

DNSOP Working Group  
Internet-Draft  
Updates: [1034](#), [1035](#) (if approved)  
Intended status: Standards Track  
Expires: April 17, 2019

D. Lawrence  
Oracle + Dyn  
W. Kumari  
P. Sood  
Google  
October 14, 2018

## **Serving Stale Data to Improve DNS Resiliency draft-ietf-dnsop-serve-stale-02**

### Abstract

This draft defines a method for recursive resolvers to use stale DNS data to avoid outages when authoritative nameservers cannot be reached to refresh expired data. It updates the definition of TTL from [\[RFC1034\]](#) and [\[RFC1035\]](#) to make it clear that data can be kept in the cache beyond the TTL expiry and used for responses when a refreshed answer is not readily available.

### Ed note

Text inside square brackets ([ ]) is additional background information, answers to frequently asked questions, general musings, etc. They will be removed before publication. This document is being collaborated on in GitHub at <https://github.com/vttale/serve-stale>. The most recent version of the document, open issues, etc should all be available here. The authors gratefully accept pull requests.

### 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 April 17, 2019.

## Copyright Notice

Copyright (c) 2018 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

<a href="#">1.</a>	<a href="#">Introduction</a>	<a href="#">2</a>
<a href="#">2.</a>	<a href="#">Notes to readers</a>	<a href="#">3</a>
<a href="#">3.</a>	<a href="#">Terminology</a>	<a href="#">3</a>
<a href="#">4.</a>	<a href="#">Background</a>	<a href="#">4</a>
<a href="#">5.</a>	<a href="#">Standards Action</a>	<a href="#">4</a>
<a href="#">6.</a>	<a href="#">EDNS Option</a>	<a href="#">5</a>
<a href="#">6.1.</a>	<a href="#">Option Format Proposal 1</a>	<a href="#">5</a>
<a href="#">6.2.</a>	<a href="#">Option Usage</a>	<a href="#">6</a>
<a href="#">6.3.</a>	<a href="#">Option Format Proposal 2</a>	<a href="#">7</a>
<a href="#">7.</a>	<a href="#">Example Method</a>	<a href="#">7</a>
<a href="#">8.</a>	<a href="#">Implementation Caveats</a>	<a href="#">9</a>
<a href="#">9.</a>	<a href="#">Implementation Status</a>	<a href="#">10</a>
<a href="#">10.</a>	<a href="#">Security Considerations</a>	<a href="#">11</a>
<a href="#">11.</a>	<a href="#">Privacy Considerations</a>	<a href="#">11</a>
<a href="#">12.</a>	<a href="#">NAT Considerations</a>	<a href="#">11</a>
<a href="#">13.</a>	<a href="#">IANA Considerations</a>	<a href="#">11</a>
<a href="#">14.</a>	<a href="#">Acknowledgements</a>	<a href="#">11</a>
<a href="#">15.</a>	<a href="#">References</a>	<a href="#">11</a>
<a href="#">15.1.</a>	<a href="#">Normative References</a>	<a href="#">11</a>
<a href="#">15.2.</a>	<a href="#">Informative References</a>	<a href="#">12</a>
	<a href="#">Authors' Addresses</a>	<a href="#">12</a>

## [1. Introduction](#)

Traditionally the Time To Live (TTL) of a DNS resource record has been understood to represent the maximum number of seconds that a record can be used before it must be discarded, based on its description and usage in [\[RFC1035\]](#) and clarifications in [\[RFC2181\]](#).

This document proposes that the definition of the TTL be explicitly expanded to allow for expired data to be used in the exceptional



circumstance that a recursive resolver is unable to refresh the information. It is predicated on the observation that authoritative server unavailability can cause outages even when the underlying data those servers would return is typically unchanged.

We describe a method below for this use of stale data, balancing the competing needs of resiliency and freshness. While this is intended to be immediately useful to the installed base of DNS software, we also propose an [[RFC6891](#)] EDNS option for enhanced signalling around the use of stale data by implementations that understand it.

