# PTP @ Meta

Oleg Obleukhov
Production Engineer
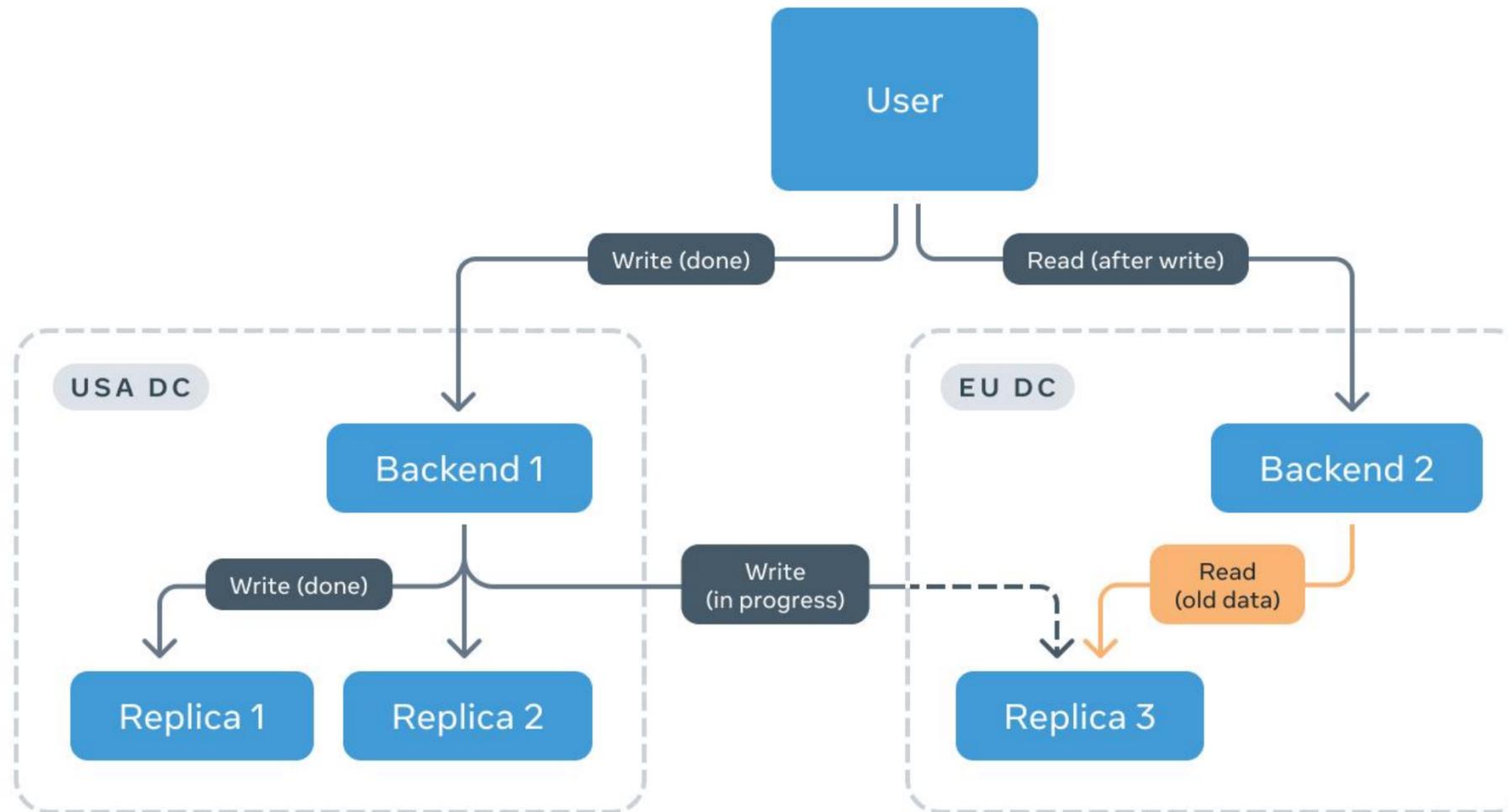
∞ Meta

