# Online and transparent self-adaptation of stream parallel patterns

**Adriano Vogel[1,2]** · **Gabriele Mencagli[2]** · **Dalvan Griebler[1,3]** ·
**Marco Danelutto[2]** · **Luiz Gustavo Fernandes[1]**

## Abstract

Several real-world parallel applications are becoming more dynamic and long-running, demanding online (at run-time) adaptations. Stream processing is a representative scenario that computes data items arriving in real-time and where parallel executions are necessary. However, it is challenging for humans to monitor and manually self-optimize complex and long-running parallel executions continuously. Moreover, although high-level and structured parallel programming aims to facilitate parallelism, several issues still need to be addressed for improving the existing abstractions. In this paper, we extend self-adaptiveness for supporting autonomous and online changes of the parallel pattern compositions. Online self-adaptation is achieved with an online profiler that characterizes the applications, which is combined with a new self-adaptive strategy and a model for smooth transitions on reconfigurations. The solution provides a new abstraction layer that enables application programmers to define non-functional requirements instead of hand-tuning complex configurations. Hence, we contribute with additional abstractions and flexible self-adaptation for responsiveness at run-time. The proposed solution is evaluated with applications having different processing characteristics, workloads, and configurations. The results show that it is possible to provide additional abstractions, flexibility, and responsiveness while achieving performance comparable to the best static configuration executions.

**Keywords** Autonomic systems · Parallel programming · Parallelism abstractions · Self-adaptive software · Stream processing

✉ Adriano Vogel
adriano.vogel@edu.pucrs.br

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

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

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

**Mathematics Subject Classification**  68W10 · 68T05

## 1 Introduction

Large amounts of data are being generated due to the proliferation of devices (e.g., sensors, cameras) used to sense the external world. In this scenario, parallel computing is relevant for processing fast enough the high amount of data being generated [13,18]. Moreover, the continuous data arrival requires stream processing applications to run for long or infinite periods, where such long executions are often subject to fluctuations in the environment (e.g., resource availability) and at the application level (input rate, workload) [24]. Hence, self-adapting entities (e.g., degree of parallelism, cores, and their frequencies) during the execution is important for achieving responsiveness [4,5,15,23].

From a programming perspective, structured parallel programming facilitates the task of parallelizing applications. Programmers can instantiate high-level pattern constructors and combine them in compositions [1,3]. In this scenario, online changing the pattern composition configurations[1] has been proposed for abstracting complexities from application programmers and increasing the adaptation space to provide flexibility [20]. This can reduce the burden on application programmers in such a way that configurations are adapted transparently [19]. However, the configuration space is usually large, which is challenging to find optimal configurations at run-time.

Moreover, more generic strategies for self-adaptive decision-making are needed, and dynamic changes can have detrimental effects on the Quality of Services (QoS). Hence, in previous work [24], we contributed with mechanisms for online self-adaptiveness of parallel patterns. The solution was integrated with a C++ programming framework (FastFlow [2]) and experimentally evaluated. In this extended version, we provide the following contributions:

– An autonomous self-adaptive strategy that avoids suboptimal configurations, which encompasses a lightweight online profiler of the application stages and an optimized decision-making for accuracy. The new strategy also supports latency as a new Service Level Objective (SLO). SLO refers to a metric of interest and its proper value to be enforced [10,12].
– A model for smooth transitions between the parallel pattern configurations. A smooth transition is important because changing the configurations can have a critical impact on the QoS of applications (see Sect. 3.3).
– Extended validation of the proposed solution, including new scenarios and applications. Noteworthy, we provide a custom version of the PARSEC's Ferret application to regulate the Input Rate (IR) and support user-defined SLOs (throughput, latency).

This paper is organized as follows. Section 2 shows the motivational context of this work. Then, Sect. 3 presents the proposed solution and Sect. 4 provides an experimental

---

[1] The term composition refers to the application topology (a.k.a. stream graph, graph topology). Here, the terms composition and configuration are used interchangeably.

evaluation. Moreover, Sect. 5 discusses related approaches and Sect. 6 concludes this paper.

## 2 Problem statement and motivation

The use of high-level parallel programming methodologies is a potential alternative to provide coding abstractions for application programmers [3], which can reduce the application programmers' burden. High-level parallel programming aims at reducing the programming efforts while ensuring proper performance. In this vein, pattern-based parallel programming provides composable recurrent structures instantiated by programmers, combining the patterns creating different configurations. Stream processing applications running in multicore machines can use high-level parallel programming frameworks, where Intel Threading Building Blocks (TBB) [25] is an example from the industry and FastFlow [2], GrPPI [7], and SPar [8] are academic frameworks.

In TBB, the application programmers can create a parallel pipeline by declaring each function as a filter and defining if a stage is parallel or sequential. In FastFlow, the user can also create pipelines and replicate (run in parallel) specific stages using the Farm pattern.

TBB creates tasks that are scheduled to a pool of threads in a runtime system's perspective, where dynamic scheduling controls thread oversubscription by avoiding context switching and time-sliced execution. However, TBB can incur scheduling overheads with fine-grained tasks and I/O blocking operations. FastFlow, on the other hand, avoids these issues with a runtime where nodes are fixedly mapped onto threads, and the runtime can statically merge the nodes without changing the user functions. Nevertheless, FastFlow has a rigid execution model that may not suit stream processing with more irregular and dynamic applications. This model may increase the demand for resources without guaranteeing performance gains. Hence, we argue that there is a need to support adaptation in existing programming frameworks.

In previous work [24], we evinced that streaming applications computing data in real-time require the programmers to create a configuration of sequential or replicated stages. However, maintaining such a configuration for the entire execution can be limited due to the dynamic nature of applications. For instance, it was demonstrated with a video stream processing application executing under fluctuations in the IR that the best configuration combining replicated, sequential, and merged stages varies with different scenarios that occur at run-time and from one programming framework to another.

