IBM

# Introduction to Parallel Sysplex

Created By: Angelo Corridori
Presented By: Riaz Ahmad
IBM Washington Systems Center
Gaithersburg, Maryland

---

IBM

# Trademarks

The following are trademarks of International Business Machines Corporation.

| | | |
|---|---|---|
| ACF/VTAM | Enterprise System/4381 | Open Blueprint |
| AD/Cycle | Enterprise System/9000 | OpenEdition* |
| ADSM | Enterprise Systems Connection Architecture | OSA |
| Advanced Function Printing | ES/3090 | OSA 1 |
| AFP | ES/4381 | OSA 2 |
| AIX* | ES/9000 | OS/2* |
| AIX/ESA | ESA/370 | OS/390 |
| AOEXPERT/MVS | ESA/390 | OS/400* |
| Automated Operations Expert/MVS | ESCON | Parallel Sysplex |
| CICS/ESA | FASTService* | Power Prestige |
| DataHub | FlowMark | PR/SM |
| DATABASE 2 | Hardware Configuration Definition | PS/2* |
| DataTrade | Hiperbatch | Processor Resource/Systems Manager |
| DB2* | Hipersorting* | RISC System/6000 |
| DFDSM | Hiperspace | S/360 |
| DFSMS | IBM* | S/370 |
| DFSMS/MVS | IBM S/390 Parallel Enterprise Server | S/390 |
| DFSMdfp | IBM S/390 Parallel Enterprise Server - Generation 3 | SAA |
| DFSMSdss | IMS/ESA | SAP R3 |
| DFSMShsm | LANRES | Sysplex Timer |
| DFSMSrmm | Micro Channel* | System/370 |
| Distributed Relational Database | MQ Series | System/390 |
| Architecture | MVS/DFP | Systems Application Architecture* |
| DRDA | MVS/ESA | SystemView |
| Enterprise Systems Architecture/370 | NetView* | VM/ESA |
| Enterprise Systems Architecture/390 | NQS/MVS | VSE/ESA |
| Enterprise System/3090 | OPC | VTAM |
| | | 3090 |

Note: Performance is in Internal Throughput Rate (ITR) ratio based on measurements and projections using standard IBM benchmarks in a controlled environment. The actual throughput that any user will experience will vary depending upon considerations such as the amount of multiprogramming in the user's job stream, the I/O configuration, the storage configuration, and the workload processed. Therefore, no assurance can be given that an individual user will achieve throughput improvements equivalent to the performance ratios stated here.

Actual performance and environmental costs will vary depending on individual customer configurations and conditions.
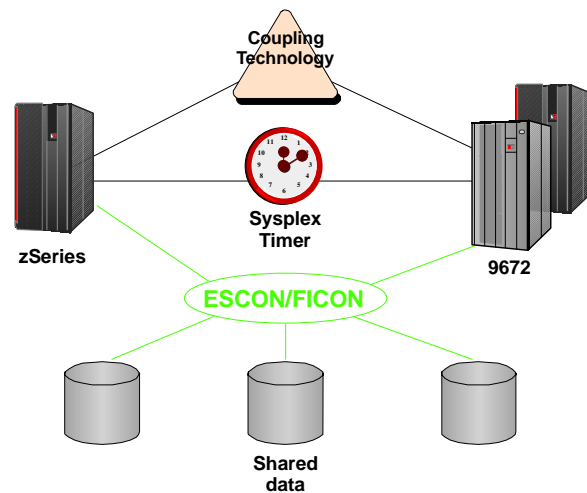
Note: IBM hardware products are manufactured from new parts, or new and used parts. Regardless, our warranty terms apply.

## Agenda

- **Parallel Sysplex Overview**
- **System Structure**
- **Coupling Facility and Link Technology**
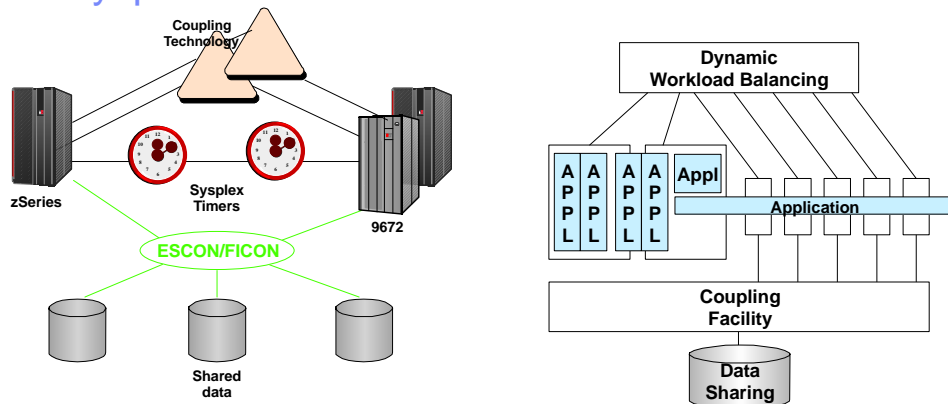- **Parallel Sysplex Software**
- **Summary**

---

## Parallel Sysplex - What is it?

- **Hardware**
  - **Timer**
  - **I/O Connectivity**
  - **Coupling Facility**

- **Software**
  - **XCF/XES**
  - **WLM**

- **Microcode**
  - **CFCC**
  - **Processor u-code**



Coupling Technology

Sysplex Timer

zSeries

9672

ESCON/FICON

Shared data

Parallel Processing!

