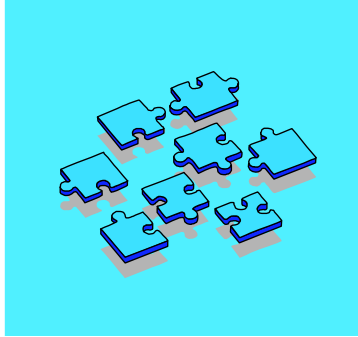# SHARE

## Open Shortest Path First (OSPF):
## An Architectural Tutorial for z/OS

*Gwen Dente, IBM Advanced Technical Skills, gdente@us.ibm.com*

**March 3, 2011 (Thursday)**
**8:00 - 9:00 AM**

**Anaheim Convention Center**
**Room 212A**

© IBM Corporation 2004 - 2011

# Abstract

- Prerequisites:
  - Knowledge of Basic Routing Concepts for Static IP Routing
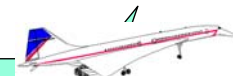  - Experience with implementing MVS or z/OS-based or z/VM TCP/IP
- Abstract:  OSPF is an ideal dynamic routing update protocol to implement in the data center to maximize application availability. It is relatively simple to design and implement ** IF ** you understand, first, the basic architecture of OSPF, and second, if you understand how to code for OSPF on the relevant platform.
  - OSPF as part of OMPROUTE is the strategic direction for implementing a dynamic IP routing protocol in z/OS CS.  The link-state algorithms of OSPF bring great advantages to an IP network over those offered by the distance vector algorithms of RIP.  As a result, many customers are choosing to move away from static routing and from dynamic RIP routing to an OSPF implementation with OMPROUTE, available since OS/390 V2R6.  To do this can be relatively simple if you understand the basic architecture of OSPF and if you understand how this architecture influences your implementation and coding choices in z/OS.
  - This session provides a tutorial of the basic concepts and considerations that need to be understood when considering an OSPF implementation.
  - The basic OSPF architectural concepts apply to any platform; some platforms have implemented extensions, but once you understand the architecture of OSPF, you can probably deal with any platform implementation.
- Acknowledgements:  Many thanks to Alan Packett and Mike Fox of the IBM Design and Development groups in Raleigh for their suggestions and several visuals.

# Agenda

1. **Why Dynamic Routing and Why OSPF?**
2. **OSPF Tutorial**
   1. General Terminology
   2. OSPF and Area Types
   3. Route Aggregation
   4. Link State Routing Protocol
   5. OSPF Packet Types
   6. Designated Router
   7. Adjacencies
   8. Flooding
   9. Route Computation
   10. Link State Advertisements (LSAs)
3. **Basic Coding Examples of OSPF with OMPRoute**
   1. URL for Full Coding Examples for OSPF with OMPROUTE
4. **Coding Enhancements:   Variables, INCLUDE Statements**
5. **Bibliography**

Warning: Fast-moving presentation!
This is a pre-requisite for the session on Coding and Designing for OSPF.

© IBM Corporation 2004 - 2011

---

1. This presentation moves at a brisk pace.  Do not get discouraged -- we know you may not be able to absorb everything during the time allocated to this topic, but that is why we have included thorough notes!!  The notes make it easy for you to review this presentation at your leisure at a later point in time.
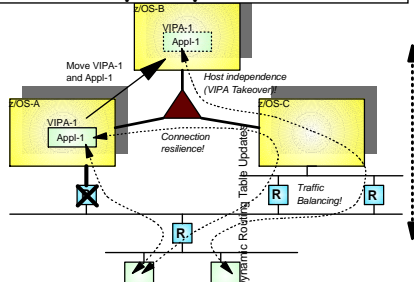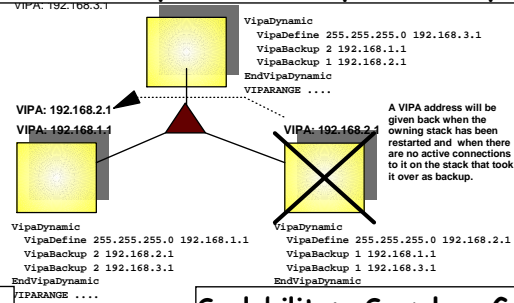
# 1. Why Dynamic Routing?

Where does Dynamic Routing Help?

1. IP Non-disruptive Reroute
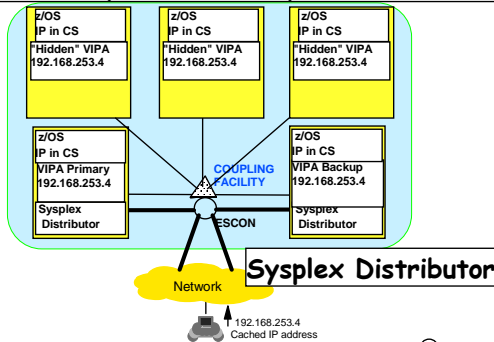2. VIPA Dynamics with VIPA Takeover (V2R8) and non-disruptive Giveback in V2R10.
3. Scalability to add or subtract OS/390 images "seamlessly" - Dynamic routing allows route discovery to the new image.
4. "Sysplex Distributor" is the name given to an enhancement in V2R10 which builds on VIPA Takeover and Sysplex technology to perform IND-like function on an OS/390 TCP stack. This function will require additional storage and CPU capacity on the distributing stack, but will give additional benefits over existing IND implementations:
   1. WLM can be consulted real-time on every connection request, rather than polling periodically, so that up-to-date load information is used to select the target node/stack.
   2. Beginning in z/OS V1R2, Sysplex Distributor may be combined with the MultiNode Load Balancing (MNLB) of CISCO by coding a Service Manager in the Sysplex Distributor node.
5. There are distance vector dynamic routing protocols, like RIPV1 and RIPV2.
   1. Every 30 seconds, a router will send its full set of distance vectors to neighboring routers, i.e., a full table worth of routing data. When a router receives a set of distance vectors, it will have to recalculate the shortest path to each destination. A distance of 16 constitutes an invalid/unreachable destination in RIP.
6. There are link-state routing protocols, as with OSPF.
   1. Every router has a full map of the entire network topology, that is, all routers have an identical copy of the network database. Changes to the network topology (due to link-state changes or outages) cause a database update for those changes only (not for an entire table of data) to be propagated throughout the network. When a change is received at a router, the router recomputes its shortest path to all destinations. There is virtually no limit to the network design with OSPF.
      1. The next page's notes show you a table comparing various dynamic routing protocols.

# Comparisons of Dynamic Routing Protocols

| Characteristics | RIP V1 | RIP V2 | OSPF | EIGRP |
|---|---|---|---|---|
| Algorithm | Distance vector | Distance vector | Link state | Distance |
| Network load (note 1) | May be high | May be high | Low | Low |
| CPU processing requirement (note 1) | Low | Low | May be high | Low |
| IP network design restrictions | Many | Some | 0 - Few | 0 - Few |
| Convergence time | LT/EQ 180 sec. | LT/EQ 180 sec. | 0 - immed. | 0 - 5 sec. |
| Information exchange | Broadcast | Broadcast / Multicast | Multicast | Multicast |
| Multipath Outbound | No | No | Yes | Yes |
| Support on OS/390 or z/OS | TCP/IP V2 | TCP/IP V3 | OS/390 V2R6 | N/a |

Note 1: Depends on network size and (in)stability.

1. These three are the most widely used interior routing protocols today.  Most TCP/IP products support RIP Version 1.
2. Link-state:  only the information and state about the links is transmitted through the network; therefore network traffic is low.
    1. Routing table must be computed, leading to higher CPU consumption for OSPF Dykstra information.
3. Distance Vector:  entire route table is transmitted through the network; network traffic load is high.
    1. No need for Routing table computation; therefore lower CPU consumption for distance vector algorithms.
4. There are a few, very special situations, where a distance vector algorithm relies on the ability to 'count to infinity' in order to fully converge after a topology change.
5. Most current RIP implementations have implemented techniques to minimize the probability for those situations, but they cannot be fully eliminated.
6. When 'counting to infinity' occurs, the convergence time may be higher than 180 seconds for RIP.
7. RIP V1 is class-full
8. RIP V2 is class-less.

# 2. OSPF Tutorial

1. Both RFC 2328 and RFC 1583 define OSPF.  However, CS for OS/390 OMPROUTE supports 1583+ and not RFC2328. Another RFC, #1587, describes a special type of OSPF area called the Not So Stubby Area.  OMPROUTE does not support RFC 1587 or the draft document that proposes additional capability for NSSAs:  draft-ietf-ospf-nssa-update-10.txt.
2. This  tutorial focuses on OSPF as described in RFC 1583.
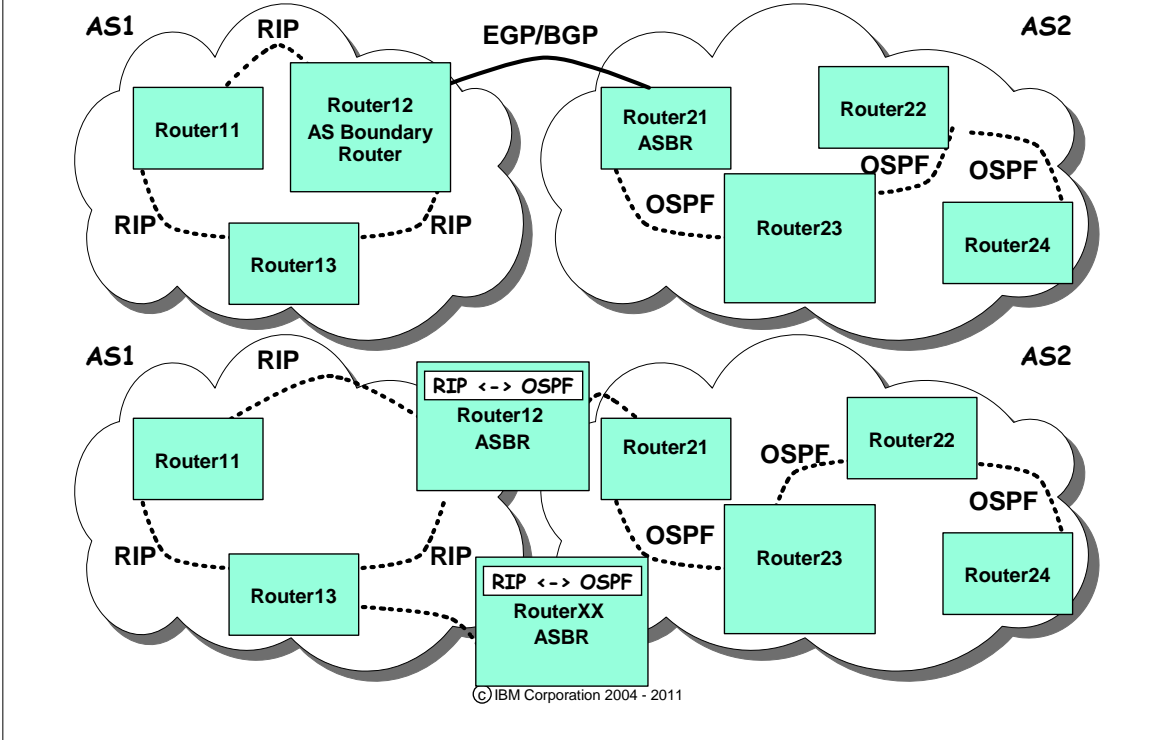
# General Terminology: AS, BGP, ASBR



| | |
|---|---|
| **AS** = Autonomous System | **A collection of routers running a common routing protocol** |
| **IGP** = Interior Gateway or Routing Protocol | **A common routing protocol that is used among routers of an autonomous system** |
| **EGP** = Exterior Gateway Protocol or **BGP** = Border Gateway Protocol | **Protocol that may be used among multiple Autonomous Systems; alternatively use a router that speaks both IGPs. Not supported in CS for OS/390.** |
| **ASBR** = AS Boundary/Border Router | **A router that speaks multiple IGP routing protocols and can import/export or redistribute routes among differing IGPs, or a router that uses BGP to communicate with another router using BGP. An ASBR may sit in each AS if it is using BGP.** |

© IBM Corporation 2004 - 2011

1. We begin here a discussion of the terminology of OSPF (Open Shortest Path First). We start with a description of only a few terms and proceed to more explanations on subsequent pages.
2. Autonomous System (AS)
   1. A group of routers running a common routing protocol. The routers connect multiple networks.
      1. Within an AS, all the routers run the same IGP.
3. ASs are interconnected by means of routers running what were formerly called External Gateway Protocols (EGPs). An enhancement to EGP, called Border Gateway Protocol (BGP), was introduced in the 1990s to provide routing loop avoidance and other improvements. The routers that interconnect Autonomous Systems are called AS Boundary Routers or AS Border Routers. (NOTE: AS Boundary Router is the preferred term so as not to be confused with another router type called an "Area Border Router.")
   1. OMPROUTE in CS for OS/390 does not support EGP or BGP; it uses another technique to interconnect Autonomous Systems. (See next page.)
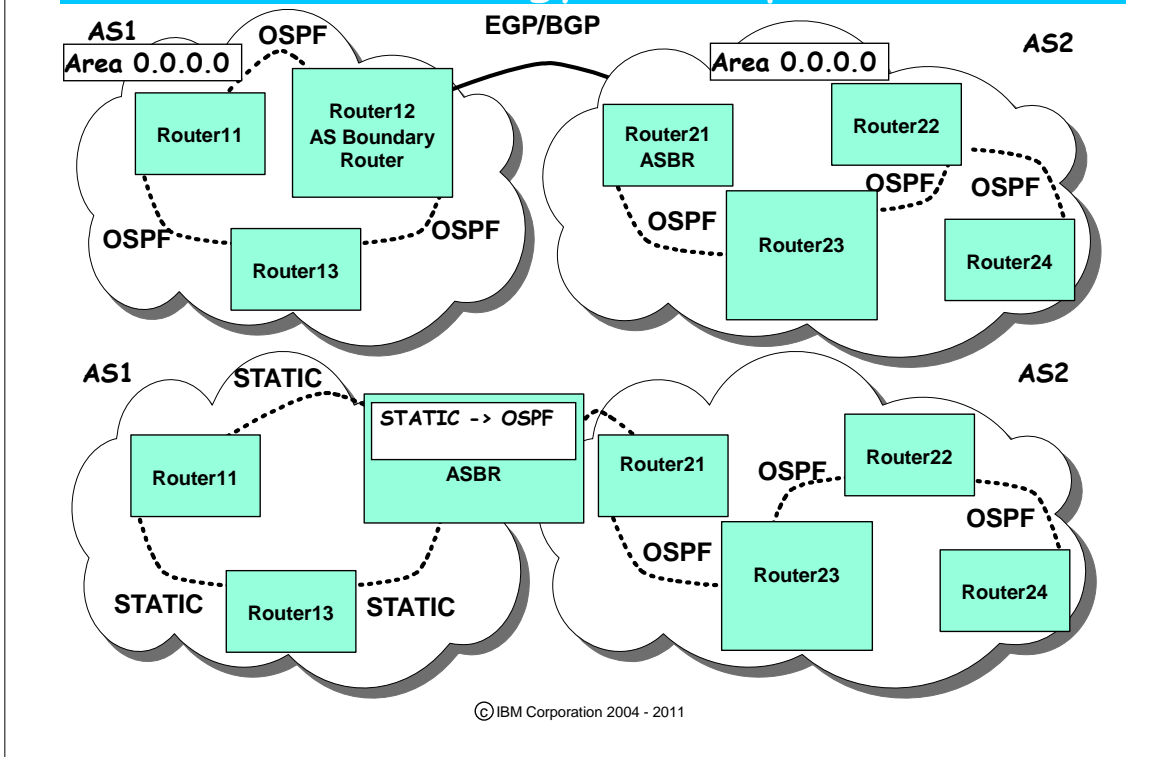
Terminology: Example 1

1. Common Interior Gateway Protocols in use today are:
   1. RIP
   2. OSPF
   3. EIGRP
   4. IGRP
   5. HELLO protocol, a protocol different from the HELLO protocol of OSPF
2. There are various flavors of BGPs in use today.  BGP may interconnect two ASs of different Protocols (as you see here with RIP and OSPF) or two ASs of the same Protocol (e.g., OSPF ASx with OSPF ASy).  We show you an example of this type of configuration on a subsequent page.  OMPROUTE does not support the use of EGPs or BGPs, so they will not be covered in this pitch.
3.  More commonly a router that communicates in several routing protocols is used to span the Autonomous Systems.  A single AS Boundary Router that "redistributes" or "imports"  the routing protocol of one IGP into the routing protocol of another can interconnect two AS's.  This is the approach used by OMPROUTE when OMPROUTE represents an ASBR.
   1. The diagram  in the bottom half of the visual depicts such a scenario.
      1. In this multiprotocol scenario, RIP routes with their metrics are imported into OSPF; OSPF routes with their Costs are imported into RIP.
4. Since metrics and costs mean different things, you must establish a route precedence methodology to determine which type of route should be preferred in the case of multiple equal-cost routes of different flavors.
5. It goes beyond the scope of this tutorial to examine these issues.   For more information, please consult the IP Configuration Guide and the Information APARs on RETAIN regarding "route precedence" and "Type 1" and Type 2" external routes.
6. If you have redundant ASBRs, as you see in the bottom half of this visual, configure them to include route tags so that the receiving AS does not attempt to re-distribute the same routes back into the AS  from which they came.  The route tag is set in the LSA Type 5 (External) packet.  Route tags are not buit by OMPROUTE, but they are passed on with the OSPF packets for the benefit of other ASBRs that may need to interpret them.

Terminology: Example 2

1. In this example you see two diagrams.
    1. The one in the top half of the visual shows how OSPF can be the IGP in each AS that is interconnected with AS Boundary Routers (ASBRs).
        1. In such a scenario, a network must employ ASBRs in each AS which is communicating with another using a Border Gateway Protocol. The single ASBR model that spans both AS's does not apply to a scenario in which multiple AS's all use OSPF as the IGP.
    2. The second diagram in the bottom half of the visual depicts an Autonomous System (AS) that uses Static Routing definitions to interface with an Autonomous System (AS) with OSPF dynamic routing protocol.
        1. A router that sits between such AS's is as much an ASBR as any other ASBR we have seen in previous diagrams.

# General Terminology: Area, ABR

**Area**

A subdivision of an AS
- reduces the router overhead required to keep all routers in the AS informed of the network topology;
- saves storage, memory, cycles in routers that do not span the areas

**ABR** = Area Border Router

A router that spans multiple areas of an AS

© IBM Corporation 2004 - 2011

---

1. We now continue our discussion of more OSPF terminology.
2. We start with the concept of "OSPF Area."
   1. If an AS is large, it may be subdivided into multiple Areas. An area is identified by a 4-octet number. You see this concept depicted in AS2: there is an Area 0.0.0.0 and another Area 1.1.1.1.
      1. Areas are interconnected by means of Area Border Routers (ABRs) which run Interior Gateway Protocols (IGPs). In our example, Router 23 represents the ABR between Area 0.0.0.0 (also called simply "Area 0"), the backbone area, and Area 1.1.1.1, a non-backbone area.
      2. The topology of an area is hidden from the rest of the Autonomous System, thus reducing routing overhead. Fewer routing updates need to be sent and smaller routing trees need to be computed and maintained, leading to a reduction in storage and CPU consumption.
      3. We discuss the different types of areas on subsequent pages.

