A vertical red bar on the left side of the slide contains various white and dark red icons representing technology: a cloud with a keyhole, a database cylinder, a server rack, a monitor, a network diagram with arrows and nodes, and a person silhouette.

Using bpftrace with Performance Co-Pilot & Grafana

Andreas Gerstmayr
Software Engineer





Performance Co-Pilot

system performance analysis toolkit



Grafana

analytics and monitoring solution



bpftrace

high-level tracing language for eBPF

Performance Co-Pilot

Toolkit for collecting, analyzing, visualizing and responding to the status and performance of servers, applications and networks.

92

Agents

Performance Metrics Domain Agents export performance data from the kernel, services (e.g. PostgreSQL) and other instrumented applications

60000+

Metrics

Metrics have associated metadata:

Semantics: instant, counter, discrete

Unit: Kilobytes, bits/sec, etc.

Type: float, (un)signed int, string, ...

Instances: e.g. sda1, sda2

Grafana



Visualization

Visualization of time series, heatmaps, tabular data, etc.

Customizable

Customizable dashboards, drag & drop, live preview

Multiple data sources

PCP, Elasticsearch, Prometheus, PostgreSQL, ...

Extensible

Custom data sources & panels

bpftrace

tracing language for eBPF

```
// read.bt file
tracepoint:syscalls:sys_enter_read
{
    @start[tid] = nsecs;
}

tracepoint:syscalls:sys_exit_read / @start[tid] /
{
    @times = hist(nsecs - @start[tid]);
    delete(@start[tid]);
}
```

bpftrace

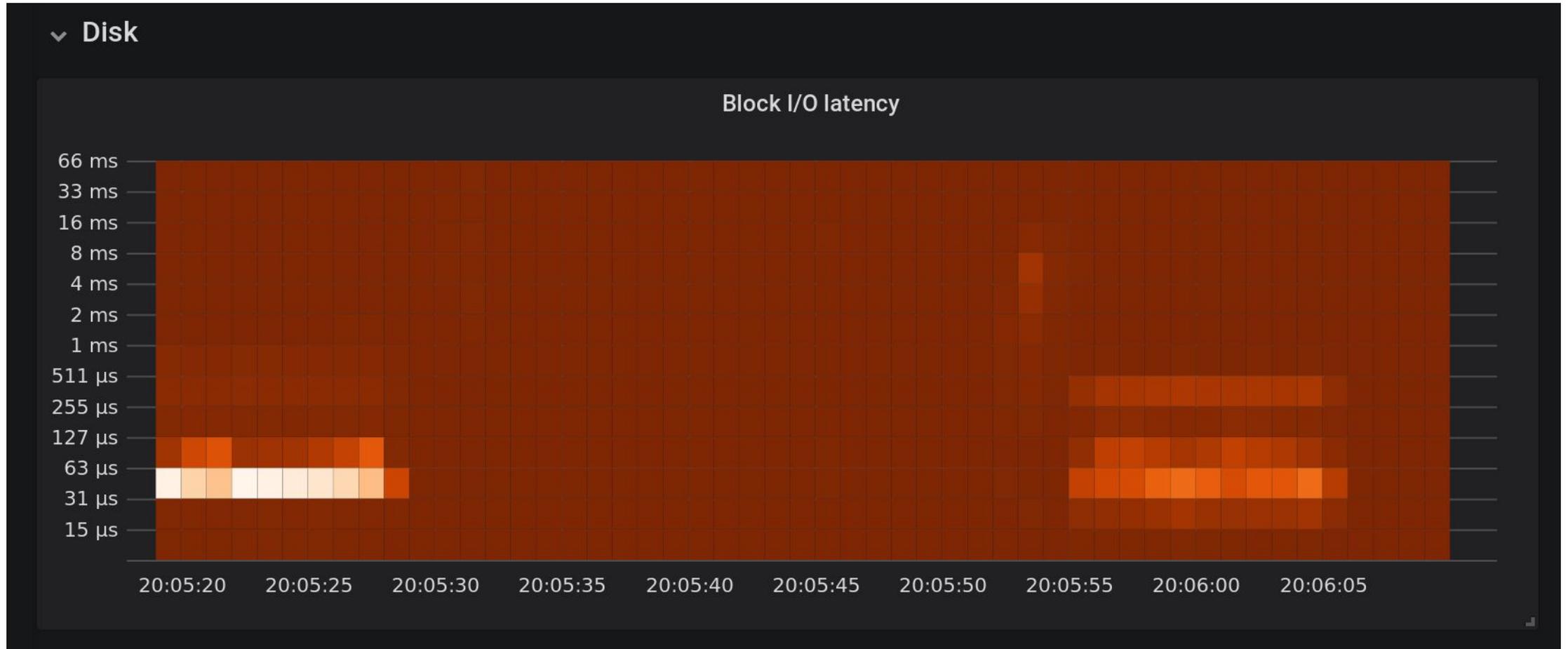
tracing language for eBPF

```
$ bpftrace read.bt  
Attaching 2 probes...  
^C
```

