

Document: [draft-cheshire-dnsext-multicastdns-05.txt](#)
Category: Standards Track
Expires 7th December 2005

Stuart Cheshire
Marc Krochmal
Apple Computer, Inc.
7th June 2005

Multicast DNS

<[draft-cheshire-dnsext-multicastdns-05.txt](#)>

Status of this Memo

By submitting this Internet-Draft, each author represents that any applicable patent or other IPR claims of which he or she is aware have been or will be disclosed, and any of which he or she becomes aware will be disclosed, in accordance with [Section 6 of BCP 79](#). For the purposes of this document, the term "[BCP 79](#)" refers exclusively to [RFC 3979](#), "Intellectual Property Rights in IETF Technology", published March 2005.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

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."

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/1id-abstracts.html>

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>

Abstract

As networked devices become smaller, more portable, and more ubiquitous, the ability to operate with less configured infrastructure is increasingly important. In particular, the ability to look up DNS resource record data types (including, but not limited to, host names) in the absence of a conventional managed DNS server, is becoming essential.

Multicast DNS (mDNS) provides the ability to do DNS-like operations on the local link in the absence of any conventional unicast DNS server. In addition, mDNS designates a portion of the DNS namespace to be free for local use, without the need to pay any annual fee, and without the need to set up delegations or otherwise configure a conventional DNS server to answer for those names.

The primary benefits of mDNS names are that (i) they require little or no administration or configuration to set them up, (ii) they work when no infrastructure is present, and (iii) they work during infrastructure failures.

Table of Contents

1.	Introduction.....	3
2.	Conventions and Terminology Used in this Document.....	4
3.	Multicast DNS Names.....	5
4.	Source Address Check.....	8
5.	Reverse Address Mapping.....	9
6.	Querying.....	9
7.	Duplicate Suppression.....	13
8.	Responding.....	15
9.	Probing and Announcing on Startup.....	18
10.	Conflict Resolution.....	22
11.	Resource Record TTL Values and Cache Coherency.....	23
12.	Special Characteristics of Multicast DNS Domains.....	28
13.	Multicast DNS for Service Discovery.....	30
14.	Enabling and Disabling Multicast DNS.....	30
15.	Considerations for Multiple Interfaces.....	31
16.	Multicast DNS and Power Management.....	32
17.	Multicast DNS Character Set.....	33
18.	Multicast DNS Message Size.....	34
19.	Multicast DNS Message Format.....	35
20.	Choice of UDP Port Number.....	38
21.	Summary of Differences Between Multicast DNS and Unicast DNS..	39
22.	Benefits of Multicast Responses.....	40
23.	IPv6 Considerations.....	41
24.	Security Considerations.....	42
25.	IANA Considerations.....	43
26.	Acknowledgments.....	43
27.	Copyright Notice.....	43
28.	Normative References.....	44
29.	Informative References.....	44
30.	Authors' Addresses.....	45

Expires 7th December 2005

Cheshire & Krochmal

[Page 2]

1. Introduction

When reading this document, familiarity with the concepts of Zero Configuration Networking [[ZC](#)] and automatic link-local addressing [[RFC 2462](#)] [[RFC 3927](#)] is helpful.

This document proposes no change to the structure of DNS messages, and no new operation codes, response codes, or resource record types. This document simply discusses what needs to happen if DNS clients start sending DNS queries to a multicast address, and how a collection of hosts can cooperate to collectively answer those queries in a useful manner.

There has been discussion of how much burden Multicast DNS might impose on a network. It should be remembered that whenever IPv4 hosts communicate, they broadcast ARP packets on the network on a regular basis, and this is not disastrous. The approximate amount of multicast traffic generated by hosts making conventional use of Multicast DNS is anticipated to be roughly the same order of magnitude as the amount of broadcast ARP traffic those hosts already generate.

New applications making new use of Multicast DNS capabilities for unconventional purposes may generate more traffic. If some of those new applications are "chatty", then work will be needed to help them become less chatty. When performing any analysis, it is important to make a distinction between the application behavior and the underlying protocol behavior. If a chatty application uses UDP, that doesn't mean that UDP is chatty, or that IP is chatty, or that Ethernet is chatty. What it means is that the application is chatty. The same applies to any future applications that may decide to layer increasing portions of their functionality over Multicast DNS.

2. Conventions and Terminology Used in this Document

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in "Key words for use in RFCs to Indicate Requirement Levels" [[RFC 2119](#)].

This document uses the term "host name" in the strict sense to mean a fully qualified domain name that has an address record. It does not use the term "host name" in the commonly used but incorrect sense to mean just the first DNS label of a host's fully qualified domain name.

A DNS (or mDNS) packet contains an IP TTL in the IP header, which is effectively a hop-count limit for the packet, to guard against routing loops. Each Resource Record also contains a TTL, which is the number of seconds for which the Resource Record may be cached.

In any place where there may be potential confusion between these two types of TTL, the term "IP TTL" is used to refer to the IP header TTL (hop limit), and the term "RR TTL" is used to refer to the Resource Record TTL (cache lifetime).

When this document uses the term "Multicast DNS", it should be taken to mean: "Clients performing DNS-like queries for DNS-like resource records by sending DNS-like UDP query and response packets over IP Multicast to UDP port 5353."

This document uses the terms "shared" and "unique" when referring to resource record sets.

A "shared" resource record set is one where several Multicast DNS responders may have records with that name, rrtype, and rrclass, and several responders may respond to a particular query.

A "unique" resource record set is one where all the records with that name, rrtype, and rrclass are under the control or ownership of a single responder, and at most one responder should respond to any given query. Before claiming ownership of a unique resource record set, a responder MUST probe to verify that no other responder already claims ownership of that set, as described in [Section 9.1](#) "Probing".

Strictly speaking the terms "shared" and "unique" apply to resource record sets, not to individual resource records, but it is sometimes convenient to talk of "shared resource records" and "unique resource records". When used this way, the terms should be understood to mean a record that is a member of a "shared" or "unique" resource record set, respectively.

Expires 7th December 2005

Cheshire & Krochmal

[Page 4]

3. Multicast DNS Names

This document proposes that the DNS top-level domain ".local." be designated a special domain with special semantics, namely that any fully-qualified name ending in ".local." is link-local, and names within this domain are meaningful only on the link where they originate. This is analogous to IPv4 addresses in the 169.254/16 prefix, which are link-local and meaningful only on the link where they originate.

Any DNS query for a name ending with ".local." MUST be sent to the mDNS multicast address (224.0.0.251 or its IPv6 equivalent FF02::FB).

It is unimportant whether a name ending with ".local." occurred because the user explicitly typed in a fully qualified domain name ending in ".local.", or because the user entered an unqualified domain name and the host software appended the suffix ".local." because that suffix appears in the user's search list. The ".local." suffix could appear in the search list because the user manually configured it, or because it was received in a DHCP option, or via any other valid mechanism for configuring the DNS search list. In this respect the ".local." suffix is treated no differently to any other search domain that might appear in the DNS search list.

DNS queries for names that do not end with ".local." MAY be sent to the mDNS multicast address, if no other conventional DNS server is available. This can allow hosts on the same link to continue communicating using each other's globally unique DNS names during network outages which disrupt communication with the greater Internet. When resolving global names via local multicast, it is even more important to use DNSSEC or other security mechanisms to ensure that the response is trustworthy. Resolving global names via local multicast is a contentious issue, and this document does not discuss it in detail, instead concentrating on the issue of resolving local names using DNS packets sent to a multicast address.

A host which belongs to an organization or individual who has control over some portion of the DNS namespace can be assigned a globally unique name within that portion of the DNS namespace, for example, "cheshire.apple.com." For those of us who have this luxury, this works very well. However, the majority of home customers do not have easy access to any portion of the global DNS namespace within which they have the authority to create names as they wish. This leaves the majority of home computers effectively anonymous for practical purposes.

To remedy this problem, this document allows any computer user to

elect to give their computers link-local Multicast DNS host names of the form: "single-dns-label.local." For example, a laptop computer may answer to the name "cheshire.local." Any computer user is granted the authority to name their computer this way, provided that the chosen host name is not already in use on that link. Having named

their computer this way, the user has the authority to continue using that name until such time as a name conflict occurs on the link which is not resolved in the user's favour. If this happens, the computer (or its human user) SHOULD cease using the name, and may choose to attempt to allocate a new unique name for use on that link. These conflicts are expected to be relatively rare for people who choose reasonably imaginative names, but it is still important to have a mechanism in place to handle them when they happen.

The point made in the previous paragraph is very important and bears repeating. It is easy for those of us in the IETF community who run our own name servers at home to forget that the majority of computer users do not run their own name server and have no easy way to create their own host names. When these users wish to transfer files between two laptop computers, they are frequently reduced to typing in dotted-decimal IP addresses because they simply have no other way for one host to refer to the other by name. This is a sorry state of affairs. What is worse, most users don't even bother trying to use dotted-decimal IP addresses. Most users still move data between machines by copying it onto a floppy disk or similar removable media.

In a world of gigabit Ethernet and ubiquitous wireless networking it is a sad indictment of the networking community that the preferred communication medium for most computer users is still the floppy disk.

Allowing ad-hoc allocation of single-label names in a single flat ".local." namespace may seem to invite chaos. However, operational experience with AppleTalk NBP names [[NBP](#)], which on any given link are also effectively single-label names in a flat namespace, shows that in practice name collisions happen extremely rarely and are not a problem. Groups of computer users from disparate organizations bring Macintosh laptop computers to events such as IETF Meetings, the Mac Hack conference, the Apple World Wide Developer Conference, etc., and complaints at these events about users suffering conflicts and being forced to rename their machines have never been an issue.

Enforcing uniqueness of host names (i.e. the names of DNS address records mapping names to IP addresses) is probably desirable in the common case, but this document does not mandate that. It is permissible for a collection of coordinated hosts to agree to maintain multiple DNS address records with the same name, possibly for load balancing or fault-tolerance reasons. This document does not take a position on whether that is sensible. It is important that both modes of operation are supported. The Multicast DNS protocol allows hosts to verify and maintain unique names for resource records where that behavior is desired, and it also allows hosts to maintain multiple resource records with a single shared name where that

behavior is desired. This consideration applies to all resource records, not just address records (host names). In summary: It is required that the protocol have the ability to detect and handle name conflicts, but it is not required that this ability be used for every record.

3.1 Governing Standards Body

Note that this use of the ".local." suffix falls under IETF jurisdiction, not ICANN jurisdiction. DNS is an IETF network protocol, governed by protocol rules defined by the IETF. These IETF protocol rules dictate character set, maximum name length, packet format, etc. ICANN determines additional rules that apply when the IETF's DNS protocol is used on the public Internet. In contrast, private uses of the DNS protocol on isolated private networks are not governed by ICANN. Since this proposed change is a change to the core DNS protocol rules, it affects everyone, not just those machines using the ICANN-governed Internet. Hence this change falls into the category of an IETF protocol rule, not an ICANN usage rule.

