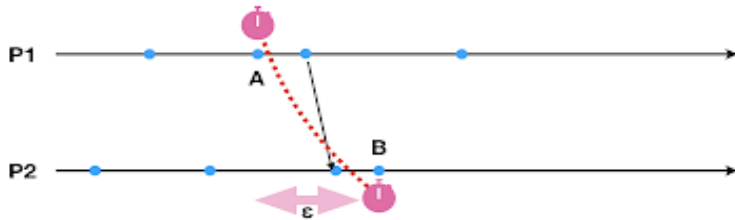# Distributing Time Synchronization in the Datacenter

Francisco Girela – Safran Navigation and Timing

February 2024

# Benefits of Precise Time

## Coherency

- Ensure that the data are the same on distributed devices

- Reduce the number of data replicas



## Efficiency

- Pre-schedule tasks to handle known low latencies

- Pipelined assignments to improve efficiency

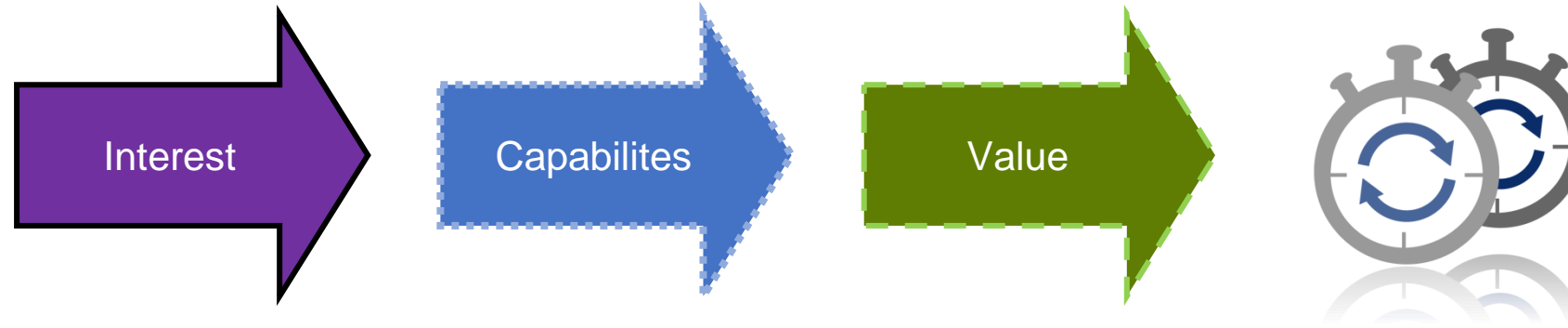- Reduce overload to ensure coherency (ε uncertainty bound)
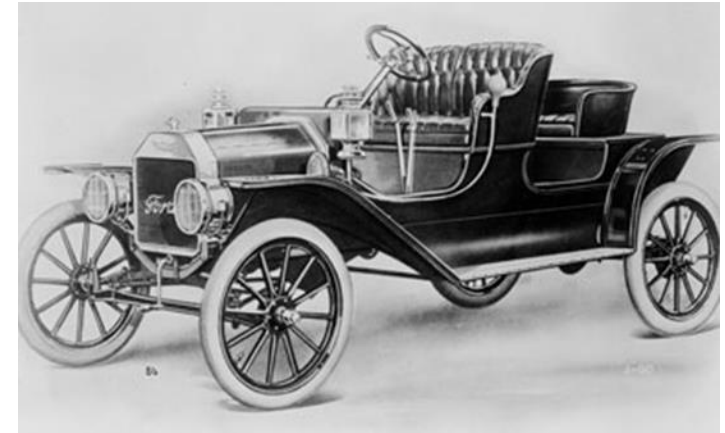
↓

Reduce CPU cycles and energy costs

## Visibility

- Have a clear view of the real order of events

- Measure latency to control bottlenecks

- Carefully allocated resources to avoid any problems

NANOG™

# Precise Time in Datacenters



| Interest | → | Capabilites | → | Value | → |
|----------|---|-------------|---|-------|---|

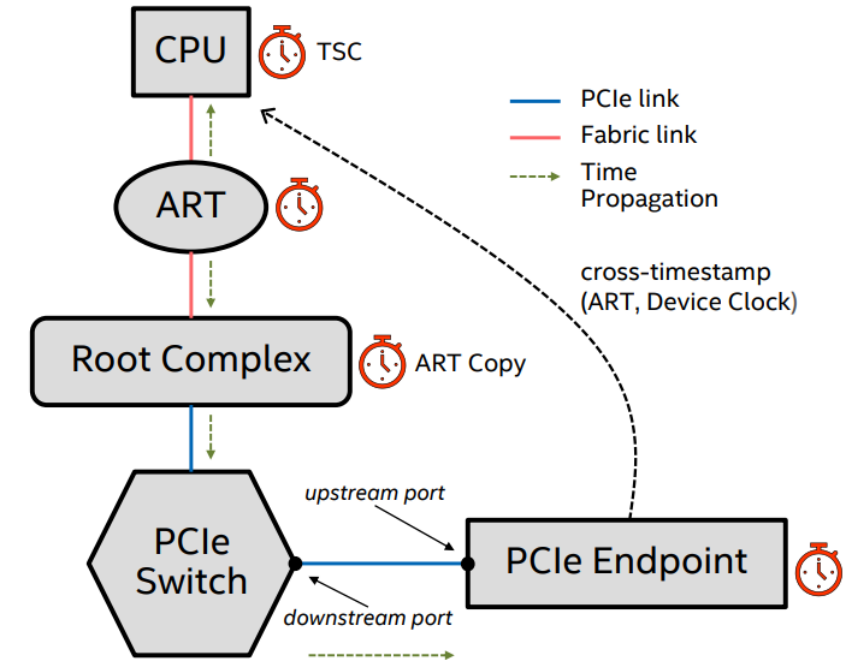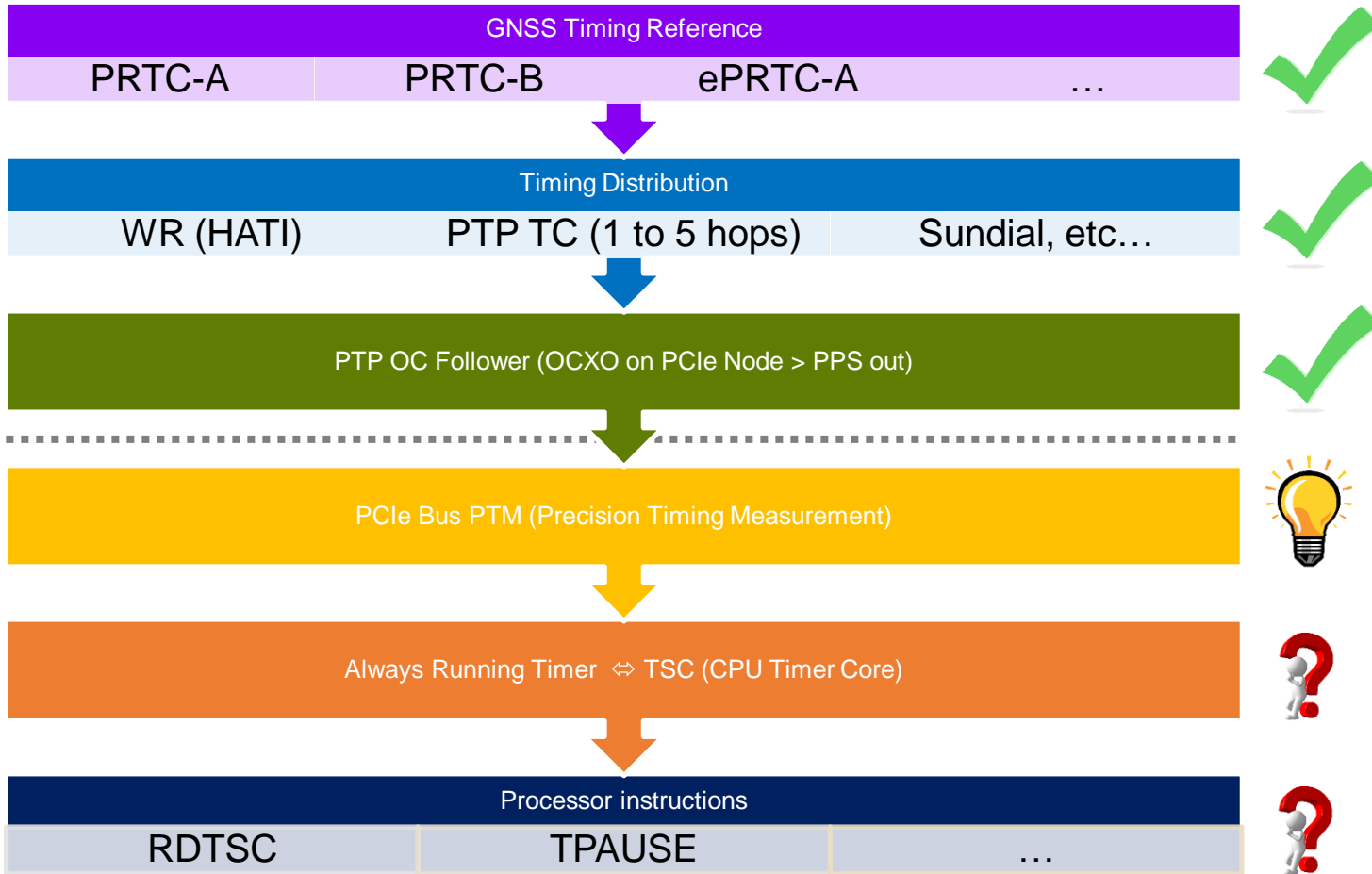Latency is one of the fundamental values

