

O Z : オブジェクト指向開放型分散システムアーキテクチャ

- L L Cタイプ3を活用する通信アーキテクチャ、実装、およびその評価 -

塚本 享治	吉江 信夫	近藤 貴士
電子技術総合研究所	住友電気工業	シャープ

水谷 功	田中 伸明
住友電気工業	松下電器産業

O Zはオブジェクト間の相互関係を分散システム全体で維持しながら複写・移動することを基本とする分散処理システムである。通信し合うオブジェクトの数が非常に多く、相互参照関係が逐次変化する上に、さらに通信相手が移動する。そのため、通信システムには、オブジェクト群を符号化した不定長バケットを、効率よく、しかも確実に宛先に送達する機能が要求される。そこで、まず、L L Cタイプ3を活かすO S Iコネクションレス型基本標準を使った7層構成の『確認付きコネクションレス型プロフィール』を実現した。次に、これを利用して大量のデータの送達を保証する『高信頼バルクデータ転送プロトコル (R B T 3)』を実現し、さらに、宛先が移動した場合にも最終移動先に大量のデータを確実に送達することのできる『ヒンティング付き高信頼バルクデータ転送プロトコル (R B T H 3)』を実現した。これらのアーキテクチャ、実装、実測性能、および評価について述べている。

Implementation and Evaluation of a Network System Based on LLC Type 3
for O Z : Object-Oriented Open Distributed System

Michiharu TSUKAMOTO	Electrotechnical Laboratory
Nobuo YOSHIE	Sumitomo Electric Industries, Ltd.
Takashi KONDO	Sharp Corporation
Isao MIZUTANI	Sumitomo Electric Industries, Ltd.
Nobuaki TANAKA	Matsushita Electric Industrial Co, Ltd.

In O Z system, distributed application programs consist of objects connected with each other. Objects are copied or moved from the caller to the callee keeping the relationships among objects. Therefore, the communicating system must transfer bulk data reliably to the destination object even if the object is migrated to another host. This paper describes our solution to these problems. At first, we developed the reliable "Acknowledged Connectionless Profile" based on LLC type 3 and connectionless OSI standards. Then, we developed "RBT3: Reliable Bulk Data Transfer Protocol, Type 3" which transfer bulk data reliably using the acknowledged connectionless profile. Further, RBT3 was extended to "RBTH3: Reliable Bulk Data Transfer Protocol with Hinting Facility, Type 3" in order to transfer bulk data reliably to the migrated destinations.

ッサである。各エンドシステムはシステムバスで接続されたLANボードを介してメディアに接続される。

エンドシステムのソフトウェアは3レベルからなる。最上位はネットワークワークソフトウェア、最上位は分散型アプリケーション・ジョーシングプログラムを搭載するためのドメインと呼ぶ仮想計算機であり、その中間には仮想計算機が相互に通信しあう環境を実現する分散カーネルである。ドメインは、システムの種々の管理を行う管理系、システムとユーザに対して共通資源へのアクセスを提供するサーバ系、オブジェクトをロードして実行する実行系に分類される。

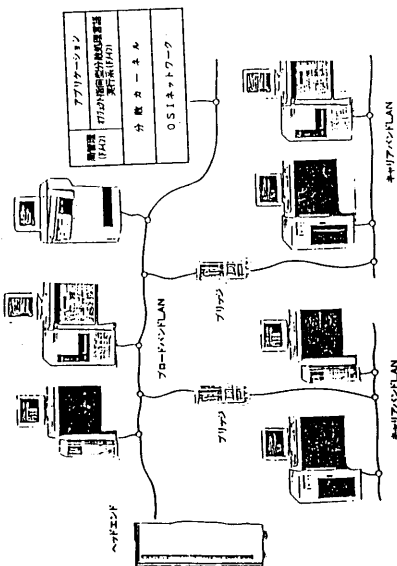


図1 OZ実験システムの構成

3. 通信システムが満たすべき条件

3.1 オブジェクトの表現

ドメインにおいては、すべての実体はオブジェクトという形式で表現される。オブジェクトがオブジェクトを引数としてオブジェクトに処理を依頼し、処理の結果をオブジェクトとして受け取る形で処理が進行する。処理効率を上げるために、オブジェクトはローカルオブジェクトとグローバルオブジェクトの2種類に分類されている。グローバルオブジェクトは分散される単位となるオブジェクトであり、分散システム全体で共有され、参照することができ、いっぽう、ローカルオブジェクトは必ずしもグローバルオブジェクトに所属し、所属するグローバルオブジェクトおよび同じグローバルオブジェクトに所属する他のローカルオブジェクト(仲間のオブジェクト)から参照されるだけである。オブジェクトは、オブジェクトの本体であるセル、およびオブジェクトの識別情報と参照関係を記憶するプロキシ、の2つの部分に分けられている。プロキシは、オブジェクトの種類を記憶する属性、グローバルオブジェクトの場合に自身を識別する識別子UID(Unique ID)、ドメイン内にセル部があるときはそのアドレスを記憶するがドメインにセル部があるときはそのドメインの識別子DID(Domain ID)を記憶するアドレス、の3つ

