

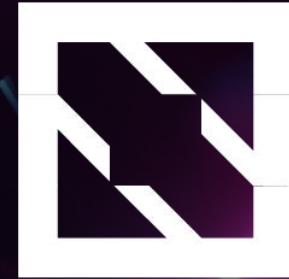


KubeCon

THE LINUX FOUNDATION



China 2024



CloudNativeCon





KubeCon



CloudNativeCon



China 2024

Unlocking LLM Performance with : Optimizing Training and Inference Pipelines

向阳，云杉网络

Outline



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



Open Source Dev & ML Summit

China 2024

1. 背景：训练和推理的效率挑战
2. 现状：传统解决方案和工具的问题
3. 方法：eBPF 构建零侵扰可观测性
4. 实践：PyTorch 全栈剖析和追踪

LLM 训练开销大、效率低



KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMIT

China 2024

	GPT-4	Llama-3.1
Size	1.8T	405B
GPU	25K A100	16K H100
Days	90~100	54
MFU	32%~36%	38%~43%

- [Everything We Know About GPT-4 - Klu.ai](#)
- [GPT4- All Details Leaked](#)
- [The Llama 3 Herd of Models](#)

训练时间长：数月

GPU 数量多：数万

模型参数大：万亿

GPU 利用率低：~40%

Component	Category	Interruption Count	% of Interruptions
Faulty GPU	GPU	148	30.1%
GPU HBM3 Memory	GPU	72	17.2%
Software Bug	Dependency	54	12.9%
Network Switch/Cable	Network	35	8.4%
Host Maintenance	Unplanned Maintenance	32	7.6%
GPU SRAM Memory	GPU	19	4.5%
GPU System Processor	GPU	17	4.1%
NIC	Host	7	1.7%
NCCL Watchdog Timeouts	Unknown	7	1.7%
Silent Data Corruption	GPU	6	1.4%
GPU Thermal Interface + Sensor	GPU	6	1.4%
SSD	Host	3	0.7%
Power Supply	Host	3	0.7%
Server Chassis	Host	2	0.5%
IO Expansion Board	Host	2	0.5%
Dependency	Dependency	2	0.5%
CPU	Host	2	0.5%
System Memory	Host	2	0.5%

Table 5 Root-cause categorization of unexpected interruptions during a 54-day period of Llama 3 405B pre-training. About 78% of unexpected interruptions were attributed to confirmed or suspected hardware issues.

$$\frac{148}{54*365/16384} = 6\%$$
$$\frac{(148+72+19+17+6+6)}{54*365/16384} = 11\%$$

GPU 年化故障高：
6%~11%

故障并不是训练低效的唯一原因



KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMIT

China 2024

Classification of the 706 low-GPU-utilization issues discovered across 400 deep learning jobs.

Category	Description	No.	Ratio
Interactive Job	The execution of a job entails regular interaction with its owner.	15	2.12%
GPU Oversubscription	A job requests more GPUs than it actually utilizes.	6	0.85%
Unreleased Job	A job does not terminate promptly after completing its computation.	9	1.27%
Non-DL Job	A job is unrelated to deep learning and does not utilize GPUs at all. For example, the job performs data analysis using CPUs solely.	4	0.57%

Improper Batch Size	Improper values of the batch size are used, which decrease the GPU computation of deep learning operators.	181	25.64%
Insufficient GPU Memory	The GPU memory is not sufficient to support more GPU computation.	22	3.12%
Model Checkpointing	A job saves model checkpoints synchronously to the distributed data store.	116	16.43%

Inefficient Host-GPU Data Transfer	The data transfer between main memory and GPU memory is not efficient.	197	27.90%
Data Preprocessing	Raw input data is preprocessed using CPUs before model training.	28	3.97%
Remote Data Read	A job opens and reads input data directly from the distributed data store.	18	2.55%
External Data Usage	A job accesses input data or model files directly from external sites.	18	2.55%
Intermediate Result Upload	A job saves intermediate training results synchronously to the distributed data store.	14	1.98%
Data Exchange	The GPUs of a distributed job continually exchange data, such as gradients and output tensors, among one another.	50	7.08%

Long Library Installation	The installation of dependent libraries takes too much time (at least 10 minutes).	12	1.70%
API Misuse	API usage violates assumptions.	16	2.27%

- Yanjie Gao (Microsoft Research) et al, ACM ICSE 2024,
[An Empirical Study on Low GPU Utilization of Deep Learning Jobs.](#)
- Yanjie Gao (Microsoft Research) et al, ACM ICSE 2023,
[An Empirical Study on Quality Issues of Deep Learning Platform.](#)

如何知晓你的训练任务是否存在这些问题？

内核计算

```
1 from torch.utils.data import DataLoader
2
3 train_batch_size = 32 384
4 eval_batch_size = 32 448
5 train_loader = DataLoader(dataset = train_data, batch_size =
6   ↪ train_batch_size, shuffle = True)
6 eval_loader = DataLoader(dataset = eval_data, batch_size =
7   ↪ eval_batch_size, shuffle = True)
```

Figure 2: A simplified example of Improper Batch Size issues. “train_batch_size” and “eval_batch_size” are specified arguments for model training and evaluation, respectively. The fix is to increase their values independently (lines 3–4).

显存拷贝

```
1 import torch
2 from torch.utils.data import DataLoader
3
4 train_loader = DataLoader(train_set, ..., num_workers=8,
5   ↪ pin_memory=True)
5 eval_loader = DataLoader(eval_set, ..., num_workers=8,
6   ↪ pin_memory=True)
7
8 for epoch in range(num_epochs):
9   ...
10  for _, data in enumerate(train_loader, 0):
11    # get the inputs
12    inputs, labels = data
13    inputs = inputs.to(device, non_blocking=True)
14    labels = labels.to(device, non_blocking=True)
```

Figure 4: A simplified example of Inefficient Host-GPU Data Transfer issues. The fix is to enable automatic memory pinning by setting the pin_memory parameter to True (lines 3–4). We further set the non_blocking parameter to True (lines 11–12), which tries asynchronous data transfer if possible.

In a distributed job, GPUs consistently share data such as gradients and output tensors. This data exchange between GPUs, whether through the network or PCI Express bus, frequently interrupts their designated tasks and causes a sudden drop in GPU utilization to zero. Within the Data Exchange category, there are 50 (7.08%) issues. A common fix is to enhance communication efficiency by minimizing the frequency of data exchange (e.g., using large batch sizes) and enabling compressed communication [34, 39]. For users of Horovod [67], enlarging the backward_passes_per_step parameter⁵ can help accumulate more local gradient updates and transmit them simultaneously. Developers can also leverage Microsoft DeepSpeed’s 3D (data, model, and pipeline) parallelism, whose 1-bit Adam and 0/1 Adam [39] optimizers demonstrate significant reductions in communication volume.

```
1 import torch
2 import threading
3 import shutil
4 import os
5
6 def main():
7   ...
8   for epoch in range(start_epoch, num_epochs):
9     if step > num_training_steps:
10       break
11
12     for i, batch in enumerate(tqdm(train_loader)):
13       ...
14       logits = model(input_data)
15
16       optimizer.zero_grad()
17       loss.backward()
18       optimizer.step()
19
20       f1, pred = evaluator.evaluate(val_loader, model, step)
21
22       if f1 > max_f1:
23         max_f1 = f1
24         torch.save(model, remote_path)
25         local_path = make_tmp_path(epoch, i)
26         torch.save(model, local_path)
27
28       + def save_file(local_path, remote_path):
29         + shutil.copyfile(local_path, remote_path)
30         + os.remove(local_path)
31
32       + threading.Thread(target = save_file, args = [local_path,
33         ↪ remote_path]).start()
34
35 if __name__ == '__main__':
36   main()
```

Figure 3: A simplified example of Model Checkpointing issues. The fix is to save the checkpoint to a local temporary file (lines 25–26), followed by an asynchronous copy to the remote data store in a separate thread (lines 28–32).

网络传输

LLM 推理开销大、时延高



China 2024

Llama	8B	70B	405B
FP32	36GB	267GB	1.48TB
FP16	20GB	135GB	758GB
INT8	12GB	70GB	382GB
INT4	8GB	37GB	193GB

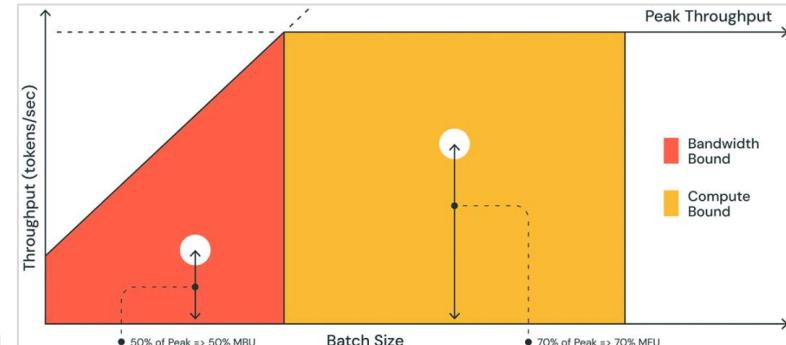
LLM Memory Requirements

Inference Training

Total Inference Memory: 1.48 TB

- Model Weights: 1.47 TB
 - KV Cache: 2.00 GB
 - Activation Memory: 3.56 GB
-
- [LLM Memory Requirements](#)
 - [LLM Inference Performance Engineering: Best Practices](#)

80GB: 1 GPU
640GB: 1 Node x 8 GPU
1.28TB: 2 Node x 8 GPU

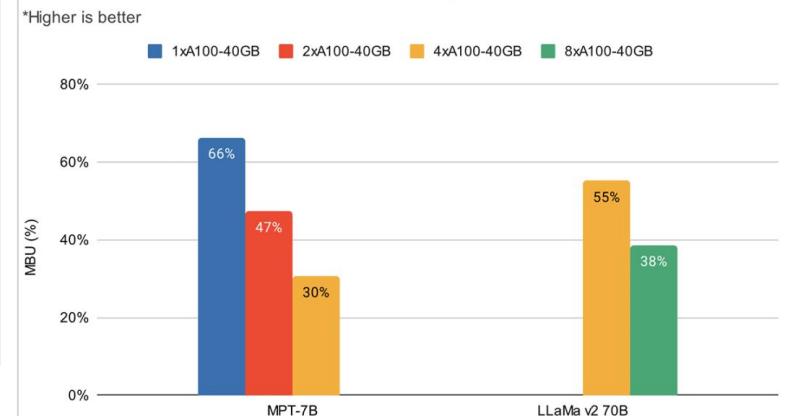


Time To First Token (TTFT)
Time Per Output Token (TPOT)
Model Bandwidth Utilization (MBU)

GPU 并非越少越好

GPU 越少，则每个 GPU 需要加载更多的模型参数。

Observed MBU across various tensor-parallelized modes



GPU 并非越多越好

GPU 越多，则通信越复杂，内存碎块越多。

没有银弹，可观测性是优化的前提。



Sig

Help Needed from vLLM team on profiling pytorch cuda memory

- Biz
- vLLM
- PyTorch
- Python / C++



youkaichao1

1 Mar 16

Hi, I'm working on [GitHub - vllm-project/vllm: A high-throughput and memory-efficient inference and serving engine for LLMs](#) (7), and the recently release of pytorch 2.2.0 caused some trouble to me. I came here for help in profiling cuda memory usage.

The basic story: vLLM tries to allocate as much memory as possible for KV Cache to accelerate LLM inference. In order to do so, it first profiles the memory usage, guess the maximum size of memory available for KV Cache, and also leave some for storing activation during inference.

What's more, vLLM uses cuda graph to reduce Python overhead.