Using time requires software modification
> new layers must be written



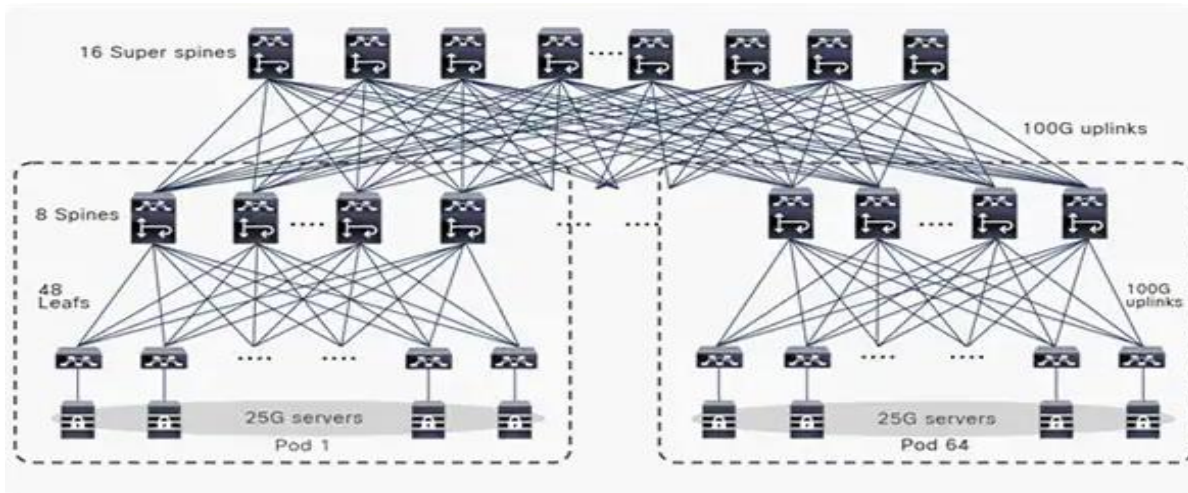*"If I had asked people **what they wanted**, they would have said **faster horses**." - Henry Ford (?)*

NANOG™

# Timing To Applications



| GNSS Timing Reference | | | | ✓ |
| PRTC-A | PRTC-B | ePRTC-A | … | |

| Timing Distribution | | | ✓ |
| WR (HATI) | PTP TC (1 to 5 hops) | Sundial, etc… | |

| PTP OC Follower (OCXO on PCIe Node > PPS out) | ✓ |

| PCIe Bus PTM (Precision Timing Measurement) | 💡 |

| Always Running Timer ⇔ TSC (CPU Timer Core) | ❓ |

| Processor instructions | | | ❓ |
| RDTSC | TPAUSE | … | |



PCIe PTM: Timing in the Last Inch, Intel, OCP-TAP
Precise Time Application, Intel, OCP-TAP

NANOG™

# Reference Architechtures

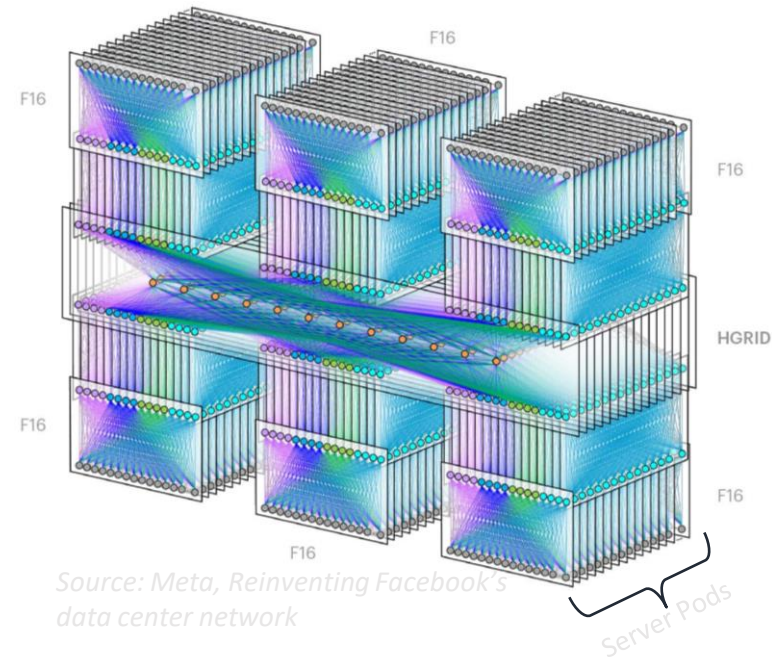1. Cisco 3-levels leaf/spine
2. Meta DC-Fabric (F16)



