

CBP: A New Parallelization Paradigm for Massively Distributed Stream Processing

Qingsong Guo¹(✉) and Yongluan Zhou²

¹ North University of China, Taiyuan, China
qingsongg@gmail.com

² University of Southern Denmark, Odense, Denmark
zhou@imada.sdu.dk

Abstract. Resource efficiency is essential for distributed stream processing engines (DSPEs), in which a streaming application is modeled as an operator graph where each operator is parallelized into a number of instances to meet the low-latency and high-throughput requirements. The major objectives of optimizing resource efficiency in DSPEs include minimizing the communication cost by collocating the tasks that transfer a lot of data between each other, and by dynamically configuring the systems according to the load variations at runtime. In the current literature, most proposals handle these two optimizations separately, and a shallow integration of these techniques, such as performing the two optimizations one after another, would result in a suboptimal solution. In this paper, we present component-based parallelization (CBP), a new paradigm for optimizing the resource efficiency of DSPEs, which provides a framework for a deeper integration of the two optimizations. In the CBP paradigm, the operators are encapsulated into a set of non-overlapping components, in which operators are parallelized consistently, i.e., using the same partitioning key, and hence the intra-component communication is eliminated. According to the changes of workload, each component can be adaptively partitioned into multiple instances, each of which is deployed on a computing node. We build a cost model to capture both the communication cost and adaptation cost of a CBP plan, and then propose several optimization algorithms. We implement the CBP scheme and the optimization algorithms on top of Apache Storm, and verify its efficiency by an extensive experiment study.

1 Introduction

Real-time big data analysis requires processing of *continuous queries* (CQ) over fast streaming data with low latency. Usually, distributed stream processing engines (DSPEs) [1, 18, 22] organize CQs as an operator graph as shown in Fig. 1(a). To handle the deluge of data, one can resort to massive parallelization that each operator is cloned with a number of instances and its inputs are

The author from North University of China is supported by National Natural Science Foundation of China (61602427) and Natural Science Foundation of Shanxi(201601D202037).

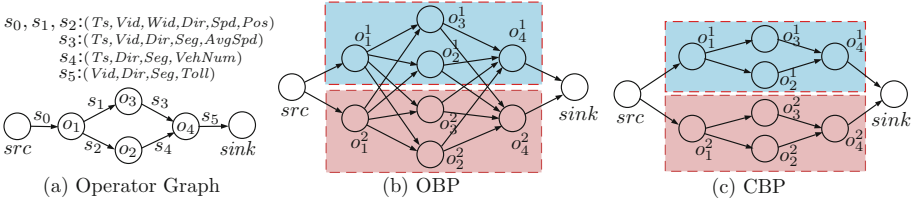


Fig. 1. Paradigms for parallelizing operator graph. This is a query that calculates the tolls of vehicles based on the source stream s_0 containing data of vehicles' speeds and positions. It consists of four operators: (1) a stateless operator o_1 that filters and partitions input data to the next operators; (2) two operators o_2 and o_3 calculate the average speed $AvgSpd$ and the traffic volume respectively; and (3) o_4 calculates the toll of each vehicle, which is a function of $AvgSpd$ and $SegNum$. The format of streams are specified in figure (a) and each operator o_i is designated with key k_i for partitioning the streams, where $k_1 = \{Ts, Vid, Spd, Dir, Seg, Pos\}$, $k_2 = \{Vid, Dir, Seg\}$, $k_3 = \{Dir, Seg\}$, and $k_4 = \{Dir, Seg\}$. We use two ways to parallelize the query: (1) in figure (b), the input streams of the operators are partitioned with different keys; and (2) in figure (c), the input streams of the four operators are partitioned consistently with the same key $\{Dir, Seg\}$.

partitioned into disjoint substreams. For the sake of resource efficiency, there are in general two critical optimizations to be considered:

1. **Runtime resource reconfiguration.** Load variations caused by the changes of data distribution and input rate are ubiquitous in the streaming context [22, 24, 25]. It is essential to provide adaptive data partitioning to achieve load balancing and to scale the number of parallel instances of each operator to avoid over-provisioning or under-provisioning.
2. **Communication cost minimization.** A large amount of data has to be continuously transmitted among the neighboring operators. Data transfer not only consumes bandwidth but also incurs significant computation overhead, including serializing and de-serializing the transmitted data. Optimizing the allocation of operator instances can to minimize cross-node communication can significantly reduce the resource consumption in a DSPE.

In existing solutions, the two problems are addressed separately. For example, M.A. Shah et al. [22] studied how to dynamically partition the input data at runtime to balance the workload across the parallel instances of an operator, while Y. Ahmad et al. [4] and P. Pietzuch et al. [19] investigated the operator placement to minimize the bandwidth usage by implicitly assumed assumption that operators do not need to be parallelized.

One can simply combine these methods to provide a complete solution. For example, we can first determine the parallelism for each operator [1], and transform the operator graph into a graph of operator instances. Thereafter, we can optimize the deployment by applying an operator placement algorithm, such as [4, 19]. Suppose we have two nodes in the cluster, Fig. 1(b) shows a possible parallelization and task allocation plan for the operator graph in Fig. 1(a).

Dynamic reconfigurations, such as re-scaling and load balancing, can be performed on each operator independently. However, such a shallow integration would provide suboptimal performance. As shown in Fig. 1(b), if the 4 operators are not parallelized consistently, e.g., partitioning the input on the same key, then each operator instance may have to transfer data to all its downstream instances. This limits the opportunity to minimize communication cost by collocating the instances that communicate with each other.

On the contrary, if we can parallelize the operators consistently using a common partitioning key, then we could have a plan as shown in Fig. 1(c), which minimizes cross-node communication. Although this idea may sound simple, it is nevertheless challenging to implement in a DSPE supporting runtime reconfiguration. First of all, dynamic data repartitioning makes it difficult or even impossible to achieve consistent parallelization of multiple operators given that the operators could be reconfigured at runtime independently. Secondly, dynamic scaling and data repartitioning involve a lot of state movements [22, 23]. The overhead of moving the states around operator instances in order to maintain the consistency of data partitioning may offset the benefits of collocating their communicating instances. Therefore we need a new parallelization framework that can optimize the parallelization of operators such that the total cost is minimized, including the communication and reconfiguration cost.

