# Extending libbpf for Kubernetes
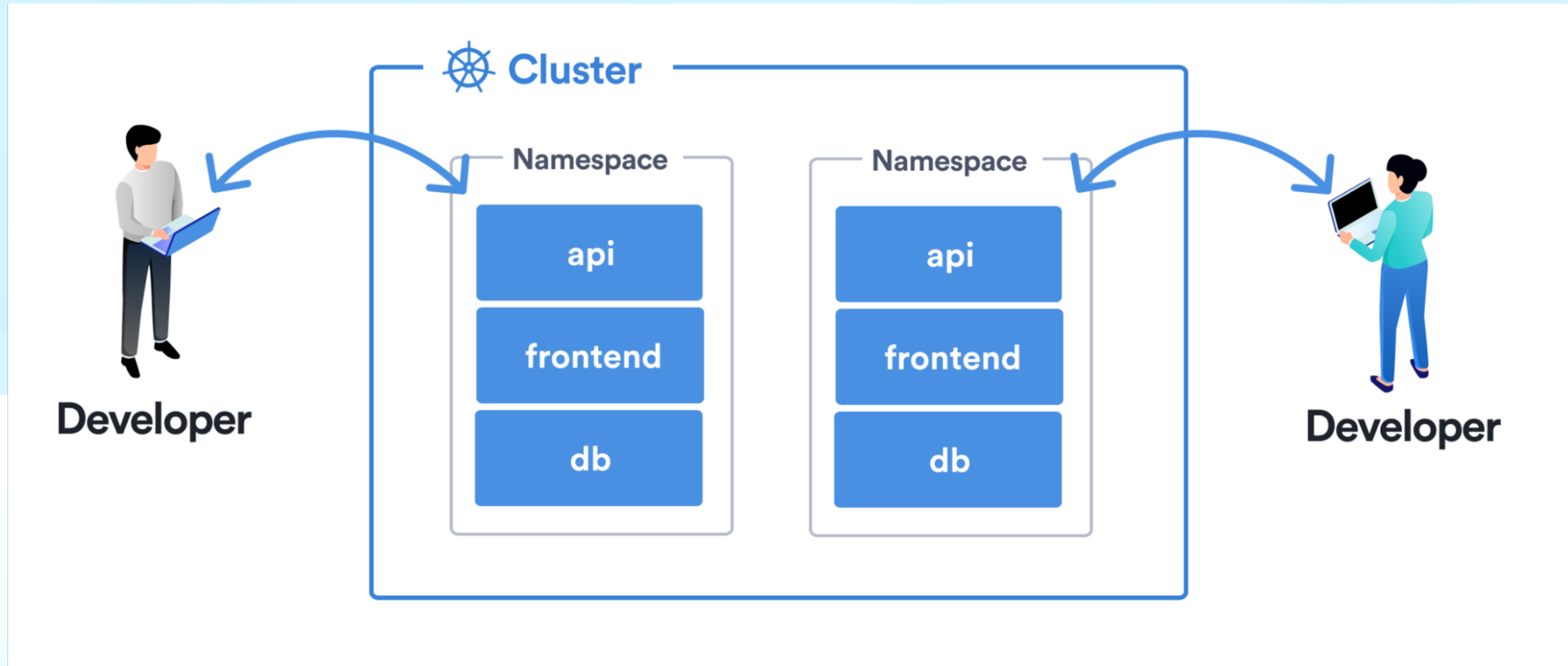
**Vinay Kulkarni**

*@iSkibum*

# Agenda

- eBPF to Proactively Do Stuff

    - Use Cases

- Problem at Hand …

- … How eBPF Came to the Rescue

- Pain Points

- Potential Solution

- Summary & Discussion

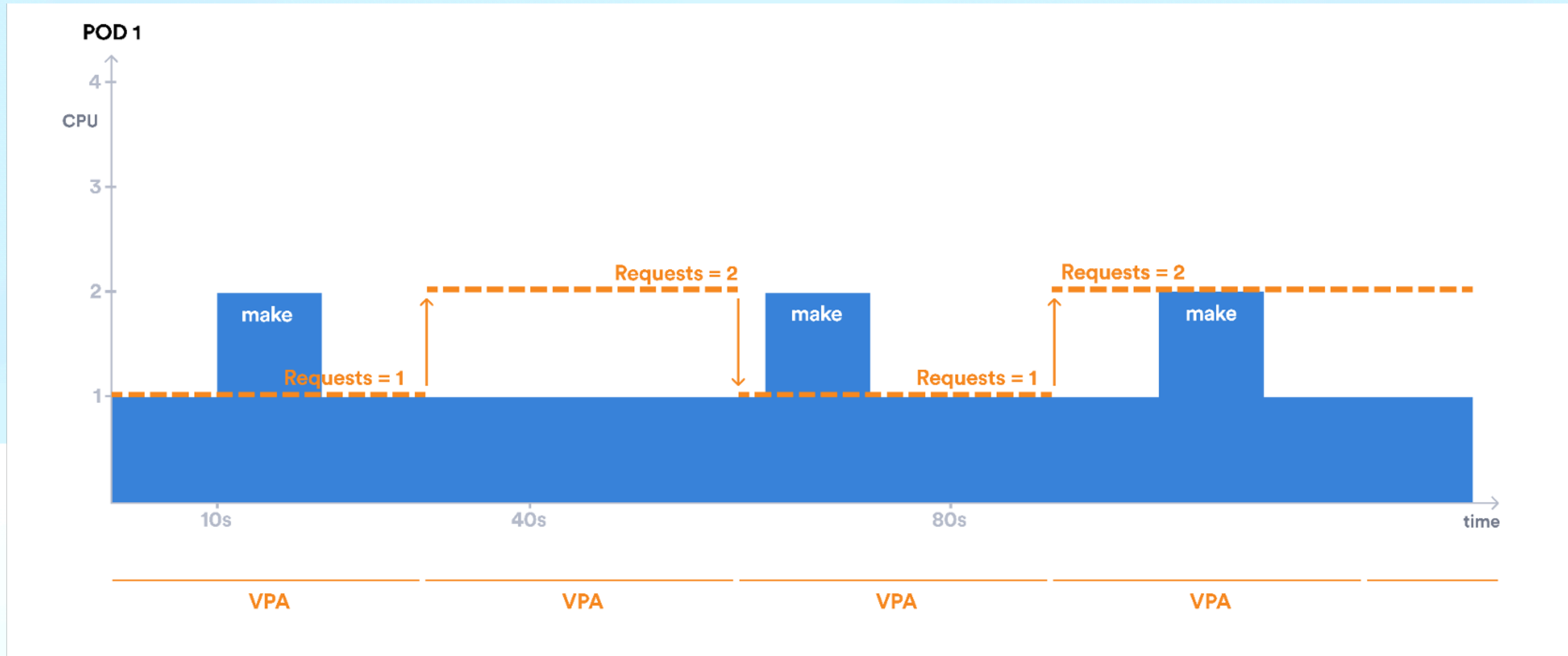# eBPF + Proactive: Use Cases

- Cloud-based Development Environment

    - Developers / dev teams write code with local IDE

    - Build code or run a battery of tests is Kubernetes pods

    - <u>Need</u>: Resize up build/test pod quickly

- Containerized Java application with high startup CPU requirements

    - Running time CPU usage is 1/10th the startup time CPU needs

        - Allocating too little CPU -> long startup time

        - Allocating startup requirements -> underutilized cluster

        - <u>Need</u>: Resize down pod quickly after startup

- eBPF network stats program attachment to pod veth
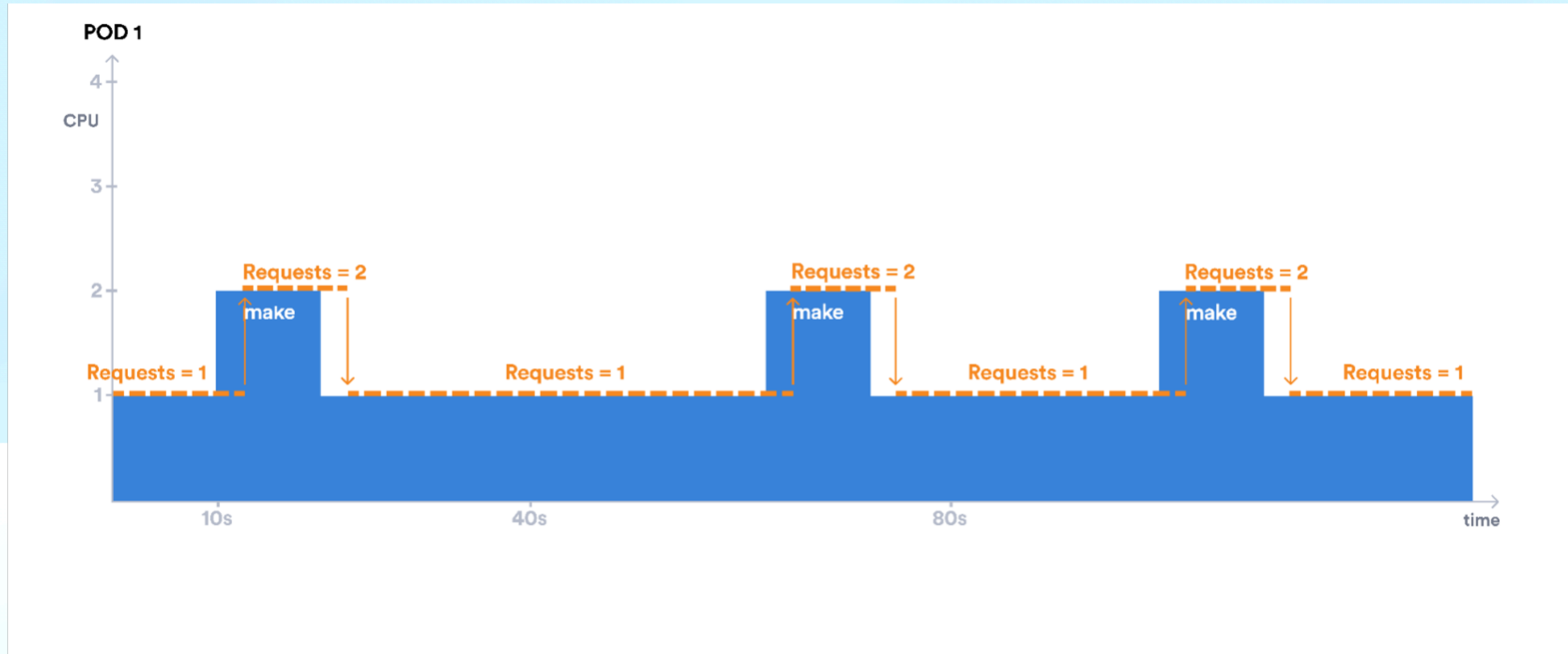
# Cloud-based Development Environment



- Write code locally -> rsync -> build or run tests in K8s pod

