

EE382N (20): Computer Architecture - Parallelism and Locality  
Fall 2011

## **Lectures 8 & 9 – Parallelism in Hardware**

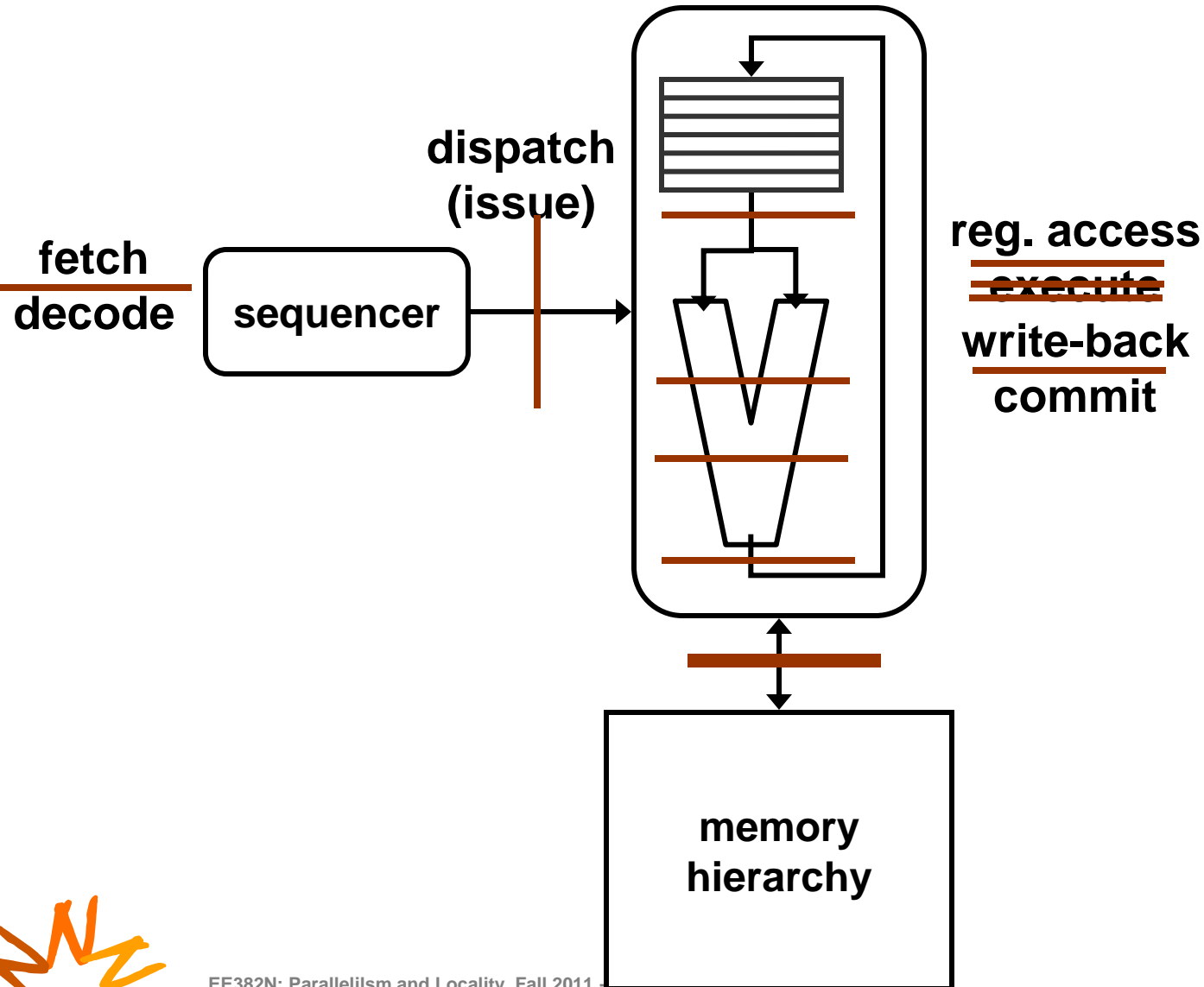
Mattan Erez



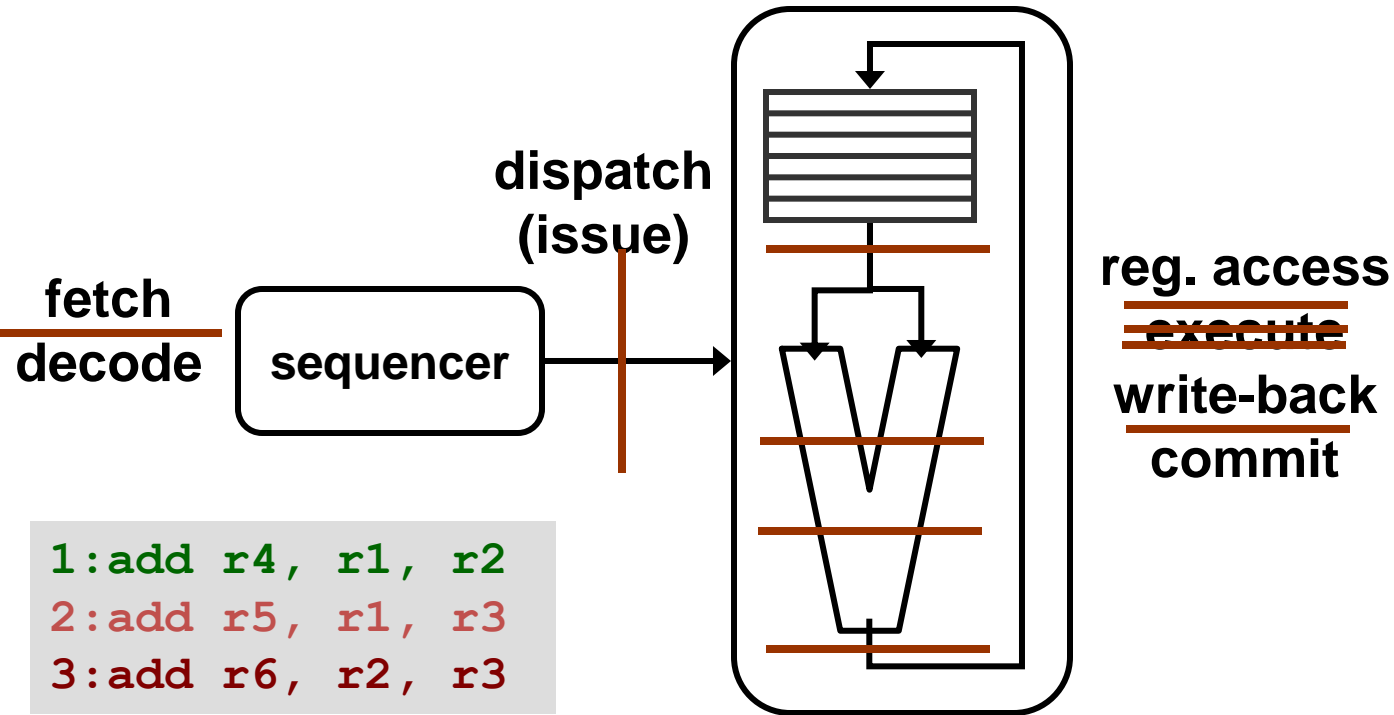
The University of Texas at Austin



# Simplified view of a pipelined processor



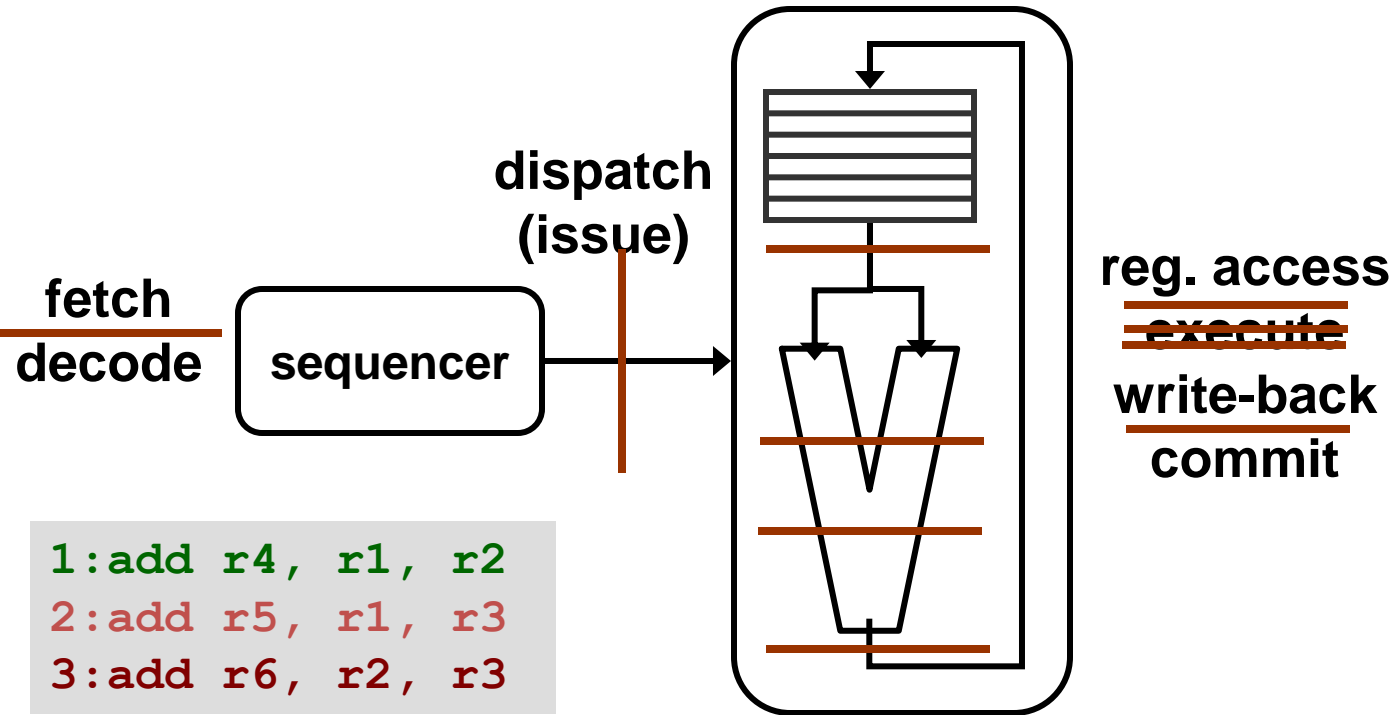
# Simplified view of a pipelined processor