1. まえがき

アプリケーションが高度になるほど、複雑なデータ構造が必要になることが多い。しかし、通信によってしか情報交換できない分散システム上で、自由に高度なデータ構造を定義し操作することは簡単ではない。筆者らは、分散システム上で相互に関係を持ち通信しあうプログラム全体をデータ要素がさまざまな方法で連結された1つの分散型データ構造とみなし、分散処理とは隣接するデータ要素間においてデータ要素群を複写・移動することだと考えた。さらに、送信側で自由に定義・作成したデータ要素群を受信側が正しく解釈できるようにするには、データにそれを解釈する手続きを付随させたオブジェクトを基本のデータ要素とすべきであると考えた。このような発想のもとに、分散システム上のすべてのものをオブジェクトで表現し、オブジェクト群を複写・移動することにより処理を進める方式の『オブジェクト指向開放型分散システムOZ』の開発を進めてきた^{1,2)}。

OZでは、オブジェクトの数は非常に多く、しかも移動可能であるため、分散されたオブジェクト間で確実にオブジェクトを転送するには、コネクション型転送を採用するの難しい。そこで、コネクションを設定しないで転送することのできるコネクションレス型転送方式を採用することにした。しかし、通常のコネクションレス型転送では、送達が保証されず、また一度に転送できるデータの長さにも制限がある。そこで、OSIコネクションレス型基本標準だけを使って、パケット長には制限があるが送達が保証される確認つきコネクション型プロファイルを開発した。次に、これを活用して不定長パケットの送達を保証するコネクションレス型マルチデータ転送プロトコルRBT3(Reliable Bulk-data Transfer, Type 3)を開発した。さらに、宛先が移動した場合にも最終宛先への送達を保証するために、RBT3にヒンティングと呼ぶ機能を付加したRBT3H3(Reliable Bulk-data Transfer with Hinting Facility, Type3)を開発した。

コネクションレス型プロトコルとしては、開発がほぼ終了したOSIの基本標準群、LANを対象としてアカデミックな分野で研究開発されてきたBLAST³⁾やNETBL⁴⁾、最近ANSIで開発中のXTP⁵⁾などがある。RBT3はOSIの基本標準をベースにしてLANの特長を活かしたプラストプロトコルである。また、RBT3はこれまでにない新しい機能を提供するプロトコルである。

本稿では、まずOZの全貌と通信システムに要求される条件について述べる。続いて、通信アーキテクチャ、実装、および実測性能について述べ、最後に考察する。

2. システム構成の概要

OZ実験システムのネットワークは3セグメントで構成されるLANであり、各セグメントはブリッジで接続されている(図1)。幹線と支線の各セグメントにはISO8802-4⁶⁾に準拠するプロトコルバス方式とキャリアバウンドトーンバス方式の4つのを採用した。ISO8802-4を採用した理由として2点があげられる。第1は、OZは生産省の大規模プロジェクト『電子計算機相互運用データベースシステム』の一環として進められており、同大規模プロジェクトでは昨年度までの前期においては、LANはISO8802-4のプロトコルバス方式を使用するという合意があったためである。第2は、これらの2方式の組み合わせは、工場用LANとして注目されているMAPにも用されており、適用分野がきわめて広いためである。

現在、このLANに接続されているマシンは、国産ワークステーション(UNIXシステムV)2機種6台と電総研が開発したデュアルCPUボード数枚からなるマルチプロ

ワールドで構成され、オブジェクトへの参照はすべてこのプロクシを介した間接参照とな
っている。UIDは、そのIDを付したドメインのDIDとDIDで付した通番LIDか
ら構成され、UIDのドメインにLIDを問い合わせると詳細が知れる。また、アドレス
フィールドがDIDのときは、そのDIDに対応するドメインに目的のセルに至るプロク
シが存在する。プロクシを集めた表をオブジェクト表と呼び、オブジェクト表だけで、そ
のドメインに存在するオブジェクト、およびそのドメインから参照しているオブジェク
トの所在などが分かる(図2)。

3. 2 オブジェクトの転送

仲間のオブジェクト間で処理を依頼し結果を受け取るさいには、オブジェクトへの参照
関係(ポインタ)を渡す。これに対して、異なる仲間に属するオブジェクト間では、指定
されたオブジェクトから参照関係を次々にたどり、ローカルオブジェクトは実体を複写
し、グローバルオブジェクトは参照関係だけを複写する。異なるドメイン間でこのことが
必要になると、オブジェクトの相互参照関係を保存するだけでなく、送信側で実行が中断
されていたオブジェクトが受信側で再開できるようにレジスタやスタックの内容も修復で
きる形で転送しなければならぬ。このために、転送するオブジェクトの量は不定長とな
り、さらに、セル本体だけでなくプロクシも転送されるために、受信側ドメインのどのブ
ロクシもそれまで参照していなかったドメインに対して新たな参照が生じることがある。
参照されているプロクシに対して、新たに参照するプロクシが生れたことを通知するのは
効率が悪いし、参照される側でその情報を維持するのも困難なため、プロクシには参照に
必要な情報しか含めていない。

