

Parallelizing Stream Compression for IoT Applications on Asymmetric Multicores (Technical Report)

Xianzhi Zeng
ISTD Pillar

Singapore University of Technology and Design
Singapore, Singapore
xianzhi_zeng@sutd.edu.sg

Shuhao Zhang
ISTD Pillar

Singapore University of Technology and Design
Singapore, Singapore
shuhao_zhang@sutd.edu.sg

Abstract—Data stream compression attracts much attention recently due to the rise of IoT applications. Thanks to the balanced computational power and energy consumption, asymmetric multicores are widely used in IoT devices. This paper introduces *CStream*, a novel framework for parallelizing stream compression on asymmetric multicores to *minimize energy consumption* without violating the user-specified *compressing latency constraint*. Existing works cannot effectively utilize asymmetric multicores for stream compression, primarily due to the non-trivial asymmetric computation and asymmetric communication effects. To this end, *CStream* is developed with the following two novel designs: 1) *fine-grained decomposition*, which decomposes a stream compression procedure into multiple fine-grained tasks to better expose the task-core affinities under the asymmetric computation effects; and 2) *asymmetry-aware task scheduling*, which schedules the decomposed tasks based on a novel cost model to exploit the exposed task-core affinities while considering asymmetric communication effects. To validate our proposal, we evaluate *CStream* with five competing mechanisms of parallelizing stream compression algorithms on a recent asymmetric multicore processor. Our extensive experiments based on a benchmark consisting of three algorithms and four datasets, show that *CStream* outperforms alternative approaches by up to 53% lower energy consumption without compressing latency constraint violation.

Index Terms—Stream compression, Edge Computing and IoT, Asymmetric Hardware

I. INTRODUCTION

Data stream compression, i.e., continuously compressing data attracts much attention recently [1], [2], [3], [4], [5], [6], [7], due to the rise of IoT applications [8], [9]. Figure 1 demonstrates a smart city use case [10], [11], [12] where stream compression is a highly attractive technique. In this application, real-time data streams (e.g., air qualities, wind speeds) from sensors are continuously gathered by the memory-limited, battery-powered patrol drones (i.e., IoT devices). The drone may continuously compress gathered data streams before uploading to the cloud center to reduce data transmission overhead. However, adopting compression does not guarantee “plug-and-play” performance benefits due to the additional compressing latency and hardware resource constraints such as the battery capacity of IoT devices.

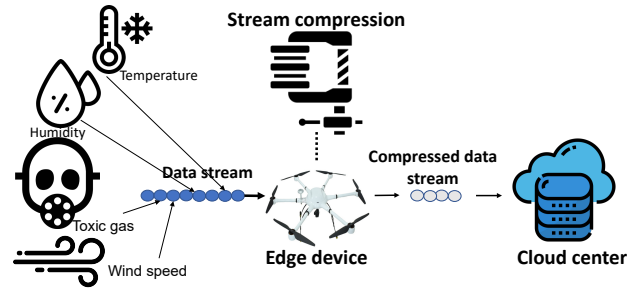


Fig. 1. Real-time data gathering and stream compression at the patrol drone.

According to a 2018 survey [13], modern ARM machines with asymmetric multicores are typical choices for IoT devices [13]. The key to asymmetric multicores is to couple relatively energy-saving and slower cores (i.e., ‘little cores’) and relatively more powerful and power-hungry cores (i.e., ‘big cores’) under the same Instruction Set Architecture (ISA). For instance, an ARM rk3399 processor can be composed with both the in-order A53 ‘little cores’ [14] and the out-of-order A72 [15] ‘big cores’. Such a novel asymmetric architecture balances computational power and energy consumption but brings non-trivial asymmetric computation and asymmetric communication effects. The asymmetric computation effect stands for that the computational power of ‘big cores’ is greater than that of ‘little cores’, and the asymmetric communication effect stands for that the communication latency from ‘big cores’ to ‘little cores’ is larger than the reverse direction.

In this paper, we propose *CStream*¹, a novel framework for parallelizing stream compression on asymmetric multicores. Different from file compression or database compression where all data to be compressed are readily presented, stream compression is an incremental procedure of handling continuously arriving data streams. Specifically, a data stream is a list of tuples chronologically arriving at the system, and

¹All of our code, data, and scripts can be found at <https://github.com/intellistream/CStream>

each tuple needs to be compressed with low latency. When stream compression is conducted in IoT devices, such as in the example in Figure 1, energy consumption is another important factor to be considered. Based on asymmetric multicores, *CStream* parallelizes a stream compression algorithm to compress data streams, such that it *minimizes energy consumption* while satisfying a user-specified *compressing latency constraint*.

Our work is linked to the literature on both parallel data compression algorithms [16], [17], [18], [19], [20], [21], [22], [23] and workload scheduling on asymmetric multicores [24], [25], [26], [27], [28], [29], [30], [31], [32], [33]. They provide highly valuable techniques and mechanisms, but none of them is able to answer the question of “*how to achieve energy efficient and low latency stream compression on asymmetric multicores*”. On the one hand, existing compression algorithms [34], [35], [36], [1], [37], [36] either assume symmetric multicore architectures [34], [36], or exploit heterogeneous hardware like GPU [37] or FPGA [22] which introduces different programming and execution patterns. Furthermore, most of them target on reducing compressing latency without taking care of energy consumption. On the other hand, existing workload scheduling on asymmetric multicores concerns the level of operating system [38], [39], [40], [41], web browser [33] and neural network inference [32]. They typically schedule the entire procedure, such as the matrix multiplication process with partitioned datasets to different cores [32]. Our experimental results will demonstrate that directly adopting such a coarse-grained scheduling mechanism for stream compression on asymmetric multicores can lead to frequent compressing latency violation and much energy dissipation.

Compared to existing works, the superiority of *CStream* is achieved from two novel designs: *fine-grained decomposition* and *asymmetry-aware task scheduling*. First, *CStream* decomposes the entire stream compression procedure (i.e., running of a stream compression algorithm on a batch of data stream, Definition 1) into fine-grained tasks with different *operational intensity* (i.e., instructions per memory access) [42], [24], [41], [43]. Tasks with higher (resp. lower) operational intensity prefer ‘big cores’ (resp. ‘little cores’). Such an approach better exposes task-core affinity and offers opportunities for better utilization of asymmetric multicores compared to existing mechanisms [41], [33], [32]. Second, *CStream* schedules the decomposed fine-grained tasks on asymmetric multicores according to their exposed task-core affinity. The involved asymmetric communication overheads require one to carefully align the communication pattern among tasks. To this end, we propose a novel cost model to guide the scheduling by considering both task-core affinity and asymmetric communication effects. Specifically, our model accurately predicts both the energy consumption and compressing latency of each decomposed task given a scheduling plan (Definition 2). Based on the model, *CStream* searches for the optimal scheduling plan by enumerating all possible plans with dynamic programming.

For a comprehensive comparison, we have implemented and evaluated the parallelization of three representative stream compression algorithms in *CStream* based on a recent ARM processor rk3399 with asymmetric multicores. The evaluation based on both real-world and synthetic datasets confirm the superiority of *CStream*. In particular, it outperforms the alternative mechanisms [41], [34] by up to 53% more energy consumption reduction without violating the strict compressing latency constraints of 11 ~ 26 microseconds for compressing each byte of data stream under varying workload characteristics. In summary, this paper makes the following contributions.

- We develop *CStream*, a novel framework for parallelizing various stream compression algorithms on asymmetric multicores to minimize energy consumption while ensuring the compressing latency is within a user-defined constraint. The design overview of *CStream* is presented in Section III.
- We propose a fine-grained decomposition mechanism (Section IV) to decompose a stream compression procedure into fine-grained tasks based on the compression behavior (i.e., varying operational intensity) of a stream compression algorithm. This allows *CStream* to better expose task-core affinities on asymmetric multicores.
- We propose an asymmetry-aware task scheduling mechanism (Section V), which schedules the decomposed tasks on asymmetric multicores based on a novel cost model to exploit the exposed task-core affinities while taking asymmetry communication effects into account.
- To the best of our knowledge, this work is also the first comprehensive study to compare various competing mechanisms to parallelize stream compression algorithms on asymmetric multicores using both real-world and synthetic datasets (Section VI and Section VII).

II. PRELIMINARIES

In this section, we give a preliminary background of stream compression for IoT applications, followed by reviewing the asymmetric multicore architecture. We summarize the terminologies used in our work in Table I.

A. Data Stream Compression for IoT

A *data stream* is a list of tuples chronologically arriving at the system. Each *tuple* represents an event including timestamp and payloads. *Data stream compression* is an incremental procedure to compress continuously arriving data streams with low latency. It can be conducted solely based on the *current tuple* (i.e., that arrives most recently) or additionally based on the *past tuples* (i.e., those arrive earlier than the current tuple). We classify the former as *stateless stream compression*, which ignore the past tuples, and the latter as *stateful stream compression*, which utilizes a state (e.g., a dictionary [44]) to keep the information of past tuples.

Our work aims to provide framework-level support to optimize stream compression (including stateless and stateful)

TABLE I
SUMMARY OF TERMINOLOGIES

Type	Notation	Description
Workload Specifications	B	The size of batch of data stream to compress
	L_{set}	User specified compressing latency constraint
Device Specifications & Roofline Model	j	A specific AMP core
	C_j	Maximum executable instructions within unit time of core j
	$L_{j',j}^{comm}$	Worst unit communication latency from core j' to j
Model Outputs	L_{est}	Estimated compressing latency
	E_{est}	Estimated total energy consumption
Cost Model Terms	t_i	Task i decomposed from a stream compression job
	p	A possible scheduling plan
	p_{opt}	The optimal scheduling plan
	i_i	Input data size of t_i
	e_i	Energy consumption of t_i
	l_i	Compressing latency of t_i
	l_i^{comp}	Computation latency of t_i
	l_i^{comm}	Communication latency of t_i
	η_i	Instructions per unit time of t_i
	ζ_i	Instructions per unit energy of t_i
	κ_i	Instructions per unit memory access (i.e., operational intensity) of t_i

for IoT applications [9]. In particular, *CStream* parallelizes each stream compression procedure on asymmetric multicores, defined as follows.

Definition 1 (Stream Compression Procedure). A *stream compression procedure* is the process of executing a stream compression algorithm (stateless or stateful) on a batch of data stream, where the batch size (B) is tunable. In this work, we assume the batch size B is pre-determined by applications. We use a pair of Algorithm-Dataset to describe a workload procedure in the following.

Note that, the compressibility of a specific compression algorithm is not a concern of this work. Instead, we focus on the following two strict design requirements for adopting stream compression for IoT applications.

Design Requirements of Stream Compression for IoT:

- (R1) Low Latency Stream Compression: Data streams generated from the IoT applications are often real-time constrained, requiring a low latency compressing to meet the quality-of-service (QoS) goal [45], [41].
- (R2) Low Energy Consumption: Stream compression for IoT needs to achieve low energy consumption as the available energy budget at IoT devices is often quite limited compared with that in a data center [8]. For instance, the devices may be solar or battery-powered and far away from a constant power source.

B. Asymmetric Multicore Architecture

The asymmetric multicore architecture is designed to balance computational power and energy consumption [46], [47], [42], and is increasingly deployed for IoT devices [9], especially due to the world wide rising concern for energy consumption and carbon emissions [48], [49], [50]. Figure 2 depicts a recent 6-core rk3399 processor [51], which couples

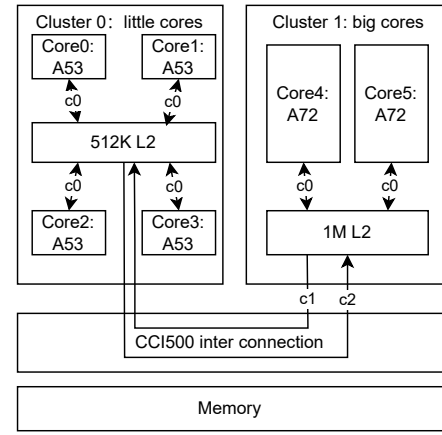


Fig. 2. The 6-core rk3399 processor.

different types of cores (A53 and A72) on the same chip. Such architecture is also often called “big-Little” architecture.

