# PMU and Performance Analysis Methods With Proposals for RISC-V in Cloud

Jiankang Chen
Shuai Xue
Huaixin Chang

Alibaba Cloud
6/26/2025

# CONTENTS

- Challenges and Methods of Performance Analysis
- CPU Bottleneck Analysis and Multi-architecture Support
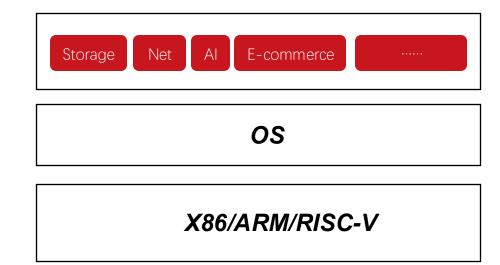- Uncore PMU Analysis and Unified PMU Support

# Challenges and Methods of Performance Analysis

# Challenges and Methods of Performance Analysis

Challenges of cloud scenario performance analysis

- Need to analyze complex business
  - Storage
    - Memory and IO intensive
  - E-commerce
    - Hybrid deployment
- The business program calling process is complex
  - The software stack involves many modules and components.
- Multi-architecture hardware
  - X86, ARM, RISC-V, GPU, etc
- Need low-overhead real-time monitoring of performance metrics

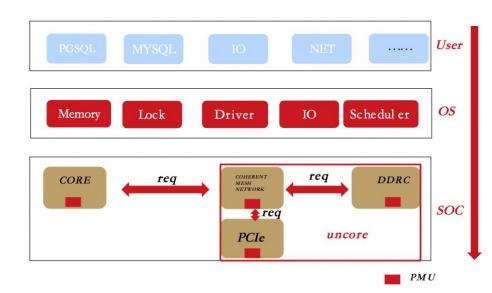| Storage | Net | AI | E-commerce | ...... |
|---|---|---|---|---|

| *OS* |
|---|

| *X86/ARM/RISC-V* |
|---|

# Challenges and Methods of Performance Analysis

Performance analysis methods for cloud computing scenarios

- From applications to OS, and then to hardware self-directed performance analysis methods
  - Perform full stack analysis through topdown

- Low overhead monitoring via hardware pmu in cloud

- Provides full-architecture analysis tools to solve multi-architecture problems

- Find bottlenecks that affect business performance through underlying hardware metric analysis
  - Use Core pmu to analyze the bottleneck points in the core
  - Analyze application bottlenecks at the soc architecture level through uncore pmu
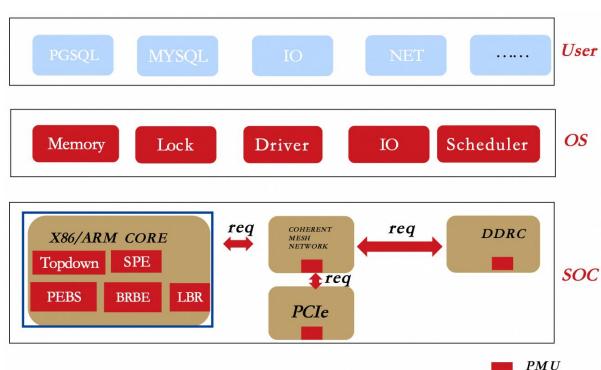
# CPU Bottleneck Analysis and Multi-architecture Support

# CPU Bottleneck Analysis and Multi-architecture Support

We use the following methods to analyze:
- **Applying the topdown analysis model**
  - Topdown for performance monitoring and analysis by perf and pas tool;

- **Multi-platform analysis framework:**
  - Essential for diverse cloud environments(X86, AMD, ARM, RISC-V)
  - PAS

- **Using advanced core pmu features:**
  - **BRBE**: BRBE is used to perform autoFDO to optimize database services by reducing frontend bound.
  - **SPE/PEBS**: Fine-grained metric analysis, such as instruction behavior and time latency, etc.
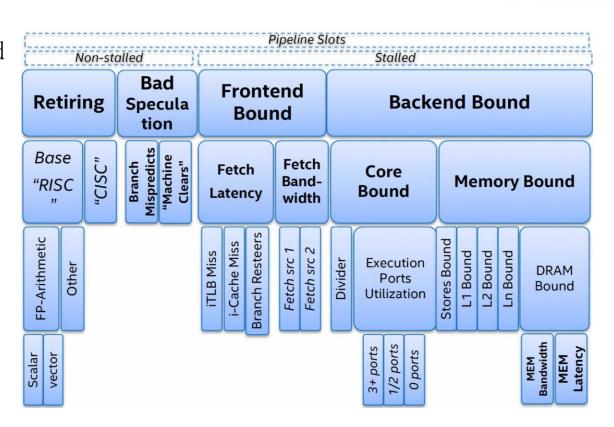


