ARMD Working Group                                      M. Karir
Internet Draft                                 Merit Network Inc.
Intended status: Informational Track                     Ian Foo
Expires: January 2012                        Huawei Technologies

                                                October 24, 2011


                   Data Center Reference Architectures
                draft-armd-datacenter-reference-arch-01.txt


Status of this Memo

Abstract

   The continued growth of large-scale data centers has resulted in a
   wide range of architectures and designs.  Each design is tuned to
   address the challenges and requirements of the specific applications
   and workload that the data is being built for.  Each design evolves
   as engineering solutions are developed to workaround limitations of
   existing protocols, hardware, as well as software implementations.

   The goal of this document is to characterize this problem space in
   detail in order to better understand if there is any gap in making
   address resolution scale in various network designs for data
   centers.  In particular it is our goal to peel back the various
   optimization and engineering solutions to develop generalized
   reference architectures for a data center.  We also discuss the
   various factors that influence design choices in developing various
   data center designs.

Conventions 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 RFC-2119 0.

Table of Contents

1. Introduction

   Data centers are a key part of delivering Internet scale
   applications.  Data center design and network architecture is an
   important aspect of the overall service delivery plan.  This
   includes not only determining the scale of physical and virtual
   servers but also optimizations to the entire data center stack
   including in particular the layer 3 and layer 2 architectures.
   Depending on the particular application requirements and scale, data
   centers can be designed in variety of ways.  Each design is often a
   representation of which aspects of the problem were and were not
   relevant to the purpose of that data center.  In this document we
   attempt to generalize the various design optimizations into a common
   generic architecture to facilitate the discussion of potential
   issues under a common framework.

2. Terminology

   ARP:      Address Resolution Protocol

   ND:       Neighbor Discovery

   Host:     Application running on a physical server or a virtual
             machine. A host usually has at least one IP address and at
             least one MAC address.

   Server:   a physical computing machine

   ToR:       Top of Rack Switch

   EoR:       End of Row

   VM:       Virtual Machines. Each server can support multiple VMs.

3. Generalized Data Center Design

   There are many different ways in which data centers might be
   designed.  The designs are usually engineered to suit the particular
   application that is being deployed in the data center.  For example,
   a massive web sever farm might be engineered in a very different way
   than a general-purpose multi-tenant cloud hosting service.  However
   in most cases the designs can be abstracted into a typical three-
   layer model consisting of the Access Layer, the Aggregation Layer
   and the Core.  The access layer generally refers to the Layer 2
   switches that are closest to the physical or virtual severs, the
   aggregation layer refers to the Layer 2 - Layer 3 boundary.  The
   Core switches connect the aggregation switches to the larger network
   core.  Figure 1 shows a generalized Data Center design, which
   captures the essential elements of various alternatives.

```
        +-----+-----+      +-----+-----+
        |   Core0   |      |   Core1   |        Core
        +-----+-----+      +-----+-----+
            /     \          /       /
           /       \----------\     /
          /    /--------/       \  /
        +------+            +------+
      +/------+ |         +/-----+ |
      | Aggr11| + --------|AggrN1| +        Aggregation Layer
      +---+---+/          +------+/
         /     \            /      \
        /       \          /        \
      +---+    +---+      +---+      +---+
      |T11|... |T1x|      |T21| ... |T2y| Access Layer
      +---+    +---+      +---+      +---+
      |   |    |   |      |   |      |   |
      +---+    +---+      +---+      +---+
      |   |... |   |      |   | ... |   |
      +---+    +---+      +---+      +---+
      |   |... |   |      |   |     |   |   Server racks
      +---+    +---+      +---+     +---+
      |   |... |   |      |   | ... |   |
      +---+    +---+      +---+      +---+
      |   |... |   |      |   | ... |   |
      +---+    +---+      +---+      +---+
```

              Figure 1: Typical Layered Architecture in DC

3.1. Access Layer

   The Access switches provide connectivity directly to/from physical
   and virtual servers.  The access switches might be placed either on
   top-of-rack (ToR) or at end-of-row(EoR) physical configuration. A
   server rack may have a single uplink to one access switch, or may
   have dual uplinks to two different access switches.

