



6 iSCSI と FCoE による ストレージ構築 ～ストレージネットワークの進化～

小口正人 ●お茶の水女子大学

ストレージはクラウドを支える重要な存在である。ストレージシステムの進化に伴い、クラウド内も含め、サーバとストレージの間をネットワークで接続する方式がいくつかの形態で登場し、利用されてきている。本稿ではイーサネットをベースにサーバとストレージが接続される iSCSI (Internet Small Computer System Interface) と FCoE (Fibre Channel over Ethernet) の 2 方式を概観する。

きない。また管理の手間がかかるため、数が増えると管理コストが高くなる。そこで図-1の右側のようにストレージへのアクセス部分をネットワークとし、この上を SCSI プロトコルが流れるような仕組みが作られた。これが SAN (Storage Area Network) である。すなわちサーバ間通信(フロントエンド)は LAN で接続され、ストレージアクセス(バックエンド)は SAN で接続される形態が、現在のサーバとストレージシステムの標準的な形となった。

iSCSI

* iSCSI 誕生の背景

ストレージシステムは従来、図-1の左側のようにサーバとストレージシステムが密に接続され、SCSI (Small Computer System Interface) などのプロトコルでアクセスが行われていた。この形態を DAS (Direct Attached Storage) と呼ぶ。しかしサーバが特定のストレージにしかアクセスできないと、偏りが生じストレージ全体を効率よく使うことがで

SAN としては FC (Fibre Channel) が広く用いられている。FC は文字通り、光ファイバをベースとしたネットワークである。しかし FC はスイッチもケーブルも高価なものであり、また光ファイバであるため銅線で接続されたイーサネットなどの LAN と比べると扱いにくい。さらに接続距離に制限があるなどといった問題もあり、機器の種類が豊富でコストパフォーマンスがよい、汎用のネットワーク技術をストレージ接続に活かすことが望まれている。

そのような要望に応じて iSCSI は登場し、2003 年 2

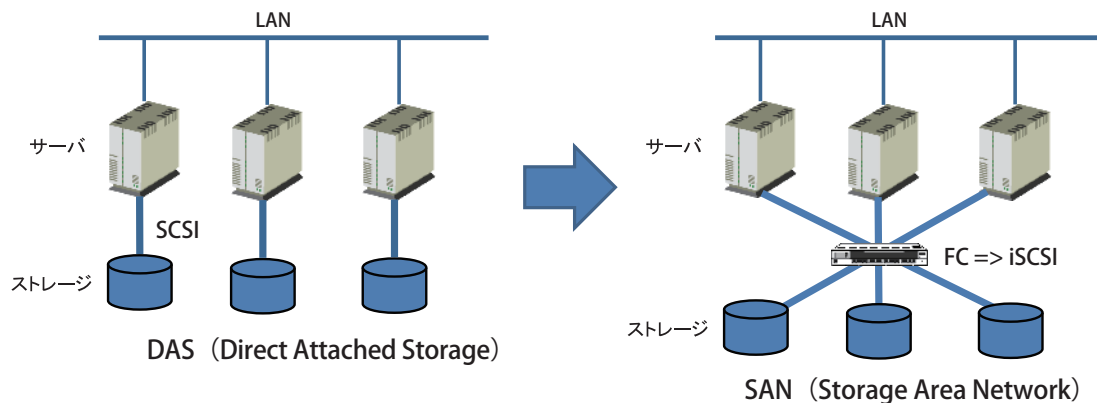


図-1 DAS から SAN への進化



月に IETF により規格が承認された¹⁾。iSCSI は SCSI プロトコルを TCP/IP ネットワーク上に載せたものであり、汎用ネットワーク環境であるイーサネットとその上の TCP/IP が動いていれば利用することができる。

* iSCSI アーキテクチャ

iSCSI は SCSI コマンドや読み書きするデータを、TCP/IP パケットにカプセル化してネットワーク越しに転送する規格である。iSCSI のプロトコルスタックを図-2 に示す。ストレージアクセスを行うサーバ側(イニシエータ)では、アプリケーションがローカルストレージへアクセスするかのよう SCSI コマンドを発行するが、iSCSI 層がこれを受け取って TCP/IP のパケットに詰め込み、イーサネット等のネットワークを介してストレージ側(ターゲット)へ送る。ターゲットでは TCP/IP パケットから SCSI コマンドを取り出してストレージアクセスを行い、結果を逆の経路でイニシエータへ返す。

