

卒業論文 2019 年度 (令和元年)

惑星規模の分散システムのための試験環境の設計と構築

慶應義塾大学 環境情報学部
廣川昂紀

惑星規模の分散システムのための試験環境の設計と構築

本研究では、惑星規模の分散システムのための試験環境の構築手法を提案する。本研究における惑星規模の分散システムとは、分散システムの中でも世界中に地理的に分散したコンピュータによって構成されるものと定義する。惑星規模の分散システムを支える技術として、P2P があげられる。P2P システムは、中央集権的なサーバを必要とせず、お互いに対等な関係をもつコンピュータが協調動作することによって成り立つ。ブロックチェーンはその一例である。惑星規模の分散システムにおける試験では、地理的に分散することによるネットワークでの通信の遅延を考慮する必要がある。よって、試験環境は地理的に分散配置されたサーバで構成されるべきである。

既存の惑星規模の分散システムの試験環境としては、PlanetLab やパブリッククラウドサービスの活用、BSafe.network があげられるが、それぞれ課題があると考ええる。PlanetLab はネットワークサービスの開発を支援する研究ネットワークであり、世界中の 717 地域 1353 のサーバに対して SSH を通して操作を行えるが、OS や CPU、メモリなどのサーバの環境を柔軟に変更できない。パブリッククラウドサービスでは、リージョンと呼ばれるデータセンターの地域を指定することでサーバを分散配置できるが、リージョンが限定的であることに加え、パブリッククラウドサービスのデータセンター周辺はネットワーク性能が高く公の実稼働環境に比べ通信の遅延が少ない。BSafe.network は 32 の大学によって構成されるブロックチェーン技術の研究を行うためのネットワークであり、各大学が保有するサーバを用いて実験が行えるが、各サーバの管理権限が分かれているため複数の大学間で共同研究を行う場合はオペレータの手作業が必要である。これらの課題を解決するため、本研究では OpenVPN と Kubernetes を組み合わせた統合的試験環境の構築手法を提案する。本システムでは、サーバがどこに分散配置されていてもネットワークでの通信の遅延を考慮した試験が可能である。加えて、仮想化技術を活用することでひとつのサーバ内に複数のアプリを様々な環境下で動作させることができ、各サーバが別々の管理下にある場合でも統合可能である。

本システムの実装後、すべてのサーバの統合管理が可能であるか、ならびに通信の遅延を考慮した上で本システムが正常に動作するかを検証した。

本研究は、惑星規模の分散システムの試験をより公の実稼働環境に近い形で柔軟に行うことを可能とし、BSafe.network のように各サーバが別々の管理下にある場合でもオペレータの手作業を省いた統合的試験環境の構築が可能である。

キーワード:

1. 惑星規模の分散システム, 2. 試験環境, 3. OpenVPN, 4. Kubernetes

慶應義塾大学 環境情報学部
廣川昂紀

Design and Implementation an experiment environment for planetary-scale distributed systems

In this study, we propose a system design for an experiment environment for planetary-scale distributed systems. The planetary-scale distributed system in this study is a distributed system consists of geographically distributed computers. P2P is an example of planetary-scale distributed systems. P2P systems do not require centralized servers and consist of computers that have an equal relationship with each other to cooperate. Blockchain is one of such P2P systems. Experiments for planetary-scale distributed systems should consider network latency. Therefore, experiment environments for planetary-scale distributed systems should consist of geographically distributed computers.

There are several services such as PlanetLab, public cloud services, and BSafe.network can be used as an experiment environment for planetary-scale distributed systems, but all of them have shortcomings. PlanetLab is a research network that supports experiments of network services. It has 1353 servers that are able to be controlled via secure shell in 717 areas around the world, but can not flexibly change the environments of the servers such as OS, CPU, and memory. In a public cloud service, servers can be geographically distributed by specifying one area called region. However, the scope of regions are limited and it is hard to consider control delay, because network performances around regions are high. BSafe.network is a network for researching blockchain technology that consists of thirty-two universities around the world. In BSafe.network, developers can use the servers owned by each university to experiment. However, the manual operations are needed for joint research between universities, because the management authority of the each server is different. To solve these issues, we propose a system design for an integrated experiment environment by combining two softwares: OpenVPN and Kubernetes. In this system, it is possible to run experiments even with control delay, regardless of where servers are geographically distributed. Besides, several applications can be run on a single server by using virtualization technology, and can be integrated even if each server is under separate management.

We verified whether integrated management of servers is possible, and whether this system considers control delay.

This research makes it possible to run experiments for planetary-scale distributed systems more flexibly in a similar manner to a public production environment. Even if each server is under separate management such as BSafe.network, It is possible to implement an integrated experiment environment without any manual operators.

Keywords :

1. Geographically Distributed System, 2. Staging Environment, 3. OpenVPN, 4. Kubernetes

Keio University Faculty of Environment and Information Studies
Koki Hirokawa

目 次

第 1 章	序論	1
1.1	背景	1
1.1.1	惑星規模の分散システム	1
1.1.2	惑星規模の分散システムの発達	2
1.1.3	試験環境	2
1.2	課題と目的	3
1.2.1	課題	3
1.2.2	目的	3
1.3	仮説	3
1.4	提案手法	4
1.5	本論文の構成	4
第 2 章	背景	5
2.1	惑星規模の分散システム	5
2.1.1	分散システム	5
2.1.2	惑星規模の分散システム	5
2.2	惑星規模の分散システムにおける使用技術と参考例	6
2.2.1	P2P	6
2.2.2	Winny	9
2.2.3	Gnutella	9
2.2.4	Bitcoin	9
2.3	試験環境	10
2.3.1	惑星規模の分散システムの試験環境	10
2.4	コンテナオーケストレーションシステム	11
2.4.1	コンテナ	11
2.4.2	Kubernetes	15
2.5	OpenVPN	18
2.5.1	OpenVPN	19
第 3 章	本研究における課題定義と仮説	21
3.1	課題定義	21
3.2	課題解決における要件	21
3.2.1	サーバ環境の柔軟性	22

3.2.2	ネットワークでの通信の遅延の考慮	22
3.2.3	異なる組織間での試験における手作業の軽減	22
3.2.4	地理的に分散したサーバの統合的な管理	22
3.3	先行研究	22
3.3.1	P2P アプリケーションの開発と性能評価のための統合開発環境の提案	23
3.3.2	プロセスレベルの仮想化を用いた大規模分散システムテストベッド	23
3.4	仮説	23
3.4.1	サーバ環境の柔軟性	23
3.4.2	ネットワークでの通信の遅延の考慮	24
3.4.3	異なる組織間での試験における手作業の軽減	24
3.4.4	地理的に分散したサーバの統合的な管理	24
3.5	提案システム概要	24
第 4 章	実装	26
4.1	実装環境	26
4.1.1	ハードウェアおよびソフトウェア	26
4.1.2	物理サーバの準備	26
4.1.3	ネットワーク構成	27
4.1.4	VM の配置	27
4.1.5	ルーターの配置	30
4.1.6	OpenVPN の設定	30
4.1.7	Kubernetes クラスタの構築	31
4.2	システム全体	34
第 5 章	評価	36
5.1	地理的に分散したサーバに対する統合管理	36
5.2	オペレータの手作業の軽減	37
5.3	通信の遅延による本システムへの影響	38
5.4	評価のまとめ	39
第 6 章	結論	41
6.1	まとめ	41
6.2	課題と展望	42
第 7 章	謝辞	43
	謝辞	43

目 次

2.1	分散システム	6
2.2	クライアントサーバ型と P2P	7
2.3	試験環境	10
2.4	アプリケーションのひとつがリソースを大幅に占有した場合	12
2.5	ひとつの物理サーバでひとつのアプリケーションを動作させる解決策	12
2.6	VM の相互独立性	13
2.7	VM の構成	13
2.8	コンテナ型仮想化	14
2.9	Docker のロゴ	15
2.10	Docker でのコンテナ作成手順	15
2.11	Kubernetes のロゴ	16
2.12	Kubernetes でのコンテナデプロイ	17
2.13	VPN	19
2.14	OpenVPN のロゴ	19
3.1	システム概要図	25
4.1	ネットワーク構成	28
4.2	OpenVPN によるサイト間接続	32
4.3	Kubernetes マスタークラスターの構築	33
4.4	マスターノードとワーカーノードの関係性	33
4.5	ネットワーク構成図	35
5.1	tc コマンドによるマスターノード・ワーカーノード間の擬似的な遅延の発生	39

表 目 次

4.1	使用したハードウェアおよびソフトウェア	26
4.2	ESXi の IP アドレス	27
4.3	Vlan ID と対応するアドレスプレフィックス	27
4.4	設置した VM の詳細	29
4.5	設置したルーターの詳細	30
4.6	OpenVPN 設定前の各サーバの疎通性	30
4.7	OpenVPN 設定前の各サーバの疎通性	31
5.1	通信の遅延によるコンテナ起動時間の変化	39

第1章 序論

本研究では，惑星規模の分散システムのための試験環境の設計と構築を行う．

本章では，惑星規模の分散システムを定義し，本研究の背景である惑星規模の分散システムの発達とシステムの試験環境について概説する．本研究の課題を明らかにした上で，目的を明確化し，目的を達するための仮説を示す．最後に仮説を裏付けるための提案手法を示し，本研究の概要を示す．

1.1 背景

本節では，本研究の背景について述べる．

初めに，本研究が対象とする惑星規模の分散システムと試験環境について概説する．分散システムについて概説した上で惑星規模の分散システムを定義し，惑星規模の分散システムの発達について述べる．加えて，惑星規模の分散システムに求められる試験環境について説明する．

1.1.1 惑星規模の分散システム

本節では，本研究における惑星規模の分散システムを定義する．

分散システムは，複数の構成要素が組み合わさって動作するシステムである．各構成要素は独立しており，互いに協調動作することによってシステム全体が成り立っている．

惑星規模の分散システムは，分散システムの中でも世界中に地理的に分散したコンピュータによって構成されるものと定義する．惑星規模の分散システムの基盤技術としては，P2Pがあげられる．

P2Pは“Peer to Peer”の略記であり，中央集権的なサーバを必要とせず，各コンピュータが互いに対等な関係を築き協調動作することで成り立つシステムモデルまたはその技術自体を指す．P2Pは，システム内に明確な役割分担と主従関係のあるクライアントサーバモデルとは対照的なシステムモデルである．クライアントサーバモデルは中央主権的なシステムであり，通信において常にクライアントとサーバで一対一の関係が成り立つ．クライアントはサーバに対し処理や情報を要求し，要求を受け取ったサーバは特定の処理を行なってクライアントへ返答する．対してP2Pシステムでは，システムを構成する各コンピュータの役割は状況に応じて柔軟に変化し，通信において多対多の関係が成り立つ．クライアントとして他のコンピュータに対し要求する場合もあれば，他のコンピュータか

らの要求に対して応答する場合もある。クライアントサーバモデルに比べて、拡張性（スケーラビリティ）ならびに耐障害性において優れているのが特徴的である。

ビットコインの中核技術であるブロックチェーンは、P2P システムの一例である。ブロックチェーンのような P2P システムでは、地理的に分散することによるネットワークでの通信の遅延を考慮した上で協調動作可能であることを開発者は意識しなければならない。

1.1.2 惑星規模の分散システムの発達

2000 年代初頭、Winny [1] や Gnutella [2] といった惑星規模の分散システムが頭角を現した。どちらのサービスも P2P 技術を基盤としており、それまでシステムモデルとして一般的であったクライアントサーバモデルとは異なる形態を採用したことで注目が集まった。P2P 技術が研究分野で取り上げられる頻度も多くなり、サービスとしても今後一層幅を広げていくと思われたが、クライアントサーバモデルに置き換わるまでの隆盛はなく後退していった。しかし、2008 年に Satoshi Nakamoto により Bitcoin のために開発されたブロックチェーン技術が登場することによって、再度 P2P 技術が脚光を浴びるようになり、開発や研究の勢いが再び盛んになってきている。

1.1.3 試験環境

試験環境とは、公の実稼働環境での運用をする前にシステム全体の試験を行うための環境である。開発者が実際に開発を行う開発環境と公の実稼働環境では、環境の差異から動作の違いが生じ、開発環境で正常に動作していたものが公の実稼働環境に反映した途端動作しなくなるといった事象が度々発生する。そのような事態を防ぐために開発環境と公の実稼働環境の間に試験環境を構築し、公の実稼働環境への適用前に試験環境にてシステム全体の試験をすることで予想外の障害が発生する可能性を低減できる。惑星規模の分散システムにおいても、システムの不具合を早期に発見するために試験環境が必要である。

先に述べたように、惑星規模の分散システムでは地理的に分散することによるネットワークでの通信の遅延が発生する。よってシステムの試験は、通信の遅延を考慮した上で正常に協調動作することを確認するため、試験環境は地理的に分散したサーバによって構成されたものでなければならない。

