



# Flexible and Efficient Memory System for a High Performance Processing Pipeline

Mario Baldi  
Distinguished Technologist

Mike Galles  
ASIC Architect

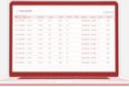
Pensando Systems

PENSANDO

# Pensando Distributed Services Platform

Centrally Managed

Policy and Services Manager



Policy



Orchestration & Provisioning

Automation



REST API

Observability



Troubleshooting & Security

Ecosystem



Compute, Analytics, IT Ops

Stateful Firewall



Encryption & TLS Offload



Micro Segmentation



Load Balancer



Networking



Storage Services



Telemetry



Distributed Services Cards



P4 Programmable Processor Pipelines



PENSANDO

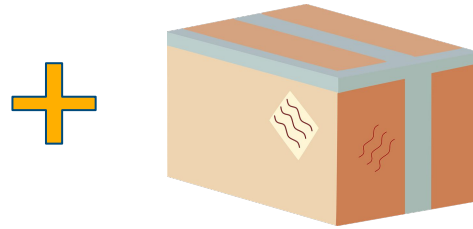
Why Pipeline? [Why P4]

# Performance

# What is specific of a pipeline?

Each stage executes a limited task

- Small number of instructions
  - E.g., an action in P4
- Simple data structures
  - E.g., a table in P4



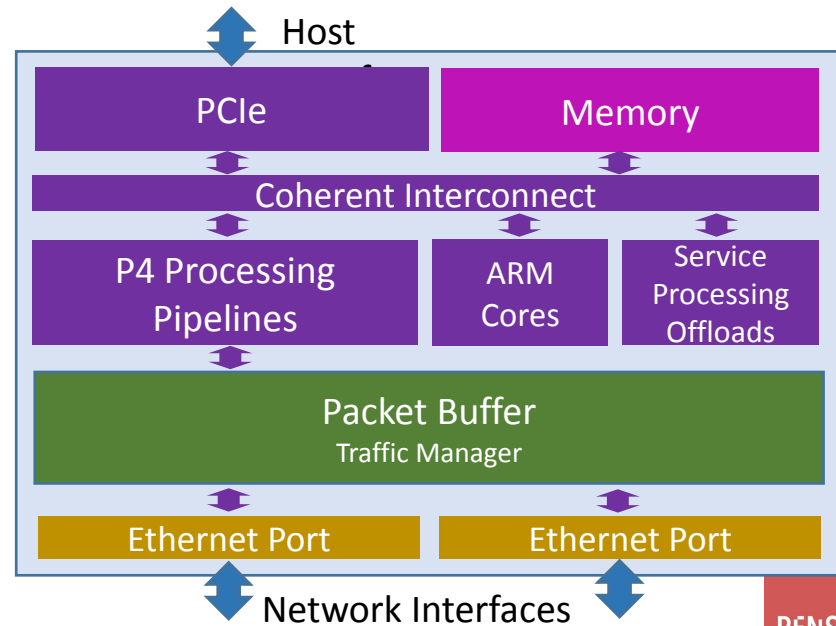
- Similar processing
  - Limited set of operations
- “Well-defined” data access

# How do we get high performance?

- **Optimize execution** of the limited set of operations
- **Specialized hardware** for table lookup
- **Fast local instruction memory** for small programs
- **Fast local data memory** for small tables

# Memory is key to achieve expected **performance**

- Pipeline
- P4 programs
- [Overall solution]



# Memory Design Options

# Local Memory

Each stage has a dedicated amount of memory

- Instruction and data (table) memory
- Possibly writeback memory (e.g., registers)

Simple architecture

(Constant) low latency

Scalability

With stages and pipelines

Limited size

Tables are possibly split across stages

Possibly inefficient

Unused memory in one stage cannot be used by others

No shared tables



# Central Memory

Stages share a common memory

- Instruction and data (table) memory
- Possibly writeback memory (e.g., registers)
- Possibly off-chip



