

Load Balancing for Skewed Streams on Heterogeneous Clusters

Muhammad Anis Uddin Nasir^{s#}, Hiroshi Horii^s, Marco Serafini^{*}, Nicolas Kourtellis[‡]
Rudy Raymond^s, Sarunas Girdzijauskas[#], Takayuki Osogami^s

[#]Royal Institute of Technology, Sweden ^sIBM Research Tokyo, Japan ^{*}Qatar Computing Research Institute [‡]Telefonica Research
anisu@kth.se, horii@jp.ibm.com, mserafini@qf.org.qa, nicolas.kourtellis@telefonica.com
rudyhar@jp.ibm.com, sarunasg@kth.se, osogami@jp.ibm.com

Abstract—Streaming applications frequently encounter skewed workloads and execute on heterogeneous clusters. Optimal resource utilization in such adverse conditions becomes a challenge, as it requires inferring the resource capacities and input distribution at run time. In this paper, we tackle the aforementioned challenges by modeling them as a load balancing problem. We propose a novel partitioning strategy called *Consistent Grouping* (CG), which enables each processing element instance (PEI) to process the workload according to its capacity. The main idea behind CG is the notion of small, equal-sized virtual workers at the sources, which are assigned to workers based on their capacities. We provide a theoretical analysis of the proposed algorithm and show via extensive empirical evaluation that our proposed scheme outperforms the state-of-the-art approaches, like key grouping. In particular, CG achieves 3.44x better performance in terms of latency compared to key grouping.

I. INTRODUCTION

Distributed stream processing engines (DSPES) have recently gained much attention due to their ability to process huge volumes of data with very low latency on clusters of commodity hardware. DSPES enable processing information that is produced at a very fast rate in a variety of contexts, such as IoT applications, software logs, and social networks. For example, Twitter users generate more than 380 million tweets per day¹ and Facebook users upload more than 300 million photos per day².

Streaming applications are represented by directed acyclic graphs (DAGs), where vertices are called *processing elements* (PEs) and represent operators, and edges are called *streams* and represent the data flowing from one PE to the next. For scalability, streams are partitioned into sub-streams and processed in parallel on replicas of PEs called *processing element instances* (PEIs).

Applications of DSPES, especially in data mining and machine learning, typically require accumulating state across the stream by grouping the data on common fields [4, 5]. Akin to MapReduce, this grouping in DSPES is usually called *key grouping* (KG) and is implemented using hashing [27]. KG allows each source PEI to route each message solely via its key, without needing to keep any state or to coordinate among PEIs. However, KG is unaware of the underlying skewness in the input streams [22], which causes a few PEIs to sustain a significantly higher load than others, as demonstrated in Figure 1 with a toy example. This sub-optimal load balancing leads to poor resource utilization and inefficiency.

The problem is further complicated when the underlying resources are heterogeneous [21, 33] or changing over time

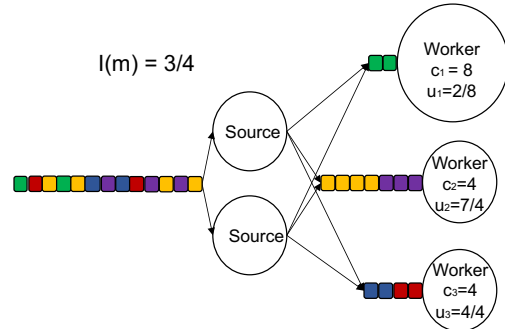


FIG. 1: Example showing that key grouping generates imbalance in the presence of a heterogeneous cluster. The capacity and the resource utilization of the i -th worker is represented by c_i and u_i respectively. Each key ($j \in \mathcal{K}$) is represented with different color box. Imbalance $I(m)$ is the difference between the maximum and the average resource utilization (see section IV for details).

[40, 35]. For various commercial enterprises, the resources available for stream mining consist of dedicated machines, private clouds, bare metal, virtualized data centers and commodity hardware. For streaming applications, the heterogeneity is often invisible to the upstream PEIs and requires inferring the resource capacities in order to generate a fair assignment of the tasks to the downstream PEIs. However, gathering statistics and finding optimal placement often leads to bottlenecks, while at the same time microsecond latencies are desired [17].

Alternatively, stateless streaming applications, like interaction with external data sources, employ *shuffle grouping* (SG) to break down the stream load equally to each of the PEIs, i.e., by sending a message to a new PEI in cyclic order, irrespective of its key. SG allows each source PEI to send equal number of messages to each downstream PEI, without the need to keep any state or to coordinate among PEIs. However, similarly to KG, SG is unaware of the heterogeneity in the cluster, which can cause some PEIs to sustain unpredictably higher load than others. Further, SG typically requires more memory to express stateful computations [27, 19].

In this present work, we study the load balancing problem for a streaming engine running on a heterogeneous cluster and processing non-uniform workload. To the best of our knowledge, we are the first to address both challenges together. We envision a light-weight and fair key grouping strategy for both stateless and stateful streaming applications. Moreover, this strategy must limit the number of workers processing each key, which is analogous to reducing the memory footprint and aggregation cost for the stateful computation [27, 19]. Towards this goal, we propose a novel grouping strategy

¹<http://www.internetlivestats.com/twitter-statistics/>

²<http://www.businessinsider.com/facebook-350-million-photos-each-day-2013-9>

called *Consistent Grouping* (CG), which handles both the potential skewness in input data distribution, as well as the heterogeneity in resources in DSPES. CG borrows the concept of virtual workers from the traditional consistent hashing [12, 13] and employs rebalancing to achieve fair assignment, similar to [34, 11, 3, 7, 35]. In summary, our work makes the following contributions:

- We propose a novel grouping scheme called Consistent Grouping to improve the scalability for DSPES running on heterogeneous clusters and processing skewed workload.
- We provide a theoretical analysis of the proposed scheme and show the effectiveness of the proposed scheme via extensive empirical evaluation on synthetic and real-world datasets. In particular, CG achieves bounded imbalance and generates almost optimal memory footprint.
- We measure the impact of CG on a real deployment on Apache Storm. Compared to key grouping, it improves the throughput of an example application on real-world datasets by up to 2x, reduces the latency by 3.44x.

II. OVERVIEW OF THE APPROACH

Consistent grouping relies on the concept of virtual workers and allows variable number of *virtual workers* for each PEI. The main idea behind CG is to assign the input stream to the virtual workers in a way that each virtual worker receives approximately the same number of messages. Later, these virtual workers are assigned to the actual workers based on their capacity. We refer to downstream PEIs as workers and to upstream PEIs as sources throughout the paper. Similar approaches have been considered in the past in the context of distributed hash tables [12, 13]. CG allows an assignment of tasks to PEIs based on the capacity of the PEIs. Thus, the powerful PEIs are assigned more work compared to less powerful PEIs. Next, we provide an overview of CG’s components.

First, we propose a novel strategy called *power of random choices* (PoRC), which assigns the incoming messages to a set of equal sized virtual workers. The basic idea behind this scheme is to introduce the notion of capacity for the virtual workers. In particular, we set the capacity of each virtual worker to the average load $\times (1+\epsilon)$, for some parameter ϵ . Note that the capacity is calculated at run time using the average load. Given a sequence of virtual workers for each key, PoRC maps a key to the first virtual worker with a spare capacity. PoRC allows the heavy keys to spread across the other virtual workers, thus reducing the memory footprint and the aggregation cost. The ϵ parameter in the algorithm provides the trade off between the imbalance and memory footprint.

Second, CG takes a radically new approach towards load balancing and allows PEIs to decide their workload based on their capacities. We call this component as *worker delegation*. Each worker monitors its workload and sends a binary signal (increase or decrease workload) to the sources in case it experiences excessive workload. This simple modification changes the distributed load balancing problem to a local decision problem, where each PEI can choose its share of workload based on its current capacity. Moreover, worker delegation provides the flexibility to implement various application-specific requirements at each PEI. The sources react to the signals by moving virtual workers from one PEI to another. Note that it is required that sources receive the signal and operate in a consistent manner, performing the same routing of messages.

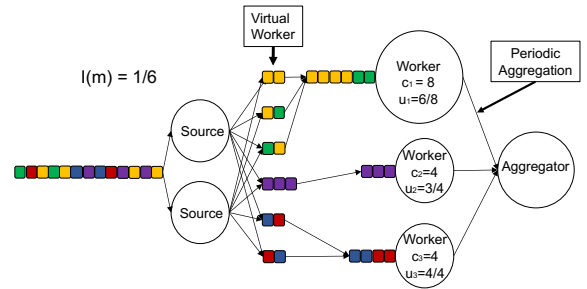


FIG. 2: Example showing that consistent grouping improves the imbalance in the presence of heterogeneous cluster, compared to key grouping. The capacity and the resource utilization of the i -th worker is represented by c_i and u_i respectively. Also, each key ($j \in \mathcal{K}$) is represented with different color box. Imbalance $I(m)$ is the difference between the maximum and the average resource utilization.

Such an operation might negatively impact the performance of a streaming application, as it requires one-to-many (from one worker to all sources) broadcast messages across the network. To overcome this challenge, we relax the consistency constraint in the DAG and allow sources to be eventually consistent. Specifically, we propose *piggybacking* that allows encoding the binary signals along with the acknowledgment message to avoid extra communication overhead.

