

ATR: Additional Truncated Packets for Large DNS Response

Linjian Song, Shengling Wang

ABSTRACT

Internet works on reliable naming and IP package transmission. As complexity added into Internet infrastructure due to security and stability consideration, some mismatch and unexpected fragility are exposed. One case is in the field of DNS (DNSSEC) and IPv6. There are some public evidence and concerns on IPv6 fragmentation issues due to larger DNS payloads over IPv6. Different from other measurement work, in this paper we proposed "glueless" measurement mechanism to analyze the end-to-end users experience worldwide. It is based on a live, real and global Ads system and shows that the IPv6 large packet drop rate is up to 38%. Our analysis shows the combination of DNSSEC, UDP and IPv6 is really not going to work very well. In the meanwhile, we advance a solution called ATR (Additional Truncated Response) as a fix on DNS server which requires no modification on users side (DNS resolver) and support incremental deployment. We also conducted a measurement which shows more than 68% impacted users can be relieved on DNS latency and failures due to large DNS response.

1 INTRODUCTION

Large DNS response is identified as a issue for a long time. There is an inherent mechanism defined in [RFC1035] to handle large DNS response (larger than 512 octets) by indicating (set TrunCation bit) the resolver to fall back to query via TCP. Due to the fear of cost of TCP, EDNS(0) [RFC6891] was proposed which encourages server to response larger response instead of falling back to TCP. However, as the increasing use of DNSSEC and IPv6, there are more public evidence and concerns on user's suffering due to packets dropping caused by IPv6 fragmentation in DNS due to large DNS response.

It is observed that some IPv6 network devices like firewalls intentionally choose to drop the IPv6 packets with fragmentation Headers[I-D.taylor-v6ops-fragdrop]. [RFC7872] reported more than 30% drop rates for sending fragmented packets. Regarding IPv6 fragmentation issue due to larger DNS payloads in response, one measurement [IPv6-frag-DNS] reported 35% of endpoints using IPv6-capable DNS resolver can not receive a fragmented IPv6 response over UDP. Moreover, most of the underlying issues with fragments are unrevealed due to good redundancy and resilience of DNS. It is hard for

DNS client and server operators to trace and locate the issue when fragments are blocked or dropped. The noticeable DNS failures and latency experienced by end users are just the tip of the iceberg.

Depending on retry model, the resolver's failing to receive fragmented response may experience long latency or failure due to timeout and retries. One typical case is that the resolver finally got the answer after several retries and it falls back to TCP after decreasing the payload size in EDNS0. To avoid that issue, some authoritative servers may adopt a policy ignoring the UDP payload size in EDNS0 extension and always truncating the response when the response size is large than a expected one. However one study [Not-speak-TCP] shows that about 17% of resolvers in the samples can not ask a query in TCP when they receive truncated response. It seems a dilemma to choose hurting either the users who can not receive fragments or the users without TCP fallback capacity. There is also some voice of "moving all DNS over TCP". But It is generally desired that DNS can keep the efficiency and high performance by using DNS UDP in most of time and fallback as soon as possible to TCP if necessary for some corner case.

To relieve the problem, this memo introduces an improvement on DNS responding process by replying an Additional Truncated Response (ATR) just after a normal large response which is to be fragmented. Generally speaking ATR provides a way to decouple the EDNS0 and TCP fallback in which they can work independently according to the server operator's requirement. One goal of ATR is to relieve the hurt of users, both stub and recursive resolver, from the position of server, both authoritative and recursive server. It does not require any changes on resolver and has a deploy-and-gain feature to encourage operators to implement it to benefit their resolvers.

2 RELATED WORK(TBD)

Fragmentation Considered Harmful in IPv4 (and IPv6)[9]
IP Fragmentation Considered Fragile

Many network operators filter all IPv6 fragments.[11]
two researchers utilized a measurement network to measure fragment filtering. They sent packets, fragmented to the minimum MTU of 1280, to 502 IPv6 enabled and reachable probes. They found that during

any given trial period, ten percent of the probes did not receive fragmented packets. [12]

Applications That Rely on Fragmentation (draft-bonica-6man-frag-deprecate) [13]

3 MEASUREMENT STUDY ON LARGE DNS RESPONSE ISSUE

We begin by study IPv6 large DNS response issue using a novel measurement approach. Basically, we want to use our measurements answer the simple question: How IPv6 DNS works worldwide with large DNS response?

