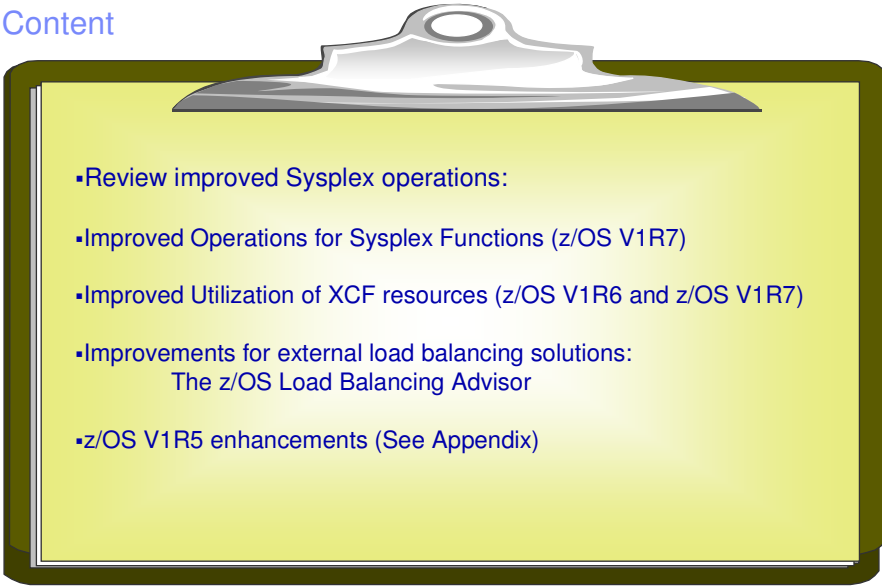# Improved Sysplex Operations
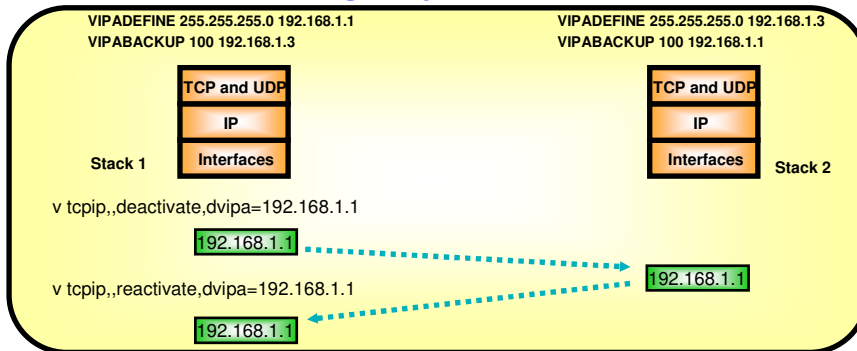
IBM Software Group, Enterprise Networking and Transformation Solutions

## Content

- Review improved Sysplex operations:
- Improved Operations for Sysplex Functions (z/OS V1R7)
- Improved Utilization of XCF resources (z/OS V1R6 and z/OS V1R7)
- Improvements for external load balancing solutions:
  The z/OS Load Balancing Advisor
- z/OS V1R5 enhancements (See Appendix)

## Improved operations: operator-initiated movement of individual stack-managed dynamic VIPA addresses

**VIPADEFINE 255.255.255.0 192.168.1.1**
**VIPABACKUP 100 192.168.1.3**

**VIPADEFINE 255.255.255.0 192.168.1.3**
**VIPABACKUP 100 192.168.1.1**

| TCP and UDP |
| IP |
| Interfaces |

**Stack 1**

| TCP and UDP |
| IP |
| Interfaces |

**Stack 2**

v tcpip,,deactivate,dvipa=192.168.1.1

192.168.1.1

192.168.1.1

v tcpip,,reactivate,dvipa=192.168.1.1

192.168.1.1

➢ **Deactivate**
  ⟩ DVIPA is deactivated and a configured backup stack will takeover the DVIPA
  ⟩ Backup DVIPA can be deactivated also removing eligibility as a backup
➢ **Reactivate**
  ⟩ Original owner can regain ownership
  ⟩ Can also reactivate a backup DVIPA that's been deactivated
  ⟩ Prior to these commands, Vary obey files were needed to cause a DVIPA takeover
  ⟩ These commands can not be used for DVIPAs created from VIPARANGE with bind, ioctl(), or the Modvipa utility

