

BGP と障害と運用と

Matsuzaki 'maz' Yoshinobu

<maz@ij.ad.jp>

2017/08/25

The screenshot shows a web browser window with the address bar displaying "enog.jp/archives/1679". The page title is "Echigo Network Operators' Group". The navigation menu includes "About ENOG", "Mailing List", "Meeting", "Echigo-IX", "お知らせ", and "活動記録". The main content area features the announcement "ENOG46 Meeting 開催のお知らせ" by "masakazu" on "2017年7月7日" with "0 Comments". The announcement details include the date and time "2017/8/25(金) 15:00~17:00" and the venue "柏崎市文化会館アルフォーレ 大会議室 (新潟県柏崎市日石町4番32号)". A list of links is provided, including "http://www.artfore.jp/". A sidebar on the right contains a search box and a "最近の投稿" (Recent Posts) section with links to "ENOG49 Meeting 開催のお知らせ", "ENOG48 Meeting を開催しました", "ENOG48 Meeting 開催のお知らせ", "ENOG47 Meeting を開催しました", and "ENOG47 Meeting 開催のお知らせ".

観測されている概要

- 2017/08/25 12:22JST頃
 - AS15169が他ASのIPv4経路をトランジット開始
 - 日頃流通しない細かい経路が大量に広報
 - これによりトラヒックの吸い込みが発生
 - 国内の各ASで通信障害を検知
- 2017/08/25 12:33JST頃
 - AS15169がトランジットしていた経路を削除

観測された問題のBGP経路概要

- 経路数
 - 全体で約11万経路 (日本分が約25000経路)
 - /10から/24まで幅広い経路(半数程度が/24)
 - 通常流れていない細かい経路が多かった
- AS PATHは概ね “701 15169 <本来のAS PATH>”
 - 広報元AS番号は正しそう
 - 各ASが直接AS15169と張っているBGP接続では今回の経路広報は観測されていない
- 対象AS
 - 全体で約7000 AS程度 (日本分が約89 AS)

