

分散型計算機システムの性能予測ツールの開発

本山博司, 大町一彦 (日立システム開発研究所)

1. 序論

分散型計算機システムがいよいよ普及しようとしている。分散システムには大別して地理的分散システムと機能分散システムとがあるが、いずれにしても、分散システム設計者は複数のノードとそれをつなぐリンクをどのように構成するか、その構成によって所期の性能が出るか否かを決定し、判定しなければならない。しかもその判定は未だシステムが稼働しないうちに行なう必要がある。すなわちシステム設計者はどうしても性能予測をしなければならない。この事自体はシングルノードの計算機システムの設計者にとっても同様ではある。しかし分散システムはシングルノードのシステムに比べてノードの数は複数がであり、しかもメッセージがノード間を渡り歩くなど従来無い要素と加わって格段に複雑になって来た。そこで性能評価ツールが必要になるのだが、分散システム用のしかも設計初期段階での性能評価ツールが、シングルノードの詳細を極めたシミュレータと同じ発想で造られてよいだろうか。設計初期段階では、後期段階と異なり、回線バッファ長や、こまかな通信制御方式上の相違、CPU走行ステップ数の微細な変化等は不明である。設計初期段階ではノードとリンクの接続のしかたや、ノード内の各資源の基本性能が、これによりのかが知りたく、もしよくならぬのであれば、別のシステム構成でどうか等々、茫洋とした設計選択枝の中からできるだけ早く、しかも誤りなく正しい幹子を見出したいのである。したがってこの設計段階で用いるべきツールは、少々の精度は落ろして、少数の基本的設計データのみに力すれば、素速く結果をユーザに示せるようなものが必要ではない。

筆者等は最近、以上のような段階における性能設計のあり方を吟味し、それに適しいツールを開発し設計者の支援に供しているが、ここに報告したい。なお本報告は、第22回大会の報告[7]を發展させたものである。

この報告では、まず第2章で、分散システム性能評価における問題点を挙げ、その解決策について既存の研究成果を検討する。第3章で、既存成果で不十分な点に関する解決策を示し、その効果に言及する。第4章では、仮想的かつ現実的な分散システムに関する実際の評価例を示す。第5章では結論と今後の課題を記して全体を終る。

2. 分散型計算機システム評価の問題点と既存の研究成果

2.1 評価時間の長大化

性能予測の手法にはシミュレーションと解析手法の二種あるが、モデリング能力が高い点から、従来多くシミュレーションが行なわれてきた。シングルノードの計算機システムをシミュレートする際、応答時間の他の評価項目の値が十分安定するまでには、通常数10分のCPU時間を要する。従ってこれと同程度の細かさで分散システムをモデル化し、シミュレーションを実行すると、少なくとも、このノード数倍は時間がかかる。ちなみに、ノード当たり20分を要するモデルで、10ノードある分散システムを評価するには、少なくとも3時間20分のCPU時間が必要である。

この長さは、評価を1ケースだけ行なう場合には、禁止的ではない。しかし設計時には、多数のシステム構成代替案を評価したい α が実情であり、この場合は評価時間を十分短かくしておかなければならない。

シミュレーション時間を短縮する方法として、モジュラ分解 (decomposition)がある。分解の考え方は COURTOIS [1]により提出された。[1]では分解を解析的手法について適用しているが、基本的にはシミュレーションモデルを分解してよい。分解し、上下二段階にしたモデルをさらに高速化するため、一方を解析モデルで評価する方式、いわゆるハイブリッド方式がある。SCHWETMAN [2]は、単純なハイブリッドモデルについて、この方式の有効性を示し、CHIU and CHOW [3]、木下 [8]は、具体システムについて、ハイブリッド方式で評価し、容量計画 (capacity planning) に適用し効果も上げた。しかし分散システム性能評価を行なうには、さらに高速化が必要であり、本報告ではこの点につき新しい工夫をした。

2.2 分散システム特有のトランザクション動作のモデル化

分散システムでは、あるノードに入ってきたトランザクションを全てそのノード内で処理するとは限らない。処理の一部は、他ノードで行なわれ、その結果を用いて、後々処理を行なう。このような動作は地理的分散システム、機能分散システム双方に見られるが、集中システムには無い。この点が自由にモデル化できるしくみを、設計者に提供しなければならない。またノードとリンクの接続法(ネットワークトポロジー)と、簡単に定義できるようにしなければならない。