## 2. Notes to readers

[ RFC Editor, please remove this section before publication!  
Readers: This is conversational text to describe what we've done, and will be removed, please don't bother sending editorial nits. :-) ]

Due to circumstances, the authors of this document got sidetracked, and we lost focus. We are now reviving it, and are trying to address and incorporate comments. There has also been more deployment and implementation recently, and so this document is now more describing what is known to work instead of simply proposing a concept.

Open questions:

- o The EDNS option that we propose for debugging is relatively full-featured to identify which RRsets are stale. It could be simplified to just indicate that an answer is stale, or even removed entirely in favour of an Extended Error result that indicates staleness.
- o What TTL value to recommend be set in stale answers returned by recursive resolvers.

## 3. Terminology

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [[RFC2119](#)] when, and only when, they appear in all capitals, as shown here.

For a comprehensive treatment of DNS terms, please see [[RFC7719](#)].



## 4. Background

There are a number of reasons why an authoritative server may become unreachable, including Denial of Service (DoS) attacks, network issues, and so on. If the recursive server is unable to contact the authoritative servers for a name but still has relevant data that has aged past its TTL, that information can still be useful for generating an answer under the metaphorical assumption that, "stale bread is better than no bread."

[RFC1035] [Section 3.2.1](#) says that the TTL "specifies the time interval that the resource record may be cached before the source of the information should again be consulted", and [Section 4.1.3](#) further says the TTL, "specifies the time interval (in seconds) that the resource record may be cached before it should be discarded."

A natural English interpretation of these remarks would seem to be clear enough that records past their TTL expiration must not be used, However, [RFC1035] predates the more rigorous terminology of [RFC2119] which softened the interpretation of "may" and "should".

[RFC2181] aimed to provide "the precise definition of the Time to Live" but in [Section 8](#) was mostly concerned with the numeric range of values and the possibility that very large values should be capped. (It also has the curious suggestion that a value in the range 2147483648 to 4294967295 should be treated as zero.) It closes that section by noting, "The TTL specifies a maximum time to live, not a mandatory time to live." This is again not [RFC2119]-normative language, but does convey the natural language connotation that data becomes unusable past TTL expiry.

Several major recursive resolver operators currently use stale data for answers in some way, including Akamai (via both Nomimum and Xerocole), Knot, OpenDNS, and Unbound. Apple can also use stale data as part of the Happy Eyeballs algorithms in mDNSResponder. The collective operational experience is that it provides significant benefit with minimal downside.

## 5. Standards Action

The definition of TTL in [RFC1035] Sections [3.2.1](#) and [4.1.3](#) is amended to read:

TTL a 32 bit unsigned integer number of seconds in the range 0 - 2147483647 that specifies the time interval that the resource record MAY be cached before the source of the information MUST again be consulted. Zero values are interpreted to mean that the RR can only be used for the transaction in progress, and should



not be cached. Values with the high order bit set SHOULD be capped at no more than 2147483647. If the authority for the data is unavailable when attempting to refresh the data past the given interval, the record MAY be used as though it is unexpired.

[ Discussion point: capping values with the high order bit as being max positive, rather than 0, is a change from [RFC2181](#). Also, we could use this opportunity to recommend a much more sane maximum value like 604800 seconds, one week, instead of the literal maximum of 68 years. ]

## 6. EDNS Option

While the basic behaviour of this answer-of-last-resort can be achieved with changes only to resolvers, explicit signalling about the use of stale data can be done with an EDNS [RFC6891](#) option. This option can be included from a stub or forwarding resolver to a recursive resolver, explicitly signalling that it does not want stale answers, or for learning that stale data was in use. It is expected that this could be useful for debugging.

[ NOTE: This section will be fleshed out a bit more thoroughly if there is interest in pursuing the option. Here are two potential options that could be used, one more fully-featured to indicate which RRsets are stale and one much more simple to indicate that stale data is present. These are proposed as mutually exclusive; the final document will have one or zero such options. We're especially soliciting feedback on this from the working group. ]

### 6.1. Option Format Proposal 1

The option is structured as follows:

	+0 (MSB)	+1 (LSB)
0:		
	OPTION-CODE	
2:		
	OPTION-LENGTH	
4:		
	STALE-RRSET-INDEX 1	
6:		
8:		
	TTL-EXPIRY 1	
	: ... additional STALE-RRSET-INDEX / TTL-EXPIRY pairs ... :	





OPTION-CODE: 2 octets per [[RFC6891](#)]. For Serve-Stale the code is TBD by IANA.