既存の惑星規模の分散システムの試験環境としては、PlanetLab やパブリッククラウドサービスの活用、BSafe.network があげられる。PlanetLab は、ネットワークサービスの開発を支援する研究ネットワークであり、世界中の 717 地域 1353 のサーバを利用することができ、パブリッククラウドサービスでは、リージョンと呼ばれるデータセンターの地域を指定することでサーバを分散配置することが可能である。最後に、BSafe.network は 32 の大学によって構成されるブロックチェーン技術の研究を行うためのネットワークであり、世界中の各大学が保有するサーバを用いて開発を行える。BSafe.network はブロックチェーンの実験環境であるが、ブロックチェーンの試験を行う場でもあるため以後試験環境と捉える。

1.2 課題と目的

本節では、本研究の課題と目的を述べる。

まず、1.1 章で述べた本研究の背景を元に、惑星規模の分散システムの試験環境における課題を示す。その上で本研究の目的を明確にする。

1.2.1 課題

本研究では、既存の惑星規模の分散システムの試験環境の課題について指摘する。

惑星規模の分散システムの試験では、地理的に分散配置されることによるネットワークでの通信の遅延を考慮した上で各コンピュータが正常に協調動作を行えるかを確認する必要がある。試験環境では各サーバを実際に地理的に分散配置する必要性があることは 1.1 章の背景で述べた通りである。さらに既存の試験環境として、Planet Lab やパブリッククラウドサービスのリージョンの活用、BSafe.network があげられるが、それぞれに課題があると考ええる。Planet Lab では世界中に分散したサーバを利用することが可能だが、OS や CPU、メモリなどのサーバの環境を柔軟に変更することができない。パブリッククラウドサービスでは、サーバの環境を自由に変更可能であり、リージョンを活用して地理的に分散した場所にサーバを設置することができるが、リージョンが限定的であり、データセンター周辺はネットワーク性能が高く公の実稼働環境に比べ通信の遅延が少ない。BSafe.network では、世界中の 32 の大学が保有するサーバを用いてシステムの実験を行えるが、各サーバの管理権限が各大学のオペレータに委ねられているため、大学間での共同研究を行う場合にオペレータの手作業が必要である。

本研究では、このように惑星規模の分散システムの試験環境の構築手法が整備されていないことを課題とする。

1.2.2 目的

本研究では、惑星規模の分散システムのための試験環境の構築手法を提案することを目的とする。

1.3 仮説

1.2.1 章で述べた課題を解決するため、本研究では地理的に分散したサーバを統合管理可能な試験環境の構築が必要であると考えた。惑星規模の分散システムの試験環境は、

- OS や CPU、Memory といったサーバ環境を柔軟に変更可能であること
- 公の実稼働環境を想定したネットワークでの通信の遅延を考慮できること
- 異なる管理権限下にある各サーバに対し統合的管理が可能であり、各オペレータの手作業を軽減できること

の三点を満たさなければならない。

本研究では、上記の必要要件を満たすことで 1.2.1 で述べた課題点を解決し、1.2.2 で述べた惑星規模の分散システムの試験環境の構築手法の提案を達成できると考えた。

1.4 提案手法

本研究では、1.3 章で述べた必要要件を満たすため、OpenVPN と Kubernetes を組み合わせた惑星規模の分散システムの統合的試験環境を提案する。

Kubernetes はコンテナオーケストレーションツールであり、コンテナ化仮想技術によってコンテナ化されたアプリケーションのデプロイやスケーリングを自動化し、統合管理するためのシステムである。Kubernetes では複数のサーバでクラスタを構成しており、クラスタリングを行うためには各サーバが互いに IP レベルで疎通可能な状態になければならない。よって、IP レベルでの疎通が取れない別々のセグメントに配置されたサーバ間では Kubernetes クラスタを構築することはできない。

そこで地理的に分散し異なるセグメントに配置されたサーバ間を繋ぐ OpenVPN オーバーレイネットワークを構築することで、各サーバを互いに IP レベルで疎通可能にする。OpenVPN は、VPN ネットワークの構築をソフトウェアで実現するために開発されたオープンソースソフトウェアである。

本研究では、OpenVPN と Kubernetes を組み合わせ、地理的に分散した拠点間で形成した OpenVPN オーバーレイネットワーク上で Kubernetes クラスタを構築した。本システムが本研究における課題点を解決できているか推定することで、要件を満たせることを確認した。

1.5 本論文の構成

本論文における以降の構成は次の通りである。

2 章では、惑星規模の分散システムと試験環境ならびに本研究で使用する技術について概説し、本研究の背景を明確化する。3 章では、本研究における課題を明確化し、課題を解決するための要件、仮説と手法について概説する。4 章では、本研究で提案する試験環境の構築方法について述べる。5 章では、3 章で述べた課題に対しての評価を行い、考察する。6 章では、本研究のまとめと今後の課題についてまとめる。

第2章 背景

本章では，本研究の背景について概説する．

本研究における惑星規模の分散システムを定義し，惑星規模の分散システムの基盤技術である P2P について概説する．加えて，本研究で対象とする試験環境の定義と惑星規模の分散システムに必要な試験環境について述べる．最後に，本研究の提案手法で用いるコンテナオーケストレーションシステムと OpenVPN について概説する．

2.1 惑星規模の分散システム

本節では，本研究が対象とする惑星規模の分散システムを定義する．惑星規模の分散システムについて概説する前に，分散システムについて詳細を述べ，違いを明らかにした上で本研究における惑星規模の分散システムを定義する．

2.1.1 分散システム

分散システムとは，複数の構成要素が協調動作することによって成り立つシステムを指す．各構成要素は独立して異なる役割を持っており，それらが組み合わさり互いに協調動作することによってシステム全体が動作する．本研究で使用する Kubernetes も，コンテナオーケストレーションに必要な機能を構成要素毎に分割した分散システムである．Kubernetes に関する説明は 2.4.2 章で詳しく行う．分散システムは複雑に思われるが，機能が細分化され各構成要素の役割が明確化されるメリットがある．機能同士の依存関係が希薄になるため細かい粒度での試験が可能となり，障害時の原因特定も行いやすい．チーム開発において開発者同士で担当範囲が重複する可能性も下がるため，開発スピードが向上するというメリットもある．

2.1.2 惑星規模の分散システム

本研究における惑星規模の分散システムとは，分散システムの中でも世界中に地理的に分散したコンピュータが協調動作することによって成り立つシステムを指す．近年注目を集めているブロックチェーンや 2000 年代初頭に登場した Winny といったサービスが，惑星規模の分散システムの一例である．惑星規模の分散システムを支える技術や例についての概説は 2.2 にて行う．惑星規模の分散システムでは，システムを構成するコンピュータ

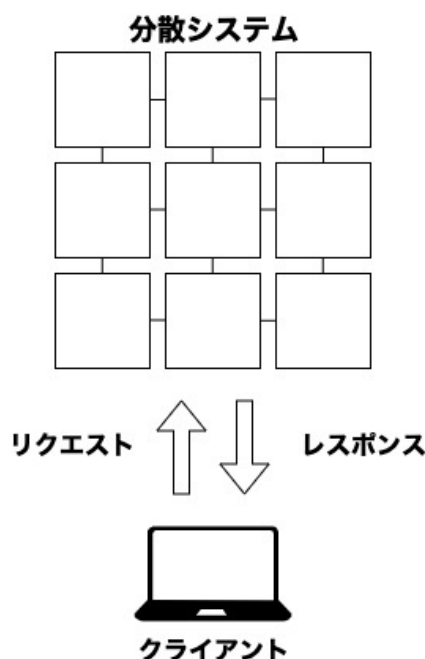


図 2.1: 分散システム

が地理的に分散していることから、開発者はネットワークでの通信の遅延を考慮する必要がある。

2.2 惑星規模の分散システムにおける使用技術と参考例

本節では、2.1.2 章で概説した惑星規模の分散システムを支える技術とサービス例について概説する。惑星規模の分散システムの基盤技術として、P2P があげられる。P2P システムは中央集権的なサーバを必要としないクライアントサーバモデルとは対照的なシステムモデルである。P2P について概説したのち、P2P システムを採用するサービスとして、Winny, Gnutella, Bitcoin を紹介する。

2.2.1 P2P

P2P は “Peer to Peer” の略記である。P2P は、クライアントサーバモデルのシステムのように中央集権的な役割を担うサーバを必要とせず、コンピュータ同士が対等な関係を築く主従関係のないシステムモデルである。またはそれを実現する技術を指す。

クライアントサーバモデルでは、通信においてサーバとクライアントで常に一対一の関係性が成り立つ。また、システム内では処理や情報を要求するクライアントと要求に対し応答するサーバで明確な役割分担がある。サーバはクライアントから要求が送られてきた際は、特定の処理を行いクライアントに対し情報を返すが、それ以外では待機状態とな

る。対してクライアントは、情報を要求したり変更してもらう必要が生じた際のみサーバと通信を行う。よってクライアントサーバモデルにおける通信は、基本的にクライアントが起点となって行われる。

P2P では各コンピュータが互いに対等な関係性を築くため、クライアントサーバモデルのような明確な役割分担がシステム上ない。P2P では、各コンピュータが状況に応じてサーバとクライアントの役割を担う。各コンピュータが臨機応変にサーバとして要求に応答し、クライアントとして処理や情報を要求する動的システムが特徴としてあげられる。

クライアントサーバモデルでは、要求する側をクライアント、対して要求に応じる側をサーバと呼んでいる。P2P では前述した通り各コンピュータは動的に役割を変化させ、サーバとしてもクライアントとしても動作することからサーバントと呼ばれる。または単にノードと呼ばれることもある。

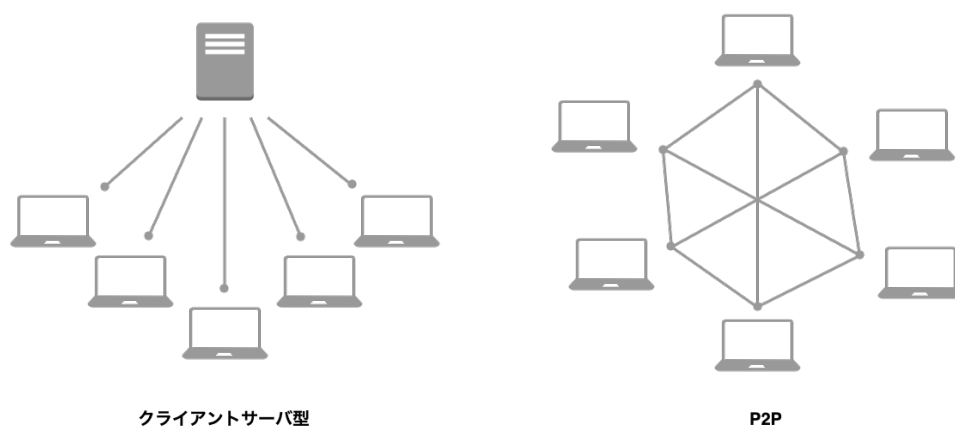


図 2.2: クライアントサーバ型と P2P

P2P の特徴

P2P では各コンピュータがサーバにもクライアントにも成り得るため、クライアントサーバモデルとは内部の実装も異なる。

第一に、データを保持する中央集権的なサーバが存在しないためアプリケーション上で必要になるデータは各コンピュータが保持する。アプリケーションの実装方式によっても異なるが、各コンピュータがデータを分割して保持する場合もあれば全てのコンピュータが同じデータを保持する場合もある。例えばファイル共有システムである Winny では、各コンピュータの保持しているデータは異なるため、データを参照する際はどのコンピュータが目的のデータを保持しているか検索し、対象となるコンピュータを決定してから通信を行う必要がある。また、ブロックチェーンでは各コンピュータが全てのデータを保持しており（全てのデータを持たない場合もある）、データを相互に検証し合うことによってデータの改竄耐性を向上させ、堅牢性を担保している。

次に、システムのメインプログラムを各コンピュータが保持し動作させなければいけない点でもクライアントサーバモデルとは異なる。クライアントサーバモデルでは、システムのメインプログラムの実行はサーバの役割であるため、サーバのみがシステムのメインプログラムを保持しておけば良い。対して、各コンピュータが状況に応じて役割を変える P2P ではシステムのメインプログラムを各々で保持する必要がある。クライアントとして他のコンピュータが保持しているデータを参照したり、情報を要求してきたコンピュータに対して応答をしなければならないからである。

