

# OSPF を用いたEBGP 接続障害の高速復旧手法 に関する提案

渡里 雅史<sup>†</sup> 屏 雄一郎<sup>†</sup> 阿野 茂浩<sup>†</sup> 山崎 克之<sup>††</sup>

<sup>†</sup> 株式会社 KDDI 研究所 <sup>††</sup> 長岡技術科学大学

Border Gateway Protocol 4 (BGP) を用いた Autonomous System (AS) 間の経路制御では、安定した接続性の確保のため、各 AS は隣接する他の AS と予め複数の代替経路を確保する。しかしながら、実際の EBGP 経路障害では、障害検知から復旧までに秒オーダーのダウンタイムを伴うと共に、大量発生する経路制御メッセージがルータの処理能力を圧迫し、ネットワーク全体の安定性を著しく低下させる問題がある。本問題を解決するため、筆者らは、ネットワーク全体の安定性を低下することなく EBGP 接続障害から高速に復旧させる手法を提案する。提案手法は、障害に伴う経路変動の伝播範囲を特定のルータに限定させた上で、経路収束の速い Open Shortest Path First (OSPF) を用いて代替経路を確保することを特徴とする。本稿では、提案手法の詳細設計ならびに汎用 PC を用いたプラットフォームでのプロトタイプ実装とその評価について述べる。

## Proposal of Fast EBGP Link Restoration Method using OSPF

Masafumi WATARI<sup>†</sup> Yuichiro HEI<sup>†</sup> Shigehiro ANO<sup>†</sup> Katsuyuki YAMAZAKI<sup>††</sup>

<sup>†</sup> KDDI R&D Laboratories Inc.

<sup>††</sup> Nagaoka University of Technology

The Border Gateway Protocol 4 (BGP) allows an Autonomous System (AS) to establish multiple external BGP connections with other ASes to ensure availability of alternative paths. However, recovery through these connections upon an external BGP link failure causes seconds of downtime. The BGP Update messages produced by the failure also heavily stresses the BGP speakers making the whole AS unstable. This paper presents a novel approach of providing fast restoration for external BGP link failures using OSPF. The method maintains stability of the network during restoration by propagating the failure to limited routers. The detail extensions made to BGP/OSPF and its evaluation through prototype implementation are presented.

### 1 はじめに

近年、IP 電話をはじめ、オンラインゲーム、株取引、企業向け IP-VPN など、安定性およびリアルタイム性を重視する通信サービスが広く普及しつつある。これらの通信サービスは、IP 転送品質劣化の影響を受け易いため、サービス品質の保証には、IP パケット単位の品質（遅延、ジッタ、ロスなど）などの確保に加え、一定の品質を保つため、安定した接続性の確保が重要である。特にこれらの通信サービスの利用形態には、ユーザ同士またはサーバが異なる管理ドメインに接続する場合が想定されるため、ドメイン間での安定した接続性確保が重要となる。

インターネットにおけるドメイン間の経路制御には、Border Gateway Protocol 4 (BGP)<sup>1)</sup> が広く利用されている。BGP におけるドメインは Autonomous System (AS) と呼ばれ、各 AS は同一または異なる他の AS と予め複数の External BGP (EBGP) 接続を確立することで、より信頼性の高い接続性を確保できる。EBGP 接続障害時には、別の EBGP 接続に切り替えることで復旧する。しかしながら、実際の EBGP 接続障害では、障害発生から復旧まで秒オーダーのダウンタイムを伴う。また、大量発生する経路制御メッセージがルータの処理能力を圧迫し<sup>2)</sup>、ネットワーク全体の安定性を著しく低下させる問題がある<sup>3)</sup>。

そこで、これらの課題を解決するため、本稿では、ネットワーク全体の安定性を低下することなくEBGP接続障害から高速に復旧する手法を提案する。提案手法は、障害に伴う経路変動の伝播範囲を特定のルータに限定させた上で、経路収束の速いOpen Shortest Path First (OSPF)<sup>4)</sup>を用いて代替経路を確保することを特徴とする。本稿では、提案手法の詳細設計ならびに汎用PCを用いたプラットフォーム上でのプロトタイプ実装とその評価について述べる。

本論文の構成は以下の通りである。2章でEBGP接続障害の現状について説明し、3章で関連研究について述べる。4章で提案する高速切替手法について述べ、5章で詳細設計について説明する。6章で提案方式の動作検証ならびに切替時間に関する評価結果を示し、7章で結論を提示する。