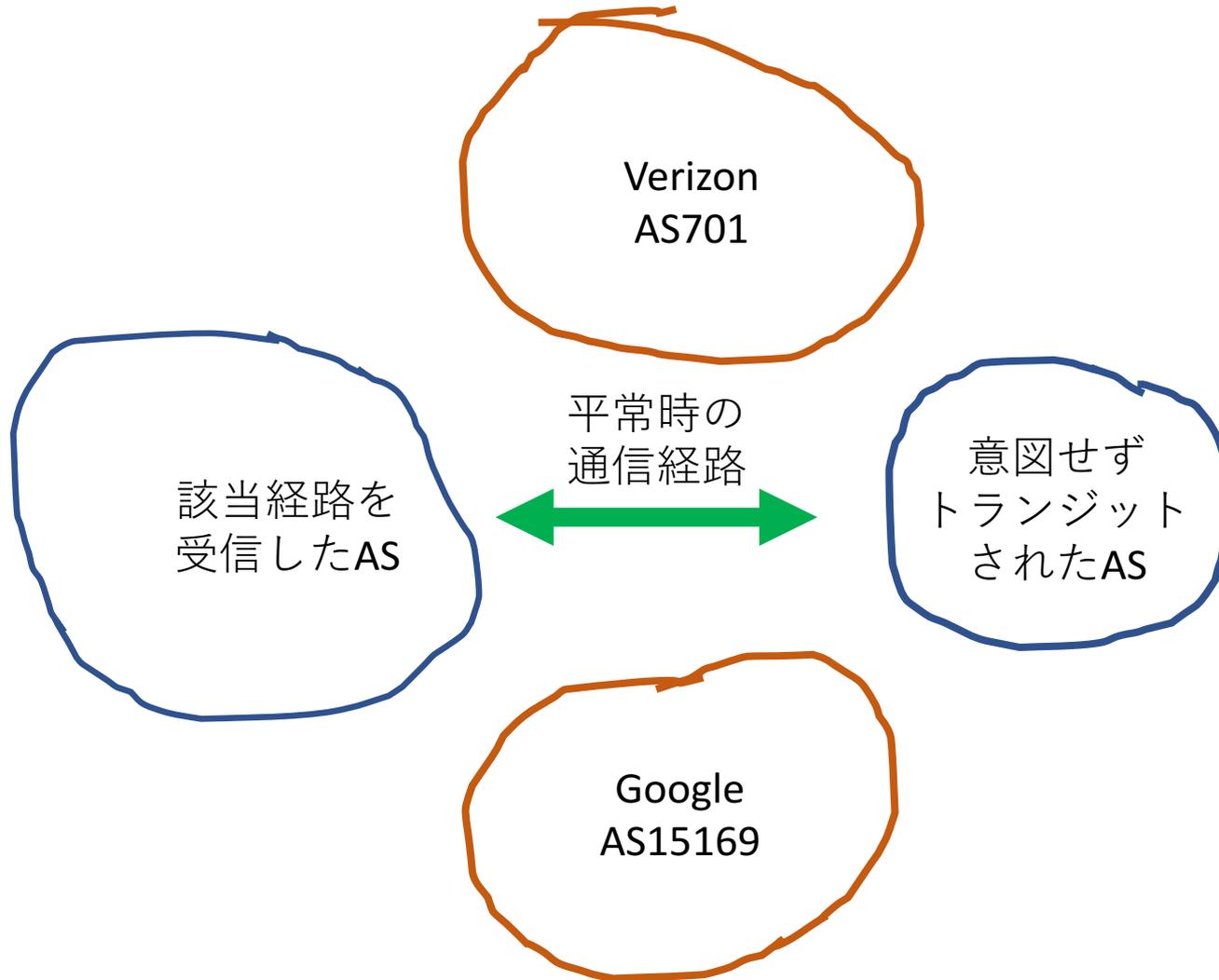
BGPは観測点によって見える情報が異なるのでご注意

トランジットされちゃったAS

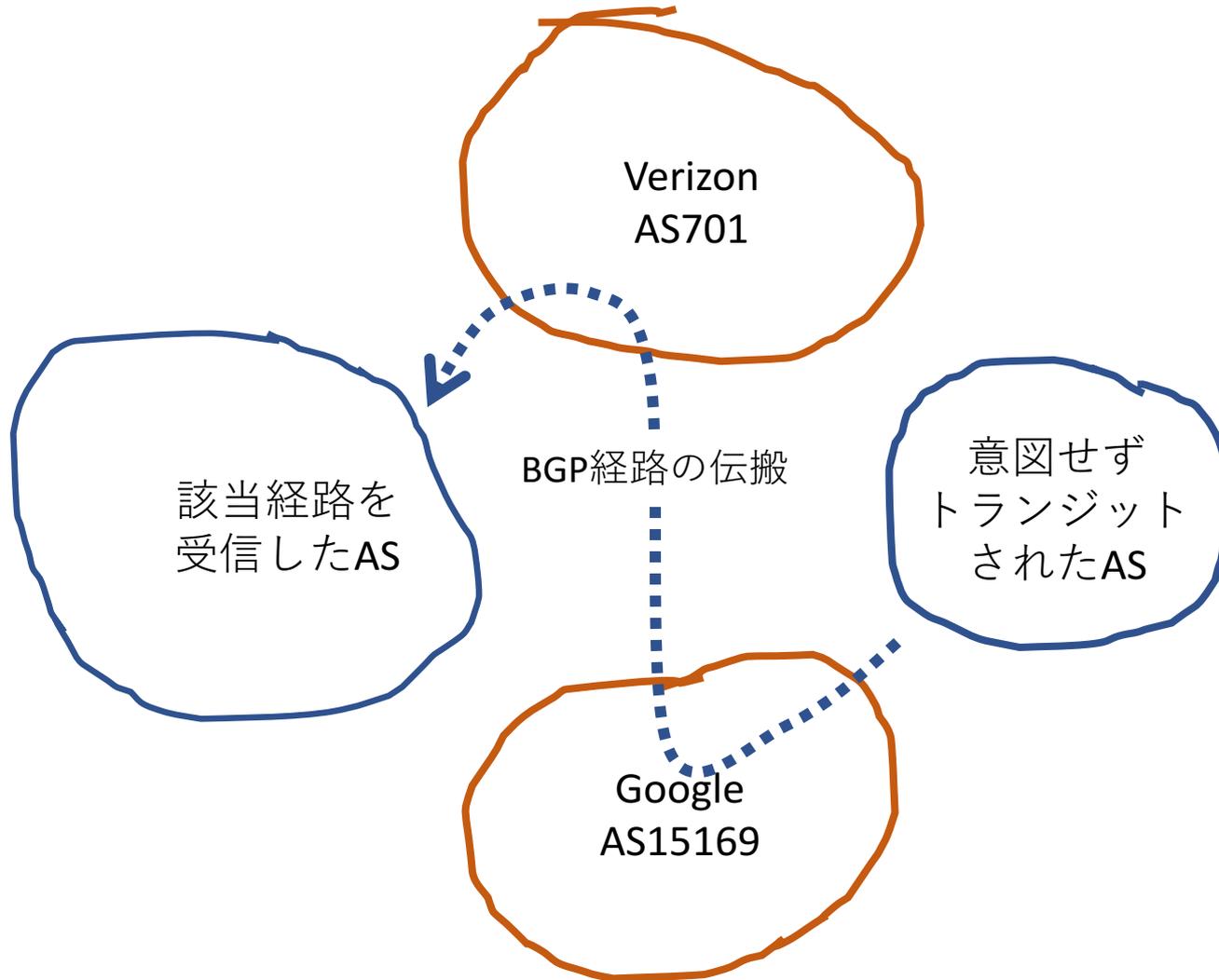
- 世界でおよそ7000 AS程度
 - 内、日本(JPNIC管轄)のものが 89 AS
- 広報されたprefix数のAS別順位
 - OCN/AS4713が大きな影響を受けている

AS番号	prefix数
4713/OCN	24381
7029/WINDSTREAM	7837
8151/UNINET	4639
9121/Turk Telecom	4606
1659/TANet	3106
9394/CTTNET	2137