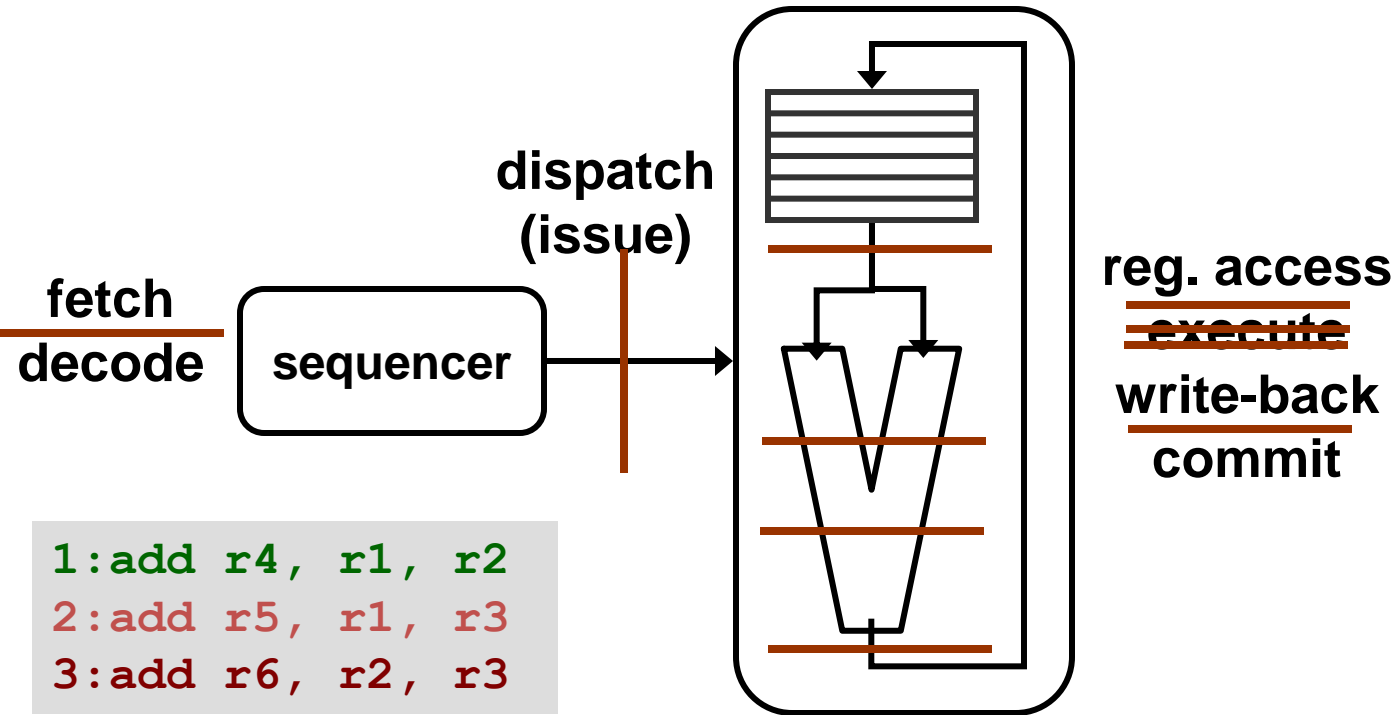
|   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|
| 1 | F | D | I | R | E | E | E | W | C |
|---|---|---|---|---|---|---|---|---|---|



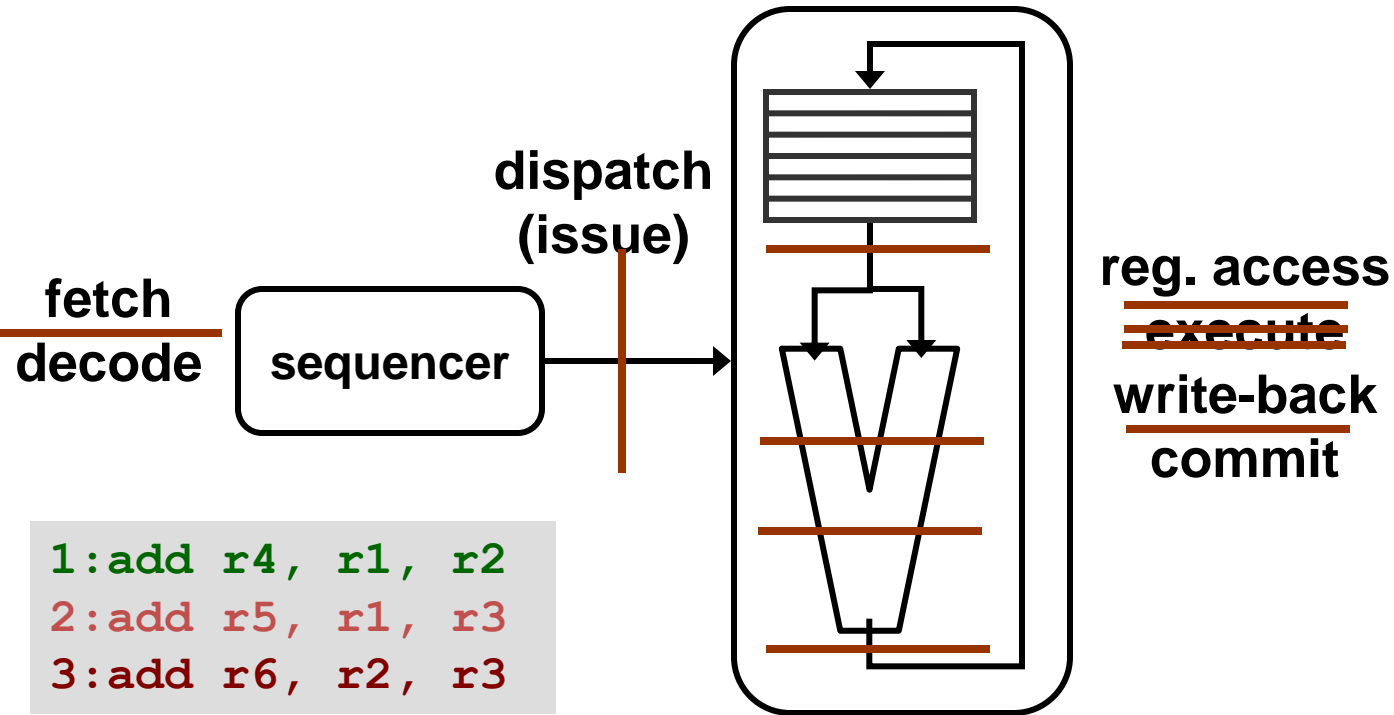
# Simplified view of a pipelined processor



# Simplified view of a pipelined processor



# Simplified view of a pipelined processor



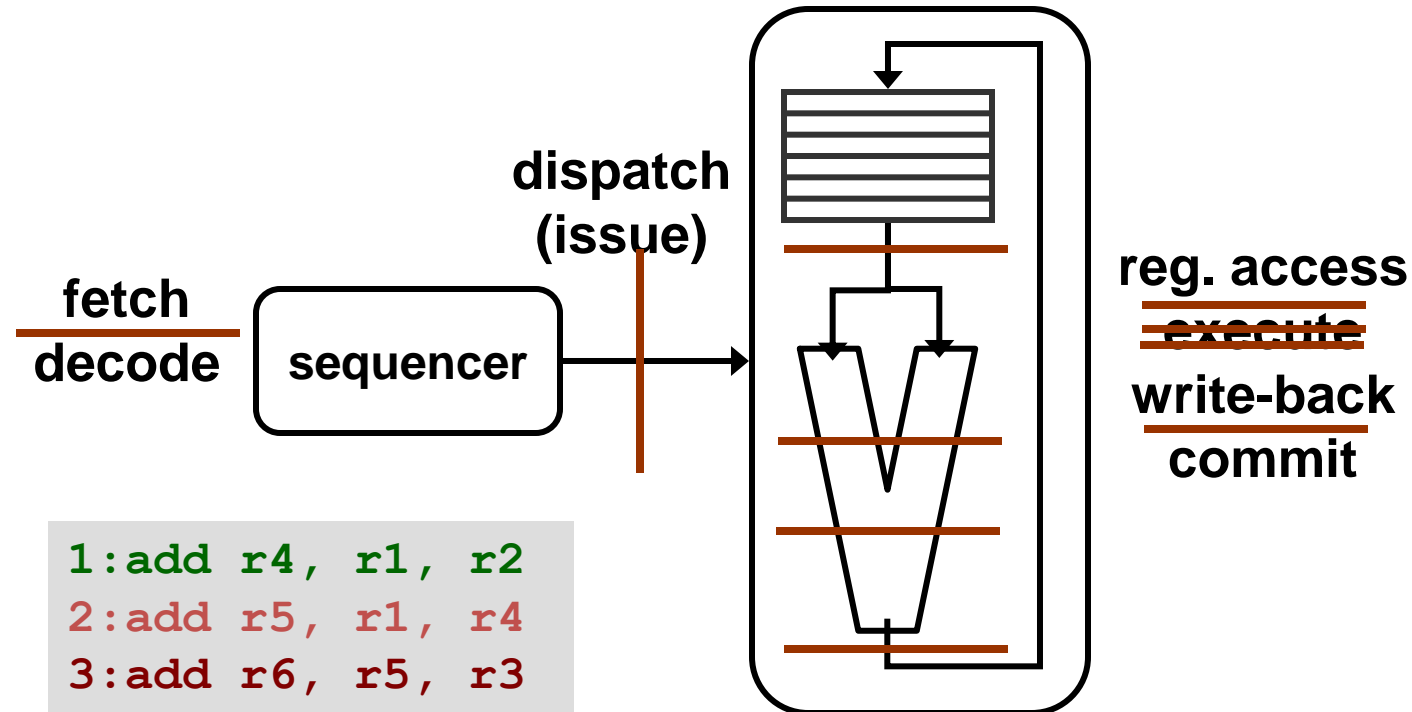
```

1: add r4, r1, r2
2: add r5, r1, r3
3: add r6, r2, r3
    
```

