# Introduction to FPGA, FINN and Brevitas

**Dr. Mario Ruiz**
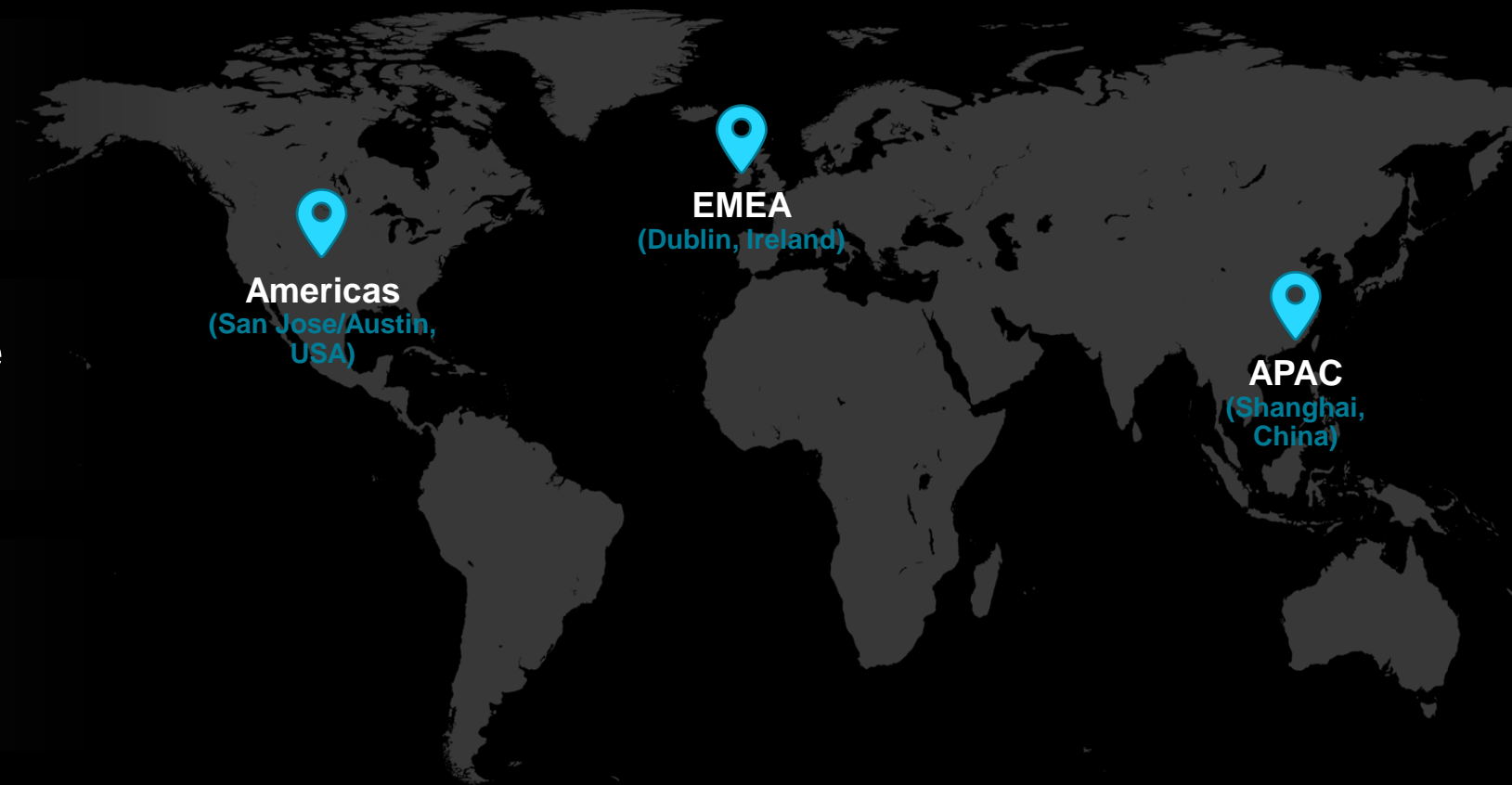**AMD University Program**

**AMD**
together we advance_

# AUP Vision

Empower academics with AMD technology to enhance teaching and learning experiences and advance state-of-the-art research.

AMD△
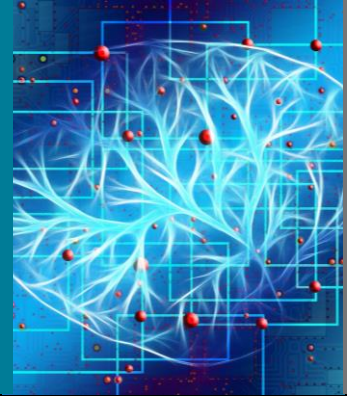together we advance_

# Our Team

Dedicated world-wide technical team

Supporting High Performance and Adaptive Compute

25+ years experience working with academia

Americas
(San Jose/Austin, USA)

EMEA
(Dublin, Ireland)

APAC
(Shanghai, China)

AMD
together we advance_

# What We Offer


**Research Programs**


**Donation Program**


**Teaching Resources**


**Training**


**Academic Solutions**


**Support**

**AMD**
together we advance_

# HACCs: Heterogeneous Accelerated Compute Clusters

Remote access to Adaptive Compute hardware
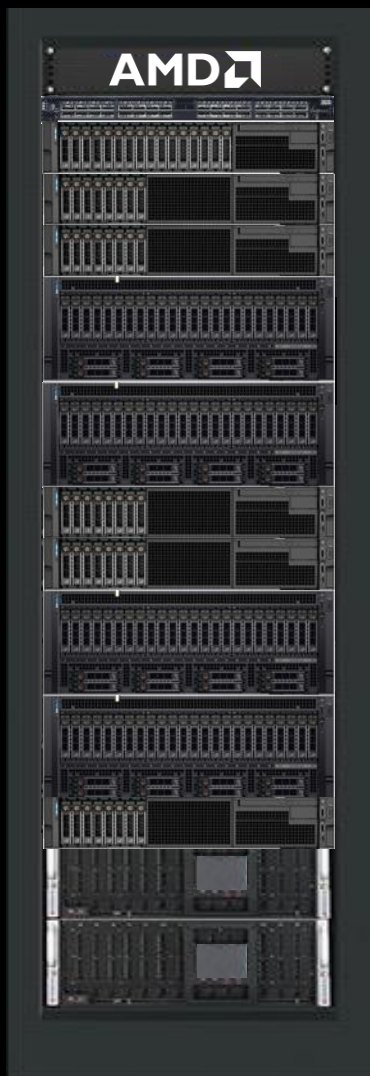
HACC user group meetings

Access to AMD researchers

Collaboration opportunities



PADERBORN UNIVERSITY

ETH zürich

ILLINOIS

UCLA

Indian Institute of Science

NUS
National University of Singapore

**AMD EPYC**

**AMD INSTINCT**

**AMD ALVEO**

**AMD VERSAL**

*www.amd-haccs.io*

★ Newest HACC at IISc, Bangalore

AMD
together we advance_

# HACC Adaptive Computing Hardware

- HACC hardware consists of:
  - Compute and Alveo™ nodes (initially U250 and U280 with HBM)
  - Latest heterogeneous nodes (SMC 4124GS) include:
    - 2 EPYC™ 3rd generation CPUs
    - 4 AMD Instinct™ MI210 GPUs
    - 2 Alveo U55C FPGA with HBM
    - 2 VCK5000 Versal Adaptive SoC with AIEs
    - Run-time via AMD ROCm™, XRT
    - SW development via HIP, Vitis, frameworks
  - 100G network

- Community hub for researchers
  - Support from in-house AMD research groups
  - Reproducible results & experiments

