

Internet Engineering Task Force  
Internet-Draft  
Intended status: Informational  
Expires: September 9, 2019

L. Song  
Beijing Internet Institute  
S. Wang  
Beijing Normal University  
March 8, 2019

ATR: Additional Truncation Response for Large DNS Response  
draft-song-atr-large-resp-03

Abstract

As the increasing use of DNSSEC and IPv6, there are more public evidence and concerns on IPv6 fragmentation issues due to larger DNS payloads over IPv6. This memo introduces an simple improvement on DNS server by replying an additional truncated response just after the normal fragmented response. It can be used to relieve users suffering on DNS latency and failures due to large DNS response. An ATR Experiment was done to show how well it works and some operational issues are discussed in this memo as well.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <https://datatracker.ietf.org/drafts/current/>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on September 9, 2019.

Copyright Notice

Copyright (c) 2019 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect

to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

## Table of Contents

1. Introduction . . . . .	2
2. The ATR mechanism . . . . .	3
3. Experiment on how well ATR works . . . . .	5
4. Operational considerations . . . . .	6
4.1. ATR timer . . . . .	6
4.2. ATR payload size . . . . .	7
4.3. Less aggressiveness of ATR . . . . .	8
5. Security Considerations . . . . .	8
6. IANA considerations . . . . .	8
7. Acknowledgments . . . . .	9
8. References . . . . .	9
Appendix A. Considerations on Resolver awareness of ATR . . . . .	11
Appendix B. Revision history of this document . . . . .	11
B.1. draft-song-atr-large-resp-01 . . . . .	11
B.2. draft-song-atr-large-resp-02 . . . . .	12
B.3. draft-song-atr-large-resp-03 . . . . .	12
Authors' Addresses . . . . .	13

## 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 [DNSSEC-impact] 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. Depending on retry model, the resolver's failing to receive fragmented response may experience long latency or failure due to timeout and retries. And, most of the underlying issues with

fragments are unrevealed due to good redundancy and resilience of DNS and dual-stack network.

Generally speaking there are two approaches for this issue. One is to make the DNS response as small as possible, for example, using ECC instead of RSA to shorten the size of Key and signature. However, few zones are signed by ECC for the time being. In addition there is an uncertainty in the algorithm rollover from RSA to ECC. Another approach is to fall back to TCP by setting on either server side or client side. For resolver it is to set EDNS0 bufsize below a certain number. For authoritative servers it is to set their maximum UDP response size small enough.

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 case.

To relieve the problem, this memo introduces a small improvement on DNS responding process by replying an Additional Truncated Response (ATR) just after a normal large response which is to be fragmented. It is a hybrid approach of using UDP when we can, and TCP only when we must. 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.

[REMOVE BEFORE PUBLICATION] Note that ATR is not just a proposed idea. Some advocates of ATR implemented it based on BIND9 ([https://gitlab.isc.org/isc-projects/bind9/merge\\_requests/158](https://gitlab.isc.org/isc-projects/bind9/merge_requests/158)). And some verify it based on a large-scale experiment platform of APNIC lab Section 3 which is introduced in this memo.

## 2. The ATR mechanism

The ATR mechanism is very simple that it involves an ATR module in the responding process of current DNS implementation. As shown in the following diagram the ATR module is right after the truncation loop if the packet is not going to be fragmented.

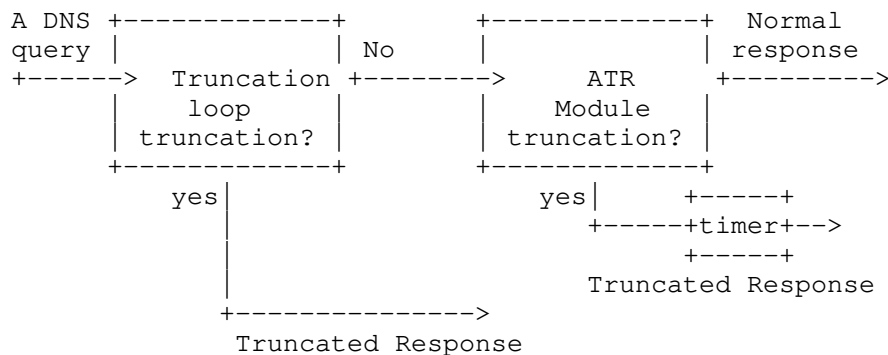


Figure 1: High-Level Testbed Components

The ATR responding process goes as follows:

- o When an authoritative server receives a query and enters the responding process, it first go through the normal truncation loop to see whether the size of response surpasses the EDNS0 payload size. If yes, it ends up with responding a truncated packets. 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 coin 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 send 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 accepts 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 surpass 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.

