

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第5022412号
(P5022412)

(45) 発行日 平成24年9月12日(2012.9.12)

(24) 登録日 平成24年6月22日(2012.6.22)

(51) Int. Cl. F I
H O 4 L 12/56 (2006.01) H O 4 L 12/56 I O O Z

請求項の数 10 (全 25 頁)

| | | | |
|-----------|------------------------------|-----------|-----------------------------------|
| (21) 出願番号 | 特願2009-180328 (P2009-180328) | (73) 特許権者 | 000004226 |
| (22) 出願日 | 平成21年8月3日(2009.8.3) | | 日本電信電話株式会社 |
| (65) 公開番号 | 特開2011-35686 (P2011-35686A) | | 東京都千代田区大手町二丁目3番1号 |
| (43) 公開日 | 平成23年2月17日(2011.2.17) | (74) 代理人 | 100064414 |
| 審査請求日 | 平成23年8月1日(2011.8.1) | | 弁理士 磯野 道造 |
| | | (74) 代理人 | 100127720 |
| | | | 弁理士 大石 恵 |
| | | (74) 代理人 | 100162374 |
| | | | 弁理士 中村 新二 |
| | | (72) 発明者 | クリステル ベルサー |
| | | | 東京都千代田区大手町二丁目3番1号 日 本電信電話株式会社内 |
| | | (72) 発明者 | 増田 暁生 |
| | | | 東京都千代田区大手町二丁目3番1号 日 本電信電話株式会社内 |

最終頁に続く

(54) 【発明の名称】 経路情報管理システム、経路情報管理方法、およびプログラム

(57) 【特許請求の範囲】

【請求項1】

I P (Internet Protocol) ネットワークを構成する A S (Autonomous System) 内で、他の A S から該他の A S の I P アドレスの集合であるプリフィクスを含む外部経路情報を取得し、かつパケットを他の A S に転送する複数のルータと、同一のプリフィクスを含む外部経路情報を前記ルータから収集して経路計算を実行する 2 以上の経路情報管理装置とが相互に通信可能に接続される経路情報管理システムであって、

前記経路情報管理装置は、

前記ルータから収集した前記プリフィクスごとに、そのプリフィクスを含む外部経路情報を用いて、前記 A S 内の 2 以上のルータの一まとまりを単位として前記 A S 内に複数設定されるルータ集合に属する少なくともいずれか 1 つのルータにおける、パケットの転送先のプリフィクスに至る転送経路を計算する経路計算を実行する経路計算部と、

前記経路計算部によって経路計算されたルータ集合ごとの経路計算の結果を含む経路計算結果情報を送信する当該ルータ集合に属する所定数のルータを選択する格納先選択部と

、
当該経路計算結果情報を、前記格納先選択部によって選択された当該所定数のルータに送信する入出力部と、

を備え、

前記ルータは、

転送すべきパケットを受信したときに、前記経路情報管理装置から受信した前記経路計

算結果情報を参照して、該パケットの転送先を決定する転送経路決定部と、

前記転送経路決定部によって決定された転送先に前記転送すべきパケットを送信する入出力部と、
を備える

ことを特徴とする経路情報管理システム。

【請求項 2】

前記経路情報管理装置は、さらに、

前記ルータ集合ごとに、少なくとも 1 以上のルータのルータ ID として、当該ルータの属する前記ルータ集合のルータ集合 ID のビット列の後にハッシュ関数によって生成されたハッシュ値のビット列を連結したビット列を記憶する格納先ノード情報を格納する記憶部と、

10

前記ルータ集合ごとに、前記経路計算結果情報の格納先の検索に用いる検索キーとして、当該ルータ集合のルータ集合 ID のビット列、前記プリフィクスのビット列、前記プリフィクスに対するマスク値のビット列、およびパディングとしてすべて「0」のビット列を連結したビット列を生成する検索キー生成部と、
を備え、

前記格納先選択部は、

前記格納先ノード情報に記憶しているルータ ID と前記検索キーとの距離を、当該ルータ ID と前記検索キーとのビットごとの排他的論理和を示すビット列において、左から見て最初に「1」が出現したビット位置について右から数えたビットの桁数の数値として算出し、その算出した距離の小さい順に所定数のルータを当該ルータ集合から送信先として選択する

20

ことを特徴とする請求項 1 に記載の経路情報管理システム。

【請求項 3】

前記経路計算結果情報は、前記検索キーと、その検索キーに含まれるプリフィクスに対応する前記ルータ集合ごとの経路計算の結果とを関連付けた情報である

ことを特徴とする請求項 2 に記載の経路情報管理システム。

【請求項 4】

前記ルータは、さらに、

そのルータの属するルータ集合内の自身を除く少なくとも 1 以上のルータのルータ ID を記憶する取得先ノード情報と前記経路計算結果情報とを記憶する記憶部と、

30

転送すべきパケットを受信したとき、当該ルータの属するルータ集合 ID のビット列、当該パケットの宛先 IP アドレスのビット列、前記宛先 IP アドレスに対するマスク値としてすべて「1」のビット列、およびパディングとしてすべて「0」のビット列を連結して形成される、取得キーを生成する取得キー生成部と、

前記取得先ノード情報に記憶しているルータ ID と前記取得キーとの距離を、前記ルータ ID と前記取得キーとのビットごとの排他的論理和を示すビット列において、左から見て最初に「1」が出現したビット位置について右から数えたビットの桁数の数値として算出し、その算出した距離の小さい順に所定数のルータを当該ルータ集合から取得先として選択する取得先選択部と、

40

を備え、

前記転送経路決定部は、

自身の前記記憶部に記憶する前記経路計算結果情報の検索キーと前記取得キーとが一致するか否かを判定し、一致するものがある場合には、その検索キーに関連付けられた前記ルータ集合ごとの経路計算の結果に基づいてパケットの転送先を決定し、一致するものがない場合には、前記取得先選択部によって取得先として選択されたルータに前記取得キーを送信して、当該ルータの記憶する前記経路計算結果情報の検索キーと前記取得キーとが一致する前記ルータ集合ごとの経路計算の結果を取得して、その取得した前記ルータ集合ごとの経路計算の結果に基づいてパケットの転送先を決定する

ことを特徴とする請求項 3 に記載の経路情報管理システム。

50

【請求項 5】

前記転送経路決定部は、さらに、前記経路計算結果情報の検索キーと前記取得キーとが一致するものがない場合、

前記取得キーの宛先 IP アドレス部分の最後のゼロでないビットをゼロに変更するとともに、前記取得キーのマスク値部分の当該宛先 IP アドレス部分の変更位置を含む右側のビットをゼロに変更し、その変更したビット列を新たな取得キーに設定し、その新たな取得キーを用いて、自身の前記取得先選択部によって選択された取得先のルータから、さらにその取得先のルータの前記取得先選択部によって選択された取得先のルータと経路計算結果情報とを取得することを繰り返し、取得されたすべての前記経路計算結果情報の検索キーと前記新たな取得キーとが一致するか否かを判定して、前記経路計算結果情報の検索キーと前記取得キーとが最長一致となるときの前記ルータ集合ごとの経路計算の結果に基づいてパケットの転送先を決定することを、前記取得キーの宛先 IP アドレス部分のビット列がすべてゼロになるまで繰り返し実行する

ことを特徴とする請求項 4 に記載の経路情報管理システム。

【請求項 6】

前記取得先ノード情報は、前記ルータ ID を複数の範囲に区分し、自身のルータのルータ ID から前記距離の近いルータ ID ほどその区分の範囲を狭く設定し、前記距離の遠いルータ ID ほどその区分の範囲を広く設定して、それぞれの前記範囲ごとにその範囲に属するルータのルータ ID を記憶し、

前記取得先選択部は、前記取得キーと前記範囲の境界のルータ ID とを比較して、前記取得キーを含む前記範囲を抽出し、前記取得キーと前記抽出した範囲に属するルータのルータ ID との前記距離を算出し、前記距離の近い順に所定数のルータを取得先として選択する

ことを特徴とする請求項 4 に記載の経路情報管理システム。

【請求項 7】

請求項 4 に記載の経路情報管理システムにおいて用いられる経路情報管理装置の経路情報管理方法であって、

前記経路情報管理装置は、

前記 AS 内の 2 以上のルータの一まとまりを単位として前記 AS 内に複数設定されるルータ集合ごとに、少なくとも所定数以上のルータのルータ ID として、当該ルータの属する前記ルータ集合のルータ集合 ID のビット列の後にハッシュ関数によって生成されたハッシュ値のビット列を連結したビット列を記憶する格納先ノード情報を格納する記憶部を備え、

前記ルータから収集した前記プリフィクスごとに、そのプリフィクスを含む外部経路情報を用いて、前記 AS 内の 2 以上のルータの一まとまりを単位として前記 AS 内に複数設定されるルータ集合に属する少なくともいずれか 1 つのルータにおける、パケットの転送先のプリフィクスに至る転送経路を計算する経路計算を実行する経路計算ステップと、

前記ルータ集合ごとに、前記経路計算ステップによって経路計算されたルータ集合ごとの経路計算の結果を含む経路計算結果情報の格納先の検索に用いる検索キーとして、当該ルータ集合のルータ集合 ID のビット列、前記プリフィクスのビット列、前記プリフィクスに対するマスク値のビット列、およびパディングとしてすべて「0」のビット列を連結したビット列を生成する検索キー生成ステップと、

前記格納先ノード情報に記憶しているルータ ID と前記検索キーとの距離を、当該ルータ ID と前記検索キーとのビットごとの排他的論理和を示すビット列において、左から見て最初に「1」が出現したビット位置について右から数えたビットの桁数の数値として算出し、その算出した距離の小さい順に所定数のルータを当該ルータ集合から送信先として選択する格納先選択ステップと、

当該経路計算結果情報を、前記格納先選択ステップによって選択された当該所定数のルータに送信する入出力ステップと

を実行することを特徴とする経路情報管理装置の経路情報管理方法。

10

20

30

40

50

【請求項 8】

請求項 4 に記載の経路情報管理システムにおいて用いられるルータの経路情報管理方法であって、

前記ルータは、

前記 A S 内の 2 以上のルータの一まとまりを単位として前記 A S 内に複数設定されるルータ集合ごとに、少なくとも所定数以上のルータのルータ ID として、当該ルータの属する前記ルータ集合のルータ集合 ID のビット列の後にハッシュ関数によって生成されたハッシュ値のビット列を連結したビット列を記憶する取得先ノード情報と、前記経路情報管理装置から送信される検索キーと前記ルータ集合ごとの経路計算の結果とを関連付けた情報を示す経路計算結果情報とを記憶する記憶部を備え、

10

転送すべきパケットを受信したとき、当該ルータの属するルータ集合 ID のビット列、当該パケットの宛先 IP アドレスのビット列、前記宛先 IP アドレスに対するマスク値としてすべて「1」のビット列、およびパディングとしてすべて「0」のビット列を連結して形成される、取得キーを生成する取得キー生成ステップと、

前記取得先ノード情報に記憶しているルータ ID と前記取得キーとの距離を、前記ルータ ID と前記取得キーとのビットごとの排他的論理和を示すビット列において、左から見て最初に「1」が出現したビット位置について右から数えたビットの桁数の数値として算出し、その算出した距離の小さい順に所定数のルータを当該ルータ集合から取得先として選択する取得先選択ステップと、