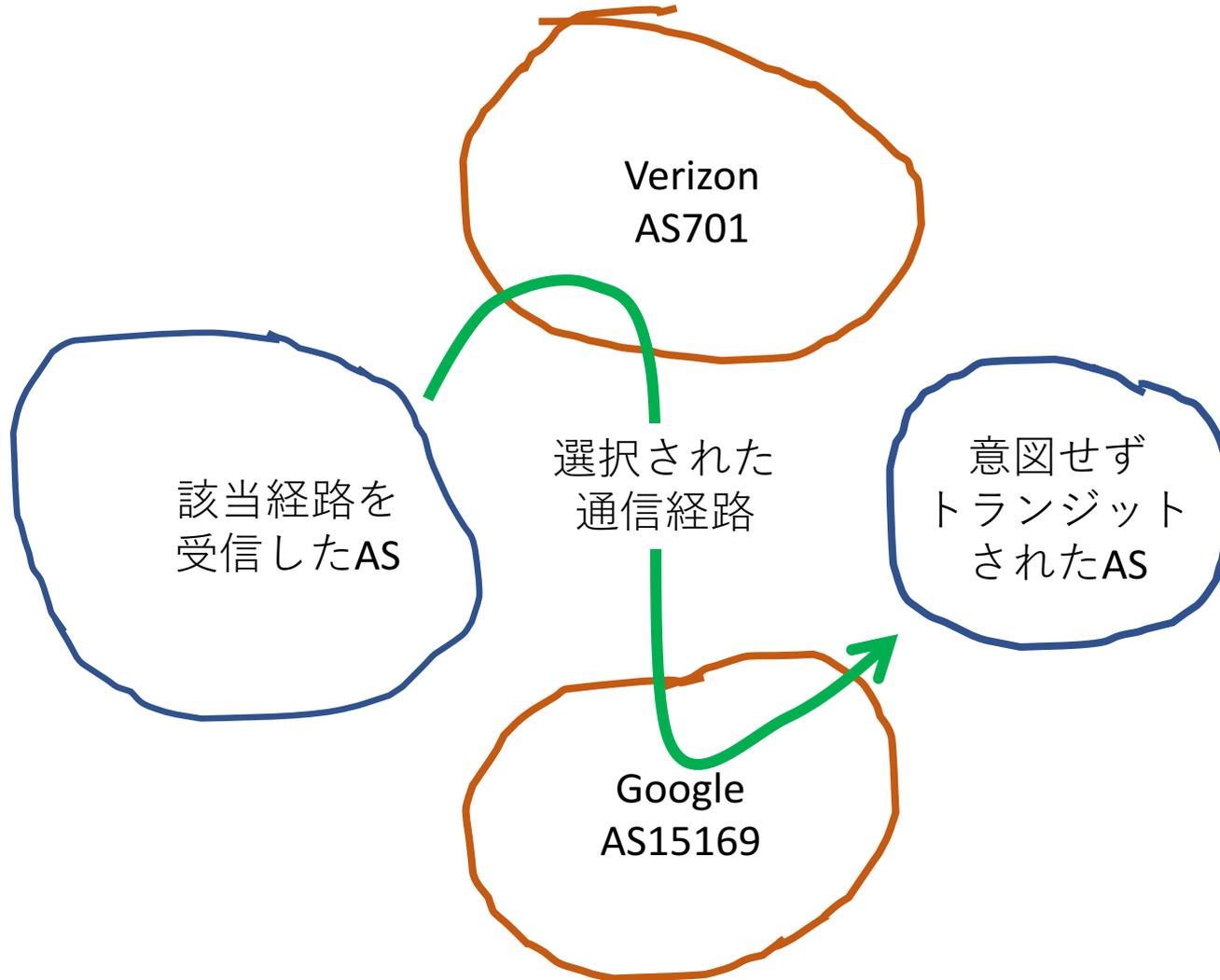
概要図 1：平常時



概要図 1：平常時



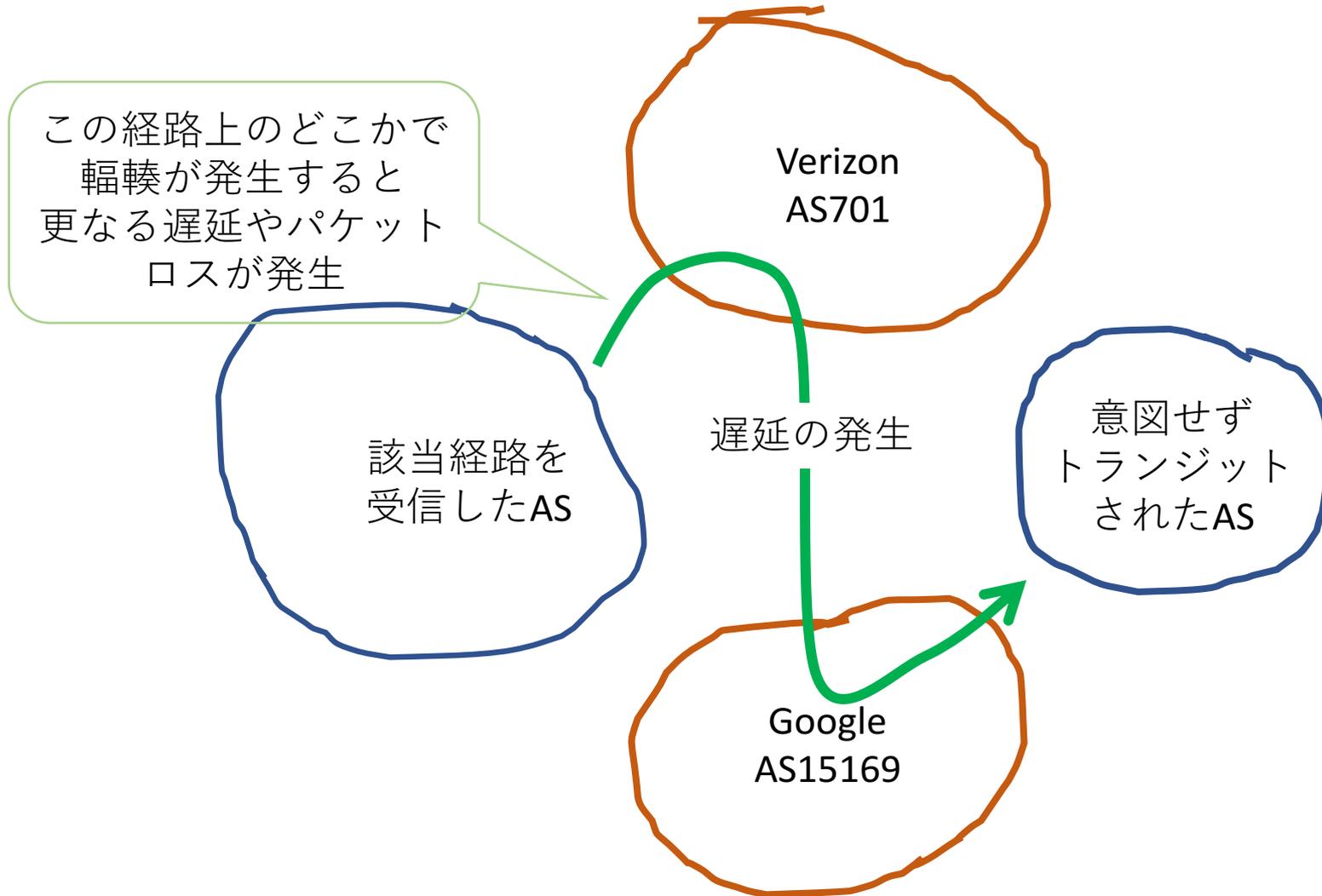
概要図 1: 平常時



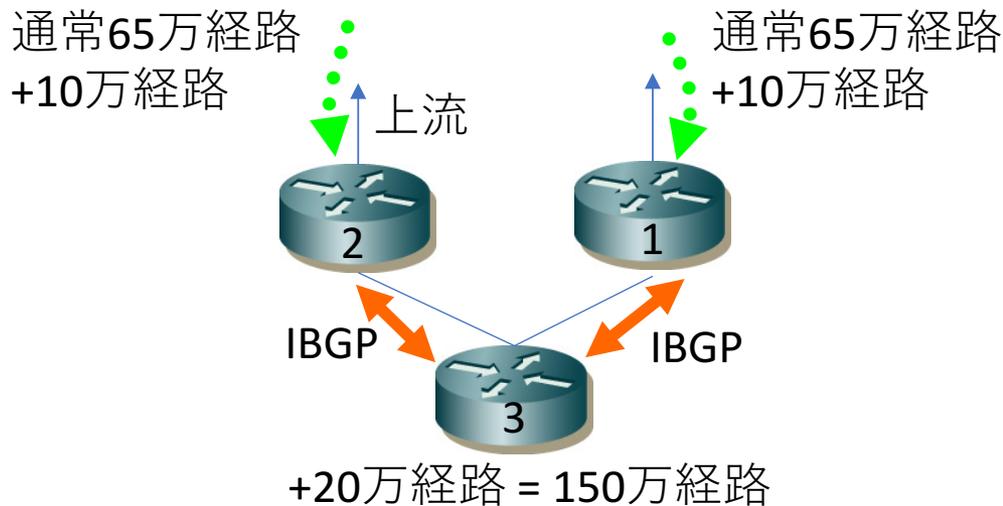
推定される障害や影響

- 広報された宛先向けの通信で遅延が発生した
 - 通信が米国経由になったため
 - 経路上に十分な帯域がない場合は輻輳の発生
- 大量の経路広報で機器が不安定になった
 - 負荷上昇でルータが不安定になった
 - RIB/FIB溢れでルータが不安定になった
 - 何らか機器のbugを踏んだ
- IXP越しの通信が迂回したかも
 - 発生条件が特殊なので、恐らく今回は発生していないのではと考えている

