

Resource Management and Scheduling in Distributed Stream Processing Systems: A Taxonomy, Review and Future Directions

XUNYUN LIU and RAJKUMAR BUYYA, The University of Melbourne, Australia

Stream processing is an emerging paradigm that continuously handles incoming data in memory, powering latency-critical application such as fraud detection, algorithmic trading, and health surveillance. Though there are a variety of Distributed Stream Processing Systems (DSPSs) that facilitate the development of streaming applications, the problem of resource management and task scheduling is not automatically handled by the DSPS middleware and remains a heavy load for the application providers. The major challenge is that there needs to be a holistic framework of resource management and scheduling to enable self-adaptive, SLA (Service Level Agreement) -aware, and cost-efficient deployment. Particularly, as the advent of cloud computing has supported customised deployment on rented resources, it is of great interest to investigate novel resource management mechanisms that host streaming systems in clouds satisfying the Quality of Service (QoS) while minimising the processing cost. In this paper, we introduce the hierarchical structure of streaming systems, define the scope of the resource management problem, and then present a comprehensive taxonomy covering critical research topics such as resource provisioning, operator parallelisation, and task scheduling. We also review the existing work following the taxonomy structure, facilitating a deeper understanding of their research contribution and method features. Finally, we propose the open issues and future research directions towards realising an automatic, QoS-aware resource management framework.

CCS Concepts: • **Software and its engineering** → *Cloud computing*; • **Networks** → *Cloud computing*; • **Computer systems organization** → *Cloud computing*;

Additional Key Words and Phrases: Resource Management, Stream Processing, Data Stream Management Systems, Task Scheduling

ACM Reference Format:

Xunyun Liu and Rajkumar Buyya. 2018. Resource Management and Scheduling in Distributed Stream Processing Systems: A Taxonomy, Review and Future Directions. *ACM Comput. Surv.* 1, 1 (March 2018), 40 pages. <https://doi.org/0000001.0000001>

1 INTRODUCTION

With the popularisation of the Internet of Things (IoT), the number of intelligent devices used for monitoring, managing and servicing has rapidly increased. These interconnected data sources generate fresh data continuously, forming a large number, or a massive flow of data streams that will eventually overwhelm the traditional data management systems. Meanwhile, the ever-growing data generation has been accompanied by the escalating demands for low-latency data processing. Time-critical applications such as fraud detection [99], algorithmic trading [1] and health surveillance [150] are gaining increasing popularity, all of which rely heavily on the low-latency guarantee to deliver meaningful results. The desire of fast data analysis gives birth to the emergence of stream processing, a new in-memory processing paradigm that allows for the collection, analysis, and visualisation of streaming data with only seconds or milliseconds latencies.

Authors' address: Xunyun Liu; Rajkumar Buyya, The University of Melbourne, The Cloud Computing and Distributed Systems (CLOUDS) Laboratory, School of Computing and Information Systems, Parkville, VIC, 3010, Australia, xunyunl@student.unimelb.edu.au, rbuyya@unimelb.edu.au.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Copyright held by the owner/author(s). Publication rights licensed to Association for Computing Machinery.

0360-0300/2018/3-ART \$15.00

<https://doi.org/0000001.0000001>

Unlike the traditional store-first, process-later batch paradigm, stream processing continuously digests incoming data to provide immediate insights before the value quickly diminishes with time. The incoming data are handled upon arrival, with the results being incrementally updated while the data flow through the system. Presented with only limited resources to handle continuous inputs, stream processing has no random access to the whole stream. Instead, it installs processing logic over time- or buffer-based windows, conducting lightweight and independent computations over recently arriving data. In this way, the strict latency requirement can be met by proper workload balancing and processing parallelisation on a host of distributed resources.

Building a distributed streaming application from scratch is a tedious job and error-prone – developers have to write code for collecting input data, wiring processing logic, and reporting the value of insights with low latency. This is further exaggerated with the burdens of dynamic scaling and failure handling which are common requirements for distributed computation. Over the recent years, various Distributed Stream Processing Systems (DSPSs) have been proposed to facilitate the development of streaming applications. From a structural perspective, a DSPS works as the middleware of a distributed system, offering unified stream management, imperative programming APIs, and a set of streaming primitives to simplify the application implementation. The state-of-the-art distributed DSPSs, such as Apache Storm [143] and Apache Flink [16], further provide transparent fault-tolerance, horizontal scalability, and state management for the upper layer applications, while abstracting away the complexity of coordinating distributed resources. A typical streaming system is thus a three-tier structure comprising user-applications, a DSPS, and the underlying infrastructure.

Though the adoption of DSPSs makes it easier to develop streaming applications, it remains a challenging and labour-intensive task to deploy a streaming system in a distributed environment satisfying certain Quality of Service (QoS) requirements with minimal resource cost. In this paper, the context of deployment problem is mainly derived from the application provider’s perspective, which involves three major research topics: (1) resource provisioning – determining the composition of the processing infrastructure, (2) operator parallelisation – configuring the degree of parallelism for streaming logic, and (3) task scheduling – deciding the placement of streaming tasks on distributed resources. The subtle interplay between these aspects plays a vital role for the deployed system to meet its functional and non-functional design requirements.

Cloud computing has offered a scalable and elastic resource pool to enable a new level of freedom in system deployment. Its customers can unilaterally provision computing capabilities as needed through an automatic-measured, subscription-oriented model, where the monetary cost is billed on a pay-as-you-go basis. The advent of cloud computing also makes it harder to manage resources for streaming systems due to a combination of influencing factors, such as the sensitive application requirements, dynamic workload characteristics, various cloud resource types, and diverse pricing models. Improper resource management and task scheduling directly affect the system performance on clouds. For example, over-provisioning and under-provisioning of resources lead to extra operational cost and Service Level Agreement (SLA) breaches, respectively. Acquiring resources from a suboptimal location introduces additional communication latency and network traffic, and inappropriate parallelisation of operators results in either overload streaming tasks or excessive overhead of context switching. Last but not least, misplacing streaming tasks to the underlying infrastructure leads to inefficient stream routing and resource contention that impair the system stability.

There have been quite a few surveys and taxonomies being conducted in the realm of distributed stream processing systems. Some of them have reviewed the whole stream processing landscape. Cugola et al. [30] wrote a seminal survey on information flow processing that aims to merge the results produced by the data stream processing model and the complex event processing mode. Compared to our work, their survey stands at a higher viewpoint covering data and processing models, the language used to express the processing logic, and the runtime system architecture. Kamburugamuve et al. [75] conducted a survey on distributed stream processing systems with a focus on fault tolerance and comparison among different DSPS implementations. Hirzel et al.’s

work [3] presents a catalog of optimisations to improve the performance of stream processing systems, but only part of the optimisation techniques, such as task scheduling and load balancing, are relevant to the deployment process for not changing the application graph and altering the processing semantics. Dayarathna et al. [33] investigated on system architecture characteristics of various event processing platforms, summarising the advancements made on open research topics such as event ordering, system scalability, event processing languages, and the use of heterogeneous devices. Their survey is on general data stream processing covering both systems and use cases, while our review has a narrower focus on deploying stream processing systems on cloud with SLA-awareness and cost-efficiency.

Some other surveys are conducted in the resource management and scheduling contexts, but each of them has a more specific focus in this area without holistically covering the deployment problem. Zhao et al. [163] surveyed various types of stream processing systems and discussed the default methods for resource management in different DSPSs. Dias de Assunção et al. [37] surveyed the state-of-the-art stream processing engines with a focus on the enabling mechanisms for resource elasticity. Hummer et al. [66] also provided an overview of stream processing and explain the key concepts pertaining to runtime adaptivity and cloud-based elasticity, but SLA-aware resource management is not included in their survey. There are also some surveys that have discussed the patterns and infrastructure to run stream processing systems elastically [53, 54, 60, 118, 127], but they emphasise more on the resource provisioning problem and lack sufficient discussion on operator parallelisation and task scheduling.

As the research in this area advances, there is a long overdue effort to define the scope of the resource management and scheduling problem in the stream processing context, then comprehensively analyse the recent progress and identify the main challenges to achieve better SLA-awareness and cost-efficiency. In this paper, we aim to bridge this gap by **proposing** a taxonomy of resource management and scheduling techniques, surveying existing work with regard to the taxonomy architecture, and discussing the open issues and challenges worth pursuing in the future.

The rest of the paper is organised as follows: we first introduce the hierarchical structure of a distributed streaming system as background, using it to organise the research topics that are involved in the deployment process and covered in this review. Then we present a taxonomy of resource management and scheduling to classify the key properties of existing work. In light of the taxonomy, the surveyed works are mapped into different categories for better comparison of their strengths and weaknesses. A thorough analysis of existing work also shed lights on the promising future directions towards an SLA-aware, cost-efficient and self-adaptive resource management and scheduling framework, which we discuss before closing this article.

2 BACKGROUND

The resource management and task scheduling problem is part of the deployment process to ensure that the pre-defined service level agreements (SLAs) are met and that the resource cost is minimised. This section first introduces the motivation and challenges of the targeted problem, and then outlines how SLA-awareness can be achieved through a self-managing and self-adaptive resource management process.

2.1 Motivation and Issues

To better understand the problem scope, Fig. 1 presents the hierarchical structure of an example streaming system. Sitting on the topmost level is the abstraction of the streaming logic, which in this case consists of four operators standing on incoming data streams. These inter-connected operators constitute a directed acyclic graph (DAG) called *topology*, representing a streaming application that produces incremental results on-the-fly unless being explicitly terminated. Each operator encapsulates certain streaming logic such as data filtering, stream aggregation or function evaluation, while the edges denote the data paths between operators as well as the sequence of operations conducted on the data streams. In most cases, the DAG of operators has been properly defined upon the completion of application development. Once entering the deployment phase, one

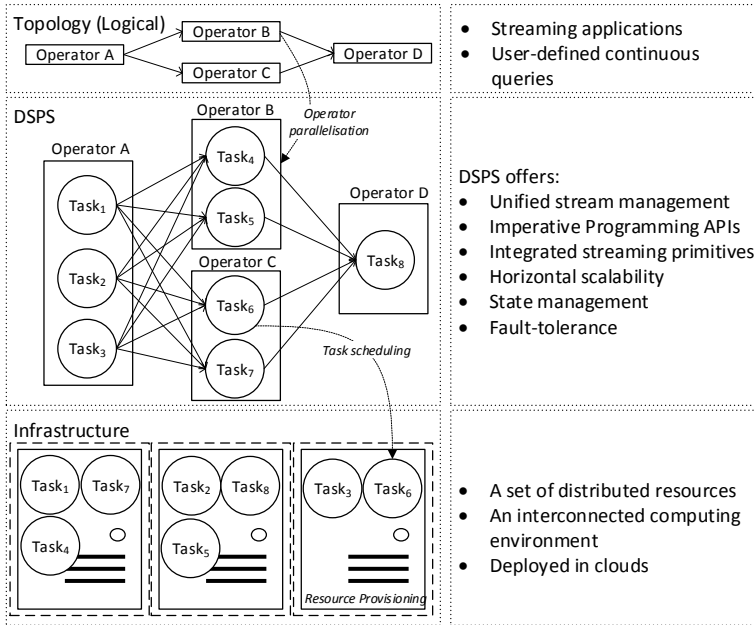


Fig. 1. The sketch of the hierarchical structure of a streaming system

has to decide where and how these streaming logic are executed in a live distributed environment to cater for the continuous and possibly fluctuating workload with SLA-awareness.

A **Distributed Stream Processing System (DSPS)** is positioned in the middle of the system structure, which serves a role similar to the **Data Base Management Systems (DBMSs)** in the conventional data processing context. In general, DSPSs expose a set of imperative programming APIs and streaming primitives to developers, encapsulating low-level implementation details such as stream routing, data serialisation, and buffer management in a unified streaming model. Developers can thus focus on the implementation of streaming logic without having to reinvent the wheels for routine data management. DSPSs also provide abstractions for parallel and distributed computing, allowing applications to exploit horizontal scalability and fault-tolerance without code changes. During the deployment phase, the parallel operators in the topology can scale with a given parallelism degree, generating multiple replicas, known as *tasks*, to execute simultaneously on top of distributed resources. As illustrated in Fig. 1, *Operator B* is parallelised into *Task₄* and *Task₅* due to *operator parallelisation*. Afterwards, a process called *task scheduling* dynamically assigns the streaming tasks to distributed resources, e.g. *Task₆* of *Operator C* is mapped to the computing node at the right end of Fig. 1 for execution. Conveniently, the DSPS guarantees the semantic correctness of parallelisation with built-in mechanisms like automatic stream splitting and tuple tracking. Heinze et al. [54] classified the existing DSPSs into three generations, among which we mainly focus on the third generation which is highly distributed and even applicable to heterogeneous environments such as edge and fog clouds. **Notable implementations falling into this generation include S4 (Simple Scalable Streaming System) [2], Apache Storm, Twitter Heron [84], Apache Flink, Samza [113] and Spark Streaming [160], etc.**

The underlying infrastructure level represents the physical view of a stream processing system, which is an interconnected computing environment created by a process called *resource provisioning*. In this paper, we only consider the Infrastructure-as-a-Service (IaaS) model for provisioning resources in clouds. This model visualises the physical infrastructure as separate service components such as computing, storage and network, where users can deploy their applications with the finest control over the entire software stack, including operating systems, middleware and applications.

There are also streaming services available in the form of the Platform as a Service (PaaS) model and the Software as a Service (SaaS) model. Notable examples include Silicus¹, Google Dataflow² and Microsoft Azure Stream Analytics³. However, the deployment of streaming applications on these services is usually managed by the service owner rather than the application provider, making it harder for the stakeholders to directly manage resources for boosting performance and cost-efficiency.

As shown in Fig. 1, deploying a streaming system can be regarded as a decision and configuration problem that constructs the hierarchical system structure in a distributed environment, where the higher tier is mapped to and hosted on the lower tier to be concrete and runnable. The primary motivation of having a resource management and scheduling framework is to free the application providers from the burden of manual tuning. By applying a collection of profiling, modelling, and decisioning techniques, the framework, which will be discussed in Section 2.2, can be trusted to ensure the deployed system meet its SLA-requirements with minimal resource consumption.

The scope of resource management and scheduling is broken down into main three sub-research topics. We explain each topic by highlighting its peculiar problem domain, as well as discussing the issues and challenges faced by developers to achieve SLA-awareness and cost-efficiency.

Resource Provisioning. Resource provisioning describes the activities to estimate, select and allocate appropriate resources from the service provider to constitute the interconnected stream processing environment.

- Resource estimation: estimating the type and amount of resources needed by the system to meet its performance and cost targets articulated in the SLA. Such estimation can be derived from the analysis of historical data as well as the prediction of future workload, but its accuracy is often affected by the instantaneous, unexpected fluctuation of inputs and system performance variations due to the dynamic nature of data streams.
- Resource adaptation: the real resource demands can fluctuate along with the varying workload, or remain vague and unclear even after the system is brought online. Finding the right point in time to scale in/out and choosing the right adaptation scheme remains a huge challenge. In addition, the profitability of adaptation is affected by a number of factors such as the selected billing model and the geographical distribution of resource pools. Take the latter case as an example, the non-negligible network latency must be taken into consideration when performing system adaptation in a distributed manner [17, 20, 116].

Operator Parallelisation. Operator parallelisation divides a parallel operator into several functionally equivalent replicas, each handling a subset of the whole operator inputs to accelerate data processing.

- Parallelism calculation: this would require accurate profiling of stream workload and probing the processing capability of each task. The details of the infrastructure also matter – the number of cores/threads in a CPU confines the maximum degree of runtime parallelism and the hardware implementation determines the cost of thread scheduling and context switching.
- Parallelism adjustment: over-parallelisation and under-parallelisation can occur at runtime as a result of workload change or resource adaptation. It remains a major challenge to monitor and profile streaming tasks at a fine-grained level in order to reveal the true performance bottleneck of the application. Another challenge is transparent state management during adjustment – stateful operators need to repartition and migrate their states properly among the constituent tasks to make parallelism adjustment transparent to developers.

¹<https://www.silicus.com/iot/services/stream-processing-and-analytics.html>

²<https://cloud.google.com/dataflow/>

³<https://azure.microsoft.com/en-gb/services/stream-analytics/>

- **Balancing data source⁴/sinks⁵**: the parallelism degree of an operator reflects the adequacy of access to the distributed resources. While making parallelisation decisions, the balance between data sources and data sinks needs to be fine-tuned as their performances are correlated due to the producer and consumer communication model in the streaming system. An overly powerful data source may cause severe backlogs in data sinks, whereas an inefficient data source would starve the subsequent operators and encumber the overall throughput [96].

Task Scheduling. Task scheduling dynamically maps streaming tasks to horizontally-scaled resources, such that data streams are partitioned and processed at different locations simultaneously and independently. In this review, we assume that online load balancing of stream routing relies on the DSPS to properly partition data streams among the streaming tasks belonging to the same operator. Also, we assume that there is a state migration mechanism, like the one in [22], to handle the migration of internal task state to the node where the stateful task is being rescheduled.

- **Minimising inter-node communication**: inter-node communication is much costlier than intra-node communication as the former involves time-consuming operations such as message serialisation and network transfer. It is therefore preferable to place communicating tasks on the same node as long as it does not cause resource contention. If the infrastructure consists of geographically distributed resources, it becomes even more prominent to reduce large data transmissions on remote and error-prone data links with limited bandwidth.
- **Mitigating resource contention**: one of the leading causes of performance deterioration is the competition for computational and network resources among collocated tasks. There is a great interest in designing a resource-aware scheduler that makes sure the accrued resource demands of collocated tasks do not exceed the node's capacity.
- **Performance-oriented scheduling**: the scheduling of tasks should be optimised towards the specific application performance targets defined in the SLA, regardless of the interference brought by workload fluctuations, VM performance variations, and the multi-tenancy mechanism at the infrastructure and the DSPS tier.

Task scheduling for stream processing systems is similar to workflow scheduling for batch processing systems — both of them are concerned with the assignment of tasks to the previously provisioned resources. However, they also differ from each other as the objects being scheduled exhibit distinct properties. Streaming tasks possess indefinite lifespan and share a great deal of intercommunication due to the producer and consumer model inherent in stream processing, while batch tasks only exist for a limited time with the dependency of execution hinging on the sequence of completion rather than the continuous and intermediate streams. These differences are naturally reflected in the scheduling process. The former involves real-time decision making considering the dynamic nature of input data, such as the varying intensity and complexity of the data flow, whereas the latter can be done prior to the processing of batch jobs based on the priori knowledge of data, tasks and the execution environment.

