

HPC and AI Together for Science Campaigns

Wes Brewer Arjun Shankar Feiyi Wang Junqi Yin

Advanced Technologies Section

National Center for Computational Sciences (NCCS)/ Oak Ridge Leadership Computing Facility (OLCF)

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



Agenda

- Opportunities for AI integration with HPC
 - Interleaved HPC and AI
 - In Situ: Co-Operating Mod-Sim and Al Models
 - Ex Situ: Edge-to-Exascale AI workflow vignette
- Specific techniques integrating scalable operations of HPC Mod-Sim and AI
 - Client-Server for Al inference supported campaigns
 - Scaling performance considerations: tight-to-loose coupling performance characteristics



Design Pattern: Interleaved Mod-Sim + ML/DL/AI at Scale





3

Converged Platform for HPC and AI with Analogs in Learning and HPC Simulation Stacks



CAK RIDGE

Junqi Yin, et al., Comparative evaluation of deep learning workloads for leadership-class systems,

Bench Council Transactions on Benchmarks, Standards and Evaluations, 2021, https://doi.org/10.1016/j.tbench.2021.100005d AI Together for Science Campaigns, 4/26/2023

HPC and AI are often treated separately







A Survey of Al4Science Workflow Applications, Middleware, and Performance, In preparation, A. Gainaru, S. Jha, F. Wang, F. Suter, W. Brewer, M. Emani, A. Shankar, et al.

Examples:

- ٠
- High-throughput virtual laboratory for drug discovery: <u>http://dx.doi.org/10.1177/10943420211001565</u> Intelligent Resolution: Integrating Cryo-EM with AI-driven Multi-resolution Simulations to Observe the SARS-CoV-2 ٠ Replication-Transcription Machinery in Action; https://doi.org/10.1101/2021.10.09.463779



Al Opportunities to Accelerate Campaigns

• Simulation steering

Example: L. Ward et al., "Colmena: Scalable Machine-Learning-Based Steering of Ensemble Simulations for High Performance Computing," 2021 IEEE/ACM Workshop on Machine Learning in High Performance Computing Environments (MLHPC), St. Louis, MO, USA, 2021, pp. 9-20

 Physics Integrated Deep Learning

Example: A predictor-corrector deep learning algorithm for high dimensional stochastic partial differential equations; Zhang et al., <u>https://arxiv.org/abs/2208.09883</u>

Fourier Neural Operator for Parametric Partial Differential Equations, Li et al., https://arxiv.org/abs/2010.08895

Commerical tools: NVIDIA SimNet

CAK RIDGE

National Laboratory

Domain-specific surrogates, data reduction

Example: DeePMD: DeePMD-kit: A deep learning package for many-body potential energy representation and molecular dynamics, https://arxiv.org/abs/1712.03641

Motif	Definition	Example
fault detection	detect algorithmic or other failure in execution, send signal for auto-	detect simulation defect caused by execution error
	matic or manual remediation	
math/cs algorithm	ML is used to enhance some mathematical (non-science-proper) com-	solver's linear system dimension is reduced based
	putation	on machine-learned parameter
submodel	a (proper) subset of a science computation is replaced by an ML model.	physics-based radiation model in a climate code
	molecular dynamics (MD) potentials as special case	replaced by ML model
steering	automatic steering of the direction of a computation for some internal	ML method to guide Monte Carlo sampling to in-
	process	clude undersampled regions
surrogate model	full science model replaced by ML approximation that captures impor-	data from tokamak simulation runs used to train
_	tant aspects, used for speed or science understanding	surrogate model
analysis	results from modeling and simulation (modsim) runs are analyzed by	use graph neural networks to analyze results of MD
	a human using ML methods	simulation
ML + modsim loop	both ML and traditional modsim, coupled	MD in loop used to refine deep learning model via
		active learning
classification	"pure" ML with little or no modsim used to classify some phe-	deep neural network inference to detect rare astro-
	nomenon; includes some other methods like reinforcement learning	physical event
various	umbrella project with multiple unrelated subprojects using possibly	CAAR/ESP/NESAP application readiness
	different kinds of AI/ML	
undetermined	manner of AI/ML use is undetermined	project is exploring AI/ML use but gives no details



