



Unified Posit/IEEE-754 Vector MAC Unit for Transprecision Computing

Luís Crespo; Pedro Tomás; Nuno Roma; Nuno Neves

INESC-ID and Instituto Superior Técnico,
Universidade de Lisboa, in Lisbon, Portugal

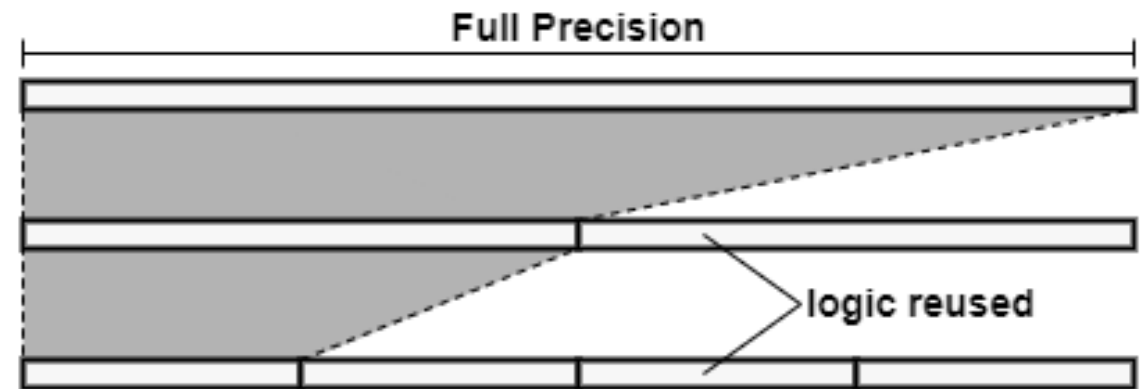
2022 IEEE International Symposium on Circuits and Systems
May 28- June 1, 2022 Hybrid Conference

Outline

- Introduction
- Proposed Architecture
 - Overview
 - Vector Structures
 - Unified Decode and Encode
 - Quire Scale and Accumulate
- Implementation Results
- Conclusion

Introduction

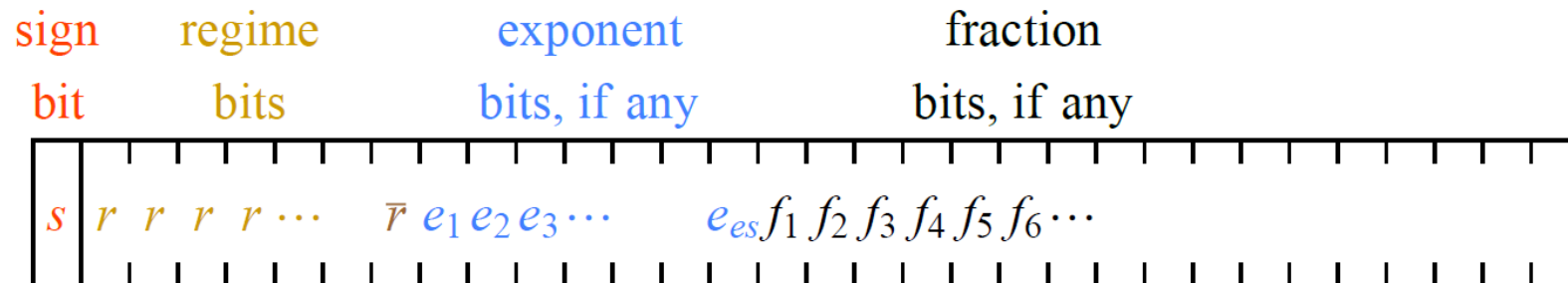
- Transprecision for performance and energy efficiency demands
- Different precisions by instantiating multiple arithmetic modules
- Vectorized Datapath
 - Different precisions with the same hardware resources
 - SIMD capabilities
 - Limited by the IEEE-754 standard



Introduction

- Posit Format

- Parameterizable precision and dynamic range $\langle n, es \rangle$
- low-precision and fused operations (quire)
- $(-1)^S \times 2^{e+k2^{es}} \times 1.f$

