



FWD NXT

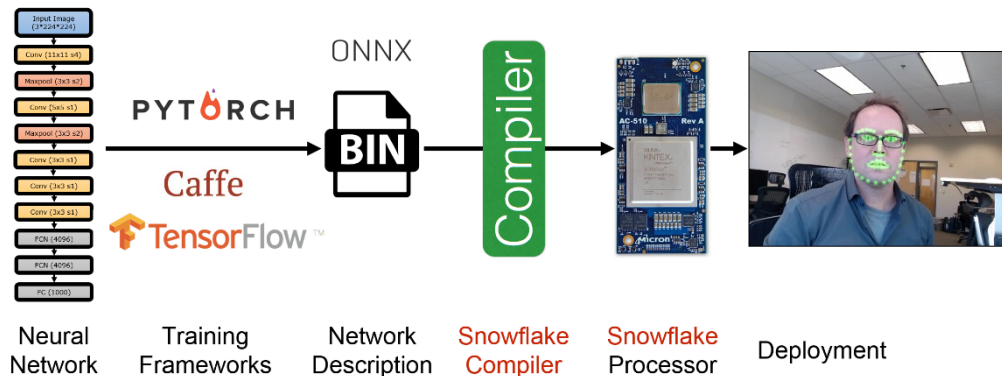
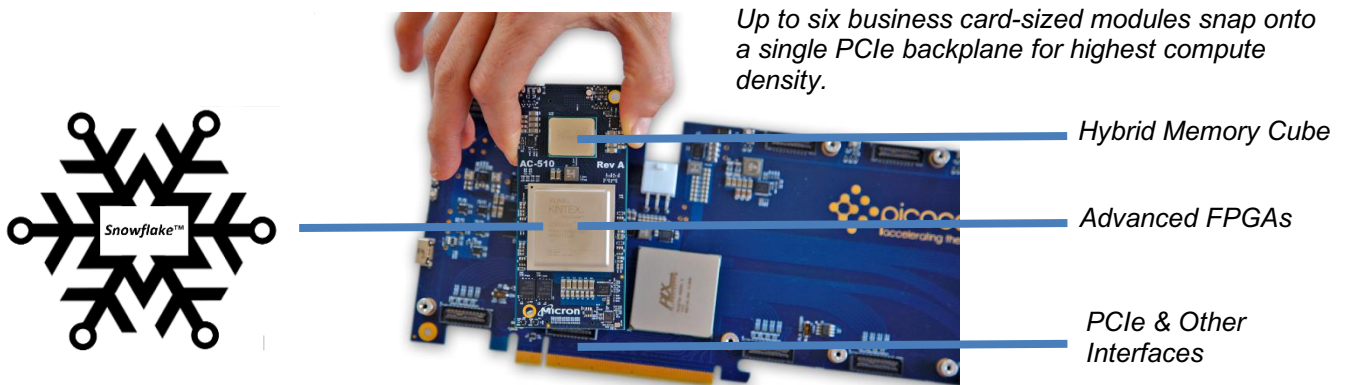
# Snowflake™

## Deep Neural Network Accelerator

### Machine Learning—Without the Need for Programming!

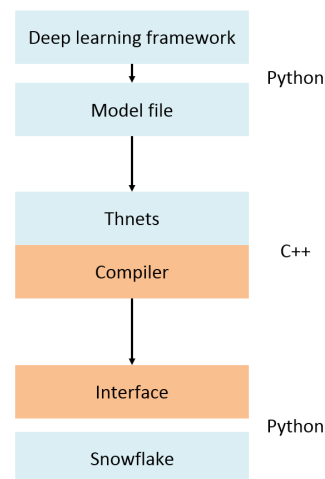
- Best performance/Watt of any deep learning solution
- Framework-agnostic development pathways—PyTorch, TensorFlow, Caffe
- Easy-to-use compiler requires no Verilog/VHDL expertise
- Most efficient scalability with highest compute density
- Exploits the power of FPGAs with the ease of GPUs

Our state-of-the-art deep learning solutions comprise a modular FPGA-based architecture with Micron’s Hybrid Memory Cube running Forward Next’s high-performance Snowflake neural network IP. Our fully integrated SDK takes trained neural network files and compiles them directly into the accelerator—*with no need for any programming*—enabling direct, rapid deployment from framework to application.



## OPERATION: COMPILER STEPS

- Model creation
- Parsing
- Partition + assignment
- Code generation
- Execute



## SYSTEMS SUMMARIES

Snowflake model	512	1K	2K	3K
<b>FPGA</b>	Micron AC510	2x Micron AC510	4x Micron AC510	6x Micron AC510
<b>Accelerator cores</b>	512	1024	2048	3072
<b>Clock Freq.</b>	187 MHz	187 MHz	187 MHz	187 MHz
<b>Peak Throughput</b>	191 G-ops/s	383 G-ops/s	766 G-ops/s	1,148 G-ops/s
<b>Memory</b>	4 GB HMC	8 GB HMC	16 GB HMC	24 GB HMC
<b>Memory B/W</b>	60 GB/s	120 GB/s	240 GB/s	360 GB/s
<b>Power system</b>	24 W	48 W	96 W	144 W

Snowflake model	256	1K	2K	3K
<b>FPGA</b>	1x Xilinx ZC706	2x Micron AC510	4x Micron AC510	6x Micron AC510
<b>Accelerator cores</b>	256	1024	2048	3072
<b>Clock Freq.</b>	250 MHz	187 MHz	187 MHz	187 MHz
<b>Peak Throughput</b>	128 G-ops/s	383 G-ops/s	766 G-ops/s	1,148 G-ops/s
<b>Memory</b>	1GB DDR3	8GB HMC	16GB HMC	24GB HMC
<b>Memory B/W</b>	4.2 GB/s	18 GB/s	36 GB/s	54 GB/s
<b>Power 1 board</b>	12 W	24 W	24 W	24 W
<b>Power system</b>	12 W	48 W	96 W	144 W
<b>Efficiency</b>	99.3%	98.2%	98.2%	98.2%