

OLCF Training AI on Frontier

Junqi Yin

Analytics and AI Methods at Scale

Feb 16, 2023

ORNL is managed by UT-Battelle, LLC for the US Department of Energy

Outline

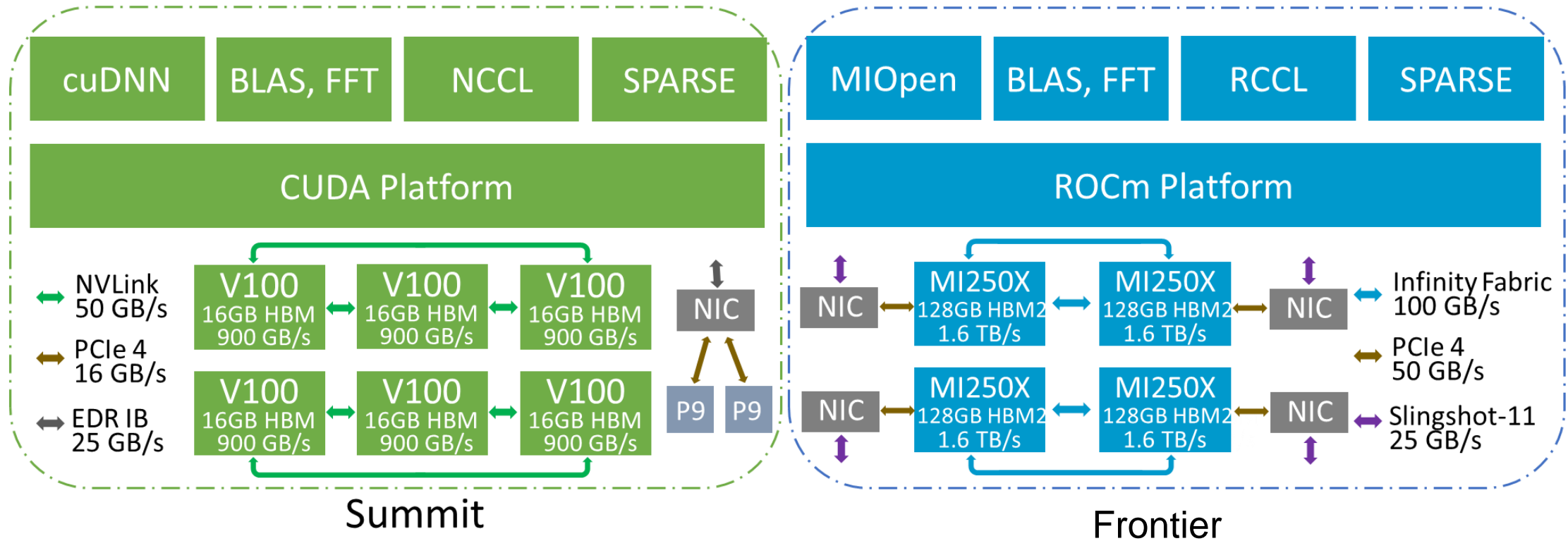
- Frontier DL Environment
- Preliminary performance numbers
 - Kernels: GEMM/CONV/LSTM
 - Models: CNN/RNN
 - Applications: ResNet50, STEMDL
- Simulation-ML integration

Deep Learning Stacks

Framework

TensorFlow, PyTorch, ...