3.1 A Case Study of Large DNS response

To illustrate this situation, here are two DNS queries, both made by a recursive resolver to an authoritative name server, both using UDP over IPv6.

Query 1:

```
$ dig +bufsize=4096 +dnssec 000-4a4-000a-000a0000-
b9ec853b-241-1498607999-2a72134a.ap2. dotnxdomain.net.
@8.8.8.8 139.162.21.135 (MSG SIZE rcvd: 1190)
```

Query 2:

```
$ dig +bufsize=4096 +dnssec 000-510-000a-000a-0000-
b9ec853b-241-1498607999-2a72134a.ap2. dotnxdomain.net.
@8.8.8.8 status: SERVFAIL (MSG SIZE rcvd: 104)
```

What we see here are two almost identical DNS queries that have been passed to Google's Public DNS service to resolve. In the first case, the DNS response is 1,190 octets long, and in the second case the response is 1,346 octets long. The DNS server is an IPv6-only server, and the underlying host of this name server is configured with a local maximum packet size of 1,280 octets. Therefore, in the first case the response being sent to the Google resolver is a single, unfragmented IPv6 UDP packet, and in the second case the response is broken into two fragmented IPv6 UDP packets. And it is this single change that triggers the Google Public DNS Server to provide the intended answer in the first case, but to return a SERVFAIL failure notice in response to the fragmented IPv6 response. When the local Maximum Transmission Unit (MTU) on the server is lifted from 1,280 octets to 1,500 octets, the Google resolver returns the server DNS response in both cases.

The only difference in these two responses is IPv6 fragmentation, but there is perhaps more to it than that.

IP fragmentation in both IPv4 and IPv6 "raises the eyebrows" of firewalls. Firewalls typically use the information provided in the transport protocol header of the IP packet to decide whether to admit or deny the packet. For example, you may see firewall rules admitting packets using TCP ports 80 and 443 as a way of allowing web

traffic through the firewall filter. For this process to work, the inspected packet needs to contain a TCP header and use the fields in the header to match against the filter set. Fragmentation in IP duplicates the IP portion of the packet header, but the inner IP payload, including the transport protocol header, is not duplicated in every ensuing packet fragment. Thus trailing fragments pose a conundrum to the firewall. Either all trailing fragments are admitted, a situation that has its own set of consequent risks, or all trailing fragments are discarded, a situation that also poses connection issues[5].

IPv6 adds a further factor to the picture. In IPv4 every IP packet, fragmented or not, contains IP fragmentation control fields. In IPv6 these same fragmentation control fields are included in an IPv6 Extension Header that is attached only to packets that are fragmented. This 8-octet Extension Header is placed immediately after the IPv6 packet header in all fragmented packets, meaning that a fragmented IPv6 packet does not contain the Upper Level Protocol Header starting at octet offset 40 from the start of the IP packet header. But in the first packet of this set of fragmented packets, the Upper Level Protocol Header is chained off the fragmentation header, at byte offset 48, assuming that there is only a Fragmentation Extension Header in the packet. The implications of this fact are quite significant. Instead of always looking at a fixed point in a packet to determine its upper-level protocol, the packet-handling device needs to unravel the Extension Header chain, raising two rather tough questions. First, how long is the device prepared to spend unraveling this chain? And second, would the device be prepared to pass on a packet with an Extension Header that it did not recognize?

