# Integrated Matrix Extension (IME)

TG Kick-off Meeting

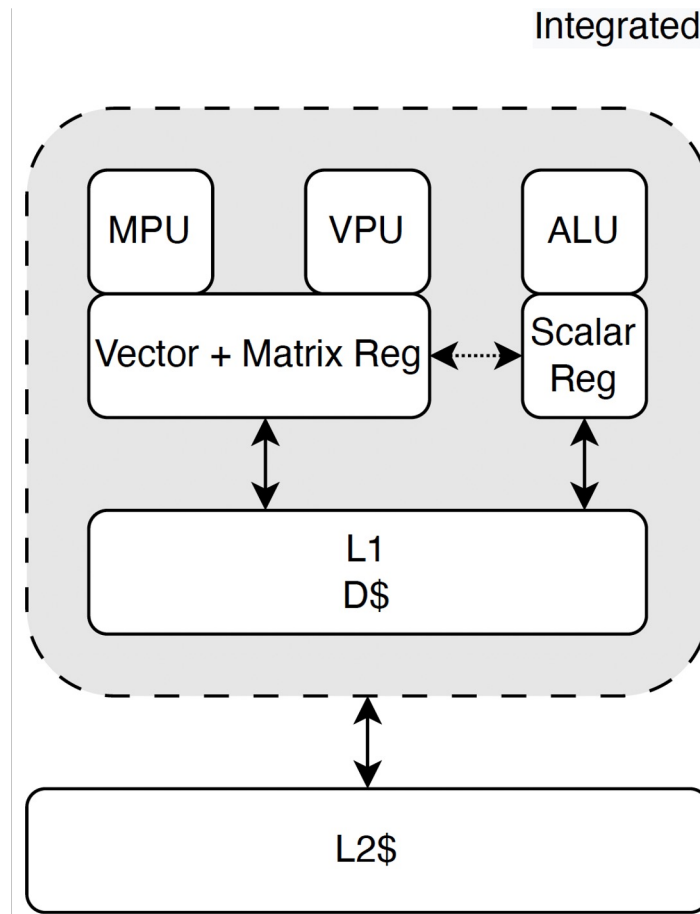**Guido Araujo**
**Jose Moreira**

02/12/24

# Agenda

- IME preliminary proposals

- Qualitative vs Quantitative approaches
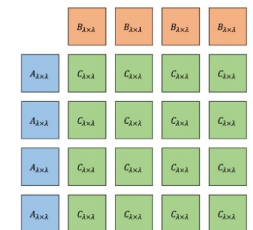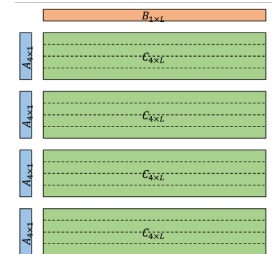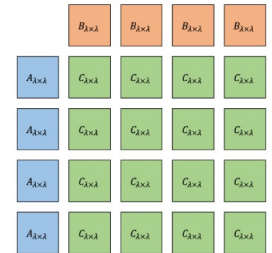
- Metrics and workloads

# Agenda

- IME preliminary proposals

- Qualitative vs Quantitative approaches

- Metrics and workloads
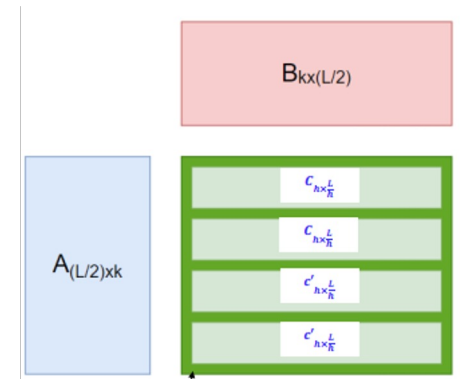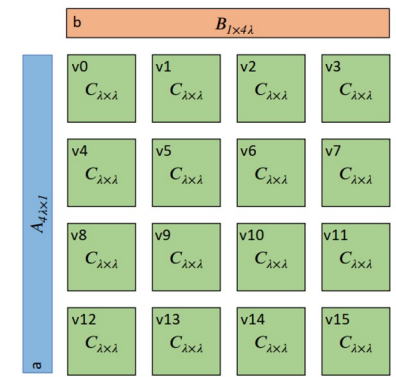
# IME - Integrated