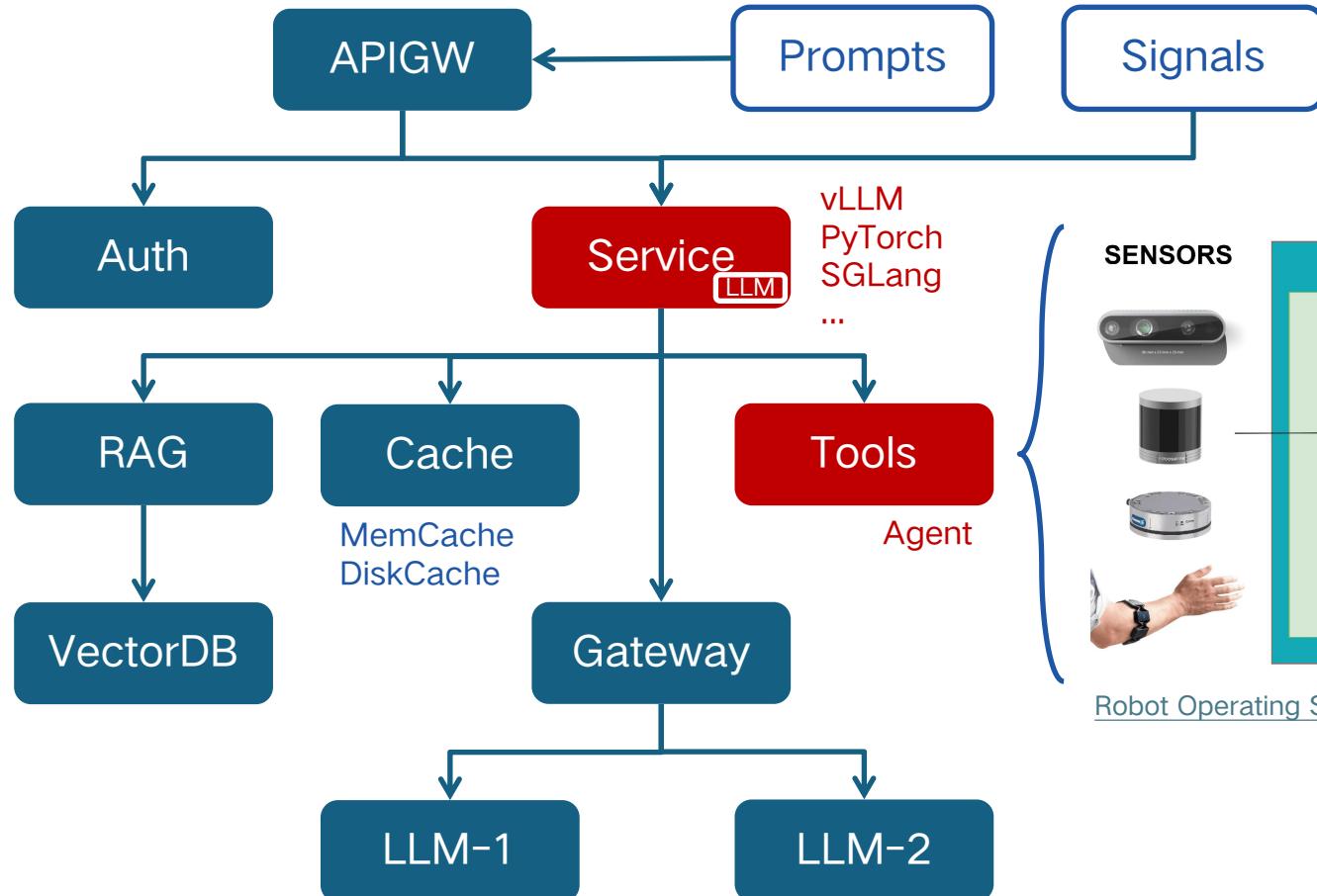
When PyTorch upgrades from 2.1.2 to 2.2.0 , there seems to be some internal change of memory allocator, and the memory that can be used is decreased. It can cause OOM error during inference.

Here are the diagnoses data (produced by `torch.cuda.memory_stats` and `torch.cuda.memory._dump_snapshot`), collected from a server with 2 L4 GPUs:

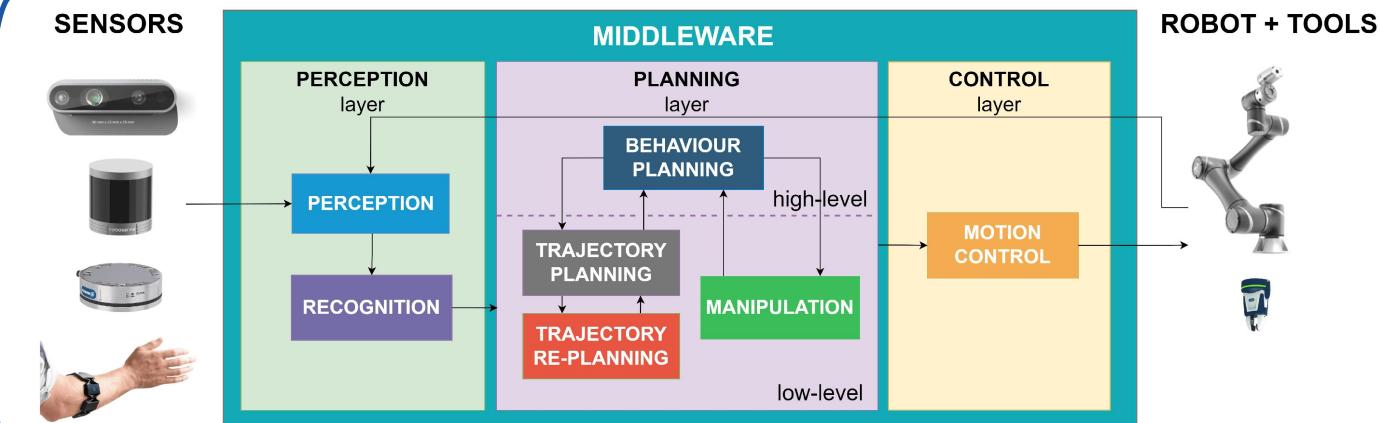
从推理应用到在线 LLM 推理服务



China 2024



Perception-Planning-Control
BEV-OCC-Transformer



端 - 自动驾驶、具身智能（ROS2）的
端到端低时延和高稳定性要求

云 - 在线推理服务是一个复杂的分布式服务
TTFT、TPOT、时延、吞吐

从大模型到小模型：消费级 GPU、CPU 协同



China 2024

Accelerating Model Training in Multi-cluster Environments with Consumer-grade GPUs

Hwijoon Lim
KAIST
Daejeon, Republic of Korea
hwijoon.lim@kaist.ac.kr

Juncheol Ye
KAIST
Daejeon, Republic of Korea
juncheol@kaist.ac.kr

Sangeetha Abdu Jyothi
UC Irvine & VMware Research
Irvine, California, USA
sangeetha.aj@uci.edu

Dongsu Han
KAIST
Daejeon, Republic of Korea
dhan.ee@kaist.ac.kr

[Accelerating Model Training in Multi-cluster Environments with Consumer-grade GPUs, SIGCOMM 2024.](#)

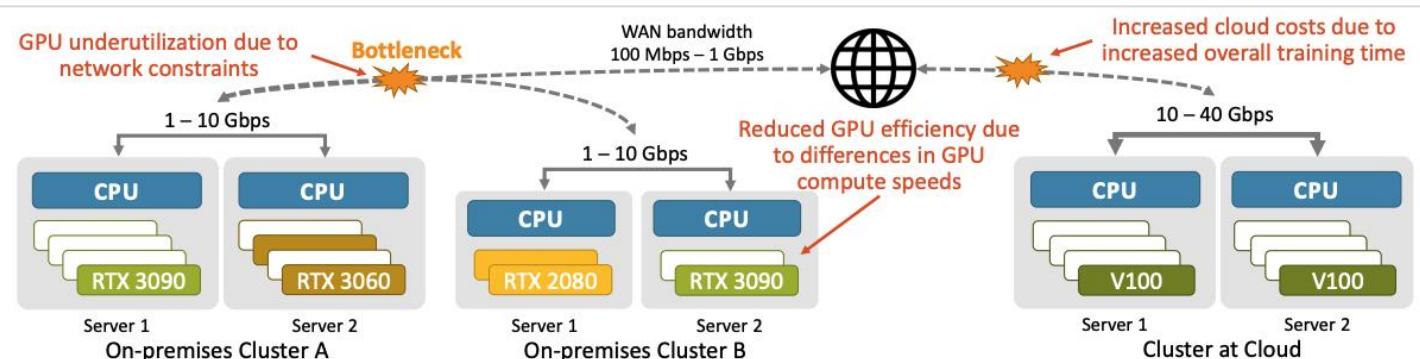
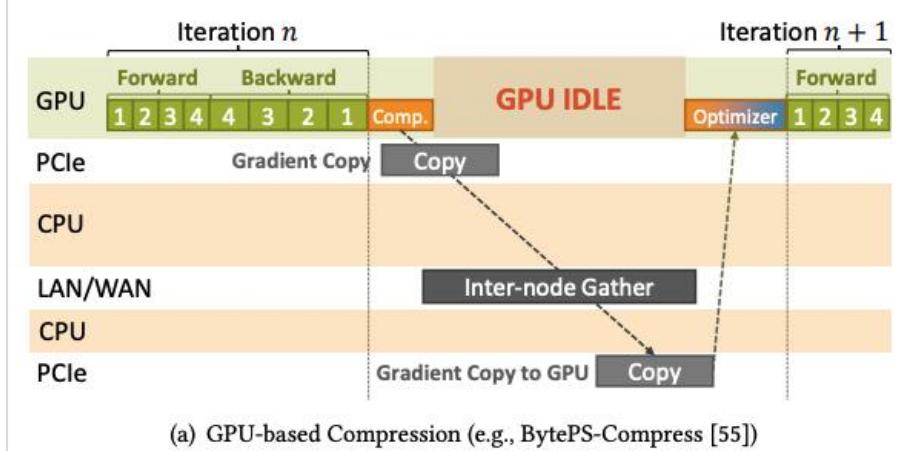
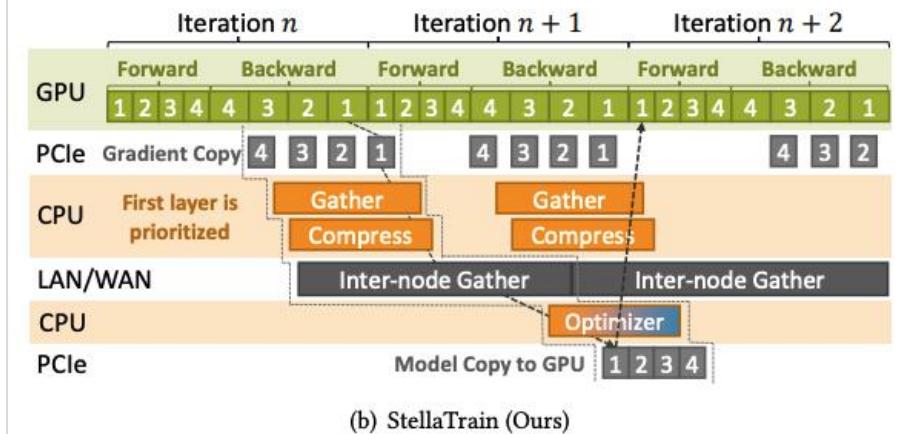


Figure 1: A multi-cluster environment with two on-premises lab clusters and a cloud cluster.



(a) GPU-based Compression (e.g., BytePS-Compress [55])



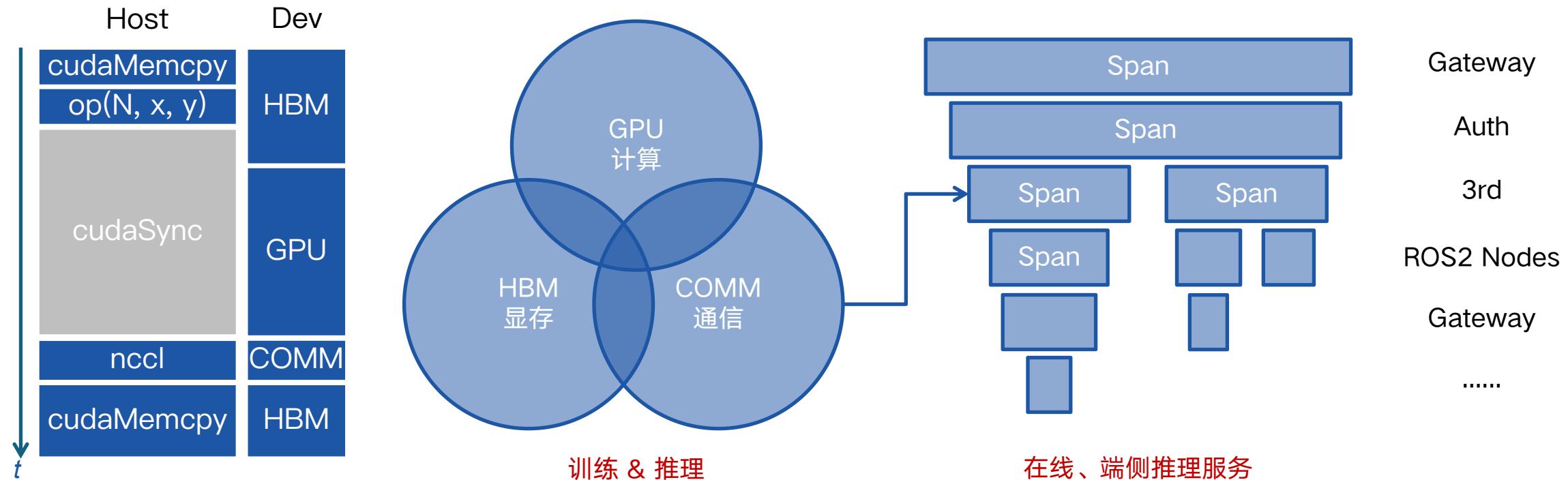
(b) StellaTrain (Ours)

Figure 4: Comparison of training pipelines.