In some cases, implementers of IPv6 equipment have found it simpler to just drop all IPv6 packets that contain Extension Headers. Some measurements of this behavior are reported in RFC 7872[6]. This document reports a 38% packet-drop rate when sending fragmented IPv6 query packets to DNS Name servers. But the example provided previously is in fact the opposite case to that reported in RFC 7872, and the example illustrates a more conventional case. It's not the queries in the DNS that can readily grow to sizes that require packet fragmentation, but the responses. The relevant question concerns the anticipated probability of packet drop when sending fragmented UDP IPv6 packets as responses to DNS queries. To rephrase the question slightly, how do DNS recursive resolvers fare when the IPv6 response from the server is fragmented?

3.2 How widespread is this problem?

For a start, it appears from the example cited here that Google's Public DNS resolvers experienced some packet-drop problem when they passed a fragmented IPv6 response (this problem was noted in mid-2017, and Google has subsequently corrected it). But was this problem limited to just one or two DNS resolvers, or do many other DNS resolvers experience a similar packet-drop issue? How widespread is this problem? We tested this question using three approaches.

3.3 Repairing Missing "Glue" with Large DNS packets

The measurement technique we are using is based on scripting inside online ads[TBD]. This allows us to instrument a server and get the endpoints who are executing the measurement script to interact with the server. However, we cannot see what the endpoint is doing. For example, we can see from the server when we deliver a DNS response to a client, but we have no clear way to confirm that the client received the response. Normally the mechanisms are indirect, such as looking at whether or not the client then retrieved a web object that was the target of the DNS name. This measurement approach has some degree of uncertainty, as there are a number of reasons for a missing web object fetch, and the inability to resolve the DNS name is just one of these reasons. Is there a better way to measure how DNS resolvers behave?

The first approach we've used here is so-called "Glueless" delegation, a novel measurement mechanism by manipulating dynamic name on DNS name servers. The basic approach is to remove the additional section from the "parent" DNS response that lists the IP address of the authoritative name servers for the delegated "child" domain(Example TBD). A resolver, when provided with this answer must suspend its effort to resolve the original DNS name and instead resolve the name server name. Only when it has completed its task can it resume the original name resolution task. We can manipulate the characteristics of the DNS response from the name server name, and we can confirm if the resolver received the response by observing whether it was then able to resume the original resolution task and query the child name server.

We tested the system using an IPv6-only name server address response that used three response sizes:

Small: 169 octets Medium: 1,428 octets Large: 1,886 octets

The local MTU on the server was set to 1,280 octets, so both the medium and large responses were fragmented.

This test was loaded into an online advertising campaign.

Results - I

68,229,946 experiments 35,602,243 experiments used IPv6-capable resolvers

Small: $34,901,983 / 35,602,243 = 98.03\% = 1.97\%$
Drop Medium: $34,715,933 / 35,666,726 = 97.335\% = 2.67\%$
Drop Large: $34,514,927 / 35,787,975 = 96.44\% = 3.56\%$ Drop

The first outcome from this data is somewhat surprising. While the overall penetration of IPv6 in the larger Internet is currently some 15% of users, the DNS resolver result is far higher. Some 52% of these 68M experiments directed their DNS queries to recursive resolvers that were capable of posing their DNS queries over a IPv6 network substrate.

That's an interesting result:

Some 52% of tested endpoints use DNS resolvers that are capable of using IPv6.

Interpreting the packet drop probabilities for the three sizes of DNS responses is not so straight forward. There is certainly an increased probability of drop for the larger DNS responses, but this is far lower than the 40% drop rate reported in RFC 7872.

It seems that we should question the experimental conditions that we used here. Are these responses actually using fragmentation in IPv6?

We observed that a number of recursive resolvers use different query options when resolving the addresses of name servers, as distinct from resolving names. In particular, a number of resolvers, including Google's public DNS resolvers, strip all EDNS(0) query options from these name server address resolution queries.

When the query has no EDNS(0) options, and in particular when there is no UDP Buffer size option in the query, then the name server responds with what will fit in 512 octets. If the response is larger, and in our case this includes the Medium and Large tests, the name server sets the Truncated Response flag in its response to indicate that the provided response is incomplete. This Truncated Response flag is a signal to the resolver that it should query again, using TCP this time.

