

# tcp/53接続を舐めて痛い目にあった話

---

2023/06/23

DNS Summer Days 2023 LT

株式会社大塚商会

たよれーるマネジメントサービスセンター システム運用課

高野 遼

# 自己紹介

- 名前：高野 遼
- 会社：株式会社大塚商会（9年目）  
各種広告等から、たのめーる（オフィス用品の通販サイト）やコピー機販売の会社と思われがちかもしれませんが、他にもWebホスティング、SOC、IaaSなど、インターネットにまつわる様々なサービスを展開しています。
- 所属：たよれーるマネジментサービスセンター システム運用課  
たよれーるマネジментサービスセンターは約300人が所属する組織。大塚商会が提供する一部サービスの開発、運用を担当。
- 自身の担当業務：所属組織向けにインフラサービスを提供（この一つにDNS（キャッシュおよびコンテンツ）が存在）および所属組織の情シス  
DNSは入社以来担当を変わらず。DNS Summer Daysには2015年以來参加を継続。

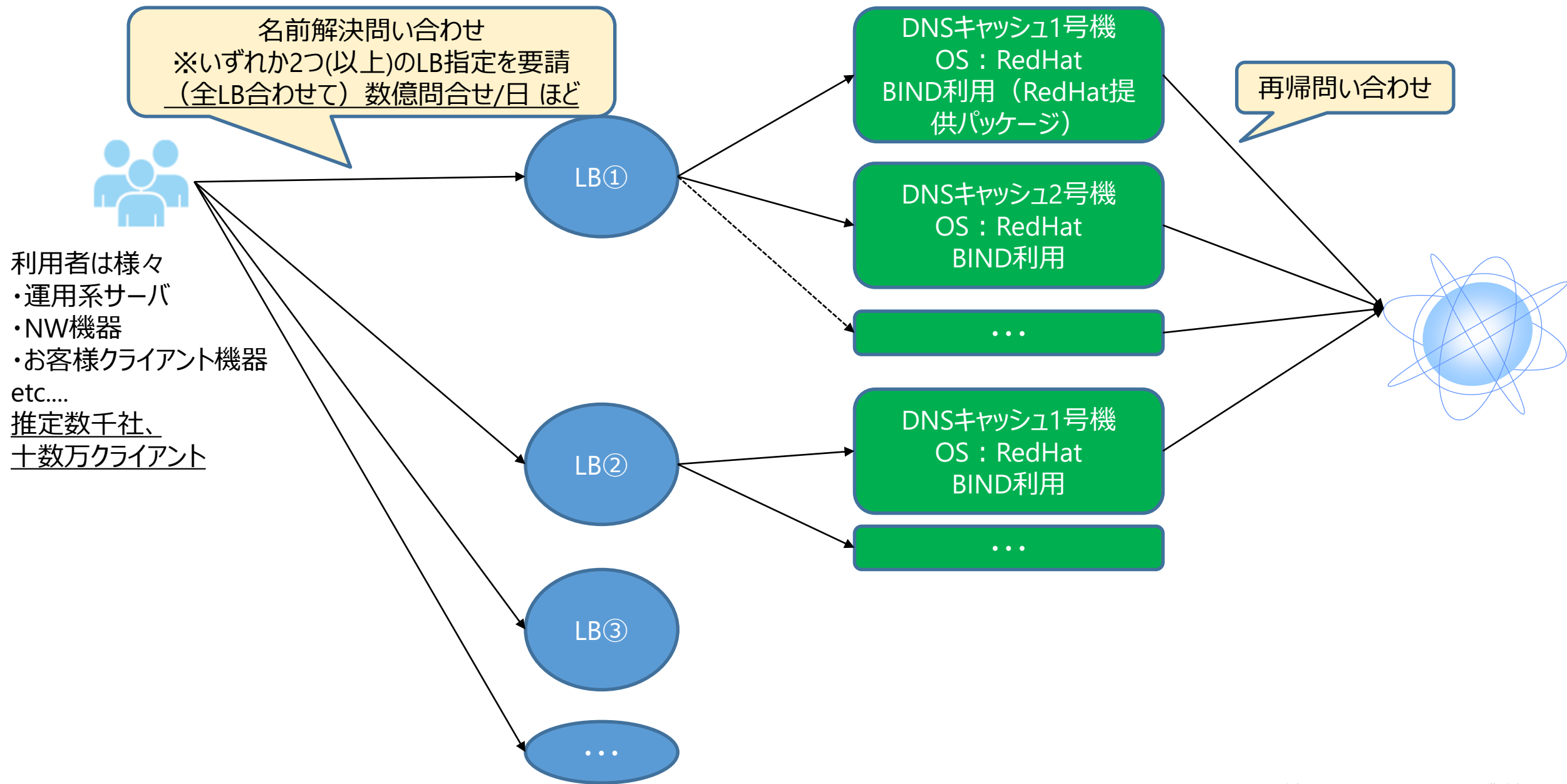
# 前置き

- 本セッションはあくまで事例紹介です  
tcp/53問い合わせが増えた（増えている）推察は記載していません。  
気になる方は、本年（2023年）のDNS Summer Days午前中のIIJ様の資料をご参照ください。
- 後にある通り、本件は現在進行形で対応中です。アドバイス等あれば、是非お願いしたい限りです