AI 训练和推理的可观测性需求



China 2024



Outline



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



Open Source Dev & ML Summit

China 2024

1. 背景：训练和推理的效率挑战
2. 现状：传统解决方案和工具的问题
3. 方法：eBPF 构建零侵扰可观测性
4. 实践：PyTorch 全栈剖析和追踪

DCGM Prometheus Exporter



China 2024

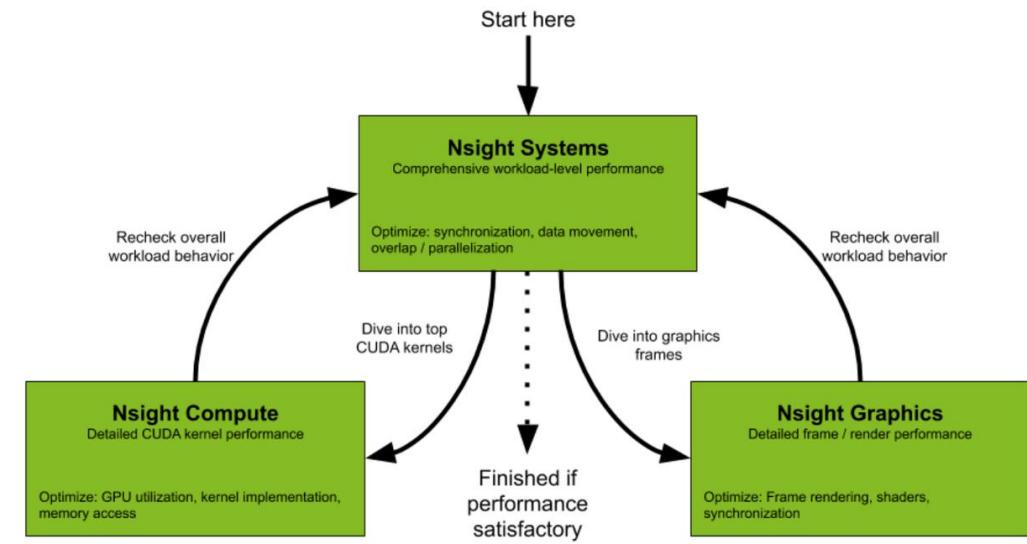


GPU: Nvidia Nsight、PyTorch Profiler



China 2024

Nsight packages



Nsight 的问题：

需要重启进程、缺少 CPU Context。

```
def trace_handler(p):
    output = p.key_averages().table(sort_by="self_cuda_time_total", row_limit=10)
    print(output)
    p.export_chrome_trace("/tmp/trace_" + str(p.step_num) + ".json")
```

```
with profile(
    activities=[ProfilerActivity.CPU, ProfilerActivity.CUDA],
    schedule=torch.profiler.schedule(
        wait=1,
        warmup=1,
        active=2),
    on_trace_ready=trace_handler
) as p:
    for idx in range(8):
        model(inputs)
        p.step()
```

1. Parameter `skip_first` tells profiler that it should ignore the first 10 steps (default value of `skip_first` is zero);
2. After the first `skip_first` steps, profiler starts executing profiler cycles;
3. Each cycle consists of three phases:
 - idling (`wait=5` steps), during this phase profiler is not active;
 - warming up (`warmup=1` steps), during this phase profiler starts tracing, but the results are discarded; this phase is used to discard the samples obtained by the profiler at the beginning of the trace since they are usually skewed by an extra overhead;
 - active tracing (`active=3` steps), during this phase profiler traces and records data;
4. An optional `repeat` parameter specifies an upper bound on the number of cycles. By default (zero value), profiler will execute cycles as long as the job runs.

PyTorch Profiler 的问题：

只能用于 PyTorch；性能影响大；需要改代码、重启进程。

需要手工精心打造
插桩、开销

RDMA 网络：网卡/交换机指标、拨测



KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI dev
Open Source Dev & ML Summit

China 2024

RDMA network behind AI Training Cluster

Enables two networked hosts to exchange data in main memory without relying on the processor(CPU)



SCALE AI
TRAINING
WORKLOADS

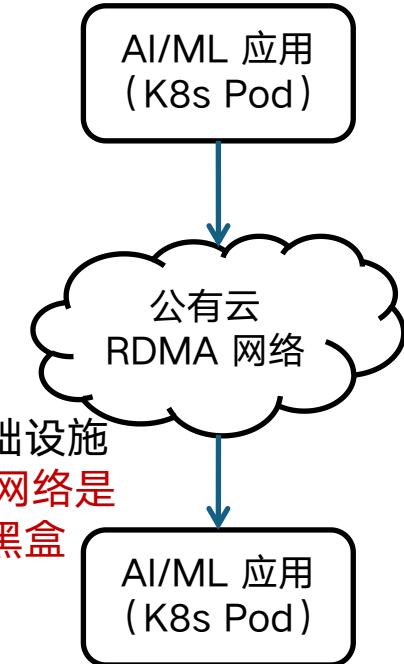


THE PICK
TO KEEP
GPUS BUSY



CAN AFFECT
TRAINING
EFFICIENCY

Meta: Network
Observability for
AI/HPC Training
Workflows
私有基础设施
网卡/交换机指标
粒度粗



公有基础设施
RDMA 网络是
性能黑盒

北京邮电大学 ByteDance 字节跳动 紫金山实验室

Beijing University of Posts and Telecommunications ByteDance Purple Mountain Laboratories

R-Pingmesh: A Service-Aware RoCE Network Monitoring and Diagnostic System

Kefei Liu, Zhuo Jiang, Jiao Zhang, Shixian Guo, Xuan Zhang, Yangyang Bai, Yongbin Dong, Feng Luo, Zhang Zhang, Lei Wang, Xiang Shi, Haohan Xu, Yang Bai, Dongyang Song, Haoran Wei, Bo Li, Yongchen Pan, Tian Pan, and Tao Huang

SIGCOMM 2024

ByteDance 字节跳动 北京邮电大学 紫金山实验室

Beijing University of Posts and Telecommunications Purple Mountain Laboratories

Hostping: Diagnosing Intra-host Network Bottlenecks in RDMA Servers

Kefei Liu, Zhuo Jiang, Jiao Zhang, Haoran Wei, Xiaolong Zhong, Lizhuang Tan, Tian Pan and Tao Huang

NSDI 2023

字节&北邮
聚焦 RDMA
网络主动拨测

在线推理服务的可观测性：分布式追踪



KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI dev
Open Source Dev & ML Summit

China 2024

```
Python TypeScript

from openai import OpenAI
from traceloop.sdk import Traceloop
from traceloop.sdk.decorators import workflow

Traceloop.init(app_name="joke_generation_service")

@workflow(name="joke_creation")
def create_joke():
    client = OpenAI(api_key=os.environ["OPENAI_API_KEY"])
    completion = client.chat.completions.create(
        model="gpt-3.5-turbo",
        messages=[{"role": "user", "content": "Tell me a joke about opentelemetry"}],
    )

    return completion.choices[0].message.content
```

OpenLLMetry

```
Python TypeScript

import openai
from langsmith.wrappers import wrap_openai
from langsmith import traceable

# Auto-trace LLM calls in-context
client = wrap_openai(openai.Client())

@traceable # Auto-trace this function
def pipeline(user_input: str):
    result = client.chat.completions.create(
        messages=[{"role": "user", "content": user_input}],
        model="gpt-3.5-turbo"
    )
    return result.choices[0].message.content

pipeline("Hello, world!")
# Out: Hello there! How can I assist you today?
```

LangSmith

```
python

@openlit.trace
def generate_one_liner():
    completion = client.chat.completions.create(
        model="gpt-3.5-turbo",
        messages=[
            {
                "role": "system",
                "content": "Return a one liner from any movie for me to guess",
            }
        ],
    )
```