概要図 1: 平常時



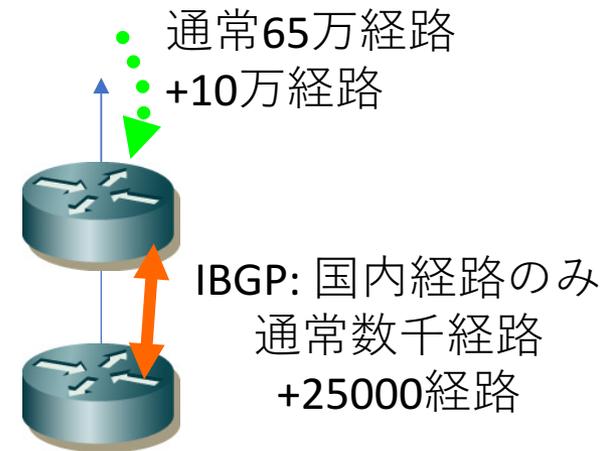
大量の経路追加



- 現状ざっくり約**65万経路**
 - 何もしていないと内部のRT3は65万x2で**130万経路**
- 今回、**10万x2追加で150万経路受信していたかも**
 - 構成によっては更に多い場合も

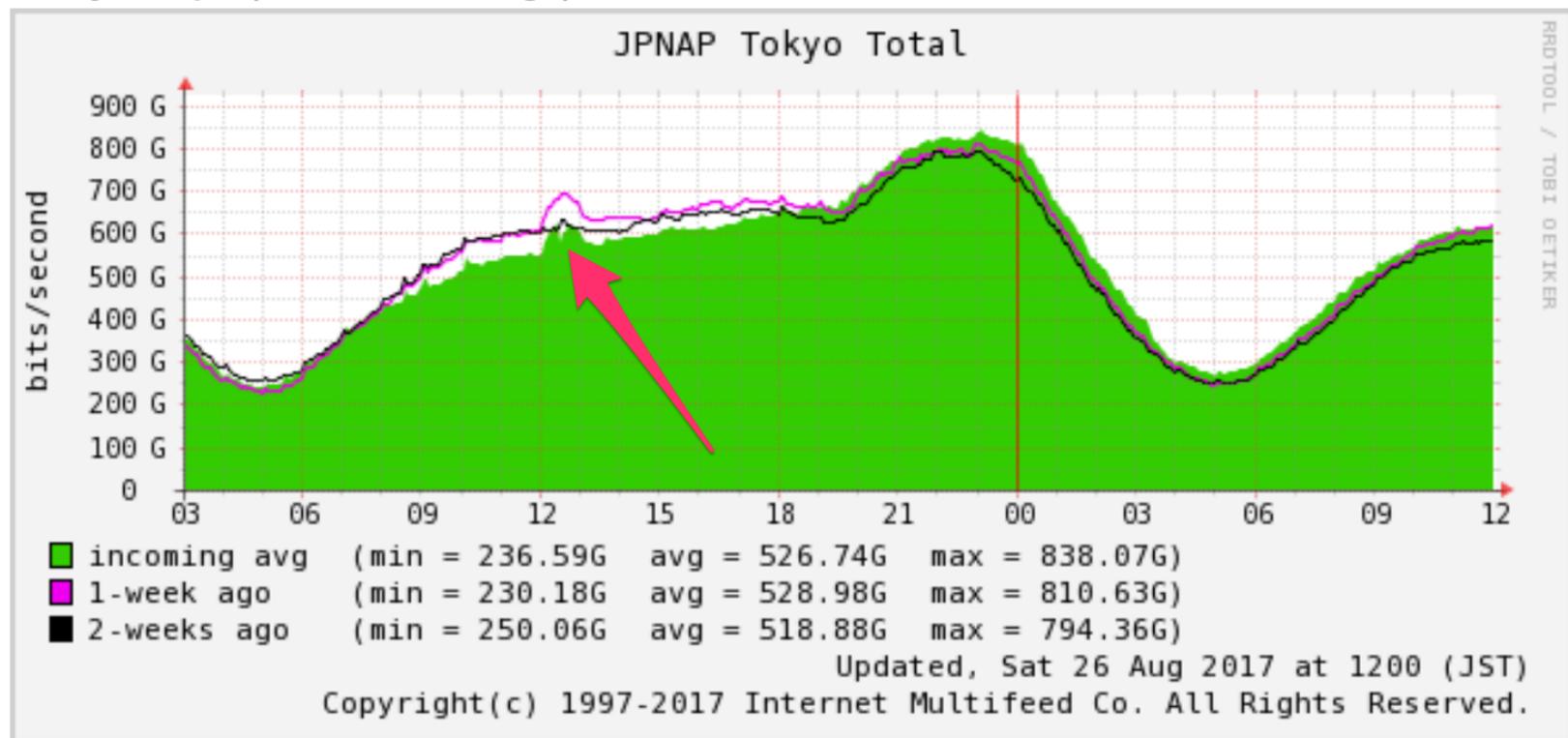
経路削減を適用してても

- 非力なルータで運用するため国内経路のみを内部ルータに渡している場合
- AS PATH(4713 等)で国内経路を識別していた場合、追加で約25000経路
 - 構成によってはもっと多い
 - 通常時の5倍から10倍の経路数が追加された可能性がある
- これら非力なルータが過負荷になるなどの障害が発生した可能性がある



