

宇宙機搭載用コンピュータのためのフォールトトレラントOS

金川 信康[†], 井原 廣一[‡], 加藤 肇彦[‡]

†(株)日立製作所 日立研究所 ‡(株)日立製作所 宇宙技術推進本部

フォールトトレランス機能を持つオペレーティングシステム(OS)を開発した。本OSでは冗長系間の処理の時間ずれから生じるデータ処理結果の矛盾の発生を防ぐためにメールボックスによるバッファリング等のITRON仕様OS, HI68Kに備わっているタスク間通信同期機能のほかに、新たに開発したランデブー方式などを用いて冗長系間の同期を採る機能を実現した。タスク間のデータの交換方式には汎用性を持たせ、カーネルとしてのITRON仕様OSの採用と相俟って種々の用途への応用を可能にし開発成果の再利用性を高めた。

Operating System for Fault-Tolerant Onboard Space Computers

Nobuyasu KANEKAWA[†], Hirokazu IHARA[‡], and Hatsuhiko KATOH[‡]

†Hitachi Res. Lab., Hitachi, Ltd.,
4026 Kuji-cho, Hitachi -city, IBARAKI 319-12 Japan
Phone: 0294-52-5111 ext.3821
E-mail:kanekawa@hrl.hitachi.co.jp

‡Space Systems Div., Hitachi, Ltd.

The Authors developed the operating systems(hereinafter OS) with fault-tolerance. In this OS, newly developed skill called rendezvous in addition to the ITRON spec. OS - HI68K's inter-task synchronization and communication function such as data buffering in mailbox for consistency. In addition to usage of ITRON spec. OS as its kernel, standarization of data exchange among tasks accomplished high versatility.

1. まえがき

交通管制や電力システム等の社会の中核ともいえる機能をコンピュータが担うにつれて、益々信頼性がコンピュータに要求されている。

また、宇宙や海洋などの極限環境下で人間に代わって作業し、あるいは人間の活動を支援する為のコンピュータが登場している。これらのコンピュータには高い信頼性と共に耐環境性が要求されている。¹⁾⁻³⁾

高い信頼性を持つフォールトトレラントコンピュータを構成するために、同一の機能をもつモジュールを複数用意しておく(冗長化)、それらの中から適切な出力を選択する方法が一般に採られている。この場合には冗長系間のデータのコンシステンシーの維持が重要な課題となる。報告者らは、フォールトトレランスのための処理、コンシステンシーの維持及びこれに必要な同期のための機能を68000用のITRON仕様OS、HI68Kをカーネルにして開発したのでここに報告する。

2. フォールトトレランスOSの機能

(1) フォールトのマスク

フォールトトレランスOSの最大の機能は、フォールト発生の影響がシステム全体の動作に及ばないようにし、システムが誤った信号を出力しないようにすることである。その為には、図2.1に示すように同一の機能を持った冗長系を複数用意しておく、それらの出力の中から適切な出力を選択する方法が採られる。

図2.1のように冗長系間で交換した処理結果を比較する相互診断の結果と、冗長系間で交換した自己診断結果を総合的に用いて各冗長系の出力の信頼度を推定し、最も信頼度の高い出力を選択するSNV(Stepwise Negotiating Voting)方式を報告者らは提案している。^{4),5)}

*TRONは"The Real Time Operating System Nucleus"の略である。

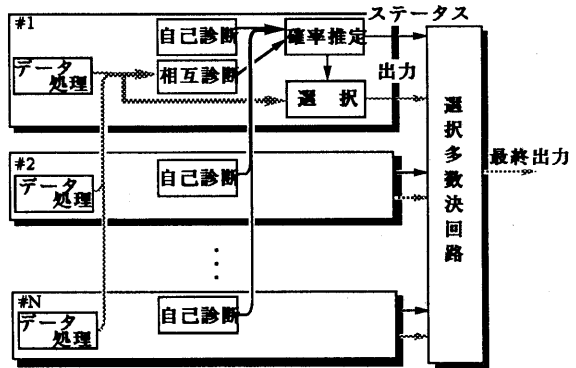


図2.1 ハードウェア構成

(2) 冗長系間の同期

データ処理が正常に進行している場合でも、異なる時刻のデータ同志の一致は保証されない。

例えば、図2.2のように冗長系Iではデータ処理が終了したが、冗長系IIではまだ終了していない場合、冗長系Iの処理済みのデータと冗長系IIの未処理のデータの間で不一致が生じてしまう。そこで、本システムのような冗長系ごとにクロックが独立している疎結合系ではフォールト・トレランスのために使うデータを冗長系間で交換する際にそれぞれの冗長系の間で、処理の同期を採り、冗長系間でデータのコンシステンシーを維持する必要がある。