Lastly, CG ensures that each message is processed in a consistent manner by discarding the message migration phase. When a source receives a request to change (increase or decrease) the workload, CG relocates virtual workers assigned to the overloaded worker, thus, only affecting the future routing of the messages. CG follows the same programming primitive as Partial Key Grouping (PKG) [27] for stream partitioning, supporting both stateless and stateful map-reduce like applications. We propose *periodic aggregation* to support map-reduce like stateful applications, which leverages the existing DAG and imposes a very low-overhead in the stream application. Figure 2 provides an example using CG for the DAG in Figure 1.

III. BACKGROUND ON STREAM PARTITIONING

Load Balancing is one of the very well-studied problems in distributed systems. Also, it is very extensively studied in theoretical computer science [25]. Next, we provide a discussion on various ways load balancing has been addressed in distributed systems, as well as state-of-art partitioning strategies to assign load to workers in such systems.

A. Load Balancing in Distributed Systems

In graph processing systems, load balancing is often found along with balancing graph partitioning, where the goal often is to minimize edge-cut between different partitions [14, 23]. Further, several systems have been proposed specifically to solve the load balancing problem, e.g., Mizan [20], GPS [32], and others. Most of these systems perform dynamic load rebalancing at runtime via vertex migration [39].

Load balancing and scheduling often appears in a similar context in map-reduce like systems, where the goal is to schedule the jobs to set of machines in order to maximize the resource utilization [15, 37]. Sparrow [29] is a stateless distributed job scheduler that exploits a variant of the power of two choices [30]. Ahmad et al. [1] improves the load balance

for map-reduce in heterogenous environment by monitoring and scheduling the jobs based on communication patterns.

Dynamic Load balancing in database systems is often implemented using rebalancing, similar to all the other systems [31]. Also, online load migration is effective for elasticity in the database systems [35, 36]. Lastly, dynamic load balancing is considered in the context of web servers [6], GPU [8], and many others.

B. Existing Stream Partitioning Functions

Messages are sent between PEIs by exchanging messages over the network. Several primitives are offered by DSPEs for sources to partition the stream, i.e., to route messages to different workers.

Key Grouping (KG). This partitioning ensures the messages with the same key are handled by the same PEI (analogous to MapReduce). It is usually implemented through hashing. KG is the perfect choice for *stateful* operators. It allows each source PEI to route each message solely via its key, without the need to keep any state or to coordinate among PEIs. However, KG does not take into account the underlying skewness in the input distribution, which causes a few PEIs to sustain a significantly higher load than others. This suboptimal load balancing leads to poor resource utilization and inefficiency.

Partial Key Grouping (PKG). PKG [27, 28, 26] adapts to the traditional power of two choices for load balancing in map-reduce like streaming operators. PKG guarantees nearly perfect load balance in the presence of bounded skew using two novel schemes: key splitting and local load estimation. The local load estimation enables each source to predict the load of workers leveraging the past history. However, similar to KG, PKG assumes that each worker has the same resources and the service time for the messages follows a uniform distribution, which is a strong assumption of many real-world use cases.

Shuffle Grouping (SG). This partitioning forwards messages typically in a round-robin fashion. It provides excellent load balance by assigning an almost equal number of messages to each PEI. However, no guarantee is made on the partitioning of the key space, as each occurrence of a key can be assigned to any PEI. It is the perfect choice for *stateless* operators. However, with *stateful* operators one has to handle, store and aggregate multiple partial results for the same key, thus incurring additional memory and communication costs.

C. Consistent Hashing

Consistent Hashing is a special form of a hash function that requires minimal changes as the range of the function changes [18]. This strategy solves the assignment problem by systematically producing a random allocation. It relies on a standard hash function that maps both messages and workers into unit-size circular ID space, i.e., $[0, 1) \subseteq \mathbb{R}$. Further, each task is assigned to the first worker that is encountered by moving in the clockwise direction on the unit circle. Consistent Hashing provides load balancing guarantees across the set of workers. Assuming n are the number of available workers, and given that the load on a node is proportional to the size of the interval it owns, no worker owns more than $O\left(\frac{\log n}{n}\right)$ of the interval (to which each task is mapped), with high probability [18].

One common solution to improve the load balance is to introduce *virtual workers*, which are copies of workers,

corresponding to points in the circle. Whenever, a new worker is added, a fixed number of virtual workers is also created in the circle. As each worker is responsible for an interval on the unit circle, creating virtual workers spreads the workload for each worker across the unit circle.

Similar to other stream partitioning functions, consistent hashing does not take into account neither the heterogeneity in the cluster or the skewness in the input stream, which restricts its immediate applicability in the streaming context. A way to deal with both heterogeneity and skewness is to employ hash space adjustment for consistent hashing [16]. Such schemes require global knowledge of the tasks assigned to each worker to adjust the hash space for the workers, i.e., movement of tasks from the overloaded worker to the least loaded worker. Even though such schemes provide efficient results in terms of load balance, their applicability in stream context incurs additional overhead due to many-to-many communication across workers. On the other hand, if implemented without global information, these schemes may produce unpredictable imbalance due to random task movement across workers.

Consistent Hashing with Bounded Load (CH). Independent from our work, Mirrokni et al. [24] proposed a novel version of consistent hashing scheme that provides a constant bound on the load of the maximum loaded worker. The basic idea behind their scheme is to introduce the notion of capacity for each worker. In particular, set the capacity of each bin to the average load times $(1 + \epsilon)$, for some parameter ϵ . Further, the tasks are assigned to workers in the clockwise direction with spare capacity.

D. Other Approaches

Power of Two Choices (PoTC). PoTC achieves near perfect load balance by first selecting two bins uniformly at random and later assigning the message to the least loaded of the two bins. For PoTC, the load of each bin is solely based on the number of messages. Using PoTC, each key might be assigned to any of the workers. Therefore, the memory requirement in worst case is proportional to the number of workers, i.e., every key appearing on all the workers.

Rebalancing. Another way to achieve fair assignment is to leverage *rebalancing* [34, 3, 7, 35, 10]. Once load imbalance is detected, the system activates a rebalancing routine that moves some of the messages and the state associated with them, away from an overloaded worker. While this solution is easy to understand, its applicability in the streaming context requires answering challenging questions: How to identify the imbalance and how to plan the migration. The answers to these questions are often application-specific, as they involve a trade-off between imbalance and rebalancing cost that depends on the size of the state to migrate. For these reasons, rebalancing creates a difficult engineering challenge, which we address in our paper.

IV. PRELIMINARIES & PROBLEM DEFINITION

This section introduces the preliminaries that are used in the rest of the paper.

We consider a DSPE running on a cluster of machines that communicate by exchanging messages following the flow of a DAG. For scalability, streams are partitioned into sub-streams and processed in parallel on a replica of the PE called *processing element instance* (PEI). Load balancing across the whole DAG is achieved by balancing along each edge independently.

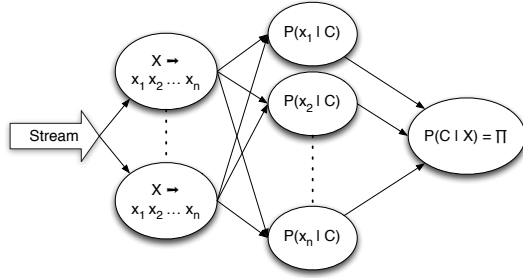


FIG. 3: Naïve Bayes implemented via key grouping (KG).

Each edge represents a single stream of data, along with its partitioning scheme. Given a stream under consideration, let \mathcal{S} be the set of sources, \mathcal{W} be the set of workers, and their sizes be $|\mathcal{S}| = s$ and $|\mathcal{W}| = n$.

Each PEI $w \in \mathcal{W}$ is deployed on a machine with a limited capacity $c_w \in \mathcal{C}$. For simplicity, we assume that there is a single important resource on which machines are constrained, such as storage and processing. Moreover, each PEI ($w \in \mathcal{W}$) has an unbounded input queue (\mathcal{Q}_w).

The input to the DAG is a sequence of messages $z = \langle i, j, v, t_i \rangle$ where i is the identifier, $j \in \mathcal{K}$ is the message key, v is the value, and t_i is the timestamp at which the message is received. The messages are presented to the engine in ascending order by timestamp. Upon receiving a message with key $j \in \mathcal{K}$, we need to decide its placement among the workers. We assume one message arrives per unit of time.

We employ queuing theory as the cost model to define the delay and the overhead at each worker. In the model, a sequence of messages arrives at a worker $w \in \mathcal{W}$. If the worker is occupied, the new message remains in the queue until it can be served. After the message is processed, it leaves the system. We represent the finish time for a message i using ϕ_i . The difference between the arrival time and the ϕ_i represents the latency of executing the message.

We define a *partitioning* function $\mathcal{H} : \mathcal{K} \rightarrow \mathcal{W}$, which maps each message into one of the PEIs. This function identifies the PEI responsible for processing the message. Each PEI is associated with one or more keys. The goal of the partitioning function is to generate an assignment of messages to the set of workers in a way that average waiting time is minimized.