IXPでトラフィック減を観測

`Daily' Graph (5 Minute Average)



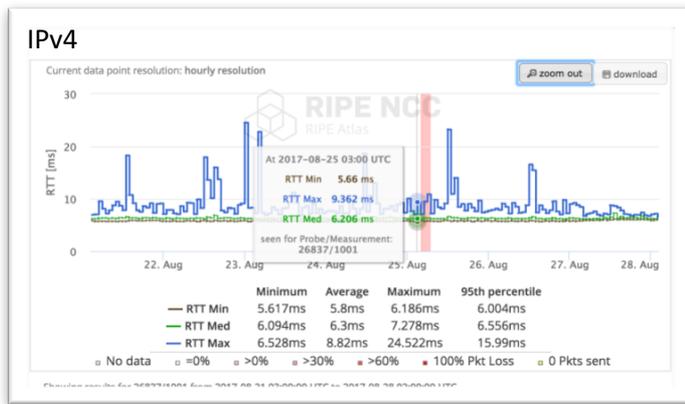
<http://www.jpnap.net/service/traffic.html>

RIPE Atlas Probeで通信影響測定

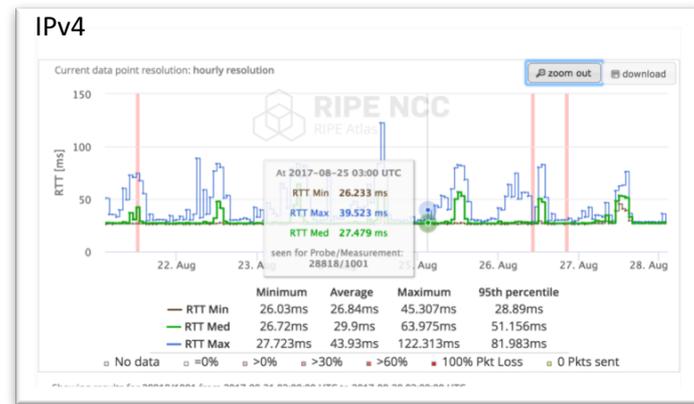
- RIPE NCCのプロジェクト
 - 世界にProbeを配っていて、エンドユーザ視点での計測が可能
- AS4713のprobeを抽出し、宛先別に影響を推定
 - OCN内で通信が完結する宛先: k.root-servers.net
 - 国内で今回の影響を受けた宛先: m.root-servers.net
 - 海外で今回の影響を受けた宛先: ctr-ams02.atlas.ripe.net