It is worth noting that there is an increasing number of DSPSs that support a lambda architecture spanning stream and batch [16, 84, 113]. Due to the strict latency constraint, batches in DSPSs are rather small and processed at small intervals, which gives birth to a new concept called micro-batch. Comparing to the canonical streaming model that handles each new piece of data when it arrives, the micro-batch model divides the stream into small batches of a fixed duration and processes them at individual batch windows, gaining better reliability and exactly-once semantics at the cost of higher latencies. The micro-batch model also plays a vital role in striking a balance between throughput and latency — the two most significant performance indicators in stream processing. For instance, enlarging the batch size in micro-batching may deteriorate end-to-end latencies due to the increased buffer filling time, but it helps improve batch throughput by reducing the frequency of small communication calls. From the scheduling perspective, the micro-batch model introduces

⁴A data source is an operator responsible for injecting data into the stream processing graph

⁵A data sink is a peripheral operator that does not have any outgoing edges and only consume data in the stream processing graph

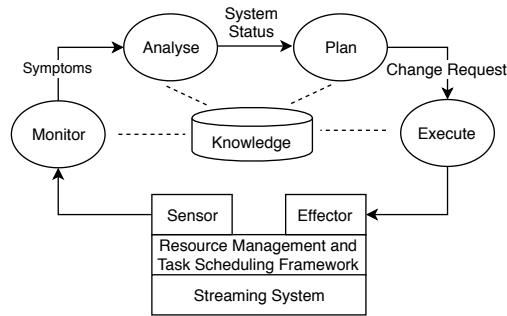


Fig. 2. The Monitor, Analyze, Plan, Execute, and Knowledge (MAPE-K) architecture for resource management and task scheduling

a new level of scheduling called job scheduling – with each batch considered as a job, the job scheduler receives the periodically generated jobs and decides when and how to schedule them in DSPS for execution. There is also a number of scheduling parameters such as batch size, job parallelism and resource shares among jobs to be tuned to cope with the variations in the workload and system conditions. However, due to the page limit, the scope of our review is confined to task scheduling alone which applies to both of the canonical streaming model and the micro-batch model.

2.2 SLA-aware Resource Management

The umbrella term SLA in cloud computing generally refers to a service-based agreement offered by the cloud service provider to its clients, encompassing the minimal standard of many measurable metrics such as system availability, Mean Time Between Failures (MTBF) and service performance, as well as what measures to take for SLA monitoring, SLA violation detection and SLA enforcement. As a paradigm with strict QoS requirements, stream processing further incorporates specific SLAs on end-to-end latency, sustained stream throughput, and processing semantic guarantee to cope with the dynamic nature of input streams and the shared nature of the infrastructure.

During the design phase, the semantic guarantee in SLA is a fundamental design requirement determining whether or not data drops are acceptable, data replays are allowed and idempotent computations are required. During the deployment phase, resource management and task scheduling care more about the SLA compliance with respect to application performance, which is obtained through constantly monitoring the fulfilment of throughput and latency constraints and dynamically making adaptation on infrastructure and DSPS at runtime.

In order to achieve SLA-awareness, the resource management and task scheduling framework needs to be implemented as an autonomic and adaptive system that efficiently fulfils the QoS requirements without human involvement. One representative solution is through a Monitor, Analyze, Plan, Execute, and Knowledge (MAPE-K) architecture that is shown in Fig. 2.

This architectural concept was first introduced by IBM to design autonomic systems with self-* capabilities, such as self-managing, self-healing and self-adapting [77]. In this architecture, SLA-awareness is enabled by the repeated execution of MAPE-K loops. Each loop starts with the monitoring component which is responsible for collecting various types of metrics for SLA monitoring. Data being collected include application metrics, task metrics, and OS (Operating System) metrics. The analysing phase in the MAPE-K loop can be driven by a sophisticated model, or implemented as simple as a boundary checker on the collected metrics to determine whether or not the performance SLA is respected. The planning phase comes into play when it receives the signal from the analysing phase reporting abnormal system status, in which case the passed-in metrics are used as inputs to generate a new plan with viable amendments. Once a new adaption plan is made, the executor in the MAPE-K loop takes the responsibility to put the new plan into effect.

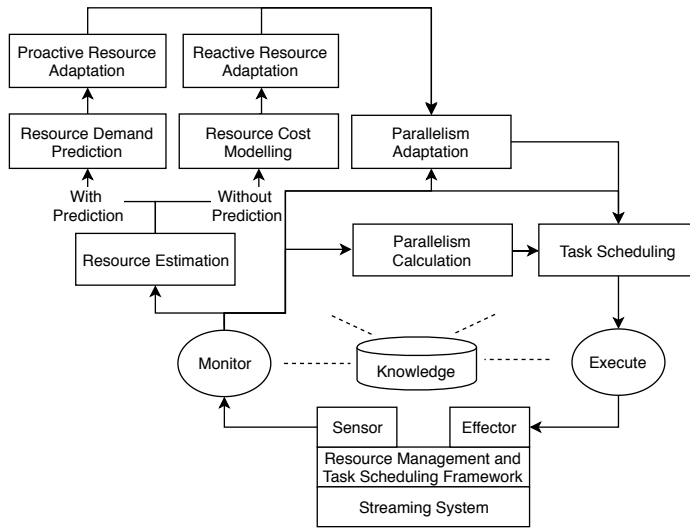


Fig. 3. The intertwined relationship among the surveyed topics in resource management and task scheduling

Meanwhile, the Knowledge component of the MAPE-K loop is an abstract module that represents the data and logic shared among the monitoring, analysing, planning and execution functions.

With the help of the MAPE-K architecture, the rest of the section explains how the subtopics introduced in Section 2.1 are intertwined with each other to achieve SLA-awareness. By replacing the analyse module and plan module with the concrete procedures of resource management and task scheduling, we can adapt Fig. 2 into Fig. 3.

Fig. 3 illustrates that the subtopics mentioned in Section 2.1 can be partially covered in a resource management and task scheduling framework to achieve SLA-awareness, but there is a dependency graph on which topics are bundled to be covered at a time. For instance, it is not compulsory for a framework targeting at SLA-aware task scheduling to alter operator parallelism and the underlying infrastructure, though such a decision would limit its ability to reduce SLA-violations. But things are different the other way around – the process of task scheduling must be invoked by the updated operator parallelism as the number of tasks being scheduled has changed. A similar relationship is also found between resource provisioning and operator parallelisation. Once there are new resources being acquired or existing resources being released, the parallelism levels of streaming operators need to be updated as well to readjust the application scale. This further triggers task scheduling in order to make sure that the streaming application can leverage the benefit of resource adaptation. In reverse, changing the operator parallelism does not necessarily modify the underlying resource configuration, in which case the resource cost of stream processing is not affected. However, through extensive experiments, Lombardi et al. [102] reported that operator parallelism and the number of resources should mutually interact and jointly scaled for the global benefit of the system.

Another point that Fig. 3 captures is that resource provisioning, operator parallelism and task scheduling all rely on the monitoring module of the MAPE-K loop to provide metrics at different levels as inputs, and the executor model to turn outputs – the proposed adaptation into a reality.

3 TAXONOMY

Figure 4 presents a taxonomy of resource management and scheduling in distributed stream processing systems. It divides the scope of the problem into seven different categories in response to the research topics and issues identified in Section 2. Each category of the taxonomy is further extended with several branches to distinguish the specific research targets and method properties for better comparison of similar work. In particular, our taxonomy covers the following aspects:

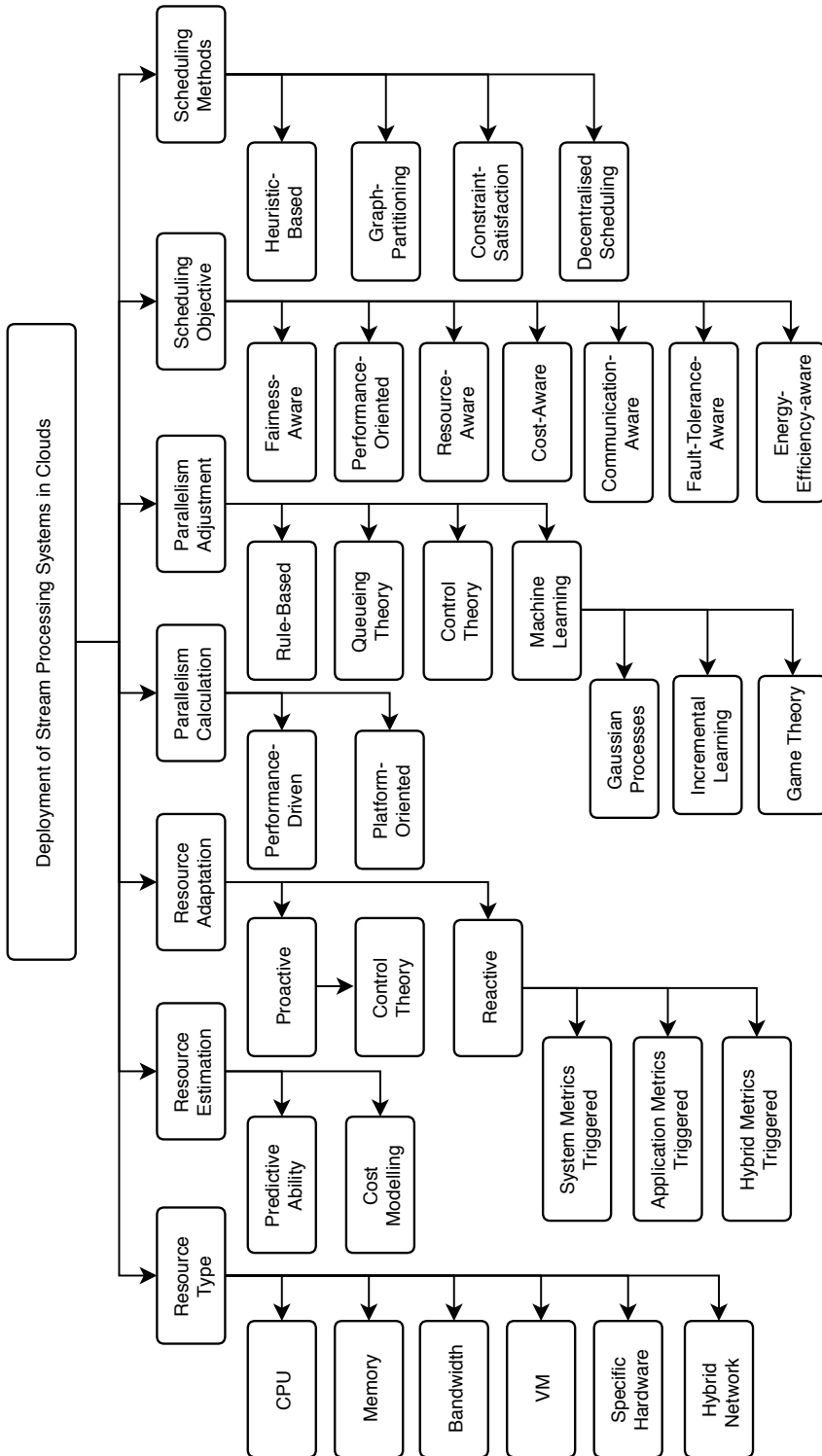


Fig. 4. The taxonomy of resource management and scheduling in distributed stream processing systems

- Resource Type: the various resource types involved in the resource management process to compose the stream processing infrastructure.
- Resource Estimation: the estimation and modelling of resource costs for a streaming system to satisfy its SLA requirements.
- Resource Adaptation: the adaptation of resource allocation to the changes of workload volume and application performance.
- Parallelism Calculation: the profiling and calculation of parallelism degree for the parallel operators in the application topology.
- Parallelism Adjustment: the adaptation of operator parallelism in response to workload variations and internal system changes.
- Scheduling Objective: the various objectives of task scheduling and the rationale behind these objectives to achieve the overall deployment target.
- Scheduling Methods: the various methods used for task scheduling.

Though being divided into different categories, the above-mentioned aspects are only conceptually distinguished in the taxonomy to help us survey the existing work related to resource management and task scheduling from different perspectives. [As discussed in Section 2.2, the activities of resource management and scheduling are often tightly correlated and conducted in a bundle to fulfil a holistic deployment target.](#) Since a surveyed research may stretch across multiple subcategories for completeness, the following sections (Section 4 ~ Section 10) may cover the same work multiple times with a different focus on different aspects of the taxonomy.

4 RESOURCE TYPE

Resource describes any physical or virtual component of limited availability within a computer system. However, depending on the actual context, the same term could contain diverse meanings and refer to various resource types at different levels of abstraction. For deployment in IaaS clouds, resource generally refers to the computing and network facilities that are available to rent through usage-based billing, such as Virtual Machines (VMs), IP addresses, and Virtual Local Area Networks (VLANs). However, when a streaming system is to be deployed in a more hybrid and geographically distributed environment, the concerned resource types also include other infrastructural components such as specific hardware and hybrid networks.

In this section, we identify the various resource types involved in the deployment and resource management process. It is worth noting that the storage resources such as block storage, file or object storage are omitted in our classification due to the rare discussion in the literature. This is credited to the fact that saving stream data to an off-site storage system is often prohibitive, which would block the dynamic data flow and cause unsustainable processing latency.

4.1 Resource Abstractions

Resource abstractions such as CPU, memory, and network bandwidth quantify the resource requirements of a streaming system, regardless of the hardware differences at the infrastructure layer. [In addition, network latency is a major constraining factor that determines if the streaming application can fulfil its low-latency SLA.](#) From the end-users' perspective, the measurement of resource abstractions is intuitive and straightforward. CPU resources can be counted by the number of used CPU cores, with loads measured by Million Instructions Per Second (MIPS) or percentage utilisations. Memory usage is quantified by Megabytes (MB). Network bandwidth consumption is gauged by Megabytes per second (MB/s) or Kilobytes per second (KB/s) and [latency is often measured by milliseconds.](#)

However, ignoring the particularity of the underlying infrastructure also means that the measurement of resource abstractions reflects the general system state rather than yielding actual resource provisioning plans. The results would be susceptible to countless modelling nuances and hardware discrepancies. Instead of being used to directly construct the infrastructure, resource abstractions

are more commonly seen in rule-based approaches (Section 6.2) to approximate the resource cost and shed lights on the direction of adjustments.

4.2 Virtual Machines

Virtual machine (VM) is an emulation of a computer system customisable to meet the specific user needs. In a cloud environment, virtual machine is the most common resource type that encapsulates the computing power and serves as the host of streaming tasks in a distributed environment.

Provisioning VMs from a particular cloud platform is a mixed problem of considering the VM price model, the location of data centres, and the network capacity of inter-connections. The actual VM configurations and placement are determined by the specific computation and communication needs of the streaming system to meet its performance and cost SLA. For the generality of discussion, this survey also includes resource management techniques that originally apply to the on-premise cluster environment, as the proposed resource estimation and adaptation methods would also benefit the VM management in clouds to prevent resource leaks and contentions.

4.3 Specific Hardware

The infrastructure of streaming systems may require specific hardware to boost performance, improve manageability, and deal with particular streaming scenarios. Due to the scarcity of supply and the indispensability of functionality, provisioning of these critical resources is often prioritised over other common computing and network resources in clouds.

Chen et al. [29] proposed a GPU-enabled extension on Apache Storm, exploiting the massively parallel computing power of the Single Instruction Multiple Data (SIMD) architecture to accelerate the processing of stream data. Similarly, Espeland et al. [41] processed distributed real-time multimedia data on GPUs with support for transparent scaling and massive data- and task-parallelism.

FPGA is reconfigurable hardware designed to enable hardware-accelerated computations. The use of FPGA as central data processing elements allows exploiting low-level data and functional parallelism in streaming applications. To facilitate the application of FPGA for stream processing, Auerbach et al. [5] presented a Java-compatible language as well as the associated compiler and run-time system to integrate the streaming paradigm into a mainstream programming environment. Neuendorffer et al. [112] from Xilinx discussed the design tools required for the fast implementation of streaming systems on FPGAs, and Sadoghi et al. [125] investigated how to map multiple streaming applications to FPGA hardware using Hardware Description Language (HDL) code.

[Specific hardware can not only be found at the centralized cloud, but also at different locations of the network to facilitate the timely processing of stream data.](#) In some use cases, deploying the streaming system requires specific sensors to collect input data or monitor the current processing state such as network transmission and power consumption. For instance, data collection sensors are employed by Zhu and Vijayakumar [146, 166] to aggregate stream data from the satellites and environmental monitoring facilities in real time. Kamburugamuve et al. [74] proposed a hybrid platform to connect smart sensors and cloud services, with the data processing logic deployed in the centralised cloud servers to enable new real-time robotics applications such as autonomous robot navigation. Traub et al. [144] optimised communication costs on sensor networks by sharing sensor reads among streaming applications, so that the amount of data transfer is reduced by a combination of data stream sampling and tailoring techniques. Also, power meters such as Watts Up are employed by Shen et al. [133] and Mashayekhy et al. [106] in their streaming systems to get live power readings from the host machines.

4.4 Hybrid Network

Traditionally, streaming systems are deployed in a single cluster or cloud environment as most of the data streams to be processed are collected from web analytic applications. However, there is an ongoing trend that the deployment migrates to a more heterogeneous and geographically distributed setting to process the huge data streams generated by the IoT applications. In this

process, novel network elements and hybrid network structures have been employed to enhance the infrastructure connectivity and create new application paradigms.

The emergence of programmable networking hardware and the expressive data plane programming languages motivate the idea of in-network computation [8, 9, 92, 126]. In-network stream processing delegates part of the computing operations to the network devices such as switches and smart Network Interface Cards (NIC), in order to reduce the increasing data traffic to data centres. However, this brings new challenges to the design of the scheduler to decide what type of computation can be done in-network. The streaming applications are sensitive to the varying network performance yet there are strict end-to-end latency constraints to respect, so the scheduler must make a careful selection of streaming tasks that can work with the limitations of the network architecture and the confined computing power of programmable devices. As well, the scheduler needs to consider maintaining scalability with the rapidly increasing communication cost as the stream data are routinely moved in many-to-many patterns.

Collaborative Fog, Edge, and IoT networks are also gaining popularity in stream processing for the ability to offload a substantial amount of control, computation and management workload to the network gateways close to data sources, thus reducing data transmission and bandwidth consumption. Papageorgiou et al. [114] identified that the low latency requirement is often challenged at the edge of the application topology due to the frequent communication with external IoT entities, so they built new decision modules to place selected tasks on edge devices at runtime using resource descriptors. Hochreiner et al. [61] discussed the distributed deployment of streaming applications over a hybrid cloud, with a threshold-based resource elasticity mechanism to deal with the variation of IoT streams. Cardellini et al. [20] also investigated distributed deployment of streaming systems over a geographically distributed Fog infrastructure, in which they focused on the design and implementation of a QoS-aware and decentralised scheduling policy. A distributed IoT network developed by Ralf et al. [117] tackles aggregation and processing of streaming data in smart city applications, which is capable of enriching input streams with semantic annotations and utilising stream reasoning techniques to allow real-time intelligence with event detection.

Mobile devices have also taken part in the network infrastructure of a streaming system to move computation closer to the data sources. To deploy stream processing application directly on smartphones, Wang et al. [151] proposed a new check-pointing method to mask the simultaneous failure of mobile devices and employed a segmented, UDP-based data transmission method to reduce the cellular network overhead. Similarly, Morales et al. [110] relied on mobile devices to pre-process data streams, and they also proposed a new check-pointing method that is both connectivity-aware and energy-aware. Yang et al. [158] discussed how to enable mobile devices to work in partnership with VMs provisioned in clouds, with a focus on the dynamic partitioning of data streams between mobile devices and data centres to achieve higher throughput and scalability.

