

Thinking about Serve Stale

Cathy Almond

2025-02-06

<https://www.isc.org>



All content © Internet Systems Consortium, Inc.

This is a combination of talk and primer covering:

- What is Serve Stale?
- A tutorial on how to configure and use it
- Things you need to be aware of if you're using it
- And a final question to you all

I'm Cathy Almond,
I work for ISC

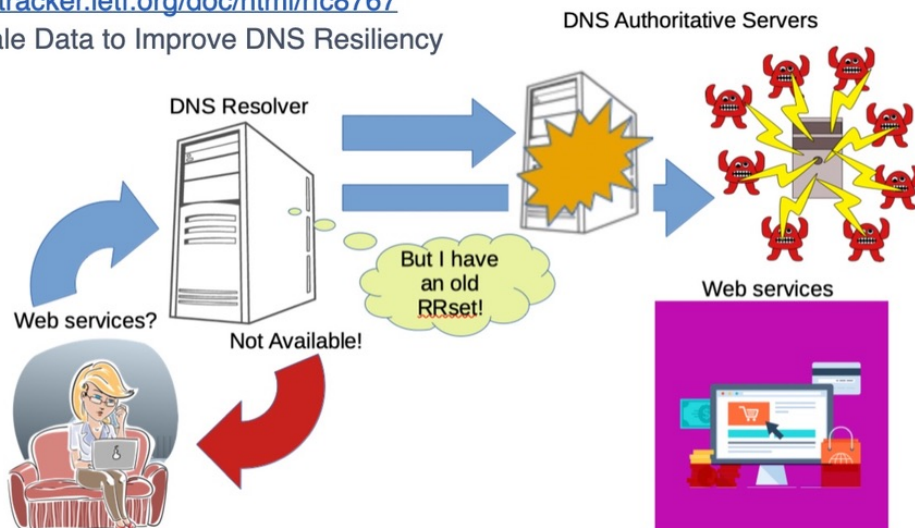
We make, maintain and distribute Open Source BIND.

But with a bit of help from my friends, I'm going to tell you about Serve Stale in Unbound, Knot Resolver and PowerDNS Recursor too

What is Serve Stale anyway?

<https://datatracker.ietf.org/doc/html/rfc8767>

Serving Stale Data to Improve DNS Resiliency



Most Internet users rely on DNS without even knowing how important it is

- I chose my T-shirt today specifically for this talk by the way!

If DNS resolution for clients wanting access to services doesn't work, then even if those services are up and running, they're not accessible

DNS is important!

But things can go wrong for authoritative DNS servers

- physical disasters
- DDoS
- People make configuration and operation mistakes

And then then DNS resolvers can't look-up what they need to to provide query responses to their clients.

It started with a big disaster back in 2016

- Successful DDoS against DynDB
- Many popular sites effectively 'down' because Resolvers couldn't reach the authoritative servers
- Management clamoured for a mitigation for this situation

These resolvers - WHAT IF they already had the answers already- DNS resolvers have caches

Resolvers are not supposed to use OLD answers - the DNS protocol includes the Time To Live (think of it as the use-by date on the food you buy).

In fact, think of your resolver as a refrigerator and the records from authoritative servers as the food you bought from the grocery store ...

But slightly old answers are very likely (most of the time) going to be correct and better than no answers at all, if something has Gone Wrong.

So why not use them as the answer of last resort, and keep the services they point to from becoming inaccessible?

If you're hungry, stale bread is better than no bread at all...

And so, Serve Stale for DNS Resolvers was born...

There's an Internet Standard (RFC) on this - written after several implementations already existed.

It's actually very well put together and easy to read and understand - and I'll cover the general principles in the next couple of slides.

How might this work?

- How long do you want to keep 'stale' cache content?
- How long do you keep a client waiting before providing a stale answer?
- How often do you try to refresh stale content?
- Should you give up eventually and stop serving stale RRsets?
- Should you indicate to the client that the answer is old data?



Resolvers have caches

Caches already have retention and eviction strategies for their content

- Ideally keep popular content available all the time if possible
- Include TTL and 'last accessed' information
- Maybe prefetch to refresh RRsets before they expire?

BUT now we're planning to potentially use stale RRsets as a backup

- Should we therefore keep content for longer?
- How much longer?
- What's this going to do for cache memory consumption by the way?

AND how to decide when to use stale RRsets - how long to keep the client waiting?