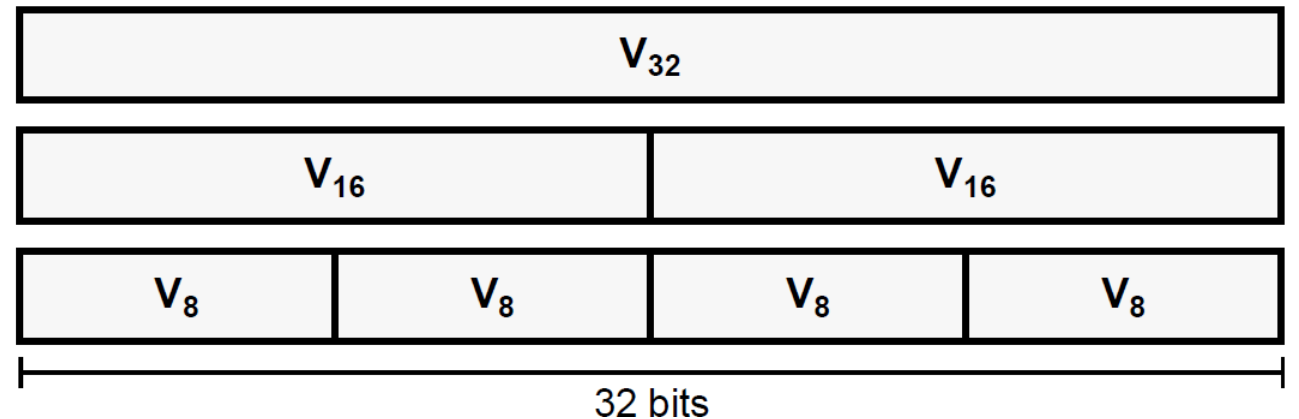


- Prohibitive overheads with quire
- Maintain compatibility with IEEE-754

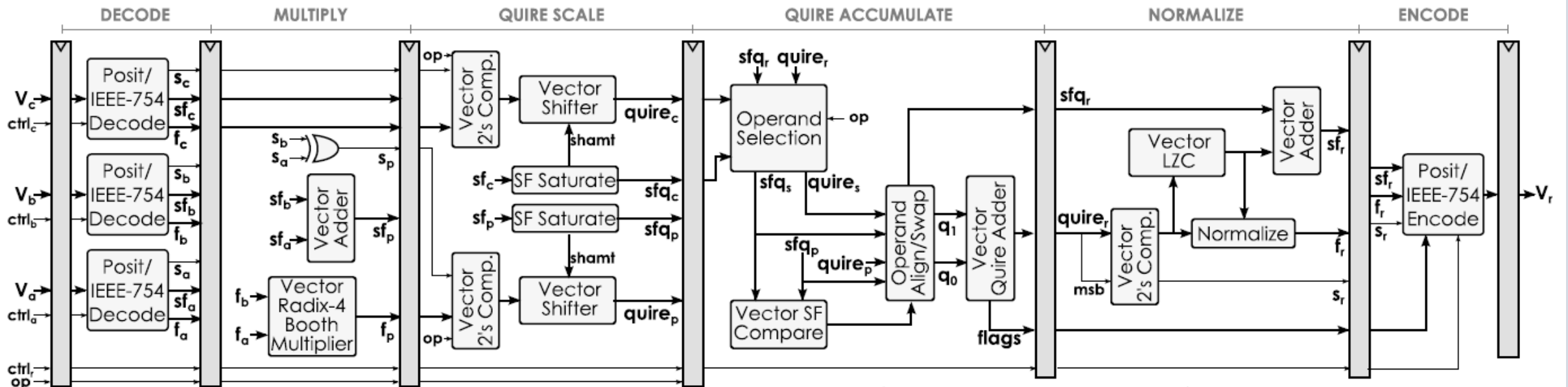
Proposed Architecture - Overview

- Floating-Point format unification
 - $(-1)^S \times 2^{sf} \times 1.f$
- Variable-precision MAC architecture with dynamic vectorization
 - 1x32-bit, 2x16-bit, and 4x8-bit vector operations
- Variable-exponent Posit configuration
 - Reduced quire

- **Dynamically configurable**



Proposed Architecture - Overview



Control Signals:

op - operation
 $ctrl_{x,pre}$ - precision (8/16/32 bits)
 $.vec$ - vector or scalar
 $.fmt$ - format (Posit/IEEE-754)
 $.es$ - posit exp. size

Supported Operations:

single
 $V_r = V_a * V_b$
 $V_r = V_a + V_c$
 $V_r = V_a - V_c$

fused
 $V_r = V_a * V_b + V_c$
 $V_r = V_a * V_b - V_c$
 $V_r += V_a * V_b$
 $V_r -= V_a * V_b$

Legend:

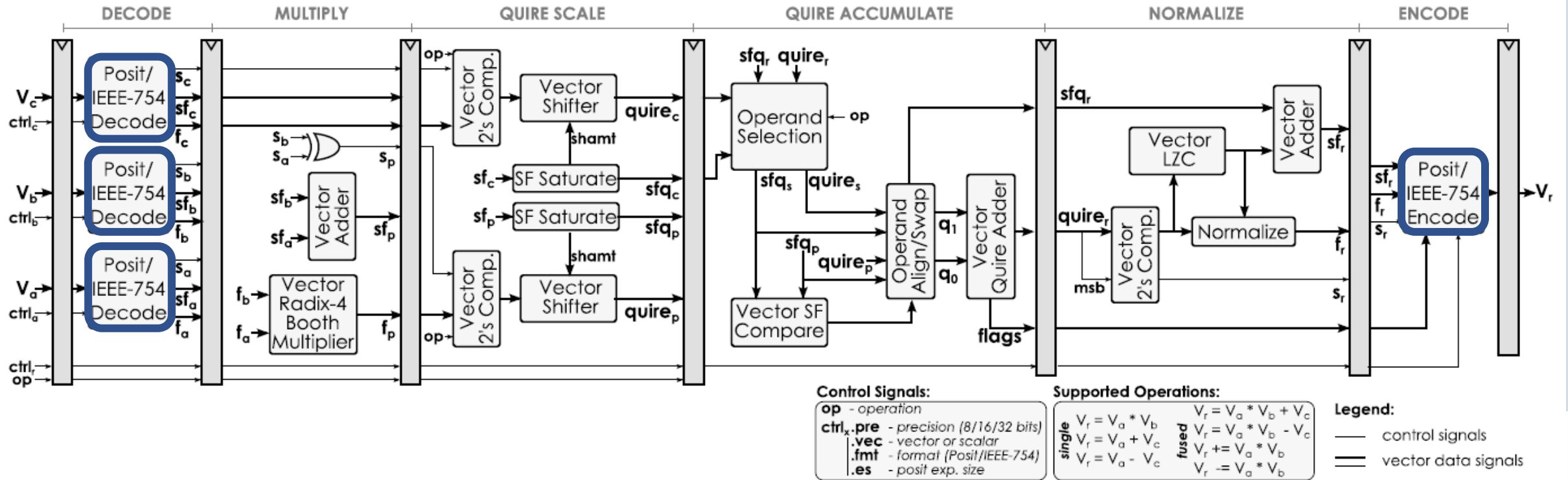
— control signals
 — vector data signals

- Precision
- Operation
- Vector or scalar
- Format
- Posit exponent size

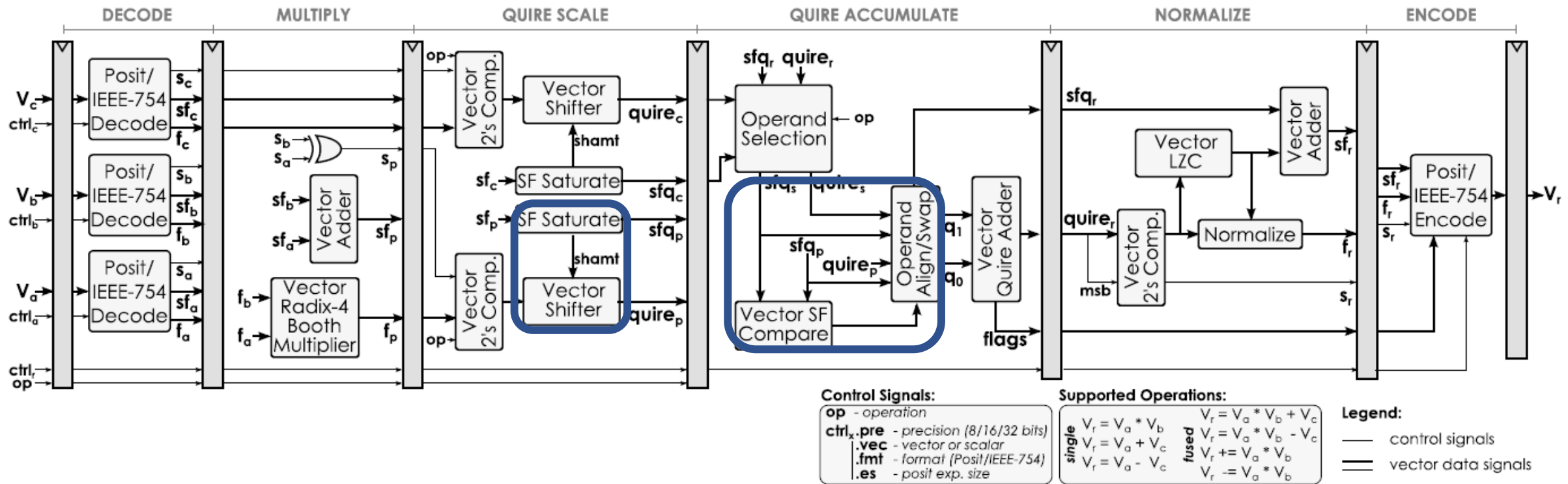
Proposed Architecture - Vector Structures

