

## Article

# Top-Down Models across CPU Architectures: Applicability and Comparison in a High-Performance Computing Environment

Fabio Banchelli \* , Marta Garcia-Gasulla  and Filippo Mantovani \* 

Barcelona Supercomputing Center, Plaça Eusebi Güell, 1-3, 08034 Barcelona, Spain; marta.garcia@bsc.es

\* Correspondence: fabio.banchelli@bsc.es (F.B.); filippo.mantovani@bsc.es (F.M.)

**Abstract:** Top-Down models are defined by hardware architects to provide information on the utilization of different hardware components. The target is to isolate the users from the complexity of the hardware architecture while giving them insight into how efficiently the code uses the resources. In this paper, we explore the applicability of four Top-Down models defined for different hardware architectures powering state-of-the-art HPC clusters (Intel Skylake, Fujitsu A64FX, IBM Power9, and Huawei Kunpeng 920) and propose a model for AMD Zen 2. We study a parallel CFD code used for scientific production to compare these five Top-Down models. We evaluate the level of insight achieved, the clarity of the information, the ease of use, and the conclusions each allows us to reach. Our study indicates that the Top-Down model makes it very difficult for a performance analyst to spot inefficiencies in complex scientific codes without delving deep into micro-architecture details.

**Keywords:** performance models; top-down model; HPC applications; MareNostrum 4; A64FX; Power 9; Zen 2



**Citation:** Banchelli, F.; Garcia-Gasulla, M.; Mantovani, F. Top-Down Models across CPU Architectures: Applicability and Comparison in a High-Performance Computing Environment. *Information* **2023**, *14*, 554. <https://doi.org/10.3390/info14100554>

Academic Editors: Lenore Mullin, John L. Gustafson and Hamid R. Arabnia

Received: 8 August 2023

Revised: 15 September 2023

Accepted: 5 October 2023

Published: 10 October 2023



**Copyright:** © 2023 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction and Related Work

Diversity in CPU architectures is the new reality in the HPC industry. The five top spots in the Top500 list of November 2022 [1] include three different CPU architectures: x86, Arm, and IBM POWER. To increase the diversity of the panorama, each of these CPU architectures has different implementations by the vendors (e.g., x86 by Intel and AMD, Arm by Fujitsu, Nvidia, and Huawei). With such diversity, it becomes increasingly difficult to establish cross-platform methods to evaluate the efficient use of these hardware systems. The complexity may not arise with benchmarks such as HPL or HPCG but with production scientific applications.

Multiple performance analysis models try to (i) measure performance, and (ii) identify performance bottlenecks. Furthermore, some models can also give hints on how to circumvent said bottlenecks.

For example, the Roofline model [2,3] plots performance in relation to arithmetic intensity (or operational intensity). This model gives feedback on whether a particular code is compute-bound or memory-bound and also tells how far the execution is from the theoretical peak. Depending on the level of analysis, the model can define different roofs (e.g., accounting for cache levels and main memory) [4]. The arithmetic intensity of a code can be either defined (i) theoretically, by counting the number of logical operations that the algorithm requires, or (ii) empirically, by measuring the operations during execution. The second method yields different results because it includes the modifications that the compiler might introduce during the implementation of the algorithm [5]. For codes far from the theoretical peak, the main limitation of the Roofline model is to tell where the performance is lost. For memory-bound codes, the performance will eventually be bounded by the memory subsystem. However, that does not mean the current implementation or execution shares the same bottleneck.

The question is, *where are the execution cycles being lost?* The Top-Down model tries to answer this question.

### 1.1. Top-Down Model

By counting and classifying the execution cycles, the Top-Down model points to the current limiting factor of the application. Intel originally described this model [6] and later established it as the Intel TMAM methodology [7]. The model is structured as a tree of metrics that organizes cycles depending on what they spent (e.g., compute resources, memory resources, lost due to stalls, etc.) The proposed evaluation methodology is to drill down the path where most cycles are lost. For Intel CPUs, there is an official definition for each level of the hierarchy. The deeper the hierarchy goes, the more micro-architecture-specific it becomes. This might make it difficult to compare results between CPUs.

In contrast to Intel CPUs, there is no official Top-Down model defined by AMD. There have been past efforts trying to map the model from Intel to the AMD 15h, Opteron, and R-Series Family processors [8]. This work highlights the limitations of mapping due to differences in micro-architecture and the available hardware counters. Although not labeled Top-Down, there are also tree-like hierarchies for other architectures. Each vendor has defined these and may not match the definitions of TMAM.

The reader should note that the Top-Down model is always defined in terms of cycles. This means that CPU (and memory) clock frequency is not taken into account.

In this work, we leverage previous experience with a production CFD code named Alya [9] to explore the insights the Top-Down model can provide. We define a Top-Down model for the AMD Zen 2 CPU based on the original Intel model and implement a workflow to measure and compute metrics for each system under study. We also analyze the applicability of the models in different CPU architectures and if the model hierarchies are comparable.

### 1.2. Contributions

1. Define the Top-Down model in AMD Zen 2.
2. Implement the Top-Down model in Intel Skylake, AMD Zen 2, A64FX, Power9, and Kunpeng 920 CPUs.
3. Apply the Top-Down model to study the effect of code modifications in a production HPC code across different CPU architectures.
4. Compare the Top-Down model across systems with different CPU architectures.

We worked on five different clusters with three different architectures: x86, aarch64, and ppc64le. For one of the clusters (the one based on Intel Skylake CPUs), the Top-Down model has already been defined [7]. For two of them, there was no model, so we had to define the Top-Down model; (i) for AMD Zen 2 we based it on a previous CPU generation [8], while (ii) for Kunpeng 920 we leveraged the model provided by the CPU manufacturer (not publicly available until now as far as we know). For the clusters A64FX and Power9, there was some model definitions similar to the Top-Down model that we had to adapt to be comparable.

## 2. HPC Systems and Their Top-Down Model

In this section, we introduce the hardware under study. We explain how the Top-Down model of each machine was constructed and how to interpret the most relevant metrics. We use the Top-Down model of TMAM as a baseline to compare with other machines. Please refer to Appendix A for a detailed listing of how to compute each metric.

Table 1 summarizes the cluster configurations of all the HPC systems under study. We include some relevant hardware features as well as the system software stack. We also show a general summary of the Top-Down model on each machine.

**Table 1.** Hardware and software configurations for MareNostrum 4, CTE-AMD, CTE-Arm, CTE-Power, and CTE-Kunpeng.

	MareNostrum 4	CTE-AMD	CTE-Arm	CTE-Power	CTE-Kunpeng
Cluster architecture					
Number of nodes	3456	33	192	52	16
CPU Model	Xeon Platinum 8160	EPYC 7742	FX1000	Power9 8335-GTH	Kunpeng 920
Architecture	x86_64	x86_64	aarch64	ppc64le	aarch64
CPUs per node	2	1	1	2	2
Cores per CPU	24	64	48	20	64
Frequency (MHz)	2100	2250	2200	3000	2600
Multi-Threading	No	No	No	No	No
Floating-point performance					
Vector/SIMD extension	AVX512	AVX2	SVE/NEON	VSX	NEON
Vector/SIMD size (B)	8	4	8/2	2	2
Peak performance (GFlop/s)	67.20	54.16	70.40/17.60	24.0	10.40
Memory subsystem					
L1 Cache (KiB)	32 private	32 private	64 private	32 private	64 private
L2 Cache (MiB)	1 private	0.5 shared	32 shared	0.5 shared	0.5 private
L3 Cache (MiB)	33 shared	16 shared	-	1 shared	32 shared
Main Memory (GB)	96	1024	32	512	256
System software					
Kernel	Linux/4.4.120	Linux/4.18.0	Linux/4.18.0	Linux/4.14.0	Linux/4.14.0
OS	SUSE/12.2	Rocky Linux/8.5	Red Hat/8.1	Red Hat/7.5	CentOS/7
Compiler	intel/2020.1	intel/2018.4	arm/20.3	pgi/20.4	gcc/11.2.0
PAPI Library	papi/6.0.0	papi/6.0.0.1	papi/git-2020-10-08	papi/6.0.0	papi/6.0.0.1
MPI Library	impi/2018.4	impi/2018.4	openmpi/4.0.5	openmpi/3.0.0	openmpi/4.1.3
Tracing Library	extrae/3.8.3	extrae/3.8.3	extrae/3.8.3	extrae/3.8.3	extrae/3.8.3
Top-Down model					
Hierarchy levels	2	1	3	6	3
Metrics relative to	Total Slots	Total Slots	Parent metric	Parent metric	Total Slots (Except for memory metrics)
Parameters	Pipeline_Width	Pipeline_Width, Mispredict_Cost	-	-	Pipeline_Width
Required event sets	2	2	3	10	2

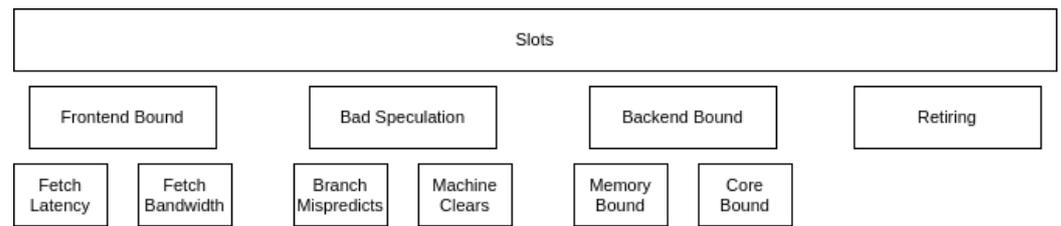
### 2.1. MareNostrum 4 General Purpose

MareNostrum 4 is the flagship Tier-0 supercomputer hosted at the Barcelona Supercomputing Center (BSC). The general purpose partition (from here on, simply MareNostrum 4) has 3456 nodes housing two Intel Xeon Platinum 8160 CPUs. This partition ranked 29th in the Top500 in June 2019.