3.2. Aggregation Layer

   In a typical data center, aggregation switches interconnect many ToR
   switches. Usually there are multiple parallel aggregation switches,
   serving the same group of ToRs to achieve load sharing. It is no
   longer uncommon to see aggregation switches interconnecting hundreds
   of ToR switches in large data centers.

3.3. Core

   Core switches connect multiple aggregation switches and act as the
   data center gateway to external networks or interconnect to
   different PODs within one data center.

3.4. Layer 3 / Layer 2 Topological Variations

3.4.1. Layer 3 to Access Switches

   In this scenario the L3 domain is extended all the way to the Access
   Switches.  Each rack enclosure consists of a single Layer 2 domain,
   which is confined to the rack.  In general in this scenario there
   are no significant ARP/ND scaling issues as the Layer 2 domain
   cannot grow very large.  This topology is ideal for scenarios where
   servers (or VMs) under one access switch don't need to be re-loaded
   with applications with different IP addresses or hosts don't need to
   be moved to other racks which are under different access switches.
   A small server farm or very static compute cluster might be best
   served via this design.

3.4.2. L3 to Aggregation Switches

   When Layer 3 domain only extends to aggregation switches, hosts in
   any of the IP subnets configured on the aggregation switches can be
   reachable via Layer 2 through any access switches if access switches
   enable all the VLANs. This topology allows for a great deal of
   flexibility as servers attached to one access switch can be re-
   loaded with applications with different IP prefix and VMs can now
   migrate between racks without IP address changes.  The drawback of
   this design however is that multiple VLANs have to be enabled on all

access switches and all ports of aggregation switches. Even though
layer 2 traffic are still partitioned by VLANs, the fact that all
VLANs enabled on all ports can lead to broadcast traffic on all
VLANs to traverse all links and ports, which is same effect as one
big Layer 2 domain.  In addition, internal traffic itself might have
to cross different Layer 2 boundaries resulting in significant
ARP/ND load at the aggregation switches.  This design provides the
best flexibility/Layer 2 domain size trade-off.  A moderate sized
data center might utilize this approach to provide high availability
services at a single location.

3.4.3. L3 in the Core only

In some cases where wider range of VM mobility is desired (i.e.
greater number of racks among which VMs can move without IP address
change), the Layer 3 routed domain might be terminated at the core
routers themselves.  In this case VLANs can span across multiple
groups of aggregation switches, which allow hosts to be moved among
more number of server racks without IP address change. This scenario
results in the largest ARP/ND performance impact as explained later.
A data center with very rapid workload shifting may consider this
kind of design.

3.4.4. Overlays

There are several approaches regarding how overlay networks can make
very large layer 2 network scale and enable mobility. Overlay
networks using various Layer 2 or Layer 3 mechanisms enable interior
switches/routers not to see the hosts' addresses. The Overlay Edge
switches/routers which perform the network address
encapsulation/decapsulation still however see host addresses.

When a large data center has tens of thousands of applications which
communicate with peers in different subnets, all those applications
send (and receive) data packets to their L2/L3 boundary nodes if the
targets are in different subnets. The L2/L3 boundary nodes have to
process ARP/ND requests sent from originating subnets and resolve
physical addresses (MAC) in the target subnets. In order to allow a
great number of VMs to move freely within a data center without re-
configuring IP addresses, they need to be under the common Gateway
routers. That means the common gateway has to handle address
resolution for all those hosts.  Therefore, the use of overlays in
the data center network can be a useful design mechanism to help
manage a potential bottleneck at the Layer 2 / Layer 3 boundary by
redefining where that boundary exists.

4. Factors that Affect Data Center Design

4.1. Traffic Patterns

   Expected traffic patterns play an important role in designing the
   appropriately sized Access, Aggregation and Core networks.  Traffic
   patterns also vary based on the expected use of the Data Center.
   Broadly speaking it is desirable to keep as much traffic as possible
   on the Access Layer in order to minimize the bandwidth usage at the
   Aggregation Layer.  If the expected use of the data center is to
   serve as a large web server farm, where thousands of nodes are doing
   similar things and the traffic pattern is largely in/out a large
   access layer with EoR switches might be of the most use as it
   minimizes complexity, allows for servers and databases to be located
   in the same Layer 2 domain and provides for maximum density.