## Parallel Sysplex Value

Coupling Technology

zSeries

Sysplex Timers

9672

ESCON/FICON

Shared data

Dynamic Workload Balancing

APPL APPL APPL APPL Appl

Application
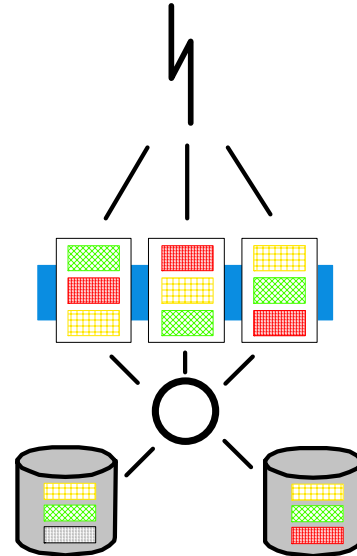
Coupling Facility

Data Sharing

- Continuous Availability
- Flexible Growth
- Scalability
- Reduced Cost
- Leverage S/390 Investment

---

## Agenda

- **Parallel Sysplex Overview**
- **System Structure**
- **Coupling Facility and Link Technology**
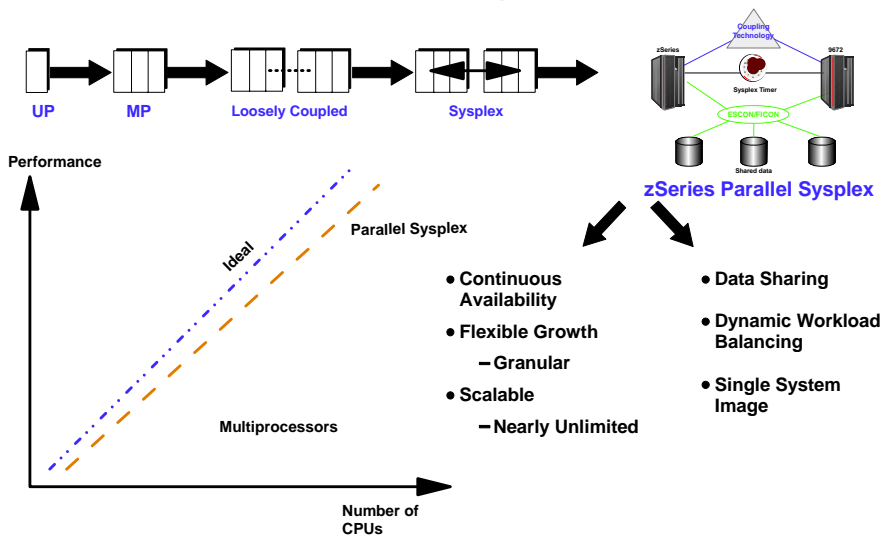- **Parallel Sysplex Software**
- **Summary**

# The zSeries Parallel Sysplex Solution

- Shared data
- Dynamic workload balancing
- Continuous application availability
- Incremental Non-disruptive growth
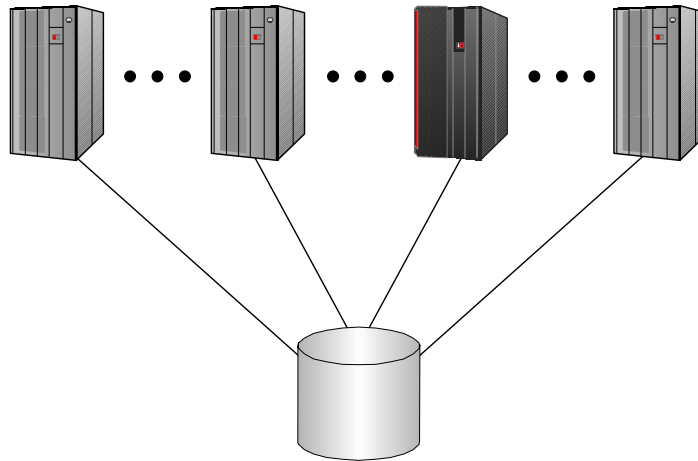
Coupling
Technology

---

# Evolution in Information Processing

**UP**   **MP**   **Loosely Coupled**   **Sysplex**

Coupling
Technology

zSeries   9672

Sysplex Timer

ESCON/FICON

Shared data

**zSeries Parallel Sysplex**

Performance

Ideal

Parallel Sysplex

Multiprocessors

Number of
CPUs

- **Continuous Availability**
- **Flexible Growth**
  - **Granular**
- **Scalable**
  - **Nearly Unlimited**

- **Data Sharing**
- **Dynamic Workload Balancing**
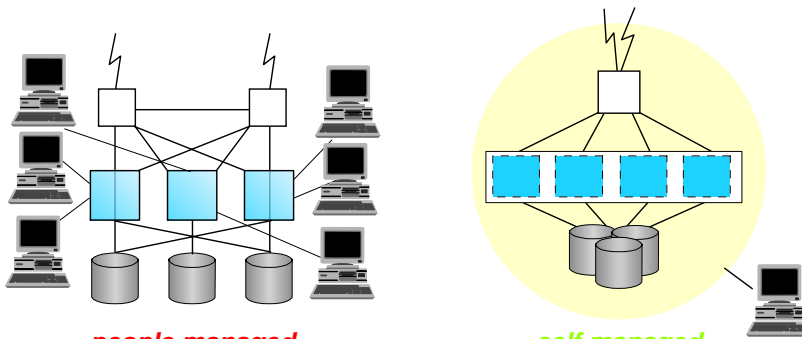- **Single System Image**

