

(19) 日本国特許庁 (JP)

(12) 特 許 公 報 (B2)

(11) 特許番号

特許第4119971号
(P4119971)

(45) 発行日 平成20年7月16日 (2008. 7. 16)

(24) 登録日 平成20年5月9日 (2008. 5. 9)

(51) Int. Cl.

F I

G 0 6 F 13/14 (2006.01)

G 0 6 F 13/14 3 2 O A

G 0 6 F 13/14 3 2 O H

請求項の数 19 (全 27 頁)

(21) 出願番号 特願平10-185331
 (22) 出願日 平成10年6月30日 (1998. 6. 30)
 (65) 公開番号 特開平11-120117
 (43) 公開日 平成11年4月30日 (1999. 4. 30)
 審査請求日 平成17年6月24日 (2005. 6. 24)
 (31) 優先権主張番号 08/885149
 (32) 優先日 平成9年6月30日 (1997. 6. 30)
 (33) 優先権主張国 米国 (US)

(73) 特許権者 597004720
 サン・マイクロシステムズ・インコーポレ
 ーテッド
 Sun Microsystems, Inc.
 アメリカ合衆国 カリフォルニア州 95
 054 サンタ クララ ネットワーク
 サークル 4150
 (74) 代理人 100064621
 弁理士 山川 政樹
 (72) 発明者 スリニバサン・ヴィスワナサン
 アメリカ合衆国・94536・カリフォル
 ニア州・フレモント・レキシントン スト
 リート・38725・アパートメント 3
 39

最終頁に続く

(54) 【発明の名称】 グローバル・ファイル・システムを基礎とした、クラスタ上の装置を可視的にするシステムおよび方法

(57) 【特許請求の範囲】

【請求項 1】

物理装置が接続されている複数のノードを含むクラスタ上の物理装置をグローバルに可視的にするシステムであって、

クラスタ上で実行するグローバル・ファイル・システムと、

それぞれが物理装置の対応するサブセットを管理するものであって、それぞれがグローバルに一意のメジャー番号に関連づけられている複数の装置ドライバと、

各装置に対して、その装置がどのノードからもアクセスされるように、少なくとも1つのグローバルに一意の識別子を確立するよう構成された装置レジストラとを備え、

そのグローバルに一意の識別子が、

クラスタのユーザが装置を識別するグローバルに一意の論理名と、そして

グローバル・ファイル・システムが装置を識別するグローバルに一意の物理名と

を含み、

グローバルに一意の物理名が、

装置を管理するために使用される装置ドライバに対応する第1の部分と、そして

装置を管理するために使用される装置ドライバに関してグローバルに一意であるマイナ
 ー名に対応する第2の部分と

を含み、

それぞれの装置の論理名と物理名との間の1対1のマッピングが装置レジストラによっ

て確立される、
システム。

【請求項 2】

さらに、

前記装置レジストラによって維持され、クラスタ内の装置の物理的関連を表す装置情報 (dev__info) データ構造であって、その各物理的関連は、前記グローバル・ファイル・システムによって維持されている装置ファイルの物理名に対応している、装置情報 (dev__info) データ構造を備え、

前記装置レジストラが、

接続された装置に対してグローバルに一意の装置タイプ (dev__t) 値を決定し、
dev__info データ構造内にエントリを作成するとともに、接続された装置に対する
対応する物理名を作成し、

接続した装置に対して、dev__t 値に基づく論理名と対応する物理名を作成し、
接続された装置の dev__t 値を、接続された装置を表す装置ファイルに関連付ける、
ように構成されている請求項 1 に記載のシステム。

【請求項 3】

前記 dev__info データ構造が、クラスタ内の装置の物理的関連を表す dev__info ツリーを含み、ルート・ノードとツリーのリーフ・ノードが、グローバル・ファイル・システムによって維持されている装置の一つを表す、装置のグローバルな一意の物理名に対応している請求項 2 記載のシステム。

【請求項 4】

それぞれのノードが装置レジストラの少なくとも一部と dev__info データ構造とのホストとなる請求項 2 記載のシステム。

【請求項 5】

各前記装置ドライバが、ローカル・マイナー番号を管理する各接続されたローカル装置に割り当てられ、

前記装置レジストラは、前記装置ドライバによって装置の特定の一つに割り当てられたローカル・マイナー番号がグローバルに一意であるかどうかを判定し、そのローカル・マイナー番号がグローバルに一意でないときそのローカル・マイナー番号をグローバル・マイナー番号に変換し、さらに、特定の装置に対応する装置ドライバに対してはメジャー番号を用い、特定の装置に対してはグローバル・マイナー番号を用いて特定の装置のためのグローバルに一意の dev__t 値を生成するように構成された装置ドライバ・インタフェース (DDI) を備えている、

請求項 1 に記載のシステム。

【請求項 6】

前記装置レジストラは、

グローバルに一意の dev__t 値をグローバルに一意の物理名にマッピングするように構成された装置ドライバインターフェースと、

グローバルに一意の物理名をグローバルに一意の論理名にマッピングするように構成されたリンク生成器と

を備えている請求項 2 記載のシステム。

【請求項 7】

装置が所定の装置クラスに分類され、これらのクラスが次のうちの少なくとも 1 つ、すなわち、

個別に列挙されているノード上の特定のドライバによって管理される、少なくとも 1 回のオカレンスを有する装置を指定する「dev__enumerate」と、

ローカルにアクセスされ、各ノード上の管理側装置ドライバと 1 対 1 の関係を有する、各ノード上で使用できる装置を指定する「dev__nodespecific」と、

そのようなノード上の装置ドライバによってアクセスできる装置を指定する「dev__global」と、

特定のノード上のドライバからアクセスされ、装置ドライバと１対１の関係を有する、「`dev__nodebound`」と
のうち少なくとも１つを含んでいる請求項 6 に記載のシステム。

【請求項 8】

クラスタ内の各装置に対して、装置を管理する装置ドライバのメジャー番号と、ローカル・マイナー番号と、グローバル・マイナー番号と、装置のホストとなるノードの `id` とを含んでいる永続的な装置構成システム (DCS) データベースを維持する、クラスタ・ノードの 1 つでホストとなる装置構成システム (DCS) をさらに備えており、DCS と DCS データベースへの要求によって、DDI がグローバルに一意の識別子を生成する請求項 7 に記載のシステム。

10

【請求項 9】

DCS が `dev__enumerate` タイプの各装置に対して、同じドライバに関連する全ての装置インスタンスのための、それぞれのグローバルに一意のグローバル・マイナー番号を生成するように構成され、
特定のノードと関係されていないそれぞれの一意のグローバル・マイナー番号に関連させられている `dev__enumerate` タイプの各装置に対して、DCS が、第 1 DCS データベース記録と、ノードの番号のいくつかに対応する複数の第 2 DCS データベース記録を生成するように構成されており、各第 2 DCS データベース記録はそれぞれの一意のグローバル・マイナー番号に関連しており、第 1 DCS データベース記録によってユーザが `dev__enumerate` タイプ装置にローカルにアクセスでき、第 2 DCS データベース記録によってユーザがあらゆるノードの `dev__enumerate` タイプ装置にアクセスできるようになっている請求項 8 に記載のシステム。

20

【請求項 10】

装置クラスのそれぞれ用の 1 つの DSO をさらに備え、それぞれの DSO を用いて、それぞれの装置クラスの装置インスタンスにグローバル・マイナー番号を割当て、DCS がグローバル・マイナー番号を割り当てるために DSO を使用するように構成されている請求項 9 に記載のシステム。

【請求項 11】

関連づけられた物理装置を有する少なくとも 1 つのサブセットの複数のクラスタ・ノードの各ノードにおける使用のためのコンピュータ・プログラム製品であって、

30

それぞれが物理装置の対応するサブセットを管理するものであって、それぞれがグローバルに一意のメジャー番号に関連づけられている複数の装置ドライバと、そして

各装置に対して、その装置がどのノードからでもアクセスされるようにする、少なくとも一つのグローバルに一意の識別子を確立するように構成された装置レジストラとを備え、

そのグローバルに一意の識別子が、

クラスタのユーザが装置を識別するグローバルに一意の論理名と、そして

グローバル・ファイル・システムが装置を識別するグローバルに一意の物理名と

を含み、

グローバルに一意の物理名が、

40

装置を管理するために使用される装置ドライバに対応する第 1 の部分と、そして

装置を管理するために使用される装置ドライバに関してグローバルに一意であるマイナー名に対応する第 2 の部分と

を含み、

それぞれの装置の論理名と物理名との間の 1 対 1 のマッピングが装置レジストラによって確立される、

コンピュータ・プログラム製品。

【請求項 12】

さらに、

前記装置レジストラによって維持され、クラスタ内の装置の物理的関連を表す装置情報 (

50

d e v _ i n f o) データ構造であって、その各物理的関連は、前記グローバル・ファイル・システムによって維持されている装置ファイルの物理名に対応している、装置情報 (d e v _ i n f o) データ構造を備え、

前記装置レジストラが、

接続された装置に対してグローバルに一意の装置タイプ (d e v _ t) 値を決定し、
d e v _ i n f o データ構造内にエントリを作成するとともに、接続された装置に対する
対応する物理名を作成し、

接続した装置に対して、d e v _ t 値に基づく論理名と対応する物理名を作成し、
接続された装置の d e v _ t 値を、接続された装置を表す装置ファイルに関連付ける
ように構成されている請求項 1 1 に記載のコンピュータ・プログラム製品。

10

【請求項 1 3】

前記 d e v _ i n f o データ構造が、クラスタ内の装置の物理的関連を表す d e v _ i n f o ツリーを含み、ルート・ノードとツリーのリーフ・ノードが、グローバル・ファイル・システムによって維持されている装置の一つを表す、装置のグローバルな一意の物理名に対応している請求項 1 2 に記載のコンピュータ・プログラム製品。

【請求項 1 4】

それぞれのノードが装置レジストラの少なくとも一部と d e v _ i n f o データ構造との
ホストとなる請求項 1 2 に記載のコンピュータ・プログラム製品。

【請求項 1 5】

各前記装置ドライバが、ローカル・マイナー番号を管理する各接続されたローカル装置
に割り当てられ、

20

前記装置レジストラは、前記装置ドライバによって装置の特定の一つに割り当てられたローカル・マイナー番号がグローバルに一意であるかどうかを判定し、そのローカル・マイナー番号がグローバルに一意でないときそのローカル・マイナー番号をグローバル・マイナー番号に変換し、さらに、特定の装置に対応する装置ドライバに対してはメジャー番号を用い、特定の装置に対してはグローバルマイナー番号を用いて特定の装置のためのグローバルに一意の d e v _ t 値を生成するように構成された装置ドライバ・インタフェース (D D I) を備えている請求項 1 4 に記載のコンピュータ・プログラム製品。

【請求項 1 6】

前記装置レジストラは、
グローバルに一意の d e v _ t 値をグローバルに一意の物理名にマッピングするように構成された装置ドライバインタフェースと、
グローバルに一意の物理名をグローバルに一意の論理名にマッピングするように構成されたリンク生成器と備えている請求項 1 2 に記載のコンピュータ・プログラム製品。

30

【請求項 1 7】

装置が所定の装置クラスに分類され、これらのクラスが次のうちの少なくとも 1 つ、すなわち、

個別に列挙されているノード上の特定のドライバによって管理される、少なくとも 1 回の
オカレンスを有する装置を指定する「d e v _ e n u m e r a t e」と、

ローカルにアクセスされ、各ノード上の管理側装置ドライバと 1 対 1 の関係を有する、各
ノード上で使用できる装置を指定する「d e v _ n o d e s p e c i f i c」と、