転送すべきパケットを受信したときに、自身の前記記憶部に記憶する前記経路計算結果情報の検索キーと前記取得キーとが一致するか否かを判定し、一致するものがある場合には、その検索キーに関連付けられた前記ルータ集合ごとの経路計算の結果に基づいてパケットの転送先を決定し、一致するものがない場合には、前記取得先選択部によって取得先として選択されたルータに前記取得キーを送信して、当該ルータの記憶する前記経路計算結果情報の検索キーと前記取得キーとが一致する前記ルータ集合ごとの経路計算の結果を取得して、その取得した前記ルータ集合ごとの経路計算の結果に基づいてパケットの転送先を決定する転送経路決定ステップと、

20

前記転送経路決定ステップにおいて決定された転送先に前記パケットを転送する入出力ステップと

を実行することを特徴とするルータの経路情報管理方法。

30

【請求項 9】

請求項 7 に記載の経路情報管理方法を、コンピュータとしての経路情報管理装置に実行させるためのプログラム。

【請求項 10】

請求項 8 に記載の経路情報管理方法を、コンピュータとしてのルータに実行させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、IP (Internet Protocol) ネットワークにおいて、インターネットの経路情報を A S (自律システム; Autonomous System) 内で分散保持する技術に関する。

40

【背景技術】

【0002】

インターネット等の IP ネットワークは、異なる組織によって運営される多数の A S が接続されて形成されている。A S 間は、それぞれの A S の境界ルータによって接続されている。A S は、その A S の外部を宛先とするパケットを転送する必要があるため、インターネットの経路情報を A S 間で B G P (Border Gateway Protocol) を用いて交換している。そして、A S 内のすべてのルータは、A S 外部の外部経路情報をすべて保持する必要がある。具体的には、A S の境界ルータ同士は、e B G P (external BGP) によって、A S 外部の宛先プリフィクス(prefix) (IP アドレスの集合) に対する外部経路情報 (A S

50

外部のどの境界ルータにパケットを転送すれば宛先に到達するかの情報)を交換している。そして、その外部経路情報は、i B G P (internal BGP) によって、A S 内のルータに広告される。ただし、i B G P を用いてルータ間で広告するためには、I G P (Interior Gateway Protocol) によってルータ間が通信可能に接続されている必要がある。

【 0 0 0 3 】

広告された外部経路情報を受信したルータは、受信した外部経路情報に基づいて、宛先プリフィクスごとに、どの境界ルータに転送すれば宛先に到達できるのかを示した転送経路、すなわち、最適な転送経路を計算して保持していた。したがって、ルータは、宛先プリフィクスの数に対応する外部経路情報および経路計算によって算出した転送経路情報を記憶する必要があった。また、ルータは、宛先プリフィクスの数に対応する外部経路情報をメッセージとして広告する必要があった。

【 0 0 0 4 】

近年は、インターネットの急速な拡大にともなって、インターネット上に存在する宛先プリフィクスの数が増大してきている。そのため、インターネットサービスプロバイダ等のネットワーク事業者は、ネットワーク設備であるルータ等における、保持すべき経路情報の記憶容量の拡張や経路計算の処理負荷量の増大に対処せざるを得なくなり、設備更改のコスト負担増を強いられている。また、ルータ一台当りの、保持すべき経路数の増大、送信する必要のあるメッセージ数の増大、および経路計算の処理負荷量の増大は、転送すべきパケットを受信してから B G P 転送経路を決定するまでの時間の増大をもたらすため、ネットワーク制御の安定性に大きな影響を及ぼす虞がある。

【 0 0 0 5 】

従来技術では、ルータの処理負荷の低減のために、ルートリフレクタ (R R) (非特許文献 1 参照) または A S コンフェデレーション (非特許文献 2 参照) を導入することによって、i B G P を用いて広告する i B G P セッション数の削減が図られている。ルートリフレクタ (R R) は、ピアを確立したルートリフレクタクライアントとの間だけで外部経路情報を転送し、ピアを確立していないルートリフレクタクライアント同士では転送を行わない。そのため、ルートリフレクタ (R R) は、A S 内に確立すべきピア数を、i B G P を用いて A S 内の各ルータ間をフルメッシュで接続する場合より減少することが可能である。また、A S コンフェデレーションは、A S 内にいくつかのサブ A S を構成することによって、ピア数をフルメッシュの場合より減少することが可能である。このように、ルートリフレクタ (R R) または A S コンフェデレーションでは、ピア数を減少できるため、ルータの処理負荷の増大を抑制することが可能である。

【 0 0 0 6 】

また、ルータの記憶容量を低減するために、Compact Routing (非特許文献 3 , 4 参照) が提案されている。Compact Routing は、パケットを一旦経路の集約点となるルータに向けて転送し、その集約点となるルータから宛先アドレスへ転送する技術である。すなわち、Compact Routing は、経路を集約することによって、各ルータの経路情報の保持に要する記憶容量を低減している。

【 先行技術文献 】

【 非特許文献 】

【 0 0 0 7 】

【 非特許文献 1 】 T.Bates, E.Chen, and R.Chandra, " BGP Route Reflection : An Alternative to Full Mesh Internal BGP (IBGP) " , RFC4456 , IETF , 2006年4月

【 非特許文献 2 】 P.Traina, D.McPherson, and J.Scudder , " Autonomous System Confederations for BGP " , RFC5065 , IETF , 2007年8月

【 非特許文献 3 】 D.Krioukov, kc claffy, K.Fall, and A.Brady , " On Compact Routing for the Internet " , SIGCOMM Comput. Commun. Rev. , Vol.37, No.3, 2007年7月, p.41 52

【 非特許文献 4 】 P.Francis, X.Xu, and H.Ballani , " FIB Suppression with Virtual Aggregation " , internet Draft, draft ietf grow va 00.txt , 2009年5月

10

20

30

40

50

【発明の概要】

【発明が解決しようとする課題】

【0008】

しかしながら、非特許文献1に記載のルートリフレクタ(RR)を用いる方法、および非特許文献2に記載のASコンフェデレーションを用いる方法は、iBGPトポロジがフルメッシュでなくなり、また外部経路情報を交換する際にiBGPセッションにおいて宛先プリフィクス(以降、プリフィクスともいう。)ごとに1つの経路情報しか広告しない。そのため、iBGPセッション数が削減されて、ルータの処理負荷の低減や、経路情報の記憶容量の低減にある程度の効果は認められる。しかし、ルータが一部の経路情報しか記憶しないため、不完全な外部経路情報に基づいて経路計算が実行されてしまうので、経路計算の結果が必ずしも最適とならず、また計算された経路がループとなって収束しないことがあるという問題が発生する。

【0009】

また、非特許文献3,4に記載のCompact Routingを用いる方法は、パケットを一旦経路の集約点となるルータに集めるため、パケットの転送経路が必ずしも最適(最短)とならず、ネットワークのリソースを無駄に使用するという問題が発生する。

【0010】

そこで、本発明の課題は、前記した問題を解決しつつ、すなわち、すべての経路情報に基づいて経路計算を実行して最適な転送経路を算出することを満足しつつ、インターネットの今後の拡大に対応できる新しいルーティングアーキテクチャとして、ルーター台当りに保持すべき経路数(記憶容量)の増大、送信するメッセージ数の増大、および処理負荷の増大を抑制する技術を提供することを目的とする。

【課題を解決するための手段】

【0011】

本発明は、IP(Internet Protocol)ネットワークを構成するAS(Autonomous System)内で、他のASから該他のASのIPアドレスの集合であるプリフィクスを含む外部経路情報を取得し、かつパケットを他のASに転送する複数のルータと、同一のプリフィクスを含む外部経路情報を前記ルータから収集して経路計算を実行する2以上の経路情報管理装置とが相互に通信可能に接続される経路情報管理システムであって、前記経路情報管理装置が、前記ルータから収集した前記プリフィクスごとに、そのプリフィクスを含む外部経路情報を用いて、前記AS内の2以上のルータの一まとまりを単位として前記AS内に複数設定されるルータ集合に属する少なくともいずれか1つのルータにおける、パケットの転送先のプリフィクスに至る転送経路を計算する経路計算を実行する経路計算部と、前記経路計算部によって経路計算されたルータ集合ごとの経路計算の結果を含む経路計算結果情報を送信する当該ルータ集合に属する所定数のルータを選択する格納先選択部と、当該経路計算結果情報を、前記格納先選択部によって選択された当該所定数のルータに送信する入出力部と、を備え、前記ルータが、転送すべきパケットを受信したときに、前記経路情報管理装置から受信した前記経路計算結果情報を参照して、該パケットの転送先を決定する転送経路決定部と、前記転送経路決定部によって決定された転送先に前記転送すべきパケットを送信する入出力部と、を備えることを特徴とする。

【0012】

このような構成によれば、経路情報管理装置がプリフィクスに係る外部経路情報を格納するため、ルータに保持すべき経路数(記憶容量)を削減することができる。また、経路情報管理装置は、同一のプリフィクスに係る外部経路情報についてすべてを収集して経路計算を実行するため、経路計算結果が最適とならないという問題は発生しない。また、AS内に複数のルータ集合を設けた場合、そのルータ集合内のBGP経路はプリフィクスに対して共通になるため、経路情報管理装置は、ルータ集合内の一つのルータに対してだけ経路計算を実行すれば良くなる。そのため、経路情報管理装置は、すべてのルータに対して経路計算を実行する必要はなく処理負荷を低減することができる。また、ルータは、経路計算を実行しなくて良いため、処理負荷を低減することができる。さらに、経路情報管

理装置によって算出されたルータ集合ごとの経路計算の結果は、ルータ集合内の所定数のルータにしか送信しないため、経路計算の結果をメッセージとして送信するときのメッセージ数を削減することができる。

【 0 0 1 3 】

本発明は、前記経路情報管理装置が、さらに、前記ルータ集合ごとに、少なくとも1以上のルータのルータIDとして、当該ルータの属する前記ルータ集合のルータ集合IDのビット列の後にハッシュ関数によって生成されたハッシュ値のビット列を連結したビット列を記憶する格納先ノード情報を格納する記憶部と、前記ルータ集合ごとに、前記経路計算結果情報の格納先の検索に用いる検索キーとして、当該ルータ集合のルータ集合IDのビット列、前記プリフィクスのビット列、前記プリフィクスに対するマスク値のビット列、およびパディングとしてすべて「0」のビット列を連結したビット列を生成する検索キー生成部と、を備え、前記格納先選択部が、前記格納先ノード情報に記憶している前記ルータIDと前記検索キーとの距離を、前記ルータIDと前記検索キーとのビットごとの排他的論理和を示すビット列において、左から見て最初に「1」が出現したビット位置について右から数えたビットの桁数の数値として算出し、その算出した距離の小さい順に所定数のルータを当該ルータ集合から送信先として選択することを特徴とする。