- https://kccncna2022.sched.com/event/182HU

# Problem is …



- Kubernetes v1.27 enables in-place restart-free resize of pod resources (CPU , memory)

- Vertical Pod Autoscaler can resize pod resources on seeing increased usage

  - Reactive - may not be good enough!

# Ideally …



- Pod resources are resized before it becomes a problem

  - Proactive

# eBPF makes it possible!

```yaml
1   apiVersion: v1
2   kind: Pod
3   metadata:
4     name: kube-build-pod
5   spec:
6     containers:
7     - name: kube-build-ctr
8       image: skiibum/kube-build-arm64:v1.25
9       imagePullPolicy: IfNotPresent
10      command: ["tail", "-f", "/dev/null"]
11      resources:
12        limits:
13          cpu: "5"
14          memory: "50Mi"
15        requests:
16          cpu: "4"
17          memory: "50Mi"
```

```python
1   # podsnoop.py: Prototype eBPF program that snoops on pod exec activity
2   #                requires linux-headers, bpfcc-tools, kubectl
3   #                To run: sudo python3 podsnoop.py
4   import os
5   from bcc import BPF
6
7   POD_SNOOP_eBPF_CODE = r"""TRACEPOINT_PROBE(syscalls, sys_enter_execve){
8       char task_cmd[32];
9       bpf_get_current_comm(&task_cmd, sizeof(task_cmd));
10      bpf_trace_printk("Launching program: %s\n", task_cmd);
11      return 0;
12  }"""
13  bpf = BPF(text = POD_SNOOP_eBPF_CODE)
14
15  while True:
16      try:
17          (task, pid, cpu, flags, ts, msg) = bpf.trace_fields()
18          if str.__contains__(msg.decode("utf-8"), "make"):
19              pod_name = os.popen("nsenter -t %s -u hostname 2>/dev/null" % pid).read().strip()
20              if pod_name == "kube-build-pod":
21                  patch_str = '{"spec":{"containers":[{                                    \
22                                        "name":"kube-build-ctr",                           \
23                                        "resources":{"requests":{"memory":"5Gi"},"limits":{"memory":"5Gi"}}}\
24                                    }]}}'
25                  patch_cmd = "kubectl patch pod %s --patch '%s' > /dev/null" % (pod_name, patch_str)
26                  os.system(patch_cmd)
27      except ValueError:
28          continue
```