Compared to conventional symmetric multicores, asymmetric multicores involve two non-trivial asymmetry effects: 1) *Asymmetric Computation Effect*. As shown in Figure 2, four A53 cores (i.e., *Core0* ~ *Core3*) and two A72 cores (i.e., *Core4* ~ *Core5*) are coupled in one chip. Although they share the same ISA (i.e., ARM V8), the A53 cores (‘little cores’) are in-order, relatively battery-saving and slower. In contrast, the A72 cores (‘big cores’) are out-of-order, power-hungry and more powerful; and 2) *Asymmetric Communication Effect*. A53 cores and A72 cores are placed in different clusters (*Cluster0* ~ *Cluster1*) resulting in different types of cross-core communication patterns [52], [53], i.e., *inter-cluster* and *intra-cluster* communication. The *inter-cluster* communication needs to go through a slow CCI500 interconnection channel, while *intra-cluster* communication involves L2 cache only.

III. MOTIVATION AND DESIGN OVERVIEW

In this section, we present the motivation for the design of our proposed framework – *CStream* for stream compression on asymmetric multicores, followed by its design overview. The detailed implementation of task decomposition and scheduling with a cost model are presented in Sections IV and V, respectively.

A. Motivations

Parallelizing stream compression on asymmetric multicores can potentially satisfy the aforementioned two design requirements of adopting stream compression for IoT applications. However, the involved asymmetry effects require a careful system design. A poor design can incur both severe compressing latency constraint violation and energy dissipation. It is a particular challenge to support varying stream compression algorithms (e.g., stateless and stateful compression), varying datasets, and varying compressing latency constraints. Our design of *CStream* is motivated by the following observations.

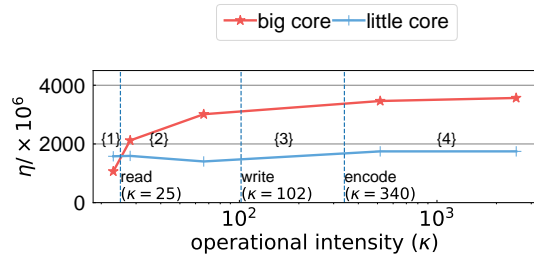


Fig. 3. The four-segment roofline model of asymmetric multicores rk3399. Dashed lines denote the κ of different stream compression steps.

Observation 1: there are varying task-core affinities in different parts of stream compression procedure. The roofline model [54], [55] reveals that the instructions per unit time (η) or energy (ζ) on hardware grows with the *operational intensity* (i.e., instructions per memory access, denoted as κ) of software, before reaching a maximum value (i.e., the so-called “roof”). A stream compression procedure is generally composed of three steps (*read*, *encode*, and *write*) with different κ . Figure 3 shows the roofline model of a ‘big core’ and a ‘little core’ of the rk3399 asymmetric multicores, and we adopt a benchmark by Lo et al. [56] for profiling the operational intensity κ and instructions per unit time η . There are two key takeaways. First, the roofline model on asymmetric multicores is more complex than the original and common assumption of “linear growth” (Section II-B). Specifically, there are four distinct segments marked as $\{1\} \sim \{4\}$ for the rooflines on ‘big core’ and ‘little core’, separated by the dotted vertical line. Each segment involves different levels of κ , putting different pressure on the L1 and L2 cache and thus leads to different κ – η relationship. In particular, η even decreases with increasing κ from 30 to 70 (i.e., the $\{2\}$ segment) on ‘little core’. We observe that this is primarily due to the increasing L1-I cache misses. As the ‘little core’ is an in-order processor, it stalls until instructions become available, severely affecting its performance. Second, we use tcomp32 algorithm [57] as an example and mark the κ of each aforementioned stream compression steps as dashed vertical lines in Figure 3. We can see that when κ is larger than 25, it is more and more cost-effective (i.e., lead to more performance gain) to run tasks on ‘big cores’. Consequently, different steps of the same stream compression procedure should be scheduled independently as they may be better scheduled to different cores due to their different κ as shown in the figure.

Observation 2: there are large differences of communication costs among asymmetric cores. To exploit task-core affinities, decomposed tasks may be scheduled to different cores. Therefore, different core communication paths (i.e., $c0$, $c1$, and $c2$ in Figure 2) may be involved. They have large differences in terms of bandwidth and latency as illustrated in Table II, measured by the STREAM benchmark [58]. The intra-cluster communication has much higher bandwidth and lower latency than inter-cluster communication, due to the slower CCI500 interconnection. More interestingly, inter-

TABLE II
BANDWIDTH AND LATENCY OF CROSS-CORE COMMUNICATION IN RK3399

Path	Bandwidth	Latency
intra-cluster $c0$	2.7 GB/s	70.4 ns
inter-cluster $c1$	0.7 GB/s	142.4 ns
inter-cluster $c2$	0.4 GB/s	420.8 ns

cluster communications of different directions (i.e., $c1$ and $c2$) are not involving the same cost. This is because of the additional *synchronization* and *hand-shaking* cycles [59] when sending data from little cores to big cores. Those asymmetric communication effects make the exploiting of task-core affinities a non-trivial quest.

In summary, these observations challenge the direct adoption of existing schemes for parallelizing stream compression on asymmetric multicores: First, existing mechanisms [41], [33], [32] consider the coarse-grained scheduling of the entire workload and do not expose the fine-grained task-core affinities in the workload. Second, previous studies on utilizing asymmetric multicores [53], [52] surprisingly overlook the different costs of $c1$ and $c2$, which is important to consider when scheduling decomposed tasks from a stream compression procedure that involves heavy inter-task communications. To achieve better utilization of asymmetric multicores in parallelizing stream compression, we should have 1) a fine-grained and complete exposition of task-core affinity in different steps of a stream compression procedure, and 2) a precise cost modeling to guide the proper scheduling of decomposed tasks under the asymmetric communication effects.

B. Design Overview

We propose a novel framework to parallelize data stream compression algorithms on asymmetric multicores, called *CStream*. Figure 4 depicts the overall workflow of *CStream*. First, *CStream* applies fine-grained decomposition of a stream compression procedure driven by the operational intensities of compression steps (e.g., *read*, *encode*, and *write*) to better expose task-core affinities. The decomposition results in several *tasks*, which can run independently and communicate with each other via message passing. For instance, one task may just conduct the *encode* step instead of running the whole steps. To increase concurrency, a task may be further replicated to multiple replicas (e.g., t_0 and t_1) handling subsets of data streams. Second, the decomposed tasks are scheduled to asymmetric multicores to minimize total energy consumption (E) without violation of user-specified compressing latency constraint (L_{set}), guided by a novel cost model. The cost model estimates both energy consumption (e_i) and compressing latency (l_i) of each task: 1) the e_i is estimated by the operational intensity (κ_i) of each task, and 2) the l_i is the summation of computation latency (l_i^{comp}) and communication latency (l_i^{comm}) of the each task t_i . In particular, l_i^{comm} varies depending on where the task and its upstream tasks are scheduled. For instance, in the example

scheduling plan shown in Figure 4, t_2 needs to fetch data from t_0 and t_1 via the slow CCI500 interconnection channel.

IV. FINE-GRAINED DECOMPOSITION

In this section, we discuss the fine-grained decomposition of a stream compression procedure in detail.

A. Stream Compression Procedure

We generally classify existing stream compression algorithms into *stateless* and *stateful* catalogs depending on whether they utilize *states*. *CStream* supports the parallelization of both *stateless* and *stateful* algorithms on asymmetric multicores. In the following, we illustrate the code template and concert examples of both types of stream compression algorithms.

Stateless Stream Compression. Algorithm 1 depicts the high-level idea of how a stateless stream compression algorithm (e.g., *tcomp32* algorithm [57]) works. It involves three steps for every batch of data streams: $s0$, $s1$, and $s2$. First, it *reads* a batch of tuples in step $s0$ before compressing. This step is mostly about memory copy, so it has low operational intensity. Second, it *encodes* each tuple in $s1$ by finding its compressible parts. This step typically involves arithmetic and logical operations in searching compressibility and therefore leads to higher operational intensity than $s0$. Third, it *writes* the compressed data to the output stream in $s2$ according to what $s1$ has encoded. This step involves both integer/float operation and memory access, and typically has a middle level of operational intensity compared with $s0$ and $s1$.

We use the *tcomp32* as a concrete example for stateless stream compression, as illustrated in Algorithm 2. The *tcomp32* conducts bit-level null suppression [57] over each non-overlapping continuous 32-bit in the data stream, and it considers leading zero bits of the integer as *compressible parts*. The incompressible part (i.e., which is the $n - \text{bit}$ number observed after cutting off leading zeros) and its length indicator are encoded thereafter and written to the output.

Algorithm 1: Stateless stream compression

Input: input stream *inData*

Output: output stream *outData*

```

1 while inData is not stopped do
2   ( $s0$ ) read the tuples from inData ;
3   ( $s1$ ) encode by finding the compressible parts;
4   ( $s2$ ) write compressed data to outData ;
5 end
```

Stateful Stream Compression. Algorithm 3 depicts the high-level idea of how a stateful stream compression algorithm works (e.g., *lz4* algorithm [60]). It involves five steps as $s0 \sim s4$. The read (i.e., $s0$) and write (i.e., $s4$) steps are the same as those in a stateless compression algorithm. In contrast, the encode step is now based on state, and can be further partitioned into three steps: $s1 \sim s3$. First, it preprocesses

Algorithm 2: *tcomp32* algorithm - a stateless stream compression algorithm

Input: input stream *inData*

Output: output stream *outData*

```

1 while not reach the end of inData do
2   /*  $s0$  */
3   number  $\leftarrow$  read next 32-bit from inData ;
4   /*  $s1$  */
5   if number = 0 then
6     n  $\leftarrow$  1 ;
7   else
8     n  $\leftarrow$  ceil(log2(number+1)) ;
9     // e.g.,  $n=2$  for number=3
10  end
11  /*  $s2$  */
12  write 5-bit  $n - 1$  to outData ;
13  write  $n$ -bit number to outData ;
14 end
```

some values before accessing the state (such as an index of a dictionary [60]) in $s1$. Second, it updates the state in-cache (or in-memory) in $s2$ with 1) the current tuple and 2) $s1$ -produced value. Third, it finally achieves encoding by state reference in $s3$. Due to the state write or read, $s2$ and $s3$ lead to lower operational intensity than $s1$.

We use *tdic32* and *lz4* as two instances of stateful stream compression, as presented in Algorithm 4 and 5, respectively. They similarly use a hash-table-based dictionary in their state to record previously compressed data, while *lz4* additionally requires a *bytePointer*, the so-called *literal* [60] part and a supplemental *buffer* to track more detailed information in its state, such as the length of matched patterns (*matchLength*). Such different state implementations between them also lead to different encoding patterns.

Algorithm 3: Stateful stream compression

Input: input stream *inData*

Output: output stream *outData*

```

1 while inData is not stopped do
2   ( $s0$ ) read the tuples from inData ;
3   ( $s1$ ) pre-process ;
4   ( $s2$ ) state update ;
5   ( $s3$ ) state-based encoding;
6   ( $s4$ ) write compressed data to outData ;
7 end
```

B. Parallelizing Stream Compression Procedure

CStream explores both pipelining and data parallelism for parallelizing a stream compression procedure as follows.

Exploring Pipelining Parallelism. First, *CStream* achieves pipelining parallelism by executing the aforementioned stream compression steps (Section IV-A) in a pipeline fashion. Specifically, each step in Algorithm 1 and 3 can run concurrently as an independent task. In this way, their varying operational intensity is exposed for further exploitation. When

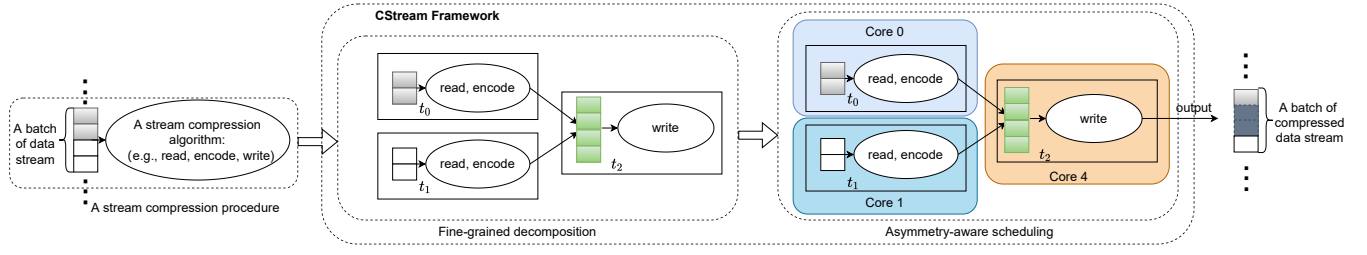


Fig. 4. The workflow of *CStream*. Using a stateless stream compression procedure as an example.