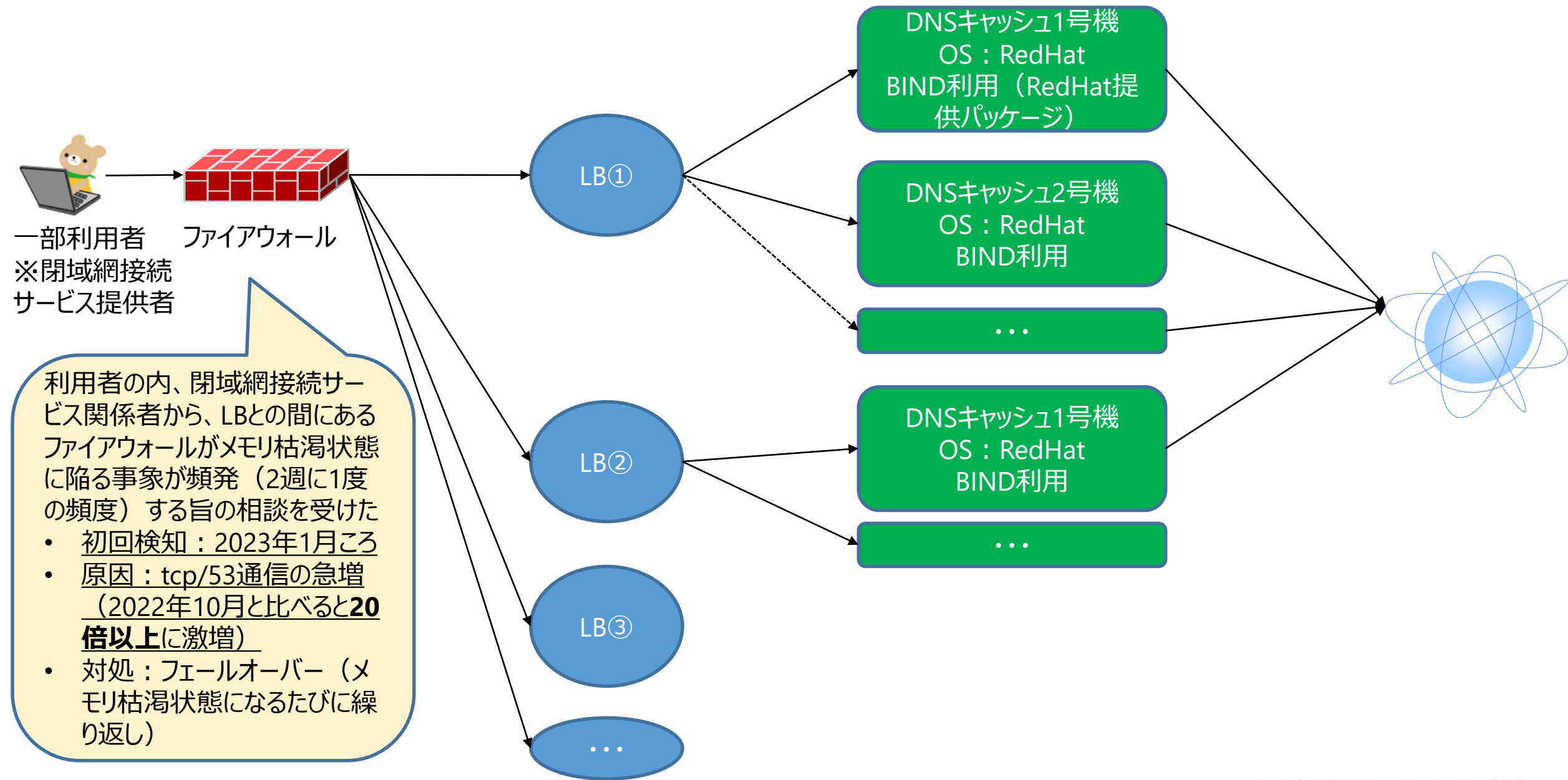
# アジェンダ

- 管理DNSキャッシュ 簡易構成図
- 問題発覚
- 調査
- 原因と対処
- 最初からtcp/53による問い合わせ自体は、RFC違反ではない（問題ない）
- LTで伝えたかったこと
- 参考①：あるDNSキャッシュサーバのtcp/53の接続数推移
- 参考②：あるDNSキャッシュのtcp問い合わせレコードTOP10 昨年との比較
- 参考③：TIME\_WAITに関して（※RedHat利用の場合）