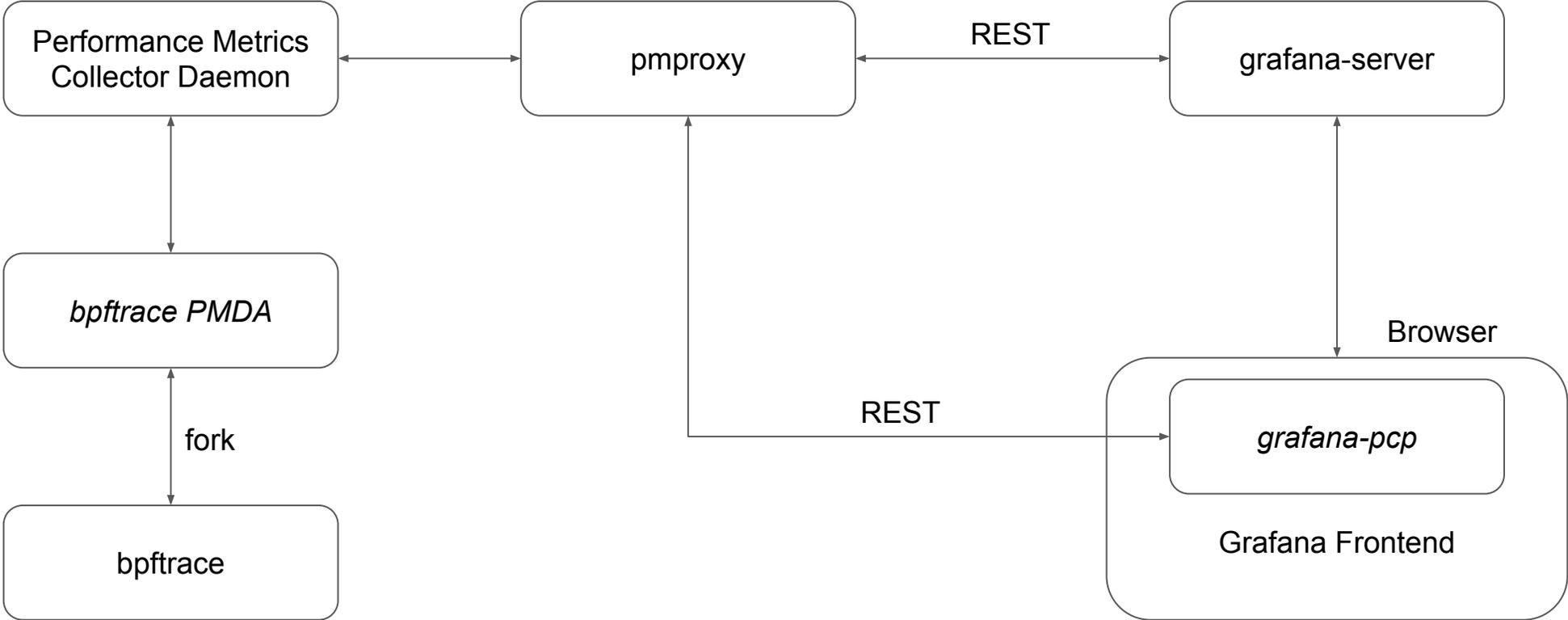
```
@times:
```

[256, 512)	326	@	
[512, 1k)	7715	@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@	
[1k, 2k)	15306	@@	
[2k, 4k)	609	@@	
[4k, 8k)	611	@@	
[8k, 16k)	438	@	
[16k, 32k)	59		
[32k, 64k)	36		
[64k, 128k)	5		