P2P のメリット

本節では、P2P のメリットについて概説する。P2P システムの利点としては、拡張性（スケーラビリティ）・耐障害性があげられる。

第一に拡張性に関しては、クライアントサーバモデルの場合、利用者が増大するとシステムを中心であるサーバへアクセスが集中し、サーバやその周辺のネットワークへの負荷が高くなり、システム的な弱点になる。システム運用者は拡張性を高めるため、ネットワーク機器のスペックをあげたり、負荷が増大した際に自動でサーバの数を増やすオートスケーリングなどの対策を取らなければならない。それに対して P2P の場合、コンピュータ同士は相互に通信を行うためアクセスは分散されやすくなる。その点で P2P は拡張性に長けている。

次に耐障害性である。クライアントサーバモデルの場合、何らかの原因でサーバが落ちるとシステム自体が停止してしまいサーバが構造上の単一障害点となる。しかし、P2P ではあるコンピュータが停止した場合でも、正常なコンピュータ同士で新たなネットワークを形成することで滞りなくシステムの動作を継続することができる。構造上の単一障害点が存在しないため、障害性に長けている。

P2P のデメリット

本節では、P2P のデメリットについて概説する。

第一に情報伝達における遅延があげられる。P2P では接続先のコンピュータが常に決まっていないため、状況に応じて接続先を変更する必要がある。すなわち、目的の情報を保持しているコンピュータを探し出したり、そもそもネットワーク上で近い距離に他のコンピュータが存在しない場合、情報の取得や送信に遅延が生じてしまう。全てのコンピュータで同じデータを保持するブロックチェーンのようなシステムにおいては、コンピュータ同士がバケツリレーのようにデータを受け渡さなければならず、端から端までデータを伝えるまでに時間が掛かってしまう問題点がある。

次にシステム全体での管理のしにくさがあげられる。P2P システムでは各コンピュータでアプリケーションを動作させるため、中央集権的なサーバと異なり、管理は各コンピュータ管理者に委ねられることになる。よって、システムに問題点が見つかり開発者が修正を含んだ更新版を配布した場合でも、実際に動作しているアプリケーションが更新されるかどうかは保証されない。同様にシステム全体の監視を行うことも困難である。

2.2.2 Winny

Winny [1] はソフトウェアエンジニア金子勇氏が開発し、2002 年に発表されたファイル共有ソフトである。システム上で中央集権的なサーバを保持せず、コンピュータ同士が相互に接続することで実現される P2P アプリケーションとして注目を浴びた。ユーザはコンピュータ内に保持されたファイルを他のコンピュータと共有することができるため、任意のファイルをアップロードしたり、逆に他のコンピュータが保持しているファイルをダウンロードすることができる。Winny では、受信ファイルの送信元や送信ファイルの宛先をユーザが確認することはできず、バックグラウンドでの処理はユーザに見せないように高い秘匿性が担保されていた。クライアントサーバモデルのシステムアーキテクチャとは打って変わって出た新しい形のアプリケーションであったが、高い匿名性も起因して、一部のユーザが違法な音楽ファイルや動画ファイル、コンピュータウイルスを Winny にアップロードしたことで著作権法違反が問われた。開発者である金子氏にも疑いがかけられ 2004 年に逮捕、その後画期的な発明であった Winny も衰退していった。なお、金子氏は裁判を経て 2011 年に無罪となった。

2.2.3 Gnutella

Winny に同じく Gnutella [2] も中央集権型サーバに依存せず、コンピュータ間の通信のみでファイルの送受信を行うファイル共有アプリケーションである。ファイルと言えど、Gnutella では主に音楽ファイルが共有されていた。Gnutella は AOL（アメリカ・オンライン）社のチームが開発したものである。著作権保護の観点から、法的に違法性を問われ公開も開発もすぐに中止されてしまった。Gnutella では、最初のプログラムの起動時には、接続先を自動で認識できないためファイル交換や検索機能は使用できない。Gnutella のシステムへ参加したい場合は、メールや掲示板を通して他の Gnutella サーバの IP アドレスとポート番号を教えてもらい、自分の Gnutella サーバへ設定することで他のサーバとの通信を確立できる。通信確立後は、最初に接続した Gnutella サーバを通して他のサーバとも連携を取ることが可能となり、音楽ファイルの交換や検索を行うことができる。

2.2.4 Bitcoin

Bitcoin [3] は 2008 年に Satoshi Nakamoto と名乗る人物によって論文にて提唱されたものである。2009 年にはソフトウェアとして公開されており、今では多くのユーザに使用されている上、仮想通貨の先駆けとして他の仮想通貨を生む大きな起点となった。同時に、2000 年代後半に勢いを失っていた P2P システムの存在を再度世に知らしめ、開発の促進を促す起爆剤の役割を果たしたと考えられる。Bitcoin は基盤技術のひとつとして Winny や Gnutella と共通する P2P ネットワークを採用している。参加するコンピュータはそれぞれがシステム上のデータを保持し相互にデータを検証しあうことで、第三者的監視機関を必要とせずにデータの堅牢性を担保することが可能である。

2.3 試験環境

本節では，本研究で着目する試験環境について概説する．

試験環境とは，システムの試験を行うための環境である．開発したシステムを公の実稼働環境へ適用する前に，試験環境にて期待する動作を正常に行うか試験する．試験環境でシステム上の不備を発見した場合，実稼働環境への適応はせずに開発環境にて修正を行う．対して，試験環境でのシステムの正常な動作を確認できた場合は，実稼働環境への適用作業へ移行する．開発環境と実稼働環境の間に試験環境を設けることで，システムの予期せぬ不具合を早期に発見できる．

以下，惑星規模の分散システムの試験に必要な試験環境について概説する．

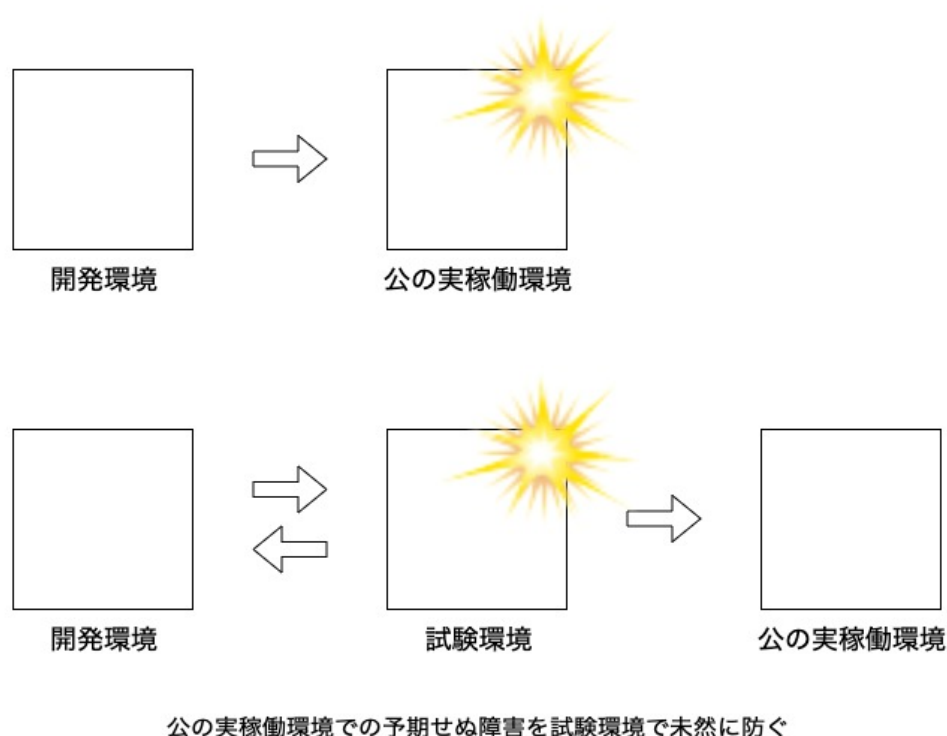


図 2.3: 試験環境

2.3.1 惑星規模の分散システムの試験環境

惑星規模の分散システムは，地理的に分散した構成要素によって成り立つシステムである．独立した構成要素が互いに協調動作することによってシステム全体が動作している．分散システムでは，構成要素単位での試験とシステム全体での協調動作の試験が必要である．なお惑星規模の分散システムの試験においては，システムを構成するコンピュータが地理的に分散していることによるネットワークでの通信の遅延の影響を考慮する必要がある．

る。通信の遅延を加味した上で分散システムの協調動作が正常に働いていることを確認できれば、公の実稼働環境においても障害発生の可能性が低くなる。よって、惑星規模の分散システムの試験環境では各コンピュータを地理的に分散させて配置させる必要がある。

惑星規模の分散システムの試験環境としては、PlanetLab やパブリッククラウドサービスの活用、BSafe.network があげられる。

PlanetLab は、ネットワークサービスの開発を支援する公の研究ネットワークであり、世界中の 717 地域 1353 のサーバを ssh を通して利用することができる。各サーバは仮想マシンを提供しており、研究者に割り当てられる仮想マシンは Slice と呼ばれる。2003 年から指導し、1000 人以上の研究者が分散ストレージ、ネットワークマッピング、P2P システム、分散ハッシュテーブルなどの新たな技術の開発のために PlanetLab を利用している。

GCP や AWS、Azure といったパブリッククラウドサービスでは、リージョンと呼ばれるデータセンターの地域を指定することでサーバの地理的な分散配置が可能である。

BSafe.network は 32 の大学で構成されるブロックチェーン技術の研究を行うためのネットワークである。BSafe.network は政治的かつ経済的に中立的な大学のみで構成される。研究を行う場合は、世界中の各大学が保有するサーバを用いて実験を行うことができる。

2.4 コンテナオーケストレーションシステム

コンテナオーケストレーションシステムは、コンテナ型仮想環境を統合管理するためのプラットフォームおよびツールを指す。2010 年代半ばから脚光を浴びようになり、今では世界的に数々のプロジェクトで実稼働環境に適用されている。サービスの立ち上げや運用過程において必要となる機能が数多く搭載されており、開発者は素早くかつ効率的に開発を進められる。コンテナは Virtual Machine（以下、VM）のデメリットを考慮して作られており、今後 VM の代わりを担う次世代の技術としてより一層注目されていく技術であると考えられる。

本研究では、コンテナオーケストレーションシステムとして Kubernetes を、CRI（コンテナ・ランタイム・インターフェース）として Docker を使用した。本節では、コンテナおよびコンテナオーケストレーションの概説と実際に使用した Kubernetes や Docker といったツールについて紹介する。

2.4.1 コンテナ

本節では、コンテナおよび Docker について概説する。

コンテナ型仮想化は、ひとつのコンピュータ上で仮想的に別のコンピュータを動作させる技術である。ホスト OS の上で動いている別のコンピュータをひとつひとつをコンテナと呼ぶ。

コンテナについて説明するにあたり、VM や物理マシンと比較しながら特徴を示していく。コンテナは、挙動としては VM と似ており、どちらも同じ課題を解決している。VM が登場する前、開発者はひとつのサーバ上で複数のアプリケーションを動作させることに

頭を悩ませていた。何故なら、アプリケーションのうちのひとつがサーバのリソースを大幅に占有した場合、他のアプリケーションのパフォーマンスが低下してしまうからである。

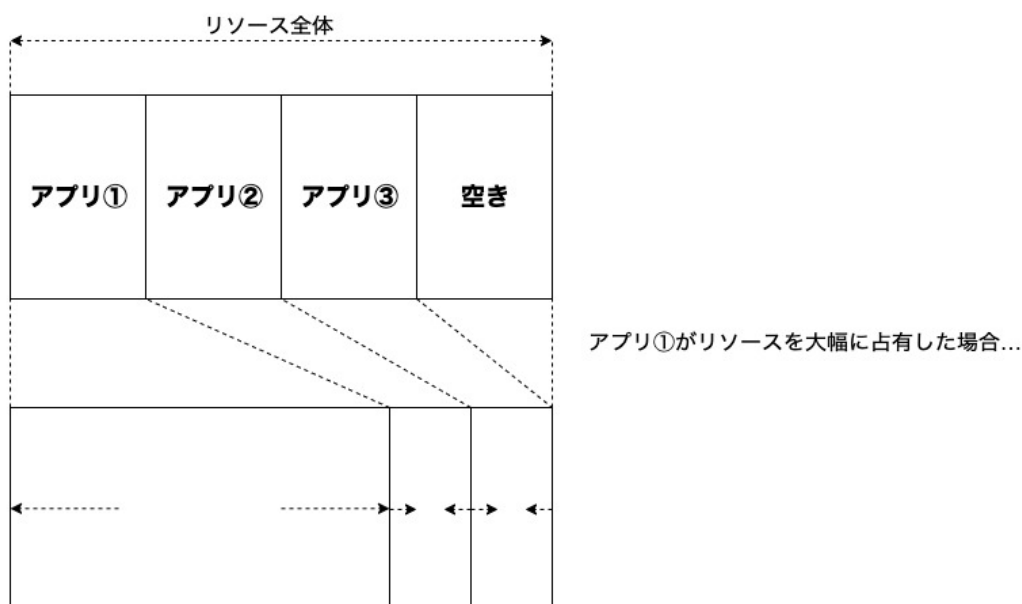


図 2.4: アプリケーションのひとつがリソースを大幅に占有した場合

解決策のひとつとして、アプリケーション毎に別々のサーバ上で動作させるものがあったが、デメリットとして維持費が嵩むことと使用されない無駄なリソースが生まれてしまうことがあった。

