

分散 Database Jungle に関する研究

A Study of
distributed Database Jungle

平成25年度 学位論文(修士)



琉球大学大学院 理工学研究科
情報工学専攻

大城 信康

要旨

スマートフォンやタブレット端末の普及により、ウェブサービスの利用者が増加している。それに伴い、サーバサイド側への負荷も増大しウェブサービスがダウンする事態が出てきている。そのため、スケーラビリティはウェブサービスにおいて重要な要件の1つとなっている。スケーラビリティとは、ある複数のノードから構成される分散ソフトウェアがあるとき、その分散ソフトウェアに対して単純にノードを追加するだけで性能を線形に上昇させることができる性質である。ウェブサービスにスケーラビリティを持たせるには、まずデータベースがスケーラビリティを持たなければならない。

そこで、当研究室ではスケーラビリティのあるデータベースとして、非破壊的木構造を用いたデータベース“Jungle”を設計・開発している。非破壊的木構造とは、データの編集の際に一度木構造として保存したデータを変更せず、新しく木構造を作成してデータの編集を行うことを言う。

本研究では、Jungle に分散データベースと永続性の実装を行った。データ分散部分には当研究室で開発中である並列分散フレームワークである Alice を使用した。学科の並列環境を用いて複数のサーバノード間でデータの分散を行うことを確認した。また、例題アプリケーションとして簡易掲示板プログラムの作成を行った。Jungle と Cassandra により作成した掲示板プログラムに対して読み込みと書き込みの負荷をかけ比較を行った。結果、Cassandra 以上の性能を確認することができた。残された課題としてデータ分割の実装、Merge アルゴリズムの設計、分断耐性の実装、過去のデータの掃除についての課題が確認された。

Abstract

Smartphone and tablet pc are widely used, thereby users of the Web services is increasing. It has caused the webservice is down. Therefore, scalability is important software factor today. Scalability in distributed system is able to increase performance linearly when just added new node to system. In order to make provide scalability to Web services, database must have scalability.

We are designing and developing a database Jungle for study of scalable database. It is use non-destructive tree structure. Non-destructive tree structure is not the destruction of data. Editing of data is done creating by new tree.

Jungle was designed as a distributed database. But data distribution and persistent has not yet been implemented in the Jungle.

In this paper, we develope persistent and distributed database on jungle. Distributed data on Jungle is developing using parallel distributed framework Alice. We implemented distributed on Jungle and confirmed with our cluster system that data is distributed between the server node. Also, we developed simple bulletinboard system with Jungle and key-value store database Cassandra. We compared Jungle and Cassandra using simple bulletinboard. As a result, we got better performance then Cassandra.

目次

第 1 章	序論	1
1.1	本論文の構成	3
第 2 章	既存のデータベース	4
2.1	Relational Database	4
2.2	NoSQL データベース	4
2.3	CAP 定理	5
2.4	既存の NoSQL データベース	6
2.4.1	memcached	6
2.4.2	MongoDB	7
2.4.3	Neo4j	8
2.4.4	Cassandra	9
第 3 章	木構造データベース Jungle	11
3.0.5	破壊的木構造	11
3.0.6	非破壊的木構造	12
3.1	Jungle におけるデータへのアクセス	15
3.2	Jungle におけるデータ編集	15
3.2.1	NodeOperation	15
3.2.2	TreeOperationLog	16
第 4 章	分散データベース Jungle の設計	18
4.1	分散バージョン管理システムによるデータの分散	18
4.1.1	分散管理システムの API	18
4.1.2	Merge によるデータ変更衝突の解決	19
4.2	Jungle のネットワークトポロジーの形成	21
4.2.1	ツリートポロジーの形成とルーティング	21
4.2.2	トポロジーの形成手段	21
4.3	並列分散フレームワーク Alice	22
4.3.1	MessagePack によるシリアライズ	23
4.4	Jungle のデータ分散	23
4.5	ログによるデータの永続性	23

第 5 章	木構造データベース Jungle の分散実装	24
5.1	Alice のトポロジーマネージャの利用	24
5.1.1	トポロジーマネージャの起動	24
5.1.2	アプリケーション側の記述	25
5.2	Alice を用いての分散実装	26
5.2.1	Alice によるプログラミング	26
5.2.2	他サーバノードの DataSegment へアクセス	28
5.2.3	独自クラスのインスタンスの送受信	29
5.3	ログのシリアルライズ	30
5.3.1	TreeOperationLog のシリアルライズ	30
5.4	ログに対する情報の追加	30
5.4.1	NetworkTreeOperationLog の実装	30
5.4.2	ログの送信部分	31
5.4.3	他サーバノードへのログの送信	31
5.4.4	ログの受信とデータ反映	32
5.5	永続性の実装	34
5.6	Merge の実装	35
5.7	掲示板プログラムにおけるデータ衝突	35
5.7.1	掲示板プログラムにおける Merge	37
第 6 章	分散木構造データベース Jungle の評価	38
6.1	実験方法	38
6.1.1	Torque Resource Manager	39
6.1.2	weighttp	39
6.1.3	掲示板プログラム	39
6.1.4	実験環境	40
6.2	実験結果 1	42
6.2.1	実験結果 1 の考察	43
6.3	実験結果 2	44
6.3.1	実験結果 2 の考察	45
第 7 章	結論	46
7.1	まとめ	46
7.2	今後の課題	46
7.2.1	pull/push 方式による分断耐性の実装	46
7.2.2	Merger アルゴリズムの設計	47
7.2.3	過去のデータの掃除について	47
	謝辞	48
	参考文献	49

発表文献

51

目 次

1.1	DEOS プロセス	1
2.1	memcached のデータ分散	6
2.2	Sharding	7
2.3	マスターとスレーブによるクラスタ	8
2.4	Consistency hashing による ring 型トポロジーの形成	9
2.5	ConsistencyLevel QUORUM による書き込み	10
2.6	ConsistencyLevel QUORUM による読み込み	10
3.1	破壊的木構造の編集	11
3.2	非破壊的木構造の編集	12
3.3	非破壊的木構造の編集手順 1	13
3.4	非破壊的木構造の編集手順 2	13
3.5	非破壊的木構造の編集手順 3	13
3.6	非破壊的木構造の編集手順 4	14
3.7	非破壊的木構造による利点	14
3.8	Node の attribute と NodePath	15
3.9	TreeOperationLog の具体例	17
4.1	分散管理システム	19
4.2	衝突の発生しないデータ編集	20
4.3	自然に衝突を解決できるデータ編集	20
4.4	衝突が発生するデータ編集	20
4.5	ツリー型の Network Topology	21
4.6	リング型のトポロジー	22
4.7	メッシュ型のトポロジー	22
4.8	DataSegment と CodeSegment によるプログラムの流れ	22
5.1	Alice によるネットワークトポロジー形成	25
5.2	DataSegment と CodeSegment によるプログラムの例	28
5.3	トポロジーの形成	28
5.4	Jungle による掲示板プログラムのデータ保持方法	35
5.5	他サーバノードの編集データ反映による整合性の崩れ 1	36
5.6	他サーバノードの編集データ反映による整合性の崩れ 2	36

6.1	実験 1 複数のクライアントからサーバ 1 台への負荷	38
6.2	実験 2 複数のクライアントから複数のノードへの負荷	39
6.3	複数のクライアントから一台への負荷	42
6.4	複数のクライアントから一台への負荷	43
6.5	分散環境下における読み込みベンチマーク結果	44
6.6	分散環境下における書き込みベンチマーク結果	45

表 目 次

6.1	簡易掲示板システムで利用した Jetty と Cassandra のバージョン	40
6.2	掲示板プログラムを実行させる VMWare クラスタの仕様 (クライアントにも利用)	40
6.3	クライアントを実行させる KVM クラスタの仕様	40
6.4	サーバノードとして利用するブレードサーバの使用	41

第1章 序論

ITシステムが巨大化していくにつれ、障害発生事例が社会に与える影響もより大きな物となる。それに伴い、ITシステムにおけるディペンダビリティへの注目が増している。

そこで、DEOS プロジェクトはITシステムにおけるディペンダビリティを担保する技術体系をまとめ、制度化、さらには事業化を目指している。DEOS プロジェクトは2006年に独立行政法人科学技術機構 (JST) はCREST プログラムの1つとして始まったプロジェクトである。DEOS プロジェクトは、変化し続ける目的や環境の中でシステムを適切に対応させ、継続的にユーザが求めるサービスを提供することができるシステムの構築法を開発することを目標としている [1]。DEOS プロジェクトではそれらの技術体系を「オープンシステムディペンダビリティ」として定義し、それをDEOS プロセスとしてまとめた (図 1.1)。

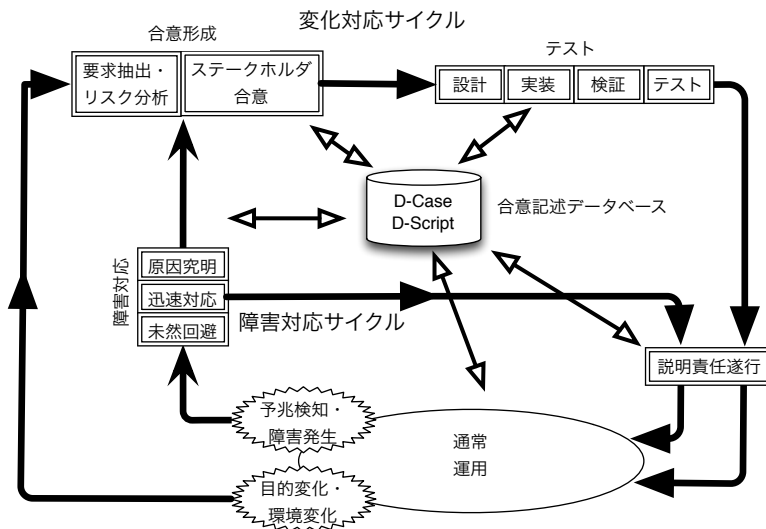


図 1.1: DEOS プロセス

DEOS プロセスには変化対応サイクルと障害対応サイクルの2つのサイクルがある。変化対応サイクルは上流プロセスにおける対象システムのオープンシステムディペンダビリティを担保するためのプロセスである。障害対応サイクルは対象システムの運用時に必要とされるプロセスである。DEOS プロセスは企画や設計という上流からシステム運用までの対象システムのライフサイクル全体に係わる。対象システムの運用は、ステークホルダからの対象システムに対する要求に関する合意を始めとしたあらゆる議論に関する合意

をベースに行われる。

DEOS プロセスを構成している各 (サブ) プロセスからアクセスが行われるデータベースとして D-ADD (DEOS Agreement Description Database) がある [2]。D-ADD はステークホルダ合意と対象システムに存在するプログラム・コード, 及び対象システムの運用状態との間の一貫性を常に保つための機構を提供する。このようなデータベースは様々なデータを柔軟に格納する必要があり, データベーススキーマの頻繁な変化に対応する必要がある。これらのデータベースは, ウェブからアクセスされることも想定される。そのため, DEOS はウェブサービスとして捉えることができる。

ウェブサービスにとってデータベースは必須であり, ウェブサービスの規模に比例してデータベースへの負荷も大きなものとなっている。そのため, データベースの処理能力の高さはそのままウェブサービスの質にも繋がってくる重要な項目となっている。データベースの処理能力をスケールさせる方法としてスケールアップ (垂直スケール) とスケールアウト (水平スケール) の考えがある。スケールアップはハードウェア的に高価な単一のマシンを用意することで処理性能を上げることをさす。スケールアウトは汎用的なマシンをいくつも用意し, 処理を分散させることで処理性能を上げることをさす。単純に処理能力をあげたいのならば, スケールアップは有効である。しかし, 単一のマシンを高性能にしても限界がある。いずれはそのマシン単体が処理できない程負荷が増えるかもしれないし, ネットワーク帯域といったマシン以外のハードウェアでボトルネックが発生するかもしれない。なによりコストがかかってしまう。

それに対しスケールアウトは, 処理が重くなるにつれて汎用的なマシンを順次追加していくことで性能をあげていく。1つのマシンに対して負荷が重くなるということ为了避免するため, ハードウェア的に高性能なマシンを要求せずすみ, 柔軟な対応をとることができる。ウェブサービスに使用されるデータベースの性能をあげる方法としては, このスケールアウトが求められている。

本研究で扱うスケーラビリティとはこのスケールアウトのことをさす。そういう意味では最も使われているデータベースである Relational Database はスケーラビリティを持つのは困難である。Relational Database にはないスケーラビリティを持つデータベースとして NoSQL と呼ばれるデータベースがある。NoSQL データベースは Consistency hashing や Sharding といった方法を使いデータを分散させスケーラビリティを得ている。データベースにおいてデータを分散させスケーラビリティを上げることはもはや必須となっている。

本論文では, スケーラビリティのあるデータベースを目指して木構造データベース Jungle の提案する。すなわち, Jungle に分散と永続性の実装をう。既存の分散データベースである Cassandra との比較を行うため, 簡易掲示板を作成し並列環境から負荷をかけることで性能比較を行った。

1.1 本論文の構成

本論文では、初めに分散データベースについて既存の製品を例に挙げながら述べる。RDB と NoSQL について説明し、既存の NoSQL データベースの特徴について説明する。第 3 章では、Jungle の持つデータ構造、非破壊的木構造について説明し、Jungle におけるデータ編集について述べる。第 4 章では、非破壊的木構造による Jungle の基本設計と、分散バージョン管理システムを参考にした分散設計について述べる。第 5 章では、第 4 章で行った設計を元に Jungle へ行った分散実装について述べる。第 6 章では、第 5 章で実装した分散データベース Jungle の評価を行うため、簡易掲示板プログラムを実装する。この掲示板プログラムは Jungle と Cassandra それぞれのデータベースを使うものを用意した。学科の並列環境上で開発した掲示板プログラムを複数のノードで実行させ、負荷をかけることで Jungle と Cassandra の性能比較を行う。第 6 章では、本研究におけるまとめと今後の課題について述べる。