# Determining current status of DVIPAs

Example of a Netstat VIPADCFG/-F report after a DVIPA is deactivated via:
  VARY TCPIP,,SYSPLEX,DEACT,DVIPA=197.11.221.1

```
netstat vipadcfg
MVS TCP/IP NETSTAT CS V1R7       TCPIP Name: TCPCS          19:52:38
Dynamic VIPA Information:

  VIPA Define:
    IpAddr/PrefixLen: 197.11.221.2/24
      Moveable: Immediate  SrvMgr: No

Deactivated Dynamic VIPA Information:

  VIPA Define:
    IpAddr/PrefixLen: 197.11.221.1/24
      Moveable: Immediate  SrvMgr: Yes

  VIPA Distribute:
  Dest:         197.11.221.1..20
    DestXCF:   ALL
      SysPt:   No   TimAff: No    Flg: BaseWLM
  Dest:         197.11.221.1..21
    DestXCF:   ALL
      SysPt:   No   TimAff: No    Flg: BaseWLM
```
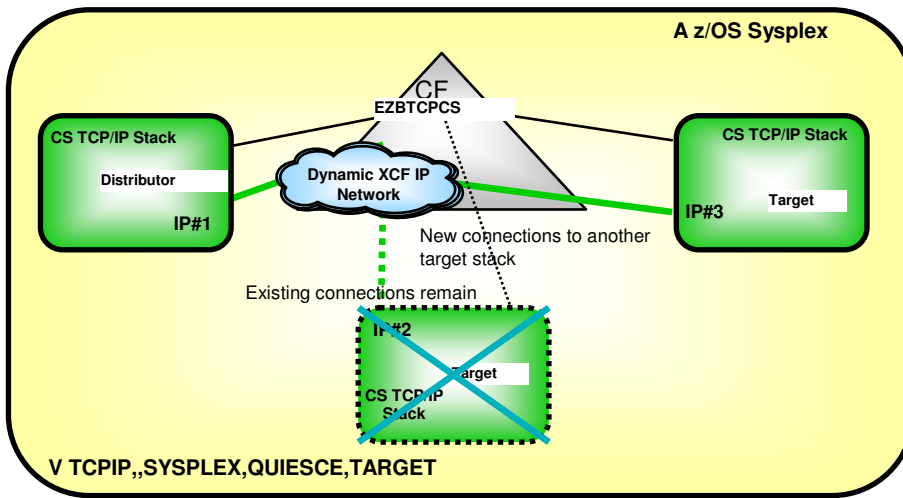
## Improved operations: operator-initiated quiesce/resume of target systems/applications

**A z/OS Sysplex**

CF
EZBTCPCS

**CS TCP/IP Stack**

Distributor

IP#1

Dynamic XCF IP Network

**CS TCP/IP Stack**

Target

IP#3

New connections to another target stack

Existing connections remain

IP #2

Target

CS TCP/IP Stack

**V TCPIP,,SYSPLEX,QUIESCE,TARGET**

# Improved operations: operator-initiated quiesce and resume of individual server applications or full target systems

➢ **Ability to quiesce a system or an application prior to shutdown**
  ⌐ Planned maintenance scenarios of system or application
    – Allows existing systems or applications to drain work queue prior to shutdown
  ⌐ Relieve temporary constraints of resources on target system
  ⌐ Temporary - Does not affect Sysplex Distributor's permanent configuration
  ⌐ Supports applications in SHAREPORT group (must be targets for SD)
  ⌐ Issued on target system being affected
  ⌐ Only way to achieve similar capability earlier was via temporary configuration changes based on OBEYFILE commands