# OSPF and Areas: Overview

Area 5.5.5.5

Rm

OSPF

Area 4.4.4.4
Totally Stubby
Area

OSPF

Rd

ASn

Virtual Link

Rc

OSPF

Re

OSPF

OSPF

Rg

OSPF

Ra

OSPF

Rb

Area 0.0.0.0
Backbone Area

OSPF

Rf

Area 1.1.1.1

Area 2.2.2.2

OSPF

Virtual Link

Rh

Rk

OSPF

Rj

ASx

RIP

Area 3.3.3.3
Stub Area

Ry

Rz

© IBM Corporation 2004 - 2011

1. This diagram depicts an OSPF Autonomous System (ASn) and a RIP Autonomous System (ASx).
2. There are many details in this diagram that we will discuss separately on the following pages in order to explain the OSPF protocol.

# OSPF Areas: Backbone & Non-Backbone



Area 5.5.5.5 Non-Backbone — Rm — OSPF

OSPF — Ra — Virtual Link — Rb

Area 2.2.2.2 Non-Backbone

Area 0.0.0.0 Backbone Area — Rc — Re — OSPF

OSPF — Rf — Rk

Area 1.1.1.1 Non-Backbone — Rg — OSPF

ASn

ASx

© IBM Corporation 2004 - 2011

1. A backbone area is a requirement in an OSPF AS. Depending on the size of the routing trees, OSPF can be compute-intensive. It is common to populate only non-backbone areas with application servers.
2. Areas are identified by an ID that is four octets in length. The backbone area is always assigned ID 0.0.0.0. Each Area number must be unique within an AS.
3. All non-backbone areas must be connected to the backbone.
4. In this subset of the OSPF AS we depict a backbone area (0.0.0.0) and three non-backbone areas (1.1.1.1, 2.2.2.2, 5.5.5.5). Note how Areas 1.1.1.1 and 2.2.2.2 are **physicall**y connected via Area Border Routers Rf and Rb to the backbone. Note, however, that Area 5.5.5.5 is not physically connected to the backbone; it is physically connected to Area 2.2.2.2, a non-backbone area. This appears to violate the rule that an area must be connected to the backbone. Yet this configuration is a valid OSPF network design. Why? Because we have **logically** or **virtually** connected the non-backbone area to the backbone! We have accomplished this feat by defining a **virtual link** between the ABR Ra and the ABR Rb, which is **physically** connected to the backbone area.
5. Although not depicted, virtual links are also used to connect discontiguous areas to make them appear logically contiguous.
6. How big should an area be? This is a vendor-specific decision. It depends on the capacity of the routers to handle the link-state database and the computation of the routing trees. Network stability, number of links, IP network address assignment design all influence the size of the Area.
7. Some customers have areas up to 200 routers in size; others have areas of 50 routers in size.

# Advertising: To Backbone & Non-Backbone

Area 5.5.5.5
Non-Backbone

<<<OSPF intra-area destinations,OSPF inter-area Summaries, RIP destinations>>>

Rm

ASn

OSPF

OSPF

Rc

OSPF

Re

Ra

Virtual Link

Rb

Area 0.0.0.0
Backbone Area

Rg

OSPF

Rf

Area 1.1.1.1
Non-Backbone

Area 2.2.2.2
Non-Backbone

OSPF

Rk

ASx   RIP

Ry      Rz

© IBM Corporation 2004 - 2011

1. We have added back into this picture the RIP AS, "ASx." It is connected to the backbone area via ASBR Rk. The purpose of this visual is now to describe how routing advertisements are shipped throughout the areas so that all destinations -- both OSPF and RIP -- may be reached from anywhere in the network.
2. All routers within an area ("intra-area routers" or "interior routers") maintain a copy of the topology of the area, including information about intra-area OSPF destinations. They also know about inter-area destinations due to advertisements sent to them by the ABRs. They also know about how to reach RIP destinations, either explicitly or through a Default route.
3. The Area Border Router maintains a copy of the database for each area to which it is connected.
   1. An area border router summarizes all OSPF links into other areas and advertises these to adjacent areas.
   2. An ABR can be customized to use address ranges for route summarization. We depict route summarization later in this presentation.
4. ASBRs ( in this case, Rk) import (or "redistribute") external destinations from other AS's into OSPF advertisements and originate what are called "External Link State Advertisements" known as Type 5 LSAs. Our ASBR, Rk, can thus originate information about RIP destinations and send this information into the OSPF AS.
   1. External Link State Advertisements (Type 5 LSAs) flow freely across all OSPF areas, except for Stub Areas.
   2. (External Links are RIP, Static, and Direct.)
   3. ASBRs may also be configured to send or "originate" default routes into the area or into the other AS.
5. All the OSPF Areas depicted in this subset diagram have knowledge of routes to the RIP Autonomous System.
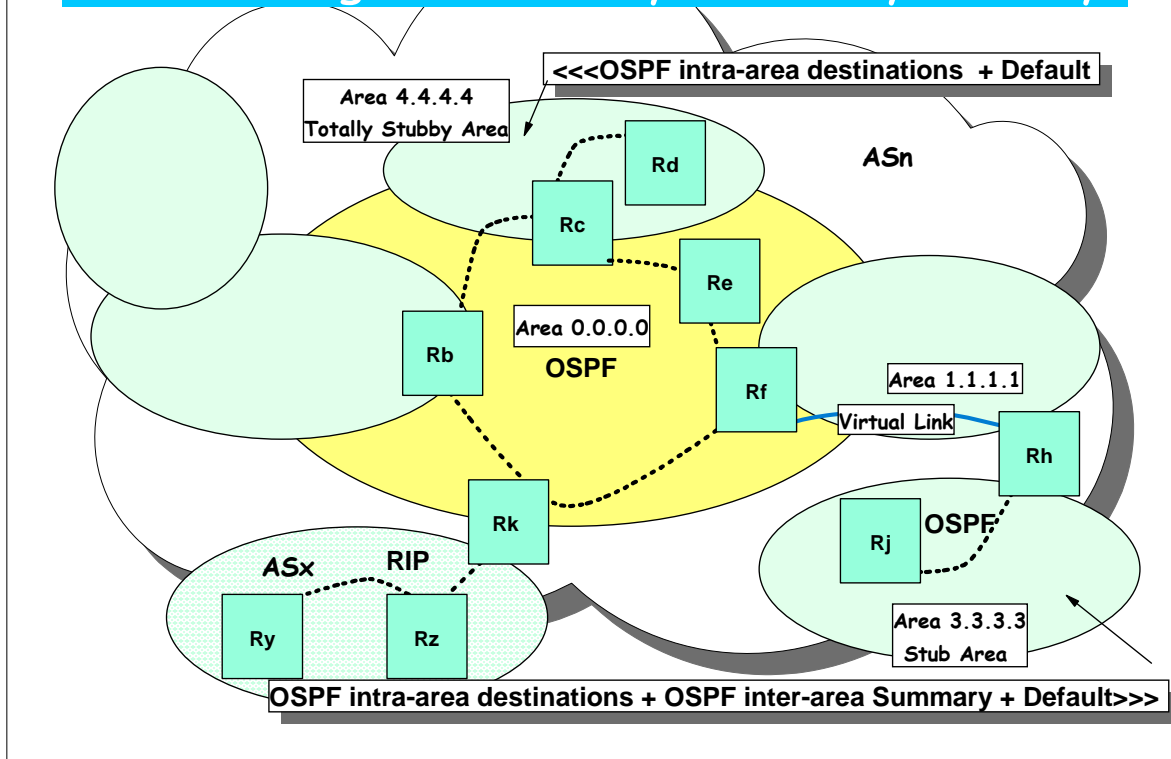
OSPF Areas: Stubby & Totally Stubby

1. This subset of our original diagram now focuses on special non-backbone areas called "Stub Areas" and "Totally Stubby Areas."
2. Stub Areas minimize storage and CPU utilization at nodes that are part of the Stub Area because they maintain less knowledge about the topology of the AS than do other types of non-backbone routers. They maintain knowledge only about intra-area destinations, summaries of inter-area destinations, and a default route in order to reach external destinations in other AS's.
    1. A Stub Area is defined with STUB=YES in the Area definition.
    2. A stub area can be adjacent to the backbone (Area 4.4.4.4) or not adjacent to the backbone (Area 3.3.3.3, which is adjacent to a non-backbone area).
    3. Stub Area 3.3.3.3 is connected via a virtual link to the backbone area.
        1. As you already know, every area must be adjacent to the backbone. If non-adjacent physically, there must be a virtual link to the backbone. In our diagram, Rh becomes part of the Backbone because of the virtual link.
        2. Stub Areas can be the end-points of a virtual link, but they cannot be on the intermediate path of a virtual link.
3. A Totally Stubby Area receives less routing information than a Stub Area; it receives only Default Routes. Designing with a Totally Stubby Area minimizes the compute-intensive operations necessary to build routing trees and minimizes the storage requirements for maintaining the topology database. A Totally Stubby Area is coded inside CS for OS/390 with an Import_Summaries=NO on the AREA statement.
4. CS for OS/390 OSPF cannot be an NSSA (Not So Stubby Area), an area that is not discussed in this brief tutorial. (Such an area sends Type 7 LSAs, not supported in CS for OS/390 or z/OS.) An NSSA is described in IETF RFC 1587 and in the IETF draft document: draft-ietf-ospf-nssa-update-10.txt. CS in OS/390 or z/OS can be a Stub Area or a "Totally Stub Area," as we will see later in this series of presentations.
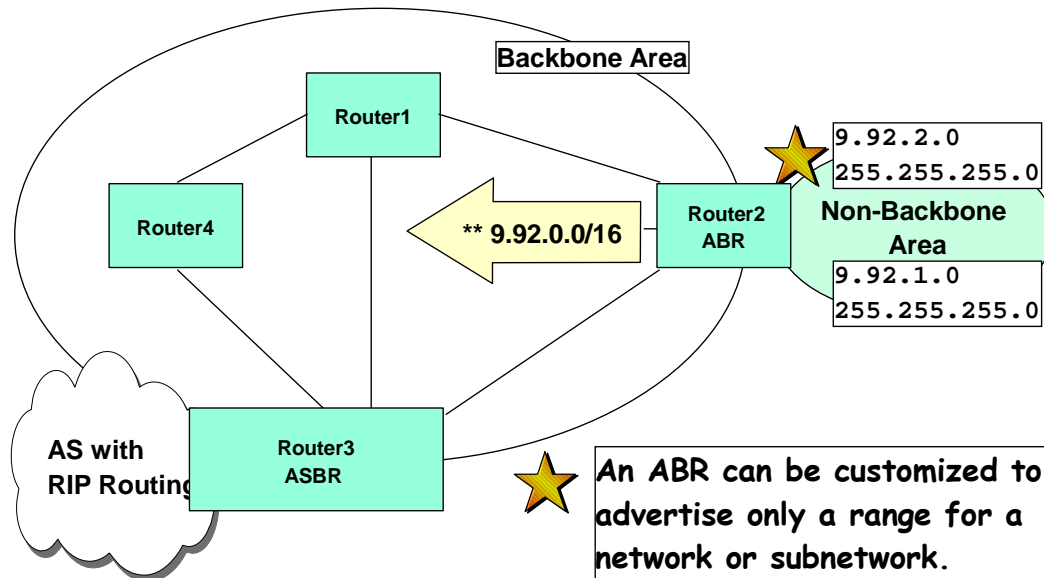
# Advertising: To Stubby & Totally Stubby

1. This subset of our original diagram now focuses on the means by which Stub Areas learn of intra-area, inter-area, and external destinations.
   1. A Stub Area receives Type 3 Summary LSAs for inter-area destinations, including the Default Route packaged inside a Type 3 Summary LSA.
   2. Stub Areas can have multiple default routes from different attached Area Borders. (One for each.)
   3. External Link State Advertisements (e.g., RIP, static, Direct routes described in Type 5 LSAs) cannot flow into a Stub Area. The Stub Areas can reach hosts in the RIP AS only by means of a default route made known to them by their Area Border Router.
2. What do these statements mean? Recall the function of the ABRs we just looked at:
   1. An area border router summarizes all OSPF links into other areas and advertises these to adjacent areas.
   2. An ABR can be customized to use address ranges for route summarization. We depict route summarization later in this presentation.
   3. You see in this subset diagram that an Area Border Router can interface with a Stub Area. So, an ABR summarizes inter-area destinations and sends Link State Advertisements into Stub Areas in the form of Type 3 Summary LSAs. It sends Default destinations into the Stub Areas, also in the form of Type 3 Summary LSAs. It does not, however, forward any External Destinations it has learned from ASBRs into the Stub Areas. (The External Destinations are described in Type 5 LSAs.) It does not need to, because the Stub Area is given a Type 3 Summary LSA Default way to get to any part of the network it may need to reach. The Area Border Router says, "I'm not telling you the external links, so here's a default route for you if you need to reach an external route."
   4. So, in summary, all OSPF links and Type 3 Summary LSA, including the LSA for a Default, can flow into a Stub Area. (More about Type 3 Summary LSAs later....)
3. A Totally Stubby Area receives only Default Routes. That is, just as with a Stub Area, it receives no External Routes (Type 5 External LSAs). But unlike a Stub Area, it also does not receive Type 3 Summary LSAs to reach inter-area destinations, with one exception: the Default Route that is packaged inside a Type 3 Summary LSA. It is coded inside CS for OS/390 with an Import_Summaries=NO on the AREA statement.
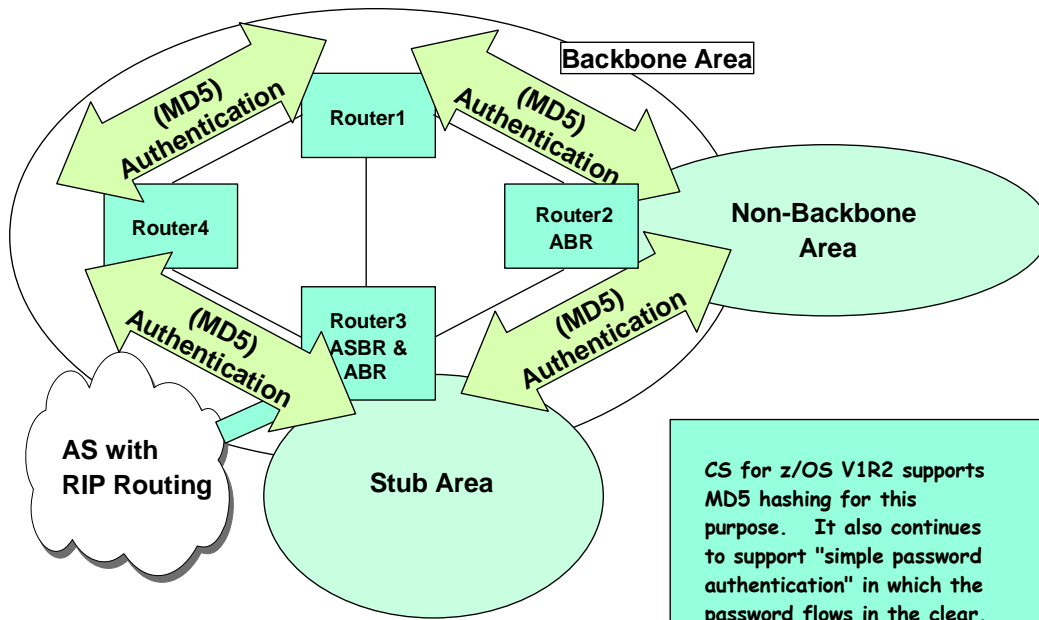
Route Aggregation

** = Router2 Advertising 9.92.0.0/16

Backbone Area

Router1

Router4

** 9.92.0.0/16

Router2 ABR

9.92.2.0
255.255.255.0

Non-Backbone Area

9.92.1.0
255.255.255.0

AS with RIP Routing

Router3 ASBR

An ABR can be customized to advertise only a range for a network or subnetwork.

© IBM Corporation 2004 - 2011

1. In a well-designed network in which areas have been carved out according to the networks or subnetworks they serve, the ABRs can be customized to advertise a range, thus reducing the number of LSAs for inter-area routes.
2. The term "Classless InterDomain Routing" (CIDR) or supernetting is also applied to this type of Route Aggregation.
3. Within the subnet, the actual subnet mask in use may be quite different, for example, 255.255.255.0. The advertised value for the mask (255.255.0.0) specifies the mask that is common to all the subnets within 9.92.0.0 (9.92.1.0 and 9.92.2.0) when the actual mask is 255.255.255.0.
   1. CS in OS/390 allows you to specify a RANGE and "Advertise=YES/NO" in order to provide Route Summarization in the LS Advertisements.
4. Without Route Aggregation (Ranges), the LSA Summaries that are sent by Router 2 about the Non-Backbone Area are:
   1. 9.92.1.0
   2. 9.92.2.0
5. With Route Aggregation (Ranges), the LSA Summary that is sent by Router 2 about the Non-Backbone Area is:
   1. 9.92.0.0
6. In aggregating routes, you must take care to choose an aggregated subnet mask that allows for a contiguous clustering of subnets in an area behind one or a set of ABRs. In other words, you cannot use a subnet mask for aggregation that will make the aggregated subnets look disjointed. So ... if you had 9.92.5.0 and 9.92.6.0 on the other side of the network behind ABR4, you could not use route aggregation if you continued to use an aggregated mask of 255.255.255.0 for the "cluster" or "pocket" of subnets numbered 9.92.1.0 and 9.92.2.0. But ... if you changed the aggregated mask to 255.255.252.0, then you could aggregate the subnets behind ABR2 as 9.92.0.0/22 and the subnets behind ABR4 as 9.92.4.0/22. The "real " subnet masks in both disjointed areas could remain 255.255.255.0.
7. Some implementations call Route Aggregation "Route Summarization" or "Summarizing Routes." Do not confuse "route summarization" with LSA Summaries, as LSA Summaries may contain either both of the non-aggregated destinations from our example (9.9.2.20 and 9.92.1.0) or may contain the aggregated destination (9.92.0.0).

# Authenticating Routers in an Area

**Backbone Area**

Router1

Router4

Router2 ABR

Router3 ASBR & ABR

**(MD5) Authentication**

**(MD5) Authentication**

**(MD5) Authentication**

**(MD5) Authentication**

**Non-Backbone Area**

**AS with RIP Routing**

**Stub Area**

CS for z/OS V1R2 supports MD5 hashing for this purpose.  It also continues to support "simple password authentication" in which the password flows in the clear.

© IBM Corporation 2004 - 2011

1. Within an area the link advertisements can be authenticated so that a "rogue" router  may not send routing advertisements that should not be accepted.
   1. Only trusted routers are allowed to introduce any changes.
   2. This capability protects more against the accidental introduction of a router into the area than it does against the malicious insertion of an alien router into the area.
2. zOS V1R2 Communications Server adds MD5 support.  MD5 is not supported in V2R6/8 or V2R10, although there is a security mechanism with "simple password authentication."
3. IBM provides a utility, "pwtokey," that can be used to generate an MD5 authentication key.  However, be aware that some routers -- Cisco being one of them -- may generate keys for MD5 authentication using nonstandard methods, so that care must be taken to ensure that Cisco and OMPROUTE define the same key.  The character value entered on the Cisco side of a network connection must be converted to HEXADECIMAL in ASCII format on the IBM side of the connection and must not have been generated with the "pwtokey" utility.  Please see examples of coding in the OSPF presentation that covers OMPROUTE coding for OSPF.
4. The key id is a one-byte constant that identifies the key for MD5 authentication.  You should really consider it to be part of the key, and it must match on both sides of a connection.  Some platforms (including Cisco) support multiple keys.  The key id identifies which key is being used.  IBM supports only one key, but we provide the key id for compatibility, permitting it to be defined to match what Cisco is using, for example.

