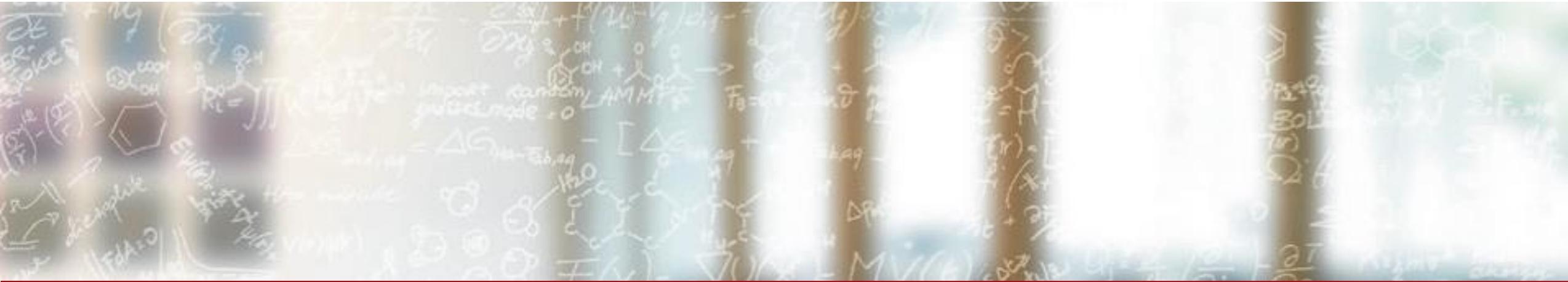




CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich



A glimpse of YAUULT

Victor Holanda Rusu and Massimo Benini, CSCS

May 5th, 2025

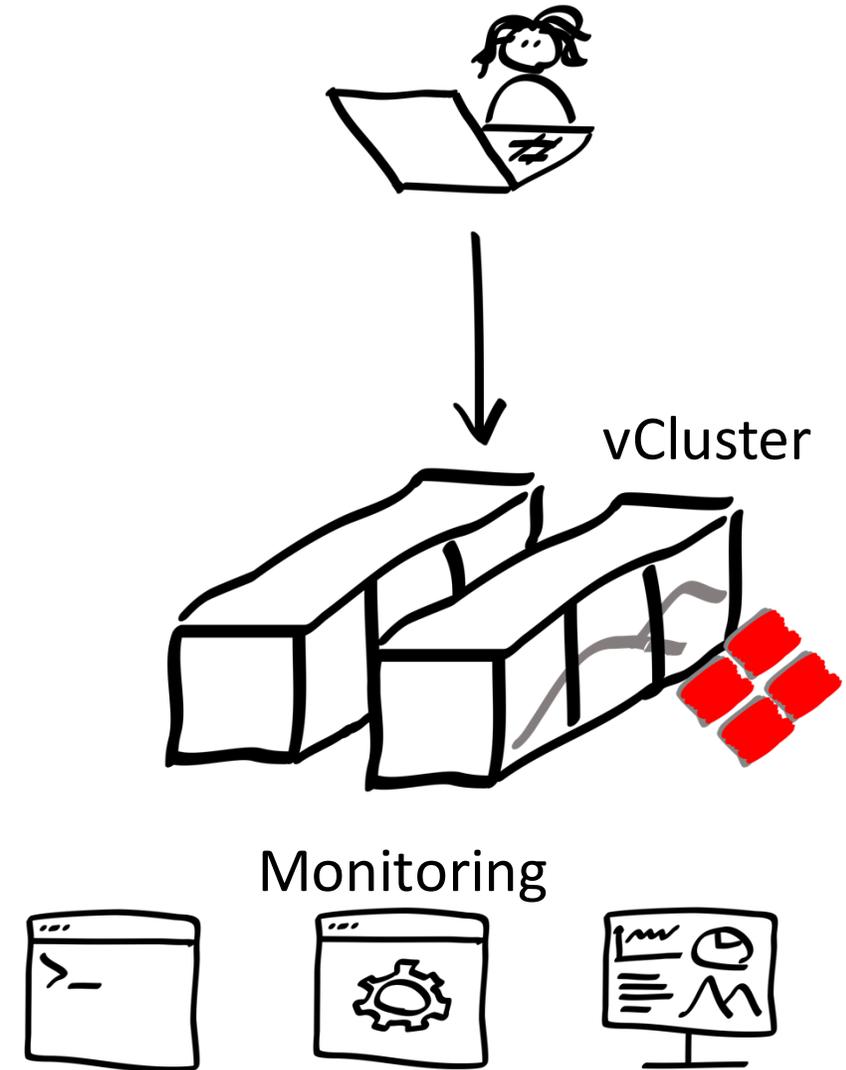
Computing Horizons: CUG 2025

Beyond a System Monitoring

50'000 Feet View Above Our Monitoring

- Infrastructure monitoring
 - Networking
 - Storage
 - Nodes
- Slurm job
 - Jobcomp
 - Slurm DB
- Specific user commands execution
 - auditd
 - syscalls
 - sudo

Can we do more?



Beyond a System Monitoring

50'000 Feet View Above User Application Monitoring

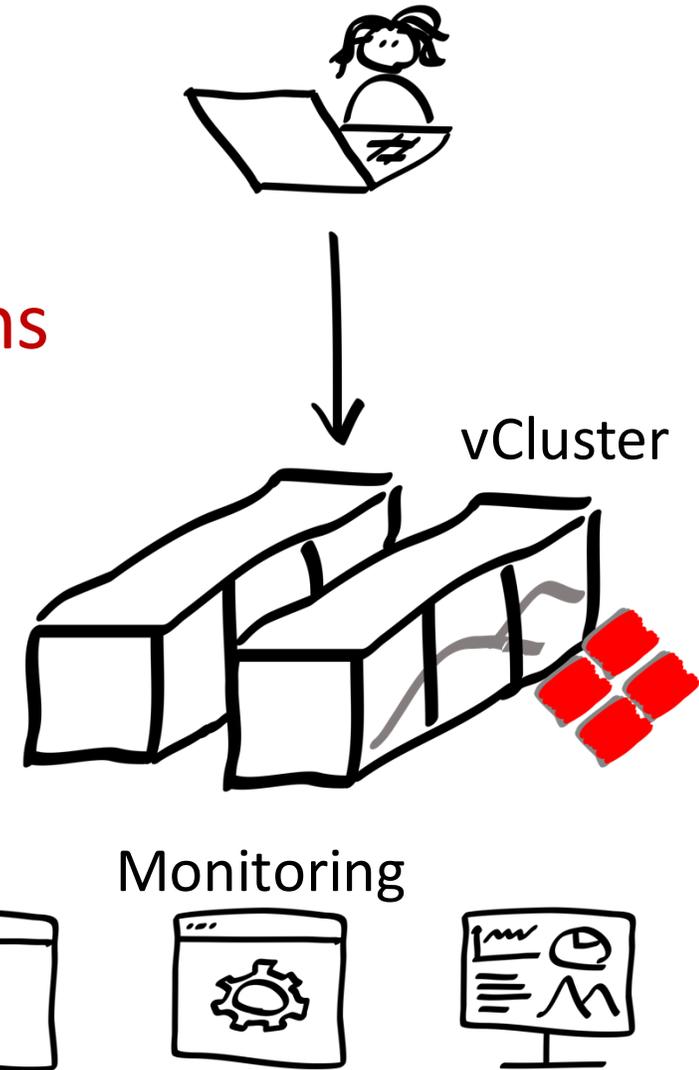
- Kernel namespace or username monitoring
- Usernamespace monitoring generally uses ld preload techniques
 - ELF tricks
 - Able to execute programs as the user
- Kernel namespace monitoring
 - Daemons reading information provided by the kernel
 - Custom kernel modules
 - Application or Kernel tracing
 - EBPF (since 2015)
- Limited by the environment
 - Mandatory Access control
 - Chrooted (container)
 - User owns their own environment
- Limited by the Kernel context

Beyond a System Monitoring

50'000 Feet View Above Our Monitoring

- Infrastructure monitoring
 - Networking
 - Storage
 - Nodes
- Slurm job
 - Jobcomp
 - Slurm DB
- Specific User commands execution
 - auditd
 - syscalls
 - sudo
 - eBPF

Can we identify the applications used by the users?



