

基於RISC-V 向量指令集之 Softmax 運算加速

Softmax Acceleration on RISC-V Vector Extension

學生: 黃士昕

指導老師: 賴謹峰 博士

Student: Shih-Hsin Huang

Advisor: Dr. Chin-Feng Lai



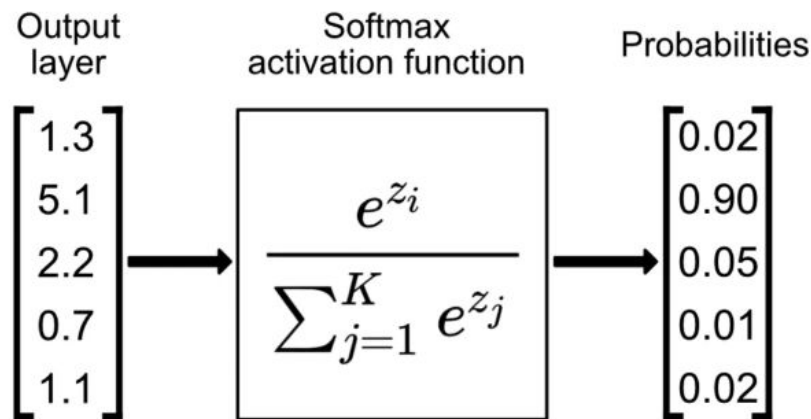
Outline

- Introduction
- Related Work
- Method
- Experiment and Evaluation
- Conclusion and Future Work

Introduction

- Research Motivation
- Research Objectives

Research Motivation - Softmax



CNN

DNN

RNN

TRANSFORMER

Challenges:

- High Computation Cost
- Numerical Instability

Research Motivation - Edge AI

Algorithm 3 Safe softmax

1: $m_0 \leftarrow -\infty$		
2: for $k \leftarrow 1, V$ do		
3: $m_k \leftarrow \max(m_{k-1}, x_k)$	Pass 1	▸ Pass 1: read X
4: end for		
5: $d_0 \leftarrow 0$		
6: for $j \leftarrow 1, V$ do	Pass 2	
7: $y_i = EXP(X_j - m_v)$		▸ Pass 2: read X, Write Y
8: $d_j \leftarrow d_j + y_i$		
9: end for		
10: for $i \leftarrow 1, V$ do	Pass 3	
11: $y_i \leftarrow \frac{y_i}{d_v}$		▸ Pass 3: read Y, Write Y
12: end for		



Challenges:

- Memory Bandwidth Bottleneck

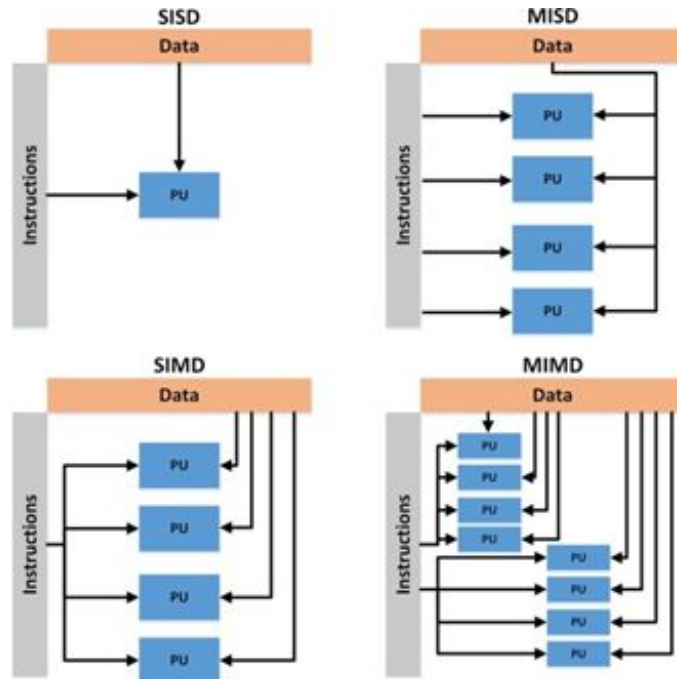
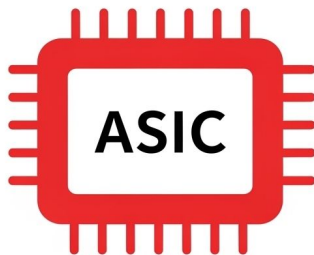
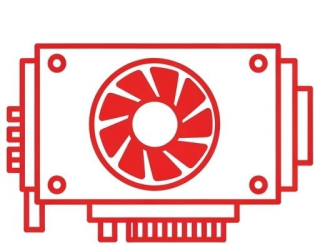
- Limited Bandwidth
- Power Sensitive

Research Motivation - RISC-V

Open Standard



Proprietary / Fragmented



RVV: Write Once, Run Anywhere

Research Objectives

- **Algorithm Design and Innovation:** To propose a Conditional Online Softmax algorithm
- **High-Performance Implementation on RISC-V Vector Extension:**
 - Softmax Function
 - Exponential Function
- **Comprehensive Performance Evaluation and Bottleneck Analysis**

Related Work

- Softmax and Numerical Stability
- A Review of Optimization Strategies and Their Limits
 - Online Softmax
 - Exponential Function
- Programming with the RISC-V Vector Extension

Background: Safe Softmax

Algorithm 3 Safe softmax

1: $m_0 \leftarrow -\infty$

2: **for** $k \leftarrow 1, V$ **do**

3: $m_k \leftarrow \max(m_{k-1}, x_k)$

4: **end for**

5: $d_0 \leftarrow 0$

6: **for** $j \leftarrow 1, V$ **do**

7: $y_j = \text{EXP}(x_j - m_v)$

8: $d_j \leftarrow d_j + y_j$

9: **end for**

10: **for** $i \leftarrow 1, V$ **do**

11: $y_i \leftarrow \frac{y_i}{d_v}$

12: **end for**

3 pass

▸ Pass 1: read X

▸ Pass 2: read X, Write Y

▸ Pass 3: read Y, Write Y

$$y_i = \frac{e^{x_i - \max_{k=1}^V x_k}}{\sum_{j=1}^V e^{x_j - \max_{k=1}^V x_k}}$$