# Link State Cost

<--Cost=1

Cost=3-->

Cost=2

Cost=3

Cost=1

Cost=2

**Router1** — C, B, A

**Router4** — J, K, L

**Router2** — G, H, I

**Router3** — D, F, E

## Area Topology

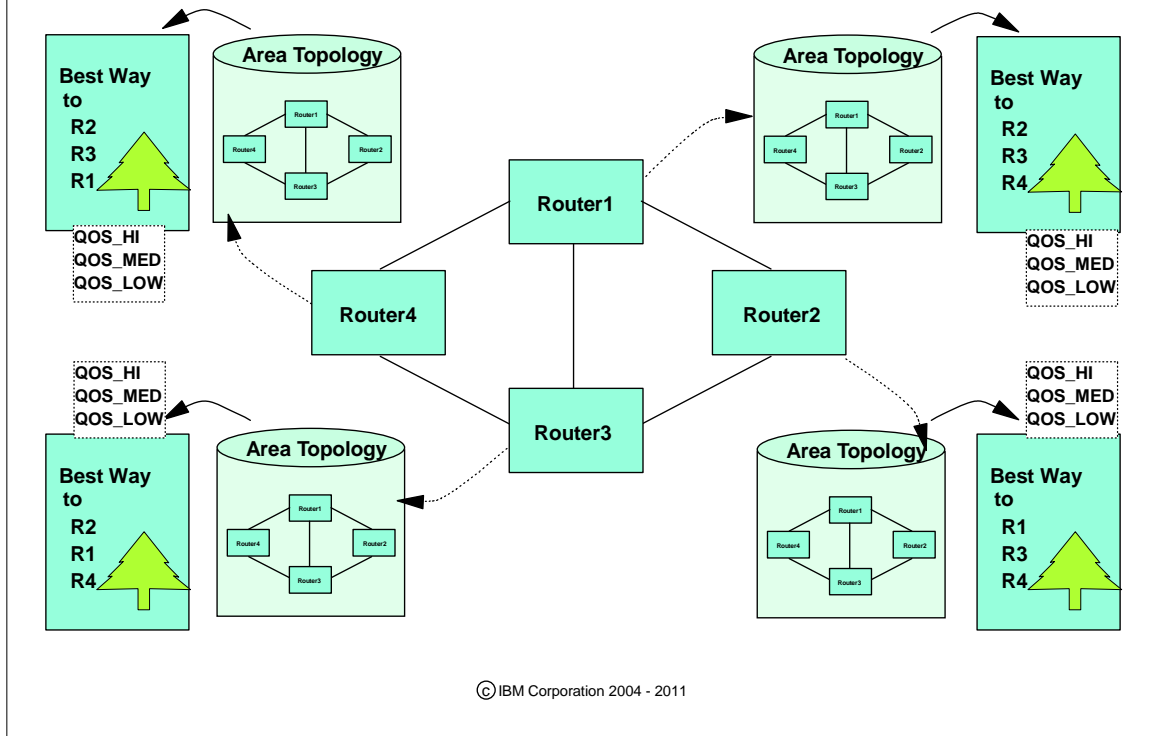| From | To/Cost ------------------------------> | | |
|------|------------|------------|------------|
| R1   | R4J - 1    | R2G - 2    | R3D - 3    |
| R2   | R1B - 2    | R3E - 2    | R4K - 3    |
| R3   | R2I - 2    | R4L - 1    | R1A - 3    |
| R4   | R1C-3      | R3F - 1    | R2H - 3    |

**Identical copy of database in each router of the Area**
**Links are really to IP Addrs/Subnets - not to Routers**
**(A, B, C, and so on are IP Addresses.)**

© IBM Corporation 2004 - 2011

1. Every link is associated with a Cost or metric that is used to compute the route to any given destination.
   1. The lower the cost, the better the route.
   2. The cost may be assigned according to any of several metrics: largest throughput, lowest delay, lowest cost, highest speed, best reliability, and so on.
2. Each router has a map of the topology of the Area. The map contains the cost of each link in the area.
   ★ The diagram/table has been simplified. Although it implies that the links extend to Routers, in fact they extend to IP addresses or Subnets in a Router. From Router 1 (R1) there is a link to one of the IP Addresses in R4 with a cost of 1. There is a link to one of the IP addresses in R2 with a cost of 2 and one to an IP address in R3 with a cost of 3.
3. The network administrator can assign a cost to every link in the network.
   1. Alternatively, some implementations of OSPF use the link speed to assign a default cost.
   2. Although not a requirement, it is usual to assign the same cost to a link in both directions. Consistency in the network helps to avoid unpredictable routes between two points.
      1. You may assign different costs in either direction if you are intending to influence traffic flow inbound and outbound to Router1.
      2. Coordinate link costs with the other router implementers in your network.
   3. In CS/390, the default link cost is 1.
4. Router 1 can reach Router 3 in three ways
   1. via Router 2 for a combined cost of 4
   2. via Router 4 for a combined cost of 2
   3. directly to Router 3 for a cost of 3

Basics of Link-State Routing Protocol

© IBM Corporation 2004 - 2011

1. Each router maintains a database describing the topology of the area.  The database is identical at each router.
   1. Each record in the Database represents one link in the network.
   2. A cost is associated with the link.  This cost is platform-specific
      1. Some platforms assign costs according to throughput or link speed.
      2. Other platforms require that the costs be assigned manually.
      3. A route to a destination is the sum of the costs of all links to the destination.  A  lower cost represents a better route.
         1. There is virtually no limit on the total cost of a route, although in reality a cost of 65534 is the highest metric that is available to a valid route.  (65535 indicates infinity, meaning that the route is not available.)  (RIP, on the other hand, has a limit of 15 hops/metric count for valid routes; 16 represents an unreachable destination.)
2. Each router uses the flooding protocol to inform the other routers in the AS of the state of its interfaces and of its reachable neighbors.
3. External routes, that is, routes learned from RIP, learned statically, or otherwise, are maintained separate from the OSPF topology.
   1. OSPF and External routes are communicated throughout the Autonomous Area unless there are Stub Areas.
      1. Stub Areas do not receive External Link advertisements; they reach areas identified in External Link advertisements by sending traffic to a default route through an Area Border Router.
4. Each router computes a "tree" of best routes to each location using itself as the root.
   1. Although OMPROUTE in CS for OS/390 or CS in z/OS currently (as of z/OS V1R2) does not compute a "tree" of best routes for each Quality of Service (QOS), the RFCs for OSPF do define this capability.

# Originating LSAs: Overview

**IBM**

A = Originating Intra-area OSPF Routing Info
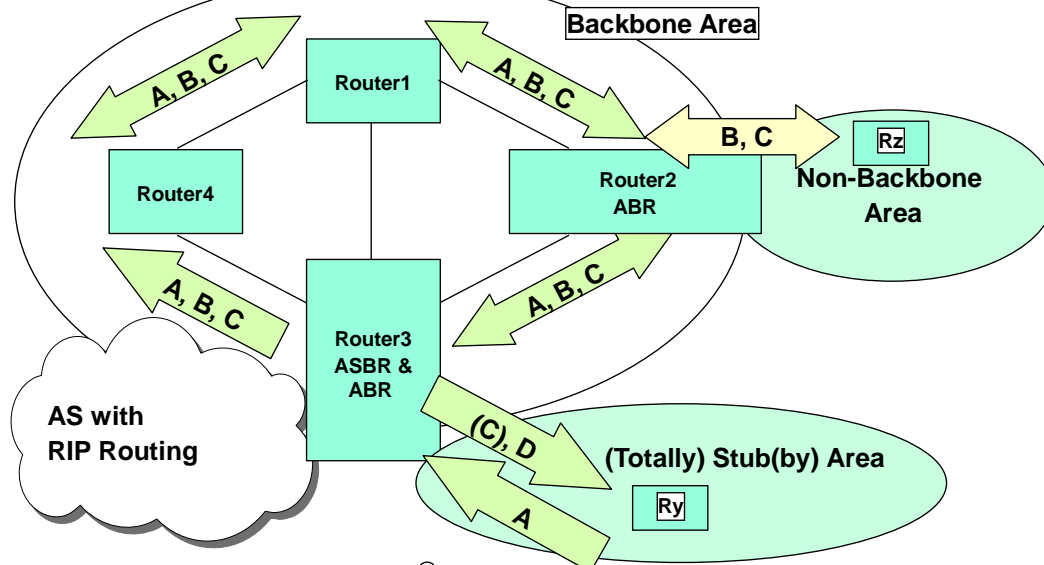B = Originating External Routing Info
C = Originating Inter-area Summary OSPF Routing Info
D = Originating Default Routing Info

**Backbone Area**

Router1

Rou...

Router2
ABR

A, C

A

Rz

**Non-Backbone Area**

A

Router3
ASBR &
ABR

A, B, C

A, C

A, B, C

**AS with RIP Routing**

A, C, D

**Stub Area**

Ry

A

An ABR can be advised by the Stub Area to omit the Type 3 Summary LSAs into the Stub Area= "Totally Stubby"

©IBM Corporation 2004 - 2011

1. This visual uses a legend to indicate what types of Link State Advertisements are sent to neighbors. Although there are actually five different types of LSAs in OSPF, the legend condenses the concepts into four groups in order to explain simply how routers synchronize their databases.
2. The whole idea of using link state routing protocols is to maintain a synchronized copy of the link state database in all nodes of an area.
3. Routers either originate route advertisements or they propagate the route advertisements received from other routers. The latter is called "flooding."
4. A routing advertisement, called a Link State Advertisement (LSA) is sent under any of the following conditions:
    1. A new link is activated or deactivated
    2. The metric of an existing link changes
    3. Every 30 minutes to refresh the databases of the neighbors
5. The ASBR (Router3) has redistributed the RIP Routes into OSPF. The redistributed (imported) routes are communicated to the other routers in the OSPF Area as External Link State Advertisements. OSPF routes are communicated as any of several types of OSPF LSAs.
6. External LSAs are not sent into a Stub Area. The Area Border Router (ABR, also Router 3) informs the Stub Area of a default route to reach any areas that are connected via RIP routes. OSPF route advertisements are also sent into the Stub Area, but they are summarized (as TYPE 3 LSAs) to indicate all the destinations that can be reached by the ABR. If the Stub Area is to be a "Totally Stubby Area," it should specify "IMPORT SUMMARIES=NO" in its AREA statement.
    1. NOTE: In reality the OSPF protocol provides for two types of Summary LSAs: Type 3 Summaries for IP destinations in other areas and Type 4 Summaries about ASBRs in the Area. LSA Type 4 Summaries pertaining to ASBRs are never advertised into Stub Areas.
7. Router 2, an ABR between the Backbone Area and the Non-Backbone Area, sends OSPF Route Summaries and External Routes into the Non-Backbone Area depicted.
8. The Stub Area could also originate a DEFAULT route, which would mean that it would send a "D." However, this would not be a good network design; the Stub Area should receive its default route from the ABR.

## Flooding LSAs: Overview

A = Flooding Intra-area OSPF Routing Info
B = Flooding External Routing Info
C = Flooding Inter-area Summary OSPF Routing Info
D = Flooding Default Routing Info

**Backbone Area**

A, B, C

Router1

A, B, C

B, C → Rz
**Non-Backbone Area**

Router4

Router2
ABR

A, B, C

A, B, C

Router3
ASBR &
ABR

**AS with RIP Routing**

(C), D

A

**(Totally) Stub(by) Area**
Ry

© IBM Corporation 2004 - 2011

1. Once databases have been synchronized with neighbors, Route Advertisements are Flooded throughout the area.  The synchronization of databases is also called "bringing up adjacencies."
2. Each LSA has a timer on it.  It is flooded again after 30 minutes in order to refresh the database.
3. Recall that the Stub Area can be customized so as to receive only the Default Route.  Omitting the inter-area route advertisements (Type 3 Summary LSAs) makes the Stub Area a "Totally Stubby Area."  The Stub Area "AREA" statement simply needs to  indicate "IMPORT SUMMARIES=NO."
    1. NOTE:  A Stub Area is never sent Type 4 Summary LSAs.  (Type 4 LSAs indicate LSAs to reach ASBRs.)

# How OSPF Operates: Details

- ➤ Discover the Neighbors
- ➤ Elect a Designated Router (DR)
- ➤ Form Adjacencies
- ➤ Synchronize the DataBases
- ➤ Compute the Routing Table (Tree)
- ➤ Advertise the Link States

Hello Protocol

Exchange Protocol

Flooding Protocol

1. The OSPF protocol has its own protocol number in IP (number 89). It uses Raw Sockets and so, unlike TCP or UDP, does not have the concept of "Port." There is no specified port for OSPF.
2. The OSPF protocol is designed to align databases at all nodes within an area. This synchronization of databases is called "bringing up adjacencies."
3. The OSPF protocol consists of three subprotocols:
   1. Hello
   2. Exchange
   3. Flooding
4. Hello is used to check that links are operational and to select what is called a designated router. Hello packets are sent at regular intervals that are customized by the network administrator on each OSPF router.
5. The Exchange protocol is used during startup of a router where it receives a copy of the network database from a neighboring router.
6. The Flooding protocol is used to propagate link state changes through the network.

# IP Datagram Format (OSPF)

| Physical Network Header | | IP Datagram as Data | |
|---|---|---|---|
| | | IP Header | Data |

| 0 ............... 7 | 8 | 15 | 16 | 19 | 31 |
|---|---|---|---|---|---|
| VERS    LEN | Type of Service | | Total Length | | |
| Identification | | | Flags | Fragment Offset | |
| TTL | PROTOCOL = 89 | | Header checksum | (20 bytes) | |
| source IP address | | | | | |
| destination IP address | | | | | |
| Options | | | padding | | |
| data | | | | | |

Ⓒ IBM Corporation 2004 - 2011

---

1. This is the IP datagram format.
   1. The IBM manuals do not show the layout of IP datagrams and TCP segments since the layouts are located in the TCP/IP RFCs and also in TCP/IP tutorials available as IBM Redbooks or as other non-IBM publications.
   2. If you have worked with SNA in the past, this fact about TCP/IP documentation and formats will be new to you. SNA is an architecture owned and controlled by IBM. IBM publishes Formats and Protocols manuals as well as Data Areas Manuals to describe internals of SNA.
2. OSPF packets are sent as IP Datagrams with a Protocol Type of 89.

# OSPF Packet Format

| Physical Network Header | IP Datagram as Data | | |
|---|---|---|---|
| | IP Header | OSPF Header | Data |

```
0..............................7 8................15 16............23 24................31
```

| Version # = 2 | Type | Packet Length | |
|---|---|---|---|
| Router ID | | | |
| Area ID | | | (24 bytes) |
| Checksum | | Autype | |
| Authentication | | | |
| Authentication | | | |
| data bytes = format depends on packet type | | | |

© IBM Corporation 2004 - 2011

1. OSPF Version # is  2 (RFC 1583).
2. OSPF Packet Types
   1. Type 1 = Hello Packet
   2. Type 2 = DataBase Description Packet
   3. Type 3 = Link State Request Packets
   4. Type 4 = Link State Update Packets
   5. Type 5 = LSA Acknowledgement Packets
3. The packet length in  bytes includes the OSPF header.
4. The Router ID is the the ID of the router originating the packet.
5. The Area ID is the area  that the packet is being sent into.
6. The Checksum  is the standard IP checksum of the entire contents of the packet, excluding the 64-bit authentication field.
7. The AuType identifies either 0 for no authentication  or 1 for plain text 64-bit password or 2  for  MD5 authentication.
8. Authentication is used by the authentication scheme.

# Types of OSPF Packets

➤ Type 1 = Hello Packets
➤ Type 2 = DataBase Description Packets
➤ Type 3 = Link State Request Packets
➤ Type 4 = Link State Update Packets
➤ Type 5 = Link State Acknowledgement Packets
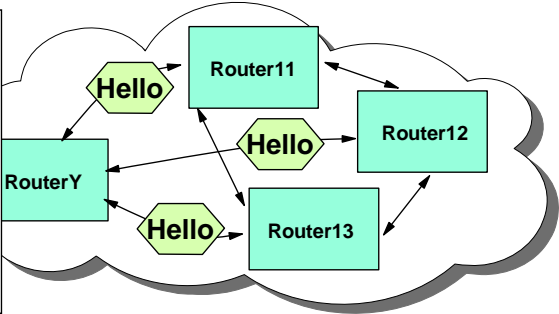
# Discovering the Neighbors

**IBM**

| Router1 | | Router2 |

**Hello (Type 1):  I am here & Here is what I know.**

**Hello (Type 1):  I am here & Here is what I know.**

**We are now "merely neighbors."**

| OSPF Packet Header, Type 1 | | |
|---|---|---|
| Network Mask | | |
| Hello Int. | Options | RTR Priority |
| Dead Router Interval | | |
| Designated Router | | |
| Backup Designated Router | | |
| RouterIds of All Neighbors on This Medium | | |

© IBM Corporation 2004 - 2011

1. The HELLO packet functions as a "keepalive" packet to determine whether a neighbor is still operating; it also sets up the initial negotiation of the connection, negotiating who the Designated Router is on the network.
   1. Connections between neighbors must maintain the same Hello_ Interval  in order to establish a neighbor relationship.  They must  also maintain the same Dead_Router_Interval.
   2. The Dead_Router_Interval is recommended to be at least 4x  the value of the Hello Interval.  If a neighbor has been unavailable for the Dead_Router_Interval, the neighbor relationship is marked as "down."
   3. These intervals have different defaults depending on the platform implementation.  Console error messages at OS/390 indicate whether these values are out of synch.
2. The Options field indicates whether the router supports the TOS field or is capable of sending External Routes.
3. The Priority  field indicates with a non-zero value that the router is eligible to be a Designated Router.
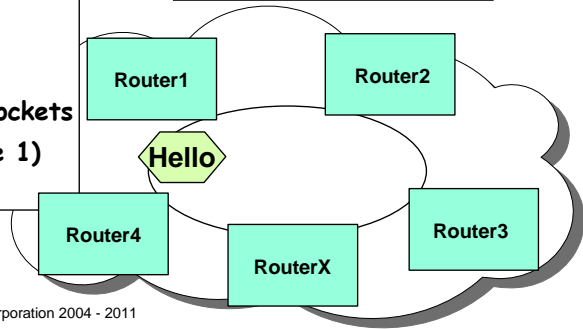
# Discovering Neighbors & DR Eligibility

**Non-Broadcast Network Multi-Access (NBMA)**
- Neighbors must be pre-defined
- X.25, Frame Relay, Native ATM
- Unicast to Interface IP Addresses: Hello (Type 1)
- Requires DR election

HELLO: RTR Priority
OSPF Packet: RouterID ...

