

# コンピュータ・ネットワークにおける NCP実現上の問題点

鍛治 勝三

(財) 日本情報処理開発センター

## 1. はじめに

1960年後半にコンピュータと遠隔地にある端末装置を結ぶコミュニケーション・システム概念が出現して以来、数多くのタイム・シェアリング・システムやオンライン・システムが開発された。

1968年に開発を開始したARPAネットワークが成功したのをきっかけに各国にかなり多くのコンピュータ・ネットワークが出現した。これらのコンピュータ・ネットワークの多くは、データ伝送と交換能力を持つサブネットを採用したが、急激に発達したため、その国の電電公社が構築するデータ交換システムがまに含めなかった。その結果として、これらのネットワークの設計者自身が自らの交換システムを構築している。

商用のバケット交換ネットワークとしては、1971年にスペインの電電公社CTNE (COMPANIA TELEFONICA NACIONAL DE ESPAÑA) が最初のサービスを開始した。Trans-Canada Telephone SystemがCCITTの勧告に準拠した方式のDatapacシステムを1976年にサービスを開始することを発表している。

ARPAネットワークが構築されてから6年経った。この間にバケット交換ネットワークに関してはかなり研究も進み、実用上問題がなくなってきた。しかし、ネットワークに接続するためにHostに存在するNetwork Control Program (NCP) に関しては、まだ初期段階を出ていない。

以下では、Hostとサブネット間結合とNCPの検討及び測定について述べる。ここでの調査対象のコンピュータ・ネットワークあるいはバケット交換ネットワークとしては、表1に示した代表的なコンピュータ・ネットワークあるいはその一部であるバケット交換ネットワークの7つである。

## 2. HostとPSPとの接続方法

まず最初に、Hostを既存のPacket Switching Processor (PSP) にどのようにして接続するかを検討する。接続媒体としては、実験用ネットワークではPSPがHostと同一構内にあるのでチャンネル間結合を多く用いられ、一方商用ネットワークではPSPが遠く離れた電話局内に設置されているので回線結合が用いられる傾向が強い。

ここではチャンネル間結合について考えてみる。Hostの標準システムとしてチャンネル間結合がサポートされていれば問題はあまりないが、サポートされていない場合には特別なI/Oルーチンを新たに作成しなければならぬ。I/Oルーチンのプログラミングは非常に複雑である。チャンネル間結合アダプタ (CCA, Channel-to-Channel Adaptor) が標準品でない (特に他メーカー) 場合、障害時の障害原因の切分けなどの保守作業が困難である (図1 (1) 参照)。

一方、Hostを直接PSPに接続するのではなく、間にミニコンピュータ (N

表1 HOST/パケット交換プロセッサ(PSP)。

ネットワーク名 比較項目	JIPNET	ARPANET	
		Local/Distant Host	Very Distant Host
稼働開始年	1974	1969	
地域	JIPDEC構内	米国全土およびヨ-	ロッパの一部
開発主体	JIPDEC	米国国防省	ARPA
目的	JIPDECの3Hostを結ぶ実験的研究用ネットワーク	全米の主要大学, 研究ネットワーク, ハワ張	実験室を結ぶ研究用ネットワーク
パケット交換プロセッサ(PSP)	NEAC 3200/50	DDP 316/5 新プロセッサ開発中	16
ノード数	3	58 (内TIP 24台)	
Host数	3	88	
PSP間伝送速度(kbit/sec)	48	50	
接続媒体	ケーブル	ケーブル	回線
通信方式	半2重(コンテンション方式)	全2重	全2重
伝送速度(Kbit/sec)	1,000	100	1.5~230.4(通常:9.6)
伝送の最小単位(bit)	8	1(IMP padding convention:16)	16
伝送方式	8ビット並列	ビット直列	ビット直列
同期方式	ブロック	非同期(bit-by-bit)	同期
誤り制御	無	無	CRC (24bit) Positive ACK
データリンク確立コマンド	Unlock Program Switch	Ready Indicator On	Line Passive State
データリンク切断コマンド	Lock Program Switch	Ready Indicator Off	-
伝送制御手順	-	-	Particular Binary Synchronous Communication
最大メッセージ長(bit)	9,216	8,095*	8,096(1,008分割)
ヘッダ長(bit)	32	72*	72
備考		*1: 先頭の32ビットはleaderであり、先頭の72ビットはHost-to-Hostメッセージのヘッダである。	HostとIMPはパケット単位で送受を行なう。第1パケットは16ビットのControl Word×32ビットのleaderから成る。第2パケット以降は、先頭に16ビットのControl Wordがあり、次に最大1008ビットのセグメントから成る。メッセージは最大9パケットに分割されて授受される。