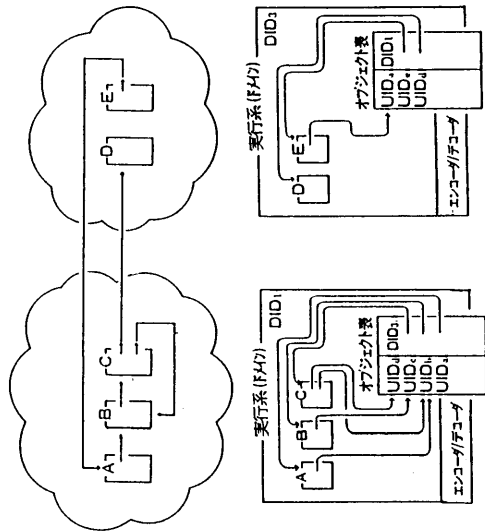


図2 オブジェクトの参照関係の表現

オブジェクトが移動するさいには、プロクシとセルを転送し、転送元ドメインにはアド
レスフィールドを移動先DIDに変えたプロクシだけを残す。移動時の処理は簡単である
が、オブジェクトを受信してもプロクシが存在するにもかかわらず、目的のオブジェク
トは移動してしまっているという状況がおきる。

通信システムに対しては、通信しあう対象の組み合わせの数が非常に多く、参照関係が
逐次変化する。しかも対象が移動するという案件のもとで、不定長のデータを効率よく確実
に宛先に送達する手段を提供することが要請される。

4. 通信アーキテクチャ

オブジェクト間における通信は、プロクシ間で行なわれるという見方とドメイン間で行
なわれるという見方の2通りが考えられる。いずれの場合も、通信相手が非常に多く、し
かも参照関係が逐次変化するため、通信対象間に固定的な通信路を設定する、いわゆるコ
ネクシオン型モデルを採用することはできない。そこで、コネクシオンを設定することな
く伝送できるコネクシオンレス型モデルを採用することにした。

4. 1 ネットワーク - 確認つきコネクシオンレス型プロファイル

通信アーキテクチャの実現にあたり、次の3つの基本方針をたてた。まず、オブジェク
トを転送するという考え方は密結合の考え方であるため、密結合が実際に実現できて有効
性を発揮するLANを対象とする。次に、独自のプロトコルを開発することを避け、OS
Iのコネクシオンレス型基本標準だけを使ってプロファイルを開発する。さらに、このプロ
ファイル単独でも一般的に使用できることを目指す。これらの方針にもとづき、コネクショ
ンレス型の欠点である送達を保証しないという問題を解決するLAN用のプロファイルを検
討することにした。

LANに関連するOSI下位2層のうち、物理層とMAC副層はINTAP(情報処理
相互運用技術協会)のLAN下位層実装規約^{*)}に準拠し、モトローラ製ETBC
(Enhanced Token Bus Controller)を使って実現された製品を使用することにした。その
ため、検討する必要があるものは、LLC副層からACSEまでである。設計したプロフ
イルを表1に示す。LLC3の送達確認機能を利用したプロファイルであるところから、確
認つきコネクシオンレス型プロファイルと名付けた。

表1 確認つきコネクシオンレス型プロファイル

層	基本標準	オプション選択
ACSE	CLACSE(DISI0035)	バージョンあり、実装情報なし
ルータ/コネクシオン	CLPPP(DIS9576)	バージョンあり、セレクタ固定なし
トラフィックネットワーク	CLSP(DIS9548)	チェックサムなし
ネットワーク	CLNP(IS06602)	インアクティブ
LLC	タイプ3(IS08802-2/Pdad2)	送達確認(reply, exchange使用付)
MAC	トークンバス(IS08802-4)	即時応答、非即時応答使い分け

但し、(1)各層の最大PDU長はMACの最大フレーム長で制限される。
(2)送信局が宛先局と同一セグメント上であれば、MACの即時応答機能の使用を推奨する。
(3)プリリッジはMACの優先度を保存する。

4. 2. 1 RBT3

コネクションレス型サービスを使ったバルクデータ転送プロトコルの多くは、受信側からの確認をとることなく一方的に連続送信し、一連の送信終了後に受信側からの通知によって正しく受信されなかったパケットを再送する、いわゆるブラスト方式と呼ばれるものがある。RBT3の設計においては、ブラスト方式による再送機能と送信側複数LSAPを用いたウィンドウを使ったフロー制御機能を基本とし、次の4種類のサービスプリミティブを提供する。

- (1) BT-DATA: ユーザデータを分割し、連続して転送するが、誤り回復やユーザの結果の通知はしない。
- (2) BT-DATA-ACK: ユーザデータを分割し、ウィンドウが許すかぎり、連続して転送する。誤り回復やユーザへの結果通知機能を持つ。エンド間での送達確認が行なわれるが誤り回復はしない。
- (3) BT-DATA-STATUS: BT-DATA-ACK 要求の結果を通知する。このサービス発行の有無は、BT-DATA-ACK のパラメータによる。
- (4) BT-ABORT: 転送の中断・放棄を通知する。