- Unique Design of IBM Hardware and Software
  - Base for Future Enhancements

---

AC42098.PRZ-7-8

# The Challenge

Single System Image/Systems Management
No Application Changes

---

# Simplified Systems Management

*people managed*

**Complex!**

*self managed*

**Simple!**

Design Objective

**A Parallel Sysplex must**
- **Look like a single image**
- **Be managed more easily than today's single image**
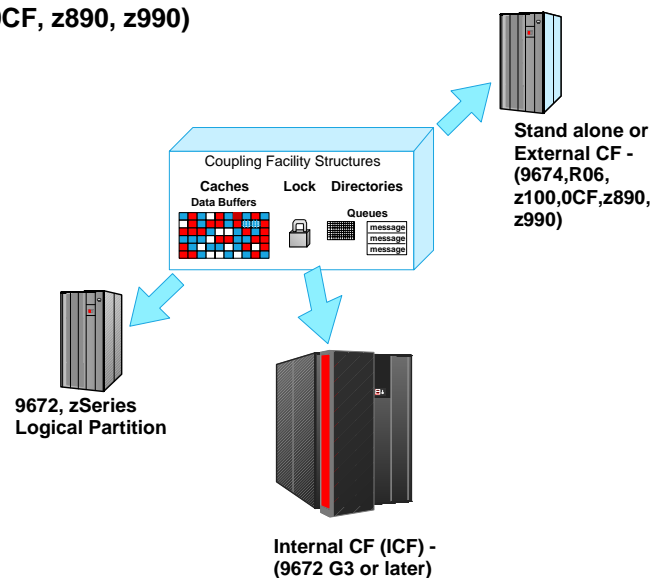
**Through**
- **Elimination of tasks**
- **Reduction of complexity**
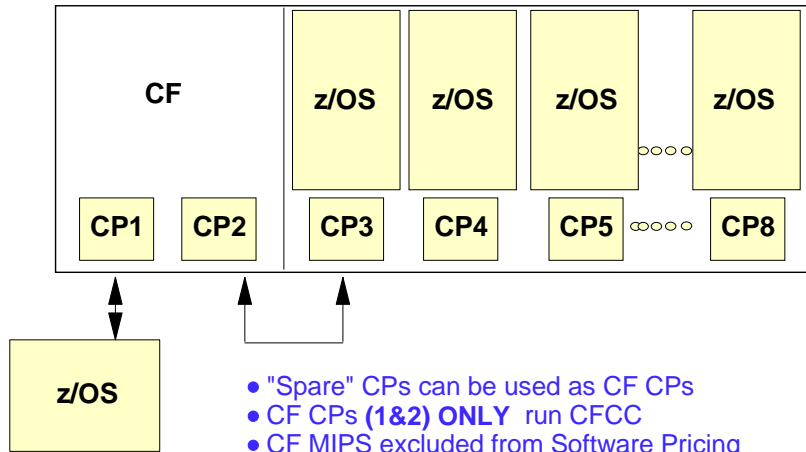- **Automation**
- **Cloning Systems**

# Agenda

- **Parallel Sysplex Overview**

- **System Structure**

- **Coupling Facility and Link Technology**

- **Parallel Sysplex Software**

- **Summary**

---

# Coupling Facility Options

- **Coupling Facility Control Code (CFCC) microcode creates a CF**
  - **Stand Alone (9674, R06, z100, 0CF, z890, z990)**
  - **Integrated (G3 +)**
  - **In an LPAR**

- **Integrated Coupling Migration Facility (ICMF)**
  - **For Test & Migration**
  - **9672 only**

Coupling Facility Structures

**Caches**
**Data Buffers**
**Lock**
**Directories**
**Queues**
message
message
message

**Stand alone or External CF - (9674,R06, z100,0CF,z890, z990)**

**9672, zSeries Logical Partition**

**Internal CF (ICF) - (9672 G3 or later)**

# Internal Coupling Facility

**CF**   **z/OS**   **z/OS**   **z/OS**   **z/OS**

**CP1**  **CP2**  **CP3**  **CP4**  **CP5**  **CP8**

**z/OS**
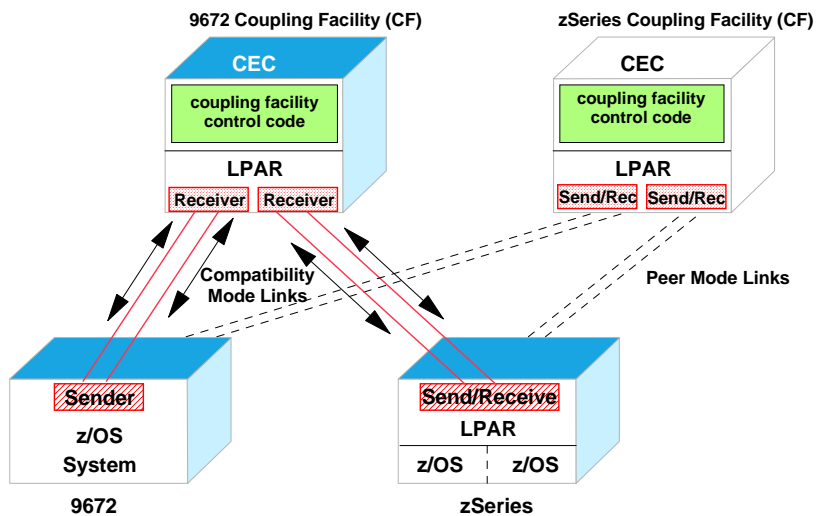
- "Spare" CPs can be used as CF CPs
- CF CPs **(1&2) ONLY** run CFCC
- CF MIPS excluded from Software Pricing
- Can run as ICMF (non-zSeries) or accessed via Coupling Links
- Available on G3 and later systems