Source: Cisco, Massively scalable data center network fabric



Source: Meta, Reinventing Facebook's data center network

- 1 building
- 16 superspines
- 64 pods → 48 x racks/pod
- ~140K server/DC

- 1 Region → 6 buildings (F16)
- 16 fabric planes
- 48 pods → 48 x racks/pod
- ~100K servers/DC → ~600K servers/region

# Time Synchronization technologies

## White Rabbit

- Sub nanosecond accuracy and precision
- Inter- and intra- datacenter sync
- In-built failover
- Extremely scalable
- Pre-calibrated

- Dedicated infrastructure required

## PTP

- Tens of nanosecond accuracy
- Can share existing network
- Standard and widely accepted at industrial level

- Susceptible to accuracy variations during high traffic patterns
- Many different implementations / tuning parameters
- Dedicated HW required

## GNSS / PPS

- Highly available
- Tens of nanoseconds accuracy

- Limited distribution capabilities
- Susceptible to outages
- Custom cabling and infrastructure
- Not possible in all locations

## NTP

- Globally available and free reference services over the Internet
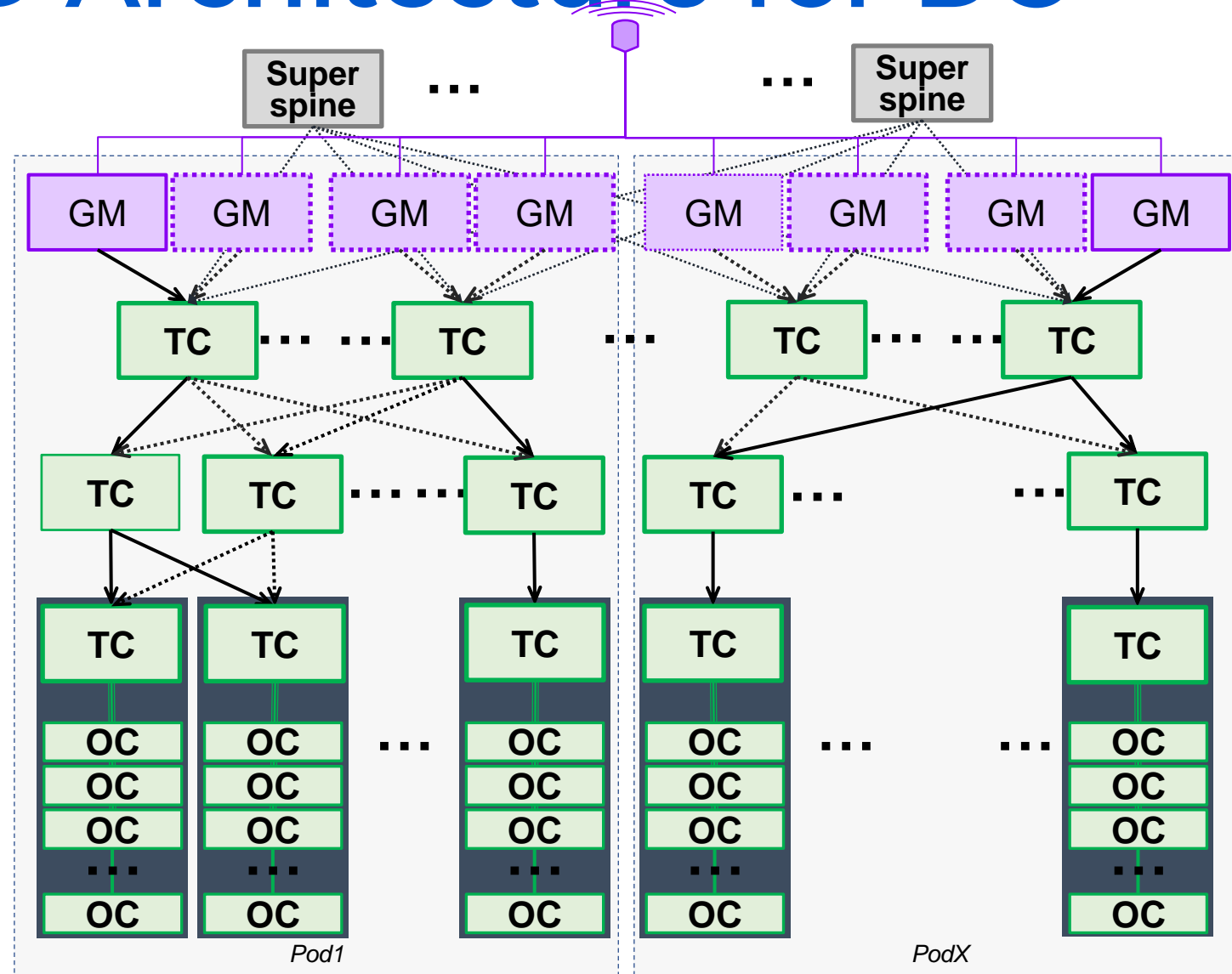
- Low levels of accuracy and precision (microseconds)

NANOG™

# Reference POD Architecture for DC Profile

*4X Open Time Servers with ART+NIC cards (GM) per pods*



- ⊕ Simple solution to put in place
  - ➤ Reduce the number of hops
  - ➤ GM handles between 5-15K clients
- ⊖ Many GNSS receivers to handle
  - ➤ Complex RF installation: Splitter, Amplifiers
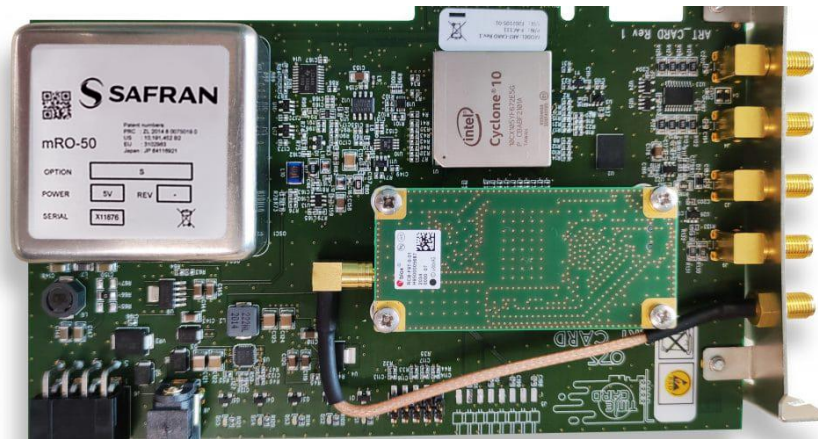  - ➤ Many references can diverge ±100ns + calibration issues

# Atomic Reference Time (ART) card

**Atomic Reference Time (ART) card**

Developed in the framework of the OCP Time Appliances Project (TAP), the ART card will provide time reference (from GNSS) to the PTP Grandmaster NIC card.