Exponent	fraction = 0	fraction \neq 0	Equation
00 _H = 00000000 ₂	\pm zero	subnormal number	$(-1)^{\text{sign}} \times 2^{-126} \times 0.\text{fraction}$
01 _H , ..., FE _H = 00000001 ₂ , ..., 11111110 ₂	normal value		$(-1)^{\text{sign}} \times 2^{\text{exponent}-127} \times 1.\text{fraction}$
FF _H = 11111111 ₂	\pm infinity	NaN (quiet, signaling)	



Background: Online Softmax

Algorithm 4 Safe softmax with online normalizer calculation

1: $m_0 \leftarrow -\infty$

2: $d_0 \leftarrow 0$

3: **for** $j \leftarrow 1, V$ **do**

4: $m_j \leftarrow \max(m_{j-1}, x_j)$

5: $d_j \leftarrow d_{j-1} \times [e^{m_{j-1}-m_j} + e^{x_j-m_j}]$

6: **end for**

7: **for** $i \leftarrow 1, V$ **do**

8: $y_i \leftarrow \frac{e^{x_i-m_V}}{d_V}$

9: **end for**

2 pass

► Pass 1: read X

Data
Dependency

► Pass 2: read X, Write Y

3N expf calls

This can be proven by induction.

Background: Exponential Function

Step 1: Range Reduction

$$e^x = 2^{x \cdot \log_2(e)} \quad t = x \cdot \log_2(e)$$

$$t = i + f, i = \lfloor t \rfloor, f = t - i$$

Step 2: Divide and Conquer

$$e^x = 2^i \cdot 2^f$$

Bit Manipulation

Polynomial Approximation

Look up table

`glibc expf()`

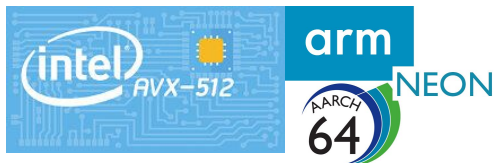
Step 3: Reconstruction



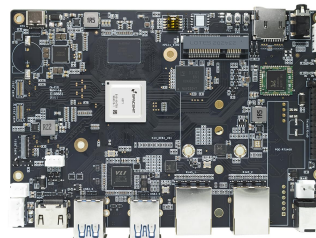
Background: RISC-V Vector Extension

RISC-V "V" Vector Extension

Version 1.0



vlen = 256b



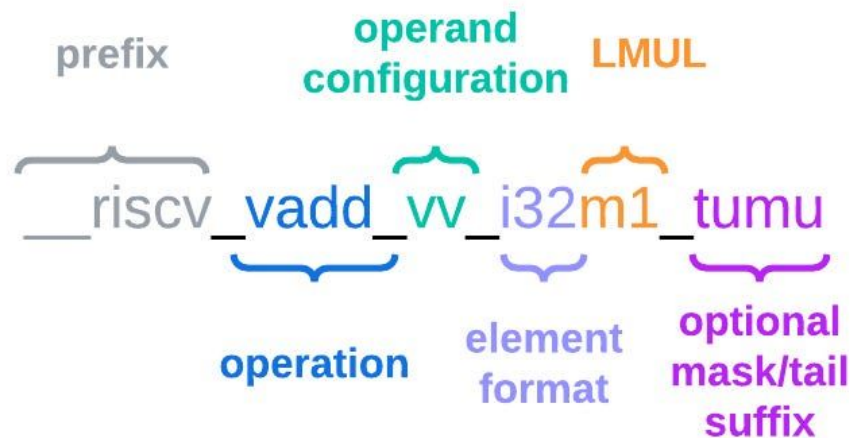
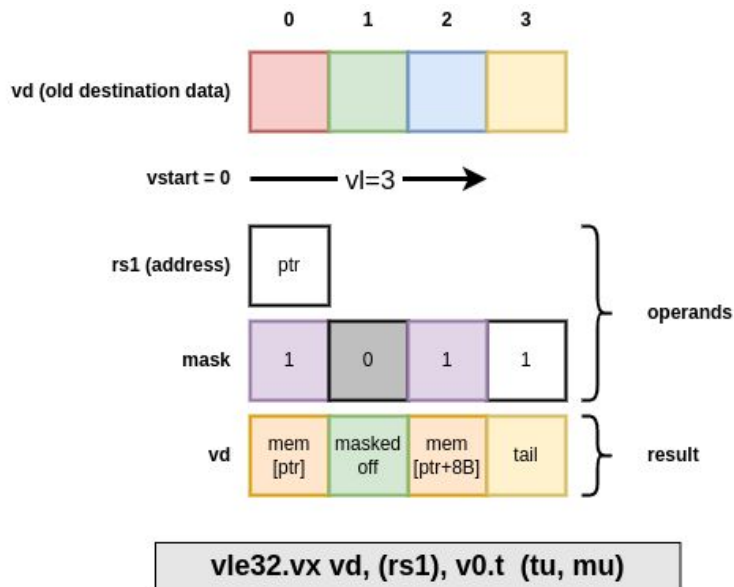
vlen = 512b