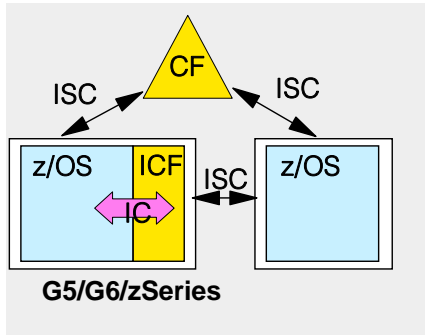**Single-System Sysplex for Software Continuous Operations**

# Coupling Links & Adapters

**9672 Coupling Facility (CF)**          **zSeries Coupling Facility (CF)**

CEC     CEC

coupling facility control code          coupling facility control code

LPAR     LPAR

Receiver   Receiver          Send/Rec   Send/Rec

**Compatibility Mode Links**          **Peer Mode Links**

Sender          Send/Receive

z/OS System          LPAR

z/OS   z/OS

**9672**          **zSeries**

Integration of zSeries Hardware & Software for Optimum Efficiency

AC42098.PRZ-13-14

# Link Technology

ISC    CF    ISC

z/OS   ICF   ISC   z/OS
IC

**G5/G6/zSeries**

ICB   CF   ICB
7 meters      7 meters

z/OS   CF   ICB   z/OS
IC   7

**G5/G6/zSeries** meters **G5/G6/zSeries**

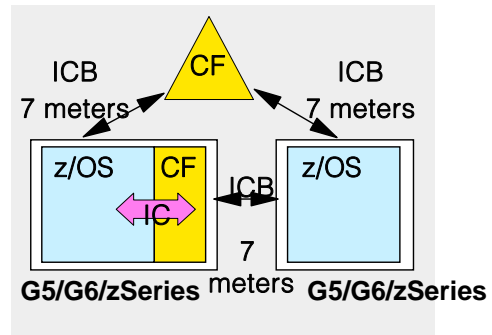*Inter System Channel (ISC)*
- ★ Original Coupling Channel
- ★ Provides long distance connections

*Internal Coupling (IC) Channel*
- ★ Fastest coupling performance
- ★ Reduced complexity
- ★ Increased reliability
- ★ Standard on 9672 G5 and up
  - ▪ **Cluster Technology Scales with Processor Speed**
  - ▪ **Peer mode available between zSeries servers/CFs**

*Integrated Cluster Bus (ICB)*
- ★ High bandwidth link for short distance
- ★ Fastest interconnection link
- ★ Improved processor utilization

---

# zSeries Coupling Technology

z990 Model 300    z900 Model 100    z800 Model 0CF

- ★ 64 Bit Architecture
- ★ Dedicated or Shared CPs
- ★ Up to 15 LPs
- ★ Up to 32 GB of storage
- ★ zSeries Peer Channels
  - ‣ InterSystem Channels-3 (ISC3)
  - ‣ Integrated Cluster Bus-3 (ICB3)
  - ‣ Integrated Cluster Bus-4 (ICB4)
- ★ zSeries Compatibility Channels
     zSeries to 9672/9674
  - ‣ InterSystem Channels (ISC)
  - ‣ Integrated Cluster Bus (ICB)
- ★ Dynamic CF Dispatch
- ★ Dynamic ICF Expansion

**z990 Model 300**
- – Up to 32 ICFs
- – Up to 16 ICB-3
- – Up to 32 ISC-3
- – Up to 16 ICB-4
- – Up to 64 GB/book
- – Upgrade to z990
- – Cannot upgrade directly from z100

**z900 Model 100**
- – Up to 9 ICFs
- – Up to 16 ICB-3
- – Up to 32 ISC-3
- – Upgrade to z900
- – Upgrade from R06

**z800 Model 0CF**
- – Up to 4 ICFs
- – Up to 6 ICB-3
- – Up to 24 ISC-3
- – Upgrade to z800

z890...
- – Up to 4 ICFs
- – Up to 16/8 ICB-3/4
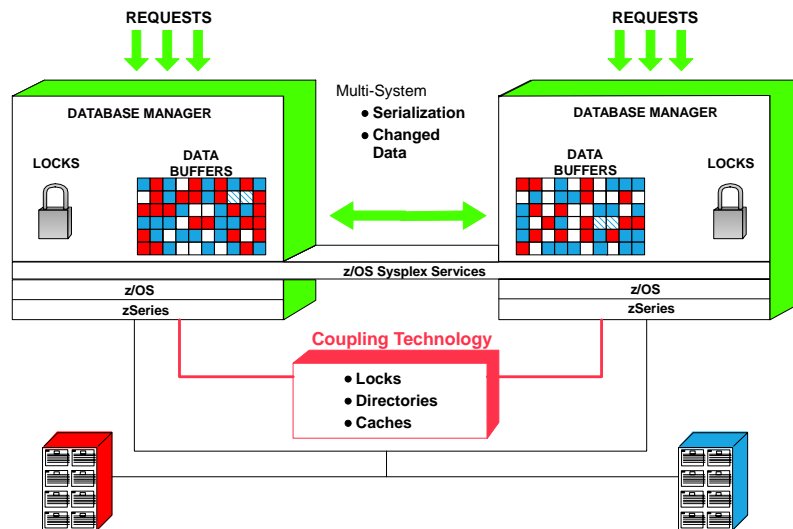- – Up to 48 ISC-3
- – Upgrade to z990

