

# Solaris IP Duplicate Address Detection

James Carlson  
Sun Microsystems, Inc.

version 1.1 2005/09/12

## Abstract

This document describes the current design of the Solaris IP Duplicate Address Detection (DAD) logic, its problems, and a new design that fixes these problems.

## 1 Background

If two nodes on a broadcast-type network share the same IP (v4 or v6) address, then chaos erupts. In general, both nodes are harmed in the process, as the other nodes on the network become confused about which node is the “real” instance of that address. The symptoms range from subtle (some applications may still work, while other mysteriously fail) to hard-to-notice (obscure messages are tucked away in system logs, safe from prying eyes).

Diagnosing the root cause of these application failures is often complex, and requires the use of network monitors and a deep understanding of the protocols. Users often mistake the symptoms as wider system or application failures, leading to ineffectual “cures,” such as replacing cables or other hardware, and frustration.

Duplicate IP addresses can result from simple administrative misconfiguration in static addressing, or the failure of protocols that are intended to provide address assignment, such as might result from the accidental deployment of multiple DHCP servers on a single subnet.

Duplicate address detection (“DAD”) is a means to detect this situation. It is necessarily an imperfect process, as any events that might disrupt network connectivity (e.g., a temporary switch failure that causes network partitioning) will make it impossible to detect all error cases, but the goal is to detect enough of these cases to be useful in common deployments.

Once a duplicate address has been detected, it’s necessary to resolve the problem. Resolution typically consists of disabling all but (at most) one of the network interfaces involved, and notifying administrators on all affected systems. Where possible, the problem is avoided by detecting the configuration error before the network interface is enabled.

## 2 Related Projects

Duplicate address detection is a small part of Apple’s “Rendezvous” technology. It allows for link-local addressing in IPv4 for ad-hoc networking. (The bulk of Rendezvous builds on DNS.)

It’s also related to the IETF’s “Detecting Network Attachment” (DNA) effort, which is aimed primarily at mobile systems. Like DAD, DNA uses specially-crafted ARP messages to do its work.

The Surya project is merging ARP into IP. This will likely render parts of the ARP/IP signaling used in the initial DAD design (described in this document) obsolete. Currently, scheduling appears to have Surya arriving in Solaris Nevada much later than this project, and thus the design described here will hold. If this changes, and Surya integrates before this project, a new signaling mechanism will be adopted once Surya bits are available in Nevada.

This project does not itself implement the IPv4 Link-Local Addressing (LLA) mechanism from RFC 3927 [2]. It implements only the duplicate detection and address defense mechanisms. A separate project will implement LLA.

## 3 Current Status and Problems

Solaris today (prior to this proposed project) has at best uneven support for DAD. In IPv4, duplicate addresses are checked only if the interface is

configured via DHCP, as that protocol, according to the RFC, explicitly requires DAD. Manual configuration (e.g., `/etc/hostname.*` and `zonecfg`) and configuration by other means (e.g., RARP) simply bypass this check. Unfortunately, this is all too common, as too many misguided customers associate DHCP with “dynamic” addressing and thus don’t use it on systems classed as “servers.”

For non-point-to-point interfaces, Solaris IPv4 uses an older mechanism called “gratuitous ARP.” When the network administrator marks a network interface as “up,” the ip driver sends an `AR_ENTRY_ADD` message to the arp module, and that causes arp to broadcast multiple (default of 5; tuned via `ndd arp_publish_count`) `ares_op$REQUEST` messages at intervals of two seconds (tuned by `ndd arp_publish_interval`) with `ar$spa` and `ar$tpa` set to the local IP address, `ar$sha` set to the local hardware address, and `ar$tha` set to the broadcast address. This gratuitous ARP process has two effects:

1. It causes other ARP caches on the network to be updated immediately. If there was an existing user of the address, all of the packets that should go to him are now redirected to the newly-configured system. (And any running applications are immediately damaged, as the new system doesn’t share the application state.)
2. It causes the previous user of the address to emit this message into his system log:

```
IP: Hardware address 'xx:xx:xx:xx:xx:xx' trying to
be our address yyy.yyy.yyy.yyy!
```

He otherwise takes no action to fix the problem. The previous user, now a victim, backs off and sulks.

The new user (the one who accidentally configured the duplicate address) isn’t even necessarily notified that anything has gone amiss. He has no way of knowing that he’s even made a mistake, let alone that he needs to do something to correct it. As the user who made the mistake may not have administrative access to the victim system, and vice versa, neither can likely fix the problem.

In IPv6, Solaris support for DAD is present but incomplete. The process is handled within `ifconfig` and other utilities. Although helpful

in that it moves the control path out to user space where it ordinarily belongs, with the data path in the kernel, this design has several unfortunate side-effects:

1. If the interface transitions (e.g., is unplugged from one network and plugged into another network), the DAD probing process is not restarted. Not only is this an operational hazard, it's also a violation of RFC 2462 [5]. (See sections 5.3 and 5.4.)
2. When we implement a DHCPv6 client on Solaris, we'll need to duplicate the DAD functionality currently in `ifconfig` and `in.ndpd` yet again, as DHCPv6 also needs to run DAD. In fact, all projects that offer a means of configuring IPv6 interfaces (e.g., SNMP) need to have an embedded DAD implementation in the current design.
3. Non-Sun code running on Solaris that configures interfaces must also fend for its own in terms of duplicate address detection. The system provides no common support. It's quite likely that most such applications do not deal with the problem and are thus out of step with the RFCs.
4. Implementors of new features in Solaris sometimes forget that they must handle DAD on their own. For example, the `ifconfig`-like features in `zoneadmd` (which configure network interfaces within non-global zones) lack DAD, and configuring and running two zones on the same network with the same IP address will cause network failures.