The complexity increases robust applications that have several sequential or replicated stages. For instance, Fig. 1 shows the structure of the Ferret [17] application from the PARSEC suite, where the four middle stages are thread pools that can run in parallel. However, the profitability of running them in parallel can vary from how balanced the stages are and according to the expected performance and QoS. Ferret has a robust structure modeled with different shapes by composing and nesting parallel patterns [6]. Figure 2 provides performance results with different Ferret's configurations, where the setup is described in Sects. 4.1 and 4.2.2. This new streaming version
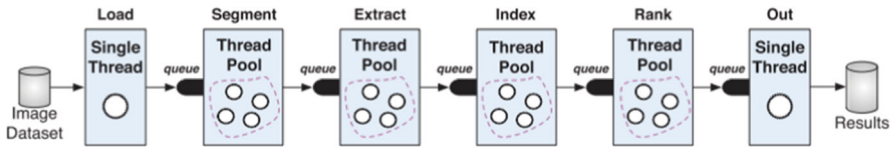
**Fig. 1** PARSEC's Ferret pipeline structure. Extracted from [6]



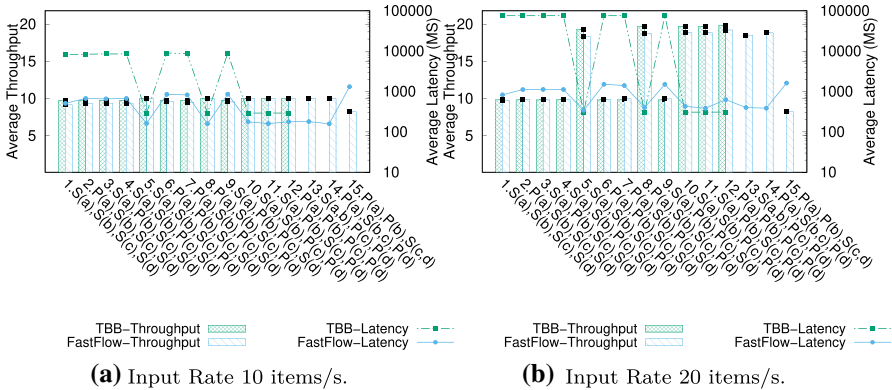**(a)** Input Rate 10 items/s.          **(b)** Input Rate 20 items/s.

**Fig. 2** Ferret stages configurations. Latency is on logarithmic scale

of Ferret computes data at a given IR and provides stream processing performance metrics like throughput and latency. In Fig. 2a the IR is 10 items per second, and 10 is a suitable throughput (items/s) for sustaining the IR. Latency is another relevant metric that corresponds to the time taken to compute a given item, where low latency is a constraint for many applications [4].

Here we show the performance from 15 configurations described in Sect. 4.2.2. There are no results from configurations 13, 14, and 15 with TBB because merging functions would require refactoring the business logic code. Configuration 12 corresponds to the native implementation using Pthreads where all stages are parallel[2], demanding more resources and not performing so well in terms of latency. In fact, the best performance is when the last stage is parallel. Although declaring such a configuration is trivial, it is not intuitive for application programmers and significantly impacts QoS.

Notably, the best configuration to be used at run-time can vary due to stream processing fluctuations [15,24], e.g., the IR can change due to network fluctuations or variations in the number of devices producing data. Consequently, flexible adaptations are expected to be applied at run-time. But, several aspects correlate nonlinearly, where we argue that the programmers should not be expected to online hand-tune the configurations. In Sect. 3 we present the proposed solution tackling the aforementioned challenges and providing abstractions with dynamic self-adaptation at run-time.

---

[2] The Pthreads version is not structured pattern-based, such results are not shown here for the sake of visual clarity. The performance of FastFlow and TBB is comparable with the native implementation. The reader interested in such a comparison can refer to reference [6].

## 3 Proposed solution

Here we describe the solution that we envision for providing additional abstractions to application programmers and to enable flexible online self-adaptation at run-time. In our perspective, a feasible solution is to support users/programmers to set only high-level goals like throughput or latency SLOs, which can be offered as parameters or attributes [10].

### 3.1 General design goals and requirements

An effective approach for dynamic parallel pattern compositions is expected to meet the following goals:

– Relieve application programmers from the burden to find the best configuration of parallel patterns. The programmers should set SLOs instead of hand-tuning configurations. This can be achieved with autonomous strategies that can provide a suitable QoS (achieves the SLO consuming fewer resources).
– Enable dynamic and flexible reconfigurations at run-time to avoid the need to recompile and rerun long-running stream processing applications.
– Respond at run-time to changing conditions, where dynamic reconfigurations can enable self-adaptiveness.

There are also important **requirements** for ensuring QoS and efficiency:

– No application downtime. The adaptation should not interrupt the data processing and output provisioning.
– Smooth transition on reconfigurations. Change from one composition to another can be necessary. However, this is expected to be stable without intrusiveness and overheads like latency glitches and throughput spikes [20].
– A suitable solution is expected to be lightweight and execute without demanding a significant extra amount of resources.
– Efficiency. An optimal configuration is one that meets user/programmer goals and that requires fewer computing resources. Consuming fewer resources increases the system efficiency, reduces energy consumption, and costs less (i.e., pay-per-use environments).

### 3.2 Decision-making strategy

A decision-making strategy is the core of a self-adaptive strategy responsible for deciding the best actions to be enforced. However, assumptions are necessary for designing a flexible and generalizable decision-making strategy. The rationale for such assumptions is to abstract technicalities that have to be implemented for each specific scenario. In Sect. 3.4 we show an example of implementation in a C++ programming framework. The main necessary assumptions are:

– Runtime system's mechanisms are available supporting dynamic changes to configurations from one configuration to another.

– The strategy receives alternative configurations to be considered at run-time. Such configurations could be defined by a user or by an expert system.
– The strategy receives information for making decisions, which can be provided by external monitoring entities.
– The data to be processed comes at a given IR and the strategy is alerted in case the IR changes.

