



Title	高度インテリジェントネットワークにおけるサービス制御ノード構成技術に関する研究
Author(s)	平野, 正則
Citation	大阪大学, 1999, 博士論文
Version Type	VoR
URL	https://doi.org/10.11501/3155636
rights	
Note	

The University of Osaka Institutional Knowledge Archive : OUKA

<https://ir.library.osaka-u.ac.jp/>

The University of Osaka

高度インテリジェントネットワークにおける
サービス制御ノード構成技術に関する研究

1998年10月

平野 正則

論文内容の要旨

博士論文題名：高度インテリジェントネットワークにおけるサービス制御ノード構成技術
に関する研究

学位申請者：平野正則

要旨：

フリーダイヤルや携帯電話等、高機能な電話サービスを効率的に実現することを狙いとした高度インテリジェントネットワークにおいて、カスタマ毎のデータを抱え、呼の接続処理の中心的な役割を果たすサービス制御ノードには、需要に応じた柔軟な構成の実現や、カスタマデータの管理機能の高度化、高い信頼性の達成等が要求される。本論文は、このような要求に的確に対応可能な構成を明らかにするため、多数のモジュールを用いた分散型メモリデータベース構成による、高性能、高信頼で経済的なサービス制御ノード構成技術として、その研究成果をまとめたものであり、以下の8章より構成した。

第1章では、本研究の背景と目的、要求条件と課題、従来研究との関係を述べた。第2章では、本研究で対象とするサービス制御ノードのシステム構成やデータベース配備について述べた。第3章では、サービス制御ノードのデータベースのアクセス特性、モジュールの負荷やその偏り、トランザクションの処理量などを総合的に考慮したデータベースの分散構成方式の評価方法を明らかにし、ファイル共用のデータベース構成に比べ、コスト性能比が2～3倍優れた分散型メモリデータベース構成法を提案した。第4章では、カスタマ毎のサービスデータに対するトラヒックが極めて大きくばらつく場合のトラヒック量の近似方法や、複数モジュールにサービスデータを分散配備した場合のモジュール間の負荷の偏りの評価方法を提案し、サービス制御ノードでは負荷の偏りが1～2割程度に収まることを明らかにした。第5章では、モジュール間の結合機構を機能単位に二重化し、50モジュール程度を接続した大規模分散構成においてもシステムダウン時間を年間数分以下にできる高信頼構成法や、モジュールの通信処理能力を向上させる周期制御方式における最適な周期時間の決定方法を提案した。第6章では、障害発生からデータベースが回復するまでの時間と、データベースの単位容量当たりのトラヒック量に着目し、メモリデータベースのリカバリ方式の評価技術を明らかにした。これに基づき、サービス制御ノードのリカバリ方式として、半導体ディスク装置にメモリデータベースのログとチェックポイント時点でのデータベースを取得する効率的なリカバリ方式を提案した。第7章では、多数のノードを保守チームが集中して保守する場合について、ノード内、ノード間で相互に冗長化されたシステムの信頼度の時間的な変化に着目し、駆付けの緊急度合いをシステム内の障害装置数や保守チームの到着の有無によりレベル分けすることを提案した。これに基づき、駆付け保守の方法を設定し、信頼性や駆付け回数の評価を通して、急いで駆け付けなければならない緊急駆付け回数の少ない保守方法を明らかにした。第8章では、本研究で得られた結論を総括した。

高度インテリジェントネットワークにおける サービス制御ノード構成技術に関する研究

目 次

謝辞	1
第1章 序論	2
1.1 背景と目的	2
1.2 研究の位置付け	3
1.3 要求条件と課題	7
1.3.1 サービス制御ノード実現に向けた要求条件	7
1.3.2 サービス制御ノード構成の課題	10
1.4 従来研究と本論文の関係	12
1.4.1 データベースの分散構成技術との関係	12
1.4.2 データベースの高信頼化技術との関係	17
1.4.3 高信頼化運転保守技術との関係	18
1.5 論文の構成と各章の概要	20
第2章 サービス制御ノードにおけるシステム構成	23
2.1 データベースの配備	23
2.2 呼の接続処理	25
2.3 多種類のサービスの実現方法	27
第3章 サービス制御ノードの分散構成技術	29
3.1 緒言	29
3.2 ノード構成条件と機能配備	30
3.2.1 ノード構成条件	30

3.2.2 機能配備	30
3.3 データベースの分散構成	33
3.3.1 サービス制御ノードにおけるデータベースの特性	33
3.3.2 データベースの格納媒体	36
3.3.3 データベース分散方法	36
3.3.3.1 複数モジュールへのデータベースの配備法	36
3.3.3.2 モジュール間の負荷の偏り	39
3.3.3.3 ダイナミックステップ数の評価	40
3.3.3.4 コスト評価	45
3.4 分散処理ノードの高信頼化	46
3.4.1 データベースのリカバリ方法	46
3.4.2 冗長化構成法	55
3.4.2.1 冗長化構成	55
3.4.2.2 制御方式	56
3.5 分散処理ノード構成法	58
3.5.1 アーキテクチャ	58
3.5.2 ハードウェア構成	58
3.6 結言	61
 第4章 モジュール間の負荷の偏りの評価技術	63
4.1 緒言	63
4.2 対象とする分散処理モデル	64
4.2.1 処理方式	64
4.2.2 評価モデル	66
4.3 サービスデータへのトラヒック量	68
4.3.1 サービスデータへのトラヒック量の近似	68
4.3.2 サービスデータの母集団の設定	70

4.3.3	近似式の特性	71
4.4	モジュール間の負荷の偏り	75
4.4.1	負荷偏り率	75
4.4.2	サービスデータ数と負荷偏り率の関係	77
4.4.3	偏り係数と負荷偏り率の関係	77
4.5	分散構成への適用	81
4.5.1	負荷の偏りの影響	81
4.5.2	負荷平準化の可能性	84
4.6	結言	87
第5章 分散処理によるサービス制御ノードのモジュール間結合技術		89
5.1	緒言	89
5.2	サービス制御ノード構成条件	90
5.2.1	ノード構成条件	90
5.2.2	性能条件	92
5.2.3	信頼度条件	94
5.3	モジュール間結合方式	95
5.3.1	モジュール間結合方式の選択	95
5.3.2	ATM結合機構の高信頼化構成法	96
5.3.3	ATM結合機構の全体構成	101
5.4	通信制御方式	103
5.4.1	制御方式	103
5.4.2	通信制御チャネルの応答時間条件	106
5.4.3	評価モデル	107
5.4.4	通信制御チャネルの応答時間評価	107
5.5	評価方法	113
5.5.1	性能測定方法	113

5.5.2	性能測定プログラムの走行方法	115
5.5.3	性能測定結果	115
5.6	結言	117

第6章 サービス制御ノードにおけるメモリデータベースのリカバリ技術 119

6.1	緒言	119
6.2	前提条件	120
6.3	評価モデルの設定	122
6.4	性能評価	128
6.4.1	スループット, レスポンスタイム	128
6.4.2	リカバリ時間	131
6.4.3	データベース容量評価	133
6.4.3.1	収容可能なデータベース容量	133
6.4.3.2	ハードウェア性能向上の影響	138
6.5	コスト評価	140
6.6	信頼性評価	141
6.6.1	不揮発化	141
6.6.2	高信頼化構成	143
6.6.2.1	モジュールとしての耐障害性	143
6.6.2.2	冗長化構成とプロセッサ稼働率	146
6.7	実現方式	147
6.8	結言	151

第7章 相互にバックアップされたサービス制御ノードの高信頼化運転保守技術

7.1	緒言	152
7.2	評価対象システムの構成と動作	153
7.2.1	評価対象システムの信頼度構成	153

7.2.2 モジュールの動作方法	156
7.3 保守方法とシステムの状態	158
7.3.1 保守方法	158
7.3.2 システムの状態	159
7.4 信頼性の評価尺度と保守パラメータの設定	161
7.4.1 信頼性の評価尺度	161
7.4.2 信頼性目標と保守パラメータの設定	163
7.5 駆付けの緊急度と保守モデル	165
7.6 信頼性と保守方法の関係	167
7.6.1 システム信頼性の評価	167
7.6.2 駆付け回数の評価	177
7.6.3 保守形態の及ぼす影響	177
7.6.4 システム、保守センタの配置方法	180
7.7 結言	183
 第 8 章 結論	184
 参考文献	189
 付録	195
 略号一覧	208
 発表論文一覧	214

謝 辞

本論文をまとめるに際して、懇切なる御指導ならびに御助言、御尽力を賜った大阪大学大学院工学研究科教授・池田博昌博士に心からの感謝の意を表します。

本論文をまとめる過程に際し、丁重なる御指導、御教示を賜った大阪大学大学院工学研究科教授・前田肇博士に厚く感謝の意を表します。

また、本論文に対して有益なる御討論、御助言を頂いた大阪大学大学院工学研究科教授・森永規彦博士、同教授・小牧省三博士、同教授・児玉裕治博士、同教授・塙澤俊之博士、大阪大学産業科学研究所教授・元田 浩博士ならびに前大阪大学教授・長谷川晃博士（現・高知工科大学教授）に感謝致します。

本研究は、筆者が日本電信電話株式会社（NTT）・ネットワークサービスシステム研究所（NS研）・高機能処理プロジェクト（N処P）において、高度インテリジェントネットワークに適用するサービス制御ノードの研究開発の一環として担当したものである。本研究開発の遂行に当たり、御指導、御鞭撻を頂いた、NTT・鈴木滋彦取締役（研究開発本部副本部長）、NTTドコモ・弓場英明取締役、NS研・研究開発企画部・花澤隆マネージャ、NS研・N処P・田中公紀マネージャ、N処P・IN装置グループ・今川仁グループリーダ、NTTマルチメディアネットワーク研究所・鈴木孝至部長に厚くお礼申し上げます。本研究の全般に渡って、多くの御討論と御助力を頂いた、東京情報大学・木ノ内康夫教授（当時NS研）に深く感謝致します。また、本研究を推進するに当たり、御討論頂いた、NTTアドバンステクノロジ・重松直樹部長（当時NS研）、NTTドコモ・上坂久一主幹技師（当時NS研）、NTT研究開発本部・研究開発推進部・寺中勝美部門長（当時情報通信研究所）、NS研・研究開発企画部・吉見正信主幹研究員、NTT光ネットワーク研究所・塙澤恒道主任研究員（当時NS研）、NTTグループ企業本部・山根道広担当課長（当時NS研）、NTT情報通信研究所・芳西崇主幹研究員、NTT入出力システム研究所・山崎幹夫主任研究員、拓殖大学・小林正光助手（当時NS研）、NTT東京支社・櫻井秀紀主査（当時NS研）、NTT-TETE九州・林誠治氏（当時NS研）に感謝致します。

最後に、これまで様々な面から御指導、御鞭撻を頂いた、大阪大学・橋本昭洋教授、国際電気通信基礎技術研究所・酒井保良副社長、東京工科大学・松永俊雄教授、NTTアドバンステクノロジ・平松琢弥部長、NTTエレクトロニクス・多嶋清次郎副事業部長、NTT情報通信研究所・武井安彦部長に心からお礼申し上げます。

第1章 序論

1.1 背景と目的

高度情報化社会に向け通信網の高度化、多様化が急速に進展してきている。電話サービスにおいても、電話機の設置場所に係わらず設定された論理的な番号に基づいて接続するサービスや、携帯電話のように電話機の位置を絶えず追跡して位置情報をネットワーク内に登録しておき、電話機の位置に係わらず通話ができるサービスなどが普及してきている。このようにネットワークでカスタマ毎のデータを保持し、この内容に基づいて多様なサービスを実現する網がインテリジェントネットワーク（IN：Intelligent Network）である。我が国では、1985年にサービスが開始されたフリーダイヤルを始めとして、携帯電話、ダイヤルQ²、仮想私設網（VPN：Virtual Private Network）、P H S（Personal Handy-phone System）等のサービスが提供され、インテリジェントネットワークを用いた電話サービスの多様化が進んでいる。

高度インテリジェントネットワーク（高度IN）は、このような状況に対応するため、多様なサービスをタイムリかつ効率的に提供することや、カスタマの要望に応じたサービス内容の変更を容易とすること等を狙いとした次世代のINであり、世界各国で研究、開発が進められている^{(1)～(4)}。この高度INでカスタマデータを抱え、呼の接続処理の中心的な役割を果たすノードがサービス制御ノード（SCP：Service Control Point）であり、需要に応じた柔軟な構成の実現や、カスタマデータの管理機能の高度化、高い信頼性の達成等が要求される⁽⁵⁾。

本論文は、このような要求に的確に対応可能な構成を明らかにするため、需要の変動や機能要求へ柔軟に対応する分散処理化、サービス毎カスタマ毎のデータを効率的に管理するデータベース化、運転保守の効率化とのバランスを考慮した経済的な高信頼化の各面からアプローチしたものであり、多数のモジュールを用いた分散型メモリデータベース構成による高性能、高信頼で経済的なサービス制御ノード構成技術として、その研究成果をまとめたものである。

1.2 研究の位置付け

電話サービスを行うための電話網は、家庭や事務所等に置かれる電話機を接続収容する交換機、および交換機間を相互に接続する伝送路等から構成される。日本全国にはおおよそ四千台の交換機が設置されている。電話機に付与される電話番号は、それを収容する交換機の位置や、交換機内での収容端子等の物理的な位置により一義的に決められている。

電話の接続は、発信者が着信者の電話番号をダイヤルすることにより行われる。まず、発信者がダイヤルを行うと、発信者の電話機を収容している交換機がダイヤルされた番号を受け取る。この番号から着信者の電話機を収容している交換機が一義的に決まるため、発信側の交換機から着信側の交換機までの伝送路を確保し、ダイヤル番号を着信側の交換機に渡す。着信側の交換機では、ダイヤルされた番号に対応する端子に接続された電話機に呼び出しをかけ、着信者が受話器を取ると、通話が出来るようになる。このように、従来の電話サービスの場合は、電話番号により、どの交換機のどの端子位置に接続された電話機に対する呼であるかが一義的に決まる。

一方、フリーダイヤルサービスの場合は、電話機を収容する交換機の設置位置や端子位置によって決まる物理的な番号ではなく、論理的な番号に基づいて接続する必要がある。このため、ダイヤルされた論理番号を物理番号に変換するデータを蓄え、検索するためのデータベース機能が必要とされる。これらの機能を四千台の交換機に持たせることは維持管理が大変であるため、交換機とは別のサービス制御ノードで実現する方法がとられている。交換機や伝送路からなる従来の電話網に、データベース機能を実現したサービス制御ノードを付加した網をインテリジェントネットワークと言う。インテリジェントネットワークでは、交換機や伝送路を伝達レイヤ、サービス制御ノードを高機能レイヤと言う。

インテリジェントネットワークを用いた最初の電話サービスは、1985年にサービスが開始されたフリーダイヤルである⁽⁶⁾ ⁽⁷⁾。フリーダイヤルの初期には、高機能レイヤはサービス制御ノード（SCP）のみから構成され、サービス処理に必要な情報等は直接、端末からSCPへ投入する方法で行われていた。また、顧客情報や空

き番号等は、支店や営業所の独自のデータベースで管理されていた。このため、サービス申し込みを受け付けた時に割り当てた空き番号等、サービス処理に必要な情報を、サービス開始に合わせてSCPへ投入する必要があり、カスタマの申し込みからサービス開始までの業務が煩雑であった^{(8)～(10)}。その後、カスタマ管理の効率化のため、1989年にサービス管理ノード(SMS)が導入され、カスタマ情報のデータベース化が行われた。これにより、SMSに接続された端末から空き番号の割当てや、サービス開始に合わせて、呼処理に必要な情報のSCPへのダウンロードが自動的に行われるようになった⁽¹¹⁾。また、SMSの導入によりカスタマ自身によるサービス条件の変更も可能となった。この時点で、カスタマ情報やサービス条件等をサービス管理ノードで一元的に管理し、呼処理に必要な番号変換情報等がSMSからSCPへダウンロードされるようになり、カスタマ管理はSMS、呼処理はSCPで機能分担する形態が確立された。その後、携帯電話やダイヤルQ²等、サービスの多様化に伴い、サービス毎にSMSとSCPが構築されていった⁽¹²⁾。この段階では、SCPはデータベースアクセス機能のみを有し、サービス制御機能は交換機側にあり、サービスの新規導入には交換機側のソフトウェアの変更が生じ多大の工数と期間が必要であった^{(8) (13) (14)}。

高度インテリジェントネットワーク(高度IN)は、今後より一層のサービスの多様化や高度化に柔軟に対応するため、需要の異なる多種多様なサービスを効率的に提供することや、新しいサービスを短期間に提供することを狙いとした次世代のインテリジェントネットワークである。高度INでは、複数のサービスでSCPやSMSを共用できることや、交換機とSCPはサービスに依存しない標準インターフェースで接続し、新サービスに伴う交換機側での機能追加を不要とすることが要求される。また、ハードウェア構成上は多様なサービスに柔軟に対応するため、サービスに依存しないSCP構成や、プロセッサの追加により処理能力が容易に拡張できること、さらに高い信頼性の達成等が望まれる^{(5) (8) (14)}。

インテリジェントネットワークの初期から高度インテリジェントネットワークに向けた進展の状況を図1.1に示す。

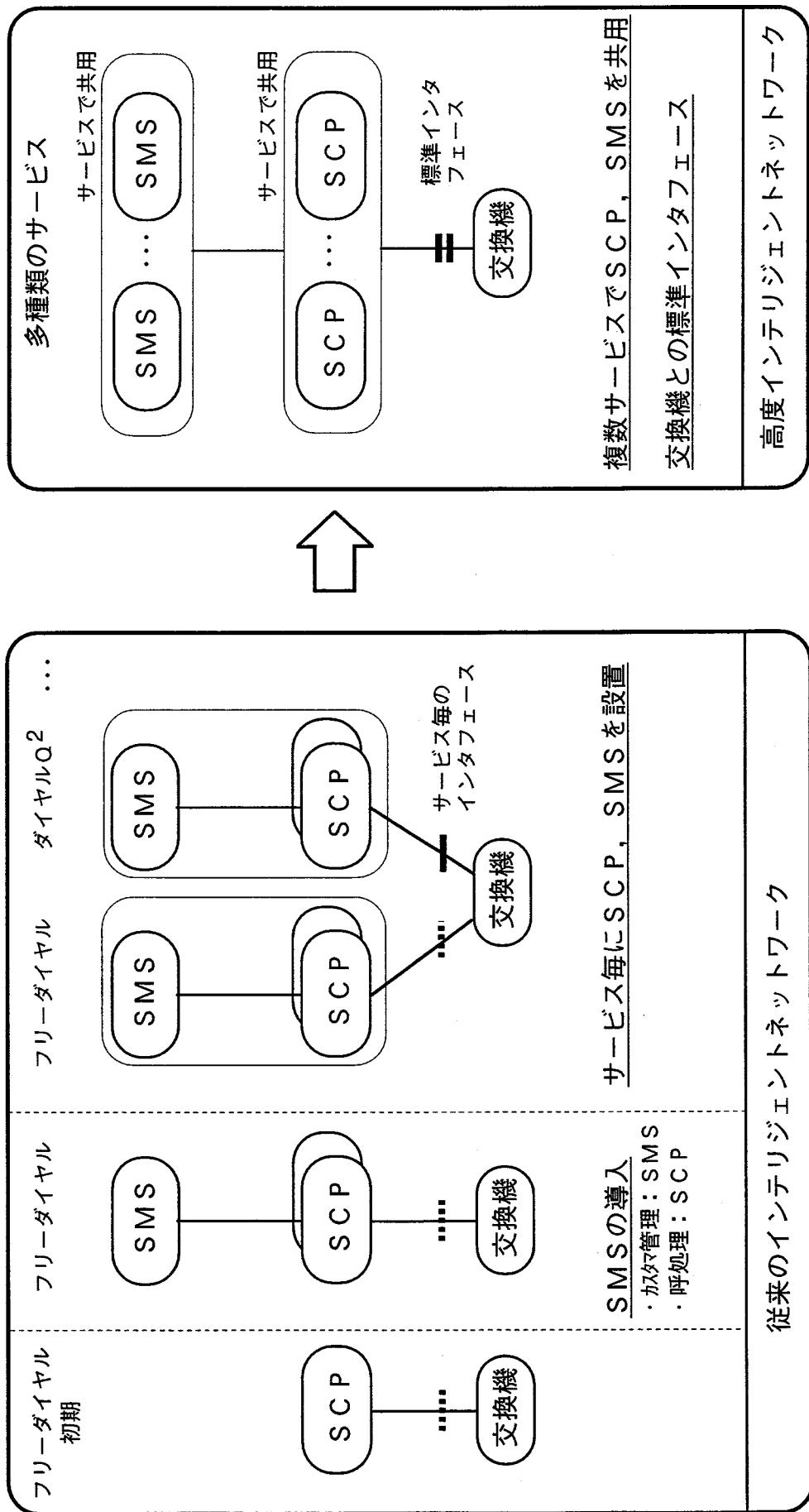


図1.1 インテリジェントネットワークの進展状況

SCP : サービス制御ノード, SMS : サービス管理ノード

高機能な電話呼の接続処理がインテリジェントネットワークでどのように行われるか、図1.1を用いて、以下で簡単に説明する。まず、電話機から投入された電話番号は交換機で解読される。その結果、電話番号が0120で始まるものであればフリーダイヤルであることが判明し、フリーダイヤルに加入しているカスタマのサービスデータを収容しているサービス制御ノード（SCP）に対して0120の後に続く6桁の論理番号を送る。1台のSCPにカスタマデータを収容できない場合は、複数のSCPに分散収容されることとなり、論理番号を送るとき、当該カスタマのサービスデータがどのSCPに収容されているかの判断が必要となる。一般に、カスタマデータの収容に当たっては、6桁の論理番号の上位1～2桁によりどのSCPに収容するかが決められており、上位1～2桁の解読によりどのSCPに収容されているかが判明する。SCPでは論理番号と物理番号の対応表をデータベースとして格納しており、このデータベースにアクセスし、論理番号を物理番号に変換して交換機に送る。交換機では送られてきた物理番号に基づいて電話呼の接続処理を行う。このとき、フリーダイヤルサービスであることから、通常の電話呼とは異なり課金は着信側に行う。交換機で解読した電話番号が0990で始まるものであればダイヤルQ²であると判断し、ダイヤルQ²に加入しているカスタマのサービスデータを収容しているSCPに問い合わせをかけ、同様に物理番号を得て、当該物理番号に基づいて接続処理を行う。このように、サービス毎にどのSCPに問い合わせを行うか等の制御や、サービスに基づく課金処理等は交換機側で行っており、新しいサービスの追加やサービス内容の変更およびSCPの増設等に伴い、交換機側のプログラムの変更が必要となる。

一方、高度インテリジェントネットワークでは、入力された電話番号を解読するところまでは同様に行うが、電話番号が0120や0990で始まるなど、高機能な電話サービスであることを検出すると、具体的なサービス内容が何であるかを判断することなく、あらかじめ決められたSCPに問い合わせをかける。SCP側では、サービス種別を判断すると共に、当該サービスを受けているカスタマのサービスデータがどのSCPに収容されているかを判断し、当該SCPにこの問い合わせを転送する。以下、呼の制御はカスタマのサービスデータを収容しているSCPで行われる。SC

P から交換機への指示も、たとえば、物理番号に接続し、課金を着信側に行う等、サービスに依存しない指示で行われる。このため、新しいサービスの追加やサービス内容の変更等が生じたとしても交換機側のプログラムの修正が不要となる。これにより新しいサービスを短期間で提供することや、サービス内容をカスタマ毎に設定することが容易となる。

高度インテリジェントネットワーク（高度 IN）においても、従来のインテリジェントネットワークと同様に、伝達レイヤと高機能レイヤの二階層で構成され、サービス制御ノードと交換機は共通線信号網で接続される。高機能レイヤはサービス制御ノード、サービス管理ノード、およびオペレーションノードから構成される⁽¹⁵⁾。サービス制御ノードは交換機と接続され、電話呼の接続処理に直接係わるノードであり、リアルタイム性の高いデータベースアクセス機能が要求される。サービス管理ノード（SMS）はカスタマの情報や、そのカスタマが受けるサービス内容等をデータベースとして管理し、サービスオーダの受け付けや、サービス内容の変更等を行う。また、SMSで管理する情報のうち呼処理に必要な情報のみを SCP へダウンロードする⁽¹⁶⁾。オペレーションノードは運転保守を行うノードであり、多数の SCP の障害処理や、ソフトウェアの変更に伴うファイル更新処理等を行う⁽¹⁷⁾。

サービス制御ノードの実現技術としては、ハードウェアの構成技術、ソフトウェアによる制御技術があるが、本論文では、サービス制御ノードのハードウェア構成技術を対象とする。本研究の位置付けを図 1.2 に示す。

1.3 要求条件と課題

1.3.1 サービス制御ノード実現に向けた要求条件

サービスの多様化やサービスのカスタマイズ化に柔軟に対応するため、サービス制御ノードに要求される条件は以下の通りである。

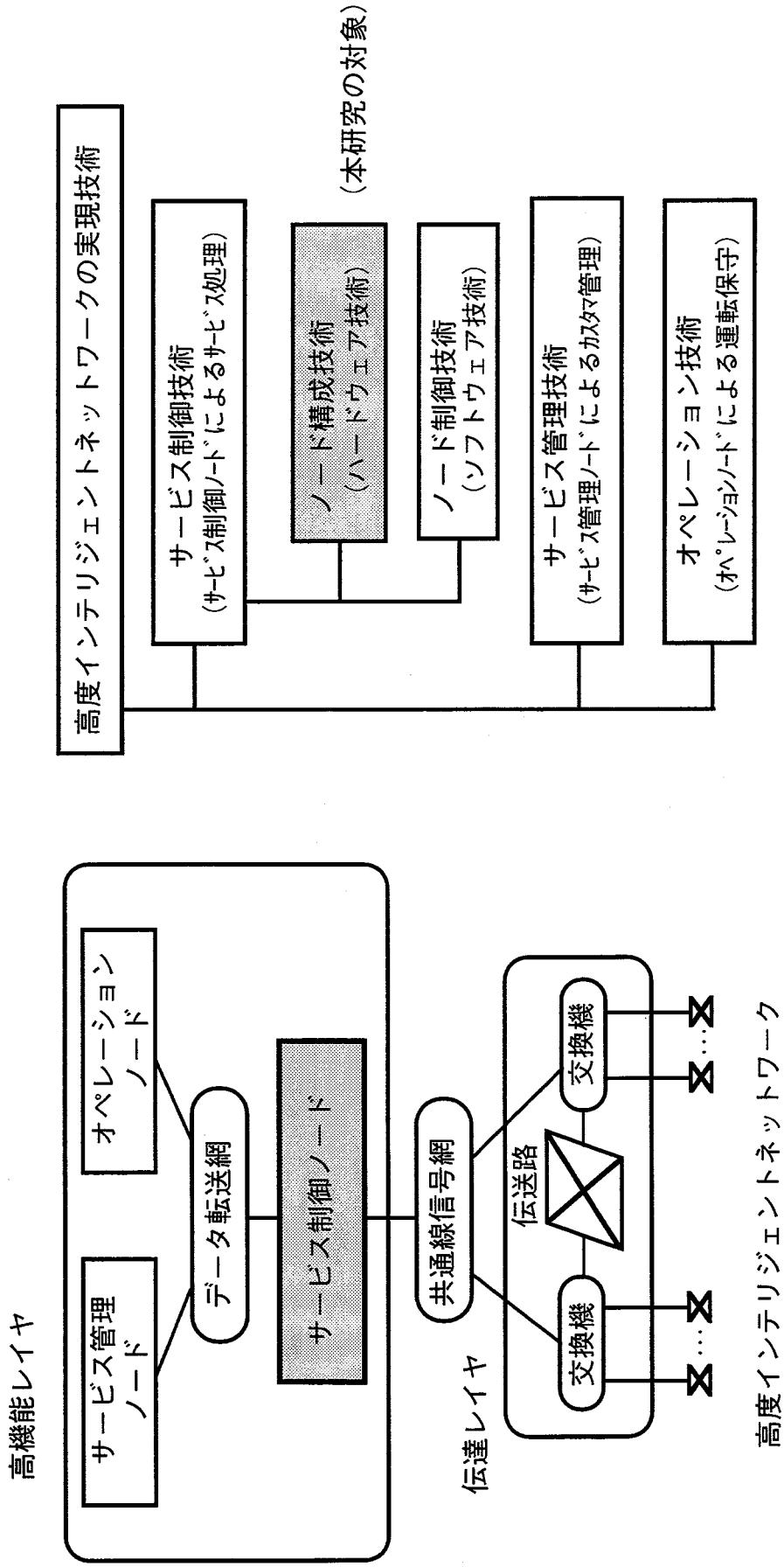


図1.2 研究の位置付け

(1) 需要に応じた柔軟な構成の実現

高度INでは、多様な内容のサービスが展開されるだけでなく、ノード規模は小規模から大規模まで幅広く分布する。またサービス毎の規模も事前に予測することが難しく、状況に応じ、即座に処理能力が拡大されなければならない。このため、分散処理構成による拡張性の高いシステムの実現が重要となる。

(2) データ管理機能の高度化

呼処理の過程で多様なデータを管理する高度INでは、今後、サービスの高度化、サービス種類の増加等によりデータの管理は一層複雑となる。これに対応するためにはデータベース技術の適用が必要となる。銀行のオンラインシステムやクレジット照会システム等のオンライントランザクション処理（OLTP：On-Line Transaction Processing）システムの分野ではデータベース技術は既に用いられているが、応答時間条件の厳しい高度INに導入するに当たっては、高速化のためデータベースの全てを主メモリや半導体ディスク装置等、半導体メモリ素子を用いた記憶媒体に格納することが必要となる。OLTPシステムにおいても半導体メモリ素子の高集積化の進展と相まって、データベースの全てを主メモリ等に格納する方法の研究は活発に行われているが、分散処理時の望ましい構成は明らかにされていない。

(3) 信頼されるサービスの実現

高度INサービスは伝達レイヤと高機能レイヤにより実現され、両レイヤが同時に機能するため、通常の電話サービスに比べて、一般的にアベイラビリティは低下する。しかし、その低下度合いは、ユーザからみて気づかない程度とすることが必須である。また、大規模なシステムになるほど、システムダウン時の影響が広く大きくなるため、一層の高い信頼性が要求される。一般的に、分散構成下では冗長化によりアベイラビリティを向上させることは容易であるが、同時にコストアップを招きやすい。このため、経済的に高信頼化をはかることが重要となる。

さらに、高度INでは、時々刻々変化するカスタマのデータをもとにサービスが行

われるため、データの喪失、破壊を防ぐこと、すなわち、ハードウェア障害、ソフトウェアのバグに対するデータの回復を保証することが併せて重要となる。

本論文では以上の要求に対応するため、分散処理化、データベース化、高信頼化の各方面からアプローチし、拡張性に優れ、信頼性の高いサービス制御ノードを経済的に実現するハードウェアシステム構成技術の確立を目的としている。本研究へのアプローチ方法を図1.3に示す。

1.3.2 サービス制御ノード構成の課題

望ましいサービス制御ノードを実現するための主要な課題は以下の通りである。

(1) データベースの分散構成技術

効率の良い分散型のサービス制御ノードを実現する上では、どのようにデータベースを配置、分散するかが重要な課題となる。なかでもデータベースを格納する記憶媒体や、複数のモジュールからのデータベースの共用形態等は、応答時間、処理の効率、さらにはデータの信頼性に影響するため、これらの選択は大きな課題となる。また、カスタマデータを分散する場合、それぞれのデータへのトラヒックのばらつきに起因してモジュール間の負荷の偏りが発生する。分散処理システムの方式設計や設備設計を効率的に行う上で、この偏りの評価や、偏りを平準化しモジュールの稼働率を高めることが課題となる。さらに、複数のモジュールを相互に結合した分散処理システムでは、結合機構がシステム全体の信頼性や性能を決定することになる。このため、結合機構の高信頼化、高性能化が重要な課題となる。

(2) データベースの高信頼化技術

データベースを効率よく高速に処理するためには、データベースの全てを主メモリ上に常駐するメモリデータベースの適用が有効である。しかし、メモリデータベース

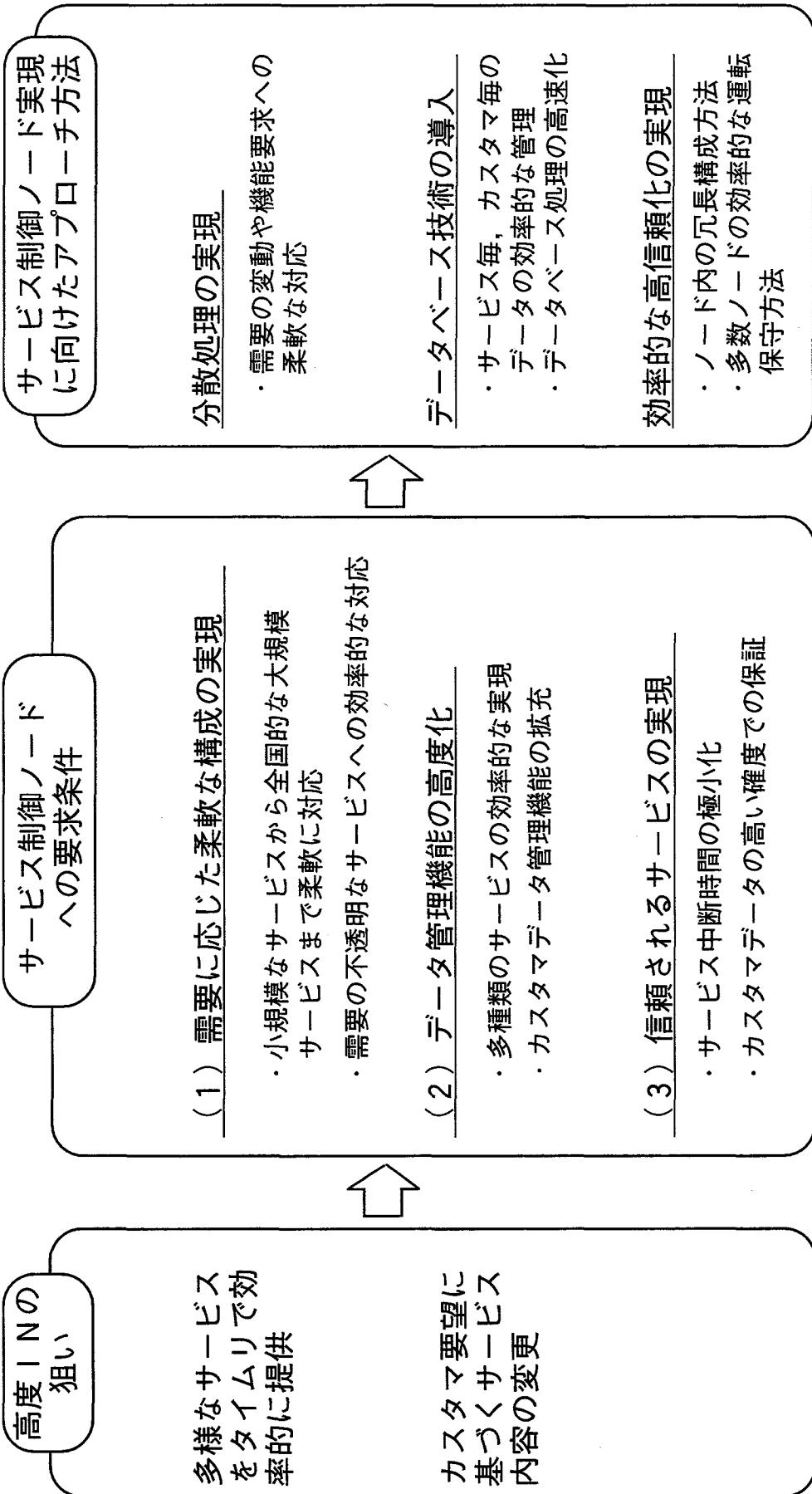


図 1.3 研究へのアプローチ方法

はソフトウェアのバグ、ハードウェア障害及び電源障害により、主メモリ上のデータベースが破壊、喪失されやすいという欠点を併せ持っている。このため、あらかじめデータベースの回復に必要な情報を障害に強い記憶装置上に格納しておき、障害時にはこの記憶装置上の情報をもとに障害直前のデータベースを効率よく復元する技術が要求される。

（3）高信頼化運転保守技術

全国的な高度INサービスを展開するためには、ネットワーク内に多数のサービス制御ノードを広域的に配置する必要がある。また、地震、水害等の大規模災害に対応するためには地理的に離れた地点に設置されたノード間での相互バックアップが必須となる。このように相互にバックアップされたサービス制御ノードに要求される信頼性を確保した上で、多数のノードの運転保守を可能な限り効率化、経済化する方法が要求される。

本研究における課題の概要を図1.4に示す。

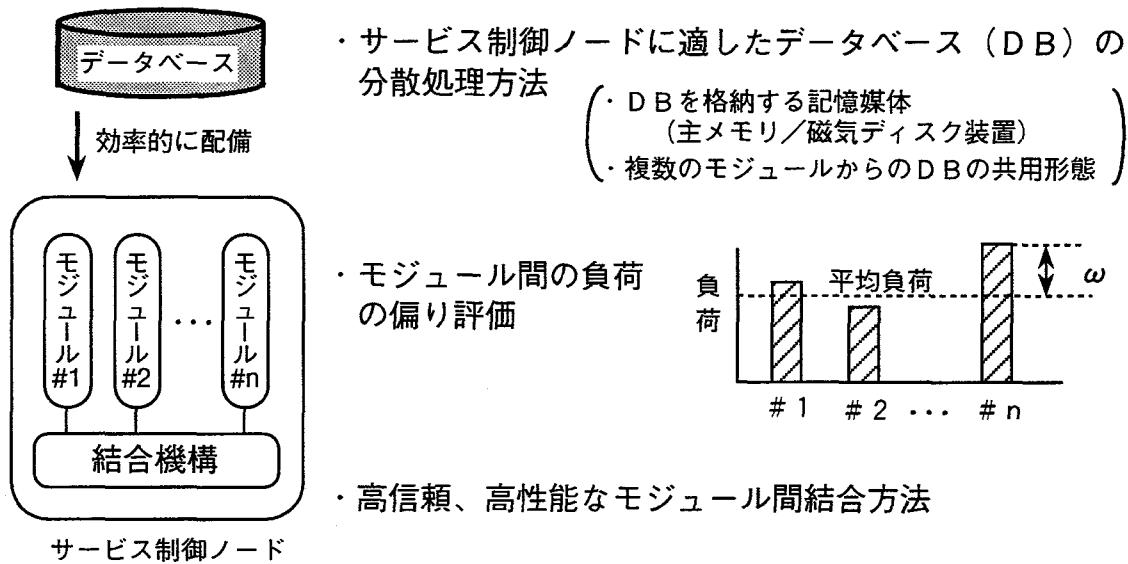
1.4 従来研究と本論文の関係

本節では、本論文で扱う課題に対して、従来の研究との関係を示す。

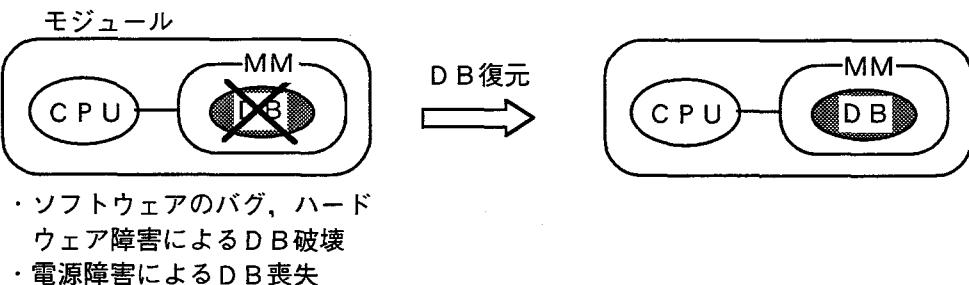
1.4.1 データベースの分散構成技術との関係

多数のモジュールによるデータベースの分散構成を効率的に実現するためには、高性能で高信頼なデータベースの分散処理方法や、モジュール間の負荷の偏りの評価方法、高信頼で拡張性の優れたモジュール間の結合方法を明らかにすることが重要となる。それについて、従来技術との関係について述べる。

(1) データベースの分散構成技術



(2) データベースの高信頼化技術



(3) 高信頼化運転保守技術

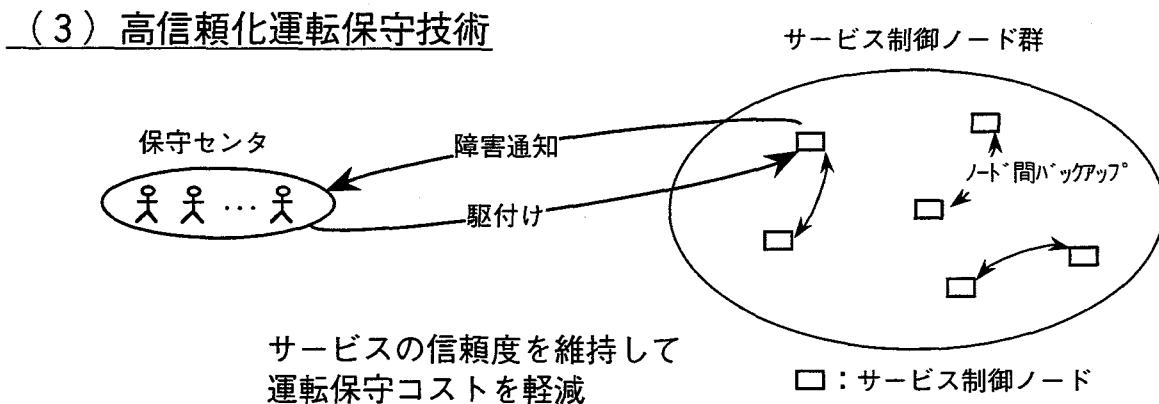


図1.4 課題の概要

(1) データベースの分散処理技術との関係

高度 IN用 SCP の実現法について、トラヒックの増大に対しては負荷分散により対応する構成の提案⁽¹⁸⁾ ⁽¹⁹⁾ はあるが、機能分散と負荷分散を統合した分散構成の提案や、分散処理を採用するまでの課題とその対応策については明らかにされていない。分散構成の実現に当たっては、高度 INに適したデータベースの分散処理の方法、高信頼化構成法を明らかにすることが主要な課題となる。OLTP (On-Line Transaction Processing) システムの分野ではデータベースを磁気ディスク装置に配備し、複数のモジュールで分散処理する方法はすでに実現されている^{(20) ~ (22)}。しかし、トランザクション処理の高速化のため、データベースのすべてを主メモリに常駐したメモリデータベースについての分散処理の実現法に関する報告例はない。また、高信頼化に関しては、地震、火災等の災害に対してもサービス中断を生じないように、地理的に離れた地点でのバックアップ方法の提案⁽²³⁾ ⁽²⁴⁾ はあるが、ノード内を分散構成とし、さらにノード間で相互バックアップしたシステムの効率的な冗長化構成法についての提案はない。

本論文では、小容量でアクセス頻度（単位時間当たりのアクセス回数）が極めて高いこと等、サービス制御ノードのデータベースの特性を明らかにし、それに基づいて、複数のモジュールからのデータベースの共用の可否、主記憶／半導体ディスク装置いずれかの格納媒体の選択の 2 点に注目して構成案を設定した。これら案について、モジュール間の負荷の偏り、データベース処理でのダイナミックステップ数に対する評価を加え、コスト、性能面から総合的に比較を行った。結果をサービス制御ノードに適した分散型メモリデータベース構成として提案する。また、モジュール間バックアップとノード間バックアップを組み合わせた効率的な冗長構成法を提案する。

(2) モジュール間の負荷の偏り評価技術との関係

高度 INサービスでは、カスタマ毎のトラヒックが大きくばらつく。トラヒックの高いカスタマとトラヒックの低いカスタマとでは、数千倍から数万倍のトラヒックの差がある。また、あるカスタマへのトランザクションは、当該カスタマのサービスデ

ータを用いて処理され、トランザクションとサービスデータが結びついている。このようにトラヒックのばらつきの大きいサービスデータを多数のモジュールに分散配備して、トランザクションを分散処理する場合、トランザクション処理の実行はサービスデータにアクセス可能なモジュールに限定されるため、モジュール間の負荷の偏りが問題となる。

一般に、分散処理システムの負荷の偏りに関する研究は、モジュールに配置されるカスタマ毎のサービスデータとトランザクション処理の結びつきの有無によって大きく2つに分類できる。サービスデータとトランザクション処理が結びつかない例としてパケット交換などの通信処理や音声、画像などのメディア変換処理があり、1つのトランザクションの処理量（ダイナミックステップ数）のばらつきは大きいが、どのモジュールでも処理が可能である。このようなシステムでは、モジュールにいかに効率よくトランザクションを割振るかが大きな課題であり、これに関しては、国内外で種々研究がなされている^{(25) (26)}。一方、サービスデータとトランザクション処理が結びつく例として、高度INシステムや、銀行のオンラインシステム、クレジット照会システムのような多くのOLTPシステムがあり、カスタマ毎のサービスデータへのトラヒックのばらつきが大きく、トランザクション処理の実行はサービスデータにアクセス可能なモジュールに限定される。従来から、サービスデータを磁気ディスク装置などの外部記憶装置上に集中して格納し、各モジュールで共用する方法がある⁽²⁷⁾。この場合は、どのモジュールも、全てのサービスデータに対するトランザクション処理が可能であり、モジュール間の負荷の偏りは生じない。しかし、近年、メモリ素子の高集積化を背景にトランザクション処理の高速化やスループット向上のため、サービスデータからなるデータベースの全てを主メモリ上に常駐するメモリデータベースが採られるようになってきた^{(28) (29)}。メモリデータベースではサービスデータの集中による共用は困難である。多数のモジュールでデータベースを非共用とした場合にはモジュール間の負荷の偏りが生ずるが、これに関する研究はほとんどなされていない。

本論文では、以上のようにトラヒックのばらつきの大きいサービスデータを多数の

モジュールに分散配備し、トランザクションを分散処理するシステムを対象として、カスタマ毎のサービスデータに対するトラヒック比が数千から数万倍と大きくばらつく場合のトラヒック量の近似方法を提案する。これに基づき、複数のモジュールにサービスデータを分散配備した場合の、モジュール間の負荷の偏りの評価方法、上限値の推定方法を明らかにする。

(3) モジュール間の結合技術との関係

高度 IN の SCP では、小規模なサービスから全国的な大規模サービスまで柔軟に対応することや需要の不透明なサービスへ効率的に対応することが要求される。このような需要の変動にタイムリに対応可能とするためには、多数のモジュールからなる大規模な分散構成により SCP を実現することが有効である⁽⁵⁾。分散処理システムを拡張性良く高信頼に実現するためには、モジュール間を接続する結合機構のデータ転送能力に制約が生ぜず高い信頼性が達成されることが要求される。また、結合機構の高信頼化に加えて、モジュール個々には高い通信処理能力が要求される。トランザクション処理の場合は、モジュール間で送受されるデータ長は数 100 バイトと短いが、送受を高頻度で行う必要がある。これを実現するため、通信処理を専用に行うハードウェア（以下、通信制御チャネルと記す）を適用し、この通信制御チャネルの単位時間当たり処理可能な通信回数を向上させることが重要となる。従来から分散処理システムのモジュール間の結合方法として、独自の LAN を用いたものや⁽²²⁾、業界標準のイーサネット、FDDI（Fiber Distributed Data Interface）を用いたものが OLT P システム等で提案されている。これらのシステムで用いられている LAN の最大データ転送速度は 10 ~ 100 Mbps 程度であり、今後のプロセッサ性能の向上や接続するモジュール数の増加を考えると必ずしも十分とは言えない。モジュールの追加により柔軟に処理能力を拡張可能とするためには、高いデータ転送能力を有するスイッチを適用する方法が有効となる。しかし、スイッチを適用する場合は、その信頼性がノード全体の信頼性を決定することとなり、スイッチの高信頼化が必須となる。パソコン等を相互に接続した企業内の LAN を構築する場合のバックボーンとして、

既にスイッチが適用されている⁽³⁰⁾が、高い信頼性が要求されるオンライントランザクション処理システムへの適用例は少ない。

本論文では、分散構成を採るSCP内のモジュール間を結合する方式として、十分なデータ転送能力を確保することからスイッチ方式とし、モジュールとのインターフェースはATMを適用する分散処理システムを対象として、結合機構の高信頼化構成法や、結合機構に接続されるモジュールの通信処理能力を向上させる制御方式について提案する。

1.4.2 データベースの高信頼化技術との関係

サービス制御ノードのデータベースは容量が比較的小さく、アクセス頻度が高い特徴を有し、トランザクション処理性能の向上のためにはデータベースの全てを主メモリ上に常駐するメモリデータベースの適用が有効である。しかし、メモリデータベースは、ソフトウェアのバグ、ハードウェア障害および電源障害により、主メモリ上のデータベースが破壊、喪失されやすいという欠点を併せ持っている。このため、あらかじめデータベースの復元に必要な情報を障害に強い記憶装置上に格納しておき、障害時にはこの記憶装置上の情報をもとに障害直前のデータベースを復元する方式を併せて採用する必要がある。この実現のため、Lehman、Eichらは磁気ディスク装置上にデータベースの更新ログ(LOG)とチェックポイント時点のデータベース(CPDB)を取得し、障害時にはLOGとCPDBから障害直前のデータベースを復元する方式を提案している^{(31)～(33)}。磁気ディスク装置の入出力ボトルネックを解消するため、不揮発化したバッファメモリにLOGをある程度バッファリングし、一括して磁気ディスク装置に書込む方式を探っている。しかし、この方式では電源障害時にバッファリングされているLOGを保持することは可能でも、ソフトウェアのバグによりバッファメモリ上のLOGが破壊される可能性がある。一方、高倉らは磁気ディスク装置を用いず、読み出しのデータ転送速度がより高速で不揮発性を持つフラッシュメモリ上にLOGとCPDBを取得し、データベースの復元を高速に行う

方式を提案している⁽³⁴⁾。主メモリを二重化するとともに、一時的に更新内容を保持するバッファを設けるなど主メモリを専用化した構造としている。このため、汎用的な装置の適用による経済性の追求やソフト制御の簡潔化が期待しにくくなるという問題がある。

本論文では、これらの要因を回避するため、汎用性の高い装置構成を前提に、LOGを主メモリでバッファリングせずトランザクション毎に外部の記憶装置に書込むこととしたメモリデータベースの復元方式（以下、リカバリ方式と記す）を検討する。この場合、外部の記憶装置の形態はLOGを単に記憶するパッシブな方式とLOGをプロセッサで処理可能なアクティブな方式に2大別される。前者を代表し高速な入出力が可能な半導体ディスク装置を適用する方式を設定し、後者については別のモジュールのプロセッサに付属する主メモリを適用する方式を代替案として設定する。2つの代替案について、障害発生からデータが回復するまでの時間と、データベースの単位容量当たりのトラヒック量に着目し、メモリデータベースのリカバリ方式の評価方法を明らかにする。これに基づき、サービス制御ノードのメモリデータベースのリカバリ方式として半導体ディスク装置を用いた方式を提案する。また、劣化判定機能を有する保守の容易なバッテリを用いた、半導体ディスク装置の効率的な不揮発化方法について提案する。

1.4.3 高信頼化運転保守技術との関係

高い信頼性が要求される交換機等では、従来からノード内で二重化する冗長構成が一般にとられている。しかし、データベースを用いて全国を対象にサービスを行う高度INのサービス制御ノードでは、ノード内の二重化のみでは、地震、火災等の大規模災害に対しては十分でなく、地理的に離れた2つの地点にノードを設置して、ノード相互でデータベース上の情報を含めてバックアップする冗長構成がとられるようになってきている⁽²³⁾。同様に、キャッシングシステムやクレジット照会システム等のオンライントランザクション処理システムにおいても、サービス内容の多様化、高

度化に伴って、ノード内で二重化し、さらに地理的に離れて設置された2つのノード間で相互バックアップする方式が採用されてきている^{(24) (35)}。

一方、システムの信頼性は、このような高信頼化のための冗長方式とならんで保守方式により大きく変動する。ノードに保守チームが常駐している場合は、障害検出と同時に修理に取り掛かることが可能である。しかし、ノード毎に保守チームを常駐させると保守コストの増大を招くため、通常、保守センタに保守チームを集中し、複数のノードを少数のチームで保守する方法が採られている⁽³⁶⁾。この場合、サービス停止につながるような障害では、常駐の保守チームにより緊急に駆付けて対応し、サービス停止に至らない軽度の障害では、ある程度時間的な余裕を持って平日勤務の保守チームにより対応する等、システムの緊急度合いに応じて駆付け時間を調整することが多い。渡辺は、二重化した交換機と宅内装置からなる加入者系を対象として、駆付け保守による修理遅延の影響について検討を行っている⁽³⁷⁾。しかし、ノード内で二重化し、さらに2つのノード間で相互バックアップを行ったシステムへの拡張や、システムの緊急度合いに応じて駆付け時間を二段階に分けた評価等は行っていない。一般に、二重化された單一ノードの場合は、経験的に、駆付けの緊急度合いはシステムダウンしているか否か等により判断される場合が多い。しかし、二重化された2つのノードが更に相互にバックアップされる場合には、駆付け要求が発生する状態が多く、どのようなときに緊急駆付けとするかの判断が重要となる。地理的に離れたノード間での相互バックアップは大規模災害への対応が本来の目的であるが、通常の装置障害に対しては、冗長度が上がるため、高い信頼性が容易に達成可能となる。許容される範囲での信頼性低下のもとで、駆付け時間の長時間化を図ることにより、保守の効率化とシステムの信頼性をバランス良く実現することが期待できる。

本論文では、データベースを持ち、ノード内は二重化されたモジュールが相互にバックアップし、さらに地理的に離れた地点に配置された2つのノードが相互にバックアップする高度INにおけるサービス制御ノードを対象として、保守の効率とシステム信頼性をバランスよく実現できる駆付け保守の方法を提案する。具体的には、システム信頼度の時間的な変化に着目して、駆付けの緊急度合いをシステム内の障害装置

数や、保守チームが到着しているか否かの状態によりレベル分けし、このレベルにより駆付け保守の方法を設定する。これに基づいて、システムの信頼性や駆付け回数を総合的に評価し、非緊急駆付けについては複数の障害をスケジュール化可能で、直ちに駆付けなければならない緊急駆付け回数の少ない効率的な保守方法を明らかにする。

1.5 論文の構成と各章の概要

本論文は、上記で述べたように高度インテリジェントネットワークにおけるサービス制御ノードのハードウェアシステム構成技術に関する研究成果をまとめたものであり、以下の8章から構成される。技術的課題と論文の構成との関係を図1.5に示す。

第1章では、本研究の背景と目的を示し、それに関連する研究状況、問題の所在を明らかにする。

第2章では、本研究で対象とするサービス制御ノードのシステム構成を検討する上で重要な、データベースの配備形態や、呼の接続処理、および複数のサービス制御ノードによる多種類のサービスの実現方法を示す。

第3章では、サービス制御ノードのデータベースに適した分散型メモリデータベース構成を提案する。サービス制御ノードのデータベースは容量が小さく、アクセス頻度が極めて高い特徴を有し、オンライントランザクション処理システムで一般に用いられている磁気ディスク装置にデータベースを格納する方法は適用できない。データベースを多数のモジュールの主メモリに分散して格納し、トランザクションを効率的に分散処理する方法を明らかにする。分散型メモリデータベース構成により、高性能で経済的なサービス制御ノードを実現できる。

第4章では、データベースの分散配置とカスタマ毎のトラヒックのばらつきに起因して起こるモジュール間の負荷の偏りの評価方法や、負荷を平準化してモジュールの

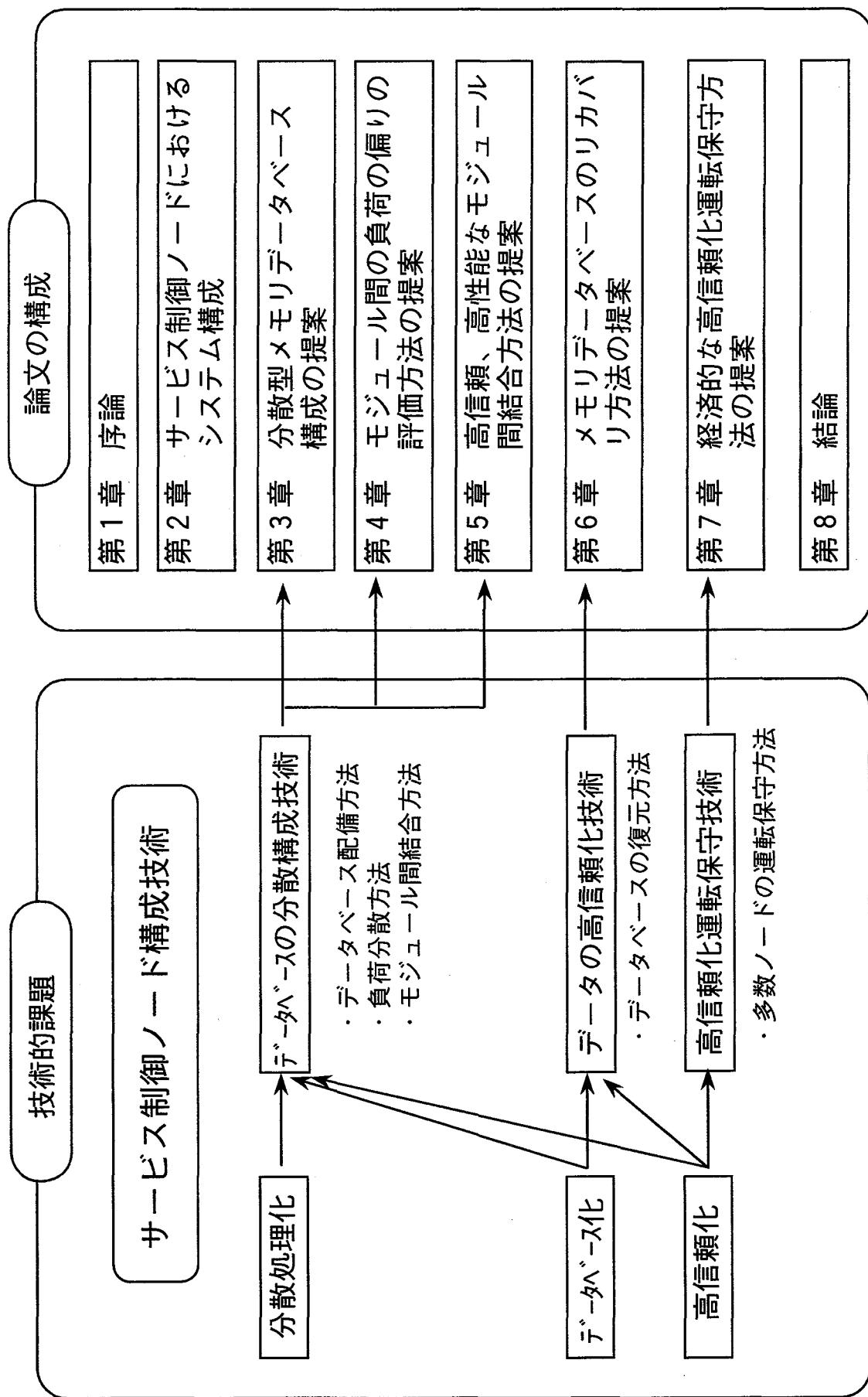


図 1.5 技術的課題と論文の構成

稼働率を高める方法を提案する。カスタマ毎のトラヒック量については、トラヒックの高いものからの順位により近似する方法を示し、これを用いてモジュール間の負荷の偏りや、偏りの上限値の推定方法を明らかにする。また、負荷の平準化については、データの配置先を示すディレクトリの検索処理とデータベース処理を組合わせることにより、動的に平準化する方法を明らかにする。負荷の偏り評価や平準化技術は、分散処理によるサービス制御ノードの方式設計、設備設計に有効である。