3.2 Private DNS Namespaces

Note also that the special treatment of names ending in ".local." has been implemented in Macintosh computers since the days of Mac OS 9, and continues today in Mac OS X. There are also implementations for Linux and other platforms [[dotlocal](#)]. Operators setting up private internal networks ("intranets") are advised that their lives may be easier if they avoid using the suffix ".local." in names in their private internal DNS server. Alternative possibilities include:

- .intranet
- .internal
- .private
- .corp
- .home

Another alternative naming scheme, advocated by Professor D. J. Bernstein, is to use a numerical suffix, such as ".6." [[djbd1](#)].

3.3 Maximum Multicast DNS Name Length

[RFC 1034](#) says:

"the total number of octets that represent a domain name (i.e., the sum of all label octets and label lengths) is limited to 255."

This text implies that the final root label at the end of every name is included in this count (a name can't be represented without it), but the text does not explicitly state that. Implementations of Multicast DNS MUST include the label length byte of the final root label at the end of every name when enforcing the rule that no name may be longer than 255 bytes. For example, the length of the name "apple.com." is considered to be 11, which is the number of bytes it takes to represent that name in a packet without using name compression:

| 0x05 | a | p | p | l | e | 0x03 | c | o | m | 0x00 |

Expires 7th December 2005

Cheshire & Krochmal

[Page 7]

4. Source Address Check

All Multicast DNS responses (including responses sent via unicast) SHOULD be sent with IP TTL set to 255. This is recommended to provide backwards-compatibility with older Multicast DNS clients that check the IP TTL on reception to determine whether the packet originated on the local link. These older clients discard all packets with TTLs other than 255.

A host sending Multicast DNS queries to a link-local destination address (including the 224.0.0.251 link-local multicast address) MUST only accept responses to that query that originate from the local link, and silently discard any other response packets. Without this check, it could be possible for remote rogue hosts to send spoof answer packets (perhaps unicast to the victim host) which the receiving machine could misinterpret as having originated on the local link.

The test for whether a response originated on the local link is done in two ways:

- * All responses sent to the link-local multicast address 224.0.0.251 are necessarily deemed to have originated on the local link, regardless of source IP address. This is essential to allow devices to work correctly and reliably in unusual configurations, such as multiple logical IP subnets overlayed on a single link, or in cases of severe misconfiguration, where devices are physically connected to the same link, but are currently misconfigured with completely unrelated IP addresses and subnet masks.
- * For responses sent to a unicast destination address, the source IP address in the packet is checked to see if it is an address on a local subnet. An address is determined to be on a local subnet if, for (one of) the address(es) configured on the interface receiving the packet, $(I \& M) == (P \& M)$, where I and M are the interface address and subnet mask respectively, P is the source IP address from the packet, '&' represents the bitwise logical 'and' operation, and '==' represents a bitwise equality test.

Since queriers will ignore responses apparently originating outside the local subnet, responders SHOULD avoid generating responses that it can reasonably predict will be ignored. This applies particularly in the case of overlayed subnets. If a responder receives a query addressed to the link-local multicast address 224.0.0.251, from a source address not apparently on the same subnet as the responder, then even if the query indicates that a unicast response is preferred (see [Section 6.5](#), "Questions Requesting Unicast Responses"), the responder SHOULD elect to respond by multicast anyway, since it can

reasonably predict that a unicast response with an apparently non-local source address will probably be ignored.

5. Reverse Address Mapping

Like ".local.", the IPv4 and IPv6 reverse-mapping domains are also defined to be link-local.

Any DNS query for a name ending with "254.169.in-addr.arpa." MUST be sent to the mDNS multicast address 224.0.0.251. Since names under this domain correspond to IPv4 link-local addresses, it is logical that the local link is the best place to find information pertaining to those names. As an optimization, these queries MAY be first unicast directly to the address in question, but if this query is not answered, the query MUST also be sent via multicast, to accommodate the case where the machine in question is not answering for itself (for example, because it is currently sleeping).

Likewise, any DNS query for a name ending with "0.8.e.f.ip6.arpa." MUST be sent to the IPv6 mDNS link-local multicast address FF02::FB, with or without an optional initial query unicast directly to the address in question.

6. Querying

There are three kinds of Multicast DNS Queries, one-shot queries of the kind made by today's conventional DNS clients, one-shot queries accumulating multiple responses made by multicast-aware DNS clients, and continuous ongoing Multicast DNS Queries used by IP network browser software.

A Multicast DNS Responder that is offering records that are intended to be unique on the local link MUST also implement a Multicast DNS Querier so that it can first verify the uniqueness of those records before it begins answering queries for them.

6.1 One-Shot Queries

An unsophisticated DNS client may simply send its DNS queries blindly to the 224.0.0.251 multicast address, without necessarily even being aware what a multicast address is.

Such an unsophisticated DNS client may not get ideal behavior. Such a client may simply take the first response it receives and fail to wait to see if there are more, but in many instances this may not be a serious problem. If a user types "http://cheshire.local." into their Web browser and gets to see the page they were hoping for, then the protocol has met the user's needs in this case.

6.2 One-Shot Queries, Accumulating Multiple Responses

A more sophisticated DNS client should understand that Multicast DNS is not exactly the same as unicast DNS, and should modify its behavior in some simple ways.

As described above, there are some cases, such as looking up the address associated with a unique host name, where a single response is sufficient, and moreover may be all that is expected. However, there are other DNS queries where more than one response is possible, and for these queries a more sophisticated Multicast DNS client should include the ability to wait for an appropriate period of time to collect multiple responses.

A naive DNS client retransmits its query only so long as it has received no response. A more sophisticated Multicast DNS client is aware that having received one response is not necessarily an indication that it might not receive others, and has the ability to retransmit its query an appropriate number of times at appropriate intervals until it is satisfied with the collection of responses it has gathered.

A more sophisticated Multicast DNS client that is retransmitting a query for which it has already received some responses, **MUST** implement Known Answer Suppression, as described below in [Section 7.1](#). This indicates to responders who have already replied that their responses have been received, and they don't need to send them again in response to this repeated query. In addition, the interval between the first two queries **SHOULD** be one second, and the intervals between subsequent queries **SHOULD** double.

6.3 Continuous Querying

In One-Shot Queries, with either a single or multiple responses, the underlying assumption is that the transaction begins when the application issues a query, and ends when all the desired responses have been received. There is another type of operation which is more akin to continuous monitoring.

Macintosh users are accustomed to opening the "Chooser" window, selecting a desired printer, and then closing the Chooser window. However, when the desired printer does not appear in the list, the user will typically leave the "Chooser" window open while they go and check to verify that the printer is plugged in, powered on, connected to the Ethernet, etc. While the user jiggles the wires, hits the Ethernet hub, and so forth, they keep an eye on the Chooser window, and when the printer name appears, they know they have fixed whatever the problem was. This can be a useful and intuitive troubleshooting

technique, but a user who goes home for the weekend leaving the
Chooser window open places a non-trivial burden on the network.

With continuous querying, multiple queries are sent over a long period of time, until the user terminates the operation. It is important that an IP network browser window displaying live information from the network using Multicast DNS, if left running for an extended period of time, should generate significantly less multicast traffic on the network than the old AppleTalk Chooser. Therefore, the interval between the first two queries SHOULD be one second, the intervals between subsequent queries SHOULD double, and the querier MUST implement Known Answer Suppression, as described below in [Section 7.1](#). When the interval between queries reaches or exceeds 60 minutes, a querier MAY cap the interval to a maximum of 60 minutes, and perform subsequent queries at a steady-state rate of one query per hour.

When a Multicast DNS Querier receives an answer, the answer contains a TTL value that indicates for how many seconds this answer is valid. After this interval has passed, the answer will no longer be valid and SHOULD be deleted from the cache. Before this time is reached, a Multicast DNS Querier with an ongoing interest in that record SHOULD re-issue its query to determine whether the record is still valid, and if so update its expiry time.

To perform this cache maintenance, a Multicast DNS Querier should plan to re-query for records after at least 50% of the record lifetime has elapsed. This document recommends the following specific strategy:

The Querier should plan to issue a query at 80% of the record lifetime, and then if no answer is received, at 85%, 90% and 95%. If an answer is received, then the remaining TTL is reset to the value given in the answer, and this process repeats for as long as the Multicast DNS Querier has an ongoing interest in the record. If after four queries no answer is received, the record is deleted when it reaches 100% of its lifetime.

To avoid the case where multiple Multicast DNS Queriers on a network all issue their queries simultaneously, a random variation of 2% of the record TTL should be added, so that queries are scheduled to be performed at 80-82%, 85-87%, 90-92% and then 95-97% of the TTL.

[6.4](#) Multiple Questions per Query

Multicast DNS allows a querier to place multiple questions in the Question Section of a single Multicast DNS query packet.

The semantics of a Multicast DNS query packet containing multiple questions is identical to a series of individual DNS query packets containing one question each. Combining multiple questions into a

single packet is purely an efficiency optimization, and has no other semantic significance.

A useful technique for adaptively combining multiple questions into a single query is to use a Nagle-style algorithm: When a client issues its first question, a Query packet is immediately built and sent, without delay. If the client then continues issuing a rapid series of questions they are held until either the first query receives at least one answer, or 100ms has passed, or there are enough questions to fill the Question Section of a Multicast DNS query packet. At this time, all the held questions are placed into a Multicast DNS query packet and sent.

6.5 Questions Requesting Unicast Responses

Sending Multicast DNS responses via multicast has the benefit that all the other hosts on the network get to see those responses, and can keep their caches up to date, and detect conflicting responses.

However, there are situations where all the other hosts on the network don't need to see every response. One example is a laptop computer waking from sleep. At that instant it is a brand new participant on a new network. Its Multicast DNS cache is empty, and it has no knowledge of its surroundings. It may have a significant number of queries that it wants answered right away to discover information about its new surroundings and present that information to the user. As a new participant on the network, it has no idea whether the exact same questions may have been asked and answered just seconds ago. In this case, triggering a large sudden flood of multicast responses may impose an unreasonable burden on the network. To avoid this, the Multicast DNS Querier SHOULD set the top bit in the class field of its DNS question(s), to indicate that it is willing to accept unicast responses instead of the usual multicast responses. These questions requesting unicast responses are referred to as "QU" questions, to distinguish them from the more usual questions requesting multicast responses ("QM" questions).