On the other hand, High-performance Computing (HPC) network has also been employed in stream processing to enable advanced interconnectivity and better scalability than the conventional Ethernet connection. Recently, Kamburugamuve et al. [76] discussed the use of Infiniband and Intel Omni-Path to improve the performance of stream processing applications, where a new Storm extension is proposed exploiting the native function of high performance interconnects to achieve significantly lower latencies and improved throughputs.

5 RESOURCE ESTIMATION

Based on the information retrieved, recorded or derived from the present and the past system states, resource estimation calculates the minimal amount of resources required by the streaming system to fulfil its SLA. The accuracy of resource estimation determines the cost-efficiency of resource provisioning, which plays a key role in a quick converge to optimal deployment and avoiding over- and under- resource utilisation.

Our taxonomy covers the following characteristics of a resource estimation method:

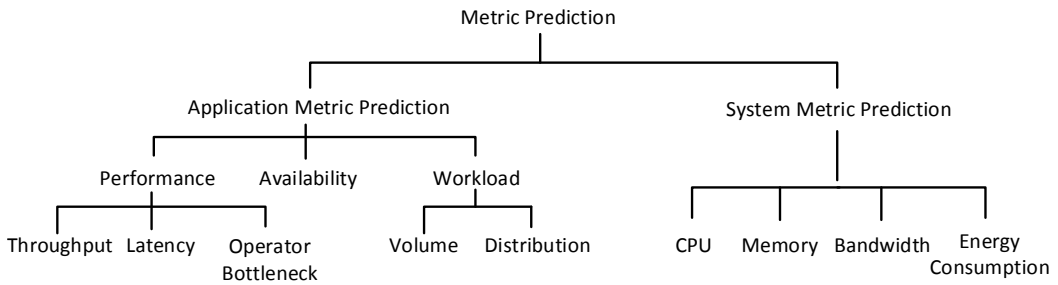


Fig. 5. The classification of predicted metrics used for resource estimation

- **Predictive Ability:** whether the resource estimation method can predict future application and system metrics, such as workload size, resource utilisation, and application performance.
- **Resource cost modelling:** how it models the resource costs based on the predicted or collected metrics, and what criteria in SLA determine the minimal amount of resource requirements.

5.1 Predictive Ability

Prediction of future application and system metrics allows active speculation of future resource demands rather than assuming a constant resource consumption pattern. Fig. 5 illustrates the classification of predicted metrics based on the level of which they are collected from the software stack. Metric of different granularities contain different information and thus contributing to resource estimation in different ways. The prediction of system metrics normally leads to a direct estimation of future resource requirements, which is oblivious to the particularity of the hosted applications; whereas the prediction of application metrics leads to an indirect estimation of resource demands, which requires further resource cost modelling to suggest the minimal resource requirement without violating the SLA requirements.

From the methodology perspective, time series analysis and queuing theory are identified as the two prominent approaches for metric prediction.

5.1.1 Time Series Analysis. A time series is a sequence of data records collected at successive points in time, and time series analysis is an umbrella term that describes a variety of models and methods on time series to find repeating patterns in the historical data. In the context of stream processing, time series analysis can work on the past system resource usages and application metrics, leading to direct and indirect estimation of future resource requirements, respectively.

Direct Resource Estimation by Predicting System Metrics. CloudScale [133] is an elastic resource scaling system built on Xen hypervisor⁶ that directly predicts the short-term resource demands based on the recent history of system metrics. They adopted a hybrid time series analysis approach combining both Fast Fourier Transform (FFT) and a discrete-time Markov chain to balance between high estimation accuracy and low overhead. The light-weight FFT is tried first for fast identification of repeating patterns in the previous time series. If not found, the heavier Markov chain model performs multi-step analysis on the metric history to provide coarse-grained and long-term resource estimations.

The same approach is also seen in the group's previous work [164], with more details revealed on the prediction process. Fast Fourier Transform (FFT) identifies the dominant frequencies of variation in the observed resource-usage time series, followed by a discrete-time Markov chain model that unveils the deeper-hidden patterns through calculating the feature value distribution for the collected resource metrics. The combination of these two methods leads to a fast yet accurate estimation model, provided that there are patterns concealed in the resource usage history.

⁶<https://www.xenproject.org/>

OrientStream [148] is a recent work on dynamic resource provisioning of stream processing systems. It features an online resource prediction module that employs an ensemble regression model on the past system metrics to suggest future resource usages. The prediction process is essentially a weighted vote of four independent regression models, reaping the benefit of reducing the overall Relative Absolute Error (RAE).

Dai et al. [32] presented VM provisioning as a multi-objective optimisation problem, which they solve with an auto-regressive model that learns and predicts the utilisation of each VM as well as the bandwidth consumption between routers. With further consideration on power management, Liu et al. [93] applied deep reinforcement learning over a linear combination of system metrics such as total power consumption, VM latency, and reliability metrics to synthetically predict future system states. Based on the forecast, a hierarchical resource provisioning model is proposed that saves energy consumption without significantly impacting application performance and availability.

Indirect Resource Estimation by Predicting Application Metrics. Time series analysis can also work on the historical application metrics to indirectly suggest the future resource demands with the help of resource cost modelling.

The first category of works predicts the future state of operators. Hidalgo et al. [59] applied a Markov chain model on the workload time series to predict whether an operator's future state would be overloaded, underloaded, or stable. Based on the state predictions, resource cost is modelled by checking the minimal amount of resources needed for the placement of tasks. Kombi et al. [83] adopted a similar method to forecast operator bottlenecks, which circumvents the rigorousness of queuing theory while still being able to estimate resource demands at the operator level.

The second category tries to predict future workload. Balkesen et al. [7] applied exponential smoothing on the periodic observations of the input stream rate to forecast the volume of future workloads. Their rate-forecasting heuristic solves the bin packing problem formulation, suggesting the future resource usages based on the stream distribution and the placement of operators. Analogously, Ishii et al. [69] employed Sequentially Discounting AutoRegression (SDAR) to predict future input rates. They formulated an optimisation problem on resource provisioning and solved it with linear programming to find the minimal resource requirement without violating the application latency SLA. HoseinyFarahabady et al. [63] also predicted the changes in the input traffic with an Auto Regressive Integrated Moving Average (ARIMA) model, which lays the foundation for a resource provisioning algorithm that causes less QoS detriments over all available servers.

Mayer et al. [108] predicted the workload distribution and its parameters with a hybrid approach of distribution moments and maximum likelihood method. The predicted workload distribution feeds into the calculation of operator parallelism and then sheds light on the resource cost by counting the number of processor cores required for task execution. Imai et al. [68] trained a linear regression model on the performance data collected in an experimental environment, in order to predict the maximum sustainable throughput of the streaming application running on a larger number of VMs. Therefore, the cost model is built by directly linking the desired application performance with the number of VMs provisioned in the infrastructure.

5.1.2 Queueing Theory. Queueing theory is a set of mathematical models studying the waiting lines and queues to describe or predict the waiting time and queue lengths. In stream processing, queueing theory is often applied to application performance metrics – especially operator latency – to shed light on the possible data flow bottlenecks.

By modelling the operator as a $G/G/1$ queueing system, De Matteis et al. [34] regarded each operator task as a single server queue where both inter-arrival times and service times have a general distribution. The Kingman's formula is then used to approximate the mean waiting time at the operator level, which sheds light on the runtime adaptation of the number of used cores as well as the CPU frequency. The same modelling and solving technique are also found in a variety of literature [25, 35, 36, 87, 100], which proves that the Kingman's formula is widely accepted for latency modelling because of its accuracy and generality applying to arbitrary distributions of the inter-arrival time and service time.

Differently, HoseinyFarahabady et al. [63] modelled the operator as a $G/G/k$ queue (k is the number of processors for the target operator) and employed Allen-Cunneen approximation to give an upper-bound of the sojourn time experienced by each tuple. Since there is no exact formula known for the $G/G/k$ -model, Allen-Cunneen approximation provides asymptotically exact results under heavy traffic and is particularly suitable for streaming applications with highly-utilised operators. There are also two papers [45, 135] that formulate the operator as a $M/M/k$ system, where M indicates Poisson distribution for arrival and Exponential distribution for service time. Accordingly, Erlang formula is applied to estimate the expected value of the total tuple sojourn time in the application. This is different to the $G/G/1$ and $G/G/k$ modelling at the operator level as the whole application topology is modelled as a Jackson open queueing network, which increases the rigorousness of the queueing model but is capable of providing more accurate latency estimations for the whole application if the model assumptions are met.

5.2 Resource Cost Modelling

With the domain knowledge of stream processing, a resource cost model summarises the various metrics collected at the runtime stack, suggesting an overall estimation of resource requirements for the streaming system to satisfy its particular SLA requirements.

Resource cost modelling is generally influenced by the application logic and the desired deployment target. For example, a filtering operator may show a resource usage pattern linear to its workload volume, while a window operator may exhibit periodical resource requirement peaks as the window slides or executes. For some applications that prefer higher reliability and availability, extra resources are required to provide fault-tolerance or confine the utilisation rate of each node in certain bounds. However, some applications are more sensitive to resource cost and communication overhead, so the deployment of these applications is consolidated to as fewer nodes as possible to reduce resource usages and inter-node communications. Proper resource cost modelling is the key to deal with these variations and suggest the overall resource requirements accordingly.

For the convenience of discussion, our taxonomy categorises different resource cost models based on the intended SLA optimisation, i.e. which SLA requirement is more critical to determine the system resource cost in general.

Minimal Cost Model. This model intends to achieve the targeted performance requirement with minimal resource cost and provisions no spare resources to improve reliability and availability. Bin packing is the most common strategy to model the minimal resource cost based on the compact task placement. Setty et al. [132] used bin-packing formulation to determine the minimal number of VMs needed by the placement of topic-subscriber pairs, with a greedy heuristic to optimise cost while respecting the constraint that the application communication must not exceed the VM bandwidth capacity. Heinze et al. developed FUGU, an elastic data stream processing prototype, to evaluate different scaling policies [57] and optimise the scaling parameters [58]. In both works, the resource requirements are estimated using a bin-packing model solved by a First Fit and Best Fit heuristic. Balkesen et al. [7] applied bin packing to dynamically re-assign data streams to different nodes, resulting in runtime adjustments to the previous round of stream assignment rather than re-optimisation from scratch to balance between result optimality and the overhead of stream redirections. Bin packing is also employed by Liu et al. [94, 95, 97], Xu et al. [156], Nardelli et al. [111] and Ghaderi et al. [49] to suggest minimal resource cost under a certain SLA requirement.

Reliability-oriented Model. This model provisions additional resources for state management and failure recovery. Madsen et al. [104, 105] proposed a Storm extension that replicates the same operator state across different nodes, allowing faster state migration to transparently scale and recover stateful operators. The resource cost is thus calculated by the needs of state management to maintain semantic correctness and fault-tolerance. [Similarly, Castro Fernandez et al. \[23\] designed a set of state management primitives to expose internal operator states to the DSPS for transparent failure handling and scaling.](#) This leads to a resource cost model built on state management with

considerations on the extra computation and communication overhead introduced by failure recovery and periodical state check-pointing.

Contention-aware Model. Model of this type permits certain resource allowances to handle random workload bursts when needed. HoseinyFarahabady et al. [63] proposed a resource cost model that tracks and confines the CPU utilisation level of each node within an accepted range, and a similar approach is also found in Thamsen et al.'s work [141]. In order to limit the memory usage and CPU consumption within a certain bound, Cammert et al. [15] proposed a cost model to estimate resource utilisation of continuous queries based on the stream characteristics such as the average inter-arrival time and the average validity of tuples. The proposed fine-grained cost model is customised to a variety of operator types and streaming logic, making it possible to even quantify the impact of (re)optimisations on query plans.

Load-balancing-oriented Model. This model focuses on the fair utilisation of available resources and is the opposite of the compact task placement which is commonly seen in the minimal cost model. Fischer et al. [42], Eskandari et al. [40] and Jiang et al. [72] regard the operator placement as a graph partitioning problem, so that they explicitly spread the streaming tasks across all available resources at the infrastructure for better load balancing. This cost model is also employed by the round-robin scheduler that is used as default by a variety of DSPSs, which favours even load distribution over participating computing nodes.

Distribution-based Model. Also, depending on the nature of the cost model, the result of resource requirement may be a probability distribution rather than a definitive value. Khoshkbarforousha et al. [80, 81] employed Mixture Density Networks (MDN), a statistical machine learning model combining Gaussian mixture models and feed-forward neural networks, to estimate the whole spectrum of resource usage as probability density functions. Modelling the resource usage as a distribution rather than a single point value captures the possible variances caused by resource contentions and interferences from parallel workloads.

6 RESOURCE ADAPTATION

In this review, we refer to resource adaptation in stream processing specifically as horizontal scaling, i.e. adding or removing VMs within the infrastructure to alter the scale of distributed computation. However, there is also vertical scaling that resizes the existing VMs and adjusts the capacity of existing hardware in terms of CPU, memory, and network resources. Vertical scaling is less preferred for scaling up the streaming system because the maximal scalability is limited by the size of the server, and stop-free vertical scaling has received limited support from the mainstream cloud provider such as Amazon, Microsoft and Google. The consequence of bringing down the whole streaming system for maintenance can be unacceptable in the presence of continuous inputs and strict latency SLA.

In some cases, vertical scaling can be combined with horizontal scaling for right-sizing the VMs for the current load, with a typical use case being workload consolidation [63]. For example, it is profitable to reduce inter-VM latency by consolidating 8 medium VMs at 50% load into 4 medium VMs at 100% load. In addition, vertical scaling can be employed to optimise energy consumption through Dynamic Voltage and Frequency Scaling (DVFS). DVFS is a power management technique that allows processors to dynamically change power states, lowering and raising CPU frequency and voltages on the fly according to the resource demands from virtual machines. It is used by Matteis et al. [34, 36] to explicitly regulate the CPU frequency, by Sun et al. [140] to model the power-to-frequency relationship, and by Shen et al. [133] to turn unused resources into energy savings without affecting application SLA.

In the rest of this section, we categorise horizontal scaling techniques into two major categories based on how they select the proper scaling time. (1) Proactive approaches that adjust resource provisioning according to the prediction of workload pressure and system behaviour in the future

time horizon. (2) Reactive approaches that scale the infrastructure only when necessary as indicated by some threshold breaches or changes of system state.

The choice of proactive or reactive approaches much depends on the predictability of workload pattern and system behaviour. In some cases, the input stream exhibits gradual and repetitive variations in volume and composition, so it is preferable to learn from the history and apply the obtained knowledge to adjust resource provisioning proactively before the application requirement changes. In other cases, the arriving data stream contains random bursts and drastic workload changes with no clear pattern, leaving the prediction of future system state no longer a viable option [11, 142]. Hence, reactive approaches are required to deal with the bursty load on the best effort basis.

6.1 Proactive Adaptation

Proactive adaptation regards the infrastructure tier as a controllable system requiring certain corrective actions from time to time, e.g. acquiring more resources to tackle under-provisioning or relinquishing over-provisioned resources for cost-efficiency. Therefore, there are continuous controlling loops that monitor the various inputs and outputs of resource management and actively suggest optimal adjustments without delay or overshoot.

A typical workflow of a controlling loop is as follows: (1) the resource estimation module predicts the future system state such as the workload arrival rate and the average input processing latency in the prediction horizon⁷. (2) The system model captures the relationship of various QoS variables, assessing the system's capability to maintain the articulated SLA. (3) the control algorithm solves an optimisation problem to find the best resource allocation for the next loop. (4) perform resource adaptation and adjust the operator parallelism accordingly to avoid data skew and load imbalance.

Based on how the optimisation problem is solved, we generally categorise the proactive adaptation methods into two groups. The first one is *loop-wise control*, which regards each prediction horizon as an independent control interval and derives proactive adjustments by applying the predefined scaling rules to the estimation of the next control loop. Methods falling this group are intuitive and straightforward to implement, but they may suffer from the problem of adjusting for short-term benefits while ignoring the long-term future.

To mitigate this, *Model Predictive Control* (MPC) optimises resource provisioning in a receding prediction horizon that consists of multiple control intervals. At each control interval, the controller solves an optimisation problem to obtain the optimal reconfiguration trajectory over the prediction horizon. However, when it comes to execution, only the first element of the optimal reconfiguration trajectory would be employed to steer the resource adaptation, while the whole trajectory is re-evaluated at the beginning of the next control interval to exploit the updated forecast in the shifted prediction horizon. De Matteis et al. [34–36] employed MPC to achieve QoS-aware and energy-efficient resource adaptation, formulating the optimisation problem as a minimisation of QoS cost, resource cost and adaptation cost. The search space of the optimisation problem is described as a tree structure and the Branch & Bound methods (B&B) are employed to prune the search tree and reduce the runtime overhead of MPC in a latency-sensitive environment. Meanwhile, HoseinyFarahabady et al. [62, 63] employed MPC to proactively alleviate the resource contention between collocated applications, in which the optimisation problem is solved by Particle-Swarm Optimization (PSO) with the execution time capped to 1% of the control interval to limit its computational overhead.

6.2 Reactive Adaptation

Based on the metric classification shown in Fig. 5, we also categorise different reactive methods by the nature of the triggering metric.

System Metrics Triggered. The system metrics such as CPU utilisation, memory usage, and bandwidth consumption contain raw information on system performance and resource utilisations, thus reflecting the need for adaptation when some metrics have breached certain thresholds. The

⁷Prediction horizon: the period in which the future values of the interested metrics are predicted

common problem associated with this type of methods is that the system metrics may not faithfully reflect the application performance. For example, a higher CPU utilisation rate does not necessarily mean higher application throughput and lower processing latency. Instead, it may imply that the current resource provisioning is not sufficient for handling the incoming workload. On the other hand, methods falling into this category are versatile and easy to implement for being application-agnostic – the simplest example would be monitoring the CPU utilisation at each host, with an upper and lower bound defined to trigger scaling in and out actions [22, 23, 57, 145]. The memory threshold method is also found in Liu et al.’s work [91].

Application Metrics Triggered. The application metrics include not only the application performance perceived by the end-user, but also internal metrics from the DSPS that include the service time, the arrival rate, and the length of input/output queue for individual operators.

Lohrmann et al. [100] present a reactive scaling strategy that reacts to latency constraint violations with appropriate scaling actions, which minimises the total resource consumption under a varying load scenario. The same approach is also employed in their Nephelē [101] implementation. Xu et al. [157] defined a metric named Effective Throughput Percentage (ETP) for each operator, which captures the state of congestion and estimates the impact of operator output towards the application throughput. The operator with the highest ETP will be given more parallelism and assigned to a new VM for scaling out.

By monitoring the input stream rates and the current processing rates within the DSPS, Cervino et al. [24] detect overload conditions in the operator buffer and then scale the number of used VMs accordingly to maintain the required throughput. Similarly, Vijayakumar et al. [146] defined a derived metric describing the difference between the processing time per data block and the average time interval of receiving one block, so that the adaptation is triggered by the calculated buffer-overflow. Kleiminger et al. [82] monitored the lengths of the input and output queues for stream processors, so that the computation can scale out from an on-premise cluster to clouds when needed. Satzger et al. [128] determined if an operator is overloaded by analysing the length of its incoming message queue, with thresholds hard-coded in the scaling logic to trigger adaptations.