# IME preliminary proposals

- Option A
  - One matrix per register vector

- Option B
  - One matrix in multiple register vectors

- Option C
  - Multiple matrices in one register vector

# IME preliminary proposals (cont.)

- Option D
  - Streaming buffers for the input matrices



- Option E
  - Variable matrix representation

# Agenda

- IME preliminary proposals

- <span style="color:red">Qualitative vs Quantitative approaches</span>

- Metrics and workloads

# Qualitative vs. Quantitative approaches

- Qualitative
    - ISA analytical evaluation
    - Analytical back-of-envelope estimates for metrics
    - Technical discussions on the merits of each approach
    - Faster progress, but ad-hoc

- Quantitative
    - Generate workloads
    - Run QEMU model, provided by proponent (open)
    - Run traces on simulator, provided by proponent (closed)
    - Measure metrics and evaluate
    - Slower progress, but precise
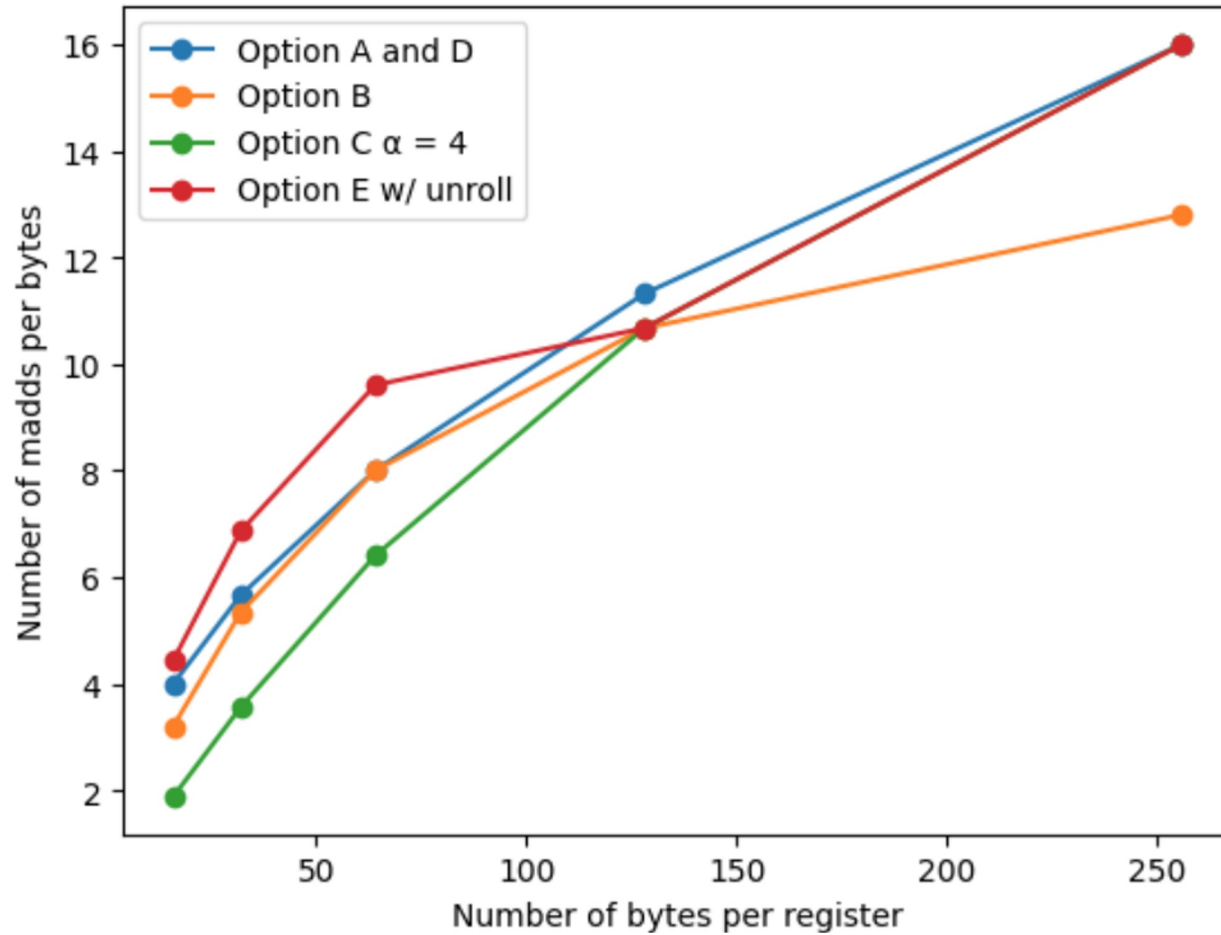
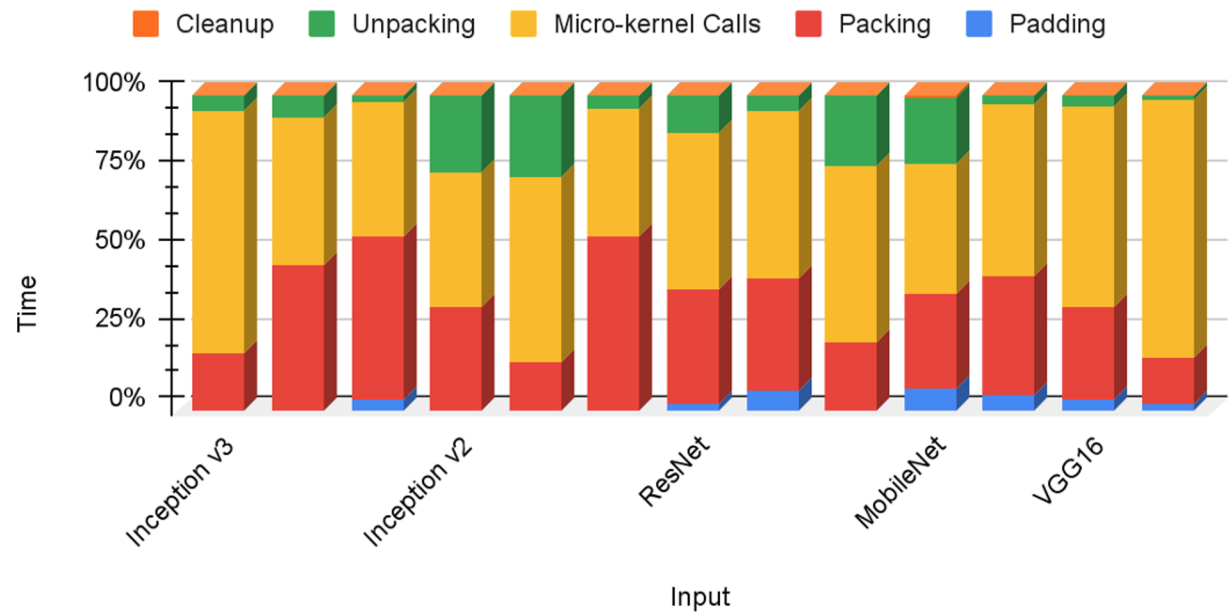# Quantitative Approach

