

Integrated Matrix Extension (IME)

Task Group Meeting

Guido Araujo
Jose Moreira

04/22/24

Agenda

- Revisiting what we have achieved so far
- Definition of gaps, agenda and working groups

What we have achieved so far

- Many presentations
- ISA modeling and analytical evaluation
 - IBM and Esperanto: Architectural options
- Technical discussions on the merits of each approach
- Analytical estimates for metrics
 - Andes and Unicamp: Computational intensity evaluation
 - Alibaba and Unicamp: Preliminary workload evaluation
 - Andes and Ventana: Register port analysis (uArch)
 - Unicamp: Packing cost evaluation*
 - Unicamp: POWER10 vector-matrix transfers*

Many presentations



SiFive



VENTANA

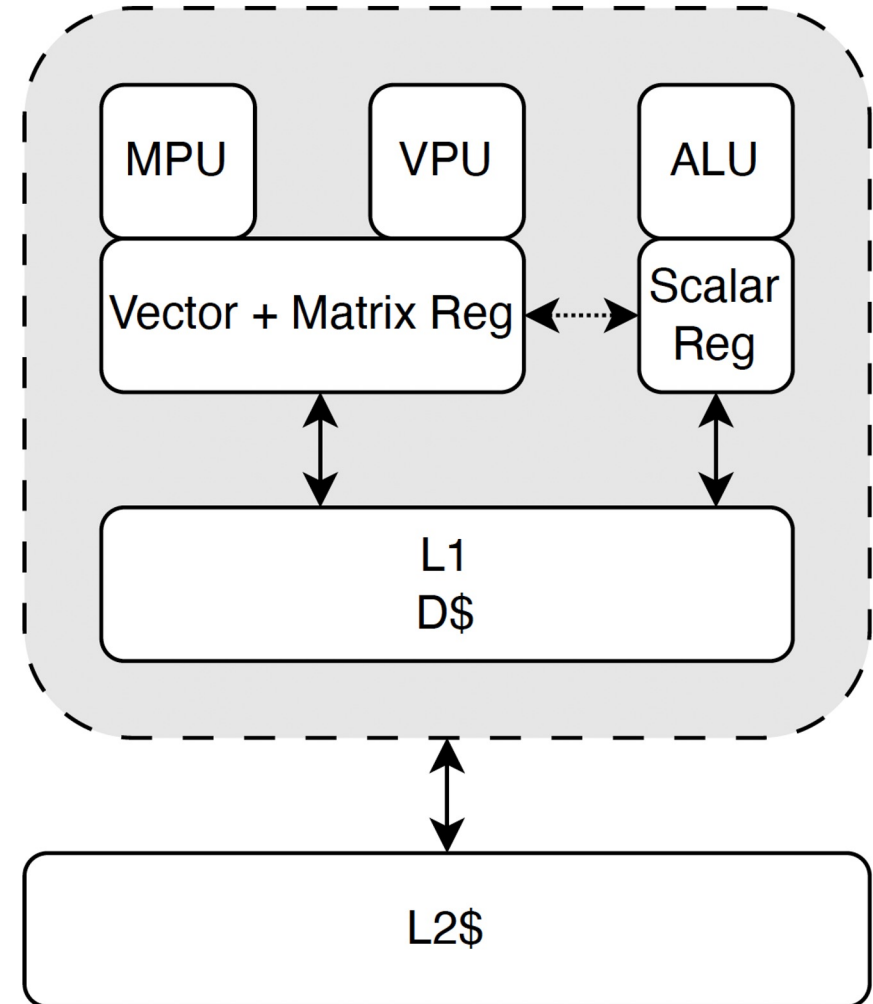