PSP: Packet Switching Processor

インターネット比較一覧表

CYCLADES (CIGALE)	CTNE	DATAPAC	
		HDLC	BSC
1974	1971	1976 (予定)	
フランス	スイス	カナダ	
IRIA	CTNE(電電公社)	Trans-Canada Telephone System	
フランス政府, 大学 研究機関等を結ぶ実 験的研究用ネット	商用パケット交換ネ ットワーフ	商用パケット交換ネ ットワーフ	
MITRA-15	UNIVAC418 III, IBM 3968	-	
12 (含 Concentrator 6台)	2	-	
20	8		
48, 19.2, 9.6, 4.8	4.8	-	
回線	回線	回線	回線
全2重	全2重, 半2重	全2重	半2重
4.8, 4.8	4.8	-	-
8	8	-	8
ビット直列	ビット直列	ビット直列	ビット直列
同期	同期	同期	同期
CRC (16bit)	CRC (16bit) 又は LRC+VRC Alternate ACK	CRC (16bit) Control Field (8bit)	CRC (16bit) Alternate ACK
-	Open Command	Request Connect	ENQ bid (Contention basis)
-	Close Command	Request Disc.	EOT
Transparent Binary Synchronous Mode Procedure, CCITT V24	Transparent Procedure or Non-transparent Procedure	HDLC (High- level Data Link Control)	BSC (Binary Synchronous Communication)
2040	2040	2000*1, 2040*2	2000*1, 2040*2
96	88	80*1, 40*2	80*1, 40*2
		*1: Virtual Callの場合 *2: Datagramの場合	*1: Virtual Callの場合 *2: Datagramの場合

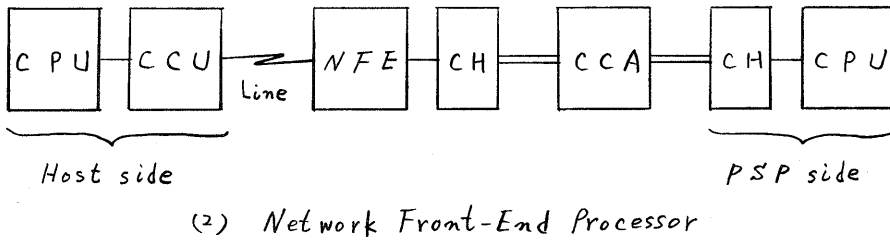
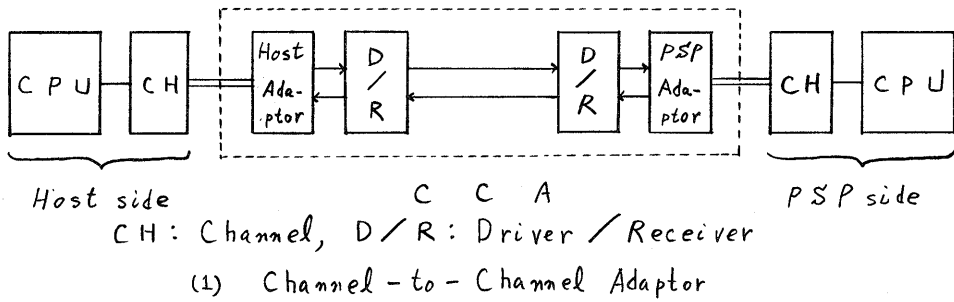


Fig.1 The Interconnection of a Host and a PSP

FE (Network Front-End processor) を介して接続する方法がある。既存のコンピュータ・システムでは既に通信制御装置 (CCU, Communication Control Unit) が設置され、その I/O ルーチンもサポートされているので、CCU-NFE 間は回線接続、そして NFE-PSP 間は CCA 接続となる (図 1 (2) 参照)。NFE として 5000 ドル位の非常に安いマイクロ・コンピュータを使えば、CCA サポートの特別な I/O ルーチンを作るよりも安くなり、ソフトウェアの保守も不必要となり、信頼性もあがり、同時にハードウェア保守も簡単になる。

