

# Integrated Matrix Extension (IME)

Task Group Meeting

**Guido Araujo**  
**Jose Moreira**

03/11/24

# Agenda

- Chair/Vice-chair selection update
- Moving forward on qualitative analysis
  - Computational Intensity
  - Locality evaluation
- Presentation
  - Matrix Tile Extension: Portable ISA For Vector-Integrated Matrix Unit
  - Erich Focht (NEC) and Marc Casas (BSC)

# Agenda

- Chair/Vice-chair selection update
- Moving forward on qualitative analysis
  - Computational Intensity
  - Locality evaluation
- Presentation
  - Matrix Tile Extension: Portable ISA For Vector-Integrated Matrix Unit
  - Erich Focht (NEC) and Marc Casas (BSC)

## Chair/Vice-chair selection update

- Call started: February 23rd.
- Call Closure: Friday, March 8th.
- Completion of Voting/Selection: Friday, March 29th

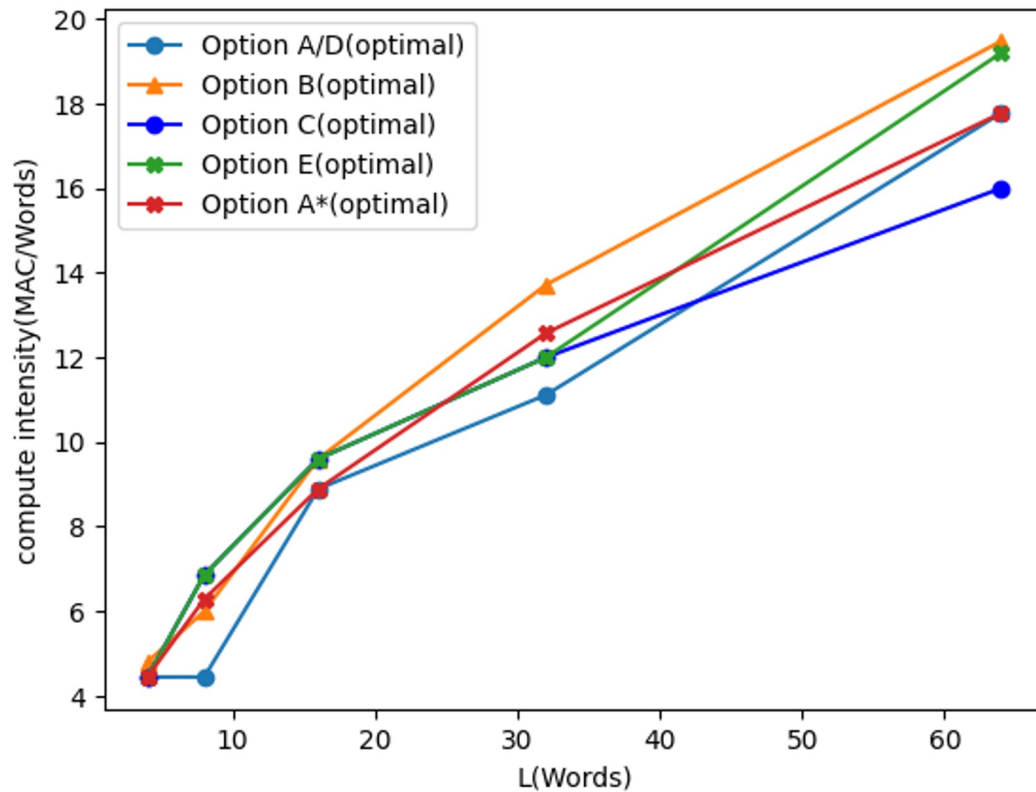
# Agenda

- Chair/Vice-chair selection update
- **Moving forward on qualitative analysis**
  - Computational Intensity
  - Locality evaluation
- **Presentation**
  - Matrix Tile Extension: Portable ISA For Vector-Integrated Matrix Unit
  - Erich Focht (NEC) and Marc Casas (BSC)

