

Title	ネットワークシミュレータSSFNetの分散プログラミングフレームワークNekoへの統合
Author(s)	松下, 誠和
Citation	
Issue Date	2005-03
Type	Thesis or Dissertation
Text version	author
URL	<a href="http://hdl.handle.net/10119/1919">http://hdl.handle.net/10119/1919</a>
Rights	
Description	Supervisor:Defago Xavier, 情報科学研究科, 修士

修 士 論 文

ネットワークシミュレータ SSFNet の  
分散プログラミングフレームワーク Neko への統合

北陸先端科学技術大学院大学  
情報科学研究科情報システム学専攻

松下 誠和

2005 年 3 月

修士論文

# ネットワークシミュレータ SSFNet の 分散プログラミングフレームワーク Neko への統合

指導教官 Defago Xavier 特任助教授

審査委員主査 Defago Xavier 特任助教授

審査委員 片山卓也 教授

審査委員 鈴木正人 助教授

北陸先端科学技術大学院大学  
情報科学研究科情報システム学専攻

310105 松下 誠和

提出年月: 2005年2月

## 概要

現在、ネットワーク化された計算機は、毎日の生活の最も必要な関連として提供されている。そして、社会の必要不可欠な部分は、ネットワーク化された計算機を構成する分散システムに頼って形成されている。例えば、トランザクション処理を有するような金融機関のATMシステムや航空機会社の時間の制約が厳しい発券機システムなどである。従って、各分散システムの安全性と耐故障性を構成することは、最つとも高い優先度とされる。

このような分散システムを構築するためには、Replication 技術やその技術を達成するために用いられる Atomic Broadcast のような分散アルゴリズムの開発、性能評価、および、研究が必要である。分散アルゴリズムの開発および性能評価を行うことは、複雑かつ多くの時間を要する作業である。なぜなら、複雑なシステムを想定してアルゴリズムを開発および性能評価を行わなければならず、容易に評価環境を構築できないからである。性能評価を行うためには、シミュレーション環境でアルゴリズムを実行、および、実システム環境でアルゴリズムを計測などが挙げられるがいずれの方法も、評価したいアルゴリズムをプログラミング言語を用いて、研究者自身で実装しなければならない。この複雑な作業による分散アルゴリズムの研究者への負担は、分散アルゴリズムを開発、および性能評価するすべての行程をサポートする単一の環境が存在しないことが主な要因である。

本研究の対象である Neko は、分散アルゴリズムを開発および性能評価を行うための分散プログラミングフレームワークである。Neko 単体で分散アルゴリズムの開発、および、性能評価に必要なすべての行程を行うことが可能である。つまり、Neko では、Neko のフレームワークを用いて実装された分散アルゴリズムを Neko 上に実装されているシミュレーション環境および実システムの環境の両方で性能評価可能であり、開発者が異なる複数の環境を用いる必要がない。また、Neko のフレームワークは、分散アルゴリズムを実装するために必要な機能（単純なメッセージインタフェース、異なるアルゴリズムを階層的に構築する構造 etc）を提供している。開発者は、このような Neko のフレームワークを用いることで、簡単に分散アルゴリズムを実装することが可能である。しかし、Neko のシミュレーション環境を用いて性能評価を行う場合、重大な問題点がある。問題点とは、Neko 上に実装されたシミュレーション環境は、実世界で用いられるような計算機ネットワークをモデル化して実装されていないので、現実的なシミュレーションネットワーク環境での性能評価を行えないことである。

この問題を解決するために本研究では、現実的な計算機ネットワークをモデル化し実装しているネットワークシミュレータである SSFNet を分散プログラミングフレームワークである Neko への統合を行う。具体的な統合方法は、Neko と SSFNet の間に Neko が SSFNet のシミュレーション環境を利用するためのインタフェースを構築する。統合する上で実装上の問題がいくつかあり、Neko と SSFNet のスケジューラの同期、アドレス方式の違い、メッセージ形式の違いなどである。これらの問題を解決することによって、Neko のネットワークシミュレーション環境として SSFNet を用いるが可能となった。このような方法

で統合し、現実的なネットワークシミュレーション環境上で性能評価を行えることが可能となった Neko 上で、分散アルゴリズムの 1 つである Atomic Broadcast を用いて性能評価を行った。Atomic Broadcast とは、分散システム中の複数のプロセスがネットワーク上からの要求（メッセージ）を全順序で目的となる計算機に配達するアルゴリズムである。Atomic Broadcast は、Consensus とよばれる分散アルゴリズムを用いて実現する種類が存在する。本研究では、2 種類の良く知られた Consensus を Atomic Broadcast の性能評価基準を用いて比較した。1 つ目のアルゴリズムは、集中的な通信よりメッセージの配達する順番を決定（合意）を用いており ([8])、2 つ目のアルゴリズムは、集中的な通信を用いていない ([7]) 2 種類の Consensus の上に構成される Atomic Broadcast を Neko 上で実装されている Ethernet を単純にモデル化したシミュレーションネットワークと、本研究で統合した Neko を用いて性能評価し、比較を行った。また、Neko と SSFNet のシミュレーションネットワークのモデルを比較することにより、統合した Neko の有効性を示した。本研究により、分散アルゴリズムの研究者へ適切な分散アルゴリズムの開発と性能評価のための有用な結果を得られるような高性能な分散プログラミングフレームワーク Neko の提供が可能となった。

# 目次

<b>第1章</b>	<b>はじめに</b>	<b>1</b>
<b>第2章</b>	<b>Neko と SSFNet について</b>	<b>3</b>
2.1	Neko とは . . . . .	3
2.1.1	Neko の問題点 . . . . .	7
2.2	SSFNet とは . . . . .	7
2.2.1	SSFNet のアーキテクチャ . . . . .	10
<b>第3章</b>	<b>統合方法</b>	<b>13</b>
3.1	統合デザインと実装上の問題の解決 . . . . .	14
3.1.1	SSFNet Networks . . . . .	17
3.2	動作テスト . . . . .	18
<b>第4章</b>	<b>How to use Neko with SSFNet</b>	<b>20</b>
4.1	インストール方法 . . . . .	20
4.2	コンフィグファイルの設定 . . . . .	21
4.3	実行 . . . . .	26
<b>第5章</b>	<b>性能評価</b>	<b>29</b>
5.1	Atomic Broadcast とは . . . . .	29
5.1.1	Atomic Broadcast の定義 . . . . .	29
5.1.2	Atomic Broadcast アルゴリズム . . . . .	31
5.2	Evaluation Neko with SSFNet by Two Consensus Algorithm . . . . .	35
5.2.1	定義 . . . . .	36
5.2.2	2つの Consensus Algorithms . . . . .	36
5.2.3	ベンチマーク . . . . .	38
5.2.4	シミュレーションモデルの比較 . . . . .	39
5.3	評価結果 . . . . .	41
<b>第6章</b>	<b>まとめ</b>	<b>44</b>
6.1	謝辞 . . . . .	44

# 第1章 はじめに

分散アルゴリズムとは、分散システムにおける耐故障性を保証するために用いられ、複数の計算機を利用して問題を解決する方法である。例えば、データベース等のサーバへデータの変更を行うような多くのシステムには、耐故障性を高めるために Replication 技術が用いられている。Replication とは、クライアントからの要求を複数のプロセスに送信することにより、単一のプロセスでシステムを保持するよりも耐故障性を高める技術のひとつである。単一のプロセスでクライアントからの要求を処理する場合に、そのプロセスが停止すれば、そのクライアントからの要求を処理できないだけでなく、システム全体が停止してしまいサービスを提供できなくなる。Replication は、複数のプロセスの中から、要求に対する処理や返答を主に行うリーダープロセスを決めておく。もし、そのリーダープロセスが何らかの故障でクラッシュしてしまった場合、他のプロセスが代わりにリーダープロセスとなり、システムを続行させる。しかし、システムを続行させるには、リーダープロセスと同じ要求を同じ順番で受信しておき、かつ、同じ順番で処理しなければならない。このような技術を実現するためには、クライアントの要求をすべてのプロセスが同じ順番（全順序）で処理する必要がある。

全順序で要求を処理するために Atomic Broadcast という分散アルゴリズムを用いて、すべての要求を全順序になるようにプロセスに送信する。Atomic Broadcast は、2つの性質を要求される。1つ目の性質は、送信されたの要求がすべてのサーバに届くか、あるいはどのサーバに届かないかのどちらかを保証することである。2つ目性質は、すべての要求が全順序ですべてのプロセスに届けられることである。この分散アルゴリズムを実現するためには、非常に複雑な処理をシステムが行わなければならない。本研究では、この Atomic Broadcast を用いて、第5章で性能評価を行っている。その章で、Atomic Broadcast の詳細について解説する。

分散アルゴリズムの研究者が、このような Replication 技術や Atomic Broadcast などの分散アルゴリズムを開発する場合、他の同じ種類の分散アルゴリズムとの有効性を示すためには性能評価が必要がある。しかし、性能評価を行うためには、実システム上にプログラミング言語などを用いて、分散アルゴリズムを実装し評価しなければならない多くの時間を要してしまう。さらに、多くの分散アルゴリズムが開発されているが、それらの性能評価を上記のような理由で、多く行われていない。研究者の負担を軽減するために、我々は分散アルゴリズムを容易に性能評価するためのツールを提供する必要がある。

本研究の目的は、NEKO のネットワークシミュレーションの問題点を SSFNet と統合することで解決し、拡張された NEKO の有効性を示すことである。NEKO と SSFNet 間

の問題であるスケジューラの同期、ネットワークアドレスやメッセージ形式の違いを解決し、NEKO と SSFNet のインタフェースを作成し統合する。そして、分散アルゴリズムである Atomic Broadcast を拡張された NEKO に実装し性能評価を行うことによって有効性を示す。

ここで、この論文の構成を解説しておく。第 2 章で、本研究の対象である Neko と SSFNet について解説し、第 3 章では Neko と SSFNet の統合について解説する。第 4 章では、第 3 章で統合したツールの利用の仕方を紹介し、第 5 章でツールと Atomic Broadcast の性能評価を行う。



## 第2章 Neko と SSFNet について

本章では、Neko と SSFNet について紹介する。Neko と SSFNet の構成、基本的な機能および、特徴について詳細に説明する。Neko については、その特徴的なレイヤー構造、メッセージを通信する方法や Neko の構成から得られる主要な部分について解説する。SSFNet については、SSFNet のコンセプト、SSFNet で利用可能プロトコルやサンプル DML ファイルを用いての簡単なパラメタの解説する。

### 2.1 Neko とは

NEKO[1] とは、Java 言語で実装された分散アルゴリズムを性能評価および開発するためのフレームワークである。離散的なイベント駆動のシミュレータを持っており、NEKO の API を用いて分散アルゴリズムを実装することにより論理時間 [9] や実時間での経過時間、レイテンシやスループットを計測することが可能である。NEKO の API を用いて実装された同一のソースコードで複数の計算機を用いた実ネットワークでの実験と単一の計算機でのシミュレーションが可能である。つまり、分散アルゴリズムの性能評価に必要な行程を Neko だけを用いることのみでできる。シンプルなメッセージパッシング方式により、分散アルゴリズムを実装することが可能であり、分散アルゴリズムを実装するために必要としている機能を API として用意している。これらの研究者は容易に分散アルゴリズムの実装と性能評価が行うことが可能である。

NEKO のアーキテクチャ (図 2.1) は、アプリケーションレベルとネットワークレベルの二つの主要な部分に分けることができる。アプリケーションレベルでは、NEKO の API を用いて分散アルゴリズムの振る舞いを実装する。複数のプロセスがシンプルなメッセージパッシング用のインタフェースを用いて、メッセージを非同期な伝播が可能である。メッセージを送信するプロセスは、send という基本的な操作を用いてメッセージをネットワーク上に送り出す。そして、ネットワークは、そのメッセージを受信するプロセスへ、deliver という操作を用いて届ける。このようなプロセスは、複数のレイヤーによってプログラミングされる。

Neko では、メッセージを伝搬するための基盤であるネットワークは、いくつかの方法を用いて制御することが可能である。1つ目は、ネットワークは、あらかじめ定義されたネットワークである実ネットワークかシミュレートされたネットワークを用いてインスタンス化できるということ。2つ目は、Neko は複数のネットワークを並行に適切に利用できるということ。最後に、Neko に実装されたネットワークは、簡単にプログラミング

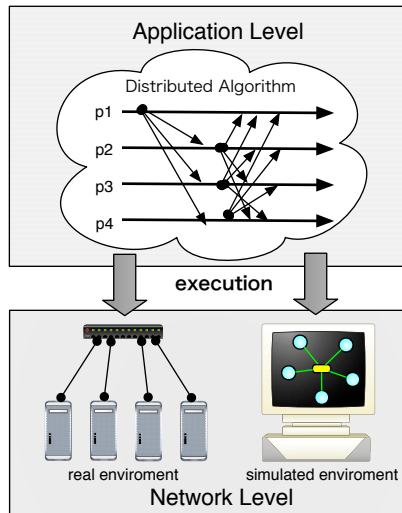


図 2.1: Architecture of Neko

でき、かつ新しいユーザー自身が定義したネットワークを追加可能である。次に、Neko のアプリケーションレベルとネットワークレベルの詳細について紹介する。

**アプリケーションレベル:** Neko のアプリケーションは、階層的なレイヤーとして構築される。送信されるメッセージは、send メソッドを用いて下位のレイヤーに渡される。そして、配達されるメッセージは、deliver メソッドを用いて上位のレイヤーに渡される。レイヤーは、*Active* と *Passive* のどちらか一方の種類を用いる。Passive レイヤー (図 2.2) は、メッセージをコントロールするためのスレッドを持ち合わせていない。つまり、下位のレイヤーから、メッセージが Passive レイヤーに渡されるということは、下位のレイヤーは、その Passive レイヤーの deliver メソッドを呼んでいることになる。一方、Active レイヤー (図 2.3) は、メッセージをコントロールするためのスレッドを持ち合わせている。このレイヤーは、メッセージを下位のレイヤーから、receive メソッドを用いて、能動的に受け取ることができる。Active レイヤーは、FIFO 形式のメッセージキューを持っており、下位のレイヤーからメッセージが Active レイヤーに渡されるということは、その Active レイヤーのメッセージキューへ追加していることになる。Active レイヤーが receive メソッドを呼ぶと、メッセージがメッセージキューに到着するまでブロックする。さらに、receive メソッドは、タイムアウトによりブロックする時間を制御可能である。receive メソッドへの引数として、0 を与えることは、ブロッキングされないことである。つまり、メッセージの到着を待たないことである。Active レイヤーは、deliver メソッドを用いることで Passive レイヤーのように振る舞うことができる。これは、Active レイヤーのメッセージキューを迂回する方法として用いることと同じである。