We define the *load* of a worker using the number of messages that are assigned to the worker at time t :

$$L_w(t) = |\{i : (\mathcal{H}(i, j) = w) \wedge (t_i < t)\}|, \text{ for } w \in \mathcal{W}$$

Also, we define *normalized load* at time t as the ratio between the *load* and the capacity of the worker.

$$\mathcal{U}_w(t) = \frac{L_w(t)}{c_w}$$

We use a definition of *imbalance* similar to others in the literature (e.g., Flux [34] and PKG [27]). We define *imbalance* at time t as the difference between the maximum and the average normalized load:

$$I(t) = \max_w \{\mathcal{U}_w(t)\} - \text{avg}_w \{\mathcal{U}_w(t)\}, \quad w \in \mathcal{W}.$$

Further, the memory footprint of a worker w is the number of unique keys assigned to the worker:

$$M_w(t) = |\{k_i : (\mathcal{H}(i, j) = w) \wedge (\phi_i < t)\}|, \text{ for } w \in \mathcal{W}$$

Problem. Given the definition of imbalance, we consider the following problem in this paper.

Problem 4.1: Given a stream of messages drawn from a heavy-tailed distribution \mathcal{K} and a set of workers $w \in \mathcal{W}$ with capacities $c_w \in \mathcal{C}$, find a partitioning function \mathcal{H} that minimizes memory footprint while keeping the imbalance ($I(t)$) bounded by a constant factor at any time instance t .

Memory Cost. One simple solution to address problem 4.1 is to employ round robin assignment as in SG, which provides an imbalance of at most one in case of a homogenous cluster. This load balance comes at the cost of memory, as messages with the same key might end up on all the workers. Also, the round robin assignment produces a higher aggregation cost [27, 28, 19], which represents the communication cost for accumulating the partial results from the set of workers.

Example. To make the discussion more concrete, we introduce a simple application that will be our running example: the *naïve Bayes classifier*. A naïve Bayes classifier is a probabilistic model that assumes independence of features in the data (hence the naïve). It estimates the probability of a class C given a feature vector X by using Bayes' theorem:

$$P(C|X) = \frac{P(X|C)P(C)}{P(X)}.$$

The answer given by the classifier is then the class with maximum likelihood

$$C^* = \arg \max_C P(C|X).$$

Given that features are assumed independent, the joint probability of the features is the product of the probability of each feature. Also, we are only interested in the class that maximizes the likelihood, so we can omit $P(X)$ from the maximization as it is constant. The class probability is proportional to the product

$$P(C|X) \propto \prod_{x_i \in X} P(x_i|C)P(C),$$

which reduces the problem to estimating the probability of each feature value x_i given a class C , and a prior for each class C . In practice, the classifier estimates the probabilities by counting the frequency of co-occurrence of each feature and class value. Therefore, it can be implemented by a set of counters, one for each pair of feature value and class value.

V. SOLUTION PRIMITIVES

In this section, we discuss our solution and its various components. Given a set of sources and a set of workers, the goal is to design a grouping strategy that is capable of assigning the messages to the workers proportionally to their capacity, while dealing with the messages' embedded skew.

Overview. In our work, we propose a novel grouping scheme called consistent grouping (CG). Our scheme borrows the concept of virtual workers from the traditional consistent hashing [12, 13] and employs rebalancing to achieve fair assignment, similar to [34, 11, 3, 7, 35]. CG allows variable number of *virtual workers* for each PEI. The main idea behind CG is to assign the input stream to the virtual workers in a way that each virtual worker has approximately equal number of messages. Later, these virtual workers are assigned to the

actual workers based on their capacity. One of the challenges is to bound the load of each virtual worker, as it implies that moving a virtual worker from one worker to another actually increases the receiving worker’s load. For this, we propose a novel grouping strategy called *power of random choices* (PoRC) that is capable of providing bounded imbalance while keeping the memory cost low. Further, we propose three efficient schemes within CG: *worker delegation*, *piggybacking* and *periodic aggregation*, which enable efficient integration of our proposed scheme into standard DSPEs. Consistent grouping follows the same programming primitive as PKG for stream partitioning. We refer to [27] for the examples of common data mining algorithms that benefit from CG.

A. Power of Random Choices

PoRC assigns the incoming messages to the set of virtual workers in a way that the imbalance is bounded and the overall memory footprint of the keys on the virtual workers is low. The basic idea behind PoRC is to introduce the notion of continuous capacity, which is a function of average load. In particular, we set the capacity of each virtual worker to the average load times $(1 + \epsilon)$, for some parameter ϵ . Note that the definition of capacity is based on the average load, rather than a hard constraint. Given a sequence of virtual workers for a key, PoRC maps the key to the first virtual worker with the spare capacity. The sequence of virtual workers for a key are produced by using a single hash function and concatenating the *salt* with the key to produce a new assignment³. We refer to the first virtual worker in the sequence as the *principal virtual worker*. The rationale behind this approach is that the heavy keys in the skewed input distribution overload their principal worker. Therefore, we allow the heavy keys to spread across the other virtual workers, which reduces the memory footprint compared to other schemes, e.g., round robin. The ϵ parameter in PoRC provides the trade off between the imbalance and memory footprint. Algorithm 1 provides the pseudocode. PoRC provides an efficient and generalized solution for the fundamental problem of load balancing for the skewed stream in streaming settings, while minimizing the memory footprint. In our work, we adapt PoRC for fair load balancing for streaming applications, which shows its effectiveness and applicability.

Algorithm 1 Pseudocode for Power of Random Choices.

Require: key, hash-function, #messages, #workers, load-vector, imbalance-factor
Ensure: $S^* \in \{1 \dots n\}$
1: **procedure** GETWORKER($j, \mathcal{H}, m_t, n, load, \epsilon$)
2: $salt \leftarrow 1$
3: $S^* \leftarrow \mathcal{H}(j + salt)$
4: **while** $(load[S^*] \geq (1 + \epsilon) \frac{m_t}{n})$ **do**
5: $salt \leftarrow salt + 1$
6: $S^* \leftarrow \mathcal{H}(j + salt)$
7: $load[S^*] \leftarrow load[S^*] + 1$
8: **return** S^*

Discussion. To show the effectiveness of PoRC, we compare its performance with KG, PKG, PoTC, SG, and CH [24] in terms of imbalance and memory footprint (see section III-A for the description of these schemes). We leverage a zipf-based dataset with different skews for this experiment (see Section VII for the description of the dataset). The top row of Figure 4 reports

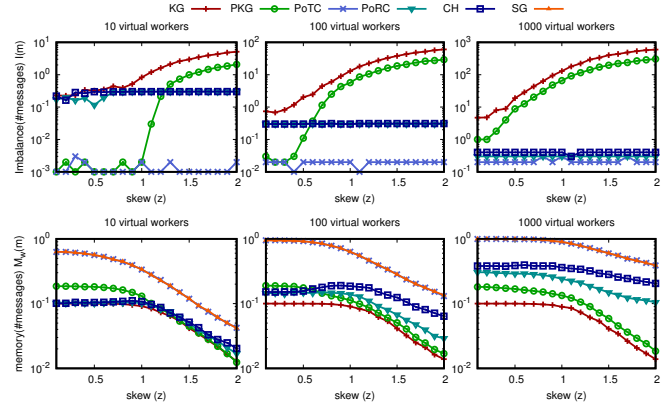


FIG. 4: Normalized imbalance and memory overhead for different schemes using zipf distribution with different skew and number of virtual workers.

the imbalance for different schemes for different number of virtual workers, i.e., 10, 100, and 1000. Results show that both key grouping and partial key grouping generate high imbalance as the skew and the number of virtual workers increase. However, the other schemes perform fairly well in terms of imbalance. Additionally, we report the memory overhead for all the schemes in Figure 4. The memory cost is calculated using the total number of unique keys that appear at each virtual worker. Results verify our claim that load balance is achieved at the cost of memory.

B. Consistent Grouping

We propose a novel grouping strategy called *Consistent Grouping* (CG), inspired by consistent hashing. CG borrows the concept of virtual workers from traditional consistent hashing and allows variable number of virtual workers for each PEI [12, 13]. It is a dynamic grouping strategy that is capable of handling both the heterogeneity in the resources and the variability in the input stream at runtime. CG achieves its goal by allowing the powerful workers to acquire additional virtual workers, which leads to ‘stealing’ work from the other workers. Moreover, it allows overloaded workers to gracefully revoke some of their existing virtual workers, which is equivalent to giving up on some of the allocated work.

CG is a lightweight and distributed scheme that allows assignment of messages to the workers in a streaming fashion. Moreover, it leverages PoRC for assignment of keys to each virtual worker in a balanced manner, which allows it to bound the load of each virtual worker. CG is able to balance the load across workers based on their capacities, which allows the DSPE to operate under realistic scenarios like heterogeneous clusters and variable workloads.

Time Slot. We introduce the notion of *time slot* (t_0), which represents the minimum monitoring time period for a PEI. t_0 is an administrative preference that can be determined based on workload traffic patterns. If workloads are expected to change on an hourly basis, setting t_0 on the order of minutes will typically suffice. For slower changing workloads t_0 can be set to an hour. Time slot guarantees that workers have enough sample of the input stream to predict their workload.