第2章 既存のデータベース

本章では、まずデータベースの種類である Relational Database と NoSQL について述べる。次に、分散データシステムにおいて重要な CAP 定理について触れる。最後に、既存の NoSQL データベースとして memcached, MongoDB, Neo4j, Cassandra の特徴について述べる。

2.1 Relational Database

Relational Database (RDB) は、列と行からなる 2 次元のテーブルにより実装されるデータベースである。データ型として文字列、数値、日付、BOOL 型がありシステムによりデータに型が強制される。RDB はスキーマの決まったデータを扱うことに長けている。構造化言語問い合わせ言語として SQL がある。

RDB はデータベースの中でも長年主要な立ち位置にあるデータベースだが、苦手としていることがある。それは、スキーマレスなデータの扱いやマシンの台数を増やして処理速度をあげることである。垂直分割や水平分割といった方法によりデータを分けることはできるが、分割を行うほどデータの扱いは複雑になっていく。

2.2 NoSQL データベース

NoSQL は Not Only SQL の略で、SQL を使わないデータベースのことを指す。NoSQL データベースは RDB とは違いスキーマがない。そのため、扱おうとしているデータの形が決まっていなくても気軽に使うことができる。

後述する Consistency Hashing や Sharding といった方法で複数ノードでデータの分散を行うことでスケーラビリティの確保を行う。単純なノードの追加により負荷分散を行うことができる。

2.3 CAP 定理

分散データシステムにおいては、次の 3 つを同時に保証することはできない。

- 一貫性 (Consistency) 全てのノードはクエリが同じならば同じデータを返す。
- 可用性 (Availability) あるノードに障害が発生しても、機能しているノードにより常にデータの読み書きが行える。
- 分断耐性 (Partition-tolerance) ネットワーク障害によりノードの接続が切れてもデータベースは機能し続けることができる。

これは CAP 定理 [3] と呼ばれる。利用するデータベース選ぶ場合、この CAP 定理を意識しなければならない。一貫性と可用性を重視したデータベースが RDB である。分断耐性を必要とする場合は NoSQL データベースとなる。そして NoSQL の場合、分断耐性に加えて、一貫性か可用性のどちらを重視しているかで用途が変わってくる。

分散データシステムを考える場合は、この CAP 定理を意識していなければならない。

2.4 既存の NoSQL データベース

ここでは既存の NoSQL データベースに説明する。それぞれの特徴を述べながら、どのような方法でスケーラビリティを確保しているのかについて述べる。

2.4.1 memcached

memcached は揮発性の分散型キャッシュである。Key-Value ストアとなっている。RDB とともに連携して使うことができ、その場合メモリの中にデータを保持させることでディスクへのアクセスを減らし処理性能を上げることができる。LRU (Least Recently Used) のため、メモリの容量がなくなると一番古いデータはメモリから削除されてしまう。memcached は永続性は考慮していない。

また、分散を行う機能はサーバ側に備わっておらず、クライアント側の実装に任せている。クライアント側ではノードのリストを保持している。データの読み書きの際には、クライアント側で実装されている分散アルゴリズムに従って読み書きをするノードが決定される (図 2.1)。

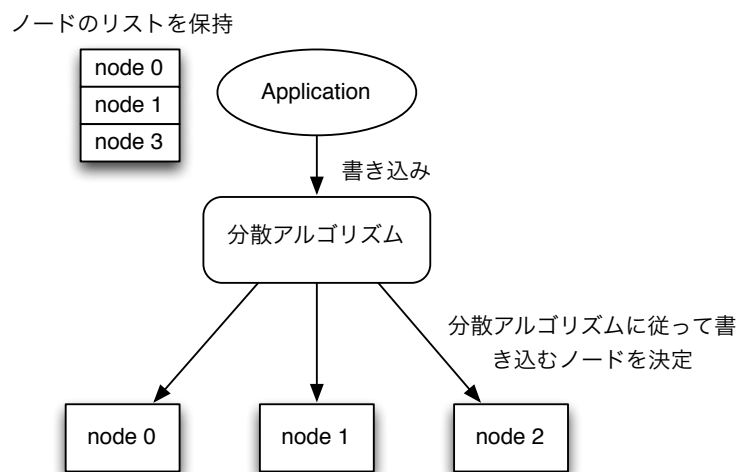


図 2.1: memcached のデータ分散

2.4.2 MongoDB

MongoDB は 2009 年に公開された NoSQL のデータベースである。JSON フォーマットのドキュメントデータベースであり、これはスキーマが無いリレーショナルテーブルに例えられる。スキーマが無いため、事前にデータの定義を行う必要がない。そのためリレーショナルデータベースに比べてデータの追加・削除が行いやすい。

MongoDB は保存したデータを複数のサーバに複製をとる。これは Replication と呼ばれる。また、1 つのサーバが全てのデータを持つのではなく、ある範囲の値を別々のサーバに分割させて保持する。これを Sharding という。MongoDB は Replication と Sharding により分断耐性と一貫性を持つ (図 2.2)。

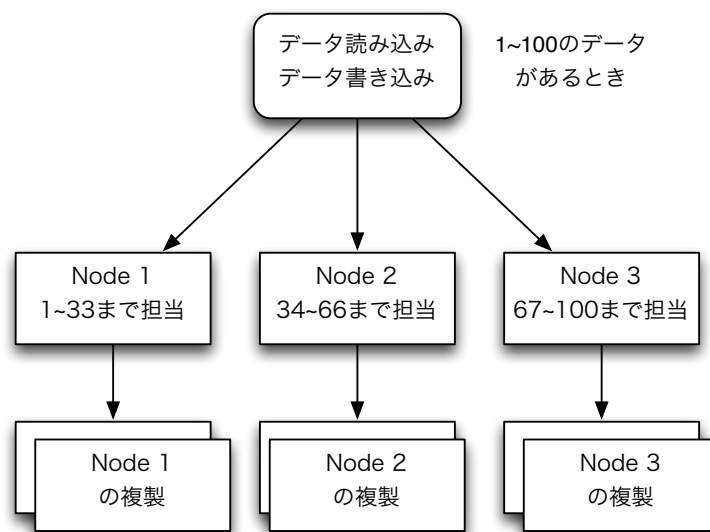


図 2.2: Sharding

2.4.3 Neo4j

Neo4j は、グラフデータベースと呼ばれる NoSQL のデータベースである。データをグラフとして保存する。グラフはノードと Relationship により表され、それぞれがプロパティを持つことができる。Relationship はグラフでいうところのエッジにあたる。ノードから Relationship を辿り、各プロパティをみることでデータの取得を行うことができる。通常データベースでは、データの取り出しに値の結合や条件の判定を行う。しかし、グラフデータベースグラフはどれだけデータが大きくなろうとノードからノードへの移動は 1 ステップですむ。そのため、どれだけデータが大きくなろうと、データが小さい時と同じ計算量でデータの取得が行える。

Neo4j はマスターとスレーブの関係になるクラスタを構成することで分散データベースとして機能する。マスターに書かれたデータはスレーブに書き込まれるが、すぐに全てのスレーブに書き込まれるわけではない。したがって、データの整合性が失われる危険がある。スレーブサーバは現在保持しているデータを返すことができる。そのため Neo4j は高い読み取り性能の要求に答えることができる可用性と分断耐性を持つ (図 2.3)。

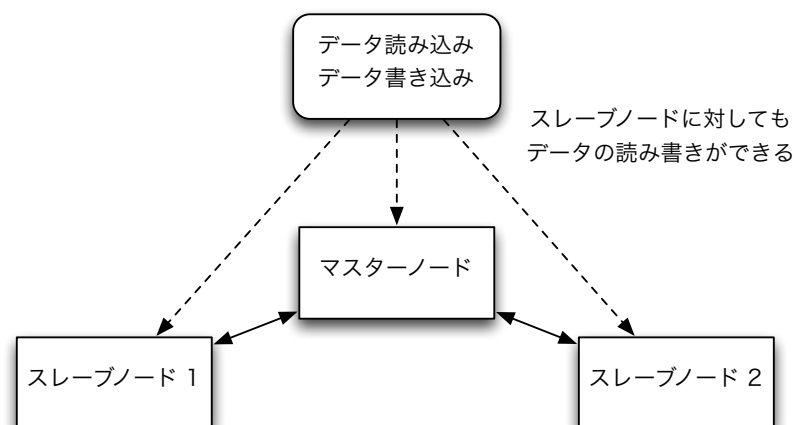


図 2.3: マスターとスレーブによるクラスタ

2.4.4 Cassandra

Cassandra[4] は 2008 年 7 月に Facebook によってオープンソースとして公開された Key-Value なデータベースである。Amazon の Dynamo[5] という分散 Key-Value データベースの影響を受けて作られている。スキーマレスな NoSQL データベースになる。

Cassandra はサーバノードの配置に Consistent hashing アルゴリズムを用いる。Consistent hashing によりノードは論理的にリング上に配置される。リングには数値で表される位置がある。データを書き込む際には、キーとなるハッシュ値に従いそのリングの位置から時計回りに近いサーバノードへと書き込まれる。Consistent hashing を用いることで、ノードの数が増減した場合に、再配置をしなくてもよいという利点がある。データの偏りにより少数のサーバへの負荷が大きい場合に、負荷が高いハッシュ値が指すリング上に新たなノードを追加することで負荷を下げるといった手段もとれる。Consistency Hashing によるリングの形成を図 2.4 に示す。

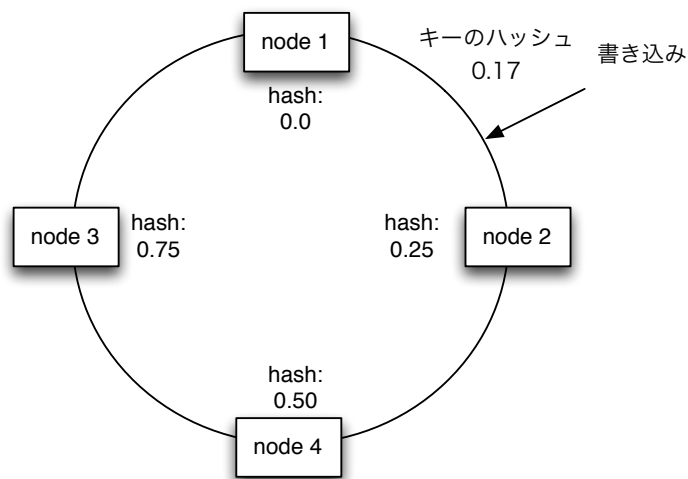


図 2.4: Consistency hashing による ring 型トポロジーの形成

Cassandra はデータを最大どれだけ配置するかを示す Replication factor と、データの読み書きをいくつのノードから行うのかを決める Consistency Level の設定が行える。Consistency Level には主に ONE, QUORUM, ALL がある。Replication factor の数値を N とした場合、ONE は 1 つのノード、QUORUM は $N/2 + 1$ のノード、ALL は N のノードへと読み書きを行う。Replication factor と Consistency Level の設定により、Cassandra は最新のデータを取得したいときとそうでないときで読み込みと書き込みの速度をあげることができる。一貫性が重要なデータに関しては QUORUM により書き込み読み込みを行うことで常に最新のデータを取得することができる。多少データが古くてもよい場合は ONE などを使用することでレスポンスを早くすることができる。ConsistencyLevel QUORUM の時のデータ書き込みと読み込みについて図 2.5 と図 2.6 に示す。Consistency

ハッシング, Replication factor と Consistency レベルの設定により Cassandra は高い可用性と分断耐性を持つ.

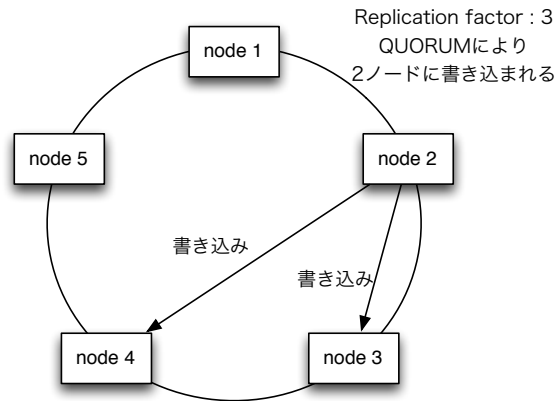


図 2.5: ConsistencyLevel QUORUM による書き込み

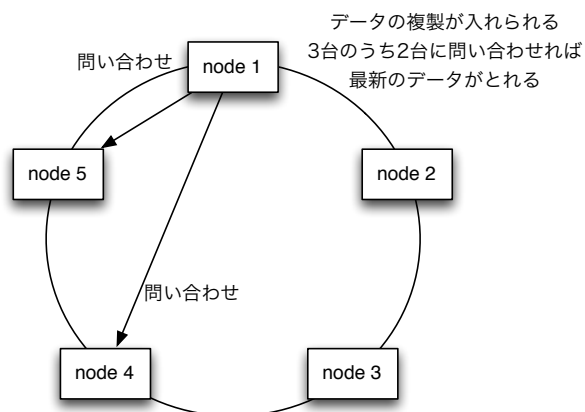


図 2.6: ConsistencyLevel QUORUM による読み込み

第3章 木構造データベースJungle

Jungle はスケーラビリティのある CMS の開発を目指して当研究室で開発されている非破壊的木構造データベースである。一般的なコンテンツマネジメントシステムではブログツールや Wiki・SNS が多く、これらのウェブサイトの構造は大体木構造であるため、データ構造として木構造を採用している。現在 Java と Haskell によりそれぞれ言語で開発されており本研究で扱うのは Java 版である。