以下ではこのカプセル化が具体的にどのように実現されているか、少々細かく説明する。SCSI コマンドを TCP/IP パケットに詰め込むために作られる iSCSI のブロックを iSCSI PDU (Protocol Data

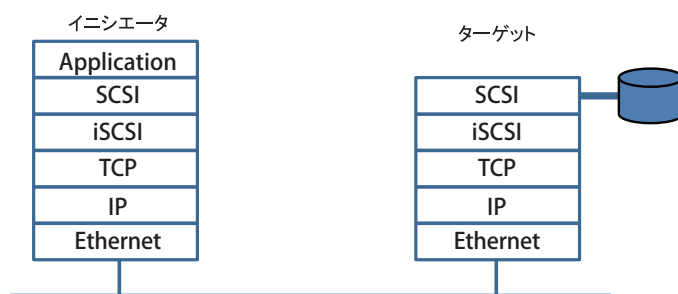


図-2 iSCSI のプロトコルスタック

Unit) と呼ぶ。この構成を図-3 に示す。この図では 1 つの iSCSI PDU が 1 つの TCP/IP パケットにカプセル化されているが、大きさが TCP の MSS (Maximum Segment Size) を超える場合には分割されて、各々のセグメントに TCP/IP のヘッダが付く。図-3 には Read の SCSI コマンドメッセージがカプセル化された場合の例が示されている。

iSCSI PDU は 48 バイトの BHS (Basic Header Segment) といくつかのオプションフィールドからなる。SCSI プロトコルによりストレージとの間で Read/Write のデータが転送される場合には、Data Segment フィールドが使われる。

SCSI コマンドメッセージをカプセル化する場合、BHS の後ろ 16 バイトのフィールドには、SCSI CDB

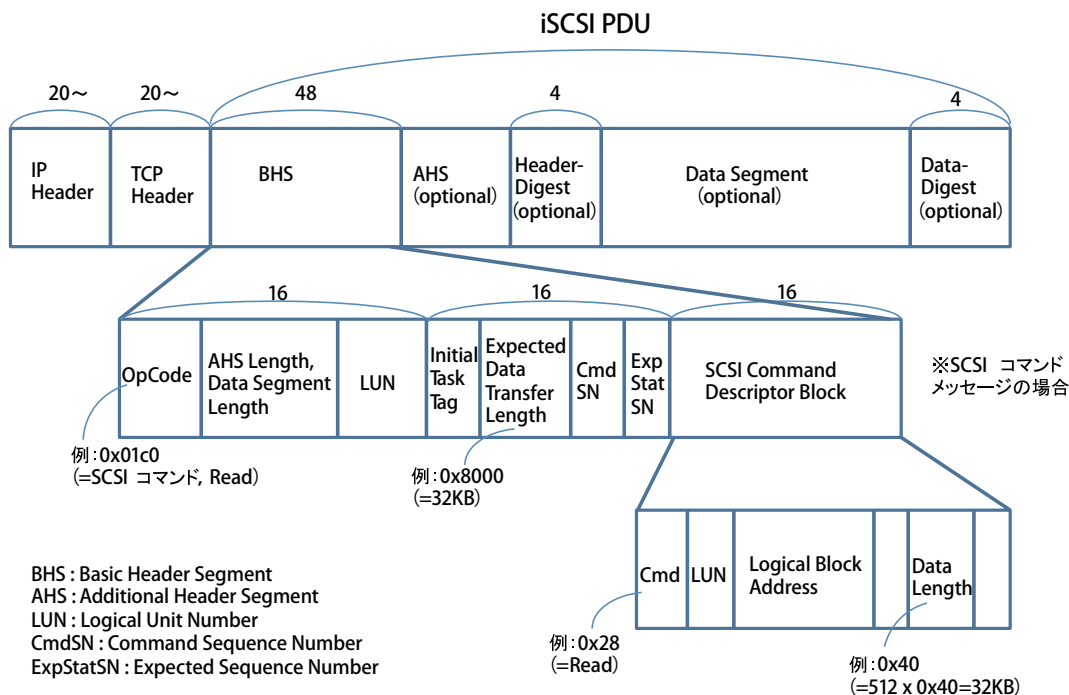


図-3 iSCSI PDU (Protocol Data Unit)