DO you try to refresh those RRsets on every client query?

And really - when is the old stuff JUST TOO OLD to use?

And ... psst! ... should you tell the client that you've given them old answers?

RFC 8767 says:

Four notable timers drive considerations for the use of stale data:

- 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.
- A query resolution timer, which caps the total amount of time a recursive resolver spends processing the query.
- A failure recheck timer, which limits the frequency at which a failed lookup will be attempted again.
- A maximum stale timer, which caps the amount of time that records will be kept past their expiration.



RFC 8767 recommends that Resolvers should have these timers ...

Actually, apart from the last one, they would have had them all anyway in some form or other

- You don't want your clients to have to wait too long (even though they usually retry their queries)
- You don't want your resolver to be doing more work that it needs to

So the practical difference with serve-stale enabled is that the client (if possible) will get old content in a query response instead of SERVFAIL

SERVE STALE instead of SERVFAIL ??

Tempered of course by not serving content so old that it's going to give the clients virtual food poisoning!

RFC 8767 also talks about:

- What TTL should you put on a stale RRset?
- Should you ever serve stale RRsets that were originally received with TTL=0 ('don't cache')?
- Also use stale content in-path to getting the answers for the client query (such as NS records and addresses?)
- Does receipt of a non-authoritative error (e.g. FORMERR, SERVFAIL, REFUSED) count as a failed refresh?
- Can you signal that stale answers have been provided



Oooh - some new things to think about we didn't touch upon before

- You're giving your client old expired stuff - already past its TTL. How long do you want the client to use it for before coming back to ask you again?

- Some authoritative servers send 'one time use only' RRsets - what should you do with those?

(ASIDE: Authoritative and Resolver operators alike pay attention to the received TTL=0 scenario later - whether you like it or not, there ARE resolver handling differences!)

- What is a failure anyway (not all auth problems are down to physical disasters and deliberate attack - sometimes people just mess things up...)

Phew! Let's talk about actually using it!

Two main things for operators to consider when enabling Serve Stale:

1. Cache - how long can/should you keep stale RRsets in cache and will there be any memory consumption implications of this for your server(s)
2. User experience - how long should clients wait before being served with stale content, and how often should they expect their resolver to attempt to refresh the old content



This slide speaks for itself:

1. If you enable Serve Stale, are you going to run into memory consumption problems? Need to plan ahead?

PARTICULARLY on the memory consumption thing - bear in mind that it's not just the popular content that is going to be retained for longer, but the one-time-use content too, and this could cause cache management issues in some implementations.

2. How do you want things to be for your clients/users?

How to enable and configure Serve-Stale

What you need to think about and how to configure Serve Stale in:

- BIND
- Unbound
- Knot Resolver
- PowerDNS Recursor



Disclaimer: I know a lot about BIND - obviously!

Hopefully I won't make too many mistakes when explaining how the other three implementations work!

(But they did all approve the slides - so use those afterwards as the source of truth rather than what I say!)

====

NOW, upfront I am going to say 'there is a lot of data in the slides that follow'

I'm NOT going to explain all the options and things to think about down to the tiniest detail

The approach I've taken, is that the slideset is your PRACTICAL TAKEAWAY, to use later.

So instead of bullet points on the slides that the speaker elaborates upon...

I'm elaborating on the slides, but giving you bullet points in my talk over!

So I recommend that you listen rather than read (unless you can multi-task effectively)

Enabling Serve Stale in BIND

- Enable the stale cache (off by default in all currently supported versions):

```
stale-cache-enable yes;
```

- Configure how long to retain stale content (default 1 day):

```
max-stale-ttl 24H;
```

- Are you going to enable serving of stale answers now, or turn the feature on/off when required?

```
stale-answer-enable yes;
```

or

```
rndc serve-stale on|off|reset
```

Note that you can't enable stale answers without first starting BIND with stale cache enabled - no stale cache means no stale answers available.



BIND has two separate things to enable

- Retaining of stale content (in case it's needed)
- Activating the processes to use stale content if it can't be refreshed

Neither are enabled by default

It's possible to enable stale cache and then turn on stale answers dynamically afterwards

How would you know you needed to do this though - and probably by the time you realised you needed to, there could be a lot of unhappiness abounding?

(Optional) Configuring stale answers in BIND

These options are only effective when both `stale-cache-enable` and `stale-answers-enable` are set to 'yes'.