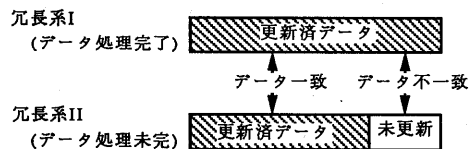


図2.2 タイミングずれによるデータ不一致

また本システムは自律分散の概念を取り入れ、システム全体の動作を司るマネージャの

役割をする部分がないシステム構成を採っている。そのため、競合及びデッドロック(3すくみ)現象の対策が必要である。

(3) フォールトからの回復(システム再構成)

フォールトの発生により一部または全部の冗長系におけるデータ処理が影響を受けた場合、速やかに冗長系の動作を停止し(セーフティ・シャットダウン)し、速やかにダウン状態から回復し正常動作に戻す機能である。他の冗長系が正常に動作している場合にはその冗長系より必要なデータを受け取り、そのデータを用いて自分自身を立ち上げ、システムを再構成し正常動作に移る。

(4) 自己の保護

他の冗長系の異常動作により自系が障害を受けないように保護をする。他の冗長系の異常動作の特に深刻な影響は、頻繁な通信要求により割込み処理に処理時間を取られ、他の処理をする時間が無くなり自分自身がシステムダウンに致ることである。このような影響を防ぐため、他の冗長系の動作を監視し、異常な動作(頻繁な通信要求など)を検出した場合、一定時間当該冗長系との接続をハードウェア的に切断する。

3. フォールトトレランスOSの開発

3.1 開発方針

フォールトトレランスOSの開発方針を以下に挙げる。

(1) 既存のOSの活用

フォールトトレランスOSの標準化と開発工程の削減のためにカーネル部分には既存のOSを活用することにした。新たにカーネルを開発するより、これらの再利用頻度の高いリアルタイムOSをカーネルに使用した方が製造上の誤り(バグ)を容易に排除でき、信頼性を向上できよう。

本開発では、わが国独自の設計思想、入手の容易性、わが国産業界への影響を考慮して

ITRON仕様リアルタイムOS、HI68Kを使用した。

(2) トランスペアレンシー

システムにアプリケーションタスクからのトランスペアレンシーを持たせ、アプリケーションタスクのフォールトトレラント機能の実現への関与を最小限にする。このため、フォールトトレラント機能は全てOSに組み込んでしまい、アプリケーション・タスクでは全くフォールトトレラントのことを考慮しなくても済むようにする。冗長系間の同期への関与は、アプリケーション・タスク実行のために必要なチェックポイントの設定のみとした。

(3) 汎用性

標準性と再利用可能性向上のために、種々のアプリケーションタスクに対応して使用できるようにした。アプリケーション・タスクの種類に依存せず使用できるようにするため、アプリケーション・タスクとフォールトトレランスOSとの間のデータの授受などのインタフェースについては特に汎用性を考慮して開発した。

ソフトウェアが再利用可能であるということは、単にソフトウェアの生産性が高いということに留まらず、ソフトウェアの信頼性の向上につながる。再利用可能なソフトウェアは繰り返し使用され、その度に使用実績が蓄積されてゆく。もしもソフトウェアをハードウェアと同様に十分に試験された高信頼性部品(ソフトウェアパッケージ)を組み合わせることができればソフトウェアの信頼性は飛躍的に向上しよう。