# Coupling Facility CFCC Levels

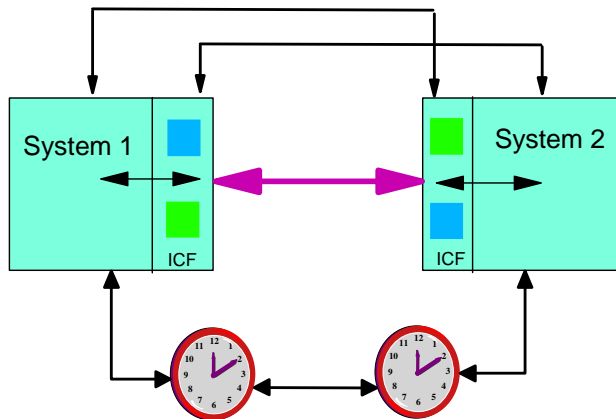| CF Level | Function | G3 | G4 | G5 | G6 | z800 | z900 | z890/990 |
|---|---|---|---|---|---|---|---|---|
| **14** | **CFCC Dispatcher Restructure** | | | | | | | x |
| **13** | **DB2 castout processing performance enhancements** | | | | | x | x | x |
| 12 | 64-bit CFCC addressability | | | | | X | X | X |
| | Message Time Ordering | | | | | X | X | X |
| | DB2 Performance | | | | | X | X | X |
| | SM Duplexing support for zSeries CFs | | | | | X | X | X |
| | Toleration for LPAR id >15 on z9xx | | | | | X | X | X |
| 11 | SM Duplexing support for 9672 G5/G6/R06 | | | X | X | | | |
| | Toleration for LPAR id >15 on z9xx | | | X | X | | | |
| 10 | z900 GA2 Level | | | | | | X | X |
| 9 | Intelligent Resource Director | | | | | X | X | X |
| | IC3 / ISC3 / ICB3 peer mode | | | | | X | X | X |
| | MQSeries Shared Queues | | | X | X | X | X | X |
| | WLM Multi-System Enclaves | | | X | X | X | X | X |
| 8 | Dynamic ICF Expansion into shared ICF pool | | | X | X | X | X | X |
| | Systems-Managed Rebuild | X | X | X | X | X | X | X |
| 7 | Shared ICF partitions on server models | | | X | X | X | X | X |
| | DB2 Delete Name optimization | X | X | X | X | X | X | X |
| 6 | ICB & IC | | | X | X | X | X | X |
| | TPF support | X | X | X | X | X | X | X |
| 5 | DB2 cache structure duplexing | X | X | X | X | X | X | X |
| | DB2 castout performance improvement | X | X | X | X | X | X | X |
| | Dynamic ICF expansion into shared CP pool | X | X | X | X | X | X | X |
| 4 | Performance optimization for IMS & VSAM RLS | X | X | X | X | X | X | X |
| | Dynamic CF Dispatching | X | X | X | X | X | X | X |
| | Internal Coupling Facility | X | X | X | X | X | X | X |
| | IMS shared message queue extensions | X | X | X | X | X | X | X |
| 3 | IMS shared message queue base | X | X | X | X | X | X | X |
| 2 | DB2 performance | X | X | X | X | X | X | X |
| | VSAM RLS | X | X | X | X | X | X | X |
| | 255 Connectors / 1023 structures for IMS Batch DL1 | X | X | X | X | X | X | X |
| 1 | Dynamic Alter support | X | X | X | X | X | X | X |
| | CICS temporary storage queues | X | X | X | X | X | X | X |
| | System logger | X | X | X | X | X | X | X |

**Details available in "Coupling Facility Level (CFLevel) Consideration" at URL   ibm.com**/servers/eserver/zseries/pso/cftable.html

TLLBPARz130

---

# Role of the Coupling Facility



AC42098.PRZ-17-18

# System Managed CF Structure Duplexing



OS:   z/OS v1.2 or later
ICFs:  zSeries G5, G6
CFs:   R06 or zSeries
CFCC: Level 11 (G5/G6) or Level 12 or higher  (zSeries)

- - Automatic Rebuild for planned reconfiguration
- - Automatic switchover for unplanned outages
- - Automatic duplexing re-establishment
- - Overlapped requests for high performance

- - Consistent Recovery Mechanism
  ► Reduced complexity
- - Faster than structure rebuild
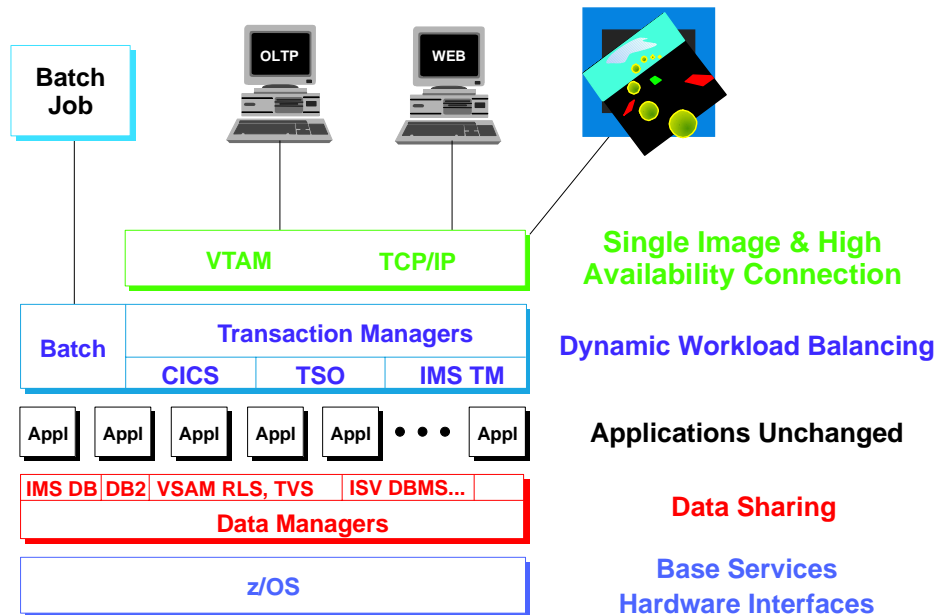- - Enables a robust "all-ICF" configuration