To address the challenges, we present **component-based parallelization (CBP)**—a new operator parallelization paradigm that considers both dynamic reconfiguration and resource optimization. In CBP, an operator graph is first decomposed into non-overlapping *components*, each being a connected subgraph. The operators in a component should have partitioning keys “compatible” with each other, i.e., sharing common attributes, and thus they can be parallelized using the same key. Each component acts as a singleton that is parallelized into a set of instances and the parallelism can be adapted at runtime in accordance with the load variations. This strategy simplifies the optimization of parallel stream processing and localizes the side-effect of reconfiguration within each component. In general, in the CBP paradigm, the more operators are grouped into a component, the less communication cost there would be, with a probable increase of the component’s reconfiguration cost. This is because every time we have to re-scale or re-balance one operator within a component, we have to trigger repartitioning of all the operators within the component. Therefore, a good trade-off should be found to minimize the total cost of a CBP plan.

We propose a cost-based optimizer to compute an optimized CBP plan for a given query graph. We develop a novel cost model that integrates the reconfiguration overhead into the optimization. We formally define the optimization problem as a **MINIMUM-COST-COMPONENT-BASED-PARALLELIZATION** problem (MCCBP). We prove that MCCBP is NP-hard, and then present two heuristic algorithms to solve it. All the techniques have been implemented on top of Apache Storm [1]. We compare our solutions with the operator placement algorithm by using both synthetic workload and an extension of the Linear Road Benchmark [6]. The experiments show that our methods can save the network

communication by up to 40%. Furthermore, our solutions can reduce the average end-to-end data latency by about 10% to 30%.

2 Background

2.1 Parallel Stream Processing

Continuous queries (CQs) [17] over streaming data are usually organized as an operator graph [1, 13, 18] in a *distributed stream processing engine* (DSPE). DSPEs like Flux [22] and StreamCloud [14] exploit data parallelism [11] to cope with the deluge of data, in which an operator is cloned into a set of independent instances each working on a partition of the input data. The number of partitions can be determined according to the input rates to achieve high throughput.

Operators can be categorized as stateless and stateful. For a stateless operator, the input tuples can be processed independently by any instance of it. While the stateful operators, such as **join** and **group-by aggregate**, are “context-sensitive”, so tuples with the same keys should be processed by the same instance to guarantee correctness. *Stream grouping* specifies the way how a stream of tuples is grouped and dispatched to the consumer operator instances. We consider two primitives: (1) *shuffle grouping*, where the input tuples are randomly routed to the operator instances; (2) *key grouping*, in which tuples are partitioned into a number of substreams based on a specified set of keys. Shuffle grouping is often the optimal choice for stateless operators since the load can be easily balanced, while key grouping is necessary for stateful operators.

Challenges of load variation. Usually, one can easily observe two kinds of variations over streaming data: (1) the fluctuation of input rates [22, 24], and (2) change of data value distributions [22, 24, 25]. If an SPE does not react to the variations, applications can run into problems:

- **Unmatched provision:** the over-provisioning or under-provisioning caused by the fluctuation of the input rates can result in low system utilization, high operational cost (e.g., using pay-as-you-go cloud services), and system failures.
- **Load imbalance:** the load distribution is skewed due to the change of data value distribution. For example some stream grouping keys become more popular than the others so that some operator instances are over-loaded while the others are under-loaded. Load imbalance can harm the processing latency and system throughput if the skewness is not resolved soon.

To handle the above problems, we resort to adaptation techniques including *dynamic scaling* [18] and *load balancing* [25]. CQs use the concept of *sliding windows* of tuples over a stream to specify the operational context of an operator. For instance, to perform a windowed join, we need to buffer the tuples within the current window(s) as the context of the join operation on the newly incoming tuples. This kind of context is called as *processing state* [8]. While processing an adaptation, the substreams should be reassigned around operator instances, and the processing states needs to be reallocated accordingly. This process is

called *state movement*. Note that both scaling and load balancing involve state movements, which consume both significant CPU and network bandwidth and thus cannot be ignored [22,23].

2.2 System Model

Data model. A data stream s is an unbounded and append-only sequence of tuples $(\dots, t_{i-1}, t_i, t_{i+1}, \dots)$. Each tuple $t = (\tau, \alpha)$ has a timestamp $\tau \in \mathbb{T}$ and a set of attributes $\alpha = (a_1, \dots, a_k)$. We assume that the attribute set α of every stream conforms to a relational schema. For simplicity, τ is assumed to be unique. In practice, if τ is not unique, existing systems usually use a unique sequence number to identify each tuple.

Operator model. A CQ is composed of a number of *operators*, each implementing a certain computation logic, such as *join*, *aggregate*, *filter*, or *user-defined functions*. An operator o is a 6-tuple, $(\text{IN}_o, \text{OUT}_o, K_o, F_o, W_o, \text{PS}_o)$, where IN_o and OUT_o are the input and output streams respectively. K_o is the *key*, a subset of attributes of the input streams IN_o , which used for partitioning IN_o . F_o defines the processing logic, where its operating context, i.e. the *processing state* PS_o , is specified by the sliding window W_o . For stateless operators like *map* and *filter*, $\text{PS} = \emptyset$.

We organize CQs as an operator graph $\mathbf{G} = (\mathcal{O}, \mathcal{S})$, which is a directed acyclic graph of the operator set \mathcal{O} and the stream set \mathcal{S} . A stream $s \in \mathcal{S}$ is represented as a directed arc (u_s, d_s) , $u_s, d_s \in \mathcal{O}$, where u_s and d_s are its producer and consumer respectively. Two special operators, *Src* and *Sink*, are responsible for spouting source streams and collecting the final results respectively. An operator graph is also referred to as a *topology* and these two terms are interchangeable throughout this paper.

Physical execution. The operator graph is executed on a cluster of identical nodes. The *execution graph* is a physical realization of the query in which each operator o is parallelized into multiple instances $\mathcal{I} = \{o^1, \dots, o^\pi\}$, where $\pi \in \mathbb{N}^+$ is the *parallelism*. For an input stream s of o , each tuple is a key-value pair $\langle k, v \rangle$, where v is the tuple and $k = t.K_o$. A partitioning function split the domain of K_o into p groups, where $p \gg \pi$. Then, the tuples of s , according their key values, form a number of substreams $\mathcal{S} = \{s^1, \dots, s^p\}$. An assignment $\mathcal{F} : \mathcal{S} \rightarrow \mathcal{I}$ allocate the processing of each substream to a unique operator instance. The degree of parallelism π and the assignment \mathcal{F} are adapted at runtime to handle load variations.

