



Simplifying and implementing service level objectives for stream parallelism

Dalvan Griebler^{1,3} · Adriano Vogel¹ · Daniele De Sensi² · Marco Danelutto² · Luiz G. Fernandes¹

Published online: 5 June 2019

© Springer Science+Business Media, LLC, part of Springer Nature 2019

Abstract

An increasing attention has been given to provide service level objectives (SLOs) in stream processing applications due to the performance and energy requirements, and because of the need to impose limits in terms of resource usage while improving the system utilization. Since the current and next-generation computing systems are intrinsically offering parallel architectures, the software has to naturally exploit the architecture parallelism. Implement and meet SLOs on existing applications is not a trivial task for application programmers, since the software development process, besides the parallelism exploitation, requires the implementation of autonomic algorithms or strategies. This is a system-oriented programming approach and requires the management of multiple knobs and sensors (e.g., the number of threads to use, the clock frequency of the cores, etc.) so that the system can self-adapt at runtime. In this work, we introduce a new and simpler way to define SLO in the application's source code, by abstracting from the programmer all the details relative to self-adaptive system implementation. The application programmer specifies which parts of the code to parallelize and the related SLOs that should be enforced. To reach this goal, source-to-source code transformation rules are implemented in our compiler, which automatically generates self-adaptive strategies to enforce, at runtime, the user-expressed objectives. The experiments highlighted promising results with simpler, effective, and efficient SLO implementations for real-world applications.

Keywords Parallel programming · Stream processing · Self-adaptive · Domain-specific language · Power-aware computing

✉ Dalvan Griebler
dalvan.griebler@acad.pucrs.br; dalvangriebler@gmail.com

¹ School of Technology, Pontifical Catholic University of Rio Grande do Sul (PUCRS), Porto Alegre, Brazil

² Department of Computer Science, University of Pisa (UNIP), Pisa, Italy

³ Laboratory of Advanced Research on Cloud Computing (LARCC), Três de Maio Faculty (SETREM), Três de Maio, Brazil

1 Introduction

Service-oriented approaches influenced new system models like cloud computing one [7]. Since a service can be represented by software components, functions, or a sequence of commands, this can help to improve the expressiveness and methodology of parallel software design. The service behavior can also be evaluated from different perspectives through the concept of quality of service (QoS), which identifies non-functional attributes. Performance metrics like throughput are used to measure QoS or to establish requirements between providers and clients. The programmer may know about the requirements of the components and how to improve the QoS of them. Consequently, the programmer can define service level objectives (SLOs) for each one of these components so that they behave as expected in a service level agreement (SLA) or in a contract [33].

Since the new- and next-generation computing systems are intrinsically offering parallel architectures, the software has to naturally exploit the architecture parallelism. Implement and meet SLOs on existing applications is not a trivial task for application programmers, because the usual software development process, besides the parallelism exploitation, also requires the implementation of autonomic algorithms or strategies. This is a system-oriented programming approach and requires the management of multiple knobs and sensors (e.g., the number of threads to use, the clock frequency of the cores, etc.) simultaneously so that the system self-adapts at runtime.

In stream processing applications, parallelism is typically exploited by using linear or nonlinear pipeline pattern compositions [27]. To this purpose, parallel programming framework such as StreamIt [35], Intel TBB [30], and FAST-FLOW [1, 2] provides different programming approaches and interfaces with a reasonable performance scalability for this domain. Although these frameworks are equipped with high-level pattern implementations to express the parallelism, they are still closer to expert system programmers rather than to the application domain programmers. Seeking to provide domain-specific and suitable abstractions for stream parallelism, SPAR [17] was created. It improves application programmers' productivity through a C++11 annotation-based language, which does not require to rewrite/restructure the sequential source code [17]. Moreover, it is important to highlight that stream processing applications are usually characterized by unpredictable load fluctuations and uncertain end of execution (may never end) [4]. However, none of these parallel programming alternatives automatically deals with this service-oriented behavior, since they are not able to guarantee SLOs due to the static resource assignment (e.g., a fixed amount of threads).

Besides the need for improving performance through the efficient exploitation of the multi-core parallelism, there are also other major concerns such as power-aware computing and efficient resource usage [16, 25]. To address these needs, the NORNIR framework was created, providing runtime support to dynamically and automatically control the resources allocated for the application according to the user needs [12]. However, NORNIR, like most existing self-adaptive solutions, only works on parallel applications and requires sequential code refactoring.

To simplify the specification of SLOs in sequential stream processing applications, differently from NORNIR and state-of-the-art parallel programming frameworks, we proposed a set of SLO attributes which can be inserted along with SPAR's stream parallelism annotations in the sequential code. In addition to that, we introduce a programming methodology where the programmer specifies which source code regions can be parallelized and the requirements that should be enforced. We implemented source-to-source code transformation rules in the SPAR compiler to automatically generate the self-adaptive strategies that enforce the user-expressed objectives at runtime. The proposed energy-aware SLO attributes were implemented using NORNIR's runtime support and studied in the previous work [18]. Our approach could also be applied to other frameworks, for example to the framework designed by the REPARA project,¹ which provides a set of C++11 attributes to introduce generic parallelism [11]. Moreover, we implement some SLOs on NORNIR and others on top of FASTFLOW, to demonstrate that this can be applied to different runtimes, with a different implementation complexity according to the used abstraction. The major contributions of this paper are summarized as follows:

- A new set of SLO attributes semantically defined by using the standard C++11 and SPAR annotations.
- Design and implementation of new self-adaptive strategies using the FASTFLOW framework.
- The implementation of the new SLOs in the SPAR compiler with source-to-source transformation rules, targeting self-adaptive strategies with NORNIR and FASTFLOW back-ends.
- An experimental evaluation with real-world applications, comparing our implementations with some state-of-the-art solutions.

We structured our paper as follows. We first present the related work in Sect. 2. Sect. 3 describes SPAR. Section 4 details the proposed SLOs for stream parallelism and its implementation. In Sect. 5, a set of experiments are analyzed and discussed. Finally, Sect. 6 makes the conclusions of the paper.

2 Related work

In the literature, there are different studies targeting power consumption, throughput, and system utilization objectives. Among them, the approach of Maggio et al. [25] monitors generic applications and supports the specification of a target performance (throughput) in the parallel code. It efficiently manages the CPU cores, adapting the amount of resource usage needed. However, it supposes that the parallel application has already been implemented and does not provide any mechanism to introduce SLO in sequential programs.

¹ <http://repara-project.eu/>.

Some existing algorithms do not explicitly model the power consumption of applications, thus only providing the possibility to specify performance SLOs [14, 24]. In some cases, it is not even possible to enforce a specific performance requirement, but only to run the application in the most efficient [32] or the most performing configuration [10, 29, 34]. Other works provide SLO on power consumption and/or application performance by acting on mechanisms different from those considered in this work, such as caches [37] or network interfaces [36]. Another alternative approach to the presented problem is to dynamically change the accuracy of the results computed by the application according to the user SLOs. This technique is known as *approximate computing* [5, 22, 38] and can be used to trade an increase on performance (or a decrease on power consumption) for a decrease on the quality of the results computed by the application. However, these approaches are usually focused on algorithms and techniques to provide SLOs rather than on programming abstractions to express SLOs, as we do in this work.

Concerning stream-parallel processing for real-time data analytic, Floratou et al. [16] introduced the notion of self-regulation in Twitter's Heron framework, called Dhalion. The user defines a target throughput as an SLO parameter for Dhalion. The self-regulator engine handles the number of processes and number of instances in a cloud infrastructure to provide the specified throughput. In the experiments, the results revealed that the system can dynamically adapt resources and automatically reconfigure to meet SLOs. We differently proposed six target SLOs to be expressed in sequential source codes.

Some works focus on high-level abstractions for energy saving on data parallelism [3, 31], by providing compiler directives for expressing energy consumption and performance objectives in OpenMP. While Shafik et al. [31] can minimize energy consumption on both sequential and parallel applications, they do not provide any means to explicitly control the performance of the application. On the other hand, in Alessi et al. [3], OpenMPE is proposed adding a new construct and two clauses (objectives) for OpenMP. Their solution was implemented using a source-to-source compiler, which recognizes the new directives and controls the number of threads used by OpenMP and applies DVFS to satisfy the SLOs expressed by the user. This is probably the closest work to the approach we are proposing in this work. The main difference is that, while Alessi et al. [3] target batch applications (i.e., applications for which all the input data is already available in memory) implemented through OpenMP, we provide support for stream processing applications, exposing ad hoc SLOs for these applications such as system utilization.