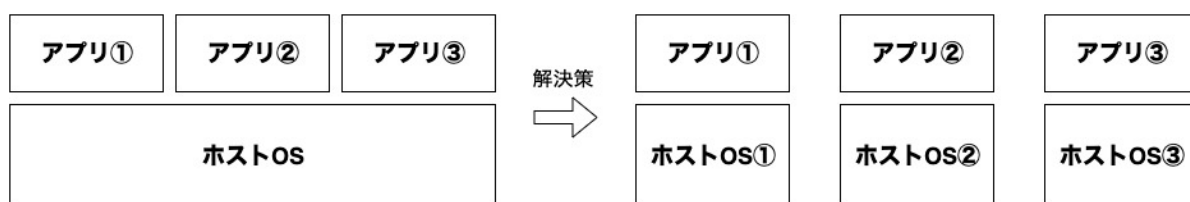


図 2.5: ひとつの物理サーバでひとつのアプリケーションを動作させる解決策

これを解決するために開発されたのが VM である。VM はソフトウェアによって仮想的に物理マシンを実現する技術であり、ひとつの物理マシン CPU 上で複数の VM を動作させることが可能である。アプリケーションはそれぞれ独立しておりお互いに不可侵な関係性であるため、ひとつのアプリケーションがリソースを占有することはなく、よりリソースを効率的に使用できる。スケーラビリティにも長けており、開発者はいつでもアプリを追加・削除でき、ハードウェアコストの削減にも貢献している。しかし、VM は処理におけるオーバーヘッドが大きく起動時間が長いなどデメリットも存在する。

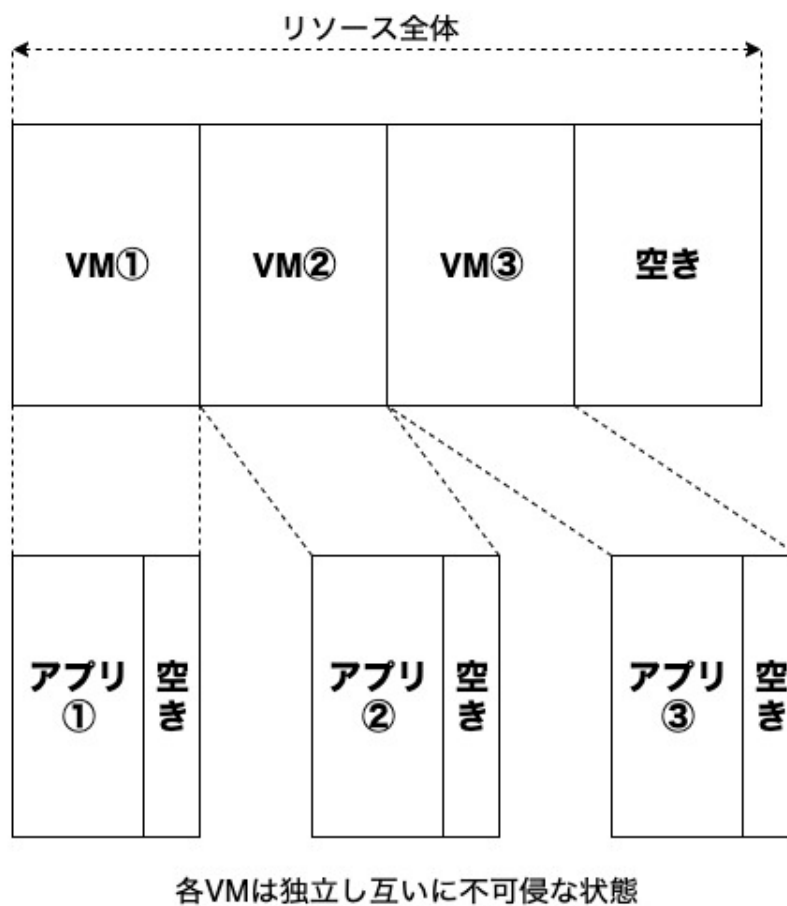


図 2.6: VM の相互独立性

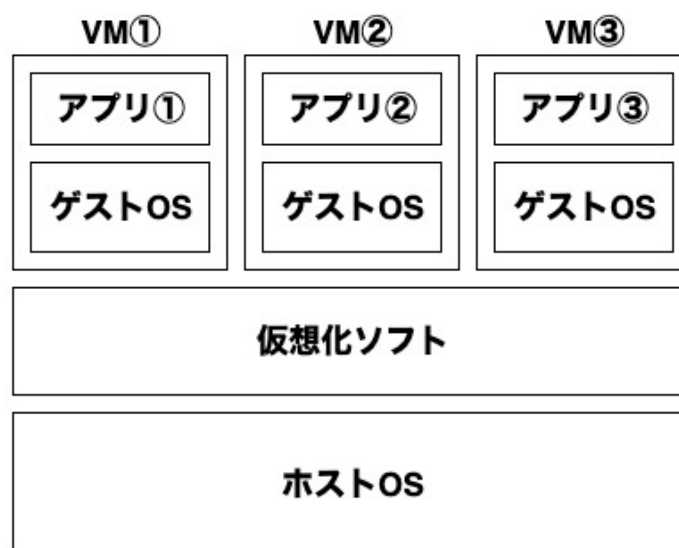


図 2.7: VM の構成

VM の後に登場した技術がコンテナ型仮想化である。コンテナ型仮想化では、各アプリケーションはひとつのホスト OS を共有するため、VM より軽量で起動時間も短い。コンテナはコンテナイメージから作成され、イメージは宣言的なファイルに基づいて生成される。これによって開発者はより簡単かつスピーディに開発を進めることが可能である。“Build Once, Run Anywhere”というコンセプトが掲げられており、一度生成されたイメージはどの環境でも動作し冪等性が担保される。一方、ホスト OS を共有するためセキュリティ面では課題が見られる。

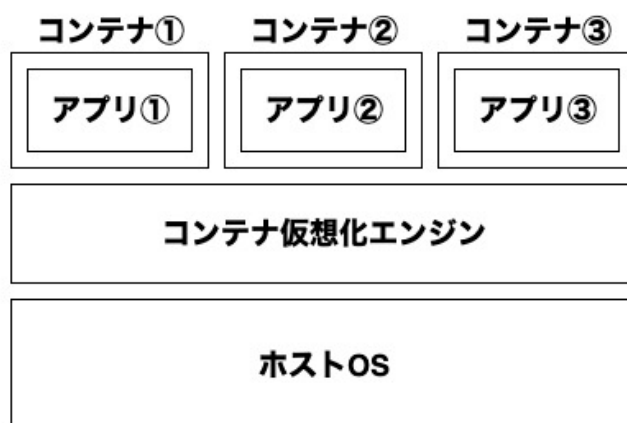


図 2.8: コンテナ型仮想化

コンテナ仮想環境を構築するためのランタイムである CRI には、dockershim (Docker), containerd, cri-o, Frakti, rktlet (rkt) などが挙げられる。本研究では、CRI のデファクトスタンダードである Docker を採用している。

Docker

Docker [4] はコンテナ型仮想環境を実現するためのプラットフォームおよびツールである。前述したように Docker では、宣言的なファイルから生成したコンテナイメージを元にコンテナを起動する。設計書となる宣言的なファイルは Docker ファイルと呼ばれる。Docker ファイルでは、ベースとなるイメージをインポートしたり、特定のコマンドの実行やファイルのコピーを行うためのコマンドが提供されている。ミドルウェアや各種環境設定をコード化して管理することができ (Infrastructure as Code), 別の環境で何度実行しても同じ結果が保証される。Docker イメージをバージョン毎に管理するための Docker Hub というサービスがあり、開発者は自身のレポジトリにイメージをプッシュしたり、他のレポジトリからイメージを取得することも可能である。

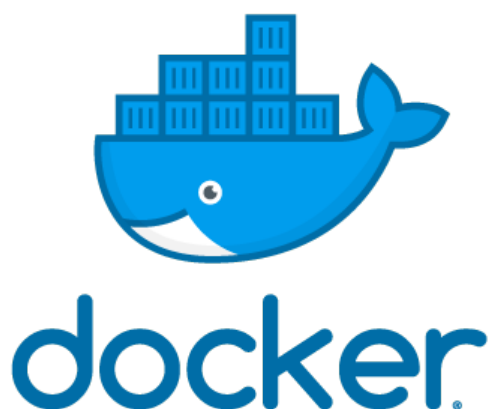


図 2.9: Docker のロゴ

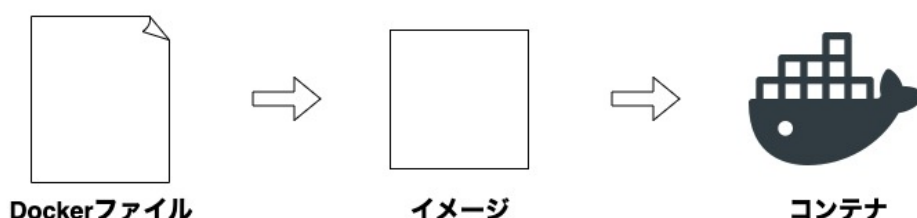


図 2.10: Docker でのコンテナ作成手順

2.4.2 Kubernetes

Kubernetes [5] はコンテナオーケストレーションエンジンであり，コンテナ化されたアプリケーションのデプロイやスケーリングなどの管理を自動化するためのプラットフォームである．

もともと Google 社内で利用されていたコンテナクラスタマネージャの「Borg」を基盤にして作られたオープンソースソフトウェアであるため信頼性が高く，現時点でコンテナオーケストレーションシステムのデファクトスタンダードとなっている．Kubernetes では，複数の Kubernetes Node の管理やコンテナのローリングアップデート，オートスケーリング，死活監視，ログ管理などサービスを実稼働環境で動かす上で必要不可欠となる機能を備えている．Docker 同様，デプロイするコンテナとその周辺のリソースは YAML 形式や JSON 形式で記述した宣言的なコードによって管理する．Infrastructure as Code に則っているため，実行環境に左右されず毎回常に同じコンテナが起動される．

GCP を筆頭にクラウド環境でもサポートされるようになり，現時点で AWS と Azure においても提供されている．そのため Kubernetes は徐々に注目を集めるようになり，今では多くの企業の実稼働環境で取り入れられている．

Kubernetes は，複数のサーバを束ねたクラスタの上で動作する．サーバの役割は二つに分かれており，システム全体を統合管理するサーバをマスターノード（コントロールプレーン），実際にコンテナを起動させるサーバをワーカーノードと呼ぶ．マスターノー



図 2.11: Kubernetes のロゴ

ドはシングルでも動作するが、基本的には冗長性や耐障害性を考慮して複数のマスターノードをクラスタリングすることが多い。クラウド環境を用いた場合、クリックひとつで Kubernetes クラスタを用意することができる。状況に応じてワーカーノードを追加・削除でき、自由にスケーリング出来る点も強みである。クラウドの種類によっては特定の条件に合わせて自動でノードのオートスケーリングを行うこともできるが、オンプレ環境では自前で実装する必要がある。

Kubernetes 自体は、多数のコンポーネントによって構成されるマイクロサービスアーキテクチャを採用している。すべてのコンポーネントが kube-apiserver と呼ばれる Kubernetes 内の API サーバを中心として動作しており、殆ど全ての処理は kube-apiserver を通して実行される。kube-apiserver はマスターノードに含まれる。他にもマスターノード内で動作するコンポーネントとしては、Kubernetes クラスタのすべての情報を保持する etcd、コンテナを起動させるノードをスケジューリングする kube-scheduler、ノード上で動作するコンテナを監視し必要に応じてコンテナを追加・削除するよう指示する kube-controller-manager などが挙げられる。対してワーカーノードで動作する主なコンポーネントには、kubelet などがある。kubelet を含め、本研究の実装で用いた kubeadm, kubectl に関しては以下で詳細に説明する。