When retransmitting a question more than once, the 'unicast response' bit SHOULD be set only for the first question of the series. After the first question has received its responses, the querier should have a large known-answer list (see "Known Answer Suppression" below) so that subsequent queries should elicit few, if any, further responses. Reverting to multicast responses as soon as possible is important because of the benefits that multicast responses provide (see "Benefits of Multicast Responses" below).

When receiving a question with the 'unicast response' bit set, a responder SHOULD usually respond with a unicast packet directed back to the querier. If the responder has not multicast that record recently (within one quarter of its TTL), then the responder SHOULD instead multicast the response so as to keep all the peer caches up

to date, and to permit passive conflict detection.

Unicast replies are subject to all the same packet generation rules as multicast replies, including the cache flush bit (see [Section 11.3](#), "Announcements to Flush Outdated Cache Entries") and randomized delays to reduce network collisions (see [Section 8](#), "Responding").

6.6 Suppressing Initial Query

If a query is issued for which there already exist one or more records in the local cache, and those record(s) were received with the cache flush bit set (see [Section 11.3](#), "Announcements to Flush Outdated Cache Entries"), indicating that they form a unique RRSet, then the host SHOULD suppress its initial "QU" query, and proceed to issue a "QM" query. To avoid the situation where a group of hosts are synchronized by some external event and all perform the same query simultaneously, a host suppressing its initial "QU" query SHOULD impose a random delay from 500-1000ms before transmitting its first "QM" query for this question. This means that when the first host (selected randomly by this algorithm) transmits its "QM" query, all the other hosts that were about to transmit the same query can suppress their superfluous query, as described in "Duplicate Question Suppression" below.

7. Duplicate Suppression

A variety of techniques are used to reduce the amount of redundant traffic on the network.

7.1 Known Answer Suppression

When a Multicast DNS Querier sends a query to which it already knows some answers, it populates the Answer Section of the DNS message with those answers.

A Multicast DNS Responder SHOULD NOT answer a Multicast DNS Query if the answer it would give is already included in the Answer Section with an RR TTL at least half the correct value. If the RR TTL of the answer as given in the Answer Section is less than half of the true RR TTL as known by the Multicast DNS Responder, the responder MUST send an answer so as to update the Querier's cache before the record becomes in danger of expiration.

Because a Multicast DNS Responder will respond if the remaining TTL given in the known answer list is less than half the true TTL, it is superfluous for the Querier to include such records in the known answer list. Therefore a Multicast DNS Querier SHOULD NOT include records in the known answer list whose remaining TTL is less than half their original TTL. Doing so would simply consume space in the packet without achieving the goal of suppressing responses, and would therefore be a pointless waste of network bandwidth.

A Multicast DNS Querier MUST NOT cache resource records observed in the Known Answer Section of other Multicast DNS Queries. The Answer Section of Multicast DNS Queries is not authoritative. By placing information in the Answer Section of a Multicast DNS Query the

querier is stating that it **believes** the information to be true. It is not asserting that the information **is** true. Some of those records may have come from other hosts that are no longer on the network. Propagating that stale information to other Multicast DNS Queriers on the network would not be helpful.

7.2 Multi-Packet Known Answer Suppression

Sometimes a Multicast DNS Querier will already have too many answers to fit in the Known Answer Section of its query packets. In this case, it should issue a Multicast DNS Query containing a question and as many Known Answer records as will fit. It MUST then set the TC (Truncated) bit in the header before sending the Query. It MUST then immediately follow the packet with another query packet containing no questions, and as many more Known Answer records as will fit. If there are still too many records remaining to fit in the packet, it again sets the TC bit and continues until all the Known Answer records have been sent.

A Multicast DNS Responder seeing a Multicast DNS Query with the TC bit set defers its response for a time period randomly selected in the interval 400-500ms. This gives the Multicast DNS Querier time to send additional Known Answer packets before the Responder responds. If the Responder sees any of its answers listed in the Known Answer lists of subsequent packets from the querying host, it SHOULD delete that answer from the list of answers it is planning to give, provided that no other host on the network is also waiting to receive the same answer record.

Previous versions of this draft specified a delay of 20-120ms before answering queries with multi-packet Known Answer lists. However, operational experience showed that, while this works well on Ethernet, on very busy 802.11 networks, it is not uncommon to observe consecutively sent packets arriving separated by as much as 200-400ms.

7.3 Duplicate Question Suppression

If a host is planning to send a query, and it sees another host on the network send a query containing the same question, and the Known Answer Section of that query does not contain any records which this host would not also put in its own Known Answer Section, then this host should treat its own query as having been sent. When multiple clients on the network are querying for the same resource records, there is no need for them to all be repeatedly asking the same question.

7.4 Duplicate Answer Suppression

If a host is planning to send an answer, and it sees another host on the network send a response packet containing the same answer record, and the TTL in that record is not less than the TTL this host would have given, then this host should treat its own answer as having been

sent. When multiple responders on the network have the same data, there is no need for all of them to respond.

This feature is particularly useful when multiple Sleep Proxy Servers are deployed (see [Section 16](#), "Multicast DNS and Power Management"). In the future it is possible that every general-purpose OS (Mac, Windows, Linux, etc.) will implement Sleep Proxy Service as a matter of course. In this case there could be a large number of Sleep Proxy Servers on any given network, which is good for reliability and fault-tolerance, but would be bad for the network if every Sleep Proxy Server were to answer every query.

8. Responding

When a Multicast DNS Responder constructs and sends a Multicast DNS response packet, the Answer Section of that packet must contain only records for which that Responder is explicitly authoritative. These answers may be generated because the record answers a question received in a Multicast DNS query packet, or at certain other times that the responder determines that an unsolicited announcement is warranted. A Multicast DNS Responder **MUST NOT** place records from its cache, which have been learned from other responders on the network, in the Answer Section of outgoing response packets. Only an authoritative source for a given record is allowed to issue responses containing that record.

The determination of whether a given record answers a given question is done using the standard DNS rules: The record name must match the question name, the record rrtype must match the question qtype (unless the qtype is "ANY"), and the record rclass must match the question qclass (unless the qclass is "ANY").

A Multicast DNS Responder **MUST** only respond when it has a positive non-null response to send. Error responses must never be sent. The non-existence of any name in a Multicast DNS Domain is ascertained by the failure of any machine to respond to the Multicast DNS query, not by NXDOMAIN errors.

Multicast DNS Responses **MUST NOT** contain any questions in the Question Section. Any questions in the Question Section of a received Multicast DNS Response **MUST** be silently ignored. Multicast DNS Queriers receiving Multicast DNS Responses do not care what question elicited the response; they care only that the information in the response is true and accurate.

A Multicast DNS Responder on Ethernet [[IEEE802](#)] and similar shared multiple access networks **SHOULD** have the capability of delaying its responses by up to 500ms, as determined by the rules described below. If multiple Multicast DNS Responders were all to respond immediately to a particular query, a collision would be virtually guaranteed. By

imposing a small random delay, the number of collisions is dramatically reduced. On a full-sized Ethernet using the maximum cable lengths allowed and the maximum number of repeaters allowed, an Ethernet frame is vulnerable to collisions during the transmission of its first 256 bits. On 10Mb/s Ethernet, this equates to a vulnerable

time window of 25.6us. On higher-speed variants of Ethernet, the vulnerable time window is shorter.

In the case where a Multicast DNS Responder has good reason to believe that it will be the only responder on the link with a positive non-null response, it SHOULD NOT impose any random delay before responding, and SHOULD normally generate its response within at most 10ms. In particular, this applies to responding to probe queries. Since receiving a probe query gives a clear indication that some other Responder is planning to start using this name in the very near future, answering such probe queries to defend a unique record is a high priority and needs to be done immediately, without delay. A probe query can be distinguished from a normal query by the fact that a probe query contains a proposed record in the Authority Section which answers the question in the Question Section (for more details, see [Section 9.1](#), "Probing").

To generate immediate responses safely, it MUST have previously verified that the requested name, rrtype and rrclass in the DNS query are unique on this link. Responding immediately without delay is appropriate for things like looking up the address record for a particular host name, when the host name has been previously verified unique. Responding immediately without delay is **not** appropriate for things like looking up PTR records used for DNS Service Discovery [[DNS-SD](#)], where a large number of responses may be anticipated.

In any case where there may be multiple responses, such as queries where the answer is a member of a shared resource record set, each responder SHOULD delay its response by a random amount of time selected with uniform random distribution in the range 20-120ms.

In the case where the query has the TC (truncated) bit set, indicating that subsequent known answer packets will follow, responders SHOULD delay their responses by a random amount of time selected with uniform random distribution in the range 400-500ms, to allow enough time for all the known answer packets to arrive.

Except when a unicast reply has been explicitly requested via the "unicast reply" bit, Multicast DNS Responses MUST be sent to UDP port 5353 (the well-known port assigned to mDNS) on the 224.0.0.251 multicast address (or its IPv6 equivalent FF02::FB). Operating in a Zeroconf environment requires constant vigilance. Just because a name has been previously verified unique does not mean it will continue to be so indefinitely. By allowing all Multicast DNS Responders to constantly monitor their peers' responses, conflicts arising out of network topology changes can be promptly detected and resolved.

Sending all responses by multicast also facilitates opportunistic

caching by other hosts on the network.

To protect the network against excessive packet flooding due to software bugs or malicious attack, a Multicast DNS Responder MUST NOT multicast a given record on a given interface if it has previously

multicast that record on that interface within the last second. A legitimate client on the network should have seen the previous transmission and cached it. A client that did not receive and cache the previous transmission will retry its request and receive a subsequent response. Under no circumstances is there any legitimate reason for a Multicast DNS Responder to multicast a given record more than once per second on any given interface.

8.1 Legacy Unicast Responses

If the source UDP port in a received Multicast DNS Query is not port 5353, this indicates that the client originating the query is a simple client that does not fully implement all of Multicast DNS. In this case, the Multicast DNS Responder **MUST** send a UDP response directly back to the client, via unicast, to the query packet's source IP address and port. This unicast response **MUST** be a conventional unicast response as would be generated by a conventional unicast DNS server; for example, it **MUST** repeat the query ID and the question given in the query packet.

The resource record TTL given in a legacy unicast response **SHOULD NOT** be greater than ten seconds, even if the true TTL of the Multicast DNS resource record is higher. This is because Multicast DNS Responders that fully participate in the protocol use the cache coherency mechanisms described in [Section 13](#) to update and invalidate stale data. Were unicast responses sent to legacy clients to use the same high TTLs, these legacy clients, which do not implement these cache coherency mechanisms, could retain stale cached resource record data long after it is no longer valid.