Eventually, many existing solutions are either simulated or validated on post-mortem data (i.e., they are executed after the application finished its execution, in a *what-if* analysis fashion) [15, 24, 28, 34]. We believe that, although a simulation may provide a first approximation about the precision that the algorithm could have in enforcing the required SLOs, it would not take into account the runtime overhead of these methods. Differently from these works, in this paper, we describe and implement a solution which has been validated by controlling real applications throughout their execution. Finally, stream processing applications vary during the execution without a pre-defined end, which makes such approaches unfeasible to apply in our application domain.

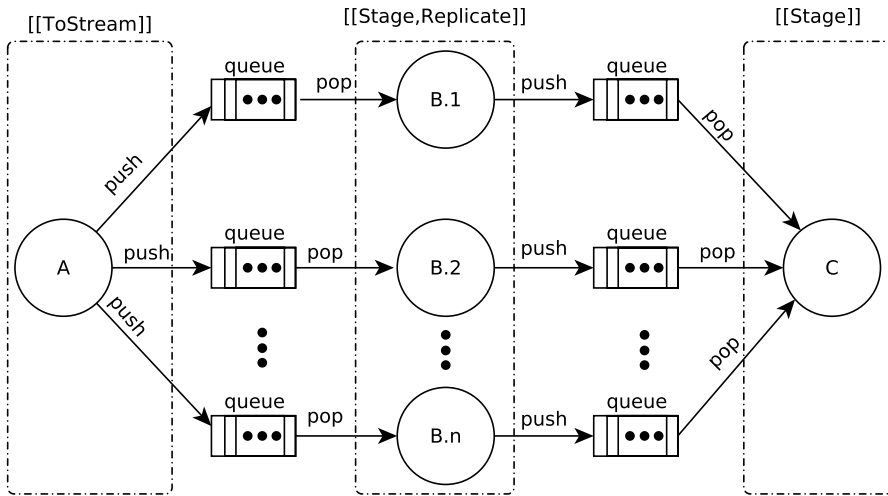


Fig. 1 SPAR runtime: activity graph and communication queues

3 SPAR: high-level stream parallelism

SPAR² is a domain-specific language (DSL) designed to support high-level stream parallelism for application programmers [17]. With SPAR, instead of rewriting the source code, the programmer introduces C++ annotations (standard C++-11 [26]) using five attributes, representing the main properties of stream processing applications. The `ToStream` attribute identifies the beginning of a stream region, which can be viewed as an assembly line. The `Stage` attribute marks a workstation in this assembly line, which can be composed by as many as necessary. Auxiliary attributes can be used inside the attribute list of an annotation sentence. The `Input` and `Output` attributes are to specify the input and output stream items, respectively, while the `Replicate` attribute is for replicating stateless stages to increase the degree of parallelism.

Listing 1 provides a short code example annotated with SPAR attributes. This example represents a typical use case of stream parallelism, where there is a sequence of operations to be performed on each stream element. The parallel activity graph produced by the SPAR compiler for Listing 1 is shown in Fig. 1. SPAR generates the parallel code with the `FASTFLOW` library [1], which implements different parallel patterns [27] for stream processing computations. The SPAR compiler parses the code of Listing 1 and represents the code with an abstract syntax tree (AST) [17]. Traversing the AST, it performs a semantic analysis of the attributes to further make the source-to-source transformations. In this step, the SPAR compiler finds the best parallel pattern that meets the parsed annotation schema. In the case of Listing 1, it will generate parallel code with three stages, where the middle one is

² SPAR website: <https://gmap.pucrs.br/spar>.

replicated. Moreover, different compositions with sequential or replicated stages can be achieved. By default, elements are scheduled from the `ToStream` stage to the `Stage` in a round-robin way. However, it is possible to use an on-demand policy by specifying the `-spar_ondemand` flag to the SPAR compiler. If the data needs to be received from the last stage in the same order it was produced by the `ToStream` stage, the programmer can specify the `-spar_ordered` flag to the SPAR compiler.

```

1| [[ spar :: ToStream ]] while(1){
2|   frame f = read_frame();
3|   if(f.empty()) break;
4|   [[ spar :: Stage, spar :: Input(f), spar :: Output(f), spar :: Replicate(n) ]]
5|   for (int i=0; i<f.length(); i++) {
6|     f[i] = convert(f[i]);
7|   }
8|   [[ spar :: Stage, spar :: Input(f) ]]{
9|     write_frame(f);
10|  }
11| }

```

Listing 1 SPAR example: image processing representation with stream parallelism.

Note that the `Replicate` attribute applies the replication role over the `Stage` in Fig. 1. Each replicated stage has its own input and output lock-free queues. The first stage executes the code inside the `ToStream` region, which generates stream items for the subsequent stages. In the default configuration of SPAR runtime, the stages actively try to push or pop stream items from the queues. If the queue is full or empty, the stage thread executes an active loop, trying to push or pop until it eventually succeeds. Every time that a given stage fails in performing push or pop, the stage generates a push or pop lost event. This may generate an extra overhead for coarse-grain computations. Therefore, users may set the SPAR runtime to behave in a blocking mode through the `spar_blocking` compiler flag. In this case, the stage thread will not stay in a loop, it will wait until it can perform push or pop in the shared queue.

4 Service level objective for stream parallelism

Service level objectives (SLOs) are traditionally included in service level agreements (SLAs), which are contracts to manage the quality of service (QoS) established between customers and providers [33]. An SLA contract defines the level of service which is acceptable by the user and attainable by the provider. The SLO is a target value or a range of values for a certain level of service to be delivered. The level of service is measured by a service level indicator (SLI). A typical structure of SLO can be written $SLI \leq target$ or $lower_bound \leq SLI \leq upper_bound$ [6]. When an SLO is violated, the system should react to guarantee the quality of service and SLA. Our design goal is to simplify the usability of SLO in stream processing applications.

Figure 2 depicts our proposed methodology to express SLOs in the application source code. The first step in the developing process is to code the stream

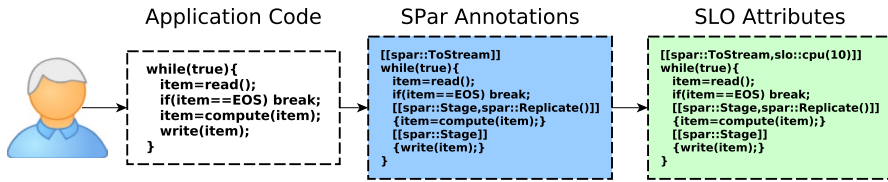


Fig. 2 Our methodology to define SLOs for stream parallelism

processing application (not needed for legacy applications). After, the programmer inserts the SPAR annotations to express the stream parallelism. This can be done following the recommendations of SPAR's annotation methodology [17]. Lastly, the programmer can insert SLO attributes along with SPAR's annotations in the source code. Therefore, the only requirement is to choose the SLO metric and its initial target value. No extra details must be provided by the application programmers, which can spend most of their time in coding the sequential application. Consequently, in this work, we support the application programmers with an opportunity to express stream parallelism with SPAR and define a target QoS through SLO attributes.

The SLO attributes are proposed to be used along with a `ToStream` annotation, which identifies the beginning of a stream parallelism region. Therefore, the SLO is applied to this particular region. Listing 2 demonstrates the definition of a power consumption SLO of 60 W in line. It is worth noting that besides the `slo::Power` attribute, no other modification is required with respect to the original SPAR code (Listing 1). While multiple SLOs attributes could be used together, there are only a few meaningful combinations. Usually, the user may need to express one SLO on performance and one SLO on power consumption. Other combinations are possible (e.g., `slo::Throughput` and `slo::Utilization` at the same time), but since they are two different representations of performance, they could conflict between each other. Additionally, adding too many constraints could lead to situations where there would be no feasible solutions, and it may be complex for the application programmer to find the right SLO values. Table 1 describes the SLO attributes proposed in this work. The attributes belong to the `slo` namespace and accept one argument, which is a value defining the target SLO.

```

1| [[ spar::ToStream, slo::Power(60) ]] while(1){
2|   frame f = read_frame();
3|   if(f.empty()) break;
4|   [[ spar::Stage, spar::Input(f), spar::Output(f), spar::Replicate(n) ]]
5|   for (int i=0; i<f.length(); i++) {
6|     f[i] = convert(f[i]);
7|   }
8|   [[ spar::Stage, spar::Input(f) ]] {
9|     write_frame(f);
10|   }
11| }

```

Listing 2 SPAR code example with power consumption SLO.

