

データベースサーバシステムの 制御方式に関する一考察

根岸 和義

(株)日立製作所 システム開発研究所

データベースサーバシステムにおいてそのプロトコルの正当性の検証は重要な課題である。本報告では、有限オートマトン(FSA)により記述された複数のプロトコルマシンで定義されているプロトコルの検証方式に関して述べる。各FSAの組合せのFSAを分析して、不当状態(コミット、ロールバックの混在した状態)や行き止まり状態(所期の最終状態に遷移できない状態)が発生しないことを検証し、何台のサーバマシンを用意すれば実機における動作確認が可能であることを示した。FSAの組合せに際して状態の直積でなく、異なる状態の組合せのみを考えることにより、サーバマシンの台数とは独立にすべての台数のサーバに対して上記の検証を可能とした。

A Study of Control Methods for
Database Server Systems

Kazuyoshi Negishi

Systems Development Lab., Hitach Ltd.
1099 Ohzenji Asao-ku Kawasaki-si 215, JAPAN

Verification of communication protocol is important for database server systems. This paper describes a verification method of the protocol. That protocol is defined by number of finite state automata(FSA). We examine these FSA to exclude invalid states which include both commit state and rollback state, and to exclude halt states which has no transition path to final states. We show how many server machines are needed for verification of the protocol with real machines. Reduced FSA is used to analyze the protocol independently of the number of database servers.

1. まえがき

近年、大型システムの発達とともに、計算機システムの処理能力に対するニーズも増大している。これに対処する手段の一つとして、複数の計算機システムを通信回線により連結し、互いのデータベースを相互にアクセス可能としたデータベースサーバシステムに対するニーズも高まっている。また、異機種種の計算機システムによりデータベースサーバシステムを構成する際に、その相互接続を保証するプロトコルの標準化がOSIのTP⁽²⁾等として推進されている。

従来、このプロトコルの定式化には、主に有限状態オートマトン(FSA)によるプロトコルマシンとその直積によるFSAが用いられてきた⁽¹⁾。そこでは、 n を固定値とした、1対 n のモデルが使用されている。

従来の方式の問題点は、1対 n のモデルに対するプロトコルの検証能力が不十分であった点である。 n を固定した場合、それをFSAとして、不正状態(各ノードのコミット、ロールバックが混在した終了状態)を生じないこと、および行き止まり状態(終了状態への遷移ができない状態)を生じないことをモデルで検証し、実機で確認することはできるが、すべての n について、FSAによる検証、および実機による確認を行うことは困難である。また、プロトコルの検証において n 台の個々の状態をすべて検証する必要はなく、状態の組合せ方の各々に着目して検証すれば十分な場合が多い。

そこで我々は、FSAによりモデル化し、プロトコルマシンの状態の直積でなく、状態の種類を考え、状態の組合せによるFSAを解析する方策を考案した。この方式によれば、ノード数 n とは独立にプロトコルの検証を行うことができる。

我々の対象とするシステムは、クライアント、

サーバモデルにより表現されるデータベースサーバシステムである。クライアントのマシンは、全体を制御するメインマシンとサーバ対応の個別マシンから成る。各ノードは、独立にジャーナルを取得し、ノード障害に備えるシステムである。このシステムにおいて、各プロトコルマシンをFSAの状態遷移により、モデル化した。プロトコルマシンのうちサーバマシン対応に設定される部分をまずサーバ直積マシンとして統合し、このマシンとシステム全体を統括するプロトコルマシンとのメッセージのやりとりを組合せマシンとしてモデル化し、不当状態への遷移を組合せマシンによりチェックする方策を取った。

本報告のシステムの特徴は、1対 n のプロトコルマシンの状態組合せマシンによるシステムのモデル化と、その検証方策の策定である。なお、具体的システムにおける1対 n のプロトコルマシンの定義方法、および障害回復に着目したその検証法則に関してはすでに報告済みである⁽⁴⁾。本報告では、任意台数のサーバにおけるプロトコルの検証の統一の方策を示す。

2. モデルの概要

我々は評価方式を明確にするため、サーバ、クライアント間のトランザクション処理に着目し、下記のように簡略化したモデルを検討した。本報告で想定するプロトコルマシンの構成を図1に示す。システムの詳細はすでに報告済み⁽⁴⁾であるので、ここでは概要のみを示す。システムはクライアントノード及び n 台のサーバノードより構成される。クライアントノードには、全体を統括するメインマシンと、サーバノードに対するインタフェースとなる個別マシンがある。一方サーバノードにはサーバマシンがある。各プロトコルマシンはメッセージにより互いに

