



# Trident

*Automated system tool for collecting and analyzing hardware performance counters*

Servesh Muralidharan & David Smith

IT-DI-WLCG, UP Team

CERN

Dec 2018



# Introduction

## ❑ Existing tools

- ❑ Execution and profiling are separate steps
- ❑ Focus more on HPC Applications
  - ❑ Fine grained analysis
  - ❑ Architecture specific optimizations

## ❑ Typical HEP Applications (**As it is today!!!**)

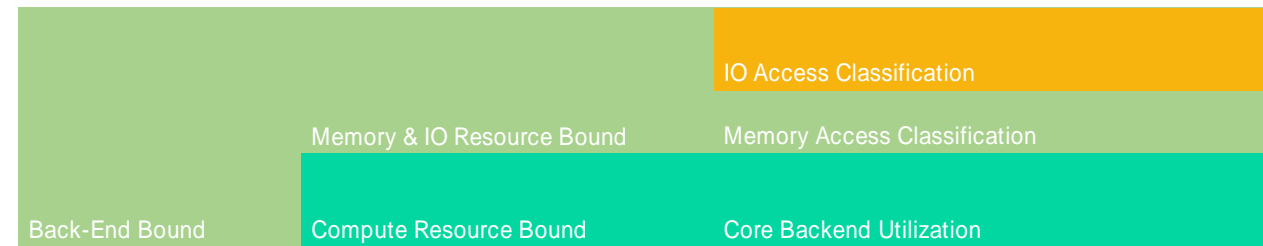
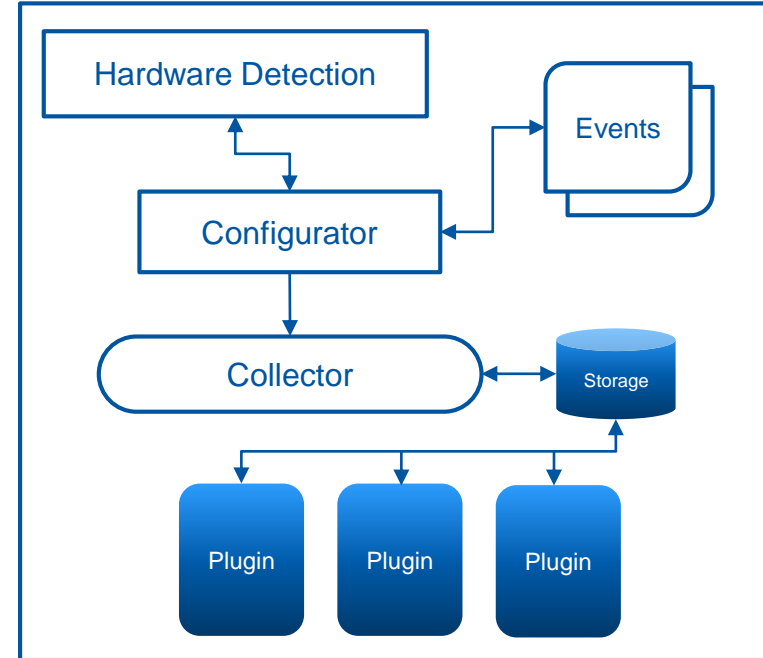
- ❑ A cluster of several hundred algorithms
- ❑ Complex framework interconnecting these algorithms
- ❑ Linear instruction spread (**no hotspots**)
- ❑ Non existence of classical numerical loops
- ❑ Average runtime of several hours

## ❑ Proposed Solution: Trident

- ❑ Continuous performance analysis
- ❑ System perspective
- ❑ Node level measurements
- ❑ Non-virtualized environment
- ❑ Low overhead and lightweight
- ❑ Utilizes hardware counters to measure Core, Memory and IO metrics
- ❑ Presents it in a HEP-human readable form

# Trident Architecture

- ❑ Detects the events supported by the node
- ❑ Extensible event files for future architectures
- ❑ Configures metrics based on these events
- ❑ These metrics are sampled at preconfigured rate
- ❑ Stored temporarily in local disk
- ❑ Plugins perform transformations to the data
- ❑ Simple low overhead online analysis
- ❑ Data exported to centralized storage for offline analysis



Extended Top Down Analysis

# Trident - Core Performance Analysis

## Instruction Per Cycle (IPC)

- Denotes ratio of parallel instructions executed
- Modern processors like Intel Haswell EP can do 4 IPC

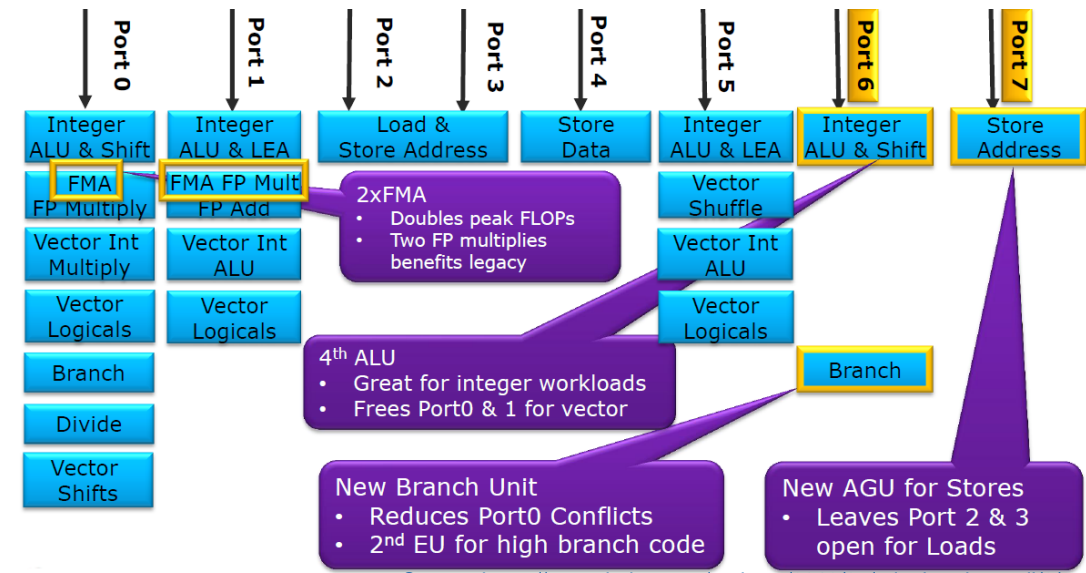
## Top down characterization

- Identifies the resources dominated by workload
- Front-End – fetch and decode program code
- Back-End – monitor and execution of uOP once the dependent data operands availability
- Retiring – Completion of the uOP
- Bad speculation – uOPs that are cancelled before retirement due to branch misprediction

## Execution Unit Port Utilization

- Determines how many cycles the port was busy
- Identifies broadly the pressure from different types of uOPs
- INT & FP operations, address calculation, etc.,

## Intel Haswell EP



Source: <https://arstechnica.com/gadgets/2013/05/a-look-at-haswell/2/>

Category	Expected Range of Pipeline Slots in This Category, for a Hotspot in a Well-Tuned:		
	Client/Desktop Application	Server/Database/Distributed application	High Performance Computing (HPC) application
Retiring	20-50%	10-30%	30-70%
Back-End Bound	20-40%	20-60%	20-40%
Front-End Bound	5-10%	10-25%	5-10%
Bad Speculation	5-10%	5-10%	1-5%

Source: <https://software.intel.com/en-us/vtune-amplifier-help-tuning-applications-using-a-top-down-microarchitecture-analysis-method>



# Trident - Memory Performance Analysis

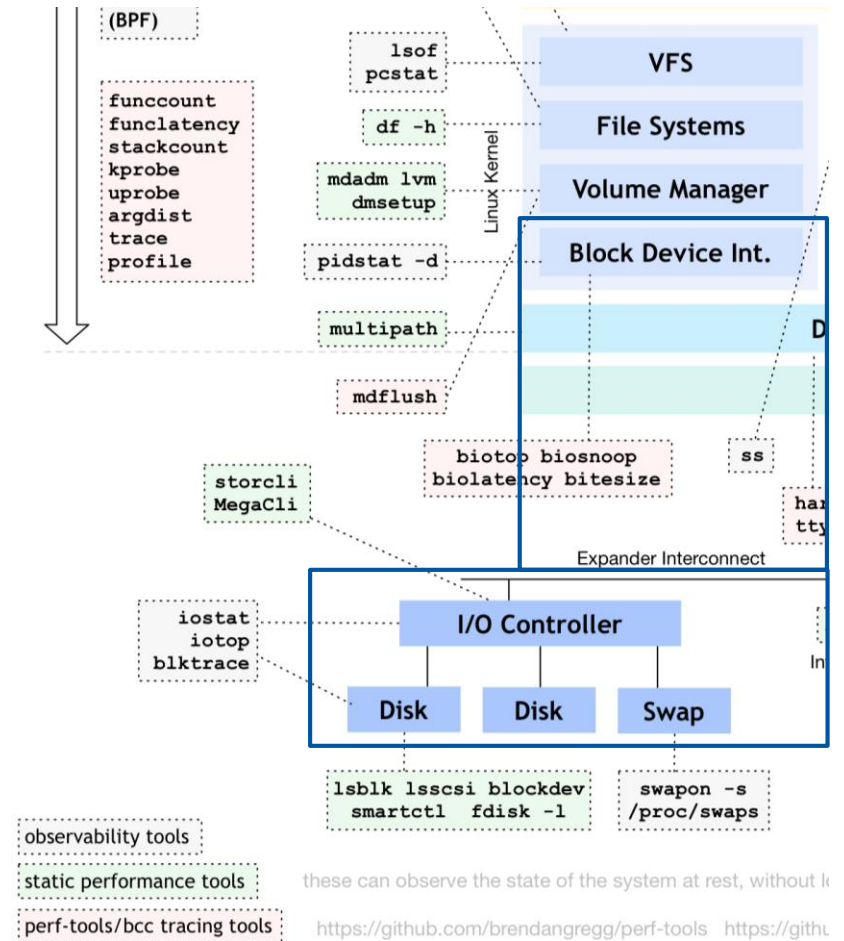
- ❑ Memory bandwidth usage
  - ❑ Amount of data read and written to main memory
- ❑ Memory transaction classification
  - ❑ Memory transaction occurs through pages from memory bank
  - ❑ Page-Hit
    - ❑ Memory bank in open state
    - ❑ Lowest access latency (Open)
    - ❑ Sequential memory access usually have high page hits
  - ❑ Page-Empty
    - ❑ Memory bank is idle and needs to be activated
    - ❑ Moderate access latency (Usually 2x of page hit)
    - ❑ Random memory access usually have high page empty counts
  - ❑ Page-Miss
    - ❑ Memory to be accessed requires closing of a page in the same bank
    - ❑ Worst access latency (Usually 3x of page hit)