Table 1 SLO attributes for SPAR

Name	Argument	Description
slo::Throughput	(min-items)	The user can specify the minimum throughput required in items per second. The respective environment variable is SLO_THROUGHPUT
slo::Power	(max-watts)	The user can specify the maximum power consumption in Watts. The respective environment variable is SLO_POWER
slo::Utilization	(min-%)	The user can specify the minimum runtime system utilization required in percentage (from 1 to 100). In our case, it represents the percentage of time that the system is active (i.e., actively processing input elements) over a time interval. The respective environment variable is SLO_UTILIZATION
slo::Latency	(max-time)	The user can specify the maximum latency in milliseconds. This latency refers to the time taken for an item passing from a stage to another one. The respective environment variable is SLO_LATENCY
slo::CPU	(max-%)	The user can specify the maximum CPU utilization in percentage (from 1 to 100%). The respective environment variable is SLO_CPU

4.1 SLO implementations

SLO attributes in SPAR were initially proposed in our previous work [18]. In this paper, we extend that work by providing new self-adaptive strategies and SLO attributes (`slo::CPU`, `slo::Latency`). Moreover, we add a new algorithm for enforcing the `slo::Throughput` SLO when it is not combined with power consumption SLO. The compiler will decide which strategy to generate when performing the source-to-source code transformations.

We first explain the self-adaptive strategies to meet the so-called energy-aware SLOs, which rely on the NORNIR runtime support [12]. NORNIR monitors the application throughout its entire execution, dynamically changing the number of resources used by the application to satisfy the requirements expressed by the user. For example, NORNIR may decide to reduce the number of replicated stages of the application to decrease its power consumption, or to increase the clock frequency of the cores to increase the application throughput.

Moreover, NORNIR can rely on different algorithms to decide how many resources to add/remove, either based on machine learning techniques [13] or on heuristics. When machine learning techniques are used, when the application starts, NORNIR performs a lightweight training phase by testing different configurations and collecting application data/performance indicators. The results collected are used to build prediction models which are used to find the optimal configuration according to the objectives specified by the user. If no feasible solution is found, NORNIR selects the resources configuration characterized by performance and power consumption as close as possible to the user requirements.

Besides providing the possibility to control existing parallel applications (by inserting instrumentation calls in the existing code), NORNIR can also be used as a programming framework (by relying on the FASTFLOW framework) for implementing stream-parallel applications with embedded self-adaptation support. We exploited this second possibility so that SPAR can translate sequentially annotated code into self-adaptive NORNIR parallel code. All details relative to the use of NORNIR are abstracted and made simple along with the stream parallelism.

In addition to the NORNIR's strategies, we also provide new self-adaptive strategies for `slo::CPU`, `slo::Latency`, and `slo::Throughput`, which rely on the SPAR runtime system (which is built on top of FASTFLOW), as illustrated in Fig. 3. Observe that this activity graph is a way of simplification from the SPAR runtime system presented in Fig. 1. The strategy follows the MAPE approach [23] with a feedback closed-loop [21]. The Monitor entity periodically collects data from the sensors, which can be originated from the application or from the operating system. These data are used by the Analyze phase, which interprets the data and extracts relevant statistics.

Afterward, the Plan phase decides if the SPAR runtime system must be adapted to meet the specified SLO. It may be impossible for a strategy to achieve a given SLO. In such a case, the strategy attempts to reach the SLO as close as possible. Since stream processing applications may have load fluctuations, unnecessary adaptations should be avoided. In our strategies, we used a threshold, which is a percentage number that can be tolerated when the actually monitored metric is higher but

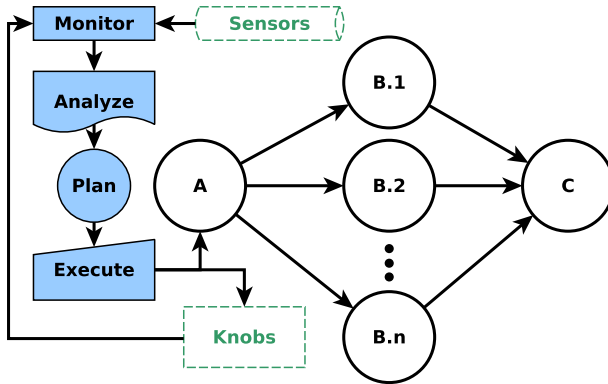


Fig. 3 SLO implementation by using self-adaptive strategies

close to the target SLO. The default threshold is 20%, such value was ascertained in [39] as a suitable one for stream processing applications. Moreover, for the sake of flexibility, users may customize the threshold value using the `SLO_THRESHOLD` environment variable. The Execute entity receives the planned action and applies the adaptation by sending instructions to the SPAR runtime system (e.g., tasks/items distribution) and system knobs (e.g., adapt the number of active replicas). Although we used this same idea for implementing all new SLOs, each SLO has its specific self-adaptive strategy (i.e., Analyze and Plan phases), described in the following section.

4.1.1 `slo::Latency`

In a previous work [39], we have shown the possibility to manage the latency by adapting the number of replicas. In this work, we extend the previous study by implementing it in the SPAR compiler as well as providing an SLO option. During our study, we have seen how the number of replicas affects the latency of stream items. Additionally, it is a presumably difficult task for programmers to manually adapt their software at runtime based on the latency SLO constraints and on the actual application latency. As a consequence, we aim at abstracting from programmers the impact of the number of replicas in latency.

We implemented a strategy for the SPAR's runtime that monitors and manages the latency of stream items by autonomously adjusting the number of replicas. Considering the representation in Fig. 3, the stage A adds a timestamp to the stream items and the Monitor entity collects from a sensor that is inside the stage C, where the latency of the stream items is measured. In the Analyze and Plan phases, the latency information from the Monitor entity is compared to the SLO, and the Plan phase decides whether to change or not the number of replicas, based on the tolerated threshold. Eventually, Execute entity sends instructions to control knobs which changes the number of active replicas and stage A which distributes the items among the active replicas from the stage B.

4.1.2 `slo::Throughput`

Our self-adaptive strategy for the `slo::Throughput` SLO is based on the number of items processed per second. The self-adaptive strategy is able to increase or decrease the throughput with the number of replicas control knob. In this SLO implementation, accordingly to Fig. 3, the Monitor entity periodically gets the number of items per second from the sensor that we installed in stage C. For each iteration, the throughput is the result of dividing the number of processed items by the time that it was taken. Throughput rates are then stored and accessed by the Analyze phase, which provides useful data statistics to the Plan phase decide if an adaptation is required.

The Plan phase also has a maximum value for the number of replicas, which is defined according to the machine's processing capabilities, gathered by another sensor that extracts hardware information. The Execution entity is updated by the Plan phase in order to send information to the control knob, which increases or decreases the number of replicas depending on the need. In addition, the Execution entity will inform the stage A (Fig. 3), which implements the task scheduler in the SPAR runtime system.

4.1.3 `slo::CPU`

In the stream processing domain, several SLOs can be relevant for defining performance/efficiency objectives. This occurs because in stream processing applications, differently from other application domains, the maximum amount of resources available are not always used nor needed. Continuously using the total resources capacity tends to reduce the efficiency of the system. Also, using the maximum resources does not actually mean that a stream processing application will achieve the best performance [13]. Performance is complex in the stream processing domain because the workload trend varies in a timely fashion according to variable input rates, volumes, resources availability, and different performance objectives. Consequently, we are employing efforts to enforce performance goals for enabling a customizable execution of stream processing applications and their unique characteristics. It is also relevant to allow programmers to define objectives regarding the consumption of resources.

Therefore, we provide an option to define the CPU utilization (`slo::CPU`) SLO when running a given application. Although there are available OS-level tools for controlling the CPU usage (e.g., `CPULimit` [9]), such tools are arguably not flexible. Considering the dynamic nature of stream processing applications, we expect to adapt the degree of parallelism of the application at runtime for optimizing the CPU utilization and meet the target SLO. The implemented self-adaptive strategy follows the schema sketched in Fig. 3.