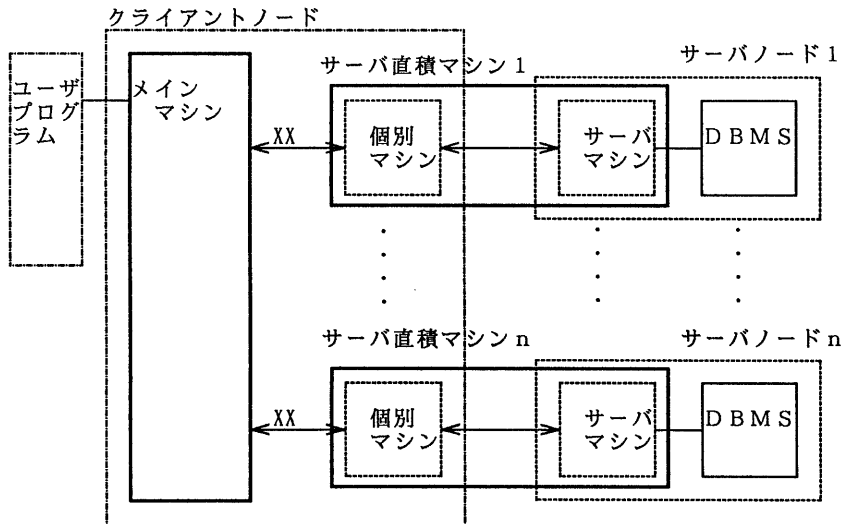


図1 プロトコルマシンの関連図

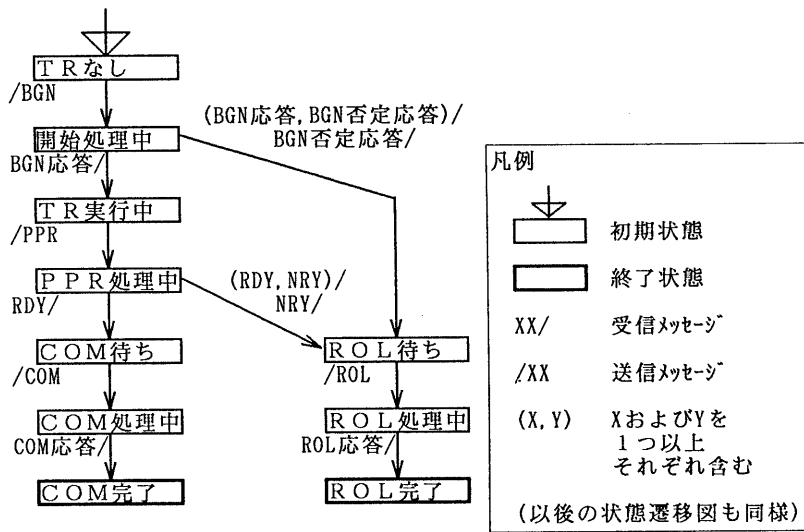


図2 メインマシンの状態遷移図

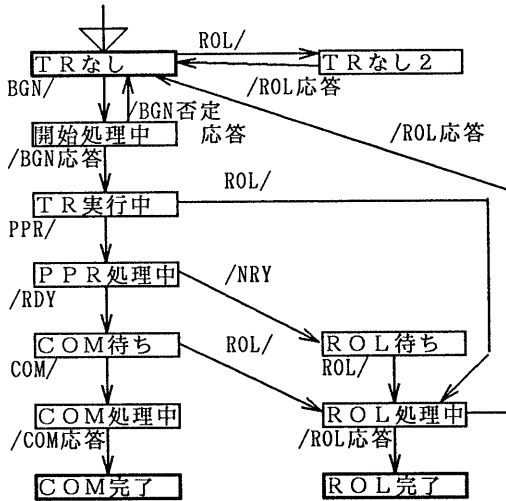


図3 サーバ直積マシンの状態遷移図

連絡する。

2. 1 メッセージの伝達モデル

メインマシンと個別マシン、個別マシンとサーバマシン間のメッセージは各プロトコルマシン間で同時には一つのみ存在する（半二重通信方式を想定）。また、メッセージの消失、追越は想定しない。ノード障害の発生時、受信されていないメッセージは削除される。