3.2 構成

(1) システム構成

ソフトウェアは、図3.1に示すようにカーネル、準カーネル、アプリケーションから構成される。

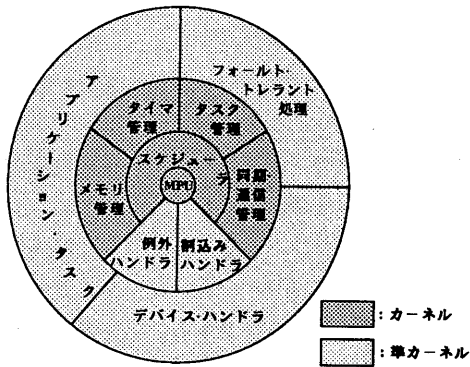


図3.1 ソフトウェアの構成

カーネルはオペレーティング・システムの中核をなすもので、本開発では開発費の削減、開発期間の短縮のためにカーネルは既存のオペレーティング・システムのものを使用した。準カーネルは例外ハンドラ、割込みハンドラ、デバイス・ハンドラ、フォールト・トレラント・タスクから構成される。これらの準カーネルの構成モジュールは新たに開発した部分である。

(a) 例外ハンドラ、割込みハンドラ

例外ハンドラはバスエラー、アドレスエラーなどのMPUの例外割込みに対する処理や、システムコールのエラー処理を実行する。

割込みハンドラは、外部割込み発生時に起動されるハンドラで、割込み処理を実行したり、他のタスクに割込み発生を通知する。

(b) デバイス・ハンドラ

外部との入出力を管理するハンドラで、例えばバブルメモリーハンドラや、サブシステム間通信ハンドラなどがある。デバイス・ハンドラは外部から入出力要求割込みを受けた割込みハンドラや、外部との入出力を要求する他のタスクによって起動される。

(c) フォールト・トレラント・ハンドラ

各種のフォールト・トレラント処理を実行するハンドラである。本ハンドラを構成するタスクは起動方法によって次のように分類される。

○フォールト・トレラント被起動タスク
アプリケーションタスクにより起動されるタスクである。起動された時点をチェック・ポイントにして、サブシステム間で各種データの照合・更新などをする。SNV方式は本タスクで実行される。

○フォールトトレラント周期起動タスク
周期的に起動するタスクである。タスク実行中には他のタスクは動作を停止する。本タスクには以下のように周期的に実行する必要がある処理、他タスクの動作の影響を受けると正常動作しない処理が含まれている。

3.3 同期

図3.2に主要なタスクの動作を示す。

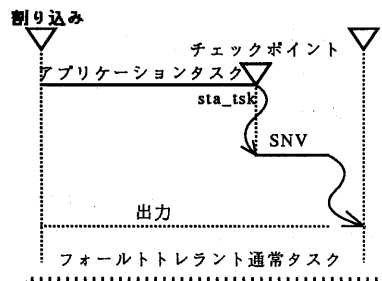


図3.2 主要なタスクの動作

フォールト・トレラント機能の存在を意識しなくともアプリケーション・タスクを作成できるように、アプリケーション・タスクの冗長系間の同期への関与は、チェックポイントでのフォールト・トレラント被起動タスクの起動のみとした。その他の冗長系間の同期はフォールトトレランスOS内部でメールボックスによるバッファリングやランデブー方式を用いてフォールトトレランスOS内部で採るようにした。表3.1にコンシステンシの維持、同期の方法をまとめる。

表3.1 コンシステンシの維持，同期の方法

| 用途 | 同期の方法 | チェックポイント | メールボックス | ランデブー | 他タスク動作抑制 | 外部割込 |
|-------------|----------------|----------|---------|-------|----------|------|
| アプリケーションタスク | | ○ | ↓ | | | |
| F | ソフトウェアポーティング | ↑ | ↓ | | | |
| | システム再構成 | ○ | ↓ | | | |
| T | 履歴データ照合訂正 | | ↓ | | | |
| O | フィードバックデータチェック | | ↓ | ○ | | |
| S | 冗長系間通信 | | ○ | | | |
| | 多重化データ照合 | | | | ○ | |
| | 入力 | | | | | ○ |
| | 出力 | | | | | ○ |