開発者は、階層的な構造のようにアプリケーションを構築することを義務づけられていないわけではない。つまり、レイヤーは違う方法によって結合させることができる。例え

ば、いくつかのレイヤーの1つのNekoメッセージのメッセージタイプに基づいて作成されたチャンネルから、メッセージがあるレイヤーに伝播されたとする。そのメッセージタイプによって、レイヤーは、開発者が定義したメソッドにより相互にメッセージを受け渡すことができる。それらのメソッドは、sendメソッドやdeliver/receiveメソッドによる制限を受けない。つまり、開発者は、すべてのタイプのJavaオブジェクトを定義し、Neko上で利用することが可能である。

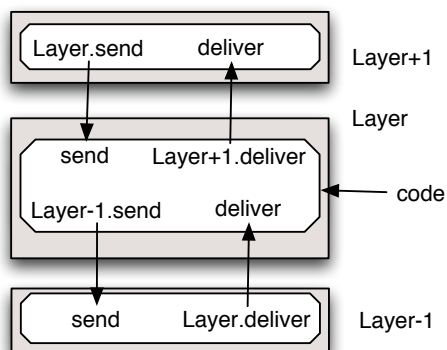


図 2.2: Passive Layer

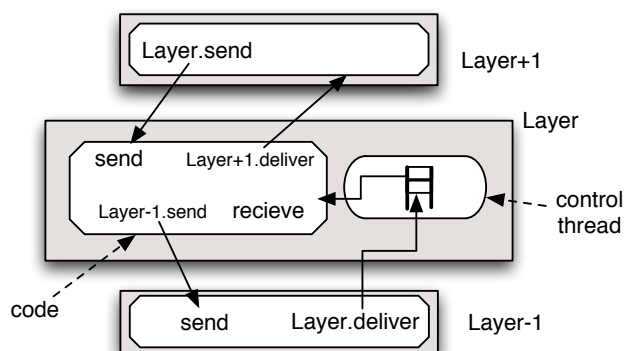


図 2.3: Active Layer

**Neko プロセス**：分散アプリケーションの各プロセスは、関連づけられた Neko プロセスを持っている。その Neko プロセスは、ネットワークとアプリケーションのレイヤーに設置されており、次に3つの重要な振る舞いに関するルールがある。

1. プロセスは、各プロセスのアドレスのような広い情報を持っている。そのプロセスのすべてのレイヤーは、すべてのプロセスにアクセスすることができる。それは、一般的な Single Multiple Data (SPMD) プログラミングに似ており、あるプログラム

は、いくつかのプロセスにより動作しており、Neko プロセスから得られるアドレスによって分岐している。

2. Neko プロセスは、プロセスによる送信や受信の履歴を取得する。
3. もし、アプリケーションが並行にいくつかのネットワークを使うなら、Neko プロセスは、適切なネットワークによるメッセージを処理する。(例えば、2つの違う物理的なネットワーク上で通信を行う場合や、違う二つのプロトコルで同じ物理的なネットワークを利用する場合)

**Neko メッセージ:** すべての基本的な通信 (send、deliver、receive) は、Neko メッセージのインスタンスによって転送される。あるメッセージは、ユニキャストまたは、マルチキャストのメッセージである。各メッセージは、Java のオブジェクトから成る content と いわれる部分と下記に示すようなヘッダーによって構成されている。

**Addressing(source (送信元アドレス) , destinations (目的地アドレス) ):** アドレスの情報は、メッセージを送信したプロセスのアドレスと目的のプロセスのアドレスから成り立っている。アドレスは、小さな整数であり、0 から開始され番号づけされる。

**Network :** Neko メッセージがいくつかのネットワークで用いられるとき、各メッセージは転送に用いられるべきネットワークの ID に運ばれる。これは、メッセージが送信される時に、はっきりと区別される。

**Message type** 各メッセージは、開発者が定義したタイプの整数によるフィールドを持っている。それは、違うプロトコルに属するメッセージとしてを見分けるために使うことができる。

## ネットワークレベル

Neko のネットワークは、低い下位のレイヤー (図 2.1) を構成する。開発者は、アプリケーションのコードを変更することなく、コンフィグファイルによってネットワークを区別することができる。

**Real Network** 実ネットワークは、Java ソケットから構築される (あるいは、他のネットワークライブラリー)。そのネットワークは、Java による連続化された Neko メッセージの content を用いる。Neko 上に実装されている TCPNetwork は、高信頼なメッセージ伝播方式である TCP/IP により構築される。TCP のコネクションは、各 1 組のプロセスの間に起動時に確立される。UDPNetwork は、信頼性の低いメッセージ伝播方式である UDP/IP により構築される。これらの実ネットワークのパフォーマンスについては、Urban らによって、[1] で性能評価を行っている。

## Simulated Network

現在、Neko は単純なバージョンの Ethernet、および、FDDI のネットワークをシミュレート可能である。複雑な現象における Ethernet 上のコリジョンなどは、モデリングされていない。しかし、異なったモデリングによるネットワークを単純なネットワークイン

タフェースによって差し替え可能である。違う種類のシミュレーションネットワークは、分散アルゴリズムをデバック時に有用なものを提供できる。Neko に実装されているネットワークは、ランダムな時間を計算するか、または、定数による時間の経過の後にメッセージを伝播する。これは、モデリングした実際のパフォーマンスは、現実的な結果ではない。

### 2.1.1 Neko の問題点

NEKO のネットワークレベルには、複雑なネットワークモデルのシミュレーションができないという大きな問題点がある。NEKO のシミュレーションネットワークは、計算機同士が直接接続された完全結合のトポロジしか想定されていない。(図 2.4) さらに、ネットワーク間の帯域幅や遅延や NIC のレイテンシーなどを柔軟に設定できない。メッセージの伝播時間は、定数によってあらかじめ決められるか、または、ランダムに生成された遅延時間を用いている。ネットワークのトラフィックおよびバッファリングによる遅延などを実ネットワークを想定したシミュレーションを行っていない。よって、NEKO では現実的なネットワークシミュレート環境上で分散アルゴリズムを性能評価できない。

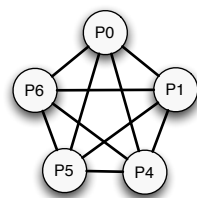


図 2.4: Neko Topology

しかし、Neko は、このような問題点は把握しており、開発者が簡単に他のネットワークを利用できるように設計されている。本研究では、この Neko の問題点である現実的なシミュレーションネットワーク環境を構築ができないことを SSFNet というネットワークシミュレータが Neko のネットワークとして振る舞うことで解決する。

## 2.2 SSFNet とは

SSFNet[2] (図 2.5) とは、Java 言語で実装されたネットワークシミュレータで、離散的なイベント駆動のシミュレータで動作している。SSFNet は、ネットワークとインターネットプロトコルをシミュレーションするために SSF (Scalable Simulation Framework) ベースのモデルである。SSFNet は、Java 言語以外に C ++ 言語でも実装されており、SSFNet

のパッケージは、大規模なネットワークモデルを構築するために、直接または、継承ができるクラスとして提供されている。その SSFNet のライブラリは、ネットワークの要素(ホスト、ルーター、ネットワーク、インタフェースカード、LAN) のためのコンポーネントモジュールとネットワークプロトコル (IP、UDP、TCP、BGP、OSPF、ICMP、HTTP など) 含んでいる。SSFNet はいくつか 大規模なモデルをサポートすることについてのいくつかの難しい方法を組み入れており、それらによって、非常に高いパフォーマンスを手に入れている。1つ目は、モデルの環境設定ための高度なサポートは、適切な結果を得るために重要なことである。2つ目は、大規模なネットワークモデルは、簡単モデルなどを高度に組み立てることにより、複数の要素による構成を用いた指向で最適に管理、構築されなければ成らない。3つ目は、高いシミュレーション性能を成し遂げることは、シミュレートされたコンポーネント間の相互作用を理解することと、SSF のパフォーマンスモデルに基づく集成的な計算結果と通信の要求を解釈することに深く依存する。

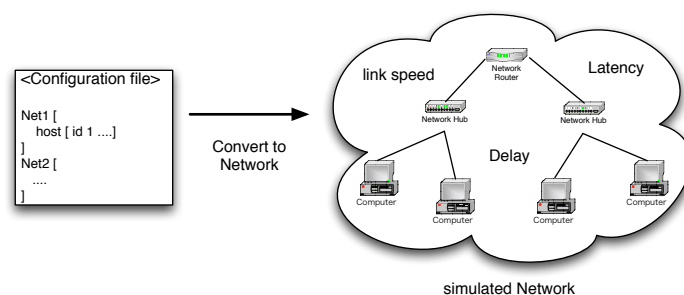


図 2.5: SSFNet

**self-configurable model** SSFNet のモデルは、開発者自身が自由に構成することが可能である。つまり、各 SSFNet クラスのインスタンスがクエリ形式のパラメータデータベースを用いることにより、それ自身を自動的に構成可能であるということである。コンフィグレーションのデータは、親のエンティティにより、自分自身の構成情報の一部を子のエンティティに渡すことで階層的に構築される。そのコンフィグレーションデータベースは、階層的なスキーマを持っている。スキーマの情報は、利用に耐えうる変数とデータベースから情報を検索するような機能を蓄えた入れ子状の基盤である。

シミュレートされた 100 万のコンポーネントにおよぶ規模は、必要なコンフィグレーション管理ためにオブジェクトデータベース技術を用いている。計画方針を除いて、それは、大規模なモデルへのパラメータへ一貫性と効率を満たすことは極端に難しいものである。構築されたデータベースは、典型的なコンフィグレーションデータと主要なモデリングエラーの原因を除くことにも対応する。

### Model composition from components

SSF のモデルは、各大規模なモデルに組み立てられたコンポーネントのための木構造の

コンフィグレーションパラメータを区別するために、単純な階層的な属性の木構造による DML(Domain Modeling Language) 表記法を用いている。

開発者は、基本的な定義のデータベースを簡単に参考にすることで、任意のネットワークコンフィグレーションをすばやく構成可能である。各ネットワークコンフィグレーションデータベースは、各コンポーネントとして実装された実行可能な C++ また、Java を用いることにより価値ある成果になる。

### Self-configurable models

DML の仕様は、1 組の key/value から成り立っている。value は、属性のリストから成っている。最上位レベルのノードは、Net で SSFNet ライブラリの Net クラスと一致する。Net は、自分自身で構成する。つまり、Net は特有のインタフェースを実装しており、クエリを用いてそれ自身のパラメータを獲得することデータベースとして機能する。

Net は、データベースを読み込み、そして、Router と Host のインスタンスを生成する属性のリストを記憶する。また、Router と Host も基本的に自分自身で SSFNet のクラスとして構成する。そして、NIC(ネットワークインタフェースカード) がそれらに適切に接続される。いったん構築されると、各インスタンスは、それらの特有のコンフィグレーションデータベースの一部の参照を受け取る。そして、また、順番に自身で構成することができる。この作業は、全体のモデルが周期的に構築、構成されるまで、子のエンティティーがさらに下の子のエンティティーを構成するのと同時に続けられる。

### Network Protocol Graphs and Network Interface Models

Router と Host は、x-kernel のデザインパターンを用いて、特有のネットワークプロトコルのグラフによりパラメータ化される。各ホストはプロトコルグラフを含んでいる。プロトコルグラフは、ここのプロトコルセッションによって組み立てられ、アプリケーションのアクティビティー、または、その Host の NIC にパケットが到着し応答があったとき、上位のプロトコルにメッセージを渡したり、下位のプロトコルに渡したりする。図 2.7 は、図 2.6 におけるプロトコルスタックの論理的な構造を示したものである。

SSFNet は、新しい TCP の実装を含んでいる。それは、RFC の要求に基づいたすべての TCP の状態、フロー制御とネットワークの混雑をコントロールする技術である。また、それらは、slow start, Jacobson's RTT estimation, Karn/Partridge アルゴリズム, exponential retransmit timer backoff, adaptive acknowledgment, および fast retransmit を含んでいる。それらは、簡単に継承し、デザインされており変更している。例えば、send-window, receive-window, および incoming message demultiplexing は、すべて明瞭なクラスとして提供されている。

SSFNet では、実世界のインターネットと同じ振る舞いとして、TCP は IP の上に働き、さらに、IP は、仮のプロトコルが表した構成された NIC の上で働く。各 NIC は、その世界による IP パケットのイベントを交換するための蓄積された SSF チャンネル (Buffer) を維持している。NIC は、同じリンクタイプの違う NIC と接続される。そして、特徴的な物理的なリンク、パケットフォロワーのバッファリング、および flow interleaving と scheduling

のためのオプションを自信で構成することをサポートしている。

### 2.2.1 SSFNet のアーキテクチャ

SSFNet のアーキテクチャ(図 2.8) は、プロトコルレベルとネットワークレベルに分けることができる。プロトコルレベルでは、実際のプロトコルスタックを構成するように、プロトコルを実装することが可能である。プロトコルスタックの最上位層の部分は、実際に JAVA で実装するように、計算機同士でどのようなやり取りをするか(プロトコル)を開発者自身が実装する。最上位層以外の層は、DML を用いて SSFNet で用意されたクラスを用いて構成する。開発者は、自由にプロトコルを定義することが可能であり、SSFNet の API を用いることでプロトコルを実装することが可能である。