AMD
together we advance_

# Contact Us

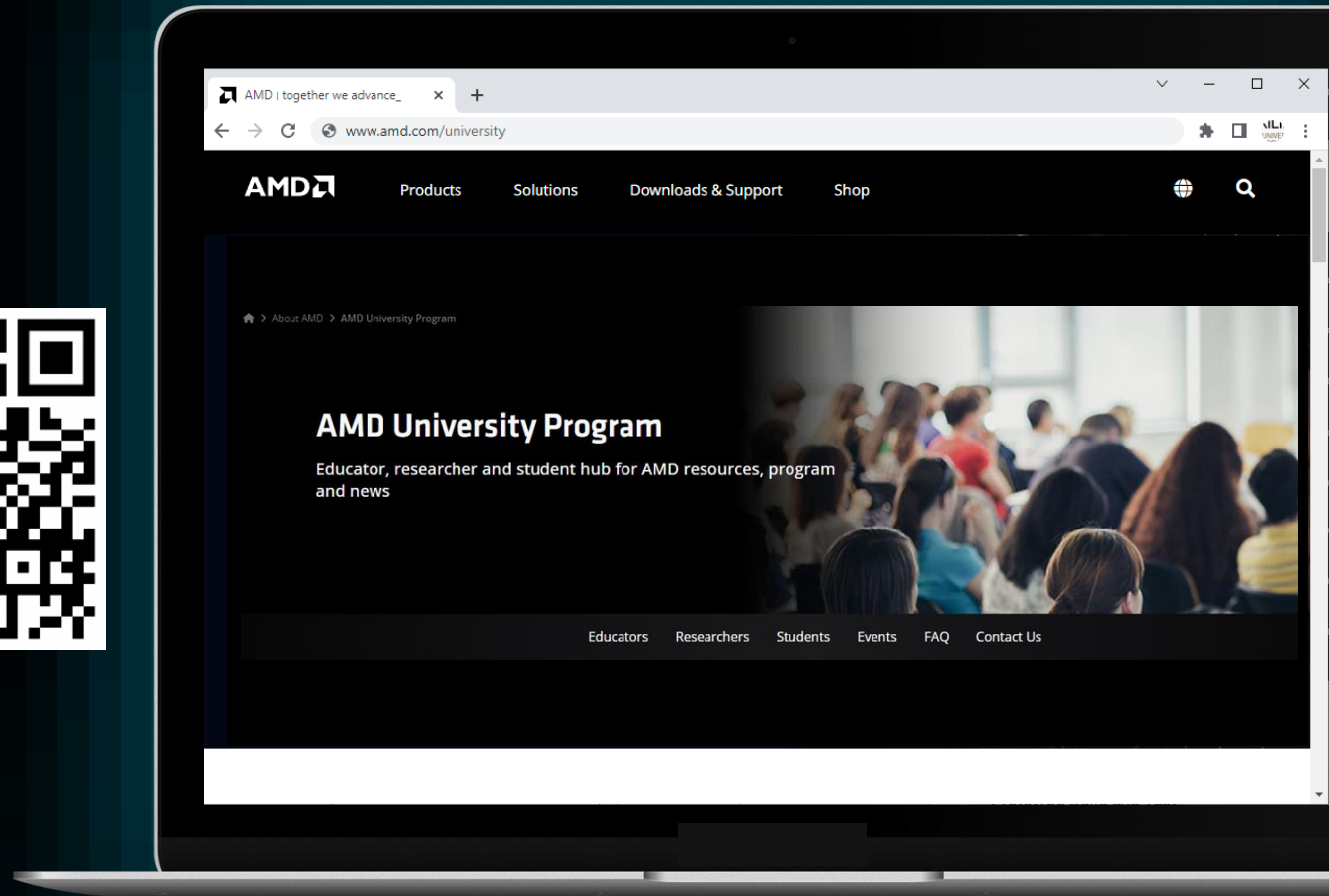## Visit our website to:

- Discover our research programs

- Access educational resources

- Submit a donation request

- Find training & other events

## Email us:
aup@amd.com

**www.amd.com/AUP**

AMD
together we advance_

# What is Adaptive Computing?

**Optimize for the Workload**
Domain-Specific Architecture for your exact requirements, accelerating the whole application
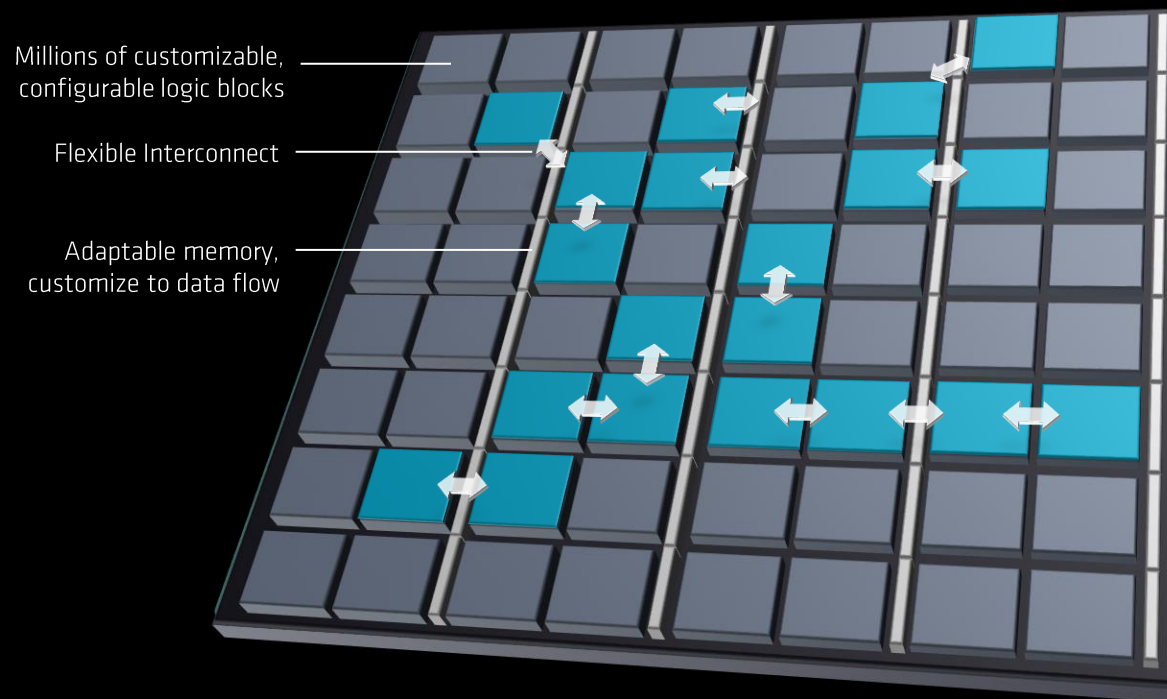
**Adapt as Algorithms Change**
Re-implement the silicon after deployment, adapting to evolving use cases
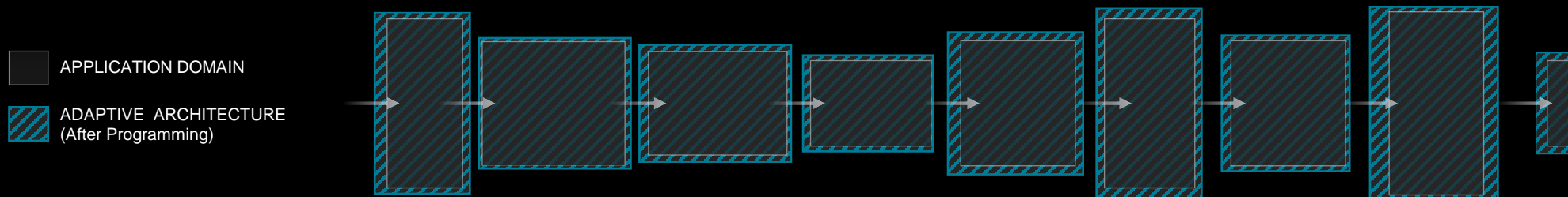
**Accelerate Pace of Innovation**
Keep pace with fast moving markets and rapid innovation cycles, e.g., AI algorithms

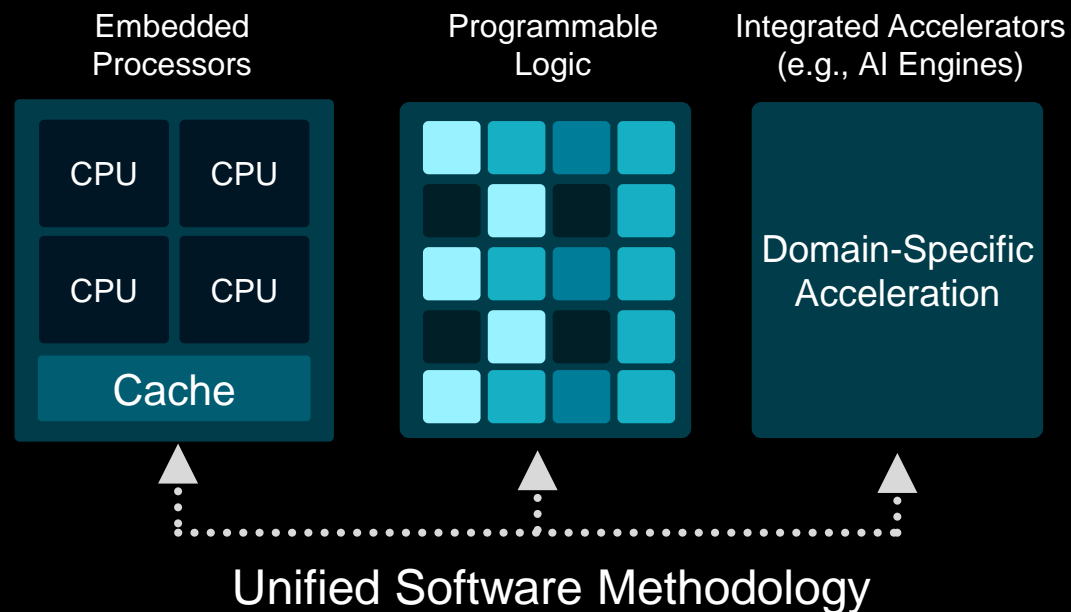Adaptive Hardware ("FPGA")
Conceptual Representation