Having sent this unicast response, if the Responder has not sent this record in any multicast response recently, it **SHOULD** schedule the record to be sent via multicast as well, to facilitate passive conflict detection. "Recently" in this context means "if the time since the record was last sent via multicast is less than one quarter of the record's TTL".

8.2 Multi-Question Queries

Multicast DNS Responders **MUST** correctly handle DNS query packets containing more than one question, by answering any or all of the questions to which they have answers. Any (non-defensive) answers generated in response to query packets containing more than one question **SHOULD** be randomly delayed in the range 20-120ms, or 400-500ms if the TC (truncated) bit is set, as described above. (Answers defending a name, in response to a probe for that name, are not subject to this delay rule and are still sent immediately.)

8.3 Response Aggregation

When possible, a responder SHOULD, for the sake of network efficiency, aggregate as many responses as possible into a single Multicast DNS response packet. For example, when a responder has several responses it plans to send, each delayed by a different interval, then earlier responses SHOULD be delayed by up to an additional 500ms if that will permit them to be aggregated with other responses scheduled to go out a little later.

9. Probing and Announcing on Startup

Typically a Multicast DNS Responder should have, at the very least, address records for all of its active interfaces. Creating and advertising an HINFO record on each interface as well can be useful to network administrators.

Whenever a Multicast DNS Responder starts up, wakes up from sleep, receives an indication of an Ethernet "Link Change" event, or has any other reason to believe that its network connectivity may have changed in some relevant way, it MUST perform the two startup steps below.

9.1 Probing

The first startup step is that for all those resource records that a Multicast DNS Responder desires to be unique on the local link, it MUST send a Multicast DNS Query asking for those resource records, to see if any of them are already in use. The primary example of this is its address record which maps its unique host name to its unique IP address. All Probe Queries SHOULD be done using the desired resource record name and query type T_ANY (255), to elicit answers for all types of records with that name. This allows a single question to be used in place of several questions, which is more efficient on the network. It also allows a host to verify exclusive ownership of a name, which is desirable in most cases. It would be confusing, for example, if one host owned the "A" record for "myhost.local.", but a different host owned the HINFO record for that name.

The ability to place more than one question in a Multicast DNS Query is useful here, because it can allow a host to use a single packet for all of its resource records instead of needing a separate packet for each. For example, a host can simultaneously probe for uniqueness of its "A" record and all its SRV records [[DNS-SD](#)] in the same query packet.

When ready to send its mDNS probe packet(s) the host should first

wait for a short random delay time, uniformly distributed in the range 0-250ms. This random delay is to guard against the case where a group of devices are powered on simultaneously, or a group of devices are connected to an Ethernet hub which is then powered on, or some

other external event happens that might cause a group of hosts to all send synchronized probes.

250ms after the first query the host should send a second, then 250ms after that a third. If, by 250ms after the third probe, no conflicting Multicast DNS responses have been received, the host may move to the next step, announcing. (Note that this is the one exception from the normal rule that there should be at least one second between repetitions of the same question, and the interval between subsequent repetitions should double.)

If any conflicting Multicast DNS responses are received, then the probing host **MUST** defer to the existing host, and **MUST** choose new names for some or all of its resource records as appropriate, to avoid conflict with pre-existing hosts on the network. In the case of a host probing using query type T_ANY as recommended above, any answer containing a record with that name, of any type, **MUST** be considered a conflicting response and handled accordingly.

If fifteen failures occur within any ten-second period, then the host **MUST** wait at least five seconds before each successive additional probe attempt. This is to help ensure that in the event of software bugs or other unanticipated problems, errant hosts do not flood the network with a continuous stream of multicast traffic. For very simple devices, a valid way to comply with this requirement is to always wait five seconds after any failed probe attempt.

If a responder knows by other means, with absolute certainty, that its unique resource record set name, rrtype and rrclass cannot already be in use by any other responder on the network, then it **MAY** skip the probing step for that resource record set. For example, when creating the reverse address mapping PTR records, the host can reasonably assume that no other host will be trying to create those same PTR records, since that would imply that the two hosts were trying to use the same IP address, and if that were the case, the two hosts would be suffering communication problems beyond the scope of what Multicast DNS is designed to solve.

9.2 Simultaneous Probe Tie-Breaking

The astute reader will observe that there is a race condition inherent in the previous description. If two hosts are probing for the same name simultaneously, neither will receive any response to the probe, and the hosts could incorrectly conclude that they may both proceed to use the name. To break this symmetry, each host populates the Authority Section of its queries with records giving the rdata that it would be proposing to use, should its probing be

successful. The Authority Section is being used here in a way analogous to the Update Section of a DNS Update packet [[RFC 2136](#)].

When a host that is probing for a record sees another host issue a query for the same record, it consults the Authority Section of that

query. If it finds any resource record there which answers the query, then it compares the data of that resource record with its own tentative data. The lexicographically later data wins. This means that if the host finds that its own data is lexicographically later, it simply ignores the other host's probe. If the host finds that its own data is lexicographically earlier, then it treats this exactly as if it had received a positive answer to its query, and concludes that it may not use the desired name.

The determination of 'lexicographically later' is performed by first comparing the record class, then the record type, then raw comparison of the binary content of the rdata without regard for meaning or structure. If the record classes differ, then the numerically greater class is considered 'lexicographically later'. Otherwise, if the record types differ, then the numerically greater type is considered 'lexicographically later'. If the rrtype and rrclass both match then the rdata is compared.

In the case of resource records containing rdata that is subject to name compression, the names **MUST** be uncompressed before comparison. (The details of how a particular name is compressed is an artifact of how and where the record is written into the DNS message; it is not an intrinsic property of the resource record itself.)

The bytes of the raw uncompressed rdata are compared in turn, interpreting the bytes as eight-bit UNSIGNED values, until a byte is found whose value is greater than that of its counterpart (in which case the rdata whose byte has the greater value is deemed lexicographically later) or one of the resource records runs out of rdata (in which case the resource record which still has remaining data first is deemed lexicographically later).

The following is an example of a conflict:

```
cheshire.local. A 169.254.99.200
cheshire.local. A 169.254.200.50
```

In this case 169.254.200.50 is lexicographically later (the third byte, with value 200, is greater than its counterpart with value 99), so it is deemed the winner.

Note that it is vital that the bytes are interpreted as UNSIGNED values, or the wrong outcome may result. In the example above, if the byte with value 200 had been incorrectly interpreted as a signed value then it would be interpreted as value -56, and the wrong address record would be deemed the winner.

[9.3](#) Announcing

The second startup step is that the Multicast DNS Responder MUST send a gratuitous Multicast DNS Response containing, in the Answer Section, all of its resource records (both shared records, and unique

records that have completed the probing step). If there are too many resource records to fit in a single packet, multiple packets should be used.

In the case of shared records (e.g. the PTR records used by DNS Service Discovery [[DNS-SD](#)]), the records are simply placed as-is into the Answer Section of the DNS Response.

In the case of records that have been verified to be unique in the previous step, they are placed into the Answer Section of the DNS Response with the most significant bit of the rrclass set to one. The most significant bit of the rrclass for a record in the Answer Section of a response packet is the mDNS "cache flush" bit and is discussed in more detail below in [Section 11.3](#) "Announcements to Flush Outdated Cache Entries".

The Multicast DNS Responder MUST send at least two gratuitous responses, one second apart. A Responder MAY send up to ten gratuitous Responses, provided that the interval between gratuitous responses doubles with every response sent.

A Multicast DNS Responder SHOULD NOT continue sending gratuitous Responses for longer than the TTL of the record. The purpose of announcing new records via gratuitous Responses is to ensure that peer caches are up to date. After a time interval equal to the TTL of the record has passed, it is very likely that old stale copies of that record in peer caches will have expired naturally, so subsequent announcements serve little purpose.

A Multicast DNS Responder MUST NOT send announcements in the absence of information that its network connectivity may have changed in some relevant way. In particular, a Multicast DNS Responder MUST NOT send regular periodic announcements as a matter of course.

Whenever a Multicast DNS Responder receives any Multicast DNS response (gratuitous or otherwise) containing a conflicting resource record, the conflict MUST be resolved as described below in "Conflict Resolution".

[9.4](#) Updating

At any time, if the rdata of any of a host's Multicast DNS records changes, the host MUST repeat the Announcing step described above to update neighboring caches. For example, if any of a host's IP addresses change, it MUST re-announce those address records.

In the case of shared records, a host MUST send a 'goodbye' announcement with TTL zero (see [Section 11.2](#) "Goodbye Packets") for the old rdata, to cause it to be deleted from peer caches,

before announcing the new rdata. In the case of unique records, a host SHOULD omit the 'goodbye' announcement, since the cache flush bit on the newly announced records will cause old rdata to be flushed from peer caches anyway.

A host may update the contents of any of its records at any time, though a host SHOULD NOT update records more frequently than ten times per minute. Frequent rapid updates impose a burden on the network. If a host has information to disseminate which changes more frequently than ten times per minute, then it may be more appropriate to design a protocol for that specific purpose.

10. Conflict Resolution

A conflict occurs when a Multicast DNS Responder has a unique record for which it is authoritative, and it receives, in the Answer Section of a Multicast DNS response another record with the same name, rrtype and rrclass, but inconsistent rdata. What may be considered inconsistent is context sensitive, except that resource records with identical rdata are never considered inconsistent, even if they originate from different hosts. This is to permit use of proxies and other fault-tolerance mechanisms that may cause more than one responder to be capable of issuing identical answers on the network.

A common example of a resource record type that is intended to be unique, not shared between hosts, is the address record that maps a host's name to its IP address. Should a host witness another host announce an address record with the same name but a different IP address, then that is considered inconsistent, and that address record is considered to be in conflict.

Whenever a Multicast DNS Responder receives any Multicast DNS response (gratuitous or otherwise) containing a conflicting resource record in the Answer Section, the Multicast DNS Responder MUST immediately reset its conflicted unique record to probing state, and go through the startup steps described above in [Section 9](#). "Probing and Announcing on Startup". The protocol used in the Probing phase will determine a winner and a loser, and the loser MUST cease using the name, and reconfigure.

It is very important that any host receiving a resource record that conflicts with one of its own MUST take action as described above. In the case of two hosts using the same host name, where one has been configured to require a unique host name and the other has not, the one that has not been configured to require a unique host name will not perceive any conflict, and will not take any action. By reverting to Probing state, the host that desires a unique host name will go through the necessary steps to ensure that a unique host is obtained.

The recommended course of action after probing and failing is as follows:

- o Programmatically change the resource record name in an attempt to

find a new name that is unique. This could be done by adding some further identifying information (e.g. the model name of the hardware) if it is not already present in the name, appending the digit "2" to the name, or incrementing a number at the end of the name if one is already present.

