

高速な障害復旧に必要な 思いやり

ソフトバンクテレコム

松嶋 聡

背景::かれこれ3年前

- 新しいネットワークをつくりなさい
 - とあるアプリケーション用です
 - 高速な障害復旧, 可能ならば sub-second で
 - メトロエリアのネットワークを紡いで、全国をカバーしなさい
 - 数千のルータ
 - 数百のメトロエリア
- 当時の前提条件
 - L1, L2 の故障検出メカニズムは使えません
 - え? broadcast media なんすか?
 - 経路制御プロトコルは BGP のみです
 - なんで?
 - BFD有史以前です
 - あっても、「BFD? そんな高い箱使えません!」
 - ああ、要は安くあげろと...

ネットワーク設計してみる（1）

Physical view

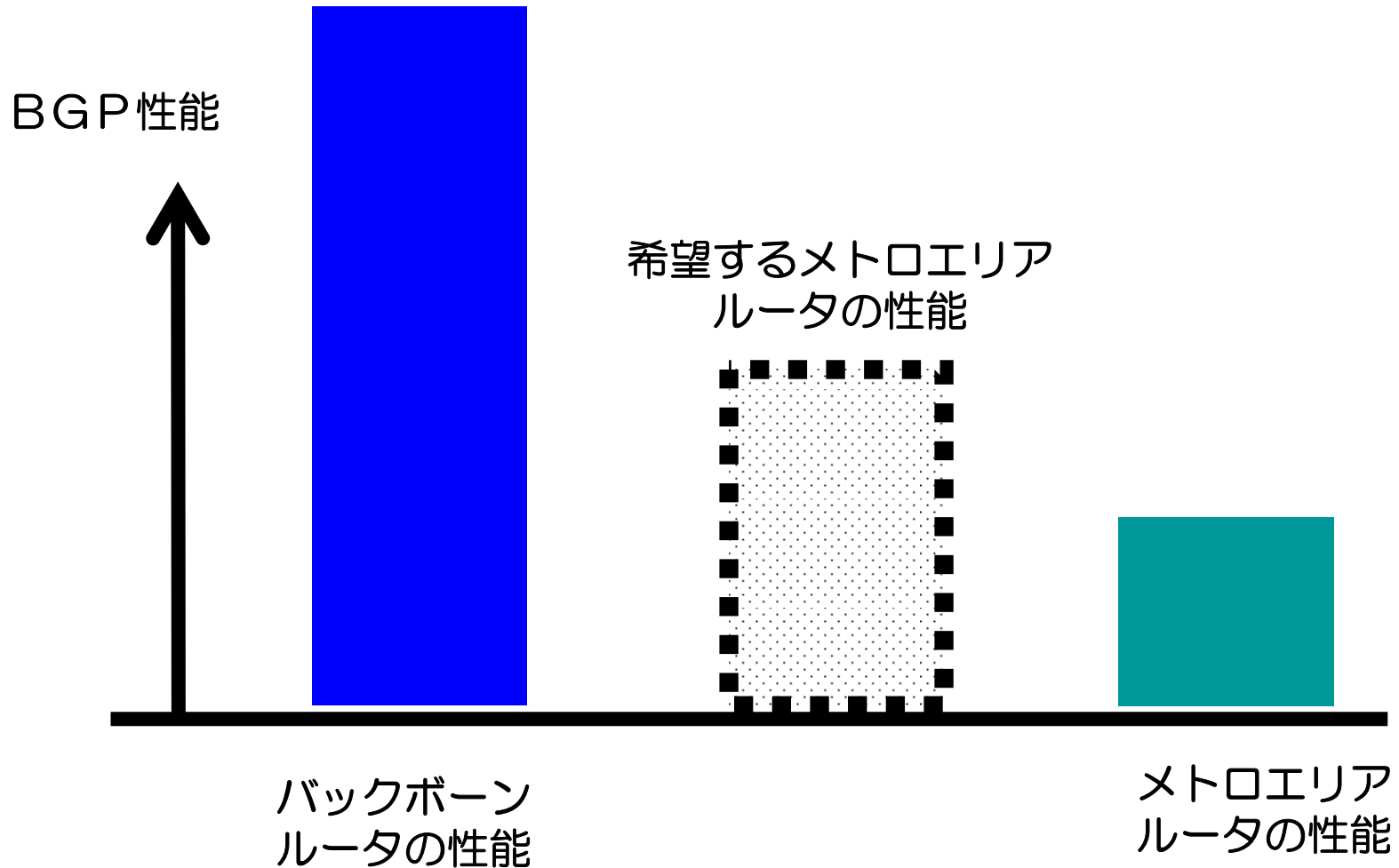
Logical view



ネットワーク設計してみる（2）

- 設計パラメータ::BGP
 - 1バックボーンルータあたり、数百ピア必要
 - ネットワーク全体で数万経路が必要
 - 経路（遅延）最適化のため、経路集約効率が限定される
- 設計パラメータ::ルータ能力
 - バックボーンルータのスペックはそれなり
 - メトロエリアルータは、しょぼしょぼ

ネットワーク設計してみる (3)

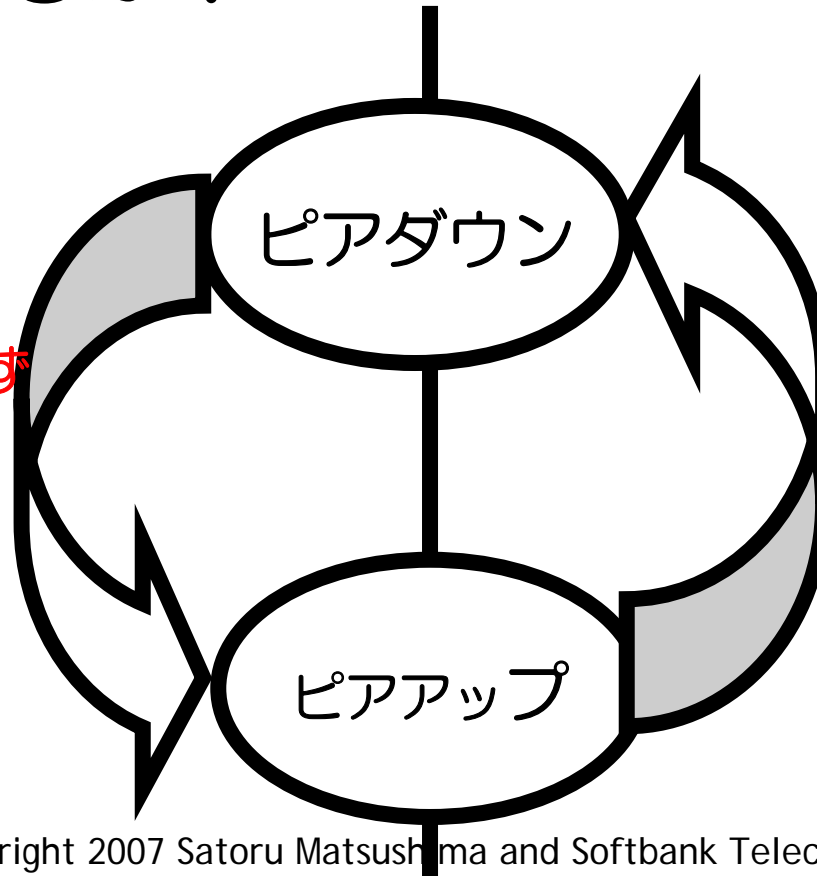


ネットワーク設計してみる（４）

- 極端な性能差のBGPピアが暴れだしたらどうなるか？

小さいルータ

- CPUパンパン
- 経路受け切れない
- ピアアップしきれず
また落ちる



大きいルータ

- ピアアップで全速
アップデート
- 落ちた経路を都度
アップデート

こんなマシン・機能があったら... (1)

- BGP性能の差をカバーするやさしい箱
 - keep-alive/hold-timer は default で、ピアの up/down 負荷をかけずに経路迂回したい
- でも、そんなの作ってくれないよねえ
 - C30さん！
 - J5rさん！

こんなマシン・機能があったら... (2)

- BFDでなくていいから、ソコソコ速い障害検知メカニズム
 - マルチポイントメディアでも動作して
 - 相手ルータがサポートしていなくても使える一般的なプロトコルでプローブ

こんな機能・マシンがあったら...