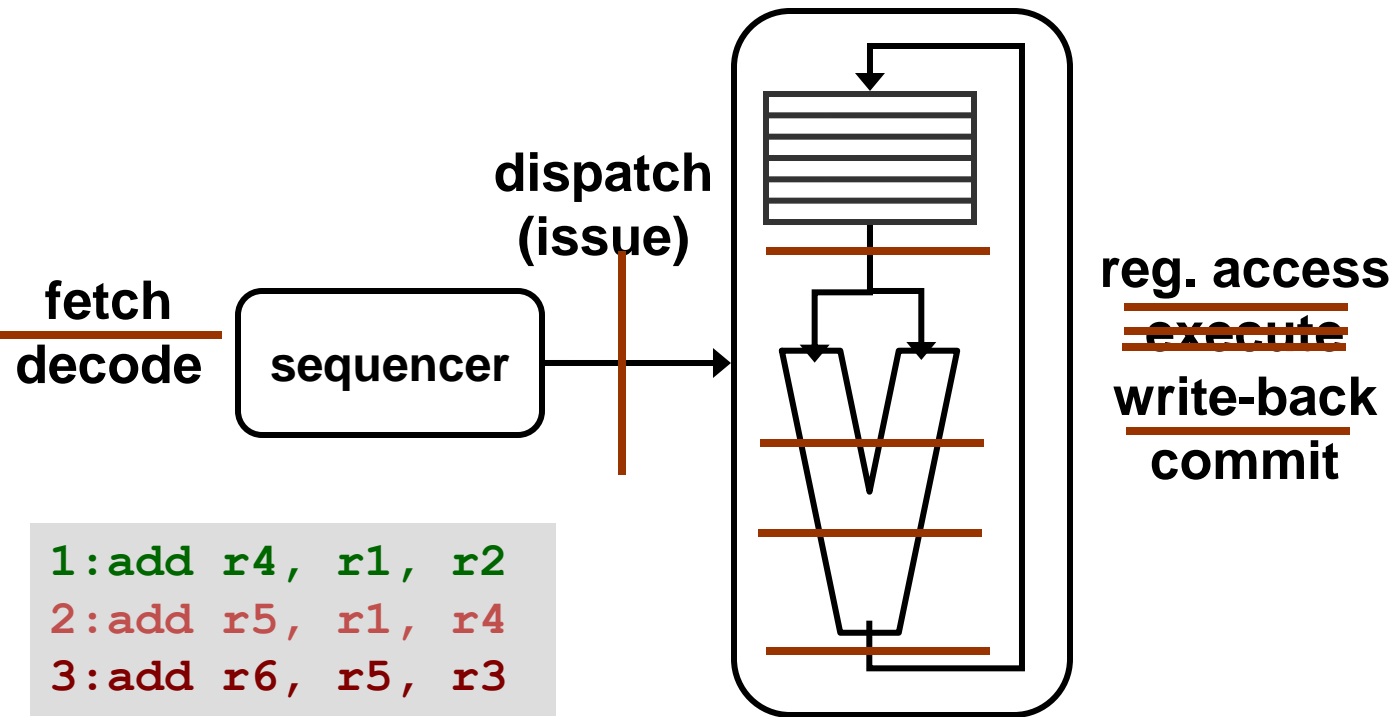
|   |   |   |   |   |   |   |   |   |   |   |   |  |  |  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|--|--|--|--|--|--|--|--|--|
| 1 | F | D | I | R | E | E | E | W | C |   |   |  |  |  |  |  |  |  |  |  |
| 2 |   | F | D | I | R | E | E | E | W | C |   |  |  |  |  |  |  |  |  |  |
| 3 |   |   | F | D | I | R | E | E | E | W | C |  |  |  |  |  |  |  |  |  |

What are the parallel resources?

# Simplified view of a pipelined processor



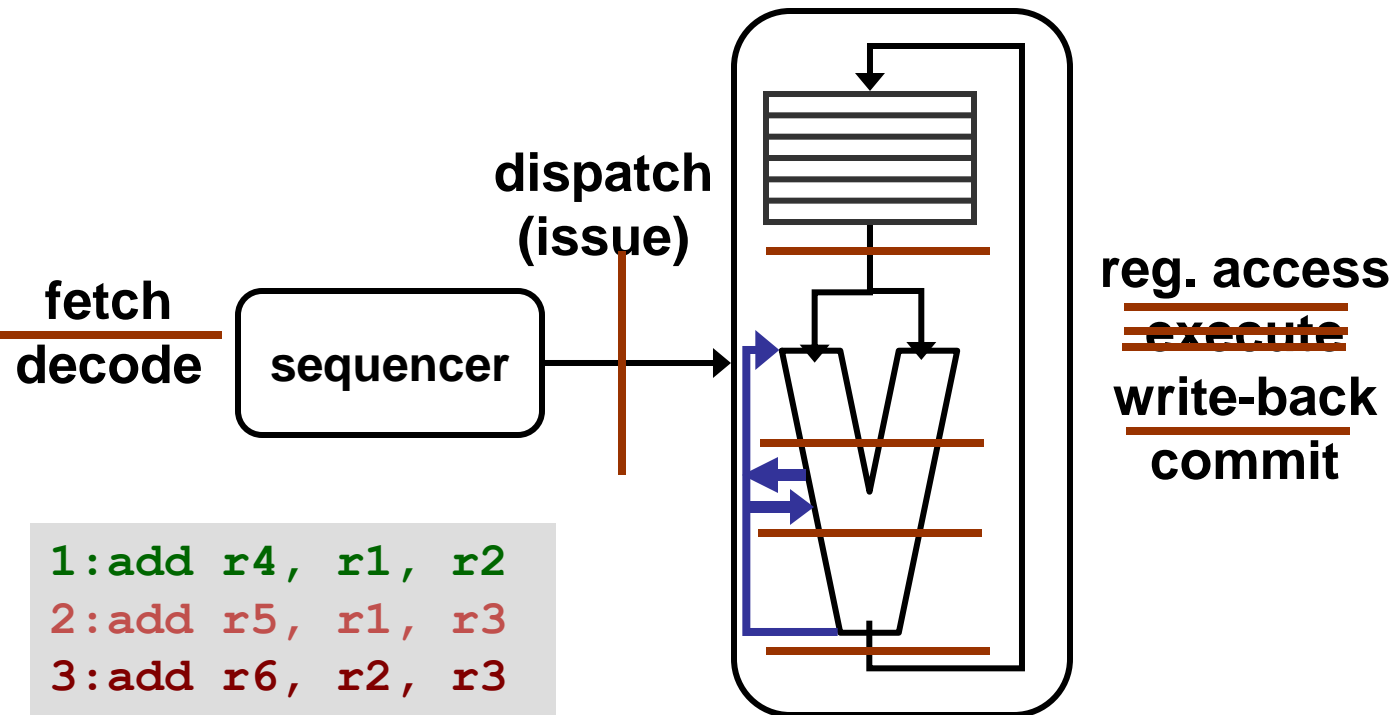
# Simplified view of a pipelined processor



|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|--|--|--|--|--|
| 1 | F | D | I | R | E | E | E | W | C |   |   |   |   |   |   |   |   |  |  |  |  |  |
| 2 |   | F | D | I |   |   |   | R | E | E | E | W | C |   |   |   |   |  |  |  |  |  |
| 3 |   |   | F | D | I |   |   |   |   |   |   | R | E | E | E | W | C |  |  |  |  |  |
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |  |  |  |  |  |



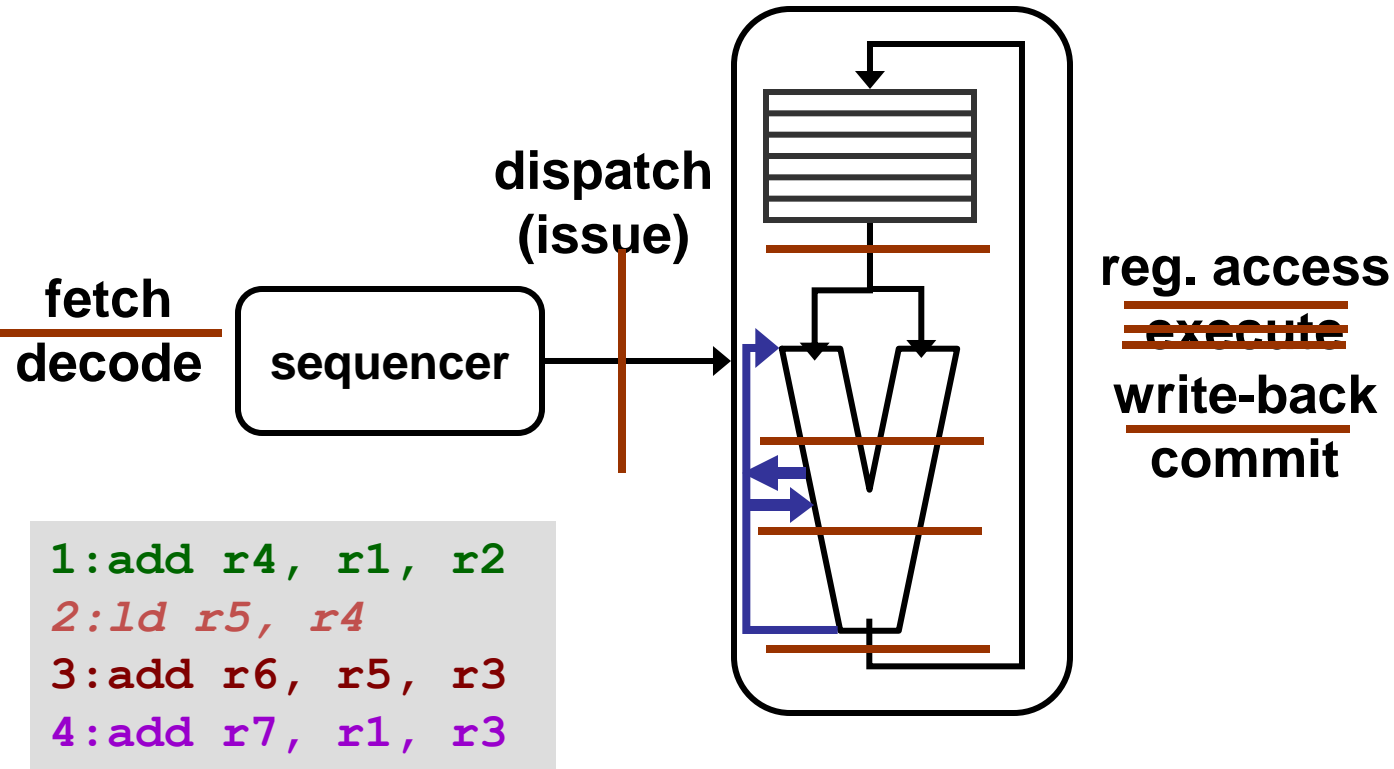
# Simplified view of a pipelined processor



|   |   |   |   |   |   |   |   |   |   |   |   |  |  |  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|--|--|--|--|--|--|--|--|--|
| 1 | F | D | I | R | E | E | E | W | C |   |   |  |  |  |  |  |  |  |  |  |
| 2 |   | F | D | I | R | E | E | E | W | C |   |  |  |  |  |  |  |  |  |  |
| 3 |   |   | F | D | I | R | E | E | E | W | C |  |  |  |  |  |  |  |  |  |