In the case of this experiment we saw some 18,571,561 medium-size records resolved using TCP and 19,363,818 large-size records resolved using TCP. This means that the observed rate of failure to resolve the name is not necessarily attributable to an inability to receive fragmented IPv6 UDP packets.

Perhaps if we remove all those instances that use TCP to retrieve the large DNS response, then what do we have left?

UDP-only queries:

Small: $34,901,983 / 35,602,243 = 98.03\% = 1.97\%$
 Drop Medium: $16,238,433 / 17,095,165 = 92.21\% = 5.01\%$
 Drop Large: $15,257,853 / 16,424,157 = 93.90\% = 7.10\%$ Drop

There is certainly a clearer signal here in this data - some 5% to 7% of experiments used DNS resolvers that appeared to be incapable of retrieving a fragmented IPv6 UDP DNS response, as compared to the “base” loss rate as experienced by the control small response of 2%. Tentatively, we can propose that a minimum of 3% of clients use DNS resolvers that fail to receive fragmented IPv6 packets.

However, in doing this we have filtered out more than one half of the tests, and perhaps we have filtered out those resolvers that cannot receive IPv6 fragmented packets.

3.4 Large DNS Packets and Web Fetch

Our second approach was to use a large response for the ‘final’ response for the requested name.

The way in which this has been done is to pad the response using bogus DNSSEC signature records (RRSIG). These DNSSEC signature records are bogus in the sense that the name itself is not DNSSEC-signed, so the content of the digital signature will never be checked, but as long as the resolver is using EDNS(0) and has turned on the DNSSEC OK bit, which occurs in some 70% of all DNS queries to authoritative name servers, then the DNSSEC signature records will be added to the response.

We are now looking at the web fetch rate, and looking for a variance between the web fetch rates when the DNS responses involve UDP IPv6 fragmentation. We filtered out all experiments that did not fetch the small DNS web object, all experiments that did not set the DO bit in their query, and all experiments that used TCP for the medium and large experiments. In this case, we are looking for those experiments where a fragmented UDP IPv6 response was passed and testing whether or not the endpoint retrieved the web object.

Results - II

68,229,946 experiments 25,096,961 experiments used UDP IPv6-capable resolvers and had the DO bit set in the query

Medium: $13,648,884 / 25,096,961 = 54.38\% = 45.62\%$
 Drop Large: $13,476,800 / 24,969,527 = 53.97\% = 46.03\%$ Drop

This is a result that is more consistent with the drop rate reported in RFC 7872, but there are a number of factors at play here, and it is not clear exactly how much of this drop rate can be directly attributed to the issue of packet fragmentation in IPv6 and the network’s handling

of IPv6 packets with Extension Headers. Again, there is also the consideration that in only looking at a subset of resolvers, namely those resolvers that use IPv6, use EDNS(0) options and set the DO bit in these queries.

3.5 Fragmented Small DNS Packets

Let’s return to the first experiment, as this form of experiment has far less potential sources of noise in the measurement. We are wanting to test whether a fragmented IPv6 packet can be received by recursive DNS resolvers, and our use of a large fragmented response is being frustrated by DNS truncation.

What if we use a customized DNS name server arrangement that gratuitously fragments the small DNS response itself? While the IPv6 specification specifies that network Path MTU sizes should be no smaller than 1,280 octets, it does not specify a minimum size of fragmented IPv6 packets.

The approach we’ve taken in this experiment is to use a user level packet processing system that listens on UDP port 53 and passes all incoming DNS queries to a back-end DNS server. When it receives a response from this back-end server it generates a sequence of IPv6 packets that fragments the DNS payload and uses a raw device socket to pass these packets directly to the device interface.

We are relying on the observation that IPv6 packet fragmentation occurs at the IP level in the protocol stack, so the IPv6 driver at the remote end will reassemble the fragments and pass the UDP payload to the DNS application, and if the payload packets are received by the resolver, there will be no trace that the IPv6 packets were fragmented.

As we are manipulating the response to the query for the address of the name server, we can tell if the recursive resolver has received the fragmented packets if the resolver resumes its original query sequence and queries for the terminal name.

