

Building your Cross-Platform Application with RHEEM

Sanjay Chawla¹ Bertty Contreras-Rojas¹ Zoi Kaoudi¹
 Sebastian Kruse^{2*} Jorge-Arnulfo Quiané-Ruiz¹

¹Qatar Computing Research Institute, Hamad Bin Khalifa University

²Hasso Plattner Institute, University of Potsdam

<http://da.qcri.org/rheem/>

ABSTRACT

Today, organizations typically perform tedious and costly tasks to juggle their code and data across different data processing platforms. Addressing this pain and achieving automatic cross-platform data processing is quite challenging because it requires quite good expertise for all the available data processing platforms. In this report, we present RHEEM, a general-purpose cross-platform data processing system that alleviates users from the pain of finding the most efficient data processing platform for a given task. It also splits a task into subtasks and assigns each subtask to a specific platform to minimize the overall cost (e.g., runtime or monetary cost). To offer cross-platform functionality, it features (i) a robust interface to easily compose data analytic tasks; (ii) a novel cost-based optimizer able to find the most efficient platform in almost all cases; and (iii) an executor to efficiently orchestrate tasks over different platforms. As a result, it allows users to focus on the business logic of their applications rather than on the mechanics of how to compose and execute them. RHEEM is released under an open source license.

1. INTRODUCTION

The pursuit of comprehensive, efficient, and scalable data analytics as well as the *one-size-does-not-fit-all* dictum have given rise to a plethora of data processing platforms (*platforms* for short). These specialized platforms include DBMS, NoSQL, and MapReduce-like platforms. In fact, just under the umbrella of NoSQL, there are reportedly over 200 different platforms¹. Each excels in specific aspects allowing applications to achieve high performance and scalability. For example, while Spark supports `Select` queries, Postgres can execute them much faster by using indices. However, Postgres is not as good as Spark for general purpose batch processing where parallel full scans are the key performance factor. Several studies have shown this kind of performance differences [20, 32, 36, 50, 57].

Moreover, today's data analytics is moving beyond the limits of a single platform. For example: (i) IBM reported that North York hospital needs to process 50 diverse datasets, which run on a dozen different platforms [35];

(ii) Airlines need to analyze large datasets, which are produced by different departments, are of different data formats, and reside on multiple data sources, to produce global reports for decision makers [9]; (iii) Oil & Gas companies need to process large amounts of diverse data spanning various platforms [19, 34]; (iv) Several data warehouse applications require data to be moved from a MapReduce-like system into a DBMS for further analysis [27, 53]; and (v) Using multiple platforms for machine learning improves performance significantly [20, 36].

To cope with these new requirements, developers (or data scientists) have to write ad-hoc programs and scripts to integrate different platforms. This is not only a tedious, time-consuming, and costly task, but it also requires knowledge of the intricacies of the different platforms to achieve high efficiency and scalability. Some systems have appeared with the goal of facilitating platform integration [2, 4, 10, 12]. Nonetheless, they all require a good deal of expertise from developers, who still need to decide which processing platforms to use for each task at hand. Recent research has taken steps towards transparent cross-platform execution [15, 28, 32, 43, 55, 56], but lacks several important aspects. Usually these efforts do not automatically map tasks to platforms. Additionally, they do not consider complex data movement (i.e., with data transformations) among platforms [28, 32]. Finally, most of the research focuses on specific applications [15, 43, 55].

Therefore, there is a clear need for a systematic approach to enable efficient *cross-platform data processing*, i.e., use of multiple data processing platforms. The Holy Grail would be to replicate the success of DBMSs for cross-platform data processing. Users simply send their tasks expressing the logic of their applications, and the cross-platform system decides on which platform(s) to execute each task with the goal of minimizing its cost (e.g., runtime or monetary cost). In other words, users focus on the high level details and the cross-platform system takes care of the low level details.

Building a cross-platform system is challenging on numerous fronts: (i) a cross-platform system not only has to effectively find all the suitable platforms for a given task, but also has to choose the most efficient one; (ii) cross-platform settings are characterized by high uncertainty as different platforms are autonomous and thus one has little control over them; (iii) the performance gains of using multiple platforms should compensate the added cost of moving data across platforms; (iv) it is crucial to achieve inter-platform parallelism to prevent slow platforms from dominating execution time; and (v) the system should be extensible to new

*Work partially done while interning at QCRI.

¹<http://db-engines.com>

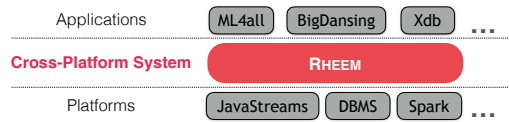


Figure 1: Rheem in the data analytics stack.

platforms and application requirements.

In this report, we present RHEEM², the first general-purpose cross-platform system to tackle all of the above challenges. The goal of RHEEM is to *enable* applications and users to run data analytic tasks *efficiently* on one or more data processing platforms. To do so, it decouples applications from platforms as shown in Figure 1. Applications issue their tasks to RHEEM, which in turn decides where to execute them. As of today, RHEEM supports a variety of platforms: Spark, Flink, JavaStreams, Postgres, GraphX, GraphChi, and Giraph. We are currently testing RHEEM in a large international airline company and in a biomedical research institute. In the former case, we aim at seamlessly integrating all data analytic activity governing an aircraft; In the latter case, we aim at reducing the effort scientists need for building data analytic pipelines while at the same time speeding up the running time. Several papers show different aspects of RHEEM: the vision behind it [17]; its optimizer [39]; its inequality join algorithm [38]; and a couple of its applications [36, 37]. A couple of demo papers showcase the benefits of RHEEM [16] and its interface [44]. This report aims at presenting the complete design of RHEEM and how all its pieces work together.

In summary, we identify four situations in which applications require support for cross-platform data processing in Section 2. For each case, we use a real application to show experimentally the benefits of cross-platform data processing using RHEEM. In Section 3, we present the data and processing model of RHEEM and show how it shields users from the intricacies of the underlying platforms. RHEEM provides flexible operator mappings that allow for better exploiting the underlying platforms. Also, its extensible design allows users to add new platforms and operators with very little effort. Then, in Section 4, we discuss the key components of RHEEM that make it novel: among them a cost-based cross-platform optimizer that considers data movement costs; a progressive optimization mechanism to deal with inconsistent cardinality estimates; and a learning tool that alleviates users from the burden of tuning the cost model. We present the RHEEM interfaces whereby users can easily code and run a data analytic task in Section 5. In particular, we present a data-flow language (RheemLatin) and a visual integrated development environment (RHEEM Studio). In Section 6, we show in detail three examples of real RHEEM plans to better illustrate how developers can build their applications using these interfaces. Section 8 outlines the limitations of RHEEM. Finally, we discuss related work in Section 9 and conclude with some open problems in Section 10.

2. CROSS-PLATFORM PROCESSING

We identified four situations in which an application requires support for cross-platform data processing [51]. Figure 2 illustrates these four cases.

²RHEEM is open source under the Apache Software License 2.0 and can be found at <https://github.com/rheem-ecosystem/rheem>.

(1) *Platform-independence*. Applications run an entire task on a single platform but may require switching platforms for different input datasets or tasks usually with the goal of achieving better performance (Figure 2(a)).

(2) *Opportunistic cross-platform*. Applications might also benefit performance-wise from using multiple platforms to run one single task (Figure 2(b)).

(3) *Mandatory cross-platform*. Applications may require multiple platforms because the platform where the input data resides, e.g., PostgreSQL, cannot perform the incoming task, e.g., a machine learning task. Thus, data should be moved from the platform it resides to another platform (Figure 2(c)).

(4) *Polystore*. Applications may require multiple platforms because the input data is stored on multiple data stores (Figure 2(d)).

In contrast to existing systems [28, 29, 32, 55, 58], RHEEM helps users in *all* above cases. The design of our system has been mainly driven by four applications: a data cleaning application, *BigDancing* [37]; a machine learning application, *ML4all* [36]; a database application, *xDB*; and an end-to-end data discovery and preparation application, *Data Civilizer* [31]. We use these applications to showcase the benefits of performing cross-platform data processing, instead of single-platform data processing, in terms of both performance and ease of use.

2.1 Platform Independence

Applications are usually tied to a specific platform. This may not constitute the ideal case for two reasons. First, as more efficient platforms become available, developers need to re-implement existing applications on top of these new platforms. For example, Spark SQL [14] and MLlib [13] are the Spark counterparts of Hive [6] and Mahout [7]. Migrating an application from one platform to another is a time-consuming and costly task and hence it is not always a viable choice. Second, for different inputs of a specific task, a different platform may be the most efficient one, so the best platform cannot be determined statically. For instance, running a specific task on a big data platform for very large datasets is often a good choice, while single-node platforms with only little overhead costs are often a better choice for small datasets [20]. Thus, enabling applications to seamlessly switch from one platform to another according to the input dataset and task is important. RHEEM *dynamically determines the best platform to run an incoming task*.