Kubeadm

Kubeadm [6] は、Kubernetes クラスタを構築するためのベストプラクティスを提供するツールである。Kubeadm が提供するコマンドをいくつか以下に示す。

kubeadm init

クラスタの最初のコントロールプレーンとなるノードを起動する。

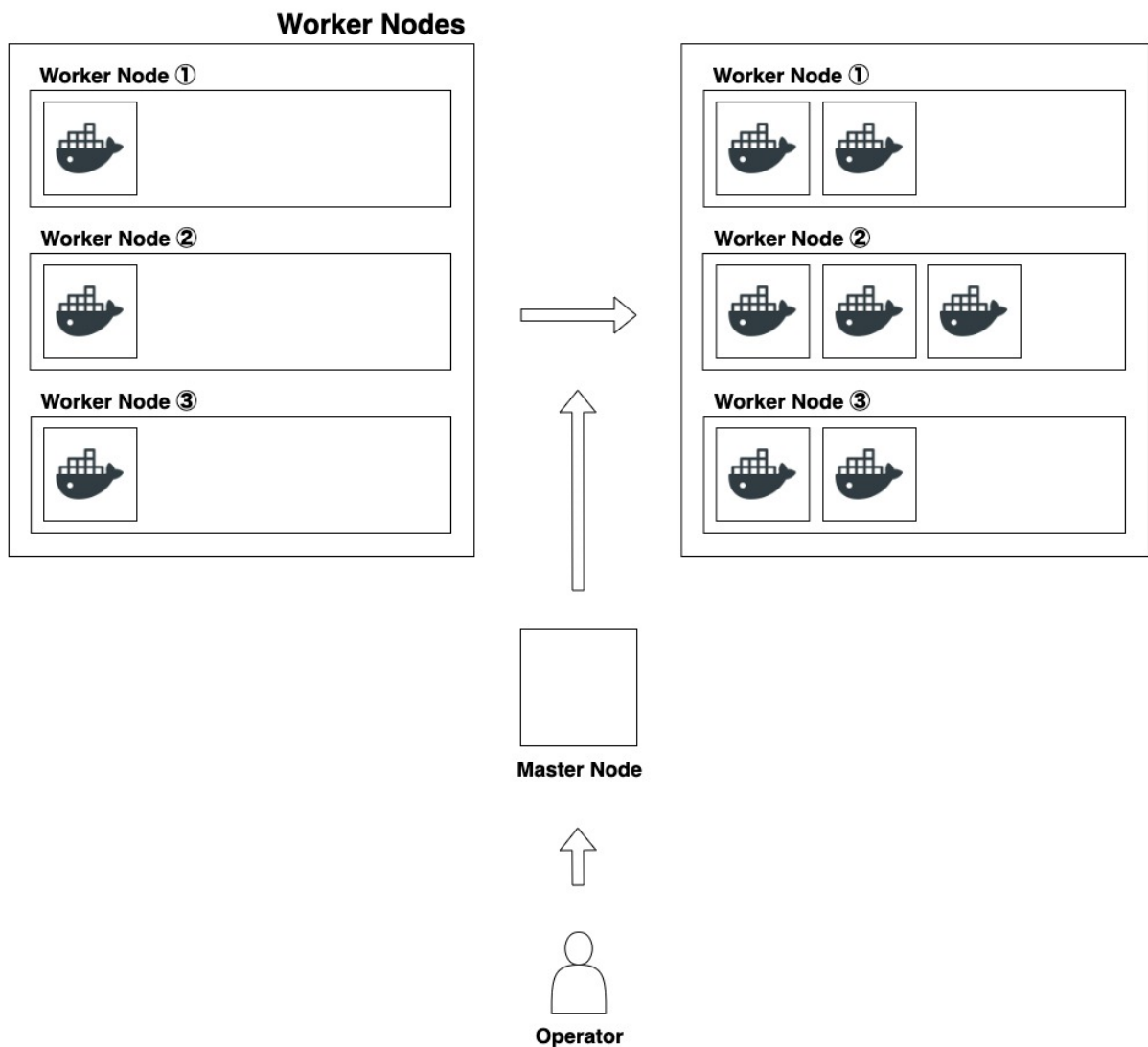


図 2.12: Kubernetes でのコンテナデプロイ

kubeadm join

クラスタに追加のコントロールプレーンまたはワーカーノードを参加させる。

kubeadm upgrade

クラスタのバージョンを最新へアップグレードする。

kubeadm reset

kubeadm init や kubeadm join によって生じた変更を取り消す。

kubelet

kubelet [7] は, Kubernetes クラスタ内の各ワーカーノードで動作するコンポーネントである。kubectl は, Docker などの CRI と連携して実際にコンテナを起動・停止する役割をもつ。具体的には etcd の情報を監視して, 自身のノードに割り当てられてまだ起動していないコンテナがあれば起動する。etcd に格納された情報は, kube-apiserver や kube-controller-manager によって kube-apiserver を通して書き換えられ, 実際のコンテナの操作に関しては kubelet が担うといった役割分担がされている。2.4.2 章の kubeadm, ならびに 2.4.2 章の kubectl は, Kubernetes クラスタ構築時や操作時に用いるコマンドツールであるのに対して, kubelet はコンテナの管理を行うデーモンとして動作する。

kubectl

kubectl [8] は, Kubernetes クラスタをコントロールするためのツールである。新規コンテナのデプロイや削除, アップデートから, 動作中のコンテナやクラスタを構成するノードの情報の取得など, サービスの運用を支援する API が提供されている。kubectl が提供するコマンドをいくつか以下に示す。

kubectl get nodes

クラスタに参加するノードのステータスやロール (役割), IP アドレス等を取得する。

kubectl get pods

ポッドの名前やステータス, 再起動の回数等を取得する。

kubectl apply

ポッドに新たな設定を反映させる。

2.5 OpenVPN

本節では, 本研究で使用した VPN 技術ならびにソフトウェア VPN である OpenVPN について概説する。

VPN とは, “Virtual Private Network” の略で, 日本語では “仮想専用線” と呼ばれる。VPN は, インターネット上に擬似的なプライベートネットワークを実現する技術, またはそのネットワーク自体を指す。VPN を使用することで, インターネット上の異なるセグメント同士であっても, あたかも専用線で接続されているかのように通信することが可能である。VPN には L2VPN と L3VPN があり, L2VPN は異なるセグメント同士を接続しひとつの擬似的な LAN を構築するものであり, L3VPN では IP プロトコルでの通信が可能となる。セキュリティ面においては, 通信内容をカプセル化することでパケットの中身の覗き見や改竄のリスクを低減することができる。

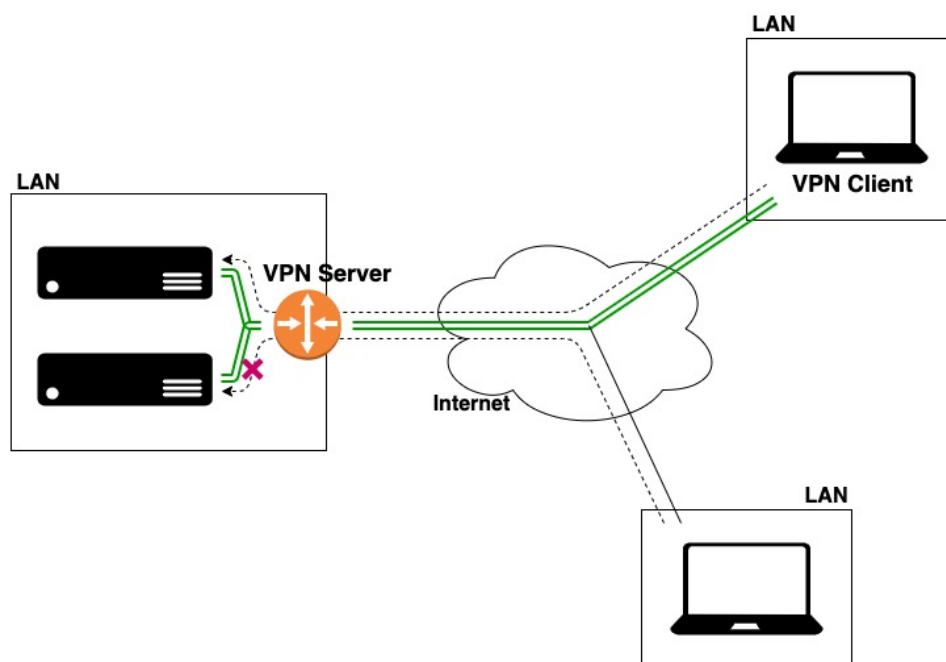


図 2.13: VPN

2.5.1 OpenVPN

OpenVPN [9] は, OpenVPN Technologies, inc. が中心になって開発しているオープンソースの VPN ソフトウェアである.



図 2.14: OpenVPN のロゴ

OpenVPN は幅広い OS でサポートされており, Window, Linux, Mac OS, iOS, Android で利用可能である. 異なる OS 間でも利用可能であるため, ひとつの VPN ネットワーク内に異なる OS が混在していても正常に動作する. OpenVPN では, 対応した OS のサーバがひとつでもあれば簡単に VPN サーバを構築することが可能である. VPN ネット

ワークに参加するためには認証が必要であり，OpenVPN では静的鍵による認証や証明書認証，ID/パスワード認証，二要素認証といった複数の認証方法から任意のものを選択できる．接続方法としては，ルーティングとブリッジが提供され，ルーティングは L3VPN，ブリッジは L2VPN に対応する．VPN ネットワーク内でブロードキャストを行いたい場合など，擬似的な LAN を構築したい場合以外は基本的に L3VPN を用いる．L3VPN にあたるルーティングでは，クライアントサーバ接続とサイト間接続が提供されている．クライアントサーバ接続では，各クライアントに認証設定が必要となり，接続の準備としてサーバ側での認証情報の生成が必要となる．認証情報を共有後，クライアント側では接続のために設定ファイルを用意し，コマンドやアプリケーションを用いて接続を行う．対するサイト間接続では，VPN の設定は各セグメントの VPN サーバで完結する．VPN サーバ間での接続が確立できれば，各拠点に配置されたサーバはお互いに疎通が可能となる．本研究の提案手法においては，ルーティング形式のサイト間接続を活用した．

第3章 本研究における課題定義と仮説

本章では，2章で述べた背景より，本研究における課題とその要件について議論し，先行研究および提案システムを概説することで本研究で用いるアプローチについて述べる．

3.1 課題定義

本研究では，惑星規模の分散システムのための試験環境の構築手法が整備されていないことを課題とする．

2.3.1章で述べたように，惑星規模の分散システムは地理的に分散したコンピュータで構成されるため，ネットワーク上の通信の遅延がシステム全体に影響を及ぼす．よって試験環境では，通信の遅延を考慮した上で地理的に分散したコンピュータの協調動作を試験しなければならない．試験環境としては PlanetLab やパブリッククラウドサービスの活用，BSafe.network が既に知られているが，これらにはそれぞれ課題があると考えられる．PlanetLab では，世界中の 717 地域 1353 のサーバという地理的に幅広い選択肢が用意されているが，サーバの OS や CPU，メモリ等を柔軟に変更できない．パブリッククラウドサービスでは，リージョンを活用することで地理的に離れた多数のデータセンターを利用することができるが，リージョンが限定的である．BSafe.network では，世界中の 32 の大学が保有するサーバを用いてブロックチェーン技術に関する実験を行えるが，各サーバの管理権限が各大学のオペレータに委ねられていることで，試験を行う際に各オペレータの手作業が必要となる．

3.2 課題解決における要件

本節では，課題解決における要件を定義する．

3.1章で述べた課題を踏まえて，惑星規模の分散システムの試験環境の構築における必要要件は，

- OS や CPU，Memory といったサーバ環境を柔軟に変更可能であること
- 公の実稼働環境を想定したネットワークでの通信の遅延を考慮できること
- 異なる管理権限下にある各サーバに対し統合的管理が可能であり，各オペレータの手作業を軽減できること

であると考えた．

3.2.1 サーバ環境の柔軟性

惑星規模の分散システムの試験では、OS や CPU、メモリ等のサーバの環境を柔軟に変更可能である必要がある。何故なら惑星規模の分散システムにおいて、システムを構成するコンピュータの環境はユーザ依存であり開発者が一意に決めることはできないからである。よって、異なる OS を搭載したコンピュータによる協調動作や、CPU やメモリといった資源を制限した場合の試験を行う必要がある。

3.2.2 ネットワークでの通信の遅延の考慮

惑星規模の分散システムは、地理的に分散したコンピュータによって構成される分散システムである。そのため、各コンピュータが地理的に分散することによるネットワーク上の通信の遅延がシステムの動作に影響を与える可能性がある。試験環境は、これらを考慮した上で試験が行える環境でなければならない。

3.2.3 異なる組織間での試験における手作業の軽減

BSafe.network のように複数の組織間で惑星規模の分散システムの試験を行う場合、各サーバの管理権限が異なることにより、試験の遂行において各オペレータの手作業が必要となる。各オペレータの手作業が要求されることによって、作業中の人為的ミスが発生する可能性が高まるとともに、手戻りによる作業の遅れが発生しやすくなる。作業が各オペレータ依存になるため、サーバ構築における冪等性が担保されない問題もある。試験環境では、属人的な手作業を軽減し試験が素早く正確に行わなければならない。