The Top-Down model in MareNostrum 4 has been constructed following Intel's TMAM definition. Figure 1 shows a schematic view of the model. This version of the model focuses on the *Slot* occupation at the boundary between the front-end and the back-end of the processor pipeline. One slot can be either consumed by one micro-operation ( $\mu$ Op) coming from the front-end or get lost because the corresponding resource in the back-end is busy. In the case of the Skylake CPU, there are four slots available per cycle (up to four  $\mu$ Ops can be dispatched per cycle). These slots can fall under one of the following categories:

- Frontend Bound: the slot was lost due to insufficient  $\mu$ OPs to execute.
- Bad Speculation: the slot was used, but to execute a speculative instruction that was later cleared.
- Back-end Bound: the slot was lost due to the back-end resources being occupied by an older  $\mu$ OP.
- Retiring: the slot was used for an instruction that eventually retired. A high number in this metric means that the pipeline has not lost slots due to stalls. This does not mean that the hardware resources are being utilized efficiently. It only means that work is flowing into the pipeline.

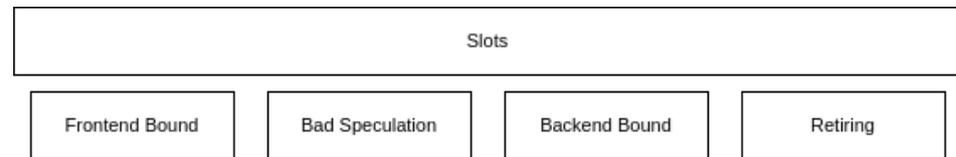
Each category is further divided into more detailed metrics. For example, the *Back-end Bound* category is divided into *Memory Bound* and *Core Bound*. In this version of the Top-Down model, child metrics add up to their parent's value. All values represent the portion of the total slots of the execution. The sum of all metrics in the first level adds up to one. With the hardware counters available in MareNostrum 4, we were able to construct two levels of the Top-Down model for the Skylake CPU. These two levels are sufficiently generic to be compared to other CPUs, even if based on different architectures.



**Figure 1.** Top-Down model hierarchy in MareNostrum 4.

## 2.2. CTE-AMD

CTE-AMD is part of the CTE clusters deployed at the BSC. It has 33 nodes made up of AMD Rome processors, similar to the CPUs of the Frontier supercomputer that was installed in 2021 at Oak Ridge National Laboratory and ranked first in the Top500 in June 2022. Our work extends the previously defined mapping of the R-Series Family processors to the EPYC 7742 CPU hosted in the CTE-AMD. Due to the limitations of the system software (i.e., hardware counters and PAPI library), we were only able to map the first level of metrics. Figure 2 shows a schematic view of the model.



**Figure 2.** Top-Down model hierarchy in CTE-AMD.

The Top-Down model in CTE-AMD shares the first level of metrics with Intel's model. However, despite sharing the same name, there are some differences between the metrics in Intel and AMD, mainly:

- Frontend Bound is based on the counter `UOPS_QUEUE_EMPTY`, which does not consider whether the back-end is stalled or not.
- Bad Speculation requires a micro-architectural parameter, *Mispredict\_Cost*, which represents the average cycles lost due to a misprediction. We define this constant as 18 based on the publicly available experimental data [10].

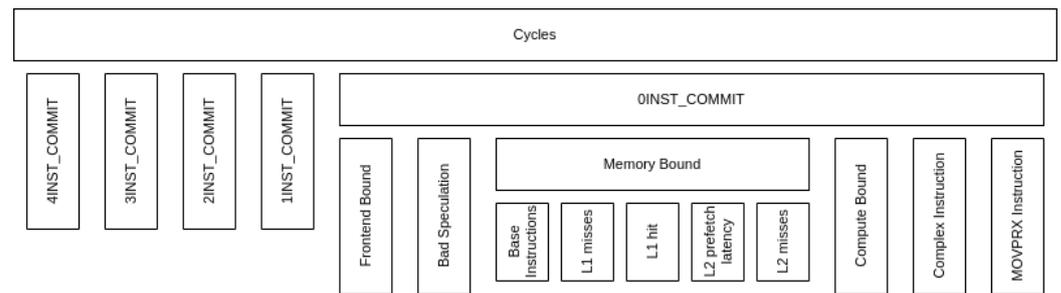
Like with MareNostrum 4, the metrics of each level are between zero and one and represent a portion of slots of the total execution.

## 2.3. CTE-Arm

CTE-Arm is another cluster of the CTE systems. It is powered by the A64FX chip developed by Fujitsu and based on the Arm-v8 instruction set. The cluster's architecture is the same as the one of the Fugaku supercomputer, which ranked first in the Top500 in June 2020.

The Top-Down model in CTE-Arm has been constructed based on the official micro-architecture manual published by Fujitsu [11]. The manual uses the term Cycle Accounting instead of Top-Down model. Figure 3 shows a schematic view of the model. The model is presented as a tree of hardware counters instead of metrics. It does not require any micro-architectural parameter.

The hierarchy is up to five levels deep. Almost all nodes in the hierarchy correspond to a hardware counter, and each node is the sum of all its children nodes. This means that there is some redundancy between counters of the hierarchy, but it also means that it is not necessary to construct the whole hierarchy to study the first and second levels. Some metrics marked as *Other* correspond to the difference between the parent node and the aggregation of all children nodes.



**Figure 3.** Top-Down model hierarchy in CTE-Arm.

In contrast to the Top-Down model in MareNostrum 4 and CTE-AMD, the metrics in CTE-Arm focus on the cycles lost during the commit stage of the instructions (not in the boundary between the frontend and back-end). The first level of the hierarchy classifies cycles depending on how many instructions were committed. The best case scenario is four instructions, while the worst is no instructions committed. Since the situation of zero commits in a cycle is the most critical, the Top-Down model hierarchy focuses on this path.

In contrast to MareNostrum 4 and CTE-AMD, the metrics in CTE-Arm are relative to the parent. This means that the metric represents a portion of the cycles of the parent metric instead of the total execution.

#### 2.4. CTE-Power

CTE-Power is part of the CTE clusters and spans across 52 nodes housing two IBM Power9 8335-GTH CPUs. Summit and Sierra are two supercomputers based on the same CPU as CTE-Power, ranked first and third in the Top500 in June 2018.

The Top-Down model in CTE-Power has been constructed based on the official PMU user guide published by IBM [12]. Figure 4 shows a schematic view of the model. It does not require any micro-architectural parameters. The manual includes a tree diagram of the first two levels of the model hierarchy. Lower levels are only referenced using their respective hardware counters.

- ICT empty no instruction to complete (the pipeline is empty). Similar to *Frontend Bound* in other models.
- Issue hold next-to-complete instruction (i.e., oldest in the pipeline) is held in the issue stage.
- Pipeline stall similar to the *Back-end Bound* category in previous models.
- Thread blocked the next-to-complete instruction is held because an instruction from another hardware thread is occupying the pipeline.
- Instruction latency cycles waiting for the instruction to finish due to the pipeline latency.
- Completion cycles cycles in which at least one instruction has been completed. Similar to the *Retiring* category in previous models.

The hierarchy is up to six levels deep. Starting from the second level, some nodes of the hierarchy do not represent an actual hardware counter that is available in the machine but the aggregation of the counters below it. This means that the Top-Down model in CTE-Power requires the construction of the whole tree if the study needs to go beyond the first level of the hierarchy. In this work, we only study metrics up to the third level and follow the *Memory* path since it is the most relevant for the application under study. Like with the CTE-Arm, the metrics in CTE-Power always represent a portion of the parent's cycles, not the whole execution.

From the official documentation, it is unclear whether the categories of the model are disjointed or not. Furthermore, the *Instruction latency* category (labeled as *Finish-to-completion* in the manual) is mentioned once when listing the top level of the model but not mentioned later on. Without more details, it is not possible to drill down the path of *Instruction latency* if it were to be the main limiting factor of the application under study.

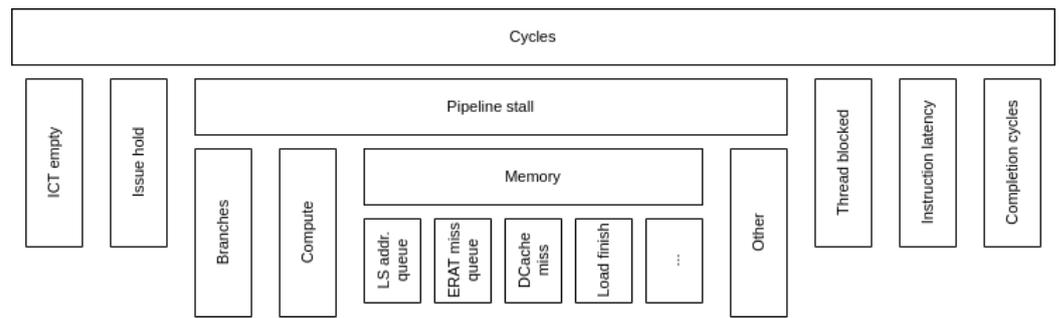


Figure 4. Top-Down model hierarchy in CTE-Power.