Size can be large



Efficient and flexible usage



Shared by stages and devices



High latency



Potential bottleneck

# Central Memory with Local Cache

Mask access latency and limit contention

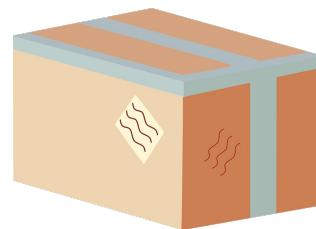


possibly best of breed

In traditional, general purpose computing

- Spatial locality
- Temporal locality

In pipelined packet processing, also



## Functional locality

# Functional Locality

A given function accesses a limited number of instructions and amount of data



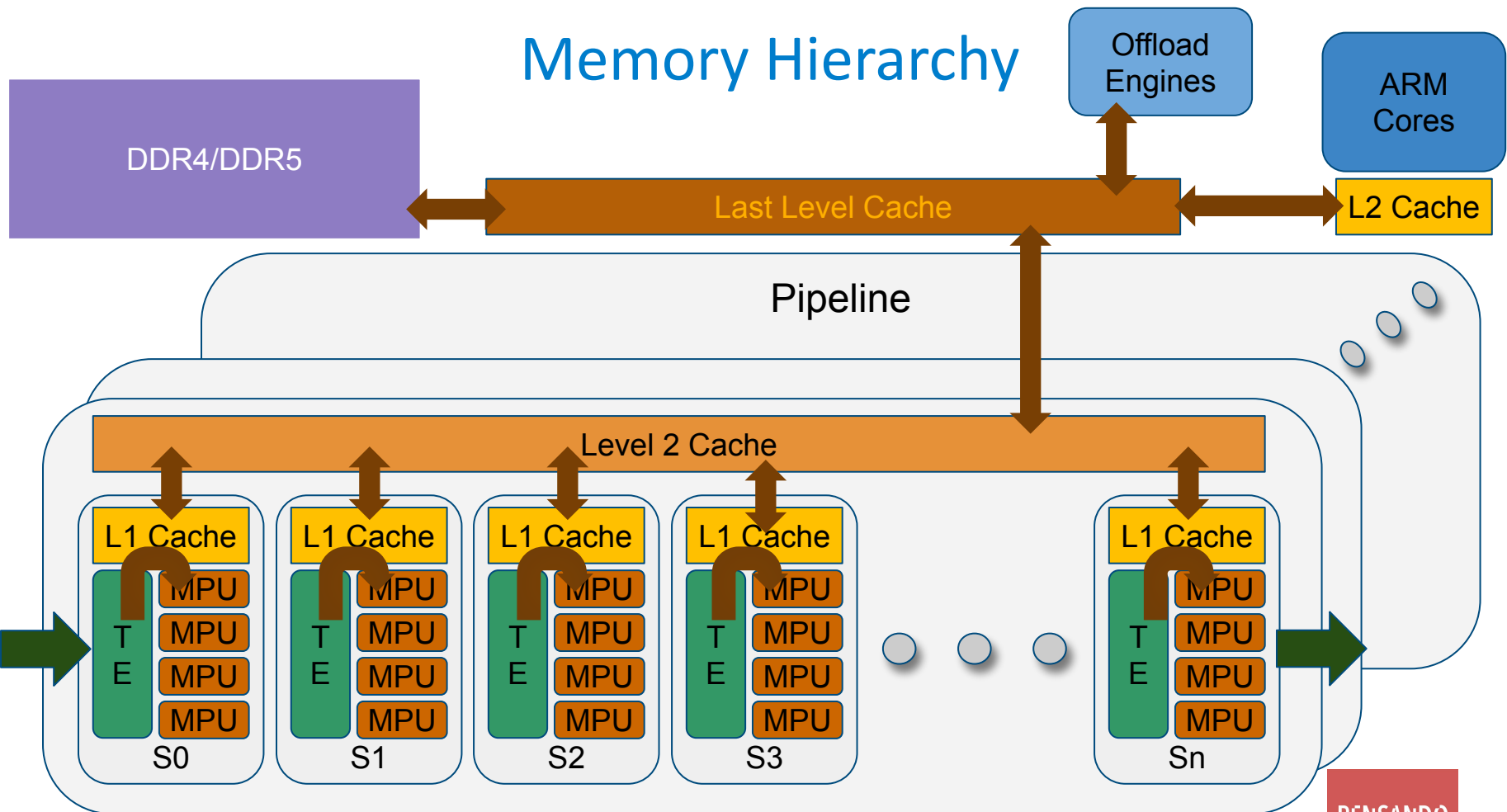
Instructions and data for the particular function are in the cache of the stage where execution takes place

- Short program at each stage
- Small tables are common

## With larger tables

- 👍 Packets of the same flow access the same entries and are processed close in time
  - Temporal locality
- 👎 Large number of active flows may result in misses

