



# 14th Workshop on the Use of High Performance Computing in Meteorology

**A cost-effective redundancy scheme for real-time data production systems through the usage of virtualization.**

**Martin Dillmann, EUMETSAT**

# Outline



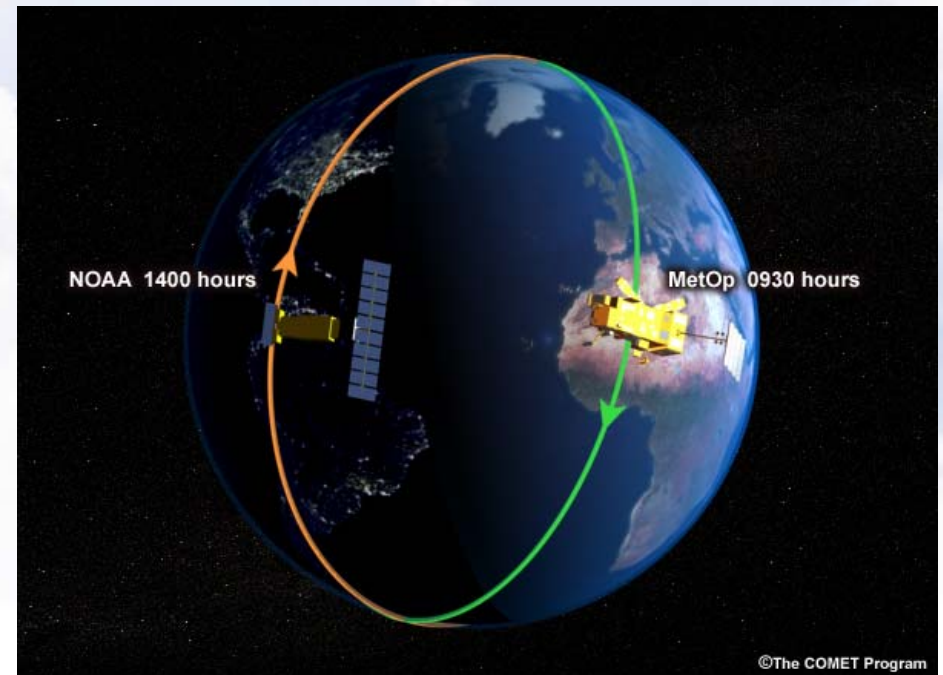
- **EPS Background**
- **Production Timing**
- **Environments and Redundancy**
- **Main Challenges**
- **Virtualisation**
- **New Redundancy Concept**
- **Project Status**
- **Outlook**



# EPS – EUMETSAT Polar System



- EUMETSAT and NOAA have agreed to provide an operational polar-orbiting service until at least 2019.



- EPS is a part of the Global Operational Satellite Observation System (GOSOS)
- It contributes to the Initial Joint Polar System (IJPS) under a cooperation agreement between EUMETSAT and NOAA



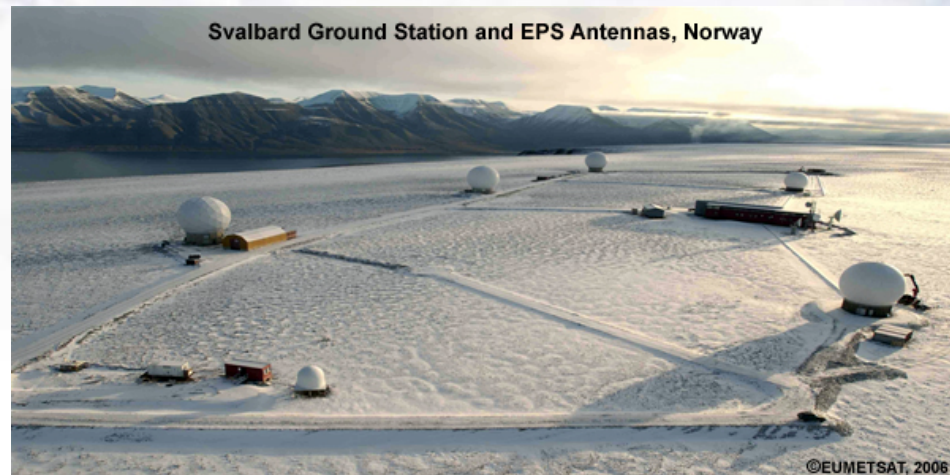


# EPS – System Components

- **3 Spacecrafts**
- **Ground Station in Svalbard (and Antarctica)**
- **M&C and data production in Darmstadt**
- **Redundant S/C control centre in Madrid**
- **R/T data distribution via EUMETCAST**