   A Data Center that is expected to host a multi-tenant cloud hosting
   service might have completely different requirements where in order
   to isolate inter-customer traffic smaller Layer 2 domains are
   preferred and though the size of the overall Data Center might be
   comparable to the previous example, the multi-tenant nature of the
   cloud hosting application requires a smaller more compartmentalized
   Access layer.  A multi-tenant environment might also require the use
   of Layer 3 all the way to the Access Layer ToR switch.

   Yet another example of an application with a unique traffic pattern
   is a high performance compute cluster where most of the traffic is
   expected to stay within the cluster but at the same time there is a
   high degree of crosstalk between the nodes.  This would once again
   call for a large Access Layer in order to minimize the requirements
   at the Aggregation Layer.

4.2. Virtualization

   Using virtualization in the Data Center further serves to increase
   the possible densities that can be achieved.  Virtualization also
   further complicates the requirements on the Access Layer as that
   determines the scope of server migrations or failover of servers on
   physical hardware failures.

   Virtualization also can place additional requirements on the
   Aggregation switches in terms of address resolution table size and
   the scalability of any address learning protocols that might be used
   on those switches. The use of virtualization often also requires the
   use of additional VLANs for High Availability beaconing which would
   need to span across the entire virtualized infrastructure.  This

would require the Access Layer to span as wide as the virtualized
infrastructure.

4.3. Impact of Data Center Design on L2/L3 protocols

When a L2/L3 boundary router receives data packets via its L3
interfaces destined towards hosts under its L2 domain, if the target
address is not present in the router's ARP/ND cache, it usually
holds the data packets and initiates ARP/ND requests towards its L2
domain to make sure the target actually exists before forwarding the
data packets to the target. If no response is received, the router
has to send the ARP/ND multiple times. If no response is received
after X number ARP/ND requests, the router needs to drop all those
data packets. This process can be very CPU intensive.

When a local host under the L2/L3 Router's L2 domain needs to send a
data frame to external peers, it usually sends ARP/ND requests to
get the physical address (i.e. MAC) of the L2/L3 routers. Many hosts
repetitively send ARP/ND requests to their default L3 gateway
routers to refresh its ARP/ND cache. This requires default routers
to process great number of ARP/ND requests when the number of hosts
under its L2 domains is very large. For IPv4, gateway routers
frequently sending out gratuitous ARP for all the hosts under its L2
domain to refresh their ARP cache for the default gateway's MAC
address can mitigate this pain point. However, for IPv6 hosts need
to validate bi-direction communication with the gateway router
before sending any data frames. Therefore, unsolicited neighbor
announcement from gateway router can't prevent hosts from sending ND
repetitively.

When hosts in two different subnets under the same L2/L3 boundary
router need to communicate with each other, the L2/L3 router not
only has to initiate ARP/ND requests to the target's Subnet, it also
has to process the ARP/ND requests from the originating subnet. This
process is even more CPU intensive.

5. Conclusion and Recommendation

In this document we have described a generalized Data Center network
design.  Our goal is to distill the essence of different designs
into a common framework in an attempt to structure the discussion
regarding various scaling issues that might appear in different
scenarios.  Different application needs such as traffic patterns,
and the role for which the data center is being designed determine
various design choices, which result in various scaling issues with
regards to port density, ARP/ND, VM mobility, and performance.  As

expected, engineering solutions serve to tune a given design to the
particular needs of the data center at the expense of other factors.

6. Manageability Considerations

   This document does not add additional manageability considerations.

7. Security Considerations

   This document has no additional requirement for security.

8. IANA Considerations

   None.

9. Acknowledgments

   We want to acknowledge the following people for their valuable
   discussions related to this draft: Kyle Creyts, Alexander Welch and
   Michael Milliken

   This document was prepared using 2-Word-v2.0.template.dot.

10. References

   [ARP]    D.C. Plummer, "An Ethernet address resolution protocol."
             RFC826, Nov 1982.

   [ND]     T. Narten, E. Nordmark, W. Simpson, H. Soliman, "Neighbor
             Discovery for IP version 6 (IPv6)." RFC4861, Sept 2007.