Learning to Scale the Summit: AI for Science on a Leadership Supercomputer, Joubert W., et al., 2022, IPDPSW; <u>https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9835678</u>

OLCF User Call – HPC and Al Together for Science Campaigns, 4/26/2023

Example: Al Surrogates Reducing Costly Calculations MC-based Exploration of High-Entropy Alloy System (MoNbTaW)

- Probability of N atoms in configuration X at temperature T follows Boltzmann's distribution $exp(-E(X)/k_{\rm B}T)$ where E is the total configuration energy and $k_{\rm B}$ is the Boltzmann constant.
- Replica of the alloy systems at various *T* are simulated via replica exchange Monte Carlo simulations with transition probability between replica *m* and *n*:

 $W(\{X_m, T_m\} | \{X_n, T_n\}) = \min[1, \exp(-\triangle)],$

where

 $\triangle = (1/k_B T_n - 1/k_B T_m)(E(X_m) - E(X_n))$

 Atoms (i,j) are exchanged with acceptance probability P_{i,i} proportional to:

CAK RIDGE National Laboratory

 $\min\left[1, \exp(-(E(\{x_i, x_j\}) - E(\{x_j, x_i\}))/k_BT)\right]$



https://doi.org/10.1038/s43588-021-00139-3

DEEP LEARNING SURROGATE MODELS FOR THE ENERGY EVALUATION OF MONBTAW ALLOY. THE ARCHITECTURE IS GIVEN BY THE NUMBER OF NODES IN EACH HIDDEN LAYER, AND IT IS THE SAME FOR EACH OF THE FOUR ELEMENT. THE R^2 SCORE IS THE AVERAGED MODEL PERFORMANCE OF ALL ELEMENTS.

name	architecture	# parameters	$\stackrel{R^2 \text{ score}}{\text{FP32} \text{ mixed}}$	
	1	1		
tiny	24-24	5,257	0.991	0.988
small	128x2-64x2 -24x2	55,817	0.991	0.992
mdedium	200x11	440,801	0.994	0.992
large	512x6-256x6 -128x6-64x6-32	2,019,009	0.994	0.993

J. Yin, F. Wang, A. Shankar, Strategies for Integrating Deep Learning Surrogate Models with HPC Simulation Applications, IPDPSW 2022. https://ieeexplore.ieee.org/document/9835386

Deployment Architectures

Tightly coupled



Loosely coupled



Semi tightly coupled







9

OLCF User Call – HPC and AI Together for Science Campaigns, 4/26/2023

HPC-AI for Inference on Summit



Brewer, Wesley, et al. "Production deployment of machine-learned rotorcraft surrogate models on HPC." MLHPC'21.

https://code.ornl.gov/whb/osmi-bench



OLCF User Call – HPC and AI Together for Science Campaigns, 4/26/2023

TF Serving on Summit Export TF Model

from tensorflow.keras.models import Model
model = Model(...)
model.save("/path/to/mymodel/1")

Start TF Serving

- module load open-ce/1.1.3-py38-0
- module load cuda/11.0.2
- tensorflow_model_server [--port 8500]
 --model_base_path=/path/to/mymodel
 --model_name="mymodel"
- Isof -i :8500 # check that TF Serving is running and listening to port 8500

Test implementation

- saved_model_cli show --all --dir 1
- python client.py

client.py

https://code.ornl.gov/whb/osmi-bench

import grpc import tensorflow as tf from tensorflow_serving.apis import predict_pb2 from tensorflow serving apis import prediction service pb2 grpc # set parameters batch size = 32input_shape, output_shape = (batch_size, 8, 48), (batch_size, 2, 12) mname, sname = 'mymodel', 'serving default' # names of model and signature iname, oname = 'inputs', 'dense_3' # names of input and output layers # create some random data and perform inference data = np.array(np.random.random(input shape, dtype=np.float32]) channel = grpc.insecure channel('localhost:8500') stub = prediction_service_pb2_grpc.PredictionServiceStub(channel) request = predict_pb2.PredictRequest() # request proto data structure request.model_spec.name = mname request.model_spec.signature_name = sname request.inputs[iname].CopyFrom(tf.make_tensor_proto(data, shape=input_shape) result = stub.Predict(request) # send request # extract results result array = np.reshape(result.outputs[oname].float val, output shape)

Approaches for Scaling Up

HAProxy Load Balancer



> singularity pull docker://haproxy

Brewer, Wesley, et al. "Production deployment of machine-learned rotorcraft surrogate models on HPC." *MLHPC*²1.