Hybrid Metrics Triggered. Since leveraging system metrics or application metrics alone may not faithfully reflect the actual application performance and resource utilisation, there are some works collecting hybrid metrics to comprehensively trigger reactive resource adaptations. The most common combination is to monitor the operator throughput and the resource utilisation at each host node, in order to deduce the average processing cost per tuple at the operator granularity. Liu et al. [95] applied this method to trigger reactive resource adaptation, so that the overall application throughput can be maintained at a pre-defined level regardless of the initial allocation of resources. In another work of the same group, scale-in is performed when the input load decreases, and so does the resource consumption of each operator [94]. The scale of adaptation is derived from the monitored load difference and a comprehensive metric of per-tuple processing cost.

Apart from stream processing with a DSPS as middleware, there is a trending serverless and event-driven architecture called FaaS (Function as a Service) that utilises a docker-based runtime to scale up or down automatically in response to demand. It provides a programming model to allow developers to write functional logic, which is completely autonomous and independent of the event sources, to be dynamically scheduled and run in response to associated events from external sources. Such built-in elasticity also means that the resource adaptation is managed by the docker runtime internally and the resource cost is billed by the workload activity rather than per hour of VM utilization. The typical examples of this architecture include Apache OpenWhisk and IBM Cloud Functions, which we will discuss more in Section 11 to shed lights on how containers facilitate the resource management of stream processing systems.

7 PARALLELISM CALCULATION

Parallelism calculation answers the question of how many streaming tasks are required for an operator to sustain its assigned workload without causing congestions to the whole application

topology. We have identified two prominent approaches in the literature. The first approach is called performance-driven parallelisation – the resulting parallelism degree is a divisor of the operator input size by the anticipated capacity of each streaming task. The second approach is platform-oriented parallelisation – it first checks the maximum number of parallelism units supported by the provisioned platform, and then distributes them as resources among different operators to ensure that the platform is not over-utilised by an excessive amount of processes and threads.

7.1 Performance-driven Parallelisation

In this approach, direct calculation of operator parallelism hinges on the accurate profiling of both operator inputs and the capacity of each streaming task, the latter of which is defined as the maximum number of tuples that a single task can sustainably handle per time unit [95]. There are direct and indirect methods to measure the volume of inputs for an individual operator. The direct methods install a metric collector at the task entrance that automatically gauges the flow traffic and regularly reports to the calculation logic [70], while the indirect method relies on the producer and consumer model to infer the input volume of a particular operator by examining the selectivity⁸ of its upstream operators [129]. The state-of-the-art DSPSs are now exposing metrics reporting APIs for light-weight stream monitoring and management⁹, so the hurdle of directly measuring operator inputs has been lowered with the abundance of collected metrics.

In addition to measuring operator input, task profiling is another piece of the puzzle to achieve performance-driven parallelisation. There are a bunch of monitoring and sampling techniques that profile the task performance from different perspectives. The most commonly profiled metrics include the average processing latency per tuple [27, 96], the idleness of task execution [72, 152], and the resource usages of a task entity [94, 95]. The relationship between task capacity and the first two metrics is readily established – a task reaches its maximum capacity when fully occupied with tuple processing under the wall clock time. On the other hand, estimating task capacity with the last metric relies on the assumption that this task is hosted by a single thread, which means its peak performance is also limited by the maximum CPU utilisation of a single CPU core.

7.2 Platform-oriented Parallelisation

The rationale of platform-oriented parallelisation is twofold – to avoid over-utilising the available resources with excessive operator parallelism, and to help incorporate some rules of thumb suggested by the DSPS developers to make full use of the parallel processing capability. Take Apache Storm as an example, it is suggested that the operator parallelism is a multiple of the number of machines deployed in the platform, and the parallelism of data source is a factor of the number of partitions of the message queue, as such configuration empirically facilitates load balancing between different hosts [96, 138]. As for Apache Flink, the official training guide suggests that using 1 CPU per slot and setting the operator parallelism as a multiple of the number of slots would help achieve balanced slot sharing¹⁰.

Platform-oriented parallelisation is commonly used in industrial deployment settings as reported by Goetz et al. [50]. Specifically, there is a concept of parallelism unit to describe the parallel processing capability of the platform, which essentially multiplies the number of nodes in the platform by the number of cores available on each node. For instance, there are 160 parallelism units available in a cluster consisting of 10 worker nodes with each incorporating 16 cores. The calculated parallelism units are then regarded as a special type of resources that can be distributed among parallel operators in the topology – the slower the task is in terms of the processing latency, the larger parallelism it gets from the resource pool of parallelism units. They also considered the fact that some tasks may exhibit a higher processing latency because of having intensive

⁸Selectivity: an operator metric that describes the number of data tuples produced as outputs per tuple consumed in inputs.

⁹Apache Storm: http://storm.apache.org/releases/2.0.0-SNAPSHOT/metrics_v2.html Apache Flink: <https://ci.apache.org/projects/flink/flink-docs-stable/monitoring/metrics.html> Apache Samza: <https://samza.apache.org/learn/documentation/0.7.0/container/metrics.html>

¹⁰<https://www.slideshare.net/dataArtisans/apache-flink-training-deployment-operations>

communications, so the number of parallelism units can be enlarged 10 to 100 times depending on the number of I/O bound operators present in the topology. This is to ensure that there are enough streaming tasks for the communication-intensive operator to split the workload and perform I/O operations.

8 PARALLELISM ADJUSTMENT

The direct calculation of operator parallelism may not be feasible in some user cases due to the lack of pilot run or monitoring facilities. Also, the results of calculation are prone to profiling errors that adversely affect the system performance. Therefore, an iterative adjustment process is needed to dynamically adapt the parallelism degree in response to the continuous variations of workload and system performance.

8.1 Rule-based Approaches

Rule-based approaches have attracted extensive research attentions due to the simplicity of implementation and effectiveness of adjustments. The core of the method is made of a collection of scaling rules that define the triggering thresholds as well as the corresponding scaling actions. In most cases, the scaling actions are greedy-based, which favour direct mitigation of the threshold violation and converging to suitable parallelism quickly at the expense of optimality. It also means that the resulting parallelism may be trapped in the local optimum and a proper backtrack mechanism is required to search for the global optimum [6].

Rule-based approaches can be generally classified as either static or dynamic in terms of execution.

Static Single Threshold. A static threshold is pre-defined in the scaling logic to trigger parallelism adjustments in a single direction. For example, the threshold on processing latency is one-sided – when the monitored latency exceeds the SLA requirement, the operator parallelism is increased to amortise the processing workload by adding more streaming tasks to the fleet. Besides, Humayoo et al. [64] assessed the necessity of adjustment with a utility threshold to evaluate if the probability of obtaining positive gain outweighs that to incur a loss. Gulisano et al. [51] defined an upper imbalance threshold to ensure the standard deviation of load distribution is below a pre-defined limit. Though setting a single threshold statically makes it fairly easy to implement the adjustment logic, expert knowledge on application characteristics and the platform specification are still required to properly decide the threshold value and the corresponding scaling actions. Furthermore, methods falling into this category lack the ability to scale reversely nor being self-adaptive as the employed threshold is fixed during the complete runtime of the system.

Static Multiple Thresholds. Multiple static thresholds are set in pairs to maintain the concerned parameters within certain upper and lower bounds. For instance, Fernandez et al. [23] defined two thresholds on the average CPU usages of each node to trigger parallelism adjustment from the perspective of local resource utilisation. This approach is also seen in Veen et al.'s work [145]. Kombi et al. [83] divided the estimated amount of operator input by the estimated capacity of a streaming task, where two performance thresholds are defined delimiting a low and a high activity level to trigger the corresponding scaling action. The major challenge for this type of methods is oscillation, where opposite scaling operations are conducted continuously due to the poorly configured thresholds or overreacting changes [48]. Therefore, a configuration of cooling time is set in practice to conservatively limit the frequency of adjustments and mitigate oscillation.

Dynamic Thresholds. With the knowledge acquired from the evaluation of the previous adjustment results, dynamic thresholds improve the method adaptivity by updating the triggering thresholds and refining the adjustment behaviours at runtime. It also helps mitigate oscillation as the parameters of scaling are dynamically updated with regard to the previous run history. Heinze et al. [57] applied reinforcement learning to reward effective adjustments and punish unnecessary changes caused by inappropriate thresholds. Bilal et al. [10] examined whether a change of parameter value has an

overall positive or negative impact on latency and throughput, where the dynamic thresholds are defined as the best performance monitored in the execution history.

8.2 Queuing Theory

The anticipation of operator congestion using queuing theory is not only useful for the estimation and adaptation of resource provisioning, but also for deciding the relevant parallelism requirement. Mayer et al. [107, 108] built an adaptive data parallelisation middleware that deduces a stationary distribution of the queue length under a certain parallelisation degree, so that the operator parallelism is adjusted accordingly to make sure that the message buffer's limit is not exceeded with a high probability. Liu et al. [96] employed a queuing network to infer the throughput distribution among operators considering their selectivity and communication pattern, based on which the operator parallelism is scaled in batch ensuring that the capability of the data source and data sink is balanced.

A predictive operator latency model is built on queuing theory and employed by Lohrmann et al. [100] to formulate a linear objective function on the minimisation of total parallelism. They applied a gradient descent search to find the optimal degree of parallelism for each operator that reduces resource footprints while enforcing the latency constraints. Similarly, Fu et al. [45] formulated a latency model based on queuing theory to determine the number of nodes that each operator needs to be placed on; however, their approach is dedicated to computationally intensive applications with no regards to the possible communication overhead and network delays. Cardellini et al. [19, 22] searched for the optimal parallelism by jointly considering operator replication and task placement within an integer linear programming formulation, and this process relies on modelling the underlying computing node as an M/M/1 queue to estimate the response time of a particular operator subject to its parallelism, service rate, and incoming load.

8.3 Control Theory

The versatile control theory also applies to the adjustment of operator parallelism. In Section 6.1, we have discussed various MPC-based algorithms that explore the optimal configuration of the target application under ever-changing operational conditions. The parallelism degree of each operator is part of configuration which is updated at the beginning of each control interval [34–36, 62, 63]. In addition, Gedik et al. [46, 48] investigated the profitability of parallelism adjustment with respect to the changes in workload volume and the availability of resources, where a control algorithm is proposed to manage the operator throughputs and congestion with appropriate parallelism. In Li et al.'s work [88], the operator parallelism is controlled by the comparison of congestion degrees¹¹ that are measured on the operator's receiving and sending queue, where the strength of intervention could be tweaked by an adjustment coefficient. Floratou et al. [44] presented a throughput-oriented policy that automatically configures the parallelism degree to ensure satisfactory throughput and alleviate backpressure. Similarly, Stela [157] also relies on monitoring throughput changes to make control decisions – the control algorithm increases the parallelism of the most congested and most influential operator to make full use of the newly added machines during scaling out. In Sun et al.'s work [137, 138], the parallelism degree of each operator is determined in proportion to its computational complexity, which is monitored and measured by the unit of MIPS, i.e. Millions of Instructions Per Second.

8.4 Machine Learning and Game Theory

The adjustment of operator parallelism can also resort to a variety of machine learning techniques. Gaussian processes (GP) is employed by Zacheilas et al. [159] to analyse historical data of workload volume and processing latency, so that the parallelism degree can be proactively adjusted to augment the system's performance. By applying incremental learning techniques to different query

¹¹The congestion degree for a particular operator queue refers to the ratio of the size of the queued messages to the overall queue buffer size.

workloads as training sets, Wang et al. [149] predicted the operator resource usages under several manually supplied candidate configurations. The optimal parallelism is then selected to minimise resource usages while considering the current query requests and stream properties. Game theory is also explored to formulate the elastic parallelism scaling problem as a non-cooperative game, with each operator regarded as an independent agent performing a local control strategy. The operator parallelism is thus determined as the system reaches the agreement of Nash equilibrium [109].

Having introduced a decentralized approach for parallelism adjustment, It is also beneficial to compare it with centralized approaches in terms of flexibility and performance. General speaking, centralized approaches have a single component conducting a streamlined decision-making process, which allows for enhanced controllability over various operators and also becomes susceptible to single point of failures. On the other hand, decentralized approaches break down the adjustment logic into local control strategies that are run by each operator. The result of adjustment may be less efficient but the robustness of decision-making is improved.

9 SCHEDULING OBJECTIVES

Scheduling plays an important role in successful deployment as it determines the resource allocation and communication pattern of streaming tasks over distributed hosts. Hirzel et al.'s survey argues that scheduling trades communication cost against resource utilisation, and that the safety of scheduling hinges on the prerequisite that each host has the right resources enough for all the streaming tasks placed on it [3].

While making scheduling decisions, it is also critical to take the cost of adjustment into consideration as there is a severe impact on the latency of the streaming application while the dataflow is being rescheduled. After all, the stability of the system should be prioritized and the oscillating placement of tasks needs to be eliminated to avoid SLA breaches.

Similar to resource cost modelling in Section 5.2, some scheduling targets have competing requirements and thus trade-offs must be made in deployment practice to satisfy particular application requirements. For example, communication-reduction and load-balancing are two conflicting targets that require task consolidation and task spreading over distributed resources, respectively. So are energy-efficiency and fault-tolerance — the former tends to remove any under-utilised components for energy conservation while the latter purposely introduces redundancies to improve reliability. It is up to the developer to choose which scheduling objective best suits the system's need. In order to evaluate and compare different scheduling policies within the same scope, we have classified various scheduling objectives into six major categories.

9.1 Fairness-aware Scheduling

The meaning of fairness is twofold when it comes to the scheduling of streaming tasks. Firstly, the amount of workload assigned to each node should be fair, avoiding load-unbalance where part of the computing infrastructure is over-utilised while the other part is under-utilised. This is mainly achieved by scheduling at runtime by placing streaming tasks on different hosts if they tend to experience load spikes at the same time [3]. Secondly, the resources allocated to each streaming application should be fair, preventing the multi-tenancy mechanism in the mainstream DSPSs from causing application starvation and resource competition. However, it should be noted that being fair in load distribution and resource allocation does not necessarily guarantee a streaming application can meet its SLA requirements [73].

Fairness-aware scheduling is adopted by many open-source DSPSs as their default scheduling strategy. For example, the default scheduler of Apache Storm assigns streaming tasks to computing nodes in a round-robin fashion, achieving coarse-grained load balancing by placing roughly the same number of tasks to each node. The FAIR scheduler of Spark supports the grouping of jobs into pools and setting different scheduling options (e.g. weight) for each pool, ensuring the fairness of resource assignment at different granularities. For instance, the FAIR scheduler can group the jobs by the pertaining user, giving each user an equal share of resources rather than giving each

job an equal share. Similarly, the default scheduler of Apache Flink endeavours to make sure that the task slots, each of which run one pipeline of parallel tasks, are utilised in a fair manner.

9.2 Performance-oriented Scheduling

Throughput and latency are the two dominant metrics measuring the performance of a streaming application from the end-user's perspective. Maintaining throughput at the required level is of vital importance to the stability of a streaming system. In a streaming environment, the data sources usually work independently and asynchronously with respect to the other parts of the streaming system. So if the processing facility lags behind in sustaining the required throughputs, the message buffer between the data source and the deployment platform will be overwhelmed by the backlogs which eventually lead to the system crash [96]. On the other hand, the importance of reducing processing latency stems from the fact that streaming applications are latency-sensitive in nature.

Performance-oriented scheduling used to be *platform-centric* in a cluster environment, which aims at producing better performance in a fixed deployment platform by optimising the resource utilisation or reducing the network communication of streaming tasks [4, 27, 40, 42, 71, 115, 137, 156]. However, as cloud computing has enabled dynamic resource provisioning during runtime, performance-oriented scheduling has become *SLA-centric* that focuses on meeting the pre-defined performance targets with elastic scaling on resource and operator parallelism [45, 68, 96].

9.3 Resource-aware Scheduling

Resource-aware scheduling matches the resource demands of streaming tasks to the capacity of distributed nodes, so that the total amount of resources required by the fused streaming tasks can be accommodated by the resources of distributed hosts [3]. Being resource-aware offers the opportunity to consume less computing and network resources to achieve the same performance target [94]. Apache Storm, for example, has a built-in resource-aware scheduler that is derived from [115]. In practice, the resource demands and capacity are described as a multi-dimensional vector, with each element representing a particular resource type, such as CPU, memory and bandwidth [95, 115]. The scheduling process is thus finding a mapping of tasks to machines such that the overall resource consumption is minimised and the resource constraints are satisfied. To be more specific, the resource constraints state that the accumulated vector of resource demands requested by the collocated tasks can not exceed the vector of resource availability on that node.

In addition, the need for resource-aware scheduling is driven by the ever-growing use of heterogeneous resources in the streaming infrastructure. The computing nodes could range from energy-constrained mobile devices to powerful virtual machines, which possess different computing powers and connection capabilities. Hence, it is of crucial importance to ensure that the workload assignment does not exceed the node's capacity and the resulting task communications can be sustained by the network facilities connecting to it. Furthermore, the task scheduling on specific hardware such as GPU and FPGA should be optimised accordingly to unlock the potential of the heterogeneous hardware [123, 124].

9.4 Cost-aware Scheduling

Cost-aware scheduling and resource-aware scheduling are strongly related since they all cut back unnecessary resource consumption for cost saving. However, they also differ from each other as the behaviour of VMs, with their startup time and billing intervals, means that reducing resource usage may not reduce the costs. Cost-aware scheduling has an ultimate goal of minimising the overall monetary cost for hosting the streaming system. In the context of stream processing, the cost optimisation problem is easily complicated by the strict latency requirement, the heterogeneity of resource types, and the diversity of billing models. For example, in a computing cloud with heterogeneous resources, the billing schemes for CPUs, GPUs and FPGAs can be vastly different and so are the programming efforts that are required to utilise them [52]. The scheduler needs to be-aware of the infrastructure, knowing the performance characteristics of different computing nodes while conducting different types of computations, and the characteristics of incoming data to

make scheduling plans that reduce the overall costs [124]. Similarly, when the deployment involves multiple geo-distributed data centres, or collaborative Fog, Edge, and IoT networks, the cost of data transmission is non-negligible and must be taken into consideration when making scheduling decisions [28].

9.5 Communication-aware Scheduling

From the perspective of implementation, inter-node communication triggers a cumbersome process involving serialisation, message queueing and network transmission, while intra-node communication can be reduced to passing an object's pointer in memory or expedited by the use of a concurrent programming framework like Disruptor¹². As inter-node data transmission incurs much higher resource consumption and significant network latency, it is preferable to place communicating task pairs on the same node as long as it does not lead to resource contention. This also implies that communication-aware scheduling is a special type of resource-aware scheduling with a focus on minimising inter-node communication [4, 27, 40, 42, 43, 72, 156]. For example, SPADE [47] has an optimising compiler that automatically maps applications to distributed resources with a goal of minimising total inter-node communication, while at the same time exploiting available operator parallelism.

To be communication-aware, the scheduler needs to monitor the task communication pattern as well as the resource usage at each computing node. The communication pattern can be represented by a weighted directed graph of streaming tasks, in which the weights associated with vertices denote the task resource requirement and the weights on edges represent the instantaneous throughput of internal streams or the accumulated volume of data transmission. On the other hand, the deployment infrastructure is also regarded as a weighted directed graph of computing nodes, where the weights on vertices denote the node's resource availability and the weights on edges represent the bandwidth capacity of network connection. Therefore, communication-aware scheduling is to find a proper mapping of these two graphs at runtime in order to minimise the number of messages sent between machines while respecting the constraints on computation and network resources.