第5章では、多数のモジュールによる分散構成を効率的に実現するため、高信頼で、拡張性の良いモジュール間の結合方法について提案する。インターフェースとしてはATMを適用し、結合機構を構成する機能ブロック毎に二重化することによりモジュール数が増加した場合でも高い信頼性が確保できる結合機構の構成法を明らかにする。また、モジュールの通信処理能力を向上させる通信処理方式を明らかにする。これにより、高信頼なサービス制御ノードを拡張性良く実現できる。

第6章では、主メモリ上に格納されたデータベースがソフトウェアのバグ、ハードウェア障害、電源障害により破壊、喪失された場合、データベースを効率的に回復する方法を提案する。呼の処理と並行して、データベースの更新ログと数分毎のチェックポイントデータベースを半導体ディスク装置に取得し、主メモリ上のデータベースが破壊、喪失した場合には、この半導体ディスク装置に取得した情報から障害直前のデータベースを高速に回復する。これにより、メモリデータベースの高信頼化を、通常の呼処理への影響が少なく実現できる。

第7章では、広域的に分散設置された多数のサービス制御ノードの効率的な運転保守方法を提案する。地震、火災等の大規模災害に対応するためには異なる地域に設置された2つのノード間での相互バックアップが必要であるが、これに伴うサービス制御ノードの信頼度の向上を運転保守の軽減に活用する方法を明らかにする。これにより、多数のサービス制御ノード群を効率的に運転保守することが可能となる。

第8章では、本研究で得られた結論について総括する。

第2章 サービス制御ノードにおけるシステム構成

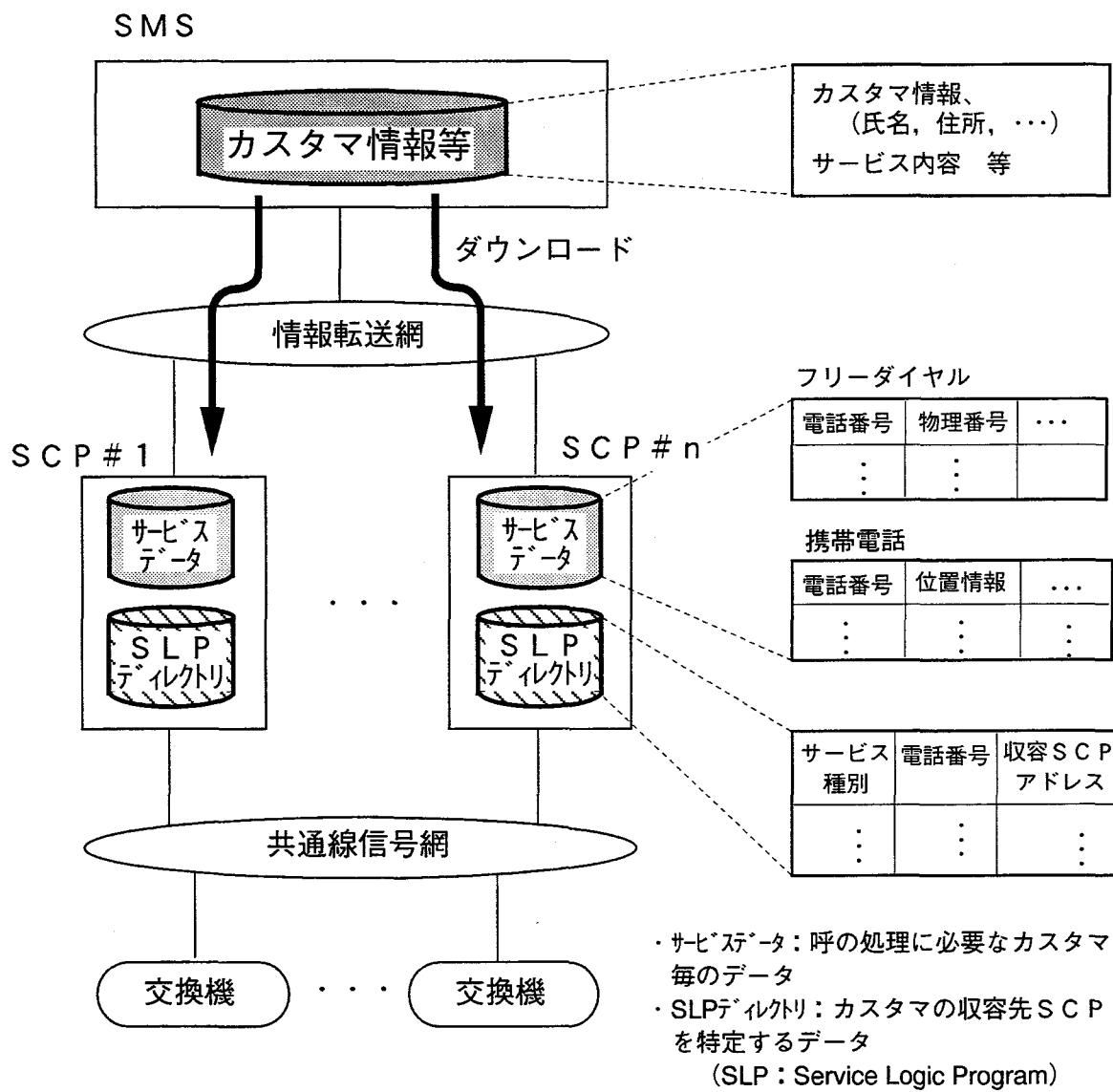
本章では、本研究で対象とするサービス制御ノードのシステム構成を考える上で重要な、データベースの配備、呼の接続処理、および今後のサービスの多様化に柔軟に対応可能とするための多種類のサービスの実現方法について述べる。

2.1 データベースの配備

高度インテリジェントネットワーク（高度IN）ではカスタマ毎に多様なデータベースを維持管理する必要がある。これらのデータベースはカスタマ管理に必要なデータベースと呼の処理に必要なデータベースに分類できる。前者は、カスタマの氏名や住所等のカスタマ情報やカスタマが受けるサービス条件等のデータベースからなる。一方、呼処理に必要なデータベースは電話番号や位置情報などからなる。高度INにおけるデータベースの配備形態を図2.1に示す⁽¹⁵⁾。

サービス管理ノード（SMS）では、カスタマの情報や、そのカスタマが受けるサービス内容等がデータベースとして管理される。ここで管理される情報のうち、実際の呼処理に必要な情報のみが、サービス制御ノード（SCP）にダウンロードされる。たとえば、フリーダイヤルの場合は、論理的な電話番号と物理番号の対などがデータベースとしてダウンロードされる。これらは、呼の接続サービスを行うために必要な情報であり、サービスデータと言われている。

SCP内に保持されるデータベースとしては、サービスデータの他にSLP（Service Logic Program）ディレクトリがある。SLPはカスタマが受けるサービス（たとえば、フリーダイヤルの場合、時間帯により接続先の電話機を変更するとか、発信者の位置により接続する電話機を変更する等のサービス）を実現するため、カスタマ毎に作成されるプログラムである。このSLPを書換えることにより、サービスのカ



- ・カスタマ管理に必要なデータベース（DB）はサービス管理ノード（SMS）に配備
- ・呼処理に必要なデータベースはサービス制御ノード（SCP）に配備

図2.1 高度インテリジェントネットワークにおけるデータベースの配備

スタマイズ化が容易に実現できる。SLPディレクトリとは、複数のSCPにカスタマ毎のSLPを分散収容した場合、どのカスタマのSLPがどのSCPに収容されているかを特定するためのデータベースである⁽⁵⁾。

2.2 呼の接続処理

電話機からの呼は一般の電話呼か高機能レイヤを利用する電話呼（以下、IN呼と記す）かを交換機で判別し、IN呼の場合は共通線信号網を介して処理要求がサービス制御ノードに送られる。サービス制御ノードはサービス毎、カスタマ毎にどのような内容のサービスを行うかのデータベースを持っており、この内容に従って交換機と協調して高機能なサービスを実現する。高度インテリジェントネットワークにおける呼処理の流れを図2.2に示す。

図2.2において、電話機から投入された電話番号は交換機で解読され、一般の電話呼かIN呼かが判別される。一般の電話呼の場合は、投入された電話番号に基づいて接続処理が行われる（図2.2の破線の①と⑤）。IN呼の場合、たとえば最初の電話番号が0120（フリーダイヤル）や0990（ダイヤルQ²）等で始まる呼の場合には、これらの電話番号の検出をトリガとして、あらかじめ決められたSCPに問い合わせがかけられる。この時、トリガとなる電話番号に加えて、後続の電話番号がSCPに送られる。SCPでは送られてきた電話番号をキーとしてデータベース処理を行い、論理番号から物理番号への変換を行う。フリーダイヤルの場合、カスタマによっては時間帯により接続先を変更したり、発信地点により着信先を振り分ける等の付加サービスを受けられている場合もあり、物理番号への変換に当たってはこのような条件も含めて行われる。次に、SCPはどのような処理を行うかの応答を交換機に返し、交換機はこの応答に基づいてIN呼の接続処理を行う（図2.2の実線の①～⑤）。

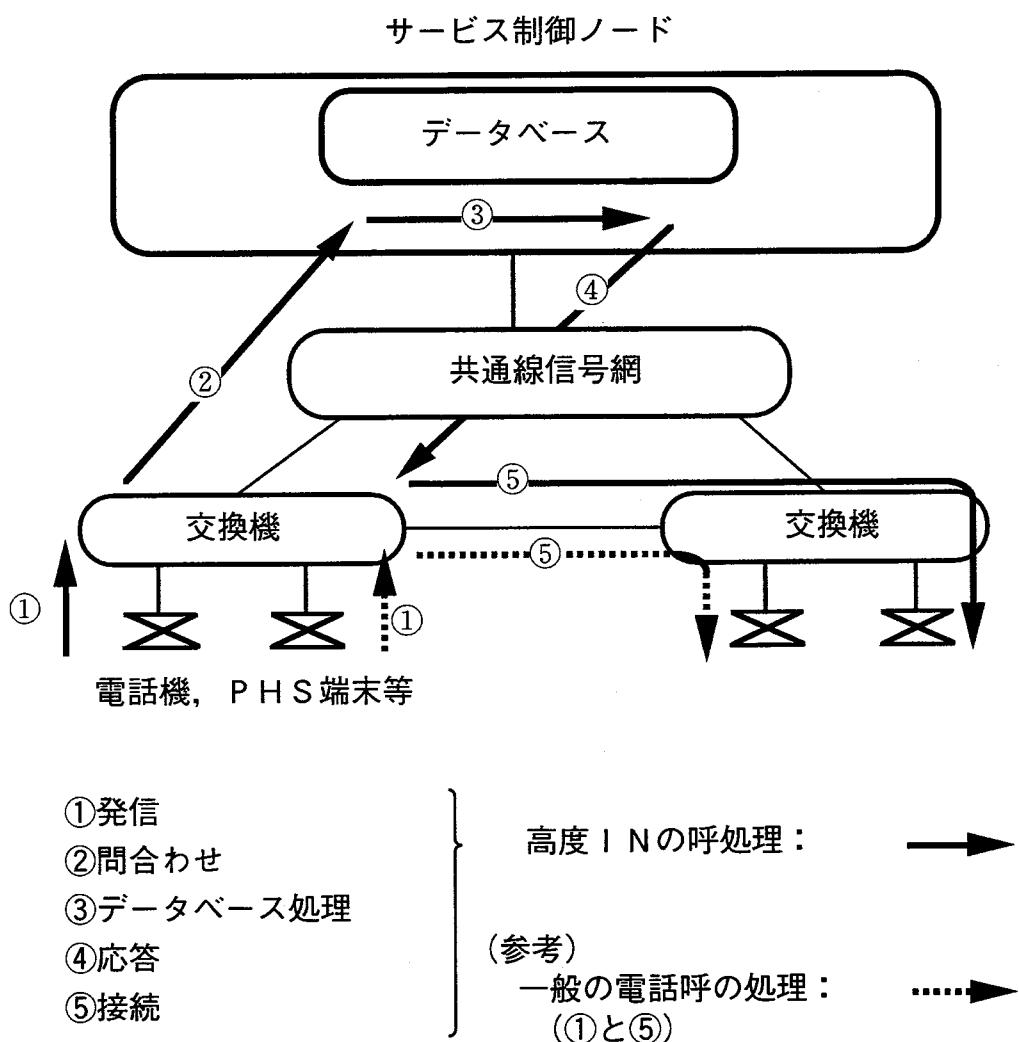
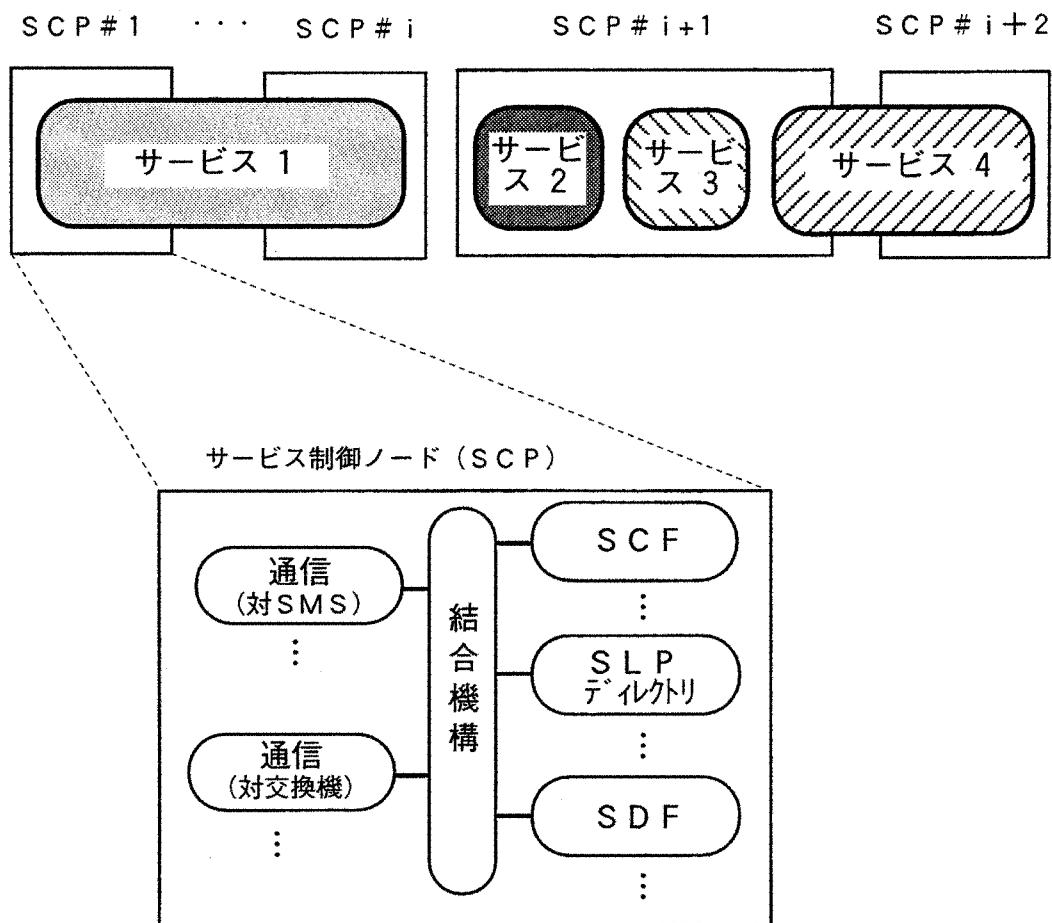


図2.2 高度インテリジェントネットワークにおける呼処理

2.3 多種類のサービスの実現方法

需要の異なる多種多様なサービスを効率的に実現するためには、大規模なサービスについては、複数のサービス制御ノードで分散することや、小規模なサービスについては、1つのサービス制御ノードで複数のサービスを実現することが必要となる。このためには、カスタマ毎のSLPやサービスデータを複数のSCPに分散収容し、カスタマ毎のトランザクションをSLPやサービスデータの配備されているSCPに割振ることが必要となる。トランザクションの振り分けは、前節で述べたSLPディレクトリにより行われる。複数のサービス制御ノードによる多種類のサービスの実現方法を図2.3に示す。

また、サービス開始後の需要の増加や、サービス内容の高度化に伴ってプロセッサの処理量が増加した場合に、柔軟に対応するためには、サービス制御ノードを複数のモジュールで構成し、モジュールの追加により対応できるノード構成が必要となる⁽⁵⁾。サービス制御ノードの分散構成方式については次章以下で考察する。



多種類のサービスを需要に応じて柔軟に分散

- ・大規模サービス → 複数ノードで分散
- ・小規模サービス → 1ノードで複数サービス
- ・サービス開始後の需要の増加、サービス内容の高度化
→ サービス制御ノードを複数のモジュールで構成し、モジュールの追加により対応

図 2.3 複数のサービス制御ノードによる多種類のサービスの実現方法

第3章 サービス制御ノードの分散構成技術

3.1 緒言

高度INのサービス制御ノード(SCP)では、幅広いレンジにわたる多様な内容のサービスを、小規模から大規模まで、機能間の分離性を良好に保ちながら処理できること、多数のカスタマの情報から構成されるデータベースを効率よく高速に処理できることが重要となる。また、障害に対しても、これらデータベースの内容が十分な確度で保証されるとともに、サービスとしては高い可用性(アベイラビリティ)の達成が必須となる。

本章では、このような要求条件に効率的に対応可能な構成として分散構成に注目し、高度INのSCPを多数のモジュールからなる大規模な分散構成により実現する方法を提案する^{(5) (38) ~ (40)}。基本的な分散構成としては、ITU-Tで標準化⁽⁴¹⁾されている、サービス制御機能(SCF)とサービスデータ機能(SDF)を別モジュールで実現し、さらにこれら機能の負荷を複数のモジュールで処理可能とした機能分散、負荷分散の統合形態としている。分散構成の実現に当たっては、高度INに適したデータベースの分散配備の方法を明らかにすることが主要な課題となる。はじめに、小容量でアクセス頻度が極めて高いこと等、高度IN用データベースの特性を明らかにし、それに基づいて、複数のモジュールからのデータベースの共用の可否、主記憶／半導体ディスク装置いずれかの格納媒体の選択の2点に注目して構成案を設定した。これら案について、モジュール間の負荷の偏り、データベース処理でのダイナミックステップ数に対する評価を加え、コスト、性能面から総合的に比較を行った。結果を高度INに適した分散型メモリデータベース構成として提案する。データベースのリカバリについては、バンキングシステム等のOLTP(On-Line Transaction Processing)系システムで、データ内容の厳密な復元を重

要視することとは対照的に、通信システムとしての高度INでは、短いリカバリ時間の実現が最も重要となる。これを考慮し、データ毎の特性に応じたりカバリレベルを設定するとともに、数分毎のチェックポイントデータベースの取得と小容量の更新ログデータの組み合わせによる、効率的なリカバリ方法を明らかにする。また、一般に、分散構成下では、冗長化による高いアベイラビリティの実現は容易であるが、同時にコストアップを招きやすい。ノード内のモジュール間バックアップとノード間にわたるバックアップを効率的に組み合わせ、高いアベイラビリティを経済的に実現する冗長構成法を提案する。

3.2 ノード構成条件と機能配備

3.2.1 ノード構成条件

本章では、パーソナル通信など全国規模のサービス、カスタマオリエンティッドな多様なサービスへの適用を考え、従来ノード⁽¹⁰⁾に比べ大規模化し、100万カスタマ程度を収容可能なノードを想定する。また、ノードにかかるトラヒックは2000呼／秒程度を想定する。

3.2.2 機能配備

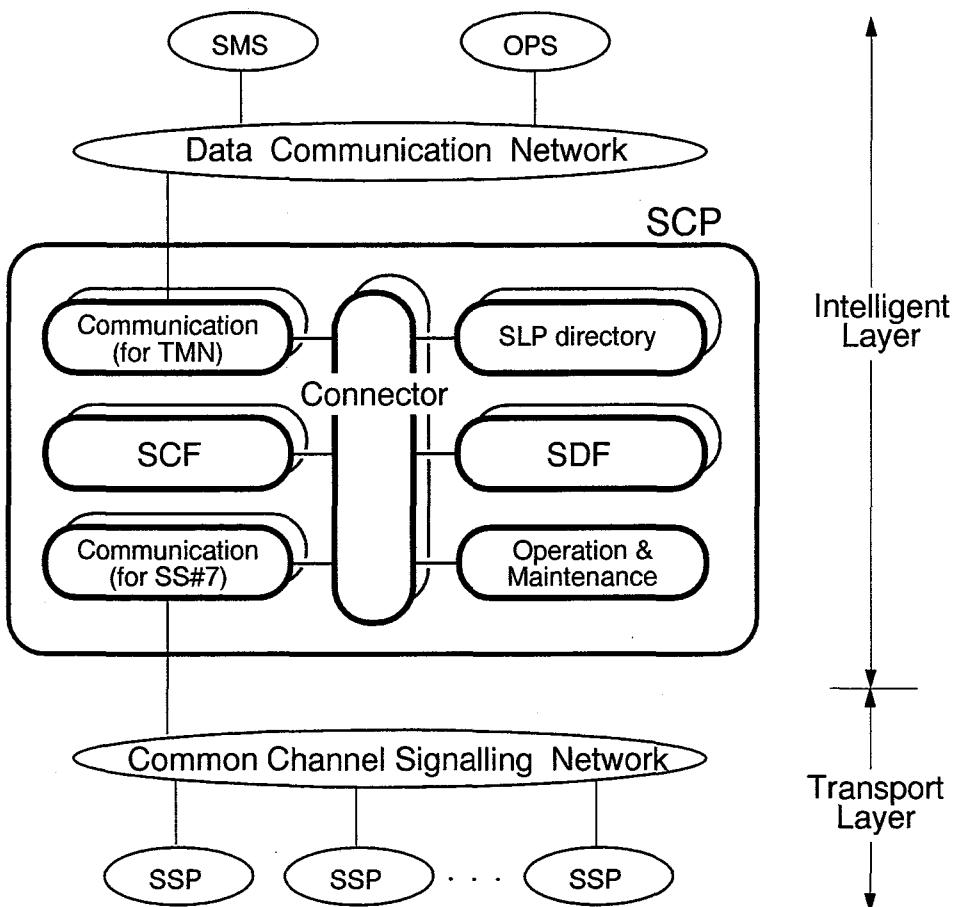
高度INに適用するSCPを実用的なシステムとして構成するためには、以下の機能が必要である。

- (1) サービス制御機能(SCF)：カスタマ毎に作成されたサービスロジックプログラム(SLP)によりサービスを実行する機能。

- (2) サービスデータ機能（SDF）：サービスを実行するために必要なサービスデータを保持し、SLPからの問い合わせに従ってサービスデータへのアクセスを行う機能。
- (3) SLPディレクトリ検索機能：ネットワーク内に複数のSCPを配備し、SLPを分散収容した場合、SLPがどのSCPに配備されているかを決定する機能。
- (4) 他ノードとの通信機能：共通線信号網を介した交換機との通信、パケット網を介したサービス管理システム（SMS）との通信機能。
- (5) 運転・保守機能：上位ノードからの指示に基づく運転操作やノード全体の障害管理を行う機能。

これら（1）から（5）の機能のモジュールへの配備に関しては、機能間の分離性がよく、負荷変動に対して機能単位に柔軟に対応できること、実現する機能に最適化した技術を適用できることから、機能毎に別々のモジュールに配備し、負荷の増大に対しては機能単位で負荷分散することとした。共通線信号網及びパケット網を介した他ノードとの通信機能については、各網の将来の発展形態が異なることから、それぞれ別モジュールで分担することとした。SCPのネットワーク内の位置付け及びSCP内の機能配備を図3.1に示す。

なお、本章では、SLPディレクトリはSCPに配備した。交換機側はカスタマによらず、前もって定めたSCPにアクセスする。そのSCPでSLPディレクトリを検索し、当該カスタマのSLPが配備されているSCPを知る。その結果に基づき、呼を当該SCPへ転送する。本方式のほかに、交換機側にSLPディレクトリを配備する方法も考えられるが、SLPの追加毎に数100台に及ぶ交換機上でのSLPディレクトリの変更が必要となる。ネットワーク管理上のオーバヘッド



SMS: Service Management System
 OPS: Operations System
 SSP: Service Switching Point

SCF : Service Logic Function
 SDF : Service Data Function
 SLP : Service Logic Program

図 3.1 SCP の機能配備

が極めて大きいため対象外とした。

S C F の負荷分散については、S C F を分担する全てのモジュールに S L P を重複して配備した。これにより、S C F を分担するどのモジュールでも S L P を実行可能とし、特定のカスタマに負荷が集中した場合においても均等な負荷分散を可能とした。他ノードとの通信機能の負荷分散については、モジュール内の回線増設とモジュール増設を組み合わせて対応することとした。S D F 、S L P ディレクトリ検索のデータベース処理の負荷分散については、負荷の偏りや、どのような記憶媒体にデータベースを格納するか等の問題を併せて検討するため、次節で考察する。なお、運転・保守機能については、ノードが大規模化されても、1台のモジュールで対応可能との見通しを得たため、負荷分散は行わないものとした。

3.3 データベースの分散構成

高度 I N におけるデータベースの特性について考察し、次にその特性にあった複数モジュールへのデータベースの配備、分散方法を明らかにする。

3.3.1 サービス制御ノードにおけるデータベースの特性

S L P ディレクトリとサービスデータから構成されるデータベースの容量、アクセス特性等を表3.1に示す。各データベースの主要な特徴は以下の通りである。

(1) ノード当たりのデータベース容量は数百M B から 1 G B であり、O L T P 系システム、たとえばバンキングシステムで必要とされるデータベース容量に比べ、おおよそ 1 / 1 0 ~ 1 / 1 0 0 と小さい。一方、単位時間、単位容量当たりの I / O アクセス回数（以下、I / O アクセス密度と記す）は 1 0 0 0 ~ 1 0 0 0 0 回 / 秒 / G B であり、O L T P 系システムに比べて 1 0 0

表3.1 サービス制御ノードにおけるデータベースの特性

	用途	データベース容量	データベースアクセス頻度	入出力アクセス頻度	入出力アクセス密度
S L P ディレクトリ (参照系データ)	収容先 S C P を特定	~1GB	2 k回/s	4 k回/s	10 k回/s/GB
サービスデータ (更新系データ)	呼処理	~1GB	5~10k回/s	10~20k回/s	10~20 k回/s/GB

- (注) ・条件：ノード当たりのトラヒックは2000呼／秒を想定
 　　(100万カスタマ収容時)
 　・入出力アクセス頻度：単位時間当たりの入出力アクセス回数
 　　(データベースが外部記憶装置に格納されている場合)
 　・磁気ディスク装置の入出力アクセス密度=数10回/s/GB

～1000倍大きい。すなわち、高度IN用データベースでは少ない容量のデータベースに極めて高頻度のアクセスがある。

- (2) 高度INのカスタマ毎のレコード（データベースの最小操作単位）への単位時間当たりのアクセス回数（以下、アクセス頻度と記す）はサービス種別にもよるが、一般的に大きくばらつく。少数のカスタマが高いアクセス頻度を有し、かつ全体のアクセス頻度の相当程度を占める。アクセス頻度の高いカスタマとアクセス頻度の低いカスタマのアクセス頻度の比は1000倍から10000倍と大きい。
- (3) データベースへのアクセスは、電話機から入力されたダイヤル番号（フリーダイヤルサービスの場合、「0120」の後に続く6桁の論理番号）等をユニークキー（データベースで、あるデータ項目に着目した時、同一の内容が存在せず、カスタマ毎のレコードを一意に認識できる項目）とし、そのキーに一致する1レコードをアクセスする形態がほとんどである。
- (4) データベースへのアクセスには、呼処理のアクセス（交換機からのアクセス）とサービス管理のアクセス（新規カスタマの追加、サービス条件等の変更に伴うSMSからのアクセス）がある。サービス管理のアクセスでは情報がSMSよりダウンロードされデータベースの書き換えが行われる。呼処理のアクセスでは、SLPディレクトリに対しては読み出しのみが行われ、サービスデータに対してはサービス内容に応じて読み出し／書き込みが行われる。呼処理のアクセスに比べ、サービス管理のアクセスは2～3桁アクセス頻度が少なく、各データベースへのアクセス特性としては呼処理のアクセスのみを考えれば十分である。従って、SLPディレクトリは読み出し専用データベース（以下、参照系データベースと記す）、サービスデータは読み出し／書き込みデータベース（以下、更新系データベースと記す）として扱う。

3.3.2 データベースの格納媒体

データベース格納媒体として、磁気ディスク装置、半導体ディスク装置、モジュール内プロセッサの主記憶（以下、メモリと記す）がある。本章では、以下の理由により、メモリと半導体ディスク装置を対象とした。

- (1) 一般にOLTP系システム⁽²⁰⁾⁽²¹⁾に適用されている磁気ディスク装置、半導体ディスク装置の単位時間当たり処理可能なI/O回数はそれぞれ40～60回／秒、1000～2000回／秒である。高度INでのノード当たり要求能力が10000～20000回／秒（表3.1）であるため、最大1GB程度のデータベースを格納するために、磁気ディスク装置は数100台を並べる必要があり適用性はない。一方、半導体ディスク装置は10台程度にI/Oアクセスを分散させることにより適用の可能性がある。
- (2) データベースの容量が比較的小さく、アクセス頻度の高い分野のデータベースでは、データベース全体をメモリに常駐させるメモリデータベースが有効であると言われている⁽²⁸⁾⁽²⁹⁾。高度INのデータベースは最大でも1GB程度と容量が小さく、アクセス頻度が非常に高いためメモリデータベースの適用性がある。ただし、ハード、ソフト障害時にデータベースが破壊されやすいため、これに対応可能なりカバリ方式を併せて採用する必要がある。

3.3.3 データベース分散方法

3.3.3.1 複数モジュールへのデータベースの配備法

複数モジュールからのデータベース(DB)の共用の可否、格納媒体の違い（メモリ、半導体ディスク装置(SD)）に注目し、以下の3方式を構成上の代替案と

して比較評価する。

【分割方式 (Partitioned Method) : DB 非共用 / メモリ格納】 データベースをレコード単位にモジュール台数分に分割し、各モジュールのメモリに格納する方式。カスタマ毎のレコードへのアクセス頻度の偏りに起因してモジュール間の負荷の偏りが生じやすい。また、高度 IN のデータベース処理はユニークキーによる 1 レコードアクセスがほとんどで、モジュール間に跨った処理は極めて少ない。このため、本方式のようにデータベースを分割配備したことによるデータベース処理のオーバヘッドは通常無視できる。しかし、データベースアクセスに先だって、どのレコードがどのモジュールに格納されているかを判別するための処理が必要となる。

【多重方式 (Duplicated Method) : DB 共用 / メモリ格納】 全てのデータベースを各モジュールのメモリに重複して格納する方式。メモリを重複して持つため、データベース格納コストが他の方式に比べて高い。更新系データベースに適用した場合、データベースの一致を保証するため、あるモジュールがデータベースを書き換えると、他の全てのモジュールのデータベースも書き換える必要がある。また、データベースを共用するため排他制御が必要となる。

【SD 共用方式 (Shared Method) : DB 共用 / SD 格納】 半導体ディスク装置に全てのデータベースを格納し、全てのモジュールで共用する方式。データベースが半導体ディスク装置上にあるため、入出力処理のオーバヘッドが生ずる。また、更新系データベースに適用した場合は、多重方式と同様に排他制御が必要となる。

各方式の構成を表 3.2 に示す。なお、データベースをレコード単位にモジュール台数分に分割し、各モジュールの半導体ディスク装置に格納する SD 非共用方式

表3.2 データベースの分散構成

略称	分割方式	多重方式	S D 共用方式
データベースの共用の可/否	非共用	共用	共用
データベースの格納媒体	メモリ	メモリ	半導体ディスク装置(S D)
構成	 P : プロセッサ DB : データベース	 Main Memory	 Semiconductor disk
モジュール間の負荷の偏り	△	有り	○ 無し
プロセッサの処理量 (ダイナミックステップ数)	△	要	○ 不要
	○	不要	×
	○	不要	△ 要
	○	無し	○ 無し
データベースの格納コスト	○ 安い	△ 高い	○ 安い

(DB非共用／SD格納)は、DB非共用／メモリ格納の分割方式に比べ、入出力オーバヘッドの発生などで明らかにコスト性能比が劣るため除外した。

これら3方式について、まず分散構成でのモジュール間の負荷の偏り、ダイナミックステップ数について考察する。次に、これら要因と性能とを総合し、システムコストとして比較評価を行う。本章では、ノードのスループットを一定とし、レスポンスタイムは一定値以下の条件でのコスト評価を行う。スループットはモジュール台数を通じて、レスポンスタイムはモジュールのプロセッサ使用率を通じてコストに反映される。

3.3.3.2 モジュール間の負荷の偏り

多重方式、SD共用方式では、データベースを共用しており、どのモジュールからも任意のデータベースへのアクセスが可能である。SCFからラウンドロビン方式により各モジュールへデータベースアクセスを割り振ることにより、短時間でのばらつきは生ずるとしても、各モジュールへのアクセス頻度は一様となる。

一方、分割方式では、データベースアクセスはカスタマのレコードを格納しているモジュールでのみ処理可能である。このため、モジュール間のアクセス頻度に偏りが生じやすい。

アクセス頻度の最も高いモジュールのアクセス頻度の平均値 (λ_0) からの偏りの割合を負荷偏り率 (ω) とおく。負荷偏り率 ω については、第4章で詳細に述べるが、カスタマ毎のアクセス頻度が1万倍程度と大きくばらついたとしても、モジュールに収容するカスタマ数が数万程度になれば、負荷偏り率 ω は0.2程度以内に収まる。

負荷の最も高いモジュールのアクセス頻度は $(1 + \omega) \lambda_0$ で表せる。このモジュールにおいても、要求されるレスポンスタイムを保証するためには、あらかじめ偏りを考慮して平均アクセス頻度の $(1 + \omega)$ 倍の負荷を見込んだ設備設計をする必要がある。

3.3.3.3 ダイナミックステップ数の評価

分散構成を採ったときの1呼当たりのダイナミックステップ数（以下、DSと記す）を表3.3に示す要因に従い算出した。特に、表3.3中のI/Oオーバヘッド（D_I）とは、データベースが半導体ディスク装置上に格納されていることによって生ずるオーバヘッドであり、データベースを半導体ディスク装置上に格納した場合とメモリ上に格納した場合のDSの差を示す。具体的には、半導体ディスク装置とメモリとの間の入出力処理のDSはもちろんあるが、さらに入出力処理のためにデータベースをページ単位で管理するためのDSを含む。従来、データベースを外部記憶装置に格納する場合は、入出力処理の容易化のためにデータベースをページ単位で管理してきた。しかし、メモリデータベースの場合は、管理単位をページにとらわれることなく構成することが可能となる。これによりDSの少ないデータベースアクセスが実現できる。図3.2に示すように、データベース処理のDSを約1/2に削減できるとの結果を得た。なお、図中のDSには、分散構成を採ることにより生ずるモジュール選択処理、排他制御等のオーバヘッドは含めていない。以上のことから、I/Oオーバヘッドは表3.3のD₀と同程度とした。

参照系データベース、更新系データベースについてDSの比較結果を図3.3、図3.4に示す。図は1呼当たりのデータベースアクセス回数（Q）が変化した場合のDSを示している。図3.3より、参照系データベースでは、データベース一致処理、排他制御が不要であるため、多重方式のDSが最も少ない。一方、更新系データベースでは、図3.4より、データベース一致処理のDSがモジュール台数Nに比例して大きくなるため、多重方式のDSが最も大きい。分割方式とSD共用方式を比べた場合、分割方式はメモリデータベースのため、半導体ディスク装置格納のデータベースであるSD共用方式に比べてDSが少ない。

表 3.3 分散構成をとったときのDSの要因

要 因	内 容
データベース 処理 (D ₀)	<ul style="list-style-type: none"> ・データベースの格納構成に依存せず各方式に共通にかかるDS ・データベースアクセス回数に比例
モジュール 選択処理 (D _S)	<ul style="list-style-type: none"> ・分割方式で、どのレコードがどのモジュールに格納されているかの判別のために必要なDS ・全てのモジュールがレコードの配備先を示すディレクトリを持ちSCFからのデータベースアクセスに対しては、まず当該ディレクトリを検索し、この結果に従ってレコードを格納しているモジュールにアクセスする方式を採用
データベース 一致処理 (D _W)	<ul style="list-style-type: none"> ・多重方式で、かつ更新系データベースの場合、あるモジュールがデータベースを更新するたびに他の全てのモジュール(n-1)台のデータベースを更新するために必要なDS
排他制御 (D _L)	<ul style="list-style-type: none"> ・多重方式及びSD共用方式でかつ更新系データベースの場合、複数のモジュールがデータベースを共用するための排他制御に必要なDS
I/O オーバーヘッド (D _I)	<ul style="list-style-type: none"> ・SD共用方式の場合、データベースが半導体ディスク装置に格納されており、半導体ディスク装置とメモリとの間の入出力処理及び入出力が容易なページ単位でのデータ管理処理のために必要なDS

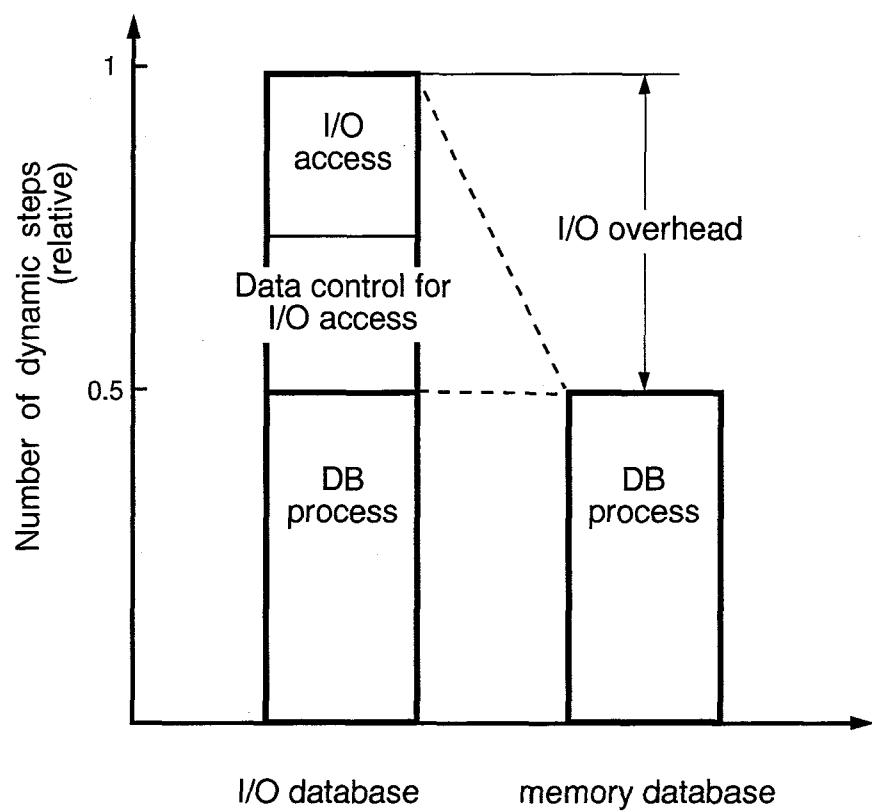


図 3.2 メモリデータベースの効果

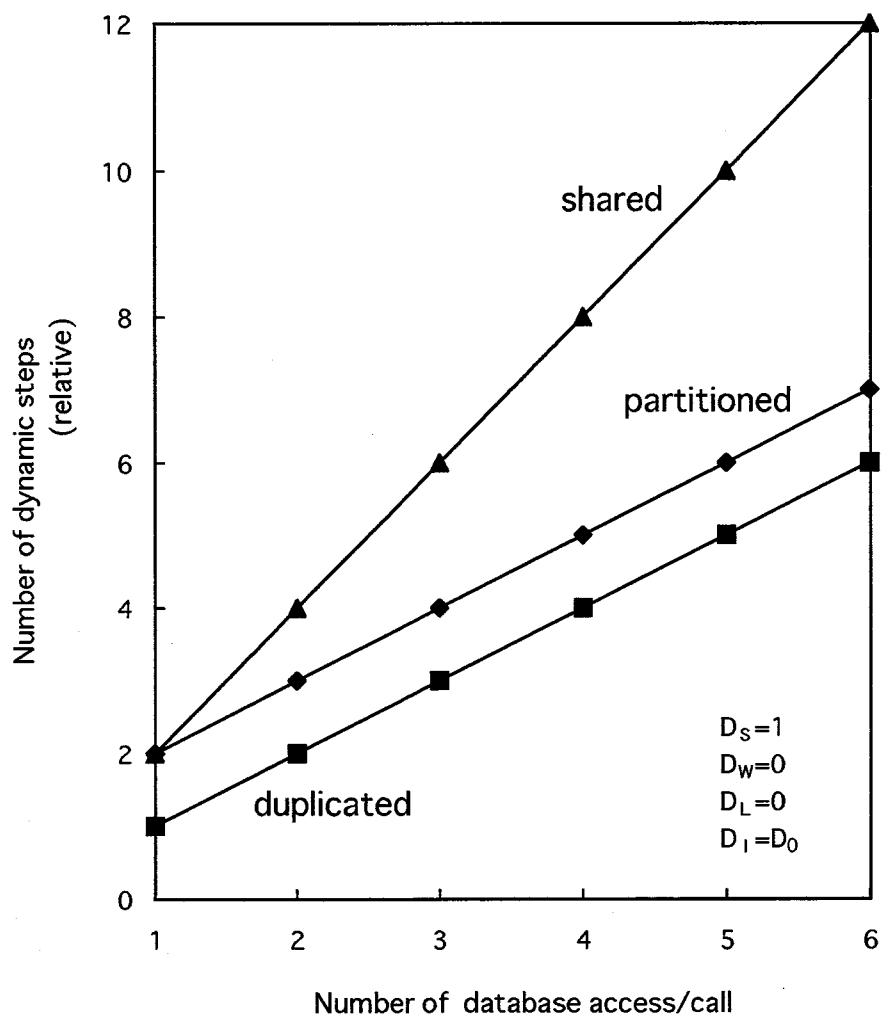


図3.3 D Sの比較（参照系データベース）

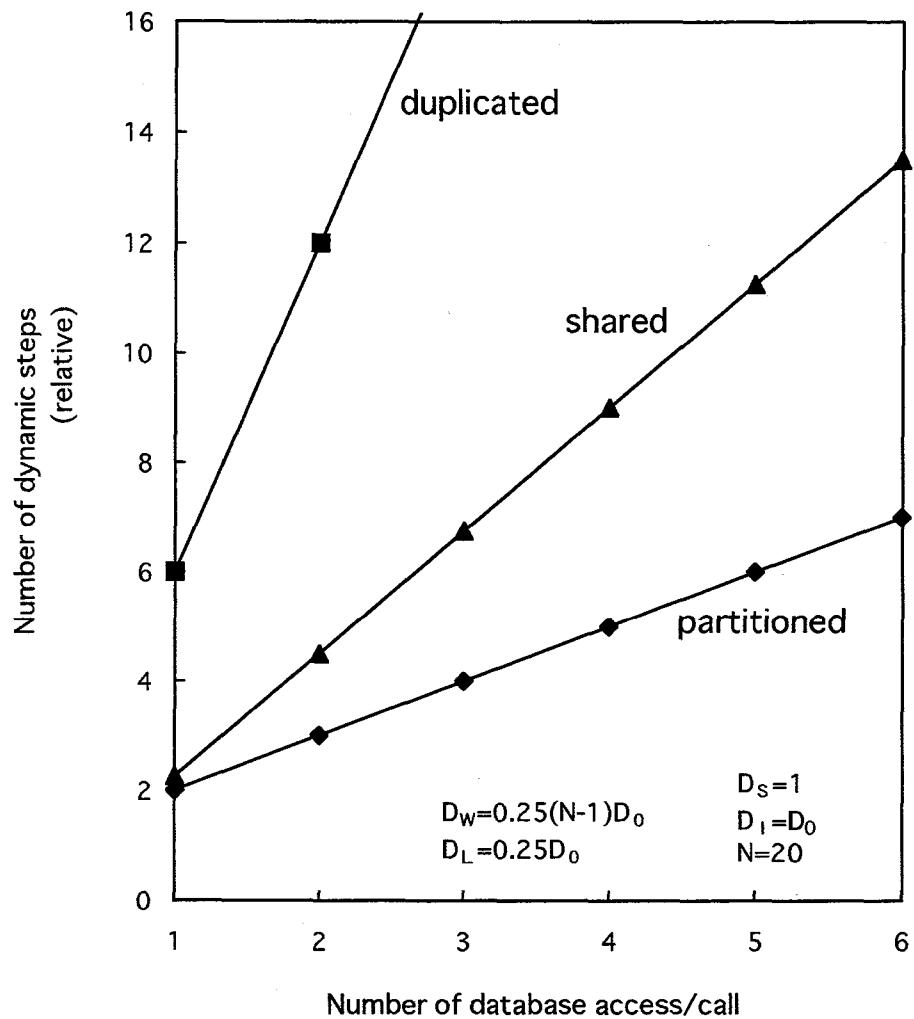


図 3.4 DS の比較 (更新系データベース)

3.3.3.4 コスト評価

各方式の優劣を D S の大小、ハードウェア性能を考慮し、総合的に比較するため、スループットを同一として、各方式でのノードコストを比較する。ここで、ノードコストとは、ノード実現に必要なハードウェアコストとした。分割方式、多重方式、S D 共用方式のノードコストをそれぞれ F_{PA} 、 F_{DU} 、 F_{SH} また、D S を D_{PA} 、 D_{DU} 、 D_{SH} としたとき、ノードコストを以下の式で表す。

$$F_{PA} = \lambda k_{CPU} D_{PA} (1 + \omega\eta) / \rho + k_{MEM} M \quad (3.1)$$

$$F_{DU} = \lambda k_{CPU} D_{DU} / \rho + k_{MEM} M \frac{\lambda D_{DU}}{\rho P} \quad (3.2)$$

$$F_{SH} = \lambda k_{CPU} D_{SH} / \rho + k_{SF} M \quad (3.3)$$

P : プロセッサ性能

ρ : プロセッサ使用率

M : データベース容量

λ : ノード全体にかかるトラヒック

k_{CPU} : 命令当たりのプロセッサコスト

k_{MEM} : 単位容量当たりのメモリコスト

k_{SF} : 単位容量当たりの半導体ディスク装置のコスト

ここで、 D_{PA} は表 3.3 に示した D_0 と D_s の和で表せる。 D_0 はモジュール間で負荷の偏りが生ずるが、 D_s はどのモジュールでも処理可能であるため、負荷の偏りは生じない。 η は D_{PA} のうち、負荷の偏りが生ずる D S の割合 ($\eta = D_0 / D_{PA}$) である。分割方式については、負荷の偏りを次式によりモジュール台数増に反映した。

$$\begin{aligned} N_{PA} &= \frac{\lambda(1+\omega)D_{PA}\eta + \lambda D_{PA}(1-\eta)}{\rho P} \\ &= \frac{\lambda D_{PA}(1+\omega\eta)}{\rho P} \end{aligned} \quad (3.4)$$

式(3.1)～(3.3)の第1項はプロセッサコストであり、第2項はデータベース格納コストである。スループットを同一としているため、 F_{PA} 、 F_{DU} 、 F_{SH} で表されるノードコストが各方式のコスト／性能を表す。

参照系データベース、更新系データベースについて、データベース格納構成のコスト比較結果を図3.5、図3.6に示す。プロセッサ、メモリのコスト、性能データは、現時点での市販品を参考に設定した。

図3.5より、データベース容量が数GB以下では、プロセッサコストに比べデータベース格納のためのコストは小さく、多重方式が最も有利である。従って、参照系データベースのSLPディレクトリでは多重方式をとることとした。これに対し、更新系データベースの場合は、図3.6より、分割方式とSD共用方式の優劣は、分割方式の負荷偏り率 ω に依存する。前節で述べたように ω は0.2程度である。 ω がこの程度であれば分割方式が有利である。従って、更新系データベースのサービスデータについては分割方式をとることとした。

3.4 分散処理ノードの高信頼化

本節では、データベースのリカバリ方法と分散構成の特質を生かした冗長化構成法について明らかにする。

3.4.1 データベースのリカバリ方法

ハード／ソフト障害によりメモリ上のデータベースは、外部記憶装置格納のデータベースに比べて、破壊されやすい。メモリ上のデータベースが破壊された場合のリカバリ方法として、OLTP系システムで一般に用いられている方法は以下の通りである⁽²⁰⁾。

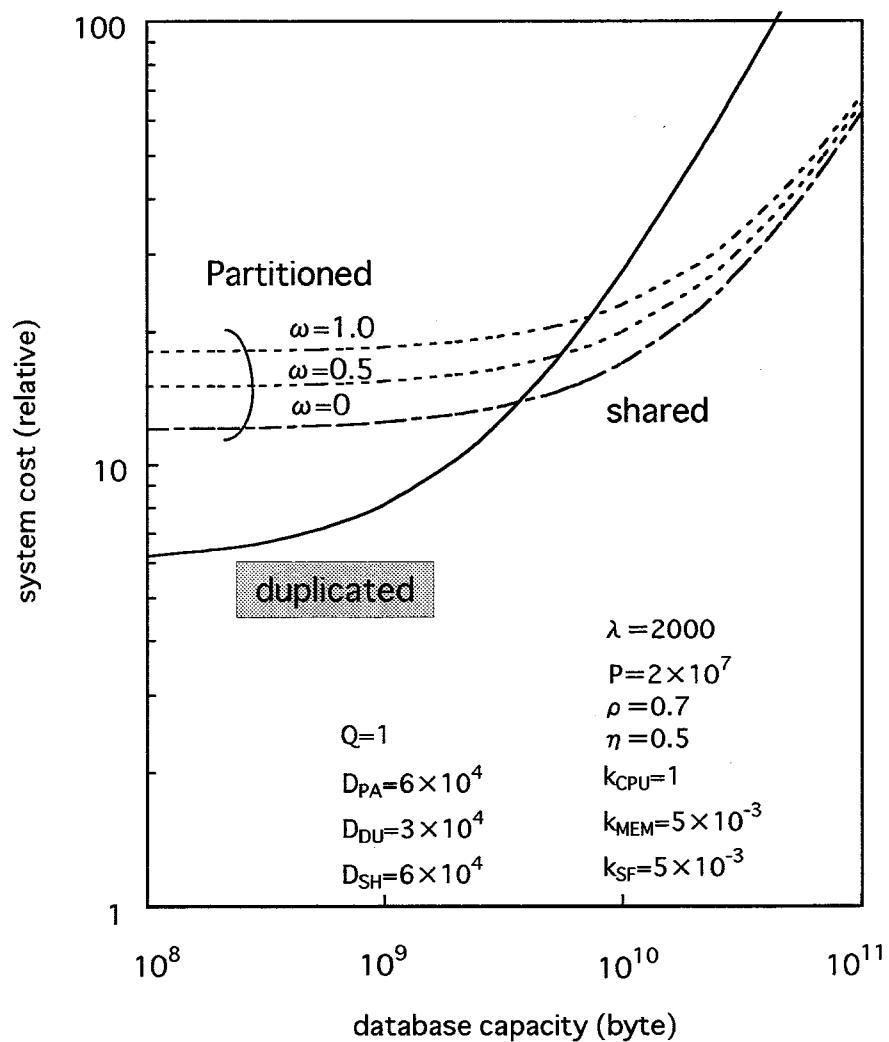


図 3.5 コスト比較（参照系データベース）

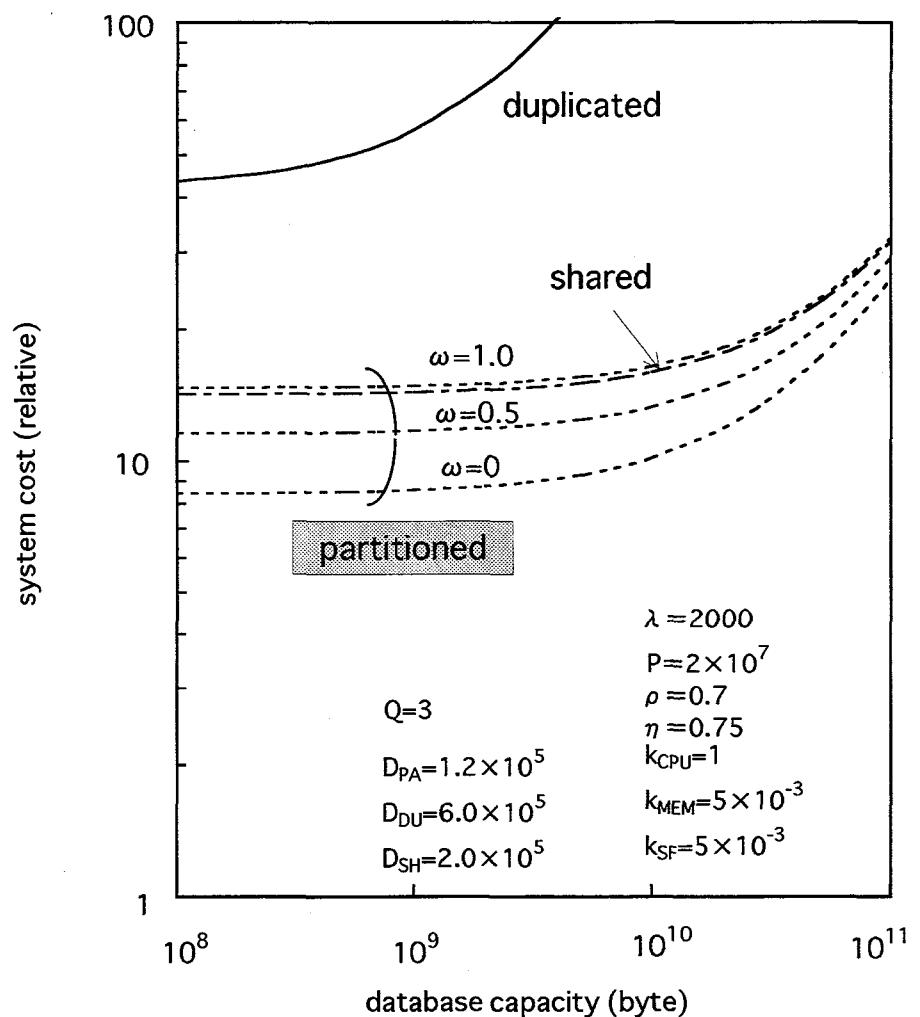


図 3.6 コスト比較（更新系データベース）

- (1) 外部記憶装置（磁気ディスク装置又は半導体ディスク装置）にチェックポイント（C P）時点のデータベース（C P D B）とチェックポイント以降のデータベースの更新ログ（L O G）を格納しておく。
- (2) 障害検出時には外部記憶装置からC P D BとL O Gを読み出し、C P D BにL O Gを上書きすることにより障害直前のデータベースを復元する。

しかし、バンキングシステムをはじめとするO L T P系システムと高度I Nとではデータベースのリカバリに関する条件は大きく異なる。バンキングシステムでは扱うデータは金銭のような重要なデータがほとんどであり、障害やバグの種類を問わずデータベースの復元が必須とされる。通常のハード／ソフト障害に対して、データベースを完全に復元することはもとより、発生頻度は極めて少ないが、アプリケーションプログラム（以下、A P Lと記す）等の論理バグに起因するデータベース破壊に対しても対処が必要となる。このため、長時間に渡ってL O Gを取得しておき、リカバリ時間は長くなるとしても、厳密に復元することを目指している。

一方、高度I Nで扱うデータには、パーソナル通信の位置情報のようにデータを喪失した場合、サービスに直接影響するデータと、回線が使用中か否かを表示するフラグのようにリカバリ時に初期値に戻せばサービスを行う上で支障が生じないデータがある。さらに、位置情報はカスタマの移動とともに更新されるため、リカバリに時間がかかると、たとえデータベースを復元したとしても実際の位置との不一致が生ずる。また、内容が一部喪失したとしても、位置情報の管理エリアをカスタマが移動した場合などには、再登録が行われ、それ以降は支障無くサービスが継続できる。このため、データの重要度に応じたデータベースのリカバリが要求されるとともに、データベースの完全な復元よりもサービス再開を迅速に行うことが優先される。

西原らは、情報の重要度に応じてデータベースの復元レベルを、①初期値に戻す（レベル0）、②C P値に戻す（レベル0.5）、③最新値に戻す（レベル1）に分

類し、LOG取得をレベル1の情報のみに絞る提案を行っている⁽⁴²⁾。ここでは、上述した高度INのリカバリに対する要求条件の分析を通して、レベル1の情報に対し、さらに復元の厳密さとリカバリ時間の2点から考察を加え、表3.4に示すように情報の重要度を新たに設定した。ここで、復元の厳密さとリカバリ時間に関しては、LOG取得時間すなわちCPから障害検出までの時間(T_{CP})と、障害発生から障害検出までの時間(T_B)に着目して考察した。これらの時間関係を図3.7に示す。

一般に、ハード／ソフト障害には、発生後数分以内に検出される障害（以下、障害レベルAと記す）と検出までに数分以上かかる障害（以下、障害レベルBと記す）がある。障害レベルAはハードの2重化チェック、パリティチェックなどで検出されるハード障害、ウォッチドッグタイマ等で検出されるソフト暴走等の障害で、通常のほとんどの障害がこれに該当する。障害レベルBは発生頻度は極めてまれであるが、APL等の論理バグに起因する障害、ハードの多重障害のためチェック機構により検出できなかった障害であり、利用者やAPLレベルでの矛盾発生により検出される障害である。障害発生以降も障害検出まではトランザクション処理は継続され、LOGは取得される。障害発生以降はデータベースの一部が破壊されている可能性がある。しかし、いつ破壊されているかは不明であり、障害発生から障害検出までを対象にリカバリを行う必要がある。リカバリを確実に行うためには、CPDBが保証されていることが必要であり、 $T_{CP} > T_B$ を満足するCPDBからのリカバリが必要となる。障害レベルBからのデータベースの復元を厳密に行うためには、 T_{CP} の長いCPDBから復元することとなり、リカバリに長い時間がかかる。

以上より、以下に示す方法によりデータベースのリカバリを行うこととした。

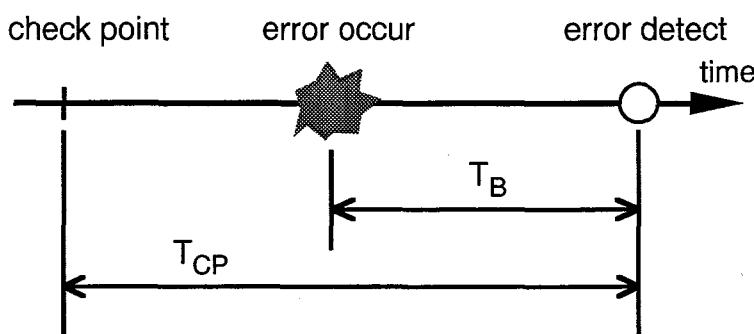
(1) 障害レベルAからのデータベースの復元を行う。障害レベルBは復元の対象としない。

