

Kekuatan Super eBPF

Giri Kuncoro

Senior Software Engineer at **Gojek**

Jakarta, August 21, 2021

Platinum sponsor :



Gold sponsor :



Silver sponsor :



Custom sponsor :



About Me



Giri Kuncoro

Senior Software Engineer at Gojek



Foundation sponsor:



Hosted by:



OpenStack Indonesia
Indonesia OpenStack Foundation Community
www.openstack.id



eBPF

 **INDONESIA**
OpenInfra Days

 **Open Infrastructure**
FOUNDATION

 **GiriKuncoro**

 **BRI**

 **EasyStack**
DATA • HYPEROPERATING

 **Biznet**
GioCloud

 **boer**
technology

 **NVIDIA.**



eBPF

“Superpowers have finally come to Linux”

- Brendan Gregg, Netflix



eBPF

“Superpowers have finally come to Linux”

“eBPF does to Linux what Javascript does to HTML”

- Brendan Gregg, Netflix



eBPF

**Run code in the kernel
without having to write a kernel module**



man bpf

man bpf

The `bpf()` system call performs a **range of operations** related to extended Berkeley Packet Filters. Extended BPF (or eBPF) is similar to the original ("classic") BPF (cBPF) used to filter network packets.

man bpf

The `bpf()` system call performs a **range of operations** related to extended Berkeley Packet Filters. Extended BPF (or eBPF) is similar to the original ("classic") BPF (cBPF) used to filter network packets.

For both cBPF and eBPF **programs**,

man bpf

The `bpf()` system call performs a **range of operations** related to extended Berkeley Packet Filters. Extended BPF (or eBPF) is similar to the original ("classic") BPF (cBPF) used to filter network packets.

For both cBPF and eBPF **programs**, the kernel statically analyzes the programs before loading them, in order to ensure that they **cannot harm the running system**.

man bpf

The `bpf()` system call performs a **range of operations** related to extended Berkeley Packet Filters. Extended BPF (or eBPF) is similar to the original ("classic") BPF (cBPF) used to filter network packets.

For both cBPF and eBPF **programs**, the kernel statically analyzes the programs before loading them, in order to ensure that they **cannot harm the running system**.

eBPF extends cBPF in multiple ways, including the ability to call a fixed set of **in-kernel helper functions** (via the `BPF_CALL` opcode extension provided by eBPF)

man bpf

The `bpf()` system call performs a **range of operations** related to extended Berkeley Packet Filters. Extended BPF (or eBPF) is similar to the original ("classic") BPF (cBPF) used to filter network packets.

For both cBPF and eBPF **programs**, the kernel statically analyzes the programs before loading them, in order to ensure that they **cannot harm the running system**.

eBPF extends cBPF in multiple ways, including the ability to call a fixed set of **in-kernel helper functions** (via the `BPF_CALL` opcode extension provided by eBPF)

and **access shared data structures such as eBPF maps**.

man bpf

eBPF programs can be written in a **restricted C** that is compiled (using the **clang** compiler) into **eBPF bytecode**. Various features are omitted from this restricted C, such as loops, global variables, variadic functions, floating-point numbers, and passing structures as function arguments.



clang & LLVM

The LLVM Project is a collection of modular and reusable **compiler and toolchain technologies**. Despite its name, LLVM has little to do with traditional virtual machines. The name "LLVM" itself is not an acronym; it is the full name of the project.

clang & LLVM

The LLVM Project is a collection of modular and reusable **compiler and toolchain technologies**. Despite its name, LLVM has little to do with traditional virtual machines. The name "LLVM" itself is not an acronym; it is the full name of the project.

Clang is an "LLVM native" C/C++/Objective-C compiler, which aims to deliver amazingly fast compiles.

man bpf

The kernel contains a just-in-time (JIT) compiler that **translates** **eBPF bytecode into native machine code** for better performance.