## 2 EBGP 接続障害の現状

EBGP接続障害時における復旧プロセスは、大きく二つに分けられる。一つは、障害発生から障害検知である。BGPでは、通常、Keepaliveメッセージを用いた接続状態の監視とホールドタイムを用いた障害検知を行う。障害検知に掛かる時間は、本メッセージの送信間隔とホールドタイムの値に大きく依存する。多くのBGPの実装では、最小3秒のホールドタイムが指定可能である。しかしながら、リアルタイム性を要するアプリケーションでは許容できないパケットロスに伴うため、近年では、BGPとは独立の接続状態監視プロトコルとして、Bidirectional Forwarding Detection (BFD)<sup>5)</sup>技術の標準化が進められている。BFDでは、独自のKeepaliveメッセージを用いた接続状態の監視により、数10ミリ秒から数100ミリ秒での高速な障害検知が可能となるため、今後は標準機能として広く実装されていく可能性がある。

二つ目のプロセスは、障害検知から障害復旧である。BGPでは、障害を検知すると最適経路の再計算を行い、代替経路が存在する場合は、Announceメッセージを用いて広報し、存在しない場合は、Withdrawメッセージを用いて障害経路の取り下げを行う。バックボーンルータでは、フルルートと呼ばれるインターネット全体の経路を扱うため、EBGP接続障害時には多くのメッセージの生成と伝搬が必要となる。図1は、Quagga Routing Software Suite (Quagga)<sup>6)</sup>のBGP実装を用い

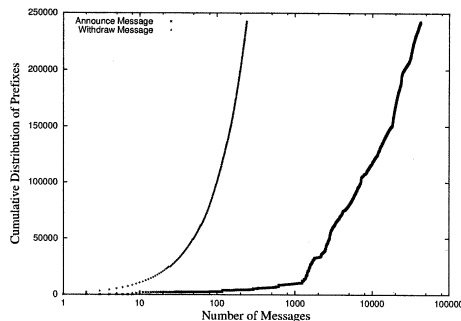


図.1 フルルート交換時の受信メッセージ数と累積プレフィックス数の関係

て、Internal BGP (IBGP) スピーカ間でフルルートの受信に必要としたメッセージ数とその累積プレフィックス数の関係を示す。計測は、Route Views Project<sup>7)</sup>のroute-views2.oregon-ix.netで観測した242,397<sup>1)</sup>のプレフィックスとASパスリストをエミュレーションした上で行った。全プレフィックスの取り下げには238個のWithdrawメッセージを要し、同プレフィックスの広報には42,411個のAnnounceメッセージを要した。また、この時のtcpdumpによる全メッセージの伝播時間は、それぞれ0.809秒と5.169秒であった。

これらの値は、BGPの実装やルータの処理性能によって異なるが、経路数の増加に伴い、今後より多くのメッセージ数が必要となることが想定される。特にAS内にRoute Reflector (RR)<sup>8)</sup>を設置している場合は、RRへWithdrawメッセージが伝播された後、RRから各クライアントに対してAnnounceメッセージが伝播されるため、メッセージ伝播だけにさらに多くの時間を要する。実際には、RRと接続するクライアント数およびBGPセッション数に応じて処理が増加し、最適経路の再計算と経路表の更新を伴うため、復旧までにより多くの時間を要する。また、これらのメッセージ伝播時は、各ルータの負荷を高めるため、AS内の安定性を一時的に低下させる問題もある。

本稿では、この障害検知から障害復旧までのプロセスに着眼し、上述の問題を解決する新たな復

<sup>1</sup> 2008年3月1日00時00分のRIBデータを使用。  
<http://archive.routeviews.org/bgpdata/2008.03/RIBS/rib.20080301.0000.bz2>

旧手法について提案する。

### 3 関連研究

EBGP 接続障害時の復旧時間短縮には、交換する Update メッセージ数を減らす方法が考えられる<sup>9, 10)</sup>。例えば、フルルートの取り下げであれば、障害したネクストホップの IP アドレスのみを通知し、受信ルータ側で、同ネクストホップに一致するすべての経路を取り下げる方法が考えられる。この場合では、実質 1 つの Withdraw メッセージでフルルートの取り下げが可能となる。しかし、実際のルータでは、メッセージ受信後に BGP Routing Information Base (RIB) の更新、最適経路の再計算、Forwarding Information Base (FIB) の更新など、経路数に応じて多くの処理を伴う。そのため、例えば 1 つの RIB/FIB エントリの更新に 146  $\mu$  秒<sup>11)</sup>を要する場合は、24 万エントリの更新には約 35 秒を要する。このため、障害復旧時間の短縮には、交換するメッセージ数の削減に加え、更新する RIB/FIB エントリ数を減らす工夫が必要となる。