- Configure the TTL used for stale RRsets in query responses (default 30s):

```
stale-answer-ttl 30s;
```

- How long to serve stale before reattempting to refresh (default 30s):

```
stale-refresh-time 30s;
```

- How long do clients have to wait for a stale answer instead of the usual recursion timeout and `SERVFAIL`?

```
stale-answer-client-timeout <duration>;
```

*The duration above can be either zero ('0') or disabled/off (default off). Recursion duration is controlled by option 'resolver-query-timeout', default 10s. **ISC support recommends not altering this!***



You've got the knobs you need (if you want them) for tweaking the standard client experience - but the defaults are sane and most users of Serve Stale won't need to adjust them.

And then there's an interesting knob: `stale-answer-client-timeout`

Oh we have had some fun with this one ...

The original idea was that you could decide how long the client had to wait before getting a stale answer (but still leaving the cache refresh to complete as normal - successfully or otherwise).

Where we're at now is that you either just have to let the client wait as usual until we replace `SERVFAIL` with Serve Stale

OR

If there's a stale cache hit, then serve that *immediately* and leave the refresh to do its thing afterwards.

Would you want to do that? (More on that at the end (BIND is not the only implementation that allows this) ...

BIND and stale answers - be aware that:

- BIND's cache content management strategy is two-fold. Opportunistic cache cleaning means that in-passing, TTL-expired content is removed. To enable stale cache, most RRsets that would otherwise be candidates for eviction, are now retained (if possible) for the period `max-stale-ttl`. Therefore your cache memory consumption may increase significantly, up to the point that memory-based (least recently used) cache cleaning is triggered.
- Content received with TTL=0 is not retained for use in stale answers, *unless* option `min-cache-ttl` has been used to override small or zero TTLs on received RRsets.
- NXDOMAIN RRsets are not retained for use in stale answers (this differs from other implementations).
- BIND automatically adds EDE options 3 (Stale Answer) or 19 (Stale NXDOMAIN) to query responses with stale content (although see above!).
- A SERVFAIL answer due to fetch-limits may instead be replaced with stale content but doesn't trigger `stale-refresh-time`.



The biggest takeaway point is the cache consumption one. BIND's routine housekeeping of expired content is less processing-intensive than the memory-based mechanism that is triggered when max-cache-size is reached.

Monitor and provision more memory/cache and/or reduce how long you retain stale content for if you need to and you will be happier...

And yes, we do do EDE. Both EDEs - except that we decided later not to apply Serve Stale to NXDOMAINs, so you'll never see EDE 19 in current BIND with Serve Stale enabled.

Enabling Serve Stale in Unbound

- Enable the serving of stale content (disabled by default):
serve-expired: yes
- Enable prefetch (to help keep popular cache up-to-date):
prefetch: yes
- How long do you want stale RRsets to continue to be used for serving stale answers (default is all content is eligible indefinitely; recommendation is 1 day)?
serve-expired-ttl: 86400 # 1 day (in seconds)
- Decide whether or not you want automatic extension of use of expired stale RRsets (default no)
serve-expired-ttl-reset: yes|no

This setting extends/resets the “use this stale until” timestamp added to stale RRsets in cache when a refresh attempt fails. Without it, content is only used stale until the limit specified in serve-expired-ttl.



The original default behaviour with “serve-expired: yes” AND NOTHING ELSE CONFIGURED is that all content in cache is eligible to be served stale if needed.

This works because Unbound has no expired-TTL cleaning mechanism, so as long as there is still room in cache, the content remains available.

Alternatively, you can say ‘only use these old RRsets until a given limit and never again thereafter

And on top of that, you can also opt to reset that ‘only for so long’ timer - but that requires a full attempt to refresh before the old content gets sent to the client.

(Optional) Configuring stale answers in Unbound

- Configure the TTL used for stale RRsets in query responses (default 30s):
`serve-expired-reply-ttl: 30 # 30 seconds (in seconds)`
- How long do clients have to wait for a stale answer instead of the usual recursion timeout and SERVFAIL?

`serve-expired-client-timeout: 1800 #1.8 seconds (in milliseconds)`

The above example is deliberately just shorter than most client-side query response timeouts of 2 seconds, but the default is zero - respond stale first and attempt to refresh afterwards.

- Configure whether or not you want to send Extended DNS Errors (EDE) in query responses containing stale RRsets:

`ede: yes`

`ede-serve-expired: yes`



Again, and similar to BIND, in Unbound you've got the knobs you need (if you want them) for tweaking the standard client experience - but the defaults are sane and most users of Serve Stale won't need to adjust them.