Boyer, Mathew, et al. "Scalable Integration of Computational Physics Simulations with Machine Learning." *AI4S'22.*



Example: Rotorcraft Aerodynamics

1E6 inferences per second

Extract velocities **CFD Simulation** Inject source terms into CFD —O— Throughput/Small —O— Throughput/Medium —O— Throughput/Large grid via Gaussian projection from CFD grid T2NI/Medium ···· T2NI/Small ····• T2NI/Large 1E+6 600 GPU (s) 500 1E+5 Throughput (samples/s) C_L , C_D , C_M 400 1E+4 Compute force coefficients 1E+3 300 to N=32768 infe Map 2D output to 1E+2 200 3D blade geometry MLS Model (CNN/Surf) 100 Dense-FCN 1E+1 Time 1E+0 0 Latent Nodes: 1 16 32 8 Space 1D Vector $[\dot{x}, \dot{y}, \dot{z}, u, v, w]$ 2D-CNN Encoder $[C_p]$ GPUs: 6 192 2D-CNN Decoder 12 24 48 96

Workflow

Weak scaling performance on Summit

Brewer, Wesley, et al. "Production deployment of machine-learned rotorcraft surrogate models on HPC." MLHPC 21.

Strategy: Maximize batch size, Concurrency = 2



13

Example: Machine-learned Boundary Conditions



Fig. 4. a) Machine-learned boundary condition implementation design. b) Flowchart for machine-learned boundary condition implementation in Orchard.

Fig. 5. Per-GPU throughput of TF Serving (TFS) and RedisAI (RAI) C++ clients for a SimNet fully connected neural network with about 100k parameters.

Boyer, Mathew, et al. "Scalable Integration of Computational Physics Simulations with Machine Learning." AI4S'22.



Rankings for RedisAl vs. TF Serving

TF SERVING VS. REDISAI VS. EMBEDDED FOR BOTH C++ AND PYTHON CLIENTS ON T4 GPU (BATCH SIZE 1024).

Rank	Approach	Avg. Latency (s)	99th Perc. Latency (s)	Throughput (samples/s)	Distributed
1	Embedded C++	0.90	0.90	1137	No
2	TF Serving with C++ client	1.11	1.14	922	Yes
3	Embedded Python	1.16	1.20	881	No
4	RedisAI with C++ client	1.25	1.26	809	Yes
5	RedisAI with Python client	1.31	1.34	768	Yes
6	TF Serving with Python client	2.50	2.53	410	Yes

Brewer, Wesley, et al. "Production deployment of machine-learned rotorcraft surrogate models on HPC." MLHPC'21.



OLCF User Call – HPC and Al Together for Science Campaigns, 4/26/2023

Application Example: DL-accelerated Monte Carlo

- Architecture
 - Application
 - Framework
 - DL stack
- Application layer is independent of DL stack, and built upon
 - TF C++
 - SmartRedis

	VAE modeling	Parallel Tempering		
Application	Wang-Landau Sampling			
	rning sal	Deep Learning Proposal		
Framework	martRedis	TensorFlow SmartRedis		
DL Stack	platform	CUDA, ROCm platform		
	PU	GPU, C		

https://doi.org/10.5281/zenodo.6612174, IPDPS23



SmartRedis Implementation

• Launch server:

redis-server --port --loadmodule

- Model can be dynamically added
- Data save to in-memory DB

Implementation 3 Use SmartRedis/RedisAI

- 1: function ENERGY
- 2: $SSDB \leftarrow redis_host_ip : port$
- SmartRedis::Client* clients[Nelements]
- 4: for $i \leftarrow 0$, Nelements do
- 5: Client_i.set_model_from_file(model_i, path_i, "TF", "GPU", batch_size, mini_batch_size, "hea", {"x"}, {"Identity"})
- 6: end for