# Agenda

# The case for PTP
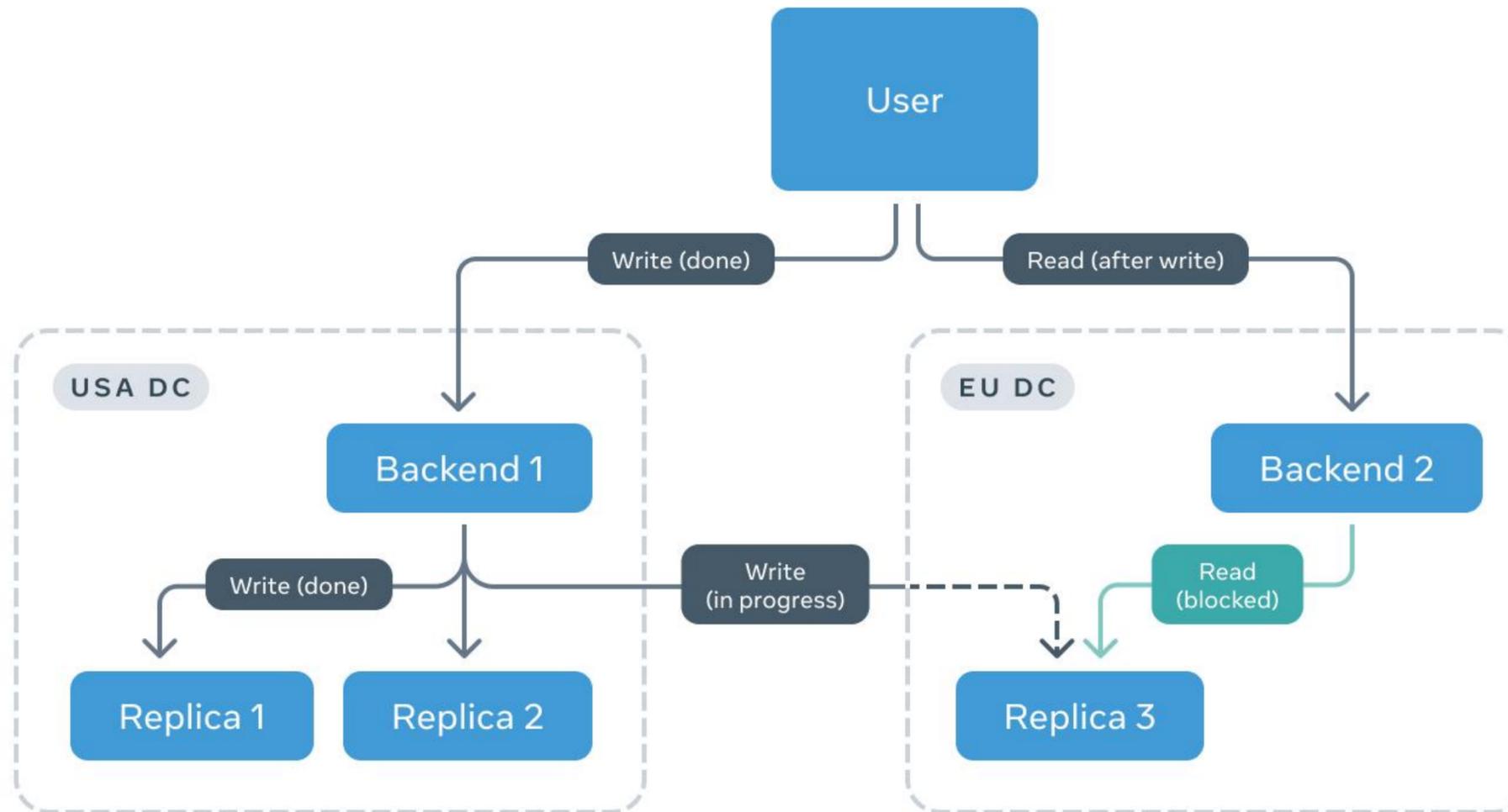
**Schematic representation of read returning outdated information**

**Schematic representation of read returning outdated information**



Commit-wait ensuring consistency guarantee (linearizability)

User

Write (done) — Read (after write)

USA DC

Backend 1

Write (done) — Write (in progress)

EU DC

Backend 2

Read (blocked)