Similar to consistent hashing, CG initializes with the same number of virtual workers for each worker, i.e., $O(\log n)$. CG manages a unit-size circular ID space, i.e., $[0, 1) \subseteq \mathbb{R}$, and

³<https://datarus.wordpress.com/2015/05/04/fighting-the-skew-in-spark/>

maps the virtual workers and keys on the unit-size ID space. We would like a scheme that is capable of monitoring the load at each worker throughout the lifetime of a streaming application and adjust the load according to the available capacity of the workers. In doing so, we introduce a novel scheme called *pairing virtual worker*.

Pairing virtual workers. The load of a worker equals to the sum of loads of the assigned virtual workers. Further, the load of each virtual worker equals to the load that is induced by the mapped messages. Ideally, we would like to assign one of the virtual workers from the overloaded worker to one of idle workers. However, it is not trivial until this point on how one can achieve such an assignment. To enable such an assignment, we propose to maintain two first-come-first-serve (FCFS) queues: *idle* and *busy*. These queues maintain the list of idle and busy workers in the DSPE and allow CG to pair any removal and addition with the opposite to balance the number of virtual workers throughout the execution. For instance, when a worker is overloaded, it sends a message to the sources. Further, the source only removes the virtual workers of the corresponding worker if it is able to pair it with an addition on an idle worker. This simple scheme ensures that the number of virtual workers in the system are balanced throughout the execution and the load of each virtual worker is bounded, which enables CG to perform fair assignment. Note that mapping the virtual workers with similar keys to the same worker might reduce the memory footprint. However, this requires maintaining all the unique keys in each virtual worker and each worker. Therefore, we opt for FCFS mapping of virtual workers to workers.

C. Integration in a DSPE

While consistent grouping is easy to understand, its applicability to the case of real world stream processing engines is not trivial. We package CG with few efficient strategies that enable its applicability in a variety of DSPEs.

Worker Delegation. First, we propose an efficient scheme called *worker delegation*, which pushes the load balancing problem to the workers and allows them to decide their workload based on their capacity. Each worker requires monitoring its workload and needs to take the decision based on their current workload and the available capacity. The decision can either be to increase the workload or to decrease the workload. The intuition behind this approach is that it is often the case that the cluster consists of a large number of workers and collecting the statistics periodically from the workers creates an additional overhead for a streaming application.

The worker delegation scheme allows the workers to interact with sources by sending binary signals: (1) increase workload and (2) decrease workload. Each worker monitors its workload and tries to keep the workload between two thresholds, i.e., if the workload exceeds the upper threshold, the worker sends a decrease signal to the sources, and if the workload is below the lower threshold, the worker sends the increase signal to the sources. This simple modification comes along with the benefit that it gives the flexibility to the workers to easily adapt to the complex application-specific requirements, i.e., processing, storage, service time and queue length.

Piggybacking. Each worker requires updating all the sources in case of experiencing undesirable (low or high) workload. Note that it is required that sources receive the signal and operate in a consistent manner, performing the same routing

of messages. Such deployment might negatively impact the performance of a streaming application, as it requires one-to-many broadcast messages across the network. To overcome this challenge, we propose to relax the consistency constraint in the DAG and allow operators to be eventually consistent. We propose to *encode* the binary signals from the workers along with the acknowledgment messages. During the execution, the sources only receive the signal from the worker as a response to its messages. This means that the worker might continue receiving the messages with the same key even after triggering the decision.

Periodic Aggregation. When the sources receive a request to increase the workload, they move one of the virtual worker from the overloaded worker to an idle worker. During the change of routing, we need to ensure that the messages that are pending in the queue of the workers must be processed in a consistent manner.

CG ensures that each message is processed in a consistent manner by discarding the message migration phase. Concretely, each worker processes the messages that are assigned to it, and any change in the routing only affects the messages that arrive later.

As a message with the same key might be forwarded to different workers, CG performs *periodic aggregation* of partial results from the workers to ensure that the state per key is consistent. Periodic aggregation leverages the same DAG and imposes a very low overhead in the stream application. Particularly, CG follows the same programming primitive as PKG for stream partitioning, supporting both stateless and stateful map-reduce like applications.

VI. ANALYSIS

We proceed to analyze the conditions under which CG achieves good load balance. Recall from Section IV that we have a set \mathcal{W} of n workers at our disposal. Each worker $w \in \mathcal{W}$ has a limited capacity, which is represented by $c_w \in \mathcal{C}$. Capacities are normalized so that the average capacity is $\frac{1}{n}$, that is $\sum_{w \in \mathcal{W}} c_w = 1$. We assume that they are ordered in decreasing order of capacity, i.e., $c_1 \geq c_2 \geq c_3 \dots \geq c_n$.

The input to the engine is a sequence of messages $z = \langle i, j, v, t_i \rangle$ where i is the identifier, $j \in \mathcal{K}$ is the message key, v is the value, and t_i is the timestamp at which the message is received. Upon receiving a message with key $j \in \mathcal{K}$, we need to decide its placement among the workers. We assume one message arrives per unit of time. The messages arrive in ascending order by timestamp.

Key distribution. We assume the existence of an underlying discrete distribution \mathcal{D} supported on \mathcal{K} from which keys are drawn, i.e., k_1, \dots, k_m is a sequence of m independent samples from \mathcal{D} ($m \gg n$). We represent the average arrival rate of messages as p_j and the cardinality of set \mathcal{K} as c , i.e., $c = |\mathcal{K}|$. We assume that they are ordered in decreasing order of average arrival rate, $p_1 \geq p_2 \dots \geq p_c$, and $\sum_{j \in \mathcal{K}} p_j = 1$. We model the load distribution as a zipf distribution with values of exponent z between 0 and 2.0. The probability mass function of the zipf distribution with z is

$$f(r, c, z) = \frac{1/r^z}{\sum_{x=1}^c (1/x^z)},$$

where r is the rank of each key, and m is the total number of elements.

Our goal is to design an algorithm to solve Problem 4.1. In the analysis of CG, we assume that t_0 represents the time slot which corresponds to the minimum time period that each worker waits after sending a signal to the workers. Also, as we are not considering elasticity, we assume that the system is well provisioned, i.e., $\frac{\sum_{j \in \mathcal{K}} P_j}{\sum_{w \in \mathcal{W}} c_w} < 1$.

A. Imbalance with Consistent Grouping

For simplification, we divide the analysis of CG into two parts: dividing the workload into small equal-sized virtual workers and assigning the virtual workers to workers based on their capacities. Assume that $\alpha > 1$ represents the number of virtual workers assigned to each worker at initial time. Then, for n heterogeneous workers, we have $\alpha \times n$ homogeneous virtual workers. Each virtual worker has the same capacity (hence, homogeneous), and the capacity is guaranteed to be at most the capacity of the worker with the lowest capacity. The sources do not know the capacity of each worker. However, since all virtual workers are homogeneous, the sources can balance the load of each worker by assigning equal number of messages to each virtual worker, and by keeping the number of virtual workers assigned to each worker proportional to its capacity.

1) *Chromatic Balls and Bins*: We model the first problem using the framework of balls and bins processes, where keys correspond to colors, messages to colored balls, and virtual workers to bins. Choose d independent hash functions $\mathcal{H}_1, \dots, \mathcal{H}_d: \mathcal{K} \rightarrow [\alpha n]$ uniformly at random. Define the Greedy- d scheme as follows: at time t , the t -th ball (whose color is k_t) is placed on the bin with minimum current load among $\mathcal{H}_1(k_t), \dots, \mathcal{H}_d(k_t)$, i.e., $P_t(k_t) = \operatorname{argmin}_{i \in \{\mathcal{H}_1(k_t), \dots, \mathcal{H}_d(k_t)\}} L_i(t)$. We define the *imbalance* as the difference between the maximum and the average load across the bins, at time t .

Observe that when $d = 1$, each ball color is assigned to a unique bin so no choice has to be made; this models hash-based key grouping. At the other extreme, when $d \gg n \ln n$, all n bins are valid choices, and we obtain shuffle grouping.

PKG [27] considers the case of $d = 2$, which is same as having two hash functions $\mathcal{H}_1(j)$ and $\mathcal{H}_2(j)$. The algorithm maps each key to the sub-stream assigned to the least loaded worker between the two possible choices, that is: $P_t(j) = \operatorname{argmin}_i (L_i(t) : \mathcal{H}_1(j) = i \vee \mathcal{H}_2(j) = i)$.

Lemma 6.1: Suppose we use n bins and let $m \geq n^2$. Assume a key distribution \mathcal{D} with maximum probability $p_1 \leq \frac{1}{5n}$. Then, the imbalance after m steps of the Greedy- d process is $O\left(\frac{\ln \ln n}{\ln d}\right)$, with high probability [27].

Observe that the imbalance in case of PKG is only guaranteed for the case when $p_1 \leq \frac{1}{5n}$. However, in the case when $p_1 > \frac{1}{5n}$, the imbalance grows proportional with the frequency of the most frequent key and number of workers.