---

# Parallel Sysplex Hardware Cluster Technology

| Hardware Component | Function |
|---|---|
| Sysplex Timer (9037) | Consistent Multi-system Time Reference |
| Coupling Links | High Performance sysplex communications |
|  - multi-mode ISC | 50 MB/sec |
|  - Single mode ISC, ISC-3 (peer) | 100 MB/sec, 200 MB/sec |
|  - HiPer Links | 100MB/sec (w/improved adapters) |
|  - ICB, ICB-3, ICB-4 (peer) | 333 MB/sec, 1000 MB/sec, 2000 MB/sec |
|  - IC, IC-3 (peer) | 700 MB/sec, 1250 MB/sec |
| Coupling Facility (7th generation) | High performance processor |
| CFCC (14th level) | CF structures (list, lock, cache) and operations (high performance contention detection, etc. ) |
| ESCON/FICON I/O Architecture and Directors | Flexible, high availability I/O connectivity |
| I/O Fencing | Failure Isolation |
| PPRC Freeze | Data Consistency for Disaster Recovery |
| IRD | CP, I/O balancing across workloads |
| CF Structure Duplexing | High availability CF data and faster failover |

## Agenda

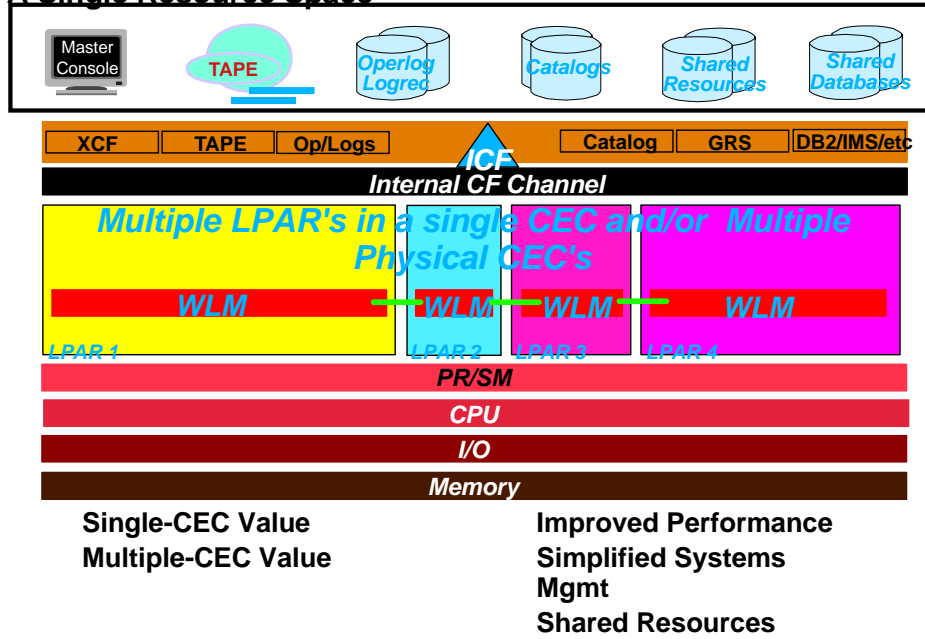- **Parallel Sysplex Overview**
- **System Structure**
- **Coupling Facility and Link Technology**
- **Parallel Sysplex Software**
- **Summary**

---

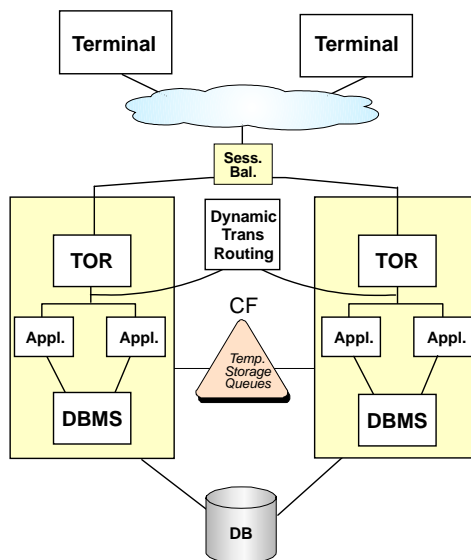## Parallel Sysplex OLTP Software Structure

| Batch Job | OLTP | WEB | |
|---|---|---|---|

**VTAM**  **TCP/IP** — Single Image & High Availability Connection

| Batch | Transaction Managers | | |
|---|---|---|---|
| | CICS | TSO | IMS TM |

Dynamic Workload Balancing

Appl · Appl · Appl · Appl · Appl · • • • · Appl — Applications Unchanged

| IMS DB | DB2 | VSAM RLS, TVS | ISV DBMS... |
|---|---|---|---|
| | Data Managers | | |