- o Probe again, and repeat until a unique name is found.
- o Record this newly chosen name in persistent storage so that the device will use the same name the next time it is power-cycled.
- o Display a message to the user or operator informing them of the name change. For example:

The name "Bob's Music" is in use by another iTunes music server on the network. Your music has been renamed to "Bob's Music (G4 Cube)". If you want to change this name, use [describe appropriate menu item or preference dialog].

How the user or operator is informed depends on context. A desktop computer with a screen might put up a dialog box. A headless server in the closet may write a message to a log file, or use whatever mechanism (email, SNMP trap, etc.) it uses to inform the administrator of other error conditions. On the other hand a headless server in the closet may not inform the user at all -- if the user cares, they will notice the name has changed, and connect to the server in the usual way (e.g. via Web Browser) to configure a new name.

The examples in this section focus on address records (i.e. host names), but the same considerations apply to all resource records where uniqueness (or maintenance of some other defined constraint) is desired.

11. Resource Record TTL Values and Cache Coherency

As a general rule, the recommended TTL value for Multicast DNS resource records with a host name as the resource record's name (e.g. A, AAAA, HINFO, etc.) or contained within the resource record's rdata (e.g. SRV, reverse mapping PTR record, etc.) is 120 seconds.

The recommended TTL value for other Multicast DNS resource records is 75 minutes.

A client with an active outstanding query will issue a query packet when one or more of the resource record(s) in its cache is (are) 80% of the way to expiry. If the TTL on those records is 75 minutes, this ongoing cache maintenance process yields a steady-state query rate of one query every 60 minutes.

Any distributed cache needs a cache coherency protocol. If Multicast DNS resource records follow the recommendation and have a TTL of 75 minutes, that means that stale data could persist in the system for a little over an hour. Making the default TTL significantly lower

would reduce the lifetime of stale data, but would produce too much extra traffic on the network. Various techniques are available to minimize the impact of such stale data.

11.1 Cooperating Multicast DNS Responders

If a Multicast DNS Responder ("A") observes some other Multicast DNS Responder ("B") send a Multicast DNS Response packet containing a resource record with the same name, rrtype and rrclass as one of A's resource records, but different rdata, then:

- o If A's resource record is intended to be a shared resource record, then this is no conflict, and no action is required.
- o If A's resource record is intended to be a member of a unique resource record set owned solely by that responder, then this is a conflict and MUST be handled as described in [Section 10](#) "Conflict Resolution".

If a Multicast DNS Responder ("A") observes some other Multicast DNS Responder ("B") send a Multicast DNS Response packet containing a resource record with the same name, rrtype and rrclass as one of A's resource records, and identical rdata, then:

- o If the TTL of B's resource record given in the packet is at least half the true TTL from A's point of view, then no action is required.
- o If the TTL of B's resource record given in the packet is less than half the true TTL from A's point of view, then A MUST mark its record to be announced via multicast. Clients receiving the record from B would use the TTL given by B, and hence may delete the record sooner than A expects. By sending its own multicast response correcting the TTL, A ensures that the record will be retained for the desired time.

These rules allow multiple Multicast DNS Responders to offer the same data on the network (perhaps for fault tolerance reasons) without conflicting with each other.

11.2 Goodbye Packets

In the case where a host knows that certain resource record data is about to become invalid (for example when the host is undergoing a clean shutdown) the host SHOULD send a gratuitous announcement mDNS response packet, giving the same resource record name, rrtype, rrclass and rdata, but an RR TTL of zero. This has the effect of updating the TTL stored in neighboring hosts' cache entries to zero, causing that cache entry to be promptly deleted.

Clients receiving a Multicast DNS Response with a TTL of zero SHOULD NOT immediately delete the record from the cache, but instead record

a TTL of 1 and then delete the record one second later. In the case of multiple Multicast DNS Responders on the network described in [Section 11.1](#) above, if one of the Responders shuts down and incorrectly sends goodbye packets for its records, it gives the other

cooperating Responders one second to send out their own response to "rescue" the records before they expire and are deleted.

Generally speaking, it is more important to send goodbye packets for shared records than unique records. A given shared record name (such as a PTR record used for DNS Service Discovery [[DNS-SD](#)]) by its nature often has many representatives from many different hosts, and tends to be the subject of long-lived ongoing queries. Those long-lived queries are often concerned not just about being informed when records appear, but also about being informed if those records vanish again. In contrast, a unique record set (such as an SRV record, or a host address record), by its nature, often has far fewer members than a shared record set, and is usually the subject of one-shot queries which simply retrieve the data and then cease querying once they have the answer they are seeking. Therefore, sending a goodbye packet for a unique record set is likely to offer less benefit, because it is likely at any given moment that no one has an active query running for that record set. One example where goodbye packets for SRV and address records are useful is when transferring control to a Sleep Proxy Server (see [Section 16](#), "Multicast DNS and Power Management").

11.3 Announcements to Flush Outdated Cache Entries

Whenever a host has a resource record with potentially new data (e.g. after rebooting, waking from sleep, connecting to a new network link, changing IP address, etc.), the host **MUST** send a series of gratuitous announcements to update cache entries in its neighbor hosts. In these gratuitous announcements, if the record is one that is intended to be unique, the host sets the most significant bit of the rrclass field of the resource record. This bit, the "cache flush" bit, tells neighboring hosts that this is not a shared record type. Instead of merging this new record additively into the cache in addition to any previous records with the same name, rrtype and rrclass, all old records with that name, type and class that were received more than one second ago are declared invalid, and marked to expire from the cache in one second.

The semantics of the cache flush bit are as follows: Normally when a resource record appears in the Answer Section of the DNS Response, it means, "This is an assertion that this information is true." When a resource record appears in the Answer Section of the DNS Response with the "cache flush" bit set, it means, "This is an assertion that this information is the truth and the whole truth, and anything you may have heard more than a second ago regarding records of this name/rrtype/rrclass is no longer valid".

To accommodate the case where the set of records from one host constituting a single unique RRSset is too large to fit in a single packet, only cache records that are more than one second old are flushed. This allows the announcing host to generate a quick burst of packets back-to-back on the wire containing all the members

of the RRSets. When receiving records with the "cache flush" bit set, all records older than one second are marked to be deleted one second in the future. One second after the end of the little packet burst, any records not represented within that packet burst will then be expired from all peer caches.

Any time a host sends a response packet containing some members of a unique RRSets, it SHOULD send the entire RRSets, preferably in a single packet, or if the entire RRSets will not fit in a single packet, in a quick burst of packets sent as close together as possible. The host SHOULD set the cache flush bit on all members of the unique RRSets. In the event that for some reason the host chooses not to send the entire unique RRSets in a single packet or a rapid packet burst, it MUST NOT set the cache flush bit on any of those records.

The reason for waiting one second before deleting stale records from the cache is to accommodate bridged networks. For example, a host's address record announcement on a wireless interface may be bridged onto a wired Ethernet, and cause that same host's Ethernet address records to be flushed from peer caches. The one-second delay gives the host the chance to see its own announcement arrive on the wired Ethernet, and immediately re-announce its Ethernet interface's address records so that both sets remain valid and live in peer caches.

These rules apply regardless of *why* the response packet is being generated. They apply to startup announcements as described in [Section 9.3](#), and to responses generated as a result of receiving query packets.

The "cache flush" bit is only set in records in the Answer Section of Multicast DNS responses sent to UDP port 5353. The "cache flush" bit MUST NOT be set in any resource records in a response packet sent in legacy unicast responses to UDP ports other than 5353.

The "cache flush" bit MUST NOT be set in any resource records in the known-answer list of any query packet.

The "cache flush" bit MUST NOT ever be set in any shared resource record. To do so would cause all the other shared versions of this resource record with different rdata from different Responders to be immediately deleted from all the caches on the network.

The "cache flush" bit does apply to questions listed in the Question Section of a Multicast DNS packet. The top bit of the rrclass field in questions is used for an entirely different purpose (see [Section 6.5](#), "Questions Requesting Unicast Responses").

Note that the "cache flush" bit is NOT part of the resource record

class. The "cache flush" bit is the most significant bit of the second 16-bit word of a resource record in the Answer Section of an mDNS packet (the field conventionally referred to as the rrclass field), and the actual resource record class is the least-significant

fifteen bits of this field. There is no mDNS resource record class 0x8001. The value 0x8001 in the rrclass field of a resource record in an mDNS response packet indicates a resource record with class 1, with the "cache flush" bit set. When receiving a resource record with the "cache flush" bit set, implementations should take care to mask off that bit before storing the resource record in memory.

11.4 Cache Flush on Topology change

If the hardware on a given host is able to indicate physical changes of connectivity, then when the hardware indicates such a change, the host should take this information into account in its mDNS cache management strategy. For example, a host may choose to immediately flush all cache records received on a particular interface when that cable is disconnected. Alternatively, a host may choose to adjust the remaining TTL on all those records to a few seconds so that if the cable is not reconnected quickly, those records will expire from the cache.

Likewise, when a host reboots, or wakes from sleep, or undergoes some other similar discontinuous state change, the cache management strategy should take that information into account.

11.5 Cache Flush on Failure Indication

Sometimes a cache record can be determined to be stale when a client attempts to use the rdata it contains, and finds that rdata to be incorrect.

For example, the rdata in an address record can be determined to be incorrect if attempts to contact that host fail, either because ARP/ND requests for that address go unanswered (for an address on a local subnet) or because a router returns an ICMP "Host Unreachable" error (for an address on a remote subnet).

The rdata in an SRV record can be determined to be incorrect if attempts to communicate with the indicated service at the host and port number indicated are not successful.

The rdata in a DNS-SD PTR record can be determined to be incorrect if attempts to look up the SRV record it references are not successful.

In any such case, the software implementing the mDNS resource record cache should provide a mechanism so that clients detecting stale rdata can inform the cache.

When the cache receives this hint that it should reconfirm some

record, it MUST issue two or more queries for the resource record in question. If no response is received in a reasonable amount of time, then, even though its TTL may indicate that it is not yet due to expire, that record SHOULD be promptly flushed from the cache.

The end result of this is that if a printer suffers a sudden power failure or other abrupt disconnection from the network, its name may continue to appear in DNS-SD browser lists displayed on users' screens. Eventually that entry will expire from the cache naturally, but if a user tries to access the printer before that happens, the failure to successfully contact the printer will trigger the more hasty demise of its cache entries. This is a sensible trade-off between good user-experience and good network efficiency. If we were to insist that printers should disappear from the printer list within 30 seconds of becoming unavailable, for all failure modes, the only way to achieve this would be for the client to poll the printer at least every 30 seconds, or for the printer to announce its presence at least every 30 seconds, both of which would be an unreasonable burden on most networks.