BT-DATA とBT-ABORTは付加的なものなので説明は省略し、目的の高信頼バルクデータ転送に関連するBT-DATA-ACKとBT-DATA-STATUS (図4-1) について述べる。確認つきコネクションレス型プロフィールでは、受信側LLC3からの送達確認を送信側が受信すれば、正しく転送されたものとみなすことができる。そこで、これをブラスト方式の送達確認として使用する。バルクデータ受信側LLC3から応答されたパケットの送達確認がすべて成功であっても、エンド間の送達は保証されていない。このために、バルクデータ受信側が最終パケットを受信すると、バルクデータ受信の成否がバルクデータ送信側に通知される。バルクデータ送信側RBT3は、これをユーザに指示する。いっぽう、バルクデータ受信側RBT3は、その通知の送達確認を正しく受信したのちに、バルクデータ受信をユーザに指示する。

RBT3の誤り回復は、送信側が主導権を持っており、受信側から誤り通知を積極的に行なうことはない。メディア上での伝送誤りにより宛先MACに届かない時は送信側MAC、宛先LLC3から送達確認が届かない場合や失敗 (宛先におけるバッファ不足など) の送達確認を受信したときは送信側LLC3で再送する。それでも回復できないときは、バルクデータ送信側RBT3ユーザにそのことを指示する (図4-2)。

RBT3で分割したパケットは、空きのLSAPがあるかぎり、ACSEを通じてLLC3に送信を要求できる。宛先LLC3から成功の送達確認を受信するとLSAPが空くので、次の送信を要求できる (図4-3)。

4. 2. 2 RBTH3

RBT3により、不定長のユーザデータを確実に転送することができる。しかし、ユーザデータを届ける宛先のオブジェクトが移動した場合には対応できない。移動先にフォワードイングする方法は、転送要求元と最終移動先間の送達を保証するのが難しく、オブジェクトが次々と移動した場合に間接参照の段数を減らすことができない。同報機能を使って移動先を放す方法も、すべてのドメインに確実に通知するのが困難である。そこで、信頼性の高い対向通信を保証するRBT3のBT-DATA-ACKを次のように拡張する (図5)。ユーザデータを分割したパケットのうち、最初のものだけをまず宛先に送信し、その応答を待つ。宛先では、目的のオブジェクトが存在すれば成功と応答し、移動してい

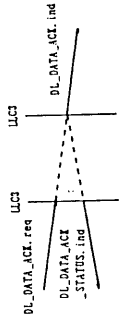


図3 LLC3送達確認

LLC3副層の基本標準には3種のタイプが存在する。LLC1はコネクションレス型で送達保証されない。LLC2はコネクション型であり、再送やフロー制御を行ない送達も保証されるが、コネクション数が小數に限られる。LLC3^{*)}はコネクションレス型でありながら送達保証される。LLC1を使用するとACSEまででは信頼性を保証する手立がないため、OZではLLC3を使用することにした。LLC3には、送達確認機能 (図3) 以外に、相手LSAPからのリードとLSAP間でのデータ交換の機能があるが、リードと交換の機能はネットワーク層以上でこれを使用する基本標準がないので採用せず、送達確認機能だけを採用することとした。

ネットワーク層のコネクション型プロトコルCLNPは、分割・組立機能、経路選択・指定機能、などを持つが送達を保証しないので採用せず、標準で許されている空機能 (inactive) で使用することにした。トランスポート層のコネクションレス型プロトコルCLTPとセッション層のコネクション型プロトコルCLSPは、アドレスセレクト機能 (selective) をもち、実質的には機能はないので、そのまま採用することにした。プレセクション層のコネクション型プロトコルCLPPはアドレスセレクトと符号化・復号化機能をもち、RBT3で使用する際には複数のパケットにわたって1つの符号化・復号化を行なうことになり、符号化と復号化が難しくなるので、OZではオクテット列だけのOZマシネクストを設けることとした。ACSEは送信の相手を指定するための種々のパラメータをもっており、OZでは応用エンティティ識別子でRBTプロトコルマシネクストを識別し、応用プロセスマシネクストでドメインを識別することにした。

LLC3をより効率よく活用するために、次のような工夫を行なった。LLC3の下位に位置するISO802-4MAC副層の基本標準には、受信側がトークンを保持したままでも即時に応答する機能 (即時応答機能) がある。これは、同一セグメント上の局間でしか使えないが、きわめて効率がよい。そこで、送信宛先局が同一セグメントにあれば即時応答機能を使用し、宛先局が同一セグメント上にならないときは即時応答機能を使わないこととした。この判別を効率よく行なうために、MACアドレスのSSI (Subnet Segment ID) をセグメントのIDに対応させることとした。LSAP以上のアドレスに関しては、Lセレクタは最大16個、CLNPは空なのでなし、Tセレクタは1つに固定、Sセレクタでアプリケーションを振り分け、Pセレクタは1つに固定した。1つのLSAPでは、応答が返ってくるまで、次のフレームが送信できないため、非即時応答では効率が悪くなる。そこで、複数のLSAPを使って同時に送信できるようにした。

CLNPからACSEまでの基本標準には送達確認を渡すサービスはないが、これはプロトコル上には表われないので、実装事項として導入した。CLNPを空機能で使用することにしたため、異なるセグメント間では通信できない。これを解決するために、MACブリッジを導入することにした。

4. 2 分散カーネル - ヒンティングつきバルクデータ転送

上記プロフィールを使って希望のサービスのサービスを提供するプロトコルを次の2段階に分けて検討する。まず、MACフレーム長によって制限されない不定長パケットが効率よく確実に転送できる方式RBT3を検討し、続いて、宛先が移動しても不定長パケットが移動先に確実に転送できるように拡張した方式RBTH3を検討する。

