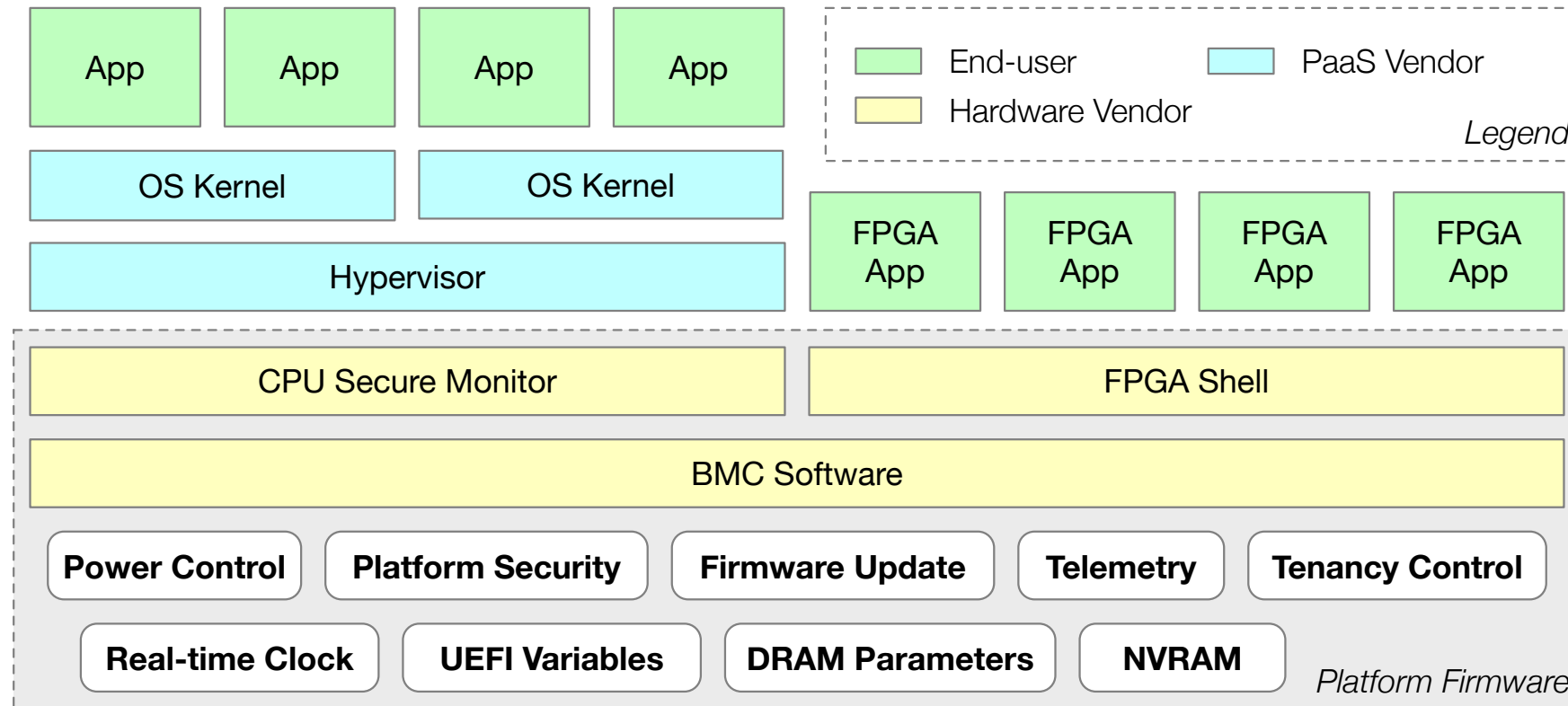
What is firmware?

- **Firmware:** a class of software that provides *low-level control* for a device's hardware
 - Provide services in the form of **actions** on **resources**



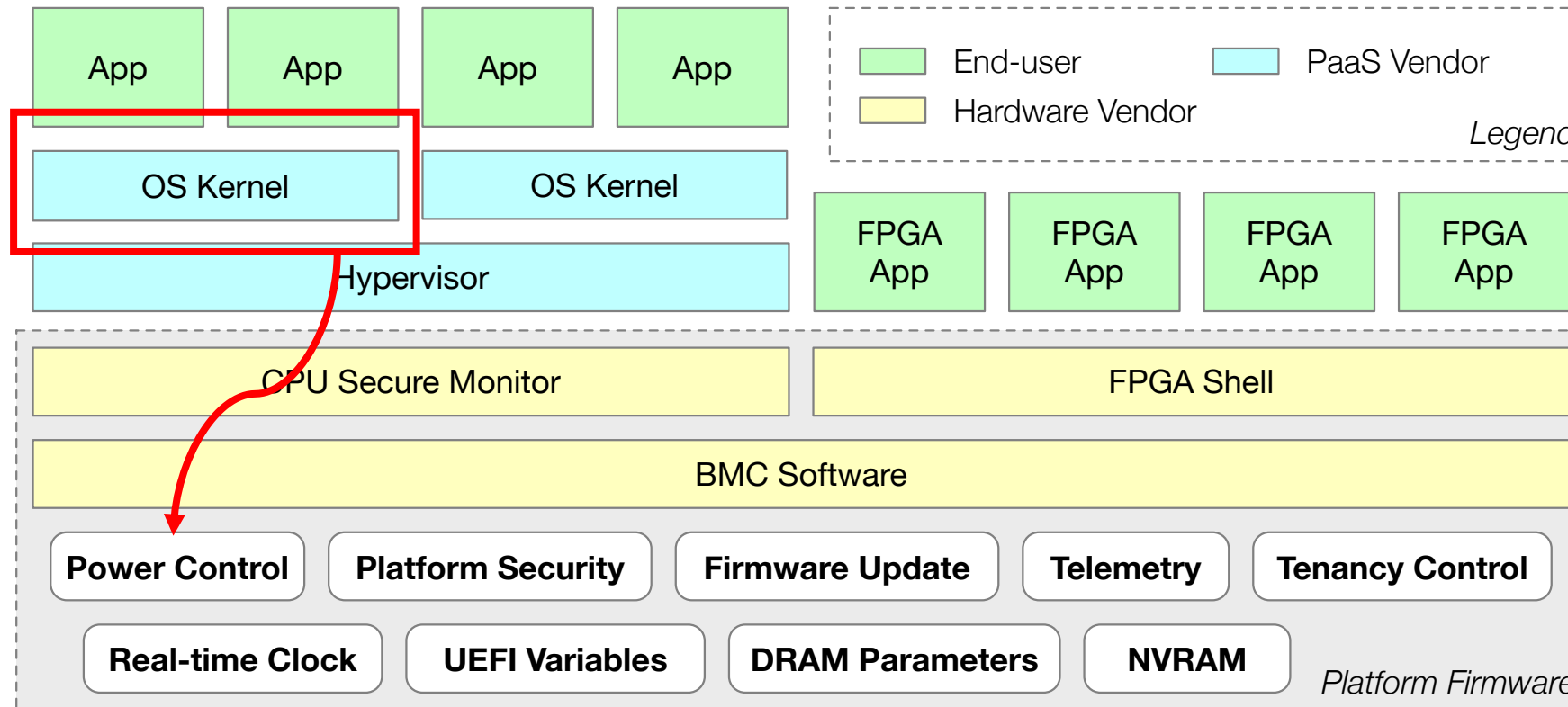
What is a firmware protocol?

- **Firmware protocol:** semantics of how to access firmware resources
 - From different places in the system



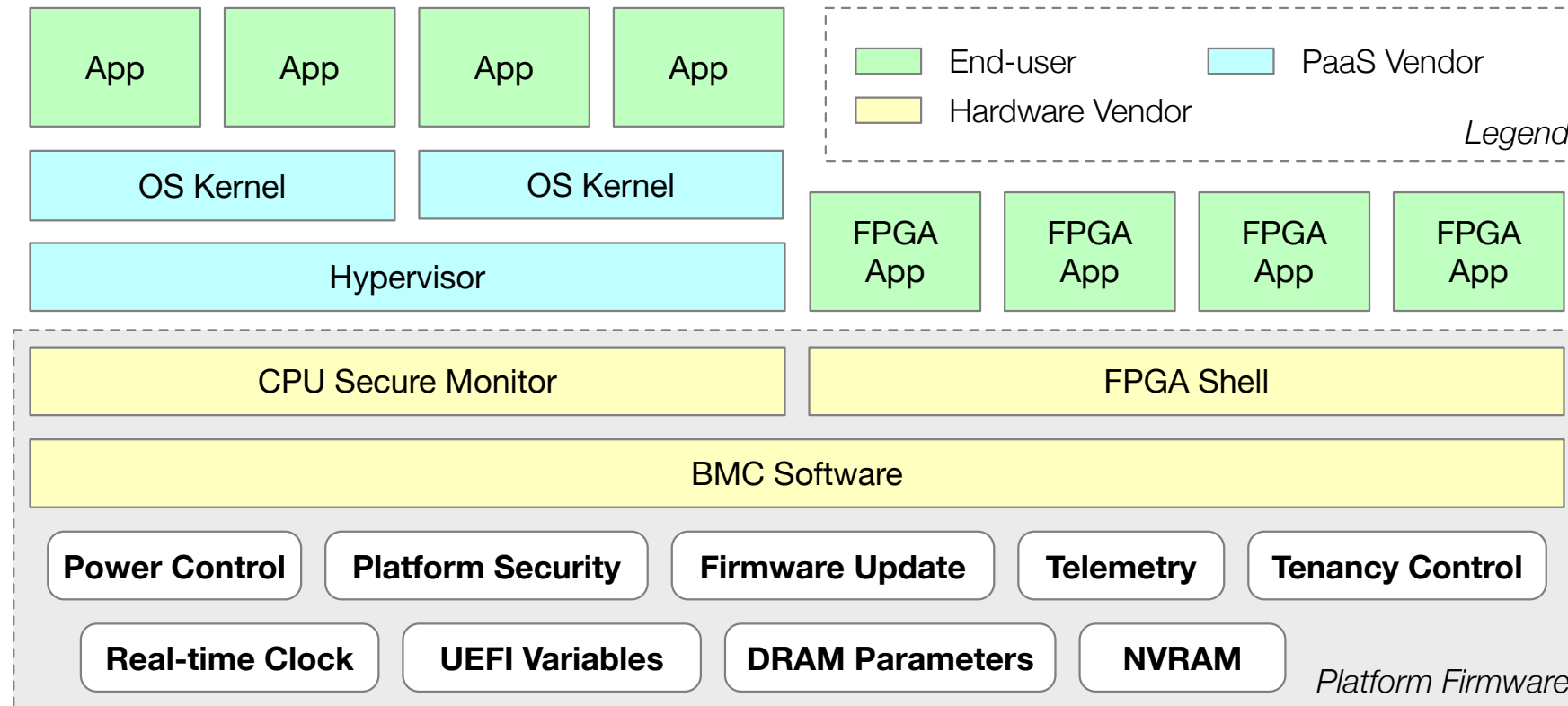
What is a firmware protocol?

- **Firmware protocol:** semantics of how to access firmware resources
 - From different places in the system



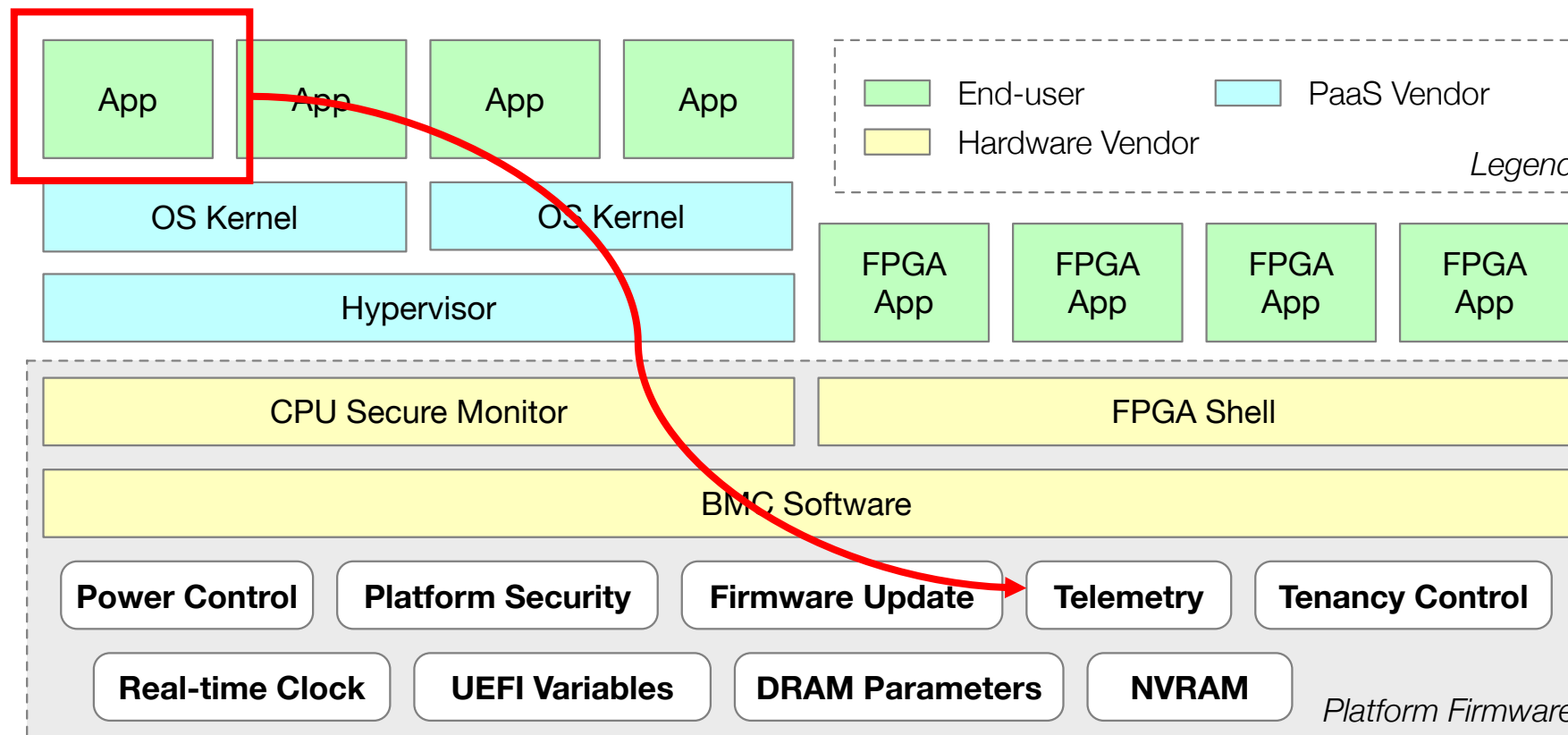
What is a firmware protocol?

- **Firmware protocol:** semantics of how to access firmware resources
 - From different places in the system



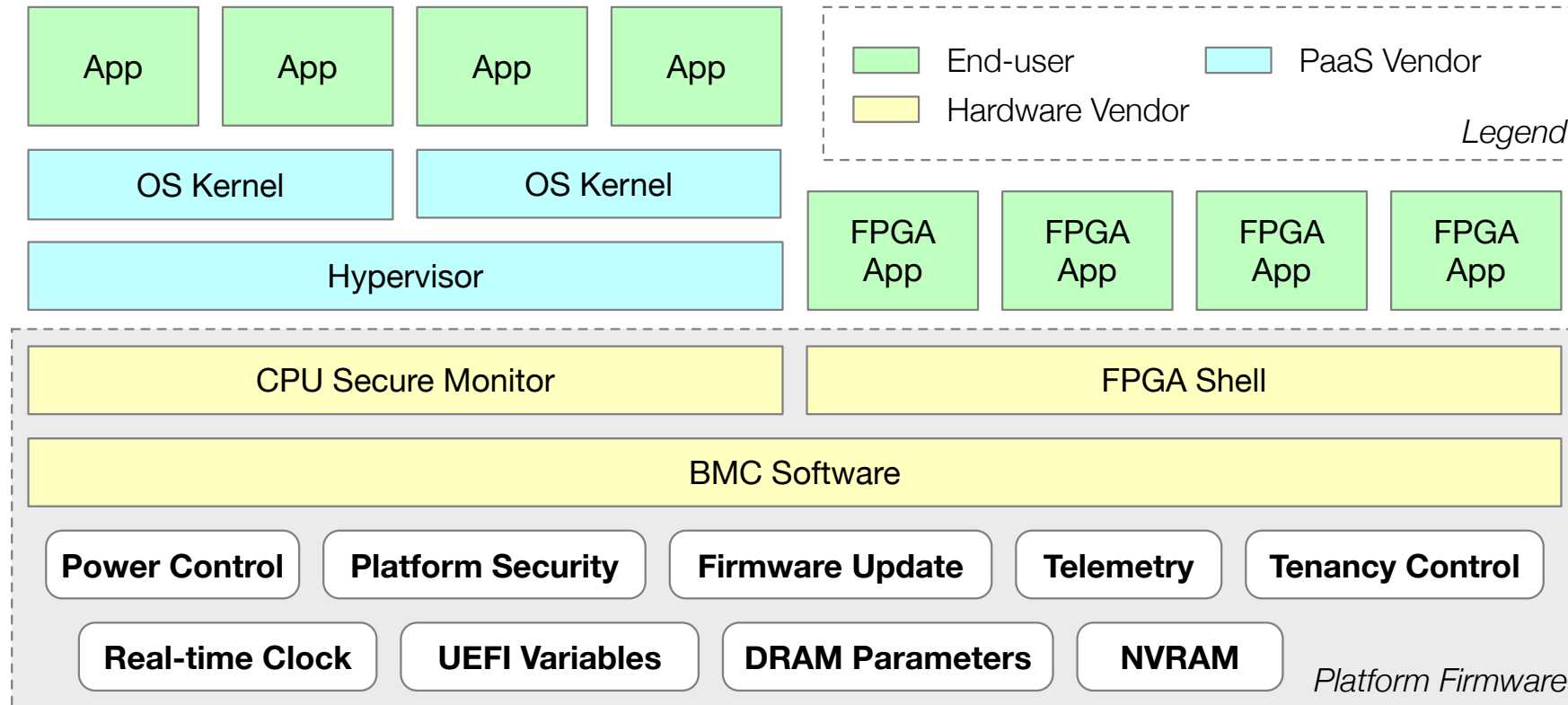
What is a firmware protocol?

- **Firmware protocol:** semantics of how to access firmware resources
 - From different places in the system



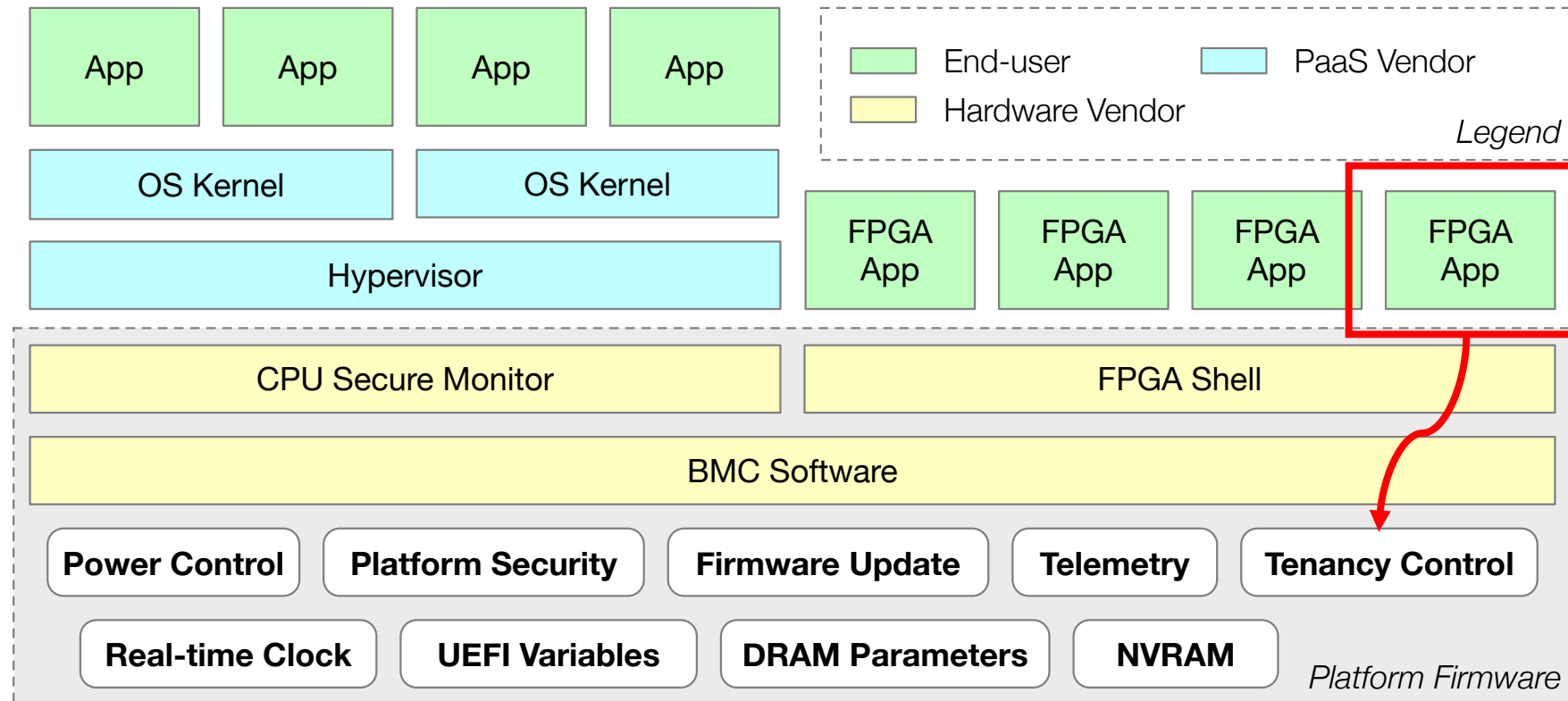
What is a firmware protocol?