クラウドを支えるデータストレージ技術

(Command Descriptor Block), すなわち SCSI コマンドがそのまま入る。ただしその中のいくつかの情報は読み出されて、BHS の前の方のフィールドに埋められる。たとえば BHS の先頭は OpCode というフィールドで、ここには SCSI CDB の中身は何であるか(たとえば SCSI Read コマンドなど)が読み出されて書き込まれる。Expected Data Transfer Length フィールドにも、SCSI CDB 中で指定している Read/Write のデータの大きさが読み出されてバイト単位で記録される(SCSI CDB 中では 512 バイト単位)。

* 代表的な iSCSI 実装の紹介

前述のように iSCSI は、イーサネットとその上の TCP/IP という汎用のネットワーク環境があれば利用することができる。最新の Linux や Windows においては、iSCSI は容易に利用できるようになっている。

Linux では当初、UNH (University of New Hampshire) により作られた参照実装が多く使われていたが、その後いくつかの実装が現れて改良が進められた。イニシエータドライバでは、Linux-iSCSI (sfnet) と Open-iSCSI がよく使われるようになってきたが、この両プロジェクトは合併し、Linux2.6.16 以降で動作する Open-iSCSI ドライバとして利用されている²⁾。またターゲットドライバとしては、iSCSI Enterprise Target が広く用いられている³⁾。一方 Windows 7 や Windows Server 2008 では、標準で Microsoft iSCSI イニシエータドライバが搭載されている。ターゲットドライバは Windows Storage Server 2008 に搭載されている。

このような導入の敷居の低さにより、iSCSI の利用は拡大してきている。一方、性能面を考慮して iSCSI 用のハードウェアも現れて利用されている。ターゲットとしては iSCSI インタフェースを持ったストレージシステムが、比較的高額なものから、量販店で売られるような割合安価なモデルまで幅広く販売されている。またイニシエータは、TCP/IP の処理を NIC にハードウェア実装する TOE (TCP Offload Engine) や、iSCSI の処理までハードウェア実装した iSCSI HBA (Host Bus Adaptor) が存

在し、目的に応じて使い分けられている。

現在の状況を見ると、iSCSI の技術はすでに発展のフェーズに乗ったと言えよう。今後はクラウドを始めとしたサーバシステムにおける利用から、エンドユーザの個人利用まで、幅広く用いられていくものと考えられる。



* FCoE 誕生の背景

「iSCSI 誕生の背景」で述べたように、現在のサーバシステムはフロントエンドを LAN で接続され、バックエンドを SAN で接続される形が一般的になった。LAN のアーキテクチャはイーサネットであり、ギガビットイーサネットから 10Gbps イーサネットへと移行しつつある。一方 SAN として主流なのは、やはり FC である。

FCoE は、INCITS (International Committee for Information Technology Standards) の FC を担当する T11 技術委員会により、FC-BB-5 の一部として標準化が行われた^{4), 5)}。FCoE は FC のフレームをカプセル化し、仕様に従って拡張された 10Gbps イーサネット上で運ぶというものである。

FCoE の目的を一言でいうと FC とイーサネットの統合であり、FCoE によって FC のフレームとイーサネットのフレームを、同じインタフェースを通し同じネットワークで運ぶことを目指した。これはユニファイド I/O と呼ばれる。

現在のサーバは、インタフェースをいくつも持っている。たとえば図-4 の左側のように、制御用のイーサネット、データ転送用のイーサネット、ストレージに繋がる FC があり、それぞれが二重化されているとすると、インタフェースは全部で 6 枚となり、6 本のケーブルがサーバから出てネットワークに繋がる形となる。これに対し、これらの FC とイーサネットをすべて FCoE で統合すると、図-4 の右側のように、二重化も含めて 2 枚のインタフェース CNA (Converged Network Adaptor) と 2 本のケーブルで済む。