Data Sharing

**z/OS**

Base Services
Hardware Interfaces

# Parallel Sysplex Resouce Sharing

**A Single Resource Space**

| Master Console | TAPE | Operlog Logrec | Catalogs | Shared Resources | Shared Databases |

| XCF | TAPE | Op/Logs | | Catalog | GRS | DB2/IMS/etc |

ICF

**Internal CF Channel**

*Multiple LPAR's in a single CEC and/or Multiple Physical CEC's*

| WLM | WLM | WLM | WLM |

LPAR 1    LPAR 2    LPAR 3    LPAR 4

**PR/SM**

**CPU**

**I/O**

**Memory**

**Single-CEC Value**
**Multiple-CEC Value**

**Improved Performance**
**Simplified Systems Mgmt**
**Shared Resources**

---

# Transaction Management

| Terminal | Terminal |

Sess. Bal.

Dynamic Trans Routing

TOR          TOR

CF

Appl.   Appl.        Appl.   Appl.

Temp. Storage Queues

DBMS          DBMS

DB

- **Single System Image**

- **Dynamic Session Balancing**

- **Dynamic Transaction Routing**

- Applications
  - *Any transaction can run anywhere*
  - *All data used by the application can be shared*

# Websphere MQ Series



---

# Parallel Sysplex Software Cluster Technology

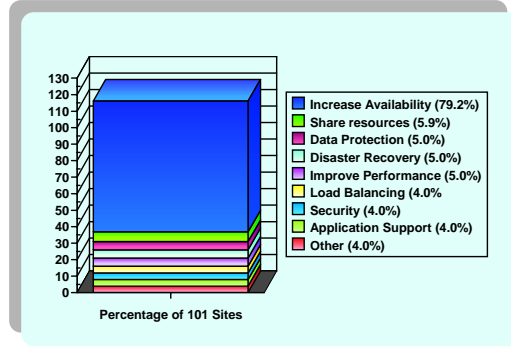| Software Component | Function |
|---|---|
| XCF | Sysplex Communication/Status Monitoring/Group Services |
| ARM | Subsystem restart (within CEC or cluster) |
| CFRM | CF Resource Management Policy |
| System Logger | High performance logging, Merged logs |
| WLM | Goal oriented unit of work management |
| WLM Enclaves | Mult-system unit of work |
| VTAM Generic Resource | Network Single System Image |
| VTAM MNPS | High Availability Network Connection |
| TCP/IP VIPA | Network Single System Image |
| TCP/IP VIPA take over/take back | High Availability Network Connection |
| CICSPlex/SM, IMS and MQ SMQ | Transaction routing/balancing |
| DB2 Sysplex Query Parallelism | SQL Query de/re-composition |
| Batch PipePlex | Cluster I/O Piping |
| ESCON Manager | ESCON I/O Systems Mangement |
| DB2, VSAM TVS, IMS/DB | Full read/write data sharing |
| IRLM | Sysplex database locking |
| Base Operating System Exploitation | Resource Sharing |
| Additional Subsystem Exploitation | Resource/Data Sharing |

## Parallel Sysplex Performance Implications

"Typical" Observed Performance (all IBM HW)

- **Multisystem Management - 3%**
- **Resource Sharing - 3%**
- **Application data sharing - <10%**
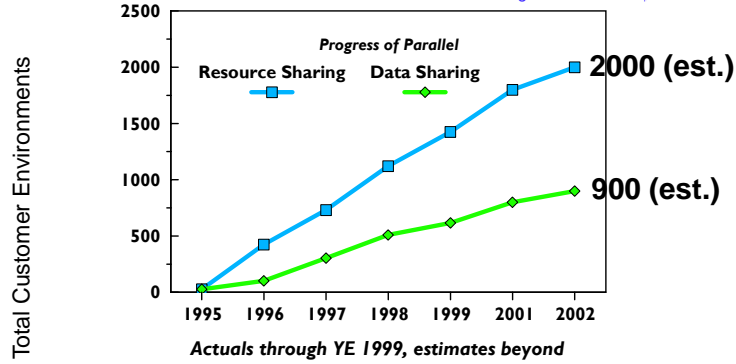- **Incremental cost of adding an image - 1/2%**

## Agenda

- **Parallel Sysplex Overview**
- **System Structure**
- **Coupling Facility and Link Technology**
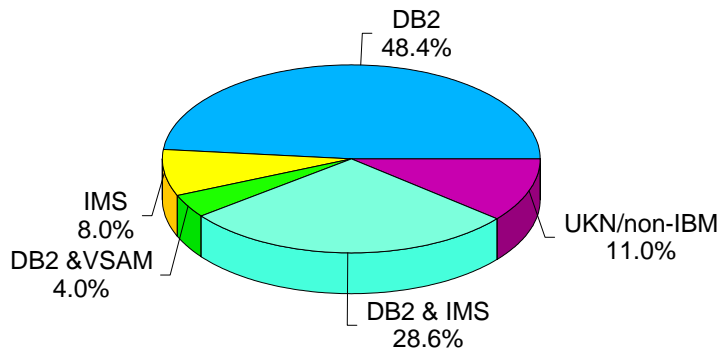- **Parallel Sysplex Software**
- **Summary**

# Parallel Sysplex Status

- Increase Availability (79.2%)
- Share resources (5.9%)
- Data Protection (5.0%)
- Disaster Recovery (5.0%)
- Improve Performance (5.0%)
- Load Balancing (4.0%
- Security (4.0%)
- Application Support (4.0%)
- Other (4.0%)