Algorithm 4: *tdic32* algorithm - a stateful stream compression algorithm

```

Input: input stream inData
Input: hash table tb with  $2^n$  entries
Output: output array outData
1 while not reach the end of inData do
    /* s0 */
2     number  $\leftarrow$  read next 32 bits from inData ;
    /* s1 */
3     index  $\leftarrow$  calculate hash value of number ;
    /* s2 */
4     result  $\leftarrow$  tb[index];
5     tb[index].word  $\leftarrow$  number ;
    /* s3 */
6     if result.word = number then
7         encoded  $\leftarrow$  (index < 1)|1 ;
8         bits  $\leftarrow$  (n+1) ;
9     else
10        encoded  $\leftarrow$  (number < 1)|0 ;
11        bits  $\leftarrow$  (33) ;
12    end
    /* s4 */
13    write (bits)-bit encoded to outData ;
14 end

```

communication latency (l_i^{comm}) of a task t_i is greater than the computation latency (l_i^{comp}) of this task or that ($l_{i'}^{comp}$) of its upstream task ($t_{i'}$), we fuse the tasks t_i and $t_{i'}$ to reduce communication overheads. For example, read (s_0) and encode (s_1) are fused, while write (s_2) runs separately as shown in Figure 4.

Exploring Data Parallelism. Second, when compressing latency constraints L_{set} is unable to meet through pipelining solely, we can replicate tasks to further explore data parallelism. We follow topologically sorted iterative scaling optimization [61] to replicate the bottleneck task. For each iteration, the bottleneck task t_i with the highest compressing latency l_i is replicated. Task replication iteration ends when compressing latency constraints can be met or hardware resources are saturated. We discuss the estimation of l_i to identify bottleneck tasks under varying scheduling plans in the next section.

The replication of all steps of both types of algorithms is straightforward, except s_2 of Algorithm 2, as it requires the manipulation of states. In our current implementation, we let each thread maintain its own state to avoid concurrent access conflicts. To validate the effectiveness of doing so, we parallelize the *tdic32* algorithm to compress the Rovio data, let each one of six working threads in s_2 step share the common concurrent dictionary (*share*) or maintain

Algorithm 5: *lz4* algorithm - a stateful stream compression algorithm

```

Input: input stream inData
Input: hash table tb
Input: maximum searching length ml
Output: output array outData
1 bytePointer  $\leftarrow$  0;
2 literal [60]  $\leftarrow$  [] ;
3 buffer  $\leftarrow$  [] ;
  // initially empty literal and buffer
4 while not reach the end of inData do
    /* s0 */
5     append 8 bits to buffer, reading from inData ;
    /* s1 */
6     number  $\leftarrow$  read the newest 32 bits from buffer ;
7     bytePointer  $\leftarrow$  bytePointer + 1 ;
    /* s2 */
8     index  $\leftarrow$  calculate hash value of number ;
    /* s3 */
9     result  $\leftarrow$  tb[index];
10    tb[index].offset  $\leftarrow$  bytePointer ;
11    tb[index].word  $\leftarrow$  number ;
12    clear contents in buffer older than bytePointer - ml ;
13    if result.word=number and bytePointer - result.offset <
        ml then
14        | expand searching in buffer and get the matchLength;
15    end
16    update literal ;
    /* s4 */
17    generate lz4-Token [60] from literal and matchLength ;
18    write the lz4-Token to outData ;
19 end

```

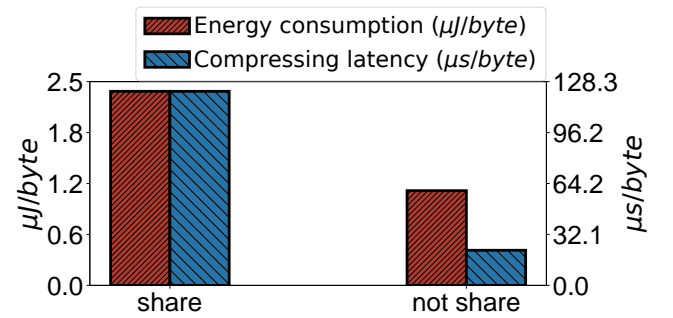


Fig. 5. Energy consumption and compressing latency comparison between whether or not share the state in *tdic32* algorithm and Rovio dataset.

its own private dictionary (*not share*), and compare the resulting energy consumption and compressing latency in Figure IV-B. Compared to sharing states with locks, this approach leads to a slight reduction (0.03 lower *compression ratio* [1], [62]) in compressibility, but significantly lower

(51%) energy consumption and lower (82%) compressing latency. Nevertheless, more advanced concurrent stateful stream processing mechanisms [63] may be applied in *CStream* to better parallelize stateful stream compression algorithms. Due to its considerable complexity, we leave it as future work.

V. ASYMMETRY-AWARE SCHEDULING

In this section, we introduce how *CStream* decides the optimal scheduling plan by the guidance of a novel cost model. We first give formal definitions to the scheduling problem, then illustrate the design of the cost model, followed by the procedures of model-guided task scheduling.

A. Problem Formulation

The goal of our scheduling is to minimize energy consumption while ensuring the user-assigned compressing latency constraint in parallelizing stream compression on asymmetric multicores. Specifically, we define a scheduling plan as follows.

Definition 2 (Scheduling Plan). *A scheduling plan is of mapping each decomposed task t_i to a specific core of asymmetric multicores. Given n tasks decomposed (i.e., $\{t_0, \dots, t_i, \dots, t_{n-1}\}$), there is an n -element array $p = \{j_0, \dots, j_i, \dots, j_{n-1}\}$ to describe a scheduling plan and j_i is the specific core of t_i to be mapped.*

We look for the optimal scheduling plan p_{opt} for a given stream compression procedure. The problem can be mathematically formulated as Equation 1~3.

$$\text{minimize}(E_{est} = \sum_i e_i) \quad (1)$$

$$\text{s.t., } \forall t_i \forall j,$$

$$L_{set} \geq L_{est} = \max(l_i) = \max(l_i^{comp} + l_i^{comm}), \quad (2)$$

$$C_j \geq \sum_{t_i \text{ at core } j} \eta_i \quad (3)$$

We refer to the energy consumption of each task t_i using the symbol e_i , and its compressing latency as l_i . Minimizing total energy consumption ($E_{est} = \sum_i e_i$) in Equation 1 is the major reasons for adopting asymmetric multicores in IoT. As the formulas show, we consider two categories of constraints that the optimization algorithm needs to make sure the scheduling plan satisfies. Constraint in Equation 2 enforces that compression latency should never exceed the user-specified constraint, due to real-time restrictions in stream analysis [8]. Due to the pipeline execution, the estimated compressing latency L_{est} is constrained by the maximum latency of each task, i.e., $L_{est} = \max(l_i)$. The l_i can be further decomposed into computation latency (l_i^{comp}) and communication latency (l_i^{comm}) as $l_i = l_i^{comp} + l_i^{comm}$. Constraint in Equation 3 enforces that the aggregated demand of instructions requested to any core ($\sum_{t_i \text{ at core } j} \eta_i$) must be smaller than its computation capacity (C_j). Task oversubscription has been studied in previous work [64], and is not the focus of this paper.

B. Cost Model

Our model estimates the energy consumption (e_i), instructions per unit time and energy (η_i and ζ_i , respectively), and compressing latency (l_i) of each task.

1) *Estimation of e_i* : The estimation of e_i is non-trivial as various properties of data and algorithms are involved. As shown in Equation 4, we estimate e_i as a proportional relationship to the instructions per unit time (η_i) and the latency (l_i), and an inverse proportional relationship to the instructions per unit energy (ζ_i).

$$e_i = \frac{\eta_i \times l_i}{\zeta_i} \quad (4)$$

2) *Estimation of η_i and ζ_i* : The instructions of unit time (η_i) is an intrinsic property of the asymmetric multicores. According to observations in Section III-A (especially the Figure 3), we use operational intensity (κ_i) to estimate η_i , by refining and formulating the aforementioned roofline graph [54]. Specifically, we estimate η_i as a four-region piece-wise linear function to κ_i , with L1 and L2 cache-awareness [65] for better modeling on asymmetric multicores, as shown in Equation 5.

Profiling of κ_i : Due to the single ISA property of asymmetric multicores, the operational intensity (κ_i) of each task t_i is not changing under varying scheduling plans. By definition, κ_i is related to the number of instructions and memory accesses of t_i . We feed each task t_i with a moderate size of data (which is large enough to prevent randomness and also within the memory capacity), and then profile the total number of instructions with *perf* [66]. The memory access is statically analysed by the specific step of stream compression algorithm.

Estimation of η_i : We can estimate η_i by using κ_i in Equation 5, and there are four segments as shown previously in Figure 3. First, when κ_i is low and does not put significant pressure on L1D (i.e., within the first boundary κ_{L1}), η_i grows relatively fast with a growth rate a_{L1} and intercept b_{L1} . Second, condition $\kappa_{L1} < \kappa_i \leq \kappa_{L2}$ means a higher κ_i that increases L1D missing and has significant impacts on the core. In this case, η_i grows slower with growth rate a_{L2} and intercept b_{L2} . Third, with κ_i increasing into $\kappa_{L2} < \kappa_i \leq \kappa_{roof}$, the L2 missing has major effects and therefore changes the a_{L2} and b_{L2} into $a_{ExceedL2}$ and $b_{ExceedL2}$, respectively. Fourth, after κ_i reaches the last boundary κ_{roof} , the η_i stays at the maximum value η_{max} instead of increasing with κ_i . When running on a certain core j , η_{max} is equal to its computational capacity (C_j). The specific value of the aforementioned parameters varies in ‘big cores’ and ‘little cores’ due to asymmetric computation effects, and we can use piece-wise linear fitting [67], [68] on the list of κ and its resulting η to acquire them.

$$\eta_i = \begin{cases} \kappa_i * a_{L1} + b_{L1} & , \kappa_i \leq \kappa_{L1} \\ \kappa_i * a_{L2} + b_{L2} & , \kappa_{L1} < \kappa_i \leq \kappa_{L2} \\ \kappa_i * a_{ExceedL2} + b_{ExceedL2} & , \kappa_{L2} < \kappa_i \leq \kappa_{roof} \\ \eta_{max} & , \kappa_i > \kappa_{roof} \end{cases} \quad (5)$$

Estimation of ζ_i : Another intrinsic property of asymmetric multicores, i.e., instructions of unit energy consumption (ζ_i) can also be estimated in a piece-wise linear form like Equation 5. Typically, the estimation of ζ_i involves different parameter values including the boundary of regions (i.e., κ_{L1} and κ_{L2}), the growth rate (i.e., a), and the intercept (i.e., b).

3) *Estimation of l_i :* The compression latency (l_i) of a task t_i is the sum of two non-overlapping components l_i^{comm} and l_i^{comp} .

Estimation of l_i^{comp} : We use l_i^{comp} to denote the general computation time spent for executing the task t_i . For tasks that have a constant workload characteristic, it is simply determined by the input size (i.e., i_i), and the mapped core j_i . As a result, Equation 6 depicts the simple linear relationship used to estimate it, with a system overhead ω_{j_i} of the mapped core j_i and a growth rate λ that are constant.

$$l_i^{comp} = \lambda \times i_i + \omega_{j_i}, \quad (6)$$

l_i^{comp} (under both big core and little core) can be acquired by a dry-run profiling. We can also use machine learning techniques (e.g., logistic regression) to train a prediction model to predict l_i^{comp} .