本章ではまず破壊的木構造と、非破壊的木構造の説明をし、Jungle におけるデータ分散の設計について述べる。

3.0.5 破壊的木構造

破壊的木構造の編集は、木構造で保持しているデータを直接書き換えることで行う。図 3.1 は破壊的木構造におけるデータ編集を表している。

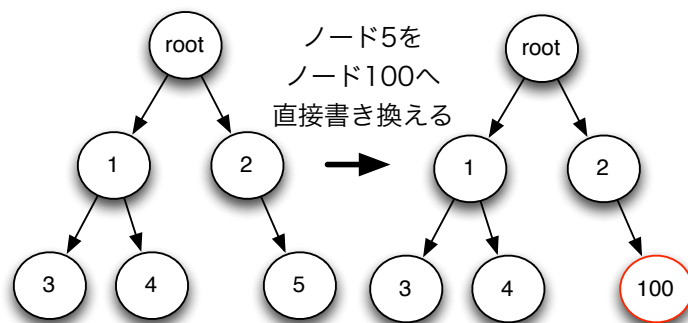


図 3.1: 破壊的木構造の編集

破壊的木構造は、編集を行う際に木のロックを掛ける必要がある。この時、データを受け取ろうと木を走査するスレッドは書き換えの終了を待つ必要があり、閲覧者がいる場合は木の走査が終わるまで書き換えをまたなければならない。これではロックによりスケーラビリティが損なわれてしまう。

3.0.6 非破壊的木構造

非破壊的木構造は破壊的木構造とは異なり、一度作成した木を破壊することはない。非破壊的木構造においてデータの編集は、ルートから編集を行うノードまでコピーを行い新しく木構造を作成することで行われる。図 3.2 は非破壊的木構造のデータ編集を示している。

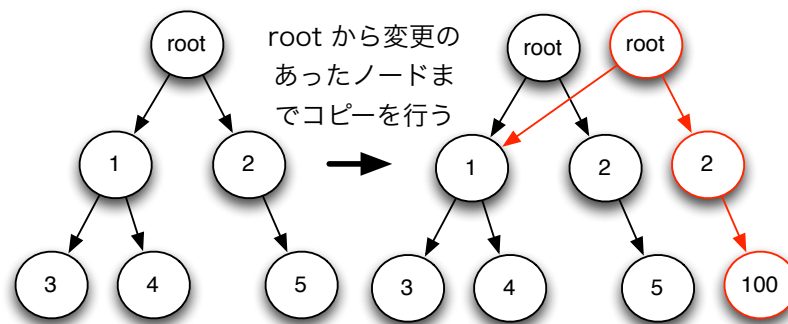


図 3.2: 非破壊的木構造の編集

非破壊的木構造におけるデータ編集の手順を以下に示す。

1. ルートから編集を行うノードまでのパスを調べる (図 3.3).
2. 編集を行うノードのコピーをとる。コピーをとったノードへデータの編集を行う (図 3.4).
3. 調べたパスに従いルートからコピーしたノードまでの間のノードのコピーをとり、繋げる (図 3.5).
4. コピーしたルートノードは編集を行っていないノードへの参照を貼り、新しい木構造を作る (図 3.6).

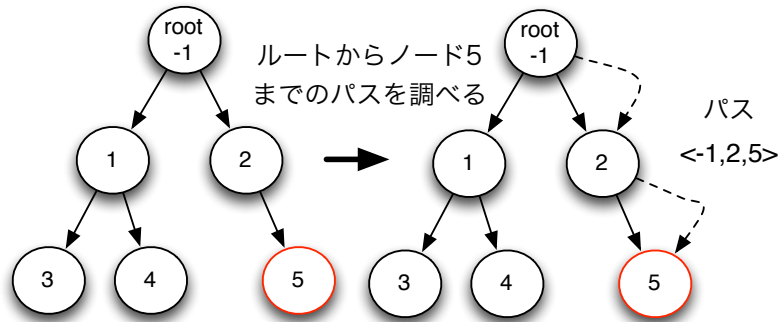


図 3.3: 非破壊的木構造の編集手順 1

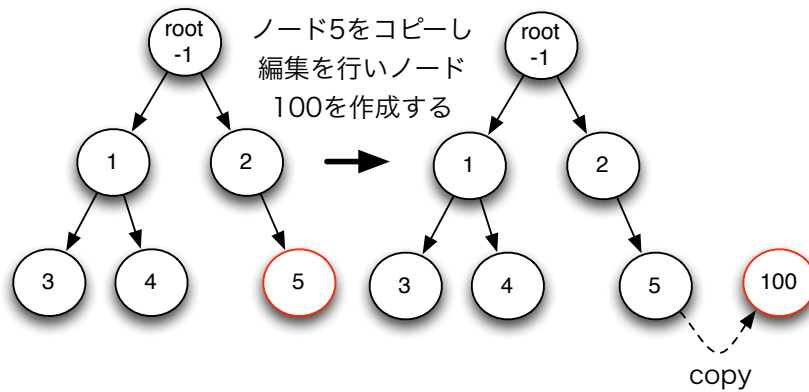


図 3.4: 非破壊的木構造の編集手順 2

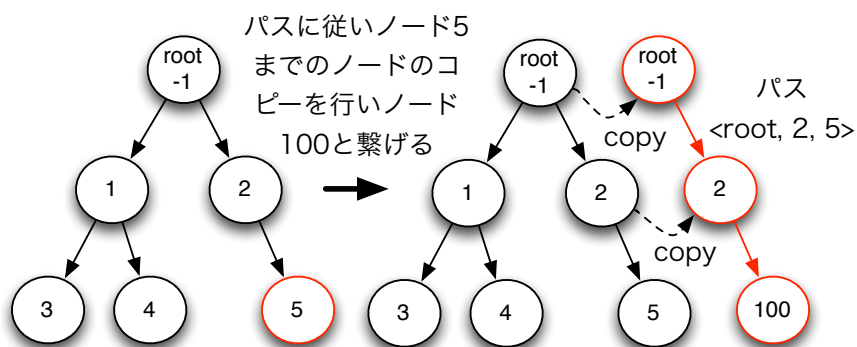


図 3.5: 非破壊的木構造の編集手順 3

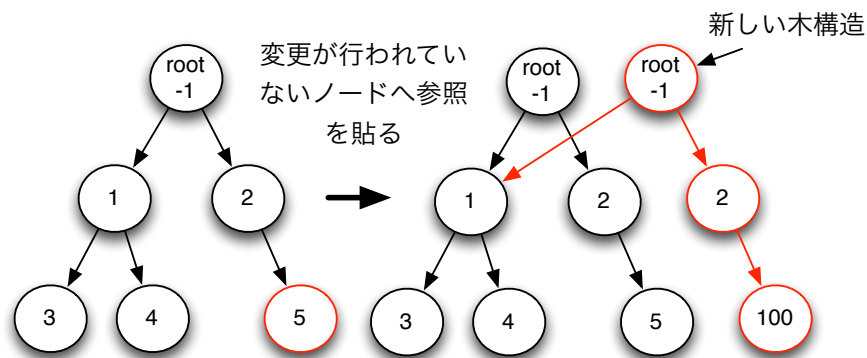


図 3.6: 非破壊的木構造の編集手順 4

非破壊的木構造においてデータのロックが必要となる部分は、木のコピーを作終りえた後にルートノードを更新するときだけである。データ編集を行っている間ロックが必要な破壊的木構造に比べ、編集集中においてもデータの読み込みが可能である (図 3.7)。そのため、破壊的木構造に比べてスケールアウトがしやすくなっている。

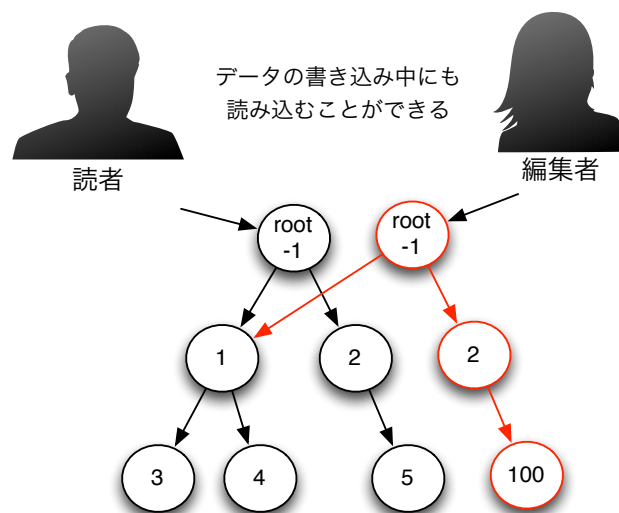


図 3.7: 非破壊的木構造による利点

3.1 Jungle におけるデータへのアクセス

Jungle におけるデータアクセス手段について述べる。Jungle ではそれぞれの Node が attribute を保持する。attribute は Key-Value によりデータを保持される。Key は String 型で Value は ByteBuffer を使用している。Jungle でデータへのアクセスは、この Node へのアクセスをさす。Node へのアクセスは、木の名前と Node を指すパスにより行える。このパスは NodePath と呼ばれる (図 3.8)。

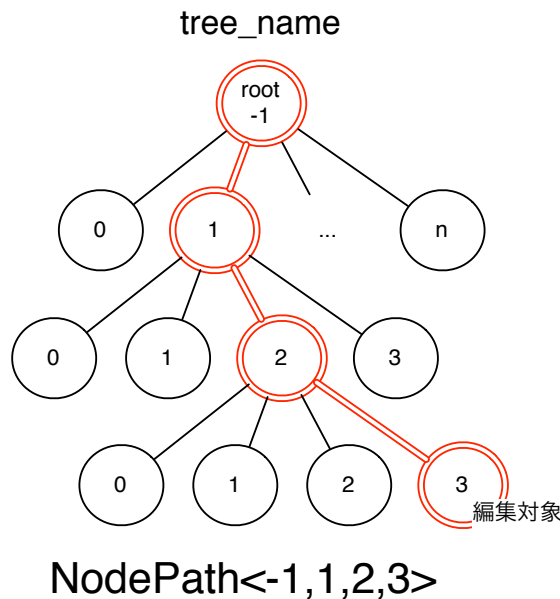


図 3.8: Node の attribute と NodePath

3.2 Jungle におけるデータ編集

3.2.1 NodeOperation

Jungle による最小のデータ編集は Node の編集を指す。Node 編集のために API が用意されており、この API は NodeOperation と呼ばれる。NodeOperation には次の 4 つの API が用意されている。

- addNewChild(NodePath _path, int _pos)
NodePath で指定された Node に子供となる Node を追加する API である。pos で指定された番号に子供として追加を行う。
- deleteChildAt(NodePath _path, int _pos)
NodePath と pos により指定される Node を削除する API である。

- `putAttribute(NodePath _path, String _key, ByteBuffer _value)`
Node に attribute を追加する API である。NodePath は attribute を追加する Node を指す。
- `deleteAttribute(NodePath _path, String _key)`
_key が示す attribute の削除を行う API である。NodePath は Node を示す。

NodeOperation は NodePath とセットで扱わなければならない、このセットを TreeOperation という。TreeOperation が 1 つのデータ編集の単位になるが、これはあくまで最小のデータ編集の単位である。Jungle によるデータの編集は TreeOperation が複数集まった単位で commit されていく。この TreeOperation の集まりを TreeOperationLog という。

3.2.2 TreeOperationLog

Jungle 内部では TreeOperation は順次ログに積まれていき、最終的に commit されることで編集が完了する。この時、ログに積まれた複数の TreeOperation は TreeOperationLog として扱われる。TreeOperationLog の仕様をソースコード 3.1 に示す。

ソースコード 3.1: TreeOperationLog の仕様

```

1 public interface TreeOperationLog extends Iterable<TreeOperation>
2 {
3     public TreeOperationLog add(NodePath _p, NodeOperation _op);
4     public TreeOperationLog append(TreeOperationLog _log);
5     public int length();
6 }
    
```

Iterable<TreeOperation>を継承しているため Iterator により TreeOperation を取り出せるようになっている。add や append メソッドを使って TreeOperation を積み上げていくことができる。

次にデータ編集により発生する TreeOperationLog の具体的な例を示す (ソースコード 3.2)。

ソースコード 3.2: トポロジーマネージャーの利用

```

1 [APPEND_CHILD:<-1>:pos:0]
2 [PUT_ATTRIBUTE:<-1,0>:key:author,value:oshiro]
3 [PUT_ATTRIBUTE:<-1,0>:key:mes,value:hello]
4 [PUT_ATTRIBUTE:<-1,0>:key:timestamp,value:0]
    
```

このログはルートノードに対し子ノードを追加し、追加した子ノードに attribute を 3 つ追加する際に図れるログである (図 3.9)。

大文字の英字は実行した NodeOperation の種類を表す。<> により囲まれている数字は NodePath を示す。NodePath の表記以降は Node の position や attribute の情報を表している。

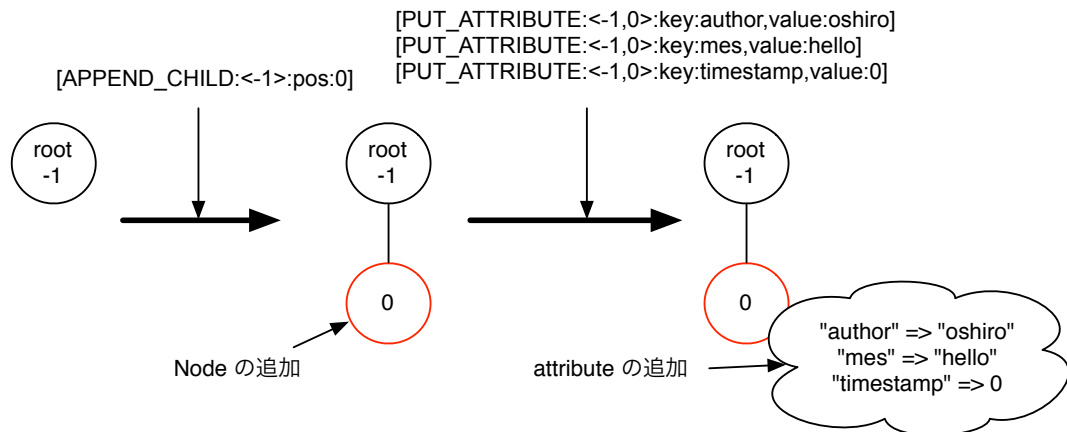


図 3.9: TreeOperationLog の具体例

ログの動作を表している図 3.9 の説明を行う。まず、APPEND_CHILD:<-1>:pos:0 により Root Node の 0 番目の子供となる Node の追加を行う。次に、追加を行った Node に対して PUT_ATTRIBUTE<-1,0> により attribute の情報を持たせていく。attribute の内容に作者の情報を表す author, メッセージの内容を表す mes, そしてタイムスタンプを timestamp とそれぞれキーにすることで追加される。

