

GPU Resource Pooling and the benefit of deploying CUPTI

Alibaba Group

Lingling Jin, Lingjie Xu

Alibaba is an AI company







- E-commerce
 - Search
 - Advertisement
- Financing & payment
- Video
- Logistics

All these services are backed up by Alibaba cloud



GPU pool



• GPU accelerator powered data center

IDC data center

GPU servers are managed elastically



Kubernetes

However GPU utilization is still low



KVM





Agenda

- GPU accelerated inference service pool: multi-tenant
- GPU resource pool over ethernet

GPU inference service pool





- GPUs are allocated in single card granularity
- Capacity is designed for worst case traffic
- Mixing different inference possible?



Experiment setup

- Single GPU V100 system
- Model: resnet50
- Start inference server, capture NVML GPU utilization at backend
- Client performance test
 - multiple concurrent request,
 - max allowed latency
 - desired batchsize
 - Calculate server throughput



NVML GPU utilization





GPU_utilization: Percent of time over the past sample period during which one or more kernels was executing on the GPU

CUPTI



- The CUDA Profiling Tools Interface (CUPTI) enables the creation of profiling and tracing tools that target CUDA applications.
- all event/metrics available in NVPROF
 - sm_efficiency
 - achieved_occupancy
 - dram_read/write
 - Single/double/half_precision_fu_utilization
- Collect GPU metrics during runtime
- 1%-10% overhead

CUPTI SM Efficiency





SM efficiency: The percentage of time at least one warp is active on a multiprocessor averaged over all multiprocessors on the GPU. CUPTI is what we want.

Explore two processes GPU sharing





C-) Alibaba Cloud

Worldwide Cloud Services Partner

Explore multi-process GPU sharing

2 Concurrent Tensorflow Alexnet Inference application share GPU

C-) Alibaba Cloud

Worldwide Cloud Services Partner



Explore multi-process GPU sharing



2 Concurrent Tensorflow Alexnet Inference application share GPU



Running two inference servers



- Single GPU V100 system
- Model
 - resnet50
 - Inception
- Batchsize 4

Performance with two inference



servers







© 2019 Alibaba Group

Perf overhead

Nvidia inference server





- Each model is implemented as GPU stream, can be run concurrently on GPU
- better control with scheduler

Nvidia inference server



GPU Activity Over Time



1 server vs 2 server performance comparison

C-) Alibaba Cloud Services Partner

Inception + resnet50, batchsize8, one instance



■ Inception ■ resnet50



Agenda

- GPU accelerated inference service pool: multi-tenant
- GPU resource pool over ethernet





- Each AI application has different CPU/GPU ratio
- many types GPU servers



 Single GPU card and GPU server is becoming so powerful, cannot be fully utilized

Distributed storage system-big success





- seperated CPU(client) and storage(chunk server)
- Stable
- High performance
- Easy to maintain
- Low cost



distributed acceleration





GPU pool over ethernet



Components can be generated automatically with minimum manual work

Components can be shared by GPU/FPGA/etc.

Challenge 1-A lot of CUDA API



- A intermediate and auto-generated YAML file to define the APIs that need to be forwarded to remote server.
- Need to add some information on pointer arguments manually



Challenge 2-Host memory



- With remote GPU, we need to detect such memories in CUDA kernel parameters to make sure the function correctness
 - cudaMallocHost

Challenge 3 - Hide API forwarding Latency

- Changing blocking API calls to asynchronous calls
- Benefit: 2x speedup



Before









- website: https://aimatrix.ai
- Code repo: https://github.com/alibaba/ai-matrix

CNN Training Performance



Tensorflow CNN performance comparison

Baseline(local run) Remote GPU





Summary

- GPU accelerated inference service pool: multi-tenant
- GPU resource pool over ethernet
- All results are preliminary and we are actively working on it
- Discussions and collaborations are welcome
- Email: l.jin@Alibaba-inc.com



Worldwide Cloud Services Partner

Thanks! Q&A