PYTORCH AMPERE® OPTIMIZED FRAMEWORK Documentation

AMPERE COMPUTING



Table of Contents

RELEASE NOTES	1
OVERVIEW	3
PYTORCH FRAMEWORK	3
Versions Compatibility	
PYTHON	3
CONFIGURATIONS	3
QUICKSTART	4
Launching Docker Container	5
Running Examples AMPERE OPTIMIZED PYTORCH PROGRAMMING GUIDE	
Overview	
Supported Inference Ops	7
PyTorch JIT Trace	
Torch Compile (beta)	9
Threading	
Programming Tips	

RELEASE NOTES

V1.9.0:

- libampere-aio updated to 0.10.0
 - New upsampling2d_nearest kernel to speed up YOLO models
- Experimental support of int8 FC to speed up models like DLRM
- Improved support of in place operators by the Pytorch-aio
- Updated to PyTorch 2.1
- Framework now is setting AIO_SKIP_MASTER_THREAD=1 env var by default, no need to specify it

V1.8.0:

- libampere-aio updated to 0.9.0
 - Resolve an issue when input shape changes frequently
 - Performance enhancements of some NLP models

V1.7.0:

- libampere-aio updated to 0.8.0
- Bug fixes and performance enhancements
- Improved memory management
- Improved model compilation times
- Improved algorithm calculating graph handled by AIO
- New operators supported: Baddbmm, sub, slice, max (elementwise), min (elementwise), neg, index (some cases), split_with_sizes, NumToTensor, Float, Int
- Options custom argument in AIO torch.dynamo backend (see below)

V1.6.0:

- libampere-aio updated to 0.7.0
- PyTorch updated to 2.0.0
- Bug fixes and performance enhancements
- New operator supported: LogSoftmax
- Torch.compile() supported (see section about it)

V1.5.2:

- Note: v1.5.2 is a bug fix release to v1.5.1 and v1.5.0. It fixes an issue related to YOLO models. Please discard v1.5.0 and v1.5.1 you've installed.
- libampere-aio updated to 0.6.1
- Bug fixes and performance enhancements
- New operators supported: Split, Chunk, Sqrt, Rsqrt, Exp, Log, Zeros_Like, Embedding, Mean
- TorchScript loops are supported.
- Improved lifetime handling of Torchscript models

V1.4.0:

- libampere-aio updated to 0.5.0
- Pytorch framework updated to 1.12.1 from 1.11.0
- Support of FP16 ops (automatic mode)
- New operators supported: deconv2d, embedding bag
- Improved memory management
- Bug fixes: Instance Norm op fix, thread safety

V1.3.0:

- Binary integer operations support.
- libampere-aio updated to 0.4.0
- New operators supported: Reshape, Squeeze, Unsqueeze, Flatten, PixelShuffle, GroupNorm, InstanceNorm.
- Using custom compiled OpenBLAS, as Pytorch BLAS backend.
- Bug fixes

V1.2.0:

- libampere-aio updated to 0.3.0
- New optimized operators: Gelu, Silu, Softmax, Div, Binary ops between Tensor and Scalar, Permute, View, Layer Norm, Size, Pow, Tanh, Sigmoid
- Improved Concat support
- Graph optimizations
- Various bugfixes

V1.1.0:

- Libampere-aio updated to 0.2.1
- Batch Matmul supported (enhancing DLRM performance)
- Adaptive Avg Pool supported
- LeakyRelu supported
- AIO_NUM_THREADS no longer needed to set Ampere Optimized PyTorch threads, inherits Pytorch intra-op thread count.

OVERVIEW

Ampere Optimized PyTorch inference acceleration engine is fully integrated with the PyTorch framework. PyTorch models and software written with the PyTorch API can run as-is, without modifications.

PYTORCH FRAMEWORK

Python is installed with Ampere Optimized PyTorch and all dependencies. No additional installation steps are needed.

Versions Compatibility

This release is based on Pytorch 2.0.0 and comes with the compatible Torchvision 0.15.1 installed.