(3)

- プローブと経路制御の連携どうしよう...
 - 検出したら telnet/snmp で config 変更？
 - 誰が？ 人間？ NMS？ ひえー
- そうだ！ルートリフレクタ使おう！
 - RRなら、ルータでがんばらなくても、ネットワークに外付けできるちゃん！
 - 経路制御に直接作用させられるからプローブと連携させやすいかも！

林さん！
よろしくお願ひします

お願ひ1：普通のプロトコルでプローブ

お願ひ2：ルータに負荷をかけない高速障害復旧

高速な障害復旧に必要な 思いやり

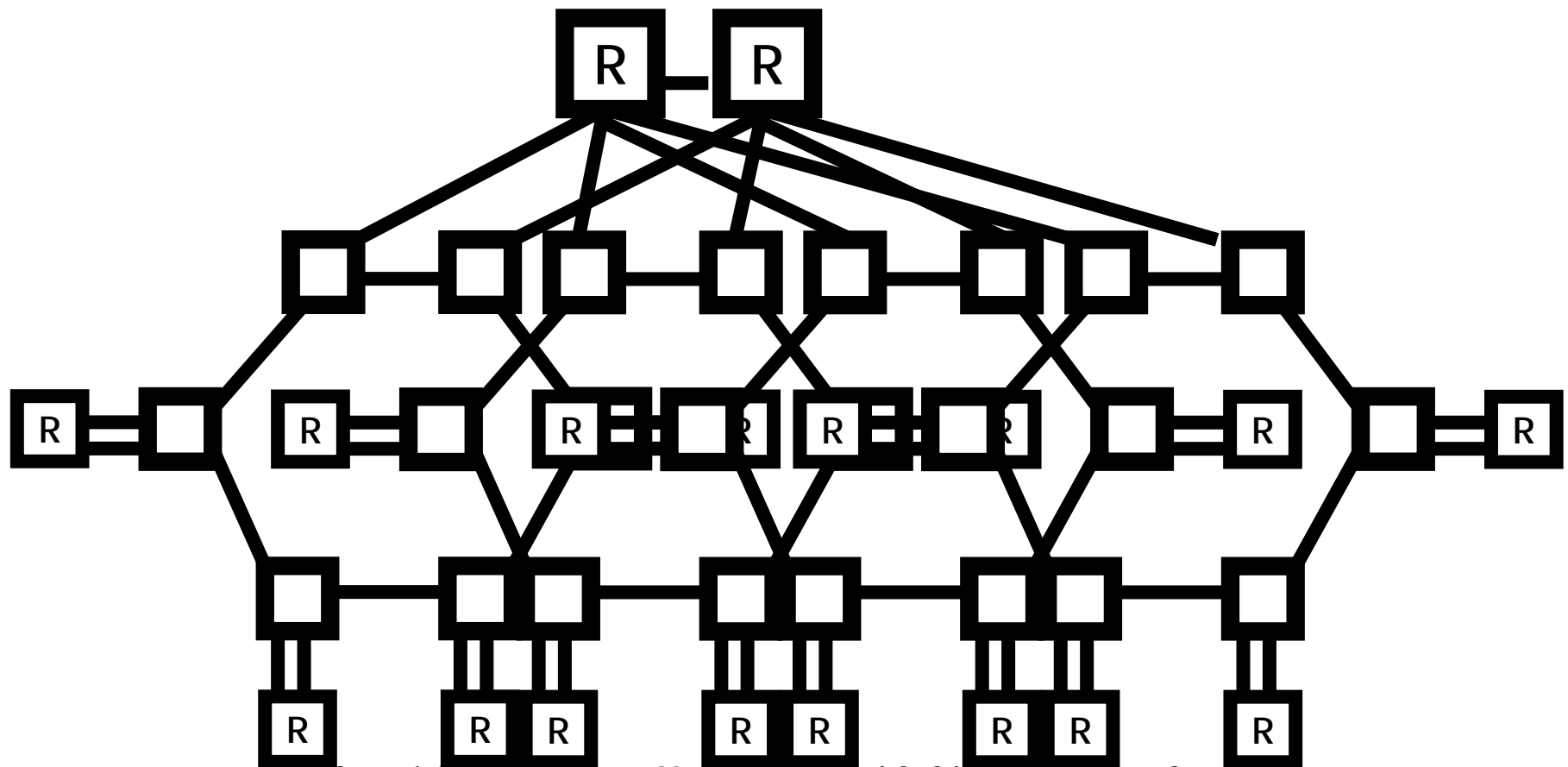
～運用編～

思いやりのある運用を目指そう！

- 誰に対して？
 - はじっこにいる、少し小さいルータたち
 - 数千～数万経路のピアを短時間で上げ下げするのは、かなりつらい
- 小さいルータに負荷をかけない
 - でも、障害検出と復旧はすばやく！
- RRでどこまでできるか？
 - keepalive/hold-timer は default のままで（3分!）
 - プローブの間隔と保護時間の調整をやり
 - ピアをあげたままで、BGPメッセージ数を最小限にキープし
 - 高速な故障検出と復旧（迂回）を実現する！