DL Stack



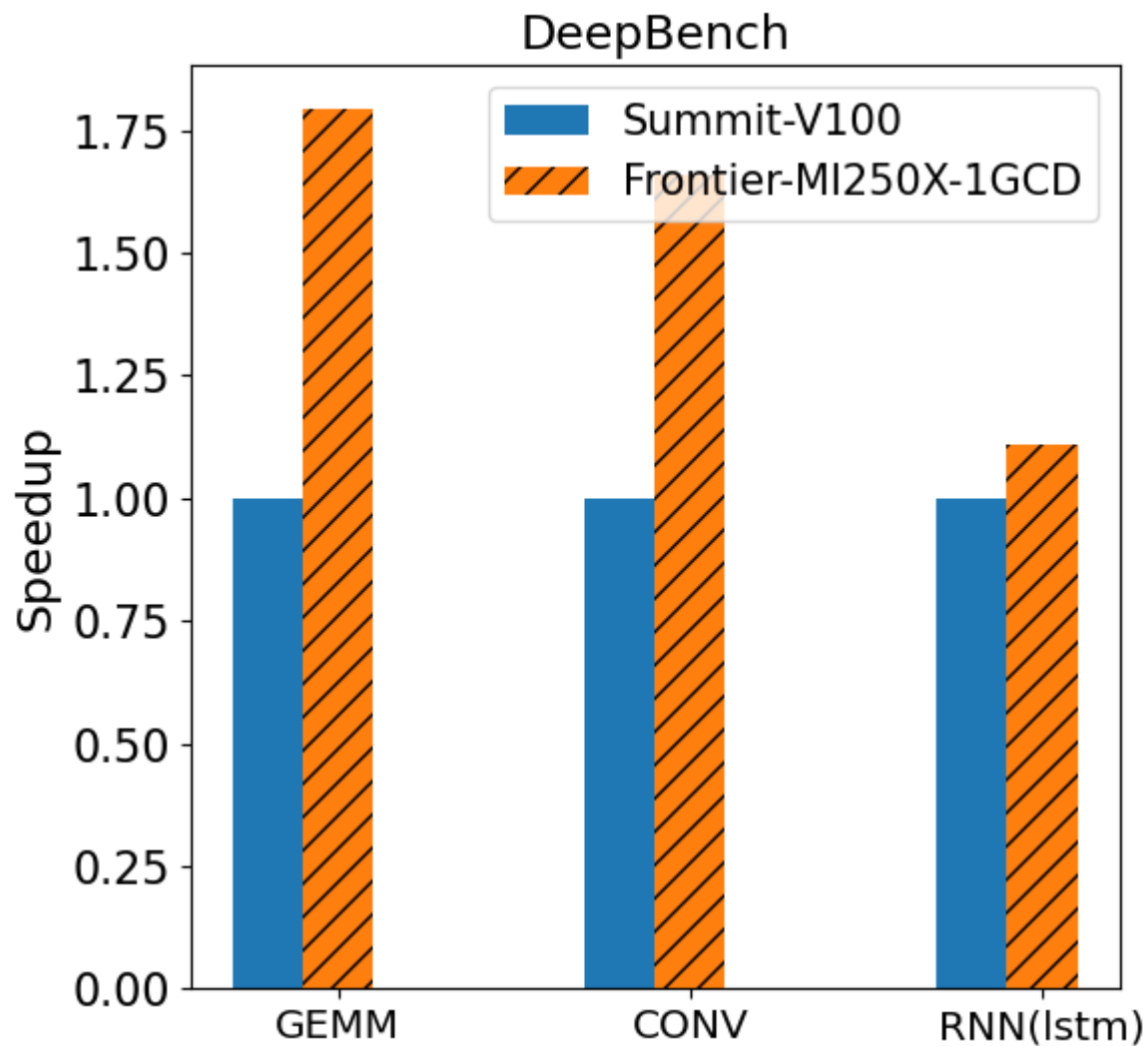
- Most DL codes work on Frontier without changes

Frontier DL Environment

- What's working?
 - TensorFlow and PyTorch
 - Third-party libraries
 - Cray DL-plugin
 - Horovod
 - DeepSpeed
 - PyG ...
- Peculiarities
 - rocm-smi
 - MIOpen cache
 - RCCL + libfabric
- Resources: <https://github.com/ROCmSoftwarePlatform>