Estimation of l_i^{comm} : Task t_i involves communication delay l_i^{comm} to fetch the data from its upstream task t_u . l_i^{comm} is determined by the fetched data size (i.e., the i_i) and the relative distance to upstream task t_u . If a scheduling plan p maps t_i and t_u respectively at core j_i and j_u , l_i^{comm} can be estimated by Equation 7.

$$l_i^{comm} = \begin{cases} 0, & \text{if } j_i = j'_i \\ \frac{i_i \times L_{j'_i, j_i}^{comm}}{\text{cache line size}} + \omega_{j'_i, j_i}, & \text{otherwise} \end{cases} \quad (7)$$

, where j_i and j'_i are determined by p .

When task t_i is collocated with its upstream $t_{i'}$, or t_i is just at the beginning step (i.e., *read*), the communication latency l_i^{comm} is 0. Otherwise, it experiences the cross-core ($j_i \neq j_{i'}$) communication as discussed in Section III-A. Formula 7 estimates the communication latency based on the total size of data to be transferred (i_i), cacheline size, static overhead ($\omega_{j_{i'}, j_i}$) between core j_i and $j_{i'}$, and the worst unit communication latency ($L_{j_{i'}, j_i}^{comm}$) between core j_i and $j_{i'}$. It is worth noting that $L_{j_{i'}, j_i}^{comm} \neq L_{j_i, j_{i'}}^{comm}$ and $\omega_{j_{i'}, j_i} \neq \omega_{j_i, j_{i'}}$ if core j_i and $j_{i'}$ are located in different clusters, due to the previous observations from Table II. For each possible j_i and $j_{i'}$, their $L_{j_{i'}, j_i}^{comm}$ and $\omega_{j_{i'}, j_i}$ can be dry-run measured by setting up a producer thread at $j_{i'}$ and a consumer thread at j_i . Note that, the accurate estimation of l_i^{comm} is difficult as the fetched data size may vary depending on the input dataset size and

compressibility. In this work, we assume that the fetched data size does not vary much in a short period of time, given the same dataset and the same compression algorithm.

C. Model-guided Scheduling

We search p_{opt} by enumerating all possible plans with the cost model. We adopt dynamic programming [69] to speed up the plan searching as different plans may overlap for the scheduling of subsets of tasks. For each plan p_{enum} enumerated, we first predict the compressing latency (l_i) and instruction per unit time (η_i) on all of its tasks (t_i), according to Equations 5 ~ 6, and then check whether constraints in Equations 2 and 3 are met. If so, we continue to predict energy consumption (e_i) of all t_i by Equation 4 and get the total estimated energy consumption E_{est} of p_{enum} ; if not, we just ignore p_{enum} and continue search. The plan with minimal E_{est} and meet all constraints is p_{opt} .

D. Adaptive to Dynamic Environment

Our model is initially instantiated with a small number of input data (10 ~ 100 batches). However, data stream characteristics, such as dynamic range and entropy of data symbols, can vary over time, and *CStream* needs to be re-optimized in response to workload changes. To adapt to dynamic scenarios, we adopt a feedback-based regulation [70] in *CStream*. Specifically, we periodically (i.e., every 50 ms) measure the compressing latency and its predicted value. If the difference is larger than a threshold, we collect the subsequent batches of data to calibrate the cost model and replan the scheduling.

For a calibration iteration k , a model parameter x^k is tuned to x^{k+1} according to its absolute error e_a^k between measured (x_{mes}^k) and estimated (x_{est}^k) value, i.e., $e_a^k = x_{mes}^k - x_{est}^k$. Subsequently, an attempt to replan the scheduling is conducted by migrating from the previous scheduling plan based on the updated model with refreshed parameter x^{k+1} . Such parameter tuning process from x^k to x^{k+1} is handled by the PID controller [71], [72] with controller settings P , I , and D , as defined in Equation 8.

$$x_{est}^{k+1} = x_{est}^k + \delta^k \quad (8)$$

where $\delta^k = P(e_a^k - e_a^{k-1}) + I(e_a^k) + D(e_a^k - 2e_a^{k-1} + e_a^{k-2})$

It is worth noting that we use the incremental form [73], [74] of PID instead of position-based PID [71], [72], as position-based PID is more costly and less robust due to its notorious *integral saturation* effect [75]. The incremental PID calibration process is ended when relative error $|e_a^k/x_{est}^k|$ is small enough, and the new optimal scheduling plan is then achieved under the eventually calibrated parameter x . For simplicity, we consider computational latency l_i^{comp} and operational intensity κ_i as the calibratable x , adjust them independently and assume all other changes as marginal noises. This is because other variable parameters are either irrelevant to workload dynamics (e.g.,

the roofline model parameters describing the hardware) or hard to precisely acquire (e.g., the communication latency l_i^{comm}). The overhead of performing such feedback-based dynamic regulation is negligible as our cost model involves solving simple linear equations and rescheduling is incrementally conducted. However, its response may be lagged when facing a bursting workload and at least 3 times of calibration should be involved (as indicated by $k-2$, $k-1$ and k in Equation 8). More sophisticated controllers [76] that monitor workload statistical information in the datastream may achieve an even better response to workload changes but beyond the scope of this work.

VI. METHODOLOGY

In this section, we first introduce the examined competing mechanisms, followed by benchmark workloads including both varying compression algorithms and datasets. Then, we discuss the instrument of targeting performance metrics.

A. Competing Mechanisms

We compare *CStream* with the following five competing mechanisms in parallelizing the stream compression procedures: *OS*, *CS*, *RR*, *BO*, and *LO* as follows.

- (i) **Operating System (*OS*)**. The replicated tasks of the whole stream compression procedure (without decomposition) are scheduled by the Linux 5.10 kernel at thread level with the *energy-aware-scheduling* (*EAS*) [39] strategy.
- (ii) **Coarse-grained Scheduling (*CS*)**. Following prior work of coarse-grained workload scheduling on asymmetric multicores [32], we can schedule each replica of the entire stream compression procedure as a single task with our asymmetry-aware scheduling scheme.
- (iii) **Round Robin (*RR*)**. Under *RR*, we apply fine-grained decomposition of *CStream*, and the decomposed tasks are scheduled in a round-robin manner, i.e., sequentially mapped to each core.
- (iv) **Big-core Only (*BO*)**. Under *BO*, the decomposed tasks are randomly scheduled to the big cores of rk3399 (core4 to core5), and little cores (core0 to core3) are left idle.
- (v) **Little-core Only (*LO*)**. Under *LO*, the decomposed tasks are randomly scheduled to the little cores of rk3399 (core0 to core3), and big cores (core4 to core5) are left idle.

Under *OS* and *CS*, the stream compression procedure is replicated into multiple tasks in achieving data parallelism (without decomposition and pipelining parallelism). The tasks are subsequently scheduled by Linux kernel under *OS*, or with asymmetry-aware scheduling of *CStream* under *CS*. Under *RR*, *BO*, and *LO*, the stream compression procedure is decomposed in a fine-grained manner as *CStream* does, but task scheduling is not asymmetry-aware.

B. Input Workloads

CStream supports varying stream compression algorithms and datasets. We select a wide range of algorithms and datasets with distinct characteristics for a comprehensive evaluation.

1) *Algorithms*: We focus on parallelizing the following three lightweight stream compression algorithms in our evaluation. Nevertheless, *CStream* can be easily extended to support other stream compression algorithms. 1) *tcomp32* cuts off unused bits of each 32-bit symbol, which is stateless stream compression following Algorithm 1's abstraction. 2) *lz4* [60] is a popular LZ77-based [44] stateful compression (i.e., Algorithm 3). It uses a hash table as its state to replace the traditional *dictionary* of LZ77. 3) *tdic32* is a simplified variable length coding created by combining the two above, which is also stateful. Specifically, it borrows the hash table from *lz4* and employs a memory I/O pattern similar to *tcomp32* (i.e., byte-unaligned encoding for each 32-bit single symbol).

2) *Datasets*: We use three real-world and one carefully designed synthetic dataset in our evaluation. These datasets cover varying statistical properties, such as 1) *vocabulary duplications* [44], 2) *dynamic range of the symbol* [62], [57], and 3) *symbol entropy* [19]. Since the *tcomp32*, *lz4* and *tdic32* share a 32-bit reading of data, we define data within 32-bit as a *symbol*, while more than 32-bit of data is treated as a *vocabulary*. The datasets are as follows. 1) *Sensor* [77] represents a type of full-text streaming data that is generated by automated sensors. Its most compressible part is the *symbol entropy*, which is packed in an XML format with only ASCII code. Furthermore, the XML pattern can result in partial *vocabulary duplication*. We let each 16 ASCII characters in *Sensor* form one 128-bit tuple in our evaluation. 2) *Rovio* [78] continuously monitors the user actions of a given game to ensure that their services work as expected, and is packed in (64-bit key, 64-bit payload). Its high key duplication leads to significant *vocabulary duplication* [79]. 3) *Stock* [80] is a real-world stock exchange dataset packed in (32-bit key, 32-bit payload) binary format. Unlike *Rovio*, its key duplication is much lower. 4) *Micro* is a synthetic dataset used to easily evaluate the impact of varying workload properties. Each tuple in *Micro* is a 32-bit plain value.

C. Instrument of Performance Metrics

Throughout this study, we focus on two important performance metrics. The first is *compressing latency constraint violation* (*CLCV* for short). For each test, we repeat the measurement 100 times, and *CLCV* refers to fraction of measurements violating L_{set} and the total measurements. The second is *energy consumption*, denoted as E_{mes} . We let E_{mes} refer to the overhead for compressing each unit of data (i.e., in $\mu J/byte$). The system overhead of each mechanism, such as profiling and scheduling in *CStream* are included in E_{mes} . We develop an energy meter as shown in Figure 6 to measure E_{mes} . It consists of the Texas Instrument's INA226 [81] chip as the sensor for current/voltage, and the Espressif's ESP32S2 [82] micro control unit (MCU) for data pre-processing and USB-2.0 communication with targeting asymmetric multicores. In this way, it provides accurate measurement with low overhead.

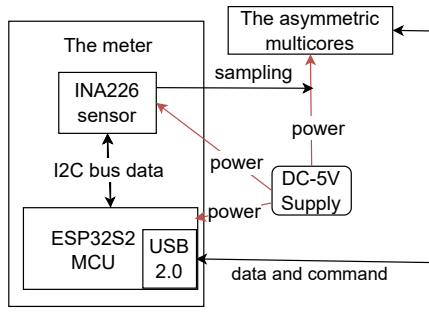


Fig. 6. The design of energy meter.

TABLE III
SPECIFICATION OF THE ASYMMETRIC MULTICORES BOARD.

Model	Radxa Rockpi 4a
Processor	4 A53 'little core' and 2 A72 'big core'
Memory	2GB LPDDR4, 64 bit dual channel @ 3200Mb/s
Software	build root [84] 2021-08-rc2, with Linux 5.10 [85] kernel

VII. EVALUATION

We use the Radxa Rockpi 4a [83] for evaluation. This platform is equipped with a rk3399 asymmetric multicores processor. Detailed specifications are shown in Table III. By default, each core runs at its highest frequency (i.e., 1.8GHz for A72, and 1.416GHz for A53). To eliminate the impact of network transmission overhead, the input datasets are first populated (synthetic dataset) or loaded (real datasets) in memory. We set the batch size (B) as 932,800 bytes and L_{set} as $26\mu s/byte$ unless otherwise stated.

A. End-to-end Comparison

In this section, we show the end-to-end comparison between *CStream* and five competing mechanisms.

Energy Consumption Comparison. Figure 7 reports the energy consumption of different mechanisms on handling different datasets. In general, we can see that *CStream* always leads to the least energy consumption. For example, *CStream* can save up to 53.23% energy consumption on the *lz4-Stock* procedure compared with *BO*. *CStream* determines the core-mapping of decomposed stream compression tasks according to their varying operational intensity. In contrast, *LO* and *BO* underutilized either the big cores or little cores on asymmetric multicores, while operational-intensity-unconscious parallelization (i.e., *OS*), or coarse-grained parallelization (i.e., *CS*) fails to explore the suitable task-core mapping between stream compression procedure and asymmetric multicores.