*To solve the bottlenecks of the program on the core*

# Topdown Analysis Model – X86

Frontend Bound:This metric is the percentage of total slots that were stalled in the frontend of the processor.

- Fetch Latency: CPU was stalled due to front-end latency issues;
  - Itlb miss, icache miss, branch resteers
- Fetch Bandwidth: CPU was stalled due to front-end bandwidth issues;
  - Inefficiencies in the instruction decoders;
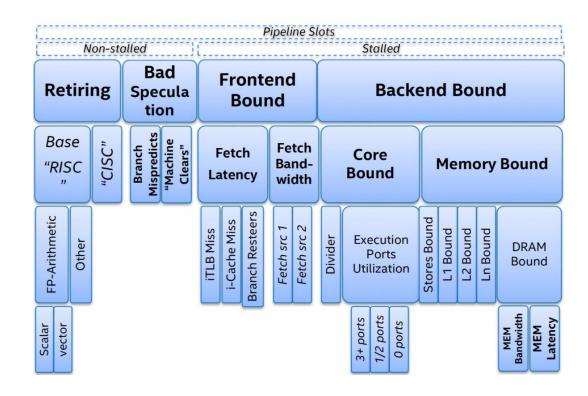  - Code restrictions for caching in the DSB;

# Topdown Analysis Model – X86

Backend Bound: No uOps are being delivered due to a lack of required resources for accepting new uOps in the Back-End.

- Core Bound: Shortage in hardware compute resources, or dependencies software's instructions
  - Certain execution units are overloaded;
  - Dependencies in program's data or instruction flow are limiting the performance
- Memory Bound:
  - Store bound: CPU was stalled on store operations
  - Cache miss bound: L1, L2, L3;
  - Dram bound
    - Memory bandwidth
    - Memroy latency

Pipeline Slots

Non-stalled | Stalled

| Retiring | Bad Speculation | Frontend Bound | Backend Bound |

| Base "RISC" | "CISC" | Branch Mispredicts | "Machine Clears" | Fetch Latency | Fetch Bandwidth | Core Bound | Memory Bound |

FP-Arithmetic | Other | iTLB Miss | i-Cache Miss | Branch Resteers | Fetch src 1 | Fetch src 2 | Divider | Execution Ports Utilization | Stores Bound | L1 Bound | L2 Bound | Ln Bound | DRAM Bound

Scalar | vector | 3+ ports | 1/2 ports | 0 ports | MEM Bandwidth | MEM Latency
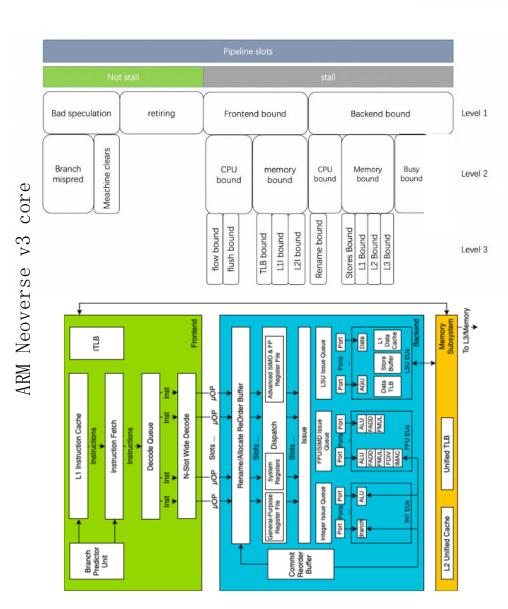
# Topdown Analysis Model – ARM

Frontend Bound: This metric is the percentage of total slots that were stalled in the frontend of the processor.

- Memory Bound: Cycle stalled by TLBi miss, L1i miss and L2i miss;

- CPU Bound: Due to frontend core resource constraints not related to instruction fetch latency issues caused by memory access components.