また CCA 接続で注意しなければならないことは、ネットワーク・プロセス間で授受されるメッセージが数十バイト以上でも、Host-PSP 間ではメッセージを分割してセグメント単位に取り扱う場合には最終セグメントは 1 バイトになる可能性があり、Host 側のチャネルも I/O スーパーバイザは 1 バイトでも転送できなければならない。もしできないならば CCA による padding や Host-PSP プロトコル及び Host-Host プロトコルでの対策を考慮されなければならない。

### 3. Host と PSP との接続事例

JIPNET の IMP<sup>\*1</sup> と FACOM 230-75 Host は CCA で接続されている。CCA の仕様を以下に示す。

通信方式: 半二重; 転送方式: 8 ビット並列・ブロック転送; 転送速度: 1 Mビット/秒; 転送制御用コマンド: Write Data with Interrupt (WDI), Write Data (WD), Read Data (RD), Sense, No Operation (NOP), Lock Program Switch (LK), Unlock Program Switch (ULK), その他

\*1: Interface Message Processor の略で、PSP に相当する。

CCAは半二重であるが、両系に対し特に制御局/従属局という関係をもたず、コンテンション方式を採用している。コンテンション方式であるので、双方が同時にWDI/WDコマンドを発信した場合、コマンドの衝突が起こる。この場合JIPNETのHost-IMPプロトコルでは、IMP側が送信権を持ち、再度WDIコマンドを発信し、一方Host側は受信状態に入り、IMP側よりのAttentionを受け付けてからRDコマンドを発信し、実際のデータ転送が行なわれる。ここで問題になるのは、IMP側のリトライが非常に早く行なわれるので、NCPがAttention受付け用のECB(Event Control Block)を再度登録する<sup>\*1</sup>のがまにあわず、IMP側のAttention到達のみはI/Oスーパーバイザによって検出されてしまい、NCPにはそのAttentionがpostされない<sup>1</sup>ので両系はデッドロックの状態に入ってしまう。IMP側のAttentionとNCP側のWDIコマンドがモニタ内でずれ違えば、すでにAttentionがpostされている状態でコマンドの不整合となり、その時点でpostされていたAttentionと次のIMPからのリトライによるAttentionとが続けてくることになり、Attentionが1つ余分に入って来たことになってしまう。コマンド不整合後のこれらの対策としては、IMP側のリトライを数ミリ秒遅らせるようにしている。コマンドの衝突はモニタ内部、チャネルーアダプタ内、アダプタ内部、アダプターIMP間などの各所で起こるが、JIPNETの第一バージョンではコマンドの衝突時のステータスが十分にチェックできていなかった。またコマンド衝突の場合の大部分において、I/Oスーパーバイザはハード障害とみなし、エラー・ロギングを収集しており、衝突が頻繁に起こるとエラー原因をディスクへ書き出すのが向に合わなくなり、コア・バンパーが無くなり、システムをダウンさせていた。これらの障害に対しては、第二バージョンでは対策がとられている。

#### 4. NCPの常駐場所

NCPをHost内に常駐せると既存OSとのインタフェースが非常に複雑になるので、ほとんどのすべてのネットワークに肉連するホスト・ソフトウェアをHostコンピュータそれ自身よりもむしろNFEに常駐させた方がよい。またNCPをHostに常駐させることによりHost側に約100Kバイト位のメモリ容量と処理時間のオーバーヘッドが生じる。JIPNETの測定によるとNCPがフル稼働している場合はCPUバウンドの状態になり、NCPのCPU優先順位が高いのでネットワークに肉連のないジョブに対してCPUリソースがディスパッチされない。

NCPをNFEに常駐させ、ネットワークに加入しているすべてのHostが同じ(標準)NFEを採用すれば、ネットワーク・ソフトウェアは一通になり、開発コストが安くなる。