5. 実装

5. 1 LANボード

コネクション型OSIプロトコルのトランスポート層までを実装する目的で開発されたLANボードを利用し、そのプログラムを入れ替えることによって目的のものを実現した。その構成を図6に示す。ボード内には、ROM、ローカルメモリ、デュアルポートメモリの3種類のメモリが設けられており、MC68000によって制御される。ROMには基本的なボードの制御プログラムが格納される。ローカルメモリはETBCが送受信専用で使用するものであり、制御CPUがデュアルポートメモリとの間のコピーを行なう。デュアルポートメモリは、ホスト計算機のバスアドレスにマッピングされており、ホスト計算機と制御CPU間でデータの授受に使用される。このボードでは、LLC3までの処理を行なっている。ホストから制御レジスタを介して送信が指示されると制御CPUはデュアルポートメモリ上の256バイト単位で連結されたバッファからローカルメモリの連続領域にコピーし、LLC3送信コンポーネント(プロトコルマシジン)処理を行なって、ETBCに送信を指示する。制御CPUがローカルメモリに受信バッファを設けてETBCに受信を指示すると受信可能になる。受信時にはLLC3受信コンポーネント処理がETBCで行なわれたのち、制御CPUに受信が指示される。制御CPUは256バイト単位のバッファをデュアルポートメモリから獲得して分割コピーを行ない、そのち、ホストに受信割り込みを通知する。デュアルポートメモリに空きがないと、ローカルメモリからコピーできない。この場合には、ローカルメモリがいずれ枯渇するのでETBCの受信エラーが発生し、送信元へ受信エラー通知が返される。このようにして資源不足による背圧が伝播する。

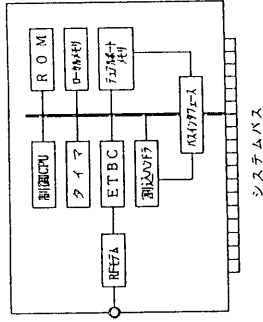


図6 LANボードの構成

5. 2 ホスト

UNIXマシンの場合には、CLNPからRBTH3までをUNIXドライバとして実装した。その構成を図7に示す。ドメインとRBTH3間におけるプロクシの登録・変更用アプリケーションは設けず、共有メモリ上にオブジェクト表の一部を設け、UNIXプロセスとして実現されるドメインと共有させた。

ユーザは、ユーザ空間内の送信バッファ列にデータを入れて、ioctlを用いてRBTH3を呼ぶ。RBTH3では、送信プロトコルマシジンを作ったのち、ユーザ空間内の送信バッファを共有メモリにコピーし、ACSEが送信可能な最大パケット長単位に分割し、ヒンティング処理を行なったのち、残りのパケットの送信を要求する。ACSEではデュアルポートメモリ内の256バイト単位のバッファ列を獲得してパケットを作り、空きウィンドウを探してLANボードに送信を要求する。現在の実装では簡略化のために、UI

ば移動先情報をつけて失敗と応答する。失敗応答を受けるとユーザに移動先を指示するとともに、RBTH3が管理するプロクシを移動先に変更して移動先に再送する。通知成功応答を受けたさいの2番目以降のバケットの転送はBT-DATA-ACKと同様である。各ユーザは内部に所持するオブジェクトのプロクシをあらかじめRBTH3に知らせておかなければならない。宛先から移動先のヒントをもらうことから、このメカニズムをヒンティングと名付け、この機能を付加したRBTH3をRBTH3と呼ぶ。

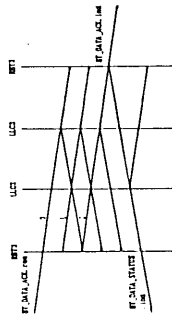


図4-1 RBTH3の動作タイミング (正常)

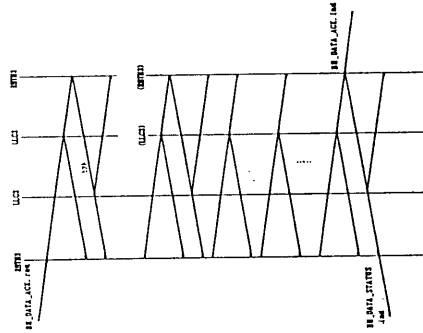


図4-2 RBTH3の動作タイミング (誤り回復)

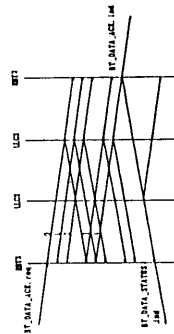


図4-3 RBTH3の動作タイミング (70-制御)

図5 RBTH3の動作タイミング

ンドウに空きがない場合はスリープして、ウインドウが空くのを待つ。
ユーザは宛先の指示に DID と UID を用いる。指示された DID がドメイン表に存在しない場合には、RBTH3 からステーション管理にアドレス獲得を指示し、得られたアドレスをドメイン表に登録して送信処理を再開する。