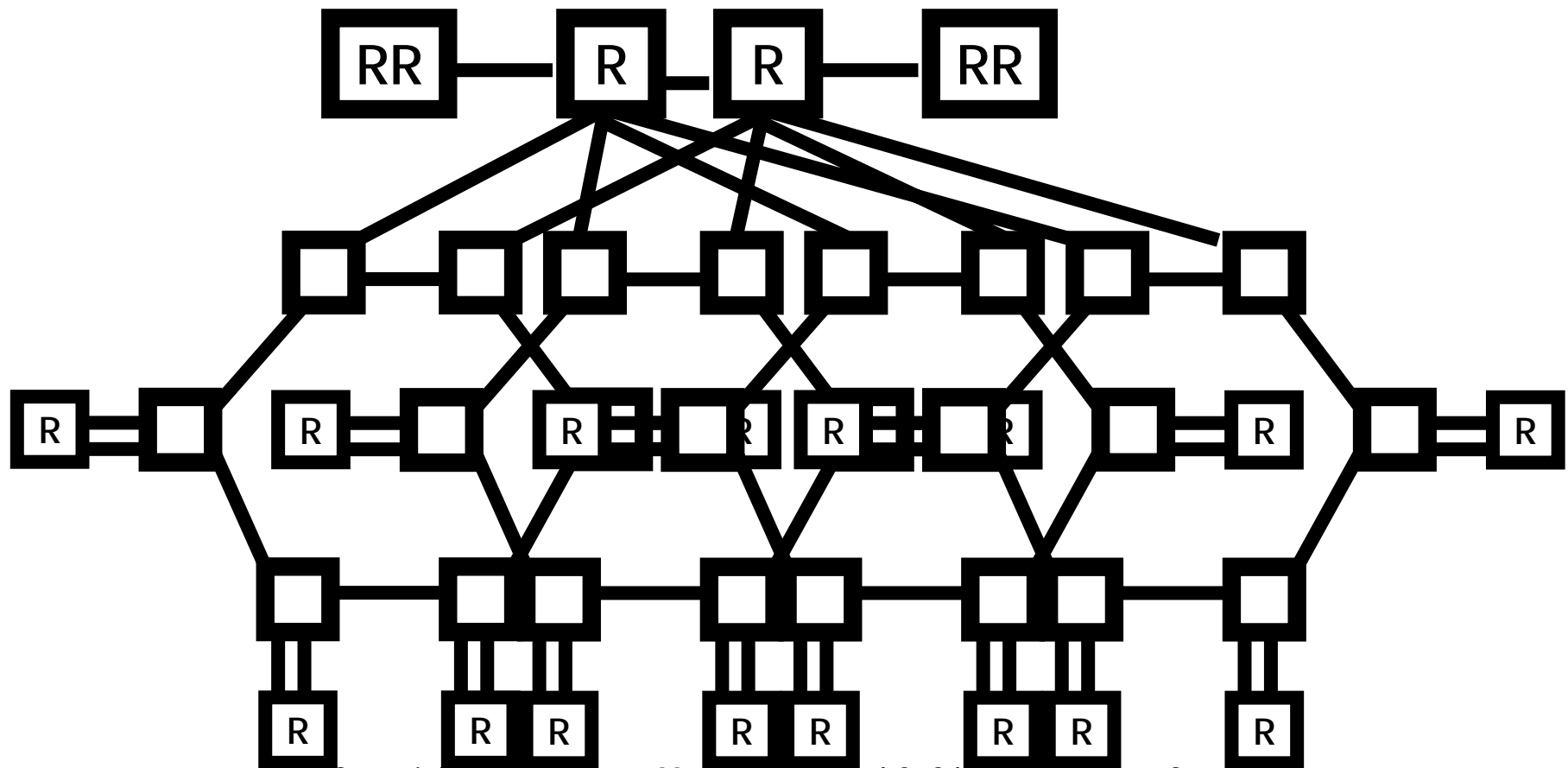
実際に構築してみよう

- メトロエリア集約ネットワークとRR



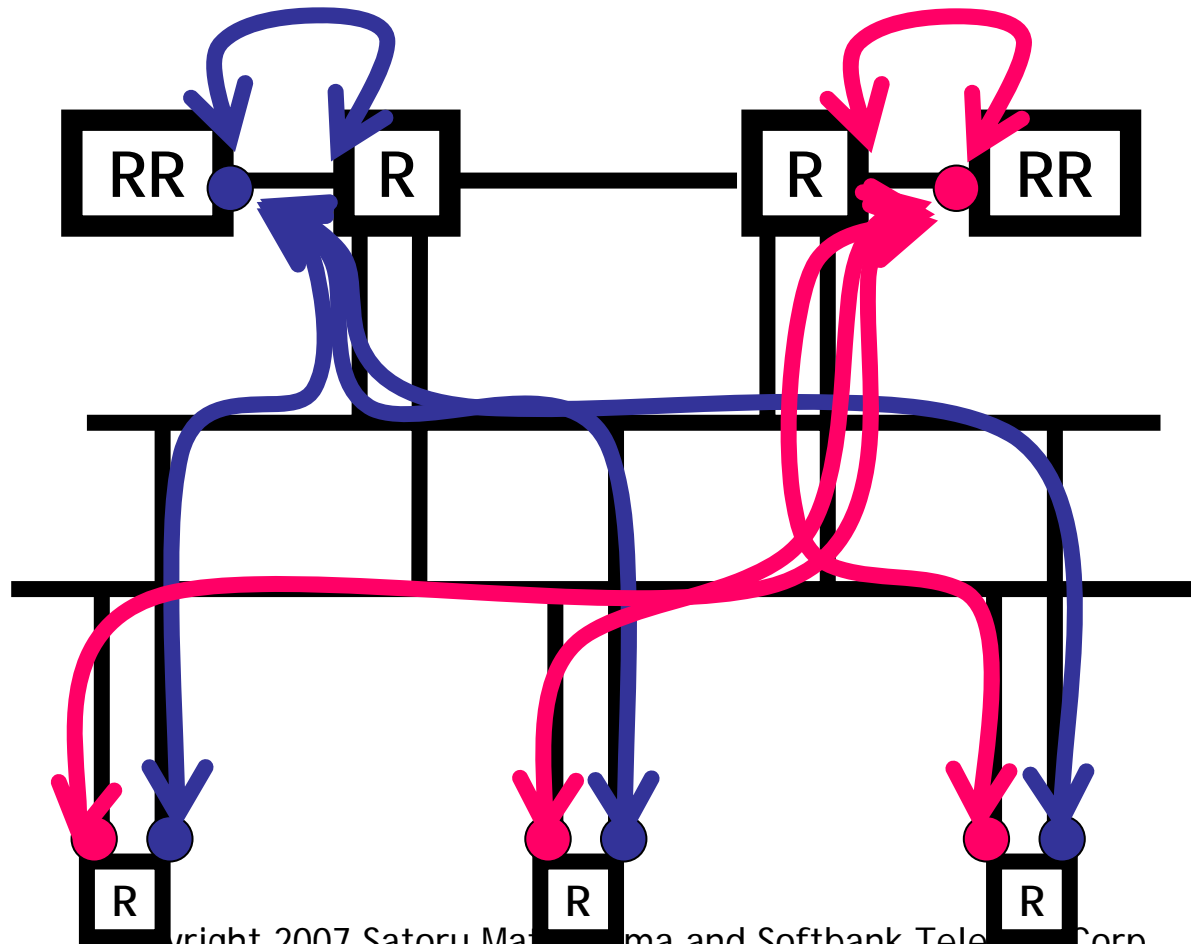
実際に構築してみよう

- メトロエリア集約ネットワークとRR