Millions of customizable, configurable logic blocks

Flexible Interconnect

Adaptable memory, customize to data flow

Matching the Architecture to the Application
*Custom Data Flow, Custom Memory Hierarchy, Custom Precision*

APPLICATION DOMAIN

ADAPTIVE ARCHITECTURE
(After Programming)

AMD
together we advance_

# Evolution to Heterogeneous Platforms

- From FPGAs to adaptive SoCs → matching the engine to the workload

- Balancing diverse technologies for domain-specific requirements

## Domain Specific Optimization

| Embedded Processors | Programmable Logic | Integrated Accelerators (e.g., AI Engines) |
|---|---|---|
| CPU CPU / CPU CPU / Cache | | Domain-Specific Acceleration |

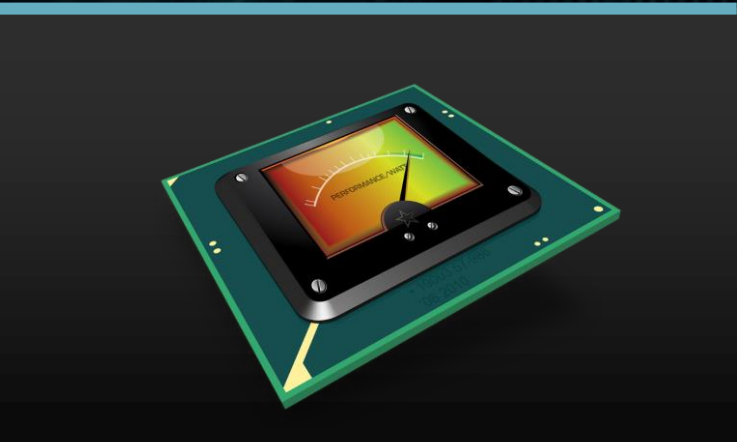Unified Software Methodology

AMD
together we advance_

# Field Programmable Gate Array (FPGA)

- Semiconductor devices
- Programmed and reprogrammed by a user
  - Configuration attributes manipulated after manufacturing
  - Matrix of configurable logic blocks (CLBs)
  - Dedicated specialized logic
  - Flexible programmable interconnects
- Ideal fit for many different workloads
  - Massive parallelism
- Hardware adaptability is a unique differentiator from CPUs and GPUs

- Invented in 1985

**Applications**
- Automotive
- Broadcast & Pro AV
- Consumer Electronics
- Data Center
- High Performance Computing and Data Storage
- Industrial
- Medical
- Video & Image Processing
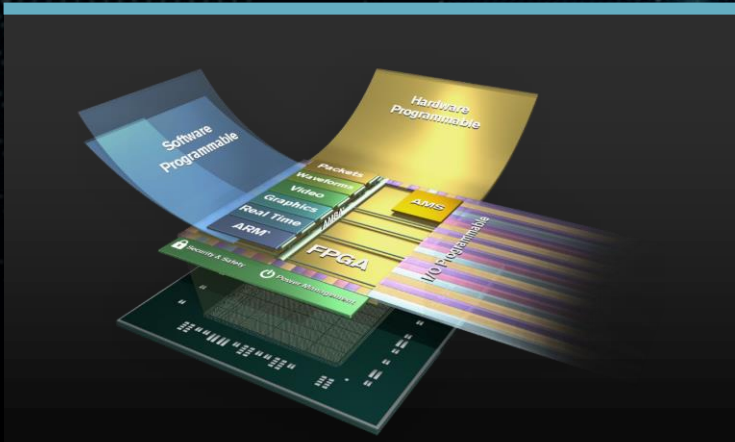- Wired Communications
- Wireless Communications

AMD
together we advance_

# Core Adaptable Hardware Technologies



## FPGAs

From high-bandwidth connectivity to massive compute engines
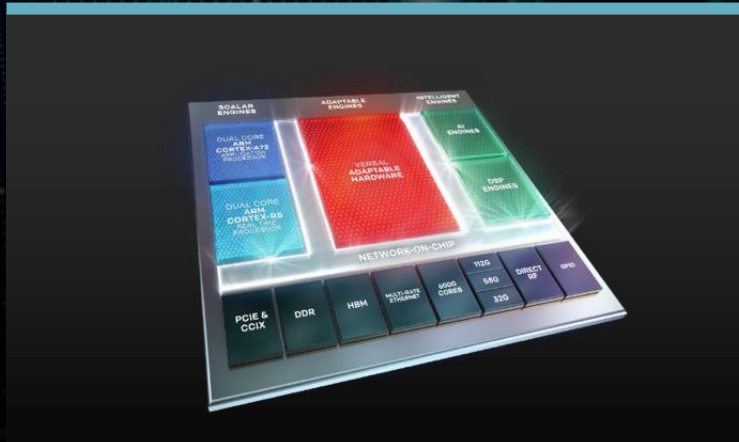
**AMD SPARTAN**   **AMD ARTIX**   **AMD KINTEX**   **AMD VIRTEX**

## SoCs

Multi-processing subsystem with Arm® cores and integrated FPGA logic

**AMD ZYNQ**

## Adaptive SoCs

Adaptive Compute Acceleration Platforms for any application, any developer

**AMD VERSAL**

**AMD** together we advance_

# Three Ages of FPGAs

- A Retrospective on the First Thirty Years of FPGA Technology
- S. M. Trimberger, "Three Ages of FPGAs: A Retrospective on the First Thirty Years of FPGA Technology," in Proceedings of the IEEE, vol. 103, no. 3, pp. 318-331, March 2015, DOI: 10.1109/JPROC.2015.2392104

# FPGA: 7-Series Architecture

- Logic elements distributed on regular columns
  - Scalability from low-cost to high-performance
- High-speed IO
- Clock management
- Interconnect matrix
  - Routing resources

Legend:
- CLB
- BRAM
- I/O
- CMT
- FIFO Logic
- BUFG
- DSP
- BUFIO & BUFR
- MGT

Artix-7 Architecture Overview

AMD
together we advance_

# Configurable Logic Block (CLB)

- Primary resource for design in AMD FPGAs
  - Combinatorial functions
  - Flip-flops
- CLB contains two slices
- Connected to switch matrix for routing to other FPGA resources
  - Carry chain runs vertically

AMD
together we advance_

# Two Types of CLB Slices

- SLICEM: Full slice
  - Can be used for logic, memory and shift register LUT
  - Has wide multiplexers and carry chain
- SLICEL: Logic and arithmetic only
  - LUT can only be used for logic (not memory)
  - Has wide multiplexers and carry chain

# Slice Resource

- Four six-input Look-Up Tables (LUT)

- Multiplexers

- Carry chains

- Four flip-flops/latches
  - Four additional flip-flops

- The implementation tool will pack multiple slices in the same CLB if certain rules are followed

# 6-Input LUT with Dual Output

- LUTs can be two 5-input LUTs with common input
  - Minimal speed impact to a 6-input LUT
  - One or two outputs
- Any combinatorial function of six variables or two functions of five variables

AMD

together we advance_

# Slice Flip-Flops and Flip-Flop/Latches

- Each slice has four flip-flop/latches (FF/L)
  - Can be configured as either flip-flops or latches
- Each slice also has four flip-flops (FF)

# Slice Flip-Flop Capabilities

- All flip-flops are D type
  - Q output
- All flip-flops have a single clock input (CK)
- All flip-flops have an active high chip enable (CE)
- All flip-flops have an active high SR input
  - Input can be synchronous or asynchronous
  - Sets the flip-flop value to a pre-determined