【 0 0 1 4 】

このような構成によれば、ルータIDのビット列にはルータ集合IDとハッシュ値とを含み、検索キーのビット列にはルータ集合IDとプリフィクスとを含み、ルータIDと検索キーとの距離の小さい順に所定数のルータが選択される。したがって、ルータ集合IDが一致して、プリフィクスに距離の近いハッシュ値を有するルータが選択される。そのため、選択される経路計算結果情報の格納先のルータがほぼ均等化される。すなわち、格納先のルータがルータ集合内でほぼ均等に選択されるようになり、ルータの記憶容量の平均化が実現できる。言い換えると、一台当りのルータの記憶容量は、すべての経路計算結果情報を記憶していた従来の一台中のルータの記憶容量に比べて、低減できる。

【 0 0 1 5 】

本発明は、前記経路計算結果情報が、前記検索キーと、その検索キーに含まれるプリフィクスに対応する前記ルータ集合ごとの経路計算の結果とを関連付けた情報であることを特徴とする。

【 0 0 1 6 】

このような構成によれば、経路計算結果情報が、検索キーと、その検索キーに含まれるプリフィクスに対応するルータ集合ごとの経路計算の結果とを関連付けた情報となっている。そのため、検索キーを媒介として、検索キーのプリフィクスに距離の近いノードIDのルータに、当該プリフィクスに対応するルータ集合ごとの経路計算の結果を記憶させることができる。したがって、あるプリフィクスに係る経路計算結果情報を探索する場合、当該プリフィクスを含む検索キーを用いれば、その検索キーに距離の近いノードIDのルータおよびそのルータに記憶されている当該プリフィクスに対応するルータ集合ごとの経路計算の結果を取得することができる。すなわち、一台当りのルータの記憶容量を抑制するために分散して記憶されている経路計算結果情報の中から、検索キーによって、その検索キーに含まれるプリフィクスに対応するルータ集合ごとの経路計算の結果を検索することが可能となる。

【 0 0 1 7 】

本発明は、前記ルータが、さらに、そのルータの属するルータ集合内の自身を除く少なくとも1以上のルータのルータIDを記憶する取得先ノード情報と前記経路計算結果情報とを記憶する記憶部と、転送すべきパケットを受信したとき、当該ルータの属するルータ集合IDのビット列、当該パケットの宛先IPアドレスのビット列、前記宛先IPアドレスに対するマスク値としてすべて「1」のビット列、およびパディングとしてすべて「0」のビット列を連結して形成される、取得キーを生成する取得キー生成部と、前記取得先ノード情報に記憶しているルータIDと前記取得キーとの距離を、前記ルータIDと前記取得キーとのビットごとの排他的論理和を示すビット列において、左から見て最初に「1

」が出現したビット位置について右から数えたビットの桁数の数値として算出し、その算出した距離の小さい順に所定数のルータを当該ルータ集合から取得先として選択する取得先選択部と、を備え、前記転送経路決定部が、自身の前記記憶部に記憶する前記経路計算結果情報の検索キーと前記取得キーとが一致するか否かを判定し、一致するものがある場合には、その検索キーに関連付けられた前記ルータ集合ごとの経路計算の結果に基づいてパケットの転送先を決定し、一致するものがない場合には、前記取得先選択部によって取得先として選択されたルータに前記取得キーを送信して、当該ルータの記憶する前記経路計算結果情報の検索キーと前記取得キーとが一致する前記ルータ集合ごとの経路計算の結果を取得して、その取得した前記ルータ集合ごとの経路計算の結果に基づいてパケットの転送先を決定することを特徴とする。

10

【 0 0 1 8 】

また、本発明は、前記転送経路決定部が、さらに、前記経路計算結果情報の検索キーと前記取得キーとが一致するものがない場合、前記取得キーの宛先IPアドレス部分の最後のゼロでないビットをゼロに変更するとともに、前記取得キーのマスク値部分の当該宛先IPアドレス部分の変更位置を含む右側のビットをゼロに変更し、その変更したビット列を新たな取得キーに設定し、その新たな取得キーを用いて、自身の前記取得先選択部によって選択された取得先のルータから、さらにその取得先のルータの前記取得先選択部によって選択された取得先のルータと経路計算結果情報とを取得することを繰り返し、取得されたすべての前記経路計算結果情報の検索キーと前記新たな取得キーとが一致するか否かを判定して、前記経路計算結果情報の検索キーと前記取得キーとが最長一致となるとき

20

【 0 0 1 9 】

このような構成によれば、ルータは、転送すべきパケットを受信したとき、まず、自身の記憶部に記憶されている経路計算結果情報の検索キーと取得キーとが一致する経路計算結果情報を抽出する。そして、一致するものがない場合には、他のルータに記憶されている経路計算結果情報を探索して、経路計算結果情報の検索キーと取得キーとが一致する経路計算結果情報を抽出することができる。すなわち、一台当りのルータの記憶容量を抑制するために分散管理されている経路計算結果情報を、取得キーと検索キーとを照合して、検索することが可能となる。また、検索キーと取得キーとが一致する経路計算結果情報が見つからない場合の処理として、取得キーの宛先IPアドレス部分のビットの有効な桁を短くしていくことによって、検索キーと取得キーとが最長一致となる経路計算結果情報を自身および他のルータから抽出することができる。したがって、転送すべきパケットを受信したルータ自身が必要な経路計算結果情報を記憶していなくても、他のルータから転送により適した経路計算結果情報を取得することができる。すなわち、ルータ集合内のルータによって経路計算結果情報が分散管理されるため、一台当りのルータの記憶容量は、すべての経路計算結果情報を記憶していた従来の一台のルータの記憶容量に比べて、低減できる。

30

【 0 0 2 0 】

本発明は、前記取得先ノード情報が、前記ルータIDを複数の範囲に区分し、自身のルータのルータIDから前記距離の近いルータIDほどその区分の範囲を狭く設定し、前記距離の遠いルータIDほどその区分の範囲を広く設定して、それぞれの前記範囲ごとにその範囲に属するルータのルータIDを記憶し、前記取得先選択部が、前記取得キーと前記範囲の境界のルータIDとを比較して、前記取得キーを含む前記範囲を抽出し、前記取得キーと前記抽出した範囲に属するルータのルータIDとの前記距離を算出し、前記距離の近い順に所定数のルータを取得先として選択することを特徴とする。

40

【 0 0 2 1 】

このような構成によれば、ルータIDを複数の範囲に区分して記憶しているため、取得キーを含む範囲に属するルータIDと取得キーとの距離計算を行えば良くなる。したがっ

50

て、ルータに記憶されているすべてのルータIDと取得キーとの距離計算を行う必要がなく、計算処理量（処理負荷）を低減することが可能となる。

【0022】

本発明は、前記経路情報管理システムにおいて用いられる経路情報管理装置の経路情報管理方法であって、前記経路情報管理装置が、前記AS内の2以上のルータの一まとまりを単位として前記AS内に複数設定されるルータ集合ごとに、少なくとも所定数以上のルータのルータIDとして、当該ルータの属する前記ルータ集合のルータ集合IDのビット列の後にハッシュ関数によって生成されたハッシュ値のビット列を連結したビット列を記憶する格納先ノード情報を格納する記憶部を備え、前記ルータから収集した前記プリフィクスごとに、そのプリフィクスを含む外部経路情報を用いて、前記AS内の2以上のルータの一まとまりを単位として前記AS内に複数設定されるルータ集合に属する少なくともいずれか1つのルータにおける、パケットの転送先のプリフィクスに至る転送経路を計算する経路計算を実行する経路計算ステップと、前記ルータ集合ごとに、前記経路計算ステップによって経路計算されたルータ集合ごとの経路計算の結果を含む経路計算結果情報の格納先の検索に用いる検索キーとして、当該ルータ集合のルータ集合IDのビット列、前記プリフィクスのビット列、前記プリフィクスに対するマスク値のビット列、およびパディングとしてすべて「0」のビット列を連結したビット列を生成する検索キー生成ステップと、前記格納先ノード情報に記憶しているルータIDと前記検索キーとの距離を、当該ルータIDと前記検索キーとのビットごとの排他的論理和を示すビット列において、左から見て最初に「1」が出現したビット位置について右から数えたビットの桁数の数値として算出し、その算出した距離の小さい順に所定数のルータを当該ルータ集合から送信先として選択する格納先選択ステップと、当該経路計算結果情報を、前記格納先選択ステップによって選択された当該所定数のルータに送信する入出力ステップとを実行することを特徴とする。

【0023】

また、本発明は、前記経路情報管理システムにおいて用いられるルータの経路情報管理方法であって、前記ルータが、前記AS内の2以上のルータの一まとまりを単位として前記AS内に複数設定されるルータ集合ごとに、少なくとも所定数以上のルータのルータIDとして、当該ルータの属する前記ルータ集合のルータ集合IDのビット列の後にハッシュ関数によって生成されたハッシュ値のビット列を連結したビット列を記憶する取得先ノード情報と、前記経路情報管理装置から送信される検索キーと前記ルータ集合ごとの経路計算の結果とを関連付けた情報を示す経路計算結果情報とを記憶する記憶部を備え、転送すべきパケットを受信したとき、当該ルータの属するルータ集合IDのビット列、当該パケットの宛先IPアドレスのビット列、前記宛先IPアドレスに対するマスク値としてすべて「1」のビット列、およびパディングとしてすべて「0」のビット列を連結して形成される、取得キーを生成する取得キー生成ステップと、前記取得先ノード情報に記憶しているルータIDと前記取得キーとの距離を、前記ルータIDと前記取得キーとのビットごとの排他的論理和を示すビット列において、左から見て最初に「1」が出現したビット位置について右から数えたビットの桁数の数値として算出し、その算出した距離の小さい順に所定数のルータを当該ルータ集合から取得先として選択する取得先選択ステップと、転送すべきパケットを受信したときに、自身の前記記憶部に記憶する前記経路計算結果情報の検索キーと前記取得キーとが一致するか否かを判定し、一致するものがある場合には、その検索キーに関連付けられた前記ルータ集合ごとの経路計算の結果に基づいてパケットの転送先を決定し、一致するものがない場合には、前記取得先選択部によって取得先として選択されたルータに前記取得キーを送信して、当該ルータの記憶する前記経路計算結果情報の検索キーと前記取得キーとが一致する前記ルータ集合ごとの経路計算の結果を取得して、その取得した前記ルータ集合ごとの経路計算の結果に基づいてパケットの転送先を決定する転送経路決定ステップと、前記転送経路決定ステップにおいて決定された転送先に前記パケットを転送する入出力ステップとを実行することを特徴とする。