Further info: <https://www.anandtech.com/show/3851/everything-you-always-wanted-to-know-about-sdram-memory-but-were-afraid-to-ask/5>



# Trident - IO Performance Analysis

- ❑ Data recorded from ProcFS
- ❑ Transfer Rate Analysis
  - ❑ Amount of data read and written to storage
  - ❑ Limited by interface and type of memory
  - ❑ Sustained Vs Bursts
- ❑ Operation Rate Analysis
  - ❑ Amount of operations performed
  - ❑ Limited by controller for fast memory
  - ❑ Limited by disk head for HDDs
  - ❑ Random Vs Sequential accesses



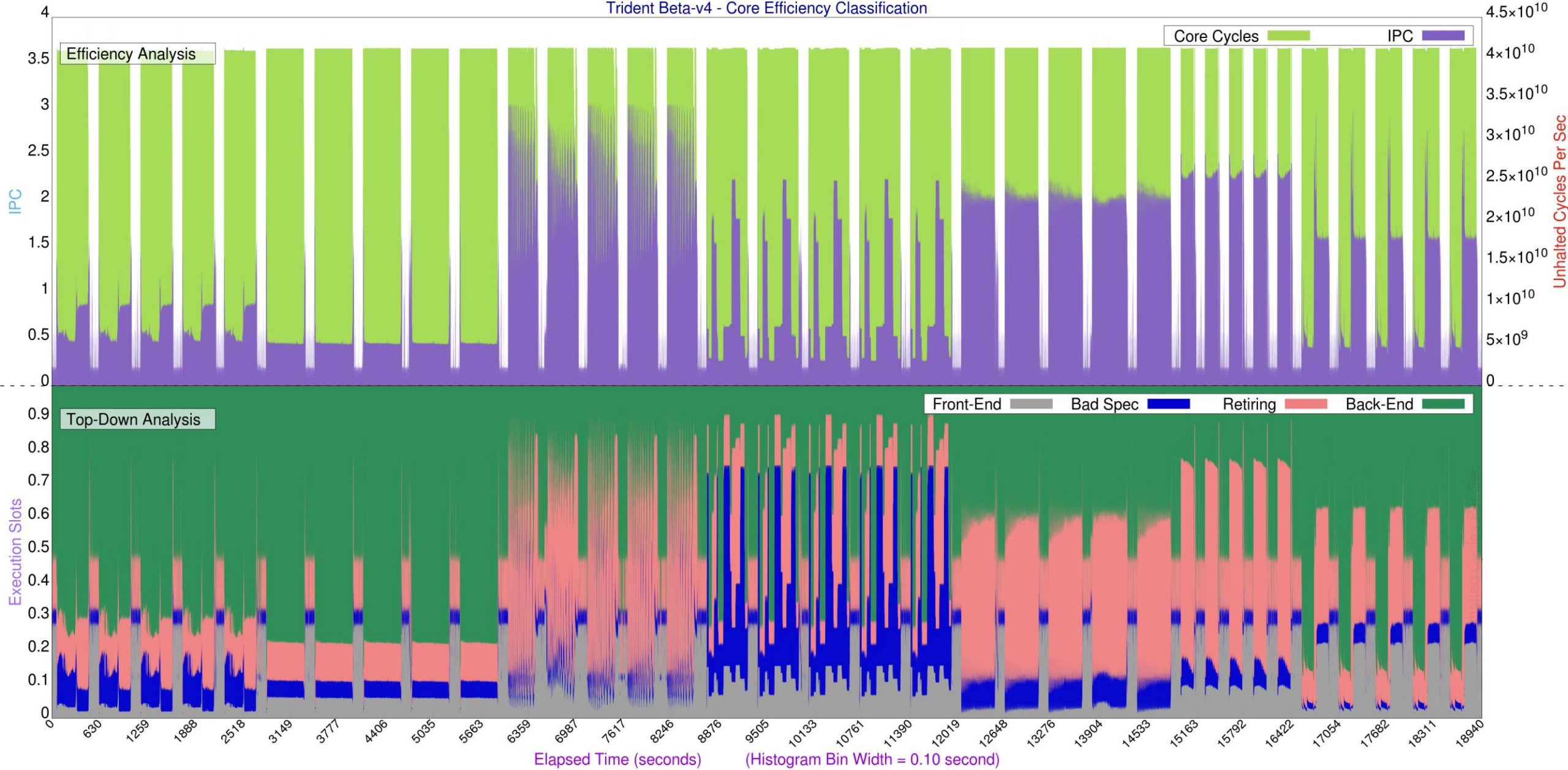
# Experiment setup

- ❑ Test System : 2x Intel(R) Xeon(R) E5-2630 v3@2.40GHz / 64GB DDR4-2133 RAM / Centos 7 (3.10.0-862.el7.x86\_64)
- ❑ Workloads
  - ❑ **HEPSPEC06 Benchmark Suite**
    - ❑ 450.soplex,471.omnetpp,447.dealll,473.astar,444.namd,453.povray,483.xalancbmk
  - ❑ ATLAS Job 1 - Geant4 MC simulation (CPU Intensive)
    - ❑ Executes: Sim\_tf (Geant4)
    - ❑ Input: EVNT files from event generation / Output: HITS files
  - ❑ ATLAS Job 2 - MC digitization and reco (CPU+I/O intensive)
    - ❑ Executes: Reco\_tf with several sub-steps
      - ❑ HITtoRDO (Digitization)
      - ❑ RDOtoRDOTrigger (Trigger simulation)
      - ❑ RAWtoESD + ESDtoAOD (MC reconstruction)
      - ❑ POOLMergeAthenaMPAOD0 (merge some small outputs)
    - ❑ Input: HITS files from MC simulation + low and high pile-up HITS for digitization / Output: AOD files
  - ❑ ATLAS Job 3 - Derivation production (I/O Intensive)
    - ❑ Executes: Reco\_tf in AODtoDAOD mode
      - ❑ AODtoDAOD
      - ❑ DAODs are used as the basis of ~all physics analysis in ATLAS
      - ❑ DAODs produced from AOD by removing whole events (skimming), whole reconstructed objects (thinning) and variables (slimming)
    - ❑ Input: AOD / Output: multiple DAODs in a train