Considering that self-adaptation is expected to encompass relevant properties, **SASO** (stability, accuracy, settling time, and overshoot) properties [12] are relevant ones to be considered when designing the proposed solution. *Stability* refers to the capacity of producing the same output under a given condition. *Accuracy* is related to achieving the control goal with sufficiently good decision-making, and *Short settling times* are desired for reaching fast enough an optimal state. Moreover, *overshooting* should be avoided by using only the amount of resources needed. In this vein, the designed decision-making for the pattern composition configurations is described in a manner that abstracts lower implementation details. However, the description provided is expected to be sufficient for replicating the proposed solution. The decision-making is built out of the following steps:

1. *Online profiling step* Lightweight instrumentation gathers execution statistics from each stage, which helps in characterizing how intensive and balanced they are. The profiling step measures the actual processing capacity of each stage and ranks them by less to more intensive. Moreover, a given stage is tagged if its average service time (time spent computing the tasks) is at least 20% higher than all other service time of stages. This step is executed at the beginning of the execution with the first configuration provided. It can be repeated at any time, such as when a given application enters a new processing phase. For increasing the profiling accuracy, it is recommended that the first configuration executes all stages sequentially.
2. *Evaluation* Assesses if the defined SLO is satisfied. If positive, goes to step 6 with the current configuration. If not, goes to step 3. The decision-making strategy infers that two values are significantly different when they contrast higher than 20% (a threshold ), which is a configurable parameter that the used value of 20% was ascertained in previous work [23].
3. *Shortlisting configurations* Previous work [24] applied experimental runs with all configurations. In practice, this can affect QoS because bad configurations could be used. Considering that the new proposed strategy aims to reduce the settling time without limiting the configuration space, the profiling step's information is used to shortlist the potentially optimal configurations. If more than one configuration can be optimal, the strategy goes to step 4, or if only one is suitable, sets this one as active and go to step 6. Also, a configuration with the most replicated stages is set if the SLO is not being achieved and there are no bottlenecks or optimal configurations.
4. *Trial phase* Activates each suitable configuration candidate for a given time interval and gathers execution statistics. The rationale for executing each shortlisted configuration is to increase the accuracy by finding the configurations that perform better for the specific application, workloads, and environments. Considering that the time interval in which each configuration is tested is a relevant parameter that

expert users can customize, 5 s is the default value ascertained from empirical results. The previous implementation from [24] tested all configurations for only 1 second because it did not profile and shortlisted the best ones. In practice, we have seen that one second as time interval is too low and subject to unpredictable variations during the training step. In the current optimized version, 5 seconds is the proper time interval because a suboptimal configuration will not be tested as these do not pass the shortlisting step.

5. *Selects the best configuration* This phase evaluates which configurations from step 5 achieved the desired SLO. If no configuration achieved the goal, enforces the one with the best value. On the other hand, if more than one is optimal, select the one with light stages merged and fewer replicated stages. This decision is to enforce the most optimal configuration that maintains QoS and at the same time consumes fewer resources for avoiding overshooting.

6. *Steady state* Stabilizes in a configuration and periodically evaluates if the SLO is being satisfied. In practice, every 10 s, the current status is verified. This comparison uses the same threshold from step 2 to avoid the instability caused by fluctuations. Steps 3 to 5 are repeated if the data gathered indicates changes or if the SLO is violated. In this case, it is searched for additional bottlenecks and potentially optimal configurations.

It is important to note that the decision-making strategy is not employing an exhaustive search. In fact, up to 20 alternative configurations are supported. However, only the suitable ones are to be activated and tested at run-time. Moreover, Sect. 3.3 addresses the relevant aspect of how to achieve safety and stability when transitioning from one configuration to another.

### 3.3 Transitioning between configurations

Reconfigurations at run-time should be smooth such that they do not compromise the QoS. One possible solution is to employ a draining phase that flushes all the tasks from the configuration to be stopped before activating the new one. From a theoretical standpoint, a flush is relevant for avoiding that two configurations run simultaneously, which would cause unpredictable performance variability or losses (throughput spikes and latency glitches [20]).

One may think that the draining phase is a trivial problem solved by simply waiting for a random time. However, we have seen that choosing for how long to wait for the draining to complete is a non-trivial value in practice. On the one hand, not waiting for enough causes performance and resource fluctuation, influencing the training step and QoS. On the other hand, waiting for too long on reconfigurations can also hurt QoS and the designed goal of avoiding application downtime. Consequently, we tackled this challenge by developing an autonomous model that automatically estimates how long to wait. Such a model is mainly expected to find a balance value being accurate, generic, and lightweight. The draining time estimation considers the following samples/entities subject to variance on parallel applications:

*Number of items buffered* This aspect refers to buffer sizes used in the runtime system and the number of computing elements (*e.g.,* nodes) that use buffers for com-

municating in a given composition. For a generalization purpose, we assume that the runtime system provides mechanisms for collecting this value or provides parameters for limiting the buffer's sizes.

*Computations' service time* Considering that applications have significant contrasts in terms of grain and tasks computational intensiveness, the service time is expected to be a broad metric and flexible for different applications. A monitor can gather data and feed the model with the information of the average service time of the tasks being processed at a given moment, which corresponds to a given active configuration.

*Processing capacity* Refers to the computation capacity of the active configuration to process the tasks buffered and finalize the draining phase. We have discussed in Sect. 2 that each configuration and programming framework has specific processing capacities in terms of the number of nodes and the mapping to threads. Consequently, only considering the service time and the number of buffered tasks would be suboptimal because the actual computational capacity of each configuration varies. Generally, the processing capacity considers the number of computing elements that compute a given application's business logic code. Additional nodes/elements that do not process business logic code necessary in the programming framework should not be included in the processing capacity.

From the provided description, it is possible to note that the model is not simple and must consider the variability of service time, runtime parameters, and processing capacity. Moreover, the model must continuously measure and accurately estimate the time to drain. Considering the potential overhead of the machinery to collect and process data at run-time, in Sect. 4.3, we characterize the transitioning between configurations using this model.