➢ **VARY TCPIP,,SYSPLEX,QUIESCE,options**
  ⌐ TARGET - Quiesces all applications on target stack.
  ⌐ PORT=xxx - Quiesce all applications bound to the specified port on this stack
    – JOBNAME=jobname - Allows quiesce of a single application in SHAREPORT group
    – ASID=asid - Further qualify job being quiesced (such as when dealing with duplicate jobnames)
  ⌐ No new TCP connections sent to the quiesced target (stack or application)
    – For all Distributed DVIPAs that the entity is a target for
  ⌐ Existing TCP connections are maintained (or in other words, the process is non-disruptive)

➢ **VARY TCPIP,,SYSPLEX,RESUME,options**
  ⌐ TARGET|PORT|JOBNAME|ASID
  ⌐ Allows identified target stacks and/or applications to once again be targets for distribution

## When is it safe to recycle a quiesced application without impact to end users?

```
NETSTAT ALL (PORT 21
Client Name: FTPD1                    Client Id: 0000002C
  Local Socket: ::..21
  Foreign Socket: ::..0
    BytesIn:          00000000000000000000
    BytesOut:         00000000000000000000
    SegmentsIn:       00000000000000000000
    SegmentsOut:      00000000000000000000
    Last Touched:      15:49:47         State:         Listen


>>>> more data here

    ConnectionsIn:        0000003121      ConnectionsDropped: 0000000000
    CurrentBacklog:       0000000000      MaximumBacklog:     0000000050
    CurrentConnections:   0000000000      SEF:                100
    Quiesced: Dest
```

Quiesce status

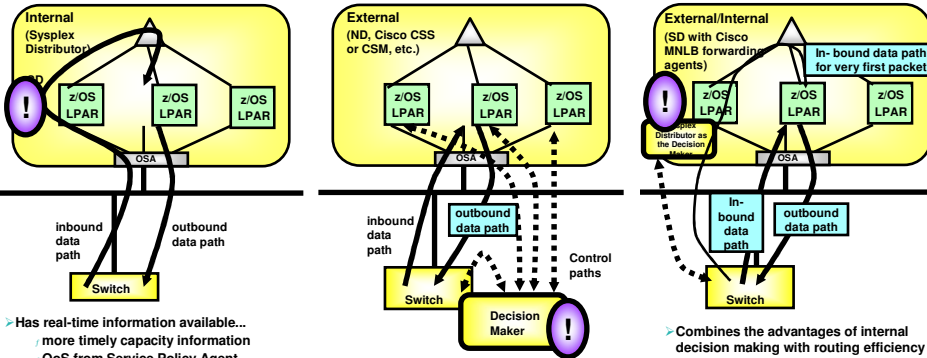No connections being processed and no connections in the backlog

# Improvements for external load balancing solutions:
## The z/OS Load Balancing Advisor

Different ways to perform load balancing in Sysplex

Each has different benefits and drawbacks. Which one you choose to use depends on which of the advantages and disadvantages are most important to you
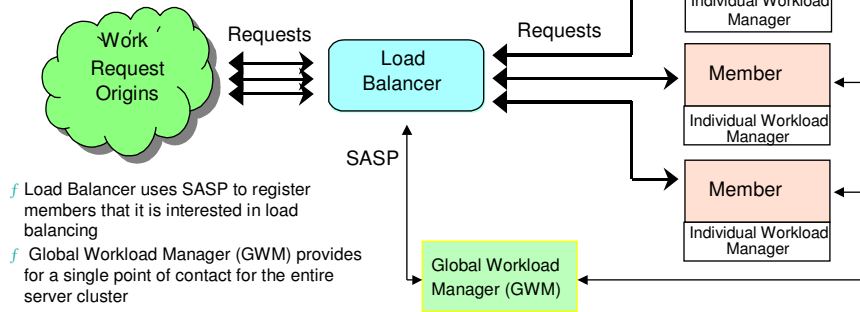
# Problem Statement

➢ Let's assume that you have selected an external IP load balancing solution for your z/OS Sysplex environment
  ƒ Some possible reasons:
    • Prefer to have a single load balancing solution across multiple platforms in your environment
    • Administration of the load balancing functions belongs to network administration domain (not z/OS administrators)
    • Requirements for Content based load balancing
      ▪ Need to perform load balancing/routing decisions based on data content (inspection of url, session IDs, cookies, etc.)
      ▪ This is often combined with additional functions
        ★ SSL offloading functions (Need to decrypt data prior to inspection)
        ★ Web caching functions
➢ But what if the external load balancing solution has no awareness of the sysplex environment?
  ƒ Is the sysplex treated just like any other server cluster?
  ƒ Is it aware of the current/changing workload conditions on the various systems in the cluster?
  ƒ Is it aware of the health and status of applications and/or systems?
➢ Making the external load balancing solution "sysplex aware" can help answer many of these questions
  ƒ The z/OS Load Balancing Advisor is a key component that allows any external load balancing solution to become "sysplex aware".