**Broadcast Network**
- Neighbors dynamically discovered
- TR, Ethernet, FDDI, LANE, HiperSockets
- Multicast to 224.0.0.5:  Hello (Type 1)
- Requires DR election

© IBM Corporation 2004 - 2011

1. This visual may appear to describe the Hello Packet from the perspective of only one router in the network.  In fact, all the routers in the NBMA and Broadcast  networks send each other Hello  Packets in a fully-meshed fashion.
2. Technically, OSPF requires multicast capability on all interfaces over which it receives OSPF information.  Otherwise, packets are discarded.
3. However, certain types of interfaces are incapable of multicast.  Either they do not have the microcode to support it or the architecture of the medium does not include multicast capability.  For such interfaces, OSPF will work with unicast if the network is defined as "non-broadcast=yes" and all DR neighbors and Non-DR neighbors are defined.  In such instances, all devices on such a network must be defined as unicast even if some of them are multicast-capable.  Since the definition burden for unicast can be onerous, the recommendation is to use multicast wherever possible.
4. A new  type of Broadcast network was introduced in z/OS V1R2: HiperSockets, also known as iQDIO.  It enables an internal connection within a CEC among z/OS images or among z/OS images together with LINUX images.  You must ensure that the HiperSockets network or LAN has at least one Designated Router and preferably a Backup Designated Router.
5. Depending on the type of network in which the neighbors are participating, you may need to define who the eligible Designated Routers are inside your OSPF definitions.  In CS/390, these definitions are maintained in the OMPROUTE Configuration file.
    1. The Priority field of the HELLO packet indicates with a non-zero value that a router is eligible to become a DR.
6. The HELLO protocol determines who the Designated Router (DR) will be.
    1. Role of the DR:
        1. It is adjacent to all other routers on the network.
        2. It generates and floods the network links advertisements on behalf of the network.  In a broadcast network this reduces the amount of router protocol traffic that is generated, as only the DR is responsible  for flooding the network with the information.
        3. It is responsible for maintaining the network topology database that is replicated at all other routers on the network and in the same area.
    2. Role of the Backup DR:
        1. It takes over should the DR fail.
7. The router with the higher Router_Priority becomes the DR on a broadcast or non-broadcast multiaccess network.  If there is a tie, the router with the higher Router_ID becomes the DR.  For point-to-point links do not specify  Router_Priority; it is irrelevant.   If your OS/390 or z/OS system is to be used primarily for production or test work and not for routing, consider setting Router_Priority to 0 for all CS interfaces so that the CS node is ineligible to become the DR.
    1. There is a wait period in the election of the DR - so, if Router 12 and 13 have communicated and Router 11 and RouterY have communicated, nevertheless, only one of them would become the DR.  This wait is the Dead_Router_Interval, during which time the router declares itself eligible for DR (with a non-zero router priority); however, until the expiration of the DR interval, no one router declares itself the DR.  This wait is introduced so that as many eligible routers as possible have a chance to come up and announce themselves before a DR is selected.
8. If more than one router declares itself the Designated Router (DR), the one with the higher RouterID becomes the DR.
9. If a DR has been declared and another router with a higher RouterID joins the network, the original DR maintains its role as Designated Router.
10. There must be a Backup DR; the original election of the DR leaves the router with the next higher RouterID as the Backup DR.
11. If the Router ID is not specified, the IP address of one of the OSPF interfaces will be used as Router ID.  For CS in OS/390 or z/OS, we recommend that you use the IP address of a STATIC VIPA or a physical interface for the RouterID to avoid the default selection of a Dynamic VIPA which could move.
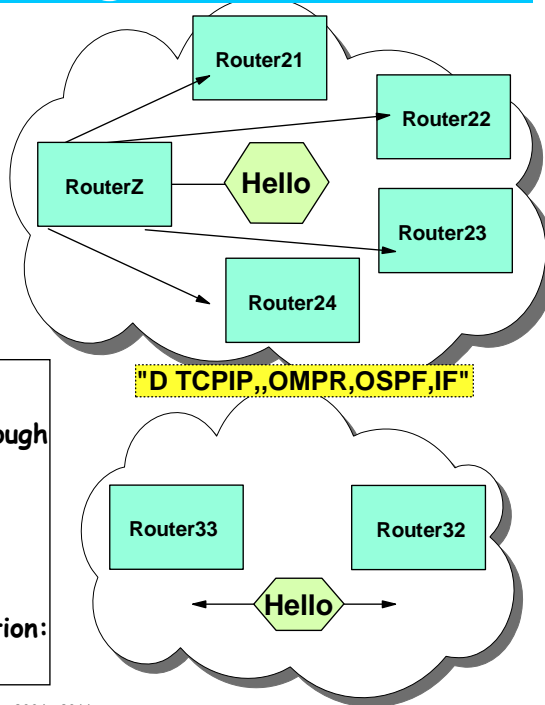
## Discovering Neighbors

**Special Types of Point-to-Multipoint Networks**

- Neighbors dynamically discovered through special signalling protocols
- z/OS MPC, XCF, IUTSAMEH
- Unicast to Each Interface: Hello (Type 1)

**Point-to-Point Networks**

- Neighbors dynamically discovered through HELLO protocol (CTCA, CLAW)
- Link-state Updates multicast to 224.0.0.5
- Link-state Update retransmissions multicast to station at end of connection: Hello (Type 1)

Router21
Router22
RouterZ
Hello
Router23
Router24

**"D TCPIP,,OMPR,OSPF,IF"**

Router33
Router32
Hello

© IBM Corporation 2004 - 2011

1. S/390 and z/Series with IP in CS/390 can dynamically discover neighbors on an XCF or an IUTSAMEH network.
   1. OSPF definition is thus simplified, as we do not need to predefine the routers that are DR-eligible in the network in the OMPROUTE Configuration file.
2. This special type of Pt-to-Multipoint network does not require a DR, as we do not need to predefine the neighors in the network in the OMPROUTE Configuration File.
3. CLAW and CTCA are point-to-point in CS for OS/390
   1. Cisco Routers describe themselves as point-to-multipoint when they are configured with the following interfaces to S/390:
      1. CLAW (NOTE: An APAR released in July of 2001 permits the definition of a CLAW interface as point-to-multipoin with a parameter on the LINK statement. See APAR PQ48766.)
      2. ESCON EMIF MPC+
4. An OMPROUTE command will display the type of network that is available to the connection you have defined:
   1. "D TCPIP,,OMPR,OSPF,IF" displays the interface types.
5. Networks that do not require DR election like Point-to-Point and Point-to-Multipoint networks may be set up as Demand Circuits and may use Hello-Suppression to minimize OSPF traffic: with Demand_Circuit=YES, LSAs are not periodically refreshed over the interface; only LSAs with real changes are advertised and LSAs flooded over this interface never age out.
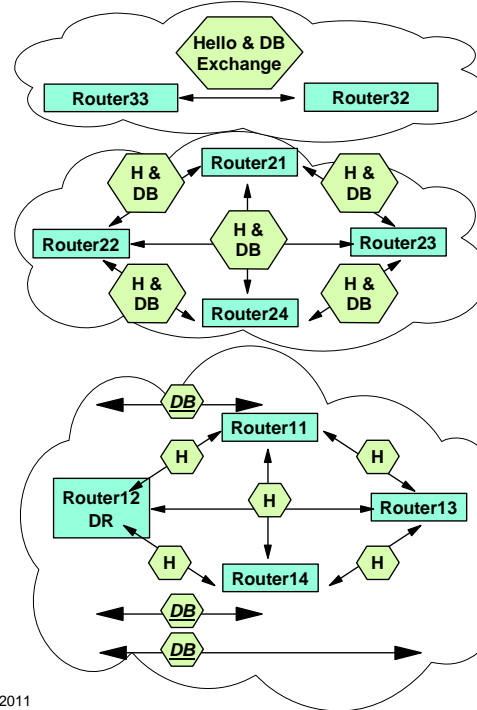
Introducing the Designated Router

**No Designated Router:**

**Point-to-Multipoint Networks &**

**Point-to-Point Networks**

- Hello Protocol with all neighbors where hardware permits
- Data Base Exchange & Adjacency with all neighbors where hardware permits (aka "full adjacency")

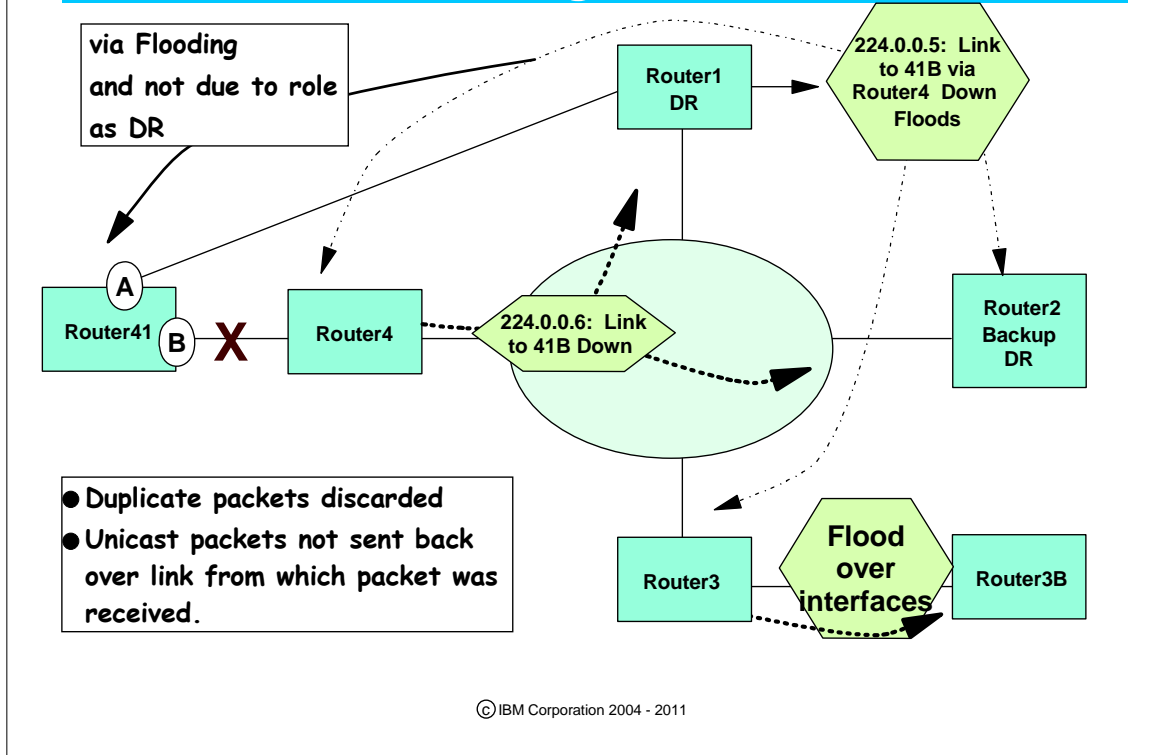**Designated Router (DR):**

**Broadcast and NBMA Networks**

- Hello Protocol with all neighbors
- Data Base Exchange & Adjacency only between DR|Backup DR and each non-DR (aka -also known as- "full adjacency")
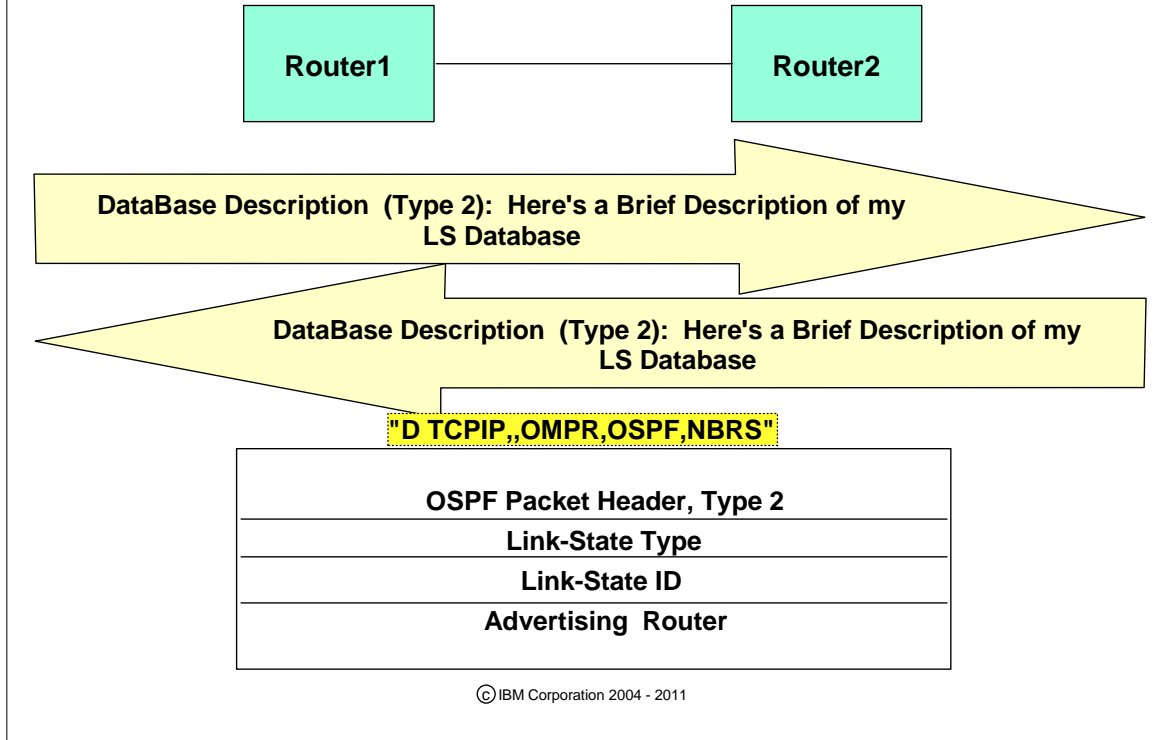
© IBM Corporation 2004 - 2011

1. No matter what the network type, the HELLO exchange occurs among all members of a network. However the Database Synchronization and Update process may not need to involve exchanges of each member with every other member. A fully meshed exchange of database information could contribute heavy traffic loads to a network and should be avoided where possible. The concept of Designated Router (DR) minimizes the exchange of topology information among network members without compromising the OSPF principle of synchronized databases for members of an OSPF area.
2. The top half of the visual shows you a Point-to-Point network and one type of Point-to-Multipoint network to illustrate the flow of HELLO packets and DATABASE synchronization packets. NOTE: The notation of "H and DB" in the hexagon of the diagram stands for "Hello and DB Exchange."
3. In Point-to-Multipoint and Point-to-Point networks routers synchronize their databases with their neighbors. In a Point-to-Point network, each router has only one neighbor on any single connection and the router must become "fully adjacent" with that neighbor. This means that the "fully adjacent" neighbor must synchronize its database with its other "fully adjacent" neighbor for the area in which the connection resides.
4. In a Point-to-Multipoint network, you could see a few different scenarios. The physical topology of one type of Point-to-Multipoint network may allow one router to be the "hub" or converging "point" of the network; this router would thus have multiple "fully adjacent" neighbors for the database exchange. On the other hand each of its neighbors sees only one neighbor with which it must synchronize databases.
5. There are also special types of Point-to-Multipoint networks like XCF which have a different physical topology from the one just described. For example, with XCF each neighbor must synchronize databases with all other members of the XCF networks. This presents a meshed network appearance with many network flows. Router21 has established adjacencies with Routers 22-24. Likewise, Router 22 has established adjacencies with Routers 21, 23, and 24.
6. The bottom half of the visual shows you a broadcast or NBMA network to illustrate the flow of HELLO packets and DATABASE synchronization packets.
7. All routers synchronize their databases with that of the Designated Router (and the Backup Designated Router) in an NBMA or Broadcast network. That is, all non-DR routers are "fully adjacent" only to the DR and the Backup DR. This minimizes routing update traffic, eliminating the meshed exchange of database information. In our network diagram, the Designated Router has only three "full adjacencies" for the purposes of database exchange: with Router11, with Router13, and with Router14. Router11 has only one adjacency --with the DR (Router12)--as do Router13 and Router14. However, note that Router11 would see Router13 as a "neighbor," although not as a "fully adjacent" neighbor.

# Role of Designated Router

via Flooding and not due to role as DR

224.0.0.5: Link to 41B via Router4 Down Floods

Router1 DR

Router41

A

B

X

Router4

224.0.0.6: Link to 41B Down

Router2 Backup DR

- Duplicate packets discarded
- Unicast packets not sent back over link from which packet was received.

Router3

Flood over interfaces

Router3B

© IBM Corporation 2004 - 2011

1. All routers synchronize their databases with that of the Designated Router.
    1. When a link-state advertisement is required due to a change, the router multicasts this LSA to the DR and its Backup DR using the special multicast address of 224.0.0.6 (for "All Designated Routers").
       1. If this is a non-broadcast network, the unicast address and not the multicast address is used
    2. The DR then floods this information to all its interfaces and to the network using the multicast address of 224.0.0.5 ("All OSPF Routers").
       1. If this is a non-broadcast network, the unicast address and not the multicast address is used.  Furthermore, the advertisement is not sent back to the same interface from which it was received.  (That is, Router 4 does not hear back from the DR that the link to 41 is down if it is attached to an NBMA network using unicast.)
2. If duplicate packets are received, as could be the case when there are point-to-point links in the network or alternate connections to the DR,  or if there is router redundancy, the sequence number in the packet let the recipient know  this is a duplicate.  The recipient then discards the duplicate.
3. Routers that hear updates from their DR flood the information over their other interfaces to other attached routers.  (See Router3-to-Router3B communication.)
    1. Each LSA must be acknowledged, either implicitly or explicitly with a Link State  Acknowledgment packet (OSPF packet Type 5).
       1. An implicit acknowledgment can be perceived due to an update that was intended for the flooding router anyway.
    2. Router 1 floods the information to Router 41A, since this is a pt.-to-pt. connection.  Router1 is not acting in the capacity of a DR here.

**Forming Adjacencies: DataBase Exchange**

IBM

Router1 — Router2

DataBase Description (Type 2): Here's a Brief Description of my LS Database

DataBase Description (Type 2): Here's a Brief Description of my LS Database

**"D TCPIP,,OMPR,OSPF,NBRS"**

| OSPF Packet Header, Type 2 |
|---|
| Link-State Type |
| Link-State ID |
| Advertising Router |

ⓒ IBM Corporation 2004 - 2011

1. Once the HELLO packets have been exchanged there is two-way communication and the routers are "merely neighbors."  They then know  with whom they must establish adjacencies.  (Remember:  bringing up an adjacency means that the databases are synchronized.)
   1. If the network is multiaccess, all routers are adjacent to the Designated and Backup Designated Routers.  The adjacencies are only with the DR or Backup DR.
   2. If the network is point-to-point, or special point-to-multipoint like MPC, XCF and IUTSAMEH,  the router forms an adjacency with the partner at the other end of the connection.
2. The routers next use the Exchange Protocol to exchange a description of their link state  database  (via OSPF Type 2 packets) with their adjacent partners. If they determine that they are fully synchronized, they are considered "fully adjacent."
3. To display the status of adjacencies, use the OMPROUTE command...
   1. "D TCPIP,,OMPR,OSPF,NBRS"

## Synchronizing the Databases

**Router1** ——— **Router2**

DataBase Description  (Type 2): Brief Description LS Database →

← DataBase Description  (Type 2): Brief Description LS Database

Link State Request  (Type 3):  I don't have these records - Send! →

← Link State Request  (Type 3):  I don't have these records - Send!