  - Flow Bound: Decode unit is awaiting input from the branch prediction unit.

  - Flush Bound: Recovering from a pipeline flush caused by bad speculation
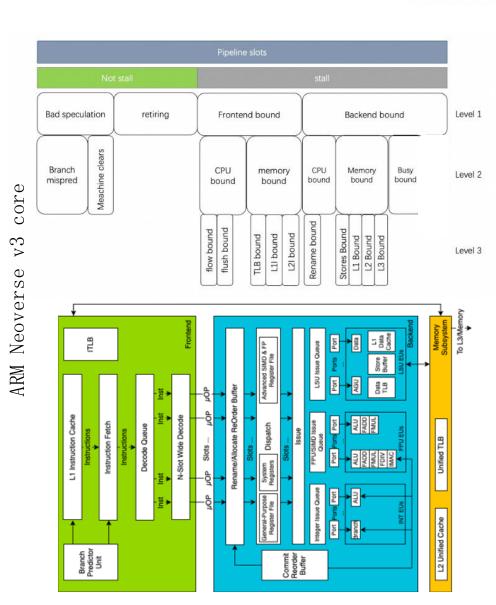


ARM Neoverse v3 core

# Topdown Analysis Model – ARM

Backend Bound: No uOps are being delivered due to a lack of required resources for accepting new uOps in the Back-End

- Memory Bound: Data cache miss  and slc miss;

- CPU Bound: Because of unavaliable rename regs to produce stall

- Busy Bound: Due to issue queues being full to accept operations



ARM Neoverse v3 core

# Multi-platform Analysis Framework -- PAS

Why need pas in cloud?

▪ Multiple architectures exist in cloud

▪ Need a lightweight analytical framework, perf is too heavy

▪ Requires fine-grained analysis at the instruction level

▪ Need to integrate different analytical capabilities to meet multiple business needs

PAS is designed to support multiple platforms

• Use Standard System Interfaces to get pmu data with low overhead

• Supports ARM/Intel/AMD, expecting RISC-V platforms

• FrameScope is the On-CPU analysing tool with improved observability by providing function-level topdown;

• RTRadar is the tool to achieve OS-level full observability, which conduct offcpu analysis.

• SysAdvisor is a tool to check and compare performance configuarations

| Micro-benchmark | MySQL | Java | | | | GPU | |
| | | Specjbb | Spark | Taobao | Nginx | training | inference |
| | | Dubbo | | | Redis | | |
| PAS | | | | | | GPU Analysis | |
| SysAdvisor | | FrameScope | | RTRadar | | Computing Analysis | |
| | | | | | | IO/Network Analysis | |
| System Interface | | | | | | | |
| Sysfs | bootcmdline | Kconfig | Perf record | Perf probe | Java Agent | Lttng | |
| Hardware | | | | | | | |
| ARM | | Intel | | AMD | | GPU | |
| Topdown | SPE | Topdown | PEBS | Topdown | IBS | DPU | |
| CoreSight | | Intel PT | | AMD ITT | | RISC-V | |

# Use Cases for PAS in Cloud



Both functions call the flush_tlb_page function:

```c
static inline void flush_tlb_page(struct vm_area_struct *vma,
                unsigned long uaddr)
{
    flush_tlb_page_nosync(vma, uaddr);
    dsb(ish);
}


static inline void flush_tlb_page_nosync(struct vm_area_struct *vma,
                unsigned long uaddr)
{
        unsigned long addr;

        dsb(ishst);
        addr = __TLBI_VADDR(uaddr, ASID(vma->vm_mm))

        __tlbi(vale1is, addr);
        __tlbi_user(vale1is, addr);
}
```
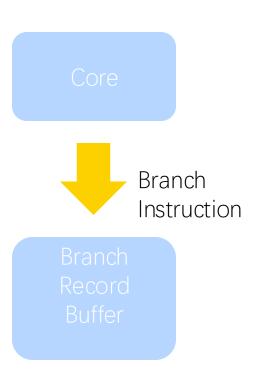
- For the instruction: __tlbi(vale1is, addr)
  - va: virtual address
  - l: last level, refers to PTE。
  - is: inner shareable, invalidate all cores TLB in inner shareable
- Ptep_clear_flush has a high backend bound due to tlb flushing;

# Advanced Core PMU Features -- Branch Record Buffer Extension

- We will use BRBE to optimize the frontend bound of database applications with autoFDO.