Compressing Latency Constraint Violation Comparison. Figure 8 presents the *CLCV* of different mechanisms on handling different datasets. We can see that *CStream* can always avoid any compressing latency constraint violations for the evaluated workloads. We find that it is mainly because

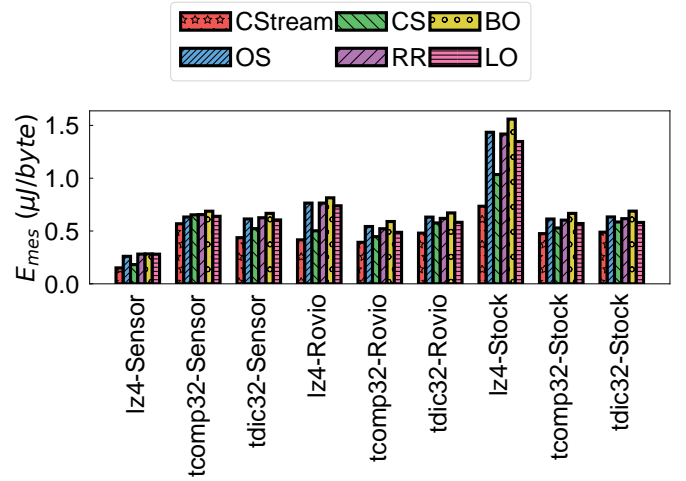
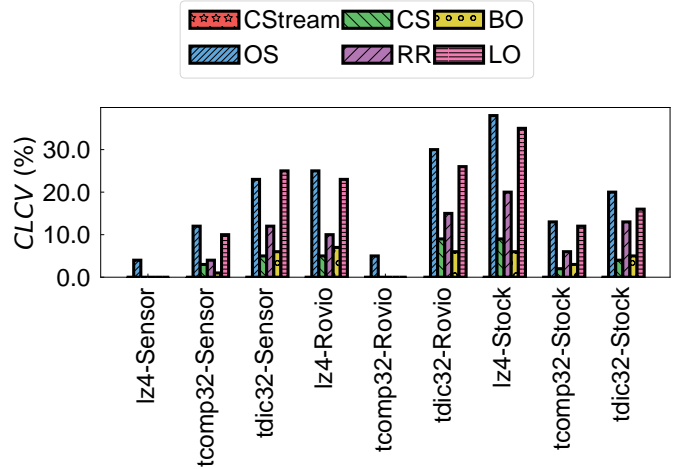
Fig. 7. Energy consumption (E_{mes}) comparison.

Fig. 8. Compressing latency constraint violation (*CLCV*) comparison. The *CLCV* of *CStream* is always zero, and the *CLCV* of *CS*, *BO*, *RR* and *LO* under *tcomp32-Rovio* and *lz4-Sensor* and are also zero.

the statistical characteristics of datasets in our experiments are relatively stable over time. We show the evaluation of more dynamic workloads shortly later. In contrast, *LO* and *RR* fail to preserve the compressing latency constraints due to the low utilization of high performance 'big cores' in conducting high operational intensity steps (e.g., the s_2 in Algorithm 1) of the stream compression procedure. On the contrary, *CS* and *BO* underutilize 'little cores' in conducting low operational intensity steps (e.g., the s_0 in Algorithm 1). In particular, *CS* tries to firstly schedule as much as possible to 'big cores' before utilizing the 'little cores', as the operational intensity of the whole stream compression procedure is relatively high (e.g., about 200 for *tcomp32-Rovio*). *OS* fails for two reasons. First, its frequent migration of tasks leads to extra overhead. Our further investigation reveals that *OS* scheduler involves about 60,000 context switches while *CStream* only involves about 10 context switches for compressing every

megabyte of input data. Second, *OS* treats stream compression tasks as a black box, and it hence leads to a relatively inaccurate estimating of latency and causes compressing latency constraint violations.

Dynamic Workloads. The stream properties (e.g., entropy or dynamic range of symbols) may change on the fly. We now evaluate the adaptivity of *CStream* to the dynamic workload by using the *tcomp32-Micro* procedure. We set the dynamic range of symbols to 500 at the beginning, and increase it to 50000 immediately after the whole fifth batch is compressed. We configure the aforementioned incremental PID controller (Equation 8) with $[P, I, D]$ as $[0.1, 0.85, 0.05]$ under the guidance of well-known *PSO tuning* [86]. We allow a 0.1 maximum value of relative error $|e_a^k/x_{est}^k|$. The energy consumption and compressing latency constraint violation of *CStream* with and without feedback-based regulation (Section V-D) are demonstrated in Figure 9. We observe that the compressing latency constraint will be violated after workload change if there is no feedback-based regulation, as the previous profiling leads to an inaccurate estimation of compressing latency. With the feedback-based regulation, *CStream* is able to adapt to the changed workload eventually from the ninth batch and switch to the new optimal scheduling plan, which requires higher energy consumption to avoid compressing latency violation. There are some vibrations during the adapting process, such as unnecessarily high energy consumption at the eighth batch, due to the continuous calibration process with the PID-based controller (Section V-D). It is worth noting that the overhead of such an adapting process is marginal, i.e., 0.9% energy consumption and 6% processing time of compressing each batch of data.

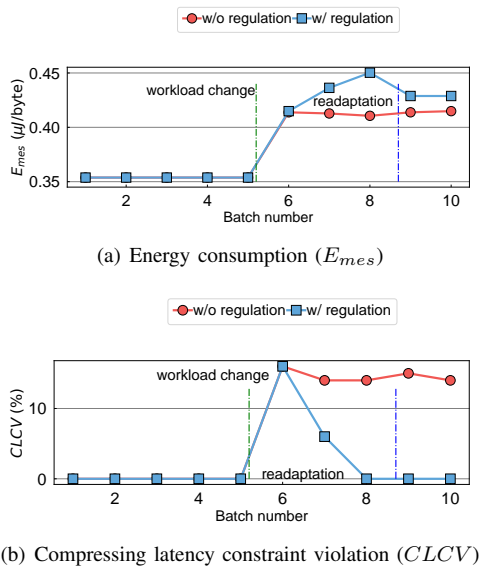


Fig. 9. Adaptation to dynamic workload. Workload changes after the 5th batch and readaptation finish at the 9th batch.

B. Workload Sensitivity Study

In this section, we show the superiority of *CStream* under varying workload characteristics, including procedure settings and data statistic properties.

1) *Procedure Settings:* We vary the compressing latency constraint (L_{set}) and batch size (B) of *tcomp32-Rovio* stream compression procedure as follows.

Varying Compressing Latency Constraint (L_{set}). We show the impact of varying compressing latency constraints shown in Figure 10. While *OS*, *RR*, *BO*, and *LO* have constant energy consumption, *CStream* and *CS* achieve more on energy saving under a larger L_{set} , and generate a lower latency scheduling plan under a smaller L_{set} . However, the coarse-grained way of *CS* leads to 1) higher energy consumption at each L_{set} when comparing with *CStream* and 2) underutilization of ‘little cores’ and failure of meeting the tight settings (i.e., smaller) of L_{set} .

Varying Batch Size (B). We study the energy consumption under different sizes of batch (i.e., B) in Figure 11. We observe that the energy consumption remains nearly stable when there is a large enough batching of data (i.e., $B > 10^3 byte$). A small batch size may increase energy consumption slightly due to the frequent cache thrashing [79]. In general, energy consumption is determined $\frac{\eta_i}{\zeta_i}$ proportion (Equation 4) of each task t_i , and both η_i (instructions per unit time) and ζ_i (instructions per unit energy) are determined by the task’s operational intensity (κ_i). *CStream* can select the minimal value of $\frac{\eta_i}{\zeta_i}$ for each task to minimize energy consumption, by utilizing 1) the fine-grained decomposed operational intensity (κ_i) when compared with *CS*, and 2) the asymmetry-aware scheduling when comparing with *OS*, *RR*, *LO* and *BO*.

2) *Data Statistic Properties:* We tune the data statistic properties including vocabulary duplication, symbol duplication, and dynamic range. Specifically, we use the synthetic dataset *Micro* to study their impacts.

Vocabulary Duplication. We first conduct *lz4-Micro* procedure under varying vocabulary duplication. The *lz4* is most sensitive to vocabulary duplication, and operational intensity (κ_i) of *lz4* tasks changes differently with vocabulary duplication due to their different functions. For instance, the κ_i of tasks conducting *state update* (*s2* in Algorithm 3) decreases with increasing vocabulary duplication, as the hash table of state can be updated less. However, the κ_i of tasks conducting *state-based encoding* (*s3* in Algorithm 3) will increase, as *lz4* is more likely to conduct ‘backward searching’ for expanding match [60]. In Figure 12, we vary the vocabulary duplication from low to high and have two observations. First, *CStream* is able to decompose the *lz4* procedure and select the best scheduling plan of tasks no matter how vocabulary duplication changes. Second, each mechanisms has its highest energy consumption when duplication is moderate, this is because the aforementioned different tendencies are reconciled under moderate duplication and leads to a total maximum energy consumption.

Symbol Duplication. We next conduct the *tdic32-Micro* procedure while varying the duplication of the 32-bit symbol.

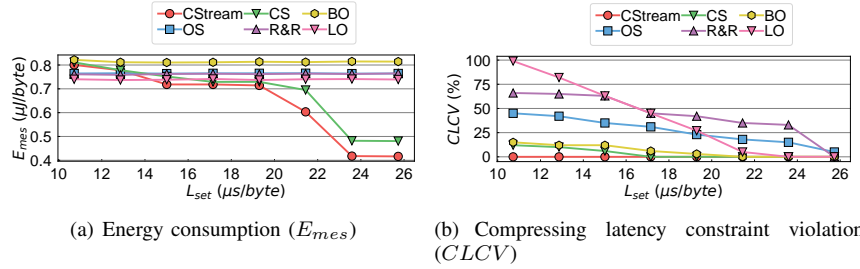
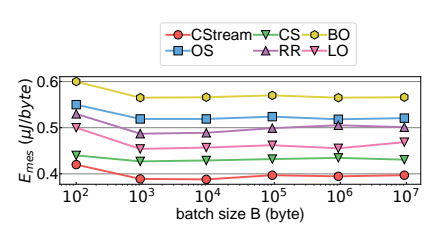
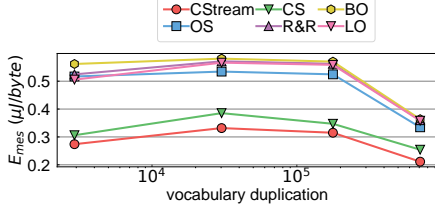
Fig. 10. Impacts of varying L_{set} .Fig. 11. Impacts of varying batch size B .

Fig. 12. Impacts of varying vocabulary duplication.

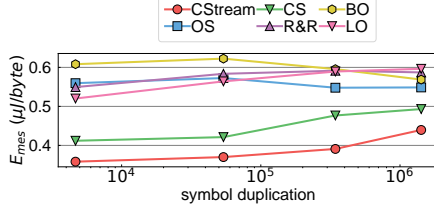


Fig. 13. Impacts of varying symbol duplication.

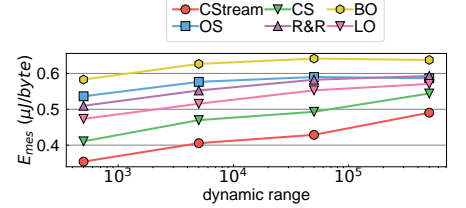


Fig. 14. Impacts of different dynamic range.

The `tdic32` is most sensitive to symbol duplication, and operational intensity (κ_i) of most `tdic32` tasks decreases with the increasing symbol duplication. While tasks conducting s_0 and s_1 remain nearly the same, higher symbol duplication leads to 1) less state update for tasks conducting s_2 and 2) shorter average encoding and writing for tasks about s_3 and s_4 . The results of parallelizing `tdic32` are shown in Figure 13. We can see that *LO* has an increasing energy consumption with increasing symbol duplication. This is because when symbol duplication increases, more `tdic32` tasks have their operational intensity κ_i dropped to the $30 \sim 70$ region, which suffers from the in-order stall in L1-I cache as depicted in Figure 3. In contrast, *BO* becomes more energy efficient with increasing symbol duplication as ‘big cores’ are out-of-order. *CStream*, *OS*, *CS* and *RR* utilize both big and little cores, and their different utilization on the asymmetric multicores lead to different energy consumption at changing symbol duplication. Nevertheless, *CStream* is always able to achieve the least energy consumption, which reaffirms its superiority compared to competing mechanisms.