この問題に対して Bonaventure et al.<sup>12)</sup>は、FIB 内にネクストホップテーブルを導入し、更新する FIB エントリ数を減らす方法を提案している。本テーブルは、各ネクストホップに対して、予めプライマリとバックアップの二つの出力インタフェースを設定し、プライマリの接続状態に応じてバックアップへの切替えを行う。障害時は、本テーブル内の 1 エントリの更新のみとなるため、障害検知から 50 ミリ秒以内での復旧が可能となる。しかしながら、AS 内に RR が介在する場合は、RR と各クライアントに対して通常の BGP におけるフルルート交換が必要となる。また、プライマリとバックアップへの切り替えは、本テーブル内のフラグに基づいて行われるため、設定可能なバックアップインタフェースが一つに限られる。複数のバックアップに対応するためには、FIB 構造のさらなる拡張、各バックアップの接続性監視、優先順位付けなど、より複雑な仕組みが必要となる。

### 4 IGP を用いた高速障害復旧手法

EBGP 接続障害は、ネクストホップへの到達性障害であると考えられる。一方、AS におけるネクストホップへの到達性は、Interior Gateway Protocol (IGP) によって確保されるため、ネクストホップ

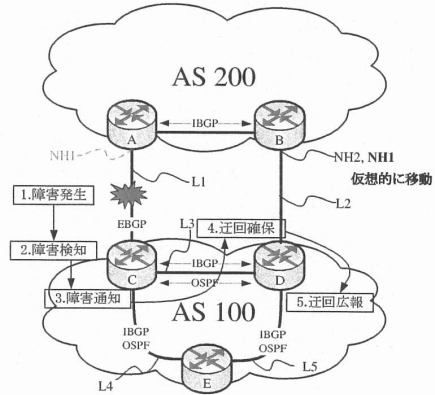


図. 2 提案手法の概要

への到達性障害は、IGP における経路障害とも考えられる。このため、IGP により障害箇所を迂回し、ネクストホップへの到達性を再確保することで、EBGP 接続障害の復旧が可能となる。そこで、本稿では、同一 AS へ複数の EBGP 接続がある構成を前提に、EBGP 接続障害からの復旧手法として、IGP を用いてネクストホップへの到達性を再確保する手法を提案する。提案手法は、同じネクストホップへの到達性を復旧させるため、BGP のフルルート交換が必要なく、また、IGP を用いて単一の FIB エントリのみ更新するため、高速かつ安定した障害復旧が可能となる。

図 2 に提案手法の概要を示す。AS100 および AS200 における AS 間の接続は、EBGP であり、AS 内は IBGP および IGP として OSPF を想定する。AS100 から AS200 へのネクストホップは、それぞれ NH1 および NH2 とする。ここで、リンク L1 にて障害が発生した場合、通常の BGP では、ルータ C において障害を検知すると共に、AS 内の BGP ルータ D および E に対して Withdraw メッセージを伝播する。一方、提案手法では、ルータ C において障害を検知すると、まず EBGP 接続を有するルータ D に対して、障害を通知する。本通知を受信したルータ D では、同 AS へのリンクである L2 を用いて、仮想的に NH1 への到達性を確保する。すなわち、ルータ C に対して NH1 への到達性を広報すると共に、自身の NH1 へ向けの経路を NH2 に設定する。これにより、ルータ C は障

Router C (Before)			Router D (Before)		
Destination	NextHop		Destination	NextHop	
BGP Prefix X	NH1		BGP Prefix X	NH2	
OSPF NH1	Connected (L1)		OSPF NH1	Router C	Connected (L3)
OSPF Router D	Connected (L3)		OSPF Router C	NH2	Connected (L2)
OSPF NH2	Router D		OSPF NH2	Connected (L5)	
OSPF Router E	Connected (L4)		OSPF Router E	Connected (L5)	

Router C (After)			Router D (After)		
Destination	NextHop		Destination	NextHop	
BGP Prefix X	NH1		BGP Prefix X	NH2	
OSPF NH1	Router D		OSPF NH1	Connected (L2)	
OSPF Router D	Connected (L3)		OSPF Router C	Connected (L3)	
OSPF NH2	Router D		OSPF NH2	Connected (L2)	
OSPF Router E	Connected (L4)		OSPF Router E	Connected (L5)	

