

グリッドにおける計算資源と光パスネットワーク資源の コアロケーション実験

竹房 あつ子^{†1} 林 通 秋^{†2} 長 津 尚 英^{†3}
中 田 秀 基^{†1} 工 藤 知 宏^{†1} 宮 本 崇 弘^{†2}
大 谷 朋 広^{†2} 田 中 英 明^{†2} 鮫 島 康 則^{†3}
今 宿 互^{†3} 神 野 正 彦^{†3} 滝 川 好 比 郎^{†3}
岡 本 修 一^{†4} 田 中 良 夫^{†1} 関 口 智 嗣^{†1}

産業技術総合研究所, KDDI 研究所, 日本電信電話は共同で, ネットワーク資源の事前予約を行うための標準ウェブサービスインタフェースを定める G-lambda プロジェクトを推進している. このインタフェースの有効性を確認するために, 計算資源と光パスネットワーク資源を事前予約によりコアロケーションする実験を, 情報通信研究機構の協力を得て行った. 本稿では, 前提とするコアロケーションシステムモデル, G-lambda プロジェクトにおけるグリッドでのネットワークサービスインタフェースの規定, 産総研が開発した計算・ネットワーク資源を同時予約するスケジューラの概要, KDDI 研が開発した GMPLS で制御される光パスネットワークの提供機構の概要と, これらにより実現されたコアロケーションシステムプロトタイプの iGrid2005 での実証実験について報告する.

An Experiment of Co-allocating Computing and Lambda Path Resources on the Grid

ATSUKO TAKEFUSA,^{†1} MICHIAKI HAYASHI,^{†2} NAOHIDE NAGATSU,^{†3}
HIDEMOTO NAKADA,^{†1} TOMOHIRO KUDOH,^{†1}
TAKAHIRO MIYAMOTO,^{†2} TOMOHIRO OTANI,^{†2} HIDEAKI TANAKA,^{†2}
YASUNORI SAMESHIMA,^{†3} WATARU IMAJUKU,^{†3} MASAHIKO JINNO,^{†3}
YOSHIHIRO TAKIGAWA,^{†3} SHUICHI OKAMOTO,^{†4} YOSHIO TANAKA^{†1}
and SATOSHI SEKIGUCHI^{†1}

AIST, KDDI R&D Laboratories and NTT are promoting collaborative project G-lambda to establish a standard Web services interface between Grid and network resources. At iGrid2005, we conducted an experiment in cooperation with NICT, where our Grid scheduling system co-allocated computing and network resources with advance reservation through the interface. In this report, our co-allocation system model, the Web services interface (GNS-WSI) between GRS (Grid Resource Scheduler) and NRM (Network Resource Management System), overview of GRS which AIST has developed and NRM which KDDI R&D laboratories developed, and overview of the experiment are described.

1. はじめに

グリッド基盤技術の成熟により, 異なる組織により管理される地理的に分散した計算資源上での大規模高

性能計算が実行可能になった. しかしながら, 現状では分散した計算資源上で実行される大規模高性能計算の多くは, 通信遅延や帯域の不安定さにより著しい性能低下を引き起こすことのないパラメタ探索型やマスター・ワーカー型のアプリケーションプログラムに限られる. これは, 以下により広域分散グリッド環境で高性能計算を実行するのは困難なためである.

インターネットにより容易に広域計算が実行できるが, ベストエフォート型の通信を行うため, 各アプリケーションプログラムの求める通信性能を保証できない. 光ネットワーク技術を用いることで高品質通信が

^{†1} 産業技術総合研究所 National Institute of Advanced Industrial Science and Technology (AIST)
^{†2} 株式会社 KDDI 研究所 KDDI R&D Laboratories Inc.
^{†3} 日本電信電話株式会社 NTT Network Innovation Laboratories
^{†4} 情報通信研究機構 National Institute of Information and Communications Technology (NICT)

提供可能であるが、現状では計算に必要なネットワーク資源の構成・提供は自動化されておらず、E-mailやFAX等により人同士で事前に交渉し、設定する必要がある。また、グリッド“スーパースケジューラ”は異なる組織間の分散した計算資源を必要に応じて確保・提供するものであるが、まだ実用段階に至っていない。現状では、他の組織の計算資源を利用するにはネットワーク資源同様、E-mail等による事前の交渉が不可欠である。