In MareNostrum 4, the *Retiring* category branches into other metrics that take into account the pipeline latency (we do not include this branch in our work because we are limited by the available counters). One could argue that the *Instruction latency* in CTE-Power could be included under *Completion cycles* similar to the model in MareNostrum 4.

### 2.5. CTE-Kunpeng

CTE-Kunpeng is a cluster powered by the Arm-based Kunpeng 920 CPU, which was deployed in 2021 at the BSC as a result of a collaboration between the BSC and Huawei. The model in CTE-Kunpeng aims to mimic the same structure and naming convention as in MareNostrum 4. The names of the hardware counters differ from those of the Skylake CPU, and it is unclear whether they cover the same situation exactly. Nonetheless, the formulation of the metrics emulates the original definition by Intel. The model in CTE-Kunpeng defines three levels. Figure 5 shows a schematic view of the model.

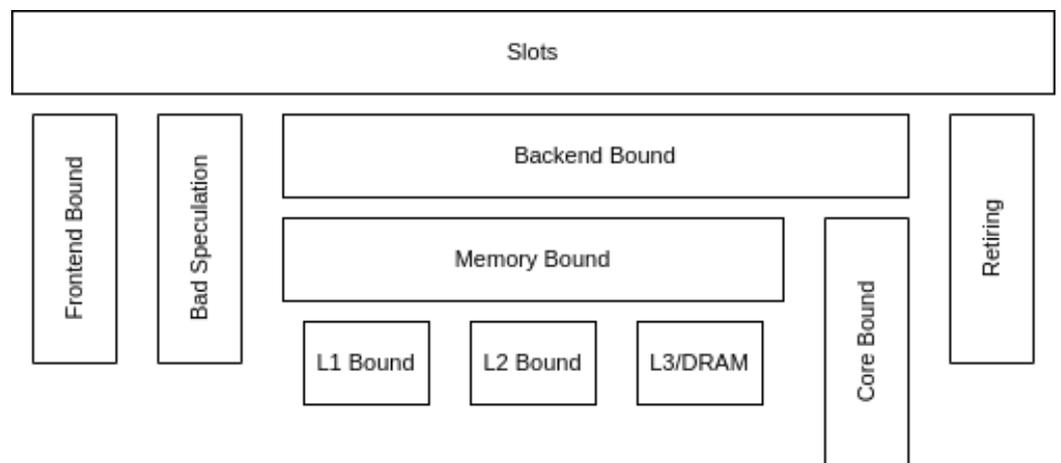


Figure 5. Top-Down model hierarchy in CTE-Kunpeng.

The first level of the hierarchy includes metrics which represent a portion of the total cycles. On the other hand, the *Memory Bound* branch in the second and third levels defines metrics that are relative to the cycles in which the core was stalled waiting for memory. CTE-Kunpeng is the only case we cover in this document where metrics are relative to different quantities depending on the level of the hierarchy. This makes it more difficult to compare metrics between models.

### 3. HPC Application: Alya

Alya is a Computational Fluid Dynamics code developed at the Barcelona Supercomputing Center and part of the PRACE Unified European Applications Benchmark Suite [13]. The application is written in Fortran and parallelized with the MPI programming model. In this work, we run a version and input of Alya that we previously studied [9]. This version implements a compile time parameter `VECTOR_SIZE`, which changes the packing of mesh

elements to expose more data parallelism to the compiler. Our previous study explored how the execution time, executed instructions, and IPC evolved while increasing `VECTOR_SIZE`. We measured that the elapsed time curve has a *U* shape, with `VECTOR_SIZE` 32 being the best configuration. Our study was limited to MareNostrum 4 and compared different compilers. In this work, we leverage our previous knowledge of Alya to run in multiple clusters and construct the Top-Down model for each one of them.

#### General Structure

Our use case is the simulation of a chemical combustion. Thus, it includes computing the chemical reaction, the reaction temperature, and the fluid velocity.

For our particular input, the execution of Alya is divided into five integration steps (or timesteps). Each step is further divided into phases. Figure 6 shows a timeline of one timestep in MareNostrum 4. The *x*-axis represents time, while the *y*-axis represents MPI processes. The timeline is color-coded to show the different execution phases. In this work, we explore the Nastin matrix assembly phase, which represents the longest time in a timestep (blue regions of the timeline) and corresponds to the non-compressible fluid.

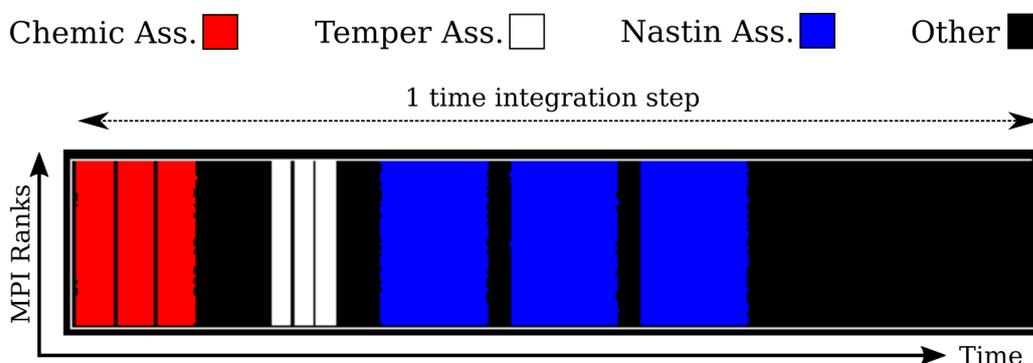


Figure 6. Timeline of one timestep in Alya.

## 4. Model Implementation

In this section, we (i) list the tools used to measure the hardware counters in each cluster and (ii) explain the methodology that chains all the tools together to compute the Top-Down model metrics according to the model definitions. All tools are publicly available except for Daltabaix, which was developed specifically for this work. Access can be provided upon request to the authors.

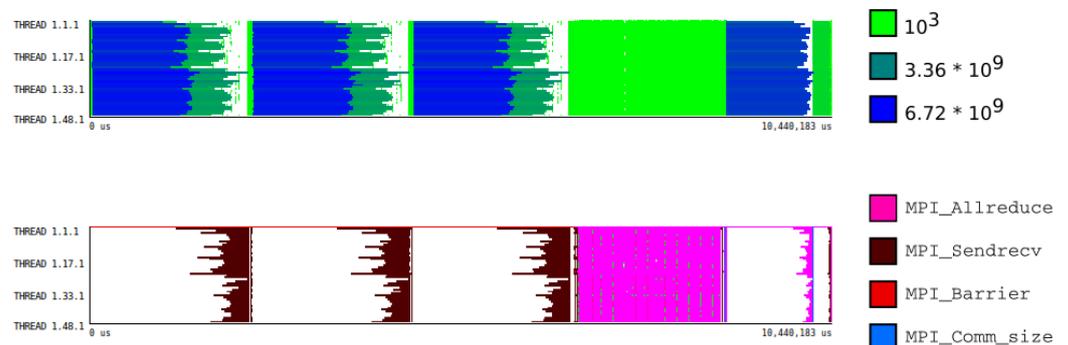
### 4.1. Tools

PAPI (Performance Application Programming Interface) [14] is a library that leverages the portability of `perf` and has an easy programming interface. It also defines a list of generic counters called presets that should be available on most CPUs. However, even with the same name, PAPI counters may not measure the same event when read in different systems. For example, `PAPI_VEC_INS` may read all issued vector instructions on a CPU while only measuring issued arithmetic vector instructions on another system.

Extræ [15] is a tracing tool that intercepts MPI calls and other events during the execution of an application and collects runtime information. The information gathered includes (i) performance counters via PAPI, (ii) which MPI primitive was called, and (iii) which processes were involved during the communication. All this information is stored in a file called a trace. Each record in a trace file has an associated timestamp. Tracing MPI applications with Extræ does not require recompiling the application.

Paraver [16] is a visualization tool that helps navigate the traces generated by Extræ. A common visualization mode when using Paraver is the timeline. Timelines represent the evolution of a given metric across time. Figure 7 shows two examples of timelines. The *x*-axis represents time, and the *y*-axis represents the MPI processes. Each burst is color-coded. For quantitative values (top), the color scale goes from dark blue (high values)

to light green (low values). For qualitative values (bottom), each color represents a different concept (e.g., MPI primitive).



**Figure 7.** Paraver timeline examples. **Top:** quantitative representation of number of instructions. **Bottom:** MPI primitive calls.

Daltabaix is a data processing tool we implemented that takes Paraver traces and computes the metrics of the Top-Down model for a given CPU. The list of metrics, which counters are necessary, and how to combine them are stored in a configuration file for each supported CPU model. Since the number of hardware counters that can be polled with a single PAPI event set is not enough to construct the whole Top-Down model in each machine, Daltabaix collects counters from traces of executions with different event sets. For parallel executions, Daltabaix takes the sum of all measurements across ranks to compute the metrics of the Top-Down model. The reason behind this method is that we define a budget of *Clocks* or *Slots* (depending on the system) that the application has consumed regardless of which cycle belongs to which core.