Results - III

10,851,323 experiments used IPv6 queries for the name server address 6,786,967 experiments queried for the terminal DNS name

Fragmented Response: $6,786,967 / 10,851,323 = 62.54\% = 37.45\%$ Drop

This is our second result:

Some 37% of endpoints used IPv6-capable DNS resolvers that were incapable of receiving a fragmented IPv6 response.

We used three servers for this experiment, on serving Asia Pacific, a second serving the America and the

third serving Eurasia and Africa. There are some visible differences in the drop rate:

Asia Pacific: 31% Drop Americas: 37% Drop Eurasia & Africa: 47% Drop

Given that this experiment occurs completely in the DNS, we can track each individual DNS resolver as they query for the name server record then, depending on if they receive the fragmented response, query for the terminal name. There are approximately 2 million recursive resolvers in today's Internet, but only some 15,000 individual resolvers appear to serve some 90% of all users. This implies that the behavior of the most intensively used resolvers has a noticeable impact on the overall picture of capabilities of DNS infrastructure for the Internet.

We saw 10,115 individual IPv6 addresses used by IPv6-capable recursive resolvers. Of this set of resolvers, we saw 3,592 resolvers that consistently behaved in a manner that was consistent with being unable to receive a fragmented IPv6 packet. The most intensively used recursive resolvers which exhibit this problem are shown in the following table.

Resolver Hits	AS	AS Name	CC	2405:200:1606:672::5
4,178,119	55836	RELIANCEJIO-IN	Reliance Jio	IN
2402:8100:c::8	1,352,024	55644	IDEANET1-IN	Idea Cellular Limited
IN 2402:8100:c::7	1,238,764	55644	IDEANET1-IN	Idea Cellular Limited
IN 2407:0:0:2b::5	938,584	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:2a::3	936,883	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:2a::6	885,322	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:2b::6	882,687	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:2b::2	882,305	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:2a::4	881,604	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:2a::5	880,870	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:2a::2	877,329	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:2b::4	876,723	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:2b::3	876,150	4761	INDOSAT-INP-AP	INDOSAT
ID 2402:8100:d::8	616,037	55644	IDEANET1-IN	Idea Cellular Limited
IN 2402:8100:d::7	426,648	55644	IDEANET1-IN	Idea Cellular Limited
IN 2407:0:0:9::2	417,184	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:8::2	415,375	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:8::4	414,410	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:9::4	414,226	4761	INDOSAT-INP-AP	INDOSAT
ID 2407:0:0:9::6	411,993	4761	INDOSAT-INP-AP	INDOSAT

This table is slightly misleading in so far as very large recursive resolvers use resolver "farms" and the queries are managed by a collection of query 'slaves'. We can group these individual resolver IPv6 addresses to their common Origin AS, and look at which networks use resolvers that show this problem with IPv6 Extension Header drops.

The second table (below) now shows the preeminent position of Google's Public DNS service as the most heavily used recursive resolver, and its Extension Header drop issues, as shown in the example at the start of this article, is consistent with its position at the head of the list of networks that have DNS resolvers with this problem.