RISC-V RVA23—A
Major Milestone



Background: RISC-V Vector Extension



Background: Research Gap and Our Contributions

The Research Gap		Our Contributions	
Algorithmic Dilemma	Safe Softmax: <ul style="list-style-type: none">• High Memory Access Cost(3-Pass)	Conditional Online Softmax	<ul style="list-style-type: none">• Balance the trade-off between memory access and computational cost
	Online Softmax: <ul style="list-style-type: none">• High Computation Cost($3N \expf()$),• Data Dependency		
Implementation Bottleneck	<ul style="list-style-type: none">• Compiler auto-vectorization ineffective	Manual Vectorization	<ul style="list-style-type: none">• RVV Intrinsics manual Vectorization
	<ul style="list-style-type: none">• Slow $\expf()$ Implementation		<ul style="list-style-type: none">• High-performance $\text{vexpf}()$ function



Methodology

- Condition Online Softmax
- Branch-free Conditional Online Softmax
- Vectorize Exponential Function

Methodology: Condition Online Softmax

Algorithm 5 Condition Online Softmax

```
1:  $m_0 \leftarrow -\infty$ 
2:  $d_0 \leftarrow 0$ 
3: for  $j \leftarrow 1, V$  do
4:   if Updated the maximum value then
5:      $d_j \leftarrow d_j \times e^{m_{j-1}-x} + 1$ 
6:   else
7:      $d_j = d_j + e^{x-m_{j-1}}$ 
8:   end if
9: end for
10: for  $i \leftarrow 1, V$  do
11:    $y_i \leftarrow \frac{e^{x_i-m_V}}{d_V}$ 
12: end for
```

An algorithmic principle: 'Each time a new maximum is found, the accumulated sum must be rescaled.'

► Pass 1: read X

2N expf() calls!

► Pass 2: read X, Write Y

Methodology: The Record-Breaker Problem

The Record-Breaker Problem is a classic model in probability theory that examines the expected number of record-breaking events within a sequence of random variables. An event is defined as a "record-breaker" if its observed value is greater than all values that have preceded it in the sequence

$$d_j = \begin{cases} (d_{j-1} \times e^{m_{j-1}-m_j} + 1), & \text{if } x_i > m_{i-1}, \\ (d_{j-1} \times 1 + e^{x_j-m_j}), & \text{otherwise.} \end{cases} \leftarrow$$

向量長度 (n)	$H_n \approx \ln n + \gamma$	更新比例
10	2.88	0.3
100	5.18	0.04
1000	7.47	0.007
10000	9.79	0.0009
100000	12.08	8.0e-5
1000000	14.40	1.3e-5
10000000	16.70	1.6e-5

18

Methodology: Branch-free Conditional Online Softmax

► Phase 2: Cross-Lane Reduction

35: **function** PHASE2_CROSSLANEREDUCTION(v_{\max} , v_{sum} , active_lanes)

36: **Input:** Per-lane results v_{\max} , v_{sum} , and the number of active lanes.

37: **Output:** The final global maximum m_{final} and normalizer d_{final} .

38:

39: $m_{\text{final}} \leftarrow v_{\max}[0]$

40: $d_{\text{final}} \leftarrow v_{\text{sum}}[0]$

41: **for** $i \leftarrow 1$ **to** active_lanes - 1 **do**

42: $m_{\text{next}} \leftarrow v_{\max}[i]; \quad d_{\text{next}} \leftarrow v_{\text{sum}}[i]$

43: $m_{\text{new}} \leftarrow \max(m_{\text{final}}, m_{\text{next}})$

44: $d_{\text{final}} \leftarrow d_{\text{final}} \cdot \text{EXP}(m_{\text{final}} - m_{\text{new}}) + d_{\text{next}} \cdot \text{EXP}(m_{\text{next}} - m_{\text{new}})$

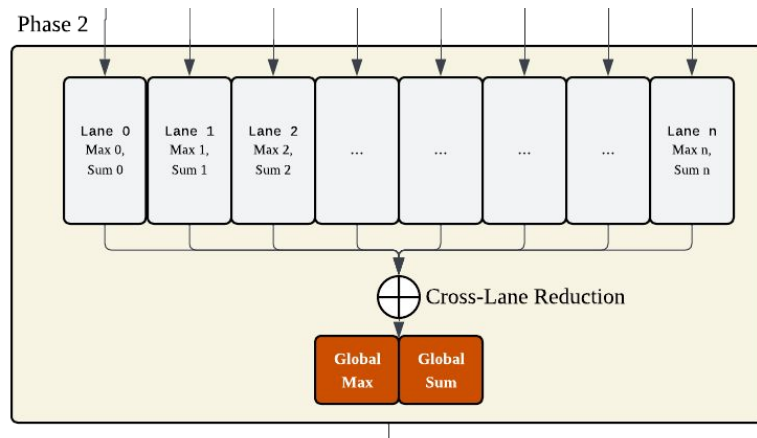
45: ► Execute the \oplus merge operation

46: $m_{\text{final}} \leftarrow m_{\text{new}}$

47: **end for**

48: **return** ($m_{\text{final}}, d_{\text{final}}$)

49: **end function**



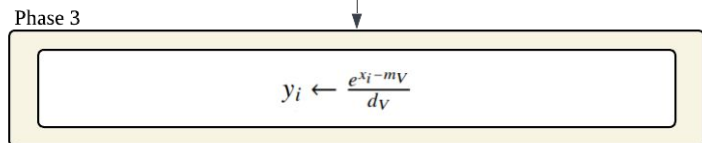
$$\begin{bmatrix} m_i \\ d_i \end{bmatrix} \oplus \begin{bmatrix} m_j \\ d_j \end{bmatrix} = \begin{bmatrix} \max(m_i, m_j) \\ d_i \cdot \exp(m_i - \max(m_i, m_j)) + d_j \cdot \exp(m_j - \max(m_i, m_j)) \end{bmatrix}$$