【0024】

10

20

30

40

50

このような構成によれば、経路情報管理装置がプリフィクスに係る外部経路情報を格納するため、ルータに保持すべき経路数（記憶容量）を削減することができる。また、経路情報管理装置は、同一のプリフィクスに係る外部経路情報についてすべてを収集して経路計算を実行するため、経路計算結果が最適とならないという問題は発生しない。また、AS内に複数のルータ集合を設けた場合、そのルータ集合内のBGP経路はプリフィクスに対して共通に定まるため、経路情報管理装置は、ルータ集合内の1つのルータに対してだけ経路計算を実行すれば良くなる。そのため、経路情報管理装置の処理負荷を低減することができるとともに、ルータでの経路計算が不要となる。また、経路情報管理装置によって算出されたルータ集合ごとの経路計算の結果は、ルータ集合内で分散管理されるため、ルータ集合ごとに所定数のルータに格納（送信）されれば良い。そのため、経路計算の結果をメッセージとして送信するときのメッセージ数を削減することができる。また、ルータは、転送すべきパケットを受信したとき、まず、自身の記憶部に記憶されている経路計算結果情報の検索キーと取得キーとが一致する経路計算結果情報を抽出する。そして、一致するものがない場合には、他のルータに記憶されている経路計算結果情報を探索して、経路計算結果情報の検索キーと取得キーとが一致する経路計算結果情報を抽出することができる。すなわち、一台当りのルータの記憶容量を抑制するために分散管理されている経路計算結果情報を、取得キーと検索キーとを照合して、検索することが可能となる。したがって、転送すべきパケットを受信したルータ自身が必要な経路計算結果情報を記憶していなくても、他のルータからより転送に適した経路計算結果情報を取得することができる。すなわち、ルータ集合内のルータによって経路計算結果情報が分散管理されるため、一台当りのルータの記憶容量は、すべての経路計算結果情報を記憶していた従来の一台のルータの記憶容量に比べて、低減できる。

【0025】

本発明は、前記経路情報管理方法を、コンピュータとしての経路情報管理装置またはルータに実行させるためのプログラムとした。

【0026】

このようなプログラムをインストールされたコンピュータは、このプログラムに基づいた機能を実現することができる。

【発明の効果】

【0027】

本発明によれば、すべての経路情報に基づいて経路計算を実行して最適な転送経路を算出することを満足しつつ、ルータ一台当りに保持すべき経路数（記憶容量）の増大、送信するメッセージ数の増大、および処理負荷の増大を抑制する技術を提供することができる。

【図面の簡単な説明】

【0028】

【図1】本実施形態の概要を表す図である。

【図2】本実施形態におけるルートサーバの機能の一例を示す図である。

【図3】本実施形態におけるルータの機能の一例を示す図である。

【図4】ルートサーバの決定方法の例を示す図である。

【図5】本実施形態における距離の定義を示す図である。

【図6】ルートサーバの格納先ノードDBの一例を示す図である。

【図7】ルータの取得先ノードDBの一例を示す図である。

【図8】ルータの経路計算結果DBの一例を示す図である。

【図9】本実施形態における経路計算結果の格納処理の流れを示す図である。

【図10】本実施形態における他のルータから転送経路情報を取得する処理の流れを示す図である。

【図11】取得キーkdのビット変更時の操作を示す図である。

【発明を実施するための形態】

【0029】

10

20

30

40

50

次に、本発明を実施するための形態（以降「本実施形態」と称す）について、適宜図面を参照しながら詳細に説明する。

【0030】

本実施形態の概要

初めに、本実施形態の概要について、図1を用いて説明する。図1に示すように、経路情報管理システム30は、AS1内に、外部経路情報の収集およびルーティングの管理を行う2以上のルートサーバ（経路情報管理装置）10A, 10B（10）および境界ルータを含む複数のルータ20（20A, 20B, ..., 20L）を備えている。そして、各ルートサーバ10は、ルータ20と通信可能に接続される。各ルータ20間は、相互に通信可能に接続される。また、各ルータ20は、一点鎖線で示される、いずれかのPOP（Point of Presence）に属している。POPは、例えば、同じ市や同じ建物等の同一のロケーションに存在するルータの集合（ルータ集合）を表している。なお、図1では、AS1は、POP1, POP2, POP3, POP4によって構成されているが、POPの数は2以上であれば、4つに限られない。また、ルートサーバ10の台数、ルータ20の台数は図1に示した台数に限られない。また、ルータ20間の接続は、図1に示した実直線に限られない。

【0031】

ここで、経路情報管理システム30において、ルータ20の記憶容量、送信する必要があるメッセージ数、および処理負荷を抑制するために設けた3つの特徴について説明する。

（特徴1）プリフィクスに係る外部経路情報を、ルートサーバ10に集約する。

（特徴2）ルートサーバ10は、POPごとにプリフィクスに係る経路計算を実行し、その経路計算の結果を、POPごとにそのPOPに属する一部のルータ20に格納（送信）する。

（特徴3）POP内では、そのPOPに属するルータ20間で、経路計算結果を分散管理し、パケットの転送に必要な経路計算結果を相互に取得する。

【0032】

（特徴1について）

図1に示すように、AS1内のルータ20Aは、AS2からプリフィクスP1に係る外部経路情報をeBGPによって受信したものとす。そして、受信したプリフィクスP1に係る外部経路情報は、ルートサーバ10Aに収集される。また、ルータ20Kが、AS3からプリフィクスP7に係る外部経路情報をeBGPによって受信したものとす。そして、受信したプリフィクスP7に係る外部経路情報は、ルートサーバ10Bに収集される。すなわち、各ルートサーバ10A, 10Bは、それぞれ、異なるプリフィクスに係る外部経路情報を収集する。なお、どのルートサーバ10が、どのプリフィクスに係る外部経路情報を収集するかは、予め決められている。そのプリフィクスを収集するルートサーバの決定方法については、後記する。

【0033】

このようにして、ルートサーバ10は、いずれかのプリフィクスについてはすべての外部経路情報を収集できるので、すべての外部経路情報に基づいて経路計算を行うことが可能となる。つまり、背景技術の欄で説明したルトリフレクタ（RR）を用いる方法やASコンフェデレーションを用いる方法に見られるような、一部の外部経路情報に基づいて経路計算を実行してしまうことにもなう、経路計算結果が最適とならないというケースや、経路計算が収束しないというケースは起こらない。また、各ルートサーバ10がプリフィクスごとに外部経路情報を分散して担当するので、AS1内のルータ20の台数の増加に対応してルートサーバ10の台数を増加することによって、ルートサーバ10一台当りの記憶容量の増加や処理負荷の増加を抑制することが可能となる。

【0034】

また、従来の境界ルータは、外部経路情報の受信リストおよび送信リストを記憶するRIB（Routing Information Base）と、パケットの転送先を記憶するFIB（Forwarding

10

20

30

40

50

Information Base)とを備えていた。しかし、本実施形態の境界ルータ20では、外部経路情報の受信リストおよび送信リストは、ルートサーバ10に集約されるため、RIBが不要となる。すなわち、ルータ20に保持すべき経路数(記憶容量)を削減することができる。

【0035】

(特徴2について)

ルートサーバ10は、プリフィクスごとにPOP内の1つのルータについて、経路計算を行えば良い。この理由は、以下の通りである。すなわち、通常は、AS1内のPOP間に張られるリンクのIGPリンクコストは、POP内のリンクコストより大きい値に設定されている。そして、経路計算においては、経路を構成するリンクごとに割り当てられたリンクコストを合計した値が最も小さい経路が優先して選択される。このようにして、同じPOP内を宛先とするパケットは、そのPOPから外に出ないようにして、POP間を行ったり来たりするような無駄な転送を防いでいる。そのため、同じPOP内のルータは、同じプリフィクスを宛先とするBGP経路については同一の経路(ネクストホップ)をとることになる。つまり、POP内のすべてのルータにとって、AS外部向けの転送経路は共通になる。そこで、同じプリフィクスを宛先とする経路計算は、POP内の1つのルータに対してだけ行えば良いことになる。以降、これを、"POPごとの経路計算"と呼ぶことにする。なお、POPごとの経路計算は、ルートサーバ10が、ルータ20における経路計算の代わりに行うもので、公知のリンクコストを用いた経路計算手法を用いることができる。ただし、このPOPごとの経路計算は、ルートサーバ10が、自身を起点として経路計算を実行するのではなく、POP内の1つのルータ20を起点として実行する点において、公知の経路計算手法と異なるだけである。したがって、POPごとの経路計算は、AS1内のすべてのルータ20について経路計算が行われていた従来の方法に比べると、処理負荷の低減になっていると言える。

【0036】

また、ルートサーバ10は、経路計算結果を各POPの一部のルータに格納(送信)する(図1の点線)。これは、後記する特徴3が備わっているために可能となっている。そして、すべてのルータ20に対して経路計算結果をメッセージとして配信していた従来に比べて、本実施形態では、メッセージ数を削減することができる。

【0037】

(特徴3について)

境界ルータ20は、転送すべきパケットを受信して、そのパケットの転送に必要な経路計算結果を自身の記憶部23(図3参照)に記憶していないとき、自身の属するPOP内の他のルータ20から、必要な経路計算結果を取得する(図1の実線矢印)。そして、境界ルータ20は、取得した経路計算結果に基づいて、受信したパケットを転送することができる。言い換えると、POP内のルータ20は、ルートサーバ10から格納(送信)された経路計算結果を分散管理している。すなわち、各ルータ20は、一部の経路計算結果しか記憶しなくてもよいので、従来より記憶容量を低減することができる。この分散管理の方法の詳細については、後記する。

【0038】

ルートサーバ10の機能

次に、本実施形態の経路情報管理システム30に用いられるルートサーバ10の機能について、図2を用いて説明する。なお、図2は、本実施形態のルートサーバ10の一部の機能のみを示し、BGPやIGP等に係る機能については記載を省略している。

【0039】

図2に示すように、ルートサーバ10は、外部経路情報や経路計算結果等の各種情報の入出力を司る入出力部11と、外部経路情報に基づいて経路計算する処理部12と、記憶部13とを備える。入出力部11は、他のルータ20等との情報を入出力する。入出力部11における情報の入出力は、処理部12からの指示に基づいて行われる。

【0040】

10

20

30

40

50

処理部 12 は、ルートサーバ 10 が備える図示しない CPU (Central Processing Unit) と図示しないメインメモリとによって構成されている。そして、この CPU がメインメモリにアプリケーションプログラムを展開して、それを実行することにより種々の演算処理を具現化する。なお、アプリケーションプログラムは、記憶部 13 に格納されている。また、処理部 12 は、外部経路情報収集部 121、経路計算部 122、格納先選択部 123、および検索キー生成部 124 を備えている。

【0041】

外部経路情報収集部 121 は、境界ルータ 20 から広告されてくる外部経路情報を入出力部 11 を介して受信する。なお、同一のプリフィクスに係る外部経路情報は、すべて、同じルートサーバ 10 に収集される。そして、外部経路情報収集部 121 は、受信した外部経路情報を記憶部 13 の外部経路情報 DB 131 に記憶する。経路計算部 122 は、外部経路情報 DB 131 に新たに記憶された外部経路情報を読み出して、プリフィクスごとに前記した "POP ごとの経路計算" を実行し、POP ごとの経路計算の結果、すなわち最適な転送経路および次善の転送経路を算出する。予め次善の転送経路を算出しておく理由は、最適な転送経路が輻輳や故障等によって利用できない場合に、直ちにこの次善の転送経路に切り替えて用いられるようにするためである。

【0042】

格納先選択部 123 は、格納先ノード DB 132 に記憶されているルータ 20 のルータ ID と、後記するプリフィクスを含む検索キーとを用いて、経路計算部 122 によって算出された POP ごとの経路計算の結果の格納 (送信) 先である POP 内のルータ 20 を選択する。そして、格納先選択部 123 は、入出力部 11 を介し、POP ごとの経路計算の結果を、POP ごとに選択したルータ 20 に格納 (送信) する。また、検索キー生成部 124 は、POP ごとの経路計算の結果を格納するとき用いる検索キーを生成する。なお、格納方法および検索キーの生成については、後記する。