AS Hits	% of Total AS	Name,CC
15169	7,952,272	17.3% GOOGLE - Google Inc., US 4761 6,521,674 14.2%
55644	4,313,225	9.4% INDOSAT-INP-AP INDOSAT, ID 55644 4,313,225 9.4%
22394	4,217,285	9.2% IDEANET1-IN Idea Cellular Limited, IN 22394 4,217,285
55836	4,179,921	9.1% CELLCO - Cellco Partnership DBA Verizon Wireless, US 55836 4,179,921 9.1%
10507	2,939,364	6.4% RELIANCEJIO-IN Reliance Jio Infocomm Limited, IN 10507 2,939,364 6.4%
5650	2,005,583	4.4% SPCS - Sprint Personal Communications Systems, US 5650 2,005,583 4.4%
2516	1,322,228	2.9% FRONTIER-FRTR - Frontier Communications, US 2516 1,322,228 2.9%
6128	1,275,278	2.8% KDDI KDDI CORPORATION, JP 6128 1,275,278 2.8%
32934	1,128,751	2.5% CABLE-NET-1 - Cablevision Systems Corp., US 32934 1,128,751 2.5%
20115	984,165	2.1% FACEBOOK - Facebook, US 20115 984,165 2.1%
9498	779,603	1.7% CHARTER-NET-HKY-NC - Charter Communications, US 9498 779,603 1.7%
20057	438,137	1.0% BBIL-AP BHARTI Airtel Ltd., IN 20057 438,137 1.0%
17813	398,404	0.9% ATT-MOBILITY-LLC-AS20057 - AT&T Mobility LLC, US 17813 398,404 0.9%
2527	397,832	0.9% MTNL-AP Mahanagar Telephone Nigam Ltd., IN 2527 397,832 0.9%
45458	276,963	0.6% SO-NET So-net Entertainment Corporation, JP 45458 276,963 0.6%
6167	263,583	0.6% SBN-AWN-AS-02-AP SBN-ISP/AWN-ISP, TH 6167 263,583 0.6%
8708	255,958	0.6% Cellco Partnership DBA Verizon Wireless, US 8708 255,958 0.6%
38091	255,930	0.6% RCS-RDS 73-75 Dr. Staicovici, RO 38091 255,930 0.6%
18101	168,164	0.4% HELLONET-AS-KR CJ-HELLOVISION, KR 18101 168,164 0.4%
		Reliance Communications DAKC MUMBAI, IN

However, one conclusion looks starkly clear to me from these results. We can't just assume that the DNS as we know it today will just work in an all IPv6 future Internet. We must make some changes in some parts of the protocol design to get around this current widespread problem of IPv6 Extension Header packet loss in the DNS, assuming that we want to have a DNS at all in this all-IPv6 future Internet.

4 ATR MECHANISM

4.1 Methodology of ATR

As we stated in related work, we could move the DNS away from UDP and use TCP instead. That move would certainly make a number of functions a lot easier, including encrypting DNS traffic on the wire as a means of assisting with aspects of personal privacy online as well as accommodating large DNS responses.

However, the downside is that TCP imposes a far greater load overhead on servers, and while it is possible

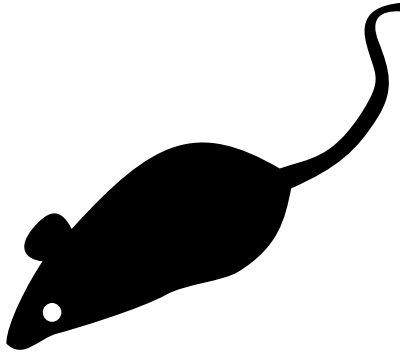


Figure 1: ATR behavior

to conceive of an all-TCP DNS environment, it is more challenging to understand the economics of such a move and to understand, in particular, how name publishers and name consumers will share the costs of a more expensive name resolution environment.

If we want to continue to UDP where it's feasible, and continue to use TCP only as the 'Plan B' protocol for name resolution, then can we improve the handling of large response in UDP? Specifically, can we make this hybrid approach — of using UDP when we can, and TCP only when we must — faster and more robust?

An approach to address this challenge is first introduced in IETF draft draft-song-atr-large-resp-00. The approach described in this draft is simple: If a DNS server provides a response that entails sending fragmented UDP packets, then the server should wait for a 10 ms period and also back the original query as a truncated response. If the client receives and reassembles the fragmented UDP response, then the ensuing truncated response will be ignored by the client's DNS resolver as its outstanding query has already been answered. If the fragmented UDP response is dropped by the network, then the truncated response will be received (as it is far smaller), and reception of this truncated response will trigger the client to switch immediately to re-query using TCP. This behavior is illustrated in Figure 1.

4.2 ATR process or Algorithm

As shown in Figure 2 the ATR module can be implemented right after truncation loop if the packet is not going to be fragmented.