### 3.4 Implementation

C++ frameworks and libraries available were considered for implementing the proposed solution. There are industry and academic solutions such as Intel TBB [25], FastFlow [2], and SPar [10]. Considering the support for performing adaptations at run-time, TBB has mechanisms only for dynamic task distribution and load balancing, where other mechanisms have to be implemented by hand. Considering that we are interested in higher-level abstractions, FastFlow is more flexible by supporting dynamic adaptation on several aspects like the parallelism degree and communication queues' concurrency modes [5,22].

FastFlow provides building blocks components to increase the flexibility in the parallelism exploitation. While parallel patterns are highly specialized structures that usually provide less flexible implementations, the building blocks approach incorporated in FastFlow [1,22] provides reusable components that can be combined and customized by programmers. Considering that building blocks are at a lower level, a drawback may consist of fewer parallelism abstractions available. It is believed that higher abstractions and flexibility can be provided by enabling building blocks compositions to be self-adapted at run-time, which can be achieved with practical implementations using FastFlow's mechanisms.
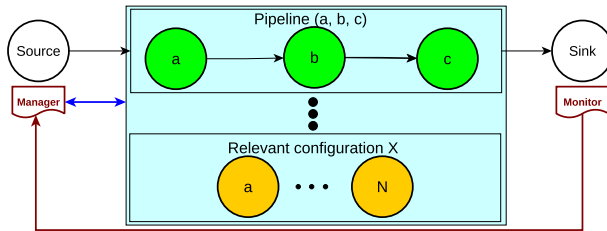
**Fig. 3** Implementation in FastFlow

Abstracting specific implementation technicalities, the proposed strategy was implemented in FastFlow as a ready-to-use C++ library, which works by default in FastFlow's blocking mode. Figure 3 provides a representation where the *Manager* is the entity that implements the self-adaptive strategy, which is embedded in the data source and uses the runtime system's mechanisms for autonomously applying changes. The *Monitor* is implemented as another entity within the Sink stage that gathers data and feds the *Manager*. Figure 3 represents a scenario where a pipeline with 3 stages is active and any other configuration that could declared would remain inactive. Moreover, the lower part of Fig. 3 demonstrates the achievable flexibility because several other configurations can be created and activated at run-time if necessary.

In the current implementation, the applications' business logic code is reused by alternative configurations. For using the proposed solution, the application programmers have to include self-adaptive strategy headers and add two extra code lines for calling the *Manager* and *Monitor*. Moreover, the higher accuracy of the profiling step demands on each application stage that a timer is initialized and that the *Monitor* is called.

The programmer can declare custom configurations using the FastFlow interface. For instance, the three staged pipeline used in Fig. 3 can be declared and added with two C++ code lines, and other compositions can be declared and included with similar coding productivity. The characteristics of the configurations (e.g., buffer sizes, stages are sequential, parallel, or merged) also have to be defined at a higher parametric description. The application programmer can be assisted with tools for designing additional configurations and coding, such as RPL [14] and SPar's compiler [8].

## 4 Evaluation

### 4.1 Experimental setup

We used a multicore machine equipped with two Intel Xeon E5-2620 processors (a total of 12 cores-24 threads), 32 GB of memory for running experiments, Ubuntu Server 16.04 as operating system, and G++ compiler (7.5.0) with O3 compilation flag. The FastFlow runtime system's buffer sizes were set to 1.

The strategy is characterized in a scenario simulating IR changes. The performance is evaluated with static configuration executions using the same configurations as a baseline. We call static configuration the executions where a given configuration is

compiled and maintained during the entire execution. The execution's correctness was checked by hashing the outputs. In Sect. 4.2 we describe the applications and configurations tested. Then, in Sect. 4.3 the decision-making is characterized with the different SLOs supported and application characteristics, Sect. 4.4 evaluates the performance of self-adaptive executions compared to baseline static executions, and Sect. 4.5 provides an overview of the results.

## 4.2 Applications and configurations

The evaluation of the proposed solution covers different applications and configurations. Considering that each application has a specific number of stages, workload pattern, and balance between stages, for each application we created a scenario of relevant configurations to be available for the self-adaptive strategy to use (or not) at run-time. In this evaluation, configurations using parallel stages use the default value of 2 replicas (parallelism degree) per stage.

### 4.2.1 Synthetic application

"Synthetic" is an application where 10,000 tasks with a total service time of 24 milliseconds (ms) are computed. This application has three functions where different configurations can be composed with a sequential or parallel stage. *Configuration 1* has the three sequential stages representing scenarios where the performance demand is not high, and the stages are balanced. *Configuration 2* has the first stage computing in parallel and stages 2 and 3 are sequential. Such a configuration can be suitable when the stages are not balanced, and the parallel stage is the bottleneck. *Configuration 3* has the second stage computing in parallel and the stages 1 and 3 sequential and *Configuration 4* has the third stage computing in parallel and the stages 1 and 2 sequential.

In *Configuration 5* stages 1 and 2 execute in parallel and the third stage is sequential, such a configuration is relevant when the performance demand is higher and the sequential stage is lighter than the others. *Configuration 6* has stages 2 and 3 execute in parallel and the first stage is sequential and in *Configuration 7* all stages execute in parallel, which can be relevant when the performance demand is higher and the stages are balanced. It is important to note that *Configuration 7* tends to consume more resources.

The configurations that are suitable vary from application characteristics and the performance demand. In this synthetic application, many other configurations could have been declared and made available for the self-adaptive strategy. However, these 7 are representative enough for evaluating the accuracy and performance of the proposed solution. Additionally, this synthetic application allows flexible customizations of load balancing between the stages. Two application versions were created for evaluating the self-adaptive strategy: one where the stages are balanced and the other that has unbalanced stages. In the balanced version, if we attribute a total computing weight of 6 each stage would have a weight of 2, meaning they are perfectly balanced. With the balanced stages, the optimal configurations are 1 and 7. On the other hand, the

unbalanced version has also a total stages' weight of 6. However, the first stage has a weight of 1, the second weight of 3, and the third stage has a weight of 2. In this case, the major bottleneck is the second stage and if the performance demand is high the third stage can become the second bottleneck.