Benefits. We use BigDancing [37] to demonstrate the benefits of providing platform independence. Users specify a data cleaning task with five logical operators: **Scope** (identifies relevant data), **Block** (defines the group of data among which an error may occur), **Iterate** (enumerates candidate errors), **Detect** (determines whether a candidate error is indeed an error), and **GenFix** (generates a set of possible repairs). RHEEM maps these operators to RHEEM operators to decide the best underlying platform. We show the power of supporting cross-platform data processing by running an error detection task on a widely used Tax dataset [30]. The task is based on the denial constraint $\forall t_1, t_2, \neg(t_1.\text{Salary} > t_2.\text{Salary} \wedge t_1.\text{Tax} < t_2.\text{Tax})$, which states that there is an inconsistency between two tuples representing two different persons if one earns a higher salary but pays a lower tax. We considered NADEEF [24], a data cleaning tool, and SparkSQL, a

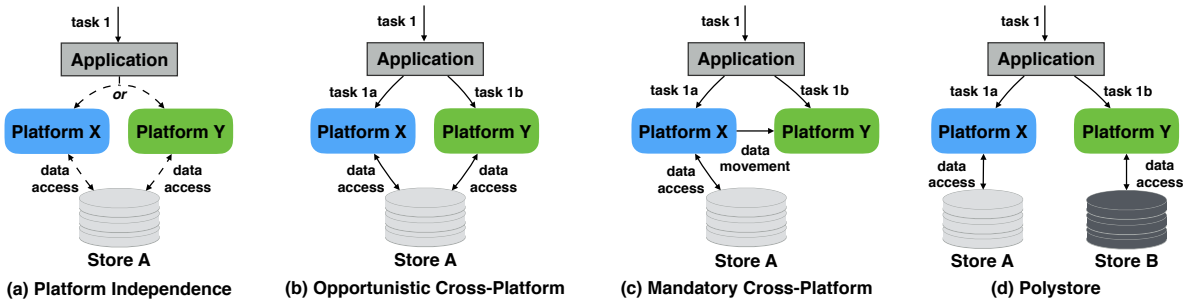


Figure 2: Cross-platform cases.

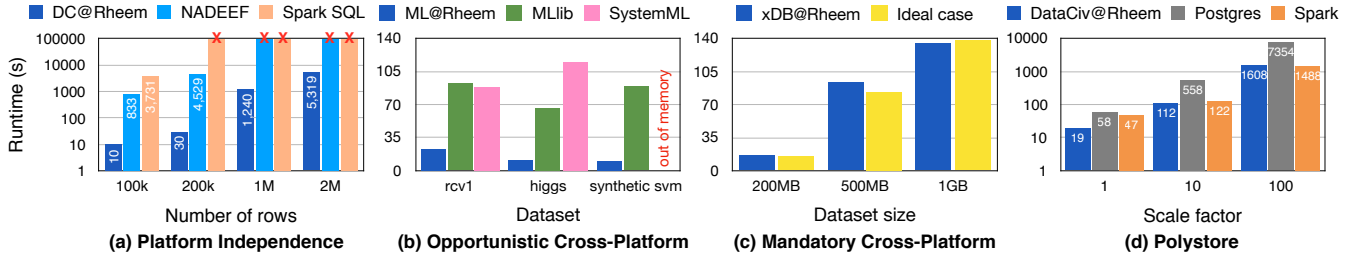


Figure 3: Benefits of the cross-platform data processing approach (using Rheem).

general-purpose framework, as baselines and forced RHEEM to use either Spark or JavaStreams per run.

Figure 3(a) shows the results³. Overall, we observe that RHEEM (DC@Rheem) allows data cleaning tasks to scale up to large datasets and be at least three orders of magnitude faster than baselines. One order of magnitude gain comes from the ability of RHEEM to automatically switch platforms. RHEEM used JavaStreams for small datasets speeding up the data cleaning task by avoiding Spark’s overhead, while it used Spark for the largest datasets. Furthermore, in contrast to SparkSQL that cannot process inequality joins efficiently, RHEEM’s extensibility allowed us to plug in a more efficient inequality-join algorithm [38], thereby further improving over these baselines. In a nutshell, BigDancing benefited from RHEEM because of its ability to effectively switch platforms and because of its extensibility to easily plug optimized algorithms. We demonstrated how BigDancing benefits from RHEEM in [16].

2.2 Opportunistic Cross-Platform

While some applications can be executed on a single platform, there are cases where their performance would be sped up by using multiple platforms. For instance, users can run a gradient descent algorithm, such as SGD, on top of Spark relatively fast. Still, we recently showed that mixing it with JavaStreams significantly improves performance [36]. In fact, opportunistic cross-platform processing can be seen as the execution counter-part of *polyglot persistence* [52], where different types of databases are combined to leverage their individual strengths. However, developing such cross-platform applications is difficult: developers must know all the cases where it is beneficial to use multiple platforms and how exactly to use them. These opportunities are often very hard (if not impossible) to spot. Even worse, like in the platform independence case, they usually cannot be determined

a priori. RHEEM finds and exploits opportunities of using multiple processing platforms.

Benefits. Let us now take our machine learning application, *ML4all* [36], to showcase the benefits of using multiple platforms to perform one single task. ML4all abstracts three fundamental phases (namely preparation, processing, and convergence) found in most machine learning tasks via seven logical operators which are mapped to RHEEM operators. In the preparation phase, the dataset is prepared appropriately along with the necessary initialization of the algorithm (**Transform** and **Stage** operators). The processing phase computes the gradient and updates the current estimate of the solution (**Sample**, **Compute**, and **Update** operators) while the convergence phase repeats the processing phase based on the number of iterations or other criteria (**Loop** and **Converge** operators). We demonstrate the benefits of using RHEEM with a classification task over three benchmark datasets, using Stochastic Gradient Descent (SGD).

Figure 3(b) shows the results. We observe that, even though all systems use the same SGD algorithm, RHEEM allows this algorithm to run significantly faster than competing Spark-based systems. This is because of two main reasons. First, this comes from opportunistically running the **Compute**, **Update**, **Converge**, and **Loop** operators on JavaStreams, thereby avoiding some of the Spark’s overhead. RHEEM runs the rest of the operators on Spark. MLlib and SystemML do not avoid such overhead by purely using Spark for the entire algorithm. Second, ML4all leverages RHEEM’s extensibility to plug an efficient sampling operator, resulting in significant speedups. We demonstrated how ML4all further benefits from RHEEM in [16].

2.3 Mandatory Cross-Platform

There are cases where an application needs to go beyond the functionalities offered by the platform on which the data is stored. For instance, a dataset is stored on a relational database and a user needs to perform a clustering task on particular attributes. Doing so inside the relational database

³The red cross means we stopped the execution after 40 hrs.

might simply be disastrous in terms of performance. Thus, the user needs to move the projected data out of the relational database and, for example, put it on HDFS in order to use Apache Flink [3], which is known to be efficient for iterative tasks. A similar situation occurs in complex data analytics applications with disparate subtasks. As an example, an application that extracts a graph from a text corpus to perform subsequent graph analytics may require using both a text and a graph analytics system. The required integration of platforms is tedious, repetitive, and particularly error-prone. Nowadays, developers write ad-hoc programs to move the data around and integrate different platforms. RHEEM *not only selects the right platforms for each task but also moves the data if necessary at execution time*.

Benefits. We use xDB⁴, a system on top of RHEEM with database functionalities, to demonstrate the benefits of performing cross-platform data processing for the above situation. It provides a declarative language to compose data analytic tasks, while its optimizer produces a plan to be executed in RHEEM. We evaluate the benefits of RHEEM with the cross-community pagerank⁵ task, which is not only hard to express in SQL but also inefficient to run on a DBMS. Thus, it is important to move the computation to another platform. In this experiment, the input datasets are on Postgres and RHEEM moves the data into Spark.

Figure 3(c) shows the results. As a baseline, we consider the ideal case where the data is on HDFS and RHEEM simply uses either JavaStreams or Spark to run the tasks. We observe that RHEEM allows xDB (xDB@Rheem) to achieve similar performance with the ideal case in all the situations, while fully automating the process. This is a remarkable result as RHEEM needs to move data out of Postgres to perform the tasks, in contrast to the ideal case.

2.4 Polystore

In many organizations, data is collected in different formats and on heterogeneous storage platforms (*data lakes*). Typically, a data lake comprises various DBMSs, document stores, key-value stores, graph databases, and pure file systems. As most of these stores are tightly coupled with an execution engine, e.g., a DBMS, it is crucial to be able to run analytics over multiple platforms. For this, users perform not only tedious, time-intensive, and costly data migration, but also complex integration tasks for analyzing the data. RHEEM *shields the users from all these tedious tasks and allows them to instead focus on the logic of their applications*.

Benefits. A clear example that shows the benefits of cross-platform data processing in a polystore case is the Data Civilizer system [31]. Data Civilizer is a big data management system for data discovery, extraction, and cleaning from data lakes in large enterprises [26]. It constructs a graph that expresses relationships among data existing in heterogeneous data sources. Data Civilizer uses RHEEM to perform complex tasks over information that spans multiple data storages. We measure the efficiency of RHEEM for these polystore tasks with TPC-H query 5. In this experiment, we assume that the data is stored in HDFS (LINEITEM and ORDERS), Postgres (CUSTOMER, REGION, and SUPPLIER), and a local file system (NATION). Thus, this task

⁴<https://github.com/rheem-ecosystem/xdb>

⁵This task basically intersects two community-DBpedia datasets and runs pagerank on the resulting dataset.

