# Towards unified secure on- and off-line analytics at scale

P. Coetzee *, M. Leeke, S. Jarvis

*Department of Computer Science, University of Warwick, United Kingdom*

**A B S T R A C T**

Data scientists have applied various analytic models and techniques to address the oft-cited problems of large volume, high velocity data rates and diversity in semantics. Such approaches have traditionally employed analytic techniques in a streaming or batch processing paradigm. This paper presents CRUCIBLE, a first-in-class framework for the analysis of large-scale datasets that exploits both streaming and batch paradigms in a unified manner. The CRUCIBLE framework includes a domain specific language for describing analyses as a set of communicating sequential processes, a common runtime model for analytic execution in multiple streamed and batch environments, and an approach to automating the management of cell-level security labelling that is applied uniformly across runtimes. This paper shows the applicability of CRUCIBLE to a variety of state-of-the-art analytic environments, and compares a range of runtime models for their scalability and performance against a series of native implementations. The work demonstrates the significant impact of runtime model selection, including improvements of between $2.3\times$ and $480\times$ between runtime models, with an average performance gap of just $14\times$ between CRUCIBLE and a suite of equivalent native implementations.

© 2014 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY license (http://creativecommons.org/licenses/by/3.0/).

## 1. Introduction

To derive insight and provide value to organisations, data scientists must make sense of a greater volume and variety of data than ever before. In recent years this challenge has motivated significant advances in data analytics, ranging from streaming analysis engines such as IBM's InfoSphere Streams [1], Backtype's Storm [2] or Yahoo!'s S4 [3], to an ecosystem of products built on the MapReduce [4] framework.

When data specialists set out to perform analyses they are typically faced with a decision: they can opt to receive continuous insight but limit analytic capabilities to a functional or agent-oriented streaming architecture, or make use of a bulk data paradigm but risk batch analyses taking hours or even days to complete. It is, of course, possible to maintain systems that target streamed and batch paradigms separately, though this is less desirable than having a single system with the semantics to account for those paradigms in a unified manner. The need to support multiple methodologies presents a further challenge; ensuring analyses are correct and equivalent across platforms. These issues are further complicated by deployment scenarios involving multi-tenant cloud systems or environments with complex access control requirements.

Our research seeks to alleviate many of the issues highlighted above through the development of CRUCIBLE, a framework consisting of (i) a domain specific language (DSL) for describing analyses as a set of communicating sequential processes, (ii) a common runtime model for analytic execution in multiple streamed and batch environments, and (iii) an approach which automates the management of cell-level security labelling uniformly across runtimes. In particular, this paper demonstrates

how CRUCIBLE (named after the containers used in chemistry for high-energy reactions) can be used across multiple data sources to perform highly parallel distributed analyses of data simultaneously in streaming and batch paradigms, efficiently delivering integrated results whilst making best use of existing cloud infrastructure.

The specific contributions of this paper are as follows:

- A high-level DSL and suite of runtime environments, adhering to a common runtime model, that provide consistent execution semantics across on- and off-line data. This is the first DSL designed specifically to target the execution of on- and off-line analytics with equal precedence.
- The development of a new primitive in the developed DSL that permits a single analytic to be run equivalently over multiple data sources: locally, over Accumulo data, and over files in the Hadoop Distributed File System (HDFS).
- A novel framework for the semi-automated management of cell-level security, applied consistently across runtime environments, enabling the management of data visibility in on- and off-line analysis.
- An evaluation of the performance of CRUCIBLE on a set of best-in-class runtime environments, demonstrating framework optimisations that result in an average performance gap of just 14× when compared to a suite of native implementations.

The remainder of this paper is structured as follows: Section 2 provides a summary of related work. Section 3 introduces the CRUCIBLE system and describes its abstract execution model. Section 4 presents a performance analysis and discussion of the CRUCIBLE runtimes and associated optimisations. Finally, Sections 5 and 6 provide avenues for further research and conclude the paper.

## 2. Related work

Large-scale data warehousing technologies abound in the literature, many of which are based on Google's MapReduce [4] and Bigtable [5], such as Hadoop [6] and HBase[7], as well as NSA's Accumulo[8], which added cell-level security, increased fault tolerance (FATE), and a novel server-side processing paradigm [9]. Tools such as Google's Drill [10], and the Apache Software Foundation implementation Dremel [11], promise SQL-like interactive querying over these Bigtable-backed frameworks. Hive [12] and Pig [13] both aim to permit definition of analytics over arbitrarily formatted data in Hadoop [6], while Cascading [14] takes a slightly more engineer-centric approach to definition of analytics over Hadoop.

Some of the more common projects in the streaming analytics space are IBM's InfoSphere Streams [1], and the open source Storm [2], developed by BackType and now an Apache Incubator project. Others include Yahoo!'s S4 [3] (also in the Apache Incubator), which offers an agent-based programming model. This makes deployment scenarios and performance prediction somewhat more challenging than Storm and Streams, which offer a lower-level abstraction. Esper [15] provides a cross-platform API for Java and .NET, and Microsoft's StreamInsight [16] product offers tight integration with Microsoft SQL Server.

Most of these technologies facilitate execution of an analytic over a single paradigm, be it online or offline. Recently, researchers have begun to translate offline analytics into an online paradigm. SAMOA [17] aims to enable Machine Learning using a streaming processing paradigm to both validate and update models in near-real-time. AT&T Research, as part of their Darkstar project [18], have constructed a hybrid stream data warehouse, DataDepot [19]. This uses online techniques to perform analysis on data as it arrives at the data warehouse, updating the contents of the bulk data store in the process. The closest research to CRUCIBLE to date has been in IBM DEDUCE [20], which defines code for MapReduce using SPADE (Stream Processing Application Declarative Engine), the programming language used in early versions of InfoSphere Streams. This permits a unified programming model (*e.g.*, allowing use of common operators), but does not offer any direct execution equivalence between a MapReduce job and an equivalent Streams SPADE job. Furthermore, SPADE is now deprecated in favour of SPL (Stream Processing Language).

CRUCIBLE builds on the most desirable attributes of these approaches in order to offer a single framework for developing secure analytics to be deployed at scale on state of the art multi-tenancy on- and off-line data processing platforms. It offers a similar programming model and approach to task parallelism as the likes of Storm and Streams, while offering consistent execution semantics across both on- and off-line data. It includes a semi-automated framework for management of security labels, and permits the application of these labels equivalently across data sources.

## 3. CRUCIBLE system

In order to facilitate the creation of advanced analytics for on- and off-line distributed execution, the CRUCIBLE DSL makes use of a higher level language abstraction than typical analytic frameworks, such as [2], [11], or [14]. This enables a degree of portability that is not typically achievable under other schemes; an engineer may write their analytic once, in a concise high-level language, and execute across a variety of paradigms without knowledge of runtime-specific implementation details. In addition, the user is afforded the ability to exploit an array of best-in-class runtime models for the execution of CRUCIBLE code.

Furthermore, this approach seeks to free domain specialists from concerns of correctness and security. The CRUCIBLE runtimes are responsible for ensuring that analytics are run with equivalent execution semantics, through adherence to CRUCIBLE's execution model, thus providing assurances of cross-platform correctness. The domain-specific nature of the CRUCIBLE language permits the user a greater degree of confidence that the analytic they *intend* is the analytic they have