#### 4.2. Measurement Methodology

In each machine, we run Alya using one full node mapping one MPI rank to each core. We configured the simulations to run with five timesteps. We manually instrumented the code to perform measurements at each timestep's beginning and end of the Nastin matrix assembly phase. We verified that the variability of the hardware counters is under 5% between runs, so we assume that we can operate between counters that belong to different executions.

Figure 8 shows a schematic view of the workflow of our study. From left to right: (i) we execute Alya with Extrae instrumentation. (ii) Extrae calls the PAPI interface, which will poll the hardware counters. (iii) At the end of the execution, Extrae generates a Paraver trace, which is fed into Daltabaix. (iv) Following the definition of the configuration file, Daltabaix extracts the relevant metrics for the Top-Down model by calling the Paraver command line tool. Metrics are printed out in tabular form and also stored in a separate file.

Given that there is no common ground between Top-Down models in each CPU, Daltabaix needs to work differently for each cluster based on the definition of the Top-Down model and its architectural parameters. The added value of this tool is that it provides a single interface to unify all the Top-Down-like models. Nonetheless, the reader should note that it is not necessary to use the same tools we present to reproduce our results. The tables presented in Appendix A are a recipe to construct the same models for each cluster, while the data gathered with hardware counters (regardless of the tool used to query them) are the ingredients. Daltabaix simply acts as the cook that follows the model definitions and outputs the data in a digestible format.

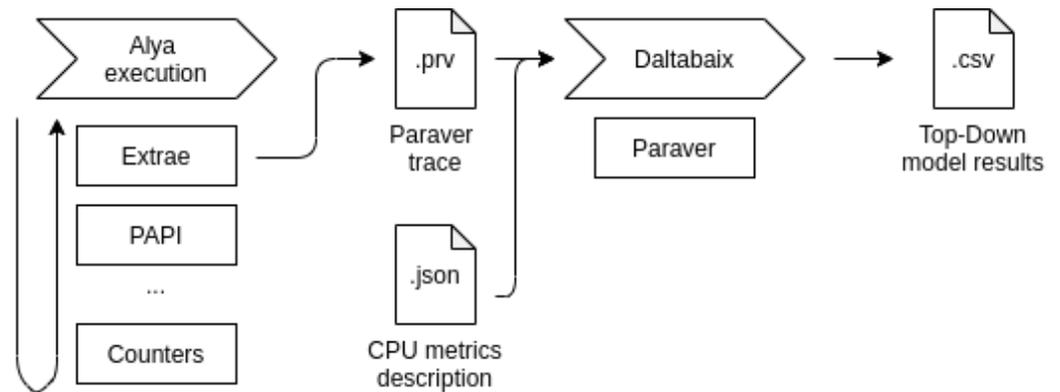


Figure 8. Workflow from application execution to model results.

### 5. Results

In this section, we present the Top-Down model results for Alya in each system under study. We lay down the model hierarchy as tables. Rows represent metrics, while each column represents runs with a given VECTOR\_SIZE. The metrics above the dotted line are not part of the Top-Down model but provide complementary information that we consider relevant to analyze the metrics as a whole. Each cell contains the value of a given metric for a given run and is color-coded with a gradient from red (metric equals zero) to green (metric equals one). The reader should keep in mind that depending on the machine, the metrics represent a portion of the total amount of consumed slots (or cycles), while for others, they represent a portion of the parent metric.

#### 5.1. MareNostrum 4

Figure 9 shows the Top-Down model for Alya in MareNostrum 4. We observe that the cycles decrease until a VECTOR\_SIZE of 32, after which they bounce back up. The number of executed instructions and the IPC also coincide with our previous knowledge of the application. The Top-Down model tells us that Alya in MareNostrum 4 is affected by different factors when increasing VECTOR\_SIZE.

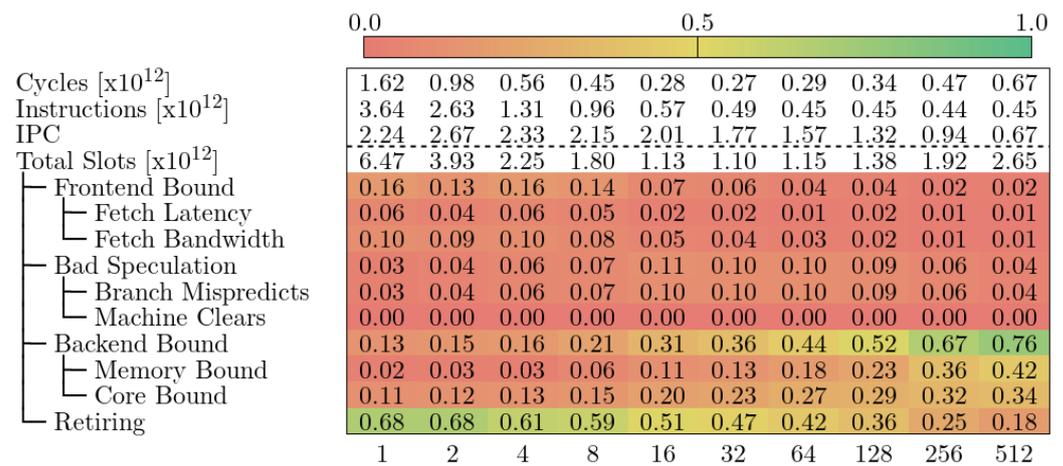


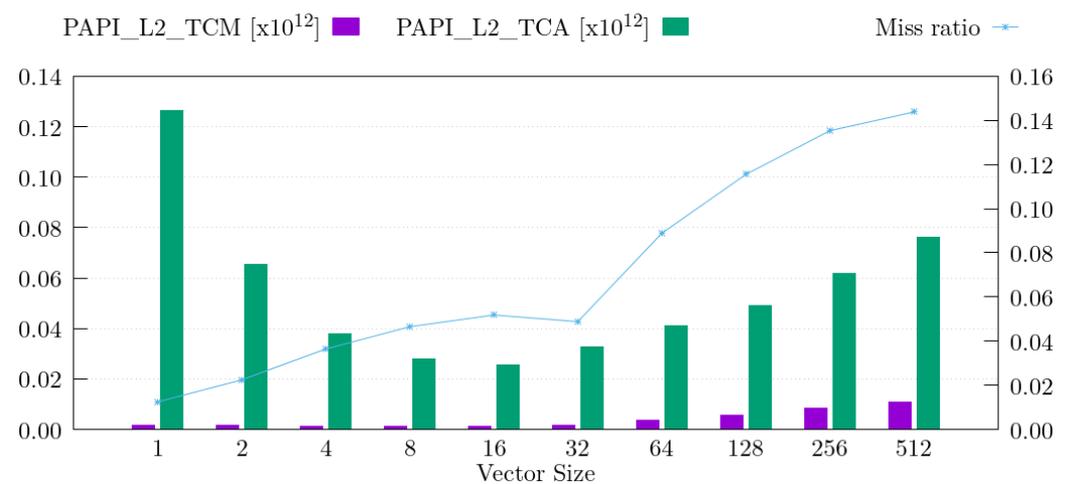
Figure 9. Top-Down model for Alya in MareNostrum 4.

1 to 32: The majority of the cycles are categorized as *Retiring*, which means that the pipeline is not stalled. The reader should note that the model does not provide insight into whether the instructions that go through the pipeline are useful calculations or not.

32–128: The code is *Core Bound*, meaning the pipeline is stalled due to the computational resources. The official definition of TMAM defines metrics under *Core Bound*, but the PAPI installation in MareNostrum 4 does not provide access to the necessary counters. An instruction mix or pipeline usage breakdown could give more information about which type of operation is clogging the CPU. Note, however, that the *Memory Bound* metric is gradually increasing. This means that the pipeline stalls due to memory access are becoming more impactful.

128–512: The code falls under *Memory Bound*, which means that slots are lost because the pipeline is waiting for memory resources. Again, the original model describes detailed metrics we cannot measure in MareNostrum 4. Furthermore, we do not know which actions the programmer could take to reduce memory pressure.

Since we could not construct further levels of the Top-Down model in MareNostrum 4, we cannot further investigate the *Memory Bound* category. However, we know from our previous study that the issue with high values of VECTOR\_SIZE is related to the memory, so we can study the number of accesses and misses to the different levels of the memory hierarchy. In the Skylake CPU, the L2 cache is the last level of private cache. It is the highest level in the memory hierarchy where we can measure memory accesses per core. Figure 10 shows the evolution of L2 cache accesses (green bars, measured with PAPI\_L2\_TCA), L2 cache misses (purple bars, measured with PAPI\_L2\_TCM), and the L2 miss ratio (line, measured as PAPI\_L2\_TCM/PAPI\_L2\_TCA).