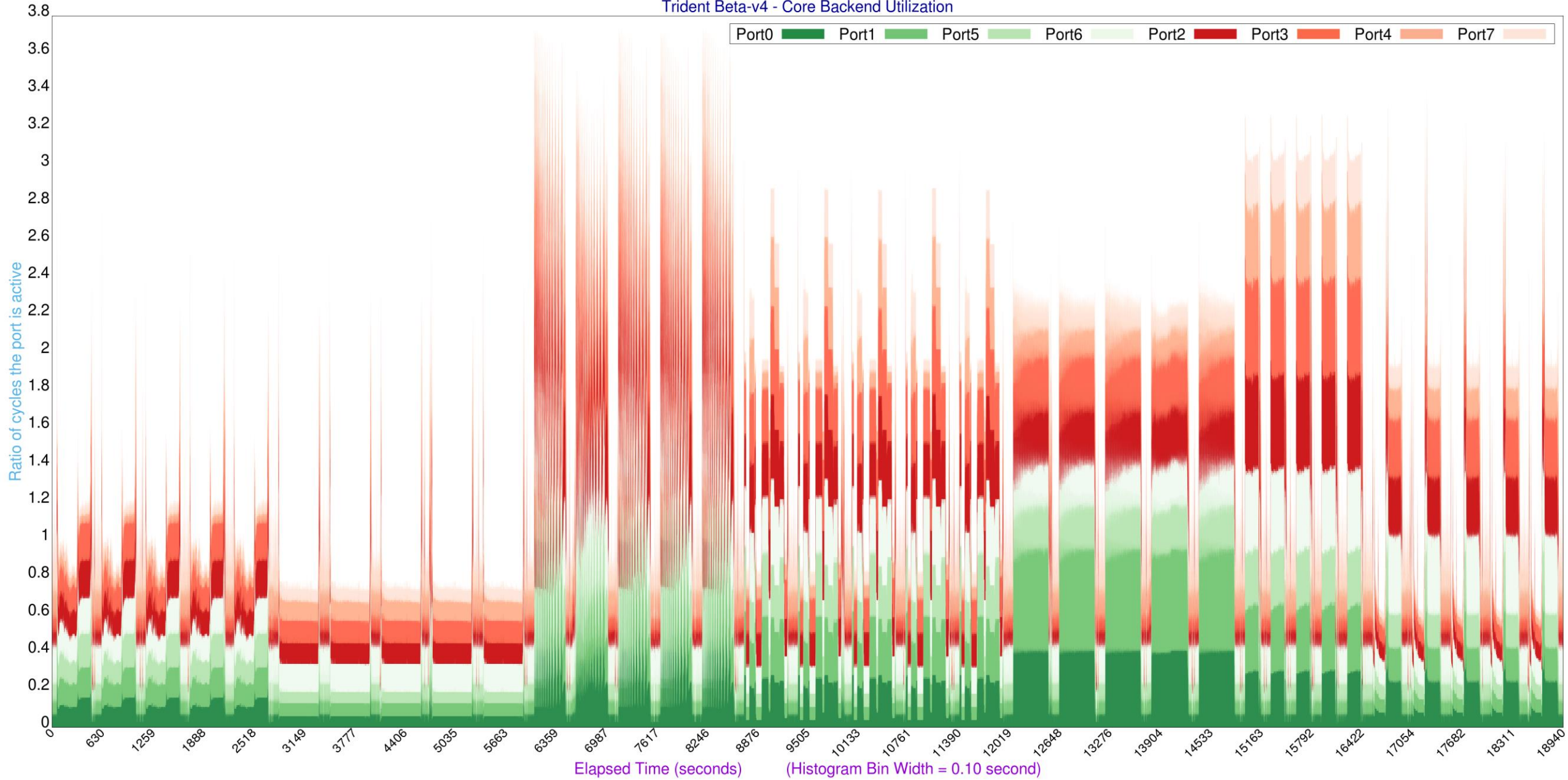


# HS06

## Trident Beta-v4 - Core Efficiency Classification

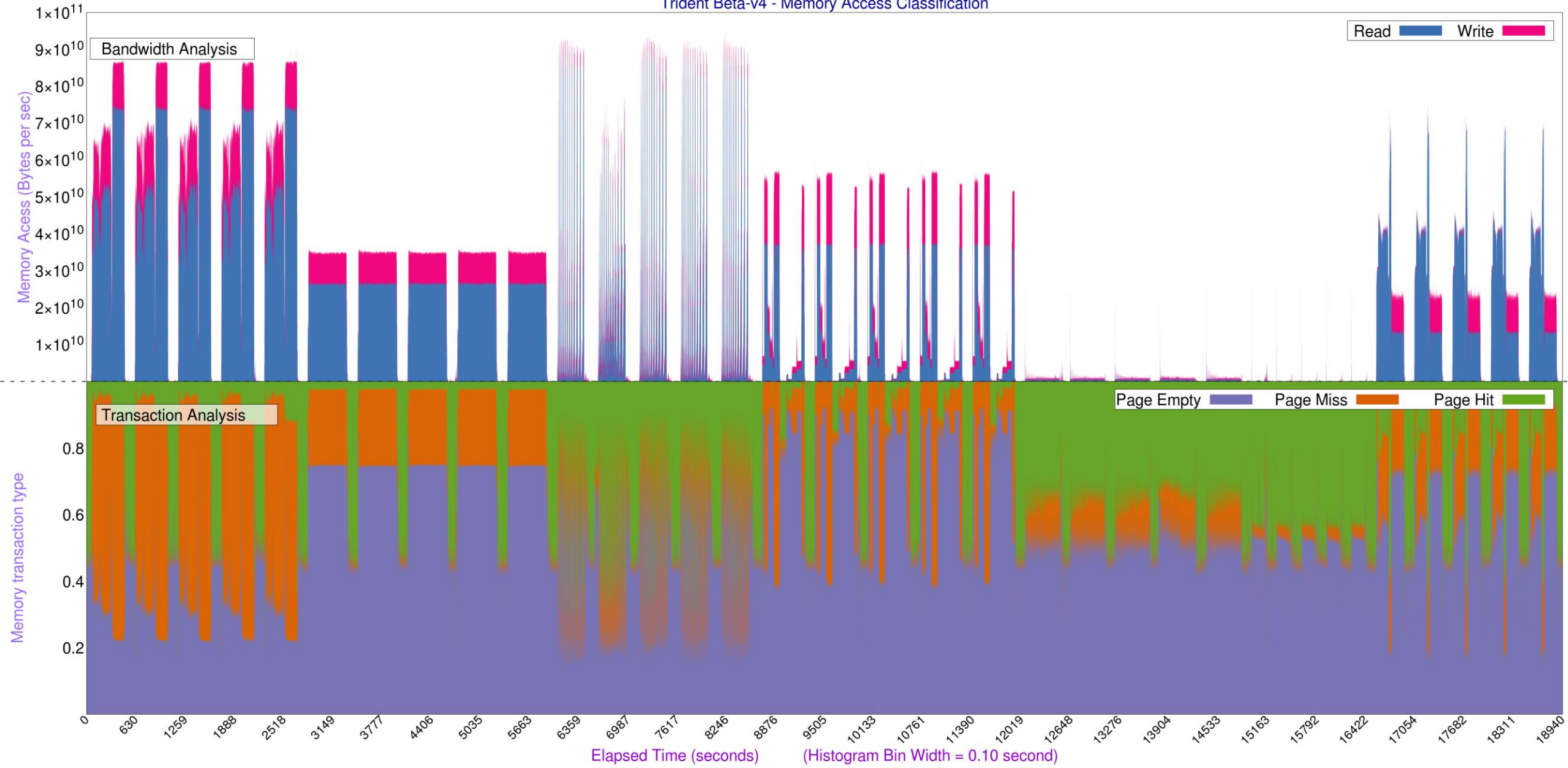


Trident Beta-v4 - Core Backend Utilization



# HS06

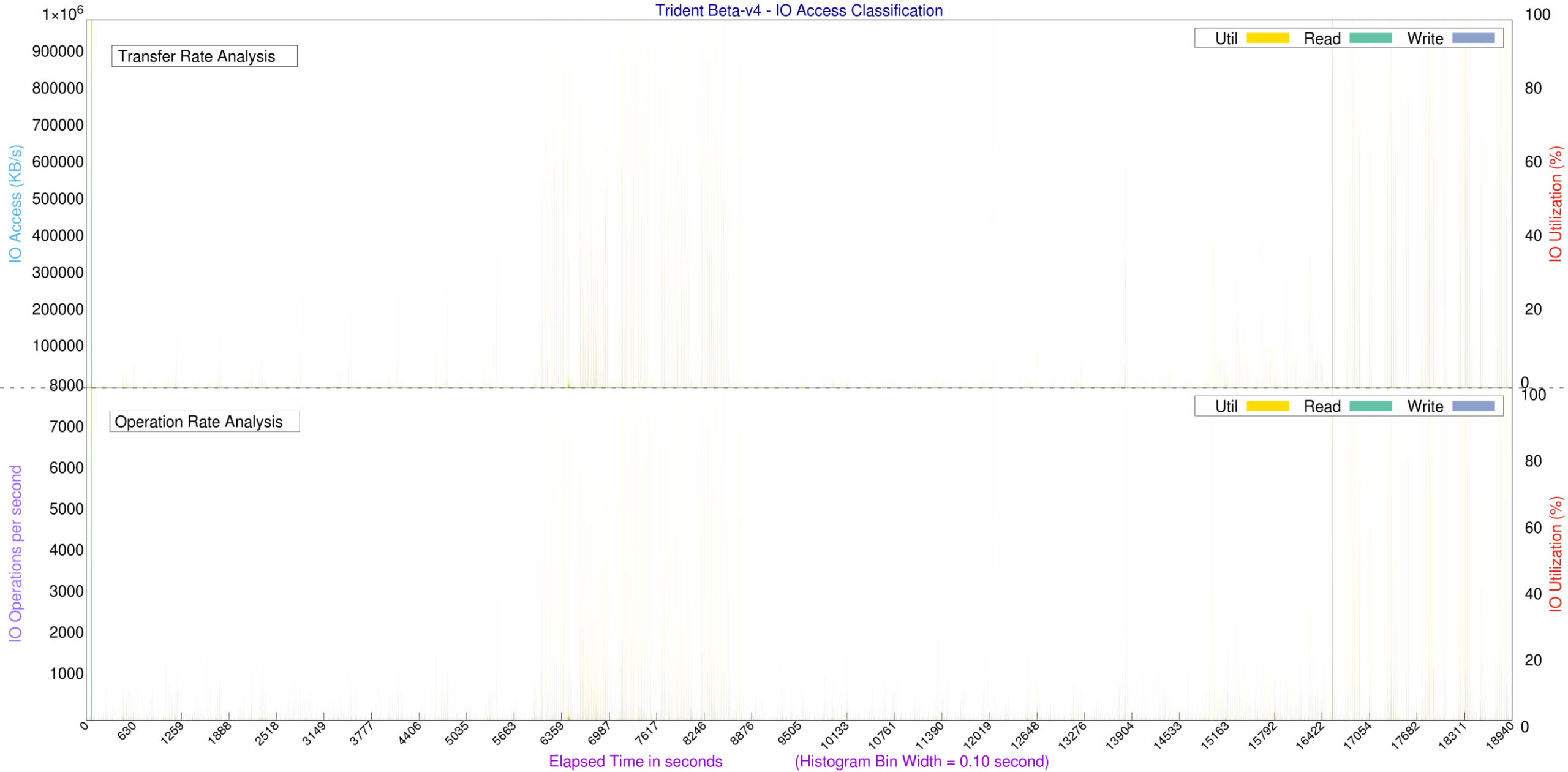
Trident Beta-v4 - Memory Access Classification





