

# データ転送量を考慮した Hadoop 性能改善方式の提案と評価

## A Study of Data Transfer Aware Hadoop Task scheduling

学籍番号：201221609

氏名：渡邊 飛雄馬

Hyuma WATANABE

現在広く使われている大規模データ分散処理フレームワークである Hadoop は MapReduce と呼ばれる分散処理の仕組みを採用している。MapReduce の Task は大量のデータ入出力を発生させるため、TaskTracker の HDD I/O やネットワーク I/O の輻輳を引き起こし、しばしばその処理速度が低下する。ところが、既存の Hadoop の Task scheduler は TaskTracker の I/O 使用率に拘わらず、空き Task slot の有無のみで Task 割り当ての判断を行ってしまう。このため、クラスタ内で I/O リソース使用率の不均衡が発生し、Hadoop の性能が低下する。

本論文では、既存の Task scheduler を用いた予備実験を行い、Hadoop の Job 実行において、データ I/O を原因として Job 実行速度が低下する仕組みを示す。

さらに、その結果に基づき、TaskTracker のネットワーク I/O 使用率と HDD I/O 使用率を考慮した Task 割り当てを行う事で Job 実行時間を短縮する Hadoop Task scheduling 方式を複数提案する。また、研究室内に構築した小規模な Hadoop クラスタと Amazon EC2 上に構築した中規模な Hadoop クラスタを用いて提案方式の性能を評価する。提案した Task scheduling 方式を組み合わせることで、平均 6% の Job 実行時間短縮を実現した。

研究指導教員：川原崎 雅敏  
副研究指導教員：森嶋 厚行