Power of Two Choices (PoTC) [2] leverages two random numbers $\mathcal{R}1(m)$ and $\mathcal{R}2(m)$. The algorithm maps each message m to the sub-stream assigned to the least loaded worker between the two possible choices, that is: $P_t(k) = \operatorname{argmin}_i (L_i(t) : \mathcal{R}1(m) = i \vee \mathcal{R}2(m) = i)$. The above random numbers can be generated by using hash functions with messages as arguments. In this case, note that the PoTC is different from the PKG in the sense that two hashes are applied to the messages, rather than the keys. The procedure is identical to the standard Greedy- d process of Azar et al. [2], therefore the following bounds hold.

Lemma 6.2: Suppose we use n bins and let $m \geq n^2$. Then, the imbalance after m steps of the Greedy- d process is $O\left(\frac{\ln \ln n}{\ln d}\right)$, with high probability [2].

Note that these bounds can be generalized to the infinite process in which n balls leave the system in each time unit (one from each worker) and the number of balls entering the system is less than n . In such cases, the relative load remains the same, therefore the bound holds. PoRC generate imbalance that is bounded by the factor ϵ , i.e., $I(m) \leq \epsilon \cdot \left(\frac{m}{n}\right)$.

2) *Fair Bin Assignment*: Given that m messages are assigned to set of n workers using PoRC, our goal is to show that consistent grouping is able to perform fair assignment to messages to the workers over time. We achieve our goal by showing that consistent grouping reduces the imbalance $I(t)$ (if it exists) over time. To make the discussion more concrete, we define the notion of busy worker using a threshold $\theta_b > 1$. In particular, we say that a worker w is busy if the load $L_w \geq \theta_b \cdot c_w$. Similarly, we define the notion of idle worker using the threshold $\theta_i < 1$. We say that a worker w is idle if its load $L_w \leq \theta_i \cdot c_w$.

Assume that α represents the average number of virtual workers per worker, i.e., the total number of virtual workers equal $\alpha \times n$. Also, assume that α_w^* represents the optimal number of virtual workers for w -th worker, namely, $\alpha_w^* = c_w n \alpha$. Clearly, $\frac{\theta_i \cdot c_w}{\alpha_w^*} \leq \frac{1}{n \cdot \alpha} \leq \frac{\theta_b \cdot c_w}{\alpha_w^*}$.

Thanks to the load balancing mechanisms such as PKG or PoTC, each virtual bin is guaranteed to have load at most $1/(\alpha n) + \gamma$ with high probability, where γ denotes the imbalance factor of the load balancing mechanism used. For PKG and PoTC, the value of γ is at most $(\ln \ln \alpha n / (m \ln d))$ as implied by Lemmas 6.1 and 6.2 (notice that the denominator m is due to the normalization of the capacity in this paper). Therefore, the expected load of a worker w having α_w virtual workers is bounded above by

$$\mathbf{E}[L_w] \leq \alpha_w \cdot \left(\frac{1}{n\alpha} + \gamma\right)$$

Now, consider that the worker w is overloaded, i.e., $\mathbf{E}[L_w] \geq \theta_b \cdot c_w$. This implies:

$$\alpha_w \cdot \left(\frac{1}{n\alpha} + \gamma\right) \geq \theta_b \cdot c_w$$

We can rearrange the above equation to have:

$$\gamma \geq \frac{\theta_b \cdot c_w}{\alpha_w} - \frac{1}{n\alpha},$$

which implies that when the worker is overloaded, it must have an imbalance that is lower bounded by the above equation. However, such an imbalance is guaranteed to be small ϵ by the load balancing mechanism used, i.e., $\epsilon \leq (\ln \ln \alpha n / (m \ln d)) \ll 1/(\alpha n)$ for PoTC and PKG, when $m \geq n^2$.

Therefore, we know that for an overloaded worker, it must hold that:

$$\epsilon \geq \gamma \geq \frac{\theta_b \cdot c_w}{\alpha_w} - \frac{1}{n\alpha}$$

Now, by solving for α_w , we get:

$$\alpha_w \geq \frac{\theta_b \cdot c_w}{\frac{1}{n\alpha} + \epsilon} \geq \theta_b \cdot c_w n \alpha (1 - \epsilon n \alpha),$$

where we use the Bernoulli's inequality $(1 + \epsilon n \alpha)^{-1} \geq (1 - \epsilon n \alpha)$ to obtain the above second inequality.

Notice that the above inequality gives the lower bound on the number of virtual workers assigned to an overloaded worker. Since its optimal number of virtual workers is $\alpha_w^* = c_w n \alpha$, we can see that $\alpha_w / \alpha_w^* \geq \theta_b (1 - \epsilon n \alpha)$, which is close to 1 since $\epsilon \ll 1 / (\alpha n)$. This gives an interesting property that once we know a worker is overloaded, we can be sure that its number of virtual workers is close to the optimal allocation. Thus, the sources can probe the capacity of workers by assigning virtual workers (taken from overloaded workers) to workers that have not reported becoming overloaded, or if there is no such one, to those that reported becoming overloaded least recently. Also notice that by letting $\theta_b = (1 + \epsilon n \alpha)$, we can guarantee that the overloaded workers are having at least the optimal number of virtual workers they should have. However, when ϵ is large (due to bad load balancing mechanisms), or when αn is large (due to having many small virtual workers), $\theta_b > 1$ will become large. This will burden the overloaded workers because they can only broadcast the overloaded cases when the threshold $\theta_b \cdot c_w$ is surpassed. This illustrates the tradeoff of load balancing mechanisms, with small imbalance factor ϵ , and the right number of virtual workers (too many is not good) in our consistent grouping strategy.

B. Memory with Consistent Grouping

KG generates the optimal memory footprint by forwarding each key to exactly one worker. Similarly, PKG produces nearly optimal memory overhead by allowing at most two workers per key. On the other end, PoTC and SG might assign each key to all the workers in the worst case, producing the memory footprint proportional to the number of workers. Assume that X_i is a random variable representing the minimum number of bins required for a ball i with color k_i . Further, assume a random variable X representing the sum of number of bins required for the balls, i.e., $X = \sum_{i=1}^c X_i$. A trivial upper bounded for X in case of shuffle grouping is given by:

$$\mathbb{E}[X] = \sum_{i=1}^c \min(\lceil p_i \cdot m \rceil, n) \quad (1)$$

PoRC allows a tradeoff between imbalance and memory using the parameter ϵ . To analyze the memory footprint of PoRC, we answer a very simple question: *What is the probability that a key is replicated on all the workers?* For this to happen, the load of $n - 1$ workers should exceed by $(1 + \epsilon)$ of the average load. Only then a key is replicated on all the workers. However, for a sufficiently large value of ϵ , i.e., $\epsilon > \frac{1}{n-1}$ this can not happen. A trivial lower bound on the number of bins required for a ball i with color k_i is $\mathbb{E}[X_i] = \lceil \frac{p_i \cdot n}{(1 + \epsilon)} \rceil$. Then,

$$\mathbb{E}[X] = \sum_{i=1}^c \lceil \frac{p_i \cdot n}{(1 + \epsilon)} \rceil \quad (2)$$

This discussion provides the basic intuition on why the memory overhead of PoRC is lower than SG and PoTC. We plan to consider the detailed analysis in future work.

TABLE I: Summary of the datasets used in experiments. Note: Percentage of messages having the most frequent key (p_1).

Stream	Symbol	Messages	Keys	p_1 (%)
Wikipedia	WP	22M	2.9M	9.32
Twitter	TW	1.2G	31M	2.67
Zipf	ZF	10M	100k	$\propto \frac{1}{\sum x^{-z}}$

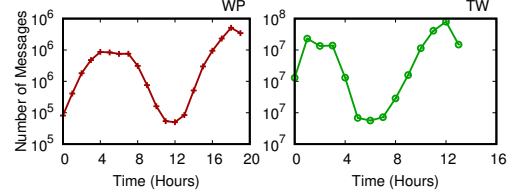


FIG. 5: Number of messages per hour for WP and TW datasets.

VII. EVALUATION

We assess the performance of our proposal by using both simulations and a real deployment. In so doing, we answer the following questions:

- Q1:** What is a good set of values for the parameters of CG?
- Q2:** How does CG perform compared to other schemes?
- Q3:** How does CG adapt to changes in input stream and resources?
- Q4:** What is the overall effect of CG on applications deployed on a real DSPE?

A. Experimental Setup

Datasets. Table I summarizes the datasets used. In particular, our goal is to be able to produce skewness in the input stream. We use two main real data streams, one from *Wikipedia* and one from *Twitter*. These datasets were chosen for their large size, different degree of skewness, and different set of applications in Web and online social network domains. The Wikipedia dataset (WP)⁴ is a log of the pages visited during a day in January 2008. Each visit is a message and the page's URL represents its key. The Twitter dataset (TW) is a sample of tweets crawled during July 2012. We split each tweet into words, which are used as the key for the message. Figure 5 reports the ingestion rate of the streams in terms of number of messages per hour. Lastly, we generate synthetic datasets with keys following Zipf distributions with exponent in the range $z = \{0.1, \dots, 2.0\}$ with 100k unique keys.