Writing Hello World





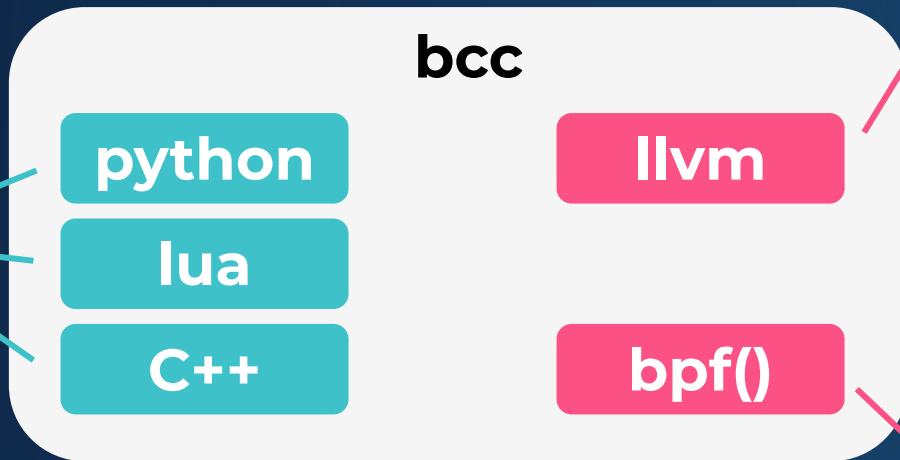
bcc

BCC **makes BPF programs easier to write**, with kernel instrumentation in C (and includes a C wrapper around LLVM), and frontends in Python and Lua.

 **bcc**

Compiles eBPF program

Language support



Wrapper for bpf() syscalls

Attaching eBPF programs to events



Triggering eBPF programs

eBPF programs can be attached to different events

- ▶ **Kprobes**
- ▶ **Uprobes**
- ▶ **Tracepoints**
- ▶ **Network packets**
- ▶ **Perf events**

bcc function names

```
b = BPF(text="""  
    int kprobe__sys_clone(void *ctx) {  
        bpf_trace_printk("Hello, OpenInfra!\\n");  
        return 0;  
    }  
""")
```

b.trace_print()



eBPF Maps



eBPF maps

Maps are a **generic data structure for storage of different types of data**. They allow sharing of data between eBPF kernel programs, and also between kernel and user-space applications.

Each map type has attributes:

- type
- max number of elements
- etc.

eBPF maps

Maps are a **generic data structure for storage of different types of data**. They allow sharing of data between eBPF kernel programs, and also between kernel and user-space applications.

Each map type has attributes:

- **type** →
- **max number of elements**
- **etc.**

BPF_MAP_TYPE_UNSPEC
BPF_MAP_TYPE_HASH
BPF_MAP_TYPE_ARRAY
BPF_MAP_TYPE_PROG_ARRAY
BPF_MAP_TYPE_PERF_EVENT_ARRAY
BPF_MAP_TYPE_PERCPU_HASH
BPF_MAP_TYPE_PERCPU_ARRAY
BPF_MAP_TYPE_STACK_TRACE
BPF_MAP_TYPE_CGROUP_ARRAY
BPF_MAP_TYPE_LRU_HASH

```
clang -O2 -emit-llvm -c bpf.c -o - | llc -march=bpf -filetype=obj -o bpf.o
```