OPTION-LENGTH: 2 octets per [[RFC6891](#)]. Contains the length of the payload following OPTION-LENGTH, in octets.

STALE-RRSET-INDEX: Two octets as a signed integer, indicating the first RRSet in the message which is beyond its TTL, with RRSet counting starting at 1 and spanning message sections.

TTL-EXPIRY: Four octets as an unsigned integer, representing the number of seconds that have passed since the TTL for the RRset expired.

## **6.2. Option Usage**

Software making a DNS request can signal that it understands Serve-Stale by including the option with one STALE-RRSET-INDEX initialized to any negative value and TTL-EXPIRY initialized to 0. The index is set to a negative value to detect mere reflection of the option by responders that don't really understand it.

If the request is made to a recursive resolver which used any stale RRsets in its reply, it then fills in the corresponding indices and staleness values. If no records are stale, STALE-RRSET-INDEX and TTL-EXPIRY are set to 0.

If the request is made to an authoritative nameserver, it can use the option in the reply to indicate how the resolver should treat the records in the reply if they are unable to be refreshed later. A default for all RRsets in the message is established by setting the first STALE-RRSET-INDEX to 0, with optional additional STALE-RRSET-INDEX values overriding the default. A TTL-EXPIRY value of 0 means to never serve the RRset as stale, while non-zero values represent the maximum amount of time it can be used before it MUST be evicted. [ Does anyone really want to do this? It adds more state into resolvers. Is the idea only for purists, or is there a practical application? ]

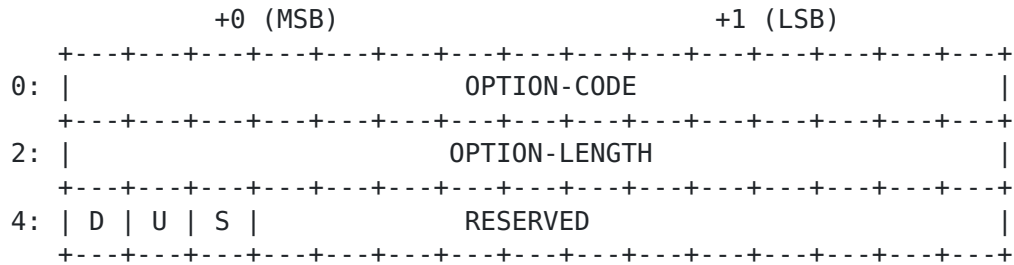
No facility is made for a client of a resolver to signal that it doesn't want stale answers, because if a client has knowledge of Serve-Stale as an option, it also has enough knowledge to just ignore any records that come back stale. [ There is admittedly the issue that the client might just want to wait out the whole attempted resolution, which there's currently no way to indicate. The absolute value of STALE-RRSET-INDEX could be taken as a timer the requester is willing to wait for an answer, but that's kind of gross overloading it like that. Shame to burn another field on that, but on the other



hand it might also be nice if a client could always signal its impatience level - "I must have an answer within 900 milliseconds!" ]

### 6.3. Option Format Proposal 2

The option is structured as follows:



OPTION-CODE: 2 octets per [\[RFC6891\]](#). For Serve-Stale the code is TBD by IANA.

OPTION-LENGTH: 2 octets per [\[RFC6891\]](#). Contains the length of the payload following OPTION-LENGTH, in octets.

D Flag: If set, the client explicitly does NOT want stale answers. If clear, the client would like an indication of whether any data in the response is stale.

U Flag: This indicates that the server understands the Serve-Stale EDNS option, and more information is communicated via the S flag. It exists to get around the issue of some authoritative servers simply echoing back EDNS options it does not understand.

S Flag: If set, this indicates that the response contains stale data. If clear, no data in the response has reached its TTL expiry.

RESERVED: Reserved for future use. Should be set to zero on send and ignored on receipt.

### 7. Example Method

There is conceivably more than one way a recursive resolver could responsibly implement this resiliency feature while still respecting the intent of the TTL as a signal for when data is to be refreshed.