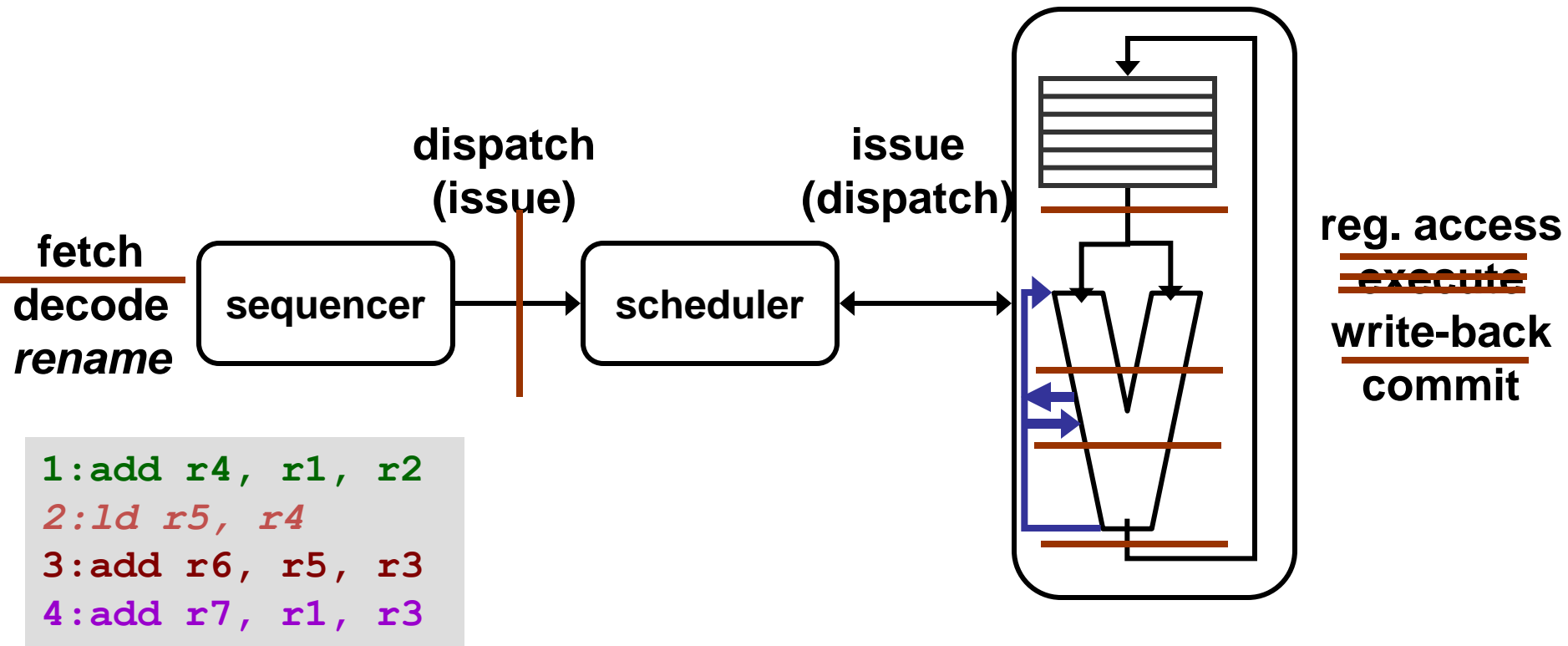
Communication and synchronization mechanisms?

# Simplified view of a pipelined processor



|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | F | D | I | R | E | E | E | W | C |   |   |   |   |   |   |   |   |   |   |   |   |
| 2 |   | F | D | I | R | L | L | L | L | L | L | L | L | W | C |   |   |   |   |   |   |
| 3 |   |   | F | D | I |   |   |   |   |   |   |   |   | R | E | E | E | W | C |   |   |
| 4 |   |   |   | F | D |   |   |   |   |   |   |   |   |   | I | R | E | E | E | W | C |

# Simplified view of a OOO pipelined processor



|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|--|--|
| 1 | F | D | I | d | R | E | E | E | W | C |   |   |   |   |   |   |   |   |  |  |
| 2 |   | F | D | I | d | R | L | L | L | L | L | L | L | L | L | L | W | C |  |  |
| 3 |   |   | F | D | I | d | R | E | E | E | W | C |   |   |   |   |   |   |  |  |

Communication and synchronization mechanisms?

# Pipelining Summary

- Pipelining is using parallelism to hide latency
  - Do useful work while waiting for other work to finish
- Multiple parallel components, not multiple instances of same component
- Examples:
  - Execution pipeline
  - Memory pipelines
    - Issue multiple requests to memory without waiting for previous requests to complete
  - Software pipelines
    - Overlap different software blocks to hide latency:  
**computation/communication**



# Parallel Execution

- Concurrency
  - what are the multiple resources?
- Communication
  - and storage
- Synchronization
  - Implicit?
  - Explicit?
- what is being **shared** ?
- What is being **partitioned** ?



# Resources in a parallel processor/system

- Execution
  - ALUs
  - Cores/processors
- Control
  - Sequencers
  - Instructions
  - OOO schedulers
- State
  - Registers
  - Memories
- Networks

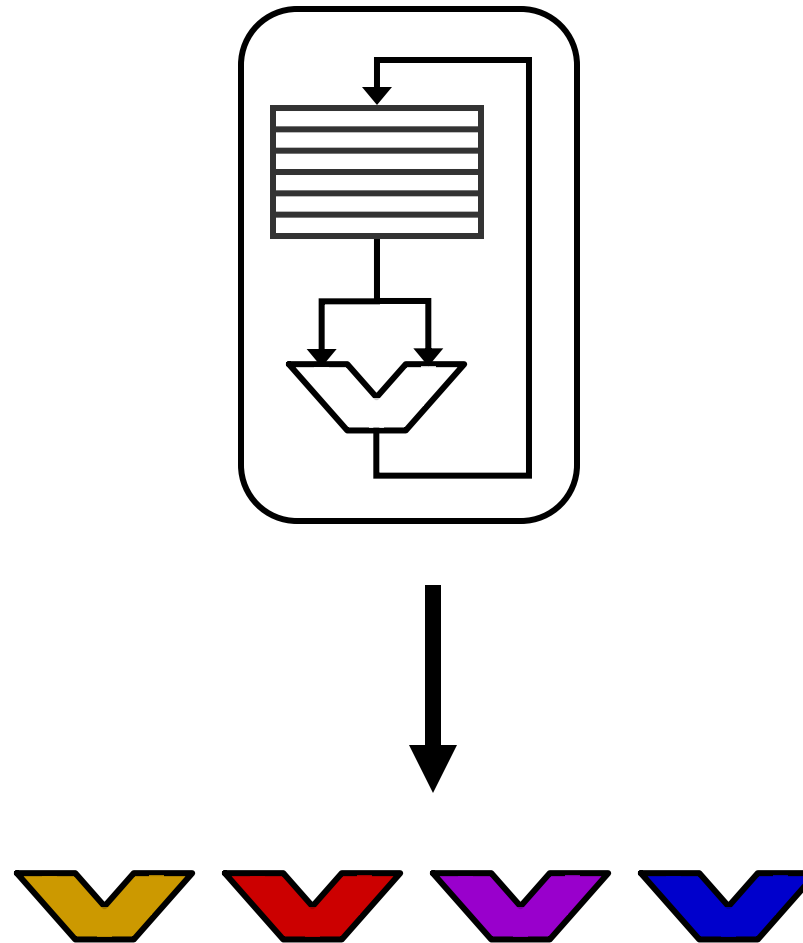


# Communication and synchronization

- Synchronization
  - Clock – explicit compiler order
  - Explicit signals (e.g., dependences)
  - Implicit signals (e.g., flush/stall)
    - More for pipelining than multiple ALUs
- Communication
  - Bypass networks
  - Registers
  - Memory
  - Explicit (over some network)