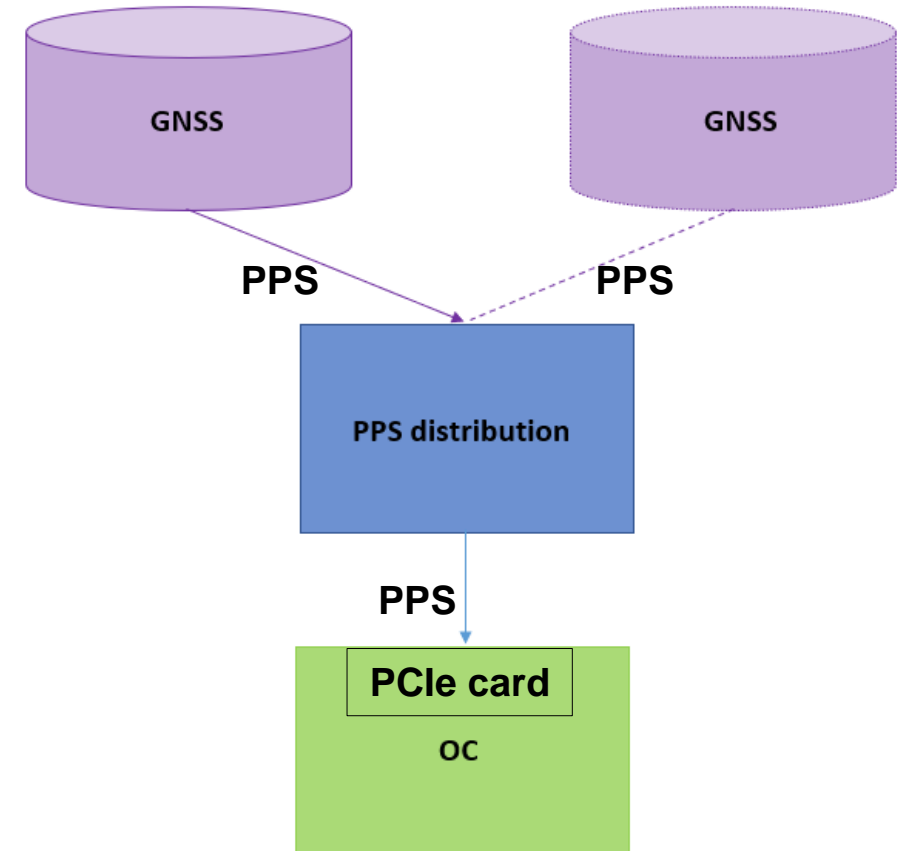**Providing time reference to an Open Time Server**

- First PCIe card **including an atomic reference** (mRO-50)

- **Software included** to monitor the synchronization of the mRO-50 on the GNSS

- **Detection** of the GNSS signal quality **to switch to holdover mode**

# PPS to the Servers through PCIe

Using a PPS distribution network and PCIe timing cards to synchronize the servers.
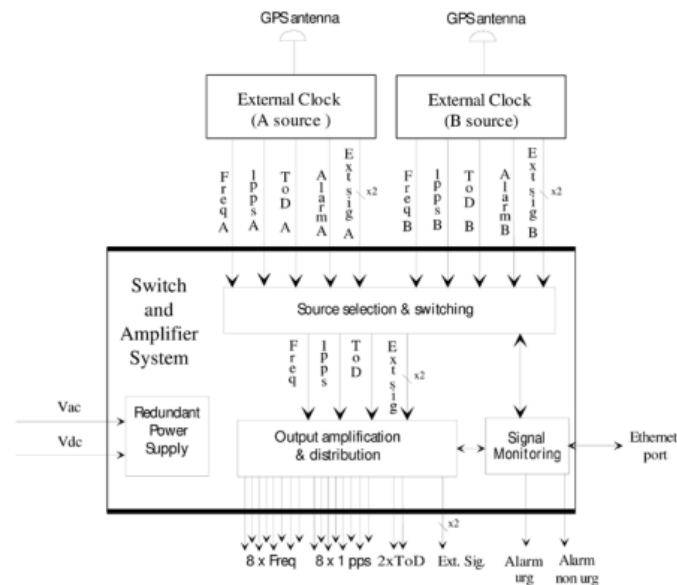
**+ Very simple and cost-effective solution**

**+ Only 1 or 2 GNSS receivers to install**

# PPS to the Servers through PCIe

**PPS distribution**

PPS distributors provide a cost-effective way to extend the distribution of time and frequency signals (pulse, low phase noise frequency signal or time of day), as a signal amplifier.



**End nodes**

Time code processors are complete synchronization systems on circuit cards ready for easy integration into servers.
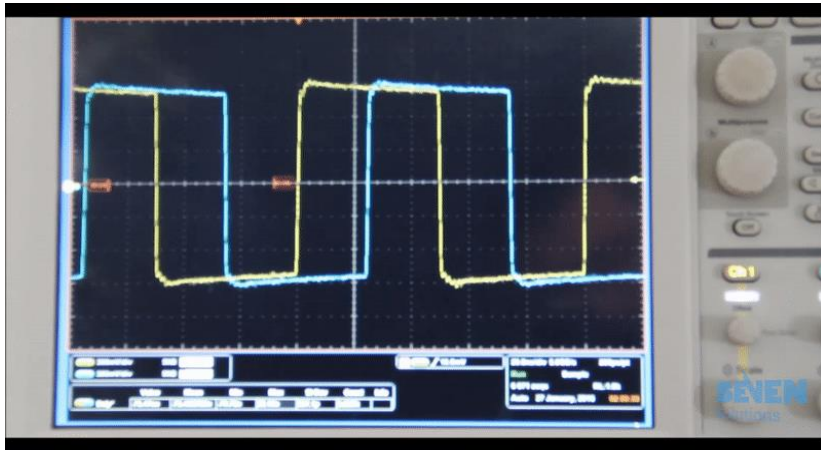
- Onboard clock/oscillator that can phase-lock to a wide **variety of external timing references** (GNSS, PPS, IRIG-B)

- The user can **prioritize multiple references** so if one is lost the unit will automatically switch to the next

- For applications where accuracy in this "holdover" conditional is essential, an **upgrade to a higher precision ovenized crystal oscillator (OCXO)** is available
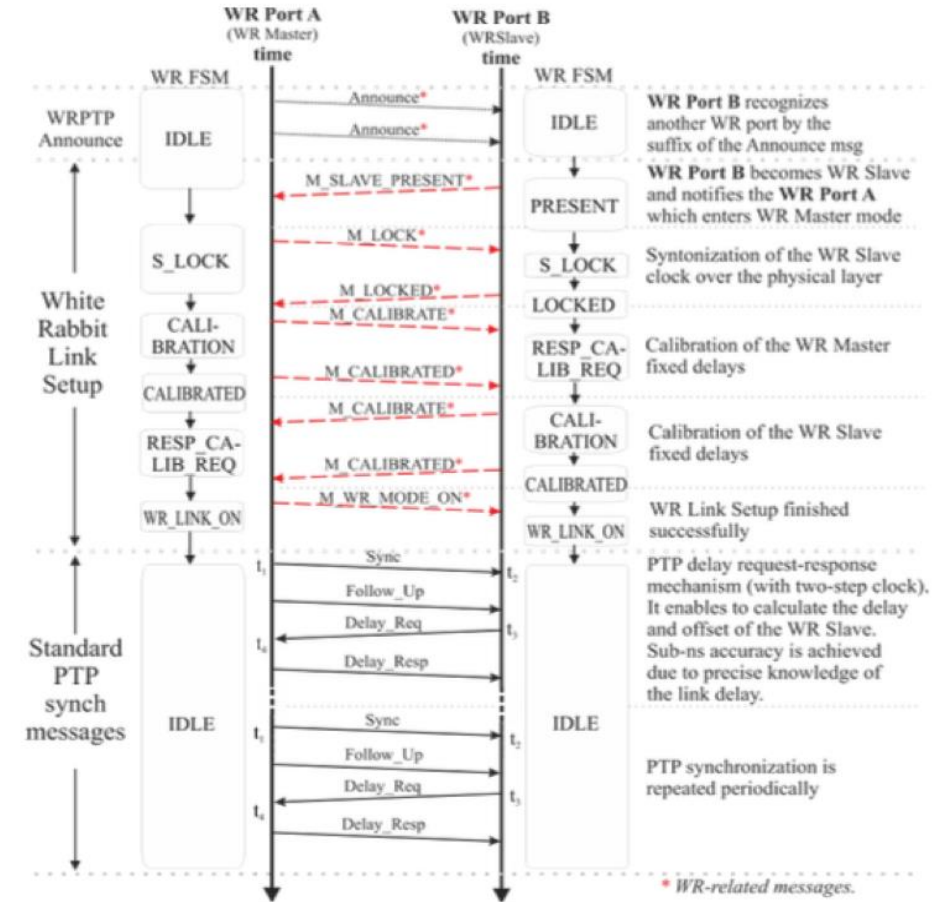
# Time transfer: WR-PTP