以上の点を考慮したコンピュータネットワーク専用のシミュレータには、田中、中村 [4], [5] および CHIU and CHOW [3] がある。これらは、性能評価の焦点を通信系に置き、各通信制御方式固有の動きについて十分な精度でモデル化している。このためネットワーク定義は十分に表現できるが、ノード内での処理については大巾に簡略化している。本報告では焦点をネットワークからノード内処理の方へ移し、通信系を上記方式より簡略化し、ノード内処理及び、トランザクション内部でのノード間渡り処理について、より詳細化した。

3. 問題点の解決法とその効果

3.1 評価時間の短縮

(1) 利用率の事前チェック

システム内に1個でも利用率が100%を超えると α が有れば、シミュレーションは安定せず、待ち行列内のトランザクション数は次第に増加する。従って特にこのような状態における過渡応答を知りたいのでもないが、シミュレーションを行なう意味は無い。このようなムダを避けるため、シミュレーションに入る前に、利用率 α の計算を行なう。利用率は、トランザクションの到着率 λ と、トランザクション当りの資源使用量が与えられれば閉ネットワークについては、原理的に簡単に求まる。しかしながら、後述するようにトランザクションの内部が複数個の処理に分かれ、その各々について資源使用量を指し示していく場合には、利用率の計算は、まず各々の処理毎に行ない、次にそれらの総計を行なうという二段階を経なければならぬので、複雑である。従って利用率の計算をツール内で自

動的に行なう事には意味がある。

3.2 モデルの階層化、ハイブリッド化、静的インタフェース化

モデルの分解により、上下2段階に階層化する事が有効である事は、SWETMAN[2]で示されている事は既に述べた。[2]ではシステム資源を主メモリ、磁気テープ装置等の「長期間資源 (long-term resource)」と、CPU、チャネル等の「短期間資源 (short-term resource)」とに分け、前者をシミュレーションで、後者を解析的に評価している。その理由は短期間資源は時間当りの使用回数が多く、これを解析化すればより評価時間短縮効果が高いからである。我々も基本的にはこの方式を取る事にした。CHIU and CHOW [3]も同じくこの方式を実施している。

しかしながら、[2],[3]と異なり、階層間の情報の受授が動的である点が、報告者等の静的インタフェース方式と異なる。以下この点について説明しよう。

(a) 動的インタフェース

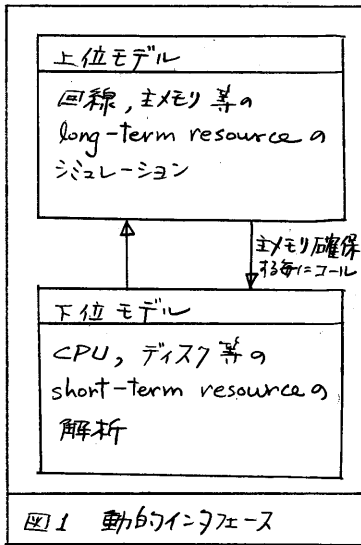


図1 動的インタフェース

上位モデルと下位モデルの動的インタフェースとは、上位モデル中でトランザクションが主メモリの占有/解放毎に下位モデルを動作させ、そのタスクの処理終了時刻を予測し、その予測時刻でタスク終了事象を発生させるよう上位モデルがシミュレートする方式を言う。すなわち主メモリ占有時間を、トランザクション毎に評価するわけであるから、上位モデルと下位モデルは交互に頻りに動作する。下位モデルの機能を数式で書けば、次のようになる。

$$T = R(N) \quad \text{----- ①}$$

但し、 $T = (T_1, T_2, \dots, T_i, \dots, T_m)$
 T_i : トランザクションタイプ i の処理終了時刻
 $N = (N_1, N_2, \dots, N_j, \dots, N_m)$
 N_j : トランザクションタイプ j の主メモリ内多重度
 R : N より T を評価する機能