(limited) C

ELF object file

- eBPF opcodes
- eBPF maps

ELF object file

- eBPF opcodes
- eBPF maps

user space

`bpf()` system calls

kernel

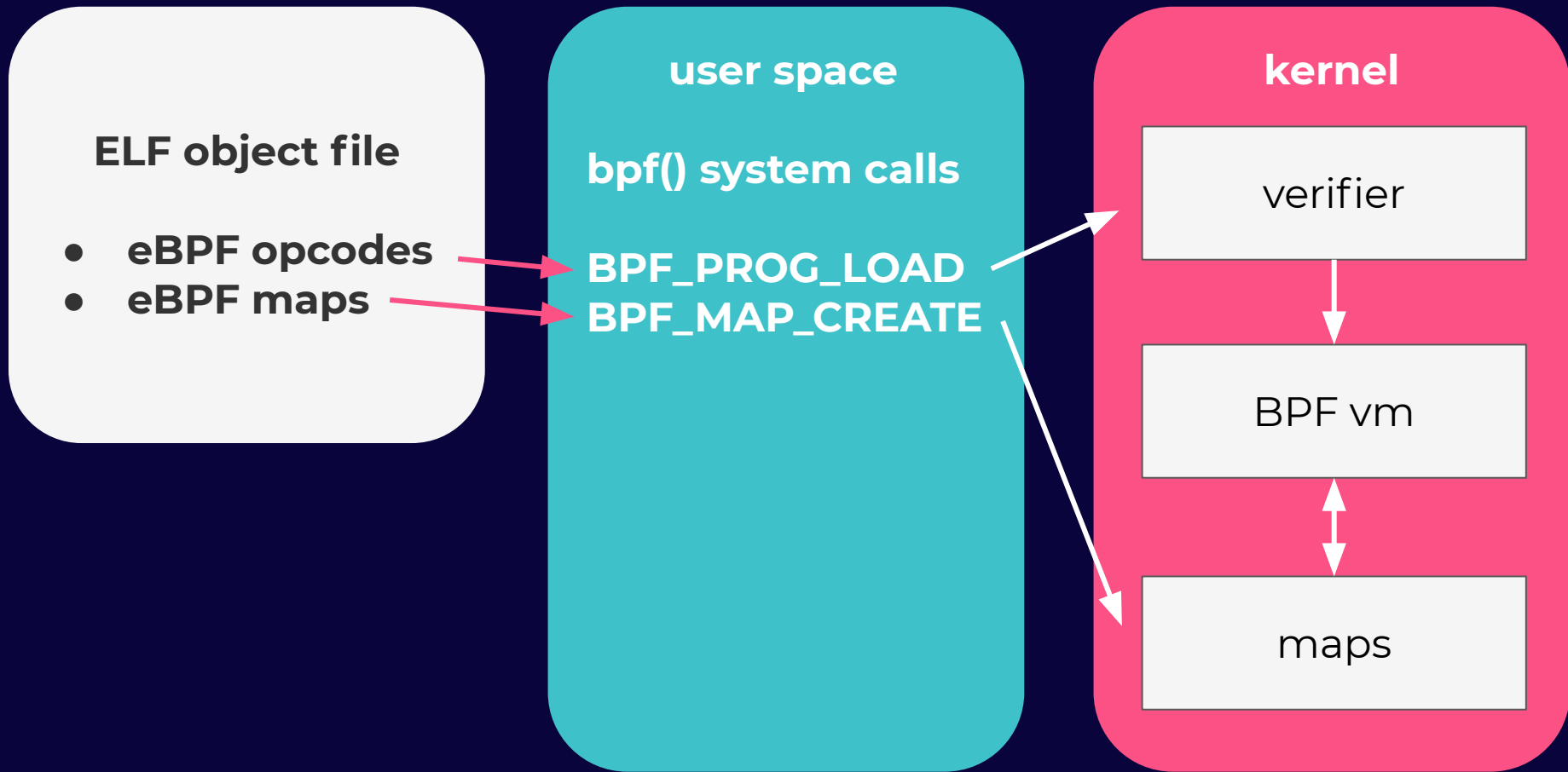
verifier



BPF vm



maps



ELF object file

- eBPF opcodes
- eBPF maps

user space

`bpf()` system calls

`BPF_PROG_LOAD`
`BPF_MAP_CREATE`

kernel

verifier

BPF vm

maps

ELF object file

- eBPF opcodes
- eBPF maps

user space

bpf() system calls

**BPF_PROG_LOAD
BPF_MAP_CREATE**

*Attach BPF
program to event*

kernel

verifier



BPF vm



maps

ELF object file

- eBPF opcodes
- eBPF maps

user space

bpf() system calls

**BPF_PROG_LOAD
BPF_MAP_CREATE**

*Attach BPF
program to event*

Read / write maps

**BPF_MAP_GET_NEXT_KEY
BPF_MAP_LOOKUP_ELEM
BPF_MAP_UPDATE_ELEM
BPF_MAP_DELETE_ELEM**

kernel

verifier



BPF vm



maps



eBPF helper functions

These helpers are used by eBPF programs to **interact with the system, or with the context in which they work**. For instance, they can be used to print debugging messages, to get the time since the system was booted, to interact with eBPF maps, or to manipulate network packets.

eBPF helper functions

These helpers are used by eBPF programs to **interact with the system, or with the context in which they work**. For instance, they can be used to print debugging messages, to get the time since the system was booted, to interact with eBPF maps, or to manipulate network packets.

```
bpf_trace_printk()  
bpf_map_*_elem()  
bpf_get_current_pid_tgid()
```

Verifier

Each eBPF program is a set of instructions that is **safe to run until its completion**. An **in-kernel verifier** statically determines that the eBPF program terminates and is safe to execute.