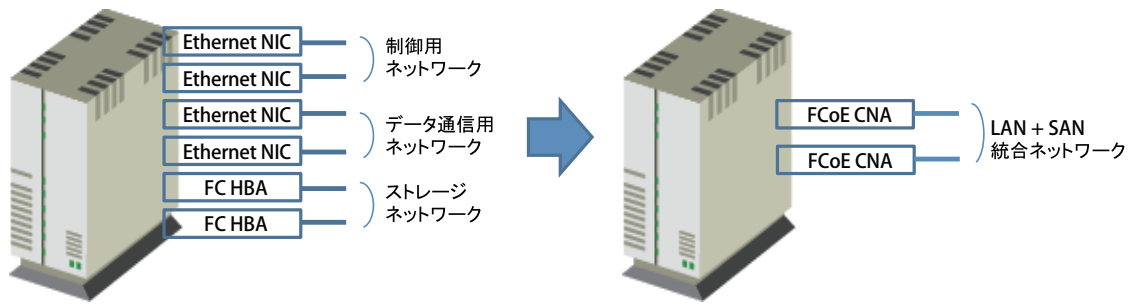


図-4 FCoE による FC とイーサネットの統合

FCoE が最も期待されている活躍の場はデータセンタである。データセンタでは、FC とイーサネット、さらにはインフィニバンドなど複数の独立したネットワークが運用されており、サーバは複数のインタフェースを持ち、複数のネットワークと接続する形態となっている。しかしながら複数のネットワークが別々に存在するのは、導入コストの面でも管理コストの面でも望ましくない。さらに、ラックマウントやブレードタイプのサーバの場合、物理的な制約から複数のインタフェースを持つことができない場合もあり、その結果、共有ストレージに直接接続されたサーバは、全体のうち一部のみに限られることとなっている。すなわち FCoE を用いれば、統合によりコストが下げられるだけでなく、機器の可用性を上げることができる。さらには FCoE を用いた統合でインタフェースやスイッチの数を減らすことによる省電力効果も期待されている。

FCoE の階層構造を 図-5 に示す。FCoE は iSCSI と異なり、TCP/IP の上に FC を載せるわけではない。トランスポート層 (TCP, UDP) もインターネット層 (IP) も用いず、第 1 層と第 2 層にあたるイーサネット部分のみを用い、その上に FC の FC-2 層以上を載せる。ただしこの 1, 2 層も既存のイーサネットそのままではなく、拡張された仕様に基づく CEE (Converged Enhanced Ethernet) が用いられる。この CEE については次節で述べる。

FCoE において、FC-2 層以上は変更を加えることなく、既存の FC を用いることができる。したがって FC のネームサービスやゾーニングなどの管理機能や SAN のアプリケーションをそのまま利用できる。これは FC フレームをイーサネットフレームにカプセル化することにより実現されている。

以下ではこのカプセル化がどのように実現されているか、少々細かく説明する。イーサネット、FC および FCoE のフレームを 図-6 に示す。

VLAN 対応のイーサネットフレームは、宛先と送信元 MAC アドレス各 6 バイト、VLAN タグ 4

* FCoE アーキテクチャ

イーサネットと FC、そしてこれを統合した

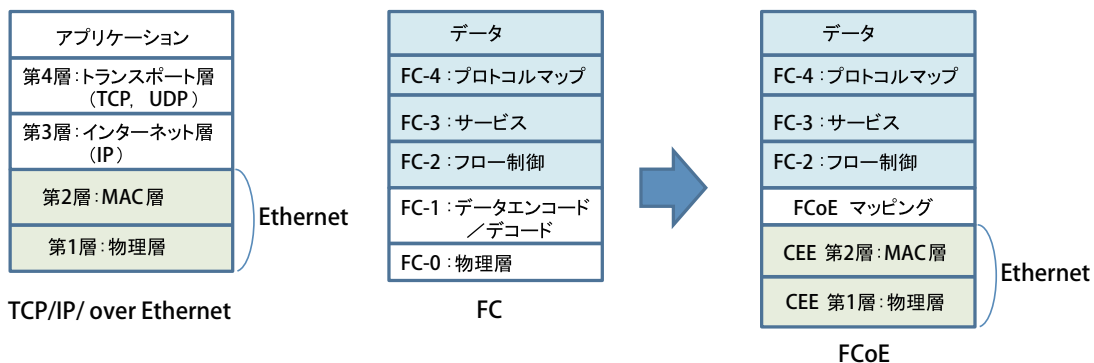


図-5 イーサネットと FC を統合した FCoE のプロトコルスタック



クラウドを支えるデータストレージ技術

バイト、タイプ/フレーム長フィールド2バイトによりヘッダが構成されている。データの後ろにはトレーラとしてFCS (Frame Check Sequence) が4バイト続く。そしてFCoE ではこのデータ部分にFCのデータが埋め込まれる。タイプ/フレーム長フィールドはFCoEを表す0x8906となっており、これによりFCoEフレームであることが識別できる。データの最大長は2,112バイトとなっており、ヘッダ等も含めたFCoEフレームは最長2,148バイトとなる。したがってジャンボフレームなどに対応するイーサネット環境が必要になる。