ネットワークレベルでは、DML を用いて現実的なネットワークモデルを定義する。DML では、ネットワークのトポロジ、リンク間の遅延、ネットワークカードの Latency および各計算機のバッファなどを詳細に定義することが可能である。DML 上で、各計算機へ振る舞いを実装したクラスを割り当てることで、各計算機の振る舞いを定義することができる。



```

Net [
  frequency 1000000000
  router [
    id 255
    graph [ProtocolSession [name ip use SSF.OS.IP]]
    interface [ id 0 buffer 16000 bitrate 100000000 latency 0.0001 ]
  ]
  host [ id 1
    nhi_route [dest default interface 0 next_hop 255(0)]
    interface [id 0 bitrate 100000000 latency 0.0001]
    graph [
      ProtocolSession [name Server use SSF.OS.TCP.Server port 10]
      ProtocolSession [name socket use SSF.OS.Socket.socketMaster]
      ProtocolSession [name tcp use SSF.OS.TCP.tcpSessionMaster]
      ProtocolSession [name udp use SSF.OS.UDP.udpSessionMaster]
      ProtocolSession [ name ip use SSF.OS.IP]
    ]
  ]
  host [ id 2
    nhi_route [dest default interface 0 next_hop 255(0)]
    interface [id 0 bitrate 100000000 latency 0.0001]
    graph [
      ProtocolSession [name client use SSF.OS.TCP.Client message-size 8]
      ProtocolSession [name socket use SSF.OS.Socket.socketMaster]
      ProtocolSession [name tcp use SSF.OS.TCP.tcpSessionMaster]
      ProtocolSession [name udp use SSF.OS.UDP.udpSessionMaster]
      ProtocolSession [ name ip use SSF.OS.IP]
    ]
  ]
  link [attach 255(0)
    attach 1(0)
    attach 2(0)
    delay 0.001
  ]
]

```

Figure 2.6: Sample SSFNet DML configuration file

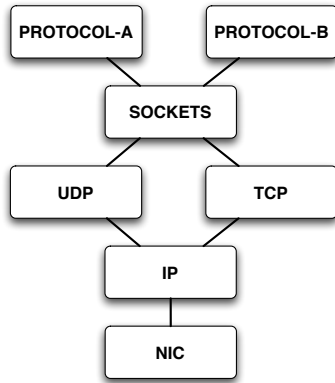


图 2.7: SSFNet protocol graph

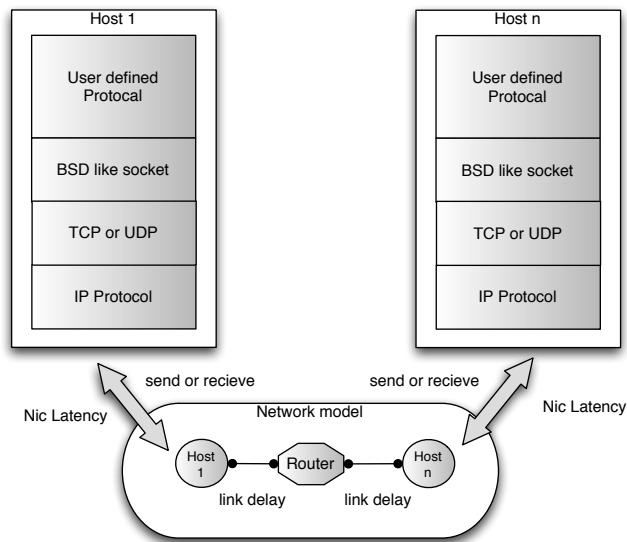


图 2.8: Architecture of SSFNet

## 第3章 統合方法

Neko では、Neko の柔軟なインタフェースを利用して、開発者自身がモデル化したシミュレートネットワークのクラスを Neko へ実装し利用することが可能である。SSFNet では、開発者が DML ファイルを用いてプロトコルスタックを自由に構成することができ、プロトコルスタックの最上位層に、開発者自身がプロトコルをモデル化し実装したクラスを利用可能である。本研究では、Neko と SSFNet を統合することが目的である。つまり、Neko のシミュレートネットワークとして、SSFNet を用いることで、分散アルゴリズムを現実的なネットワークでシミュレートできるようにすることである。具体的な方法は、Neko と SSFNet の間にインタフェースを作成することである。Neko 側のインタフェースとして、SSFNetNetwork というシミュレートネットワークのクラスを作成する。SSFNet 側のインタフェースとして、NekoProxy という SSFNet のプロトコルスタックの最上位層のクラスを作成する。SSFNetNetwork が NekoProxy を介して、SSFNet のシミュレーションネットワークを用いてメッセージを通信するために、SSFNet の TCP および UDP ソケットインタフェースを利用できるようにする。これは、あたかも、Neko の実ネットワークにおける振る舞いと非常に似ている。実ネットワークを Neko が利用する場合、Java のソケットクラスを用いて、異なる計算機上のプロセスと通信を行い分散アルゴリズムを実行している。つまり、SSFNet のネットワークを実ネットワークとして利用していることである。実際、NekoProxy のクラスの中には、Neko が実ネットワークを用いる場合に使われているクラスである、TCPNetwork と UDPNetwork を継承しているクラスがある。

しかし、統合する上で下記のようないくつかの Neko と SSFNet 間の特有の問題がある。

- 2つのイベント駆動のスケジューラの同期
- アドレスのマッピング方式
- メッセージ形式の違い
- コンフィグファイルの変更

この章では、次のような3つの統合する上で考慮した問題について解説していく。

- Neko と SSFNet の全体的な統合のデザインと、デザイン中のクラスとその役割によって分けられる個別の部分。
- 上記に示した Neko と SSFNet 間の特有の問題の詳細と解決方法。

- Neko 側のインタフェースとして実装した SSFNetNetwork には、TCPNetwork と UDPNetwork の 2 つの異なった種類のシミュレートネットワークが存在する。2 つの種類の SSFNetNetwork の利用方法と詳細。

### 3.1 統合デザインと実装上の問題の解決

SSFNet の全体的な統合のデザインと、実際に実装したクラス名を用いて、そのクラスの役割を説明し、Neko と SSFNet の統合方法について解説する。Neko と SSFNet を統合する上で下記のようないくつかの実装上の制限がある。

- Neko が SSFNet をコントロールする、つまり、Neko のシミュレーション時間で SSFNet を駆動させる。
- Neko と SSFNet の両方とも、システムの根幹を成すスケジューラの様な中心となる実装をいっさい変更しない。つまりインタフェースとなる一部のパッケージにより統合する。これは、SSFNet のバージョンが上がるような大きな変更にも耐えるためである。
- Neko のスケジューラと SSFNet のスケジューラを同期的に並行に動作させる。

基本的な設計指針は、Neko と SSFNet の両方にインタフェースとなるクラスを作成し、両方のシミュレータのスケジューラを並行に動作させながら、シミュレーションを行うことである。全体的な設計デザインは、図 3.1 のように表すことができる。メッセージがプロセス  $P_i$  から  $P_j$  へ伝播される例をもとに各クラスの役割について解説する。最上位で Neko プロセスが動作し、開発者が実装した分散アルゴリズムを実行している。Neko プロセス  $P_i$  は、Neko 上で Network Layer でシミュレートネットワークとして認識されている SSFNet Networks を介して、SSFNet 上の Host である NekoProxy にメッセージを伝播している。SSFNet Networks を用いて、SSFNet の Host である NekoProxy にメッセージを仲介して貰う時、Neko の ID と SSFNet の IP アドレスを 1 対 1 に対応させた Mapping Table を用いる。(アドレス方式の違い) NekoProxy は、あらかじめ作成された TCP か UDP のソケットを用いて通信する。しかし、この時、Neko 上と SSFNet 上の両方のスケジューラのシミュレーション時間の一貫性が取れていない。つまり、スケジューラの同期がここで問題になる。この問題の解決方法については、**スケジューラの同期**で詳しく述べる。そして、SSFNet のスケジューラでメッセージの伝播をシミュレートし、SSFNet によってメッセージが他の NekoProxy に送信される。メッセージが受信されると、Neko が NekoProxy からメッセージを取得し、SSFNet Networks を用いて Neko プロセス  $j$  へメッセージを伝播しようと試みる。しかし、この時、分散アルゴリズムのメッセージとして実装されているクラスである Neko メッセージは、複数の宛先 ID を保持している。送信時には、この宛先 ID を SSFNet の IP アドレスと対応させて送信することは可能である。しかし、SSFNet Networks を用いてメッセージを宛先の Neko プロセスへ伝搬するとき、そのメッセージは

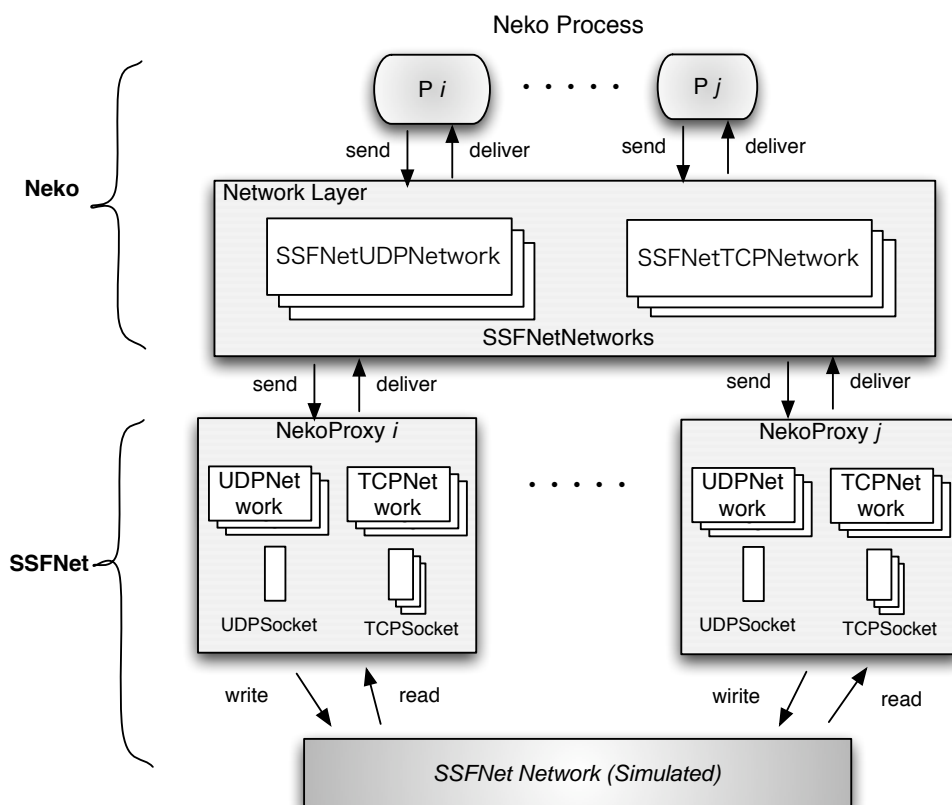


図 3.1: Integration Design

一意指定された Neko プロセスへのメッセージである。つまり、複数の宛先を含んでいると、そのメッセージを一意指定された Neko プロセスへ伝播することができない。この問題の解決方法は、**ネットワークアドレスの違いとメッセージ形式の違い**で詳しく述べる。最後に、Neko プロセス  $P_j$  は、自分自身に伝搬されたメッセージを受信し、そのメッセージにより分散アルゴリズムを実行していく。

このような順序で、メッセージが SSFNet をシミュレートネットワークとして伝播していき分散アルゴリズムが展開されていく。Neko 側の役割は、Neko プロセスからのメッセージを SSFNet へ渡すことで、これは、send イベントを生成していることに対応する。SSFNet 側の役割は、Neko からの送信イベントを処理し、そのイベントに関するメッセージを Neko へ渡すことであり、これは、deliver イベントを生成していることに対応する。つまり、シミュレートネットワークとして SSFNet を用いる場合、Neko 側の役割は、send イベントを生成することで、SSFNet 側の役割は、deliver イベントを生成することである。2つのスケジューラを用いて、この2つのイベントが一貫性を持って処理されるように、つまり、メッセージの伝播が適切に行われるように Neko へ SSFNet の統合を行った。

## スケジューラの同期

2つの離散的なイベント駆動のシミュレータをイベント発生におけるシミュレーション時間の一貫性を持って、並行に動作させることは容易ではない。2つのスケジューラを単純に動作させた場合、どちらかのスケジューラがシミュレーション時間上、先に進んでいるか、もしくは、まだその時間に達していないので、メッセージを各プロセスへ伝播できない。つまり、これは、まったくシミュレータとして機能していないことに相当する。

この問題は、Neko 側の send イベントと SSFNet 側の deliver イベントを一貫性を持って、処理できるように、スケジューラを下記のようなルールに従って、交互に動作させることで解決した。

- Neko のスケジューラは、send イベントをその send イベントが起こる時間付きで、SSFNet へ渡し（SSFNet のソケットを用いてメッセージを SSFNet へ渡していること）停止した後、SSFNet のスケジューラを動作させる。
- また、Neko のスケジューラは、SSFNet からの deliver イベントを実行したら（Neko プロセスへメッセージを伝播すること）停止し、SSFNet のスケジューラを動作させる。
- SSFNet のスケジューラは、Neko から渡された send イベントをいっしょに渡された時間に従って、実行したら（メッセージの伝播を SSFNet でシミュレートし deliver イベントを生成）deliver イベントか、Neko スケジューラ中の次のイベントの時間まで動作する。deliver イベントの場合は、そのイベントを Neko へそのイベントが起こる時間付きで渡し停止した後、Neko のスケジューラを動作させる。

この同期方法を単純に説明すると、Neko の次のイベントのシミュレーション時間まで SSFNet のスケジューラを動作させる機会を与えることである。このルールに基づいてスケジューラを動作させることで、シミュレーション時間上の一貫性を持って、メッセージの伝播を処理することが可能となる。