Hostに常駐しているNCPがユーザ・プログラムではなく制御プログラムとして作成されている場合は、NCPの障害が即OSの障害になることが多く、作成、デバッグ、そして保守・改良が困難である。一方NCPがNFEに常駐していれば作成も容易で、保守や改良もHostシステムのローカル運用に対し何ら影

\*1: Attention受付け用ECBをSET ECBマクロ命令により登録し、Attentionが入るとそのECBにpostされる。PostされたECBは再度登録しないと次のAttentionが受け取れない。

響を与えずに行うことができる。

ネットワーク・ソフトウェアが一通りにになると、ネットワーク・プロトコルの改定もずつとやりやすくなる。如入している Host がたくさんになり、NCP がすべての Host 内に常駐している状態では、ネットワーク・プロトコルの改定は不可能となり、複雑で不効率なプロトコルに甘んじなければならなくなってくる。

NCP を NFE に常駐させた場合は、新たに Host - NFE プロトコルが必要となってくるが、これは Host - PSP プロトコルや Host - Host プロトコルに比べてずつと簡単である。しかしながらプロトコルの重複を増加させる欠点が生じる。

現在 JIPNET で作成されている NCP はすべて Host 内に常駐している。したがって、ここでは検討をより具体的にやるため、NCP は Host 内に常駐していると仮定する。

### 5. 既存 OS の検討事項

NCP の作成においては、特に既存 OS とのインタフェース機能を中心に検討し、NCP の肥大化による Host の過度の負担を防ぎ、かつネットワークに因り係せぬローカル・ジョブへの影響を最少限におさえなければならぬ。

NCP も 1 つのプロセスであるので、既存 OS に各々独立したジョブステップの向での通信機能が備わっているが否が、また備わっているとしたり、それはユーザプロセスで使うことができるか否かによって、NCP の位置づけは、次の通りが考えられる。

- ① 制御プログラムの一部とする
- ② ユーザ・プログラムとする
- ③ ①と②の中向的なサブモジュールとする

既存 OS にジョブステップ間通信機能が無い場合には、新たにこの機能をサポートするプログラムを作成しなければならぬ。ロール・イン/ロール・アウト機能やバッチ・メモリを採用した OS では、受信側ジョブステップが実主記憶上に存在した場合があるので、その処理が非常に複雑であり、この機能をユーザが独自で開発するのは不可能に近いと考えられる。

これらの内、①はある意味では理想であるが、OS の設計時に NCP の組み込みが考慮されていない場合はあまり良いものとはできなると考えられる。ネットワーク結合のために増加した制御プログラムが、その Host のみによるローカル・ジョブのオーバヘッドの増加をもたらし可能性がある。また既存 OS 部分に改造を加えねばならぬ可能性も強く、いわゆる標準 OS ではなくなるため、その後のメーカーの OS メンテナンス・サポートが期待できぬ恐れがある。

これに対し、②のユーザ・プログラムとして作成する形はユーザとしてはもっとも作成が容易であるが、独立したプロセスの起動・停止や制御・監視ができない可能性があり、ネットワーク・サービスに関連する既存 OS 内の制御表を参照するなど、かなり無理が生じる可能性が強い。

以上のような点を総合すると、結果的には③のサブモジュール形式で NCP を作成することも望ましい。サブモジュール形式を採用した OS では、あるサブモジュールの障害が別のサブモジュールへ影響を与えず、完全に自立している。

しかしこれら3つの選択は既存 OS の機能および性格によって決定づけられ

てしまい、選択の余地がない可能性がある。今後新たに開発されるOSは他のコンピュータとのコミュニケーションを前提として設計されなければならない。

ユーザ・プロセスは直接 Host-Host プロトコルのコマンドを使うのではなく、マクロなコマンド (NCP サービス・マクロと呼ばれる) を通じて NCP にサービスを依頼するのが普通である。NCP サービス・マクロを既存OSのマクロ命令と同じようにスーパーバイザ・コール (SVC) で実現できることが望ましい。そのためには既存OSの設計時に SVC コードの何番から何番までをユーザに解放するのを完全に決めておき、SVC ルーチン作成の詳細マニュアルを準備しておかなければならない。OS/360<sup>[1]</sup> では SVC コード 255 ~ 200 まではユーザのために予約されている。