40

そのようなノード上の装置ドライバによってアクセスできる装置を指定する「d e v _ g l o b a l」と、

特定のノード上のドライバからアクセスされ、装置ドライバと 1 対 1 の関係を有する、「
d e v _ n o d e b o u n d」と

のうち少なくとも 1 つを含んでいる請求項 1 6 に記載のコンピュータ・プログラム製品。

【請求項 1 8】

クラスタ内の各装置に対して、装置を管理する装置ドライバのメジャー番号と、ローカル・マイナー番号と、グローバル・マイナー番号と、装置のホストとなるノードの i d とを含んでいる永続的な装置構成システム (D C S) データベースを維持する、クラスタ・

50

ノードの1つでホストとなる装置構成システム(DCS)をさらに備えており、DCSとDCSデータベースへの要求によって、DDIがグローバルに一意の識別子を生成する請求項17に記載のコンピュータ・プログラム製品。

【請求項19】

DCSがdev__enumerateタイプの各装置に対して、同じドライバに関連する全ての装置インスタンスのための、それぞれのグローバルに一意のグローバル・マイナー番号を生成するように構成され、

特定のノードと関係されていないそれぞれの一意のグローバル・マイナー番号に関連させられているdev__enumerateタイプの各装置に対して、DCSが、第1DCSデータベース記録と、ノードの番号のいくつかに対応する複数の第2DCSデータベース記録を生成するように構成されており、各第2DCSデータベース記録はそれぞれの一意のグローバル・マイナー番号に関連しており、第1DCSデータベース記録によってユーザがdev__enumerateタイプ装置にローカルにアクセスでき、第2DCSデータベース記録によってユーザがあらゆるノードのdev__enumerateタイプ装置にアクセスできるようになっている請求項18記載のコンピュータ・プログラム製品。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は一般に、ファイル・システムを通じて装置アクセスを提供するシステムおよび方法に係り、特に、クラスタ上の装置をグローバルに可視的にするシステムおよび方式に関する。

【0002】

【従来の技術】

複数のコンピュータを含むクラスタ上でUnixベースのコンピュータ・アプリケーションを使用することがますます一般的になっている。クラスタ・オペレータのターゲットは、クラスタの動作を、単一のコンピュータの場合と同様にアプリケーション/ユーザに対して透過的なものにすることである。たとえば、クラスタは通常、ユーザが、ファイルがどこに存在するかにかかわらずクラスタ上のすべての従来型のファイルを見てそれにアクセスできるようにするグローバル・ファイル・システムを用意している。しかし、上記透過性はクラスタ上の装置アクセスには及ばない。

【0003】

通常、Unixベースのシステム上の装置アクセスは、装置をファイルとして扱う特殊ファイル・システム(たとえば、SpecFS)を通して行われる。この特殊ファイル・システムは、単一のノード上でのみ動作する。すなわち、このファイル・システムでは、特定のノード上の装置を見てそれにアクセスすることができるのはそのノード上のユーザだけであり、これは、クラスタ上のグローバル装置可視性のターゲットに反する。このような制限は、様々なノード上で動作する特殊ファイル・システムの間で調整が行われないことと、装置のグローバル可視性に対処する装置命名方式がないことによるものである。次に、従来技術の装置アクセス・システムのこのような点を図1ないし4を参照して説明する。

【0004】

図1には、中央演算処理装置(CPU)102と、高速メモリ104と、複数の物理的装置106と、CPU102がデータを制御しメモリ102および物理的装置106と交換できるようにする一群の物理的装置インタフェース108(たとえばバスまたは他の電子式インタフェース)とを含む従来型のコンピュータ・システムのブロック図が示されている。メモリ102は、ランダム・アクセス・メモリ(RAM)でも、あるいはキャッシュ・メモリでもよい。

【0005】

物理的装置106には、高可用性装置112や、プリンタ114や、カーネル・メモリ116や、通信装置118や、記憶装置120(たとえば、ディスク装置)を含めることが

できるが、これらに限らない。プリンタ 114 および記憶装置 120 は良く知られている。高可用性装置 112 には、関連する二次装置を有する記憶装置やプリンタが含まれる。そのような装置は、二次装置が、そのそれぞれの一次装置が故障したときに一次装置に代替することができるので可用性が高い。カーネル・メモリ 116 は、システム性能累計・報告統計メモリ 102 のプログラム済み領域である。通信装置 118 はモデムと、ISDN インタフェース・カードと、ネットワーク・インタフェース・カードと、その他の種類の通信装置とを含む。装置 106 は、擬似装置 122、すなわち実際の物理的装置に関連付けられないソフトウェア装置を含むこともできる。

【0006】

コンピュータ 100 のメモリ 104 は、オペレーティング・システム 130 と、アプリケーション・プログラム 150 と、データ構造 160 を記憶することができる。オペレーティング・システム 130 は、コンピュータ 100 が動作可能であるかぎり CPU 102 内で実行され、プロセッサ 102 と、CPU 102 内で実行中のアプリケーション 150 とのためのシステム・サービスを行う。オペレーティング・システム 130 は、Sun ワークステーション上で使用される SolarisTM オペレーティング・システムの v. 2.6 に基づいて構成され、カーネル 132 と、ファイル・システム 134 と、装置ドライバ 140 と、装置ドライブ・インタフェース (DDI) フレームワーク 142 とを含む。Solaris および Sun はそれぞれ、サンマイクロシステムズ社 (Sun Microsystems, Inc.) の商標および登録商標である。カーネル 116 は、メモリ 104、ファイル・システム 134、または装置 106 へのアクセスを求める要求など、アプリケーション 150 からのシステム呼出しを処理する。ファイル・システム 134 およびその装置 106 および装置ドライバ 140 との関係について図 2A および 2B を参照して説明する。

【0007】

図 2A を参照するとわかるように、Solaris オペレーティング・システムの v. 2.6 およびそれ以前のバージョンによって使用されるファイル・システム 134 の高レベル表現が示されている。Solaris では、ファイル・システム 134 は、すべてのファイル、装置 106、ネットワーク・インタフェース (コンピュータ 100 はネットワーク化されていると仮定する) がアクセスされる媒体である。この 3 つの異なる種類のアクセスはそれぞれ、ファイル・システム 134 の 3 つの構成要素、すなわち Unix ファイル・システム 138u (UFS)、特殊ファイル・システム 138s (SpecFS)、ネットワーク・ファイル・システム 138n (NFS) によって行われる。

【0008】

Solaris では、アプリケーション 150 は最初、ファイル、装置、またはネットワーク・インタフェース (本明細書では「ターゲット」と呼ぶ) に対するオープン要求をカーネル 132 を介してファイル・システム 134 に発行することによってそのターゲットにアクセスする。ファイル・システム 134 は次いで、必要に応じて UFS 138u、SpecFS 138s、または NFS 138n へ要求を送る。ターゲットが首尾良くオープンされた場合、UFS、SpecFS、または NFS は、要求されたファイル、装置、またはネットワーク・ノードにマップされている v ノード・オブジェクト 136 をファイル・システム 134 に返す。ファイル・システム 134 は次いで、v ノード・オブジェクト 136 をファイル記述子 174 にマップし、そのファイル記述子はカーネル 132 を介してアプリケーション 150 に返される。要求側アプリケーションはその後で、ファイル記述子 174 を使用して、返された v ノード・オブジェクト 136 に関連付けられた対応するファイル、装置、またはネットワーク・ノードにアクセスする。

【0009】

v ノード・オブジェクト 136 は、カーネル 132 とファイル・システム 134 との間のインタフェースとして働く v ノード/VFS インタフェースまたはレイヤ (VFS) 172 に従って一般的な 1 組のファイル・システム・サービスを行う。Solaris は、v ノード・オブジェクト 136 から継承し、UFS、SpecFS、NFS に関連付けられ

10

20

30

40

50

た各ターゲットタイプについてカスタマイズされたメソッドとデータ構造も含む i ノード・オブジェクト 136 i、s ノード・オブジェクト 136 s、r ノード・オブジェクト 136 r も備える。これらのクラス 136 i、136 s、136 r は、v ノード 136 とそのそれぞれのターゲットとの間の低レベル・インタフェースを形成する。したがって、UFS、SpecFS、または NFS が v ノード・オブジェクトを返すとき、そのオブジェクトは、実際のターゲット動作を実行する対応する i ノード、s ノード、または r ノードに関連付けられる。Solaris ファイル・システムの一般的な性質について論じたが、次に、Solaris で使用されるファイル・ベースの装置アクセス方法について論じる。

【0010】

図 2B を参照するとわかるように、Solaris アプリケーション 150 は通常、それがオープンする必要がある装置の論理名 166 を使用してファイル・システム 134 (カーネル 132 を介して) に装置アクセス要求を発行する。たとえば、アプリケーション 150 はコマンド `open(/dev/dsk/disk_logical_address)` を用いて SCSI 装置へのアクセスを要求することができる。論理名 `/dev/dsk/disk_logical_address` は、オープンする装置が特定の論理アドレスにあるディスクであることを示す。Solaris では、SCSI ディスクの論理アドレスは「c0t0d0sx」であってよく、この場合、「c0」は SCSI コントローラ 0 を表し、t0 はターゲット 0 を表し、d0 はディスク 0 を表し、sx は特定のディスクの x 番目のスライスを表す (SCSI ディスク・ドライバは最大で 8 つのスライスを有することができる)。

【0011】

論理名は、1 つのリンク・ジェネレータ 144、すなわち DDI フレームワーク 142 の 1 つのユーザ空間拡張子によって割り当てられ、装置の接続時に装置のドライバ 140 から供給される情報と DDI フレームワーク 142 によって生成される装置の対応する物理名とに基づくものである。特定の装置ドライバ 140 のインスタンスがノード 100 に接続されると、DDI フレームワーク 142 はそのドライバ 140 の接続ルーチン呼び出す。ドライバ 140 は次いで、そのインスタンスに関連付けることのできる各装置について DDI フレームワーク 142 の `ddi_create_minor_nodes` メソッド 146 に固有のローカル識別子を割り当てそれを呼び出す。通常、この固有のローカル識別子はマイナー名 (たとえば、「a」) およびマイナー番号 (たとえば、「2」) を構成する。`ddi_create_minor_nodes` メソッド 146 は、呼び出されると、所与の装置を表すリーフ・ノードを `DevInfo` ツリー 162 内に作成する。たとえば、SCSI ドライブ (すなわち、インスタンス) は最大で 8 つのスライス (すなわち、デバイス) を有することができるので、ローカル SCSI ドライブ 140 は 8 つのスライスのそれぞれに固有のローカル識別子を割り当て、ローカル識別子を用いて `ddi_create_minor_nodes` メソッド 146 を最大で 8 回呼び出す。

【0012】

各装置 106 には、ターゲット装置 106 に関する構成情報を与える UFS ファイル 170 も関連付けられる。特定の UFS ファイル 170 i の名前は、コンピュータ上の装置の物理的位置から導かれる物理名 168 i と同じである。たとえば、SCSI 装置は物理名 `168/devices/iommu/sbus/esp1/sd@addr:minor_name` を有することができ、この場合、`addr` は装置ドライバのアドレスであり、`minor_name` は装置インスタンスのマイナー名である。この名前は装置ドライバ `sd` によって割り当てられる。物理名がどのように割り当てられるかについては以下で図 3 を参照して説明する。

【0013】

ファイル・システム 134 は、それ自体が、ターゲット装置の論理名が与えられた場合にターゲット装置をオープンできるようにするために、論理ファイル名 166 を物理ファイル名 168 にマップする論理名空間データ構造 164 を使用する。装置 106 の物理名は

10

20

30

40

50

、コンピュータ・システム 100 に関連付けられた装置タイプ、バス接続、コントローラ、ドライバ、装置の階層を表す装置情報 (DevInfo) ツリー 140 (図 1 に示されている) 内の装置の位置から導かれる。物理名 168 によって識別される各ファイル 170 は、その属性として、ターゲット装置に固有に関連付けられた識別子または dev__t (装置タイプの短縮形) を含む。この dev__t 値は、SpecFS 138s を介して正しいターゲット装置にアクセスするためにファイル・システム 134 によって使用される。次に、図 3 を参照して、dev__t 値がどのように割り当てられ、DevInfo ツリー 140 がどのように DDI フレームワーク 142 によって維持されるかについて説明する。