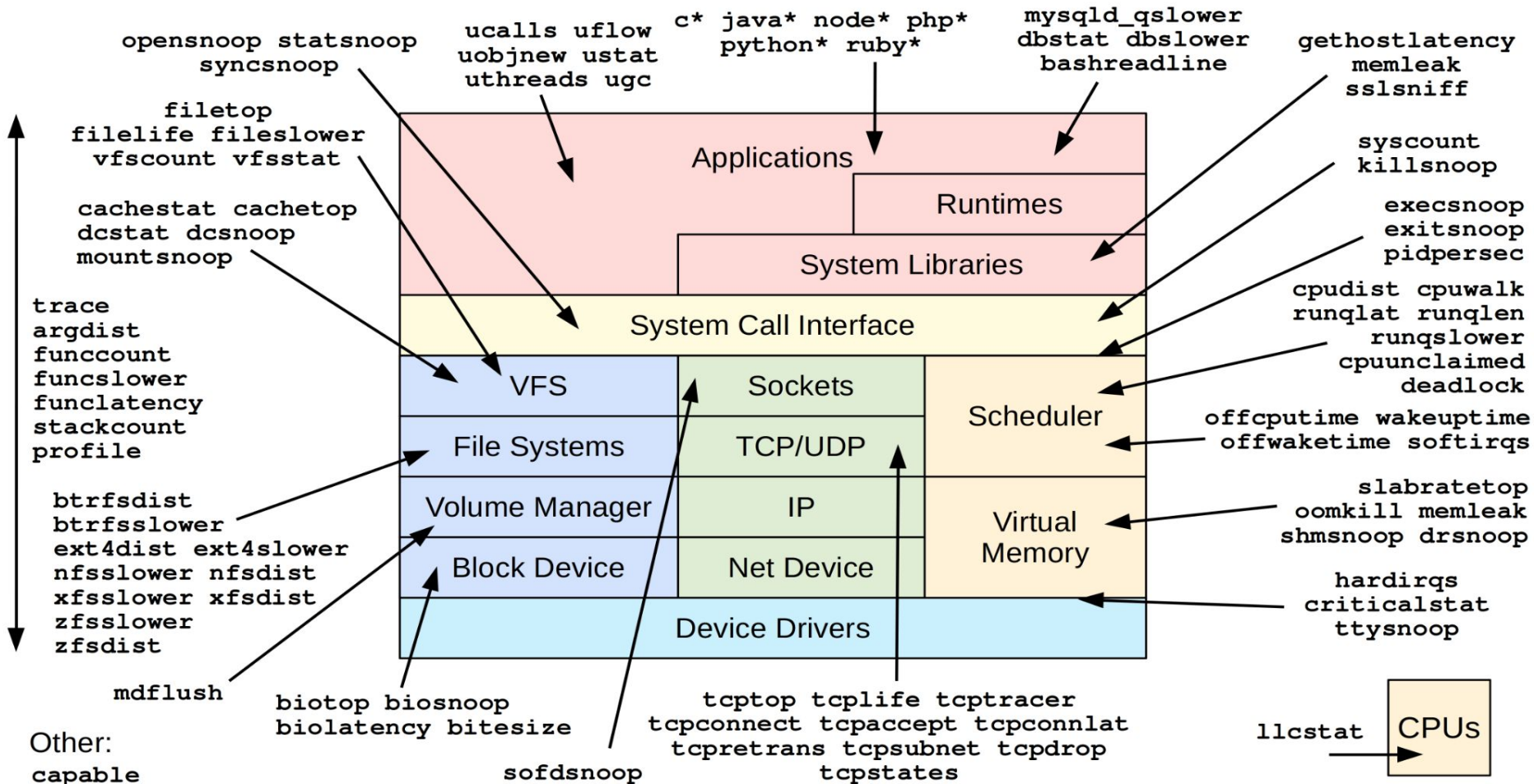
- No loops
- No bad pointer dereferences
- Restricted program size
- Always exits