OpenLIT

支持的语言有限、需要修改代码

Outline



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE
SUMMIT

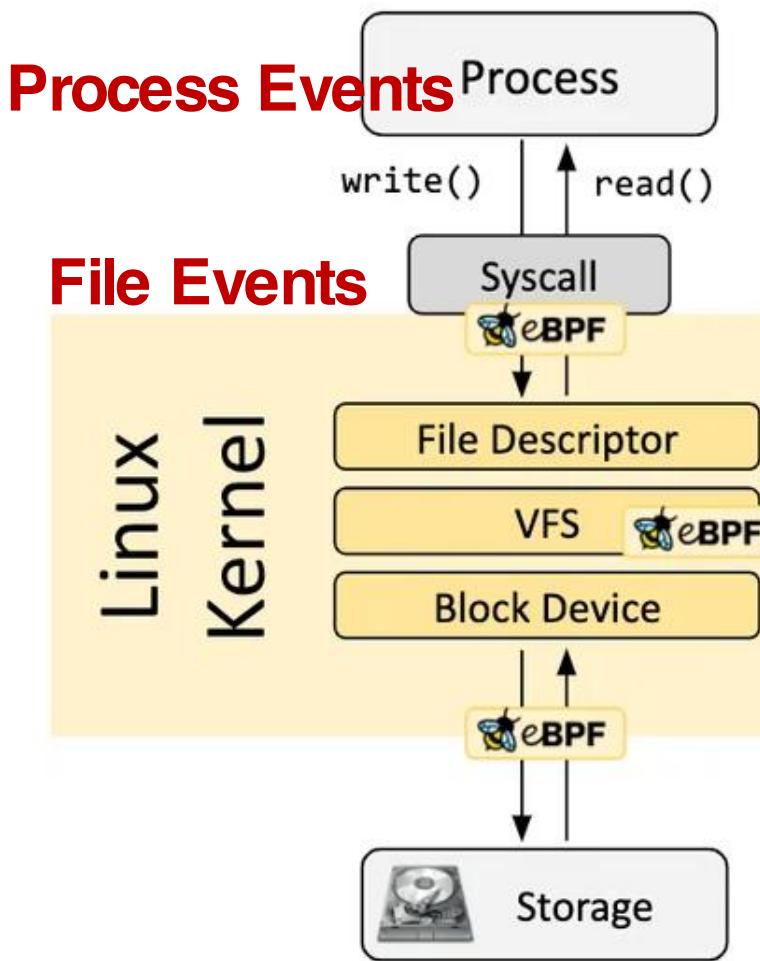


AI_dev
Open Source Dev & ML Summit

China 2024

1. 背景：训练和推理的效率挑战
2. 现状：传统解决方案和工具的问题
3. 方法：eBPF 构建零侵扰可观测性
4. 实践：PyTorch 全栈剖析和追踪

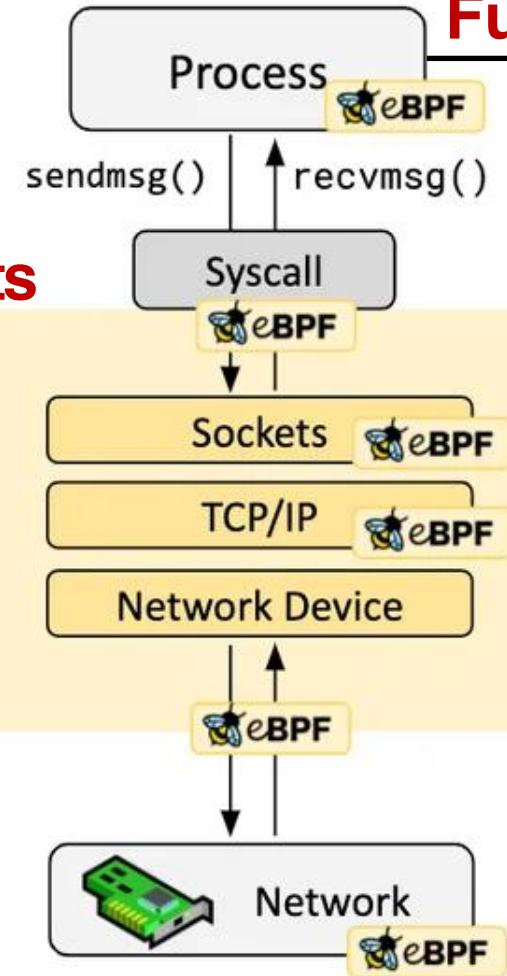
eBPF 的可观测性能



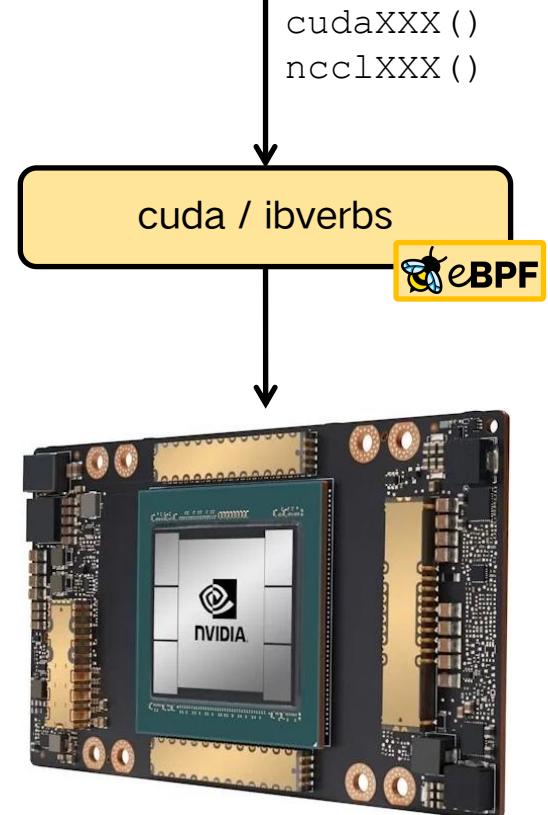
Socket Events

Kernel Events

HW Events



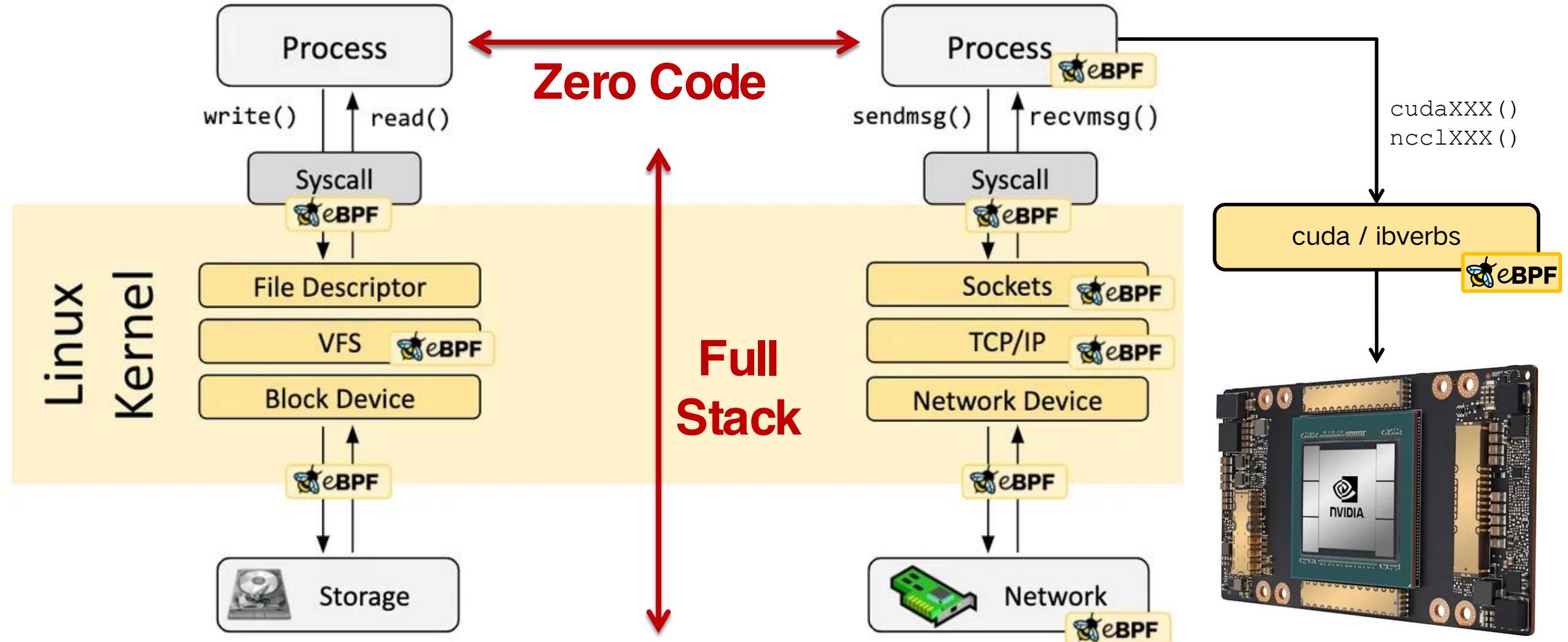
Perf Events Function Call Events



使用 eBPF 实现可观测性的优势



China 2024



他山之石：eBPF Profiling & Tracing



KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMIT

China 2024



GPU Profiling with eBPF at Meta

How Meta is building and using a low overhead GPU observability tool set using eBPF

Riham Selim

@ebpfsummit

Meta: eBPF GPU Profiling

The slide features a background graphic of a hexagonal lattice structure with circuit board elements. In the center, there is a logo for "eBPF Summit" with a bee icon and the text "powered by ISOVALENT". A small portrait of Riham Selim is in the top right corner.

Trace-enabled Timing Model Synthesis for ROS2-based Autonomous Applications

Hazem Abaza^{*†}, Debayan Roy*, Shiqing Fan[‡], Selma Saidi[†] and Antonios Motakis^{*}
[†]Technische Universität Dortmund, ^{*}Huawei Dresden Research Center, [‡]Huawei Munich Research Center
{hazem.abaza, selma.saidi}@tu-dortmund.de, {debayan.roy6, shiqing.fan, antonios.motakis}@huawei.com

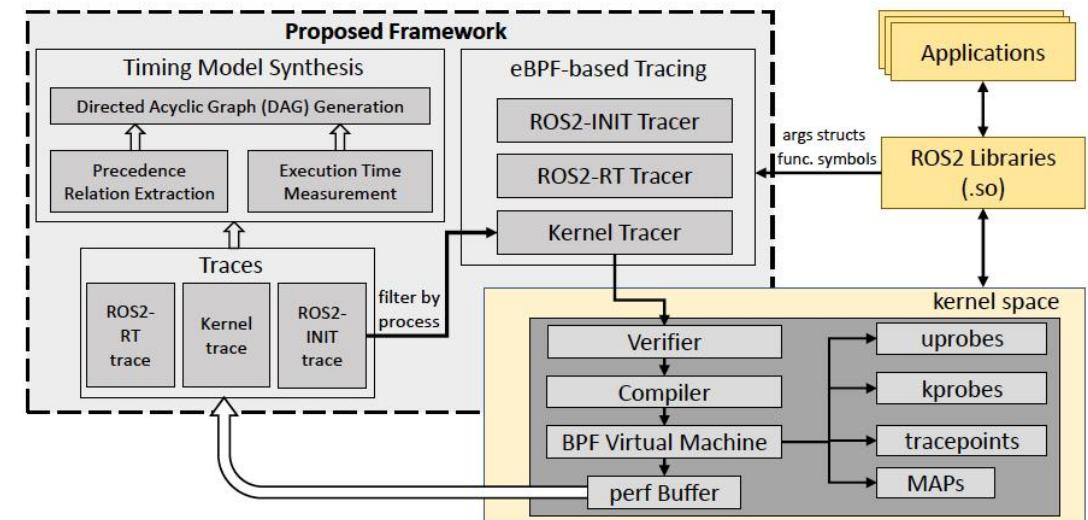


Fig. 1: Proposed trace-enabled timing model synthesis framework.

华为：eBPF Tracing + ROS2

使用 eBPF 实现持续剖析的技术挑战



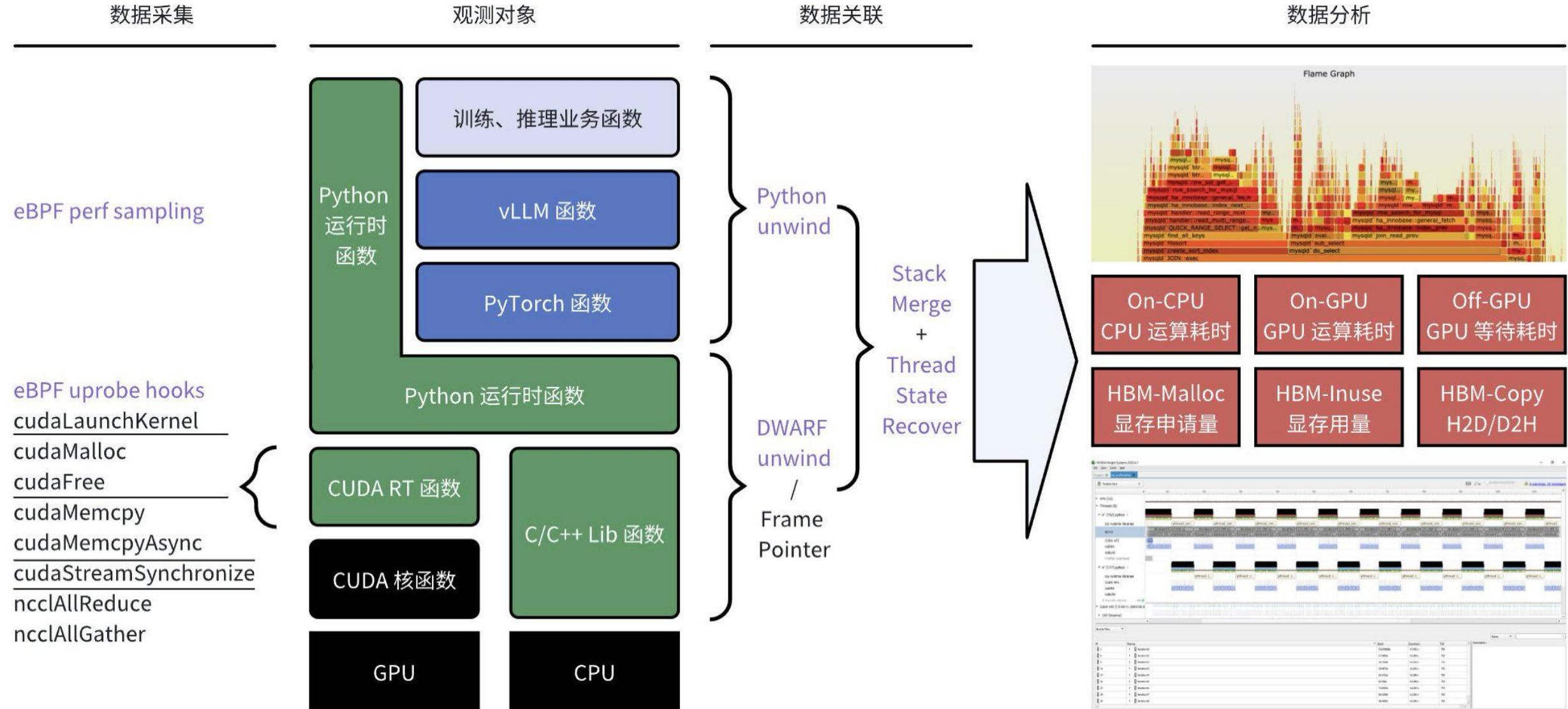
KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI_dev
Open Source GPU & ML Summit