In this example method three notable timers drive considerations for the use of stale data, as follows:



- o A client response timer, which is the maximum amount of time a recursive resolver should allow between the receipt of a resolution request and sending its response.
- o A query resolution timer, which caps the total amount of time a recursive resolver spends processing the query.
- o A maximum stale timer, which caps the amount of time that records will be kept past their expiration.

Recursive resolvers already have the second timer; the first and third timers are new concepts for this mechanism.

When a request is received by the recursive resolver, it SHOULD start the client response timer. This timer is used to avoid client timeouts. It SHOULD be configurable, with a recommended value of 1.8 seconds.

The resolver then checks its cache for an unexpired answer. If it finds none and the Recursion Desired flag is not set in the request, it SHOULD immediately return the response without consulting the cache for expired records.

If iterative lookups will be done, it SHOULD start the query resolution timer. This timer bounds the work done by the resolver, and is commonly around 10 to 30 seconds.

If the answer has not been completely determined by the time the client response timer has elapsed, the resolver SHOULD then check its cache to see whether there is expired data that would satisfy the request. If so, it adds that data to the response message and SHOULD set the TTL of each expired record in the message to 30 seconds. The response is then sent to the client while the resolver continues its attempt to refresh the data. 30 second was chosen because historically 0 second TTLs have been problematic for some implementations, and similarly very short TTLs could lead to congestive collapse as TTL-respecting clients rapidly try to refresh. 30 seconds not only sidesteps those potential problems with no practical negative consequence, it would also rate limit further queries from any client that is honoring the TTL, such as a forwarding resolver. [ NOTE: we're looking for further working group feedback on this value. ]

The maximum stale timer is used for cache management and is independent of the query resolution process. This timer is conceptually different from the maximum cache TTL that exists in many resolvers, the latter being a clamp on the value of TTLs as received from authoritative servers. The maximum stale timer SHOULD be



configurable, and defines the length of time after a record expires that it SHOULD be retained in the cache. The suggested value is 7 days, which gives time to notice the resolution problem and for human intervention to fix it.

This same basic technique MAY be used to handle stale data associated with delegations. If authoritative server addresses are not able to be refreshed, resolution can possibly still be successful if the authoritative servers themselves are still up.

## **8. Implementation Caveats**

Answers from authoritative servers that have a DNS Response Code of either 0 (NOERROR) or 3 (NXDOMAIN) MUST be considered to have refreshed the data at the resolver. In particular, this means that this method is not meant to protect against operator error at the authoritative server that turns a name that is intended to be valid into one that is non-existent, because there is no way for a resolver to know intent.

Resolution is given a chance to succeed before stale data is used to adhere to the original intent of the design of the DNS. This mechanism is only intended to add robustness to failures, and to be enabled all the time. If stale data were used immediately and then a cache refresh attempted after the client response has been sent, the resolver would frequently be sending data that it would have had no trouble refreshing. As modern resolvers use techniques like pre-fetching and request coalescing for efficiency, it is not necessary that every client request needs to trigger a new lookup flow in the presence of stale data, but rather than a good-faith effort have been recently made to refresh the stale data before it is delivered to any client. The recommended period between attempting refreshes is 30 seconds.

It is important to continue the resolution attempt after the stale response has been sent, until the query resolution timeout, because some pathological resolutions can take many seconds to succeed as they cope with unavailable servers, bad networks, and other problems. Stopping the resolution attempt when the response with expired data has been sent would mean that answers in these pathological cases would never be refreshed.

Canonical Name (CNAME) records mingled in the expired cache with other records at the same owner name can cause surprising results. This was observed with an initial implementation in BIND when a hostname changed from having an IPv4 Address (A) record to a CNAME. The version of BIND being used did not evict other types in the cache when a CNAME was received, which in normal operations is not a





significant issue. However, after both records expired and the authorities became unavailable, the fallback to stale answers returned the older A instead of the newer CNAME.

[ This probably applies to other occluding types, so more thought should be given to the overall issue. ]

Keeping records around after their normal expiration will of course cause caches to grow larger than if records were removed at their TTL. Specific guidance on managing cache sizes is outside the scope of this document. Some areas for consideration include whether to track the popularity of names in client requests versus evicting by maximum age, and whether to provide a feature for manually flushing only stale records.