## 4 Architectural Details

For both IPv4 and IPv6, DAD logic will be placed in the kernel. See Appendix A on page 22 for details on this decision.

The DAD logic will be triggered when the administrator sets the `IFF_UP` bit on the network interface and `IFF_RUNNING` is already set, when the `IFF_RUNNING` flag is set by the system and `IFF_UP` has already been set, or when the IP address is changed. This covers the three administrative actions: reconfiguration and network recabling. See Appendix B on page 23 for details on this decision.

To maintain compatibility and consistency with existing applications, those that read the flags directly with `ioctl`s (`SIOCGIFFLAGS` and `SIOCGLIFFLAGS`) will see `IFF_UP` even while DAD is running. If it were not done this way, then a typical application that reads the current flags, changes one flag, and then writes the flags back would inadvertently clear the `IFF_UP` bit and shut down the interface. (And it would be impossible to tell the difference between an application doing that by accident, and one that actually intends to shut down the network interface.)

A new `IFF_DUPLICATE` flag will be set by the in-kernel DAD logic if a failure occurs. This flag cannot be set or cleared by the user. It is automatically removed when either (a) `IFF_UP` is again set by the administrator and the interface is successfully established (i.e., the other system actively using the address has been disabled, removed, or reconfigured) or (b) the interface address is changed by the administrator (even if the interface is not brought up).

The routing socket interface changes slightly. Instead of reporting the new interface immediately (normally a sequence of messages consisting of `RTM_IFINFO` noting that the link has gone up, `RTM_NEWADDR` noting the local address, and `RTM_ADD` noting the local network), IP will hold off and report the new interface only after DAD completes. This is done to avoid having routing daemons and other applications that listen for new interfaces attempting to use and advertise the network before it has been validated. It also provides a simple way for applications to wait for DAD completion, if they are interested in it.

Although the routing socket timing is changed, this does not prevent those applications from seeing the flags via the `ioctl`s. If they do, they may see interfaces that go “up” briefly and then down again in the event of a duplicate address. It’s assumed that this does not happen often, but that if it does happen, it’s not catastrophic. The routing socket behavior described in this design is thus just an optimization.

Note also that there’s a possible failure mode here: if the user modifies the interface after setting `IFF_UP` but before DAD completes, routing socket messages will be sent up prematurely that show the interface as being up. To fix this, we could add another flag to indicate “DAD in progress.” However, the effect of this failure seems minimal (any real DAD failure will still cause the interface to be torn down and additional routing socket messages to be sent), the window is very small (1 to 5 seconds by default), the chance of it happening is slight (it requires multiple

administrators touching the same new interface), and the expense in fixing it is high (IFF\_\* flags are a precious resource). Thus, we won't fix this hole today, but may consider it in the future if it proves to be an actual problem.

While DAD is running, a small number of packets transmitted through IP by forwarding or applications may be buffered by the stack for transmission once DAD probing to validate the address completes successfully.

Failure of DAD will cause IFF\_UP to be reset by the system, as though “ifconfig [interface] down” were executed by an administrator, and the IFF\_DUPLICATE flag to be set, and will result in the generation of a log message and a sysevent<sup>1</sup> recording the problem.

The current support in DHCP (for v4), ifconfig (for v6), in.ndpd, and libinetcfg (v6 only) will be removed. These utilities will instead rely on the kernel to do the check when the 'up' flag is set.

No tuning is required or expected for most systems, but several undocumented ndd variables will be added in case we learn differently in deployment. Five ndd parameters are added for /dev/arp, as described in Table 1 below.

<i>Parameter</i>	<i>Default</i>	<i>Units</i>	<i>RFC name</i>
arp_probe_delay	1000	milliseconds	PROBE_WAIT
arp_probe_interval	1500	milliseconds	PROBE_MIN/MAX
arp_probe_count	3	packets	PROBE_NUM
arp_fastprobe_delay	100	milliseconds	PROBE_WAIT
arp_fastprobe_interval	150	milliseconds	PROBE_MIN/MAX
arp_fastprobe_count	3	packets	PROBE_NUM
arp_defense_interval	300000	milliseconds	see text
arp_publish_interval	2000	milliseconds	ANNOUNCE_INTERVAL
arp_publish_count	5	packets	ANNOUNCE_NUM

Table 1: New ARP NDD Parameters

Five ndd parameters are added for /dev/ip, as shown in Table 2 below.

The exact behavior of these tunables is described in the “Protocol Operation” sections below.

<sup>1</sup>To be determined later

<i>Parameter</i>	<i>Default</i>	<i>Units</i>	<i>RFC name</i>
ip_ndp_defense_interval	300000	milliseconds	
ip_max_temp_idle	86400	seconds	
ip_max_temp_defend	1	count	
ip_max_defend	3	count	
ip_defend_interval	30	seconds	DEFEND_INTERVAL

Table 2: New IP NDD Parameters

DAD probing for IPv4 can be disabled (if desired) by setting the ARP probe counts and defense interval to zero. When DAD is disabled, Solaris will launch straight into announcing the newly-configured address via gratuitous ARP, just as Solaris has historically done. Note, however, that ongoing conflict detection and resolution cannot (and should not) be disabled.

