

## Introduction to ML Frameworks

Alessandro Fanfarillo, PhD Frontier Center of Excellence SMTS Software System Design Eng. December 2022

### **CAUTIONARY STATEMENT**

This presentation contains forward-looking statements concerning Advanced Micro Devices, Inc. (AMD) such as the features, functionality, performance, availability, timing and expected benefits of AMD's current products, future products and markets, which are made pursuant to the Safe Harbor provisions of the Private Securities Litigation Reform Act of 1995. Forward-looking statements are commonly identified by words such as "would," "may," "expects," "believes," "plans," "intends," "projects" and other terms with similar meaning. Investors are cautioned that the forward-looking statements in this presentation are based on current beliefs, assumptions and expectations, speak only as of the date of this presentation and involve risks and uncertainties that could cause actual results to differ materially from current expectations. Such statements are subject to certain known and unknown risks and uncertainties, many of which are difficult to predict and generally beyond AMD's control, that could cause actual results and other future events to differ materially from those expressed in, or implied or projected by, the forward-looking information and statements. Investors are urged to review in detail the risks and uncertainties in AMD's Securities and Exchange Commission filings, including but not limited to AMD's most recent reports on Forms 10-K and 10-Q.

AMD does not assume, and hereby disclaims, any obligation to update forward-looking statements made in this presentation, except as may be required by law.















# CDNA/CDNA2 Architectures

# **COMPUTE GPU ARCHITECTURE ROADMAP**



### 2<sup>nd</sup> GENERATION MATRIX CORES

**OPTIMIZED COMPUTE UNITS FOR SCIENTIFIC COMPUTING** 



DOUBLE PRECISON (FP64) MATRIX CORE THROUGHPUT REPRESENTATION

### 2ND GENERATION CDNA ARCHITECTURE TAILORED-BUILT FOR HPC & AI





# AMDA ROCM 5.0 Democratizing exascale for all

EXPANDING	OPTIMIZING	ENABLING
SUPPORT & ACCESS	PERFORMANCE	DEVELOPER SUCCESS
<ul> <li>Support for Radeon Pro</li></ul>	<ul> <li>MI200 Optimizations: FP64</li></ul>	<ul> <li>HPC Apps &amp; ML Frameworks</li></ul>
W6800 Workstation GPUs	Matrix ops, Improved Cache	on AMD InfinityHub
<ul> <li>Remote access through the</li></ul>	<ul> <li>Improved launch latency and</li></ul>	<ul> <li>Streamlined and improved</li></ul>
AMD Accelerator Cloud	kernel performance	tools increasing productivity

02

# Machine Learning, Deep Learning & Al



### Machine Learning, Deep Learning, AI, RL, ...



**Artificial Intelligence** (AI): techniques enabling computers to mimic human behavior

**Machine Learning** (ML): techniques that give computers the ability to learn without being explicitly programmed (learning from data)

**Neural Networks** (NN): models and techniques that teach computers to process data like the human brain (able to model non-linear distributions)

**Deep Learning** (DL): use of NNs with three or more layers (able to model very complex distributions)

**Reinforcement Learning** (RL): use models (ML/DL) to learn how to achieve a goal in unexplored/unpredictable environments (learning based on rewards and penalties)

## **AI/ML ECOSYSTEM SUPPORT & ADOPTION**

**UPSTREAMED SOURCE & BINARY SUPPORT** 



# **ML FRAMEWORKS & LIBRARIES**

UPSTREAMED SOURCE & BINARY SUPPORT ALLOW SCIENTISTS TO EASILY USE EXISTING CODE

	Source	Container	PIP Wheel
TensorFlow	TensorFlow GitHub	Infinity Hub	pypi.org
<sup>с</sup> РуТогсh	PyTorch GitHub	Infinity Hub	pytorch.org
ONNX RUNTIME	ONNX-RT GitHub	Docker Instructions	onnxruntime.ai
JAX	GitHub public fork	Docker Hub	Est 2022
DeepSpeed	DeepSpeed GitHub	Docker Hub	deepspeed.ai
CuPy	<u>cupy.dev</u>	Docker Hub	<u>cupy.dev</u>



## Focused on Targeted ML Use-Cases

Most common models on HuggingFace supported on AMD platforms today



#### **VIDEO & IMAGE RECOGNITION**

Optimized Models Resnet, VGG, Inception GoogleNet, ResNext, Detectron2, RetinaNet, Mask R-CNN