図. 3 障害前後における経路表

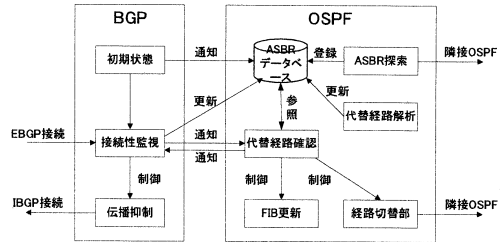


図. 4 機能構成

害した NH1 への到達性をルータ D 経由で確保できると判断し、BGP における Withdraw メッセージの送信を抑制し、OSPF を用いて迂回経路を確保する。

図 3 に、障害前後におけるルータ C および D の経路表を示す。ルータ C における障害前の Prefix X へのネクストホップは NH1 であるが、障害後のルータ C におけるネクストホップはルータ D となる。ルータ C は、FIB にて NH1 を recursive に検索することで、宛先ネクストホップを解決する。また、ルータ D は、NH1 へのネクストホップを自身のリンクである L2 に向けることで、AS200 へのパケット転送が可能となる。

## 5 システム設計

本章では、提案手法に関わる詳細設計について述べる。図 4 に提案手法の主要機能構成を示す。BGP では、EBGP 接続状態の監視と OSPF への通知、また、OSPF による復旧時は、Withdraw メッセージ伝播の抑制を行う。OSPF では、BGP からの通知に基づき、新たなデータベースの管理と迂回経路の確保、また、隣接 OSPF ルータとネクストホップ情報の交換を行う。

### 5.1 代替経路の管理

本稿では、EBGP 接続を持つルータを AS Border Router (ASBR) と呼ぶ。ASBR は、EBGP 接続障害時の代替経路を確保するため、AS 内の各 ASBR と予めネクストホップ情報を共有し、ASBR データベースに管理する。ネクストホップ情報の共有には、ASBR の探索を兼ねるため、OSPF を拡張し、新たに代替経路通知 LSA を定義した。各 ASBR は、OSPF 起動時にネイバに対して本 LSA をフラッシングし、EBGP 接続のネクストホップア

表 1 ASBR データベース

AS 番号	ネクストホップ	ルータ ID	Pref.
200	172.0.1.1	10.0.0.1	20
200	172.0.2.1	10.0.0.2	10
300	172.0.3.1	10.0.0.3	30
300	172.0.4.1	10.0.0.4	10

ドレスと AS 番号を通知する。本 LSA を受信した ASBR では、送信元のルータ ID と共にこれらの情報を ASBR データベースに格納する。表 1 に本データベースの例を示す。ASBR データベースには、使用する代替経路の優先順位を決めるため、プレファレンス値を設けた。

なお、本 LSA は、通常ルータとの互換性を維持しながら AS 内の全 ASBR に通知するため、Opaque LSA<sup>13)</sup> を用いた。通常ルータで特別な処理をすることなく、全 ASBR への通知が可能となる。

### 5.2 提案手法に関わる BGP 拡張

通常の BGP における Finite State Machine (FSM)<sup>1)</sup> では、ホールドタイムのタイムアウトに伴い、Established から Idle 状態へ遷移し、Withdraw メッセージの伝播を行う。一方、提案手法では、IGP による障害復旧とするため、本 FSM を拡張し、新たに Fallback 状態を定義した。障害を検知した BGP ルータは、Fallback 状態へ遷移し、OSPF に対して EBGP 接続障害の通知を行う。OSPF にて代替経路を確保できない場合は、通常の BGP における Idle 状態へと遷移する。OSPF において代替経路が確保できた場合は、Withdraw メッセージ伝播を抑制し、復旧は OSPF に委ねる。また、障害箇所が復旧した場合は、Fallback 状態から通常の BGP における OpenSent 状態へと遷移し、OSPF



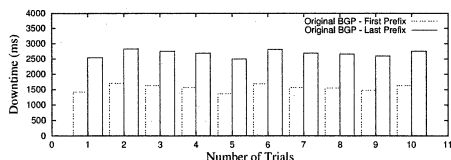


図. 7 通常の BGP におけるダウンタイム

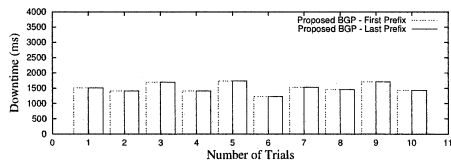


図. 8 提案手法におけるダウンタイム

提案手法では、最後のプレフィックスの復旧時間は平均 15.10 秒であった。また、最初と最後のプレフィックス間での差は 1 パケットであり、提案手法ではプレフィックス数に依存しない障害復旧が可能であることを確認した。なお、回復時はどちらもパケットロスがないことを確認した。