esperanto.ai

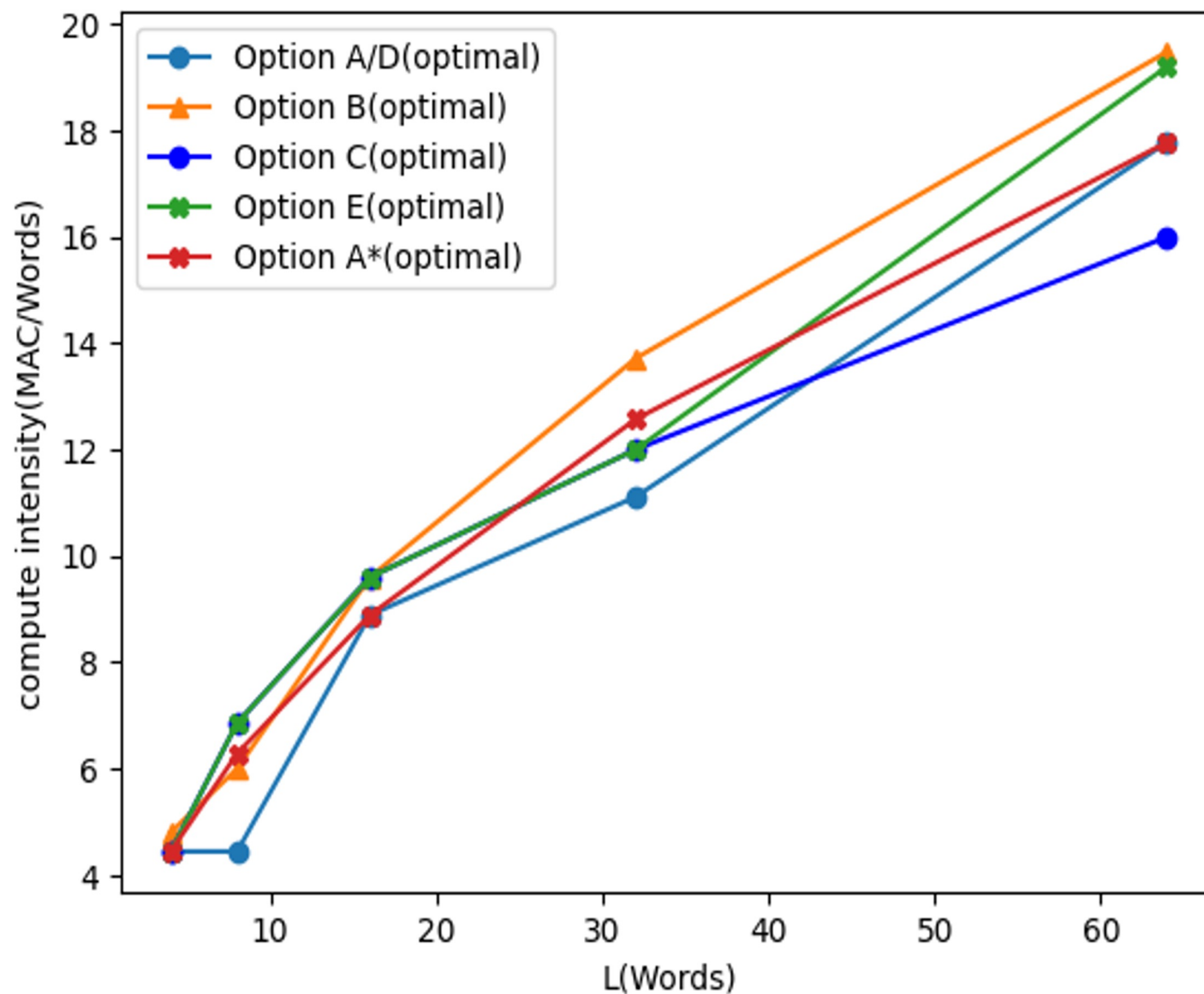
NEC



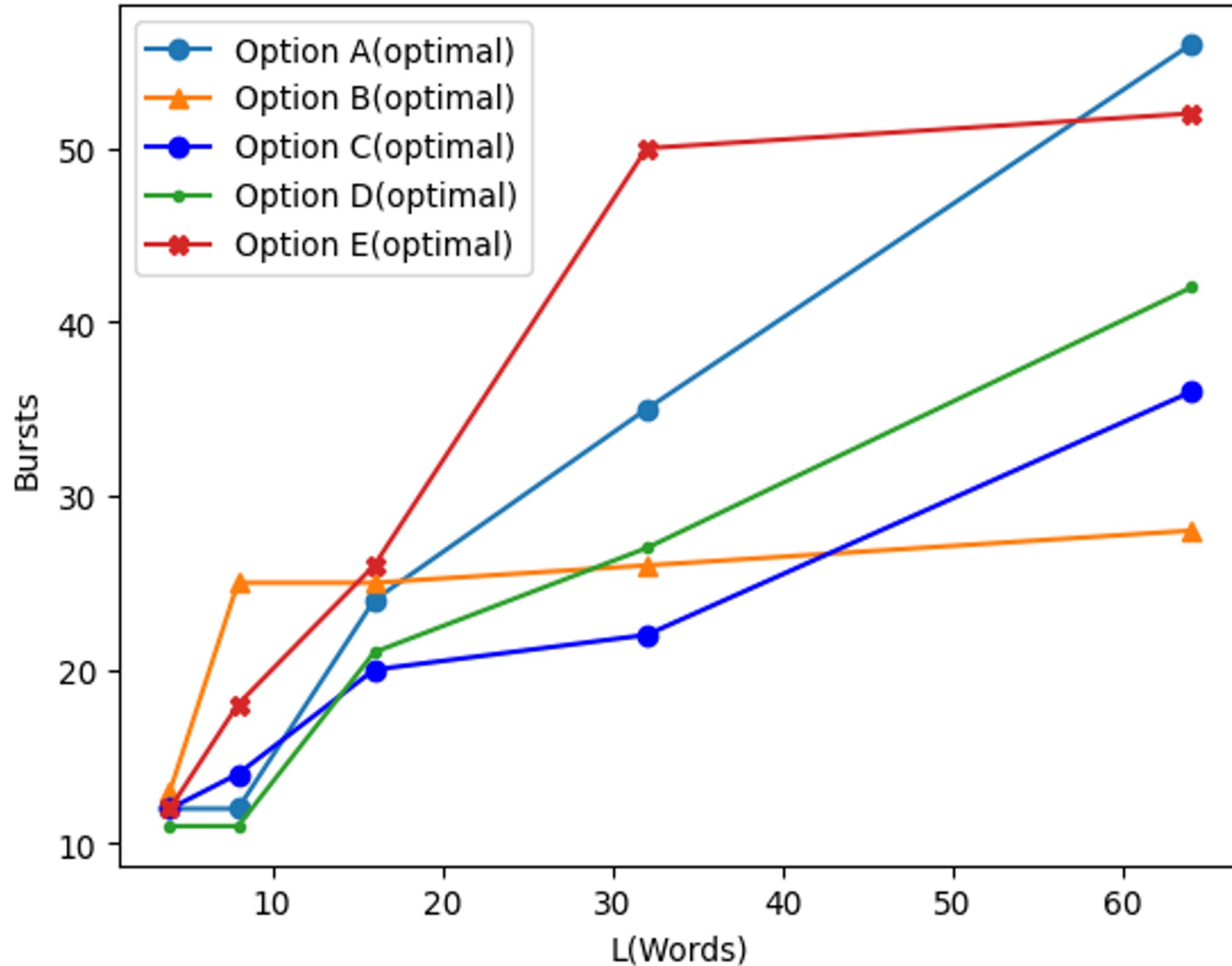
Integrated



Computational Intensity



Evaluating impact on packing*



Charter Criteria Discussion

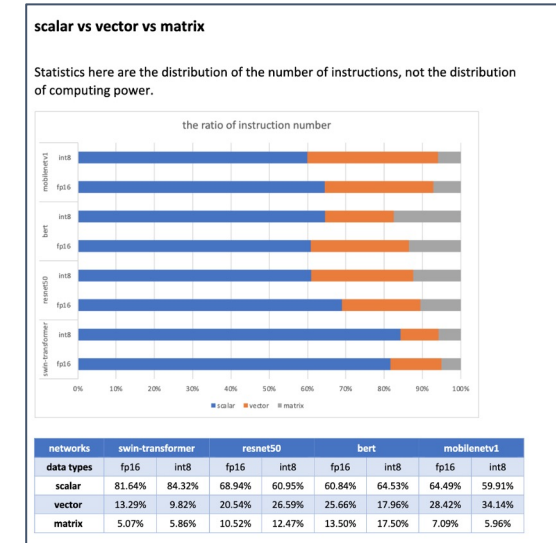
- Based on proposing Matrix-TG → 9 explicitly guide-lines
- Suggest to add Metrics for Performance/uArch Cost