Methodology: Branch-free Conditional Online Softmax

► Phase 3: Final Normalization

```
50: function PHASE3_FINALNORMALIZATION( $X, Y, N, m_{final}, d_{final}$ )  
51:   while  $N > 0$  do  
52:      $v_l \leftarrow \text{SETVECTORLENGTH}(N)$   
53:      $v_{data} \leftarrow \text{VECTORLOAD}(X, v_l)$   
54:      $v_{data} \leftarrow v_{data} - m_{final}$   
55:      $v_{data} \leftarrow \text{VECTOREXP}(v_{data}, v_l)$   
56:      $v_{result} \leftarrow v_{data} / d_{final}$   
57:      $\text{VECTORSTORE}(Y, v_{result}, v_l)$   
58:      $X \leftarrow X + v_l; \quad Y \leftarrow Y + v_l; \quad N \leftarrow N - v_l$   
59:   end while  
60: end function
```

Phase 3


$$y_i \leftarrow \frac{e^{x_i - m_V}}{d_V}$$

Methodology: Vectorize Exponential Function

Algorithm 7 Vectorize Exponential Function

```

1: function VECTOREXPFAST(x,vl)
2:   Input: An input vector x of length vl.
3:   Output: An approximate result vector for  $e^x$ .
4:   Constants:
       $C_{\log 2e} \leftarrow 1.442695041$ 
       $P_A, P_B, P_C \leftarrow 0.3371..., 0.6576..., 1.0017...$ 
       $\triangleright$  Coefficients for polynomial approximation
  
```

```

5:    $t \leftarrow x \cdot C_{\log 2e}$   $\triangleright$  Phase 1: Range Reduction
6:    $i \leftarrow \text{VECTORFLOOR}(t)$   $\triangleright$  Convert base from  $e$  to 2
7:    $f \leftarrow t - i$   $\triangleright$  Get integer part  $i = \lfloor t \rfloor$ 
    $\triangleright$  Get fractional part  $f = t - i$ , where  $f \in [0, 1)$ 
  
```

```

8:    $\triangleright$  Phase 2: Polynomial Approximation for  $2^f$ 
    $\triangleright$  Use 2nd-degree polynomial evaluated with Horner's method:  $(A \cdot f + B) \cdot f + C$ 
9:    $v_A \leftarrow \text{VECTORBROADCAST}(P_A, vl)$ 
10:   $v_B \leftarrow \text{VECTORBROADCAST}(P_B, vl)$ 
11:   $v_C \leftarrow \text{VECTORBROADCAST}(P_C, vl)$ 
12:   $p_1 \leftarrow \text{VECTORFMA}(f, v_A, v_B, vl)$   $\triangleright$  Parallel Fused-Multiply-Add:  $f \cdot A + B$ 
13:   $\text{exp2}_f \leftarrow \text{VECTORFMA}(p_1, f, v_C, vl)$   $\triangleright$  Parallel FMA:  $(f \cdot A + B) \cdot f + C$ 
  
```

```

14:   $\triangleright$  Phase 3: Bit-level Reconstruction for  $2^i \cdot 2^f$ 
    $\triangleright$  This phase computes the result by adding  $i$  to the exponent field of  $\text{exp2}_f$ 
15:   $\text{bits}_f \leftarrow \text{REINTERPRETASINT}(\text{exp2}_f)$   $\triangleright$  Get the integer bit representation of  $2^f$ 
16:   $\text{bits}_i \leftarrow i \ll 23$   $\triangleright$  Represents multiplication by  $2^i$  for float32
17:   $\text{bits}_{\text{final}} \leftarrow \text{bits}_f + \text{bits}_i$   $\triangleright$  Merge the exponents via integer addition
18:   $\text{result} \leftarrow \text{REINTERPRETASFLOAT}(\text{bits}_{\text{final}})$   $\triangleright$  Convert final bit pattern back to float
19:
20:  return result
  
```

Experiment and Evaluation

- Experimental Setup
- Results and Discussion

Experiment and Evaluation: Hardware Setup

Main Experimental Platform	Performance Reference Platform
RISC-V Platform	x86 Platform
Board: Banana Pi BPI-F3	CPU: Intel Core i7-13700
CPU: Spacemit X60 (8 Cores @ 1.6 GHz)	Cores: 8 P-cores + 8 E-cores (@ up to 5.2 GHz)
ISA: RV64GCV 1.0	ISA: x86-64 (with AVX2)
Vector Spec: VLEN = 256-bit	Vector Spec: AVX2 (256-bit)
Role in Study: <ul style="list-style-type: none">- Algorithm Development and Validation- Vectorization Performance Analysis	Role in Study: <ul style="list-style-type: none">- Industry-Standard Performance Reference- Cross-Microarchitecture Behavior Comparison