* CEE について

FCoEはFCフレームがイーサネットヘカプセル化されているため、FCから見ると下がイーサネットになっていることは認識できない。しかしながらFCはパケットが喪失せず宛先に届くことを前提として設計されており(ロスレス)、パケットロスがあることを前提にトランスポート層で到達保証を行うTCP/IP/ over イーサネット環境とは異なる。したがってFCをイーサネット上に載せようとした場合、ロスレスのイーサネットが必要となってくる。その

ための規格がCEEである。CEEはIEEEにおいて、DCB (Data Center Bridging) という名前で標準化が進められている⁶⁾。DCBの用途はFCoEだけではなく、FCoEはその上位に載る有力なアプリケーションといえる。

DCBの主要な規格は以下の3つである。

- 802.1Qau : Congestion Notification (CN) : 輻輳通知
- 802.1Qaz : Enhanced Transmission Selection (ETS) : 拡張送信選択
- 802.1Qbb : Priority-based Flow Control(PFC) : 優先度ベースフロー制御

CNはフレームの転送経路において輻輳が発生した際に、これを送信元へ知らせることで送信レートの制御を行う。輻輳を検出したスイッチは輻輳通知メッセージを送信元へ送り、これを受け取った送信元は、フレーム送出を抑制することにより輻輳を緩和させる。

ETSではトラフィックの優先度クラスを設け、各クラスの保証帯域を規定してフレームの配信を行う枠組みである。優先度をプライオリティグループに割り当て、LANやSAN、プロセス間通信など、異なる種類のトラフィックをグループごとに区別し