# HS06

Trident Beta-v4 - IO Access Classification

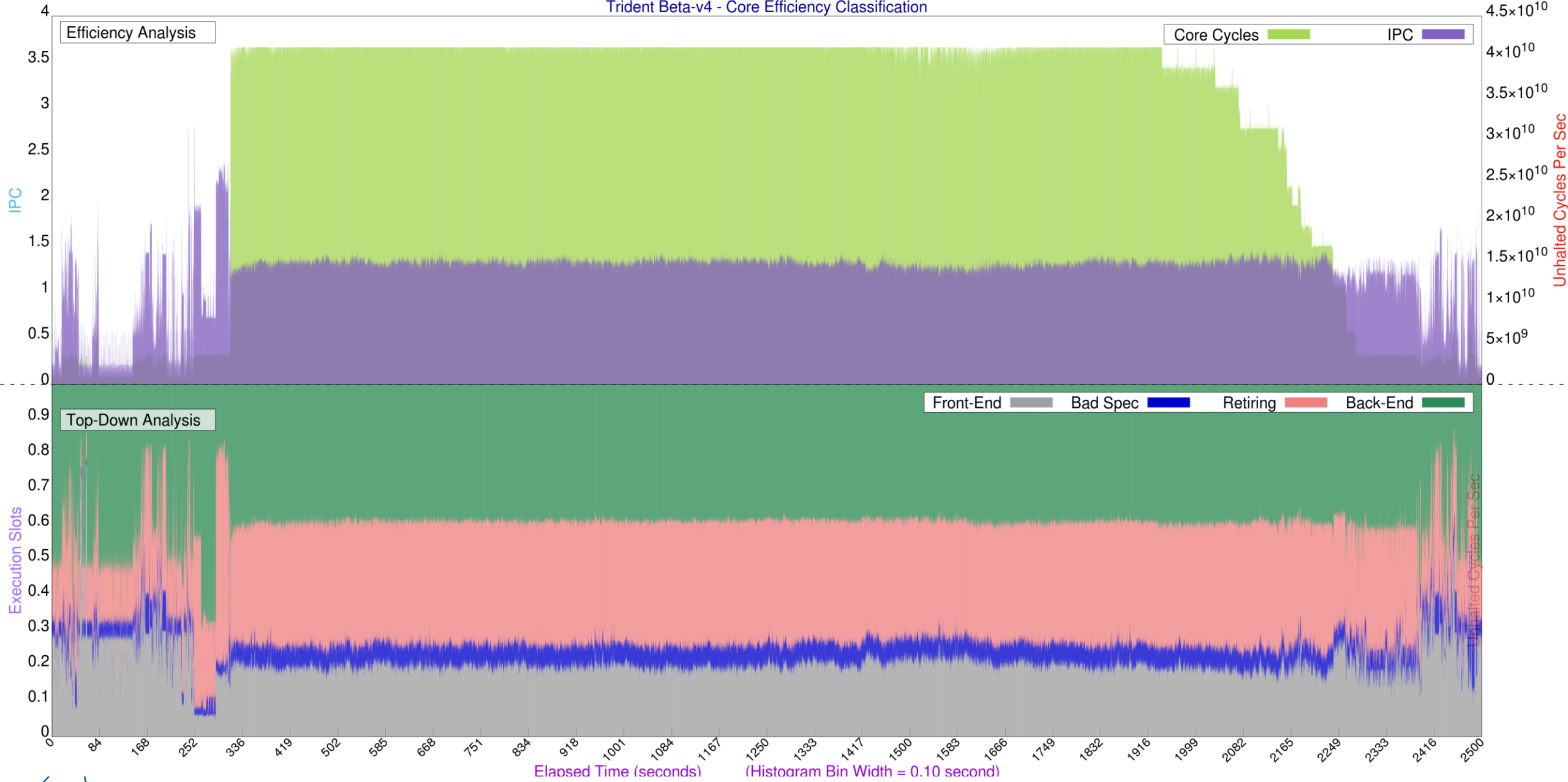


# Experiment setup

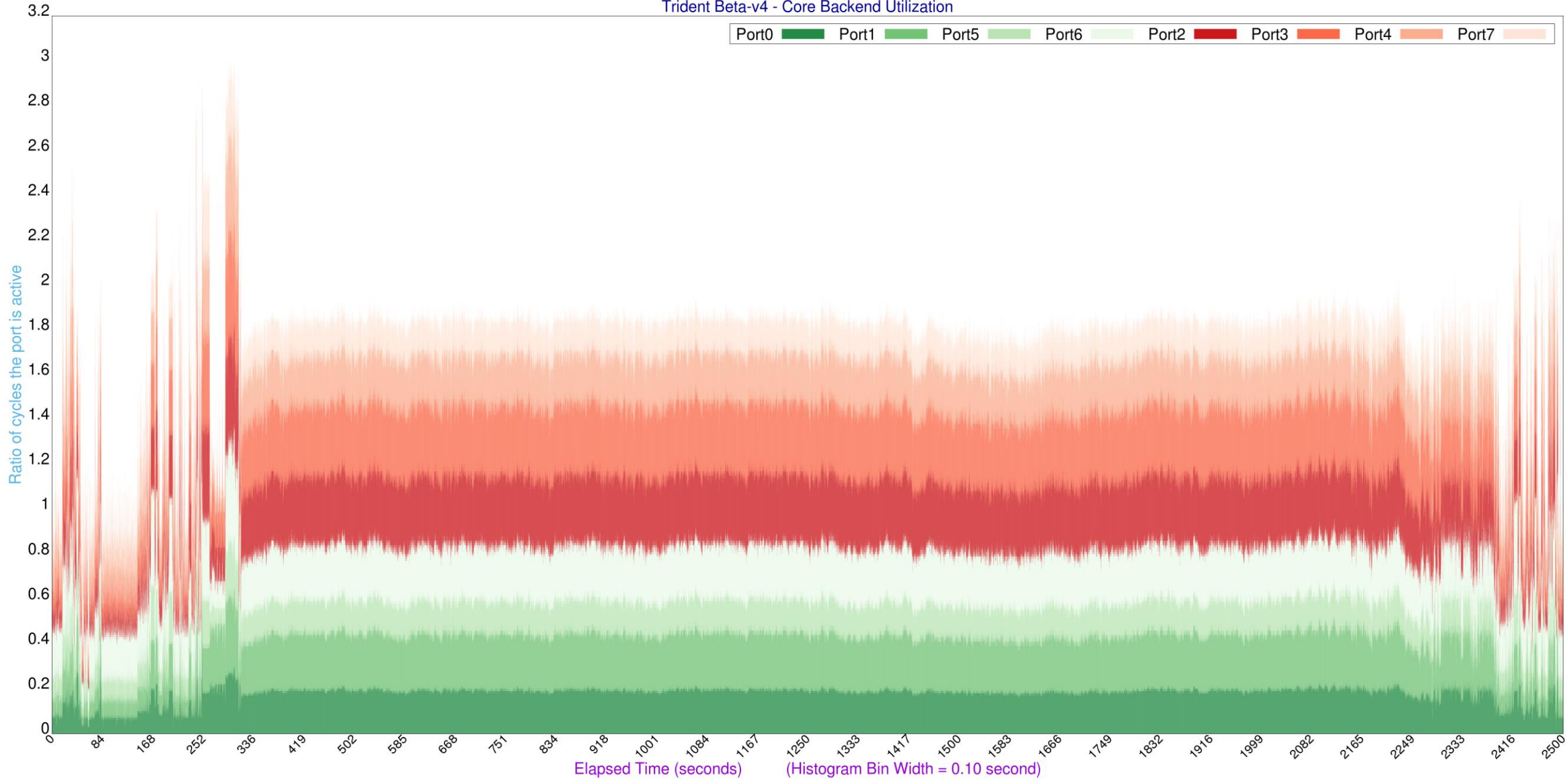
- ❑ Test System : 2x Intel(R) Xeon(R) E5-2630 v3@2.40GHz / 64GB DDR4-2133 RAM / Centos 7 (3.10.0-862.el7.x86\_64)
- ❑ Workloads
  - ❑ HEPSPROC06 Docker Container – Cloud Benchmark Suite (Thanks Domenico Giordano!!!)
    - ❑ 450.soplex,471.omnetpp,447.dealll,473.astar,444.namd,453.povray,483.xalancbmk
  - ❑ ATLAS Job 1 - Geant4 MC simulation (CPU Intensive)
    - ❑ Executes: Sim\_tf (Geant4)
    - ❑ Input: EVNT files from event generation / Output: HITS files
  - ❑ ATLAS Job 2 - MC digitization and reco (CPU+I/O intensive)
    - ❑ Executes: Reco\_tf with several sub-steps
      - ❑ HITtoRDO (Digitization)
      - ❑ RDOtoRDOTrigger (Trigger simulation)
      - ❑ RAWtoESD + ESDtoAOD (MC reconstruction)
      - ❑ POOLMergeAthenaMPAOD0 (merge some small outputs)
    - ❑ Input: HITS files from MC simulation + low and high pile-up HITS for digitization / Output: AOD files
  - ❑ ATLAS Job 3 - Derivation production (I/O Intensive)
    - ❑ Executes: Reco\_tf in AODtoDAOD mode
      - ❑ AODtoDAOD
      - ❑ DAODs are used as the basis of ~all physics analysis in ATLAS
      - ❑ DAODs produced from AOD by removing whole events (skimming), whole reconstructed objects (thinning) and variables (slimming)
    - ❑ Input: AOD / Output: multiple DAODs in a train