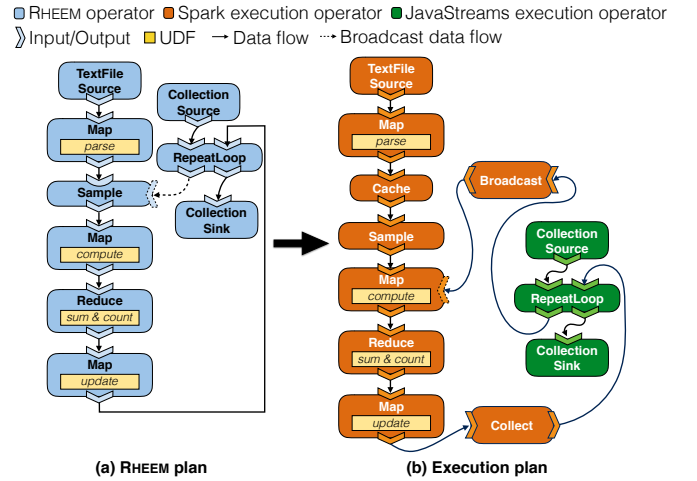


Figure 4: SGD example.

performs join, groupby, and orderby operations across three different platforms. In this scenario, the common practice is to move the data into the database to enact the queries inside the database [27, 53] or move the data entirely to HDFS and use Spark. We consider these two practices as the baseline. For a fairer IO comparison, we also set the “parallel query” and “effective IO concurrency” features of Postgres to 4.

Figure 3(d) shows the results. RHEEM (DataCiv@Rheem) is significantly faster, namely up to 5 \times , than the current practice. We observed that loading data into Postgres is already approximately 3 \times slower than it takes RHEEM to complete the entire task. Even when discarding data migration times, RHEEM can still perform quite similarly to the parallel version of Postgres. The pure execution time in Postgres for scale factor 100 amounts to 1,541 sec compared to 1,608 sec for RHEEM, which exploits Spark for data parallelism. We also observe that RHEEM has negligible overhead over the case where the developer writes ad-hoc scripts to move the data to HDFS for running the task on Spark. In particular, RHEEM is twice faster than Spark for scale factor 1 because it moves less data from Postgres to Spark.

3. RHEEM MODEL

First of all, let us emphasize that RHEEM is *not* yet another data processing platform. On the contrary, it is designed to work between applications and platforms (as shown in Figure 1), helping applications to choose the right platform(s) for a given task. RHEEM is the first general-purpose cross-platform system that shields users from the intricacies of the underlying platforms and let them focus only on the logic of their applications. We define the RHEEM data and processing models in the following.

Data Quanta. The RHEEM data model relies on *data quanta*, the smallest processing units from the input datasets. A data quantum can express a large spectrum of data formats, such as database tuples, edges in a graph, or the full content of a document. This flexibility allows applications and users to define a data quantum at any granularity level, e.g., at the attribute level rather than at the tuple level for a relational database. This fine-grained data model allows RHEEM to work in a highly parallel fashion, if necessary, to achieve better scalability and performance.

Rheem Plan. RHEEM accepts as input a RHEEM *plan*:

a directed data flow graph whose vertices are RHEEM operators and whose edges represent data flows among the operators. A RHEEM operator is a platform agnostic data transformation over its input data quanta, e.g., a **Map** operator transforms an individual data quantum while a **Reduce** operator aggregates input data quanta into a single output data quantum. Only **Loop** operators accept feedback edges, which allows iterative data flows to be expressed. Users or applications can refine the behavior of operators with a UDF. Optionally, applications can also attach the selectivities of the operators through a UDF. RHEEM comes with default selectivity values in case they are not provided. A RHEEM plan must have at least one source operator, i.e., an operator reading or producing input data quanta, and one sink operator per branch, i.e., an operator retrieving or storing the result. Intuitively, data quanta are flowing from source to sink operators, thereby being manipulated by all inner operators. As our processing model is based on primitive operators, RHEEM plans are highly expressive. This is in contrast to other systems that accept either declarative queries [32, 58] or coarse-granular operators [28].

EXAMPLE 1. Figure 4(a) shows a RHEEM plan for the stochastic gradient descent algorithm (SGD). Initially, the dataset containing the data points is read via a **TextFileSource** operator and parsed using a **Map** operator while the initial weights are read via a **Collection** source operator. After the **RepeatLoop** operator, the weights are fed to the **Sample** operator, where a set of input data points is sampled. Next, **Map(compute)** computes the gradient for each sampled data point. Note that as **Map(compute)** requires all weights to compute the gradient, the weights are broadcasted at each iteration to the **Sample** operator (denoted by the dotted line). Then, the **Reduce** operator computes the sum and count of all gradients. The next **Map** operator uses these sum and count values to update the weights. This process is repeated until the loop condition is satisfied. The resulting weights are output in a collection sink.

Execution Plan. Given a RHEEM plan as input, RHEEM uses a cost-based optimization approach to produce an *execution plan* by selecting one or more platforms to efficiently execute the input plan. The cost can be any user-specified cost, e.g., runtime or monetary cost. The resulting execution plan is again a data flow graph, where the vertices are now *execution operators*. An execution operator implements one or more RHEEM operators with platform-specific code. For instance, the **Cache** Spark execution operator in RHEEM implements the **Cache** RHEEM operator by calling the `RDD.cache()` operation of Spark. An execution plan may also comprise additional execution operators for data movement (e.g., data broadcasting) or data reuse (e.g., data caching). Additionally, each execution operator has attached a UDF where its cost is specified. RHEEM learns such costs from execution logs using machine learning. We discuss more details in Section 4.5.

EXAMPLE 2. Figure 4(b) shows the SGD execution plan produced by RHEEM when Spark and *JavaStreams* are the only available platforms. This execution plan exploits high parallelism for the large dataset of input data points and avoids the extra overhead incurred by big data processing platforms for the smaller collection of weights. Note that the execution plan also contains three execution operators for transferring (**Broadcast**, **Collect**) and making data quanta reusable across the platforms (**Cache**).

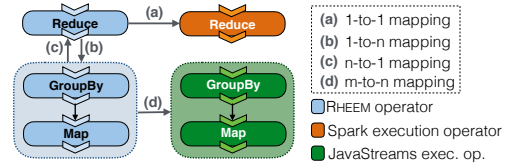


Figure 5: Operator mappings.

Operator Mappings. To produce an execution plan, RHEEM relies on flexible *m-to-n mappings* to map RHEEM operators to execution operators. Supporting *m-to-n mappings* is particularly useful as it allows to map whole subplans of RHEEM operators to subplans of execution operators. Additionally, a subplan of RHEEM (or execution) operators can map to another subplan of RHEEM (respectively execution) operators. As a result, we can handle different abstraction levels among platforms, e.g., to emulate RHEEM operators that are not natively supported by a specific platform. This is not possible in other systems, such as [28].

EXAMPLE 3. Figure 5 illustrates the mapping for the **Reduce** RHEEM operator. This operator directly maps to the **Reduce** Spark execution operator via a 1-to-1 mapping (mapping (a)). However, it does not have a direct mapping to a *JavaStreams* execution operator. Instead, it maps to a set of RHEEM operators (**GroupBy** and **Map**) via a 1-to-n mapping (mapping (b)) and vice-versa (n-to-1 mapping (c)). In turn, this set of RHEEM operators maps to a set of *JavaStreams* execution operators (**GroupBy** and **Map**) via an m-to-n mapping (mapping (d)).

Data movement. Data flows among operators via *communication channels* (or simply *channels*). A channel can be any internal data structure within a data processing platform (e.g., **RDD** for Spark or **Collection** for *JavaStreams*), or simply a file. In the case of two execution operators of different platforms connected within a plan, it is necessary to convert the output channel of one to the input channel of the other (e.g., from **RDD** to **Collection**). These conversions are handled by conversion operators, which in fact are regular execution operators. For example, we can convert a Spark **RDD** channel to a *JavaStreams* **Collection** channel using the **SparkCollect** operator (see Figure 4(b)). We represent the space of data movement paths across all platforms as a *channel conversion graph*, where the channels form its vertices and the *conversion operators* form its directed edges connecting one source channel to a target channel. Unlike other approaches [28, 32], developers do not need to provide conversion operators for all combinations of source and target channels. It is thus much easier for developers to add new platforms to RHEEM.

Extensibility. We designed RHEEM to address extensibility as a first-class citizen rather than as “nice-to-have” feature. Users add new RHEEM and execution operators by merely extending or implementing few abstract classes/interfaces. RHEEM provides template classes to facilitate the development for different operator types. Users also add operator mappings by simply implementing an interface and specifying a graph pattern that matches the RHEEM operator. As a result, users can plug a new platform by providing: (i) its execution operators and their mappings; and (ii) the communication channels that are specific to the new platform (e.g., **RDDChannel** for Spark). Users neither have to modify the RHEEM code nor integrate the newly added platform with all the already supported platforms.

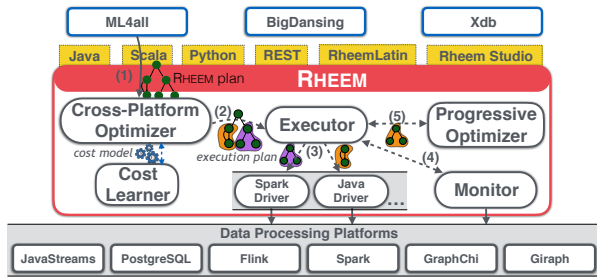


Figure 6: Rheem’s ecosystem and architecture.