Link State Update  (Type 4):  Here are the requested records →

← Link State Update  (Type 4):  Here are the requested records

◄— Explicit Link State Acknowledgements (OSPF Type 5) or Implicit in LS Update Packet —►

**We are now "fully adjacent."**

© IBM Corporation 2004 - 2011

1. If the router databases are not fully synchronized, they request more information from the adjacent router (Link State Request/OSPF packet Type 3); then they exchange Link State Updates (OSPF Type 4 packets) to get in synch so that they may become "fully adjacent."
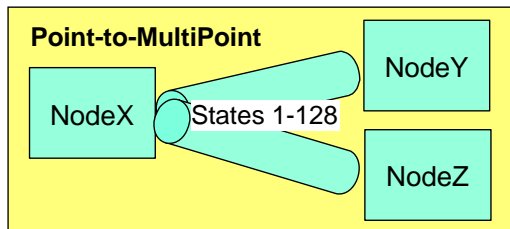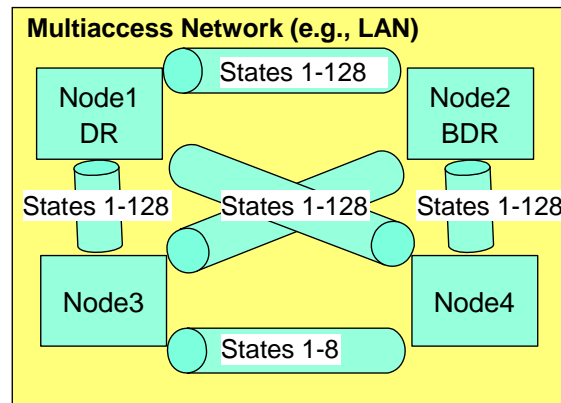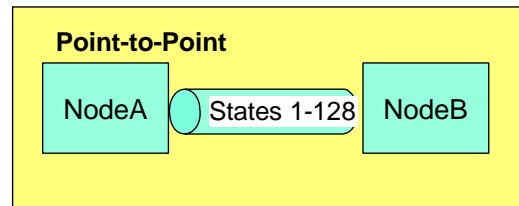   1. The Exchange Protocol is more complicated than what is depicted here.  One router (the one with the higher RouterID) assumes the master role and the other the slave role.  The master sends its database descriptions, one at a time.  The slave acknowledges each one and includes in the acknowledgement its database descriptions.  If the new description indicates that this record is newer than what the recipient already has in its database, this description is saved as  a record "of interest."
      1. The records are compared according to "type," "advertising router," and "link state ID."  A sequence number in the record reveals whether the record is newer or older.
   2. Once all descriptions have been received, the neighbors send out database requests  for more complete information about the records "of interest," which are then followed with the actual update packets.  The Link State Update packets (OSPF Packet type 4) implement the "flooding" procedure.
   3. Each update packet must be acknowledged, either explicitly with a DataBase Acknowledgement packet (OSPF Type 5) or implicitly in the LS Update packets (OSPF Type 4).
2. It may not be until this point that the routers are "fully adjacent."  This means that the link state databases are fully synchronized.  Recall that whether or not a router interface is "fully adjacent" with every neighbor it knows about depends upon the network type to which the interface is attached.  If the interface is part of a network that requires the election of a Designated Router and the router itself is neither the DR nor the Backup DR, it might have many neighbors (neighbor state 8), but it will be fully adjacent only with the Designated Router and the Backup Designated router (neighbor state 128).  If the Database synchronization process is quite lengthy, you might consider setting the DB_Exchange_Interval higher than the default length (equivalent to the Dead_Router_Interval).  If the database synchronization fails due to insufficient time, the neighbor process will never stablize beyond Neighbor State of 8 (and may even continually cycle through states 8, 16, and 32 as indicated in Messages EZZ7919I and EZZ7921I, event 12), and the DB exchange and loading process will fail with event code 15:  "failure to thrive" in Message EZZ7921I OSPF ADJACENCY FAILURE.  DB_Exchange_Interval was introduced with APARs PQ45413, PQ37048, and PQ39760 for V2R7, V2R8 and V2R10 respectively.
3. Again, to display the status of adjacencies, use the OMPROUTE command...
4. "D TCPIP,,OMPR,OSPF,NBRS"

# Summary:  Neighbor States

| STATE | | Router1 | Router2 | | STATE | |
|---|---|---|---|---|---|---|

**"D TCPIP,,OMPR,OSPF,NBRS"**

| STATE (Router1) | | Router2 STATE |
|---|---|---|
| DOWN 1 | | DOWN 1 |
|  | Hello (Type 1):  I am here & Here is what I know. → | INIT 4 |
| INIT 4 | ← Hello (Type 1):  I am here & Here is what I know. | |
| 2-WAY 8 | We are now "merely neighbors." | 2-WAY 8 |
| EXSTART 16 | ← DB Desc. (Type 2):  EMPTY:  No DB Description → | EXSTART 16 |
| EXCHANGE 32 | DB Desc. (Type 2):  Brief Description of my DB → | |
|  | ← DB Desc. (Type 2):  Brief Description of my DB | EXCHANGE 32 |
| LOADING 64 | LS Req. (Type 3):  I don't have these - send them → | LOADING 64 |
|  | ← LS Req. (Type 3):  I don't have these - send them | |
|  | LS Upd. (Type 4):  Here are records you want → | |
| FULL 128 | ← LS Upd. (Type 4):  Here are records you want | FULL 128 |
|  | We are now "fully adjacent." | |

**"D TCPIP,TCPIP,OMPR,OSPF,DATABASE,AREAID=0.0.0.0  "**

1. This page represents a summary of what we have seen in the protocol exchanges on the previous pages.  It depicts the state of the conversation that is transpiring between neighboring routers in the OSPF network.
2. It also indicates the state changes that the routers (or CS nodes) go through during each phase of the protocol exchange.
3. Depending on the platform, the command displays for monitoring the routing in the network might show either the verbal status ("DOWN," INIT," etc.) or a number that  can be equated with that state (DOWN=1, INIT=4, and so on).
4. 1 (DOWN):  Initially the routers are down and have no contact with each other.
5. 2 (Attempt):  In NonBroadcast MultiAccess (NBMA) networks a different state may be indicated even when a router is marked down.  "Attempt" indicates that no contact has been made but the Hello packet will continue to send the packets to "attempt" to make contact.
6. 4 (Initialize):  The Hello packet has been identified/received but no exchange of information has taken place between neighboring routers.  Contact between the neighboring routers has been made, however.
7. 8 (2-Way):  The Hello packets have been received and acknowledged by both neighboring routers, as indicated by the presence of the router itself in the neighbor's Hello packet.  That is, each router sees itself in the neighbor's Hello packet.  The designated router (DR) is selected and thereafter follows the selection of the Backup DR.
8. 16 (ExStart):  Neighbor routers form adjacencies between themselves.  Neighbor routers' communication is more advanced in this state and the routers decide who is the "master" and who is the "slave" and what is the initial DataBase Sequence Number. The transmission of the link state database can begin.
9. 32 (Exchange):  The neighbors send their link state database to their adjacent routers.  The link state database records describes the characteristics of the database and each must be acknowledged.  Link state request packets requesting the neighbor's recent LSA's status may be sent to the neighboring routers.  In this state the neighbors are capable of receiving and sending all types of OSPF routing protocol packets.
10. 64 (Loading):  The link state request packets are being transmitted and received from neighboring routers requesting the most recent LSAs.
11. 128 (Full):  If the neighboring routers are displayed with this state, it  means that they are fully adjacent. l The adjacent routers exchange LSAs and appear in their neighbors' router-LSAs.
12. The synchronization of databases can be verified with a command to display the OSPF database; there is a field called CHECKSUM TOTAL, which can be used to compare if 2 routers have synched their DBs.
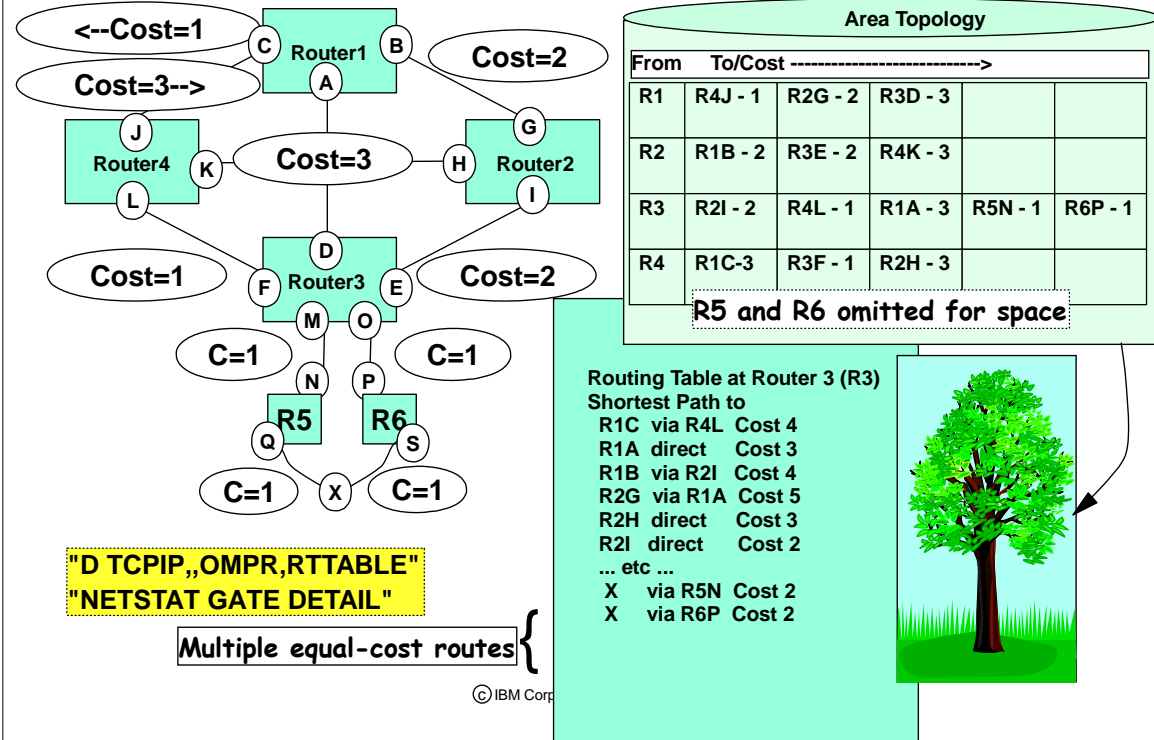    1. D TCPIP,TCPIP,OMPR,OSPF,DATABASE,AREAID=0.0.0.0

# Neighbor States

| Nbr. State # | Meaning |
|---:|---|
| 1 | Down |
| 2 | Attempt |
| 4 | Init |
| 8 | 2-way |
| 16 | ExStart |
| 32 | Exchange |
| 64 | Loading |
| 128 | Full |

**Point-to-Point**

NodeA — States 1-128 — NodeB

**Multiaccess Network (e.g., LAN)**

Node1 DR — States 1-128 — Node2 BDR

States 1-128    States 1-128    States 1-128

Node3    States 1-8    Node4

**Point-to-MultiPoint**

NodeX — States 1-128 — NodeY / NodeZ

© IBM Corporation 2004 - 2011

1. The Neighbor State Codes are described in RFC 1583.
2. A multiaccess network is a network in which all the attached hosts share one transport medium. Examples of multiaccess networks are LANs, ATM, etc.
3. Although you will frequently hear that the best neighbor state to be in is State 128 (Full Adjacency), in fact, on a Multiaccess network the only Full Adjacency is with the DR and the BDR. In the broadcast network depicted, Node3 and Node4 are perfectly fine with Neighbor State of 8 (2-way communication), as they need not synchronize their databases with each other. They synchronize only with the DR and the BDR. They simply need to know that they are available for communication with each other in the subnetwork.

# Computing the Routing Table

**IBM**

<--Cost=1

**C** Router1 **B**  Cost=2

Cost=3-->  **A**

**J**
Router4 **K**  Cost=3  **H** Router2  **G**
**L**  **I**

Cost=1  **F** Router3 **E**  Cost=2
**D**
**M** **O**

C=1  **N** **P**  C=1

**R5** **R6**
**Q** **S**

C=1 **X** C=1

**"D TCPIP,,OMPR,RTTABLE"**
**"NETSTAT GATE DETAIL"**

Multiple equal-cost routes {

© IBM Corp

## Area Topology

| From | To/Cost ----------------------------> | | | | |
|------|------------|------------|------------|------------|------------|
| R1 | R4J - 1 | R2G - 2 | R3D - 3 | | |
| R2 | R1B - 2 | R3E - 2 | R4K - 3 | | |
| R3 | R2I - 2 | R4L - 1 | R1A - 3 | R5N - 1 | R6P - 1 |
| R4 | R1C-3 | R3F - 1 | R2H - 3 | | |

R5 and R6 omitted for space

Routing Table at Router 3 (R3)
Shortest Path to
R1C  via R4L  Cost 4
R1A  direct     Cost 3
R1B  via R2I  Cost 4
R2G  via R1A  Cost 5
R2H  direct     Cost 3
R2I   direct     Cost 2
... etc ...
X    via R5N  Cost 2
X    via R6P  Cost 2

1. All routers in the area have identical copies of the AS topology database.
2. Each router computes its own routing table using a "spanning tree" algorithm.
   1. This is also called the "Dijkstra" algorithm for computing "Shortest Path FIrst."
   2. Every time a link-state update is received the entire tree is recomputed.
      1. Christian Huitema quotes a time of 200 milliseconds to compute the routes for a network with 200 nodes.  (p. 107 of "Routing in the Internet," Prentice Hall.
3. If there are multiple equal-cost routes, the routing tree retains them.
   1. If there are multiple routes to the same destination, the routing tree selects only the shortest path for retention.
4. D TCPIP,,OMPR,RTTABLE shows the cost of the routes in the network.  TSO NETSTAT GATE DETAIL also shows the costs of the links to each gateway attached to the router.
5. The RFC provides for the OSPF protocol to compute a different routing tree for each type of service (represented by the TOS bits - bits 3, 4, and 5 -  of the TOS Byte in the IP Header).  However, no known implementation of OSPF provides for multiple trees.

# Advertising Link States: Flooding Review

A = Flooding Intra-area OSPF Routing Info
B = Flooding External Routing Info
C = Flooding Inter-area Summary OSPF Routing Info
D = Flooding Default Routing Info

**Backbone Area**

A, B, C

Router1

A, B, C

B, C → Rz

Router4

Router2 ABR

**Non-Backbone Area**

A, B, C

A, B, C

Router3 ASBR & ABR

**AS with RIP Routing**

(C), D

A

**(Totally) Stub(by) Area**

Ry

ⓒ IBM Corporation 2004 - 2011

1. You have seen this visual before. It is just a reminder that there are several types of LSAs that can be sent into the network. These LSA Types are summarized on the next pages.
2. LSAs can be updated only every 5 seconds. If an LSA with a different sequence number arrives within 5 seconds of a previous LSA for the same route, it is ignored.
3. If the age of the LSA in the DB is less than 1 second, another LSA for the same route canot be accepted.
4. A router is required to acknowledge LSAs, even if they are duplicates. However, it forwards only one copy of the LSA if it has received duplicates.
5. Every 30 minutes the Designated Routers must exchange their link state databases.

# LSA Header Format

| Physical Network Header | IP Datagram as Data | | |
|---|---|---|---|
| | IP Header | OSPF Header | Data |

```
0................7  8..............15  16..............23  24..............31
```

**(24 bytes)**

| LS Age | Options | LS Type |
|---|---|---|
| Link State ID | | |
| Advertising Router | | |
| LS Sequence Number | | |
| LS Checksum | LS Length | |
| data bytes = format depends on LS Type | | |

20 bytes

© IBM Corporation 2004 - 2011

1. OSPF in OMPROUTE supports five types of LSAs.
2. An LSA Type 7 is described in the RFCs. It is used to advertise a Not So Stubby Area (NSSA). OMPROUTE does not support this type of LSA.

# Types of Link State Advertisements

➤ Type 1 = Router Links
  ➤ Describes all router's interfaces into an area and their states (Intra-area destinations)
➤ Type 2 = Network Links
  ➤ Describes DR connection to medium in broadcast or NBMA nets
➤ Type 3 = Summary Links, created by Area Border Routers
  ➤ Describe inter-area routes
➤ Type 4 = Summary Links, created by Area Border Routers
  ➤ Describe routes to the AS Boundary Routers
➤ Type 5 = AS External Links, created by AS Boundary Routers
  ➤ Describe routes to destinations outside the AS

```
"D TCPIP,,OMPR,OSPF,DATABASE,AREA=..."
"D TCPIP,,OMPR,OSPF,DBSIZE"
"D TCPIP,,OMPR,OSPF,LSA,LSTYPE=..."
```

© IBM Corporation 2004 - 2011

1. More about Type 1:  Router LSA
   1. Indicates whether originating router has Virtual Links, is an ASBR or an ABR
   2. Indicates for each link whether it is pt-to-pt, connection to a transit network, connection to a Stub area, or Virtual Link.
2. More about Type 2:  Network LSA
   1. Originated by DR in broadcast and NBMA networks
   2. Contains list of all routers that are attached (adjacent) to the DR
3. More about Type 3:  Summary LSA
   1. Originated by ABR
   2. Contains IP Destination Address inside the area
   3. One advertisement per destination
   4. Network Mask is mask of IP Address Destination
   5. Used to describe Default Route to Stub Area
      1. LS Id Field and Network Mask both set to 0.0.0.0
4. More about Type 4:  Summary LSA about ASBR
   1. Originated by A=ABR
   2. Contains IP Address of ASBR
   3. Network Mask is 0.
5. More about Type 5:  External LSAs
   1. Originated by ASBRs
   2. Describe Destinations outside the AS
      1. Can be default destination
         1. LS ID Field is 0.0.0.0
      2. Can be network destination
         1. LS ID Field contains network address
      3. Identifies whether to use Type 1 metric for destination
         1. Metric is comparable to a link-state metric
      4. Or... Identifies whether to use Type 2 metric
         1. Metric is greater than link-state metric
6. Notes on Type 7 LSAs:  Not So Stubby Area
   1. This optional OSPF feature is described in IETF RFC 1587 and in Internet Draft draft-ietf-ospf-nssa-update-10.tx.  NSSAs are similar to the existing OSPF Stub Area.  However, unlike an OSPF Stub Area, they may import AS external routes in a limited fashion.
7. The following OMPROUTE commands display information about the LSAs that have reached this router:
   1. D TCPIP,,OMPR.,OSPF,DATABASE,AREA=...
   2. D TCPIP,,OMPR,OSPF,DBSIZE
   3. D TCPIP,,OMPR,OSPF,LSA,LSTYPE=...

# Aging the Link State Records

| RouterA | | RouterB | | RouterC | |
|---|---|---|---|---|---|
| New Record: Age = 0 | D=2 | Record: Age = 2 | Record: Age = 3 | D=4 | Record: Age = 7 | Record: Age = 8 |

**"D TCPIP,,OMPR,OSPF,RTTABLE"**

➤ A new record has an **Age** of "0."

➤ The record age is incremented by the Transmission Delay (D value in diagram) that is configured for the interface it is sent over.

➤ It is also incremented for every second it is maintained.

➤ Maximum **Age** of a record is one hour, at which point routers will no longer use the record in their route table calculation. The record will eventually be removed from the database.

➤ A router must retransmit a record *at least* once every 30 minutes.

1. D TCPIP,,OMPR,RTTABLE displays the age of information about connections.

# Appendix: Basic Coding Examples of OSPF with OMPRoute in z/OS

Please consult the OSPF presentations available at www.ibm.com/support/techdocs for more examples of:
- coding and designing for OSPF on z/OS and for
- avoiding common errors and reviewing hints and tips.

Or consult the Extended Coding Examples in the Appendix of this document.