# Geant4

## Trident Beta-v4 - Core Efficiency Classification



Trident Beta-v4 - Core Backend Utilization

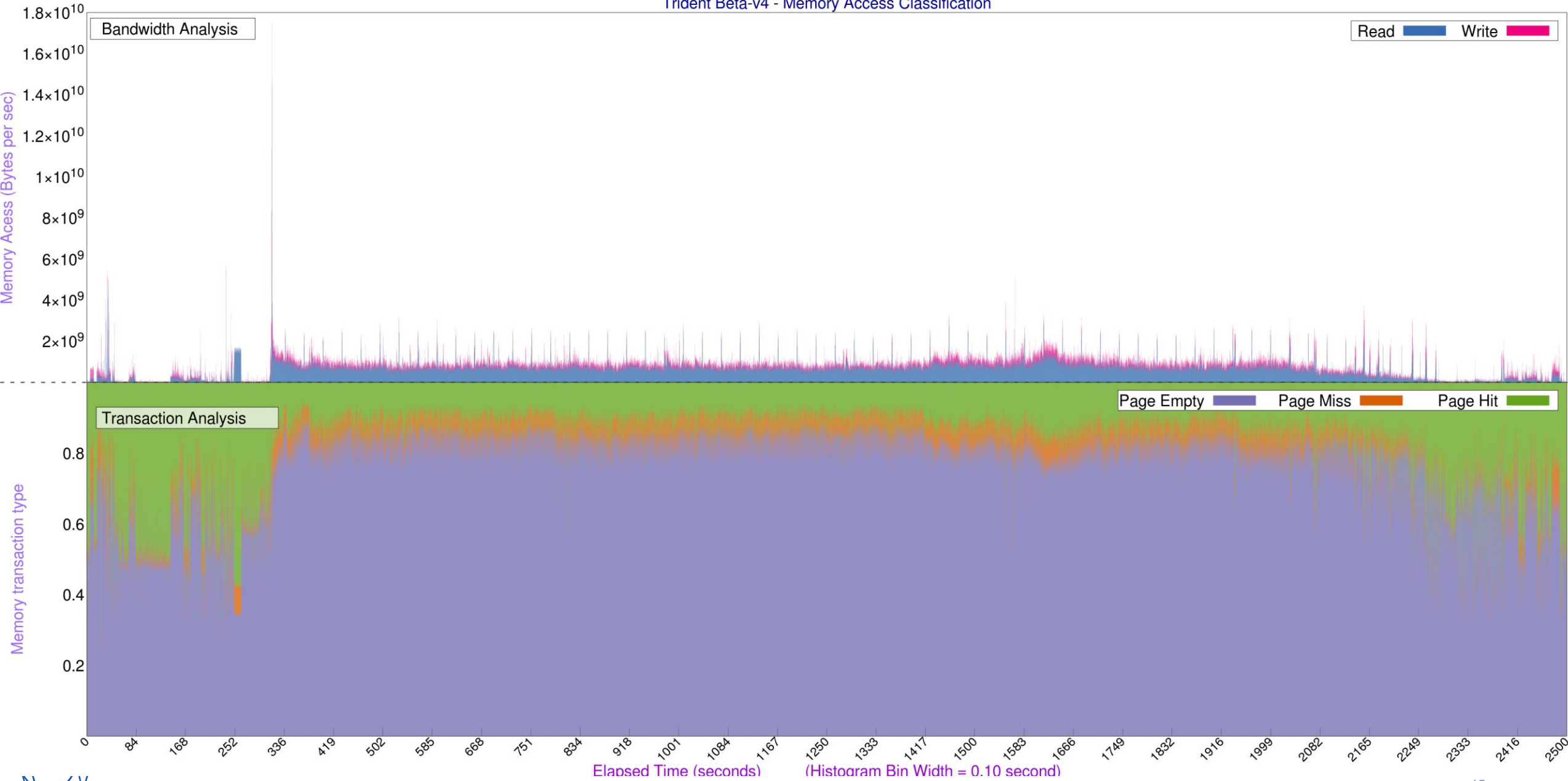


Elapsed Time (seconds) (Histogram Bin Width = 0.10 second)

