

# GOTC

## 全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

# OPEN SOURCE , OPEN WORLD #

### 「CNCF云原生」专场

本期议题：eBPF在旷视PaaS平台的落地实践

王续 2021年07月10日

旷视PaaS平台的演进

初识eBPF

eBPF实践之路

## 愿景:

面向全体研发, 在**混合云**场景下, 定义基于Kubernetes的**Application**概念, 提供CI和周边基础设施, 形成内部完善的、面向**云原生的**、**GitOps**驱动的**PaaS**平台

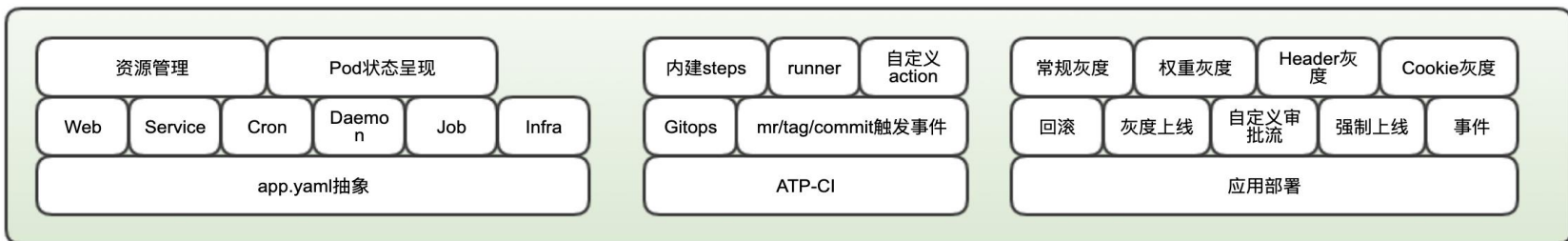
## 现状:

**900+**个应用, **10+**个集群, 每日**200+**次上线

## 应用管理



## 应用生命周期管理



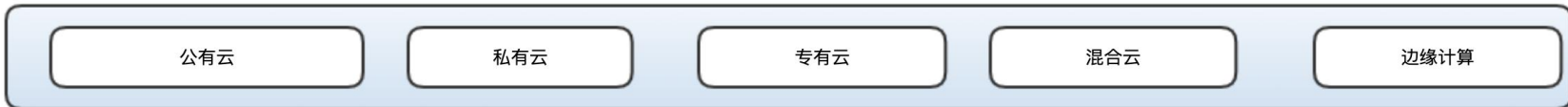
## 应用资源



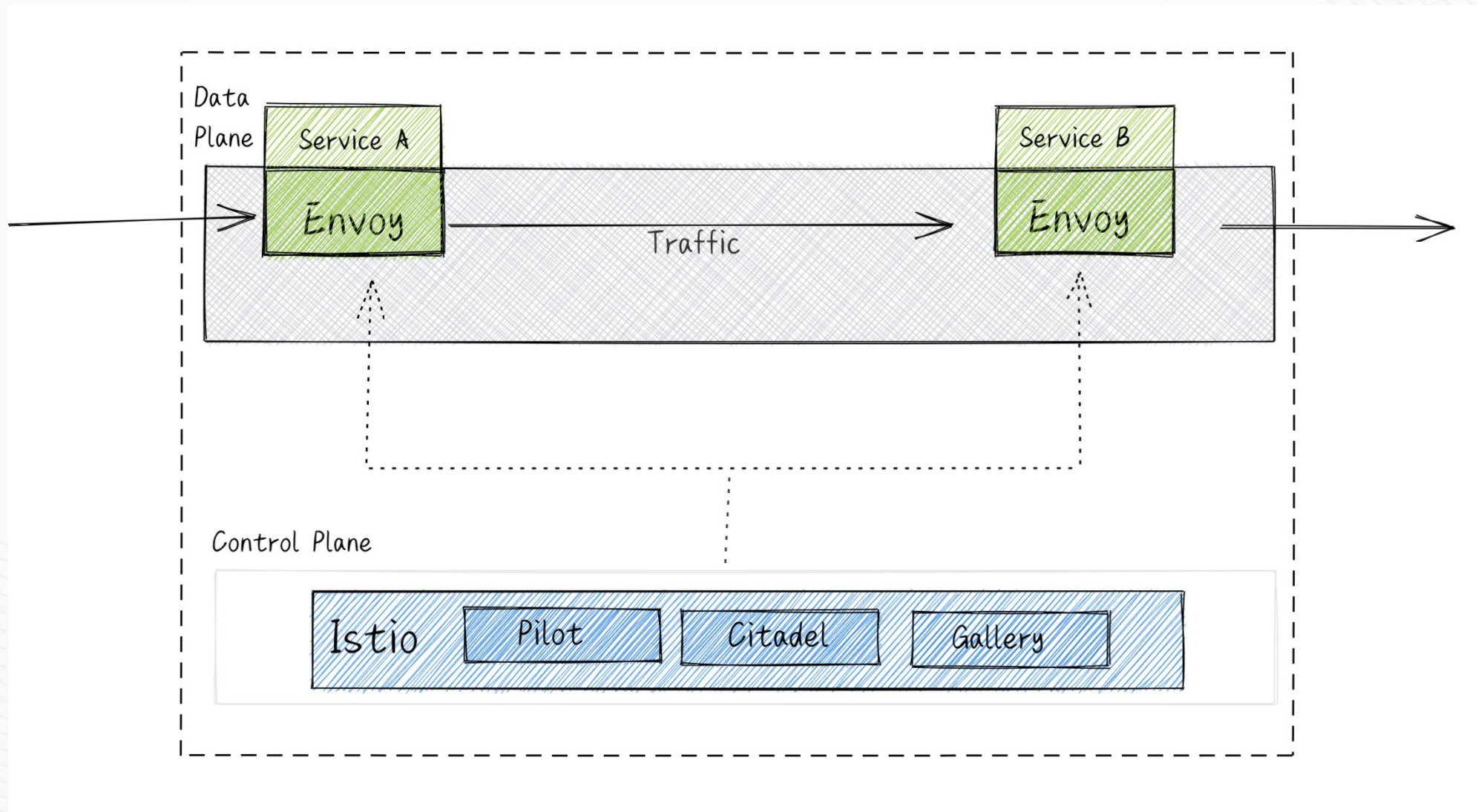
## K8S集群



## 场景支持







再优雅一点?