(b) 静的インタフェース

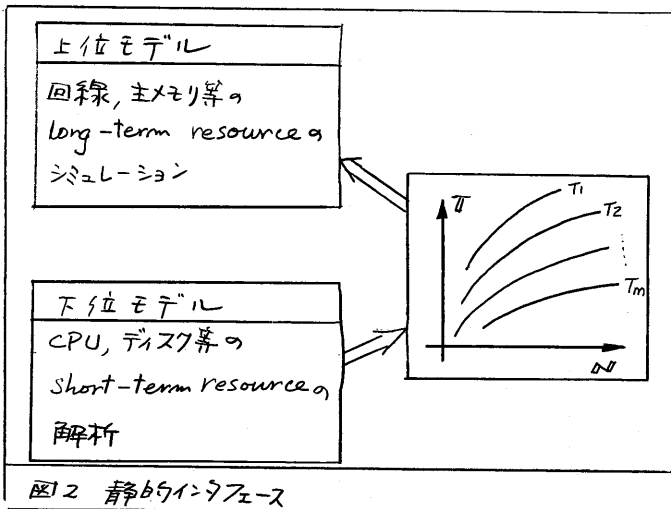
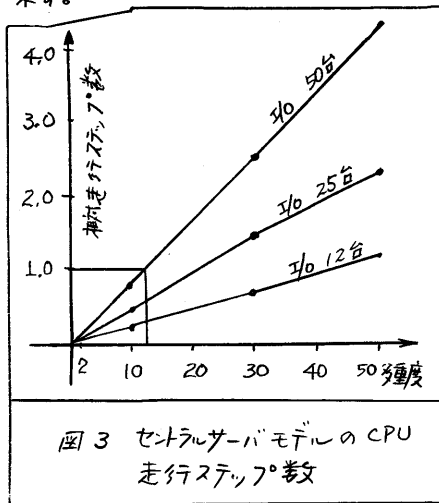


図2 静的インタフェース

静的インタフェースとは①式の $T = R(N)$ を、毎回計算するかわりに、左図のように計算結果をテーブルの形で持つておく方式を言う。従って下位モデルは、このテーブルを作るために、1回だけ動作すればよく、上位モデルはこのテーブルを引くのみで、タスク終了時刻が予測できる。

静的インタフェースを取る事により、下位モデルの走行回数は、次のように短縮された。下位モデルとして BUZEN [6] のセントラルサーバモデルを用いた場合、I/O 台数 12~50, 多重度 2~50 における CPU 処理ステップ数実測値を図 3 に示す。



これはトランザクションタイプが1種類のみの場合である。但し、値は I/O 50台, 多重度 12 の時の走行ステップ数を 1 単位として表わした。以下 I/O 50台, 多重度 12 のケースについて考えるとする。いま、シミュレーションを行い、全部で m 回トランザクションが、主メモリを占有したとすれば、メモリの解放も同様に m 回だから計 m 回確保/解放される。動的インタフェース方式では、この場合

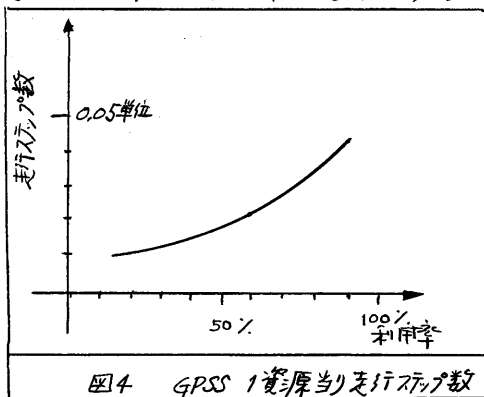
$$2 \times m \times 1.0 \text{ 単位}$$

の走行ステップ数が下位モデルで費されるのに対し、静的インタフェース方式で費されるのは、

$$1.0 \text{ 単位}$$

のみであり、通常 m は $10^3 \sim 10^5$ はあるのだから、その差は $2m$ 単位としてよい。すなわち静的の方が、動的より $2m$ 単位ステップ数だけ走行ステップ数が短かくて済む。しかし静的インタフェース方式はテーブルサイズに限界があるため、トランザクションタイプが増え、各タイプ内の多重度が大きくなると、実行できない

次に上下モデル間のインタフェースは上記と同じように静的ではあるが、下位モデルがセントラルサーバモデルのような解析手法ではなく、GPSS で組んだシミュレーションモデルにした場合を想定し、下位モデルの走行ステップ数を予測してみよう。モデルは上記の場合同様 I/O 50台, 多重度 12 とする。1 トランザクション当たり平均 P 回の I/O 要求を出すとすると、CPU は P 回使用する事になるが計 $(2P+1)$ 回資源を使用する。一般にシミュレーションでは資源の利用率により、資源当たりシミュレーションに要する走行ステップ数は異なる。GPSS による M/M/1 待行列の走行ステップ数を実測すると、図 4 のようになる。