### 4.2.2 Ferret

Ferret is a stream-parallel benchmark that searches for similarities on data items like audio, images, and video [6,17]. For the evaluation, we modified the original ferret version to a streaming version. This streaming version computes data items at a fixed speed instead of reading the data as fast as possible from the disks, simulating a scenario where the data comes in real-time from the network at a given speed to be computed. The streaming version also covers the instrumentation to collect stream processing metrics like throughput and latency, instead of the execution time. We used the PARSEC native as the input set, which is a representative workload.

Ferret can be modeled with several configurations. In this evaluation, we created 15 alternative configurations for challenging the self-adaptive strategy to find the best ones at run-time considering different scenarios. The four parallel stages from Fig. 1 are represented as user functions a,b,c,d. Moreover, the configurations from 1 to 12 explore possible combinations of sequential (S) and parallel (P) stages, whereas configurations 13, 14 and 15 cover the merging of sequential stages. Merging can be relevant when the stages are unbalanced and the lighter ones can be merged. Importantly, the self-adaptive strategy has a profiling step for characterizing the stages and their workload.

### 4.2.3 Person recognition

The Person Recognition is a stream processing application [9] where we used a customized version that has three functions to detect and verify people in video streams. It receives a video input and applies a denoising step for improving the quality. Then, it detects and marks the faces with a red circle. These faces are compared with the training set of faces. The experiments were run using as input a 30 s video with a resolution of 260 pixels.

In the Person Recognition, we used 5 alternative configurations from reference [24] that cover sequential, parallel, and merged stages. In *Configuration 1* all application functions are merged in a sequential stage (1S.). *Configuration 2* separates the functions into two stages (Pipe-2S.), whereas *Configuration 3* runs with one more stage(Pipe-3S.). Considering that some applications or performance goals are not suitable for sequential stages, *Configuration 4* shows an example of a pipeline with a parallel stage (P.S.1) running all functions, which in FastFlow is a Farm parallel pattern. Considering that functions can be decomposed into multiple parallel stages, *Configuration 5* provides a variation of *Configuration 4* where two parallel stages (P.S.2) are employed, which can be useful for applications that are not embarrassingly parallel.

### 4.3 Self-adaptive strategy characterization

This section characterizes the decision-making process of the self-adaptive strategy. The first results to characterize the solution are from the synthetic application. The proposed solution is compared to the previous one called PDP21 [24]. Figure 4a shows the results of balanced application stages where the defined SLO is to have a throughput (items/s) outcome equal to the IR, where there are two changes in the IR representing fluctuations that can occur at run-time in stream processing.

The PDP21 strategy started trying all configurations. Considering that the SLO was being achieved, the new strategy avoided the unnecessary training in step 2 of the decision strategy (see Sect. 3.2). By Reacting to the IR change around the second 30, the new strategy accurately went on one step to configuration 7 that executes all stages in parallel, inferred by the profiling step that detected balanced stages. By contrast, the PDP21 strategy tested all configurations again, resulting in lower throughput and higher latency for several seconds due to testing suboptimal configurations. Then, after the second 50, the IR dropped, and the executions went back to configuration 1 that sustained the SLO without demanding additional resources.

Figure 4b shows a distinct outcome in a scenario with unbalanced stages. The execution starts with a throughput lower than the IR. Therefore, the new strategy searches for better configurations, which results in shortlisting and entering the trial phases with configurations 3, 5, 6, and 7. The rationale behind such as decision is that the profiling correctly detected the second stage as the bottleneck (Sect. 4.2.1) and shortlisted the configurations where the second stage is parallel as an attempt to overcome the
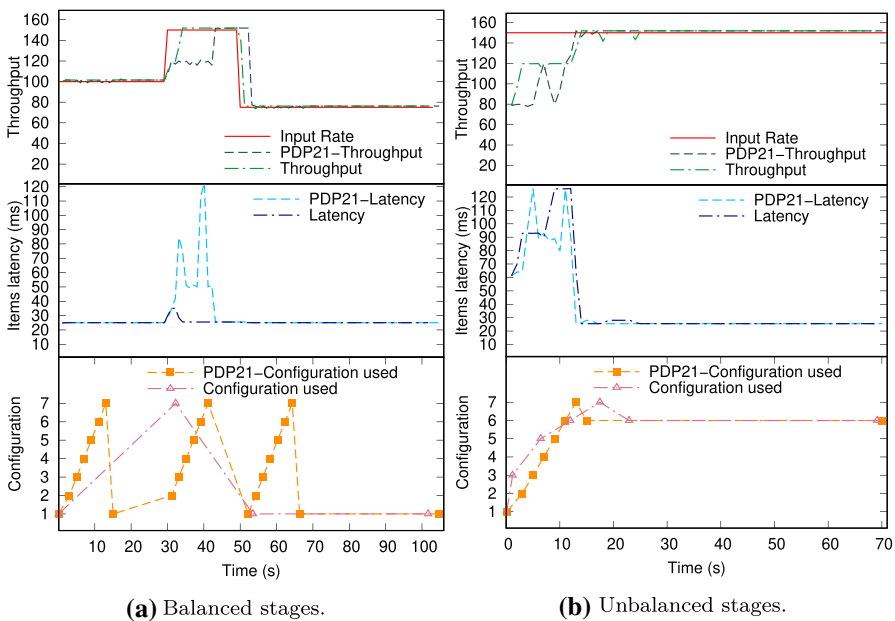


