

科学研究費助成事業 研究成果報告書

平成 26 年 6 月 6 日現在

機関番号：62615
研究種目：若手研究(A)
研究期間：2011～2013
課題番号：23680005
研究課題名(和文)大規模DNSトラフィックの統計的解析と異常検出

研究課題名(英文)An analysis of large scale DNS traffic

研究代表者

福田 健介 (Fukuda, Kensuke)

国立情報学研究所・アーキテクチャ科学研究系・准教授

研究者番号：90435503

交付決定額(研究期間全体)：(直接経費) 12,900,000円、(間接経費) 3,870,000円

研究成果の概要(和文)：本研究では、国内ccTLDサーバ(jp DNSサーバ)へ送信された全てのクエリデータの収集・解析することで、以下に示す研究成果を得た。(1) DNSキャッシュリゾルバによるDNSサーバ選択アルゴリズムのマクロレベル効率性解析。(2) DNSSEC検証を行うキャッシュリゾルバの推定アルゴリズムの開発。(3) キャッシュリゾルバでのDNSソフトウェア推定アルゴリズムの開発。(4) DNSSEC検証失敗時の影響評価。

研究成果の概要(英文)：We analyzed large scale DNS queries measured at JP ccTLD servers (JP DNS servers) in order to quantify interactions between DNS authoritative server and cache resolvers. The main results are as follows: (1) A macroscopic analysis of efficiency of DNS server selection algorithms in cache resolvers. (2) A development of a technique for counting DNSSEC validators. (3) A development of a technique for identifying DNS software at end-host. (4) An analysis of macroscopic effect of DNSSEC validation failure on cache resolver performance.

研究分野：計算基盤

科研費の分科・細目：情報ネットワーク

キーワード：インターネット トラフィック DNS

1. 研究開始当初の背景

DNSはホスト名からIPアドレスへの変換等を行うインターネットワイドでの分散データベースであり、現在のインターネットにおいて必要不可欠の機能である。DNSは階層構造を持ち、その頂点にあるルートサーバ群は世界的に分散配置されており、その下の階層では例えばccTLDの一つである“.jp”は日本のレジストリ機関によって管理されている(jpサーバ)。これらのjpサーバは実際には複数のサーバおよびそのレプリケーションから構成されている。クライアントは経験的なアルゴリズムによって、より近いサーバへアクセスすることで効率的な応答性能を得ることができる。

(1) DNSサーバの返答は直接ユーザの感じる遅延に影響を及ぼすことから、サーバは低遅延・低負荷を実現するために自律分散的に管理・運用される。しかし、サーバソフトウェアの最適化等のサーバのミクロな効率については改良が行われているが、DNSの大規模環境でのマクロレベルの効率や異常トラフィックに関する振る舞いについてはほとんどわかっておらず、複製配置・チューニングは経験的に行われているのが現状である。

(2) また、2010年10月よりDNSセキュリティ機能拡張であるDNSSECが利用可能となったが、新しい機能でありトラフィックの増加を伴うことから、その影響を知ること、サーバインフラ等の構築運用を行うDNSオペレーション事業者にとっても重要である。これらの大規模なトラフィックデータは未だに世界的にも得られていないものである。

2. 研究の目的

(1) DNSの大規模環境でのマクロレベルの効率や異常トラフィックに関する振る舞いについては、データ収集の困難さもあり、いまだに十分な理解を得られていない。そのため、世界中からDNSサーバへ送信される全クエリデータをサーバ再度で収集・解析することで、国レベルでのマクロなDNSの挙動を明らかにする必要がある。これには、サーバ側でのデータ収集を行った後に、クライアントでのサーバ選択アルゴリズムのマクロな評価が必要となる。また、クライアント側でのDNSソフトウェアの実装による違いについて把握する必要がある。