Replica 1    Replica 2

Replica 3

**Commit-wait reads issued against PTP and NTP backed clusters**

# The PTP architecture

# Regional PTP architecture

# The PTP rack

# GNSS



GNSS antenna in one of the Meta Data Center location



Huber-Suhner GNSS-over-fiber technology tested in Meta Dublin office

# The Time Card

## Facebook Time Card

**The Time Card**



Frequency vs. Tempertature

Frequency to temperature ratio of an atomic clock

# The Time Card



Short cable between PPS-out of the Time Card and PPS-in of the NIC



## Offset between the Time Card and the Network Card PHC

# ptp4u

# C4U architecture



```
$ cat /etc/ptp4u.yaml
clockaccuracy: 34
clockclass: 6
draininterval: 30s
maxsubduration: 1h0m0s
metricinterval: 1m0s
minsubinterval: 1s
utcoffset: 37s
```

# The PTP network

# Two-step PTP exchange

**Transparent clock**

PTP Transparent Clock

HW Timestamping NIC

## Transparent Clock and Correction Field

| PTP OC | ────── | *SW | ────── | *SW | ────── | PTP TA |

TC2                              TC1

**T2**

T1, T2, CF$_A$

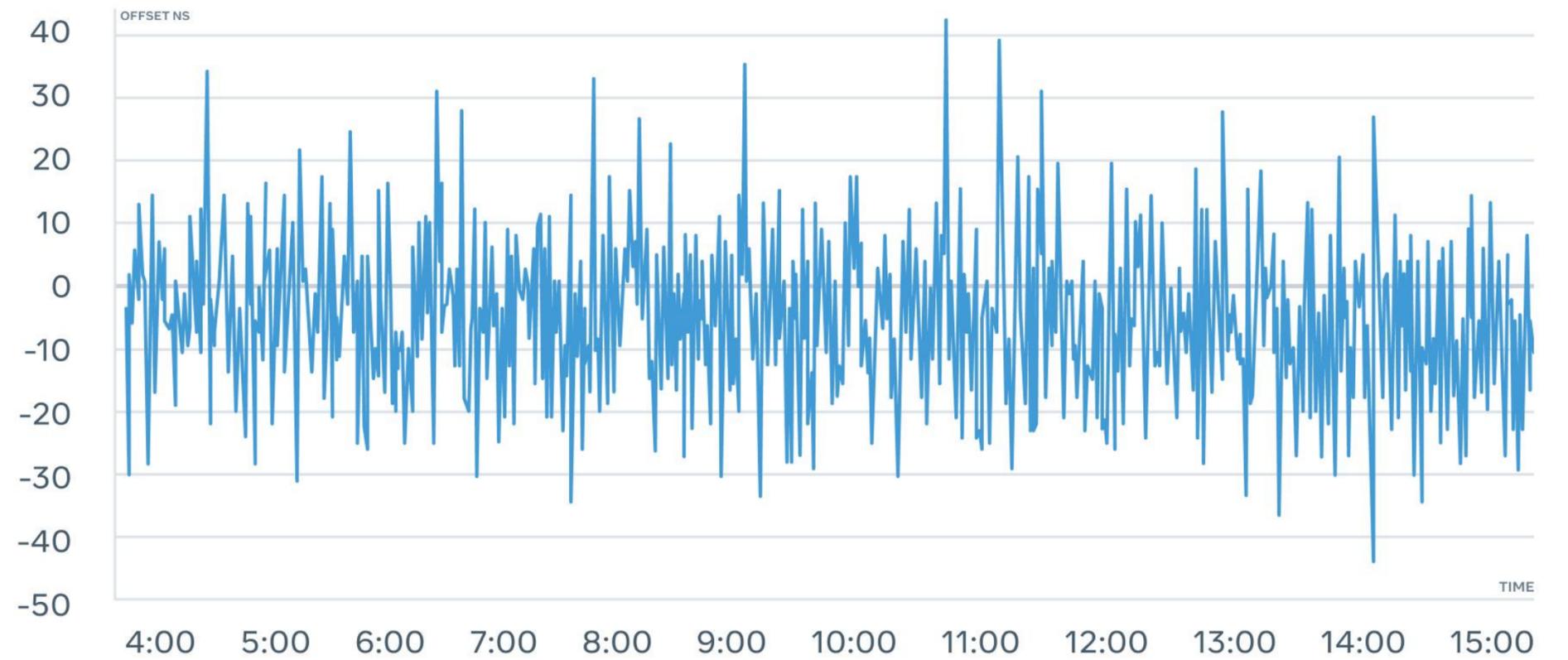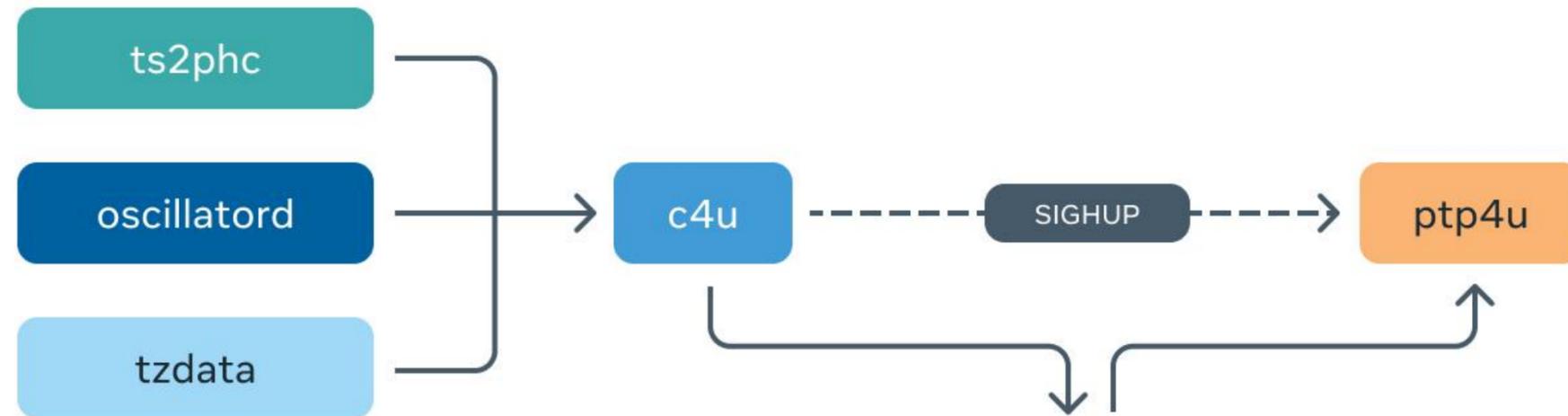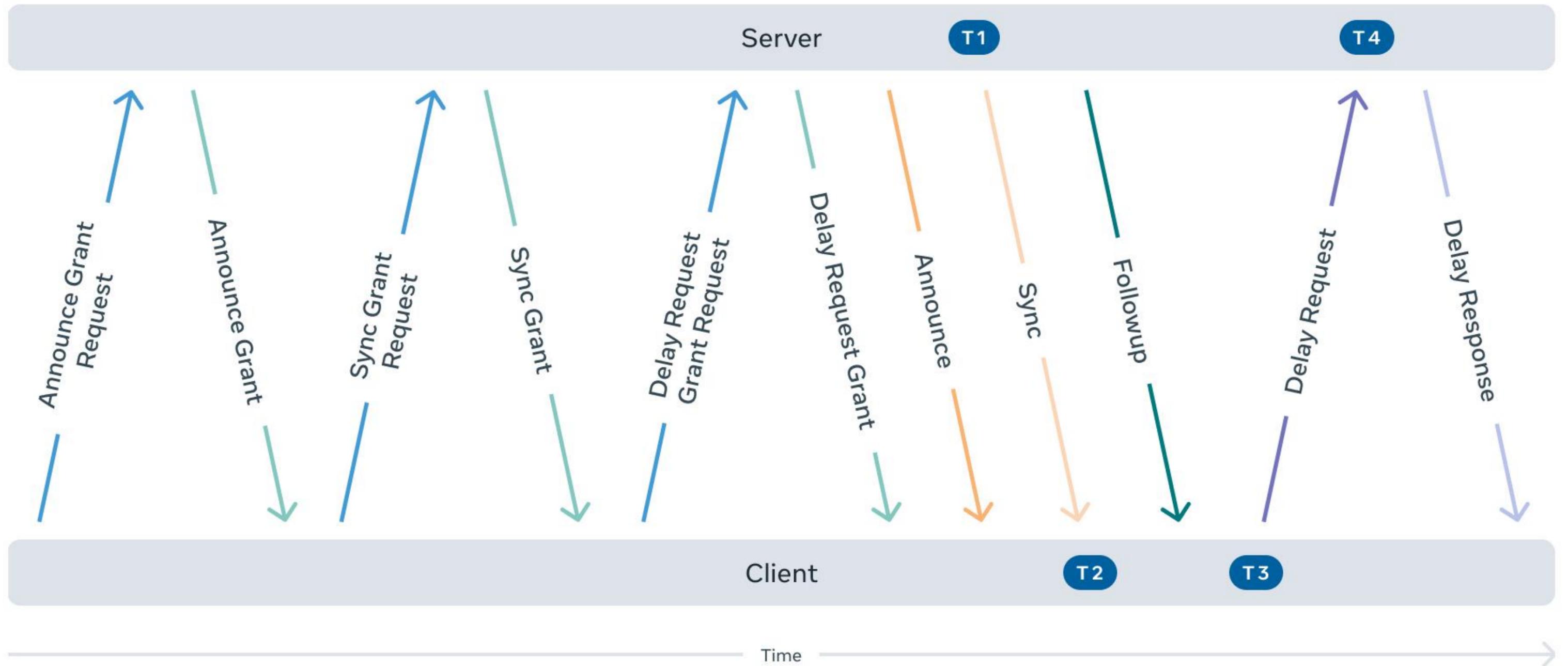| Origin timestamp = T1 CF$_A$ = 864 ns | **SYNC** | + | Origin timestamp = T1 CF$_A$ = 423 ns | **SYNC** | + | Origin timestamp = T1 CF$_A$ = 0 ns | **SYNC** |