【0043】

記憶部 13 は、外部経路情報 DB 131 および格納先ノード DB 132 を格納している。外部経路情報 DB 131 は、境界ルータ 20 から広告されたプリフィクスに係る外部経路情報を記憶している。例えば、外部経路情報は、プリフィクスと、そのプリフィクスに関連付けられたパス属性情報とを含んでいる。パス属性情報は、AS path やネクストホップに至るリンクに割り当てられたリンクコスト等である。なお、AS path は、経由してきた AS を記録しておくためのものであり、AS を経由するごとに AS 番号が追加される。そして、AS path に記録された AS 番号の数が少ない方の経路が優先して用いられる。外部経路情報 DB 131 に記憶されている外部経路情報は、経路計算部 122 によって、経路計算を実行するとき用いられる。また、格納先ノード DB 132 は、ルータ 20 を識別可能なルータ ID を記憶している。このルータ ID は、AS の管理者によって予め設定される。そして、格納先ノード DB 132 は、格納先選択部 123 が格納先のルータ 20 を選択するとき参照される。格納先ノード DB 132 の具体例については、後記する。

【0044】

ルータ 20 の機能

次に、本実施形態のルータ 20 の機能について、図 3 を用いて説明する。なお、図 3 は、本実施形態のルータ 20 の一部の機能のみを示し、BGP や IGP 等に係る機能については記載を省略している。

【0045】

図 3 に示すように、ルータ 20 は、ルートサーバ 10 や他のルータ 20 との間で送受信される情報の入出力を制御する入出力部 21 と、処理部 22 と、記憶部 23 とを備える。入出力部 21 における情報の入出力は、処理部 22 からの指示に基づいて行われる。

【0046】

処理部 22 は、ルータ 20 が備える図示しない CPU (Central Processing Unit) と図示しないメインメモリとによって構成されている。そして、この CPU がメインメモリにアプリケーションプログラムを展開して、それを実行することにより種々の演算処理を具

現化する。なお、アプリケーションプログラムは、記憶部 2 3 に格納されている。また、処理部 2 2 は、ルートサーバ選択部 2 2 1、転送経路決定部 2 2 2、取得先選択部 2 2 3、カプセル化処理部 2 2 4、および取得キー生成部 2 2 5 を備えている。

【 0 0 4 7 】

ルートサーバ選択部 2 2 1 は、記憶部 2 3 のルートサーバ I D D B 2 3 1 に記憶されているルートサーバ I D 情報を参照して、後記するルートサーバ決定方法に基づいて、他の A S から受信した外部経路情報を広告するルートサーバ 1 0 を選択する。そして、ルートサーバ選択部 2 2 1 は、入出力部 2 1 を介して、他の A S から受信した外部経路情報を、選択したルートサーバ 1 0 へ広告する。

【 0 0 4 8 】

転送経路決定部 2 2 2 は、転送すべきパケットを受信したときに、記憶部 2 3 の経路計算結果 D B 2 3 2 に記憶している経路計算結果情報を探索し、そのパケットの転送先を決定する。しかし、経路計算結果 D B 2 3 2 は、経路計算に必要な経路計算結果情報の一部しか記憶していない場合がある。その場合には、転送経路決定部 2 2 2 は、取得先選択部 2 2 3 によって選択された他のルータ 2 0 から不足している経路計算結果情報を取得する。そして、転送経路決定部 2 2 2 は、その取得した経路計算結果情報に基づいて、パケットの転送経路（転送先）を決定する。

【 0 0 4 9 】

取得先選択部 2 2 3 は、転送すべきパケットを受信したときに、取得先ノード D B 2 3 3 に記憶しているルータ 2 0 に係る情報を参照して、不足する経路計算結果情報を保持するルータ 2 0 を選択する。また、カプセル化処理部 2 2 4 は、転送すべきパケットに対して、出口となる境界ルータ 2 0 のアドレスのヘッダでカプセル化する処理を実行する。また、出口の境界ルータ 2 0 では、カプセル化処理部 2 2 4 は、カプセル化されたパケットから、その出口の境界ルータ 2 0 のアドレスのヘッダを取り外す処理を実行する。また、取得キー生成部 2 2 5 は、経路計算結果情報の取得先のルータ 2 0 を選択するときに用いる取得キーを生成する。なお、取得先の選択方法および取得キーの生成方法については、後記する。

【 0 0 5 0 】

記憶部 2 3 は、ルートサーバ 1 0 から配信される経路計算結果情報を格納する経路計算結果 D B 2 3 2、ルートサーバ 1 0 の I D 番号を格納するルートサーバ I D D B 2 3 1、および経路計算結果情報を分散管理するルータ 2 0 に係る情報を格納する取得先ノード D B 2 3 3 を格納する。経路計算結果 D B 2 3 2 は、転送経路決定部 2 2 2 の経路計算において参照される。ルートサーバ I D D B 2 3 1 は、ルートサーバ選択部 2 2 1 が、外部経路情報を広告するルートサーバ 1 0 を選択するときに参照される。取得先ノード D B 2 3 3 は、取得先選択部 2 2 3 が必要な経路計算結果情報の取得先のルータ 2 0 を選択するときに参照される。経路計算結果 D B 2 3 2、ルートサーバ I D D B 2 3 1、および取得先ノード D B 2 3 3 の具体例については、後記する。

【 0 0 5 1 】

ルートサーバ決定方法

ここで、ルータ 2 0 におけるルートサーバ決定方法について、図 4 を用いて説明する（適宜図 3 参照）。図 4 は、印で示すルートサーバ 1 0 とそれをカバーするハッシュ空間を表している。図 4 中の印の中の数字は、ルートサーバ 1 0 の I D 番号を示している。そして、図 4 は、ハッシュ空間に、3 つのルートサーバ 1 0 A、1 0 B、1 0 C が存在する場合を表している。例えば、ルートサーバ 1 0 A の I D 番号が 5、ルートサーバ 1 0 B の I D 番号が 1 0、ルートサーバ 1 0 C の I D 番号が 1 5、であるものとして説明する。なお、これらのルートサーバ 1 0 A、1 0 B、1 0 C の I D 番号に係る情報が、図 3 の記憶部 2 3 のルートサーバ I D D B 2 3 1 に記憶されることになる。

【 0 0 5 2 】

まず、next(m)という関数を導入する。next(m)はハッシュ空間においてハッシュ空間を時計回り（つまり I D を大きくする方向）で次のノードの I D を返す。例えば、next(2

10

20

30

40

50

)=5、next(6)=10、next(15)=15となる。すなわち、next(m)は、ハッシュ空間のある値m以上の範囲において最小の値となるIDを返す。ここで、mは、プリフィクスの値をハッシュ関数に代入して算出されたハッシュ値である。ただし、図4に示すハッシュ空間では、 $0 < m \leq 15$ である。

【0053】

そして、図4では、ハッシュ空間におけるルートサーバ10Aの分担範囲は0を超え5以下、ルートサーバ10Bの分担範囲は5を超え10以下、ルートサーバ10Cの分担範囲は10を超え15以下となる。このようにして、ルートサーバ10のID番号を予め決めておけば、プリフィクスに関する経路情報を送信するルートサーバ10を特定することが可能となる。

【0054】

経路計算結果情報の分散管理

前記したように、POP内のルータ20がBGP経路については同一の経路をとるので、経路計算結果情報をPOP内のルータ20によって分散管理することが可能となる。そして、経路計算結果情報を分散管理することによって、ルータ20が保持すべき経路数(記憶容量)を抑制することが可能となる。以下に、分散管理の方法について説明する。

【0055】

(フレームワーク)

まず、分散管理のためのフレームワークについて説明する。AS内のルータ20には、そのルータを識別するために、ASの管理者によってユニークなルータIDが割り当てられる。また、経路計算結果情報は、分散管理が可能ないように、検索キーと経路計算結果の転送経路情報(最適な転送経路および次善の転送経路)とが紐付けられて形成される。ただし、ルータIDと検索キーとは、同じビット長で表現される。

【0056】

まず、検索キーは、ルートサーバ10において、経路計算結果情報を、どのルータ20に格納(送信)するかの選択に用いられる。そして、ルートサーバ10は、後記する検索キーとルータIDとの距離に基づいて、検索キーに距離の近いルータIDのルータを、距離の近い順にk個選択する。k個選択する意味は、冗長性の確保のためである。したがって、冗長性を考えなければ、 $k = 1$ であっても良い。そして、ルートサーバ10は、選択したk個のルータ20に、経路計算結果情報を格納(送信)する。

【0057】

ここで、距離の定義について、図5を用いて説明する。図5では、説明を簡単にするために、ルータIDおよび検索キーのビット長をそれぞれ5ビットで表している。各ルータa, b, c, dのルータIDがそれぞれ10011, 10101, 10110, 11000であったとする。そして、検索キーが、10111であったとする。この場合、ルータIDと検索キーとのExclusive ORを算出する。そして、Exclusive ORの演算結果のビット列において、左から見て最初に「1」が出現した位置について右から数えた桁数の数値を、「距離」と定義する。すなわち、各ルータa, b, c, dに対して、ルータIDと検索キーとの距離は、それぞれ3, 2, 1, 4となる。したがって、検索キーはルータcに最も近いと判定される。そして、冗長性を考慮して $k = 2$ とした場合には、ルータcの次に距離が近いのは、ルータbと判定される。このようにして、検索キーおよびその検索キーに紐付けられた転送経路情報によって形成される経路計算結果情報は、ルータb, cに格納されることになる。この距離を用いることによって、距離の近いものほど、左から数えたビット列の一致する長さが長くなり、距離が0の場合は、完全に一致することを表す効果がある。

【0058】

次に、ルータ20が、他のルータ20から必要な経路計算結果情報を取得するときに、取得先のルータ20を選択する方法について説明する。この場合、経路計算結果情報を格納するときと同じ距離を用いて、必要な経路計算結果情報を格納したルータ20を特定することができる。以下に、ルータID、検索キー、および転送経路情報の構成と、格納先ノードDB132、取得先ノードDB233、および経路計算結果DB232の一例につ

10

20

30

40

50

いて、図6, 7, 8を用いて説明する。なお、本実施形態では、例えば、ビット長は160ビットとして説明する。

【0059】

ルータIDの構成を式(1)に示す。

<ルータID> ::= <POPのID> <ハッシュ値> …式(1)

ただし、POPのIDは、そのルータ20が属するPOPを識別するIDであり、32ビットとする。POPのIDは、AS管理者によって設定されるユニークな識別情報である。また、ハッシュ値は、例えば、ハッシュ関数の一つである、ランダム値を基にしたMD5(Message Digest 5)により算出したMD5ハッシュ値であって、128ビットとする。

10

【0060】

経路計算結果情報の格納先のルータ20を選択するとき用いる検索キーkrの構成を式(2)に示す。

<kr> ::= <POPのID> <プリフィクス値> <マスク値> <パディング> …式(2)