China 2024

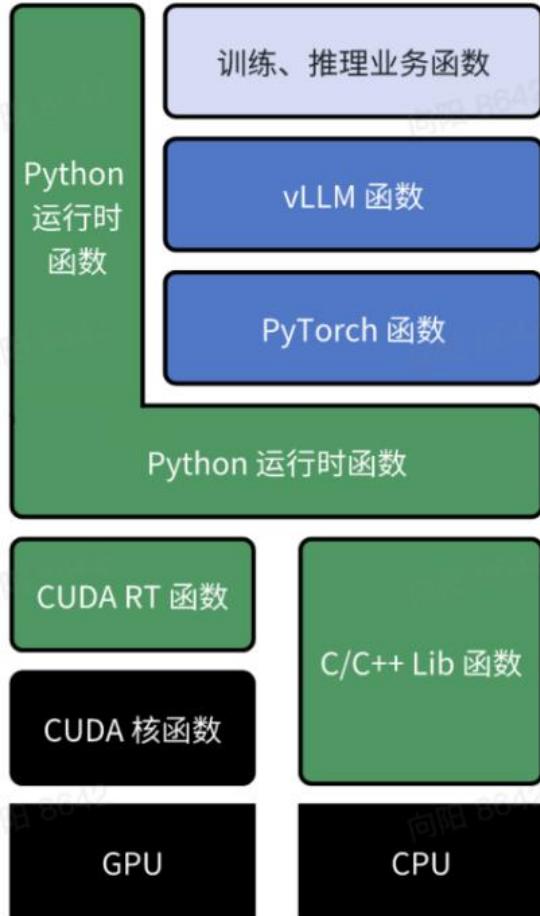


如何合并 Python Stack 和 C/C++ Stack



China 2024

栈合并目的：全栈剖析



```
def baz():
    time.sleep(100)

def bar():
    baz()

def foo():
    bar()

def main():
    foo()

main()
```

```
#0 __select_nocancel ()
#1 pysleep
#2 time_sleep
#3 call_function
#4 PyEval_EvalFrameEx
#5 fast_function
#6 call_function
#7 PyEval_EvalFrameEx
#8 fast_function
#9 call_function
#10 PyEval_EvalFrameEx
#11 fast_function
#12 call_function
#13 PyEval_EvalFrameEx
#14 fast_function
#15 call_function
#16 PyEval_EvalFrameEx
#17 _PyEval_EvalCodeWithName
#18 PyEval_EvalCodeEx
#19 PyEval_EvalCode
```

<https://github.com/deepflwio/deepflow>

例：剖析显存申请和使用量



China 2024

① eBPF uprobe
Hook `cuda_malloc`
获取显存申请调用栈

```
SEC("uprobe/cuda_malloc")
int uprobe_cuda_malloc(struct pt_regs *ctx) {
    __u64 id = bpf_get_current_pid_tgid();
    __u32 tgid = id >> 32;

    void *address = (void *) PT_REGS_PARM1(ctx);
    __u64 size = (__u64) PT_REGS_PARM2(ctx);
    malloc_data_t *data = cuda_malloc_info_lookup(&tgid);
    __u64 call_time = bpf_ktime_get_ns();

    if (data) {
        data->address = address;
        data->size = size;
        data->call_time = call_time;
        data->rip = PT_REGS_IP(ctx);
    } else {
        malloc_data_t newdata = { .address = address, .size = size, .call_time = call_time, .rip = PT_REGS_IP(ctx) };
        cuda_malloc_info_update(&tgid, &newdata);
    }

    return 0;
}
```

③ eBPF uprobe
Hook `cuda_free`
获取释放的显存地址

```
SEC("uprobe/cuda_free")
int uprobe_cuda_free(struct pt_regs *ctx) {
    __u64 id = bpf_get_current_pid_tgid();

    void *addr = (void *) PT_REGS_PARM1(ctx);

    __u32 zero = 0;
    unwind_state_t *state = heap_lookup(&zero);
    if (state == NULL) {
        return 0;
    }
    __builtin_memset(state, 0, sizeof(unwind_state_t));

    struct stack_trace_key_t *key = &state->key;
    key->tgid = id >> 32;
    key->pid = (__u32) id;

    /*
     * CPU idle stacks will not be collected.
     */
    if (key->tgid == key->pid && key->pid == 0) {
        return 0;
    }

    key->cpu = bpf_get_smp_processor_id();
    bpf_get_current_comm(&key->comm, sizeof(key->comm));
    key->timestamp = bpf_ktime_get_ns();

    key->mem_addr = (__u64) addr;
    bpf_perf_event_output(ctx, &NAME(cuda_memory_output), BPF_F_CURRENT_CPU, &state->key, sizeof(state->key));

    return 0;
}
```

④ 计算当前
显存消耗

```
SEC("uretprobe/cuda_malloc")
int uretprobe_cuda_malloc(struct pt_regs *ctx)
{
    __u64 id = bpf_get_current_pid_tgid();
    __u32 tgid = id >> 32;

    long ret = PT_REGS_RC(ctx);
    if (ret != 0) {
        return 0;
    }

    malloc_data_t *data = cuda_malloc_info_lookup(&tgid);
    if (data == NULL) {
        return 0;
    }

    cuda_malloc_info_delete(&tgid);

    __u32 zero = 0;
    unwind_state_t *state = heap_lookup(&zero);
    if (state == NULL) {
        return 0;
    }
    __builtin_memset(state, 0, sizeof(unwind_state_t));

    struct stack_trace_key_t *key = &state->key;
    key->tgid = id >> 32;
    key->pid = (__u32) id;

    /*
     * CPU idle stacks will not be collected.
     */
    if (key->tgid == key->pid && key->pid == 0) {
        return 0;
    }

    key->cpu = bpf_get_smp_processor_id();
    bpf_get_current_comm(&key->comm, sizeof(key->comm));
    key->timestamp = bpf_ktime_get_ns();

    bpf_probe_read_user(&key->mem_addr, sizeof(__u64), (void *)data->address);
    key->mem_size = (__u64) data->size;

    // add one frame for cudaMalloc
    // native unwinding start from uretprobe, which has the reg state after cudaMalloc return
    add_frame(&state->stack, data->rip);

    state->regs.ip = PT_REGS_IP(ctx);
    state->regs.sp = PT_REGS_SP(ctx);
    state->regs.bp = PT_REGS_FP(ctx);
    add_frame(&state->stack, state->regs.ip);
    bpf_tail_call(ctx, &NAME(progs_jmp_uprobe_map), PROG_DWARF_UNWIND_IDX);

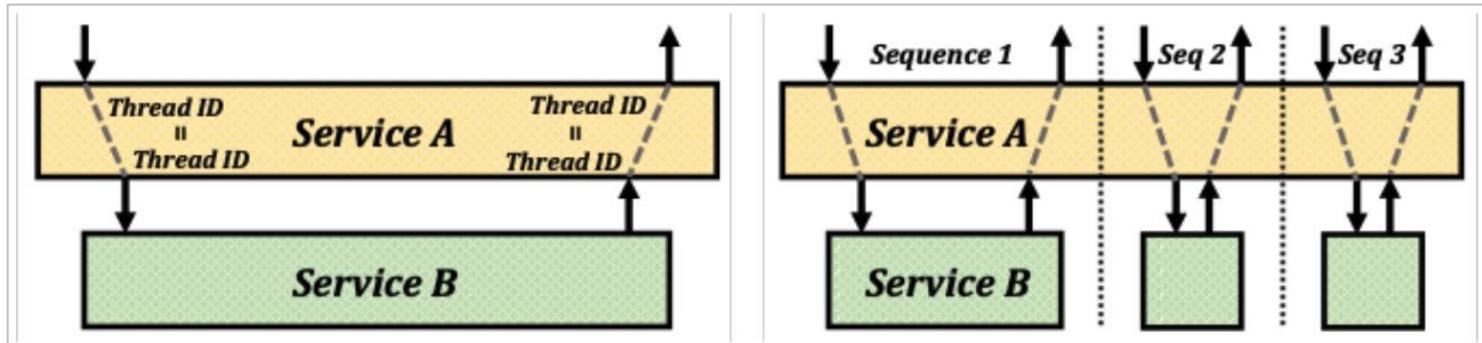
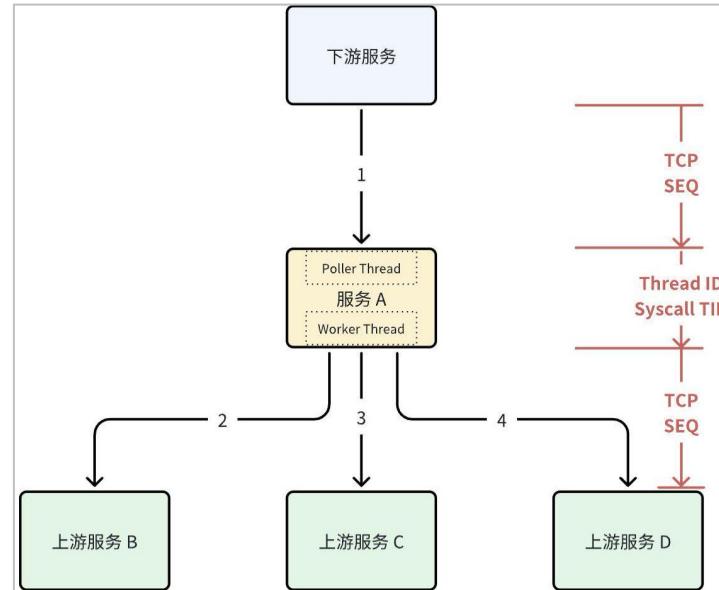
    return 0;
}
```

② eBPF uretprobe
Hook `cuda_malloc`
获取申请的显存地址

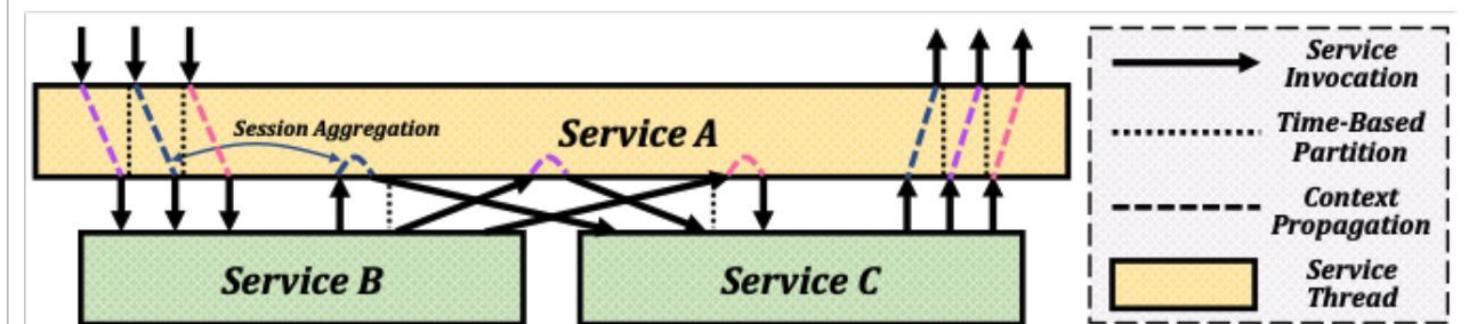
使用 eBPF 实现分布式追踪的技术挑战



China 2024



(a) Implicit context propagation via thread IDs.
(b) Using time sequences to partition thread-reusing spans.



(c) Association for multiple requests or responses.

Figure 7: Intra-component causal association.