メインマシンと個別マシンの間では、メインから個別には常に同一の内容を全個別マシンへ

送信する。これは、トランザクションの処理開始時に関連する全ノードへトランザクション開始を連絡する方式を想定している。また、トランザクションに関連しないノードに関してはモデル化の対象外とする。また、個別マシンからメインマシンへのメッセージは、下記のいずれかとする。

(1) 全個別マシンから同一のメッセージを受信する。

(2) 全個別マシンから二種類のメッセージのいずれかを受信する。ただし、各メッセージを

少なくとも一つの個別マシンから受信するものとする。

2.2 プロトコルマシンの詳細

メインマシンの状態遷移図を図2に示す。個別マシンとサーバマシンの直積を考え、メインマシンとのやりとりに関連のない部分を統合したものが、サーバ直積マシンである。図3にサーバ直積マシンの状態遷移図を示す。この図において、開始状態は「TRなし」状態であり、終了状態は「COM完了」状態、「ROL完了」状態、「TRなし」状態のいずれかである。

3. プロトコルの正当性検証

次に、メインマシンとサーバ直積マシンとの組合せマシンを考える。図4の例でこれを説明する。まず、初期状態は、メインマシン、サーバ直積マシンとも「TRなし」状態であるから、組合せマシンの初期状態は（「TRなし」、（「TRなし」））である。ここで、（ S_0, S_1, \dots, S_m ）は、 S_0 がメインマシンの状態、 S_1, \dots, S_m が m 通りのサーバ直積マシンの状態の組合せを表す。メインマシンからサーバ直積マシンへBGNメッセージを同時に送信すると組合せマシンの状態は、（「開始処理中」、「開始処理中」）となる。次に、サーバマシンがダウンすることなく開始処理を終了すると

全てのサーバ直積マシンからBGN応答が返り、組合せマシンは状態（「TR実行中」、 $\left(\begin{matrix} \text{「TR実行中」} \\ \text{「TRなし」} \end{matrix} \right)$ ）に遷移する。しかし、いくつかのサーバノードでダウンが発生した場合、当該サーバ直積マシンはBGN否定応答を返すから、組合せマシンの状態は、（「ROL待ち」、 $\left(\begin{matrix} \text{「TR実行中」} \\ \text{「TRなし」} \end{matrix} \right)$ ）に遷移する。また、すべてのノードでダウンが発生した場合、すべてのノードがBGN否定応答を返すので、組合せマシンの状態は、（「ROL待ち」、 $\left(\begin{matrix} \text{「ROL待ち」} \\ \text{「ROL完了」} \end{matrix} \right)$ ）に遷移する。以下同様にして、図6が得られる。ここで、組合せマシンの終了状態は、 $\left(\begin{matrix} \text{「COM完了」} \\ \text{「COM完了」} \end{matrix} \right)$ 、 $\left(\begin{matrix} \text{「ROL完了」} \\ \text{「ROL完了」} \end{matrix} \right)$ 、 $\left(\begin{matrix} \text{「ROL完了」} \\ \text{「ROL完了」} \end{matrix} \right)$ 、 $\left(\begin{matrix} \text{「ROL完了」} \\ \text{「TRなし」} \end{matrix} \right)$ 、 $\left(\begin{matrix} \text{「ROL完了」} \\ \text{「TRなし」} \end{matrix} \right)$ のいずれかである。従って、「TRなし」と「ROL完了」を同一視すれば、この組合せマシンにおいて、プロトコルの規約に矛盾する不正状態への遷移は発生しない。アルゴリズム1がこの組合せマシンを生成する手順を示す。また、定理1よりこのマシンが、不正状態（正常でない状態）を生じないなら、対応する直積マシンにおいても任意の n に対して不正状態への遷移は発生しない。