ただし、POPのIDは、前記と同様である。つまり、POPのIDによって、どのPOPのための検索キーであるかが特定され、同じPOPのIDを含むノードIDを見つけやすくすることができる。プリフィクス値は、外部経路情報に含まれるプリフィクスであり、32ビットのIPv4(Internet Protocol version 4)である。マスク値は、プリフィクス値に対するマスク値であり、32ビットである。このマスク値は、IPv4で用いられるサブネットマスクに相当する。パディングは、その値をすべて「0」とし、64ビットである。なお、「::=」は、左辺のビット列と右辺のビット列とが同じになることを表し、右辺は<>のビット列をこの順で連結することを表している。

20

【0061】

経路計算結果情報を形成する、検索キーkrに紐付けられる転送経路情報valueの構成を式(3)に示す。

<value> ::= <バージョン番号> <n> <NH₁> <優先度1> …
… <NH_n> <優先度n> …式(3)

ここで、バージョン番号は、ルートサーバ10が経路計算を実行するごとに付与される。nは、valueに格納する経路計算の結果である転送経路(NH_{__}*)の数である。NH_{__}*は、プリフィクスに対応する転送経路(ネクストホップ情報)を表す。優先度は、転送経路を用いる優先順位を指定するもので、優先度の高いものから順に用いられる。なお、転送経路情報valueは、複数の転送経路を格納しているので、選択された転送経路の障害時には、次に優先度の高い転送経路を直ちに用いることができる。

30

【0062】

次に、ルートサーバ10の記憶部13の格納先ノードDB132(図2参照)の一例を、図6を用いて説明する。格納先ノードDB132は、POPごとに、そのPOPに属するルータ20の、ルータ名、IPアドレス、およびルータIDを記憶する。なお、POP内では経路計算結果情報が分散管理されるため、格納先ノードDB132は、AS内のすべてのルータ20に係る情報について記憶しなくても良い。

40

【0063】

また、ルータ20の記憶部23の取得先ノードDB233(図3参照)の一例を、図7を用いて説明する。なお、図7に示す取得先ノードDB233は、ある一つのルータ20に記憶されているものを表している。rangeは、ルータIDをいくつかに分けて示している。例えば、図4に示すハッシュ空間がルータIDとして用いられるID全体を示すものとする、ルートサーバ10の分担範囲がrangeに相当する。そして、その範囲は、自身のルータIDに距離の近い範囲は狭く、距離の遠い範囲は広く設定される。図7では、冗長性をもたせて、取得先のルータの台数k=2とした場合を表しており、rangeごとに、2つの取得先のルータに係る情報を記憶した状態を表している。

【0064】

50

また、ルータ20の記憶部23の経路計算結果DB232(図3参照)の一例を、図8を用いて説明する。経路計算結果DB232は、POP内の一部の経路計算結果情報を記憶している。なお、経路計算結果情報は、検索キーkrと転送経路情報valueとの組(kr, value)によって表される。

【0065】

(ルートサーバ10の格納処理の流れ)

次に、本実施形態における、ルートサーバ10による格納処理の流れについて、図9を用いて説明する(適宜図2参照)。ステップS901では、経路計算部122が、外部経路情報DB131を参照して、各プリフィクスに対してPOPごとの経路計算を実行し、転送経路情報valueを生成する。ステップS902では、検索キー生成部124が、ルータ20から受信したプリフィクスについて、そのプリフィクスを含む検索キーKrを生成する。ステップS903では、格納先選択部123が、POPごとに、検索キーkrとルータIDとの距離の近い(小さい)順にk個(予め決められている所定数)のルータIDを抽出する。そして、ステップS904では、格納先選択部123が、抽出したルータIDのルータ20に、そのルータ20が記憶している取得先のルータに係る情報をk個取得する要求find_node(kr)を送信する。なお、find_node(kr)を受信したルータ20の取得先選択部223は、検索キーkrとrange(図7参照)の境界のルータIDとを比較して、検索キーkrを含むrangeを抽出し、検索キーkrと抽出したrangeに属するルータ20それぞれのルータIDとの距離を算出し、距離の近い順にk個の取得先のルータに係る情報を、ルートサーバ10に返信する。

【0066】

ステップS905では、格納先選択部123は、find_node(kr)の送信先のルータ20から、取得先のルータに係る情報をk個取得する。ステップS906では、格納先選択部123は、ステップS905で取得したルータの中で、find_node(kr)を未送信のルータがあるか否かを判定する。そして、未送信のルータがある場合(ステップS906でYes)、ステップS907で、格納先選択部123が、未送信のルータに、find_node(kr)を送信する。そして、処理は、ステップS905へ戻る。また、未送信のルータがない場合(ステップS906でNo)、ステップS908では、格納先選択部123は、POPごとに、検索キーkrに距離の近い順にk個のルータIDを選択し、選択したk個のルータIDのルータに、経路計算結果情報(kr, value)を格納(送信)する。なお、ステップS908で格納された経路計算結果情報(kr, value)は、ルートサーバ10が経路計算を再実行した場合には、削除される。また、ルートサーバ10が、find_node(kr)の送信先から新しく取得したルータIDは、格納先ノードDB132に格納し、更新していく。

【0067】

(他のルータ20から転送経路情報を取得する処理)

次に、他のルータ20から転送経路情報を取得する処理について、その概要を説明する(適宜図3参照)。まず、ルータ20の取得キー生成部225は、転送すべきパケットが到着すると、転送先の決定のために、転送すべきパケットの宛先IPアドレスを含む取得キーkdを生成する。なお、この取得キーkdの詳細については、後記する。そして、転送経路決定部222は、自身の記憶する経路計算結果DB232を参照して、取得キーkdと一致する、経路計算結果情報に含まれる検索キーkrを探索する。取得キーkdに一致する検索キーkrが見つかった場合には、その検索キーkrに紐付けられた転送経路情報valueが読み出され、パケットの転送に用いられる。

【0068】

また、取得先選択部223は、転送すべきパケットが到着すると、取得キーkdに距離の近いk個のルータIDのルータ20を選択する。なお、この距離は、前記した距離と同様の定義である。そして、取得先選択部223は、その選択したルータ20に、そのルータ20が保持する取得キーkdに距離の近いk個のルータに係る情報と、取得キーkdに一致する検索キーkrの経路計算結果情報(kr, value)がある場合には経路計算

結果情報 (k r , v a l u e) とを取得するための要求 f i n d _ v a l u e (k d) を送信する。そして、転送経路決定部 2 2 2 は、取得した複数の経路計算結果情報 (k r , v a l u e) の中から、取得キー k d に最長一致 (例えば、複数の検索キー k r と取得キー k d との距離を算出したとき、最も距離の近いビット列を見つけること) となる検索キー k r を含む経路計算結果情報 (k r , v a l u e) を選択する。次に、転送経路決定部 2 2 2 は、選択した該経路計算結果 (k r , v a l u e) から転送経路情報 v a l u e を読み出して、パケットの転送に用いる。

【 0 0 6 9 】

(取得キー k d の構成)

ここで、取得キー k d の構成を式 (4) に示す。

< k d > : : = < P o P の I D > < I P アドレス > < マスク値 > < パディング >

・ ・ 式 (4)

ただし、P o P の I D は、転送すべきパケットを受信したルータの属する P o P の I D である。つまり、P o P の I D によって、どの P o P のための取得キーであるかが特定される。I P アドレスは、転送すべきパケットの宛先 I P アドレスであり、3 2 ビットである。マスク値は、プリフィクス値に対するマスク値であり、3 2 ビットである。このマスク値は、I P v 4 で用いられるサブネットマスクに相当し、すべて「 1 」に設定する。パディングは、その値をすべて「 0 」とし、6 4 ビットである。

【 0 0 7 0 】

このように、境界ルータ 2 0 では、転送すべきパケットが到着すると、転送先の決定のために、分散管理されている転送経路情報 v a l u e を取得する必要がある。そのため、境界ルータ 2 0 は、転送経路情報 v a l u e との組を形成している検索キー k r を生成できればよいが、転送すべきパケットからは、宛先 I P アドレスしか取得できない。そして、その宛先 I P アドレスを示すビット列と、宛先のプリフィクスを示すビット列とは完全に一致するとは限らない。すなわち、境界ルータ 2 0 は、取得キー k d を生成する場合、経路計算結果情報 (k r , v a l u e) を格納するとき用いた検索キー k r と同じビット列を生成できるとは限らない。そこで、ルータ 2 0 が、取得キー k d を用いて行う転送経路情報の取得処理の流れについて、図 1 0 を用いて説明する (適宜図 3 参照) 。

【 0 0 7 1 】

(転送経路情報の取得処理の流れ)

図 1 0 に示すように、ステップ S 1 0 0 1 では、取得キー生成部 2 2 5 が、転送すべきパケットの宛先 I P アドレスを含む取得キー k d を生成する。ステップ S 1 0 0 2 では、取得先選択部 2 2 3 が、取得キー k d と取得先ノード D B 2 3 3 に記憶されているルータ I D とを比較して、取得キー k d に距離の近い順に k 個のルータを選択する。ステップ S 1 0 0 3 では、取得先選択部 2 2 3 が、選択したルータに、取得キー k d に距離の近い k 個のルータに係る情報と、取得キー k d に一致する検索キー k r が存在すれば、その検索キー k r を含む経路計算結果情報 (k r , v a l u e) の転送経路情報 v a l u e とを取得するための要求 f i n d _ v a l u e (k d) を送信する。

【 0 0 7 2 】

ステップ S 1 0 0 4 では、転送経路決定部 2 2 2 は、f i n d _ v a l u e (k d) の送信先のルータ 2 0 から、取得キー k d に距離の近い k 個のルータ 2 0 に係る情報と、取得キー k d に一致する検索キー k r が存在すれば、その検索キー k r を含む経路計算結果情報 (k r , v a l u e) の転送経路情報 v a l u e とを取得する。なお、f i n d _ v a l u e (k d) を受信したルータ 2 0 の取得先選択部 2 2 3 は、取得キー k d と r a n g e (図 7 参照) の境界のルータ I D とを比較して、取得キー k d を含む r a n g e を抽出し、取得キー k d と抽出した r a n g e に属するルータ 2 0 それぞれのルータ I D との距離を算出し、距離の近い順に k 個の取得先のルータに係る情報を、送信元のルータ 2 0 に返信する。

【 0 0 7 3 】

ステップ S 1 0 0 5 では、転送経路決定部 2 2 2 が、転送経路情報 v a l u e が取得さ

10

20

30

40

50

れたか否かを判定する。そして、転送経路情報 `value` が取得された場合 (ステップ S 1005 で `Yes`)、ステップ S 1006 では、転送経路決定部 222 が、転送経路情報 `value` に含まれる転送経路 (式 (3) における `NH_*`) と、現在使われている転送経路とのどちらが、転送すべきパケットの IP アドレスと最長一致かを比較し、最長一致の方の転送経路を現在使われている転送経路に設定する。ただし、ステップ S 1005 で取得された転送経路情報 `value` に含まれる転送経路および現在使われている転送経路のいずれもが、転送すべきパケットの IP アドレスとの距離が同じ場合には、転送経路情報 `value` のバージョン番号が最新の方に、現在使われている転送経路を設定する。そして、ステップ S 1007 では、転送経路決定部 222 が、現在使われている転送経路を不図示の `FIB` に設定する。そして、処理は、ステップ S 1008 へ進む。なお、ステップ S 1005 で、転送経路情報 `value` が取得されなかった場合 (ステップ S 1005 で `No`)、処理は、ステップ S 1008 へスキップする。