Differently from the previous SLOs, the Monitor entity is periodically getting the current CPU utilization from the sensor, which is reading it from the operating system. The Analyze phase calls the Monitor entity for providing CPU utilization statistics to the Plan phase, which aims to decide whether the number of replicas should be increased or decreased. To avoid oscillation and instability regarding the

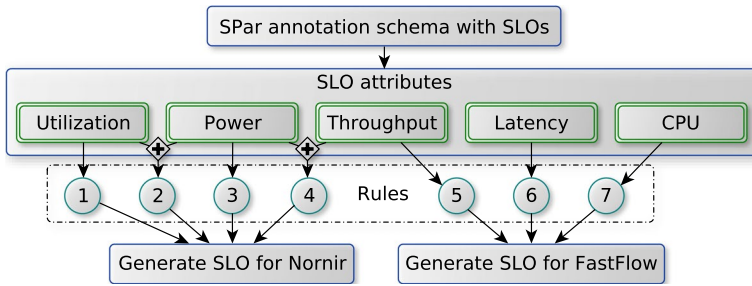


Fig. 4 Source-to-source transformation for the SLO attributes with SPAR

replicas reconfiguration, a threshold (described in Sect. 4.1) value is used so that the number of replicas is not increased when the utilization is close to the SLO. Finally, the Execution entity simply sends this information to the system knob apply an action (increase, decrease, or stay as it is) as well as to the task scheduler in the stage *A* to manage stream item in a correct manner.

4.2 Source-to-source transformations

Self-adaptive strategies for each SLO attribute are automatically generated during the program compilation. We used the SPAR compiler to implement the source-to-source code transformations. This required to add a new compilation step inside the compiler when performing the transformations from the SPAR annotations to parallel patterns. In this step, the compiler builds the SPAR runtime system with the communications, scheduling, and synchronizations. Based on the semantics previously specified, we added semantics-checking for the SLO attributes to ensure correct code generation. All transformations and analysis are performed in the AST, and the parallel code generation is based on transformation rules [17]. In addition to that, we also implemented transformation rules to generate the appropriate self-adaptive strategy for each SLO attribute annotated in the source code.

Figure 4 depicts a high-level representation of the source-to-source transformations targeting the proposed SLOs. This occurs after the semantic analysis where the annotation abstract syntax tree (AAST) is built from the source codes. The AAST also contains an internal representation of the SLO attributes. The compiler checks if SLO attributes were specified in the specific `spar::ToStream` node. In this case, the appropriate rules are applied as shown in Fig. 4.

As presented in Fig. 4, we implemented seven transformation rules to implement the proposed SLO attribute declarations. The first four rules will generate parallel code with SLO strategies to be executed with the NORNIR framework's back-end. On the other hand, the last three rules will generate parallel code with the new SLO strategies (proposed in this paper) for the FASTFLOW framework. When using FastFlow's back-end, the task scheduler thread also hosts the self-adaptive strategies. On the other hand, in NORNIR, this is managed by an extra thread. Besides, NORNIR and FASTFLOW have different programming interfaces, the

parallel patterns are conceptually similar. Modifications were only necessary to accommodate the proper routine names. Considering that parallel patterns were implemented in the previous source-to-source transformation step based on the rules already designed in [17], here we concentrate on the transformation rules related to the SLOs. It is important to note that here we describe only the meaning of the transformations required; the implementation details are arguably not relevant for presenting a simplified description. The transformations performed are the following:

1. Implement the following transformation steps for the `slo::Utilization` attribute: (a) insert the routine which implements the SLO utilization strategy in the `NORNIR` library before the declaration of `spar::ToStream`; and (b) give as a parameter the attribute argument to be the target SLO for the strategy routine.
2. Implement the following transformation steps for the `slo::Utilization` and `slo::Power` attributes: (a) insert the routine which implements the SLOs utilization and power strategies in the `NORNIR` library before the declaration of `spar::ToStream`; and (b) give as parameters the attribute arguments to be the target SLO for the strategy.
3. Implement the following transformation steps for the `slo::Power` attribute: (a) insert the routine which implements the SLO power strategy in the `NORNIR` library before the declaration of `spar::ToStream`; and (b) give as a parameter the attribute argument to be the target SLO for the strategy routine.
4. Implement the following transformation steps for the `slo::Throughput` and `slo::power` attributes: (a) insert the routine which implements the SLOs throughput and power in the `NORNIR` library before the declaration of `spar::ToStream`; and (b) give as parameters the attribute arguments to be the target SLOs for the strategy routine.
5. Implement the following transformation steps for the `slo::Throughput` attribute: (a) insert the routine which implements the SLO throughput strategy in the `FASTFLOW` library before the declaration of `spar::ToStream`; and (b) give as a parameter the attribute argument to be the target SLO for the strategy routine.
6. Implement the following transformation steps for the `slo::Latency` attribute: (a) insert the routine which implements the SLO latency strategy in the `FASTFLOW` library before the declaration of `spar::ToStream`; and (b) give as a parameter the attribute argument to be the target SLO for the strategy routine.
7. Implement the following transformation steps for the `slo::CPU` attribute: (a) insert the routine which implements the SLO CPU utilization strategy in the `FASTFLOW` library before the declaration of `spar::ToStream`; and (b) give as a parameter the attribute argument to be the target SLO for the strategy routine.

After the SLO transformation rules were applied, the parallel pattern generated in the AST is built with these rules' configurations, either for the `NORNIR`'s self-adaptive manager or for the implemented `SPAR`'s manager. The `SPAR`'s manager runs a MAPE feedback closed-loop in the generated `SPAR`'s task scheduler,

which is on top of the FastFlow library. We also support the `-spar_blocking` and `-spar_ordered` compilation flags that are natively supported in SPAR (see Sect. 3). These compilation flags were used for the experiments in the next section.

5 Experiments

In this section, we first introduce the considered real-world applications. Then, we will compare the code generated by SPAR with handwritten parallel implementations for these applications, both regarding maximum performance achieved and productivity. Also, we analyze the self-adaptation capabilities automatically generated by the SPAR compiler under different scenarios. The experiments have been executed in the following two machines:

- **M1** is a machine equipped with 32 GB of RAM memory and two Intel(R) Xeon(R) CPU E5-2620 v3 2.40 GHz processors (12 cores-24 hardware threads). The operating system used was Ubuntu Server 64 bits with the kernel 4.4.0-59-generic. The GCC version used was the 5.4.0 using the compiler `-O3` flag.
- **M2** is a dual-socket NUMA machine with two Intel Xeon E5-2695 Ivy Bridge CPUs running at 2.40 GHz featuring 24 cores (12 per socket). The machine exposes 13 frequency levels, ranging from 1.2 to 2.4 GHz, at steps of 0.1 GHz. Each core has 2-way hyper-threading, 32 KB private L1, 256 KB private L2 and 30 MB of L3 shared with the cores on the same socket. The machine has 64 GB of DDR3 RAM. We used Linux 3.14.49 x86_64 shipped with CentOS 7.1 and gcc version 4.8.5. For all our experiments, we disabled the hyper-threading feature.

5.1 Applications

We briefly describe the real-world application set, input loads, and parallel implementations. For a detailed description of how *Lane Detection* and *Person Recognition* have been parallelized by using SPAR, please refer to [19], while for *Pbzip2* more details can be found in [20].

Lane Detection is a video processing application to detect road lanes, implemented by using the OpenCV library. To introduce parallelism in the sequential code, it is annotated with SPAR by identifying three stages: (1) a first stage which reads the frames; (2) another stage, replicated a number of times, which processes the frames in parallel; (3) the last stage which displays the frames in the proper order, with the lanes properly marked. As input workload, we used a 22 MB MPEG-4 video (640 × 360 pixels).

Person Recognition is an application used to recognize people in a video. The parallel structure of this application is similar to *Lane Detection*, with the middle stage detecting the faces from the crowd and searching in an image database to

Table 2 Performance improvement with respect to a handwritten implementation using FASTFLOW for the energy-aware SLOs

Power, throughput, utilization	Pbzip2	Lane detection	Person recognition
Performance improvement (%)	+ 0.48%	− 1.45%	− 0.92%
LOC reduction (%)	− 15.86%	− 21.51%	− 24.49%

Negative percentages are the overheads (means slower) added by the SPAR and SLO abstractions. LOC Reduction, negative values mean that SPAR with SLO attributes is more concise than the handwritten one

classify each face detected. As input workload, we used a 4.8 MB MPEG-4 video (640 × 360 pixels) along with a training set of 10 face images of 150 × 150 pixels.

Pbzip application is a parallel implementation of the `bzip2` block-sorting files compressor.³ This is a very coarse-grained application characterized by a stream-parallel programming model. The SPAR version is annotated with three stages, where the middle stage is replicated. The input file to compress that we used for our experiments is a 6,3GB file containing a dump of all the abstract present on the English Wikipedia on 01/12/2015.