What's eBPF?

TL;DR Edition

- It is a technology that allows programs to run code within a virtual machine inside the Linux kernel
- No need to add additional modules or modify the kernel source code
- Sandboxed environment
- Limited C-based DSL
 - No memory allocation
 - No infinite loops
 - No sleeps
- Library interface for several languages, e.g. Python, Go, C/C++, etc.
- Write modules using a AWK-like language (bpftrace)

What's eBPF?

How does the code look like?

```
1 #ifndef BPFTRACE_HAVE_BTFF
2 #include <linux/sched.h>
3 #endif
4
5 BEGIN
6 {
7     printf("%-15s %-7s %-7s %s\n", "TIME", "PID", "PPID", "ARGS");
8 }
9
10 tracepoint:syscalls:sys_enter_exec*
11 {
12     $task = (struct task_struct *)curtask;
13     printf("%15s %-7d %-7d ", strftime("%H:%M:%S.%f", nsecs), pid, $task->real_parent->pid);
14     join(args.argv);
15 }
```

```
1 from bcc import BPF, USDT
2 import sys
3
4 bpf_text = """
5 #include <uapi/linux/ptrace.h>
6
7 struct data_t {
8     u64 timestamp;
9     char client_ip;
10    char file_path;
```

```
    }, char *ipstrptr, char *pathptr, u32 file_size) {
```