GOTC

应用无感

高性能

技术成本低

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

# The BSD Packet Filter: A New Architecture for User-level Packet Capture\*

Steven McCanne<sup>†</sup> and Van Jacobson<sup>†</sup>  
Lawrence Berkeley Laboratory  
One Cyclotron Road  
Berkeley, CA 94720  
mccanne@ee.lbl.gov, van@ee.lbl.gov

December 19, 1992

## Abstract

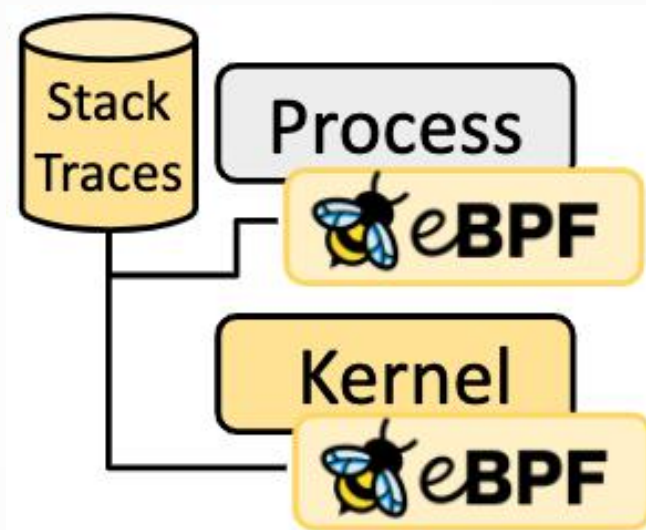
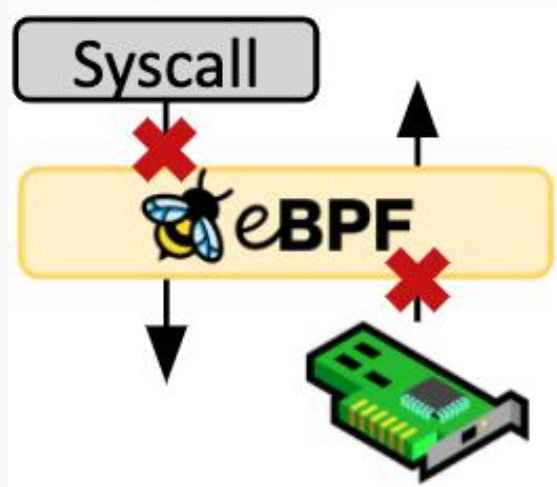
Many versions of Unix provide facilities for user-level packet capture, making possible the use of general purpose workstations for network monitoring. Because network monitors run as user-level processes, packets must be copied across the kernel/user-space protection boundary. This copying can be minimized by deploying a kernel agent called a *packet filter*, which discards unwanted packets as early as possible. The original Unix packet filter was designed around a stack-based

SunOS, the Ultrix Packet Filter[2] in DEC's Ultrix and Snoop in SGI's IRIX.

These kernel facilities derive from pioneering work done at CMU and Stanford to adapt the Xerox Alto 'packet filter' to a Unix kernel[8]. When completed in 1980, the CMU/Stanford Packet Filter, CSPF, provided a much needed and widely used facility. However on today's machines its performance, and the performance of its descendents, leave much to be desired — a design that was entirely appropriate for a 64KB PDP-11 is simply not a good match to a 16MB Sparcstation

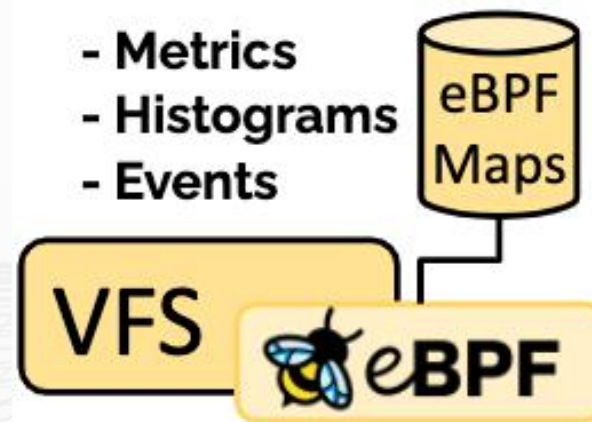
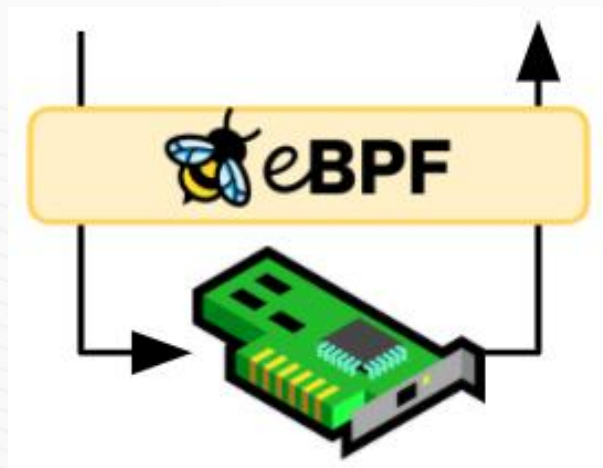


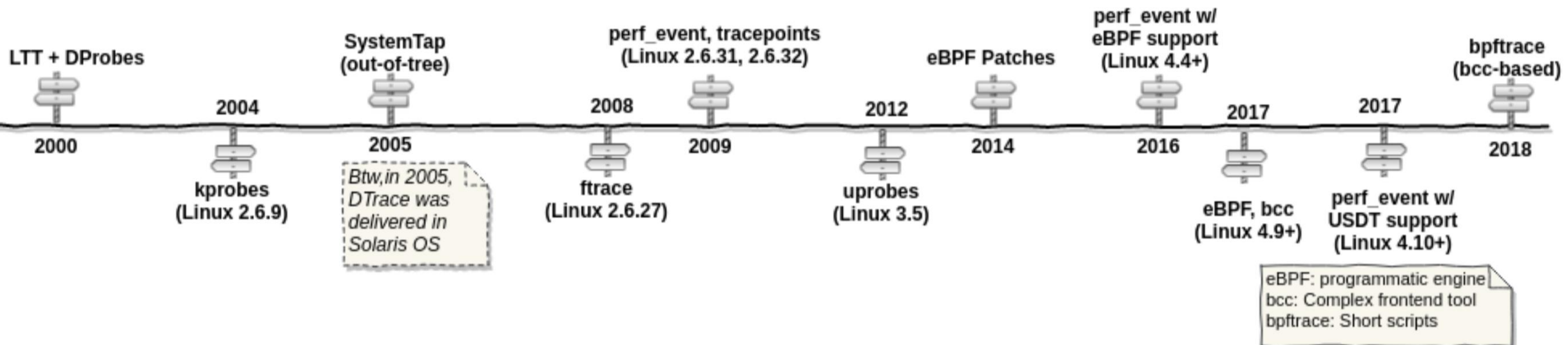
维度	cBPF	eBPF
内核版本	Linux 2.1.75 (1997年)	Linux 3.18 (2014年) [4.x for kprobe/uprobe/tracepoint/perf-event]
寄存器数目	2个: A, X	10个: R0-R9, 另外 R10 是一个只读的帧指针
寄存器宽度	32位	64位
存储	16 个内存位: M[0-15]	512 字节堆栈, 无限制大小的 "map" 存储
限制的内核调用	非常有限, 仅限于 JIT 特定	有限, 通过 bpf_call 指令调用
目标事件	数据包、seccomp-BPF	数据包、内核函数、用户函数、跟踪点 PMCs 等