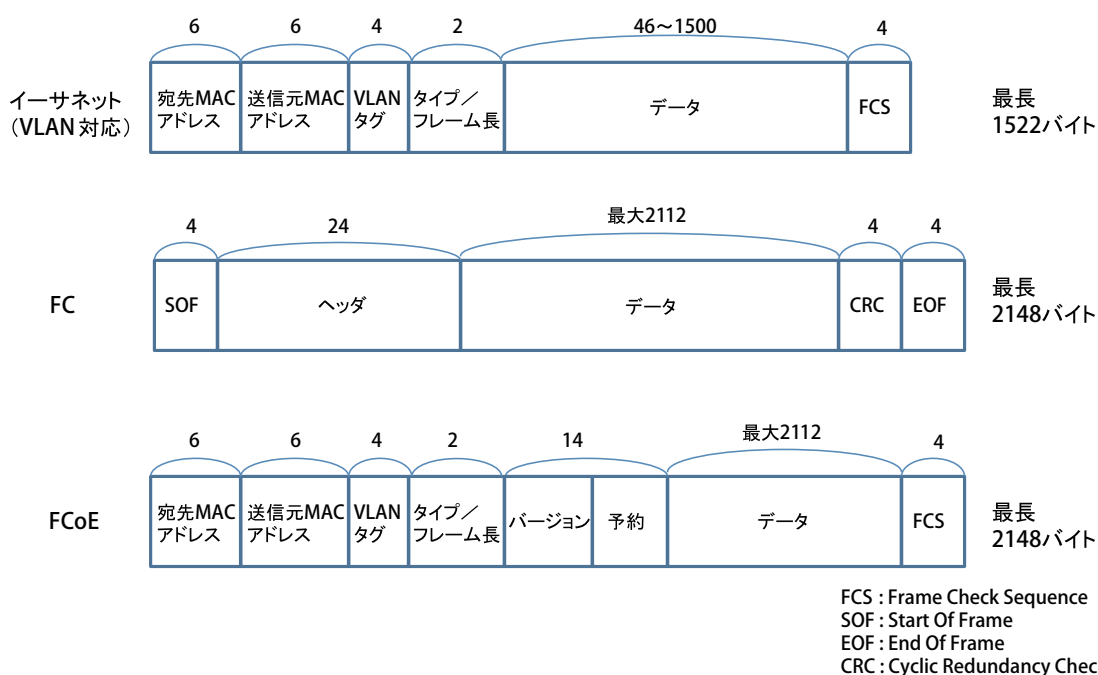


図-6 イーサネット, FC, FCoEのフレーム

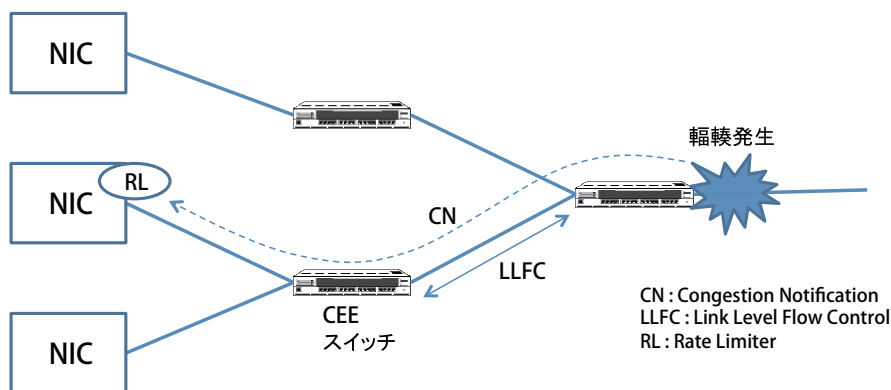


図-7 CN と PFC の動作の概念図

で扱うことができる。

PFCはスイッチにおいて輻輳が発生した際に、その1つ手前のスイッチまたは端末に対してPAUSEフレームを送り、転送を一時的に停止させる。元々802.3x PAUSEという隣接ノード間でフレーム転送を一時的に停止させるフロー制御の規格があるのに対し、これに優先度クラスを対応させて、優先度別にフレーム転送を一時的に停止できるようにしたものである。

フロー制御が行われている様子を図-7に示す。LLFC (Link Level Flow Control) は隣接ノード間のフロー制御で、PFCにより実現される。しかし深刻な輻輳が発生した場合、PFCだけでは一時的な対処しかできないため、CNで元栓を絞ることによって、ロスレスのイーサネットを実現している。

* FCoE の発展方向

FCoEの導入は、まず最初にFCoEのCNAをインタフェースに持つサーバが導入され、FCoE対応のスイッチにより既存のLAN/SANの枠組みに接続されるところから進んでいくであろう。FCoEに対応したストレージが導入されて、エンド・ツー・エンドがFCoEで統合されるようになるまでには、まだ少し時間がかかりそうである。

FCoEによる統合のメリットは前述の通りであるが、果たして統合は進むのであろうか。統合を阻害する要因もいくつか考えられる。まずメリットとして考えられているコスト削減が思惑通り進むかとい

うことである。これは技術が成熟して安定した安価な製品が出てくるかどうかにかかっているであろう。また性能面で考えると、現在8Gbpsが主流となっているFCに対し、10Gbpsイーサネットを用いれば物理的な転送速度は向上するが、FCも16Gbpsへと移行し始めたところである。たとえばSSDのような高速ス

トレージが用いられる場合には、接続ネットワークの速度が性能に直結するため、少しでも速いことが求められる可能性がある。さらに、LANとSANは異なる組織が管理しているようなケースが多いため、これらを統合するには組織を変える必要も出てくるかもしれない。これらの阻害要因が今後どのように変わっていくかが、FCoE普及の鍵といえるであろう。

参考文献

- 1) RFC 3720 : Internet Small Computer Systems Interface (iSCSI) (<http://tools.ietf.org/html/rfc3720>).
- 2) Open-iSCSI (<http://www.open-iscsi.org/>).
- 3) iSCSI Enterprise Target (<http://iscsitararget.sourceforge.net/>).
- 4) INCITS Working Draft Proposed American National Standard for Information Technology, Fibre Channel Backbone-5 (FC-BB-5) Rev. 2.00 (<http://www.t11.org/ftp/t11/pub/fc/bb-5/09-056v5.pdf>).
- 5) INCITS 462-2010, American National Standard, Fibre Channel Backbone-5 (FC-BB-5) (http://www.techstreet.com/standards/INCITS/462_2010?product_id=1724386).
- 6) Data Center Bridging, IEEE 802 Tutorial (http://www.ieee802.org/802_tutorials/07-November/Data-Center-Bridging-Tutorial-Nov-2007-v2.pdf).

(平成23年1月22日受付)

小口正人(正会員) ■oguchi@computer.org

平成2年慶大・理工・電気卒。平成7年東大大学院工学系研究科電子工学専攻博士課程修了。博士(工学)。学術情報センター中核的研究機関研究員、東大生産技術研究所特別研究員、中央大学研究開発機構助教授、お茶の水女子大学理学部情報科学科助教授を経て、平成18年より同教授。ネットワークコンピューティング・ミドルウェアに関する研究に従事。IEEE、ACM、電子情報通信学会各会員。