The ATR responding process goes as follows:

- o When an authoritative server receives a query and enters the responding process, it first goes through the normal truncation loop to see whether the size of response surpasses the EDNS0 payload size. If yes, it ends

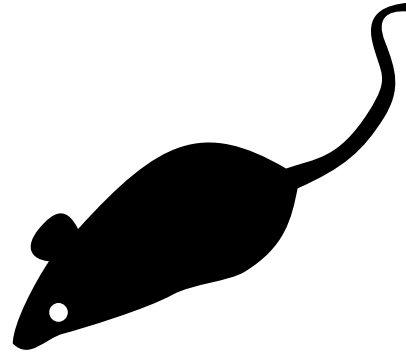


Figure 2: ATR module in DNS response process

up with responding a truncated packet. If no, it enters the ATR module.

- o In ATR module, similar like truncation loop, the size of response is compared with a value called ATR payload size. If the response of a query is larger than ATR payload size, the server firstly sends the normal response and then coins a truncated response with the same ID of the query.

- o The server can reply the coined truncated response in no time. But considering the possible impact of network reordering, it is suggested a timer to delay the second truncated response, for example 10 50 millisecond which can be configured by local operation.

Note that the choice of ATR payload size and timer SHOULD be configured locally. And the operational consideration and guidance is discussed in Section 4.2 and Section 4.1 respectively.

There are three typical cases of ATR-unaware resolver behavior when a resolver sends query to an ATR server in which the server will generate a large response with fragments:

- o Case 1: a resolver (or sub-resolver) will receive both the large response and a very small truncated response in sequence. It will happily accept the first response and drop the second one because the transaction is over.

- o Case 2: In case a fragment is dropped in the middle, the resolver will end up with only receiving the small truncated response. It will retry using TCP in no time.

- o Case 3: For those (probably 30%*17% of them) who can not speak TCP and sitting behind a firewall stubbornly dropping fragments. Just say good luck to them!

In the case authoritative server truncated all response surpasses certain value, for example setting IPv6-edns-size to 1220 octets, ATR will be helpful for resolver with TCP capacity, because the resolver still has a fair chance to receive the large response.

4.3 ATR timer

TBD

4.4 ATR payload size

Regarding the operational choice for ATR payload size, there are some good input from APNIC study [scoring-dns-root] on how to react to large DNS payload for authoritative server. The difference in ATR is that ATR focuses on the second response after the ordinary response.

For IPv4 DNS server, it is suggested the study that do not truncate and fragment IPv4 UDP response with a payload up to 1472 octets which is Ethernet MTU(1500) minus the sum of IPv4 header(20) and UDP header(8). The reason is to avoid gratuitously fragmenting outbound packets and TCP fallback at the source.

In the case of ATR, the first ordinary response is emitted without knowing it be to fragmented or not on the path. If a large value is set up to 1472 octets, payload size between 512 octets and the large value size will probably get fragmented by aggressive firewalls which leads losing the benefit of ATR. If ATR payload size set exactly 512 octets, in most of case ATR response and the single unfragmented packets are under a race at the risk of RO.

Given IPv4 fragmentation issue is not so serious compared to IPv6, it is suggested in this memo to set ATR payload size 1472 octets which means ATR only fit large DNS response larger than 1500 octets in IPv4.

For IPv6 DNS server, similar to IPv4, the APNIC study is suggested that do not truncate IPv6 UDP packets with a payload up to 1,452 octets which is Ethernet MTU(1500) minus the sum of IPv6 header(40) and UDP header(8). 1452 octets is chosen to avoid TCP fallback in the context that most TCP MSS in the root server is not set probably at that time.

In the case of ATR considering the second truncated response, a smaller size: 1232 octets, which is IPv6 MTU for most network devices(1280) minus the sum of IPv6 header(40) and UDP header(8), should be chosen as ATR payload size to trigger necessary TCP fallback. As a complementary requirement with ATR, the TCP MSS should be set 1220 octets to avoid Packet Too Big ICMP message as suggested in the APNIC study.

In short, it is recommended that in IPv4 ATR payload size SHOULD be 1472 octets, and in IPv6 the value SHOULD be 1232 octets.