(2) CPDBは数分おきに取得し、CPDB取得が完了した時点で、それ以前

表 3.4 データの重要度に基づく分類

重要度 レベル	復元 レベル	情報例	リカバリ条件		
			LOG の要否	LOG 取得時間	C P D B の要否
1.5	すべての障害に対し障害直前に戻す	なし *1	要	数分以上	要
1	ほとんどのハード／ソフト障害に対し障害直前に戻す	位置情報	要	数分以下	要
0.5	C P 値に戻す	(位置情報)	否	—	要
0	初期値に戻す	回線の使用状況表示フラグ	否	—	否

*1 バンキングシステムの”預金額”が相当する



T_B : エラー発生から障害検出までの時間

T_{CP} : チェックポイントデータベース取得から
障害検出までの時間

データベースのリカバリ条件

$$T_{CP} > T_B$$

図 3.7 T_{CP} , T_B の関係

に取得したLOG及び旧のCPDBは破棄する。

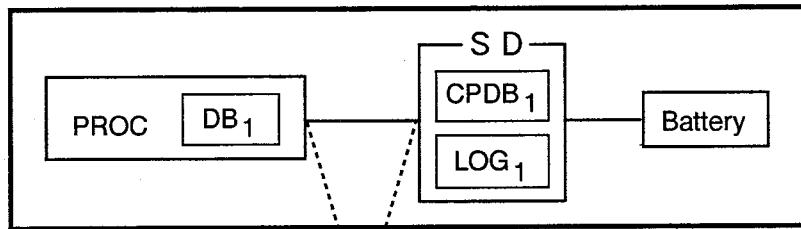
- (3) 障害発生時にはCPDBとLOGを用いて障害直前のデータベースを約1分以内でメモリ上に復元する。
- (4) CPDB及びLOGの格納媒体については、二重化された半導体ディスク装置を用い、バッテリによりデータの不揮発化を図る。

データベースのリカバリのためのハード構成を図3.8に示す。データベースのバックアップ方式には、通信によりLOGをバックアップする別モジュールへ送り込む方式がある⁽¹⁸⁾。ここで提案した入出力による方式は、通信による方式に比べて、以下の点で優れている。

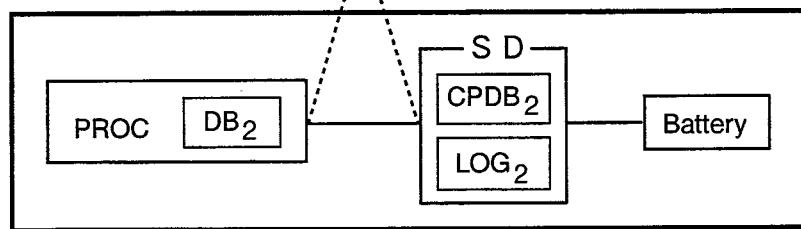
- (1) 入出力処理のためのDSはLOG送受信のためのDSに比べて数分の1ですむ。
- (2) 入出力のデータ転送速度は、通信処理のデータ転送速度に比べて高速であるため、スループットが大きく、レスポンスが短い。
- (3) 不揮発化された半導体ディスク装置にデータベースが格納されており、ノード全体での電源断に対してもリカバリが可能である。

メモリデータベースのリカバリ方式に関し、入出力による方式と通信による方式の適用領域や実用化に向けた課題とその対応策については第6章で述べる。

Module #1



Module #2



PROC : Processor
SD : Semiconductor disk
DB : Current Database
CPDB : Database at checkpoint
LOG : Log data since checkpoint

図 3.8 データベースのリカバリとモジュール間相互バックアップ

3.4.2 冗長化構成法

3.4.2.1 冗長化構成

冗長化の方法はモジュールの機能に応じて異なるため、以下のようにモジュールを2つのタイプに分類した。

【タイプI】同一機能に属するモジュール群の中では、どのモジュールでも処理可能なモジュール。SCFを実現するモジュールや、参照系データベースに多重方式を適用したSLPディレクトリ検索機能を実現するモジュールがこれに該当する。

【タイプII】同一機能のモジュールであっても、特定のモジュールのみが処理可能なモジュール。更新系データベースに分割方式を適用したSDFを実現するモジュールがこれに該当する。

以下、モジュールの各タイプ毎に、ノード内及びノード間での冗長化構成法について述べる。

(1) ノード内冗長化

タイプIについては、同一機能を有するモジュール全体が障害となったモジュールの処理を肩代わりする。すなわち、同一機能に属するモジュール全体として、プロセッサの使用率が許す範囲でグループとしてバックアップする。

タイプIIについては、1対1のモジュールの組を構成し、相互にバックアップする。具体的には、障害となったモジュールの半導体ディスク装置を別のモジュールが引継ぎ、データベース処理を再開する。モジュールと半導体ディスク装置間のチャネル数が少なく、モジュール、半導体ディスク装置の構成が容易であること、及びソフト制御が簡単なことから1対1のバックアップ構成とした。モジュール間バ

ックアップの概略構成を図3.8に示す。

(2) ノード間冗長化

地震、火災等の重度な障害及びノード内の複数モジュール障害に対してもアベイラビリティの向上をはかるため、地理的に離れて設置された2つのノード間で相互にバックアップを行うこととした⁽²³⁾ ⁽²⁴⁾ ⁽⁴³⁾。図3.9に2つのノード間でのバックアップ構成を示す。これら2つのノードは全く同じモジュール構成を取っている。

タイプIについては、2つのノード間でも同一機能を有するモジュール全体で相互にバックアップする構成をとる。タイプIIでは、ノード内ではモジュール相互に1対1のバックアップ構成をとるが、さらに、2つのノード間で対応するモジュール相互に1対1のバックアップ構成を取っている。2つのノード間の対応するモジュール間では正常時には、DBのLOGを相互に送出し、バックアップのためのDBを更新している。

3.4.2.2 制御方式

全てのモジュールについて、正常時には、プロセッサ能力の半分を使用し、残りの半分はバックアップ用とした。タイプIの場合、同一機能のモジュール群の中で半分以上のモジュールが障害となり、要求されるトラヒックを処理できなくなった時点で、障害となったノードのサービスはバックアップする側のノードに全て移される。タイプIIの場合は、ノード内で1台のモジュールが障害となった時は、ノード内のバックアップするモジュールが障害となったモジュールのトラヒックを引き継ぐ。さらに、バックアップしているモジュールも障害となった時点で、障害となったノードのサービスはバックアップする側のノードに全て移される。

地震、火災等の重度な障害でノード全体が一度にサービスができなくなった場合は、その時点で、障害となったノードのサービスはバックアップする側のノードに

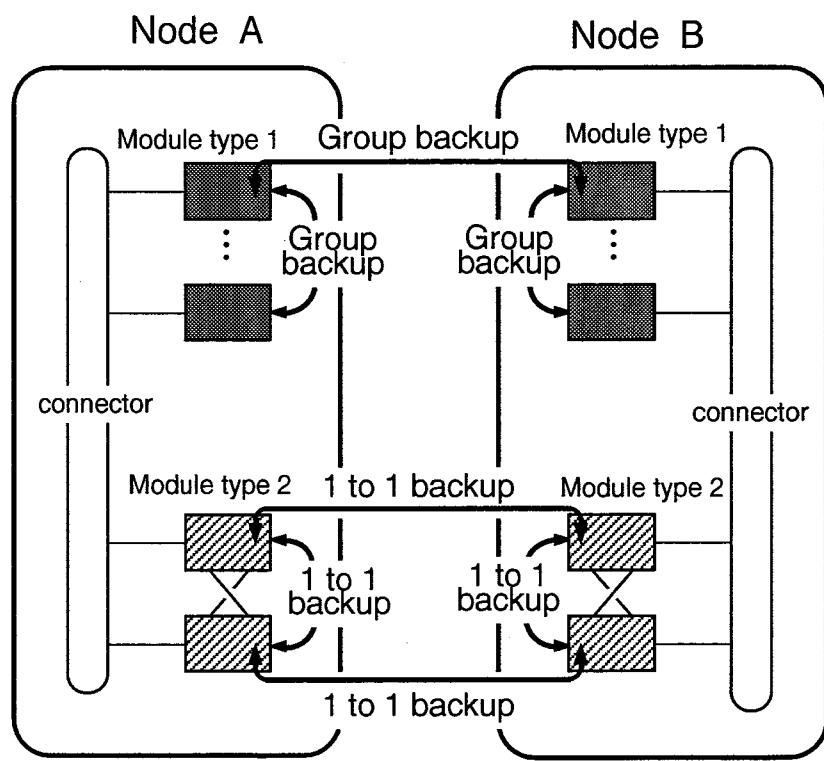


図 3.9 ノード間バックアップ

全て移される。

モジュール相互のバックアップとノード相互のバックアップが同時に発生する確率は極めて低いことから、ノード内と2つのノード相互のバックアップのためのプロセッサ能力の余裕を共用させ、経済的な冗長構成を可能とした。ノード内のモジュール相互のバックアップと2つのノード相互のバックアップを行うことにより、1分／20年以下の不稼働率を達成できる見通しを得ている。

3.5 分散処理ノード構成法

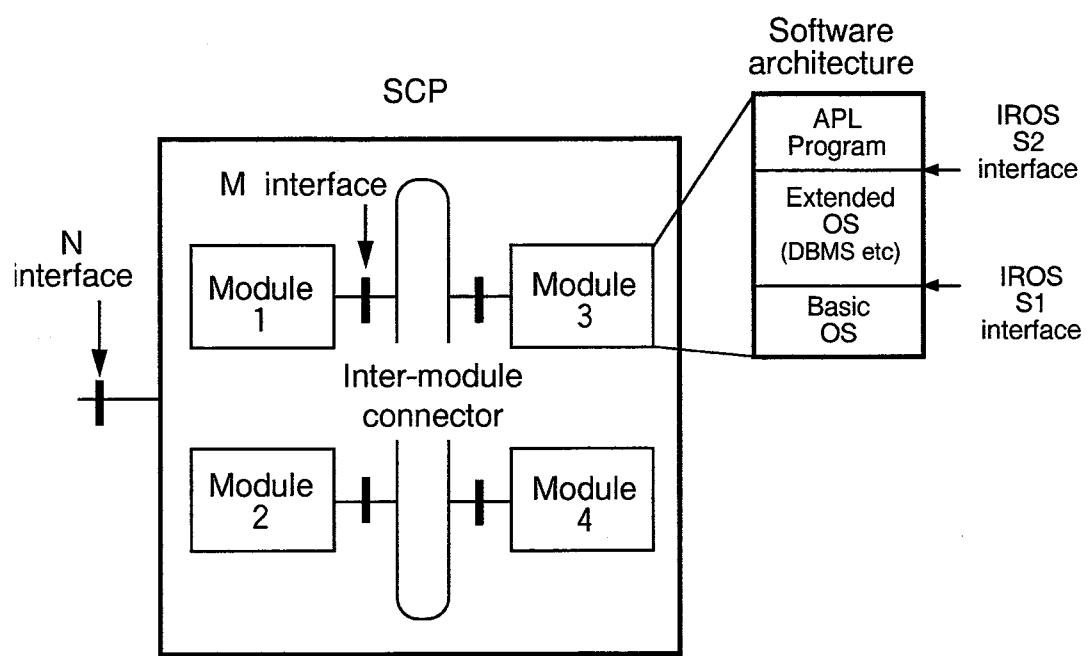
前節までに述べた方式を組み込んで試作した高度IN用SCPのアーキテクチャ及びハードウェア構成について述べる。

3.5.1 アーキテクチャ

交換機との通信はTCP（Transaction Capabilities Application Part）、SMSとの通信はTMN（Telecommunication Management Network）で行う。ノード内はATMレイヤをベースとし、ノード間通信と同一にするため、呼処理系はTCP、保守運用系はTMNを使用したモジュール間インタフェースを規定し、各モジュール間はこのインタフェースを介して通信を行う。また、各モジュール内のソフトウェアについては、階層化を進め、基本OS（S1）、拡張OS（S2）インターフェースにそれぞれIROS（Interface for Realtime Operating Systems）⁽⁴⁴⁾ ⁽⁴⁵⁾を採用した。SCPのアーキテクチャを図3.10に示す。

3.5.2 ハードウェア構成

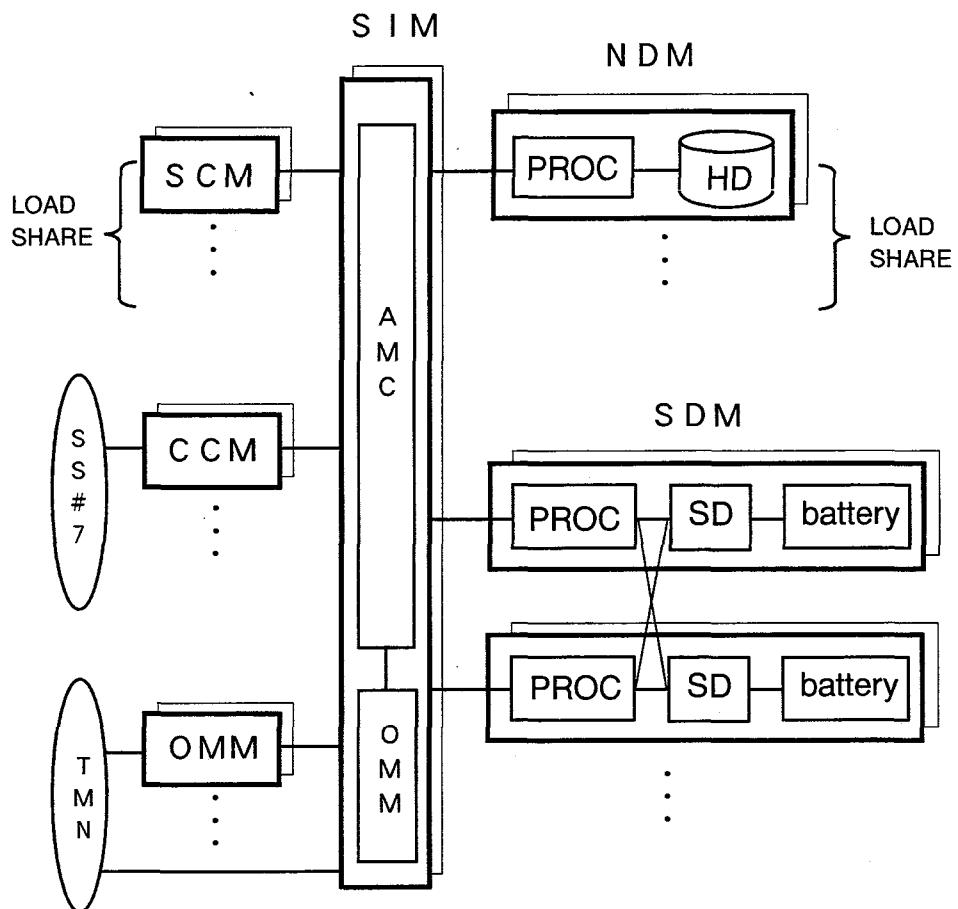
SCPのハードウェア構成を図3.11に示す。SCFはサービス制御モジュール



N interface : Inter-Node Interface

M interface : Inter-Module Interface

図3.1.0 SCPのアーキテクチャ



N D M : network directory module
 S D M : service data module
 S C M : service control module
 C C M : communication control module
 O MM : operation and maintenance module
 S I M : system interface module
 A M C : ATM connector

図 3.1.1 SCP のハードウェア構成

(S C M)、S D Fはサービスデータモジュール (S D M)、S L Pディレクトリ検索機能はネットワークディレクトリモジュール (N D M)で実現した。通信処理には共通線信号網を介した交換機との通信とパケット網を介したS M Sとの通信があるが、各網の将来の発展形態が異なることから、それぞれ通信制御モジュール (C C M)、オペレーション管理モジュール (O M M)で分担した。モジュール間の結合機構としては、大規模ノードに十分なスループットを確保することからA T M結合機構 (A M C)を採用した。A T Mスイッチには、1入力当たり156M b p sのデータ転送能力を有する 16×16 のスイッチを用いた。A M C全体では2.4G b p sのデータ転送速度を持つ。運転・保守機能を実行するためのプロセッシング能力は小さくて済むため、O M Mと共に用化し、A M CとO M Mを一体化してシステムインターフェースモジュール (S I M)で実現した。

3.6 結言

サービス制御ノードのデータベースへのアクセス特性、モジュールの負荷やその偏り等を総合的に考慮した、データベースの分散構成方式の評価方法を明らかにした。具体的には、複数のモジュールからのデータベースの共用の可否、主メモリ／半導体ディスク装置いずれかの格納媒体の選択に注目して構成案を設定した。これら代替案に対して、モジュール間の負荷の偏り、データベース処理でのダイナミックステップ数に対する評価を加え、コスト、性能面から総合的に評価する方法を明らかにした。これに基づき、サービスデータについては分割方式、S L Pディレクトリについては多重方式による分散型メモリデータベース構成を提案した。データベースの高信頼化に関しては、ソフトウェアのバグやハードウェア障害によりデータベースが破壊された場合、データベースの復元の厳密さと復元時間に着目してデータ毎の重要度に基づいたリカバリレベルを設定し、サービス制御ノードのメモリデータベースに適した効率的なリカバリ方法を提案した。また、ノード内のモジュ

ール間バックアップとノード間バックアップを効率的に組み合わせ、高いアベイラビリティを経済的に実現する方法を明らかにした。さらに、モジュール間の結合機構としてATMを用い、これらの技術を総合して試作したサービス制御ノードの全体構成を示した。

第4章 モジュール間の負荷の偏りの評価技術

4.1 緒言

本章では、分散構成における方式設計や設備設計を効率的に行うため、データベースを多数のモジュールに分散配備した場合のモジュール間の負荷の偏りの評価方法を提案する^{(46)～(48)}。

フリーダイヤル等の全国規模の高度INサービスを展開するためには、サービスデータを複数のモジュールに分散配置し、大量のトランザクションを分散処理することが必要となる。モジュール間の負荷の偏りの評価方法を明らかにすることは、少数のサービスデータへトラヒックが集中する傾向にあるサービス制御ノードのデータベースの分散構成を経済的に実現する上で有用である。特に、新しくシステムを設計する段階では、通常、サービス内容も大幅に更新されており、ユーザ特性が明確でない場合が多い。このような中でシステムの基本的な構造となる分散構成を決定してゆくためには、モジュール間の負荷の偏りについて大まかにでも事前に予測できる方法が必要となる。

本章では、カスタマ毎のサービスデータに対するトラヒック比が数千倍から数万倍と大きくばらつく場合のトラヒック量を大きいものから順位付けし、その順位によりカスタマ毎のトラヒック量を近似する方法を提案する。これに基づき、複数のモジュールにサービスデータを分散配備した場合のモジュール間の負荷の偏りの評価方法や、上限値の推定方法を明らかにする。また、現状のプロセッサ性能、高度INサービスにおけるトランザクションの処理量などから負荷の偏りが1～2割程度になることを明らかにする。さらに、この程度の偏りに対しては、SLPディレクトリ検索処理とサービスデータに基づくデータベース処理を組み合わせることにより、負荷を動的に平準化することが有効であることを示す。

4.2 対象とする分散処理モデル

フリーダイヤルサービス、パーソナル通信サービスなどの全国的な高度INサービスを対象とする。システム構成としてはネットワーク内のNカ所にノードを広域的に配置し、1つのノードは更にM台のモジュールで構成する。分散システム構成を図4.1に示す。

4.2.1 処理方式

高度INサービスを受ける呼（以下、IN呼と記す）はユーザからの発信（フリーダイヤルサービスの場合、「0120*****」のダイヤル）により生起される。発信されたIN呼は交換機（SSP）で受信され、共通線信号網を介してトランザクションがサービス制御ノード（SCP）に送られる。SCPではサービスデータを用いてトランザクション処理を行い、その結果をSSPに送信する。SSPではそれに基づき呼の接続処理、切断処理等を行う。SCPへのトランザクションとしては、IN呼を接続、切断するため交換機から送られるトランザクションの他にサービス条件等を変更するためサービス管理ノード（SMS）から送られるトランザクションがある。交換機からのトランザクションに比べてSMSからのトランザクションはトラヒック量が2～3桁少なく、SCPへのトラヒックとしては交換機からのトランザクションのみを考えれば十分である。本章ではIN呼の接続処理、切断処理のトランザクションのみを対象とする。

サービスデータを複数のモジュールに分散して配置する場合、トランザクション処理の前にサービスデータが配置されているモジュールを特定する必要がある。本章では、サービスデータの配置先モジュールを示すディレクトリをすべてのSCPに配置する。SSPは前もってSSP毎に定められたSCPにアクセスする。そのSCPでディレクトリ検索処理を行い、サービスデータが配置されているSCPとSCP内のモジュールを特定し、そこへトランザクションを転送する。ディレクトリ検索処理

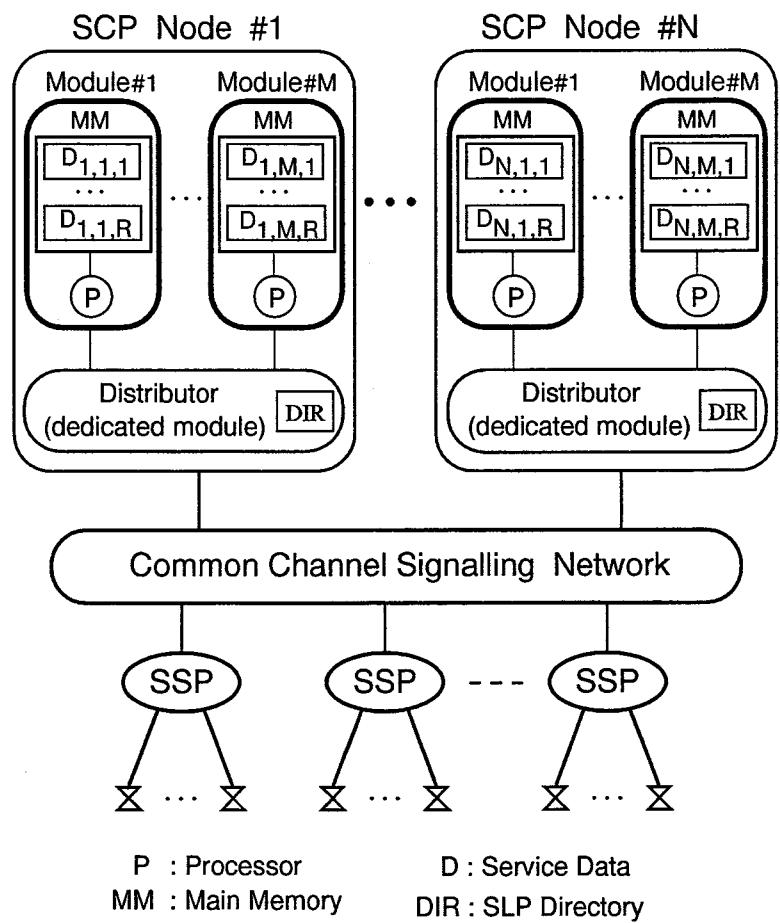


図 4.1 分散システム構成

はIN呼当たり1回行う。IN呼を処理するため複数のトランザクション処理が必要な場合、ディレクトリ検索処理は最初のトランザクションのみを行い、2回目以降のトランザクションは直接サービスデータが配置されているモジュールに送る⁽⁵⁾。

1つのIN呼を処理するためのトランザクション数はサービス内容にも依存するが、フリーダイヤルサービスの場合は2~3回程度であり、フリーダイヤルサービスのIN呼毎のばらつきは小さい。

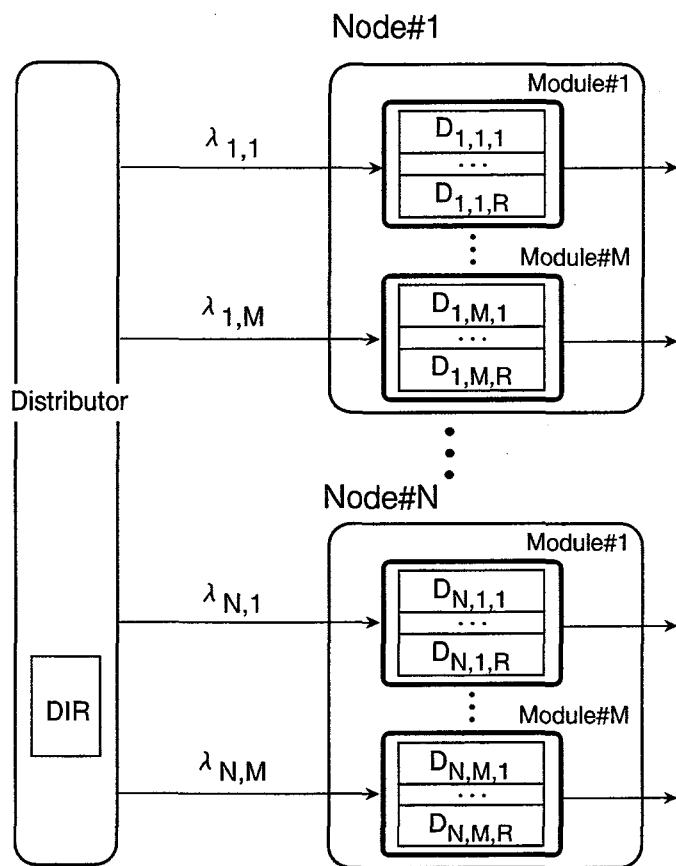
モジュールの負荷はトランザクション当たりの処理量とモジュールに配置したサービスデータへのトラヒック量の総和との積となる。1つのトランザクションの処理はサービスデータへの参照、更新などの処理がほとんどであり、トランザクション毎の処理量の差は小さい。本章ではすべてのトランザクションの処理量は一定とする

フリーダイヤルサービスを例にとると、1つの回線に接続するIN呼の数がカスタマ毎に大きく異なる。さらに、ユーザからのIN呼を受信する回線は1回線のところが大半であるが、大量のIN呼を同時に受け付け可能とするため、カスタマによっては数10から数100の回線を設置する場合がある。これらの回線へIN呼を接続するためのトランザクションはモジュールに配置されたサービスデータを用いて処理され、そのデータへのトラヒック量は非常に高くなる。このようにフリーダイヤルサービスなどの企業を対象としたサービスでは、極く少数のサービスデータへトラヒックが集中する傾向がある。

4.2.2 評価モデル

サービスデータへのトラヒック量のばらつきに起因するモジュール間の負荷の偏りの評価モデルを図4.2に示す。

各モジュールには一律にR個のサービスデータを無作為に抽出して配置する。サービスデータはどれか1つのモジュールにのみ配置される。サービスデータへのトラヒック量の予測などに基づき、モジュールに配置するサービスデータ数を調整する方法も考えられるが、サービスデータ配置操作の自動化が困難となること、季節的要因や



N : Number of nodes
 M : Number of modules
 R : Number of service data allocated in each module
 DIR : SLP Directory
 $D_{i,j,k}$: Service Data
 $y_{i,j,k}$: Amount of traffic for the service data
 $\lambda_{i,j}$: Amount of traffic for the module

$$\lambda_{i,j} = \sum_{k=1}^R y_{i,j,k}$$

(i : Node number , j : Module number , k : Service data number)

図 4.2 評価モデル

経済的要因に基づく変動を予測することは難しいことから、無作為に抽出して配置するとした。

発生したトランザクションは、それを処理するために必要なサービスデータを配置しているモジュールに振り分ける必要がある。評価モデルでは振り分けのためのディレクトリ検索処理は別の専用モジュールで処理されるとした。評価モデルの各モジュールではサービスデータと関係するトランザクション処理のみを実行する。

4.3 サービスデータへのトラヒック量

4.3.1 サービスデータへのトラヒック量の近似

フリーダイヤルサービスにおけるカスタマの1カ月間のIN呼数（IN呼数は度数に比例するとした）の測定結果をもとに、その分布例を図4.3に示す。図4.3の横軸はIN呼数の高いカスタマからの順位を示し、縦軸はIN呼数を示す。図中の黒い点はIN呼数の高い順にカスタマを4つのグループに分割したときの各グループ内のカスタマの順位、IN呼数の平均を示す。点線は各グループ毎のカスタマの占める範囲を示す。図4.3中の直線により、x番目のカスタマのIN呼数Kは、おおよそ

$$K = Cx^{-\alpha} \quad (1 \leq x \leq h) \quad (4.1)$$

の関係で近似できる。ここで、hはカスタマの数である。 α は図4.3中で直線の傾きを示し、カスタマ間のIN呼数の偏りを示す係数と見なすことができる。以後、 α を偏り係数と記す。Cはサービス特性等に依存して決まる比例定数である。サービスデータへのトラヒック量を式(4.1)を用いて近似した場合、図4.3中のフリーダイヤルサービスの4つの点の近似としては、おおよそ $\alpha = 0.8$ と見なすことができる。また、 $\alpha = 0.8$ の場合、 $x = 1$ と $x = 10^5$ のKの比は約 10^4 であり、IN呼数の高いカスタマと低いカスタマとのIN呼数の比は約 10^4 となる。

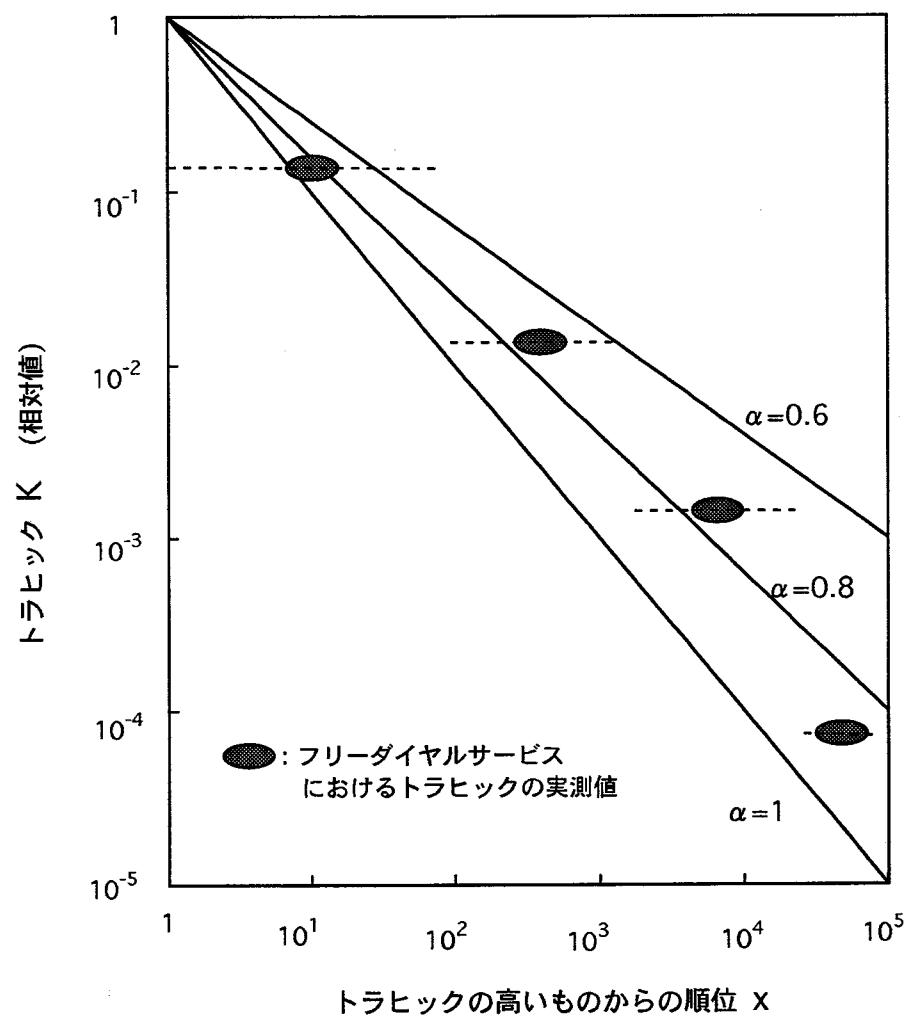


図 4.3 カスタマ毎の IN 呼数

時々刻々変化するサービスデータへのトラヒック量を直接測定することは容易でないため、以下のように仮定する。

(1) IN呼当たりのトランザクション数は一定とする。

(2) 一般にフリーダイヤルサービスなどのIN呼の発生はある時間帯（たとえば、会社が営業を行っている時間など）に集中する傾向がある。どのカスタマについても、IN呼が集中する時間帯および集中するIN呼の割合は同一とし、かつこの時間帯内ではIN呼は一様に分布するものとする。

ノードやモジュールの設備容量の設計ではトラヒックが集中する時間帯を重点的に考慮するため、上記仮定に基づいて、サービスデータへのトラヒック量 y をトラヒック量の高いものからの順位 x を用いて

$$y = cx^{-\alpha} \quad (1 \leq x \leq h) \quad (4.2)$$

の関係で近似する。ここで、最もトラヒック量の高いサービスデータ ($x = 1$) と最もトラヒック量の低いサービスデータ ($x = h$) とのトラヒック量の比（以下、最大トラヒック量比と記す）を g とおくと、 $h = g^{\frac{1}{\alpha}}$ となる。

語彙論において、ある文学作品で使用されている言葉の使用度数 y とその使用度数の大きい方からの順位 x の関係を表すものとして式 (4.2) と同様の関係式が使用されている。これはジップの法則として知られている^(4.9)。

4.3.2 サービスデータの母集団の設定

サービスデータの取りうる x の値は 1 から h までの整数値である。これを 0 から 1 までに正規化するため、

$$z = \frac{x-1}{h-1} \quad (0 \leq z \leq 1) \quad (1 \leq x \leq h) \quad (4.3)$$

とおき、式 (4.2) に式 (4.3) 、 $h = g^{\frac{1}{\alpha}}$ を代入すると次式を得る。

$$y = c \left\{ (g^{\frac{1}{\alpha}} - 1)z + 1 \right\}^{-\alpha} \quad (0 \leq z \leq 1) \quad (4.4)$$

ここまで展開では、 z は等間隔に分布する h 個の離散値をとる。しかし、一般には通信サービスを受けるカスタマ数は非常に多く、 $z = 0$ から $z = 1$ までの間に無数のサービスデータが稠密に分布していると考えられる。式 (4.4) は順位を 0 から 1 までに正規化したサービスデータの母集団の特性を示すものとみなすことができる。以下、式 (4.4) で表されるサービスデータの母集団を $U(\alpha, g)$ で表す。式 (4.4) はパレート分布と言われているものと同一の形態をとる。パレート分布は所得とそれを得る人数の関係を近似する方法として Pareto によって最初に提案されたものである⁽⁵⁾。パレート分布は高額所得者の分布とも言われ、高い所得を有する人数は少数であるが、全体の所得に占める割合が高い不均等の分布を良く近似できると言われている。

4.3.3 近似式の特性

母集団 $U(\alpha, g)$ から任意に 1 つのサービスデータを抽出した場合のトラヒック量の期待値 y_0 は次式で表せる（付録 4.1 参照）。

$$y_0 = \int_0^1 y dz = \int_0^1 c \left\{ (g^{\frac{1}{\alpha}} - 1)z + 1 \right\}^{-\alpha} dz$$

$$= \begin{cases} \frac{c(g^{\frac{1}{\alpha}-1} - 1)}{(1-\alpha)(g^{\frac{1}{\alpha}} - 1)} & (\alpha \neq 1) \\ \frac{c \log g}{(g-1)} & (\alpha = 1) \end{cases} \quad (4.5)$$

式 (4.5) の $\alpha \neq 1$ の式において、 α を 1 に近づけた場合の極限値は $\alpha = 1$ の式と同一となる。このため、以後、 α の値に係わらず式 (4.5) の $\alpha \neq 1$ の式を用いることとする。

次に、サービスデータへのトラヒック量の大きい方から z までの累積値が期待値に占める割合 η は次式で表せる（付録4.2参照）。

$$\begin{aligned}\eta &= \frac{\int_0^z y dz}{y_0} \\ &= \frac{(z - zg^{-\frac{1}{\alpha}} + g^{-\frac{1}{\alpha}})^{1-\alpha} - g^{1-\frac{1}{\alpha}}}{(1 - g^{-\frac{1}{\alpha}})}\end{aligned}\quad (4.6)$$

z と y との関係を図4.4、 z と η との関係を図4.5に示す。図4.4より、 $\alpha = 1$ の場合、 $z = 0$ から $z = 0.1$ の間で、トラヒック量は1から 10^{-3} まで急激に低下する。一方、 $z = 0.1$ から $z = 1$ の間では、トラヒック量は 10^{-3} から 10^{-4} まで緩やかに低下する。

α が1より大きくなるに従って、トラヒック量の高い領域での低下傾向は緩やかとなり、 z の全領域に渡ってトラヒック量の低下の度合いが一定となってくる。これは、 α が無限大となった場合、 y が次式に収束することから裏付けられる（付録4.3参考）。

$$\begin{aligned}\lim_{\alpha \rightarrow \infty} y &= ce^{-(l_{og}g)z} \\ &= cg^{-z}\end{aligned}\quad (4.7)$$

一方、 α が1より小さくなるに従って、トラヒック量の高い領域では一段と急激に低下する。但し、全体としてトラヒック量の高い領域が減少し、ほとんどがトラヒック量の小さいサービスデータのみの一様な分布に近づく。これは、 α が0となった場合、次式に示すように η が z に比例して増加することから裏付けられる。

$$\lim_{\alpha \rightarrow 0} \eta = z \quad (4.8)$$

以上のことから、式(4.4)は α を変化させることにより、一様な分布から指數関数的に低下の度合いが一定となる分布まで幅広く表現できることがわかる。

全体として、 α が極めて小さい領域を除けば少数のサービスデータが大部分のトラヒック量を占める。たとえば $g = 10^4$ の場合、図4.5より $\alpha = 1$ 、 $\alpha = \infty$ ともにトラヒック量の高いサービスデータの上位約2割で総トラヒック量の約8割を占める。モ

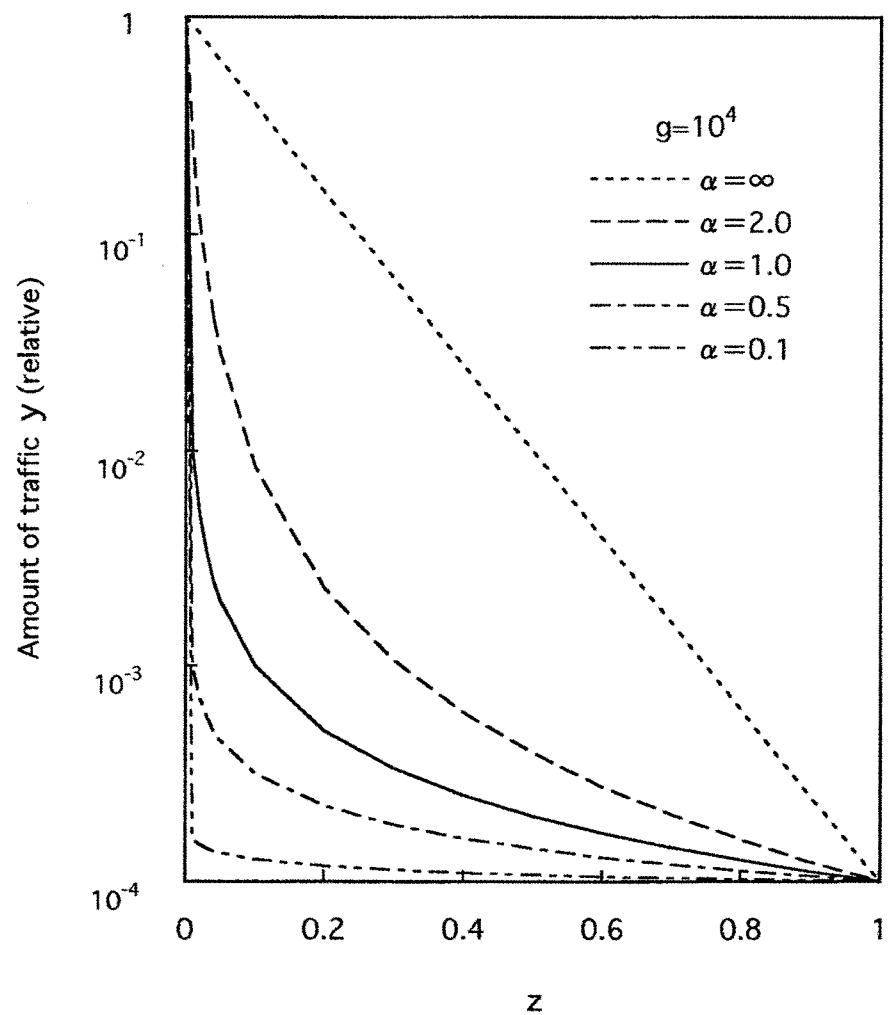


図4.4 トラヒック量

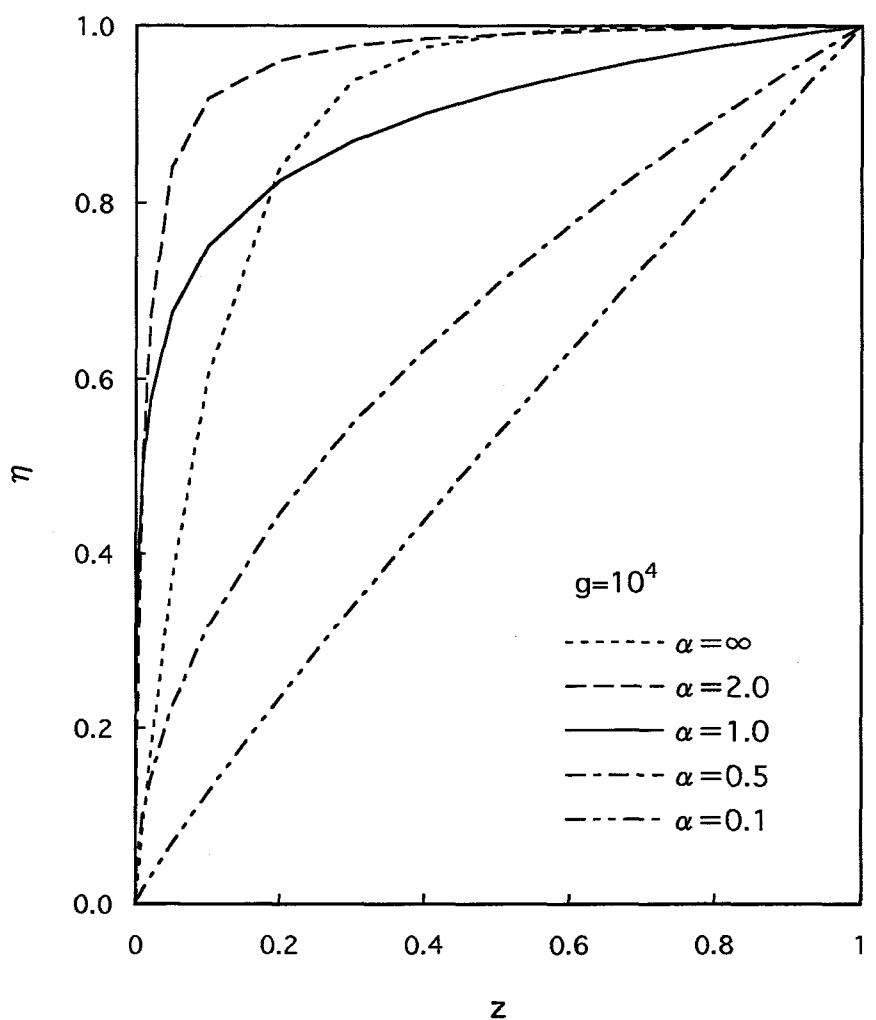


図 4.5 z と η の関係

ジユール間の負荷の偏りに影響を及ぼすのはトラヒック量の高い少数のカスタマのばらつきである。図4.5では z の小さい領域での分布が明確でないため、横軸を対数として図4.6に示す。図4.6より $z = 0.01$ の場合、 $\alpha = 1$ では $\eta = 0.5$ であるが、 $\alpha = \infty$ では $\eta = 0.1$ となる。 $\alpha = 1$ の分布は $\alpha = \infty$ に比べて、極く少数のサービスデータへトラヒックが集中する分布を適切に近似できる。以上の考察に基づき、式(4.4)により、サービスデータへのトラヒック量を近似し、これを用いてモジュール間の負荷の偏りについて評価する。

4.4 モジュール間の負荷の偏り

4.4.1 負荷偏り率

$U(\alpha, g)$ の中から無作為にサービスデータを R 個抽出し、モジュールに配置する。 R 個のサービスデータへのトラヒック量をそれぞれ $cx_1^{-\alpha}, cx_2^{-\alpha}, \dots, cx_R^{-\alpha}$ とおくと、総トラヒック量 λ は、

$$\lambda = cx_1^{-\alpha} + cx_2^{-\alpha} + \dots + cx_R^{-\alpha} \quad (4.9)$$

となる。 λ の期待値を λ_0 とおき、その偏りを $\frac{\lambda}{\lambda_0}$ で表す。このとき標準偏差 σ は次式で表せる（付録4.4参照）。

$$\begin{aligned} \sigma &= \left\{ \frac{1}{(h-1)^R} \int_1^h \cdots \int_1^h \left(\frac{\lambda}{\lambda_0} - 1 \right)^2 dx_1 \cdots dx_R \right\}^{1/2} \\ &= \frac{A(\alpha, g)}{\sqrt{R}} \end{aligned} \quad (4.10)$$

$$\text{但し、 } A(\alpha, g) = \left\{ \frac{(1-\alpha)^2 (g^{\frac{1}{\alpha}} - 1)(g^{\frac{1}{\alpha}-2} - 1)}{(1-2\alpha)(g^{\frac{1}{\alpha}-1} - 1)^2} - 1 \right\}^{1/2} \quad (4.11)$$

モジュールの最大負荷を $(1 + \omega)\lambda_0$ で表し、 ω を 3σ で近似する。 λ の分布が

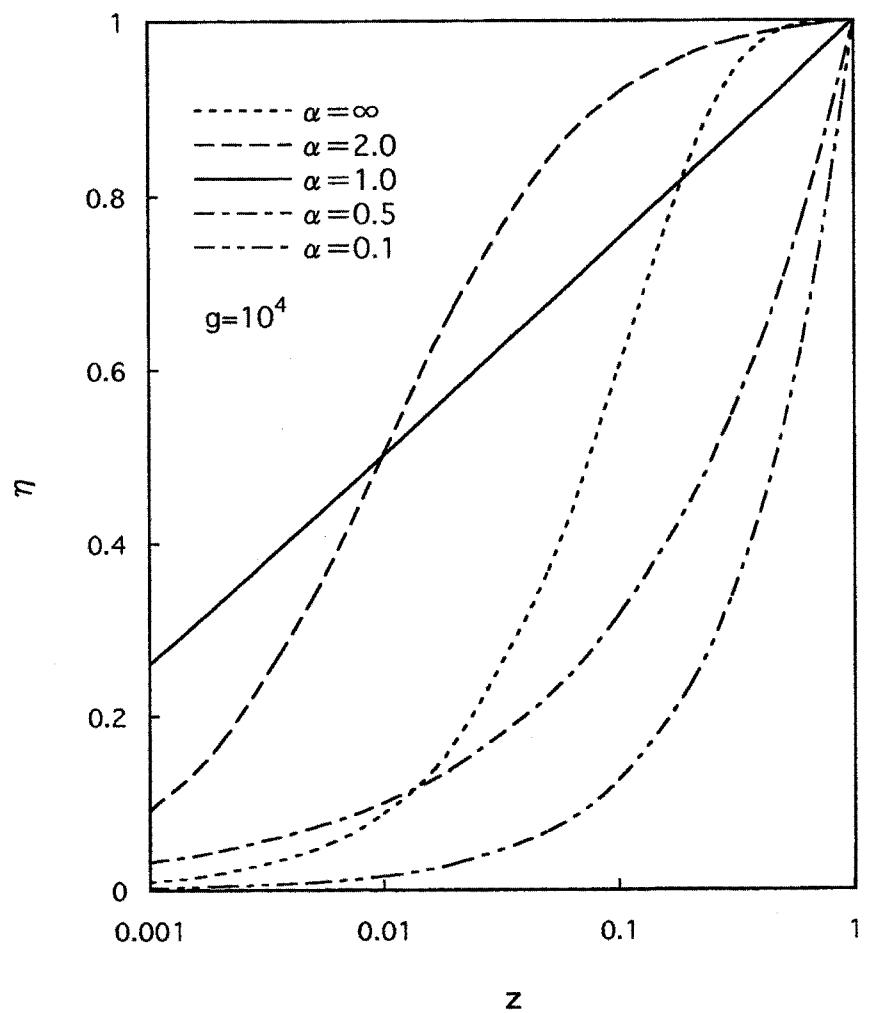


図 4.6 z と η の関係

正規分布で近似できるとした場合、約 99.95% のモジュールがこの最大負荷以内に収まる。数百モジュールが設置されたとしても、ほとんどのモジュールの負荷がこの最大負荷をこえることはない。以下、 ω を負荷偏り率と記す。

$$\omega \approx 3\sigma = \frac{3A(\alpha, g)}{\sqrt{R}} \quad (4.12)$$

ω は \sqrt{R} に反比例し、その係数 A は α 、g に依存して決まる。

4.4.2 サービスデータ数と負荷偏り率の関係

モジュールに格納するサービスデータ数 R と負荷偏り率 ω との関係を図 4.7 に示す。なお、図 4.7 では最大トラヒック量比 g は 10^4 としている。図 4.7 より、 $\alpha = 0.8$ の場合、R と g が同程度、すなわち R が 1 万の場合は $\omega \approx 0.3$ と大きい。R が g の 10 倍になると $\omega \approx 0.1$ 程度となり、さらに、R が g の 100 倍になると ω は数% 以下となり、無視できる程度となる。

g と ω の関係を図 4.8 に示す。g が小さくなるに従って ω は小さくなる。たとえば、R が 10 万の場合、 $g = 10^4$ では $\omega \approx 0.1$ であるが、 $g = 10^2$ になると $\omega \approx 0.02$ に低下する。

図 4.3 のフリーダイヤルサービスの場合には $g = 10^4$ であったが、パーソナル通信サービスのように個人レベルの通話の場合には g は小さくなり、それに伴って ω は低下する。たとえば、g がフリーダイヤルサービスに比べて 1/100 程度になったとした場合、 ω は数% 以下となり、無視できる程度となる。

4.4.3 偏り係数と負荷偏り率の関係

負荷偏り率 ω は、式 (4.12) に示す通り、 \sqrt{R} に反比例し、係数 A に比例する。サービスデータへのトラヒックの偏り係数 α と負荷偏り率 ω の係数 A との関係を図 4.9 に示す。図 4.9 からわかる通り、g の値に係わらず $\alpha = 1$ で A は極大値をとる。これは A を α で微分し、 $\alpha \rightarrow 1$ とすると

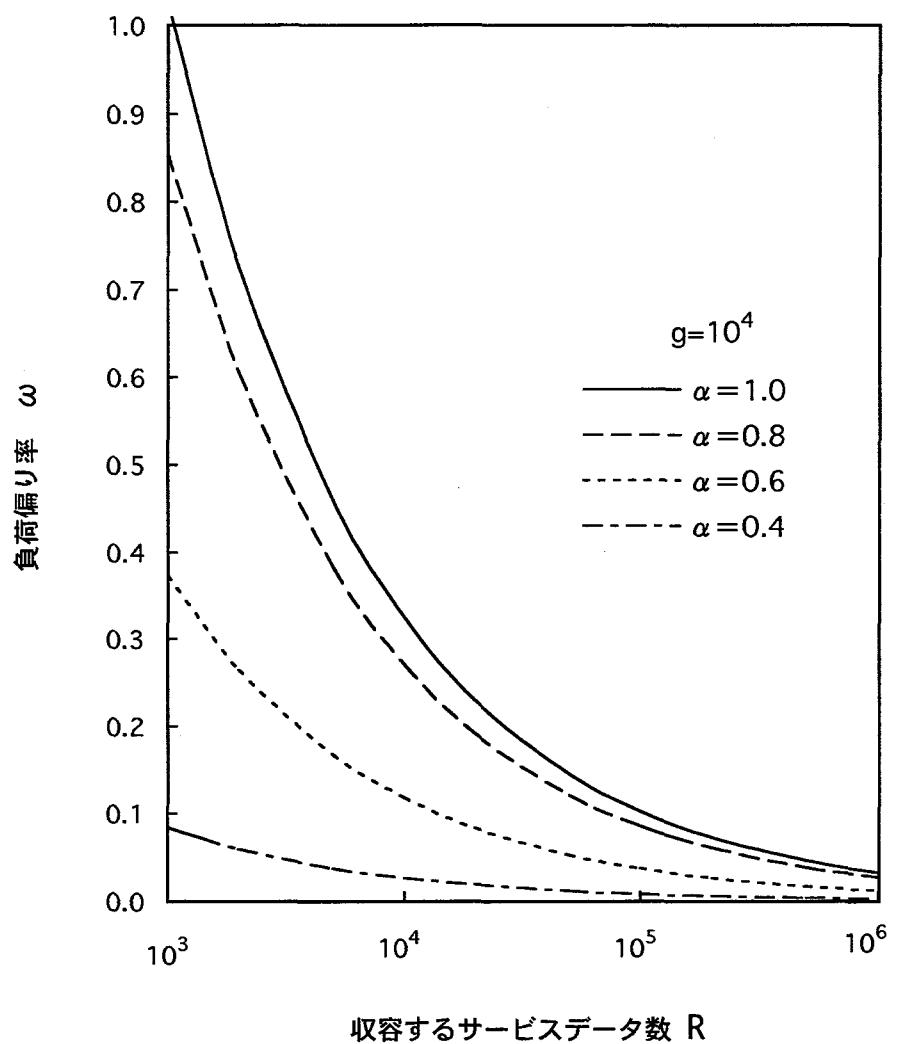


図 4.7 R と ω の関係

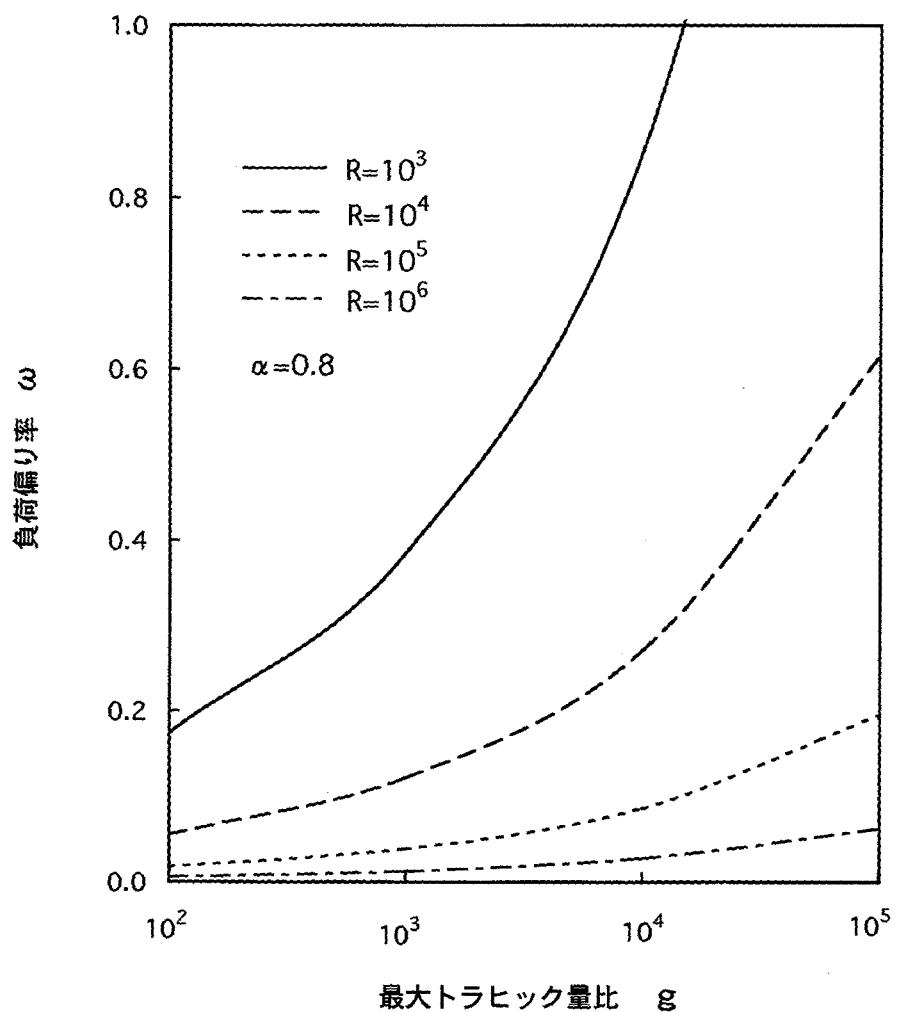


図4.8 g と ω の関係

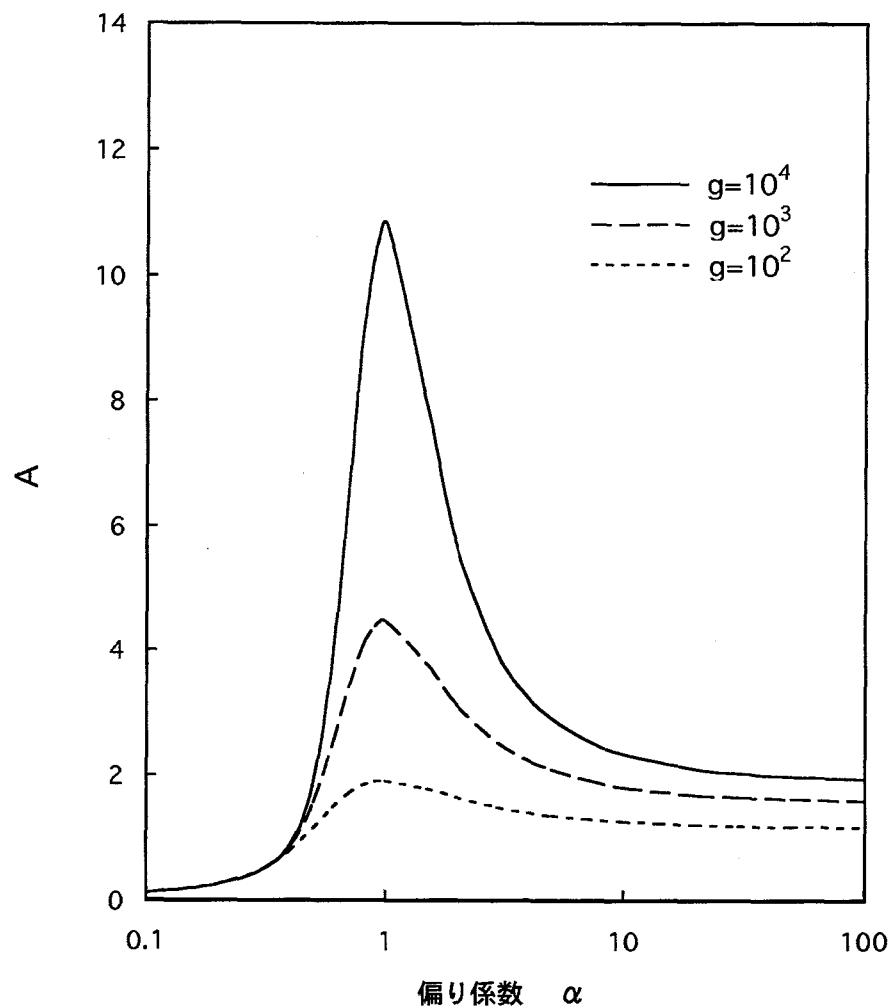


図 4.9 α と A の関係

$$\lim_{\alpha \rightarrow 1} \frac{dA(\alpha, g)}{d\alpha} = 0 \quad (4.13)$$

となる（付録4.5参照）ことから裏付けられる。 α が0から1に近づくに従って、トラヒック量の小さな大多数のサービスデータ間でのトラヒック量の差が拡大するため、Aは増加する。一方、 α が1より大きくなると、高いトラヒック量をとるサービスデータ間のばらつきが縮小する傾向にあり、Aは減少する。 $\alpha = 1$ と $\alpha = \infty$ でのAの比($A(1, g)/A(\infty, g)$)は、 $g = 10^4$ の場合、約6となる。