No.	Guides	Option A	Option B	Option C	Option D	Option E
1	VLEN agnostic at binary level	Not Disclosed while $LI=\lambda^2$	To Be Discussed (seems support?)		Not Disclosed while $LI=\lambda^2$	(1) Source Compliant (2) Binary Compliant for vmul/fused ISA
2	Deterministic Result (FMAC rounding/ordering)	Shall support	Shall support	Shall support	Shall support	Support if finalized
3	Re-producible result with plain scalar/vector	Shall support	Shall support	Shall support	Shall support	Support (BIT TRUE test)
4	Near peak (~90%) performance for GEMM kernels is possible					GEMM kernels Near Peak U-rate
5	Higher (~2X) performance than vector					Over (>3X) enhancement than RVV
6	Maximization computation intensity for GEMM kernels					GEMM kernels Near Peak U-rate
7	Minimization additional architecture state	None new state	None new state	None new state	Not Support (New Streaming buffer for A/B)	None new state (not considering ZOB)
8	Live-migration with larger vector registers	Not Support for different lambda	To Be Discussed (seems support?)	Not Support for different vector element type	Not Support for different lambda	Under Working (AMM 2.0)
9	Proper Support for packing/reformat data	May Need Additional handling while $LI=\lambda^2$	Support (No additional interleaving/shuffle required)	Support (No additional interleaving/shuffle required)	May Need Additional handling while $LI=\lambda^2$	Support (No additional interleaving/shuffle required)
10	Metrics for Performance/uArch cost (Suggest to consider)	Feasible uArch cost(VRF R/W) for Specific MAC Performance	High uArch cost(VRF R/W) for MAC Performance [†]	feasible uArch cost(VRF R/W) for MAC Performance	Feasible uArch cost(VRF R/W) for MAC Performance	Feasible uArch cost(VRF R/W) for MAC Performance



Workloads and Benchmarking

- [Alibaba] Analysis of matrix and vector instruction distribution in ML Models



- [Unicamp] Analysis of vector and matrix transfers in Polybench (HPC benchmark)

	Name	Algorithm	vector = vector op vector		matrix = vector op matrix		vector = vector op matrix		matrix = matrix op matrix	
			lines	%time	lines	%time	lines	%time	lines	%time
3	2mm	$D := \alpha A^T B^T C + \beta A^T D$							2 mac: (83-89) (90-96) with scalar	
4	3mm	$G := (A^T B)^T C^T D$							3 mac: (79-85) (87-93) (95-101)	
5	atax	$y := (x^T A)^T A + y$					2 mac: (68-72) (68-75)			
6	bicg	$s := r^T A + s$ $q := A^T p + q$					2 mac: (77-85) (77-85)			
7	cholesky	Algorithm to find L, in which $A := LL^T$ (triangular application)	2mac: (70-71) (76-77)							
8	doligen	R matrix multiplications in cubic buffers							1 mac: (69-74) 1 mac: (75-76)	
9	gemm	$C := \alpha A^T A^T B + \beta A^T C$							1 mac: (77-83) with scalar	
10	gemver	$A := A + u1.v1 + u2.v2$ $x := \beta A^T A^T y + z$ $w := \alpha A^T A^T x$	1 add: (100-101)		2 mac: (92-94)		1 mac: (96-98) with scalar 1 mac: (103-105) with scalar		2 add (92-94)	
11	gesummv	$y := \alpha A^T A^T x + \beta A^T B^T x$			2 macs (76-84) (76-84) with scalar					
12	mvt	$x1 := x1 + A^T y1$ $x2 := x2 + A^T y2$					2 mac: (75-77) (79-81)			
13	symm	$C := \alpha A^T A^T B + \beta A^T C$					1 mac: (81-84) with scalar			
14	sy2k	$C := \alpha A^T A^T B + \alpha A^T B^T A + \beta A^T C$							2mac: (79-85) with scalar	
15	syk	$C := \alpha A^T A^T A + \beta A^T C$							1mac: (75-78) with scalar	

Agenda

- Revisiting what we have achieved so far
- Definition of gaps, agenda and working groups

Gaps and agenda (14 meetings)

- Today: Definition of working groups
- 2 Meetings: Consolidation of the feature table
- 1 Meetings: Definition of workloads and benchmarking
- 2 Meetings: Discussion and definition of two candidates
- 3 Meetings: Quantitative analysis
- 2 Meetings: Definition of the final candidate
- 4 Meetings: Discussion and writing of the specification

Working groups

Groups	Coordinator	Members
Option A and A*		
Option B		
Option C		
Option D		
Option E		
Workloads and Benchmarking	Guido	

Quantitative Approach

Polybench
Convbench
Other



Options'
Simulators

Metrics

- Option A

- Option B

- Option C

- Option D

- Option E