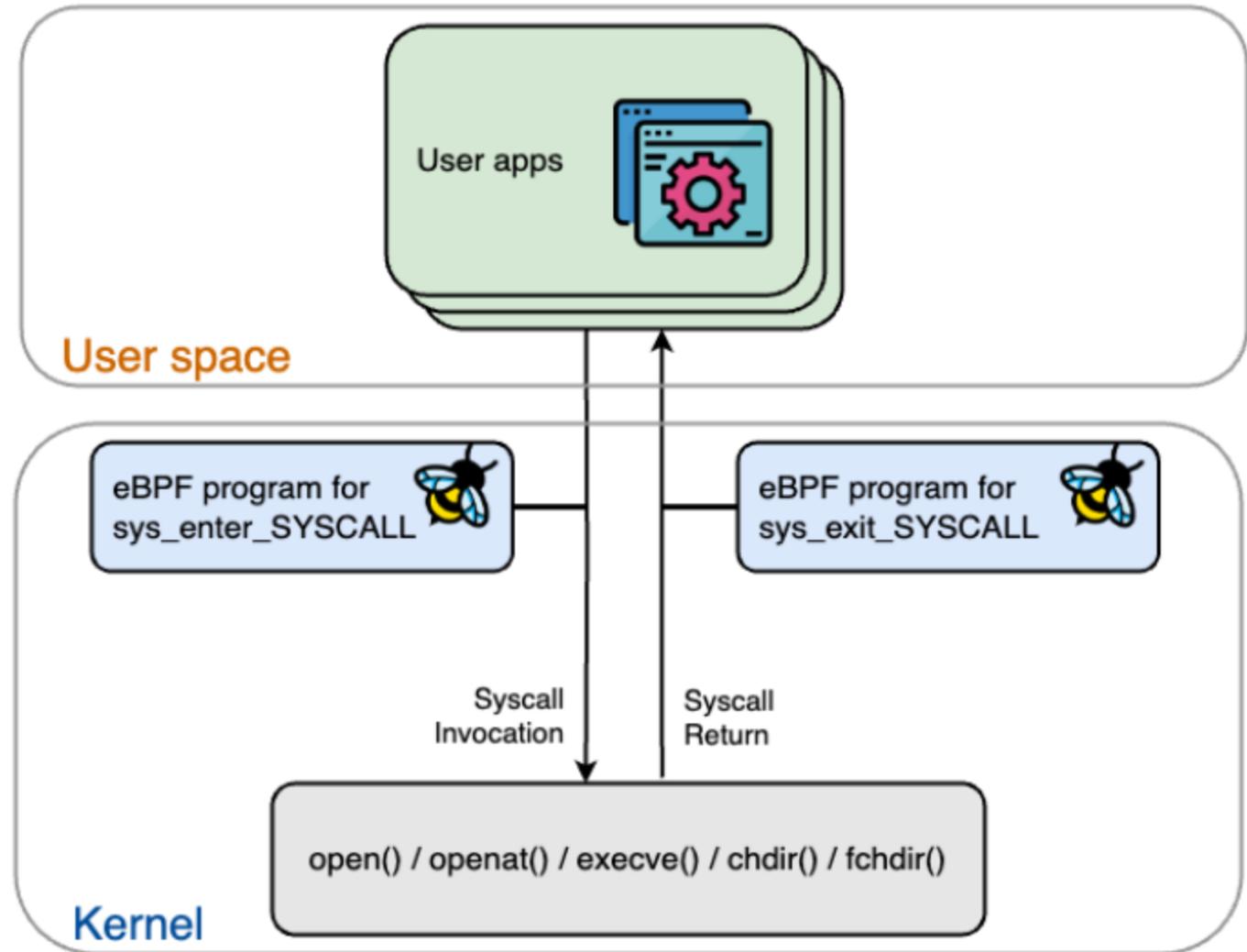
```
        lata.client_ip), (void *)ipstrptr);
        lata.file_path), (void *)pathptr);
        ita));
```

```
28 u.enable_probe(probe="file_transfer", fn_name="trace_file_transfers")
29 b = BPF(text=bpf_text, usdt_contexts=[u])
30
31 def print_event(cpu, data, size):
32     event = b["events"].event(data)
33     print("{0}: {1} is downloading file {2} ({3} bytes)".format(
34         event.timestamp, event.client_ip, event.file_path, event.file_size))
35
36 b["events"].open_perf_buffer(print_event)
37 while True:
38     try:
39         b.perf_buffer_poll()
40     except KeyboardInterrupt:
41         exit()
```

What's YAULT?