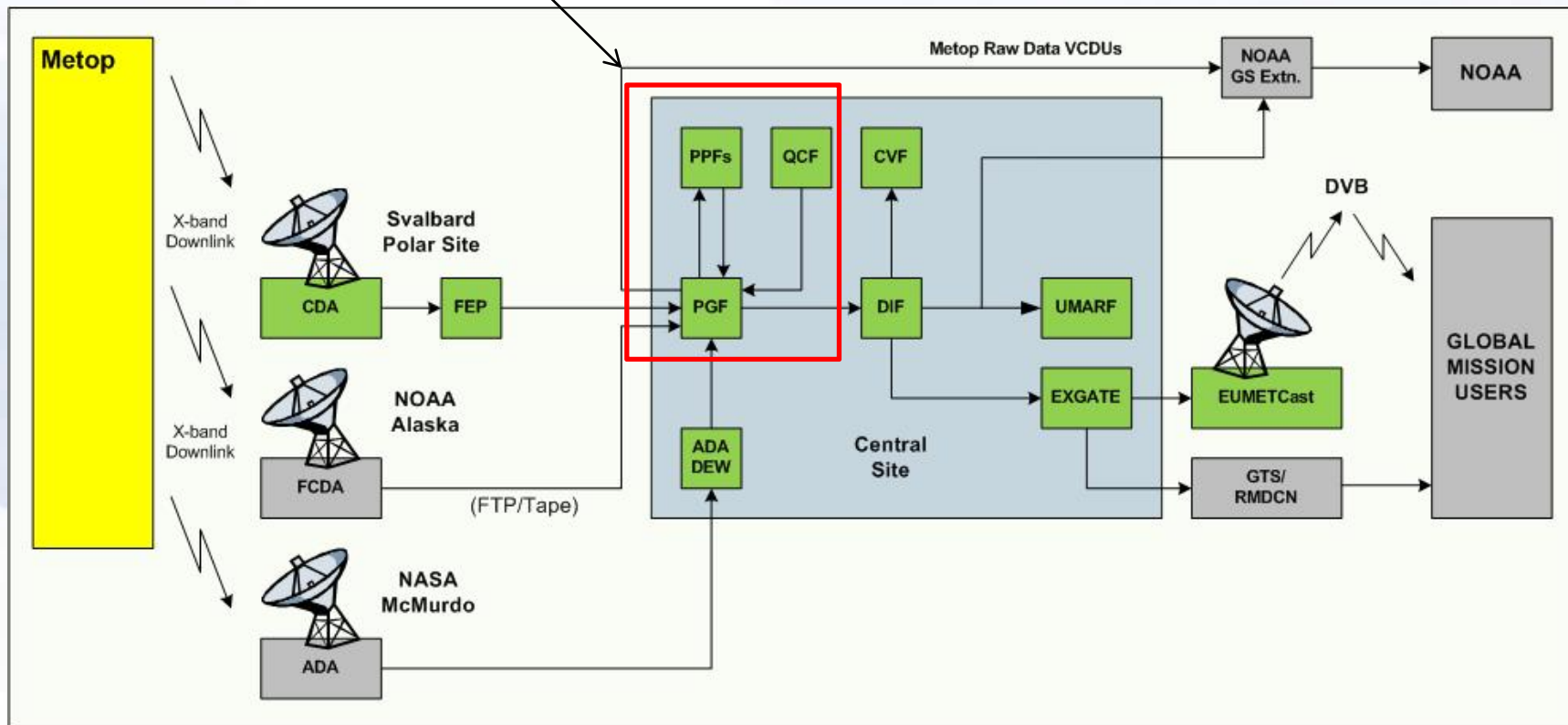
Metop-A:  
Launched: 19.10.2006  
Operational: 15.5.2007





# EPS System Overview

## Data Production





# EPS: ADA

## Antarctic Data Acquisition (ADA) :

- Metop mission data downlink to be supported by NASA/NSF 10m antenna at McMurdo Sound (77°S)
- MGS support will only be utilised for X-band, although antenna is S-band capable (may be used for auto-track)
- Mission data to be sent by combined satellite/land link to Darmstadt via Australia
- Phases : Demonstration Feb 2011 to Feb 2014 (average of 9 passes/day)  
Operational Feb 2014 onwards (all passes of operational satellite)
- Support to be focussed on prime Metop, with option to support backup Metop
- Metop mission data : raw data transfer to NOAA nearer to sensing time than present