以上が掲示板プログラムにおける 1 つの書き込みで発生する TreeOperationLog である。

第4章 分散データベースJungleの設計

前章では木構造データベース Jungle の仕様について述べた。非破壊的木構造によりデータを破壊せずに保持する Jungle は、過去のデータもみることができる。また、データ編集をした際には、TreeOperationLog として履歴が残る。Jungle で分散設計をするにあたり、参考になっているシステムとして分散バージョン管理システムがある。本章では、分散バージョン管理システムとは何かを説明し、Jungle の分散データベースの設計について述べる。

4.1 分散バージョン管理システムによるデータの分散

Jungle の分散設計は Git や Mercurial といった分散バージョン管理システム（以下、分散管理システム）の機能を参考に作る。分散管理システムとは、多人数によるソフトウェア開発において変更履歴を管理するシステムである。分散管理システムでは開発者それぞれがローカルにリポジトリのクローンを持ち、開発はこのリポジトリを通すことで進められる（図 4.1）。ローカルのリポジトリは独立に存在し、サーバ上にあるリポジトリや他人のリポジトリで行われた変更履歴を取り込みアップデートにかけることができる。また、逆にローカルのリポジトリに開発者自身がかけたアップデートを他のリポジトリへと反映させることもできる。分散管理システムでは、どれかリポジトリが壊れたとしても、別のリポジトリからクローンを行うことができる。ネットワークに障害が発生しても、ローカルにある編集履歴をネットワーク復旧後に伝えることができる。そのため、可用性と分断耐性が高いと言える。

4.1.1 分散管理システムの API

分散管理システムでは、主に次の 3 つの API を使用することで、データの分散を行っている。

- commit
前のバージョンのデータに変更を加えたことをリポジトリに登録する。
- push
ローカルのリポジトリで行った変更履歴を別リポジトリへと伝える。
- pull
他のリポジトリの変更履歴をローカルのリポジトリに受け取る。

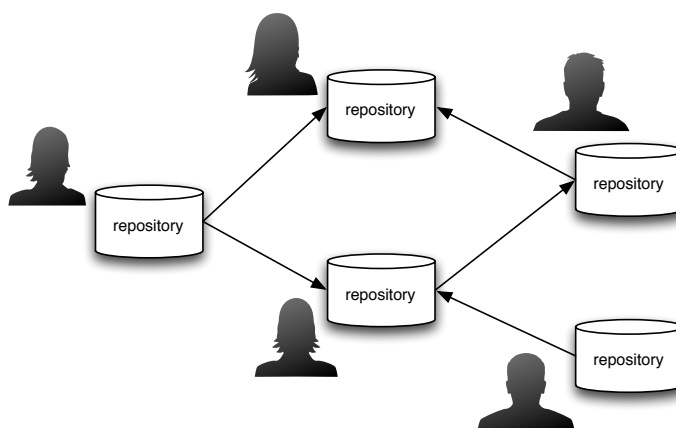


図 4.1: 分散管理システム

commit によりローカルのリポジトリに変更履歴が積み重なり, push をすることで他のリポジトリにローカルで行ったデータ変更を伝える. また, pull に好きなタイミングでデータの更新履歴を受け取ることができる.

分散管理システムはデータ変更の履歴を自由に受け取ることができるが, その変更履歴をそのまま適応できないときがある. それはお互いが同じデータに対して編集を行っている場合である. この状態を衝突という. この衝突を解決する方法として Merge がある. Merge は衝突を起こしたデータ両方の編集履歴を適用したデータを作成することである. 分散管理システムではこの衝突を解決する方法が必要になる.

4.1.2 Merge によるデータ変更衝突の解決

分散管理システムでは, データの更新時において衝突が発生する時がある. それは, 分散管理システムを参考にしている Jungle においても起こる問題である. データの変更を行うときには, 元のデータに編集が加えられている状態かもしれない. Jungle はリクエストがきた場合, 現在もっているデータを返す. そのためデータは最新のものであるかは保証されない. この場合, 古いデータに編集が加えられ, それを更に最新のデータへ伝搬させなければならない. このように他のリポジトリにより先にデータ編集が行われており, データの伝搬が素直にできない状態を衝突という. この衝突を解決する手段が必要である. 分散管理システムでは, 衝突に対して Merge と呼ばれる作業で解決をはかる. Merge は, 相手のリポジトリのデータ編集履歴を受け取り, ローカルにあるリポジトリの編集と合わせる作業である. データ衝突に対して Jungle はアプリケーションレベルでの Merge を実装することで解決をはかる.

以下に Merge が必要な場合とそうでない場合のデータ編集についての図を示す (図 4.2,4.3,4.4).

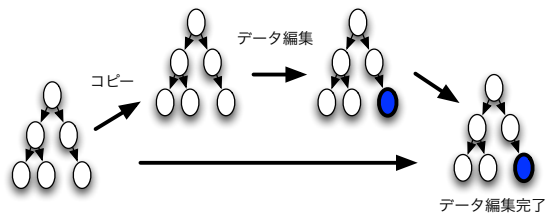


図 4.2: 衝突の発生しないデータ編集

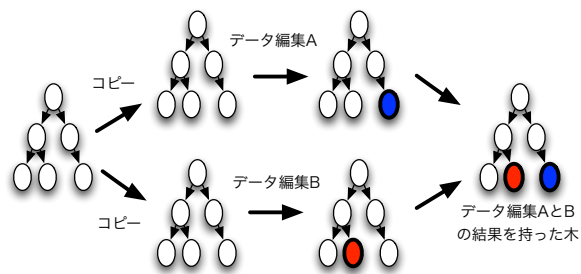


図 4.3: 自然に衝突を解決できるデータ編集

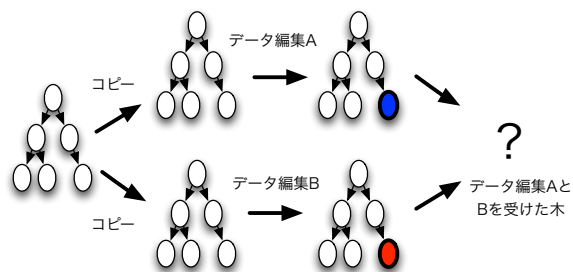


図 4.4: 衝突が発生するデータ編集

4.2 Jungle のネットワークトポロジーの形成

分散管理システムを参考に Jungle でもそれぞれのデータベースが独立に動くようにしたい。そのために必要なことはトポロジーの形成と、サーバノード間でのデータアクセス機構である。また、データ分散のために形成したトポロジー上で扱うデータを決めなければならぬ。

4.2.1 ツリートポロジーの形成とルーティング

分散データベース Jungle で形成されるネットワークトポロジーはツリー構造を想定している。ツリー構造ならば、データの整合性をとる場合、一度トップまでデータを伝搬させることで行える。トップもしくはトップまでの間にあるサーバノードでデータ伝搬中に衝突が発生したら Merge を行い、Merge の結果を改めて伝搬すればよいからである (図 4.5)。また、リング型 (図 4.6)、メッシュ型 (図 4.7) ではデータ編集の結果を他サーバノードに流すとき流したデータが自分自身にくることにより発生するループに気をつける必要がある。ツリー構造の場合は、サーバノード同士の繋がりで閉路が無い。そのため、自分自身が行ったデータ編集の履歴を繋がっているノードに送信するだけですむ。このルーティングの方式はスプリットホライズンと呼ばれるものである。

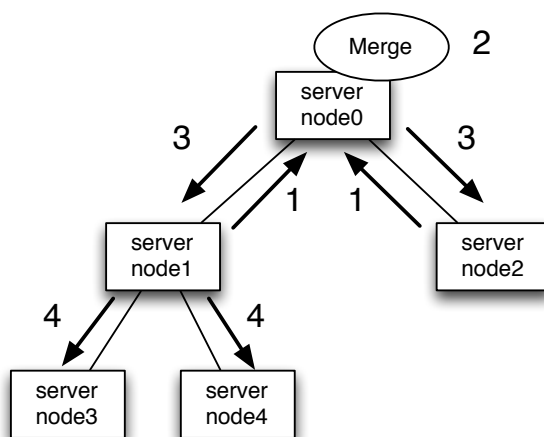


図 4.5: ツリー型の Network Topology

4.2.2 トポロジーの形成手段

Jungle で使用するネットワークトポロジーはツリー型を考えているが、リング型やメッシュ型といった他のネットワークトポロジーによる実装に関して試す余地はある。そのため、ツリーだけでなく、自由にネットワークトポロジーの形成を行えるようにしたい。

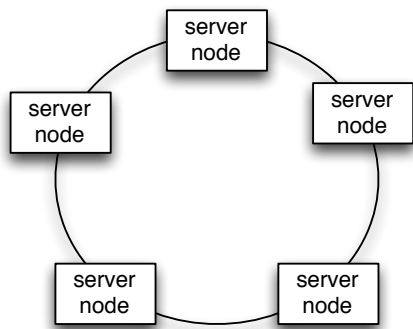


図 4.6: リング型のトポロジー

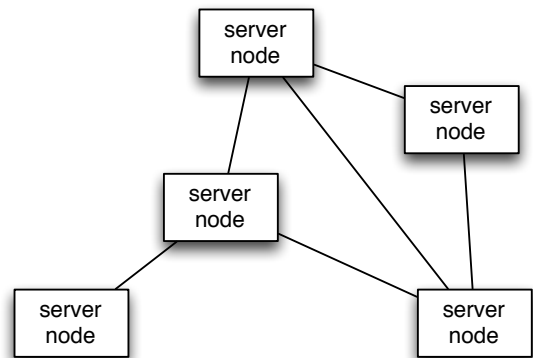


図 4.7: メッシュ型のトポロジー

そこで当研究室で開発を行っている並列分散フレームワークである Alice を使用する。Alice はユーザが望んだマシンへの接続や必要なデータへのアクセスを行う機構と、ネットワークトポロジー形成機能を提供している。

4.3 並列分散フレームワーク Alice

Alice は当研究室で開発している並列分散フレームワークである。Alice はデータを DataSegment, タスクを CodeSegment という単位で扱うプログラミングを提供している。タスクの部分となる CodeSegment は, 計算に必要なデータである DataSegment が揃い次第実行が行われる (図 4.8)。CodeSegment の結果により出力される新たなデータでは, 別の CodeSegment が実行されるための DataSegment となる。DataSegment と CodeSegment の組み合わせにより並列・分散プログラミングの依存関係が表されるようになっている。

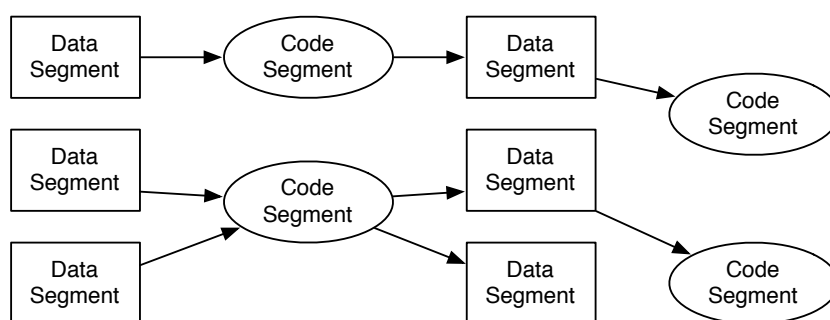


図 4.8: DataSegment と CodeSegment によるプログラムの流れ

4.3.1 MessagePack によるシリアルライズ

Alice では DataSegment のデータ表現に MessagePack[6] を利用している。MessagePack はオブジェクトをバイナリへと変換させるシリアルライズライブラリである。Alice によりネットワークを介してデータにアクセスするときは、そのデータが MessagePack でシリアルライズが行えることが条件である。

4.4 Jungle のデータ分散

Alice によりトポロジーの形成とデータアクセスの機構が提供される。後はデータ分散の為にどのデータをネットワークに流すのか決めなければならない。そこで選ばれたのが TreeOperationLog である。TreeOperationLog はデータ編集の履歴になる。どの Node にどのような操作をしたのかという情報が入っている。この TreeOperationLog を Alice を使って他サーバノードに送り、データの編集をしてもらうことで同じデータを持つことが可能となる。Alice を用いるため、この TreeOperationLog は MessagePack によりシリアルライズ可能な形にすることが必要である。

4.5 ログによるデータの永続性

Jungle は非破壊でオンメモリにデータを保持している。だが、オンメモリのままでは電源が落ちた際にデータが失われてしまう。ディスクからデータを読み込むことでデータの復旧を行えるようにしたい。そこで、ログによるデータの永続性の実装を行う。

Jungle の永続性実装の余地として Journal という機能が元々用意されている。この Journal にはディスクへ書き出すためのクラスとして Writer が用意されている。Jungle はデータの編集が完了した際にこの Writer クラスの write 関数を呼び出すようになっている。