# Agenda

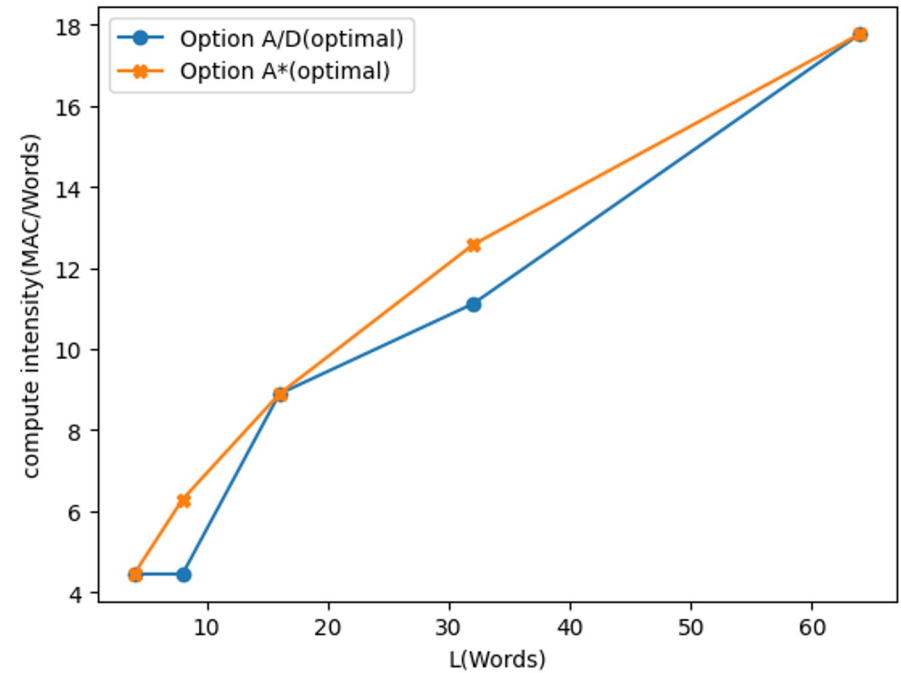
- Chair/Vice-chair selection update
- Moving forward on qualitative analysis
  - Computational Intensity
  - Locality evaluation
- Presentation
  - Matrix Tile Extension: Portable ISA For Vector-Integrated Matrix Unit
  - Erich Focht (NEC) and Marc Casas (BSC)

# Computational Intensity (Adding A\*)



Thanks

CN-Ke  
BING



# 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   | <b>VLEN agnostic</b> at binary level                             | Not Disclosed while $L! = \lambda^2$                      | To Be Discussed (seems support?)                      |   | Not Disclosed while $L! = \lambda^2$                | (1) Source Compliant<br>(2) Binary Compliant for vmul/fused ISA |
| 2   | <b>Deterministic Result</b> (FMAC rounding/ordering)             | Shall support   | Shall support   | Shall support   | Shall support                                       | Support if finalized  |
| 3   | <b>Re-producible result</b> with plain scalar/vector             | Shall support   | Shall support   | Shall support   | Shall support                                       | Support (BIT TRUE test)   |
| 4   | <b>Near peak (~90%) performance</b> for GEMM kernels is possible |   |   |   |   | GEMM kernels Near Peak U-rate                                   |
| 5   | <b>Higher (~2X) performance</b> than vector                      |   |   |   |   | Over (>3X) enhancement than RVV                                 |
| 6   | Maximization <b>computation intensity</b> for GEMM kernels       |   |   |   |   | GEMM kernels Near Peak U-rate                                   |
| 7   | Minimization additional <b>architecture state</b>                | None new state  | None new state  | None new state  | Not Support (New Streaming buffer for A/B)          | None new state (not considering ZOB)                            |
| 8   | <b>Live-migration</b> 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 <b>packing/reformat data</b>                  | May Need Additional handling while $L! = \lambda^2$       | Support (No additional interleaving/shuffle required) | Support (No additional interleaving/shuffle required) | May Need Additional handling while $L! = \lambda^2$ | Support (No additional interleaving/shuffle required)           |
| 10  | <b>Metrics for Performance/uArch cost</b> (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                |