FDRE

D

Q

CE

C

R

UG474_c3_05_102910

AMD
together we advance_

# 7-Series FPGA I/O

- Wide range of voltages
  - 1.2V to 3.3V operation
- Wide I/O standards support
  - Single ended and differential
  - Referenced voltage inputs
  - 3-state capability
- Very high performance
  - Up to 1600 Mbps LVDS
  - Up to 1866 Mbps single-ended for DDR3
- Easy memory interfacing
  - Hardware support for QDRII+ and DDR3
- Digitally controlled impedance
- Power reduction features

# 7-Series Block RAM and FIFO

- Fully synchronous operation
  - Outputs are latched
- Optional internal pipeline register
  - Higher frequency operation
- Two independent ports access common data
  - Individual address, clock, write enable, clock enable
  - Independent data widths for each port

AMD
together we advance_

# 7-Series Block RAM and FIFO

- Multiple configuration options
  - True dual-port, simple dual-port, single-port
- Integrated cascade logic
- Byte-write enable in wider configurations
- Integrated control for fast and efficient FIFOs
- Integrated 64/72-bit Hamming error correction

Each block RAM block can be used as

36 Kb BRAM / FIFO

or

18 Kb BRAM

18 Kb BRAM / FIFO

(1) 36 Kb BRAM
OR
(1) 36 Kb or FIFO

(2) independent 18 Kb block RAMs
OR
(1) 18 Kb FIFO + (1) 18 Kb block RAM

AMD
together we advance_

# 7-Series DSP48E1 Slice



- 25x18 signed multiplier
- 48-bit add/subtract/accumulate
- 48-bit logic operations
- Pipeline registers for high speed
- Pattern detector
- SIMD operations (12/24 bit)
- Cascade paths for wide functions
- Pre-adder

# 7-Series FPGAs Clock Management

- Global clock buffers
  - High fanout clock distribution buffer
- Low-skew clock distribution
  - Regional clock routing
- Clock regions
  - Each clock region is 50 CLBs high and spans half the device
- Clock management tile (CMT)
  - One Mixed-Mode Clock Managers (MMCMs) and one Phase Locked Loop (PLL) in each Clock
  - Performs frequency synthesis, clock de-skew, and jitter-filtering
  - High input frequency range

**AMD**
together we advance_

# Programming Model

**Hardware Description Languages (HDL)**

- Verilog

- VHDL

- System Verilog

- Closer to the metal
  - Low level abstraction
  - Describe the behaviour

**High-Level Synthesis (HLS)**

- C/C++

- High level of abstraction
  - Write algorithms

- Vitis HLS generates the architecture
  - Guided by user directives

**AMD Vivado**

**AMD Vitis**

**AMD**
together we advance_

# VHLD/Verilog counter

## VHDL

```vhdl
library IEEE;
use IEEE.STD_LOGIC_1164.ALL;
use IEEE.STD_LOGIC_UNSIGNED.ALL;
entity counter is
    Port ( clk: in std_logic;
           rst: in std_logic;
           cout: out std_logic_vector(3 downto 0)
    );
end counter;
architecture rtl of counter is
signal counter_up: std_logic_vector(3 downto 0);
begin
    process(clk)
    begin
    if(rising_edge(clk)) then
        if(rst='1') then
            counter_up <= x"0";
        else
            counter_up <= counter_up + x"1";
        end if;
    end if;
    cout <= counter_up;
    end process;
end rtl;
```

## Verilog

```verilog
module counter(
    input clk,
    input rst,
    output reg [7:0] count
    );

  always @(posedge(clk)) begin
    if (rst)
        count <= 0;
      else
        count <= count + 1;
    end
endmodule
```

AMD
together we advance_

# Vitis HLS Vector addition

```cpp
void vadd(const int* in1, // Read-Only Vector 1
          const int* in2, // Read-Only Vector 2
          int* out,       // Output Result
          int elements    // Number of elements
          ) {
// Simple vector addition kernel.
vadd1:
    for (int i = 0; i < elements; i++) {
        out[i] = in1[i] + in2[i];
    }
}
```

AMD
together we advance_

# What is AMD Vitis™ HLS and HLS Benefits

Structured C/C++

RTL Code

Automated C/C++ to RTL Conversion

Allows Significantly Faster Design Iterations

Significantly Accelerates Simulation – Important For Wireless, Video Applications

AMD
together we advance_

# AI on FPGA

# DNNs and their Potential

**1.** Requires little domain expertise

**2.** NNs are a "universal approximation function"

**3.** If you make it big enough and train it long enough
- Can outperform humans and existing algorithms on specific tasks

---

**Will not only increasingly replace other algorithms, but also…**



Nature, Oct 2021

**… solve previously unsolved problems**

- ChatGPT, Copilot
- Stable diffusion
- Protein folding



*Stable Diffusion Prompt: "Pencil sketch of an international group of semiconductor research scientists, studio Ghibli"*

**AMD**
together we advance_

# Spectrum of ML use case with very different requirements

Sensor Intelligence
**10k LUTs**
**10s infps**
**Latency sensitive**

Cyber security
Communications
**100s M infps**
**microsec latency**

**Orders of Magnitude**

Recommender, Computer Vision or
Natural Language Processing related DNNs
**10s k infps**

High energy particle physics
**nsec latency**

**infps = inferences per second**

**AMD**
together we advance_

# DNN Compute Requirements are Outpacing Moore's Law



AlexNet to AlphaGo Zero: A 300,000x Increase in Compute

Log Scale    Linear Scale

**Doubling every 3.4months Since 2012**

**Moore's Law Doubling every 18months**

Compute Requirements Petaflops/day (log scale)

Time

*Source: https://blog.openai.com/ai-and-compute*

# Innovation is needed to provide the necessary performance scalability

# Specialization Is #1 Industry Approach to Achieve Performance Scalability and Energy Efficiency

# Adaptive Computing or Dedicated Silicon for DPUs



- With increasing specialization of the device, potential sales volume decreases
  - Hard to amortize the increasing NRE costs involved in building ASSPs
  - FPGAs become more attractive
- Increasing specialization scales performance for both ASSPs and FPGAs

- **The opportunity for FPGAs lies in their ability to specialize**

35 *ASSP: Application Specific Standard Product*

# Vitis AI - ML in general

# Customization levels on Adaptive Computing

**MPE**



Custom Dataflow | Quantization | Sparsity

Customized for ML in general

Customized for Specific Topologies

Customized in Datatypes

Customized in Connectivity

Specialization/Performance/Efficiency

**AMD**
together we advance_

# Popular Approach: Matrix of Processing Engines (MPEs) Specializing for AI in general

- Popular layer-by-layer compute
- Batching to achieve high compute efficiency
  - At latency cost (latency ~ batch size)
- Specialized processing engines
  - Operators
  - ALU types
    - tensor-, matrix- or vector-based

- **Customized for ML in general**
  - **Designed to run any DNN**
- **Works really well for computer vision and natural language processing (10s kinfps)**
- **Popular approach: Vitis AI (FPGA or AIE) as well as majority of AI accelerators**

**DNN**

Matrix of Processing Engines (MPE)
"layer-by-layer" compute

PE

Buffer

AMD
together we advance_

# AMD Vitis™ AI Integrated Development Environment

*A Complete AI Stack for Adaptable AMD Targets*

TensorFlow

PyTorch

**Vitis™ AI Tools & Components**

| Model Zoo | Community or User Models | Optional 3rd Party Framework Enablement |
|---|---|---|
| Optimizer | | ONNX RUNTIME    TensorFlow Lite |
| Quantizer | | |
| Libraries  Profiler  Compiler | | tvm |
| Runtime | | |

**Domain-Specific Architectures**

| Embedded Deep Learning Processing Units | Data Center Deep Learning Processing Units |
|---|---|

**Supported AMD Targets**

VCK190    ZCU102    ZCU104    Kria™ K26 SOM    U50/C/LV    U200 / U250    U280    VCK5000

Your Platform

39

AMD together we advance_

[Public]

# AI Model Zoo – Expanding to Diverse AI Applications

- A comprehensive AI model repository
  - Open and free to download for any user
  - State-of-the-art models from Pytorch, TF & TF2
  - Retrainable, appliable to various data set & scenario
  - Deployable on AMD FPGA and Versal Adaptive SoC
- New models in each release



40

# Extensive Application Coverage

## Classification

- Inception
- Mobilenet
- Resnet
- VGG
- EfficientNet
- MLPerf ResNet50
- OFA ResNet
- Vision Transformer
- Car Type classification
- Car Color classification