いま、CPU と I/O の利用率がともに 60% であると仮定すると、図 4 より走行ステップ数は 0.02 単位である。また、シミュレーションは m 個のトランザクションについて行なうとすると、結局、下位モデルのシミュレーションに要する走行ステップ数は、

$$(2P+1) \times 0.02 \times m \text{ 単位ステップ}$$

となる。

ミニで上位のシミュレーションモデルの評価に要する実行ステップ数を、トランザクション当り、S単位ステップとして、以上の結果を図5にまとめる。

上位モデル	VS (単位ステップ)	
下位モデル	解析	方式I 1.0
		方式II 2V
	シミュレーション	方式III (2P+1) × 0.02 × m

単位: 「単位ステップ」すなわち下位モデル
αステップ数を1.0とした時の比
V: 全完了トランザクション数
S: 上位モデルの1トランザクション
当りの実行ステップ数

図5. 各方式の実行ステップ数比較

全体の実行ステップ数を各方式について比較してみよう。便宜上図5の如く各方式を方式I, II, IIIと呼ぶ。

(a) 方式Iと方式IIの比較

方式Iのステップ数はVS+1.0, 方式IIはVS+2Vであるから、方式IとIIの比は、下のようになる。

$$\begin{aligned} & (I) / (II) \\ &= (VS+1.0) / (VS+2V) \\ & \approx VS / (VS+2V) \\ &= S / S+2 \quad \text{----- (3)} \end{aligned}$$

すなわち、S ≈ 1.0で、上位モデルの実行ステップ数が下位モデルと同程度ならば方式Iは方式IIの1/3の実行ステップ数で評価を完了できる。逆に、上位シミュレーションモデルで複雑なシミュレーションを行ない、S ≈ 10となると、(3)は10/12 = 83%となり、両者の差はほとんど無くなる。しかしながら、常にIの方がIIより速い。

(b) 方式IIと方式III

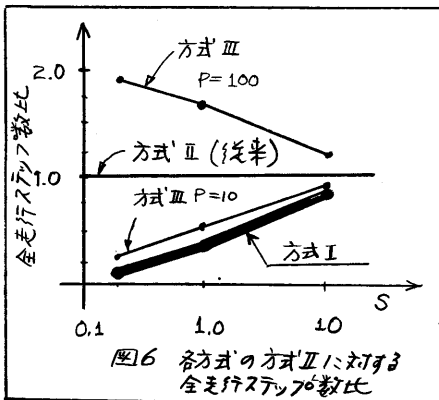
上記(1)における比較と同様にして、

$$\begin{aligned} & (III) / (II) = (VS + (2P+1) \times 0.02 \times m) / (VS+2V) \\ &= (S + (2P+1) \times 0.02 \times m / V) / (S+2) \quad \text{----- (4)} \end{aligned}$$

ミニで、上位シミュレーション、下位シミュレーションが安定するに要するトランザクション量V, mが、ほとんど同程度(通常10³~10⁵)とすると、

$$(4) = (S + (2P+1) \times 0.02) / (S+2) \quad \text{----- (5)}$$

となり、トランザクション当りI/O使用回数P = 24.5を境界にして、Pがこれより小さければ方式IIIの方が実行ステップ数が低くなる。方式IIIをつかって、下位モデルを、さらに複雑にしていって、1トランザクション当りの実行ステップ数が、2単位を超えれば、方式IIよりIIIの方が実行時間が長くなる(2)が、それがわりシミュレータの高いモデリング能力をひき出す事ができる。

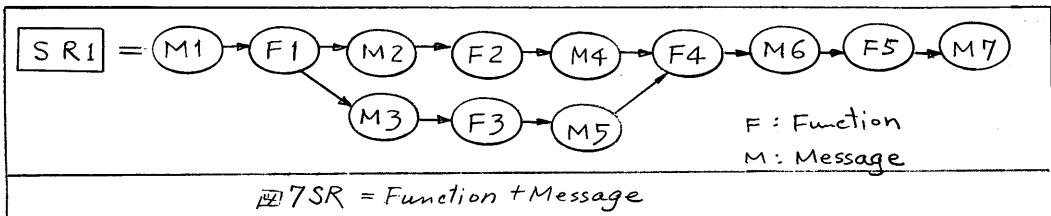


(3), (4)式を、上位シミュレーションモデルの1トランザクション当りステップ数Sについてプロットしたものを図6に示す。図6は方式II(従来αハイブリッド方式)を1.0として、他の方式と比較したもので、1.0より上になれば従来方式より遅く、下になれば従来方式より速い事を意味する。