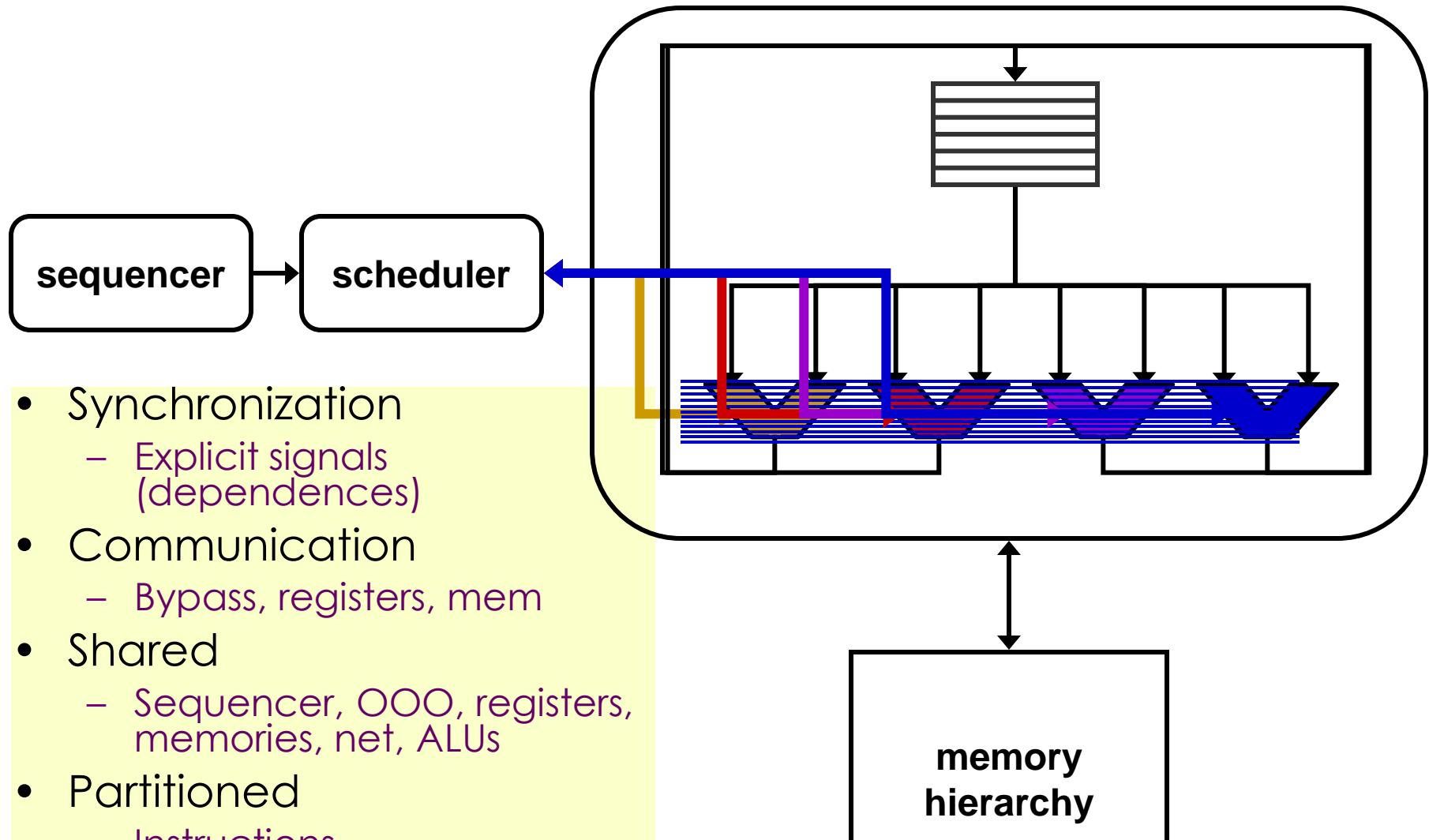


# Organizations for ILP (for multiple ALUs)





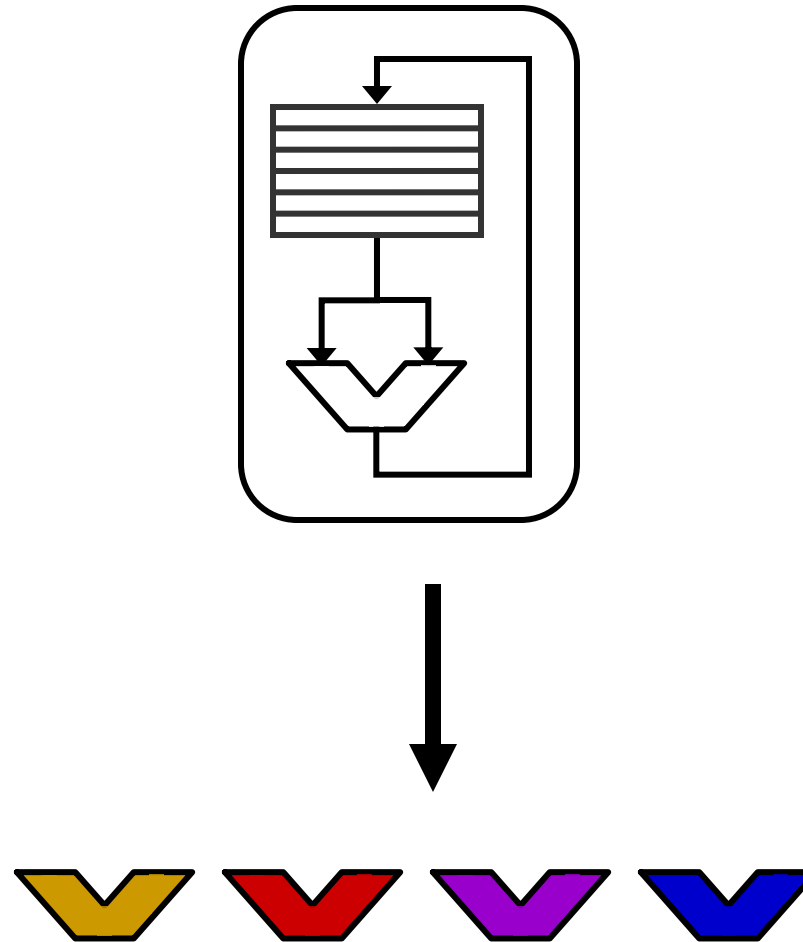
# Superscalar (ILP for multiple ALUs)



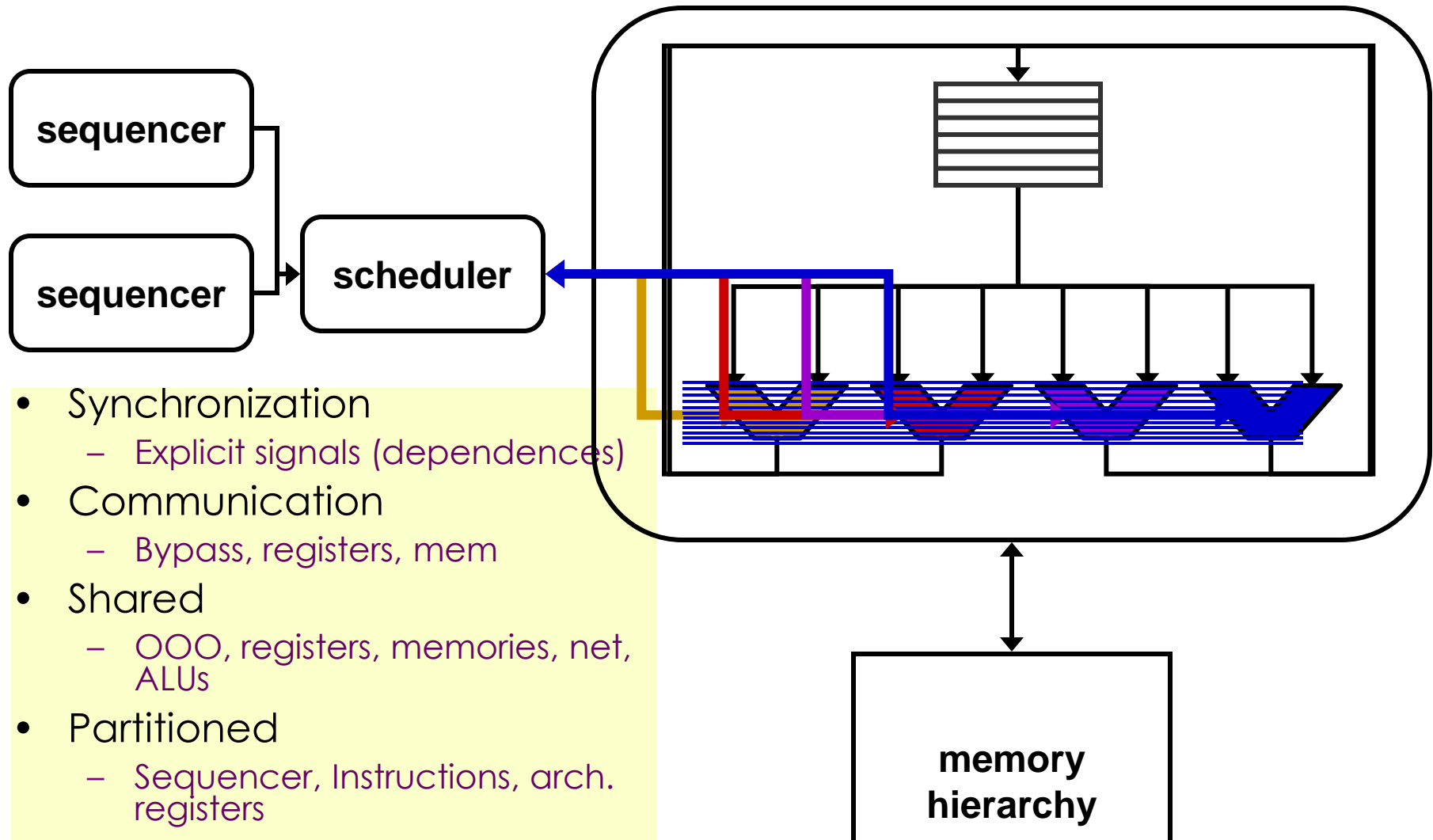
- Synchronization
  - Explicit signals (dependences)
- Communication
  - Bypass, registers, mem
- Shared
  - Sequencer, OOO, registers, memories, net, ALUs
- Partitioned
  - Instructions