PYTHON

Pytorch 2.0.0 is built for Python 3.10, supporting Ubuntu 22.04. Regarding other Python versions, please contact your Ampere sales representative. If you are using the software through a third party, contact their customer support team for help. You can also contact the AI team at ai-support@amperecomputing.com.

CONFIGURATIONS

Ampere Optimized PyTorch inference engine can be configured by a set of environment variables for performance and debugging purposes. They can be set in the command line when running Pytorch models (e.g., AIO_NUM_THREADS=16 python run.py -p fp32) or set in the shell initialization script.

AIO_PROCESS_MODE

This variable controls whether the Ampere Optimized PyTorch inference engine is used to run the Pytorch model:

- 0: disabled.
- 1: enabled (Default).

AIO_CPU_BIND

Enables core binding. If enabled, each Ampere Optimized PyTorch thread will bind itself to a single core:

- 0: Core binding disabled.
- 1: Core binding enabled (Default).

AIO_MEM_BIND

Binds memory to NUMA (Non-uniform memory access) node 0. For optimal performance, numactl (https://linux.die.net/man/8/numactl) is preferred. numactl bind will affect both the Pytorch framework and the optimized framework buffers, while the optimized framework is unable to affect buffers allocated by the Pytorch framework:

- 0: Membind disabled.
- 1: Membind to node 0 (Default).

AIO_NUMA_CPUS

Select the cores that Ampere Optimized PyTorch should bind to (if CPU_BIND is enabled):

- Not set: use the first N cores of the machine, excluding hyper-threaded (Default).
- Set: use N first cores from the list of cores for N threads. The list is in space separated, 0-based number format. For example, selecting cores 0 to 1: AIO_NUMA_CPUS="0 1".

AIO_DEBUG_MODE

Control the verbosity of debug messages:

- 0: No messages
- 1: Errors only
- 2: Basic information, warnings, and errors (Default)
- 3: Most messages
- 4: All messages

QUICKSTART

The following instructions run on Altra/Altra Max Linux machines installed **with Docker**. When you are already using a virtual machine pre-installed with the version of Ampere Optimized PyTorch (e.g. on a cloud service provider) that you need, you can skip the following step of launching Docker container.

Note: This docker image is developed for benchmarking and evaluation purpose, not for deployment into production environment. We will provide required Debian, RPM and Python packages as needed for your production deployment.

Launching Docker Containermas

Pulling Docker Image from Docker Hub repository

\$ docker pull amperecomputingai/pytorch:1.9.0

Launching Docker Container

\$ docker run --privileged=true --rm --name pytorch-aio --network host -it amperecomputingai/pytorch:1.9.0

Warning: This user has, by default, root privileges with Docker. Please limit permission according to your security policy.

Running Examples

You can try Ampere Optimized PyTorch by either running the Jupyter Notebook examples or Python scripts on the CLI level.

To run the Jupyter Notebook QuickStart examples follow the instructions below:

Set AIO_NUM_THREADS to the requested value first.

\$ export AIO_NUM_THREADS=16; export OMP_NUM_THREADS=16
\$ cd /workspace/aio-examples/
\$ bash start_notebook.sh

If you run the Jupyter Notebook Quickstart on a cloud instance, make sure your machine has port 8080 open and on your local device run:

\$ ssh -N -L 8080:localhost:8080 -I <ssh_key> your_user@xxx.xxx.xxx

Use a browser to point to the URL printed out by the Jupyter Notebook launcher.

You will find Jupyter Notebook examples (examples.ipynb) under the /classification and /object detection folders.

The examples run through several inference models, visualize results they produce, and present the performance numbers.

To use CLI-level scripts:

Set AIO_NUM_THREADS to the requested value first.

\$ export AIO_NUM_THREADS=16; export OMP_NUM_THREADS=16
\$ cd /workspace/aio-examples/

Go to the directory of choice, e.g.