9.6 Fault-Tolerant Scheduling

Due to the large size of deployment, faults in a stream processing system are not only considered as exceptions but rather normal events. This implies that fault-tolerance should be made a first-class citizen in the scheduling phase to allow fast and efficient error-handling. In a data streaming system, the consequences of faults can range from a single tuple failure to cascading node crashes [65]. A tuple failure affects the timely delivery of messages, which could be caused by the package discarding on overloaded networks. A node crash, on the other hand, impairs the proper functioning of stream operators that are allocated to this node. In general, we categorise various fault-tolerance techniques into two groups: (1) state management, which allows stateful operators to survive from possible node crashes, and (2) event tracking, which ensures that messages are delivered with regard to the desired semantic. Schedulers that are fault-tolerance-aware can alleviate the overhead of state management, reduce the risk of event replay, and expedite the recovery process by taking the possible failures into consideration during the placement of streaming tasks [86, 131, 139, 147, 162]. For example, the frequency of state check-pointing can be reasonably decreased by being availability-aware [19]: stateful tasks can be scheduled on more reliable computing nodes while stateless tasks that are fail-fast and easy to recover can be assigned to nodes with relatively lower availability. Also, placing communicating tasks in the vicinity and making sure that the bandwidth of network link is not over-utilised can help reduce the risk of message delivery errors [97].

9.7 Energy-Efficient Scheduling

Reducing the total energy consumption is of great interests to the scheduling process [137, 140]. The total energy consumption is unnecessarily increased by the under-utilised computing nodes, so

¹²<https://lmax-exchange.github.io/disruptor/>

it is preferable to perform workload consolidation periodically in order to put the low-load nodes into shut-down or low-power mode [94]. Another critical source of energy consumption is the continuous communication among different streaming tasks. Depending on the distance of data transfer as well as the implementation of the underlying network infrastructure, the actual energy consumption of conveying a tuple over a message channel can vary significantly. This implies that the scheduler should also be aware of energy consumptions when deciding the stream routing, putting a large volume of internal streams on wired and reliable network connections rather than channels that are susceptible to interferences to reduce the possibility of retransmission.

10 SCHEDULING METHODS

The previous section covers the various objectives of scheduling but provides little explanation on how these targets are achieved. In this section, we categorise different scheduling methods into four groups and explain the design and implementation of associated schedulers in details.

10.1 Heuristic-based Scheduling

The scale of the scheduling problem increases exponentially along with the growing application and platform complexity. Since finding the optimal schedule in such a huge solution space is an NP-hard problem, heuristic methods are preferred over exact algorithms to trade off optimality, completeness, and accuracy for speed. Aniello et al. [4] pioneered the dynamically scheduling of streaming tasks to improve application performance at runtime, where a greedy heuristic is applied to minimise inter-node traffic and avoid load imbalances among all the nodes. T-Storm [156] extended their work by allowing hot-swapping of scheduling algorithms and fine-grained control over worker node consolidation. The proposed traffic-aware scheduling algorithm has a greedy-based heuristic in its kernel that keeps trying to assign streaming tasks to available nodes with minimum incremental traffic load. Chatzistergiou et al. [27] also proposed an improved heuristic that utilises the domain-specific group-wise communication pattern between streaming tasks to minimise the communication cost, which guarantees to produce a schedule in linear-time outperforming the existing quadratic-time solutions in practical cases. Similarly, Rizou et al. [120, 121] came up with a task placement heuristic to minimise the network load which is calculated as the bandwidth-delay product of data streams between operators. Sun et al. [140] proposed an energy-efficient heuristic that differentiates the scheduling of critical and non-critical operators to minimise the response time and system fluctuations. R-Storm modelled the scheduling problem as a multi-dimensional Knapsack problem, for which they proposed a heuristic algorithm to put communicating tasks in proximity while ensuring no resource constraints on CPU and memory are violated [115]. The list of heuristic-based schedulers goes on with works done by Cammert et al. [14], Sun et al. [139], Heinze et al. [55, 56] and Jiang et al. [72].

It is also worth mentioning that heuristic can play a complementary role alongside the exact algorithms for better execution efficiency. The SODA scheduler [153] for System S, a proprietary DSPS developed at IBM, uses a local search heuristic as a backup solution to the main approach of mixed-integer optimisation. The heuristic method steps in when the CPLEX-based solution fails or becomes too slow to converge. In addition, meta-heuristic has been employed in the scheduling process to improve method adaptivity. Smirnov et al. [136] investigated the use of genetic algorithms to yield throughput improvement as compared to the greedy heuristics, where the task placement is adapted as an evolutionary process utilising the performance statistics gathered at runtime.

10.2 Graph-Partitioning based Scheduling

As we have discussed in Section 9.5, the scheduling process can be regarded as a graph partitioning problem where the communication graph is divided into smaller components hosted on different computing nodes. The quality of partitioning is often measured by the total amount of inter-partition communications, the degree of load balance across the platform, and the execution time required to work out a partition plan. By assuming the streaming tasks cannot move after their initial placement, Xing et al. [155] employed a static partitioning method to select an operator placement

plan that is resilient enough to withstand different input rate combinations. For dynamic scheduling, Fischer et al. [42] collected the communication behaviour of applications, built the communication graph at runtime, and then set a partitioning objective function in the METIS software to reduce network loads and balance the CPU usage and bandwidth consumption over the platform. Similarly, Khandekar et al. [79] proposed a minimum-ratio cut subroutine to achieve hierarchical partitioning of the operator graph in System S. Eskandari et al. [40] also discussed hierarchical scheduling of streaming tasks with METIS, proposing a two-phase approach that improves on the traditional k -way partitioning method by allowing to dynamically compute the number of computing nodes required in the platform. Ghaderi et al. [49] employed a randomised scheduling algorithm with a theoretically provable guarantee on low-complexity, which enables a smooth trade-off between the cost of approaching the optimal partitioning and the queueing performance. In Li et al.'s work [88], the streaming tasks are firstly partitioned based on the dependency graph of communication, while determining the actual task assignment further involves joint optimisation on the topology structure, inter-node traffic and worker node load-balancing.

The theoretical aspect of graph partitioning in the context of streaming task scheduling has been investigated by Eidenbenz et al. [39]. They proved that optimal partitioning is an NP-hard problem and proposed an approximation algorithm that deterministically achieves a constant-factor approximation under a few assumptions on resource provisioning and processing cost.

10.3 Constraint-Satisfaction based Scheduling

Constraint satisfaction problems (CSPs) regard the entities of interest as set of objects whose state must satisfy a number of constraints or limitations. Thinking the placement of tasks as objects, task scheduling in stream processing can be naturally considered as a constraint satisfaction problem subject to various resource and SLA constraints and requiring efficient search methods to be solved in a reasonable time. When comparing to the heuristic-based scheduling discussed in Section 10.1, constraint-satisfaction based scheduling emphasises more on the result optimality and tends to traverse a large area of the solution space to maximise the objective function.

Cardellini et al. [18, 21] formulated an optimal scheduling problem considering the application and resource heterogeneity. The objective function is to minimise migration costs, and the constraints are modelled as the satisfaction of the application SLA. The problem is then solved by CPLEX, a widely used integer programming toolkit. Jiang et al. [73] also formulated a mixed integer program on scheduling to achieve max-min fairness in resource allocation for multiple streaming applications, where the non-convex constraints are converted to several linear constraints using linearisation and reformulation techniques. Schneider et al. [130] proposed a scheduling algorithm for the ordered streaming runtime to minimise synchronisation, global data and access locks, which allows any thread to execute any operator while maintaining the constraints of tuple order in operator communication. Load-balancing is added as an implicit constraint by Zhang et al. [161] to ensure more task assignment will be assigned to the node with the lowest CPU and memory consumptions. For a similar purpose, Liu et al. [98] proposed a runtime-adaptive scheduler that assigns tasks loads in proportion to the processing capacity of nodes. By dynamically migrating tasks assignment from slow nodes to fast nodes, the latency difference between the fastest and slowest nodes is mitigated. Buddhika et al. [13] formulated a resource-constrained problem on scheduling to reduce interference that adversely impacts the performance of streaming computations. They proposed a proactive scheduling algorithm that accounts for the changes in the stream packet arrivals and cluster resource utilisations, which utilises a new data structure of prediction ring to track the amount of workload expected in a given time window.

Constraint satisfaction problems can also be solved by exhausted search. Li et al. [89] trained a model with Support Vector Regression (SVR) on a collection of monitored features to predict metrics like the average latency of tuple processing and the average size of tuple transfer. The resulting scheduler algorithm is essentially an exhaust search algorithm that traverses the whole solution space to find the optimal schedule with the minimised end-to-end latency.

10.4 Decentralised Scheduling

A decentralised scheduler is not a tangible entity that collects global information from the deployment platform and makes holistic scheduling decisions for the whole streaming system. Instead, it offloads the scheduling logic to the individual streaming operator or computing node, regarding each as an independent agent that collaborates with each other to converge to a feasible scheduling plan. The first prominent benefit of decentralised scheduling is robustness, which eliminates the single point of failure and allows graceful degradation in the presence of computing node crashes – the nodes that are not actively cooperating will be excluded from the scheduling resource pool. The second merit of this design is that it can base the scheduling decision on the accurate prediction of communication latency between different hosts, which is of crucial importance for dealing with streaming systems that are geographically-distributed on Edge and Fog cloud.

Specifically, the Vivaldi algorithm [31] – a decentralised approach that has linear complexity with respect to the number of network locations – is often employed to calculate accurate coordinates of distributed nodes in a latency network. Pietzuch et al. [116] pioneered the use of the Vivaldi algorithm to make continuous optimisation in stream processing scheduling without the global knowledge of the system. In their work, a stream-based overlay network is proposed to map the upper streaming system and the underlying physical network, so that the task placement is determined by searching in a multi-dimensional cost space in a decentralised manner. Cardellini et al. [17] presented a distributed and self-adaptive QoS-aware scheduler based on the Vivaldi algorithm, which can deal with infrastructure with non-negligible latencies. Rizou et al. [122] employed the Vivaldi algorithm to form a continuous latency space, and the proposed scheduler ensures that the QoS guarantee on latency is fulfilled while the network load incurred is reduced.

Repantis et al. [119], on the other hand, designed a set of fully distributed algorithms to discover and evaluate the reusability of data streams and processing components, enabling sharing-aware component composition while being consistent with QoS requirements. Chaturvedi et al. also studied the reusability of distributed streams in the context of Storm to improve resource efficiency, proposing dataflow reuse algorithms that identify the intersection of reusable tasks and streams to collaboratively reuse the outputs of overlapping dataflow [26]. Zhou et al. [165] proposed a decentralised and asynchronous scheduling algorithm that improves load balancing by dynamically migrating operators from overloaded nodes to lightly loaded ones.

Unless otherwise stated, the schedulers surveyed in the other subsections are centralised designed, which are often collocated on the master node of the deployment platform for the convenience of metric collection and scheduling coordination.

At the end of our review, we present a tabular comparison of key works regarding resource management and scheduling in distributed streaming systems. It summarises the state-of-the-art in the area on what methods the researchers are using to continuously adapt the system deployment, accommodating load fluctuations by provisioning more resources and adjusting the application scale. It is the general understanding of the community that such adaptation should satisfy the SASO properties – stability, accuracy, short settling time and no overshoot. Scaling actions can be executed either in a reactive or proactive fashion, as long as the system satisfies the performance expectation and works in a correct configuration with the minimal amount of resources being wasted. For proactive approaches, there needs to be a prediction module that forecasts variations in the workload and periodically checks if the resource configuration and application parallelism need to be scaled-in/out. In this process, stability is achieved by introducing a parameter that controls the interval of assessment, which helps mitigate oscillations, while accuracy and short settling time are obtained by some sort of resource cost modelling that computes the expected resource consumption with given input load and configuration. For reactive approaches, a monitoring system at task level is essential that collects fine-grained information on resource usage at runtime to determine whether and how the system should be scaled. And a resource cost model can also help adapting the system to the workload in a single reconfiguration that guarantees accuracy and immediacy.

Table 1. A Review and comparison of key works regarding resource management and scheduling in distributed streaming systems

Work	Resource Type	Prediction Method	Cost Modelling	Resource Adaptation	Parallelism Calculation	Parallelism Adjustment	Scheduling Objective	Scheduling Methods
De Matteis et al. [36]	CPU	Time series analysis	Minimal cost	Proactive	–	MPC	–	–
De Matteis et al. [35]	CPU	Time series analysis	Minimal cost	Proactive	–	MPC	–	–
HoseinyFarahabady et al. [63]	CPU, Mem	Time series analysis	Contention-aware	Proactive	–	MPC	Resource-aware	Control theory
Imai et al. [68]	VM	Time series analysis	Minimal cost	Proactive	Platform-ori.	–	Performance-oriented	Heuristic
Cardellini et al. [22]	VM	–	Reliability-aware	Reactive	–	Queueing theory	Communication-aware	Heuristic
Xu et al. [157]	VM	–	Minimal cost	Reactive	–	Control theory	Performance-oriented	Heuristic
Khoshkbarforoushha et al. [81]	CPU	Time series analysis	Distribution-based	–	–	–	–	–
Wang et al. [148]	CPU, Mem	Ensemble regression	Minimal cost	Proactive	–	Rule-based	–	–
Thamsen et al. [141]	CPU, Mem	Time series analysis	Contention-aware	Proactive	–	Rule-based	–	–
De Matteis et al. [34]	CPU	Time series analysis	Minimal cost	Proactive	–	MPC	–	–
Cardellini et al. [19]	CPU, Mem	–	Reliability-aware	Reactive	–	Queueing theory	Communication-aware	Heuristic
Kombi et al. [83]	CPU, Mem	Time series analysis	Minimal cost	Proactive	–	Rule-based	Resource-aware	Heuristic
Hidalgo et al. [59]	CPU, Mem	Markov chain	Minimal cost	Proactive	–	MPC	Fairness-aware	Round-robin
Shieh et al. [134]	CPU	–	Minimal cost	Reactive	–	Rule-based	Fairness-aware	Round-robin
HoseinyFarahabady et al. [62]	CPU, Mem	ARIMA	Contention-aware	Proactive	–	MPC	Performance-oriented	MPC
Mencagli et al. [109]	VM	–	Minimal cost	Reactive	–	Machine learning	–	–
Smirnov et al. [136]	VM	–	Minimal cost	Reactive	–	–	Resource-aware	Heuristics
Jiang et al. [73]	VM	–	Load-balancing	Reactive	–	Rule-based	Fairness-aware	CSP-based
Sun et al. [139]	CPU, Mem	–	Reliability-aware	Reactive	Performance-ori.	Control theory	Fault-tolerant	Heuristics
Shukla et al. [135]	VM	Queueing theory	Minimal Cost	Proactive	–	Control theory	Resource-aware	Heuristics
Cardellini et al. [21]	CPU, Mem	–	Reliability-aware	Reactive	–	Queueing theory	Communication-aware	Heuristic
Buddhika et al. [13]	CPU, Mem, BW	Time series analysis	Contention-aware	Proactive	–	–	Performance-oriented	CSP-based
Li et al. [88]	CPU, Mem,	–	Minimal cost	Reactive	–	–	Fairness-aware	Graph-based
Schneider et al. [130]	CPU	–	Contention-aware	Reactive	–	Rule-based	Resource-aware	CSP-based
Liu et al. [98]	CPU,	–	Load-balancing	Reactive	–	Rule-based	Fairness-aware	CSP-based
Ghaderi et al. [49]	VM	–	Minimal cost	Reactive	–	–	Resource-aware	Graph-based
Zhang et al. [161]	CPU, Mem	–	Load-balancing	Reactive	–	–	Communication-aware	CSP-based
Li et al. [89]	VM	SVR	Minimal cost	Proactive	–	–	Performance-oriented	CSP-based
Eskandari et al. [40]	VM	–	Load-balancing	Reactive	–	–	Performance-oriented	Graph-based
Sun et al. [137]	CPU, Mem	Time series analysis	Reliability-aware	Proactive	Performance-ori.	–	Fault-tolerant	Heuristics
Eidenbenz et al. [39]	VM	–	Cost minimal	Reactive	–	–	Communication-aware	Graph-based
Lohrmann et al. [100]	CPU	Queueing theory	Load-balancing	Reactive	–	Queueing theory	Fairness-aware	Round-robin
Heinze et al. [58]	CPU	Queueing theory	Minimal cost	Reactive	–	Queueing theory	–	–
Lin et al. [91]	BW	–	Minimal cost	Reactive	–	Rule-based	–	–
Veen et al. [145]	VM	–	Minimal cost	Reactive	–	Rule-based	Fairness-aware	Round-robin
Heinze et al. [57]	VM	–	Minimal cost	Reactive	–	Rule-based	Resource-aware	Heuristics
Madsen et al. [104]	VM	–	Reliability-aware	Reactive	–	Rule-based	Fault-tolerant	Heuristics
Setty et al. [132]	VM	–	Minimal cost	Reactive	–	–	Fairness-aware	Heuristics

11 GAP ANALYSIS AND FUTURE DIRECTIONS

Although many research efforts have investigated the resource management and scheduling in distributed streaming systems, there exist theoretical and technical gaps to the prospect of an SLA-aware and cost-efficient framework that relieves the deployment burden for application providers. In this section, we discuss the identified gaps and shed light on the future directions on this front.

11.1 Fine-Grained Profiling

Accurate profiling of application and system metrics plays an important role in the decision-making process as they reflect the current state of the streaming system and indicate whether the desired SLA requirements have been satisfied. However, most of the existing work based their deployment decisions on coarse-grained metrics such as application throughput, end-to-end latency, operator capacity and the volume of internal streams. These metrics collected at the operator or application level are too general to reveal the actual bottleneck of the data stream, so that the amendments can only be made on a best-effort basis with little guarantee on the adjustment effects. In order to capture the real culprit that throttles the application performance, a fine-grained profiling mechanism is required to fulfil the following expectations. (1) It should be installed at the task level to obtain fine-grained information such as the lengths of input/output queue, the task capacity on different infrastructure, and the average resource cost for processing a single tuple. (2) the application metric collected from the DSPS tier should be cross-validated with the system metrics to identify the probable cause and the severity of the processing bottleneck, allowing accurate amendments to be made in the next adjustment cycle. (3) proper sampling and quantisation techniques should be employed to reduce the profiling overhead while providing strong enough guarantee on result accuracy.

11.2 Straggler Mitigation

A straggler is a slow running entity that adversely impacts the performance of the whole streaming system. It could be a streaming task enduring severe resource contention or data skew, or a computing node that is over-utilised or affected by the performance variation of the host cloud. In either case, the local performance degradation caused by the straggler will soon propagate throughout the topology structure due to the the producer and consumer communication model. The first path of propagation is through the operator DAG — with a straggler, the upstream operator will be throttled by the accumulated backlogs, and the downstream operators will stagnate without receiving sufficient inputs. The second path of propagation is through the performance correlation of sibling tasks belonging to the same operator. If one of these tasks becomes a straggler and performs significantly worse than the others, the logic of tuple emitting will reduce the volume of data stream sent to the other sibling tasks in order to not overwhelm the straggler. This could lead to under-utilisation on other nodes as the healthy sibling tasks could have been placed in different places processing more inputs.