LANポートから受信割り込みが発生するとデモンが起動され、共有メモリ中に 256 バイト単位のバッファを獲得して、これにデュアルポートメモリからコピーのち、プロトコルを解析して RBTH3 に渡す。RBTH3 では、それが最初のパケットならば受信側プロトコルマシンの作成したあと、オブジェクト表を参照して受信処理を連結する。一方、宛先オブジェクトがドメインに存在すれば、後続のパケットを受信して連結する。連理のパケットを受信し、組み立てると宛先ドメインに受信を指示する。受信を指示されたドメインのプロセスはドライバを呼びだし、共有メモリをユーザメモリにコピーする。

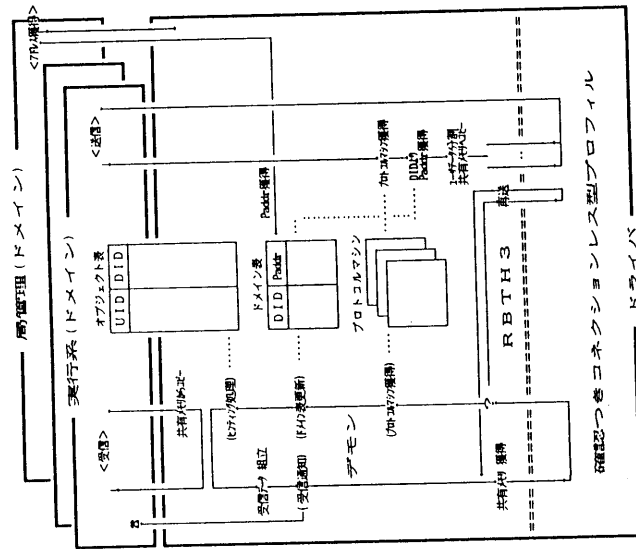


図7 ホスト計算機上のRBTH3ソフトウェア構成

5.3 ブリッジ

ブリッジは、5.1 で述べた LAN ポード複数枚を UNIX マシンにバスで接続し、UNIX マシン上のソフトウェアによって実現されている。その構成を図8に示す。ブリッジ用の LAN ポードでは、MAC 副層の送受信処理しか行わない。中継するパ

ケットはデュアルポートメモリに記憶される。UNIX マシン上のソフトウェアは、LAN ポードとインターフェースをとるためのドライバ、パケットの中継を行なう中継処理部、ブリッジを管理するブリッジ管理部、およびオペレータとインターフェースをとるブリッジモニター部から構成されており、ブリッジの管理情報は共有メモリに記憶されている。

ブリッジでは、可能な限り無用なパケットを中継しないように配慮しなければならぬ。そこで、次の2種類の手法を採用した。

(1) ETBC のフィルタリング：OZ のアドレス体系は SSI を利用している。そこで、各 ETBC のアドレスフィルタには、それが接続されているセグメントの SSI を設定し、その SSI を持つパケットは受信しないようにしている。

(2) 中継処理部の学習・フィルタリング：中継処理部の受信プロセッサは、ドライバから、受信した LAN ポードの番号と受信したパケットの PCI を受取ると、LAN ポードの番号に対応するポート番号、宛先 SSI と宛先 SSI を求める。次に、そのポート側に宛先 SSI のセグメントが存在することをフィルタリングテーブルに記憶し(学習)、さらに、宛先 SSI がどのポートの先に存在するかをフィルタリングテーブルから求める(フィルタリング)。そのうち、求めたポートに対応する送信プロセッサに送信を指示する。対応するポートがなければすべての送信プロセッサに送信を指示する。LAN の構成を変更したさいには、以前の学習結果を消さなければならぬ。そのため、フィルタリングテーブルの各エントリについてタイムアウトしており、一定時間以上アクセスされないとエントリは削除される。

このように、SSI をフィルタリングに用いることにより、必要なパケットだけをブリッジにとりこみ、さらにフィルタリングテーブルのエントリ数を少量にすることにより、無用な負荷がブリッジにかかるのを避けている。さらに、中継処理の効率を上げるため、受信したパケットは LAN ポードのデュアルポートにそのまま残し、送信プロセッサからドライバに送信が指示されたときに、ドライバ内で受信側バッファから送信側バッファへコピーしている。

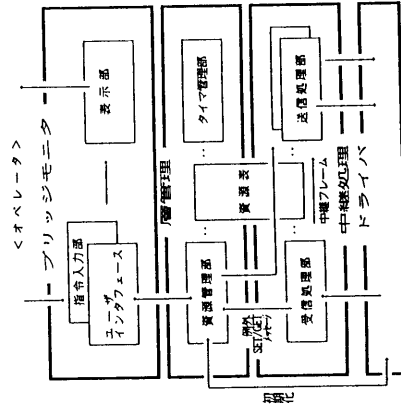


図8 ブリッジの構成

6. 性能測定実験

OZの基本部分の開発が完了したので、それを改良したOZ+の開発に着手することにした。その設計に役立てるために、OZの通信に関する部分の性能測定実験を行なった。その実験性能を図9にまとめ示す。