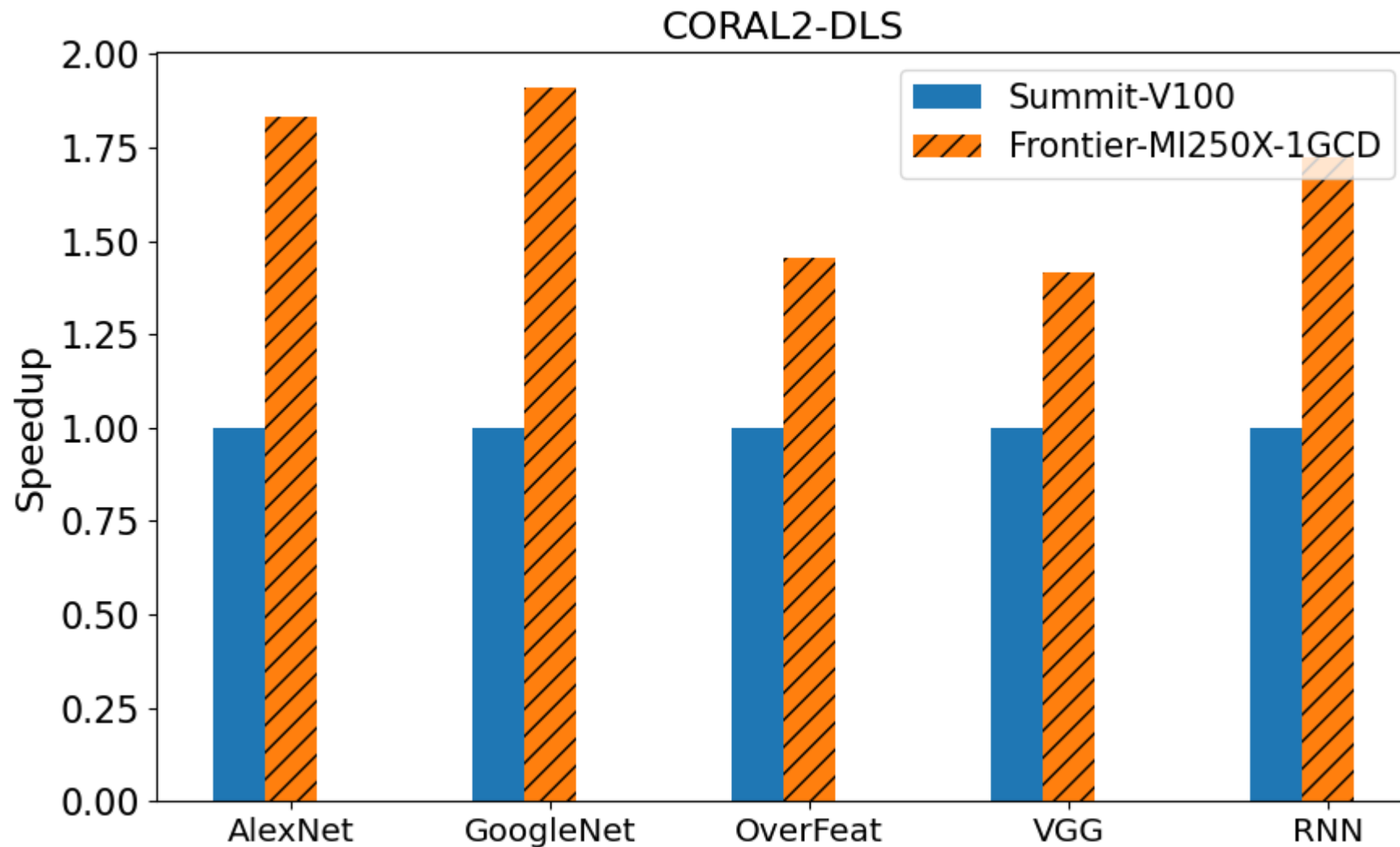
Performance baselines: Kernels

- Kernel Ops (fp32)
 - GEMM ~ 1.7x
 - CONV ~1.6x
 - LSTM ~1.1x



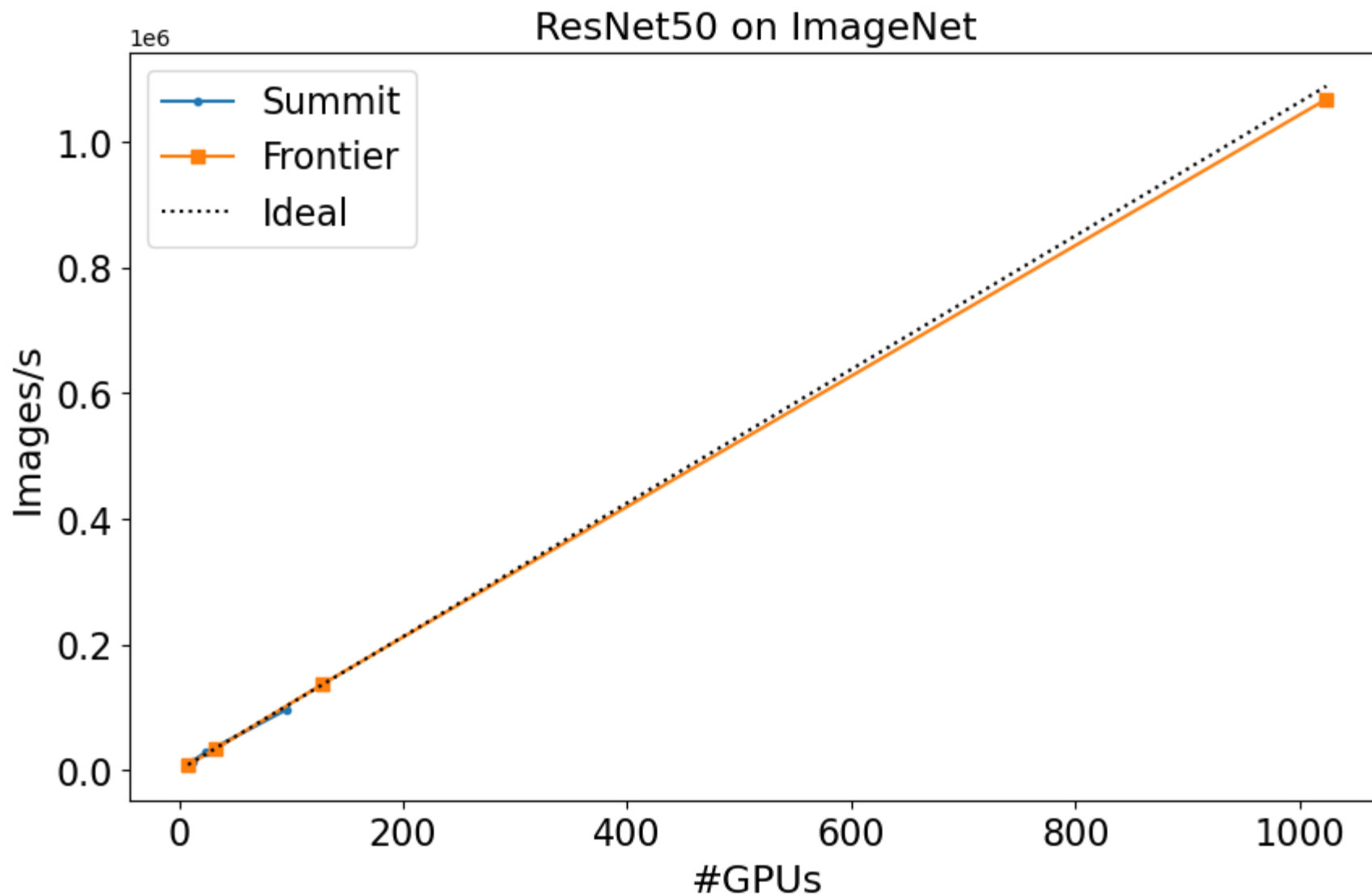
Performance baselines: Models

- CNN (fp32)