- **Firmware protocol:** semantics of how to access firmware resources
 - From different places in the system



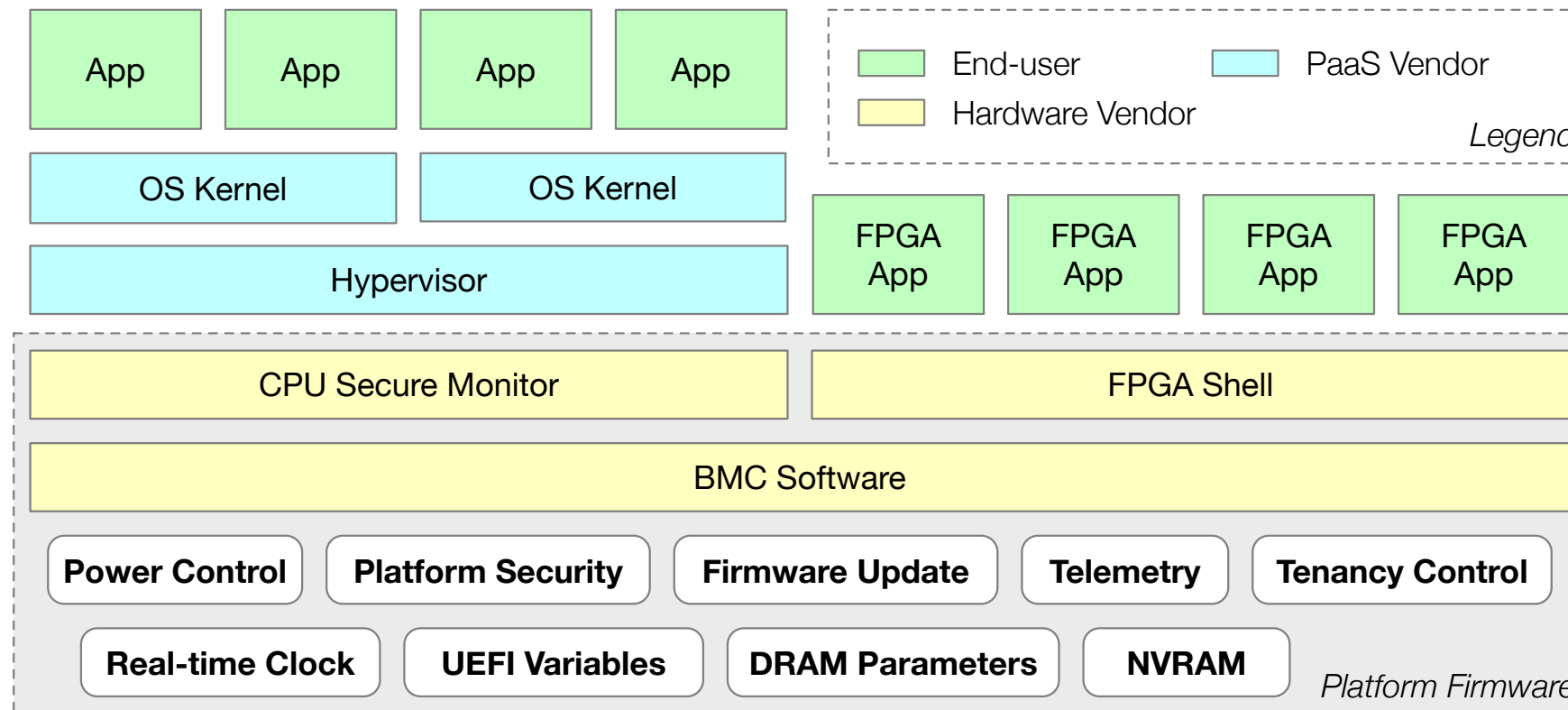
What is a firmware protocol?

- **Firmware protocol:** semantics of how to access firmware resources
 - From different places in the system

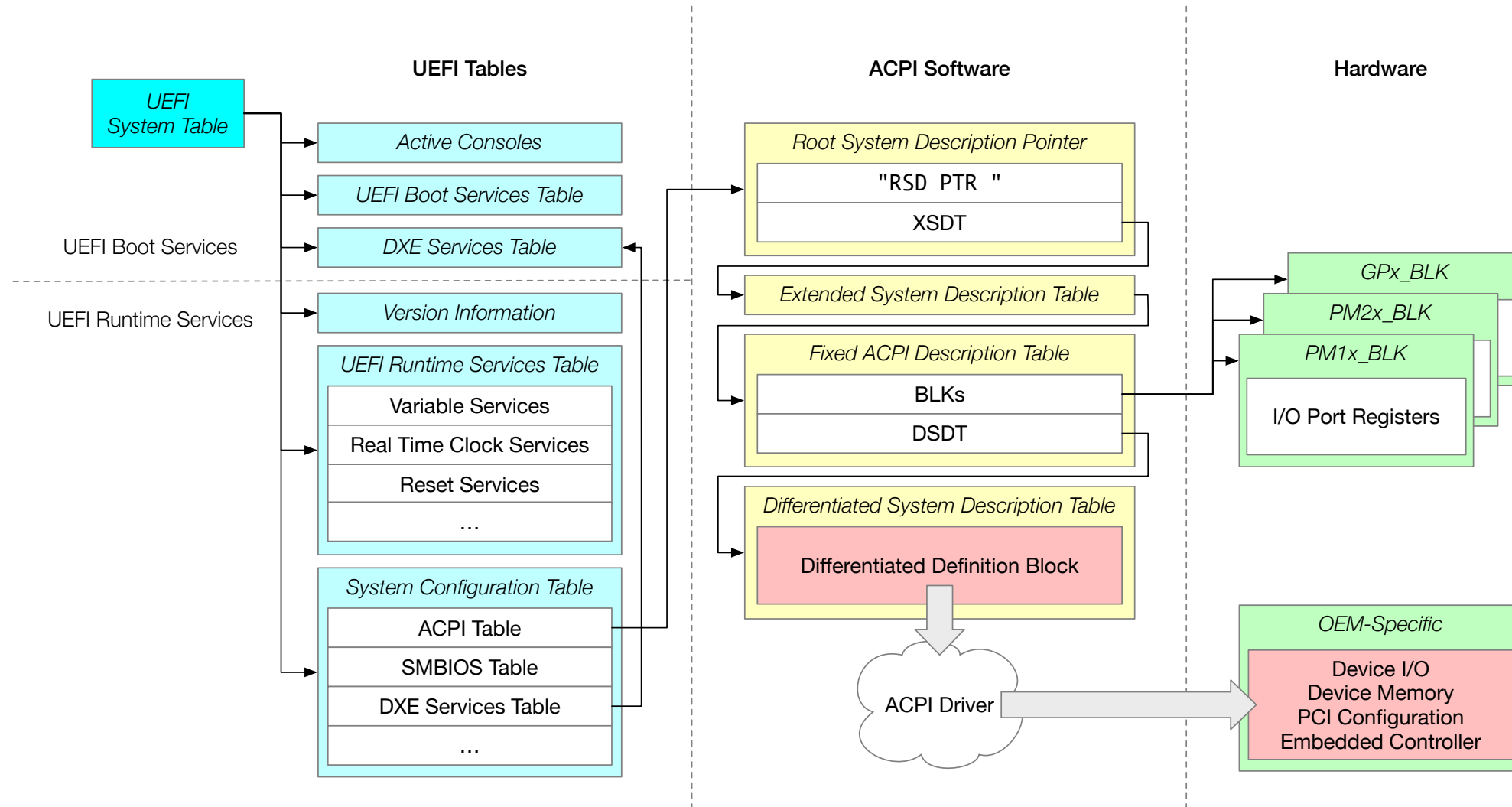


What is a firmware protocol?

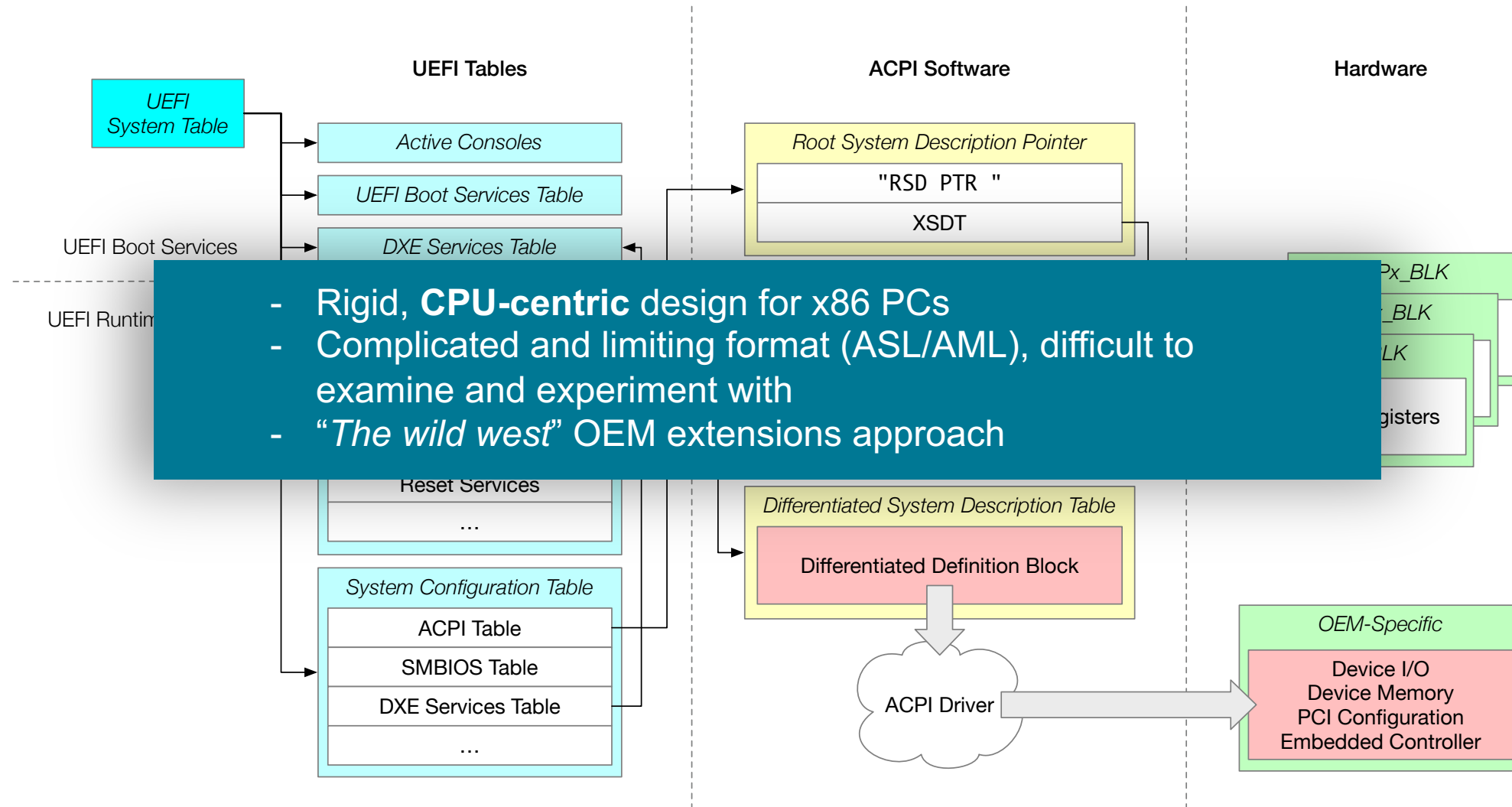
- **Firmware protocol:** semantics of how to access firmware resources
 - From different places in the system