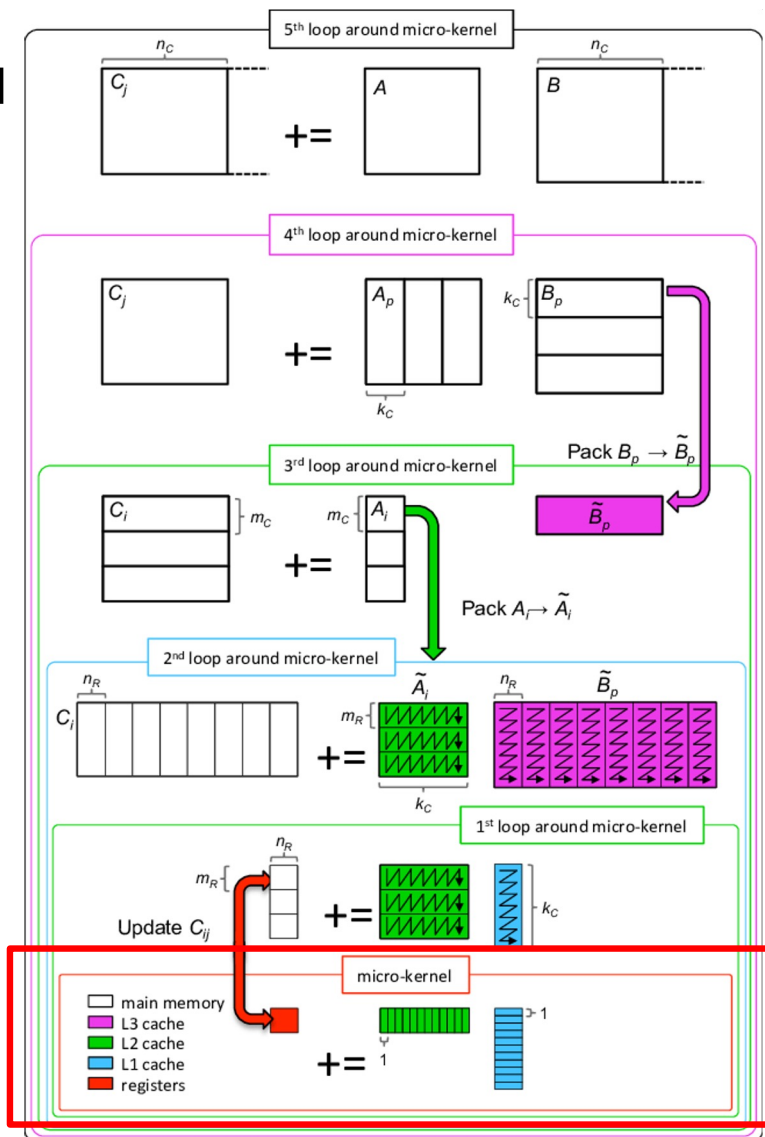




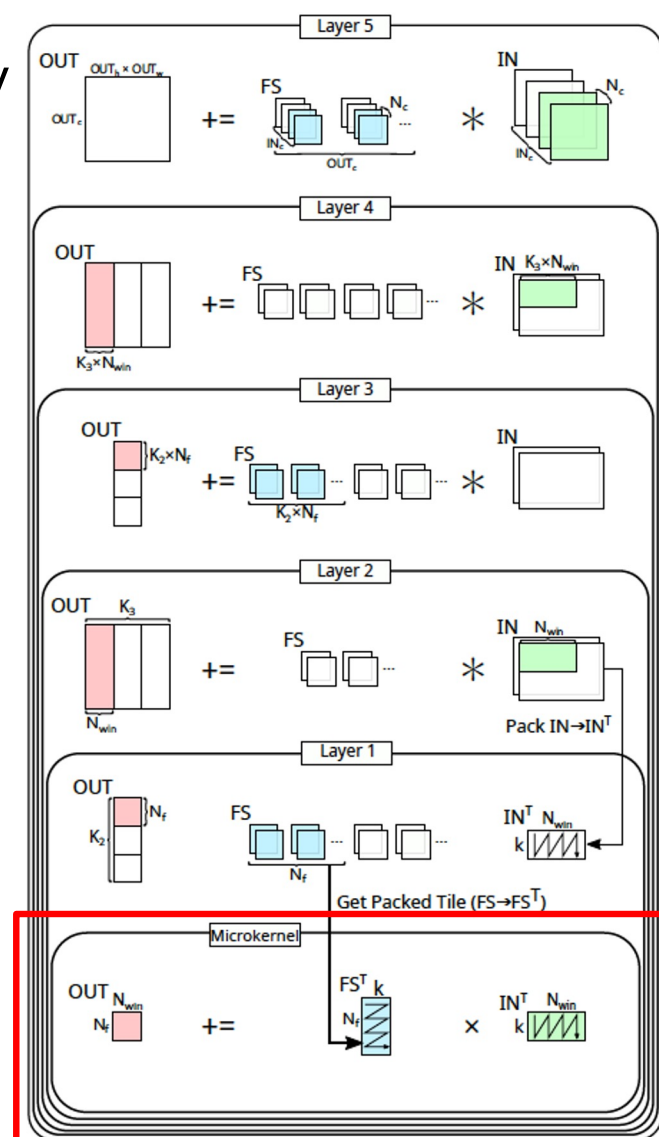
# Agenda

- Chair/Vice-chair selection update
- Moving forward on qualitative analysis
  - Computational Intensity
  - Locality evaluation