$\alpha = 1$ でAが極大値、すなわち ω が極大値をとるため、 α の値が明確でない場合でも、 g について推定できるときは $\alpha = 1$ で評価することにより、 ω の上限を把握することができる。なお、図4.3に示したフリーダイヤルサービスの場合、 $\alpha = 1$ と $\alpha = 0.8$ のときの ω の差は約20%程度である。

4.5 分散構成への適用

4.5.1 負荷の偏りの影響

前節までの評価方法を用いて、現状のプロセッサ性能とトランザクション処理のダイナミックステップ数を基に分散構成への影響を明らかにする。

最も負荷の高いモジュールの負荷を L_{MAX} 、プロセッサ使用率を ρ_{MAX} 、プロセッサ性能をPとすると次式が成り立つ。但し、ここで負荷はモジュールへのトラヒック量とダイナミックステップ数(D)との積とした。

$$L_{MAX} = (1 + \omega)Ry_0 D \quad (4.14)$$

$$\rho_{MAX} = \frac{L_{MAX}}{P} \quad (4.15)$$

最も負荷の高いモジュールについても要求されるレスポンスタイムを保証するためにはRを少なくしてプロセッサ使用率 ρ_{MAX} を0.8程度以下に抑える必要がある。Rは式(4.12)、(4.14)、(4.15)より次式となる。

$$R = \rho_{MAX} \frac{P}{y_0 D} - \left(\sqrt{(3A)^2 \rho_{MAX} \frac{P}{y_0 D} + \frac{(3A)^2}{4}} - \frac{(3A)^2}{2} \right) \quad (4.16)$$

式 (4.16) の第 1 項は負荷の偏りが無い場合の配置可能なサービスデータ数 (R_0) 、第 2 項は負荷の偏りに対応するため減少すべきサービスデータ数を表している。

なお、サービスデータはサービス内容にも依存するが最大でも数 100 バイト～1000 バイト程度と評価しており、サービスデータ数が数万～数十万程度までは、主メモリ容量からの制約はほとんど無い。

P と R の関係を図 4.1.0 に示す。図 4.1.0 ではモジュール間の負荷の偏りが無い場合、すなわち $\alpha = 0$ のときの収容可能なサービスデータ数 R_0 との相対値 (R/R_0) で表している。現状のプロセッサ性能、ダイナミックステップ数およびサービスデータへのトラヒック特性から、 $D = 100$ キロステップ、 $P = 20$ MIPS、 $g = 10^4$ 、 $y_0 = 3 \times 10^{-3}$ トランザクション／秒／サービスデータとすると、モジュール間の負荷の偏りに関しては以下のことが言える。

(1) $\alpha = 0.8$ の場合、 $R \approx 4.8$ 万、 $\omega \approx 0.12$ となる。サービスデータ間の最大トラヒック量比 g が 10^4 と大きく、かつ各サービスデータへのトラヒック量の予測などを行わず無作為に抽出して同数配置したとしても、 R が 5 万程度になるとモジュール間の負荷の偏りを 1～2 割程度に抑えることができる。

$\alpha = 1$ の場合はモジュール間の負荷の偏りが極大となるが、 $R = 4.7$ 万であり、 $\alpha = 0.8$ の場合と比べてほとんど差はない。

(2) ノード全体を 1 つの単位として見たとき、ノード間の負荷の偏りもモジュール間の負荷の偏りと同様に評価することができる。ノード間の負荷偏り率を ω_n とすると、ノード内には $M \cdot R$ 個のサービスデータが配置されていることから、

$$\omega_n = \frac{3A(\alpha, g)}{\sqrt{M \cdot R}} \quad (4.17)$$

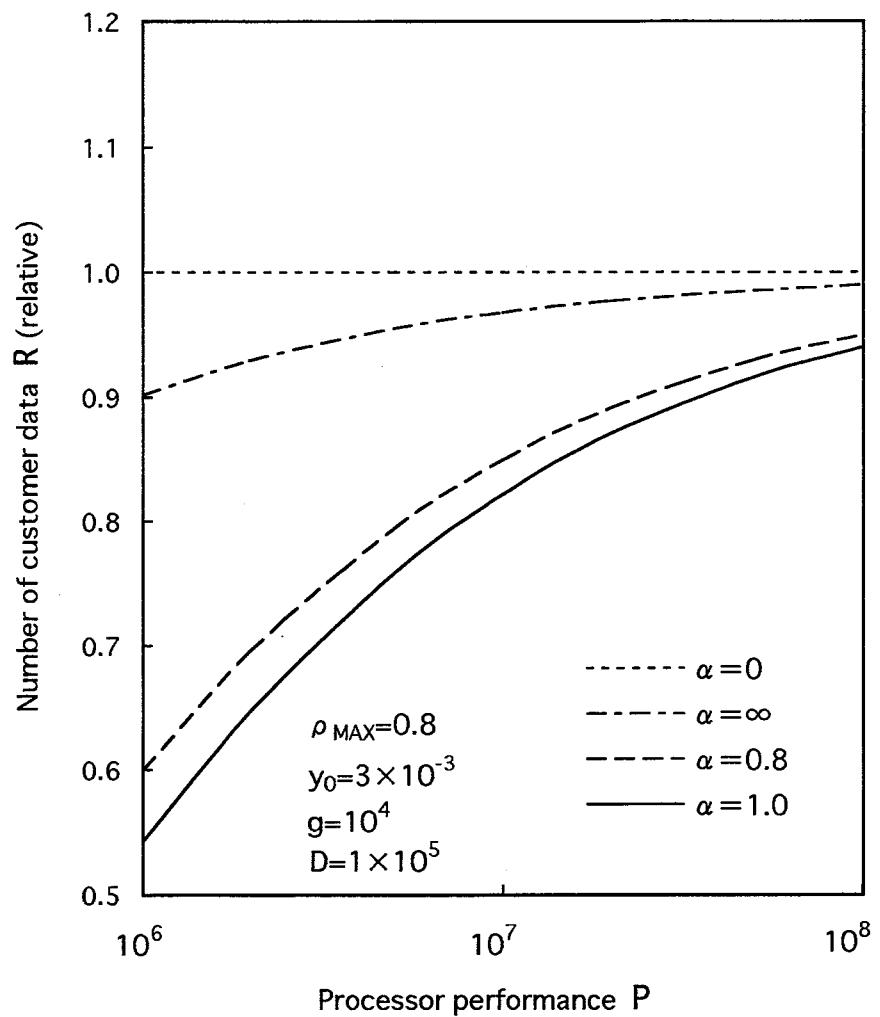


図 4.1.0 P と R の関係

と表せる。Rを5万程度、ノード内のモジュール数を20程度として、ノード内のサービスデータ数を100万 ($M \cdot R = 10^6$) とすると、 ω_n は数%以下となり、ノード間の負荷の偏りは無視できる程度となる。

以上、フリーダイヤルサービスを中心にモジュールに格納できるRについて考察した。式(4.16)はサービスデータからなるデータベースを複数のモジュールに分散収容したシステムであればサービスに依存せず適用できる。 R/R_0 は $P/y_0 D$ に依存し、Pが大きくなるに従って大きくなり、 $y_0 D$ が大きくなるに従って小さくなる。現状のプロセッサ性能を想定すると、ダイナミックステップ数の大きい、より高度なサービスや、平均的なトラヒック量の高いサービスの場合には、 $y_0 D$ が大きくなり R/R_0 が減少する。また、将来プロセッサ性能が高くなった場合でもサービスの高機能化に伴い、ダイナミックステップ数の増加がプロセッサ性能の向上度より大きい場合には R/R_0 が減少する。

4.5.2 負荷平準化の可能性

前節まで無作為に抽出したサービスデータでの負荷の偏りを評価してきたが、次にモジュール間での負荷の平準化の可能性について考察する。

負荷の平準化を実施するとした場合、サービス開始前にカスタマ毎の事前の予測に基づいて負荷が平準化するように、サービスデータを配置するモジュールを選定する方法とサービス開始後、動的に負荷を平準化する方法とが考えられる。前者については、4.2節で述べたように、データ配置操作の自動化が困難となること、季節的要因や経済的要因に基づく変動を予測することが難しいことから実用的ではない。モジュール間の負荷を動的に平準化する方法としては以下の2案が考えられる。

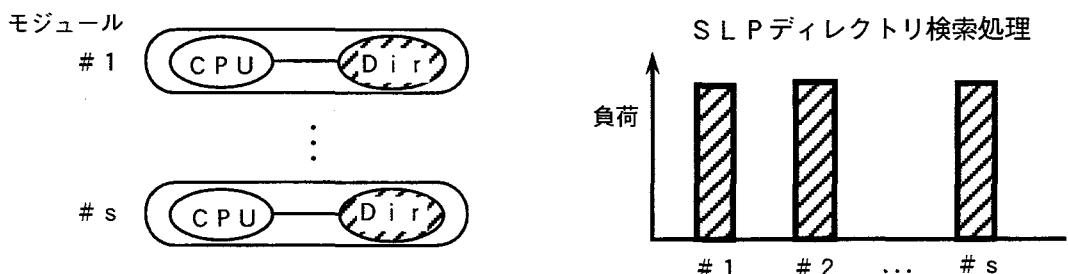
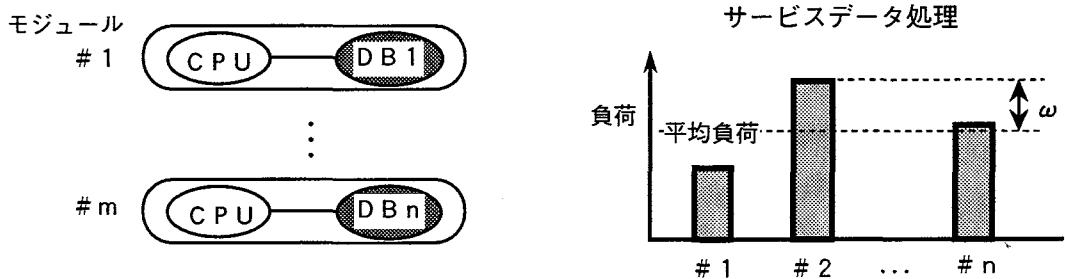
【方法1】ノード内のモジュール間で負荷の高いモジュールから負荷の低いモジュールへサービスデータそのものを移動し、モジュール間の負荷を平準化する。

【方法 2】ノード内にサービスデータに結びつく処理だけでなく、サービスデータと独立に実行可能な処理を共存させ、後者をその時点で負荷の低いモジュールに多く割り当てるることにより、モジュール間の負荷を平準化する。

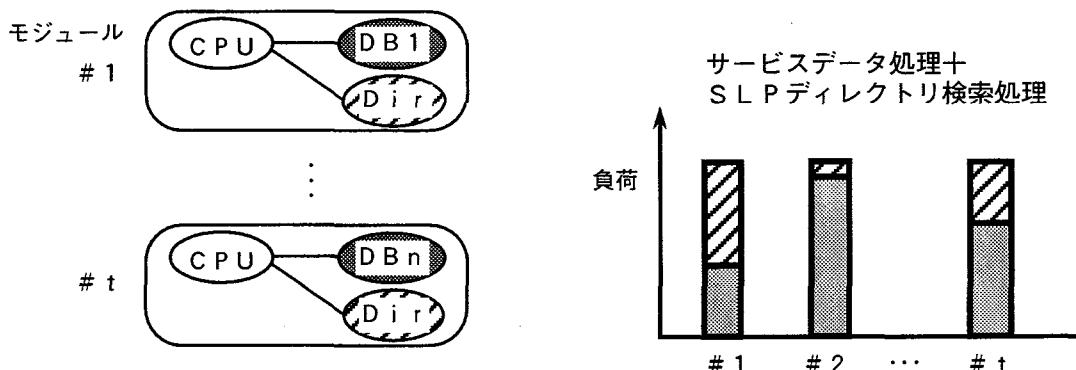
方法 1 は確実に負荷を平準化できるが、サービスデータの移し換えのためには、当該サービスデータの一時的な閉塞とモジュール間の転送が必要になるなど、プロセッサのオーバヘッドが大きく、頻繁な実行や高負荷時の実行には問題が多い。また、サービスデータの再配置により各 SCP が持つ SLP ディレクトリの変更処理も必要となる。このため、夜間、負荷が低いときなどに実施する必要があり、負荷の動的な変化には十分に対応できない。

方法 2 でどの程度までモジュール間の負荷を平準化できるかは、サービスデータと独立に実行可能な、どのモジュールにも割り当てできる処理量に依存する。負荷平準化のためのプロセッサオーバヘッドは小さく、負荷の動的な変化には十分対応できると考えている。

方法 2において、サービスデータと独立に実行可能な処理として、4.2 節で述べた SLP ディレクトリ検索処理がある。図 4.2 の評価モデルでは別の専用モジュールで SLP ディレクトリ検索処理を行い、トランザクションの振り分けを行うとして分離したが、この処理を各モジュールのトランザクション負荷に応じて動的に割り当てる考えを考慮する。具体的には、SLP ディレクトリを各モジュールに重複して配置し、任意のモジュールが SLP ディレクトリ検索を実行可能とする。負荷の平準化方法を図 4.11 に示す。図 4.11 で機能別処理とは、サービスデータ処理と SLP ディレクトリ検索処理をそれぞれ別のモジュールで処理するとした方法である。SLP ディレクトリ検索処理はどのモジュールでも可能であり、ラウンドロビン方法等によりモジュールの負荷を一様とすることができる。しかし、サービスデータ処理は最も負荷の高いモジュールに合わせてモジュールの稼働率を設定する必要があり、全体としてモジュールの稼働率が低下する。負荷の平準化を行うためには、1 モジュールにサービスデータと SLP ディレクトリの両方を格納する必要があり、メモリやファイ



(1) 機能別処理



(2) 負荷平準化処理

平準化方法：サービスデータ処理と SLP ディレクトリ検索処理を組合わせ、サービスデータ処理の小さいモジュールに SLP ディレクトリ検索処理を多く割当てる。

: サービスデータ
(更新系データベースで分割して各モジュールに配備)

: SLP ディレクトリ
(参照系データベースで全てのモジュールに重複して配備)

図 4.1.1 負荷の平準化方法

ル容量は多く必要となるが、モジュールの稼働率を高くでき、機能別処理に比べて少ないモジュール台数でサービスデータ処理と SLP ディレクトリ検索処理を行うことができる。負荷の平準化がどの程度行えるかはサービスデータ処理と SLP ディレクトリ検索処理の処理量や、サービスデータ処理でのモジュール間の負荷偏り率に依存する。SLP ディレクトリ検索処理は IN 呼毎に必要であり、そのダイナミックステップ数はサービスデータに関するトランザクション処理のダイナミックステップ数と同程度である。IN 呼当たりのトランザクション数を 2~3 回とすると、SLP ディレクトリ検索処理の処理量はサービスデータに関するトランザクションの総処理量の 2~3 割程度となる。一方、モジュール間の負荷偏り率 ω は 1~2 割程度であり、方法 2 を採用した場合、SLP ディレクトリ分のメモリ量は増加するが、モジュール間の負荷を十分に平準化できると考えている。

4.6 結言

カスタマ毎のサービスデータへのトラヒックが数千倍から数万倍と大きくばらつく場合、カスタマ毎のトラヒックを高にものからの順位を用いて近似する方法を提案した。これに基づき、複数のモジュールにサービスデータを分散して配置し、大量のトランザクションを分散処理するシステムを対象として、モジュールに収容するサービスデータ数とモジュール間の負荷の偏りの関係や、その偏りの上限値の推定方法を明らかにした。これらの結果を用い、現状のプロセッサ性能、高度 IN サービスにおけるトランザクションのダイナミックステップ数などから、トラヒック量のばらつきが数千から数万倍と大きいサービスデータを無作為に抽出して各モジュールに配置したとしても、モジュール間の負荷の偏りが 1~2 割程度に収まることを示した。さらに、この偏りに対しては、SLP ディレクトリ検索処理とサービスデータに基づくデータベース処理を効率的に組み合わせることにより、負荷を動的に平準化できることを示した。

モジュール間の負荷の偏りの評価技術は、高度 IN のサービス制御ノードに限らず、OLTP システム等、データベースを分散処理するシステムの方式設計や設備設計の効率化に有効である。

なお、本章では、フリーダイヤルサービスのデータに基づいて負荷の偏りについて評価した。しかし、現時点では必ずしもこれらのデータが十分でないため、今後、サービス毎、カスタマ毎のトラヒック量の測定を拡充し、提案した近似式の検証をさらに進めてゆく必要がある。また、どのカスタマについてもトラヒックが集中する時間帯は同一とし、かつこの時間帯についてはトラヒックは一様に分布すると仮定して検討を進めた。しかし、すべてのカスタマについてトラヒックの集中する時間帯が同一になるわけではなく、サービス毎、カスタマ毎の時間変動を含めたトラヒック量の測定も拡充してゆく必要がある。これらについては、今後の課題と考えている。

第5章

分散処理によるサービス制御ノードの モジュール間結合技術

5.1 緒言

高度インテリジェントネットワーク（高度IN）のサービス制御ノード（SCP）では、小規模なサービスから全国的な大規模サービスまで柔軟に対応することや需要の不透明なサービスへ効率的に対応することが要求される。このような需要の変動にタイムリに対応可能とするためには、多数のモジュールからなる大規模な分散構成によりSCPを実現することが有効である⁽⁵⁾。分散処理システムを拡張性良く高信頼に実現するためには、モジュール間を接続する結合機構のデータ転送能力に制約が生ぜず高い信頼性が達成されることが要求される。また、結合機構の高信頼化に加えて、モジュール個々には高い通信処理能力が要求される。トランザクション処理の場合は、モジュール間で送受されるデータ長は数100バイトと短いが、送受を高頻度で行う必要がある。これを実現するため、通信処理を専用に行うハードウェア（以下、通信制御チャネルと記す）を適用し、この通信制御チャネルの単位時間当たり処理可能な通信回数を向上させることが重要となる。

本章では、分散構成を採るSCP内のモジュール間を結合する方式として、十分なデータ転送能力を確保することからスイッチ方式とし、モジュールとのインターフェースにはATMを適用する分散処理システムを対象として、結合機構の高信頼化構成法や、結合機構に接続されるモジュールの通信処理能力を向上させる制御方式について提案する^{(51)～(54)}。はじめに、ATM結合機構を各モジュールと接続する機能、信号を多重化／分離化する機能、スイッチングを行う機能により実現し、機能単位に二重化してモジュール数が増加した場合でも信頼度の低下を抑えるとともに、接続するモ

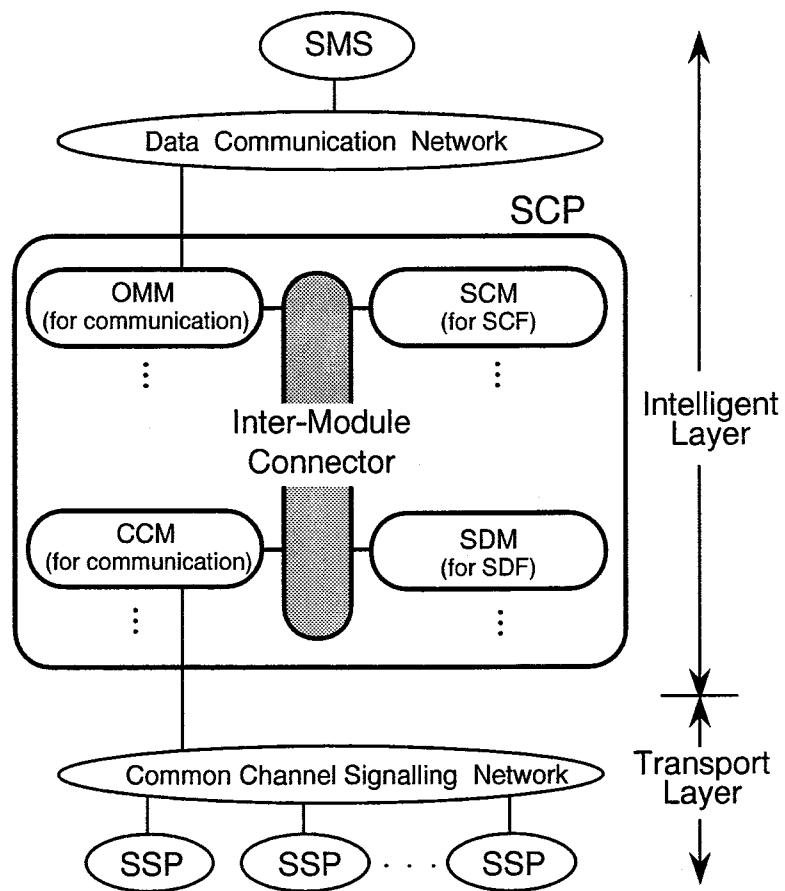
ジユール数に応じて機能単位に増設可能とした、高信頼で経済的な構成法を明らかにする。次に、モジュールの通信処理能力を向上させる制御方式として、ソフトウェアからの通信要求を主メモリにキューイングし、通信制御チャネルが一定周期毎にまとめて処理する方式を取り上げ、周期時間と応答時間の関係を明らかにし、高度 IN の SCP では数 10 ミリ秒の周期で処理することが有効であることを示す。さらに、以上の検討結果を踏まえて試作した分散処理システムの評価を行うため、モジュール間での通信に要する時間を測定する性能測定プログラムの構成法や測定結果について示す。

5.2 サービス制御ノード構成条件

5.2.1 ノード構成条件

高度 IN に適用する SCP を構成するためには、ITU-T で標準化^{(4) (1)} されているサービス制御機能 (SCF) やサービスデータ機能 (SDF) に加えて、交換機、サービス管理ノード (SMS) との通信機能が必要となる。機能間の分離性がよく、負荷変動に対して機能単位に柔軟に対応できること、実現する機能に最適化した技術を適用できることから、これらの機能を別々のモジュールに配備し、負荷の増大に対しては機能単位で負荷分散するノード構成を想定する⁽⁵⁾。すなわち、SCF はサービス制御モジュール (SCM)、SDF はサービスデータモジュール (SDM)、共通線信号網を介した交換機との通信は通信制御モジュール (CCM)、サービス管理ノードとの通信は運転、保守モジュール (OMM) で実現し、負荷の増大に対しては各機能を実現するモジュールの追加により対応する。SCP のハードウェア構成を図 5.1 に示す。

SCP の規模としては、パーソナル通信など全国規模のサービス、カスタマオリエンティッドな多様なサービスへの適用を考え、従来ノード⁽¹⁰⁾ に比べ大規模化し、



SMS : Service Management System
 SCP : Service Control Point
 SSP : Service Switching Point
 SDF : Service Data Function
 SCF : Service Control Function
 SDM : Service Data Module
 SCM : Service Control Module
 CCM : Communication Control Module
 OMM : Operation & Maintenance Module

図 5.1 SCP 構成

100万カスタマ程度を収容可能とする。また、SCPにかかるトラヒックは2000呼／秒程度を想定する。接続されるモジュール数は呼当たりのダイナミックステップ数やプロセッサ性能に依存し、約50モジュール程度を想定する。

5.2.2 性能条件

交換機に収容される電話機から発信された呼は一般的の電話呼かインテリジェントネットワークを利用する高機能な電話サービスを受ける電話呼（以下、IN呼と記す）かを判別され、IN呼の場合は共通線信号網を介してトランザクションがSCPに送られる。SCPはサービス毎、カスタマ毎にどのような内容のサービスを行うかのデータベースを有し、このデータベースの内容に従い、交換機と協調して高機能な電話サービスを実現する。本章で想定する機能分散と負荷分散を統合したSCPでは、交換機から送られてきたトランザクションは、まずCCMで受信され、結合機構を介してSCMへ転送される。SCMでは必要に応じて、SDMにアクセスし、データベースの参照および更新を行い、トランザクション処理結果をCCMを経由して交換機に送信する。

1つのIN呼を処理するためには少なくとも発信時と切断時にそれぞれ1回のトランザクションが発生する。また、サービス内容によっては、呼の接続処理中に交換機からSCPへの問合せのトランザクションが必要となることもある。IN呼当たりのトランザクション数としては2～4Tr／呼を想定する。また、1つのトランザクションを処理するために必要なモジュール間の通信回数は、サービス内容にも依存するが、5回／Tr程度である。従って、IN呼当たりのモジュール間の通信回数は10～20回／呼となり、単位時間当たり2000呼を処理するために必要なモジュール間の総通信回数は20000～40000回／秒程度となる。呼の処理のためにモジュール間で送受されるメッセージのデータ長は約200B程度であり、モジュール間の総通信量は30～60Mbps程度となる。SCPにかかるトラヒックとモジュール間の通信量の関係を図5.2に示す。

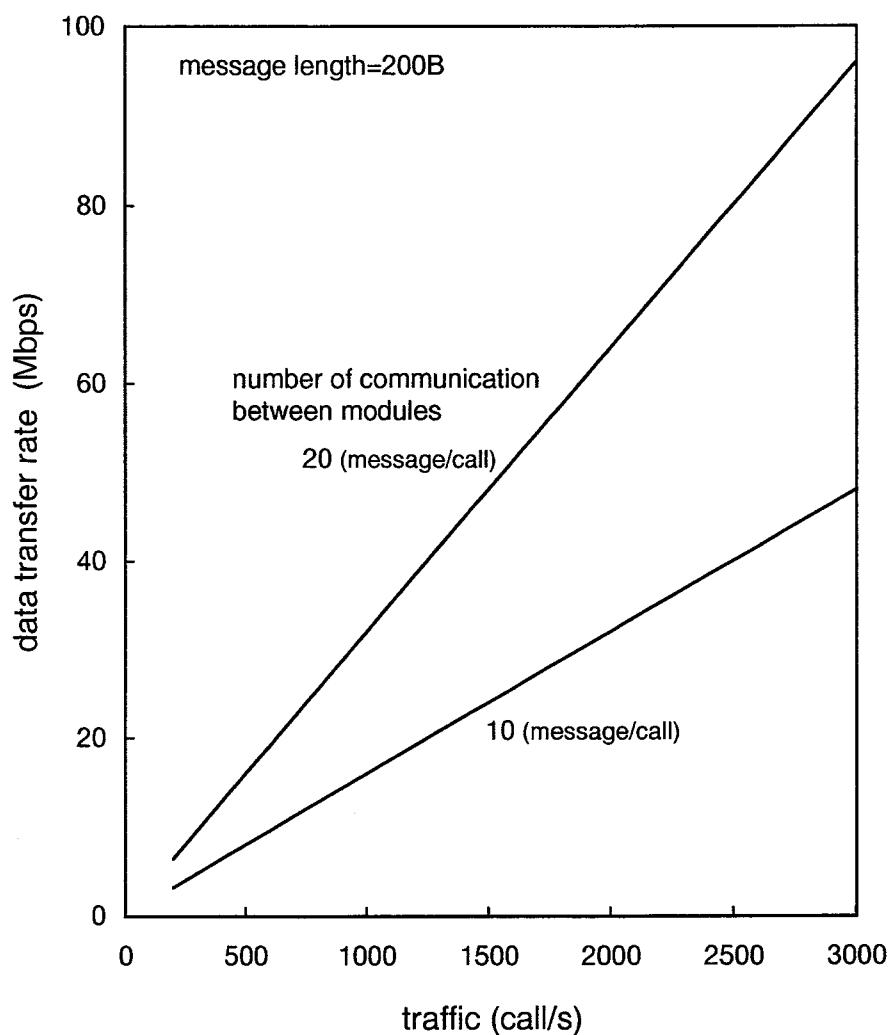


図5.2 トライックとモジュール間のデータ転送量

モジュールに要求される通信処理能力は、入出力を合わせると、総通信回数の2倍をモジュール台数50で割った値となり、1000～2000回／秒程度となる。

一般の電話呼の場合、ダイヤル終了から呼出し音が返るまでの時間が接続品質として規定されている。着信先の位置等の条件による変動はあるがおおよそ数秒から10秒程度である。高度INサービスの場合は交換機からSCPへの問合せが必要であり、一般の電話呼に比べて、呼出し音が返るまでに要する時間が、この問合せの分だけ余計にかかる。高度INサービスの接続品質を考慮すると、SCPへの問合せの遅延時間は1秒程度以下に抑える必要がある。この時間より交換機とSCP間の共通線信号網の遅延時間を差し引き、SCPノード内で許容される時間は300ミリ秒程度を想定する。

モジュール間の通信には、呼処理のための通信以外に、SMSからカスタマのデータをSCP内のモジュールにダウンロードするための通信や運用情報をSMSに送るための通信等がある。このようなダウンロードや運用のためのトラヒックは、呼処理のためのトラヒックに比べ2～3桁小さいため、モジュール間の総通信量や総通信回数としては、呼処理のための通信のみを考えておけば十分である。このため、本章ではダウンロード等のトラヒックについては対象外とした。

5.2.3 信頼度条件

高度INサービスは図5.1に示すように交換機の上位にSCPを配置する形で実現され、交換機とSCPのそれぞれの不稼働率が加算されてサービスに影響する。高度INサービスを一般の電話サービスと同程度の不稼働率でユーザに提供するためにはSCPの不稼働率が交換機の不稼働率に比べてはるかに小さく、無視できる程度であることが望ましい。交換機の信頼度条件としては20年間に1時間程度のシステムダウン（不稼働率 $\approx 6 \times 10^{-6}$ ）が設定されている⁽⁵⁵⁾。本章では、SCPの不稼働率を交換機の不稼働率より1～2桁低く想定し、 $10^{-7} \sim 10^{-6}$ 以下（20年間に1～10分以下のシステムダウン）を目標とする。SCPを構成するモジュールは相互バッ

クアップ等により信頼度上無視できる程度とすることができます。分散処理システムの信頼度のボトルネックはモジュール間の結合機構であり、結合機構の不稼働率の目標は $10^{-7} \sim 10^{-6}$ 以下とする。

5.3 モジュール間結合方式

5.3.1 モジュール間結合方式の選択

モジュール間の結合方式は、大きくバス／リング結合とスイッチ結合に大別できる。SCPのモジュール間の結合方式としては、以下の理由によりスイッチ結合とし、モジュールとスイッチ間はATMにより接続することとした。

OLT Pシステム等で用いられているバス／リング結合の代表的なものとしてFDDIがある。FDDIの最大データ転送速度は100Mbpsであるが、実効データ転送能力は電文長が短くなるに従って小さくなり^(5,6)、200B程度の短電文では10～30%程度である。SCPとしての所要データ転送能力(30～60Mbps)を確保するためには、複数のFDDIリンクを実現することが必要となり、負荷の分散や障害時のソフトウェア処理が複雑となる。また、FDDIリンクを制御する通信制御チャネルをモジュール毎に複数実装する必要がありコストアップにもつながる。一方、スイッチ結合の場合は高いデータ転送能力を実現することは容易である。たとえば、1入力当たり156Mbpsの転送能力を有する8×8のATMスイッチの場合、1.2Gbpsのデータ転送能力があり、転送能力に対する制限はほとんど無いと考えて良い。

また、モジュール間の通信をATMで実現しておくと、将来、共通線信号網を介したSSPとの通信や情報転送網を介したSMSとの通信等のノード間の通信がATM化された場合、モジュール間の通信とノード間の通信を同様に扱えるメリットがある。

5.3.2 ATM結合機構の高信頼化構成法

FDDIはモジュール間で通信を行うためのスイッチング機構が各モジュールに分散されており、1台のモジュール内のスイッチング機構が障害になったとしてもSCPノード全体の障害とはならない。一方、スイッチ結合の場合は、スイッチング機構が一ヵ所に集中するため、結合機構の障害が即座にSCPノード全体の障害につながる。このため、結合機構の高信頼化が必須となる。

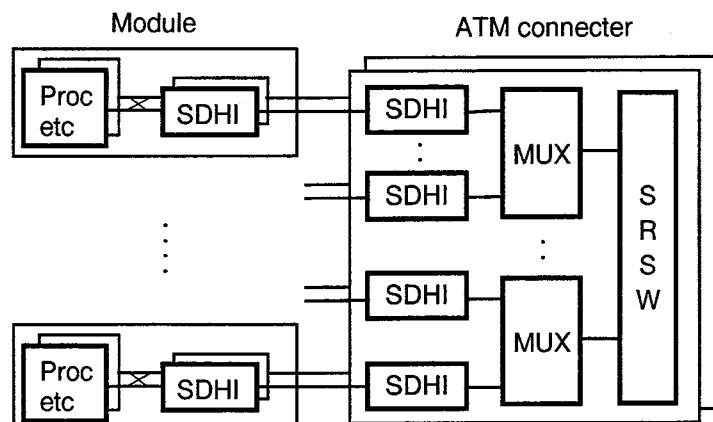
ATM結合機構は各モジュールと接続する光インターフェース部(SDH1)、各モジュールからのノードへの信号を多重化／分離化する部分(MUX)、ATMセルのルーティングを行うスイッチ部(SRSW)の機能ブロックから構成することができる。SDH1は接続するモジュール毎に必要であり、モジュール台数Nに比例して数が増加する。MUXは信号を多重化するモジュール数mに依存し、モジュール台数がNの場合、MUXの台数Kは次式で表せる。

$$K = \lceil (N-1)/m \rceil + 1 \quad (5.1)$$

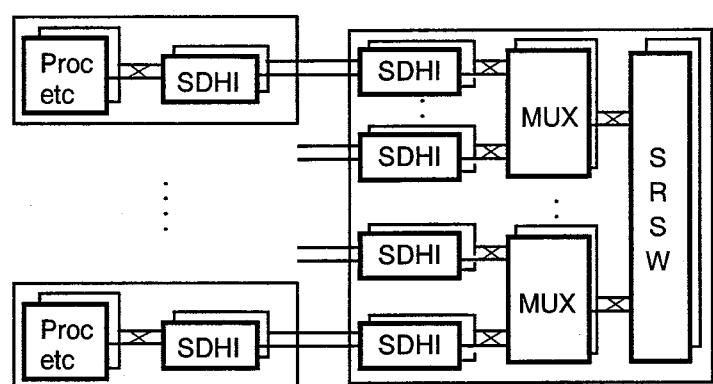
SRSWの増設は難しく、あらかじめ最大のノード規模に見合ったものを1台用意する必要がある。

ATM結合機構の高信頼化をはかるためには各機能ブロックを二重化する必要があるが、二重化の方法として、ATM結合機構全体として二重化する方法(方法1)とATM結合機構を構成する機能ブロック毎に二重化し、機能ブロック間で交差を設ける方法(方法2)を考えられる。両方法のATM結合機構の構成を図5.3に示す。ATM結合機構に接続するモジュールはプロセッサとSDH1から構成され、それぞれ二重化されているとする。また、プロセッサとSDH1間には交差が設けられており、現用系のプロセッサが障害となり予備系に切替わったとしてもSDH1やATM結合機構の切替えは不要とする。

方法1は多数のパソコン等を接続した企業内のLAN等で幅広く使用されている一重化構成のATM結合機構を2台用いて全体として信頼度を向上させる方法に相当する。方法1では現用系のSDH1、MUX、SRSWのどれか1台が故障した場



Method 1



Method 2

SDHI : Synchronized Digital Hierarchy Interface

MUX : Multiplexing and Demultiplexing

SRSW: Self-Routing Switch

図 5.3 ATM結合機構の信頼度構成

合でも ATM 結合機構の切替えが生ずる。また、モジュール側の SDHI の故障に対しても ATM 結合機構の切替えが必要となる。方式 2 の場合は、機能ブロック間で交差が設けられており、現用系の、ある機能ブロックが故障した場合は、当該機能ブロックのみ予備系に切替えることとなる。方法 1、2 の信頼度構成モデルを図 5.4 に示す。

方法 1、2 の不稼働率をそれぞれ UA_1 、 UA_2 とおくと、次式で表せる。

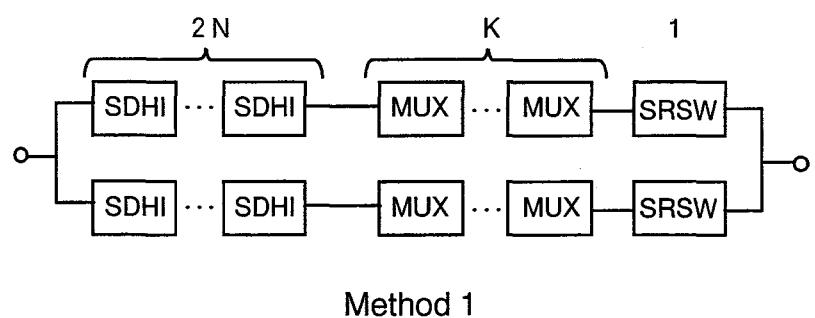
$$UA_1 = \frac{\lambda_1^2}{(\mu + \lambda_1^2)} \quad (5.2)$$

$$UA_2 = \frac{N(2\lambda_{SDHI})^2}{(\mu + 2\lambda_{SDHI})^2} + \frac{K\lambda_{MUX}^2}{(\mu + \lambda_{MUX})^2} + \frac{\lambda_{SRSW}^2}{(\mu + \lambda_{SRSW})^2} \quad (5.3)$$

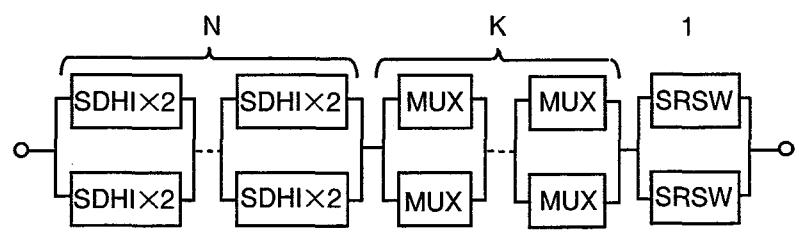
$$\text{但し、 } \lambda_1 = 2N\lambda_{SDHI} + K\lambda_{MUX} + \lambda_{SRSW} \quad (5.4)$$

である。ここで、 λ_{SDHI} 、 λ_{MUX} 、 λ_{SRSW} はそれぞれ SDHI、MUX、SRSW の故障率、 μ は修理率である。モジュール台数と不稼働率の関係を図 5.5 に示す。図 5.5において、試作したハードウェア量から λ_{SDHI} 、 λ_{MUX} 、 λ_{SRSW} はそれぞれ 2.5×10^{-5} 、 5×10^{-5} 、 10^{-4} とした。また、従来の交換機の修理状況を考慮し μ は 1 とした。方法 1 は図 5.5 からわかるように接続するモジュール台数の増加とともに不稼働率が増加し、要求される信頼度条件 ($10^{-7} \sim 10^{-6}$ 以下) のもとでは、接続可能なモジュール数は数台～10数台に制限され、大規模ノードを構成することは難しい。一方、方法 2 は機能ブロック間で交差を設けるため、モジュール台数が増加しても ATM 結合機構の不稼働率の増加は小さく、50台程度のモジュールを接続した大規模システムにおいても要求される不稼働率を満足する。

方法 1、方法 2 ともに結合機構を構成する機能ブロックは同一であり、機能ブロックのハードウェア量の差はほとんどない。しかし、方法 2 の場合は 0 系と 1 系の機能ブロック間に交差を設けるため、機能ブロック間を結合するケーブル類が多く必要となる。このため、結合機構のコストは方法 2 が方法 1 に比べて若干高くなる。しかし、システム全体で見るとモジュールのコストが支配的であり、結合機構のコスト差に基づくシステムコストの差はほとんどない。また、結合機構のスループットやレスポンスタイムの性能については方法 1 と方法 2 で差はない。



Method 1



Method 2

図5.4 ATM結合機構の信頼度ブロック図

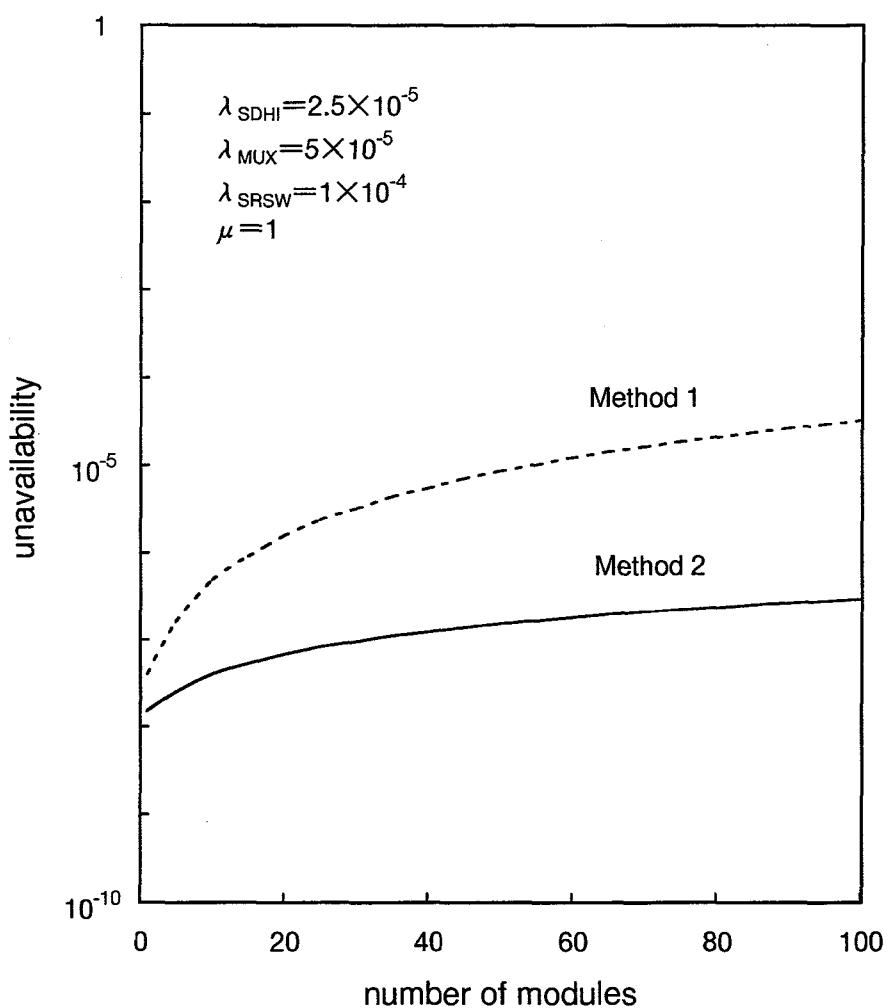


図 5.5 ATM 結合機構の不稼働率

以上、信頼度面では方法2は大規模構成においても要求される信頼度条件を満足するが、方法1はモジュール数が10数台以上になると要求条件を満足しない。一方、コスト面では方法1が方法2に比べて若干有利である。一般に規模が大きくなるに従ってシステムとして要求される信頼度条件は高くなる^(5,7)。このことから、ATM結合機構としては、機能単位に交差を設け、接続するモジュール台数が増加した場合でも信頼度の低下度合いが小さい方法2を採用することとした。

方法2において、機能ブロック間の交差の実現方法を図5.6に示す。図5.6はSDH IとMUX間の交差の例を示しており、実線で囲った機能ブロックは現用系、破線で囲った機能ブロックは予備系を示す。現用系、予備系のSDH Iから信号が両系のMUXに供給され、MUXでは両方の信号のどちらか一方を選択する回路（図5.6のS）を有し、通常は現用系のSDH Iからの信号を選択する。現用系のSDH Iが故障した場合には予備系のSDH Iの信号へ切替える。また、MUX内に両系のSDH Iからの信号の途絶や同期ずれを検出する監視回路（モニタ）を設け、現用系のSDH Iからの信号に異常を検出した場合にはモニタが自律的に選択回路を予備系のSDH Iからの信号に切替える。図5.6には示していないが、MUXとSR SW間の交差も同様の方法で実現した。

5.3.3 ATM結合機構の全体構成

高度INサービスは小規模から大規模まで幅広く分布する。このため、要求に応じてモジュールを追加し、柔軟に処理能力を拡張できることが要求される。併せて、接続するモジュール数が少ない場合、結合機構のコスト負担が少ないと要求される。このことから、ATM結合機構を機能ブロック単位に増設可能とし、SDH Iについては1モジュール単位に、MUXについては8モジュール単位に、接続するモジュール台数に応じて増設可能とした。

システムを立上げる場合は、ATM結合機構を構成する機能ブロックの初期化やルーティング情報の設定を行う必要がある。また、モジュールを増設した場合にもルーチ

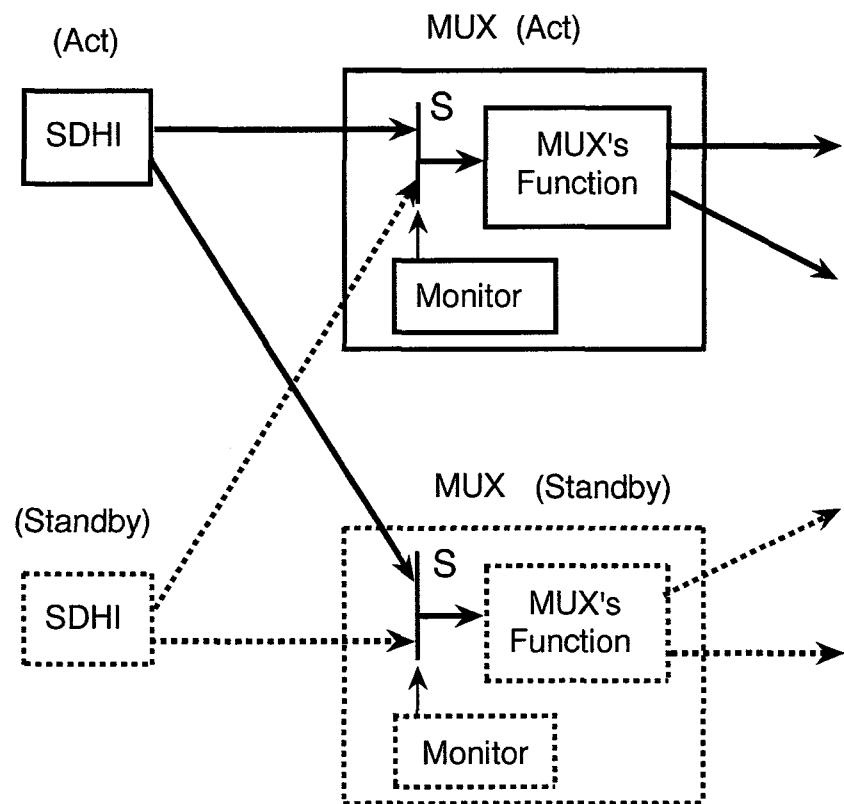


図 5.6 機能ブロック間の交差の実現法

ング情報の再設定が必要となる。このような処理は一度行えば変更されることは無いため、通常の呼処理でプロセッサに負荷がかかることはない。このため、これらの処理を行うプロセッサは、SMSとの通信を行うモジュール（OMM）と共にすることとし、ATM結合機構をこのモジュールと一体化し、システムインターフェースモジュール（SIM）として実現した。ATM結合機構およびプロセッサを含むSIMの全体構成を図5.7に示す。

SIMの実装に関しては、一つの架の中にプロセッサやSRSW、MUX、SDHIを搭載し、この架で24モジュールを接続可能とした。また、より多くのモジュールを接続可能とするため、別の架にMUXやSDHIを搭載し、最大で88モジュールまで接続可能とした。SIMの実装構成を図5.8に示す。各モジュールとATM結合機構とのインターフェース速度は、200B程度のメッセージを1000～2000回／秒の頻度で送受した場合でも十分対応可能とするため6Mbpsとした。

5.4 通信制御方式

ATM結合機構に接続され、モジュールの通信処理を専用に行う通信制御チャネルの制御方式や応答時間について考察する。

5.4.1 制御方式

通信制御チャネルの制御方式として即時処理と周期処理がある。即時処理とは、ソフトウェアでデータ転送の要求が発生すると同時に通信制御チャネルに起動をかけ、通信制御チャネルでは起動内容に応じた処理を即座に実行し、処理終了と同時にソフトウェアへ割込みにより通知する方式である。応答時間条件の厳しい磁気ディスク装置等を接続する入出力チャネルでは一般に即時処理が採用されている。一方、周期処理とは、ソフトウェアでデータ転送の要求が発生すると、その転送要求をメモリ上の

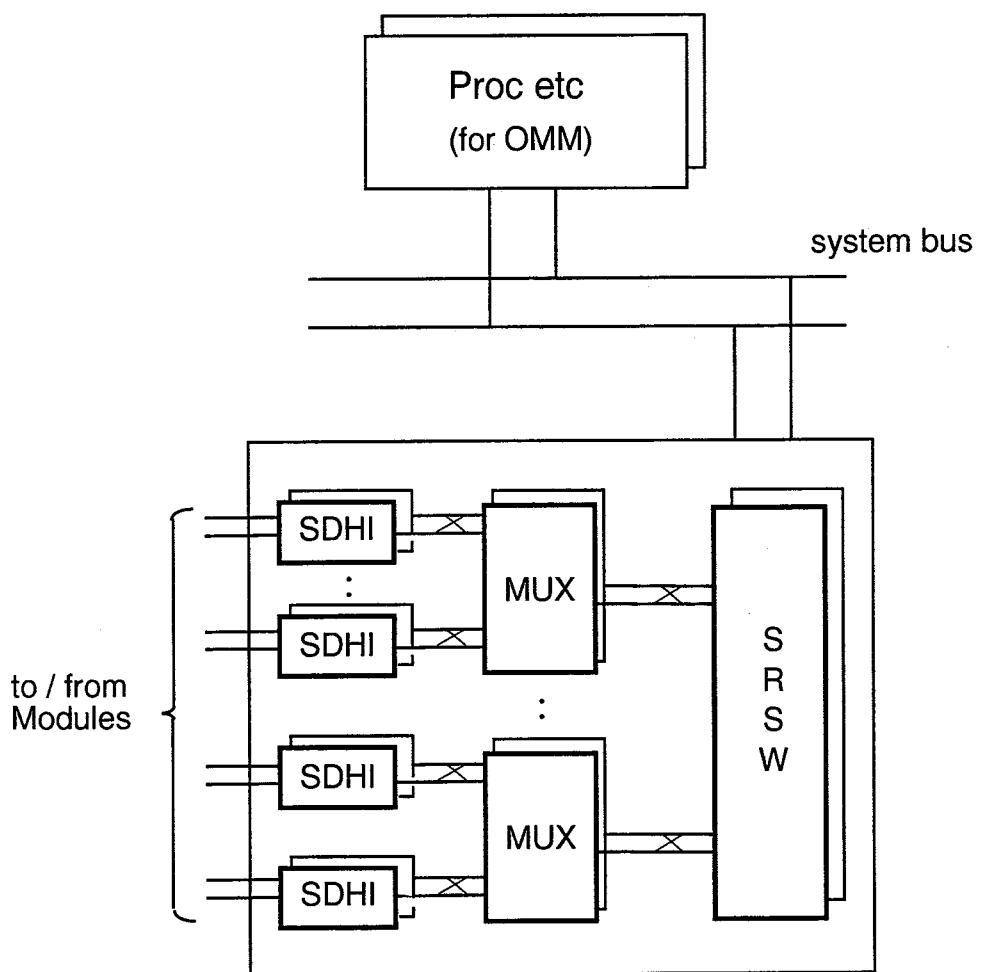


図5.7 SIMの全体構成

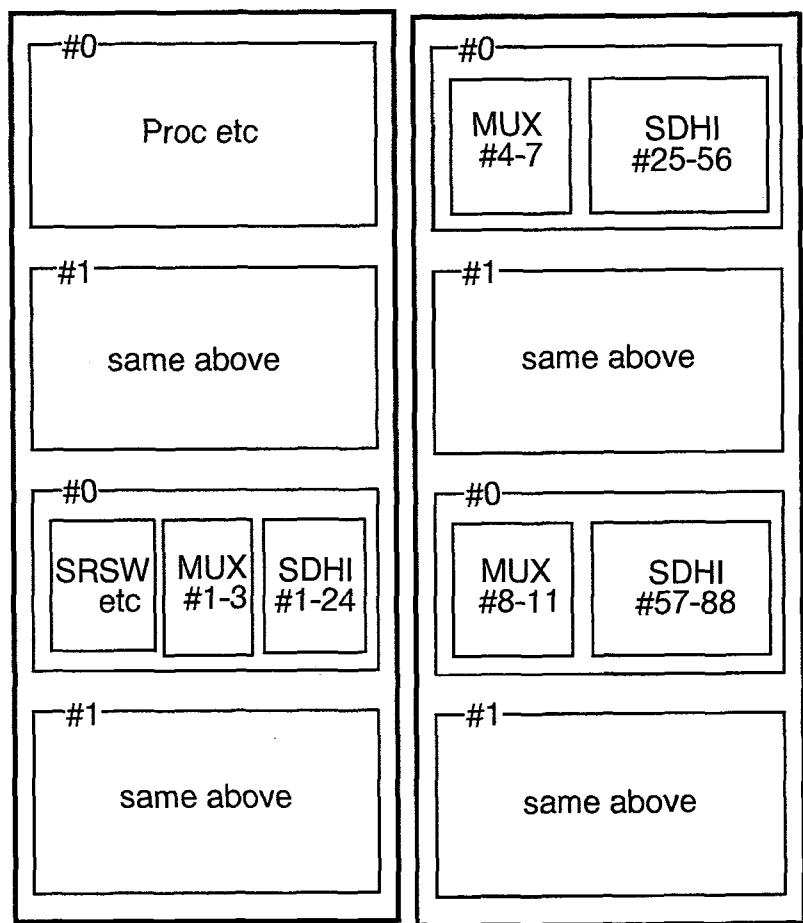


図5.8 SIMの実装構成

キューに格納するのみで、通信制御チャネルには起動をかけない方式である。通信制御チャネルは、あらかじめ決められた周期でメモリ上のキューをルックインし、処理要求の有無を調べ、それに基づいて処理する。処理が終了するとメモリ上のキューにその旨を書き込む。ソフトウェアから非同期に要求される起動処理やソフトウェアへの終了割込み処理が不要となること、メモリ上のキューから複数の転送要求を読出しが可能であることから、スループットの向上が期待できる。本章で対象とするモジュールの通信制御チャネルには1000～2000回／秒の高い通信処理能力が必要とされるため、通信制御チャネルの制御方式として、周期処理を採用することとした。周期処理の場合、周期時間を長く設定すると一度に多くの通信要求を処理できるためスループットは向上するが、転送要求がメモリ上のキューに書込まれてからルックインするまでは待たされるため、ソフトウェアからの転送要求から通信制御チャネルでの処理が終了するまでの時間（以下、応答時間と記す）が長くなることが懸念される。周期時間の設定に当たっては応答時間の評価が必要となる。これに関しては、まず、SCPに要求されるトランザクション処理時間から通信制御チャネルの応答時間への要求条件を明らかにし、次に簡単な单一キュー、単一窓口の待ちモデルにより、周期時間と応答時間の関係を以下で評価する。

5.4.2 通信制御チャネルの応答時間条件

SCPノード内でのトランザクション処理時間 T_{RES} は、プロセッサの処理時間 T_{CPU} 、モジュール間の通信時間 T_{COM} 、トランザクション当たりのモジュール間通信回数 f により、次式で表せる。

$$T_{RES} = T_{CPU} + f \cdot T_{COM} \\ = \frac{D}{(1-\rho)P} + f \cdot (2T_{CCH} + T_{ATM}) \quad (5.5)$$

ここで、Dはトランザクション当たりのダイナミックステップ数、Pはプロセッサの性能、 ρ はプロセッサの使用率、 T_{CCH} は通信制御チャネルの応答時間、 T_{ATM} はメッセージがATM結合機構を通過する時間である。なお、モジュール間の通信に要する

時間 T_{com} は図 5.9 に示すように、送信モジュール側の通信制御チャネルの応答時間、メッセージが ATM 結合機構を通過する時間、受信モジュール側の通信制御チャネルの応答時間の和で表せるとした。また、通信制御チャネルの応答時間は、送信、受信で同一とした。

5.2 節で述べたように、 T_{res} に許容される時間は約 300 ミリ秒程度である。 ρ を 80 %、D を 100 k ステップ、P を 20 MIPS とすると、 T_{cpu} は 25 ミリ秒となる。 T_{atm} はメッセージのデータ長（約 200 B）、モジュールと ATM 結合機構とのインターフェース速度（6 Mbps）で決まり 1 ミリ秒程度である。トランザクション当たりのモジュール間通信回数 f は 5 回 / T_r であり、 T_{cch} に許容される時間は 25 ミリ秒程度となる。

5.4.3 評価モデル

通信制御チャネルへの通信要求はポアソン到着とし、窓口は一定周期 T 毎に開き、キューから周期時間内に処理可能なメッセージ数 b を取り出して処理する。キュー内のメッセージ数が周期時間内に処理可能なメッセージ数より少ない場合は、キューからすべてのメッセージを取り出して処理する。このとき、取り出した全てのメッセージの処理が終了したとしても周期途中で再度窓口を開いて、新たに発生したメッセージを処理することはしない。キュー内のメッセージ数が周期時間内に処理可能なメッセージ数より多い場合は b 個のメッセージのみ処理し、残りは次の周期で処理する。周期 T 内に発生したメッセージ数を n 、処理残りのメッセージ数を h とすると、周期の終了時点ではキューの中に $h + n$ 個のメッセージがある。評価モデルを図 5.10 に示す。

5.4.4 通信制御チャネルの応答時間評価

周期の始めに行うキュー内のメッセージの検索等周期内で処理するメッセージの

Module #A Module #B

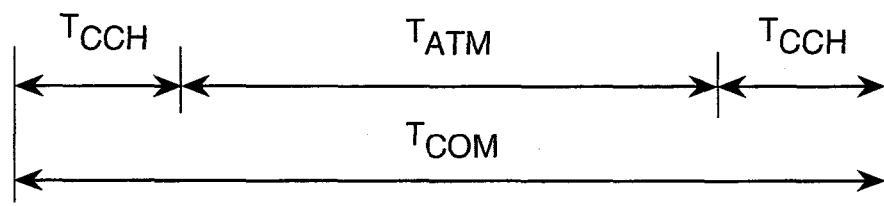
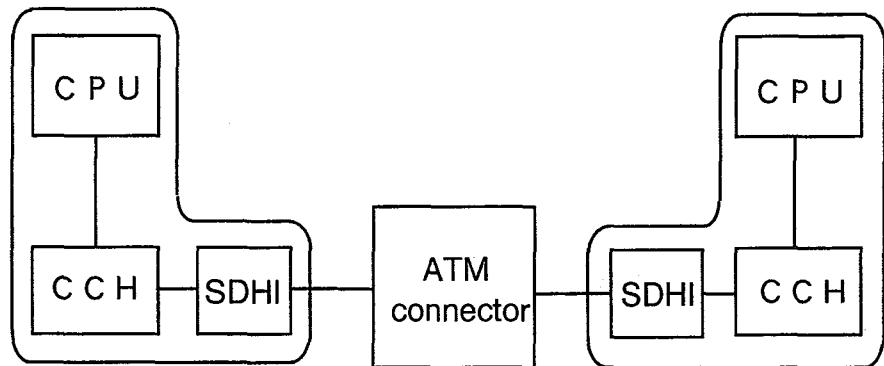