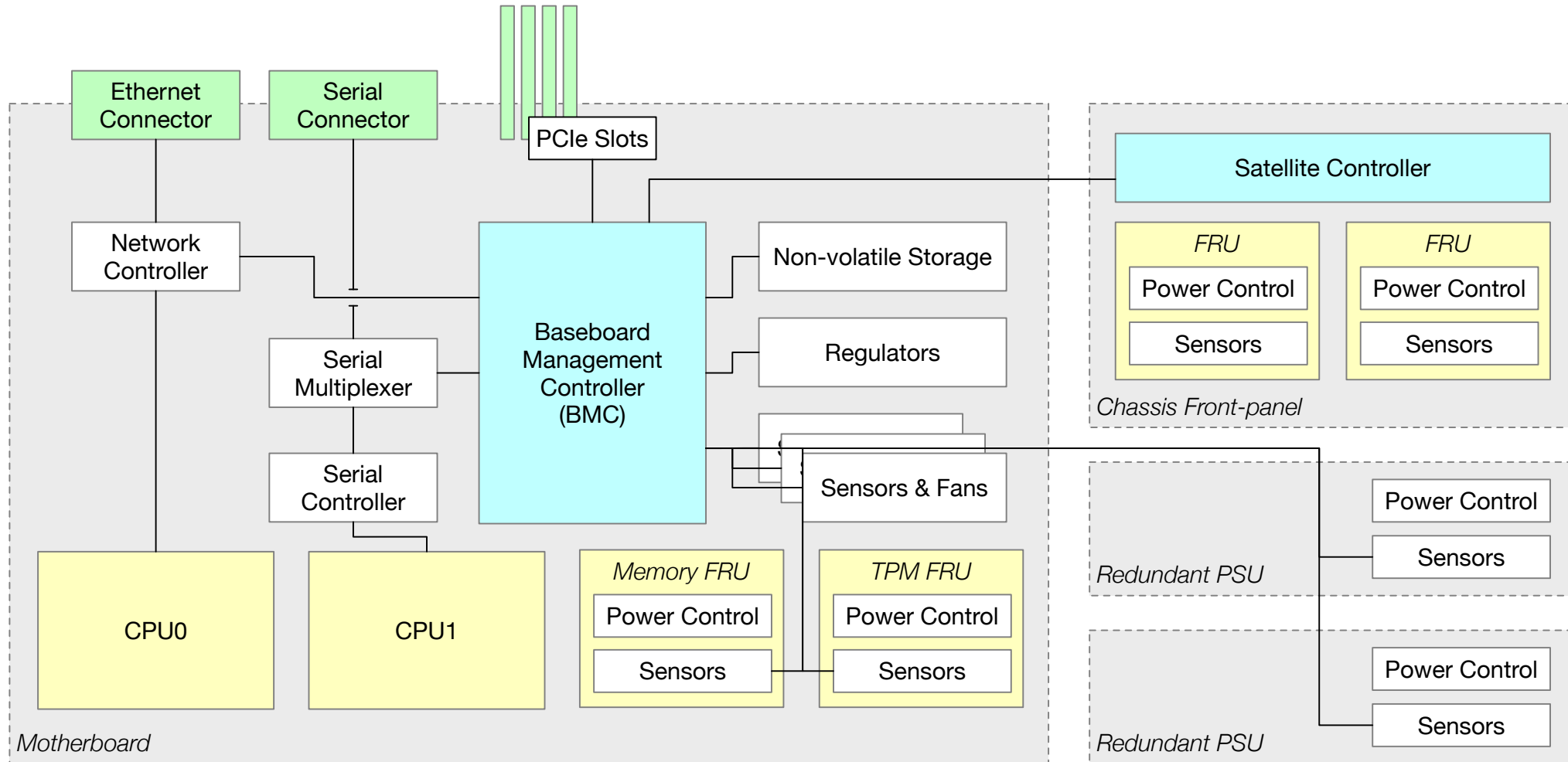
Existing firmware protocols – UEFI/ACPI



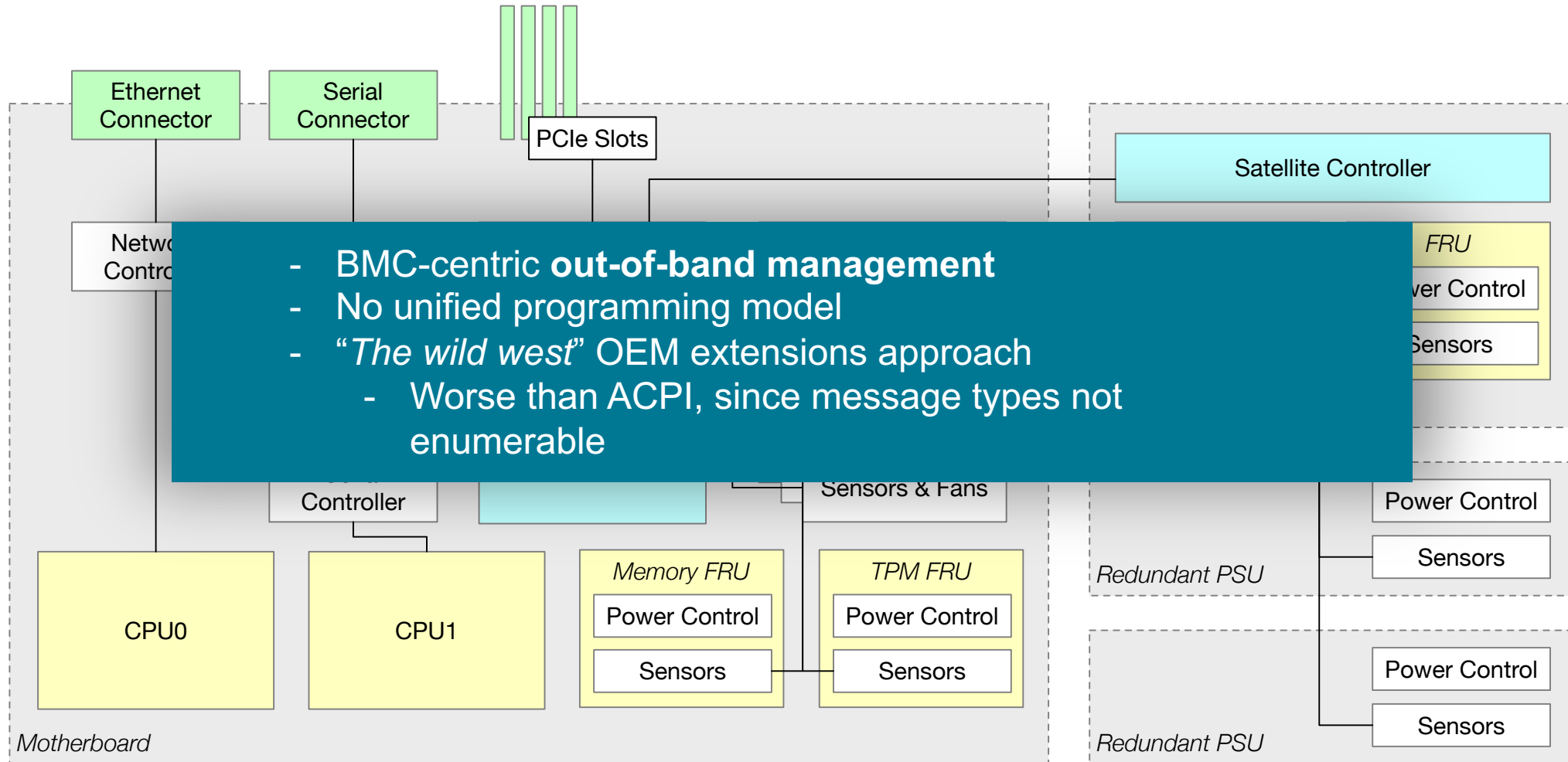
Existing firmware protocols – UEFI/ACPI



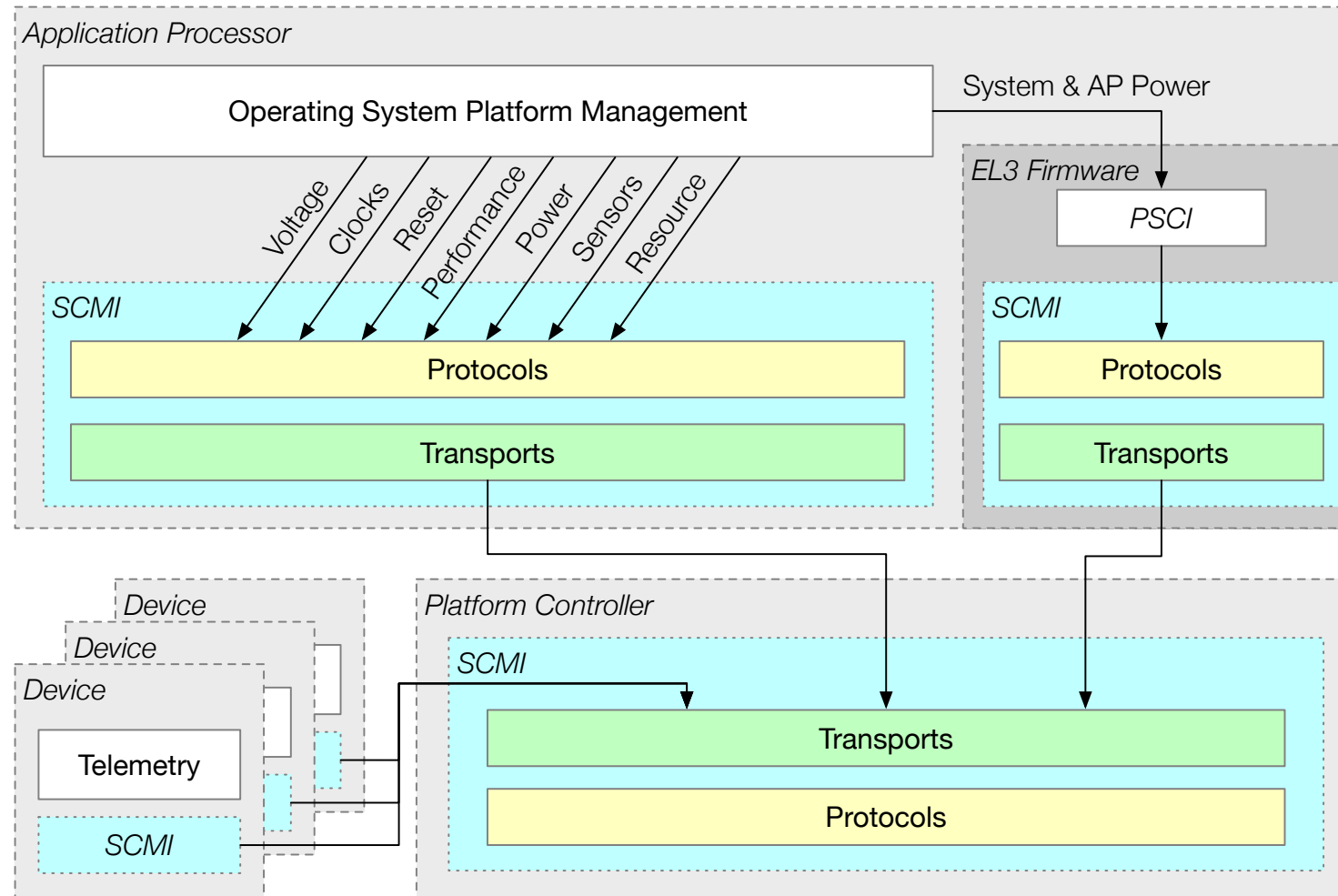
Existing firmware protocols – IPMI



Existing firmware protocols – IPMI



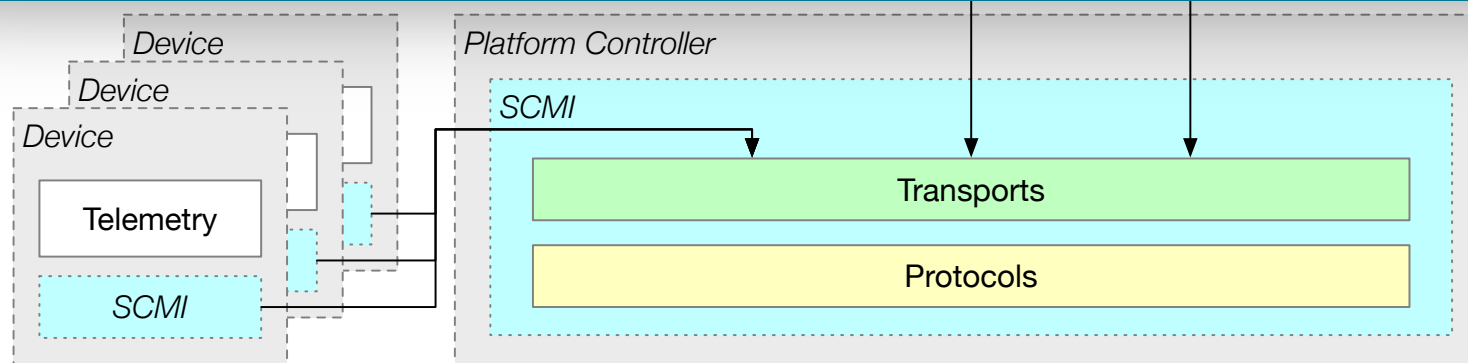
Existing firmware protocols – SCMI



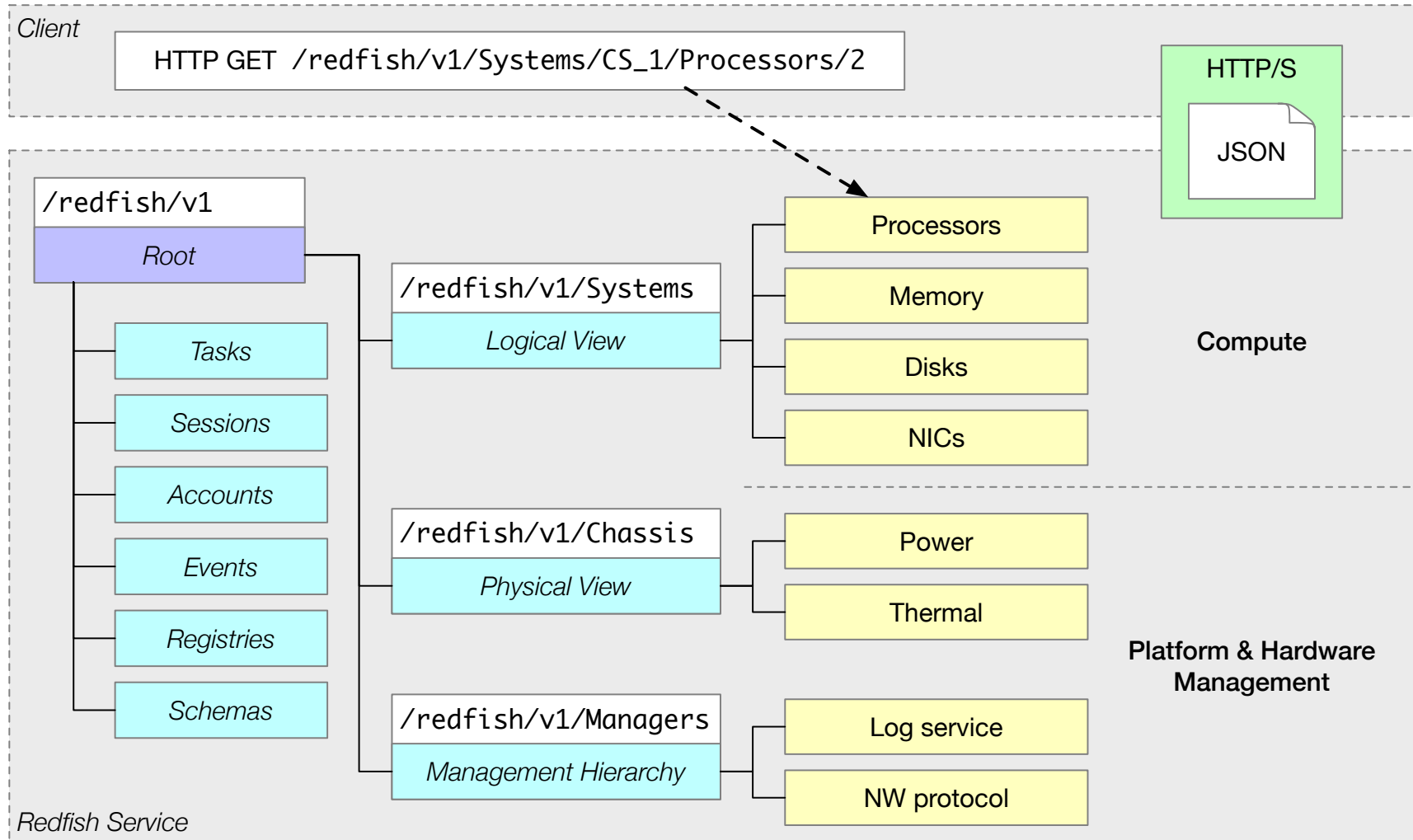
Existing firmware protocols – SCMI



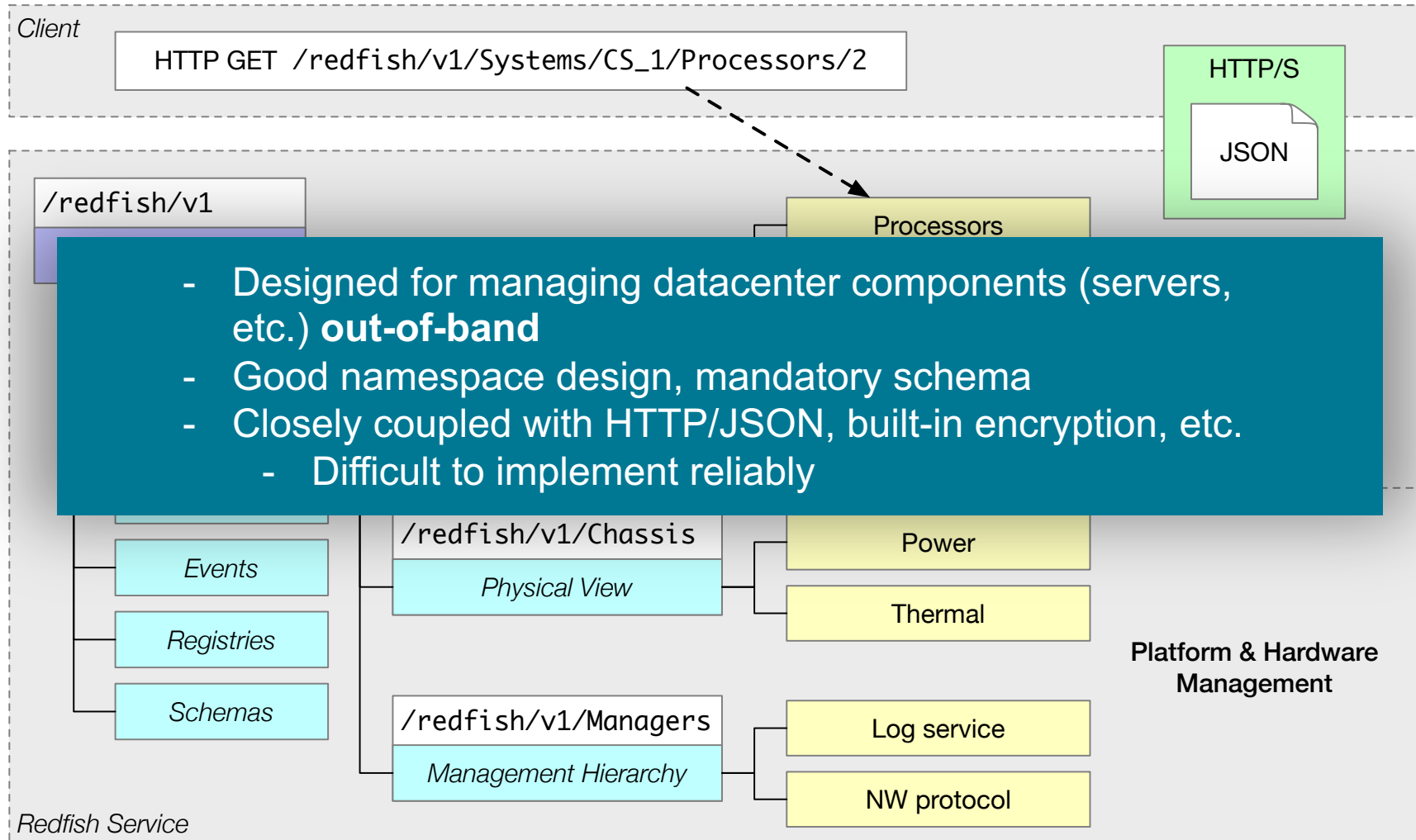
- Targets SoC platforms with clear role delegations
- Relatively new (2017) thus few adopters
 - Little documentation, restrictive access
- Requires mapping from resource ID to actual resource
 - Device tree model on Linux

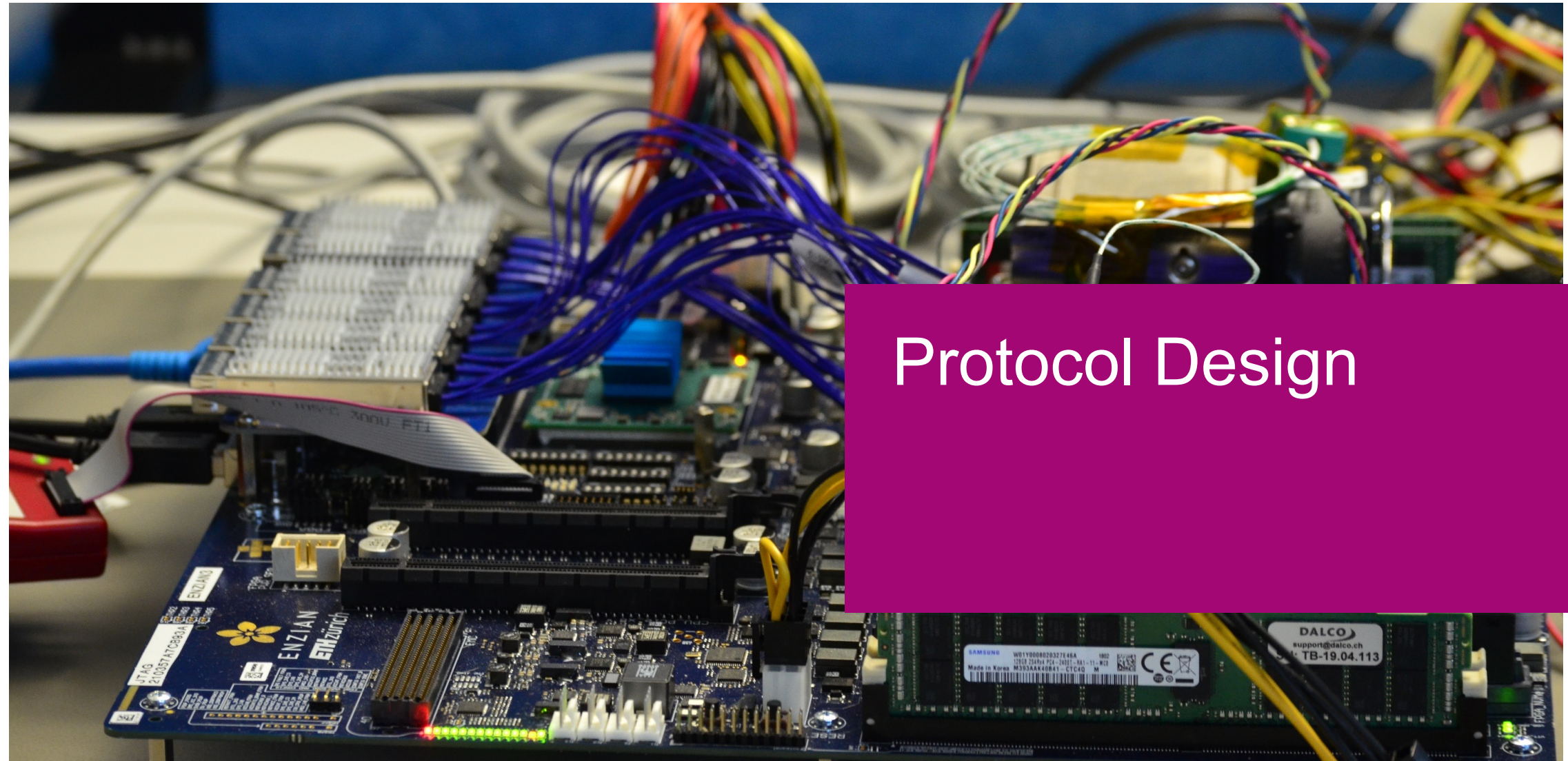


Existing firmware protocols – DMTF Redfish



Existing firmware protocols – DMTF Redfish





Protocol Design

Requirements for a suitable firmware protocol for Enzian



DINFK

Recall that Enzian is a **research** *heterogeneous* system...

Requirements for a suitable firmware protocol for Enzian

Recall that Enzian is a **research** *heterogeneous* system...