(2) DNSの新しい機能拡張であるDNSSECの導入に伴って、新しいタイプのトラフィックが生成され始めている。また、DNSSECは選択的に導入されているため、どのクライアントが実際にDNSSECを使用しているのかは定かではない。そのため、DNSSEC使用クライアントを同定するアルゴリズムの開発および、DNSSECを導入することによって生じる問題点を明らかにする必要がある。

3. 研究の方法

研究に用いるデータセットとしては、DITL (Day in the life of the Internet) キャンペーン時に収集された全jpサーバへのクエリデータを用いる。このデータには、クライアントであるキャッシュリゾルバのIPアドレス、クエリ内容、およびサーバでの受信タイムスタンプが含まれる。前述のように、jpサーバは広域分散しているため、全サーバ数は20以上となる。DITLでは2日間の計測を行っており、データは48時間分利用可能である。解析には、2009-2013年のデータを使用した。

4. 研究成果

(1) 「クライアントによるサーバ選択アルゴリズムの効率性解析」

クライアントであるキャッシュリゾルバでは、DNSソフトウェアが適切なサーバをサーバ選択アルゴリズムによって選択する。このアルゴリズムはソフトウェアの実装によって異なるため、ミクロな実装の違いはわかっているものの、マクロな効率に関して情報はない。とりわけ、DNSは処理の軽いプロトコルとして設計されているため、サーバ側にて、サーバ・クライアント間の遅延を測定する方法がない。そのため、サーバサイドでクライアントへの遅延を近似的に計算するために、2点間の測地線距離を用いる手法を開発した。測地線距離を計算するためには、ランドマークとなる都市情報および、ランドマーク間の物理的なネットワークポロジ情報が必要となる。本課題では、インターネット上で使用可能な、光ファイバ敷設マップを用いることで、世界的な都市間の測地線距離を定義した。また、各サーバ・クライアントの位置座標はIPアドレスを緯度経度に変換するデータベースであるGeoIPデータベースを使用した。測地線距離を遅延の代替とすることは直感的には正しいが、実際にどの程度の誤差が含まれるかを調べるために、東京に設置されたサーバから、トラフィックログに現れるクライアントのIPアドレスの一部にpingコマンドを利用して遅延を測定した。その結果、直接的な距離と比べて測地線距離を用いることで遅延との相関が高くなることを確認した。

クライアントごとにjpサーバへのサーバ選択の効率性を計算するために、2つの効率性を表す指標を導入した。1つ目は、実際にアクセスのあったサーバのうちで最も測地線距離が小さくなるサーバアクセスを行った際の測地線距離の総和と、実際にアクセスを行ったサーバへの測地線距離の総和の比である。これは、現状のサーバ構成において物理的に可能なアクセスと実際に生じたアクセスの比であることから、現状の構成を用いた際の最適値とのずれを表している。2つ目は、実際にはアクセスがなくても最も近いサーバへのアクセスを行った際の測地線距

離の総和と、実際にアクセスのあった測地線距離の総和を比較するものである。これは、理想状態であり、IP ルーティングの複雑さを無視したものとなる。実際のトラフィックデータを用いて、サーバ選択の効率性を計算したところ、日本に存在する 75%のキャッシュリゾルバ、40-60%の海外に存在するキャッシュリゾルバは、サーバ選択アルゴリズムにより、アクセスがあったサーバの中で最も近いサーバにアクセスしていることが確認できた。つまり 25%もしくは 60%のサーバでは、サーバ選択アルゴリズムの改良によりより近いサーバを選択することで効率を改善できる可能性を示した。同様に、15-35%の海外に存在するキャッシュリゾルバでは、理想的なサーバへのアクセスが実際に行われていることがわかった。この効率性はそれほど高いものとは言えないが、これらの改善にはサーバの増設や、インターネットトポロジの改善が必要であることが示された。

