T2R2東京工業大学リサーチリポジトリ Tokyo Tech Research Repository

論文 / 著書情報 Article / Book Information

論題(和文)	TSUBAME2.0におけるMulti-rail InfiniBandネットワークの性能評価
Title(English)	
著者(和文)	野村 哲弘, 遠藤 敏夫, 松岡 聡
Authors(English)	Akihiro Nomura, Toshio Endo, SATOSHI MATSUOKA
出典(和文)	
Citation(English)	, Vol. 2012-ARC-194/HPC-137, ,
発行日 / Pub. date	2012, 12
権利情報 / Copyright	ここに掲載した著作物の利用に関する注意:本著作物の著作権は(社))情報処理学会に帰属します。本著作物は著作権者である情報処理学 会の許可のもとに掲載するものです。ご利用に当たっては「著作権法」ならびに「情報処理学会倫理綱領」に従うことをお願いいたします 。 The copyright of this material is retained by the Information Processing Society of Japan (IPSJ). This material is published on this web site with the agreement of the author (s) and the IPSJ. Please be complied with Copyright Law of Japan and the Code of Ethics of the IPSJ if any users wish to reproduce, make derivative work, distribute or make available to the public any part or whole thereof.

TSUBAME2.0 における Multi-rail InfiniBand ネット ワークの性能評価

野村 哲 $\mathbf{U}^{1,2,\mathrm{a}}$ 遠藤 敏夫¹ 松岡 聡¹

概要:TSUBAME2.0のネットワークは Fat tree トポロジであるものの,大規模実行時に集団通信性能が 劣化することが観測されている.本稿では想定される原因としてスイッチ間リンクにおけるパケット衝突 とスイッチ間リンクの性能劣化に着目し,それぞれの問題を緩和するネットワーク設定を提示し,バンド 幅および集団通信性能への影響を示す.ネットワーク設定の改善の結果,通信の確率的な遅延の発生をほ ぼなくすことができ,大規模実行時のインジェクションバンド幅において 16.0%~39.5%の性能向上を確 認した.

1. 背景

2010年に東京工業大学に導入された TSUBAME2.0 スー パーコンピュータ [1] は, 日本初の 1PFlops を超える計 算機であり, 4000 基以上の GPU を備える GPU スパコ ンという面で知られているが,2つの独立した2段Fat tree トポロジの InfiniBand QDR ネットワークで構成され る Full-bisection Multi-rail ネットワークで 1408 台の計算 ノードが接続されている点も特徴である.Fat tree ネット ワークでは, Edge スイッチ-ノード間のリンクの本数と同 数以上のリンクで Edge スイッチ–Core スイッチ間を接続 することで,スイッチ間リンクが通信ボトルネックとな ることを防いでいる.一方, Multi-rail ネットワークとは, 全ノードが複数のネットワークに所属することで,ネット ワークの数に比例した通信性能を得ることができるもので あり, TSUBAME2.0 では I/O ノードへの接続の有無を除 いて相似形である2つの InfiniBand QDR ネットワークを 同時に使用することによって,任意の2ノード間で理論上 は80Gbpsの通信性能を得ることができる.これらの性質 から, TSUBAME2.0 は全体全通信やランダムな1対1通 信が多く発生するアプリケーションにおいて Torus などの 他のトポロジを持つコンピュータに比べて高い性能を示す ことが期待される.

しかしながら,実際の TSUBAME2.0 における大規模実 行時には,集団通信の性能が想定している理論性能よりも

© 2012 Information Processing Society of Japan

低くなる現象および, 普段は1秒以下で終了するサイズの 集団通信において確率的に5秒以上遅延する現象が経験的 に知られており, 通信性能がボトルネックとなっているア プリケーションにおいてネットワークが致命的な実効性能 低下要因となっている.

2. TSUBAME ネットワークの通信性能評価

上記のように経験的に知られていた通信性能の劣化につ いて,2種類のマイクロベンチマークを用いて通信性能を 計測することで,実際にどの程度の通信性能劣化が発生し ているかを確認した.以下全ての実験は,表1に示す環境 (TSUBAME2.0のThin Node)において,以下のような条 件下で通信性能を評価した.

- 1ノードあたり1プロセスを配置した.
- ランク番号の割り当てはノード番号順とした.同一 ラック,同一 Edge スイッチ配下のノードのランク番 号は連続している.
- 環境変数 MV2_NUM_HCAS を設定することで, MVA-PICH2 に複数 rail を使った通信を行うよう指定した.
- 2.1 ランダムペアの Sendrecv 性能

1つ目のベンチマークとして,MPI_Alltoall 関数の1 フェーズの通信を模してバンド幅を計測するベンチマーク を実行した.通信する全プロセスをランダムに組み合わせ, 各ペアにおいて一定メッセージサイズのMPI_Sendrecvを 同時に行うことにより,バイセクション通信性能を測定し た.各ノード間で通信の実行時間にばらつきが発生する が,集団通信をこのようなフェーズごとの通信として実装 した場合には,フェーズ間では一番遅いプロセスの遅延が

東京工業大学 学術国際情報センター Global Scientific Information and Computing Center, Tokyo Institute of Technology
² JST, CREST

a) nomura.a.ac@m.titech.ac.jp

ノード数	1408 (うち 1300 台を使用)		
ネットワーク	Dual Rail Infini Band 4x QDR (40Gbps x 2)		
トポロジ	2 段 Fat Tree		
スイッチあたりのノード数	14 もしくは 16		
CPU	Intel Xeon E5670 x 2 (2.93GHz)		
OS	SuSE Linux Enterprise Server		
MPI ライブラリ	MVAPICH2 1.8		
OFED ドライバ	MLNX_OFED_LINUX-1.5.3-3.0.0		
サブネットマネージャ	OpenSM 3.3.9.MLNX_20111006_e52d5fc-0.1		

表 1 評価環境の緒元 (TSUBAME Thin ノード)

各プロセスに伝播していき,最終的な集団通信の実行時間 に対する支配項となるため,今回の実行においては一番実 行に時間がかかったペアの実行時間を通信実行時間として 定め,実行結果は各ノードのインジェクションバンド幅に 換算した.なお,すべての実験において疑似乱数シードを 固定することで同じノード数の実験における通信相手の組 み合わせを固定し,通信相手の組み合わせによる性能への 影響を排除している.

図1にメッセージサイズを512MiBとしてプロセス数を 変えたときの最低インジェクションバンド幅の推移を示す. バンド幅が高いほど通信性能が良いと言える.40ノードま での実行ではほぼ理論性能である7.5GB/sの性能が出てい るが,ノード数が増えるにつれて最良時の最低インジェク ションバンド幅が5.2GB/s,3.9GB/s,3.1GB/s,2.6GB/s と段階的に低下している様子が観測される.これは,Edge スイッチをまたぐ通信が支配的になるにつれて,スイッチ 間リンクにおける通信の衝突が発生しやすくなるためと思 われる.性能低下が離散的である原因は,1つのスイッチ 間リンクの通信の多重度が離散的であるためと思われる.

参加ノード数が600を超えると、最低インジェクション バンド幅が500MB/sを下回る試行が出現するようになり、 ノード数が増えるにつれてその頻度は増加していること が観察された.本稿ではメッセージサイズ512MiBの試行 のみを図示しているが、他のメッセージサイズにおける実 行時にも傾向は変化しなかった.実験時間が限られていた ため、以降の本ベンチマークについてはメッセージサイズ 512MiBの場合のみを計測している.

2.2 Alltoall 性能

2 つめのベンチマークとして MPI_Alltoall の実行時間 を測定した.図2に1通信あたりのメッセージサイズを 1MiB としたときの MPI_Alltoall の実行時間の推移を示 す.実行時間が短いほど通信性能が良いと言える.最低イ ンジェクションバンド幅ベンチマークと同様に,参加ノー ド数が 600 ノードを超えると異常値(通信が通常時の数倍 かかるケース)が発生することが分かる.これらの異常値 を排除しても(図3),通信性能にばらつきが大きいことが わかる.また,最低インジェクションバンド幅ベンチマー











クと同様にノード数が増えると実効バンド幅が低下してい ることが読み取れる.



図 4 Fat Tree におけるルーティングの衝突

3. 性能劣化の要因

今回,性能劣化の要因としてルーティング戦略によるものおよび,不調リンクの存在によるものの2つを想定した.以下にそのそれぞれについて性能が低下する原因を述べる.

3.1 ルーティング戦略

我々は,段階的な性能劣化の原因として,図4に示す ように,2段Fat TreeのEdgeスイッチとCoreスイッチ 間のリンクが有効に使われておらず,通信の衝突が起こっ ていると推量した.理想的には図5のように,各通信が 上流のリンクの間で完全にバランスされてリンク速度を使 い切る通信ができるのであるが,InfiniBandのルーティン グはパケットの送信先ごとに次ホップのスイッチを固定す る静的ルーティングであるため,すべての通信パターンに おいて通信が衝突しないルーティングを行うことは不可能 である.また,集団通信に頻出する通信パターンにおいて のみルーティングを最適化することも考えられるが,実際 の運用では故障ノードの発生によって歯抜けとなるノード が出現するため,この方法は大規模計算機環境では破綻す る.実際に今回の実験中にも複数台の計算ノードがダウン して,実行対象のノードリストから取り除かれている.

InfiniBand でのルーティングテーブルはサブネットマ ネージャが管理しており,TSUBAME2.0 では OpenSM 3.3.9 の UpDn ルーティング戦略に基づいて決定されてい る.今回は,実験時に最新であった OpenSM 3.3.15 に附 属する以下のルーティング戦略を用いて,ルーティング戦 略による通信性能の変化を観察した.

- UpDn: TSUBAME で通常利用されている戦略である.ツリー状のトポロジを共通する祖先に到達するまで送信元および送信先から辿り,得られた経路のうち最短のものを採用する.
- MinHop: 送信元と送信先を結ぶ最短経路のうち任意の経路を次ホップの転送先として選択する
- Fat Tree: 完全 Fat Tree を仮定して通信の衝突を避けるように次ホップの転送先を送信先に応じて順に割り



Vol.2012-ARC-202 No.3 Vol.2012-HPC-137 No.3

2012/12/13

図 5 Fat Tree における理想的なルーティング



図 6 UpDn および DFSSSP におけるネットワーク利用効率のシ ミュレーション結果

振ることで通信の衝突を回避する.

 DFSSSP(Deadlock Free Single-source Shortest Path)[2]: 全ホストで全対全の通信経路を作成した時 に負荷が完全にバランスするように経路を構成する.

図 6 はネットワークシミュレータ ORCS[3] による UpDn および DFSSSP におけるバイセクションバンド幅のシミュ レーション結果である.この結果より,ルーティング戦略 を変化させることで通信性能が向上することが期待できる.

3.2 不調リンク

TSUBAME2.0の大規模ネットワークにおいてはスイッ チ間リンクにケーブル異常やスイッチポート異常に起因す る不調なリンクが発生していることが,ibdiagnet コマン ドにおけるポートの速度およびパフォーマンスカウンタの 値から判明した.これらのエラーは一時的なもの(一旦リ ンクを再起動することによって復帰する)と恒久的なもの の両方があり,場合によっては全く通信できなくなるわけ ではなく,速度低下や大量のパケットロスを起こすものの 通信できてしまうものがある.そのようなリンクがネット ワーク上に存在すると,そのリンクを用いる通信だけが遅 延することにより,全体の通信のボトルネックとなってし まうことが推察される.以下にTSUBAME2.0で観察され た主な不調リンクの症状を示す.

 速度低下: TSUBAME2.0 のネットワークにおける 1 ポートあたりの通信速度は 4x QDR(QDR データレー トのリンク 4 本分)の 40Gbps であるが, 1x QDR や

Vol.2012-ARC-202 No.3
Vol.2012-HPC-137 No.3
2012/12/13

表 2	ベンチマーク条件
-----	----------

ラベル	ルーティング戦略	不調リンク無効化
minhop	MinHop	N
updn	UpDn	N
ftree	Fat Tree	N
hetero	Fat Tree / DFSSSP	Y
updn2	UpDn	Y

4x SDR(いずれも 10Gbps) に縮退してしまっているリ ンクが発生した .

 異常パケットの発生: TSUBAME2.0のネットワーク 全域で全く通信を行っていない状態においても, symbol_error_counterや port_rcv_errorsのようにパケット 破損が発生しているときに上昇するカウンタが秒間数 百パケットのオーダで上昇していることを確認した.

後者については何らかの原因でパケットが無限ループ している,制御パケットが異常発生している,もしくはパ ケットを正常にエンコード・デコードできなくなったこと などが原因と考えられる.いずれにせよ,当該リンクを通 過するパケットは(場合によっては確率的に)遅延もしくは 破損して正常に届かなくなると考えられる.そこでネット ワーク中のこれらの以上リンクを個別に停止することで, Fat Tree のトポロジを多少犠牲にして安定したネットワー クを構成した.

4. 性能劣化の検証と改善

前節で述べた2つの原因および解決策を実際に実行して 2節に述べたそれぞれのベンチマークにおける性能の変化 を観察した.本来はそれぞれの解決策を切り分けて実行す べきであるが,実験時間の制約や,上記の結論に至った過 程および不調リンクの無効化が不可逆な操作である点を理 由に個別に切り分けた実験を行うことはできなかった.ま た,DFSSSPはTSUBAME2.0ネットワークの1st railに 適用した際に,正常な性能が発揮できず,全通信にかかる 時間が10倍以上となりネットワークが不安定になってし まったため,1st railにdfssspを用いた実験は行っていな い.表2に今回実行した実験の条件を示す.

heteroにおいては、1st rail に Fat Tree 戦略、2nd rail に DFSSSP と、異なる戦略を採用した.TSUBAME2.0の Multi-rail ネットワークは I/O ノードへの接続の有無を除 いて相似であり、同じルーティング戦略を用いることで相 似形のルーティングテーブルが作成され、性能低下する通 信パターンも同様のものとなることが考えられる.そのた め、各 rail におけるルーティング戦略を違えることによっ て弱点となる通信パターンが分散して相互に通信性能を補 完しあうことが期待される.

全ベンチマークの実行結果を図 7,図 9および図 10 に 示す.ルーティング戦略の選択にかかわらず,3.2 節に示 したような異常値が観察されていることがわかる.不調リ



ンクの排除後の実験のみを示した図 8,図 11 と比較すれ ばわかるように,速度低下やカウンタ異常を起こしている リンクを排除することによってこれらの異常値のほとんど が発生しなくなり,図3と比較して異常値以外の部分にお いても性能が安定するようになった.

また,図7および図8に示す通り,hetero戦略において他のルーティング戦略と比べて16.0%~39.5%最低インジェクションバンド幅の向上がみられた.グラフの形状を観察することにより,hetero以外のトポロジにおいて,ノード数およびノードの組み合わせによって性能低下の幅が振動していたものが,hetero戦略の採用によってその分の性能低下を防げるようになったと推論することができる.

図 10 は, MPI_Alltoall の性能比較である. Alltoall 通 信においては通信相手を切り替えながら Sendrecv 相当の 通信を (ノード数-1) 回行っている.そのため,不調リンク 排除前の実験では通信時間のぶれが蓄積し,通信性能が安 定していないことが分かる.他方,図 11 に示す不調リンク 排除後の実験では,わずかな異常値を除いて通信性能は極 めて安定していることがわかる.Alltoall の性能比較では, ルーティング戦略の選択に起因する性能の差は確認できな かった.この点について,何故最低インジェクションバン ド幅と違う傾向が見られるか解明することは今後の課題で ある.なお,ノード数が 600 を超える部分での Alltoall の バンド幅 1.74GB/s であり,最低インジェクションバンド 幅の約半分である.ノード数 200 付近の hetero における 性能劣化を含めて,何故 600 ノード超で性能劣化が起こる かの解明も今後の課題である.

5. おわりに

我々は,TSUBAME2.0 における大規模実行時の通信性 能低下がどの程度発生しているかを計測し,スイッチ間リ ンクの混雑と性能劣化に着目し,これらを解消することで 通信速度の向上を図った.理論性能からの平均性能の乖離 については rail ごとに異なるルーティング戦略を用いるこ とでマイクロベンチマークにおいて最悪インジェクション バンド幅の平均値で 16.0% ~ 39.5%の性能向上を,確率的



図 8 最低インジェクションバンド幅の分布 (不調リンク排除後)







図 10 図 9 から異常値を除いたもの



図 11 Alltoall 通信実行時間の分布 (不調リンク排除後)

な性能劣化についてはネットワーク上の不調リンクを遮断 することで観測されなくなるという結果を得た. 今後は今回の実験時に発生したバグのために実験の継続 を断念したルーティング戦略の評価を行い,より理論性能 に近い実行性能を得るとともに,シミュレーション上の性 能と実際の実行時の性能のギャップおよび,マイクロベン チマークと実際の集団通信の性能のギャップについて原因 を明らかにする必要がある.また,実アプリケーションで の性能の変化も比較して,これらの処置がアプリケーショ ンの性能向上に資することを示す必要がある.

なお,東京工業大学学術国際情報センターでは今回の実 験の成果をもとに,TSUBAME2.0の運用において不調リ ンクの検出を強化して通信性能の劣化を未然に防ぐ運用を 2012年9月より行っている.

謝辞 TSUBAME2.0 における InfiniBand の性能調査 のため,2012年8月8日~10日の間の48時間,TSUB-AME2.0 のネットワークを占有しての実験を行わせていた だきました.本期間中TSUBAME2.0の利用を控えていた だいた全てのユーザに感謝いたします.

また,本実験の実施時には,NEC,Mellanox,Torsten Hoefler 博士および,Jens Domke 氏に多数の助言およびご 協力をいただきました.ここに感謝いたします.

参考文献

- [1] 東京工業大学学術国際情報センター: TSUBAME 計算サー ビス, http://tsubame.gsic.titech.ac.jp/.
- [2] Domke, J., Hoefler, T. and Nagel, W.: Deadlock-Free Oblivious Routing for Arbitrary Topologies, *Proceedings* of the 25th IEEE International Parallel & Distributed Processing Symposium (IPDPS), IEEE Computer Society, pp. 613–624 (2011).
- [3] Schneider, T., Hoefler, T. and Lumsdaine, A.: ORCS: An Oblivious Routing Congestion Simulator, Technical Report 675, Indiana University (2009).