3 Component-Based Parallelization

3.1 CBP Abstraction

In essence, CBP decomposes an operator graph into a set of non-overlapping *components*, which act as the parallelization unit. In particular, CBP relies on

two essential properties: *compatibility* and *connectivity*. Compatibility concerns if some operators can be parallelized consistently. A set of operators $\{o_1, \dots, o_k\}$ is compatible iff the intersection of their keys is not empty, i.e., $K_{o_1} \cap \dots \cap K_{o_k} \neq \emptyset$. Note that the compatibility property is not *transitive*. For example, suppose we have three operators o_1 , o_2 , and o_3 with keys $K_1 = \{a_1, a_2\}$, $K_2 = \{a_2, a_3\}$, and $K_3 = \{a_1, a_3\}$ respectively. Even though any pair of them are compatible, they as a whole are incompatible because $K_1 \cap K_2 \cap K_3 = \emptyset$. The rationale of assembling the topology into components is to reduce the communication cost. One can benefit from placing compatible operators into a node only if they are connected by streams.

Formally, we can define a component as follow.

Definition 1 (Component). A *component* $\mathcal{C} = (\mathcal{O}_{\mathcal{C}}, \mathcal{S}_{\mathcal{C}})$ is an induced sub-graph of the operator graph $\mathcal{G} = (\mathcal{O}, \mathcal{S})$, where \mathcal{C} is connected and the operators in $\mathcal{O}_{\mathcal{C}}$ are compatible.

Let $\text{IN}(\mathcal{C})$ be the set of all input streams of the operators in component \mathcal{C} , then $\text{IN}(\mathcal{C}) = \cup_{o \in \mathcal{O}_{\mathcal{C}}} \text{IN}_o$. Assuming $\mathcal{O}_{\mathcal{C}} = \{o_1, \dots, o_{|\mathcal{C}|}\}$. The streams of $\text{IN}(\mathcal{C})$ can be grouped by a partition function over the key $K = K_{o_1} \cap \dots \cap K_{o_{|\mathcal{C}|}}$, which is the intersection of the keys of all the operators in \mathcal{C} . Since $K \neq \emptyset$, all the streams of $\text{IN}(\mathcal{C})$ can be partitioned uniformly into p substreams. For the convenience of discussion, we regard the streams in $\text{IN}(\mathcal{C})$ as a *composite stream* cs , which is partitioned into a set of substreams $\mathcal{CS} = \{cs^1, \dots, cs^p\}$. In addition, each component \mathcal{C} is parallelized into a number of instances $\mathcal{CI} = \{ci^1, \dots, ci^\pi\}$, where π is the *parallelism* of \mathcal{C} and each instance has a clone of the computation logic of each operator in \mathcal{C} . The parallel processing of the composite stream \mathcal{CS} is specified by an assignment $\mathcal{F}_{\mathcal{C}} : \mathcal{CS} \rightarrow \mathcal{CI}$, which is adapted at runtime to handle load variations.

4 MCCBP

4.1 Metrics

The cost of a CBP plan can be put into three parts: (1) *Processing cost* \mathcal{PC} , which is the CPU usage of the computation, (2) *Communication cost* \mathcal{CC} , which is the CPU and network usages of data transmission, and (3) *Adaptation cost* \mathcal{AC} , which is the CPU and network usages of carrying out adaptations.

In particular, we assume that \mathcal{PC} keeps the same regardless of the physical execution, and thus it can be disregarded in our cost model. In addition, we categorize data communication into *inter-component communication* and *intra-component communication*. The first one involves three sequential steps: (1) data serialization, (2) network propagation, and (3) de-serialization. Steps (1) and (3) consume CPU cycles and step (2) occupies network bandwidth. In contrast, the intra-component communication is realized via local memory access, whose overhead is negligible. Therefore, we only take the overhead of inter-component communication into account.

Statistics measurements. The cost calculation relies on the statistics of execution of the operator graph. In our implementation, the statistics are measured periodically over a sequence of time intervals of length Δ , which are called as *statistics windows*. Suppose the historical data spans m statistics windows that start at the time instance $\tau = 0$, then the timespan of historical data is $[0, m\Delta]$. The following discussions are confined within the timespan $[0, m\Delta]$.

For the input stream $s \in \mathbf{S}$ of a component that is split into p partitions, the statistics are represented as a sequence of histograms $\mathbf{Y}(s) = (Y_1, \dots, Y_m)$, where the histogram $Y_r = (y_{1,r}, \dots, y_{p,r})^T$, $r = 1 \dots m$, is a vector recording the data rate of the p partitions over the r -th statistics window. In other words, the data distribution of s at the r -th window can be approximated with Y_r . With \mathbf{Y} , we can derive other statistics on demands. For instance, denote $s = (o_i, o_j)$, then the load l_{ij} of s during $[0, m\Delta]$ is $l_{ij} = \sum_{r=1}^m \sum_{k=1}^p y_{kr}$.

The adaptation cost is closely related to the adaptation frequency f , where $\Delta = 1/f$. For simplicity, we assume that SPE performs an adaptation at each window. Let ψ_i^r be the number of state movements in the r -th adaptation of component C_i , then $\mathcal{AC} = \sum_{i=1}^{|\mathcal{C}|} \psi_i$, where $\psi_i = \sum_{r=1}^m \psi_i^r$ is the adaptation cost of C_i .

4.2 Problem Formulation

Consider an operator graph that is grouped into a set of disjoint components $\mathcal{C} = \{C_1, C_2, \dots\}$, it is called a CBP plan if $\cup_{i=1}^{|\mathcal{C}|} \mathbf{0}_{C_i} = \mathbf{0}$ and $\mathbf{0}_{C_i} \cap \mathbf{0}_{C_j} = \emptyset$ for any two components of \mathcal{C} . Let \mathbf{X} be the streams interconnecting components in \mathcal{C} . Let $w(C_i)$ be the adaptation cost of C_i and $c(s)$ be the communication cost incurred by stream s . Since \mathcal{PC} is independent on the CBP plan, the cost of a CBP plan \mathcal{C} , denoted as $cost(\mathcal{C})$, is measured by the sum of the communication cost \mathcal{CC} and adaptation cost \mathcal{AC} . That is,

$$cost(\mathcal{C}) = \mathcal{CC} + \mathcal{AC} = \sum_{s \in \mathbf{X}} c(s) + \sum_{C_i \in \mathcal{C}} w(C_i) \tag{1}$$

We introduce a constraint on the adaptation cost, $w(C_i) \leq \beta$, to prevent any component from being the bottleneck. Consequently, the objective of optimizing a CBP plan is to minimize $cost(\mathcal{C})$. We denote this problem as MINIMUM COST COMPONENT-BASED PARALLELIZATION (MCCBP), which is a variant of graph partitioning problem under constraints of connectivity and compatibility. Formally, it is stated as follow.