4. RHEEM INTERNALS

In this section, we give the details of the RHEEM internals. Figure 6 depicts the RHEEM ecosystem, i.e., the RHEEM core architecture together with three main applications built on top of it. Users provide a RHEEM plan to the system (Step (1) in Figure 6), using Java, Scala, Python, REST, RheemLatin, or RHEEM Studio API (yellow boxes in Figure 6). The *cross-platform optimizer* compiles the RHEEM plan into an execution plan (Step (2)), which specifies the processing platforms to use; the *executor* schedules the resulting execution plan on the selected platforms (Step (3)); the *monitor* collects statistics and checks the health of the execution (Step (4)); the *progressive optimizer* re-optimizes the plan if the cardinality estimates turn out to be inaccurate (Step (5)); and the *cost learner* helps users in building the cost model offline. In the following, we explain each of these components using the pseudocode in Algorithm 1, which shows the entire data processing pipeline.

4.1 Optimizer

The cross-platform optimizer (Line 1 in Algorithm 1) is responsible for selecting the most efficient platform for executing each single operator in a RHEEM plan. One might think of a rule-based optimizer for selecting the right platforms to perform a given RHEEM plan. However, while a rule-based optimizer could determine how to split and execute a plan, e.g., based on its processing patterns [32, 58], it is neither practical nor effective. First, by setting rules, one may make only very simplistic decisions based on the different cardinality and complexity of each operator. Second, the cost of a task on any given platform depends on many input parameters, which hampers a rule-based optimizer’s effectiveness as it oversimplifies the problem. Third, as new platforms and applications emerge, maintaining a rule-based optimizer becomes cumbersome.

We thus pursue a more flexible cost-based approach: *we split a given RHEEM plan into subplans and determine the best platform for each subplan so that the total plan cost is minimized.* Figure 4(b) shows how the RHEEM plan of Figure 4(a) was split into two subplans to be executed in JavaStreams and Spark. Below, we give the four main phases of the optimizer, namely *plan inflation*, *cardinality and cost estimation*, *data movement planning*, and *plan enumeration*. Technical details about these can be found in [39].

At first, the optimizer passes the RHEEM plan through an inflation phase. That is, it applies a set of operator mappings as described in Section 3. The optimizer then annotates the inflated plan with the cost of each execution operator. RHEEM represents cost estimates as intervals with a confidence value, which allows it to perform on-the-fly re-optimization as we will see in Section 4.4. The

Algorithm 1: Cross-platform data processing

Input: RHEEM plan *rheemPlan*

```

1 exPlan ← Optimize(rheemPlan)
2 monitor ← StartMonitor(exPlan)
3 finished ← ExecuteUntilCheckpoint(exPlan, monitor)
4 while ¬finished do
5   updated ← UpdateEstimates(exPlan, monitor)
6   if updated then exPlan ← ReOptimize(exPlan)
7   finished ← ResumeExecution(exPlan, monitor)

```

cost (e.g., wallclock time or monetary cost) of an execution operator depends on (i) its resource usage (CPU, memory, disk, and network) and (ii) the unit costs of each resource (e.g., how much one CPU cycle costs). While the unit costs depend on hardware characteristics, the resource usage of each execution operator depends on its input cardinality. Next, the optimizer looks for the best way to move data quanta among execution operators of different platforms. As noted earlier, we model the problem of finding the most efficient communication path among execution operators as a graph problem, which we proved to be NP-hard. Our solution to this problem relies kernelization and can discover all ways to connect execution operators of different platforms via a sequence of communication channels. After the best data movement strategy is found, the optimizer attaches the data movement cost to the inflated plan. At last, it determines the optimal way of executing a RHEEM plan based on the cost estimates of its inflated plan. For this, it must consider the previously computed data movement costs as well as the start-up costs of data processing platforms. Thus, instead of taking a simple greedy approach that neglects data movement and platform start-up costs, we follow a principled approach: we use an enumeration algebra together with a lossless pruning technique. Our pruning technique is guaranteed to not prune a subplan that is part of the optimal execution plan. As a result, the optimizer can output the optimal execution plan without an exhaustive enumeration.

4.2 Executor

The executor receives an execution plan from the optimizer to run it on the selected data processing platforms (Lines 3 and 7 in Algorithm 1). For example, the optimizer selected the Spark and JavaStreams platforms for our SGD example in Figure 4(a). Overall, the executor follows well-known approaches to parallelize a task over multiple compute nodes, with only few differences in the way it divides an execution plan. In particular, it divides an execution plan into *stages*. A stage is a subplan where (i) all its execution operators are from the same platform; (ii) at the end of its execution, the platforms need to give back the execution control to the executor; and (iii) its terminal operators materialize their output data quanta in a data structure, instead of being pipelined into the next operator.

In our SGD example of Figure 4(b), the executor divides the execution plan into six stages as illustrated in Figure 7. Note that Stage3 contains only the RepeatLoop operator as the executor must have the execution control to evaluate the loop condition. This is why the executor also separates Stage1 from Stage5. Then, it dispatches the stages to the relevant platform drivers, which in turn submit the stages as a job to the underlying platforms. Stages are connected by data flow dependencies so that stages with no dependencies

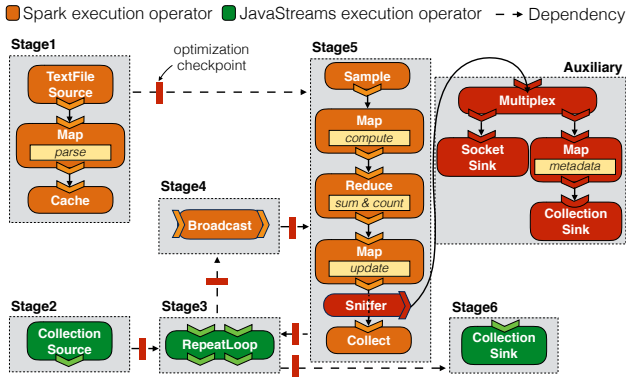


Figure 7: Stage dependencies for SGD.

(e.g., Stage1 and Stage2) are dispatched first in parallel and any other stage is dispatched once its input dependencies are satisfied (e.g., Stage3 after Stage2).

Data Exploration. As data exploration is a key piece in the field of data science, the executor *optionally* allows applications to run in an *exploratory mode* where they can pause and resume the execution of a task at any point. Achieving this in a cross-platform setting is very challenging, because most platforms, such as Spark, Flink, Giraph, Postgres, and Hadoop, do not support pausing task computations at all – let alone resuming a task from an intermediate state. Thus, the challenge resides in enabling the underlying platforms to support data exploration efficiently. RHEEM achieves this by injecting *sniffers* into execution plans and attaching *auxiliary* execution plans. A sniffer is an execution operator that duplicates intermediate results and sends them to an auxiliary execution plan. For example, the user would like to keep track of the weights at each iteration of SGD and thus a sniffer is necessary right after updating the weights (Stage5 in Figure 7). The sniffer sends the weights to an auxiliary plan that is responsible for reporting them back to the user (the socket sink operator in Figure 7). This auxiliary plan is also responsible for computing and storing additional metadata for efficient task resumption (the map and collection sink operators of the auxiliary plan in Figure 7). When resuming a task, the executor performs the task by re-using as much as possible from the previously computed metadata. For instance, if the user pauses the SGD task at iteration i and resumes it later on, the executor fetches the previously computed weights of iteration i and resumes the task.

4.3 Monitor

Recall that the cross-platform optimizer operates in a setting that is characterized by high uncertainty. For instance, the semantics of UDFs and data distributions are usually unknown because of the little control over the underlying platforms. This uncertainty can cause poor cardinality and cost estimates and hence can negatively impact the effectiveness of the optimizer [42]. To compensate this uncertainty, RHEEM registers the execution of a plan with the monitor (Line 2 in Algorithm 1). The monitor collects light-weight execution statistics for the given plan, such as data cardinalities and operator execution times. It is also aware of lazy execution strategies used by the underlying platforms and assigns measured execution time correctly to operators. RHEEM uses these statistics to improve its cost model and re-optimize ongoing execution plans in case of poor cardi-

nality estimates. Additionally, the monitor is responsible for checking the health of the execution. For instance, if it finds a large mismatch between the real output cardinalities and the estimated ones, it pauses the execution plan and sends it to the progressive optimizer.

4.4 Progressive Optimizer

To mitigate the effects of bad cardinality estimates, RHEEM employs a *progressive query optimization* approach. The key principle is to re-optimize the plan whenever the cardinalities observed by the monitor greatly mismatch the estimated ones [45]. Applying progressive query optimization in our setting comes with two main challenges. First, we have only limited control over the underlying platforms, which makes plan instrumentation and halting executions difficult. Second, re-optimizing an ongoing execution plan must efficiently consider the results already produced.