**(a)** Balanced stages.                    **(b)** Unbalanced stages.

**Fig. 4** Characterization with the synthetic application

**(a)** Ferret with IR 10.
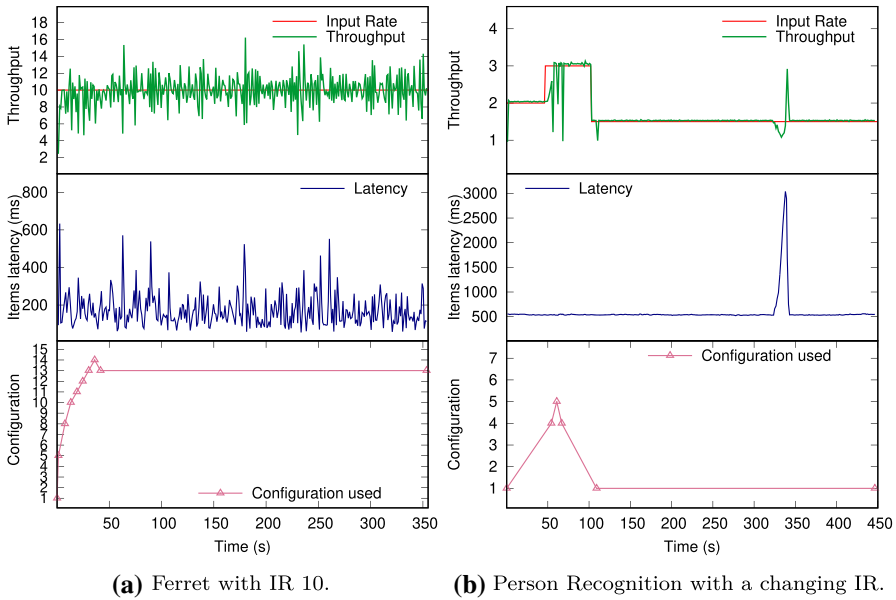
**(b)** Person Recognition with a changing IR.

**Fig. 5** Throughput (items/s) Characterization

bottleneck. Even during the trial phase, it is noticeable that the performance improved in terms of throughput and latency. Then, the strategy stabilized with configuration 6 that provides QoS and demands fewer threads than configuration 7.

The PDP21 strategy had to apply all configurations to find the best one. By contrast, the new strategy inferred the best configuration with fewer steps, which is very relevant for real-world applications [15]. The strategies used different time intervals for testing each configuration, the new strategy uses the default value of five seconds, and PDP21 tests configurations for one second. Although the time interval can be customized for specific application's characteristics, five seconds is expected to be a suitable value for a wide range of applications. Another relevant aspect evinced in Fig. 4 concerns the transitioning model. Notably, the transitions between configurations are smooth without throughput drops or latency glitches.

The new strategy is also characterized with more realistic applications, where we only show results from the new strategy for the sake of visual clarity. Figure 5 provides results with two applications. Figure 5a shows Ferret where is notable that the metrics collected in real-time present fluctuations due to the application's processing characteristics. Importantly, the self-adaptive strategy's profiling step detected the *Rank* stage as the bottleneck and shortlisted configurations where this stage executes in parallel. Then, after the trial phases, it stabilized with configuration 13 that presented a suitable performance, and that consumes fewer resources with the first stages merged.

The results from the Person Recognition application shown in Figure 5b emphasize the accuracy of the decision-making, which chooses the best configuration according to IR changes. In a scenario with SLO violations, configurations 4 and 5 were shortlisted and tried to achieve higher performance. Hence, the strategy applied configuration 4.
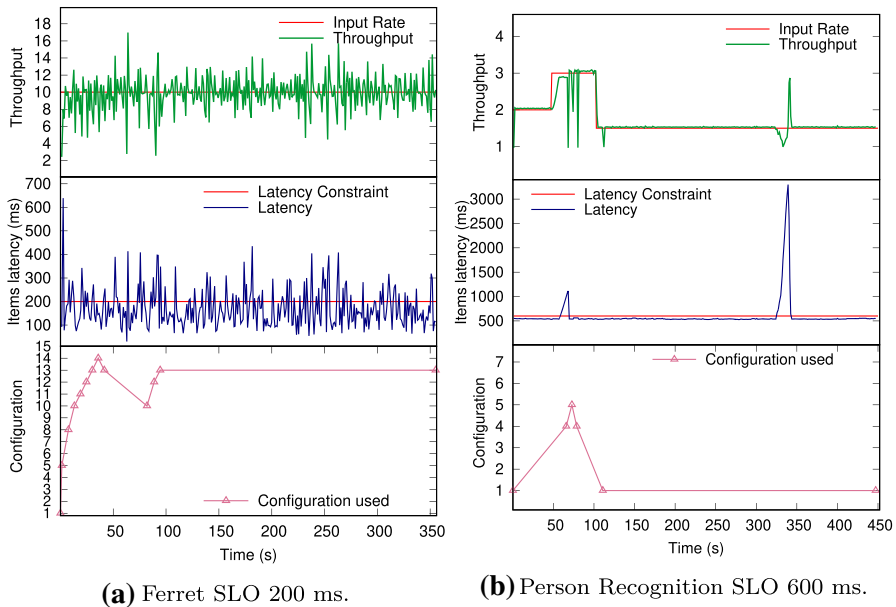
**(a)** Ferret SLO 200 ms.

**(b)** Person Recognition SLO 600 ms.

**Fig. 6** Latency characterization

Under a lower IR, the self-adaptive strategy returned to configuration 1 to increase efficiency by demanding fewer resources. Although the throughput was reduced during some reconfigurations, the transitioning model showed accuracy because there was no application downtime.

Figure 6 shows results from the self-adaptive strategy with latency SLO as a constraint. Figure 6a evinces Ferret with an SLO of 200 ms with fluctuations due to Ferret's characteristics. The strategy stabilizes with configuration 13, overcoming the bottleneck on stage *Rank*. Near 100 s time, a significant application fluctuation increased the latency. Hence, the strategy detected an SLO violation and searched for a better configuration because some change could have occurred. The third pool stage (*Vec*) was detected as an additional bottleneck, where the strategy shortlisted and tried configurations 10, 12, and 13, where the two bottlenecks are executed in parallel. However, the strategy returned to configuration 13 that remained the most suitable configuration.