11.6 Passive Observation of Failures

A host observes the multicast queries issued by the other hosts on the network. One of the major benefits of also sending responses using multicast is that it allows all hosts to see the responses (or lack thereof) to those queries.

If a host sees queries, for which a record in its cache would be expected to be given as an answer in a multicast response, but no such answer is seen, then the host may take this as an indication that the record may no longer be valid.

After seeing two or more of these queries, and seeing no multicast response containing the expected answer within a reasonable amount of time, then even though its TTL may indicate that it is not yet due to expire, that record MAY be flushed from the cache. The host SHOULD NOT perform its own queries to re-confirm that the record is truly gone. If every host on a large network were to do this, it would cause a lot of unnecessary multicast traffic. If host A sends multicast queries that remain unanswered, then there is no reason to suppose that host B or any other host is likely to be any more successful.

The previous section, "Cache Flush on Failure Indication", describes a situation where a user trying to print discovers that the printer is no longer available. By implementing the passive observation described here, when one user fails to contact the printer, all hosts on the network observe that failure and update their caches accordingly.

12. Special Characteristics of Multicast DNS Domains

Unlike conventional DNS names, names that end in ".local.", "254.169.in-addr.arpa." or "0.8.e.f.ip6.arpa." have only local significance. Conventional DNS seeks to provide a single unified namespace, where a given DNS query yields the same answer no matter where on the planet it is performed or to which recursive DNS server the query is sent. (However, split views, firewalls, intranets and the like have somewhat interfered with this goal of DNS representing a single universal truth.) In contrast, each IP link has its own private ".local.", "254.169.in-addr.arpa." and "0.8.e.f.ip6.arpa." namespaces, and the answer to any query for a name within those domains depends on where that query is asked.

Multicast DNS Domains are not delegated from their parent domain via use of NS records. There are no NS records anywhere in Multicast DNS Domains. Instead, all Multicast DNS Domains are delegated to the IP addresses 224.0.0.251 and FF02::FB by virtue of the individual organizations producing DNS client software deciding how to handle those names. It would be extremely valuable for the industry if this special handling were ratified and recorded by IANA, since otherwise the special handling provided by each vendor is likely to be inconsistent.

The IPv4 name server for a Multicast DNS Domain is 224.0.0.251. The IPv6 name server for a Multicast DNS Domain is FF02::FB. These are multicast addresses; therefore they identify not a single host but a collection of hosts, working in cooperation to maintain some reasonable facsimile of a competently managed DNS zone. Conceptually a Multicast DNS Domain is a single DNS zone, however its server is implemented as a distributed process running on a cluster of loosely cooperating CPUs rather than as a single process running on a single CPU.

No delegation is performed within Multicast DNS Domains. Because the cluster of loosely coordinated CPUs is cooperating to administer a single zone, delegation is neither necessary nor desirable. Just because a particular host on the network may answer queries for a particular record type with the name "example.local." does not imply anything about whether that host will answer for the name "child.example.local.", or indeed for other record types with the name "example.local."

Multicast DNS Zones have no SOA record. A conventional DNS zone's SOA record contains information such as the email address of the zone administrator and the monotonically increasing serial number of the last zone modification. There is no single human administrator for any given Multicast DNS Zone, so there is no email address. Because

the hosts managing any given Multicast DNS Zone are only loosely coordinated, there is no readily available monotonically increasing serial number to determine whether or not the zone contents have changed. A host holding part of the shared zone could crash or be

disconnected from the network at any time without informing the other hosts. There is no reliable way to provide a zone serial number that would, whenever such a crash or disconnection occurred, immediately change to indicate that the contents of the shared zone had changed.

Zone transfers are not possible for any Multicast DNS Zone.

13. Multicast DNS for Service Discovery

This document does not describe using Multicast DNS for network browsing or service discovery. However, the mechanisms this document describes are compatible with (and support) the browsing and service discovery mechanisms proposed in "DNS-Based Service Discovery" [[DNS-SD](#)].

14. Enabling and Disabling Multicast DNS

The option to fail-over to Multicast DNS for names not ending in ".local." SHOULD be a user-configured option, and SHOULD be disabled by default because of the possible security issues related to unintended local resolution of apparently global names.

The option to lookup unqualified (relative) names by appending ".local." (or not) is controlled by whether ".local." appears (or not) in the client's DNS search list.

No special control is needed for enabling and disabling Multicast DNS for names explicitly ending with ".local." as entered by the user. The user doesn't need a way to disable Multicast DNS for names ending with ".local.", because if the user doesn't want to use Multicast DNS, they can achieve this by simply not using those names. If a user *does* enter a name ending in ".local.", then we can safely assume the user's intention was probably that it should work. Having user configuration options that can be (intentionally or unintentionally) set so that local names don't work is just one more way of frustrating the user's ability to perform the tasks they want, perpetuating the view that, "IP networking is too complicated to configure and too hard to use." This in turn perpetuates the continued use of protocols like AppleTalk. If we want to retire AppleTalk, NetBIOS, etc., we need to offer users equivalent IP functionality that they can rely on to, "always work, like AppleTalk." A little Multicast DNS traffic may be a burden on the network, but it is an insignificant burden compared to continued widespread use of AppleTalk.

15. Considerations for Multiple Interfaces

A host should defend its host name (FQDN) on all active interfaces on which it is answering Multicast DNS queries.

In the event of a name conflict on **any** interface, a host should configure a new host name, if it wishes to maintain uniqueness of its host name.

A host may choose to use the same name for all of its address records on all interfaces, or it may choose to manage its Multicast DNS host name(s) independently on each interface, potentially answering to different names on different interfaces.

When answering a Multicast DNS query, a multi-homed host with a link-local address (or addresses) should take care to ensure that any address going out in a Multicast DNS response is valid for use on the interface on which the response is going out.

Just as the same link-local IP address may validly be in use simultaneously on different links by different hosts, the same link-local host name may validly be in use simultaneously on different links, and this is not an error. A multi-homed host with connections to two different links may be able to communicate with two different hosts that are validly using the same name. While this kind of name duplication should be rare, it means that a host that wants to fully support this case needs network programming APIs that allow applications to specify on what interface to perform a link-local Multicast DNS query, and to discover on what interface a Multicast DNS response was received.

16. Multicast DNS and Power Management

Many modern network devices have the ability to go into a low-power mode where only a small part of the Ethernet hardware remains powered, and the device can be woken up by sending a specially formatted Ethernet frame which the device's power-management hardware recognizes.

To make use of this in conjunction with Multicast DNS, we propose a network power management service called Sleep Proxy Service. A device that wishes to enter low-power mode first uses DNS-SD to determine if Sleep Proxy Service is available on the local network. In some networks there may be more than one piece of hardware implementing Sleep Proxy Service, for fault-tolerance reasons.

If the device finds the network has Sleep Proxy Service, the device transmits two or more gratuitous mDNS announcements setting the TTL of its relevant resource records to zero, to delete them from neighboring caches. The relevant resource records include address records and SRV records, and other resource records as may apply to a particular device. The device then communicates all of its remaining active records, plus the names, rrtypes and rrclasses of the deleted records, to the Sleep Proxy Service(s), along with a copy of the specific "magic packet" required to wake the device up.

When a Sleep Proxy Service sees an mDNS query for one of the device's active records (e.g. a DNS-SD PTR record), it answers on behalf of the device without waking it up. When a Sleep Proxy Service sees an mDNS query for one of the device's deleted resource records, it deduces that some client on the network needs to make an active connection to the device, and sends the specified "magic packet" to wake the device up. The device then wakes up, reactivates its deleted resource records, and re-announces them to the network. The client waiting to connect sees the announcements, learns the current IP address and port number of the desired service on the device, and proceeds to connect to it.

The connecting client does not need to be aware of how Sleep Proxy Service works. Only devices that implement low power mode and wish to make use of Sleep Proxy Service need to be aware of how that protocol works.

The reason that a device using a Sleep Proxy Service should send more than one goodbye packet is to ensure deletion of the resource records from all peer caches. If resource records were to inadvertently remain in some peer caches, then those peers may not issue any query packets for those records when attempting to access the sleeping device, so the Sleep Proxy Service would not receive any queries for

the device's SRV and/or address records, and the necessary wake-up message would not be triggered.

The full specification of mDNS / DNS-SD Sleep Proxy Service is described in another document [not yet published].

17. Multicast DNS Character Set

Unicast DNS has been plagued by the lack of any support for non-US characters. Indeed, conventional DNS is usually limited to just letters, digits and hyphens, with no spaces or other punctuation. Attempts to remedy this for unicast DNS have been badly constrained by the need to accommodate old buggy legacy DNS implementations. In reality, the DNS specification actually imposes no limits on what characters may be used in names, and good DNS implementations handle any arbitrary eight-bit data without trouble. However, the old rules for ARPANET host names back in the 1980s required names to be just letters, digits, and hyphens [[RFC 1034](#)], and since the predominant use of DNS is to store host address records, many have assumed that the DNS protocol itself suffers from the same limitation. It would be more accurate to say that certain bad implementations may not handle eight-bit data correctly, not that the protocol doesn't support it.

Multicast DNS is a new protocol and doesn't (yet) have old buggy legacy implementations to constrain the design choices. Accordingly, it adopts the simple obvious elegant solution: all names in Multicast DNS are encoded using precomposed UTF-8 [[RFC 3629](#)]. The characters SHOULD conform to Unicode Normalization Form C (NFC): Use precomposed characters instead of combining sequences where possible, e.g. use U+00C4 ("Latin capital letter A with diaeresis") instead of U+0041 U+0308 ("Latin capital letter A", "combining diaeresis"). Some users of 16-bit Unicode have taken to stuffing a "zero-width non-breaking space" character (U+FEFF) at the start of each UTF-16 file, as a hint to identify whether the data is big-endian or little-endian, and calling it a "Byte Order Mark" (BOM). Since there is only one possible byte order for UTF-8 data, a BOM is neither necessary nor permitted. Multicast DNS names MUST NOT contain a "Byte Order Mark". Any occurrence of the Unicode character U+FEFF in a Multicast DNS name MUST be interpreted as a zero-width non-breaking space.

For names that are restricted to letters, digits and hyphens, the UTF-8 encoding is identical to the US-ASCII encoding, so this is entirely compatible with existing host names. For characters outside the US-ASCII range, UTF-8 encoding is used.

Multicast DNS implementations MUST NOT use any other encodings apart from precomposed UTF-8 (US-ASCII being considered a compatible subset of UTF-8).