この Journal と Writer クラスを新たに用意し、ディスクへと書き込みが行える機能を実装する。このとき、ログをどのようなデータ表現でハードディスクへと書きだすかという問題が発生するが、これは Alice を使うことで解決している。Alice を用いるため MessagePack によりシリアルライズ可能な TreeOperationLog ができる。このシリアルライズ可能な TreeOperationLog をそのままハードディスクへ書き込むことでログの永続性ができる。

第5章 木構造データベースJungleの分散実装

前章では Jungle のアーキテクチャと分散設計について説明した。トポロジーの形成と他サーバノードのデータのアクセス方法には Alice を使用する。また, Jungle ではデータ編集のログとして TreeOperationLog がある。この TreeOperationLog を Alice により他サーバノードへ送ることでデータの分散を行う。

本章では Jungle に行った分散実装について述べる。まず, Alice によるトポロジー形成の仕方を述べ, Alice を利用したプログラミングの方法について説明を行う。Alice の提供する分散プログラミングにおけるデータアクセス機構について述べ, それらを使用する上での注意点を述べる。次に, それらをふまえて Jungle の分散実装を行っていく。データ分散に必要なクラスを用意することで他サーバノードとの通信を行う。最後に, データ更新時の衝突において具体的な Merge の例として掲示板プログラムにおける Merge について述べる。

5.1 Alice のトポロジーマネージャの利用

5.1.1 トポロジーマネージャの起動

Alice を用いてサーバノードでトポロジーの形成を行う方法を述べる。Alice のトポロジーマネージャの起動はソースコード 5.1 の様に行う。

ソースコード 5.1: Alice によるネットワークトポロジーマネージャの起動

```
1 % java -cp Alice.jar alice.topology.manager.TopologyManager -p 10000 -conf ./topology/tree5.dot
```

-p オプションはトポロジーマネージャが開くポートの番号, -conf オプションには dot ファイルのパスを渡す。

ポート番号は Alice により記述された並列分散プログラムの起動時に渡す必要がある。dot ファイルには, トポロジーをどのように形成するかが書かれている。以下に, サーバノード数 5 で, 2 分木ツリー構造を形成する dot ファイルの例を示す (ソースコード 5.2)。

ソースコード 5.2: ネットワークトポロジ設定用 dot ファイル

```

1 % cat tree5.dot
2 digraph test {
3   node0 -> node1 [label="child1"]
4   node0 -> node2 [label="child2"]
5   node1 -> node0 [label="parent"]
6   node1 -> node3 [label="child1"]
7   node1 -> node4 [label="child2"]
8   node2 -> node0 [label="parent"]
9   node3 -> node1 [label="parent"]
10  node4 -> node1 [label="parent"]
11 }
    
```

node0 や node1 はサーバノードの名前を示す。サーバノードの間にはラベルがあり、Alice 上ではこのラベルに指定される文字列(キー)を使うことで他のサーバノードのデータへアクセスすることができる。node0 -> node1 はサーバノード同士の繋がりを示している。次に続く label="child1" は、node0 が node1 のデータに"child1"という文字列を使うことでアクセスできることを示す。

dot ファイルを読み込んだ Alice のトポロジーマネージャーに対して、サーバノードは誰に接続を行えばよいかを訪ねる。トポロジーマネージャーは訪ねてきたサーバノードに対してノード番号を割り振り、dot ファイルに記述している通りにサーバノード同士が接続を行うよう指示をだす。

トポロジーマネージャーは接続要求先を聞いてくるサーバノードに対して名前を割り振り、接続相手を伝える。dot ファイル 5.2 により形成されるトポロジを図 5.1 に示す。

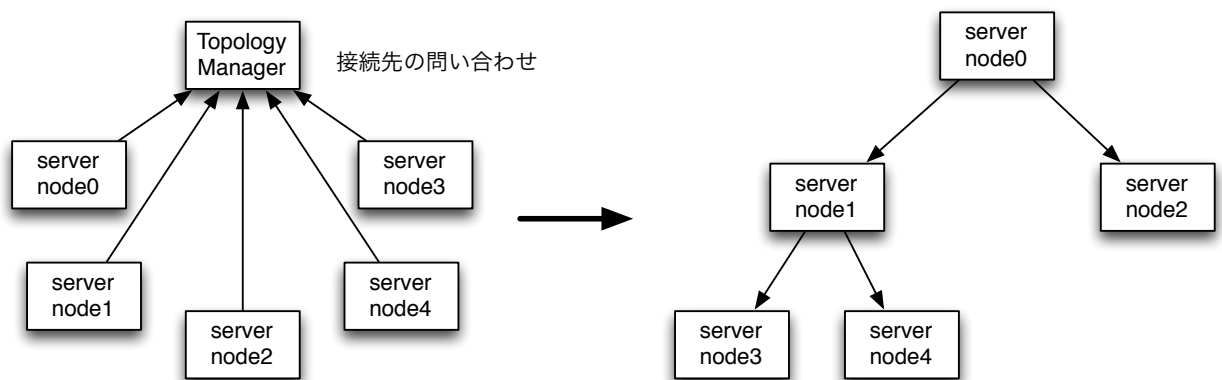


図 5.1: Alice によるネットワークトポロジ形成

5.1.2 アプリケーション側の記述

次は Jungle 側のプログラムが最初に Alice のトポロジーマネージャーと通信を行うようにする。そのためには Alice の TopologyNode クラスに必要な情報を渡してインスタンスを生成する(ソースコード 5.3)。

ソースコード 5.3: アプリケーションの起動

```
1 public static void main( String[] args ) throws Exception
2 {
3     RemoteConfig conf = new RemoteConfig(args);
4     new TopologyNode(conf, new StartJungleCodeSegment(args, conf.bbsPort));
5 }
```

TopologyNode クラスは第 2 引数として CodeSegment を受け取る。TopologyNode のインスタンスはまず初めにトポロジーマネージャーへ接続を行う。次にトポロジーマネージャーから受け取った情報を元に別のサーバノードとトポロジーマネージャーの形成を行う。その後、第 2 引数で渡された StartJungleCodeSegment の実行を行う。StartJungleCodeSegment には通常のアプリケーションの処理が書かれる。

アプリケーションの起動時にはコンフィグの情報として、トポロジーマネージャーが動いているサーバのドメインとポート番号を渡す必要がある。例えば、mass00.cs.ie.u-ryukyu.ac.jp というサーバ上でポート番号 10000 を指定してトポロジーマネージャーを起動した場合は次のようになる (ソースコード 5.4)。

ソースコード 5.4: トポロジーマネージャーの利用

```
1 % java Program -host mass00.cs.ie.u-ryukyu.ac.jp -port 10000
```

5.2 Alice を用いての分散実装

Alice のポロジーマネージャー形成と他のサーバのデータへのアクセスする機構を用いるためには、Alice が提供するプログラミングスタイルに沿わなければならない。それは DataSegment (データ) と CodeSegment (タスク) によるプログラムである。ここではまず、DataSegment と CodeSegment によるプログラムの方法について説明し、他サーバとの通信部分の実装について述べる。

5.2.1 Alice によるプログラミング

Alice は DataSegment (データ) と CodeSegment (タスク) 単位でプログラミングを行うことを述べた。CodeSegment には計算に必要な DataSegment が登録される。そして DataSegment が準備され次第 CodeSegment による計算が実行される。DataSegment の取得は文字列のキーを使うことで行える。ソースコード 5.5 に CodeSegment の例を示す。

ソースコード 5.5: CodeSegment の実行

```
1 public class TestCodeSegment extends CodeSegment {
2     public Receiver arg1 = ids.create(CommandType.TAKE);
3
4     public TestCodeSegment() { }
5
6     public void run() {
7         int count = ds.asInteger();
8         count++;
9         System.out.println("count_□=□"+count);
10        if(c > 10) { exit(0); }
11        CodeSegment cs = new TestCodeSegment();
12        cs.setKey("count");
13        ods.put("local", "count", c);
14    }
15
16    public static void main(String[] args) {
17        CodeSegment cs = new TestCodeSegment();
18        cs.arg1.setKey("local", "count"); // setKey API
19        cs.ods.put("local", "count", 0);
20    }
21 }
```

コードの説明を行う。このプログラムは、数字を 1 から 10 まで出力を行い終了するプログラムである。17 行目から 19 行目の処理が最初に行われる。まず `TestCodeSegment` という `CodeSegment` のインスタンス `cs` を生成する。 `cs` は `arg1` という `Receiver` クラスのフィールドを保持しており、 `Receiver` クラスは `DataSegment` を受けとるためのクラスである。 `arg1` に対し `setKey` API を使うことで、使用したい `DataSegment` のキー”count”を登録することができる。これによりキー”count”に対してデータが登録された場合、そのデータを受け取り `cs` の計算が自動で始まる。 `setKey` API の第一引数に渡している”local”はどのマシンの `DataSegment` にアクセスするのかを指定している。この場合は自分自身を表す”local”になる。

データの登録は `ods.put` により行える。上記のコード 19 行目では `put` により”count”をキーとして数値の 0 を登録している。 `put` がされると `cs` の計算が始まり別スレッドにより 8 行目からの処理が行われる。

`put` によりキー”count”に登録された数値 0 は `Receiver` である `ds` を使って取ることができる。7 行目から 13 行目では `ds.asInteger()` により、”count”に登録したデータの中身を受け取りインクリメントし出力する。そして、最後には `ods.put` を行っている。新たな `TestCodeSegment` も生成しており、これはインクリメントされた”count”が `put` されることで実行される。この一連の処理を”count”の数値が 10 以上になるまで行う。

`DataSegment` ヘデータの追加と `CodeSegment` の実行について表した図 5.2 になる。

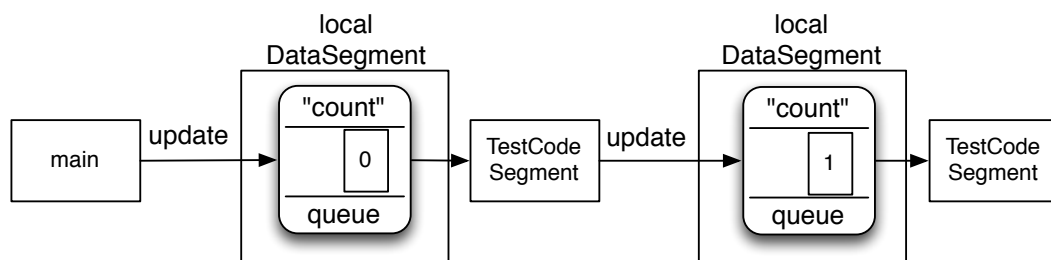


図 5.2: DataSegment と CodeSegment によるプログラムの例

5.2.2 他サーバノードの DataSegment へアクセス

Alice における基本的なプログラミングは述べた。次はネットワークを介して他サーバノードの DataSegment にアクセスするプログラムについて述べる。

まず, Alice により 2 分木 3 ノードのトポロジーが形成された場合を想定する。その時に実際に作られるトポロジーは図 5.3 となる。

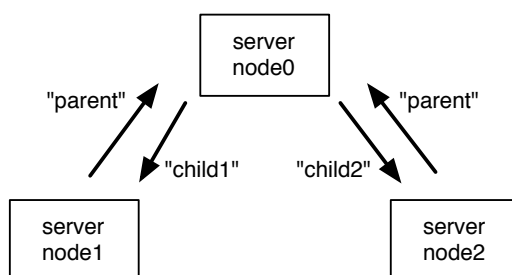


図 5.3: トポロジーの形成

ネットワークを介した DataSegment へのアクセスはそのサーバノードを示す文字列のキーを追加することで行える。他サーバノードを示す文字列のキーとは図 5.3 に矢印の隣に書かれている文字列 "parent", "child1", "child2" のことを指す。例えば, server node0 が server node1 の DataSegment に入っている "count" というデータをを使用したい場合は, 次のように setKey を行えばよい (ソースコード 5.6)。

ソースコード 5.6: CodeSegment で他サーバノードの DataSegment を使用する

```
1 CodeSegment cs = new RemoteCodeSegment();
2 cs.arg1.setKey("child1", "count");
```

また, 他サーバノードの DataSegment にデータを送りたい場合は, put を行うときにサーバノードへのキーを追加すればよい。例として, server node1 や server node2 が server node0 の DataSegment に "message" というキーでデータを追加したい場合は, 次のようになる (ソースコード 5.7)。

ソースコード 5.7: 他サーバーノードの DataSegment にデータを追加する

```
1 ods.put("parent", "message", "Hello_parent");
```

5.2.3 独自クラスのインスタンスの送受信

最後に、独自クラスのインスタンスの DataSegment での扱い方について述べる。Alice では MessagePack を用いてシリアライズを行い他サーバノードへと送信している。MessagePack はインスタンス単位でシリアライズを行うことができる。そのため、Alice でもプリミティブな型に限らずクラスのインスタンスを DataSegment として扱うことができる。

MessagePack によりシリアライズとなるクラスはいくつか制限がある。それはそのクラスに @Message アノテーションを付けることと、そのクラスが保持するフィールドが MessagePack によりシリアライズ可能であることである。例えば次のようなクラスである (ソースコード 5.8)。

ソースコード 5.8: MessagePack によりシリアライズ可能なクラス 1

```
1 import org.msgpack.annotation.Message
2
3 @Message
4 public class Student {
5     String name;
6     int age;
7 }
```

上記の Student クラスはプリミティブ型しか保持していない。そのためシリアライズが可能であるまた、次のようなクラスもシリアライズ可能な型となる (ソースコード 5.9)。

ソースコード 5.9: MessagePack によりシリアライズ可能なクラス 2

```
1 import org.msgpack.annotation.Message
2
3 @Message
4 public class Class {
5     List<Student> studentList;
6 }
```

この場合、フィールドはプリミティブな型でない Student クラスのフィールドを保持している。しかし、Student クラスはシリアライズ可能な形で作成しているため、クラスのフィールドとして保持しても問題はない。