Simulation and Real Deployment. We process the datasets by simulating the DAG presented in Figure 3. The stream is composed of timestamped keys that are read by multiple independent sources (\mathcal{S}) via shuffle grouping, unless otherwise specified. The sources forward the received keys to the workers (\mathcal{W}) downstream. In our simulations we assume that the sources perform data extraction and transformation, while the workers perform data aggregation, which is the most computationally expensive part of the DAG. Thus, the workers are the bottleneck in the DAG and the focus for the load balancing. Note that for simulation, we ignore the network latency. The selected workloads represent a variety of streaming applications. In particular, any application that performs reduce-by-key or group-by operation follows a similar pattern.

⁴http://www.wikibench.eu/?page_id=60

TABLE II: Notation for the algorithms tested.

Symbol	Algorithm
KG	Key Grouping
PKG	Partial Key Grouping
PoTC	Power of Two Choices
PORC	Power of Random Choices
CH	Consistent Hashing with Bounded Load
SG	Shuffle Grouping
CG	Consistent Grouping

TABLE III: Metric used for evaluation of the algorithms.

Metric	Description
Memory Cost	Replication cost of the keys
Queue Length	Number of messages in the queue
Resource Utilization	Ratio between number of messages and capacity of worker.
Imbalance	Difference between the maximum and the average resource utilization.
Execution Latency	Difference between arrival and finish time.
Throughput	Number of messages processed per second.

Algorithms. Table II defines the notations used for the different algorithms. We use a 64-bit Murmur hash function for implementation of KG to minimize the probability of collisions. Unlike the algorithms in Table II, other related load balancing algorithms [34, 9, 38, 3, 7] require the DSPE to support operator migration. Many top DSPEs, such as Apache Storm, do not support migration. Thus, we omit these algorithms from the evaluation.

Metrics. Table III defines the metrics used for the evaluation of the performance of different algorithms.

Monitoring Performance. For CG, each worker requires monitoring its resource utilization that enables the fair message assignment. In case of simulations, we define the resource utilization as the ratio between the number of assigned messages and the capacity of a worker. We define the notion of idle and busy worker using the $\mathcal{U}_w(t) < 0.75 \cdot c_w$ and $\mathcal{U}_w(t) > 0.85 \cdot c_w$ thresholds respectively. For the real-world experiments, we suggest using the queue length as a parameter for monitoring the resource utilization. In particular, the resource utilization is defined by:

$$\mathcal{U}_w(t) = \frac{\#\text{tuples in the queue}}{\text{input queue capacity}} = \frac{L_w(t)}{c_w}$$

The choice of the parameter was motivated by its availability in the standard Apache Storm distribution (ver 1.0.2).

B. Experimental Results

Q1: In the first experiment, we simulate the CG scheme by varying the value of ϵ and fixing the number of sources to 1 and the number of workers to 10. Each worker is homogeneous and the number of virtual workers per worker are set to 10. We select the WP dataset and simulate CG for different values of ϵ . We leverage KG, SG and PORC for task assignment to the virtual workers. Figure 6 reports the imbalance and the memory overhead for the experiment. The results verify our claim that epsilon provides a trade-off between imbalance and memory. In particular, CG generates low imbalance at lower values of epsilon and produces low memory footprint for higher values of epsilon. Also, the experiment shows that CG is

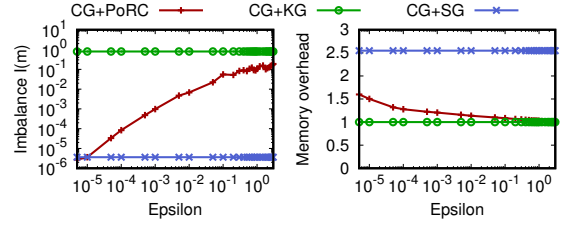


FIG. 6: Experiment reporting the imbalance and the memory overhead for different values of epsilon. The setup includes 10 workers, each having 10 virtual workers mapped to it.

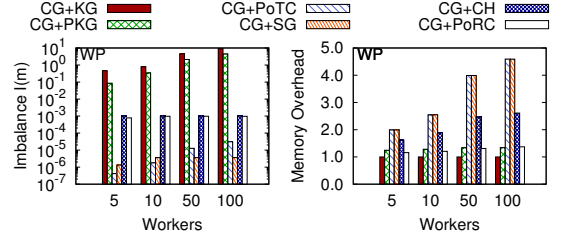


FIG. 7: Normalized imbalance and memory overhead comparing several assignment strategies along with consistent grouping on a homogeneous cluster with 5, 10, 50 and 100 workers using WP dataset. Each worker spawns 10 virtual workers and the $\epsilon = 0.01$ for CH and PORC.

able to interpolate well between the KG and SG schemes. Based on this experiment, we use the value of $\epsilon=0.01$ henceforth as it provides a middle ground between memory and imbalance.

Next, we analyze the allocation strategies, i.e., KG, PKG, PoTC, PORC, CH and SG. We simulate an experiment on a homogeneous cluster with 5, 10, 50 and 100 workers using the WP dataset. The number of virtual workers per worker are set to 10, i.e., equivalent to splitting the keys into 50, 100, 500 and 1000 bins. For CH and PORC, we set $\epsilon = 0.01$. Figure 7 shows the imbalance after the assignment of the streams. Results show that KG and PKG generate high imbalance, whereas PoTC and SG generate nearly perfect load balance. Both CH and PORC bound the imbalance close to a constant factor from the value of ϵ . The imbalance in case of KG and PKG grows linearly with the increase in the number of workers. This behavior is due to the fact that both these schemes restrict a single key to a constant number of workers. CH and PORC bound the imbalance upto a constant factor for each bin. PoTC and SG achieve near perfect imbalance by exploiting all the possible workers. Interestingly, PORC achieves bounded imbalance while keeping the memory footprint as low as PKG, as shown in Figure 7. In particular, PORC generates nearly perfect memory footprint and operates very close to KG. The gain in memory footprint depends on the distribution of the workload and the size of the deployment, and achieving gains in orders of magnitude is not always possible. Henceforth, we leverage PORC for the next experiments and analyze consistent grouping.

Q2: To answer this question, we compare the imbalance and the memory overhead of CG with KG, PKG, PoTC, CH and SG. We simulate the DAG using the WP dataset and report the value of imbalance measured at the end of the simulation. The cluster consists of different number of workers, i.e., 5, 10, 50 and 100 workers. Each experiment considers a cluster of homogeneous machines. For CG and CH, we set the value of epsilon equal to 0.01. Figure 8 reports the imbalance and

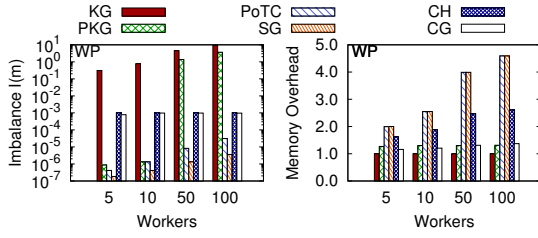


FIG. 8: Normalized imbalance and memory overhead comparing different grouping strategies on a homogeneous cluster with 5, 10, 50 and 100 workers using the WP dataset. Each worker spawns 10 virtual workers and the $\epsilon = 0.01$ for CH and PoTC.

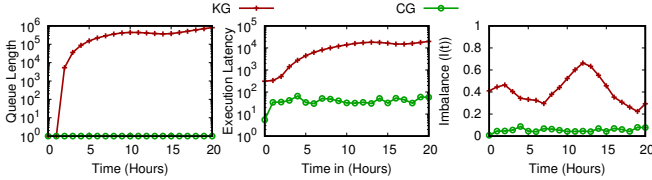


FIG. 9: Effect on queue length, execution latency and resource utilization on a homogeneous cluster with 10 workers for KG and CG using the WP dataset. Each worker spawns 10 virtual workers and the ϵ equals 0.01 for CG.

the memory overhead for different schemes (note the log scale). Results show that KG performs the worst in terms of the imbalance while generating the optimal memory footprint. PKG on the other hand provides nearly perfect imbalance and optimal memory footprint for smaller deployments, i.e., 5 and 10 workers. However, the imbalance grows as the number of workers increase. PoTC and SG provide very similar performance, i.e., provide nearly perfect imbalance and generate higher memory footprint. CH provides bounded imbalance and reduces the memory footprint compared to PoTC and SG. CG provides the bounded imbalance and improves the memory footprint compared to CH. This behavior is due to the fact that CG leverages randomness to redistribute the messages once the principal worker reaches the capacity, whereas CH always choose the next worker in the ring.

Additionally, we report the queue length, execution latency and the resource utilization among the workers by setting the capacity of the workers in a way that each worker operates at 80% of the capacity using shuffle grouping. We report each metric as a difference between the maximum and minimum value. Note that the difference between the maximum and minimum resource utilization represents the imbalance. Due to space restriction, we only report the results for 10 workers. For comparison, we also simulate and report KG and CG. Note that as PKG, PoTC and SG provide nearly perfect load balance, we do not report their results (the different between maximum and average queue length, execution latency and resource utilization equals 0). We simulate the WP dataset, set the value of ϵ equal to 0.01 and set the number of virtual workers per worker equal to 10 for CG. Figure 9 shows the results of the experiment over time. Results show that the difference between the maximum and minimum queue length and execution latency increases over time using KG, whereas CG keeps both queue length and execution latency very low. Also, the imbalance is high for KG, whereas CG keeps the imbalance to close to zero.