RIPE Atlasで見る: OCN内

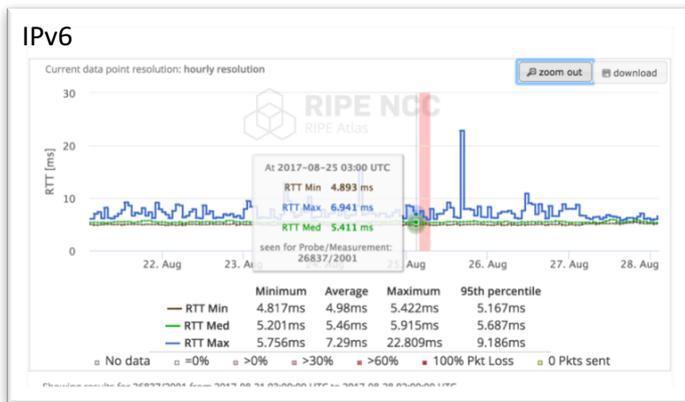
Probe26837



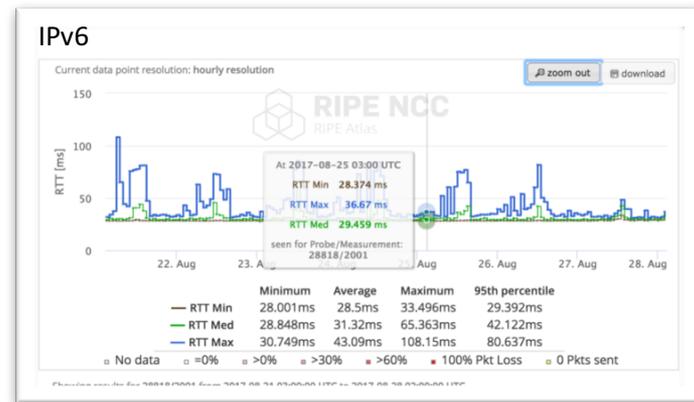
Probe28818



IPv6



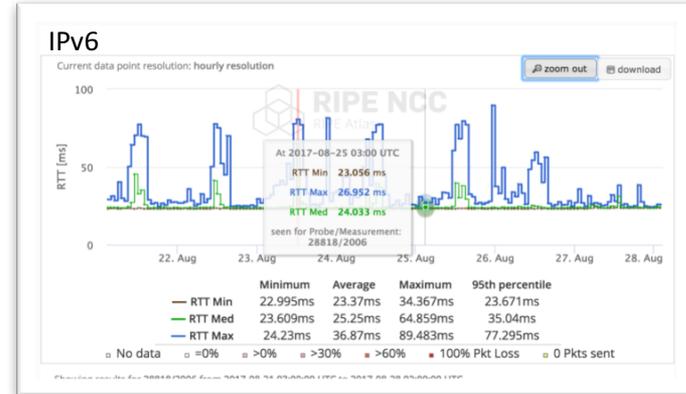
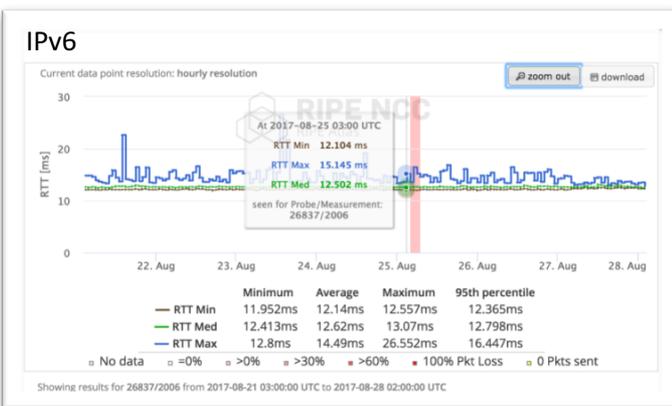
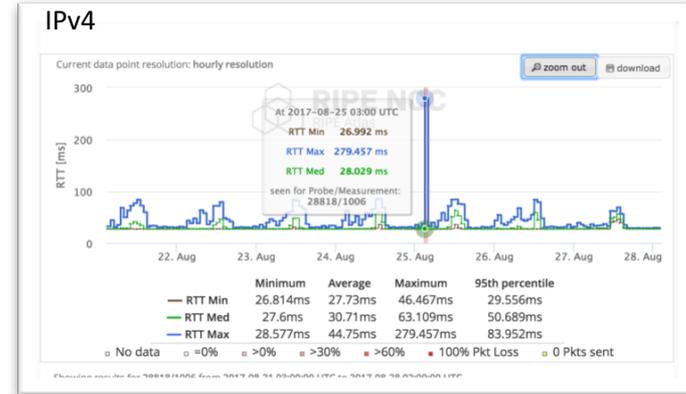
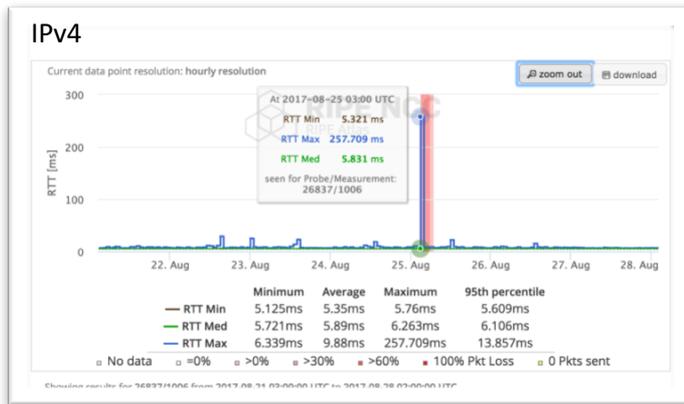
IPv6



RIPE Atlasで見る: OCNと国内

Probe26837

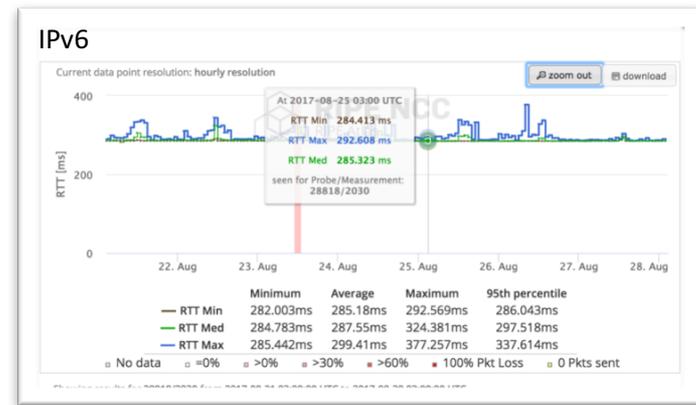
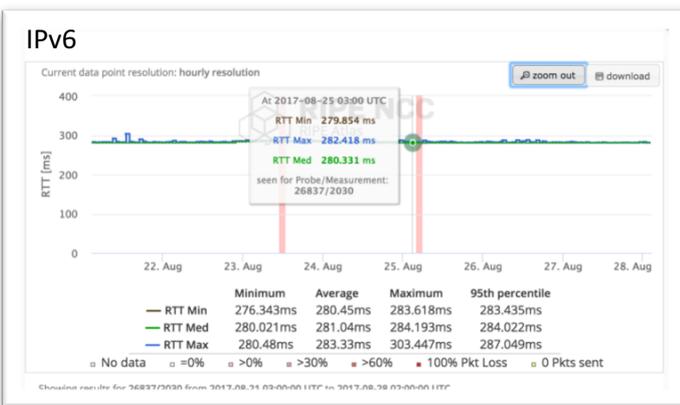
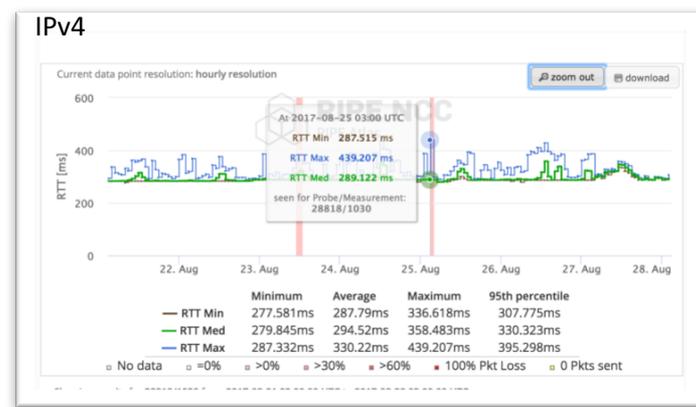
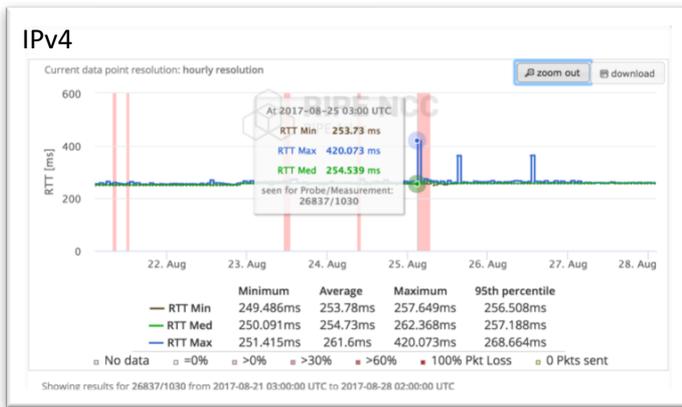
Probe28818



RIPE Atlasで見る: OCNと海外

Probe26837

Probe28818



RIPE Atlasから見えること

- 該当Probeでは国内、海外のIPv4通信に遅延の増加やパケットロスを観測
- IPv6へ直接の影響はほとんどなかった模様
 - IPv4のBGP経路が対象であったため
 - probe26837ではIPv4/IPv6で宛先に寄らずパケットロスが観測されているためProbe近傍のどこかで輻輳が発生していたかもしれない