Definition 2 (MCCBP). *Given an operator graph $G = (\mathbf{0}, \mathbf{S})$ and a positive constant β , the MCCBP problem is to find a CBP plan, which is a partition of G into a set of disjoint components $\mathcal{C} = \{C_1, C_2, \dots\}$, to achieve the following objective:*

$$\begin{aligned} & \text{minimize } cost(\mathcal{C}) \\ & \text{subject to } \cup_{i=1}^{|\mathcal{C}|} \mathbf{0}_{C_i} = \mathbf{0} \\ & \qquad \qquad w(C_i) \leq \beta \end{aligned}$$

MCCBP can be proved to be NP-hard by simplifying it to a *Minimum-Capacity-Graph-Partitioning* (MCGP) problem, which has been shown to be NP-hard.

5 Computing CBP Plans

5.1 Greedy Algorithm

A straightforward idea is to obtain an initial CBP plan \mathcal{C}_0 in advance, and then make improvement incrementally. The algorithm, as shown in Algorithm 1, begins with the initial plan \mathcal{C}_0 (Line 2) and makes improvement step by step (Line 9–21). The initial plan \mathcal{C}_0 is generated by a depth-first search (DFS) of the operator graph. The traversal is tracked by an operator stack OS . In each iteration, we peek an operator from OS . Let o be the current operator being visited and $\mathcal{C}(o)$ be the component containing o . Then o will be popped out from OS if it has no unvisited child or is a leaf node. Otherwise we choose an unvisited child v of o and then check the compatibility between v and $\mathcal{C}(o)$. If they are compatible, v will be added into component $\mathcal{C}(o)$; Otherwise, a new component \mathcal{C}_i containing operator v is created.

The essence of Algorithm 1 is to reduce the cost by moving operators around components. Let $move(\mathcal{C}_i, \mathcal{C}_j, o_k)$ be the *potential movement* that attempts to move o_k from \mathcal{C}_i to \mathcal{C}_j . It is *admissible* if $o_k \in \mathcal{C}_i$ and the new operator set $\mathcal{C}_j \cup \{o_k\}$ can form a component. Given a CBP plan \mathcal{C} , the execution of the potential movement $move(\mathcal{C}_1, \mathcal{C}_2, o_k)$ gives rise to a new plan \mathcal{C}' if it is admissible. The admissibility of it is checked in Line 8.

The movement results in the following change of costs: (1) the change of communication cost between \mathcal{C}_1 and \mathcal{C}_2 , and (2) the change of adaptation costs of \mathcal{C}_1 and \mathcal{C}_2 . Hence the profit $\delta_{12}(o_k)$ of $move(\mathcal{C}_1, \mathcal{C}_2, o_k)$ consists of two parts: the changes of the communication cost and adaptation cost, denoted as $\delta_{12}^1(o_k)$ and $\delta_{12}^2(o_k)$ respectively. Let $\varphi_1(o_k)$ be the data rate between o_k and \mathcal{C}_1 . Then, $\varphi_1(o_k) = \sum_{\substack{(o_k, o_t) \in \text{SV}(o_t, o_k) \in \mathcal{S} \\ o_t \in \mathcal{C}_1}} l_{kt}$. $\varphi_1(o_k)$ does not contribute to \mathcal{CC} if $o_k \in \mathcal{C}_1$, otherwise it does. After the movement, $\varphi_i(o_k)$ contributes to \mathcal{CC} , but $\varphi_j(o_k)$ does not contribute to \mathcal{CC} . Therefore, the gain on communication cost is $\delta_{12}^1(o_k) = \varphi_2(o_k) - \varphi_1(o_k)$. Let $\psi(o_k)$ be the new adaptation cost of a component, then we have $\delta_{12}^2(o_k) = (\psi_1 + \psi_2) - (\psi_1(o_k) + \psi_2(o_k))$.

Summing all together, we get the overall profit of the movement, $\delta_{12}(o_k) = \delta_{12}^1(o_k) + \delta_{12}^2(o_k)$. In each run, we choose an admissible movement with the maximum positive profit to execute. Suppose that $\delta_{12}(o_k)$ is the best movement in the current run, then the load and state statistics of \mathcal{C}_1 and \mathcal{C}_2 should be changed after the execution of $move(\mathcal{C}_1, \mathcal{C}_2, o_k)$ (Line 14). The movement also causes changes of the profits of any admissible movement involving \mathcal{C}_1 or \mathcal{C}_2 . To prepare the next iteration, we should recompute the profits of these admissible movements (Line 15).

Algorithm 1. Greedy Algorithm

Input: Operator graph $G = (O, S)$, load statistics $\{Y(s_1), Y(s_2), \dots\}$, state statistics $\{Z(o_1), Z(o_2), \dots\}$

Output: CBP plan \mathcal{C}

```

1  $\mathcal{C} \leftarrow \text{InitialPartition}(G)$ ;
2 compute load statistics  $\mathcal{Y}(\mathcal{C}_i)$ , state statistics  $\mathcal{Z}(\mathcal{C}_i)$ , and adaptation cost  $\psi_i$  for
  each component  $\mathcal{C}_i \in \mathcal{C}$ ;
3  $\delta \leftarrow 1.0$ ;
4 while  $\delta > 0$  do
5   foreach  $o_k \in O$  do
6      $\mathcal{C}_i \leftarrow$  get the component containing  $o_k$ ;
7     foreach  $\mathcal{C}_j \in |\mathcal{C}|$  and  $j \neq i$  do
8       if  $a_{jk} \neq -1$  and  $\mathcal{C}_j \cup \{o_k\}$  is compatible then
9          $\delta_{ij}^1(o_k) \leftarrow \ell_{jk} - \ell_{ik}$ ;
10         $\delta_{ij}^2(o_j) \leftarrow (\psi_i + \psi_j) - (\psi_i(o_k) + \psi_j(o_k))$ ;
11         $\delta_{ij}(o_k) \leftarrow \delta_{ij}^1(o_k) + \delta_{ij}^2(o_k)$ ;
12       $\delta \leftarrow \max\{\delta_{ij}(o_k)\}$ ;
13      move  $o_k$  from  $\mathcal{C}_i$  to  $\mathcal{C}_j$ ;
14      update the load and state statistics of  $\mathcal{C}_1$  and  $\mathcal{C}_2$ ;
15      recompute the profits for any admissible movement involves  $\mathcal{C}_i$  or  $\mathcal{C}_j$ ;
16 return  $\mathcal{C}$ ;
```