(1) メールボックスによる同期

例えば冗長系間での処理の進行状況に差が生じ，図3.3のように冗長系Iが冗長系IIより早く処理Aを終了した場合には以下の様にしてメールボックスにより同期を採る。

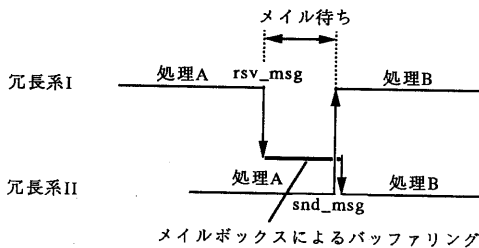


図3.3 メールボックスによる同期

冗長系Iでは処理Aを終了した後，冗長系IIから処理Aの結果のデータがメールとして送られてくるのを待ち，データを受け取った後に処理Bを開始する。このメール待ちにはタイムアウトを設け，冗長系IIがダウンしてメールが送られてこなかった場合でも，タイムアウト設定時間後には処理Bを開始して，冗長系Iのダウンを防止している。一方，冗長系IIでは冗長系Iからデータを処理A終了前に受け取っているため，処理Aが終了するまでメールボックスによりデータをバッファリングする。

ここでは，rev_msg，snd_msgシステムコー

ルを使用している。

(2) ランデブー

冗長系間で，同時に同一の処理を開始するためにランデブー方式により同期を採る。

図3.4は，同期が不要な処理Aのあと，ランデブー方式により同期を採り，同期が必要な処理Bを開始する様子を示したものである。

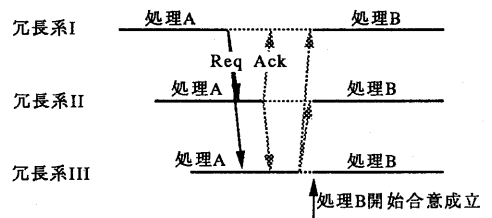


図3.4 ランデブー方式による同期

先に処理Aを終えた冗長系Iはランデブーのためのマスタになり，冗長系II，冗長系IIIに処理B開始のリクエスト信号(Req)を送る。冗長系II，冗長系IIIでは，処理Aを終えランデブー処理に入ると，冗長系Iからのリクエストに対してアクナレッジ信号(Ack)を返送する。所定の数のアクナレッジを受け取った時点で処理B開始のための同意が成立し，処理Bが開始される。

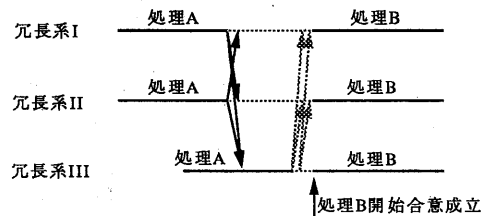


図3.5 ランデブー方式による同期

(2つの冗長系がマスタ)

また、先に処理Aを終えた冗長系が2つ以上ある場合にはそれらの冗長系がマスタになる。図3.5に2つの冗長系がマスタになった場合のランデブーの様子を示す。

なお、自分が先に処理Aを終えたかどうか（マスタかどうか）は、ランデブー処理の開始時に既に他の冗長系よりリクエストが到着しているかどうかで判定できる。

合意成立に必要なアクナレッジの数とリクエストの数との間には次のような関係がある。

(a) 自分がマスタであるとき

$$\text{Nack} = (N - N_{\text{req}} - 1) (N_{\text{req}} + 1) \dots (3-1)$$

ただし、Nack：アクナレッジの数

Nreq：リクエストの数

N：全サブシステム数

(b) 自分がマスタでないとき

$$\text{Nack} = (N - N_{\text{req}} - 1) N_{\text{req}} \quad (3-2)$$

なお、これらの式の導出については、付録で述べる。

一部の冗長系が障害などのためにアクナレッジを送れない場合には、他の冗長系はタイムアウトにより処理Bを開始する。この場合、図3.6のように全ての正常な冗長系からアクナレッジが到着した後指定したタイムアウト時間後に処理B開始の合意が成立する。