3.3 分散特有動作のモデル化支援機能

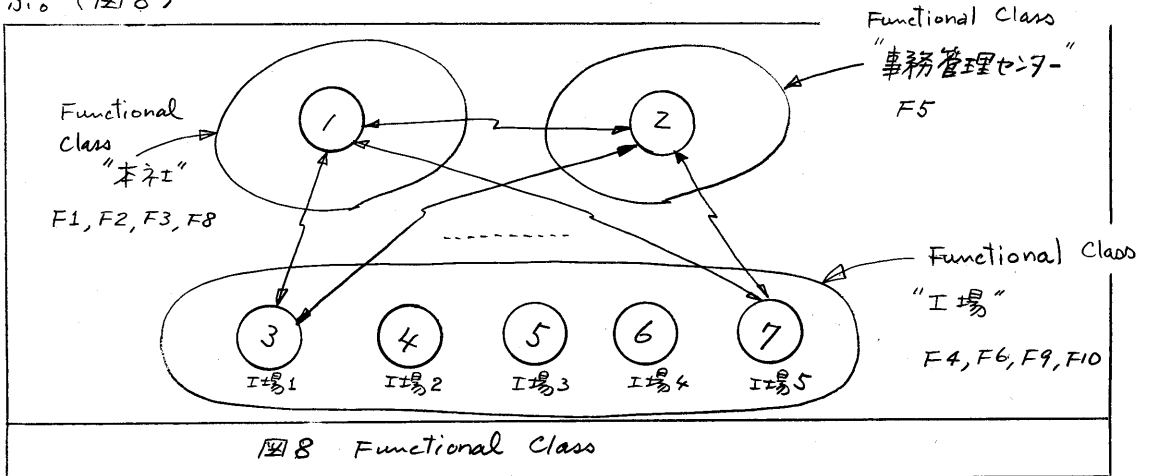
(1) Service Request = Function + Message (図9)

端末から入力されるオンラインランザクション、あるいはTSSのコマンドあるいはバッチジョブ等、システムに対する処理要求は様々なるが呼ばれている。これから以後Service Request, SRと総称する。ツール内であつた)処理要求は、全てSRであり、その他の形式は無い。分散システムでは、1個のSR内で、複数のノードにまたがって処理するものが必ずある。(もしよければ、分散システムではなく孤立したシステムが単に複数個有るにすぎない)このようなSRの性格をモデル化しやすきようは、"Function"を作った。FunctionはSRの内部に有り、SRの構成要素である。FunctionはSRの内部に複数個あつてよく、各Functionは、必ず1個のノード内で処理される。1個のFunctionの一部分があるノードで処理され、他の部分があるノードで処理されるような事は無い。しかし、SR内の複数のFunctionが、別々のノードで処理されてよい。Functionに入力されたり、Functionが出力したりする情報をMessageと呼ぶ。Functionは必ずMessageを入力し、必ず出力する。

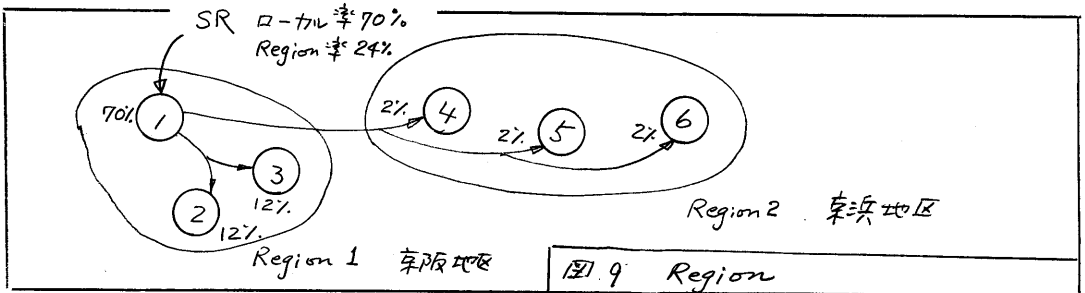


(2) Function Class / Region

分散システムを構成する各々のノードは、システムに入ってくる処理要求の全てを処理できるわけではない。個々のノードで処理できる業務内容は決まっている。この事をモデル化しやすきようは、各ノードをグループに分け、それらのグループ内のノードはどれも、共通のFunctionを処理できるが、それら以外のFunctionは処理できないものとした。このグループを"Functional Class"と呼ぶ。(図8)