5.2 MWSC

We proceed to consider an alternative that transforms MCCBP into the *minimum weighted set cover problem* (MWSC). Let $\Omega = \{\mathcal{C}_1, \mathcal{C}_2, \dots\}$ be a set containing all the possible components of O . Let N be the cardinality of Ω , i.e., $N = |\Omega|$. A CBP plan $\mathcal{C} = \{\mathcal{C}_i | \mathcal{C}_i \in \Omega\}$ is a subset of Ω . It is apparent that the plan \mathcal{C} is a set cover of O , since $\bigcup_{i=1}^{|\mathcal{C}|} \mathcal{C}_i = O$ and $\mathcal{C}_i \cap \mathcal{C}_j = \emptyset$ for $\forall \mathcal{C}_i, \mathcal{C}_j \in \mathcal{C}$. Therefore, MCCBP is equivalent to find a subset \mathcal{C} of Ω such that \mathcal{C} is a partition of O . We attempt to optimize this problem by enumerating all the feasible components and finding the optimal CBP plan from them.

Each component associates with adaptation cost ψ_i and intra-component communication cost ϕ_i , where $\phi_i = \sum_{o_i, o_j \in \mathcal{C}} l_{ij}$. For each component $\mathcal{C}_i \in \Omega$, we assign a weight w_i to it such that $w_i = \psi_i + l - \phi_i$, where l is the overall load, $l = \sum_{i=1}^n \sum_{j=1}^n l_{ij}$. It is obvious that $\psi_i > 0$ and $l - \phi_i \geq 0$.

Let x_i be a decision variable that indicates whether component \mathcal{C}_i is chosen in the set cover \mathcal{S} , where $x_i = 1$ if \mathcal{C}_i is picked, or $x_i = 0$ otherwise. Then the MCCBP is transformed to the weighted set cover problem. A set cover \mathcal{S} of O has some redundant operators, for example $\mathcal{C}_i \cap \mathcal{C}_j = o_k$. Denote \mathcal{S}' as the new set cover by discarding o_k . Since $\psi_i > 0$ and $l - \phi_i \geq 0$, the cost of \mathcal{S}' is definitely smaller than that of the former one, i.e., $w(\mathcal{S}') < w(\mathcal{S})$. Finally, we can get the minimum set cover of O by removing all the redundant operators.

Algorithm 2. MWSC

Input: Operator graph $\mathbf{G} = (\mathbf{0}, \mathbf{S})$, load statistics $\{\mathbf{Y}(s_1), \mathbf{Y}(s_2), \dots\}$, state statistics $\{\mathbf{Z}(o_1), \mathbf{Z}(o_2), \dots\}$

Output: CBP plan \mathcal{C}

```

1  $l \leftarrow \sum_{i=1}^n \sum_{j=1}^n l_{ij}$  ; /* overall loads */
2  $\Omega \leftarrow \text{Enumerate}(\mathbf{G}, k)$  ;
3 foreach component  $\mathbf{C}_i$  in  $\Omega$  do
4     compute the adaptation cost  $\psi_i$  ;
5      $\phi_i \leftarrow \sum_{o_i, o_j \in \mathbf{C}_i} l_{ij}$  ;
6      $w_i \leftarrow \psi_i + l - \phi_i$  ; /* weight of  $\mathbf{C}_i$  */
7  $\mathcal{S} \leftarrow$  compute the MWSC of  $\mathbf{0}$  over  $\Omega$  ;
8  $\mathcal{C} \leftarrow \mathcal{S}$  ;
9 return  $\mathcal{C}$  ;
```

Definition 3 (MWSC). Given a universe $\mathbf{0}$ and a family Ω of subsets of $\mathbf{0}$, the minimum weighted set cover of $\mathbf{0}$ can be expressed as an integer linear programming:

$$\begin{aligned}
 & \text{minimize } w(\mathcal{S}) = \mathbf{w}^T \mathbf{x} & (2) \\
 & \text{subject to } \sum_{\mathbf{C}_i: o \in \mathbf{C}_i}^N x_i \geq 1 \quad \text{for each operator } o \in \mathbf{0}, \\
 & \quad \quad \quad x_i \in \{0, 1\}
 \end{aligned}$$

where $\mathbf{w} = (w_1, \dots, w_N)$ is the weight vector and $\mathbf{x} = (x_1, \dots, x_N)$ is the decision vector for Ω respectively.

Apparently, a MWSC is a partition of $\mathbf{0}$. Thus,

$$w(\mathcal{S}) = \sum_{i=1}^N x_i \psi_i + |\mathbf{S}|l - \sum_{i=1}^N x_i \phi_i = \underbrace{\sum_{i=1}^N x_i \psi_i}_{\mathcal{A}\mathcal{C}} + \underbrace{\left[l - \sum_{i=1}^N x_i \phi_i\right]}_{\mathcal{C}\mathcal{C}} + \underbrace{(|\mathbf{S}| - 1)l}_{\text{constant}} \quad (3)$$

where $|\mathbf{S}|$ is the number of edges of $\mathbf{G} = (\mathbf{0}, \mathbf{S})$.

Comparing to the cost model Eq. (1), we have the first component $\sum_{i=1}^N x_i \psi_i$ and the second $l - \sum_{i=1}^N x_i \phi_i$ of Eq. (3) equal to the adaptation cost $\mathcal{A}\mathcal{C}$ and communication cost $\mathcal{C}\mathcal{C}$ respectively. As the third component $(|\mathbf{S}| - 1)l$ is a constant, the best solution of MWSC is equivalent to the optimal CBP plan.

The idea is depicted in Algorithm 2. We first enumerate all the possible components of \mathbf{G} (Line 2). Then we compute the adaptation cost ψ_i and the load ϕ_i of each component \mathbf{C}_i , and assign a weight to each component (Line 3–6). Finally, we compute a solution \mathcal{S} of MWSC and take it as a CBP plan by discarding all the redundant operators (Line 7–8). MWSC can be solved exactly with a MIP solver like Gurobi [2] when N is not too large. But we also implement a greedy routine to solve MWSC (Line 7) according to the description

in [9, Chap. 35]. The greedy routine is a useful option when N is large. Since the set cover \mathcal{S} obtained through the greedy routine might not be a CBP plan, we have to remove the redundant operators to get the final solution \mathcal{C} .