# Server Application State Protocol (SASP) Architecture

## Server Cluster

Work Request Origins

Requests

Load Balancer

Requests

SASP

Member

Individual Workload Manager

Member

Individual Workload Manager

Member

Individual Workload Manager

Global Workload Manager (GWM)

ƒ Load Balancer uses SASP to register members that it is interested in load balancing

ƒ Global Workload Manager (GWM) provides for a single point of contact for the entire server cluster

  ƒ Provides Load Balancer with recommendations on how work should be balanced across servers

  ƒ Based on current capacity and server availability

ƒ SASP: Open protocol

  ƒ currently being pursued as an individual Internet Draft RFC submitted to the IETF

# Products Supporting SASP

➢ IBM products supporting SASP

   ƒ EWLM (Enterprise Workload Manager)
- Part of IBM Virtualization Engine 1.0
  - Supported Platforms:
    - ★ IBM AIX 5L™ Version 5.2
    - ★ Microsoft® Windows® 2000 Advanced Server, 2000 Server, 2003 Enterprise Edition, 2003 Standard Edition
    - ★ Sun Microsystems Solaris 8 (SPARC Platform Edition), 9 (SPARC Platform Edition)
    - ★ z/OS V1R6 (part of IBM Virtualization Engine Enterprise Workload Manager for z/OS V1.1.0 )

   ƒ z/OS Load Balancing Advisor (became available 4Q2004)
- Part of z/OS Communications Server (z/OS V1R4 and higher)
  - V1R4 (APAR PQ90032) - V1R5 and V1R6 (APAR PQ96293)
    - ★ Documentation of function available at:
      http://www-1.ibm.com/support/docview.wss?uid=swg27005585
- Will be part of the base z/OS V1R7 Communications Server
- Can be used within scope of EWLM z/OS solution or in a z/OS WLM environment

➢ Load Balancers
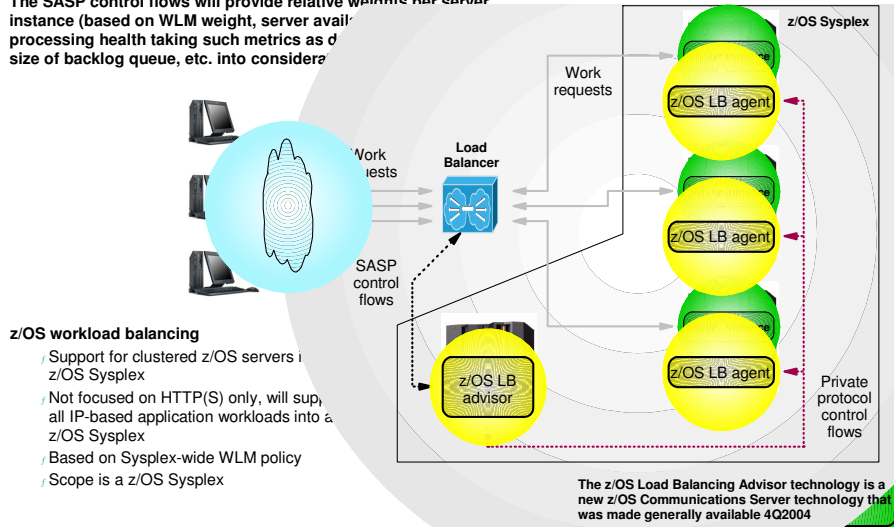
   ƒ CISCO CSM level 4.1 (2.5)
   ƒ Nortel Alteon
   ƒ Other vendors possibly in the future