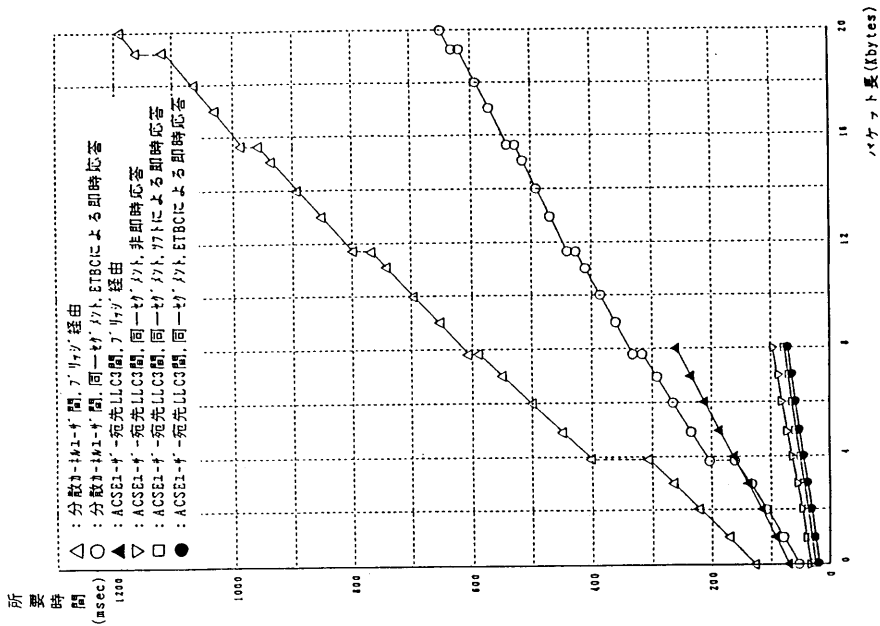


図9 性能測定値

6. 1 確認つきコネクショントラック型プロファイル

LLC3送達確認の実現法には、即時応答を使用しない方法、ソフトウェアで即時応答を実現する方法、チップ(ETBC)によって即時応答を実現する方法、の3通りがある。図9における▽、□、および●は、UNIXユーザープロセスからドライバ内に実装されたACSEに対してデータを渡し、宛先LLC3からの送達確認がユーザープロセスに戻るまでの時間を測ったものである。即時応答の場合には、送信側がトークンを保持している時間内に送達確認を返す必要がある。ETBCでは送達確認応答がハードウェアで行なわれるためにスロットタイムを短くすることができるが、ソフトウェアで実現する場合は短くできない。この実験では、ETBCの場合は、X'40'オクテットタイム(X'40'オクテットの送信時間)、ソフトウェアの場合はほぼ近いX'120'オクテットタイムとした。

ETBC即時応答では、パケット受信から送達確認送信までの時間の実測値は数10μsec以下であるため、0バイトのユーザーデータ(PCI長は108バイト)の場合の立ち上がり時間16msecが、送信側におけるプロトコル処理とモジュール間インタラクションに要したオーバヘッドといえる。ソフトウェアで即時応答を実現する場合のオーバヘッドはさらに4msec大きく、即時応答機能を使わない場合のオーバヘッドはまたさらに5msec大きい。ETBC即時応答の場合の傾き8μsec/バイトから、メディア上の所要時間約1μsec/バイトを減じた7μsec/バイトが、送信側でユーザー空間からデュアルポートメモリを経てローカルメモリまでの間のバッファ確保とコピーに要する1バイトあたりの時間となる。即時応答機能を使わない場合の傾きがそれ以外に比べて大きいのは、受信側LANボードにおいてローカルメモリからデュアルポートメモリにコピーしたのちに送達確認を戻しているためである。別の見方をすれば、この傾きの差2μsec/バイトがローカルメモリからデュアルポートメモリへのバッファ獲得・コピーに要する時間といえる。

6. 2 RBTH3

図9における○は、分散カーネルユーザー(すなわちドメイン)からドライバ内のRBTH3に渡されたデータが即時応答付きLLC3を使用したACSEによって宛先のドメインに渡され、エコーバックされて送信側ドメインに戻されるまでの時間を半分にしたものである。

RBTH3の立ち上がりは50msecである。MACの最大フレーム長は約4Kバイトであり、4KBまでの傾き35μsec/バイトが、受信側側におけるバッファ獲得とコピーの性能である。ACSEで転送されたデータを宛先側でコピーする速度が送信側のコピーと同じと考え、35μsec/バイトから、両側におけるACSEからETBCまでの速度2×7μsec/バイトとメディア上の速度1μsec/バイトを引いた、20μsec/バイトが両側RBTH3におけるバッファ獲得とコピーのおおよその性能と言えよう。

送信データが4Kバイトを越えると、RBTH3によって分割して、複数のLAPをウインドウとして使って転送し、宛先RBTH3で組立が行なわれる。ちょうど4Kバイトごとに段階的に性能が低下するが、4パケット以上になるとウインドウ効果によって段階的な性能低下は認められなくなる。

6. 3 ブリッジ

ブリッジを介して、6.1と6.2に対応する性能を測定したものが、図9における▲

8. むすび

LLC3をベースにした確認つきコネクショナル型プロファイル、これを有効に適用する高信頼性バルクデータ転送プロトコルRBT3、および宛先が移動した場合でも宛先に送達することを保証する機能をさらに付加したRBTH3を設計して実現した。未だ十分な性能を達成してはいないが、いくつかの分野では方式的な正しさが実証できたといえる。ブリッジを介したバルクデータ転送時に見られた性能上の問題に関しては、早急に検討を行いたい。OZを改良したOZ+では、その検討をもとに、よりよいプロトコルとするよう改良を加えたい。