産業技術総合研究所、KDDI研究所、日本電信電話は共同で、ネットワーク資源の事前予約を行うための標準ウェブサービスインタフェースを定めるG-lambdaプロジェクト¹⁾を推進している。このインタフェースの有効性を確認するために、計算資源と光バスネットワーク資源を事前予約によりコアロケーションする実験を、情報通信研究機構の協力を得て行った。

本稿では、前提とする事前予約のためのコアロケーションシステムモデル、G-lambdaプロジェクトにおけるグリッドでのネットワークサービスインタフェースの規定、産総研が開発した計算・ネットワーク資源を同時予約するスケジューラの概要、KDDIが開発したGMPLSで制御される光バスネットワークの提供機構の概要と、これらにより実現されたコアロケーションシステムプロトタイプ*iGrid2005*での実証実験について報告する。

分散する計算資源間には高品質かつ安定した通信を可能にする光バスネットワークを用い、IETFで標準化が進められているGMPLS (Generalized Multi-Protocol Label Switching)²⁾で制御することにより通信環境を柔軟に構成して提供する。GMPLSは分散した光クロスコネクタ(OXC)等のネットワーク機器を制御するプロトコル群である。必要に応じて計算資源間にユーザ専用の光バスネットワークを張ることにより、通信性能を保証することができる。

次に、G-lambdaプロジェクトにおいてグリッドで柔軟に光バスネットワークを構築するため、ウェブサービスベースのネットワークサービスインタフェースGNS-WSI (Grid Network Service - Web Services Interface)を規定し、必要なネットワークの予約、修正、解放、予約状況の取得を可能にした。

また、分散した計算資源とネットワーク資源の事前予約によるコアロケーションを実現するため、GNS-WSIを用いた先行予約ベースのスケジューリングシステムを開発し、2005年9月に米国サンディエゴで開催された*iGrid2005*においてそのプロトタイプシステムの実証実験を行った。コアロケーションシステムは、主に両資源のコアロケーションを行うグリッド資源スケジューラ(GRS: Grid Resource Scheduler)と光バスネットワーク資源を管理するネットワーク管理システム(NRM: Network Resource Management system)からなる。実証実験では、ユーザがGUI(Graphical

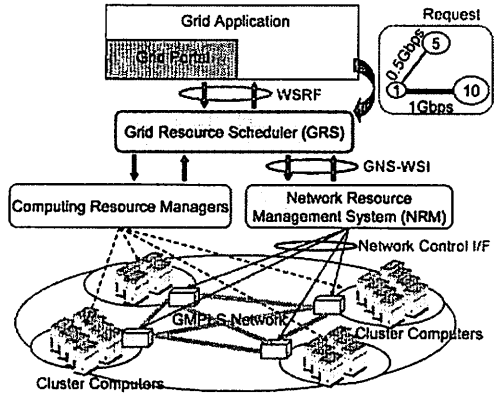


図1 コアロケーションシステムの概要。

User Interface)から簡単な操作で計算機とネットワークに関する要求をGRSに送り、GRSが適宜NRMに問い合わせながら適切な計算・ネットワーク資源を同時に事前予約し、予約された資源上でユーザの分子動力学アプリケーションを実行した。実験環境には、日本の6都市に配備されたPCクラスタとその間のGMPLSネットワークを用いた。本実験により、事前予約に基づくコアロケーションシステムとそのインタフェースの有効性・実現可能性を示した。また、資源予約サービスに関する技術的課題を明らかにした。

2. ウェブサービスベースのコアロケーションアーキテクチャ

図1にウェブサービスベースのグリッドコアロケーションシステム概念図を示す。コアロケーションシステムはグリッド資源スケジューラ(GRS)、ネットワーク資源管理システム(NRM)、計算資源管理システム(CRM)、分散する実計算資源と計算資源間を結ぶGMPLSで制御されるネットワークで構成される。GRSがユーザの要求に応じ、各サイトで計算資源を管理するOpenPBS³⁾やN1 GridEngine⁴⁾等のCRMと、ネットワーク資源を管理するNRMにウェブサービスを介して連携し、適切な計算資源と資源間の光バスネットワーク帯域を選択し、事前予約する。