次に、組合せマシンの各状態に対して直積マ

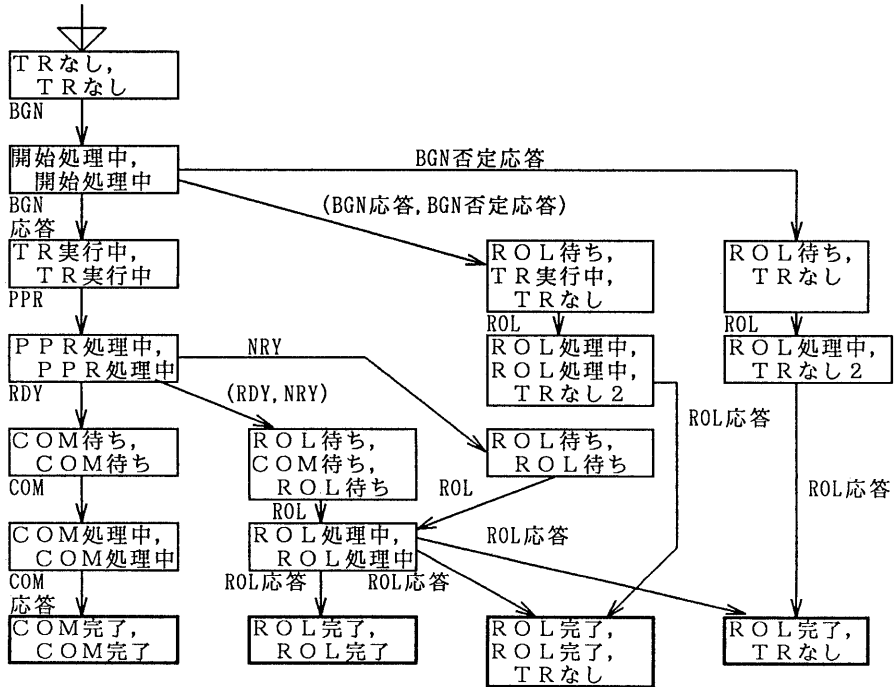


図4 組合せマシンの状態遷移図

[注] (n_1, \dots, n_m) :
 対応する組合せマシンの状態 (q^0, q_1, \dots, q_m)
 に対し、 $(q_1 \times n_1, \dots, q_m \times n_m)$ が DN_i の唯一の
 要素であることを示す。
 $[n_1, \dots, n_m]$:
 同上、 DM_i の唯一の要素であることを示す。

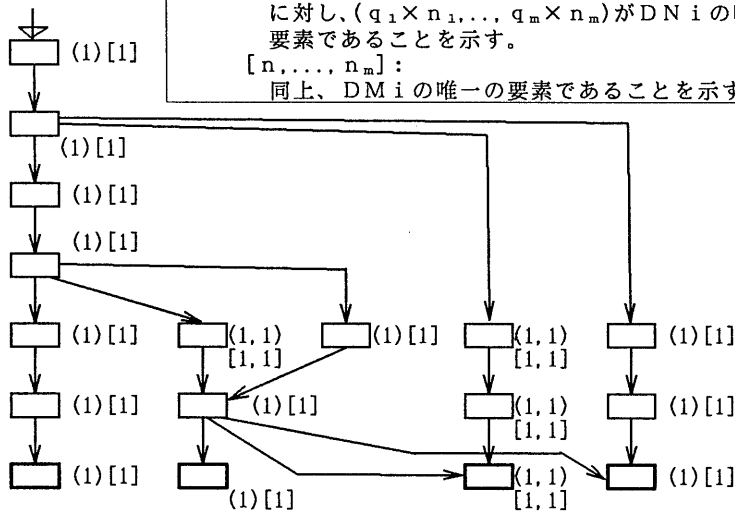


図5 アルゴリズム2の評価結果

シン上で当該状態に対応する状態に遷移するために最小限何台のサーバマシンが必要かをアルゴリズム2を使用して分析する。アルゴリズム2に従い、各ノードの

DN_i : ノード i に対応する直積マシンの状態で初期状態から遷移可能な状態の集合

DM_i : ノード i に対応する直積マシンの状態で終了状態へ遷移可能な状態の集合

を求める。

for all $q \in DN_i$

$\exists q' \in DM_i (q \supset q')$

ならノード i は行き止まり状態ではない。ノード i における DN_i の要素のうち最小ノード数を d_i とした時、 $\text{Max}(d_i)$ が実機でのテストに必要なノード数である。

本例題に上記のアルゴリズム1を適用した結果を分析の結果を図5に示す。図5ではステップ2の DN_i を $()$ により表し、 DM_i を $[\]$ により表す(記述法は図5の注記の通り)。この例では、行き止まり状態がないことがわかる。図5の DN_i の状態数の最大値が2であることから本プロトコルマシンでは最小限2台のサーバマシンがあれば全ての組合せを検証できることがわかる。