   [STUDY]  Rees, J., Karir, M., "ARP Traffic Study." MANOG52, June
             2011. URL
             http://www.nanog.org/meetings/nanog52/presentations/Tuesda
             y/Karir-4-ARP-Study-Merit Network.pdf

   [DATA1]  Cisco Systems, Data Center Design - IP Infrastructure ,
             October 2009. URL
             http://www.cisco.com/en/US/docs/solutions/Enterprise/Data_
             Center/DC_3_0/DC-3_0_IPInfra.html

   [DATA2]  Juniper Networks, Government Data Center Network Reference
             Architecture, 2010. URL
             www.juniper.net/us/en/local/pdf/reference-
             architectures/8030004-en.pdf

Authors' Addresses

   Manish Karir
   Merit Network Inc.
   1000 Oakbrook Dr, Suite 200
   Ann Arbor, MI 48104, USA
   Phone: 734-527-5750
   Email: mkarir@merit.edu

   Ian Foo
   Huawei Technologies
   2330 Central Expressway
   Santa Clara, CA 95050, USA
   Phone: 919-747-9324
   Email: Ian.Foo@huawei.com


Intellectual Property Statement

                        Problem Statement for ARMD
                  draft-ietf-armd-problem-statement-00

Abstract

   This document examines issues related to the massive scaling of data
   centers.  Our initial scope is relatively narrow.  Specifically, we
   focus on address resolution (ARP and ND) within the data center.

Status of this Memo

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

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

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

   This Internet-Draft will expire on April 20, 2012.

Table of Contents

1.  Introduction

   This document examines issues related to the massive scaling of data
   centers.  Specifically, we focus on address resolution (ARP in IPv4
   and Neighbor Discovery in IPv6) within the data center.  Although
   strictly speaking the scope of address resolution is confined to a
   single L2 broadcast domain (i.e., ARP runs at the L2 layer below IP),
   the issue is complicated by routers with many interfaces (on which
   address resolution is performed) or with IEEE 802.1Q domains, where
   individual VLANs form their own broadcast domains.  Thus, the scope
   of address resolution spans both the L2 link and the devices attached
   to those links.

   This document is intended to support the ARMD WG identify potential
   future work areas.  The scope of this document intentionally starts
   out relatively narrow, mirroring the ARMD WG charter.  Expanding the
   scope requires careful thought, as the topic of scaling data centers
   generally has an almost unbounded potential scope.  This document
   aims to list "pain points" that are being experienced in current data
   centers.  It is separate exercise to determine which (if any) of
   these pain points should lead to specific protocol work, whether in
   ARMD or some other WG.


2.  Terminology

   Application:  a service that runs on either a physical or virtual
      machine, providing a service (e.g., web server, database server,
      etc.)

   Broadcast Domain:  The set of all links and switches that are
      traversed in order to reach all nodes that are members of a given
      L2 domain.  For example, when sending a broadcast packet on a
      VLAN, the domain would include all the links and switches that the
      packet traverses when broadcast traffic is sent.

   Host (or server):  Physical machine on which a system is run.  A
      system can consist of an application running on an operating
      system on the "bare metal" or multiple applications running within
      individual VMs on top of a hypervisor.  Traditional non-
      virtualized systems will have a single (or small number of) IP
      addresses assigned to them.  In contrast, a virtualized system
      will use many IP addresses, one for the hypervisor plus one (or
      more) for each individual VM.

Hypervisor:  Software running on a host that allows multiple VMs to
    run on the same host.

L2 domain:   IEEE802.1Q domain supporting up to 4095 VLANs.  The
    notion of an L2 broadcast domain is closely tied to individual
    VLANs.  Broadcast traffic (or flooding to reach all destinations)
    reaches every member of the specific VLAN being used.

Virtual machine (VM):  A software implementation of a physical
    machine that runs programs as if they were executing on a bare
    machine.  Applications do not know they are running on a VM as
    opposed to running on a "bare" host or server.