Next, we mimic the heterogeneity in the cluster by assuming a cluster consisting of n machines in which y machines are z

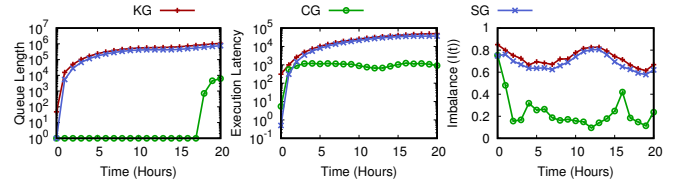


FIG. 10: Effect on queue length, execution latency and resource utilization due to heterogeneity in the cluster for KG, CG and SG.

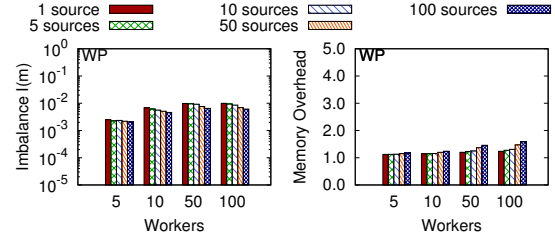


FIG. 11: Normalized imbalance and memory overhead on a homogeneous cluster with 5, 10, 50 and 100 workers with using 1, 10, 50, and 100 sources using WP dataset. Each worker spawns 10 virtual workers and the ϵ equals 0.01 for CG.

times more powerful than rest of the machines. In particular, we vary the value of z between 2 to 10 and vary the value of y between 1 to $n - 1$. For instance, when $y = 1$ and $z = 2$, a machine in a 10-machine cluster has twice the capacity than all the other nine machines. We simulate the KG, SG, CG for comparison and use the value of epsilon equal to 0.01. In case of CG, each worker is initialized with 10 virtual workers. We observe similar behavior in all the configurations and report only a single iteration with $y = 3$ and $z = 5$. Figure 10 reports the queue length, execution latency and resource utilization for the three approaches. Results show that queue length and execution latency grow for KG and SG. Similarly, the imbalance is pretty high for these approaches. On the other hand, CG provides the lower queue length and execution latency. Also, it keeps the imbalance close to zero. Note that there is a spike after 17 hours for queue length, which is due to the fact that we leverage the resource utilization as a metric to segregate between idle and busy workers.

Q3: Further, we evaluate the performance of CG by increasing the number of sources. In particular, we compare the performance of different deployments using 1, 10, 50, and 100 sources. For assignment of messages to sources, we use SG. Figure 11 reports the performance of CG in terms of imbalance and memory overhead. Results show that both imbalance and memory footprint almost remain the same on a log scale by both increasing the number of workers and number of sources. Therefore, we can conclude that CG is able to provide similar performance even under higher number of sources and workers.

Further, we study the behavior of CG on the number of virtual workers. We reuse the configuration for the experiment reported in Figure 10 and report the queue length, execution latency and resource utilization for CG, i.e., $y = 3$ and $z = 5$. We perform the experiment using number of virtual workers equal to 5, 10, 20, 50, 100 and 1000. Figure 12 shows that setting the number of virtual workers to a value of 5 does not provide desired results. This is due to the fact that there are not enough virtual workers to move around the workers. Similarly,

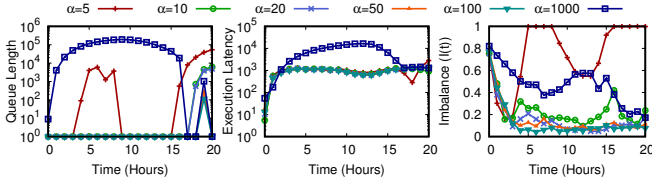


FIG. 12: Queue length, execution latency and resource utilization for different number of virtual workers.

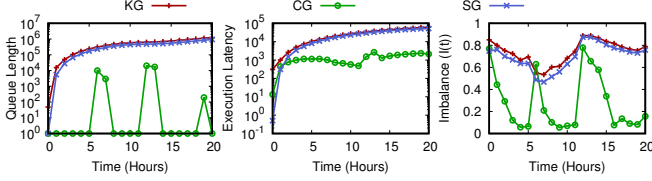


FIG. 13: Queue length, execution latency and resource utilization for when resources are changing over time. The resources change after processing 6M and 12M messages.

when the number of virtual workers are equal to 1000, the system takes longer time to converge, hence impacting the performance. Executions using 10 and 20 virtual workers provide similar performance. Lastly, the execution using 100 virtual workers generate the best results.

Next, we study the performance of CG by dynamically changing the resources over time. To initialize the resources, we reuse the configuration from the previous experiment and change the capacity of resources twice during the execution, i.e., after processing 6M and 12M messages. Concretely, we change the values of y and z (represented as $\{y, z\}$) after 6M and 12M messages to $\{5, 4\}$ and $\{2, 10\}$ respectively. We execute the experiment for 100 virtual workers and change the resources in a way that the sum of resources remains the same. Also, we report the results of KG and SG for comparison. Figure 13 reports the queue length, execution latency and resource utilization of the experiment. Results show that CG adapts very efficiently to the change in resources.

Q4: Lastly, we study the effect of CG on streaming applications deployed on an Apache Storm cluster running in a private cloud. The storm cluster consists of 8 medium sized machines with 2 virtual CPUs and 4 GB of memory each. Moreover, a Kafka cluster with 8 partitions is used as a data source. We perform experiments to compare CG, PKG, KG, and SG on the TW dataset. The parameters are selected in a way that the number of sources and workers match the number of executors in the Storm cluster. In this experiment, we use a topology configuration with 8 sources and 24 workers. We report overall throughput, end-to-end latency and memory footprint.

In the first experiment, we evaluate the performance of the algorithms in a homogeneous cluster. We emulate different levels of CPU consumption per key, by adding a fixed delay to the processing. We prefer this solution over implementing a specific application to control better the load on the workers. We choose a range that can bring our configuration to a saturation point, although the raw numbers would vary for different setups. Even though real deployments rarely operate at saturation point, CG allows better resource utilization, therefore supporting the same workload on a smaller number of machines, but working on a higher overall load point each. In this case, the minimum delay (0.1ms) corresponds approximately to reading 400kB sequentially from memory, while the

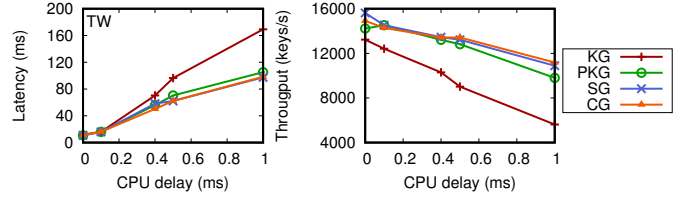


FIG. 14: Throughput and latency for TW dataset on a homogenous Storm cluster.

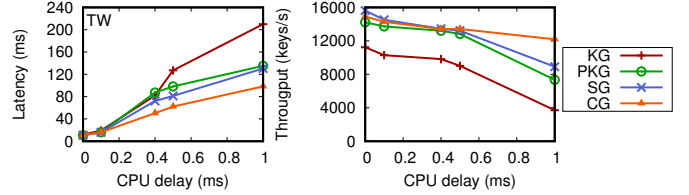


FIG. 15: Throughput and latency for TW dataset on a heterogenous Storm cluster.

maximum delay (1ms) to $\frac{1}{10}$ -th of a disk seek.⁵ Nevertheless, even more expensive tasks exist: parsing a sentence with NLP tools can take up to 500ms.⁶

Figure 14 reports the throughput and end-to-end latency for the TW dataset on the homogenous cluster. Also, during the experiment, KG was consuming 7% of memory in the cluster vs. 8.5% for PKG and CG, and 14% for SG. Results shows that KG provides low memory overhead but coupled with low throughput and high execution latency. Alternatively, PKG, SG and CG provide superior performance in terms of throughput, latency and memory consumption.

Further, we evaluate the performance of CG in the presence of heterogeneity in the cluster. We use the `cpulimit` application to change the resource capacity over time and monitor the behavior of different approaches in terms of throughput and end-to-end latency. In particular, we limit the cpu resources of two of the executors to 30% of the available CPU resources to mimic the heterogeneity in the cluster. During the experiment, we give the system 10 minutes grace period to reach a stable state before collecting the statistics. Figure 15 reports the throughput and the end-to-end latency of the experiment. Results show that CG outperforms the other approaches both in terms of throughput and end-to-end latency. In particular, and compared to KG, it provides up to $2\times$ better end-to-end latency and $3.44\times$ better performance in terms of throughput.

Overall, we observe that CG is a very competitive solution with respect to KG, PKG and SG, performing much better with respect to throughput and end-to-end latency and imposing a small memory footprint, while at the same time tackling the problem of heterogeneity of available resources at the workers in the cluster.