影響時間の推定

- 迂回による遅延や輻輳の影響
 - 今回広報された経路に関わる通信
 - 経路やRIPE Atlasなどの記録からすると**20分以内**には収束したはず
- 大量の経路広報で不安定になった機器の影響
 - 影響した機器があった場合、そこを通過する通信
 - 収束時間は**その運用組織の対応に依存**

僕がやってたこと 08/25(金)

- 14時台
 - ENOG46会場に向かい中、某社SOCより連絡
 - 調査開始して、関係各所に情報共有
- 15時台
 - JANOG MLに投稿 [janog:14007]
 - ENOG46で発表
- 17時台
 - AS15169 NOCに情報共有
 - Twitterに投稿
- 18時台
 - ENOG46懇親会



Yoshinobu Matsuzaki
@maz_zzz

08/25 12:22JSTから発生した経路障害の件。実はまだ詳細は調査中なのだけど、米国のAS15169が他ネットワークを誤ってトランジットしてしまった様に見える。この影響で国内であっても一部通信が米国経由になって遅延やパケットロス、通信障害が発生した模様

17:28 - 2017年8月25日

僕がやってたこと 08/28(月)

- その日のIA研@IJJで事象を概説することになり、慌てて資料作成
- 何とかできたので公開



 **Yoshinobu Matsuzaki**
@maz_zzz

08/25の通信障害の件、これまでに調べたことを資料にまとめてみたよ。
[attn.jp/maz/p/t/pdf/20 ...](https://www.attn.jp/maz/p/t/pdf/20170825-routeleakage.pdf)

16:32 - 2017年8月28日



電子情報通信学会 研究会発表申込システム
研究会 開催プログラム

インターネットアーキテクチャ研究会 (IA) [schedule] [select]

専門委員長 飯田 耕吉 (北大)
副委員長 新 藤 (IJJ), 大崎 博之 (関西学院大), 渡久 智樹 (阪大)
幹事 作元 雄輔 (京都大東院), 岸 健一郎 (NICT)
幹事補佐 大平 健司 (徳島大), 坂野 遼平 (NTT), 渡辺 俊貴 (NEC)

日時	2017年 8月28日(月) 12:30 - 19:30
議題	インターネット運用・管理, 一般 (JANOG協催)
会場名	IIJセミナールーム (飯田橋グランブルーム 13階)
住所	〒102-0071 東京都千代田区富士見2-10-2 飯田橋グランブルーム
他の共催	◆JANOG協催
お知らせ	◎会場に入るためには事前に入館登録が必要となります。 8/25までに以下のフォームから参加申込みをお願いします。 https://docs.google.com/forms/d/e/1FAIpQLSczmXOcYPMYsaIRwplSAE7qPBvChv88mdIR4gJkte_x50kIQ/viewform
参加費に ついて	この開催は「技術完全電子化」研究会です。参加費についてはこちらをご覧ください。参加費のお支払いが必要な研究会はJA研究会です。

8月28日(月) 午後 開会

<https://www.attn.jp/maz/p/t/pdf/20170825-routeleakage.pdf>

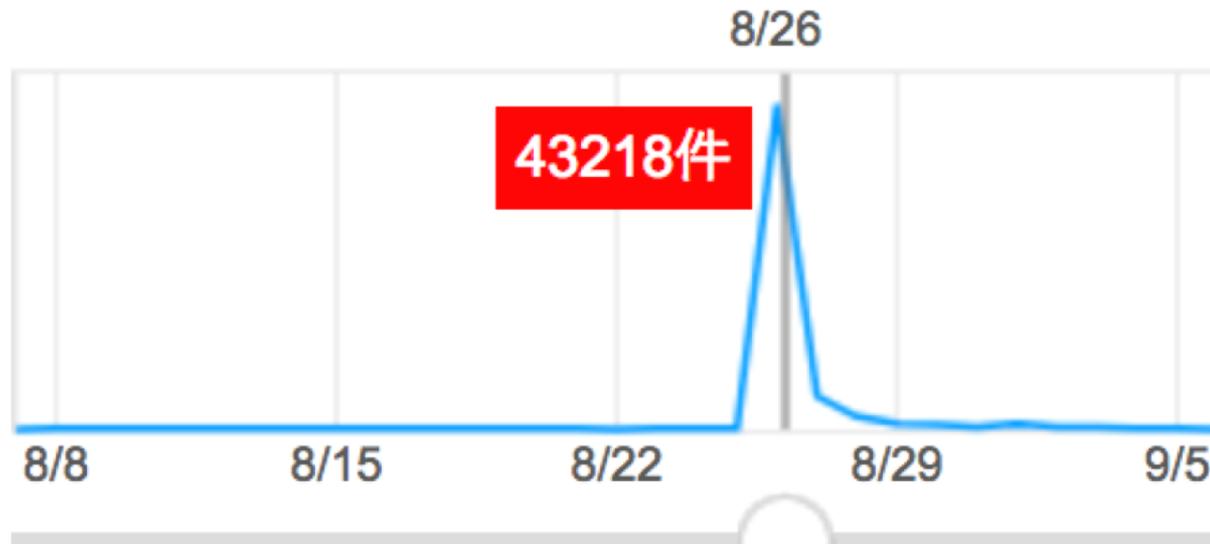
通信障害の分析グラフ

24時間

7日間

30日間

ツイート数の推移



<https://search.yahoo.co.jp/realtime/search?p=%E9%80%9A%E4%BF%A1%E9%9A%9C%E5%AE%B3&ei=UTF-8>