3.2.4 地理的に分散したサーバの統合的な管理

惑星規模の分散システムの試験環境では、すべてのサーバを統合管理する必要がある。試験環境内のすべてのサーバに対し一斉に指示できるようにすることで、試験を迅速に進めることができるからである。なお惑星規模の分散システムでは各サーバが地理的に分散しており、ネットワーク上で別のセグメントに配置されている可能性がある。よって、試験環境では異なるセグメントに配置されたサーバに対して統合的な管理を行う必要がある。

3.3 先行研究

本節では、惑星規模の分散システムの試験環境の構築手法として提案された先行研究について述べる。先行研究では、各コンピュータを操作するためのデバッグエージェントを独自実装したものや、ネットワークエミュレータを用いた仮想ネットワークによる手法が提案されている。本研究ならびに PlanetLab、パブリッククラウドサービス、BSafe.network

では実際に地理的に分散したコンピュータを用いているのに対し、Emulab では仮想的にネットワークを再現することで惑星規模の分散システムの試験を可能としている。

3.3.1 P2P アプリケーションの開発と性能評価のための統合開発環境の提案

既存の提案手法として、地理的に分散したコンピュータを統合管理するためにデバッグエージェントと呼ばれる遠隔操作のアプリケーションを独自で実装したものがある。デバッグエージェントは、試験対象のアプリケーションに対して命令を送信したり通信内容をログとして抽出する。デバッグエージェントを通して、P2P アプリケーションの開発、性能評価を支援するための統合開発環境の構築に成功している。しかし、デバッグエージェントは予め各サーバで動作させる必要がある。また、デバッグエージェントは試験対象のアプリケーションに依存するため、アプリケーションの変更をする場合はデバッグエージェントにも修正が必要になる。

3.3.2 プロセスレベルの仮想化を用いた大規模分散システムテストベッド

分散システムや分散ネットワークを、ネットワークエミュレータによって構築された仮想的なネットワーク上で開発する取り組みである。仮想ネットワークの構築には Emulab [10] を活用している。Emulab は大規模なソフトウェアシステムであり、仮想ネットワーク内に点在するマシン同士の接続環境を自由に設定することが可能である。数台のコンピュータ上に数千台の仮想環境をプロセスレベルで構築し、それらをネットワークシミュレータにより相互接続することによって、擬似的なネットワーク環境における試験を可能にするものである。

3.4 仮説

本研究では 3.2 章で述べた惑星規模の分散システムの試験環境の構築における必要要件を満たす手法を提案したい。そこで OpenVPN と Kubernetes を組み合わせることで、本研究の課題解決のための必要要件を満たした試験環境を構築できるのではないかと考えた。各要件に対して、本研究で提案するシステムによる実現が可能であると考えられる点を本節では述べる。

3.4.1 サーバ環境の柔軟性

本研究のシステムでは、コンテナ型仮想化技術である Docker を活用する。Docker において、Docker ファイルから生成されるコンテナイメージは一度ビルドされれば他の OS でも動作し冪等性が担保される。コンテナは CPU やメモリといった資源も仮想的に制限す

ることができるため、Docker を活用することにより様々なサーバ環境下での試験を行うことが可能となる。

3.4.2 ネットワークでの通信の遅延の考慮

本研究で提案するシステムは、公のネットワーク上で構築することを前提としている。OpenVPN を活用することによる VPN オーバーレイネットワークを構築することにより、ネットワークにおいて別々のセグメントに配置されたサーバ間の疎通性を確保する。公のネットワーク上で地理的分散システムの試験環境を構築することで、通信の遅延を考慮した試験を行うことが可能である。

3.4.3 異なる組織間での試験における手作業の軽減

本研究の提案手法で活用する Kubernetes では、複数のサーバをクラスタ化することにより、オペレータがマスターノードへ指示を送ることで、すべてのノードの操作が可能である。Kubernetes クラスタの構築において、サーバの管理権限が要求されるのは最初のクラスタ構築時のみである。クラスタリングの完了後は、各サーバの管理権限を必要とせずサーバの統合管理が可能のため、異なる組織間での試験においてオペレータの手作業を軽減することが可能である。

3.4.4 地理的に分散したサーバの統合的な管理

本研究では、OpenVPN と Kubernetes を組み合わせた惑星規模の分散システムのための試験環境を提案する。地理的に分散することによりネットワーク上で互いに通信のできないサーバ間で疎通性を確保した上で、Kubernetes クラスタを構築する。よって、惑星規模の分散システムの試験環境における統合管理が可能だと考える。

3.5 提案システム概要

本節では、提案システムの概要を述べる。惑星規模の分散システムを構成するサーバ間を、OpenVPN オーバーレイネットワークを構築することにより疎通可能な状態にする。OpenVPN オーバーレイネットワーク上で地理的に分散したサーバをクラスタリングすることで、Kubernetes クラスタを構築する。OpenVPN と Kubernetes を組み合わせた本研究の提案システムは、オペレータが Kubernetes クラスタ内のマスターノードを介してすべてのサーバに対し操作を可能とする。



図 3.1: システム概要図

第4章 実装

本章では提案手法の実装について述べる。

4.1 実装環境

本節では，本研究で構築した実装環境について概説する．環境構築は実験用の LAN で行い、LAN をパブリックなインターネットと仮定する。

4.1.1 ハードウェアおよびソフトウェア

本研究で使用したハードウェアおよびソフトウェアとそのバージョンを以下に示す。

表 4.1: 使用したハードウェアおよびソフトウェア

ハードウェア/ソフトウェア	機種/バージョン
Server	FUJITSU PRIMERGY S6 (12 CPUs, Memory 48GB)
VMWare ESXi	6.5
VyOS	1.2.1
OpenVPN	2.3.4
Ubuntu	18.04
kubeadm	1.16.3
kubelet	1.16.3
kubectrl	1.16.3
HA-Proxy	1.8.8

4.1.2 物理サーバの準備

本研究では，実装において複数のセグメントおよび Kubernetes クラスタの構築に複数のサーバが必要であったため，それらを仮想的に作成できる VMWare ESXi（以下，ESXi）を導入した．使用したのは，ESXi6.5である．ESXiはホスト OS を必要とせず，直接ハードウェアにインストールさせて動作させるハイパーバイザー型であるため，まず初めに

ESXi インストーラのブータブルイメージを USB メモリに書き込み，FUJITSU サーバにインストールした．計二台の FUJITSU サーバに ESXi をインストールし，それぞれ以下の IP アドレスを設定した．

表 4.2: ESXi の IP アドレス

名前	IP アドレス
1 台目	10.4.0.13
2 台目	10.4.0.14

4.1.3 ネットワーク構成

本研究で構築したネットワーク構成について説明する．

まず初めに，ESXi の仮想スイッチと VLAN を用いて二つの ESXi サーバ上に新たに計三つの論理セグメントを構築した．Vlan によって論理的にセグメントを分割することで，お互いに通信不可能な環境とした．以下に，Vlan ID と対応するアドレスプレフィックスを示す．なお，10.4.0.0/16 のアドレスプレフィックスは Vlan ID 0 に対応している．

表 4.3: Vlan ID と対応するアドレスプレフィックス

Vlan ID	アドレスプレフィックス
0	10.4.0.0/16
10	192.168.10.0/24
20	192.168.20.0/24
30	192.168.30.0/24

4.1.4 VM の配置

ネットワーク構築後，Kubernetes クラスタの構築に必要なサーバを VM として立ち上げた．それぞれの VM の OS には Ubuntu18.04 を採用した．以下に構築したサーバの詳細を示す．

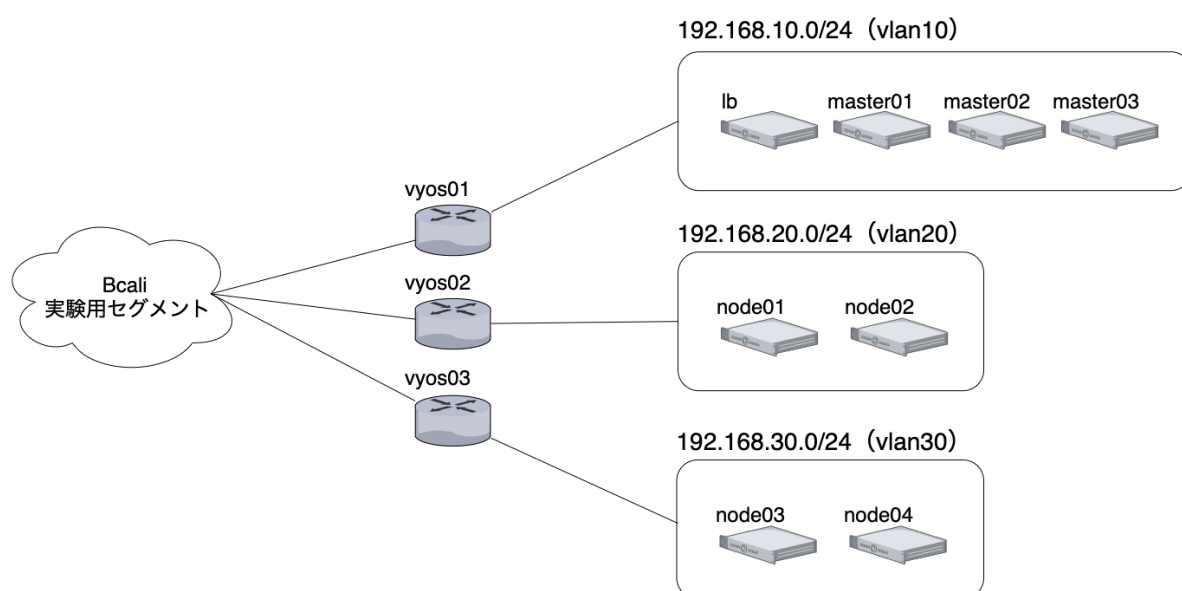


図 4.1: ネットワーク構成

表 4.4: 設置した VM の詳細

名前	ネットワークインターフェース名	Vlan ID	IP アドレス	役割
lb	ens160	10	192.168.10.253	マスターノードのロードバランサー
master01	ens160	10	192.168.10.101	マスターノード
master02	ens160	10	192.168.10.102	マスターノード
master03	ens160	10	192.168.10.103	マスターノード
node01	ens160	20	192.168.20.101	ワーカーノード
node02	ens160	20	192.168.20.102	ワーカーノード
node03	ens160	30	192.168.30.101	ワーカーノード
node04	ens160	30	192.168.30.102	ワーカーノード

4.1.5 ルーターの配置

次に各拠点に OpenVPN の設定をするルーターを設置した．ルーターの OS には VyOS 1.2.1, OpenVPN はバージョン 2.3.4 を採用した．以下にルーターのネットワーク情報を示す．

表 4.5: 設置したルーターの詳細

名前	ネットワークインターフェース名	Vlan ID	IP アドレス
vyos01	eth0	0	10.4.0.90
vyos01	eth1	10	192.168.10.1
vyos02	eth0	0	10.4.0.91
vyos02	eth1	20	192.168.20.1
vyos03	eth0	0	10.4.0.92
vyos03	eth1	30	192.168.30.1

全てのルーターはお互いに疎通可能である．さらに，各拠点に設置されたサーバと疎通できるよう eth1 のネットワークインターフェースには別の IP アドレスを設定した．この時点での各サーバの疎通性は以下の通りである．

表 4.6: OpenVPN 設定前の各サーバの疎通性

	lb	master01	master02	master03	node01	node02	node03	node04
lb		○	○	○	×	×	×	×
master01	○		○	○	×	×	×	×
master02	○	○		○	×	×	×	×
master03	○	○	○		×	×	×	×
node01	×	×	×	×		○	×	×
node02	×	×	×	×	○		×	×
node03	×	×	×	×	×	×		○
node04	×	×	×	×	×	×	○	

4.1.6 OpenVPN の設定

4.1.5 で示したように，OpenVPN の設定をする前ではすべてのサーバはお互いに疎通可能な状態にはない．Kubernetes は，クラスタに参加するサーバのすべてが疎通可能，厳密には IP reachable な環境下にある必要がある．そこで OpenVPN を用いて，複数の分離した LAN を仮想的に接続し Kubernetes の要件を満たそうと試みた．本実装では，OpenVPN の site-to-site モードを採用した．client-server モードを採用しなかった理由としては以下の二点が挙げられる．

1. Kubernetes は通信時にデフォルトゲートウェイに設定したネットワークインターフェースを使用するため，サーバ毎に OpenVPN を設定する client-server モードではトンネルインターフェースを通して通信ができない．
2. サーバ毎に証明書と鍵の管理が必要なため扱いづらい．

対して，site-to-site モードでは以下の利点が挙げられる．

1. ルーティングはルーターに任せられるため，サーバは通信時にデフォルトゲートウェイに設定されたネットワークインターフェースを使用できる．
2. OpenVPN の設定は LAN 内のルーターのみ．