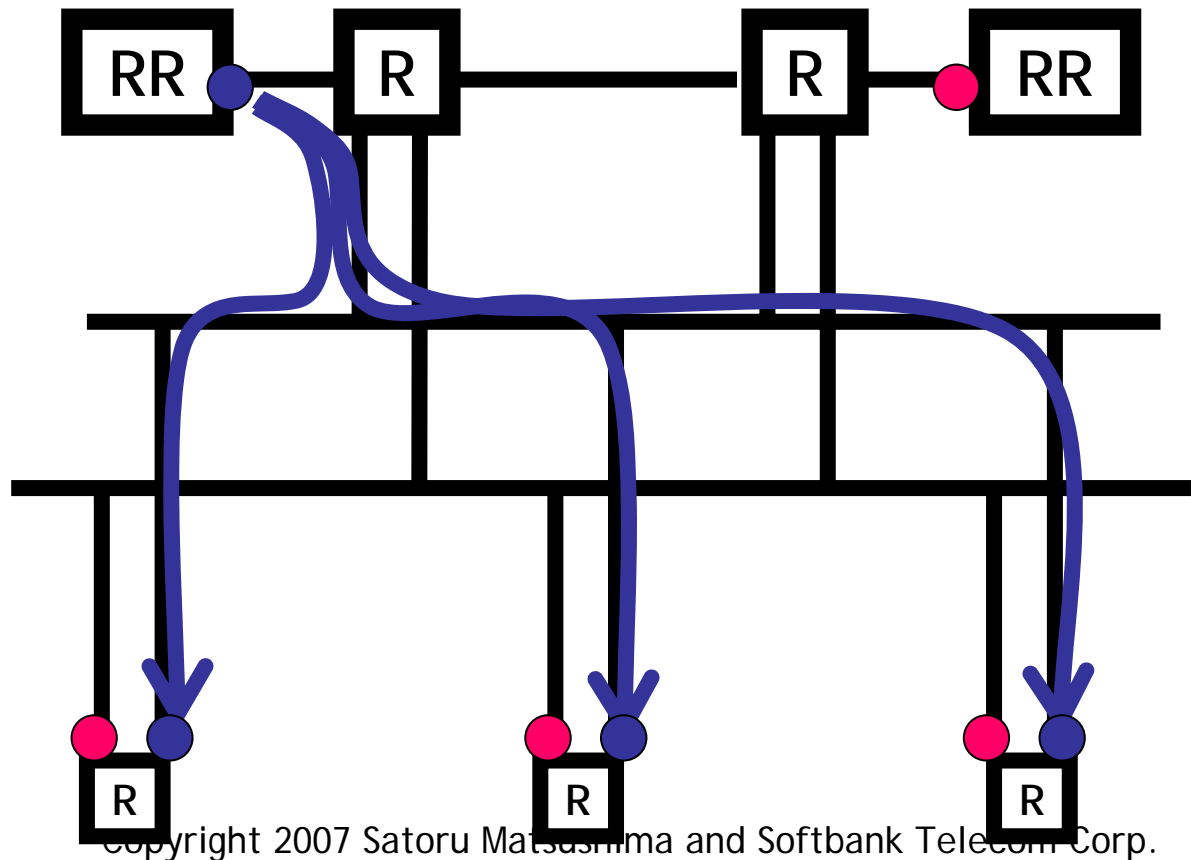
実際に構築してみよう

- BGPピアの設定



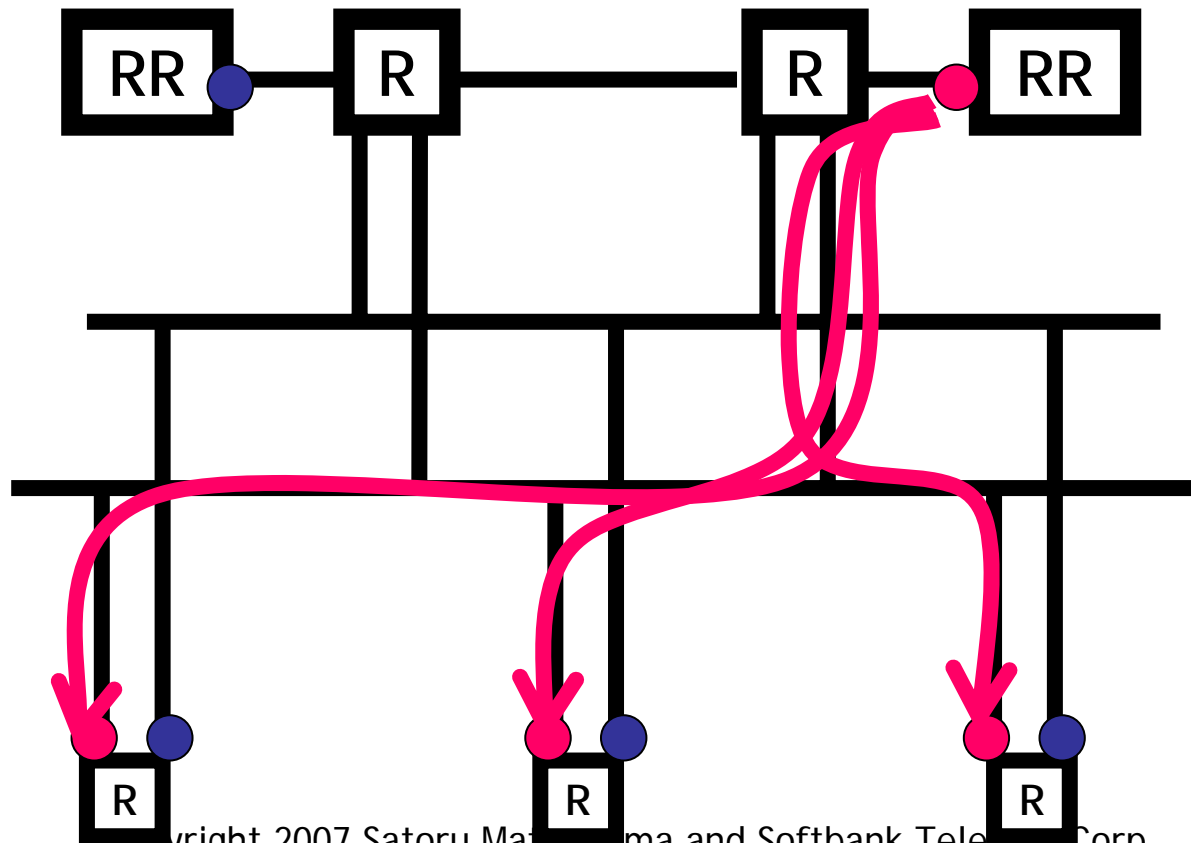
実際に構築してみよう

- プローブのかけ方