 - AlexNet
 - GoogleNet
 - OverFeat
 - VGG
- RNN



Performance baselines: Apps

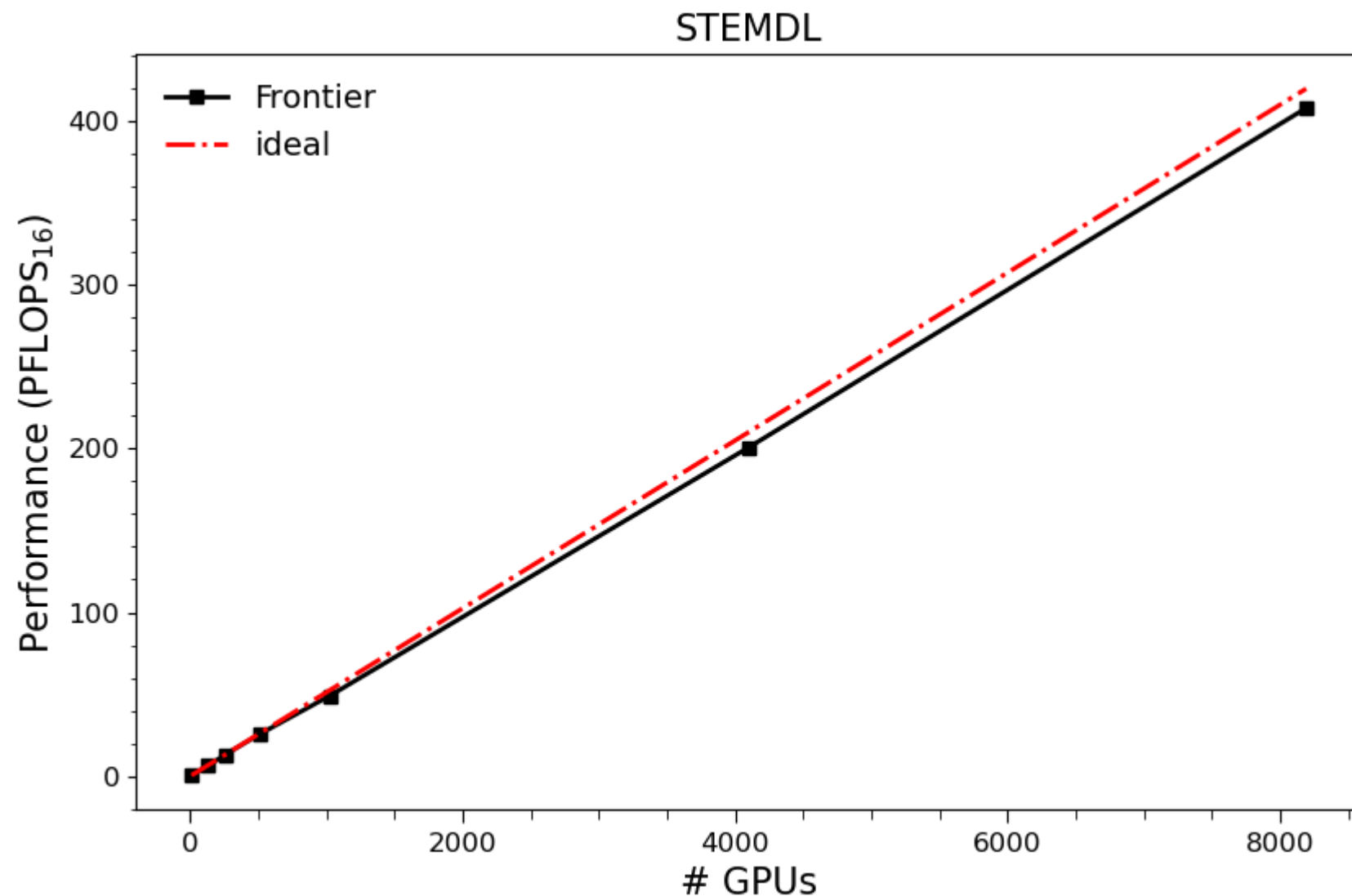
- ResNet50
 - Mixed
 - ~1.0x per GCD
 - 98% at 1024



Performance baselines: Apps

- STEMDL