Dynamic Range. Finally, the results of `tcomp32-Micro` procedure under varying dynamic range of the 32-bit symbol, are shown in Figure 14. The `tcomp32` is most sensitive to dynamic range, as increasing dynamic range makes its arithmetic computation in s_1 and writing output data in s_2 more costly. Therefore, operational intensity (κ_i) and compressing latency (l_i) of most tasks t_i in the procedure are increased, resulting in a higher energy consumption according to Equations 4. *CStream* always outperforms others, but its energy saving becomes less significant when the dynamic range is high, as there is less room for optimization.

C. System Configuration Analysis

In this section, we evaluate the system sensitivity by tuning core frequency statically and dynamically. The `tcomp32-Rovio` procedure is used as illustration example.

Static Frequency Regulation. We vary core frequency statically and evaluate how it affects the measured compressing energy consumption (E_{mes}). The results are shown in Figure 15. Obviously, the *CStream* outperforms other mechanisms under varying frequency settings. It is also worth noting that low frequency doesn’t imply lower energy consumption. Although lower frequency results in low power (i.e., measured in Watts), their increased latency in stream compression may lead to even higher energy consumption, especially when tasks run on ‘little cores’.

Dynamic Frequency Regulation. We also consider to dynamically regulate the frequency by using the DVFS [38], [40], [87]. we report how each of the six mechanism cooperates with different DVFS strategies in Figure 16. We use the “default” strategy for comparison by fixing each core at its highest frequency. The “conservative” and “on-demand” are two DVFS methods trying to reduce energy consumption by frequency reconfiguration on the fly, and their major difference is that “conservative” strategy changes frequency less when compared with “on-demand”.

We have three observations here. First, *CStream* always achieves the least energy consumption and least compressing latency constraint, regardless of using what kind of DVFS strategy. Second, the “conservative” strategy can further reduce energy consumption for all mechanisms compared with their default conditions. However, the compressing latency constraint violation for all mechanism are increased by using such strategy. This because the “conservative” DVFS only offers *relative coarse-grained guarantee of meeting the compressing latency constraints*. Specifically, there is

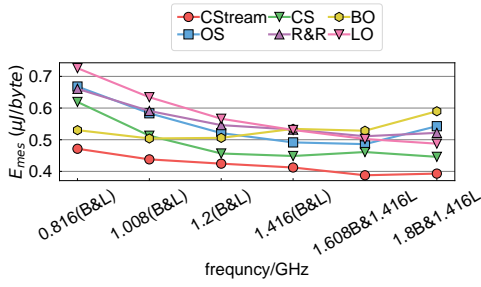


Fig. 15. Impacts of statically varying core frequency. “B” denotes the big cores while “L” denotes the little cores.

uncertain extra overhead in dynamic frequency regulation, leading to more compressing latency constraint violations. Third, the “on-demand” strategy leads to no improvement but more compressing latency constraint violation and energy consumption. This is because it changes the frequency too often and involves much extra overhead in the frequency switching.

D. Break-down Analysis

Since *CStream* incorporates fine-grained decomposition and asymmetry-aware (i.e., both computation and communication) task scheduling as two key contributions, we conduct a break-down analysis to study their impacts respectively. We use *tcomp32-Rovio* procedure as the illustration example, and the following break-down factors will be studied.

- *simple* refers to following a symmetric-multicore-aware parallel data compression [88], which only exploits the data parallelism. Specifically, the whole procedure is treated as a task t_{all} . As t_{all} fails to meet the compressing latency constraint, it is then replicated into two equivalent tasks (each one is denoted as a t_{re}).
- *+decom.* adds the fine-grained decomposition on stream compression and enables the basic exposition of task-core affinity as discussed in Section IV. Specifically, the procedure is decomposed into two tasks namely t_0 (conducting s_0 and s_1 in Algorithm 1) and t_1 (conducting s_2 in Algorithm 1). Both t_0 and t_1 will be randomly scheduled to asymmetric multicores.
- *+asy-comp.* adds the asymmetric-computation-awareness to schedule t_0 and t_1 , including all modeling in Section V-B but ignoring the asymmetric communication effects. Specifically, the unit communication latency L_{j_i', j_i}^{comm} and $L_{j_i, j_i'}^{comm}$ in Equation 7 are treated the same for any core j_i' and j_i that non-equal.
- *+asy-comm.* adds the consideration of asymmetric communication to schedule t_0 and t_1 , and is the fully functional *CStream* as introduced in Section III.

Impacts of Fine-grained Decomposition. By comparing *+decom.* with *simple*, we found that the fine-grained decomposition provides opportunities for better utilization of both big cores and little cores, which can reduce energy consumption a lot. To further comprehend the effects of

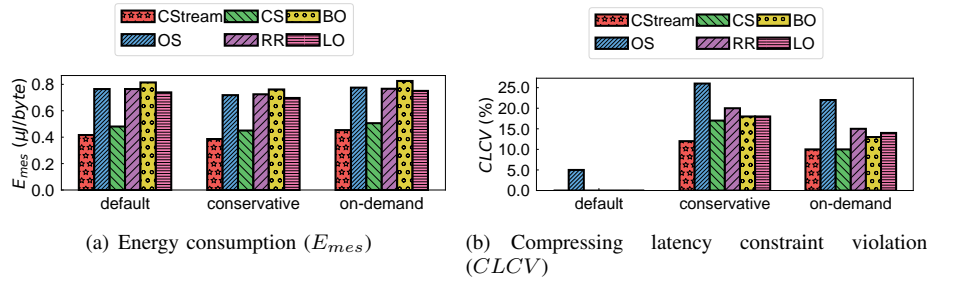


Fig. 16. Impacts of DVFS strategies.

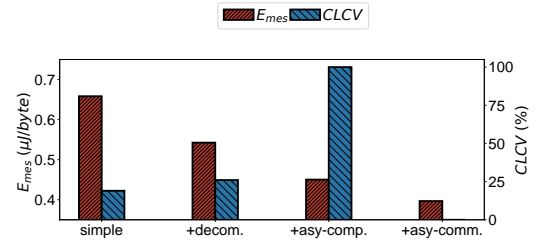


Fig. 17. Factor analysis about energy dissipation consumption (E_{mes}) and compressing latency constraint violation ($CLCV$).

decomposition, we compare the operational intensity (κ), compressing latency (l) and energy consumption (e) of the decomposed tasks t_0/t_1 , single thread all procedure (t_{all}), and the replicated version of t_{all} (i.e., $t_{re} \times 2$) in Table IV. There are three takeaways. First, t_0 is better to be scheduled to ‘big cores’ than other tasks due to its highest operational intensity. Specifically, it can reduce 53% compressing latency when scheduled to a ‘big core’ instead of a ‘little core’, with only 8% increased energy consumption. Second, t_{all} and $t_{re} \times 2$ lead to much under utilization of asymmetric multicores, as they simply *reconcile* the operational intensity of t_0 (i.e., 340) and t_1 (i.e., 102) to a medium value 220, ignoring the fact that they have large difference and should be treated differently. Third, only applying *+decom.* is far from optimal. Specifically, both t_0 and t_1 should be scheduled to the “right place” (i.e., big core and little core respectively as shown in table) in achieving energy saving or compressing latency reduction, but *+decomp.* schedules them randomly to asymmetric multicores.

Impacts of Asymmetry-aware Scheduling. The asymmetry-aware scheduling includes the awareness of both asymmetric computation and asymmetric communication, and we can observe their impacts respectively from *+asy-comp.* and *+asy-comm.* First, *+asy-comp.* adds asymmetric computation awareness to guide the scheduling of fine-grained tasks t_0 and t_1 . It can always correctly determine the varying task-core affinity among them according to their different operational intensity (Table IV), which reduces energy consumption. Nevertheless, due to the ignorance of asymmetric communication effects, it is too aggressive

TABLE IV

COMPARISON AMONG DECOMPOSED TASK t_0/t_1 , SINGLE THREAD ALL PROCEDURE t_{all} , AND REPLICATED t_{all} (DENOTED AS $t_{re} \times 2$).

Task	operational intensity κ	compressing latency l ($\mu s/byte$)		energy consumption e ($\mu J/byte$)	
		big core	little core	big core	little core
t_0	320	15.0	32.6	0.29	0.27
t_1	102	13.5	21.7	0.32	0.10
t_{all}	220	28.3	53.2	0.59	0.34
$t_{re} \times 2$	220	15.0	27.1	0.75	0.51

for energy saving and violates the latency constraint (L_{set}) frequently. Second, *CStream* further adds asymmetric communication awareness (i.e., *+asy-comm.*) compared to *+asy-comp.* It can model and plan well before executing t_0 and t_1 , and thus prevent the compressing latency constraint violation while ensuring the least energy consumption.

As the asymmetry-aware scheduling is guided by cost model, we evaluate its correctness here. We show the relative error rate associated with estimated compressing latency (L_{est}) or energy consumption (E_{est}) per procedure by our model. Specifically, we define $relative_error^L = \frac{|L_{pro} - L_{est}|}{L_{pro}}$ and $relative_error^E = \frac{|E_{pro} - E_{est}|}{E_{pro}}$, where L_{pro} and E_{pro} are the measured compressing latency and energy consumption of the tested procedure, respectively. Table V shows the model accuracy of all evaluated algorithms under their optimal scheduling plans in compressing Rovio. Overall, our estimation approximates the measurement well for the L_{est} and E_{est} of all three algorithms. It is therefore able to determine the optimal scheduling plan as shown previously. The inaccuracy is mainly caused by the difficulty in accurately estimating the unit overhead of communication ($L_{j',j}^{comm}$) on the fly, as it is affected by multiple factors such as memory access patterns and hardware prefetcher units.

TABLE V

MODEL CORRECTNESS UNDER OPTIMAL SCHEDULING PLANS.

Estimation object	method	lz4	tcomp32	tdic32
	variable			
Compressing Latency	$L_{est}(\mu s/byte)$	25.5	23.2	23.3
	$L_{pro}(\mu s/byte)$	23.6	21.7	25.3
	$relative_error^L$	0.08	0.07	0.08
Energy Consumption	$E_{est}(\mu J/byte)$	0.47	0.43	0.44
	$E_{pro}(\mu J/byte)$	0.42	0.40	0.48
	$relative_error^E$	0.14	0.08	0.09

VIII. RELATED WORK

In this section, we review the related work and reveal the limitations that motivate this work.

Parallel Data Compression. A lot of compression algorithms have been proposed since 1950s [89], focusing on improving the theoretical compressibility [16] and reducing compressing complexity [18], [19]. On the FPGA, Milward et al. [90] implemented a novel dictionary-based parallel

compression implementation, Sano et al. [91] achieved float-point data compression, Tian et al. [92] proposed parallel compression for scientific data, and Bark et al. [22] parallelized the lz4 algorithm. On the GPU, the parallelization of several LZ algorithms [44] have been conducted such as [93] and [94], and Huang et al. [95] have also parallelized trajectory-specific compression algorithms. Due to the significant difference in hardware architectures, it is unclear how can those works be applied to parallel compression on asymmetric multicores. Researcher have also utilized the Symmetric Multicore Processors (SMPs) for parallel compression: to compress the float-point data [35] and to drill the inner parallelism (in a SIMD-like manner) of LZ77 [96]. Recently, Knorr et al. [97] achieved high-throughput of large volumes of scientific data scaling up to 24 threads, and Dua et al. [36] compressed the hyperspectral images on supercomputer with a hyper-cube structure. Our work differ significantly from existing works from three aspects: 1) hardware architecture (i.e., we focus on asymmetric multicores), 2) compression algorithms (i.e., we conduct data stream compression), and 3) performance metrics (i.e., we consider both compressing latency constraint violation and energy consumption).

Efficient Utilization of asymmetric multicores.