6 Evaluation

6.1 Experimental Setup

Evaluation metrics—We use the following metrics in the evaluation:

- **Communication cost** counts the number of tuples transmitted through inter-component communication.
- **Adaptation cost** counts the number of state movements in an adaptation process.
- **End-to-end latency** indicates the time completing the processing of a source tuple. It includes the time spent on processing, adaptation, and communication, and thus it is a overall metric to reflect the effectiveness of CBP.

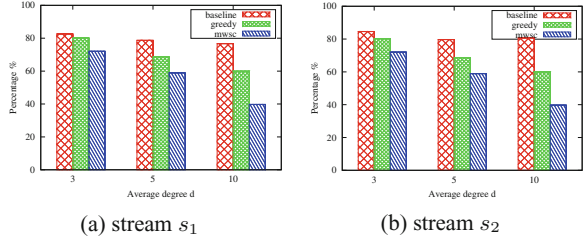
Tested solutions—We implement the sparse-cut algorithm, a graph partition algorithm used in COLA [15], to compare with our solutions. Note that the objective of baseline is merely to minimize communication cost. In general, we evaluated the following three algorithms: (1) *greedy algorithm*, (2) the *MWSC algorithm*, and (3) the *baseline algorithm* which implements an OBP-based operator placement algorithm in [15].

We implement our algorithms in Java and integrate them with Apache Storm [1] by extending it with runtime adaptation. Part of the experiments are conducted via simulation, while the rest are conducted on Amazon’s EC2 with medium VM instances (*m1.medium*), where each has 1.7 GB of RAM, moderate IO performance and one EC2 compute unit (approximately equivalent to a 1.2 GHz 2007 Xeon CPU). While these VMs have low processing capabilities, they are representatives of public cloud VMs.

6.2 Simulation Result

In the test, we used a randomized topology $\mathbf{G} = (\mathbf{0}, \mathbf{S})$. In the topology \mathbf{G} , each operator o , except *src*, maintains computing states and randomly forwards the received data to downstream operators according to the selectivity $\delta(o)$. The specific setting of \mathbf{G} is summarized in Fig. 2. Operator *src* generates two synthetic streams s_1 and s_2 to simulate two types of variations, where the key values of s_1 and s_2 follow the uniform distribution and Zipf respectively. Therefore, s_1 only results in scaling. In contrast, data distribution of s_2 is skewed and thus the adaptations involve both scaling and load balancing. Each operator of \mathbf{G} randomly chooses two attribute of *sch* as the partition key. The data arrivals of s_1 and s_2 follow a Poisson process $X(t) : P[N(t + \tau) - N(t) = k] = (k!)^{-1} e^{-\lambda\tau} (\lambda\tau)^k$, where τ is set to 1 s and $\lambda = 10,000$. Both s_1 and s_2 conform to

Parameters	Settings
Random graph	$G = (O, S, d)$
Number of operators	$ O = 100$
Average degree d	$d = \{3, 5, 10\}$
Selectivity $\delta(o)$	$\delta(o) \sim N(0.5, 1.0)$
Size of states $ PS_o $	$ PS_o = X(t)$

Fig. 2. Setting of parameters

Fig. 3. Comparison of communication costs

the schema: $\text{SynStream}(ts:\text{Unix timestamp}, a_1:\text{int}, a_2:\text{int}, a_3:\text{int}, a_4:\text{int})$, in which each attribute has 4 Bytes.

We measured the communication cost and state movements by varying the average degree d and the adaptation frequency f . Let N_1 be the number of tuples processed by all the operators and N_2 be the number of tuples in the states of all the operators at every adaptation. We calculated the percentages, $\frac{100n_c}{N_1}$ and $\frac{100n_a}{N_2}$, of tuples involved in the communication and state movement, where n_c and n_a are the communication cost and adaptation cost respectively.

Comparison of communication costs—Figure 3a and b show the percentages achieved by three algorithms. We can observe that the baseline algorithm can save the cost by at most 20%, but the CBP solutions can reduce the cost by at least 20%. In particular, the greedy algorithm saves about 20% when $d = 3$, and it increases to 40% when $d = 10$. MWSC outperforms the greedy algorithm. It significantly reduces the communication cost by about 27.8% when $d = 3$ and by nearly 60% when $d = 10$. The baseline algorithm deploys the operator graph based on a placement plan, which is generated in advance by graph partitioning. Since the operators are incompatible, the physical topology of the query changes as adaptation process. The parallelization plan is no longer optimal when the physical topology has been changed. Therefore, we cannot optimize the communication cost efficiently with operator placement.

The intra-component communications of a CBP plan are eliminated completely regardless of the change of physical topology. This is confirmed by Fig. 3. By comparing Fig. 3a and b, the communication costs of CBP solutions keep the same regardless of the difference of s_1 and s_2 . However, the costs of baseline algorithm is slightly different over s_1 and s_2 , where the cost is about 3% higher over s_2 than that over s_1 . The frequency f shows similar impact to the algorithms.

Comparison of adaptation costs—Figures 4 and 5 shows the impact of load variation and the adaptation frequency. In this experiment, the frequency f is varied by changing the length of adaptation window from 1 min to 10 min, i.e., $1/f = \{1, 2, 5, 10\}$. Figure 4 plots the adaptation cost of each algorithm when $1/f = 1$. It is clear that CBP has larger adaptation costs than the baseline algorithm. Moreover, the adaptation cost over a skewed stream, s_2 in Fig. 4b, is higher than the uniformly distributed stream, s_1 in Fig. 4a. We observe similar

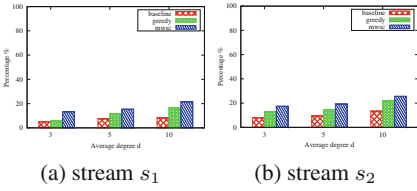


Fig. 4. Comparison of adaptation costs

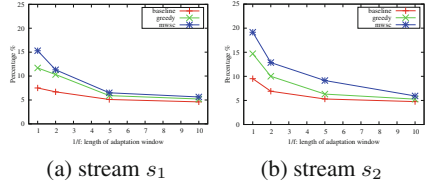


Fig. 5. Adaptation costs with respect to $1/f$

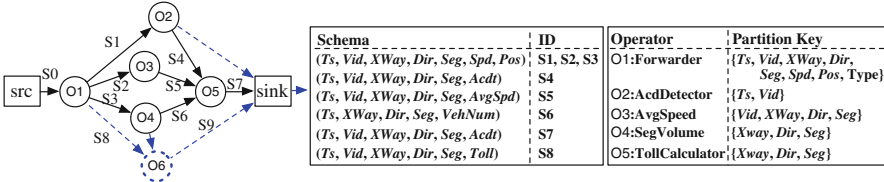


Fig. 6. Operator graph for LRB

results when $f = \{2, 5, 10\}$. The results also justify an implicit assumption in this paper that the adaptation cost is normally higher when we assemble operators into components.