【0014】

図 3 を参照すると、コンピュータ・システム 100 に関する仮説的な DevInfo ツリー 162 の図が示されている。DevInfo ツリー 162 の各ノードは、コンピュータ 100 に関連付けられた装置システムの物理的構成要素に対応する。それぞれの異なるレベルは、それぞれの異なるレベルの装置階層に対応する。上位のノードに直接接続されたノードは、上位のレベルのオブジェクトのインスタンスであるオブジェクトを表す。したがって、DevInfo ツリーのルート・ノードは常に「I」ノードであり、この下に装置階層全体が存在する。中間ノード (すなわち、リーフ・ノードおよびリーフ親ノード以外のノード) は、ネクサス装置と呼ばれ、コントローラ、バス、ポートなどの中間構造に対応する。DevInfo ツリーの最も低いレベルの次に装置ドライバがあり、それぞれ、1 つまたは複数の装置をエクスポートまたは管理することができる。リーフ・レベルには、それぞれ、装置タイプに応じていくつかの装置インスタンスをエクスポートできる実際の装置がある。たとえば、SCSI 装置は最大 7 個のインスタンスをもつことができる。

【0015】

図 3 に示した仮説的な DevInfo ツリー 162 は、物理アドレス addr0 にあるメモリ・マップされた入出力装置 (iommu) 用の入出力コントローラを含むコンピュータ・システム 100 を表す。iommu は、アドレス addr1 にあるシステム・バス (sバス) およびアドレス addr2 にある PCI バスなどの高速バスに接続された入出力装置と CPU の対話を管理する。それぞれのアドレス addr3 および addr4 にある 2 つの SCSI コントローラ (esp1 および esp2) は、アドレス addr5 にある非同期転送モード (ATM) コントローラと共に s バスに結合される。第 1 の SCSI コントローラ esp1 は、4 つの SCSI 装置インスタンス (dev0、dev1、dev2、dev3) を管理するアドレス 0 (@0 として表される) にある SCSI 装置ドライバ (sd) に関連付けられる。これらの装置インスタンスはそれぞれ、単一の物理装置 106 のそれぞれのスライスに対応する。第 1 の SCSI コントローラ esp1 は、他の物理装置 106 の複数の SCSI 装置インスタンス (図示せず) を管理するアドレス 1 にある SCSI 装置ドライバ (sd) にも関連付けられる。

【0016】

コンピュータ・システム 100 と共に使用できる各装置ドライバ・タイプには所定の固有のメジャー番号が割り当てられる。たとえば、SCSI 装置ドライバ sd にはメジャー番号 32 が割り当てられる。各装置には、単一の装置ドライバによって管理される一群の装置内で固有のマイナー番号が割り当てられる。たとえば、アドレス 0 にあるドライバ sd に関連付けられた装置 dev0、dev1、dev2、dev3 はそれぞれ、マイナー番号 0、1、2、3 とマイナー名 a、b、c、d とを有する。同様に、アドレス 1 にあるドライバ sd によって管理される装置は、装置 dev0 ないし dev3 に関連付けられたマイナー番号とは異なるマイナー番号 (たとえば、マイナー番号 4 ないし 7 を有する 4 つ) を有する。マイナー番号およびマイナー名は、新しい各装置インスタンスごとに親装置ドライバ 140 (図 1) によって割り当てられる (SCSI インスタンスが、特定の SCSI ドライブでよく、SCSI 装置がそのドライブの特定のスライスでよいことを想起されたい)。これによって所与の装置ドライバによってエクスポートされる各装置は固有のマ

10

20

30

40

50

イナー番号とマイナー名とを有する。すなわち、ドライバはマイナー番号 - 名空間を管理する。

【 0 0 1 7 】

各マイナー番号は、親ドライバのメジャー番号と組み合わせられると、各装置を固有に識別する `dev__t` 値を形成する。たとえばアドレス 0 にあるドライバ `sb` によって管理される装置 `dev 0`、`dev 1`、`dev 2`、`dev 3` は、それぞれの `dev__t` 値 (3 2 , 0)、(3 2 , 1)、(3 2 , 2)、(3 2 , 3) を有する。 `SpecFS 1 3 8 s` は `dev__t` 値の対応する装置へのマッピングを維持する。そのため、 `SpecFS` へのすべての装置オープン要求は、固有の `dev__t` 値を使用してオープンすべき装置を識別する。

【 0 0 1 8 】

装置への `DevTree` パスはその装置の物理名を形成する。たとえば、装置 `dev 0` の物理名は文字列 `/ devices / i o m m u @ a d d r 0 / s b u s @ a d d r 1 / e s p 1 @ a d d r 3 / s d @ 0 : a` によって与えられ、この場合、 `sd@0 : a` は、アドレス 0 にある `sd` ドライバによって管理されマイナー名が `a :` である装置、すなわち装置 `dev 0` を指す。この物理名は、対応する装置にアクセスするのに必要なすべての情報を保持する (`s` ノードに対応する) 特殊ファイル 1 7 0 (図 2 に示されている) を識別する。特に、各特殊ファイル 1 7 0 の属性は、対応する装置に関連付けられた `dev__t` 値を保持する。

【 0 0 1 9 】

前述のように、リンク・ジェネレータ 1 4 4 は、リンク・ジェネレータによって管理される装置に適用できる 1 組の規則に従って装置の物理名から装置の論理名を生成する。たとえば、アドレス 0 にあるドライバ `sd` によって管理される装置 `dev 0` の場合、 `SCSI` 装置のリンク・ジェネレータは論理名 `/ dev / d s k / c 0 t 0 d 0 s 0` を生成することができる。この場合、 `c 0` はコントローラ `esp 1 @ a d d r 3` を指し、 `t 0` は `sd@0` ドライバによって管理される物理ディスクのターゲット `id` を指し、 `d 0` は `sd@0` ドライバを指し、 `s 0` はマイナー名 `a` とマイナー番号 0 とを有するスライス指定する。 `sd@1` ドライバに関連付けられた装置 `dev 0` には、同じリンク・ジェネレータ 1 4 4 によって論理名 `dev / d s k / c 0 t 1 d 1 s 4` が割り当てられる。2 つの `dev 0` 装置が、ターゲット値、ディスク値、スライス値の差によって区別される論理名を有することに留意されたい。次に、図 4 を参照して、現在 `Solaris` でこのインフラストラクチャをどのように使用して、コンピュータ 1 0 0 上に存在する特定の装置をオープンできるようにしているかを説明する。

【 0 0 2 0 】

図 4 を参照すると、アプリケーション 1 5 0 の要求に応じて装置をオープンする際に様々なオペレーティング・システム構成要素によってコンピュータ 1 0 0 のメモリ 1 0 4 で実行される動作の流れ図が示されている。メモリ 1 0 4 は、アプリケーション 1 5 0 が実行されるユーザ空間 1 0 4 U と、オペレーティング・システム構成要素が実行されるカーネル空間 1 0 4 K に分割される。この図では、動作が行われる順序と、各動作の送信元およびターゲットである装置が 1 組のラベル付き矢印で示されている。必要に応じて、破線は、参照が渡されるオブジェクトを示す。メモリ 1 0 4 の図の隣に、ラベル付き矢印に関連する各動作が定義されている。動作はメッセージまたは関数呼出しとして定義され、その場合、メッセージ名の後に、処理すべきデータまたは受信側エンティティによって返されるデータが続く。たとえば、メッセージ (4 - 1) 「 `open (logical__name)` 」は、「 `logical__name` 」によってユーザ空間 1 0 4 U 内に表された装置をオープンすることをカーネル 1 3 2 に要求する、アプリケーション 1 5 0 から発行されるメッセージである。この特定の例では、アプリケーションは装置 `dev 2` をオープンすることを要求する。

【 0 0 2 1 】

オープン・メッセージ (4 - 1) を受信した後、カーネル 1 3 2 はメッセージ (4 - 2) 「 `get__vnode (logical__name)` 」をファイル・システム 1 3 4 に発

10

20

30

40

50

行する。このメッセージは、カーネル 132 がオープン動作を完了するのに必要な装置 dev2 の v ノードを返すようファイル・システム 134 に要求する。これに応答して、ファイル・システム 134 は、論理名空間 164 を使用して論理名 166 を対応する物理名 168 に変換する。ファイル・システム 134 は次いで、物理名によって指定されたファイルを見つけ、そのファイルの属性から対応する装置の dev__t 値を判定する。ファイル・システム 134 は、dev__t 値を得た後、メッセージ (4-3) 「get__vnode (dev__t)」を SpecFS 138s に発行する。このメッセージは、装置 dev2 にリンクされた v ノードの参照を返すことを SpecFS 138s に要求する。SpecFS 138s は、メッセージ (4-3) を受信した後、要求された v ノード 136 と、v ノード 136 を装置 dev2 にリンクする s ノード 136s とを作成し、v ノード 136 の参照 (4-4) をファイル・システム 134 に返す。ファイル・システム 134 は次いで、v ノード参照をカーネルに返す (4-5)。

【0022】

カーネル 132 は、v ノード参照を得た後、v ノード 136 に関連付けられた装置 dev2 をオープンすることを求める要求 (4-6) を SpecFS 138s に発行する。SpecFS 138s は、装置 dev2 を管理することが知られているドライバ 2 にオープン・コマンド (4-7) を発行することによってこの要求を満たそうとする。ドライバ 2 は、装置 dev2 をオープンできる場合、オープン動作が成功したことを示す open__status メッセージ (4-8) を返す。そうでない場合、ドライバ 2 は同じメッセージ (4-8) で失敗表示を返す。SpecFS 138s は次いで、同様な状況メッセージ (4-9) をカーネル 132 に直接返す。メッセージ (4-9) で「成功」が返されたと仮定した場合、カーネル 132 は、ファイル記述子、すなわち装置 dev2 にリンクされた v ノード 136 のユーザ空間表現をアプリケーション 150 に返す (4-10)。アプリケーション 150 は、ファイル記述子を得た後、ファイル・システム動作を使用してカーネル 132 およびファイル・システム 134 を介して装置 dev2 にアクセスすることができる。たとえば、アプリケーション 150 は、返されたファイル記述子に関する読取り要求を発行することによって装置 dev2 からの入力データを実行する。このようなファイル・システム・コマンドは次いで、SpecFS 136s と、装置 dev2 を管理する v ノード・オブジェクト 136 および s ノード・オブジェクト 136s とによって実際の装置コマンドに変換される。

【0023】

したがって、Solaris では、コンピュータ・システム 100 のユーザは、そのシステム 100 上の装置に比較的容易にアクセスすることができる。しかし、Solaris によって使用される方法では、いくつかの異なるコンピュータがクラスタの一部として構成されているときでも、ユーザがコンピュータを介して装置に透過的にアクセスすることはできない。すなわち、第 1 のコンピュータ上で実行されているアプリケーションが、Solaris を使用して第 2 のコンピュータ上の装置を透過的にオープンすることはできない。

【0024】

Solaris の現行のバージョンがマルチコンピュータの場合に透過的な装置アクセスを行うことができない理由は、dev__t 番号およびマイナー番号が装置が接続されるときに割り当てられる現行の方法に関するものである。再び図 3 を参照するとわかるように、コンピュータ 100 に装置が接続されるたびに、装置の関連するドライバによって制御される 1 組の装置内で固有であり、したがって、ドライバのメジャー番号と組み合わせられたときにコンピュータ 100 の固有の dev__t 値にマップすることのできるマイナー番号をドライバがその装置に割り当てる。しかし、同じ装置およびドライバが第 2 のコンピュータ上に設けられている場合、そのドライバおよび装置には、同一ではないにしても同様な 1 組のメジャー番号およびマイナー番号ならびに dev__t 値が割り当てられる。たとえば、両方のコンピュータが SCSI ドライバ sd (メジャー番号 = 32) と、SCSI ドライバ sd によって管理される 4 つの SCSI 装置インスタンスを有する場合、各