提供するネットワーク資源には、通信性能が保証可能な光バスネットワークを用いる。また、光バスネットワークの制御にはIETFで標準化が進められているGMPLS技術を用いる。GMPLSを用いることにより、分散した計算機や記憶装置等を要求されるバンド幅で適宜接続することができるため、計算効率と利用効率の向上が期待できる。

グリッドにおいて提供される資源は地理的に分散しているため、資源間を繋ぐネットワークは複数の都市や国を跨って接続され、複数のネットワークオペレータ

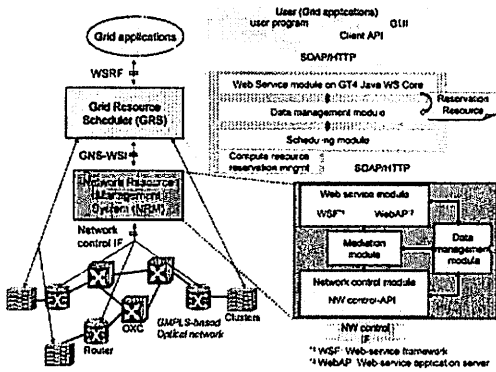


図2 コアロケーションシステムアーキテクチャ

が存在することが予想される。よって、GRSとNRM間の連携が必要とされる情報交換には、ネットワークオペレータ間で共通のインタフェースを用いることが重要である。G-lambdaプロジェクトでは、グリッドネットワークサービスウェブサービスインタフェース(GNS-WSI)としてネットワーク資源の確保のための基本的な機能を定義する。

2.1 GNS-WSI

2拠点間を接続する光パスネットワークの提供を前提とした初期GNS-WSIとして、ウェブサービス記述言語WSDLにより以下のポーリングベースのオペレーションを定義した。

netResourceReservation GRSから要求された仮想パスネットワークをOXC間の物理リンクに割り当て、事前予約する。GRSはエンドポイントの拠点ID、バンド幅、予約期間、遅延や耐故障性に関する要求をNRMに送信する。複数のパスを同時に要求することもできる。

netResourceModification 予約されたパスのバンド幅や予約期間等の予約情報を変更する。

netResourceRelease 予約されたパスを解放する。

netResourceReservationStatusQuery 要求された／予約されたパスの現在のステータスを返す。予約、変更、解放要求の結果は、このオペレーションにより取得することができる。また、確保されたパスが予約した時刻に実際に利用可能になったか否かの情報も取得することができる。

netResourceQuery GRSの指定したパスが指定した期間に利用可能かどうかを返す。

netAvailableResourceQuery GRSが指定した期間に利用可能なパスをすべて返す。

2.2 グリッド資源スケジューラ (GRS)

GRSはユーザの要求に応じて適切な高品質計算環境を提供するため、計算資源は各サイトのCRMに、ネットワーク資源はGNS-WSIを介してNRMに問い合わせる適切な資源を同時に確保する。GRSでは

事前予約、解放、予約情報取得のためのオペレーションを提供している。

GRSのプロトタイプとして、Globus Toolkit 4 (GT4)^{5),6)}を用いてWSRF(Web Services Resource Framework)⁷⁾仕様に基づく計算・ネットワーク資源を同時に確保するスケジューラを開発した。WSRFはウェブサービスにおけるステートフルリソースのモデリングとアクセスのための仕様であり、OASISでその標準化が進められている。

コアロケーションシステムを構成する主要モジュールGRSとNRMの詳細を図2に示す。GRSは次のモジュールで構成される。

Client API, GUI ユーザにプログラミング／グラフィカルインタフェースを提供する。

ウェブサービスモジュール WSRF に基づく予約サービスインタフェースをユーザに提供する。このモジュールで個々の資源予約要求に対する予約リソースを生成し、データ管理モジュールへ格納する。ユーザにはこの予約リソースにアクセスするためのエンドポイントリファレンスを返す。これにより、ユーザは現在の予約状況や予約結果を取得することができる。