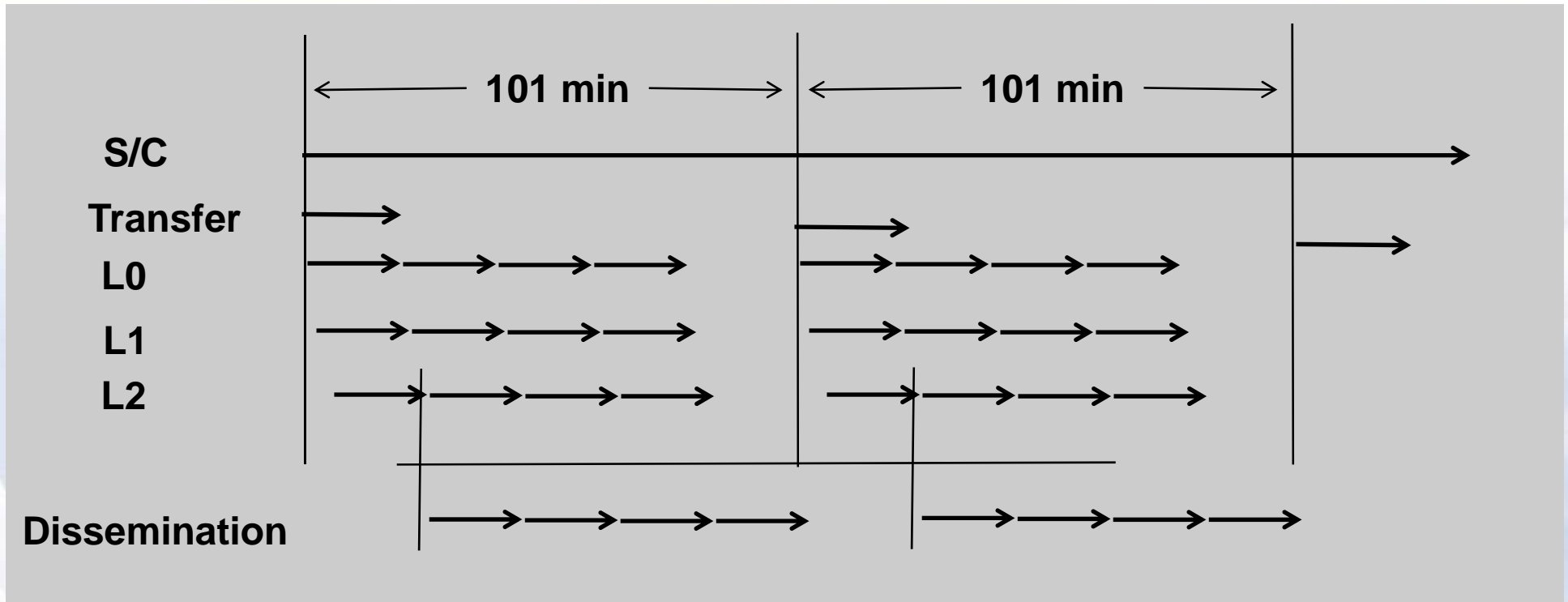


Level-1 product timeliness to improve from max. 135 to 65 mins

(sensing time to product dissemination), current avg = 115



# Production Timing - Scheme



## Requirements:

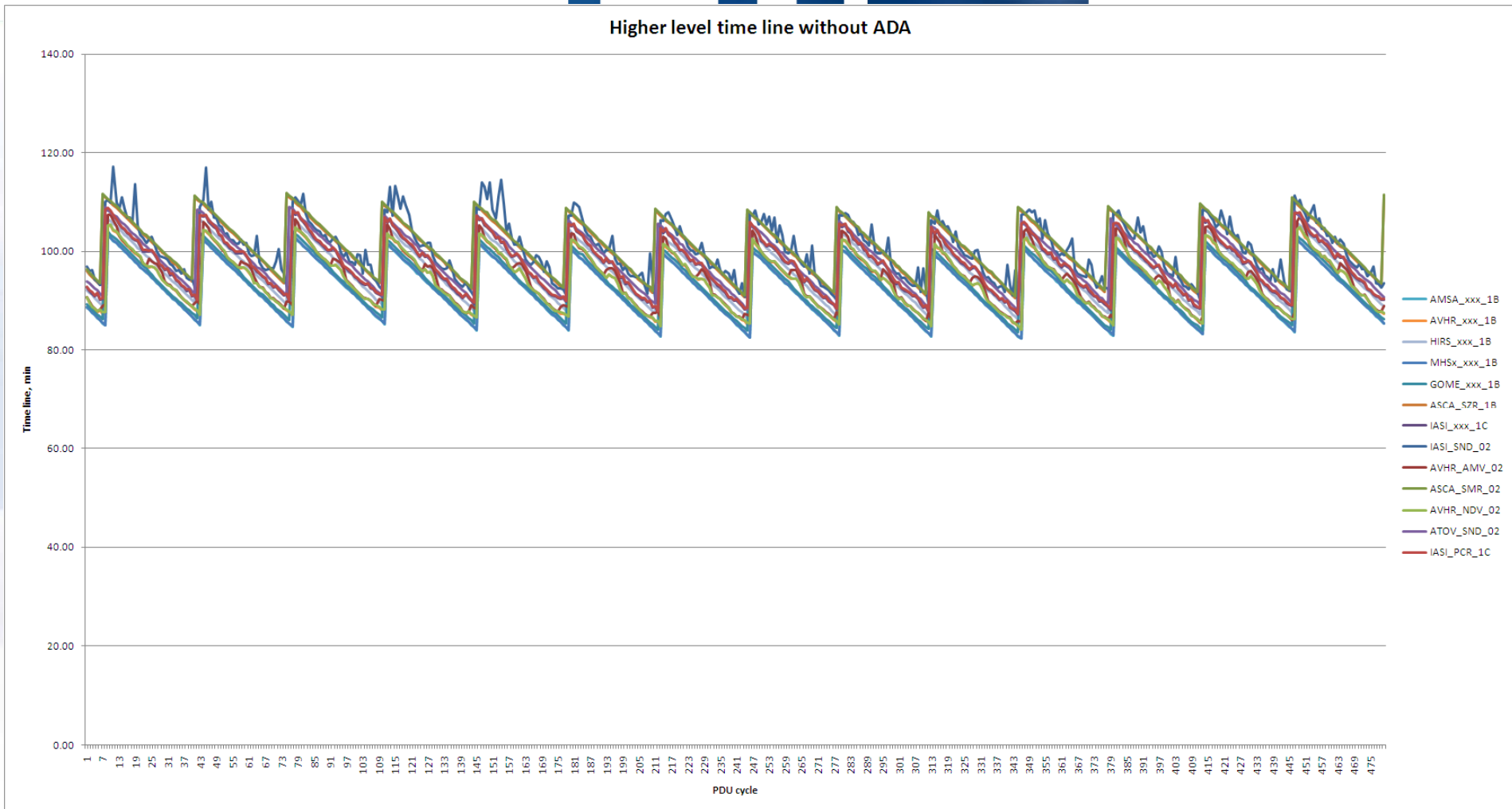
- **Global Level 1 data with an average latency of 2 hours 15 minutes**
- **Selected global Level 2 products with an average latency of 3 hours.**