## IndustrialVision/Robotics

- FADNet
- PSMNet
- PMG
- Superpoint
- HFNet

## Detection

- ssd_mobilenet
- Yolov3
- Yolov4
- YoloX
- Refinedet
- EfficientDet
- Pointpillars
- Centerpoint
- CLOCs
- Pointpainting
- Multi-taskv3
- OFA-Yolo

## Medical Image

- RCAN
- SESR
- OFA-RCAN
- DRUnet
- SSR
- C2D2lite

## Segmentation

- ENet
- Semantic FPN
- Salsanext
- Salsanextv2
- SOLO
- Mobilenetv2
- 2D-Unet
- FPN-ResNet18
- Unet-Chaos-CT
- HardNet
- Sa-Gate

## NLP

- Bert-base
- Sentiment detection
- Customer satisfaction
- Open-information-extraction

## Video Analytics

- Face Recognition
- Face Quality
- Face ReID
- Person ReID
- FairMOT
- FaceMask Detection
- MoveNet

## Text-OCR

- Textmountain, OCR

41

**AMD**
together we advance_

# Compiling for DPU - an XIR-based Toolchain

- Xilinx Intermediate Representation (XIR)
  - Graph-based intermediate representation of the AI algorithms
  - Designed for compilation and efficient deployment of the DPU on the FPGA platform.
- XIR-based compilation flow
  - First, transform the input models to XIR format
  - Breaks up computing graph to subgraphs
  - Execute DPU subgraph to a compiled xmodel file

```
┌──────────────┐        ┌──────────┐              ┌──────────┐             ┌──────────┐
│ Deep learning│───────▶│  Parser  │─────────────▶│ Compiler │────────────▶│ Runtime  │
│  frameworks  │        │          │              │          │             │          │
└──────────────┘        └──────────┘              └──────────┘             └──────────┘
                                      .xmodel                    .xmodel
                                    (XIR-based)                (XIR-based)
```

AMD
together we advance_

# Techniques for Further Specialization with
## Adaptive Compute Architectures

# Specialization beyond MPEs

**MPE**

Custom Dataflow

Quantization

Sparsity

Buffer

Customized for
ML in general

Customized for
Specific Topologies

Customized in
Datatypes

Customized in
Connectivity

**AMD**
together we advance_

# Dataflow - Specializing for Individual Topologies

- Hardware instantiates the topology as a dataflow architecture
  - Customize everything to the **specifics of the given DNN**, any operation, any connectivity
- Benefits:
  - Improved efficiency
  - Low fixed latency
- Scale performance & resources to meet the application requirements
  - If resources allow, we can completely unfold to create a circuit that inferences at clock speed and thereby meet these new throughput requirements



**DNN**

allocated resource ~ compute requirement per layer

200Minfps

**FPGA**

**Dataflow can scale performance to meet the application requirements**

AMD
together we advance_

# Specialization beyond MPEs



Customized for ML in general

Custom Dataflow

Customized for Specific Topologies

Quantization

Customized in Datatypes

Sparsity

Customized in Connectivity

**AMD**
together we advance_

# Customizing Arithmetic to Minimum Precision



C= f(size of accumulator, size of weight,size of activation)

- Popular approach which reduces bits in the data representation of weights and activations while preserving accuracy

- Reducing precision shrinks hardware cost/ scales performance
  - Instantiate n-times more compute within the same fabric, thereby scale performance n-times

- Reduces memory footprint
  - NN model can stay on-chip => no memory bottlenecks

- With dataflow: every layer has dedicated compute resources, we can mix and match precision across layers
  - Exploit custom arithmetic at a greater degree than MPEs

| Precision | Model size [MB] (ResNet50) |
|-----------|----------------------------|
| 1b        | 3.2                        |
| 8b        | 25.5                       |
| 32b       | 102.5                      |

**Reducing precision saves resources/ scales performance, and reduces memory**
**However, it requires quantization support in the training software**

5b  1b  2b  3b

AMD
together we advance_

# Specialization beyond MPEs

MPE

Buffer

Custom Dataflow

Quantization

Sparsity

Customized for
ML in general

Customized for
Specific Topologies

Customized in
Datatypes

Customized in
Connectivity

AMD
together we advance_

# Sparsity

- DNNs are naturally sparse

- Sparse topologies result in irregular compute patterns which are difficult to accelerate on vector- or matrix-based execution units

- With streaming dataflow architectures, where every neuron and synapse is represented in the hardware, we can fully exploit this

**FPGA**

**Optimized Dataflow on FPGA**

**AMD**
together we advance_

# Taking it to the Extreme: **LogicNets**

# Specialization beyond MPEs



Custom Dataflow

Quantization

Sparsity

Full co-design

Customized for Specific Topologies

Customized in datatypes

Customized in connectivity

**Specialize the DNN to suit hardware**

AMD
together we advance_

# LogicNets with Adaptive Computing



thresholding activation

6x1-bit inputs → Σ → batch norm → ⊓ → 1x1-bit output

float weights

convert (enumerate inputs)

6:1 LUT

**Traditional**

Pick DNN topology → Train → Optimize hardware → Deploy

Adjust the parameters of DNN while iterating on training dataset until accuracy

**LogicNets**

Design a circuit (=unrolled DNN) → Train → Deploy

Typically sparse topologies

Adjust the parameters of DNN (=LUT contents) while iterating on training dataset until accuracy

Maximum performance by design (classification at clock rate) [5]

Compared to unrolled DF: sparse to suit the interconnect

AMD together we advance_

# Unique Opportunity for Adaptive Computing

- FPGAs can scale DNN performance through extreme specialization

- Reduced precision arithmetic
  - Arbitrary bitwidth
  - Mix & match bitwidths between layers

- Fine-grained sparsity

- Scalable, layer-parallel streaming dataflow

*Reducing precision*
*f: < 8bit operations*

*Mix & Match precisions*

*Reducing connectivity to suit the interconnect on our devices*

*Map each layer to parallel hardware*

**High degree of specialization doesn't make sense for ASSPs**

FPGA

**AMD**
together we advance_

# How much do we get out of the different specializations?

# Deep Network Intrusion Detection System (NIDS)

Traditional: Hand-coded rules
**Emerging trend: Neural networks**

Network interface L1-L3 → Packet processing/Feature extraction → **Traffic Classification** Identify malware such as DOS, fuzzers, worms, etc → Packet filter drop/pass → Network interface L1-L3

**Network processing system**

FPGA

**Goal:** Implement **NN-based traffic classifier** delivering 100G **line-rate** throughput = 150 Mips
Latency sensitive (buffer 10s of MB/msec)

Dataset: UNSW-NB15 Moustafa, Nour, and Jill Slay. "UNSW-NB15: a comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set)." 2015 military communications and information systems conference (MilCIS). IEEE, 2015.

AMD
together we advance_

# Results – Implementations



Specialization

| Generic ML architecture | | Specialized for topology<br>2 different folding factors | | Fully<br>Co-designed | |
|---|---|---|---|---|---|
| | **MPE (VitisAI)** | | **FINN (fold 8, fold 1)** | | **LogicNet** |
| **Topology / #layers / #OPs** | MLP / 3 / 92KOPs | | MLP / 3 / 92KOPs | | **Circuit / 4 / 15.4KOPs** |
| **#inputs / neuron** | 64 | | 64 | | **7** |
| **#bits / weight & activation** | 8b & 8b | | 2b & 2b | | 2b & 32b |
| **Accuracy** | 92.3% | | 91.9% | | 91.3% |

Same network

Reduced precision

Sparsity

New sparse topology with 7 Inputs/LUT 6x less OPs

4 implementations with varying degree of specialization

**AMD**
together we advance_

# Results – Throughput and Latency

Performance

## Throughput & Latency

471Mips @ 9ns

300Mips @ 18ns

25Mips @ 240ns

22Kips @ 26us

| | | | |
|---|---|---|---|
| 100000000 | | | |
| 1000000 | | | |
| 10000 | | | |
| 100 | | | |
| 1 | | | |

Fold 8

Fold 1

**Vitis AI MPE**

**FINN DF**

**LogicNet**

Performance scaling
And latency reduction
through specialization
~1000x

Further unrolling
~8x

DF unrolling reduces
latency further

Full co-design ~ 2x
sparse connectivity suits the
interconnect

Sparse topology reduces
#pipeline stages => latency

Classification ~ clock rate

**Specialization scales performance and reduces latency
by orders of magnitude if application is amenable**

AMD
together we advance_

# Resource Cost - Compute, Memory

**Efficiency**

## Resources

■ kLUTs   ■ DSPs   ■ BRAM   ■ URAM

Reducing precision reduces LUT & DSP cost, even when further parallelized (1000x speedup)