Figure 6b evinces a latency constraint of 600 ms, where a fluctuating IR varies from 1.5 to 3 FPS. A reconfiguration may be needed when the IR changes because not sustaining the IR increases the buffering and latency. This occurred after the second 50 when the IR increased. The active configuration did not sustain the IR, which increased the number of items buffered and the latency. Hence, the strategy detected the latency violation and self-adapted to configuration 4. After the second 300, there is a fluctuation (also seen in Fig. 6b) that caused the throughput to decrease and the latency to increase. This fluctuation was not long enough for a reconfiguration because the latency SLO was being achieved when the self-adaptive entered the training step. Notably, the transition between configurations occurs without application downtime, which shows that the model's estimation is accurate.

## 4.4 Performance evaluation

In this section, we compare the final performance of the self-adaptive executions to static ones using real-world applications. The results represent an average of 10 runs and we also show the standard deviation, which is difficult to visualize in the figures because it is very low.

Figure 7 shows results from Ferret, where the self-adaptive strategy was able to effectively adapt and find the best configuration (13) for achieving a performance competitive with the best static configurations. The best throughput in FastFlow, the runtime system of the self-adaptive solution, was with configuration 12 where the self-adaptive throughput was 6.3% lower. However, in the latency metric, the self-adaptive was 39.7% better than static FastFlow with configuration 12.

Figure 8 provides results from Person Recognition, where a notable outcome is that the self-adaptive executions have a good performance competitive with the best static scenarios. This is due to the accuracy of the self-adaptive strategy, especially the profiling, trial, and transitioning steps.



**(a)** Throughput with IR 10.          **(b)** Throughput with IR 20.

**Fig. 7** Performance comparison with Ferret. Latency on logarithmic scale



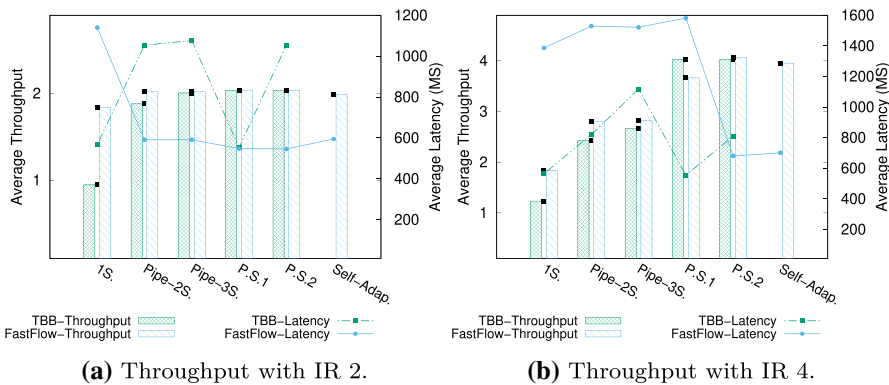**(a)** Throughput with IR 2.          **(b)** Throughput with IR 4.

**Fig. 8** Performance comparison: Person Recognition Application

## 4.5 Results summary

The evaluation provided here shows that our solution for online self-adapting the parallel patterns:

- has effective mechanisms for reconfiguring and maintaining program's executions correctness;
- accurately characterizes the applications for finding bottleneck stages;
- can transparently react to unpredictable fluctuations (e.g., IR, workload) that occur at run-time;
- locates in a few steps, the best configuration according to different SLOs (throughput, latency) and that demands fewer resources.
- provides a transitioning model that is sufficiently accurate as no application downtime neither latency glitches occurred due to reconfigurations (Sect. 4.3);
- has a negligible overhead of instrumentation by the fact of achieving a competitive performance (Sect. 4.4).

## 5 Related work

Considering that modern software applications are executed in dynamic environments and subject to variations at run-time, the literature has different types and levels of adaptation for coping with real-world challenges [13]. For instance, strategies can apply adaptations actions in the environment/resources and at the application software level. When designing self-adaptiveness, control theory [21] is one of the most accepted principles.

In parallel computing, many entities can be adapted at run-time [13], such as batches size [19] for increasing throughput, the cores frequency for reducing energy consumption [5] and dynamically changing the degree of parallelism [4,15,23]. However, these optimizations are arguably not flexible enough for the adaptations that real-world stream processing applications demand. Hence, changing the application structure was proposed as a more powerful and flexible entity to be dynamically adapted [20,24].

Concurrent recompilation has been proposed for reducing application downtime [20]. However, the techniques needed for controlling downtime are intrusive, which can affect the computing (ordering, throughput) and consume additional resources. In practice, we have seen that this approach is hard to generalize to other applications and programming frameworks.

There are approaches employing profiling for guiding the deployment of stream applications [16]. By contrast, we focus on more critical requirements and scenarios with online profiling, where reconfigurations are online without pausing the applications and restarting their executions from scratch.

The study of [11] proposed *Grizzly*, a solution that encompasses adaptive compilation to change the executions at run-time, which is a reaction to changes in applications' data characteristics. Grizzly's decision-making performs only speculative optimizations where in practice the accuracy can vary.

Contrasting with the related approaches, we provide a strategy that uses online profiling and tries only the suitable candidate configurations. This design is an attempt to be generic enough to a wide range of applications and frameworks. In this vein, it is possible to avoid the demand to rerun the applications to apply changes and not demanding intrusive techniques (e.g., input duplication, resource throttling, output smoothing). Recompilation is unnecessary because multiple configurations are created, and the best ones for SLO and efficiency are found and applied at run-time.

Considering that our solution increases the self-adaptation space, we can combine it with other specific entities, e.g., self-adapting the degree of parallelism when high throughput is desirable in configurations with parallel stages. To the best of our knowledge, this is the first approach that provides an accurate decision-making strategy for choosing at run-time the best parallel pattern configuration to be employed and that online self-adapts when it is necessary.