3.  Background

   Large, flat L2 networks have long been known to have scaling
   problems.  As the size of an L2 network increases, the level of
   broadcast traffic from protocols like ARP increases.  Large amounts
   of broadcast traffic pose a particular burden because every device
   (switch, host and router) must process and possibly act on such
   traffic.  In addition, large L2 networks can be subject to "broadcast
   storms".  The conventional wisdom for addressing such problems has
   been to say "don't do that".  That is, split large L2 networks into
   multiple smaller L2 networks, each operating as its own L3/IP subnet.
   Numerous data center networks have been designed with this principle,
   e.g., with each rack placed within its own L3 IP subnet.  By doing
   so, the broadcast domain (and address resolution) is confined to one
   Top of Rack switch, which works well from a scaling perspective.
   Unfortunately, this conflicts in some ways with the current trend
   towards dynamic work load shifting in data centers and increased
   virtualization as discussed below.

   Workload placement has become an issue within data centers.  Ideally,
   it is desirable to be able to move workloads around within a data
   center in order to optimize server utilization, add additional
   servers in response to increased demand, etc.  However, servers are
   often pre-configured to run with a given set of IP addresses.
   Placement of such servers is then subject to constraints of the IP
   addressing restrictions of the data center.  For example, servers
   configured with addresses from a particular subnet could only be
   placed where they connect to the IP subnet corresponding to their IP
   addresses.  If each top of rack switch is placed within its own
   subnet, a server can only be connected to the one top of rack switch.
   This same constraint occurs in virtualized environments, as discussed
   next.

   Server virtualization is fast becoming the norm in data centers.

With server virtualization, each physical server supports multiple
virtual servers, each running its own operating system, middleware
and applications.  Virtualization is a key enabler of workload
agility, i.e., allowing any server to host any application and
providing the flexibility of adding, shrinking, or moving services
among the physical infrastructure.  Server virtualization provides
numerous benefits, including higher utilization, increased data
security, reduced user downtime, and even significant power
conservation, along with the promise of a more flexible and dynamic
computing environment.

The discussion below focuses on VM placement and migration.  Keep in
mind, however, that even in a non-virtualized environment, many of
the same issues apply to individual workloads running on standalone
machines.  For example, when increasing the number of servers running
a particular workload to meet demand, placement of those workload may
be constrained by IP subnet numbering considerations.

The greatest flexibility in VM and workload management occurs when it
is possible to place a VM (or workload) anywhere in the data center
regardless of what IP addresses the VM uses and how the physical
network is laid out.  In practice, movement of VMs within a data
center is easiest when VM placement and movement does not conflict
with the IP subnet boundaries of the data center's network, so that
the VM's IP address need not be changed to reflect its actual point
of attachment on the network from an L3/IP perspective.  In contrast,
if a VM moves to a new IP subnet, its address must change, and
clients will need to be made aware of that change.  From a VM
management perspective, management is simplified if all servers are
on a single large L2 network.

With virtualization, a single physical server can host 10 (or more)
VMs, each having its own IP (and MAC) addresses.  Consequently, the
number of addresses per machine (and hence per subnet) is increasing,
even when the number of physical machines stays constant.  Today, it
is not uncommon to support 10 VMs per physical server.  In a few
years, the number will likely reach 100 VMs per physical server.

In the past, services were static in the sense that they tended to
stay in one physical place.  A service installed on a machine would
stay on that machine because the cost of moving a service elsewhere
was generally high.  Moreover, services would tend to be placed in
such a way as to facilitate communication locality.  That is, servers
would be physically located near the services they accessed most
heavily.  The network traffic patterns in such environments could
thus be optimized, in some cases keeping significant traffic local to
one network segment.  In these more static and carefully managed
environments, it was possible to build networks that approached

scaling limitations, but did not actually cross the threshold.

Today, with the proliferation of VMs, traffic patterns are becoming more diverse and less predictable.  In particular, there can easily be less locality of network traffic as services are moved for such reasons as reducing overall power usage (by consolidating VMs and powering off idle machine) or to move a virtual service to a physical server with more capacity or a lower load.  In today's changing environments, it is becoming more difficult to engineer networks as traffic patterns continually shift as VMs move around.

In summary, both the size and density of L2 networks is increasing. In addition, increasingly dynamic workloads and the increased usage of VMs is creating pressure for ever larger L2 networks.  Today, there are already data centers with 120,000 physical machines.  That number will only increase going forward.  In addition, traffic patterns within a data center are changing.


4.  Representative Data Center Designs

   This section outlines some general data center designs and how they impact address resolution.  These designs may only approximate what happens in real data centers, but it is hoped that they can serve as a useful vehicle for describing pain points that are being experienced today in current data centers.