# Production Timing – without ADA

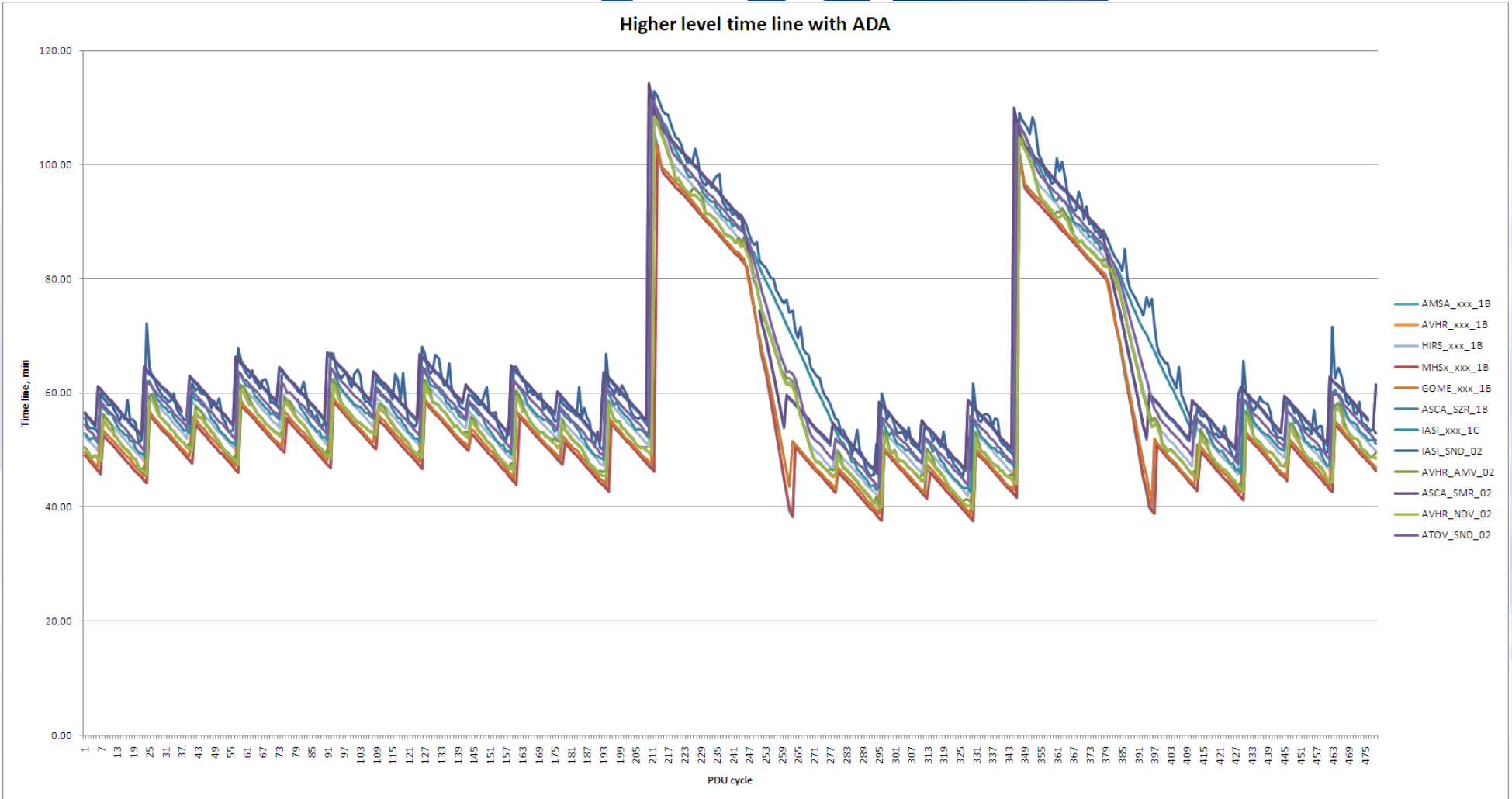
Higher level time line without ADA





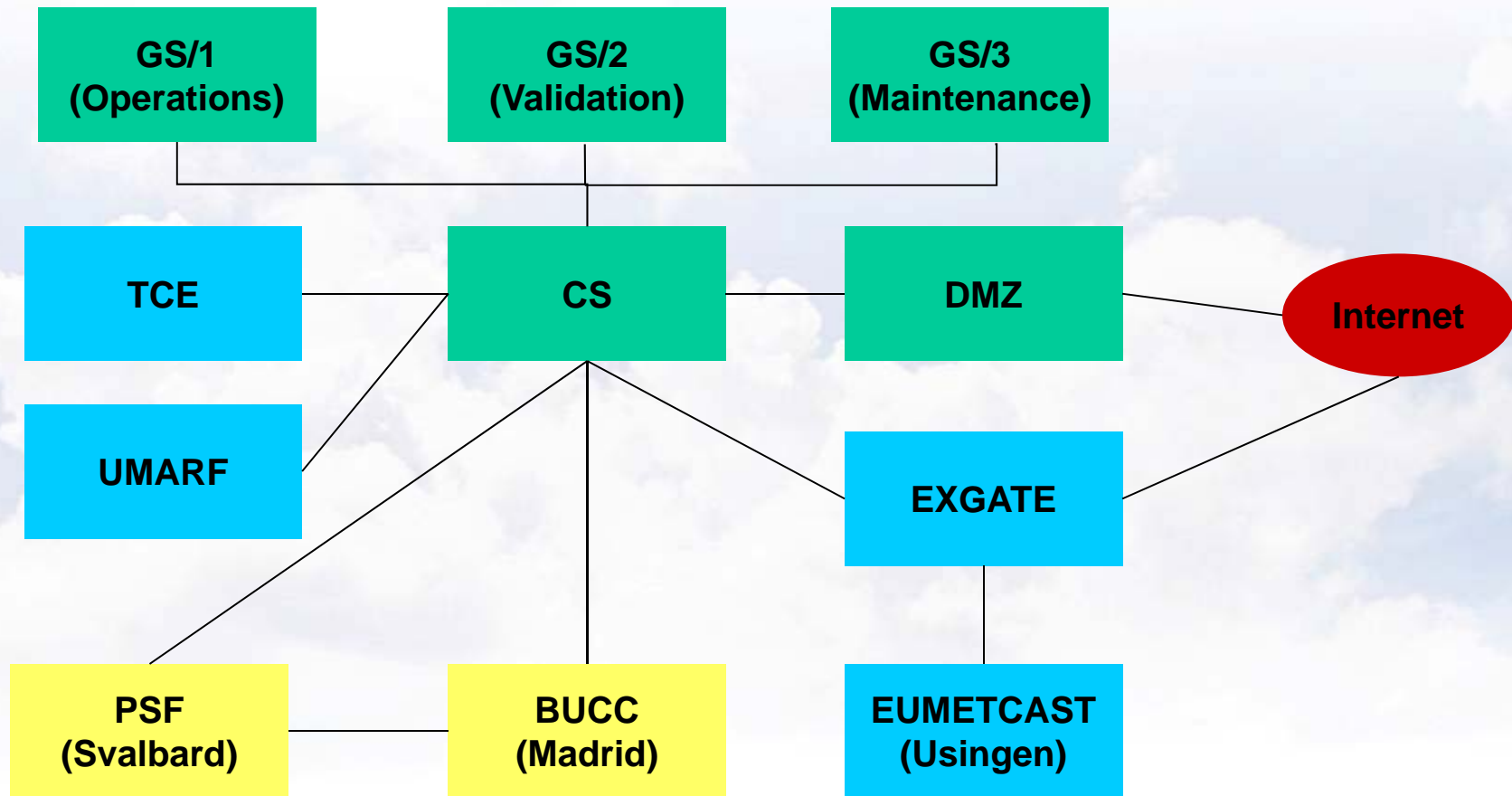


# Production Timing – with ADA (initial phase)





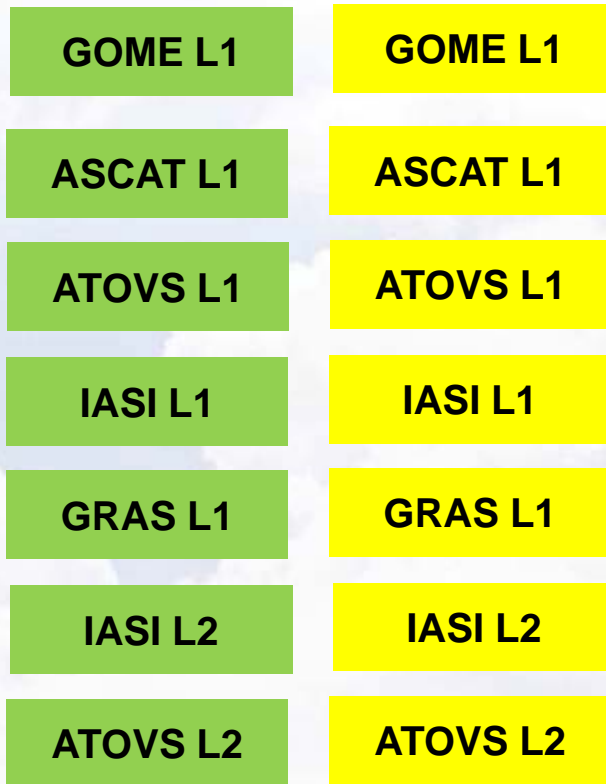
# Environments and Redundancy



# Environments and Redundancy - Classical

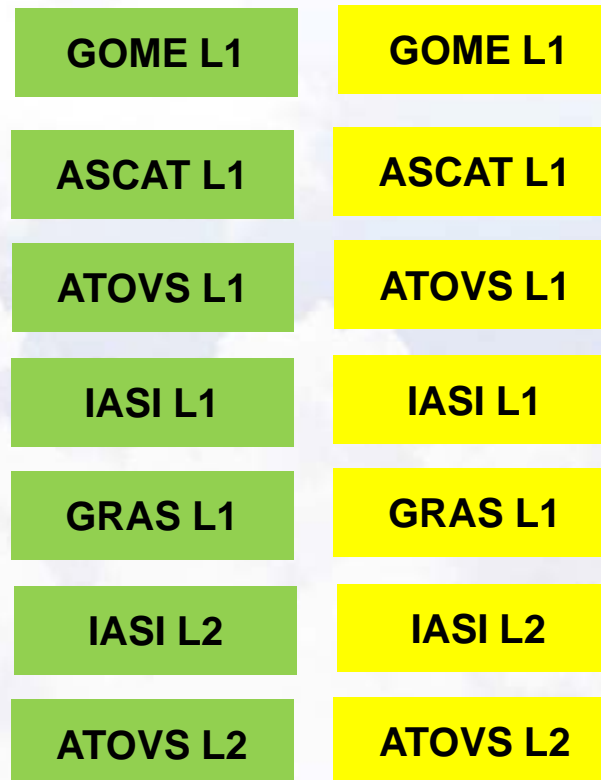
## Redundancy Approach

### GS/1



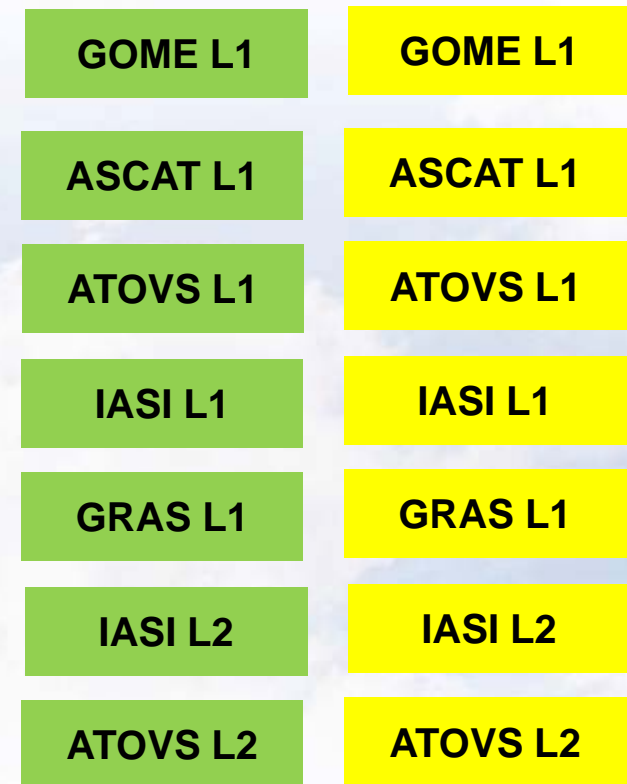