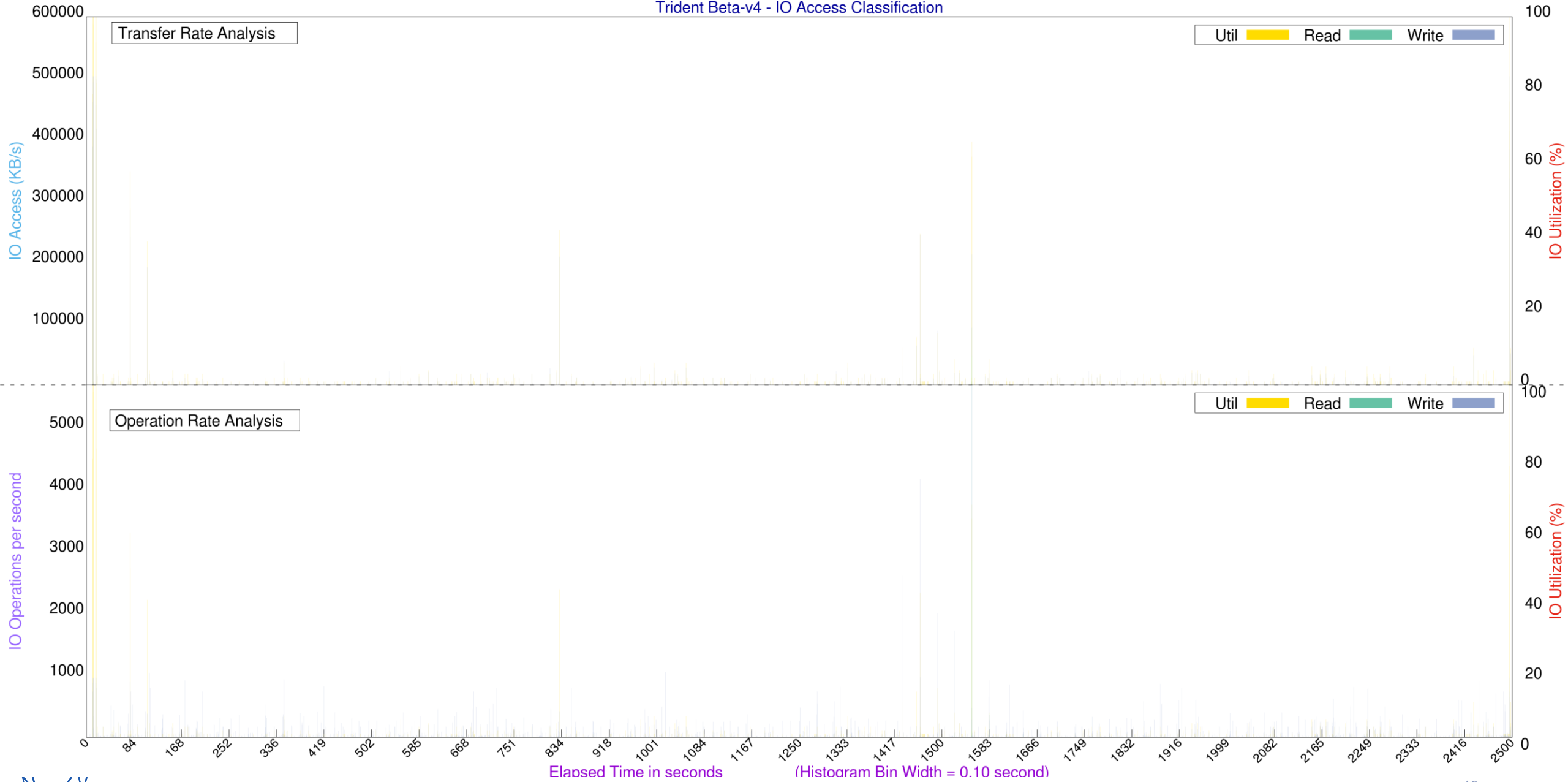


# Geant4

## Trident Beta-v4 - Memory Access Classification







# Experiment setup

- ❑ Test System : 2x Intel(R) Xeon(R) E5-2630 v3@2.40GHz / 64GB DDR4-2133 RAM / Centos 7 (3.10.0-862.el7.x86\_64)
- ❑ Workloads
  - ❑ HEPSPROC06 Docker Container – Cloud Benchmark Suite (Thanks Domenico Giordano!!!)
    - ❑ 450.soplex,471.omnetpp,447.dealll,473.astar,444.namd,453.povray,483.xalancbmk
  - ❑ ATLAS Job 1 - Geant4 MC simulation (CPU Intensive)
    - ❑ Executes: Sim\_tf (Geant4)
    - ❑ Input: EVNT files from event generation / Output: HITS files
  - ❑ ATLAS Job 2 - MC digitization and reco (CPU+I/O intensive)
    - ❑ Executes: Reco\_tf with several sub-steps
      - ❑ HITtoRDO (Digitization)
      - ❑ RDOtoRDOTrigger (Trigger simulation)
      - ❑ RAWtoESD + ESDtoAOD (MC reconstruction)
      - ❑ POOLMergeAthenaMPAOD0 (merge some small outputs)
    - ❑ Input: HITS files from MC simulation + low and high pile-up HITS for digitization / Output: AOD files
  - ❑ ATLAS Job 3 - Derivation production (I/O Intensive)
    - ❑ Executes: Reco\_tf in AODtoDAOD mode
      - ❑ AODtoDAOD
      - ❑ DAODs are used as the basis of ~all physics analysis in ATLAS
      - ❑ DAODs produced from AOD by removing whole events (skimming), whole reconstructed objects (thinning) and variables (slimming)
    - ❑ Input: AOD / Output: multiple DAODs in a train

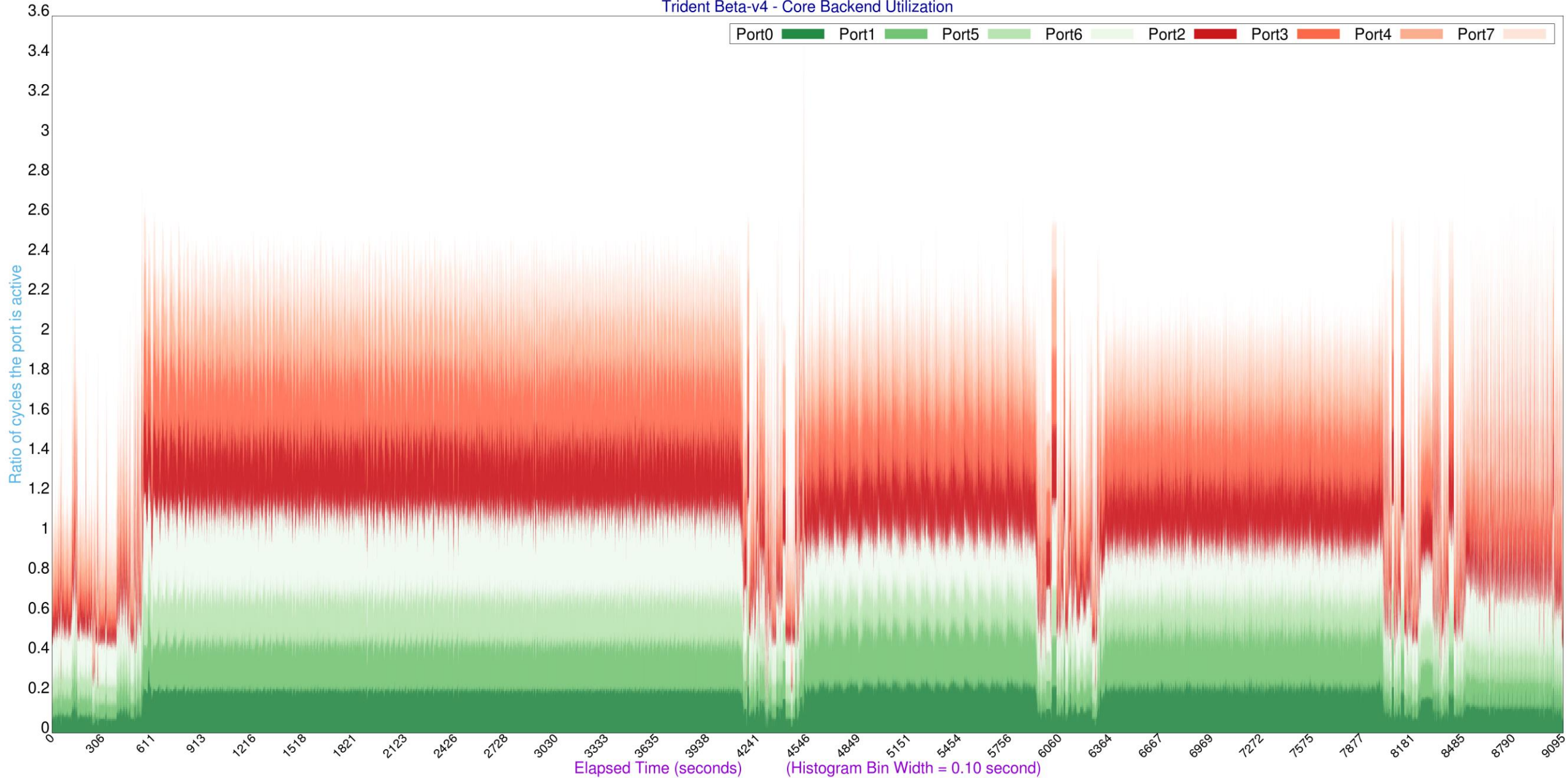
# MC Digi + Reco

## Trident Beta-v4 - Core Efficiency Classification



# MC Digi + Reco

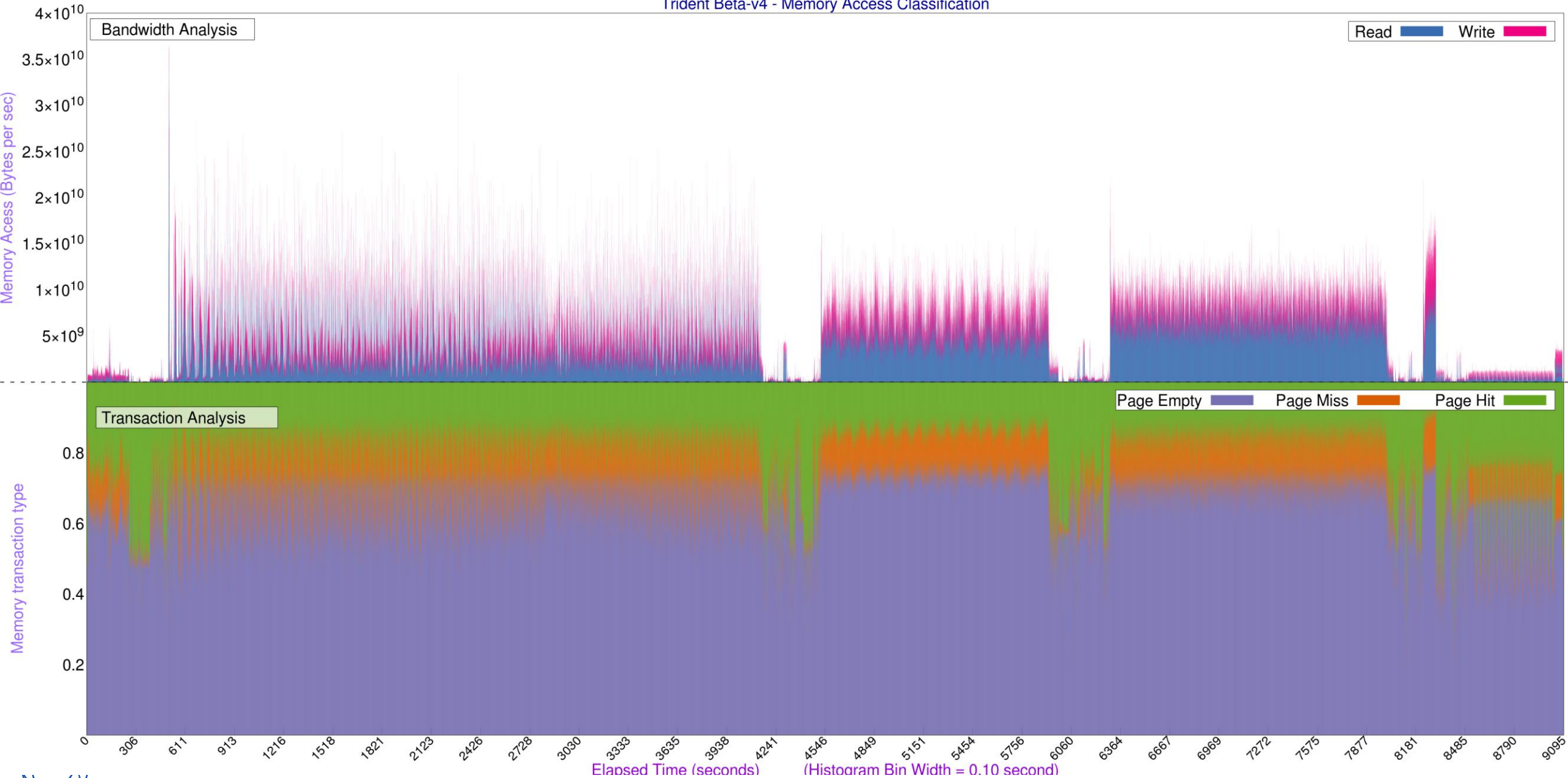
Trident Beta-v4 - Core Backend Utilization