   Many data centers build their L2 networks using a two-tier approach consisting of access and aggregation switches.  Servers connect to access switches (e.g., top-of-rack switches) and access switches in turn are interconnected via aggregation switches.  In the following, we describe two common layouts.

4.1.  Scenario 1: L3 Terminates at the Access Link

   In Scenario 1, the L3 network extends all the way to the access switches, with the L2 broadcast domain terminated at the access switch.  All servers attached to an access switch are part of the same L2 broadcast domain and the same IP subnet.  Each access switch terminates its own L2 broadcast domain, and machines connected to different access switches are numbered out of different IP subnets. This approach works well from an address resolution perspective because the overall number of machines (physical and virtual) in a single L2 domain is relatively small, e.g., in the low hundreds.

   The main disadvantage to this scenario is that VMs cannot easily be moved from a server attached to one access switch to a server on a different access switch, as doing so requires changing the VM's IP

address, or taking additional steps at the IP routing level to ensure
that traffic continues to reach the VM at its new location, even
though its IP address no longer matches the subnet configuration of
the physical network.

4.2.  Scenario 2: L3 Terminates at the Aggregation Switch

In Scenario 2, the L3 network extends only to the aggregation
switches (or perhaps to routers that connect to the aggregation
switches).  The aggregation switches (or the routers that connect to
multiple aggregation switches) could terminate multiple distinct IP
subnets (e.g., one per VLAN) or one large IP subnet.  In order to let
hosts belonging to different IP subnets be placed under any access
switches, it is necessary for access switches to enable multiple
VLANs and aggregation switches to enable some VLANs (or subnets) over
many physical ports.  This configuration breaks the confinement of
the VLAN's broadcast domain and makes it equivalent to all the access
switches being part of the same L2 broadcast domain (and IP subnet).
Thus, this configuration allows VMs to be moved to servers connected
to other access switches, but increases the size of the L2 broadcast
domain, which can lead to difficulties outlined below.


5.  Address Resolution in IPv4

In IPv4, ARP provides the function of address resolution.  To
determine the link-layer address of a given IP address, a node
broadcasts an ARP Request.  The request is delivered to all portions
of the L2 network, and the node with the requested IP address replies
with an ARP response.  ARP is an old protocol, and by current
standards, is sparsely documented.  For example, there are no clear
requirement for retransmitting ARP requests in the absence of
replies.  Consequently, implementations vary in the details of what
they actually implement [RFC0826][RFC1122].

From a scaling perspective, there are a number of problems with ARP.
First, it uses broadcast, and any network with a large number of
attached hosts will see a correspondingly large amount of broadcast
ARP traffic.  The second problem is that it is not feasible to change
host implementations of ARP - current implementations are too widely
entrenched, and any changes to host implementations of ARP would take
years to become sufficient deployed to matter.  That said, it may be
possible to change ARP implementations in hypervisors, L2/L3 boundary
routers, and/or ToR access switches, to leverage such techniques as
Proxy ARP and/or OpenFlow infused directory assistance approaches.
Finally, ARP needs to take steps in order to flush out stale or
changed entries.  However, the existing standards do not provide
clear implementation guidelines for how to do this.  Consequently,

some implementations are "chatty" in that they just periodically
flush caches every few minutes and rerun ARP.


6.  Problem Itemization

   This section articulates some specific problems or "pain points" that
   are related to large data centers.  It is a future activity to
   determine which of these areas can or will be addressed by ARMD or
   some other IETF WG.

6.1.  ARP Processing on Routers

   One pain point with large L2 broadcast domains is that the routers
   connected to the L2 domain need to process "a lot of" ARP traffic.
   Even though the vast majority of ARP traffic may well not be for that
   router, the router still has to process enough of the ARP request to
   determine it can safely be ignored.  The ARP algorithm specifies that
   a recipient must update its ARP cache if it receives an ARP query
   from a source for which it has an entry [RFC0826].

   A common router architecture has ARP processing handled in a "slow
   path" software processor rather than directly by a hardware ASIC as
   is the case when forwarding packets.  Such a design significantly
   limits the rate at which ARP traffic can be processed.  Current
   implementations today can support in the low thousands of ARP packets
   per second.