Outline



KubeCon



CloudNativeCon



THE LINUX FOUNDATION
OPEN SOURCE
SUMMIT



China 2024

1. 背景：训练和推理的效率挑战
2. 现状：传统解决方案和工具的问题
3. 方法：eBPF 构建零侵扰可观测性
4. 实践：PyTorch 全栈剖析和追踪

DeepFlow 中的 eBPF AutoProfiling



KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI_dev
Open Source Dev & ML Summit

China 2024



1. GPU Profiling



KubeCon



CloudNativeCon

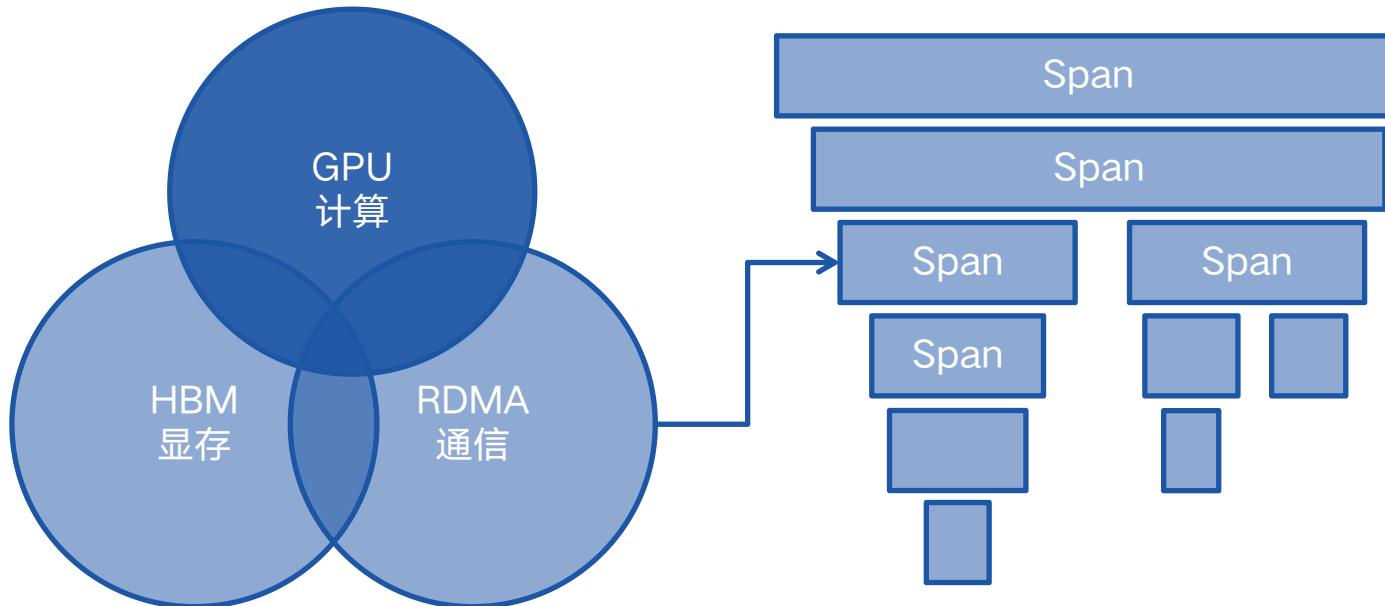


THE LINUX FOUNDATION
OPEN SOURCE
SUMMIT



AI_dev
Open Source Dev & ML Summit

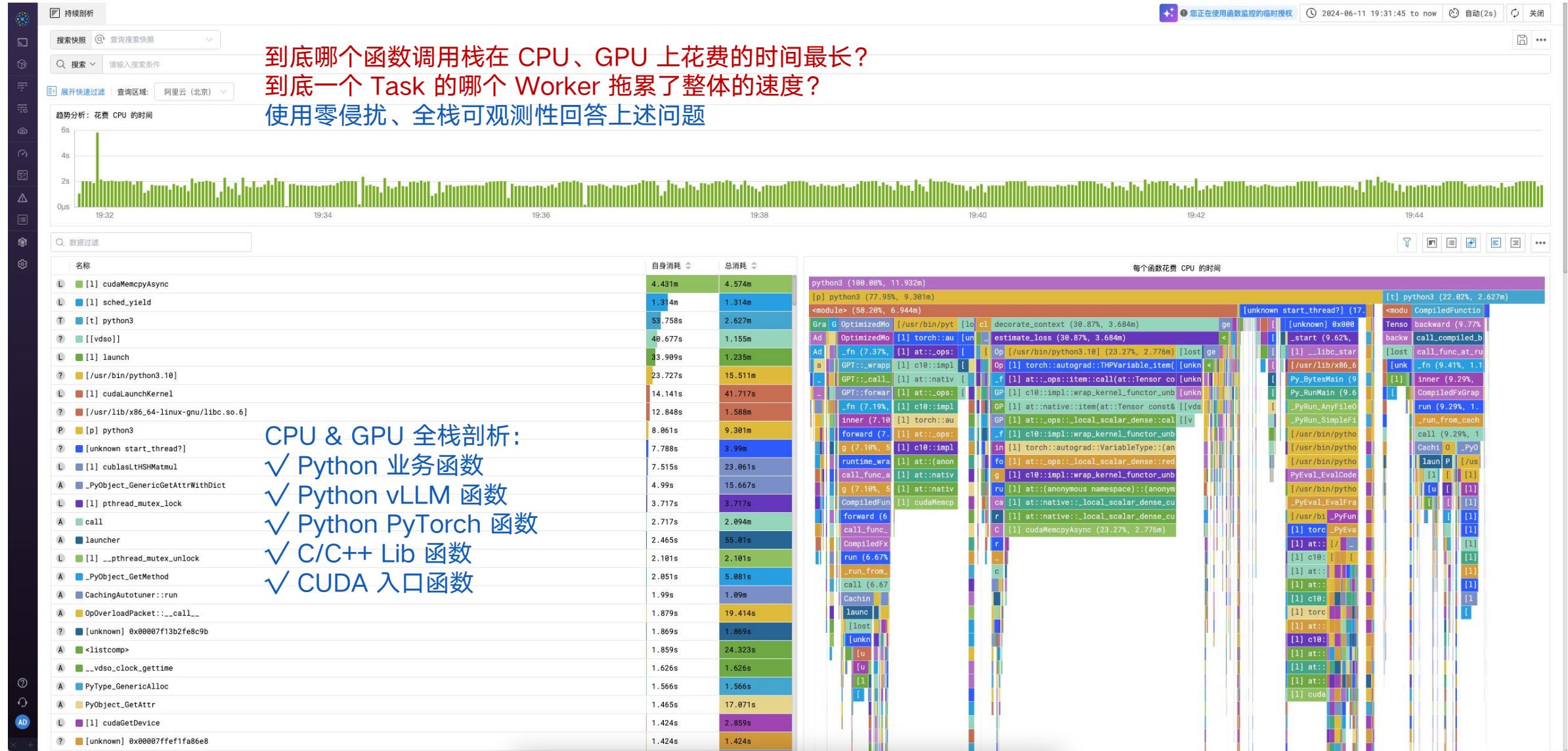
China 2024



PyTorch + nanoGPT CPU & GPU 火焰图



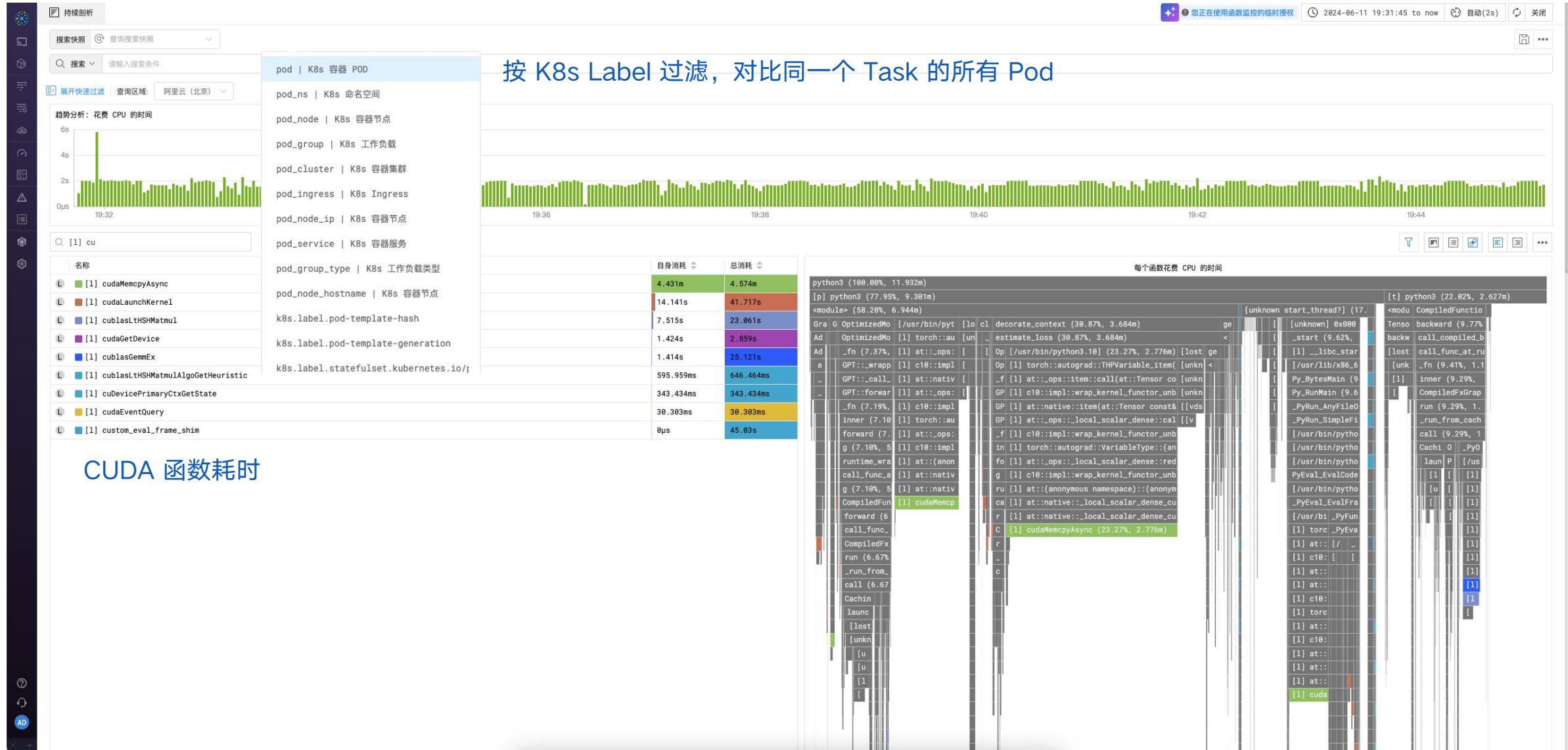
China 2024



PyTorch + nanoGPT CPU & GPU 火焰图



China 2024



2. HBM Profiling



KubeCon

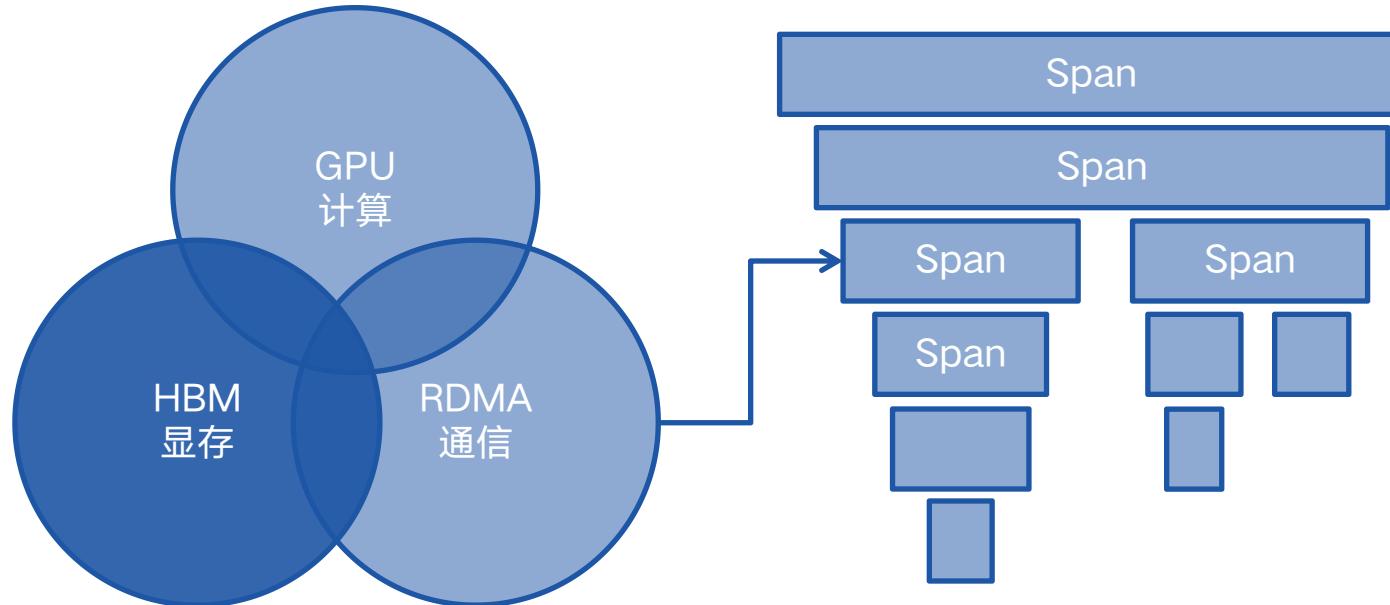