ネットワーク・プロセスがコネクションを切断せずに異常終了した場合、既存OSのターミネータがそのことをNCPに通知してくれば、このことをリモートNCPに通知し、コネクションを自動切断することが出来る。ターミネータが通知してくれず、シガモタイム・アウト処理がなされていない場合にはリモートプロセスは永久に待たされることになる。

以上をまとめると、最低限以下の機能が既存OSに備わっていることが望ましい。

- ① NCP とプロセスとの間で非同期に情報の受け渡しが可能なること
- ② ビット・ストリング・データのトランスペアレントな送受が可能なること
- ③ NCP がプロセスの発生、監視および強制終了をおこなえること
- ④ プロセスの実行に必要なファイルを動的に割付け・返却が可能なること
- ⑤ 通信制御プログラム中の端末制御プログラムが伝送制御手順毎に完全に独立したモジュールであること
- ⑥ アカウンティングのためにリソース使用情報をNCPが取り出せること
- ⑦ ローカル・リソース全体の使用状況がNCPから参照できること

## 5. パフォーマンスの測定

### (1) ファイル転送のパフォーマンス測定

1 ブロックが 576 バイトのファイルを 1 hop<sup>\*1</sup> 離れた 2 つの Host 間で転送した場合の測定によると、JIPNET の FACOM 230-75 から HITAC 8450 へのスループットは 25.7 K ビット/秒、round trip time<sup>\*2</sup> (RTT) は 179 ミリ秒である。F-H 間では別に 2 hop の迂回ルートもあるので、それも併用するとスループットは 30.4 K ビット/秒に増加し、RTT は 153 ミリ秒に減少し、adaptive routing の効果が現われている。

同様な実験<sup>[2]</sup> が 1970 年ごろの ARPANET を使って SRI の XDS 940 から Utah 大学の PDP-10 で行なわれており、スループットは 26.3 K ビット/秒、RTT が 175 ミリ秒が測定されている。その後 IMP が改善されており、現在測定すればスループットは 30.3 K ビット/秒、RTT は 152 ミリ秒になることが推定される<sup>[3]</sup>。

2 つのネットワークの測定結果を比較するとほぼ等しくなっている。これは両ネットワークとも同じ規模の回線、IMP コンピュータとして Host コンピュー

\*1: 通過する回線の個数

\*2: ここでは、プロセスが X-メッセージ送信要求を出してから、それに対する acknowledgment がプロセスに返されるまでに経過した時間と定義する。

タを使っているため、プロトコルの差異は影響してないことがわかる。

(2) ネットワーク・パフォーマンスの測定

JIPNETを使って図2のような迂回ルートのあるトポロジで2つのHost間へ2つの異なるメッセージ長に対してコネクション数を変化させて総合パフォーマンスを測定した(図3参照)。

最初のケースは、パケットのオーバーヘッドが最小になるようにメッセージ長を最大パケット長(full packet)の1152ビットにした。そのようなメッセージに対する最大スループットはルートR<sub>1</sub>のみ使用した場合は約41 Kビット/秒であり、active connectionが4個でほぼ上限に達している。

これはサブネットのフロー・コントロールが原因であり、発信地IMPと目的地IMP間にパイプという論理的なバスを張り、そのパイプの容量を制限して輻輳を防止している。サブネット内ではメッセージを単一パケット・メッセージ(1152ビット以下)と複数パケット・メッセージ(9216ビット以下)とに分けて、単一パケット用のパイプ(S-パイプ)を4本と複数パケット用のパイプ(L-パイプ)を1本設けている。各パイプはコンカレントにメッセージを流すことができ、複数パケット・メッセージは目的地IMPにバッファを確保した後L-パイプを渡し、一時に1個だけ流すことができる。