   To further reduce the ARP load, some routers have implemented
   additional optimizations in their ASIC fast paths.  For example, some
   routers can be configured to discard ARP requests for target
   addresses other than those assigned to the router.  That way, the
   router's software processor only recieves ARP requests for addresses
   it owns and must respond to.  This can significantly reduce the
   number of ARP requests that must be processed by the router.

   Another optimization concerns reducing the number of ARP queries
   targeted at routers, whether for address resolution or to validate
   existing cache entries.  Some routers can be configured to send out
   periodic gratuitous ARPs, helping to reduce the number of ARP queries
   they receive.  The gratuitous ARP pre-populates the ARP caches on
   neighboring devices, or refreshes the "last validated" timestamp on
   such entries, reducing the number of ARP queries they send to the
   router.

   Finally, another area concerns how routers process IP packets for
   which no ARP entry exists.  Such packets must be held in a queue
   while address resolution is performed.  Once an ARP query has been

resolved, the packet is forwarded on.  Again, the processing of such
packets is handled in the "slow path".  This effectively limits the
number of ARP "cache misses" that a router can process and is viewed
as a problem in some networks today.

Although address-resolution traffic remains local to one L2 network,
some data center designs terminate L2 subnets at individual
aggregation routers (i.e., Scenario 2).  Such routers can be
connected to a large number of interfaces (e.g., 100).  While the
address resolution traffic on any one interface may be manageable,
the aggregate address resolution traffic across all interfaces can
become problematic.

Another variant of Scenario 2 has individual routers servicing a
relatively small number of interfaces, with the individual interfaces
themselves serving very large subnets.  Once again, it is the
aggregate quantity of ARP traffic seen across all of the router's
interfaces that can be problematic.  This "pain point" is essentially
the same as the one discussed above, the only difference being
whether a given number of hosts are spread across a few large subnets
or many smaller ones.

6.2.  MAC Address Table Size Limitations in Switches

L2 switches maintain L2 MAC address forwarding tables for all sources
and destinations traversing through the switch.  These tables are
populated through learning and are used to forward L2 frames to their
correct destination.  The larger the L2 domain, the larger the tables
have to be.  While in theory a switch only needs to keep track of
addresses it is actively using, switches flood broadcast frames
(e.g., from ARP), multicast frames (e.g., from Neighbor Discovery)
and unicast frames to unknown destinations.  Switches add entries for
the source addresses of such flooded frames to their forwarding
tables.  Consequently, MAC address table size can become a problem as
the size of the L2 domain increases.  The table size problem is made
worse with VMs, where a single physical machine now hosts ten (or
more) VMs, since each has its own MAC address that is visible to
switches.

In Scenario 1, the size of MAC address tables in switches s not
generally a problem.  In Scenario 2, however MAC table size
limitations can be a real issue. [xxx: do we have numbers?  For what
size L2 broadcast domains do we start seeing problems? ]

7.  Summary

This document has outlined a number of problems or "pain points"

related to address resolution in large data centers.

8.  Open Issues

   1.  The document concentrates on ARP, but the same analysis needs to
       be performed for IPv6's Neighbor Discovery.

9.  Acknowledgments

   This document has been significanlty improved by comments from Linda
   Dunbar and Sue Hares.  Igor Gashinsky deserves addition credit for
   highlighting some of the ARP-related pain points and for clarifying
   the difference between what the standards require and what some
   router vendors have actually implemented in response to operator
   requests.

10.  IANA Considerations

   This document makes not request of IANA.

11.  Security Considerations

   This documents lists existing problems or pain points with address
   resolution in data centers.  This document does not create any
   security implications nor does it have any security implications.
   The security vulnerabilities in ARP are well known and this document
   does not change or mitigate them in any way.

12.  Informative References

   [RFC0826]  Plummer, D., "Ethernet Address Resolution Protocol: Or
              converting network protocol addresses to 48.bit Ethernet
              address for transmission on Ethernet hardware", STD 37,
              RFC 826, November 1982.

   [RFC1122]  Braden, R., "Requirements for Internet Hosts -
              Communication Layers", STD 3, RFC 1122, October 1989.

Author's Address

    Thomas Narten
    IBM

    Email: narten@us.ibm.com