White Rabbit (WR) is an IEEE 1588 (PTP) implementation that achieves sub-nanosecond accuracy.

Basis for the High Accuracy profile in IEEE 1588-2019.



White Rabbit uses the information collected by the exchange of timestamped packets for correcting the constant offset between nodes
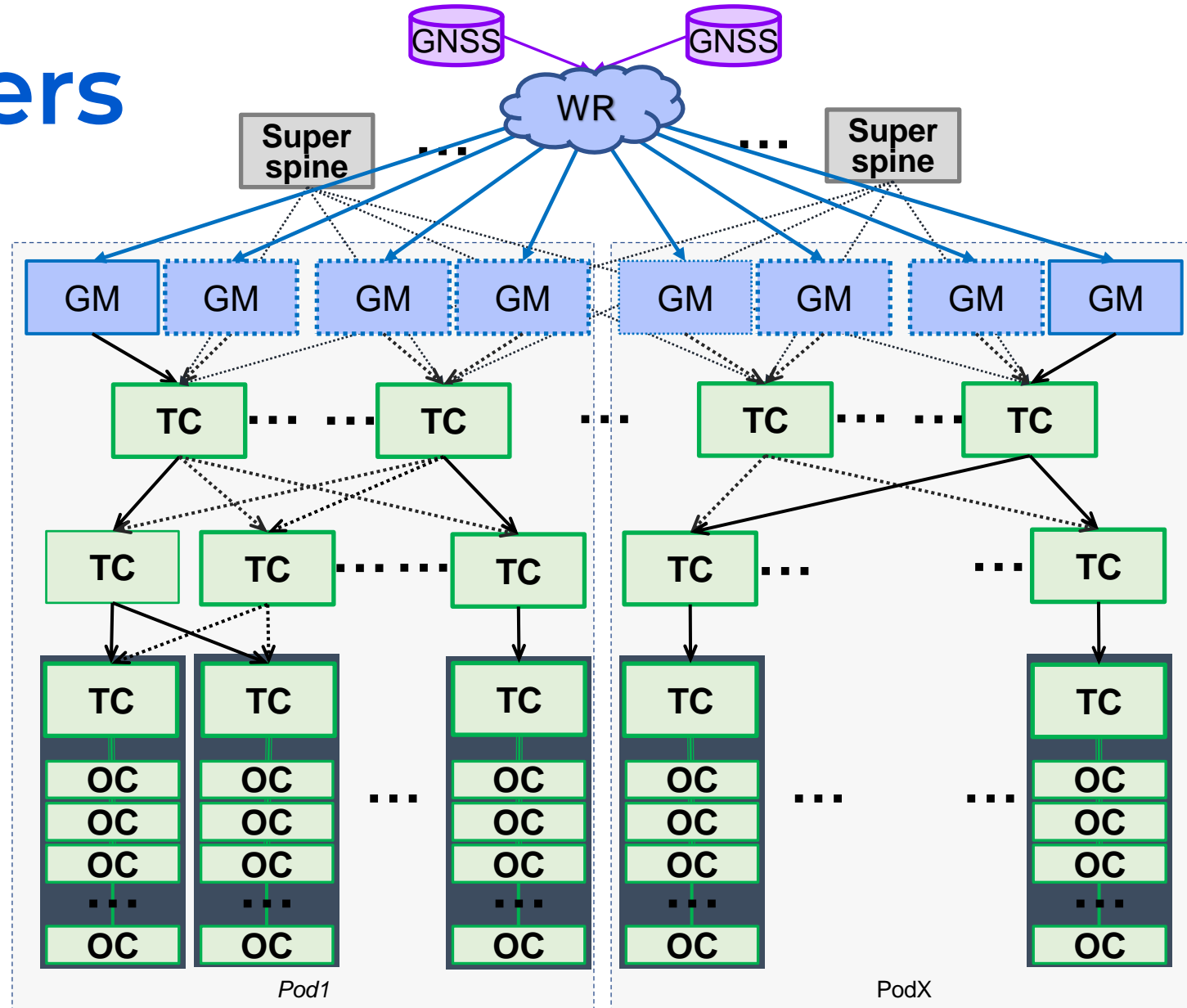


NANOG

# WR Time Servers

*Using WR at the core of DC to synchronize all PTP GM (Open Time Server) at each pods*

- Simple solution to put in place
  - Reduce the number of hops
  - GM handles between 5-15K clients
- Sharing a common clock (<1ns accuracy)
  - Linked clocks increase resiliency and accuracy
  - Solution for intra-DC and inter-DC
  - Only 1 or 2 GNSS receivers to install
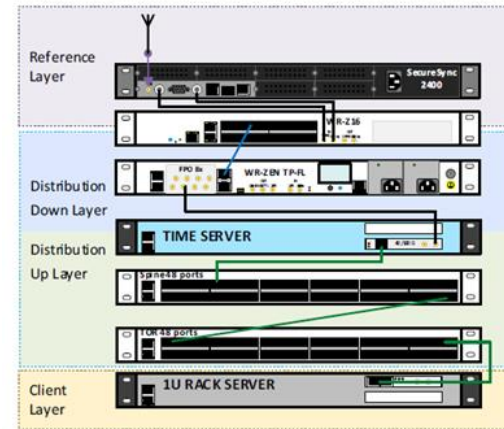  - Relative accuracy is reduced by ±100ns

# White Rabbit for Time Servers

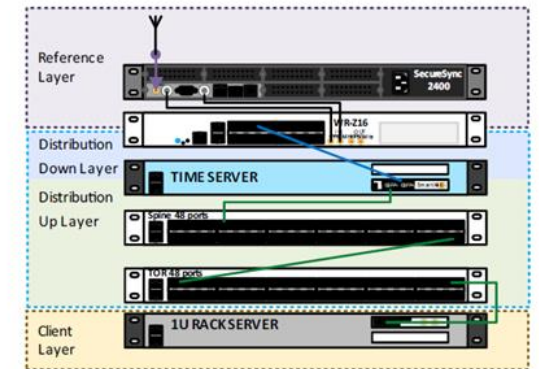*Combining GNSS time servers with High Accuracy timing distribution*