 - Tiramisu network
 - 220M parameters
 - 97% at 8192 GPUs

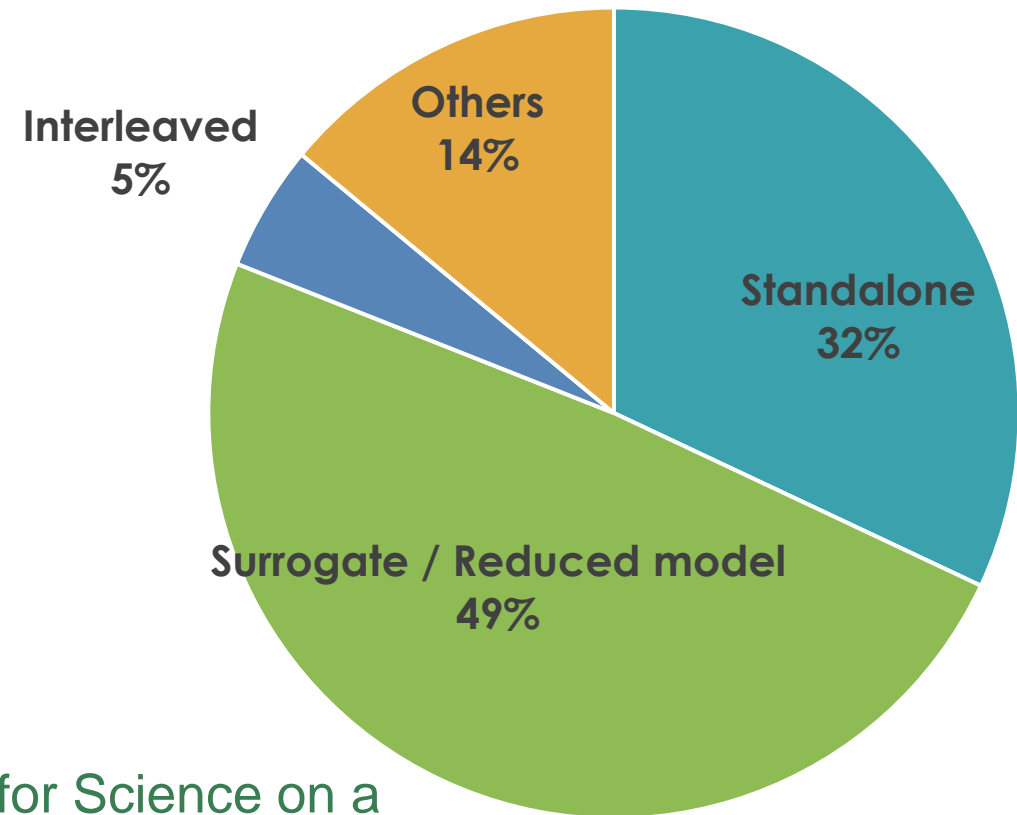
[Accelerating Collective Communication in Data Parallel Training across Deep Learning Frameworks.](#)
[USENIX NSDI'22, 2022](#)