## 5 Protocol Operation – IPv4

### 5.1 Mechanisms

I have found two mechanisms for IPv4 DAD. One was originally described in Stuart Cheshire’s expired draft [1] and has found its way in modified form into RFC 3927 [2]. The other is buried in Microsoft’s Windows CE documentation [3]. See Appendix C on page 25 for details explaining why the RFC 3927 mechanism is the better choice.

See [2] for detailed operation of the protocol. In short, we need to be able to send broadcast “probe” packets, which look like this:

ar\$op	ares_op\$REQUEST
ar\$sha	local hardware L2 address
ar\$spa	INADDR_ANY
ar\$tha	all zeros
ar\$tpa	local interface IP address

Table 3: ARP Probe Contents

The above probe is crafted such that it elicits a response from any existing user of the address (which will see ar\$tpa as a request), but, unlike gratuitous ARP, will not corrupt the ARP caches of other nodes (which will see the empty ar\$spa and ignore the message).

## 5.2 Timers

The RFC recommends a random initial delay (up to one second), followed by default one-to-two second spacing of three probes, followed by a 2 second silence period before announcements begin. This causes a 4- to 7-second delay before an interface becomes usable.

As I'm enhancing an existing system where the expectation is that the IP interface is usable as soon as configured, this delay seems excessive and likely to cause service calls. Thus, I'll interpret the section 2.3 "Shorter Timeouts" of the RFC to apply to all drivers that respond to the Sun proprietary DLPI DL\_NOTIFY\_REQ message for DL\_NOTE\_LINK\_\* (i.e., those that can properly report link up/down events) and where the IP address involved is not in the special 169.254.0.0/16 link-local space. When both of these conditions apply, as they will for most interfaces, the system will shorten the timer for those interfaces by a factor of ten, resulting in at most a 700ms delay.

The "slow" behavior is tuned using `arp_probe_interval` for the timer and `arp_probe_count` for the number of probes needed. The "fast" behavior is tuned using `arp_fastprobe_interval` and `arp_fastprobe_count`.

In order eliminate self-synchronizing behavior, the timers mentioned here (and elsewhere in this document) are subjected to "fuzz." Instead of the min/max values used in the RFC, we define a central point (the timer interval), and then add or subtract a small random value from that interval. As a rough estimate, we will use +/- 20%, except for the initial delay, which is just a straight 0..N range.

Thus, the PROBE\_MIN and PROBE\_MAX constants from the RFC are both represented by `arp_probe_interval` in the table.

Note that MAX\_CONFLICTS and RATE\_LIMIT\_INTERVAL from the RFC aren't needed for this project, because they're specific to the randomized address selection needed for IPv4 Link Local addresses. This project doesn't implement support for those addresses.

### 5.3 Conflicts

There are two distinct types of conflicts that can be detected. In all cases, conflicts are at least logged for possible administrative action.

The first type of conflict is with an existing user of the address: the conflict is detected when any ARP message is received with `arp_spa` set to our IP address and `arp_sha` set to some other hardware address.

The second kind is a conflict with another system currently doing DAD, and is detected when a probe is received that has `arp_spa` set to our IP address and `arp_sha` set to some other hardware address.

If we detect a conflict of either kind during DAD probing, before we begin actively using the address, then we give up immediately. This is done by sending a new `AR_CN_FAILED` message to IP, and having IP tear down any interfaces that use the address. IP sets the `IFF_DUPLICATE` flag, clears `IFF_UP`, and sends routing socket messages.

If we are actively using the address after DAD probing completes and detect the second kind of conflict, one with someone sending DAD probes, then we defend our address by sending a gratuitous ARP. He should back off. This is done within ARP itself, which sees that probe as a regular query. The peer sees a conflict as in the previous paragraph.

If we are actively using the address and detect a conflict of the first kind, one with another active user of the address, then we have a choice to make: we need to determine whether to defend the address we've got or tear it down. As a heuristic, I will assume that if `IFF_DHCP_RUNNING` or `IFF_TEMPORARY` is set and the local address hasn't been used in a long time (by default, one day, but tunable using `ip_max_temp_idle`), then the system shouldn't try too hard to defend it (just once per `ip_defend_interval`, by default, and tunable using `ip_max_temp_defend`), as it's safe to give up and let `dhcpcagent` or `in.ndpd` give us a new address to try.

Otherwise, if the address is manually configured or if we've got active connections, we defend the address more aggressively, but give up if additional conflicts are detected as described in the RFC.

### 5.4 Ongoing Defense

As described above, it's possible for conflicting hosts to fail to detect each other during initial DAD probing, due to connectivity problems at the time

DAD runs or due to mishandling of ARP probes by Certain Router Vendor Software. Thus, addresses need to be defended after establishment.

In order to make this on-going conflict detection work, and be robust against network partition and repair, we must announce our addresses periodically. The RFC requires implementations to broadcast all ARP messages – both queries and replies. We will intentionally ignore this requirement. It's of questionable correctness for the limited usage within link-local addressing, but clearly wrong in all other contexts. Broadcasting an ARP association that you just broadcasted 10 milliseconds ago in response to back-to-back queries and forcing all other hosts to process this flood of broadcast replies is just foolish.