## アドレス方式の違い

NEKO では、プロセスを識別するためにシーケンシャルな ID を用いている。一方、SSFNet では、計算機の識別に IP アドレスを用いている。NEKO では、Neko メッセージに含まれている宛先 ID を用いて、メッセージの伝播を行っている。しかし、IP アドレスを指定して送信できないので SSFNet を利用して、他の NEKO プロセスに送信することができない。NEKO の ID と SSFNet の IP アドレスを 1 対 1 に対応させるマッピングテーブルを作成した。NEKO からの送信時と SSFNet からの受信時にこのマッピングテーブルを用いることで、アドレス方式の違いを解決した。

## メッセージ形式の違い

NEKO のメッセージは、1 つ以上の宛先 ID と文字列をメッセージとして送信する。一方、SSFNet は、送信時にメッセージサイズを指定してオブジェクトを送信する。よって、

SSFNet では、NEKO メッセージをオブジェクトとして、変更を加えることなく、送信することが可能である。しかし、通常、NEKO ではメッセージを送信する場合、NEKO メッセージの複数の宛先 ID を指定してメッセージの送受信を行っている。SSFNet のソケットを用いて、Neko メッセージを送信する場合は、マッピングテーブルを用いて従来の送信方法を用いて、IP アドレスを指定し、そのメッセージを宛先の SSFNet 上の Host へ送信可能である。しかし、SSFNetNetworks を用いて、Neko メッセージを Neko プロセスに渡す場合、そのメッセージは、一意に指定したプロセスへの伝播されなければならない。しかし、Neko メッセージの宛先 ID を用いると、その宛先 ID に対応したすべての Neko プロセスへ送信してしまうことになり、不要なメッセージの複製を行ってしまう。

この問題を解決するために、NEKO プロセスを一意に識別できる情報を付加した SSFNet メッセージクラスを作成した。SSFNet メッセージは、SSFNet のソケットを用いてメッセージを伝播する時に、宛先の Host の IP アドレスを用いて割り出した（マッピングテーブルより）Neko プロセスの ID を付加する。Neko は、SSFNet Networks を用いて Neko プロセスにメッセージを渡す時に、Neko メッセージとして伝播されてきた SSFNet メッセージから、最終的な宛先である Neko プロセスの ID を用いて、そのメッセージを一意に指定された Neko プロセスへ渡す。

### 3.1.1 SSFNet Networks

SSFNet を Neko のシミュレートネットワークとして用いるためのクラスである SSFNetNetworks (図 3.1) について詳しく説明する。シミュレートネットワークとして SSFNet を指定したい場合、Neko のコンフィグファイル上に、ある SSFNetNetworks の種類のクラスをパラメタとして指定することで利用可能である。

SSFNetNetworks は、TCP と UDP の 2 つの種類のクラスに分類している。TCP クラスのネットワークは、SSFNetTCPNetwork と呼ばれ、メッセージを伝播するためのプロトコルとして信頼性の高い TCP/IP を用いている。ネットワークシミュレータとして、SSFNetTCPNetwork が指定されると、Neko の初期化時に、SSFNet 上の各 Host がプロセス数に応じて、受信用と送信用の TCP ソケット生成し、各 Host 同士がメッセージを伝播するための TCP コネクションを張る。

一方、UDP クラスのネットワークは、SSFNetUDPNetwork と呼ばれ、メッセージを伝播するためのプロトコルとして、信頼性の低い UDP/IP を用いている。このネットワークは、Neko の初期化時に、1 つだけ UDP ソケットを生成する。

これらのネットワーククラスは、Neko のコンフィグファイルを用いて、複数指定することができる。例えば、TCP のクラスを 2 つ選択した場合は、各 Host 間に 2 組の送受信用のソケットが生成され、コネクションが張られる。UDP のクラスを選択した場合も同様である。また、TCP と UDP のクラスを混合して、送信することが可能であり、開発者は、ネットワークモデルに応じて複数のソケットを利用することができる。

## 3.2 動作テスト

本項では、前述の方法を用いて、統合を行った Neko の動作テストを行う。シミュレーションを行うアルゴリズムは、アトミックブロードキャストの1つである Lamport[9] のアルゴリズムを用いる。このアルゴリズムは、イベント数がプロセス数  $n$  に対して、 $O(n^2)$  の計算量で増加し、同一の時間に複数のイベントを処理する。また、すでに Neko を用いてクラスとして実装されており、テストケースとして適している。既存の Neko に実装されているシミュレーションネットワークである BasicNetwork および MetricNetwork と、本研究で実装したネットワークである SSFNetTCPNetwork および SSFNetUDPNetwork を用いてテストを行う。計測する対象は、既存の Neko 単体と SSFNet を用いた場合のシミュレーションにおける実時間の比較と、SSFNet を用いた場合、最大どれくらいのプロセス数をシミュレート可能であるかの2つである。実時間の比較方法は、プロセス数が10で、Lamport のアルゴリズムを1000回連続で実行した時の実時間を計測する。最大プロセス数の計測は、プロセス数の値を変更しながら、Lamport のアルゴリズムを1回だけ実行し、いくらのプロセス数で実行できないかの境界を計測する。評価環境は、Cpu clock が2.6GHz、Memory が約1GByte である。経過時間の計測には、Unix のコマンドである time コマンドを用いる。

	Simulated Network class	execution time
Neko	BasicNetwork	5.0 [sec]
	MetricNetwork	65.5 [sec]
SSFNet	SSFNetUDPNetwork	63.5 [sec]
	SSFNetTCPNetwork	126.2 [sec]

表 3.1: Unit test

Neko の BasicNetwork と SSFNet を用いた場合を比較すると、実時間レベルで最大13倍長い。これ2つの理由がある。1つ目は、Neko のスケジューラのスレッドと SSFNet のスケジューラを交互に動作させているために、スレッドを切り替えるオーバーヘッドが生じている。2つ目は、SSFNet は、ネットワークを離散的なイベント駆動でシミュレートしているため、Neko のネットワークシミュレータよりもシミュレーションに多くの時間を有するためである。Neko の MetricNetwork と SSFNet のネットワークを比較すると、実時間レベルで最大2倍長い。MetricNetwork は、BasicNetwork と比べ複雑なシミュレートを行っており、SSFNetUDPNetwork とほぼ同じ経過時間でシミュレート可能なことがわかる。SSFNet を用いた場合の最大プロセス数は、Neko 単体と比較して、かなり制限があることがわかる。SSFNetUDPNetwork を用いた場合のシミュレート可能な最大プロセス数は253である。この理由は、Neko の実ネットワーク環境で実験を行う場合に用いられている UDPNetwork クラスを継承して用いているからである。実ネットワークで実験を行う時に UDPNetwork クラスは、コンフィグレーションファイルの実 IP アドレスを利用



している。SSFNet では、IP アドレスとして、実際に IPv4 など用いられているアドレスのフォーマットを用いておらず、順番に 1 から番号付けしているために制限が起きてしまう。また、SSFNetTCPNetwork を用いた場合は、シミュレート可能な最大プロセス数は 53 である。この理由は、スレッドを多く生成しすぎて、Java の実行環境 “Memory out of error” により停止してしまいますためである。SSFNetTCPNetwork では、各プロセスが送信用と受信用のソケットをプロセス数だけ生成する。そのソケットの数だけ、スレッドを用いて送信と受信を行っているので、プロセス数が増える度にスレッドを生成してしまう。

SSFNetUDPNetwork の問題は、Neko 上の UDPNetwork クラスをリファクタリングを行うことで、シミュレート可能なプロセス数を増やすことは可能あると考える。SSFNetTCPNetwork の問題は、JAVA の実行環境の問題か OS のスレッド生成の上限に関する問題かのいずれかであると考えられる。この問題についての解決方法は、現在調査中である。

## 第4章 How to use Neko with SSFNet

本章では、ネットワークシミュレーション環境としてSSFNetを用いてシミュレーションを行う方法を解説する。はじめに、Neko上でSSFNetを利用する場合のインストール方法について解説する。次に、サンプルアルゴリズムを用いて、NekoとSSFNetのコンフィグレーションファイルの設定方法とおよび実行方法について解説する。

### 4.1 インストール方法

はじめに、SSFNetとNekoをダウンロードしなければならない。SSFNetはライセンス規約により、自由に配布できないので、ユーザ自身が<http://www.ssfnet.org/>より、ダウンロードしなければならない。Nekoは、<http://lsrwww.epfl.ch/neko/>から取得できる。次にSSFNetを任意のディレクトリにインストールする。インストール方法は、SSFNetのwebページに記載されているので、ここでは省略する。つぎに、Nekoを任意のディレクトリに展開する。nekoというが生成されるので、そのディレクトリは移動して、下記のようにコマンドを入力する。

```
例:SSFNetがインストールされたディレクトリ (/home/***/work/ssfnet)

$ ./configuer -java='Java コマンドのパス' -ssfnet='/home/***/work/ssfnet/lib'
$ ../dest_neko/bin/ant
```

上記のような手順で、インストールすることにより、SSFNetのクラスパスがNekoのコマンド群に設定される。<sup>1</sup>

---

<sup>1</sup>Nekoのコマンド群を使うために、パス変数にNekoがインストールされたディレクトリを設定しなければならない(詳細は<http://lsrwww.epfl.ch/neko/>)

## 4.2 コンフィグファイルの設定

シミュレーションを実行するためのコンフィグファイルの設定方法について解説する。基本的には、ネットワークシミュレータとして SSFNet を用いる場合、従来の Neko のコンフィグファイルに SSFNet 特有のパラメタを追加するだけで実行可能である。SSFNet 特有のパラメタとは、SSFNet の基本となる DML ファイル、開発者がネットワークモデルを記述した DML ファイル、メッセージサイズおよびネットワークモデルの Host 数である。まずはじめに、Neko のコンフィグファイルの記述方法を解説する。つぎに、SSFNet のコンフィグファイルである DML の記述方法とサンプルネットワークモデルの解説をする。シミュレーションを行うアルゴリズムは、すでに Neko 上に実装されている Lamport の Atomic Broadcast を用いる。

### Neko コンフィグファイルの設定

Neko のコンフィグファイルの記述方法は、`パラメタ名 = 値 or 文字列 or 真偽値` の単純な形式である。開発者は、自分自身で実装したアルゴリズムに、このコンフィグファイルからパラメタを設定することも可能であり、自由なパラメタが設定可能である。Neko では、コンフィグファイルからパラメタを取得する方法として、Apache の org.java.util パッケージの configurations クラスを利用しており、パラメタ名を複数指定することで配列としてパラメタを取得可能である。

図 4.1 を参照しながら各パラメタについて説明する。一行目の simulation は、false の場合は、分散アルゴリズムを実行する環境として、実ネットワークを用いる。実ネットワークを指定した場合、プロセスが実行している計算機の IP アドレスやポート番号などを記述しなければならない。本研究では、実ネットワーク特有のパラメタについては解説しない。5 行目の process.num は、アルゴリズムを実行するプロセスの数である。この値を変更することにより、実行するプロセスの数を変更することが可能だが、シミュレーションネットワークとして SSFNet を用いる場合、後に説明する host.num の値以下でなくてはならない。なぜなら、SSFNet では、DML ファイルにより、ネットワーク上で利用できる計算機の数記述されており、その数以上のプロセスは、SSFNet の計算機に配置できないからである。7 行目の host.num は、SSFNet のシミュレーションネットワーク上にシミュレートされる計算機の数である。Neko のプロセスは、SSFNet の計算機と 1 対 1 に対応しており、その計算機上で Neko のプロセスが実行されることになる。8 行目の algorithm は、プロセス上で実行される分散アルゴリズムのクラス名である。9 行目の process.initializer は、Neko のプロセスの振る舞いであるアルゴリズムを階層的に構築するためのクラス名である。このクラスが algorithm で指定されたアルゴリズムのクラス名を動的に読み込み、階層的にアルゴリズムを構築する。10 行目の network は、シミュレーションネットワークのクラス名である。このパラメタを変更することにより、シミュレーションネットワークの変更が可能である。このサンプルコンフィグファイルでは、SSFNetUDPNetwork を指定しているが、他のネットワークや、複数のネットワークを指定することが可能である。この値を用いて、ネットワークの切り替えを簡単に可能にして

いる。12行目からは、シミュレーションネットワークとしてSSFNetを用いた場合の固有のパラメタである。13行目のmessage.sizeは、Nekoのプロセスが通信に用いる1つあたりのメッセージのサイズであり、単位はbyteである。14行目のmy.dmlは、開発者が記述したDMLファイルの絶対パスである。DMLファイルは、複数のファイルから構成することができるため、複数のDMLファイルを指定したい場合は、複数のmy.dmlパラメタを用いてすべてのDMLファイルの絶対パスを記述しなければならない。16行目のssfnet.dmlは、SSFNetが記述したDMLファイルの基本となる構成を記述したファイルである。このファイルは、SSFNet単体で実行させる場合でもかならず指定しなければならない。18行目のdirect.mappingは、SSFNetのある計算機に対して、Nekoのプロセスを指定して対応させたいときに用いる。パラメタにtrueの値を指定すると、20行目から26行目のようなssfnet.processを複数記述して、NekoのプロセスのIDとSSFNetのNHIアドレスを用いて記述する。記述形式は、'NekoのプロセスID / SSFNetのNHIアドレス'である。

以上が、Nekoのコンフィグファイルの記述方法と各パラメタについての説明である。このサンプルコンフィグファイルは、必要最低限の実行に必要な設定である。この他のパラメタとして、log4jを用いたログの取得方法やその出力ファイル名の指定など、さまざまなパラメタが存在する。さらに詳細な情報を知りたい場合は、インストールの章で紹介したNekoのwebページで調べてほしい。

#### SSFNet DML ファイルの設定

SSFNetのコンフィグファイルであるDMLファイルの記述方法とその各パラメタについて説明する。サンプルのDMLファイルで記述されたネットワークモデルのトポロジと各パラメタは、図4.2のように表すことができる。単純なLocal Area Networkをモデリングしており、1つのルータに対して、1つのハブを中継して7つの計算機が接続されている。各計算機とルータにID番号が割り振られているが、これはIPアドレスではなく、SSFNet固有のNHIアドレスである。すべてのリンクスピード、リンク間のパケット遅延およびNIC(Network Interface Card)上の1つのパケットに対する遅延であるレイテンシは、同じ値にモデリングしており対象である。各計算機上にあるNICのレイテンシーは、1パケットごとの送信遅延を表現しており100msであり、NICのBufferは、8000byteである。ルータのインタフェースのレイテンシーも100msであり、Bufferは、16000byteである。ハブは、SSFNetではシミュレートされていないが、複数の計算機をルータに接続する場合において、DMLファイルの属性としてハブを用いることにより簡単にルータに各計算機を接続する状態を記述可能となる。もし、このハブを用いないとルータは、接続される各ホストの数だけインタフェースを用意しなければならないし、それらが接続されていることをDMLで記述しなければならないため、多くの手間がかかる。