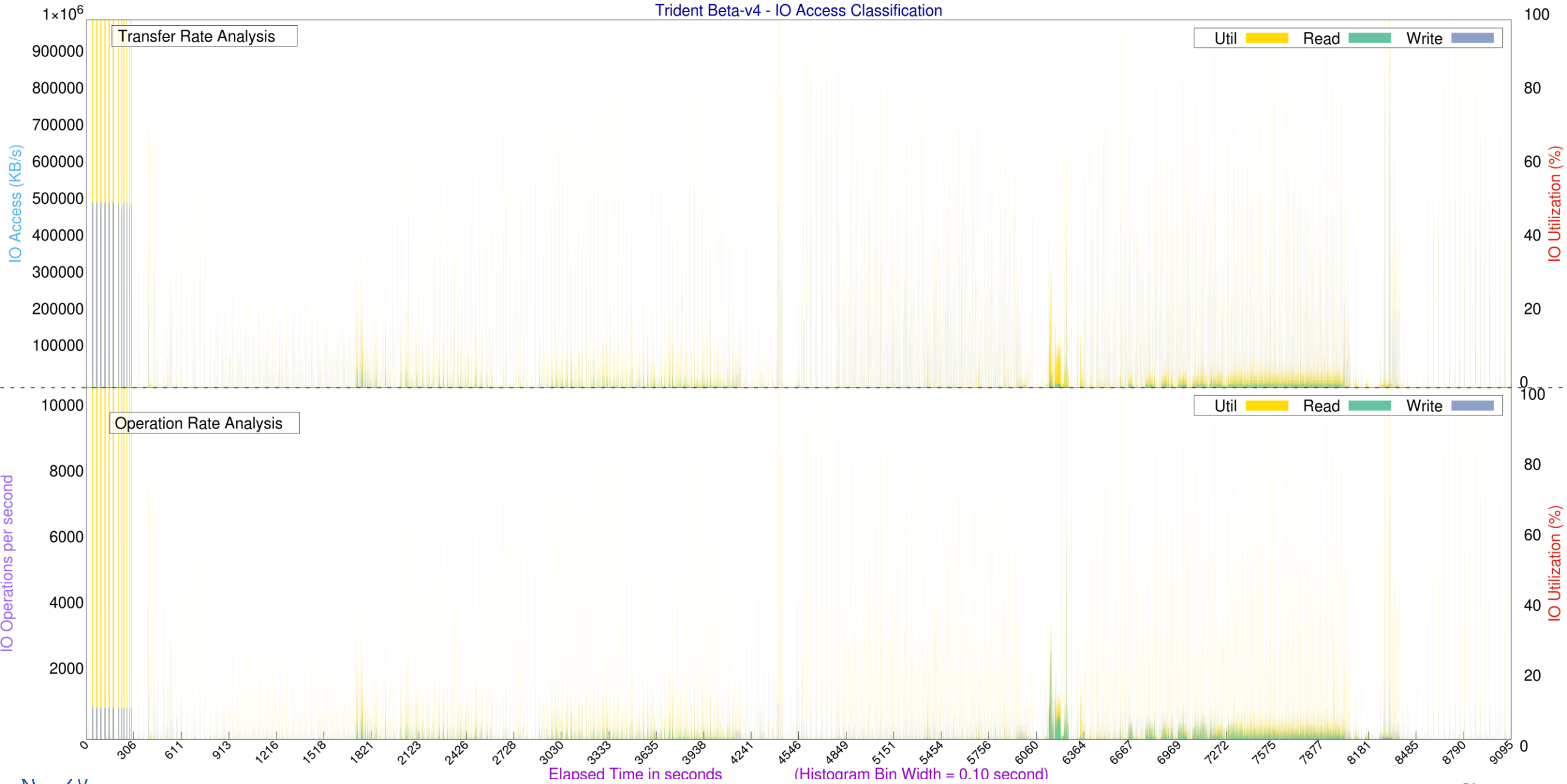
# MC Digi + Reco

Trident Beta-v4 - Memory Access Classification



# MC Digi + Reco

## Trident Beta-v4 - IO Access Classification

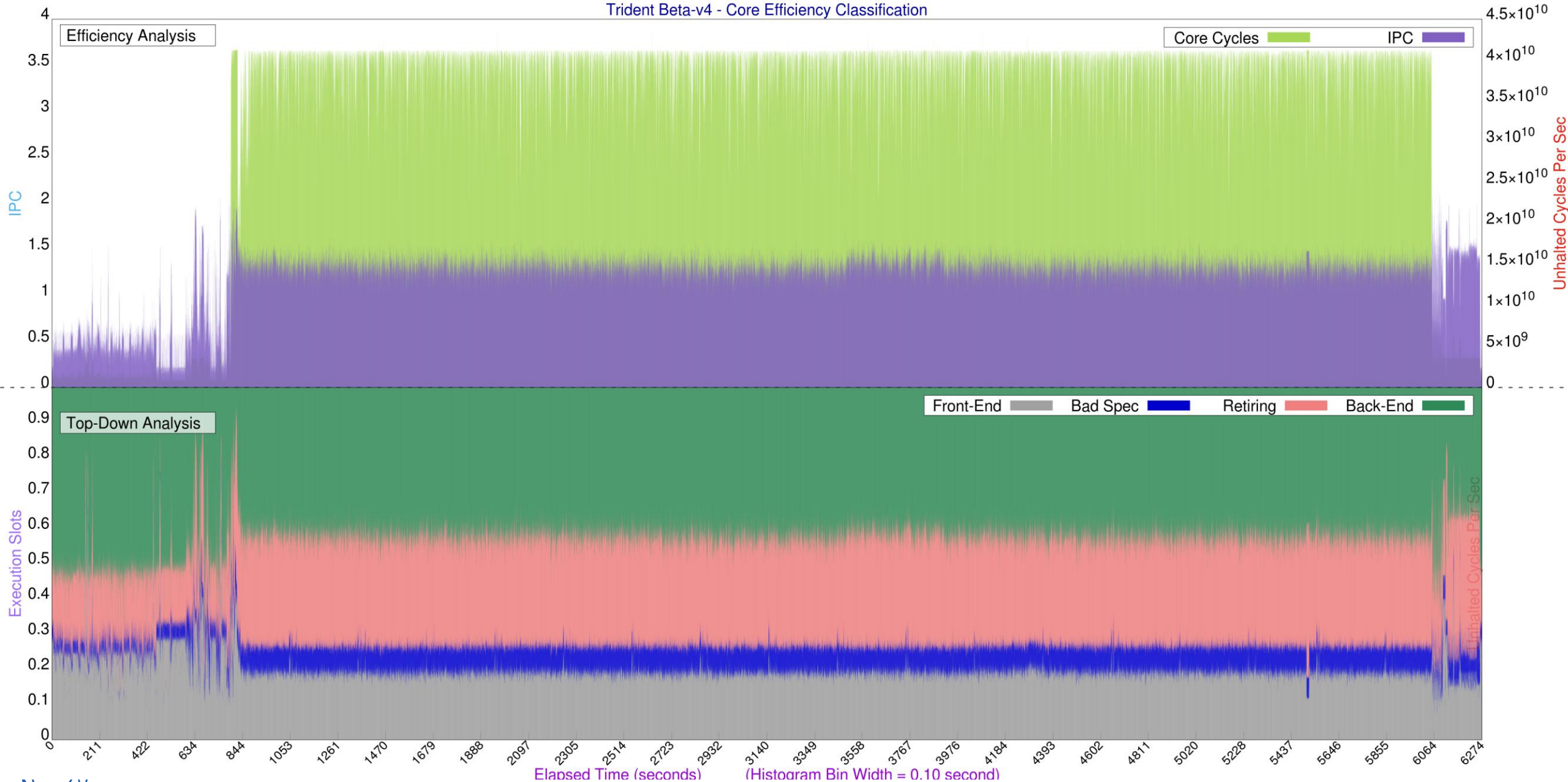


# Experiment setup

- ❑ Test System : 2x Intel(R) Xeon(R) E5-2630 v3@2.40GHz / 64GB DDR4-2133 RAM / Centos 7 (3.10.0-862.el7.x86\_64)
- ❑ Workloads
  - ❑ HEPSPROC06 Docker Container – Cloud Benchmark Suite (Thanks Domenico Giordano!!!)
    - ❑ 450.soplex,471.omnetpp,447.dealll,473.astar,444.namd,453.povray,483.xalancbmk
  - ❑ ATLAS Job 1 - Geant4 MC simulation (CPU Intensive)
    - ❑ Executes: Sim\_tf (Geant4)
    - ❑ Input: EVNT files from event generation / Output: HITS files
  - ❑ ATLAS Job 2 - MC digitization and reco (CPU+I/O intensive)
    - ❑ Executes: Reco\_tf with several sub-steps
      - ❑ HITtoRDO (Digitization)
      - ❑ RDOtoRDOTrigger (Trigger simulation)
      - ❑ RAWtoESD + ESDtoAOD (MC reconstruction)
      - ❑ POOLMergeAthenaMPAOD0 (merge some small outputs)
    - ❑ Input: HITS files from MC simulation + low and high pile-up HITS for digitization / Output: AOD files
  - ❑ ATLAS Job 3 - Derivation production (I/O Intensive)
    - ❑ Executes: Reco\_tf in AODtoDAOD mode
      - ❑ AODtoDAOD
      - ❑ DAODs are used as the basis of ~all physics analysis in ATLAS
      - ❑ DAODs produced from AOD by removing whole events (skimming), whole reconstructed objects (thinning) and variables (slimming)
    - ❑ Input: AOD / Output: multiple DAODs in a train