#### <u>Markets</u>

Automotive/Self Driving Cars Healthcare/Medical Imaging Public Safety



#### LANGUAGE PROCESSING

<u>Optimized Models</u> GNMT, BERT, GPT-2, BART, DeBERTa, DistilBERT, RoBERTa, T5

<u>Markets</u> Customer Service Web Services/E-Commerce



#### **RECOMMENDATION ENGINE**

Optimized Models DLRM

#### <u>Markets</u> Web Services/E-commerce SaaS

# AMD INFINITY HUB

Ready-to-Use HPC/ML Containers

#### AMD Instinct<sup>™</sup> MI200 Support

Over 15 key applications & frameworks on Infinity Hub & a catalogue supporting over 50 applications, frameworks & tools

#### Accelerating Instinct<sup>™</sup> Adoption

Over 5000 application pulls since launch last year

#### PERFORMANCE RESULTS

Published Performance Results for Select Apps / Benchmarks

AMD.com/InfinityHub Instinct Application Catalog

AMDA PRODUCTS - SOLUTIONS - SHOP - DRIVER	S & SUPPORT
AMD Infinity Hub	
AMDR         PRODE         ROCM" LEARNING CENTER	Computational Science Starts Here The AMD Infinity Hub contains a collection of advanced GPU software containers ar researchers, scientists and engineers to speed up their time to science.
Amber Amber is a suite of biomolecular simulation programs. It began in the late 1970's, and is maintained by an active development community; see our history page and our contributors pa	Chroma The Chroma package supports data-parallel programming constructs for lattice field theory and in particular lattice QCD. It uses the SciDAC QDP++ data-parallel programming (in C++) that
MORE INFO	MORE INFO PULL TAG
	GROMACS
<b>GRID</b> Grid is a library for lattice QCD calculations that employs a high-level data parallel approach while using a number of techniques to target multiple types of parallelism. The library currently	GROMACS GROMACS is a versatile package to perform molecular dynamics, i.e. simulate the Newtonian equations of motion for systems with hundreds to millions of particles.

### AMD Instinct<sup>™</sup> GPUs & ROCm<sup>™</sup> SOFTWARE ECOSYSTEM USEFUL WEB RESOURCES

- ▲ AMD Instinct GPUs:
  - ▲ AMD Instinct<sup>™</sup> MI100 GPU page: <u>https://www.amd.com/en/products/server-accelerators/instinct-mi100</u>
  - ▲ AMD Instinct<sup>™</sup> MI210 GPU page: <u>https://www.amd.com/en/products/server-accelerators/instinct-mi210</u>
  - ▲ AMD Instinct<sup>™</sup> MI Series Product Page: <u>https://www.amd.com/en/graphics/instinct-server-accelerators</u>
  - ▲ AMD Instinct<sup>™</sup> HPC Solutions Page: <u>https://www.amd.com/en/graphics/servers-instinct-mi-powered-servers</u>
  - ▲ AMD Instinct<sup>™</sup> Machine Learning Solutions Page: <u>https://www.amd.com/en/graphics/servers-instinct-deep-learning</u>
  - AMD CDNA2 Architecture: <u>https://www.amd.com/en/technologies/cdna2</u>
  - CDNA2 WP: <u>https://www.amd.com/system/files/documents/amd-cdna2-white-paper.pdf</u>
- ▲ AMD ROCm<sup>™</sup> open software platform:
  - ▲ AMD ROCm<sup>™</sup> pages: <u>https://www.amd.com/en/graphics/servers-solutions-rocm</u>
  - AMD Infinity Hub: <u>https://www.amd.com/en/technologies/infinity-hub</u>
  - ▲ AMD ROCm<sup>™</sup> Deep Learning: <u>https://docs.amd.com/bundle/ROCm-Deep-Learning-Guide-v5.2/page/Deep\_Learning\_Training.html</u>
  - ROCm Information Portal (DOCs & Learning Ctr.): <u>https://docs.amd.com/</u>
- ▲ HPC & AMD page: <u>www.AMD.com/HPC</u>





# TensorFlow

### **TensorFlow and Keras**

- TensorFlow (TF) is an open-source library for solving problems of ML/DL/AI
- Focus on training and inference in neural networks
- Functional and OOP capabilities
- Support for multiple accelerators (<u>https://www.tensorflow.org/tutorials/distribute/keras</u>)
- Keras is a high-level API for building NNs (adopted in 2017)
- Two options to install a ROCm-aware version of TF:
  - 1. Via docker image
  - 2. Via wheel package (Python 3.7+ is required)