サンプルDMLファイル(図4.3および図4.4)について説明する。このDMLファイルは、Nekoのコンフィグファイルの14行目のmy.dmlで指定している`evaluation-net-p7.dml`のファイルに記述されている。図4.3は、主にネットワークのトポロジを記述しており、図4.4は、ネットワークで用いられる各属性の値を記述している。

```
1: simulation = true
2: #
3: # The number of communicating processes.
4: #
5: process.num = 5
6: # The number of Host on DML (Network model)
7: host.num = 7
8: algorithm = lse.neko.examples.basic.Lamport
9: process.initializer = lse.neko.examples.basic.TestInitializer
10: network = lse.neko.sim.ssfnet.udp.SSFNetUDPNetwork
11: #network = lse.neko.sim.ssfnet.tcp.SSFNetTCPNetwork
12: #User defined parameter on SSFNet's DML file
13: message.size = 200
14: my.dml = /home/m-tomo/work/nekoTest/finalTest/evaluation-net-p7.dml
15: #ssfnet dml
16: ssfnet.dml = /home/m-tomo/usr/ssfnet/examples/net.dml
17: # correspondent Neko process to SSFNet Host.
18: direct.mapping = false
19: #direct.mapping = true
20: ssfnet.process = 0/7
21: ssfnet.process = 1/6
22: ssfnet.process = 2/5
23: ssfnet.process = 3/4
24: ssfnet.process = 5/3
25: ssfnet.process = 6/2
26: ssfnet.process = 7/1
```

☒ 4.1: Sample Neko configuration file

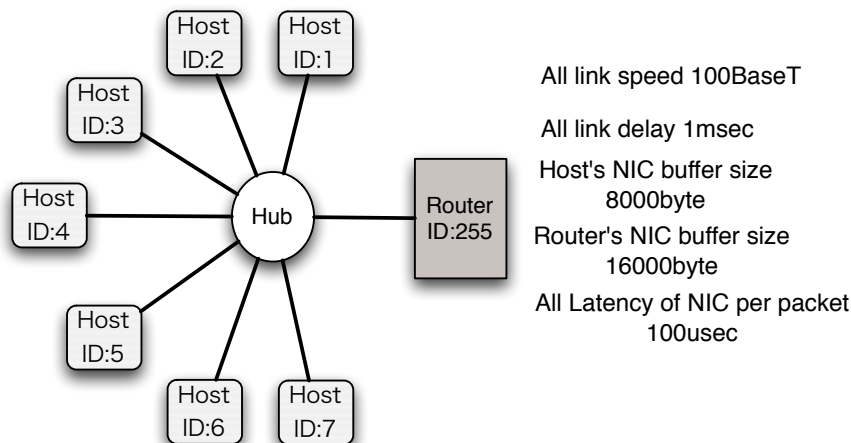


図 4.2: Network model of sample SSFNet DML file

はじめに、図 4.3 から説明する。一行目の Net は、Net の中に Net を記述することも可能であり、実ネットワークの世界でいうサブネットワークを定義している。2 行目の frequency は、最上位の Net に記述することが可能であり、最小のシミュレーション時間の単位を定義している。つまり、frequency は、1GHz なので 1ns である。この値は、各インタフェースの bitrate と関係しており、もし、bitrate が 100,000,000 (100MbaseT) ならば、100msec 以上でなければならない。よって、この frequency は、1GbaseT のリンクスピードまでシミュレート可能である。3 行目の router は、TCP/IP 参照モデルのネットワーク層に分類されるルータと呼ばれる機器を定義している。NHI アドレスとして、255 番の ID 番号を割り振っており、graph によりプロトコルとして IP を用いている。もし、LAN 以上の大きなネットワークモデルを記述する場合は、graph に OSPF や BGP のプロトコルを IP プロトコルの上位に記述しなければならない。interface は、計算機が外部つまり、ネットワークにパケットを送信するためのインタフェースを記述する属性である。すべての計算機 (Host,router,etc) は、graph と interface を持たなければならない。interface にも、NHI アドレスが割り振られており、このインタフェースの NHI アドレスは `255:0` となり、このアドレスを用いて、他の機器は接続する。\_extend は、他の DML ファイル名か同一ファイル内の他のトップレベルの属性の下位の属性のパラメタをすべての参照することができる。つまり、図 4.4 の dictionary の 100BaseT の 2 つの値を参照していることになる。9 行目の host は、計算機を定義しており、idrange を用いることにより、一度に複数の host を定義することが可能である。この場合は、ID1 から ID7 までを指定しているため、合計 7 つの計算機を定義している。nhi\_route は、ハブを中継して計算機 (このモデルの場合、ルーター) に接続すること定義している。つまり、このように定義することにより、簡単にルータに接続している状態を定義可能である。この host も router と同様に、dictionary の process 属性の値をすべてを参照している。interface について記述されており、さらに、host の graph は、複数のプロトコルから記述されている。graph は、そ

の計算機のプロトコルスタックを定義している。計算機のプロトコルスタックの構成と各属性については、後の図 4.4 の説明時に詳細に解説する。1 3 行目の link は、TCP/IP 参照モデルのデータリンク層を定義しており、NHI アドレスを用いて各計算機とルータの接続を定義している。\_find を用いて、1 つだけ他の属性を参照することが可能であり、リンク間の遅延である delay を参照している。以上がネットワークのトポロジの部分についての説明と記述方法である。

```
1: Net [
2:     frequency 1000000000
3:     router [
4:         id 255
5:             graph [ProtocolSession [name ip use SSF.OS.IP]]
6:             interface [ id 0 buffer 16000 _extends .dictionary.100BaseT ]
7:         ]
8:     #####all latency of nic is 0.0001 [sec]####
9:     host [ idrange [from 1 to 7]
10:         nhi_route [dest default interface 0 next_hop 255(0)]
11:         _extends .dictionary.process
12:     ]
13:     link [attach 255(0)
14:         attach 1(0)
15:         attach 2(0)
16:         attach 3(0)
17:         attach 4(0)
18:         attach 5(0)
19:         attach 6(0)
20:         attach 7(0)
21:         _find .dictionary.link_delay.delay
22:     ]
23: ]
```

図 4.3: Sample SSFNet DML file (part of Net)

次に、ネットワークの各パラメータが定義された dictionary (図4.4) について説明する。3行目で参照されてある12行目の `NekoProxy_standard.graph` で定義されてあるプロトコルスタックについて説明する。最上位の `ProtocolSession` クラスつまりプロトコルは、本研究で作成した `Neko` と `SSFNet` のインタフェースである `NekoProxy` クラスを指定している。DML ファイル中の計算機のプロトコルスタックを構築する定義するとき、最上位層に `NekoProxy` クラスを指定し、`Neko` のコンフィグファイルの `network` に `SSFNetNetwork` クラスを指定するだけで、`Neko` のシミュレーションネットワークとして、`SSFNet` を用いることができる。次の下位層である `socketMaster` クラスは、ソケットインタフェースをシミュレートするクラスであり、下位層のプロトコルである `tcpSesseionMaster` クラスか `udpSesseionMaster` クラスを用いて、TCP ソケット、UDP ソケットを生成する。最下位層の IP クラスは、TCP/IP 参照モデルのトランスポート層をシミュレートしている。この graph つまりプロトコルスタックは、`SSFNet` を `Neko` のシミュレートネットワークとして用いる場合は、常に同様の定義である。従って、計算機のプロトコルスタックを定義する場合は、必ずこの記述を用いなければならない。7行目から11行目は、各計算機のインタフェースの仕様を定義したもので、例として記述している。28行目の `tcpinit-default-ssfnet` は、TCP ソケットを用いる場合の TCP プロトコルの仕様を定義したものである。ここに定義した仕様は、`SSFNet` がデフォルトで定義してある仕様をそのまま用いた。開発者は、自分自身で定義したネットワークモデルに従って、この仕様を変更することができる。最後に46行目の `udpinit` は、UDP プロトコルの仕様を定義しており、変更できるパラメータとして、一度に送信可能なパケットサイズを定義できる。

### 4.3 実行

設定したコンフィグファイルを用いて、シミュレーションする方法とその出力結果について説明する。シミュレートするアルゴリズムは、Lamport の Atomic Broadcast で、シミュレートネットワークとして UDP ソケットを用いて通信を行う `SSFNetUDPNetwork` を用いる。`Neko` のコンフィグファイル名は、`sample.config` とする。ターミナルで下記のように、`neko` コマンド<sup>2</sup>に `Neko` のコンフィグファイルを引数として実行する。

```
$ neko sample.config
```

`SSFNet` がネットワークモデル初期かし、そのネットワークモデルの情報を標準出力した後、`Neko` のシミュレーションが開始される。そして、各プロセスがアルゴリズムを終了したことを出力し、最後にアルゴリズムによって、経過した時間を出力する。出力結果は下記の通りである。

---

<sup>2</sup>`dest_neko/bin` にパスが通していること



SSFNet の初期化情報が出力される

```
.  
.br/>p0 finished  
p1 finished  
p2 finished  
p4 finished  
finishing at 2.3568
```

通常は、この出力結果を性能評価に用いるわけではない。Neko のコンフィグファイルに、メッセージの通信履歴を取るために log4j クラスを用いるように設定する。もしくは、アルゴリズム中に異なった形式をファイルに出力するように実装することもできる。開発者は、いずれの場合も、その通信履歴を用いてグラフ化などを行い性能評価に利用する。この章では、SSFNet を Neko のシミュレートネットワークとして用いるために、コンフィグファイルの設定方法と実行方法および実行結果について説明した。もっと詳細な方法を知りたいければ、インストールの説明時に紹介した Neko と SSFNet の web サイトで入手可能である。

```

1: dictionary [
2:     process [
3:         interface [id 0 buffer 8000 _extends .dictionary.100BaseT]
4:         _find .dictionary.NekoProxy_standard.graph
5:     ]
6:     # standard network interface
7:     10BaseT [bitrate 10000000 _find .dictionary.nic_latency.latency]
8:     100BaseT [bitrate 100000000 _find .dictionary.nic_latency.latency
9: ]
10:     1GBaseT [bitrate 1000000000 _find .dictionary.nic_latency.latency ]
11:     nic_latency [latency 0.0001]
12:     link_delay [delay 0.001]
13:     NekoProxy_standard [graph [
14:         ProtocolSession [
15:             name NekoProxy use lse.neko.sim.ssfnet.NekoProxy
16:         ]
17:         ProtocolSession [name socket use SSF.OS.Socket.socketMaster]
18:         ProtocolSession [
19:             name tcp use SSF.OS.TCP.tcpSessionMaster
20:             #_extends .dictionary.tcpinit-default-ssfnet
21:         ]
22:         ProtocolSession [
23:             name udp use SSF.OS.UDP.udpSessionMaster
24:             #_extends .dictionary.udpinit
25:         ]
26:         ProtocolSession [ name ip use SSF.OS.IP]
27:     ]]
28:     tcpinit-default-ssfnet[
29:         ISS 0          # initial sequence number
30:         MSS 1024       # maximum segment size
31:         RcvWndSize 16  # receive buffer size in units of MSS
32:         SendWndSize 16 # maximum send window size in units of MSS
33:         SendBufferSize 16 # send buffer size in units of MSS
34:         MaxConWnd 64   #
35:         MaxRexmitTimes 12 # maximum number of retransmission times
36:         TCP_SLOW_INTERVAL 0.5 # granularity of TCP slow timer, sec
37:         TCP_FAST_INTERVAL 0.2 # granularity of TCP fast(delay-ack), sec
38:         MSL 3          # maximum segment lifetime, sec
39:         MaxIdleTime 600 # maximum idle time, sec
40:         delayed_ack true # delayed ack option, true or false
42:         fast_recovery 60 # use fast recovery algorithm, true/false
43:         show_report false # print summary connection report, true/false
44:         debug false
45:         # if true, dump verbose TCP diagnostics to files
46:     ]
47:     udpinit [
48:         max_datagram_size 10000 # max UDP datagram size (payload bytes)
49:         debug false 28 # print verbose UDP diagnostics
50:     ]

```

⊗ 4.4: Sample SSFNet DML file (part of dictionary)

## 第5章 性能評価

本研究で SSFNet と統合した Neko を用いて、分散アルゴリズムの 1 つである Atomic Broadcast の性能評価を行う。Atomic Broadcast については、後に詳細を説明する。性能評価を行う目的は、複雑な分散アルゴリズムをシミュレートネットワークとして SSFNet を選択した場合、つまり、本研究で実装し統合した Neko 上でシミュレートが可能であることを示すことで、その統合された Neko の有効性を示すことである。また、その有効性を示しつつ、Atomic Broadcast の性能評価も同時に行う。本研究で性能評価を行う Atomic Broadcast は、すでに Neko 上で実装されているアルゴリズムである。Urban ら [5] によって行われた評価結果を参考にし、Atomic Broadcast に必要不可欠な分散アルゴリズムである Consensus の異なった 2 種類の性能評価を行う。2 つの異なったアルゴリズムの特徴を浮かび上がらせることで、システム的设计者により、システムに適切なアルゴリズムの選択を提供する判断材料とさせる。しかし、分散アルゴリズムをこのように評価するためには、計算機とネットワークのリソースを考慮にいれなければならない。従って、シミュレートネットワークとして SSFNet を用いることで、より現実的な性能評価結果が得られる。つまり、システム設計者は、より適切な評価に基づいた判断ができる。