図 5.9 モジュール間の通信時間

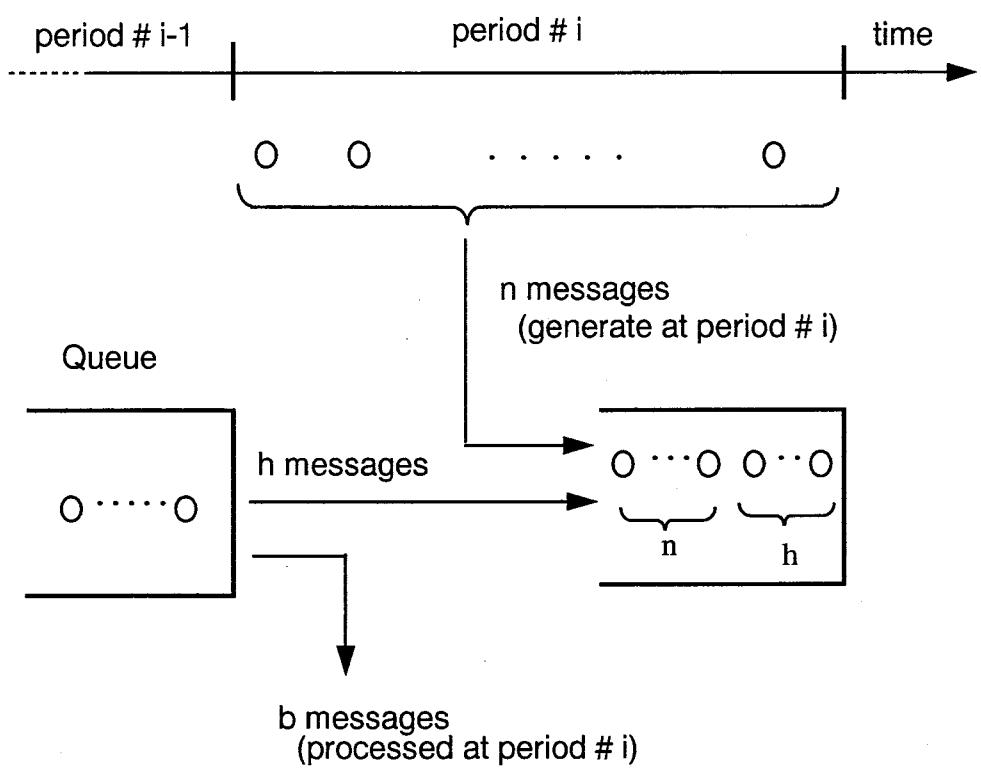


図 5.1.0 周期駆動モデル

数に依存せず必要な固定分の処理時間を t_0 、メッセージの処理に必要な処理時間を t_1 とおくと、周期 T 内に処理可能なメッセージ数 b は

$$b = \frac{T - t_0}{t_1} \quad (5.6)$$

となる。固定分 t_0 を周期内で処理するメッセージで一様に分担するとしたときのメッセージの平均処理時間 s は

$$s = \frac{T}{b} = t_1 + \frac{t_0 \cdot t_1}{T - t_0} \quad (5.7)$$

となる。以下では簡単化のため、周期 T 内においてメッセージは平均処理時間 s で処理されるものとして評価する。T を大きくすることにより s は小さくなりスループットは向上するが、発生したメッセージは周期終了まで待たされたため応答時間は長くなる。n 個のメッセージについて、発生から処理が終了するまでの平均時間（以下、平均応答時間と記す）を $t_{CCH}(n, h)$ とする。n 個のメッセージは周期の終了時点まで待たされ、次の周期では、まず、処理残りの h 個のメッセージの処理が終了した後、キューに登録された順に処理されるとする。平均応答時間 $t_{CCH}(n, h)$ はメッセージの発生時点から周期終了までの平均待ち時間 T_q 、処理残りの h 個のメッセージの処理時間 T_h 、n 個のメッセージの平均処理終了時間 T_n の和となり、次式で表せる。

$$t_{CCH}(n, h) = T_q + T_h + T_n \quad (5.8)$$

ここで、n 個のメッセージは周期 T 内でランダムに発生するため T_q は $T/2$ 、1 つのメッセージの処理時間は s であり T_h は $s \cdot h$ となる。また、n 個のメッセージのうち、1 番目のメッセージの処理が終了する時間は s、2 番目のメッセージの処理が終了する時間は $2s$ 、n 番目のメッセージの処理が終了する時間は $n \cdot s$ であり、平均処理終了時間 T_n は $(n+1)s/2$ となる。従って、 $t_{CCH}(n, h)$ は次式となる。

$$t_{CCH}(n, h) = \frac{T}{2} + s \cdot h + \frac{(n+1)s}{2} \quad (5.9)$$

周期内に n 呼発生する確率を P(n)、処理残り呼数が h である確率を X(h) とおくと、すべてのメッセージの平均応答時間 T_{CCH} は次式で表せる。

$$T_{CCH} = \sum_{h=0}^{\infty} \sum_{n=0}^{\infty} \{n \cdot X(h) \cdot P(n) \cdot t_{CCH}(h, n)\} / \sum_{n=0}^{\infty} \{n \cdot P(n)\} \quad (5.10)$$

式 (5.10) に式 (5.9) を代入し次式を得る。

$$T_{CCH} = \frac{T}{2} + s \cdot \bar{X} + \frac{s}{2a} \cdot \bar{P} \quad (5.11)$$

$$\text{ここで、 } \bar{P} = 1 \cdot 2 \cdot P(1) + 2 \cdot 3 \cdot P(2) + 3 \cdot 4 \cdot P(3) + \dots \quad (5.12)$$

$$\bar{X} = X(1) + 2X(2) + 3X(3) + \dots \quad (5.13)$$

である。単位時間当たりの平均メッセージ数を A とおくと、周期 T 内で発生する平均メッセージ数 a および P (n) は次式で表せる。なお、X (h) の算出方法は付録 5.1 に示す。

$$a = A \cdot T \quad (5.14)$$

$$P(n) = \frac{e^{-a} \cdot a^n}{n!} \quad (5.15)$$

A と T_{CCH} の関係を図 5.1.1 に示す。図 5.1.1 において、試作した通信制御チャネルで走行するファームウェアのダイナミックステップ数等から t₀ は 0.5 ミリ秒、t₁ は 0.5 ミリ秒とした。図 5.1.1 より、周期 T を長くすることにより、A は向上するが、向上度合いは T が長くなるに従って飽和する。たとえば、T が 16 ミリ秒以上では A はほぼ飽和する。一方、T_{CCH} は T が長くなるに従って長くなる。T が 16 ミリ秒では、A が 1800 (呼/秒) で T_{CCH} は約 20 ミリ秒程度であり、A をこれ以上大きくすると T_{CCH} は急激に増加する。

以上の検討結果から、通信制御チャネルの制御方式として、16 ミリ秒で周期処理することとした。このとき式 (5.5) より、モジュール間の通信時間 T_{COM} は 40 ミリ秒程度、トランザクション処理時間 T_{RES} は 225 ミリ秒程度となり、5.2 節で述べた要求条件 (300 ミリ秒) を満足する。

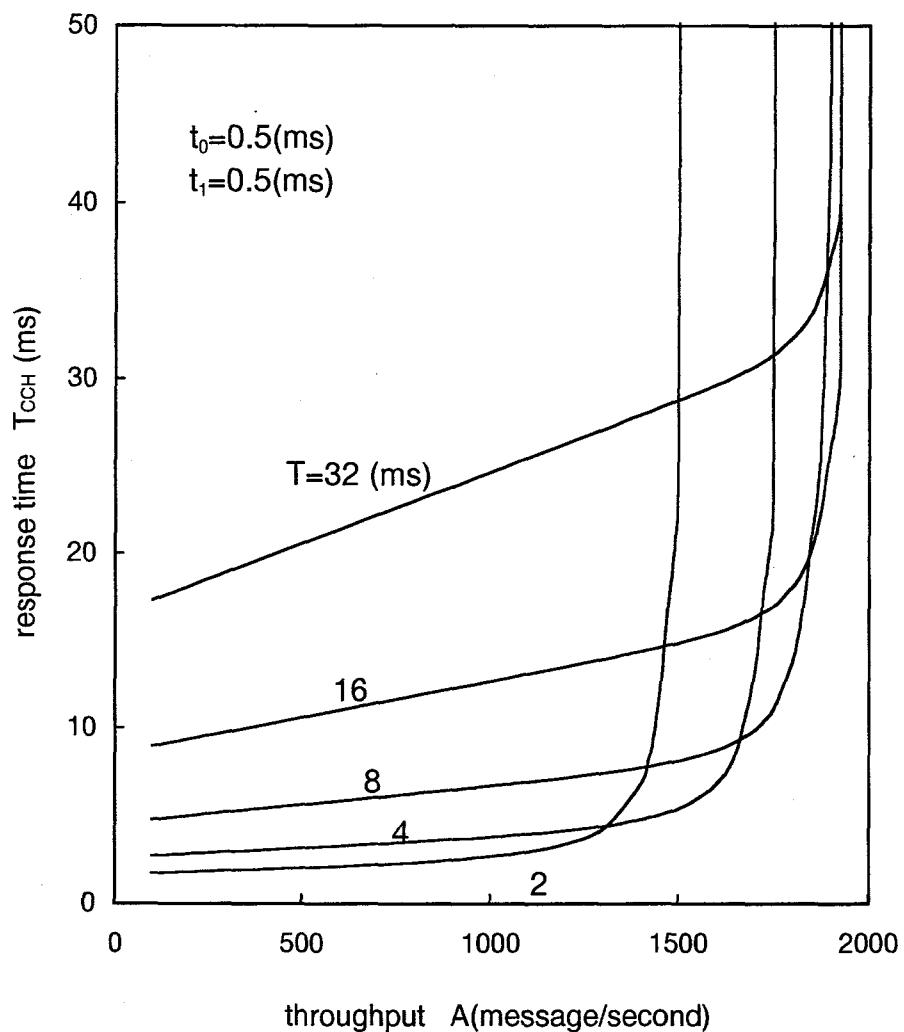


図5.1.1 スループットと応答時間

5.5 評価方法

本節では、前節での検討結果の検証を目的として作成したモジュール間の通信時間を測定する性能測定プログラムの構成法および測定結果について考察する。

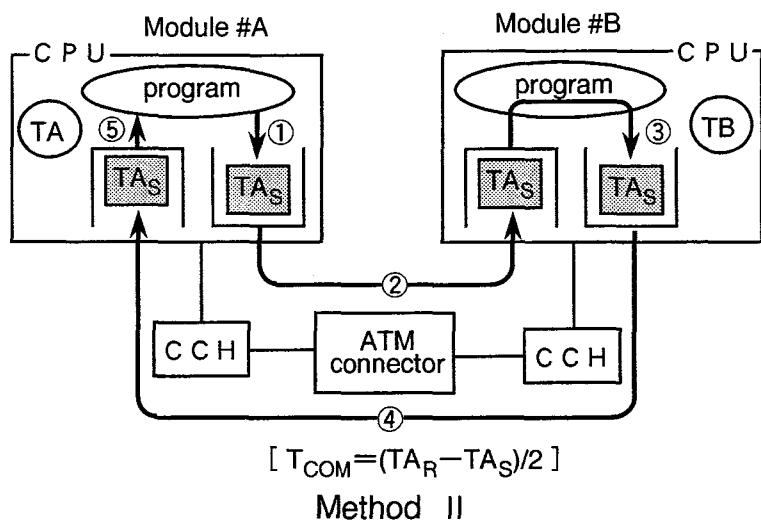
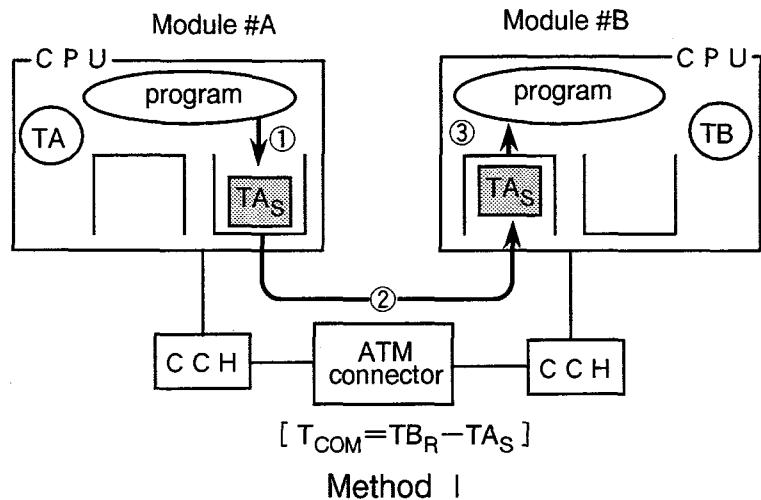
5.5.1 性能測定方法

試作した装置上で走行するプログラムによりモジュール間でメッセージの送受信を行い、通信に要する時間 T_{COM} を実測評価する。通信時間の測定方法として以下の2案が考えられる。

【方法 I】モジュールAでキューに送信コマンドを登録するとき、メッセージのデータ部に登録時刻を書き込む。このメッセージがモジュールBで受信された時刻と、メッセージに埋め込まれた時刻の差によって通信時間を測定する。

【方法 II】モジュールAでキューに送信コマンドを登録するとき、メッセージのデータ部に登録時刻を書き込む。当該メッセージがモジュールBで受信されると、当該メッセージのアドレスをモジュールA宛に変更して送り返す。モジュールAでメッセージを受信した時刻とメッセージに埋め込まれた時刻との差により通信時間を測定する。この場合、往復の時間をカウントしたこととなり実際のモジュール間の通信時間はこの値の1／2となる。

方法I、IIの測定方法を図5.12に示す。方法Iの場合は、2つのモジュールのタイマを利用して通信時間を測定するため、タイマ値の差により誤差が生ずる。モジュール間でのメッセージの通信に要する時間は数10ミリ秒程度であり、2つのモジュールのタイマ値をこれに見合った精度で一致させることは非常に難しい。一方、方法IIの場合は、1つのモジュールのタイマを利用して通信時間を測定するため、高い



TAs : Message from Module #A to Module #B

: Queue on main memory

TA : Timer in Module #A

TB : Timer in Module #B

T_{COM} : Communication time between Modules

TA_S : Time of TA when the message was sent by Module #A

TB_R : Time of TB when the message was received by Module #B

TA_R : Time of TA when the message was received by Module #A

図 5.1.2 モジュール間の通信時間測定方法

精度で測定できる。また、モジュール間でタイマ値に差が生じたとしても測定に誤差は生じない。このことから、本性能測定プログラムでは方法Ⅱを採用することとした。

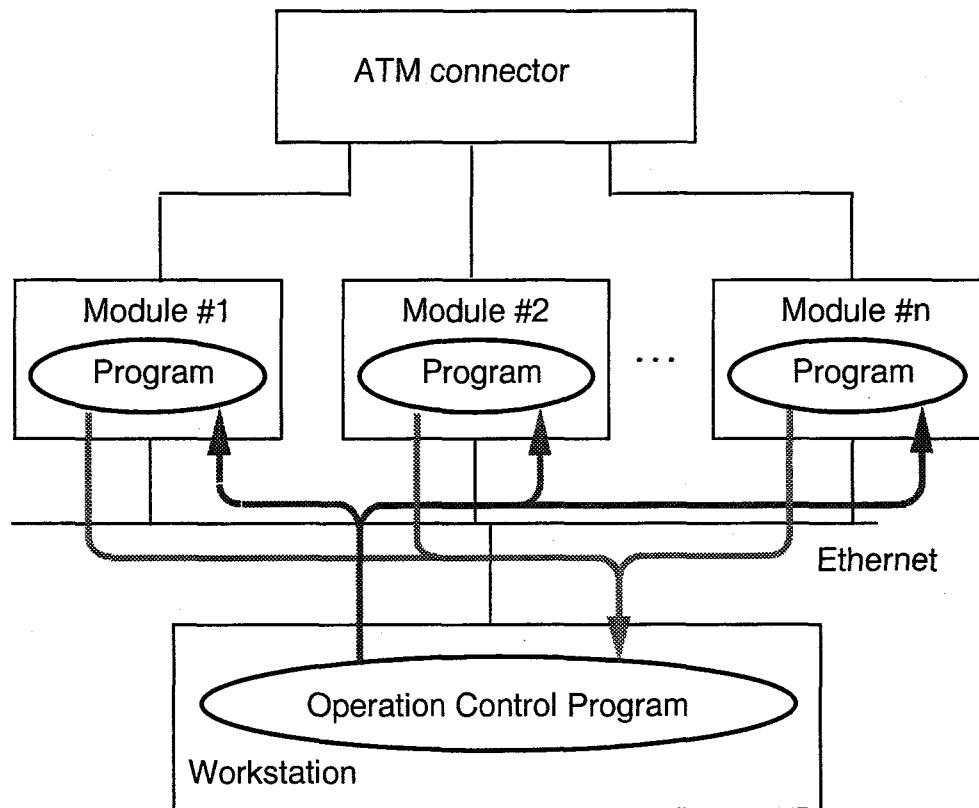
5.5.2 性能測定プログラムの走行方法

性能測定プログラムを各モジュール上で走行させるためには、各々のモジュールに対して、プログラムのローディングや走行指示が必要となる。また、各モジュール上の性能測定プログラムを同期をとって走行させる必要がある。方法として以下の2案が考えられる。

【方法A】モジュール毎のパネルやコンソールから行う方法

【方法B】ワークステーションと各モジュールを接続し、ワークステーションから行う方法

方法Aはモジュール毎に人を配置する必要があり、モジュール数が少ない場合はよいが、モジュールの増加とともに、操作が大変となる。また、モジュール毎に細かい同期を取ることが難しくなる。方法Bはワークステーション上で新たなプログラムを開発する必要があるが、モジュールが増加しても一人で走行させることが可能であり操作性がよい。また、モジュール間の同期を取ることも容易にできる。本章では、モジュール数が多い大規模ノード構成時でも一人で操作できることを考慮し、方法Bを採用することとした。ワークステーション上のプログラム規模としては、10キロライン程度で実現できた。各モジュールとワークステーションはイーサネットで接続し、一台のワークステーションから性能測定プログラムのローディング、起動および測定結果の表示等を行うこととした。性能測定プログラムの走行環境を図5.13に示す。



← · Program Loading
· Execution Control

← · Results of Execution

図 5.1.3 性能測定プログラムの走行環境

5.5.3 性能測定結果

性能測定プログラムを試作した装置上で走行させて実測した通信時間の分布を図5.14に示す。図5.14では転送するデータ量を振らせて測定した。図5.14より T_{COM} の平均は約40ミリ秒であり、ほぼ5.4節での評価通りの値となった。また、データ長が256Bでも90%以上の通信は50ミリ秒以内に収まることが判明した。データ転送量が大きくなるに従って通信時間が長くなる理由はATM結合機構の通過時間や、5.4節では考慮しなかったが通信チャネルと主メモリ間のデータ転送時間がデータ長に比例して増加するためである。

5.6 結言

大規模分散処理システムの結合機構の高信頼化構成法、通信処理能力を向上させる通信制御チャネルの制御方式およびモジュール間の通信時間の測定方法について提案した。モジュール間の接続にATMを適用し、結合機構を各モジュールと接続する機能、信号を多重化／分離化する機能、スイッチングを行う機能により実現し、機能単位に二重化した高信頼化構成法を明らかにした。これにより、結合機構に50モジュールを接続した大規模分散構成においても、不稼働率が $10^{-7} \sim 10^{-6}$ の高い信頼性を実現できた。また、各モジュールの通信制御については、ある周期毎にまとめて複数のメッセージを処理する方式を取り上げ、周期時間と通信時間の関係を明らかにし、高度INのSCPでは数10ミリ秒の周期で処理することが有効であることを示した。ここで示したモジュール間の結合方式は、高信頼で経済的大規模分散処理システムを設計する上で有用である。本章では、高度INに適用するサービス制御ノードを対象として検討を進めたが、ここで提案した方式は、適用業務の拡大から多数のモジュールによる分散処理化が進んでいるOLTPシステムにおいても、同様に適用できる。

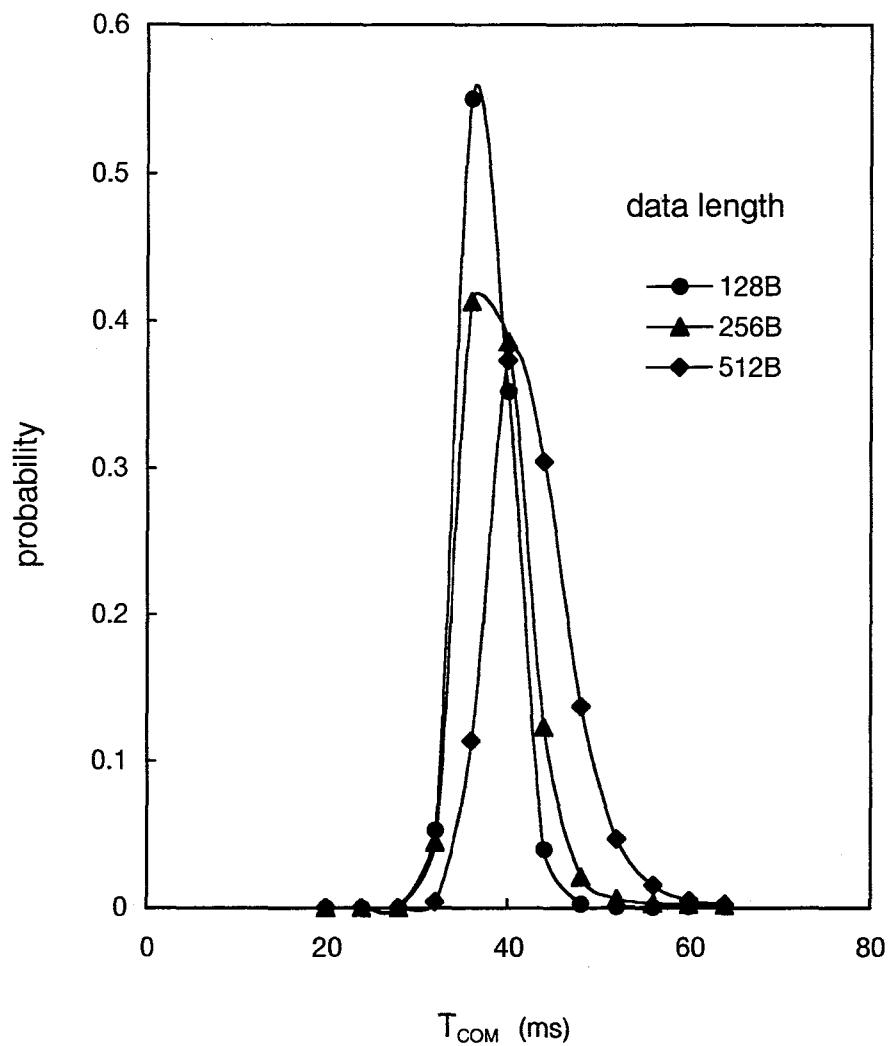


図 5.1.4 モジュール間の通信時間の分布

第6章 サービス制御ノードにおける メモリデータベースのリカバリ技術

6.1 緒言

高度INのSCPでは、多数のカスタマの情報から構成されるデータベースをもとに高度な通信サービスを実現するため、データベースを効率よく高速に処理できることが重要であるが、同時に障害に対しても、これらデータベースの内容が十分な確度で保証されることが要求される。また、高度INのデータベースは容量が比較的小さく、アクセス頻度が高い特徴を有し、トランザクション処理性能の向上のためには、第3章で述べたようにデータベースのすべてを主メモリ上に常駐するメモリデータベースの適用が有効である。

既に、銀行のオンラインシステム、クレジット照会システム等のOLTP (On-Line Transaction Processing) システムでもサービスの広域化、対象業務の拡大を背景に、より大量のトランザクションを高速に処理する必要がでてきており、主メモリの大容量化の進展と相まって、メモリデータベース化が進められている。しかし、メモリデータベースは、ソフトウェアのバグ、ハードウェア障害および電源障害により、主メモリ上のデータベースが破壊、喪失されやすいという欠点を併せ持っている。このため、あらかじめリカバリに必要な情報を障害に強い記憶装置上に格納しておき、障害時にはこの記憶装置上のリカバリ情報をもとに障害直前のデータベースを復旧する方式を併せて採用する必要がある^{(5) (28) (29)}。外部の記憶装置の形態はLOGを単に記憶するパッシブな方式とLOGをプロセッサで処理可能なアクティブな方式に2大別される。前者を代表し高速な入出力が可能な半導体ディスク装置を適用する方式を設定し、後者については別のモジュールのプロセッサに付属する主メモリを適用する方式を代替案として設定する。後者的方式では、当該の主メモリ上でLOGを逐一データベースに反映する方法をとると、リカバリ時間の大幅な短縮が期待できる。

本章では、2つの代替案を中心に、データベースへの単位時間、容量当たりのトランザクション数やりカバリ時間、モジュール毎に収容可能なデータベース容量に注目し、両方式の適用領域を明らかにする。また、電源障害時の不揮発化方法、冗長化構成法など信頼度構成面から考察し、サービス制御ノードのメモリデータベースのリカバリ方式としては、バッテリバックアップによる半導体ディスク装置を用いたりカバリ方式が有効であることを示す。ここで用いるバッテリバックアップについては、試作を通じて、保守の容易な方法を明らかにする等、全体として高度INで実用性の高いリカバリ方式を提案する^{(58)～(62)}。ここで示したリカバリ方式はメモリデータベース化が進められている多くのOLTPシステムにおいても有用と考えている。

6.2 前提条件

高度INサービスの特性に基づいて、本章で前提とした条件は以下の通りである。

- (1) 高度INのデータベースは容量が小さくアクセス頻度が極めて高い。パーソナル通信サービスを例にとると、カスタマ対応のデータとしては位置情報、暗証番号、スクリーニング情報等からなり数100バイト程度である。トランザクションとしてはカスタマからの発信、カスタマへの着信及び着信先を指定する位置登録などがあり1時間当たり3回程度、すなわち、カスタマ当たりの平均トランザクション数は 10^{-3} (Tr/秒) 程度であり、データベースへの単位時間、容量当たりのトランザクション数（以下、トランザクション密度と記す）は1～10 (Tr/MB/秒) 程度となる。データベースのサイズはモジュールに収容するカスタマ数に依存し、プロセッサ性能やトランザクション当たりのダイナミックステップ数等によってきまる。収容するカスタマ数を数10万程度とし、データベース容量としては100MB程度を想定する。

- (2) 一般の電話呼の場合、ダイヤル終了から呼出し音が返るまでの時間が接続品質として規定されている。着信先の位置等の条件による変動はあるがおおよそ数秒から 10 秒程度である。高度 IN サービスの場合は交換機から SCP への問合せが必要であり、一般の電話呼に比べて、呼出し音が返るまでに要する時間が、この問合せの分だけ余計にかかる。高度 IN サービスの接続品質を考慮すると、SCP への問合せの遅延時間は 1 秒程度以下に抑える必要がある。この時間より交換機と SCP 間の共通線信号網の遅延時間を差し引くと SCP ノード内で許容される時間は数 100 ミリ秒、データベース処理に許容される時間は 100 ミリ秒程度である。
- (3) 高度 IN でトランザクション処理中に更新されるデータとしてはパーソナル通信サービスの位置情報、回線が使用中か否かを表示するフラグなどがある。前者はいったん登録されたデータが失われると、その後の呼の接続が出来なくなりユーザサービスに大きな支障が出るため、LOG 取得による情報の高信頼化が必要となる。一方、後者はリカバリ時に初期値に戻すことでユーザサービスに支障をきたさず、網内で対処できるデータで、LOG の取得は必ずしも必要でない。
- (4) データベースのリカバリを行っている間はサービスが一時中断し、ユーザからの呼を受付けることができない。通常、ユーザは再度ダイヤルすることとなるが、1 回の試行に 10 ~ 20 秒程度かかる。ユーザの再試行を考慮すると、数秒以内での高速なリカバリは必ずしも必要とされない。また、リカバリ時間が数分以上になると再試行の呼と新たな呼が加算されて輻輳を引き起こすことが懸念される。さらに、当初、呼が受けられなかったユーザも数回ダイヤルした時点でサービスが再開できていた場合には、ほとんど障害を意識することはないと想定される。このことを考慮するとリカバリに許容される時間は 60 秒程度となる。

(5) SCPの構成としては、データベースを複数のモジュールに分割して配備し、トランザクションを分散処理する負荷分散を前提とした。1モジュールとしての能力は、高度INのトランザクション処理のダイナミックステップ数および最近のプロセッサの性能から、単位時間当たり処理可能なトランザクション数は数100(Tr/秒)とする。また、地震、火災などの災害に対しても十分な信頼度を確保するため、ノード内ではモジュールを単位とした二重化を行い、ノード間では2つのノード相互のバックアップを行う。呼処理は交換機とその上位に配置したSCPで行う。サービスオーダ処理やサービス内容の変更等のカスタマ対応はSCPの上位に配置されるSMSで行う。すなわち、カスタマ情報等の管理はSMSで行い、呼処理に必要な情報のみがSMSからSCPへダウンロードされる。

(6) 高度INのノードにおいても交換機と同等の保守で対応可能なことが必須である。交換機に用いている論理回路から構成されるパッケージの保守はモジュールからの故障の報告にもとづき予備パッケージと交換して正常性を確認することにより行われる。電源障害に対応するため半導体メモリのバックアップにバッテリを用いたとしても、定期保守を不要とし、論理回路から構成されるパッケージと同等の保守条件で対応可能とする。

6.3 評価モデルの設定

前節で述べたように単位時間当たり数100のトランザクションを処理するためには、LOGを取得する外部の記憶装置は数100(回/秒)の入出力が要求される。最近の磁気ディスク装置の単位時間当たりの入出力回数は40～60(回/秒)程度でありLOG取得への適用は難しい。トランザクション毎にLOGを取得するためには、外部の記憶装置に高い入出力能力が要求され、半導体メモリの使用が必須となる。磁

気ディスク装置を複数並べて並列動作させるアレーディスク装置については、入出力回数の面では適用の可能性もあるが、容量面での使用効率が極端に低下するため、本章では対象とはしない（付録6.1参照）。

高頻度のLOG取得が可能な外部記憶系の構成として、当該モジュールのプロセッサ配下に接続する半導体ディスク装置と、別のモジュールのプロセッサに付属する主メモリとを代替案として、次に示す方式1、2を設定する。プロセッサ外部へのデータ送受はファイル装置への入出力処理による方法とプロセッサ間通信による方法の2通りに大別されるが、両方式はこれに対応している。加えて、方式2では単にLOGを取得するだけでなく、本来のトランザクション処理と並行してLOGに関係する処理を行えるなど、この面からもリカバリ時間短縮の可能性を併せて評価することができる。なお、両方式とも半導体メモリとしてDRAMを使用することとした。このため、電源障害時にデータが消失する問題があり、バッテリによるバックアップが必要となるが、これに対しても併せて考察する。

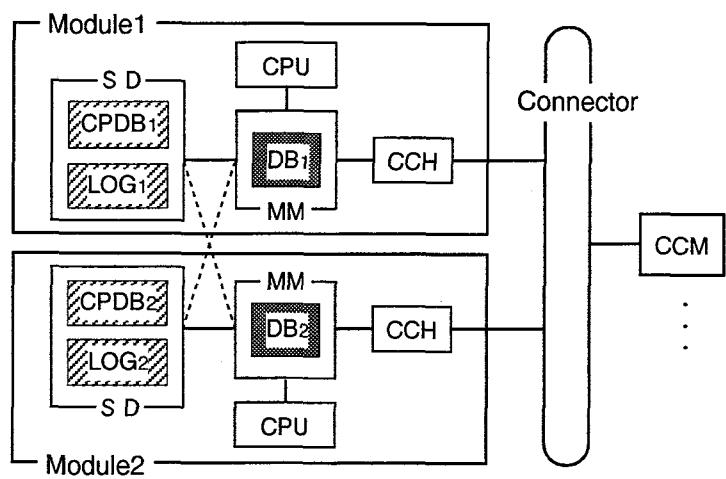
方式1：リカバリ情報の格納媒体として半導体ディスク装置を用いる。本方式ではトランザクション毎にLOGを半導体ディスク装置に書込む。また、リカバリ時間を短縮するため、トランザクション処理と並行して一定周期でCPDBを半導体ディスク装置に取得する。なお、取得したLOGを長時間に渡って保持することは不要であり、リカバリに必要なCPDBの書き込みが終了した場合は、CPDB取得開始時点より前のLOGは破棄する。

方式2：リカバリ情報の格納媒体として他のモジュールの主メモリを用いる。トランザクション毎に通信チャネルを介して他のモジュールにLOGを送る。他のモジュールではデータベースのコピー(DBC)を有し、受信したLOGをプロセッサによりDBCに上書きすることによって、常に最新のデータベースを保持する。このため、CPDBは使用しない。なお、DBCに上書き後、LOGは破棄する。

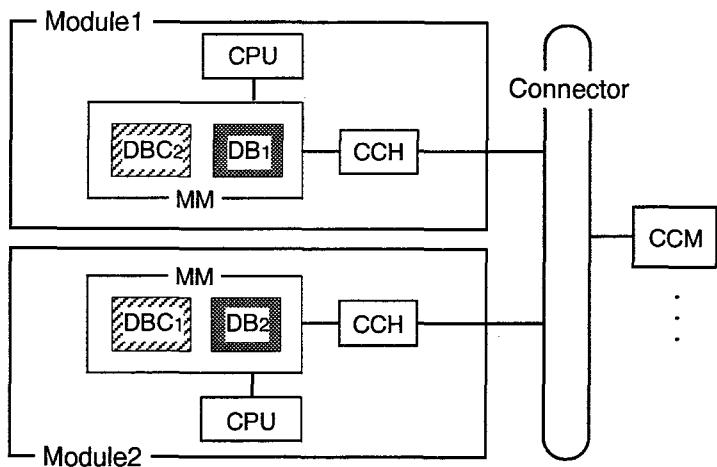
高度INではLOGを長時間に渡って保持しておくことは不要であり、本章ではリカバリに不必要となった後はLOGを破棄するとした。LOGを長時間に渡って保持する必要がある場合は、方式1ではSD上のLOGを破棄する前に大容量の記憶媒体に移す方法、方式2では相互にバックアップ関係にある双方のモジュールで、LOGを主メモリ上である程度バッファリングして大容量の記憶媒体に移す方法が考えられる。いずれにしてもLOGをバッファリングする必要があり、本章の方式1、2をベースに、LOGを大容量の記憶媒体へ書き込む処理を追加することにより容易に拡張できる。

方式1、2共にモジュールが障害となった場合、他のモジュールで処理を再開するための予備を含めたシステム構成例を図6.1に示す。2台のモジュールが1/2づつ負荷を分担し相互にバックアップする。OSやDBMS等のプログラムファイルの格納媒体としては磁気ディスク装置を使用するが、プログラムファイルはシステム起動時にのみアクセスされ、呼処理でのアクセスはない。このため、本章で対象とするメモリデータベースのリカバリ方式の選定に直接影響しないことから図6.1では省略している。

トランザクションの受け付けから応答までのプロセッサ(CPU)、半導体ディスク装置(SD)、通信チャネル(CCH)での処理の流れを図6.2に示す。方式1についてはCPDB取得処理のフローについても併せて示す。図6.2の直線上的記号はCPU、SD、CCHの処理時間である。なお、1回の入出力に要するSDの処理時間、CCHの処理時間は、アクセス時間とデータ長に比例するデータ転送時間との和で表わしている。ここで、図6.2で使用している記号の内容は表6.1に示す通りである。



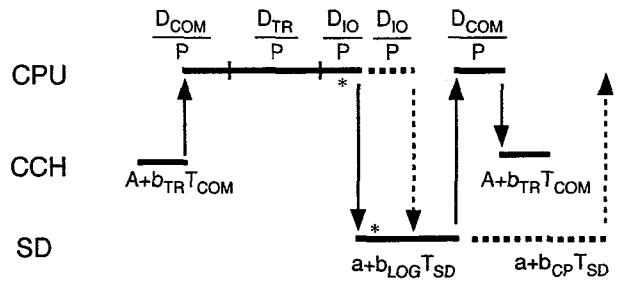
(1) Method 1



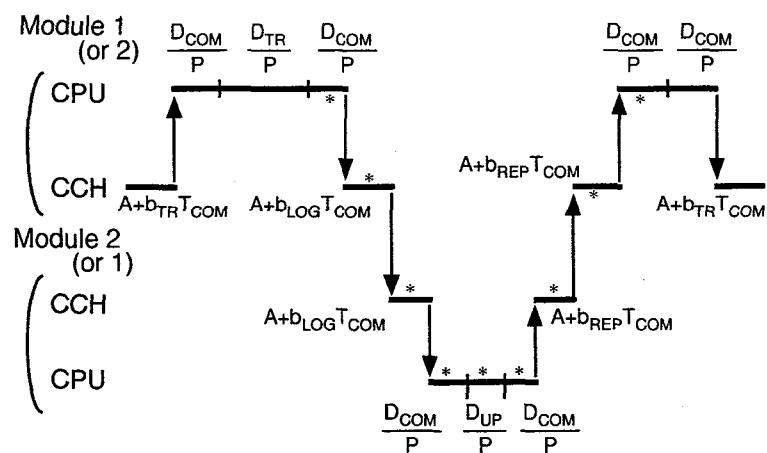
(2) Method 2

DB: Current database	CPDB: Database at the checkpoint
LOG: Log data	DBC: Database copy
MM: Main memory	SD: Semiconductor disk device
CCH: Communication channel	CCM: Communication Control Module

図 6.1 システム構成



(1) Method 1



(2) Method 2

— : Transaction Processing : Checkpoint Processing

* : Transaction without logging don't need these processing.

図 6.2 处理フロー

表 6.1 性能評価パラメータ

分類	記号	内容	設定値
D S	D _{TR}	トランザクション処理の D S	80KS
	D _{IO}	入出力処理の D S	5KS
	D _{COM}	通信処理の D S	10KS
	D _{UP}	L O G を D B に上書きする処理の D S	10KS
デ ー タ 長	b _{LOG}	L O G のデータ長	512B
	b _{TR}	トランザクションのデータ長	100B
	b _{REP}	方式 2 でバックアップ側モジュールからの L O G 受信の応答のデータ長	100B
	b _{RCV}	リカバリ時に 1 回の入力で半導体ディスク装置から読込むデータ長	64KB
	b _{CP}	C P D B 取得のため、1 回の出力で半導体ディスク装置に書込むデータ長	64KB
装 置 性 能	P	プロセッサ性能	20MIPS
	a	半導体ディスク装置のアクセス時間	1ms
	A	通信チャネルのアクセス時間	0.5ms
	T _{SD}	半導体ディスク装置のデータ転送時間	200ns/B
	T _{COM}	通信チャネルのデータ転送時間	1000ns/B

6.4 性能評価

方式1、2を性能面から比較評価する。

6.4.1 スループット、レスポンスタイム

方式1、2のCPU使用率 ρ_{CPU1} 、 ρ_{CPU2} はそれぞれ次式で表せる。

$$\rho_{CPU1} = (\lambda D_1 + \frac{MD_{IO}}{b_{CP}T_{CP}}) / P \quad (6.1)$$

$$\rho_{CPU2} = \lambda D_2 / P \quad (6.2)$$

ここで、 λ は単位時間当たりに当該モジュールで処理するトランザクション数、 T_{CP} はCPDBの取得周期、Mはデータベース容量である。 D_1 、 D_2 は方式1、2のトランザクション処理に必要なダイナミックステップ数(DS)であり次式となる。 ω はLOG取得が必要なトランザクションの割合である。

$$D_1 = D_{TR} + 2D_{COM} + \omega D_{IO} \quad (6.3)$$

$$D_2 = D_{TR} + 2D_{COM} + \omega(4D_{COM} + D_{UP}) \quad (6.4)$$

CPUの最大使用率を ρ_{CPU}^M とすると、CPUの能力に基づく、方式1、2の単位時間当たり処理可能なトランザクション数(以下、スループットと記す) λ_{CPU1} 、 λ_{CPU2} はそれぞれ次式で表せる。

$$\lambda_{CPU1} = \frac{\rho_{CPU}^M P - MD_{IO} / b_{CP}T_{CP}}{D_1} \quad (6.5)$$

$$\lambda_{CPU2} = \frac{\rho_{CPU}^M P}{D_2} \quad (6.6)$$

方式1の場合、LOG、CPDBをSDに書き込むため、SDの使用率 ρ_{SD1} は次式で表せる。

$$\rho_{SD1} = \lambda\omega(a + b_{LOG}T_{SD}) + \frac{M}{b_{CP}T_{CP}}(a + b_{CP}T_{SD}) \quad (6.7)$$

SDの最大使用率を ρ_{SD}^M とおくと、SDの能力に基づくスループット λ_{SD1} は次式と

なる。

$$\lambda_{SD1} = \frac{\rho_{SD}^M - \frac{M}{b_{CP}T_{CP}}(a + b_{LOG}T_{SD})}{\omega(a + b_{LOG}T_{SD})} \quad (6.8)$$

方式1の場合、 λ_{CPU1} と λ_{SD1} のどちらか小さい方でモジュールとしてのスループット λ_1 が決まる。方式2の場合は、 λ_{CPU2} によってモジュールとしてのスループット λ_2 が決まる。スループットを決定する要因としてはCCHのボトルネックも想定される。しかし、CCHの複数設置による負荷の分散は容易であり、本章ではボトルネックの要因としてCCHは考慮していない。

方式1はトランザクション処理と並行してCPDBを取得する必要があり、 λ_1 は T_{CP} に影響される。図6.3に T_{CP} と λ_1 の関係を示す。図6.3の算出に用いたDS、装置性能等を表6.1に示す。なお、 ω は0.7、 ρ_{CPU}^M 、 ρ_{SD}^M は0.8とした。 T_{CP} の短い領域ではSDのボトルネックが発生し、 T_{CP} はある値より短くできない。Mが100MBの場合、 T_{CP} が30～40秒ではSDにより、40秒以上ではCPUによりほぼ λ_1 が決まる。

方式1、2でトランザクション間の競合による待ちが無い場合の処理は図6.2中の実線で示した処理フローに従って進められる。トランザクション処理時間 T_{RES1}^0 、 T_{RES2}^0 はCPU、CCH、SDの処理時間の総和となり次式で表せる。

$$T_{RES1}^0 = (D_{TR} + D_{IO} + 2D_{COM})/P + (a + b_{LOG}T_{SD}) + (2A + 2b_{TR}T_{COM}) \quad (6.9)$$

$$T_{RES2}^0 = (D_{TR} + D_{UP} + 6D_{COM})/P + \{6A + 2(b_{TR} + b_{LOG} + b_{REP})T_{COM}\} \quad (6.10)$$

表6.1の条件では、 T_{RES1}^0 は約8ミリ秒、 T_{RES2}^0 は約12ミリ秒となる。ポアソン到着、指數分布サービスを想定し、使用率を0.8程度とすると、煩雑になるため式は省略するが、レスポンスタイムは競合がない場合の処理時間の約5倍程度となる。トランザクション間の競合を考慮しても、両方式のレスポンスタイムは100ミリ秒以内に収まり高度INの条件を満足する。

しかし、方式1の場合、CPDB取得処理との競合によりトランザクション処理が待たされることも考慮する必要がある。CPDB取得処理の b_{CP} を大きく設定すると

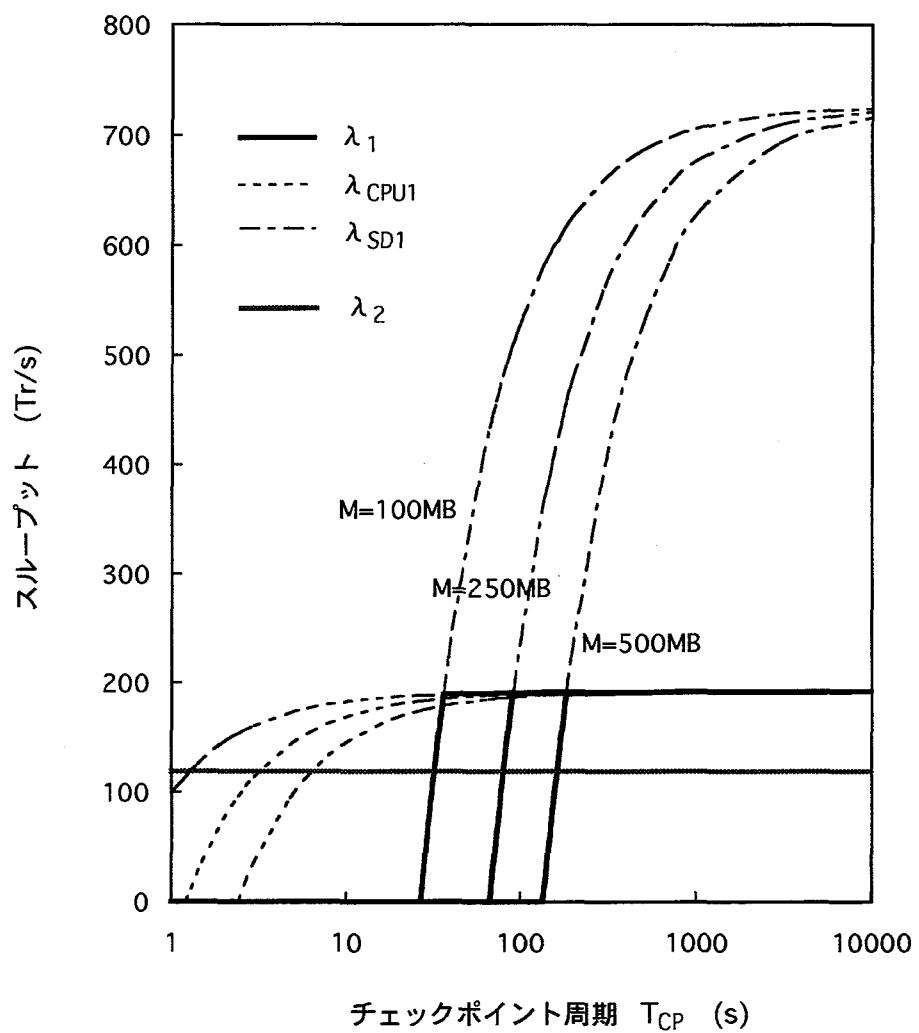


図6.3 チェックポイント時間とスループット

SDの占有時間が長くなり、CPDB取得処理と競合したトランザクションは長く待たされる。これについては、 b_{CP} を100KB以下に設定すればレスポンスタイムへの影響が無いことを確認した（付録6.2参照）。なお、 b_{CP} を小さくするとCPUの負荷が増加するため、本章では64KBに設定した。

以上、全体としてレスポンスタイムについては方式1、2の間で有意な差は生じないが、スループットに注目すると方式1が優れている。たとえば、Mを100MB、 T_{CP} を60秒とすると、 λ_1 は λ_2 に比べて約30%高い。これは一般に、方式2での通信処理は入出力処理に比べてDSが大きいこと、LOGをDBCにその都度反映する処理はCPDBを取得する処理に比べて、 T_{CP} が短くない限り、トランザクション当たりのDSが大きいことによる。

6.4.2 リカバリ時間

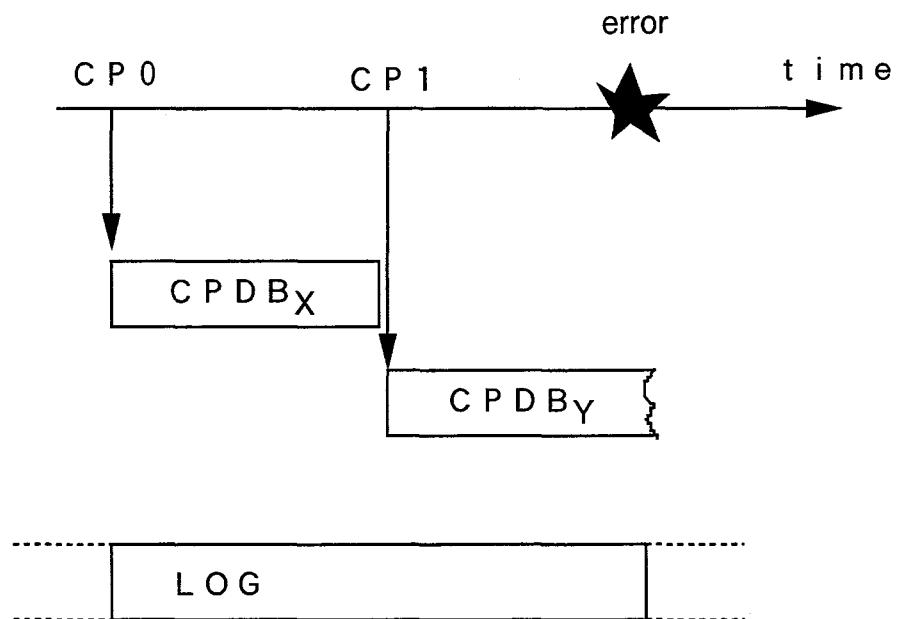
方式1の場合のリカバリ時間 T_{RCV1} は、SDからCPDBと最大で $2T_{CP}$ 分のLOGを読出す時間と、LOGをCPDBに上書きする時間の和で、次式となる。

$$T_{RCV1} = \frac{M + 2b_{LOG}\omega\lambda T_{CP}}{b_{RCV}}(a + b_{RCV}T_{SD}) + 2\omega\lambda T_{CP} \frac{D_{UP}}{P} \quad (6.11)$$

ここで、 $2T_{CP}$ とした理由は以下のリカバリを行うためである。

CPDBを取得中にエラーが発生した場合、CPDBは書換え途中であり、内容が保証されない。このため、CPDBの書き込みエリアをSD上に2面設け、交互に書き込むこととした。CPDB取得方法を図6.4に示す。リカバリを行う場合は、CPDBの取得が完了した最新のCPDBを用いる。図ではCPDB_yの取得中にエラーが発生し、CPDB_yが保証できないため、CPDB_xからリカバリする。LOGはCPOから障害時点までのすべてのLOGを読出す。LOGの量は障害発生ケースに影響されるが、確実なリカバリを保証するため最大値（ $2T_{CP}$ 分のLOG）とした。

方式2の場合は、既に、他のモジュールの主メモリ上に最新のデータベースが存在するため、リカバリ時間 T_{RCV2} は、障害でダウンしたことの検出とサービスを再開す



C P D B : チェックポイントデータベース

→取得エリアを半導体ディスク上に2面(X, Y)持ち、
交互にLOGを取得する。

L O G : データベースの更新ログ

→シーケンシャルに書き込む

図 6.4 C P D B 取得方法

る手続き等からなり数秒ですむ。

方式1の場合、Mを100MB、 T_{CP} を100秒、 λ を100(Tr/秒)とすると T_{RCV1} は約30秒となる。リカバリ時間は明らかに方式2が有利である。

6.4.3 データベース容量評価

前節で述べたように方式1はスループットで優れ、方式2はリカバリ時間で優れるという特徴があるため、その優劣を直接判断しにくい。本節ではリカバリ時間を条件としたとき、プロセッサ性能を同一としたモジュールで収容可能なデータベース(DB)容量を求め、これにより両方式の優劣を総合的に比較評価する。

6.4.3.1 収容可能なデータベース容量

カスタマ対応のデータ量をRとし、カスタマのデータを使用するトランザクションの単位時間当たりの平均件数を λ_c とおくと、トランザクション密度 ϕ は次式で表せる。

$$\phi = \frac{\lambda_c}{R} \quad (6.12)$$

モジュールに収容するカスタマ数をKとおくと、データベース容量Mは $K \cdot R$ となり、単位時間当たりのトランザクション数 λ は、

$$\lambda = \lambda_c K = \frac{\lambda_c}{R} K R = \phi M \quad (6.13)$$

となる。ここで、現状の主メモリの最大容量は数GB程度であり、 λ が数100(Tr/秒)のとき、 ϕ が約1/10(Tr/MB/秒)を下まわると容量面からHDを使用しなければならず、メモリデータベースに適さない。

方式1で収容可能な最大データベース容量を M_1 としたとき、 T_{RCV1} は式(6.1)、(6.7)を式(6.11)に代入し、また、 λ に ϕM_1 を代入することにより、次式で表せる。

$$T_{RCV1} \geq \frac{M_1(a + b_{RCV}T_{SD})}{b_{RCV}} \left\{ 1 + \frac{2b_{LOG}\omega\phi D_{IO}M_1}{b_{CP}(\rho_{CPU}^M P - \phi D_1 M_1)} \right\} + \frac{2\omega\phi D_{UP}D_{IO}M_1^2}{b_{CP}P(\rho_{CPU}^M P - \phi D_1 M_1)} \quad (6.14)$$

$$\begin{aligned} T_{RCV1} &\geq \frac{M_1(a + b_{RCV}T_{SD})}{b_{RCV}} \left\{ 1 + \frac{2b_{LOG}\omega\phi M_1(a + b_{CP}T_{SD})}{b_{CP}[\rho_{SD}^M - \omega\phi M_1(a + b_{LOG}T_{SD})]} \right\} \\ &+ \frac{2\omega\phi D_{UP}M_1^2(a + b_{CP}T_{SD})}{b_{CP}P\{\rho_{SD}^M - \omega\phi M_1(a + b_{LOG}T_{SD})\}} \end{aligned} \quad (6.15)$$

式 (6.14)、(6.15) はそれぞれ C P U に起因する条件と S D に起因する条件を示す。方式 1 ではリカバリ時間が与えられると両式を同時に満足する容量 M_1 のデータベースを収容できる。

T_{RCV1} と M_1 の関係を図 6.5 に示す。図中の実線は式 (6.14)、(6.15) の両方を満足する M_1 を表し、破線は式 (6.15) を満足する M_1 を表している。式 (6.14) を満足する M_1 はほぼ実線と同一となるため省略したが、 T_{RCV1} に比例して M_1 が増加する領域では、わずかな差はあるが S D によって M_1 が決まる。 M_1 が飽和する領域では C P U によって M_1 が決まる。これは M_1 が飽和値に近づくと、C P D B 取得のための負荷を C P U に割り当てることが難しくなり、急激な T_{CP} の延長と L O G 量の増加を引き起こすためである。

方式 2 で収容可能な最大データベース容量 M_2 は、式 (6.6)、(6.13) より、次式となる。

$$M_2 = \frac{\rho_{CPU}^M P}{\phi D_2} \quad (6.16)$$

方式 1、2 で収容可能な最大データベース容量の比を α ($= M_1 / M_2$) とおく。 α が 1 より大きい場合は方式 1 が有利であり、 α が 1 より小さい場合は方式 2 が有利となる。 ϕ と α の関係を図 6.6 に示す。図 6.6 から判るように、要求されるリカバリ時間 T_{RCV} を 60 秒とすると、 ϕ が 0.5 (Tr/MB/秒) 以上であれば方式 1 の方が、0.5 (Tr/MB/秒) 以下であれば方式 2 の方がより多くのデータベースを収容できる。