S-パイプは4本であるので、4本以上のコネクションを張ってコンカレントにメッセージを送信しても、S-パイプのバウンドで抑えられて、スループットは増大しない。

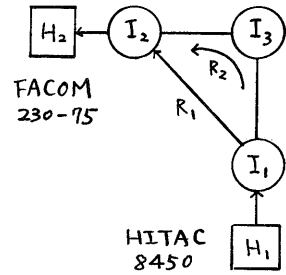


Fig.2 Network Topology

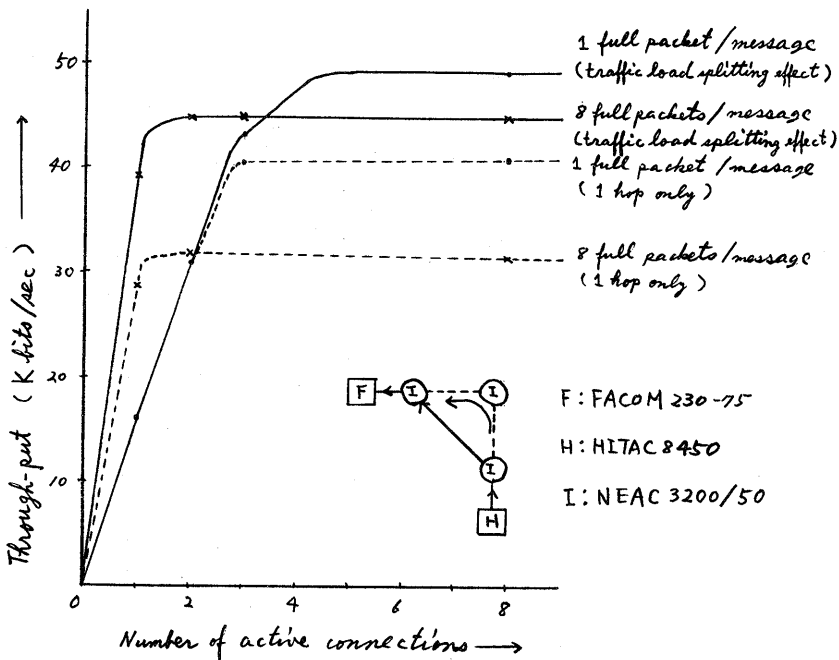


Fig.3 Maximum Through-put Measurements



ルート  $R_1$  と迂回ルート  $R_2$  を併用した場合の最大スループットは約 50 Kビット/秒に増加する。

これはサブネットの adaptive routing の効果である。Adaptive routing としては Shortest Queue + Bias を採用しており、最短ルートの出力回線の送信待ちキューが長くなると迂回ルートの出力回線が選択される。

第 2 のケースは、メッセージを最大メッセージ長 (8-full packet) の 9216 ビットにした。そのようなメッセージに対する最大スループットはルート  $R_1$  のみを使用した場合は約 32 Kビット/秒であり、非常に効率が悪い。これはレーパイプが 1 本しかないため、別のコネクションのメッセージはレーパイプが空くまで送信が待たされるためである。ルート  $R_1$  と迂回ルート  $R_2$  を併用した場合の最大スループットは約 45 Kビット/秒に増加する。これはメッセージの一部 (パケット) が迂回ルートを通り、全体としての RTT を減少させており、レーパイプが早く空き状態になるためである。

ロング・メッセージの最大スループットは 1970 年ごろの ARPANET の約 70 Kビット/秒<sup>[2]</sup> に比べると非常に悪くなっている。このころの ARPANET ではロング・メッセージを特別扱わず、リンク毎にフロー・コントロールを行っていた。しかし特定の Host への active リンクが 5 本を越えると reassemble lockup が起りうるミッドが判明して、その後フロー・コントロールが強化されたためたがんと低い値 (31.1 Kビット/秒)<sup>[3]</sup> になると推定される。

次に単一コネクションでメッセージ長を変化させて RTT を測定した (図 4 参照)。

RTT は、1 フル・パケット・メッセージの場合 1 hop では 68 ミリ秒、2 hop では 100 ミリ秒である。サブネットの測定から推定すると、両 Host の送受信処理に約 32 ミリ秒が加わっている。

8 フル・パケット・メッセージの場合 1 hop では 298 ミリ秒、2 hop では 337 ミリ秒である。1 hop の場合、2 hop の迂回ルートを用いると 217 ミリ秒に減少する。この 217 ミリ秒は 1 hop ルートに 5 パケットを流した場合に相当し、このことから 1 hop ルートへ 5 パケット、2 hop ルートへ 3 パケットを振り分けたミッドが推定される。