- Presentation
  - Matrix Tile Extension: Portable ISA For Vector-Integrated Matrix Unit
  - Erich Focht (NEC) and Marc Casas (BSC)

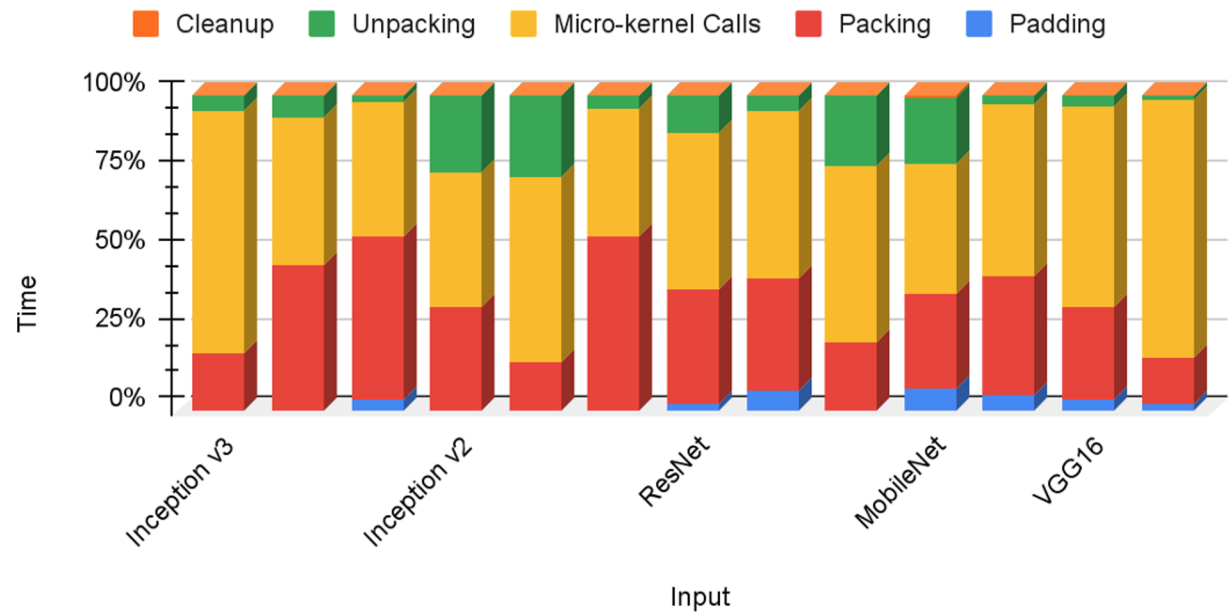
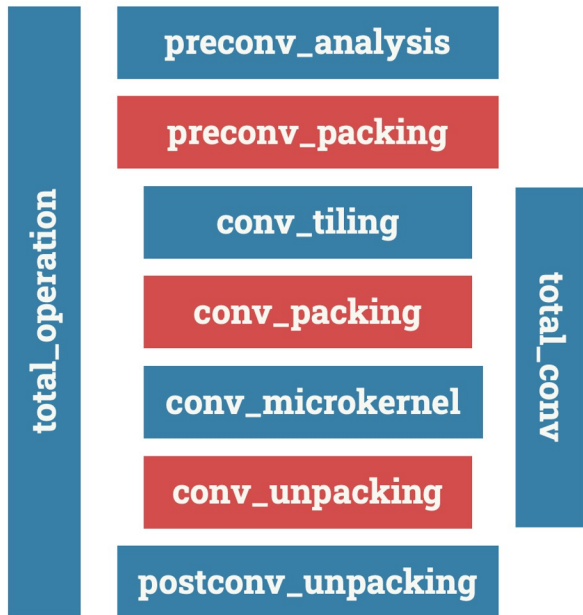
# GEMM