VIII. CONCLUSION

We studied the load balancing problem for streaming engines running in a heterogeneous cluster and processing varying workload. We proposed a novel partitioning strategy called Consistent Grouping. CG leverages two very simple, but extremely powerful ideas: *power of random choices* and *fair virtual worker assignment*. It efficiently achieves fair

⁵http://brenocon.com/dean_perf.html

⁶<http://nlp.stanford.edu/software/parser-faq.shtml#n>

load balancing for streaming applications processing skewed workloads. We provided a theoretical analysis of the proposed algorithm and showed via extensive empirical evaluation that the cg outperforms the state-of-the-art approaches. In particular, CG achieves 3.44x better performance in terms of latency compared to key grouping.

REFERENCES

- [1] F. Ahmad, S. T. Chakradhar, A. Raghunathan, and T. Vijaykumar. Tarazu: optimizing mapreduce on heterogeneous clusters. In *SIGARCH*, volume 40, pages 61–74. ACM, 2012.
- [2] Y. Azar, A. Z. Broder, A. R. Karlin, and E. Upfal. Balanced allocations. *SIAM J. Comput.*, 29(1):180–200, 1999.
- [3] C. Balkesen, N. Tatbul, and M. T. Özsu. Adaptive input admission and management for parallel stream processing. In *DEBS*, pages 15–26. ACM, 2013.
- [4] Y. Ben-Haim and E. Tom-Tov. A Streaming Parallel Decision Tree Algorithm. *JMLR*, 11:849–872, 2010.
- [5] R. Berinde, P. Indyk, G. Cormode, and M. J. Strauss. Space-optimal heavy hitters with strong error bounds. *ACM Trans. Database Syst.*, 35(4):1–28, 2010.
- [6] V. Cardellini, M. Colajanni, and S. Y. Philip. Dynamic load balancing on web-server systems. *IEEE Internet computing*, 3(3):28, 1999.
- [7] R. Castro Fernandez, M. Migliavacca, E. Kalyvianaki, and P. Pietzuch. Integrating scale out and fault tolerance in stream processing using operator state management. In *SIGMOD*, pages 725–736. ACM, 2013.
- [8] L. Chen, O. Villa, S. Krishnamoorthy, and G. R. Gao. Dynamic load balancing on single-and multi-gpu systems. In *IPDPS*, pages 1–12. IEEE, 2010.
- [9] M. Cherniack, H. Balakrishnan, M. Balazinska, D. Carney, U. Cetintemel, Y. Xing, and S. B. Zdonik. Scalable distributed stream processing. In *CIDR*, volume 3, pages 257–268, 2003.
- [10] T. Das, Y. Zhong, I. Stoica, and S. Shenker. Adaptive stream processing using dynamic batch sizing. In *Proceedings of the ACM Symposium on Cloud Computing*, pages 1–13. ACM, 2014.
- [11] B. Gedik. Partitioning functions for stateful data parallelism in stream processing. *VLDB Journal*, 23(4):517–539, 2014.
- [12] B. Godfrey, K. Lakshminarayanan, S. Surana, R. Karp, and I. Stoica. Load balancing in dynamic structured p2p systems. In *INFOCOM*, volume 4, pages 2253–2262. IEEE, 2004.
- [13] P. B. Godfrey and I. Stoica. Heterogeneity and load balance in distributed hash tables. In *INFOCOM*, volume 1, pages 596–606. IEEE, 2005.
- [14] J. E. Gonzalez, Y. Low, H. Gu, D. Bickson, and C. Guestrin. Powergraph: Distributed graph-parallel computation on natural graphs. In *OSDI*, pages 17–30, 2012.
- [15] B. Hindman, A. Konwinski, M. Zaharia, A. Ghodsi, A. D. Joseph, R. Katz, S. Shenker, and I. Stoica. Mesos: A Platform for Fine-grained Resource Sharing in the Data Center. In *NSDI*, Berkeley, CA, USA, 2011.
- [16] J. Hwang and T. Wood. Adaptive performance-aware distributed memory caching. In *ICAC*, volume 13, pages 33–43, 2013.
- [17] E. Kalyvianaki, M. Fiscato, T. Salonidis, and P. Pietzuch. Themis: Fairness in federated stream processing under overload. In *SIGMOD*, pages 541–553. ACM, 2016.
- [18] D. Karger, E. Lehman, T. Leighton, R. Panigrahy, M. Levine, and D. Lewin. Consistent hashing and random trees: Distributed caching protocols for relieving hot spots on the world wide web. In *STOC*, pages 654–663. ACM, 1997.
- [19] N. R. Katsipoulakis, A. Labrinidis, and P. K. Chrysanthis. A holistic view of stream partitioning costs. *VLDB*, 10(11), 2017.
- [20] Z. Khayyat, K. Awara, A. Alonazi, H. Jamjoom, D. Williams, and P. Kalnis. Mizan: a system for dynamic load balancing in large-scale graph processing. In *Proceedings of the 8th ACM European Conference on Computer Systems*, pages 169–182. ACM, 2013.
- [21] A. Kolioussis, M. Weidlich, R. Castro Fernandez, A. L. Wolf, P. Costa, and P. Pietzuch. Saber: Window-based hybrid stream processing for heterogeneous architectures. In *SIGMOD*, pages 555–569. ACM, 2016.
- [22] J. Lin et al. The curse of zipf and limits to parallelization: A look at the stragglers problem in mapreduce. In *7th Workshop on Large-Scale Distributed Systems for Information Retrieval*, volume 1, 2009.
- [23] G. Malewicz, M. H. Austern, A. J. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski. Pregel: a system for large-scale graph processing. In *SIGMOD*, pages 135–146. ACM, 2010.
- [24] V. Mirrokni, M. Thorup, and M. Zadimoghaddam. Consistent hashing with bounded loads. *arXiv preprint arXiv:1608.01350*, 2016.
- [25] M. Mitzenmacher, R. Sitaraman, et al. The power of two random choices: A survey of techniques and results. In *Handbook of Randomized Computing*, pages 255–312, 2001.
- [26] M. A. U. Nasir, G. D. F. Morales, D. García-Soriano, N. Kourtellis, and M. Serafini. Partial key grouping: Load-balanced partitioning of distributed streams. *arXiv preprint arXiv:1510.07623*, 2015.
- [27] M. A. U. Nasir, G. D. F. Morales, D. García-Soriano, N. Kourtellis, and M. Serafini. The power of both choices: Practical load balancing for distributed stream processing engines. In *ICDE*, pages 137–148, April 2015.
- [28] M. A. U. Nasir, G. D. F. Morales, N. Kourtellis, and M. Serafini. When two choices are not enough: Balancing at scale in distributed stream processing. In *ICDE*, pages 589–600, May 2016.
- [29] K. Ousterhout, P. Wendell, M. Zaharia, and I. Stoica. Sparrow: distributed, low latency scheduling. In *SOSP*, pages 69–84, 2013.
- [30] G. Park. A Generalization of Multiple Choice Balls-into-bins. In *PODC*, pages 297–298, 2011.
- [31] E. Rahm and R. Marek. Dynamic multi-resource load balancing in parallel database systems. In *VLDB*, volume 95, pages 11–15. Citeseer, 1995.
- [32] S. Salihoglu and J. Widom. Gps: A graph processing system. In *Proceedings of the 25th International Conference on Scientific and Statistical Database Management*, page 22. ACM, 2013.
- [33] S. Schneider, J. Wolf, K. Hildrum, R. Khandekar, and K.-L. Wu. Dynamic load balancing for ordered data-parallel regions in distributed streaming systems. In *Middleware*, page 21. ACM, 2016.
- [34] M. A. Shah, J. M. Hellerstein, S. Chandrasekaran, and M. J. Franklin. Flux: An adaptive partitioning operator for continuous query systems. In *ICDE*, pages 25–36. IEEE, 2003.
- [35] L. Suresh, M. Canini, S. Schmid, and A. Feldmann. C3: Cutting tail latency in cloud data stores via adaptive replica selection. In *NSDI*, pages 513–527, 2015.
- [36] R. Taft, E. Mansour, M. Serafini, J. Duggan, A. J. Elmore, A. Aboulnaga, A. Pavlo, and M. Stonebraker. E-store: Fine-grained elastic partitioning for distributed transaction processing systems. *VLDB*, 8(3):245–256, Nov. 2014. ISSN 2150-8097. doi: 10.14778/2735508.2735514. URL <http://dx.doi.org/10.14778/2735508.2735514>.
- [37] V. K. Vavilapalli, A. C. Murthy, C. Douglas, S. Agarwal, M. Konar, R. Evans, T. Graves, J. Lowe, H. Shah, S. Seth, et al. Apache hadoop yarn: Yet another resource negotiator. In *SCC*, page 5, 2013.
- [38] Y. Xing, S. Zdonik, and J.-H. Hwang. Dynamic load distribution in the borealis stream processor. In *ICDE*, pages 791–802, 2005.
- [39] D. Yan, J. Cheng, Y. Lu, and W. Ng. Effective techniques for message reduction and load balancing in distributed graph computation. In *WWW*, pages 1307–1317. ACM, 2015.
- [40] M. Zaharia, A. Konwinski, A. D. Joseph, R. H. Katz, and I. Stoica. Improving mapreduce performance in heterogeneous environments. In *OSDI*, volume 8, page 7, 2008.