This point bears repeating: After many years of debate, as a result of the need to accommodate certain DNS implementations that apparently couldn't handle any character that's not a letter, digit or hyphen (and apparently never will be updated to remedy this limitation) the unicast DNS community settled on an extremely baroque

encoding called "Punycode" [[RFC 3492](#)]. Punycode is a remarkably ingenious encoding solution, but it is complicated, hard to understand, and hard to implement, using sophisticated techniques including insertion unsort coding, generalized variable-length integers, and bias adaptation. The resulting encoding is remarkably

compact given the constraints, but it's still not as good as simple straightforward UTF-8, and it's hard even to predict whether a given input string will encode to a Punycode string that fits within DNS's 63-byte limit, except by simply trying the encoding and seeing whether it fits. Indeed, the encoded size depends not only on the input characters, but on the order they appear, so the same set of characters may or may not encode to a legal Punycode string that fits within DNS's 63-byte limit, depending on the order the characters appear. This is extremely hard to present in a user interface that explains to users why one name is allowed, but another name containing the exact same characters is not. Neither Punycode nor any other of the "Ascii Compatible Encodings" proposed for Unicast DNS may be used in Multicast DNS packets. Any text being represented internally in some other representation MUST be converted to canonical precomposed UTF-8 before being placed in any Multicast DNS packet.

The simple rules for case-insensitivity in Unicast DNS also apply in Multicast DNS; that is to say, in name comparisons, the lower-case letters "a" to "z" (0x61 to 0x7A) match their upper-case equivalents "A" to "Z" (0x41 to 0x5A). Hence, if a client issues a query for an address record with the name "cheshire.local", then a responder having an address record with the name "Cheshire.local" should issue a response. No other automatic equivalences should be assumed. In particular all UTF-8 multi-byte characters (codes 0x80 and higher) are compared by simple binary comparison of the raw byte values.

No other automatic character equivalence is defined in Multicast DNS. For example, accented characters are not defined to be automatically equivalent to their unaccented counterparts. Where automatic equivalences are desired, this may be achieved through the use of programmatically-generated CNAME records. For example, if a responder has an address record for an accented name Y, and a client issues a query for a name X, where X is the same as Y with all the accents removed, then the responder may issue a response containing two resource records: A CNAME record "X CNAME Y", asserting that the requested name X (unaccented) is an alias for the true (accented) name Y, followed by the address record for Y.

18. Multicast DNS Message Size

[RFC 1035](#) restricts DNS Messages carried by UDP to no more than 512 bytes (not counting the IP or UDP headers). For UDP packets carried over the wide-area Internet in 1987, this was appropriate. For link-local multicast packets on today's networks, there is no reason to retain this restriction. Given that the packets are by definition link-local, there are no Path MTU issues to consider.

Multicast DNS Messages carried by UDP may be up to the IP MTU of the physical interface, less the space required for the IP header (20 bytes for IPv4; 40 bytes for IPv6) and the UDP header (8 bytes).

In the case of a single mDNS Resource Record which is too large to fit in a single MTU-sized multicast response packet, a Multicast DNS Responder SHOULD send the Resource Record alone, in a single IP datagram, sent using multiple IP fragments. Resource Records this large SHOULD be avoided, except in the very rare cases where they really are the appropriate solution to the problem at hand. Implementers should be aware that many simple devices do not re-assemble fragmented IP datagrams, so large Resource Records SHOULD NOT be used except in specialized cases where the implementer knows that all receivers implement reassembly.

A Multicast DNS packet larger than the interface MTU, which is sent using fragments, MUST NOT contain more than one Resource Record.

Even when fragmentation is used, a Multicast DNS packet, including IP and UDP headers, MUST NOT exceed 9000 bytes.

19. Multicast DNS Message Format

This section describes specific restrictions on the allowable values for the header fields of a Multicast DNS message.

19.1. ID (Query Identifier)

Multicast DNS clients SHOULD listen for gratuitous responses issued by hosts booting up (or waking up from sleep or otherwise joining the network). Since these gratuitous responses may contain a useful answer to a question for which the client is currently awaiting an answer, Multicast DNS clients SHOULD examine all received Multicast DNS response messages for useful answers, without regard to the contents of the ID field or the Question Section. In Multicast DNS, knowing which particular query message (if any) is responsible for eliciting a particular response message is less interesting than knowing whether the response message contains useful information.

Multicast DNS clients MAY cache any or all Multicast DNS response messages they receive, for possible future use, provided of course that normal TTL aging is performed on these cached resource records.

In multicast query messages, the Query ID SHOULD be set to zero on transmission.

In multicast responses, including gratuitous multicast responses, the Query ID MUST be set to zero on transmission, and MUST be ignored on reception.

In unicast response messages generated specifically in response to a particular (unicast or multicast) query, the Query ID MUST match the

ID from the query message.

Expires 7th December 2005

Cheshire & Krochmal

[Page 35]

19.2. QR (Query/Response) Bit

In query messages, MUST be zero.

In response messages, MUST be one.

19.3. OPCODE

In both multicast query and multicast response messages, MUST be zero (only standard queries are currently supported over multicast, unless other queries are allowed by future IETF Standards Action).

19.4. AA (Authoritative Answer) Bit

In query messages, the Authoritative Answer bit MUST be zero on transmission, and MUST be ignored on reception.

In response messages for Multicast Domains, the Authoritative Answer bit MUST be set to one (not setting this bit implies there's some other place where "better" information may be found) and MUST be ignored on reception.

19.5. TC (Truncated) Bit

In query messages, if the TC bit is set, it means that additional Known Answer records may be following shortly. A responder MAY choose to record this fact, and wait for those additional Known Answer records, before deciding whether to respond. If the TC bit is clear, it means that the querying host has no additional Known Answers.

In multicast response messages, the TC bit MUST be zero on transmission, and MUST be ignored on reception.

In legacy unicast response messages, the TC bit has the same meaning as in conventional unicast DNS: it means that the response was too large to fit in a single packet, so the client SHOULD re-issue its query using TCP in order to receive the larger response.

19.6. RD (Recursion Desired) Bit

In both multicast query and multicast response messages, the Recursion Desired bit SHOULD be zero on transmission, and MUST be ignored on reception.

19.7. RA (Recursion Available) Bit

In both multicast query and multicast response messages, the Recursion Available bit MUST be zero on transmission, and MUST be ignored on reception.

19.8. Z (Zero) Bit

In both query and response messages, the Zero bit **MUST** be zero on transmission, and **MUST** be ignored on reception.

19.9. AD (Authentic Data) Bit [[RFC 2535](#)]

In query messages the Authentic Data bit **MUST** be zero on transmission, and **MUST** be ignored on reception.

In response messages, the Authentic Data bit **MAY** be set. Resolvers receiving response messages with the AD bit set **MUST NOT** trust the AD bit unless they trust the source of the message and either have a secure path to it or use DNS transaction security.

19.10. CD (Checking Disabled) Bit [[RFC 2535](#)]

In query messages, a resolver willing to do cryptography **SHOULD** set the Checking Disabled bit to permit it to impose its own policies.

In response messages, the Checking Disabled bit **MUST** be zero on transmission, and **MUST** be ignored on reception.

19.11. RCODE (Response Code)

In both multicast query and multicast response messages, the Response Code **MUST** be zero on transmission. Multicast DNS messages received with non-zero Response Codes **MUST** be silently ignored.

19.12. Repurposing of top bit of qclass in Question Section

In the Question Section of a Multicast DNS Query, the top bit of the qclass field is used to indicate that unicast responses are preferred for this particular question.

19.13. Repurposing of top bit of rrclass in Answer Section

In the Answer Section of a Multicast DNS Response, the top bit of the rrclass field is used to indicate that the record is a member of a unique RRSet, and the entire RRSet has been sent together (in the same packet, or in consecutive packets if there are too many records to fit in a single packet).

20. Choice of UDP Port Number

Arguments were made for and against using Multicast on UDP port 53. The final decision was to use UDP port 5353. Some of the arguments for and against are given below.

20.1 Arguments for using UDP port 53:

- * This is "just DNS", so it should be the same port.
- * There is less work to be done updating old clients to do simple mDNS queries. Only the destination address need be changed. In some cases, this can be achieved without any code changes, just by adding the address 224.0.0.251 to a configuration file.

20.2 Arguments for using a different port (UDP port 5353):

- * This is not "just DNS". This is a DNS-like protocol, but different.
- * Changing client code to use a different port number is not hard.
- * Using the same port number makes it hard to run an mDNS Responder and a conventional unicast DNS server on the same machine. If a conventional unicast DNS server wishes to implement mDNS as well, it can still do that, by opening two sockets. Having two different port numbers is important to allow this flexibility.
- * Some VPN software hijacks all outgoing traffic to port 53 and redirects it to a special DNS server set up to serve those VPN clients while they are connected to the corporate network. It is questionable whether this is the right thing to do, but it is common, and redirecting link-local multicast DNS packets to a remote server rarely produces any useful results. It does mean, for example, that the user becomes unable to access their local network printer sitting on their desk right next to their computer. Using a different UDP port eliminates this particular problem.
- * On many operating systems, unprivileged clients may not send or receive packets on low-numbered ports. This means that any client sending or receiving mDNS packets on port 53 would have to run as "root", which is an undesirable security risk. Using a higher-numbered UDP port eliminates this particular problem.

Continuing the previous point, since using an unprivileged port allows normal user-level code to bind, a given machine may have more than one such user-level application running at a time. Because of this, any code binding to UDP port 5353 MUST use the SO_REUSEPORT

option, so as to be a good citizen and not block other clients on the machine from also binding to that port.

21. Summary of Differences Between Multicast DNS and Unicast DNS

The value of Multicast DNS is that it shares, as much as possible, the familiar APIs, naming syntax, resource record types, etc., of Unicast DNS. There are of course necessary differences by virtue of it using Multicast, and by virtue of it operating in a community of cooperating peers, rather than a precisely defined authoritarian hierarchy controlled by a strict chain of formal delegations from the top. These differences are listed below:

Multicast DNS...

- * uses multicast
- * uses UDP port 5353 instead of port 53
- * operates in well-defined parts of the DNS namespace
- * uses UTF-8, and only UTF-8, to encode resource record names
- * defines a clear limit on the maximum legal domain name (255 bytes)
- * allows larger UDP packets
- * allows more than one question in a query packet
- * uses the Answer Section of a query to list Known Answers
- * uses the TC bit in a query to indicate additional Known Answers
- * uses the Authority Section of a query for probe tie-breaking
- * ignores the Query ID field (except for generating legacy responses)
- * doesn't require the question to be repeated in the response packet
- * uses gratuitous responses to announce new records to the peer group
- * defines a "unicast response" bit in the rrclass of query questions
- * defines a "cache flush" bit in the rrclass of response answers
- * uses DNS TTL 0 to indicate that a record has been deleted
- * monitors queries to perform Duplicate Question Suppression
- * monitors responses to perform Duplicate Answer Suppression...
- * ... and Ongoing Conflict Detection
- * ... and Opportunistic Caching

22. Benefits of Multicast Responses

Some people have argued that sending responses via multicast is inefficient on the network. In fact using multicast responses results in a net lowering of overall multicast traffic, for a variety of reasons, in addition to other benefits.