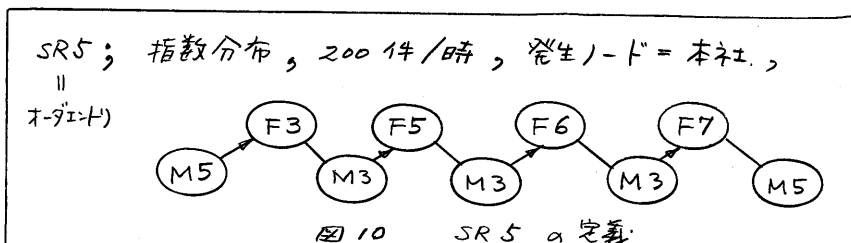
また、ファイルが分散している場合など、そのSRが入って来たノード内のファイルを多く参照し、他ノードのファイルは少ししか参照しないといった事がよくある。例えば、全国に分散したデータベースで、大阪でおきた検索アトラクションは、ほとんど関西圏内のファイル参照のみ、少数が東京のファイルを参照するのである。このような「ファイル参照の地域性」をモデル化するため「Region」の考え方を導入した。1つのRegionは、ノードの集合で、全ノードは必ずどこかのノードに属するものとする。各SR内のFunctionは、SRが到着したノード内で処理される率（ローカル率）、SRが到着したノード以外だが同一Region内のノードで処理される率（Region率）を指定しておく。シミュレーターは、この率に従って、Functionの処理を各ノードに配分していく。このように処理すべきノードは、当然そのFunctionを処理できるFunctional Classに属していなければならぬ。



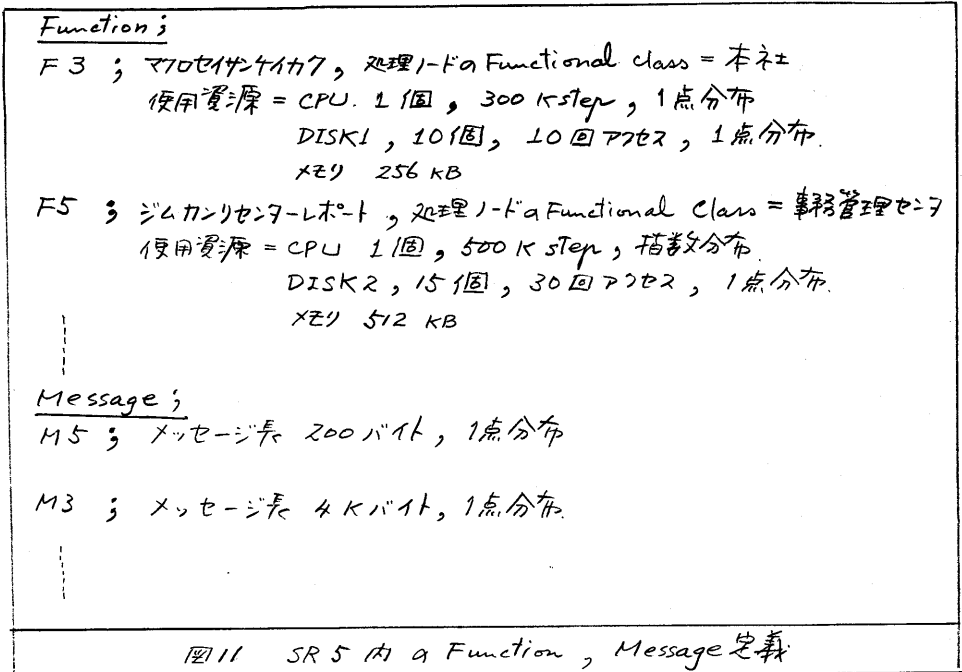
4. 評価例

以上述べてきた技法を用いて、仮想的な分散システムについて、性能評価を行なってみたい。システムは図8のように、本社、事務管理センター、工場とから成っている。本社ではオーダーリーのトランザクションを受けつけ、それにとどまらずマクロ生産計画を立て、工場に生産を指令する。工場は指令にとどまらずミクロ生産計画を立て、部品発注をし、生産を開始する。また、製品が完成に到るまで、その状態管理を行なうため、工場は独自の製品トラッキングを行なう。