**GS/1 Network  
OPE SAN**

### GS/2



**GS/2 Network  
VALSAN**

### GS/3

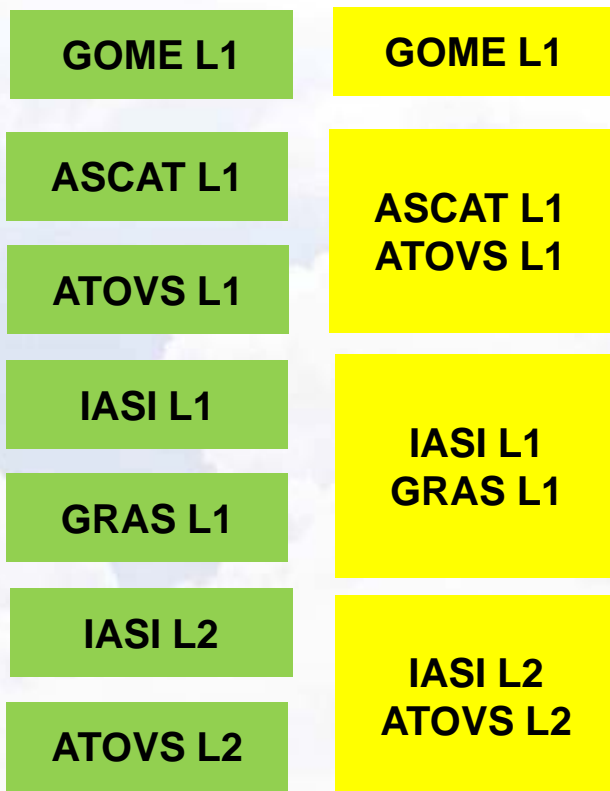


**GS/3 Network  
VAL SAN**

# Environments and Redundancy - Real Redundancy Approach

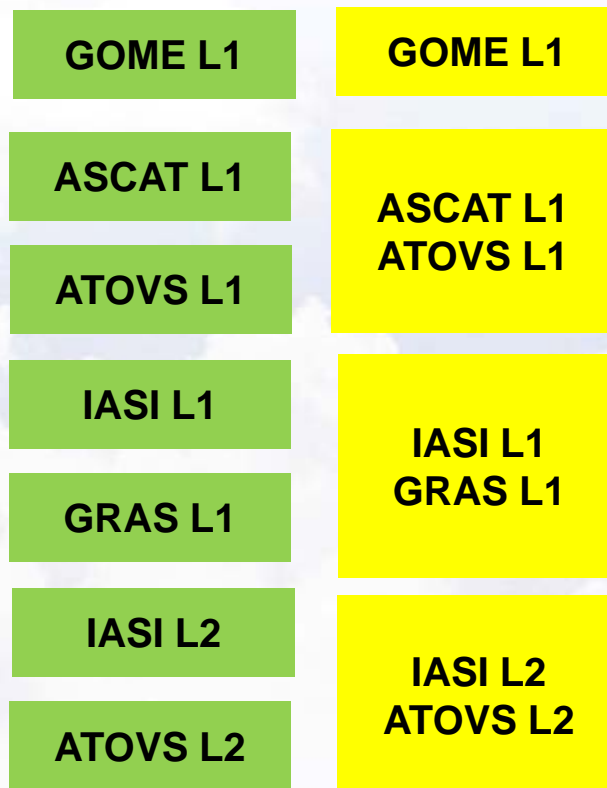


## GS/1



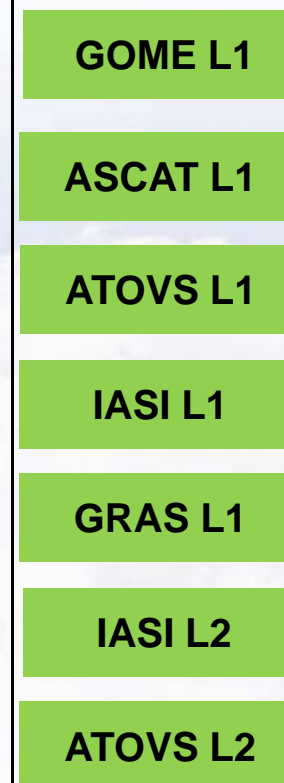
**GS/1 Network  
OPE SAN**

## GS/2



**GS/2 Network  
VAL SAN**

## GS/3



**GS/3 Network  
VAL SAN**





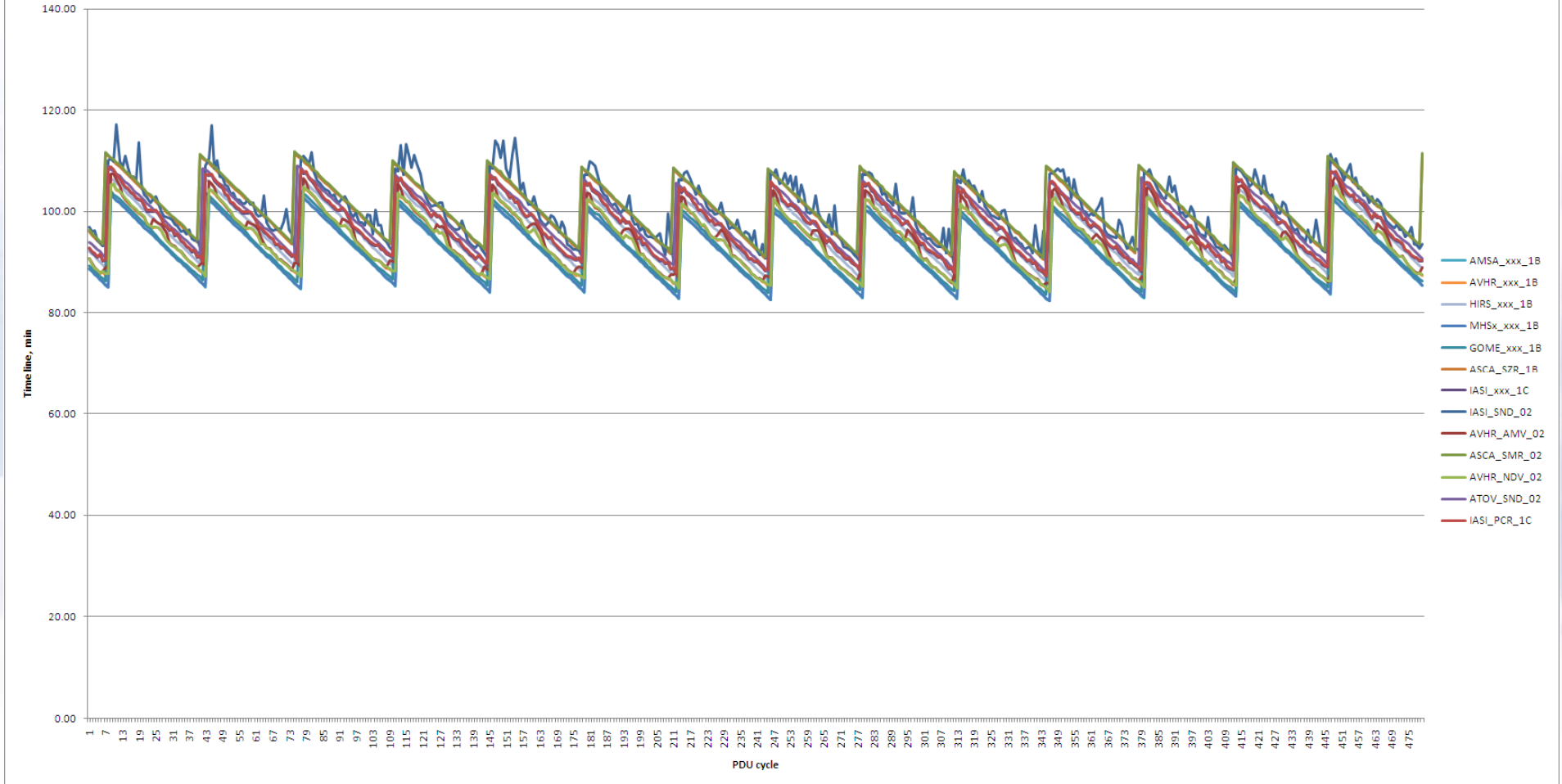
# Main Challenges

- The system has to be a **SPC (Sustained Performance Computing) System**
- Full redundancy needed (HW&SW), allowed fail-over times in the range of minutes
- Capability to verify new software versions
- Capability to scientifically validate the software
- No downtime allowed for installations
- Network Separation
- SAN Separation
- Operational independence of environments
- **BUT: Less hardware/better utilisation**