# Basic OSPF Definitions in OMProute

IBM

```
OSPF
   RouterID=10.138.165.9
```

**BEST PRACTICE:** *Specify RouterID as static VIPA* if multiple interfaces in CS.

```
   Comparison=Type2
```

**Specify Comparison to designate how to compare external routes when computing the cost of the entire route.**

```
   Demand_Circuit=YES
;
```

**BEST PRACTICE:** *Do not code* **Demand_Circuit=YES or Hello_Suppression for OSPF protocol on the z platform. Typically used for lines that are costly to run, like x.25 lines or dial-up modem lines. Included as parameters in z OMPROUTE, but only because they are part of the RFCs. Stops periodic refreshing of the LSAs over an interface.**

```
Area
   Area_Number=1.1.1.1

   STUB_AREA=NO | YES
;
```

**Specify Area Number and type of area if Stub or Totally Stubby Area.**

```
AS_Boundary_Routing
   Import_RIP_Routes = YES

   Import_Static_Routes = YES
;
```

**If CS in z/OS OSPF node is an ASBR between OSPF and a RIP AS, and between OSPF and a STATIC AS, indicate it.**

© IBM Corporation 2004 - 2011

1. Guideline for OSPF statement in omproute configuration file: The standalone statements of "RouterID," "Comparison," and "Demand_Circuit" continue to be supported. However, the OSPF statement is the preferred method for defining them, and future potential standalone parameters are to be added to this statement only. If both the OSPF statement and the standalone statements are coded, the last one coded in the configuration file takes precedence.
2. A softcopy for a valid OMPRoute Configuration file exists in SEZAINST(EZAORCFG) and in DOC APAR II11555. However, this softcopy does not contain some of the parameter enhancements that you see on the next couple of pages. Nevertheless, it is a good starting point for coding both IPv4 and IPv6 interfaces.
3. Assuming you have more than one interface in a CS node, always code the Router_ID. We recommend coding it with the value of a static VIPA so that it does not by chance default to a dynamic VIPA which could move to another node and generate confusion in the OSPF network.
4. Comparison to designate how to compare external routes when computing the cost of the entire route. Compare to type 1 or 2 externals. Valid values are Type1 (or 1) or Type2 (or 2). Comparison may be specified on either the OSPF Statement (preferred) or as a Standalone Statement (original coding technique).
   1. The comparison tells OMPROUTE where external routes fit in the IPv4 OSPF hierarchy. OSPF supports two types of external metrics. Type 1 external metrics are equivalent to the link state metric. Type 2 external metrics are greater than the cost of any path internal to the autonomous system. Use of type 2 external metrics assumes that routing between autonomous systems is the major cost of routing a packet, and eliminates the need for conversion of external costs to internal link state metrics.
5. The periodic nature of OSPF routing traffic requires a link's underlying data-link connection to be constantly open. In theory this can result in unwanted usage charges on network segments whose costs are very high, typically on x.25 links or dial-up links with modems. There are two configuration steps that can be taken to inhibit the periodic nature of the protocol, although they are not particularly useful steps to take with the OSPF protocol in z/OS Communications Server: Configure the link with Demand_Circuit or configure it with both Demand_Circuit and Hello_Suppression. Hello suppression is only meaningful if the link is a demand circuit and is either point-to-point or point-to-multipoint. Hello suppression will inhibit the periodic transmission of OSPF hello packets.
   1. IBM does not recommend specifying these parameters for z/OS CS, even though they are part of the OSPF architecture. They end up not being particularly useful in the long run because we don't have dial-up lines from the mainframe and we don't usually have X.25 lines being controlled out of the mainframe anymore. Here, however, are the specifics: Specify Demand_Circuit here and on individual OSPF_Interface statements if there is no need for constant/periodic refreshing of link state advertisements (LSAs) sent over the interface. Only LSAs with real changes will be readvertised. In addition, aging of these LSAs will be disabled such that they will not age out of the link state database.Demand_Circuit may be specified on either the OSPF Statement (preferred) or as a Standalone Statement (original coding technique).
6. If you are an ASBR, you may choose to Import routes from other protocols. Even if you are importing the Default Route from another protocol, you must specify AS_Boundary_Routing parameters.

# Basic OSPF Definitions in OMProute

**IBM**

```
;
;Sample OSPF_INTERFACE (Broadcast: token-ring, ethernet, fddi):

    OSPF_Interface
        IP_Address=9.59.101.5
        Name=TR1
        Subnet_Mask=255.255.255.0
        Attaches_To_Area=0.0.0.0
        MTU=1500
        Cost0=100

        Hello_Interval = 10
        Dead_Router_Interval = 40
        Router_Priority=0;
;
```

**For full-octet IP address, Link Name from IP Profile is required.**
**Specify Subnet Mask to avoid default of standard Class Mask.**

**Code Area that interface attaches to; Area defaults to 0.0.0.0. Specify MTU to avoid default; Specify reasonable cost.**

**Specify Hello_Interval and Dead_Router_Interval; they must be the same value for all nodes attaching to the medium.**

**BEST PRACTICE:  Code Router_Priority=0 for the z/OS platform unless you are dealing with HiperSockets; allow Router to be the DR.  A non-zero value for Router_Priority means that the router is  eligible to become a Designated Router.**

ⓒ IBM Corporation 2004 - 2011

1. If you specify the full four-octet IP interface address, you MUST code the valid Link Name used in the IP Profile.  The only exception to this rule is for Dynamic VIPAs --  you may code the full four-octet IP subnet address and then use a fictitious Link Name.
2. Always code the Subnet Mask of the subnet to which the interface attaches; otherwise OMPRoute will use the standard Class Mask.
3. If the interface does not attach to the backbone area, you MUST code the area number the interface attaches to; if it attaches to the backbone, the area defaults to 0 and coding it is unnecessary.
4. Always code the valid MTU size; otherwise you will get the default of 576 bytes.
5. Always code a reasonable cost for the type of interface.  Despite many examples in the literature to the contrary, set your highest-speed interface to a value like 50 or 100, so that when new technology is introduced, you can set the Cost of that technology even lower.
    1. The cost of an OSPF interface can be dynamically changed using the MODIFY <procname>,OSPF,WEIGHT,NAME=<if_name>,COST=<cost> command. This new cost is flooded quickly throughout the OSPF routing domain, and modifies the routing immediately.  The cost of the interface reverts to its configured value whenever OMPROUTE is restarted. To make the cost change permanent, you must reconfigure the appropriate OSPF_INTERFACE statement in the configuration file.
6. HELLO INTERVAL and DEAD ROUTERINTERVAL must be the same at both ends of  a connection!  RFC recommends that the DEAD_ROUTERINTERVAL be at least four times the value of the HELLO_INTERVAL.
7. Always try to let the router platforms that are meant to do the bulk of the routing become the Designated Routers for interfaces and media that require DR election.  It would be appropriate to make z/OS DR-ineligible on a broadcast medium like a LAN, since a standalone router could easily assume this duty.  For HiperSockets broadcast medium, it would be necessary to make the zSeries OSPF images the DR and Backup DR.  In this example, if we can assume there is a standard router on this LAN, it would be better to code z/OS CS with a Router_Priority of 0 in order to disable the ability to become DR.

## CISCO Example with OSA-E QDIO in Totally Stubby Area

**Area 1.1.1.1**   **Area 0**

CS for z/OS     192.168.130.2   CISCO 7500

192.168.130.39              192.168.51.2

**OMPROUTE.CONF for LAN**

```
;
AREA
    Area_Number=1.1.1.1
Stub_area=Yes
[Import_Summaries=No]        Optional
;

;
OSPF_INTERFACE
    IP_address=192.168.130.39
    Name=GIG0B
    Subnet_mask=255.255.255.0

    Attaches_To_Area=1.1.1.1
    MTU=1500
    Retransmission_Interval=5
    Transmission_Delay=1
    Router_Priority=0
    Hello_Interval=10
    Dead_Router_Interval=40  }
    Cost0=100
;
```

**CISCO Router**

```
!
interface Vlan130
 ip address 192.168.130.2 255.255.255.0
 ip ospf hello-interval 10  }
 ip ospf dead-interval 40   }
 ip ospf priority 10
!
interface Vlan510
 ip address 192.168.51.2 255.255.255.0
 ip ospf hello-interval 10
 ip ospf dead-interval 40
!
router ospf 100
 area 1.1.1.1 stub no-summary
 network 192.168.51.0 0.0.0.255 area 0
 network 192.168.130.0 0.0.0.255 area 1.1.1.1
 .....
```

Significant for Totally Stubby

© IBM Corporation 2004 - 2011

1. Many  examples of Cisco interfacing with CS in an OSPF network may be found in a Whitepaper on OSPF with IBM & CISCO:
   1. "OSPF Design and Interoperability Recommendations for Catalyst 6500 and OSA-Express Environments" http://www-1.ibm.com/servers/eserver/zseries/networking/pdf/ospf_design.pdf
2. Other examples may be found in the IBM Redbook:
   1. "Networking with z/OS and Cisco Routers:  An Interoperability Guide (SG24-6297)
3. This diagram shows the CISCO router behaving as an Area Border Router attached to backbone area 0 and Totally Stubby Area 1.1.1.1.  CS for z/OS resides in the Totally Stubby Area.
   1. Two areas are coded in the CISCO IOS.  Only one area number is coded inside of CS for z/OS.
   2. The coding of "stub no-summary" at the CISCO router makes area 1 a Totally Stubby area.  ("stub" alone would make area 1 merely a stub area, but NOT a totally stubby area.  Remember:  Interarea OSPF summary advertisements are not sent into totally stubby areas.)
   3. Note that CS for z/OS also indicates that Area 1 represents a Stub Area.  Since the ABR knows it is not to send summary advertisements into Area 1, it is not necessary to code anything additional at CS.  However, you see that we have indeed added the parameter "Import_Summaries=No" on the CS side.  In fact, this was not necessary; the parameter is optional because the ABR (i.e., the CISCO router) has already indicated that "no-summary" should be sent to CS.  It does not hurt anything to code "Import_Summaries=No" at a CS for z/OS Totally Stubby OSPF node -- you may regard it as documentation if you choose to code it.
   4. Note the coding for Area number in each platform.
   5. Note how the full four-octet IP address is specified in OMPRoute; therefore the name assigned to the interface is valid:  GIG0B.
   6. The Subnet mask is that of the subnet to which this interface attaches.
   7. The MTU value is coded so as to avoid the default of 576.
   8. Note how the values of the Hello and Dead_Router intervals at both ends of a network agree.
   9. Note how the Cisco ABR is the Designated Router on this Broadcast network.  Not only is the OSPF priority at CISCO set to 10, but we have actually made z/OS ineligible for DR responsibility by setting the priority to 0.

**IBM**

Appendix:  z/OS Communications
Server Enhancements to OSPF Coding:

System Variables,

INCLUDE Statements

# Combining Parameters under OSPF Statement

**IBM**

OSPF
    RouterID=<value>
    Comparison=<value>
    Demand_Circuit=<value>

The parameters indicated used to be stand-alone parameters, but should now be combined under the OSPF statement. The parameters, will, however, continue to be valid as stand-alone values.

1. OSPF:   Use the OSPF statement to specify various parameters that apply globally to IPv4 OSPF, either to all interfaces or to the overall OSPF autonomous system. This statement is intended to replace the following standalone statements:
   1. ROUTERID
   2. COMPARISON
   3. DEMAND_CIRCUIT
      1. Guideline: Those standalone statements continue to be supported. However, the OSPF statement is the preferred method for defining them, and future potential standalone parameters are to be added to this statement only. If both the OSPF statement and the standalone statements are coded, the last one coded in the configuration file takes precedence.

# OMPRoute Initialization: Configuration DD Statement (V1R9)

```
//OMPROUTE PROC
//OMPROUTE EXEC PGM=OMPROUTE,REGION=4096K,TIME=NOLIMIT,
// PARM=('POSIX(ON)',
//     'ENVAR("_CEE_ENVFILE=DD:STDENV")/-t2 -d1')
//OMPCFG DD DSN=USER1.OMPROUTE(&OMPCFG),DISP=SHR

//STDENV   DD DSN=USER1.OMPROUTE(OMPENV1),DISP=SHR

//SYSPRINT DD SYSOUT=*
//CEEDUMP  DD SYSOUT=*,DCB=(RECFM=FB,LRECL=132,BLKSIZE=132)
//*SYSMDUMP DD
DSN=(USER1.OMPROUTE.DUMP),DISP=(NEW,DELETE,CATLG),
//*         DCB=(RECFM=FBS,LRECL=4096,BLKSIZE=4096),
//*         UNIT=SYSDA,SPACE=(CYL,(100,100),RLSE),VOL=SER=IPCS08
```

Be careful about leaving tracing and debugging on.
STDENV files must be VB
Specify OMPCFG DD card to identify the Configuration

© IBM Corporation 2004 - 2011

1. Prior to z/OS V1R9 it was necessary for individual started procedures to be maintained for each OMPROUTE instance.
2. Both internal and external users have requested a way to specify an OMPROUTE configuration file name which includes an MVS system symbol in the started procedure for OMPROUTE, so that one started procedure could be shared by multiple OMPROUTE instances.
3. OMPROUTE now supports a DD:OMPCFG statement in its started procedure. This allows for MVS system symbols to be used in the name of the OMPROUTE configuration file, eliminating the necessity to maintain multiple OMPROUTE started procedures
4. As was the case prior to V1R9, beware of leaving Tracing and Debugging options active in a running OMPROUTE address space.
5. As was the case prior to V1R9, beware of specifying an OMPROUTE Environment file that is not a VB or HFS (zFS) file.
6. search order remains the same:
   1. a. DD:OMPCFG
   2. b. OMPROUTE_FILE environment variable
   3. c. /etc/omproute.conf
   4. d. hlq.ETC.OMPROUTE.CONF

## System Symbols in OMPROUTE Configuration File (V1R9)

```
Routerid=1.1.1.&VIPA1
;
OSPF_Interface
IP_ADDRESS=10.10.10.&VIPA1
SUBNET_MASK=255.255.255.0
;
; Where &VIPA1=1 in the IEASYMxx PARMLIB
; member, the above translates to:
; Routerid=1.1.1.1
;
OSPF_Interface
IP_ADDRESS=10.10.10.1
SUBNET_MASK=255.255.255.0
```

hlq.PARMLIB(IEASYMxx)

&VIPA1=1

- **OMPROUTE supports MVS system symbols in its configuration files.**
- **To confirm correct parsing, use OMPROUTE output commands or enable -t2 -d1 in OMPRoute initialization.**

© IBM Corporation 2004 - 2011

1. The ability to use the MVS system symbols in the OMPROUTE configuration file is nice in and of itself because now OMPROUTE configuration files can be shared
2. between OMPROUTE instances. It was possible to share configuration files between OMPROUTE configuration files between OMPROUTE instances prior to V1R9 by using wildcarding; however in an OSPF environment there was no way to wildcard the Routerid, so if you did share configuration files, there was no way to specify a unique routerid for each OMPROUTE instance.
3. If you need to see how a symbol was translated, turn on –t2 –d1 OMPROUTE trace and look for the text "Translated to". For each line that contained an MVS system symbol there will be a line in the trace file which shows to what the symbol was translated.

## Deleted Routes: D TCPIP,,OMP,RTTABLE,DELETED (V1R9)

- **Display TCPIP,,OMP,RTTABLE,DELETED to see all deleted routes in OMPROUTE's main routing table**

```
D TCPIP,TCPCS1,OMP,RTTABLE,DELETED
EZZ8137I IPV4 DELETED ROUTES 816
TYPE    DEST NET        MASK      COST   AGE     NEXT HOP

 DEL    10.11.0.0       FFFF0000  16     6       NONE
 DEL    10.11.2.1       FFFFFFFF  16     5       NONE
 DEL    10.61.0.2       FFFFFFFF  16     6       NONE
 …
 …
15 NETS DELETED, 2 NETS INACTIVE
```

1. This is an example of the output of the D TCPIP,OMP,RTTABLE,DELETED display.
2. The same information can be seen also in the F OMP,RTTABLE,DELETED display
3. In order to investigate deleted networks as indicated in OSPF RTTABLE display, it used to be necessary to run an OMPROUTE debug trace and analyze the EZZ8061I and EZZ7943I messages indicating each network as it is deleted, or to take a dump of the OMPROUTE address space and send it to support.

# INCLUDE Statement for OMPROUTE Configuration File (V1R10)

```
AREA
AREA_NUMBER=1.1.1.1
STUB_AREA=NO;

INCLUDE /u/user1/omproute.conf
INCLUDE //'USER1.INC10'
Include //'USER1.&SYSNAME..OMP'

OSPF_INTERFACE
IP_ADDRESS=10.9.128.128
NAME=DUMMY_SASRVA2
SUBNET_MASK=255.255.255.240
ROUTER_PRIORITY=0
ATTACHES_TO_AREA=1.1.1.1;
```

- **OMPROUTE    "Include file":**
  - Easier to share common OMPROUTE definitions within a Sysplex

```
>>__ _____><
     |_Include__ _//'fully qualified MVS dataset name'_ _|
                 |_/file system absolute pathname_____|
```

© IBM Corporation 2004 - 2011

1. Common configuration statements can be grouped into separate files and specified in the OMPROUTE configuration via the INCLUDE statement
2. Single, multiple, and nested INCLUDE statements can be used in configuring OMPROUTE
3. Rules:
   1. 1. INCLUDE statement must be the only configuration statement on the line.
   2. 2. INCLUDE statement must not end with semicolon.
   3. 3. There must be no more than 10 nested INCLUDE statements.
   4. 4. Static system symbols can be specified as part of the data set name.
   5. 5. Only 1 INCLUDE statement can be specified per line, anything else that follows the statement will be ignored.
4. If a syntax error is encountered in the final version of the configuration file after INCLUDE file(s) were processed, use debug level d1 to print a copy of the expanded configuration file to your OMPROUTE trace.

# Appendix:
# Extended Coding Examples for OMPROUTE

The examples included here are not exploiting the enhancements just presented in the previous Appendix.

# Sysplex As Non-Backbone Area: #1a

Area 1.1.1.1

HOST12          HOST13

9.67.111.0

HOST11          HOST14
QDIO            QDIO

9.67.112.0

ABR:
HOST01

HOST02    9.67.101.0

Area 0.0.0.0

Two non-backbone Area Hosts connected only via XCF connections
Two non-backbone Area Hosts connected also via OSA-E in QDIO mode
One backbone interior router (HOST02) connected to ABR

ⓒ IBM Corporation 2004 - 2011

1. These simple diagrams do NOT necessarily depict "best practices."  For example, we often recommend the use of Private IP addresses for interior links like XCF connections, because the best practice would be to implement VIPA addresses as Public Addresses in each of the LPARs.  With VIPA as our target IP addresses, there would be no need to "waste" a set of Public Addresses on the XCF connections.
2. Nevertheless, to minimize the amount of coding we need to show you in a limited space, we depict only the physical adapter addresses and not any assumed VIPA addresses.

## OMPROUTE Configuration: Hosts 12 & 13

```
Area
    Area_number=1.1.1.1;

OSPF_Interface
    IP_address=9.67.111.*
    Name=DYNXCF
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=55296;
```

Area 1.1.1.1

HOST12    HOST13

9.67.111.0

9.67.112.0

9.67.101.0    Area 0.0.0.0

1. This configuration file can be used on both of the hosts that are attached only to the XCF.
2. Area 1.1.1.1 needs to be defined.
3. By using a wildcard definition for the XCF interface, the OSPF_INTERFACE statement can be used unchanged on all of the hosts attached to the XCF. Note that the interface name is a dummy on wildcard interfaces
4. Because these two hosts only have one interface (the XCF) it is not necessary to define the ROUTER_ID -- it will be the address of the XCF interface since OMPROUTE picks one of the configured interfaces for router id and that is the only interface available.
5. If you follow a "best practices" scenario, you would code at least one VIPA address in each z/OS image and define it as the Router_ID.