なお、この研究は、通産省の大型プロジェクト『電子計算機相互運用データベースシステムの研究開発』の一環として進められているものである。その機会を与えられ、しかもいろいろ助言をいただいた電子技術総合研究所情報アーキテクチャ部長榎上昭男博士に感謝致します。また、住友電気工業、松下電器産業、シャープの関係各位のご助力にも感謝致します。さらに、筆者らの要望に応じて、LANボードを再開発していただいた住友電気工業の方々、およびブリッジを開発していただいた松下電気産業の方々にも感謝致します。

参考文献

- 1) M. Tsukamoto et al.: The Architecture of Object-Oriented Open Distributed System: OZ. Interoperable Information Systems ISIS '88. Ohmsha. pp153-166 (1988.11)
- 2) 塚本他: OZ: 対外向指向開放型分散型ネットワーク - ネット指向型分散型ネットワーク言語とゆ美装、情報学会研究会報告、7093/7言語21-4 (1989.6)
- 3) D. Theriault: BLAST, an Experimental File Transfer Protocol. MIT-LCS Computer System Research Group RFC-217 (1982.3)
- 4) D. Clark et al.: NETBLT: A High Throughput Transport Protocol. Proc of ACM SIGCOMM '87. pp353-359 (1987.8)
- 5) M. Cohn: A Lightweight Transfer Protocol for the U.S. Navy SAFENET Local Area Network Standard. Proc of 13th Conf on Local Computer Networks. pp151-156 (1988.10)
- 6) ISO8802-4: Information Processing System - Local Area Networks - Part4: Token-passing Bus Access Method and Physical Layer Specification
- 7) ISO7498-1/Add.1: Information Processing System - Basic ReferenceModel - AD1: Addendum Covering Connectionless-mode Transmission
- 8) 情報処理相互運用技術協会: LAN下位層互換規約案 S012(V1.0) (1989.3)
- 9) ISO8802-2/Pdad2: Addendum to ISO/DIS 8802-2 Logical Link Control - Acknowledged Connectionless-mode Service, Type 3 Operation

と△である。

確認つきコネクショナル型プロファイルにおける立ち上がりは、即時応答つきLLC3を使った場合のほぼ4倍強の立ち上がりとなっている。4回の転送が行なわれることを考慮するとこの値は妥当なものといえる。4Kバイト以下において、ブリッジを介した▲の傾きから、それに対応する同一セグメント上における●の傾きを引いたものがブリッジにおけるバッファ獲得とコピーの速度となる。この値は、20μsec/バイトである。いっぽう、RBTH3に関する△と○の4Kバイトまでの傾きの差は19μsec/バイトであり、上記の値とほぼ同値となる。

ウィンドウ効果が有効に働いていないのは、連続してPDUを送信するさいには、PDU中継処理中に送信側から次々とPDUを受信し、さらに送達確認PDUも受信するため、ブリッジで輻射が生じていることが原因だと思われる。

7. 考察

7. 1 性能

少量データ転送時に十分な応答速度が得られるように改善しなければならぬ。LANボードとホストドライバ間のハンドシェイクの改善、ドメイン、UNIX共有メモリ、LANボードデュアルポートメモリ、LANボードローカルメモリというメモリの種類の減少、ドメイン (UNIXプロセス) とドライバ間における記憶空間の問題とインタラクショナルの改善、などを行なう必要がある。

大量データ転送時に十分な性能が得られていない点も改善しなければならぬ。この解決方法は上記問題点の解決方法と関連するが、メモリの種類の減少やDMAの採用などが必要である。現在の実装で容易に改善できる点としては、ドメインとドライバ間でデータ交換に使用するバッファの単位を大きくすることが上げられる (現在のところバッファ単位は256バイトであり、上記実験に使用したマシンでは、ドメインからドライバにコピーインし再びコピーアウトする性能は、25μsec/バイトと38μsec/バイトである)。

7. 2 方式

確認つきコネクショナル型プロファイルは正しく動作することが実証された。MAP (Manufacturing Automation Protocol) では即時応答機能しか採用しないために複数のセグメントにまたがった通信ができないが、宛先によって即時応答機能を使用するか否かを選択するというかわれわれの方法はこの問題を解決している。RBTH3のヒンチン機能を有効にするためのものであり、コネクショナル型転送方式を活用する上で有効であると確信している。

RBTH3は、LLC3の即時応答機能を使えば、予想通りの振る舞い (性能は十分ではないが) をしていることが確かめられた。RBT3はその一部であるため、RBT3も正しく動作し、予想通りの振る舞いをするものと思われる。しかし、ブリッジを介したバルクデータ転送の場合には、輻射が起きていると思われる。この原因としては、まずブリッジの性能が十分でないことがあげられるが、LLC3送達確認がブリッジを輻射させる要因となっている可能性がある。得られた特性を使って数理的な手法やシミュレーション手法を使って、RBT3およびRBTH3とブリッジとの関連性について、究明する必要がある。