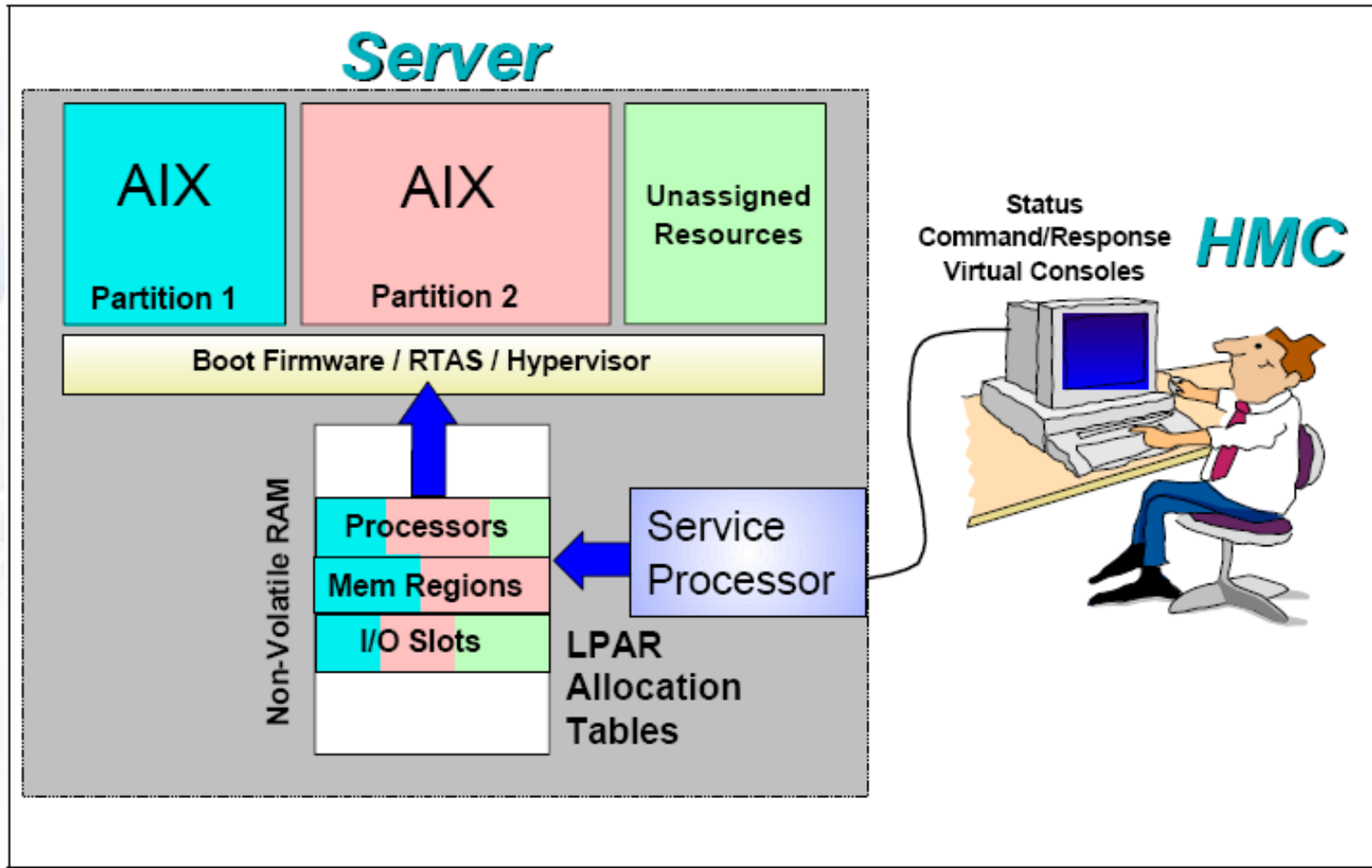
# Main Challenges – load scheme

Higher level time line without ADA





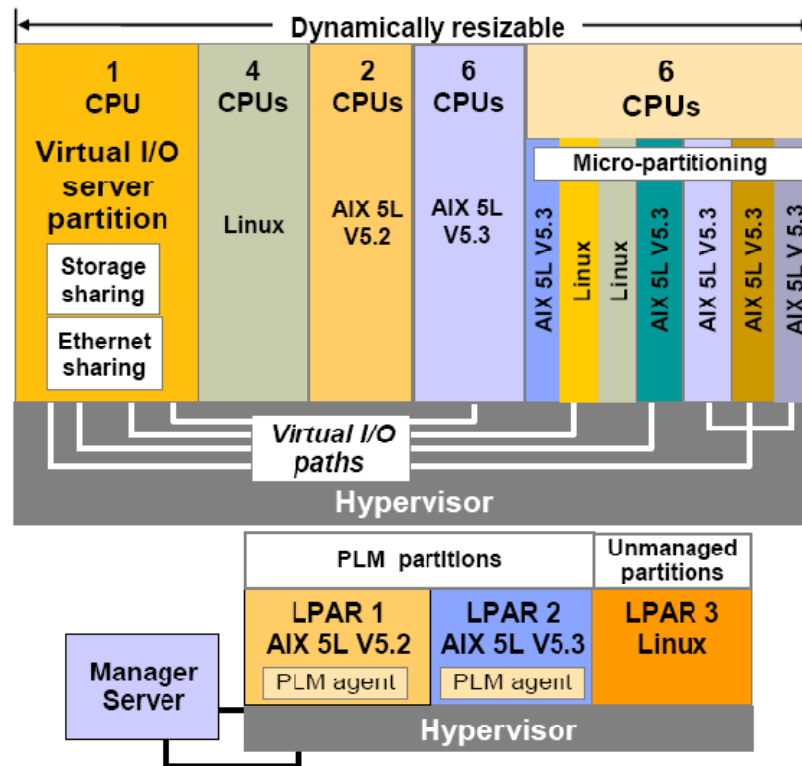
# Virtualization





# Virtualization: IBM p-series

## p5 advanced virtualization option



### Virtual I/O server

- Shared Ethernet
- Shared SCSI and Fibre Channel-attached disk subsystems
- Supports AIX 5L V5.3 and Linux\* partitions

### Micro-Partitioning

- Share processors across multiple partitions
- Minimum partition 1/10<sup>th</sup> processor
- AIX 5L V5.3 or Linux\*

### Partition Load Manager

- Both AIX 5L V5.2 and AIX 5L V5.3 supported
- Balances processor and memory request

### Managed via HMC

\* SLES 9 or RHEL AS 3





# Virtualization : IBM Hypervisor

In its current form the Hypervisor hosts the following services:

- Dynamic or static assignment of CPU time (down to 1/10 CPU) to partitions according to rules that are stored inside the Hypervisor;
- Dynamic or static assignment of memory to partitions according to rules that are stored inside the Hypervisor;
- Exclusive or shared access to I/O devices;
- Virtual SCSI devices that are served by a special partition (VIO Server, see next chapter) which connects to the real disk devices;
- Virtual networks interfaces that can be used for the transfer of data between the partitions or between partitions and real network devices being hosted by the VIO server;
- Communication with other Hypervisor for partition mobility purposes.