Simulation-ML Integration

- Best of both worlds: FP64 simulation + FP16 modeling
- Common use cases
 - Surrogate modeling
 - Reduced model
 - Interleaved

OLCF (2019 - 2021)
662 projects (147 INCITE, 72 ALCC, 62 ECP, DD)

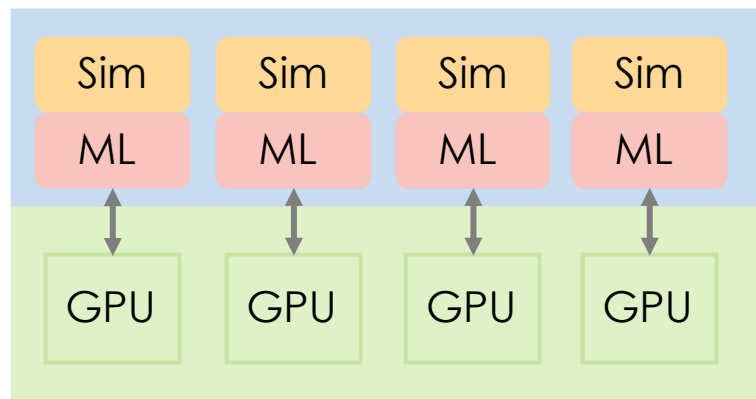


[Learning to Scale the Summit: AI for Science on a Leadership Supercomputer, IPDPSW 2022](#)

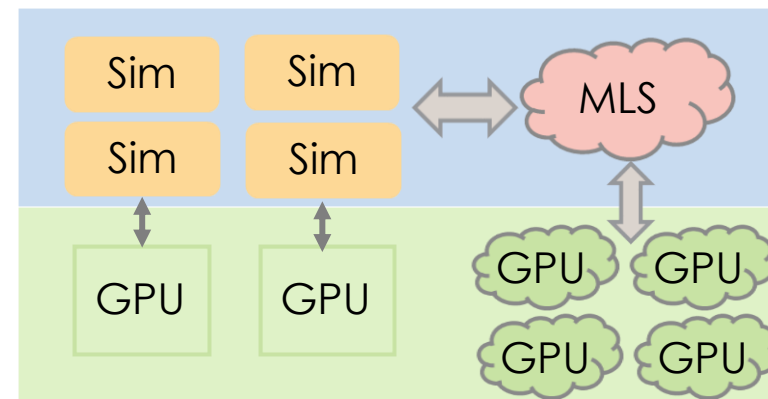
Simulation-ML Integration

- Tightly coupled
 - Single executable
 - one-to-one
- Loosely coupled
 - Different machines
 - many-to-many
- Semi-tightly
 - Separate executables
 - 1-to-1, many-to-many