以下に，OpenVPN 設定後の各サーバの疎通性を示す．

表 4.7: OpenVPN 設定前の各サーバの疎通性

	lb	master01	master02	master03	node01	node02	node03	node04
lb		○	○	○	○	○	○	○
master01	○		○	○	○	○	○	○
master02	○	○		○	○	○	○	○
master03	○	○	○		○	○	○	○
node01	○	○	○	○		○	○	○
node02	○	○	○	○	○		○	○
node03	○	○	○	○	○	○		○
node04	○	○	○	○	○	○	○	

4.1.7 Kubernetes クラスタの構築

OpenVPN による拠点間の接続を行った後，Kubernetes クラスタを構築した．本研究の実装では，Kubeadm を使用した高可用性 Kubernetes クラスタを構築するため，まず初めに複数マスターへのリクエストを振り分けるロードバランサーを設置した．ロードバランサーの構築には，HA-Proxy 1.8.8 を採用した．

```

1  frontend kubernetes
2      bind *:6443
3      option tcplog
4      mode tcp
5      default_backend kubernetes-backend
6
7  frontend etcd
8      bind *:2379

```

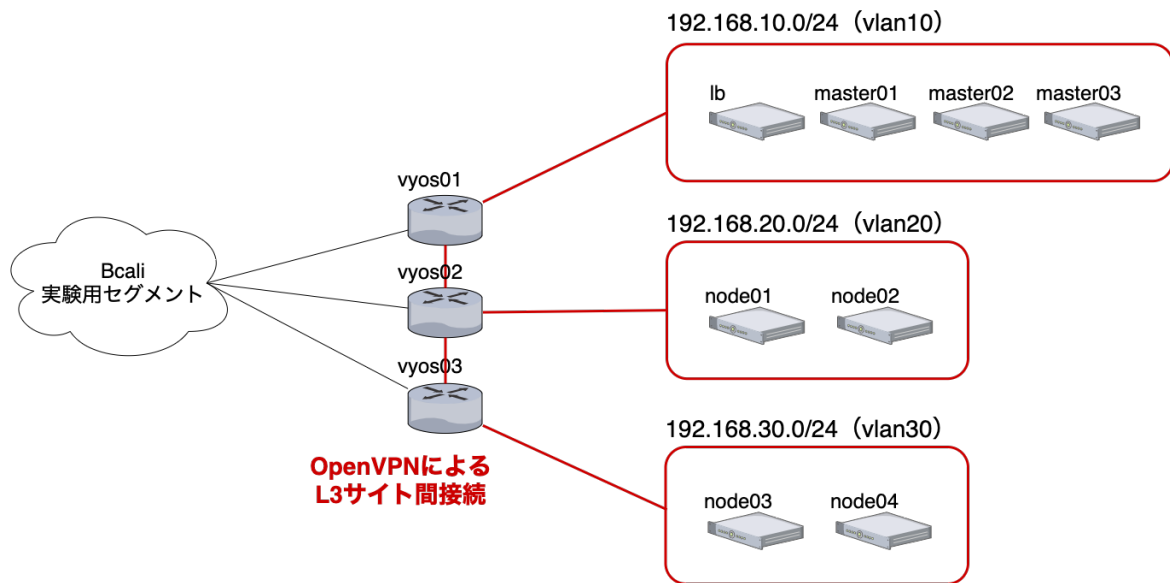


図 4.2: OpenVPN によるサイト間接続

```

9      option tcplog
10     mode tcp
11     default_backend etcd-backend
12
13     backend kubernetes-backend
14         mode tcp
15         balance roundrobin
16         option tcp-check
17         server master01 192.168.10.101:6443 check
18         server master02 192.168.10.102:6443 check
19         server master02 192.168.10.103:6443 check
20
21     backend etcd-backend
22         mode tcp
23         balance roundrobin
24         server master01 192.168.10.101:2379 check
25         server master02 192.168.10.102:2379 check
26         server master03 192.168.10.103:2379 check

```

上記の設定で、ロードバランサーのポート 6443 番とポート 2379 番へのリクエスを三台のマスターノードへと振り分けている。

次に、マスターノードとワーカーノードを立ち上げるにあたり必要なパッケージをインストールした。Kubernetes のランタイムとして使用する Docker に加え、クラスタ構築時に用いる kubeadm と kubelet、クラスタ操作時に必要な kubectl を apt によって取得した。パッケージの用意が完了したのち、マスターノードからクラスタ構築作業を行った。kubeadm ではクラスタの初期化用に init コマンドが用意されており、初めのマスターノードにて実行することでクラスタの基盤を作成可能である。初期化に成功した場合、以下のようなテキストが出力される。

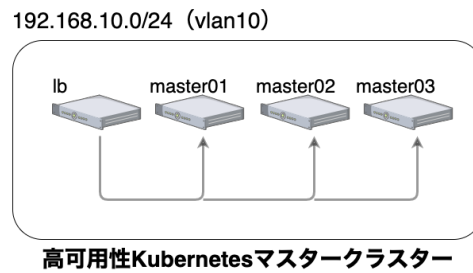


図 4.3: Kubernetes マスタークラスターの構築

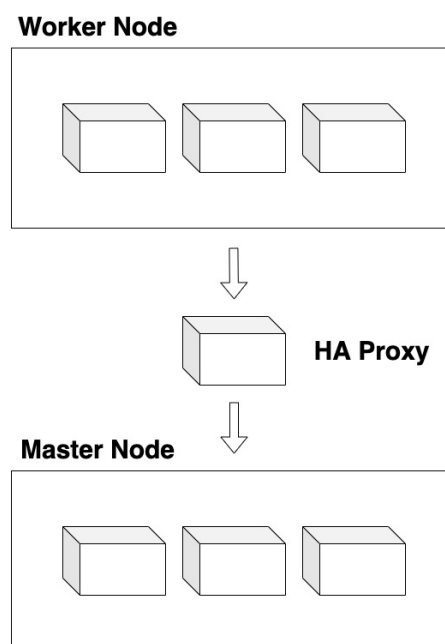


図 4.4: マスターノードとワーカーノードの関係性

```

1  You can now join any number of control-plane nodes by copying
    certificate authorities
2  and service account keys on each node and then running the
    following as root:
3
4  kubeadm join 192.168.10.253:6443 --token { token } \
5      --discovery-token-ca-cert-hash sha256:{ hash } \
6      --control-plane
7
8  Then you can join any number of worker nodes by running the
    following on each as root:
9
10 kubeadm join 192.168.10.253:6443 --token { token } \
11     --discovery-token-ca-cert-hash sha256:{ hash }

```

出力にある通り，与えられたコマンドを他のマスターノードとワーカーノードから実行することでクラスタへの参加が行える．各サーバにて上記のコマンドを実行した結果，マスターノードからクラスタが構築できていることを確認できた．

```
1  $ kubectl get nodes
2  NAME                STATUS    ROLES    AGE      VERSION
3  master01             Ready     master   58d      v1.16.3
4  master02             Ready     master   58d      v1.16.3
5  master03             Ready     master   58d      v1.16.3
6  node01               Ready     <none>    8d       v1.16.3
7  node02               Ready     <none>    4d22h    v1.16.3
8  node03               Ready     <none>    6d16h    v1.16.3
9  node04               Ready     <none>    8d       v1.16.3
```

4.2 システム全体

本研究で構築した実装環境の図を以下に示す．

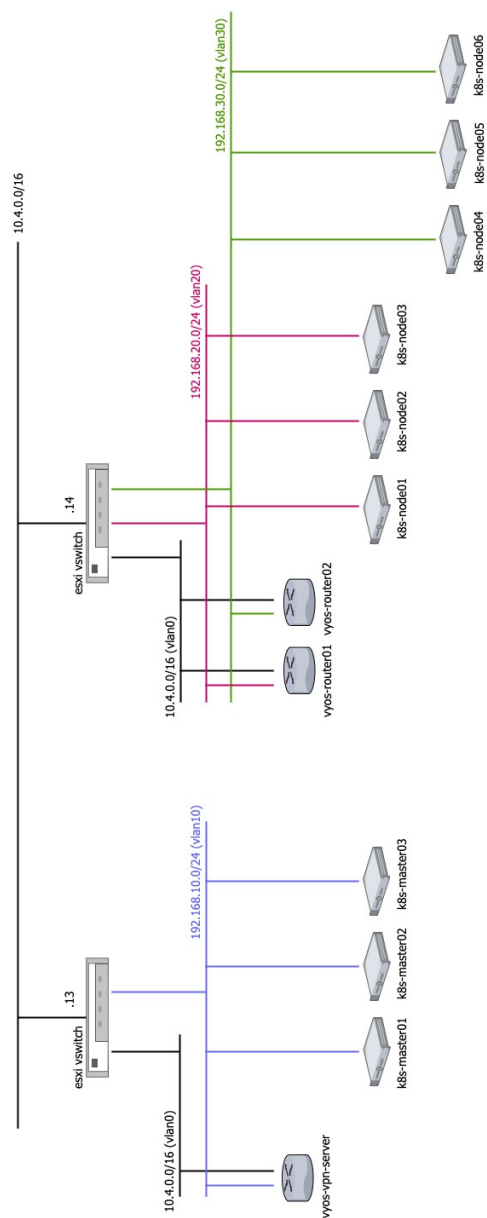


図 4.5: ネットワーク構成図

第5章 評価

本章では、本システムが 3.2 章で述べた課題解決における必要要件を満たしているか確認することで本研究の提案の評価を行う。必要要件は以下の通りである。

- OS や CPU, Memory といったサーバ環境を柔軟に変更可能であること
- 公の実稼働環境を想定したネットワークでの通信の遅延を考慮できること
- 異なる管理権限下にある各サーバに対し統合的管理が可能であり、各オペレータの手作業を軽減できること

評価では、本システムが地理的に分散したサーバに対して適用可能か、ならびに本システムを利用することによる手作業の削減の二点における定性評価、加えてネットワーク上の遅延による本システムへの影響を定量的に評価した。

5.1 地理的に分散したサーバに対する統合管理

本システムでは、互いに疎通の取れない別セグメントを OpenVPN オーバーレイネットワークで繋ぎ、その上で Kubernetes クラスタを構築した。VPN によって異なるセグメントへの疎通性が確保されていることは、ping コマンドを用いて確認した。以下は、Kubernetes クラスタ内のマスターノード (192.168.10.101) から別セグメントに配置されたワーカーノードに対して ping コマンドを実行した際の結果である。

```
1  $ ping 192.168.20.101 -c5
2  PING 192.168.20.101 (192.168.20.101) 56(84) bytes of data.
3  64 bytes from 192.168.20.101: icmp_seq=1 ttl=62 time=1.41 ms
4  64 bytes from 192.168.20.101: icmp_seq=2 ttl=62 time=1.02 ms
5  64 bytes from 192.168.20.101: icmp_seq=3 ttl=62 time=1.28 ms
6  64 bytes from 192.168.20.101: icmp_seq=4 ttl=62 time=1.06 ms
7  64 bytes from 192.168.20.101: icmp_seq=5 ttl=62 time=1.43 ms
8
9  --- 192.168.20.101 ping statistics ---
10  5 packets transmitted, 5 received, 0% packet loss, time 4005
    ms
11  rtt min/avg/max/mdev = 1.028/1.243/1.432/0.172 ms
```

パケットロスはなく、通信が行えていることが確認できた。

次に、kubectl コマンドによって異なるセグメントに位置するサーバをクラスタリングできていることを確認した。

```

1 $ kubectl get nodes -owide
2 NAME STATUS ROLES VERSION INTERNAL-IP OS
  -IMAGE KERNEL-VERSION CONTAINER-RUNTIME
3 master01 Ready master v1.16.3 192.168.10.101
  Ubuntu 18.04.3 LTS 4.15.0-70-generic docker://18.9.7
4 master02 Ready master v1.16.3 192.168.10.102
  Ubuntu 18.04.3 LTS 4.15.0-70-generic docker://18.9.7
5 master03 Ready master v1.16.3 192.168.10.103
  Ubuntu 18.04.3 LTS 4.15.0-70-generic docker://18.9.7
6 node01 Ready <none> v1.16.3 192.168.20.101
  Ubuntu 18.04.3 LTS 4.15.0-74-generic docker://18.9.7
7 node02 Ready <none> v1.16.3 192.168.20.102
  Ubuntu 18.04.3 LTS 4.15.0-74-generic docker://18.9.7
8 node03 Ready <none> v1.16.3 192.168.30.101
  Ubuntu 18.04.3 LTS 4.15.0-74-generic docker://18.9.7
9 node04 Ready <none> v1.16.3 192.168.30.102
  Ubuntu 18.04.3 LTS 4.15.0-74-generic docker://18.9.7