... except possibly `serve-expired-client-timeout`

This works similarly to the BIND option `stale-answer-client-timeout` - it's about how long to wait for refresh recursion to complete for the first time, before opting to use an existing stale RRset for the reply to the client. But in Unbound you can make more use of it, as well as opting to serve stale first and then refresh afterwards.

Also note that you will need to enable both EDE and the serving of EDE when providing stale answers to clients.

Unbound and stale answers - useful to know:

- Unbound's cache content management strategy is based on eviction of RRsets once cache memory limits are reached and on a Least Recently Used (LRU) basis. Enabling Serve Stale therefore is unlikely to change cache composition or memory consumption.
- Unbound has an additional back-end cache DB feature available which may change how Serve Stale operates - see <https://unbound.docs.nlnetlabs.nl/en/latest/manpages/unbound.conf.html#cache-db-module-options>
- Content received with TTL=0 is not retained for use in stale answers (because it is never added to cache) *unless* options cache-min-ttl and/or cache-min-negative-ttl have been used to override small or zero TTLs on received RRsets.



The takeaway from this slide is that you shouldn't need to worry overly about cache consumption in Unbound if you enable stale answers - it's not going to change cache consumption. Cache content that was least recently used will still be evicted once cache memory limits have been reached.

Also note that back-end cache DB feature - I did not dig into this and I leave its functionality for the reader to research!

Unbound and stale answers - useful to know:

- There are no options to directly control how frequently Unbound attempts to refresh stale content before using it again; client queries will continue to use the stale content without initiating a new refresh until 5 seconds have passed since the last failure.
- The client experience with `serve-expired-ttl-reset: yes` can vary because although a failed refresh will reset the `serve-expired-ttl` timer on the `RRset`, if that `RRset` was no longer eligible for stale use *before* the refresh, it won't be sent to the client until *after* the refresh attempt has finished processing and has failed (`serve-expired-client-timeout` isn't used for non-eligible stale content).



What this means is that after the first 'serve-expired-client-timeout' and then use of stale RRsets, those same stale RRsets will continue to be used immediately for subsequent client queries until 5 seconds have passed, at which point, the next client query will cause another refresh attempt to be initiated.

And 'serve-expired-ttl-reset' is complicated. What it essentially means is that IF you decide to set the limit on 'this stale content is so stale that it would now be poisonous if used' (serve-expired-ttl) instead of leaving all stale content eligible until it's evicted from cache, that you DO have the possibility of using it in desperation. HOWEVER, a thorough attempt to refresh it will take place first (serve-expired-client-timeout does NOT apply here). But at least there is the possibility of 'answers of last resort' - should you want them?

Enabling Serve Stale in Knot Resolver

- Enable the serving of stale content (disabled by default):

options:

serve-stale: true

This setting enables the use of stale RRsets in order to avoid sending SERVFAIL when cache content can't be refreshed. They are usable for up to 24 hours after TTL expiry. This stale validity period is not directly configurable.



Knot Resolver has a simplistic approach to Serve Stale, with few options. Basically, just enable it and it's 'on'!

Knot Resolver and stale answers - need to know:

- Knot Resolver's cache content management strategy is based on eviction of RRsets once cache memory limits are reached. Enabling `serve-stale` therefore is unlikely to change cache composition or memory consumption.
- Stale RRsets are returned in query responses with a 1 second TTL.
- There is no directly-configurable client query/resolver timeout.
- It's nevertheless possible to change some timers (e.g. `.timeout` and `.callback`) by directly modifying the `lua` module: https://gitlab.nic.cz/knot/knot-resolver/-/blob/master/modules/serve_stale/serve_stale.lua
- Content received with `TTL=0` is eligible to be used in stale answers *even if* Knot Resolver option `cache/ttl-min` (default 5s) is set to zero.
- Knot Resolver (version 6 and up) automatically adds EDE options 3 or 19 to query responses with stale content.



This option is not specifically for Serve Stale - it is about setting timeouts during recursion anyway. If "serve-stale: true" had not been set, then the client would receive SERVFAIL instead of stale content.

You probably don't want to adjust this.

(Optional) Configuring stale answers in Knot Resolver

- Adjust how long to wait before retrying unreachable servers (default 1000 ms):

cache:

ns-timeout: <time ms|s|m|h|d>

As long as as all nameservers are marked unreachable and have been for less than ns-timeout, clients will receive immediate stale answers without an attempt to refresh first.