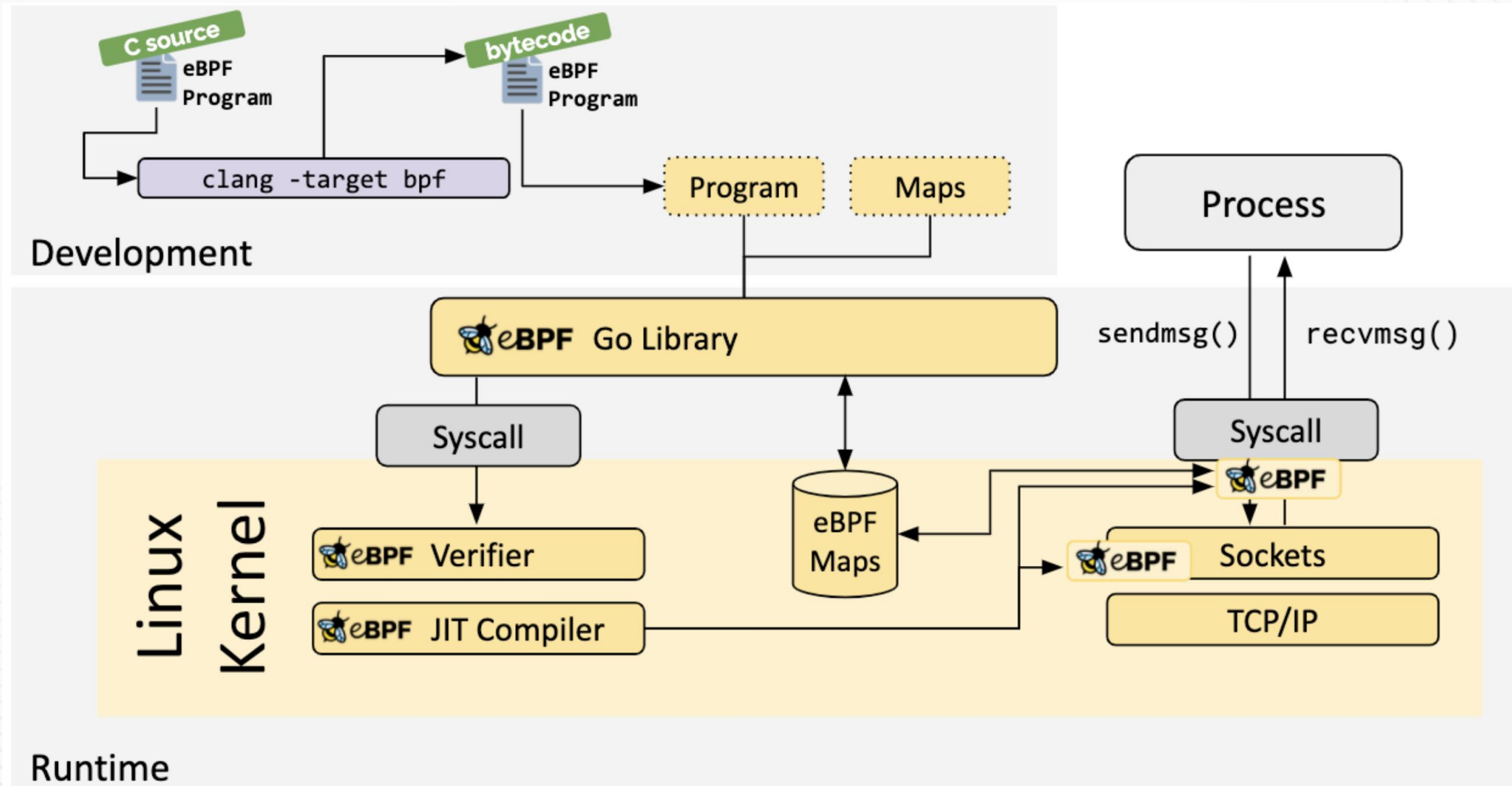


安全 性能分析

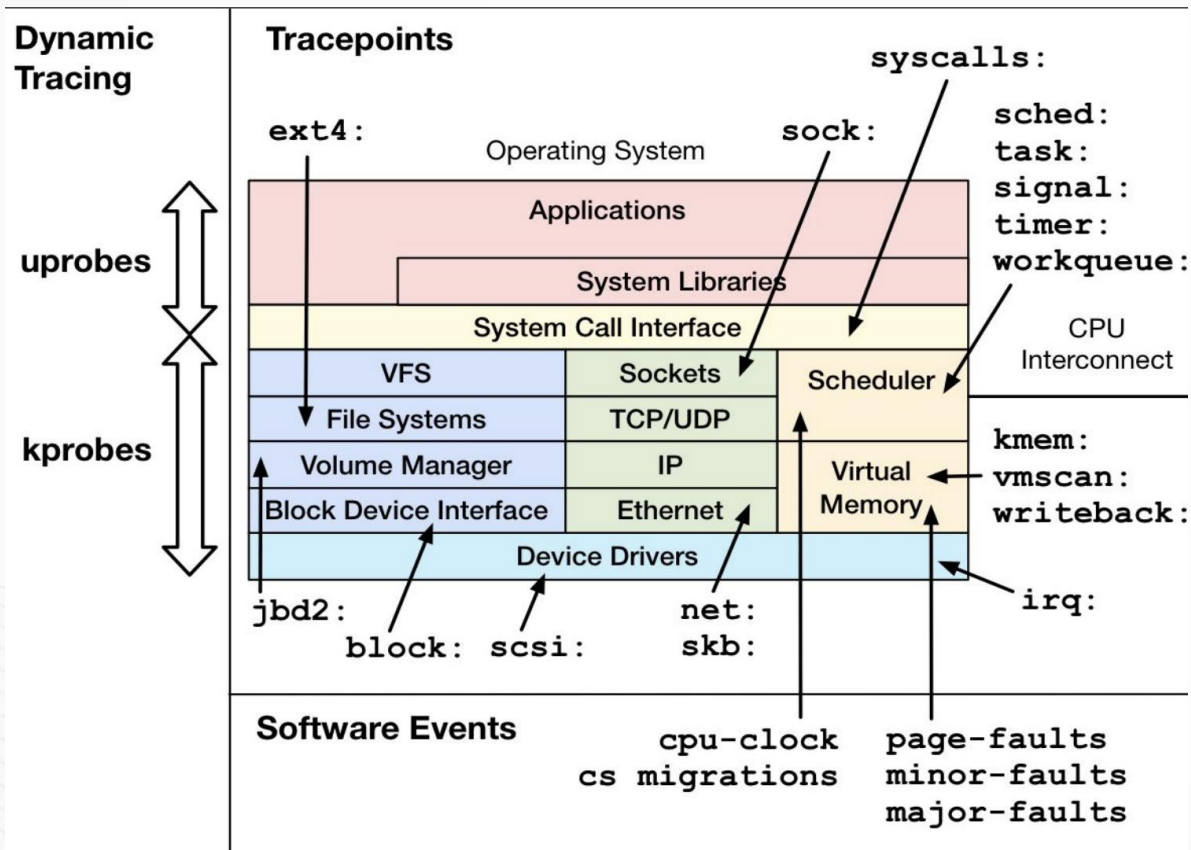
网络 观测&监控







```
# ./tcplife.bt
PID   COMM      LADDR          LPORT  RADDR          RPORT  TX_KB  RX_KB  MS
20976 ssh       127.0.0.1      56766  127.0.0.1      22      6    10584  3059
20977 sshd      127.0.0.1      22     127.0.0.1      56766  10584  6    3059
14519 monitord  127.0.0.1      44832  127.0.0.1      44444   0      0    0
4496  Chrome_IOT 7f00:6:5ea7::a00:0 42846  0:0:bb01::     443     0      3    12441
4496  Chrome_IOT 7f00:6:5aa7::a00:0 42842  0:0:bb01::     443     0      3    12436
4496  Chrome_IOT 7f00:6:62a7::a00:0 42850  0:0:bb01::     443     0      3    12436
4496  Chrome_IOT 7f00:6:5ca7::a00:0 42844  0:0:bb01::     443     0      3    12442
4496  Chrome_IOT 7f00:6:60a7::a00:0 42848  0:0:bb01::     443     0      3    12436
4496  Chrome_IOT 10.0.0.65      33342  54.241.2.241   443     0      3    10717
4496  Chrome_IOT 10.0.0.65      33350  54.241.2.241   443     0      3    10711
4496  Chrome_IOT 10.0.0.65      33352  54.241.2.241   443     0      3    10712
14519 monitord  127.0.0.1      44832  127.0.0.1      44444   0      0    0
```

