

科学技術研究向け超高速大域ネットワーク基盤

平木敬 稲葉真理 菅原豊

情報理工学系研究科コンピュータ科学専攻

概要

ネットワークの性能を引き出すためには、時として、階層型モデルによる非効率さが枷となるケースがある。2004年度、我々はレイヤー間協調方式を提案し、日本・北米大陸・ヨーロッパ大陸にまたがるネットワーク転送実験を行い、internet2によるバンド幅距離積の世界記録(Land Speed Record)を更新した。また超高速文字列照合機構をFPGAの上に実装したネットワークカードを作成し、10 Gbpsを超える処理速度を達成、レイヤー間協調方式の有用性を示した。

1. はじめに

光スイッチ技術、WDM等、近年のめざましい技術革新に伴い、ネットワーク環境が整備された。例えば、日米間、あるいは、欧米間の通信はOC192を日常的に利用できるようになってきている。このような環境でネットワークのぎりぎりの性能を引き出すためには、階層型モデルによる非効率さが枷となるケースが多々見られる。特に高速ネットワークに対応するためにハードウェア化を行う際には、ハードウェアとソフトウェアの切り分けをどう行うかの設計が重要であり、たとえば、階層をまたいだハードウェア実装を行うとで、大幅な性能向上が見られることがある。我々は、これをレイヤー間協調方式(inter-layer co-ordination)と呼ぶことにする。我々は、レイヤー間協調方式として、データ転送速度を改善するための、複数ストリーム協調方式、パケット間ギャップ調整

方式、TRC-TCP方式を、またTCPストリーム用文字列照合ハードウェアを提案してきた。本稿では、2004年度に行った仕事のうち、ハードウェアTRC-TCPの実装および日本とヨーロッパ大陸、アメリカ大陸をつなぐ遠距離実験、および、書き換え可能であるFPGAを持つ2ポート10 Gbpsイーサネットカードと、そのFPGA上に実装されたTCPストリームを認識する超高速文字列照合機構について述べる。

2. TRC-TCP 10 Gbps ハードウェア

ネットワーク通信で標準的に使われるTCPは、データの信頼性のため、送信済み受信未確認のデータ("inflight data")をバッファリングしており、このTCPウィンドウと呼ばれるバッファサイズの調整により動的にスピードの調整を行っているが、遠距離高速通信では、性能が十分に出ないことが知られている。2003年度、我々は、TCPのウィンドウサイズによって定まる速度に、ネットワークインターフェースのパケット送出速度を合致させるTRC-TCP方式で、性能が著しく改善されることを1 Gbpsのイーサネットを終端とする、約8 Gbpsの日米間の実験線を用いて実証した。本年度は急速に広まりつつある10 Gbpsのネットワークに対応するため、すなわち、10 Gイーサネットでスタンダードフレームを利用する場合パケット送出間隔が約1.3 μ 秒、に対応するため、ハードウェアTRC-TCPを既存のネットワークカード上に実装し、日米欧を

またぐ10 Gbpsのネットワークの上で実験を行った。

今回、我々は一般的に入手可能な標準的部品を組み合わせたサーバと、インテリジェント10 Gbpsイーサネットカード Chelsio T110(図1)を組み合わせることでシステムを構成した。この場合のボトルネックは、ネットワークではなく、サーバのI/Oバスとなる。ボトルネックスピードは、PCI-X 1.0 64bit-133MHzで、およそ8.5Gbpsとなる。



図1. Chelsio T110

Chelsio T110は512MBパケットバッファを持ちフルTCPオフロード機能を持つため、サーバのCPU負荷が少ないという特徴を持ち、データバスはASIC制御コードはFPGAという、プログラマブルなインテリジェントイーサネットカードで、近距離通信では7.5Gbpsを達成する。我々はT110のバッファのキューイングをシンプルにすることで、シングルストリームによるデータ転送時に再高速が達成できるように最適化したうえで、TCPによって定まる速度とイーサネットパケット送出速度の調整を行った。

3. 遠距離データ転送実験

実験はRioworks HDAMA rev.Eのマザーボードに、2.2GHz AMD Opteron248、1GB PC3200 DDR SDRAMを搭載したサーバにLinux 2.6.6 for x86_64を使用、iperf v1.7.0を使って、二点間でデータ転送を行った。MTU

1500Bのスタンダードと7832Bおよび8192Bのジャンボフレームでの実験を行った。バッファサイズは、RTT 250msecのとき、384MB、約500msecのとき、448MBとした。また複数ストリームを用いディスク上のファイルデータ転送実験も行った。実験は(1)東京から北米経由CERN(2)CERNから東京経由ピッツバーグの実際のネットワークを用いて行った。ここでは(2)CERNピッツバーグの実験結果を中心に実験結果を述べる。