We tackle these challenges by using *optimization checkpoints*. An optimization checkpoint tells the executor to pause the plan execution in order to consider a re-optimization of the plan beyond the checkpoint. The progressive optimizer inserts optimization checkpoints into execution plans wherever (i) cardinality estimates are uncertain (having a wide interval or low confidence) or (ii) the data is at rest (e.g., a Java collection or a file). For instance, the optimizer inserts an optimization checkpoint right after Stage1 as the data is at rest because of the **Cache** operator (see Figure 7). When the executor cannot dispatch a new stage anymore without crossing an optimization checkpoint, it pauses the execution and gives the control to the progressive optimizer. The latter gets the actual cardinalities observed so far by the monitor and re-computes all cardinalities from the current optimization checkpoint (Line 5 in Algorithm 1). In case of a mismatch, it re-optimizes the remaining of the plan (from the current optimization checkpoint) using the new cardinalities (Line 6). It then gives the new execution plan to the executor, which resumes the execution from the current optimization checkpoint (Line 7). RHEEM can switch between execution and progressive optimization any number of times at a negligible cost.

4.5 Cost Model Learner

Profiling operators in isolation might be unrealistic whenever platforms optimize execution across multiple operators, e.g., by pipelining. Indeed, we found cost functions derived from isolated benchmarking to be insufficiently accurate. We thus take a different approach.

Learning the Cost Model. Recall that each execution operator o is associated with a number of resource usage functions (r_o^m , where m is CPU, memory, disk, or network). For instance, the cost function to estimate the CPU cycles required by the **JavaFilter** operator is $r_{JavaFilter}^{CPU} := c_{in} \times (\alpha + \beta) + \delta$, where parameters α and β denote the number of required CPU cycles for each input data quantum in the operator itself and in its UDF, and parameter δ describes some fixed overhead for the operator start-up and scheduling. We then multiply each of these resource usage functions r_o^m with the time required per unit (e.g., msec/CPU cycle) to get the time estimate t_o^m . The total cost estimate for operator o is defined as: $f_o = t_o^{CPU} + t_o^{mem} + t_o^{disk} + t_o^{net}$. However, obtaining the parameters for each resource, such as the α, β, δ values for CPU, is not trivial. We, thus, use

execution logs to *learn* these parameters in an offline fashion and model the cost of individual execution operators as a *regression problem*. Note that the execution logs contain the runtimes of execution stages (i.e., pipelines of operators as defined in Section 4.2) and not of individual operators. Let $((o_1, C_1), (o_2, C_2), \dots, (o_n, C_n)), t$ be an execution stage, with o_i , $0 < i \leq n$, where o_i are execution operators, C_i are input and output true cardinalities, and t is the measured execution time for the entire stage. Furthermore, let $f_i(\mathbf{x}, C_i)$ be the total cost function for execution operator o_i with \mathbf{x} being a vector with the parameters of all resource usage functions (e.g., CPU cycles, disk I/O per data quantum). We are interested in finding $\mathbf{x}_{\min} = \arg \min_{\mathbf{x}} \text{loss}(t, \sum_{i=1}^n f_i(\mathbf{x}, C_i))$. Specifically, we use a *relative* loss function defined as $\text{loss}(t, t') = \left(\frac{|t-t'|+s}{t+s} \right)^2$, where t' is the geometric mean of the lower and upper bounds of the cost estimate produced by $\sum f_i(\mathbf{x}, C_i)$ and s is a regularizer inspired by additive smoothing that tempers the loss for small t . Note that we can easily generalize this optimization problem to multiple execution stages: we minimize the weighted arithmetic mean of the losses of multiple execution stages. In particular, we use as stage weights the sum of the relative frequencies of the stages' operators among all stages, so as to deal with skewed workloads that contain certain operators more often than others. Finally, we apply a genetic algorithm [47] to find \mathbf{x}_{\min} . In contrast to other optimization algorithms, genetic algorithms impose only few restrictions on the loss function to be minimized. Hence, our cost learner can deal with arbitrary cost functions. Applying this technique allows us to calibrate the cost functions with only little additional effort.

Logs Generation. Clearly, the more execution logs are available, the better RHEEM can tune the cost model. Thus, RHEEM comes with a log generator. It first creates a set of RHEEM plans by composing all possible combinations of RHEEM operators forming a particular topology. We found that most data analytic tasks in practice follow three different topologies: *pipeline* (e.g., batch tasks), *iterative* (e.g., ML tasks), and *merge* (e.g., SPJA tasks). It then generates all possible executions plans for the previously created set of RHEEM plans. Next, it creates different configurations for each execution plan, i.e., it varies the UDF complexity, output cardinalities, input dataset sizes, and data types. Once it has generated all possible plans with different configurations, it executes them and logs their runtime.

5. RHEEM INTERFACES

RHEEM provides a set of native APIs for developers to build their applications. These include Java, Scala, Python, and REST. Examples of using these APIs can be found in the RHEEM repository⁶. The code developers have to write is fully agnostic of the underlying platforms. Still, in case the user wants to force RHEEM to execute a given operator on a specific platform, she can invoke the `withTargetPlatform` method. Similarly, she can force the system to use a specific execution operator via the `customOperator` method, which further enables users to employ custom operators without having to extend the API.

Although the native APIs are quite popular among developers, many users are not proficient using these APIs.

⁶<https://github.com/rheem-ecosystem/rheem-benchmark>

Thus, RHEEM also provides two APIs that target non-expert users: a data-flow language (*RheemLatin*) and a visual IDE (*RheemStudio*). We explain these interfaces using our SGD example from Figure 4. However, for the sake of explanation, before going into the details of these two interfaces, we first show how one can implement SGD on RHEEM using one of its native APIs. The salient feature of all these APIs is that they are all platform-agnostic. It is RHEEM that figures out on which platform to execute each of the operators.

5.1 Platform-Agnostic Native API

Let us explain how users can code their applications using one of the native APIs of RHEEM. We use the Scala API and our SGD running example (see Listing 1)⁷.

```

1 val context = new RheemContext(new Configuration)
2   .withPlugin(Spark.basicPlugin)
3   .withPlugin(JavaStreams.basicPlugin)
4 val plan = new PlanBuilder(context)
5 val points = plan.readTextFile("hdfs://myData.csv")
6   .map(parsePoints)
7 val finalWeights = plan.loadCollection(createRandomWeights())
8   .repeat(50, { weights =>
9     points.sample(sampleSize).withBroadcast(weights)
10    .map(computeGradient())
11    .reduce(_ + _)
12    .map(updateWeights())
13  }).collect()

```

Listing 1: SGD task using the Scala API.

First, a user creates the RHEEM context, where she specifies the available platforms (Lines 1-3): Spark and JavaStreams in this example. She then initializes her RHEEM plan with this context (Line 4). Eventually, she creates the graph of RHEEM operators that defines the SGD task (Lines 5-13). Note that RHEEM plans must have at least one source operator (Line 5), i.e., an operator reading or producing input data quanta, and one sink operator per branch (Line 13), i.e., an operator retrieving or storing the result. Recall that a RHEEM plan must have at least one source operator (Line 5) and one sink operator per branch (Line 13). Also, observe that this code is fully agnostic of the underlying platforms. Still, in case the user wants to force RHEEM to execute a given operator on a specific platform, she can invoke the `withTargetPlatform` method. Similarly, she can force the system to use a specific execution operator via the `customOperator` method, which further enables users to employ custom operators without having to extend the API. For clarity reasons, we did not include the UDF implementations in Listing 1.

5.2 RheemLatin

RHEEM provides a data-flow language (*RheemLatin*) for users to specify their tasks [44]. Our goal is to provide ease-of-use to users without compromising expressiveness. RheemLatin follows a procedural programming style to naturally fit the pipeline paradigm of RHEEM. This is similar to the R language, which is quite popular among data scientists. It draws its inspiration from PigLatin [48] and hence it has PigLatin's grammar and supports most PigLatin's keywords. In fact, one could see it as an extension of PigLatin for cross-platform settings. For example, users can specify the platform for any part of their queries. More importantly,

⁷The complete source code of this task is available online: <https://github.com/rheem-ecosystem/rheem-benchmark>.

it provides a set of configuration files whereby users can add new keywords to the language together with their mappings to RHEEM operators. As a result, users can easily adapt RheemLatin for their applications. Listing 2 illustrates how one can express our SGD example with RheemLatin.

```

1 import '/sgd/udfs.class' as taggedPointCounter;
2 lines = load 'hdfs://myData.csv';
3 points = map lines -> {taggedPointCounter.parsePoints(lines)};
4 weights = load taggedPointCounter.createWeights();
5 final_weights = repeat 50 {
6   sample_points = sample points -> {taggedPointCounter.getSample()}
7   with broadcast weights;
8   gradient = map sample_points ->
9     {taggedPointCounter.computeGradient()};
10  gradient_sum_count = reduce gradient -> {gradient.sumcount()};
11  weights = map gradient_sum -> {gradient_sum_count.average()} with
12    platform 'JavaStreams';
13 store final_weights 'hdfs://output/sgd';

```

Listing 2: SGD task in RheemLatin.

The user starts by importing all her required UDFs (Line 1). She then parses all the data points from the input dataset (Lines 2 and 3) and initializes the weights (Line 4). Next, she proceeds to perform the core of SGD: she takes a sample of data points (Line 6), computes the gradient for each sampled data point (Line 7), updates the weights (Lines 8 and 9), and repeats the process 50 times (Line 5). She can also repeat such a core process until convergence by using `WhileLoop` instead of `Repeat`. Optionally, she can specify the platform for any part of her query. For instance, she might know that updating the weights on each iteration is a lightweight computation and hence might specify to use `JavaStreams` (Line 9). She finishes by storing the final weights on HDFS (Line 10).