予約リソース 個々の予約要求に対する予約サービス情報を保持する。

データ管理モジュール 各予約処理の状況を管理する。格納データの持続性を保証している。

スケジューリングモジュール ユーザの計算資源とネットワーク資源に関する要求と、予約期間とデッドライン等の時間に関する要求をもとに、NRM、CRMと連携して適切な資源を割り当てる。要求された資源が先行予約できれば「予約成功」、できなければ「予約失敗」のステータスをデータ管理モジュールへ格納する。このステータスはウェブサービスモジュールを介してユーザに送られる。予約が失敗した場合は、ユーザが新しい条件を再度送信することになる。

計算資源予約管理 各サイトのCRMで管理される予約情報を管理する。

2.3 ネットワーク資源管理システム (NRM)

NRMではGMPLSで制御される光パスネットワーク資源をGNS-WSIを介して提供する。KDDI研究所が開発したNRMは、下記の3つの機能モジュールから構成され、GRSに対してGMPLSによるネットワークサービスを提供する⁸⁾。

ウェブサービスモジュール GNS-WSIを介してGRS等のグリッドブローキングシステムにGMPLSで制御される光パスネットワーク資源を提供する。これにより、GRSはネットワーク資源の発見、光パスの事前予約、予約の修正等を行うことができる。

中間モジュール GRSからの要求と実際のネットワーク資源間を仲介し、GMPLS制御光ネットワークを

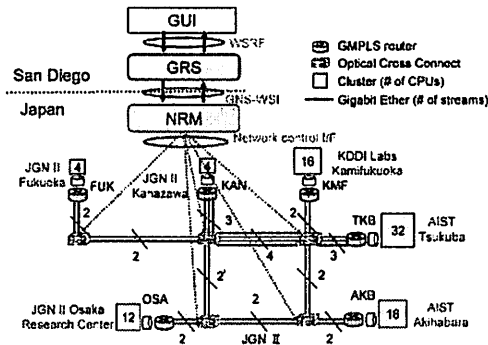


図3 iGrid2005 実験環境。

計算資源間のパスとして仮想化する。また、光パスのスケジューリングも行う。
 ネットワーク制御モジュール ネットワーク制御インタフェースを介して GMPLS ルータを制御し、要求された光パスを利用可能にしたり、それらの状況や性能をモニタする。

3. iGrid2005 実証実験の概要

2005年9月に米国カリフォルニア大学サンディエゴ校で開催されたグリッドに関するデモを中心とした国際会議 iGrid2005 において、本グリッドコアロケーションシステムの実証実験を行った。実験では、ユーザがグリッドの資源に関する要求を GUI から GRS へ送信し、GRS が NRM と連携して柔軟に適切な計算・ネットワーク資源を同時に事前予約する。予約が成功すると、予約された時刻に確保された資源を用いて科学アプリケーションプログラムを実行する。

図3に実証実験環境を示す。実験では、GRSとGUIはサンディエゴの会場のPC上で、NRMは大手町のWS上で実行した。図3に示すように、日本の6都市(つくば、秋葉原、上福岡、金沢、大阪、福岡)に配置された計算資源とそれを結ぶGMPLSネットワークを用いた。本実験では、各計算資源サイトでは事前予約機能を持たないCRMを用いたため、全ての計算資源をコアロケーションシステムで占有し、先行予約の管理をGRSで行った。GMPLSネットワークは光クロスコネクとGMPLSルータで構成されている。予約した資源上では、グリッドミドルウェア Ninf-G2⁹⁾と Globus Toolkit 2 (GT2)⁵⁾を用いて実装された分子動力学アプリケーションプログラムを実行する。実証実験にあたり、産総研、KDDI研、NTTがGNS-WSIを規定し、産総研が開発したGRS、KDDI研が開発したNRMを用いた。また、NICTが運営するJGNIIのネットワークテストベッド上にKDDI研、NTTが動作環境を構築した。

