



FWD NXT

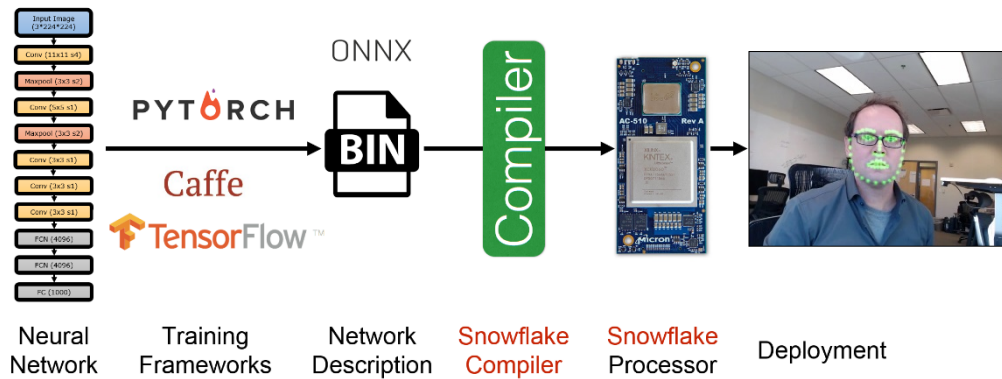
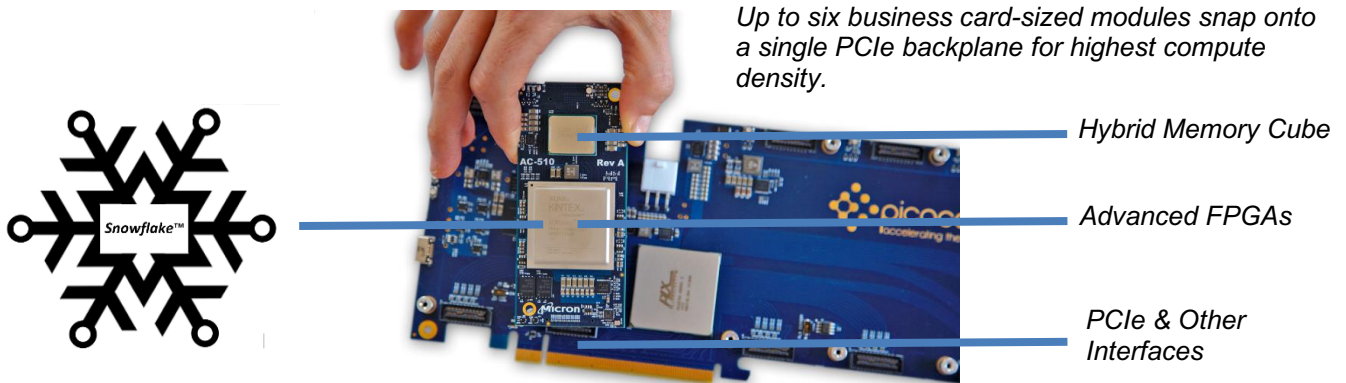
Snowflake™

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

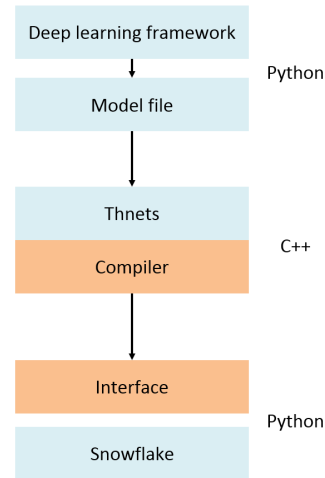
Our state-of-the-art deep learning solutions comprise a modular FPGA-based architecture with Micron’s Hybrid Memory Cube running Forward Next’s high-performance Snowflake neural network IP. Our fully integrated SDK takes trained neural network files and compiles them directly into the accelerator—*with 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	512-510	1k-511	512-852
FPGA	Micron AC510	Micron AC511	Micron SB-852
Accelerator cores	512	1024	1024
Clock Freq.	187 MHz	250 MHz	250 MHz
Peak Throughput	191 G-ops/s	512 G-ops/s	512 G-ops/s
Memory	4 GB HMC	2 GB HMC	512 GB DDR4 2 GB HMC
Memory B/W	60 GB/s	120 GB/s	120 GB/s
Power system	24 W	48 W	150 W (MAX)