# OMPROUTE Configuration: Hosts 11 & 14

```
RouterId = 9.67.111.11; [host 11]
RouterID = 9.67.111.14; [host 14]

Area
   Area_number=1.1.1.1;

OSPF_Interface
   IP_address=9.67.111.*
   Name=DYNXCF
   Subnet_mask=255.255.255.0
   Attaches_to_area=1.1.1.1
   MTU=55296;

OSPF_Interface
   IP_address=9.67.112.*
   Name=OSAGBE11
   Subnet_mask=255.255.255.0
   Attaches_to_area=1.1.1.1
   MTU=8992;
```

Area 1.1.1.1

9.67.111.0

HOST11    HOST14

QDIO    QDIO

9.67.112.0

9.67.101.0    Area 0.0.0.0

© IBM Corporation 2004 - 2011

1. Because this router has two interfaces, we will explicitly define the router id so that we know what it is.  To be consistent with the other XCF hosts, the XCF interface is used on as the router id in this example.  Wildcard addresses are not allowed on the RouterID definition, so the same configuration file can't be used on both hosts.  Except for the routerid statements, however, the configuration files are identical copies (because of the use of wildcard definitions for each interface)
2. Once again, area 1.1.1.1 must be defined.
3. The XCF interface definition is copied from the config files for hosts 12 and 13.
4. The OSA-express is configured as a Gigabit ethernet in this example

# OMPROUTE Configuration: Host 01

```
RouterID=9.67.101.1;

Area
   Area_number=1.1.1.1;

Area
   Area_number=0.0.0.0;

OSPF_Interface
   IP_address=9.67.112.*
   Name=OSAGBE11
   Subnet_mask=255.255.255.0
   Attaches_to_area=1.1.1.1
   MTU=8992;

OSPF_Interface
   IP_address=9.67.101.1
   Name=CTC1TO2
   Subnet_mask=255.255.255.252
   Attaches_to_area=0.0.0.0
   Destination_address=9.67.101.2
   MTU=65527;
```
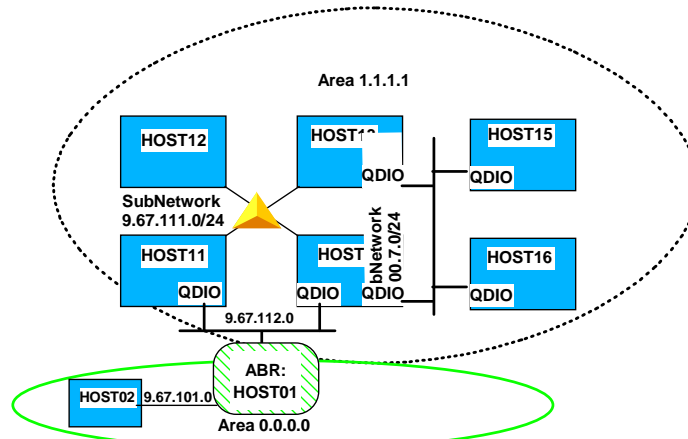
Area 1.1.1.1

9.67.111.0

QDIO   QDIO

9.67.112.0

ABR: HOST01

9.67.101.0

Area 0.0.0.0

ⒸIBM Corporation 2004 - 2011

1. This host is the area-border router between areas 0.0.0.0 (the backbone) and 1.1.1.1 (the sysplex's area). The mere presence of more than one Area statement makes this router an ABR.
2. An area is coded statement for each area
3. The OSA Express definition is copied from the other hosts on the network
4. The point-to-point link to host 2 is also defined.
5. There would probably be additional OSPF interfaces on this host in a real configuration.

# OMPROUTE Configuration: Host 02

```
RouterID=9.67.101.2;

OSPF_Interface
   IP_address=9.67.101.2
   Name=CTC2TO1
   Subnet_mask=255.255.255.252
   Attaches_to_area=0.0.0.0
   Destination_address=9.67.101.1
   MTU=65527;
```

Area 1.1.1.1

9.67.111.0

QDIO    QDIO

9.67.112.0

9.67.101.0

HOST02

Area 0.0.0.0

1. The main thing this configuration shows is that it is not necessary to define the backbone area (0.0.0.0) in routers that are in the backbone.

Sysplex As Non-Backbone Area: #1b

1. This example is the same as the previous one, except with the addition of the new Gigabit Ethernet subnetwork 9.100.7.0.  For purposes of this example, assume that we do not want that subnetwork advertised beyond area 1.1.1.1.

# OMPROUTE Configuration: Host 01

```
RouterID=9.67.101.1;

Area
    Area_number=1.1.1.1;

Area
    Area_number=0.0.0.0;

Range
    IP_Address=9.100.7.0
    Subnet_mask=255.255.255.0
    Area_number=1.1.1.1
    Advertise=No;

OSPF_Interface
    IP_address=9.67.112.*
    Name=OSAGBE11
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=8992;

OSPF_Interface
    IP_address=9.67.101.1
    Name=CTC1TO2
    Subnet_mask=255.255.255.252
    Attaches_to_area=0.0.0.0
    Destination_address=9.67.101.2
    MTU=65527;
```

Area 1.1.1.1

SubNetwork
9.67.111.0/24

9.67.112.0

ABR:
HOST01

9.67.101.0

Area 0.0.0.0

Ⓒ IBM Corporation 2004 - 2011

1. This is identical to the previous Host 01 configuration except for the addition of the Range statement.
2. In this example, we don't want the the 9.100.7.0 subnet advertised into the backbone area.  We accomplish this with the Range configuration statement.  First we define the range as all the addresses in the 9.100.7.0/24 subnet that exist in Area 1.1.1.1. Then we add the Advertise=No, which tells the Area-border router not to advertise this range into other areas.
3. This is the only type of "filtering by ip address" support that OMPROUTE has for OSPF.
    1. The RANGE statement may be used either for Route Aggregation or for Filtering.

# OMPROUTE Configuration: Host 14
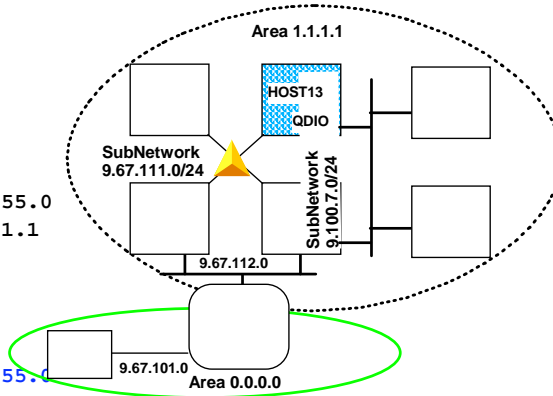
```
RouterID = 9.67.111.14;

Area
    Area_number=1.1.1.1;

OSPF_Interface
    IP_address=9.67.111.*
    Name=DYNXCF
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=55296;
OSPF_Interface
    IP_address=9.100.7.*
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=8992;
OSPF_Interface
    IP_address=9.67.112.*
    Name=OSAGBE11
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=8992;
```
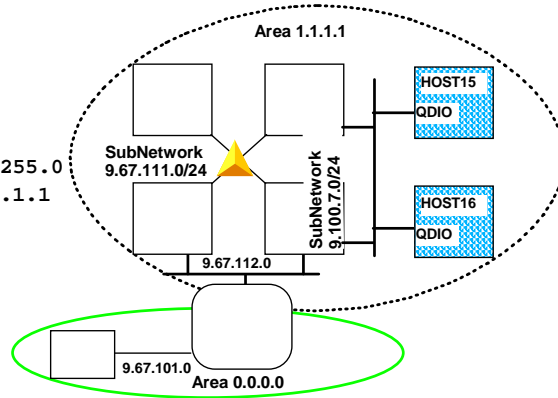
**Area 1.1.1.1**

**SubNetwork 9.67.111.0/24**

**HOST14**

**QDIO**

twork
7.0/24

QDIO

9.67.112.0

9.67.101.0

**Area 0.0.0.0**

© IBM Corporation 2004 - 2011

1. The only change from the previously shown configuration for this host is the addition of the OSPF_Interface for the Gigabit Ethernet network for subnetwork 9.100.7.0.
2. We can share this Gigabit Ethernet definition with HOST13 because we have used a wildcard in its definition. Therefore, here you see two ways of defining the Gigabit Ethernet interface to OMPROUTE: one uses a wildcard with no NAME parameter; the other uses an explicit IP address with the exact Gigabit Interface Linkname.

# OMPROUTE Configuration: Host 13

```
RouterID=9.67.111.13;

Area
    Area_number=1.1.1.1;

OSPF_Interface
    IP_address=9.67.111.*
    Name=DYNXCF
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=55296;

OSPF_Interface
    IP_address=9.100.7.*
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=8992;
```

**Area 1.1.1.1**

HOST13
QDIO
**SubNetwork 9.67.111.0/24**
**SubNetwork 9.100.7.0/24**
9.67.112.0
9.67.101.0
**Area 0.0.0.0**

1. Because host 13 now has two interfaces, it is necessary to define the OSPF router id if we want to ensure that we get the value we want.
2. Also, the OSPF_interface for the 9.100.7.0 Gigabit Ethernet network is added.   Note once again that we are sharing the same definition for this Gigabit Ethernet definition to the 9.100.7.0 network as the one we used for HOST14.
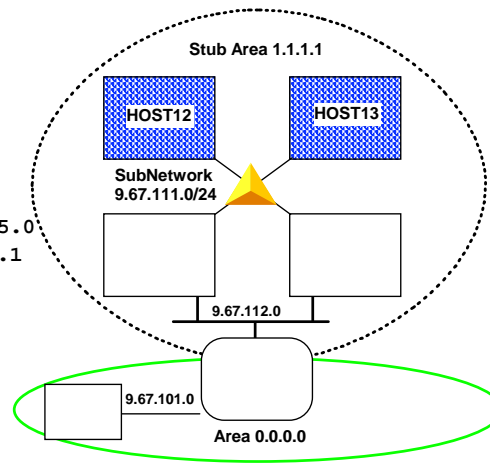
# OMPROUTE Configuration: Hosts 15 & 16

```
Area
    Area_number=1.1.1.1;

OSPF_Interface
    IP_address=9.100.7.*
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=8992;
```

**Area 1.1.1.1**

**SubNetwork**
**9.67.111.0/24**

**SubNetwork**
**9.100.7.0/24**

HOST15
QDIO

HOST16
QDIO

**9.67.112.0**

**9.67.101.0**
**Area 0.0.0.0**

1. Like hosts 12 and 13 in the previous configuration, these hosts have a relatively simple configuration for attachment via Gigabit Ethernet to the OSPF network.  Each may use a wildcarded definition.

# Sysplex Attached as Stub Area: #2

**Stub Area 1.1.1.1**

HOST12    HOST13

SubNetwork
9.67.111.0/24

HOST11    HOST14
QDIO      QDIO

9.67.112.0

ABR:
HOST01

HOST02  9.67.101.0

**Area 0.0.0.0**

**Two Stub Area Hosts connected only via XCF connections**
**Two Stub Area Hosts connected also via OSA-E in QDIO mode**
**One backbone interior router (HOST02) connected to ABR**

# OMPROUTE Configuration: Hosts 12 & 13

```
Area
    Area_number=1.1.1.1
    Stub_area=yes;

OSPF_Interface
    IP_address=9.67.111.*
    Name=DYNXCF
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=55296;
```

Stub Area 1.1.1.1

HOST12          HOST13

SubNetwork
9.67.111.0/24

9.67.112.0

9.67.101.0

Area 0.0.0.0

1. The difference between this configuration and the first configuration for these hosts is the addition of Stub=Yes to the area statement for area 1.1.1.1
2. Because 1.1.1.1 is now defined as a stub area, OSPF external routes (like static or direct routes) will not be imported into it.
3. OMPROUTE does not support Not So Stubby Area (NSSA), which is a method defined by RFC 1537 for limited importation of external routes.

# OMPROUTE Configuration: Hosts 11 & 14

```
RouterId = 9.67.111.11; [host 11]
RouterID = 9.67.111.14; [host 14]

Area
    Area_number=1.1.1.1
    Stub_area=yes;

OSPF_Interface
    IP_address=9.67.111.*
    Name=DYNXCF
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=55296;

OSPF_Interface
    IP_address=9.67.112.*
    Name=OSAGBE11
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=8992;
```

**Stub Area 1.1.1.1**

**SubNetwork 9.67.111.0/24**

**HOST11**      **HOST14**

9.67.112.0

9.67.101.0

**Area 0.0.0.0**

1. This configuration is identical to the first configuration for these two hosts except for the addition of Stub=yes on the area statement for 1.1.1.1

# OMPROUTE Configuration: Host 01

```
RouterID=9.67.101.1;

Area
    Area_number=1.1.1.1
    Stub_area=yes
    Stub_default_cost=10;

Area
    Area_number=0.0.0.0;

OSPF_Interface
    IP_address=9.67.112.*
    Name=OSAGBE11
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=8992;

OSPF_Interface
    IP_address=9.67.101.1
    Name=CTC1TO2
    Subnet_mask=255.255.255.252
    Attaches_to_area=0.0.0.0
    Destination_address=9.67.101.2
    MTU=65527;
```
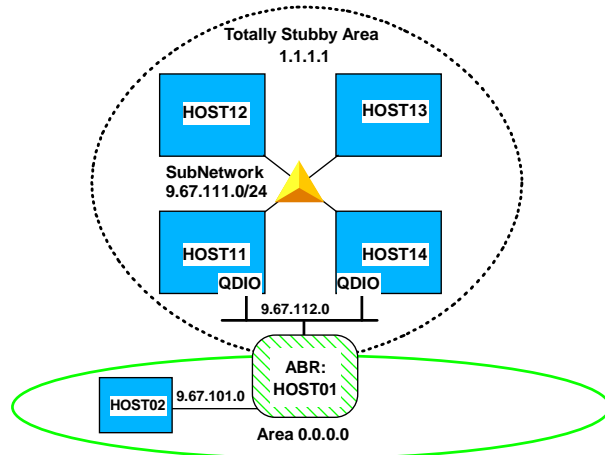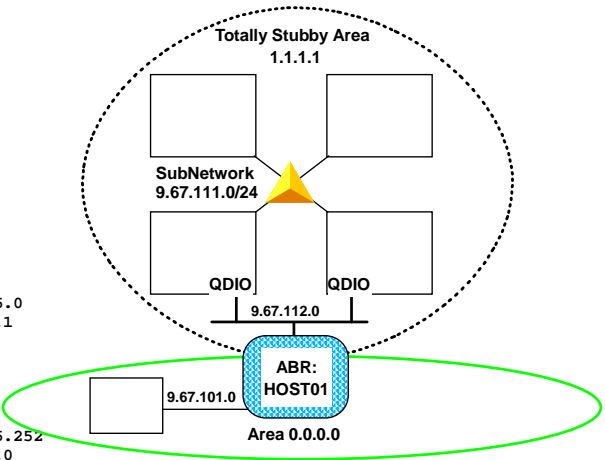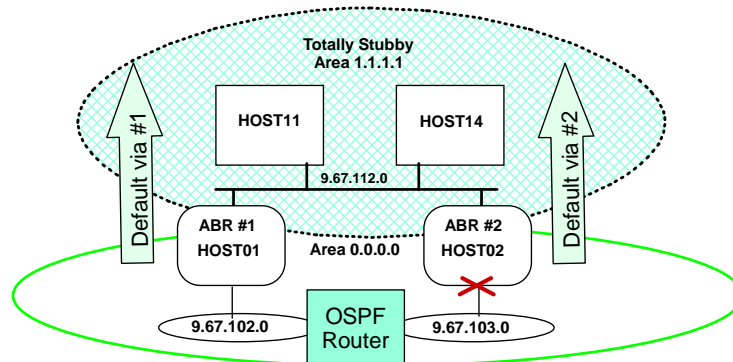
Stub Area 1.1.1.1

SubNetwork 9.67.111.0/24

9.67.112.0

ABR: HOST01

9.67.101.0

Area 0.0.0.0

1. This host is the area-border router between areas 0.0.0.0 (the backbone) and 1.1.1.1 (the sysplex's stub area)
2. An area statement for each area is required
3. This definition is the same as the first definition for this host except for the addition of the stub area information. The area definition for a stub area at an area-border router is more interesting than it is for a router that resides wholly inside the stub area.
4. An area border router advertises a default route into each stub area it attaches to. Routers inside the stub area will use the default route to reach OSPF external destinations via the area-border router. The stub_default_cost parameter specifies the OSPF cost that will be assigned to the default route that this router advertises into the stub area.
5. OMPROUTE will advertise summary LSAs for other areas into stub areas unless commanded not to. Because of this, a default OMPROUTE stub area is not what is called a "Totally Stubby Area." But note that it is also NOT a "Not So Stubby Area" (NSSA) as defined in RFC 1587, because it does not import type 7 LSAs for OSPF external destinations. (OSPF in OMPROUTE for OS/390 or z/OS does not support Type 7 LSAs.) Basically all the optimization you are getting in this case is that OSPF External LSAs are not propagated into the stub area, plus it can't be used as a transit network for virtual links, plus all the area-border routers advertise default routes into the stub area.
    1. Assume your ABR is also an ASBR and you have coded the parameters we have seen earlier in the AS_Boundary_Routing statement: "Originate_Default_Route=Yes" and "Default_Route_Cost=n." Even with these parameters, the Stub Area will receive a Type 3 Summary LSA for a default route with -- NOT the Default_Route_Cost -- but rather the ("Stub_Default_Cost" plus 1). In fact, the "Originate_Default_Route=" parameter has no effect whatsoever on the generation of the Default Route for Stub Areas. The Stub Default and its cost are derived solely from the OSPF protocol governing such exchanges and by the coding (or default setting) of "Stub_Default_Cost."

# Sysplex As Totally Stubby Area: #3

**Totally Stubby Area**
**1.1.1.1**

HOST12    HOST13

**SubNetwork**
**9.67.111.0/24**

HOST11    HOST14
QDIO      QDIO

**9.67.112.0**

**ABR:**
**HOST01**

HOST02  **9.67.101.0**

**Area 0.0.0.0**

**Two Totally Stubby Area Hosts connected only via XCF connections**
**Two Totally Stubby Area Hosts connected also via OSA-E in QDIO mode**
**One backbone interior router (HOST02) connected to ABR**

# OMPROUTE Configuration: Host 01

```
RouterID=9.67.101.1;

Area
    Area_number=1.1.1.1
    Stub_area=yes
    Stub_default_cost=10
    Import_summaries=no;


Area
    Area_number=0.0.0.0;

OSPF_Interface
    IP_address=9.67.112.*
    Name=OSAGBE11
    Subnet_mask=255.255.255.0
    Attaches_to_area=1.1.1.1
    MTU=8992;

OSPF_Interface
    IP_address=9.67.101.1
    Name=CTC1TO2
    Subnet_mask=255.255.255.252
    Attaches_to_area=0.0.0.0
    Destination_address=9.67.101.2
    MTU=65527;
```

**Totally Stubby Area**
**1.1.1.1**

**SubNetwork**
**9.67.111.0/24**

QDIO    QDIO

**9.67.112.0**

**ABR:**
**HOST01**

9.67.101.0

**Area 0.0.0.0**

©IBM Corporation 2004 - 2011

1. By adding import_summaries=no to the area statement for stub area 1.1.1.1 at the area-border router, we have changed it into a true OSPF Stub area (totally stubby area).
2. Now the ONLY route that is not wholly contained within the stub area that this area-border router will advertise into the stub area is the default route. This can be a signficant performance optimization for hosts in the stub area and will greatly simplify their routing tables.
3. The drawback is that if there is more than one area-border router serving the stub area, they will only advertise default routes into the stub area, so hosts inside the totally stubby area cannot pick an area-border router based on destination address. In other words, if there are multiple area-border routers serving a totally stubby area, they must be redundant connections.