Instead, to achieve the correct effect, we will implement two timers for our published addresses. The first timer, set to 15 seconds by default (and tunable via `arp_broadcast_interval`), controls the minimum interval between broadcasted replies. If we are transmitting a unicast ARP message with an `ar$pa` that hasn't been broadcasted recently (within this interval), then the current message will be sent as a broadcast instead and the timer is reset. If the second timer, set to 5 minutes by default (and tunable via `arp_defense_interval`), expires without a given address having been broadcasted, then we'll send a traditional gratuitous ARP reply message.

With broadcasts at a 5 minute interval, a subnetwork with an improbable 2000 hosts on it will see less than 7 packets per second of purely ARP traffic. If someone created such a network, ARP messages would probably be the least of the concerns.

By default, when a conflict on an active address is detected, we defend up to 3 times (tunable with `ip_max_defend`) within any 30 second span (tunable with `ip_defend_interval`). The RFC describes a slightly different response: it defends just once within any 10 second span. The more complicated design chosen here allows for a packet drop or two between the conflicting peers, and allows the less-aggressive peer to back off.

Note that these values are also carefully chosen to work well with older systems that don't have DAD. Older Solaris systems blindly send out 5 gratuitous ARP messages (the sender doesn't stop, even if there are conflicts) when configuring an interface. We will defend against the first three conflicts, but then give up and let the other guy have the address on the 4th. There's no way to stop him from taking the address. This allows for correct behavior with at most one packet drop, while still being robust

against stray ARP messages.

## 5.5 Future Work

As an additional feature, we could be stubborn when giving up our last address, or shutting down the last interface with a default route pointing to it. When being “stubborn,” we would choose to defend the address for longer than suggested in the RFC, but not indefinitely. This feature is not included in the proposed implementation because it does not appear to be necessary.

We can also add an interface flag indicating that a given address should be defended indefinitely. This can be used on important servers that must never back down. It’s not clear whether this needs to be in the first release, or, if implemented, whether it deserves a new flag. (`IFF_PREFERRED` might be sufficient.)

All of these “stubbornness” heuristics are based on things that only IP knows. So, I will have ARP send `AR_CN_BOGON` to IP in these cases (as it has always done) and have IP decide how or if to defend the address. It does so by sending a new `AR_ENTRY_ADD` command back up to ARP with a new `ACE_F_DEFEND` flag set to tell ARP that it needs to send just one gratuitous ARP message or `AR_ENTRY_DELETE` if it chooses to shut down.

## 5.6 Summary of RFC Differences

Below is a list of differences between the behavior described in the RFC (by section) and the proposed design.

- Section 2.2 states that hosts “MUST NOT” check the address periodically. We, however, wish to have timely detection of failures in even rarely-used addresses and will thus send infrequent gratuitous ARPs for all addresses.

Note that rarely-used addresses may represent idle fail-over interfaces or other crucial resources. Detecting configuration errors in them only when they’re pressed into service would be a serious mistake.

- Section 2.2.1 specifies an unusual mechanism for generating the necessary “fuzz” in timer values to avoid synchronizing effects among multiple hosts. Rather than use this mechanism (which relies on separate minimum and maximum constants), we will borrow the better-known RIP timer mechanism, which uses a random delta around a single central value. The two are equivalent in effect, but the latter is easier to manage and understand.
- Section 2.2.1 also uses separate timers for ANNOUNCE\_WAIT and probing. For ease of implementation, and because these timers have about the same value, we’ll set the probe timer when we send the probe, and change state when this timer expires. This means that our delay between the last probe and the first announcement will be equivalent to a probe period (1 to 2 seconds) instead of the fixed 2-second timer described in the RFC.
- Section 2.3 weakens the commitment to short timers from that found in the original drafts. In essence, it mandates the longer timers for all common interfaces.

The timer policy set for DAD in this document conforms with the RFC within the narrow confines of the IPv4 Link Local address space. It does not appear to me that the RFC’s recommendations are necessarily sound, but the implementation will comply.

- Section 2.4 specifies 2 gratuitous ARP announcements after establishing an address. Existing Solaris practice is 5 announcements. We will retain the existing Solaris behavior to avoid uncovering bugs in peers.

## 6 Protocol Operation – IPv6

### 6.1 Basics

RFCs 2461 [4] and 2462 [5] together specify a default timer of 1 second (RetransTimer) and just one probe (DupAddrDetectTransmits), after an initial delay of 0 to 1 second (MAX\_RTR\_SOLICITATION\_DELAY). The current

implementation does this during `ifconfig` time (and in `in.ndpd` for temporary addresses), but fails to repeat the process when the link goes down and back up, check periodically, and use an initial delay.

As for IPv4 above, I'll assume that a link supporting `DL_NOTIFY_REQ` should also use a 150ms timer. However, as the single transmit probe ('required' by RFC 2462 [5]) seems a bit too weak to me, I'll increase this to a default of 3 and document the rationale (namely that it's much too likely that the first transmit is just dropped, and three isn't a lot of extra traffic). Ndd tunables for the retransmit time and the number of probes will be added.

The procedure is similar to the one for IPv4: Neighbor Solicitation messages are sent with a zero source address and the target address set to the desired address ("probes"), as shown below.