5.3 Rheem Studio

Although the native APIs and RheemLatin cover a large number of users, some might still be unfamiliar with programming and data-flow languages. Also, some other users may simply desire to speed up the process of composing their data analytic tasks. To this end, RHEEM provides a visual IDE (RHEEM Studio) where users can compose their data analytic tasks in a *drag and drop* fashion [44]. Figure 8 shows the RHEEM Studio’s GUI. The GUI is composed of four parts: a panel containing all RHEEM operators, the drawing surface, a console for writing RheemLatin queries, and the output terminal. The right-side of Figure 8 shows how operators are connected for an SGD plan. The studio provides default implementations for any of the RHEEM operators, which enables users to run common data analytic tasks without writing code. Yet, expert users can provide a UDF by double-clicking on any operator.

Users can draw such a plan by simply *dragging* as many RHEEM operators as required from the left-side panel and *dropping* them on the drawing surface. They consequently connect the operators as required by their data analytic task. The right-side of Figure 8 shows how operators are connected for SGD. While connecting operators, the studio validates such connections and gives feedback to users in case that a connection cannot be established, e.g., the output and input of two connected operators are of different data types. Last but not least, the studio provides default implementations for any of the RHEEM operators, which enables users to run common data analytic tasks without writing a

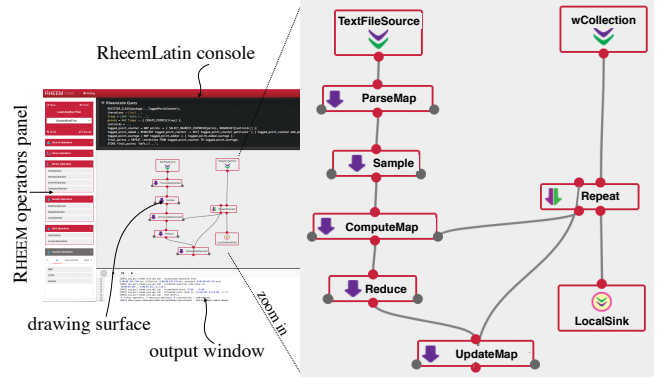


Figure 8: SGD task in the Rheem Studio.

single line of code. Yet, expert users can provide a UDF by double-clicking on any operator.

6. EXAMPLES OF RHEEM PLANS

We now provide in detail three examples of how users can implement their tasks using the Scala native API and the RheemLatin interface. For this, we consider three popular data analytic tasks: *WordCount* (a well-known aggregate task), *K-means* (a very representative iterative task), and *PolyJoin* (a common task over difference data sources).

Users start their RHEEM plans in Scala with a preamble that defines the context and the platforms to be used, as shown in Listing 3. For the sake of presentation, we do not include this preamble in our Scala code examples below.

```

1 val context = new RheemContext(new Configuration)
2   .withPlugin(Spark.basicPlugin)
3   .withPlugin(JavaStreams.basicPlugin)
4 val plan = new PlanBuilder(context)

```

Listing 3: Preamble in the Scala API.

WordCount is an aggregate task that computes the frequency with which each word appears in a dataset. Listing 4 shows the RheemLatin query for this task: Line 1 imports all the required UDFs, Line 2 loads the input data; Lines 3 and 4 parse the words and convert them into records; Line 5 computes the frequency of each word; and Line 6 stores the final word count on disk. Note that users naturally define the flow of their analytical tasks with RheemLatin. Alternatively, users can implement this task using one of the native APIs of RHEEM. Listing 5 shows the Scala code for this task. Similar to the RheemLatin query, the Scala code keeps the plan composition simple.

```

1 import '/wordcount/udfs.class' as wordcount;
2 lines = load 'hdfs://myWords.txt';
3 words = flatmap lines -> { wordcount.splitWords() };
4 tuples = map words -> { wordcount.convert2Tuple() };
5 adds = reduce tuples -> { wordcount.getWord() }, tuples -> {
6   wordcount.reduce() };
7 store adds '/output/wordcount';

```

Listing 4: Word Count task in RheemLatin.

K-means is a widely used ML task for clustering data points together according to their similarity. We show the RheemLatin query in Listing 6. In contrast to the Word-Count task, this task is iterative (Lines 4–7). We observe

```

1 val words = plan.readTextFile("hdfs://myWords.csv")
2   .flatMap(_.split("\\W+"))
3   .map(word => (word.toLowerCase, 1))
4   .reduceByKey(_+_1, (c1, c2) => (c1._1, c1._2 + c2._2))
5   .collect()

```

Listing 5: Word Count task using the Scala API.

that defining loops in RheemLatin is quite similar to coding in a high-level language (e.g., Scala), which makes it intuitive for most users. Listing 7 shows its counterpart in Scala.

```

1 lines = load 'hdfs://myPoints.txt';
2 points = map lines -> kmeans.parsePoints();
3 centroids = load 'hdfs://myInitialCentroids.txt';
4 final_centroids = repeat centroids AS current_centroid for 50 {
5   distance = map points -> kmeans.selectNearestCentroid() with
6     broadcast current_centroid;
7   centroids_sum = reduce distance -> kmeans.reduce();
8   new_centroids = map centroids_sum -> kmeans.average(); }
9 store final_centroids 'hdfs://output/kmeans';

```

Listing 6: K-means task in RheemLatin.

```

1 val points = plan.readTextFile("hdfs://myPoints.csv")
2   .map(createPoints)
3 val initialCentroids =
4   plan.loadCollection(Kmeans.createRandomCentroids(k))
5 val finalCentroids = initialCentroids .repeat( iterations , {
6   currentCentroids =>
7   val newCentroids = points.mapJava(
8     new SelectNearestCentroid,
9   )
10   .withBroadcast(currentCentroids, "centroids")
11   .reduceByKey(_+_centroidId, +)
12   .map(_average newCentroids)}
13 finalCentroids .collect()

```

Listing 7: K-means task using the Scala API.

PolyJoin is a common task in polystore scenarios, i.e., joining several datasets from different data sources. In this case, we consider the TPC-H Q5 and assume that: the region, suppliers, and customer relations are on Postgres; the nations relations is on the local file system; and the orders and lineitem relations are on HDFS. Despite the complexity of this query, we observe that the RheemLatin query (Listing 8) and the Scala (Listing 9) are still simple as they follow the logical flow of the task itself. Lines 1-7 in Listing 8 load the dataset, Lines 8-12 select and project the required tuples, and Lines 13-22 join the resulted tuples before making the group-by in Line 23.

7. RHEEM VS. MUSKETEER

We experimentally compare RHEEM with its closest competitor, Musketeer [32]. More experiments concerning the optimizer can be found in [40].

Setup. We ran our experiments on a cluster of 10 machines. Each node has one 2 GHz Quad Core Xeon processor, 32 GB main memory, 500 GB SATA hard disks, a 1 Gigabit network card and runs 64-bit platform Linux Ubuntu 14.04.05. In RHEEM we used the following platforms: Java’s Stream library (JavaStreams), Spark 1.6.0 (Spark), Flink 1.3.2 (Flink), GraphX 1.6.0 (GraphX), Giraph 1.2.0 (Giraph), a Java graph library (JGraph), and HDFS 2.6.0 to store files. We used all these platforms with their default settings and configured

```

1 import 'polyjoin/udfs.class' as polyjoin;
2 region = load 'postgres:///tpch/region';
3 suppliers = load 'postgres:///tpch/suppliers';
4 customers = load 'postgres:///tpch/customers';
5 nations = load 'file:///nations' delimiter '|';
6 orders = load 'hdfs:///orders' delimiter '|';
7 lineitems = load 'hdfs:///lineitems' delimiter '|';
8 region.filter = filter region [1] == 'ASIA';
9 region.project = map region.filter -> { polyjoin.projectRecord(0, 1) };
10 suppliers.project = map suppliers -> { polyjoin.projectRecord(0, 3) };
11 customers.project = map customers -> { polyjoin.projectRecord(0, 3) };
12 order.filter = filter orders -> { polyjoin.isBetween( 4,
13   '1994-01-01', '1995-01-01' ) };
14 join1 = join nation [2], region.project [0];
15 map.join1 = map join1 -> { polyjoin.tuple2Record(0, 0, 0, 1) };
16 join2 = join map.join1 [0], customers.project [1];
17 map.join2 = map join2 -> { polyjoin.tuple2Record(0, 0, 0, 1, 1, 0) };
18 join3 = join map.join2 [2], order.filter [0];
19 map.join3 = map join3 -> { polyjoin.tuple2Record(0, 0, 0, 1, 1, 0) };
20 join4 = join map.join3 [2], lineitems [0];
21 map.join4 = map join4 -> { polyjoin.tuple2Record(0, 0, 0, 1, 1, 2, 1,
22   5, 1, 6) };
23 join5 = join map.join4 -> { polyjoin.record2Tuple(2, 0) },
24   suppliers.project -> { polyjoin.record2Tuple(0, 1) };
25 map.join5 = map join5 -> { polyjoin.tuple2Record(0, 1, 0, 3, 0, 4) };
26 groupBy = groupby map.join5[0];
27 store groupBy 'output/polyjoin';