\$ cd classification/resnet_50_v1

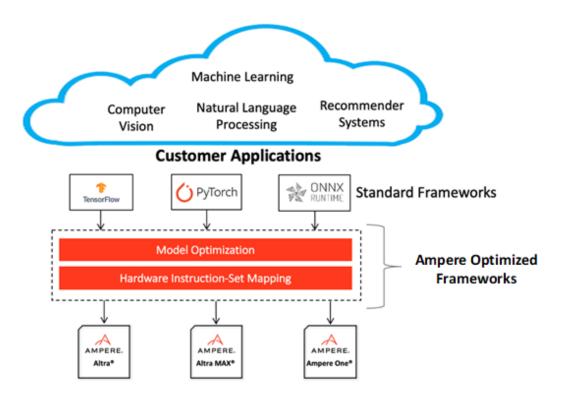
Evaluate the model.

\$ numactl --physcpubind=0-15 python3 run.py -p fp32

AMPERE OPTIMIZED PYTORCH PROGRAMMING GUIDE

Overview

Ampere Optimized PyTorch is powered by Ampere[®] AI backend that accelerates Deep Learning (DL) operations on the Ampere[®] Altra family of processors. Ampere Optimized PyTorch accelerates DL operations through model optimization, highly vectorized compute kernels and multi-thread operations that are automatically tuned to deliver the best latency and throughput on Ampere Altra processors. It delivers 2-5x gains over alternative backend solutions.



Supported Inference Ops

Ampere Optimized Pytorch accelerates the most common Pytorch ops that are used in diverse types of models. Here is a list of accelerated ops and formats (Note: non-accelerated ops will still run without a problem, at the original framework operator speed):

Layer	FP32	Explicit FP16 (Model defined)	Implicit FP16 (Automatic on- the-fly conversion)	Int8 (Automatic on-the-fly) [Experimental]	Notes
Conv2d	Y		Y		
Deconv2d	Y				Without bias
Linear	Y	Y	Y	Y	
MaxPool2d	Y				
AvgPool2d	Y				
AdaptiveAvgPool2d	Y				
Relu	Y		Y		
Relu6	Y		Y		
LeakyRelu	Y		Y		
Softmax	Y		Y		
LogSoftmax	Y		Y		
Gelu	Y		Y		
Silu	Y		Y		
Sigmoid	Y				
Tanh	Y		Y		
Transpose	Y		Y		
Permute	Y		Y		
BatchNorm	Y		Y		
LayerNorm	Y				
GroupNorm	Y				
InstanceNorm	Y				
Add	Y	Y	Y		Int version not optimized

Mul	Y	Y	Y		Int version not optimized
Div	Y	Y	Y		Int version not optimized
Pow	Y	Y	Y		Int version not optimized
Matmul	Y		Y	Y	Limited support 2 nd argument have to be 2d matrix
MM	Y		Y		
BMM	Y		Y		
PixelShuffle	Y				
View	Y	Y	Y		
Reshape	Y	Y	Y		
Squeeze	Y	Y	Y		
Unsqueeze	Y	Y	Y		
Flatten	Y	Y	Y		
Contiguous	Y		Y		
Size	Y	Y	Y		One dimension case
EmbeddingBag	Y	Y	Y		Sum mode
Embedding	Y		Y		
Split	Y		Y		
Chunk	Y		Y		
Sqrt	Y		Y		
Rsqrt	Y		Y		
Exp	Y		Y		
Log	Y		Y		
Zeros_like	Y				
Mean	Y		Y		
Baddbmm	Y		Y		

Slice	Υ	Y	
Select	Y	Y	
Neg	Y	Y	
Split with sizes	Y	Y	
Index	Y	Y	Limited support
Max	Y	Y	Elementwise
Min	Y	Y	Elementwise
Sub	Y	Y	
UpsampleNearest2 d	Y	Y	Support only constant scale factor and output size (cannot change in model lifetime)

PyTorch JIT Trace

While Pytorch Eager Execution provides excellent model building, programming, and debugging experience, it is slower than graph execution. So, Torchscript is typically used for inference deployment. In the current version of Ampere Optimized Pytorch, Torchscript mode is also accelerated.