- Synchronize to GPS, SAASM GPS, Galileo, multi-GNSS and many other timing references
- Generate virtually any time and frequency output signals
- Multiple internal oscillator options
- Highly modular (configure-to-order)
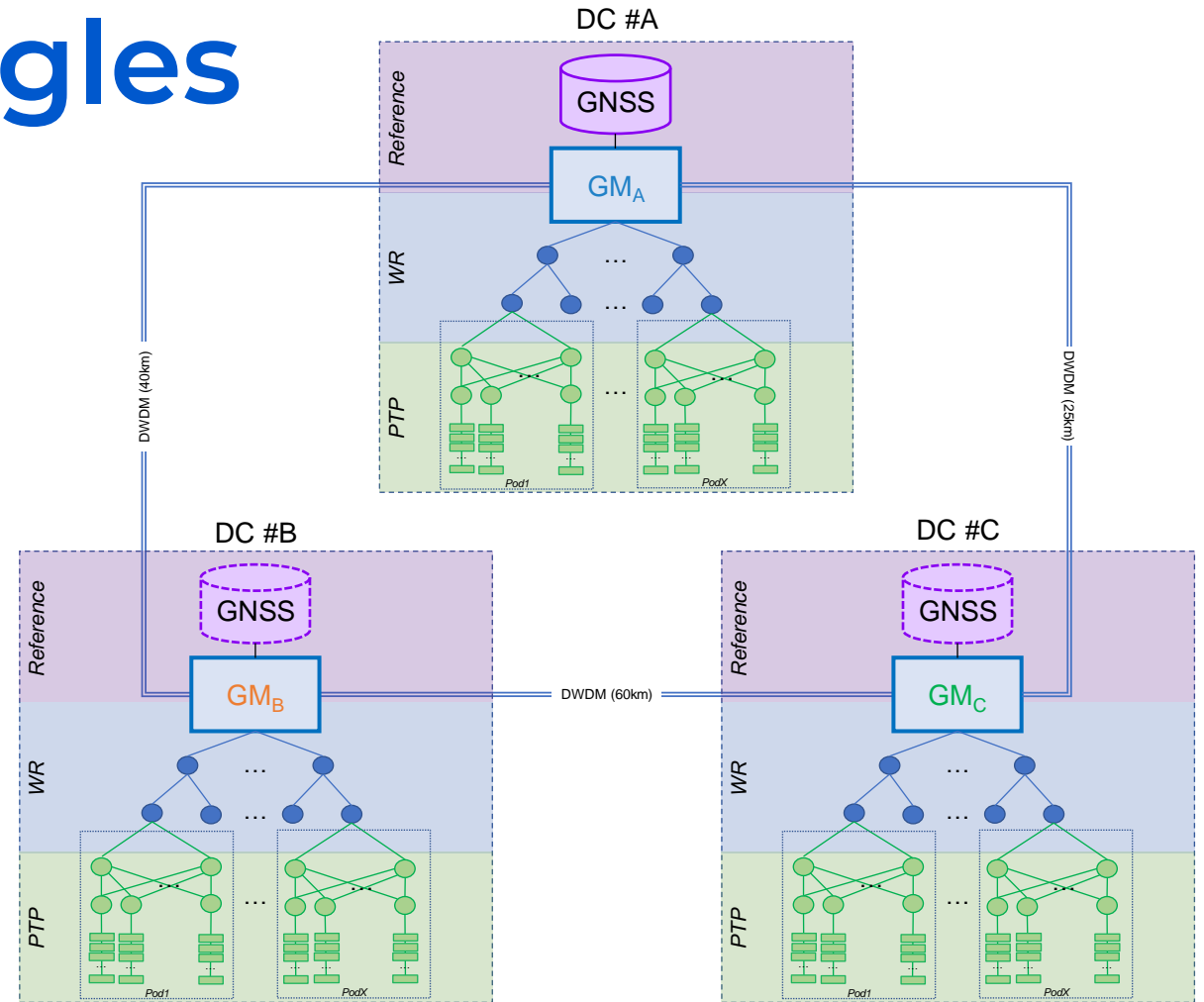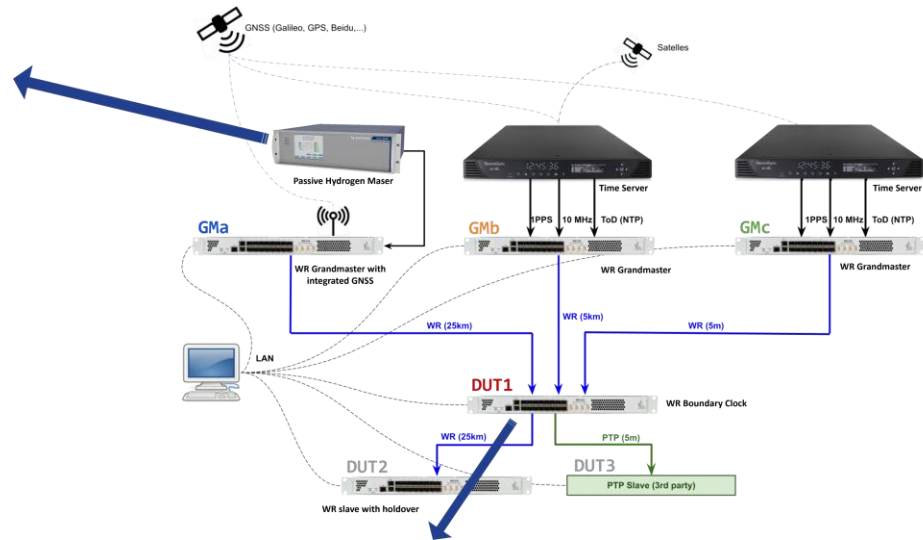
**Through PPS**



**Through WR powered by HATI IP core**



- Sub-nanosecond accuracy
- Picosecond level precision
- Interoperability (PTP, NTP, 10 MHz, 1PPS)
- Failover capability
- Holdover