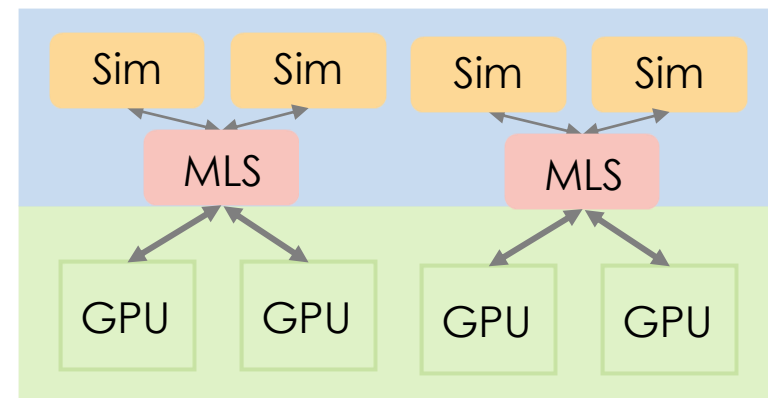
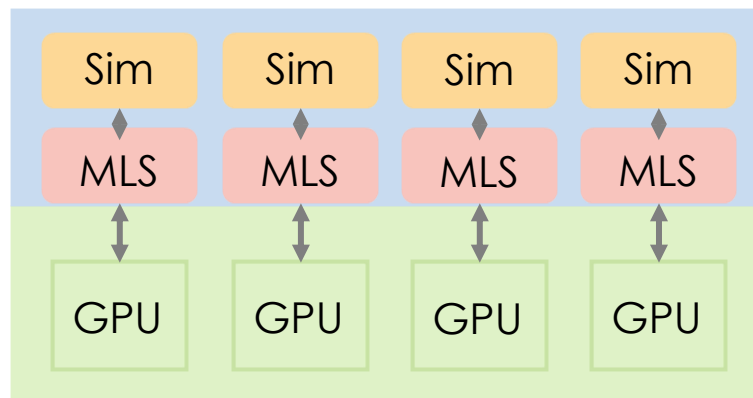
Tightly coupled



Loosely coupled



Semi tightly coupled



Simulation-ML Integration

Method	Pro	Con
Framework C++ API (TensorFlow/PyTorch C++)	<ul style="list-style-type: none">• Portable• Better latency• Easy to deploy	<ul style="list-style-type: none">• Not flexible
Framework Server (TensorFlow Serving/TorchServe)	<ul style="list-style-type: none">• Flexible• Better throughput• Options to deploy	<ul style="list-style-type: none">• High maintenance
Third-party API (SmartRedis/RedisAI)	<ul style="list-style-type: none">• Easy integration• More functionality• Model support	<ul style="list-style-type: none">• Portability

TensorFlow C++

- Assume model in TF *SavedModel* format
- Link with *libtensorflow_cc.so*
- Support *half*, *uint8*, ...

Implementation 1 Use TensorFlow C++ API

```
1: function LOADMODEL(string ModelPath)
2:   tensorflow::SessionOptions SessOpt
3:   tensorflow::RunOptions RunOpt
4:   tensorflow::SavedModelBundle Model           ▷ Get model path, and
                                                setup session and run op-
                                                tions
5:   tensorflow::LoadSavedModel(SessOpt, RunOpt,
6:     ModelPath, {"serve"}, &Model)           ▷ load pre-trained model
7:   return Model
8: end function

9: function PREDICT(tensorflow::SavedModelBundle Model, tensorflow::Tensor In-
10:  puts)
11:   tensorflow::Tensor Inputs                 ▷ input tensor
12:   vector<tensorflow::Tensor> Outputs
13:   string InputNode
14:   string OutputNode                         ▷ input and output nodes
15:   Model.GetSession().Run(data, {OutputNode}, {}, &Outputs)
16:   return Outputs
17: end function

18: function ENERGY
19:   for  $i \leftarrow 0, N_{elements}$  do
20:      $Model_i \leftarrow LOADMODEL(Path_i)$ 
21:   end for
22:   for  $i \leftarrow 0, N_{elements}$  do           ▷ Add energy by element species
23:     tensorflow::TensorShape DataShape
24:     tensorflow::Tensor Inputs(tensorflow::DT_UINT8, DataShape)
25:     for  $j \leftarrow 0, N_{neighbors}$  do
26:        $Inputs \leftarrow Atom[i][j]$ 
27:     end for
28:     Outputs = PREDICT( $Model_i$ , Inputs)
29:      $Energy \leftarrow Outputs$ 
30:   end for
31:   return Energy
32: end function
```

TensorFlow Serving

- Launch server:
`tensorflow_model_server --port --model_config_file`
- Support grpc & http