5.2 Comparison with handwritten implementations

Before evaluating the ability to satisfy SLO specified by the user, we want to show that from a performance standpoint, the code generated by SPAR is comparable with a handwritten implementation. On the other hand, we would like to show that our solution reduces the code intrusion required to transform a sequential application into a parallel one. As reference implementations for *Pbzip*, we consider the original Pthreads version, while for *Lane Detection* and *Person Recognition* applications, we consider the handwritten FASTFLOW versions described in [19].

Performance To measure the maximum performance presented in Table 2, we executed both the reference and our solution generated versions by running them with 24 threads (to have at most one thread per core). The reported results refer to machine M2. For our generated version, we did not specify any SLO, but we still monitor the application by using NORNIR. By doing so, we monitored both the overhead introduced by the interaction with the self-adaptive support and possible inefficiencies in the generated code. As shown by the results in Table 2 relative to the energy-aware SLOs, for *Lane Detection* and *Person Recognition*, the overhead is negligible (below 1.5%). For *Pbzip2*, there is a slight improvement caused by the use of FASTFLOW and its optimizations as runtime support in NORNIR, while the reference implementation was based on Pthreads.

Code Intrusion To measure the code intrusion, we rely on lines of code (LOC) metric in Table 2 for energy-aware SLOs. Despite that LOC is not universally accepted, it is commonly used to compare different implementations of the same application [40]. For our measurements, we only considered the source files

³ <http://compression.ca/pbzip2/>.

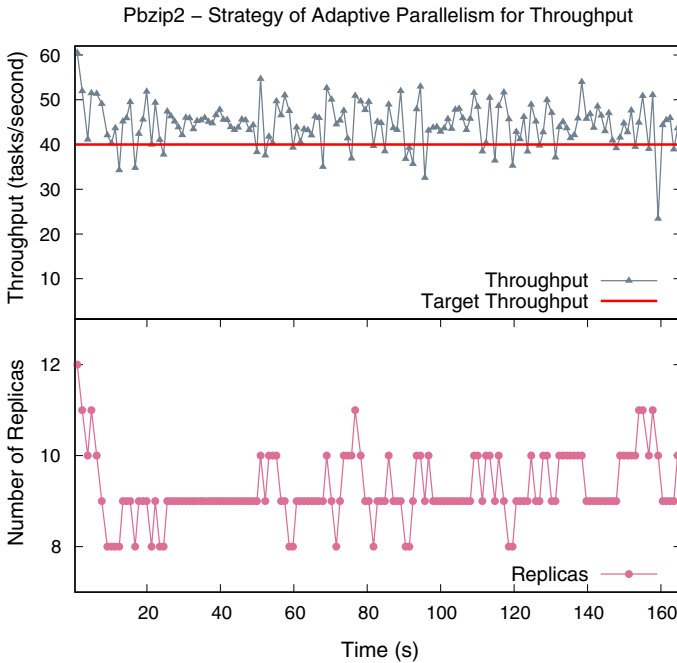


Fig. 5 M1-characterization of *Pbzip2* application with `slo::Throughput (40)`

containing the code relevant for the parallelization. In all the cases, parallelizing an application by using SPAR with SLOs requires a lower code intrusion with respect to FASTFLOW [17, 20]. Since SLOs can be defined by only inserting the objective through attributes, this practice reduces significantly the lines of code. The hand-written version increased significantly the lines of code because implementing a strategy requires implementing all the details relative to resource management and monitoring.

5.3 SLO analysis

In this section, we analyze the use of the SLO attributes laying emphasis in the generated self-adaptive strategies. The main goal is to provide a discussion regarding the adaptivity and effectiveness of the SLO strategies with a set of stream processing applications.

The `slo::Throughput` SLO can be used in any application parallelized with SPAR that has at least one replicated stage. In Fig. 5 is shown the result of *Pbzip* with a target throughput of 40 tasks per second, representing an SLO defined by the user. Figure 5 shows the measured throughput compared to SLO as well as the number of replicas used on each monitoring step. It is possible to observe that the throughput oscillated significantly during the execution, which is caused by the application and its input load characteristic. Because of the throughput oscillations, the self-adaptive

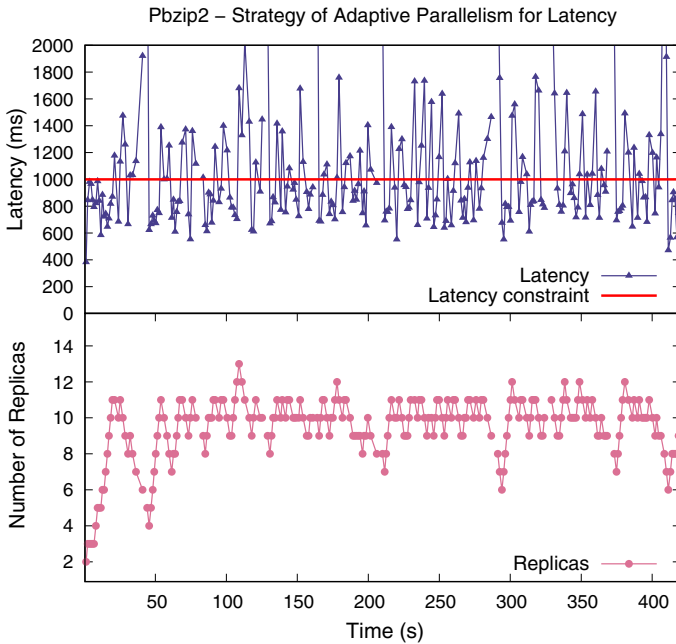


Fig. 6 M1-lane detection application with `slo::latency(1000 ms)`

strategy needed to change the number of replicas several times responding to throughput fluctuations and pursuing performance optimization. In some events, it is possible to note SLO violations caused by the execution variation. The strategy responded to such variations, but sometimes it was not fast enough since such short-duration load spikes occur randomly, and their prediction is not possible.

Regarding the `slo::Latency` attribute, we tested this SLO under different configurations to evaluate if the strategy impacts on the application performance. For instance, we tested in a video streaming application using a file as an input to simulate a typical execution. A representative outcome of this experiment is shown in Fig. 6 with a latency constraint of 1000 ms, which simulates the definition of a representative SLO by the user.

In Fig. 6, the strategy is characterized by the measured latency. We also presented the number of replicas used in different instants of the execution. Considering the results from Fig. 6, we can identify that the latency varied during the execution because some frames require more time to be processed, thus causing unpredictable variations. Despite the adaptive strategy changed the number of replicas when necessary responding to the latency oscillations, some SLO violations occurred due to such short-duration fluctuations. Under a more stable workload trend, the adaptive strategy is expected to find a suitable number of replicas and maintain this number throughout the entire execution.

Concerning the use of the `slo::CPU` SLO, Fig. 7 shows the execution of the Pbzip2 application with the attribute defining the maximum utilization to 60%, which is an empirically defined scenario, simulating an execution that could have

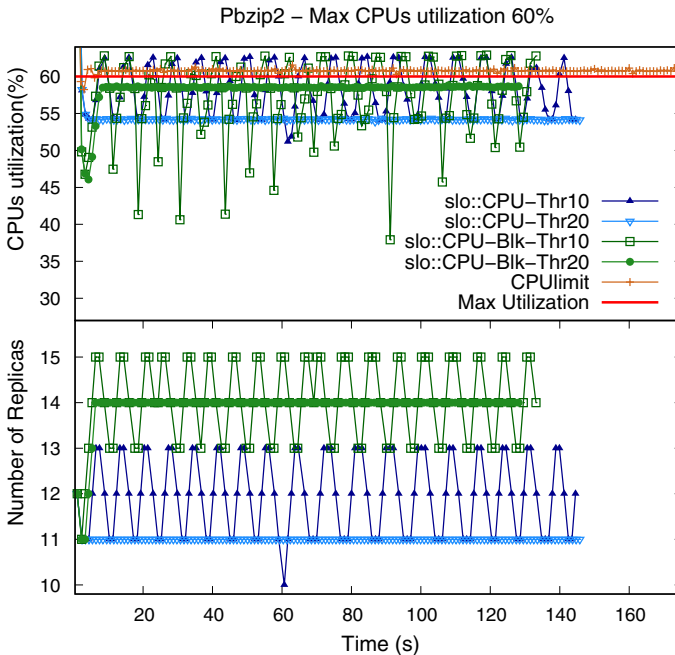


Fig. 7 M1-characterization of *Pbzip2* application with `s1o::CPU (60)`

a CPU load slightly higher than half of the machine's resources. Such a scenario is representative of applications running on shared environments. We tested this SLO strategy with two representative threshold values, using the environment variable (`SLO_THRESHOLD`): 10 and 20%. These were the most suitable thresholds for stream parallelism, as seen in [39]. We also ran one variant using the blocking mode (`-spar_blocking` compilation option in `SPAR`) that tends to consume fewer CPU resources by only distributing new tasks upon requests from the active threads. The results are compared to the `CPUlimit` utility tool, which also was set to limit the CPU usage in 60%. For the tests using `CPUlimit`, we set a number of application threads equal to the number of hardware threads, which is what is done by default in several runtimes. The self-adaptive strategy, on the other hand, uses a custom number of active threads by changing the status of the replicas at runtime according to the heuristic policy implemented (Sect. 4.1.3).