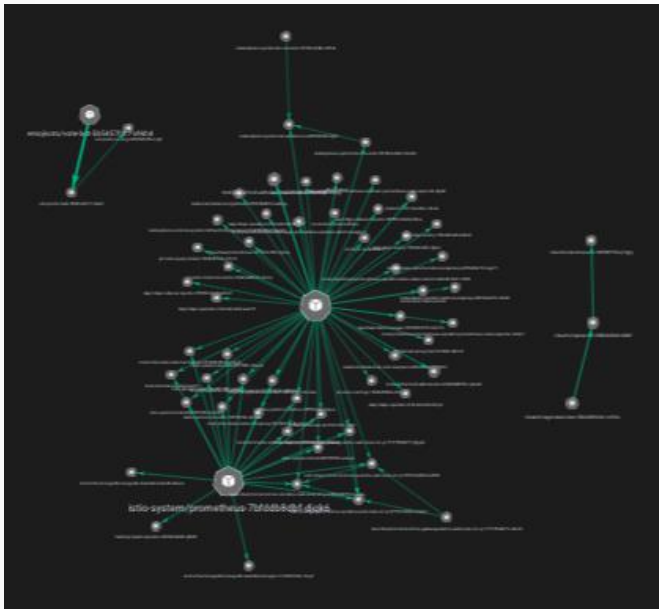


```

// session ended: calculate lifespan and print
if ($newstate == TCP_CLOSE && @birth[$sk]) {
    $delta_ms = (nsecs - @birth[$sk]) / 1e6;
    $lport = $sk->__sk_common.skc_num;
    $dport = $sk->__sk_common.skc_dport;
    $dport = ($dport >> 8) | (($dport << 8) & 0xff00);
    $tp = (struct tcp_sock *)$sk;
    $pid = @skpid[$sk];
    $comm = @skcomm[$sk];
    if ($comm == "") {
        // not cached, use current task
        $pid = pid;
        $comm = comm;
    }

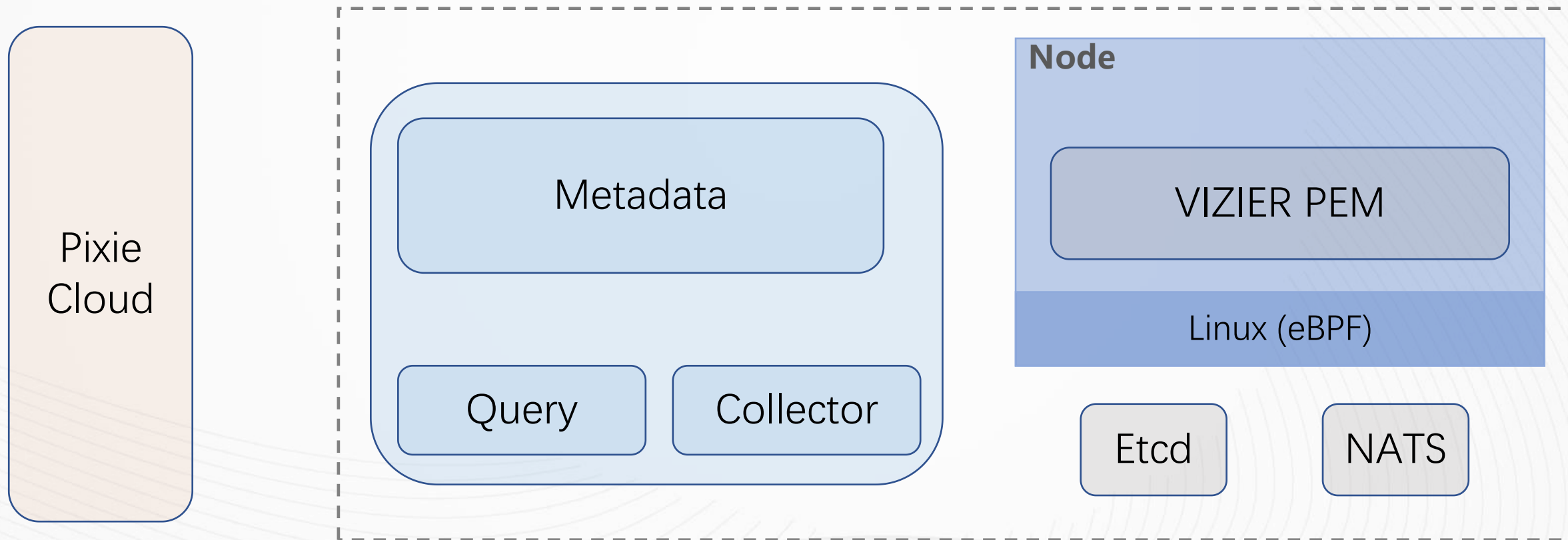
    $family = $sk->__sk_common.skc_family;
    $saddr = ntop(0);
    $daddr = ntop(0);
    if ($family == AF_INET) {
        $saddr = ntop(AF_INET, $sk->__sk_common.skc_rcv_saddr);
        $daddr = ntop(AF_INET, $sk->__sk_common.skc_daddr);
    } else {
        // AF_INET6
        $saddr = ntop(AF_INET6,
            $sk->__sk_common.skc_v6_rcv_saddr.in6_u.u6_addr8);
        $daddr = ntop(AF_INET6,
            $sk->__sk_common.skc_v6_daddr.in6_u.u6_addr8);
    }
    printf("%-5d %-10.10s %-15s %-5d %-15s %-6d ", $pid,
        $comm, $saddr, $lport, $daddr, $dport);
    printf("%5d %5d %d\n", $tp->bytes_acked / 1024,
        $tp->bytes_received / 1024, $delta_ms);

    delete(@birth[$sk]);
    delete(@skpid[$sk]);
    delete(@skcomm[$sk]);
}
    
```



PID	PROCESS NAME	ASID	TIMES...	VI...	AVERAG...
22,953	/usr/bin/promtail -config.file=/etc/promtail/promtail.yaml -client.url=http://loki-distributo...	9	2021/6/29...	6.8 GB	414.1 MB
22,953	/usr/bin/promtail -config.file=/etc/promtail/promtail.yaml -client.url=http://loki-distributo...	9	2021/6/29...	6.8 GB	414.1 MB
23,951	haproxy -W -db -f /usr/local/etc/haproxy/haproxy.cfg	9	2021/6/29...	3 GB	84.8 MB
23,951	haproxy -W -db -f /usr/local/etc/haproxy/haproxy.cfg	9	2021/6/29...	3 GB	84.8 MB
23,309	/app/mcd-admission	9	2021/6/29...	3 GB	28.3 MB
23,309	/app/mcd-admission	9	2021/6/29...	3 GB	28.3 MB
22,305	/app/mcd-admission	9	2021/6/29...	2.9 GB	28 MB
22,305	/app/mcd-admission	9	2021/6/29...	2.9 GB	28 MB
6,983	/app/mcd-admission	9	2021/6/29...	2.8 GB	27.8 MB
6,983	/app/mcd-admission	9	2021/6/29...	2.8 GB	27.8 MB