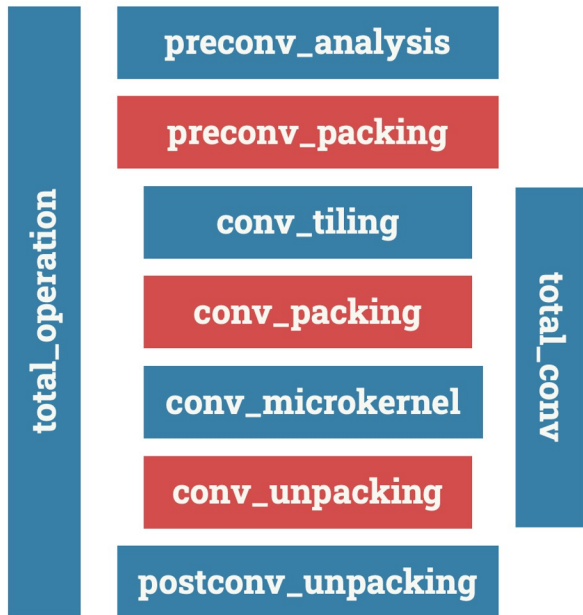| | Polybench Convbench Mobile | QEMU | Options' Simulators | Metrics |
|---|---|---|---|---|
| Option A | | | | |
| Option B | | | | |
| Option C | | | | |
| Option D | | | | |
| Option E | | | | |

# Qualitative approach

# Quantitative approach

# Agenda

- IME preliminary proposals

- Qualitative vs Quantitative approaches

- Metrics and workloads

# Performance metrics

- Computational intensity

- GEMM execution time

- CONV execution time

- ISA integration with RVV

- Cache misses

- Your favorite metric: TBD

# Domain workloads

- ML: Convbench

- HPC: Polybench

- Mobile: TBD

- Your favorite domain: TBD

# Thanks!!

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

02/12/24