Knot Resolver's cache management strategy is similar to Unbound's - so enabling Serve Stale isn't going to change much about cache consumption or cache composition.

As noted before, there's not much that you can tweak in the Knot Resolver configuration to change Serve Stale behaviour, but you could modify some timers in the lua module instead.

But please note one BIG DIFFERENCE in relation to TTL=0 content - Knot resolver already defaults to overriding the cache TTL on RRsets received with TTL < 5s. These RRsets are added to cache regardless - and they will still be added to cache if this override is disabled using cache/ttl-min. There is nothing in Knot Resolver's Serve Stale handling to exempt 'one time use' TTL=0 RRsets from use when stale, if they cannot be refreshed.

Enabling Serve Stale in PowerDNS Recursor

- Enable the serving of stale content (disabled by default) by choosing how many times to reset a stale RRset's TTL in the configuration or command line:

serve-stale-extensions <count>

This setting defines how many times a stale RRset can be 'revived' in cache by giving it a TTL that is the smaller of 30s or the original TTL of the RRset.

The 'revival' only happens after the RRset refresh attempt fails.

The 'revival' mechanism takes into account how long ago the last extension was done when doing a new extension. RRsets that are only queried once in a while get the same maximum stale period as compared with RRsets queried very frequently (assuming original TTL >= 30s).



PowerDNS Recursor has taken a slightly different approach to implementing Serve Stale.

Note that all implementations have to consider:

- What TTL to put on stale answers served to clients
- How long to continue to use stale RRsets from cache before attempting to refresh them again

PowerDNS Recursor handles this when it 'revives' stale content by giving a new TTL of 30s in cache, which then 'counts down' in the same way as the originally received TTL from the authoritative servers did - after which it can't be served again without attempting to refresh it.

There is also a 'count' that you can configure - which is how many times this 'add another 30s TTL' can be allowed to happen.

Ah but (you might say), "what if the RRset is only queried once every 5 minutes, does that mean it can survive in cache to be reused much longer than an RRset that is queried every second?" I asked. Apparently not - the <count> and calculation is done on the basis that the refresh attempt was done after each 30s, so old content that is infrequently queried doesn't get reused for longer than popular content.

Also note that the 30s TTL used when 'reviving' stale content might not be 30s if the original authoritative TTL was shorter than 30s. (But the calculations on how long this content is available are still done on <count> x 30s).

PowerDNS and stale answers - useful to know:

- PowerDNS Recursor's cache content management strategy is based both on eviction of expired content and on eviction of RRsets based on LRU once cache memory limits are reached. Enabling Serve Stale defers eviction of expired content until `serve-stale-extensions` x 30s. Therefore your cache memory consumption will increase and there may be more LRU cleaning occurring.
- Content received with TTL=0 is eligible to be used in stale answers *even if* PowerDNS Recursor option `minimum-ttl-override` (default 1s) is set to zero.
- There is no specific configurable timer for how long to wait for on a refresh attempt before serving stale content; the time to wait is the same, irrespective of whether PowerDNS would respond with `SERVFAIL` or use stale RRsets.
- The TTL on RRsets used stale is 30s or less, and counts down in cache in the same way as the original RRset's TTL would have done.
- PowerDNS does not (yet) add EDE options to query responses with stale content, but this is planned as a future feature.



PowerDNS Recursor is akin to BIND in that there are two cache eviction strategies, one based on TTL (content expiry) and the other on reaching cache memory limits, so you will need to think about cache content and cache memory consumption if you enable Serve Stale.

And coming soon - EDE options.

Other resolver implementations exist, for example:

- Appliances that use Open Source DNS 'inside the box' will almost certainly enable Serve Stale and provide some configuration options.
- Akamai DNS has a switch (default 'on') for serving stale content - it is an extension to prefetch functionality for 'popular' content
- Cloud-based resolver solutions appear to implement some types of 'failure mode' recovery, and this probably includes serving of stale content (*disclaimer - I have not researched this!*)



If you're using something other than the DNS resolver implementations I covered then please ask your provider for advice on whether or not they provided this feature and if they do, how it works and what configurable options are available.

Not covered in this presentation...