IP Source	all zeros
IP Destination	multicast
ICMP Type	135
Target Address	local interface IP address

Table 4: NDP Probe Contents

Conflicts are detected when either a Neighbor Advertisement or a probe is received with the desired address. (Solicitations with non-zero source are ignored; they're regular address resolution packets, and an interface doing DAD isn't ready for use yet.)

DAD is not performed if any of the `IFF_NOLOCAL`, `IFF_NONUD`, `IFF_ANYCAST`, or `IFF_LOOPBACK` flags is set, or if the local address is configured as unspecified ("`::`" – all zeroes).

## 6.2 Reporting Conflicts

One important wart in the implementation is that `ip_rput_v6()` currently discards the `M_PROTO` message (which contains the sender's hardware address) very early in the process. Since probe messages must not include the source link-layer address option (per RFC 2461 section 7.1.1), this means we have no easy way to report which node might be the current holder of the address. To fix this, I will change the v6 `rput` functions to

pass this mblk upwards into NDP. Probes are multicast, and thus always include M\_PROTO.

The existing `ifconfig(1M)/in.ndpd(1M)/libinetcfg` DAD code doesn't report the address of the current holder, either, so failing to do so for this project would not be a regression, but it's easy enough to do, so I'll add it.

### 6.3 The NOLOCAL Flag

In the existing code, `IFF_NOLOCAL` is set while DAD is in progress, and then reset afterwards. This, however, means that applications doing the usual `SIOCGLIFFLAGS/SIOCSLIFFLAGS` sequence to set or clear a flag may render newly-created interfaces unusable by accident, and applications will see the interface twice via routing socket messages. In the new code, the kernel will not set this flag at all during DAD, and, as with IPv4, delays the routing socket notification.

While this is a visible change in the system behavior, the combination of (a) unlikelihood that many applications rely on the behavior of the Solaris-specific `IFF_NOLOCAL` flag, (b) unlikelihood of actual duplicate addresses on any running network, and (c) unlikelihood that the up-to-down time for duplicates will seriously affect applications such that the flag would be needed means that the changes proposed ought to have no effect on existing or new applications.

Using `IFF_NOLOCAL` for DAD on all interfaces was considered and rejected. The problems include “remembering” the previous state – whether the user himself set the “nolocal” flag – so that it can be restored after DAD completes, and the lack of any apparent purpose for this usage, other than to affect kernel source address selection. Conditioning source address selection on DAD completion with internal ipif flags seems simpler.

### 6.4 Summary of RFC Differences

Below is a list of differences between the behavior described in the RFCs (by section) and the proposed design.

- RFC 2461 section 10 specifies a `RETRANS_TIMER` value of 1 second. As with IPv4, we use a much shorter timer on “reliable” media.

- RFC 2462 section 5.1 specifies a `DupAddrDetectTransmits` value of 1 packet. As with IPv4, we will use a default of 3 packets for more robust behavior.

## 7 Published Entries

Published ARP entries (other than local addresses) come from at least two sources: “`arp -s ip-address ether pub`” commands, and automatically configured proxy ARP entries from `pppd`.

The historical behavior for manually-configured static entries is that the “newest” definition of the mapping established on the network is the one that ought to be propagated. For that reason, the existing code treats updates for published non-self addresses specially: it updates the entry as though it were a regular ‘resolved’ entry. This behavior isn’t always what’s wanted by the administrator, especially for `pppd` proxy ARP entries, and that’s the source of the following RFE:

```
1253974 Please make a "permanent contents" option for
arp vs. static
```

Even the `ATF_PERM` flag doesn’t make the contents permanent, and an examination of the 4.4BSD code makes it clear that the “temp” option on the `/usr/sbin/arp` command (turning off `ATF_PERM`) refers only to the entry expiry time, not the immutability of the contents. OpenBSD, though, has a “permanent” command line option for the `arp` command that sets `RTF_PERMANENT_ARP` and makes the entry immutable. (Modern BSD uses a unified routing table with ARP entries, thus the `RTF` flag. Solaris does not have a unified table.)

Linux, as usual, appears to get this wrong, at least in the 2.4 kernel sources. It treats `ATF_PERM` as a license to ignore ARP updates. We likely cannot follow this bad example, as it would break existing configurations.

The behavior of DAD in these two cases is different. In the static entry case, no duplicate address detection is needed, and instead the entry needs to be kept up-to-date with changes on the network. In the `pppd` proxy ARP case, the entry must be treated like any other local address, with DAD and defense.

In order to reconcile the expected DAD behavior with the above traditional behaviors, it seems necessary to distinguish between the “authoritative” (e.g., PPP proxy arp) and “non-authoritative” (e.g., traditionally mutable “static” entry) behavior. The solution to that is exactly what’s needed for the above RFE (as is present in OpenBSD), and thus we’ll include a resolution to that RFE with this project, along with PPP and Mobile IP updates.

A new ATF\_AUTHORITY (authoritative entry) flag will be established for ARP entries. This will be set when the /usr/sbin/arp “permanent” keyword is used, PPP establishes proxy-arp entries, and with mipagent for Mobile IP. In the kernel, this will map to ACE\_F\_AUTHORITY, and internal use of ACE\_F\_MYADDR (by ip) will imply ACE\_F\_AUTHORITY. (Note that ACE\_F\_AUTHORITY records needn’t be published, but that ACE\_F\_MYADDR ones must be. Unpublished authoritative records might be useful for users wishing to avoid the insecurity of ARP.)