**How many ALUs?**

# SMT/TLS (ILP for multiple ALUs)

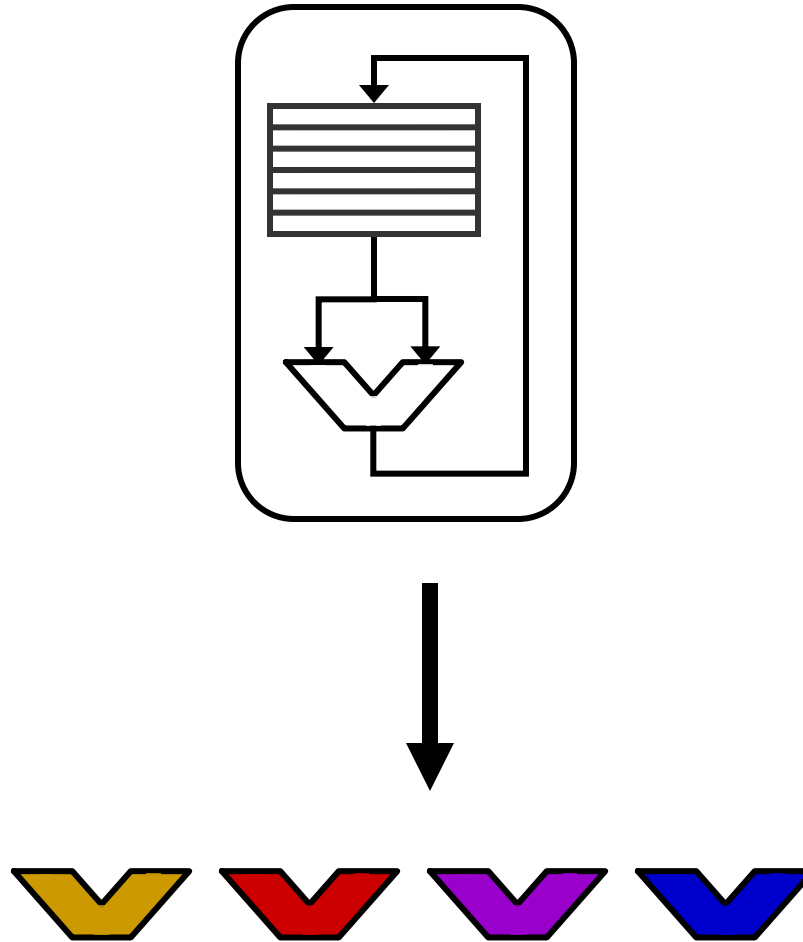


# SMT/TLS (ILP for multiple ALUs)

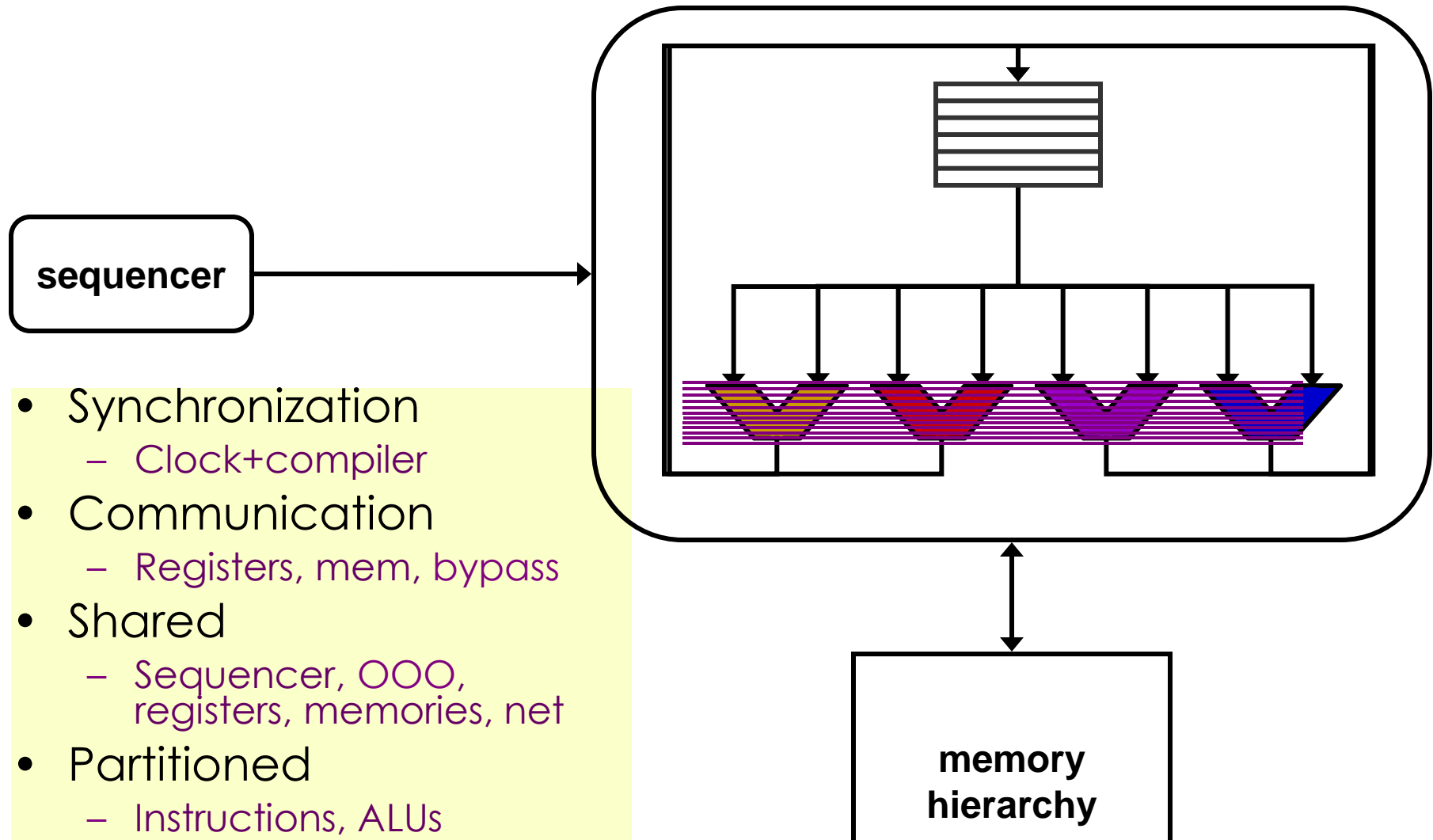


**Why is this ILP? How many threads?**

# VLIW (ILP for multiple ALUs)

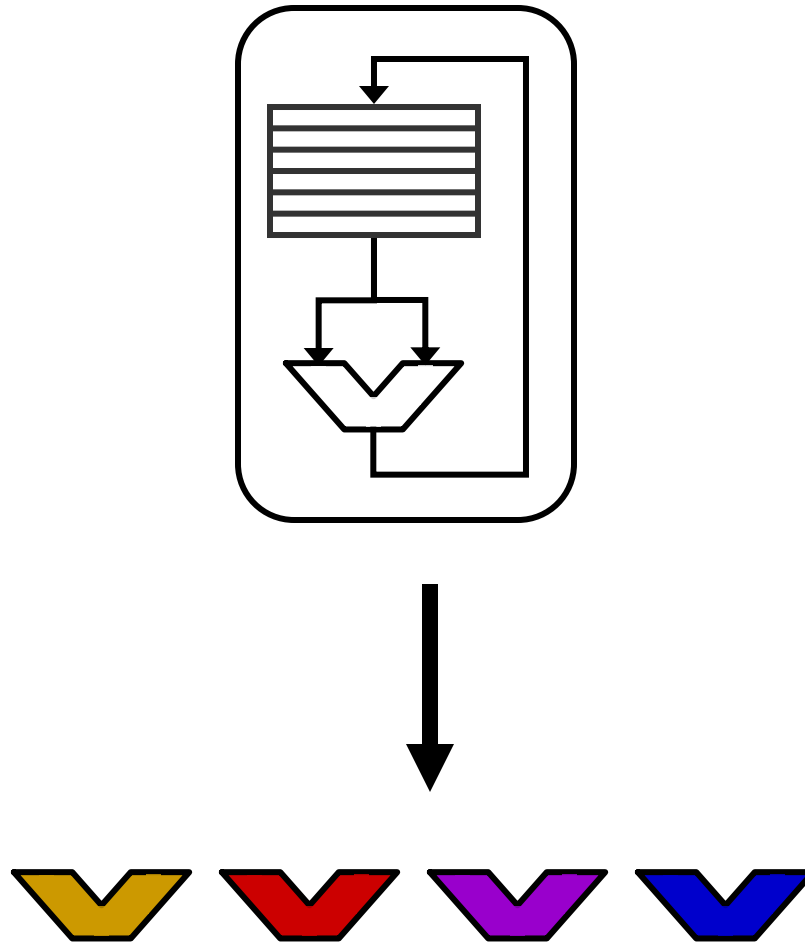


# VLIW (ILP for multiple ALUs)

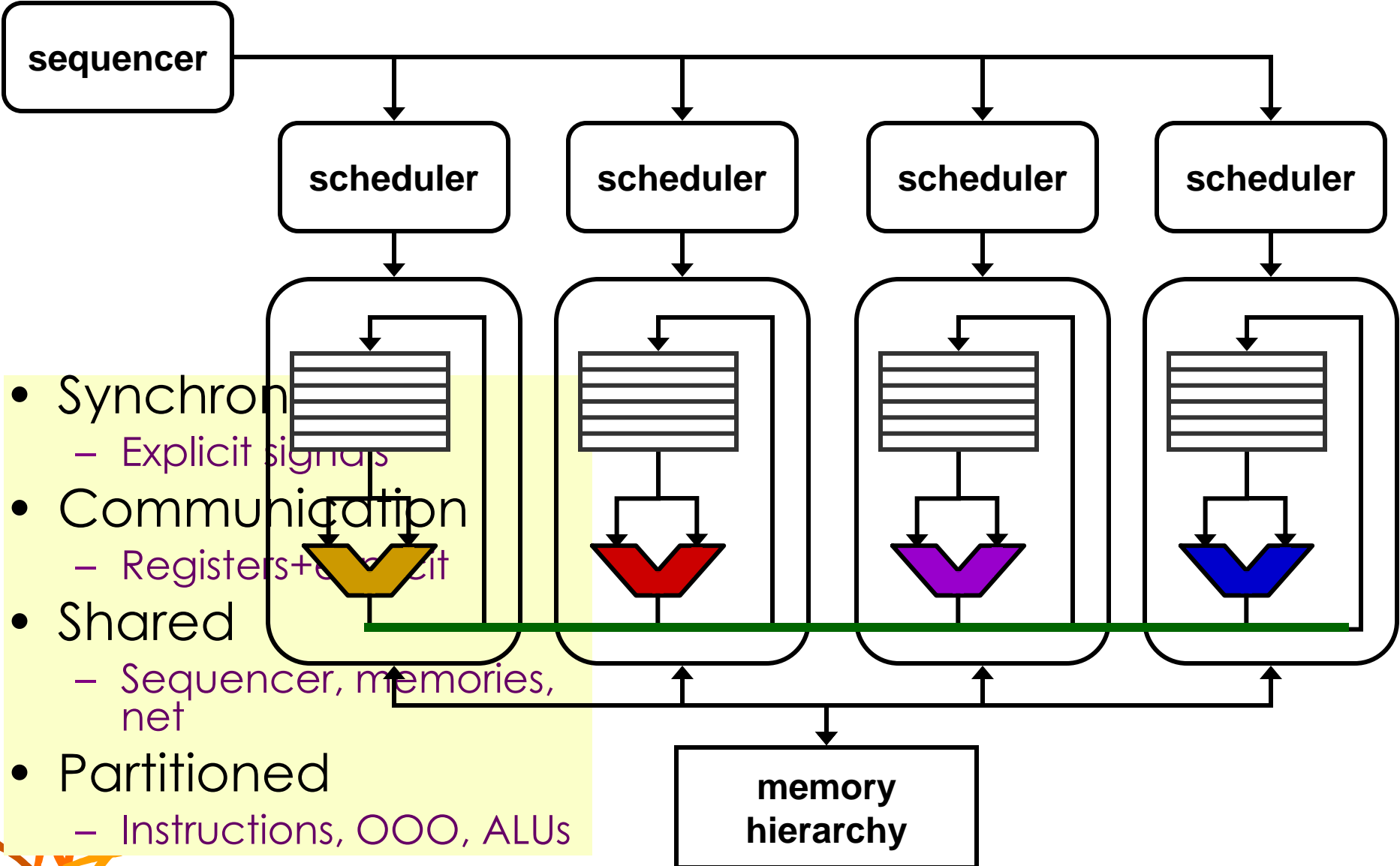


**How many ALUs?**

# Explicit Dataflow (ILP for multiple ALUs)



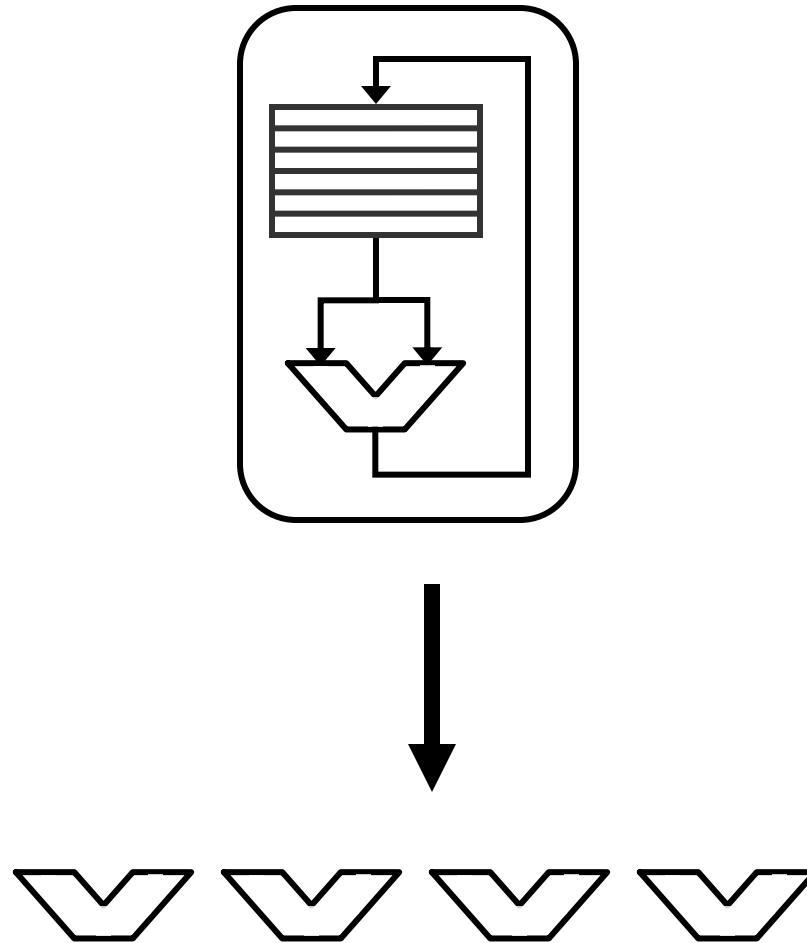
# Explicit Dataflow (ILP for multiple ALUs)



- Synchronous
  - Explicit signals
- Communication
  - Registers + multi-ported ALU
- Shared
  - Sequencer, memories, net
- Partitioned
  - Instructions, OOO, ALUs