4.5 Less aggressiveness of ATR

There is a concern ATR sends TC=1 response too aggressively especially in the beginning of adoption. ATR

can be implemented as an optional and configurable feature at the disposal of authoritative server operator. One of the idea to mitigate this aggressiveness, ATR may respond TC=1 responses at a low possibility, such as 10

Another way is to reply ATR response selectively. It is observed that RO and IPv6 fragmentation issues are path specific and persistent due to the Internet components and middle box. So it is reasonable to keep a ATR "Whitelist" by counting the retries and recording the IP destination address of that large response causing many retires. ATR only acts to those queries from the IP address in the white list.

5 ATR EVALUATION

It is worth of mentioning APNIC report[How-ATR-Work] on "How well does ATR actually work?" done by Geoff Huston and Joao Damas after 00 version of ATR draft. It was reported firstly in IEPG meeting before IETF 101 and then posted in APNIC Blog later.

It is said the test was performed over 55 million endpoints, using an on-line ad distribution network to deliver the test script across the Internet. The result is positive that ATR works! From the end users' perspective, in some 9% of IPv4 cases the use of ATR by the server will improve the speed of resolution of a fragmented UDP response by signaling to the client an immediate switch to TCP to perform a re-query. The IPv6 behavior would improve the resolution times in 15% of cases.

It also analyzed the pros and cons of ATR. On one hand, It is said that ATR certainly looks attractive if the objective is to improve the speed of DNS resolution when passing large DNS responses. And ATR is incrementally deployable in favor of decision made by each server operator. On another hand, ATR also has some negative factors. One issue is adding another DNS DDoS attack vector due to the additional packet sent by ATR, (author's note : very small adding actually.) Another issue is risk of RO by the choice of the delay timer which is discussed fully in Section 4.1.

As a conclusion, it is said that "ATR does not completely fix the large response issue. If a resolver cannot receive fragmented UDP responses and cannot use TCP to perform DNS queries, then ATR is not going to help. But where there are issues with IP fragment filtering, ATR can make the inevitable shift of the query to TCP a lot faster than it is today. But it does so at a cost of additional packets and additional DNS functionality". "If a faster DNS service is your highest priority, then ATR is worth considering", said at the end of this report

6 NETWORK RE-ORDERING CONSIDERATION

As introduced in Section 2 ATR timer is a way to avoid the impact of network reordering(RO). The value of the timer is critical, because if the delay is too short, the ATR response may be received earlier than the fragmented response (the first piece), the resolver will fall back to TCP bearing the cost which should have been avoided. If the delay is too long, the client may timeout and retry which negates the incremental benefit of ATR. Generally speaking, the delay of the timer should be "long enough, but not too long".

To the best knowledge of author, the nature of RO is characterized as follows hopefully helping ATR users understand RO and how to operate ATR appropriately in RO context.

- o RO is mainly caused by the parallelism in Internet components and links other than network anomaly [Bennett]. It was observed that RO is highly related to the traffic load of Internet components. So RO will long exists as long as the traffic load continue increase and the parallelism is used to enhance network throughput.

- o The probability of RO varies largely depending on the different tests samples. Some work shown RO

probability below 2[Tinta] and another work was above 90% agreed that RO is site-dependent and path-dependent. It is observed in that when RO happens, it is mostly exhibited consistently in a small percentages of the paths. It is also observed that higher rates smaller packets were more prone to RO because the sending inter-spacing time was small.

- o It was reported that the inter-arrival time of RO varies from a few milliseconds to multiple tens of milliseconds [Tinta]. And the larger the packet the larger the inter-arrival time, since larger packets will take longer to be transmitted.

Reasonably we can infer that firstly RO should be taken into account because it long exists due to middle Internet components which can not be avoided by end-to-end way. Secondly the mixture of larger and small packets in ATR case will increase the inter-arrival time of RO as well as the its probability. The good news is that the RO is highly site specific and path specific, and persistent which means

the ATR operator is able to identify a few sites and paths, setup a tunable timer setting for them, or just put them into a blacklist without replying ATR response.

Based on the above analysis it i