# Update on Ara

23/11/2022

**Matteo Perotti**

**Matheus Cavalcante**

**Professor Luca Benini**

**Integrated Systems Laboratory**

**ETH Zürich**
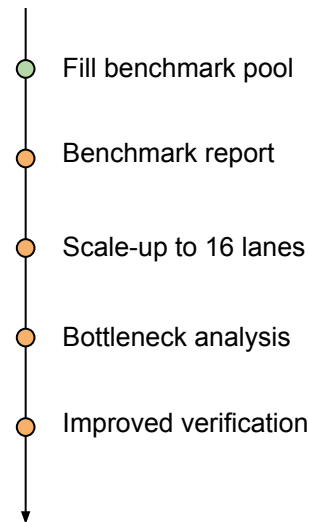
# Summary

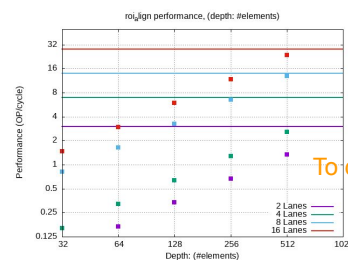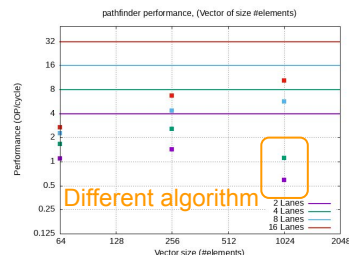- **Software**
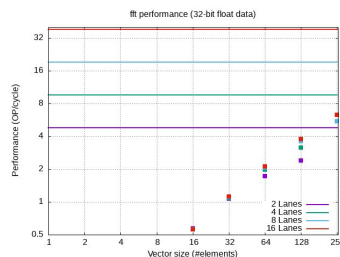  - New benchmarks upstream
  - Stall analysis

- **Hardware (RTL + Backend)**
  - 8-lanes trials
  - Mask instructions support
  - Multi-Core

Fill benchmark pool

Benchmark report

Scale-up to 16 lanes

Bottleneck analysis

Improved verification

# Longer vectors perform better

# Benchmarks

- Vector intrinsics

- How architectural choices reflect on programming model

- Pitfalls and coding guidelines



Naive intrinsics



Optimized intrinsics



Optimized ASM

# Benchmarks

- Vector intrinsics

- How architectural choices reflect on programming model

- Pitfalls and coding guidelines



Naive intrinsics

Optimized intrinsics

Optimized ASM

# Benchmarks

- Vector intrinsics

- How architectural choices reflect on programming model

- Pitfalls and coding guidelines



Naive intrinsics

Optimized intrinsics

Optimized ASM

# Scalar core problem?

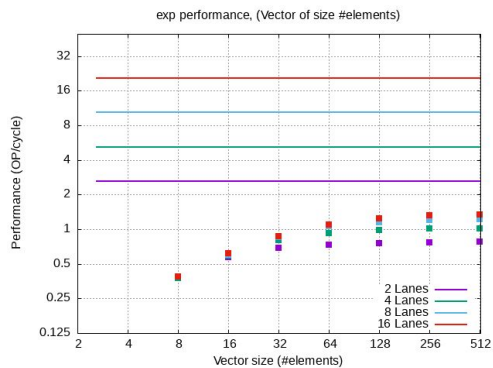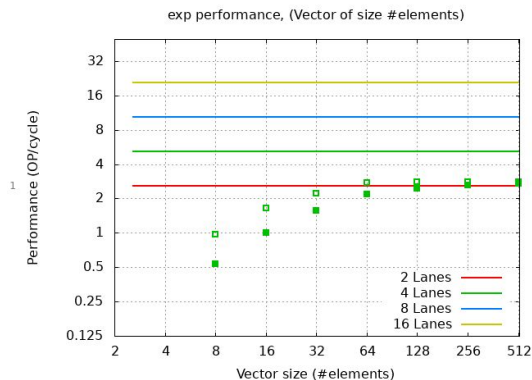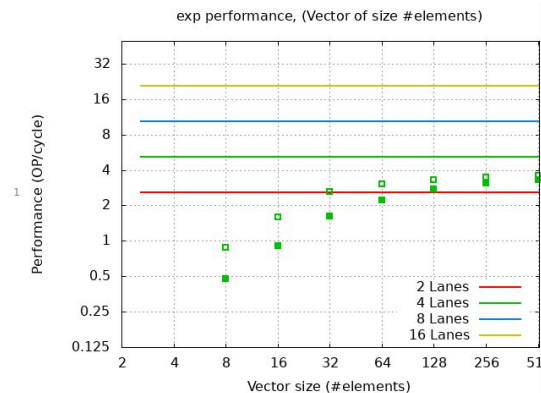**Issue rate limitation only for real system!**