In the results from Fig. 7, we can observe that `CPUlimit` was unable to enforce the required SLO. It is relevant to highlight that all executions presented a high CPU utilization in the first second. This event is caused by the application startup routines, such as threads and queues creation. The threshold of 10% introduced instability by triggering too frequent changes in the number of replicas, which also induced variation in CPU utilization. On the other hand, the threshold of 20% was the most accurate and stable one. By using the `-spar_blocking` compilation flag, it reduced the CPU utilization. Consequently, this resulted in an opportunity to use more replicas in the parallel region.

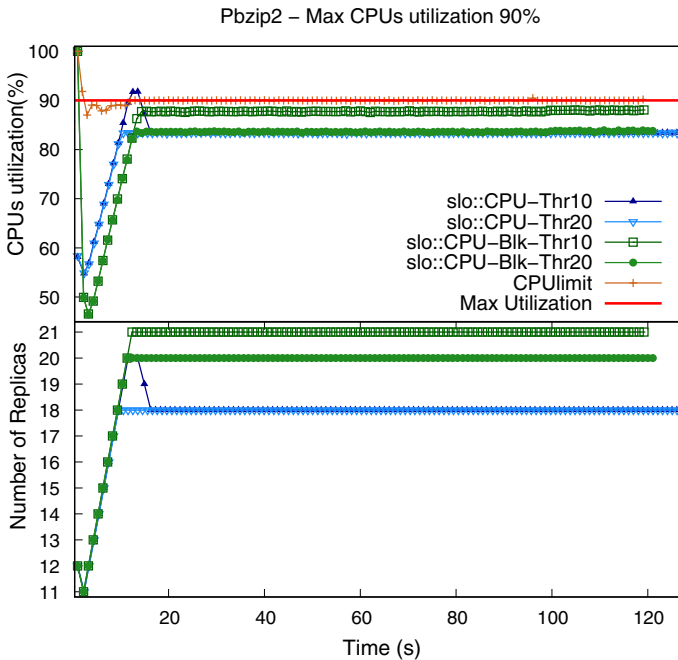


Fig. 8 M1-characterization of *Pbzp2* application with `s1o::CPU (90)`

Figure 8 introduces the results of CPU utilization with a higher SLO value of at most 90% CPU utilization. Such a scenario was tested in order to evaluate the executions when almost all the machine resources available could be used. It tends to impact in the number of replicas used by the adaptive strategy. Comparing the different versions, we can visualize that the threshold 20% again caused fewer SLO violations by reaching a stable number of replicas after the first calibration phase. The `CPUlimit` presented oscillations in the utilization while the `-spar_blocking` compilation flag again enabled the use of additional replicas and avoided SLO violations.

We now show the results obtained by running with the `s1o::CPU` SLO with all the considered applications. The results presented are an average of 10 executions. In Fig. 9, is shown the throughput of the execution considering the three representative applications and two representative `s1o::CPU` SLO configurations: 60 and 90%. It is important to note that the SLO strategies are compared to a static degree of parallelism version using the `CPUlimit` for SPAR and Intel TBB.

Considering the SLO of 60%, it is possible to identify a similar outcome regarding the different applications. In the self-adaptive executions, when using the `spar_blocking` compilation flag, it achieved a higher throughput rates than the default non-blocking execution. The self-adaptive strategy dynamically tunes the number of replicas resulting in the highest throughput rates. This result indicates that the way in which `CPUlimit` works (i.e., continuously pausing and

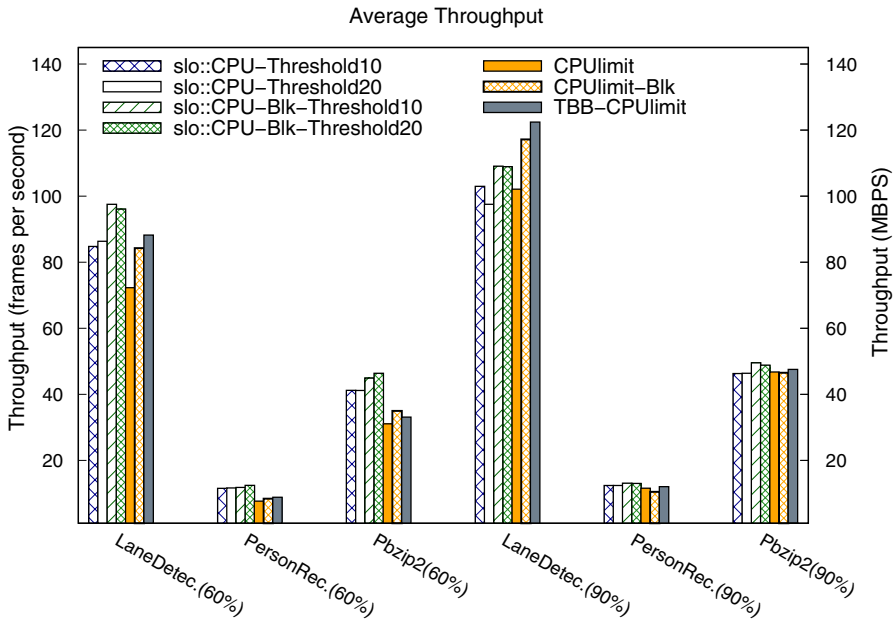


Fig. 9 M1-throughput of applications. Left side in frames per second refers to video applications. Right side in MBPS is related to the throughput of Pbzip2

resuming the target process) causes performance overhead. CPUlimit in SPAR had a lower throughput, while TBB and SPAR Blocking achieved better performance.

The result of running with `slo::CPU` in 90% showed similar results with respect to 60%. Although the contrasts between our generated self-adaptive strategy and CPUlimit were smaller, our strategy again was significantly better in most cases. In Lane Detection with 90% CPU utilization SLO, both TBB and SPAR blocking achieved the highest throughput. CPUlimit blocked significantly less the threads with the low CPU restriction of 90% CPU utilization SLO, which increased the application performance. In Lane Detection, the TBB version outperformed SPAR because TBB improves the load balancing, while in Person Recognition and Pbzip2 both versions achieved similar performance. Considering the different applications and their execution characteristics, it is possible to note that CPUlimit performed better in those applications with a more balanced load, while performed worst in the irregular processing applications (Person Recognition). This indicates that CPUlimit is not a suitable alternative for limiting CPU utilization in stream processing applications, which are usually unbalanced because of their intrinsic dynamic nature.

In order to further characterize CPUlimit, we also evaluated the impact of the number of replicas. Figure 10 presents the results on Pbzip with a representative `slo::CPU` SLO of 60%. In this test, the results from our self-adaptive strategy are compared to a static number of replicas in SPAR and TBB managed by CPUlimit. The throughput of our strategies is presented in all number of replicas because any of those numbers could be used during the execution, depending on the decisions made by the regulator algorithm. It is possible to note that the configuration using

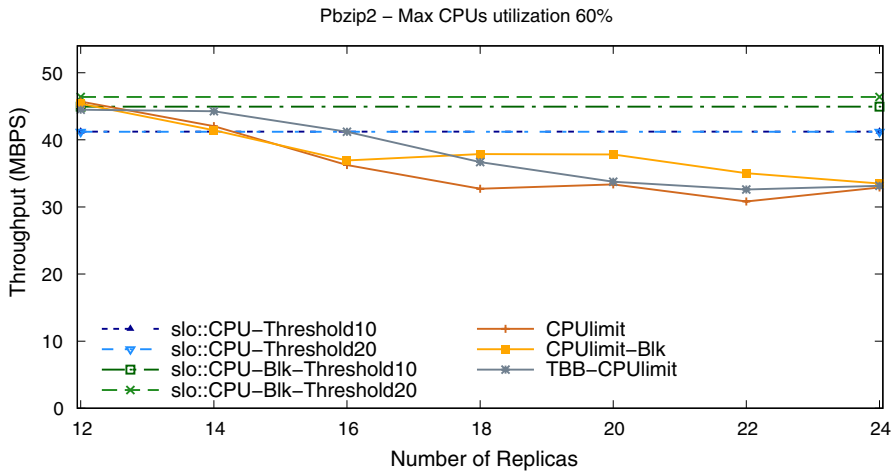


Fig. 10 M1-CPUlimit characterization with different number of replicas

12 replicas was the best CPUlimit configuration in SPAR and TBB, although the self-adaptive strategy in blocking mode still achieved the highest throughput. Regarding CPUlimit, the blocking mode only achieved a better performance in specific cases comparing to the default non-blocking mode. Comparing the results where TBB outperformed SPAR running with one application thread per hardware thread in Fig. 9, the several number of replicas in TBB only won with 14 and 16 replicas. On the other hand, SPAR with the blocking mode outperformed TBB in most cases.