(2) 「DNSSEC 検証ホストの推定アルゴリズムの開発」

DNSSEC は DNS へのセキュリティ拡張の一つとして知られているが、キャッシュリゾルバが実際に DNSSEC を使用するかどうかはキャッシュリゾルバの設定に依存している。2010 年より jp サーバでも DNSSEC が利用可能となり、DITL データを観測することで、DNSSEC クエリを用いたキャッシュリゾルバの増加を観測している。しかしながら、通常、DNSSEC 処理を行う際には、2 つのクエリ(DS, DNSKEY)をサーバへ送信する必要があるが、実際のトラフィックにはそのどちらか一方のみが送られている場合が多く存在することが明らかになった。これらのキャッシュリゾルバが実際に DNSSEC 処理を行っているかどうかは定かではないことから、キャッシュリゾルバの DNSSEC 使用パターンを解析し、DNSSEC 使用ホストを同定するアルゴリズムを開発した。メインのアイデアは、DNSSEC クエリ(DS レコード)が送信される場合には、元々必要とされるクエリ(A や AAAA レコード)が送信されるため、DNSSEC クエリと元のクエリの共起関係に着目することにある。実際には、クエリの共起関係を全てのクエリに対して調査することは計算コストが大きくなることから現実的ではないため、DNSSEC クエリ数と元のクエリ数の比率(DSratio)によって DNSSEC 処理を行っているかを判定する。

評価を行うにあたって、正解データにあたる情報が必要であることから、DITL トレース中より、1 つ以上 DNSSEC クエリを送信したキャッシュリゾルバを抽出し、その IP アドレスに対して DNSSEC 要求を付したクエリを送ることで実際に DNSSEC 処理が行われているかについて確認した。この試行に対して、ほぼ 10%のキャッシュリゾルバが返答を返し、その約 60%が DNSSEC を実際に使用していることを確認した。同定アルゴリズムは上記

DSratio の他に、5 つの DNSSEC トラフィックに関係するトラフィック特徴量をキャッシュリゾルバごとに抽出し、機械学習アルゴリズムの一つである分類木を用いて同定を行った。その結果、同定に関しては、6 つの特徴量のうち DSratio が最も大きな寄与を示すことが確認できた。また、同定精度(f-measure)は、0.8-0.9 と高い精度で分類可能であることが示された。1 つ以上 DNSSEC クエリを送った潜在的なキャッシュリゾルバの約 70%が実際の DNSSEC 処理を行っているとの推定値を得た。また、DNSSEC 処理を行っていないキャッシュリゾルバの IP アドレスを調べたところ、多くのアドレスは、パブリック DNS サービスを提供している組織であることが明らかとなった。

(3) 「キャッシュリゾルバでの DNS ソフトウェア推定アルゴリズムの開発」

DNS サーバサイドでは、キャッシュリゾルバで使用されているソフトウェア(DNS ソフトウェア)を知ることは通常困難である。これは、DNS が軽量なプロトコルであることに起因する。本研究課題では、DNS サーバに存在するキャッシュリゾルバの DNS クエリパターンより、キャッシュリゾルバで使用されている DNS ソフトウェアを推定するアルゴリズムを開発した。推定アルゴリズムを開発するにあたり、代表的な DNS ソフトウェアである、BIND、UNBOUND、Windows DNS server をバーチャルマシン上にインストールし、典型的なクエリ送信時のクエリパターンを収集し、それらの結果より、経験的な 15 の分類ルールを構築した。これらのルールは、A/AAAA レコードをクエリに含んだ場合に、追加的に NS レコードをどのようにクエリするか、A もしくは AAAA のクエリの場合、どのように AAAA、A のクエリを追加的にクエリするか等、一つのトリガとなるクエリとそこから生成される付加的なクエリの組み合わせで表現される。

これらの分類ルールを同定ツールとして実装し、国内バックボーンで収集された DNS クエリを用いて評価を行った。正解データとしては、既存手法である該当クエリへの Chaos クエリの送信の結果を使用した。この返答のうち、正しいと思われる結果についてのみを抽出した。ツールをこれらの正解データのクエリログを適用したところ、99%の精度で DNS ソフトウェアを同定することができた。また、従来手法で同定することができなかったホストのうち 78%のホストに関して、DNS ソフトウェアを同定することが可能となった。さらに、分類ルールを解析結果より再調査したところ、主要なルールは、IPv6/IPv4 に関するものであること、分類能力が高いと予想される DNSSEC に関連したルールは、DNSSEC の普及が進んでいないため、現時点では、ほとんどルールとして使用されていないことが明らかとなった。