# z/OS Load Balancing Advisor (LBA) for outboard load balancers

**The SASP control flows will provide relative weights per server instance (based on WLM weight, server avail... processing health taking such metrics as d... size of backlog queue, etc. into considera...**

Work requests

**Load Balancer**

Work requests

SASP control flows

z/OS Sysplex

z/OS LB agent

z/OS LB agent

z/OS LB advisor

z/OS LB agent

Private protocol control flows

**z/OS workload balancing**
- Support for clustered z/OS servers i... z/OS Sysplex
- Not focused on HTTP(S) only, will sup... all IP-based application workloads into a... z/OS Sysplex
- Based on Sysplex-wide WLM policy
- Scope is a z/OS Sysplex

**The z/OS Load Balancing Advisor technology is a new z/OS Communications Server technology that was made generally available 4Q2004**

## Load Balancer Advisor Configuration

Advisor IP Address and portA
Group1: Cluster_IP address, port, protocol [TCP/UDP]
- IPx, P1
- IPy, P1
- IPz, P1

Loadbalancer: LBx

SASP

IPz   IPx   IPy

PortA

TCP/IP S1   TCP/IP S2   TCP/IP S1

Server P1   Server P1

Sysplex LB Advisor

Sysplex LB agent   W L M

Sysplex LB agent   W L M

SYSA   SYSB

Agent config

IP socket communications
(Can flow over any network path)

Advisor config

Load Balancer configures...
IP address and port (PortA) of Advisor

One or more groups of homogeneous applications or systems (e.g. Group1)

The cluster IP address clients use to connect to an application (e.g. Cluster_IP address)

The protocol which the group will use (TCP or UDP)
For each group
IP address, and port of each application in the group (e.g. IPx P1, IPy P1, IPz P1)

## What information does the z/OS Load Balancing Advisor provide?

- Status of Applications/Systems
  - Is the application and/or system active?

- Weights indicating how well is each system/application doing. The weights are composed of two main elements:

  - **WLM weight**
    - The WLM weight capacity as we know from other WLM-based load balancing solutions, such as Sysplex Distributor (System or Server-specific)
      - ✓ A numeric value between 0 and 64
      - ✓ Server-specific available on z/OS V1R7

  - **Communications Server weight (Health of the application)**
    - This weight is calculated based on the availability of the actual server instances (are they up and ready to accept workload) and how well TCP/IP and the individual server instances process the workload that is sent to them.
      - ✓ Expressed as a numeric percentage value between 0 and 100
    - Purpose of calculations is to:
      - ✓ Prevent stalled server from being sent more work (accepting no new connections and new connections are being dropped due to backlog queue full condition)
      - ✓ Proactively react to server that is getting overloaded (accepting new connections, but size of backlog queue increases over time approaching the max backlog queue size)

- The final weight is calculated by combining the WLM and the CS weights into a single metric
  - Final weight = WLM weight * CS weight / 100

- For more information see charts for session 3931, "The z/OS Load Balancing Advisor: Making External IP Load Balancers Sysplex Aware"

Weights are a combination of WLM recommendations and QoS information maintained by the TCP/IP stack

# Some strategies for workload request balancing into a z/OS Sysplex

**A** **DNS/WLM as workload balancing should not be used any longer**

**B** **Where HTTP workload is to be balanced based on content of HTTP requests, an outboard load balancer that supports contents inspection must be deployed**

- If HTTPS workload is to be included, the load balancing node must be accompanied by an SSL/TLS offload technology
- Can be combined with a cache appliance for improved performance

**C** **UDP workload balancing must be deployed using an outboard load balancer - SD does not support UDP balancing**

**D** **Remaining TCP connection balancing can be deployed using either SD or an outboard load balancer:**

- SD has more real-time information available than outboard load balancers - even with outboard load balancers using the SASP protocol
- Who is to apply management control over the workload balancing function will be a major factor in deciding which solution to use
- When using an outboard load balancer the use of SASP is recommended as that will help improve the load balancing decisions (makes the outboard load balancers *"sysplex aware"*)