# Memory Hierarchy



## Variable Access Latency

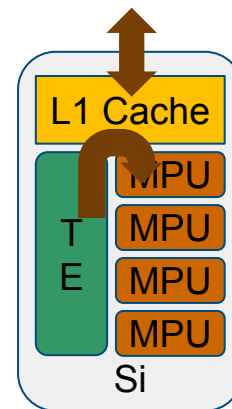
- Low on cache hits
- High on cache misses

### Options for handling high latency

- Stall the pipeline
  - Affect throughput → performance
- Process another packet

## Handling Variable Latency

- Table Engine (TE) can issue up to 28 reads in parallel
  - Possibly related to multiple packets
- As each completes, action data is provided to one of the Match Processing Units (MPUs)
  - MPU fetches instructions
    - Likely to still be in cache
- Packet ordering preserved

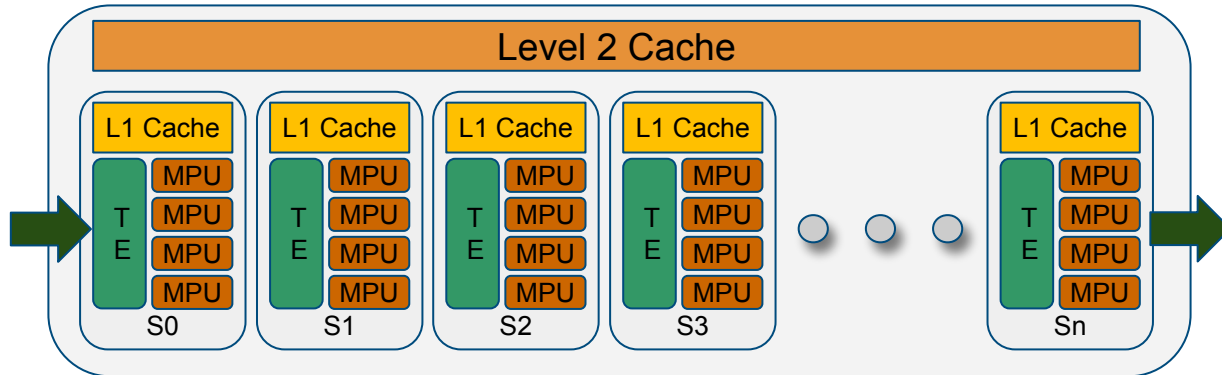


# Experimental Evaluation

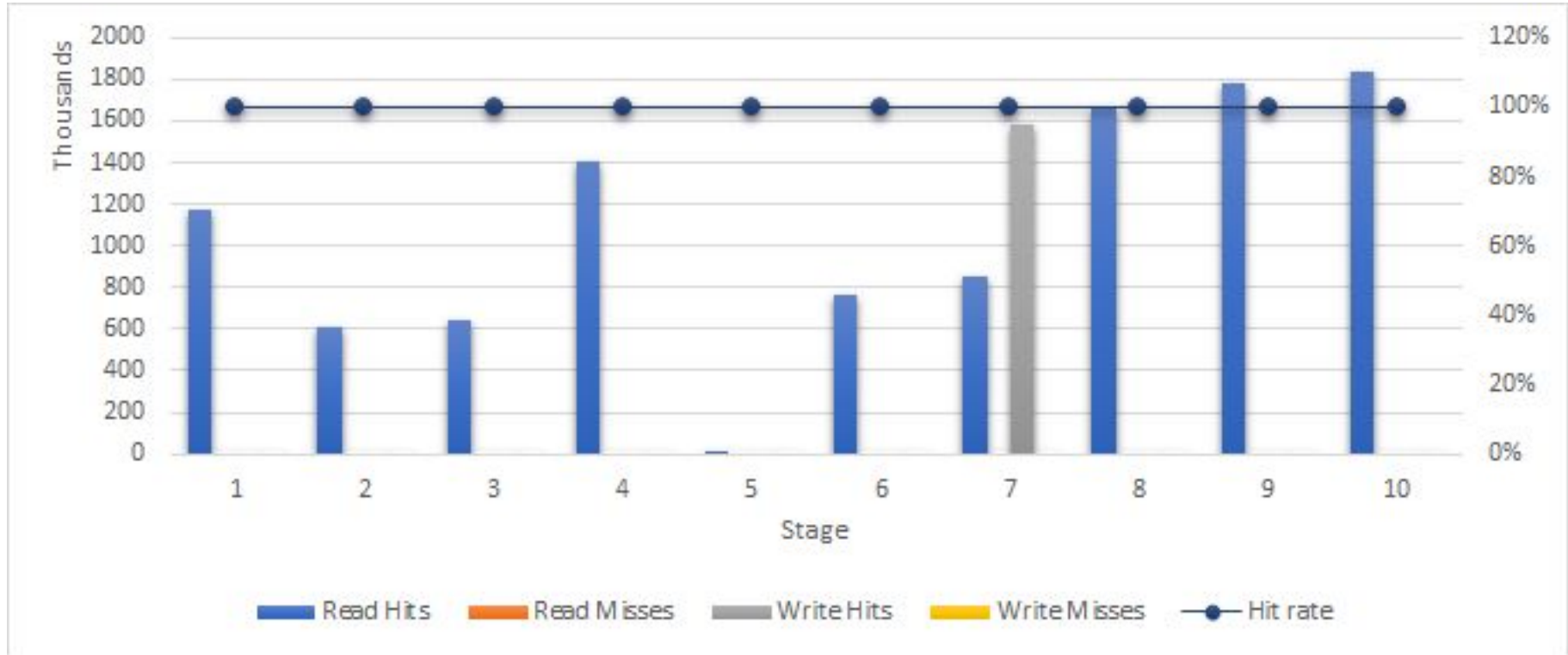


# Experimental Setup

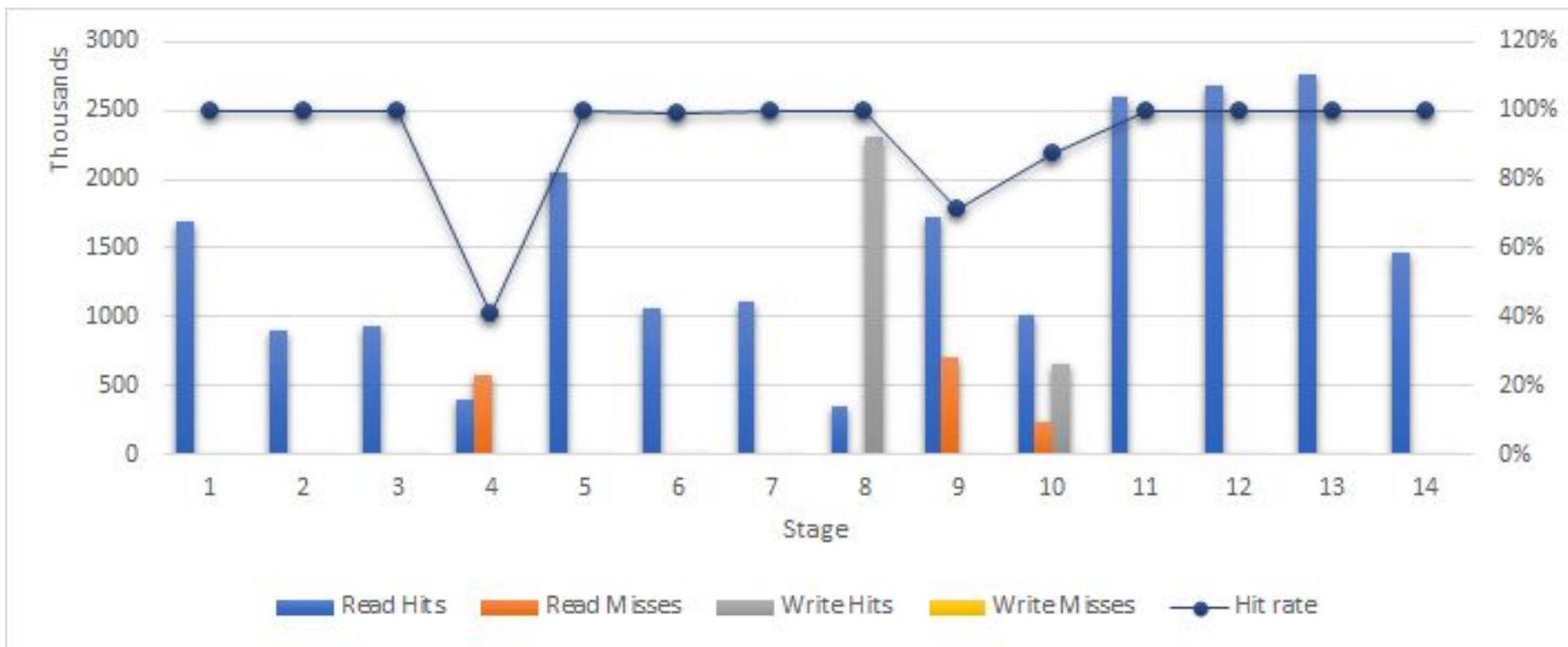
- Production P4 code for cloud data center
  - Two different P4 programs
- Various configuration and traffic scenarios
- L1 data cache hits and misses per stage
- Instruction cache hits and misses per MPU