- * One multicast response can update the cache on all machines on the network. If another machine later wants to issue the same query, it already has the answer in its cache, so it may not need to even transmit that multicast query on the network at all.
- * When more than one machine has the same ongoing long-lived query running, every machine does not have to transmit its own independent query. When one machine transmits a query, all the other hosts see the answers, so they can suppress their own queries.
- * When a host sees a multicast query, but does not see the corresponding multicast response, it can use this information to promptly delete stale data from its cache. To achieve the same level of user-interface quality and responsiveness without multicast responses would require lower cache lifetimes and more frequent network polling, resulting in a significantly higher packet rate.
- * Multicast responses allow passive conflict detection. Without this ability, some other conflict detection mechanism would be needed, imposing its own additional burden on the network.
- * When using delayed responses to reduce network collisions, clients need to maintain a list recording to whom each answer should be sent. The option of multicast responses allows clients with limited storage, which cannot store an arbitrarily long list of response addresses, to choose to fail-over to a single multicast response in place of multiple unicast responses, when appropriate.
- * In the case of overlayed subnets, multicast responses allow a receiver to know with certainty that a response originated on the local link, even when its source address may apparently suggest otherwise.
- * Link-local multicast transcends virtually every conceivable network misconfiguration. Even if you have a collection of devices where every device's IP address, subnet mask, default gateway, and DNS server address are all wrong, packets sent by any of those devices addressed to a link-local multicast destination address will still be delivered to all peers on the local link. This can be extremely helpful when diagnosing and rectifying network problems, since it facilitates a direct communication channel between client and

server that works without reliance on ARP, IP routing tables, etc. Being able to discover what IP address a device has (or thinks it has) is frequently a very valuable first step in diagnosing why it unable to communicate on the local network.

23. IPv6 Considerations

An IPv4-only host and an IPv6-only host behave as "ships that pass in the night". Even if they are on the same Ethernet, neither is aware of the other's traffic. For this reason, each physical link may have **two** unrelated ".local." zones, one for IPv4 and one for IPv6. Since for practical purposes, a group of IPv4-only hosts and a group of IPv6-only hosts on the same Ethernet act as if they were on two entirely separate Ethernet segments, it is unsurprising that their use of the ".local." zone should occur exactly as it would if they really were on two entirely separate Ethernet segments.

A dual-stack (v4/v6) host can participate in both ".local." zones, and should register its name(s) and perform its lookups both using IPv4 and IPv6. This enables it to reach, and be reached by, both IPv4-only and IPv6-only hosts. In effect this acts like a multi-homed host, with one connection to the logical "IPv4 Ethernet segment", and a connection to the logical "IPv6 Ethernet segment".

23.1 IPv6 Multicast Addresses by Hashing

Some discovery protocols use a range of multicast addresses, and determine the address to be used by a hash function of the name being sought. Queries are sent via multicast to the address as indicated by the hash function, and responses are returned to the querier via unicast. Particularly in IPv6, where multicast addresses are extremely plentiful, this approach is frequently advocated.

There are some problems with this:

- * When a host has a large number of records with different names, the host may have to join a large number of multicast groups. This can place undue burden on the Ethernet hardware, which typically supports a limited number of multicast addresses efficiently. When this number is exceeded, the Ethernet hardware may have to resort to receiving all multicasts and passing them up to the host software for filtering, thereby defeating the point of using a multicast address range in the first place.
- * Multiple questions cannot be placed in one packet if they don't all hash to the same multicast address.
- * Duplicate Question Suppression doesn't work if queriers are not seeing each other's queries.
- * Duplicate Answer Suppression doesn't work if responders are not seeing each other's responses.

* Opportunistic Caching doesn't work.

* Ongoing Conflict Detection doesn't work.

24. Security Considerations

The algorithm for detecting and resolving name conflicts is, by its very nature, an algorithm that assumes cooperating participants. Its purpose is to allow a group of hosts to arrive at a mutually disjoint set of host names and other DNS resource record names, in the absence of any central authority to coordinate this or mediate disputes. In the absence of any higher authority to resolve disputes, the only alternative is that the participants must work together cooperatively to arrive at a resolution.

In an environment where the participants are mutually antagonistic and unwilling to cooperate, other mechanisms are appropriate, like manually administered DNS.

In an environment where there is a group of cooperating participants, but there may be other antagonistic participants on the same physical link, the cooperating participants need to use IPSEC signatures and/or DNSSEC [[RFC 2535](#)] signatures so that they can distinguish mDNS messages from trusted participants (which they process as usual) from mDNS messages from untrusted participants (which they silently discard).

When DNS queries for *global* DNS names are sent to the mDNS multicast address (during network outages which disrupt communication with the greater Internet) it is *especially* important to use DNSSEC, because the user may have the impression that he or she is communicating with some authentic host, when in fact he or she is really communicating with some local host that is merely masquerading as that name. This is less critical for names ending with ".local.", because the user should be aware that those names have only local significance and no global authority is implied.

Most computer users neglect to type the trailing dot at the end of a fully qualified domain name, making it a relative domain name (e.g. "www.example.com"). In the event of network outage, attempts to positively resolve the name as entered will fail, resulting in application of the search list, including ".local.", if present. A malicious host could masquerade as "www.example.com" by answering the resulting Multicast DNS query for "www.example.com.local." To avoid this, a host MUST NOT append the search suffix ".local.", if present, to any relative (partially qualified) domain name containing two or more labels. Appending ".local." to single-label relative domain names is acceptable, since the user should have no expectation that a single-label domain name will resolve as-is.

25. IANA Considerations

IANA has allocated the IPv4 link-local multicast address 224.0.0.251 for the use described in this document.

IANA has allocated the IPv6 multicast address set FF0X::FB for the use described in this document. Only address FF02::FB (Link-Local Scope) is currently in use by deployed software, but it is possible that in future implementers may experiment with Multicast DNS using larger-scoped addresses, such as FF05::FB (Site-Local Scope).

When this document is published, IANA should designate a list of domains which are deemed to have only link-local significance, as described in [Section 12](#) of this document ("Special Characteristics of Multicast DNS Domains").

The re-use of the top bit of the rrclass field in the Question and Answer Sections means that Multicast DNS can only carry DNS records with classes in the range 0-32767. Classes in the range 32768 to 65535 are incompatible with Multicast DNS. However, since to-date only three DNS classes have been assigned by IANA (1, 3 and 4), and only one (1, "Internet") is actually in widespread use, this limitation is likely to remain a purely theoretical one.

No other IANA services are required by this document.

26. Acknowledgments

The concepts described in this document have been explored, developed and implemented with help from Freek Dijkstra, Erik Guttman, Paul Vixie, Bill Woodcock, and others.

Special thanks go to Bob Bradley, Josh Graessley, Scott Herscher, Roger Pantos and Kiren Sekar for their significant contributions.

27. Copyright Notice

Copyright (C) The Internet Society (2005).

This document is subject to the rights, licenses and restrictions contained in [BCP 78](#), and except as set forth therein, the authors retain all their rights. For the purposes of this document, the term "[BCP 78](#)" refers exclusively to [RFC 3978](#), "IETF Rights in Contributions", published March 2005.

This document and the information contained herein are provided on an "AS IS" basis and THE CONTRIBUTOR, THE ORGANIZATION HE/SHE REPRESENTS OR IS SPONSORED BY (IF ANY), THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIM ALL WARRANTIES, EXPRESS OR IMPLIED,

INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE
INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED
WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Expires 7th December 2005

Cheshire & Krochmal

[Page 43]

28. Normative References

- [RFC 1034] Mockapetris, P., "Domain Names - Concepts and Facilities", STD 13, [RFC 1034](#), November 1987.
- [RFC 1035] Mockapetris, P., "Domain Names - Implementation and Specifications", STD 13, [RFC 1035](#), November 1987.
- [RFC 2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [RFC 2119](#), March 1997.
- [RFC 3629] Yergeau, F., "UTF-8, a transformation format of ISO 10646", [RFC 3629](#), November 2003.

29. Informative References

- [dotlocal] <<http://www.dotlocal.org/>>
- [djbd1] <<http://cr.yp.to/djbdns/dot-local.html>>
- [DNS-SD] Cheshire, S., and M. Krochmal, "DNS-Based Service Discovery", Internet-Draft (work in progress), [draft-cheshire-dnsext-dns-sd-03.txt](#), June 2005.
- [IEEE802] IEEE Standards for Local and Metropolitan Area Networks: Overview and Architecture.
Institute of Electrical and Electronic Engineers,
IEEE Standard 802, 1990.
- [NBP] Cheshire, S., and M. Krochmal,
"Requirements for a Protocol to Replace AppleTalk NBP",
Internet-Draft (work in progress),
[draft-cheshire-dnsext-nbp-04.txt](#), June 2005.
- [RFC 2136] Vixie, P., et al., "Dynamic Updates in the Domain Name System (DNS UPDATE)", [RFC 2136](#), April 1997.
- [RFC 2462] S. Thomson and T. Narten, "IPv6 Stateless Address Autoconfiguration", [RFC 2462](#), December 1998.
- [RFC 2535] Eastlake, D., "Domain Name System Security Extensions", [RFC 2535](#), March 1999.
- [RFC 3492] Costello, A., "Punycode: A Bootstring encoding of Unicode for use with Internationalized Domain Names in Applications (IDNA)", [RFC 3492](#), March 2003.
- [RFC 3927] Cheshire, S., B. Aboba, and E. Guttman,
"Dynamic Configuration of IPv4 Link-Local Addresses",

[RFC 3927](#), May 2005.

[ZC] Williams, A., "Requirements for Automatic Configuration of IP Hosts", Internet-Draft (work in progress), [draft-ietf-zeroconf-reqts-12.txt](#), September 2002.

Expires 7th December 2005

Cheshire & Krochmal

[Page 44]

30. Authors' Addresses

Stuart Cheshire
Apple Computer, Inc.
1 Infinite Loop
Cupertino
California 95014
USA

Phone: +1 408 974 3207
EMail: rfc [at] stuartcheshire [dot] org

Marc Krochmal
Apple Computer, Inc.
1 Infinite Loop
Cupertino
California 95014
USA

Phone: +1 408 974 4368
EMail: marc [at] apple [dot] com