PCP + Grafana + bpftrace



Architecture



Installation

Installation on Fedora 31

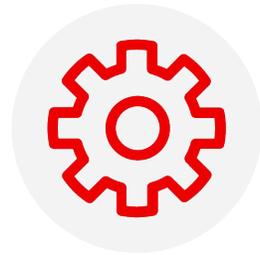
```
# bpftrace version with support for kernel 5.4+
```

```
$ dnf copr enable agerstmayr/bpftrace
```

```
$ sudo dnf install pcp pcp-pmda-bpftrace grafana grafana-pcp bpftrace
```

```
$ sudo systemctl start pmproxy grafana-server
```

bpfttrace PMDA Configuration



Configure authentication

```
$ sudo vim /var/lib/pcp/pmdas/bpfttrace/bpfttrace.conf
```



Install/Enable the PMDA

```
$ cd /var/lib/pcp/pmdas/bpfttrace
```

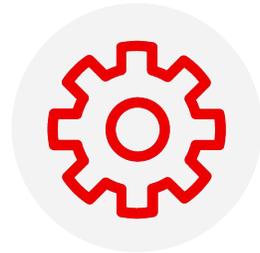
```
$ sudo ./Install
```

Grafana Configuration



Enable plugin

Go to Configuration / Plugins and enable *Performance Co-Pilot*



Configure data source

Go to Configuration / Data Sources / Add data source, select *PCP bpftrace*, enter URL of *pmproxy*, click Save

bpftrace tools

runqlat

CPU scheduler run queue latency

```
15
16 tracepoint:sched:sched_wakeup, tracepoint:sched:sched_
17 {
18     @qtime[args->pid] = nsecs;
19 }
20
21 tracepoint:sched:sched_switch
22 {
23     if (args->prev_state == TASK
24         @qtime[args->prev_pid] = n
25     }
26
27     $ns = @qtime[args->next_pid];
28     if ($ns) {
29         @usecs = hist((nsecs - $ns) / 1000);
30     }
31     delete(@qtime[args->next_pid]);
32 }
33
34 END
35 {
```

tracepoint:sched:sched_kthread_stop	probe
tracepoint:sched:sched_kthread_stop_ret	probe
tracepoint:sched:sched_migrate_task	probe
tracepoint:sched:sched_move_numa	probe
tracepoint:sched:sched_pi_setprio	probe
tracepoint:sched:sched_process_exec	probe
tracepoint:sched:sched_process_exit	probe
tracepoint:sched:sched_process_fork	probe

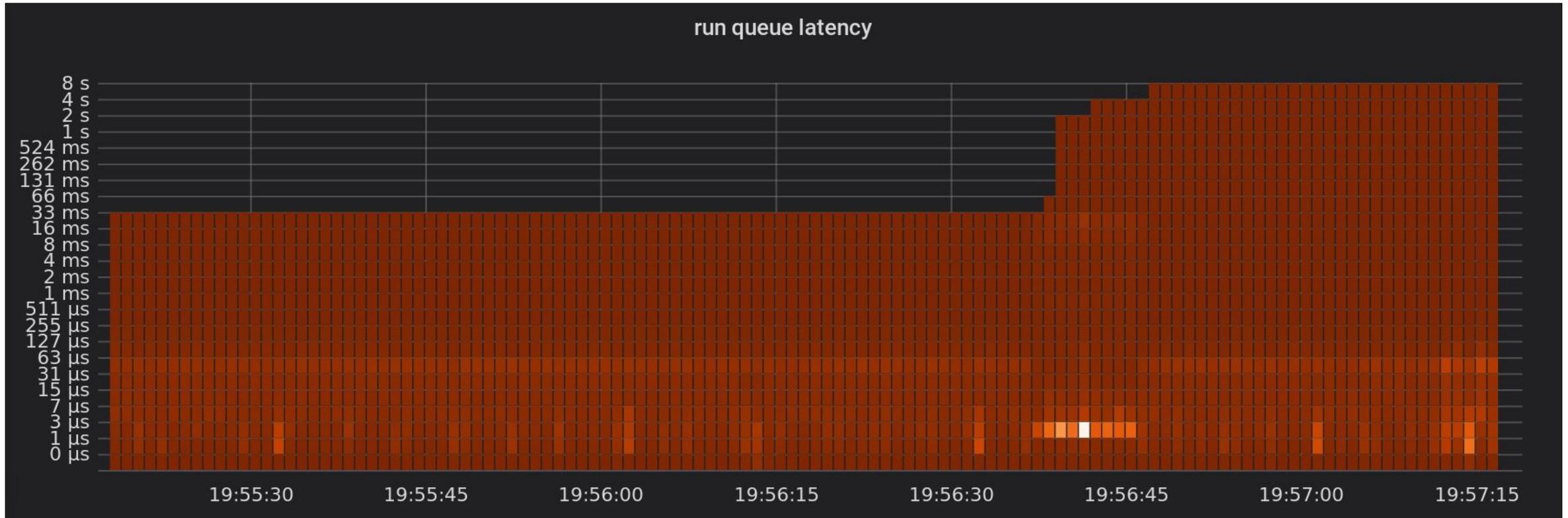
Legend legend format Format Heatmap URL override URL