ここで、 DM_i を定義する過程で、ノードの障害発生時にしか遷移しないものが含まれている可能性がある。これを避けるにはアルゴリズム1において、直積マシンの遷移として正常時の遷移のみを対象とするようアルゴリズムを変更すればよい。ここでは、簡単のためその可能性は考慮しない。

4. DN_i , DM_i の考察

DN_i の持つ意味を以下に考察する。組合せマシンのノード i について、これに対応する直積マシンの $n = 1 \sim \infty$ のすべてのノードのうち、初期状態から移行可能なノードに対応する状態の組合せの集合を考える。この集合の要素のうち包含関係 \supset に関して、他の要素を含まない要素(最小要素)の集合が DN_i である。

DM_i は逆に、終了状態に移行可能なノードに対して上記を行ったものである。

しかし、このような考え方で DN_i , DM_i を求めるのでは、間接的に直積マシンを扱っており、組合せマシンを考えた意味が少ない。そこで、付録に示すアルゴリズムでは、直積マシンを考慮せず、組合せマシンの状態を扱う方を示した。

5. まとめ

本報告では、データベースサーバシステムのプロトコルマシンモデルに対して、プロトコルの検証をノード別の状態ではなく、状態の種別に対して行うことにより、ノード数と独立にプロトコルが行き止まり状態を生じないことの検証を行うアルゴリズムを示した。さらに、実機においてプロトコルの検証に必要なノード数を求めるアルゴリズムを示した。

本方式により、データベースサーバシステムのように1対 n の関係にあるプロトコルシステムの n と独立したプロトコル検証が可能となった。また、プロトコル実証時の必要なノード数を求めることができた。

参考文献

- [1] Dale Skeen, "Formal Model of Crash Recovery in a Distributed System", IEEE Trans. on Software Eng. Vol. SE-9, No. 3 May 1983 pp219-228.
- [2] ISO/IEC DIS 10026 Information Processing Systems-Open Systems Interconnection-Distributed Transaction Processing
- [3] J.E. ホップ クラフト, J.D. ウルマン, "言語理論とオートマトン", サイエンス社(1971)
- [4] 根岸他, "分散データベースシステムの制御方式に関する一考察", 信学技法 Vol. 90, No. 146, 1990年7月, pp29-38.

付録 記号の定義とアルゴリズム、定理

文献 [3] の記法に基づき、n 台のサーバ直積マシンを下記で表す。

$$M = (X, Q, \delta, q_0, F)$$

ここで、

- X : マシンの入力記号の集合
- Q : マシンの状態の集合
- δ : マシンの状態遷移関数の集合
- q : マシンの初期状態
- F : マシンの終了状態の集合

とする。

次に、X, Q, δ にアルファベット順を仮定し、A を X, Q, δ のいずれかに固定した時、

$$(A)^n : A \text{ の要素の長さ } n \text{ の列}$$

とする。

長さ n の列 $a_1, \dots, a_n \in (A)^n$ に対して、

$$U(a_1, \dots, a_n) = (a_1, \dots, a_n \text{ のうち重複するものを除きアルファベット順に並べ換えたもの})$$

なる関数を考える。次に、

$$(A)^+ : A \text{ の要素の列の集合}$$

とした時、U を拡張して、

$$U((A)^+) = \{ U(a_1, \dots, a_n) \mid a_1, \dots, a_n \in (A)^+ \}$$

とする。これを用いて、M と同様にメインマシンを下記で表す。

$$M^0 = (X^0, Q^0, \delta^0, q_0^0, F^0)$$

ここで、

$$X^0 = U((X)^+)$$

である。

M^0 と M はその入出力メッセージをすべて入力記号と見なすことにより、FSA を定義している。以下の直積マシン、および組合せマシンでは、 M^0 と n 個の M が各々 $U(x^1, \dots, x^n)$, および x^1, \dots, x^n を入力した時、メッセージの送受信が成立したと見なす。

サーバ直積マシンとメインマシンのノード数 n における直積マシンを下記のように定義する。 $M^{pn} = (X^{pn}, Q^{pn}, \delta^{pn}, q_0^{pn}, F^{pn})$

ここで、

$$X^{pn} = (X)^n$$

$$Q^{pn} = Q^0 \times (Q)^n$$

$$\delta^{pn} : Q^{pn} \times X^{pn} \rightarrow Q^{pn}$$

$$q_0^{pn} = q_0^0 q_0^1 \dots q_0^n$$

$$F^{pn} = F^0 \times (F)^n$$

である。ただし、

$$\delta^{pn}(q_i^1 q_i^2 \dots q_i^n, x^1, \dots, x^n) = (q_j^1 q_j^2 \dots q_j^n)$$