7:	for $i \leftarrow 0$, Nelements do	▷ Add energy by element species
8:	uint8 Data[]	
9:	data_key ← "data_i"	
10:	outputs_key ← "energy_i"	
11:	for $j \leftarrow 0$, Nneighbors do	
12:	$Data \leftarrow Atom[i][j]$	
13:	end for	
14:	Client _i .put_tensor(data_key, Data	, NT[t], Nneighbors*(NE-1), SmartRe-
	dis::TensorType::flt, SmartRedis::MemoryL	ayout::contiguous))
15:	$Client_i.run_model(model_i, \{data)$	ta_key}, {outputs_key})
16:	Client _i .unpack_tensor(outputs_ke	y, Outputs.data(), NT[t], SmartRe-
	dis::TensorType::flt, SmartRedis::MemoryL	ayout::contiguous)
17:	$Energy \leftarrow Outputs$	
18:	end for	
19:	return Energy	
20:	end function	

TF C++ Implementation

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

Im	plementation 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,	
	ModelPath, {"serve"}, &Model)	▷ load pre-trained model
6:	return Model	
7:	end function	
8:	function PREDICT(tensorflow::SavedModelBundle M	Iodel, tensorflow::Tensor In-
9:	tensorflow::Tensor Inputs	⊳ input tensor
10:	vector <tensorflow::tensor> Outputs</tensorflow::tensor>	
11:	string InputNode	
12:	string OutputNode	▷ input and output nodes
13: 14: 15:	Model.GetSession().Run(data, {OutputNode}, {} return Outputs end function	, &Outputs)
16.	function ENERGY	
17.	fon i (0 Nolomonto do	
12.	$10F i \leftarrow 0, Neiements 00$ $Model \leftarrow LOADMODEL(Path)$	
18. 19:	$model_i \leftarrow LOADMODEL(Path_i)$ end for	
20.	for i (0 Nolomento do) A	ld anaray by alamant spacios
20.21.	tensorflow: TensorShane DataShane	id energy by element species
21.	tensorflowTensorInputs(tensorflowDT_LUN	TP DataShana)
22.	for <i>i</i> (0 <i>N m ci a b b m c d c</i>	(16, DataShape)
25.	for $j \leftarrow 0$, is neighbors do	
24.	$inputs \leftarrow Atom[i][j]$	
25:	end for $Outputs = DPEDICT(Model Inputs)$	
20:	$Suppose = PREDICT(Model_i, inputs)$	
27:	$Energy \leftarrow Outputs$	
20:	end for	
29:	return Energy	
30:	end function	



18

Implementation Considerations

OVERVIEW OF THE CHARACTERISTICS OF DIFFERENT DEPLOYMENT METHODS FOR SURROGATE MODELING IN SIMULATIONS.

deployment	usability		performance		
method	portability	flexibility	concurrency	latency	throughp
TensorFlow C++ API TensorFlow Serving SmartRedis/RedisAI					



HPC-AI Usage Ex Situ: Distributed System





J. Yin, G. Zhang, H. Cao, S. Dash, B. Chakoumakos, and F. Wang, "Toward an Autonomous Workflow for Single Crystal Neutron Diffraction." SMC 2022

OLCF User Call – HPC and Al Together for Science Campaigns, 4/26/2023

Demonstration



CAK RIDGE National Laboratory

21

OLCF User Call – HPC and AI Together for Science Campaigns, 4/26/2023

HPC + AI Training Resources

- Upcoming: AI and Model Deployment
 - <u>https://www.olcf.ornl.gov/calendar/ai-for-science-at-scale-intro/</u>, June 15th, 2023
 - <u>https://www.olcf.ornl.gov/calendar/smartsim-at-olcf/</u> HPE Training, July 13th, 2023
- Al Model Training
 - AI on Frontier: <u>https://olcf.ornl.gov/wp-content/uploads/2-16-23_AIonFrontier.pdf</u>
 - Data and Viz: <u>https://www.olcf.ornl.gov/calendar/data-visualization-and-analytics-</u> <u>training-series-jupyter-workflow-at-olcf/</u>
- Hyperparameter exploration
 - Link to OLCF training: <u>https://www.olcf.ornl.gov/wp-content/uploads/Jupyter_DL_Workflow.pdf</u> (RayTune)



Thank you!

We are Hiring! Visit: https://jobs.ornl.gov/

CAK RIDGE

DLCF User Call – HPC and AI Together for Science Campaigns, 4/26/2023