Further unrolling costs proportional to folding factor, however synthesis prunes!
Memory merged into logic

LogicNet low compute cost due to limited connectivity

|  | VitisAI | FINN (fold8) | FINN (fold1) | LogicNets |
|---|---|---|---|---|
| ■ kLUTs | 122 | 44 | 8-69 | 16 |
| ■ DSPs | 1124 | 0 | 0 | 0 |
| ■ BRAM | 290 | 166 | 0 | 0 |
| ■ URAM |  |  |  |  |

**Customizing arithmetic, sparse implementations and learned circuits greatly reduce resources and improve device efficiency**

AMD
together we advance_

# Deep Network Intrusion Detection System (NIDS) Results

- This example illustrates the trade-offs between specialization and performance and efficiency

- Custom arithmetic is effective to **scale performance** and **dataflow to reduce latency**
  - If application is amenable, custom arithmetic can meet extreme throughput requirements such as in NIDS

- Reduced precision, fine-granular sparsity & learned circuits can **shrink the resource** requirements despite speedup

- These are some of the opportunities which make most sense to exploit with FPGAs

**AMD**
together we advance_

# General Introduction to FINN

**AMD**

# Project Mission and Key Techniques

# FINN – Project Mission

- Custom Specialization
  - for creating **high-throughput, ultra-low-latency** DNN inference engines

- End-to-End
  - flow for the *easy* creation of **specialized hardware architectures** for FPGAs

- Open Source
  - for full **transparency and flexibility** to adapt to end user applications and
  - for easy customer interactions

# Two Key Techniques for Customization in FINN

**Streaming Dataflow Architectures**
for FPGAs

**Custom Precision:**
**Few-bit Weights and Activations**



e.g. 1-bit weights

e.g. 3-bit activations

AMD
together we advance_

# Customized Dataflow Processing versus More Generic Architectures

**Matrix of Processing Engines (MPE)
(Vitis AI, TPUs, GPUs)**

MAC, Vector
Processor or VLIW

**Dataflow Architectures
with FPGAs and FINN**

**Customized
Data path**

**AMD**
together we advance_

# Matrix of Processing Engines (MPEs) Specializing for AI in General

- Popular layer-by-layer compute
- Batching to achieve high compute efficiency
  - At latency cost (latency ~ batch size)
- Customized for ML in general
  - Designed to run any DNN
  - Specialized processing engines
    - Operators
    - ALU types
- Works really well for computer vision and natural language processing
- Popular approach: Vitis AI (FPGA or AIE) as well as majority of AI accelerators



**DNN**

Matrix of Processing Engines (MPE) "layer-by-layer" compute

PE

Buffer

AMD together we advance_

# Dataflow - Specializing for Individual Topologies

- Hardware instantiates the topology as a dataflow architecture
- Customize everything to the specifics of the given DNN, any operation, any connectivity
- Benefits
  - Improved efficiency
  - Low fixed latency
- Scale performance and resources to meet the application requirements

**DNN**

allocated resource ~ compute requirement per layer

**FPGA**

**Dataflow can scale performance to meet the application requirements**

AMD
together we advance_

# Dataflow Processing:
## *Scaling to Meet Performance and Resource Requirements*



FPGA (fold 1)

200MRps

Scaling to maximize throughput

FPGA (fold 10)

20MRps

Scaling to fit into available resources

FPGA (fold 1000)

200kRps

**1. Scale performance and resources to meet the application requirements**
**2. If resources allow, we unfold completely, creating a circuit for inference at clock speed**

AMD
together we advance_

# Customized Dataflow Processing versus More Generic Architectures

### Matrix of Processing Engines (MPE) (Vitis AI, TPUs, GPUs)

MAC, Vector Processor or VLIW



### Dataflow Architectures with FPGAs and FINN

**Customized Data path**



- Customized for typical DNN operations
  - e.g., multiply accumulate
- Lower throughput (~10KRps)
- Flexibility through programming
- Applications: CV, Speech

- Customized/adapted for specific DNN topologies
- Streaming interfaces
- Specialization -> higher efficiency
- Lower latency (no intermediate buffering)
- Higher throughput (~100MRps)
- Flexibility through reconfiguration
- Applications: radio, networking, material science, particle physics – smaller DNNs

AMD
together we advance_

# Quantization

- Reducing precision shrinks hardware cost/scales performance
  - For integer datatypes, LUT cost proportional to both bitwidths in weight and activations (e.g., INT8 : INT1 ≈ 70×)
  - n-times more compute fits into the same fabric, thereby, scaling performance n-times or shrinking hardware cost accordingly

- Energy
  - Faster execution or smaller footprint → less energy ($E = P \cdot time$)
  - Using reduced precision operators saves energy
  - Reduces memory footprint
    - ResNet50 @ 32b: 102.5 MB, ResNet50 @ 2b: 6.4 MB
    - NN model can stay on-chip → no external memory access → saves energy

| Precision | Model size [MB] (ResNet50) |
|-----------|---------------------------|
| 1b | 3.2 |
| 8b | 25.5 |
| 32b | 102.5 |



| Operation | | Picojoules per Operation | | |
|---|---|---|---|---|
| | | 45 nm | 7 | 45 / 7 |
| + | Int 8 | 0.03 | 0.007 | 4.3 |
| | Int 32 | 0.1 | 0.03 | 3.3 |
| | BFloat 16 | -- | 0.11 | -- |
| | IEEE FP 16 | 0.4 | 0.16 | 2.5 |
| | IEEE FP 32 | 0.9 | 0.38 | 2.4 |
| × | Int 8 | | 0.07 | 2.9 |
| | Int 32 | | 1.48 | 2.1 |
| | BFloat 16 | -- | 0.21 | -- |
| | IEEE FP 16 | 1.1 | 0.34 | 3.2 |
| | IEEE FP 32 | 3.7 | 1.31 | 2.8 |
| SRAM | 8 KB SRAM | 10 | 7.5 | 1.3 |
| | 32 KB SRAM | 20 | 8.5 | 2.4 |
| | 1 MB SRAM[1] | 100 | 14 | 7.1 |
| GeoMean[1] | | | | 2.6 |
| | | Circa 45 nm | Circa 7 nm | |
| DRAM | DDR3/4 | 1300[2] | 1300[2] | 1.0 |
| | HBM2 | -- | 250-450[2] | -- |
| | GDDR6 | -- | 350-480[2] | -- |

is pJ per 64-bit access.

AMD
together we advance_

# The FINN Framework

# FINN Framework: From DNN to FPGA Deployment



**Brevitas**
Training in PyTorch
Algorithmic optimizations

- Train or even learn reduced precision DNNs
- Library of standard layers
- Pretrained examples

**FINN Compiler**
Hardware Architecture
Build

- Perform optimizations
- Assemble parameterized HLS/RTL modules
- Generate a DNN hardware IP

**Deployment**

- Embed the DNN IP into an infrastructure design
- Generate a Python run-time
- Enable integration with your application
- System integration available for some embedded and Alveo platforms, including HACC

AMD
together we advance_

# Brevitas:
## *A PyTorch Library for Quantization-Aware Training*



add quantization
resize layers
change hyperparameters
retrain

FP32

INT

**Precision**
Preset or learned

**Scaling Factors**
Granularities, strategies and constraints

**Target Tensors**
Weights, activations, accumulators

**Export to ONNX**
To import into the FINN compiler

https://github.com/Xilinx/brevitas

AMD
together we advance_

# FINN Compiler
## *Transform DNN into Custom Dataflow Architecture*

QONNX representation of the quantized DNN

FINN

- Uses an ONNX-based network description as intermediate representation (IR)
- Is a Python library of graph transformations
- Generates a synthesizable description of each layer (HLS/RTL) encapsulated as an IP block
- Produces a synthesized stitched IP block representing the complete network

Stitched DNN accelerator IP

AMD
together we advance_

# FINN Compiler - Network preparation

NN exported from Brevitas

**FINN Compiler**

1. Import, streamlining, and other transformations
2. Adjust folding to suit performance/resource requirements
3. Generate IP cores and stitched IP design

Stitched IP Design

## QONNX

- Directly exported from Brevitas
- Input format to FINN compiler
- Quantization operator(s)
  - Quant, BipolarQuant, Trunc
- No tensor annotations



## FINN-ONNX

- Previously used as input format
- IR in the FINN compiler
- MultiThreshold to represent activation quantization
- Custom datatype annotations on tensor