For more info on how to install TF + ROCm please visit:

https://docs.amd.com/bundle/ROCm-Deep-learning-Guide\_5.2/page/Frameworks\_Installation.html

Tutorials available at: <u>https://github.com/tensorflow/docs</u>

More examples available at: <u>https://github.com/tensorflow/examples</u>





# PyTorch

### **PyTorch**

- PyTorch is an optimized tensor library for deep learning using GPUs and CPUs.
- Object oriented (python-like) and able to work on multiple accelerators
- Popularity increasing in the AI & research communities
- Two options to install a ROCm-aware version of PyTorch:
  - 1. Via docker image
  - 2. Via wheel package
  - 3. AMD is one of the founding members of the <u>PyTorch foundation</u>

For more info on how to install PyTorch + ROCm please visit:

https://docs.amd.com/bundle/ROCm-Deep-learning-Guide\_5.2/page/Frameworks\_Installation.html

Tutorials available at : <u>https://github.com/pytorch/tutorials</u>

Examples available at: <u>https://github.com/pytorch/examples.git</u>

# **PyTorch 1.12.0**

AMD ROCm<sup>™</sup> SUPPORT **THROUGH BINARIES** FROM PyTorch.ORG

PyTorch Build	Stable (1.12.0)		Preview (Nightly)	
Your OS	Linux		Mac	
Package	Conda	Pip		L
Language	Python			C
Compute Platform	CUDA 10.2	CUDA 11.3	CUDA 11.6	R
Run this Command:	pip3 instal	l torch torch	vision torcha	udi

chaudio --extra-index-url https://downlo ad.pytorch.org/whl/rocm5.1.1

LibTorch

C++/Java

ROCm 5.1.1

LTS (1.8.2)

Windows

Source

CPU



#### Disclaimer

The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions, and typographical errors. The information contained herein is subject to change and may be rendered inaccurate for many reasons, including but not limited to product and roadmap changes, component and motherboard version changes, new model and/or product releases, product differences between differing manufacturers, software changes, BIOS flashes, firmware upgrades, or the like. Any computer system has risks of security vulnerabilities that cannot be completely prevented or mitigated. AMD assumes no obligation to update or otherwise correct or revise this information. However, AMD reserves the right to revise this information and to make changes from time to the content hereof without obligation of AMD to notify any person of such revisions or changes.

THIS INFORMATION IS PROVIDED 'AS IS." AMD MAKES NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE CONTENTS HEREOF AND ASSUMES NO RESPONSIBILITY FOR ANY INACCURACIES, ERRORS, OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION. AMD SPECIFICALLY DISCLAIMS ANY IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY, OR FITNESS FOR ANY PARTICULAR PURPOSE. IN NO EVENT WILL AMD BE LIABLE TO ANY PERSON FOR ANY RELIANCE, DIRECT, INDIRECT, SPECIAL, OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF AMD IS EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

Third-party content is licensed to you directly by the third party that owns the content and is not licensed to you by AMD. ALL LINKED THIRD-PARTY CONTENT IS PROVIDED "AS IS" WITHOUT A WARRANTY OF ANY KIND. USE OF SUCH THIRD-PARTY CONTENT IS DONE AT YOUR SOLE DISCRETION AND UNDER NO CIRCUMSTANCES WILL AMD BE LIABLE TO YOU FOR ANY THIRD-PARTY CONTENT. YOU ASSUME ALL RISK AND ARE SOLELY RESPONSIBLE FOR ANY DAMAGES THAT MAY ARISE FROM YOUR USE OF THIRD-PARTY CONTENT.

ANSYS, CFX and any and all ANSYS, Inc. brand, product, service and feature names, logos and slogans are registered trademarks or trademarks of ANSYS, Inc. or its subsidiaries in the United States or other countries.

Eclipse® is registered to the Eclipse Foundation in the United Stated, other countries, or both

© 2022 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, ROCm, Radeon, Radeon Instinct and combinations thereof are trademarks of Advanced Micro Devices, Inc. in the United States and/or other jurisdictions. Other names are for informational purposes only and may be trademarks of their respective owners.

The OpenMP name and the OpenMP logo are registered trademarks of the OpenMP Architecture

# AMDA

### Questions?

alessandro.fanfarillo@amd.com

Ð

#