提案手法の切り替え時間を計測するため、ルータ C の障害検知から FIB 更新完了までの時間  $T_1$ 、ルータ D の LSA 受信から FIB 更新完了までの時間  $T_2$  を求めた。また、同様に回復時において、ルータ C の回復検知から FIB 更新完了までの時間  $T_1$ 、ルータ D の LSA 受信から FIB 更新完了までの時間  $T_2$  を求めた。表 2 は、それぞれ 10 回の平均を示す。障害時の復旧は、ルータ間での LSA 伝播時間を除き平均 2.3 ミリ秒であった。回復時は、ルータ D の FIB 更新に実装上の問題と考えられる多くの時間を要している。ただし、ルータ C の切り替えのタイミングで復旧するため、上述の通りパケットロスがないことを確認した。

また、FIB 更新に伴うネクストホップの recursive lookup に要した時間を計測した結果、CPU クロック数にして 50 回の平均が 248,680 クロックであり、時間にしておよそ 117  $\mu$  秒と十分に小さいことを確認した。また、2 回目以降は、キャッシュを参照するため、転送するパケット毎に recursive lookup の必要がなく、データトラヒックへの影響が少ないことを確認した。

表 2 経路切り替え時間 (ミリ秒)

	$T_1$	$T_2$	$T_1 + T_2$
接続障害時	2.1	0.2	2.3
接続回復時	2.1	221.3	223.4

## 7 今後とまとめ

本稿では、EBGP 接続障害からの復旧時間を短縮するため、IGP を用いて代替経路を確保する新たな手法を提案した。提案手法は、フルルートの交換が必要なく、また、特定のルータ間で制御するため、AS 内の安定性を維持した復旧が可能となる。実インターネットの約 24 万経路をエミュレーションした環境下で、障害検知から 3 ミリ秒以内の切り替えが可能であることを確認した。

今後は、対向 AS 双方への提案手法の導入、また、BFD を導入した場合の障害発生から復旧までの全ダウンタイムを計測し、提案手法の有用性評価を行う。

## 参考文献

- 1) Y. Rekhter, T. Li, and S. Hares, "A Border Gateway Protocol 4 (BGP-4)," RFC 4271, January 2006.
- 2) D. Meyers, L. Zhang, and K. Falls, "Report from the IAB Workshop on Routing and Addressing," RFC 4984, September 2007.
- 3) L. Wang, X. Zhao, D. Pei, R. Bush, D. Massey, A. Mankin, S. Wu, and L. Zhang, "Observation and analysis of BGP behavior under stress," ACM SIGCOMM Workshop on Internet Measurement, November 2002.
- 4) J. Moy, "OSPF Version2," RFC 2178, April 1998.
- 5) D. Katz and D. Ward, "Bidirectional Forwarding Detection," Internet Draft (Work in Progress), March 2008.
- 6) Quagga Routing Software Suite <http://www.quagga.net/>
- 7) Route Views Project <http://www.routeviews.org/>
- 8) T. Bates, E. Chen, and R. Chandra, "BGP Route Reflection: An Alternative to Full Mesh Internal BGP (IBGP)," RFC 4456, April 2006.
- 9) M. Bhatia, J. Halpern, and P. Jakma, "Advertising Multiple NextHop Routes in BGP," Internet Draft (Work in Progress), August 2006.
- 10) D. Pei, M. Azuma, N. Ngyuyen, J. Chen, D. Massey, and L. Zhang, "BGP-RCN: Improving BGP convergence through Root Cause Notification," Computer Networks and ISDN Systems, June 2005.
- 11) P. Francois, C. Filsfils, J. Evans, and O Bonaventure, "Achieving sub-second IGP convergence in large IP networks," ACM SIGCOMM Computer Communication Review, July 2005.
- 12) O. Bonaventure, C. Filsfils, and P. Francois, "Achieving Sub-50 Milliseconds Recovery Upon BGP Peering Link Failures," ACM CoNext'05, October 2005.
- 13) R. Coltun, "The OSPF Opaque LSA Option," RFC 2178, April 1998.