The outcome from Fig. 10 highlights the correlation between the number of replicas and the application throughput, showing that using a tool like CPUlimit for limiting the CPU utilization SLO is inefficient in the stream processing context. The results indicate that even if CPUlimit is used, a suitable number of replicas has still to be found. However, finding a suitable number of replicas tends to be a complex task in stream processing applications. Additionally, the number of replicas often has to be adapted during execution according to performance or efficiency goals, because this class of applications runs without a defined end of the computation. Therefore, rerun the application multiple times until a suitable number of replicas is found, it becomes unfeasible for stream processing applications. Consequently, our strategy that dynamically adapts the number of replicas in SPAR at runtime is a feasible and effective approach, which showed promising performance outcomes.

In Fig. 11, we analyze a different scenario, where the user requires a throughput as well as an energy constraint. This scenario exploits the usage of energy strategies. The defined SLO throughput (`slo::Throughput`) was 20 tasks per second and power consumption (`slo::Power`) lower than 65 W for the *Pbzip* application. In this test, we add some external noise to show that our strategy for controlling performance and energy succeeds in providing the required SLO even in the presence of unexpected behaviors.

In particular, besides the usual calibration done in the first seconds of execution, after 50 s from the start of *Pbzip*, we start another application on the same machine.

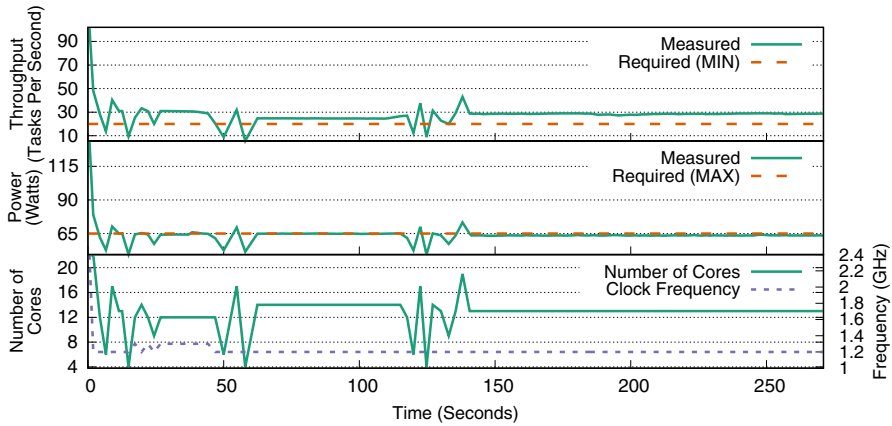


Fig. 11 M2-*Pbzip2* application with `slo::Throughput` (20) and `slo::Power` (65)

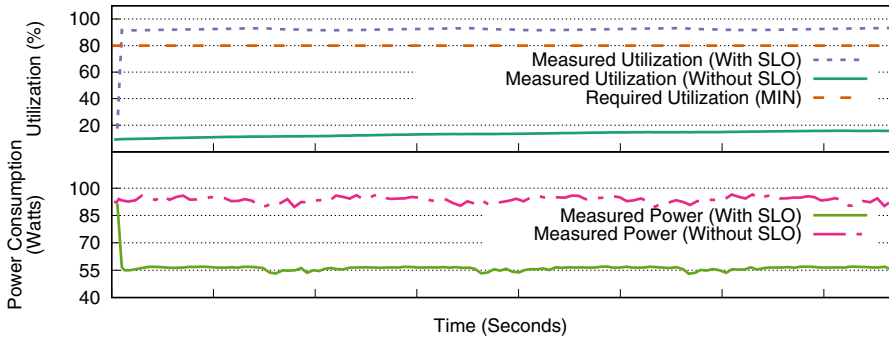


Fig. 12 M2-*lane detection* application with `slo::Utilization` (80)

Since the two applications share some resources (i.e., cores, memory, among others), the throughput of *Pbzip2* starts to decrease. In response to this issue, our generated code and the compatible runtime recompute the prediction models, now considering the presence of external interference. As a consequence, as we can see from the bottom part of Fig. 11, our generated strategy with its runtime increases the number of replicas of the middle stage from 12 to 14. When the other interfering application terminates (around 120 s from the start of *Pbzip2*), our generated strategy recomputes the models and decreases the number of replicas from 14 to 13. As we can see from the two upper parts of the figure, our generated strategy uses its runtime to satisfy the user requirements throughout the entire execution (excepts for the phases where the models are computed), independently from the presence of other applications running on the system.

In Fig. 12, we analyze the *Lane Detection* application, in a scenario where it produces no more than 50 frames/s. In such a case, using all the available resources could be inefficient, since they could be idle for most of the time. To avoid such

Table 3 Power consumption reduction obtained by a parallel application with the same throughput of the sequential one

	Pbzip2	Lane detection	Person recognition
Power consumption reduction (%)	- 9.43%	- 10.37%	- 7.39%

scenario, we set a utilization SLO (`slo::Utilization`) of 80%. In the upper part of Fig. 12, we report the utilization when an SLO is specified and when it is not specified. In the bottom part, we report the power consumption. As shown by the result when an SLO is not specified, the utilization would be around 20%. This utilization means that the threads of the application would spend 80% of the time waiting for new frames to arrive. By requiring a minimum utilization of 80%, our generated strategy decreases the number of resources allocated to the application, decreasing the power consumption from 90 to 55 W. This event occurs without decreasing the overall performance of the application. Indeed, the threads still spend some time waiting for new data, but it is reduced from 80 to 5% (the utilization is around 95%).

The target of the experiment in Table 3 is to demonstrate that parallelization is not only useful for improving the performance of an application, but it can also be used to reduce its power consumption. In a nutshell, we want to show that a parallel application with the same performance of the sequential one has lower power consumption. We were able to limit the SLO throughput by using the `slo::Throughput` combined with `slo::Power`.

The interpretation we would like to give to these results is that, even if the performance of a sequential application is satisfactory, parallelizing it may still be useful for reducing its power consumption. This effect occurs since by increasing the number of replicas (and thus the number of cores used by the application); we can reduce the clock frequency while keeping the same performance. Since the power consumption increases linearly with the number of cores but more than quadratic with the clock frequency [8], running an application on more cores at a lower frequency is usually more energy efficient than running it on fewer cores at a higher frequency. Having tools and methodologies for doing that automatically and with low code intrusion, like those we proposed through SLO attributes in this work is of paramount importance for enabling such techniques in real-world scenarios.

6 Conclusion

In this work, we presented a new and simpler way to express SLOs in sequential source codes. Our approach was designed to target stream processing applications along with its parallelization support. We validated this approach by implementing it in the SPAR language and compiler, which now recognizes the C++11 SLO attributes and automatically performs source-to-source transformations to the self-adaptive strategies implemented in the FASTFLOW and NORNIR libraries. The main advantage is that application programmers can now simply define SLOs by inserting

the attributes in the source code, and the compiler generates the appropriate self-adaptive strategy to meet the target SLO. This new approach does not require from programmers implementation expertise either system resource management.

Moreover, our implemented solution has proven to be efficient and offers many opportunities to improve the QoS in stream processing applications. We were able to reduce power consumption and increase performance in certain cases. Regarding the CPU utilization SLO, the performance was improved in most cases compared with CPUlimit. While our strategy relies on changing the number of replicas, CPUlimit works at the operating system level limiting the CPU utilization by continuously pausing and resuming the target process. Although the goals and efforts in this work were more in the abstraction of SLOs implementation, we visualize a set of future works. For instance, a deep performance validation can be conducted to cover different workloads and stream processing scenarios. We also plan to implement other self-adaptive strategies and SLOs. We would like to refine our `slo::CPU` and `slo::Latency` SLOs to combine them with power consumption SLOs. Eventually, our approach could be extended to other computing environments such as cloud or cluster architectures.