□ Ideal Dispatcher + Ara
■ CVA6            + Ara

- **Ideal dispatcher analysis + Tune miss rate/penalty**
- CVA6 + Scalar Mem Sys bottlenecks the gains
  - Many in-Ara improvements are hidden by CVA6
- **Better scalar core** does **not** fully **solve** the **problems**
- Should **repeat analysis after internal optimization** with ideal dispatcher (barber's pole, new hazard handling engine)

# Ara problem?



D     S     L     OP     V     OP     FPU     RES     V
            S     REQ     R     Q              Q       R
                          F                            F

# Ara problem?

V Instruction

D    S    L    OP    V    OP    FPU    RES    V
          S    REQ   R    Q            Q      R
                     F                        F

# Ara problem?

Stall back-pressure

1. Lower LMUL?

D   S   L   OP   V   OP   FPU   RES   V
        S   REQ  R   Q          Q     R
                 F                     F

# Ara problem?

Stall back-pressure

D  S  L
       S

OP
REQ

V
R
F

OP
Q

FPU  RES
      Q

V
R
F

1. Downstream full queues?
2. Slide stall?
3. WAR, WAW no source?
4. Exception check?
5. Feedback answer?
6. Back pressure from next stage?

*__Crucial stalls for short vectors!__*

# Ara problem?

Stall **without** back-pressure for != units

1. Data hazards?
2. Arbitration?
3. Bank conflict?

# Ara problem?

Stall **without** back-pressure for != units



D    S    L
S    OP
REQ    V
R
F    OP
Q    FPU    RES
Q    V
R
F

From here on, the vector length help hide stalls from behind!

# Ara problem?

Stall **without** back-pressure for != units

Result

D  S  L
      S

OP
REQ

V
R
F

OP
Q

FPU  RES
      Q

V
R
F

1. Bank conflict?

# Ara problem?

**Issue rate limitation only for real system!**



jacobi2d performance, (matrices of size #elements x #elements)

jacobi2d performance, (matrices of size #elements x #elements)

☐ Ideal Dispatcher + Ara
■ CVA6             + Ara

- CVA6 bottlenecks the gains
  - Many in-Ara improvements are hidden by CVA6
- **Do we need renaming? WAW-WAR analysis**
- **Why short-vectors are a problem? Stall analysis**
- **Diminishing returns** to improve short vector perf (larger buffers)
- Should **repeat analysis after internal optimization** with ideal dispatcher (barber's pole, new hazard handling engine)
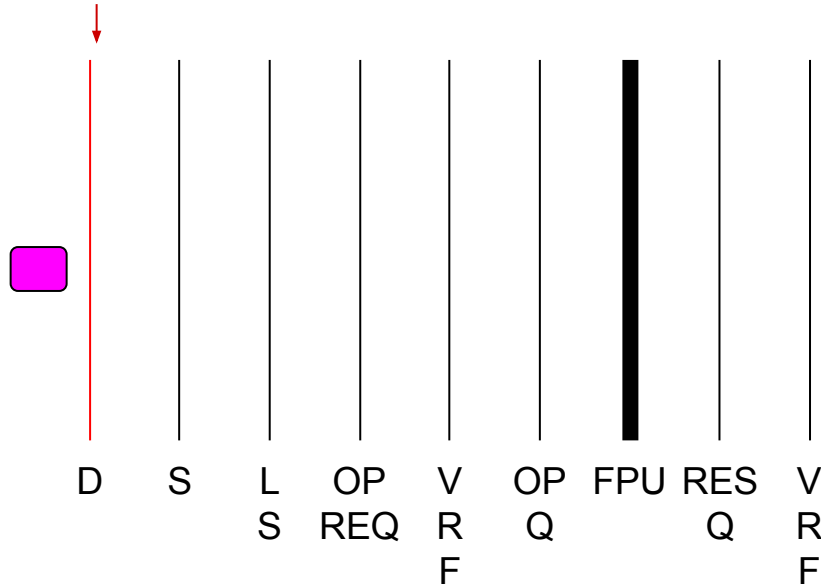
# Scale-up problem

- **HW**:
  8 lanes WIP for full closure
  16 lanes is almost infeasable:
    - `With    SLDU - Infeasable`
    - `Without SLDU - WNS: -400 ps`

- **SW**:
  Efficient use of the resources:
    - Longer vectors
    - Hard to partition a problem!