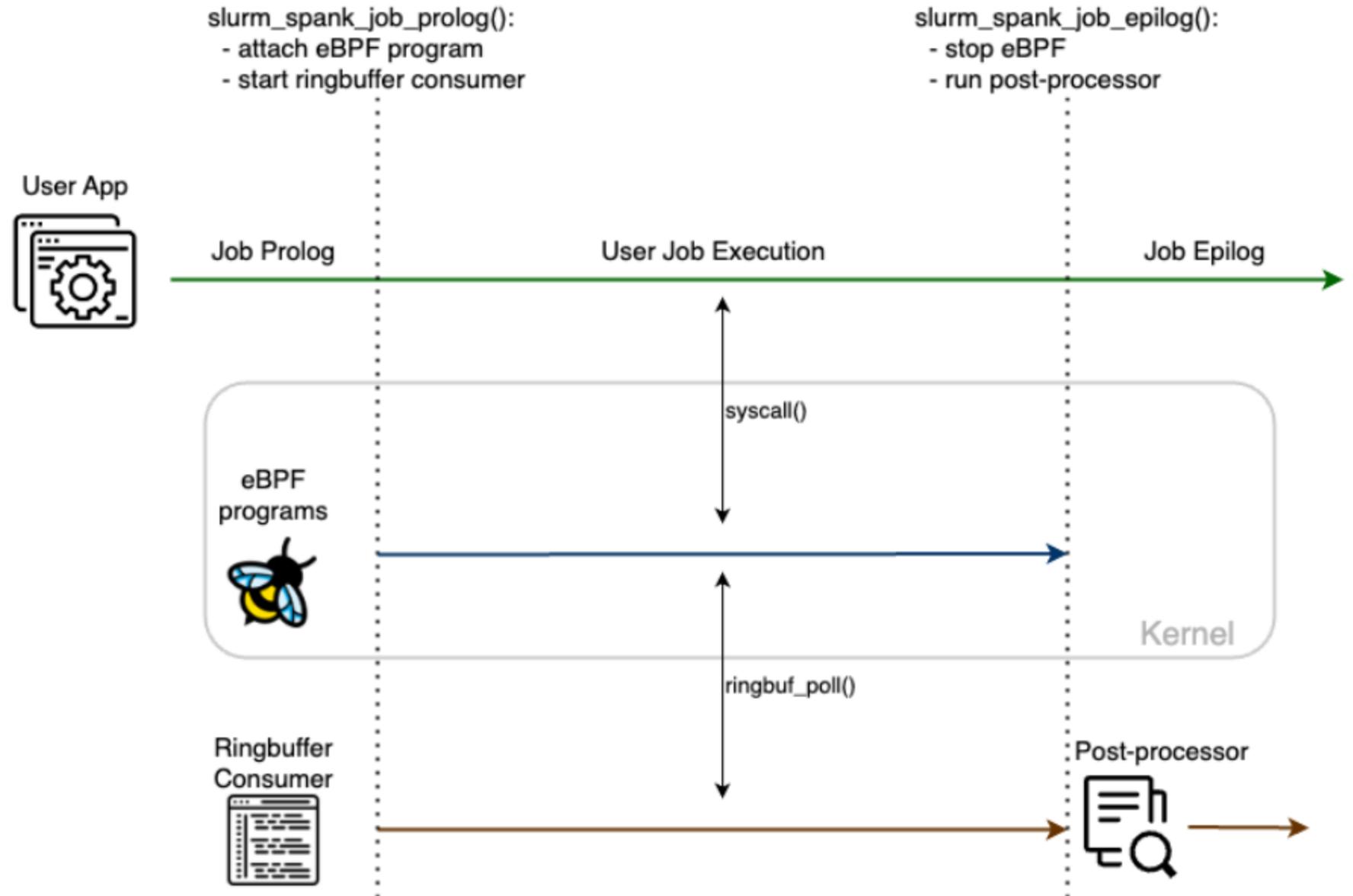
What are we thinking?

- We use eBPF to trace user application syscalls and their arguments
- Slurm plugin to associate application and Slurm job context
- Use this information to identify user application



What's YAULT?

What are we thinking?