## 9. Implementation Status

[RFC Editor: per [RFC 6982](#) this section should be removed prior to publication.]

The algorithm described in the [Section 7](#) section was originally implemented as a patch to BIND 9.7.0. It has been in production on Akamai's production network since 2011, and effectively smoothed over transient failures and longer outages that would have resulted in major incidents. The patch was contributed to Internet Systems Consortium and the functionality is now available in BIND 9.12 via the options `stale-answer-enable`, `stale-answer-ttl`, and `max-stale-ttl`.

Unbound has a similar feature for serving stale answers, but it works in a very different way by returning whatever cached answer it has before trying to refresh expired records. This is unfortunately not faithful to the ideal that data past expiry should attempt to be refreshed before being served.

Knot Resolver has an demo module here: <https://knot-resolver.readthedocs.io/en/stable/modules.html#serve-stale>

Details of Apple's implementation are not currently known.

In the research paper "When the Dike Breaks: Dissecting DNS Defenses During DDoS" [[DikeBreaks](#)], the authors detected some use of stale answers by resolvers when authorities came under attack. Their research results suggest that more widespread adoption of the technique would significantly improve resiliency for the large number of requests that fail or experience abnormally long resolution times during an attack.



## **10. Security Considerations**

The most obvious security issue is the increased likelihood of DNSSEC validation failures when using stale data because signatures could be returned outside their validity period. This would only be an issue if the authoritative servers are unreachable, the only time the techniques in this document are used, and thus does not introduce a new failure in place of what would have otherwise been success.

Additionally, bad actors have been known to use DNS caches to keep records alive even after their authorities have gone away. This potentially makes that easier, although without introducing a new risk.

## **11. Privacy Considerations**

This document does not add any practical new privacy issues.

## **12. NAT Considerations**

The method described here is not affected by the use of NAT devices.

## **13. IANA Considerations**

This document contains no actions for IANA. This section will be removed during conversion into an RFC by the RFC editor.

## **14. Acknowledgements**

The authors wish to thank Matti Klock, Mukund Sivaraman, Jean Roy, and Jason Moreau for initial review. Feedback from Robert Edmonds, Giovane Moura, Davey Song, and Ralf Weber has also been incorporated.

## **15. References**

### **15.1. Normative References**

- [RFC1034] Mockapetris, P., "Domain names - concepts and facilities", STD 13, [RFC 1034](https://www.rfc-editor.org/info/rfc1034), DOI 10.17487/RFC1034, November 1987, <<https://www.rfc-editor.org/info/rfc1034>>.
- [RFC1035] Mockapetris, P., "Domain names - implementation and specification", STD 13, [RFC 1035](https://www.rfc-editor.org/info/rfc1035), DOI 10.17487/RFC1035, November 1987, <<https://www.rfc-editor.org/info/rfc1035>>.

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC2181] Elz, R. and R. Bush, "Clarifications to the DNS Specification", [RFC 2181](#), DOI 10.17487/RFC2181, July 1997, <<https://www.rfc-editor.org/info/rfc2181>>.
- [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>>.

## **15.2. Informative References**

- [DikeBreaks]  
Moura, G., Heidemann, J., Mueller, M., Schmidt, R., and M. Davids, "When the Dike Breaks: Dissecting DNS Defenses During DDos", ACM 2018 Internet Measurement Conference, DOI 10.1145/3278532.3278534, October 2018, <<https://www.isi.edu/~johnh/PAPERS/Moura18b.pdf>>.
- [RFC7719] Hoffman, P., Sullivan, A., and K. Fujiwara, "DNS Terminology", [RFC 7719](#), DOI 10.17487/RFC7719, December 2015, <<https://www.rfc-editor.org/info/rfc7719>>.

### Authors' Addresses

David C Lawrence  
Oracle + Dyn

Email: tale@dd.org

Warren Kumari  
Google  
1600 Amphitheatre Parkway  
Mountain View CA 94043  
USA

Email: warren@kumari.net

Puneet Sood  
Google

Email: puneets@google.com