Acknowledgements This study was partially funded by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior-Brasil (CAPES)-Finance Code 001 and by the FAPERGS 01/2017-ARD Project PARAElastic (No. 17/2551-0000871-5). We would like to thank Laboratório de Alto Desempenho (LAD) from PUCRS for partially providing computing resources.

References

1. Aldinucci M, Danelutto M, Kilpatrick P, Torquati M (2014) FastFlow: high-level and efficient streaming on multi-core. In: Programming multi-core and many-core computing systems, vol 1, PDC. Wiley, p 14
2. Aldinucci M, Meneghin M, Torquati M (2010) Efficient smith-waterman on multi-core with fast-flow. In: Proceedings of the Euromicro Conference on Parallel, Distributed and Network-Based Processing, pp 195–199
3. Alessi F, Thoman P, Georgakoudis G, Fahringer T, Nikolopoulos DS (2015) Application-level energy awareness for OpenMP. In: International workshop on OpenMP. Springer, pp 219–232
4. Andrade HCM, Gedik B, Turaga DS (2014) Fundamentals of stream processing. Cambridge University Press, New York
5. Ansel J, Pacula M, Wong YL, Chan C, Olszewski M, O'Reilly U-M, Amarasinghe S (2012) Siblingrivalry. In: Proceedings of the 2012 International Conference on Compilers, Architectures and Synthesis for Embedded Systems-CASES '12. ACM Press, New York, pp 91
6. Beyer B, Jones C, Petoff J, Murphy NR (2016) Site reliability engineering. O'Reilly, Boston
7. Buyya R, Vecchiola C, Selvi T (2013) Mastering cloud computing. McGraw Hill, New York
8. Chandrakasan AP, Brodersen RW (1995) Minimizing power consumption in digital CMOS circuits. Proc IEEE 83(4):498–523
9. CPUlimit (2018) CPU Usage Limiter for Linux roadmap. <http://cpulimit.sourceforge.net/>. Last access Dec, 2018
10. Curtis-Maury M, Blagojevic F, Antonopoulos CD, Nikolopoulos DS (2008) Prediction-based power-performance adaptation of multithreaded scientific codes. IEEE Trans Parallel Distrib Syst 19(10):1396–1410
11. Danelutto M, Garcia JD, Sanchez LM, Sotomayor R, Torquati M (2016) Introducing parallelism by using REPARA C++11 attributes. In: Euromicro International Conference on Parallel, Distributed, and Network-Based Processing. IEEE, pp 354–358

12. De Sensi D, De Matteis T, Danelutto M (2018) Simplifying self-adaptive and power-aware computing with Nornir. *Future Gener Comput Syst* 87:136–151
13. De Sensi D, Torquati M, Danelutto M (2016) A reconfiguration algorithm for power-aware parallel applications. *ACM Trans Archit Code Optim* 13(4):43:1–43:25
14. Delimitrou C, Kozyrakis C (2014) Quasar: resource-efficient and qos-aware cluster management. *SIGARCH Comput Archit News* 42(1):127–144
15. Ding Y, Kandemir M, Raghavan P, Irwin MJ (2008) A helper thread based edp reduction scheme for adapting application execution in cmps. In: *IEEE International Symposium on Parallel and Distributed Processing, 2008. IPDPS 2008*, pp 1–14
16. Floratou A, Agrawal A, Graham B, Rao S, Ramasamy K (2017) Dhalion: self-regulating stream processing in heron. *Proc VLDB Endow* 10:1825–1836
17. Griebler D, Danelutto M, Torquati M, Fernandes LG (2017) SPar: a DSL for high-level and productive stream parallelism. *Parallel Proc Lett* 27(01):1740005
18. Griebler D, De Sensi D, Vogel A, Danelutto M, Fernandes LG (2018) Service level objectives via C++11 attributes. In: *Euro-Par 2018: parallel processing workshops*. Springer, Turin, pp 12
19. Griebler D, Hoffmann RB, Danelutto M, Fernandes LG (2017) Higher-level parallelism abstractions for video applications with SPar. In: *Parallel Computing is Everywhere, Proceedings of the International Conference on Parallel Computing, ParCo'17*. IOS Press, Bologna, pp 698–707
20. Griebler D, Hoffmann RB, Danelutto M, Fernandes LG (2018) High-level and productive stream parallelism for Dedup, Ferret, and Bzip2. *Int J Parallel Program* 47:253–271
21. Hellerstein JL, Diao Y, Parekh S, Tilbury DM (2004) *Feedback control of computing systems*. Wiley, New York
22. Hoffmann H, Sidiroglou S, Carbin M, Misailovic S, Agarwal A, Rinard M (2011) Dynamic knobs for responsive power-aware computing. *SIGPLAN Not* 46(3):199–212
23. Kephart JO, Chess DM (2003) The vision of autonomic computing. *Computer* 36(1):41–50
24. Li J, Martinez JF (2006) Dynamic power-performance adaptation of parallel computation on chip multiprocessors. In: *Proceedings of the International Symposium on High-Performance Computer Architecture*, pp 77–87
25. Maggio M, Hoffmann H, Santambrogio MD, Agarwal A, Leva A (2010) Controlling software applications via resource allocation within the heartbeats framework. In: *IEEE Conference on Decision and Control*. IEEE, pp 3736–3741
26. Maurer J, Wong M (2008) Towards support for attributes in C++ (revision 6). Technical report, The C++ Standards Committee
27. McCool M, Robison AD, Reinders J (2012) *Structured parallel programming: patterns for efficient computation*. Morgan Kaufmann, Burlington
28. Petrica P, Izraelevitz AM, Albonesei DH, Shoemaker CA (2013) Flicker: a dynamically adaptive architecture for power limited multicore systems. *ACM SIGARCH Comput Archit News* 41(3):13
29. Pusukuri KK, Gupta R, Bhuyan LN (2011) Thread reinforcer: dynamically determining number of threads via os level monitoring. In: *Proceedings of the 2011 IEEE international symposium on workload characterization, IISWC '11*. IEEE Computer Society, Washington, DC, pp 116–125
30. Reinders J (2007) *Intel threading building blocks*. O'Reilly, New York
31. Shafik RA, Das A, Yang S, Merrett G, Al-Hashimi BM (2015) Adaptive energy minimization of OpenMP parallel applications on many-core systems. In: *Parallel programming and run-time management techniques*, pp 19–24
32. Sridharan S, Gupta G, Sohi GS (2013) Holistic run-time parallelism management for time and energy efficiency. In: *Proceedings of the 27th International ACM Conference on International Conference on Supercomputing-ICS '13*. ACM Press, New York, pp 337
33. Sturm R, Morris W, Jander M (2000) *Foundations of service level management*. SAMS, Boston
34. Suleman MA, Qureshi MK, Patt YN (2008) Feedback-driven threading. In: *Proceedings of the 13th International Conference on Architectural Support for Programming Languages and Operating Systems-ASPLOS XIII*, vol 42. ACM Press, New York, pp 277
35. Thies W, Karczmarek M, Amarasinghe SP (2002) StreamIt: a language for streaming applications. In: *Proceedings of the International Conference on Compiler Construction*. Springer, Grenoble, pp 179–196
36. Totoni E, Jain N, Kalé LV (2015) Power management of extreme-scale networks with on/off links in runtime systems. *TOPC* 1(2):16
37. Totoni E, Torrellas J, Kale LV (2014) Using an adaptive hpc runtime system to reconfigure the cache hierarchy. In: *Proceedings of SC 2014*. IEEE Press, pp 1047–1058

38. Vassiliadis V, Parasyris K, Chalios C, Antonopoulos CD, Lalis S, Bellas N, Vandierendonck H, Nikolopoulos DS (2015) A programming model and runtime system for significance-aware energy-efficient computing. *SIGPLAN Not* 50(8):275–276
39. Vogel A, Griebler D, Sensi DD, Danelutto M, Fernandes LG (2018) Autonomic and latency-aware degree of parallelism management in SPAr. In: *Euro-Par 2018: parallel processing workshop*. Springer, Turin, pp 12
40. Weyuker EJ (1988) Evaluating software complexity measures. *IEEE Trans Softw Eng* 14(9):1357–1365

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.