# For More Information....

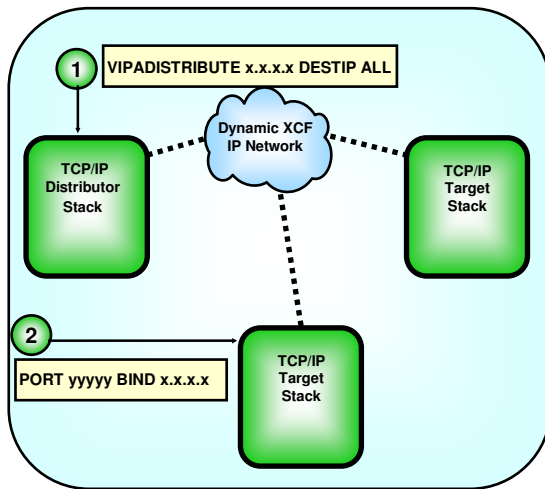| URL | Content |
| --- | --- |
| http://www.ibm.com/servers/eserver/zseries | IBM eServer zSeries Mainframe Servers |
| http://www.ibm.com/servers/eserver/zseries/networking | Networking: IBM zSeries Servers |
| http://www.ibm.com/servers/eserver/zseries/networking/technology.html | IBM Enterprise Servers: Networking Technologies |
| http://www.ibm.com/software/network/commserver | Communications Server product overview |
| http://www.ibm.com/software/network/commserver/zos/ | z/OS Communications Server |
| http://www.ibm.com/software/network/commserver/z_lin/ | Communications Server for Linux on zSeries |
| http://www.ibm.com/software/network/ccl | Communication Controller for Linux on zSeries |
| http://www.ibm.com/software/network/commserver/library | Communications Server products - white papers, product documentation, etc. |
| http://www.redbooks.ibm.com | ITSO redbooks |
| http://www.ibm.com/software/network/commserver/support | Communications Server technical Support |
| http://www.ibm.com/support/techdocs/ | Technical support documentation (techdocs, flashes, presentations, white papers, etc.) |
| http://www.rfc-editor.org/rfcsearch.html | Request For Comments (RFC) |

IBM

Appendix:
Sysplex Enhancements
z/OS V1R5

# IP Sysplex enhancements in z/OS V1R5

ƒ Increase ports on VIPADISTRIBUTE from 4 to 64 (PTFed back to z/OS V1R2 - APAR PQ65205)

ƒ Dynamic port definition for VIPADISTRIBUTE dynamic VIPA when server binds to dynamic VIPA

ƒ Increase limit of DVIPAs per stack from 256 to 1024

ƒ New round-robin distribution method in Sysplex Distributor (PTFed back to z/OS V1R4 - APAR PQ76866)
- Alternative to WLM-based distribution
- Useful where availability is more important than capacity

ƒ Sysplex Distributor affinity
- Configurable timer-based stickyness per source IP address, server DVIPA and port

ƒ Support DVIPA activation based on VIPABACKUP before VIPADEFINE ever processed

# Dynamic Port Assignment (z/OS V1R5)



Allows applications on any port to become automatic candidates for Sysplex Distributor load balancing when they bind to a Distributed DVIPA

1 If PORT is omitted from VIPADISTRIBUTE definition, the Sysplex Distributor reacts dynamically to servers binding to the distributed DVIPA and creating a listening socket, by adding a port to the list of ports for which connection workload balancing will occur.
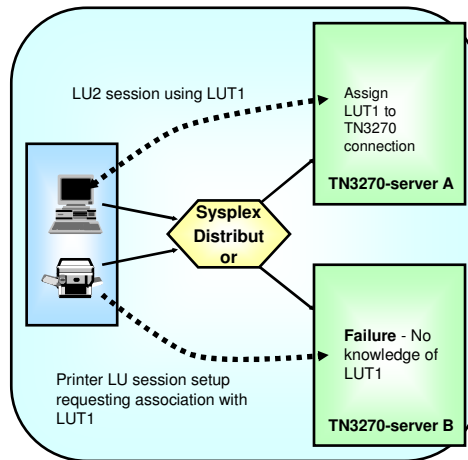
2 Servers either explicitly bind to distributive DVIPA or configure "PORT BIND" definition if binding to INADDR_ANY.