### 5.1 Atomic Broadcast とは

Atomic Broadcast (Total Order Broadcast と呼ばれる) は、分散システムにおいて基本となる問題である。Atomic Broadcast は、すべてのメッセージを同じ順序 (全順序) で届けるプリミティブな Broadcast として定義される。例えば、Process Replication (プロセスの複製技術)、複製かされたデータベースのトランザクション処理をサポートするために用いられる。Atomic Broadcast といっても、約 60 以上のアルゴリズムが存在している。Atomic Broadcast の基本的な定義と 5 つのクラスに分類した Atomic Broadcast の特徴について説明していく。

#### 5.1.1 Atomic Broadcast の定義

Atomic Broadcast は、2 つのプリミティブな操作から成り立っている。(m はメッセージ)

- A-broadcast(m)
- A-deliver(m)

さらに、前提条件2つがあり、すべてのメッセージは、一意に識別可能であることと、送信プロセスのよって与えられたその識別子はメッセージといっしょに含まれて運ばれることである。Atomic Broadcast を定義する特性について説明する。ここでいう正確なプロセスとは、クラッシュなどの理由により、まちがった内容のメッセージを送信することや、メッセージの処理が極端に遅いプロセスではなく、正常に期待した振る舞いを行うプロセスのことである。

**VALIDITY : 正確性** 正確なプロセスがメッセージ  $m$  を A-broadcast する場合、最終的にメッセージ  $m$  を A-deliver することである。

**AGREEMENT : 同意** 正確なプロセスがメッセージ  $m$  を A-deliver する場合、最終的には、すべての正確なプロセスはメッセージ  $m$  を A-deliver する

**INTEGRITY : 完全** 任意のメッセージ  $m$  のために、すべてのプロセスはメッセージ  $m$  を多くて1度しか A-broadcast しない。これは、メッセージ  $m$  が前もって、送信プロセスによりメッセージ  $m$  を A-broadcast されている場合に限る

**TOTAL ORDER : 全順序** 正確なプロセス  $p$  と  $q$  がメッセージ  $m$  とメッセージ  $m'$  を A-deliver した場合、 $p$  はメッセージ  $m'$  の前にメッセージ  $m$  を A-deliver する。これは、 $q$  がメッセージ  $m'$  の前にメッセージ  $m$  を A-deliver する場合に限る。

AGREEMENT と TOTAL ORDER にさらに Uniform 特性を加えることにより、Uniform Atomic Broadcast を定義することが可能である。

**UNIFORM AGREEMENT** 正確なプロセスがメッセージ  $m$  を A-deliver する場合、すべての正確なプロセスは、最終的にメッセージ  $m$  を A-deliver する。

**UNIFORM TOTAL ORDER** プロセス  $p$  とプロセス  $q$  (正確または、不正確なプロセス) が、メッセージ  $m$  とメッセージ  $m'$  を A-deliver する場合、プロセス  $p$  はメッセージ  $m$  を A-deliver する。これは、プロセス  $q$  がメッセージ  $m'$  を A-deliver する場合に限る。

Uniform 特性は、ただひとつのプロセスでも A-broadcast または、A-deliver できない場合は、すべてのプロセスも A-broadcast または、A-deliver しないということである。Uniform 特性は、Atomic commit または、Active Replication の用いるような特定のアプリケーションに必要とされている。例えば、複製されたプロセスがある副作用を伴うような可能性があるシステムなどに利用される。それらのシステムは、異なったプロセスからの矛盾するような振る舞いが生じる Atomic Broadcast を用いている。しかし、Uniform 特性は、すべてのアプリケーションで必要としないし、残念なことに大幅なコストがかかる。システム設計者は、Uniform または、non-Uniform な Atomic Broadcast を区別して考慮しなければならない。

## 5.1.2 Atomic Broadcast アルゴリズム

Atomic Broadcast を約 60 のアルゴリズムを 5 つのクラスに分類し、5 つのクラスについて説明する。異なる 3 つのタイプの 1 つのプロセス (sender、destinations、sequencer) によって、生成されたメッセージの順番によって分類される。その他の要因で、3 つのタイプから、下記のような 5 つのクラスに分類することができる。

1. Fixed sequencer
2. Moving sequencer
3. Privilege based
4. Communication history
5. Destinations agreement

### Fixed Sequencer

Fixed Sequencer (図 5.1) は、非常に簡単アルゴリズムである。システム内の 1 つのプロセスが Sequencer に選ばれ、すべてのメッセージの順番を決める。図 5.1 では、p1 が Sequencer である。アルゴリズムは、あるプロセス p がメッセージ m を broadcast したい場合、まずはじめに Sequencer である p1 に対して、メッセージ m を送信する。Sequencer である p1 がメッセージを受信すると、メッセージ m にシーケンス番号を割り当て、シーケンス番号といっしょにメッセージ m を他のすべてのプロセスへ送信する。non-uniform の場合は、プロセスがメッセージ m とシーケンス番号を受信すると、即座にメッセージ m を deliver する。uniform の場合は、Sequencer が ack メッセージを用いて、すべてのプロセスが安定 (正確なプロセス) していることを確認してから、メッセージ m を deliver する。このアルゴリズムは、めったに uniform であることはない。なぜなら、他のクラスのアルゴリズムより比較的成本が高いからである。

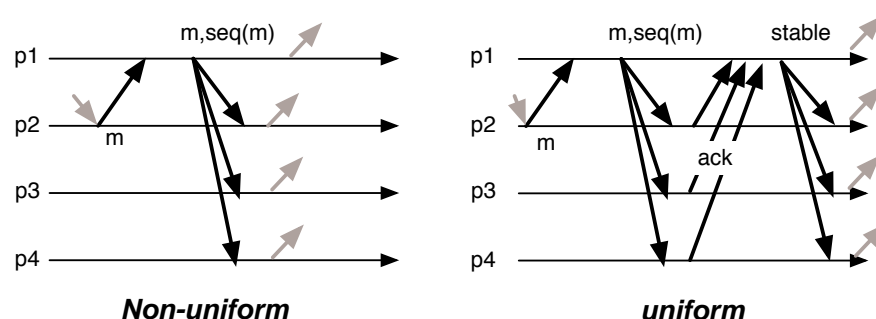


図 5.1: Fixed Sequencer

## Moving Sequenser

Moving Sequenser (図 5.2) の Sequencer の役割は、あるプロセスから他のプロセスへ渡されることである。これは、トークン (優先権番号) によって行われ、トークンはシーケンス番号を運び、プロセス間を常に移動している。

Non-uniform の場合は、プロセス p 2 がメッセージ m を broadcast する場合、p 2 はすべての他のプロセスにメッセージ m を送信する。各プロセスは、メッセージ m を受信すると receiveQueue にメッセージ m を蓄える。トークンを保持しているプロセス p1 は、receiveQueue に蓄えられている最初にメッセージである m にシーケンス番号を割り当てる。そして、トークンを次のプロセス p2 に送信し、シーケンス番号といっしょにメッセージ m を broadcast する。あるプロセスは、以下のような場合に、メッセージ m を deliver することが可能となる。

- メッセージ m が受信されている場合
- メッセージ m のシーケンス番号が受信されている場合
- メッセージ m の前の各メッセージが deliver されている場合

uniform の場合は、宛先プロセスがメッセージ m を deliver できる前に、メッセージ m が安定するまで待たなければならないことと似ている。(例えば、そのプロセスがすべてのプロセスにより応答されるまで) メッセージ m が安定することを確認するには、ack をかき集めたトークンにより行われる。トークンは、常に、プロセス間を循環している必要はなく最終的にストップする。

Non-uniform は、メッセージ m が安定することを確認する必要があるため、uniform よりも早くトークンを止めることができる。

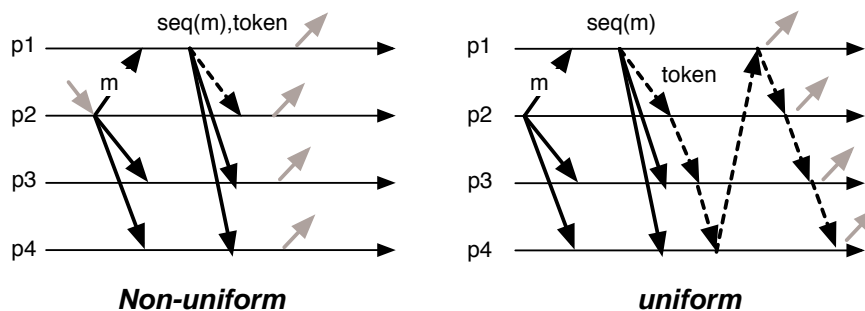


図 5.2: Moving Sequenser

## Privilege-Based

Privilege Based (図 5.3) では、delivery する順番を決定するのは、送信プロセスである。あるプロセスがメッセージ m を broadcast する場合、そのプロセスは、最初に privilege を取得しなければならない。Non-uniform の場合、プロセス p1 がメッセージ m を broadcast

する場合、プロセス p1 はトークンを受信するまで m を sendQueue に蓄えておく。トークンは、シーケンス番号を運び常にプロセス間を循環している。プロセス p1 がトークンを受信すると、sendQueue から最初のメッセージ m を取り出し、トークンによって運ばれたシーケンス番号を用いて、メッセージ m にシーケンス番号を割り当てる。そして、シーケンス番号といっしょにメッセージ m を broadcast する。次に、プロセス p1 はトークンのシーケンス番号を 1 つインクリメントし次のプロセス p2 へトークンを送信する。送信するメッセージの数を減らすために、プロセスは、単一のメッセージ中として、メッセージ m といっしょにトークンを送信することができる。プロセスは、メッセージ m を受信したときに、そのシーケンス番号によってメッセージ m を deliver することができる。uniform の場合は、トークンはさらに ack を運んでいる。メッセージ m が deliver される前に、プロセスはメッセージ m が安定するまで待たなければならない。これを実現するためには、十分なトークンの round trip を必要とする。

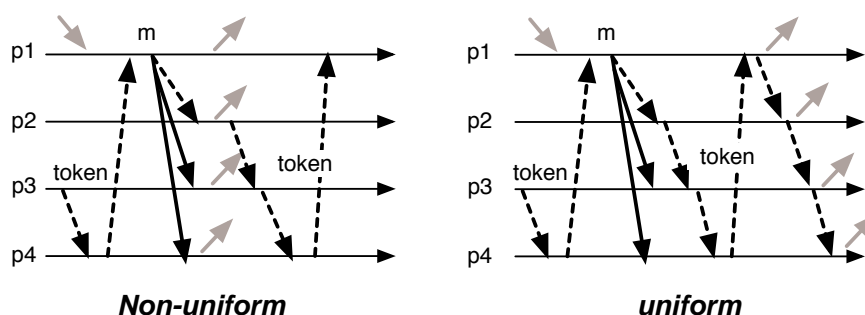


図 5.3: Privilege Based

### Communication History

Communication History 図 5.4 は、メッセージを deliver する順番を Privilege Based のように送信プロセスによって決定する。Privilege Based と違う点は、各プロセスがいつでもメッセージを送信可能なことである。宛先のプロセスは、もはや deliver するメッセージが total order (全順序) を破らないことを確認するために、他のプロセスから生成されたメッセージを監視する。ほとんどの Communication History では、メッセージの causal history によって定義された不完全な順番を用いる。そして、あらかじめ決定された方法に従って順序化されたメッセージが一致する時、この不完全な順番を total order に変換する。

Communication History のアルゴリズムを説明する。ある不完全な順番は、 $A\text{-broadcast}(m)$  イベントの論理時間を含んでいる各メッセージ  $m$  を “time stamp” するための論理クロックを用いて生成される。この不完全な順番は、後に、つながりを調停するための送信プロセスの ID 上の lexicographical order を用いて、total order に変換される。つながりとは、もし 2 つのメッセージ  $m$  と  $m'$  が同じ論理的な time stamp を持っているなら、メッセージ  $m$  は  $m'$  の前である。 ( $id(sender(m)) < id(sender(m'))$ ) この時、 $id(p)$  は、プロセス  $p$  の

IDである。プロセス  $p$  は、将来低い time stamp を運ばないメッセージ  $m'$  を受信していないことを知っているなら、1度だけあるメッセージ  $m$  を deliver することができる。

このアルゴリズムは、まだ確立されていない。なぜなら、静かなプロセス（何の反応もないプロセス）は、deliver からの他のプロセスを妨害する可能性があるからである。これを避けるためには、プロセス  $p$  が少しの時間の遅延の後、空のメッセージを broadcast すること要求される。

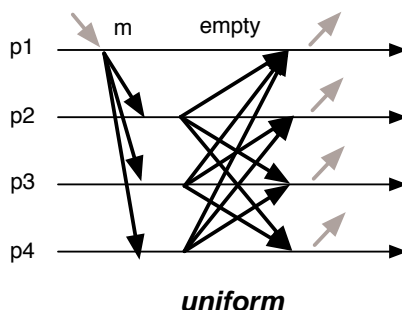


図 5.4: Communication History

### Destinations Agreement

Destinations Agreement 図 5.5 では、メッセージを deliver する順番は、下記の 2つのいずれかの方法で、宛先のプロセスによって決定される。

- 順番の情報は、決定的に結合された各プロセスによって生成される。
- 順番の情報は、例えば consensus のようなアルゴリズムを用いて、プロセス間で合意の上で得られる。

メッセージ  $m$  を broadcast するために、プロセス  $p_1$  はメッセージ  $m$  をすべてのプロセスに送ることおよび、メッセージ  $m$  の delivery のために coordinator（調整役）として振る舞う。他のプロセスは、メッセージ  $m$  を受信すると、各プロセスは、ack と logical time stamp を送信元のプロセス  $p_1$  へ送信する。プロセス  $p_1$  はすべての time stamp を収集すると、すべての受け取った time stamp から、最大値を計算し、timestamp  $TS(m)$  を決める。最後に、プロセス  $p_1$  は、 $TS(m)$  をすべてのプロセスに送信し、各プロセスはメッセージ  $m$  を  $TS(m)$  に従って deliver する。



