



# Leveraging TVM as a front-end compiler for RISC-V based custom tinyML processor

Vaibhav Verma and Mircea R. Stan High-Performance Low-Power Lab (HPLP) University of Virginia vv8dn@virginia.edu





## Need for on-device tinyML

- Latency enable real-time AI systems
- **QoS** cannot rely on connectivity in remote areas
- Security sending data over network not secure
- Privacy keep private data locally on device
- Bandwidth send "information" to cloud rather than "data"
- Cost data communication is costly





### But tinyML/Edge AI is different!

Smaller neural network models

Smaller batch sizes (≈ 1)

Edge devices are power, cost, area and size limited

Edge processors support both AI and non-AI applications

Edge processors lack support for Keras, PyTorch, MXNet etc.



#### AI-RISC

Custom RISC-V processor with ISA extensions targeting AI applications

Tightly integrated AI accelerators for fine-grained offloading of AI tasks

End-to-end hardware/software codesign solution

Support for AI and non-AI applications on the same processor





## Hardware/Software Co-design

- End-to-end design methodology
  - TVM on the frontend
  - Synopsys ASIP Designer on the backend
- Support for multiple Domain Specific Language (DSL) frontends - Pytorch, MXNet, TFLite, Tensorflow, DarkNet etc.
- ✓ Verified with both 32 and 64-bit RISC-V
- ✓ Quantization support





#### Compilers are hard!

- > 2-step compilation TVM + ASIP Designer generated C compiler
- Goal is to have no/minimum interaction with TVM generated C code
- Problem 1 ASIP Compiler is not smart enough to detect opportunities for complex instructions like VMM, GEMM etc.
  - Works well for simple instructions like MAC.
- Solution Expose new instructions to TVM
  via compiler intrinsics.





## Compiler issues with custom instructions

- **Problem 2** Breaking the convolution schedule leads to accuracy issues.
- Issue 1 Wrong allocation of operands
  - Solved by TVM buffer and strided access of operands from bigger matrix
- Issue 2 Wrong data type in quantized NN
  - Defined input/output data types in TVM hardware intrinsic call
  - 8-bit inputs and 16/32-bit accumulated result.
- Issue 3 Kernel and Input data layout
  - TVM supports only a few Input/Filter layouts with specific ISA.
  - Adding support for required Input-Filter layout combinations in TVM for C hardware target used in AI-RISC.



### More Compiler issues

- **Problem 3** TVM support for breaking the convolution computation to match custom extension kernel size is limited.
  - TVM throws random errors when breaking the computation schedule using "tensorize" schedule pass.
  - Exact same convolution works with tensorize as a standalone kernel but not as a part of neural network.
- Solution Trying to debug the exact issue but till we find a reliable solution we work with what we have.
  - Breaking the computation schedule into error-free parts and adapting the custom instructions accordingly for testing purposes.



#### TVM generated C code





### Speedup on GEMV kernel

- A Matrix  $\rightarrow$  8x8
- B Vector  $\rightarrow$  1x8
- Input datatype  $\rightarrow$  int8
- Output datatype  $\rightarrow$  int16





# Speedup on single CONV2D kernel

- Input image  $\rightarrow$  7x7
- Input channels  $\rightarrow$  8
- Filter  $\rightarrow$  2x2
- Output channels  $\rightarrow$  2
- Input datatype  $\rightarrow$  int8
- Output datatype  $\rightarrow$  int16
- data\_layout  $\rightarrow$  NHWC
- kernel\_layout  $\rightarrow$  HWIO





#### Speedup on ResNet-8 network from MLPerf Tiny





#### Questions?



• This work is funded by SRC under GRC AIHW task 2945.001.