- Vectorized adder
 - Strategically placed multiplexers
 - Sub-adder carry-in is the carry-out or the input carry-in
- Vectorized barrel shifter
 - One shifting level for each vector configuration
 - Cropped or OR'ed with the adjacent vector element
- Vectorized LZC
 - Tree-like structure
 - Capture intermediate results
- Vectorized radix-4 Booth multiplier

Proposed Architecture - Unified Decode and Encode



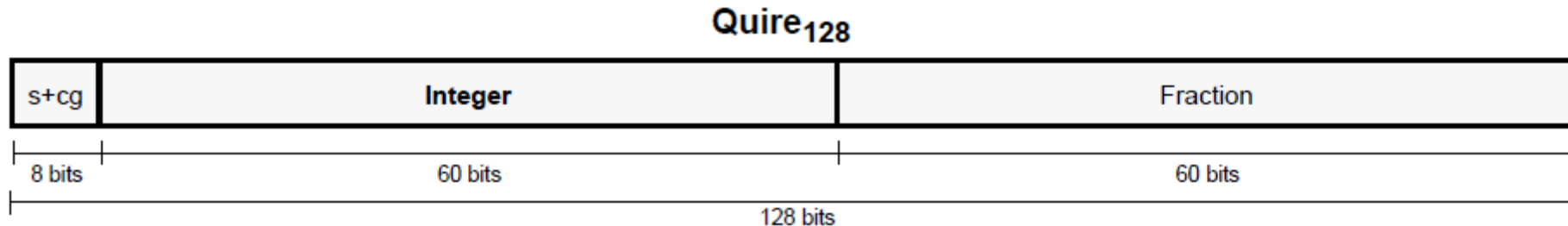
Proposed Architecture - Quire Scale and Accumulate



Proposed Architecture - Quire Scale and Accumulate

- Conversion

- Shift amount cannot be greater than the integer size
- Saturates the shifting amount
- No right shift



- Alignment

- Determine operand with lower scale factor
- Shift it by the difference

Implementation Results

UNIT	NUM. BITS	PIPEL. STAGES	ASIC TECH.	DELAY (ns)	AREA (μm^2)	POWER (mW)	PERF. (GOPS)	AREA EFF. ($\times 10^{-6}$ GOPS/ μm^2)	EDP ($\times 10^{-22}$ J.s)
Ref. Posit Std. MAC	8	5	28 nm	0.65	7598	21	1.54	202.4	0.89
Ref. Posit Std. MAC	16	5	28 nm	0.8	17384	47	1.25	71.91	3.01
Ref. Posit Std. MAC	32	5	28 nm	0.91	39767	108	1.10	27.63	8.94
Proposed VMAC	8/16/32	6	28 nm	1.5	51563	99	2.7/1.3/0.7	51.7/25.7/12.9	5.6/11.1/22.3
Posit DFMA [10]	32	5	45 nm	1.5	112350	370	0.67	5.95	83.25
FP VFMA [13]	16/32/64	3	90 nm	1.5	180610	44	2.7/1.3/0.7	14.8/7.4/3.7	2.5/4.9/9.9
Posit VMULT [14]	8/16/32	-	90 nm	2.3	91861	64	1.7/0.9/0.4	18.9/9.5/4.7	8.5/16.9/33.9

- Vs 32-bit Posit MAC
 - 30% area increase, similar power consumption
 - More area- and energy-efficient
- Vs transprecision architecture estimation
 - 50% less area and 2.9x less power

Implementation Results

UNIT	NUM. BITS	PIPEL. STAGES	RESULT SOURCE	DELAY (<i>ns</i>)	AREA (μm^2)	POWER (<i>mW</i>)
Proposed VMAC	8/16/32	6	Synthesis	1.5	51563	99
Posit DFMA [10]	32	5	Estimated	1.24	39324	266
FP VFMA [13]	16/32/64	3	Estimated	0.77	16044	21
Posit VMULT [14]	8/16/32	-	Estimated	1.18	8160	31

- Vs DFMA
 - HW requirements mitigated by the 128 bits quire
- Vs VFMA and VMULT
 - Decoding, encoding, quire (50%) and pipeline increase HW requirements
- VMAC presents a much higher functionality

Conclusion

- Variable-precision datapath with SIMD processing capabilities
- Low- and high-precision operations
- Unique support for Posit and IEEE-754
- Dynamic exponent and reduced quire
- Without requiring a prohibitive chip area size
- Higher functionality



Q&A

2022 IEEE International Symposium on Circuits and Systems
May 28- June 1, 2022 Hybrid Conference