**Figure 10.** L2 accesses and misses in MareNostrum 4.

We observe that the number of L2 cache accesses follows a *U* shape, similar to the execution cycles. In contrast, the number of L2 cache misses stays flat from VECTOR\_SIZE 1 to 32 and increases drastically for higher values. The combined view of L2 cache accesses and misses (miss ratio) shows that there is a noticeable jump starting at VECTOR\_SIZE 32. This jump coincides with the point at which the elapsed time of Alya stops decreasing, and the code modification appears to be detrimental. It also matches the results presented in Figure 9; the highest metric when VECTOR\_SIZE is between 32 and 128 is *Core Bound*, but *Memory Bound* is the metric that is consistently increasing. We conclude that the code modification of Alya is beneficial in MareNostrum 4 up to VECTOR\_SIZE 32 because it decreases the number of instructions executed and the number of L2 accesses. However, the modification reaches an inflection point, after which the L2 miss ratio becomes too high and the elapsed time bounces back up.

The Top-Down model has helped us to identify, in general terms, the part of the CPU that stalls during execution. However, more detailed metrics require access to hardware counters that are not accessible on our platform. This limits the insight that the model can provide. We can complement the information the Top-Down model gives us to compensate for the missing metrics (e.g., L2 accesses and misses). Without this complementary data,

we cannot tell how execution cycles are lost. Furthermore, the study of MareNostrum 4 shows that looking at the metric with the highest value only tells part of the story.

### 5.2. CTE-AMD

Figure 11 shows the Top-Down model for Alya in CTE-AMD. We observe a behavior similar to MareNostrum 4 in Figure 9 (i.e., the number of execution cycles decreases until VECTOR\_SIZE 32 and bounces back up). In this case, we also identify the Back-end Bound category as the main limiting factor for high values of VECTOR\_SIZE. However, the starting point of *Back-end Bound* is 56% for CTE-AMD, while it is 13% for MareNostrum 4. Moreover, the portion of slots categorized as *Retiring* is always lower in CTE-AMD compared to MareNostrum 4. Unfortunately, we cannot drill further down because we have no definition of metrics or more hardware counters available in the cluster.

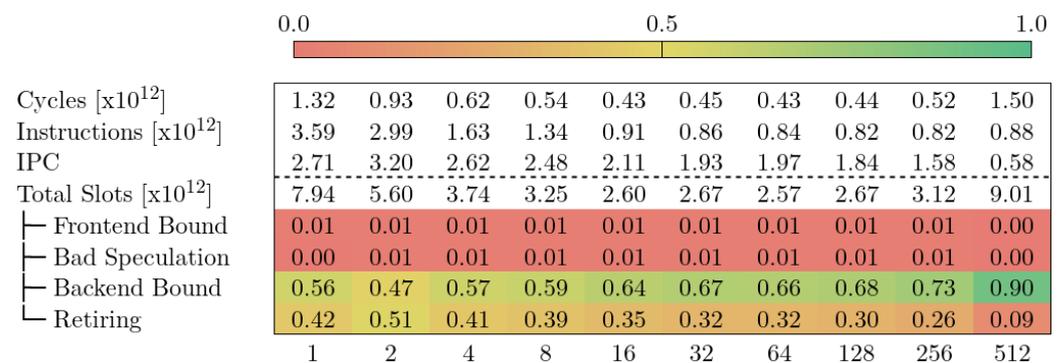


Figure 11. Top-Down model for Alya in CTE-AMD.

### 5.3. CTE-Arm

Figure 12 shows the Top-Down model for Alya in CTE-Arm. The top part represents the first three levels of the model, while the bottom represents the metrics under *Memory Bound*. Like with MareNostrum 4 and CTE-AMD, the execution cycles decrease until a certain value of VECTOR\_SIZE, and then they jump back up. We observe that in CTE-Arm, the values of VECTOR\_SIZE 128 and 256 yield the lowest amount of cycles. This is a much higher VECTOR\_SIZE compared to MareNostrum 4 and CTE-AMD, which was around 32 and 64.

The Top-Down model in CTE-Arm does not have a *Retiring* metric, like in MareNostrum 4. We can achieve a similar metric by combining  $\{4,3,2,1\}INST\_COMMIT$ . This new metric groups cycles where at least one instruction was committed, similar to the definition of *Retiring*, but includes cycles where the CPU could not commit some instruction. The key difference between the model in MareNostrum 4 in CTE-Arm that makes it impossible to have a comparable *Retiring* is that the first uses *Slots* to construct its metrics while the second uses *Cycles*.

Looking at the other end of the spectrum, in Figure 12 we observe that for VECTOR\_SIZE 1, CTE-Arm spends 75% of the execution cycles completely stalled (*0INST\_COMMIT*). This trend is observable across all values of VECTOR\_SIZE, with 4 showing the lowest (59%) value and 512 showing the highest (79%). The reader should note that the *0INST\_COMMIT* metric indicates how badly the pipeline is stalled but does not reflect the total execution time. Furthermore, the metric is relative to the total execution cycles, which means that the run with the lowest *0INST\_COMMIT* is not necessarily the fastest. As it stands, the Top-Down model in CTE-Arm does not allow us to compare metric-to-metric different runs. What it allows us to do is to walk down the hierarchy in a particular run or observe general trends across runs.

For low values of VECTOR\_SIZE, the metrics *Memory Bound* and *Compute Bound* are the main limiting factors. From VECTOR\_SIZE 32 onward, *Memory Bound* accounts for over half of the cycles lost in *0INST\_COMMIT*. We can further drill down and measure the metrics below *Memory Bound* (shown in the bottom part of Figure 11). We observe that *Base*

*instructions* is always the dominant metric, with the exception of VECTOR\_SIZE 1, where *L1 hit* is higher. The description of this metric in the official documentation by Fujitsu [11] states: *Cycles caused by instructions belonging to Base Instructions* (the set of *Base Instructions* in the Armv8 ISA contains basic scalar arithmetic and memory instructions). With only this definition, it is unclear if there is an overlap between the metric and other metrics under *Memory Bound*. However, Table 14-6 of the same document verifies no overlap. Our observation is that the proportion of cycles lost due to *L1 misses* increases while the inverse happens for *L1 hit*. With the information available, we cannot blame the cycles lost waiting for the completion of a memory access *Memory Bound* to a specific type of load operation. We can only conclude that it is a scalar instruction (not SIMD nor SVE).

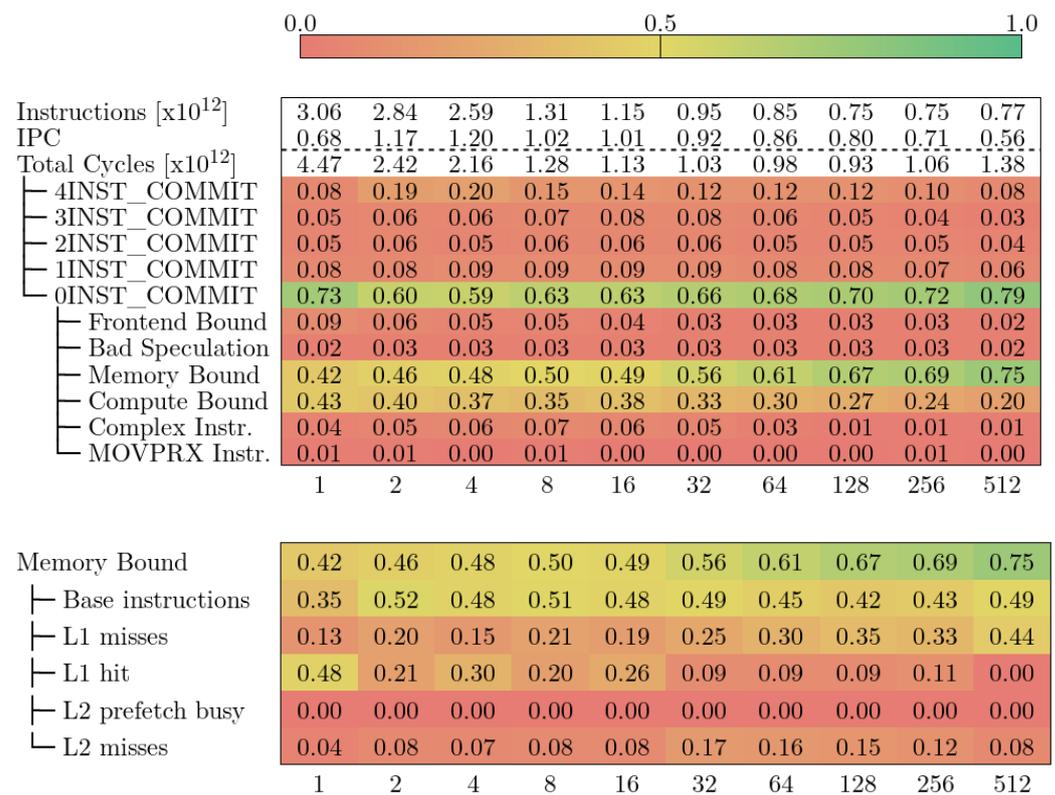


Figure 12. Top-Down model for Alya in CTE-Arm.

#### 5.4. CTE-Power

Figure 13 shows the Top-Down model for Alya in CTE-Power. Like CTE-Arm, the number of total cycles decreases while VECTOR\_SIZE increases, but they start bouncing back up when VECTOR\_SIZE is 512.

Throughout all the executions, the main limiting factor is the *Pipeline stall* category, which is comparable to the *Retiring* category in MareNostrum 4. Furthermore, the *Memory* subcategory always represents between 85% and 90% of the stalled cycles. In contrast, the *Compute* subcategory, the second highest, only accounts for 15% max.