The straggler mitigation techniques have been initially studied in batch processing systems and then ported to most stream processing systems with micro-batch paradigm. Spark Streaming, for example, has a built-in speculative straggler mitigation technique applicable to various workloads, regardless of being either CPU, disk, or network throttled. This is made possible by having speculative backup copies of slow tasks run in neighbouring nodes. Through extensive evaluation, Khan et al. [78] suggest that using mean/standard deviation instead of median for straggler detection, and that using a confidence level to decide if a task can be executed on a node with a history of abnormal behaviours rather than blacklisting that node entirely.

However, there still lacks enough research attention on detecting and mitigating stragglers for canonical stream processing systems, partially because the one-tuple-at-a-time processing model is more dynamic and less trackable. A straggler mitigation mechanism needs to quickly identify the root cause of the performance deterioration and cuts the chain of propagation with active intervention. A straggling computing node can be detected by its soaring resource usages and

slow response time, while a straggler streaming task is revealed by the extended average tuple processing time or the sudden rise in resource consumption. In the context of resource management and scheduling, the possible measures to mitigate stragglers include provisioning new resources, adding more parallelism, or rescheduling the straggler on a different node.

11.3 Transparent State Management

An integrated state management system consists of two parts: (1) State elasticity, which allows dynamically scaling up and down the operator parallelism with a state repartitioning and migration mechanism, supporting the relocation of the operator internal state and providing a guarantee on the semantic correctness during the scaling process. (2) State persistence, which backups the computational states to persistent storage or a different node in order to mask the loss of states caused by JVM or node crashes. There are some preliminary efforts from both academia and industry towards achieving transparent state management [22, 23, 97, 104]. ChronoStream [154], for example, treats the internal state as a first-class citizen and provides state elasticity to cope with workload fluctuation and dynamic resource allocation. However, significant gaps still present in the following aspects. First, there is limited support for the diverse representation of operator state. In most existing state management frameworks, the abstraction and presentation of operator states are limited to key-value mapping for the ease of implementation. But it is possible that computational states exist in other forms such as graphs, hashes and trees that can hardly be indexed by certain keys. One promising research direction would be supporting arbitrary data structure for operator state representation while keeping the repartitioning and migration process entirely transparent to the end-users. The second gap is to reduce the excessive overhead of state migration, which could be overwhelming if the adaptation of resource provisioning, operator parallelism, and task scheduling have not considered the current state placement. Particularly, there is little research on gradual, stepwise task scheduling that eventually converges to the state satisfying the SLA requirements without incurring too much state migration overhead over a short adjustment period. In contrast, most scheduling algorithms in existence determine a new task mapping from scratch by re-applying the scheduling heuristic, re-invoking a graph partitioning algorithm, or re-conducting an exhausting search in the solution space.

11.4 Resource-Availability-aware Scheduling

The existing schedulers have often falsely assumed that, once provisioned, the same amount of resources will be offered to the streaming system throughout its standing lifecycle. Therefore, few of them has considered the fluctuation of node resource availability and what implication it might have for the performance of the streaming system.

In fact, it is common and inevitable to experience fluctuation of resource availability in a distributed cloud environment thanks to two major contributing factors. (1) Multitenancy: multiple tenants of a shared platform may experience performance interference as they compete for limited resources, despite mechanisms like virtualisation and cgroups have provided a certain level of isolation for resource allocation. The temporal and spatial performance variations on Amazon EC2, as reported by Kumbhare et al. [85], can be as severe as 23% of VMs having a normalized core performance worse than 80% of the expectation. (2) Background activities: unexpected background events, such as scheduled system backup, security update, and initialisation of another collocated application could take up a portion of resources that were previously made available to the streaming system. Having a scheduler that is aware of node resource availability can help the streaming system avoid resource contentions.

Resource-availability-aware scheduling is particularly useful when there are no further spare resources for the system to scale out due to the limitation of budget or other performance constraints. In that case, we are interested in changing the mapping of tasks to underlying resources so that the local resource shortage can be amortised over the whole platform. For instance, Imai et al. [67] and Buddhika et al. [13] discussed how to optimize the usage of the available resources through remapping of tasks without expanding the resource pool. The basic idea is that, if tasks of different

operators in the topology process less workload accordingly, their resource consumption is expected to be reduced proportionally. So there is an increasing possibility to find a new task mapping that satisfies the updated resource allocation constraints affected by the fluctuation of availability. Such informed scheduling decision will allow the application performance to degrade gracefully without causing straggler problems discussed in Section 11.2.

11.5 Heterogeneity-aware Scheduling

As discussed in Section 4, streaming systems are incorporating an increasing number of heterogeneous hardware elements and network facilities to improve application performance and enhance system manageability. A typical use case of streaming processing in IoT applications exhibits the following dataflow: data collection activities take place at the edges of a network, data preprocessing and aggregation are performed over the geographically-distributed network, and the heaviest logic for processing and analysis are hosted on cloud servers. The key to maintaining low-latency is to place certain stream processing elements on micro data centres (also known as Cloudlets) which are closer to the data sources compared to centralised cloud data centres.

In this process, the coordination of different types of resources imposes new challenges to the design and implementation of a heterogeneity-aware scheduler. The first challenge is that it would be impossible to describe the computing capability and the network capacity of heterogeneous resources using a unified measure. For example, we can hardly assert that a VM provisioned in clouds is always ten times faster than a mobile device regardless of the workload characteristics. Neither can we be certain about the exact bandwidth difference between a LAN and a wireless connection across regions when routing internal streams. Secondly, since the behaviour of scheduling strongly depends on its environment, the heterogeneity-aware scheduler need to integrate multiple scheduling mechanisms to fulfil stringent performance requirements under environmental changes [103]. For example, one scheduling mechanism may make dynamic use of resources at the edge of the network to meet the low latency demand, while the other one may favour cloud resources for high stability and low control message overhead in adapting the task placement. It is a challenge to find a seamless way switching among different scheduling mechanisms.

The anticipated heterogeneity-aware scheduler also needs to incorporate a customised resource model for each component to track its current resource usage and availability individually. This means that the profiling method should work at the fine-grained task level producing results that are not necessarily shared across the whole deployment platform. A gossip-like negotiation protocol is a promising solution for different components to converge to a scheduling plan that best suits their current situation.

11.6 Energy-Efficient Scheduling

Apart from reducing the total energy consumption through active workload consolidation, it is also of great interests to cut back the proportion of brown energy consumption with task scheduling.

Over the last several years, the energy supply of the infrastructure of streaming systems has been enriched by the green power generated from renewable sources such as sun, wind, water and biomass waste. Energy-efficient scheduling intends to reduce the carbon emission and other negative impacts on our environment by scheduling computational-intensive tasks on nodes driven by green power, as well as allocating a large chunk of communication on links powered by green energy. To do this, the scheduler needs to exploit suitable forecast mechanisms to predict the supply of renewable energy in an online fashion, as renewable energy can be intermittent and much more variable than conventional energy from the grid. The scheduler is then committed to produce a task mapping that satisfies the energy supply constraints while trying to maximise the use of green energy. If the DSPS adopts a lambda architecture that span stream and batch processing, energy-efficient scheduling can postpone the execution of batch jobs, if their deadline permits, until there is enough supply of green energy.

It is also common that saving energy on computation and communication are two conflicting targets that cannot be achieved at the same time through the scheduling of streaming tasks. So a

theoretical or empirical model on energy consumption is required to evaluate and compare different scheduling plans, ensuring the overall optimal in the reduction of brown energy consumption.

11.7 Cost Efficiency with Different Pricing Models

The monetary cost of resource usages in clouds largely depends on the actual pricing and billing model chosen by the users. Apart from the on-demand pricing model that has been intensively studied in the literature, a variety of alternative pricing models are also offered by mainstream cloud service provider like Amazon, Google and Microsoft to help users tailor their choices on resource provisioning and reduce the operational cost. To start with, reserved instances with a fixed-term contract are much cheaper than the on-demand ones, which makes them a good fit to host the baseline workload while leaving the on-demand instances for scaling out when needed. Also, the bidding price model can lower the cost of resource usage significantly as these instances are much cheaper for being hosted on the spare compute capacity in the cloud. However, a streaming system using price-bidding instances needs to handle interruptions in infrastructure under a fairly short notice, which imposes great challenges for the latency-sensitive system to adapt task placement and migrate the associated computational state accordingly. A comprehensive resource provisioning and task scheduling model combining the use of on-demand, reserved, and price-bidding resources is a promising research topic that would be welcomed by industry users.

11.8 Container-based Deployment

Containerisation of clouds allows the services and applications to adapt efficiently and operate at an unprecedented scale. Containers offer a logical packaging mechanism that decouples the applications from the environment in which they actually run, so there is a clean separation of concerns by clearly differentiating the procedures of application development and deployment. The ability of containers to run virtually anywhere and the isolation of the CPU, memory, storage, and network resources at the OS-level make it profitable to host streaming applications that are dynamic in nature [12].

However, resource management and scheduling in streaming systems over containers would require an overhaul in the design and implementation of existing DSPSs. The most prominent challenge is transparent state management over the container cloud that is initially designed to host state-less micro-services. The stateful streaming tasks may have to store their computational state externally which could raise new concerns on the performance of state access. [Function as a service \(FaaS\) frameworks, for example, rely on the use of distributed key/value stores for state management. Apache OpenWhisk uses consul¹³ as a hierarchical key/value store that is accessible by every component of the system for various purposes such as dynamic configuration, feature coordination, leader election, etc.](#) Besides, the flexibility of arbitrary placement and dynamic scaling of containers makes it hard to keep track of the destination of each internal stream, so that the tuple emitting logic needs to be revised to make sure that the provisioned containers are coordinated properly in sending and receiving messages.

11.9 Integration of Different DSPSs

The diverse user requirements may require different DSPSs to be deployed at the same time to tackle different use cases. It then raises the questions of how to avoid performance interference between collocated DSPSs and how to select the appropriate middleware that best improves the user experience. There are some preliminary efforts to enable federated execution on top of different streaming engines [38, 90], however, they all lack the ability to theoretically formulate an engine selection problem for a submitted streaming application, where the objective function and the resource and performance constraints caused by DSPS collocation are clearly defined. It is also interesting to investigate how to concatenate different DSPSs together to host a single streaming

¹³<https://www.consul.io/intro/index.html>

application, where each DSPS can handle the part of workload or streaming logic that it excels at processing.

12 SUMMARY

It is of great interest to study resource management and task scheduling in distributed stream processing systems in order to satisfy the Quality of Service (QoS) requirements with minimal resource cost. This topic has received extensive attention in the literature — many have paved the way for SLA-aware, self-adaptive deployment by proposing enabling techniques such as elastic resource scaling, dynamic task scheduling, and runtime operator parallelisation. However, there are still many gaps between the state-of-the-art and the prospect that the monitoring, tuning, and adaptation burden of deployment can be completely offloaded to a comprehensive resource management and scheduling framework, which can address key challenges of dynamic workload characteristics, heterogeneous cloud resources types, and ever-changing QoS requirements without requiring user intervention.

In this paper, we summarise the achievements made on this front and identify the gaps to bridge by presenting a comprehensive review of resource management and scheduling techniques in stream processing. Our narrative starts with defining the resource management and task scheduling problem, and then organising the research topics of interest around a singular context of achieving SLA-awareness and cost-efficiency while deploying stream processing systems on cloud. We also identified the issues and challenges associated with each research topic and developed a taxonomy of existing work to differentiate the specific work properties and method features. Following the structure of the taxonomy, we discussed each research topic in details and compared the strengths and weaknesses of different methods that fall into the same category. Finally, we shed light on the promising directions to promote future research in this area.

REFERENCES

- [1] 2010. Efficient event processing through reconfigurable hardware for algorithmic trading. *Proceedings of the VLDB Endowment* 3, 1-2 (2010), 1525–1528.
- [2] 2010. S4: Distributed Stream Computing Platform. In *Proceedings of the IEEE International Conference on Data Mining Workshops*. IEEE, 170–177.
- [3] 2014. A Catalog of Stream Processing Optimizations. *Comput. Surveys* 46, 4 (2014), 1–34.
- [4] Leonardo Aniello, Roberto Baldoni, and Leonardo Querzoni. 2013. Adaptive Online Scheduling in Storm. In *Proceedings of the 7th ACM International Conference on Distributed Event-Based Systems*. ACM Press, 207–218.
- [5] Joshua Auerbach, David F Bacon, Perry Cheng, and Rodric Rabbah. 2010. Lime: A Java-Compatible and Synthesizable Language for Heterogeneous Architectures. *ACM SIGPLAN Notices* 45, 10 (2010), 89–108.
- [6] Nathan Backman, Rodrigo Fonseca, and Ugur Çetintemel. 2012. Managing Parallelism for Stream Processing in the Cloud. In *Proceedings of the 1st International Workshop on Hot Topics in Cloud Data Processing (HotCDP '12)*. ACM Press, 1–5.
- [7] Cagri Balkesen, Nesime Tatbul, and M. Tamer Özsu. 2013. Adaptive Input Admission and Management for Parallel Stream Processing. In *Proceedings of the 7th ACM International Conference on Distributed Event-Based Systems*. ACM Press, 15–24.
- [8] Anne Benoit, Henri Casanova, Veronika Rehn-Sonigo, and Yves Robert. 2009. Resource Allocation Strategies for Constructive In-Network Stream Processing. In *Proceedings of the IEEE International Symposium on Parallel & Distributed Processing*. IEEE, 1–8.
- [9] Anne Benoit, Henri Casanova, Veronika Rehn-Sonigo, and Yves Robert. 2011. Resource Allocation for Multiple Concurrent In-Network Stream-Processing Applications. *Parallel Comput.* 37, 8 (2011), 331–348.
- [10] Muhammad Bilal and Marco Canini. 2017. Towards Automatic Parameter Tuning of Stream Processing Systems. In *Proceedings of the ACM Symposium on Cloud Computing*. ACM Press, 189–200.
- [11] Ioannis Boutsis and Vana Kalogeraki. 2012. RADAR: Adaptive Rate Allocation in Distributed Stream Processing Systems under Bursty Workloads. In *Proceedings of the 31st IEEE Symposium on Reliable Distributed Systems*. IEEE, 285–290.
- [12] Antonio Brogi, Gabriele Mencagli, Davide Neri, Jacopo Soldani, and Massimo Torquati. 2018. Container-Based Support for Autonomic Data Stream Processing Through the Fog. In *Proceedings of the 23rd European Conference on Parallel Processing*, Vol. 8374. Springer, 17–28.
- [13] Thilina Buddhika, Ryan Stern, Kira Lindburg, Kathleen Ericson, and Shrideep Pallickara. 2017. Online Scheduling and Interference Alleviation for Low-Latency, High-Throughput Processing of Data Streams. *IEEE Transactions on Parallel and Distributed Systems* 28, 12 (2017), 3553–3569.

- [14] Michael Cammert, Christoph Heinz, Jurgen Kramer, Bernhard Seeger, Sonny Vaupel, and Udo Wolske. 2007. Flexible Multi-Threaded Scheduling for Continuous Queries over Data Streams. In *Proceedings of the 23rd IEEE International Conference on Data Engineering Workshop*. IEEE, 624–633.
- [15] Michael Cammert, J Kramer, B. Seeger, and S. Vaupel. 2008. A Cost-Based Approach to Adaptive Resource Management in Data Stream Systems. *IEEE Transactions on Knowledge and Data Engineering* 20, 2 (2008), 230–245.
- [16] Paris Carbone, Stephan Ewen, Seif Haridi, Asterios Katsifodimos, Volker Markl, and Kostas Tzoumas. 2015. Apache Flink: Unified Stream and Batch Processing in a Single Engine. *Bulletin of the IEEE Computer Society Technical Committee on Data Engineering* 36, 1 (2015), 28–38.
- [17] Valeria Cardellini, Vincenzo Grassi, Francesco Lo Presti, and Matteo Nardelli. 2015. Distributed QoS-Aware Scheduling in Storm. In *Proceedings of the 9th ACM International Conference on Distributed Event-Based Systems (DEBS '15)*. ACM Press, 344–347.
- [18] Valeria Cardellini, Vincenzo Grassi, Francesco Lo Presti, and Matteo Nardelli. 2016. Optimal Operator Placement for Distributed Stream Processing Applications. In *Proceedings of the 10th ACM International Conference on Distributed and Event-based Systems*. ACM Press, 69–80.
- [19] Valeria Cardellini, Vincenzo Grassi, Francesco Lo Presti, and Matteo Nardelli. 2017. Optimal Operator Replication and Placement for Distributed Stream Processing Systems. *ACM SIGMETRICS Performance Evaluation Review* 44, 4 (2017), 11–22.
- [20] Valeria Cardellini, Vincenzo Grassi, Francesco Lo Presti, and Matteo Nardelli. 2015. On QoS-Aware Scheduling of Data Stream Applications over Fog Computing Infrastructures. In *Proceedings of the IEEE Symposium on Computers and Communication*. IEEE, 271–276.
- [21] Valeria Cardellini, Francesco Lo Presti, Matteo Nardelli, and Gabriele Russo Russo. 2017. Optimal Operator Deployment and Replication for Elastic Distributed Data Stream Processing. *Concurrency and Computation: Practice and Experience* 43, 34 (2017), 4334–4353.
- [22] Valeria Cardellini, Matteo Nardelli, and Dario Luzzi. 2016. Elastic Stateful Stream Processing in Storm. In *Proceedings of the International Conference on High Performance Computing & Simulation*. IEEE, 583–590.
- [23] Raul Castro Fernandez, Matteo Migliavacca, Evangelia Kalyvianaki, and Peter Pietzuch. 2013. Integrating Scale out and Fault Tolerance in Stream Processing Using Operator State Management. In *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD '13)*. ACM Press, 725–736.
- [24] Javier Cervino, Evangelia Kalyvianaki, Joaquin Salvachua, and Peter Pietzuch. 2012. Adaptive Provisioning of Stream Processing Systems in the Cloud. In *Proceedings of the 28th IEEE International Conference on Data Engineering Workshops*. IEEE, 295–301.
- [25] Badrish Chandramouli, Jonathan Goldstein, Roger Barga, Mirek Riedewald, and Ivo Santos. 2011. Accurate Latency Estimation in a Distributed Event Processing System. In *Proceedings of the 27th IEEE International Conference on Data Engineering (ICDE '11)*. IEEE, 255–266.
- [26] Shilpa Chaturvedi, Sahil Tyagi, and Yogesh Simmhan. 2017. Collaborative Reuse of Streaming Dataflows in IoT Applications. In *Proceedings of the 13th IEEE International Conference on e-Science*. IEEE, 403–412.
- [27] Andreas Chatzistergiou and Stratis D. Viglas. 2014. Fast Heuristics for Near-Optimal Task Allocation in Data Stream Processing over Clusters. In *Proceedings of the 23rd ACM International Conference on Information and Knowledge Management (CIKM '14)*. ACM Press, 1579–1588.
- [28] Wuhui Chen, Incheon Paik, and Zhenni Li. 2016. Cost-Aware Streaming Workflow Allocation on Geo-Distributed Data Centers. *IEEE Trans. Comput.* 1 (2016), 1–14.
- [29] Zhenhua Chen, Jielong Xu, Jian Tang, Kevin Kwiat, Charles Kamhoua, and Chonggang Wang. 2016. GPU-accelerated High-throughput Online Stream Data Processing. *IEEE Transactions on Big Data* 3, 99 (2016), 1–12.
- [30] Gianpaolo Cugola and Alessandro Margara. 2012. Processing Flows of Information: From Data Stream to Complex Event Processing. *Comput. Surveys* 44, 3 (2012), 1–62.
- [31] Frank Dabek, Russ Cox, Frans Kaashoek, and Robert Morris. 2004. Vivaldi: A Decentralized Network Coordinate System. *ACM SIGCOMM Computer Communication Review* 34, 4 (2004), 15–26.
- [32] Wenyun Dai, Longfei Qiu, Ana Wu, and Meikang Qiu. 2016. Cloud Infrastructure Resource Allocation for Big Data Applications. *IEEE Transactions on Big Data* 3, 99 (2016), 1–11.
- [33] Miyuru Dayarathna and Srinath Perera. 2018. Recent Advancements in Event Processing. *Comput. Surveys* 51, 2 (2018), 1–36.
- [34] Tiziano De Matteis and Gabriele Mencagli. 2016. Keep Calm and React with Foresight: Strategies for Low-Latency and Energy-Efficient Elastic Data Stream Processing. In *Proceedings of the 21st ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*. ACM Press, 1–12.
- [35] Tiziano De Matteis and Gabriele Mencagli. 2017. Elastic Scaling for Distributed Latency-Sensitive Data Stream Operators. In *Proceedings of the 25th Euromicro International Conference on Parallel, Distributed and Network-based Processing*. IEEE, 61–68.
- [36] Tiziano De Matteis and Gabriele Mencagli. 2017. Proactive Elasticity and Energy Awareness in Data Stream Processing. *Journal of Systems and Software* 127, C (2017), 302–319.
- [37] Marcos Dias de Assunção, Alexandre da Silva Veith, and Rajkumar Buyya. 2018. Distributed Data Stream Processing and Edge Computing: A Survey on Resource Elasticity and Future Directions. *Journal of Network and Computer Applications* 103, 1 (2018), 1–17.

