

# Supporting TVM on RISC-V Architectures with SIMD Computations

**Jenq-Kuen Lee<sup>1</sup>, Chun-Chieh Yang<sup>1</sup>, Allen Lu<sup>2</sup>, P. Chen<sup>1</sup>, YM Chang<sup>1,2</sup>,  
CH Chang<sup>1</sup>, Yi-Ru Chen<sup>1</sup>, HH Liao<sup>1</sup>, Chao-Lin Lee<sup>1,2</sup>, Ssu-Hsuan Lu<sup>2</sup>,  
and Shao-Chung Wang<sup>3</sup>**

<sup>1</sup>Department of Computer Science, National Tsing Hua University, Taiwan

<sup>2</sup>Peakhills Group Corporation

<sup>3</sup>Andes Technology Corporation



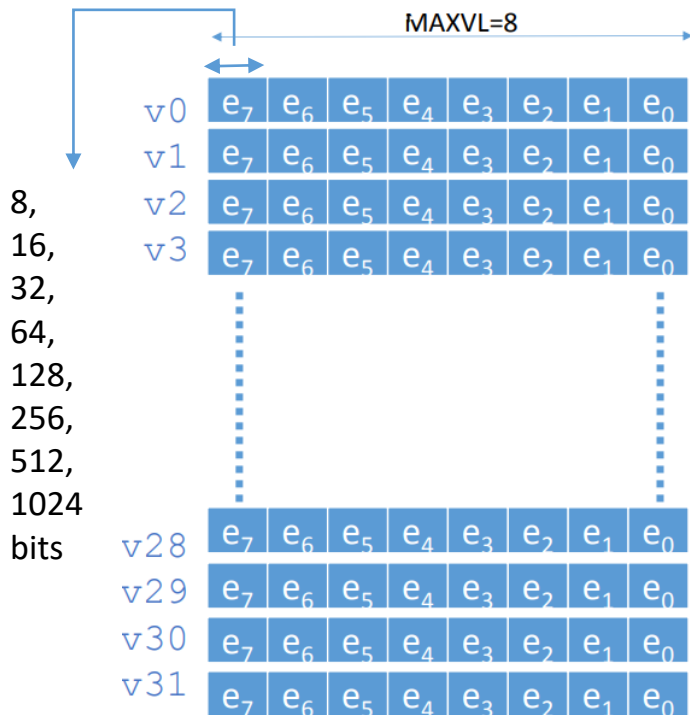
**PEAKHILLS**  
GROUP

**ANDES**  
TECHNOLOGY

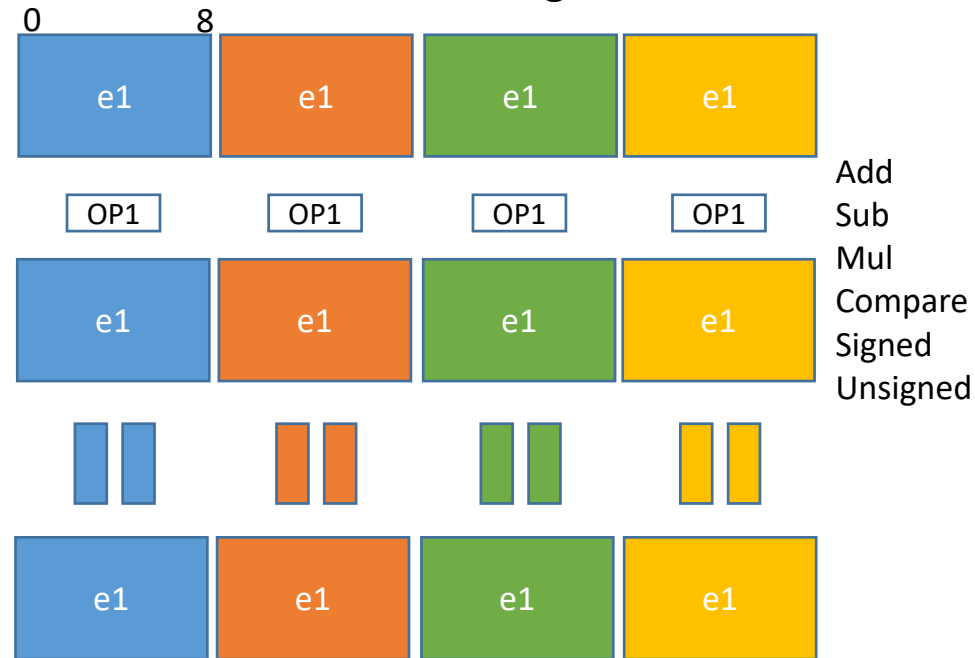
**PL** NTHU  
LAB