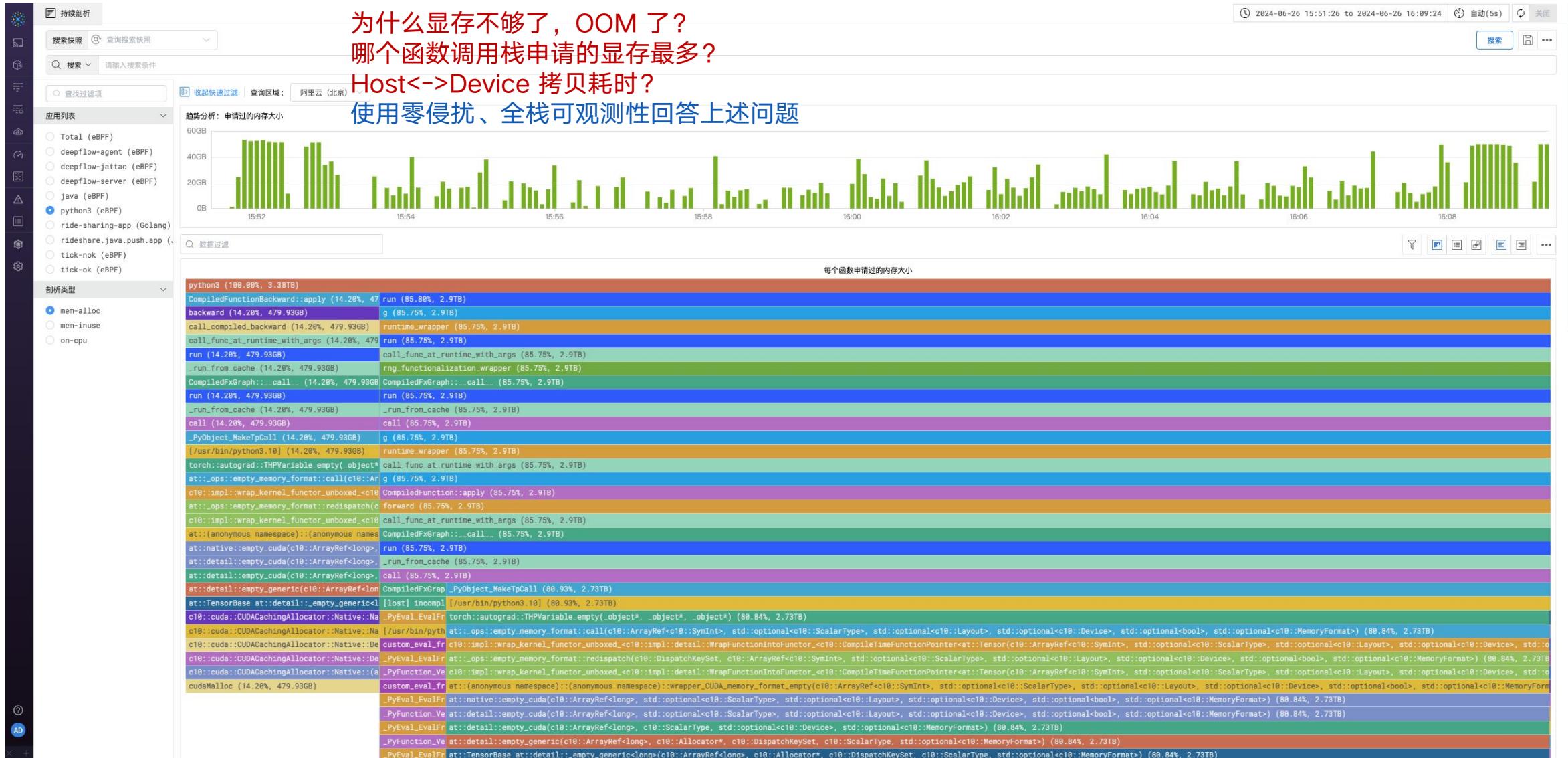
| CloudNativeCo

THE LINUX FOUNDATION
 OPEN SOURCE SUMMIT

 AI_dev
Open Source GenAI & ML Summit

— China 2024





CUDA mem-inuse 显存实时用量火焰图



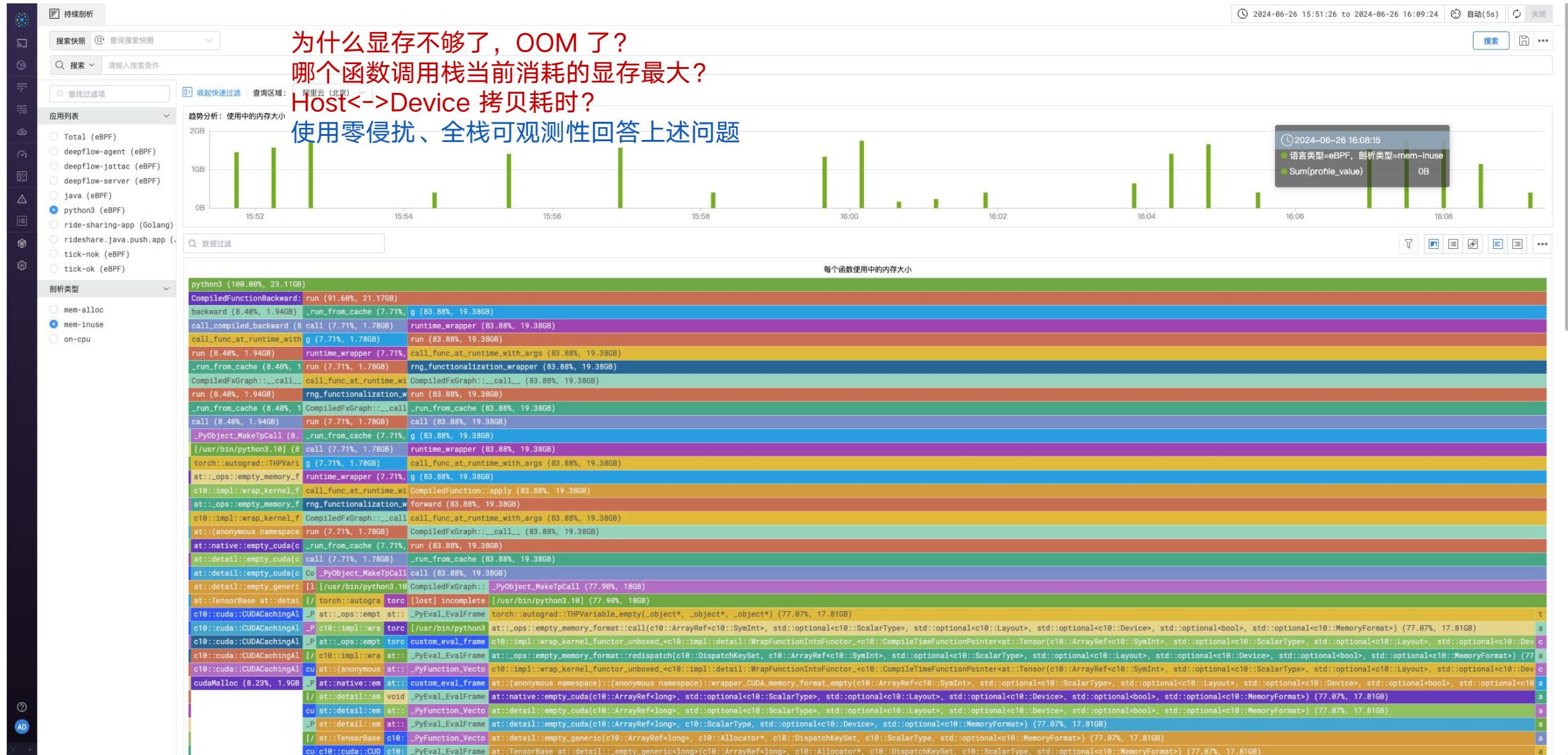
KubeCon



CloudNativeCon

THE LINUX FOUNDATION
OPEN SOURCE SUMMITAI_dev
Open Source Dev & ML Summit

China 2024



3. COMM. Profiling



KubeCon



CloudNativeCon

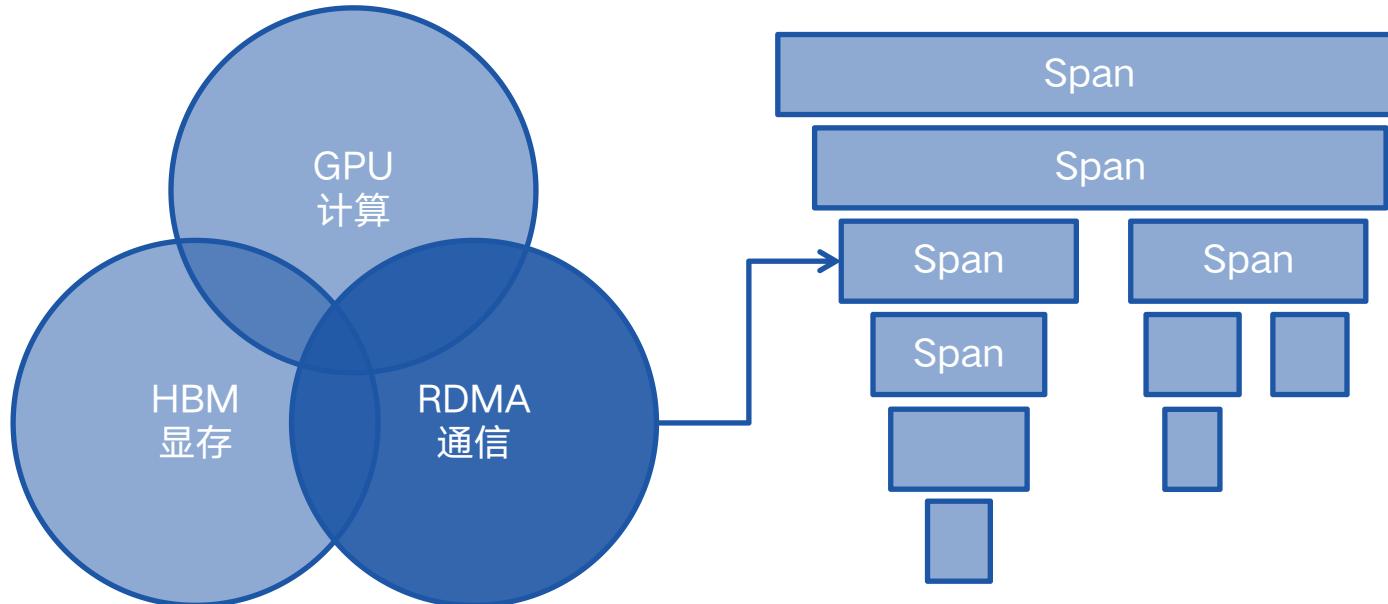


THE LINUX FOUNDATION
OPEN SOURCE SUMMIT



AI_dev
Open Source Dev & ML Summit

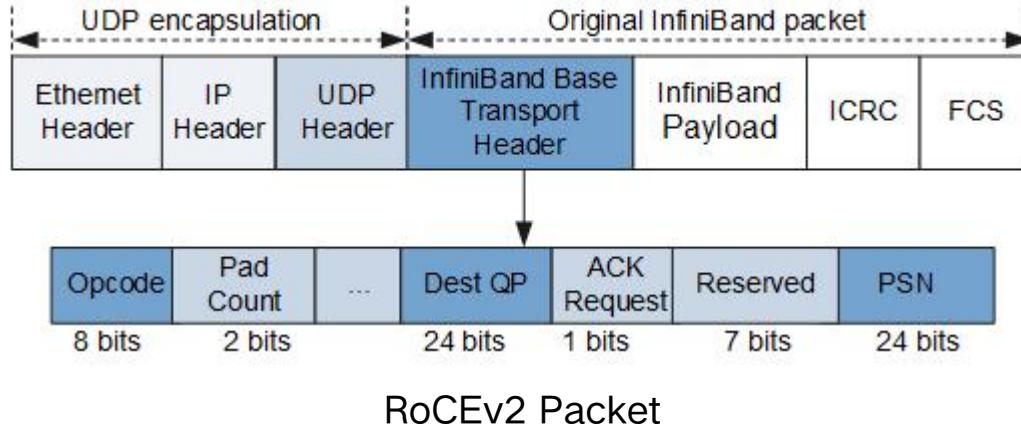
China 2024



RDMA 网络性能剖析



China 2024



eBPF Hooks

Function	作用
ibv_modify_qp	Queue Pair 状态修改
ibv_reg_mr	Memory Region 注册
ibv_dereg_mr	Memory Region 注销
ibv_post_send	buffer 发送或 RDMA 读写
ibv_poll_cq	等待 wr 完成