- action = (command == 'make') ? resize pod : have a beer ;)

# Some Rough Edges

- Ugly demo code: https://github.com/vinaykul/ebpf-playground/tree/main/ebpf-pod-resize-libbpf-rs-CO-RE

```rust
fn create_resize_bpf_map_entry(cid: String, cmd: String, rc_map: &mut Map) -> Result<()> {
    let cg_path;
    let splt = cid.split("//");
    let vec: Vec<&str> = splt.collect();
    let dir_name = vec[1];
    let mut find_dir = Command::new("find");
    find_dir.arg("/sys/fs/cgroup/unified/kubepods").arg("-type").arg("d").arg("-name").arg(dir_name);
    let cg_path_out = find_dir.output().expect("failed to execute process");
    if cg_path_out.stdout.len() == 0 {
        let mut find_dir_cgv1 = Command::new("find");
        find_dir_cgv1.arg("/sys/fs/cgroup/kubepods").arg("-type").arg("d").arg("-name").arg(dir_name);
        let cg_path_out_v1 = find_dir_cgv1.output().expect("failed to execute process");
        cg_path = String::from_utf8(cg_path_out_v1.stdout).unwrap();
    } else {
        cg_path = String::from_utf8(cg_path_out.stdout).unwrap();
    }
    let mut ls_cgid = Command::new("ls");
    ls_cgid.arg("-ladi").arg(cg_path.trim());
    let ls_cgid_out = ls_cgid.output().expect("failed to execute process");
    let ls_cg_out = String::from_utf8(ls_cgid_out.stdout).unwrap();
    let cgsplt = ls_cg_out.split(" ");
    let veccg: Vec<&str> = cgsplt.collect();
    let cgroup_id = veccg[0];
    println!("DBG: cgroup_id: '{}' , cmd: '{}'", cgroup_id,  cmd);
```

- Hacky way to find containerID <> cgroup_id mapping, isn't it?

# Some Rough Edges

```c
SEC("tracepoint/syscalls/sys_enter_execve")
int podsnoop(void *ctx) {
    u64 cgroup_id = bpf_get_current_cgroup_id();
    struct pod_command *val = bpf_map_lookup_elem(&resize_containers_map, &cgroup_id);
    if (val != NULL) {
        u8 is_equal = 1;
        struct pod_exec_event pxevent = {};
        bpf_get_current_comm(&pxevent.cgroup_cmd, sizeof(pxevent.cgroup_cmd));
        //TODO: Find a more efficient way. Maybe 'val->cmd' should be u64 hash
        for (u8 i = 0; i < BUF_SIZE; i++) {
            if (pxevent.cgroup_cmd[i] != val->cmd[i]) {
                is_equal = 0;
                break;
            }
            if (pxevent.cgroup_cmd[i] == '\0' || val->cmd[i] == '\0') {
                break;
            }
        }
        if (is_equal) {
            pxevent.cgroup_id = cgroup_id;
            u64 id = bpf_get_current_pid_tgid();
            pxevent.cgroup_pid = id & 0xFFFFFFFF;
            long rv = bpf_perf_event_output(ctx, &pod_exec_events, BPF_F_CURRENT_CPU, &pxe
            if (rv != 0) {
                bpf_printk("DBG: podsnoop call to bpf_perf_event_output failed. ErrCode: %
            }
        }
    }
    return 0;
}
```

- ~~Maybe add bpf_strncmp(...)?~~ NVM: It has already been added in libbpf

# Potential libbpf Extensions

- Add: `u64 bpf_get_container_cgroup_id(const char *container_id)`

  - How: Scan /sys/fs/cgroup for container_id (For cgroups v1, look under /sys/fs/cgroup/cpu)

    - If found, return its i-node number

    - If not, return 0

- Add: `int bpf_get_cgroup_id_container(u64 cgroup_id, const char *container_id)`

  - How: Native version of `find /sys/fs/cgroup -inum <cgroup_id>`

- Add: ~~`int bpf_strncmp(const char *s1, const char *s2, register size_t n)`~~

# Summary + Discussion + Q&A

- At least two use cases that could leverage simplified cgroup_id <> container_id helper.

- Is this the right way to do it?

  - If not, any alternative suggestions?