となるのは、

$$\delta^0(q_i^0, U(x^1, \dots, x^n)) = q_j^0 \text{ かつ、}$$

$$\delta(q_i^k, x^k) = q_j^k \text{ (} k = 1 \sim n \text{)}$$

の時である。

また、サーバ直積マシンとメインマシンの組合せマシンを下記のように定義する。

$$M^s = (X^s, Q^s, \delta^s, q_0^s, F^s)$$

ここで、

$$X^s = U((X)^+)$$

$$Q^s = Q^0 \times U((Q)^+)$$

$$\delta^s : Q^s \times X^s \rightarrow Q^s$$

$$q_0^s = (q_0^0 q_0^1)$$

$$F^s = F^0 \times U((F)^+)$$

である。ただし、

$$\delta^s(U(q_i^0 q_i^1 \dots q_i^n), U(x^1, \dots, x^n))$$

$$= U(q_j^0 q_j^1 \dots q_j^n) \text{ となるのは、あるノード数 } n \text{ において}$$

$$\delta^{pn}(q_i^0 q_i^1 \dots q_i^n, x^1, \dots, x^n)$$

$$= q_j^0 q_j^1 \dots q_j^n$$

が存在する時に限る。

組合せマシンの状態遷移を直積マシンを考えずに求めるアルゴリズムを以下に示す。

[アルゴリズム 1]

組合せマシンのすべてのノード i の状態 :

$$(q_i^0 q_{i1}^1 \dots q_{im}^m) \in Q^s$$

に対して、下記により状態遷移を定義する。

(1) q_i^0 を左辺とするクライアントマシンの状態遷移の集合を $G(q_i^0)$ とする。

(2) q_{ik}^k を左辺とするサーバ直積マシンの状態遷移の集合を $G(q_{ik}^k)$ とする。 $G(q_{ik}^k)$ のすべての部分集合 (空集合を除く) の集合を $G(q_{ik}^k)^+$ とする。

(3) クライアントマシンとサーバ直積マシンの間で記号 $X = x_1, \dots, x_r$ を入出力し、ノード j の状態 $(q_j^0 q_{j1}^1 \dots q_{jm}^m)$ に遷移する状態遷移の集合 G_{ijx} を考える。

$$G_{ijx} = \{ \{ \delta^0 \}, G^1, \dots, G^m \} \mid$$

$$\delta^0 \in G(q_i^0),$$

$$G^k \in G(q_{ik}^k)^+,$$

$$\delta^0(q_i^0, x) = q_j^0,$$

$$U(G^1, \dots, G^m \text{ のすべての記号}) = x,$$

$$U(G^1, \dots, G^m \text{ のすべての右辺})$$

$$= q_{j1}^1 \dots q_{jm}^m$$

(4) このような G_{ijx} の存在する時、

$$\delta^s(q_i^0 q_{i1}^1 \dots q_{im}^m, x)$$

$$= q_j^0 q_{j1}^1 \dots q_{jm}^m$$

とする。

ここで、本報告では、初期状態から遷移可能な直積マシンの状態においては、送信可能な記号列 x に対してはこれを受信する状態遷移が必ず存在することを仮定する。このような状態遷移が存在するかどうかをチェックするように本方式のアルゴリズムを拡張することは可能であるが、それは本報告の対象外とする。

次に、不正状態と行き止まり状態を下記のように定義する。

不正状態：個別のサーバ直積マシンの終了状態の集合を複数のグループに分けた時、直積マシンの終了状態のうちこれら複数のグループの二つ以上のグループに属する終了状態を含むものを不正状態とする。

行き止まり状態：FSAにおいて、初期状態から遷移可能な状態で、その状態から終了状態への遷移ができない状態を行き止まり状態とする。この時、以下の定理が成立する。

[定理1]

組合せマシンにおいて、不正状態（正常でない状態）を生じないなら、対応するすべての直積マシンでもこのような状態は発生しない。

[証明]