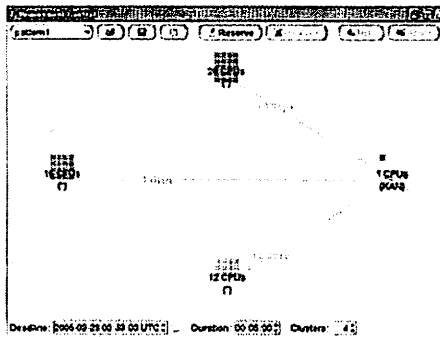


図4 資源予約入力画面。

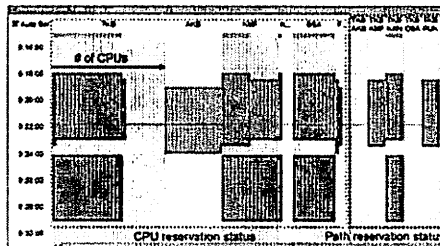


図5 資源予約タイムテーブル (左: 計算資源の予約CPU数, 右: 光パスネットワークの予約帯域)。

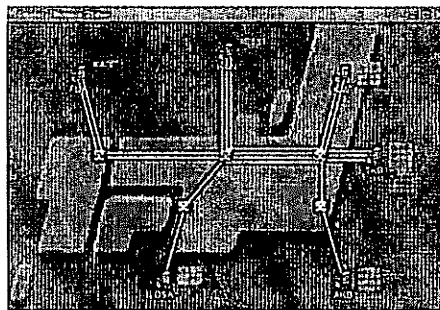


図6 予約状況を表示する地図画面。

3.1 実証実験の詳細

図4, 図5, 図6, 図7に実証実験の表示画面を示す。実証試験では、図4の資源予約入力画面から簡単な操作でクラスタ数と各クラスタのCPU数、その間のネットワーク構成と帯域、および予約期間とデッドラインに関する要求を入力し、GRSに送信する。GRSは通常要求された資源に対して実際の計算資源サイトとその間のネットワークを割当ててるが、ユーザが一部の計算資源サイトを指定することも可能である。

ユーザの要求がGUIからGRSに送信されると、GRSがNRMとウェブサービスを介して連携して適切な計算・ネットワーク資源が事前予約し、その結果

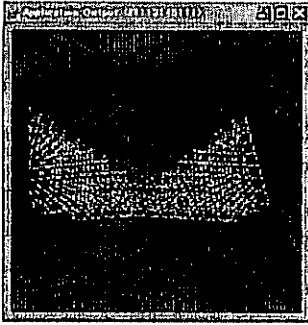


図7 分子動力学アプリケーション実行結果画面。

が図5の資源予約タイムテーブルに反映される。ここでは、画面の左側に各サイトで予約されたCPUの数、右側に予約された計算資源間の光バスネットワークバンド幅を示している。タイムテーブルの縦軸は時刻を示し、中央の横線が現在の時刻を表している。

図6は現在のネットワークの予約状況を表示する地図画面を示す。地図上に表示されるサイト間のバスネットワークは、実際の物理ネットワーク上でGMPLSの制御パケットを拾って表示させたものである。予約した時刻になると、NRMからの制御により割当てられた光バスが利用可能となり、地図上にバスが出力される。

図7は先行予約で確保され、動的に構築された高品質計算環境において分子動力学シミュレーションを実行した結果である。この実験により、地理的に分散した計算資源とその間の高品質ネットワークをGNS-WSIを介して統一的に扱い、予約ベースで必要に応じて柔軟に確保・提供できることを実証した。

3.2 議論

3.2.1 GRS-NRM間のSOAP通信遅延

GRSとNRMが連携してコアケーションする際、GRSとNRM間で複数のSOAPメッセージが送受信される。この際、通信遅延に加えてSOAP通信にかかる遅延が加わるため、本実験のように通信遅延の大きい環境では、ユーザがGUIから資源要求を送信してから結果を受け取るまでの時間が長くなってしまふ。iGrid2005では数秒程度と許容できる時間であったが、大規模環境でより複雑な資源要求がされる場合は深刻な問題となる。GRS-NRM間の通信回数をできるだけ少なくし、大規模環境にも耐えうるプロトコルを開発する必要がある。