- Objective

  - Capture a recent sequence of branches in an easy-to-consume format

- Requirements

  - Low compute and memory overheads for capture/analysis;

  - Low overheads while recording

  - Low cost of context switch and on reading

- Branch Record

  - Source VA of taken branch/exception

  - Target VA of taken branch/exception

- Branch Record Buffer

  - EL1 feature, for recording EL0 and/or EL1;

- For risc-v the smctr (Control Transfer Records extension) is the same as ARM BRBE;

Core

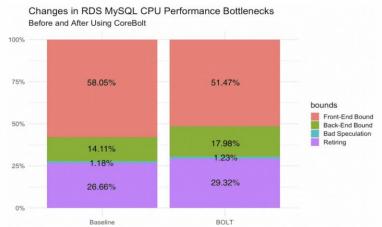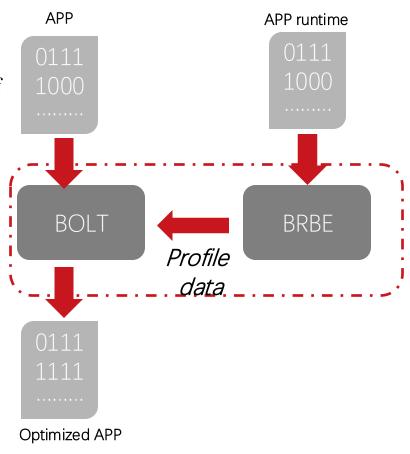Branch Instruction

Branch Record Buffer

# Practice of the BRBE In Cloud

We optimize database services by combining Bolt and BRBE
- Bolt's optimization principle
  - Bolt(binary optimization and layout tool)is a post-link optimizition tool
  - We collects the profile data( branchs track of program) of business programs during runtime;
  - Bolt using this profile data to directly modifies and optimizes the binary files;
  - Reducing the jump instructions on the critical execution path;
- Comparison of TPS of RDS MySQL:*read_only increased by 20%, read_write increased by 16%, and write_only increased by 22%.*

APP

APP runtime

0111
1000
.........

0111
1000
.........

BOLT

BRBE

*Profile data*

0111
1111
.........

Optimized APP

Changes in RDS MySQL CPU Performance Bottlenecks
Before and After Using CoreBolt



bounds
- Front-End Bound
- Back-End Bound
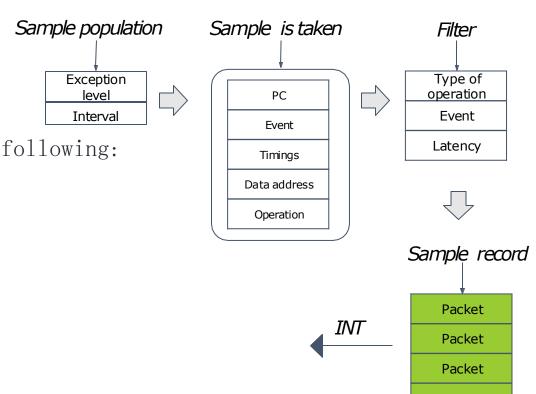- Bad Speculation
- Retiring

*Compared with before optimization, frontend bound decreased*

# Advanced Core PMU Features -- Statistical Profiling Extension

- **PMU Defects:** PMU Due to the existence of interrupts, the PMU sampling instruction will deviate from the actual instruction position

- SPE advantage: can perform instruction-level accurate sampling

  - The sampling process of SPE can be divided into the following:
    - 1. Choose an operation
    - 2. Collect data about the operation
    - 3. Optionally discard the record based on a filter
    - 4. Write the record to memory
    - 5. Interrupt when the buffer is full

- SPE trace the processor pipeline
  - Record operations (memory, exception, SVE, etc)
  - Gather associated information for the operation (PC value, data address, event type, timestamp, etc.)