10

20

30

40

50

ドライバ `s d` はそのローカルな 1 組の `SCSI` 装置に同じ 1 組のマイナー番号を割り振る (たとえば、両方の組が 0 ないし 3 のマイナー番号を有する)。したがって、装置はその `dev__t` 値に従ってアクセスされるので、第 1 のノード・アプリケーションが第 2 のノード上の `SCSI` ディスクをオープンする必要がある場合、そのアプリケーションは、どちらのコンピュータ・システム上の `SpecFS` に対しても `SCSI` ディスクを明確に識別することができない。

【0025】

【発明が解決しようとする課題】

したがって、アプリケーションが、どこで実行されているかにかかわらず、コンピュータ・クラスタのノード上に存在する装置に透過的にアクセスできるようにするファイル・ベースの装置アクセス・システムが必要である。

10

【0026】

【課題を解決するための手段】

要約すると、本発明は、クラスタ内で使用され、ローカル装置名を、クラスタ・ノードのいずれかでランしているアプリケーションから他のクラスタ・ノードに位置する装置にアクセスできる、グローバルで一意的な論理名にマッピングするシステムおよび方法にある。

【0027】

特に、好ましい実施形態には、クラスタのグローバル・ファイル・システム、装置 (デバイス)・ドライバ・インタフェース (`DDI`) および装置 (デバイス) 情報ツリーが含まれている。装置ドライバ・インタフェースは、少なくとも予め定められた装置クラスのサブセットに関連付けられたインスタンスに対して、ローカル装置名がグローバルに一意的であるかどうか判断する。ローカル装置名がグローバルに一意的ではない場合、`DDI` はそれぞれのローカル・インスタンスに対してグローバルに一意的な識別子を決める。装置ドライバ・インタフェースによって維持される装置情報ツリーは、クラスタ内の装置の物理的関連を表す。装置情報ツリーのルート・ノードとリーフ・ノードの間の各パスは、装置インスタンスのそれぞれ 1 つを表すグローバル・ファイル・システムによって維持される装置ファイルの物理名に対応する。リンク・ジェネレータは、各インスタンスに対して、そのインスタンスの対応する物理名にマッピングされるグローバルに一意的な識別子に基づいて論理名を生成し、特定のインスタンスのグローバルに一意的な識別子を、その特定のインスタンスを表す装置ファイルに関連付ける。

20

30

【0028】

このフレームワークを与えられると、クラスタのユーザは、アクセスされる装置をその論理名で識別するアクセス要求をグローバル・ファイル・システムに出すことにより、任意の装置にアクセスできる。グローバル・ファイル・システムは、論理名を対応する物理名に変換し、物理名によって識別される装置ファイルを介して装置のグローバルに一意的な識別子にアクセスし、そのグローバルに一意的な識別子によって識別される装置にアクセスして、アクセス要求を解決する。

【0029】

好ましい実施形態においては、グローバル・ファイル・システムは、1 つのネットワーク・ノード上でしか稼動しない可能性があり、各ノードが、装置ドライバ・インタフェースと装置情報ツリーのインスタンスにそれぞれ対応する。本発明の実施形態においては、グローバル・ファイル・システムとそれぞれの装置ドライバ・インタフェースとの間のリンクは、プロキシ・ファイル・システムによって提供される。好ましい実施形態はまた、各々がローカル装置名をインスタンスに割り当てることのできる装置を管理するための複数の装置ドライバを含むこともある。これらの装置ドライバは、各種の装置ドライバ・インタフェースと共設される。

40

【0030】

好ましい予め定められた装置クラスのセットは、次を含む。すなわち、

「`dev__enumerate`」・・特定のドライバによって管理される少なくとも 1 つのインスタンス (それぞれは、個別に列挙されている特定のノード上の特定のドライバに

50

よって管理される)で装置を指定。

「dev__nodespecific」・・・ローカルにアクセスされ、各ノード上の管理ドライバと1対1の関係を有している、各ノード上で利用可能な装置を指定。

「dev__global」・・・任意のノードのドライバからアクセスできる装置を指定。

「dev__nodebound」・・・特定のノード上のドライバだけから通常アクセスされ、そのドライバと1対1の関係を有している装置を指定。

【0031】

本発明は、それぞれの装置クラスのインスタンスのローカル装置名を、対応するグローバルに一意の識別子および論理名に変換するための一連の規則(ルール)を定義するものである。特に、dev__enumerate装置のローカル装置名は、装置名およびローカルに一意のローカル・マイナー番号の組合せを含み、グローバルに一意の識別子は、グローバル・マイナー番号を含み、論理名は、装置名とグローバル・マイナー番号の組合せを含んでいる。同様に、dev__nodespecific装置のローカル装置名と論理名はそれぞれ、装置名から形成され、dev__global装置のローカル装置名と論理名はそれぞれ装置名から形成される。

【0032】

本発明はまた、クラスタ・ノードの1つのノードに装置構成システム(DCS)が組み込まれており、その1つのノードには、クラスタ内の各装置に対して、装置名、論理装置を管理する装置ドライバのメジャー番号と、グローバル・マイナー番号、および論理装置のホストとなるノードのホストidとを含んでいる永続的な装置構成システム(DCS)データベースが維持されている。DDIは、DCSデータベースを使用してDDIのために装置それぞれのグローバル・マイナー番号を生成するDCSの支援によって、各論理装置のグローバルに一意の識別子、論理名および物理名を生成する。

【0033】

本発明はまた、クラスタ上の装置をグローバルに可視的にするためのシステムでもあり、クラスタには装置が接続される複数のノードが含まれる。好ましいシステムは、少なくとも1つのグローバルに一意の識別子を各装置に設定して、その装置がどのノードからもアクセスできるようにする装置レジストラを含む。このシステムはクラスタ上でランするグローバル・ファイル・システムを含むことができ、少なくとも1つのグローバルに一意の識別子には、クラスタのユーザが装置を識別するグローバルに一意の論理名と、グローバル・ファイル・システムが装置を識別するグローバルに一意の物理名とが、備えられている。好ましい実施形態においては、装置レジストラによって設定される、各装置の論理名と物理名との間の1対1のマッピングがある。

【0034】

本発明はまた、クラスタ内の装置の物理的関連を表す装置レジストラによって維持される装置情報(dev__info)データ構造も含み、当該物理的関連のそれぞれは、装置のそれぞれを表しグローバル・ファイル・システムによって維持されている装置ファイルの物理名に対応している。このフレームワークを与えられると、装置レジストラは好ましい実施形態では、接続された装置に対してグローバルに一意の装置タイプ(dev__t)値を決定し、dev__infoデータ構造にエントリおよび対応する物理名を接続された装置に対して作成し、接続された装置に対してdev__t値に基づく論理名と、対応する物理名とを作成し、そして接続された装置のdev__t値を、接続された装置を表す装置ファイルに関連付ける。

【0035】

本発明の他の目的および特徴は、以下の詳細な説明および添付の特許請求の範囲を図面と共に検討したときにより容易に明らかになる。

【0036】

【発明の実施の形態】

図5を参照すると、本発明を実施することのできるコンピュータ・クラスタ200のブロック図が示されている。クラスタ201は、関連する装置106とアプリケーション15

10

20

30

40

50

0とを有する複数のノード202を含む。図1と同様に、装置106は高可用性装置112と、プリンタ114と、カーネル・メモリ116と、通信装置118と、記憶装置120とを含むことができる。本発明の議論の都合上、グローバル・ファイル・システム206は、クラスタ201上に記憶されているすべてのファイル用の単一のグローバル・ファイル空間を維持し、1つのノード202上で実行される。グローバル・ファイル・システム206は装置106の少なくとも2つの表現をサポートする。物理名空間(PNS)表現305は、カーネル空間からアクセスできるものであり、それぞれのノード202上の装置106の物理的構成に対応する。論理名空間(LNS)表現304は物理名空間305のユーザ空間バージョンであり、すなわち論理名空間304内の各エントリは物理名空間305内の対応するエントリにマップされる。本発明は、このグローバル・ファイル・システム206の多数の点を修正し、アプリケーション150による装置106への透過的なグローバル・アクセスを可能にする。クラスタ201は、本発明の重要な構成要素である装置構成システム(DCS)208を備えるノード204も含む。

10

【0037】

他の実施形態には、それぞれ、それ自体の物理名空間および論理名空間を維持する、任意の数のグローバル・ファイル・システム206がある。そのような環境では、特定の装置は、グローバル・ファイル・システム206およびそれに関連付けられた物理名空間および論理名空間のみを通してアクセスされる。

【0038】

上記で図1ないし4を参照して説明したように、従来型のSolaris装置アクセス・システムでは、透過装置アクセスは単一のコンピュータ・システム内でしかできない。従来技術においてアクセスすべき装置のdev_t値にファイル・システムによってマップされた論理名を生成する方法のある態様は、現行の装置アクセス・システムのクラスタへの拡張には適合しない。たとえば、数組の装置106-1、106-2がそれぞれ、4つのSCSIディスク・ドライブを含むと仮定すると、現在使用されている論理命名システムでは、それぞれの異なるノード106-1、106-2上のそれぞれの異なるドライブが同じdev_t値を有する。この場合、アプリケーション150-1がノード202-2上の特定の1つのディスク・ドライブに透過的にアクセスすることは不可能である。次に、本発明がどのように、そのような透過的なグローバル装置アクセスを行うかについて説明する。

20

30

【0039】

図6を参照すると、1つのノード202と、DCS208を備えるノード204の他の詳細が示されている。この図では、ファイル・システム206は、1つの特定のノード202-2上に存在する場合と同様には示されていない。各ノード202は、内部にオペレーティング・システム(OS)ルーチン/オブジェクト240およびデータ構造300が定義されているメモリ230を含む。OSルーチン240は、オペレーティング・システム・カーネル242と、プロキシ・ファイル・システム(PxFS)244と、特殊ファイル・システム258と、装置ドライバ・フレームワーク(DDI)270と、1組のサーバ・オブジェクト(DSO)と、装置ドライバ280とを含む。

【0040】

40

前述のように、カーネル242は、メモリ230、ファイル・システム206、または装置106へのアクセスを求める要求など、アプリケーション150からのシステム呼出しを処理する。カーネル242は、グローバル装置アクセスをサポートするように本発明によって修正されているのでカーネル132(図1)とは異なる。プロキシ・ファイル・システム(PxFS)244は、SolarisPxFSファイル・システムに基づくものであるが、カーネル242と同様に、本明細書ではグローバル装置アクセスをサポートするように修正される。PxFS244は、ノード202-i内のアプリケーション150-iがいくつかの異なるノード202を介してファイル・システム206とシームレスに対話できるようにするオブジェクトの集合を含む。PxFSオブジェクトには、PxFSクラスタ246と、PxFSサーバ248と、f_obj_s(ファイル・オブジェクト)

50

250と、vノード(仮想ファイル・ノード)252と、sノード(特殊ファイル・ノード)254と、px_vノード(プロキシvノード)256とが含まれる。これらのオブジェクトはそれぞれ、ファイル・システム206の動作にตอบสนองしてPxF S 244の必要に応じて作成されるので図6では任意選択(opt)と示されている。