Drilling into the *Memory* category, we observe that two metrics evolve noticeably when increasing VECTOR\_SIZE: *Store reorder queue* and *Store finish*. While the first one represents 50% of the cycles due to memory stalls with VECTOR\_SIZE 1, its weight decreases while the weight of *Store finish* increases.

- *Store reorder queue* is defined by the counter PM\_CPLU\_STALL\_SRQ, which measures the cycles in which a store operation was stalled because the store reorder buffer (SRQ) was full (i.e., too many store operations were in-flight at the same time).

- *Store finish* is defined by the counter `PM_CMPLU_STALL_STORE_FINISH`, which measures the cycles waiting for a store operation that requires all its dependencies to be met to finish (i.e., the nominal latency of a store operation).

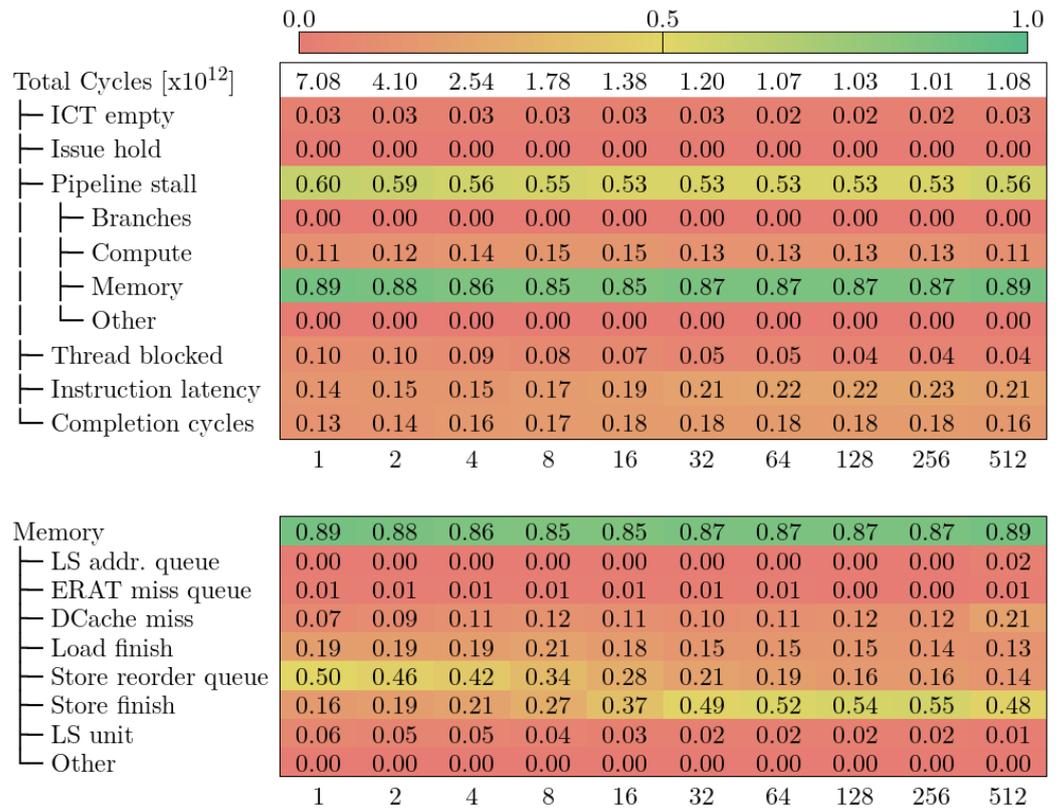


Figure 13. Top-Down model for Alya in CTE-Power.

Following the definition of both metrics, we can infer that the parameter `VECTOR_SIZE` affects the amount of store operations that are in-flight at the same time. For low values of the parameter ( $VECTOR\_SIZE \leq 16$ ), there are too many stores at once, which stalls the pipeline. From that point forward, the concurrent store operations decrease, meaning that the pipeline is still waiting for stores to complete. However, these have a minimum amount of cycles to complete (pipeline latency). The Top-Down model in CTE-Power is heavily dependent on the micro-architecture, so going deep into the hierarchy might give more insight into performance bottlenecks but makes it harder to compare against other machines.

### 5.5. CTE-Kunpeng

Figure 14 shows the Top-Down model for Alya in CTE-Kunpeng. Furthermore, the model has been defined to be very similar to the original Top-Down model from Intel. Contrary to MareNostrum 4, the results in CTE-Kunpeng show that *Core Bound* is the main limiting factor for low values of `VECTOR_SIZE`. We know that the CPU in CTE-Kunpeng has a lower single-thread performance compared to MareNostrum 4 (see Table 1). In the case of Alya, the weaker floating-point throughput of the core is reflected as cycles stalled due to the computational resources being occupied. At the time of writing, we do not have the definition of the metrics under *Core Bound*, so we cannot construct the whole hierarchy.

As `VECTOR_SIZE` increases, so does the *Memory Bound* metric. Starting at `VECTOR_SIZE` 256, Alya has become bound by the memory subsystem. We can study the different cache levels in CTE-Kunpeng for the *Memory Bound* part. This in-depth study was not possible in MareNostrum 4 since we lacked the PAPI counters to compute the metrics. At this point,

the model suggests that cycles are lost due to L3 Cache and DRAM accesses or misses. The cache hierarchy of the Kunpeng 920 CPU shares the L3 level across cores, which makes it difficult to blame accesses or misses on a particular core.

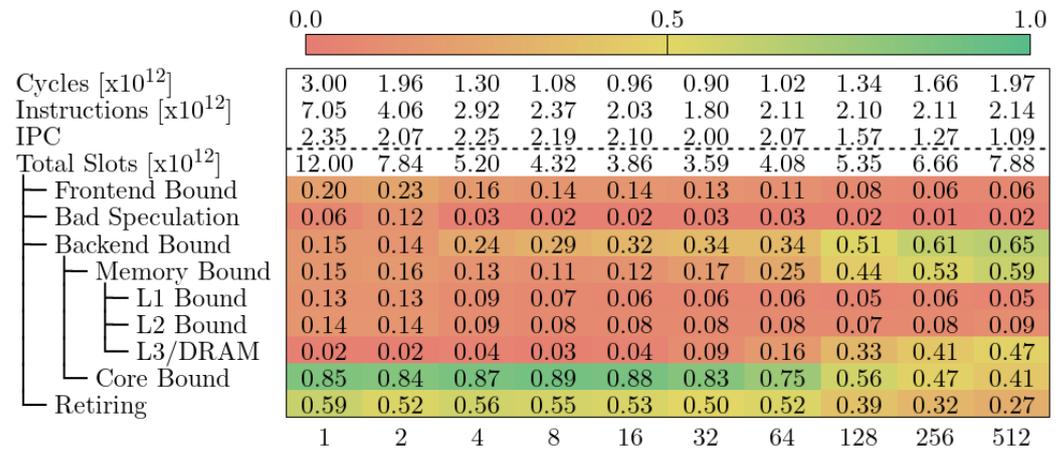


Figure 14. Top-Down model for Alya in CTE-Kunpeng.

### 6. Conclusions and Discussion

In this paper, we have studied the Top-Down model for three different ISAs (x86-64, Arm-v8, and IBM Power9) implemented by five different chip providers (Intel, AMD, Fujitsu, Huawei, and IBM) in five HPC clusters. We used Alya, a CFD application, to measure each model’s metrics and gain insights into performance bottlenecks. The results of our study can be summarized into two main categories: (i) conclusions that can be drawn when studying the Top-Down model within the same cluster and (ii) considerations that relate to the Top-Down results gathered from different clusters.

Within the same cluster we found that implementing the Top-Down model can be tricky; going deep into the model’s hierarchy means needing hardware counters, which are not always available due to limitations in tools, the maturity of the system, and the software configuration. Furthermore, some architectural details can be missing.

On a case-by-case basis, x86-64 Skylake is the most documented and with official support from Intel. The metrics are well defined and easy to understand, but they require hardware counters that are not available in MareNostrum 4. This issue is even more apparent in CTE-AMD, where we were not able to go past the first level of metrics of the Top-Down model. In the case of CTE-Power, the amount of counters needed to construct the whole hierarchy implies obtaining the data with multiple runs. Once constructed, interpreting the metrics requires a deeper understanding of the micro-architecture of the system compared to the other clusters. The model in CTE-Arm tries to strike a balance; the official documentation accompanies the *Cycle Accounting* with multiple tables of complementary metrics, so we were able to go beyond the first level of metrics. However, the model definition in CTE-Arm does not allow for metric comparisons between runs because metrics are relative to their parent in the tree and not to the root.

Once the hardware counters are mapped to the Top-Down model in a given system, we comment on the insights we can obtain from the model itself. The metrics are always defined relative to the execution cycles (or parent metric) of their respective runs. Thus, when comparing runs performed with different software configurations (e.g., changing VECTOR\_SIZE) on the same cluster, having a higher or lower value of a certain metric does not imply a better or worse execution time. The Top-Down model only shows general trends and how this evolves throughout different runs with different configurations (e.g., changing VECTOR\_SIZE). Once constructed, the interpretation of the metrics in all clusters depends on micro-architectural knowledge.