**AMD**
together we advance_

# FINN Passes - ONNX Graph Transformations



ONNX-IR ➡ **Streamlining**

ONNX-IR ➡ **...**

**Streamline()**

*simplified*

**LowerConvsToMatMul()**

*Convolutions as Im2Col + MatMul*

AMD
together we advance_

# FINN Passes - ONNX Graph Transformations



**ConvertToHWLayers and SpecializeLayers**

```
template<
    unsigned  MW,    // Width of the input matrix
    unsigned  MH,    // Height of the input matrix
    unsigned  SIMD,  // Number of input columns computed in parallel
    unsigned  PE,    // Number of output rows computed in parallel
    typename  TI,    // Input Datatype
    typename  TO,    // Output Datatype
    typename  TW,    // Weight Datatype
    typename  TA,    // Activation Datatype
>
void Matrix_Vector_Activate_Batch(
    hls::stream<hls::vector<TI>> &in,
    hls::stream<hls::vector<TO>> &out,
    TW  const &weights,
    TA  const &activation
);
```

Corresponding to
finn-hlslib function call
or finn-rtllib module

**Optimization, lowering, code generation... are all transformations**

AMD
together we advance_

# FINN Hardware Folding

NN exported from Brevitas

**FINN Compiler**

1. Import, streamlining, and other transformations
2. Adjust folding to suit performance/resource requirements
3. Generate IP cores and stitched IP design

Stitched IP Design

$$\begin{pmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} w_{00} & w_{01} & w_{02} & w_{03} \\ w_{10} & w_{11} & w_{12} & w_{13} \\ w_{20} & w_{21} & w_{22} & w_{23} \\ w_{30} & w_{31} & w_{32} & w_{33} \end{pmatrix} \cdot \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

**SIMD**

PE0: $w_{00} \cdot x_0 + w_{01} \cdot x_1$  $+w_{02} \cdot x_2 + w_{03} \cdot x_3$  $w_{20} \cdot x_0 + w_{21} \cdot x_1$  $+w_{22} \cdot x_2 + w_{23} \cdot x_3$

PE1: $w_{10} \cdot x_0 + w_{11} \cdot x_1$  $+w_{12} \cdot x_2 + w_{13} \cdot x_3$  $w_{30} \cdot x_0 + w_{31} \cdot x_1$  $+w_{32} \cdot x_2 + w_{33} \cdot x_3$

$\begin{pmatrix} y_0 \\ y_1 \end{pmatrix}$  $\begin{pmatrix} y_2 \\ y_3 \end{pmatrix}$

Clock:    0        1        2        3

**AMD**
together we advance_

# FINN HLS/RTL Library - Parameterizable Kernel Library

- Kernels representing individual layers, a.k.a. Operators
- Flexible parametrization as for
  - Degree of parallelism (output channels, input channels, kernel dimensions …)
  - Datatypes (INT8, ternary, INT2, …)
  - Behaviour (activation function)
- Instantiated and stitched by FINN compiler with AXI-Stream data path
- Implemented as synthesizable C++ (Vitis HLS) or SystemVerilog

Layer $(i - 1)$

AXI-Stream

```
template<
    unsigned  MW,   // Width of the input matrix
    unsigned  MH,   // Height of the input matrix
    unsigned  SIMD, // Number of input columns computed in parallel
    unsigned  PE,   // Number of output rows computed in parallel
    typename  TI,   // Input Datatype
    typename  TO,   // Output Datatype
    typename  TW,   // Weight Datatype
    typename  TA,   // Activation Datatype
>
void Matrix_Vector_Activate_Batch(
    hls::stream<hls::vector<TI>> &in,
    hls::stream<hls::vector<TO>> &out,
    TW  const &weights,
    TA  const &activation
);
```

AXI-Stream

Layer $(i + 1)$

AMD
together we advance_

78

# FINN Compiler: IP Generation Flow

*NN exported from Brevitas*

**FINN Compiler**

1. Import, streamlining, and other transformations
2. Adjust folding to suit performance/resource requirements
3. Generate IP cores and stitched IP design

*Stitched IP Design*

› Stream-in, stream-out FPGA IP block

» Easy "bump-in-the-wire" integration into streaming systems

» Simple data movement, fully deterministic

**AMD**
together we advance_

# Deployment with PYNQ for Python Productivity

```
# instantiate the accelerator
accel = models.cnv_w2a2_cifar10()
# generate an empty numpy array to use as input
dummy_in = np.empty(accel.ishape_normal, dtype=np.uint8)
# perform inference and get output
dummy_out = accel.execute(dummy_in)
```

- Use PYNQ-provided Python abstractions and drivers
- User provides NumPy array input, calls driver, retrieves NumPy array output
  - Internally use PYNQ DMA driver to wr/rd NumPy arrays into I/O streams

https://github.com/Xilinx/PYNQ
https://github.com/Xilinx/finn-examples

AMD
together we advance_

# FINN Dataflow Build Mode

FINN flow = Python script making calls to FINN API

Produce output files from input ONNX and config

**DataflowBuildConfig**

input ONNX → **FINN** → output files

Consists of a sequence of steps, each step

- is a Python function with a standardized interface
- consumes and produces ONNX
- may produce other files
- may be standard or custom
- may have config-dependent behavior

ONNX → step → ONNX → step → ONNX

step → output files

step → output files

Can be resumed from intermediate steps

ONNX files act as checkpoints

**AMD ⌐**
together we advance_

# FINN Infrastructure and Workflow

# The FINN Ecosystem and Software Stack

| finn-examples |
|---|

| finn |
|---|

| qonnx | finn-hlslib | brevitas |
|---|---|---|
| ONNX    ONNX Runtime | Vitis HLS | PyTorch |
| **Core infrastructure** | **Operator library** | **Frontend** |

☑ *FINN* project landing page: https://xilinx.github.io/finn

- Quick Start, Documentation, Examples (Jupyter Notebooks)
- Links to Repos

**AMD**
together we advance_

# A FINN End-to-End Flow

**Brevitas**

Trained Network in PyTorch/Brevitas

Brevitas FINN-ONNX Export

**Network Preparation**

Network of high-level ONNX layers

**Simulation and Emulation Flows**

Streamlining Transformations

Simulation using Python

Streamlined network of high-level ONNX layers

Convert to HW Layers

Network of HW layers, maximum folding

Prepare cppsim

Specialize Layers

Network of HLS/RTL layers, maximum folding

Network of HLS layers With C++ wrappers

Adjust folding to maximize performance

Run cppsim (HLS C++)

Network of HLS/RTL layers, desired folding

**Vivado HLS and IPI**

Create IP per layer

Network of HLS/RTL layers, IP per layer

Create stitched design

Network of HLS/RTL layers, stitched IP
**Ready to be integrated in Vivado IPI**

Prepare rtlsim (stitched)

Prepare rtlsim (layer by layer)

Full-network Verilator model

Network of HLS/RTL layers with Verilator models

PYNQ

Traditional HW Design RTL Simulation

Emulation (rtlsim) using PyVerilator

AMD
together we advance_

# FINN Workflow

Customization of Arithmetic

Customization of Hardware Architecture

FINN and Brevitas can be used as co-design tools to implement your DNN use case on an FPGA.

- Train a quantized neural network in PyTorch using Brevitas
- Converting trained QNN to Vivado IP
- Fine-tune model to meet resource/performance targets
- Integrate generated IP into a larger design

But you can leverage the infrastructure beyond that…

**Brevitas**

Brevitas FINN-ONNX Export

**Network Preparation**

Streamlining Transformations

Convert to HW Layers

Specialize Layers

Adjust folding to maximize performance

**Vivado HLS and IPI**

Create IP per layer

Create stitched design

Network of HLS/RTL layers, stitched IP
**Ready to be integrated in Vivado IPI**

PYNQ™

AMD
together we advance_

# Research in the FINN Ecosystem

Infrastructure to enable research on advanced quantization schemes and analysis of quantized neural networks