【0041】

DDIフレームワーク270(以下では「DDI」と呼ぶ)も、従来技術(図1)を参照して説明したDDIフレームワーク142に類似している。しかし、DDIフレームワーク270は本発明では、DCS360と対話し、いくつかの異なるノード202上でこれらのノードからアクセスすることのできる装置106との互換性を有する物理名および論理名を生成する。DDI270は、ローカル・ノード202に新しい装置が接続されるたびに呼び出される接続メソッド272を含む。従来技術の接続メソッドとは異なり、接続メソッド272は、DCS360のサービスを使用して、接続された各装置ごとにグローバルに整合する物理名を作成するように構成される。DDIフレームワーク270は、対応する物理名から固有の論理名を生成するリンク・ジェネレータ274の集合も含む。それぞれの異なる種類の装置106ごとにそれぞれの異なる種類のリンク・ジェネレータがある。したがって、接続ルーチン272およびリンク・ジェネレータ274はそれぞれ、装置106をそれぞれ、カーネル・レベルおよびユーザ・レベルでグローバルに見えるようにする物理名空間及び論理名空間を構築する。

【0042】

本発明は、それぞれ、装置106の特定のクラス312を管理する、1組のDSO290をクラスタ200の各ノード上に含む。それぞれの装置クラスは、特定の装置106をオープンすることを求めるユーザの要求を一般的には透過的なグローバル装置アクセス・システムにより、具体的にはDCS372によって満たさなければならない場合の特殊性に対処する本発明の新しい形態である。好ましい実施形態には、dev__enumerat e314、dev__node__specif ic316、dev__global318、dev__nodebound320の4つの装置クラスと、DSO__enum292、DSO__nodespec294、DSO__global296、DSO__nodebound298の4つの対応するDSO290がある。

【0043】

dev__enumerat eクラス314は、各装置が接続されるときに関連するドライバ280によって列挙される特定のノード202で複数のインスタンスを有することのできる装置106(たとえば、複数の記憶装置120)に関連付けられる。dev__node__specif icクラス316は、ノード当たり1つのインスタンスしかなく、その結果ドライバ280によって列挙されることがない装置106(たとえば、カーネル・メモリ116)に関連付けられる。dev__globalクラス318は、各ノード上に存在するドライバを使用してローカルまたはリモートにアクセスできる装置106(たとえば、通信装置118)のクラスである。dev__nodeboundクラスは、特定のノード上のドライバを使用した場合にのみアクセスすることのできる装置(たとえば、HA装置112)に使用される。

【0044】

ドライバ280は、接続中の各オブジェクトについての装置クラス情報312が得られるときにはそれを含む追加の構成情報を報告することを除いてドライバ140に類似している。

【0045】

データ構造300はDevInfoツリー302とddi__minor__nodesテーブル306とを含む。多数のOSルーチン240と同様に、データ構造300は、従来技術(図1)で使用される同様な名称のデータ構造160に類似している。しかし、各データ構造は、本発明の動作を有効にする、従来技術との重要な違いを具体化する。具体的には、DevInfoツリー302は、クラスタ200内の選択されたクラスの装置を見つけるために必要な追加の中間ノードを含む。DevInfoツリーで表される物理名空間

10

20

30

40

50

305を変更したため、論理名空間304も従来技術の論理名空間164とは異なる。最後に、ddi__minor__nodesテーブル306は、従来技術で使用されるddi__minor__nodesテーブルとは異なり追加のフィールドを含む。たとえば、本発明のddi__minor__nodesテーブルは、global__minor__numberフィールド308と、local__minor__numberフィールド310と、(device)classフィールド312(前述)を含む。従来技術のddi__minor__nodesテーブルには、フィールド308とフィールド312のどちらも含まれていない。

【0046】

ノード204は、OSルーチン/オブジェクト340およびデータ構造370が定義されているメモリ330を含む。OSルーチン/オブジェクト340は、装置構成システム(DCS)360と、DCS上のmap__minorメソッドと、前述のDSOと同一の1組のDSO290とを含む。データ構造370はDCSデータベース372を含む。

【0047】

DCS360は、従来技術とはまったく異なるものであり、少なくとも2つの重要な機能を果たす。第1に、DCS360はDDI270と共に働き、新たに接続された装置に、この装置にグローバルにかつ透過的にアクセスできるようにするグローバル・マイナー番号を割り当てる。第2に、DCS360はファイル・システム206と共に働き、PXF S244アプリケーション150が、接続された装置106に透過的にアクセスできるようにする。DCS__database372は、DCS372によって生成されたすべての重要な結果を常に保持する。DCS360のこの2つの態様について以下でそれぞれ、図7Aおよび7Bならびに図8Aおよび8Bを参照して説明する。

【0048】

図7Aを参照すると、ノード202内のDDIフレームワークおよびノード204内のDCS360が、ノード202に接続される装置380の適切なdev__t値、論理名、物理名を規定する動作を示す流れ図が示されている。DDI270、リンク・ジェネレータ274、DCS360、それらの拡張機能は全体として、クラスタ200用の装置レジストラとして働く。動作およびメッセージは図4Aの場合と同様に示されている。流れ図に表した動作を説明する前に、本発明によって管理されるいくつかの名前空間について図7Bを参照して説明する。

【0049】

図7Bを参照すると、2つのノード202-1、202-2を含む例示的なクラスタに本発明で使用されるマイナー名/番号空間307、物理名空間305、論理名空間304の概念図が示されている。後述のように、ノード202に装置106を接続するたびに、装置のドライバが装置にローカル・マイナー番号307__numおよびローカル・マイナー名307__nameを割り当てる。DDI270はこの情報を使用して、グローバルに一意のマイナー番号を生成し、装置106のグローバルに一意の物理名305__nameを形成する。物理名305__nameはクラスタの装置階層内の装置を特定する。次いで、リンク・ジェネレータ274は物理名305__nameをグローバルに一意の論理名304__nameにマップする。DDI270-1、270-2およびリンク・ジェネレータ274-1、274-2がそれぞれ、協働して共通のグローバル物理名空間305およびグローバル論理名空間304を生成することに留意されたい。これとは対照的に、各ドライバはそのノード202に対してマイナー名/番号空間しか生成しない。したがって、本発明はローカルのマイナー名/番号をグローバル物理名および論理名にマップする。これらのグローバル名前空間はファイル・システム206の一部である。したがって、任意のノード202上のアプリケーション150は、ファイル・システム206を使用して、クラスタ200上のすべての装置106を、それらが単一のコンピュータ上に配置されている場合と同様に見てそれらにアクセスすることができる。フレームワークを形成する名前空間について説明したが、次に、本発明について図7Bを参照して説明する。

【0050】

図7Aを参照するとわかるように、ノード220に装置106を接続すると、DDI270はドライバ280に接続メッセージ(7-1a)を発行する。これに対して、ドライバ280は、接続されたインスタンスに関連付けられた各装置のDDI270にcreate_ddi_minor_nodesメッセージ(7-1b)を発行する。create_ddi_minor_nodesメッセージ(7-1b)は、装置380の構成を示し、適切な装置ドライバ280によって割り当てられたローカル・マイナー番号(minor_num)382およびminor_name384と、1つのクラス312から選択されたdevice_class386とを含む。たとえば、装置がノード202に接続された第3のSCSIディスク・ドライブである場合、minor_num、minor_name、classはそれぞれ、「3」、「a」(それがその装置上の第1のスライスであることを示す)、「dev_enumerate」であってよい。

10

【0051】

create_ddi_minor_nodesメッセージ(7-1b)に回答して、DDI270は、local_minor_numフィールド310をminor_num値382(7-2)に等しく設定することによってddi_minor_nodesテーブル380を更新する。DDI270は次いで、DCS360にdc_map_minorメッセージ(7-3)を発行し、装置380の適切なグローバル・マイナー番号388を返すことをDCS360に要求する。前の文の「適切な」の意味は装置クラスに依存する。すなわち、dev_enumerate装置およびdev_nodebound装置は固有のグローバル・マイナー番号388を必要とするが、dev_global装置およびdev_nodespecific装置は必要としない。dc_map_minorメッセージ(7-3)は、(1)「gminor」、すなわちDCS360によって生成されたグローバル・マイナー番号388のリターン・フィールドと、(2)装置ドライバ280によって生成されたローカル・マイナー番号384を保持する「Iminor」と、(3)装置ドライバ280によって生成された装置クラス386を保持する「class」の3つのフィールドを含む。map_minorメッセージ(7-3)に回答して、DCS360は、メッセージ(7-3)で識別されたクラスのローカルDSO290に同様なds_map_minorメッセージ(7-4)を発行する。

20

【0052】

DSO290は特に、接続中の装置に割り当てるべきグローバル・マイナー(gmin)番号388を判定する。gmin番号をどのように割り当てるかは装置のクラス386に依存する。たとえば、列举される各装置は特定のノードでアクセスしなければならないので、dev_enumerateクラス314のDSO292は、各dev_enumerate装置にクラスタ全体にわたって一意のgmin番号388を割り当てる。これに対して、どのノードでdev_global装置にアクセスするかは重要ではないので、dev_globalクラス318のDSO296はこのような装置に同じgmin番号を割り当てる。他のクラスの場合、dev_nodespecificクラス316のDSO294は、そのクラスの各装置に同じ非ヌルgmin番号を割り当て、dev_nodeboundクラス320のDSO298は、そのクラスの各装置にクラスタ全体にわたって一意のgmin番号を割り当てる。

30

40

【0053】

DSO292、298は、まずDCSデータベース372に問い合わせ、グローバル・マイナー番号がまだあるかどうかを判定することによってグローバル・マイナー番号を割り当てる。

【0054】

DCSデータベース372は、常に保持され、クラスタ200内のすべての装置106について、メジャー番号に関するフィールド390と、グローバル・マイナー番号に関するフィールド388と、内部(またはローカル)マイナー番号に関するフィールド382と、(サーバ・クラス386と数値394とを含む)装置サーバidに関するフィールド392とを含む。マイナー名、メジャー番号、グローバル・マイナー番号、ローカル・マイ

50

ナー番号についてはすでに説明した。数値 394 は、接続中の装置のサーバであるノード 202 を識別する。第 1 のクラスのサーバの ID は無関係であり、第 2 のクラスの場合には、装置にアクセスする必要があるノードの位置と同じであるので、数値 394 の情報は、`dev_global` 装置および `dev_nodespec` 装置については任意選択である。DCS データベース 272 の例を表 1 に示す。

【表 1】

装置(ファイルでない)	メジャー 390	グローバル マイナー 388	内部 マイナー 382	装置サーバ 392:	
				サーバ クラス 386	数値 394
tcp	42	0	0	dev_global	0
kmem	13	12	12	dev_node_ spec	0,
disk c2l0d0s0	32	24	24	dev_enum	node id
kmem	13	1	12	dev_enum	node 0 id
kmem	13	2	12	dev_enum	node 1 id
kmem	13	3	12	dev_enum	node 2 id
kmem	13	4	12	dev_enum	node 3 id
HA devices	M	X1	X1	dev_nodeb ound	id

【0055】

表 1 の第 1 行は、tcp インタフェースに関するエントリを示す。tcp インタフェースは、クラスタ 200 内のあらゆるノード 202 からアクセスできるので `dev_global` 装置である。tcp 装置のメジャー番号は 42 であり、この値はすべての tcp ドライバに関連付けられる。tcp 装置のグローバル最小値 388 およびローカル最小値 382 ならびにサーバ数値 394 (すなわち、`node_id`) は 0 に設定されることに留意されたい。これは、どのノードから tcp インタフェースにアクセスするかは重要ではないからである。したがって、DCS データベースにはクラスタ 200 全体に関する tcp エントリは 1 つしかない。表 1 内の第 2 のエントリは、省略時指定ではローカルにアクセスされるカーネル・メモリ装置に関するエントリである。このため、これは `dev_nodespec` クラスの装置である。メジャー番号 13 は、kmem 装置ドライバに関連付けられる。kmem 装置は、特定のサーバでアクセスされるわけではないので、ヌル数値 394 を有し、かつ同一の非ヌルグローバル最小値および非ヌルローカル最小値 (12) を有する。これは、`dev_nodespec` 装置では、DCS 360 がローカル・マイナー番号と同じグローバル・マイナー番号を割り当てるに過ぎないからである。この例では、それぞれのノード 202 上に配置された kmem 装置どうしを区別する必要がないので DCS データベース 372 には `dev_nodespec` クラスの kmem エントリは 1 つしかない。