# RISC-V with two vector ISAs to support fall-back engine with AI Models

## Super Word Vector V Extension



## Packed Vector (SubWord SIMD) P Extension With Fixed-Point and Integer Instructions

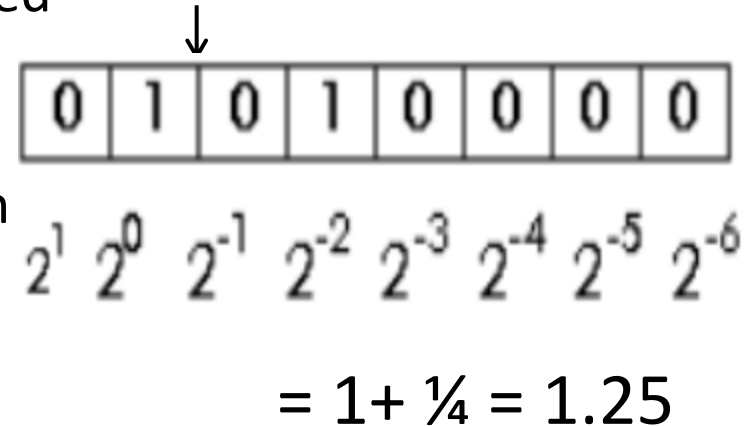


[RISC-V DSP \(P\) Extension Proposal Chuan-Hua Chang, Andes Technology Corporation](#)

Courtesy: Vector ISA, Roger Espasa, Esperanto Technologies

# [RFC] Fixed-point type implementation proposal #4446

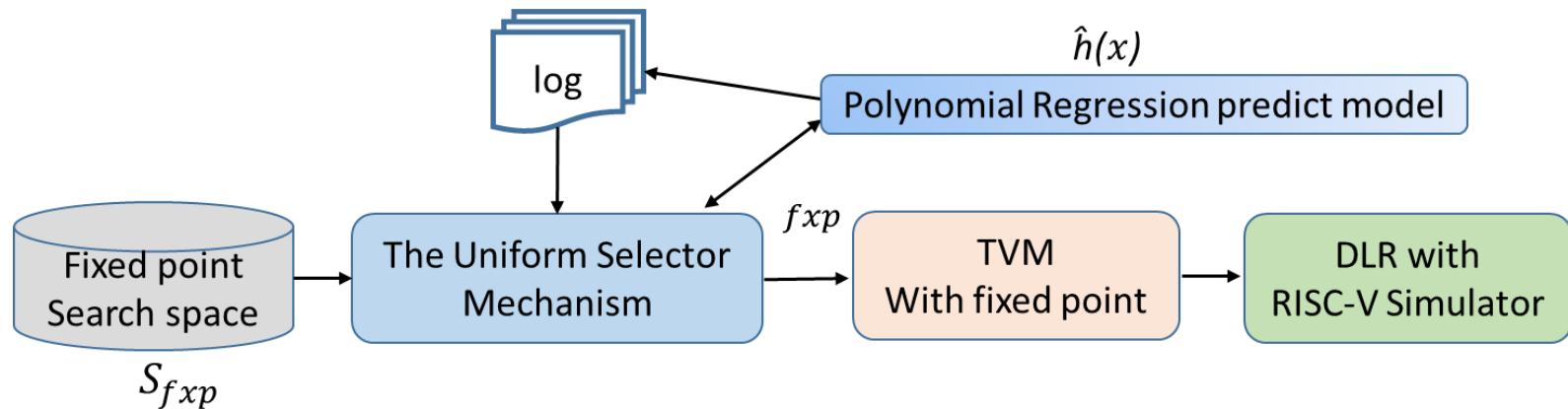
- RISC-V P extension (Subword SIMD) with fixed-point instructions.
- We refer Fxp as fixed-point value, Fp as floating-point value and PP as point position
  - $F_{xp} = F_p * \text{pow}(2, PP)$
- Support Fixed-point Type with TVM
- Compiler time with type information for the binary point position of the variable.