# 管理DNSキャッシュ 簡易構成図



# 問題発覚



- 利用者の内、閉域網接続サービス関係者から、LBとの間にあるファイアウォールがメモリ枯渇状態に陥る事象が頻発（2週に1度の頻度）する旨の相談を受けた
- 初回検知：2023年1月ころ
  - 原因：tcp/53通信の急増（2022年10月と比べると**20倍以上**に激増）
  - 対処：フェールオーバー（メモリ枯渇状態になるたびに繰り返し）

# 調査

- 試しにtcp/53で名前解決問い合わせをテスト。すると、回答は得られるものの遅延が頻発長くて数秒程度。udp/53による同一FQDNの問い合わせはすぐ回答が返るので、明らかにtcp/53が原因。
- DNSキャッシュホスト上で確認すると、大量の TCP TIME\_WAITの行がクライアントIPアドレスは一意でなく様々。

- namedログを見ると、下記のようなログが頻出

```
Nov 14 08:38:05 XXX named[2894]: 14-Nov-2022 08:38:05.143 client: warning: client @0x7f41617edfb0 (no-peer): TCP client quota reached: quota reached
```

- 「rndc status」コマンド出力結果の内 “tcp clients” が、上限の100に張り付いていた
- tcp/53による名前解決問い合わせは確かに増えているが、20倍までにはなっていない

**⇒今回の事象の原因は、DNSキャッシュサーバのnamedプロセスのTCP接続上限に達してしまい、tcp/53による問い合わせクライアントが上限が空くまで接続リトライを繰り返しかけていたため、と結論**

# 原因と対処

## (原因)

- ① DNSキャッシュサーバではTCP接続上限に関するオプション：tcp-clients の明示はしていなかった（デフォルト値の"100"が設定されていた）

このオプションの存在および内容は知っていたが、リリース当時（2020年ごろ）のtcp/53接続数（数十回/日 程度）から、このオプションのチューニングは不要と判断してしまった。

- ② TIME\_WAITが多発していた

RedHatはTIME\_WAIT発生から60秒間、TIME\_WAIT状態を保持する仕様。

## (対処)

- tcp-clients を明示的に設定し、制限を大きく緩和
- 上限に達した際に検知ができるよう、監視を強化  
→ログファイルへのTCP client quota reachedの出力を監視。
- (今後) TIME\_WAIT発生要因を究明し対処  
→LBのTCP Profileが原因とみている。調査中。



# 最初からtcp/53による問い合わせ自体は、RFC違反ではない（問題ない）

- （大）昔は、「最初はudpで問い合わせるべき」のような記載だった（tcpはサポート“すべき”レベルだった）

## RFC 1123 6.1.3.2 Transport Protocols

DNS resolvers and recursive servers **MUST** support UDP, **and SHOULD support TCP**, for sending (non-zone-transfer) queries. Specifically, a DNS resolver or server that is sending a non-zone-transfer query **MUST** send a UDP query first. **If the Answer section of the response is truncated and if the requester supports TCP, it SHOULD try the query again using TCP.**

（意訳）

DNS リゾルバーと再帰サーバーは、送信（非ゾーン転送）クエリ機能として、UDPは必須、**TCPはできる限りサポートしなければならない**。具体的には、DNS リゾルバーまたはサーバーが非ゾーン転送クエリは、最初に UDP クエリを送信する必要がある。**もし応答の回答セクションが断たれた場合かつクライアント側が TCP による問合せを可能としている場合、次にTCPを使用して問い合わせを再試行すべきである。**

- 2016年に発表されたRFC1123のアップデート版：RFC7766 で、最初からtcpによる問合せもOKとなった

## RFC7766 5. Transport Protocol Selection

Stub resolvers and recursive resolvers **MAY** elect to send either TCP or UDP queries depending on local operational reasons. **TCP MAY be used before sending any UDP queries.**

（意訳）

スタブ リゾルバと再帰リゾルバは、ローカルの運用上の理由に応じて、TCP クエリまたは UDP クエリの送信を選択してもよい。**UDPクエリ送信の前にTCPクエリを使用してもよい。**