【0056】

第 3 のエントリは、SCSI ドライバがメジャー番号 32 を有する SCSI ディスク c0

t 0 d 0 t 0 に関するエントリである。D C S 3 6 0 は、S C S I 装置にそのローカル・マイナー番号 3 8 2 (2 4) と同じグローバル・マイナー番号 3 8 8 を割り当てている。これは、D C S データベース 3 7 2 には他の S C S I 装置が表されていないからである。しかし、他の S C S I 装置 c 0 t 0 d 0 t 0 を同じローカル番号 (2 4) で異なるノードに登録した場合、D C S 3 6 0 はその S C S I に異なるグローバル番号、おそらく 2 5 を割り当てる。同じローカル番号を有する S C S I 装置どうしを区別するために、D C S データベース 3 7 2 は完全なサーバ情報を含む。この場合、数値 3 9 4 はサーバ 2 0 2 のホスト i d に設定される。

【 0 0 5 7 】

エントリ 4 ないし 7 は、d e v _ e n u m e r a t e 装置として登録された 4 つのカーネル・メモリ装置に関するエントリである。好ましい実施形態では、d e v _ n o d e s p e c i f i c 装置を登録するたびに、カーネル内のすべてのノード 2 0 2 について、D C S データベース 3 7 2 内に追加のエントリを作成することができ、これによって、ユーザはローカル・ノード以外のノード上の d e v _ n o d e s p e c i f i c 装置にアクセスすることができる。したがって、4 つのノード 2 0 2 - 1、2 0 2 - 2、2 0 2 - 3、2 0 2 - 4 があると仮定すると、D C S 2 6 0 はこれらのノードのそれぞれごとに d e v _ e n u m e r a t e クラスのカーネル・メモリ装置を登録することができる。他の d e v _ e n u m e r a t e 装置の場合と同様に、各 k m e m 装置には固有のグローバル番号が割り当てられる。d e v _ e n u m e r a t e 情報は、カーネル・メモリ装置をオープンすることを求める一般的な要求 (たとえば、o p e n (/ d e v i c e s / k m e m)) をユーザが発行するときには使用されない。d e v _ e n u m e r a t e 情報は、カーネル・メモリ装置をオープンすることを求める特定の要求をユーザが発行するとき使用される。たとえば、要求 o p e n (/ d e v i c e s / k m e m 0) によって、ユーザはノード 0 上の k m e m 装置をオープンすることができる。

【 0 0 5 8 】

最後のエントリは、一般的な高可用性 (H A) 装置が D C S データベース 3 7 2 内でどのように表されるかを示す。メジャー番号 3 9 0、グローバル・マイナー番号、ローカル・マイナー番号は、m a p _ m i n o r _ n o d e s メッセージに与えられる値 M、X 1、X 1 から取り出される。数値 3 9 4 は、特定のノードに結合された装置の i d に設定される。この「i d」はノード i d ではない。この i d は、各 H A サービスごとにクラスタ 2 0 0 に対して固有に作成される。

【 0 0 5 9 】

装置 3 8 0 のグローバル・マイナー番号 3 8 8 が判定された後、適切な D S O 2 9 0 が、D C S データベース 3 7 2 を新しい情報で更新し (7 - 5)、グローバル・マイナー番号 3 8 8 を D C S 3 6 0 に返す (7 - 6)。D C S 3 7 2 は次いで、グローバル・マイナー番号 3 8 8 を D D I 2 7 0 に返し (7 - 7)、D D I 2 7 0 が d d i _ m i n o r _ n o d e s テーブル 3 0 6 (7 - 9)、論理名空間 3 0 4、物理名空間 3 0 5、d e v _ i n f o ツリー 3 0 2 (7 - 9) を更新する。D D I 2 7 0 は、新しいグローバル・マイナー番号 3 8 8 を書き込むことによって d d i _ m i n o r _ n o d e s テーブル 3 0 6 を更新する。名前空間 3 0 4 / 3 0 5 の更新はこれよりも複雑であり、この更新について次に説明する。

【 0 0 6 0 】

第 1 に、D D I 2 7 0 は D e v I n f o ツリー 3 0 2 に新しいリーフ・ノードを追加し、D e v I n f o ツリーの構造が、上記で図 3 を参照して説明した構造から、d e v _ e n u m e r a t e が接続されたクラスタ・サイトを表す追加の「/ h o s t i d」ノード・レベルを「/ d e v i c e s」ノードのすぐ下に含めるように変更される。各ノード 2 0 2 が、そのノード上の装置を表すそれ自体の D e v I n f o ツリー 3 0 2 を有することに留意されたい。しかし、物理名空間で表されるように、D e v I n f o ツリーの集合は、この追加の / h o s t i d ノードを含む単一の表現として組み合わせられる。(たとえば、典型的な物理名は文字列 / d e v i c e s / h o s t i d / . . . から始まる)。各装

10

20

30

40

50

置はまた、リーフ・レベルで、ローカル・マイナー番号 382 ではなくグローバル・マイナー番号 388 に関連付けられる。(dev__enumerate 装置の) 必要に応じて、DevInfo ツリー 302 の各リーフ・ノードの dev__t 値は、対応する装置のグローバル・マイナー番号 388 およびそのドライバのメジャー番号 390 から導かれる。たとえば、グローバル・マイナー番号 GN と、マイナー名 MN と、ドライバ sd@addr/sbus@addr/esp@addr/sd@addr:MN

【0061】

この物理名は、メジャー番号およびグローバル・マイナー番号から導かれた dev__t 値を属性として含む、所与の装置に関する構成情報を含む UFS ファイル 170 (図 2B) の物理名に対応する。

【0062】

本発明のリンク・ジェネレータ 274 は、DevInfo パスの少なくとも一部と、DCS から返されたグローバル・マイナー番号に従って修正されたドライバから与えられたマイナー名とから装置 (および対応する UFS) の論理名を導く。

【0063】

たとえば、ノード 202-1 が、最初にドライバ・マイナー名 a ないし d およびドライバ・マイナー番号 0 ないし 3 で割り当てられた 4 つのスライスを含む 1 枚の SCSI ディスクを有し、ノード 202-2 が、マイナー名 a ないし f およびマイナー番号 0 ないし 5 で割り当てられた 6 つのスライスを含む 1 枚の SCSI ディスクを有すると仮定する。これらの装置を接続したときに、DCS 360 が第 1 の SCSI ディスクに関してはグローバル・マイナー番号 0 ないし 3 を返し、第 2 の SCSI ディスクに関してはグローバル・マイナー番号 4 ないし 9 を返したと仮定する。DDI 270 は、これらのグローバル・マイナー番号を使用して物理名 (後述) を作成し、リンク・ジェネレータ 274 は DDI 270 を使用して、以下のように物理名にマップされる論理名を作成する。

ドライバ 280 からのマイナー名	リンク・ジェネレータ 274 からの論理名
a (ノード 202-1)	/dev/dsk/c0t0d0s0
b "	/dev/dsk/c0t0d0s1
c "	/dev/dsk/c0t0d0s2
d "	/dev/dsk/c0t0d0s3
a (ノード 202-2)	/dev/dsk/c1t0d0s0
b "	/dev/dsk/c1t0d0s1
...	
f "	/dev/dsk/c1t0d0s5

【0064】

ノード 202-1 および 202-2 に割り当てられた論理名はそれぞれの異なるクラスタ値 (論理名文字列 cxt0d0sy の cx 部分。この場合、「x」および「y」は変数である) を有する。これは、論理名が装置の物理名にマップされ、クラスタ内で、それぞれの異なるノード上の装置がそれぞれの異なるコントローラに関連付けられるからである。たとえば、ノード 202-1 コントローラは c0 として表され、ノード 202-2 は c1 として表される。

【0065】

DDI 270 は、同じ gmin 情報を使用して物理名空間 305 を生成し、対応する装置の dev__t 値が属性として含まれるファイルを識別する論理名と物理名との間のマップを生成する。上記の例では、論理名空間 304 と、論理名空間と物理名空間とのマップは以下のように更新される (addr が任意のアドレスを代替することに留意されたい)。

10

20

30

40

論理名	DevInfoツリー302からの物理名
/dev/dsk/c0t0d0s0	/devices/node_202-1/iommu@addr/sbus@addr/espl@addr/sd@0:a
/dev/dsk/c0t0d0s1	" /espl@addr/sd@0:b
/dev/dsk/c0t0d0s2	" /espl@addr/sd@0:c
/dev/dsk/c0t0d0s3	" /espl@addr/sd@0:d
/dev/dsk/clt0d0s0	/devices/node_202-2/iommu@addr/sbus@addr/espl@addr/sd@0:minor
/dev/dsk/clt0d0s1	" /espl@addr/sd@0:e
/dev/dsk/clt0d0s2	" /espl@addr/sd@0:f
...	
/dev/dsk/c0t0d0s5	" /espl@addr/sd@0:i

10

【0066】

この例は、DDI270が、dev__enumerate装置、たとえばSCSI装置がそのクラスのメンバである論理名および物理名を生成することを示している。簡単に言えば、dev__enumerate装置を命名する規則では、特定のドライバ（たとえば、sd）によって列挙される各インスタンスが固有のグローバル・マイナー番号を有し、そのグローバル・マイナー番号をドライバのメジャー番号と組み合わせたときに、対応する固有のdev__t値が形成される必要がある。このような規則は、各インスタンスに関連付けられた物理名が、そのインスタンスのホストidとグローバル・マイナー番号を、他の従来の物理パス情報と共に含まなければならないことも指定する。他のクラスの他の装置を命名する規則は、上記でdev__enumerateクラスに関して説明した規則に類似している。

20

【0067】

特に、DDI270は、dev__nodespecific装置に形式/dev/dev__ice__nameの論理名および以下の形式の物理名を割り当てる。/devices/pseudo/driver@gmin:device__name
この場合、device__nameは名前384であり、pseudoは、この種の装置が擬似装置であることを示し、driverは対応するドライバのidであり、@gmin:device__nameは、dev__nodespecific装置のグローバル番号388および装置名384を示す。たとえば、カーネル・メモリ装置の論理名および物理名はそれぞれ、/dev/kmemおよびdevices/pseudo/mm@12:kmemであってよい。前述のように、kmem装置には、それに特定のノード上でアクセスできるようにする論理名を与えることができる。たとえば、DDI270は、論理名/dev/kmem0を物理名/devices/hostid0/pseudo/mm@0:kmemにマップする。

30

40

【0068】

dev__globalクラスの場合、DDIによって生成される各論理名は、ファイル・システムによってクラスタ200内の任意の装置に関連付けられる共通の物理パスを識別する。このような装置の論理名は、/dev/device__nameの形式であり、以下の形式の物理名にマップされる。

/devices/pseudo/clone@gmin:device__name

この場合、device__nameはドライバに特有の名前384であり、pseudoは、この種の装置が擬似装置であることを示し、cloneはこの装置がクローン可能で

50

あることを示し、@gmin:device_nameは、dev_global装置のグローバル番号388および装置名384を示す。たとえば、表1のtcp装置の論理名は/dev/tcpであり、物理名は/devices/pseudo/clone@0:tcpである。前述のように、本発明では、kmem装置の場合とは異なり、dev_global装置どうしを区別することができないことに留意されたい。すなわち、すべてのdev_global装置が区別不能である。