【0074】

ステップ S 1008 では、転送経路決定部 222 が、`find_value(kd)` を、未送信のステップ S 1004 で取得したルータがあるか否かを判定する。未送信のステップ S 1004 で取得したルータがある場合 (ステップ S 1008 で `Yes`)、ステップ S 1009 では、取得先選択部 223 が、未送信のルータ 20 へ、`find_value(kd)` を送信する。そして、処理は、ステップ S 1004 へ戻る。また、未送信のステップ S 1004 で取得したルータがない場合 (ステップ S 1008 で `No`)、処理は、ステップ S 1010 へ進む。

【0075】

ステップ S 1010 では、転送経路決定部 222 が、取得キー `kd` の IP アドレス部の最後 (一番右) のゼロでないビットをゼロに変更し、同時にマスク値部分の当該 IP アドレスの変更位置のビットを含む右側のビットをゼロに変更し、その変更したビット列を新たな取得キー `kd` に設定する。この取得キー `kd` のビット変更時の操作の具体例について、図 11 を用いて説明する。図 11 に示すように、ビット変更前の IP アドレス部のビット列およびマスク値部分のビット列がそれぞれ `<10101100>` `<11111111>` であったとする。IP アドレス部の最後 (一番右) のゼロでないビットは、右から 3 番目である。この右から 3 番目のビットをゼロに変更する。そして、マスク値部分については、右から 3 番目のビットを含む右側の 3 つのビットをゼロに変更する。

【0076】

ステップ S 1011 では、転送経路決定部 222 が、取得キー `kd` の IP アドレス部のホスト部分がすべてゼロかそれ以外かを判定する。取得キー `kd` の IP アドレス部のホスト部分の一部がゼロでない場合 (ステップ S 1011 で `No`)、処理は、ステップ S 1002 へ戻る。また、取得キー `kd` の IP アドレス部のホスト部分がすべてゼロの場合 (ステップ S 1011 で `Yes`)、ステップ S 1012 では、転送経路決定部 222 が、現在使われている転送経路を含む経路計算結果情報のレプリケーションを、当初の取得キー `kd` に $(k+1)$ 番目に距離の近いルータ ID のルータに対して実行する。そして、転送経路情報の取得処理を終了する。

【0077】

なお、転送すべきパケットは、ステップ S 1007 で `FIB` に設定された転送経路に基づいて転送される。この際、カプセル化処理部 224 は、転送すべきパケットを、出口の境界ルータの IP アドレスを宛先とするヘッダによってカプセル化して、`IGP` を用いて、そのパケットを転送する。このカプセル化の効果は、ルータ 20 が他の `POP` の `BGP` 経路情報を保持する必要がなくなること、および `AS` 内で経由するルータ 20 が意図しない方路へパケットを転送することを防止できること、である。

【0078】

また、ステップ S 1010 において取得キー `kd` の IP アドレス部のビット列を、最後 (一番右) のゼロでないビットをゼロに変更し、再びステップ S 1002 へ戻って検索を続ける理由は、経路計算結果情報 (`kr, value`) は、宛先 IP アドレスの集合であ

10

20

30

40

50

るプリフィクスに基づいて管理されているのに対して、取得キー k_d は宛先 IP アドレスを用いて生成されているからである。すなわち、取得キー k_d にそのまま一致する検索キー k_r が見つかる可能性は低いと考えられる。そこで、IP アドレス部のビット列を、最後（一番右）のゼロでないビットをゼロに変更しつつ、最長一致となるビット列を探索していくためである。

【0079】

また、ステップ S1006 において、最長一致となる転送経路が複数見つかった場合には、転送経路情報 `value` に含まれるバージョン番号が最新の方が選ばれる。また、ステップ S1012 において、レプリケーションする経路計算結果情報は、 $k+1$ 番目に距離の近いルータ 20 に格納されるが、 $k+1$ 番目のルータ 20 にすでに格納済みであれば、 $k+2$ 番目のルータにレプリケーションされる。そして、レプリケーションされた経路計算結果情報は、所定時間が経過した後（タイムアウト後）、消去される。なお、その所定時間は、取得キー k_d とルータ ID との距離の大きさに反比例して設定される。すなわち、タイムアウトの所定時間は、距離が遠いほど短く設定される。

【0080】

（記憶容量の削減効果）

前記したように、本実施形態では、ルータ 20 が記憶していた外部経路情報がルートサーバ 10 に記憶されるため、従来のルータに比べて、記憶容量を低減することが可能である。また、ルートサーバ 10 も、プリフィクスを分担して記憶する台数を増加することによって、一台当りの記憶容量を低減することができる。さらに、従来の各ルータでは、パケット転送のために必要な転送経路の情報（FIB に記憶される情報）については、すべてのプリフィクスに対して記憶しなければならなかったが、本実施形態では、POP 内で分散管理することによって、ルータ 20 の記憶容量を低減することが可能である。定量的には、ルータ 20 一台当りの保持する経路数は、全経路数を POP 内のルータ数で除算した値に、レプリケーションされる情報の経路数を加算して算出される。なお、インターネットにおいては、トラヒックのほとんどは少数の宛先に集中する傾向があるため、レプリケーションされる経路数は、全経路数を POP 内のルータ数で除算した値に比較して小さい。

【0081】

（メッセージ数および計算処理量の削減効果）

ルートサーバ 10 が経路計算を実行するため、各ルータ 20 は、経路計算を実行しないで良い。そのため、本実施形態のルータ 20 では、従来のルータに比較して、計算処理量が低減できる。さらに、本実施形態では、AS 内のルータ 20 は、複数の POP のいずれかに属している。そして、POP 内では BGP 経路については同一の転送経路となるので、ルートサーバ 10 は、POP ごとに 1 つのルータ 20 のための経路計算を実行すれば良く、また、その経路計算結果をメッセージとして POP 内のルータ数より少ない k 個のルータ 20 に送信すれば良い。したがって、ルートサーバ 10 が、ルータ 20 ごとにすべての経路計算結果を送信する場合に比べて、定性的には、メッセージ数および計算処理量（処理負荷）を低減することができる。

【0082】

定量的には、POP 内のルータ数の平均を R_{pop} 、ルートサーバ 10 からすべてのルータ 20 に経路計算結果（メッセージ）を送信する場合のメッセージ数を M とすると、経路計算結果の配信のために、ルートサーバ 10 が送信するメッセージ数 R_m は、 $M / R_{pop} + N_f$ となる。ここで、 N_f は、`find_node` メッセージ数である。この N_f は、初期状態では、最大で POP 内のルータ数分になることもあるが、`find_node` が繰り返し替えされることによって、ルートサーバ 10 の格納先ノード DB132 が更新され、ルートサーバ 10 が知り得ないルータ 20 が無くなっていくため、最終的には、 k （冗長性）に収束する。そのため、メッセージ数は、削減される。

【0083】

また、`find_value` メッセージ数は、転送すべきパケットの宛先の分布に依存

する。インターネットでは、パケットの宛先のほとんどが、ごく一部のプリフィクスに集中することから、レプリケーションによって、他のルータに取得要求を行う必要がなくなっていく。したがって、`find_value`メッセージ数は、ほとんどゼロに収束する。したがって、本実施形態では、POPおよびレプリケーションの導入にともなって、メッセージ数を低減することができる。

【0084】

以上、本実施形態の経路情報管理システム30によれば、ルートサーバ10がプリフィクスに係る外部経路情報を格納するため、ルータ20に保持すべき経路数(記憶容量)を削減することができる。また、ルートサーバ10は、同一のプリフィクスに係る外部経路情報についてすべてを収集して経路計算を実行するため、経路計算結果が最適とならないという問題は発生しない。また、AS内に複数のPOPを設けた場合、そのPOP内のBGP経路はプリフィクスに対して共通になるため、ルートサーバ10は、POP内の1つのルータ20に対してだけ経路計算を実行すれば良くなる。そのため、ルートサーバ10の処理負荷を低減できるとともに、ルータ20での経路計算が不要となる。また、ルートサーバ10によって算出されたPOPごとの経路計算の結果は、POP内で分散管理されるため、POPごとにk個(所定数)のルータ20に格納(送信)されれば良い。そのため、経路計算の結果をメッセージとして送信するときのメッセージ数を削減することができる。また、POP内のルータ20によって経路計算結果情報が分散管理されるため、一台当りのルータ20の記憶容量は、すべての経路計算結果情報を記憶していた従来の一台のルータ20の記憶容量に比べて、低減できる。また、必要な経路計算結果の取得のための探索は、POP内のルータ20だけに限られるため、その取得処理が早くなる。

【0085】

なお、本実施形態は、これらに限定されるものではなく、その趣旨を変えない範囲で実施することができる。例えば、取得先選択部223が検索キー`kr`や取得キー`kd`に距離の近いk個のルータ20を選択するとき、本実施形態では検索キー`kr`や取得キー`kd`を含む`range`で選択されているk個のルータを取得するものとして説明した。しかし、ルータ20の記憶部23に記憶されているルータIDと検索キー`kr`や取得キー`kd`との距離を計算して、最も距離の近いk個のルータ20を選択するようにしても構わない。また、図7では、一例として、記憶部23の取得先ノードDB233に`range`ごとにk個のルータに係る情報を格納する場合を示した。しかし、`range`ごとにk個より多くのルータ20に係る情報が記憶されている場合には、その中から最も距離の近いk個のルータ20を選択するようにしても構わない。

【0086】

また、転送経路決定部222が取得キー`kd`に一致する検索キー`kr`を探索するとき、経路計算結果情報(`kr, value`)を`range`ごとに紐付けて記憶されていれば、探索範囲を絞ることができ、探索時間を小さくすることができる。また、ルータ20は、`find_value(kd)`によって取得した転送経路情報`value`を、自身の記憶部23の図示しないキャッシュメモリに格納しておくこと、探索を早めることができる。なお、キャッシュメモリに格納された情報は、所定時間後(タイムアウト後)に消去される。なお、レプリケーションやキャッシュは、転送経路情報を格納する1つのRadix Treeによって管理される。

【0087】

また、ルータ20が、ルートサーバ10の機能を兼ねていても構わない。また、本実施形態は、AS間の通信プロトコルであるeBGPを改変するものではないので、AS内に閉じて適用することができる。そして、本実施形態を適用したASは、従来通り、他のASと接続することが可能である。

【0088】

また、本実施形態において、ルートサーバ10(図2参照)の各部11, 12の処理は、ルートサーバ10をコンピュータで実現したときに搭載されるプログラムによって実現

されてもよい。また、本実施形態において、ルータ（図3参照）の各部21, 22の処理は、ルータ20をコンピュータで実現したときに搭載されるプログラムによって実現されてもよい。このプログラムは、通信回線を介して提供することもできるし、CD-ROM等のコンピュータ読み取り可能な記録媒体に書き込んで配布することも可能である。