これらの制約にそった形で作成し DataSegment にネットワークを介してクラスのインスタンスを put することができる。DataSegment から受け取ったデータはそのままではシリアライズされたものため、一度手元で元のクラスにコンバートすることで扱う。例として、Alice における Student クラスのコンバートを次に示す (ソースコード 5.8)。

ソースコード 5.10: DataSegment

```
1 // public Receiver arg1 = ids.create(CommandType.PEEK);
2 Student s = arg1.asClass(Student.class);
```

MessagePack でシリアライズ可能な形としているため DataSegment はネットワークを介して送受信が可能である。

5.3 ログのシリアライズ

Jungle の具体的な分散実装について述べる。実装にあたり、解決しなければならない問題はまず、ログを DataSegment で扱える形にすることである。そのためには、@Message アノテーションを付けたログのクラスの作成を行わなければならない。

5.3.1 TreeOperationLog のシリアライズ

ログの実体である TreeOperationLog をシリアライズ可能な形にするにあたって気をつけなければならないのが、フィールドをシリアライズ可能な部分である。TreeOperationLog は TreeOperation をいくつも保持し、TreeOperation は NodePath と NodeOperation を保持するものであった。そのため、これら全てをシリアライズ可能な形にしなければならない。

基本的にこれらの実装は、フィールドを全てプリミティブなものだけにすればよい。MessagePack は List を扱うこともできるため、TreeOperationLog で継承されていた Iterable の挙動も List を使うことで実装を行うことができた。

5.4 ログに対する情報の追加

TreeOperationLog をシリアライズ可能な形にした後、問題が発生した。それは、TreeOperationLog 自体は木の名前を保持していないというものである。そのため、TreeOperationLog だけを受け取っても、そのログがどの木に対して行われるのかわからなかった。そこで、TreeOperationLog の情報だけでなく、木の名前と UUID、それと timestamp の情報も付与してシリアライズが可能な NetworkTreeOperationLog の実装を行った。

5.4.1 NetworkTreeOperationLog の実装

NetworkTreeOperationLog の実装の一部を以下 (図 5.11) に示す。

ソースコード 5.11: NetworkTreeOperation が持つフィールド

```
1 @Message
2 public class NetworkTreeOperationLog implements TreeOperationLog
3 {
4     public LinkedList<NetworkTreeOperation> list;
5     String treeName;
6     long timestamp;
```

List により保持している NetworkTreeOperation は TreeOperation をシリアライズ可能な形にしたものである。TreeOperationLog を implements し、木の名前と timestamp をを保持する。他サーバードへ伝える必要のある情報が増えた場合、このように NetworkTreeOperationLog に情報を付与することで対応することができる。

5.4.2 ログの送信部分

ログを送信するタイミングはいつ行うか。それは、木の編集が成功した時である。木の編集が成功した結果得られる `TreeOperationLog` を `NetworkTreeOperationLog` に変換し、`ods.put` を使って `CodeSegment` 側から利用できるようにする。

しかし、この時気をつけなければならないことがある。それは、`ods.put` の処理をレスポンスを返すスレッドの中で行うと、レスポンスが悪くなる可能性があることだ。そのため、`ods.put` を行うのは別の `Thread` にしたほうがよい。以下のコードは `commit` に成功した後に、`NetworkTreeOperation` へと変換したログを別スレッドに渡して処理させるコードである (ソースコード 5.12)。

ソースコード 5.12: `NetworkTreeOperation` を `put` するために別スレッドを立ち上げる

```
1 NetworkTreeOperationLog netLog = new NetworkTreeOperationLog(_uuid, _treeName, newLog);
2 CodeSegment cs = new LogPutCodeSegment(netLog);
3 cs.execute();
```

`LogPutCodeSegment` の実装は次のようになっている (ソースコード 5.13)。

ソースコード 5.13: `put` を行うためだけの `CodeSegment` の用意

```
1 public class LogPutCodeSegment extends CodeSegment {
2     NetworkTreeOperationLog log;
3     public LogPutCodeSegment(NetworkTreeOperationLog _log) {
4         log = _log;
5     }
6     @Override
7     public void run() {
8         ods.put("log", log);
9     }
10 }
```

上で述べた問題は、ベンチマークテストなど、大量の負荷をかけたさいに発生する。負荷とは `Jungle` のデータに変更が加わることである。多数のデータの変更により大量のログが生成される。そのため、`ods.put` により `DataSegment` の "log" にアクセスが集中してしまい、レスポンスが悪くなっていた。`ods.put` を行うタイミングには気をつけなければず、上記のコードにしても改良の余地はある。

5.4.3 他サーバノードへのログの送信

上記の実装により `DataSegment` の "log" キーにアクセスすることで `TreeOperationLog` を取得できるようになった。次はこのデータを他サーバノードへ送る部分の実装を行う。

他サーバノードに対してデータを送るときに必要なことは、そのサーバノードの `DataSegment` にアクセスするキーである。Alice では接続を行っている相手の `DataSegment` のキーのリストは "_CLIST" キーを使うことで取得することができる。"_CLIST" キーで得られるリストを使って他サーバノードへとデータを `put` するコードを次に示す (ソースコード 5.14)。

ソースコード 5.14: 他サーバノードへのログの送信部分

```

1 public class LogUpdateCodeSegment extends CodeSegment {
2     Receiver log = ids.create(CommandType.TAKE);
3     Receiver clist = ids.create(CommandType.PEEK);
4
5     public LogUpdateCodeSegment() {
6         log.setKey("log");
7         clist.setKey("_CLIST");
8     }
9
10    List<String> list = clist.asClass(List.class);
11    for (String node : list) {
12        ods.put(node, log.key, log.getVal()); // Send datasegment to other node
13    }
14    :

```

12 行目の `ods.put(node, log.key, log.getVal())` が他サーバノードへデータを送る部分になる。変数 `log` には `NetworkTreeOperationLog` が入っている。変数 `list` には `"_CLIST"` により得られたデータが入っている。それは文字列の `List` となっている。 `list` の中身を 1 つずつ受け取り、 `ods.put` する際に引数として渡してやることで、他サーバノードの `DataSegment` にアクセスが行える。

`list` の具体的な内容は、図 5.1 で組まれたトポロジでいうところの、 `"parent"`、 `"child1"`、 `"child2"` の文字列の `List` となっている。もし子どもとなるサーバノードがないときは `"parent"` だけが `List` には入れられる。

5.4.4 ログの受信とデータ反映

次は受け取ったログからデータの編集を行う部分の実装を行う。 `Jungle` はのデータを変更する手段として木構造データ毎に `TreeEditor` クラスが提供される。この `TreeEditor` を使用し、ログに入っている `TreeOperation` を 1 つ 1 つ取り出し同じ編集を行わせる。例えば次のようになる (ソースコード 5.15)。

ソースコード 5.15: ログを受け取ってのデータの反映

```

1 // Receiver log ← "log"キーから取得できるデータが入っている
2 // Jungle jungle
3 NetworkTreeOperationLog netLog = log.asClass(NetworkTreeOperationLog.class);
4 String treeName = netLog.getTreeName();
5 JungleTree tree = jungle.getTreeByName(treeName);
6 TreeEditor editor = tree.getLocalTreeEditor(); // Editor の取得
7 for (TreeOperation op : netlog) {
8     NodePath path = op.getNodePath();
9     NodeOperation nodeOp = op.getNodeOperation();
10    either = edit(editor, path, nodeOp, nodeOp.getPosition()); // データの編集を行う。
11    if(either.isA()) {
12        // エラー処理. 編集失敗
13    }
14    editor = either.b();
15 }

```

7 行目で取り出された `TreeOperation` からさらに `NodePath` と `NodeOperation` を取り出しているのが 8 行目と 9 行目になる。最後に `edit` 関数に `TreeEditor` と `NodePath`、それと

NodeOperation を引き渡している。edit 関数は次のようになる (ソースコード 5.16)。

ソースコード 5.16: edit 関数の実装

```
1 Either<Error, JungleTreeEditor> edit(JungleTreeEditor editor, NodePath path,  
2 NodeOperation nodeOp, int pos) {  
3     String key = "";  
4     Command c = nodeOp.getCommand();  
5     switch (c) {  
6         case PUT_ATTRIBUTE:  
7             key = nodeOp.getKey();  
8             ByteBuffer value = nodeOp.getValue();  
9             return editor.putAttribute(path, key, value);  
10        case DELETE_ATTRIBUTE:  
11            key = nodeOp.getKey();  
12            return editor.deleteAttribute(path, key);  
13        case APPEND_CHILD:  
14            return editor.addNewChildAt(path, pos);  
15        case DELETE_CHILD:  
16            return editor.deleteChildAt(path, 0);  
17    }
```

NodeOperation の API の種類 (Command) 毎に使用する TreeEditor の API を変えている。もし NodeOperation の API の種類を増えるようなことが合っても、上記のコードのように対応する TreeEditor の API を書くことで対応できる。

5.5 永続性の実装

次は、ログの書き出しによる永続性の実装について述べる。第 3 章で Jungle は Writer という永続性実装のための機能が元々用意されていることを述べた。永続性で考えなければならないことは、どのようなデータをどんなデータ表現で書き込むかということである。今回の実装ではログである `TreeOperationLog` を書き出す。また、`TreeOperationLog` の情報を保持しつつ、`MessagePack` でシリアライズできるクラスとして `NetworkTreeOperationLog` の実装を行った。

つまり、Writer で `NetworkTreeOperationLog` を書き出すようにすればよい。以下にログをディスクへ書き出すためのクラス `PersistentChangeListWriter` の実装を示す (ソースコード 5.17)。

ソースコード 5.17: `NetworkTreeOperation` をディスクへ書き出す

```
1 public class PersistentChangeListWriter implements ChangeListWriter {
2
3     MessagePack msgpack;
4     OutputStream out;
5
6     public PersistentChangeListWriter(OutputStream _out, MessagePack _msgpack) {
7         out = _out;
8         msgpack = _msgpack;
9     }
10    @Override
11    public Result write(ChangeList cs)
12    {
13        NetworkTreeOperationLog log
14        = new NetworkTreeOperationLog(cs.uuid(), cs.getTreeName(), cs);
15        try {
16            msgpack.write(out, log);
17            out.flush();
18            return Result.SUCCESS;
19        } catch (IOException e) {
20            // エラー処理
21        }
22        return Result.FAILED;
23    }
```

`write` 関数は Jungle のデータ編集が完了すると呼び出される関数である。引数として渡されている `ChangeList` は `TreeOperationLog` と同じく `Iterable<TreeOperation>` を継承している。また UUID や木の名前も取得することができる。そのため、`NetworkTreeOperationLog` へと変換が行える。

ログの書き出しを行いたいときはこの `PersistentChangeListWriter` を設定することで行えるようになった。これにより木の編集が行われるたびに `NetworkTreeOperationLog` が書き込まれていく。読み込みたいときは `MessagePack` を使ってディスクから読み込み、データ分散実装と同じの方法で木の編集を行っていくことができる (ソースコード 5.15, 5.16)。

5.6 Merge の実装

Jungle に分散実装を行った後の問題としてデータ衝突がある。他のサーバノードから送られてくるデータが既に手元で変更を加えた木構造を対象とした場合に発生する問題である。Jungle ではこれをアプリケーション毎に Merge を実装することで解決させる。

5.7 掲示板プログラムにおけるデータ衝突

今回分散実装を行い、例題として掲示板プログラムを用意した。掲示板プログラムに実装を行った Merge について述べる。まず Jungle を用いた掲示板プログラムのデータ保持方法を図 5.4 に示す。

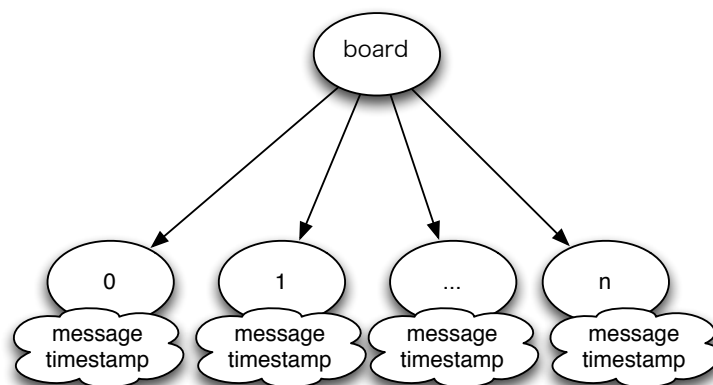


図 5.4: Jungle による掲示板プログラムのデータ保持方法

掲示板プログラムでは各掲示板毎に 1 つの木構造が作成される。掲示板への 1 つの書き込みは子ノードを 1 つ追加することに相当する。また、各子ノードは attributes として書き込みの内容である message と書き込まれた時間を表す timestamp を保持している。先に追加された順で子ノードには若い番号が割り振られる。

他サーバノードからの書き込みをそのまま子ノードの後ろに追加してしまうと、データの整合性が崩れてしまう。この時の状態を表しているのが図 5.5 と 5.6 になる。