- Logging and statistics (mostly included in software providers' documentation).
- Serve Stale in relation to DNSSEC validation failures and stale DNSSEC material.
- Use of stale cache data when following intermediate referrals.
- Expired cache containing both CNAME/DNAME and other records for the same name (due to different query responses at different times).
- “Fail” responses (as opposed to failure to respond by authoritative servers) - which of these trigger Serve Stale?

It's reasonable to assume that the software implementations handle all of the above sanely - but if you need to know - ask your software provider!



Please refer to the documentation for logging and statistics - and also take a look at both of them if you enable Serve Stale in your environment . Is it being triggered? Is it helping you?

There are a lot of nuances and edge scenarios that I didn't dig into (well, I dug into some of them while preparing this talk, mainly to make sure I understood the options correctly, but also because I was fascinated by how the underlying resolver architecture drove the Serve Stale implementation choices made by the developers).

If you have any specific “what if...?” questions, then I would encourage you to ask your software provider directly.

A challenge!

- Is the Serve Stale feature genuinely useful?
- Do operators of resolvers know for certain that it is helping their servers?
- I have never yet heard of a situation where having Serve Stale enabled 'saved the day' - but would love to hear from anyone who has one!
- One 'useful' that is perhaps passing under the radar might be that resolvers are 'smoothing over' short term authoritative server unavailability or slowdowns. *Would we even know that this is happening?*
- What do we think about 'serve stale content first, refresh afterwards' (available in some implementations) - and does that count as 'useful'?
- Analysis of Resolver logging and statistics from servers in production environments *might* be the starting point for an assessment of the practical benefits of the Serve Stale feature - a talk for OARC45?



I asked the questions above at the end of my talk. Please watch the recording if you want to hear the discussion, but essentially:

- Several operators confirmed my hypothesis that it does get triggered, but on a small scale ("smoothing over the cracks...")
- No-one has yet seen it 'save the day' on a big DNS authoritative services outage
- One hypothesis on "smoothing over the cracks" is that the cracks are BGP and other routing 'flaps'.

So the conclusion we came to, is that this feature IS genuinely useful, but not a knight in shining armour riding in to save the day - more it's a series of small bandages being applied just to hold things together from time to time on an ad-hoc basis to improve the 'client experience'.

(We didn't touch on 'serve stale content first, refresh afterwards' – which sounds tempting if you are being measured on time to respond to client queries – but do you really want to serve a stale answer when there may be a different answer available from the authoritative servers once that refresh has completed?).

Discussion time ...

RFC 8767: “Stale bread is better than no bread.”



Photo by Vincent van Zelst via Wikimedia Commons



References:

- <https://kb.isc.org/docs/serve-stale-implementation-details>
- <https://kb.isc.org/docs/changes-to-serve-stale-option-stale-answer-client-timeout-in-bind-918-and-newer>
- <https://unbound.docs.nlnetlabs.nl/en/latest/topics/core/serve-stale.html#serving-stale-data>
- <https://blog.nlnetlabs.nl/some-country-for-old-men/>
- https://websites.pages.nic.cz/knot-resolver.cz/documentation/v5.7.4/modules-serve_stale.html
- <https://www.knot-resolver.cz/documentation/latest/config-serve-stale.html>
- https://gitlab.nic.cz/knot/knot-resolver/-/blob/master/modules/serve_stale/serve_stale.lua
- <https://docs.powerdns.com/recursor/appendices/internals.html#serve-stale>
- <https://docs.powerdns.com/recursor/settings.html#>



Thanks to:

- [Wouter Wijngaards](#) and [Yorgos Thessalonikefs](#) from [NL.net Labs](#) for detailed information about Unbound's Serve Stale implementation
- [Vladimir Cunat](#) of [CZ.NIC](#) for explaining what happens with the different configuration options in Knot Resolver
- [Otto Moerbeek](#) at [PowerDNS](#) for disentangling in my mind the different approach taken to Serve Stale by [PowerDNS Recursor](#)

Also acknowledging the royalty-free image creators contributing to slide 2:

- [clickschool](#) of [Pixabay](#)
- [Megan Rexazin Conde](#) of [Pixabay](#)
- [GraphicMama-team](#) of [Pixabay](#)
- [Yayamamo](#) of [Wikimedia Commons](#)



Thank you for listening.

- ISC main website: <https://www.isc.org>
- Software downloads: <https://www.isc.org/download> or <https://downloads.isc.org>
- Presentations: <https://www.isc.org/presentations>
- Main GitLab: <https://gitlab.isc.org>