runqlat

```
1  #include <linux/sched.h>
2  tracepoint:sched:sched_wakeup,
3  tracepoint:sched:sched_wakeup_new
4  {
5      @qtime[args->pid] = nsecs;
6  }
7  tracepoint:sched:sched_switch
8  {
9      if (args->prev_state == TASK_RUNNING) {
10         @qtime[args->prev_pid] = nsecs;
11     }
12     $ns = @qtime[args->next_pid];
13     if ($ns) {
14         @usecs = hist((nsecs - $ns) / 1000);
15     }
16     delete(@qtime[args->next_pid]);
17 }
```

runqlat

CPU scheduler run queue latency



tcpconnect

Trace TCP connect()

```
$ sudo /usr/share/bpftrace/tools/tcpconnect.bt
```

```
Attaching 2 probes...
```

```
Tracing tcp connections. Hit Ctrl-C to end.
```

TIME	PID	COMM	SADDR	SPORT	DADDR	DPORT
18:17:38	22576	curl	-	45666	1.1.1.1	80
18:17:53	22598	curl	-	50966	8.8.4.4	80
18:20:40	23429	wget	-	40606	2a02:26f0:10e:1b8::d44	443
18:21:15	1746	syncthing	-	44096	2a03:b0c0:0:1010::bb:4001	443

```
^C
```

tcpconnect

Trace TCP connect()

```
// table-retain-lines: 10
```

```
BEGIN
```

```
{
```

```
    printf("%s,%s,%s,", "TIME", "PID", "COMM");
```

```
    printf("%s,%s,%s,%s\n", "SADDR", "SPORT", "DADDR",  
"DPORT");
```

```
}
```

```

1 kprobe:tcp_connect
2 {
3     $sk = ((struct sock *) arg0);
4     $inet_family = $sk->__sk_common.skc_family;
5
6     if ($inet_family == AF_INET || $inet_family == AF_INET6)
7     {
8         // [...]
9         $lport = $sk->__sk_common.skc_num;
10        $dport = $sk->__sk_common.skc_dport;
11        $dport = ($dport >> 8) | (($dport << 8) & 0x00FF00);
12
13        time("%H:%M:%S,");
14        printf("%d,%s,", pid, comm);
15        printf("%s,%d,%s,%d\n", $saddr, $lport, $daddr, $dport);
16    }
17 }

```

tcpconnect

Trace TCP connect()

trace TCP connect()						
TIME ▾	PID	COMM	SADDR	SPORT	DADDR	DPORT
18:21:15	1746	syncting	2a02:8388:xxx:xxxx:xxxx:xxxx:xxxx:xxxx	44096	2a03:b0c0:0:1010::bb:4001	443
18:20:40	23429	wget	2a02:8388:xxx:xxxx:xxxx:xxxx:xxxx:xxxx	40606	2a02:26f0:10e:1b8::d44	443
18:17:53	22598	curl	192.168.0.185	50966	8.8.4.4	80
18:17:38	22576	curl	192.168.0.185	45666	1.1.1.1	80

Sampling kernel stacks

```
profile:hz:99
{
    @stacks[kstack] = count();
}
```


Preconfigured container

```
$ sudo -H podman run -d \  
  --privileged \  
  -v /lib/modules:/lib/modules:ro \  
  -v /usr/src:/usr/src:ro \  
  -p 127.0.0.1:3000:3000 \  
  quay.io/performancecopilot/pcp-preview
```

Live Demo

 github.com/andreasgerstmayr

 linkedin.com/in/andreasgerstmayr