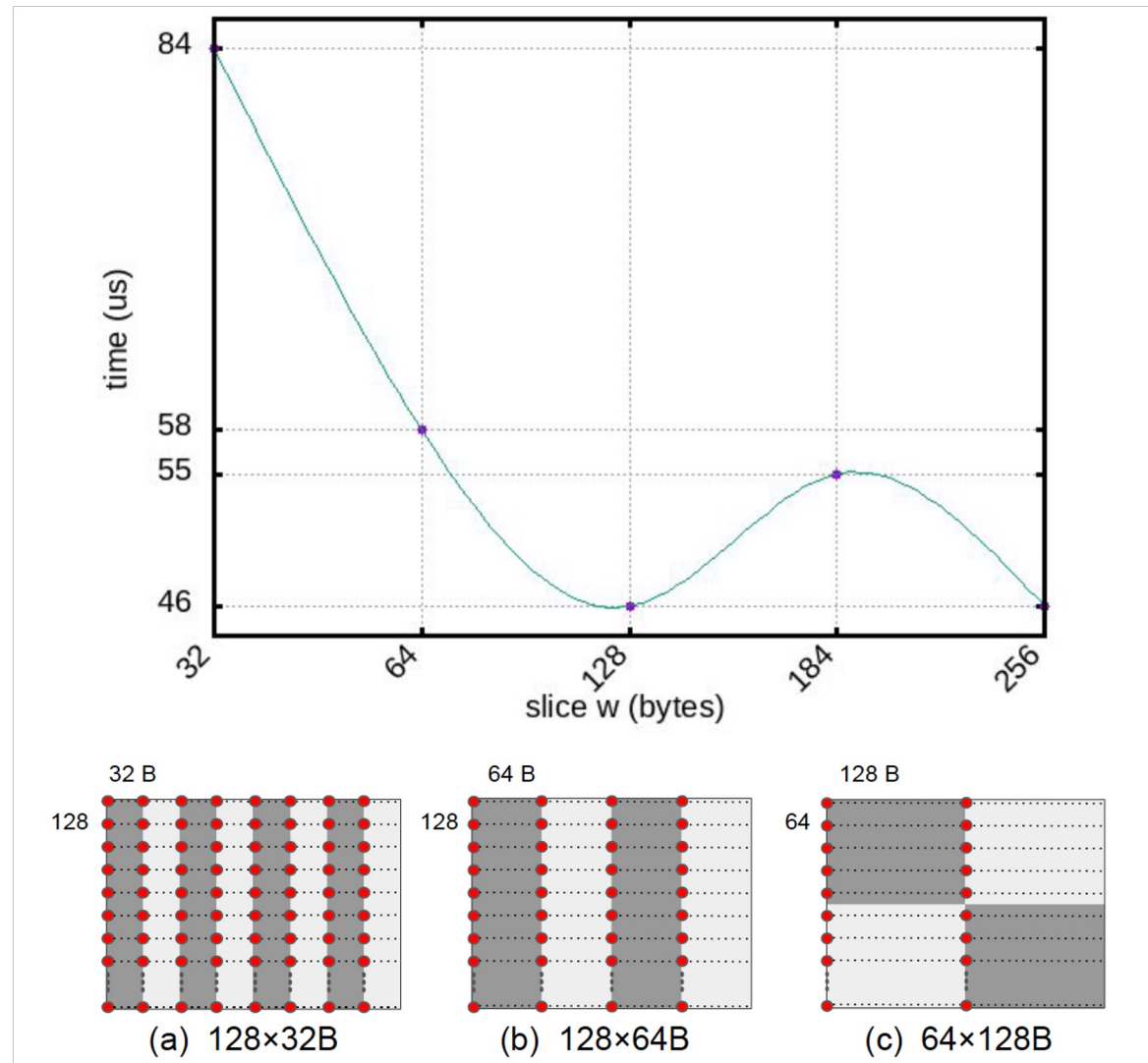
# CONV



# Locality is important for packing

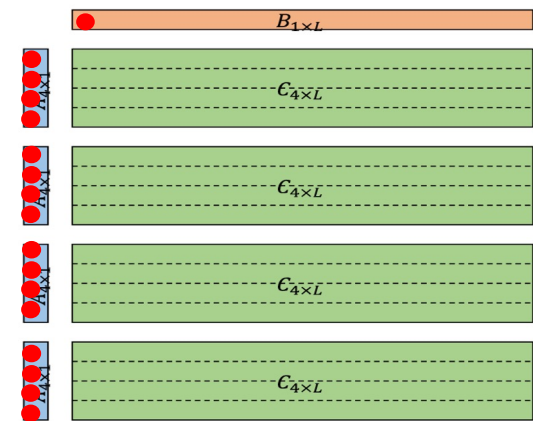
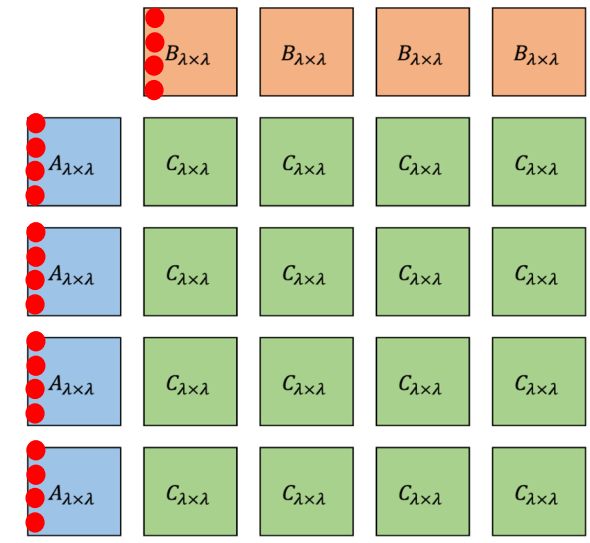


# Tiling and memory burst



# Evaluate impact on packing

- Option A
  - One matrix per register vector $5\lambda$
- Option B
  - One matrix in multiple register vectors $17$



# Agenda

- Chair/Vice-chair selection update
- Moving forward on qualitative analysis
  - Computational Intensity
  - Locality evaluation
- **Presentation**
  - Matrix Tile Extension: Portable ISA For Vector-Integrated Matrix Unit
  - Erich Focht (NEC) and Marc Casas (BSC)