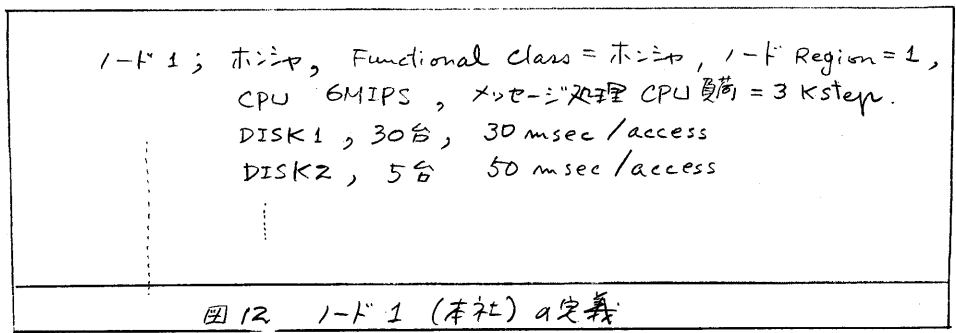
SRは全体で12種類あり（SR1～SR12）そのうち、「オーダー」と呼ばれるSR5は、次のような内訳となっている。



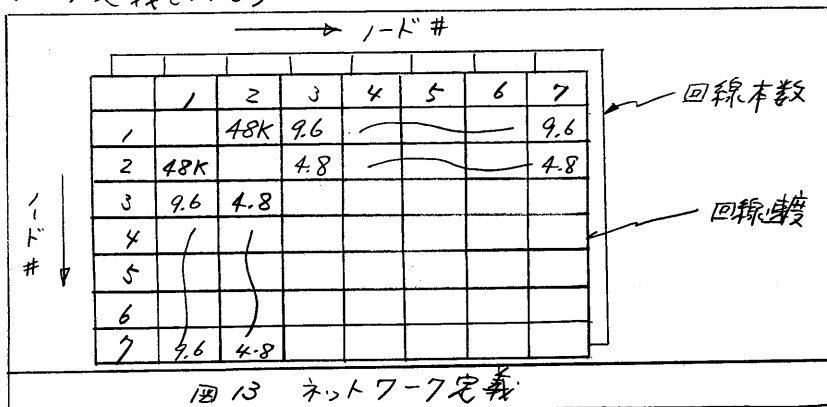
さらに、SR5内のMessage, Functionの定義とする。



次に、各ノド内にある資源の定義及び Functional Class の定義を行なう。



最後にネットワーク定義を行なう)



以上のように入力データに対し、以下のような結果を得た。(全CPU時間 8分10秒)

2. LINK UTILIZATION (%)

(ORIGIN NODE)		1	2	3	4	5	6	7
IDS.NI	1	0.0	4.1791	2.4801	1.3801	1.4121	1.7081	1.5661
	2	13.9011	0.0	4.1181	3.9801	3.8401	3.2611	3.0951
	3	23.1701	34.8041	0.0	0.0	0.0	0.0	0.0
	4	26.1691	28.7751	0.0	0.0	0.0	0.0	0.0
	5	16.4611	14.1681	0.0	0.0	0.0	0.0	0.0
	6	16.4561	28.7331	0.0	0.0	0.0	0.0	0.0
	7	17.6131	32.0171	0.0	0.0	0.0	0.0	0.0

表1

DYNAMIC MODEL (上位モデル)

による
回線利用率
評価結果

表2

STATIC MODEL による
利用率事前 check (100%以下) OK

3. RESOURCE UTILIZATION BY RESOURCE CLASS

RESOURCE CLASS (%)		1	2	3	4	5
NODE	1	10.5141	33.8081	2.9171	36.8081	0.0391
	2	1.8971	13.2501	0.0	0.0	0.0071
	3	60.3311	0.8101	48.2761	71.4251	1.1691
	4	60.3311	0.8101	48.2761	71.4251	1.1691
	5	60.3311	0.8101	48.2761	71.4251	1.1691
	6	60.3311	0.8101	48.2761	71.4251	1.1691
	7	60.3311	0.8101	48.2761	71.4251	1.1691