図 6.7 は ϕ と T_{RCV} から算出した方式 1、2 の適用領域を示す。図 6.7 より、デー

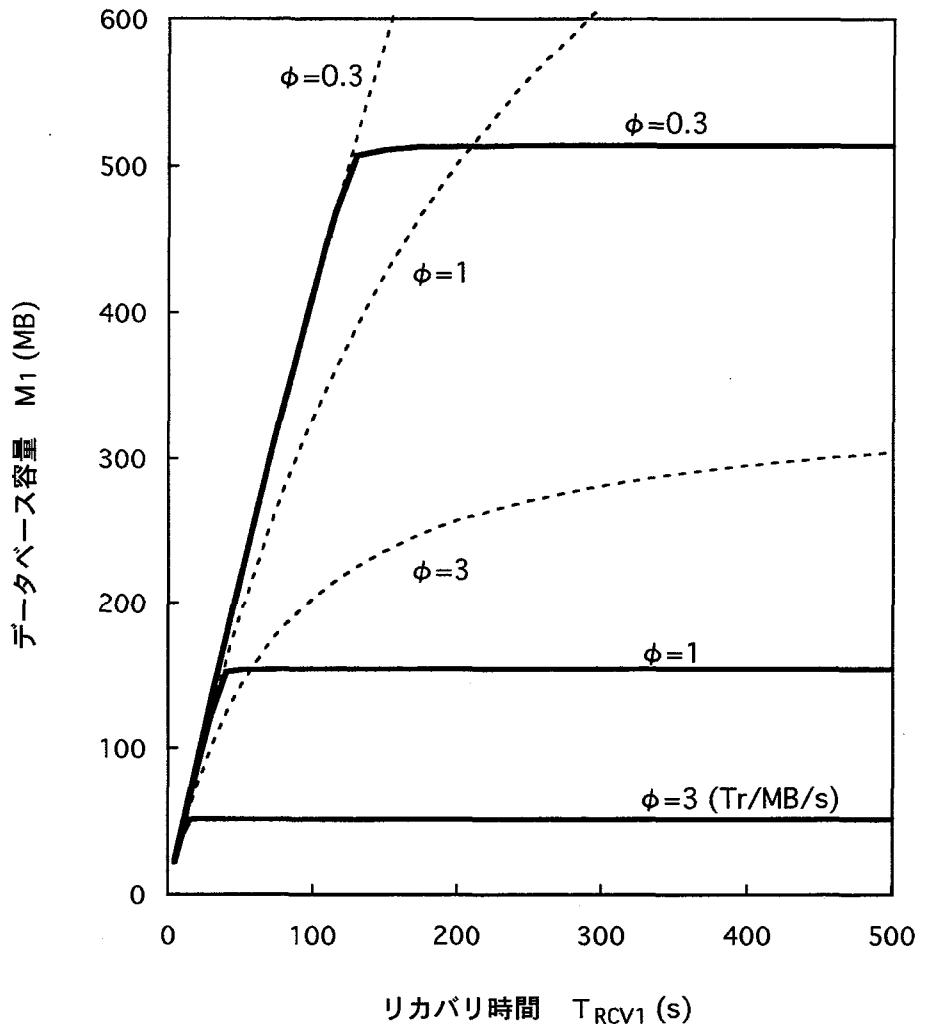


図 6.5 リカバリ時間とデータベース容量

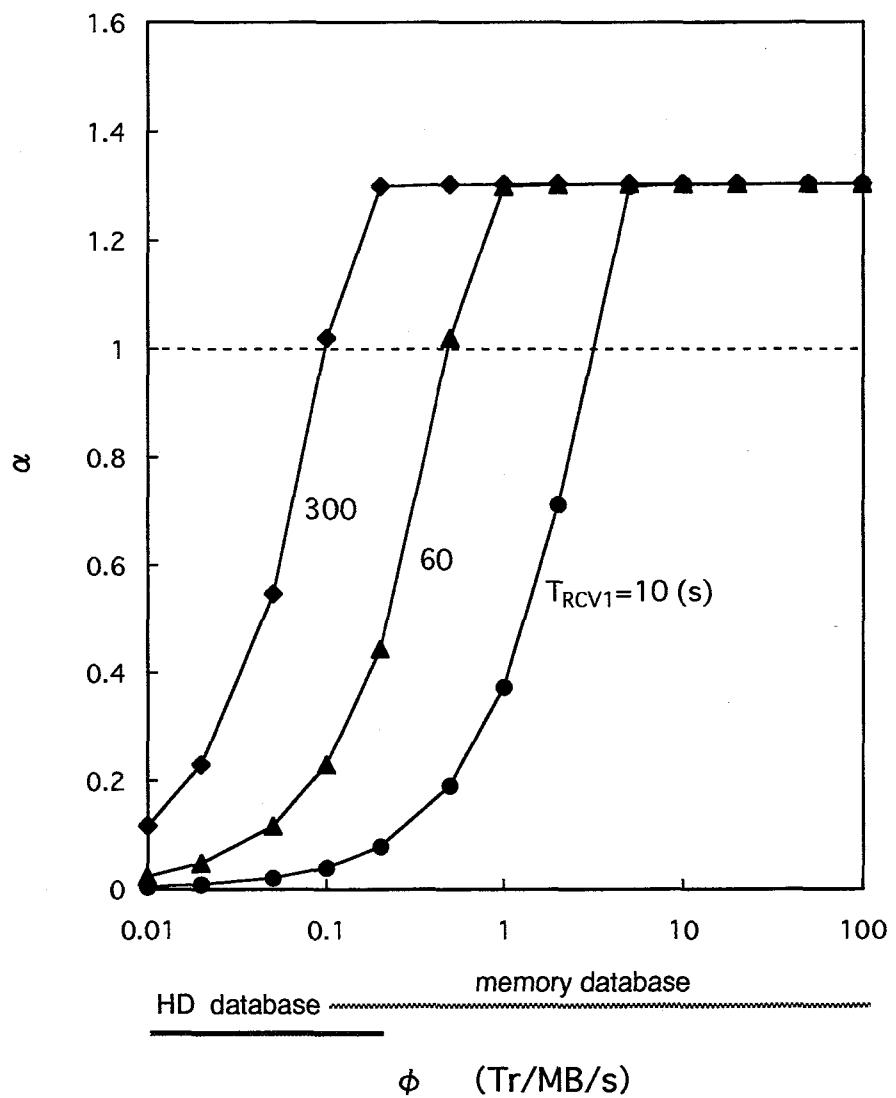


図 6.6 ϕ と α の関係

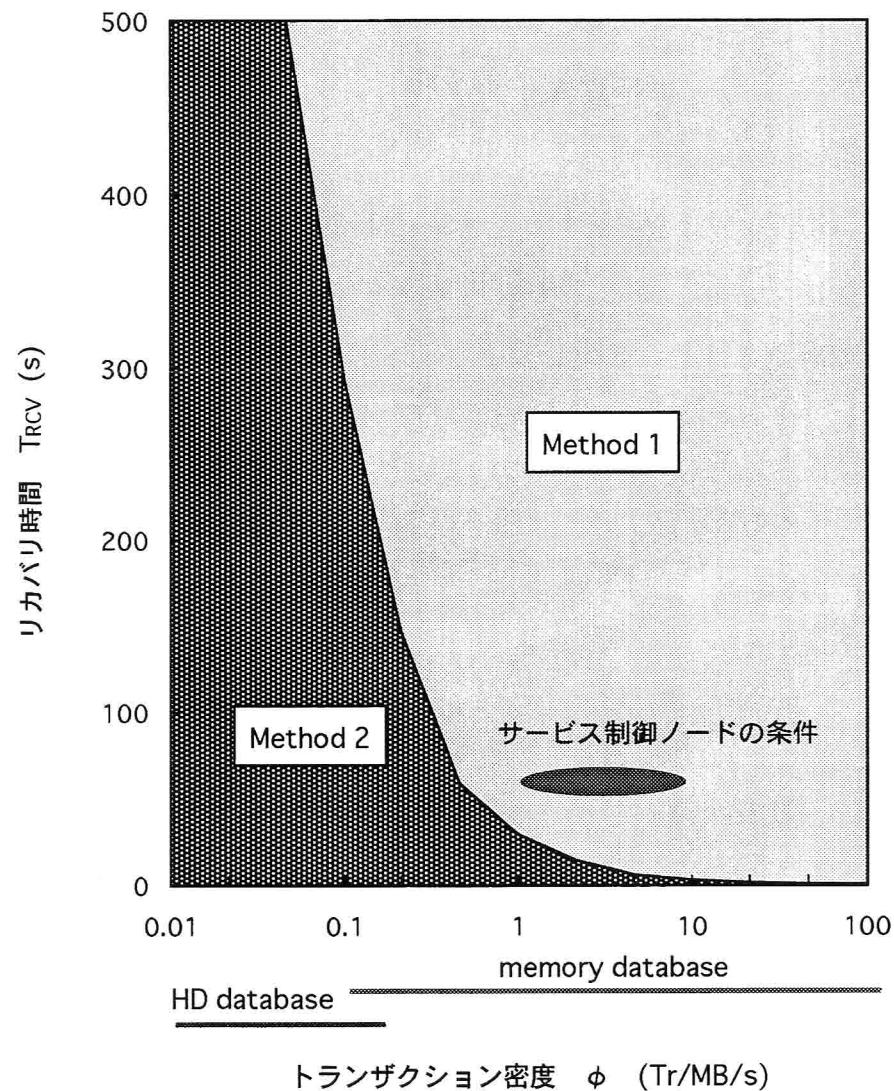


図 6.7 リカバリ方式の適用領域

タベースのトランザクション密度が高く、許容されるリカバリ時間が長いほど方式1を適用する方が有利であり、逆の場合は方式2を適用する方が有利となる。 ϕ が0.1(Tr/MB/秒)以上でなければメモリデータベースの実現は難しいことから、メモリデータベースに適する領域では広範囲にわたり方式1が有利となる。

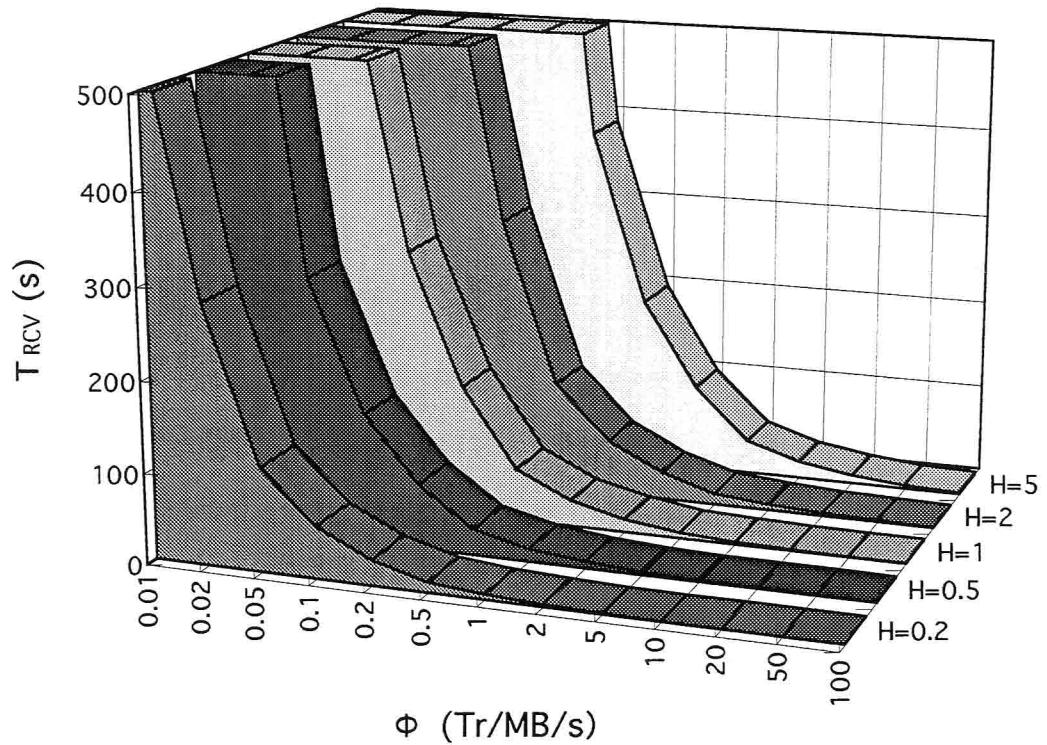
6.4.3.2 ハードウェア性能向上の影響

CPU性能、SD性能は今後とも向上することが予想される。CPU、SDの性能向上の度合いを h_{CPU} 、 h_{SD} とし、将来、CPU性能が $h_{CPU} \cdot P$ 、SD性能が a/h_{SD} 、 T_{SD}/h_{SD} となったとする。また、CPU、SDの性能向上度の比をH($= h_{CPU}/h_{SD}$)とすると、式(6.14)、式(6.15)は以下の通りとなる。

$$T_{RCV1} \geq \frac{\left(\frac{M_1}{h_{CPU}}\right)H(a + b_{RCV}T_{SD})}{b_{RCV}} \left\{ 1 + \frac{2b_{LOG}\omega\phi D_{IO}\left(\frac{M_1}{h_{CPU}}\right)}{b_{CP}\left[\rho_{CPU}^M P - \phi D_1\left(\frac{M_1}{h_{CPU}}\right)\right]} \right\} \\ + \frac{2\omega\phi D_{UP}D_{IO}\left(\frac{M_1}{h_{CPU}}\right)^2}{b_{CP}P\left\{\rho_{CPU}^M P - \phi D_1\left(\frac{M_1}{h_{CPU}}\right)\right\}} \quad (6.17)$$

$$T_{RCV1} \geq \frac{\left(\frac{M_1}{h_{CPU}}\right)H(a + b_{RCV}T_{SD})}{b_{RCV}} \left\{ 1 + \frac{2b_{LOG}\omega\phi\left(\frac{M_1}{h_{CPU}}\right)H(a + b_{CP}T_{SD})}{b_{CP}\left[\rho_{SD}^M - \omega\phi\left(\frac{M_1}{h_{CPU}}\right)H(a + b_{LOG}T_{SD})\right]} \right\} \\ + \frac{2\omega\phi D_{UP}\left(\frac{M_1}{h_{CPU}}\right)^2 H(a + b_{CP}T_{SD})}{b_{CP}P\left\{\rho_{SD}^M - \omega\phi\left(\frac{M_1}{h_{CPU}}\right)H(a + b_{LOG}T_{SD})\right\}} \quad (6.18)$$

図6.8は ϕ と T_{RCV} が与えられた場合、方式1と方式2の優劣がHによってどのように変化するかを表している。将来、CPU、SDの性能が向上したとしても、その向



shadow zone : $\alpha < 1$ (method 2)
no shadow zone : $\alpha > 1$ (method 1)

図 6.8 リカバリ方式の適用領域
(ハード性能との関係)

上度度合いが同一、すなわち $H = 1$ であれば、式 (6.14) 、 (6.15) と式 (6.17) 、 (6.18) は M_1 が (M_1 / h_{CPU}) となる以外は全く同一の式となる。これは、リカバリ時間同一とした場合、 M_1 を h_{CPU} 倍大きくできることを意味している。方式 2 では h_{CPU} に比例して M_2 が増加する。方式 1 、方式 2 共に収容可能なデータベース容量の向上度が同一となるため、現時点の評価結果がそのまま当てはまる。また、図 6.8 からわかるように、 H が大きくなるに従って方式 2 の適用領域が拡大し、逆に小さくなるに従って方式 1 の適用領域が拡大する。

6.5 コスト評価

方式 1 、 2 をコスト面から比較評価する。図 6.1 に示すように、モジュールを構成するプロセッサ、主メモリ、通信チャネル等は方式 1 、 2 に係わらず必要となる。このようにモジュールのコストを構成する主要部分は両方式で同様に必要とされる。方式 1 と方式 2 のモジュール構成上の主要な相違点は、方式 1 の場合はチェックポイントデータベース (CPDB) やデータベースの更新ログ (LOG) を取得するため半導体ディスク装置が必要であること、方式 2 の場合はバックアップする相手モジュールのデータベースのコピー (DBC) を主メモリに常駐させておくため主メモリの記憶容量が多く必要となることのみである。半導体ディスク装置と主メモリの記憶媒体の差や必要とされる記憶容量に若干の差があるとしても、モジュールのコストへの影響は少なく、方式 1 、 2 でモジュールコストはほとんど変わらない。

モジュールコストが同一であるとすれば、システムとしてのコストはモジュール台数に依存して決まる。モジュール台数はシステムに要求されるスループットを、方式 1 、 2 のモジュール 1 台当たりのスループット λ_1 、 λ_2 で割った値となる。前節の性能評価で述べたように、スループットは方式 1 が方式 2 に比べ約 30 % 高いため、システムとしてのコストは方式 1 が方式 2 に比べ約 30 % 低くなる。

6.6 信頼性評価

6.6.1 不揮発化

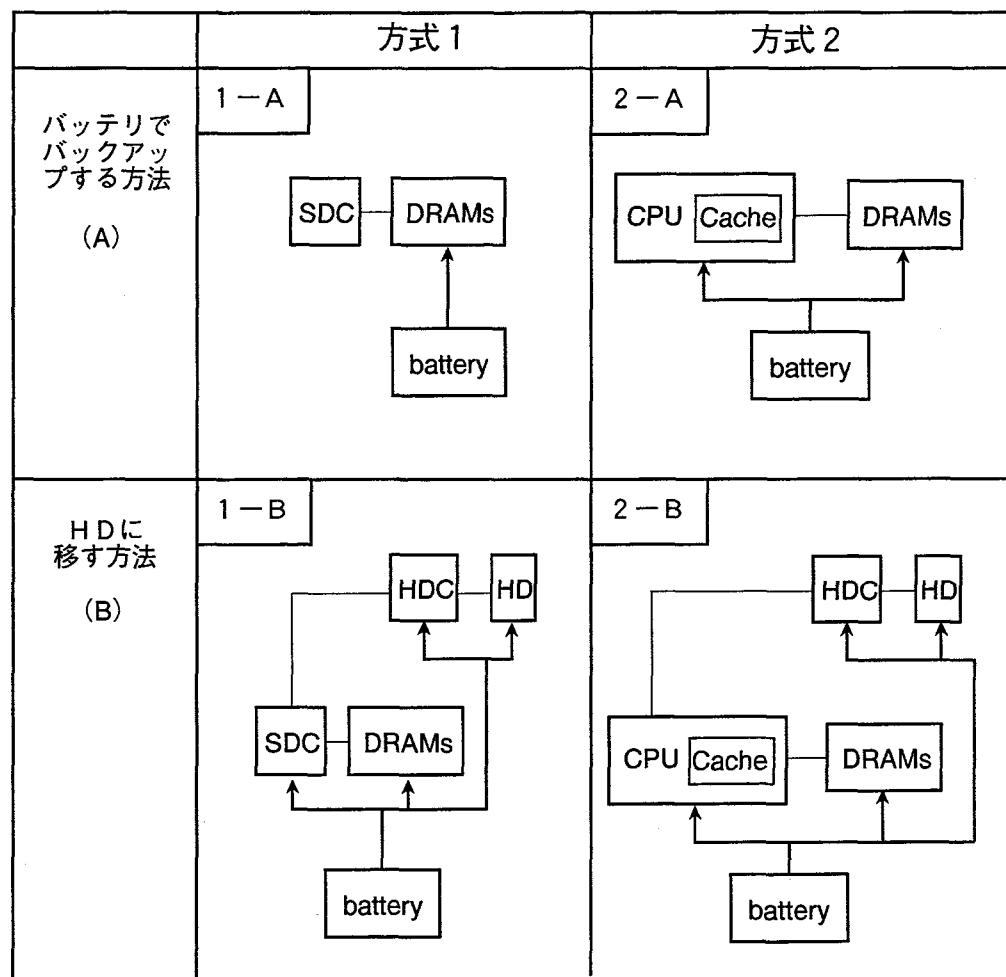
方式1、2ともにメモリ素子としてDRAMを使用するため、電源の供給が停止した場合には主メモリ上のデータベースはもとより、半導体ディスク装置または他のモジュールの主メモリに格納しているバックアップ情報も消失する。電源系の障害は、①モジュール内の電源部の障害、②モジュールへの配電系統の障害、③局舎全体の電源障害に分類できる。

方式1の場合、半導体ディスク装置および電源部を二重化し、配電系を別にすることにより①、②の部分の一重障害では情報は消失せず、情報の消失は二重障害のみとなる。方式2の場合、相互にバックアップの関係にあるモジュールの電源部は別であり、モジュールへの配電系を別にすることにより同様に①、②の一重障害では情報の消失を回避できる。しかし、両方式とも、このままでは③の障害には対応できない。局舎電源を構成する各装置についても、通常、多重化等により高信頼化がはかられているが、最近の10年間においても、落雷や増設工事中の事故等により、局舎全体の電源断が発生している。電源断の時間は、ほとんどが1～2時間以内であり、極くまれではあるが、10時間程度となった場合がある。情報の消失を防止するためには固有の電源による不揮発化が必要である。

不揮発化の方法としては、電源供給停止後、バッテリによりバックアップする方法（方法A）とDRAMの内容を磁気ディスク装置（HD）に移す方法（方法B）がある。前者はバックアップ時間Tが長くなるに従って、バッテリ量が大きくなる問題があるが、装置構成は簡単である。一方、後者は一度HDに格納すれば、バックアップ時間の制約は無い。しかし、データをHDに移す間はHDや書込み／読み出しを制御するコントローラを含め、バッテリによるバックアップが必要となる。

方式1、2の不揮発化方法を表6.2に示す。方法2-Aでは、最新の情報がキャッシュメモリ上にある場合があり、キャッシュメモリの内容を主メモリに書き出した

表 6.2 不揮発化方法



あと、主メモリのDRAMのみを不揮発化する必要がある（付録6.3参照）。各方法の消費電力は次式で表せる。記号の内容を表6.3に示す。

$$E_{1-A} = p_{MS}BST/m \quad (6.19)$$

$$E_{1-B} = (p_{MA}BS/m + p_{SDC} + p_{HDC} + p_{HD})S/M_{HD} \quad (6.20)$$

$$E_{2-A} = (p_{MA}BS/m + p_{CPU})T_0 + p_{MS}BS(T - T_0)/m \quad (6.21)$$

$$E_{2-B} = (p_{MA}BS/m + p_{CPU} + p_{HDC} + p_{HD})S/M_{HD} \quad (6.22)$$

バックアップ時間と消費電力の関係を図6.9に示す。地震、火災等の特別な障害を除けば、局舎の電源が長時間に渡って障害になることは少なく、情報のバックアップ時間は20時間程度とした。図6.9より、バックアップ時間がこの程度であればいずれの方法も消費電力の差はほとんど無い。方法BはバッテリのほかにHDおよびHD制御装置が必要であり、ハードウェア量の点では方法Aが有利である。これに加えて、HDは可動部を有しており、信頼性、保守性の面でも方法Bは劣る。全体としていずれの方法もバッテリ量の差は小さいが、可動部を持たないこと、構成が最も簡単なことから方法1-Aの実用性が高い。

要求される電力容量がこの程度まで小さくなると、一般に大量に使用されているニッカド電池を使用可能となる。ニッカド電池はシール鉛電池に比べて単位重量当たりの電力容量（重量電力密度）が大きく、装置に組み込むことが容易である。また、ニッカド電池の寿命（約5年程度）はシール鉛電池に比べ数年長い。

6.6.2 高信頼化構成

モジュールとしての耐障害性、冗長化構成の選択の2面から方式1、2について簡単に比較する。

6.6.2.1 モジュールとしての耐障害性

データベースのリカバリのために増加するメモリの容量は方式1の方がLOGの

表 6.3 不揮発化評価パラメータ

記号	内容	設定値
P_{MA}	動作時のDRAMの消費電力	0.25W
P_{MS}	スタンバイ時のDRAMの消費電力	$1.9 \times 10^{-3}W$
S	記憶容量	500MB
m	DRAMの記憶容量	16Mbit
B	1バイトを実現するために必要なビット数	10bit/Byte
P_{SDC}	半導体ディスク装置の消費電力	40W
P_{CPU}	プロセッサの消費電力	100W
P_{HDC}	磁気ディスク制御装置の消費電力	40W
P_{HD}	磁気ディスク装置の消費電力	30W
M_{HD}	磁気ディスク装置のデータ転送速度	4MB/s
T_0	キャッシュメモリの内容を主メモリに書き出す時間	1s

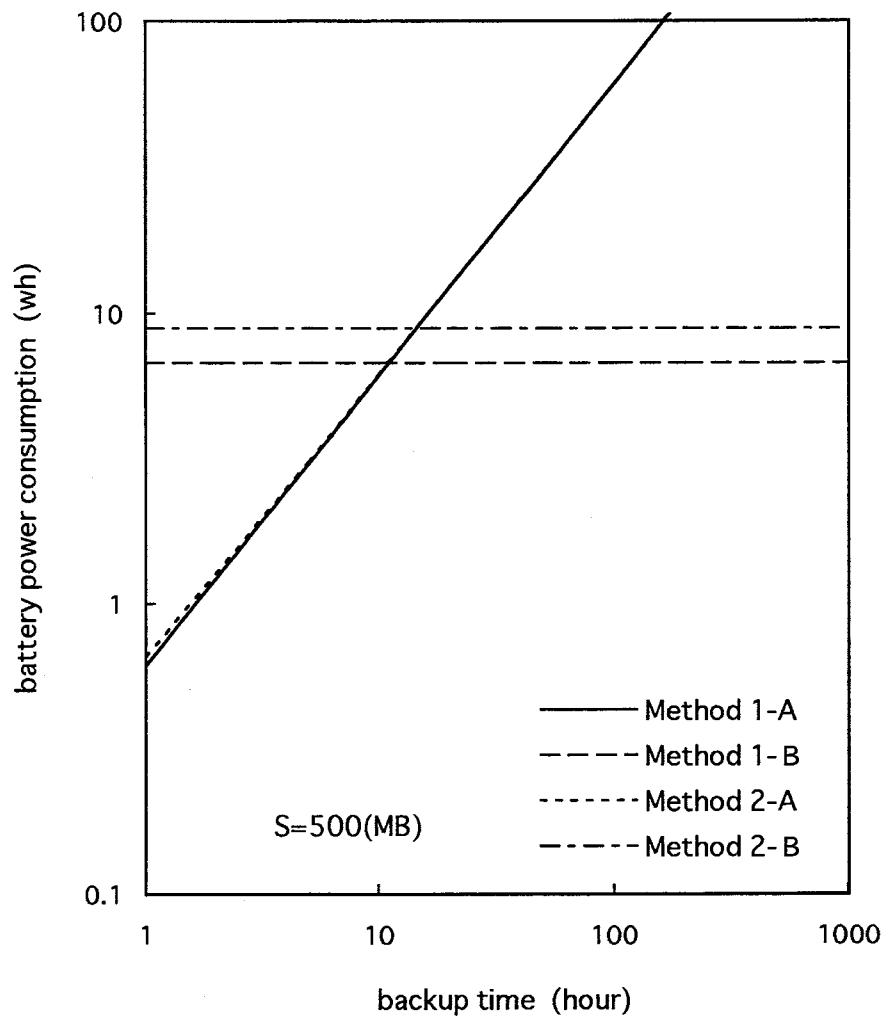


図 6.9 バックアップ時間とバッテリの消費電力

容量分多くなるが、モジュール全体で見ると、プロセッサ系がF I T数の主要部を占めるため、方式1、2でF I T数の差は小さい。このため、ハードウェア障害に起因するモジュール障害発生の度合いはほとんど同じとなる。

S C Pではシステム稼働中にファイル更新等の保守を行う必要がある。方式2ではトランザクション処理毎にバックアップ側のDBを最新の状態に更新してしまうため、稼働中のファイル更新の際のバグや操作ミスにより誤ってDBを破壊した場合には、DBを元に戻すことはできない。このため、CPDBを事前に取得する等の何らかの新しい処置が必要となる。これに対して、方式1ではCPDBとLOGを残しており、DB破壊が発生しても、CPDBやLOGを用いてDB破壊を限定し、復元することができる。また、SDの容量を増加させることでT_{cp}とは独立に、CPDBの内容を多重に必要な時期まで保存しておくことが容易である。このため、方式1の方がサービス中のソフト保守はしやすい。

6.6.2.2 冗長化構成とプロセッサ稼働率

ノード内バックアップ方法として、1対1の相互バックアップのほかに、1モジュールが障害となった場合、障害となったモジュールの負荷を他のnモジュール（n≥2）で分担する方法もある。プロセッサ稼働率はn／(n+1)となり、nの増加と共に高く設定できる。方式1ではSDを他のモジュールからアクセス可能とする必要があり、nに比例してSDのポート数が増加し、ハード構造が複雑になる。一方、方式2はモジュール間結合機構を介して他のモジュールと結合されており、ハード構成上、nの増加に対して柔軟に対応できる。しかし、1モジュールの障害を複数のモジュールがバックアップすると、障害発生時の切り替えや障害装置回復時の切り戻しのソフト制御が極めて複雑となる。このため、実用性の高いシステムとして構築するためには両方式共に1対1の相互バックアップとし、1／2のプロセッサ稼働率で動作させる方が望ましい。

加えて、地震、火災などの災害に対応するためには、地理的に離れた2つの地点に

設置されたノード間での相互バックアップが必要となる。この場合、どちらかのノードが障害となった時には、別のノードがすべての負荷を処理する必要があり、通常は $1/2$ のプロセッサ稼働率で動作させる。第3章ではノード内バックアップとノード間バックアップのためのプロセッサ稼働率の余裕を共用させることにより、効率的なバックアップ方法を提案した。プロセッサの稼働率をノード間バックアップに合わせて $1/2$ に設定し、ノード内バックアップ時にもこの $1/2$ のプロセッサ稼働率の余裕を活用する方法である。ノード間バックアップとノード内バックアップが同時に発生する確率は極めて低いため、経済的にシステムの高信頼化を達成できる。ノード間バックアップとノード内バックアップの両方を組合せた冗長構成を採る場合は、ノード内バックアップ方法に依存せず、方式1、2に係わらずプロセッサ稼働率は同一となる。

6.7 実現方式

前節までの検討を踏まえ、高度IN用SCPとして試作したモジュールの方式と構成について述べる。

リカバリ方式としては半導体ディスク装置にLOGとCPDBを取得する方式1を採用した。また、半導体ディスク装置の不揮発化はDRAMのみをバッテリによりバックアップする方法とした。これは主に次の理由による。

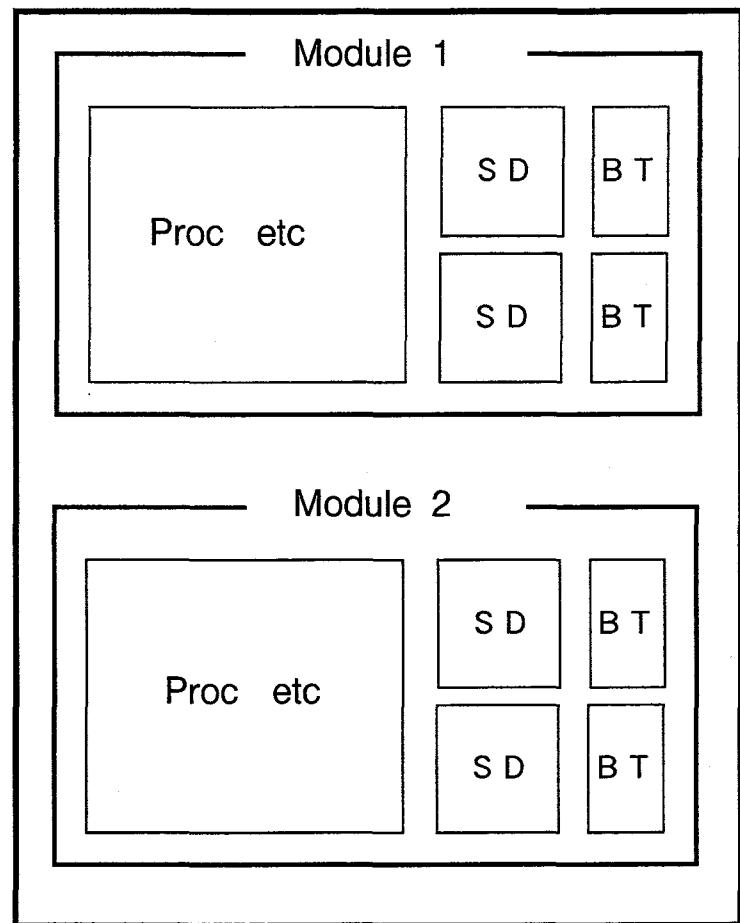
高度INのデータベースのトランザクション密度は $1 \sim 10$ (Tr/MB/秒) 程度であり、要求されるリカバリ時間は60秒程度である。この条件では、図6.7から判るように明らかに方式1の方がより多くのデータベースを収容できる。また、20時間程度の不揮発化ではSDをバッテリによりバックアップする方法が最も構成が簡単でハードウェア量が少ない。高信頼化のための冗長化構成の面からは、高度INではノード間バックアップとノード内バックアップの併用が必要であり、この場合、方式1、2でプロセッサ稼働率の差はない。加えて、ソフト障害に対処する能力は方式1

の方が高い。

システムとしてのバックアップ動作の確認等を行うためモジュールを試作した。1つの架の中に2台のモジュールを搭載し、同一架内のモジュール間で、SDを介して相互にバックアップ可能とした。バッテリはCPU、SDと同一バックボードに搭載し、全体としてコンパクトな実装が可能となった。モジュールの実装構成を図6.10に示す。試作を通して、電源断時に正常にバックアップに移行して20時間のバックアップに耐え得ること、電池の周囲温度の上昇が室温に比べて10°C以下に収まること、およびコンピュータ等に適用されているVCCIの伝導雑音や放射雑音の規定値を満足すること等を確認した。

CPUやSDの論理パッケージは故障検出回路を有しており、その保守は、故障が検出された時点での報告にもとづき予備パッケージと交換して正常性を確認することにより行われる。一方、バッテリには寿命があり、ある時間経過後には使用できなくなる。このため、通常は寿命が切れる前に定期的に交換する必要がある。モジュールを構成するCPU、SD等の保守とバッテリの保守が異なることは、2通りの保守体制をとる必要があり、保守が煩雑となる問題がある。加えて、バッテリの寿命は周囲の温度など使用状況に応じて数倍変動する。このため、定期保守する場合は、大半のバッテリではバックアップ能力がまだ残っているとしても、劣化が早く寿命が切れると予想されるごく少数のバッテリに合わせて、交換周期を短く設定する必要がある。

これらを解決するため、バッテリの劣化判定機構を開発し^(6.3)、CPU、SD等の論理回路から構成されるパッケージと同様に障害や劣化の報告を受けてから交換可能とするとともに、平均的にバッテリを長く使用できるようにした。ここで、劣化判定とは、バッテリが所望の時間バックアップに耐えられるか否かを判定し、耐えられないと判断した場合、交換が必要なことを遠隔の保守センタに報告する機能である。図6.11に劣化判定機構を組込んだバッテリの構成を示す。バッテリは単3のニッカド電池6個を1パックとし、4パック（電池総数は24個）を使用して、論理回路を搭載するパッケージ（30×33cm²）と同一のものに実装し、パッケージ2枚分の体積で実現できた。劣化判定は1パック毎に図6.11のスイッチをテスト側（T



Proc: Processor
SD: Semiconductor disk device
BT: Battery

図 6.10 モジュールの実装構成

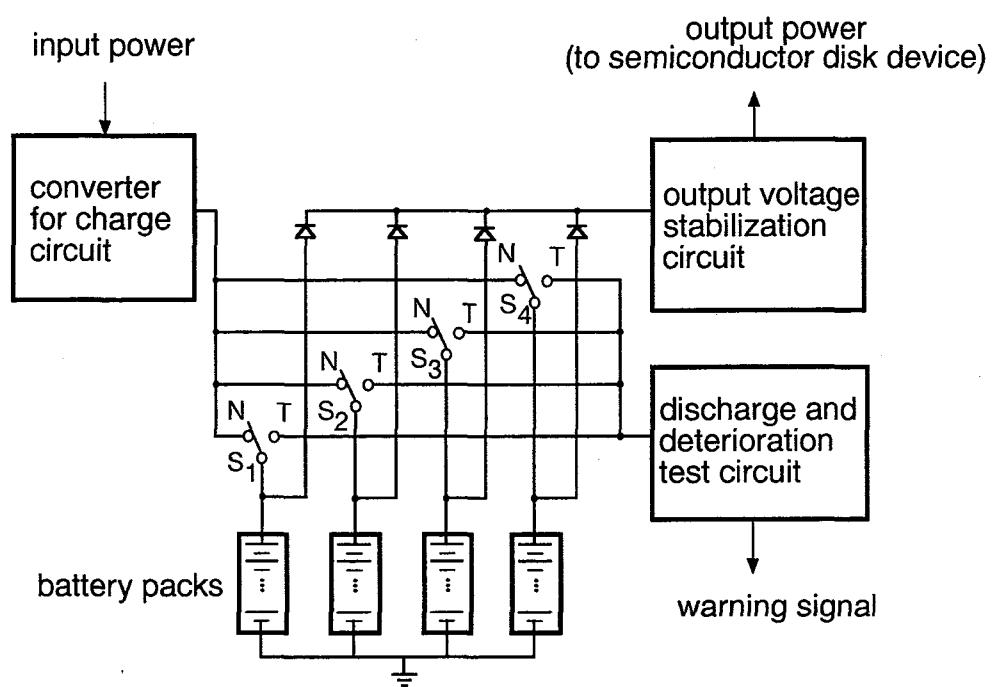


図 6.1.1 バッテリの構成

側)にセットし放電回路を通して放電させ、経過時間と電圧低下からバッテリのバッカアップ能力が20時間以下となったとき劣化と判断する。1パックの判定が終了して正常であれば、スイッチを充電側(N側)に倒して充電する。充電が終了した時点で次のパックの劣化試験を行う。これにより劣化試験中に電源障害が発生した場合でも残りの3パックで半導体ディスク装置のバッカアップを可能とした。1パックの劣化判定と再充電に約3日を要し、約半月で4パックの試験が終了する。この試験を3ヵ月周期で行うこととした。

6.8 結言

メモリデータベースのリカバリのため、トランザクション毎のLOGを取得する記憶装置として、半導体ディスク装置を適用する方式と他のモジュールの主メモリを適用する方式を代替案として設定し、性能、信頼性の2つの面から総合的に比較評価を行った。性能に関しては、スループットとリカバリ時間を併せて、モジュール毎に収容可能なデータベース容量を評価し、データベースのトランザクション密度と要求されるリカバリ時間が与えられた場合の適用領域を明らかにした。信頼性に関しては、不揮発化方法および冗長化構成法について両方式の比較評価を行った。また、バッテリバックアップについては劣化判定機構を新たに開発し、論理パッケージと同様に障害や劣化の報告を受けてから保守可能とし、バッテリ導入上の問題を軽減する方法を提案した。その結果、高度INのメモリデータベースのリカバリ方式としてバッテリバックアップを行った半導体ディスク装置による実用性の高い方式を明らかにした。なお、本章では、LOGはリカバリに必要な時点を過ぎれば破棄するとして検討を進めたが、OLTPシステムの中には数時間前の状態からリカバリを可能とするためや、リカバリ以外の用途に使用するため、LOGを長時間に渡って保持する必要のあるアプリケーションもある。この場合は、本章のモデルに、LOGを大容量の記憶媒体に取得する処理を追加することにより容易に拡張できる。

第7章

相互にバックアップされたサービス制御ノードの 高信頼化運転保守技術

7.1 緒言

高い信頼性が要求される通信システムや、バンキングシステム、クレジット照会システム等のオンライントランザクション処理（O L T P）システムでは、従来からノード内で二重化する冗長構成が一般にとられている^{(21) (55)}。しかし、データベースを用いて全国を対象にサービスを行うようなシステムでは、ノード内の二重化のみでは、地震、火災等の大規模災害に対しては十分でなく、地理的に離れた2つの地点にノードを設置して、ノード相互でデータベース上の情報を含めてバックアップする冗長構成がとられるようになってきている^{(23) (24) (35)}。インテリジェントネットワークの急激な普及や、O L T Pシステムでのサービス内容の多様化、高度化に伴って、このような形態のシステムが増加する傾向にある。

一方、システムの信頼性は、このような高信頼化のための冗長方式とならんで保守方式により大きく変動する。ノードに保守チームが常駐している場合は、障害検出と同時に修理に取り掛かることが可能である。しかし、ノード毎に保守チームを常駐させると保守コストの増大を招くため、通常、保守センタに保守チームを集中し、複数のノードを少数のチームで保守する方法が採られている。この場合、サービス停止につながるような障害では常駐の保守チームにより緊急に駆付けて対応し、サービス停止に至らない軽度の障害ではある程度時間的な余裕を持って平日勤務の保守チームにより対応する等、システムの緊急度合いに応じて駆付け時間を調整することが多い。

地理的に離れたノード間での相互バックアップは大規模災害への対応が本来の目的であるが、通常の装置障害に対しては、冗長度が上がるため、高い信頼性が容易に達成可能となる。許容される範囲での信頼性低下のもとで、駆付け時間の長時間化を図ることにより、保守の効率化とシステムの信頼性をバランス良く実現することが期

待できる。

本章では、データベースを持ち二重化された2つのノードを離れた地点に配置し、相互にバックアップするシステムを対象として、駆付け保守の方法と信頼性の関係を検討し、保守の効率化方法について考察する^{(64) (65)}。はじめに、ノード内で二重化された2つのノードが相互にバックアップするシステムを駆付けにより保守する場合のシステムの構成と信頼性算出方法を示す。既に、渡辺らはノード内で二重化した交換機と宅内装置からなる加入者系を対象として、駆付け保守による修理遅延の影響について検討を行っている⁽³⁷⁾。本章では、渡辺らの考え方を、二重化された2つのノードが相互にバックアップするような構成をとるシステムに拡張して用いている。次に、信頼度の時間的な変化に注目し、システムが障害となった場合、どの程度急いで駆付けるかをシステム内の障害装置数や保守チームが到着しているか否かの状態によりレベル分けすることを提案し、どのレベルまで緊急駆付けを行うかに応じて、信頼性や駆付け回数の変化を明らかにする。これらの結果に基づいて、システムに要求される信頼度を満足した上で、非緊急駆付けについては複数の障害をスケジュール化可能で、直ちに駆付けなければならない緊急駆付け回数の少ない効率的な保守方法を明らかにする。

7.2 評価対象システムの構成と動作

7.2.1 評価対象システムの信頼度構成

本章では、データベースを保有し、トランザクション処理や呼処理を行うモジュールを基本の構成単位とする。ここで、モジュールは本章での障害検出、修理の基本の構成単位であり、プロセッサ、メモリ、ファイル装置、通信装置等、一式から構成される。ノード内ではモジュールが二重化され、さらに地理的に離れた2つのノードで1対1に相互にバックアップされる構成を信頼性評価の単位とする。以下、これをシ

システムと記す。このような冗長構成がとられる例として、高い信頼性が要求される高度インテリジェントネットワーク（高度IN）やOLT Pがある。高度INでは、電話機から発信された呼は交換機、共通線信号網を介してデータベースを保有するサービス制御ノード（SCP）に集められる。また、OLT Pでは、端末からのトランザクションが通信回線を通じて処理センタに集められる等、ノードの設置に際して地理的な制約が少ないという特徴がある。

高度INの場合、全国的な大規模サービスを想定すると、ネットワーク内に数10から数100程度のシステムを配置する必要がある。OLT Pの場合は、通常、サービス毎には1システムから数システムで構成されるが、あるベンダが複数の顧客のシステムの保守を請負う場合には、多数のシステムの保守をまとめて行うこととなる。従って、保守対象は数10から数100の規模で配置されるシステム群とする。図7.1に評価対象システム全体の構成を示す。

高度INやOLT Pのデータベースは、呼毎、トランザクション毎に更新されるデータベースであることから、異なる地点に配置されたノード間で相互にバックアップを行う場合は、データベースそのものを相互に持ち合い更新内容をバックアップ側のデータベースに反映する必要がある⁽¹⁵⁾ ⁽⁴³⁾。本章では図7.1に示すように、2つのノードの対応するモジュール#1-1-1と#1-2-1および#1-1-2と#1-2-2がデータベースを相互に持ち合い、更新内容を報告し合うことによりバックアップ側のデータベースに反映した。

高度INでのパーソナル通信サービスのように单一サービスを多数のノードで分散して実現する必要がある場合のノード間バックアップ方法としては、1対1の相互バックアップのほかに、1ノードがダウンした場合、その負荷を複数のノードで分担する方法もある。この場合、ダウンしたノードの負荷が複数のノードに分散されるため、平均的にノードの稼働率を高く設定できるという長所がある。しかし、1ノードの障害を複数のノードがバックアップすると、通常動作時において、データベースの更新情報を常に複数のバックアップ側ノードに送信しなければならないことに加えて、ノードダウン発生時の切り替えや障害ノード回復時の切り戻しのソフトウェア制

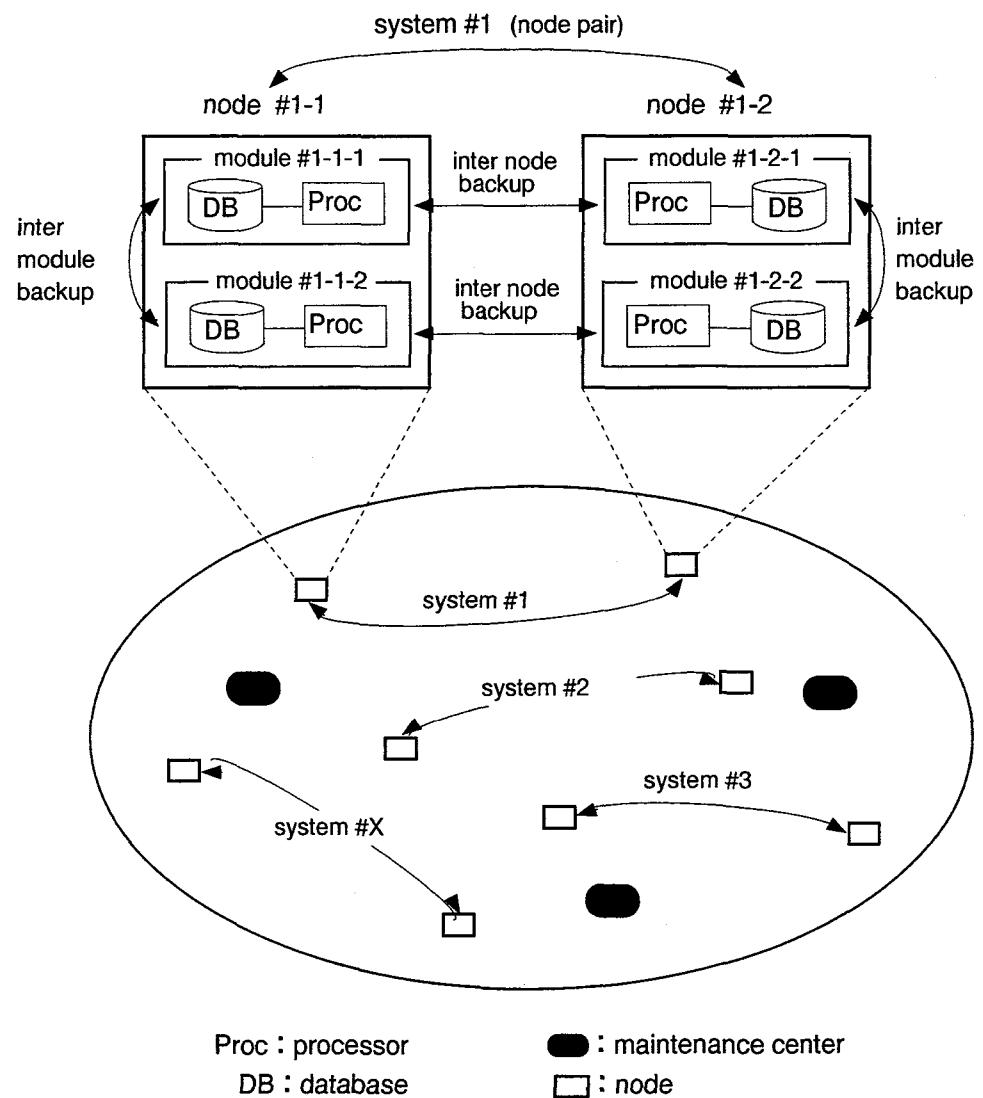


図 7.1 評価対象システムの構成

御が複数のノードでほぼ同時に行われ、システム群の制御が極めて複雑となる。実用性の高いシステムとして、本章では1対1の相互バックアップとした。システムの信頼度構成を図7.2に示す。

7.2.2 モジュールの動作方法

一般に、冗長化により信頼度は向上するが、モジュールの稼働率が低下し、コストアップを招くこととなる。本章では、ノード単体での2重化構成の場合と同様のモジュール稼働率を実現するため、通常時はシステム内の4台（2台×2ノード）のモジュールをすべて動作させ、それぞれのノードで負荷を1/2ずつ分担するとした⁽⁵⁾。1台のモジュールが障害となった場合は、同一ノード内の別のモジュールが障害となつたモジュールのデータベースを格納して負荷を処理する。同一ノード内の2台のモジュールが共に障害となつた場合は、システム内のもう一方のノードのモジュールが共に正常であれば障害となつたそれぞれのモジュールのデータベースを格納して負荷を処理する。もう一方のノードで1台のモジュールが障害となっている場合、すなわち、同一システム内で3台のモジュールが障害となった場合は、正常な1台のモジュールが2ノード分の負荷に対応し、すべてのデータベースを格納して処理する。この場合、正常な1台のモジュールに最大処理能力の2倍の負荷がかかるが、トラヒックの規制により両ノードにかかる負荷の1/2を処理するデグレード運転を行う。1台のモジュールがすべてのデータベースを保持することとなるが、このためのファイル、メモリ等の記憶容量は用意されているものとする。システム内の4台のモジュールがすべて障害となつた場合はシステムダウンとなり当該システムで行っていたサービスが停止する。

モジュールの動作方法として、各モジュールが1/2の負荷を分担する方法のほかに、ノード内の2台のモジュールのうち1台のモジュールが1の負荷を分担し、別の1台のモジュールは待機させておく方法もある。この方法では、すべてのモジュールが正常であったとしても、負荷を分担したモジュールには、常に1の負荷がかかって

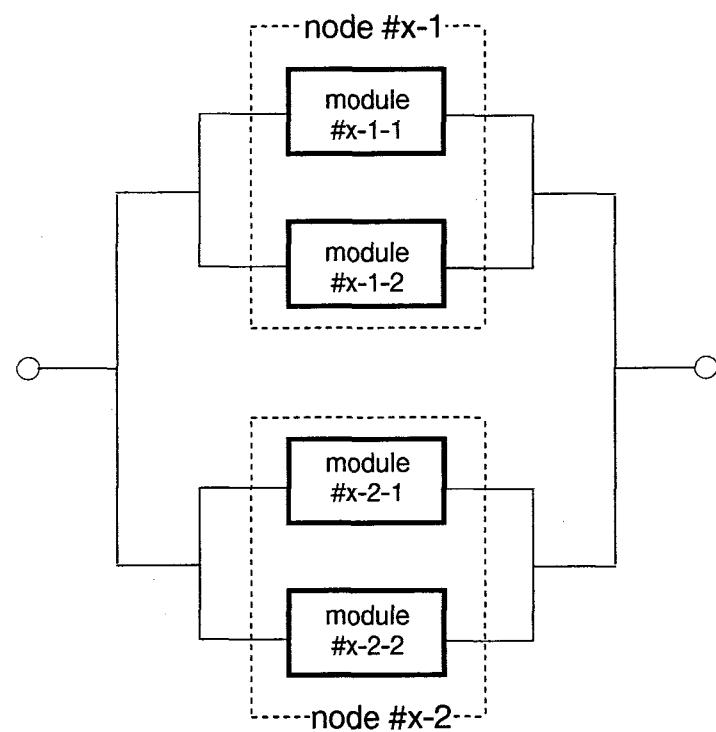


図 7.2 システムの信頼度構成

おり、発生頻度は小さいとしても、突発的な負荷の増加に対してレスポンスタイムを保証することが難しい。一方、本章で想定したように、正常時には各モジュールが1／2の負荷を分担する方法では、モジュールの稼働率に余裕があるため、突発的な負荷の増加に対しても、レスポンスタイムの保証が十分可能である。

7.3 保守方法とシステムの状態

7.3.1 保守方法

保守チームは少數の保守センタに集中配置されるとする。各ノードは、その修理を担当するいずれか1つの保守センタに所属する。各保守センタは障害が報告（以下、駆付け要求と記す）されると保守チームが当該ノードへ駆付けて修理を行う。ソフトウェアの保守は遠隔から行われるものとし、本章ではハードウェアの保守のみを対象とする。渡辺らは二重化されたモジュールを対象とし、ノード内のモジュールが正常か修理中かの状態のほかに、保守チームが駆付けているか否かの状態を加えた状態遷移を提案している⁽³⁷⁾。本章においても、保守チームの駆付けをモデルに組込むため、ノードの状態として保守チームがノードに駆付けているか否かを加えている。

当該ノードに保守チームがいない状態で障害が発生すると、保守センタへ駆付け要求が出される。当該ノードに保守チームが既にいる時は、新たな障害が発生しても、保守チームへの駆付け要求を出さず、既にいる保守チームが新たな障害を含めて修理を行うとした。但し、保守センタは担当するノードを含むシステム内の2ノードの状態を把握しており、駆付けの際、どの程度急ぐかは、次節以下に示すように2ノードの状態により判断可能とする。

修理は保守チームが到着した時点で開始する。一度あるノードに保守チームが駆付けた場合、当該ノードでの障害モジュール数に係わらず、すべてのモジュールを同時に修理できるものとする。すなわち、ノード内の2台のモジュールとも障害となった

場合、両方のモジュールを同時に修理する。保守チームはそのノードでの修理がすべて終了するまで、そのノードに留まって修理を行う。ノード内のすべての障害モジュールの修理が完了した場合は直ちに立ち去る。

7.3.2 システムの状態

システムの状態は各ノードでの障害モジュール数と保守チームが駆付けているか否かにより定義され $S_{ij,kl}$ で表す。ここで、 i 、 k は各ノードの障害モジュール数、 j 、 l は各ノードに保守チームが駆付けているか否か（0 の場合は駆付けていない、1 の場合は駆付けている）を表す。2つのノードに対して同一の方法で保守を行うすると $S_{ij,kl}$ の対称性から図 7.3 に示すようにシステムは 15 通りの状態で表すことができる。具体的には、2つの状態 $S_{i_1j_1,k_1l_1}$ と $S_{i_2j_2,k_2l_2}$ で $i_1 = k_2$ 、 $j_1 = l_2$ 、 $k_1 = i_2$ 、 $l_1 = j_2$ が成立する場合は、2つの状態を $S_{i_1j_1,k_1l_1}$ で縮退して表している。この縮退時の添字の表記としては ij,kl を 4 桁の数字と見たときの小さい方を採用している。図 7.3 の横軸はシステム内の障害モジュール数 ($= i + k$)、縦軸はシステムに駆付けている保守チーム数 ($= j + l$) である。

モジュールの平均故障率および平均修理率（駆付けてから修理に要する平均時間の逆数）はシステムの状態によらず一定とするが、平均駆付け時間はシステムの状態毎に設定する。平均故障率を λ 、平均修理率を μ 、単位時間に保守チームの駆付ける率を $\omega_{ij,kl}$ 、すなわち、平均駆付け時間を $\omega_{ij,kl}^{-1}$ と表す。いずれも指数分布とする。各状態間の遷移は図 7.3 に示す通りとなる。状態 $S_{ij,kl}$ の時刻 t における確率を $P_{ij,kl}(t)$ 、その導関数を $P'_{ij,kl}(t)$ で表すと状態遷移の方程式は以下の通りとなる。

$$P'_{00,00}(t) = -4\lambda P_{00,00}(t) + \mu P_{00,11}(t) \quad (7.1)$$

$$P'_{00,10}(t) = 4\lambda P_{00,00}(t) - (3\lambda + \omega_{00,10})P_{00,10}(t) + \mu P_{10,11}(t) \quad (7.2)$$

$$P'_{00,11}(t) = \omega_{00,10}P_{00,10}(t) - (3\lambda + \mu)P_{00,11}(t) + 2\mu P_{00,21}(t) + 2\mu P_{11,11}(t) \quad (7.3)$$

$$P'_{00,20}(t) = \lambda P_{00,10}(t) - (2\lambda + \omega_{00,20})P_{00,20}(t) + \mu P_{11,20}(t) \quad (7.4)$$

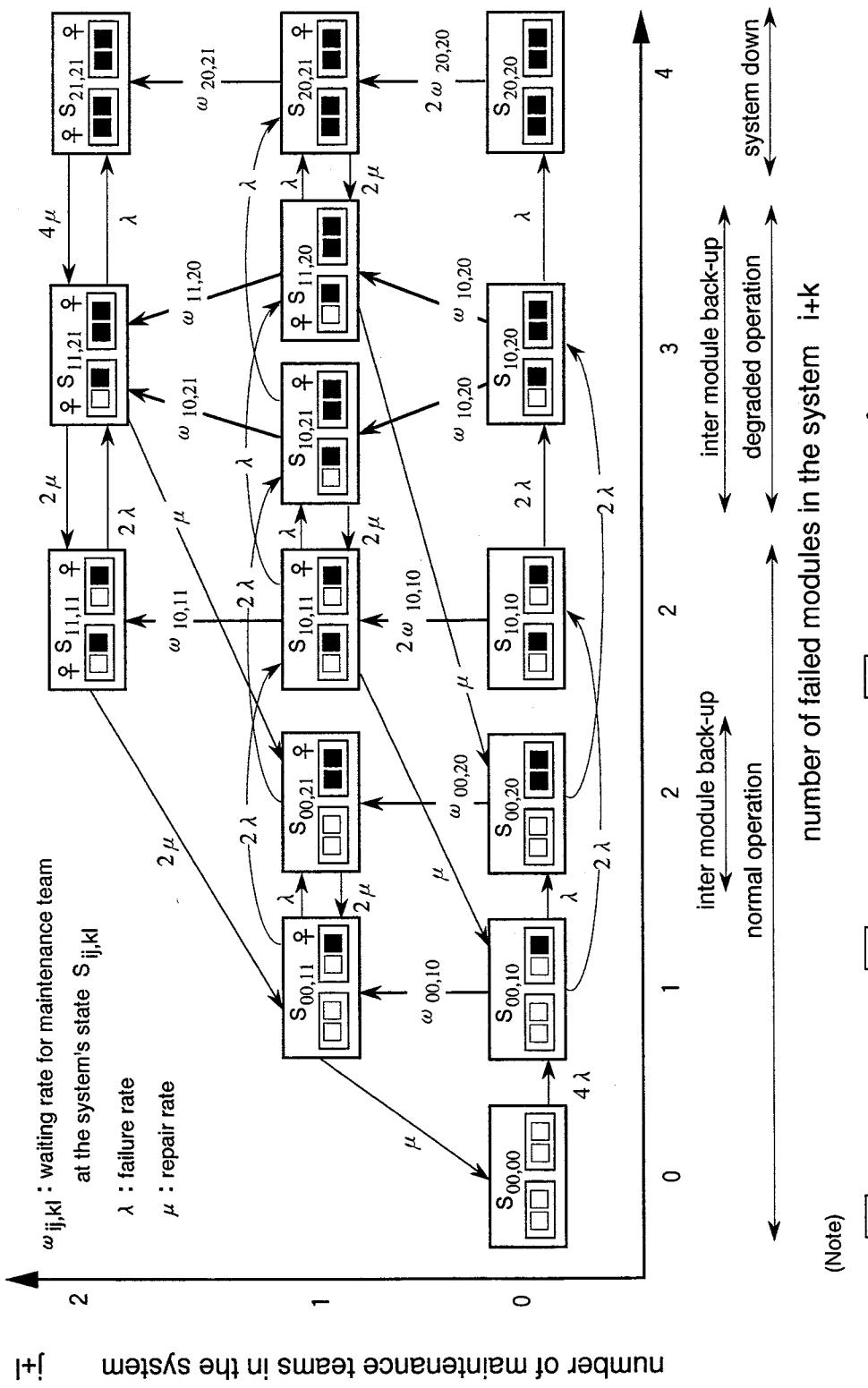


図7.3 状態遷移図

(Note)

$\square\square$: no module failed $\square\blacksquare$: one module failed $\blacksquare\blacksquare$: two module failed \varnothing : maintenance team exist at the node
 $S_{ij,kl}$: system's state. (i,k : number of failed modules in the node j,l : number of maintenance team in the node)

$$P_{00,21}(t) = \lambda P_{00,11}(t) + \omega_{00,20} P_{00,20}(t) - (2\lambda + 2\mu) P_{00,21}(t) + \mu P_{11,21}(t) \quad (7.5)$$

$$P_{10,10}(t) = 2\lambda P_{00,10}(t) - (2\lambda + 2\omega_{10,10}) P_{10,10}(t) \quad (7.6)$$

$$P_{10,11}(t) = 2\lambda P_{00,11} + 2\omega_{10,10} P_{10,10}(t) - (2\lambda + \mu + \omega_{10,11}) P_{10,11}(t) + 2\mu P_{10,21}(t) \quad (7.7)$$

$$P_{11,11}(t) = \omega_{10,11} P_{10,11}(t) - (2\lambda + 2\mu) P_{11,11}(t) + 2\mu P_{11,21}(t) \quad (7.8)$$

$$P_{10,20}(t) = 2\lambda P_{00,20}(t) + 2\lambda P_{10,10}(t) - (\lambda + 2\omega_{10,20}) P_{10,20}(t) \quad (7.9)$$

$$P_{10,21}(t) = 2\lambda P_{00,21}(t) + \lambda P_{10,11}(t) + \omega_{10,20} P_{10,20}(t) - (\lambda + 2\mu + \omega_{10,21}) P_{10,21}(t) \quad (7.10)$$

$$P_{11,20}(t) = \lambda P_{10,11}(t) + \omega_{10,20} P_{10,20}(t) - (\lambda + \mu + \omega_{11,20}) P_{11,20}(t) + 2\mu P_{20,21}(t) \quad (7.11)$$

$$P_{11,21}(t) = \omega_{10,21} P_{10,21}(t) + 2\lambda P_{11,11}(t) + \omega_{11,20} P_{11,20}(t) - (\lambda + 3\mu) P_{11,21}(t) + 4\mu P_{21,21}(t) \quad (7.12)$$

$$P_{20,20}(t) = \lambda P_{10,20}(t) - 2\omega_{20,20} P_{20,20}(t) \quad (7.13)$$

$$P_{20,21}(t) = \lambda P_{10,21}(t) + \lambda P_{11,20}(t) + 2\omega_{20,20} P_{20,20}(t) - (2\mu + \omega_{20,21}) P_{20,21}(t) \quad (7.14)$$

$$P_{21,21}(t) = \lambda P_{11,21}(t) + \omega_{20,21} P_{20,21}(t) - 4\mu P_{21,21}(t) \quad (7.15)$$

$$\begin{aligned} & P_{00,00}(t) + P_{00,10}(t) + P_{00,11}(t) + P_{00,20}(t) + P_{00,21}(t) + P_{10,10}(t) + P_{10,11}(t) + P_{11,11}(t) \\ & + P_{10,20}(t) + P_{10,21}(t) + P_{11,20}(t) + P_{11,21}(t) + P_{20,20}(t) + P_{20,21}(t) + P_{21,21}(t) = 1 \end{aligned} \quad (7.16)$$

7.4 信頼性の評価尺度と保守パラメータの設定

7.4.1 信頼性の評価尺度

サービスへ支障をきたしている状態として、4台のモジュールが障害となりサービスが停止している状態（以下、システムダウンと記す）と3台のモジュールが障害となりシステム全体の1／2の負荷を処理している状態（以下、デグレードと記す）の2通りがある。また、一方のノードがダウンしてノード間バックアップ状態となっている場合は、もう一方のノードで大規模災害が発生すると本来の目的が達成できない。また、大規模災害が生じないまでも、呼やトランザクションのノード間での転送、ノード回復後の切り戻し等、システムの運用が煩雑となるため、これらの状態について