3.1 東京 - 北米 - CERN

2004年10月 - 日本、カナダ、米国、オランダとCERN研究所(スイス・ジュネーブ)の研究者が共同して、日本からスイスまでの世界最長の10ギガビットイーサネット回線を構築した。ネットワークの全長は約18,500Km、17のタイムゾーンを通過しており、10GigabitイーサネットWANPHY技術を用い光を用いたSONET/SDH機器の相互運用により、東京大学からCERN研究所のサーバを恰も同一LAN上であるかのように接続を行った。2004年度最初の遠距離データ転送実験は、この全経路WANPHYネットワークを利用して行われた。2004年度最初の遠距離データ転送実験は、このWANPHYネットワークを利用して行われ、スタンダードフレームを用いて平均7.57Gbpsのデータ転送を達成した。

3.2 CERN 東京 ピッツバーグ

2004年11月 東京大学とWIDEプロジェクトの研究者が日本、カナダ、米国、オランダとCERN研究所(スイス・ジュネーブ)の研究者との協力により、米国ピッツバーグから日本を経由してスイス・ジュネーブまでの世界最長の10ギガビット回線を構築し、ピッツバーグのSC2004国際会議会場に設置したシステムからCERN研究所に設置したシステム間

を接続した。ネットワークの全長は約 31,248 Km あり、17 のタイムゾーンを通過している。ピッツバーグの SC2004 国際会議会場に設置された東京大学の計算機から、SCinet を経由し、Abilene, StarLight を経由し APAN/JGN2 ネットワークで東京に到達し、WIDE Project が運用する T-LEX に接続し、T-LEX から米国シアトルまでは TycoTelecommunication が IEEAF に寄贈した回線を用い、Pacific Northwest Gigapop の機器に接続した。シアトルからは、回線は CA*net4 の専用ラムダによってシカゴの StarLight まで運ばれ SURFnet のシカゴ-アムステルダム間のラムダに接続され、アムステルダムの NetherLight に接続された。NetherLight と CERN 研究所の間は、SURFnet のアムステルダム-ジュネーブのラムダによって接続されている。ネットワークの往復遅延時間(RTT)は約 437 ミリ秒である。(図 2)



図 2 ネットワーク構成図

実験では、まず一本の TCP ストリームを用いて、データ転送を行い 7.21 ギガビット/秒のデータ転送を実現した。データ転送には 1500 バイトの標準イーサネットフレーム長を用い、データ転送は 20 分間行った。図 3 に シングルストリームのデータ転送レートを示す。 定常状態においては、PCI-X の理論的限界値の 95 % 程度が継続的に達成されており、TRC-TCP が有効であることがわかる。このパ

ンド幅・距離積は 225,298 テラビットメートル/秒となる。

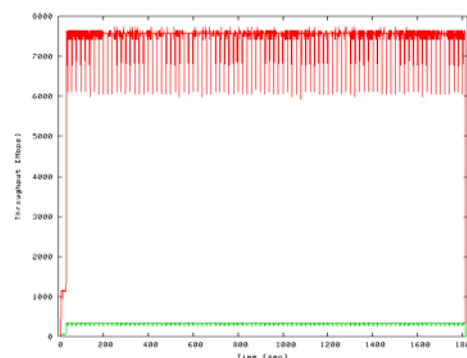


図 3 単一ストリーム転送レート

次にマルチストリームによる自己競合のあるシステムにおける、ディスク上のデータファイルの転送実験で、TRC-TCP 方式の有無による比較を示す。

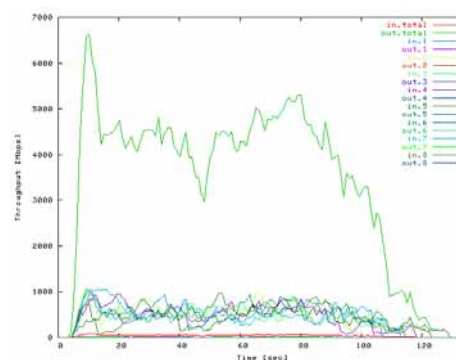


図 4 自己競合 TRC-TCP 無し

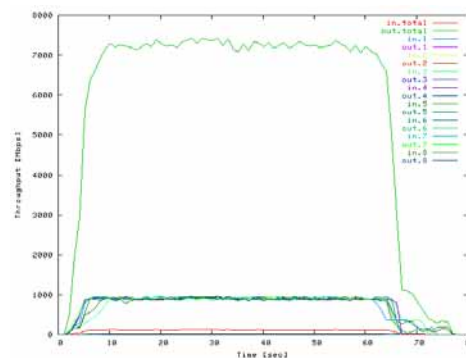


図 5 自己競合 TRC-TCP 有り

マルチストリームを用いたシステムにおいては、自己競合による性能低下が問題になることが知られている。図 4、図 5 は、TRC-TCP の無・有による、ストリームごと、そしてシステム全体のパフォーマンスを示している。

TRC-TCPにより、バースト性を取り除くことで、著しい安定化およびその結果の性能向上が見られる。ここでは、8台構成のシステムで、1.6Gbpsのディスク間データ転送を達成した。

4. プログラマブル10Gbps イーサカード

プログラマブルなFPGAを持ち複数の10Gbpsイーサネットポートを持つネットワークカードを開発し(図6)FPGAの上に(1)TCPストリームを認識する超高速文字列照合機構および(2)10Gbpsパケットヘッダロガーを実装した。ここでは文字列照合について述べる。



図6. プログラマブルネットワークカード

近年、Webデータのデータマイニングあるいはネットワークセキュリティーのために、高速文字列照合の要求が高まっている。特にセキュリティーに関しては、照合パターンを必要に応じて変更する必要がある。またエンドノードでの煩雑な処理を無くすためにネットワークの入口でフィルタリングをかければ良いが、そのためには、中間ノードでTCPストリームごとの文字列照合をワイヤレートで行う必要がある。我々はFPGA向きの文字列照合アルゴリズムSBT法を提案、実装し、合計1000文字のルールが10Gbpsで処理できることを確認した。

5. おわりに

ネットワーク高速化に従い、ハードウェア化の重要性は、日々増している。一方、技術の進歩に従いユーザ要求は日々変化することも

多く、プログラマブルなFPGAとASICの組み合わせによるシステムの重要性は今後、ますます増大すると考えられる。

6. 謝辞

ネットワーク実験はWIDEプロジェクト、CERN, SURFnet, StarLight, Neitherlight, Tyco Telecommunications, IEEAF, Pacific NorthwestGigapop, CA*net4networks, APAN富士通コンピュータテクノロジーズ、NTTコミュニケーションズ、東陽テクニカ、Foundry networks, Cisco Systems, Juniper Networks, Clear Sight社の協力のもとに行われた。

参考文献

- [1] K. Hiraki, M. Inaba, J. Tamatsukuri, R. Kurus, Y. Ikuta, H. Koga, and A. Jinzaki, "Data Reservoir: Utilization of Multi-Gigabit Backbone Network for Data-Intensive Research," Proc. Super Computing 2002, (SC2002) CD-ROM, Nov., 2002.
- [2] M. Nakamura, M. Inaba, and K. Hiraki, "Fast Ethernet is sometimes faster than Gigabit Ethernet on LFN--Observation of congestion control of TCP streams" Proc. Int. Conf. on Parallel and Distributed Computing And Systems (PDCS2003) Nov. 2003 pp.854-859
- [3] K. Hiraki, M. Nakamura, J. Senbon, Y. Sugawara, T. Itoh and M. Inaba, "End-node transmission rate control kind to intermediate routers - towards 10Gbps era" (PFLDnet 2004) Feb. 2004
- [4] T.Ito, M. Inaba, "Theoretical Analysis of Performances of TCP/IP Congestion Control Algorithm with Different Distances", Networking 2004 May 2004
- [5] M. Nakamura, H. Kamezawa, J. Tamatsukuri, M. Inaba, K. Hiraki, K. Mizuguchi, K. Torii, S. Nakano, S. Yoshita, R. Kurusu, M. Sakamoto, Y. Furukawa, T. Yanagisawa, Y. Ikuta, J. Shitami, A. Zinzaki "Long Fat Pipe Congestion Control for Multi-Stream Data Transfer" Proceedings of the International Symposium on Parallel Architectures, Algorithms and Networks, ISPAN 2004
- [6] Y. Sugawara, M. Inaba, and Kei Hiraki, "Over 10Gbps String Matching Mechanism for Multi-Stream Packet Scanning Systems", Field-Programmable Logic and Applications, 14th International Conference, FPL 2004
- [7] H. Kamezawa, M. Nakamura, J. Tamatsukuri, N. Aoshima, M. Inaba, K. Hiraki, J. Shitami, A. Jinzaki, R. Kurusu, M. Sakamoto, and Y. Ikuta, "Inter-layer coordination for parallel TCP streams on Long Fat pipe Networks", Super Computing 2004, High Performance Networking and Computing, SC2004