RE...	RE...	RE...	REQ_PATH	RESP_STATUS	RESP_BODY	LA...	SERVICE
10.0.1....	40040	POST	/_bulk	200	{ took: 81, errors: false, items: [ { index: { _index: queue_metrics_d1bd92a3b039400cbaafc60...	82.3 ms	clearml/ela...
10.0.1....	40040	POST	/_bulk	200	{ took: 72, errors: false, items: [ { index: { _index: queue_metrics_d1bd92a3b039400cbaafc60...	73.3 ms	clearml/ela...
10.0.1....	40040	POST	/_bulk	200	{ took: 67, errors: false, items: [ { index: { _index: queue_metrics_d1bd92a3b039400cbaafc60...	68.1 ms	clearml/ela...
10.0.1....	40040	POST	/_bulk	200	{ took: 94, errors: false, items: [ { index: { _index: queue_metrics_d1bd92a3b039400cbaafc60...	95.4 ms	clearml/ela...
10.0.1....	40040	POST	/_bulk	200	{ took: 52, errors: false, items: [ { index: { _index: queue_metrics_d1bd92a3b039400cbaafc60...	53.3 ms	clearml/ela...
10.0.1....	40040	POST	/_bulk	200	{ took: 62, errors: false, items: [ { index: { _index: queue_metrics_d1bd92a3b039400cbaafc60...	62.9 ms	clearml/ela...
10.0.1....	40040	POST	/_bulk	200	{"took":40,"errors":false,"items":[{"index":{"_index":"worker_stats_d1bd92a3b039400cbaafc60a7a5b1e52b_2021-0...	40.7 ms	clearml/ela...





## Table

TI...	UPID	RE...	RE...	M...	REQ_HEADERS	RE...	REQ_PATH	RE...
4/20/2...	00000032-0000-325...	192.16...	47082	1	{ Connection: cl...	GET	/catalogue?page=1&size=...	
<pre>}, resp_status: 500, resp_message: Internal Server Error, resp_body: {   error: Do: database connection error,   status_code: 500,   status_text: Internal Server Error }, resp_body_size: 98, latency: 440476, pod: px-sock-shop/catalogue-8f6fdb6d8-24x72, service: px-sock-shop/catalogue</pre>								

```
def http_data(start_time: str, num_head: int):  
    df = px.DataFrame(table='http_events', start_time=start_time)  
  
    df.namespace = df.ctx['namespace']  
    df.node = df.ctx['node']  
    df.pod = df.ctx['pod']  
    df.pid = px.upid_to_pid(df.upid)  
  
    # Remove some columns.  
    df = df.drop(['upid', 'trace_role', 'content_type', 'minor_version'])  
  
    # Restrict number of results.  
    df = df.head(num_head)  
  
    return df
```

```
constexpr DataElement kHTTPElements[] = {
    canonical_data_elements::kTime,
    canonical_data_elements::kUPID,
    canonical_data_elements::kRemoteAddr,
    canonical_data_elements::kRemotePort,
    canonical_data_elements::kTraceRole,
    {"minor_version", "HTTP minor version, HTTP1 uses 1, HTTP2 set this value to 0",
     types::DataType::INT64,
     types::SemanticType::ST_NONE,
     types::PatternType::GENERAL_ENUM},
    {"content_type", "Type of the HTTP payload, can be JSON or protobuf",
     types::DataType::INT64,
     types::SemanticType::ST_NONE,
     types::PatternType::GENERAL_ENUM,
     &kHTTPContentTypeDecoder},
    {"req_headers", "Request headers in JSON format",
     types::DataType::STRING,
     types::SemanticType::ST_NONE,
     types::PatternType::STRUCTURED},
    ...
};
// clang-format on

constexpr auto kHTTPTable = DataTableSchema("http_events", "HTTP request-response pair events",
                                             kHTTPElements, std::chrono::milliseconds{1000});
```

```
SocketTraceConnector::SocketTraceConnector(std::string_view source_name)
: SourceConnector(source_name, kTables), conn_stats_(&conn_trackers_mgr_), uprobe_mgr_(this) {
  proc_parser_ = std::make_unique<system::ProcParser>(system::Config::GetInstance());
  InitProtocolTransferSpecs();
}

{kProtocolHTTP, TransferSpec{FLAGS_stirling_enable_http_tracing,
                             kHTTPTableNum,
                             {kRoleClient, kRoleServer},
                             TRANSFER_STREAM_PROTOCOL(http)}}},
```

## 注册transfer

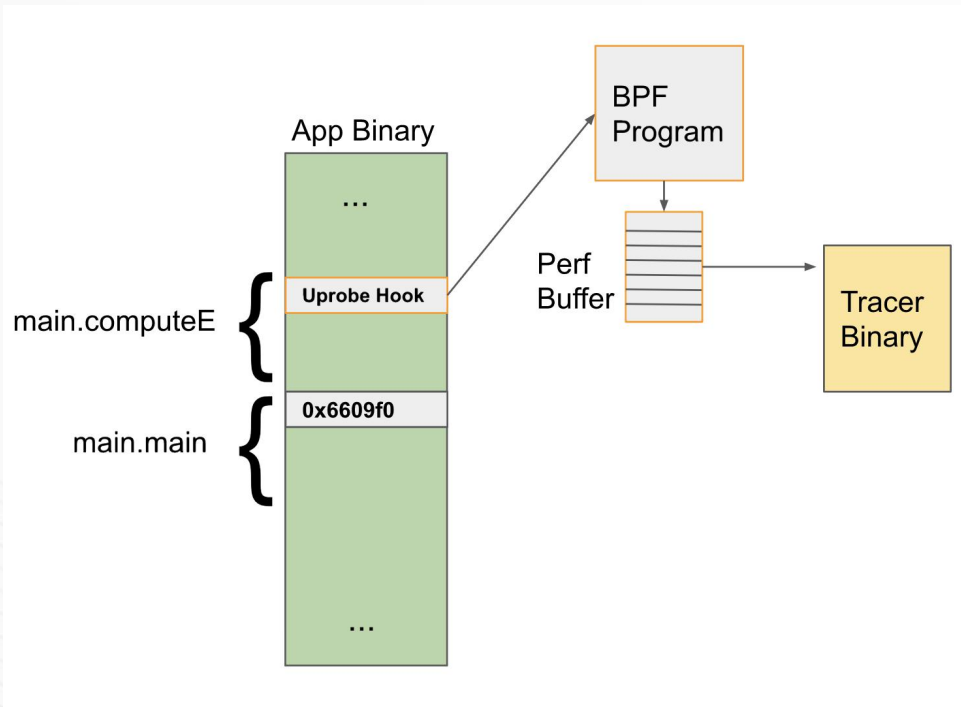
```
if (tracker->state() == ConnTracker::State::kTransferring) {
  // ProcessToRecords() parses raw events and produces messages in format that are expected by
  // table store. But those messages are not cached inside ConnTracker.
  auto result = tracker->ProcessToRecords<TProtocolTraits>();
  for (auto& msg : result) {
    AppendMessage(ctx, *tracker, std::move(msg), data_table);
  }
}
```