## 6 Conclusions

In this paper, we presented a solution for supporting self-adaptive pattern compositions for stream processing applications. A relevant implication emphasized is the importance of having well-defined building blocks components as composable and nestable objects. The building blocks enable the creation of complex, well-defined structures (e.g., patterns) that we have shown to be possible to self-adapt at run-time. Moreover, the results demonstrated that self-adaptiveness could provide new efficient abstractions and autonomous responsiveness for applications that compute data in real-time.

The components of our solution can be generalized to be used in other scenarios. For instance, the online profiler has the potential to be used for other application classes and workloads. Moreover, the self-adaptive strategy can be generic enough to be customized with other programming frameworks and execution environments. We expect that one could apply the strategy to provide self-adaptations and abstractions for regular parallel applications.

This study is limited in some aspects. We designed the solution to be generic, but mechanisms in the programming frameworks are necessary for achieving self-adaptive pattern compositions. Currently, the FastFlow framework supports such mechanisms. In future works, we intend to support additional applications and workloads in our solution. Moreover, we intend to evaluate our solution for self-adapting other parallel patterns. Considering that we focus on self-adaptation at the application level, a future approach could tackle changes in environment resources by integrating other adaptive entities.

**Availability of data and code:** The code and data analyzed in the current study are available from the corresponding author on reasonable request.

## Declarations

**Conflict of interest** The authors have no conflicts of interest to declare that are relevant to the content of this article.

## References

1. Aldinucci M, Campa S, Danelutto M, Kilpatrick P, Torquati M (2014) Design patterns percolating to parallel programming framework implementation. Int J Parallel Prog 42(6):1012–1031
2. Aldinucci M, Danelutto M, Kilpatrick P, Torquati M (2017) Fastflow: high-level and efficient streaming on multicore. Programming multi-core and many-core computing systems, parallel and distributed computing pp 261–280
3. Danelutto M, Mencagli G, Torquati M, González Vélez H, Kilpatrick P (2020) Algorithmic skeletons and parallel design patterns in mainstream parallel programming. Int J Parallel Prog 49: 1–22
4. De Matteis T, Mencagli G (2016) Keep calm and react with foresight: strategies for low-latency and energy-efficient elastic data stream processing. SIGPLAN Not 51(8):13:1–13:12
5. 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
6. De Sensi D, De Matteis T, Torquati M, Mencagli G, Danelutto M (2017) Bringing parallel patterns out of the corner: the $P^3$ARSEC benchmark suite. ACM Trans Archit Code Optim 14(4):1–26
7. del Rio Astorga D, Dolz MF, Fernández J, García JD (2017) A generic parallel pattern interface for stream and data processing. Concurr Comput 29(24):e4175
8. Griebler D, Danelutto M, Torquati M, Fernandes LG (2017) SPar: a DSL for high-level and productive stream parallelism. Parallel Process Lett 27(01):1740005
9. Griebler D, Hoffmann RB, Danelutto M, Fernandes LG (2017) Higher-level parallelism abstractions for video applications with SPar. IOS Press, Bologna, pp 698–707
10. Griebler D, Vogel A, De Sensi D, Danelutto M, Fernandes LG (2019) Simplifying and implementing service level objectives for stream parallelism. J Supercomput 76:4603–4628
11. Grulich P, Sebastian B, Zeuch S et al Grizzly: efficient stream processing through adaptive query compilation. In: ACM SIGMOD international conference on management of data, pp 2487–2503 (2020)
12. Hellerstein J, Diao Y, Parekh S, Tilbury D (2004) Feedback control of computing systems. Wiley, p 456
13. Herodotou H, Chen Y, Lu J (2020) A survey on automatic parameter tuning for big data processing systems. ACM Comput Surv 53(2):1–37
14. Janjic V, Brown C, Mackenzie K, et al (2016) RPL: a domain-specific language for designing and implementing parallel C++ applications. In: Euromicro international conference on parallel, distributed, and network-based processing, pp 288–295. IEEE

15. Kalavri V, Liagouris J, Hoffmann M et al (2018) Three steps is all you need: fast, accurate, automatic scaling decisions for distributed streaming dataflows, pp 783–798
16. Liu X, Dastjerdi AV, Calheiros RN, Qu C, Buyya R (2017) A stepwise auto-profiling method for performance optimization of streaming applications. ACM Trans Auton Adap 12(4):1–33
17. Lv Q, Josephson W, Wang Z, Charikar M, Li K (2006) Ferret: a toolkit for content-based similarity search of feature-rich data. In: ACM SIGOPS/EuroSys European conference on computer systems, pp 317–330
18. Mencagli G, Torquati M, Cardaci A et al (2021) Windflow: high-speed continuous stream processing with parallel building blocks. IEEE Trans Parallel Distrib Syst 32(11):2748–2763
19. Metzger P, Cole M, Fensch C, Aldinucci M, Bini E (2020) Enforcing deadlines for skeleton-based parallel programming. In: IEEE RTAS, pp 188–199
20. Rajadurai S, Bosboom J, Wong WF, Amarasinghe S (2018) Gloss: Seamless live reconfiguration and reoptimization of stream programs. ACM Not 53(2):98–112
21. Shevtsov S, Berekmeri M, Weyns D, Maggio M (2017) Control-theoretical software adaptation: a systematic literature review. IEEE T Software Eng 44(8):784–810
22. Torquati M (2019) Harnessing parallelism in multi/many-cores with streams and parallel patterns. Computer Science Dept. - University of Pisa, Italy (Ph.D. thesis)
23. Vogel A, Griebler D, Fernandes LG (2021) Providing high-level self-adaptive abstractions for stream parallelism on multicores. Softw Pract Exp 51(6):1194–1217
24. Vogel A, Mencagli G, Griebler D, Danelutto M, Fernandes LG (2021) Towards on-the-fly self-adaptation of stream parallel patterns. In: Euromicro international conference on parallel, distributed and network-based processing. IEEE, Valladolid, pp 89–93
25. Voss M, Asenjo R, Reinders J (2019) Pro TBB: C++ parallel programming with threading building blocks. Apress