Benefits:
- ƒ Lifts restriction on number of ports that can be distributed for a DVIPA.
- ƒ Minimizes updates to the VIPADISTRIBUTE statement
- ƒ **But** applications being load balanced may be less visible/controlled!
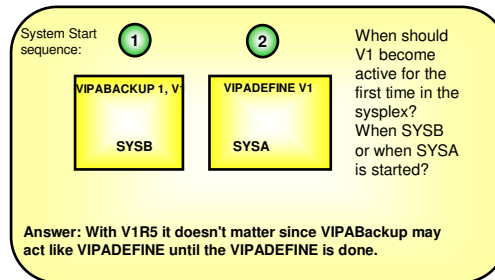
# Application Server Affinity (z/OS V1R5)

➢Optional keyword TIMEDAFFINITY on VIPADISTRIBUTE statement determines whether or not a connection from a client (as identified by source IP address) to a particular server instance of several serviced by Sysplex Distributor shall establish an affinity for future connections from the same client (IP address) to the same Distributed DVIPA and port(s).

➢The affinity duration maintained between connections differs from application to application (configurable)
  ➢The affinity timer begins after the last connection from that client has terminated
    ➢As long as new TCP connections from that client arrive within this interval, the affinity persists
➢Beware of Gateway type of applications!
  ➢All connections from the Gateway application will have affinity to a single server!

LU2 session using LUT1

Assign LUT1 to TN3270 connection

**TN3270-server A**

**Sysplex Distributor**

**Failure** - No knowledge of LUT1

**TN3270-server B**

Printer LU session setup requesting association with LUT1

# Allowing for DVIPA activation when owning system is not started

ƒ Support DVIPA activation based on VIPABACKUP before VIPADEFINE ever processed

ƒ Allows independent startup order for DVIPA owning systems and backup systems
- Backup system will retain ownership of the DVIPA until the primary owning system is started
- Ownership reverts to "VIPADEFINE" stack when that stack is activated

ƒ Requires changes to the VIPABACKUP statements
- MOVEABLE IMMEDIATE, address mask, etc. should be specified
- Also requires that VIPADISTRIBUTE statement is present (after the VIPABACKUP statement) on the Backup TCP/IP stacks

System Start sequence:

**1** **2**

VIPABACKUP 1, V1

SYSB

VIPADEFINE V1

SYSA

When should V1 become active for the first time in the sysplex? When SYSB or when SYSA is started?

**Answer: With V1R5 it doesn't matter since VIPABackup may act like VIPADEFINE until the VIPADEFINE is done.**

IBM

# Trademarks and notices

The following terms are trademarks or registered trademarks of International Business Machines Corporation in the United States or other countries or both:

| | | | |
|---|---|---|---|
| AIX7 | GDDM7 | PrintWay™ | z/Architecture™ |
| AnyNet7 | GDPS7 | PR/SM™ | z/OS7 |
| AS/4007 | HiperSockets™ | pSeries7 | z/VM7 |
| Candle7 | IBM7 | RACF7 | zSeries7 |
| CICS7 | Infoprint7 | Redbooks™ | |
| CICSPlex7 | IMS™ | Redbooks (logo)™ | |
| CICS/ESA7 | IP PrintWay™ | S/3907 | |
| DB27 | iSeries™ | System/3907 | |
| DB2 Connect™ | Language Environment7 | ThinkPad7 | |
| DPI7 | MQSeries7 | Tivoli7 | |
| DRDA7 | MVS™ | Tivoli (logo)7 | |
| e business(logo)7 | MVS/ESA™ | VM/ESA7 | |
| ESCON7 | NetView7 | VSE/ESA™ | |
| eServer™ | OS/27 | VTAM7 | |
| ECKD™ | OS/3907 | WebSphere7 | |
| FFST™ | Parallel Sysplex7 | xSeries7 | |

Cisco, Cisco Systems, the Cisco Systems logo, Catalyst, and Cisco IOS are registered trademarks or trademarks of Cisco Systems, Inc. and/or its affiliates in the United States and certain other countries.