# Vector Multi-Core

✓  **Partition** the problem

✓  **Utilization is higher** since vectors are "longer"

✓  **Scale-up** in terms of **FPUs**

➤  **Should repeat the analysis after optimization**
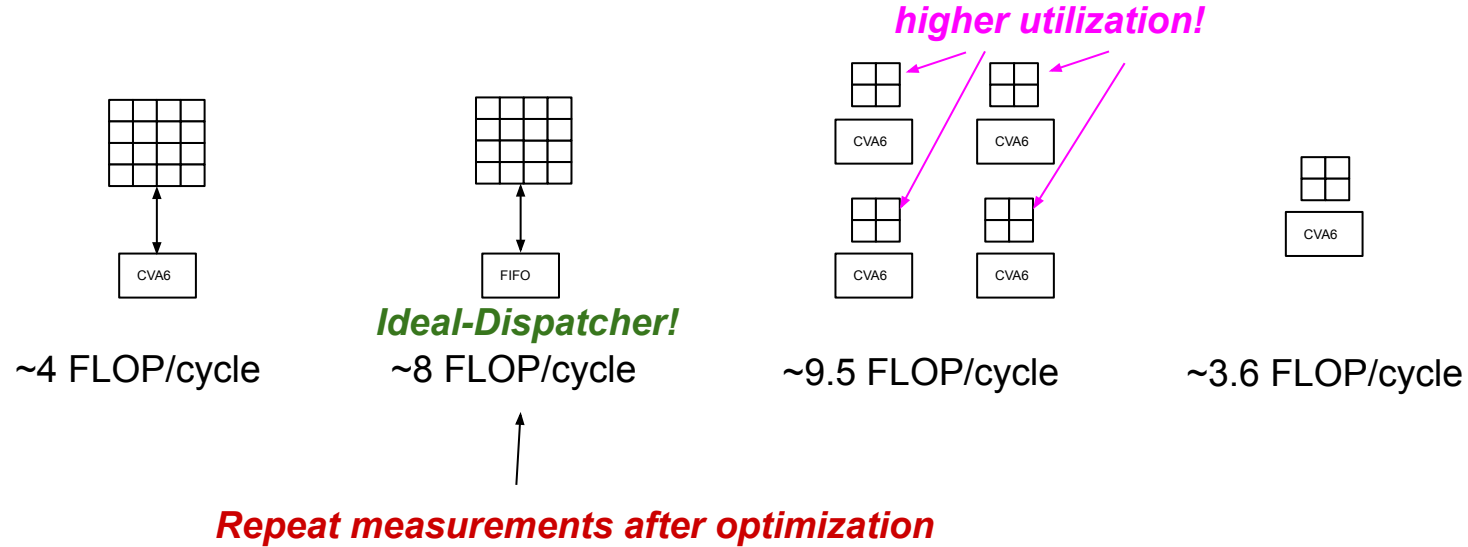
✗  **Increased CVA6 traffic to upper level of memory**

# Vector Multi-Core (8x8 `fmatmul`)

*higher utilization!*

~4 FLOP/cycle          ~8 FLOP/cycle          ~9.5 FLOP/cycle          ~3.6 FLOP/cycle

*Ideal-Dispatcher!*

*Repeat measurements after optimization*

Power and efficiency measurements - Many Ara systems in an SoC wrapper
- Memory left outside
- Interconnect left outside

# PRs - Compliance



➤ Vector complex shuffling
✓ Mask-reg instructions    (1)
✓ Mask-reg instructions    (2)
✓ Fixed-Point support      (1)
✓ Fixed-Point support      (2)
✗ Segment memory ops

✗ strncpy, strncmp
✓ AWB

# PRs - Benchmarks

[HW, SW] Add Power Analysis feature ✕
#179 opened 4 days ago by mp-17 • Draft • 3 tasks

add kernels: gemv, spmv, conjugate gradient ✕
#175 opened 12 days ago by husterZC • Changes requested

[SW] Add lavaMD benchmark ✓
#166 opened 20 days ago by mp-17 • Draft • 1 of 3 tasks

[HW] Draft PR for Implementing Ara on FPGA ✓
#146 opened on Sep 19 by hossein1387 • Draft • 3 tasks

[HW] cva6: Increase AXI data width ✓
#91 opened on Nov 18, 2021 by niwis • 2 of 3 tasks

[HW] Add support for the single-laned configuration ✕
#75 opened on Sep 28, 2021 by suehtamacv • Draft • 3 tasks

✓ `gemv, spmv, conjugate gradient`
✓ `lavaMD`

✕ `strncpy, strncmp`
✓ AWB

# Further

- **Software**
  - Stall analysis
  - WAW / WAR profiling

- **Hardware (RTL + Backend)**
  - Compliance + Verification
  - Scale up to 16 lanes

Fill benchmark pool

Benchmark report

Scale-up to 16 lanes

Bottleneck analysis

Improved verification