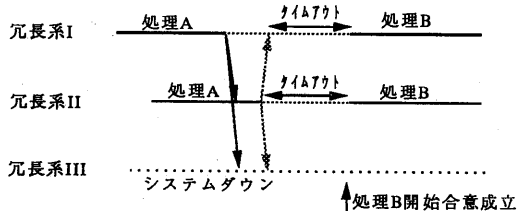


図3.6 ランデブー方式による同期

(冗長系III障害)

3.4 データ交換

(1) 冗長系間通信

自律分散の原理に従い、冗長系間に絶対的な優先順位が生じるのを避けるため、個々の冗長系は同一のアルゴリズムに従って動作する。もし、冗長系間に絶対的な優先順位が存在すると、優先順位の高い冗長系の信頼性が系全体の信頼性を左右してしまうからである。従って、冗長系間の通信でも優先順位は絶対的ではなく、常に相対的に決定される。例えばループ型のネットワークで個々の冗長系が中心に向かって右側の冗長系への通信を優先させるアルゴリズムに基づいて動作すると、冗長系間の同期が採れている場合には図3.7のように全ての冗長系が反時計廻りに通信要求を出すことになる。

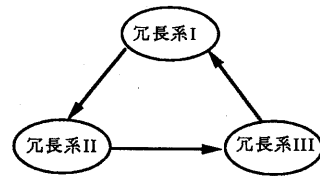


図3.7 通信要求

このとき通信方式が半2重（送信と受信を交互にする）方式であると、全ての冗長系が通信要求を送信しているために、どの冗長系も送信された通信要求を受信することができないためにデッドロックに陥ってしまう。デッドロック発生時の再試行（リトライ）の時間間隔も絶対的な優先順位付けを避けるために全ての冗長系とも同一にしなければならないため、デッドロックの発生を繰り返してしまう。このようなデッドロック発生を防ぐために、冗長系間の通信は全2重（送信と受信が同時にできる）方式とした。

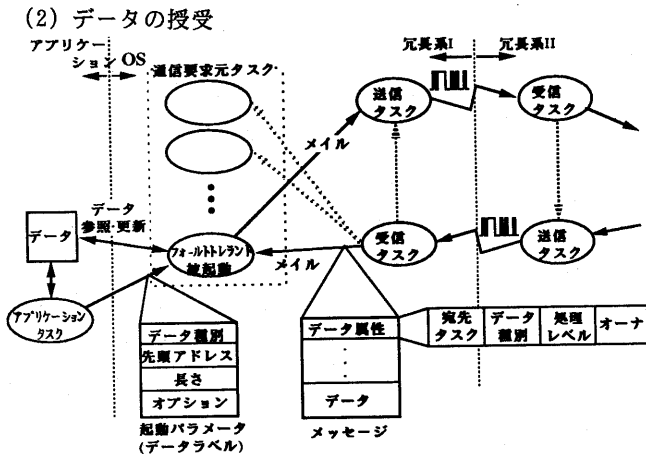


図3.9 データ交換

図3.9にデータ交換の方法を示す。アプリケーション・タスクからフォールト・トレラント被起動タスクへのデータの受け渡しは、図のようにデータ・ラベルを起動パラメータとして渡す。データ・ラベルは多数決を採るデータのデータ種別、先頭アドレス（ポイント）、データ長などを表している。フォールト・トレラント被起動タスクでは、データ・ラベルをもとに他の冗長系とデータを交換し、自身のデータが誤っている場合には、データ・ラベルにより示された領域に他の冗長系から受け取った正しいデータを書き込む。このようにデータをラベルで渡すことにより任意の形式のデータを扱うことができる。

送信、受信タスクとの間では、メールボックスを用いてデータをメールとして交換する。メールには交換すべきデータの他に、データの属性として宛先タスク名、データ種別、処理レベル、オーナー等の情報が含まれている。

4. まとめ

ITRON仕様OS, HI68Kをカーネルにしてフォールトトレランスを実現するためのOSを開発した。

アプリケーションタスクのフォールト・トレラント機能への関与の低減、タスク間で交

換するデータ形式の標準化に留意して開発したために、フォールトトレランスOSの汎用性を高めることができた。

図4.1の様にチェックポイントでフォールト・トレラント被起動タスクを起動し、その際に起動パラメータとしてデータラベルをフォールト・トレラント被起動タスクにわたせば、各種のアプリケーションにおいて容易にフォールトトレラント機能