- [38] Michael Duller, Jan S. Rellermeyer, Gustavo Alonso, and Nesime Tatbul. 2011. Virtualizing Stream Processing. In *Proceedings of the 12th International on Middleware Conference*. Springer, 269–288.
- [39] Raphael Eidenbenz and Thomas Locher. 2016. Task Allocation for Distributed Stream Processing. In *Proceedings of the 35th Annual IEEE International Conference on Computer Communications*. IEEE, 1–9.
- [40] Leila Eskandari, Zhiyi Huang, and David Eysers. 2016. P-Scheduler: Adaptive Hierarchical Scheduling in Apache Storm. In *Proceedings of the Australasian Computer Science Week Multiconference (ACSW '16)*. ACM Press, 1–10.
- [41] Havard Espeland, Paul B. Beskow, Hakon K. Stensland, Preben N. Olsen, Stale Kristoffersen, Carsten Griwodz, and Pal Halvorsen. 2011. P2G: A Framework for Distributed Real-Time Processing of Multimedia Data. In *Proceedings of the 40th International Conference on Parallel Processing Workshops*. IEEE, 416–426.
- [42] Lorenz Fischer and Abraham Bernstein. 2015. Workload Scheduling in Distributed Stream Processors Using Graph Partitioning. In *Proceedings of the IEEE International Conference on Big Data*. IEEE, 124–133.
- [43] Lorenz Fischer, Thomas Scharrenbach, and Abraham Bernstein. 2013. Scalable Linked Data Stream Processing via Network-Aware Workload Scheduling. In *Proceedings of the 9th International Conference on Scalable Semantic Web Knowledge Base Systems*. Springer, 81–96.
- [44] Avriilia Floratou, Ashvin Agrawal, Bill Graham, Sriram Rao, and Karthik Ramasamy. 2017. Dhalion: Self-Regulating Stream Processing in Heron. *Proceedings of the VLDB Endowment* 10, 12 (2017), 1825–1836.
- [45] Tom Z.J. Fu, Jianbing Ding, Richard T.B. Ma, Marianne Winslett, Yin Yang, and Zhenjie Zhang. 2015. DRS: Dynamic Resource Scheduling for Real-Time Analytics over Fast Streams. In *Proceedings of the 35th IEEE International Conference on Distributed Computing Systems*. IEEE, 411–420.
- [46] Bugra Gedik. 2014. Partitioning Functions for Stateful Data Parallelism in Stream Processing. *The VLDB Journal* 23, 4 (2014), 517–539.
- [47] Bugra Gedik, Henrique Andrade, Kun-Lung Wu, Philip S. Yu, and Myungcheol Doo. 2008. SPADE: The System S Declarative Stream Processing Engine. In *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD '08)*. ACM Press, 1123–1132.
- [48] Bugra Gedik, Scott Schneider, Martin Hirzel, and Kun-Lung Wu. 2014. Elastic Scaling for Data Stream Processing. *IEEE Transactions on Parallel and Distributed Systems* 25, 6 (2014), 1447–1463.
- [49] Javad Ghaderi, Sanjay Shakkottai, and R Srikant. 2016. Scheduling Storms and Streams in the Cloud. *ACM Transactions on Modeling and Performance Evaluation of Computing Systems* 1, 4 (2016), 1–28.
- [50] P. Taylor Goetz and Brian O'Neill. 2014. *Storm blueprints: patterns for distributed real-time computation*. Packt Pub. 1–426 pages. <https://www.oreilly.com/library/view/storm-blueprints-patterns/9781782168294/>
- [51] Vincenzo Gulisano, Ricardo Jimenez-Peris, Marta Patino-Martinez, Claudio Soriente, and Patrick Valduriez. 2012. StreamCloud: An Elastic and Scalable Data Streaming System. *IEEE Transactions on Parallel and Distributed Systems* 23, 12 (2012), 2351–2365.
- [52] Jiong He, Yao Chen, Tom Z.J. Fu, Xin Long, Marianne Winslett, Liang You, and Zhenjie Zhang. 2018. HaaS: Cloud-Based Real-Time Data Analytics with Heterogeneity-Aware Scheduling. In *Proceedings of the 38th IEEE International Conference on Distributed Computing Systems*, Vol. 1. IEEE, 1017–1028.
- [53] Thomas Heinze. 2011. Elastic Complex Event Processing. In *Proceedings of the 8th Doctoral Symposium on Middleware (MDS '11)*. ACM Press, 1–6.
- [54] Thomas Heinze, Leonardo Aniello, Leonardo Querzoni, and Zbigniew Jerzak. 2014. Cloud-Based Data Stream Processing. In *Proceedings of the 8th ACM International Conference on Distributed Event-Based Systems (DEBS '14)*. ACM Press, 238–245.
- [55] Thomas Heinze, Zbigniew Jerzak, Gregor Hackenbroich, and Christof Fetzer. 2014. Latency-Aware Elastic Scaling for Distributed Data Stream Processing Systems. In *Proceedings of the 8th ACM International Conference on Distributed Event-Based Systems (DEBS '14)*. ACM Press, 13–22.
- [56] Thomas Heinze, Yuanzhen Ji, Yinying Pan, Franz Josef Grueneberger, Zbigniew Jerzak, and Christof Fetzer. 2013. Elastic Complex Event Processing under Varying Query Load. In *Proceedings of the First International Workshop on Big Dynamic Distributed Data*. Springer, 25–30.
- [57] Thomas Heinze, Valerio Pappalardo, Zbigniew Jerzak, and Christof Fetzer. 2014. Auto-Scaling Techniques for Elastic Data Stream Processing. In *Proceedings of the 30th IEEE International Conference on Data Engineering Workshops*. IEEE, 296–302.
- [58] Thomas Heinze, Lars Roediger, Andreas Meister, Yuanzhen Ji, Zbigniew Jerzak, and Christof Fetzer. 2015. Online Parameter Optimization for Elastic Data Stream Processing. In *Proceedings of the Sixth ACM Symposium on Cloud Computing (SoCC '15)*. ACM Press, 276–287.
- [59] Nicolas Hidalgo, Daniel Wladdimiro, and Erika Rosas. 2017. Self-Adaptive Processing Graph with Operator Fission for Elastic Stream Processing. *Journal of Systems and Software* 127, 1 (2017), 205–216.
- [60] Christoph Hochreiner, Stefan Schulte, Schahram Dustdar, and Freddy Lecue. 2015. Elastic Stream Processing for Distributed Environments. *IEEE Internet Computing* 19, 6 (2015), 54–59.
- [61] Christoph Hochreiner, Michael Vogler, Stefan Schulte, and Schahram Dustdar. 2016. Elastic Stream Processing for the Internet of Things. In *Proceedings of the 9th IEEE International Conference on Cloud Computing*. IEEE, 100–107.
- [62] M. Reza Hoseiny Farahabady, Hamid R. Dehghani Samani, Yidan Wang, Albert Y. Zomaya, and Zahir Tari. 2016. A QoS-Aware Controller for Apache Storm. In *Proceedings of the 15th IEEE International Symposium on Network Computing and Applications*. IEEE, 334–342.

- [63] MohammadReza HoseinyFarahabady, Albert Y Zomaya, and Zahir Tari. 2017. QoS- and Contention- Aware Resource Provisioning in a Stream Processing Engine. In *Proceedings of the IEEE International Conference on Cluster Computing*. IEEE, 137–146.
- [64] Mahammad Humayoo, Yanlong Zhai, Yan He, Bingqing Xu, and Chen Wang. 2014. Operator Scale Out using Time Utility Function in Big Data Stream Processing. In *Proceedings of the International Conference on Wireless Algorithms, Systems, and Applications*. Springer, 54–65.
- [65] Waldemar Hummer, Christian Inzinger, Philipp Leitner, Benjamin Satzger, and Schahram Dustdar. 2012. Deriving a Unified Fault Taxonomy for Event-Based Systems. In *Proceedings of the 6th ACM International Conference on Distributed Event-Based Systems (DEBS '12)*. ACM Press, 167–178.
- [66] Waldemar Hummer, Benjamin Satzger, and Schahram Dustdar. 2013. Elastic Stream Processing in the Cloud. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 3, 5 (2013), 333–345.
- [67] Shigeru Imai, Thomas Chestna, and Carlos A. Varela. 2012. Elastic Scalable Cloud Computing Using Application-Level Migration. In *Proceedings of the Fifth IEEE International Conference on Utility and Cloud Computing*. IEEE, 91–98.
- [68] Shigeru Imai, Stacy Patterson, and Carlos A Varela. 2017. Maximum Sustainable Throughput Prediction for Data Stream Processing over Public Clouds. In *Proceedings of the 17th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*. IEEE, 1–10.
- [69] Atsushi Ishii and Toyotaro Suzumura. 2011. Elastic Stream Computing with Clouds. In *Proceedings of the 4th IEEE International Conference on Cloud Computing*. IEEE, 195–202.
- [70] Gabriela Jacques-Silva, Bugra Gedik, Rohit Wagle, Kun-Lung Wu, and Vibhore Kumar. 2012. Building User-Defined Runtime Adaptation Routines for Stream Processing Applications. *Proceedings of the VLDB Endowment* 5, 12 (2012), 1826–1837.
- [71] Jiahua Fan, Haopeng Chen, and Fei Hu. 2015. Adaptive Task Scheduling in Storm. In *Proceedings of the 4th International Conference on Computer Science and Network Technology*. IEEE, 309–314.
- [72] Jiawei Jiang, Zhipeng Zhang, Bin Cui, Yunhai Tong, and Ning Xu. 2017. StroMAX: Partitioning-Based Scheduler for Real-Time Stream Processing System. In *Proceedings of the International Conference on Database Systems for Advanced Applications*, Vol. 3882. Springer, 269–288.
- [73] Yuxuan Jiang, Zhe Huang, and Danny H. K. Tsang. 2017. Towards Max-Min Fair Resource Allocation for Stream Big Data Analytics in Shared Clouds. *IEEE Transactions on Big Data* 4, 1 (2017), 130–137.
- [74] Supun Kamburugamuve, Leif Christiansen, and Geoffrey Fox. 2015. A Framework for Real Time Processing of Sensor Data in the Cloud. *Journal of Sensors* 2015, 1 (2015), 1–11.
- [75] Supun Kamburugamuve and Geoffrey Fox. 2013. Survey of Distributed Stream Processing for Large Stream Sources. *Grids Ucs Indiana Edu* 2 (2013), 1–16.
- [76] Supun Kamburugamuve, Karthik Ramasamy, Martin Swamy, and Geoffrey Fox. 2017. Low Latency Stream processing: Apache Heron with Infiniband & Intel Omni-Path. In *Proceedings of the 10th International Conference on Utility and Cloud Computing*. ACM Press, 101–110.
- [77] J.O. Kephart and D.M. Chess. 2003. The Vision of Autonomic Computing. *Computer* 36, 1 (2003), 41–50.
- [78] Danish Khan, Kshiteej Mahajan, Rahul Godha, and Yuvraj Patel. 2015. *Empirical Study of Stragglers in Spark SQL and Spark Streaming*. Technical Report. 1–12 pages. <http://pages.cs.wisc.edu/>
- [79] Rohit Khandekar, Kirsten Hildrum, Sujay Parekh, Deepak Rajan, Joel Wolf, Kun-Lung Wu, Henrique Andrade, and Bugra Gedik. 2009. COLA: Optimizing Stream Processing Applications via Graph Partitioning. In *Proceedings of the 10th ACM/IFIP/USENIX International Conference on Middleware*. Springer, 308–327.
- [80] Alireza Khoshkbarforoushha, Rajiv Ranjan, Raj Gaire, Prem P. Jayaraman, John Hosking, and Ehsan Abbasnejad. 2015. Resource Usage Estimation of Data Stream Processing Workloads in Datacenter Clouds. *arXiv:1501.07020 [cs]* (2015). arXiv:1501.07020
- [81] Alireza Khoshkbarforoushha, Rajiv Ranjan, and Peter Strazdins. 2016. Resource Distribution Estimation for Data-Intensive Workloads: Give Me My Share and No One Gets Hurt! In *Communications in Computer and Information Science*. Vol. 393. Springer, 228–237.
- [82] Wilhelm Kleiminger, Evangelia Kalyvianaki, and Peter Pietzuch. 2011. Balancing Load in Stream Processing with the Cloud. In *Proceedings of the 27th IEEE International Conference on Data Engineering Workshops*. IEEE, 16–21.
- [83] Roland Kotto Kombi, Nicolas Lumineau, and Philippe Lamarre. 2017. A Preventive Auto-Parallelization Approach for Elastic Stream Processing. In *Proceedings of the 37th IEEE International Conference on Distributed Computing Systems*. IEEE, 1532–1542.
- [84] Sanjeev Kulkarni, Nikunj Bhagat, Masong Fu, Vikas Kedigehalli, Christopher Kellogg, Sailesh Mittal, Jignesh M. Patel, Karthik Ramasamy, and Siddharth Taneja. 2015. Twitter Heron: Stream Processing at Scale. In *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD '15)*. ACM Press, 239–250.
- [85] Alok Gautam Kumbhare, Yogesh Simmhan, and Viktor K. Prasanna. 2014. PLASStiCC: Predictive Look-Ahead Scheduling for Continuous Dataflows on Clouds. In *Proceedings of the 14th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*. IEEE, 344–353.
- [86] Myungcheol Lee, Miyoung Lee, Sung Jin Hur, and Ikkyun Kim. 2015. Load Adaptive Distributed Stream Processing System for Explosive Stream Data. In *Proceedings of the 17th International Conference on Advanced Communication Technology*, Vol. 5. IEEE, 753–757.