も考慮する必要がある。各状態の発生確率と当該状態の継続時間の2つの面から、保守方式に対する信頼性の評価尺度として以下のものを用いる。

- (1) 平均システムダウン時間と不稼働率：システムダウンの開始から終了までの時間（システムダウン時間）の平均値と、サービス時間に対するシステムダウン時間の総計の占める割合
- (2) 平均デグレード時間とデグレード率：デグレードの開始から終了までの時間（デグレード時間）の平均値と、サービス時間に対するデグレード時間の総計の占める割合
- (3) 平均ノード間バックアップ時間とノード間バックアップ率：ノード間バックアップにより、1つのノードでのみサービスが行われている時間（ノード間バックアップ時間）の平均値と、サービス時間に対するノード間バックアップ時間の総計の占める割合

システムとしての不稼働率を P_{SDN} 、平均システムダウン時間を T_{SDN} 、デグレード率を P_{DEG} 、平均デグレード時間を T_{DEG} 、ノード間バックアップ率を P_{BUP} 、平均ノード間バックアップ時間を T_{BUP} とすると、以下のように表される。ここで、定常状態、すなわち $P_{ij,kl}(t)$ が 0 のときの状態 $S_{ij,kl}$ の確率を $P_{ij,kl}$ と表す。

$$P_{SDN} = P_{20,20} + P_{20,21} + P_{21,21} \quad (7.17)$$

$$T_{SDN} = P_{SDN} / \{\lambda(P_{10,20} + P_{10,21} + P_{11,20} + P_{11,21})\} \quad (7.18)$$

$$P_{DEG} = P_{10,20} + P_{10,21} + P_{11,20} + P_{11,21} \quad (7.19)$$

$$T_{DEG} = P_{DEG} / \{2\lambda(P_{00,20} + P_{00,21} + P_{10,10} + P_{10,11} + P_{11,11}) + 2\mu P_{20,21} + 4\mu P_{21,21}\} \quad (7.20)$$

$$P_{BUP} = P_{00,20} + P_{00,21} + P_{10,20} + P_{10,21} + P_{11,20} + P_{11,21} \quad (7.21)$$

$$T_{BUP} = P_{BUP} / \{\lambda(P_{00,10} + P_{00,11}) + 2\lambda(P_{10,10} + P_{10,11} + P_{11,11}) + 2\mu P_{20,21} + 4\mu P_{21,21}\} \quad (7.22)$$

7.4.2 信頼性目標と保守パラメータの設定

G r a y らは、不稼働率の大小に基づいて、高信頼システムの簡単なクラス分けを行っている。たとえば、不稼働率が 10^{-A} のシステムは、クラス A の信頼度とし、現状の通信システムの信頼度はクラス 5 程度、すなわち不稼働率が 10^{-5} 程度としている。今後、より高信頼のシステムが要求され、その信頼度目標としてクラス 7 程度、すなわち不稼働率が 10^{-7} 程度とすることを提案している⁽⁴³⁾。また、高度 IN の場合は、交換機の上位にサービス制御ノード (SCP) を配置してサービスが行われるため、交換機と SCP のそれぞれの不稼働率が加算されてサービスに影響する。高度 IN サービスを一般の電話サービスと同程度の信頼度でユーザに提供するためには、SCP の不稼働率は交換機の不稼働率 ($\approx 6 \times 10^{-6}$)⁽⁵⁵⁾ に比べて 1 ~ 2 衡小さく、無視できる程度であることが望ましい。これらの状況を踏まえて、二重化された 2 つのノードが相互にバックアップするシステムの不稼働率の目標は 10^{-7} 以下とする。

能條らは不稼働率を極めて小さくしたとしても同時に平均システムダウン時間に制限を加えることを提案している⁽⁵⁷⁾。これは、システムダウンの発生がどんなにまれであるとしても、一度システムダウンとなった場合、それが長時間となることはサービスへの影響が大きくなるためである。ここでの提案も踏まえ、平均システムダウン時間については 1 ~ 3 時間程度とする。

デグレード状態とは、サービス停止には至っていないが、処理能力が 1 / 2 に低下している状態であり、極力発生頻度を押さえる方が望ましい。本章ではデグレード率は不稼働率に比べて 3 衡程度ゆるやかとし 10^{-4} 以下とする。これは 1 年間に 1 時間程度のデグレードに相当する。また、平均デグレード時間については、平均システムダウン時間の数倍を想定し、10 時間程度とする。

同一システム内の 1 つのノードが障害となり、残りの正常な 1 つのノードがサービスを行っているとき、サービス中の正常なノードが大規模災害で破壊されるとサービスが停止し、最悪、データベースが破壊される。この場合は、バックアップ関係にある同一システム内の別のノードを早急に修理してサービスを再開する必要がある。地

震等の大規模災害の影響を定量化することは難しいが、バックアップ中のノードが罹災する頻度を1000年に1回以下、大規模災害の発生頻度を50年に1回程度とすると、ノード間バックアップ率としては1/20以下が必要となる。

また、これらのシステムに適用され、地理的に離れた地点で相互バックアップするような比較的規模の大きい装置の運用実績、開発状況から、平均故障間隔 λ^{-1} はおよそ1000時間から5000時間程度と想定し、ここでは、議論の煩雑化を避けるため、平均的な値として λ^{-1} は2000時間とした。

次に、装置の修理については、高度INやOLT Pに適用される装置の修理状況等を考慮し、平均修理時間 μ^{-1} は1時間とした。また、駆付け保守については、短時間で対応する緊急駆付けと長時間待たせてもよい非緊急駆付けの二種類を想定した。以下、非緊急駆付け時間を T_1 、緊急駆付け時間を T_2 とする。ここで、二重化された単一ノードの場合は、経験的に、駆付けの緊急性はシステムダウンしているか否か等により判断される場合が多い。しかし、二重化された2つのノードが更に相互にバックアップされる場合には、駆付け要求が発生する状態が多く、どのようなときに緊急駆付けとするかの判断が重要となる。

非緊急駆付けの時間 T_1 としては、①1～2日以内、②1週間程度、③1カ月以上、等が考えられる。保守チームの配備に着目すると、①の場合は、夜間は除くとしても、平日の昼間はもとより、土曜や日曜にも対応できる保守チームを配備する必要がある。②は平日勤務で、夜間や土、日曜は休みの保守チームで対応できる。また、①に比べて、駆付けに際し、時間的な余裕があるため、幾つかの障害を待たせる等のスケジュール化が可能であり、少ない保守チームで対応できる。③まで駆付け時間が延びた場合は、障害の発生を契機としてノードに駆付けるのではなく定期保守等で各ノードを巡回したときに、障害となっているモジュールを修理することで対応可能となる。少ない保守チームで、複数の障害に対応可能とするためには、1～2日以内では短すぎるため、本章では非緊急駆付け時間が100時間以上許容されることを条件とする。

7.5 駆付けの緊急度と保守モデル

モジュールの障害が発生し駆付け要求が出される状態 $S_{ab,cd}$ ($i=a, j=b, k=c, l=d$ で、 ab,cd が $00,10/10,11/10,21/00,20/10,10/11,20/20,21/10,20/20,20$ の場合) は 9通りあり、どの程度緊急に駆付ける必要があるか否かの度合いに応じて、この 9通りの状態をレベル分けすることを提案し、どのレベルまで緊急駆付けを行うかに基づいて保守モデルを設定する。本章では不稼働率の時間的な変化、すなわち $P_{SDN}(t)$ に着目し、以下に示す方法により緊急度合いのレベル分けを行った。

駆付け要求が出される状態 $S_{ab,cd}$ の時刻 $t = 0$ における確率 $P_{ab,cd}(0)$ を 1 とし、経過時間 t の関数として不稼働率 $P_{SDN}(t)$ を次式で算出する。

$$P_{SDN}(t) = P_{20,20}(t) + P_{20,21}(t) + P_{21,21}(t) \quad (7.23)$$

$$\text{但し, } \omega_{ij,kl} = 0 \quad (7.24)$$

$$P_{ab,cd}(0) = 1 \quad (7.25)$$

$$P_{ij,kl}(0) = 0 \quad (i \neq a, j \neq b, k \neq c, l \neq d) \quad (7.26)$$

駆付け要求が出される 9通りの状態について、 t と $P_{SDN}(t)$ の関係を図 7.4 に示す。 $P_{SDN}(t)$ は式 (7.1) ~ (7.16) を式 (7.24) ~ (7.26) の条件のもとで数値計算した(付録 7.1 参照)。ここで、 $\omega_{ij,kl}$ を 0 としたのは、すべての駆付けが行われないとしたためである。図 7.4 で黒ぬりのマークを付した線はシステム内の 2 ノード共に保守チームがいない状態、白抜きのマークを付した線は片方のノードに既に保守チームが駆付けている状態を表している。 t が数時間以内では障害モジュール数 ($= a + c$) が多いほど $P_{SDN}(t)$ は高い。しかし、既に一方のノードに保守チームが駆付けている場合は、短時間のうちに当該ノードで修理が行われ t の増加と共に $P_{SDN}(t)$ は小さくなる。 t が 10 時間以上では、そのノードに障害モジュールが無かった場合と同様の傾向を示す。たとえば、 $S_{20,21}$ は 4 台のモジュールが障害となりシステムダウンしているが、既に一方のノードに保守チームが駆付けているため、短時間のうちに当該ノードで 2 台のモジュールの修理が行われ、 t の増加と共に $P_{SDN}(t)$ は $S_{00,20}$ と同様の値となる。図

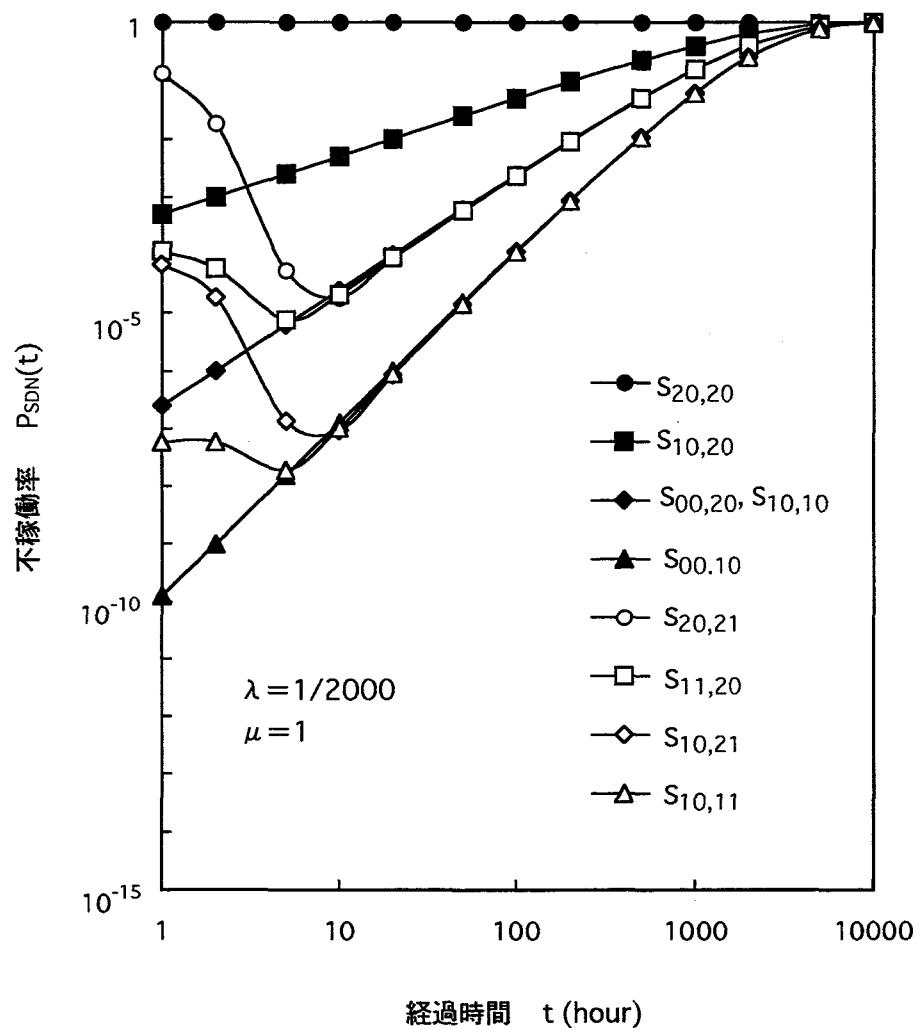


図 7.4 経過時間と不稼働率

7.4 から判るように t が 10 時間以上では大きく 4 つにグループ化できる。同一の経過時間 t で $P_{SDN}(t)$ を比較すると、この大きさは、この状態で駆付け時間を長時間化していくときのリスクの大きさを表しており、駆付けの緊急度はこのリスクに対応すると考えられる。この尺度を緊急レベル H とし、図 7.4 に基づいて 4 段階に分けると、状態 $S_{ab,cd}$ の緊急レベルは次式で表される。

$$H = (a + c) - (a \cdot b + c \cdot d) \quad (7.27)$$

式 (7.27) の第 1 項はシステム内の障害モジュール数、第 2 項は保守チームが到着しているノードでの障害モジュール数である (b, d が 1 の場合は保守チームが到着していることを示す)。すなわち、 H はシステム内で障害となっているが修理が始まっていないモジュール数となっている。

この緊急レベル H により駆付け要求が出される 9 通りの状態を分類して図 7.5 に示す。緊急駆付けをどの状態で行うか否かについて、 H を基に 3 通りの保守モデル(以下、簡単にモデルと記す)を設定する。モデル 0 は H が 4 以上で緊急駆付けを行うモデルである。モデル 1、モデル 2 はそれぞれ H が 3 以上、 H が 2 以上で緊急駆付けを行うモデルである。図 7.5 にモデル毎にどの状態で緊急駆付けを行うかを併せて示す。また、駆付け要求が出される 9 通りの状態の駆付け時間 $\omega_{ab,cd}^{-1}$ として非緊急駆付け時間 T_1 と緊急駆付け時間 T_2 のどちらを設定するかを整理して表 7.1 に示す。

7.6 信頼性と保守方法の関係

7.6.1 システム信頼性の評価

要求されるシステム信頼性を満足した上で、保守方法に影響する緊急、非緊急駆付け時間として、どのような時間設定が可能であるか評価する。なお、システム信頼性は、はじめに表 7.1 に示した駆付け時間を式 (7.1) ~ (7.16) に代入して、定常状

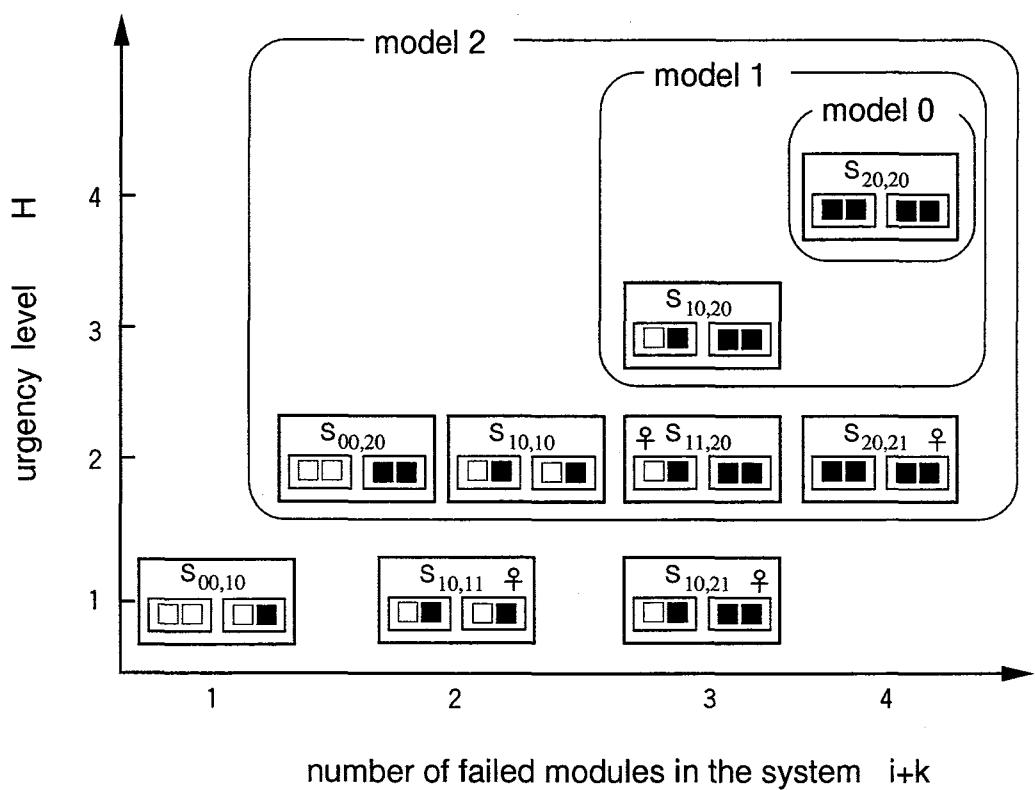


図 7.5 緊急レベルと駆付けモデル

表7.1 駆付けモデルと設定する駆付け時間

緊急 レベル	状態	駆付け時間	設定する駆付け時間		
			モデル0	モデル1	モデル2
4	$S_{20,20}$	$\omega_{20,20}^{-1}=$	T_2	T_2	T_2
3	$S_{10,20}$	$\omega_{10,20}^{-1}=$	T_1	T_2	T_2
2	$S_{20,21}$	$\omega_{20,21}^{-1}=$	T_1	T_1	T_2
	$S_{11,20}$	$\omega_{11,20}^{-1}=$	T_1	T_1	T_2
	$S_{10,10}$	$\omega_{10,10}^{-1}=$	T_1	T_1	T_2
	$S_{00,20}$	$\omega_{00,20}^{-1}=$	T_1	T_1	T_2
1	$S_{10,21}$	$\omega_{10,21}^{-1}=$	T_1	T_1	T_1
	$S_{10,11}$	$\omega_{10,11}^{-1}=$	T_1	T_1	T_1
	$S_{00,10}$	$\omega_{00,10}^{-1}=$	T_1	T_1	T_1

T_1 : 非緊急駆付け時間

T_2 : 緊急駆付け時間

態における各状態の発生確率を数値計算し、次にそれらを式 (7.17) ~ (7.22) に代入して算出した。

非緊急駆付け時間 T_1 と平均システムダウン時間 T_{SDN} の関係を図 7.6、緊急駆付け時間 T_2 と平均システムダウン時間の関係を図 7.7 に示す。平均システムダウン時間は T_1 に係わらずほぼ一定となるが、 T_2 に対してはモデルに関係なく約 $T_2/2$ となる。平均システムダウン時間として 1 ~ 3 時間を条件としたため、 T_2 は 2 ~ 6 時間以下に抑える必要がある。以下、 T_2 は 2 ~ 6 時間として、 T_1 と不稼働率、デグレード率、平均デグレード時間、ノード間バックアップ率の関係について評価し、 T_1 として保守のスケジュール化が可能な 100 時間以上許容されるか否かについて考察する。

T_1 と不稼働率 P_{SDN} の関係を図 7.8 に示す。Gray の信頼度分類によるクラス 7、すなわち、不稼働率が 10^{-7} 以下を条件としたため、モデル 0 は T_1 が 40 ~ 50 時間程度となり若干不足する。モデル 1 は 100 ~ 400 時間程度許容される。モデル 2 は非緊急駆付けを全く行わなくとも良く、過剰な保守方式となる。

T_1 とデグレード率 P_{DEG} の関係を図 7.9 に示す。デグレード率として不稼働率より 3 衡程度大きい 10^{-4} 以下を条件としたため、モデル 0 では T_1 は 50 時間程度、モデル 1 では 100 ~ 200 時間程度許容される。モデル 2 では制約が無くなる。 T_1 と平均デグレード時間 T_{DEG} の関係を図 7.10 に示す。平均デグレード時間は、モデル 1、2 の場合は、 T_1 に依存せずほぼ一定となる。しかし、モデル 0 の場合は、 T_1 にはほぼ比例して平均デグレード時間が長くなる。平均デグレード時間として 10 時間程度を条件としたため、モデル 0 では T_1 に許容される時間は 20 時間程度となり、モデル 1、2 では制約が無くなる。

T_1 とノード間バックアップ率 P_{BUP} の関係を図 7.11 に示す。ノード間バックアップ率を $1/20$ 以下とするためには、モデル 0、1 の場合は T_1 を 200 時間以下に抑える必要があるが、モデル 2 では T_1 への制約はなくなる。

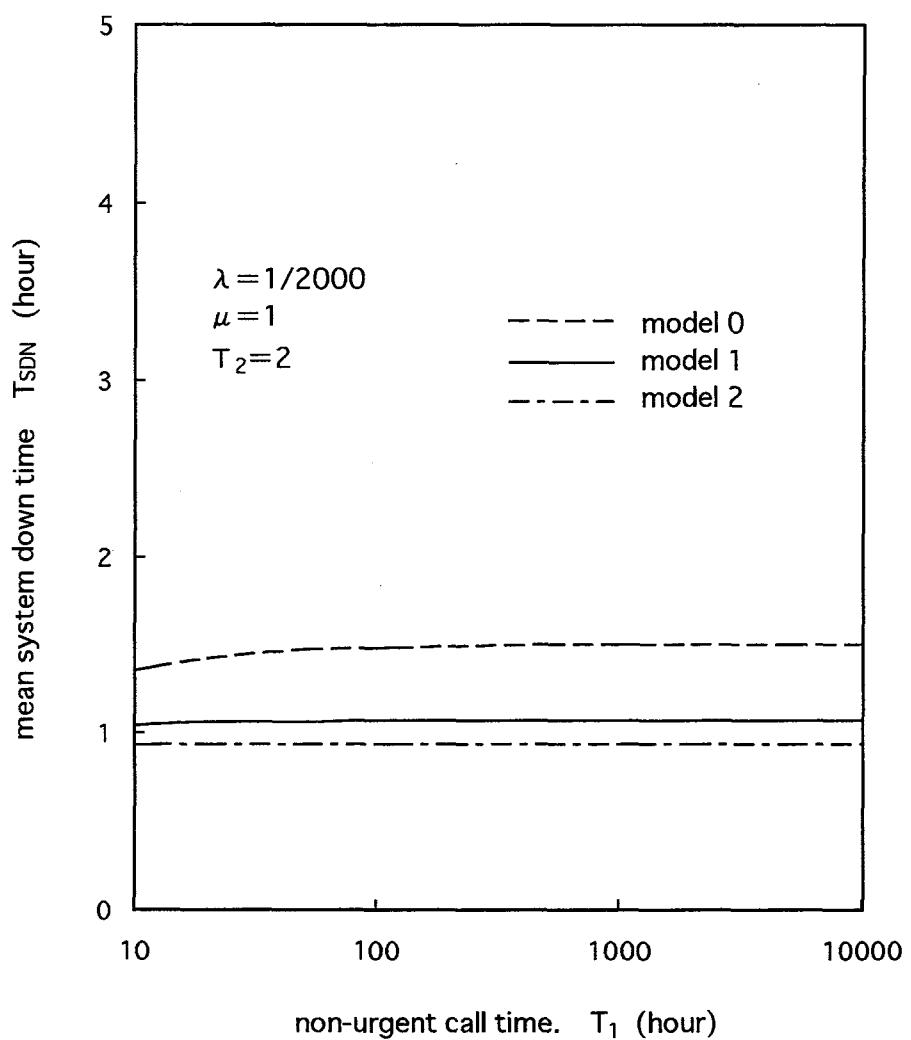


図 7.6 非緊急駆付け時間と平均システムダウン時間

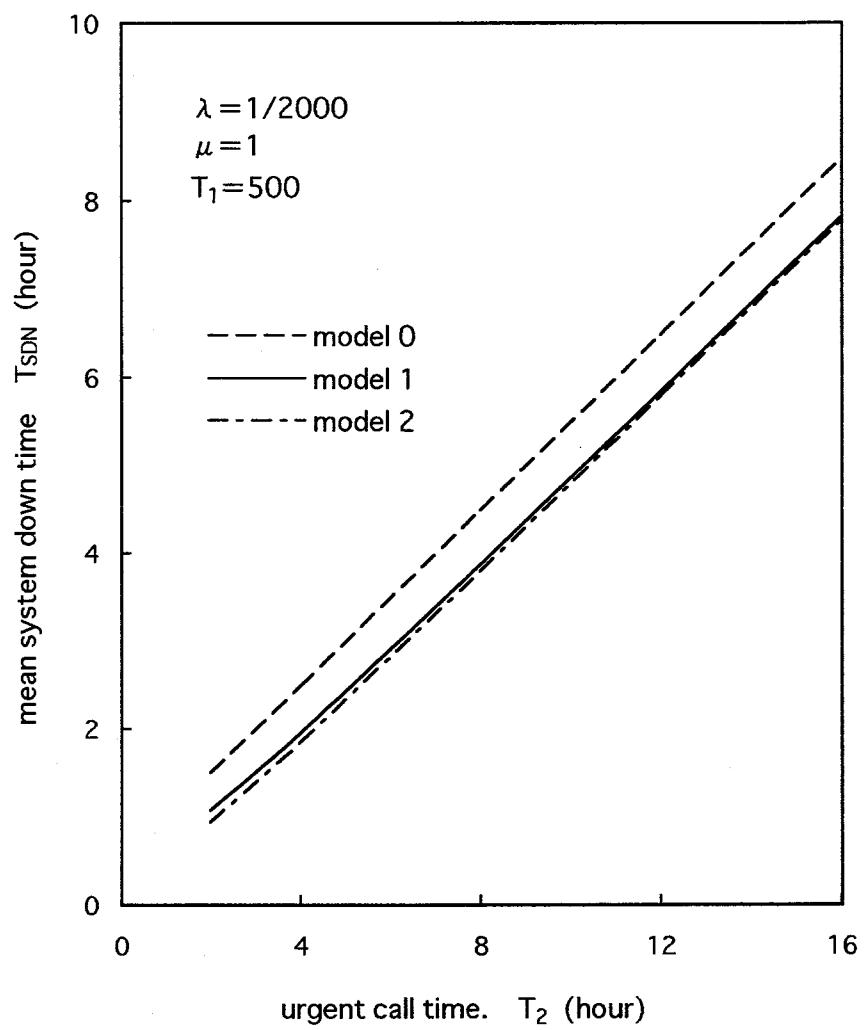


図 7.7 緊急駆付け時間と平均システムダウン時間

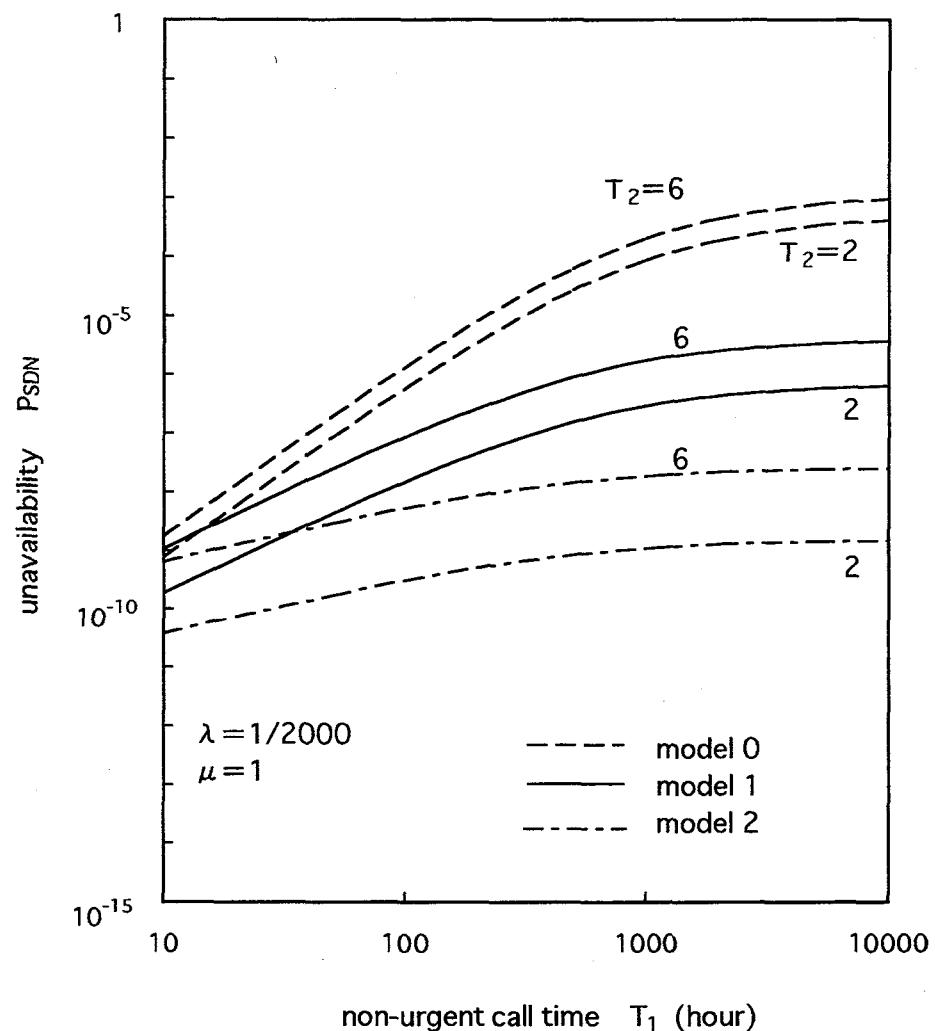


図 7.8 非緊急駆付け時間と不稼働率

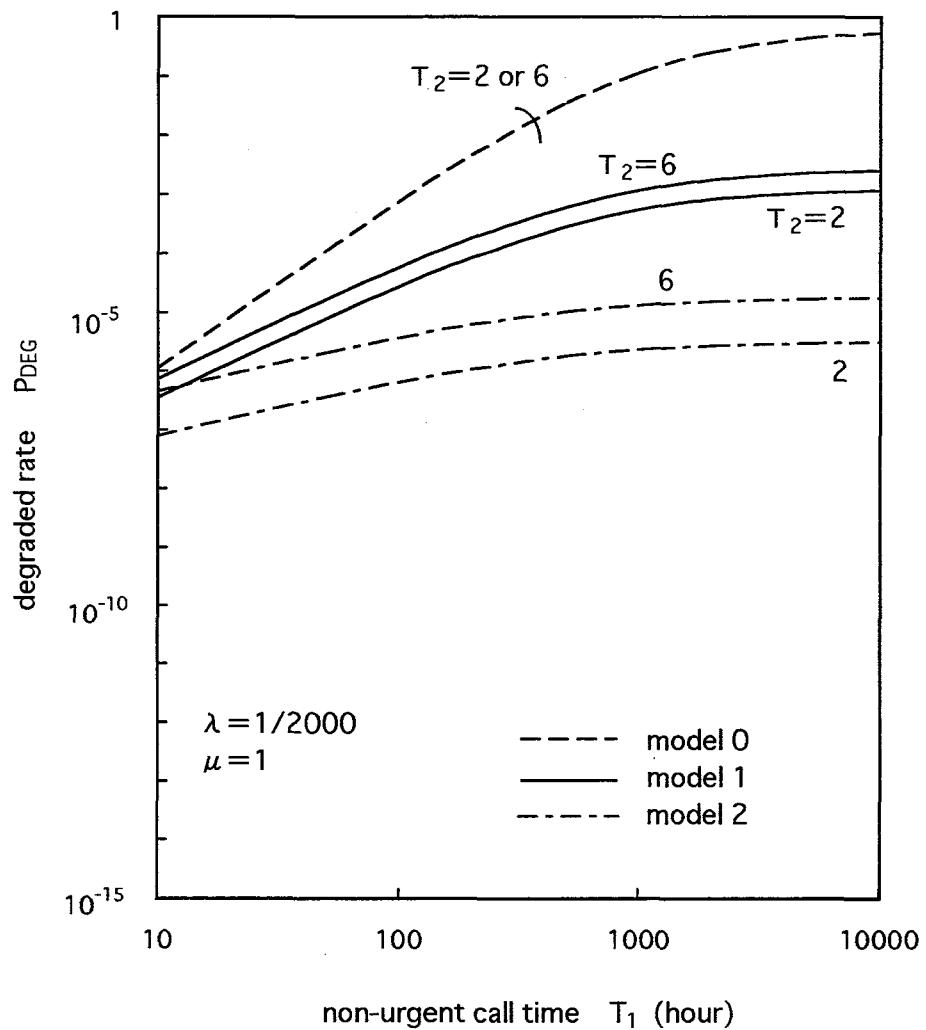


図 7.9 非緊急駆付け時間とデグレード率

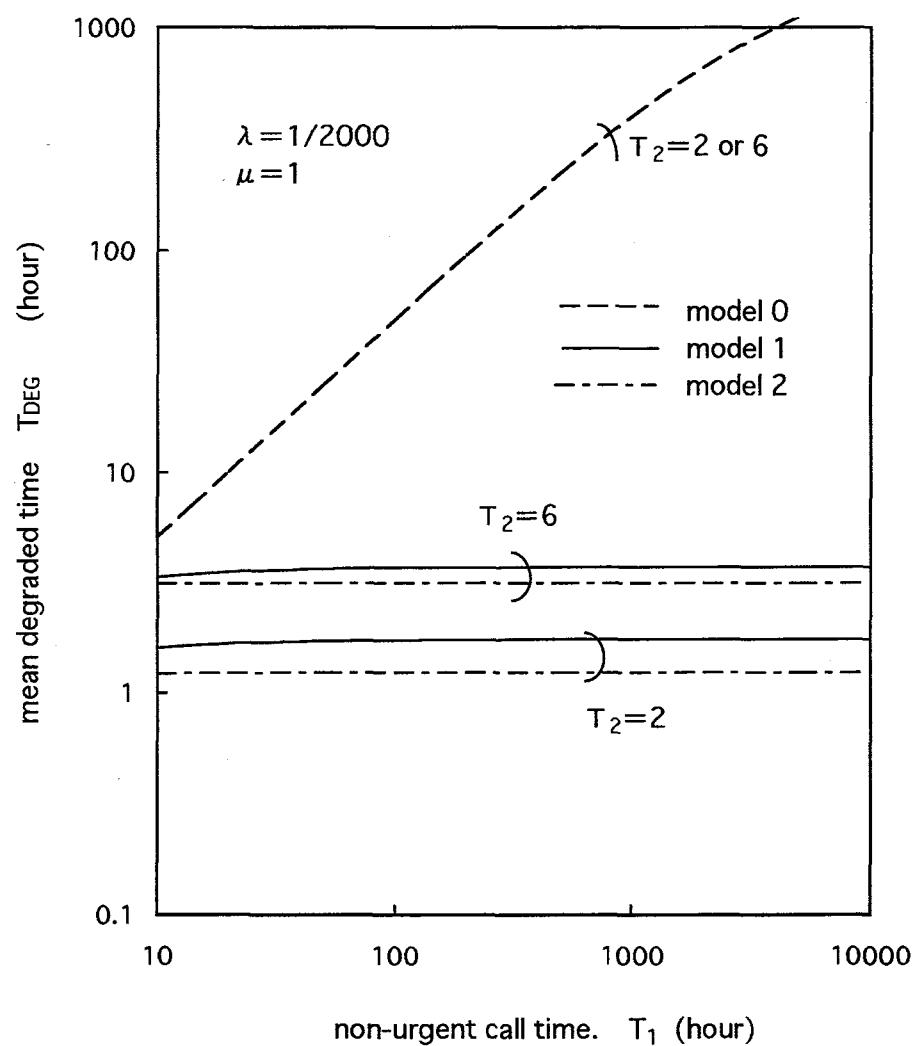


図 7.10 非緊急駆付け時間と平均デグレード時間

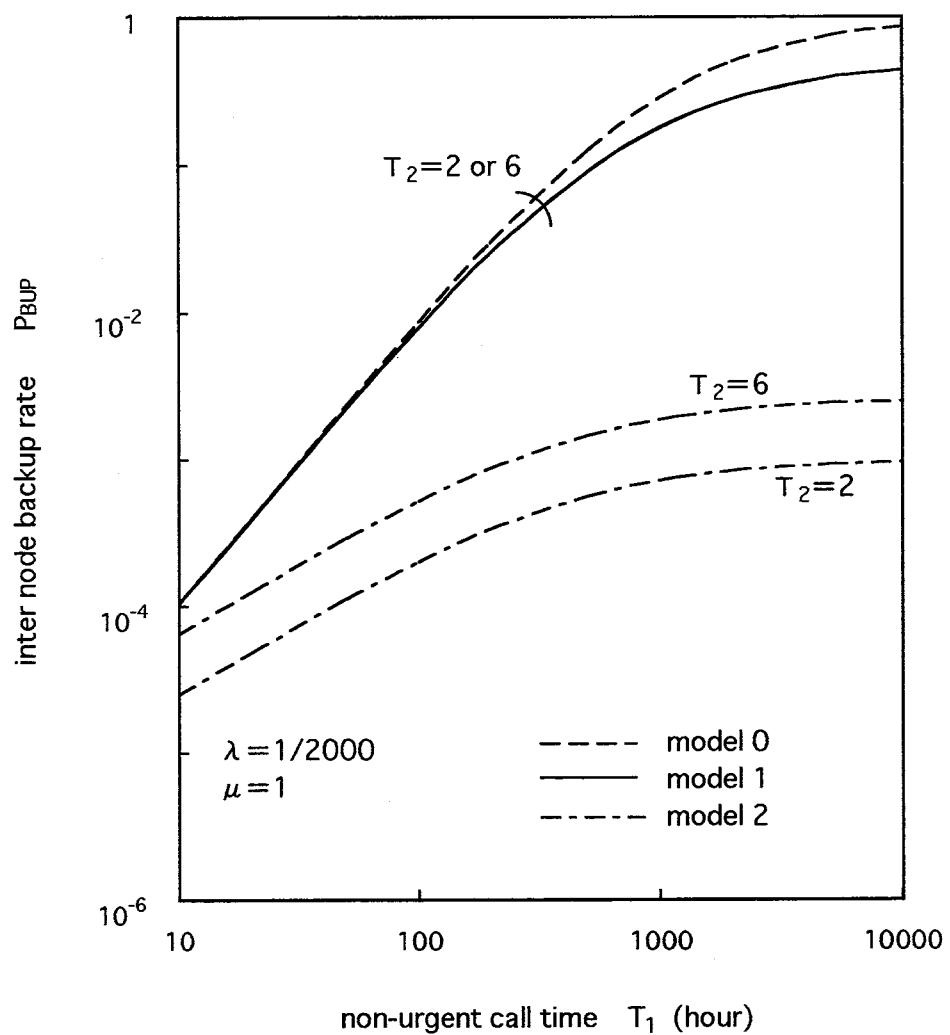


図 7.1.1 非緊急駆付け時間とノード間バックアップ率

7.6.2 駆付け回数の評価

本節では、駆付け回数の面から各モデルの評価を行う。緊急駆付け要求が出される状態の確率の総和を P_{ECL} 、非緊急駆付け要求が出される状態の確率の総和を P_{OCL} とすると、1年間で1つのシステムに対して発出される緊急駆付け回数 N_{ECL} および非緊急駆付け回数 N_{OCL} は、1年間の延べ時間が 8760 時間であることから、次式で表される。

$$N_{ECL} = 8760 \cdot P_{ECL} / T_2 \quad (7.28)$$

$$N_{OCL} = 8760 \cdot P_{OCL} / T_1 \quad (7.29)$$

ここで、総駆付け回数を N_{TCL} ($= N_{ECL} + N_{OCL}$) として、 T_1 と N_{TCL} の関係を図 7.1.2 に示す。図 7.1.2 の下部には、不稼働率を 10^{-7} 以下とした場合の各モデルの T_1 の許容範囲についても併せて示している。保守のスケジュール化が可能と考える T_1 が 100 時間程度で、総駆付け回数を比較するとモデル間での差はほとんど無い。なお、総駆付け回数はモデルに関係なく T_1 を長くすることにより減少する。これは T_1 の増加につれて必ずしも 1 台のモジュールが障害となる毎に駆付けて修理するのではなく、2 台のモジュールが障害となった時点で駆付ける場合が発生しているためである。

T_1 と緊急、非緊急駆付け回数の関係を図 7.1.3 に示す。保守の効率化の点では緊急駆付け回数が少ない方が有利である。 T_1 が 100 時間で比べると、モデル 2 はシステムあたりの緊急駆付け回数が年に 2 回程度となり、モデル 1 とモデル 0 では 0.1 回程度となる。図 7.1.2、図 7.1.3 では、 T_2 は 2 時間として評価したが、 T_2 を 6 時間としても緊急、非緊急駆付け回数は変化なかった。

7.6.3 保守形態の及ぼす影響

保守を効率化するポイントは緊急駆付け回数が少ないと、非緊急駆付けについては保守のスケジュール化が可能と考える T_1 が 100 時間程度許容されることである。 T_1 として許容される時間や緊急駆付け回数について、前節までの評価結果をまとめ

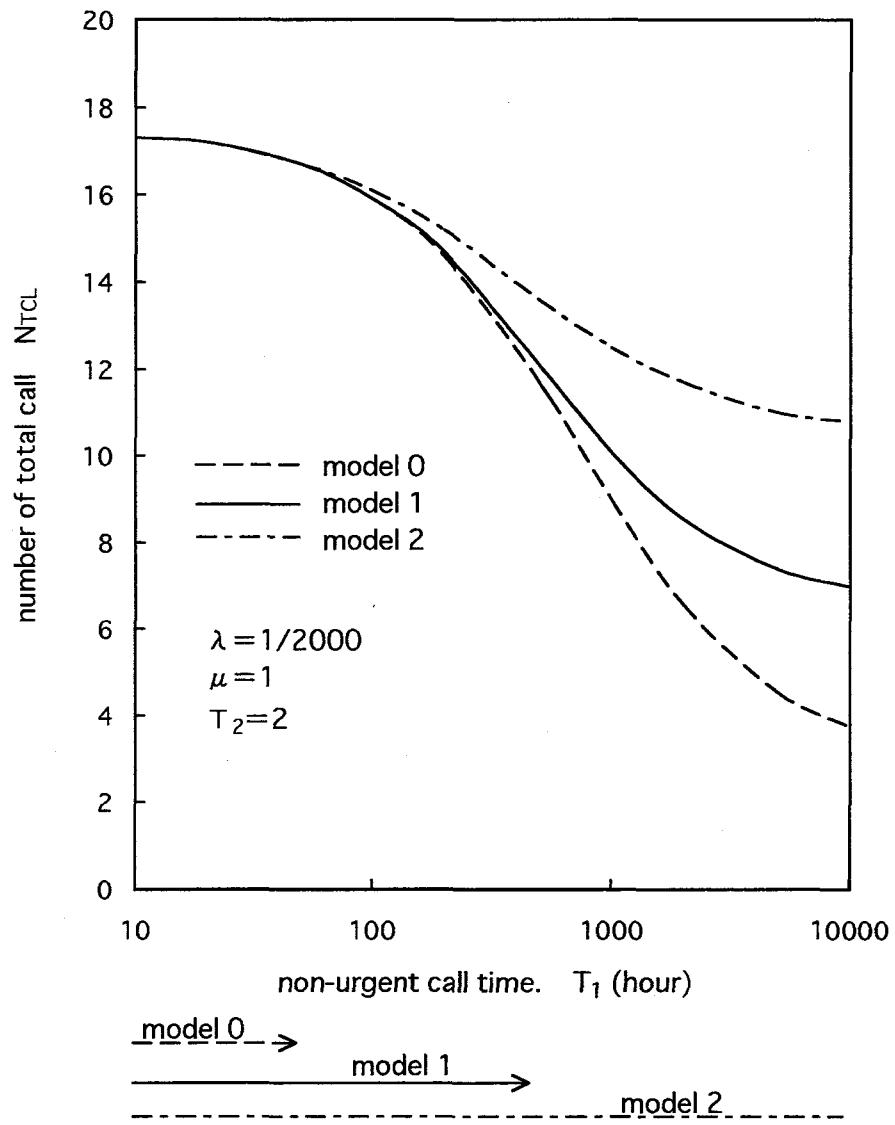


図 7.1.2 非緊急駆付け時間と総駆付け回数

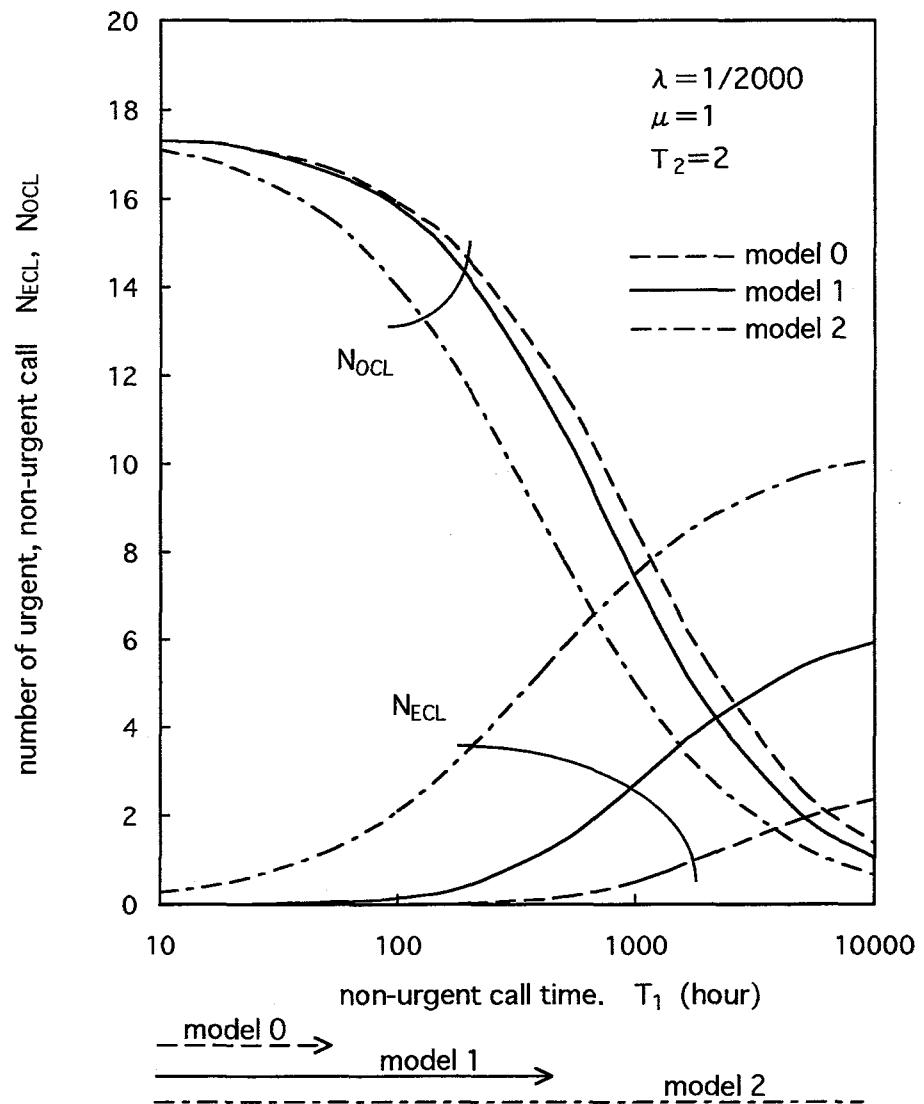


図 7.1.3 非緊急駆付け時間と緊急, 非緊急駆付け回数

て表7.2に示す。表7.2より、モデル0は平均デグレード時間が T_1 に比例して長くなるため、 T_1 を延ばすことは難しい。 T_1 として許容される時間は20時間程度と短く、保守のスケジュール化が十分行えない。モデル1では T_1 を100時間程度と大きく設定でき、保守のスケジュール化が十分可能である。モデル2は信頼性の面からは非緊急駆付け時間 T_1 への制約はないが、モデル1に比べて緊急駆付け回数が20倍程度多いという問題を持つ。

全体として、相互にバックアップした二重化ノードの保守としては、緊急レベル3以上で緊急駆付けを行うモデル1、すなわち、3~4台のモジュールが障害で、かつ両方のノードに保守チームがいない状態($S_{10,20}$ 、 $S_{20,20}$)のとき緊急駆付けを行うモデルが有用であると考える。モデル1の場合、すべてのモジュールが障害でシステムダウンしているが、既に一方のノードに保守チームが駆付けている状態($S_{20,21}$)では緊急駆付けを行わない。この状態に対して、緊急駆付けを行うとしてもシステム信頼度の向上効果はほとんどなく、7.5節で提案した緊急レベルに基づくレベル分けが有効と言える。

7.6.4 システム、保守センタの配置方法

これまでの検討結果を踏まえ、システムや保守センタの配置方法について簡単に考察する。保守センタからノードまでは数時間以内の緊急駆付けを可能とする必要がある(以下、保守センタから緊急駆付け可能なエリアをセンタエリアと記す)。このためには、たとえば、車による駆付けを行うとした場合、センタエリアは保守センタからおおよそ100Km半径のエリア内となる。また、広域的な被害を被る地震、水害等でノードのペアが同時に被害を被らないためには、ノード間を数100Km以上離す必要がある。加えて、保守センタ数を極力少なくして、保守チームの集約化を図ることも必要である。これらの点を考慮したシステムおよび保守センタの配置の一例を図7.14に示す。図において、保守センタは2カ所に配置し、保守センタ相互間は数100Km以上離す。すべてのノードはいずれかのセンタエリア内に配置され、同

表7.2 駆付けモデルの評価結果

	設定条件	モデル0	モデル1	モデル2
非緊急駆付け時間(T_1)として許容される時間(時間)	平均システムダウン時間(1~3時間)	制限なし	制限なし	制限なし
	不稼働率(10^{-7} 以下)	20 ~ 40	200 ~ 400	制限なし
	デグレード率(10^{-4} 以下)	50	100 ~ 200	制限なし
	デグレード時間(10時間)	20	制限なし	制限なし
	ノード間バックアップ率(1/20以下)	200	200	制限なし
緊急駆付け回数(回/システム/年)	~0	0.1	2.0	

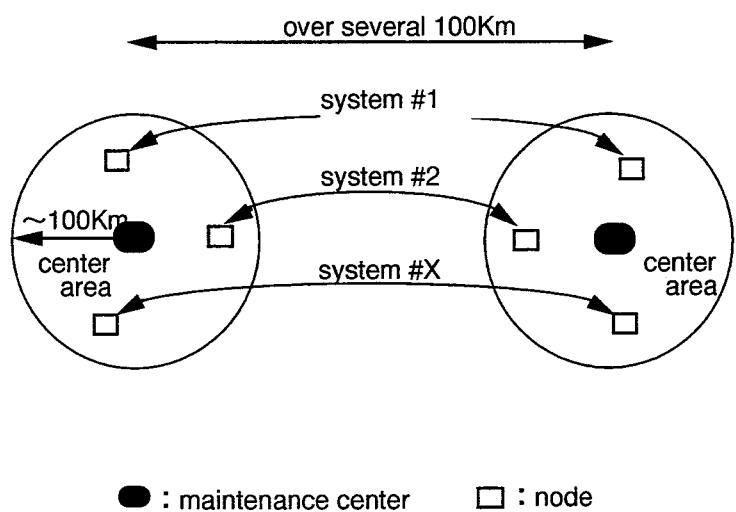


図7.1.4 システム、保守センタの配置方法

一システム内のノードのペアはそれぞれ別のセンタエリア内に配置する。

7.7 結言

データベースを持ち、ノード内は二重化され、さらに2つのノード間で相互にバックアップされたシステム群を対象として、駆付け保守とシステム信頼性の関係について明らかにした。保守の方法は、障害が発生した時点で、システムの障害状態に応じて駆付ける方式とし、緊急、非緊急の二種類とした。このとき、障害発生時のシステム信頼度の時間的な変化に注目して、駆付けの緊急度合いをシステム内の障害装置数や保守チームが駆付けているか否かの状態に基づいてレベル分けすることを提案し、どのレベルまで緊急駆付けを行うかに応じて保守モデルを設定して信頼性や駆付け回数の評価を行った。これらの結果に基づいて、システムに要求される信頼度条件を満足し、非緊急駆付けに関しては保守のスケジュール化が可能な100時間程度許容され、緊急駆付けに関しては駆付け回数の少ない効率的な保守方法を明らかにした。

高い信頼性が要求される高度INシステムや、OLTPシステムでは地震や火災等に際してもサービスを継続するため、離れた地点での相互バックアップを行う方向にある。ここで示した保守の評価方法は、このバックアップを単に冗長度の増加としてとらえるのではなく、保守の効率化に活用するという面で、これら多くの高信頼化システムの保守方法を決定する上で有用である。

第8章 結論

本論文では、高い拡張性を持ち、高性能で高信頼なサービス制御ノードを実現するためのハードウェア構成技術や、全国的な高度 IN サービスを提供するためネットワーク内に多数設置されたサービス制御ノードの効率的な運転・保守技術について研究した。具体的には、需要の変動や機能要求へ柔軟に対応する分散処理化、サービス毎カスタマ毎のデータを効率的に管理するデータベース化、運転保守の効率化とのバランスを考慮した経済的な高信頼化の各面からアプローチし、（1）多数のモジュールによるデータベースの分散構成技術、（2）モジュール間の負荷の偏りの評価技術、（3）モジュール間の効率的な結合技術、（4）メモリデータベースの効率的なりかわり技術、（5）多数のサービス制御ノードの経済的な運転保守技術、を明らかにした。また、試作を通して、提案した技術の有効性を確認した。本研究で得られた主な成果は以下の通りである。

（1）データベースの分散構成技術

データベースのアクセス特性、モジュールの負荷やその偏りなどを総合的に考慮して、データベースの分散構成方式の評価方法を明らかにした。これに基づき、サービス制御ノードとして効率的な分散型メモリデータベース構成法を提案した。

具体的には、呼処理に必要なサービスデータや、複数の SCP に SLP を分散配備したときの配備先を特定する SLP ディレクトリ等、サービス制御ノードとして必要なデータベースの容量やアクセス頻度から、データベースの格納媒体として OLTP システムで一般に用いられている磁気ディスク装置は適用できず、主メモリや半導体ディスク装置に格納する必要があることを示した。また、データベースの格納媒体や、複数のモジュールからのデータベースの共用の可否に基づいてデータベースの分散構成の代替案を設定し、モジュール間の負荷の偏り、データベース処理でのダイナミックステップ数、データベースの格納コストなどを総合的に考慮した評価方法を明ら

かにした。この結果に基づき、更新系データベースであるサービスデータについては、データベースを各モジュールに分散して配備し、主メモリに格納する分割方式を提案し、参照系データベースである SLP ディレクトリについては全てのデータベースを各モジュールに重複して配備し、主メモリに格納する多重方式を提案した。データの高信頼化に関しては、ハードウェア障害やソフトウェアのバグによりデータが失われた場合のサービスへの影響度合いに基づく、データの復元条件を明らかにした。また、アベイラビリティの向上に関しては、モジュール間バックアップとノード間バックアップを組み合わせた効率的な冗長構成法を提案した。

分散型メモリデータベース構成と、モジュール間、ノード間での相互バックアップを組み合わせた冗長構成により、高性能で高信頼なサービス制御ノードを経済的に実現することができる。

(2) モジュール間の負荷の偏りの評価技術

カスタマ毎のサービスデータに対するトラヒックが大きくばらつく場合について、トラヒック量の近似方法を提案した。また、これを用いて、複数のモジュールにサービスデータを分散配備した場合のモジュール間の負荷の偏りの評価方法を明らかにした。

具体的には、最もトラヒックの高いカスタマと最もトラヒックの低いカスタマのトラヒック比が数千倍から数万倍と大きくばらつく場合について、各カスタマのトラヒックの高いものからの順位と、トラヒックのばらつきの度合いを表す偏り係数を用いて近似する方法を提案した。この近似方法は、高度 IN サービスのようにトラヒックの高い少数のカスタマが全体のトラヒックの大半を占めるような場合に良く近似できることを示した。これを用い、カスタマ毎のサービスデータを複数のモジュールに分散配備した場合のモジュール間の負荷の偏りの評価方法を明らかにした。この負荷の偏りは、モジュールに収容するサービスデータ数の平方根に反比例すること、偏り係数が 1 のとき極大値を取ることを示した。

また、高度 IN サービスのトランザクションの処理量や、現状のプロセッサ性能か

ら、モジュールに収容できるサービスデータ数は数万から10万程度であり、このとき、モジュール間の負荷の偏りは1～2割程度に収まることを明らかにした。さらに、このモジュール間の負荷の偏りについては、SLPディレクトリ検索処理とサービスデータに基づくデータベース処理を組み合わせて、モジュール間の負荷を動的に平準化する方法を提案した。

負荷の偏りや平準化技術は分散処理によるサービス制御ノードの方式設計、設備設計に有効である。

(3) モジュール間の結合技術

大規模分散処理システムの結合機構の高信頼化構成法、通信処理能力を向上させる通信制御チャネルの制御方式およびモジュール間の通信時間の測定方法について提案した。

モジュール間の接続にはATMを適用し、結合機構を各モジュールと接続する機能、信号を多重化／分離化する機能、スイッチングを行う機能により実現し、機能単位に二重化し、接続するモジュール台数が50モジュール程度とした大規模分散処理システムにおいても、システムの不稼働率を $10^{-7} \sim 10^{-6}$ 以下にすることができた。各モジュールの通信制御については、ある周期毎にまとめて複数のメッセージを処理する方式を取り上げ、周期時間と通信時間の関係を明らかにし、高度INのSCPでは数10ミリ秒の周期で処理することが有効であることを示した。

ここで提案したモジュール間の結合技術は、高度INのサービス制御ノードに限らず、適用業務の拡大から多数のモジュールによる分散処理化が進んでいるOLTSPシステムにおいても、高信頼で経済的な分散処理システムを構築する上で有用である。

(4) メモリデータベースのリカバリ技術

障害発生からデータベースが回復するまでの時間と、データベースの単位容量当たりのトラヒック量に着目し、メモリデータベースのリカバリ方式の評価方法を明らかにした。これに基づき、サービス制御ノードのリカバリ方式として、半導体ディスク

装置にメモリデータベースのログとチェックポイント時点でのデータベースを取得する効率的なリカバリ方式を提案した。

主メモリ上のデータベースが破壊された場合、データベースを復元するためには、あらかじめ復元に必要なログを外部の記憶装置に取得しておく必要がある。外部の記憶装置の形態としては、入出力による方法と通信による方法に大きく2大別され、それぞれを代表して、半導体ディスク装置を適用する方式と、別のモジュールの主メモリを適用する方式を代替案として設定した。この2つの代替案について、スループット、リカバリ時間に対する評価を行い、モジュールに収容可能なデータベース容量に着目し、両方式の適用領域を明らかにした。これに基づき、サービス制御ノードのメモリデータベースのリカバリ方式として、半導体ディスク装置にログとチェックポイント時点でのデータベースを取得する、効率的なリカバリ方式を明らかにした。また、半導体ディスク装置に取得したログ等の消失を防止するための不揮発化に関しては、直接バッテリでバックアップする方法と、データを磁気ディスク装置に移す方法を比較評価し、サービス制御ノードとして要求される20時間程度のバックアップでは直接バックアップする方法がハードウェア量や、制御の簡単化の点で有利であることを示した。また、半導体ディスク装置をバックアップするバッテリについては、新たに劣化判定機構を開発し、保守の容易な方法を明らかにした。

メモリデータベースのリカバリ方式の評価方法は、サービス制御ノードに限らず、トランザクション処理の高性能化のためメモリデータベース化が進んでいる多くのOLTPシステムにおいても有効である。

(5) 高信頼化運転保守技術

データベースを保有し、二重化された2つのノードを離れた地点に配置して相互にバックアップするシステムを対象とし、これらがネットワーク内に多数配置された場合について、駆付け保守の方法と信頼性の関係を明らかにし、高い信頼性と平均的に長い駆付け時間の双方を満足する効率的な運転保守方式を提案した。

保守チームは保守センタに集中配置され、システムの障害状態に応じて短時間で対

応する緊急駆付けと、長時間待たせてもよい非緊急駆付けの二種類の駆付けにより保守する方式を想定した。このとき、システムの信頼性の時間的な変化に注目し、駆付けの緊急度合いをシステム内の障害装置数や、保守チームが到着しているか否かの状態によりレベル分けし、このレベルにより駆付け保守の方法を設定することを提案した。これに基づいて、システムの信頼性や駆付け回数を総合的に評価し、システム内で3～4台のモジュールが障害で、かつ両方のノードに保守チームがいない状態のとき緊急駆付けを行い、それ以外の状態では複数の障害に対して保守のスケジュール化を行うなど非緊急に対応することが、保守の効率化の点で有効であることを明らかにした。さらに、この結果を踏まえ、サービス制御ノードや保守センタの配置方法を明らかにした。

高い信頼性が要求される高度INシステムや、OLTPシステムでは地震や火災等に際してもサービスを継続するため、離れた地点での相互バックアップを行う方向にある。ここで示した保守の評価方法は、このバックアップを単に冗長度の増加としてとらえるのではなく、保守の効率化に活用するという面で、これら多くの高信頼化システムの保守方法を決定する上で有用である。

本研究により、拡張性のよい、高性能で、高信頼なサービス制御ノードを、経済的に構成することができる。また、多数のサービス制御ノードの効率的な運転保守が可能となる。本研究の成果を組み込んだ装置は、既に、1997年11月からPHSサービスやローミングサービスに使用されている⁽⁶⁶⁾。引き続き、パーソナル通信サービス等、高度インテリジェントネットワークを用いた高機能な電話サービスに用いられる予定である。高度INサービスは、今後、多くの通信通信事業者により提供されることが予想されるが、本研究により提案した技術は通信事業者の区別なく適用できる。本研究が高度インテリジェントネットワークにおけるサービス制御ノードの構成方式や運転保守方式の進展の一助となれば幸いである。

参考文献

- (1) S. Suzuki, "IN Rollout in Japan," IEEE commun. mag., vol.31, pp.48-55, March 1993.
- (2) E. Cancer, R. Mccann, and M. Aboudharam, "IN Rollout in Europe," IEEE commun. mag., vol.31, pp.38-47, March 1993.
- (3) P. A. Russo, K. Bechard, E. Brooks, R. L. Corn, R. Gove, W. L. Honig, and J. Young, "IN Rollout in the United States," IEEE commun. mag., vol.31, pp.56-63, March 1993.
- (4) 鈴木滋彦, "高度インテリジェントネットワーク," 電子情報通信学会誌, vol.77, no.4, pp.410-415, April 1994.
- (5) 平野正則, 鈴木孝至, 塩澤恒道, 芳西崇, 木ノ内康夫, "分散処理による高度 IN用サービス制御ノードの構成," 信学論(B-I), vol.J79-B-I, no.8, pp.539-550, Aug. 1996.
- (6) 石川弘, 石川秀樹, "ネットワークの高度化とサービスの展開," NTT技術ジャーナル, no.12, pp.4-8, Dec. 1989.
- (7) 関根俊彦, 佳山東一, 宮本雅昭, 渡部剛士, "NSP/NSSPのソフトウェア機能拡充," NTT技術ジャーナル, no.12, pp.70-72, Dec. 1994.
- (8) 鈴木滋彦, "新ノードシステムの開発," NTT R&D, vol.45, no.6, pp.497-506, June 1996.
- (9) 上坂久一, 網谷駿介, "INサービスの先駆け フリーダイヤルサービス," NTT技術ジャーナル, no.12, pp.21-26, Dec. 1989.
- (10) 上坂久一, 吉江金三郎, 児玉博義, 吉見正信, "高度電話網サービス制御ソフトウェア構成法," NTT電気通信研究所研究実用化報告, vol.36, no.8, pp.1059-1064, Aug. 1987.
- (11) 村瀬節雄, 池尻稔, 武井伊佐夫, "サービスのカスタマ化とカスタマコントロール," NTT技術ジャーナル, no.12, pp.16-20, Dec. 1989.
- (12) 伊藤弘, "高度INのアーキテクチャ技術," NTT技術ジャーナル, no.6, pp.13

-17, June 1993.