【0069】

本発明のクラス・ベースの命名システムの利点は、それが、Solarisの従来のバージョン向けに設計された従来型のソフトウェアとの互換性を有することである。たとえば、従来型のプログラムはopen(/dev/kmem)要求を発行することができ、この場合、本発明を具体化するSolarisのバージョンは、ローカルkmem装置のハンドルを返す。本発明は、dev_global装置とdev_enumerate装置に対して同様な結果を与える。従来技術には、dev_nodetound装置に関する概念はなかった。

【0070】

それぞれの異なるクラスの装置がクラスタ200のそれぞれの異なるノード上でアクセスできる一貫したグローバル名空間を、DDI270およびDCS360がどのように形成するかについて説明したが、次に、他のノード上の装置に対するオープン要求に応答するために本発明によって使用されるステップについて図8Aおよび8Bを参照して説明する。

【0071】

図8Aおよび8Bを参照すると、ノード202-3上に存在する装置106-2(図8B)にアクセス(オープン)するためにノード202-1上で実行されているアプリケーション150からの要求(8-1)に応答して本発明によって実行されるステップの流れ図が示されている。この例では、ファイル・システム206およびDCS360はそれぞれ、ノード202-2および204上に存在する。アプリケーション150は装置の論理名上のカーネル242にオープン要求を発行する。カーネル242は次いで、ファイル・システム206に問い合わせ装置のdev_t値を判定する。ファイル・システムがカーネル242とは異なるノード上にあるので、これは、プロキシ・ファイル・システムPxF Sを使用するマルチステップ・プロセスであり、プロキシ・ファイル・システムPxF Sの大部分の態様はすでにSolarisの現行のバージョンによって定義されている。しかし、本発明は、Solarisの従来のバージョンとはまったく異なり、PxF Sクライアント246やPxF Sサーバ248などのプロキシ・ファイル・システム要素をDCS360との対話をサポートするように修正する。PxF Sクライアント246とPxF Sサーバ248とファイル・システム206との間の対話について簡単に説明する。

【0072】

まずファイル・システム206にアクセスする必要のあるカーネル242などのオブジェクトは、そのローカルPxF Sクライアント246にアクセス要求を発行する。PxF Sクライアントは、ファイル・システム206と共に配置されたPxF Sサーバ248の参照を保持する。この参照によって、PxF Sクライアント246はカーネルの要求をPxF Sサーバ248を介してファイル・システム206に伝達することができる。ファイル・システム206は、要求されたアクセスを実行し、要求されたファイルを表すvノード・オブジェクト252を作成し、vノード・オブジェクト252の参照をPxF Sサーバ248に返す。ノード202-1および202-2がそれぞれの異なるアドレス空間であるので、vノード252の参照は、ノード202-1内のPxF Sクライアント246およびカーネル242には無用である。したがって、PxF Sサーバ248は、vノード252にリンクされたファイル転送オブジェクト(f_obj)250を作成し、f_obj250の参照をPxF Sクライアント246に返す。PxF Sクライアント246は、f_obj参照を受信した後、f_obj250にリンクされたプロキシvノード(px_vnode)256を作成する。カーネル242は次いで、単にローカルpx_vノード

10

20

30

40

50

ド 2 5 6 にアクセスすることによって、v ノード 2 5 2 で表されたファイル情報にアクセスすることができる。

【 0 0 7 3 】

カーネル 2 4 2 は、この機構を使用して、P x F S クライアント 2 4 6 に対してオープンされる装置の論理名に関する参照メッセージ (8 - 2) を発行し、P x F S クライアント 2 4 6 は同様な参照メッセージ (8 - 3) を P x F S サーバ 2 4 8 へ送る。P x F S サーバ 2 4 8 は、ファイル・システム 2 0 6 に `lookup (logical_name)` , `get_vnode` メッセージ (8 - 4) を発行し、論理シンボリック・リンクを介して `logical_name` を対応する `physical_name` にマップし、その `physical_name` で識別される U F S ファイルを表す `v_node 2 5 2` の参照を返すことを要求する。 `physical_name` がこの例と同様に装置を参照するものであるとき、装置の属性は装置の固有の `dev_t` を含む。前述のように、ファイル・システム 2 0 6 は次いで、P x F S サーバ 2 4 8 に v ノードを返し (8 - 5) 、P x F S サーバ 2 4 8 が対応する `f_obj 2 5 0` を作成し、`f_obj 2 5 0` の参照を P x F S クライアント 2 4 6 に返す (8 - 6) 。P x F S クライアント 2 4 6 は次いで、要求された装置に関する `dev_t` 情報が属性に含まれる `px_vnode 2 5 6` を作成し、その `px_vnode 2 5 6` の参照をカーネル 2 4 2 に渡す (8 - 7) 。この点で、カーネル 2 4 2 は `px_vnode 2 5 6` の P x F S クライアント 2 4 6 にオープン・メッセージ (8 - 8) を返す。P x F S クライアント 2 4 6 は、このメッセージを受信した後、`dev_t` 値を含む `px_vnode` の属性から、対応する v ノード 2 5 2 が装置を表し、従って、D C S 3 6 0 によってオープン・メッセージを処理しなければならないと判定する。 `px_vnode 2 5 6` が `dev_t` 値を含まない場合、P x F S クライアント 2 4 6 は他のチャネルを通してオープン要求 (8 - 8) を満たす。S o l a r i s の従来のバージョンの場合と同様に、装置はローカルにしかアクセスできないので、P x F S クライアントは `dev_t` 値に関する試験を実行しない。

【 0 0 7 4 】

`px_vnode 2 5 6` が `dev_t` 値 4 3 0 を含むので、P x F S クライアント 2 4 6 は、`dev_t` に対応する装置の D C S 3 6 0 に `resolve` メッセージ (8 - 9) を発行する。次に、D C S 3 6 0 がこの要求をどう扱うかについて図 8 B を参照して説明する。

【 0 0 7 5 】

図 8 B を参照するとわかるように、`resolve (dev_t)` メッセージ (8 - 9) に応答して、D C S 3 6 0 は D C S データベース 3 7 2 を参照し、`dev_t` 値に対応する装置の位置および I D を判定する。装置クラス 3 1 2 に関する前述の議論と同様に、`dev_enumerate` クラスおよび `dev_nodebound` クラスの装置は、D C S データベース 3 7 2 の数値フィールド 3 9 4 にその位置が指定された特定のノード上でアクセスされる。これに対して、`dev_global` クラスまたは `dev_node_specific` クラスの装置は、要求側アプリケーションのローカル・ノード上でアクセスされる。D C S 3 6 0 は、オープンすべき装置の位置を判定した後、要求された装置が属する装置クラスを管理し、要求されたオブジェクトを備えるノードに対してローカルな D S O 2 9 0 の参照 (`DSO_ref`) を P x F S クライアント 2 4 6 に返す (8 - 1 0) 。この例では、要求された装置 1 0 6 - 2 が `dev_enumerate` クラスの装置であり、ノード 2 0 2 - 3 上に存在すると仮定した場合、`DSO_ref` はノード 2 0 2 - 3 上の `DSO_enum` オブジェクト 2 9 2 に返される。

【 0 0 7 6 】

P x F S クライアント 2 4 6 は、メッセージ (8 - 1 0) を受信した後、装置 1 0 6 - 2 に関する `get_device_obj` 要求を、参照された D S O 2 9 2 (8 - 1 1) に発行する。これに応答して、D S O 2 9 2 は `create_specvp ()` メッセージ (8 - 1 2) を発行し、ノード 2 0 2 - 3 上の `SpecFS 4 1 0` に装置 1 0 6 - 2 の s ノードを作成し返す (8 - 1 3) ように要求する。D S O 2 9 2 は次いで、P x F S サ

10

20

30

40

50

サーバ248-2にsノードのf__obj参照を要求し(8-14a)、PxFSSサーバ248-2は、要求されたf__objを返す(8-14b)。DSO292は次いで、sノードのfobj参照をPxFSSクライアント246に返す(8-15)。クライアント246は次いで、このfobjに関するオープン要求(8-16)を発行し、この要求がPxFSSサーバ248-2を介してSpecFS410へ送られる(8-17)。

【0077】

SpecFS410は次いで、装置106-2をオープンすることを試みる。オープン動作の結果に応じて、SpecFS410は成功と失敗のどちらかを示す状況メッセージ(8-18)を返す。オープンが成功した場合、状況メッセージ(8-18)が、オープンされたsノード432の参照も含む。PxFSSサーバ248-2は、状況メッセージ(8-18)内の「成功」を受信した後、オープンされたv__node252-2に関するf__obj250-2を作成し、PxFSSクライアント246に返す(8-19)。PxFSSクライアント246は、ノードを介してf__obj250-2にリンクされたpx__vnode256-2を作成する。装置オープン動作の最後のステップとして、PxFSSクライアントはpx__vnode256-2をカーネル242に返し(8-20)、カーネル242が、対応するユーザ空間ファイル記述子(fd)434を作成する。カーネル242はこのファイル記述子をアプリケーション150-1に返し(8-21)、次いでアプリケーション150-1は、ファイル記述子434を使用して装置106-2と直接(すなわち、カーネル242、PxFSSクライアント246、px__nodeを介して)対話することができる。

【0078】

いくつかの特定の実施形態を参照して本発明について説明したが、この説明は本発明を例示するものであり、本発明を制限するものと解釈すべきではない。当業者には、添付の特許請求の範囲で定義される本発明の真の趣旨および範囲から逸脱せずに様々な修正が構想されよう。

【図面の簡単な説明】

【図1】単一のコンピュータ上の装置へのアクセスを可能にするために使用される構成要素を示す従来技術のコンピュータ・システムのブロック図である。

【図2】従来技術におけるアプリケーションとオペレーティング・システム・カーネルとファイル・システムと装置との間の関係を示すブロック図(A)と従来技術における装置論理名と物理名とファイル・システムと装置タイプ識別子(dev__t)と装置との間の関係を示すブロック図(B)である。

【図3】従来技術で使用される装置情報ツリー(DevInfo Tree)に整合する例示的な装置情報ツリーの図である。

【図4】アプリケーション150の要求に応じて装置をオープンする際に従来技術のコンピュータ・システム100のメモリ104で実行される動作の流れ図である。

【図5】本発明を実施することのできるコンピュータ・クラスタのブロック図である。

【図6】図5のクラスタの代表的なノード202および204で実施された本発明を構成するメモリ・プログラムおよびデータ構造のブロック図である。

【図7A】装置ドライバ・インタフェース(DDI)フレームワークおよび装置構成システム(DCS)が、ノード202に接続中の装置の適切なdev__t値、論理名、物理名を規定する動作を示す流れ図である。

【図7B】本発明によって規定されたローカル・マイナー名/番号と物理名と論理名との間の関係を示す図である。

【図8A】ノード202-3上に存在する装置にアクセス(オープン)するためにノード202-1上で実行されているアプリケーション150からの要求に応答して本発明によって実行される各ステップを示す流れ図である。

【図8B】ノード202-3上に存在する装置にアクセス(オープン)するためにノード202-1上で実行されているアプリケーション150からの要求に応答して本発明によって実行される各ステップを示す流れ図である。

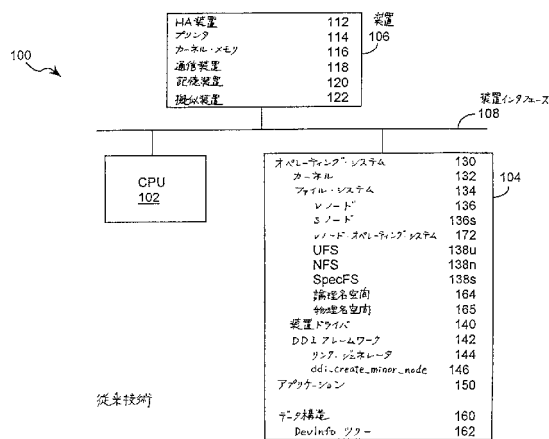
【符号の説明】

- 150 アプリケーション
 200 クラスタ
 201 クラスタ
 202 ノード
 204 ノード
 206 グローバル・ファイル・システム
 208 装置構成システム
 240 オペレーティング・システム（OS）ルーチン／オブジェクト
 242 カーネル
 244 プロキシ・ファイル・システム（P×FS）
 250 ファイル・オブジェクト
 252 仮想ファイル・ノード
 254 特殊ファイル・ノード
 256 プロキシ・ノード
 258 特殊ファイル・システム
 270 装置ドライバ・フレームワーク（DDI）
 272 接続メソッド
 274 リンク・ジェネレータ
 280 装置ドライバ