# Derivation Production

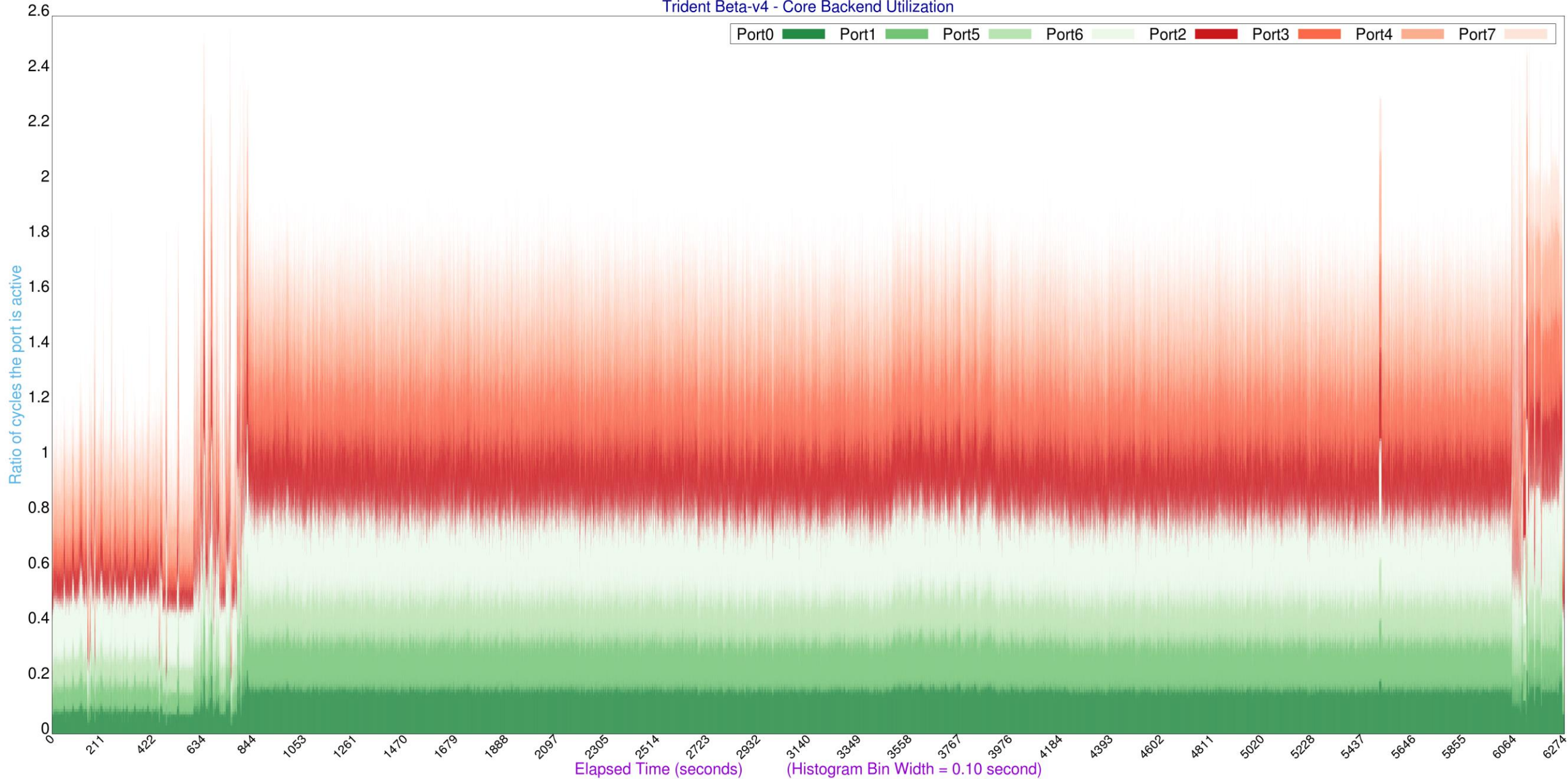
Trident Beta-v4 - Core Efficiency Classification





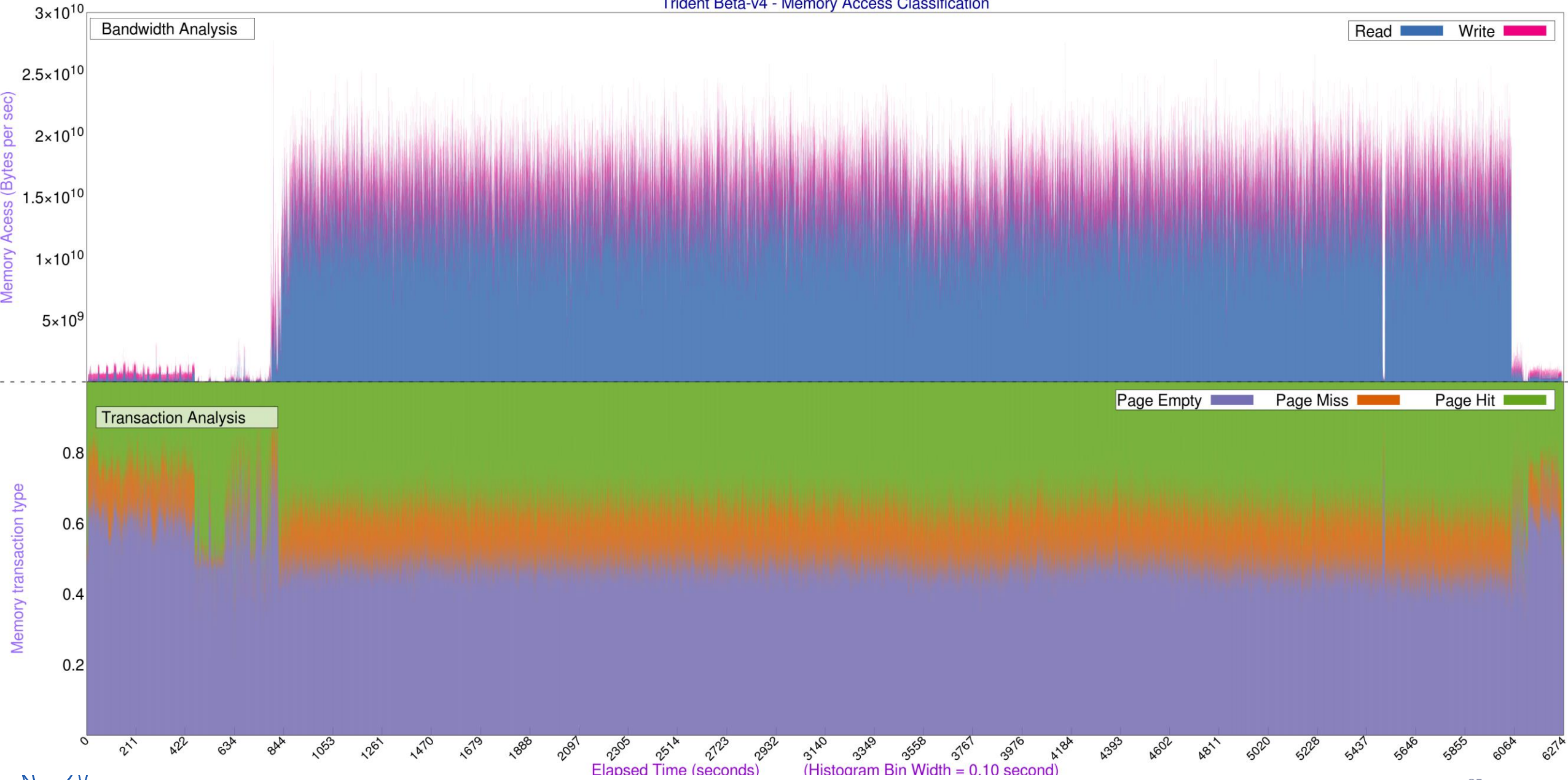
# Derivation Production

Trident Beta-v4 - Core Backend Utilization



# Derivation Production

Trident Beta-v4 - Memory Access Classification



# Derivation Production

Trident Beta-v4 - IO Access Classification