- (13) 鈴木滋彦, “インテリジェントネットワーク（IN）の高度化への展望,” NTT技術ジャーナル, no.6, pp.8 -12, June 1993..
- (14) 鈴木孝至, 上坂久一, 木ノ内康夫, “高度INのサービス制御技術,” NTT技術ジャーナル, no.6, pp.18 -22, June 1993.
- (15) 弓場英明, 佐藤清実, 今川仁, 鈴木孝至, 田中豪, “新ノード高度IN技術,” NTTR&D, vol.45, no.6, pp.559-568, June 1996.
- (16) 佐藤清実, “高度INのサービス管理技術,” NTT技術ジャーナル, no.6, pp.23 -37, June 1993.
- (17) 花澤隆, “高度INのオペレーション技術,” NTT技術ジャーナル, no.6, pp.28 -32, June 1993.
- (18) M. T. Chao, and P. B. Passero, “Using General Purpose Computing Equipment as a Base for an Advanced Intelligent Network Systems Platform,” IEEE IN'95 Workshop, May 1995.
- (19) M. Syrett, D. Skov, and A. Kristensen, “HP in IN,” IEEE IN'95 Workshop, May 1995.
- (20) J. Gray, and A. Reuter, “TRANSACTION PROCESSING : CONCEPTS AND TECHNIQUES,” MORGAN KAUFMANN PUBLISHERS, CALIFORNIA, 1993.
- (21) 渡辺栄一, J.グレイ, “オンライントランザクション処理:OLTPシステム,” 近代科学社, 1995.
- (22) 鶴保征城, 木ノ内康夫, 星子隆幸, 仲谷元, 宮川順治, “ループを用いた大規模分散処理システム,” 情報処理学会誌, vol.31, no.5, pp.686-697, May 1990.
- (23) 竹内進, 宮崎達三, 貝沢哲男, 倉野明彦, “サービスを実現する高機能レイヤの信頼性を向上,” NTT技術ジャーナル, no.9, pp.17-20, Sep. 1991.
- (24) 浜田浩平, 黒住弘明, “リアルタイムネットワークシステム,” 信学誌, vol.73, no.11, pp.1161-1166, Nov. 1990.
- (25) T. L. Casavant, and J. G. Kuhl, “A Taxonomy of Scheduling in General-Purpose Distributed Computing Systems,” IEEE Trans. on Software and Eng., vol.14, no.2,

pp.141-154, Feb. 1988.

- (26) 渡辺尚, 太田剛, 水野忠則, 中西暉, “双方向ピギーバックに基づいた動的負荷分散方式,”信学論(D-I), vol.J78-D-I, no.3, pp.302-312, March 1995.
- (27) 田中克之, “並列汎用機の全貌,” 日系B P社, Sep. 1994.
- (28) M.H. Eich, “Foreword Main Memory Databases : Current and Future Research Issues,” IEEE Trans. on Knowl. and Data Eng., vol.4, no.6, pp.507-508, Dec. 1992.
- (29) H. Garcia-Molina, and K. Salem, “Main Memory Database Systems : An Overview,” IEEE Trans. on Knowl. and Data Eng., vol.4, no.6, pp.509-516, Dec. 1992.
- (30) 泉谷建司, “E t h e r n e t と F D D I,” pp.338-342, ソフト・リサーチ・センター, 1996.
- (31) T. J. Lehman, and M. J. Carey, “A Recovery Algorithm for A High-Performance Memory-Resident Database System,” Proc. ACM SIGMOD Int. Conf. on Management of Data, pp.104-117, May 1987.
- (32) M.H. Eich, “Main Memory Database Recovery,” ACM FJCC, pp.1226-1232, Nov. 1986.
- (33) V. Kumar, and A. Burger, “Performance Measurement of Some Main Memory Database Recovery Algorithms,” Proc. of 7th Int. Conf. on Data Engineering, pp.436-443, April 1991.
- (34) 高倉弘喜, 上林彌彦, “フラッシュメモリバックアップ方式の設計と性能評価,” 信学論(D-I), vol.J76-D-I, no.10, pp.514-521, Oct.1993.
- (35) (財)金融情報システムセンター, “金融情報システム白書（平成9年版）,” 財経詳報社, 1996.
- (36) 渡辺均, 能條哲, “信頼性設計のためのシステム開発の動向,” NTT技術ジャーナル, no.3, pp.14 -18, March 1993.
- (37) H. Watanabe, “A reliability design method for maintenance strategy of telecommunication network,” Reliability and Maintainability Symposium, pp.476-483, Jan. 1993.
- (38) M. Hirano, Y. Kinouchi, and T. Suzuki, “Distributed Control Node Architecture in the Advanced Intelligent Network,” ISS '95, April 1995.

- (39) 平野正則, 塩澤恒道, 木ノ内康夫, 中村篤, 井上潮, “分散データベースシステムのコストパフォーマンス評価,” 電子情報通信学会春季大会, D104, 1992.
- (40) 大光明直孝, 林誠治, 富上幸成, 平野正則, “分散構成におけるデータベースのバックアップ方法,” 電子情報通信学会総合大会, D-6-22, 1997.
- (41) ITU-T 勧告, “Q.1200 シリーズの能力セット 1”, 1993.
- (42) 西原琢夫, 富田清次, “蓄積情報を有するリアルタイムシステム構成法,” 信学論(D-I), vol.J78-D-I, no.8, pp.770-776, Aug. 1995.
- (43) J. Gray, and D. P. Siewiorek, “High-Availability Computer Systems,” COMPUTER, pp.39-48, Sep. 1991.
- (44) Technical Requirements TR100001, “Interface for Realtime Operating Systems,” Ver. 1, NTT, 1992.
- (45) 大久保利一, 香西省治, 二神新, 大南正人, “リアルタイム性を考慮した通信処理用OSインターフェースの設計と性能評価,” 信学論(D-I), Vol.J78-D-I, No.8, pp.687-698, Aug. 1995.
- (46) 平野正則, 塩澤恒道, 木ノ内康夫, 鈴木孝至, “高度INの分散処理におけるデータ配置による負荷の偏り,” 信学論(B-I), vol.J-80-B-I, no.2, pp.87-97, Feb. 1997.
- (47) M. Hirano, T. Shiozawa, Y. Kinouchi, and T. Suzuki, “Large-Scale Distributed Control Node in the Advanced Intelligent Network,” APCC '95, June 1995.
- (48) 富上幸成, 平野正則, 大光明直孝, 吉見正信, 木ノ内康夫, “高度INにおける分散データベース処理の負荷平準化に関する一考察,” 電子情報通信学会通信ソサイエティ大会, B-623, 1996.
- (49) 水谷静夫, “朝倉日本語新講座2語彙,” pp.102-105, 朝倉書店, 1983.
- (50) 小杉肇, “統計調査論,” pp.128-129, 朝倉書店, 1976.
- (51) 平野正則, 櫻井秀紀, 今川仁, 木ノ内康夫, “分散処理による高度IN用サービス制御ノードのモジュール間結合方式,” 信学論(B-I), vol.J81-B-I, no.8, pp.519-530, Aug. 1998.

- (52) H. Sakurai, M. Ito, M. Hirano, and H. Imagawa, "Inter-module Connections in Distributed-processing Service Control Points in the Advanced Intelligent Network and Testing of Associated Hardware," ICIN '96, Nov. 1996.
- (53) 櫻井秀紀, 塩澤恒道, 平野正則, 今川仁, "A T M結合機構を用いた分散処理システムの試験方式," 電子情報通信学会通信ソサイエティ大会, B-406, 1995.
- (54) 斎藤秀一, 大光明直孝, 櫻井秀紀, 伊藤守夫, 平野正則, "分散処理による高度 IN用サービス制御ノードのモジュール間結合方式に関する一考察," 電子情報通信学会総合大会, B-695, 1996.
- (55) 脇村慶明, "デジタル交換システム," 信学誌, vol.73, no.11, pp.1167-1173, Nov. 1990.
- (56) S. Heatley, and D. Stokesberry, "Analysis of Transport Measurements Over a Local Area Network," IEEE Commun. Mag., vol.27, pp.16-22, June 1989.
- (57) S. Nojo, and H. Watanabe, "Incorporating reliability specifications in the design of telecommunication networks," IEEE Commun. Mag., vol.31, pp.40-43, June 1993.
- (58) 平野正則, 山根道広, 山崎幹夫, 木ノ内康夫, 林 誠治, "高度 IN用サービス制御ノードにおけるメモリデータベースのリカバリ方式," 信学論(B-I), vol.J80-B-I, no.8, pp.596-608, Aug. 1997.
- (59) 林誠治, 平野正則, 石橋宏純, 小林伸幸, 白石正裕, "高度 IN用データベースのリカバリ方式に関する一考察," 電子情報通信学会総合大会, B-658, 1996.
- (60) 斎藤秀一, 伊藤守夫, 平野正則, 櫻井秀紀, 林誠治, "高度 IN用半導体ファイル装置のバッテリバックアップ方法," 電子情報通信学会総合大会, B-6-118, 1997.
- (61) 川瀬克之, 秋葉昭浩, 大光明直孝, 平野正則, "高度 IN用大容量半導体ファイル装置のメモリ試験に関する一検討," 電子情報通信学会総合大会, B-6-4, 1998.
- (62) 山根道広, 秋葉昭浩, 平野正則, "半導体ファイル用バッテリバックアップ電源の高信頼化構成," 電子情報通信学会総合大会, B-6-66, 1998.

- (63) K. Takeno, M. Yamasaki, and S. Muroyama, "A New Backup Power Supply with a Battery Deterioration Test Circuit," INTELEC'95, pp.591-596, Oct. 1995.
- (64) 平野正則, 山根道広, 小林正光, 木ノ内康夫, "相互にバックアップされたデータベースを持つ二重化ノードにおける駆付け保守の評価," 信学論(B-I), vol.J81-B-I, no.10, pp.***-***, Oct. 1998.
- (65) 山根道広, 平野正則, 木ノ内康夫, "センタ間バックアップを行う高信頼化システムの平均修理時間," 電子情報通信学会春季大会, A-389, 1994.
- (66) T. Hanazawa, "NTT's Intelligent Network Today and Tomorrow," ICIN '98, May 1998.

付録

第4章

付録4.1：トラヒックの期待値の算出

$$y_0 = \int_0^1 y dz = \int_0^1 c \left\{ (g^{\frac{1}{\alpha}} - 1)z + 1 \right\}^{-\alpha} dz \quad (\text{A}\cdot 4.1.1)$$

ここで、

$$z = \frac{x-1}{h-1} \quad (\text{A}\cdot 4.1.2)$$

$$h = g^{\frac{1}{\alpha}} \quad (\text{A}\cdot 4.1.3)$$

とおき、式 (A·4.1.2) 、 (A·4.1.3) を式 (A·4.1.1) に代入し次式を得る。

$$y_0 = \frac{c}{h-1} \int_1^h x^{-\alpha} dx \quad (\text{A}\cdot 4.1.4)$$

(1) $\alpha \neq 1$ の場合

式 (A·4.1.4) から次式を得る。

$$y_0 = \frac{c}{h-1} \left[\frac{x^{1-\alpha}}{1-\alpha} \right]_1^h = \frac{c(h^{1-\alpha} - 1)}{(1-\alpha)(h-1)} \quad (\text{A}\cdot 4.1.5)$$

式 (A·4.1.5) に式 (A·4.1.3) を代入し、次式を得る。

$$y_0 = \frac{c(g^{\frac{1}{\alpha}-1} - 1)}{(1-\alpha)(g^{\frac{1}{\alpha}} - 1)} \quad (\text{A}\cdot 4.1.6)$$

(2) $\alpha = 1$ の場合

式 (A·4.1.4) から次式を得る。

$$y_0 = \frac{c}{h-1} [L_{og} x]^h = \frac{c L_{og} h}{(h-1)} \quad (\text{A}\cdot\text{4.1.7})$$

式 (A·4.1.7) に式 (A·4.1.3) を代入し、次式を得る。

$$y_0 = \frac{c L_{og} g}{(g-1)} \quad (\text{A}\cdot\text{4.1.8})$$

付録 4.2 : トラヒックの累積値の期待値に対する割合の算出

$$\eta = \frac{\int_0^z y dz}{y_0} \quad (\text{A}\cdot\text{4.2.1})$$

ここで、

$$\int_0^z y dz = \int_0^z c \left\{ (g^{\frac{1}{\alpha}} - 1)z + 1 \right\}^{-\alpha} dz \quad (\text{A}\cdot\text{4.2.2})$$

に、式 (A·4.1.2) 、 (A·4.1.3) を代入し、次式を得る。

$$\begin{aligned} \int_0^z y dz &= \frac{1}{h-1} \int_1^{(h-1)z+1} cx^{-\alpha} dx \\ &= \frac{c((hz-z+1)^{1-\alpha} - 1)}{(1-\alpha)(h-1)} \end{aligned} \quad (\text{A}\cdot\text{4.2.3})$$

式 (A·4.2.3) に式 (A·4.1.3) を代入し次式を得る。

$$\int_0^z y dz = \frac{c((g^{\frac{1}{\alpha}} z - z + 1)^{1-\alpha} - 1)}{(1-\alpha)(g^{\frac{1}{\alpha}} - 1)} \quad (\text{A}\cdot\text{4.2.4})$$

式 (A·4.2.4) 、式 (A·4.1.6) を式 (A·4.2.1) に代入し、次式を得る。

$$\eta = \frac{(z - zg^{\frac{1}{\alpha}} + g^{\frac{1}{\alpha}})^{1-\alpha} - g^{1-\frac{1}{\alpha}}}{(1 - g^{\frac{1}{\alpha}})} \quad (\text{A}\cdot\text{4.2.5})$$

付録 4.3 : α を無限大としたときのトラヒック量の算出

$$y = c \left\{ \left(g^{\frac{1}{\alpha}} - 1 \right) z + 1 \right\}^{-\alpha} \quad (\text{A}\cdot 4.3.1)$$

y は z の関数であり、 $z = 0$ のまわりでテーラ展開すると、次式で表せる。

$$y(z) = y(0) + y'(0) \frac{z}{1!} + y''(0) \frac{z^2}{2!} + \cdots + y^{(n)}(0) \frac{z^n}{n!} + \cdots \quad (\text{A}\cdot 4.3.2)$$

また、式 (A·4.3.1) より次式を得る。

$$y^{(n)}(z) = c(-1)^n \alpha(\alpha+1)\cdots(\alpha+n-1) \left(g^{\frac{1}{\alpha}} - 1 \right)^n \left\{ \left(g^{\frac{1}{\alpha}} - 1 \right) z + 1 \right\}^{-(\alpha+n)} \quad (\text{A}\cdot 4.3.3)$$

式 (A·4.3.1) 、式 (A·4.3.3) で $z = 0$ とおき、

$$y(0) = c \quad (\text{A}\cdot 4.3.4)$$

$$y^{(n)}(0) = c(-1)^n \alpha(\alpha+1)\cdots(\alpha+n-1) \left(g^{\frac{1}{\alpha}} - 1 \right)^n \quad (\text{A}\cdot 4.3.5)$$

式 (A·4.3.4) 、式 (A·4.3.5) を式 (A·4.3.2) に代入して、 $\alpha \rightarrow \infty$ とおくと次式を得る。

$$\lim_{\alpha \rightarrow \infty} y(z) = c \left(1 + \sum_{n=0}^{\infty} \left\{ \lim_{\alpha \rightarrow \infty} (-1)^n \alpha(\alpha+1)\cdots(\alpha+n-1) \left(g^{\frac{1}{\alpha}} - 1 \right)^n \right\} \frac{z^n}{n!} \right) \quad (\text{A}\cdot 4.3.6)$$

ここで、ロピタルの規則を適用し、次式が成り立つ。

$$\lim_{\alpha \rightarrow \infty} \alpha^n \left(g^{\frac{1}{\alpha}} - 1 \right)^n = \left\{ \lim_{\alpha \rightarrow \infty} \alpha \left(g^{\frac{1}{\alpha}} - 1 \right) \right\}^n = (L_{og}g)^n \quad (\text{A}\cdot 4.3.7)$$

また、 $i < n$ の場合、次式が成り立つ。

$$\lim_{\alpha \rightarrow \infty} \alpha^i \left(g^{\frac{1}{\alpha}} - 1 \right)^n = \left\{ \lim_{\alpha \rightarrow \infty} \alpha \left(g^{\frac{1}{\alpha}} - 1 \right) \right\}^i \lim_{\alpha \rightarrow \infty} (g^{\frac{1}{\alpha}} - 1)^{n-i} = 0 \quad (\text{A}\cdot 4.3.8)$$

式 (A·4.3.7) 、式 (A·4.3.8) を式 (A·4.3.6) に代入し、次式を得る。

$$\lim_{\alpha \rightarrow \infty} y(z) = c \left(1 + \sum_{n=0}^{\infty} (-L_{og}g)^n \frac{z^n}{n!} \right) = ce^{-(L_{og}g)z} = cg^{-z} \quad (\text{A}\cdot 4.3.9)$$

付録 4.4：負荷の偏りの標準偏差の算出

$$\sigma = \left\{ \frac{1}{(h-1)^R \lambda_0^2} \int_1^h \cdots \int_1^h (cx_1^{-\alpha} + \cdots + cx_R^{-\alpha} - \lambda_0)^2 dx_1 \cdots dx_R \right\}^{1/2} \quad (\text{A}\cdot 4.4.1)$$

ここで、

$$\begin{aligned} H &= (cx_1^{-\alpha} + \cdots + cx_R^{-\alpha} - \lambda_0)^2 \\ &= c^2(x_1^{-2\alpha} + x_2^{-2\alpha} + \cdots + x_R^{-2\alpha}) + 2c^2(x_1^{-\alpha}x_2^{-\alpha} + \cdots + x_{R-1}^{-\alpha}x_R^{-\alpha}) \\ &\quad - 2c\lambda_0(x_1^{-\alpha} + x_2^{-\alpha} + \cdots + x_R^{-\alpha}) + \lambda_0^2 \end{aligned} \quad (\text{A}\cdot 4.4.2)$$

とおき、各項毎に積分を行う。

【 $x_i^{-2\alpha}$ の項の積分】

$$\int_1^h \cdots \int_1^h x_i^{-2\alpha} \cdot dx_R \cdots dx_R = \frac{(h^{1-2\alpha} - 1)(h-1)^{R-1}}{(1-2\alpha)} \quad (\text{A}\cdot 4.4.3)$$

【 $x_i^{-\alpha}x_j^{-\alpha}$ の項の積分】

$$\int_1^h \cdots \int_1^h x_i^{-\alpha}x_j^{-\alpha} dx_1 \cdots dx_R = \frac{(h^{1-\alpha} - 1)^2(h-1)^{R-2}}{(1-\alpha)^2} \quad (\text{A}\cdot 4.4.4)$$

【 $x_i^{-\alpha}$ の項の積分】

$$\int_1^h \cdots \int_1^h x_i^{-\alpha} dx_1 \cdots dx_R = \frac{(h^{1-\alpha} - 1)(h-1)^{R-1}}{(1-\alpha)} \quad (\text{A}\cdot 4.4.5)$$

【定数項の積分】

$$\int_1^h \cdots \int_1^h dx_1 \cdots dx_R = (h-1)^R \quad (\text{A}\cdot 4.4.6)$$

式 (A・4.4.3) ~ 式 (A・4.4.6) 式を式 (A・4.4.1) に代入し、次式を得る。

$$\begin{aligned}\sigma &= \frac{1}{(h-1)^{R/2} \lambda_0} \left\{ \frac{c^2 R (h^{1-2\alpha} - 1) (h-1)^{R-1}}{(1-2\alpha)} \right. \\ &\quad + \frac{c^2 R (R-1) (h^{1-\alpha} - 1)^2 (h-1)^{R-2}}{(1-\alpha)^2} \frac{2c R \lambda_0 (h^{1-\alpha} - 1) (h-1)^{R-1}}{(1-\alpha)} \\ &\quad \left. + \lambda_0^2 (h-1)^R \right\}^{1/2}\end{aligned}\quad (\text{A} \cdot 4.4.7)$$

$\lambda_0 = y_0 R$ 及び本文中の式 (4.5) を代入し、次式を得る。

$$\begin{aligned}\sigma &= \frac{1}{\sqrt{R}} \left\{ \frac{(1-\alpha)^2 (h-1) (h^{1-2\alpha} - 1)}{(1-2\alpha) (h^{1-\alpha} - 1)^2} - 1 \right\}^{1/2} \\ &= \frac{1}{\sqrt{R}} \left\{ \frac{(1-\alpha)^2 (g^{\frac{1}{\alpha}} - 1) (g^{\frac{1}{\alpha}-2} - 1)}{(1-2\alpha) (g^{\frac{1}{\alpha}-1} - 1)^2} - 1 \right\}^{1/2}\end{aligned}\quad (\text{A} \cdot 4.4.8)$$

付録 4.5 : $\alpha = 1$ で A が極大値を取ることの証明

$$A(\alpha, g) = \sqrt{\left\{ \frac{(1-\alpha)^2 (g^{\frac{1}{\alpha}} - 1) (g^{\frac{1}{\alpha}-2} - 1)}{(1-2\alpha) (g^{\frac{1}{\alpha}-1} - 1)^2} - 1 \right\}} \quad (\text{A} \cdot 4.5.1)$$

$A(\alpha, g)$ を α で微分すると次式を得る。

$$\begin{aligned}\frac{dA(\alpha, g)}{d\alpha} &= \frac{1}{2A} \left\{ \frac{(2\alpha^2 (g^2 - g^{\frac{1}{\alpha}}) (1 - g^{\frac{1}{\alpha}}) + (1-2\alpha) (g^2 + 1 - 2g^{\frac{1}{\alpha}}) g^{\frac{1}{\alpha}} \log g) (\alpha-1)^2}{\alpha^2 (2\alpha-1)^2 (g - g^{\frac{1}{\alpha}})^2} + \right. \\ &\quad \left. \frac{2(g^2 - g^{\frac{1}{\alpha}}) (g^{\frac{1}{\alpha}} - 1) (\alpha-1) (g\alpha^2 - g^{\frac{1}{\alpha}}\alpha^2 + g^{\frac{1}{\alpha}} \log g - g^{\frac{1}{\alpha}}\alpha \log g)}{\alpha^2 (2\alpha-1) (g - g^{\frac{1}{\alpha}})^3} \right\}\end{aligned}\quad (\text{A} \cdot 4.5.2)$$

ここで、 $\alpha \rightarrow 1$ とすると次式を得る。

$$\lim_{\alpha \rightarrow 1} \frac{dA(\alpha, g)}{d\alpha} = \frac{1}{2A} \left\{ -g(g-1)^2(2 + \log g) \lim_{\alpha \rightarrow 1} \frac{(\alpha-1)^2}{(g-g^{\alpha})^2} + \right.$$

$$\left. 2g(g-1)^2 \lim_{\alpha \rightarrow 1} \frac{(\alpha-1)(g\alpha^2 - g^{\alpha}\alpha^2 + g^{\alpha}\log g - g^{\alpha}\alpha\log g)}{(g-g^{\alpha})^3} \right\} \quad (\text{A} \cdot 4.5.3)$$

ここで、ロピタルの規則を適用し次式が成り立つ。

$$V = \lim_{\alpha \rightarrow 1} \frac{(\alpha-1)^2}{(g-g^{\alpha})^2} = \left\{ \lim_{\alpha \rightarrow 1} \frac{\frac{d}{d\alpha}(\alpha-1)}{\frac{d}{d\alpha}(g-g^{\alpha})} \right\}^2 = \left\{ \lim_{\alpha \rightarrow 1} \frac{1}{g L_{og} g} \right\}^2 = \frac{1}{g^2 (L_{og} g)^2} \quad (\text{A} \cdot 4.5.4)$$

$$W = \lim_{\alpha \rightarrow 1} \frac{(\alpha-1)(g\alpha^2 - g^{\alpha}\alpha^2 + g^{\alpha}L_{og}g - g^{\alpha}\alpha L_{og}g)}{(g-g^{\alpha})^3}$$

$$= \frac{\lim_{\alpha \rightarrow 1} d^3(\alpha-1)(g\alpha^2 - g^{\alpha}\alpha^2 + g^{\alpha}L_{og}g - g^{\alpha}\alpha L_{og}g) / d\alpha^3}{\lim_{\alpha \rightarrow 1} d^3(g-g^{\alpha})^3 / d\alpha^3}$$

$$= \frac{3gL_{og}g(2+L_{og}g)}{6g^3(L_{og}g)^2} = \frac{2+L_{og}g}{2g^2(L_{og}g)^2} \quad (\text{A} \cdot 4.5.5)$$

なお、Wの分子、分母の3回微分の値は以下の通りである。

$$W \text{の分子} = \lim_{\alpha \rightarrow 1} \left\{ \frac{6g^{\alpha}(L_{og}g)^3}{\alpha^3} - \frac{18(g-g^{\alpha})g^{\alpha}(L_{og}g)^2(2\alpha+L_{og}g)}{\alpha^6} \right.$$

$$\left. + \frac{3(g-g^{\alpha})^2g^{\alpha}L_{og}g(6\alpha^2+6\alpha L_{og}g+(L_{og}g)^2)}{\alpha^6} \right\}$$

$$= 6g^3(L_{og}g)^3 \quad (\text{A} \cdot 4.5.6)$$

$$\begin{aligned}
W \text{の分子} &= \lim_{\alpha \rightarrow 1} \left\{ 6(g - g^{\frac{1}{\alpha}}) + \frac{3g^{\frac{1}{\alpha}} L_{og} g}{\alpha^4} (2\alpha^3 + \alpha(2-\alpha)L_{og}g + (1-\alpha)(L_{og}g)^2) \right. \\
&\quad \left. + \frac{(\alpha-1)g^{\frac{1}{\alpha}}(L_{og}g)^2}{\alpha^6} (4\alpha^2 L_{og}g - 6\alpha^2 - 6\alpha L_{og}g + (\alpha-1)(L_{og}g)^2) \right\} \\
&= 3gL_{og}g(2 + L_{og}g) \tag{A・4.5.7}
\end{aligned}$$

式 (A・4.5.4)、式 (A・4.5.5) を式 (A・4.5.2) に代入し、次式を得る。

$$\lim_{\alpha \rightarrow 1} \frac{dA(\alpha, g)}{d\alpha} = 0 \tag{A・4.5.8}$$

第5章

付録5.1 : X(h) の算出

周期 i で発生するメッセージ数を n^i 、その発生確率を $P^i(n^i)$ 、周期 i での処理残りメッセージ数を h^i 、その発生確率を $X^i(h^i)$ 、周期内で処理可能な最大メッセージ数を b とおく。周期 $i+1$ での処理残りメッセージ数 h^{i+1} は次式で表せる。

$$h^{i+1} = \begin{cases} n^i + h^i - b & (n^i + h^i \geq b) \\ 0 & (n^i + h^i < b) \end{cases} \tag{A・5.1.1}$$

周期 i で h^i 個のメッセージが残る確率 $X^i(h^i)$ 、と n^i 個のメッセージが発生する確率 $P^i(n^i)$ は独立であり、周期 $i+1$ で h^{i+1} 個のメッセージが残る確率 $X^{i+1}(h^{i+1})$ は、上記関係式が成立する h^i と n^i のすべての組合せの確率 $X^i(h^i) \cdot P^i(n^i)$ の総和となり、次式で表せる。

$$\begin{aligned} X^{i+1}(1) &= X^i(0) \cdot P(b+1) + X^i(1) \cdot P(b) + \dots \\ &\quad \dots + X^i(b) \cdot P(1) + X^i(b+1) \cdot P(0) \end{aligned} \tag{A・5.1.2}$$

$$\begin{aligned} X^{i+1}(2) &= X^i(0) \cdot P(b+2) + X^i(1) \cdot P(b+1) + \dots \\ &\quad \dots + X^i(b+1) \cdot P(1) + X^i(b+2) \cdot P(0) \end{aligned} \tag{A・5.1.3}$$

• • • • •

また、 $X^{i+1}(h^{i+1})$ の総和は 1 となることから、

$$X^{i+1}(0) = 1 - X^{i+1}(1) - X^{i+1}(2) - X^{i+1}(3) - \dots \tag{A・5.1.4}$$

となる。ここで初期値は、

$$X^0(0) = 1 \tag{A・5.1.5}$$

$$X^0(1) + X^0(2) + X^0(3) + \dots = 0 \tag{A・5.1.6}$$

である。式 (A・5.1.5) 、 (A・5.1.6) を式 (A・5.1.2) ~ (A・5.1.4) の漸化式に代入し、 i を大きくするに従って $X^{i+1}(h^{i+1})$ の平衡状態の値 $X(h)$ が算出できる。

第 6 章

付録 6.1：アレーディスク装置の適用性

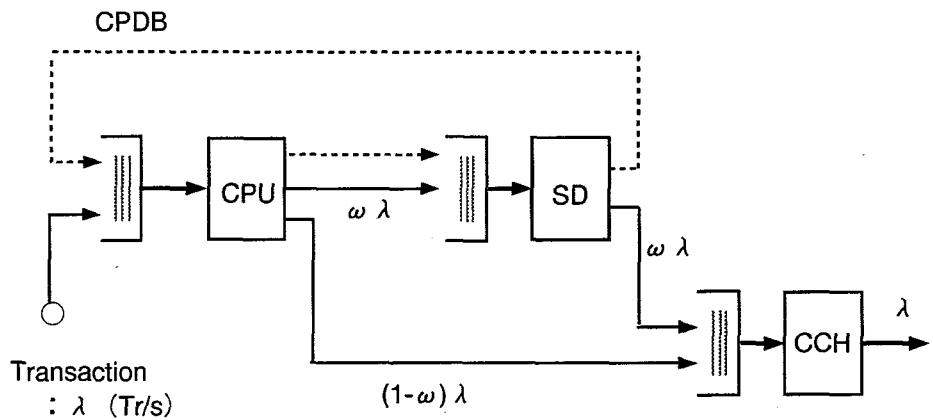
磁気ディスク装置を数台～10数台使用し、並列に動作可能としたアレーディスク装置がある。最近のアレーディスク装置の単位時間当たりの入出力回数は数100（回/秒）程度であり、LOG取得用としての適用の可能性がある。しかし、数GBの磁気ディスク装置を使用したアレーディスク装置の記憶容量は数10GB程度となり、数100MB～1GB程度のデータベースのCPDBやLOGを取得する記憶装置として使用する場合は、容量面での使用効率が極端に低下し、コスト性能比が悪い。

付録 6.2：方式 1 のレスポンスタイムの評価

C P D B は b_{CP} 単位に S D に書込む。書き込みを確実に行うため書き込み終了の確認をとる。すべての C P D B を T_{CP} 内に書き込むためには、 $T_{CP} \cdot b_{CP} / M$ の周期で書きめば良いが、ここではトランザクション処理への影響が最も厳しいケースとして、 b_{CP} の書き込み終了と同時に次の書き込みの要求を C P U キューに登録することとした。シミュレーション評価モデルを図 A・6.2.1 に示す。トランザクションはポアソン到着とし、平均発生頻度を λ とした。 b_{CP} とレスポンスタイムの関係を図 A・6.2.2 に示す。図 A・6.2.2 より b_{CP} が小さい場合には S D での待ちがほとんど問題とならないが、 b_{CP} が 100 KB より大きくなるに従って S D での待ちが急激に増大しレスポンスタイムが悪化する。

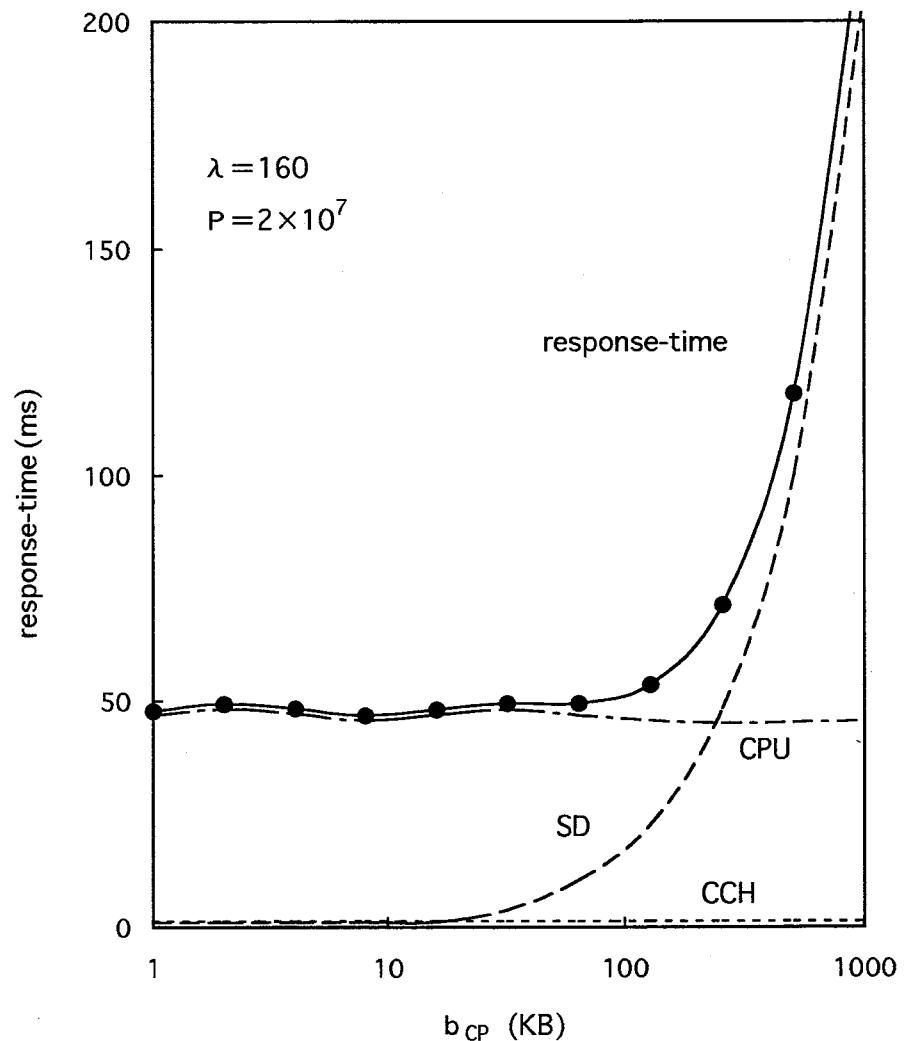
付録 6.3：キャッシュメモリ上の最新情報の書き出し

最近の高性能マイクロプロセッサではキャッシュメモリが使用される。しかも、キャッシュメモリの制御方式として、ほとんどの場合コピーバック方式がとられる。コピーバック方式では、最新のデータはキャッシュメモリ上に存在し、データをキャッシュメモリから追い出すときに初めて主メモリに反映される。従って、電源障害が発生した場合、障害時点にすべてのデータベースが主メモリ上に反映されている保証はない。このため、キャッシュメモリの内容を主メモリに書き出す必要がある。マイクロプロセッサの書き換えを即座に主メモリに反映するためには、キャッシュ制御方式としてストアスルー方式を採用する必要があるが、ストアスルー方式ではプロセッサと主メモリ間のボトルネックが発生し、マイクロプロセッサの性能を十分に引き出せない。



CPU: central processing unit
 SD: semiconductor disk device
 CCH: communication control channel
 CPDB: check-point database

図A・6.2.1 レスポンスタイムの評価モデル



図A・6.2.2 b_{CP} とレスポンスタイムの関係

第7章

付録7.1：経過時間と不稼働率の関係の算出

本文中の状態遷移の方程式を基に、経過時間と不稼働率の関係を差分方程式を用いて算出した。概要を以下に示す。

ある時間 t から微小時間 δ 経過した場合、システムの各状態の確率は次式で近似できる。

$$P_{ij,kl}(t+\delta) \approx P_{ij,kl}(t) + \delta P'_{ij,kl}(t) \quad (\text{A}\cdot7.1.1)$$

式 (A\cdot7.1.1) を本文中の式 (7.1) ~ (7.16) に代入し、さらに、駆付けを全く行わないことから、本文中の式 (7.24) を適用し次式を得る。

$$\mathbf{P}(t+\delta) = \mathbf{A} \bullet \mathbf{P}(t) \quad (\text{A}\cdot7.1.2)$$

ここで、 \mathbf{A} 、 \mathbf{P} は以下の行列で表せる。

$$\mathbf{A} = \begin{bmatrix} 1 - 4\lambda\delta & 0 & \cdots & \cdots & 0 & 0 \\ 4\lambda\delta & 1 - 3\lambda\delta & \cdots & \cdots & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & \cdots & 1 - 2\mu\delta & 0 \\ 0 & 0 & \cdots & \cdots & 0 & 1 - 4\mu\delta \end{bmatrix} \quad (\text{A}\cdot7.1.3)$$

$$\mathbf{P}(t) = \begin{bmatrix} P_{00,00}(t) \\ P_{00,10}(t) \\ \cdots \\ P_{20,21}(t) \\ P_{21,21}(t) \end{bmatrix} \quad (\text{A}\cdot7.1.4)$$

時間 t から n δ 時間経過後は以下の通りとなる。

$$\mathbf{P}(t + n\delta) = \mathbf{A}^n \bullet \mathbf{P}(t) \quad (\text{A}\cdot7.1.5)$$

ここで、 $t = 0$ とおくと次式となり、 $t = 0$ から $n \delta$ 時間経過後の各状態の確率が求まる。

$$\mathbf{P}(n\delta) = \mathbf{A}^n \bullet \mathbf{P}(0) \quad (\text{A}\cdot7.1.6)$$

なお、 $t = 0$ のときの $\mathbf{P}(0)$ は本文中の式 (7.25) 、(7.26) で与えられる。

各状態の確率から、 $n \delta$ 時間経過後の不稼働率は次式で表せる。

$$P_{SDN}(n\delta) = P_{20,20}(n\delta) + P_{20,21}(n\delta) + P_{21,21}(n\delta) \quad (\text{A}\cdot7.1.7)$$

δ を小さく設定することにより、精度が高くなる。ここでは、 δ を 1/1000 時間刻みで算出した。これ以上、 δ を小さくしても結果はほとんど変化なかった。なお、式 (A·7.1.5) の数値計算には、Mathematica を使用した。

略号一覧

A M C : ATM connector : ATMコネクタ

サービス制御ノード（SCP）を複数のモジュールで構成する場合、モジュール間をATMにより接続する結合機構。

A P L : application program : アプリケーションプログラム

サービスを実現するための応用プログラム。

A T M : asynchronous transfer mode : 非同期転送モード

広帯域の伝送／交換技術。本論文ではモジュールと結合機構との間のデータ転送に使用。

B T : battery : バッテリ

電源断においても、半導体ディスク装置上の情報を不揮発化するために使用。

C C H : communication channel : 通信チャネル

モジュール間の通信を行うためのチャネル。

C C M : communication control module : 通信制御モジュール

サービス制御ノード（SCP）を機能分散／負荷分散で構成する場合のモジュールの一種であり、共通線信号網を介して交換機（SSP）との通信機能を実現するモジュール。

C P : check point : チェックポイント

ソフトウェアのバグやハードウェア障害により、主メモリ上のデータベースが破壊された場合、障害直前のデータベースを復元するため、ある時間間隔で主メモリ上のデータベースを半導体ディスク装置にコピーするが、このコピーを取得する契機を与える時間のこと。

C P D B : check point database : チェックポイントデータベース

ソフトウェアのバグやハードウェア障害により、主メモリ上のデータベースが破壊された場合、障害直前のデータベースを復元するため、チェックポイント毎に取得したデータベースのこと。

C P U : central processing unit : 中央演算処理装置

プロセッサを構成する一要素であり、主メモリ上の命令を実行する装置。

D B : database : データベース

カスタマ毎のデータから構成されるデータの集合。

D B C : database copy : データベース コピー

ソフトウェアのバグやハードウェア障害により、主メモリ上のデータベースが破壊された場合、それを復元するため、他のモジュールの主メモリに格納しておくデータベースのコピー。

D R A M : dynamic random access memory : ダイナミック ランダムアクセスメモリ

主メモリや半導体ディスク装置で用いるメモリ素子。

D S : dynamic step : ダイナミックステップ

ある機能を実現するために処理しなければならない命令数のこと。要求される機能を実現するためのプロセッサの負荷を表す尺度として用いる。

F D D I : fiber distributed data interface : エフディィディアイ

100Mビット／秒のリング型LANの名称。

H D : hard disk : 磁気ディスク装置

O S や A P L のプログラムファイル等を格納する、磁気を用いた記憶装置。

H D C : hard disk controller : 磁気ディスク制御装置

磁気ディスク装置へのアクセスを制御する装置。

I N : intelligent network : インテリジェントネットワーク

ネットワーク内にカスタマ毎のデータを保持し、この内容に従って高度な電話サービスを実現する網。

I R O S : interface for realtime operating system : アイロス

NTTで開発している交換機やサービス制御ノード等で使用しているリアルタイムOS。

I T U - T : international telecommunication union-telecommunication sector : 国際電気通信連合標準化部門

電気通信に係わる標準化を行っている国際機関。

L A N : local area network ; ローカルエリアネットワーク

同一建物内、あるいは同一敷地内などの比較的狭い地域に分散設置されたモジュールなどを結ぶ構内ネットワークシステム。

L O G : log : ログ

トランザクションによりデータベースを更新した場合の更新内容。データベースが破壊された場合、それを復元するために用いる。

MM : main memory : 主メモリ

プロセッサを構成する一要素であり、命令やデータを格納する記憶装置。

M U X : multiplexing and demultiplexing : 多重／分離化部

AMCを構成する機能ブロックの一種であり、SDHIからの信号の多重化や、SDHIへの信号の分離化を行う機能ブロック。

N D M : network directory module : 網ディレクトリモジュール

サービス制御ノード(SCP)を機能分散／負荷分散で構成する場合のモジュールの一種であり、カスタマ毎のサービスデータが配備されているSCPを特定するための網ディレクトリを格納し、それに対するアクセスを行うモジュール。

O L T P : on-line transaction processing : オンライントランザクション処理

銀行のオンラインシステムやクレジット照会システムなど、データベースを有するホストシステムと、通信回線で接続された多数の端末とから構成されたシステムで、端末からの問合わせに対してホストシステムでデータベースの更新や参照を行い、結果を端末に返す処理のこと。

O M M : operation and maintenance module : オペレーション管理モジュール

サービス制御ノード（S C P）を機能分散／負荷分散で構成する場合のモジュールの一種であり、サービス管理ノードやオペレーションノードとの通信機能を実現するモジュール。

O P S : operation system : オペレーションノード

インテリジェントネットワークを構成する通信ノードの一種であり、多数のサービス制御ノード（S C P）の障害処理や、ソフトウェアの変更に伴うファイル更新処理等を行うノード。

P H S : personal handy-phone system : 簡易型携帯電話

携帯電話の一種。本論文で対象とするサービス制御ノードにはP H S の位置情報がデータベースとして格納される。

P roc : processor : プロセッサ

C P Uや主メモリ、チャネル等から構成される処理装置のこと。

S C F : service control function : サービス制御機能

I T U-Tで規定されているサービス制御ノードを実現するための機能の一つであり、呼の接続制御等のサービス制御を実現する機能。

S C M : service control module : サービス制御モジュール

サービス制御ノード（S C P）を機能分散／負荷分散で構成する場合のモジュールの一種であり、サービス制御機能（S C F）を実現するモジュール。

S C P : service control point : サービス制御ノード

インテリジェントネットワークを構成する通信ノードの一種であり、カスタマ毎のデータベースを保持し、交換機（S S P）と協調して、高機能な電話呼の接続処理を行うノード。

S D : semiconductor disk device : 半導体ディスク装置

半導体メモリを用いて構成し、チャネルを経由してアクセスできる記憶装置。機械的な動作を伴わないため、磁気ディスク装置に比べてアクセス時間が短い。

S D C : semiconductor disk device controller : 半導体ディスク制御装置

半導体ディスク装置へのアクセスを制御する装置。

S D F : service data function : サービスデータ機能

I T U – T で規定されているサービス制御ノードを実現するための機能の一つであり、呼の接続のために必要なカスタマ毎のサービスデータを管理し、それに対するアクセスを実現する機能

S D H I : synchronous digital hierarchy interface : 同期デジタルハイアーキインタフェース部

AMC を構成する機能ブロックの一種で、モジュールと S D H で接続する機能ブロック。

S D M : service data module : サービスデータモジュール

サービス制御ノード (S C P) を機能分散／負荷分散で構成する場合のモジュールの一種であり、サービスデータ機能 (S D F) を実現するモジュール。

S I M : system interface module : システムインターフェースモジュール

サービス制御ノード (S C P) を機能分散／負荷分散で構成する場合のモジュールの一種であり、モジュール間を接続する AMC と OMM を一体化したモジュール。

S L P : service logic program : サービス論理プログラム

カスタマが受けるサービス（たとえば、フリーダイヤルの場合、時間帯により接続先の電話機を変更するとか、発信者の位置により接続する電話機を変更する等のサービス）を実現するため、カスタマ毎に作成されたプログラム。S L P を書き換えることによりサービスのカスタマイズ化が容易に実現できる。

S M S : service management system : サービス管理ノード

インテリジェントネットワークを構成する通信ノードの一種であり、カスタマの氏名や住所等のカスタマ情報や、カスタマが受けるサービス条件等のデータベースを管理するノード。実際の呼処理に必要な情報は S M S からサービス制御ノード (S C P) にダウンロードされる。

S R S W : self-routing switch : 自己ルーティングスイッチ

AMC を構成する機能ブロックの一種であり、A T M 信号のルーティングを行う機能ブロック。

S S P : service switching point : 交換機

電話呼の接続処理を行うノード。高機能な電話呼（I N呼）の場合は、S S PがI N呼であることを検出すると共通線信号網を介してサービス制御ノード（S C P）に問い合わせをかけ、S S PとS C Pが協調してI N呼の接続を行う。

T C A P : transaction capabilities application part : トランザクション機能応用部

共通線信号網で、データベースへのアクセスなどを効率的に行うためのプロトコル。

T M N : telecommunication management network : 通信管理ネットワーク体系

オペレーションノードと、交換機、伝送装置等のネットワークエレメントとを容易に接続するための通信アーキテクチャ。

V C C I : voluntary control council for interference by information technology equipment : 情報処理装置等電波障害自主規制協議会

電子装置からの妨害電波の自主規制を行っている協議会。

V P N : virtual private network : 仮想私設網

インテリジェントネットワークを用いた電話サービスの一種で、加入電話網をあたかも社内の内線電話のように利用できるようにしたサービス。

発表論文一覧 (下線は関連発表論文を示す)

A. 論文誌

- (1) 福村好美, 平野正則, 塩澤恒道, “バス結合型マルチプロセッサのキャッシュメモリ構成方式,” 信学論(D-I), vol.J74-D-I, no.10, pp.721-728, Oct. 1991.
- (2) 木ノ内康夫, 桜井紀彦, 平野正則, 塩澤恒道, “長短データアクセス混在時のストライプアレーディスクのスループット解析,” 信学論(D-I), vol.J76-D-I, no.8, pp.417-428, Aug. 1992.
- (3) 平野正則, 鈴木孝至, 塩澤恒道, 芳西 崇, 木ノ内康夫, “分散処理による高度IN用サービス制御ノードの構成,” 信学論(B-I), vol.J79-B-I, no.8, pp.539-550, Aug. 1996.
- (4) 平野正則, 塩澤恒道, 木ノ内康夫, 鈴木孝至, “高度INの分散処理におけるデータ配置による負荷の偏り,” 信学論(B-I), vol.J80-B-I, no.2, pp.87-97, Feb. 1997.
- (5) 平野正則, 山根道広, 山崎幹夫, 木ノ内康夫, 林 誠治, “高度IN用サービス制御ノードにおけるメモリデータベースのリカバリ方式,” 信学論(B-I), vol.J80-B-I, no.8, pp.596-608, Aug. 1997.
- (6) 平野正則, 櫻井秀紀, 今川仁, 木ノ内康夫, “分散処理による高度IN用サービス制御ノードのモジュール間結合方式,” 信学論(B-I), vol.J81-B-I, no.8, pp.519-530, Aug. 1998.
- (7) 平野正則, 山根道広, 小林正光, 木ノ内康夫, “相互にバックアップされたデータベースを持つ二重化ノードにおける駆付け保守の評価,” 信学論(B-I), vol.J81-B-I, no.10, pp.***-***, Oct. 1998.

B. 國際会議

- (1) M. Hirano, Y. Kinouchi, and T. Suzuki, “Distributed Control Node Architecture in the Advanced Intelligent Network,” ISS '95, April 1995.

(2) M. Hirano, T. Shiozawa, Y. Kinouchi, and T. Suzuki, "Large-Scale Distributed Control Node in the Advanced Intelligent Network," APCC '95, June 1995.

(3) H. Sakurai, M. Ito, M. Hirano, and H. Imagawa, "Inter-module Connections in Distributed-processing Service Control Points in the Advanced Intelligent Network and Testing of Associated Hardware," ICIN '96, Nov. 1996.

C. NTT社内論文誌

- (1) 多嶋清次郎, 松本博幸, 平野正則, 小濱晴雄, "DIPS VLSI プロセッサ本体系の構成," 研究報, vol.33, no.1, pp.17-30, 1984.
- (2) S. Tajima, H. Matsumoto, and M. Hirano, "Central Processing Subsystem for DIPS VLSI Processor," Rev. of the ECL, vol.32, no.4, pp.700-708, 1984.
- (3) 平野正則, 山口利和, 東海林敏夫, 北村士守, "DIPS Vシリーズ用汎用インターフェース制御装置の構成," 研究報, vol.36, no.6, pp.761-770, 1987.
- (4) 矢沢良一, 平野正則, 山口利和, 岡田靖史, "DIPS-V30E のハードウェア構成," 研究報, vol.37, no.9, pp.523-532, 1988.
- (5) 矢沢良一, 平野正則, 小濱晴雄, 丸山正人, "DIPS-V30EX のハードウェア構成," NTT R&D, vol.38, no.8, pp.885-894, 1989.

D. その他の機関誌

- (1) 平野正則, 矢沢良一, 園田雅文, "INSネットへの直接接続を実現した分散プロセッサ," NTT技術ジャーナル, no.8, pp.55-57, 1989.
- (2) 脇村慶明, 岡田勝行, 小町谷忠芳, 平野正則, "高性能システムバスの構成技術を開発," NTT技術ジャーナル, no.6, pp.81-84, 1990.

E. 全国大会

- (1) 桑原敏, 平野正則, "マイクロプログラム制御のチャネルにおける多重処理の影響について," 昭和52年度電子通信学会総合全国大会, 1381, 1977.

- (2) 小原和博, 平野正則, 多嶋清次郎, “マルチプロセッサシステムにおけるバス制御方式の検討,” 昭和 55 年度情報処理学会第 21 回全国大会, 4G-6, 1980.
- (3) 平野正則, “アドレス変換制御部の LSI 化に関する一考察,” 情報処理学会第 22 回 (昭和 56 年前期) 全国大会, 7J-5, 1981.
- (4) 平野正則, 小原永, “VLSI 化論理設計法に関する一考察,” 電子通信学会情報・システム部門全国大会, 551, 1981.
- (5) 平野正則, “密結合マルチプロセッサのメモリアクセス方式に関する一考察,” 情報処理学会第 34 回 (昭和 62 年前期) 全国大会, 5Q-5, 1987.
- (6) 森啓, 平野正則, “スプリット制御バスに適したキャッシュ一致制御方式,” 電子情報通信学会春季全国大会, D-120, 1990.
- (7) 塩澤恒道, 平野正則, “密結合マルチプロセッサにおけるキャッシュ制御方式,” 情報処理学会第 40 回 (平成 2 年前期) 全国大会, 4L-1, 1990.
- (8) 平野正則, “密結合マルチプロセッサにおける二階層キャッシュ制御方式,” 電子情報通信学会秋季全国大会, D-80, 1990.
- (9) 平野正則, 塩澤恒道, 福村好美, “バス結合型マルチプロセッサ構成に関する一考察,” 電子情報通信学会春季全国大会, D-141, 1991.
- (10) 平野正則, 塩澤恒道, 木ノ内康夫, 中村篤, 井上潮, “分散データベースシステムのコストパフォーマンス評価,” 電子情報通信学会春季大会, D104, 1992.
- (11) 塩澤恒道, 平野正則, “排他制御方式に関する一考察,” 電子情報通信学会春季大会, D-123, 1992.
- (12) 塩澤恒道, 平野正則, “コピーバックキャッシュメモリを用いたメモリリカバリ方式に関する一考察,” 電子情報通信学会秋季大会, D-56, 1992.
- (13) 塩澤恒道, 平野正則, “メモリリカバリ機能を有するキャッシュメモリ適用法の検討,” 電子情報通信学会秋季大会, D-74, 1993.

- (14) 木ノ内康夫, 平野正則, 塩澤恒道, “連続運転の容易化に向けた時空間多相プロセッサの適用について,” 電子情報通信学会春季大会, SD-6-2, 1994.
- (15) 山根道広, 平野正則, 木ノ内康夫, “センタ間バックアップを行う高信頼化システムの平均修理時間,” 電子情報通信学会春季大会, A-389, 1994.
- (16) 櫻井秀紀, 塩澤恒道, 平野正則, 今川仁, “A T M結合機構を用いた分散処理システムの試験方式,” 電子情報通信学会ソサイエティ大会, B-406, 1995.
- (17) 林誠治, 平野正則, 石橋宏純, 小林伸幸, 白石正裕, “高度 I N用データベースのリカバリ方式に関する一考察,” 電子情報通信学会総合大会, B-658, 1996.
- (18) 斎藤秀一, 大光明直孝, 櫻井秀紀, 伊藤守夫, 平野正則, “分散処理による高度 I N用サービス制御ノードのモジュール間結合方式に関する一考察,” 電子情報通信学会総合大会, B-695, 1996.
- (19) 富上幸成, 平野正則, 大光明直孝, 吉見正信, 木ノ内康夫, “高度 I Nにおける分散データベース処理の負荷平準化に関する一考察,” 電子情報通信学会ソサイエティ大会, B-623, 1996.
- (20) 斎藤秀一, 伊藤守夫, 平野正則, 櫻井秀紀, 林誠治, “高度 I N用半導体ファイル装置のバッテリバックアップ方法,” 電子情報通信学会総合大会, B-6-118, 1997.
- (21) 大光明直孝, 林誠治, 富上幸成, 平野正則, “分散構成におけるデータベースのバックアップ方法,” 電子情報通信学会総合大会, D-6-22, 1997.
- (22) 川瀬克之, 秋葉昭浩, 大光明直孝, 平野正則, “高度 I N用大容量半導体ファイル装置のメモリ試験に関する一検討,” 電子情報通信学会総合大会, B-6-4, 1998.
- (23) 山根道広, 秋葉昭浩, 平野正則, “半導体ファイル用バッテリバックアップ電源の高信頼化構成,” 電子情報通信学会総合大会, B-6-66, 1998.