NANOG™

# Datacenters Triangles

- Multiple GNSS compared through WR links
- Voting mechanism to select the most reliable reference
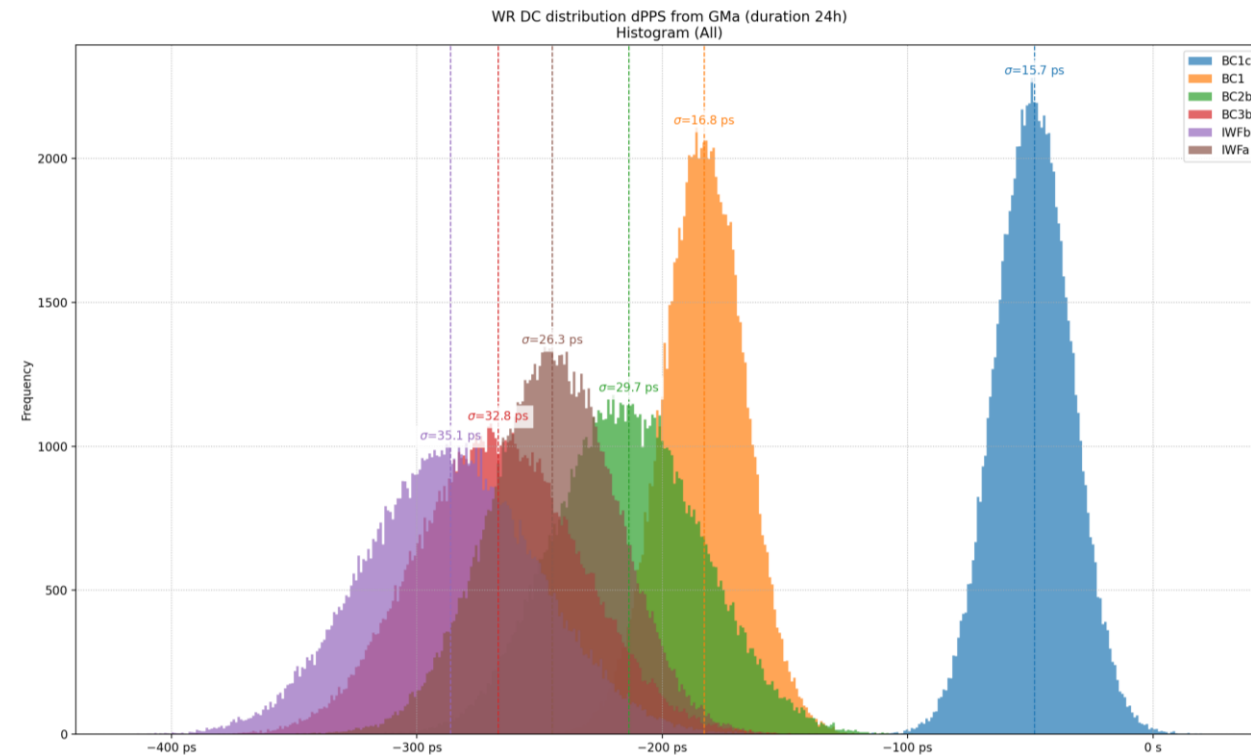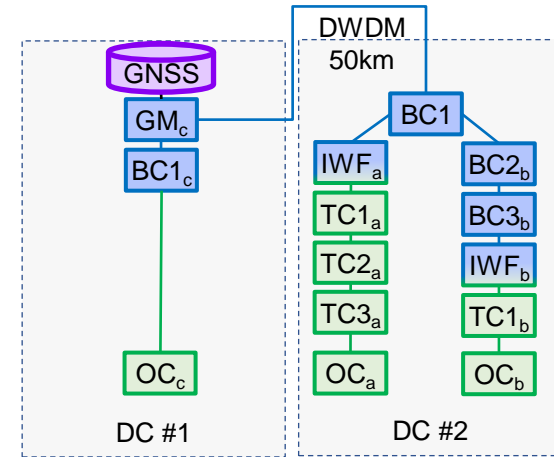- Metro-area connection using DWDM links

# PTP TC/BC vs WR BC

*PTP TC is preferred over PTP BC in datacenter*



- **⊖ PTP BC**
  - ➤ Each hops introduce an error that is propagated
  - ➤ Behaves differently depending on servo and OXCO

-------------

- **⊕ PTP TC for spine/leaf/TOR switches**
  - ➤ Simple to process, no need of specific HW
  - ➤ Well supported by more manufacturers
- **⊖ PTP TC**
  - ➤ Scalability concerns

-------------

- **⊕ WR BC**
  - ➤ Few picoseconds error introduced by each hop
  - ➤ Allows sharing common clock until a specific point and then scales



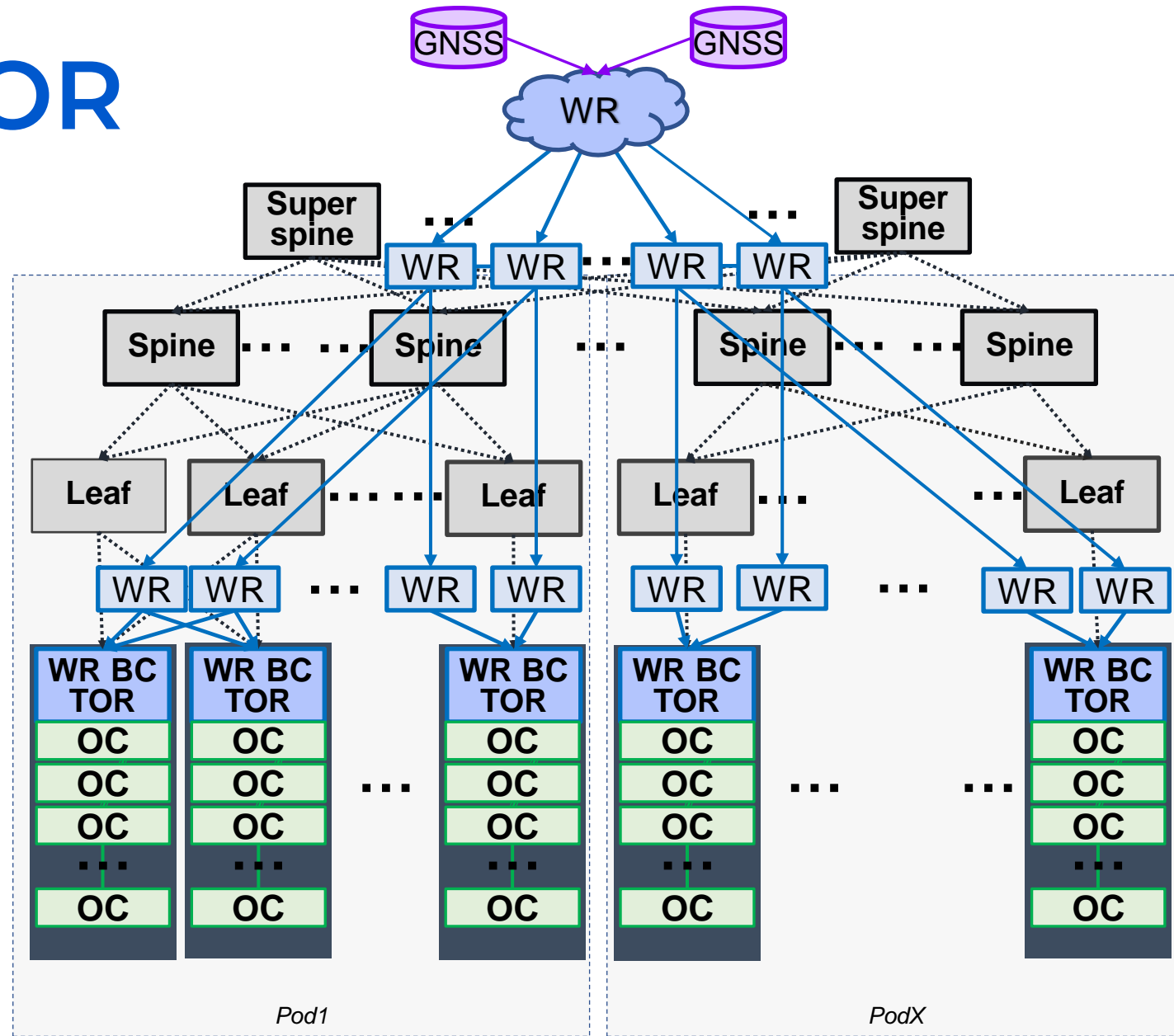WR DC distribution dPPS from GMa (duration 24h)
Histogram (All)