直積マシンにおいて、初期状態から状態 $q_i^1 q_i^2 \dots q_i^n$ への状態遷移が存在する時、組合せマシンにおいて初期状態から対応する状態 $U(q_i^1 q_i^2 \dots q_i^n)$ への状態遷移が存在する。定義より、直積マシンの状態及び状態遷移に対応して、組合せマシンの状態および状態遷移が必ず存在する。従って、直積マシンにおいて、初期状態から不正状態への遷移が存在するなら、それに対応して直積マシンにおいて、初期状態から対応する不正状態への遷移が存在する。

[証明終]

直積マシンMの状態 $q^{pn} \in Q^{pn}$ と、 $q \in Q$ に対して、

$N(q^{pn}, q) = q^{pn}$ 中の q の個数

を定義する。また、 $q^0 \in Q^0$ に対して、

$N(q^{pn}, q^0) = 1$

とする。次に、 q^{pn} と q^{pm} の包含関係 \supseteq を以下のように定義する。

$q^{pn} \supseteq q^{pm}$ とは、

for all $q \in q^{pn}$, $N(q^{pn}, q) \neq 0$

$N(q^{pn}, q) \geq N(q^{pm}, q)$

である。また、

$q^{pn} \supseteq q^{pm}$ かつ $q^{pn} \neq q^{pm}$ の時、

$q^{pn} \supset q^{pm}$ であるとする。

サーバ直積マシン任意の n 個とメインマシンを組合せた場合、組合せマシンに行き止まり状態が存在しない場合、直積マシンにも行き止まり状態が存在しないかどうかをチェックするために、 DN_i 、および DM_i を求める以下のアルゴリズムを考案した。

[アルゴリズム2]

[ステップ1]

組合せマシンのノード毎に探索済み状態集合 W_i (初期値 ϕ)、および探索対象状態集合 V_i を用意する。

[ステップ2]

組合せマシンの初期状態 0 で、 $V_0 = \{(q^0, q_0)\}$ とする。

[ステップ3]

ノード 0 を対象に以下の処理 [ステップ4] を行う。

[ステップ4]

もし、for all $q \in V_i$ ($\exists q' \in W_i$ ($q \supseteq q'$)) なら何もしないで終了。

それ以外なら、 $U = \{q \mid q \in V_i, \neg \exists q' \in W_i (q \supseteq q')\}$ を求める。

$V_i \cup W_i$ の包含関係に付いての最少の部分集合を新たな W_i とする。

ノード i から入力記号 X によるノード j への遷移各々に付き [ステップ4-1] ~ [ステップ4-3] を行う。

[ステップ4-1]

状態遷移 $i \rightarrow j$ (入力記号 X) を生成する規則集合 $G_{i,j,x}^k$ の各々につき下記の [ステップ4-1-1] を行う。

[ステップ4-1-1]

$q \in U$ なる要素 q の各々につき下記の (5-1-1-1) を行う。

[ステップ4-1-1-1]

q と $G_{i,j,x}^k$ の規則の左辺の状態を比較し、各記号毎に、

① q の同一記号の個数が $G_{i,j,x}^k$ の左辺の状態の記号の個数より多いとき、規則をすべて適用して、さらに余った記号に適用可能な規則を任意に適用してできた、すべての右辺の組合せを遷移先の状態とする。

② q の記号の個数が $G_{i,j,x}^k$ 中の規則の左辺の状態より少ないとき、当該記号については q に関係なく規則をすべて適用したものと、遷移先の状態を定める。

を行い、生成した遷移先の状態の集合を求める。

[ステップ4-2]

遷移先の状態の集合のうち包含関係 \supseteq につき最少の部分集合を V_j とする。

[ステップ4-3]

ノード j を新たなノード i として、[ステップ4] を再帰的に実施する。

[ステップ5]

上記がすべて終了すると、 $DN_i = W_i$ である。

上記の遷移規則の右辺と左辺の記号を反対にし、初期状態のかわりに終了状態の集合を用いることにより、 DM_i を求めることができる。

上記のアルゴリズム終了の保証は下記のようにして得られる。

[ステップ4] でなにもせず再帰的サーチを無限回繰り返さなければ、ノード数は有限であり、このアルゴリズムは有限ステップで終了。

[ステップ4] でさらに処理をするということとは、

$\exists q \in V$

($\neg \exists q' \in W_i (q \supseteq q')$)

であり、たは W_i に加えられる。

このような入れ換えは Q の要素数が有限であること、および個別の要素 $q \in Q$ に関する $N(q, q)$ の小さいものを追加していくことより、無限回繰り返すことはできない。従って、アルゴリズムは有限ステップで終了する。