```

Listing 8: PolyJoin task in RheemLatin.

```

1 val regions : DataQuanta[Record] =
2   plan.readTable("postgres:///tpch/region")
3   .map( createRecord())
4   .filter ((r: Record) => r.getString(1) == "ASIA")
5   .map( projectRecord(., 0, 1))
6 val suppliers : DataQuanta[Record] =
7   plan.readTable("postgres:///tpch/supplier")
8   .map[Record]( createRecord())
9   .map[Record]( projectRecord(., 0, 3))
10 val customers : DataQuanta[Record] =
11   plan.readTable("postgres:///tpch/customer")
12   .map[Record]( createRecord())
13   .map( projectRecord(., 0, 3))
14 val nations : DataQuanta[Record] = plan.readTextFile("file:///nation")
15   .map( createRecord())
16 val orders : DataQuanta[Record] = plan.readTextFile("hdfs:///order")
17   .map( createRecord())
18   .filter ( isBetween(., 4, fromDate, toDate) )
19 val lineitems : DataQuanta[Record] =
20   plan.readTextFile("hdfs:///lineitem")
21   .map(createRecord())
22 nations
23   .join( getColumn(., 2), regions, getColumn(., 0))
24   .map( tuple2Record(., 0, 0, 0, 1))
25   .join( getColumn(., 0), customers, getColumn(., 1))
26   .map( tuple2Record(., 0, 0, 0, 1, 1, 0))
27   .join( getColumn(., 2), orders, getColumn(., 1))
28   .map( tuple2Record(., 0, 0, 0, 1, 1, 0))
29   .join( getColumn(., 2), lineitems, getColumn(., 0))
30   .map( tuple2Record(., 0, 0, 1, 1, 2, 1, 5, 1, 6))
31   .join [Record, Tuple2[String, String]]( record2Tuple(., 2, 0),
32     suppliers, record2Tuple(., 0, 1))
33   .map( tuple2Record(., 0, 1, 0, 3, 0, 4))
34   .groupByKey((r: Record) => r.getField(0))
35   .collect()

```

Listing 9: PolyJoin task using the Scala API.

the maximum RAM of each platform to 20 GB. We disabled the RHEEM stage parallelization feature to have only one single platform running at any time. We obtained all the cost functions required by our optimizer as described in Section 4.5. We considered the cross-community pagerank task (CrocoPR), because the authors reported this task to be a case where Musketeer chooses multiple platforms. Note that, for fairness reasons, we perform the data preparation part of CrocoPR (i.e., union the different communities

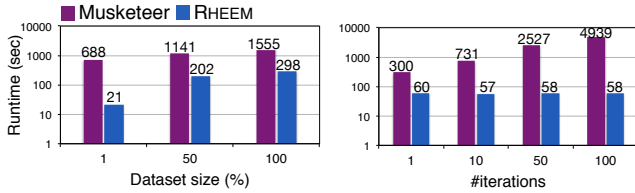


Figure 9: Rheem outperforms Musketeer by more than one order of magnitude.

pages) as a separate script for Musketeer. This is because its language (Mindi) is not optimized for dealing with UDFs, thereby it would be much slower to provide the data preparation as a UDF. In contrast, RHEEM seamlessly performs both parts (data preparation and page rank) as a single task. We used the DBPedia pagelinks dataset (20 GB).

Results. Figure 9 shows the results in log scale when varying the dataset sizes for 10 iterations and the number of iterations for 10% of the dataset. Overall, we observe the superiority of RHEEM over Musketeer, especially as the number of iterations increases: RHEEM is up to 85 times faster than Musketeer. Note that, in contrast to Musketeer, RHEEM keeps its runtime constant as the number of iterations increases. This is because: (i) Musketeer, among other things, checks dependencies, compiles and package the code, and writes the output to HDFS at each iteration (or stage), which comes with a high overhead; (ii) RHEEM executes the page rank part of the task (i.e., after the data preparation) on `JavaStreams`, which allows it to perform each iteration with almost zero overhead.

8. LIMITATIONS

As of now, RHEEM does not support any stream processing platforms. While users can easily supply new batch processing platforms, stream processing requires to extend RHEEM’s core. We plan to do so by following the lambda architecture paradigm [46]. In addition, RHEEM currently relies on the fault-tolerance of the underlying platforms and is, thus, susceptible to failures while moving data across platforms. We plan to incorporate some basic fault-tolerance mechanism at the cross-platform level. Other remaining issues include: adding methods that speed up inter-platform communications, such as the one proposed in [33], integrating RHEEM with resource managers to incorporate changes in the availability of computing resources, and supporting simultaneous execution of RHEEM jobs.

9. RELATED WORK

The research and industry communities have proposed a myriad of different data processing platforms [5, 8, 11, 18, 25, 59]. In contrast, we do not provide a data processing platform but a novel system on top of them.

Cross-platform data processing has been in the spotlight only very recently. Some works focus only on integrating different data processing platforms with the goal of alleviating users from their intricacies [1, 2, 10, 12, 29]. However, they still require expertise from users to decide when to use a specific data processing platform. For example, BigDAWG [29] requires users to specify where to run tasks via its `Scope` and `Cast` commands, which already require expertise from users. Only few works share a similar goal with us [28, 32, 43, 55, 58]. However, they substantially differ from RHEEM. Two main

differences are that they consider neither data movement costs nor progressive task optimization techniques, although both aspects are crucial in cross-platform settings. Additionally, each of these works differs from RHEEM in various ways. As Musketeer’s main goal is to decouple front-end languages (e.g., SQL and PigLatin) from the underlying platforms [32], it is not as expressive and extensible as RHEEM. Furthermore, as it maps task patterns to specific underlying platforms, it is not clear how one can efficiently map a task when having similar platforms (e.g., Spark vs. Flink or Postgres vs. MySQL). Similarly, in Myria [58], it is hard to allocate tasks when having similar platforms because it comes with a rule-based optimizer. Additionally, its rule-based optimizer also makes it hard to maintain. IRES [28] supports only 1-to-1 mappings between abstract tasks and their implementations, which limits expressiveness and optimization opportunities. Moreover, it assumes direct data movement paths between platforms, which is hard to maintain for many platforms. QoX focuses only on ETL workloads [55]. DBMS+ [43] is limited by the expressiveness of its declarative language and hence it is neither adaptive nor extensible. Other complementary works focus on improving data movement across different platforms [33] or libraries by using a common intermediate representation and executing the scripts in LLVM [49], but none of them address the cross-platform optimization problem. Tensorflow [15] follows a similar idea, but for cross-device execution of machine learning tasks and thus it is orthogonal to RHEEM. In fact, RHEEM could use TensorFlow as an underlying platform.

The research community has also studied the problem of federating relational databases [54]. Garlic [22], TSIMMIS [23], and InterBase [21] are just three examples. However, all these works significantly differ from RHEEM in that they consider a single data model and simply push query processing to where the data is. Other works integrate Hadoop with an RDBMS [27, 41], however, one cannot easily extend them to deal with more diverse tasks and platforms.

10. CONCLUSION

Given today’s data analytic ecosystem, supporting cross-platform data processing has become rather crucial in organizations. We have identified four different situations in which an application requires or benefits from cross-platform data processing. Driven by these cases, we built RHEEM, a cross-platform system that decouples applications from data processing platforms to achieve efficient task execution over multiple platforms. RHEEM follows a cost-based optimization approach for splitting an input task into subtasks and assigning each subtask to a specific platform, such that the cost (e.g., runtime or monetary cost) is minimized. Our experience while building RHEEM raised several interesting questions that need to be addressed in the future, namely: *How can we (i) reduce the inter-platform data movement costs? (ii) address the cardinality and cost estimation problem? (iii) efficiently support fault-tolerance across platforms? (iv) add new platforms automatically? and (v) improve data exploration in cross-platform settings?*

11. REFERENCES

- [1] Apache Beam. <https://beam.apache.org>.
- [2] Apache Drill. <https://drill.apache.org>.
- [3] Apache Flink. <https://flink.apache.org>.