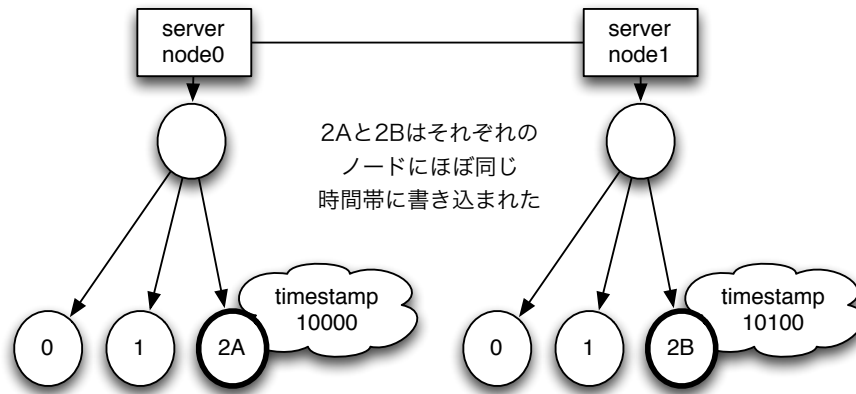


図 5.5: 他サーバノードの編集データ反映による整合性の崩れ 1

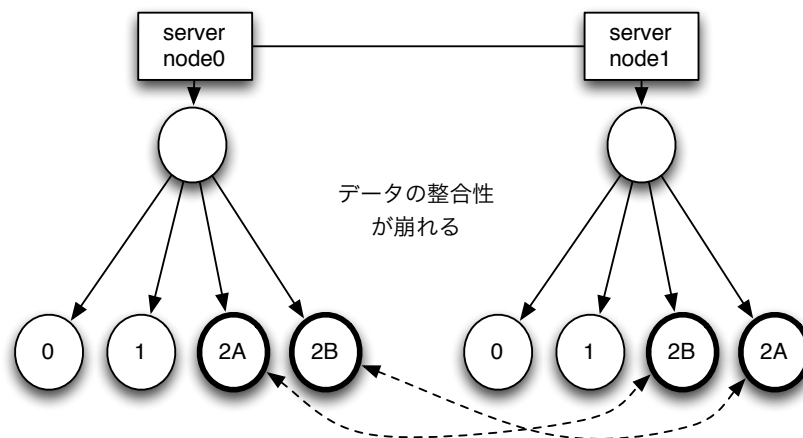


図 5.6: 他サーバノードの編集データ反映による整合性の崩れ 2

5.7.1 掲示板プログラムにおける Merge

図 5.6 の server node0 の木の状態にするのが理想である。掲示板への書き込みの表示は、書き込みされた時間が早い順に表示されるようにしたい。これを timestamp を利用することで行う。他サーバノードから来たデータに関しては、timestamp を参照し、次に自分の保持している木の子ノードの timestamp と比べていくことでデータの追加する場所を決める。これが今回実装を行った掲示板システムにおける Merge になる。

第6章 分散木構造データベース Jungle の評価

前章では Jungle における分散データベースの詳細な実装について述べた。本章では実装を行った Jungle に対して Cassandra との性能比較を行い評価をする。性能比較の為に簡易な掲示板プログラムを Jungle と Cassandra それぞれに作成した。複数のノードに繋がっている状態においても性能を測りたいため、学科が提供する VMWare の並列環境を利用する。また、我々の研究室が利用しているブレードサーバ上で動いている KVM もクライアントとして利用する。Jungle は永続性はなく分散だけ実装で測定を行っている。

6.1 実験方法

実験は同じ機能を提供している簡易掲示板プログラムを Jungle と Cassandra それぞれで動かす、HTTP リクエストにより負荷をかけて行う。レスポンスが返ってくるまでの時間をはかり、平均時間と標準偏差を求めグラフに出力する。

また、実験は2つ行う。まず行う実験は、複数のクライアントから1つのノードに負荷をかける方法である(図 6.1)。

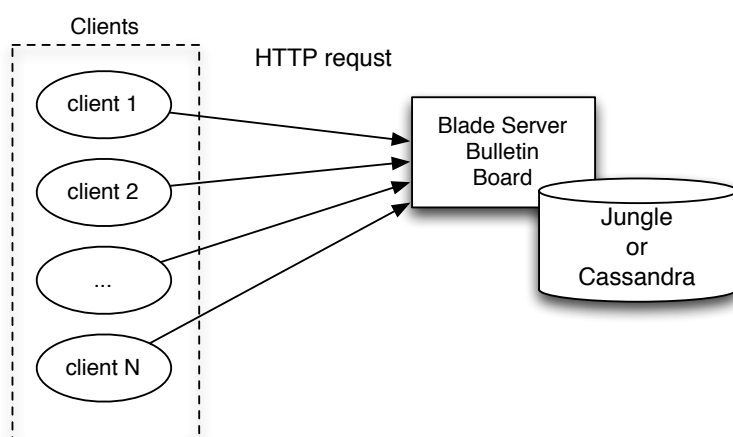


図 6.1: 実験 1 複数のクライアントからサーバ 1 台への負荷

次に行う実験は複数のノードに対し複数のクライアントから負荷をかける方法である。それぞれ大量の HTTP リクエストをだし、全てのリクエストの処理にかかる時間を測定する。

クライアントの数に比例してノードを増やすことでレスポンスを維持できるかスケールビリティを調べるためである (図 6.2)。

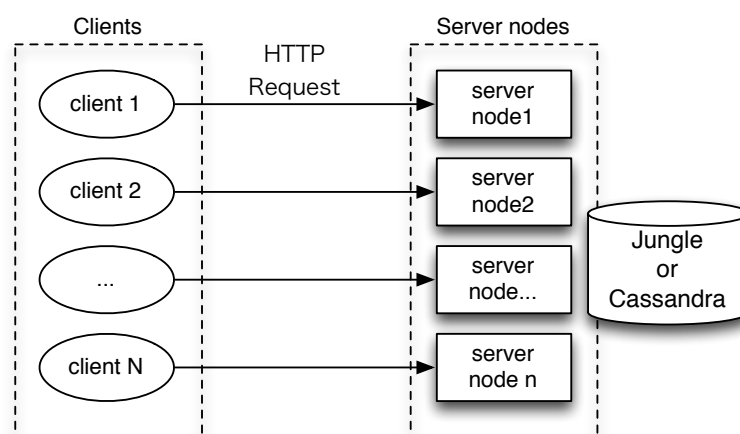


図 6.2: 実験 2 複数のクライアントから複数のノードへの負荷

6.1.1 Torque Resource Manager

並列環境下にあるマシン全てに命令を出し、タスクを実行させることは非常に大変である。そのため、今回の実験において並列環境のマシンに同時にタスクを実行させるために Torque Resource Manager を利用する。Torque は Queue によりタスクの実行順序を制御する。Queue にタスクをいれる際には、そのタスクをいくつかのノードで実行するか、いくつかのコア数を使用するかといったリソースの設定も行うことができる。

6.1.2 weighttp

最初の実験で1つのノードに負荷をかけるプログラムはウェブサーバの測定ツールである weighttp を使用する。weighttp は総リクエスト数、同時接続数、ネイティブスレッド数をオプションとして指定することができる C 言語で書かれたプログラムである。

6.1.3 掲示板プログラム

今回使用する掲示板プログラムは組み込み用ウェブサーバである Jetty をフロントエンドとして利用し、バックエンドに Jungle と Cassandra を利用している。

表 6.1: 簡易掲示板システムで利用した Jetty と Cassandra のバージョン

名前	バージョン
Jetty	6.1.26
Cassandra	2.0.4

6.1.4 実験環境

サーバノードとクライアントを実行させるサーバの仕様

使用する VMWare と KVM のクラスタの使用を以下に示す. クラスタは仕様を表 6.2 と表 6.3 に示す.

表 6.2: 掲示板プログラムを実行させる VMWare クラスタの仕様 (クライアントにも利用)

名前	概要
CPU	Intel(R) Xeon(R) CPU X5650@2.67GHz
Memory	8GB
OS	CentOS 5.8
HyperVisor	VMWare ESXi
JavaVM	Java(TM) SE Runtime Environment (build 1.7.0-b147)

表 6.3: クライアントを実行させる KVM クラスタの仕様

名前	概要
CPU	Intel(R) Xeon(R) CPU X5650@2.67GHz
Memory	8GB
OS	CentOS 5.8
HyperVisor	KVM
JavaVM	Java(TM) SE Runtime Environment (build 1.7.0-b147)

ブレードサーバの仕様

最初の実験ではブレードサーバ 1 台で掲示板プログラムを動かす, 並列環境から複数のクライアントで負荷をかける. ブレードサーバの仕様を表 6.4 に示す

表 6.4: サーバノードとして利用するブレードサーバの使用

名前	概要
CPU	Intel(R) Xeon(R) CPU X5650@2.67GHz
物理コア数	12
論理コア数	24
Memory	132GB
OS	Fedora 16
JavaVM	Java(TM) SE Runtime Environment (build 1.7.0_51-b13)

Jungle 実行時の JavaVM のオプションの設定

サーバで Jungle を実行するときは, JavaVM がデフォルトで設定している Heap サイズの容量を大きくする. Jungle では非破壊でデータを保持するため, データで使用するメモリの量が大きい. Java の Heap サイズをデフォルトのままベンチマークプログラムを走らせると, エラーの `java.lang.OutOfMemoryError: GC overhead limit exceeded` が出力されてプログラムが終了してしまう. このエラーは Full GC にかかる回数が多いか, プログラムの 98%以上 GC に使用されていると出力されるエラーである. そのため, ブレードサーバでは `-Xmx20g -Xms10g` をつけ, VM 側では `-Xmx6g -Xms4g` のオプションを付けて行う.

サーバの環境

HTTP によりノードに負荷を掛ける場合気をつけることがある. それはサーバの設定により最大コネクション数や開くことのできるファイル記述子の数に制限がかかっていることである. この 2 つの値はデフォルトでは小さなものとなっており, そのままではカーネルの設定がネックとなったベンチマーク結果がでる可能性がある. そこで次のようにコマンドを実行することでコネクション数の制限を増やすことができる.

ソースコード 6.1: コネクション数を増やす

```
1 % sudo sysctl -w net.core.somaxconn=10000
```

ファイル記述子の制限を増やす場合は次のコマンドを実行する

ソースコード 6.2: ファイル記述子の制限を増やす

```
1 % ulimit -n 10000
```

6.2 実験結果 1

複数のクライアントからサーバノード一台に対して負荷をかける実験を行った。クライアントの数は 10 台から始まり 5 台ずつ増やしていき、最大 45 台まで増える。各クライアントにおいて以下のオプションをつけた `weighttp` プログラムが実行される。

ソースコード 6.3: `weighttp` のオプション (実験 1)

```
1 weighttp -n 20000 -c 20 -t 2 -k "http://url"
```

このオプションは 2 つのネイティブスレッドを使用し、同時に 20 のコネクションを張り、通信の間、コネクションを切らずに 2 万件の HTTP request を送信することを表している。Cassandra はサーバノードが一台の為、Replication factor 1 で Consistency Level は ONE となる。

実験の結果はグラフ 6.3, 6.4 となる。横軸はクライアントノードの数を表しており、値が増えるほどリクエストの数も増え負荷が高まる。縦軸は 2 万件のリクエスト全てにレスポンスを返し終えた時間を表している (単位: 秒)。

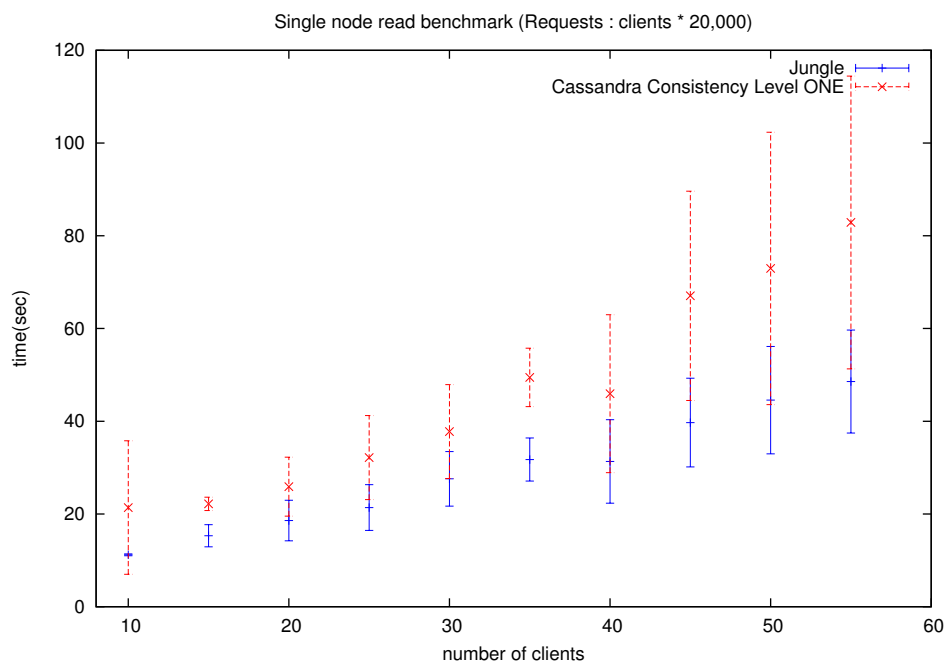


図 6.3: 複数のクライアントから一台への負荷

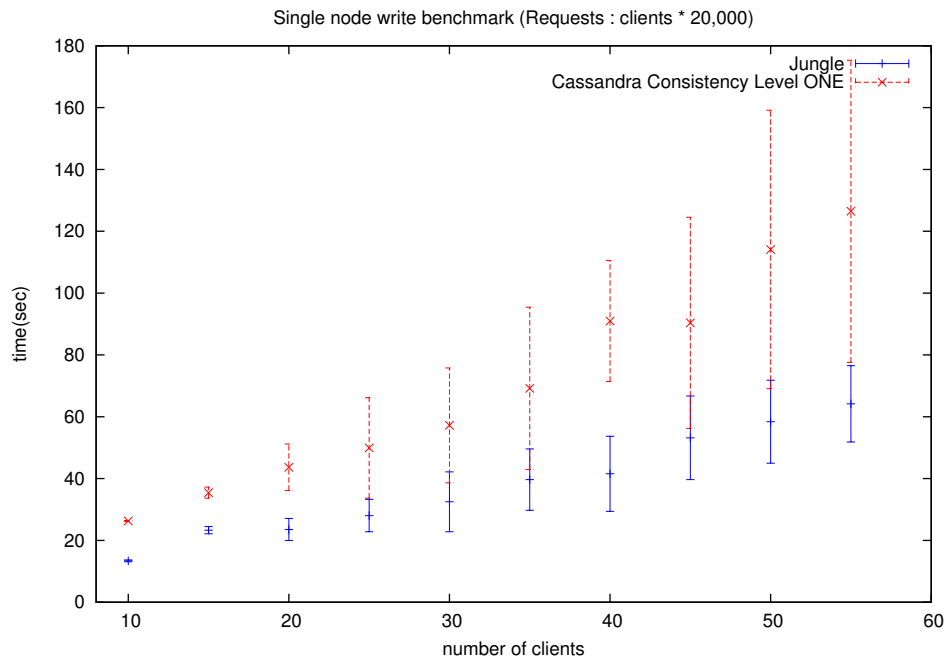


図 6.4: 複数のクライアントから一台への負荷

6.2.1 実験結果 1 の考察

読み込み, 書き込みともに Jungle のほうが良い結果となっている. 書き込みの差が大きく開いていることに関しては, Cassandra はディスクへと書きだすときもあるのも原因の 1 つと考えられる. Jungle はオンメモリであることから, やはり差はでてしまう. しかしディスクに書き出していないことは別の要因も考えられる. Jungle は非破壊的木構造なため, ロックをほとんど必要としない. 書き込み時においてもロックが必要なときは木のコピーをとりおえて, ルートノードを更新するときのみである. 書き込みの速度が早いことは Jungle のロックが少ないことも要因の 1 つとしてあげられる.