## 8 Compatibility

There are several cases to consider (“old” means a Solaris system prior to the DAD changes described in this document, while “new” means a system including the DAD changes):

1. Old system starts up with an address already in use by a new system.

The new system will detect the gratuitous ARP from the old system as a conflict with an existing address, and will choose to defend the address.

This should result in errors logged on both offending systems, and with the new system in control of the address, at least temporarily. As the old system won’t back off, the new system will stop defending the address after three tries, and thus lose.

2. New system starts up with an address already in use by an old system.

The old system responds correctly to the probes, and the new interface is shut down as a duplicate. No ARP cache pollution results, as the probes contain no source IP address.

3. New system starts up with an address already in use by another new system.

This case is documented in the Protocol Operation section. The system just starting up loses, and the existing one is unaffected.

4. New system starts up with an address already in use by some non-Solaris system.
  - (a) Linux
  - (b) Windows
  - (c) Cisco
  - (d) AIX
  - (e) HP/UX
5. Non-Solaris system starts up with an address already in use by a new Solaris system.

## 9 Interactions With Other Projects

### 9.1 IPMP

IPMP indirectly uses gratuitous ARP when it moves an address from one interface to another, so that all of the local peers are notified of the new location. This will be treated as an instance of “address defense” as described above, except that the draft’s prohibition against sending more than one gratuitous ARP within 10 seconds will be ignored for the sake of backward compatibility.

This is logically similar to the case where the local Ethernet interface has its address changed on the fly. We still own the mapping and there are no known conflicts, so we’ll send out gratuitous ARPs to announce the change as we’ve always done.

### 9.2 DHCP

The `dhcpageant(1M)` daemon contains a small amount of duplicate address detection logic, and the logic isn’t well-integrated with the daemon’s

event-driven design. (It sleeps for 5 seconds waiting on a per-interface response right in the middle of the main multi-interface polling loop, rather than treating it as a separate state.)

The current code in `cmd/cmd-inet/sbin/dhcpagent/arp_check.c` sends a single probe similar to the one described above for the new DAD procedure (using all-ones for `ar$tha` instead of all-zeros), and then waits 5 seconds for any `ares_op$REPLY` where `ar$spa` has the desired IP address. If one is received, then it's treated as a duplicate address, and the state machine is notified.

Unlike RFC 3927[2], this doesn't send multiple packets in case the first one is lost and doesn't deal with `ares_op$REQUEST` messages. It thus doesn't respond correctly to gratuitous ARPs from peers, and instead just discards those. (RFC 2131[6] recommends the use of ARP to detect duplicate addresses, but does not describe the procedure in much detail. Implementations likely vary considerably in behavior.)

In the new design, this logic is all removed. The DHCP agent needs a new state – `PRE.BOUND` – to do its work. In the `PRE.BOUND` state, we have configured an interface and marked it `IFF_UP`, but have not yet sent the DHCP ACK or set up routes. DHCP then transitions from `PRE.BOUND` to `BOUND` if the routing socket message that marks it `IFF_UP` is seen, or back to `INIT` (to try again) if `IFF_DUPLICATE` is seen.

### 9.3 Routing

The routing protocols will see “up” interfaces during DAD probing when they use `SIOCGIFCONF`, and will attempt to use them. Because the time constants used for these protocols are much larger than DAD (10 seconds for OSPF and 5 to 30 seconds for RIP), only one packet might be delayed or lost, and the effect should be minimal.

## 10 Spanning Tree

Spanning tree (STP; 802.1D[9]) is known to cause lengthy delays for hosts attached to networks via switches; on the order of tens of seconds on some networks. These delays cause the initial phase of DAD – where we attempt to discover whether there are any duplicates before going “live” with the configured address – to fail to detect any conflicts. It's thus up to on-going

conflict resolution (as in the case of repair after network partitioning) to find the problem.

This is clearly suboptimal, as on-going conflict resolution can cause disruption. Worse, there's no way for an external node to tell a switch not to create these delays; no way to tell the switch that the node is just a leaf and cannot forward. The problem can only be eliminated by an administrator disabling STP on the port, and most switches have STP enabled by default.

Once STP is configured, the port goes through Blocking, Listening, Learning, and (finally) Forwarding states strictly on a timer (the "Forward Delay Timer"). Even worse still, there's no way to know when the port is discarding packets, because there is no precise way to know when it enters or leaves Forwarding state.

(Note that Topology Change Notification BPDUs, which might provide some clue about link stability, are sent towards the Spanning Tree Root only, not towards the leaves. Simple end nodes are never the Root, and thus never see these messages.)

A switch that's using STP emits 802.3 packets to 01:80:c2:00:00:00 with SSAP/DSAP set to 42 and control set to 03. The fifth octet in the message is a set of flags. The LSB (01) is set to indicate "Topology Change." In other words, it's a packet that (in part) looks like this:

```
0000  01 80 c2 00 00 00 00 04  96 18 42 a9 00 26 42 42
0010  03 00 00 00 00 01
                        **
```

According to the standards (and in experiments with Extreme switches), the TC bit is set shortly after forwarding is enabled or disabled, and remains set for the length of the Forwarding Delay Timer. Because TC behaves as a one-shot, multiple port transitions can keep it set for extended periods of time. Thus, watching for a 0-to-1 transition is not necessarily accurate.