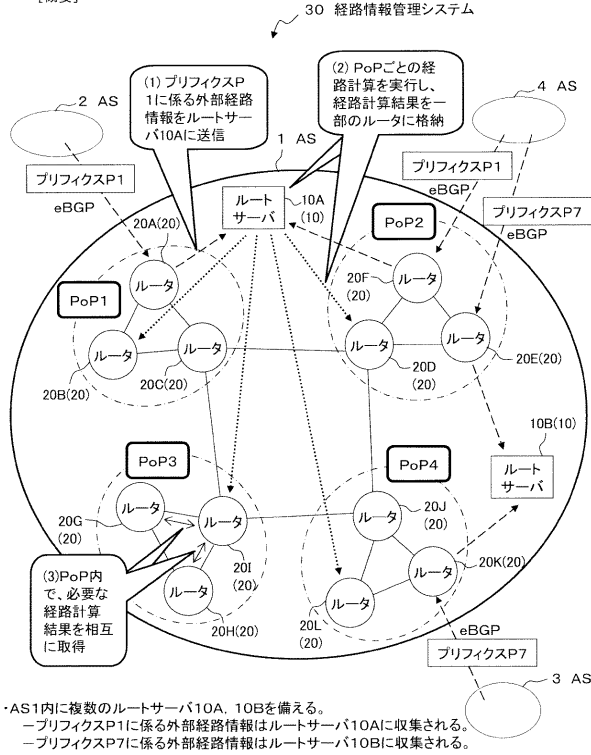
【符号の説明】

【0089】

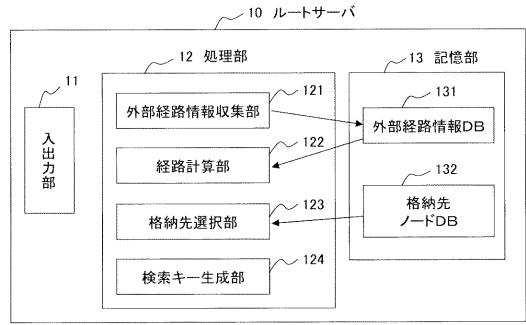
| | | |
|-----|------------------|----|
| 1 | AS | |
| 10 | ルートサーバ（経路情報管理装置） | |
| 11 | 入出力部 | |
| 12 | 処理部 | 10 |
| 13 | 記憶部 | |
| 20 | ルータ（境界ルータ、ルータ） | |
| 22 | 処理部 | |
| 23 | 記憶部 | |
| 30 | 経路情報管理システム | |
| 121 | 外部経路情報収集部 | |
| 122 | 経路計算部 | |
| 123 | 格納先選択部 | |
| 124 | 検索キー生成部 | |
| 131 | 外部経路情報DB | 20 |
| 132 | 格納先ノードDB | |
| 221 | ルートサーバ選択部 | |
| 222 | 転送経路決定部 | |
| 223 | 取得先選択部 | |
| 224 | カプセル化処理部 | |
| 225 | 取得キー生成部 | |
| 231 | ルートサーバIDDB | |
| 232 | 経路計算結果DB | |
| 233 | 取得先ノードDB | |

【図1】

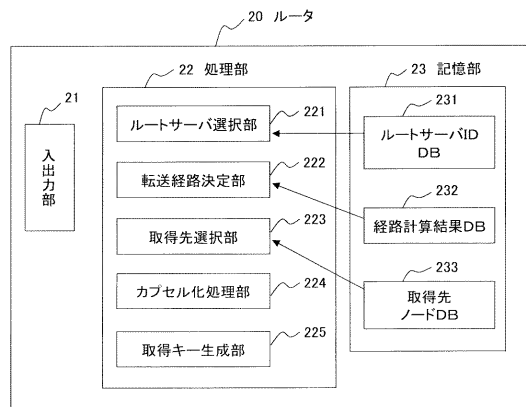
【概要】



【図2】

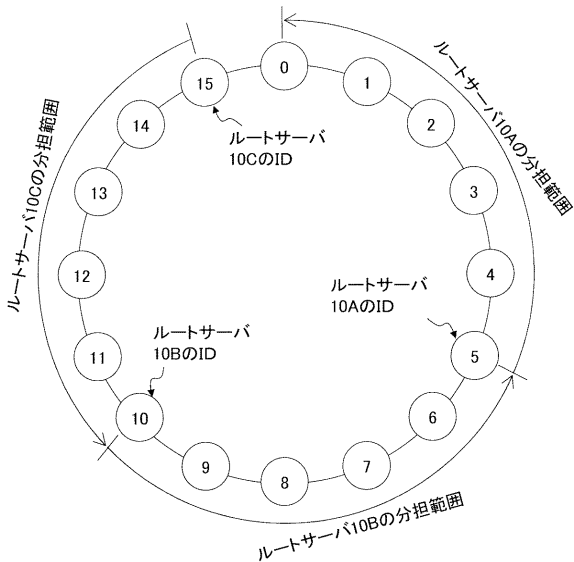


【図3】



【図4】

【ルートサーバ10の分担範囲(ハッシュ空間表示)】



- 内の数字はルートサーバのID番号を表す。
- ・ルートサーバ10Aの分担範囲: 0超え~5以下
- ・ルートサーバ10Bの分担範囲: 5超え~10以下
- ・ルートサーバ10Cの分担範囲: 10超え~15以下

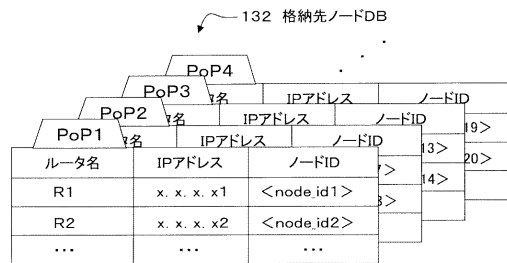
【図5】

【距離の定義】ノードIDと検索キーとのビット列のExclusive ORを演算し、Exclusive ORの演算結果のビット列において、左から見て最初に「1」が出現した位置を右からの数えた桁数の数値

(例)

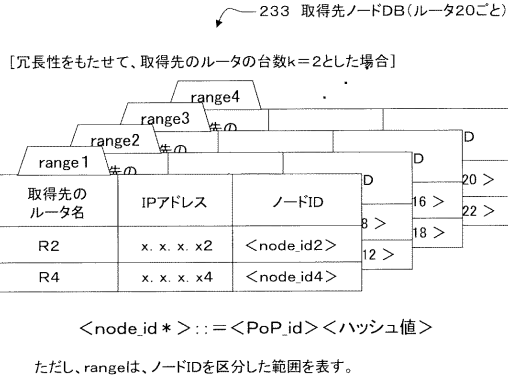
| ルータ | ノードID | 検索キー | Exclusive OR | 距離 | 近さの順位 |
|-----|-------|-------|--------------|----|-------|
| a | 10011 | 10111 | 00100 | 3 | 3 |
| b | 10101 | | 00010 | 2 | 2 |
| c | 10110 | | 00001 | 1 | 1 |
| d | 11000 | | 01111 | 4 | 4 |

【図6】

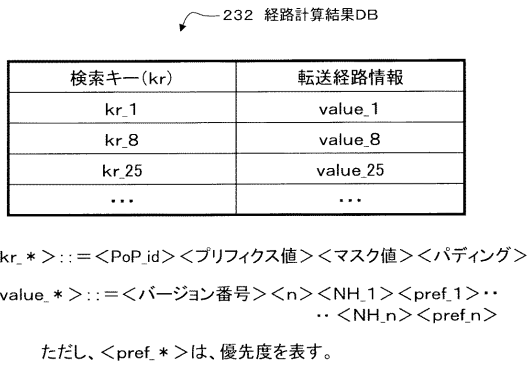


<node.id*> ::= <PoP.id> <ハッシュ値>

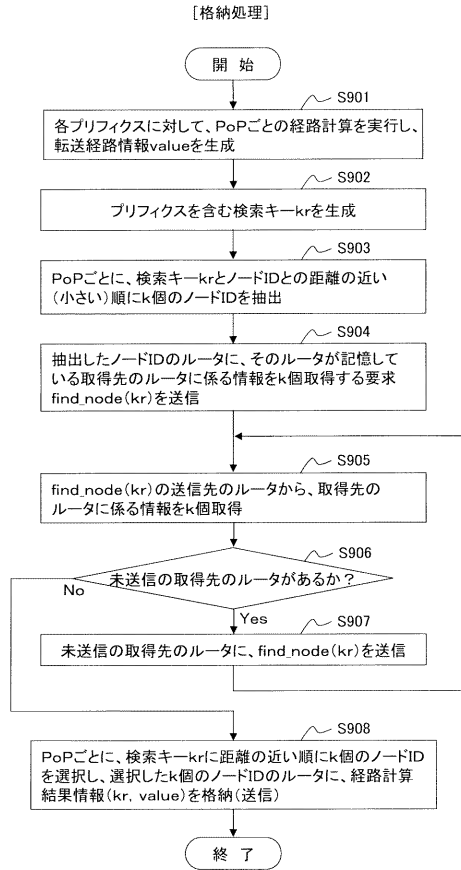
【 図 7 】



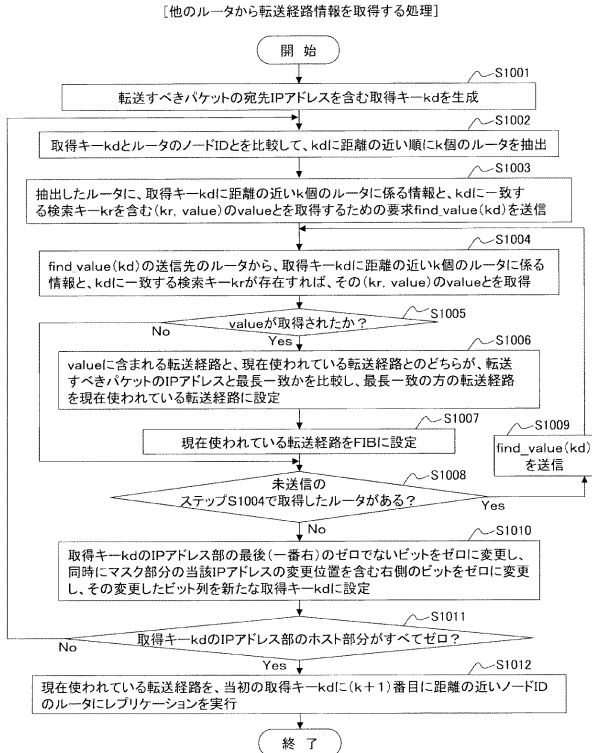
【 図 8 】



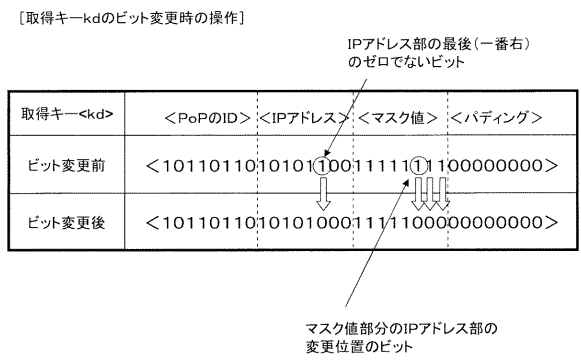
【 図 9 】



【 図 10 】



【 図 11 】



フロントページの続き

審査官 安藤 一道

(56)参考文献 特開2006-270781(JP,A)
特開2004-260655(JP,A)
特開2003-143143(JP,A)
特開2010-199800(JP,A)

(58)調査した分野(Int.Cl., DB名)
H04L 12/56