Figure 5 shows the impact of adaptation frequency f . As we can see from the figure, the cost drops greatly at the beginning when we increase $1/f$. The number of state movements is determined by two factors: (1) the adaptation frequency f , and (2) the skewness of the data. The skewness usually goes serious if we increase $1/f$, i.e., it always involves more state movements in a single adaptation. As we expected, the decline of adaptation cost is much gentle when $1/f$ is larger.

6.3 End-to-End Latency

We proceed to evaluate the end-to-end latency of the tested solutions. In this experiment, we use the *Linear Road Benchmark* (LRB) [6]. LRB models a road toll network, in which tolls depend on the level of congestion. The primitive LRB gadget only has 7 operators, which is too small to represent a large-scale computation. So we extend it by connecting a number of LRB gadgets together with a road network. The road network $G = (V, E)$ is a graph where an edge $e \in E$ stands for an expressway of LRB and a vertex $v \in V$ represents the joint of expressways. This extension has a wide range of applications. If we want to measure the traffic between two locations or track the route of a vehicle, then an LRB gadget must dispatch result to its downstream LRB gadgets. Consequently, we introduce a new operator o_6 to calculate the traffic between every pairs of vertices every 1 min. Figure 6 shows the topology of the extended LRB, where some new streams (blue dashed arcs), have been added into a LRB gadget to fulfill the requirement.

Table 1. Statistics about end-to-end latency (ms)

$1/f$	Mean				Median				95%				Maximum			
	1	2	5	10	1	2	5	10	1	2	5	10	1	2	5	10
Greedy	677	610	566	617	141	121	109	116	1501	1236	1095	1130	2825	2223	1736	2117
MWSC	583	517	534	602	131	120	97	118	1532	1333	1054	1171	3103	2703	1853	1853
Baseline	775	710	673	681	153	137	114	127	1017	928	856	889	2109	1809	1673	1681

$G = (V, E)$ is generated with the random graph presented in Sect. 6.2. In particular, $|V| = 10$, $|E| = 30$, and the average degree $d = 3$. Therefore, we have 30 LRB gadgets and 180 operators in total excluding *srcs* and *sinks*. For each LRB gadget, the data rate of the source stream is controlled with the Poisson process used in the previous section. The experiments are conducted on EC2 with 30 VMs and accomplished in two phases: (1) We first deploy it over EC2, and keep it running for two hours to collect statistics. The length of statistic window is set to 1 min. (2) With the statistics, we partition the topology into components or subgraphs with the tested algorithms. Thereafter, we deploy the partitioned topology on EC2 and run the experiments.

We measure the end-to-end latency at 4 scales of the adaptation frequency f , i.e., $1/f = \{1, 2, 5, 10\}$. The latency values are given in Table 1, where “95%” is the 95th percentile of latency. In general, the results follow what we expected. By comparing the mean, we observe that our algorithms reduce the latencies by about 10%–25%. It shows that the CBP algorithm can indeed improve the performance and thus confirms the effectiveness of CBP. We can further identify the impacts of adaptation process and load imbalance in these values. For example, the CBP algorithms are more sensitive to adaptation process and load imbalance comparing to the greedy algorithm. The maximum latency is 3103 ms for MWSC when $1/f = 1$, which is higher than the maximum latency of Baseline.

Tuples with a latency smaller than the median are less affected by the adaptation process and load imbalance. In contrast, tuples with latencies larger than 95-th percentiles are greatly affected by the adaptation process and load imbalance. We take the latency when $1/f = 1$ as an example, the medians of MWSC and Greedy are about 75% and 87% of that of Baseline. So the results confirm that CBP can save communication cost efficiently. In contrast, the 95-th percentiles for MWSC and Greedy are about 29% and 26% greater than the baseline algorithm.

During an adaptation, input tuples are buffered by the upstream operators. The tuples will be replayed to downstream after the completion of adaptation. Therefore, adaptation process increases the end-to-end latencies for a portion of tuples. As we can see from the table, the maximum latency peaks up to about 3 s.

For each algorithm, each numeral characteristic drops with the increase of $1/f$ at first and then grow with the increase of $1/f$ on the contrary. This behavior is obvious for the 95-th percentile. In terms of the 95-th percentile, it is obvious MWSC is higher than Greedy and Baseline. This phenomena confirms the

impact of adaptation process and load imbalance. The adaptation cost drops with the increase of $1/f$, but load imbalance get worse on the contrary. Thus we observe that all lines are concave. It means that the adaptation frequency is very important as it can trade off between impact of adaptation cost and load imbalance. In this experiment, $f = 1/2$ is the best choice for MWSC and $f = 1/5$ is the best choice for Greedy and Baseline.

7 Related Work

Parallel stream processing. Much work has been focused on exploiting parallelism in stream processing. The early SPEs aim at providing transparent parallelization for distributed stream processing in a shared-nothing environment. Aurora [7] and Borealis [3] supports intra-query parallelism by organizing a topology into a set of boxes and conducting parallelization via *box-splitting*.

Many SPE proposals, e.g., System S [5] and Flux [22], leverage partitioned parallelism [11] to improve scalability. They propose new “Exchange” operators between stream producers and consumers to encapsulate the adaptive state partitioning and stream routing. In recent years, many efforts have been made to improve the scalability of parallelization [12, 20, 21]. The MapReduce model [10] enables programmer to think in a *data-centric* fashion and hence provides a practical implementation for partitioned parallelism. Distributed SPEs like Apache Storm [1], Yahoo! S4 [18], and StreamCloud [14] are inspired by such a model.

Operator placement. If an application is geographically distributed, the transmission latency is sensitive to the communication channels. The SAND project [4] exploits the knowledge of the underlying network characteristics such as topology and link bandwidths to make intelligent in-network placement of query graph. In contrast, [19] develops a *stream-based overlay network* (SBON) over Borealis, which is a network-aware optimization framework that manages operator placement within a pool of wide-area overlay nodes in order to make efficient use of network bandwidth. The placement decisions are made based on the cost space that encodes multidimensional metrics such as latency and load.