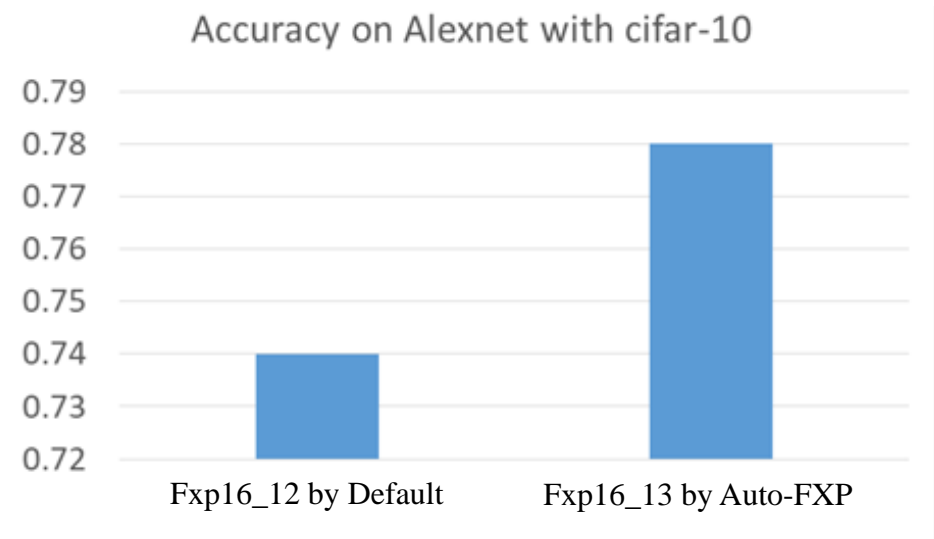
## References for Fixed-Point Type

- (1) AC fixed-Point by Mentor graphics (<https://www.mentor.com/hls-ip/downloads/ac-datatypes>)
- (2) Our early proposal to Khronos for OpenCL fixed-point feature set ([https://www.khronos.org/assets/uploads/developers/library/2018-khronos-group-openc1-embedded-outreach/Taipei-DSP-Profile-NTHU\\_Jan18.pdf](https://www.khronos.org/assets/uploads/developers/library/2018-khronos-group-openc1-embedded-outreach/Taipei-DSP-Profile-NTHU_Jan18.pdf))

# Auto-FXP with TVM on RISC-V with p Extension



- Using machine learning model to auto-tune the binary point position.
- It can find the best binary point position for fixed-point expression when we have TVM on RISC-V with p extension.
- The work extends AutoTVM and can enhance the accuracy while enjoy the low power numeric benefits.
- The tuning work is done with spike simulator incorporated with RISC-V P extension (Subword SIMD).

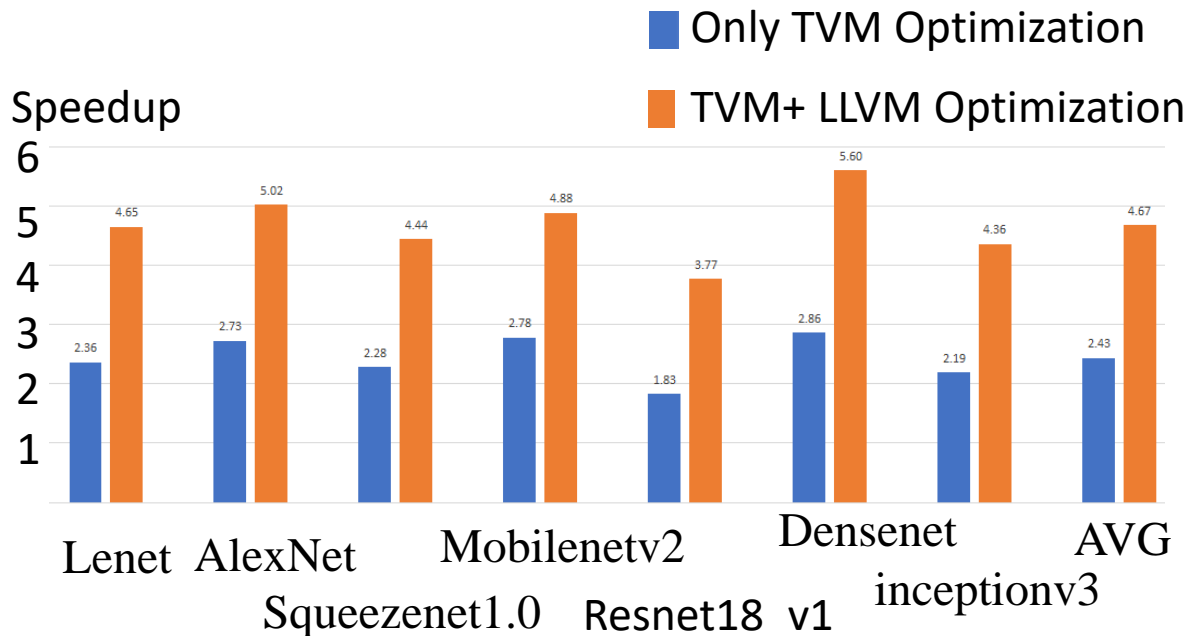


# TVM for RISC-V with V Extension (Superword SIMD)

- TVM Optimization
  - The TVM RISC-V codegen will lower SIMD computation with SIMD intrinsics into LLVM.
  - The LLVM backend will need to generate the corresponding SIMD instructions.
  - Need to tune the scheduler to provide a large loop index space for vector parallelism.
- LLVM Optimization
  - VSETVL Redundancy Elimination
  - VMulADD Resource Utilization
  - Fast Vector Initializer

Speedup based on runtime executed instructions

- Spike Simulator
- assume 512 bits vector register
- V SIMD in <4 x float32>, <8 x float32>, <16 x float32>
- Spec v0.7.0, TVM v0.6, LLVM 9.0.0
- Compare with SIMD float32 and no SIMD float32



# Summary

- Thank you AWS team help with AI model validation flow.
- Look forward to contributing codes to the TVM source trees.
- More detailed of our work can also be found in the following.
  - Experiments and AI Model Validations for Neo/TVM on RISC-V Architectures with SIMD, Allen Lu, et al, RISC-V Summit, San Jose, Dec 2019 (Poster).
  - Enabling TVM on RISC-V Architectures with SIMD Instructions, Allen Lu, Chao-Lin Lee, Yuan-Ming Chang, Piyo Chen, Hsiang-Wei Sung, Heng Lin, Shao-Chung Wang, and Jenq-Kuen Lee, RISC-V Forum, March 2019 (Oral presentation).