電気通信事故検証会議の報告

平成 29 年 8 月に発生した大規模な
インターネット接続障害に関する検証報告

平成 29 年 12 月
電気通信事故検証会議

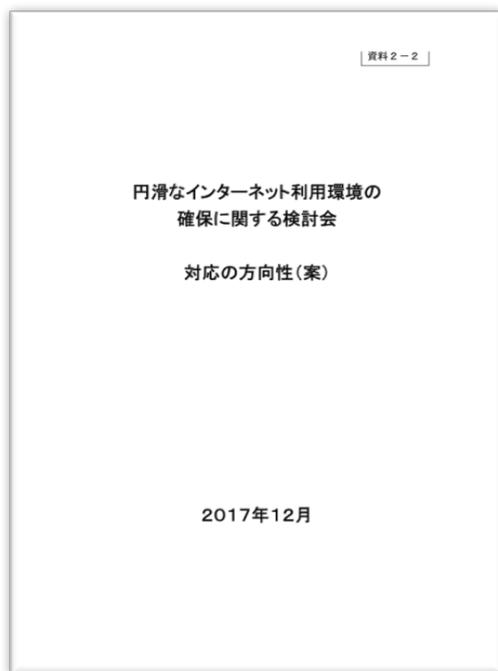
- 人為的ミスの未然防止
- 誤送信された経路情報の受信防止
及び不要な経路情報の送信防止
- 障害に関する情報の電気通信事業者間での共有
- 利用者周知

16 ベストエフォートサービスにおいては、停止が
ない限りは品質の低下と判断していないことから、
総務省は、本事案については電気通信事業法上の
事故には該当しないと判断している。

http://www.soumu.go.jp/main_content/000523153.pdf

円滑なインターネット利用環境の確保に関する検討会

http://www.soumu.go.jp/main_sosiki/kenkyu/smooth_internet/



検討会(第2回)配布資料: 対応の方向性(案)

- の通信の基盤として、適切な情報伝達手段の確保が重要である。また、インターネット利用環境の円滑な利用を確保するためには、適切なセキュリティ対策の取組が不可欠である。また、インターネット利用環境の円滑な利用を確保するためには、適切なセキュリティ対策の取組が不可欠である。
- また、インターネット利用環境の円滑な利用を確保するためには、適切なセキュリティ対策の取組が不可欠である。

http://www.soumu.go.jp/main_content/000525811.pdf

情報集約と展開？

- 日本経済新聞、1月23日紙面
 - サービスで障害が発生した場合、第一報を接続事業者から総務省に報告
 - 同省が...関係する業者に必要な情報を送る
 - 対象となる障害の度合いやどこまで報告するかは業者の意見を踏まえて今後詰める

情報共有、周知への課題

- 会社事情
 - 会社ポリシー、NDA
 - 「公式見解」による硬直化
 - 利用者対応
 - 即時性と不確かな”障害”情報
 - 報道対応
 - 報道は犯人探しになりがち
- どうすれば、うまくやれるだろうか？

経路制御でやりたいこと

- **到達性確保**
- **トラヒック制御**
- 付随して確保したいことは他にもあるけど
 - 冗長性
 - 拡張性
 - 運用容易性

自身が広報するBGP経路には届きたい範囲(スコープ)がある

- 到達性を担保する経路(トランジット向け広報)
 - グローバルに届いて欲しい
- 到達性を確保する経路(ピア向け広報)
 - 隣接ASとその配下のみに伝搬して欲しい
- トラヒックを制御する経路
 - 隣接ASとその配下のみに伝搬して欲しい
 - 場合によっては隣接ASだけでも良い
- でも、これらは自分では制御できない
 - 一旦広報した経路は、受信側のポリシーで制御される

そんなわけで

並
我

並
用

情報の集取と展開

- みんな何から情報集めてますか？
- どこで情報提供するのが嬉しいですか？
- それっぽい情報を出す元はどこが良い？
 - 個人
 - JANOG WG
 - ICT-ISAC
 - JPNIC
 - お役所

BGP運用どうします？

- フィルターやリミッターの設定？
 - リミッターはかなり覚悟が必要な設定
 - 今回の件は防げない
- 相互接続を推進、もしくは縮退？
 - 接続先を増やしてAS_PATH長で勝てる状態を作る
 - 接続先を減らして漏洩の可能性を減らす
- leak community？
 - <https://datatracker.ietf.org/doc/draft-heiz-idr-route-leak-community/>

局面に応じた情報共有

- 発生直後
 - 事象把握のために事業者間で情報交換
 - JANOG等で、信頼関係を構築しておきましょう
 - 硬直化前にできるだけ早く
- 利用者周知
 - 障害認知の公表、状況の適宜更新
 - 権威ある周知情報があると参照しやすい
- 報道対応
 - 僕たちが伝えたいことと、報道価値は別モノかも

送出経路ポリシーは基本2つ

- ピア/上流に広報する経路
 - 自身と配下の経路
- 顧客ASに広報する経路
 - フルルート/デフォルトルート
- **BGP**でトラヒック制御しようとするすると複雑になっていく
 - 例えば、特定ピアのみ細かな経路の広報などなど

BGPでトラヒック制御しない

- 到達性の維持に努める
- 必要なところに必要な帯域を用意する
- 必要なASと適宜ピアを拡充する
- ピア、上流に常に同一の経路広報を行う
 - 誰かが経路流出させても、周辺ASと多拠点でピアできていれば、流出経路はAS PATH長で採用されない可能性が高まり、影響範囲が極小となる

複雑化するBGP

- 経路に局所性が出てきている
 - トラヒック制御のため
 - 隠すため
- **BGP**は見ているポイントで経路が異なる
 - 漏れ出した経路は、特定のネットワークだけに届くかも

経路削減は延命策 not 解決策

- 経路フィルタで削減された経路数は現状のもの
 - 将来やトラブル発生時の経路数は未知
- max-prefixなどで受信経路数を制限できるけど
 - BGPピアが切断しても大丈夫？
 - alert受信して、機器が不安定になる前に対処可能？
- 十分余裕を持った運用を心がけるしかない
 - あるいは複数の安全策の併用

BGP運用でできること

- 統一した広報ポリシー
 - AS間のトラヒック制御に細かい経路を利用しない
- 受信ポリシーの検討
 - 受信経路で酷い影響が出ないようにする
- 余裕を持った機材での運用
- 事象を第三者が検証できるように、経路アーカイブ等へ経路情報を提供しておく

peerlockingの適用

- Peerlocking
 - https://www.nanog.org/sites/default/files/Snijders_Everyday_Practical_Bgp.pdf
- あるASと堅牢に相互接続できてるなら、そのAS関連の経路を他から受信しない
 - 他ASが経路を漏洩しても影響を受けない
 - 迂回路が直接接続のどれかに限られるため、事故が起こった際の影響はより酷いことになるかも
 - 運用は結構大変

peerlocking模式図

- AS_PATHにAS-YYYが含まれている経路を他のASから一切受信しない
 - 上流、ピア、顧客