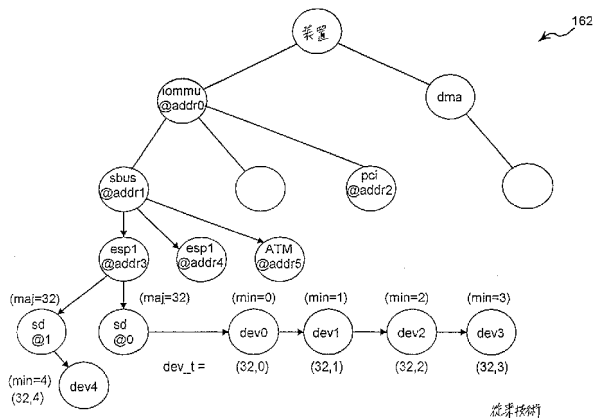
10

20

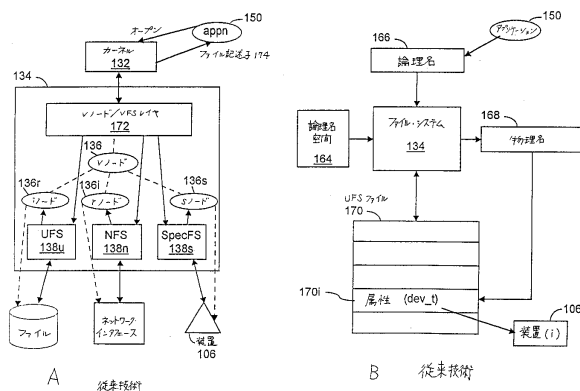
【図1】



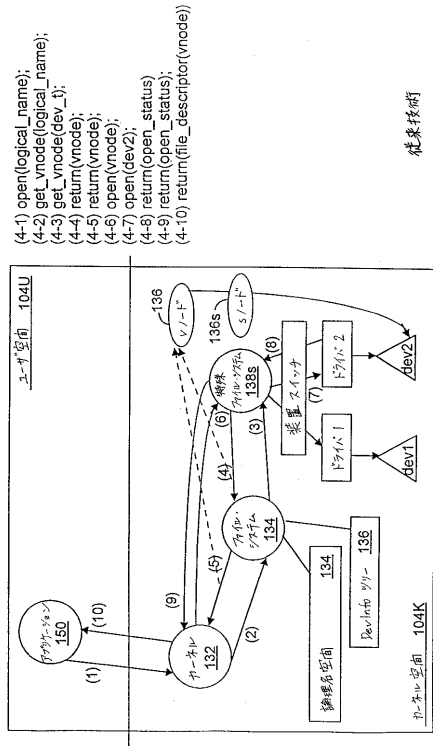
【図3】



【図2】

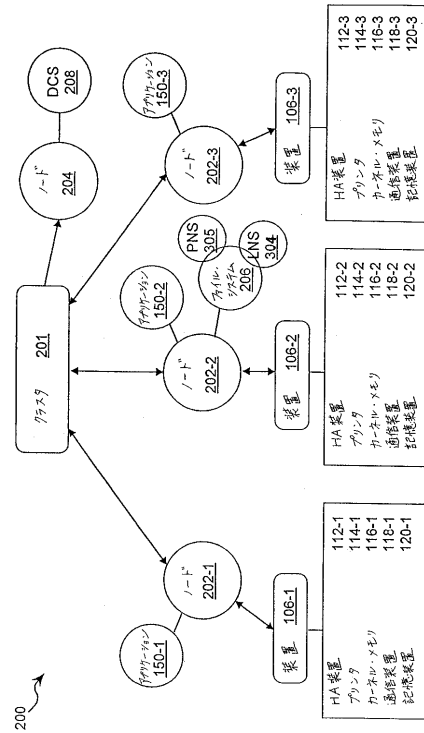


【 図 4 】

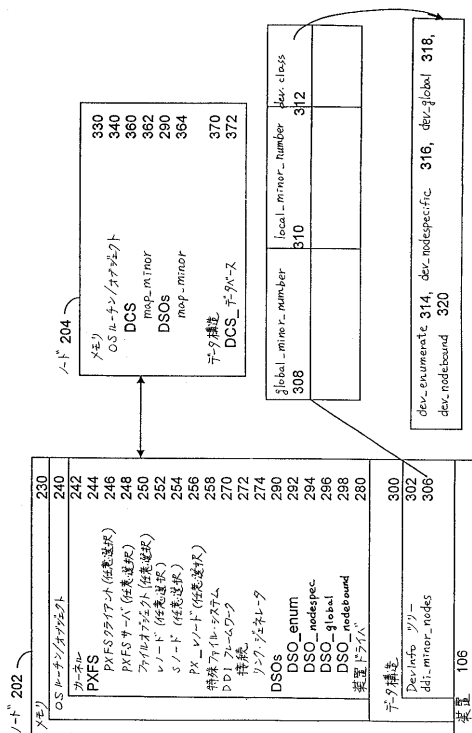


従来技術

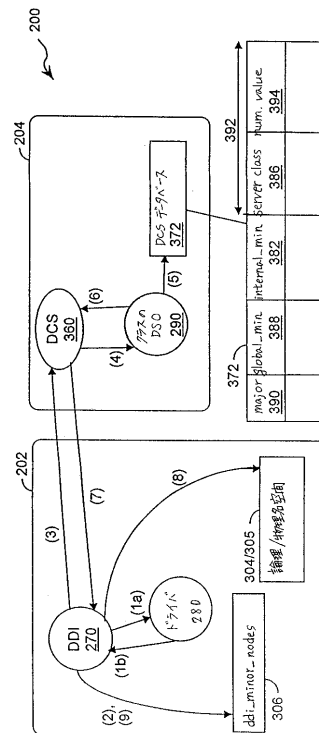
【 図 5 】



【 図 6 】



【 図 7 A 】

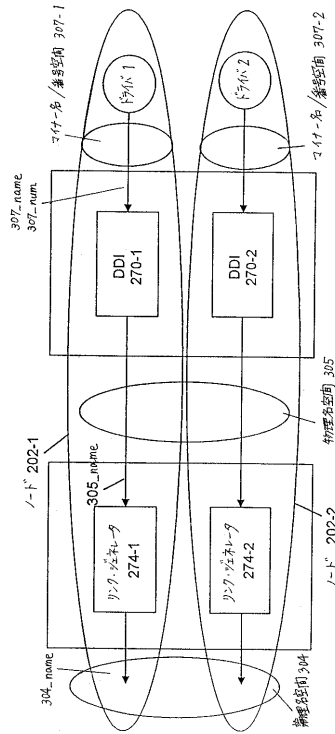


```

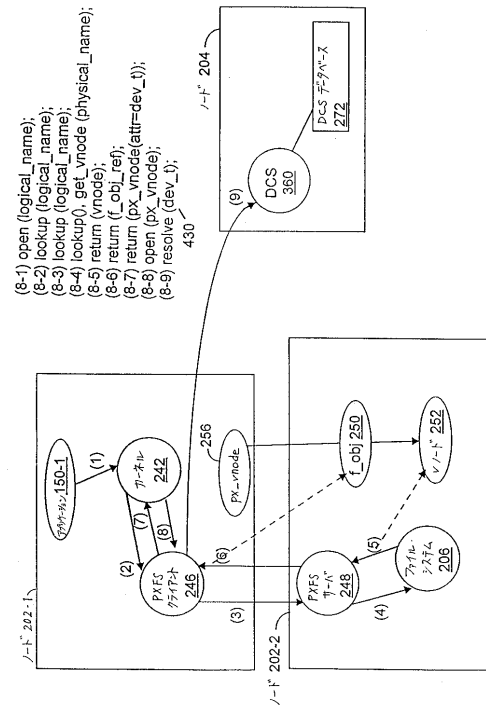
(7-1) attach();
(7-1b) ddl_create_minor_nodes (minor_num 382, minor_name 384, class 386);
(7-2) (7-5) update ddl_minor_nodes (gminor, lminor, class);
(7-3) dc_map_minor (major, lminor, gminor, class);
(7-4) ds_map_minor (major, lminor, gminor);
(7-5) return (gminor);
(7-6) update DCS_database (major, lminor, gminor, DS_type, DS_num_val);
(7-7) return (gminor);
(7-8) update filesystem;

```

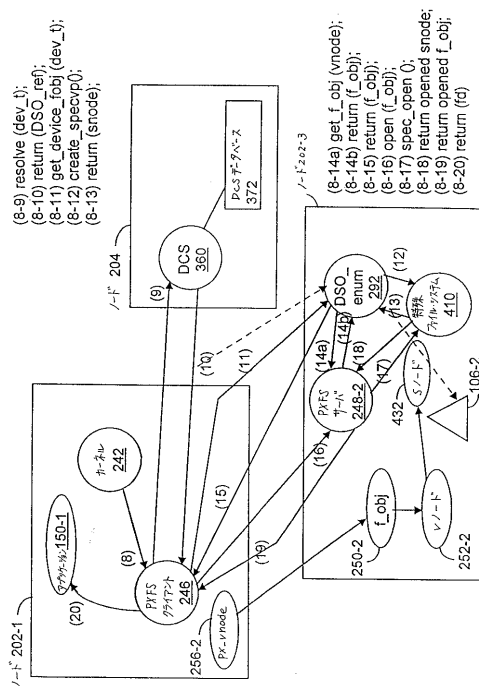
【 図 7 B 】



【 図 8 A 】



【 図 8 B 】



フロントページの続き

(72)発明者 シアマク・ナザリ

アメリカ合衆国・9 1 0 0 6・カリフォルニア州・アルカディア・オークウッド アヴェニュー・1
8 6 2

(72)発明者 アニル・スワループ

アメリカ合衆国・9 2 3 5 4・カリフォルニア州・ロマ リンダ・シエロ レーン・1 1 5 6 2

(72)発明者 ヨーセフ・カリディ

アメリカ合衆国・9 4 0 8 6・カリフォルニア州・サニーバイル・ウエスト ガーランド テラス
・6 3 3

審査官 石井 茂和

(56)参考文献 米国特許第0 5 5 5 5 4 0 1 (U S , A)

Brent Welch , A Comparison of Three Distrebuted File System Architectures: Vnode, Sprit
e, and Plan9 , COMPUTING SYSTEMS, USENIX ASSOCIATION , 米国 , the MIT Press , 1 9 9 4年
7月 8日 , VOL.7 No.2 , p175-199

Greg Thiel , LOCUS operating system, a transparent system , COMPUTER COMMUNICATIONS , NL
, ELSEVIER SCIENCE PUBLISHERS, BV. , 1 9 9 1年 7月 1日 , vol.14, no.6 , p336-346

Jim-Min Lin & Shang Rong Tsai , Supporting user mobility in the Distributed Logical Mac
hine System , MICROPROCESSING and MICROPROGRAMMING , NL , ELSEVIER SCIENCE PUBLISHERS, BV
, 1 9 9 5年 8月 1日 , vol.41, no.4 , p315-330

(58)調査した分野(Int.Cl. , D B 名)

G06F 13/14

WPI(DIALOG)

JSTPlus(JDream2)