Extensible

Organised and discoverable

Modular and composable

Simple to implement

Inspectable

Usable

Overview – Enzian Firmware Resource Interface



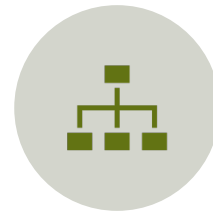
Stateless RPC



Synchronous
with async
extension



Shared actions
on resources



Hierarchical
namespaces



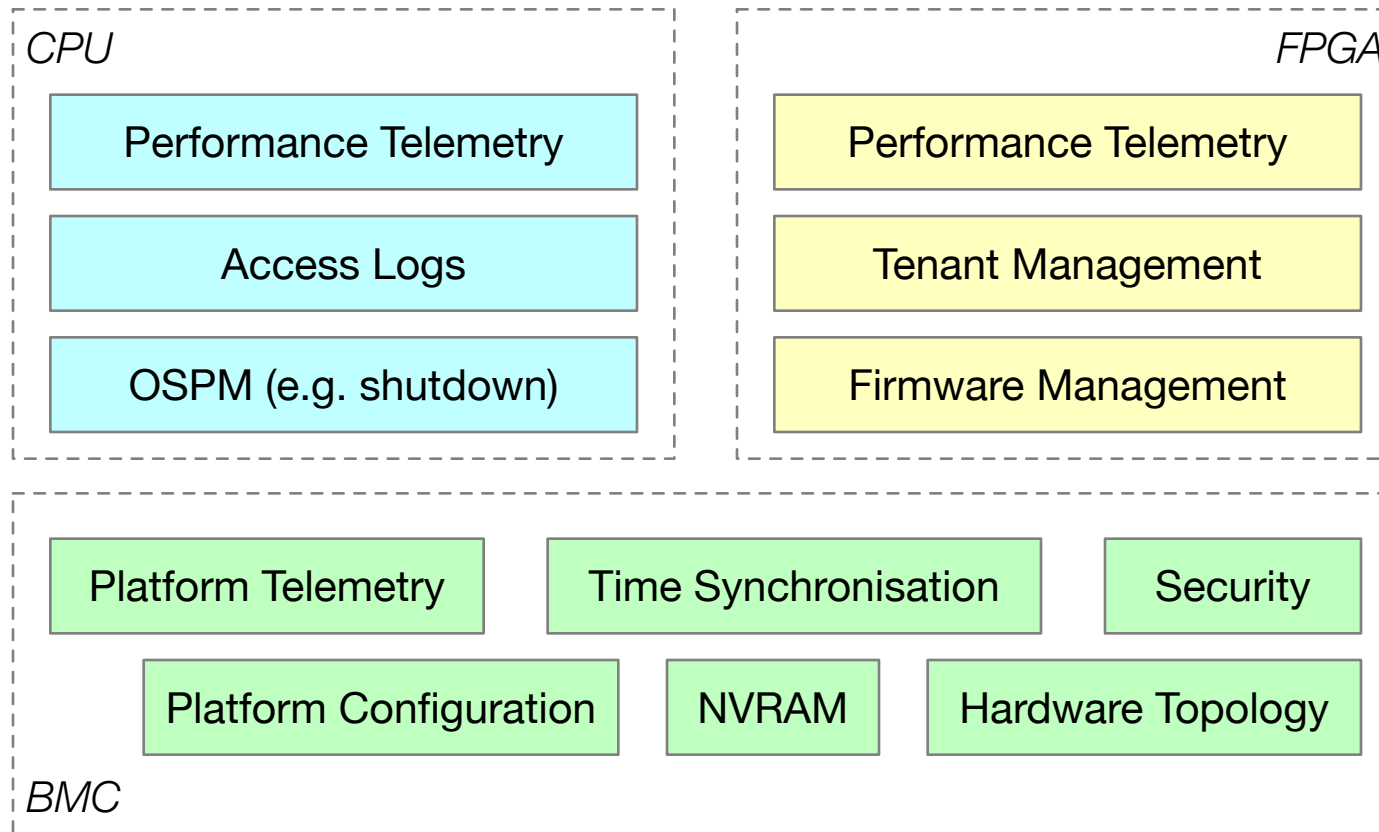
Schema-based



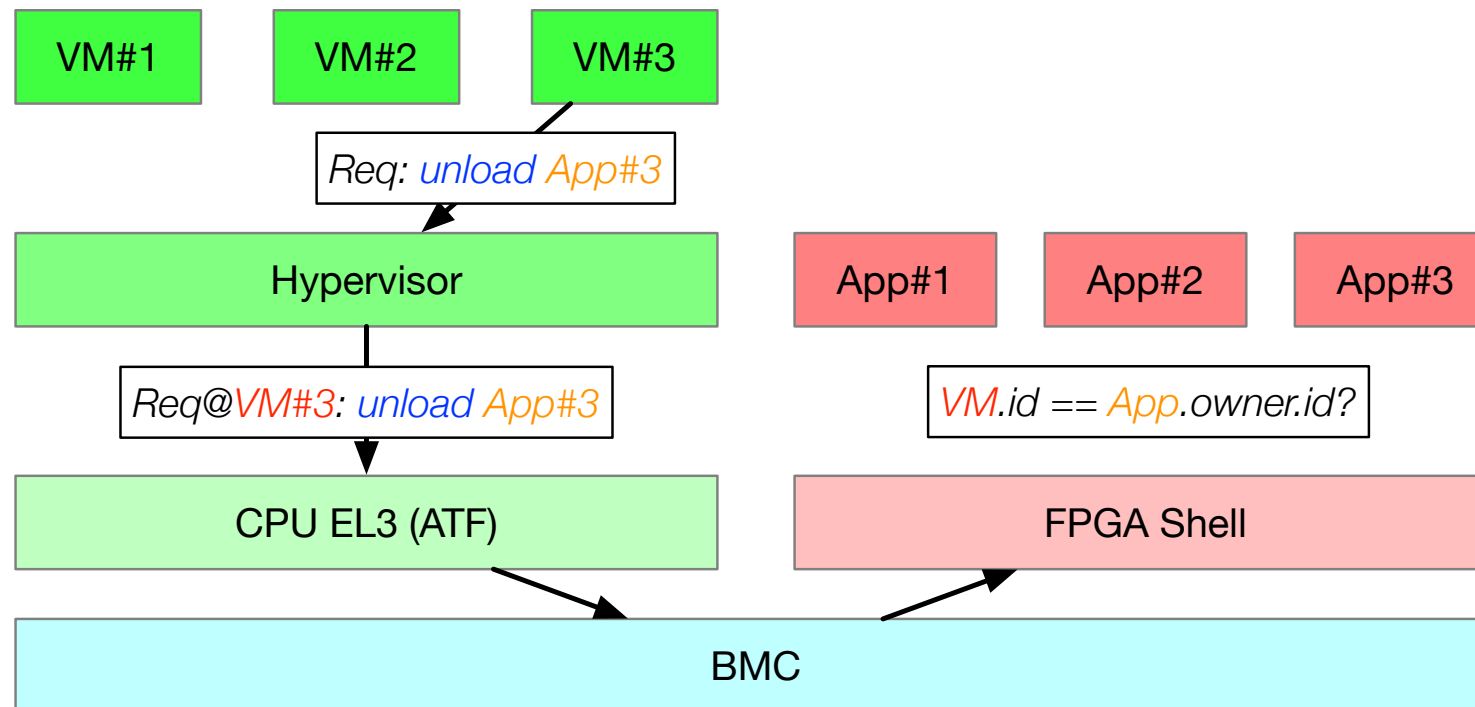
Symmetric
client/actor
model

Symmetric Client/Actor Design

- Every hardware component in the system can be both a client and an actor



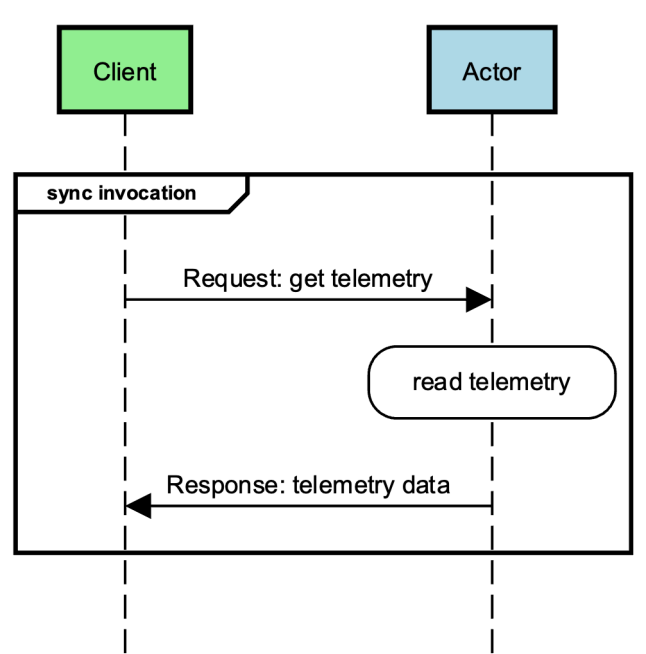
- Attack model: trust all **forwarding agents** along the way a request traverses
 - Permission check happens along this route



Asynchronous Event Mechanism (WIP)

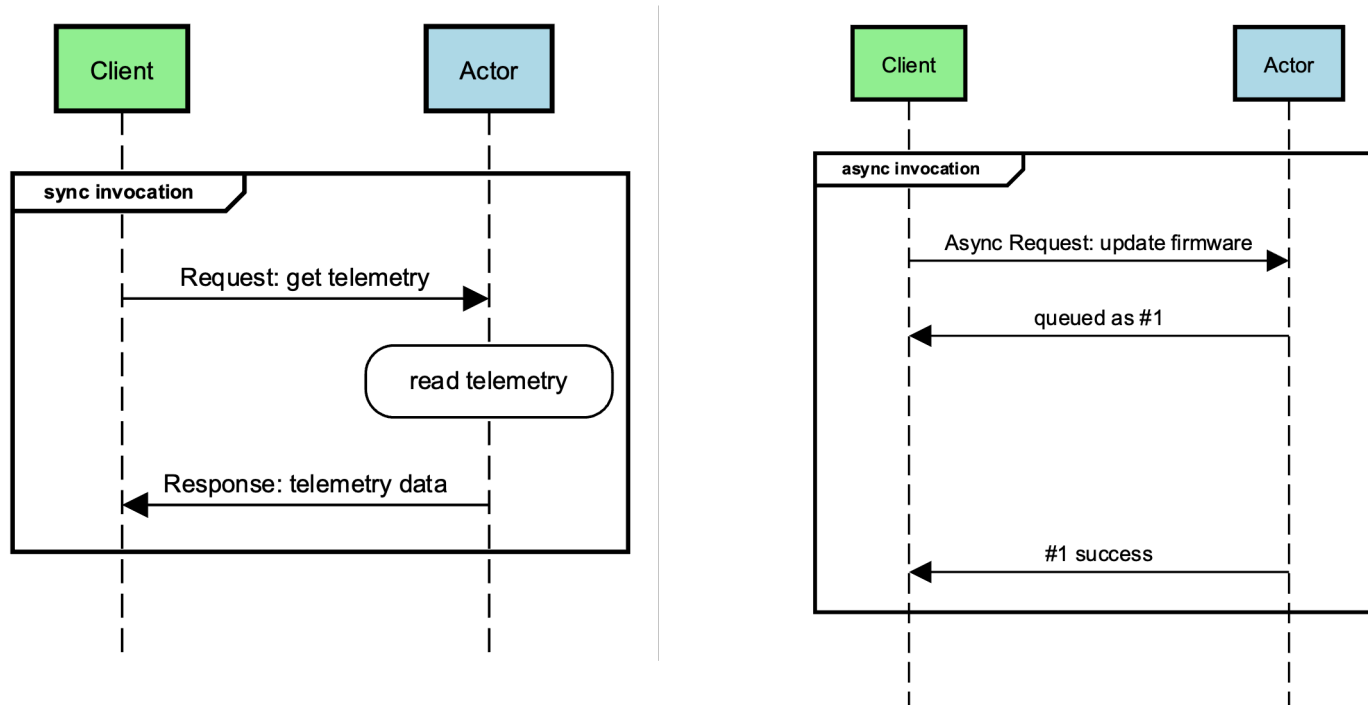
- Default synchronous due to in-band nature
 - Asynchrony useful for **slow** or **periodic** requests
 - Eliminates protocol overhead to generate request

Asynchronous Event Mechanism (WIP)



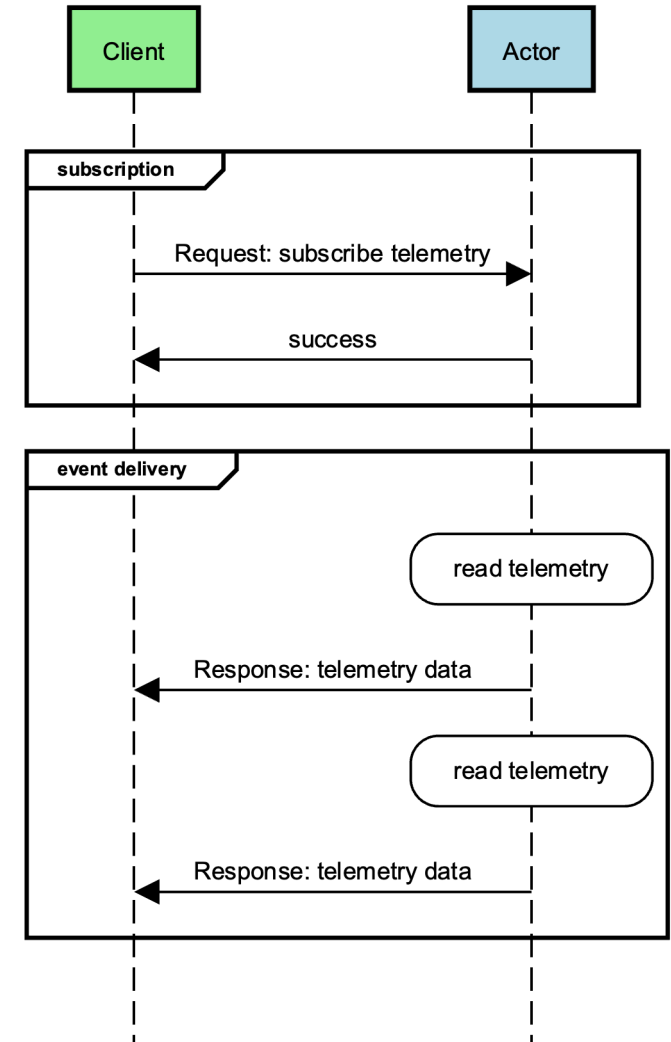
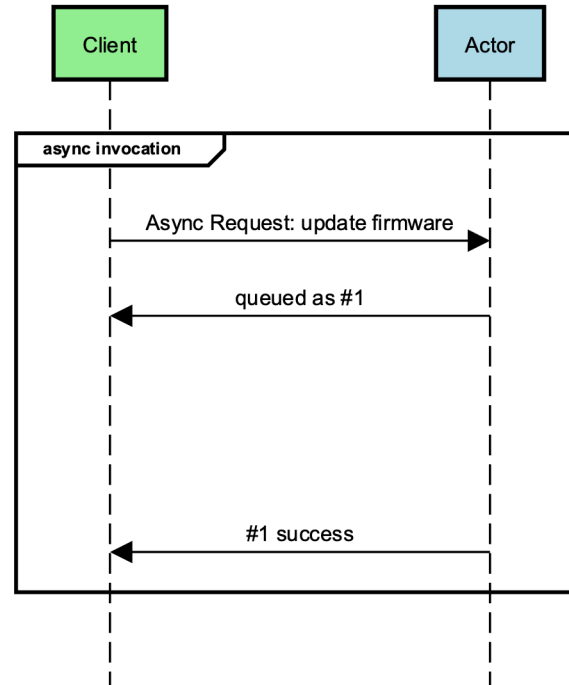
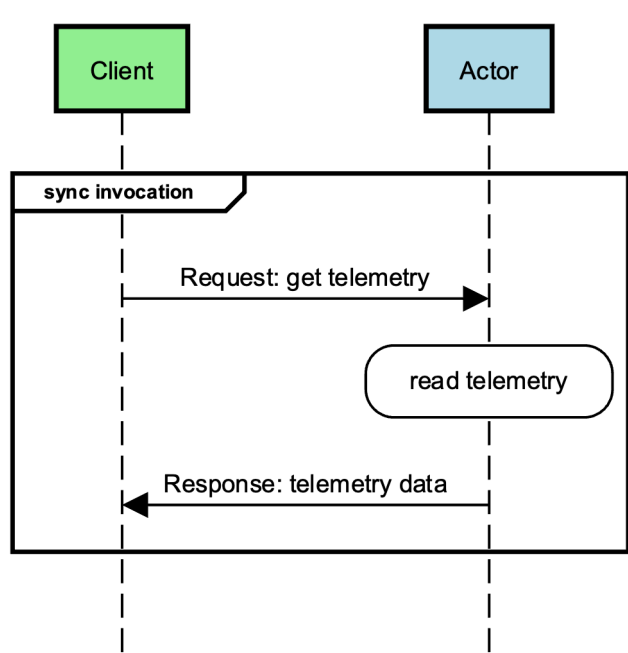
- Default synchronous due to in-band nature
 - Asynchrony useful for **slow** or **periodic** requests
 - Eliminates protocol overhead to generate request

Asynchronous Event Mechanism (WIP)

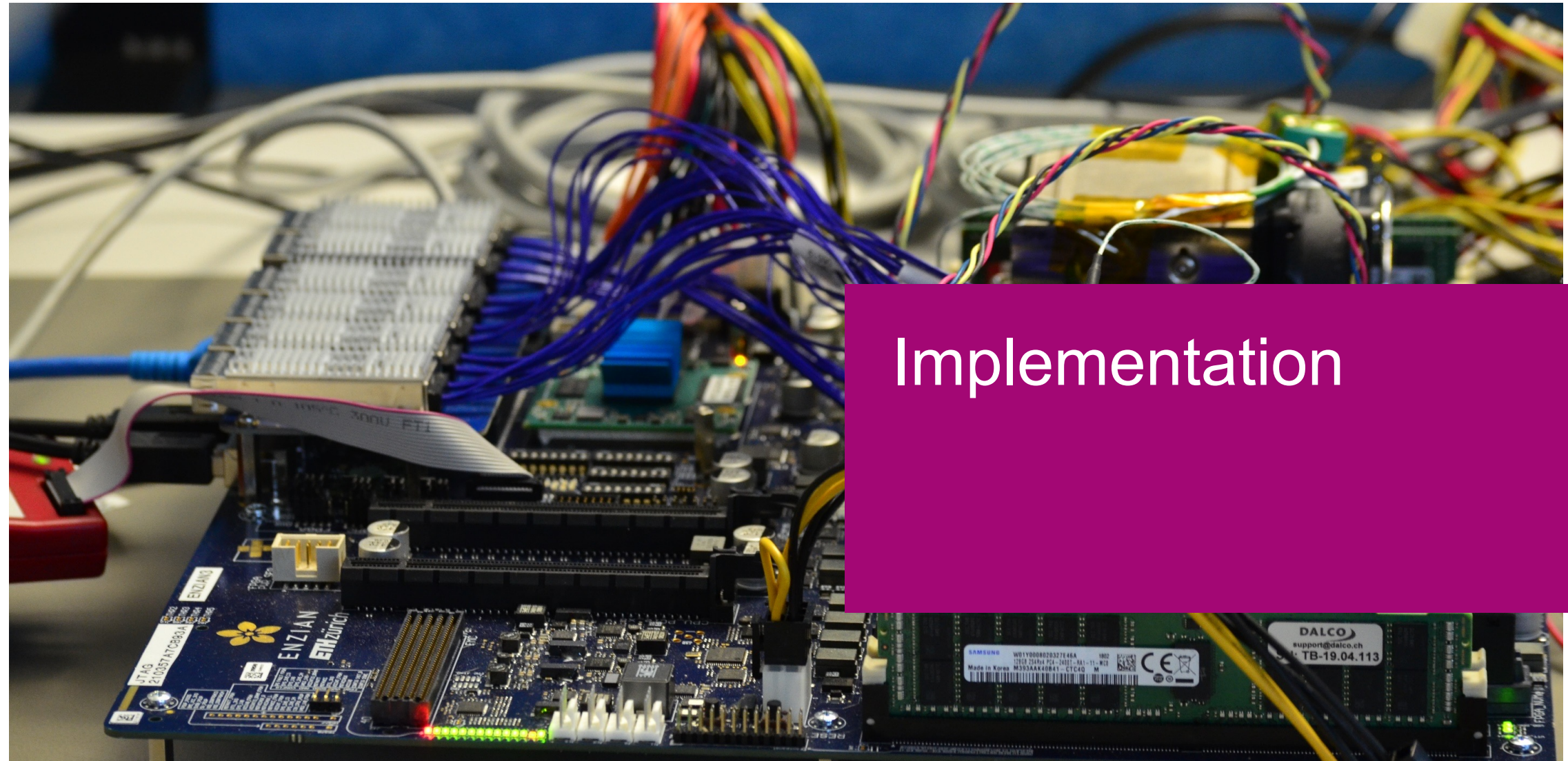


- Default synchronous due to in-band nature
 - Asynchrony useful for **slow** or **periodic** requests
 - Eliminates protocol overhead to generate request

Asynchronous Event Mechanism (WIP)

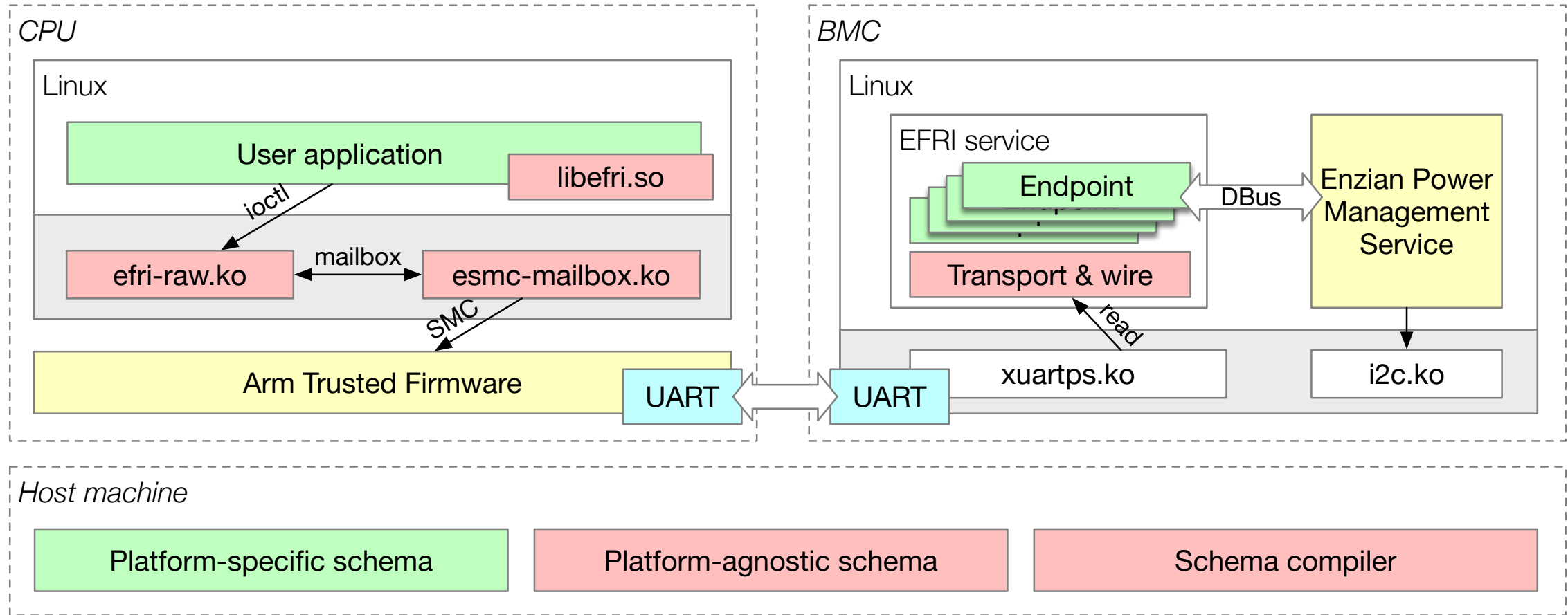


- Default synchronous due to in-band nature
 - Asynchrony useful for **slow** or **periodic** requests
 - Eliminates protocol overhead to generate request