Among different clusters, we discovered even more limitations. Having no standard naming convention for metrics across systems makes it hard to compare the results of Top-Down models gathered on different clusters. Even with the same names, metrics in different systems have subtle differences, e.g., on MareNostrum 4, CTE-AMD, and CTE-Kunpeng, the inefficiencies of the Top-Down model are gathered at the boundary between the frontend and back-end of the CPU, while for CTE-Arm and CTE-Power, the model counts the cycles lost during the commit stage of the instructions.

We have seen that on different machines, we can go deeper in the hierarchy of the model, depending on the availability of hardware counters. However, even having access to the whole Performance Monitoring Unit (where hardware counters are stored) does not magically mean having more insight because a correct interpretation of most of the counters requires a deep knowledge of the underlying micro-architecture that the scientist running a scientific code does not have. This creates a trade-off when using the Top-Down model; the deeper the detail and insight, the more it depends on previous micro-architecture knowledge and the more difficult it becomes to compare against other machines. Moreover, in some clusters, the definition of the model requires architectural parameters, making it more difficult to define the model and compare it with other clusters.

As a general and crucial remark, no metric in any of the models tells *how much* the CPU resources are being used. One could measure 100% of *Slots* in the *Retiring* category, but that does not tell if the code is running efficiently, it simply implies that there are no pipeline stalls. In the Top-Down model defined by Intel, there are some metrics under the *Retiring* category that try to indicate if there is room for improvement. However, we did not include these metrics in our work because (once more) we were limited by the available hardware counters.

While the community accepts the Top-Down model as a method for spotting inefficiencies of HPC codes, we have experienced that:

1. It requires additional information (either micro-architectural details or further hardware counters) to draw a complete picture when analyzing and improving the performance of a scientific application;
2. It does not quantify how much each of the resources within the compute node is used/saturated;
3. It does not allow us to easily compare clusters of the same architecture or different architectures.

All these limitations make the work very difficult for a performance analyst who wants to use the Top-Down model as an inspection tool for spotting inefficiencies in complex scientific codes. The dependencies on micro-architectural knowledge also make it unfeasible for it to be proposed as a co-design tool to the researcher owner of the scientific application under study.

Discussion A simpler, more insightful, and micro-architectural-independent model to measure the saturation of computational resources would be as follows: One counter for execution cycles (e.g., CPU\_CYCLES), which is already available in all clusters. One counter per pipeline that counts *active cycles* (i.e., cycles where the hardware resource is performing useful work). With this information, the performance analyst can compute the portion of cycles in which each computational resource is being used, thus measuring the efficiency of the code for a given hardware. We state that this approach is micro-architectural-independent because it only counts *active cycles*, which is a metric that can be interpreted the same across different functional units. The Roofline model can assist this study by giving more information about how far the actual performance of the code is from the theoretical peak.

In the case of the memory hierarchy, we propose to look at the miss ratio. However, the instruction miss ratio is architecture-dependent since different ISAs may perform a different number of accesses per instruction. A more accurate measurement would be the micro-operation miss ratio, but it has the downside that it is a heavily micro-architecture-dependent metric. A middle-ground would be to have a request counter and a request

hit/miss counter per level of the memory hierarchy. These counters are already available in some clusters but not all. To obtain an accurate performance metric, one would have to know the latency and bandwidth of each level of the hierarchy so that the cost of each access can be measured (again, a micro-architectural-specific feature). There is no simple and hardware-independent way to precisely measure the efficiency of memory access.

Finally, we believe that cycle aggregate metrics, as presented in the Top-Down model, have the benefit of being easy to interpret (a single number) but hide heterogeneous regions of code. In the case of Alya, we know that there is a specific region that we want to study: the Nastin matrix assembly. Nevertheless, for other codes, this previous knowledge might not be available, and the performance analyst will have to construct the Top-Down model for the whole execution of the application. A more precise and automated method would be to first categorize phases of the execution using clustering techniques [17] and then apply the Top-Down model to each one of the phases.

**Author Contributions:** Conceptualization, methodology, and implementation by F.B. Supervision by M.G.-G. and F.M. All authors contributed equally to the conclusions section. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research received no external funding.

**Institutional Review Board Statement:** Not applicable.

**Data Availability Statement:** The data presented in this study are available on request from the corresponding author.

**Conflicts of Interest:** The authors declare no conflict of interest.

### Appendix A. Event Sets and Metrics

In Table A1, we list the PAPI event sets used to construct the Top-Down model on each system under study. The left column displays the name of the event set, while the right column lists the PAPI events included in a given set. Events repeated in more than one set are omitted.

In Table A2, we list the PAPI counters and formulas used to construct the Top-Down model on each system under study. The left column displays the name of a metric, while the right column details the formula to compute the metric. Counters are written in monospace font and using the exact same name as displayed by the `papi_native_avail` command. Metrics of the model as well as helper metrics are written in *cursive*.

**Table A1.** Event set reference for MareNostrum 4, CTE-AMD, CTE-Arm, CTE-Power, and CTE-Kunpeng.

<b>MareNostrum 4</b>	
EventSet1	UOPS_RETIRED:RETIRE_SLOTS, BR_MISP_RETIRED:ALL_BRANCHES MACHINE_CLEAR:COUNT, EXE_ACTIVITY:BOUND_ON_STORES CPU_CLK_THREAD_UNHALTED, CYCLE_ACTIVITY:STALLS_MEM_ANY
EventSet2	IDQ_UOPS_NOT_DELIVERED:CORE IDQ_UOPS_NOT_DELIVERED:CYCLES_O_UOPS_DELIV_CORE UOPS_ISSUED:ANY, INT_MISC:RECOVERY_CYCLES
<b>CTE-AMD</b>	
EventSet1	UOPS_QUEUE_EMPTY, RETIRED_BRANCH_INSTRUCTIONS_MISPREDICTED RETIRED_TAKEN_BRANCH_INSTRUCTIONS_MISPREDICTED RETIRED_INDIRECT_BRANCH_INSTRUCTIONS_MISPREDICTED
EventSet2	CYCLES_NOT_IN_HALT, RETIRED_UOPS, RETIRED_INSTRUCTIONS

Table A1. Cont.

<b>CTE-Arm</b>	
EventSet1	CPU_CYCLES, 0INST_COMMIT, 1INST_COMMIT 2INST_COMMIT, 3INST_COMMIT, 4INST_COMMIT
EventSet2	LD_COMP_WAIT, EU_COMP_WAIT, BR_COMP_WAIT ROB_EMPTY, UOP_ONLY_COMMIT, SINGLE_MOVPRFX_COMMIT LD_COMP_WAIT_EX, LD_COMP_WAIT_L2_MISS
EventSet3	LD_COMP_WAIT_L2_MISS_EX, LD_COMP_WAIT_L1_MISS LD_COMP_WAIT_L1_MISS_EX, LD_COMP_WAIT_PFP_BUSY
<b>CTE-Power</b>	
EventSet1	PM_RUN_CYC, PM_CMPLU_STALL_BRU, PM_NTC_ISSUE_HELD_ARB PM_NTC_ISSUE_HELD_DARQ_FULL, PM_NTC_ISSUE_HELD_OTHER
EventSet2	PM_CMPLU_STALL_EXEC_UNIT, PM_CMPLU_STALL_NTC_DISP_FIN PM_CMPLU_STALL_SRQ_FULL, PM_ICT_NOSLOT_CYC
EventSet3	PM_CMPLU_STALL_EMQ_FULL, PM_CMPLU_STALL_LOAD_FINISH PM_CMPLU_STALL_NTC_FLUSH, PM_CMPLU_STALL_THRD
EventSet4	PM_1PLUS_PPC_CMPL, PM_CMPLU_STALL_DCACHE_MISS PM_CMPLU_STALL_LMQ_FULL, PM_CMPLU_STALL_LSU_MFSR
EventSet5	PM_CMPLU_STALL_LARX, PM_CMPLU_STALL_LRQ_FULL PM_CMPLU_STALL_LSAQ_ARB, PM_CMPLU_STALL_STORE_DATA
EventSet6	PM_CMPLU_STALL_ERAT_MISS, PM_CMPLU_STALL_HWSYNC PM_CMPLU_STALL_LHS, PM_CMPLU_STALL_LRQ_OTHER
EventSet7	PM_CMPLU_STALL_LSU_FIN, PM_CMPLU_STALL_ST_FWD PM_CMPLU_STALL_STORE_FIN_ARB, PM_CMPLU_STALL_STORE_FINISH
EventSet8	PM_CMPLU_STALL_EIEIO, PM_CMPLU_STALL_SLB, PM_CMPLU_STALL_TLBIE
EventSet9	PM_CMPLU_STALL_LWSYNC, PM_CMPLU_STALL_PASTE PM_CMPLU_STALL_STORE_PIPE_ARB
EventSet10	PM_CMPLU_STALL_STCX, PM_CMPLU_STALL_TEND
<b>CTE-Kunpeng</b>	
EventSet1	INST_RETIRED, CPU_CYCLES FETCH_BUBBLE, INST_SPEC
EventSet2	MEM_STALL_ANYLOAD, MEM_STALL_ANYSTORE, EXE_STALL_CYCLE MEM_STALL_L1MISS, MEM_STALL_L2MISS

Table A2. Top-Down modelcounter reference for MareNostrum 4, CTE-AMD, CTE-Arm, CTE-Power, and CTE-Kunpeng.