**Brevitas**
Quantization training library

**qonnx toolkit**
Infrastructure for backend-agnostic, mathematically preserving optimizations on (Q)ONNX models

Enables early design space exploration

**finn compiler**
Library of transformations/infrastructure to convert QNNs to FPGA design

**FINN library**
FPGA dataflow specific HW components

Infrastructure for research on neural network hardware design

**System integration**

Explore new optimized neural network layer implementations

together we advance_

# Status and Outlook

# Status Summary

> *"The FINN toolset is showing **huge potential using it in upcoming SICK products**.
> It is **easy to use** and with an **extraordinary performance** and very promising results.
> In the future, flexible implementations of ML in our products with FINN can be a great advantage and even replace static architectures as they are currently used.
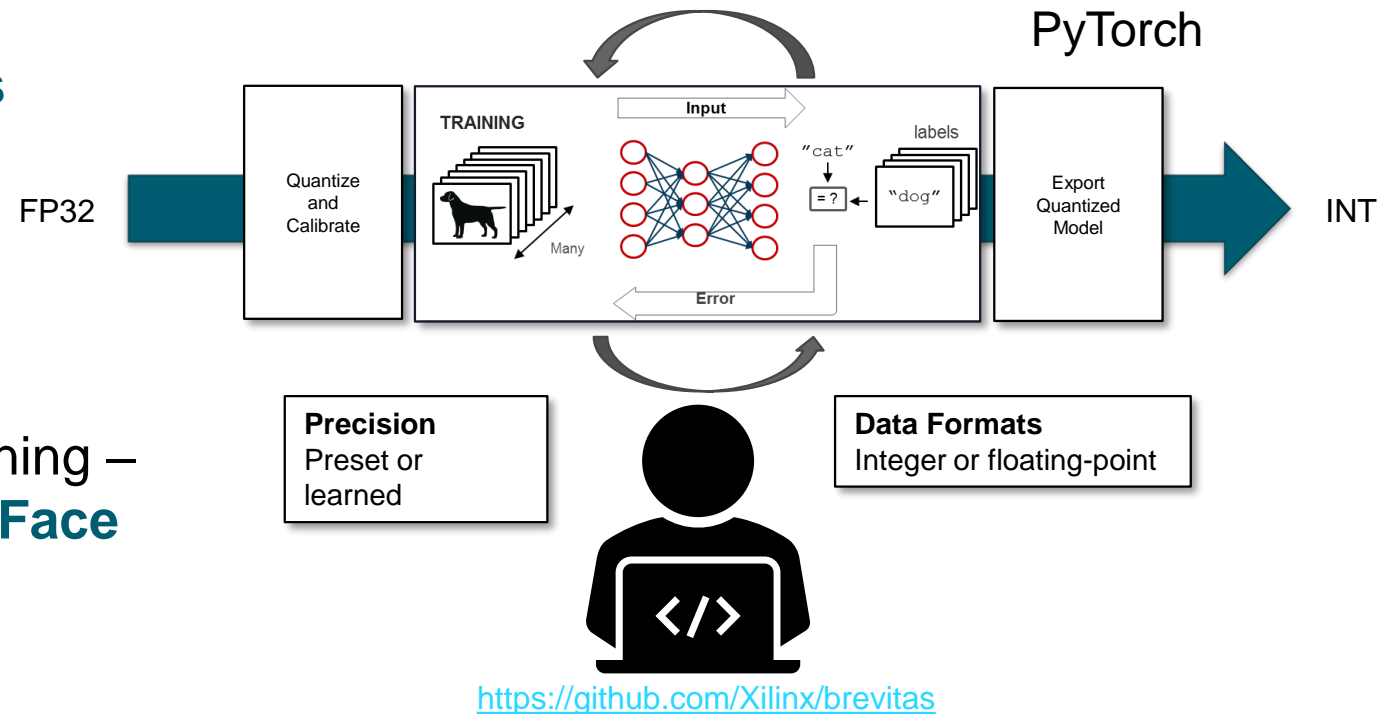> Thanks to the FINN team for the great cooperation"
> – Sick AG*

- **Open-Source Adoption**
  - ~2k+ GitHub stars summarized across repos
  - 250k+ Brevitas downloads
  - ~200k QONNX downloads
  - 17k+ FINN compiler downloads
- **Academic Results**
  - ACM TRETS 2020, FPL'2020, DFT'2019 Best Paper awards
  - 1000+ citations on original paper
- **University Classes on computer architecture for ML with FINN**
  - Stanford, UNC Charlotte, NTNU in Norway, EPFL in Switzerland
  - Regular tutorials, also available on YouTube: https://www.youtube.com/watch?v=zw2aG4PhzmA
- **Business units providing customer support**
  - Lead engineering team: Custom and Strategic Engineering, Dublin

https://github.com/Xilinx/brevitas
https://github.com/Xilinx/finn
https://github.com/Xilinx/finn-hlslib
https://github.com/Xilinx/finn-examples
https://github.com/fastmachinelearning/qonnx

**AMD**
together we advance_

# FINN Layer Support

| Layer | Current Support | Outlook |
|---|---|---|
| GEMM | ✓ | |
| Conv1D and Conv2D | ✓ | |
| - Dense | ✓ | |
| - Depthwise | ✓ | |
| - Separable (pointwise) | ✓ | |
| Elementwise (add, sub) | ✓ | others easily doable |
| Activation | ReLU, SeLU | |
| BatchNorm | ✓ (absorbed by streamlining) | |
| Pooling | ✓ | |
| Scale | ✓ (absorbed by streamlining) | |
| Concat | ✓ | |
| Reshape | ✓ (must be streamlinable) | |
| Transpose | ✓ (must be streamlinable) | |
| Clip by Value | ✓ (absorbed by streamlining) | |
| TransposeConv2D | ✓ | optimized version (WIP) |
| UpSample | ✓ | |
| DownSample | ✓ | |

AMD
together we advance_

# Brevitas Updates

- Targets the **entire AMD product range**

- First-class support for **integer datatypes**
  - prototype support for minifloats (e.g., FP8)

- Supports **PTQ and QAT**

- Out of the box support for distributed training – (e.g., DDP, interoperability with **HuggingFace Accelerate** (PP))

- Interoperability with **HuggingFace Transformers**

PyTorch

FP32

| Quantize and Calibrate | TRAINING | Export Quantized Model |

Input

labels

"cat"

= ?

"dog"

Many

Error

INT

**Precision**
Preset or learned

**Data Formats**
Integer or floating-point

https://github.com/Xilinx/brevitas

AMD
together we advance_

# FINN Compiler Updates

**FINN v0.10.1 Release**

- **Refactoring** of operator instantiation infrastructure
  - FINN compiler used to assume that hardware blocks are synthesized from HLS code
  - New class hierarchy to facilitate integration of RTL components
  - Provide users with an interface to override the compiler's choice
    for HLS vs. RTL implementation on a per-layer basis

- **RTL component** library optimizing the implementations of critical layers
  - Efficient implementation of 4-bit and 8-bit compute leveraging DSP slices
  - Efficient implementation of multi-level thresholding
  - Eradication of (regularly long) HLS synthesis times for layers with an RTL option

- **Compiler optimization pass** for accumulator and weight bit width minimization

- **Added board support** in system integration flow
  - **RFSoC 4x2** and **U55C** (contributed by University of Paderborn)

**AMD**
together we advance_

# FINN Technical Roadmap: Capabilities

- Operator Hardening

  - **Revised RTL Thresholding by binary search**

    - Ingestion of fp32 inputs

  - DSP-enabled Generalized Datatype Support

    - Efficient higher-precision integer compute: **`int4, int8`**, …, `int16`
    - Small standard floating-point formats: `float16, bfloat16`
    - Custom MiniFloats: `fp4 – fp8`

  - Internal clock pumping of DSP datapaths to increase their operational density
    We are aiming at a standard operational frequency around 500 MHz

- New Operators

  - Optimized **transposed convolution**

  - Fallback float layers to mitigate streamlining limits

AMD
together we advance_

# FINN Technical Roadmap: Ease of Use

- **FINN Library**
  - Refactoring of streamed layer interfaces
    - Packed flat `ap_uint<W>` → explicit `hls::vector<T, N>`
  - Combining HLS and RTL components into one FINN Library

- **FINN Examples**
  - MobileNet-v1 and VGG10-RadioML **with efficient DSP compute**
  - New example: German Traffic Sign Recognition Benchmark

**FINN-examples v0.0.7 Release**

**AMD**
together we advance_

# Resources

- [https://github.com/Xilinx/brevitas](https://github.com/Xilinx/brevitas)
  [https://github.com/Xilinx/finn](https://github.com/Xilinx/finn)

- [https://github.com/Xilinx/finn-hlslib](https://github.com/Xilinx/finn-hlslib)

- [https://github.com/Xilinx/finn-examples](https://github.com/Xilinx/finn-examples)

- [https://github.com/fastmachinelearning/qonnx](https://github.com/fastmachinelearning/qonnx)

- [https://amd.com/aup](https://amd.com/aup)

AMD
together we advance_

Q & A

# COPYRIGHT AND DISCLAIMER

**AMD**
together we advance_