6.3 実験結果 2

学科の並列環境クラスタを用いて分散環境下での実験を行う学科の提供する VM は 48 台だが、ブレードサーバ上で動く KVM から 12 台を利用し、合計 60 台を使用する。Jungle と Cassandra をそれぞれサーバノード 10 台、20 台、30 台で動かす、クライアントも 10 台、20 台、30 台と増やして負荷をかける。クライアントとサーバノードの数は 1:1 となるため、横軸の値の数が増えると総リクエストは増えても 1 台に与えるリクエスト数は変わらない。縦軸はリクエストに全てに対してレスポンスを返さった時間を表す (単位: 秒) KVM 側はクライアント側だけに利用する。weighttp に付ける引数は実験 1 と同じとする。各クライアントから 2 万のリクエストを送る。Cassandra は Consistency Level ONE と QUORUM の両方を計測する。QUORUM の Replication factor は 5 で設定してある。

測定は読み込みと書き込みの両方を行う。測定の結果をグラフにしたのを図 6.5, 6.6 に示す。横軸はクライアントとサーバノードの数を表す。

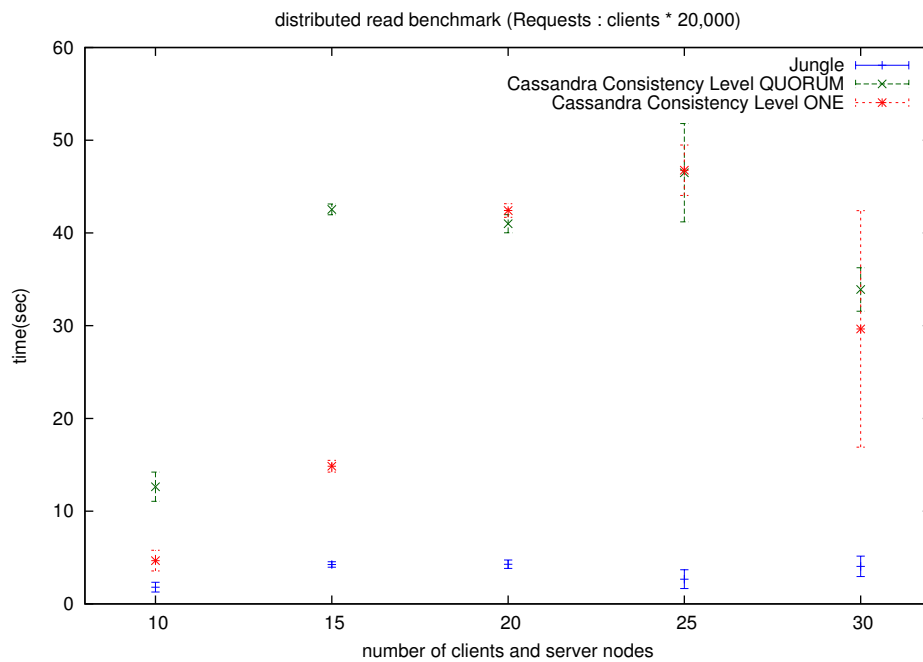


図 6.5: 分散環境下における読み込みベンチマーク結果

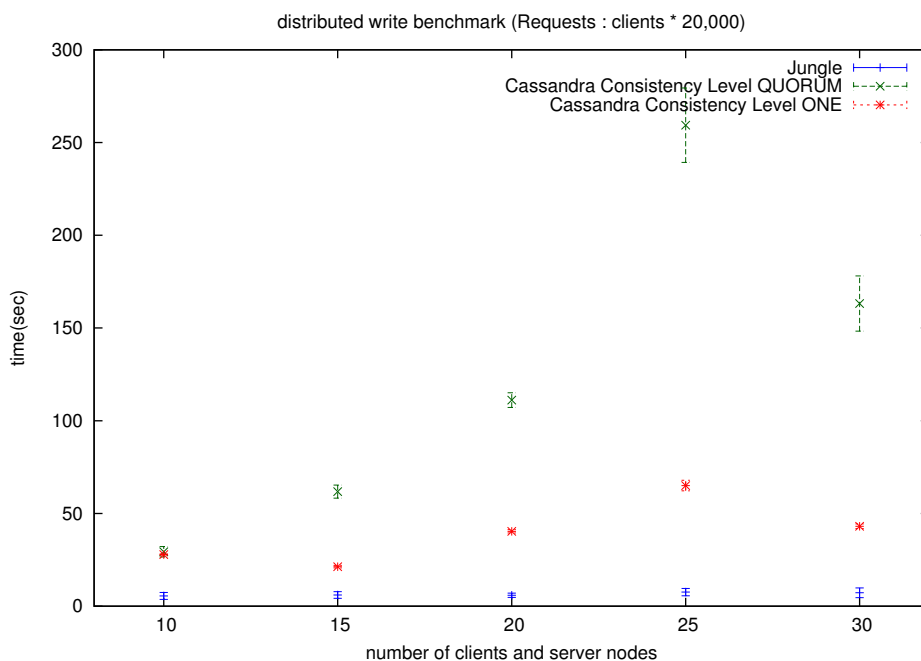


図 6.6: 分散環境下における書き込みベンチマーク結果

6.3.1 実験結果 2 の考察

こちら、Jungle のほうが Cassandra にくらべて良い結果となっている。特に読み込みに関しては Consistency Level QUORUM の場合と比べると 3 倍以上離れている場合もある。実験 1 に比べて Jungle と Cassandra の差が開いているのは Cassandra の Consistency Level が QUORUM に設定されていることが要因の 1 つとしてあげられる。今回 Cassandra の Replication factor は 5 と設定している。そのため、Consistency Level が QUORUM の場合は、書き込みは 3 つのノードに書き込まれたことを確認し、読み込みは 3 つのノードからデータを取得して最新のデータを返すのである。Jungle の結果が横軸の値が増えても横ばいになっていることにも注目したい。これは Jungle の場合、リクエストが来た際に、それぞれのノードがローカルにある木の情報をすぐに返すためである。そのため、クライアントが増え、総リクエスト数が増加しても一台に対する負荷が増えない限りは同じレスポンス速度を維持できる。

第7章 結論

7.1 まとめ

本研究では、まず初めにRDBとNoSQLの説明を行い、既存のNoSQLであるCassandra, MongoDB, Neo4jがスケーラビリティをどのように確保しているかを述べた。次に木構造データベースJungleで使われている非破壊的木構造について述べ、破壊的木構造に比べロックが少ないというメリットがあることを論じた。Jungleは非破壊的木構造により過去のデータを保持することでMergeを行うことができる等、分散管理システムと似た分散設計が行えることを述べた。

また、Jungleの分散設計では当研究室で開発している並列分散フレームワークAliceを用いた。Aliceにより自由にトポロジを組むことができ、他サーバノードへのデータアクセス機構を手に入れることができた。Jungleの分散実装ではデータの編集履歴であるTreeOperationLogをAliceが使用できるようにし、木の名前といった必要な情報を追加することでデータの分散を行った。さらに、Jungleに元々設計されていたJournalを使ってログをディスクへ書き出すことで永続性の実装を行った。

最後に簡易掲示板を作成し、Cassandraとの性能比較を行った。読み込み、書き込みの負荷をかける実験を2つ行った。1つの実験ではサーバノード1台に対し複数のクライアントから負荷をかけた。2つめの実験では複数のクライアントに対し同じ数のサーバノードを用意し数を増やしていき負荷を高めた。どちらの実験もJungleがCassandraよりも良い結果を示すことを確認した。

7.2 今後の課題

7.2.1 pull/push方式による分断耐性の実装

現在の実装のJungleは、プログラムの起動時に複数ノードが接続をしトポロジを形成する。プログラムの途中で接続がきれるとトポロジがくずれたままになる。接続がきれたJungleは単独では稼働し続けるが、トポロジへの復帰を行えるようにしたい。そのためには、今の実装で行っている非同期でログを送信する方法とは別に同期をとり差分のデータを伝搬する方法を実装する必要がある。これは、分散管理システムにおけるpull/push APIの機能にあたる。

7.2.2 Merger アルゴリズムの設計

Jungle は Merge を使うことでデータ衝突の問題を解決をはかるが、この Merge はアプリケーション毎に考えなければならない。今回、Jungle における Merge の例として掲示板プログラムにおける Merge について述べた。だが掲示板のような単純な Merge ですむアプリケーションは少ない。また、アプリケーション毎でデータの保存の仕方といったものも違って来る。そのため、アプリケーションに合った Merge アルゴリズムを設計しなければならない。

7.2.3 過去のデータの掃除について

Jungle は非破壊でデータを保持し続けるため、非常に多くのメモリを使用してしまう。ある程度の単位で過去のデータの掃除を行いたい。今回分散実装を行ったことで、多数のノードでデータが保持され、その内の数台がディスクへ書き出すといったことも可能になった。しかし、Merge の問題含め、どのタイミングで過去のデータを掃除すべきかは自明ではない。分断耐性の実装の問題とも関わってくるが、どのデータがどれだけ複製して持っているといった情報も扱う必要がでてくるかもしれない。

謝辞

本研究を行うにあたりご多忙にも関わらず日頃より多くの助言、ご指導をいただきました河野真治准教授に心より感謝いたします。

研究を行うにあたり、並列計算環境の調整、意見、実装に協力いただいた谷成 雄さん、杉本 優さん、並びに並列信頼研究室の全てのメンバーに感謝いたします。

また、本研究は、JST/CREST 研究領域「実用化を目指した組み込みシステム用ディペンドブル・オペレーティングシステム」D-ADD 研究チームとして実施された。様々な研究や勉強の機会を与えてくださった、株式会社 Symphony の永山辰巳さん、同じく様々な助言を頂いた森田育宏さんに感謝いたします。様々な研究に関わることで自身の研究にも役立てることが出来ました。

最後に、大学の修士まで支えてくれた家族に深く感謝します。

参考文献

- [1] 所眞理雄. DEOS プロジェクト研究成果集 Dependability Engineering for Open Systems, 2013.
- [2] 永山 辰巳, 横手 靖彦. オープンシステムディペンダビリティと D-Case を繋ぐリポジトリ, 2013.
- [3] Nancy Lynch and Seth Gilbert. Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. *ACM SIGACT News*, 2002.
- [4] Avinash Lakshman and Prashant Malik. Cassandra - a decentralized structured storage system. *LADIS*, Mar 2003.
- [5] Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Voshall, and Werner Vogels. Dynamo: Amazon’s highly available key-value store.
- [6] Messagepack. <http://msgpack.org/>.
- [7] 大城信康, 河野真治. Data segment の分散データベースへの応用. 日本ソフトウェア科学会, September 2013.
- [8] 玉城将士, 河野真治. Cassandra を使った cms の pc クラスタを使ったスケーラビリティの検証. 日本ソフトウェア科学会, August 2010.
- [9] 玉城将士, 河野真治. Cassandra を使ったスケーラビリティのある cms の設計. 情報処理学会, March 2011.
- [10] 玉城将士, 河野真治. Cassandra と非破壊的構造を用いた cms のスケーラビリティ検証環境の構築. 日本ソフトウェア科学会, August 2011.
- [11] Fay Chang and Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Mike Burrows, Tushar Chandra, Andrew Fikes, and Robert E. Gruber. Bigtable : A distributed storage system for structured data.
- [12] Matt Welsh. The staged event-driven architecture for highly-concurrent server applications.

- [13] Eric Brewer Matt Welsh, David Culler. Seda : An architecture for well-conditioned , scalable internet services. *SOSP*.

発表履歴

- Java による授業向け画面共有システムの設計と実装, 大城信康, 谷成雄 (琉球大学), 河野真治 (琉球大学), オープンソースカンファレンス 2011 Okinawa, Sep, 2011
- Continuation based C の GCC 4.6 上の実装について, 大城信康, 河野真治 (琉球大学), 第 53 回プログラミング・シンポジウム, Jan, 2012
- GraphDB 入門 TinkerPop の使い方, 大城信康, 玉城将士 (琉球大学), 第 15 回 Java Kuche, Sep, 2012
- ディペンダブルシステムのための木構造を用いた合意形成データベースの提案と実装, 大城信康, 河野真治 (琉球大学), 玉城将士 (琉球大学), 永山 辰巳 (株式会社 Symphony), 情報処理学会システムソフトウェアとオペレーティング・システム研究会 (OS), May, 2013
- Data Segment の分散データベースへの応用, 大城信康, 杉本優 (琉球大学), 河野真治 (琉球大学), 日本ソフトウェア科学会 30 回大会 (2013 年度) 講演論文集, Sep, 2013