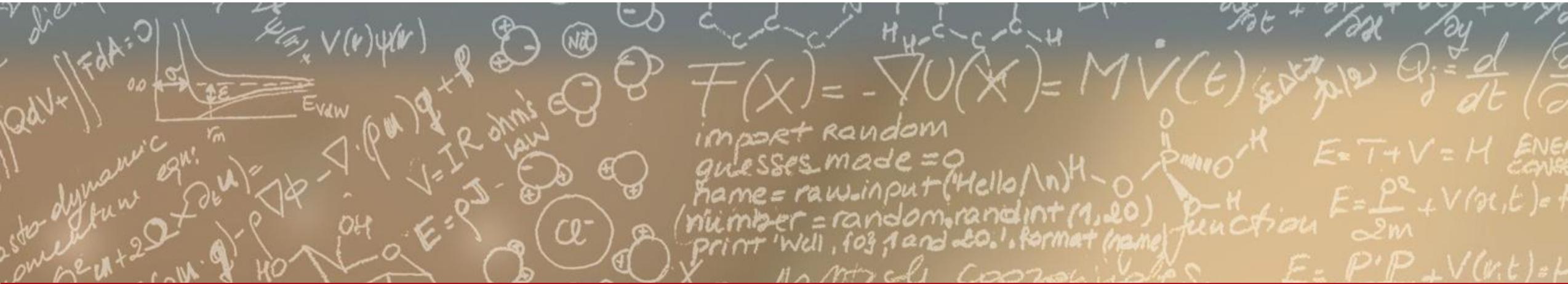




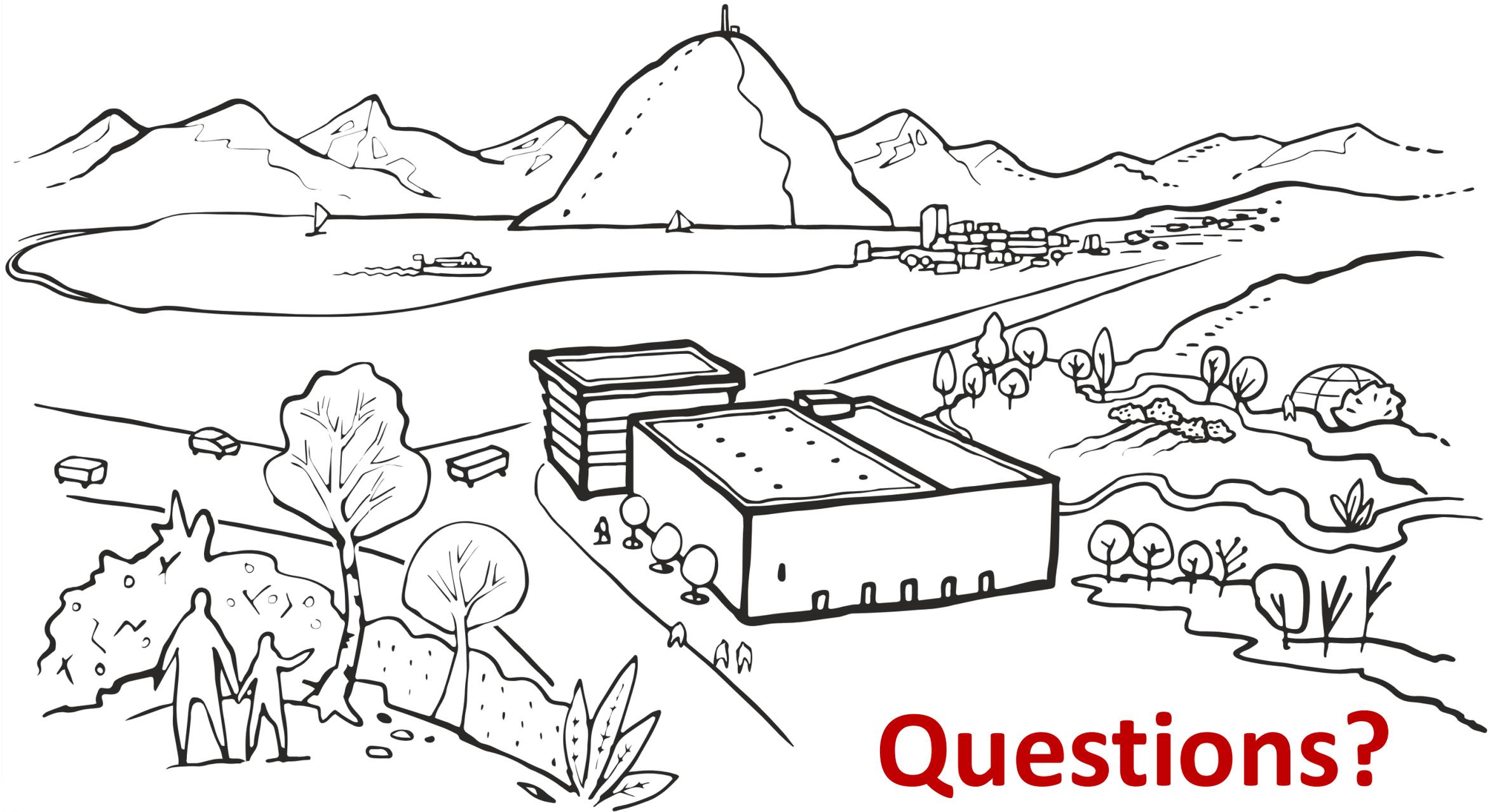
CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich



Thank you for your attention.



Questions?