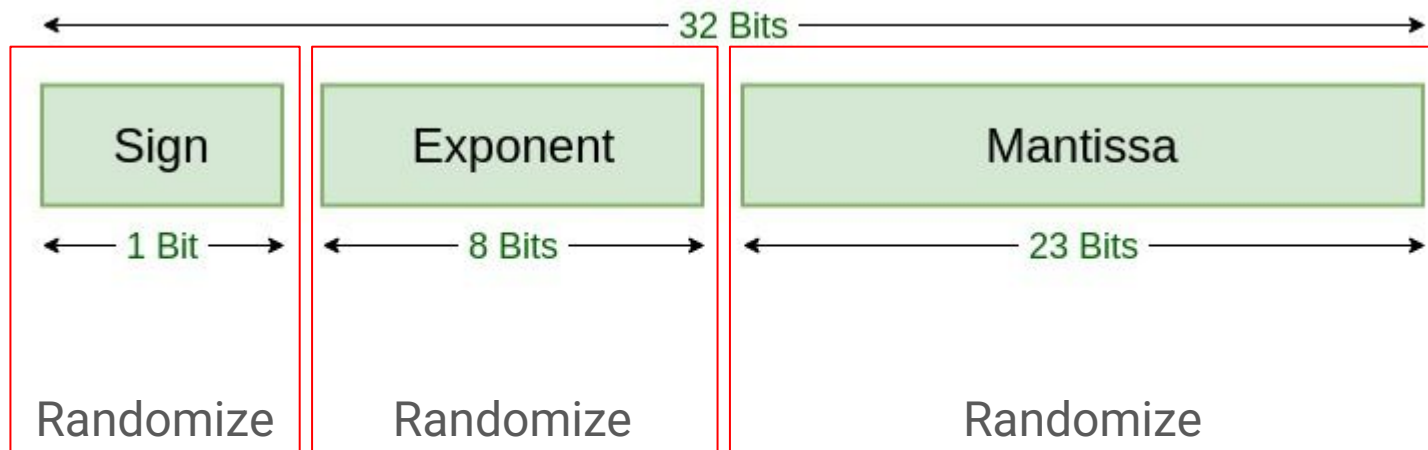
Experiment and Evaluation: Software Setup

System	Compilation	Measurement & Metrics
OS: GNU/Linux (Kernel 6.6.63)	Target ISA: -march=rv64gcv	Primary Metrics: Throughput (Elem/s), Execution Time (ms)
Compiler: RISC-V GCC 13.2.0	Optimization: -O2 (Standard), -O3	Timing Method: clock_gettime()
Execution Control: taskset (Bind to a single core)	Reproducibility: -static	Reliability: Multiple runs for each test (Ensures stable results)

Experiment and Evaluation: Data Generation

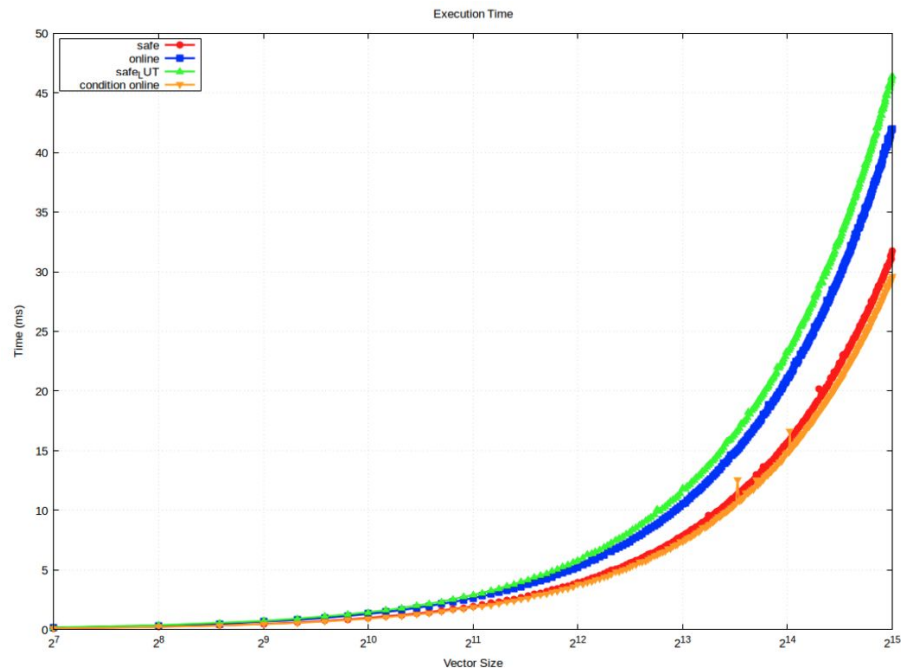
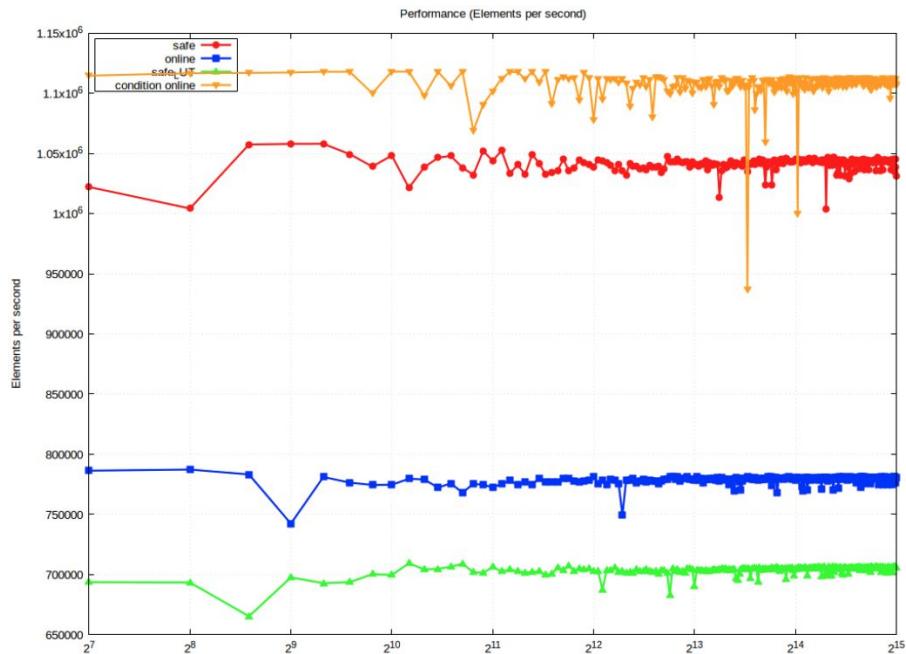
Our strategy generates 32-bit floating-point numbers by independently randomizing their constituent bit-fields, with a specific constraint to avoid special values