を実現できる。

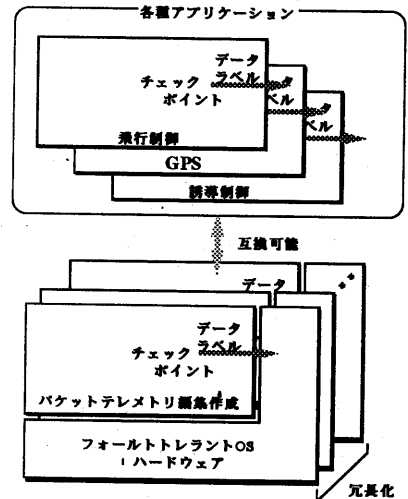


図4.1 フォールトトレランスOSの汎用性

参考文献

- 1) 山田 弘善ほか：電子機器に対する宇宙環境，信学誌，66，3，pp.254-258 (1983)
- 2) 高野 忠ほか：耐故障性を有する人工衛星搭載用コンピュータの開発，信学技報，SANE86-43，pp.25-30 (1986)
- 3) T. Takano et al.: Fault-Tolerant Onboard Computers, Proc. 16th Int'l Symp. Space Tech. and Sci., ISTS-16, Sapporo, pp.1097-1100 (1988-06)

- 4) 金川 信康ほか：“汎用性を持たせた多数決の方法に関する一提案”，信学技報，FTS88-36，pp.65-70 (1989-02)
- 5) N. Kanekawa et al.，”Dependable Onboard Computer Systems with a New Method - Stepwise Negotiating Voting，” Proc. 19th Fault-Tolerant Computing Symp.，FTCS-19，pp.13-19 (1989)
- 6) 金川 信康ほか：新しい多数決方式によるフォールトトレラントコンピュータシステム，信学論，J73-D-I，2，pp.109-116 (1990)
- 7) Tadashi Takano et al. ”Longlife Dependable Computers for Spacecrafts，” IFIP Int'l Working Conf. Dependable Computing for Critical Applications，(1989)
- 8) 加藤 肇彦，石川 知雄：マイクロコンピュータ用オペレーティングシステムの動向，情報 処理，20，6，pp.717-724 (平1-6)
- 9) H I 6 8 K ユーザーズマニュアル，(株)日立製作所 (昭62-3)
- 10) 加藤 肇彦，茶木 英明：新マイクロプロセッサH32とその開発環境，情報研究報告，89-MC-56-3，(1989)

付 録

アクナレッジの数とレクエストの数の関係

(a) 自分がマスタのとき

自分自身は，自分以外のマスタが出したリクエストを受け取るから，受け取るリクエストの数 N_{req} は，

$$N_{req} = N_m - 1 \quad (A-1)$$

ただし， N_m ：マスタの数

となる。

それぞれのスレーブ（マスタにならなかった冗長系）は，マスタから受け取ったリクエストの数だけアクナレッジを出すから，受け

取るアクナレッジの数は，

$$N_{ack} = N_s \cdot N_m \quad (A-2)$$

ただし， N_s ：スレーブの数

となる。マスタ以外の冗長系はスレーブであるから，

$$N_s = N - N_m \quad (A-3)$$

となる。

ここで，式(A-1)，(A-2)，(A-3)から，

$$N_{ack} = (N - N_{req} - 1)(N_{req} + 1) \quad (A-4)$$

となる。

(b) 自分がスレーブのとき

自分自身は全てのマスタの出したリクエストを受け取るから，

$$N_{req} = N_m \quad (A-5)$$

となる。

自分自身は，自分以外のスレーブの出したアクナレッジを受け取り，それぞれのスレーブ（マスタにならなかった冗長系）は，マスタから受け取ったリクエストの数だけアクナレッジを出すから，受け取るアクナレッジの数は，

$$N_{ack} = (N_s - 1) \cdot N_m \quad (A-6)$$

となる。ここで，式(A-5)，(A-6)，(A-3)から，

$$N_{ack} = (N - N_{req} - 1) N_{req} \quad (A-7)$$

となる。