

分散ファイルシステムによる並列データ I/O 測定 An Examination of Data I/O Speed on a Parallel Data Storage System

村田 健史^{1*}; 渡邊 英伸¹; 鶴川 健太郎²; 村永 和哉²; 鈴木 豊²; 建部 修見³; 田中 昌宏³; 木村 映善⁴
MURATA, Ken T.^{1*}; WATANABE, Hidenobu¹; UKAWA, Kentaro²; MURANAGA, Kazuya²; YUTAKA, Suzuki²; TATEBE, Osamu³; TANAKA, Masahiro³; KIMURA, Eizen⁴

¹ 情報通信研究機構, ² 株式会社 セック, ³ 筑波大学, ⁴ 愛媛大学

¹National Institute of Information and Communications Technology, ²Systems Engineering Consultants Co., LTD., ³University of Tsukuba, ⁴Ehime University

現在、多くの科学研究分野ではデータのほとんどがデジタル化され、その量および種類は大規模化の一途をたどっている。これからますます大規模化・複雑化するデータ指向型科学時代を踏まえて、ビッグデータ処理がより容易に、また一元的行うことができるクラウドシステムが求められている。

NICT サイエンスクラウドは、地球惑星科学を含む様々な科学研究データおよびソーシャルデータのためのクラウドシステムである。NICT サイエンスクラウドでは (1) データ伝送・データ収集機能、(2) データ保存・データ管理機能、(3) データ処理・データ可視化機能の3つの柱(機能)から構成されている。それぞれの機能についての基盤技術を開発するだけでなく、複数の基盤技術を組み合わせることでシステム化を行うことができる。システムを実際に科学研究に応用・適用することで、様々な分野でのビッグデータ科学・データインテンシブ科学が可能となる。

本研究では、NICT サイエンスクラウド上で科学研究のビッグデータ処理を行うための基盤技術について議論する。データサイズが大きい場合にクラウドデータ処理で解決すべき問題点の一つはデータ I/O である。例えば、100MB/sec で 100TB のデータを読み出すとすると、1,000,000 秒(約 11.5 日)かかる。すなわち、大規模科学データを処理するためには、高速 I/O 技術が不可欠である。本発表では、並列ファイルシステム (GPFS) と分散ファイルシステム (Gfarm) の2つのシステムでのデータ読み出し速度の比較を行い、それらのスケーラビリティを比較する。