- **No underflow, no overflow, no subnormal number**



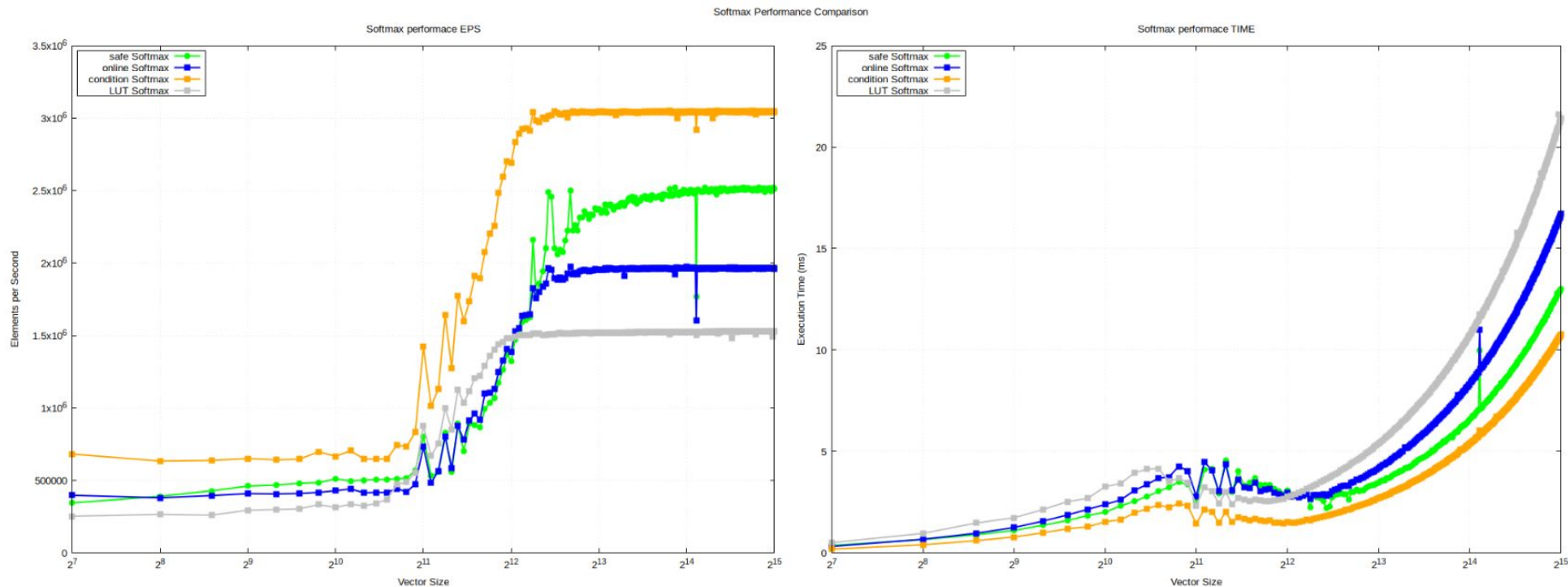
Experiment and Evaluation: Conditional Online Softmax

- RISC-V Platform



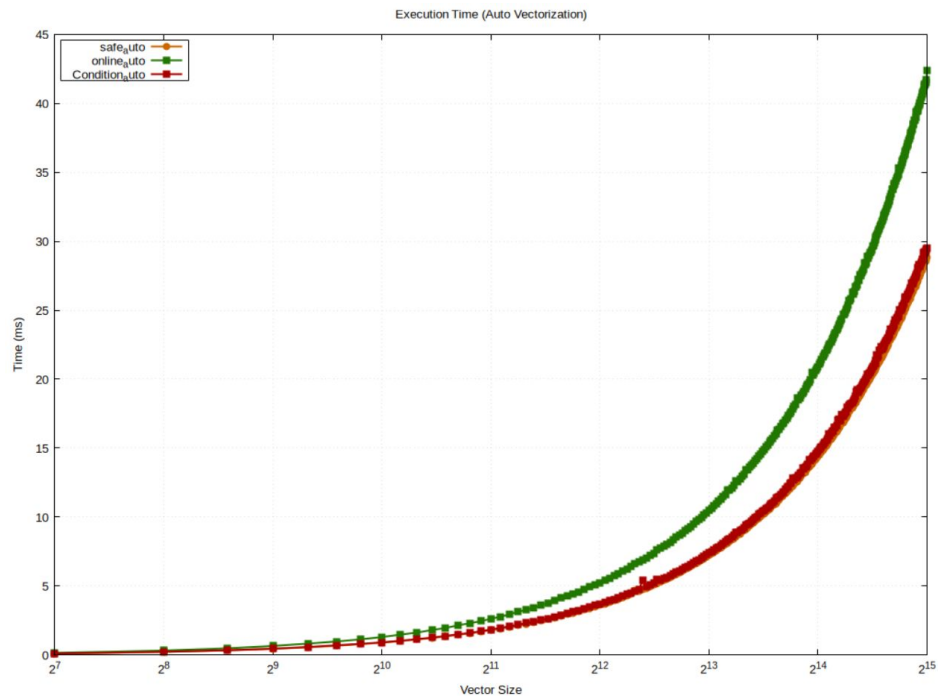
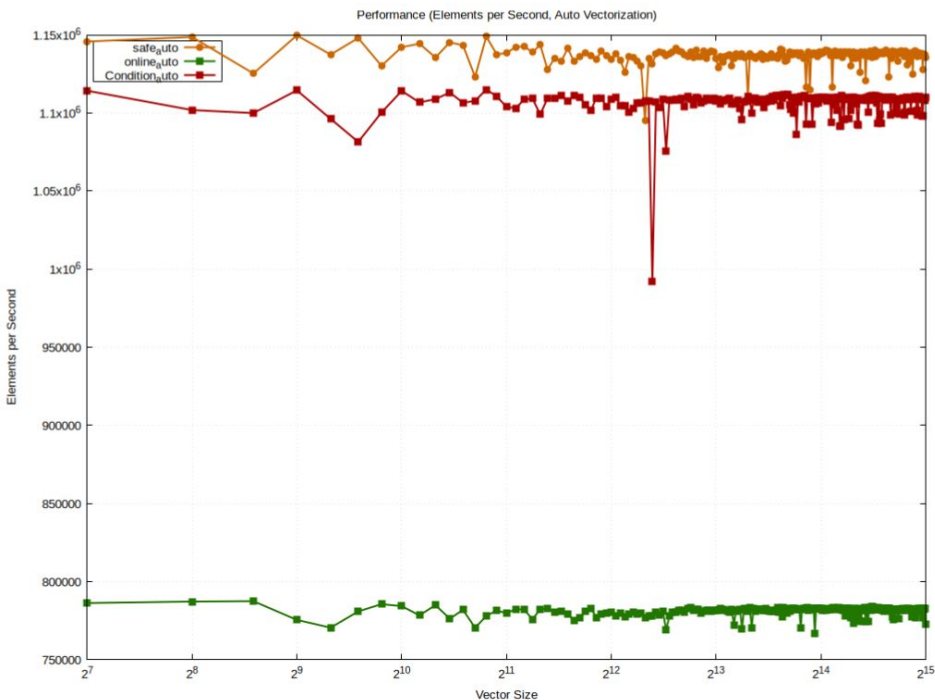
Experiment and Evaluation: Conditional Online Softmax