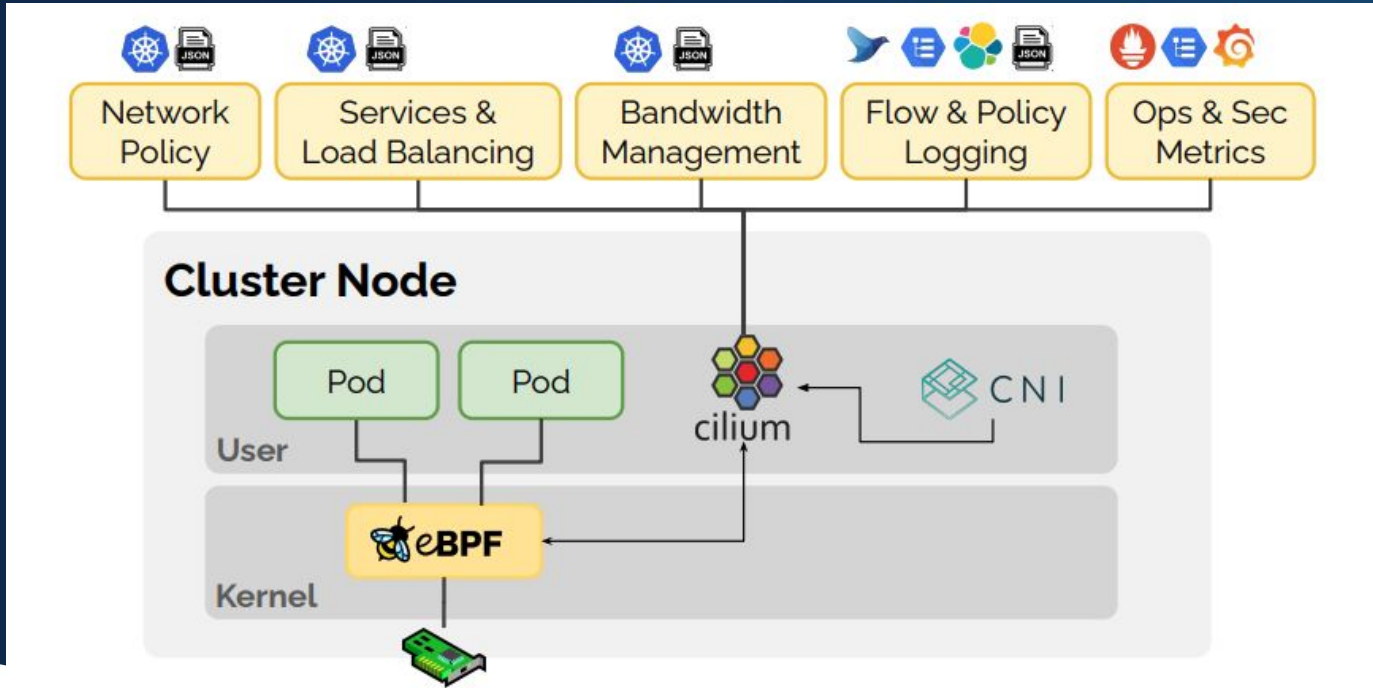
BPF Tools

Kekuatan super untuk semua



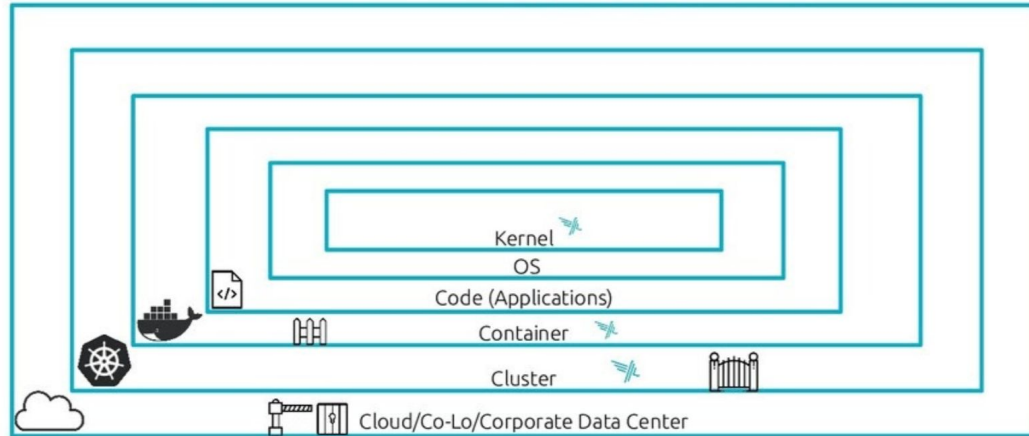
Linux bcc/BPF Tracing Tools





Prevention is not enough.

Combine with runtime detection tools. Use a *defense-in-depth* strategy.



@leodigo



Script: px/namespace --start_time=-5m --namespace=px-sock-shop
Live View: https://withpixie.ai:443/live/clusters/gke_pl-dev-infra_us-west1-a_zasgar-demo/script?namespace=px-sock-shop&script=px%2Fnamespace&start_time=-5m

work.withpixie.ai/live/clusters/gke_zasgar-demo/namespaces/px-sock-shopcon?start_time=-5m

px-sock-shop

Namespace Service Graph

Pod List

POD	RSS	VSZ
px-sock-shop/carts-75fdd5984-mblpv	171.86 MB	921.34 MB
px-sock-shop/carts-db-567bc578f-47fx	332.11 MB	1.86 GB
px-sock-shop/catalogue-64659d9654-zz894	12.71 MB	18.4 MB
px-sock-shop/catalogue-db-66dc796b66-v4zxx	250.53 MB	2.53 GB
px-sock-shop/front-end-786597bb5f-w99lm	67.53 MB	427.03 MB
px-sock-shop/load-test-55c4d5647f-2zb6s	2.77 GB	3.88 GB
px-sock-shop/orders-64cb96fcb9-flwb9	162.53 MB	895.41 MB
px-sock-shop/orders-db-5cfc884cf-zwzqc	389.62 MB	1.74 GB

Service List

SERVICE	LATENCY_MS (P99)	REQUESTS_PER_S	ERROR_RATE_PCT	INBOUND_BYTES
px-sock-shop/carts	287.47	38.37	0.25	
px-sock-shop/catalogue	67.24	25.59	0	
px-sock-shop/front-end	1605.52	105.23	0.22	
px-sock-shop/orders	212.72	12.75	1.38	
px-sock-shop/payment	1.78	12.80	0	
px-sock-shop/shipping	5.11	11.03	0	
px-sock-shop/worker	8.17	RR K4	n	