```

上記の結果から、異なる IP アドレスをもつサーバに対しクラスタリングを行えていることが確認できた。

最後に、Kubernetes クラスタ上にアプリケーションのデプロイが可能であることを確認した。Nginx を Kubernetes のワークロードである Deployment としてクラスタ上にデプロイした結果が以下である。

```

1 $ kubectl create deployment --image nginx hello-world
2 $ kubectl get pods -owide
3 NAME AGE IP NODE READY STATUS RESTARTS
  GATES NOMINATED NODE READINESS
4 hello-world-c6c6778b4-5n74d 1/1 Running 0
  d22h 10.44.0.1 node01 <none> <none> 4
5 hello-world-c6c6778b4-6mrj4 1/1 Running 0
  d22h 10.42.0.1 node03 <none> <none> 4
6 hello-world-c6c6778b4-fmnxt 1/1 Running 0
  d22h 10.47.0.1 node02 <none> <none> 4
7 hello-world-c6c6778b4-r8b5w 1/1 Running 0
  d22h 10.44.0.2 node04 <none> <none> 4

```

node01 から node04 のすべてのワーカーノードに対してコンテナのデプロイが完了したことを確認できた。これにより、異なるセグメントに位置するノードに対して統合管理が可能であることを確認した。

5.2 オペレータの手作業の軽減

本節では、本システムが異なる組織間での惑星規模の分散システムの試験においてオペレータの手作業を軽減可能であるか評価する。複数の大学によって構成される研究ネットワークで、惑星規模の分散システムを試験する場合を考える。試験に必要な作業を手作業で行った場合、以下の作業が必要であると考え。

1. 各大学のリソース（OS, CPU, Memory）の共有
2. 各大学における作業内容の確定・共有
3. 各大学のサーバ管理者とのスケジューリングの調整
4. 各大学でのデプロイ作業
5. 各大学から作業完了の連絡が来るのを待機
6. 全大学での作業完了の共有
7. 試験環境の利用開始

上記の作業は、試験環境に対する変更が必要になる度に行わなければならない。手作業に加え、オペレータ間のコミュニケーションが密に求められる。そのため、各オペレータが手作業で試験の準備をしなければいけないことによる時間の消費、属人性によるミスが発生する可能性がある。対して本研究の提案システムを用いた場合、以下の作業のみで試験の準備を行える。

1. Kubernetes の定義ファイルの作成
2. kubectl コマンドによるコンテナの起動

システムの統合管理が可能であるため、各オペレータの手作業を軽減することができることに加え、属人性を排除し作業の冪等性を担保することができる。

5.3 通信の遅延による本システムへの影響

本節では、本システムが地理的に分散したコンピュータによって構成される場合を考慮し、通信の遅延による本システムへの影響をコンテナの起動時間を計測することによって定量評価した。実験環境では、tc コマンドを用いて擬似的に通信の遅延を発生させた。これによってマスターノードとワーカーノードの通信において 100ms, 300ms, 500ms の遅延を発生させ、異なるコンテナ起動数においてすべてのワーカーノードでコンテナが起動するまでの時間を計測し 5 回の平均値を算出した。計測には kubectl コマンドのコンテナ監視機能を用いており、単位はこれに依存し秒 (s) となっている。

結果は以下ようになった。

結果から、起動するコンテナ数が少ない場合は遅延によるコンテナ起動時間の変化は見られなかったが、コンテナ数が増えると遅延時間による差が生じた。Kuberntes では、コンテナの起動には三つのステップがある。コンテナをどのワーカーノードにて起動するかを決めるスケジューリングの段階、ワーカーノードで動作する kubelet が自身のノードで起動すべきコンテナがあることを把握する段階、最後に kubelet がコンテナを起動する段階である。初めのスケジューリングはマスターノード、最後のコンテナの起動はワー

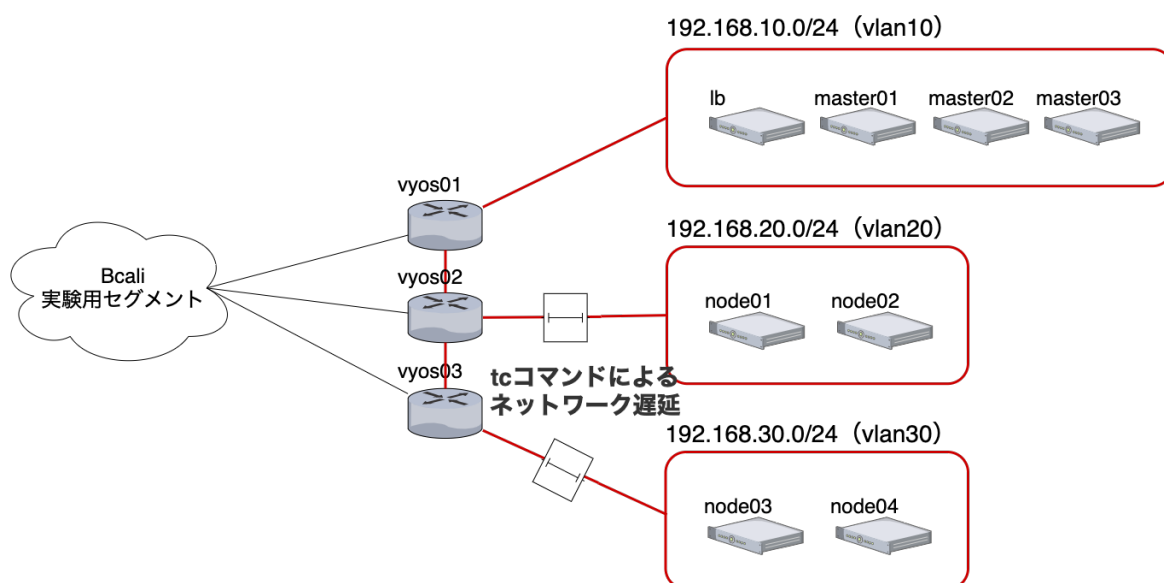


図 5.1: tc コマンドによるマスターノード・ワーカーノード間の擬似的な遅延の発生

コンテナ数	4	20	100
遅延時間			
100ms	24.9s	38.2s	125.4s
300ms	24.6s	37.4s	133s
500ms	25.5s	38s	141s

表 5.1: 通信の遅延によるコンテナ起動時間の変化

カーノードによって行われ、途中にワーカーノードがマスターノードの情報を参照する段階のみ遅延時間が影響する。今回の結果より、起動時間全体に対する遅延の影響は比較的少なく、大半の時間はスケジューリングとワーカーノードでのコンテナの起動のための時間であると考え。スケジューリングとコンテナの起動は、マスターノードとワーカーノードの性能に依存しており、各ノードの性能を向上させることでコンテナ起動時間は短縮可能である。よって、本システムが地理的に分散したコンピュータによって構成され、Kubernetes クラスタにおける通信の遅延が発生した場合においても動作可能であると考えた。

5.4 評価のまとめ

本章では、二つの定性評価とひとつの定量評価を行うことで本システムが課題解決のための必要要件を満たしているか確認した。まず初めに OpenVPN オーバーレイネットワーク上で構築した Kubernetes クラスタが正常に動作しているかを、kubectl によるクラスタの状態確認および各ワーカーノードにコンテナを起動することで確認した。すべてのノード

ドが正常に動作し、各ノードでコンテナを起動できることを確認できた。よって、コンテナ型仮想技術によりサーバ環境を柔軟に変更できるという要件定義を満たした。次に、本システムによる手作業の削減について定性評価を行なった。複数の大学で別々の管理権限下にあるサーバを用いて試験を行う際、試験の準備を手作業で行う場合に比べ、本システムを用いた場合ではオペレータの手作業とコミュニケーションを大きく削減することが可能であることを確認した。これにより、異なる管理権限下にあるサーバに対する統合管理が可能であり、各オペレータの手作業を軽減できるという要件定義を満たした。最後に、試験環境を構成するサーバが地理的に分散することによる通信の遅延が本システムに与える影響を定量評価した。結果、試験環境において起動するコンテナ数の増加に伴って通信の遅延がコンテナの起動時間に影響することを確認した。しかし、通信の遅延の影響は起動時間全体に対し小さく、起動時間は Kubernetes クラスタを構成するノードの性能に依存するため、各ノードの性能を上げることで地理的に分散した場合でも拡張性があると考えた。これにより、試験環境を地理的に分散したサーバによって構成し、公の実稼働環境を想定したネットワークでの通信の遅延を考慮するという必要要件を満たした。

第6章 結論

本章では、本研究のまとめと今後の課題を示す。

6.1 まとめ

本研究では、惑星規模の分散システムの試験を行うための試験環境の構築手法を提案した。惑星規模の分散システムは、地理的に分散したコンピュータによって構成されるため、試験においてはネットワーク上の通信の遅延を考慮する必要があることを2章で示した。3章では、既存の提案手法としてPlanetLab、クラウドサービスのリージョンの活用、BSafe.networkをあげたが、それぞれに欠点があり未だ惑星規模の分散システムの試験環境の構築手法に課題が残されていることを指摘した。加えて、惑星規模の分散システムの試験における以下四点の必要要件を明らかにした。

- OS や CPU, Memory といったサーバ環境を柔軟に変更可能であること
- 公の実稼働環境を想定したネットワークでの通信の遅延を考慮できること
- 異なる管理権限下にある各サーバに対し統合的管理が可能であり、各オペレータの手作業を軽減できること

上記の必要要件を満たすため、本研究ではOpenVPNとKubernetesを組み合わせた惑星規模の分散システムのための構築手法を提案した。試験環境を構成するサーバ間をOpenVPNオーバーレイネットワークによって繋げ、その上でKubernetesクラスタを構築することで、地理的に分散したサーバを統合的に管理することができるのではないかと考えた。4章では、提案手法の構築を行った。ESXiを利用して仮想的にネットワーク環境を構築し、疎通性のないセグメント間でKubernetesクラスタを立ち上げた。5章では、3章で明らかにした惑星規模の分散システムの試験環境における必要要件を本システムが満たしているか評価を行なった。ネットワーク上の別セグメントに位置するサーバによるKubernetesクラスタが、各サーバに対し統括的な指示ができることより、本システムが本研究の課題に対する必要要件を満たしていることを確認した。本節の冒頭でも述べたように、惑星規模の分散システムではネットワークでの通信の遅延がシステムに影響を及ぼす。そのため、通信の遅延を考慮した上で各コンピュータが協調動作できていることを試験しなければならず、惑星規模の分散システムのための試験環境を構築することは容易ではない。本研究は、地理的に分散したコンピュータによって構成される分散システムの試験環境を構築手法を提案するものであり、システムの堅牢性の向上に繋がったと考える。

6.2 課題と展望

本節では，本研究の課題と展望について述べる．本研究の実装は，ESXi によって構築した仮想ネットワーク上で異なるセグメントを繋ぐ OpenVPN オーバーレイネットワークを実装し，さらにその上で Kubernetes クラスタの構築を行なった．実装は LAN 内で行なったものであり，実際に地理的に分散したコンピュータを用いて実装が行えなかった点は課題として残された．実用に向けた次の段階としては，本研究で提案したシステムを公のインターネット上で実装し再評価する必要があると考える．Bitcoin の基盤技術であるブロックチェーンが登場したことによって，惑星規模の分散システムには今後より注目が集まると考える．今後，惑星規模の分散システムの堅牢性をより一層向上させる必要があり，そのためには試験環境の提案手法の確立が大きな課題である．

第7章 謝辞

本論文の執筆にあたり，常に優しく，最後まで見捨てずにご指導してくださった慶應義塾大学政策・メディア研究科特任准教授鈴木茂哉博士，同大学政策・メディア研究科博士課程阿部涼介氏に感謝致します．お忙しいにも関わらず，毎週のようにミーティングを設けてくださったこと，研究について一から教えてくださったこと，行き詰まっている際に親身に相談に乗ってくださったことには本当に感謝しております．

参考文献

- [1] 金子 勇. Winny の技術, 2005.
- [2] Gtk-gnutella. <http://gtk-gnutella.sourceforge.net/en/?page=news>.
- [3] Satoshi Nakamoto. Bitcoin: A peer-to-peer electronic cash system. <http://www.cryptovest.co.uk/resources/Bitcoin%20paper%20original.pdf>, 2008.
- [4] Docker. <https://www.docker.com/>.
- [5] Kubernetes. <https://kubernetes.io/ja/>.
- [6] Kubeadm. <https://github.com/kubernetes/kubeadm>.
- [7] kubelet. <https://github.com/kubernetes/kubelet>.
- [8] kubectrl. <https://github.com/kubernetes/kubectrl>.
- [9] Openvpn. <https://openvpn.net/>.
- [10] Emulab. <https://www.emulab.net/portal/frontpage.php>.