# LTで伝えなかったこと

- 名前解決問い合わせ=udp/53の考えはもう古い  
今回の事象発生までの間は、少なくとも自身はそう思っていた（tcp/53は、udp/53で解決できない場合くらいでしか使用されないと思っていた）。  
TCPに関わる設定は、DNSオプションだけでなくカーネルパラメータ、LBを利用している場合はTCP Profileなど様々。既に動いているものも、今まで意識していなければ今一度TCP周りの設定見直しを。
- TCP接続上限に達していたとしても、（きちんと監視していないと）すぐには検知されない、心配な方は一度ご確認を  
多少の遅延で済んでしまうため、名前解決問い合わせのレスポンスタイム監視だけでは引っかけられない可能性も。  
相談が無かったら、今も気づいていなかった...？
- ブラウザが重いなあ... と思ったら一度tcp/53問い合わせテストを  
tcp/53による問い合わせ自体を受け付けていないDNSサーバもいるとかいないとか。  
指定しているDNSサーバがtcp/53に対応していないことが原因の可能性も。  
心当たりのある方は一度テストをし、万が一tcp/53問い合わせが失敗したのならDNSサーバ運用者へ問い合わせを。

## 参考① : あるDNSキャッシュサーバのtcp/53の接続数推移

集計日付	tcp/53 問い合わせ数	全問い合わせでtcp/53が占める割合
2022/05/26(木)	52,019	0.38%
2022/06/23(木)	47,820	0.35%
2022/07/28(木)	168,125	1.23%
2022/08/25(木)	157,559	1.12%
2022/09/22(木)	24,319	0.17%
2022/10/27(木)	289,537	2.04%
<b>2022/11/24(木)</b>	<b>783,392</b>	<b>4.99%</b>
2022/12/22(木)	1,062,519	6.87%
<b>2023/01/27(木)</b>	<b>2,174,765</b>	<b>11.46%</b>
2023/02/22(水)	2,025,857	10.84%
2023/03/23(木)	2,127,693	10.61%
<b>2023/04/27(木)</b>	<b>1,269,677</b>	<b>6.73%</b>
2023/05/25(木)	690,440	3.81%

- 問題発覚は2023年1月だったが、それ以前の2022年11月にはtcp/53問い合わせ増加が始まっていた。
- 問題発覚の2023年1月は、問い合わせ数も率も高かった。
- 2023年4月以降はなぜか減少傾向。

## 参考②：あるDNSキャッシュのtcp問い合わせレコードTOP10 昨年との比較

2022/06/09(木)			2023/06/08(木)		
tcp/53による問い合わせレコード	TYPE	回数	tcp/53による問い合わせレコード	TYPE	回数
login.live.com	A	25012	login.live.com	A	16652
b1sync.zemanta.com	A	4556	s.yimg.jp	A	16063
assets.pinterest.com	AAAA	542	www.bing.com	A	13452
res.cdn.office.net	AAAA	475	signaler-pa.clients6.google.com	A	10768
ag.gbc.criteo.com	A	455	signaler-pa.clients6.google.com	TYPE65	10482
gem.gbc.criteo.com	A	441	www.google.com	A	10223
dns.msftncsi.com	A	239	www.bing.com	TYPE65	9407
officecdn.microsoft.com	AAAA	167	ssl.gstatic.com	A	9250
ims-prod07.adobelogin.com	A	141	ssl.gstatic.com	TYPE65	9103
login-us.microsoftonline.com	A	117	googleads.g.doubleclick.net	A	9014

- 1年の間で、tcp/53問い合わせレコードTOP10の内容が大きく変わっている。
- ただしlogin.live.comの問い合わせ回数は大差なし。2022年6月以前から、tcp/53で接続する仕様になっていると推測（なお tcp/53:udp/53=1:2 程度の比率）。回答サイズは小さいのでEDNSではない。

## 参考③ : TIME\_WAITに関して (※RedHat利用の場合)

- TIME\_WAITの保持時間は60秒。これはカーネルソース内でハードコードされているため変更不可 (らしい)
- LBのTCP ProfileがTIME\_WAIT多発の原因と睨んでいる理由  
<https://access.redhat.com/solutions/536113>  
LBの設定に依っては、TCPのタイムスタンプが不適切なものになってしまうとのこと。
- net.ipv4.tcp\_fin\_timeout はTIME\_WAITの保持時間の設定ではない  
このオプションは FIN\_WAIT2 状態にあるソケットの寿命であり、TIME\_WAITの保持時間設定で無い。  
TIME\_WAITに関するブログでいくつか見かけたので、念のため記載。
- TIME\_WAIT 状態のソケットを再利用できる net.ipv4.tcp\_tw\_recycle オプションは、RedHat8より廃止に  
<https://access.redhat.com/ja/solutions/6984791>  
なお、net.ipv4.tcp\_tw\_reuse はOutbound宛て通信でしか再利用できない。  
現在 net.ipv4.tcp\_tw\_recycle を有効化して運用されている場合は、OSリプレイス等の際に設計の見直しを。