However, it is fair to conclude the following:

- If any STP packets are seen after a link goes active, then there may be bridges that are currently partitioning the network. Any time there's a 1-to-0 transition on the TC bit, we should trigger DAD or address defense.

- There may be bridges that we don't know about. Thus, as a precaution, and observing that 15 seconds is a typical value for the Forwarding Delay parameter (and thus about 30 seconds total time to enable a port), we make sure that ongoing address defense checks are performed at 40 seconds after the port is enabled.

We will at least do the latter.

## 11 Implementation

DAD probing begins when the arp module receives `AR_ENTRY_ADD` for a new mapping. On successful completion, at the same time the first gratuitous ARP announcement is transmitted, `AR_CLIENT_NOTIFY` with `AR_CN_READY` is sent. On failure, `AR_CLIENT_NOTIFY/AR_CN_FAILED` is sent.

The arp module sends `AR_ARP_EXTEND` to ip (when it is first pushed on the stream) to tell it that DAD will be performed. This signal indicates that IP should not try to use any address on this interface until DAD completes. If the resolver doesn't send `AR_ARP_EXTEND`, this means that it's one of the "old" (ATM-related) ARP modules, and IP should assume the old gratuitous ARP based behavior is needed. IP doesn't wait for this message from ARP; it just assumes the old-style behavior unless told otherwise.

While we must assume that there are other external resolvers (at least until ATM is finally retired to the bit-bucket), we needn't assume that there are any other valid clients for the 'arp' module other than ip. In other words, the arp module can assume that the client understands the new messages. (Such compatibility, though, could be added in the future if needed by defining a response for `AR_ARP_EXTEND`.)

If a duplicate address is detected during normal operation (after DAD probing has completed), the arp module will send `AR_CLIENT_NOTIFY / AR_CN_BOGON` to ip, just as it has done in the past. In this case, the ip module makes the determination of what action to take next: either tear down the affected interface, send a new `AR_ENTRY_ADD` to defend it, or just discard the message to ignore the notification.

An important performance consideration for ARP is that we must not send a message to IP each time we get a null (same address) update, or an ARP storm can cause the system to lock up, as described in:

4653899 ARP hurricane can deny service

In the old code, this was accomplished in a rather round-about way. Since we deleted resolved entries from our own table once we told IP the resolution, we couldn't determine whether any given packet represented a change in address or just another instance of the same address. We thus deleted all but the ones that looked like gratuitous ARPs (`ar$spa == ar$tpa`).

This strategy had multiple flaws. First of all, it ignored the ARP RFC, which requires us to notice an address change on any ARP message we receive. This has caused occasional interoperability problems and complaints, as described in:

4157198 ARP cache inconsistency between arp and ip modules.

Secondly, it meant that if we ever got a storm of gratuitous ARPs, we'd be back in the same denial-of-service mess. Only part of the hurricane problem was solved.

To fix this problem, we'll keep the resolved entries around in the arp module as long as the corresponding information is known to be cached by IP. To do this, we set a 20 minute timer on each ARP-resolved entry, rather than just deleting it when resolution is complete. If the timer expires, we check with IP to see if there's still an IRE cache entry. If so, then reset the timer. If not, then delete the entry, and allow future updates to appear to be "new."

## 12 UnARP

As a part of restructuring the ARP code to keep all of the resolved entries in the table, it becomes trivial to handle RFC 1868 [8] "UnARP." As a part of this project, support for handling of these messages (with `ar$hlen` set to zero) will be added to the receive side. Transmission of UnARP messages as well (when deleting an authoritative published entry) might be possible, but isn't a goal of the design.

## A User or Kernel?

DAD could be done in either user space or kernel space. There are several clear advantages and reasons to prefer a user space implementation:

- Configuring new IP addresses is in the control path, not the data path, so it's not very performance sensitive and thus doesn't need to be inside the kernel.
- If there are flaws in the implementation, failing in user space just results in a core dump rather than a kernel panic, and patches could be done on a live system rather than with a reboot.
- Linking into administrative bits and logging is potentially easier.
- The API used could support alternative implementations, if such a thing were ever necessary.

The drawbacks, though, are:

- There is no daemon that does anything quite like this. We'd either have to add yet another almost always idle but critical system daemon, or hijack and extend an existing daemon to do the work.
- If it crashed or were disabled or removed by the administrator (misplaced security fears?), IP would just fall dead.
- The protocol itself is both quite simple and entangled with normal ARP entry maintenance (e.g., local addresses versus configured entries), IP operation (e.g., routing socket messages and interface flags), and DLPI information (e.g., link up/down notifications).

Thus, it seems that little purpose is served by placing it in a separate critical system daemon, so it will be put in the kernel. It could be moved to user space later, if desired. Doing so would require the creation of an API to control ARP responses.

## B Starting DAD Operation

There are multiple ways that the DAD feature could be accommodated. Here are the benefits and drawbacks of the many that have been discussed.

1. Start DAD when user sets IFF\_UP, but don't read back IFF\_UP until DAD is done.

This is the "conservative" choice. It assumes that applications may read the flags while DAD is in progress, and then be confused about the state of the interface, as it isn't quite up until DAD probing completes.

The key failing is the potential for an unintended-command side effect. If some application reads the current flags, changes them, and then writes the flags back, this creates an ambiguity. Did that application attempt to turn the interface off (turning off the IFF\_UP bit that hasn't yet been set), or did it just leave the IFF\_UP bit unchanged?