- [4] Apache Flume. <https://flume.apache.org/index.html>.
- [5] Apache HBase. <http://hbase.apache.org/>.
- [6] Apache Hive: A data warehouse software for distributed storage. <http://hive.apache.org>.
- [7] Apache Mahout. <http://mahout.apache.org>.
- [8] Apache Spark: Lightning-Fast Cluster Computing. <http://spark.incubator.apache.org/>.
- [9] Fortune magazine. <http://fortune.com/2014/06/19/big-data-airline-industry/>.
- [10] Luigi Project. <https://github.com/spotify/luigi>.
- [11] PostgreSQL. <http://www.postgresql.org/>.
- [12] PrestoDB Project. <https://prestodb.io>.
- [13] Spark MLlib: <http://spark.apache.org/mllib>.
- [14] Spark SQL programming guide. <http://spark.apache.org/docs/latest/sql-programming-guide.html>.
- [15] M. Abadi et al. TensorFlow: A System for Large-Scale Machine Learning. In *OSDI*, pages 265–283, 2016.
- [16] D. Agrawal, L. Ba, L. Berti-Equille, S. Chawla, A. Elmagarmid, H. Hammady, Y. Idris, Z. Kaoudi, Z. Khayyat, S. Kruse, M. Ouzzani, P. Papotti, J.-A. Quiané-Ruiz, N. Tang, and M. Zaki. Rheem: Enabling Multi-Platform Task Execution. In *SIGMOD*, pages 2069–2072, 2016.
- [17] D. Agrawal et al. Road to Freedom in Big Data Analytics. In *EDBT*, pages 479–484, 2016.
- [18] A. Alexandrov et al. The Stratosphere platform for big data analytics. *VLDB J.*, 23(6):939–964, 2014.
- [19] A. Baaziz and L. Quoniam. How to use big data technologies to optimize operations in upstream petroleum industry. In *21st World Petroleum Congress*, 2014.
- [20] M. Boehm, M. Dusenberry, D. Eriksson, A. V. Evfimievski, F. M. Manshadi, N. Pansare, B. Reinwald, F. Reiss, P. Sen, A. Surve, and S. Tatikonda. SystemML: Declarative Machine Learning on Spark. *PVLDB*, 9(13):1425–1436, 2016.
- [21] O. A. Bukhres et al. InterBase: An Execution Environment for Heterogeneous Software Systems. *IEEE Computer*, 26(8):57–69, 1993.
- [22] M. J. Carey et al. Towards Heterogeneous Multimedia Information Systems: The Garlic Approach. In *RIDE-DOM*, pages 124–131, 1995.
- [23] S. S. Chawathe et al. The TSIMMIS Project: Integration of Heterogeneous Information Sources. In *IPSJ*, pages 7–18, 1994.
- [24] M. Dallachiesa, A. Ebaid, A. Eldawy, A. K. Elmagarmid, I. F. Ilyas, M. Ouzzani, and N. Tang. NADEEF: a commodity data cleaning system. In *SIGMOD*, pages 541–552, 2013.
- [25] J. Dean and S. Ghemawat. MapReduce: Simplified Data Processing on Large Clusters. *Communications of the ACM*, 51(1), 2008.
- [26] D. Deng, R. C. Fernandez, Z. Abedjan, S. Wang, M. Stonebraker, A. K. Elmagarmid, I. F. Ilyas, S. Madden, M. Ouzzani, and N. Tang. The Data Civilizer System. In *CIDR*, 2017.
- [27] D. J. DeWitt, A. Halverson, R. V. Nehme, S. Shankar, J. Aguilar-Saborit, A. Avanes, M. Flasz, and J. Gramling. Split query processing in polybase. In *SIGMOD*, pages 1255–1266, 2013.
- [28] K. Doka, N. Papailiou, V. Giannakouris, D. Tsoumakos, and N. Koziris. Mix 'n' match multi-engine analytics. In *IEEE BigData*, pages 194–203, 2016.
- [29] A. J. Elmore et al. A Demonstration of the BigDAWG Polystore System. *PVLDB*, 8(12):1908–1911, 2015.
- [30] W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional Functional Dependencies for Capturing Data Inconsistencies. *ACM Transactions on Database Systems (TODS)*, 33(2):6:1–6:48, 2008.
- [31] R. C. Fernandez, D. Deng, E. Mansour, A. A. Qahtan, W. Tao, Z. Abedjan, A. K. Elmagarmid, I. F. Ilyas, S. Madden, M. Ouzzani, M. Stonebraker, and N. Tang. A Demo of the Data Civilizer System. In *SIGMOD*, pages 1639–1642, 2017.
- [32] I. Gog et al. Musketeer: all for one, one for all in data processing systems. In *EuroSys*, 2015.
- [33] B. Haynes, A. Cheung, and M. Balazinska. PipeGen: Data Pipe Generator for Hybrid Analytics. In *SoCC*, pages 470–483, 2016.
- [34] A. Hems, A. Soofi, and E. Perez. How innovative oil and gas companies are using big data to outmaneuver the competition. Microsoft White Paper, <http://goo.gl/2Bn0xq>, 2014.
- [35] IBM. Data-driven healthcare organizations use big data analytics for big gains. White paper, <http://goo.gl/AFIHpk>.
- [36] Z. Kaoudi, J.-A. Quiane-Ruiz, S. Thirumuruganathan, S. Chawla, and D. Agrawal. A Cost-based Optimizer for Gradient Descent Optimization. In *SIGMOD*, 2017.
- [37] Z. Khayyat, I. F. Ilyas, A. Jindal, S. Madden, M. Ouzzani, P. Papotti, J. Quiané-Ruiz, N. Tang, and S. Yin. BigDansing: A System for Big Data Cleansing. In *SIGMOD*, pages 1215–1230, 2015.
- [38] Z. Khayyat, W. Lucia, M. Singh, M. Ouzzani, P. Papotti, J. Quiané-Ruiz, N. Tang, and P. Kalnis. Lightning Fast and Space Efficient Inequality Joins. *PVLDB*, 8(13):2074–2085, 2015.
- [39] S. Kruse, Z. Kaoudi, J.-A. Quiané-Ruiz, S. Chawla, F. Naumann, and B. Contreras-Rojas. RHEEMix in the Data Jungle – A Cross-Platform Query Optimizer. arXiv: 1805.03533 <https://arxiv.org/abs/1805.03533>, 2018.
- [40] S. Kruse, Z. Kaoudi, J.-A. Quiané-Ruiz, S. Chawla, F. Naumann, and B. Contreras-Rojas. RHEEMix in the Data Jungle – A Cross-Platform Query Optimizer. arXiv: 1805.03533 <https://arxiv.org/abs/1805.03533>, 2018.
- [41] J. LeFevre, J. Sankaranarayanan, H. Hacigümüs, J. Tatemura, N. Polyzotis, and M. J. Carey. MISO: soup up big data query processing with a multistore system. In *SIGMOD*, pages 1591–1602, 2014.
- [42] V. Leis et al. How good are query optimizers, really? *Proc. VLDB Endow.*, 9(3):204–215, 2015.
- [43] H. Lim, Y. Han, and S. Babu. How to Fit when No One Size Fits. In *CIDR*, 2013.
- [44] J. Lucas, Y. Idris, B. Contreras-Rojas, J.-A. Quiané-Ruiz, and S. Chawla. Cross-Platform Data Analytics Made Easy. In *ICDE*, 2018.
- [45] V. Markl, V. Raman, D. Simmen, G. Lohman,

- H. Pirahesh, and M. Cilimdžić. Robust query processing through progressive optimization. In *SIGMOD*, pages 659–670, 2004.
- [46] N. Marz and J. Warren. *Big Data: Principles and best practices of scalable realtime data systems*. Manning, 2015.
- [47] M. Mitchell. *An introduction to genetic algorithms*. MIT press, 1998.
- [48] C. Olston, B. Reed, U. Srivastava, R. Kumar, and A. Tomkins. Pig Latin: A Not-so-foreign Language for Data Processing. In *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, SIGMOD '08, pages 1099–1110, 2008.
- [49] S. Palkar, J. J. Thomas, A. Shanbhag, M. Schwarzkopt, S. P. Amarasinghe, and M. Zaharia. Weld: A Common Runtime for High Performance Data Analysis. In *CIDR*, 2017.
- [50] A. Pavlo, E. Paulson, A. Rasin, D. J. Abadi, D. J. DeWitt, S. Madden, and M. Stonebraker. A Comparison of Approaches to Large-Scale Data Analysis. In *SIGMOD*, pages 165–178, 2009.
- [51] J.-A. Quiané-Ruiz and Z. Kaoudi. Cross-Platform Query Processing. In *ICDE (tutorial)*, 2018.
- [52] P. J. Sadalage and M. Fowler. *NoSQL distilled: A brief guide to the emerging world of polyglot persistence*. Addison-Wesley Professional, 2012.
- [53] S. Shankar, A. Choi, and J.-P. Dijkstra. Integrating Hadoop Data with Oracle Parallel Processing. Oracle White Paper, <http://www.oracle.com/technetwork/database/bi-datawarehousing/twp-integrating-hadoop-data-with-or-130063.pdf>, 2010.
- [54] A. P. Sheth and J. A. Larson. Federated Database Systems for Managing Distributed, Heterogeneous, and Autonomous Databases. *ACM Computing Surveys*, 22(3):183–236, 1990.
- [55] A. Simitsis, K. Wilkinson, M. Castellanos, and U. Dayal. Optimizing Analytic Data Flows for Multiple Execution Engines. In *SIGMOD*, pages 829–840, 2012.
- [56] M. Stonebraker. The Case for Polystores. <http://wp.sigmod.org/?p=1629>, 2015.
- [57] D. Tsoumakos and C. Mantas. The Case for Multi-Engine Data Analytics. In *Euro-Par*, pages 406–415, 2013.
- [58] J. Wang, T. Baker, M. Balazinska, D. Halperin, B. Haynes, B. Howe, D. Hutchison, S. Jain, R. Maas, P. Mehta, D. Moritz, B. Myers, J. Ortiz, D. Suciu, A. Whitaker, and S. Xu. The Myria Big Data Management and Analytics System and Cloud Services. In *CIDR*, 2017.
- [59] F. Yang, J. Li, and J. Cheng. Husky: Towards a More Efficient and Expressive Distributed Computing Framework. *PVLDB*, 9(5):420–431, 2016.