## 解析events

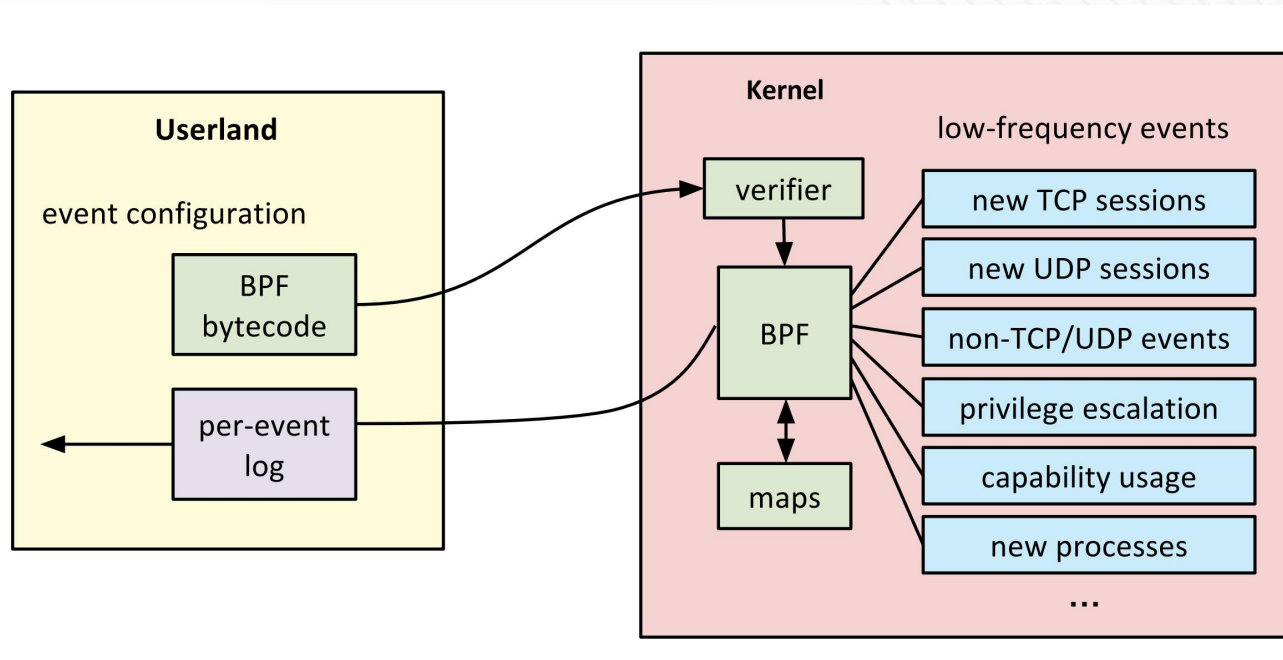
```
r.Append<r.ColIndex("remote_addr")>(conn_tracker.remote_endpoint().AddrStr());
r.Append<r.ColIndex("remote_port")>(conn_tracker.remote_endpoint().port());
r.Append<r.ColIndex("trace_role")>(conn_tracker.role());
r.Append<r.ColIndex("major_version")>(1);
r.Append<r.ColIndex("minor_version")>(resp_message.minor_version);
r.Append<r.ColIndex("content_type")>(static_cast<uint64_t>(content_type));
```