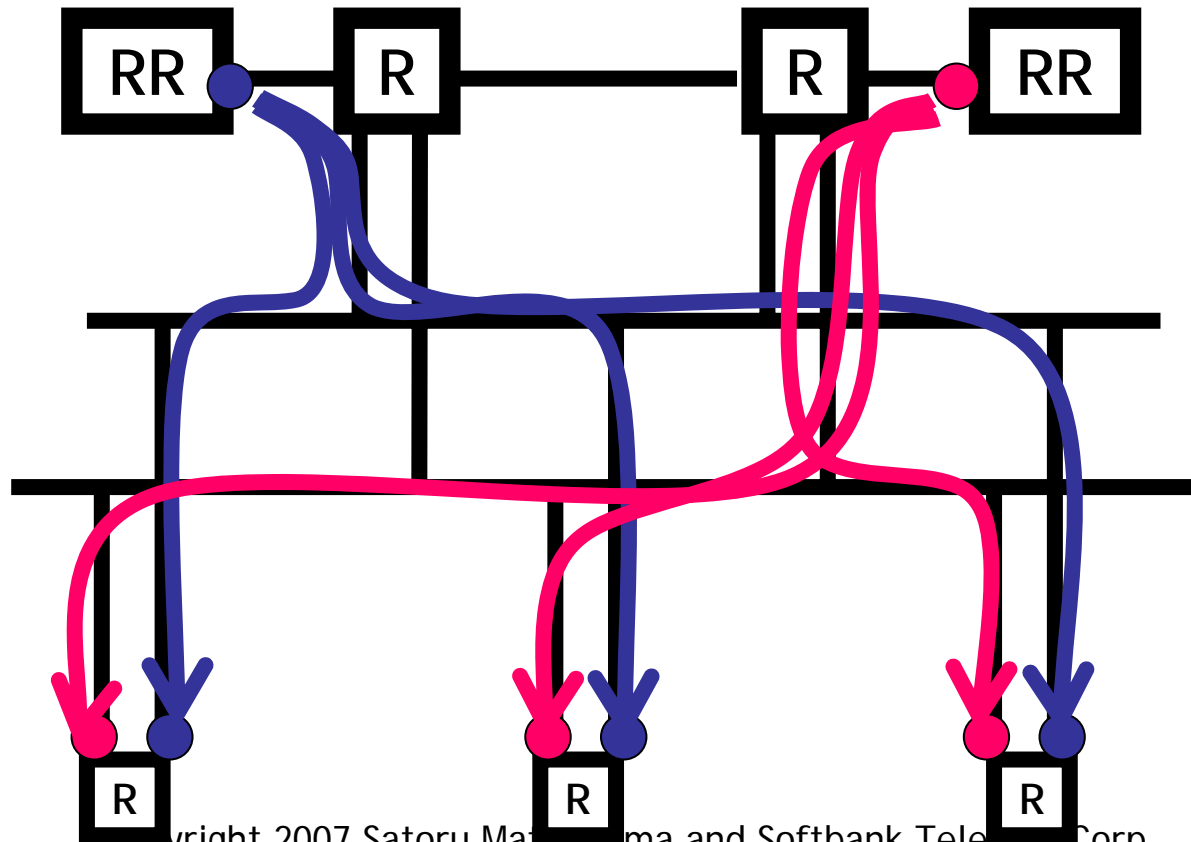
実際に構築してみよう

- プローブのかけ方



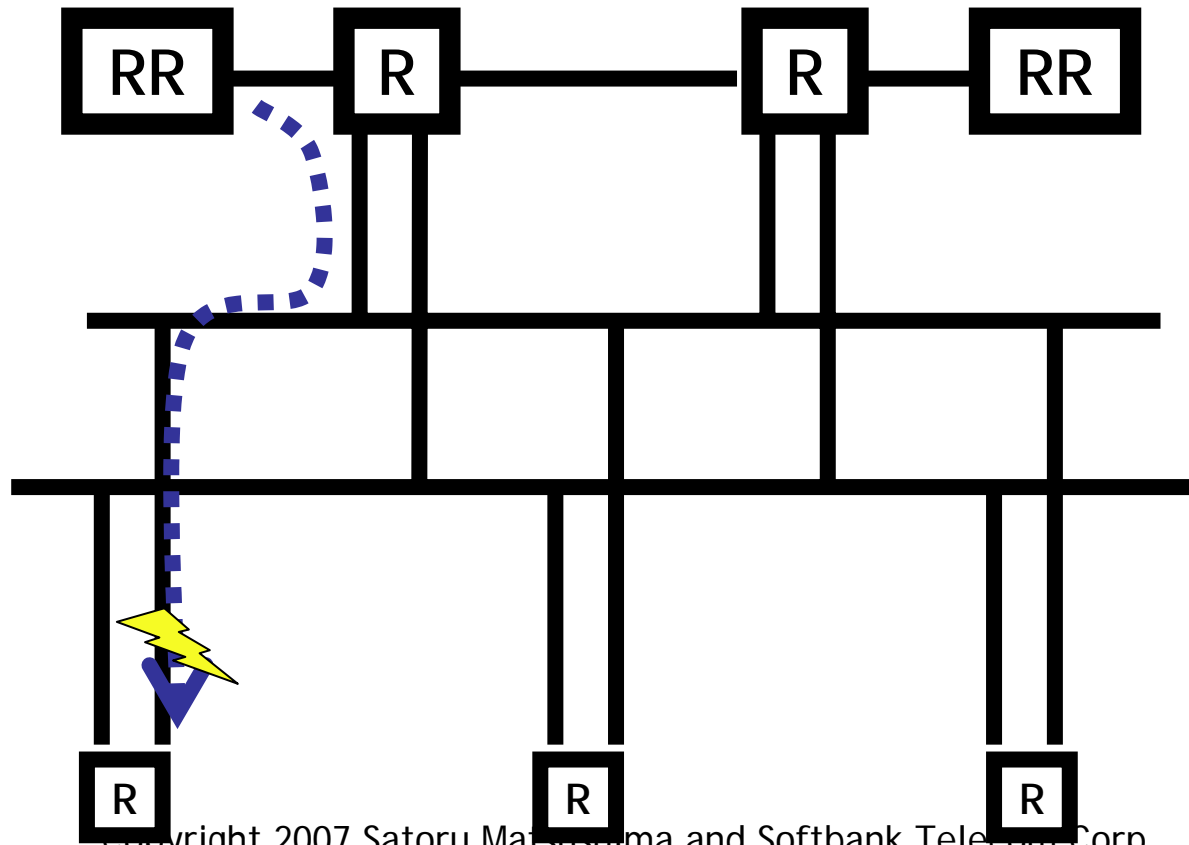
実際に構築してみよう

- プローブのかけ方 (BGPピアと同じ)



