

課題名(タイトル):

通信・ファイル I/O に関する研究開発

利用者氏名:

○石川 裕\*、高木 将通\*、堀 敦史\*、Gerofi Balazs\*、森江 善之\*、畑中 正行\*、小倉 崇浩\*、Martsinkevich Tatiana\*、亀山 豊久\*

理研における所属研究室名:

\*計算科学研究センター フラッグシップ 2020 プロジェクト システムソフトウェア開発チーム

<p>1. 本課題の研究の背景、目的、関係するプロジェクトとの関係</p>	<p>Tofu2のオフロード機能を用いて実装可能かどうかの確認を行った。まず、Tofu2 インターコネク ト・ハードウェアのオフロード機能をモデル化し、 butterfly アルゴリズムの通信の各段で通信相手 が変わってゆく場合に、通信ハードウェア資源が 0(1)になるような手法を設計した。さらにこの手 法に基づいて、butterfly アルゴリズムの barrier、allgather (recursive doubling) およ び alltoall でプロトタイプ実装を行った。</p>
<p>本課題では、FX100 が有する Tofu2 ネットワー クインターフェイスを用いてポスト「京」向けに 開発している通信・ファイル I/O 機能のプロトタ イプ実装および性能検証を実施することを目的 とする。我々は、ポスト「京」向けの OS、MPI ラ イブラリ等のシステムソフトウェアの研究開発 を行なっている。このようなシステムソフトウェ アの開発は、ターゲット計算機のハードウェア開 発と並行して進めている。ポスト「京」で採用さ れるネットワークインターフェイス TofuD は、 FX100 の Tofu2 ネットワークの後継版であり互換 性を有する。FX100 を用いることによりポスト 「京」向けの通信に関するソフトウェア実装実 験が可能となる。</p>	<p>これらの実装により、演算を伴うリダクションを 除く永続型集団通信を OFI の fi_trigger インター フェイスを用いてオフロード実装が可能であるこ とを確認した。なお、Tofu2にはバリア同期時にリ ダクション同様の演算を行う機能が提供されてお り、リダクション型集団通信はこの機能で実現で きる。</p>
<p>2. 具体的な利用内容、計算方法</p> <p>今回、ポスト「京」に向けて MPI 通信ライブラ リで規格化が進んでいる永続型集団通信の設計、 評価を行った。</p>	<p>4. まとめ</p> <p>昨年度のツリー型のアルゴリズムの実装であ った Trinaryx3 アルゴリズム に加え butterfly アルゴリズムを用いた集団通信も 0(1)の通信ハ ードウェア資源消費で実装可能であることを確 認し、省資源型のオフロード可能な永続型集団通 信関数群を提供できる見込みがたった。</p>
<p>Tofu2 および TofuD にも適用可能な Open Fabrics Interfaces (OFI)が提供する fi_trigger と呼ばれるネットワークオフロードインターフェ イスを用いて broadcast, allgatherなどを例に永 続型集団通信アルゴリズムをオフロード可能であ ることを示した。</p>	<p>5. 今後の計画・展望</p> <p>今回実装していない永続型集団通信についても同 様に設計、評価を行う。これらの設計、評価を元 にポスト「京」向けの永続型集団通信の適切な実 装を行う予定である。</p>
<p>3. 結果</p> <p>昨年度(平成29年度)開発した Trinaryx3 アル ゴリズムの broadcast 永続型集団通信の実装を、 OFI の fi_trigger インターフェイスで書き直した。 barrier、allgather、alltoallなどの永続型集 団通信は butterfly アルゴリズムで実装できる。</p>	

平成 30 年度 利用研究成果リスト

**【口頭発表】**

森江 善之, 畑中 正行, 高木 将通, 堀 敦史, 石川 裕, “Open Fabrics Interfaces オフロード API を用いた MPI 永続型集団通信の設計と評価”, 情報処理学会研究報告, 2018(HPC-168), 2018 年 3 月, 石川県加賀市