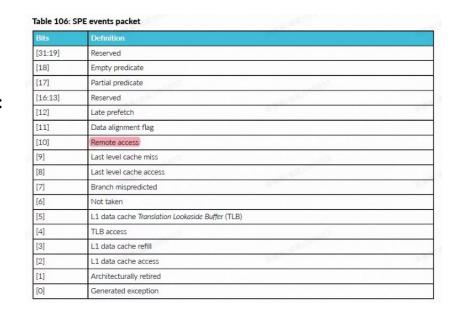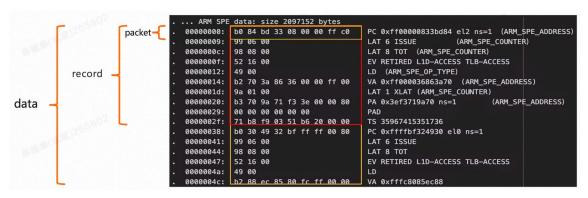*Sample population*

| Exception level |
| Interval |

*Sample is taken*

| PC |
| Event |
| Timings |
| Data address |
| Operation |

*Filter*

| Type of operation |
| Event |
| Latency |

*Sample record*

| Packet |
| Packet |
| Packet |
| …… |

*INT*

SPE Data Recording Format
- One record represents the tracking record of an operation
- One record contains multiple packets.
- Different packets represent different information, including:
  - Address: Virtual PC
  - Event Source: Where in the system the data was returned from. E.g. Level 1 cache, Last Level cache, or another socket in a multi-socket system;
  - ISSUE: Number of cycles taken from Dispatch to Issue
  - TOT: Number of cycles taken from Dispatch to Complete
  - XLAT: Number of cycles spent on pagetable walk
  - Events: Events occur:  L1D-ACCESS
  - Type: Operation type
  - Address: Virtual address that instruction access
  - Timestamp
- ARM supports multiple events in patckets;
- The Operation type in packet:
  - Load/store/branch/sve,etc instructions

Table 106: SPE events packet

| Bits | Definition |
| --- | --- |
| [31:19] | Reserved |
| [18] | Empty predicate |
| [17] | Partial predicate |
| [16:13] | Reserved |
| [12] | Late prefetch |
| [11] | Data alignment flag |
| [10] | Remote access |
| [9] | Last level cache miss |
| [8] | Last level cache access |
| [7] | Branch mispredicted |
| [6] | Not taken |
| [5] | L1 data cache *Translation Lookaside Buffer* (TLB) |
| [4] | TLB access |
| [3] | L1 data cache refill |
| [2] | L1 data cache access |
| [1] | Architecturally retired |
| [0] | Generated exception |

# Practice of SPE in Cloud

- DDRC statistics only count the latency of requests within DDRC;
- The ddr latency calculated by DDRC PMU may be low.

- Therefore, we can monitor the memory access latency by SPE;

- We know that spe can count the three latency of load instructions:;
  - Total latency(total lat);
  - Issue latency(issue lat);
  - Xlat latency (xlat);

- Memory latency equals total latency minus address translation and issue latency

  mem lat = total lat - xlat – issue lat

*DDR latency*



| 0b00000 | Total latency. Cycle count from the operation being dispatched for issue to the operation being complete. Included for all operations. |
| 0b00001 | Issue latency. Cycle count from the operation being dispatched for issue to the operation being issued for execution. This counts any delay in waiting the operation being ready to issue. Included for all operations. |
| 0b00010 | Translation latency. Cycle count from a virtual address being passed to the MMU for translation to the result of the translation being available. Included for all load, store and atomic operations. |

*Disadvantages of this approach:  The tool needs to be customized;*
*Demand : A new independent PMU is needed to count latency, for example, AMD has the XI pmu*
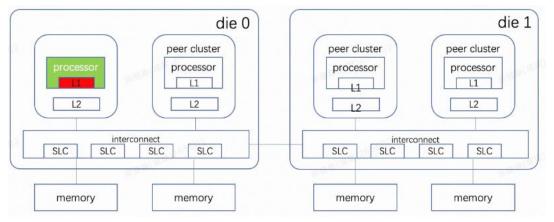
Packet Data Source：

- N2 TRM: Data source: Refers to where does core get data from;

