

発表概要

Hadoop 上で動作する Sawzall サブセットの実装

中 田 秀 基^{†1} 井 上 辰 彦^{†1,†2} 工 藤 智 宏^{†1}

(平成 23 年 1 月 20 日発表)

Sawzall は、Google が 2006 年に発表した大容量データの並列バッチ処理に適した言語である。Sawzall の計算モデルは MapReduce 型の分散演算であるが、リダクション操作を組み込みの Aggregator に限定することで、エンドユーザによる容易な記述を可能にしている。我々は現在開発中の並列データ処理機構上の言語処理系を開発するための 1 ステップとして、Scala 言語による Sawzall 言語のサブセット処理系を実装した。文法やセマンティクスに関しては明確な定義がなかったため、2006 年の論文をベースに推測した。その結果、最近公開された Sawzall 処理系とは機能的に若干の相違がある。構文解析に Scala 言語の Parser Combinator を用いることで、処理系の記述量が削減できた。現在の実行対象処理系は Hadoop である。Hadoop の Mapper 上で言語インタプリタを動作させ、Reducer 上では我々の提供する Aggregator を動作させる。Scala は Java VM 上で動作することから、Java で記述される Hadoop 上での実行は容易である。本発表では、本処理系の実装について詳しく述べる。さらに、Hadoop で直接記述した場合と、プログラム量および実行速度の点で比較を行う。比較の結果、プログラム量は大幅に小さくなる一方、実行速度の面でも一定のオーバーヘッドがあることが確認された。

An Implementation of Sawzall Subset Interpreter Working on Hadoop

HIDEMOTO NAKADA,^{†1} TATSUHIKO INOUE^{†1,†2}
and TOMOHIRO KUDOH^{†1}

Sawzall is a script language designed for batch processing of large amount of data, which is introduced by Google in 2006. The processing model of Sawzall is the MapReduce. Sawzall allows programmers only to program 'mappers' to ease the burden. Sawzall provides a set of 'built-in aggregaters', from which programmer choose mapping function. We are developing distributed data processing system for large scaled data. As a part of the project, we have implemented an interpreter for Sawzall subset in Scala language. We employed parser

combinator for lexical parsing. Currently, the system is running on Hadoop. In the paper, we provide detailed implementation of the system. We also evaluated the system with naked Hadoop in terms of program size and execution speed. We confirmed that, with Sawzall, program size is much smaller, while there are certain overhead in terms of execution speed.

†1 独立行政法人産業技術総合研究所

National Institute of Advanced Industrial Science and Technology

†2 株式会社創夢

SOUM Corporation