(4) 「DNSSEC 検証失敗時の影響評価」
DNSSEC は今後のインターネットセキュリティを考える上ではキーとなる技術である。現在の DNS では、キャッシュリゾルバから送信されたクエリが正しい権威サーバによって返答されたかどうかを検証することが困難であるが、DNSSEC では、権威サーバのツリーをルートサーバから順番にたどることによってその一連の返答が正しい権威サーバからのものであることを保証する。そのため、権威サーバ間では、検証を行うための鍵となるレコード(DS, DNSKEY)を正しく設定する必要がある。仮にこれらが不一致になったとすると、DNSSEC を用いたクエリは全て正しい返答を得ることができなくなる。実際に、これらの不一致をおこす例が世界各地で起きていることから、これらの検証失敗の影響を推定することがネットワーク管理者には必要となる。

本課題では、DITL の DNS クエリデータから実際のキャッシュリゾルバのアクセスパターンを抽出し、jp サーバのいくつかに障害が発生した際の、DNSSEC 検証失敗の影響を評価した。その結果、18%のキャッシュリゾルバ(および 70%の AS)は、DNSSEC 障害の最初の 10 分間で使用不可能となることが示された。同様に 50%のキャッシュリゾルバは 6 時間以内に使用不可能となる。しかしながら、現在の jp ゾーンにおける DNSSEC の普及割合は低いため、障害の最初の 10 分間で影響を受けるキャッシュリゾルバは 0.8%程度であると予想される。さらに、ドメイン名の数とそのクエリ頻度はべき分布でモデル化できることから、人気のあるドメイン名に関する障害はそうでないドメイン名に比べて非常に大きな影響があることがわかった。

5. 主な発表論文等

(研究代表者、研究分担者及び連携研究者には下線)

〔雑誌論文〕(計 1 件)

1. 浅井大史, 福田健介, 江崎浩, DNS 探索グラフによる IPv4/IPv6 トランスポートを考慮した DNS 委譲構造の分析, コンピュータソフトウェア, 査読有, vol.30(2), pp.135-146, 2013

〔学会発表〕(計 5 件)

1. Kensuke Fukuda, Shinta Sato, Takeshi Mitamura, Towards Evaluation of DNS Server Selection with Geodesic Distance, Proc. of IEEE/IFIP NOMS2014, 査読有, p.8, Krakow, Poland, May 7, 2014

2. Ruetee Chitpranee, Kensuke Fukuda, Towards Passive DNS Software Fingerprinting, Proc. of AINTEC2013, 査読有, pp.9-16, Chiang Mai, Thailand, Nov. 14, 2013

3. Kensuke Fukuda, Shinta Sato, Takeshi Mitamura, A Technique for Counting DNSSEC

Validators, Proc. of IEEE INFOCOM2013, 査読有, pp.80-84, Turin, Italy, Apr. 15, 2013

4. 浅井大史, 福田健介, 江崎浩, DNS 探索グラフの可視化と解析, Proc. of WIT2011, 査読有, p.8, 札幌, Jun. 3, 2011

5. Kensuke Fukuda, Shinta Sato, Takeshi Mitamura, Preliminary Evaluation of Potential Impact of Failure in DNSSEC Validation, Proc. of Workshop of DNS Health and Security (DNS-EASY2011), 査読有, p.13, Rome, Italy, Oct. 18, 2011

6. 研究組織

(1)研究代表者

福田 健介 (FUKUDA, Kensuke)

国立情報学研究所・アーキテクチャ科学研究系・准教授

研究者番号：90435503