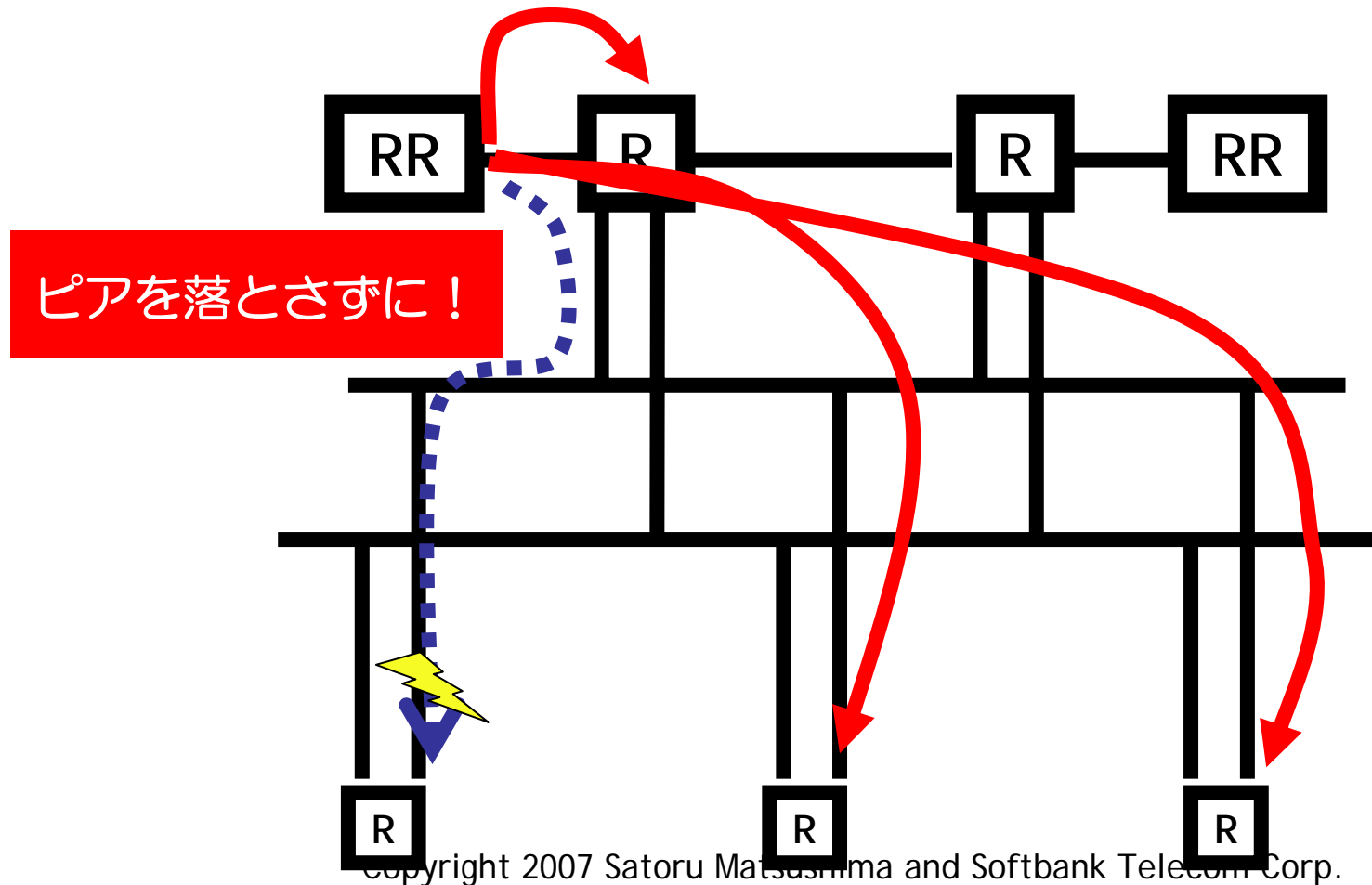
検証::障害検知時の動作

- プローブが通らなくなる



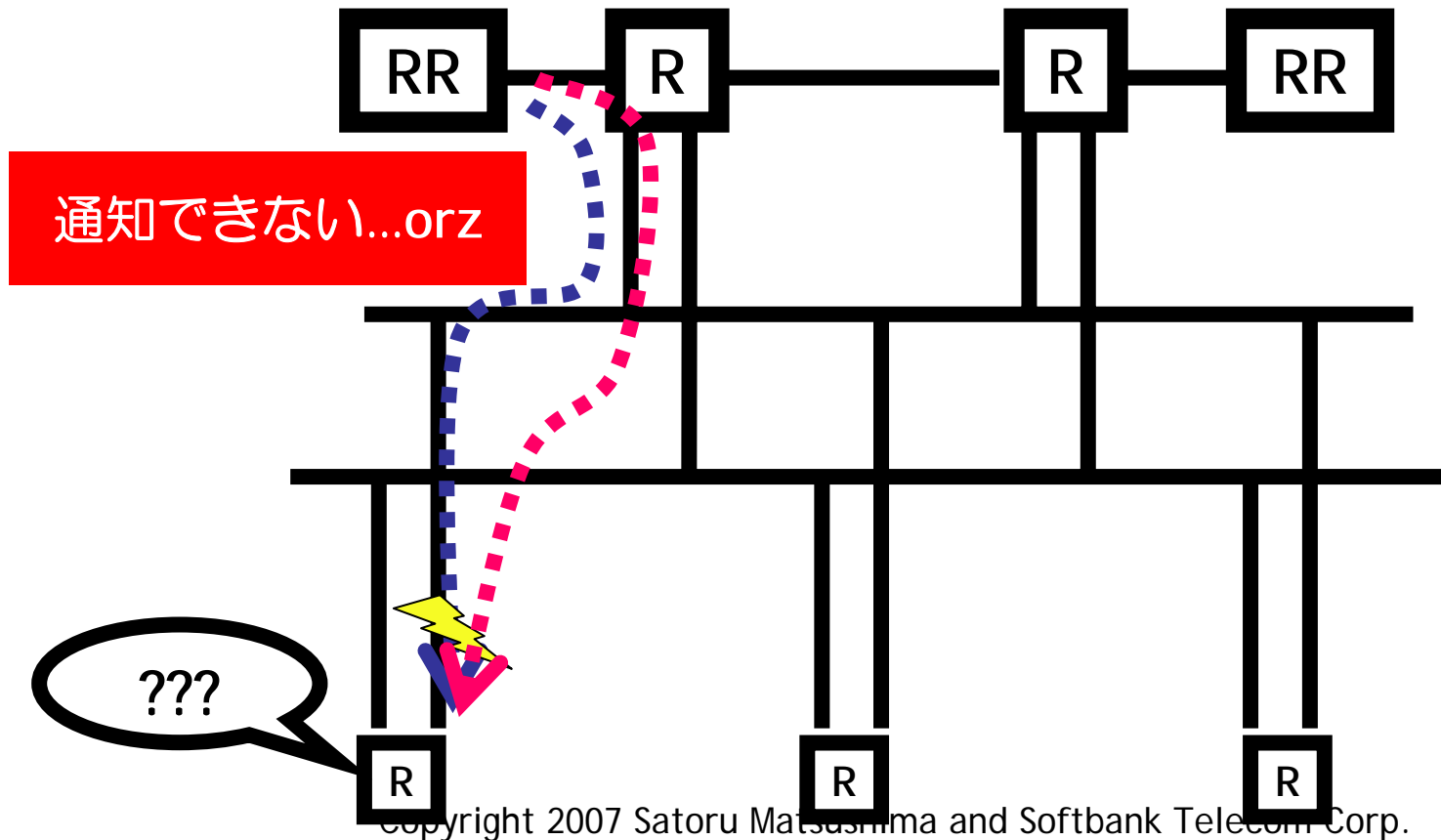
検証::障害検知時の動作

- すぐさま他のルータに通知 (withdraw)