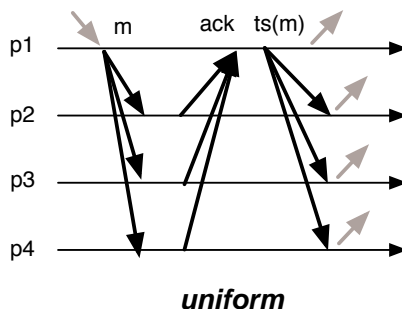


図 5.5: Destinations Agreement

## 5.2 Evaluating Neko + SSFNet with Two Consensus Algorithm

Atomic Broadcast を用いた二つの異なる Consensus のアルゴリズムについて、比較および性能評価を行う。この性能評価は、Urban らによって、[5] で、Consensus のアルゴリズムを Neko に実装し行われている。本研究で行う性能評価と Urban[5] で行った性能評価の違いは、シミュレートを行う環境であるシミュレートネットワークに SSFNet を用いることである。そして、Urban ら [5] による評価結果での評価結果と、本研究での評価結果を比較することで Neko の有効性を示す。

分散合意問題、すなわち、Consensus、Atomic Broadcast 等は、トランザクションや時間に制約があるようなアプリケーションを含んだ分散システムの耐故障性のために必要不可欠とされている。Unbun らによる [5] では、2つの知られた Consensus のアルゴリズムの比較を提供している。1つ目のアルゴリズムは、あるプロセスだけに通信が集中するようなパターンを用いており、2つ目のアルゴリズムは、通信が集中しないパターン、つまり、各プロセスに均等に通信が行われるようなアルゴリズムである。他の特筆すべき特徴は、この2つのアルゴリズムとも非常に似ている。この性能評価では、Consensus アルゴリズムの性能評価のための一般的な方法論も提案している。その方法論について簡単に説明する。この2つの Consensus アルゴリズムは、各プロセスが他のプロセスへ Atomic Broadcast を送信するようなシステムで分析されている。Atomic Broadcast アルゴリズムは、メッセージの delivery する順番を決定するために、Consensus の順番付けを実行する。Atomic Broadcast の性能評価を行うことは、現実的なシナリオの使い方として、Consensus アルゴリズムの性能評価を理解するの非常に良い方法である。この性能評価で用いられている Atomic Broadcast は、2つの Consensus アルゴリズムの内のどちらか一方を用いている。このアルゴリズムを比較するために、Neko を用いて多様な異なるシミュレートネットワークを用いている。ネットワークモデルは、Contention aware metrics [6] と呼ばれるモデルを用いている。このモデルは、メッセージを交換するために、計算機およびネットワーク上で複数の計算機（プロセス）が同じリソースを用いており、Neko のシミュレー

トネットワークの MetricNetwork として実装されている。Atomic Broadcast のための性能測定基準は、*early latency* である。この時間は、メッセージ  $m$  の送信と最も早いメッセージ  $m$  の delivery の間の経過である。Urban ら [5] は、(1) failuers および suspicions も起きない環境で実行した状態の latency の評価と、(2) プロセスがクラッシュした後の一時的な latency の評価の二つを行っている。本研究では、(1) に限定して、SSFNet を用いた評価と MetricNetwork を用いた評価との比較を行う。

### 5.2.1 定義

性能評価で用いるシステムモデルは、分散システムのための広く確立されたモデルである。そのモデルは、メッセージ交換をだけを用いて通信を行うプロセスから成り立っている。システムは、非同期で通信を行うタイミングの振る舞いなどの前提条件はない。ネットワークは、メッセージを失わないし、また、重複したメッセージを作らないといったような外見上の信頼性がある。慣例的に、これは、失ったメッセージを再送することで簡単に達成できる。

次に、分散合意問題の非公式な定義について説明する。Consensus の問題において、各プロセスは、initial 値を提案する。Uniform Consensus は、提案された値から、すべてのプロセスが同じ値を決定することを続けて行う。

Atomic broadcast は、5.1 章で説明した通り、 $A\text{-broadcast}(m)$  と  $A\text{-deliver}(m)$  の 2 つの基本的な操作により定義される。Uniform Atomic Broadcast は、2 つのことを保証する。それは、(1) もし、あるメッセージがあるプロセスによって  $A\text{-broadcast}$  されるなら、すべての正確なプロセスは、いつかはそのメッセージを  $A\text{-deliver}$  することと、(2) すべてのプロセスは、同じ順番でメッセージを  $A\text{-deliver}$  することである。

### 5.2.2 2つの Consensus Algorithms

2 つの Consensus アルゴリズムについて、それらの共通のポイントと異なるポイントを解説する。そして、Consensus の上に構築される Atomic Broadcast のついで紹介する。

#### Consensus アルゴリズム

Consensus を解決するために、Chandra-Toueg  $\diamond S$  アルゴリズム [7] と Mostefaoui-Raynalno の  $\diamond S$  のアルゴリズム [8] を用いる。これからは、前者を CT アルゴリズムと後者を MR アルゴリズムとしてそのアルゴリズムを呼ぶ。

#### 共通のポイント

このアルゴリズムは、多くの前提条件と特性を共有している。両方のアルゴリズムは、 $\diamond S$  failure detectors を用いた非同期モデル向けに設計されている。両方とも、交替制のコーディネータ（調整役）法に基づいており、各プロセスは非同期なラウンドのシーケンスを実行する。そして、各ラウンドにおいて、プロセスは、コーディネータの役割を取得する。そのコーディネータの役割は、すべてのプロセス上で決定的な値を押し付けること

である。もし、それが成功すると、そのコーディネータは終了する。このケースでは、新しいラウンドから開始される。

### ラウンドの実行と異なるポイント

Consensus の各ラウンドが実行すると、CT アルゴリズムは、集中的な通信の問い合わせを用いる。(図 5.6 中の左) 一方、MR アルゴリズムは、集中的な通信の問い合わせを用いない。(図 5.6 中の右) この 2 つの図は、各 2 つのアルゴリズムの 1 つのラウンドを実行を表している。

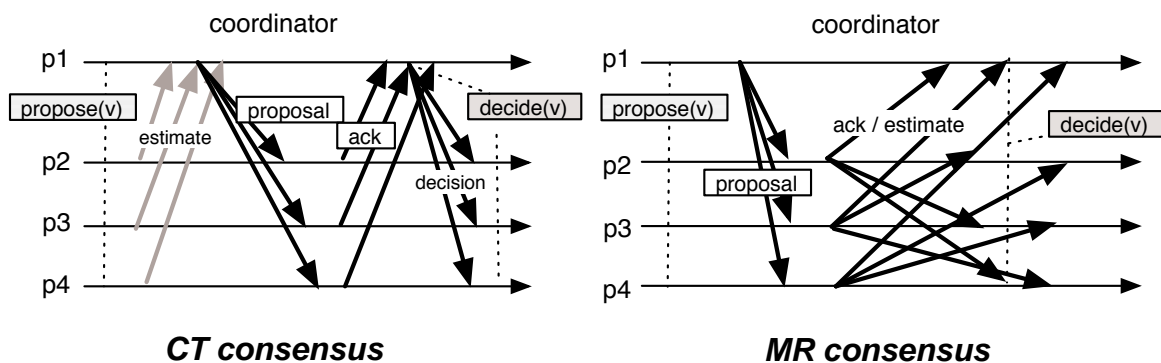


図 5.6: Example run of the both consensus algorithm

- CT アルゴリズムでは、コーディネータは、はじめにプロセスの多数決からの決定値のために estimates メッセージを集めている。このフェーズは 2 番目のラウンドとその後に一意必要なことである。
- 両方のアルゴリズムにおいて、コーディネータは proposal メッセージを送信する。
- その proposal が受信されると、プロセスは、acknowledgment(ack メッセージ)を送信する。CT アルゴリズムでは、ack メッセージはコーディネータだけに送信される。MR アルゴリズムでは、ack メッセージはすべてのプロセスに送信される。それに加えて、MR アルゴリズムのプロセスは、ack メッセージの上にプロセスの最新の estimate メッセージを便乗させる。順番に proposal を選択するための次のラウンドのコーディネータを認める。これは、MR アルゴリズムが estimate メッセージを送信するための独立した行程を要求していないからである。似た方法の便乗する estimate は、次のラウンドのコーディネータが ack メッセージを受信しないので、CT アルゴリズムでは不可能である。
- プロセスの多数決からの ack を受信すると、そのコーディネータ (CT アルゴリズム) とすべてのプロセス (MR アルゴリズム) は決心する。CT アルゴリズムのコーディネータは、その決定をすべてのプロセスに送信することを必要としている。(decision

メッセージ) これは、MT アルゴリズムでは必要ない。なぜなら、各プロセスが独立に決定しているからである。

### The Chandra-Toueg atomic broadcast algorithm

このアルゴリズムは、プロセスがすべてのプロセスに送信したメッセージによって、A-broadcast を実行する。あるプロセスが各メッセージを受信する時に、そのプロセスは delivery する順番が決定されるまで、そのメッセージをバッファにためておく。その delivery する順番は、Consensus のシーケンスが番号付け (1、2、etc) することにより決定される。その値は、はじめに提案され、各 Consensus の決定値はメッセージ ID のシーケンスである。その delivery する順番は、Consensus(1、2、etc) のようなシーケンスの連鎖によって与えられる。

このアルゴリズムは、Consensus のアルゴリズムを理解することから、システムモデルと耐故障性の保証を引き継がれていく。この性能評価では、両方の CT および MR Consensus アルゴリズムで実現された Atomic Broadcast を用いる。

## 5.2.3 ベンチマーク

この章では、性能評価基準、および workload(仕事量) から成るベンチマークについて説明する。順序よく意義のある結果を得るために、評価 (Consensus) に基づいてのコンポーネントの用語よりもむしろ、評価 (プロセスが Atomic Broadcast を送信する) に基づいてのシステムの用語を述べる。

### 性能評価基準

主なる性能評価基準は、Atomic Broadcast の early latency である。Early latency  $L$  は、以下のように単体の Atomic Broadcast のために定義される。A-broadcast(m) がプロセス  $p_j$  ある時間  $t_0$  と、A-deliver(m) が  $p_i$  により時間  $t_i$  で起こったと仮定する。( $i = 1, \dots, n$ ) すると、latency は、時間が最初のメッセージ m の A-deliver (例、 $L = (\min_{i=1, \dots, n} t_i) - t_0$ ) までの経過した時間である。この性能評価では、 $L$  が多くのメッセージを超えることと、いくつかの実行を計算している。

### Workload:仕事量

latency は、いつも一定の仕事量に基づいて計測されている。この性能評価では、単純な仕事量を選択している。それは、(1) すべての宛先プロセスは同じ一定の割合で Atomic Broadcast メッセージを送信する、および、(2) A-broadcast イベントは、stochastic process から来ることである。Atomic Broadcast の全体の割合のことを *throughput* と呼び、 $T$  がそれを示す。一般的に *throughput*  $T$  に依存して、どれくらいの latency  $L$  が決定される。

## 5.2.4 シミュレーションモデルの比較

性能評価への研究方法は、実システム上での測定を得るために適しているような一般的な結果を取得可能なシミュレーションである。それは、異なった環境の変化をシミュレートするためにネットワークモデルのパラメータを変更できることである。性能評価では、Neko を用いて分散アルゴリズムが実装されたおり、Neko のシミュレートネットワークとして実装された MetricNetwork を用いている。本研究では、さらにシミュレートネットワークとして SSFNet を用いて性能評価を行う。一方、SSFNet は、現実的なネットワークモデルを定義することが可能である。この2つのシミュレートネットワークは、共通する特性を持っており、ネットワークモデルを比較することで SSFNet の有効性を示すことができる。

### MetricNetwork model in Neko

Neko のシミュレートネットワークであり、性能評価に用いられている MetricNetwork (図 5.7) について解説する。MetricNetwork は、実世界における Ethernet を単純にモデル化したものである。モデル中の要所となる点は、リソースを複数のプロセスで共有し計算していることである。この点は、分散アルゴリズムの性能のための限定的な要素であるリソースの複数のプロセスでの共有として重要である。計算機とネットワーク自身の両方ともは、ボトルネックになる可能性もある。モデル中でのリソースの上限の種類は、メッセージの伝達手段を提供するプロセス間で共有されているネットワークリソースと、ネットワークングをするためのレイヤースタックとネットワークの制御によるプロセスの処理を提供するプロセスごとの CPU リソースである。あるメッセージ  $m$  はプロセス  $p_i$  からプロセス  $p_j$  へ転送されることは、(1) CPU  $i$ 、(2) ネットワーク、および (3) CPU  $j$  の順番にリソースを用いている。メッセージ  $m$  は、各ステージがもし混雑しているなら、メッセージキューで待つ次のステージに押し出される。その時間は、ネットワークリソース上で 1 time unit 費やされる。また、CPU リソース上で  $\lambda$  time unit 費やされ、これは、メッセージを送信と受信することが単純なコストを持っていると考えることできる。

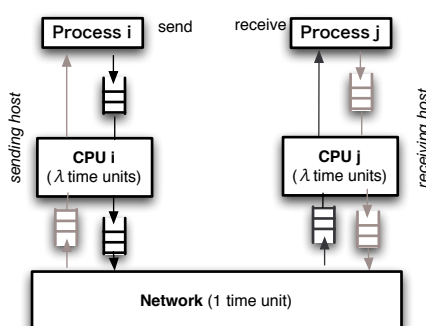


図 5.7: MetricNetwork in Neko

その $\lambda$ パラメータ ( $0 \leq \lambda$ ) は、ネットワーク上でメッセージが転送されることを計算機が比較する上でのメッセージの処理速度の関連を表している。異なった値は、異なったネットワーク環境をモデル化する。この性能評価では、 $\lambda$ のために設定の種類を用いての実験を考慮している。さらに、このネットワーク環境はユニキャストとマルチキャストによるメッセージの通信も考慮している。