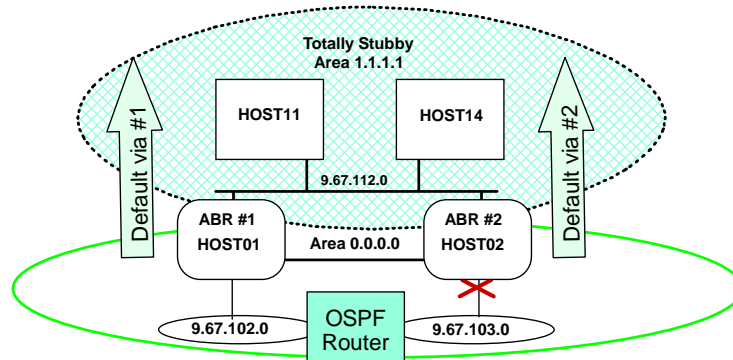How to Interconnect ABRs for Totally Stubby Area? Ex. 1

- ABRs interconnected via Link in Totally Stubby Area; two Default Gateways for each HOSTnn.
  - LSA about loss of connectivity to subnet 9.67.103.0 via ABR #2 is NOT sent into Totally Stubby Area 1.1.1.1 by either ABR #1 or ABR #2, since only defaut routes are sent into Totally Stubby Areas.
- As far as hosts in the Totally Stubby Area know, all destinations outside the area can be reached via either ABR. They will continue sending packets for 9.67.103.0 to both ABR #1 and/or ABR#2 (depending on multipath settings). The packets sent via ABR #2 wil not be delivered
- This is an undesirable way to configure a Totally Stubby Area. It violates the rule that it should never matter which ABR is used to reach any destination outside the area.
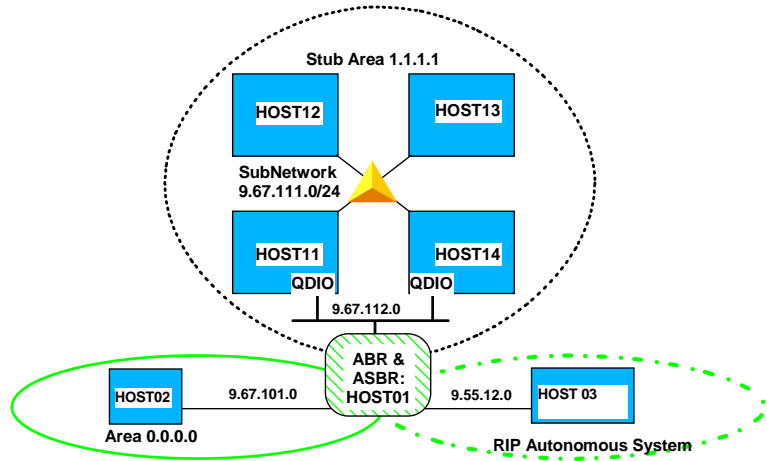
1. One ABR in a stub area does not accept the summary default route from another ABR in the same stub area. In other words, ABR #2 will not accept the stub area summary default route from ABR #1, so packets sent to the backbone via ABR #2 cannot be delivered by going to ABR #1 through the totally stubby area
2. Dead gateway processing may kick in when packets fail to be delivered via ABR #2, but only in some cases, and its recovery is not robust enough to give the redundancy that should be expected when two ABRs attach a totally stubby area to the backbone.

# How to Interconnect ABRs for Totally Stubby Area? Ex. 2

**Totally Stubby
Area 1.1.1.1**

Default via #1

HOST11

HOST14

9.67.112.0

ABR #1
HOST01

Area 0.0.0.0

ABR #2
HOST02

Default via #2

9.67.102.0

OSPF
Router

9.67.103.0

- ABRs interconnected via Link in Backbone Area 0 AND Link in Area 1; two Default Gateways for each HOST1n.
  - As before, LSA about loss of connectivity to subnet 9.67.103.0 via ABR #2 is NOT sent into Totally Stubby Area 1.1.1.1 by either ABR #1 or ABR #2 since only default routes are sent into Totally Stubby Areas
  - Packets sent to ABR#2 for 9.67.103.0 will still reach the destination, via ABR#1 and the link in Area 0
  - HOST11 and HOST14 know of two default gateways; regardless of which gateway they choose (they have coded IPCONFIG MULTIPATH), the packets can arrive at 9.67.103.0

# Sysplex Attached to ASBR/ABR #4



**Stub Area 1.1.1.1**

HOST12   HOST13

**SubNetwork 9.67.111.0/24**

HOST11   HOST14
QDIO     QDIO

9.67.112.0

**ABR & ASBR: HOST01**

HOST02   9.67.101.0

HOST 03   9.55.12.0

**Area 0.0.0.0**

**RIP Autonomous System**

ABR is now also an ASBR attached to the backbone and to the RIP autonomous system.

# OMPROUTE Configuration: Host 01

```
RouterID=9.67.101.1;

AS_Boundary_Routing
    Import_RIP_Routes=Yes
    Originate_Default_Route=Yes;

Area
    Area_number=1.1.1.1
    Stub_area=yes
    Stub_default_cost=10;

Area
    Area_number=0.0.0.0;

[OSPF Interfaces elided]

RIP_interface
    IP_address=9.55.12.1
    Name=CTC1TO3
    Subnet_mask=255.255.255.0
    RIPV2=YES
    [send_only=default]  <== see notes!
    MTU=10000;

Originate_RIP_default
    condition=always;
```
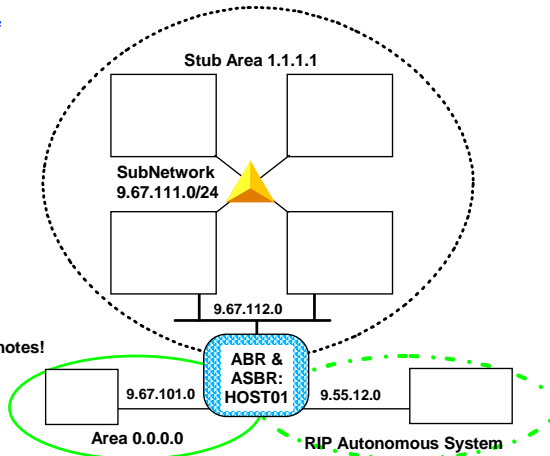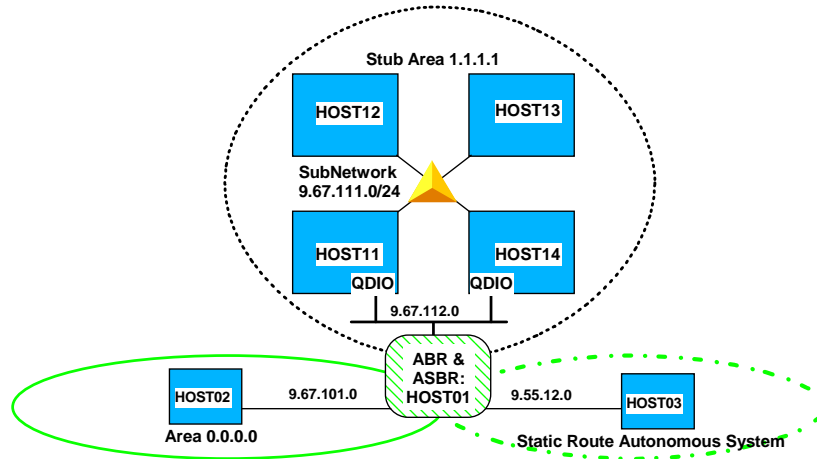
**Stub Area 1.1.1.1**

**SubNetwork 9.67.111.0/24**

9.67.112.0

**ABR & ASBR: HOST01**

9.67.101.0        9.55.12.0

**Area 0.0.0.0**        **RIP Autonomous System**

1. This host is the area-border router between areas 0.0.0.0 (the backbone) and 1.1.1.1 (the sysplex's stub area). It is now also an Autonomous System Boundary router (ASBR) between the RIP and OSPF AS's
2. The AS_Boundary_Routing statement is an OSPF statement.  It tells this router to import RIP routes into the OSPF AS, and to originate an OSPF default route into the OSPF AS.  This OSPF statement does NOT affect what OMPROUTE advertises into the RIP AS, that is defined by RIP definition statements.
3. A RIP_interface is defined for the interface that connects to the RIP AS
4. the originate_RIP_default statement tells OMPROUTE to always  advertise a default route into the RIP AS.
5. So this router advertises default routes into all networks to which it is attached. This is a pretty important router!
6. If  this host is the only ASBR for the RIP subnet and you don't want to flood your RIP network with the contents of the OSPF network (which could get quite large), you can add the send_only=default parameter to your RIP_interface statement.  This will prevent this host from importing any routes into the RIP AS except for the default route.  As a result, all hosts in the RIP AS will route to this host for destinations not known within the RIP AS. This is  NOT advised if this host is a needed router **within** the RIP AS (i.e., if any RIP hosts depend on this  host to route to other hosts within the same RIP AS).

# Sysplex Attached to ASBR/ABR #5

**Stub Area 1.1.1.1**

HOST12  HOST13

**SubNetwork 9.67.111.0/24**

HOST11  HOST14

QDIO  QDIO

9.67.112.0

**ABR & ASBR: HOST01**

HOST02  9.67.101.0  9.55.12.0  HOST03

**Area 0.0.0.0**

**Static Route Autonomous System**

ABR is now also an ASBR attached to backbone and to autonomous system of Static Routes.

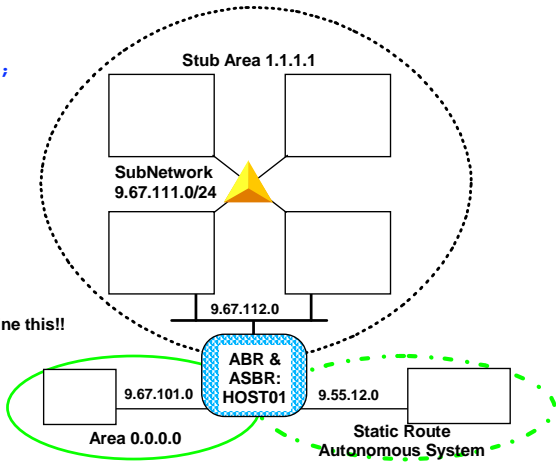# OMPROUTE Configuration:  Host 01

```
RouterID=9.67.101.1;

AS_Boundary_Routing
    Import_Static_Routes=yes
    Import_Direct_Routes=Yes
    Originate_Default_Route=Yes;

Area
    Area_number=1.1.1.1
    Stub_area=yes
    Stub_default_cost=10;

Area
    Area_number=0.0.0.0;

[OSPF Interfaces elided]

Interface    <=== don't forget todefine this!!
          (see notes)
    IP_address=9.55.12.1
    Name=CTC1TO3
    Subnet_mask=255.255.255.0
    MTU=10000;
```

**Stub Area 1.1.1.1**

**SubNetwork
9.67.111.0/24**

**9.67.112.0**

**ABR &
ASBR:
HOST01**

**9.67.101.0**

**9.55.12.0**

**Area 0.0.0.0**

**Static Route
Autonomous System**

Ⓒ IBM Corporation 2004 - 2011

1. This host is the area-border router between areas 0.0.0.0 (the backbone) and 1.1.1.1 (the sysplex's stub area). It is now also an Autonomous System Boundary router (ASBR) between the  OSPF AS and an AS that is using static routing, meaning that any routes needed to reach destinations in the static AS are defined in the TCP/IP profile in a GATEWAY or BEGINROUTES statement.  OMPROUTE is using NO routing protocol in the  static AS.
2. The AS_Boundary_Routing statement is an OSPF statement.  It tells this router to import static routes into the OSPF AS, and to originate an OSPF default route into the OSPF AS.   It also tells OMPROUTE to import direct routes into the OSPF AS.  In this case, that means that OMPROUTE will advertise a direct route to the subnet that the non-routing interface is attached to, into the OSPF network (in this example, that subnet is 9.55.12.0/24).  A direct route is one that is not learned from a routing protocol, but rather from OSPF's knowledge of an Interface.
3. An "Interface" statement  is defined for the interface that connects to the static AS.  The Interface statement indicates that OMPROUTE is not communicating **any** routing protocol over that interface.
4. **Please note** that even though OMPROUTE is not communicating a routing protocol over the CTC1TO3 interface, it should still be defined to OMPROUTE, using the Interface statement. If it is not, then OMPROUTE will use a default value for its MTU size (576) AND update the stack's MTU value with that value. Also, if the interface is not defined to OMPROUTE, OMPROUTE will use the class mask for the interface's subnet mask  (in this case 255.0.0.0), with possible undesirable results (for example, a direct route will be  to 9.0.0.0 over that interface will be advertised and added to the TCP/IP routing table!)

# References, Bibliography

# For More Information ...

**URL**

| | |
|---|---|
| http://www.ibm.com/servers/eserver/zseries | IBM Enterprise Servers (z900 & S/390) |
| http://www.ibm.com/servers/eserver/zseries/networking | z900 Networking |
| http://www.ibm.com/servers/eserver/zseries/networking/technology.html | Networking White Papers and Information |
| http://www.ibm.com/software/network | Networking & Communications Software |
| http://www.ibm.com/software/network/commserver | Communications Server |
| http://www.ibm.com/software/network/commserver/library | CS White Papers, Product Doc, etc. |
| http://www.redbooks.ibm.com | ITSO Redbooks |
| http://www.rfc-editor.org/rfcsearch.html | RFCs |
| http://www.ibm.com/support/techdocs/ | Advanced Technical Support  (Flashes, Presentations, White Papers, etc.) |

**Whitepaper on OSPF with IBM & CISCO:**
**"OSPF Design and Interoperability Recommendations for Catalyst 6500 and OSA-Express Environments"  http://www-1.ibm.com/servers/eserver/zseries/networking/pdf/ospf_design.pdf**
**(Beware of HELLO and DR Timers -- they are set too short for most production networks.)**

**Redbook on "Networking with z/OS and Cisco Routers:  An Interoperability Guide (SG24-6297)**

# Bibliography

| | |
|---|---|
| **RFCs 1583 & 2328** | OSPF V2 by John Moy:  URL http://www.ietf.cnri.reston.va |
| **GG24-3376** | TCP/IP Tutorial and Technical Reference |
| **SC31-8513** | OS/390 Communications Server:  IP Configuration |
| **Huitema, Christian** | Routing in the Internet, Prentice Hall |
| **Black, Uyless** | IP Routing Protocols:  RIP, OSPF, BGP, PNNI & CISCO Routing Protocols  (SR23-9498) |
| **Parkhurst, William R.** | Cisco Router OSPF Design and Implementation Guide, Parkhurst McGraw-Hill (SR23-8683) |
| **Cisco Systems** | CCIE Fundamentals:  Network Design and Case Studies, 2nd Edition (SR23-9437) |
| **Doyle, Jeff (Cisco Systems)** | CCIE Professional Development:  Routing TCP/IP, Volume 1 (SR23-9241) |

# Bibliography Redbooks (V1R10-V1R12)

| | |
|---|---|
| **SG24-7696 (V1R10)** <br> **SG24-7798 (V1R11)** <br> **SG24-7896 (V1R12)** | IBM z/OS VnRxx Communications Server TCP/IP Implementation, Volume 1 - Basic Functions, Connectivity, and Routing |
| **SG24-7697 (V1R10)** <br> **SG24-7799 (V1R11)** <br> **SG24-7897 (V1R12)** | IBM z/OS VnRxx Communications Server TCP/IP Implementation, Volume 2 - Standard Applications |
| **SG24-7698 (V1R10)** <br> **SG24-7800 (V1R11)** <br> **SG24-7898 (V1R12)** | IBM z/OS VnRxx Communications Server TCP/IP Implementation, Volume 3 - High Availability, Scalability, and Performance |
| **SG24-7699 (V1R10)** <br> **SG24-7801 (V1R11)** <br> **SG24-7899 (V1R12)** | IBM z/OS VnRxx Communications Server TCP/IP Implementation, Volume 4 - Security and Policy-Based Networking |

© IBM Corporation 2004 - 2011

1. These redbooks are all now available from www.redbooks.ibm.com.

# Bibliography Redbooks (V1R7, V1R8)

| | |
|---|---|
| **SG24-7169** | Communications Server for z/OS V1R7 TCP/IP Implementation, Volume 1 - Base Functions, Connectivity, and Routing |
| **SG24-7170** | Communications Server for z/OS V1R7 TCP/IP Implementation, Volume 2 - Standard Applications |
| **SG24-7171** | Communications Server for z/OS V1R7 TCP/IP Implementation, Volume 3 - High Availability, Scalability, and Performance |
| **SG24-7172** | Communications Server for z/OS V1R7 TCP/IP Implementation, Volume 4 - Security and Policy |
| **SG24-7339** | Communications Server for z/OS V1R8 TCP/IP Implementation, Volume 1: Base Functions, Connectivity, and Routing |
| **SG24-7340** | Communications Server for z/OS V1R8 TCP/IP Implementation, Volume 2 - Standard Applications |
| **SG24-7341** | Communications Server for z/OS V1R8 TCP/IP Implementation, Volume 3 - High Availability, Scalability, and Performance |
| **SG24-7342** | Communications Server for z/OS V1R8 TCP/IP Implementation, Volume 4 - Policy-Based Network Security |
| **SG24-6297-00** | Networking with z/OS and Cisco Routers: An Interoperability Guide |

1. These redbooks are all now available from www.redbooks.ibm.com.

# Bibliography Redbooks (V1R2)

| | |
|---|---|
| **SG24-5227-03** | Communications Server for z/OS V1R2 TCP/IP Implementation Guide Volume 1: Base and TN3270 Configuration |
| **SG24-5228-03** | Communications Server for z/OS V1R2 TCP/IP Implementation Guide Volume 2: UNIX Applications |
| **SG24-6516-00** | Communications Server for z/OS V1R2 TCP/IP Implementation Guide Volume 4: Connectivity and Routing |
| **SG24-6517-00** | Communications Server for z/OS V1R2 TCP/IP Implementation Guide Volume 5: Availability, Scalability, and Performance |
| **SG24-6839-00** | Communications Server for z/OS V1R2 TCP/IP Implementation Guide Volume 6: Policy and Network Management |
| **SG24-6840-00** | Communications Server for z/OS V1R2 TCP/IP Implementation Guide Volume 7: Security |
| **SG24-5229-01** | OS/390 eNetwork Communications Server **V2R7** TCP/IP Implementation Guide Volume 3: MVS Applications |
| **SG24-6297-00** | Networking with z/OS and Cisco Routers: An Interoperability Guide |

© IBM Corporation 2004 - 2011

1. These redbooks are all now available from www.redbooks.ibm.com.
2. The new Volume 4 contains information not only on implementing OSPF in z/OS, but also information on integrating with Cisco Routers, including the integration with EIGRP.
3. Note also that the MVS Applications  volume of the redbook series has not been updated since OS/390 V2R7, although there have been enhancements to several of the socket applications discussed there.  Notably, the CICS sockets application has received important enhancements in the z/OS V1R2 release.  These CICS enhancements are included in one of the presentations in this course.
4. The final redbook mentioned also contains information on using EIGRP and BGP with Cisco.

# End