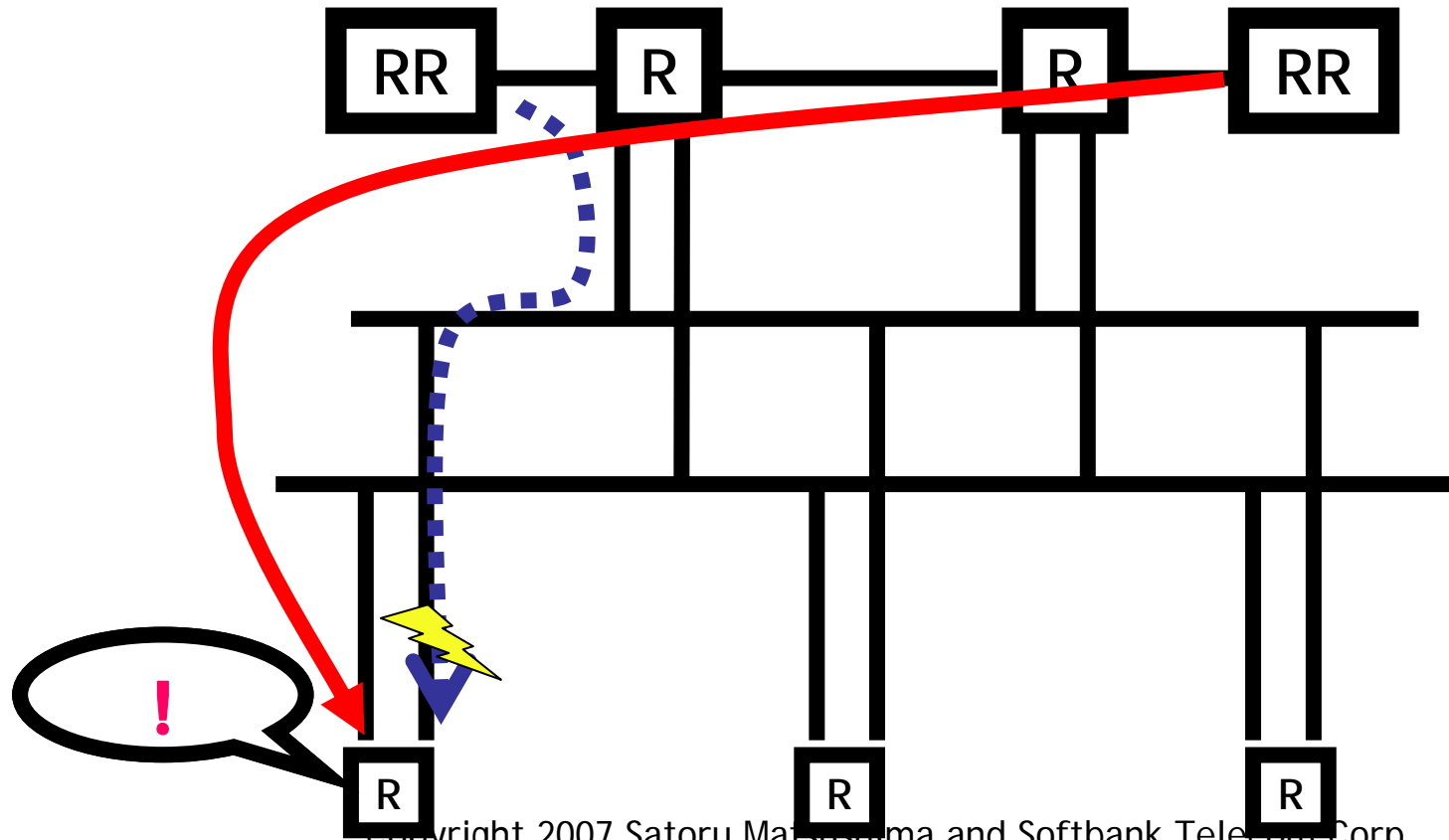
検証::障害検知時の動作

- 対象ルータが障害検知できないケース



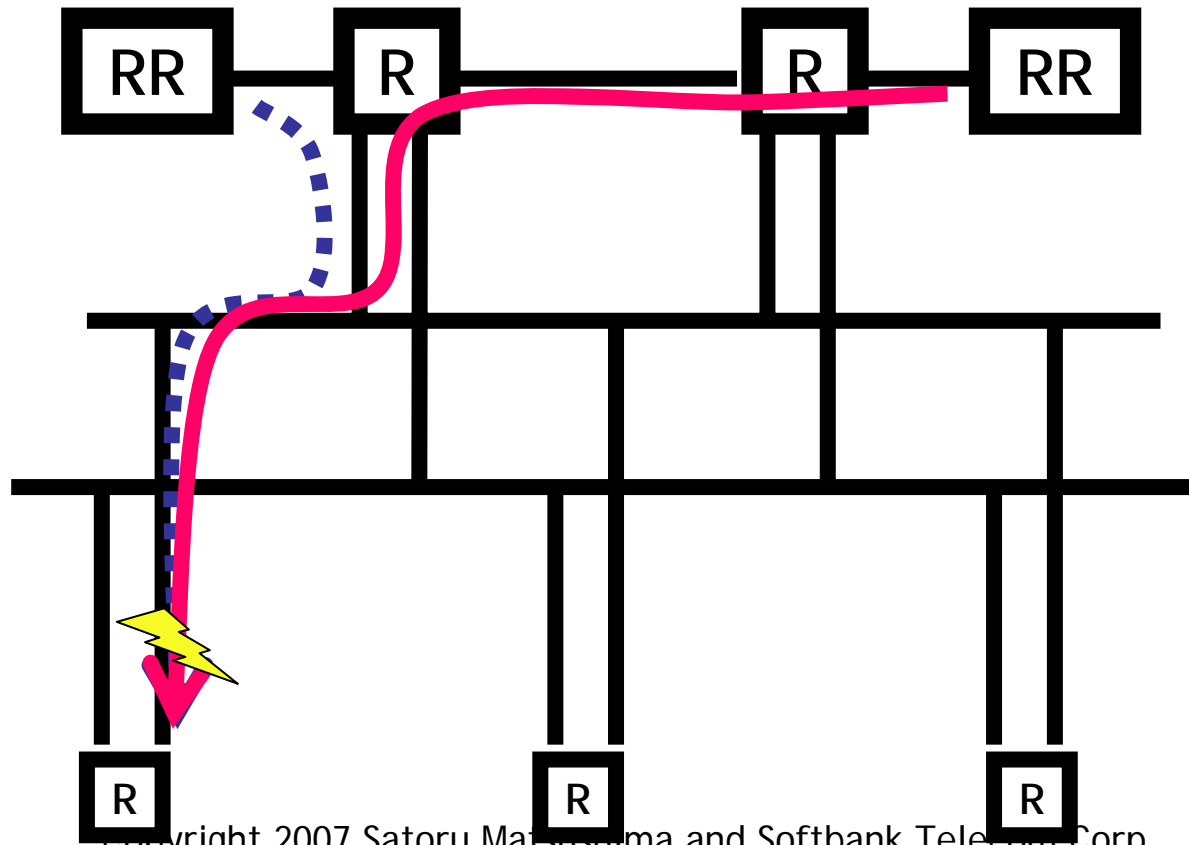
検証::障害検知時の動作

- 到達できるピアから通知してあげる
(予備経路を強める)