- x86 Platform



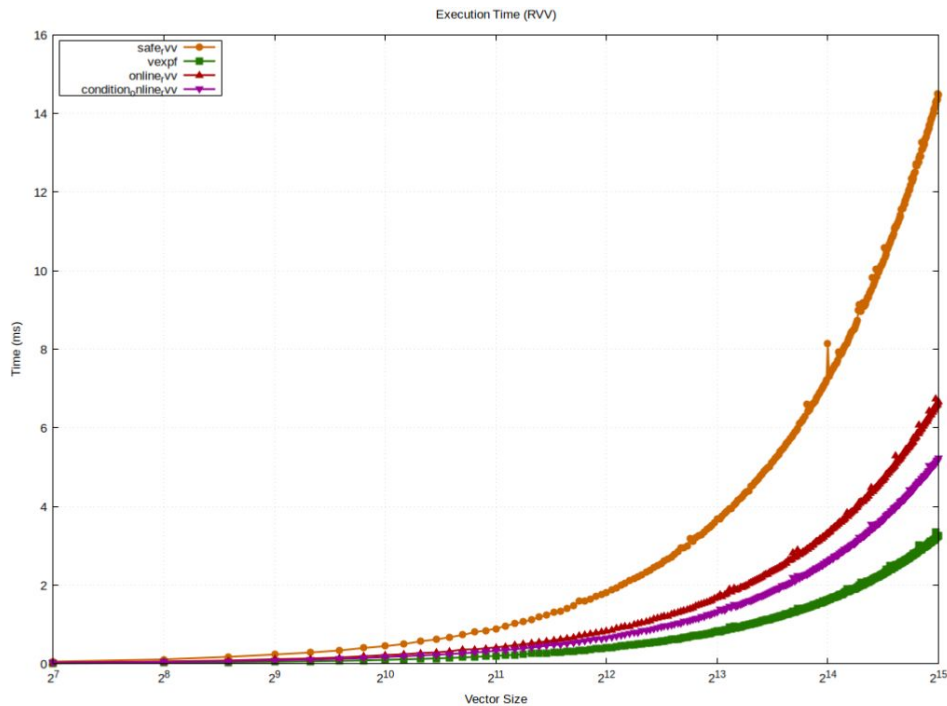
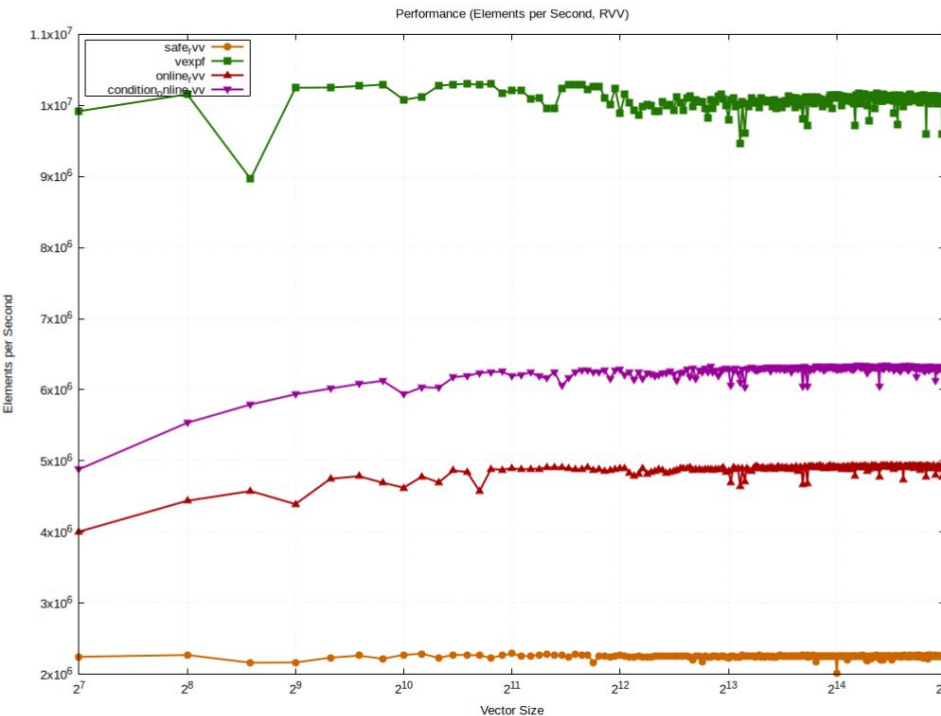
Experiment and Evaluation: Auto Vectorization