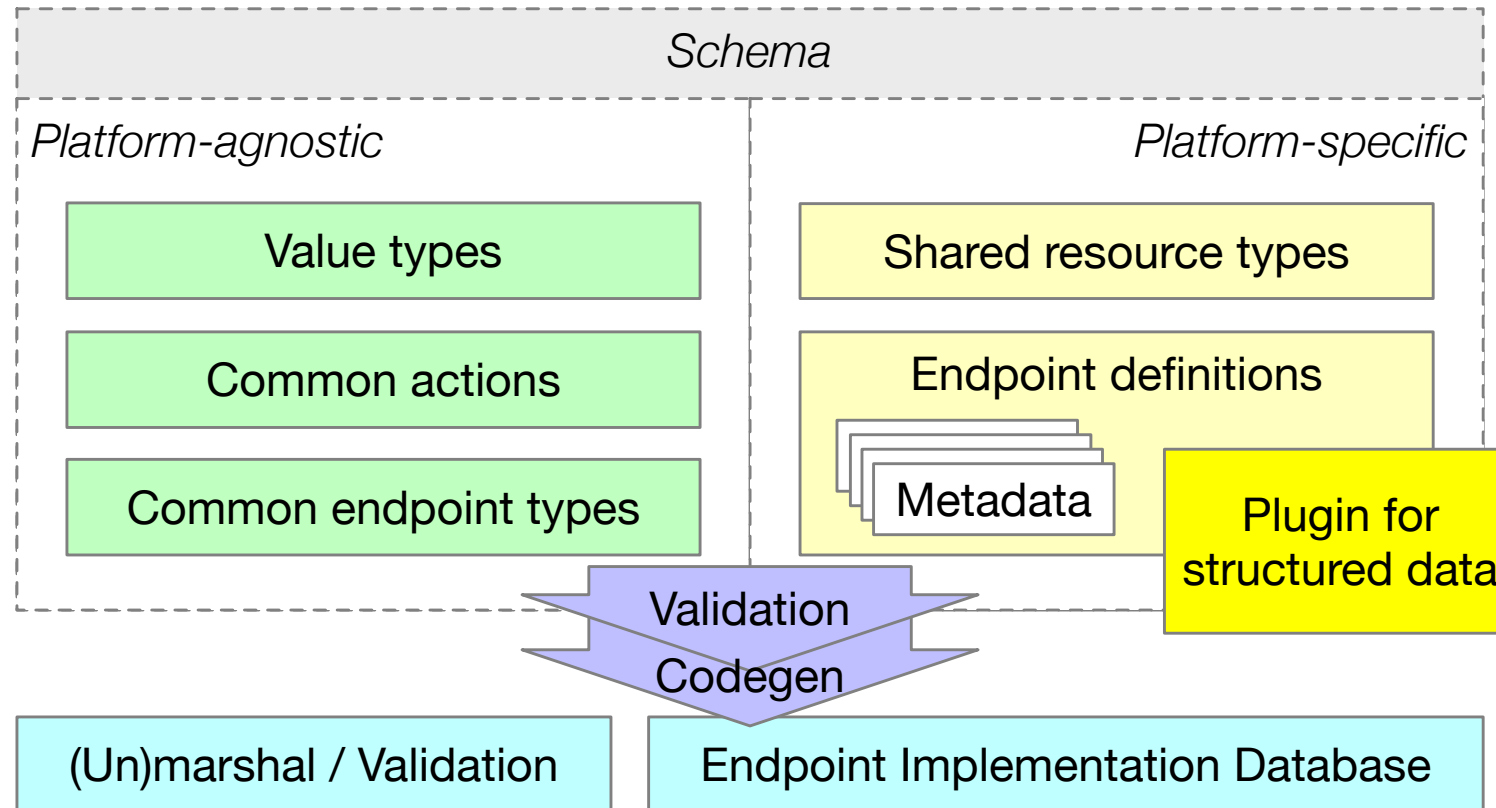
Implementation

Reference implementation overview

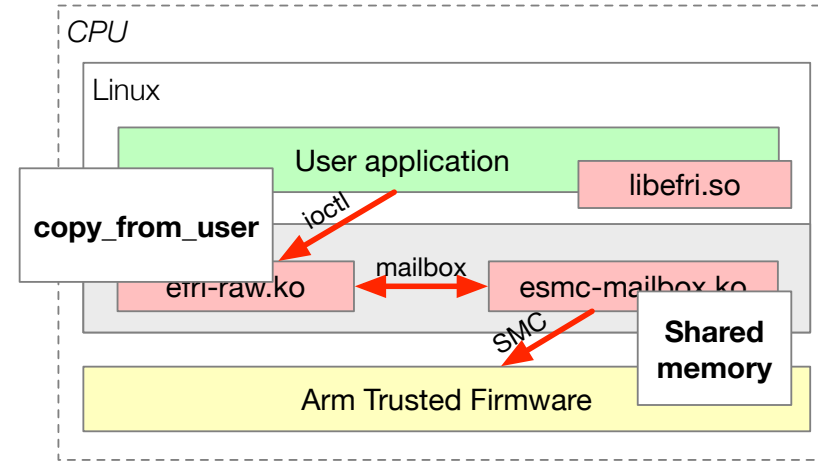
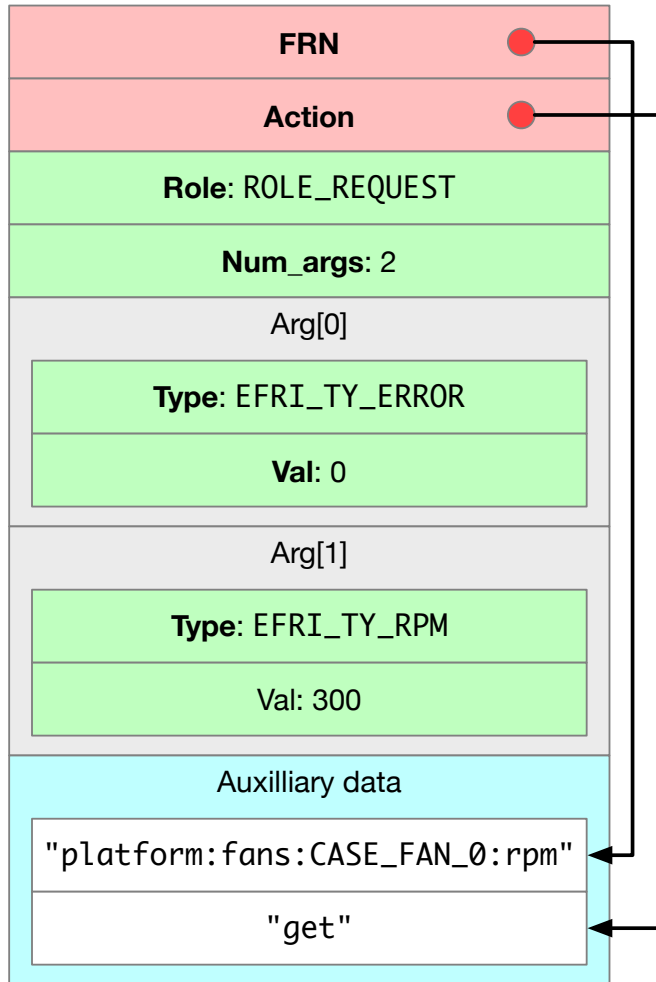


The schema compiler

- Generate boilerplate code from YAML schema



Implementation on the CPU

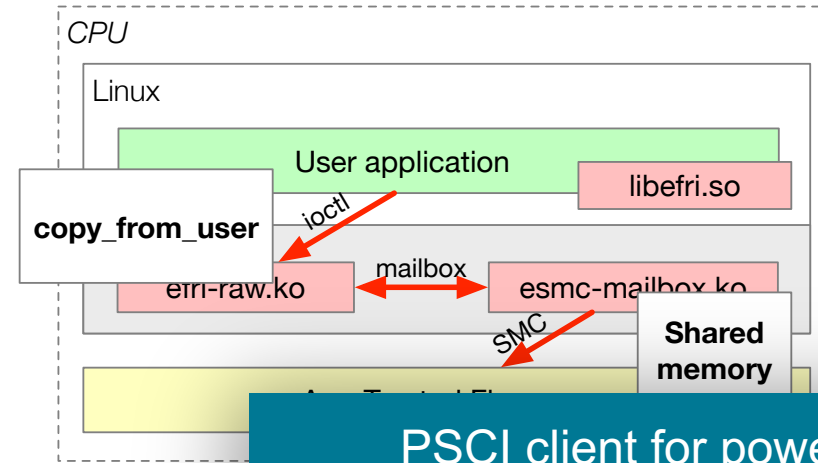
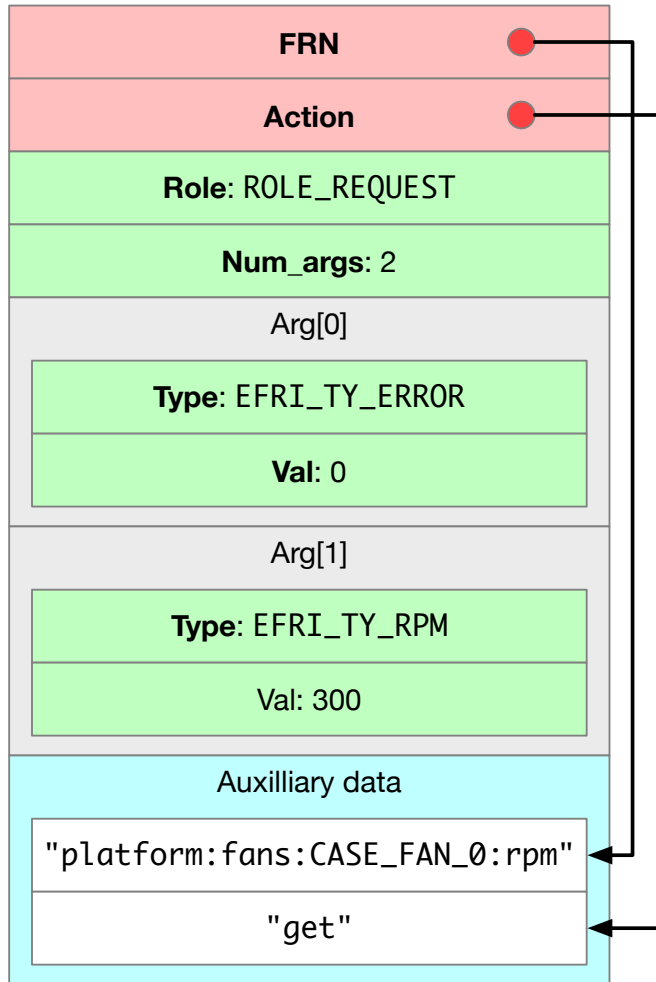


```

==EFRI== => platform:fans:CASE_FAN_0:rpm get \n
           request preamble          FRN          action

==EFRI== <= platform:fans:CASE_FAN_0:rpm get Error$0 Rpm$300#RPM \n
           response preamble          FRN          action error code value
    
```

Implementation on the CPU



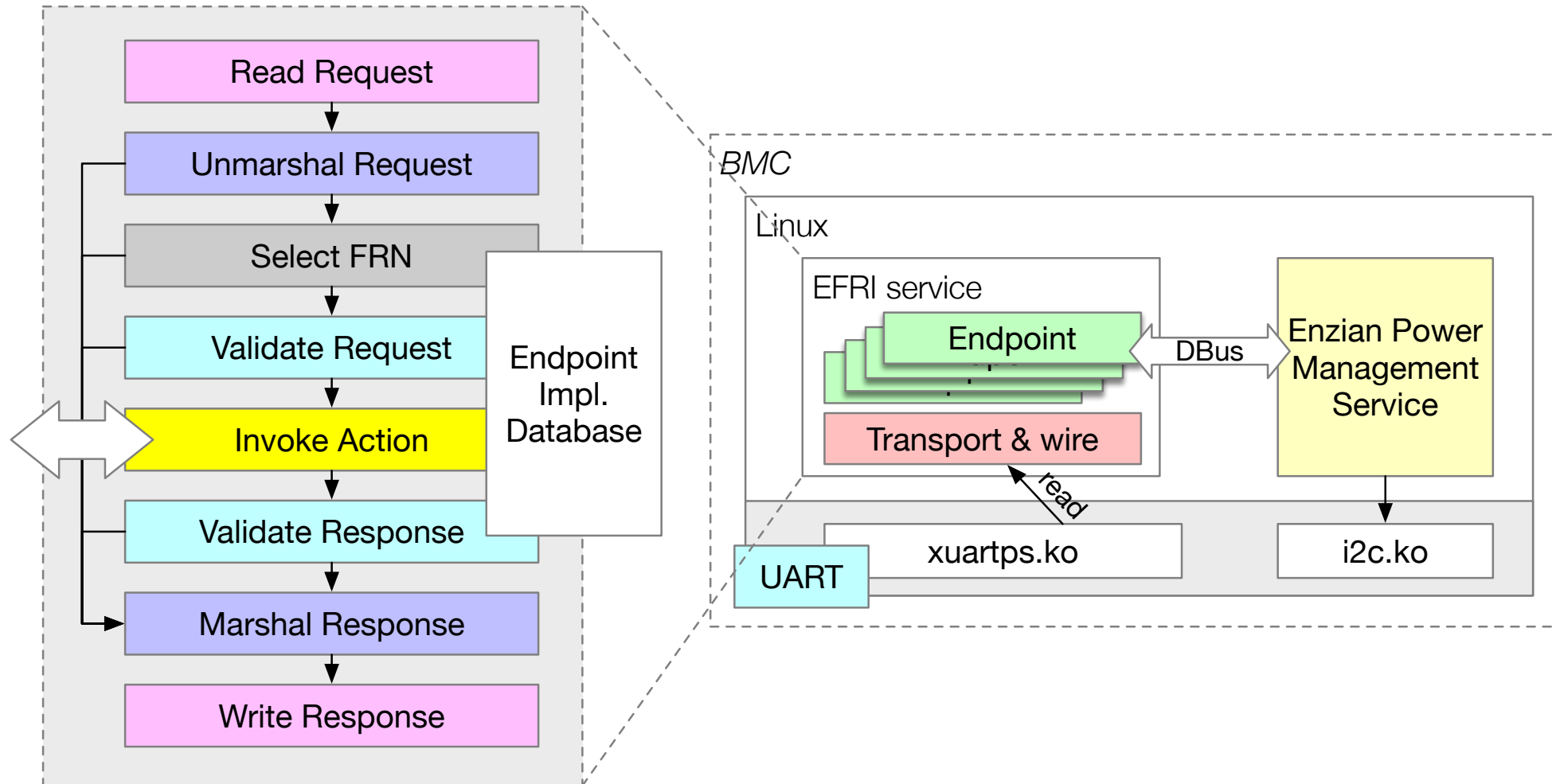
PSCI client for power-off!
DRAM parameters, etc.

```

==EFRI== => platform:fans:CASE_FAN_0:rpm get \n
           request preamble          FRN          action

==EFRI== <= platform:fans:CASE_FAN_0:rpm get Error$0 Rpm$300#RPM \n
           response preamble          FRN          action error code value
    
```

Implementation on the BMC

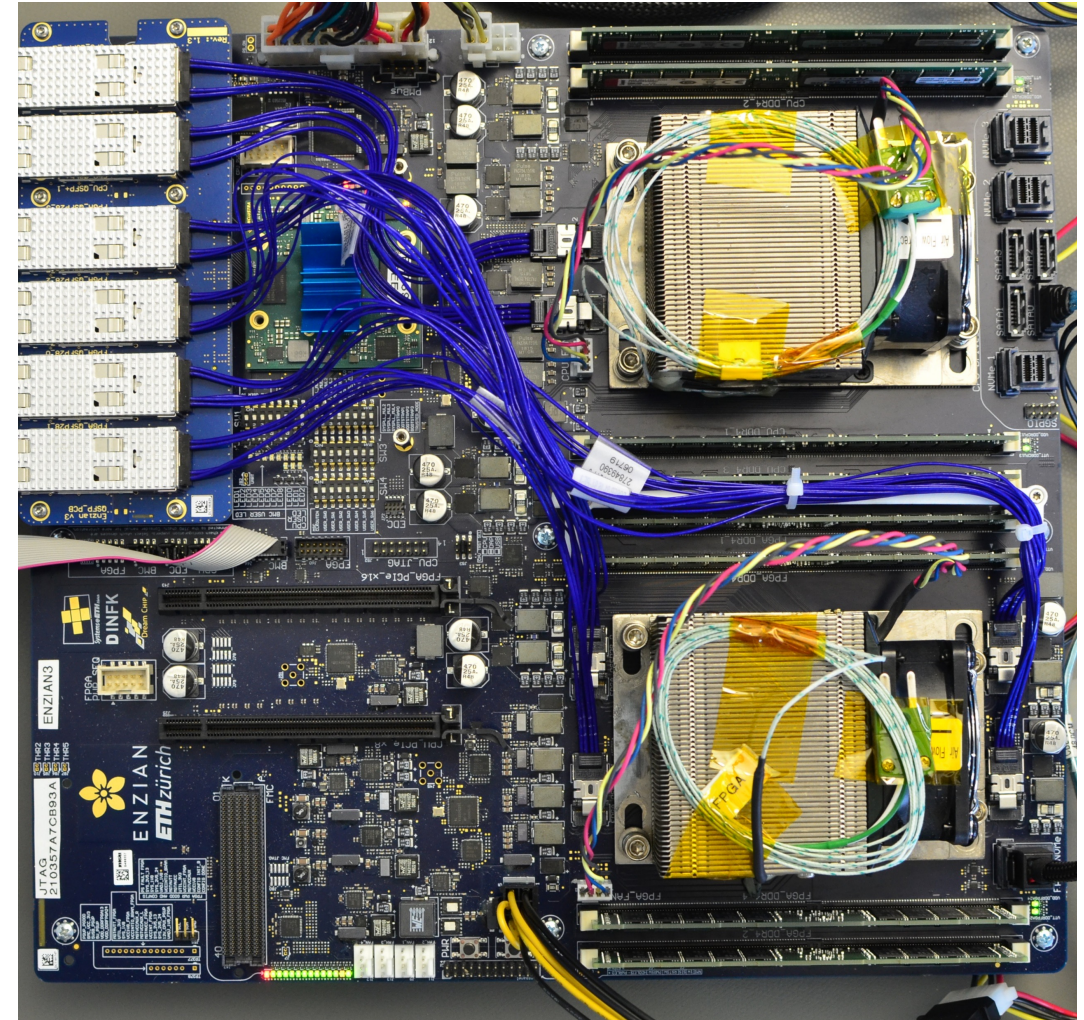




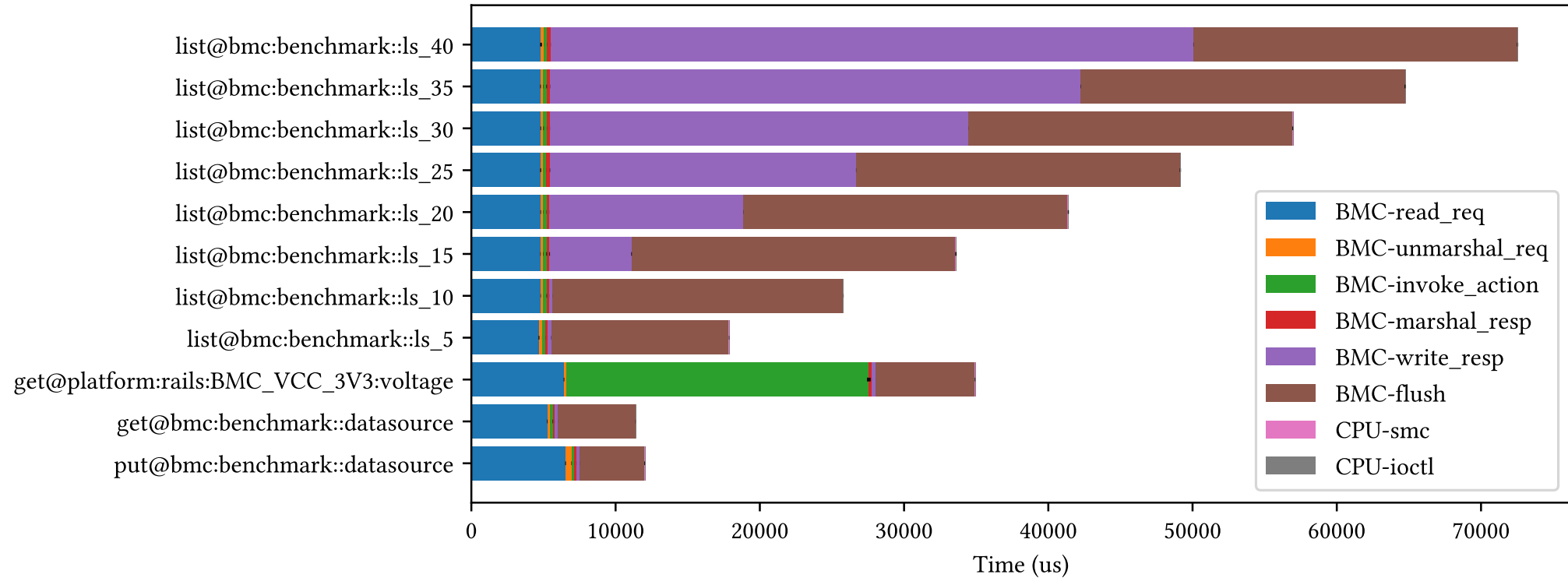
Evaluation

Experiments Platform

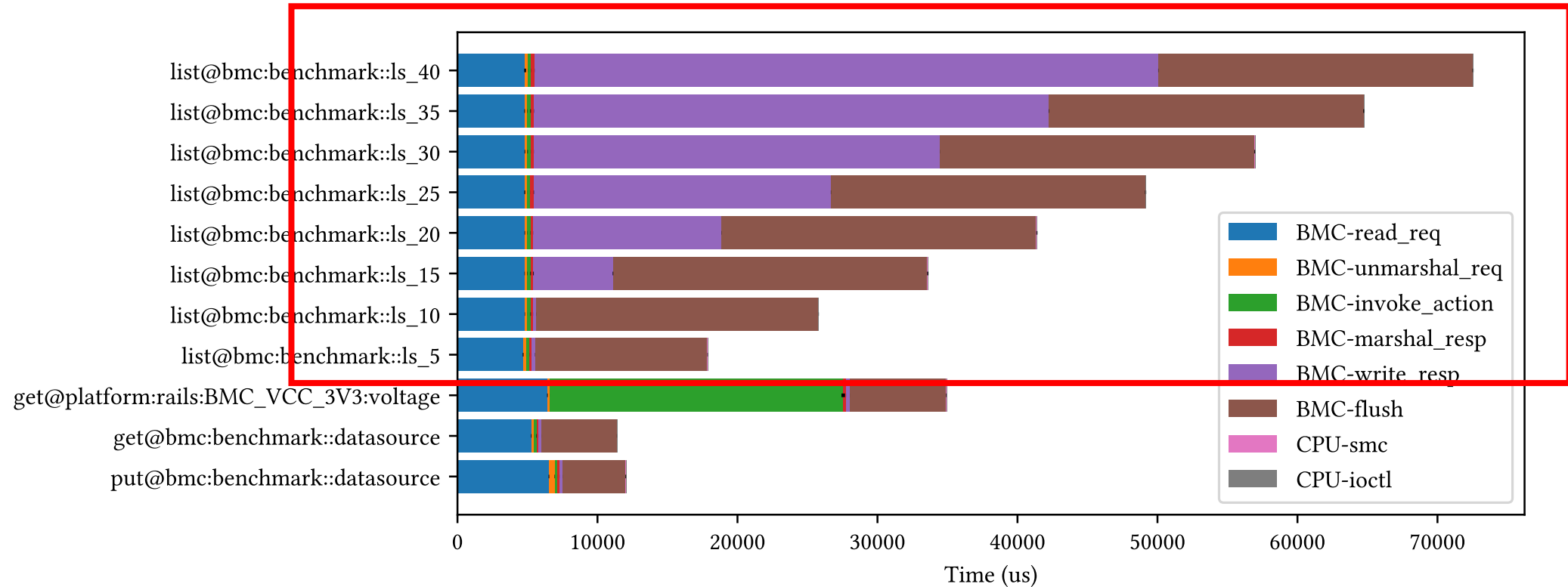
- Enzian v3 “zuestoll”
 - ThunderX-1, Zynq-7000 BMC
- 115200 8n1 for the serial link
- Toolchain:
 - Schema compiler generates ATF and BMC code
 - OpenBMC + Python 3
 - ATF codebase version unknown (from Cavium)
- Codebase augmented to collect timestamps at various stage of processing