To use Ampere Optimized Pytorch, conversion of Pytorch module to Torchscript is needed. There are two ways to convert: torch.jit.script() or torch.jit.trace(input) API calls. See https://pytorch.org/docs/stable/jit.html for more details. After converting to Torchscript user should call torch.jit.freeze() to freeze the models and enable model optimizations for inference.

Torch Compile (beta)

Ampere Optimized Pytorch support torch.compile API introduced in Pytorch 2.0 release. This is a new mode for optmizing model for infenence. To take advantage of it user must compile the model with AIO backend by using compiled_model = torch.compile(model, backend="aio", options={"modelname": "model"}). It is important to explicitly select "aio" backend and pass additional parameter named options with "modelname" field. See https://pytorch.org/get-started/pytorch-2.0/ for more information.

Note: In this release this is a beta feature. Torchscript is likely to be faster than torch.compile.

Implicit FP16 mode

Ampere Optimized PyTorch backend now provides automatic FP16 operator conversion that can boost the performance of your FP32 model on-the-fly. It automatically performs FP16 conversion and computation for certain whitelisted operators through regular expression. To take advantage of that, you can set environment variable.

\$export AIO_IMPLICIT_FP16_TRANSFORM_FILTER=".*"

This activates automatic FP16 conversion for all supported operators. It is estimated that this has minor impact on the accuracy of common models. Please contact us if you have any questions about this feature.

Implicit INT8 mixed precision mode (experimental)

Another performance boost can be achieved with implicit model quantization. It is enabled with:

\$export AIO_QUANTIZE_INT8="ALL"

It will replace supported operations with their quantized INT8 versions. At this moment only LINEAR and limited cases of MATMUL layers are supported. Please mind this is a new and experimental feature. Interface may change in future.

It is possible to use this mode along with implicit FP16 mode, however sometimes it causes slowdowns compared to the FP32/INT8 mode.

Threading

Ampere Optimized PyTorch controls the number of Ampere Optimized Pytorch intra_op threads with torch.set_num_threads(). This controls both the number of threads used for ops delegated to Ampere Optimized Pytorch as well as the ops running on default CPU backend.

Some default CPU backend ops (non-AIO) also need to set OMP_NUM_THREADS environment variable to control the intra_op threads.

Programming Tips

In the first two inference passes, Ampere Optimized Pytorch performs a runtime compilation of PyTorch script and prepares Ampere Optimized Pytorch network. So, the latency of the first two passes is expected to be longer. Subsequent passes will be accelerated.

Ampere Optimized PyTorch provides much better latency scaling as core count increases, compared to other platforms. You can easily try the optimal number of cores with the above set_num_threads() function that can give you the best price / performance, while meeting your latency requirements.

Models are optimized for the shape of the tensors that are used during the compilation phase (see above). Passing different shape tensors will work but is suboptimal. To get the best performance pad varying shape tensors when running inference.

If any issues occur, Ampere AI team is ready to help. Typically, the first step is to get more debug logs and send them to ai-support@amperecomputing.com. Please set environment variable AIO_DEBUG_MODE=4 to capture low level logs.

Limitations

Ampere Optimized PyTorch doesn't support dynamic ranks of tensors (different rank in subsequent passes). Dynamic shapes of Tensors are supported but not recommended, ideally one should pad inputs to the network to get best performance.

We can also provide more in-depth profiling of your model to help enhance performance to meet your needs.

Ampere Computing[®] / 4655 Great America Parkway, Suite 601 / Santa Clara, CA 95054 / www.amperecomputing.com

Ampere Computing, the Ampere Computing logo, Altra, and eMAG are registered trademarks of Ampere Computing.

Arm is a registered trademark of Arm Holdings in the US and/or elsewhere. All other trademarks are the property of their respective owners.

©2022 Ampere Computing. All rights reserved.

AMP 2019-0039