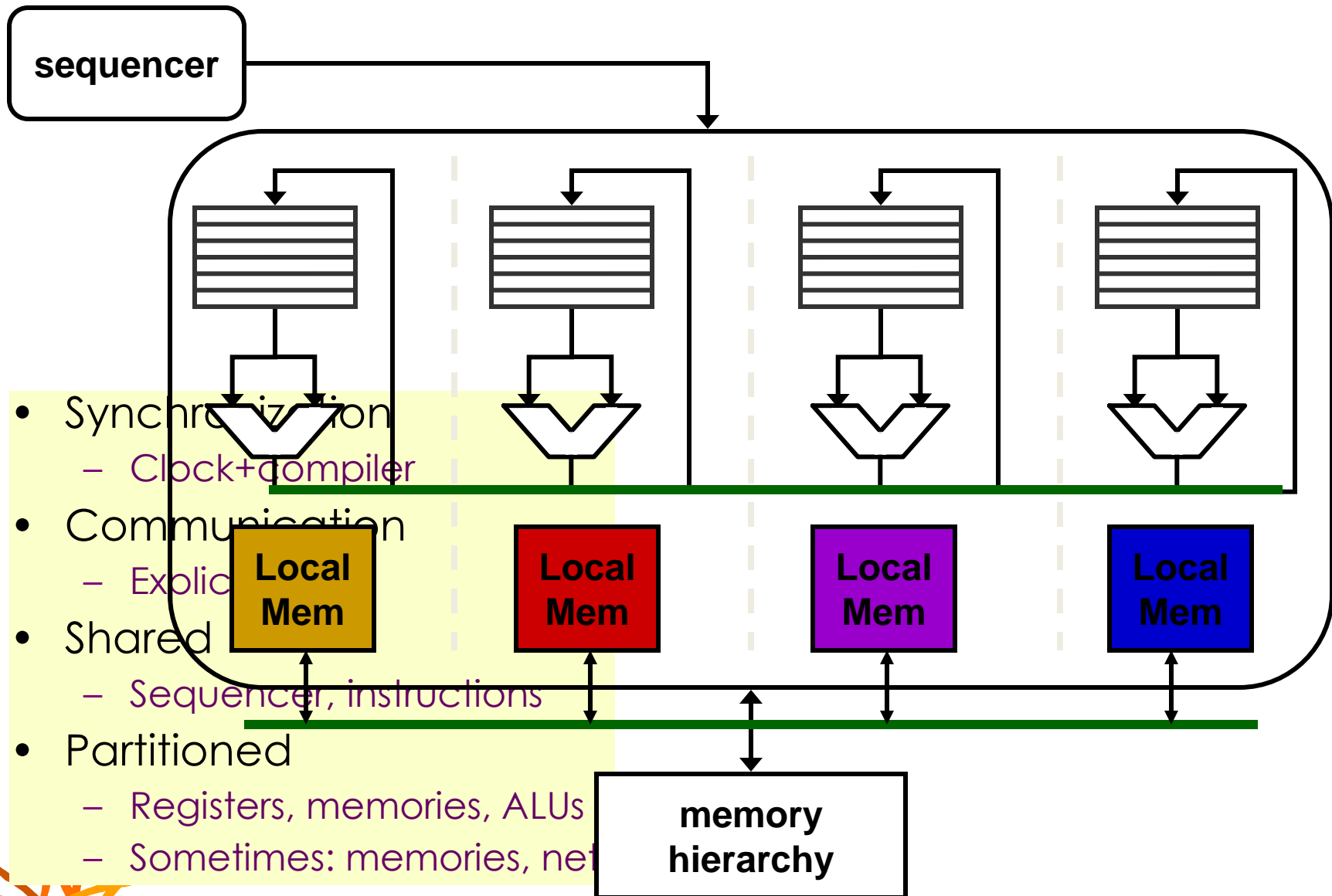
# DLP for multiple ALUs



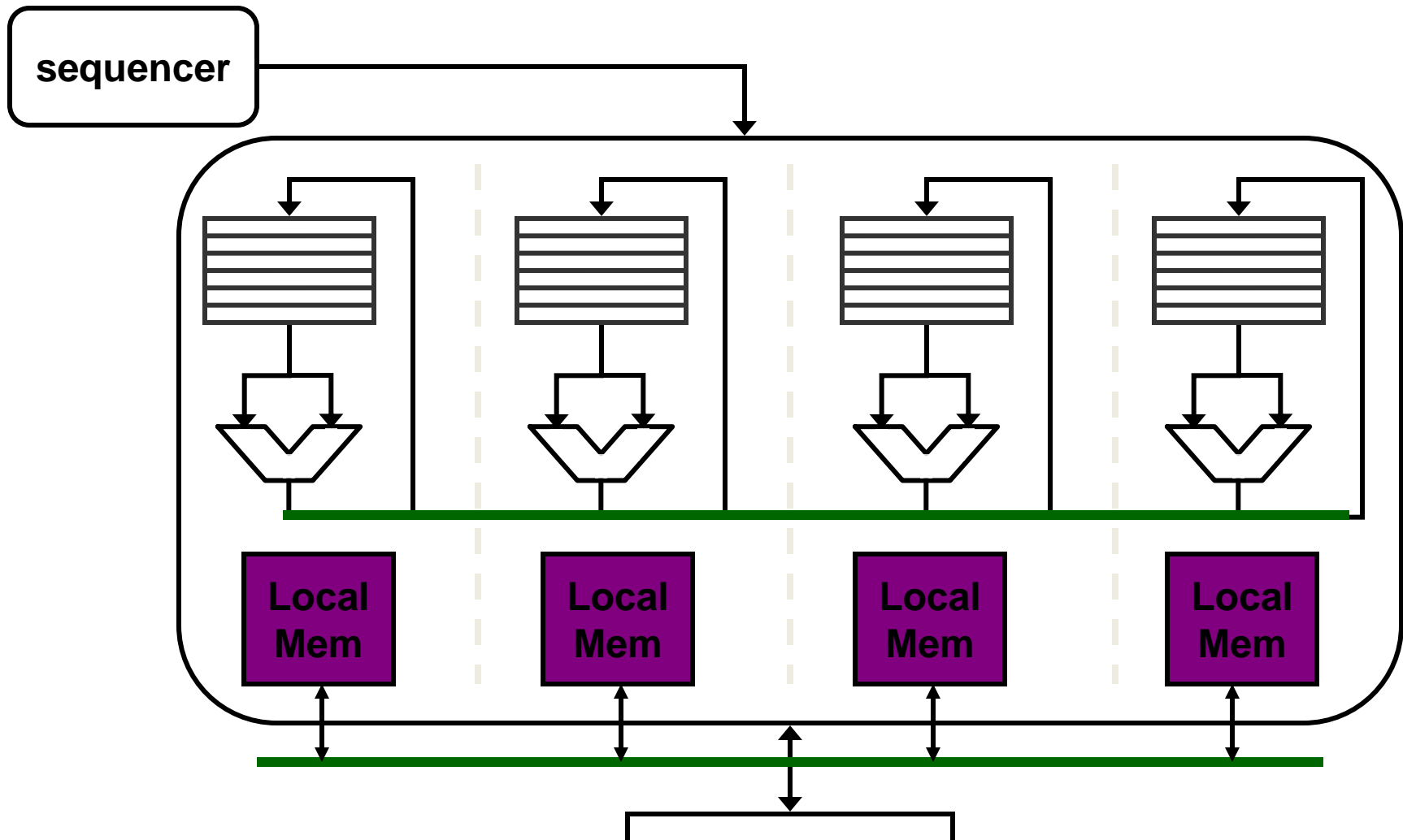
**From HW – this is SIMD**



# SIMD (DLP for multiple ALUs)

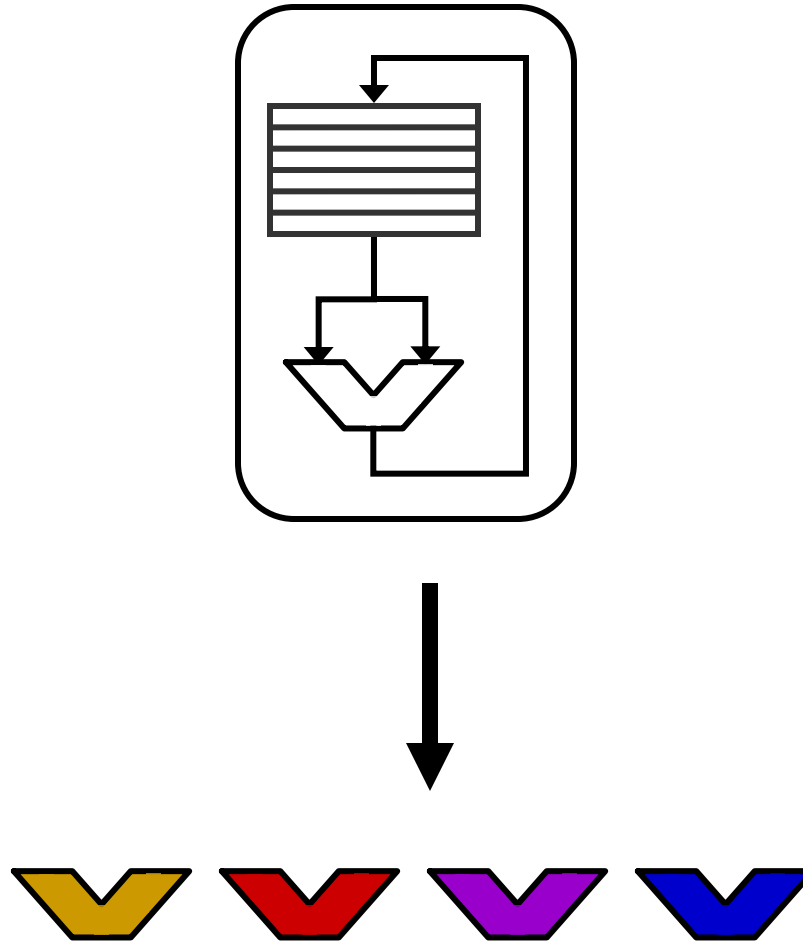


# Vectors (DLP for multiple ALUs)

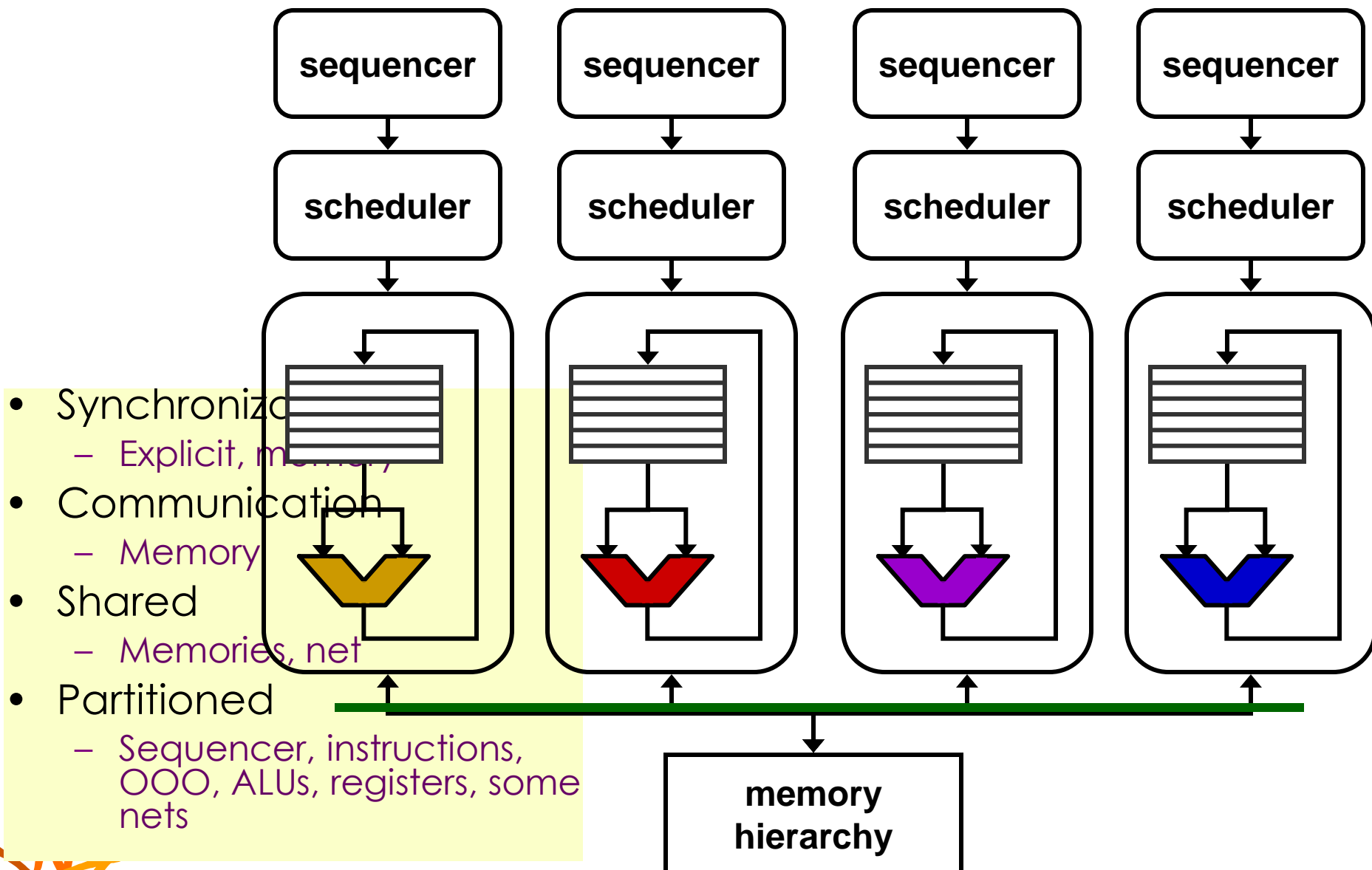


**Vectors: memory addresses are part of single-instruction and not part of multiple-data**

# TLP for multiple ALUs

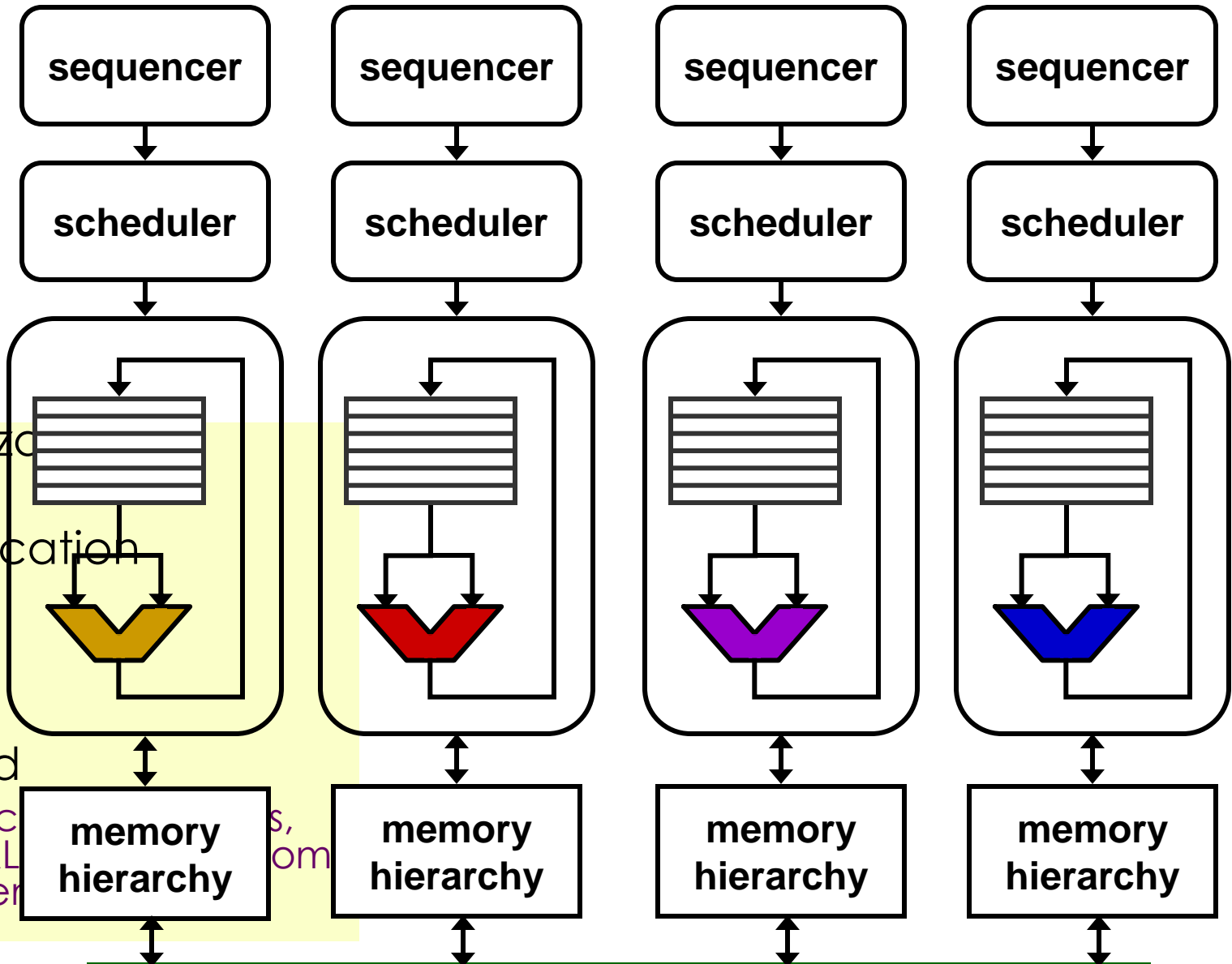


# MIMD – shared memory (TLP for multiple ALUs)



# MIMD – distributed memory

- Synchronized
  - Explicit
- Communication
  - Explicit
- Shared
  - Net
- Partitioned
  - Sequence
  - OOO, ALU
  - nets, mem



# Summary of communication and synchronization

| Style              | Synchronization | Communication |
|--------------------|-----------------|---------------|
| <b>Superscalar</b> |                 |               |
| <b>VLIW</b>        |                 |               |
| <b>Dataflow</b>    |                 |               |
| <b>SIMD</b>        |                 |               |
| <b>MIMD</b>        |                 |               |



# Summary of sharing in ILP HW

| Style       | Seq | Inst | OOO | Regs | Mem | ALUs | Net |
|-------------|-----|------|-----|------|-----|------|-----|
| Superscalar |     |      |     |      |     |      |     |
| SMT/TLS     |     |      |     |      |     |      |     |
| VLIW        |     |      |     |      |     |      |     |
| Dataflow    |     |      |     |      |     |      |     |



# Summary of sharing in DLP and TLP

| Style  | Seq | Inst | OOO | Regs | Mem | ALUs | Net |
|--------|-----|------|-----|------|-----|------|-----|
| Vector |     |      |     |      |     |      |     |
| SIMD   |     |      |     |      |     |      |     |
| MIMD   |     |      |     |      |     |      |     |