2. Start DAD when user sets IFF\_UP, and leave IFF\_UP set until (and unless) DAD fails.

This is the "optimistic" choice. It assumes that DAD failure will be quite rare – thus, the up-down flap visible in the flags in the failure case won't have significant impact – and that applications that will be confused by seeing the "UP" bit set when the interface isn't yet fully usable are unlikely to exist.

3. Create a new flag (IFF\_UP\_WANTED) that triggers DAD. Set IFF\_UP internally only when DAD is complete.

This case is a refinement on (1). It almost avoids the unintended-command problem, but causes a new one: it's no longer compatible with existing programs that may set IFF\_UP on their own. We are thus forced to solve that problem as well, in order to keep compatibility, likely using solution (2).

It only "almost" solves the problem, as the command ambiguity still exists for standard applications that use SIOCSIFFLAGS rather than the Solaris-proprietary SIOCSLIFFLAGS because no 32 bit flags are left.

4. Create some new interface entirely for requesting an interface to be marked as “up,” and rewrite everything to use that interface.

This is essentially equivalent to (3) in terms of risks and benefits.

5. Create a new `IFF_DAD_IN_PROGRESS` flag, set by the system along with `IFF_UP`. This is a refinement on (2). It allows applications to detect the DAD probing interval. However, the applications that might be possibly confused are already deployed, so the new bit likely does little good. In fact, it’s not clear who could use such a flag.

The risk of unintended commands is hard to quantify, but the consequences are quite unfortunate. It will result in interfaces that are left either up or down when the opposite was desired.

The risk of confusion due to seeing the UP bit set when the interface is not fully up is limited to applications that scan interfaces during DAD, and the consequences are limited. If the application uses these addresses for binding (e.g., a multicast applications), then it’ll bind to an interface that can’t yet send or receive.

This likely results in a small amount of packet loss, but that can be covered by buffering in IP during DAD, just as is done when regular ARP resolution is in progress. Small amounts of packet loss are likely not completely unexpected at interface start-up time. (In fact, STP makes such loss common.)

If the application uses these addresses in order to tell other systems about connectivity (e.g., routing applications), then it may erroneously report that these networks are reachable before it’s actually known whether they’re reachable. If DAD succeeds (the most likely case), it’s at most just a short blip. If it doesn’t, then the network interface will go back down, and the result is a short-lived black hole.

Of course, interfaces might go down at any time for administrative reasons, so if applications are completely confused by interfaces appearing to go up and then back down, then those applications are already in serious trouble even without the prospect of DAD failure.

Thus, as the least problematic and most compatible solution, I pick (2). To lessen the likelihood of confused applications, the routing socket messages associated with the new interface will be delayed until DAD completes. Since most applications that care about interfaces going up or down listen to routing socket messages rather than scanning the kernel

with SIOCGIFCONF periodically, this should minimize the chance that any application ever even notices the short up-but-not-quite-up interval caused by DAD.

## C Choosing DAD Standard

The documentation for the actual process followed by Windows CE is sketchy, but through various documents found on the Web, including “Introduction to IP Version 6,” I’ve pieced together the following picture.

Windows CE sends out a gratuitous ARP when the interface is configured, just like Solaris does today. This has the traditional gratuitous ARP format – both `ar$spa` and `ar$tpa` set to the local address. The difference appears to be that Windows CE doesn’t update its cache when it receives such a request message, so an all-Windows network won’t get confused.

That would allow it to use these requests as an effective DAD mechanism. However, such a thing would be in direct violation of RFC 826 [7], which says that the cache is updated before looking at `ar$op`.

More importantly for us, this leaves the backward-compatibility problem on the table. Even if we could violate RFC 826 in this way for new systems, writing a special case just for packets that are obviously gratuitous ARP (`ar$spa == ar$tpa`), it would leave us with the problem that existing RFC-correct systems will still be confused by these packets, and will still update their local caches when a duplicate address is configured.

Thus, either Microsoft’s documentation is inaccurate, and they actually use Cheshire’s method, or their method is useless for us. Either way, Cheshire’s method is preferred.

## References

- [1] Stuart Cheshire, *IPv4 Address Conflict Detection*, December 2002 (expired work-in-progress).
- [2] Stuart Cheshire, et al., RFC 3927: *Dynamic Configuration of IPv4 Link-Local Addresses*, May 2005.
- [3] Microsoft Corporation, *Microsoft Windows CE .NET 4.2 Duplicate IP Address Detection for IPv4*.

- [4] Thomas Narten, et al., RFC 2461: *Neighbor Discovery for IP Version 6 (IPv6)*, December 1998.
- [5] Susan Thomson and Thomas Narten, RFC 2462: *IPv6 Stateless Address Autoconfiguration*, December 1998
- [6] Ralph Droms, RFC 2131: *Dynamic Host Configuration Protocol*, March 1997.
- [7] David C. Plummer, RFC 826: *Ethernet Address Resolution Protocol: Or converting network protocol addresses to 48.bit Ethernet address for transmission on Ethernet hardware*, November 1982.
- [8] Gary Malkin, RFC 1868: *ARP Extension - UNARP*, November 1995.
- [9] *Part 3: Media Access Control (MAC) Bridges*, ANSI/IEEE Std 802.1D, 1998 Edition (ISO/IEC 15802-3: 1998).