COLA [15] employs graph-partitioning algorithms to compute an optimal allocation of operators with regard to a cost model that captures the communication and CPU costs. The operator graph is partitioned into processing elements (PE) at compile-time, which acts as a deployable unit. COLA aims at balancing load across the processing nodes and minimizing the communication cost of the PEs. It only measures the CPU cost incurred by processing and communicating, but ignores the network bandwidth usage. In addition, COLA does not consider how to parallel the operators. Moreover a partition plan obtained at compile-time is incapable to handle the load variations at runtime.

The essence of operator placement is to optimize an assignment of operators to the computing nodes based on an objective function. Unfortunately, the existing solutions are static and the cost of the state migration cannot be ignored in the presence of load variations. For more detailed comparisons of the placement strategies, please refer to a survey paper [16].

8 Conclusion

We present CBP, a succinct parallelization paradigm for DSPEs that leverages both the connectivity and compatibility of operators. CBP seamlessly integrates operator placement with parallelization and thereby provides a framework to integrate the optimizations of runtime resource reconfiguration and communication cost minimization. Furthermore, we introduce a cost model that captures the cost of communication and adaptation. Two algorithms are proposed to optimize the CBP plans for a given computation. The extensive experiments confirm that an optimized CBP plan can improve the resource efficiency of DSPEs significantly.

References

1. Apache Storm. <http://storm.apache.org/>
2. Gurobi Parallel MIP solver. <http://www.gurobi.com/resources/getting-started/mip-basics>
3. Abadi, D.J., Ahmad, Y., Balazinska, M., Cetintemel, U., Cherniack, M., Hwang, J.-H., Lindner, W., Maskey, A.S., Rasin, A., Ryvkina, E., Tatbul, N., Xing, Y., Zdonik, S.: The design of the borealis stream processing engine. In: CIDR 2005, Asilomar, CA, January 2005
4. Ahmad, Y., Çetintemel, U.: Network-aware query processing for stream-based applications. In: VLDB 2004, vol. 30, pp. 456–467 (2004)
5. Andrade, H., Gedik, B., Wu, K., Yu, P.S.: Scale-up strategies for processing high-rate data streams in system S. In: ICDE 2009
6. Arasu, A., Cherniack, M., Galvez, E., Maier, D., Maskey, A., Ryvkina, E., Stonebraker, M., Tibbetts, R.: Linear road: a stream data management benchmark. In VLDB 2004
7. Carney, D., Çetintemel, U., Cherniack, M., Convey, C., Lee, S., Seidman, G., Stonebraker, M., Tatbul, N., Zdonik, S.: Monitoring streams: a new class of data management applications. In: VLDB 2002, pp. 215–226 (2002)
8. Castro Fernandez, R., Migliavacca, M., Kalyvianaki, E., Pietzuch, P.: Integrating scale out and fault tolerance in stream processing using operator state management. In: SIGMOD 2013, pp. 725–736. ACM, New York (2013)
9. Cormen, T.H., Stein, C., Rivest, R.L., Leiserson, C.E.: Introduction to Algorithms, 3rd edn. The MIT Press, Cambridge (2009)
10. Dean, J., Ghemawat, S.: Mapreduce: simplified data processing on large clusters. In: OSDI 2004, vol. 6. USENIX Association, Berkeley (2004)
11. DeWitt, D., Gray, J.: Parallel database systems: the future of high performance database systems. *Commun. ACM* **35**(6), 85–98 (1992)
12. Gedik, B., Schneider, S., Hirzel, M., Wu, K.-L.: Elastic scaling for data stream processing. *IEEE Trans. Parallel Distrib. Syst.* **25**, 1447–1463 (2010)
13. Graefe, G.: Encapsulation of parallelism in the volcano query processing system. In: SIGMOD 1990, pp. 102–111. ACM (1990)
14. Gulisano, V., Jimenez-Peris, R., Patino-Martinez, M., Valduriez, P.: StreamCloud: a large scale data streaming system. In: ICDCS 2010, pp. 126–137 (2010)
15. Khandekar, R., Hildrum, K., Parekh, S., Rajan, D., Wolf, J., Wu, K.-L., Andrade, H., Gedik, B.: COLA: optimizing stream processing applications via graph partitioning. In: Bacon, J.M., Cooper, B.F. (eds.) *Middleware 2009*. LNCS, vol. 5896, pp. 308–327. Springer, Heidelberg (2009). doi:[10.1007/978-3-642-10445-9_16](https://doi.org/10.1007/978-3-642-10445-9_16)

16. Lakshmanan, G.T., Li, Y., Strom, R.: Placement strategies for internet-scale data stream systems. *IEEE Internet Comput.* **12**(6), 50–60 (2008)
17. Motwani, R., Widom, J., et al.: Query processing, resource management, and approximation in a data stream management system. In: *CIDR 2003*, pp. 245–256, January 2003
18. Neumeyer, L., Robbins, B., Nair, A., Kesari, A.: S4: distributed stream computing platform. In: *ICDMW 2010*, pp. 170–177. IEEE Computer Society, Washington, DC (2010)
19. Pietzuch, P., Ledlie, J., Shneidman, J., Roussopoulos, M., Welsh, M., Seltzer, M.: Network-aware operator placement for stream-processing systems. In: *ICDE 2006*. IEEE (2006)
20. Schneider, S., Andrade, H., Gedik, B., Biem, A., Wu, K.-L.: Elastic scaling of data parallel operators in stream processing. In: *IPDPS*, pp. 1–12 (2009)
21. Schneider, S., Hirzel, M., Gedik, B., Wu, K.-L.: Auto-parallelizing stateful distributed streaming applications. In: *PACT 2012*, pp. 53–64. ACM, New York (2012)
22. Shah, M.A., Chandrasekaran, S., Hellerstein, J.M., Franklin, M.J.: Flux: an adaptive partitioning operator for continuous query systems. In: *ICDE*, pp. 25–36 (2002)
23. Wu, S., Kumar, V., Wu, K.-L., Ooi, B.C.: Parallelizing stateful operators in a distributed stream processing system: how, should you and how much? In: *DEBS 2012*, pp. 278–289 (2012)
24. Xing, Y., Hwang, J.-H., Çetintemel, U., Zdonik, S.: Providing resiliency to load variations in distributed stream processing. In: *VLDB 2006*, pp. 775–786. VLDB Endowment (2006)
25. Xing, Y., Zdonik, S., Hwang, J.-H.: Dynamic load distribution in the borealis stream processor. In: *ICDE 2005*, pp. 791–802. IEEE Computer Society (2005)