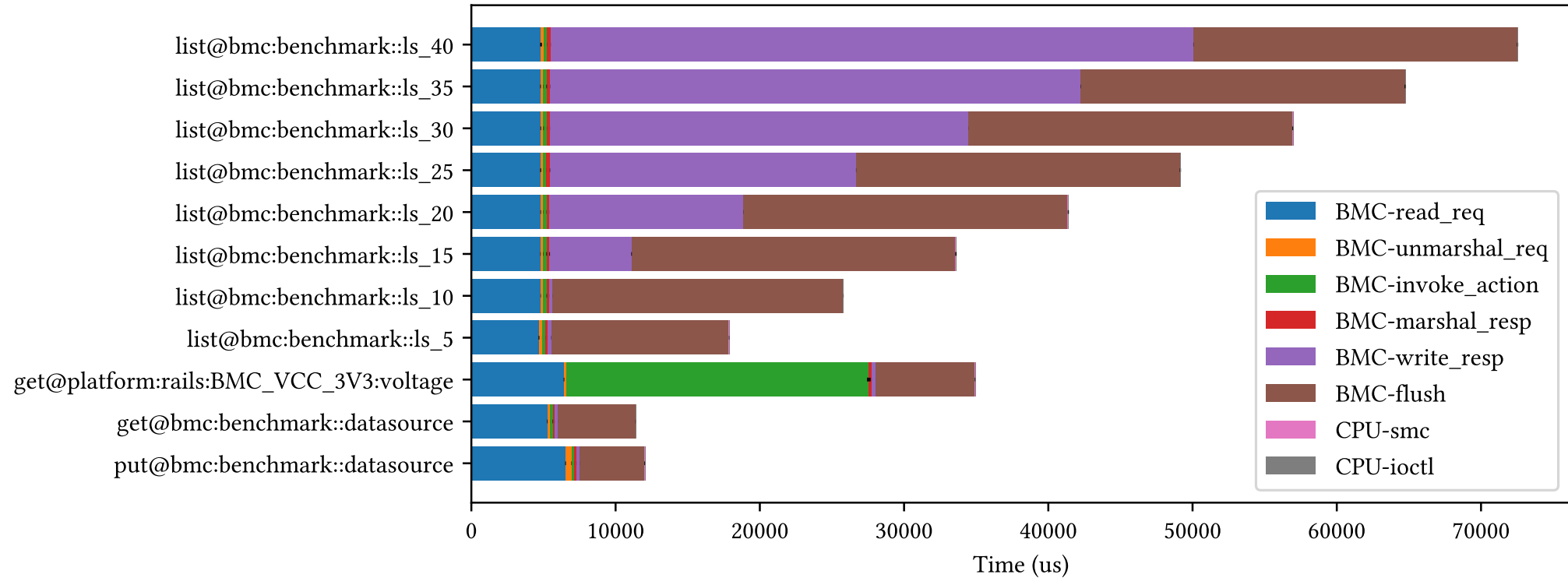
Performance – latency breakdown



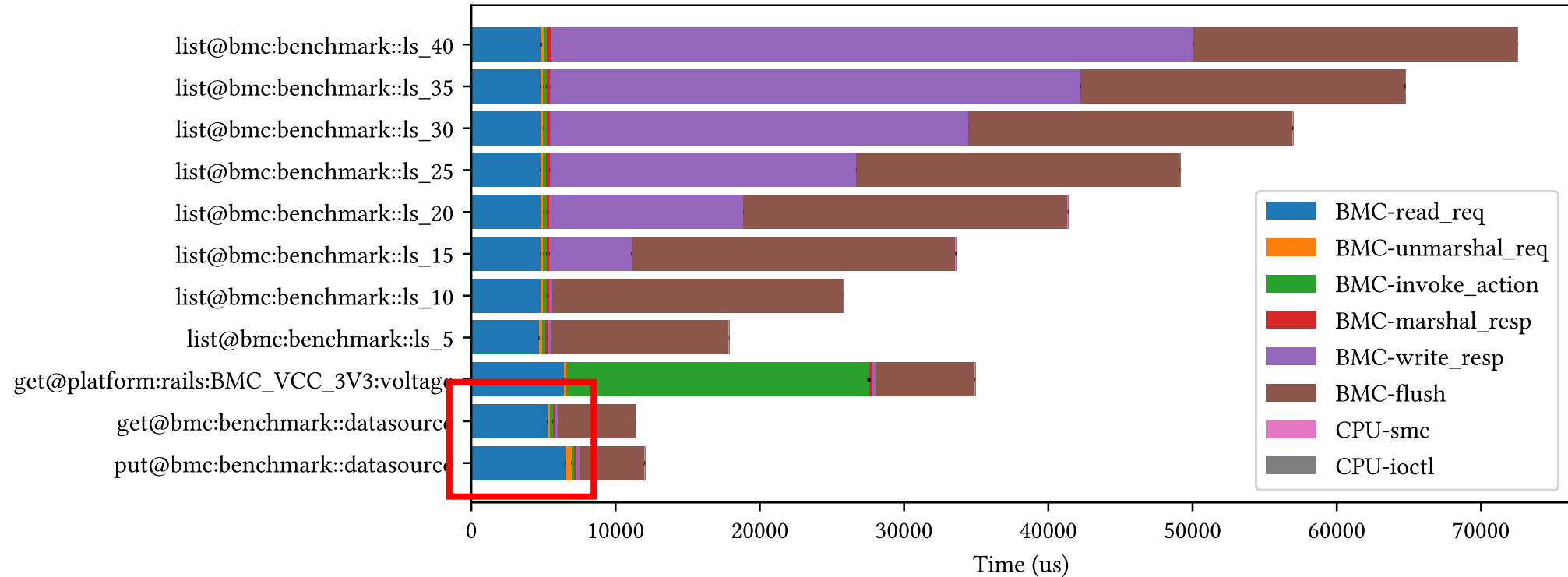
Performance – latency breakdown



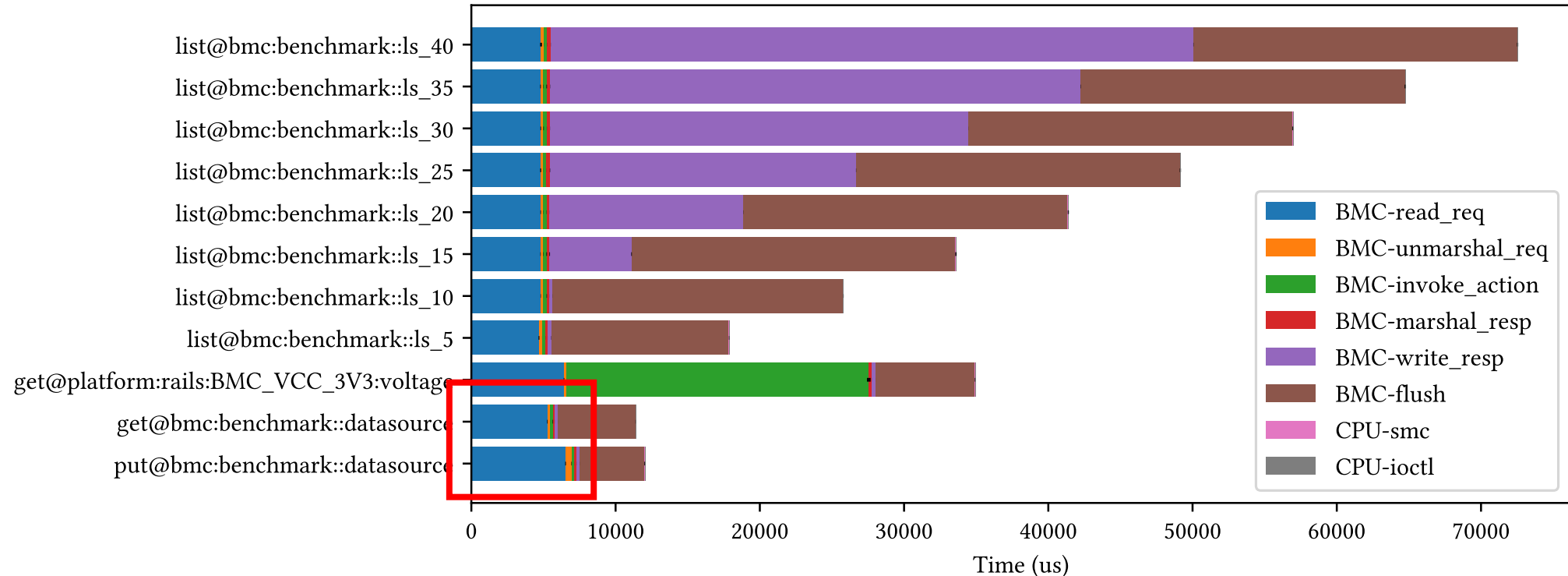
Performance – latency breakdown



Performance – latency breakdown



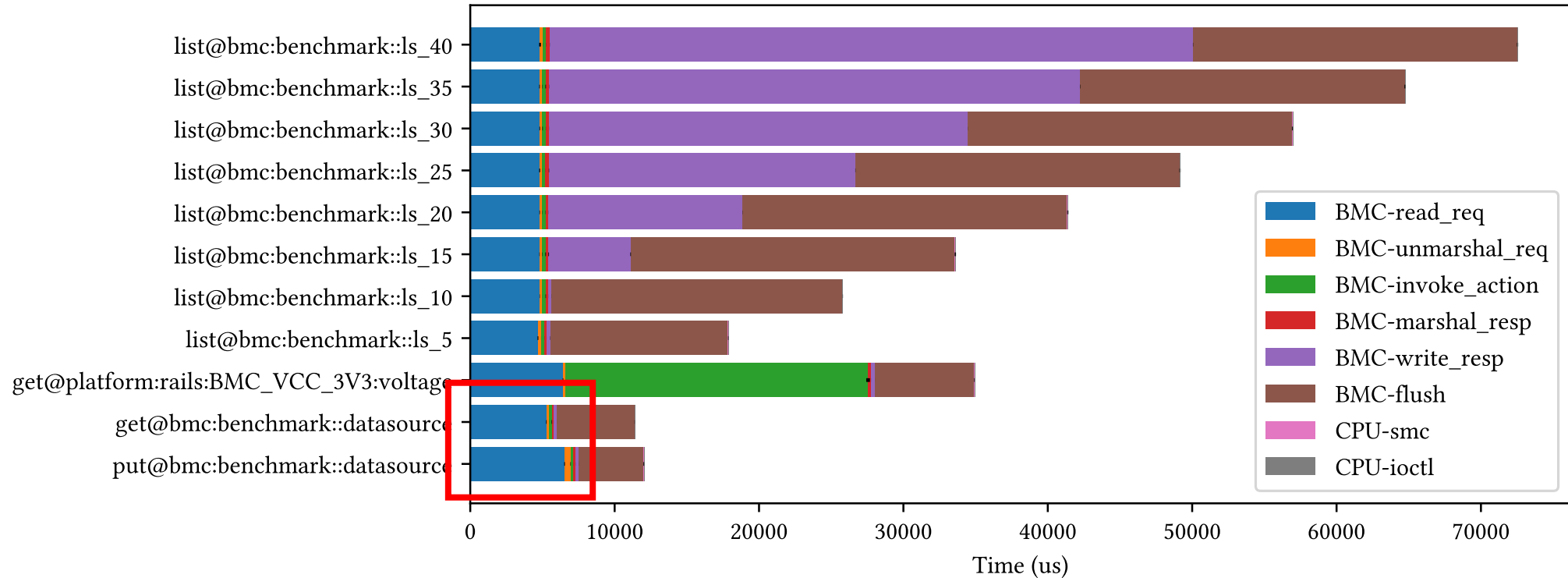
Performance – latency breakdown



==EFRI== => bmc:benchmark::datasource get\n

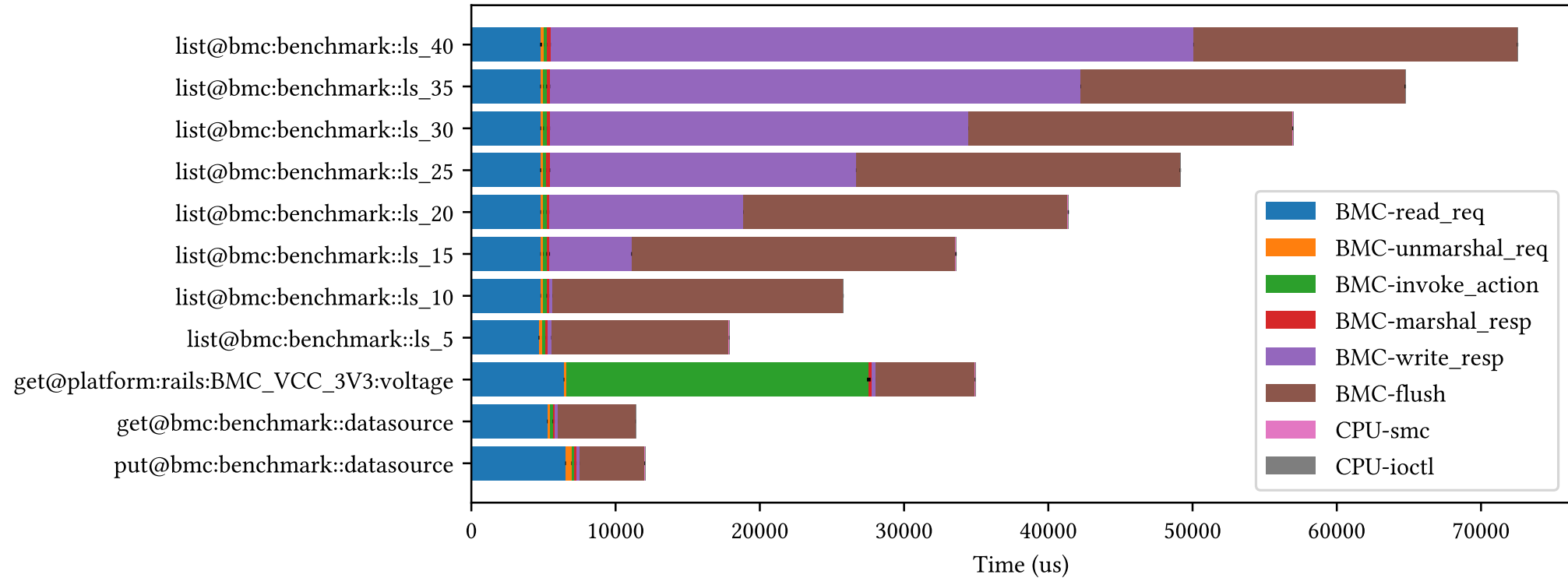
==EFRI== => bmc:benchmark::datasource put Rpm\$5#RPM\n

Performance – latency breakdown

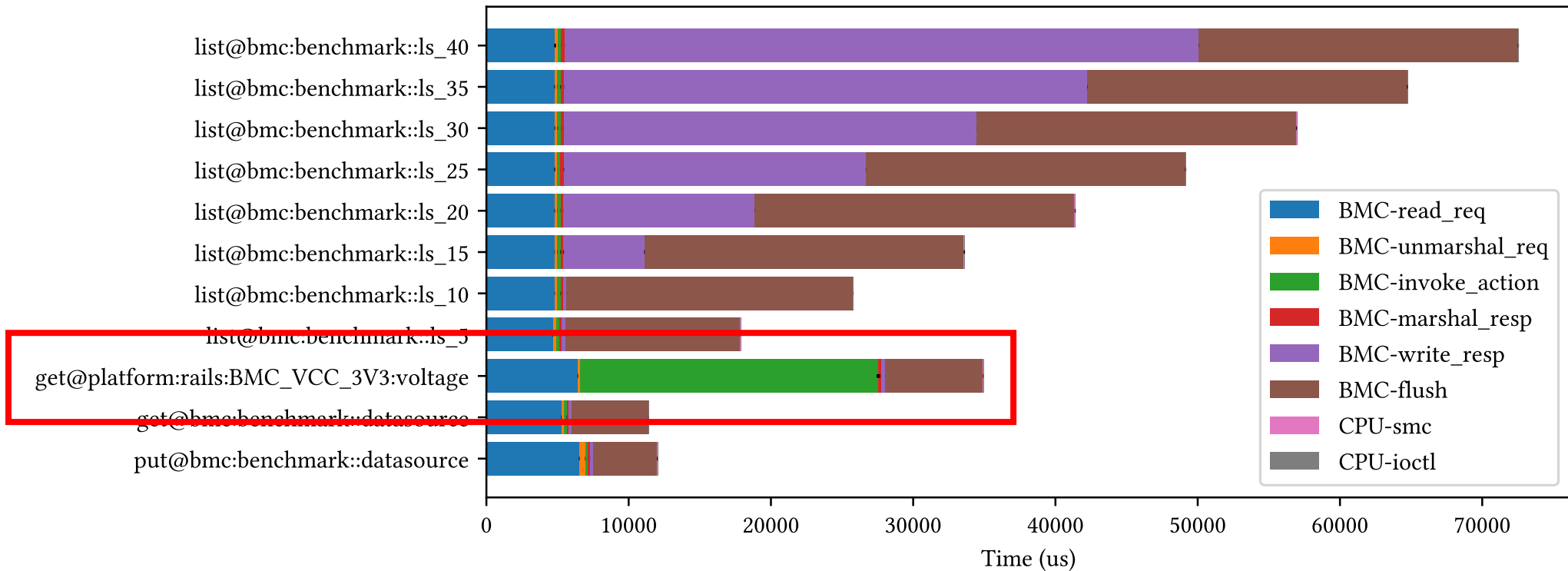


==EFRI== => bmc:benchmark::datasource get\n 42 vs 52 bytes
 ==EFRI== => bmc:benchmark::datasource put Rpm\$5#RPM\n 120us/byte (UART: 86us/byte)