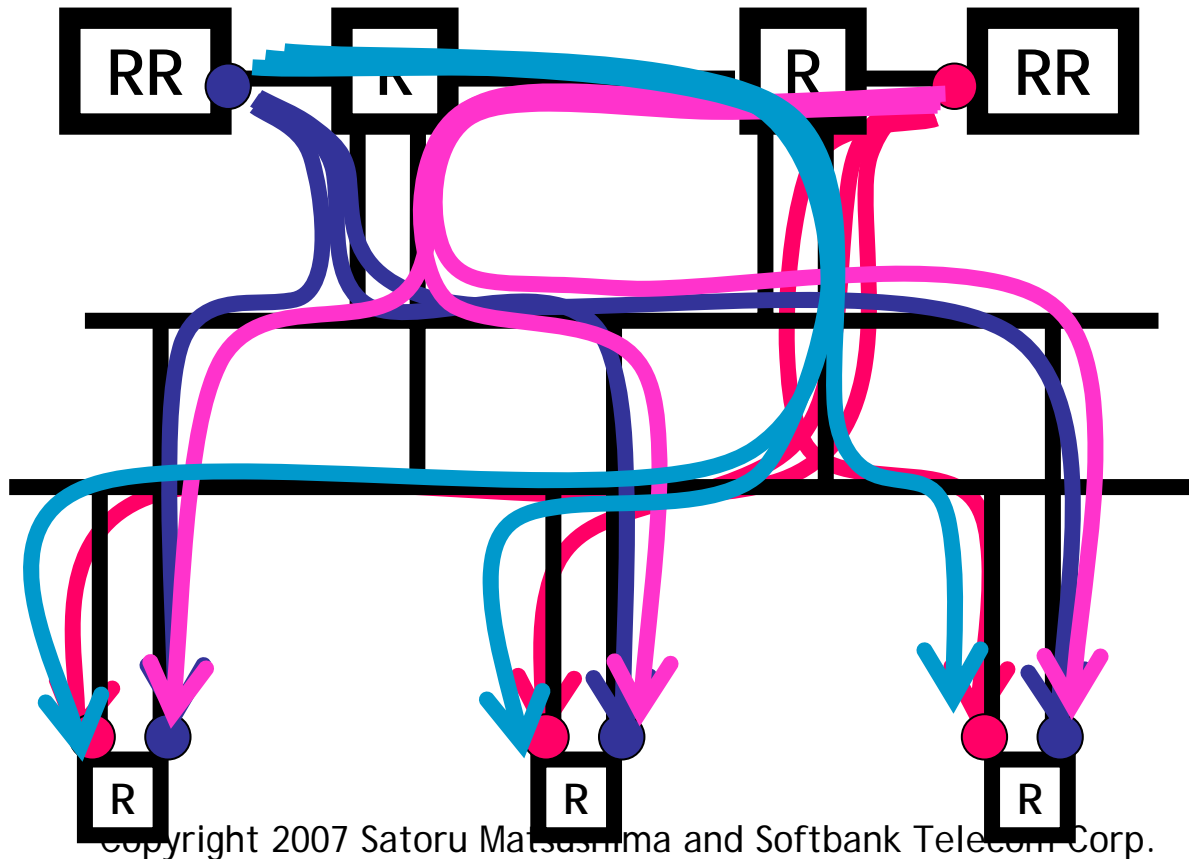
検証::障害検知時の動作

- プローブを追加してやる必要があった！



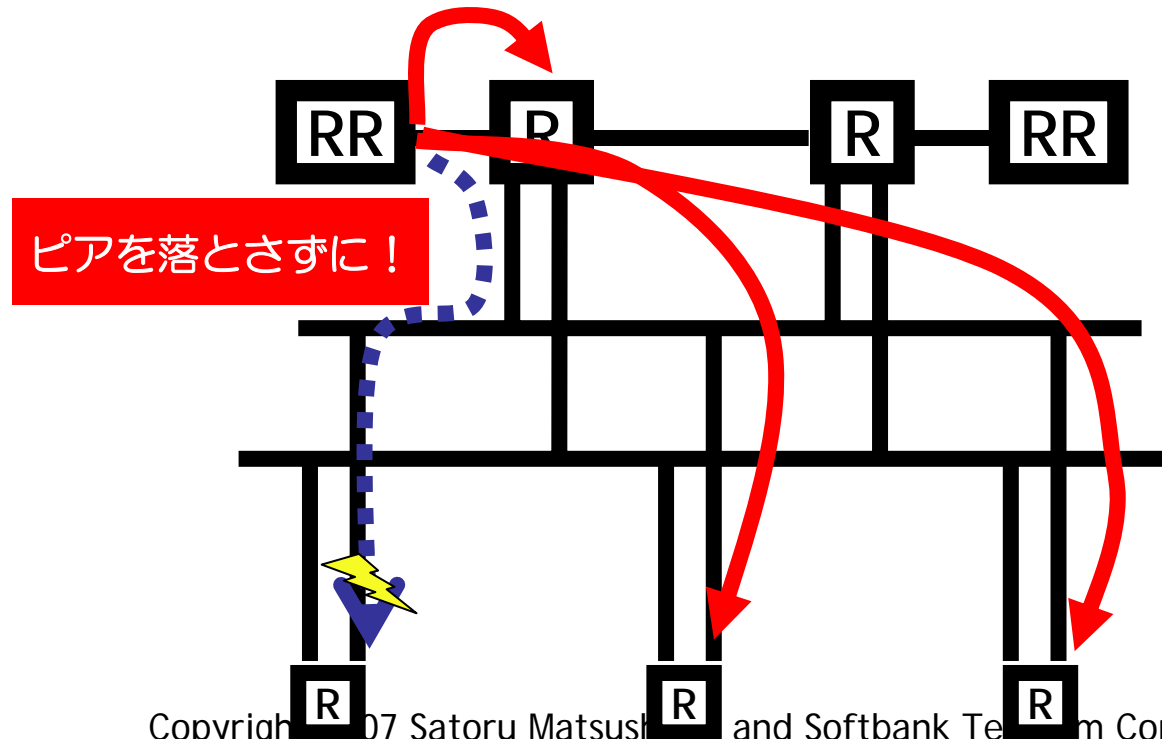
検証::障害検知時の動作

- やりなおし::プローブのかけ方



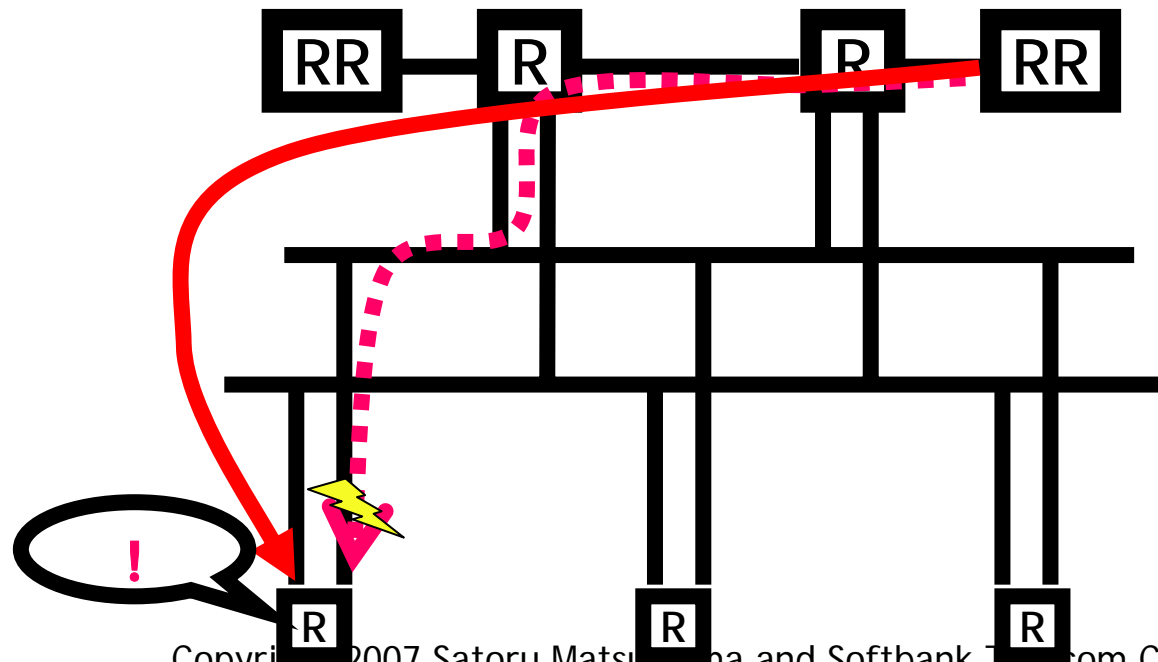
検証後の動作ポリシー

- 対象ルータ以外へwithdraw
 - BGPピアと同じパスのプロローブで検出した場合



検証後の動作ポリシー

- 対象ルータへ 予備経路を強めるupdate
 - BGPピアと異なる（別のRRの）パスのプローブで検出した場合



検証後の動作ポリシー

- プローブの保護時間
 - 障害検出するまで：
 - 障害復旧検出から実際に復旧させるまで：

**フローブを受け取る側と
フローブを送る側の
負荷をみなから慎重に**

まとめ::思いやる運用は実現できたか？

- BGP peer up/down はかなり抑制された
 - プロブがダウンしても、Peer downまで至らないケース多数
 - ICMPの負荷という話もあるが、小さい箱にハードウェアBFDが載らない限り原理は同じ
- かなりアクロバティックな経路制御に...
 - プロブとBGPピアのパスを揃えるのが大変
 - マルチホップのBFDもこんな感じなのか...
 - 予備経路がいくつもある場合、優先順位を考えて設定するのが大変
- プロブはBFDでやれる事に越したことはない
 - 実際にパケットを転送するルータ自身が、プロブをもつ
 - でもやっぱり気になる、能力差の大きいピア
 - マルチホップ環境での高速検知が必要なケースがあったら、またこれをやることになるのだろう...
- RRで経路制御ポリシーを作れるのはなかなかおもしろい
 - でも、ホントはRRはアトリビュートに変更を加えてはいけない

まとめ::終わりに

- BFD有史以前の、涙ぐましい高速障害復旧の実装・運用について紹介しました
 - 数千の小さい箱でも高速障害復旧
 - 汎用プロトコルをプローブにBFDライクな実装
 - アクロバティックな経路制御
- しかし、今回対象としたネットワークにBFDが導入されるのは果たしていつか？
 - BGPピアの能力差がはげしい環境
 - 低能カルータには厳しい経路数を扱う環境
- 本質的には、経路数爆発状態における高速障害検知技術の運用を今から考えておきなさい、ということなのかも