<b>MareNostrum 4</b>	
Pipeline_Width	4
Clocks	CPU_CLK_THREAD_UNHALTED
Slots	$Pipeline\_Width \times Clocks$
Frontend_Bound	$IDQ\_UOPS\_NOT\_DELIVERED: CORE / Slots$
Fetch_Latency	$Pipeline\_Width \times IDQ\_UOPS\_NOT\_DELIVERED: CYCLES\_0\_UOPS\_DELIV\_CORE / Slots$
Fetch_Bandwidth	$Frontend\_Bound - Fetch\_Latency$
Bad_Speculation	$(UOPS\_ISSUED: ANY - UOPS\_RETIRED: RETIRE\_SLOTS + Pipeline\_Width \times INT\_MISC: RECOVERY\_CYCLES) / Slots$
Branch_Mispredicts	$Mispred\_Clears\_Fraction \times Bad\_Speculation$
Machine_Clears	$Bad\_Speculation - Branch\_Mispredicts$
Mispred_Clears_Fraction	$BR\_MISP\_RETIRED: ALL\_BRANCHES / (BR\_MISP\_RETIRED: ALL\_BRANCHES + MACHINE\_CLEARS: COUNT)$
Backend_Bound	$1 - Frontend\_Bound - (UOPS\_ISSUED: ANY + Pipeline\_Width \times INT\_MISC: RECOVERY\_CYCLES) / Slots$
Memory_Bound	$Memory\_Bound\_Fraction \times Backend\_Bound$
Core_Bound	$Backend\_Bound - Memory\_Bound$
Memory_Bound_Fraction	$CYCLE\_ACTIVITY: STALLS\_MEM\_ANY + EXE\_ACTIVITY: BOUND\_ON\_STORES / Backend\_Bound\_Cycles$
Retiring	$UOPS\_RETIRED: RETIRE\_SLOTS / Slots$
<b>CTE-AMD</b>	
Pipeline_Width	6
Mispredict_Cost	18
Clocks	CYCLES_NOT_IN_HALT
Slots	$Pipeline\_Width \times Clocks$
Frontend_Bound	$UOPS\_QUEUE\_EMPTY / Slots$
Bad_Speculation	$Branch\_Instructions \times Mispredict\_Cost / Slots$
Branch_Instructions	$RETIRED\_BRANCH\_INSTRUCTIONS\_MISPREDICTED + RETIRED\_INDIRECT\_BRANCH\_INSTRUCTIONS\_MISPREDICTED + RETIRED\_TAKEN\_BRANCH\_INSTRUCTIONS\_MISPREDICTED$
Backend_Bound	$1 - (Frontend\_Bound + Bad\_Speculation + Retiring)$
Retiring	$RETIRED\_UOPS / Slots$

Table A2. Cont.

		CTE-Arm
Clocks	CYCLES_NOT_IN_HALT	
4_Instruction_Commit	4INST_COMMIT/Clocks	
3_Instruction_Commit	3INST_COMMIT/Clocks	
2_Instruction_Commit	2INST_COMMIT/Clocks	
1_Instruction_Commit	1INST_COMMIT/Clocks	
0_Instruction_Commit	0INST_COMMIT/Clocks	
Frontend_Bound	ROB_EMPTY /O_INST_COMMIT	
Bad_Speculation	BR_COMP_WAIT /O_INST_COMMIT	
Memory_Bound	LD_COMP_WAIT /O_INST_COMMIT	
Compute_Bound	EU_COMP_WAIT /O_INST_COMMIT	
Complex_Instructions	UOP_ONLY_COMMIT /O_INST_COMMIT	
MOVPRX_Instructions	SINGLE_MOVPRX_COMMIT/O_INST_COMMIT	
		CTE-Power
Clocks	PM_RUN_CYC	
No_Instruction_To_Execute	PM_ICT_NOSLOT_CYC/Clocks	
Instruction_Held_In_Issue	PM_ISSUE_HOLD/Clocks	
Backend_Bound	PM_CMPLU_STALL/Clocks	
Stalled_By_Other_Thread	PM_CMPLU_STALL_THRD/Clocks	
???	PM_1PLUS_PPC_CMPL/Clocks	
Completion_Cycles	1 – (No_Instruction_To_Execute + Instruction_Held_In_Issue + Backend_Bound + Stalled_By_Other_Thread + ???)	
		CTE-Kunpeng
Pipeline_Width	4	
Clocks	CPU_CYCLES	
Slots	Pipeline_Width × Clocks	
Frontend_Bound	FETCH_BUBBLE/Slots	
Bad_Speculation	(INST_SPEC – INST_RETIRED)/Slots	
Backend_Bound	1 – (Frontend_Bound + Bad_Speculation + Retiring)	
Memory_Bound	Memory_Stall_Cycles / EXE_STALL_CYCLE	
Core_Bound	(EXE_STALL_CYCLE – Memory_Stall_Cycles) / EXE_STALL_CYCLE	
Memory_Stall_Cycles	MEM_STALL_ANYLOAD + MEM_STALL_ANYSTORE	
Retiring	INST_RETIRED/Slots	

## References

- Top500 List. 2022. Available online: <https://www.top500.org/lists/top500/2022/11/> (accessed on 1 November 2022).
- Williams, S.; Waterman, A.; Patterson, D. Roofline: An Insightful Visual Performance Model for Multicore Architectures. *Commun. ACM* **2009**, *52*, 65–76. [CrossRef]
- Ofenbeck, G.; Steinmann, R.; Caparros, V.; Spampinato, D.G.; Püschel, M. Applying the roofline model. In Proceedings of the 2014 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), Monterey, CA, USA, 23–25 March 2014; IEEE: Piscataway, NJ, USA, 2014; pp. 76–85.
- Ilic, A.; Pratas, F.; Sousa, L. Cache-aware Roofline model: Upgrading the loft. *IEEE Comput. Archit. Lett.* **2014**, *13*, 21–24. [CrossRef]
- Banchelli, F.; Garcia-Gasulla, M.; Houzeaux, G.; Mantovani, F. Benchmarking of State-of-the-Art HPC Clusters with a Production CFD Code. In Proceedings of the Platform for Advanced Scientific Computing Conference, Geneva, Switzerland, 29 June–1 July 2020. [CrossRef]
- Yasin, A. A Top-Down method for performance analysis and counters architecture. In Proceedings of the 2014 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), Monterey, CA, USA, 23–25 March 2014; pp. 35–44. [CrossRef]
- Intel. *Top-Down Microarchitecture Analysis Method*; Intel: Santa Clara, CA, USA, 2022.
- Jarus, M.; Oleksiak, A. Top-Down Characterization Approximation based on performance counters architecture for AMD processors. *Simul. Model. Pract. Theory* **2016**, *68*, 146–162. [CrossRef]
- Banchelli, F.; Oyarzun, G.; Garcia-Gasulla, M.; Mantovani, F.; Both, A.; Houzeaux, G.; Mira, D. A portable coding strategy to exploit vectorization on combustion simulations. *arXiv* **2022**, arXiv:2210.11917.
- Fog, A. *The Microarchitecture of Intel, AMD, and VIA CPUs—An Optimization Guide for Assembly Programmers and Compiler Makers*; Copenhagen University College of Engineering: Ballerup, Denmark, 2022.
- A64FX Microarchitecture Manual. 2021. Available online: [https://raw.githubusercontent.com/fujitsu/A64FX/master/doc/A64FX\\_Microarchitecture\\_Manual\\_en\\_1.6.pdf](https://raw.githubusercontent.com/fujitsu/A64FX/master/doc/A64FX_Microarchitecture_Manual_en_1.6.pdf) (accessed on 7 August 2023).
- POWER9 Performance Monitor Unit User's Guide. 2018. Available online: [https://wiki.raptorcs.com/w/images/6/6b/POWER9\\_PMU\\_UG\\_v12\\_28NOV2018\\_pub.pdf](https://wiki.raptorcs.com/w/images/6/6b/POWER9_PMU_UG_v12_28NOV2018_pub.pdf) (accessed on 7 August 2023).
- Unified European Applications Benchmark Suite. Available online: <https://prace-ri.eu/training-support/technical-documentation/benchmark-suites/> (accessed on 7 August 2023).
- Terpstra, D.; Jagode, H.; You, H.; Dongarra, J. Collecting Performance Data with PAPI-C. In *Tools for High Performance Computing 2009: Proceedings of the 3rd International Workshop on Parallel Tools for High Performance Computing, September 2009, ZIH, Dresden*; Müller, M.S., Resch, M.M., Schulz, A., Nagel, W.E., Eds.; Springer: Berlin/Heidelberg, Germany, 2010; pp. 157–173.
- Extrac, 2023. Available online: <https://tools.bsc.es/extrac> (accessed on 7 August 2023).

16. Pillet, V.; Pillet, V.; Labarta, J.; Cortes, T.; Cortes, T.; Girona, S.; Girona, S.; Computadors, D.D.D. PARAVÉR: A Tool to Visualize and Analyze Parallel Code. In *Proceedings of WoTUG-18: Transputer and Occam Developments*; Technical Report; IOS Press: Amsterdam, The Netherlands, 1995.
17. Gonzalez, J.; Gimenez, J.; Labarta, J. Automatic detection of parallel applications computation phases. In Proceedings of the 2009 IEEE International Symposium on Parallel & Distributed Processing, Chengdu, China, 10–12 August 2009; pp. 1–11. [[CrossRef](#)]

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.