- There are currently 8 data sources in spec, but the number of implementations varies for different SOC.
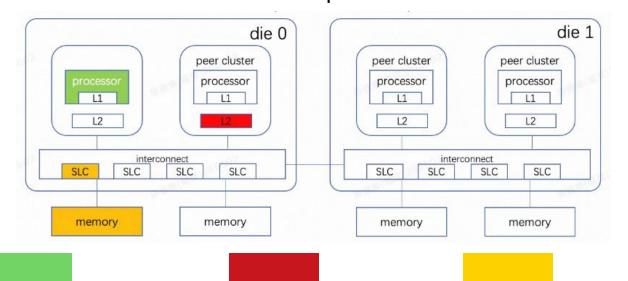
Table 22-2: SPE data source packet

| Value | Name |
|---|---|
| 0b0000 | L1 data cache |
| 0b1000 | L2 cache |
| 0b1001 | Peer core |
| 0b1010 | Local cluster |
| 0b1011 | System cache |
| 0b1100 | Peer cluster |
| 0b1101 | Remote |
| 0b1110 | Dynamic Random Access Memory (DRAM) |

- N2: One core per cluster

**data source is l1 cache**



**data source is peer cluster.**



the core requesting data    new data    old data
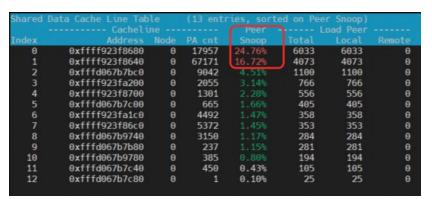
# Practice of SPE in Cloud

Another important use case is detecting false sharing in X265 by SPE data source

False sharing: multiple threads concurrently modifying variables within the same cache line lead to performance degradation;

We use perf c2c cmd to collect the sampling data of X265
- In the perf c2c statistics, each row here represents a cacheline.
- The peer snoop field describes the percentage of peer snoop.
- A high peer snoop value indicates false sharing has occurred.
- It can be seen that the addresses where false sharing occurs are mainly 0x688118 / 0x688104
- Through spe, we can find the code line where false sharing occurs and solve the problem
- After optimization, performance increased by 8.6%



perf c2c statistics          Detailed information corresponding to cache line 0xffff923f8680          Hot code

# Uncore PMU Analysis and Unified PMU Support

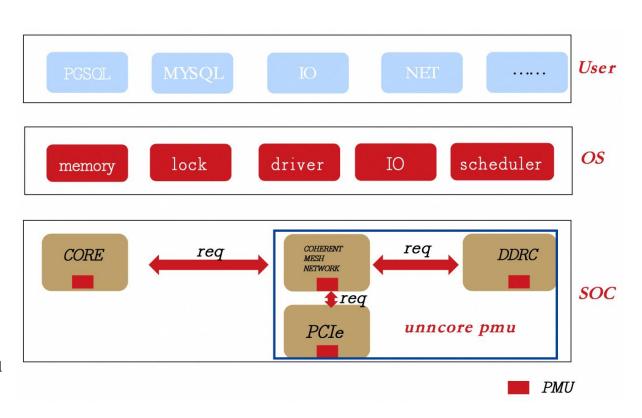# Uncore PMU Analysis and Unified PMU Support

The main usage scenarios of uncore pmu in cloud are as follows:
- **Low-overhead monitoring**
- **Bandwidth and Latency:** Cloud services focus on bandwidth and latency between components.
  - Inter-die bandwidth
  - Inter-socket bandwidth
  - Memory bandwidth
  - Memory latency
  - Some IO and network services also focus on PCIe bandwidth;
- **Last level cache pmu events**

Currently the key components of uncore pmu are:
- CMN PMU
- PCIe PMU
- DDRC PMU
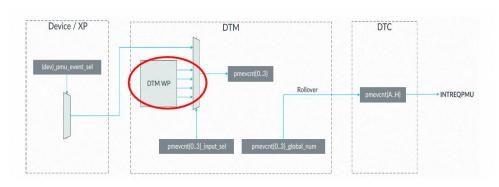
ARM propose a new uncore pmu architecture
- Use coresight pmu replace all uncore pmu
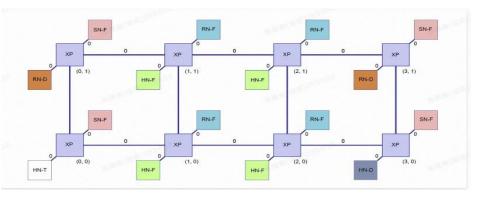


22

# Bandwidth and Cache Analysis of CMN PMU

In ARM,CMN PMU is mainly used to analyze cross-die data bandwidth and last level cache.