### 3. Experiment on how well ATR works

It is worth 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 or unanswered factors. One 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. It is also founded that the trailing UDP packet may generate ICMP Port Unreachable messages back to the server as a kind of noise (a rate of approximately 1 in 5 responses in our experiments). Note that in author's argument, it is not a big issue and the server can simply ignore it if it decides to adopt ATR.

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

#### 4. Operational considerations

There are some operational consideration on ATR, such as the parameter of the ATR timer and ATR payload size, and policies on when ATR is triggered to avoid side-effect.

##### 4.1. ATR timer

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% [Paxson] [Tinta] and another work was above 90% [Bennett]. But it is 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 is hard to provide a perfect value of ATR timer for all ATR users due to the diversity of networks. It seems OK to set the timer with a range from ten to hundreds ms, just below the timeout setting of typical resolver. It is suggested that a decision should be made as operator-specific according to the statistic of the RTT of their users. Some measurement shown [Brownlee][Liang] the mean of response time is below 50 ms for the sites with lots of anycast instance like L-root, .com and .net name servers. For that sites, delay less than 50 ms is appropriate.

#### 4.2. 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.

#### 4.3. 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 retries. ATR only acts to those queries from the IP address in the white list.

#### 5. Security Considerations

There may be concerns on DDoS attack problem due to the fact that the ATR introduces multiple responses from authoritative server. The extra packet is pretty small. In the worst case, it's 50% more packets and they are small

DNS cookies [RFC7873] and RRL on authoritative may be possible solutions

#### 6. IANA considerations

No IANA considerations for this memo



## 7. Acknowledgments

Many thanks to reviewers and their comments. Geoff Huston and Joao Damas did a testing on the question "How well does ATR actually work?". Alexander Dupuy proposed the idea to distinguish ATR responses from normal ones. Akira Kato contributed ideas on operational consideration. Shane Kerr help author with the security consideration. Stephane Bortzmeyer gave thought of happyeyeballs on resolver side.

Acknowledgments are also give to Mukund Sivaraman, Evan Hunt and Mark Andrews who implement it and maintained it in a brunch in BIND9 code base.

## 8. References

### [ATR-Github]

"XML source file and test script of DNS ATR", September 2017, <[https://github.com/songlinjian/DNS\\_ATR](https://github.com/songlinjian/DNS_ATR)>.

[Bennett] Bennett, J. C. R., "Packet Reordering is Not Pathological Network Behavior", December 1999, <<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.461.7629&rep=rep1&type=pdf>>.

### [Brownlee]

Brownlee, N., "Response time distributions for global name servers", 2002, <<http://www.caida.org/publications/papers/2002/nsrtd/nsrtd.pdf>>.

### [DNSSEC-impact]

Broek, G. V. D., "DNSSEC meets real world: dealing with unreachability caused by fragmentation", April 2014, <<https://repository.ubn.ru.nl/bitstream/handle/2066/132796/132796.pdf?sequence=1>>.

### [How-ATR-Work]

Huston, G., "How well does ATR actually work?", April 2018, <<https://blog.apnic.net/2018/04/16/how-well-does-atr-actually-work/>>.

### [I-D.taylor-v6ops-fragdrop]

Jaeggli, J., Colitti, L., Kumari, W., Vyncke, E., Kaeo, M., and T. Taylor, "Why Operators Filter Fragments and What It Implies", draft-taylor-v6ops-fragdrop-02 (work in progress), December 2013.

- [IPv6-frag-DNS] Huston, G., "Dealing with IPv6 fragmentation in the DNS", August 2017, <<https://blog.apnic.net/2017/08/22/dealing-ipv6-fragmentation-dns>>.
- [Liang] Liang, J., "Measuring Query Latency of Top Level DNS Servers", February 2013, <<https://netsec.ccert.edu.cn/duanhx/files/2013/02/latency.pdf>>.
- [Not-speak-TCP] Huston, G., "A Question of DNS Protocols", August 2013, <<https://labs.ripe.net/Members/gih/a-question-of-dns-protocols>>.
- [Paxson] Paxson, V., "End-to-End Internet Packet Dynamics", August 1999, <<https://cseweb.ucsd.edu/classes/fa01/cse222/papers/paxson-e2e-packets-sigcomm97.pdf>>.
- [RFC1035] Mockapetris, P., "Domain names - implementation and specification", STD 13, RFC 1035, DOI 10.17487/RFC1035, November 1987, <<https://www.rfc-editor.org/info/rfc1035>>.
- [RFC6891] Damas, J., Graff, M., and P. Vixie, "Extension Mechanisms for DNS (EDNS(0))", STD 75, RFC 6891, DOI 10.17487/RFC6891, April 2013, <<https://www.rfc-editor.org/info/rfc6891>>.
- [RFC7872] Gont, F., Linkova, J., Chown, T., and W. Liu, "Observations on the Dropping of Packets with IPv6 Extension Headers in the Real World", RFC 7872, DOI 10.17487/RFC7872, June 2016, <<https://www.rfc-editor.org/info/rfc7872>>.
- [RFC7873] Eastlake 3rd, D. and M. Andrews, "Domain Name System (DNS) Cookies", RFC 7873, DOI 10.17487/RFC7873, May 2016, <<https://www.rfc-editor.org/info/rfc7873>>.
- [scoring-dns-root] Huston, G., "Scoring the DNS Root Server System", November 2016, <<https://blog.apnic.net/2016/11/15/scoring-dns-root-server-system/>>.
- [Tinta] Tinta, S. P., "Characterizing End-to-End Packet Reordering with UDP Traffic", August 2009, <<https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/35247.pdf>>.