# Boundary Conditions for the usage of Virtualisation

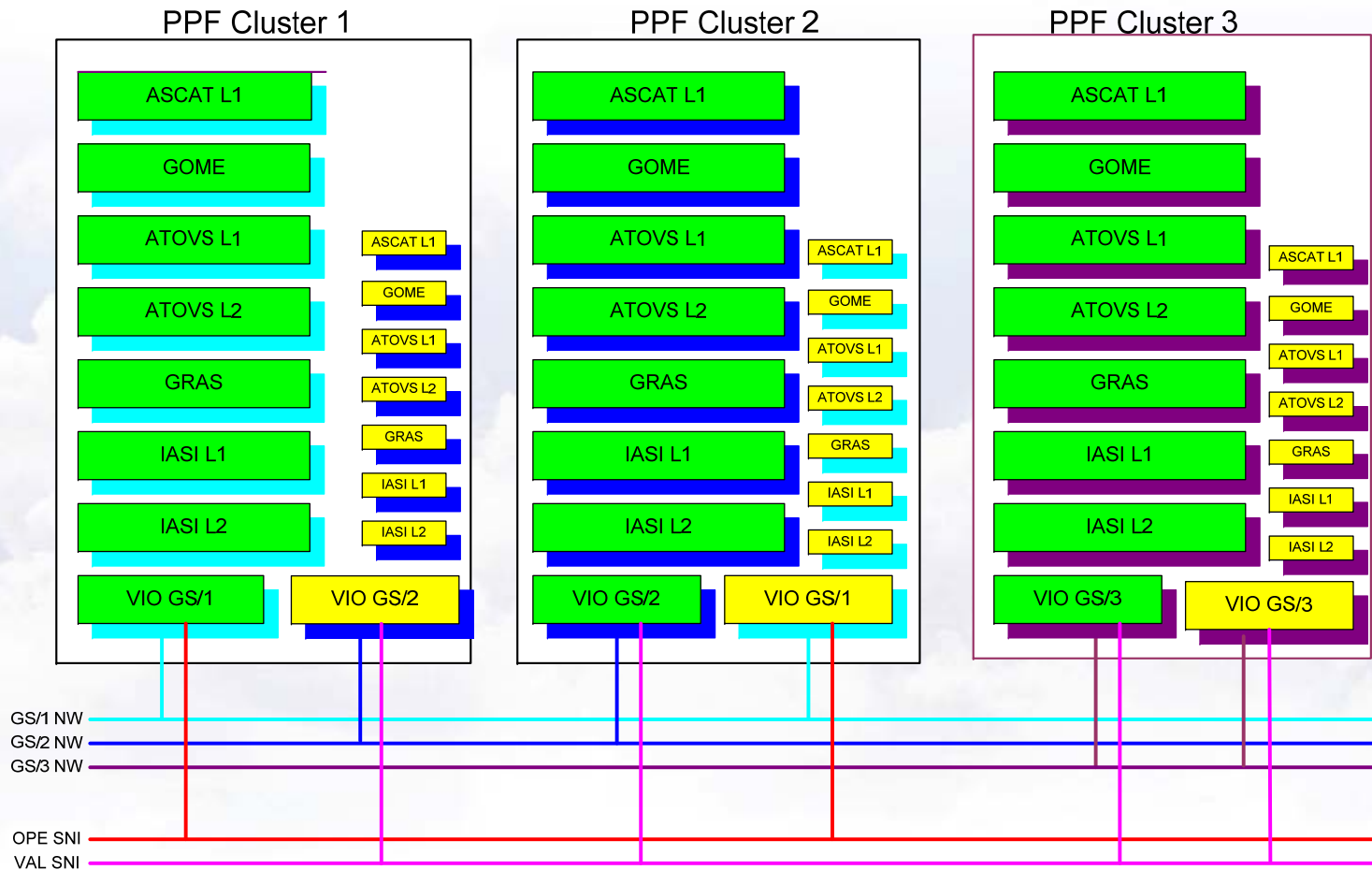
- Redundant units must not be on the same computer
- Redundant units must not be in the same infrastructure unit (rack, power, cooling)
- Operational production must not be affected under any failure scenario
- New system must be deployed in parallel to existing one
- Validation and test environments must be representative
- **BUT: Less hardware/better utilisation**

# New redundancy concept: Major failure cases

## Most likely to least likely:

- Data failures
- Software failures
- Operational errors/human errors
- Unforeseen side effects of changes
- Hardware failures
- Infrastructure failures

# New Redundancy Concept: Shared Virtualization







# New Redundancy Concept : Proposed H/W



IBM power-6 p575 HPC node:

**Processor cores:** 32 4.7 GHz POWER6 processor cores per node

**Cache:** 4 MB L2 cache per processor core  
32 MB L3 cache shared per two cores

**RAM (memory):** Up to 256 GB per node

**Internal disk:** Two SAS small form factor disks per node (73.4 GB or 146.8 GB 10K rpm)

**I/O:** Eight 1Gb Ethernet, four 4 Gb FC, two additional disks in I/O drawer (one shared per two nodes)

**Rack:** Special Water-cooled Rack, up to eight machines per rack

# New Redundancy Concept : Shared Virtualization

## Advantages:

- All the computational resources can be used in the nominal case;
- No degradation of the operational production in all failure cases;
- Because of the size of the machines (32 CPUs) and the fact that all PPFs run on one machine the maximum co-usage of resources is possible;
- Simple PPF configuration as all the PPF nodes would exist twice, no co-sharing of redundant machines needed;

## Disadvantages:

- Complex initial set-up which only would have to be done once and the complex tuning of the virtual machines.



# Deployment and Testing

- **New System installed in parallel to existing one – some power and cooling restrictions;**
- **Full parallel network integration;**
- **Storage integration of new system with new disk arrays, old system stays on the existing one, will be used as mirror after old system will have been decommissioned;**
- **System and integration level testing in parallel;**
- **Facility-level deployment and testing horizontally, each facility is installed G3-G2-G1, SW before HW or vice versa.**



# Project status

- **Full System is life since middle of the year**
- **Switch-over cost was 3 L0 PDUs plus affected L1 and L2 PDUs**
- **Full rack failover tests were working on infrastructure level but revealed software problems**



# Potential Extensions

- **More computing power for meteorological products (day-2)**
- **Usage of shared file system (IBM GPFS) to realise a data driven design**
- **Scales well as long as a single task fits in half a machine**





**Thank you**