### Network model in SSFNet

SSFNet を用いてモデル化することができるネットワークについて解説する。(図 5.8) SSFNet は、現実的なネットワークモデルを構築可能である。protocol は、実世界に存在する SSFNet が実装したプロトコル (e.g TCP/IP、UDP/IP etc) 用いることができ、そのプロトコルを用いてプロトコルスタックを構成することが可能である。SSFNet の各計算機 (Host、Router) は、NIC を持っておりパラメータとして、buffer 値や送出するパケットあたりの latency(送出するまでの遅延) を各計算機ごとに設定可能である。また、計算機同士を接続している link の送信遅延である delay と link のバンド幅を各リンクごとに設定可能である。そして、10000 台以上の計算機を同時にシミュレーション可能であり、そのネットワークトポロジの開発者が自由に定義可能である。

Neko の MetricNetwork は、単純な Ethernet をモデル化しているので、この SSFNet のモデルと似ている。それは、Neko における  $\lambda$  の値は、SSFNet における latency であり、Neko におけるネットワーク遅延は、SSFNet で定義できる link delay とバンド幅に対応することができる。

SSFNet のモデルにおいて、 $host_i$  から  $host_j$  までのパケットの伝播流れを説明する。 $host_i$  は、定義されたプロトコルに従って、自分自身の NIC を用いてパケットの送信を試みる。この時、もし NIC の buffer が混雑しているならば、そのパケットは、混雑が解消するまで待って、パケットを送出する。送出する時、NIC 上に定義された latency に従って、そのパケットの送出を送らせる。次に、接続されている Router へ、link に定義されたバンド幅に従って経過時間が決定され、定義された delay $l$  によって、経過時間に遅延を足した時間後、その link を通じて伝播される。Router の Host と同じように buffer と latency が定義されているので、その定義に従った経過時間がシミュレートされる。その後、同じように link、Host の NIC の順番で伝播され、プロトコルに従って、 $host_j$  へパケットが送信される。

Neko のシミュレートネットワークとして SSFNet を用いる場合、メッセージのサイズを定義することが可能であるが、SSFNet では、パケットに分割してそのメッセージは、目的の Host へ送信される。つまり、上記のような順序のシミュレートを SSFNet は、分割されたパケット毎に行う。これは、実ネットワークを精密にモデル化していることにより行われ、Neko のシミュレートネットワークである MetricNetwork よりも、現実的かつ有効なシミュレート結果が得られる。

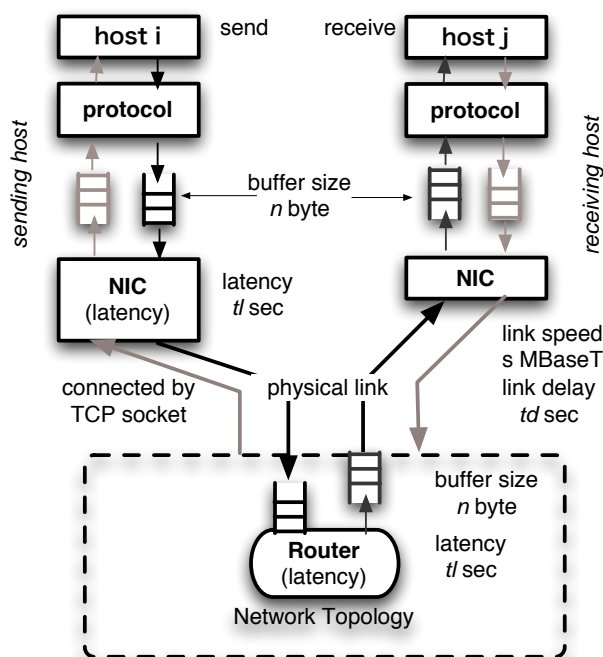


図 5.8: SSFNetNetwork in SSFNet

### 5.3 評価結果

この2つのネットワークモデルを用いた性能評価の結果を報告する。NekoのMetricNetworkで $\lambda$ : 0.1、1、10のように設定し、プロセス数 $n=7$ を用いてシミュレーションを行った結果と、SSFNetのSSFNetTCPNetworkで、すべての $latency\ of\ NIC = 100[\mu sec]$ 、 $link\ delay = 1msec$ 、および、 $link\ speed = 100MBaseT$ と、各計算機上のNICのBufferを8000byte、ルータのBufferを16000byteに設定して、プロセス数 $n=7$ を用いて、2種類のネットワークモデルで結果を取得した。SSFNetのネットワークポロジは、7つの計算機がルータに接続している単純なLocal Area Networkを定義している。この各パラメータは、4章で説明したDMLファイルのパラメータを用いている。図5.9は、MetricNetworkを用いた3つの異なった $\lambda$ による結果で、図5.10は、SSFNetTCPNetworkを用いた結果である。各図は、latency対Throughputのグラフを示している。

#### The case of using MetricNetwork

- CTアルゴリズムは、 $\lambda = 10$ の場合、最も悪いパフォーマンスである。(図5.9)この理由は、CTアルゴリズムがMTアルゴリズムよりも多くの負荷をコーディネータへ与えているからである。その要因は、CTアルゴリズムは、proposalメッセージとackメッセージを収集し、さらに、決定値を図5.6の左記のように送信しなければならないためである。
- MTアルゴリズムは、 $\lambda = 1$ の場合、最も悪いパフォーマンスである。(図5.9)パ

パフォーマンスの相違点は、throughput がおおよそ比例していることである。この振る舞いの理由は、CT アルゴリズムとは違って、CPU リソースの負荷が重要なことではないからである。従って、結論付けられる事実は、MR アルゴリズムはネットワークに負荷を多く与える。throughput の増大は、ネットワークバッファのモデルにおける高負荷な割り込み時間を導いている。

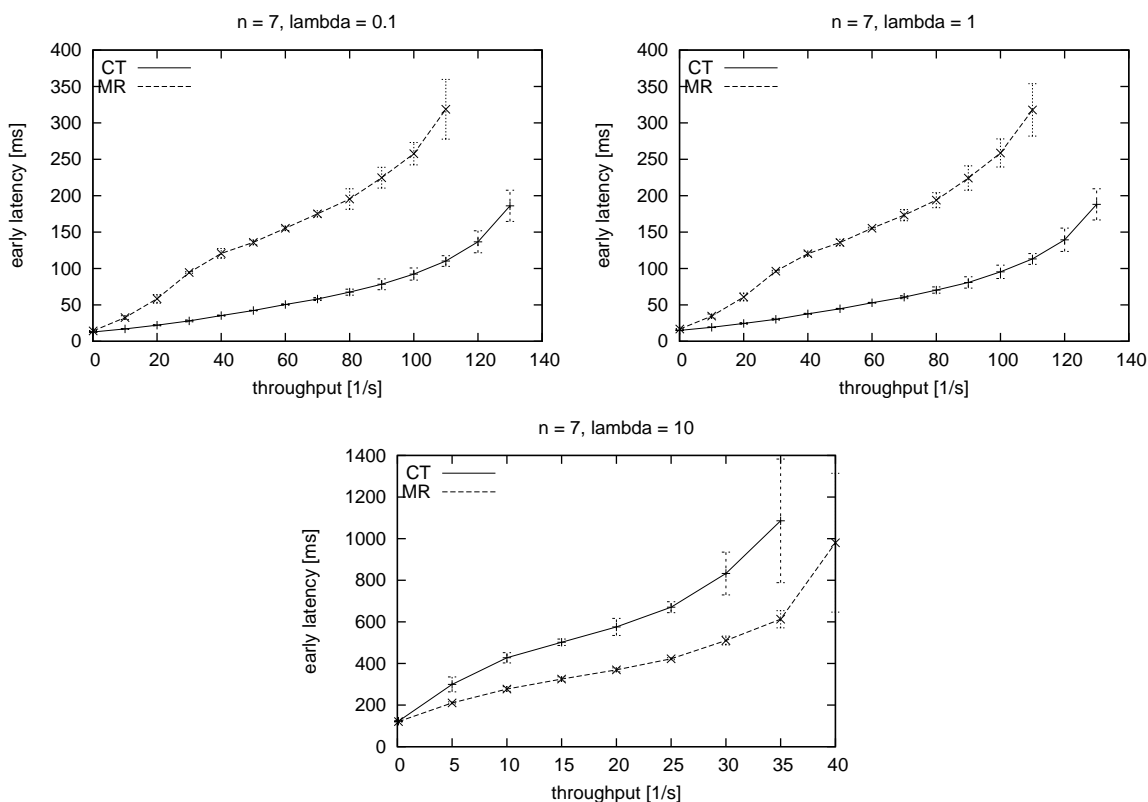


図 5.9: Result by using MetricNetwork in Neko ( $\lambda = 0.1, 1, 10$ )

### The case of using SSFNetTCPNetwork

- SSFNetTCPNetwork を用いたネットワークモデルにおいて、CT アルゴリズムは、MT アルゴリズムよりも、悪い性能である。(図5.10) このグラフは、MetricNetwork を用いた事例における  $\lambda = 10$  のグラフと非常によく似ている。MR アルゴリズムは、MetricNetwork と同じでおおよそ比例して throughput が増加している。CT アルゴリズムは、MetricNetwork の時と非常に良く似た曲線を描いている。同じような結果が得られた理由は、SSFNet で定義したネットワークモデルが MetricNetwork の  $\lambda = 10$  の場合と同じようなモデルといえるからである。SSFNet のネットワークモデルの *latency of nic* は、MetricNetwork の  $\lambda$  に当たると考えられるので、CT ア



ルゴリズムのコーディネータにおけるメッセージの送信がボトルネックとなっていることがわかる。

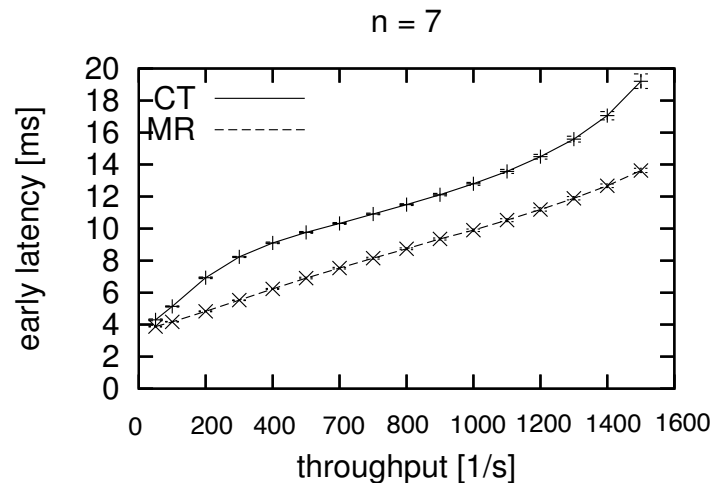


図 5.10: Result by using SSFNetTCPNetwork in SSFNet

## 第6章 まとめ

本研究では、Nekoの問題点である現実的なシミュレートネットワーク環境を持っていない点をSSFNetのシミュレーションネットワークを利用することで解決した。NekoとSSFNetの概略と特徴について、詳細に説明を行った。具体的な方法は、Neko側とSSFNet側の両方に、インタフェースだけを作成し、NekoとSSFNetの中核の実装を変更することなく統合した。実装上の問題点として、NekoとSSFNetのスケジューラの同期、アドレス方式の違い、メッセージ方式の違いおよびコンフィグファイルの仕様変更などがあつた。これらの問題を解決し実装したシミュレートネットワークとしてSSFNetが利用可能なNekoの性能評価を分散アルゴリズムのひとつであるAtomic Broadcastを用いて行った。その際、Atomic Broadcastの定義について詳細に解説し、性能評価で用いたAtomic Broadcastを実現するための分散アルゴリズムである2つのConsensusについても解説した。また、性能評価において、NekoのシミュレートネットワークであるMetricNetworkのネットワークモデルについて解説し、SSFNetのネットワークモデルとの比較を行うことで*Neko with SSFNet*の有効性を示した。

### 6.1 謝辞

In the first, I'm most grateful to Xavier Défago. I took pleasure in our steady on our laboratory which gave me wonderful environment. I would like to say "Arigato" to Péter Urbán who is collaborate researcher. Your Neko is towering tool and gave me most enjoyment. Lastly, I was most interested life with partners in Xavier Lab.

## 参考文献

- [1] Peter Urban, Xavier Defago, and Andre Schiper, “Neko: A single environment to simulate and prototype distributed algorithms.” *Journal of Information Science and Engineering*, 18(6):981-997, November 2002.
- [2] J.Cowie,D.M.Micol,A.T.Ogielski, “Towards Realistic Million-Node Internet Simulations.” *Proc Int Conf on Parallel Distributed Processing Techniques Applications (PDPTA'99)*, June, 1999.
- [3] X.Defago, A.Schiper, and P.Urban, “Comparative performance analysis of ordering strategies in atomic broadcast algorithms.” *IEICE Trans. on Information and Systems*, Vol.E86-D, No.12, pp.2698-2709, December 2003.
- [4] X.Defago, A.Schiper, and P.Urban, “Total order broadcast and multicast algorithms: Taxonomy and survey. ” *Research Report IS-RR-2003-009*, Japan Advanced Institute of Science and Technology, Ishikawa, Japan, September 2003.
- [5] P.Urban, N.Hayashibara, A.Schiper, and T.Katayama, “Performance Comparison of a Rotating Coordinator and a Leader Based Consensus Algorithm.” In *Proc, 23rd IEEE Int'l Symp, on Reliable Distributed System (SRDS)*, page 4-17, Florianopolis, Brazil, October 2004.
- [6] P.Urban, X.Defago, and A.Schiper, “Contention-aware-metrics for distributed algorithms: Comparison of atomic broadcast algorithms.” *Proc. 9th IEEE Int'l Conf. on Computer Communications and Networks (IC3N 2000)*, 2000, page 582-589.
- [7] T.D. Chandra and S.Toueg, “Unreliable failure detectors for reliable distributed systems.” *Journal of the ACM*, 43(2), 1996, 225-267.
- [8] A.Mostefaui and M.Raynal, “Solving consensus using Chandra-Toueg’s unreliable failure detectors: A general quorum-based approach.” *Proc. 13th Int'l Symp. on Distributed Computing (DISC)*, Bratislabe, Slovak Republic, 1999, 49-63.
- [9] L.Lamport, “Time, clocks, and the ordering of events in a distributed system.” *Commun. ACM*, vol.21,no.7,pp.558-565, 1978