3.2.2 光バスネットワークのモニタリング

本実験環境はGUI、GRS、NRM等のプログラムモジュールと6拠点に分散したクラスタおよびその間のGMPLSネットワークからなり、組み上げて安定した操作が可能になるまで多くの人手と時間を要した。特に、GMPLS光バスネットワーク構築におけるデバッグ手法がまだ未熟であるため、予約した時刻に光バス

ネットワークが利用可能にならない等の不具合が発生した際にその原因究明に時間を要した。しかしながら、実証実験を目的として開発したネットワークの予約状況を表示する地図画面(図6)では、GMPLSでバスネットワークを利用可能にする際にGMPLSルータ間で送受信されるメッセージパケットを直接拾って表示させていたため、比較的容易に原因箇所を特定することができた。よって、光バスネットワークレベルのモニタリングはデバッグ手法として非常に有効であり、今後その手法の改良と標準化が重要となる。

3.2.3 光バスネットワークのルーティング

GMPLSネットワークでは、OSPF(Open Shortest Path First)ベースの経路計算手法を用いて動的に光バスの経路を設定することができる。このGMPLSルーティングプロトコルの機能により、光ネットワークの障害耐力と資源利用効率を向上させることができる。しかしながら、GRS-NRM間のネットワークサービスで事前にきめられた光バスの通信遅延に関する条件が、動的なルーティングにより変化してしまうため、本実験ではこの機能を利用しなかった。よって、本実験ではGMPLSネットワークにおける各光バスの経路を明示的に指定し、固定することにした。GMPLSネットワークにおける動的なルーティング機能を有効に利用するためには、グリッドアプリケーションの通信遅延に関する要求条件を明らかにする必要がある。

3.2.4 耐故障性とエラー通知機能

資源予約の耐故障性はコアケーションシステムの実用化に向けた重要な課題である。iGrid2005の実験では、光バスの先行予約が成功したにも関わらず、設定上のバグやGMPLSメッセージングのエラーによる不具合が発生し、予約された時刻にバスネットワークが利用可能にならないことがあった。一方、Ninf-GとGT2で実装されたアプリケーションプログラムにはこのエラーが検知できないため、プログラムの実行を開始して計算資源間の通信が確立されないまま待ち続け、複数プロセスが計算資源上に残ってしまった。この問題を解決するため、実験ではアプリケーションの実行開始時間を10秒程度遅らせるとともに、GRSがNRMに予約したネットワークが予約時刻に利用可能になったかどうか問い合わせ、何らかのエラーが発生した場合はGUIからユーザに通知するモニタ機能を用意した。今後、資源予約の耐故障性を強化するとともに、致命的なエラーが発生した場合に事前にユーザに通知したり、その際にユーザに何らかの補償をするメカニズムを開発する必要がある。

4. 関連研究

グリッドで高品質なネットワークを提供する試みがいくつかある。CANARIE¹⁰⁾はカナダの先進的インターネット開発の推進を目指した組織であり、教育目

的の光バスネットワークとして CA*net4 を配備している。CA*net4 ネットワークの制御において、エンドユーザがネットワークをソフトウェアオブジェクトとして扱い、光バスの配備や再構成が容易にできるようにするため、UCLP (User Controlled LightPaths) ソフトウェアを開発している¹¹⁾。UCLP はウェブサービスを介してユーザが指定する個々の光バスネットワークを直接操作するものであり、事前予約機能や計算機とのコアロケーション機能は用意されていない。

VIOLA は本研究同様、グリッドで計算資源とネットワーク資源を同時に提供することを目的としている¹²⁾。VIOLA では GGF において標準化がすすめられている WS-Negotiation/-Agreement 仕様に基づくメタスケジューラを UNICORE¹³⁾ ベースのグリッド環境で開発することを提案している。しかしながら、実装の詳細は明らかでない。

グリッドスーパースケジューラでは Silver¹⁴⁾、CSF¹⁵⁾、NAREGI スーパースケジューラ¹⁶⁾ など、複数開発されているが、その多くが実用段階には至っていない。また、計算資源の割当てが主体であり、性能が保証されたネットワーク資源と計算資源のコアロケーションを実現しているものはない。