図 4 から RTT の式  $Tr(n, p)$  を求めると、

$$\langle 1 \text{ hop の場合} \rangle \quad Tr(1, p) = 60 + 30p \quad (1)$$

$$\langle 2 \text{ hop の場合} \rangle \quad Tr(2, p) = 100 + 30p \quad (2)$$

$$\langle 1 \neq 2 \text{ hop の場合} \rangle \quad Tr(1 \neq 2, p) = 88 + 16 \times \{p + \text{mod}(p, 2)\} \quad (3)$$

である。ここで  $n$  は hop 数、 $p$  はメッセージに含まれるパケット数である。各々

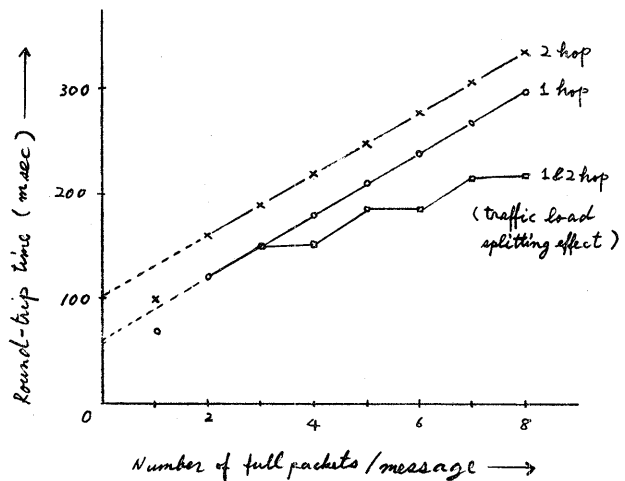


Fig. 4 Minimum Round-Trip Time Measurements

の式の定数項は第1パケットとRFA\*1などのフロー伝送制御パケットが目的地IMPへ到着するのに必要な時間の和であり、hop数の関数である。pの係数30は1フル・パケットが1hopだけ前進するのに必要な時間であり、30pは最終パケットが先頭パケットよりもどの位遅れて目的地IMPへ到着するかを表わしている。

以上から次のRTTの近似式 $Tr(h, p)$ が導き出される。

$$Tr(h, p) \doteq 60h + 30p \quad (4)$$

p=7(8064ビット)の場合は、

$$Tr(h, 7) = 60h + 210 \quad (5)$$

であり、ARPANETのフル・メッセージ(8063ビット)の場合の式〔4〕

$$Tr^A(h, 8) = 60h + 200 \quad (6)$$

とほとんど一致する。

式(3)からadaptive routingにより $p=2n$ のRTTが $p=2n-1$ のRTTに等しくなることが解かる。

## 6. おわりに

以上でNCPの実現上の向題点、検討および測定について述べたが、今後はコンピュータ・ネットワークをコンピュータ・ユーティリティ(バーチャル・ネットワーク)へと発展させるための共通ネットワーク制御言語(Network Job Control Language, NJCL), 各Hostに散在している共用データ・ベースへのアクセス(Distributed Data Base Access, DDBA), 複数Hostを併用して実行されるアプリケーション(Distributed Application)などの開発を容易にするためのサポート・プログラムが必要となってくるであろう。

本論文作成のために、JIPNETの測定に協力してくれた開発第一課の各位に深く感謝の意を表わします。

最後に、紙面の関係で本文で直接引用した文献のみを示し、参考にさせて頂いた文献の一覧表は発表当日に配布する予定である。

## 参考文献

- [1] IBM, "IBM System/360 Operating System: System Programmer's Guide," Form C28-6550-6
- [2] G. D. Cole, "Performance measurements on the ARPA computer network," Proc. Second Symposium on Problems in the Optimization of Data Communications Systems, pp. 39-45, Oct. 1971
- [3] J. M. McQuillan, et al, "Improvements in the design and performance of the ARPA network," FJCC '72, pp. 741-755
- [4] D. C. Wood and R. K. Trehan, "An assesment of the performance of packet switching networks," EUROCOMP '75, pp. 1-11

\*1: Request For Allocation の略で、発信地IMPが目的地IMPにRFAパケット(80ビット)を送り、目的地IMPにバッファを確保する。