## 构造table

## Remote Debug



## Security



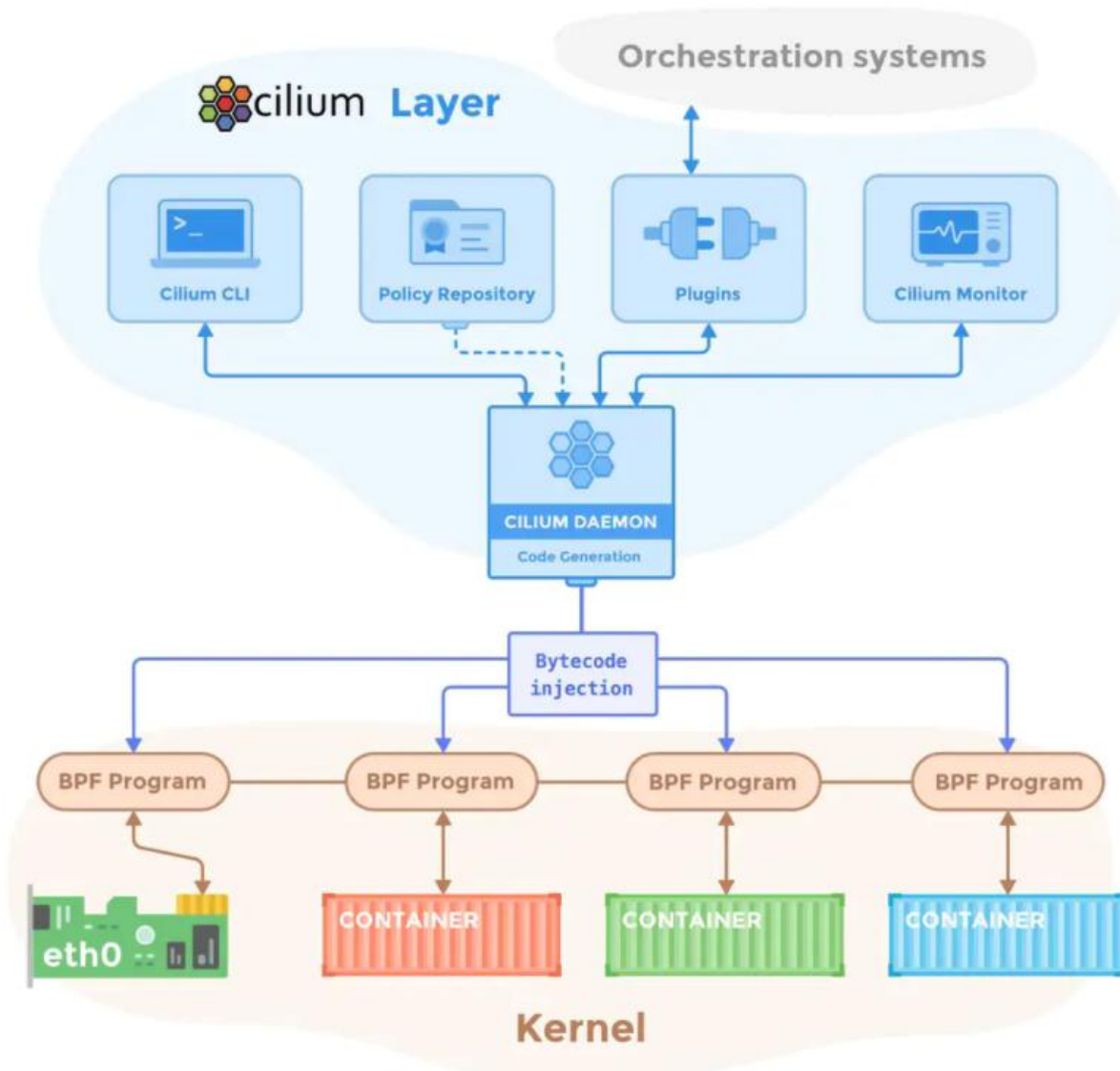
## 性能

Latency percentiles at 50k RPS ( $\approx 20\%$  of maximum RPS), ms

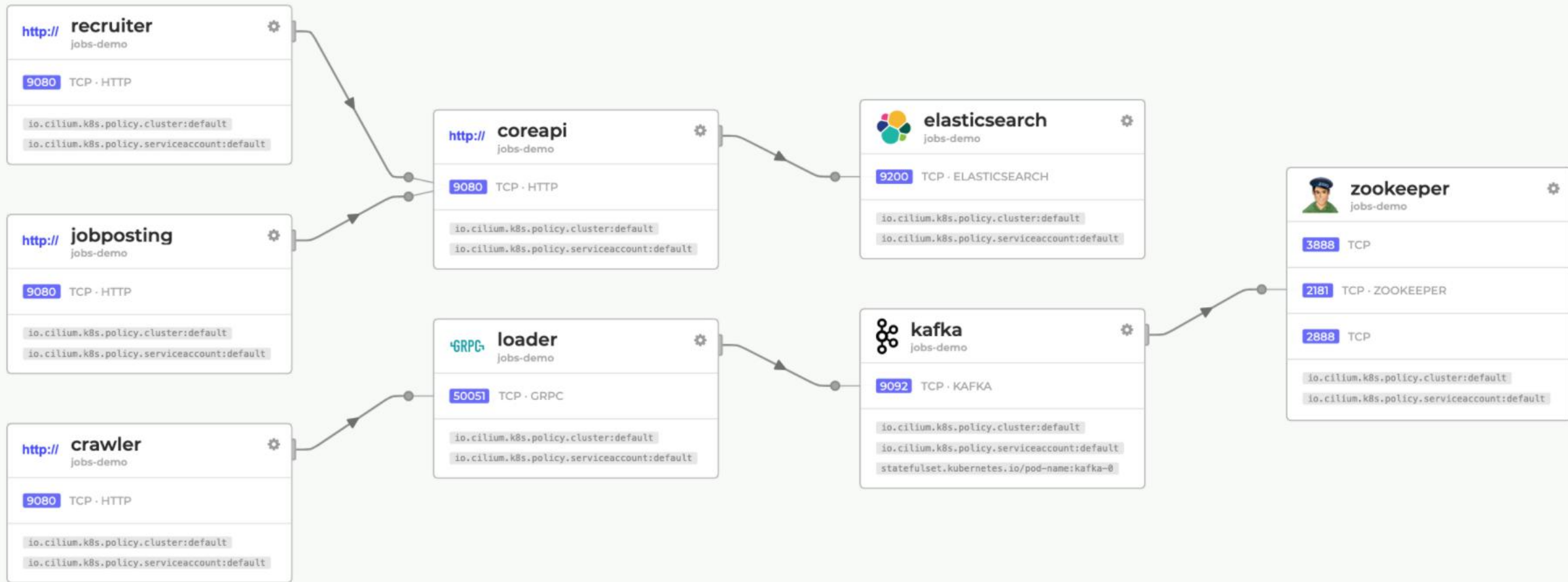
Setup	95 %ile	99 %ile	99.5 %ile	99.99 %ile	99.999 %ile	Max Latency
IPvlan	0.6	0.8	0.9	5.7	9.6	15.8
aws-vpc	0.7	0.9	1	5.6	9.8	403.1
host-gw	0.7	0.9	1	7.4	12	202.5
vxlan	0.8	1.1	1.2	5.7	201.5	402.5
-- net=host	0.5	0.7	0.7	6.4	9.9	14.8

## Network Policy

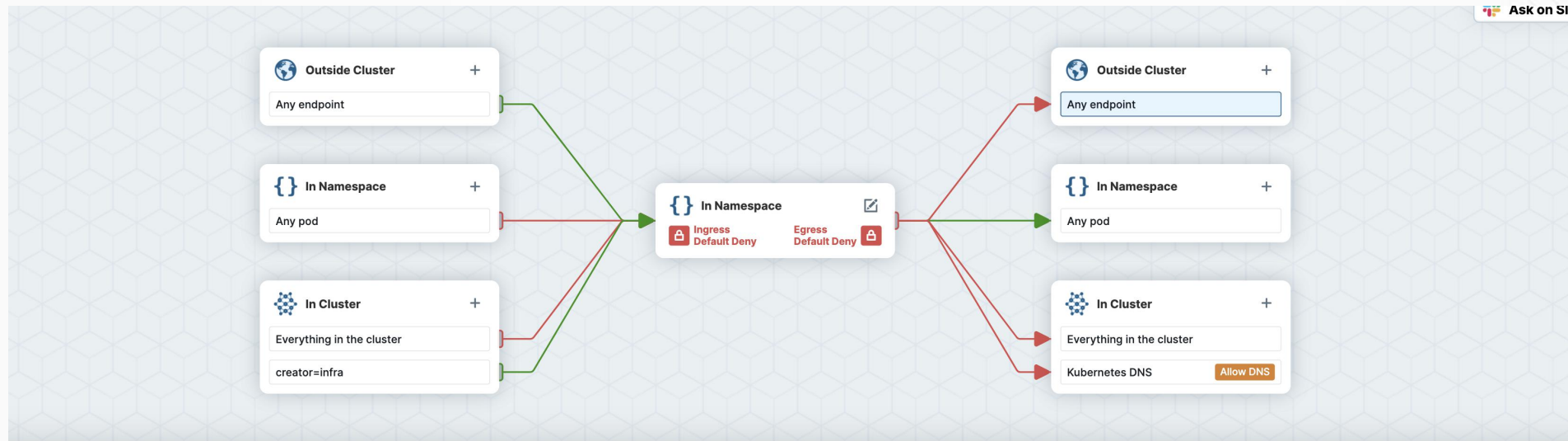
```
ingress:  
- from:  
  - namespaceSelector:  
    matchLabels:  
      user: alice  
  podSelector:  
    matchLabels:  
      role: client
```



## jobs-demo







```
apiVersion: networking.k8s.io/v1
kind: NetworkPolicy
metadata:
  name: untitled-policy
spec:
  podSelector: {}
  ingress:
    - from:
      - ipBlock:
          cidr: 0.0.0.0/0
      - namespaceSelector:
          matchLabels:
            creator: infra
  egress:
    - to:
      - podSelector: {}
```

Policy Rating ■■■■ Download [Download](#) [Share](#)