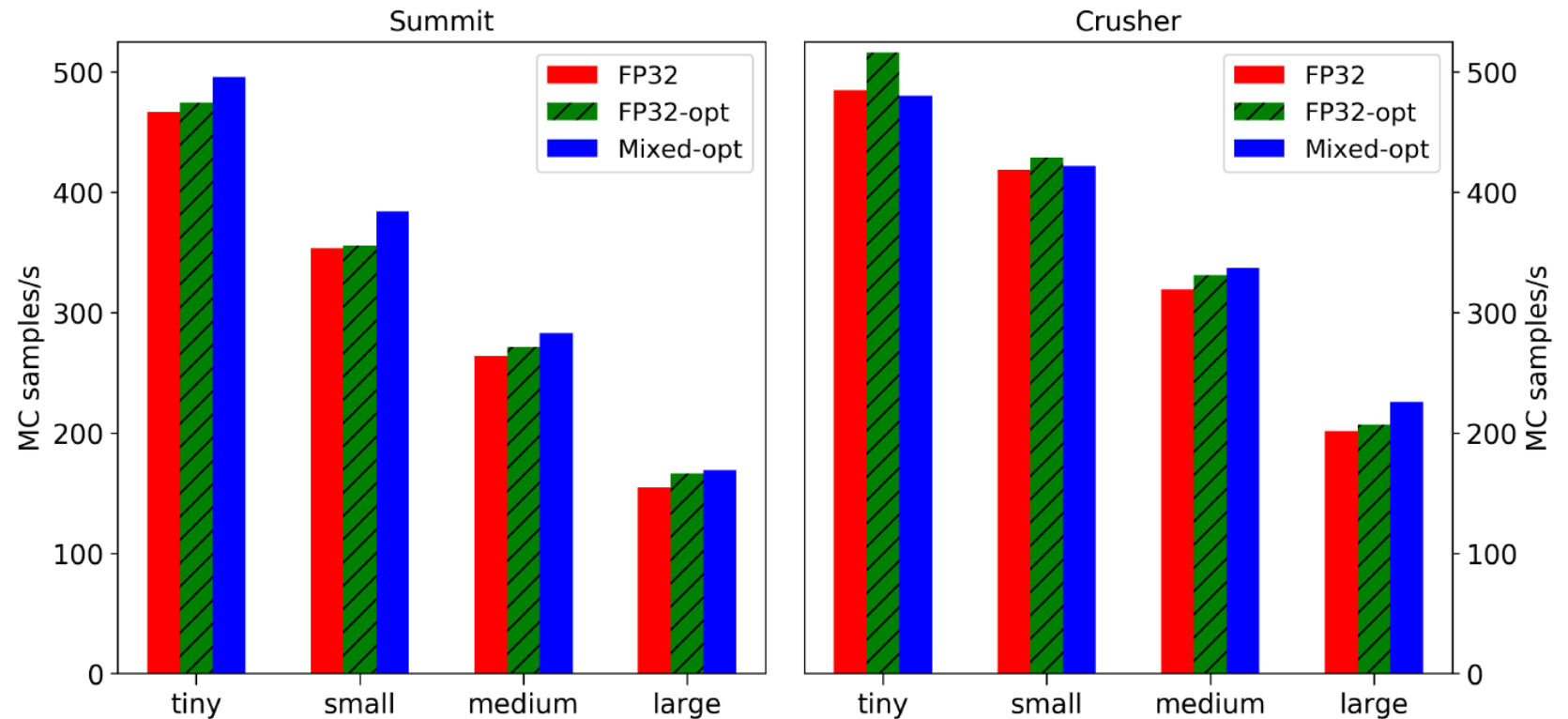
Implementation 2 Use TensorFlow Serving

```
1: function ENERGY
2:   Host ← server_ip : port
3:   for i ← 0, Nelements do
4:     Stubi ← tensorflow::serving::PredictionService::NewStub(
      grpc::CreateChannel(Host, grpc::InsecureChannelCredentials()))
5:   end for

6:   for i ← 0, Nelements do           ▷ Add energy by element species
7:     tensorflow::serving::PredictRequest predictRequest
8:     tensorflow::serving::PredictResponse Outputs
9:     grpc::ClientContext context
10:    predictRequest.mutable_model_spec().set_name(Modeli)
11:    predictRequest.mutable_model_spec().set_signature_name("serving_default")
12:    for j ← 0, Nneighbors do
13:      predictRequest.mutable_inputs() ← Atom[i][j]
14:    end for
15:    Stubi.Predict(&context, predictRequest, &Outputs)
16:    Energy ← Outputs
17:  end for
18:  return Energy
19: end function
```

Simulation-ML Integration

- 1.2x per GCD over V100



Strategies for Integrating Deep Learning Surrogate Models with HPC Simulation Applications, IPDPSW 2022

Questions ?

yinj@ornl.gov