

Preconfiguring Robust Logical Topology with Autonomous Routing for Data Center Networks

データセンターネットワークのための環境変動に対応した自律的ルーティング用トポロジ構成手法

村田研究室
下間 雄太

2014/2/14 1

データセンター

多数のサーバーとサーバー間を結ぶネットワークで構成

- サーバー間の連携により多量のデータを処理
- 分散ファイルシステム、分散コンピューティングなど
- ネットワークが処理性能に与える影響大
- サーバー間の帯域の不足により、他サーバーとのデータの連携にかかる時間が増大

2014/2/14 2

データセンターネットワークにおける問題と本研究の目的

- トラフィック変動が頻繁に発生
- 故障が頻繁に発生

↓

- 環境変動発生時でもサーバー間に十分な帯域を確保することが必要

研究の目的：
データセンターネットワークにおいて頻発する環境変動に対応してサーバー間に十分な帯域を確保する経路制御手法の確立

- トラフィック状況・故障状況に応じて、通信サーバー間に十分な帯域を確保可能な経路を確立
- 頻繁な環境変動に瞬時に対応可能な自律分散型経路制御

2014/2/14 3

環境変動に対応した自律的な経路制御手法

アプローチ

複数の論理トポロジを用いた自律的経路制御手法による輻輳・故障箇所の迂回

- 論理トポロジを複数構成
 - 各論理トポロジは物理ネットワークの全ノードと一部のリンクから構成
 - いずれの論理トポロジを用いても全ノードに転送可能
- 各機器が論理トポロジを自律的に選択することにより、経路を決定
 - 輻輳・故障で利用不可なリンクを含まない論理トポロジのうち、次ホップへのリンクの負荷がもっとも低いものを選択
 - 利用不可な論理トポロジの情報をパケットに付与して転送することにより、他のノードにも伝達

2014/2/14 4

論理トポロジを用いた自律的な経路制御の動作例

リンク 4-7 で故障が発生した場合

- ノード 4 はパケットに論理トポロジ A と C が使用できないという情報を付与
- 論理トポロジ B を用いてノード 6 に転送
- その後の中継ノードはパケットを見て使用可能な論理トポロジ B のみでパケットを転送

ノードが保持する論理トポロジの集合

2014/2/14 5

論理トポロジが満たすべき性質

- 全ノード間に十分な迂回経路が存在
 - 故障・輻輳が発生した場合にも、代替経路を用いた通信路を確保することが必要
- 故障の影響が大きい機器は存在しない
 - 特定の機器の故障により、通信帯域が激減することは避けることが必要
- 迂回経路上でトラフィックが集中しない
 - 迂回先でのトラフィック集中により、サーバー間に確保できる帯域が制限されることは望ましくない

要件を満たすような論理トポロジの集合を構成する手法を検討

2014/2/14 6

論理トポロジの設計手順

1. 論理トポロジの候補の集合を生成

論理トポロジ候補:

- 物理トポロジ内のリンクのサブセットで構成された木構造の全パターン
- 各論理トポロジに含まれるリンク数は最小
→利用不可なリンクが生じた場合に、
利用不可なリンクが各論理トポロジに含まれる確率は小

2. 論理トポロジを選択し、経路制御に用いる論理トポロジの集合へ追加

- 各候補論理トポロジ追加時の論理トポロジ集合の適切性に関する指標を計算しもっとも良い論理トポロジを選択

指標

- 十分な迂回路数: サーバー間の迂回路数
- 機器の故障の影響: 平均ノード媒介中心性の最大値
- 故障迂回時の輻輳可能性: 平均リンク媒介中心性の最大値

$$B_n^{node} = \frac{1}{G} \left| \sum_{g \in G} \left(\sum_{s,d \in S} \frac{R_g^{node}(s,n,d)}{R_g(s,d)} \right) \right|$$

$$B_l^{link} = \frac{1}{G} \left| \sum_{g \in G} \left(\sum_{s,d \in S} \frac{R_g^{link}(s,l,d)}{R_g(s,d)} \right) \right|$$

2014/2/14

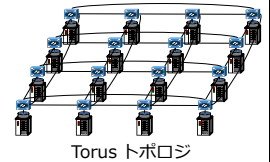
修士論文発表会

7

シミュレーション評価

評価環境

- 評価対象トポロジ
 - 4ポートのスイッチ16台で構築されたTorusトポロジ
 - 4ポートのスイッチ20台で構築されたFatTreeトポロジ
- 比較対象経路制御
 - Equal Cost Multi-Path (ECMP) [1]
- 各リンクの帯域
 - 10 Gbps
- トラヒックの発生方法
 - 10%のサーバーラック間で発生
 - 空き帯域がある限りトラヒックの送信量を増加させていく
- 故障の発生方法
 - 単一リンク故障を各リンクにおいて発生



Torus トポロジ

評価指標

- 通信サーバーラック間の通信帯域の最小値 (Gbps)

[1] C. HOPPS, "Analysis of an Equal-Cost Multi-Path Algorithm," RFC 2992, Internet Engineering Task Force, Nov. 2000. <http://tools.ietf.org/html/rfc2992>.

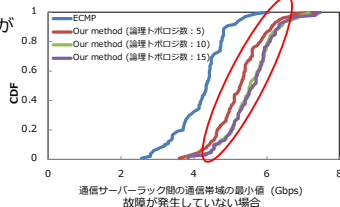
2014/2/14

修士論文発表会

8

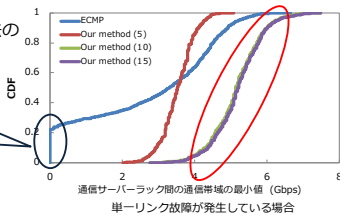
評価結果: 提案手法と ECMP の比較

- ECMP より提案手法の方が通信帯域が大きい
- 少量の論理トポロジでも ECMP より多く通信帯域を確保できている



- 故障時も同様に ECMP より提案手法の方が通信帯域が大きい

リンク故障の影響で通信に失敗したサーバー間が存在する



2014/2/14

修士論文発表会

9

まとめと今後の課題

まとめ

- 論理トポロジを用いた自律的な経路制御手法および論理トポロジの構築手法を提案
- 提案手法で ECMP より通信帯域を多く確保できることを確認
- FatTree においても Torus と同様 ECMP よりも通信帯域を多く確保できることを確認
- リンクの単一故障において全てのサーバーラック間で経路を確保できることを確認

今後の課題

- 規模の大きいネットワークでの評価

2014/2/14

修士論文発表会

10