The approaches to effectively manage the asymmetric multicores include modeling [24] and predicting [25] of performance/energy, CPI stack [43] and DVFS-collibrated scheduling [87], [98], [26]. The performance/energy model has been merged into the popular mainline Linux from version 5.0 onwards as the *Energy Aware Scheduling (EAS)* [99]. Yu et al. [41] recently proposed a collaborative OS scheduler addressing comprehensive multiple optimization objects on asymmetric multicores. However, these models treat the software running on asymmetric multicores as a black box due to the isolation by OS and system-level statistics, which overlooks many optimization opportunities as demonstrated in our experiments. There are also a few existing user-space research projects working on energy or latency optimization on asymmetric multicores, including virtual machine [31], web browser [33], game governor [100] and artificial intelligence framework [32]. However, all of them focus on scheduling the whole workload, which is relatively coarse-grained. For instance, Wang et al. [32] optimized the matrix operation tasks which are entirely computation intensive. Our work is orthogonal to them as we especially exploit the fine-grained behaviour (i.e., the different operational intensity among steps, Section III-A) of stream compression procedures.

IX. CONCLUSION

This paper introduced *CStream*, a novel framework to parallelize stream compression on asymmetric multicores. *CStream*'s superiority is gained by both fine-grained decomposition and asymmetry-aware scheduling strategy. We have experimentally demonstrated that *CStream* achieves the following desired properties: 1) when the compressing latency constraint (L_{set}) by the user is relatively loose, it

can achieve the least energy consumption; and 2) when encountering a tight L_{set} , its latency constraint violation is always minimized. In the future, we plan to further exploit *CStream* on more stream compression algorithms and on other hardware architectures such as Intel Agilex and Nvidia Jetson to achieve energy-efficient and low latency stream compression for a wide range of IoT applications.

REFERENCES

- [1] G. Pekhimenko, C. Guo, M. Jeon, P. Huang, and L. Zhou, "Tersecades: Efficient data compression in stream processing," in *2018 USENIX Annual Technical Conference (USENIX ATC 18)*. Boston, MA: USENIX Association, Jul. 2018, pp. 307–320. [Online]. Available: <https://www.usenix.org/conference/atc18/presentation/pekhimenko>
- [2] G. Theodorakis, F. Kounelis, P. Pietzuch, and H. Pirk, "Scabbard: Single-node fault-tolerant stream processing," *Proc. VLDB Endow.*, vol. 15, no. 2, p. 361–374, oct 2021. [Online]. Available: <https://doi.org/10.14778/3489496.3489515>
- [3] J. Azar, A. Makhoul, M. Barhamgi, and R. Couturier, "An energy efficient iot data compression approach for edge machine learning," *Future Generation Computer Systems*, vol. 96, pp. 168–175, 2019.
- [4] D. Zordan, B. Martinez, I. Vilajosana, and M. Rossi, "On the performance of lossy compression schemes for energy constrained sensor networking," *ACM Transactions on Sensor Networks (TOSN)*, vol. 11, no. 1, pp. 1–34, 2014.
- [5] C. J. Deepu, C.-H. Heng, and Y. Lian, "A hybrid data compression scheme for power reduction in wireless sensors for iot," *IEEE transactions on biomedical circuits and systems*, vol. 11, no. 2, pp. 245–254, 2016.
- [6] A. Ukil, S. Bandyopadhyay, and A. Pal, "Iot data compression: Sensor-agnostic approach," in *2015 data compression conference*. IEEE, 2015, pp. 303–312.
- [7] P. Shilane, M. Huang, G. Wallace, and W. Hsu, "Wan-optimized replication of backup datasets using stream-informed delta compression," *ACM Transactions on Storage (ToS)*, vol. 8, no. 4, pp. 1–26, 2012.
- [8] S. Zeuch, A. Chaudhary, B. D. Monte, H. Gavrilidis, D. Giouroukis, P. M. Grulich, S. Breß, J. Traub, and V. Markl, "The nebula stream platform for data and application management in the internet of things," in *CIDR 2020, 10th Conference on Innovative Data Systems Research, Amsterdam, The Netherlands, January 12-15, 2020, Online Proceedings*. [www.cidrdb.org](http://cidrdb.org/cidr2020/papers/p7-zeuch-cidr20.pdf), 2020. [Online]. Available: <http://cidrdb.org/cidr2020/papers/p7-zeuch-cidr20.pdf>
- [9] M. Bansal, I. Chana, and S. Clarke, "A survey on iot big data: Current status, 13 v's challenges, and future directions," vol. 53, no. 6, 2020. [Online]. Available: <https://doi.org/10.1145/3419634>
- [10] A. Lavric, V. Popa, and S. Sfichi, "Street lighting control system based on large-scale wsn: A step towards a smart city," in *2014 International Conference and Exposition on Electrical and Power Engineering (EPE)*, 2014, pp. 673–676.
- [11] V. A. Memos, K. E. Psannis, Y. Ishibashi, B.-G. Kim, and B. B. Gupta, "An efficient algorithm for media-based surveillance system (eamsus) in iot smart city framework," *Future Generation Computer Systems*, vol. 83, pp. 619–628, 2018.
- [12] P. Rizwan, K. Suresh, and M. R. Babu, "Real-time smart traffic management system for smart cities by using internet of things and big data," in *2016 international conference on emerging technological trends (ICETT)*. IEEE, 2016, pp. 1–7.
- [13] (2020) Eclipse iot working group. iot developer survey 2018. [Online]. Available: <https://blogs.eclipse.org/post/benjamin-cab%C3%A9/key-trends-iotdeveloper-survey-2018>, 2018.
- [14] (2021) Arm cortex-a53 mpcore processor technical reference manual , <https://developer.arm.com/documentation/ddi0500/j/>. Last Accessed: 2021-05-12.
- [15] (2021) Arm cortex-a72 mpcore processor technical reference manual , <https://developer.arm.com/documentation/100095/0003>. Last Accessed: 2021-05-12.
- [16] (2021) 7-zip home page, <https://www.7-zip.org/>. Last Accessed: 2021-07-25.
- [17] W. Li and Y. Yao, "Accelerate data compression in file system," in *2016 Data Compression Conference (DCC)*. IEEE Computer Society, 2016, pp. 615–615.
- [18] A. Gupta, A. Bansal, and V. Khanduja, "Modern lossless compression techniques: Review, comparison and analysis," in *2017 Second International Conference on Electrical, Computer and Communication Technologies (ICECCT)*. IEEE, 2017, pp. 1–8.
- [19] A. Moffat, "Huffman coding," *ACM Computing Surveys (CSUR)*, vol. 52, no. 4, pp. 1–35, 2019.
- [20] A. Ozsoy, M. Swamy, and A. Chauhan, "Pipelined parallel lzss for streaming data compression on gpgpus," in *2012 IEEE 18th International Conference on Parallel and Distributed Systems*. IEEE, 2012, pp. 37–44.
- [21] K. K. Yong, M. W. Chua, and W. K. Ho, "Cuda lossless data compression algorithms: a comparative study," in *2016 IEEE Conference on Open Systems (ICOS)*. IEEE, 2016, pp. 7–12.
- [22] M. Bark, S. Ubik, and P. Kubalik, "Lz4 compression algorithm on fpga," in *2015 IEEE International Conference on Electronics, Circuits, and Systems (ICECS)*. IEEE, 2015, pp. 179–182.
- [23] A. Gupta, A. Bansal, and V. Khanduja, "Modern lossless compression techniques: Review, comparison and analysis," in *2017 Second International Conference on Electrical, Computer and Communication Technologies (ICECCT)*. IEEE, 2017, pp. 1–8.
- [24] M. Pricopi, T. S. Muthukaruppan, V. Venkataramani, T. Mitra, and S. Vishin, "Power-performance modeling on asymmetric multi-cores," in *2013 International Conference on Compilers, Architecture and Synthesis for Embedded Systems (CASES)*. IEEE, 2013, pp. 1–10.
- [25] N. Mishra, C. Imes, J. D. Lafferty, and H. Hoffmann, "Caloree: Learning control for predictable latency and low energy," *ACM SIGPLAN Notices*, vol. 53, no. 2, pp. 184–198, 2018.
- [26] M. E. Haque, Y. He, S. Elnikety, T. D. Nguyen, R. Bianchini, and K. S. McKinley, "Exploiting heterogeneity for tail latency and energy efficiency," in *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*, 2017, pp. 625–638.
- [27] S. Srinivasan, N. Kurella, I. Koren, and S. Kundu, "Exploring heterogeneity within a core for improved power efficiency," *IEEE Transactions on Parallel and Distributed Systems*, vol. 27, no. 4, pp. 1057–1069, 2015.
- [28] V. Petrucci, O. Loques, and D. Mossé, "Lucky scheduling for energy-efficient heterogeneous multi-core systems," in *Proceedings of the 2012 USENIX conference on Power-Aware Computing and Systems*, 2012, pp. 7–7.
- [29] M. Pricopi and T. Mitra, "Task scheduling on adaptive multi-core," *IEEE transactions on Computers*, vol. 63, no. 10, pp. 2590–2603, 2013.
- [30] D. Lustig, C. Trippel, M. Pellauer, and M. Martonosi, "Armor: Defending against memory consistency model mismatches in heterogeneous architectures," in *Proceedings of the 42nd Annual International Symposium on Computer Architecture*, 2015, pp. 388–400.
- [31] T. Cao, S. M. Blackburn, T. Gao, and K. S. McKinley, "The yin and yang of power and performance for asymmetric hardware and managed software," in *2012 39th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2012, pp. 225–236.
- [32] M. Wang, S. Ding, T. Cao, Y. Liu, and F. Xu, "Asymo: scalable and efficient deep-learning inference on asymmetric mobile cpus," in *Proceedings of the 27th Annual International Conference on Mobile Computing and Networking*, 2021, pp. 215–228.
- [33] Y. Zhu and V. J. Reddi, "High-performance and energy-efficient mobile web browsing on big/little systems," in *2013 IEEE 19th International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, 2013, pp. 13–24.
- [34] V. Pankratius, A. Jannesari, and W. F. Tichy, "Parallelizing bzip2: A case study in multicore software engineering," *IEEE software*, vol. 26, no. 6, pp. 70–77, 2009.
- [35] M. Burtischer and P. Ratanaworabhan, "pfpc: A parallel compressor for floating-point data," in *2009 Data Compression Conference*. IEEE, 2009, pp. 43–52.
- [36] Y. Dua, V. Kumar, and R. S. Singh, "Parallel lossless hsi compression based on rls filter," *Journal of Parallel and Distributed Computing*, vol. 150, pp. 60–68, 2021.
- [37] R. A. Patel, Y. Zhang, J. Mak, A. Davidson, and J. D. Owens, *Parallel lossless data compression on the GPU*. IEEE, 2012.
- [38] W. Wolff and B. Porter, "Performance optimization on big.little architectures: A memory-latency aware approach," in *The 21st ACM SIGPLAN/SIGBED Conference on Languages, Compilers, and Tools for Embedded Systems*, ser. LCTES '20. New York, NY, USA:

- Association for Computing Machinery, 2020, p. 51–61. [Online]. Available: <https://doi.org/10.1145/3372799.3394370>
- [39] (2021) Energy aware scheduling, <https://www.kernel.org/doc/html/latest/scheduler/sched-energy.html>. Last Accessed: 2021-05-10.
- [40] H. Ribic and Y. D. Liu, “Energy-efficient work-stealing language runtimes,” *SIGARCH Comput. Archit. News*, vol. 42, no. 1, p. 513–528, Feb. 2014. [Online]. Available: <https://doi.org/10.1145/2654822.2541971>
- [41] T. Yu, R. Zhong, V. Janjic, P. Petoumenos, J. Zhai, H. Leather, and J. Thomson, “Collaborative heterogeneity-aware os scheduler for asymmetric multicore processors,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 5, pp. 1224–1237, 2020.
- [42] S. Balakrishnan, R. Rajwar, M. Upton, and K. Lai, “The impact of performance asymmetry in emerging multicore architectures,” in *32nd International Symposium on Computer Architecture (ISCA’05)*. IEEE, 2005, pp. 506–517.
- [43] K. Van Craeynest, A. Jaleel, L. Eeckhout, P. Narvaez, and J. Emer, “Scheduling heterogeneous multi-cores through performance impact estimation (pie),” in *2012 39th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2012, pp. 213–224.
- [44] J. Ziv and A. Lempel, “A universal algorithm for sequential data compression,” *IEEE Transactions on Information Theory*, vol. 23, no. 3, pp. 337–343, 1977.
- [45] X. Gu, P. S. Yu, and K. Nahrstedt, “Optimal component composition for scalable stream processing,” in *25th IEEE International Conference on Distributed Computing Systems (ICDCS’05)*. IEEE, 2005, pp. 773–782.
- [46] R. Kumar, K. I. Farkas, N. P. Jouppi, P. Ranganathan, and D. M. Tullsen, “Single-isa heterogeneous multi-core architectures: The potential for processor power reduction,” in *Proceedings. 36th Annual IEEE/ACM International Symposium on Microarchitecture, 2003. MICRO-36*. IEEE, 2003, pp. 81–92.
- [47] R. Kumar, D. M. Tullsen, P. Ranganathan, N. P. Jouppi, and K. I. Farkas, “Single-isa heterogeneous multi-core architectures for multithreaded workload performance,” in *Proceedings. 31st Annual International Symposium on Computer Architecture, 2004*. IEEE, 2004, pp. 64–75.
- [48] X.-P. Zhang and X.-M. Cheng, “Energy consumption, carbon emissions, and economic growth in china,” *Ecological economics*, vol. 68, no. 10, pp. 2706–2712, 2009.
- [49] U. Soytaş and R. Sari, “Energy consumption, economic growth, and carbon emissions: challenges faced by an eu candidate member,” *Ecological economics*, vol. 68, no. 6, pp. 1667–1675, 2009.
- [50] K. Fang, R. Heijungs, and G. R. de Snoo, “Theoretical exploration for the combination of the ecological, energy, carbon, and water footprints: Overview of a footprint family,” *Ecological Indicators*, vol. 36, pp. 508–518, 2014.
- [51] (2021) Rockchip wiki rk3399, http://opensource.rock-chips.com/wiki_RK3399. Last Accessed: 2021-05-10.
- [52] S. Z. Sheikh and M. A. Pasha, “Energy-efficient cache-aware scheduling on heterogeneous multicore systems,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 1, pp. 206–217, 2022.
- [53] A. Suyyagh and Z. Zilic, “Energy and task-aware partitioning on single-isa clustered heterogeneous processors,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, no. 2, pp. 306–317, 2019.
- [54] S. Williams, A. Waterman, and D. Patterson, “Roofline: an insightful visual performance model for multicore architectures,” *Communications of the ACM*, vol. 52, no. 4, pp. 65–76, 2009.
- [55] J. W. Choi, D. Bedard, R. Fowler, and R. Vuduc, “A roofline model of energy,” in *2013 IEEE 27th International Symposium on Parallel and Distributed Processing*. IEEE, 2013, pp. 661–672.
- [56] Y. J. Lo, S. Williams, B. Van Straalen, T. J. Ligocki, M. J. Cordery, N. J. Wright, M. W. Hall, and L. Oliker, “Roofline model toolkit: A practical tool for architectural and program analysis,” in *International Workshop on Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems*. Springer, 2014, pp. 129–148.
- [57] P. Elias, “Universal codeword sets and representations of the integers,” *IEEE Transactions on Information Theory*, vol. 21, no. 2, pp. 194–203, 1975.
- [58] (2022) Stream: Sustainable memory bandwidth in high performance computers, <https://www.cs.virginia.edu/stream/>. Last Accessed: 2022-06-29.
- [59] (2021) An introduction to amba axi, <https://developer.arm.com/documentation/102202/latest/>. Last Accessed: 2021-12-05.
- [60] (2021) lz4 source code, <https://github.com/lz4/lz4/>. Last Accessed: 2021-07-25.
- [61] S. Zhang, J. He, A. C. Zhou, and B. He, “Briskstream: Scaling data stream processing on shared-memory multicore architectures,” in *Proceedings of the 2019 International Conference on Management of Data*, 2019, pp. 705–722.
- [62] G. Pekhimenko, V. Seshadri, O. Mutlu, M. A. Kozuch, P. B. Gibbons, and T. C. Mowry, “Base-delta-immediate compression: Practical data compression for on-chip caches,” in *2012 21st International Conference on Parallel Architectures and Compilation Techniques (PACT)*, 2012, pp. 377–388.
- [63] S. Zhang, Y. Wu, F. Zhang, and B. He, “Towards concurrent stateful stream processing on multicore processors,” in *2020 IEEE 36th International Conference on Data Engineering (ICDE)*, 2020, pp. 1537–1548.
- [64] C. Iancu, S. Hofmeyr, F. Blagojević, and Y. Zheng, “Oversubscription on multicore processors,” in *2010 IEEE International Symposium on Parallel & Distributed Processing (IPDPS)*. IEEE, 2010, pp. 1–11.
- [65] A. Ilic, F. Pratas, and L. Sousa, “Cache-aware roofline model: Upgrading the loft,” *IEEE Computer Architecture Letters*, vol. 13, no. 1, pp. 21–24, 2013.
- [66] (2021) Perf wiki, https://perf.wiki.kernel.org/index.php/Main_Page. Last Accessed: 2021-11-07.
- [67] A. Magnani and S. P. Boyd, “Convex piecewise-linear fitting,” *Optimization and Engineering*, vol. 10, no. 1, pp. 1–17, 2009.
- [68] A. Toriello and J. P. Vielma, “Fitting piecewise linear continuous functions,” *European Journal of Operational Research*, vol. 219, no. 1, pp. 86–95, 2012.
- [69] R. Bellman, “The theory of dynamic programming,” *Bulletin of the American Mathematical Society*, vol. 60, no. 6, pp. 503–515, 1954.
- [70] S. Schneider, J. Wolf, K. Hildrum, R. Khandekar, and K.-L. Wu, “Dynamic load balancing for ordered data-parallel regions in distributed streaming systems,” in *Proceedings of the 17th International Middleware Conference*, ser. Middleware ’16. New York, NY, USA: Association for Computing Machinery, 2016. [Online]. Available: <https://doi.org/10.1145/2988336.2990475>
- [71] K. H. Ang, G. Chong, and Y. Li, “Pid control system analysis, design, and technology,” *IEEE transactions on control systems technology*, vol. 13, no. 4, pp. 559–576, 2005.
- [72] L. Fan, L. Xiong, and V. Sunderam, “Fast: differentially private real-time aggregate monitor with filtering and adaptive sampling,” in *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, 2013, pp. 1065–1068.
- [73] L. Shen, Z. Liu, Z. Zhang, and X. Shi, “Frame-level bit allocation based on incremental pid algorithm and frame complexity estimation,” *Journal of Visual Communication and Image Representation*, vol. 20, no. 1, pp. 28–34, 2009.
- [74] S. Tzafestas and N. P. Papanikolopoulos, “Incremental fuzzy expert pid control,” *IEEE Transactions on Industrial Electronics*, vol. 37, no. 5, pp. 365–371, 1990.
- [75] H.-B. Shin and J.-G. Park, “Anti-windup pid controller with integral state predictor for variable-speed motor drives,” *IEEE Transactions on Industrial Electronics*, vol. 59, no. 3, pp. 1509–1516, 2011.
- [76] V. Kalavri, J. Liagouris, M. Hoffmann, D. Dimitrova, M. Forshaw, and T. Roscoe, “Three steps is all you need: Fast, accurate, automatic scaling decisions for distributed streaming dataflows,” in *Proceedings of the 13th USENIX Conference on Operating Systems Design and Implementation*, ser. OSDI’18. USA: USENIX Association, 2018, p. 783–798.
- [77] (2021) Beach weather stations - automated sensors, <https://catalog.data.gov/dataset/beach-weather-stations-automated-sensors/resource/3b820f68-4dca-4ea7-8141-f37d9237734d>. Last Accessed: 2021-11-12.
- [78] (2019) Creator of the angry birds game, www.rovio.com. Last Accessed: 2021-05-10.
- [79] S. Zhang, Y. Mao, J. He, P. M. Grulich, S. Zeuch, B. He, R. T. Ma, and V. Markl, “Parallelizing intra-window join on multicores: An experimental study,” in *Proceedings of the 2021 International Conference on Management of Data*, 2021, pp. 2089–2101.
- [80] (2018) Shanghai stock exchange, <http://english.sse.com.cn/>. Last Accessed: 2021-11-12.
- [81] (2021) Ina226, <https://www.ti.com/product/INA226>. Last Accessed: 2021-11-12.

- [82] (2021) esp32s2, <https://www.espressif.com/en/products/socs/esp32-s2>. Last Accessed: 2021-11-12.
- [83] (2021) Rock pi 4 wiki, <https://wiki.radxa.com/Rockpi4>. Last Accessed: 2021-05-10.
- [84] (2021) buildroot home page, <https://buildroot.org/>. Last Accessed: 2021-05-10.
- [85] (2021) The linux kernel archives , <https://www.kernel.org/>. Last Accessed: 2021-05-12.
- [86] M. I. Solihin, L. F. Tack, and M. L. Kean, “Tuning of pid controller using particle swarm optimization (psa),” in *Proceeding of the international conference on advanced science, engineering and information technology*, vol. 1, 2011, pp. 458–461.
- [87] T. Somu Muthukaruppan, A. Pathania, and T. Mitra, “Price theory based power management for heterogeneous multi-cores,” ser. ASPLOS ’14. New York, NY, USA: Association for Computing Machinery, 2014, p. 161–176. [Online]. Available: <https://doi.org/10.1145/2541940.2541974>
- [88] J. Gilchrist, “Parallel data compression with bzip2,” in *Proceedings of the 16th IASTED international conference on parallel and distributed computing and systems*, vol. 16, no. 2004. Citeseer, 2004, pp. 559–564.
- [89] D. A. Huffman, “A method for the construction of minimum-redundancy codes,” *Proceedings of the IRE*, vol. 40, no. 9, pp. 1098–1101, 1952.
- [90] M. Milward, J. L. Nunez, and D. Mulvaney, “Design and implementation of a lossless parallel high-speed data compression system,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 15, no. 6, pp. 481–490, 2004.
- [91] K. Sano, K. Katahira, and S. Yamamoto, “Segment-parallel predictor for fpga-based hardware compressor and decompressor of floating-point data streams to enhance memory i/o bandwidth,” in *2010 Data Compression Conference*. IEEE, 2010, pp. 416–425.
- [92] J. Tian, S. Di, C. Zhang, X. Liang, S. Jin, D. Cheng, D. Tao, and F. Cappello, “Wavesz: A hardware-algorithm co-design of efficient lossy compression for scientific data,” in *Proceedings of the 25th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, 2020, pp. 74–88.
- [93] A. Ozsoy and M. Swamy, “Culzss: Lzss lossless data compression on cuda,” in *2011 IEEE International Conference on Cluster Computing*. IEEE, 2011, pp. 403–411.
- [94] A. Yang, H. Mukka, F. Hesaaraki, and M. Burtscher, “Mpc: a massively parallel compression algorithm for scientific data,” in *2015 IEEE International Conference on Cluster Computing*. IEEE, 2015, pp. 381–389.
- [95] Y. Huang, Y. Li, Z. Zhang, and R. W. Liu, “Gpu-accelerated compression and visualization of large-scale vessel trajectories in maritime iot industries,” *IEEE Internet of Things Journal*, vol. 7, no. 11, pp. 10 794–10 812, 2020.
- [96] J. Shun and F. Zhao, “Practical parallel lempel-ziv factorization,” in *2013 Data Compression Conference*. IEEE, 2013, pp. 123–132.
- [97] F. Knorr, P. Thoman, and T. Fahringer, “ndzip: A high-throughput parallel lossless compressor for scientific data,” in *2021 Data Compression Conference (DCC)*. IEEE, 2021, pp. 103–112.
- [98] B. Salami, H. Noori, and M. Naghibzadeh, “Fairness-aware energy efficient scheduling on heterogeneous multi-core processors,” *IEEE Transactions on Computers*, vol. 70, no. 1, pp. 72–82, 2020.
- [99] A. Mascitti, T. Cucinotta, and M. Marinoni, “An adaptive, utilization-based approach to schedule real-time tasks for arm big. little architectures,” *ACM SIGBED Review*, vol. 17, no. 1, pp. 18–23, 2020.
- [100] X. Li and G. Li, “An adaptive cpu-gpu governing framework for mobile games on big. little architectures,” *IEEE Transactions on Computers*, 2020.