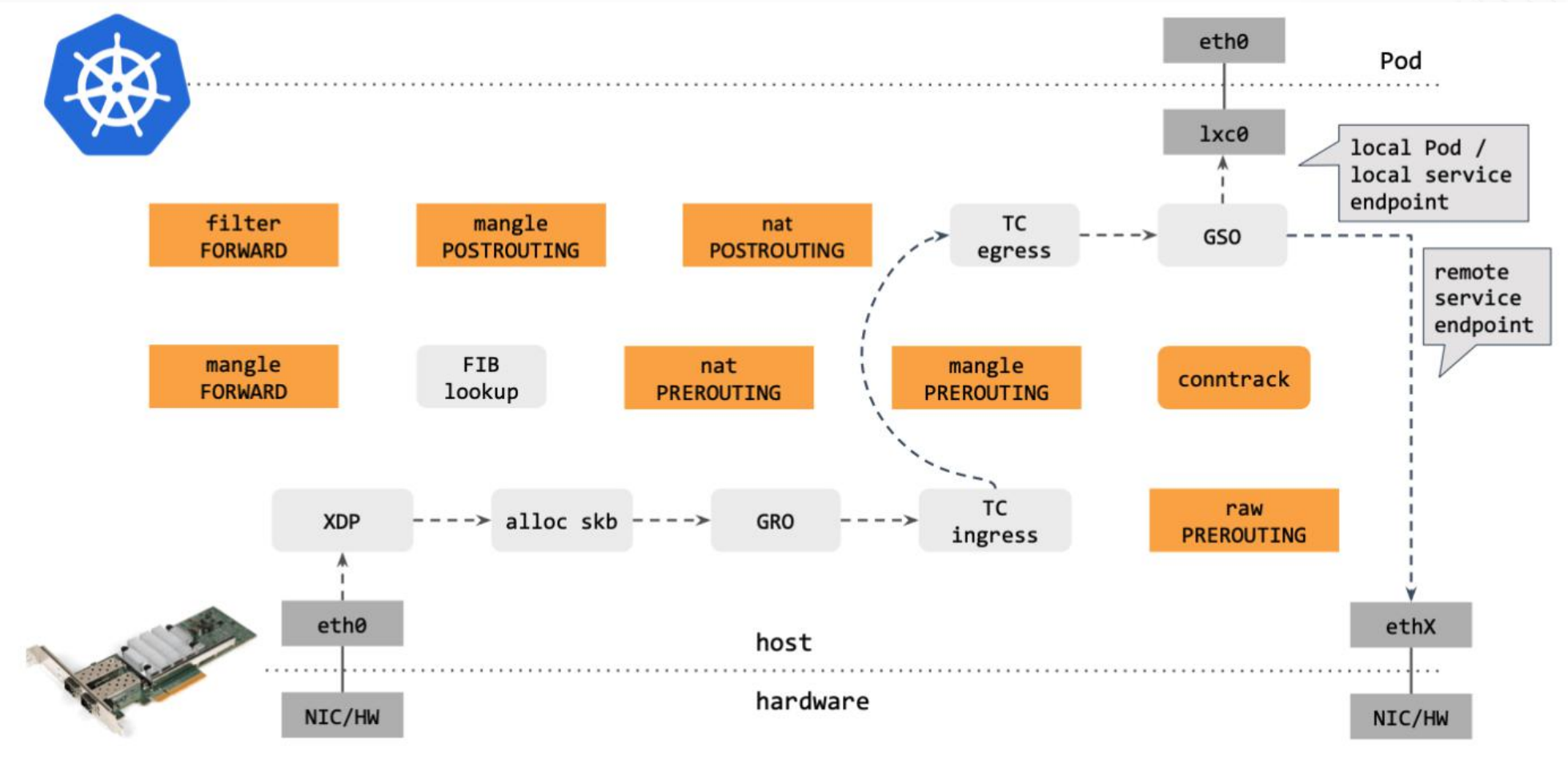
Main tutorial Flows upload

### Welcome to the Network Policy Editor! *Beta*

This tutorial will teach you how to create a network policy using Editor. It explains basic network policy concepts and guides you through the steps needed to achieve the desired least-privilege security and zero-trust concepts.

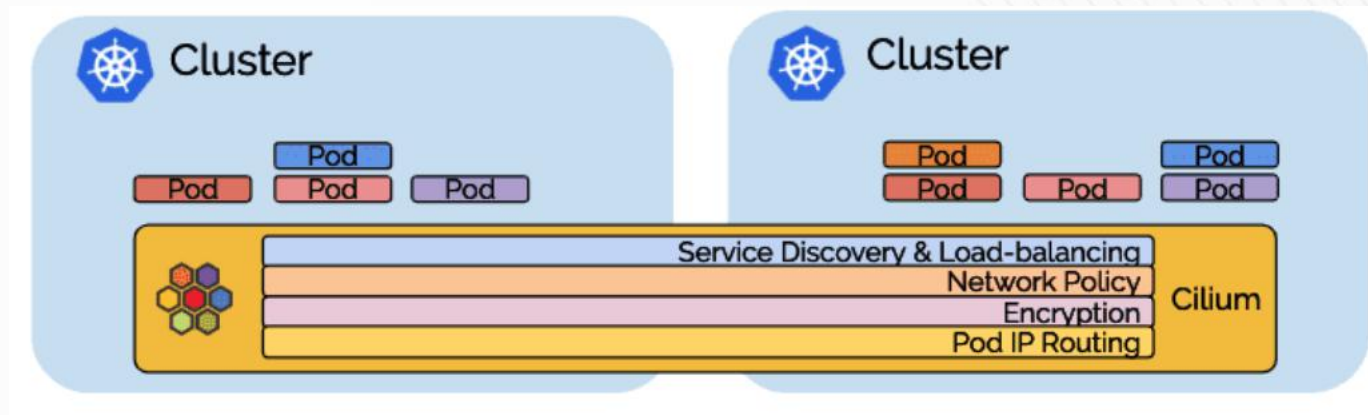
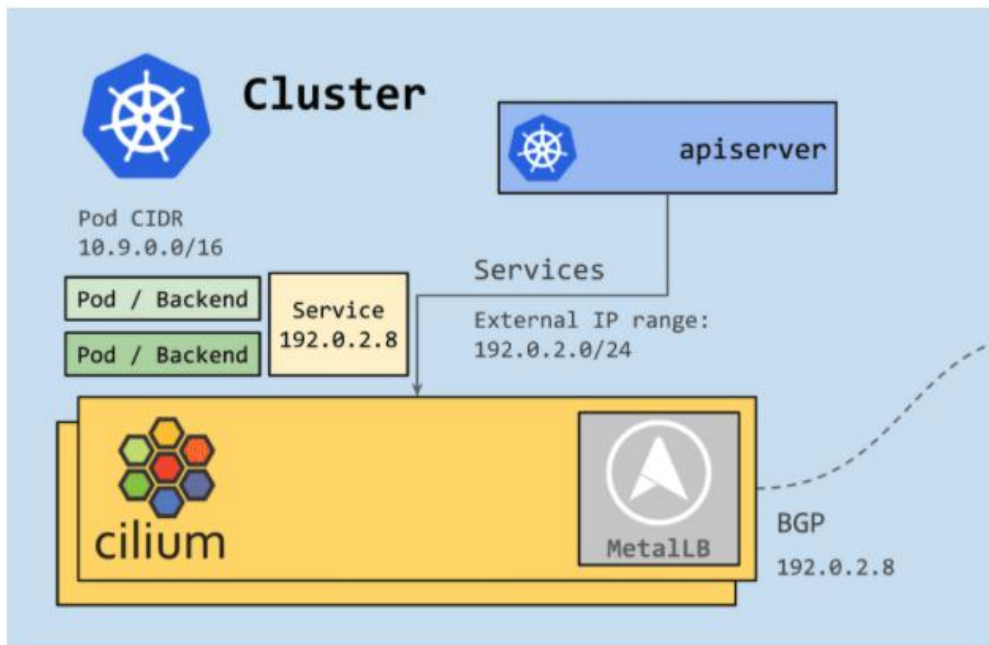
#### Step 1. What pods do you want to secure?

First, select the pods to which the policy should be applied by



## L4 LB

## ClusterMesh



# 参考资料

GOTC

<https://github.com/icopy-site/awesome-cn/blob/master/docs/awesome/awesome-ebpf.md>  
<http://machinezone.github.io/research/networking-solutions-for-kubernetes/>  
[https://github.com/iovisor/bpf-docs/blob/master/Express\\_Data\\_Path.pdf](https://github.com/iovisor/bpf-docs/blob/master/Express_Data_Path.pdf)  
<https://cilium.io/blog/2020/11/10/ebpf-future-of-networking/>  
<https://www.kernel.org/doc/html/latest/networking/filter.html>  
<https://www.tcpdump.org/papers/bpf-usenix93.pdf>  
<https://github.com/sbueringer/kubecon-slides>  
<https://developer.aliyun.com/article/779357>  
<https://nakryiko.com/posts/bpf-ringbuf/>  
<https://github.com/weaveworks/scope>  
<https://github.com/pixie-labs/pixie>  
<https://lwn.net/Articles/132196/>  
<https://github.com/cilium/cilium>  
<https://github.com/draios/sysdig>  
<https://github.com/iovisor/bcc>  
<https://ebpf.io>

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

**GOTC**

**THANKS**

**全球开源技术峰会**

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE