

## メガノードマルチキャストシミュレーションに向けての送信実験報告

三角真<sup>†</sup> 中川晋一<sup>†,††</sup> 我如古津世史<sup>††</sup> 知念賢一<sup>††</sup> 篠田陽一<sup>††</sup>

† 独立行政法人 情報通信研究機構  
〒184-8795 東京都小金井市貫井北町 4-2-1  
†† 北陸先端科学技術大学院大学  
〒923-1292 石川県能美市旭台 1-1

E-mail: †{misumi,snakagaw}@nict.go.jp, ††{tganeko,k-chinen,shinoda}@jaist.ac.jp

**あらまし** インターネットを用いて大量のセンサ情報を複数地点で同時取得する通信方法として、IP マルチキャストが合理的である。しかし、IP マルチキャストを用いて1万を越えるデータを複数地点でルーティングし取得できたという報告は世界的になされていない。評価・実証にはシミュレーション環境が必要だが、実験に供する数万の送信ノード群、高性能のルーティング機構と受信・解析のシステム構築が必要となる。我々は、これらについて送信ノードの多重化(1台の実PCで何台分の仮想ノードが実装可能か)、受信ノードの性能評価(何万の送信ノードからの通信が受信可能か)の検討を行い、ルーティング機構として最も軽いレイヤ2スイッチングを用いて実験環境 StarBED を用いて合計1メガノード(100万)の通信に成功した。実験途上で遭遇した様々な問題と解決法について報告する。

**キーワード** 性能評価, ベンチマーク, ユビキタスコンピューティング, 大規模ネットワークテストベッド

## Experiment Report of Mega-order Sensor Node Multicast Simulation

MISUMI MAKOTO<sup>†</sup>, NAKAGAWA SHIN-ICHI<sup>†,††</sup>, GANEKO TSUYOSHI<sup>††</sup>, CHINEN  
KENICHI<sup>††</sup>, and SHINODA YOICHI<sup>††</sup>

† National Institute of Information and Communications Technology  
4-2-1, Nukui-Kitamachi Koganei, Tokyo 184-8795 Japan  
†† Japan Advanced Institute of Science and Technology  
1-1, Asahidai Nomi, Ishikawa 923-1292 Japan

E-mail: †{misumi,snakagaw}@nict.go.jp, ††{tganeko,k-chinen,shinoda}@jaist.ac.jp

**Abstract** IP multicast is the most efficient transport method for many to many data, especially in case of ten thousand sending and receiving pairs. We managed to succeed as the first case in the world for Mega-order send-receive pairs using L2 IP multicast, process-multiplexing condition at 'StarBED.' In this report, some results and assessments from this experiment were investigated.

**Key words** Performance evaluation, Benchmark, Ubiquitous computing, Large-scale network test bed

### 1. はじめに

インターネットを用いたセンサーネットワークの通信手法として、IP マルチキャストが運用上最も合理的である。微小センサノードからの通信を想定し、複数の受信ノードに対する任意地点へのデータ配送においては、自由度の高い IP マルチキャストが有効である。[1], [4], [6] 従来の IP マルチキャスト通信に関する研究は、放送の様に単一の送信ホストから多数の受信ホストに対するものであり、検討されてきた内容は、主に中継ノード上でのパケットメモリコピーとマルチキャストツリー

構造に関するものなどであり、数万を越えるような打ち上げが行なわれる場合のマルチキャストネットワークの挙動や伝送に関する要件は明らかにされていない。今回、IP マルチキャストのスケールに関する諸問題を明らかにするため、大規模ネットワーク実験設備である StarBED を用いてメガオーダーマルチキャストネットワークの実実験環境を構築した。100 台程度の実 PC を用いて1メガノードを作成し、毎秒合計1メガパケットのパケットの送出を目的として実験を行なう環境構築に成功した。構築したシミュレーション環境の送信ノードの仮想化技術、および実験途上で生じた問題点について報告する。

## 2. 目的と方針

メガオーダーを超えるマルチキャストパケットを打ち上げる(senderが100万をこえる)経験は未知である。ルーチングに関する問題として、スイッチのMACテーブルルックアップの負荷、スイッチングエンジンの負荷、その他、通常の利用では予測できない要因が障害となって、受信側へのパケット到着間隔の延長(ジッタの増大)やパケットロスが生じることが考えられる。また、今回想定するIPマルチキャストは一種のflooding typeのbroadcast trafficであり、送信ノード側のインターフェイスに対しても負荷となる。さらに1個1個のパケット長が100バイト程度であるパケットが毎秒1メガ個到着する場合のデータ量は100メガバイト/秒であり、最速のハードディスク(最速で10Mbytes/sec, 2006年)を用いても記録を取ることは不可能である。従って、ルーチング、送信側インターフェイスに対する余剰トラフィックのフィルタリング、受信機構の負荷分散と記録方法の検討が必要になる。

### 2.1 通信モデル

本研究では、現在想定されている接続されたセンサからの信号を符号化、データパケットを生成、受信側の状態を確認せずにデータを伝送する様な、微小ノードを用いたデータ収集を目的としたセンサネットワークのトラフィックをIPマルチキャスト(今回は特にツリー構造を考える事の無いようにL2,IPv4)を想定した。

### 2.2 仮想化

大規模なIPマルチキャストの通信を想定した場合、1メガノードからのパケットが流れたとしても、それらのエミュレーションの規模(送信数の増大)により、伝送ネットワークにかかる負荷も変化する。本研究の目的は実験環境の構築と評価であり、一つ一つの送信ノードの送信精度、受信ノードの解析性能のベンチマークを行い、送信・受信間隔とパケットロスを尺度として、ネットワークの問題点を探索する必要がある。

特に、実験環境として、実際に100万ノードを実ノード(PC)を用いて用意することは現実的ではなく、仮想ノードを実ノードと送信間隔を基準にした送信精度の低下が許容範囲に収まる事を条件として、可能な限り一つの実ノードから多くのトラフィックを送信させることを検討した。以下、実ノードから多数のノードに通信させる事を仮想化と呼ぶ。

多重化の階層として、MAC層での多重化、IP層での多重化、UDP層での多重化、アプリケーション層での多重化が考えられる。

多重化を行う層が、下位であるほど実際のパケットの挙動に近づき、また、それに伴い、同数のパケットがネットワークに流れたとしてもネットワークの負荷が大きくなると考えられる。

今回は、最も単純に多量のパケットを送出するために、プロセス多重を用いて、UDP層での多重化を行い、パケットの送出行なした。

本実験の様な、定期的にパケットを送出する単純なモデルを仮想化する場合、複数のノードをシミュレートするプログラムを作成するよりも、実際のプログラムそのものを複数動作させ

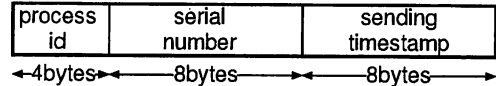


図1 送信パケットのフォーマット

る事が容易であると考えた。そのため、定期的にパケットを送出するプロセスを1台のPCで複数動かさせ、それぞれのプロセスを1つの仮想ノードと見なした。

## 3. 実験

### 3.1 本実験の目的

メガオーダーの仮想ノードを実現し、それらの仮想ノードからパケットの送出・伝送の確認を目的とし、以下の実験を行った。

### 3.2 計測プログラム概要

#### 3.2.1 send\_udp

パケットを送信するためのプログラムで、図1に示すように、データ部分に、送信間隔を調査するための送信時のタイムスタンプ、多重化した際に送信プロセスを一意に判別するためのプロセス識別子、パケットの到着順序及び消失を検知するためのパケットのシリアル番号を持ったUDPパケットを、指定した時間毎に送出する。今回の実験では、1秒毎にパケットを送出する。

なお、多重化のために、複数の送信ノードを単一のPCでシミュレートするために、send\_udpは、指定した数だけプロセスを実行する。

#### 3.2.2 sara

当初、tcpdumpや専用のプログラムを用いてパケットをHDDにダンプするという手法も試みたが、専用のプログラムの場合は、毎秒合計5000ノードからのパケットの受信で0.1%程度、毎秒合計30000ノードからのパケットの受信で1%程度のパケットドロップが認められた。そのため、パケットをキャプチャし、各ノードが送出したパケットの送信間隔と受信間隔の平均値及び標準偏差、ドロップしたパケット数を、ノード毎に算出し記憶し、終了時に各ノード毎の統計情報と、受信したパケット全てのノードの統計情報を出力する。Save and Recording Agent(sara)を作成した。

送信間隔及び受信間隔を求める際には、シリアル番号が非連続な場合、つまり、パケットドロップした場合は、図2に例を示すように、それらの間隔を平均値や標準偏差の算出に利用しない。そのため、パケットのドロップが送受信間隔の平均値と標準偏差に影響を及ぼすことは無い。

#### 3.2.3 ami

A Multicast Introducer(ami)は、受信ノードが、マルチキャストグループに参加するためのプログラム。

### 3.3 実験ネットワークトポロジ

図3に、今回メガノード実験を行った際のネットワークを示す。

Group Aに属する1台のPCをIGMP Querierとし、207

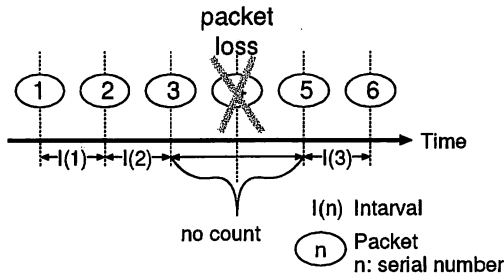


図2 パケットドロップ時の送受信間隔の扱い

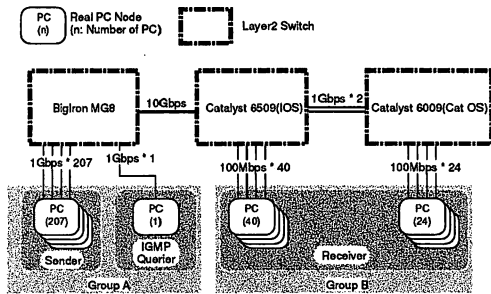


図3 実験ネットワークポロジ

表1 実験に利用したPCのスペック

	GroupA (Sender and IGMP Querier)	GroupB (Receiver)
CPU	Pentium3 1GHz	Pentium3 1GHz
Memory	512MB	512MB
NetIF	1000Base-T	100Base-T
OS	Linux 2.6.13	Linux 2.6.13

台を送信のためのPCとして利用した。また、Group Bに属する64台のPCを受信ノードとして利用した。

全てのPCは同一のLayer 2ネットワーク上に存在する。なお、スイッチ間はTagged VLANを用いて通信する。

### 3.4 計測手法

実験の駆動は、宮地らの設計開発したSpringOS [2], [3]を用いて行った。実験流れを図4に示す。

PCをコントロールするために送信PCと受信PCで、“slave”を実行し、これらのPCをコントロールするためにコントロールPCで“master”を用いる。なお、コントロールPCは、純粹に実験の駆動のみを行い、実験そのものための通信は行わない。

送信PC及び受信PCそれぞれの実験用のシナリオを受け取り、ネットワークを設定する。続いて、受信PCで、マルチキャストグループに参加するためのプログラム(ami)とパケットを受信し統計処理を行うプログラム(sara)を実行する。layer3のネットワークが正常に設定されている事を確認するために、送信PCから受信PCに対して定期的にpingを打つ。pingの応答が有った時点でネットワークの設定が終了したと見なし、“master”に“Network was ready”のメッセージを送る。送信

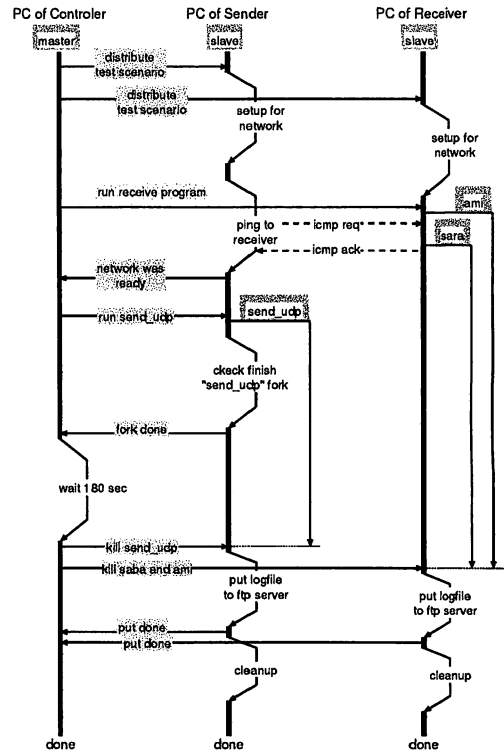


図4 実験駆動の流れ

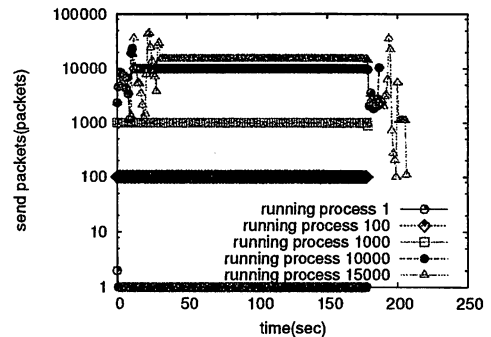


図5 実行プロセス数と1秒当りの送出パケット数

PCでノードをエミュレートするsend\_udpを実行する。図5は、プロセスを実行し始めてからの1秒当りのパケット送出数を実行プロセス数毎に示したグラフであるが、これより多くのプロセスを実行した場合は、全てのプロセスの実行完了までに暫く時間が掛ることが分かる。そのため、全ての仮想ノードが出現してから観測を始めるために、指定した数のプロセスの起動を監視し、起動の完了を“master”に報告する。全ての送信PCで、指定した数の仮想ノードを実行後、180秒間の観測を行う。その後、送信PCでsend\_udpを停止し、また、受信PCでamiとsaraを停止する。送信PCと受信PCは、ログ

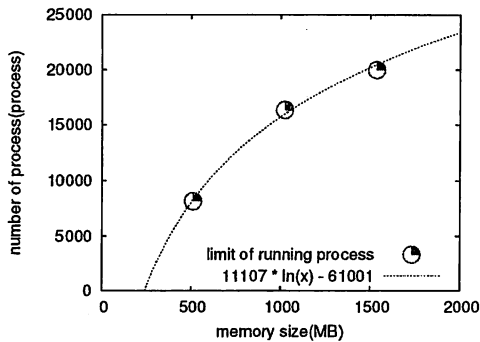


図6 実装メモリの量と最大実行可能プロセス数の関係

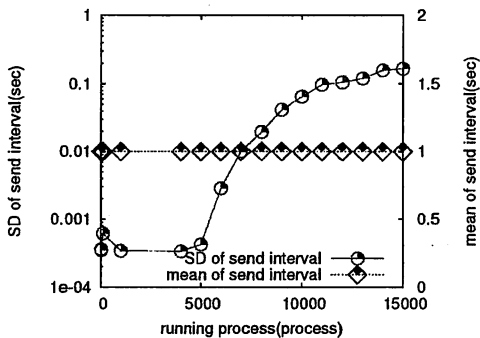


図7 プロセスの多重度と送信間隔の平均値と標準偏差の関係

ファイルを ftp サーバに put し、その作業が終了すると終了した旨を “master” に報告、そして、実験のための設定を消去し実験シナリオを終了する。

## 4. 結果と考察

### 4.1 送信ノードの条件検討

表1に、今回実験に用いた送信ノードを示す。本実験では、PC1台で1つのプロセスを実行した場合の送信間隔の平均値と標準偏差を基準とし、多重度を検討した。多重度は、複数の仮想送信ノードが実行される1台のPCと、受信ノードとなる1台のPCを用いて検討した。

PCで単一のプログラムを実行した際の送信間隔の標準偏差は、図7に示すように、数百マイクロ秒の精度であった。そこで、本実験では送信間隔の標準偏差が、想定する送信間隔の1000分の1である、1マイクロ秒以下の場合を、正常な動作をしているとみなした。また、実験によって求めたメモリの量と最大実行プロセス数を表す図6を基に、表1のGroupAのPCで多重度の探索を行い、5000の仮想送信ノード多重が可能であることを確認した。

### 4.2 受信ノードの条件検討

1台のPCで実行した受信プログラムが、処理できる送信ノード数の上限を求めた。

当初、高性能な10GbE搭載のDual Opteron PCで数十万

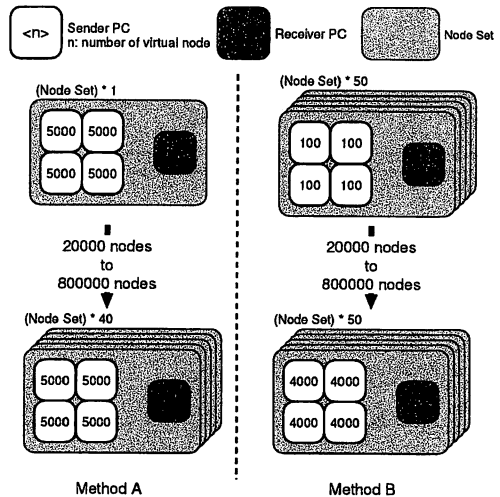


図8 実験ノードの増加方式

ノードからのパケットを処理しようと考えていた。しかし、saraを用いて受信したところ毎秒合計50000ノードからのパケットを受信した時点で、saraが0.5%のパケットを取りこぼしていることが分かった。また、調査したところ、saraの用いているパケットキャプチャライブラリである、libpcapで純粋にパケットのキャプチャを行うだけでも、0.5%のパケットの取りこぼしが発生することが分かった。当実験では、メガオーダーの送出されたパケットを十分に受信出来るノードを用意したいと考えたため、表1のGroupBを複数台用い、パケットを分散して受信する手法を用いた。

事前に、1台のPCで多重可能なノード数が5000で有ると分かっていたため、送信仮想ノードを多重化したPCを1台ずつ増やすことで、5000刻みで送信ノードを増加し、単一の受信プログラムが処理可能なノード数を探索した。パケットのドロップ数を判断基準とし、恒常的にパケットをドロップしない台数を探索した。結果、表1のGroupBでは20000台のノードまでのパケットを落とすことなく処理できる事を確認した。

### 4.3 送受信ノードセット

4.2で述べたように、1台の受信ノード(表1のGroupB)で、1メガオーダーのパケットをモニタは困難であったため、20000台の仮想ノードに対して1台の受信ノードを用いる。

1メガオーダーを実現するための仮想ノードの増加の手順として、図8の様な手法が考えられる。図8に、両手法で20000台の仮想ノードと、800000台の仮想ノードを実現する際の具体例を挙げる。

図8のMethodAの場合は、1送受信セットの、4台のPCそれぞれで5000ノードを仮想化し、仮想ノード数を変化させる場合は、このノードセットの数を増減する。実験を通して、パケットを送信するPCでは常に5000ノードが多重化されており、また、受信ノードは常に毎秒合計20000パケットを処理する。そのため、実験で利用する仮想ノードの数を変えた計画

であっても、送受信ノードの精度はそれぞれの計測と同様である。そのため、送受信のための PC の精度に関しては、仮想ノード数に依存せず同様の条件で計測が可能である。

図 8 の Method B の場合は、1 セット (4 台の送信のための PC と 1 台の受信のための PC) を、あらかじめ 1 メガ仮想ノードを実現出来る 50 セット分を利用し、それぞれの送信のための PC での多重度を変化させることで、実験に用いる仮想ノードの数を変化させる。この手法は、送信ノードが少数の場合でも、実験に用いる全ての PC を利用するため、全体の仮想ノード数を変化させ計測した場合でも、実験を通して、実験ネットワーク全体に対して均等に負荷を与えることが可能である。しかし、仮想ノードの数を変化するために、送信のための PC 上での多重度が増加すること、また、受信ノードで単位時間あたりに処理するパケット数が変化することから、計測値の微細な変化が、送信側、ネットワーク、或は受信側のいずれに起因するものかの特定が困難である。

当実験では、送受信の統計情報を取るために均質にパケットを増加したいと考えたため、MethodA を用いた。

#### 4.4 L2 ルーティングの問題

以下に挙げるような問題の対策をし、同一 L2 ネットワークでメガオーダーのパケットを送出し伝送できる系を構築した。

##### 4.4.1 フラディングパケットの送信ノードへの影響

IGMP Snooping が有効でなく、マルチキャストパケットが同一 Layer2 ネットワークにフラディングしてしまう場合、送出するパケットが少量であれば問題は発生しない。しかし、メガオーダーのノードからのパケットがフラディングすると、それらのパケットを受信することによりネットワークインタフェースに負荷が掛り、パケットの送出に影響が及ぼされることが確認された。

これより、多量のマルチキャストパケットを用いて通信する場合、受信ノードの接続されたポート以外にパケットを送出しない、IGMP Snooping の様な機構が必要となる。

##### 4.4.2 IGMP Querier の過負荷

今回、PC で mrouted を動作させ IGMP Querier とした。その際、毎秒メガオーダーのパケットが送出されると、それらのパケット全てがマルチキャストルータである IGMP Querier に対して全てのマルチキャストパケットが到着する。そのため、ネットワークカードで多量の割り込みが発生し、OS が過負荷になり IGMP Query の送出が正常に行えない状態であった。今回の PC マルチキャストルータの役割は、マルチキャストパケットを転送することではなく、IGMP Query を定期的に出すことであった。そのため、ネットワークインタフェースのドライバオプションで単位時間あたりの割り込み回数を制限することで、OS の処理する受信パケットを制限し、IGMP Querier が過負荷になることを防いだ。

#### 4.5 伝送品質の低下

1 メガノードからの毎秒合計 1 メガパケットの送出を行うことに成功した。

図 9 に受信間隔の平均値と標準偏差、図 10 に送信間隔の平均値と標準偏差を示す。これら 2 つの図は相似で有ることから、

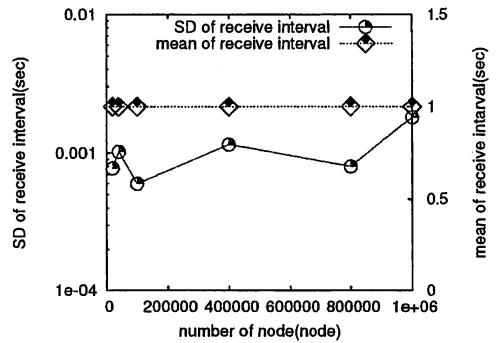


図 9 ノード数毎の受信間隔の平均値と標準偏差

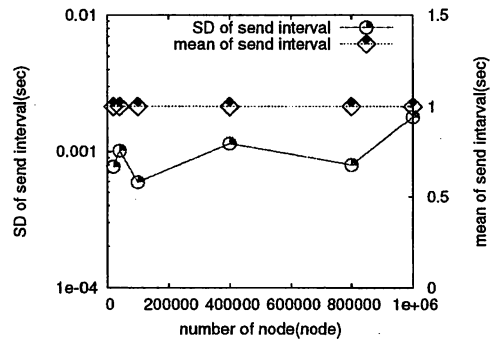


図 10 ノード数毎の送信間隔の平均値と標準偏差

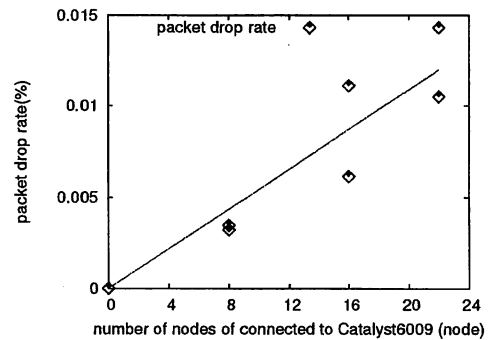


図 11 Catalyst6009 に接続された受信ノードの台数とパケットドロップ率の関係

ネットワーク上で遅延や揺らぎの影響を受けていないことが推測できる。

しかし、我々が [1] で述べたように、Catalyst6009 に接続された受信ノードを用いたことが原因であると思われる、パケットドロップが観測された。図 11 に、Catalyst6009 に接続されている受信ノード割合と、パケットロスに関するグラフを示す。このグラフは、640000 仮想ノードからパケットを送出し、32 台の PC でパケットを受信し、Catalyst6509 と Catalyst6009 に接続される受信ノードの台数を変えて測定を

行ったものである。X軸が、Catalyst6009に接続されたPCの台数で、Y軸がパケットドロップ数を基に求めたパケットドロップ率である。つまり、ネットワーク上に流れているマルチキャストパケットの数は同様であるが、受信ノードの接続されたスイッチの割合が異なる。その結果、Catalyst6009に接続されたPCの台数と、パケットドロップの割合はほぼ比例していることが分かる。これより、Catalyst6009内でのパケットドロップは、単位時間に通過するパケットの量に関わらず、一定の確率で起きていると推測される。

## 5. ま と め

送信ノードが $10^6$ 規模のマルチキャストパケットの送出・伝送を実現目的とし、これらを実現するための要件を探索するため、試行的にStarBEDを用いて、プロセス多重で1メガノードをエミュレートし、合計毎秒1メガパケットをマルチキャストで送出し受信することに成功した。

メガノードを実現するために、単一のPCでの多重度の検討を行った。また、単一の受信ノードでどれだけのパケットが受信可能であるかについても検討を行った。これらの結果から、送受信セットを作成し、セットを増減することで仮想ノード数を変更することで、一定の精度を保った送受信ノードを用意した。

複数の実PCで実現される多数のノードの計測を行う際に、全てのノードからのデータを一定以上の時間するために、プロセスの起動のばらつきを考慮し、計測のスケジューリングを行った。

数万程度のマルチキャストパケットが同一のlayer2ネットワークにフラディングすると、パケットの送出に影響が有ることが確認された。また、毎秒1メガのマルチキャストパケットがフラディングすれば、これらのパケットを浴びたPCはパケットの送信さえ行えない状態になることが確認された。よって、メガオーダーのマルチキャストパケットを用いた通信を行う際には、IGMP Snoopingの様なパケットをフラディングしないための機構が必要であることが示唆された。

送受信間隔の平均値と標準偏差やパケットのドロップ数の統計情報から、ネットワークの障害や問題、ボトルネックリンクが推測できることが示唆された。

## 謝 辞

本研究を行うにあたり、ご指導ご助言頂いた東京工業大学理工学研究科酒井善則博士、山岡克式博士、篠宮俊輔氏、情報通信研究機構大槻英樹博士、町沢朗彦氏、インテック北口善明氏、ならびに諸氏に深謝する。また、本研究は情報通信研究機構運営費交付金(情報通信部門)、JGN2-A17003、「IPマルチキャストを用いたSimple Node Administration Protocol実験」、NiCT北陸IT研究開発支援センター、平成17年度厚生労働省がん研究助成金研究総合研究「がん情報ネットワークを利用した総合的がん対策支援の具体的方法に関する研究」若尾班等の支援、Cisco Systems Japanならびに三井物産株式会社の技術協力を得て行った。関係各位に深謝する。パケット受信プログ

ラムsaraを作成するにあたり、参考にさせて頂いたENMA [8]の制作者である中村豊氏に感謝します。

## 文 献

- [1] 三角真, 中川晋一, 我如古津世史, 知念賢一, 篠田陽一: "メガノードIPマルチキャストセンサネットに関する実験的検討", 信学技報, CQ2006016(2006-4), pp.75-81, ISSN 0913-5685, 2006
- [2] Toshiyuki Miyachi, Ken-ichi Chinen and Yoichi Shinoda: "Automatic Configuration and Driving Internet Experiments On An Actual Node-base Testbed", Tridentcom2005, pp.274-282, Trento, ISBN 0-7695-2219-X, 2005
- [3] 宮地 利幸, 知念 賢一, 篠田 陽一: 「SpringOS/VM: 大規模ネットワークテストベッドにおける仮想機械運用技術」, 情報処理学会, ISSN 0919-6072, 2005
- [4] 三角真, 中川晋一, 我如古津世史, 知念賢一, 出口真人, 篠田陽一: 「メガオーダーノードセンサネットエミュレーションに関する実験的検討」, The Seventh Workshop on Internet Technology - WIT2005, 2005.11.
- [5] Takemoto, S., T. Yamamoto, A. Mukai, S. Otsuka and K. Fujimori: Crustal Strain Observation for Nine Years with a Laser Strainmeter in Kobe, Japan, Journal of Geodynamics, Vol. 35/4-5, (2003) 483 - 498.
- [6] 出口真人, 中川晋一, 三角真, 篠田陽一: マルチキャストを用いたモバイルコンピュータデバイス拡張の提案, IPSJ-SIG Technical Report 2004-MBL-31(21), pp157-162, 2004
- [7] 根日屋英之, 植竹古都美: ユビキタス無線工学と微細 RFID, 東京電機大学出版局, 2004
- [8] 中村 豊, 知念 賢一, 砂原 秀樹, 山口 英: ENMA: パケットモニタによるWWWサーバの性能計測システムの設計と実装, 電子情報通信学会 和文論文誌 D-1, p329-338, 2000.3.