# WR Down to TOR

*Timing network and data network are independent down to TOR that works as BC receiving HA (WR) but transmitting PTP*

⬇

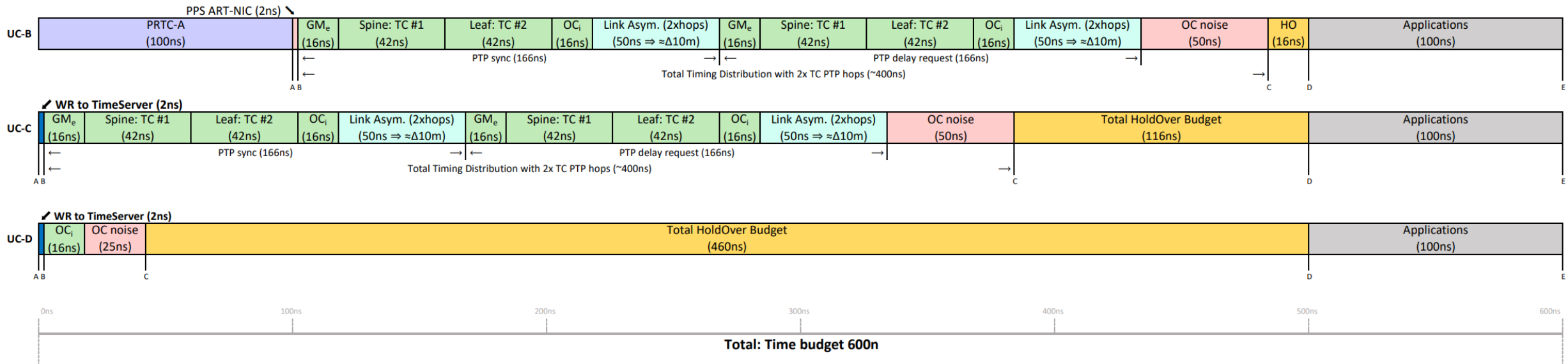➕ Accuracy @ TOR BC <1ns

➕ Only 1 hops PTP
  ➢ Accuracy @ OC Server → 10's ns

➕ Few PTP clients (<50) for TOR BC

➕ Resilient solution

➕ OC NIC can be very basic

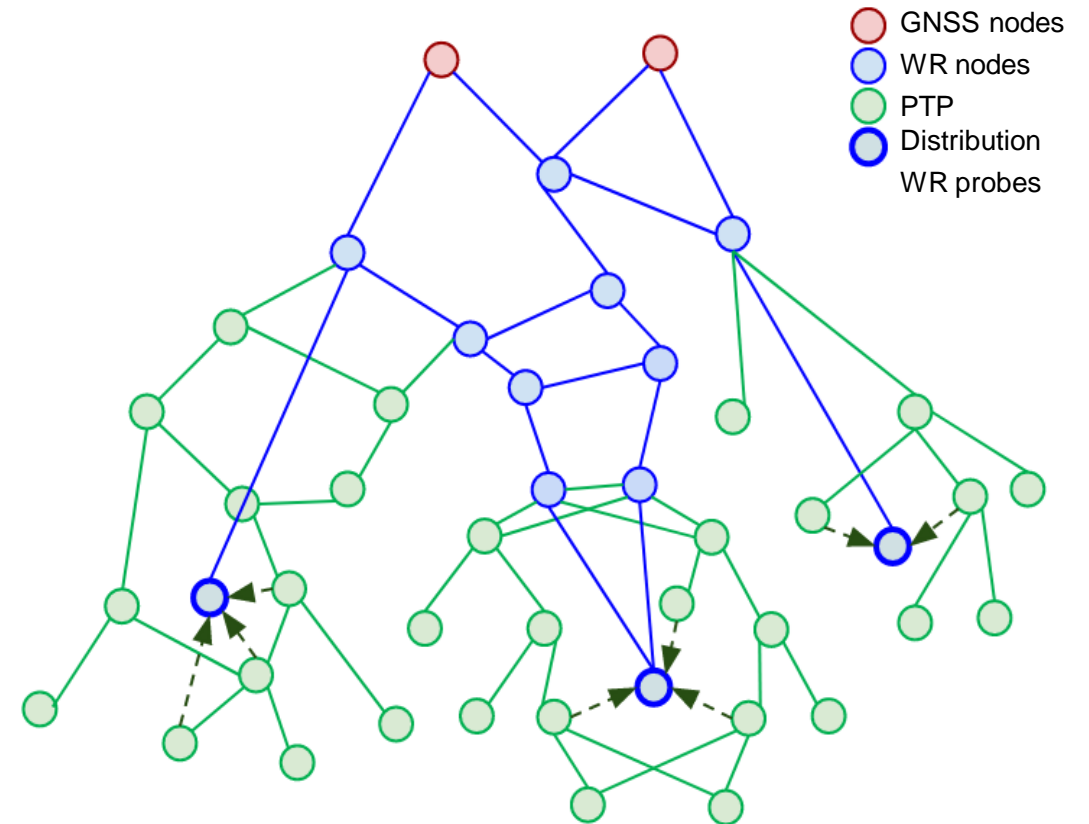➖ Adding a parallel timing network

# Time Budget Optimization

- **Improving accuracy for timing distribution increases holdover budget and thus to enhance resiliency**

- **Through WR a common clock is shared among the DC and thus it allows to:**
  - Remove PTRC-A time-error.
  - Dedicate Holdover budget to final OC node

# Supervision Network

*Using **WR as ground-truth** to monitor the timing distribution PTP DC Profile*

- A well-tested, reliable and deterministic sub-nanosecond accuracy allows one to properly monitor other timing distribution systems. Otherwise, a timing distribution network could be degraded without knowing it.

- Inserting distributed "WR probes" at strategic points allows one to measure the timing performance of "PTP distribution" network in real-time and act in case of unexpected behaviour.



Legend:
- GNSS nodes
- WR nodes
- PTP
- Distribution WR probes

NANOG™

# Wrap Up

## Linking GNSS

The accuracy of WR allows to connect and compare GNSS receiver between them to detect abnormal behaviour. It also reduces the number of GMs.

## Increase Holdover budget

By consuming negligible timing-budget with WR and reducing the number of PTP hops, the reliability is increased thanks to longer holdover budget.

## Supervision Network

Real-time multi-source timing comparison benefiting from the accuracy of WR. It allows to improved traceability and resiliency.

## Future proof solution

Targeting ultra-accurate & reliable timing allows to prepare for future applications needing smaller but still undefined error-bound ($\varepsilon$).

**NANOG**™

# Thank you

FEB-2024

Francisco.Girela@nav-timing.safrangroup.com

safran-navigation-timing.com