- CMN is composed of many different types of nodes in a cpu die
    - XP: Crosspoint - A switch or router logic module.
    - RN-F(a CPU core), HN-F (SLC cache module), SN-F(connect Memory Controller)
- The CMN PMU consists of DTM (local in xp) and DTC(global in die) components;
    - The DTM is located in each XP and has 4 local counters that can handle up to 4 events of the connected device nodes and the XP itself;
    - DTC has only 2 counters per CMN instance and 8 global counters;
    - The local counter in DTM overflows and carries into the global counter, and support overflow interrupt.
    - Each node supports different events, for example, HNF supports last level cache related events.
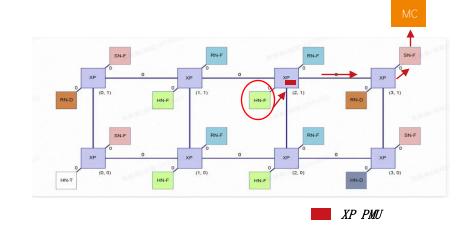
# Bandwidth and Cache Analysis of CMN PMU

- HNF contains last level cache miss rate events and some memory request events

- HNF SLC cache miss rate can be obtained through the following two events

  - PMU_HN_CACHE_MISS_EVENT (0x1)

    - Counts the total cache misses.

  - PMU_HNSLC_SF_CACHE_ACCESS_EVENT (0x2)

    - The total number of cache accesses.

  - Metric:

  *Cache miss rate (%) = Total cache misses / Total cache accesses*
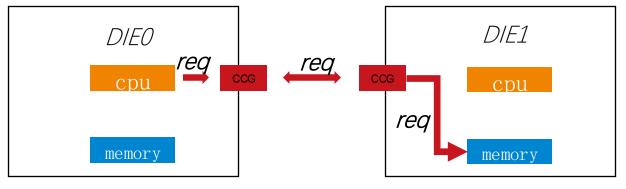


■ *XP PMU*

- If have a cache miss, hnf can send a memroy request

- If HNF resend memory requests frequently, the memory controller has become a bottleneck.

- MC message retry rate mainly include two events:

  - PMU_HN_MC_RETRIES_EVENT Number of requests that are retried to the memory controller.

  - PMU_HN_MC_REQS_EVENT Total number of requests that are sent to the memory controller.

  - metric:

  *MC message retry rate (%) = MC total messages retried / MC total messages received x 100*

# Bandwidth and Cache Analysis of CMN PMU

- CCG connects different CPU dies and is responsible for exchanging data between CPU dies.

- The most important events for CCG are to count the data bandwidth across the die.

- Data Bandwidth events metric:
  - Inter-Socket RX DATA Bandwidth = Flit_Size * CCG_WP0_RX_DATA_FLIT * CMN_Freq / PMU_CYCLE_COUNTER
  - Inter-Socket TX DATA Bandwidth = Flit_Size * CCG_WP0_TX_DATA_FLIT * CMN_Freq / PMU_CYCLE_COUNTER
  - Inter-Die RX DATA Bandwidth = Flit_Size * CCG_WP0_RX_DATA_FLIT * CMN_Freq / PMU_CYCLE_COUNTER
  - Inter-Die TX DATA Bandwidth = Flit_Size * CCG_WP0_TX_DATA_FLIT * CMN_Freq / PMU_CYCLE_COUNTER
  - Flit_Size = 32 Bytes

- CCG_WP0_RX_DATA_FLIT: The CMN PMU can count the number of data flit passing through CCG

- PMU_CYCLE_COUNTER: CPU cycle spent transferring data

# Bandwidth and Cache Analysis of CMN PMU

- In actual use, we will monitor the real-time data bandwidth across the die using the CCG XP PMU.
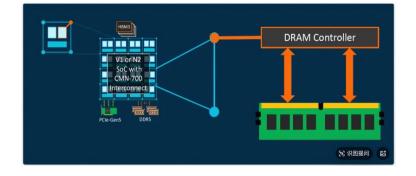- The following figure shows the real-time monitoring of the business:

# Memory bandwidth Analysis Via DDRC PMU

- For DDRC PMU, we often use it to monitor the DDR memory bandwidth of the entire machine.
- In addition, we will monitor the memory bandwidth usage of each machine in the cloud in real time online.
- And the main events we monitor include the following:
  - Memory total bandwith
  - memory read bandwidth
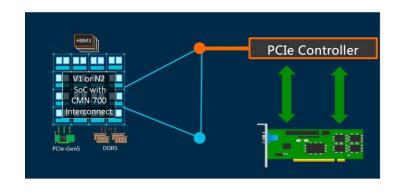  - memory write bandwidth;
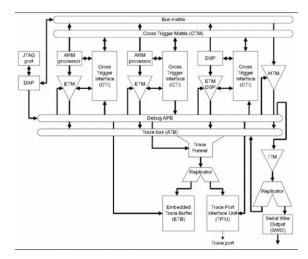
# Device Bandwidth Analysis of PCIe PMU

- In cloud, we also use PCIe PMU to count the bandwidth of IO and Net on the PCIe link.
- PCIe controller PMU provides:
  - Time Based Analysis (RX/TX data throughput and time spent in each state)
  - Lane Event counters (Error and Non-Error for lanes)
- **Time Based Analysis**: the counters are 64-bit width and measure data in two categories,
  - Percentage of time does the controller stay in each state in a configurable duration
  - Amount of data processed (Units of 16 bytes).
- **Lane Event counters**: Event (Error and Non-Error) across different layers of PCIe protocol, a 32 bit counter
  - Multiple event counters group, and each group include multiple events;
- For PCIe PMU, the events we commonly use are TLP Bandwidth, including
  - PCIe RX Bandwidth: PCIE_RX_DATA * 16B / Measure_Time_Window
  - PCIe TX Bandwidth: PCIE_TX_DATA * 16B / Measure_Time_Window

```
Performance counter stats for 'system wide':

    681,297,119      pcie_bdf_200/Rx_PCIe_TLP_Data_Payload/

   20.515887961 seconds time elapsed

>>> int("681,297,119".replace(',', '')) * 16/1000/1000
10900 MB
```

*FIO*

PCIe Controller

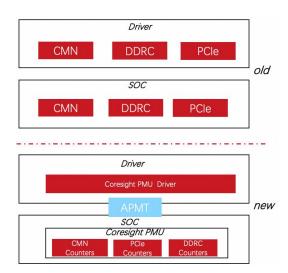V1 or N2 SoC with CMN-700 Interconnect

PCIe-Gen5   DDR5

# Use Coresight PMU to Uniformly Manage Uncore

- Traditional uncore pmu architectural features
  - In the traditional uncore pmu architecture, all pmu components are independent;
  - Poor maintainability and scalability, as every new uncore architecture requires a unique driver.

- New CoreSight PMU architecture
  - All uncore pmu events are managed by a coresight pmu hardware component;
  - There is only one coresight pmu driver in the driver layer;
  - Different uncore pmu component events can be extended in the same driver

- Introduced APMT new acpi table for extending uncore pmu nodes;
  - APMT(ARM Performance Monitoring Unit Table)
  - APMT supports multiple node types:
    - Memory controller, smmu, pcie, cpu cache



*Coresight can be used to trace branch and memory access instructions.*



*Uncore PMU new mode*

# RISC-V proposal

- RFC Proposal1:

Motivtion:
Mixed deployment runs multiple containers on the same node with both online and offline applications. Since online applications are latency-sensitive, we need a Latency PMU to monitor request latency.

Issue:
- DDRC statistics only count the latency of requests within DDRC.
- It is possible that the foreground request is blocked in CMN, while the DDR is still relatively idle. As a result, even if the overall latency is longer, the latency seen only from DDR is relatively small.

Solution:
- Add a new pmu component on the risc-v core;
- The new core pmu component needs to be able to count the latency of the entire life cycle of a request from the core.
- The new core pmu needs to support the collection of latency indicators such as memory read request latency, write request latency, and total latency.

# RISC-V proposal

- RFC Proposal2:

Motivation:
- Lack of unified PMU specifications results in numerous redundant perf drivers in the kernel with repetitive code for device discovery, interrupt registration, and counter enablement. ARM's Coresight PMU specification addresses this by providing a unified design for components like DDR, SMMU, and PCIe Root Complex.

Issue:
-  Lack of a unified uncore pmu for the risc-v architecture

Solution：
- Add an uncore unified architecture pmu adapted to the risc-v architecture
- Add new acpi table to adapt risc-v uncore pmu node;
- The pmu needs to support multi-component event statistics and add multi-component counters
- Design supports overflow interrupt

Q & A