表3

DYNAMIC MODEL による
利用率評価結果

1. RESOURCE UTILIZATION BY NODE (%)

(RESOURCE CLASS)		1	2	3	4	5
NODE	1	15.3721	34.8001	4.5001	70.5001	3.6361
	2	2.0161	14.4001	0.0	0.0	0.3221
	3	57.0751	1.2001	53.1001	59.8501	28.3121
	4	55.5001	0.7501	47.7001	62.5501	28.1861
	5	56.9251	1.0501	48.4501	64.3501	23.1191
	6	60.3751	0.7501	46.9501	73.0501	27.5911
	7	58.2001	1.0501	46.9501	68.4001	27.4511

CPU

DISK

表4

3. SERVICE REQUEST PERFORMANCE DATA				
SR.#	SERVICE REQUEST DESCRIPTION	NUMBER COMPLETED	MEAN TURNAROUND TIME (SEC)	VARIANCE OF TURNAROUND TIMES
1	セイベン ハッチユウ	41	0.4741	0.088
2	ハコハ"イ トウケイ/ハコハ"イ ケイカ	41	0.4741	0.088
3	セイザン ケイカク	41	1.4401	0.325
4	クミタテ シ"ユウシ"ヨ ケイカク	41	1.8041	0.318
5	17"ヒコ ハッチユウ	131	4.3041	0.912
6	セイベン アテ"タス	241	2.1141	1.424
7	ヒキアテ	2191	3.9241	4.132
8	1"ユウカ ヲウ"	15561	0.9681	0.032
9	1"セツケイ"チ"タテツク"ウ	41	1.6091	0.232
10	1"ケ"ン"タイ"セツケイ"1/ホカ	22141	0.5041	0.024
11	1"ケ"ン"タイ"セツケイ"2/ホカ	531	4.1061	1.070
12	17"ヒコ"チ"タ"ス"メ"カ"チ"タ"ス	241	2.1331	0.914

完了個数

平均応答時間

応答時間標準偏差

DYNAMIC MODEL (=上位モデル)

による SR 応答時間
評価結果

5. 結論

分散型計算機システムの性能評価ツールを開発した。この種のツールの課題は性能評価時間の短縮と、分散型有のモデリング機能支援の充実にある事から発し、時間短縮のためはSCHWETMANのハイブリッド方式をさらに高速化した方式を採用した。モデリング支援のため、トランザクションを複数の「ファンクション」と「メッセージ」で構成し、さらにノードを Functional class と呼ばれるグループに分け、各々の Functional class 内のノードが共通のファンクションを処理できるように定義する機能を設けた。これにより、トランザクション内の処理を、複数ノードで行なうようなモデルが簡単にできるようになった。

現在のところ、実験的分散システムの実測値により、本方式の結果の検証を行っている。また、ここで採用したハイブリッド方式の評価時間の測定と組織的に行なっている。今後この二点を課題とした。

6. 参考文献

- [1] Courtois, P. Decomposability, instabilities, and saturation in multiprogramming systems. *Comm. ACM* 18, 7 (July 1975), 371-377
- [2] Schwetman, H.D. Hybrid Simulation Models of Computer Systems. *Comm. ACM* 21, 9 (Sept. 1978) 718-723
- [3] Chiu, W. and Chow, W.M. A performance models of MVS. *IBM sys. J.* 17, 4 (1978) 444-463
- [4] 田中, 通信ネットワーク計画性能設計支援ツール SPACNET-Pの開発
第23回 情報処理大会 (予定)
- [5] 中村, 通信ネットワークの性能設計 - SPACNET-Pによるアプローチ
第23回 情報処理大会 (予定)
- [6] BUZEN, J. Computational algorithm for closed queueing networks with exponential servers. *Comm. ACM* 16, 9 (Sept. 1973) 527-531
- [7] 本山, 大町 ハイブリッド手法を用いた分散システム性能評価手法
第22回 情報処理大会
- [8] 木下, 石沢 解析手法を取り入れた計算機性能評価シミュレータの開発
第22回 情報処理大会