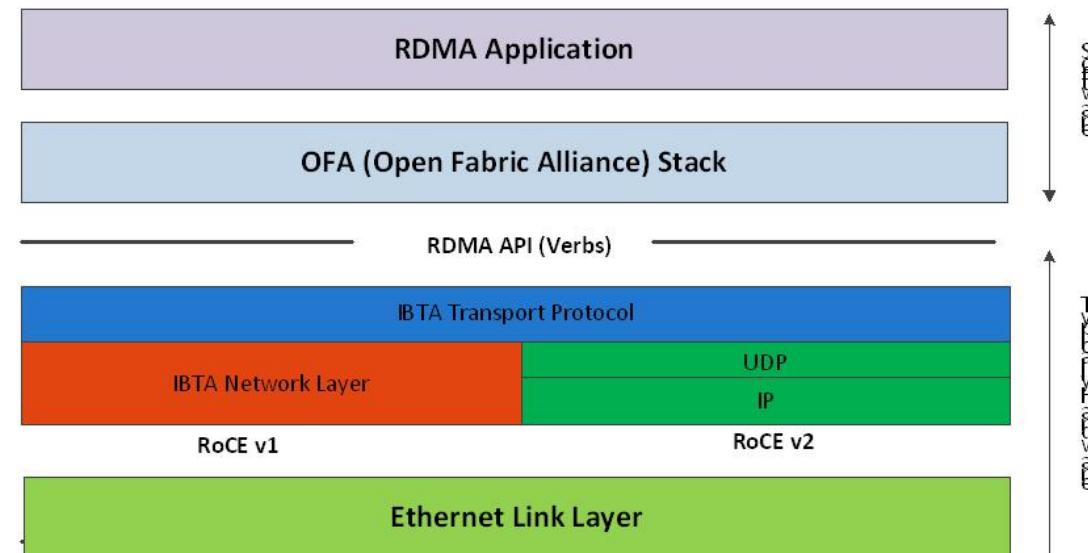
libverbs Hooks, NCCL Hooks, PyTorch Hooks, ...

关键标签:

- Client 及其关联的 K8s Pod、标签
- Server 及其关联的 K8s Pod、标签
- Client Queue Pair
- Server Queue Pair

关键指标:

- 丢包率: 即 NACK 的比例
- 时延: ACK 的时延
- 吞吐: 通信对的 bps、pps 等



4. Distributed Tracing



KubeCon



CloudNativeCon

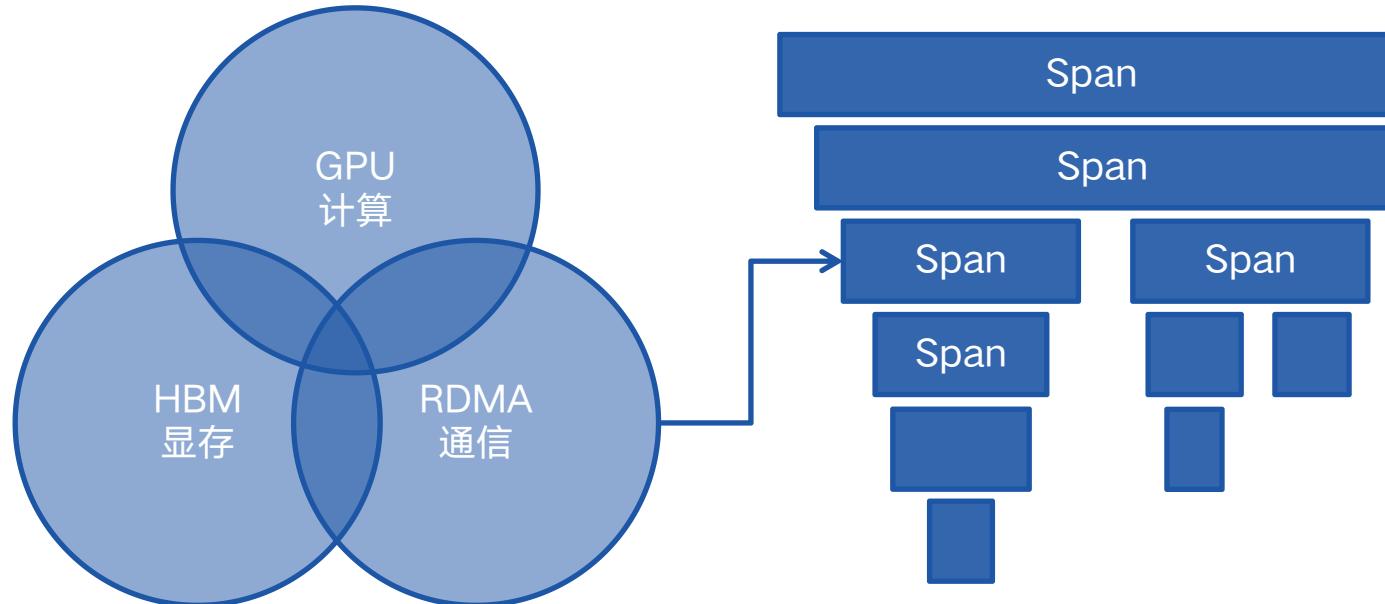


THE LINUX FOUNDATION
OPEN SOURCE SUMMIT

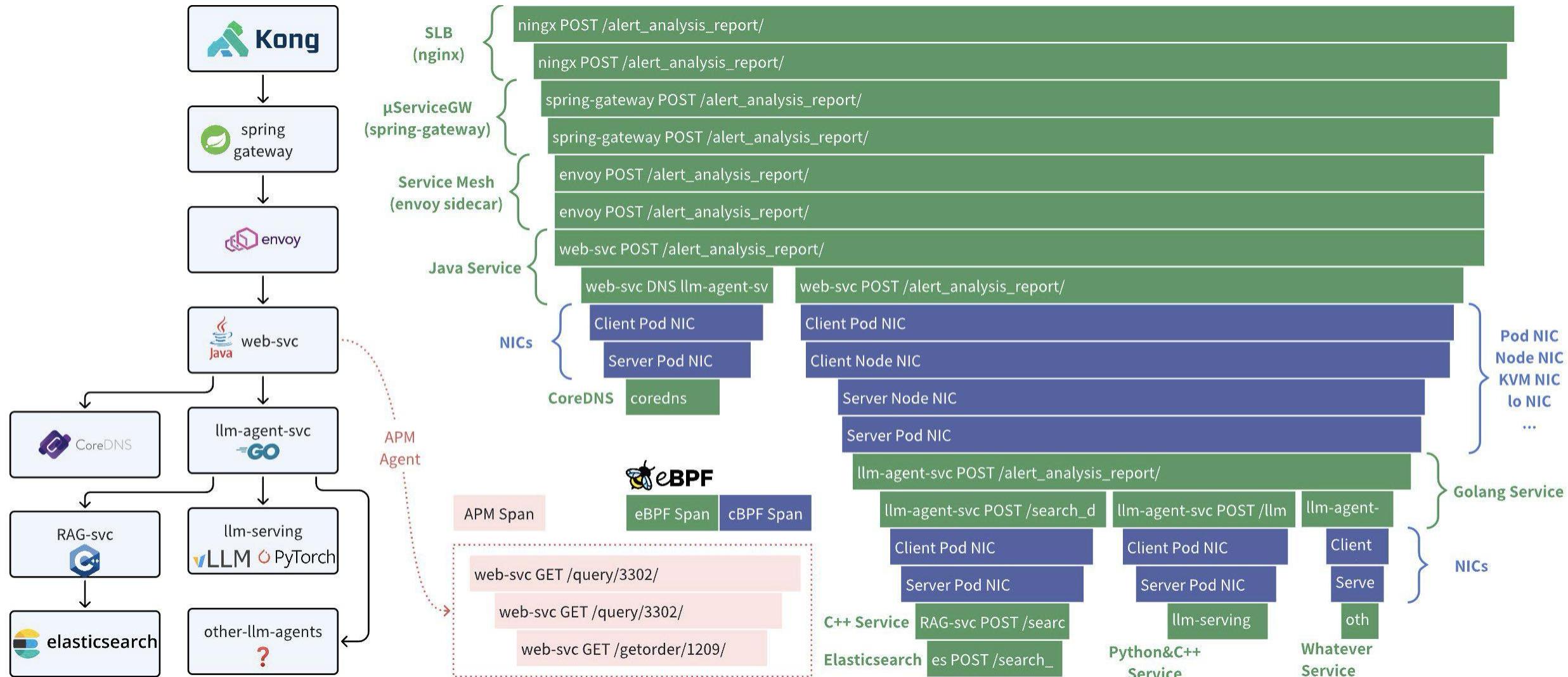


AI_dev
Open Source Dev & ML Summit

China 2024



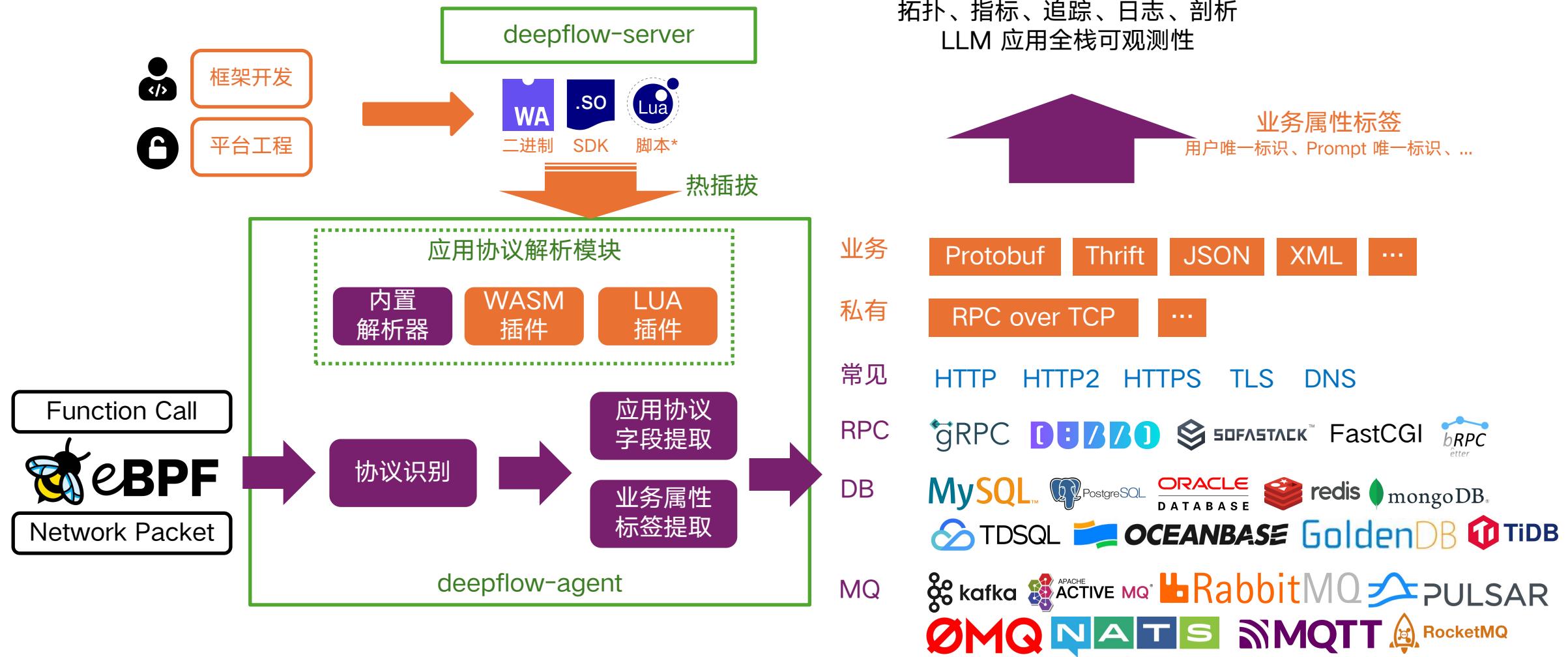
云侧、端侧的在线推理服务



内置的及可编程的协议识别能力



China 2024



DeepFlow: 零侵扰实现 AI Infra 和应用的全栈可观测性



无需修改代码，无需重启进程