**Percentage of 101 Sites**

*Source: Strategic Research Corporation 1998 Clustering Practices Profile*

**Progress of Parallel**

Resource Sharing    Data Sharing

**2000 (est.)**

**900 (est.)**

Total Customer Environments

2500
2000
1500
1000
500
0

1995  1996  1997  1998  1999  2001  2002

*Actuals through YE 1999, estimates beyond*

---

# Data Sharing Database Summary

DB2
48.4%

UKN/non-IBM
11.0%

IMS
8.0%

DB2 &VSAM
4.0%

DB2 & IMS
28.6%

|  |  | Number | % |
|---|---|---|---|
| TOTAL SITES |  | 525 | 100 |
| Sites with DB2 |  | 425 | 81.0 |
| Sites with IMS |  | 192 | 36.6 |

# Parallel Sysplex Production by Industry

## Data Sharing

- Utilities 3.0%
- Travel/Trans. 3.0%
- Telecom 8.6%
- Process 2.3%
- Manufacturing 7.2%
- Insurance 13.0%
- Health 2.5%
- Government 7.6%
- Distribution 10.3%
- Banking/Finance 39.0%
- Other 3.4%

## Resource Sharing

- Utilities 4.0%
- Travel/Trans. 3.9%
- Telecom 7.6%
- Process 3.6%
- Manufacturing 6.8%
- Insurance 11.9%
- Health 3.4%
- Government 9.5%
- Distribution 10.4%
- Banking/Finance 34.0%
- Other 4.8%

---

# zSeries Continuous Availability

### Single System

### Parallel Sysplex

1 to 32 Systems

### GDPS

Site 1          Site 2

- **Built In Redundancy**

- **Capacity Upgrade on Demand**

- **Capacity Backup**

- **Hot Pluggable I/O**

- **Addresses Planned/Unplanned HW/SW Outages**

- **Flexible, Nondisruptive Growth**
  - ► Capacity beyond largest CEC
  - ► Scales better than SMPs
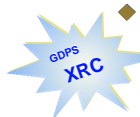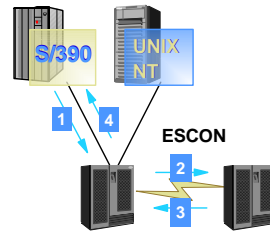
- **Dynamic Workload/Resource Management**

- **Addresses Site Failure/Maintenance**

- **Sync/Async Data Mirroring**
  - ► Eliminates Tape/Disk SPOF
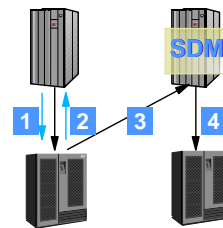  - ► No/Some Data Loss

- **Application Independent**

# GDPS/PPRC and GDPS/XRC

**GDPS PPRC**

- ◆ <u>Peer to Peer Remote Copy</u> (PPRC)
  - ◆ Synchronous data mirroring
- ◆ GDPS manages secondary data consistency
  - ◆ No or limited data loss in failover - user policy
- ◆ Production site exception condition monitoring
  - ◆ GDPS initiates and executes failover
- ◆ Distance between sites up to 40KM (fiber)
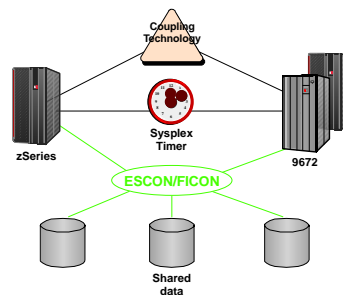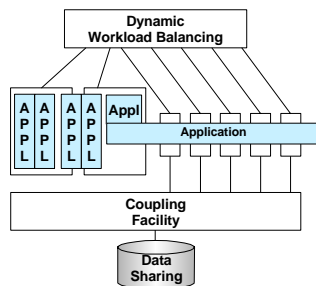- ◆ **Continuous Availability and Disaster Recovery solution**

**GDPS XRC**

- ◆ <u>eXtended Remote Copy</u> (XRC)
  - ◆ Asynchronous data mirroring
  - ◆ Limited data loss to be expected in unplanned failover
- ◆ XRC manages secondary data consistency
- ◆ GDPS executes parallel sysplex restart
  - limited user involvement
- ◆ Supports any distance
- ◆ **Disaster Recovery solution**

---

# Parallel Sysplex Value

Parallel Sysplex

**Dynamic Workload Balancing**

APPL APPL APPL APPL Appl

**Application**

**Coupling Facility**

**Data Sharing**

**Coupling Technology**

zSeries

**Sysplex Timer**

9672

**ESCON/FICON**

**Shared data**

The **best** server for diverse workloads . . .
- **Traditional**
- **New network computing**
  **continuous application availability**
  **virtually unlimited capacity**
  **leverage existing investments**
  **lowest incremental cost**
  **classic strengths!**

AC42098.PRZ-33-34

IBM Systems Group

IBM

# Additional Parallel Sysplex Information

- **www.ibm.com/servers/eserver/zseries/pso**
  - ▶ **zSeries Parallel Sysplex Cluster: What is it and what can it do for you?**
    - ─ **Business Value Overview**
  - ▶ **System-Managed CF Structure Duplexing (GM13-0103)**
  - ▶ **Configuring consoles for maximum availability**
  - ▶ **Availability Checklist**
  - ▶ **CF Configuration Options**
  - ▶ **Leveraging z/OS TCP/IP Dynamic VIPAs and Sysplex Distributor for Higher Availability**