## Appendix A. Considerations on Resolver awareness of ATR

ATR proposed in this memo is a server-side function which requires no change in resolver, so it is not required that resolver MUST recognize ATR and react accordingly. But it may be helpful for some cases where a resolver is able to recognize ATR response, for example by checking the large edns0 payload size and Truncation bit.

One case where ATR is used as a troubleshooting tool by which resolver operators are able to flag problematic name servers. The resolver operator is enabled to log cases where ATR responses are received without a (reassembled) UDP response to a query. In the case of receiving a ATR, RDNS can choose to restrict maximum EDNS to a lower value than the default 4096 that is currently used.

Another case is that when receiving a ATR response a ATR-aware resolver can adopt a "happy eyeballs" strategy by opening a separate transaction sending the query via TCP instead of falling back to TCP and closing the original UDP transaction. Listening to port 53 on both TCP and UDP port 53 will enhance the availability and reduce the latency. It will add more tolerance to network reordering issues as well. However, it should be taken into account about the balance of resolver's resource. Less priority should be given to that function when the resolver is "busy".

The awareness of ATR on resolver can also avoid sending ICMP Port Unreachable messages back to the server. In some implementations, reusing the same UDP sockets for multiple queries will not generate that ICMP noise.

However, resolver use case of ATR is currently outside of the scope of server-ATR proposal. It needs further discussion.

## Appendix B. Revision history of this document

### B.1. draft-song-atr-large-resp-01

After receiving reviews and comments, changes of 01 version are shown as follows:

- o Rewrite introduction and add another goal of ATR as a measuring tool;
- o Add section 3 indicating a ATR response. A bit in the EDNS0 OPT header is defined as an indicator of ATR response. The flag bit is called "ATR Response" (AT) bit;

- o Add Section 4 Operation considerations, which discuss ATR timer , ATR payload size, and less aggressiveness of ATR;
- o Add IANA consideration to register the AT bit;
- o Add section 7 Acknowledgments;
- o Append a list of references regarding Network reordering, and APNIC's study on IPv6 and DNS;
- o Add Appendix A, An introduce of APNIC testing work and author's comments;
- o Appendix B. Considerations on Resolver awareness of ATR;
- o Change the category="std" . It is said in RFC6891 IETF Standards Action is required for assignments of new EDNS(0) flags. So the draft should be categorized as standard track if registering AT bit is desired in this document.

Change history is also available in the public GitHub repository where this document is maintained: <[https://github.com/songlinjian/DNS\\_ATR](https://github.com/songlinjian/DNS_ATR)>.

#### B.2. draft-song-atr-large-resp-02

Changes in 02 version of ATR draft:

- o Remove the section of introduction of AT bit as well as requirement of IANA registration of that bit;
- o Change the category of this document to experimental and move the introduction of APNIC's experiment from Appendix A to section 3;
- o Add more names in Acknowledgments part after IETF102;

#### B.3. draft-song-atr-large-resp-03

Changes in 03 version of ATR draft:

- o Add related work in the introduction session;
- o Introduce ICMP noise as a finding of APNIC's experiment and propose how to avoid it in Appendix A;
- o Change the category from "exp" to "info";
- o Move to ISE for review.

- o Add one author S. Wang

Authors' Addresses

Linjian Song  
Beijing Internet Institute  
2nd Floor, Building 5, No.58 Jing Hai Wu Lu, BDA  
Beijing 100176  
P. R. China

Email: [songlinjian@gmail.com](mailto:songlinjian@gmail.com)  
URI: <http://www.biigroup.com/>

Shengling Wang  
Beijing Normal University  
Beijing Normal University, No. 19, XinJieKouWai St., HaiDian District  
Beijing 100875  
P. R. China

Email: [wangshengling@bnu.edu.cn](mailto:wangshengling@bnu.edu.cn)  
URI: <https://cist.bnu.edu.cn/>