5. まとめと今後の課題

本稿では、事前予約に基づくコアロケーションシステムモデル、産総研、KDDI 研、NTT の共同プロジェクト G-lambda におけるグリッドでのネットワークサービスインタフェースの規定、産総研が開発した計算・ネットワーク資源を同時予約するスケジューラの概要、KDDI 研が開発した GMPLS で制御される光バスネットワークの提供機構の概要と、これらにより実現されたコアロケーションシステムプロトタイプの iGrid2005 での実証実験について報告した。iGrid2005 実験では、ユーザが GUI から簡単な操作で計算機とネットワークに関する要求を入力することで、GRS が GNS-WSI を介して NRM と連携して適切な計算・ネットワーク資源を同時に事前予約し、柔軟かつ容易に高品質な計算環境を提供可能であることを実証し、その有効性を示した。また、資源予約サービスに関する技術的課題を明らかにした。

今後は G-lambda において複数 GRS と複数 NRM を考慮して GNS-WSI を洗練していくとともに、GGF や IETF 等において標準化を進める。また、海外の光ネットワークとの連携実験も検討している

謝辞 iGrid2005 実証実験にあたり、多大なご協力をいただいた谷村勇輔氏、武宮博氏、岡崎史裕氏 (産総研)、鈴木正敏氏、鶴沢宗文氏 (KDDI 研究所)、大原拓也氏、築島幸男氏 (NTT)、下條真司氏、秋山豊和氏 (NICT 大阪 JGNII RC)、小林正朋氏 (東京大学) に深く感謝いたします。

参考文献

- 1) G-lambda. <http://www.g-lambda.net/>.
- 2) GMPLS. RFC 3945.
- 3) OpenPBS. <http://www.openpbs.org/>.
- 4) N1 Grid Engine. <http://gridengine.sunsource.net/>.
- 5) Foster and Kesselman: Globus: A Metacomputing Infrastructure Toolkit., *International Journal of Supercomputer Applications*, pp. 115-128 (1997).
- 6) Foster: Globus Toolkit Version 4: Software for Service-Oriented Systems, *IFIP International Conference on Network and Parallel Computing*, Springer-Verlag LNCS 3779, pp.2-13 (2005).
- 7) OASIS Web Services Resource Framework (WSRF) TC: Web Services Resource 1.2 (WS-Resource) (2005). <http://www.oasis-open.org/committees/wsrf/>.
- 8) Hayashi, Miyamoto, Otani, Tanaka, Takefusa, Nakada, Kudoh, Nagatsu, Sameshima and Okamoto: Managing and Controlling GMPLS Network Resources for Grid Application, *Optical Fiber Communications Conference (OFC) 2006* (to appear).
- 9) Tanaka, Nakada, Sekiguchi, Suzumura and Matsuoka: NinF-G: A Reference Implementation of RPC-based Programming Middleware for Grid Computing, *Journal of Grid Computing*, pp.41-51 (2003).
- 10) CANARIE. <http://www.canarie.ca/>.
- 11) Boutaba, Golab, Iraqi, Li and Arnaud: Grid-Controlled Lightpaths for High Performance Grid Applications, *Journal of Grid Computing, Special Issue on High Performance Networking* (2004).
- 12) Barz: Dynamic allocation of network resources in VIOLA, *VIOLA workshop* (2005).
- 13) Romberg: The UNICORE Architecture: Seamless Access to Distributed Resources, *Proceedings of the Eighth IEEE International Symposium on High Performance Distributed Computing*, pp.287-293 (1999).
- 14) Moab Grid Scheduler (Silver) Administrator's Guide version 4.0. <http://www.clusterresources.com/products/mgs/docs/>.
- 15) Community Scheduler Framework. <http://sf.net/projects/gcsf>.
- 16) 畑中, 中野, 井口, 大野, 佐賀, 秋岡, 中田, 松岡: OGSA アーキテクチャに基づく NAREGI スーパースケジューラの設計と実装, *情報処理学会研究報告 2005-HPC-102*, pp.33-38 (2005).