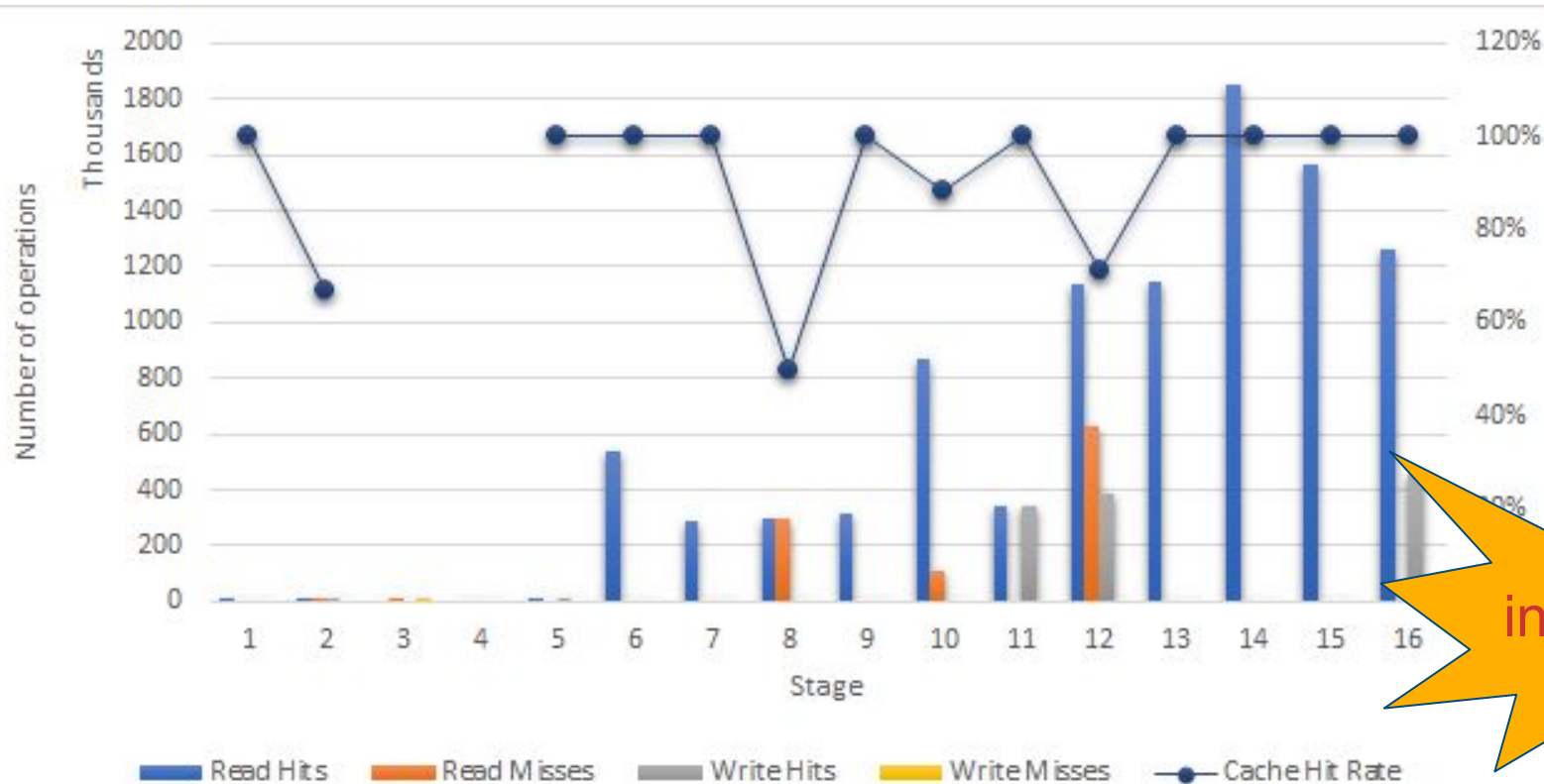
# 1,000 active flows



# 1,800,000 active flows



# Other P4 Program



100%  
instruction  
hit rate

# Instruction Cache

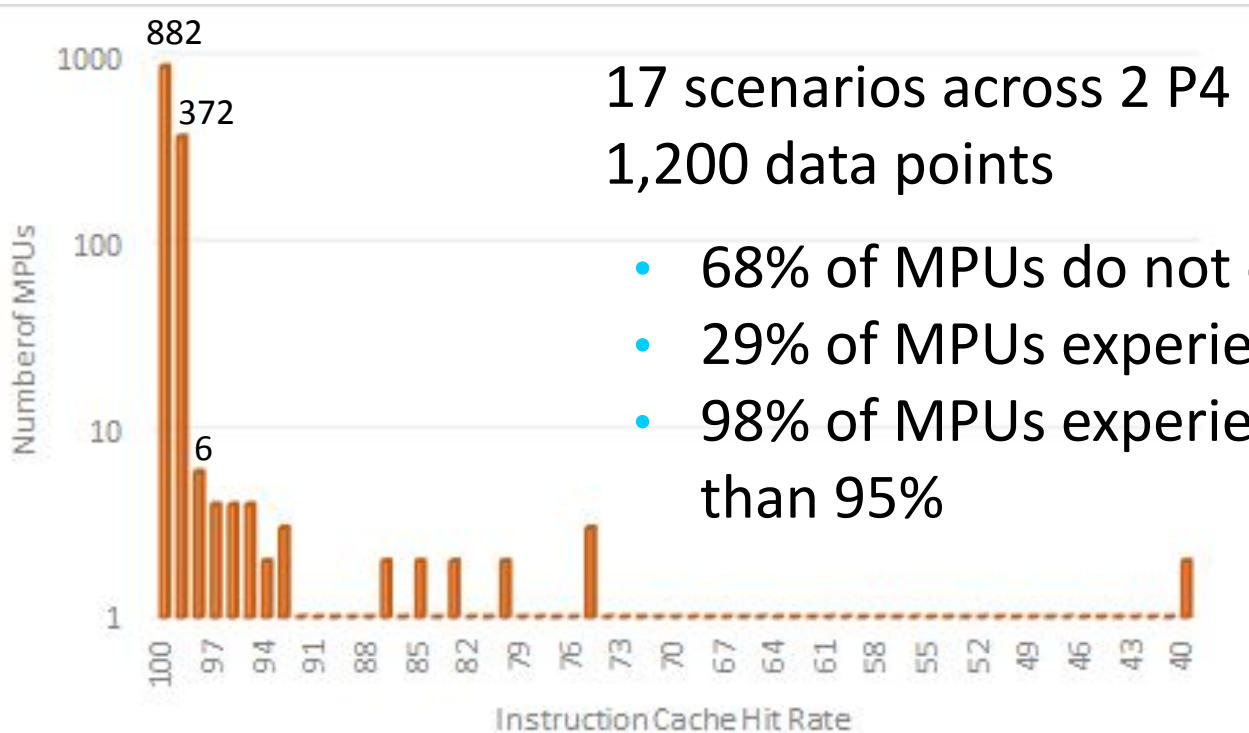
- Common actions compile to 30-50 instructions
- L1 cache can contain 2k 8-byte instructions
- Potentially 40-60 actions per stage

**Instruction misses are extremely rare**

# Instruction Cache

17 scenarios across 2 P4 programs, over 1,200 data points

- 68% of MPUs do not experience misses
- 29% of MPUs experience 1% misses
- 98% of MPUs experience hit rate higher than 95%



## Concluding Remarks

- P4 over pipelines for performance and flexibility
- Memory plays a key role in the performance and flexibility of target architectures
- A large central memory with a caching system successfully provides both
  - At least for a device at the network edge



# Thank You

[baldi\[at\]pensando.io](mailto:baldi@pensando.io)  
[pensando.io/blog](https://pensando.io/blog)

**PENSANDO**