- [87] Boduo Li, Yanlei Diao, and Prashant Shenoy. 2015. Supporting Scalable Analytics with Latency Constraints. *Proceedings of the VLDB Endowment* 8, 11 (2015), 1166–1177.
- [88] Chunlin Li, Jing Zhang, and Youlong Luo. 2017. Real-Time Scheduling Based on Optimized Topology and Communication Traffic in Distributed Real-Time Computation Platform of Storm. *Journal of Network and Computer Applications* 87, 10 (2017), 100–115.
- [89] Teng Li, Jian Tang, and Jielong Xu. 2016. Performance Modeling and Predictive Scheduling for Distributed Stream Data Processing. *IEEE Transactions on Big Data* 7790, 99 (2016), 1–12.
- [90] Harold Lim and Shivnath Babu. 2013. Execution and Optimization of Continuous Queries with Cyclops. In *Proceedings of the International Conference on Management of Data (SIGMOD '13)*. ACM Press, 1069–1072.
- [91] Qian Lin, Beng Chin Ooi, Zhengkui Wang, and Cui Yu. 2015. Scalable Distributed Stream Join Processing. In *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD '15)*. ACM Press, 811–825.
- [92] Ming Liu, Liang Luo, Jacob Nelson, Luis Ceze, Arvind Krishnamurthy, and Kishore Atreya. 2017. IncBricks: Toward In-Network Computation with an In-Network Cache. *ACM SIGARCH Computer Architecture News* 45, 1 (2017), 795–809.
- [93] Ning Liu, Zhe Li, Jielong Xu, Zhiyuan Xu, Sheng Lin, Qinru Qiu, Jian Tang, and Yanzhi Wang. 2017. A Hierarchical Framework of Cloud Resource Allocation and Power Management Using Deep Reinforcement Learning. In *Proceedings of the 37th IEEE International Conference on Distributed Computing Systems*. IEEE, 372–382.
- [94] Xunyun Liu and Rajkumar Buyya. 2017. D-Storm : Dynamic Resource-Efficient Scheduling of Stream Processing Applications. In *Proceedings of the 23rd International Conference on Parallel and Distributed Systems*. IEEE, 1–8.
- [95] Xunyun Liu and Rajkumar Buyya. 2017. Performance-Oriented Deployment of Streaming Applications on Cloud. *IEEE Transactions on Big Data* 14, 8 (2017), 1–14.
- [96] Xunyun Liu, Amir Vahid Dastjerdi, Rodrigo N Calheiros, Chenhao Qu, and Rajkumar Buyya. 2017. A Stepwise Auto-Profiling Method for Performance Optimization of Streaming Applications. *ACM Transactions on Autonomous and Adaptive Systems* 12, 4 (2017), 1–33.
- [97] Xunyun Liu, Aaron Harwood, Shanika Karunasekera, Benjamin Rubinstein, and Rajkumar Buyya. 2017. E-Storm: Replication-Based State Management in Distributed Stream Processing Systems. In *Proceedings of the 46th International Conference on Parallel Processing*. IEEE, 571–580.
- [98] Yuan Liu, Xuanhua Shi, and Hai Jin. 2016. Runtime-Aware Adaptive Scheduling in Stream Processing. *Concurrency and Computation: Practice and Experience* 28, 14 (2016), 3830–3843.
- [99] Giorgia Lodi, Leonardo Aniello, Giuseppe A. Di Luna, and Roberto Baldoni. 2014. An Event-Based Platform for Collaborative Threats Detection and Monitoring. *Information Systems* 39 (2014), 175–195.
- [100] Björn Lohrmann, Peter Janacik, and Odej Kao. 2015. Elastic Stream Processing with Latency Guarantees. In *Proceedings of the 35th IEEE International Conference on Distributed Computing Systems*. IEEE, 399–410.
- [101] Björn Lohrmann, Daniel Warneke, and Odej Kao. 2014. Nephel Streaming: Stream Processing Under QoS Constraints at Scale. *Cluster Computing* 17, 1 (2014), 61–78.
- [102] Federico Lombardi, Leonardo Aniello, Silvia Bonomi, and Leonardo Querzoni. 2018. Elastic Symbiotic Scaling of Operators and Resources in Stream Processing Systems. *IEEE Transactions on Parallel and Distributed Systems* 29, 3 (2018), 572–585.
- [103] Manisha Luthra, Boris Koldehofe, Pascal Weisenburger, Guido Salvaneschi, and Raheel Arif. 2018. TCEP: Adapting to Dynamic User Environments by Enabling Transitions between Operator Placement Mechanisms. In *Proceedings of the 12th ACM International Conference on Distributed and Event-based Systems*. ACM Press, 136–147.
- [104] Kasper Grud Skat Madsen, Philip Thyssen, and Yongluan Zhou. 2014. Integrating Fault-Tolerance and Elasticity in a Distributed Data Stream Processing System. In *Proceedings of the 26th International Conference on Scientific and Statistical Database Management (SSDBM '14)*. ACM Press, 1–4.
- [105] Kasper Grud Skat Madsen, Yongluan Zhou, and Li Su. 2016. Enorm: Efficient Window-Based Computation in Large-Scale Distributed Stream Processing Systems. In *Proceedings of the 10th ACM International Conference on Distributed and Event-based Systems*. ACM Press, 37–48.
- [106] Lena Mashayekhy, Mahyar Movahed Nejad, Daniel Grosu, Quan Zhang, and Weisong Shi. 2015. Energy-Aware Scheduling of MapReduce Jobs for Big Data Applications. *IEEE Transactions on Parallel and Distributed Systems* 26, 10 (2015), 2720–2733.
- [107] Ruben Mayer, Boris Koldehofe, and Kurt Rothermel. 2014. Meeting Predictable Buffer Limits in the Parallel Execution of Event Processing Operators. In *Proceedings of the IEEE International Conference on Big Data*. IEEE, 402–411.
- [108] Ruben Mayer, Boris Koldehofe, and Kurt Rothermel. 2015. Predictable Low-Latency Event Detection with Parallel Complex Event Processing. *IEEE Internet of Things Journal* 2, 4 (2015), 274–286.
- [109] Gabriele Mencagli. 2016. A Game-Theoretic Approach for Elastic Distributed Data Stream Processing. *ACM Transactions on Autonomous and Adaptive Systems* 11, 2 (2016), 1–34.
- [110] Jefferson Morales, Erika Rosas, and Nicolas Hidalgo. 2014. Symbiosis: Sharing Mobile Resources for Stream Processing. In *Proceedings of the IEEE Symposium on Computers and Communications*. IEEE, 1–6.
- [111] Matteo Nardelli. 2016. QoS-Aware Deployment of Data Streaming Applications over Distributed Infrastructures. In *Proceedings of the 39th International Convention on Information and Communication Technology, Electronics and Microelectronics*. IEEE, 736–741.

- [112] Stephen Neuendorffer and Kees Vissers. 2008. Streaming Systems in FPGAs. In *Embedded Computer Systems: Architectures, Modeling, and Simulation*. Springer, 147–156.
- [113] Shadi A Noghabi, Kartik Paramasivam, Yi Pan, Navina Ramesh, Jon Bringham, Indranil Gupta, and Roy H Campbell. 2017. Samza: Stateful Scalable Stream Processing at LinkedIn. *Proceedings of the VLDB Endowment* 10, 12 (2017), 1634–1645.
- [114] Apostolos Papageorgiou, Ehsan Poormohammady, and Bin Cheng. 2016. Edge-Computing-Aware Deployment of Stream Processing Tasks Based on Topology-External Information: Model, Algorithms, and a Storm-Based Prototype. In *Proceedings of the 5th IEEE International Congress on Big Data*. IEEE, 259–266.
- [115] Boyang Peng, Mohammad Hosseini, Zhihao Hong, Reza Farivar, and Roy Campbell. 2015. R-Storm: Resource-Aware Scheduling in Storm. In *Proceedings of the 16th Annual Conference on Middleware (Middleware '15)*. ACM Press, 149–161.
- [116] Peter Pietzuch, Jonathan Ledlie, Jeffrey Shneidman, Mema Roussopoulos, Matt Welsh, and Margo Seltzer. 2006. Network-Aware Operator Placement for Stream-Processing Systems. In *Proceedings of the 22nd International Conference on Data Engineering (ICDE '06)*. IEEE, 49–49.
- [117] T Ralf, Muhammad Intizar Ali, Payam Barnaghi, Sorin Ganea, Frieder Ganz, Manfred Haushwirth, Brigitte Kjærsgaard, K Daniel, Alessandra Mileo, Septimiu Nechifor, Amit Sheth, and Vlasios Tsiatsis. 2014. Real Time IoT Stream Processing and Large-scale Data Analytics for Smart City Applications. In *Proceedings of the European Conference on Networks and Communications*. IEEE, 1–5.
- [118] Rajiv Ranjan. 2014. Streaming Big Data Processing in Datacenter Clouds. *IEEE Cloud Computing* 1, 1 (2014), 78–83.
- [119] Thomas Repantis, Xiaohui Gu, and Vana Kalogeraki. 2006. Synergy: Sharing-Aware Component Composition for Distributed Stream Processing Systems. In *Proceedings of the ACM/IFIP/USENIX International Conference on Middleware*, Vol. 4290. Springer, 322–341.
- [120] Stamatia Rizou, Frank Durr, and Kurt Rothermel. 2010. Solving the Multi-Operator Placement Problem in Large-Scale Operator Networks. In *Proceedings of the 19th International Conference on Computer Communications and Networks*. IEEE, 1–6.
- [121] Stamatia Rizou, Frank Durr, and Kurt Rothermel. 2011. Fulfilling End-To-End Latency Constraints in Large-Scale Streaming Environments. In *Proceedings of the IEEE International Performance Computing and Communications Conference*. IEEE, 1–8.
- [122] Stamatia Rizou, Frank Durr, Kurt Rothermel, F Durr, and Kurt Rothermel. 2010. Providing QoS Guarantees in Large-Scale Operator Networks. In *Proceedings of the 12th IEEE International Conference on High Performance Computing and Communications*. IEEE, 337–345.
- [123] Marek Rychly, Petr Koda, and Pavel Mr. 2014. Scheduling Decisions in Stream Processing on Heterogeneous Clusters. In *Proceedings of the Eighth International Conference on Complex, Intelligent and Software Intensive Systems*. IEEE, 614–619.
- [124] Marek Rychly, Petr Škoda, and Pavel Smrž. 2015. Heterogeneity-Aware Scheduler for Stream Processing Frameworks. *International Journal of Big Data Intelligence* 2, 2 (2015), 70–82.
- [125] Mohammad Sadoghi, Rija Javed, Naif Tarafdar, Harsh Singh, Rohan Palaniappan, and Hans-Arno Jacobsen. 2012. Multi-query Stream Processing on FPGAs. In *Proceedings of the 28th IEEE International Conference on Data Engineering*. IEEE, 1229–1232.
- [126] Amedeo Sapio, Ibrahim Abdelaziz, Abdulla Aldilajjan, Marco Canini, and Panos Kalnis. 2017. In-Network Computation is a Dumb Idea Whose Time Has Come. In *Proceedings of the 16th ACM Workshop on Hot Topics in Networks*. ACM Press, 150–156.
- [127] Kai-Uwe Sattler and Felix Beier. 2013. Towards Elastic Stream Processing: Patterns and Infrastructure. In *Proceedings of the First International Workshop on Big Dynamic Distributed Data*. IEEE, 49–54.
- [128] Benjamin Satzger, Waldemar Hummer, Philipp Leitner, and Schahram Dustdar. 2011. Esc: Towards an Elastic Stream Computing Platform for the Cloud. In *Proceedings of the 4th IEEE International Conference on Cloud Computing*. IEEE, 348–355.
- [129] Scott Schneider, Martin Hirzel, Bugra Gedik, and Kun-Lung Wu. 2012. Auto-Parallelizing Stateful Distributed Streaming Applications. In *Proceedings of the 21st International Conference on Parallel Architectures and Compilation Techniques (PACT '12)*. ACM Press, 53–64.
- [130] Scott Schneider and Kun-Lung Wu. 2017. Low-Synchronization, Mostly Lock-Free, Elastic Scheduling for Streaming Runtimes. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation*. ACM Press, 648–661.
- [131] Zoe Sebepou and Kostas Magoutis. 2011. CEC: Continuous Eventual Checkpointing for Data Stream Processing Operators. In *Proceedings of the 41st IEEE/IFIP International Conference on Dependable Systems & Networks*. IEEE, 145–156.
- [132] Vinay Setty, Roman Vitenberg, Gunnar Kreitz, Guido Urdaneta, and Maarten Van Steen. 2014. Cost-Effective Resource Allocation for Deploying Pub/Sub on Cloud. In *Proceedings of the 34th IEEE International Conference on Distributed Computing Systems*. IEEE, 555–566.
- [133] Zhiming Shen, Sethuraman Subbiah, Xiaohui Gu, and John Wilkes. 2011. CloudScale: Elastic Resource Scaling for Multi-Tenant Cloud Systems. In *Proceedings of the 2nd ACM Symposium on Cloud Computing (SOCC '11)*. ACM Press, 1–14.

- [134] Ce-Kuen Shieh, Sheng-Wei Huang, Li-Da Sun, Ming-Fong Tsai, and Naveen Chilamkurti. 2017. A Topology-Based Scaling Mechanism for Apache Storm. *International Journal of Network Management* 27, 3 (2017), 1933–1952.
- [135] Anshu Shukla and Yogesh Simmhan. 2017. Model-driven Scheduling for Distributed Stream Processing Systems. *arXiv:1702.01785 [cs]* (2017). arXiv:1702.01785
- [136] Pavel Smirnov, Mikhail Melnik, and Denis Nasonov. 2017. Performance-Aware Scheduling of Streaming Applications Using Genetic Algorithm. *Procedia Computer Science* 108, 6 (2017), 2240–2249.
- [137] Dawei Sun and Rui Huang. 2016. A Stable Online Scheduling Strategy for Real-Time Stream Computing Over Fluctuating Big Data Streams. *IEEE Access* 4, 1 (2016), 8593–8607.
- [138] Dawei Sun, Hongbin Yan, Shang Gao, Xunyun Liu, and Rajkumar Buyya. 2017. Rethinking Elastic Online Scheduling of Big Data Streaming Applications over High-Velocity Continuous Data Streams. *The Journal of Supercomputing* 74, 2 (2017), 615–636.
- [139] Dawei Sun, Guangyan Zhang, Chengwen Wu, Keqin Li, and Weimin Zheng. 2017. Building a Fault Tolerant Framework with Deadline Guarantee in Big Data Stream Computing Environments. *J. Comput. System Sci.* 89, 1 (2017), 4–23.
- [140] Dawei Sun, Guangyan Zhang, Songlin Yang, Weimin Zheng, Samee U. Khan, and Keqin Li. 2015. Re-Stream: Real-Time and Energy-Efficient Resource Scheduling in Big Data Stream Computing Environments. *Information Sciences* 319 (2015), 92–112.
- [141] Lauritz Thamsen, Thomas Renner, and Odej Kao. 2016. Continuously Improving the Resource Utilization of Iterative Parallel Dataflows. In *Proceedings of the 36th IEEE International Conference on Distributed Computing Systems Workshops*. IEEE, 1–6.
- [142] Rafael Tolosana-Calasanz, José Ángel Bañares, Congduc Pham, and Omer F. Rana. 2016. Resource Management for Bursty Streams on Multi-Tenancy Cloud Environments. *Future Generation Computer Systems* 55 (2016), 444–459.
- [143] Ankit Toshniwal, Jake Donham, Nikunj Bhagat, Sailesh Mittal, Dmitriy Ryabov, Siddarth Taneja, Amit Shukla, Karthik Ramasamy, Jignesh M. Patel, Sanjeev Kulkarni, Jason Jackson, Krishna Gade, and Maosong Fu. 2014. Storm@twitter. In *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD '14)*. ACM Press, 147–156.
- [144] Jonas Traub, Sebastian Breß, Tilmann Rabl, Asterios Katsifodimos, and Volker Markl. 2017. Optimized On-Demand Data Streaming from Sensor Nodes. In *Proceedings of the ACM Symposium on Cloud Computing*. ACM Press, 586–597.
- [145] Jan Sipke van der Veen, Bram van der Waaij, Elena Lazovik, Wilco Wijbrandi, and Robert J. Meijer. 2015. Dynamically Scaling Apache Storm for the Analysis of Streaming Data. In *Proceedings of the First IEEE International Conference on Big Data Computing Service and Applications*. IEEE, 154–161.
- [146] Smita Vijayakumar, Qian Zhu, and Gagan Agrawal. 2010. Dynamic Resource Provisioning for Data Streaming Applications in a Cloud Environment. In *Proceedings of the Second IEEE International Conference on Cloud Computing Technology and Science*. IEEE, 441–448.
- [147] Rohit Wagle, Henrique Andrade, Kirsten Hildrum, Chitra Venkatramani, and Michael Spicer. 2011. Distributed Middleware Reliability and Fault Tolerance Support in System S. In *Proceedings of the 5th ACM International Conference on Distributed Event-Based Systems*. ACM Press, 335–346.
- [148] Chunkai Wang, Xiaofeng Meng, Qi Guo, Zujian Weng, and Chen Yang. 2016. OrientStream: A Framework for Dynamic Resource Allocation in Distributed Data Stream Management Systems. In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management*. ACM Press, 2281–2286.
- [149] Chunkai Wang, Xiaofeng Meng, Qi Guo, Zujian Weng, and Chen Yang. 2017. Automating Characterization Deployment in Distributed Data Stream Management Systems. *IEEE Transactions on Knowledge and Data Engineering* 29, 12 (2017), 2669–2681.
- [150] Di Wang, Elke A. Rundensteiner, Han Wang, and Richard T. Ellison. 2010. Active Complex Event Processing: Applications in Real-Time Health Care. *Proceedings of the VLDB Endowment* 3, 1-2 (2010), 1545–1548.
- [151] Huayong Wang and Li-Shiuan Peh. 2014. MobiStreams: A Reliable Distributed Stream Processing System for Mobile Devices. In *Proceedings of the 28th IEEE International Parallel and Distributed Processing Symposium*. IEEE, 51–60.
- [152] Daniel Warneke and Odej Kao. 2011. Exploiting Dynamic Resource Allocation for Efficient Parallel Data Processing in the Cloud. *IEEE Transactions on Parallel and Distributed Systems* 22, 6 (2011), 985–997.
- [153] Joel Wolf, Nikhil Bansal, Kirsten Hildrum, Sujay Parekh, Deepak Rajan, Rohit Wagle, Kun-Lung Wu, and Lisa Fleischer. 2008. SODA: An Optimizing Scheduler for Large-Scale Stream-Based Distributed Computer Systems. In *Proceedings of the 9th ACM/IFIP/USENIX International Conference on Middleware (Middleware '08)*. Springer, 306–325.
- [154] Yingjun Wu and Kian-Lee Tan. 2015. ChronoStream: Elastic Stateful Stream Computation in the Cloud. In *Proceedings of the 31st IEEE International Conference on Data Engineering*. IEEE, 723–734.
- [155] Ying Xing, Jeong-Hyon Hwang, Uğur Çetintemel, and Stanley B Zdonik. 2006. Providing Resiliency to Load Variations in Distributed Stream Processing. In *Proceedings of the 32nd International Conference on Very Large Data Bases*. IEEE, 775–786.
- [156] Jielong Xu, Zhenhua Chen, Jian Tang, and Sen Su. 2014. T-Storm: Traffic-Aware Online Scheduling in Storm. In *Proceedings of the 34th IEEE International Conference on Distributed Computing Systems*. IEEE, 535–544.
- [157] Le Xu, Boyang Peng, and Indranil Gupta. 2016. Stela: Enabling Stream Processing Systems to Scale-in and Scale-out On-demand. In *Proceedings of the IEEE International Conference on Cloud Engineering*. IEEE, 22–31.
- [158] Lei Yang, Jiannong Cao, Yin Yuan, Tao Li, Andy Han, and Alvin Chan. 2013. A Framework for Partitioning and Execution of Data Stream Applications in Mobile Cloud Computing. *ACM SIGMETRICS Performance Evaluation Review* 40, 4 (2013), 23–32.

- [159] Nikos Zacheilas, Vana Kalogeraki, Nikolas Zygouras, Nikolaos Panagiotou, and Dimitrios Gunopoulos. 2015. Elastic Complex Event Processing Exploiting Prediction. In *Proceedings of the IEEE International Conference on Big Data*. IEEE, 213–222.
- [160] Matei Zaharia, Tathagata Das, Haoyuan Li, Timothy Hunter, Scott Shenker, and Ion Stoica. 2013. Discretized Streams: Fault-Tolerant Streaming Computation at Scale. In *Proceedings of the 24th ACM Symposium on Operating Systems Principles (SOSP '13)*. ACM Press, 423–438.
- [161] Jing Zhang, Chunlin Li, Liye Zhu, and Yanpei Liu. 2016. The Real-Time Scheduling Strategy Based on Traffic and Load Balancing in Storm. In *Proceedings of the 18th IEEE International Conference on High Performance Computing and Communications*. IEEE, 372–379.
- [162] Zhe Zhang, Yu Gu, Fan Ye, Hao Yang, Minkyong Kim, Hui Lei, and Zhen Liu. 2010. A Hybrid Approach to High Availability in Stream Processing Systems. In *Proceedings of the 30th IEEE International Conference on Distributed Computing Systems*. IEEE, 138–148.
- [163] Xinwei Zhao, Saurabh Garg, Carlos Queiroz, and Rajkumar Buyya. 2017. A Taxonomy and Survey of Stream Processing Systems. In *Software Architecture for Big Data and the Cloud* (1 ed.). Elsevier, 183–206.
- [164] Zhenhuan Gong, Xiaohui Gu, and John Wilkes. 2010. PRESS: PRedictive Elastic ReSource Scaling for Cloud Systems. In *Proceedings of the International Conference on Network and Service Management*. IEEE, 9–16.
- [165] Yongluan Zhou, Beng Chin Ooi, Kian-lee Tan, and Ji Wu. 2006. Efficient Dynamic Operator Placement in a Locally Distributed Continuous Query System. In *On the Move to Meaningful Internet Systems*. Springer, 54–71.
- [166] Qian Zhu and Gagan Agrawal. 2008. Resource Allocation for Distributed Streaming Applications. In *Proceedings of the 37th International Conference on Parallel Processing*. IEEE, 414–421.

Received March 2017