1 Outgoing Traffic 2 State

10:42
work.withpixie.ai

px-sock-shop

Namespace Service Graph

Service List

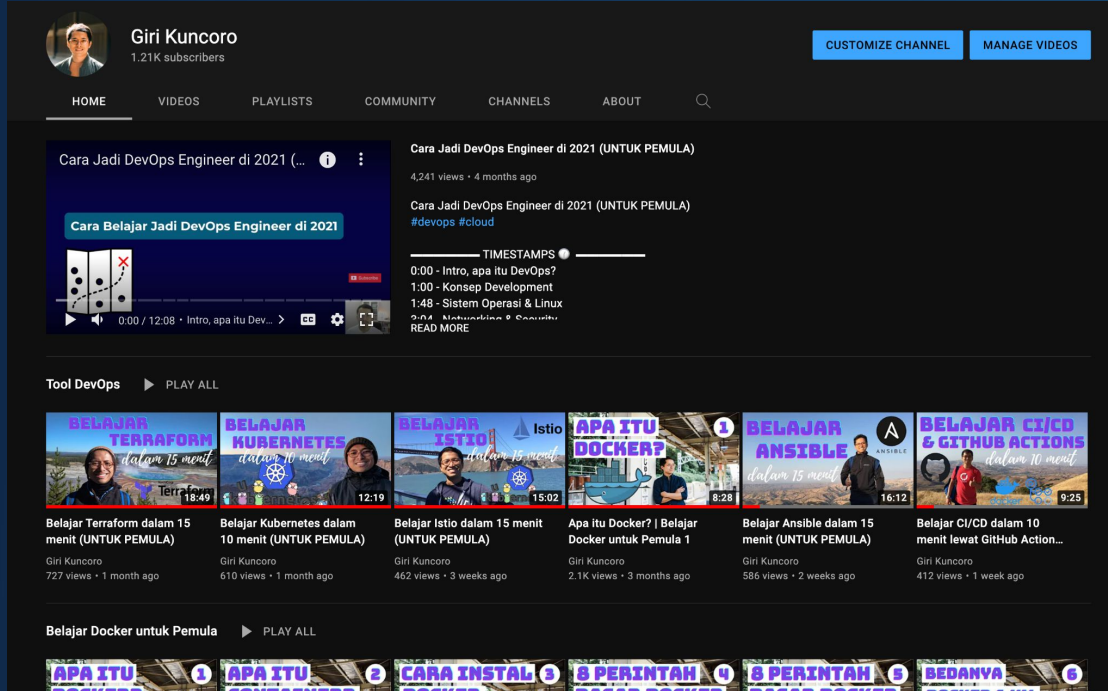
Service List Namespace Service Graph PIXIE

References

IO Visor Project (iovisor.org)

Brendan Gregg (brendangregg.com)

Youtube Channel



The screenshot shows the YouTube channel page for Giri Kuncoro, who has 1.21K subscribers. The channel is categorized under 'Tool DevOps' and features a video player for 'Cara Jadi DevOps Engineer di 2021 (UNTUK PEMULA)' with 4,241 views. Below the video player is a list of related videos, including 'Belajar Terraform dalam 15 menit (UNTUK PEMULA)', 'Belajar Kubernetes dalam 10 menit (UNTUK PEMULA)', 'Belajar Istio dalam 15 menit (UNTUK PEMULA)', 'Apa itu Docker? | Belajar Docker untuk Pemula 1', 'Belajar Ansible dalam 15 menit (UNTUK PEMULA)', and 'Belajar CI/CD dalam 10 menit lewat GitHub Action...'. The channel also has a 'PLAY ALL' button for a playlist of DevOps tutorials.

Sponsored by:



Open Infrastructure
FOUNDATION



NVIDIA[®]



boer
technology



Biznet
GioCloud



EasyStack
open cloud computing

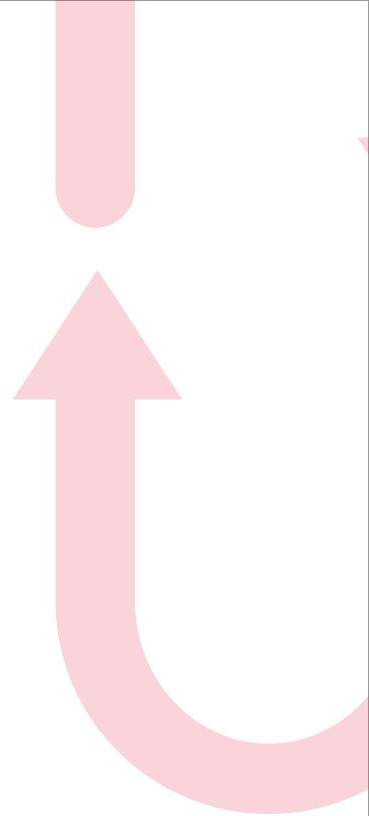


intek

IND  **CENTER**



Open Networking
Indonesia



Hosted by:



OpenStack Indonesia

Indonesia OpenStack Foundation Community
www.openstack.id

Community Partners:



Thanks!

Do you have any questions?

 @girikuncoro

 GiriKuncoro

Platinum sponsor :



Gold sponsor :



Silver sponsor: Custom sponsor :