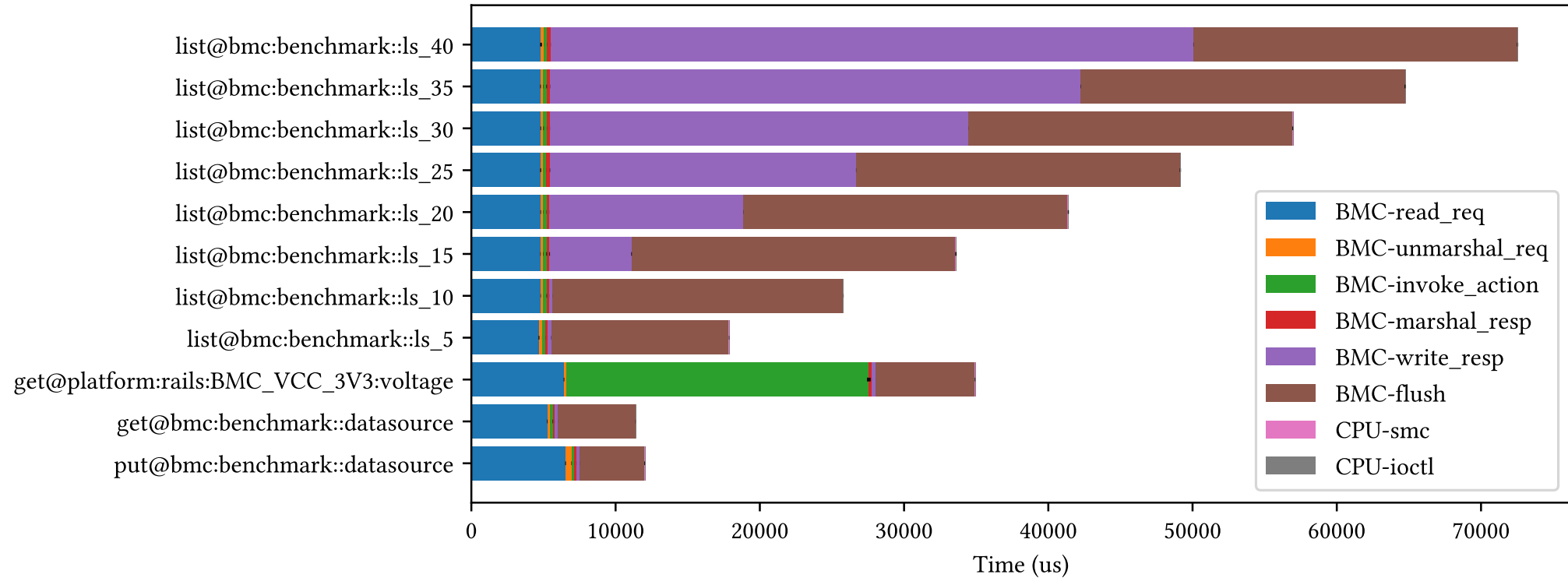
Performance – latency breakdown



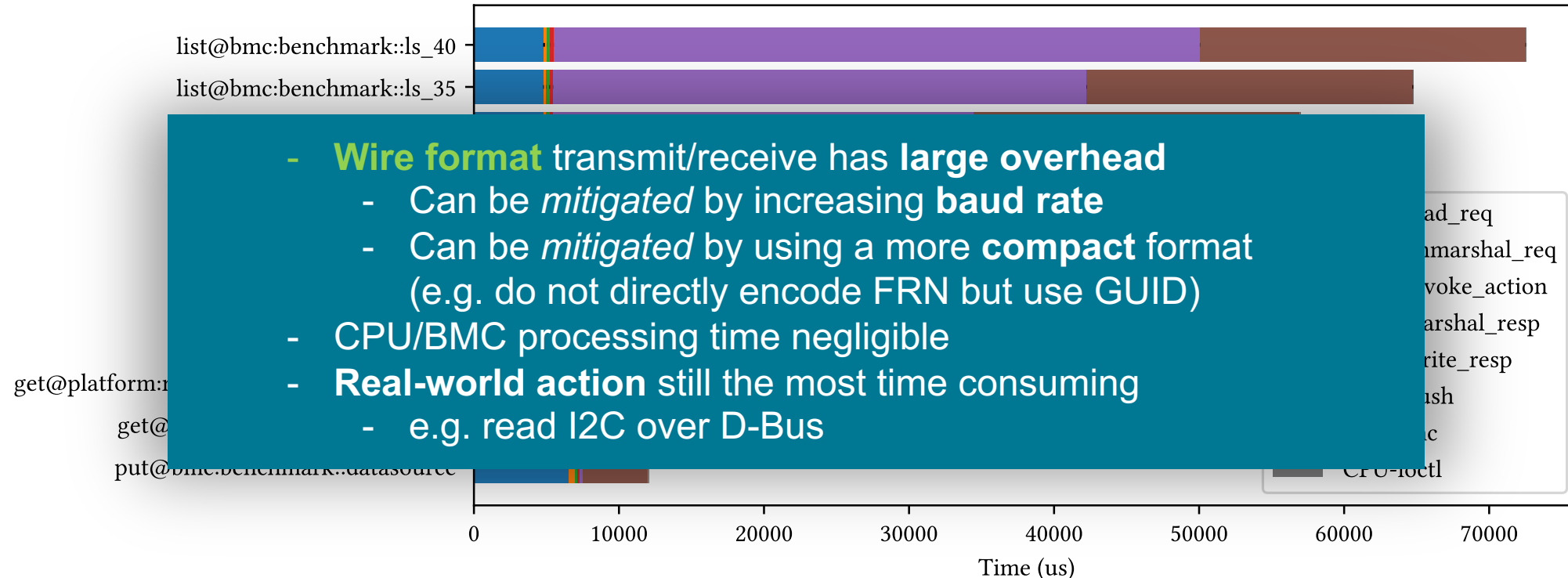
Performance – latency breakdown



Performance – latency breakdown



Performance – latency breakdown



Case studies for ease of extension



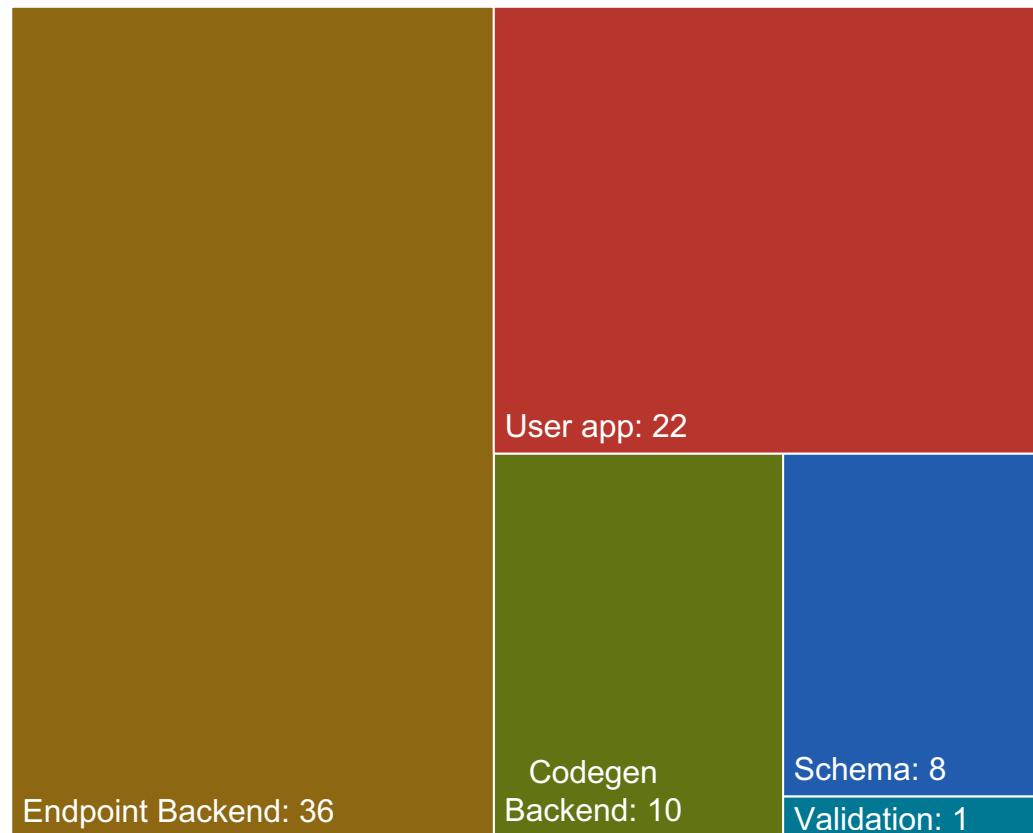
DINFK

- Measure LOC changes for implementing a specific firmware feature

Case studies for ease of extension

- Measure LOC changes for implementing a specific firmware feature

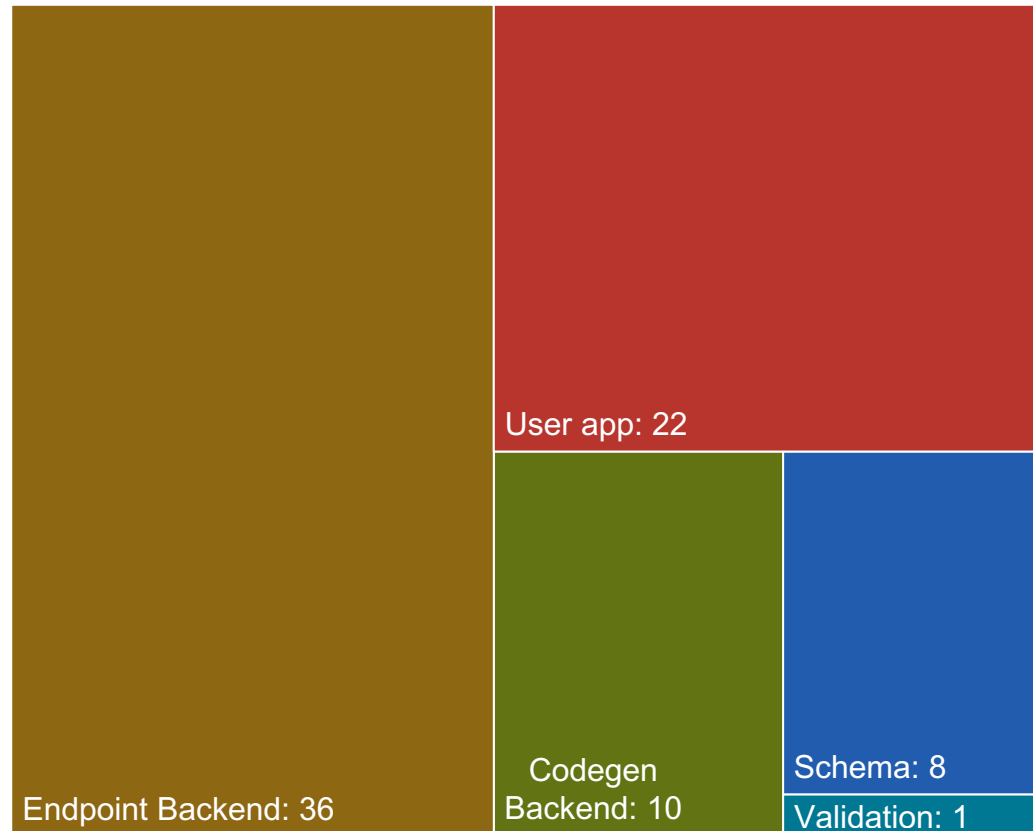
CPU power-down



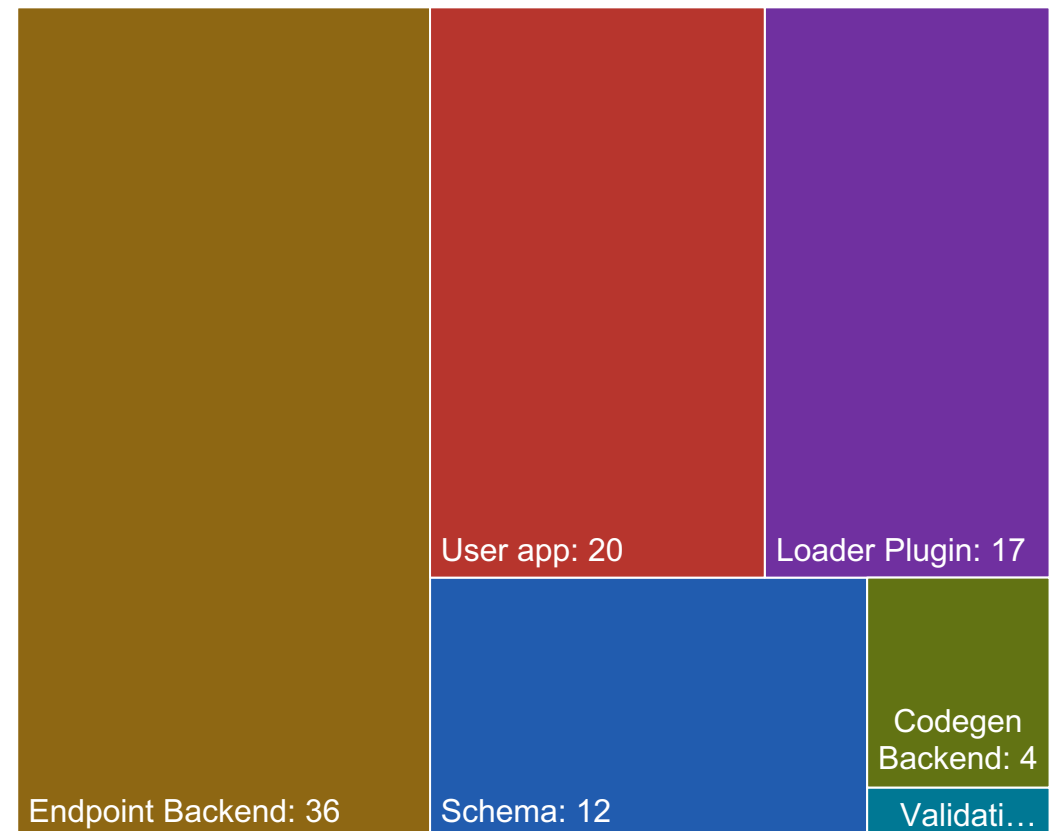
Case studies for ease of extension

- Measure LOC changes for implementing a specific firmware feature

CPU power-down



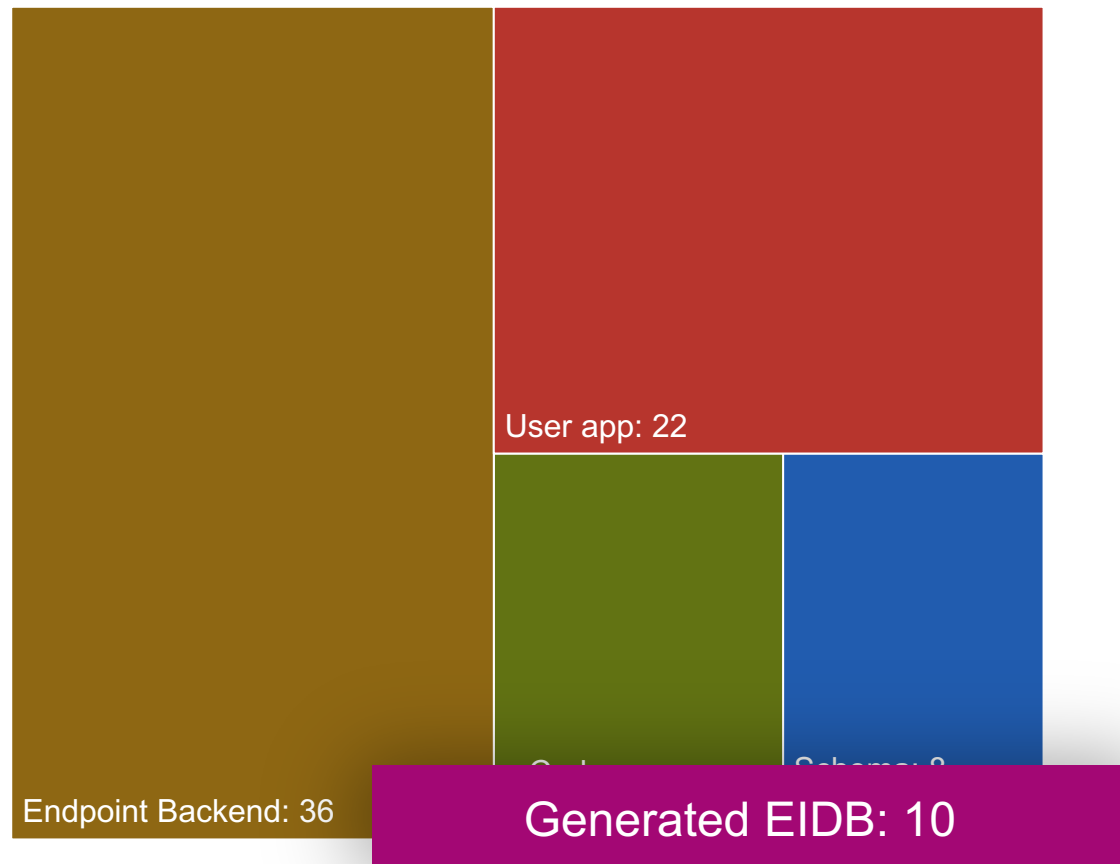
Power rail telemetry



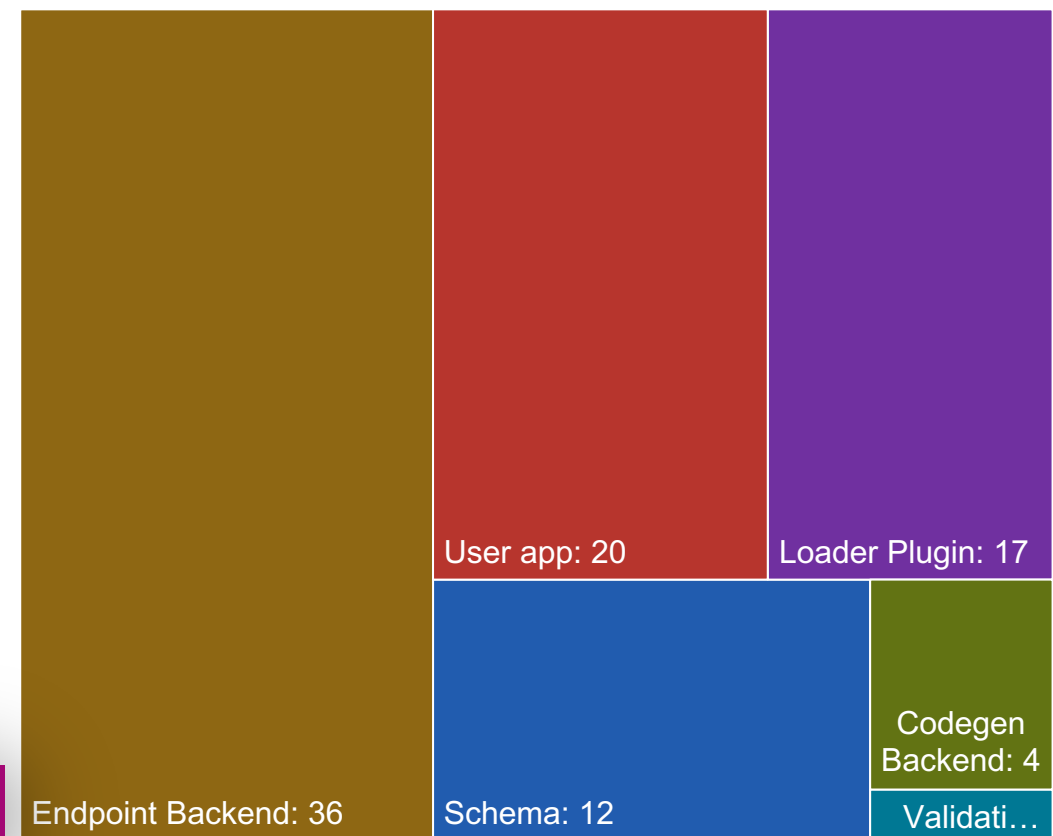
Case studies for ease of extension

- Measure LOC changes for implementing a specific firmware feature

CPU power-down



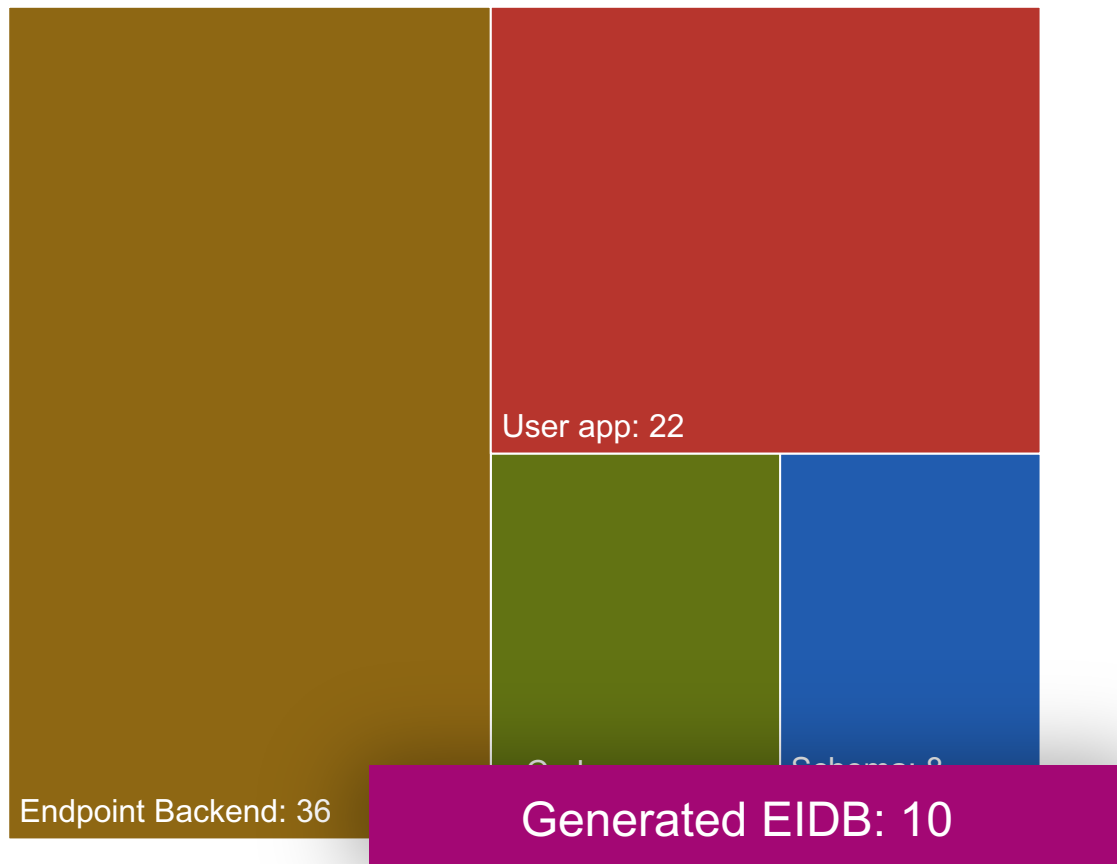
Power rail telemetry



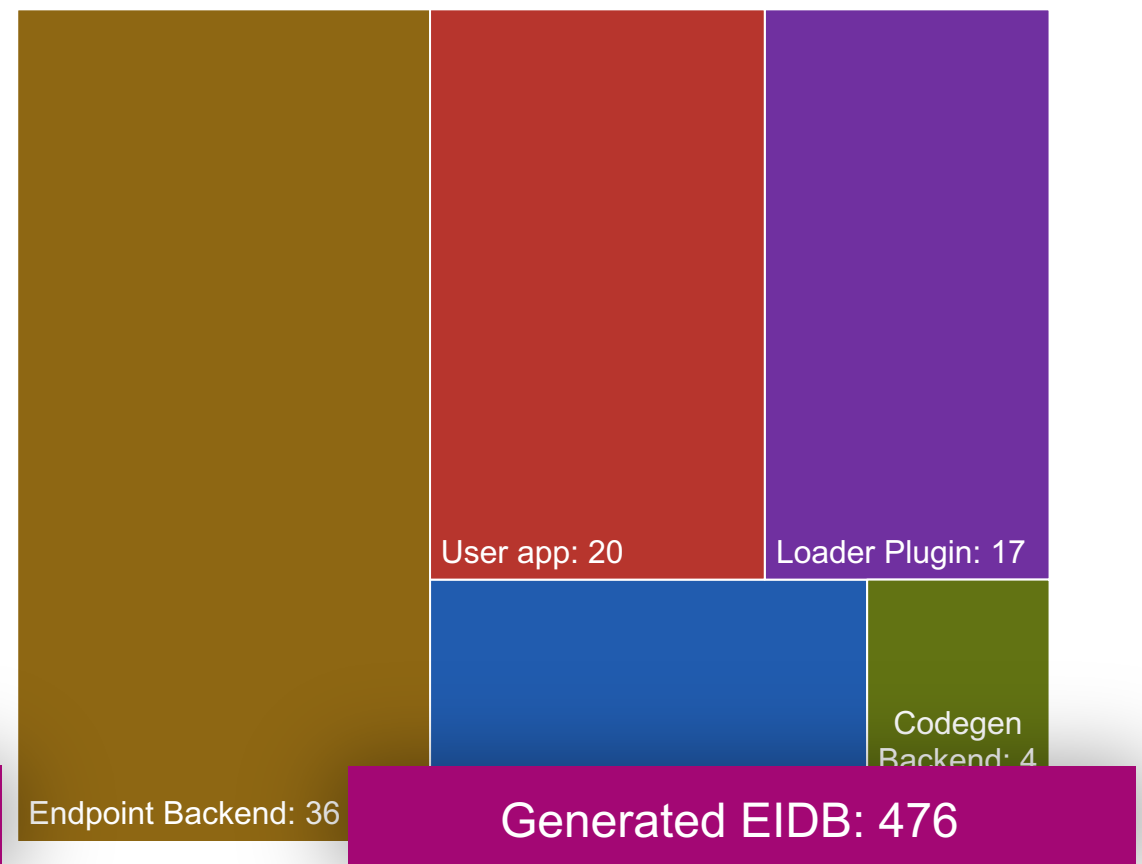
Case studies for ease of extension

- Measure LOC changes for implementing a specific firmware feature

CPU power-down

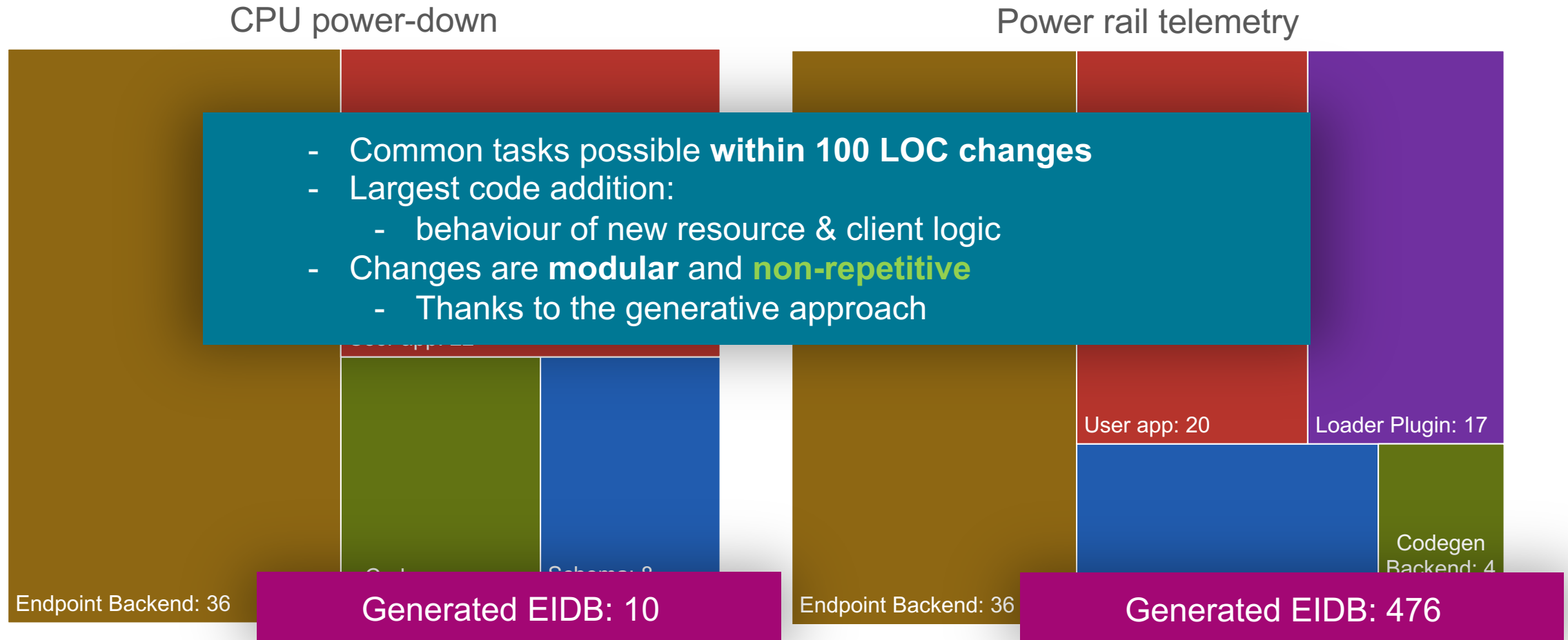


Power rail telemetry

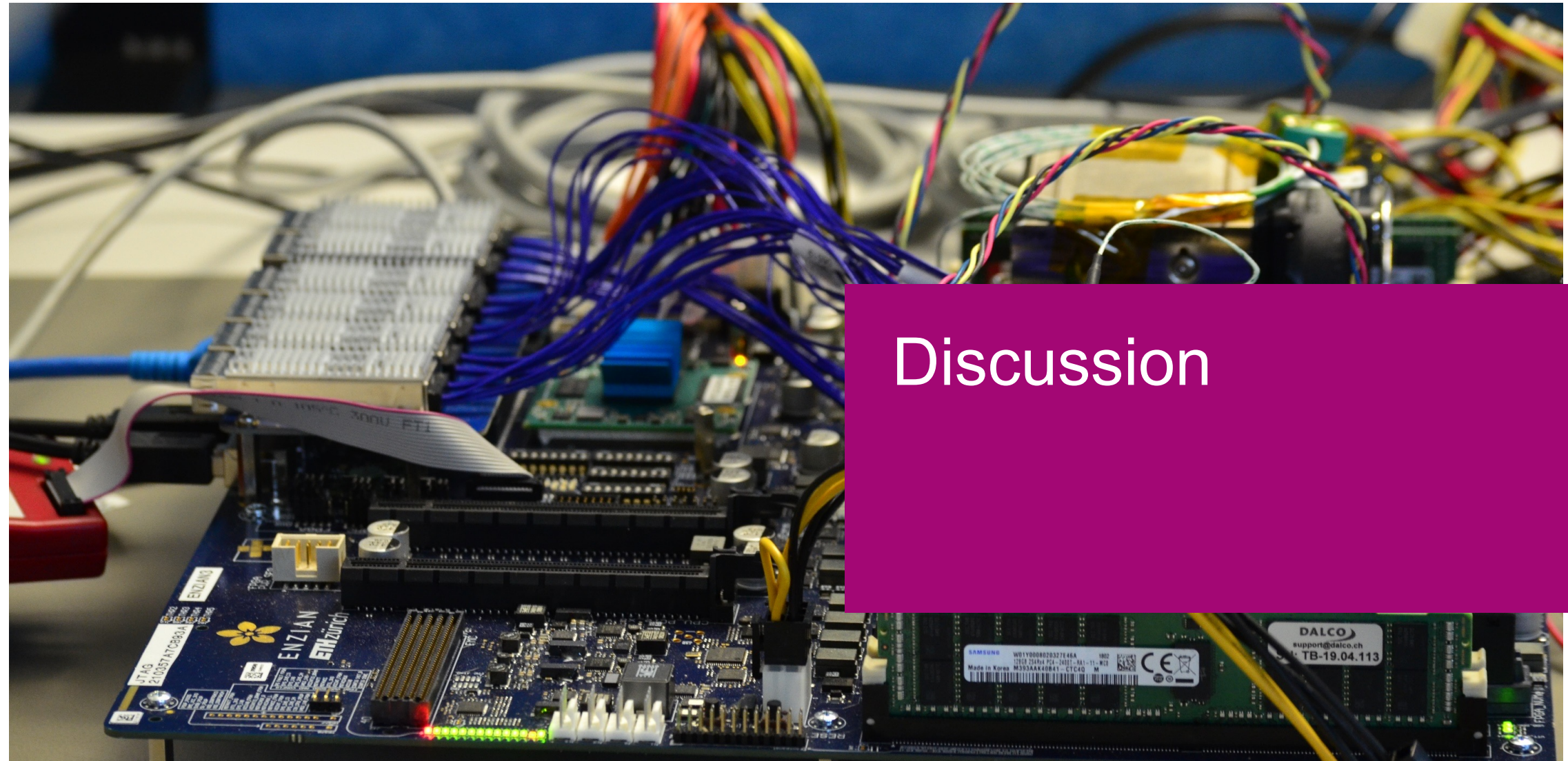


Case studies for ease of extension

- Measure LOC changes for implementing a specific firmware feature

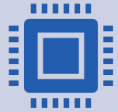


Live Demo



Discussion

- Better CPU-side API & libraries, language support, etc.
 - Upper EFRI kernel module now exposes argspage to userspace
- Better integration of benchmark framework into base protocol
 - For more predictable time estimation *in-vivo*
 - Important for cross-source profiling
- Asynchronous Event Mechanism implementation
- Security Model implementation
- Explore FPGA-side operation with EnzianShell
- Expand the use-case horizon: implement more services
- ...



A *better* firmware protocol for heterogeneous systems

Industrial standards aren't *bad*, they are designed **for a different problem**



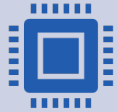
EFRI is...

Practical to implement
Easy to extend
Acceptably performant



A future where...

EFRI is the *common substrate* of heterogeneous platforms and applications



A *better* firmware protocol for heterogeneous systems

Industrial standards aren't *bad*, they are designed **for a different problem**



EFRI is...

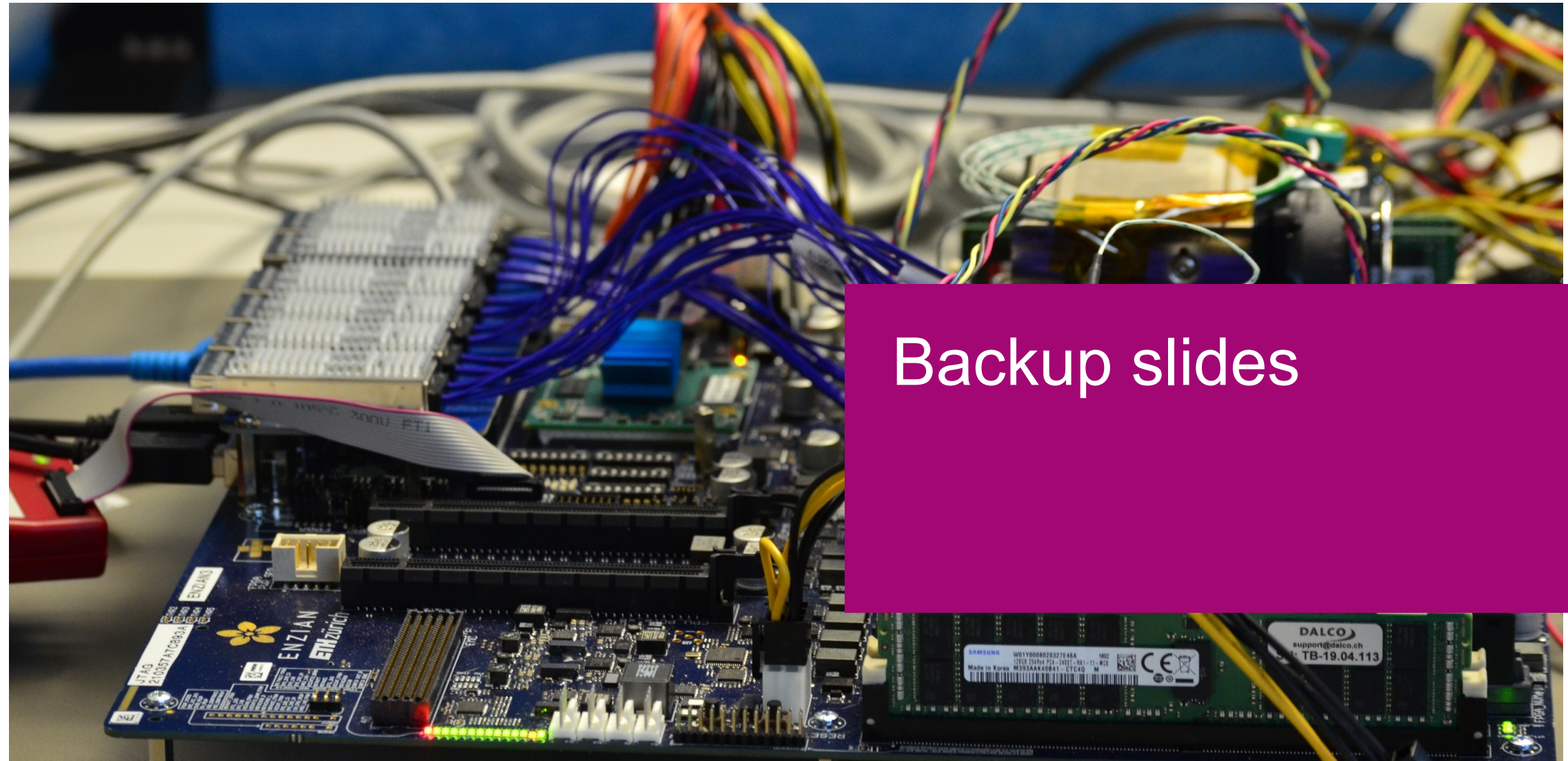
Practical to implement
Easy to extend
Acceptably performant



A future where...

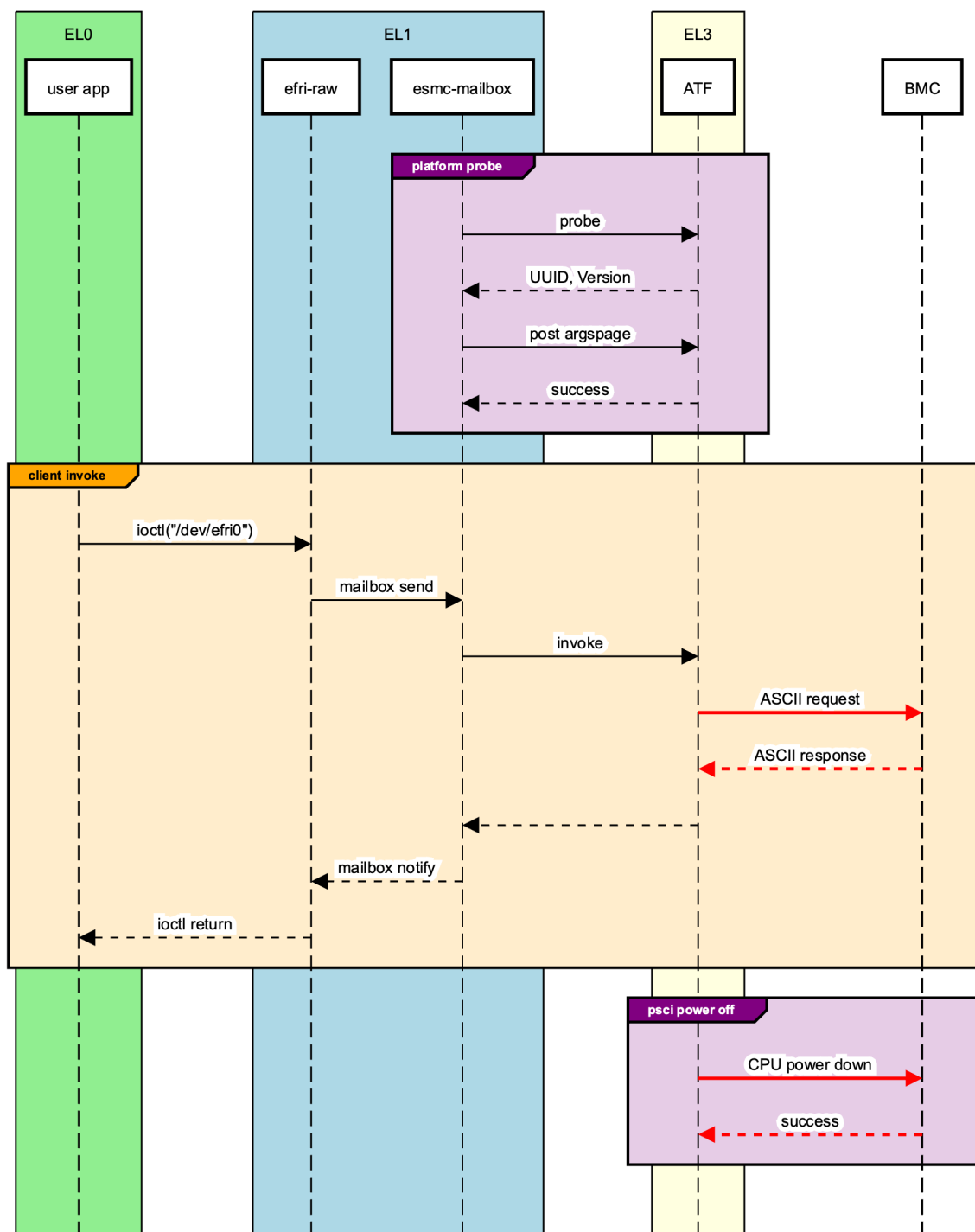
EFRI is the *common substrate* of heterogeneous platforms and applications

Questions?



Backup slides

Global sequence diagram for client invocations



Schema definitions

```
1  actions:
2      put: &put
3          name: put
4          args: [p]
5          retval: [err]
6      get: &get
7          name: get
8          args: []
9          retval: [err, T]
10 aliases:
11     rw: &rw [*put, *get]
12     ro: &ro [*get]
```

```
power_rail: &power_rail
  actions: *ro
  impl: datasource-dbus-power
  class: [platform, rails]
  subsystem: telemetry
  rail_voltage: &rail_voltage
  name: voltage
  type: voltage
  <<: *power_rail
  rail_current: &rail_current
  name: current
  type: current
  <<: *power_rail
```

Codegen backend

```
1 def gen_ep(ep: Endpoint):
2     ...
3     if ep.impl.startswith('datasource'):
4         ds_common = f'"{frn}", {ep.type.name.capitalize()}, {actions}'
5         if 'dummy' in ep.impl:
6             ep_impl = f'DummyDataSource({ds_common}, readonly={"ro" in
7                 ↪ ep.impl})'
8         elif 'dbus-power' in ep.impl:
9             # ep should have device and monitor properties set
10            assert hasattr(ep, 'device') and hasattr(ep, 'monitor'), \
11                'datasource-dbus-power requires device and monitor set'
12            ep_impl = f'DBusPowerDataSource({ds_common}, "{ep.device}",
13                ↪ "{ep.monitor}")'
14        else:
15            assert False, 'unknown datasource impl ' + ep.impl
16    ...
```


Endpoint implementation database

```
1 FRN_EP_MAP["platform:rails:VDD_DDRCPU13:voltage"] =  
  ↪ DBusPowerDataSource("platform:rails:VDD_DDRCPU13:voltage", Voltage, [  
2     ActionDesc("get", 0, 2),  
3     ActionDesc("subscribe", 2, 1)  
4 ], "pac_cpu", "VMON9")  
5 FRN_EP_MAP["platform:rails:VDD_DDRCPU13:current"] =  
  ↪ DBusPowerDataSource("platform:rails:VDD_DDRCPU13:current", Current, [  
6     ActionDesc("get", 0, 2),  
7     ActionDesc("subscribe", 2, 1)  
8 ], "ina226_dds_cpu_13", "CURRENT")
```


Endpoint backend

```
1 E_DBUS = TypedValue(Error, 6) # EFRI_INTERNAL_ERROR
2 class DBusPowerDataSource(DataSource):
3     @classmethod
4     def _get_service(cls):
5         return dbus.SystemBus().get_object('systems.enzian.Power',
6         ↪ '/systems/enzian/Power')
7     def get(self):
8         try:
9             raw_val = self._get_service().read_device_monitor(self.device,
10            ↪ self.monitor)
11         except dbus.DBusException:
12             return E_DBUS, None
13         return E_OK, TypedValue(self.ty, raw_val)
```

Userspace application

```
1  #include "libefri.h"
2  int main(int argc, char *argv[]) {
3      int fd = open("/dev/efri_raw0", O_RDWR);
4
5      // fill argspage
6      struct efri_raw_invoke_msg msg;
7      efri_frame_t *frame = &msg.frame.frame;
8      char *cur = msg.frame.aux;
9      frame->frn = push_string(msg.buf, &cur,
↪  "platform:rails:BMC_VCC_3V3:voltage");
10     frame->action = push_string(msg.buf, &cur, "get");
11     frame->role = ROLE_REQUEST;
12     frame->num_args = 0;
13     msg.dest = EFRI_DEST_BMC;
14     msg.size = cur - msg.buf;
15
16     // invoke and print result
17     ioctl(fd, EFRI_RAW_INVOKE, &msg);
18     describe_frame(frame);
19 }
```