- RISC-V Platform



Experiment and Evaluation: Branch-free Conditional Online Softmax

- RISC-V Platform



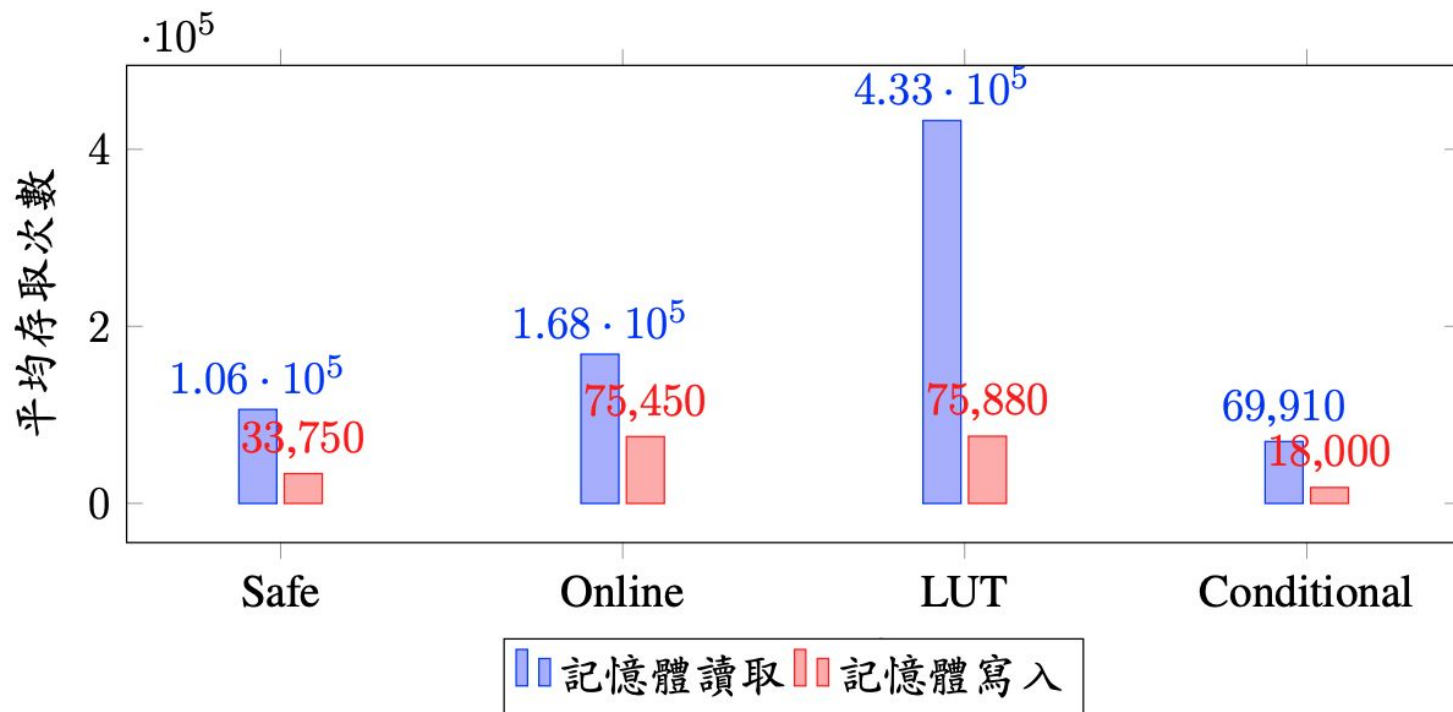
Experiment and Evaluation:

演算法	實現平台/方式	吞吐量 (Elem/s)	加速比*	記憶體訪問
Safe Softmax	Scalar (RISC-V)	1.04e6	1.00x	3R + 2W
Online Softmax	Scalar (RISC-V)	0.78e6	0.75x	2R + 1W
LUT Softmax	Scalar (RISC-V)	0.70e6	0.67x	2R + 1W
Conditional Online	Scalar (RISC-V)	1.11e6	1.07x	2R + 1W
Safe Softmax	Scalar (x86)	2.50e6	1.00x	3R + 2W
Online Softmax	Scalar (x86)	1.95e6	0.78x	2R + 1W
LUT Softmax	Scalar (x86)	1.5e6	0.6x	2R + 1W
Conditional Online	Scalar (x86)	3.10e7	1.24x	2R + 1W
Safe Softmax	Auto-Vectorized	2.20e6	1.00x	3R + 2W
Online Softmax	Manual RVV	4.90e6	2.22x	2R + 1W
Conditional Online	Manual RVV	6.4e6	2.9x	2R + 1W
Safe Softmax + vexpf	Manual RVV	1.01e7	4.59x	3R + 2W



Experiment and Evaluation: Memory Analysis

- x86 Platform



Conclusion and Future Work

Conclusion and Future Work

Conclusion:

- Algorithmic Innovation:
 - Proposed the Conditional Online Softmax algorithm with **speedup of 7%** on RISC-V and **decreased memory access with 37%** on x86 in scalar execution.
- High-Performance Implementation:
 - Achieved **speedup of up to 2.9x** on the RISC-V platform through manual vectorization and a custom vexpf library.

Future Work:

- Hardware/Software Co-design:
 - Develop a custom vector instruction for the Conditional Online update logic.
- Algorithm Integration:
 - Integrate the optimized Softmax kernel into state-of-the-art algorithms like FlashAttention.

Q&A



Thank You