**T1**

**T3**

T1, T2, T3, CF$_A$

| CF$_B$ = 0 ns | **DELAY_REQ** | + | CF$_B$ = 451 ns | **DELAY_REQ** | + | CF$_B$ = 874 ns | **DELAY_REQ** |

**T4**

T1, T2, T3, T4, CF$_A$, CF$_B$

| T4 hw timestamp when packet received CF 874 ns | **DELAY_RESP** | T4 hw timestamp when packet received CF 874 ns | **DELAY_RESP** | T4 hw timestamp when packet received CF 874 ns | **DELAY_RESP** |

**PTP HW TIMESTAMPING**

$$mean\_path\_delay = ((T4-T3) + (T2-T1) - CF_A-CF_B) / 2$$

$$clock\_offset = (T2 - T1) - mean\_path\_delay$$

## Transparent clock impact

```
ptp4l[43.662]: offset          -9 s2 freq  -12372 path delay      4114
ptp4l[44.662]: offset          17 s2 freq  -12349 path delay      4114
ptp4l[45.662]: offset          37 s2 freq  -12324 path delay      4078
ptp4l[46.662]: offset         -70 s2 freq  -12420 path delay      4153
ptp4l[47.662]: offset          95 s2 freq  -12276 path delay      4039
ptp4l[48.662]: offset      266776 s2 freq +254434 path delay      4181
ptp4l[49.662]: offset     -430864 s2 freq -363173 path delay    168255
ptp4l[50.662]: offset      -80141 s2 freq -141710 path delay    168255
ptp4l[51.662]: offset      217086 s2 freq +131475 path delay       408
ptp4l[52.662]: offset       16268 s2 freq   -4217 path delay     57459
ptp4l[53.662]: offset        8101 s2 freq   -7504 path delay     57459
ptp4l[54.662]: offset       55912 s2 freq  +42738 path delay      4776
ptp4l[56.305]: offset      -48984 s2 freq  -45385 path delay     19209
ptp4l[56.662]: offset      -37194 s2 freq  -48290 path delay     19209
ptp4l[57.662]: offset       29964 s2 freq   +7710 path delay    -12022
ptp4l[58.662]: offset        9943 s2 freq   -3322 path delay    -12022
ptp4l[59.662]: offset      -19403 s2 freq  -29685 path delay      8279
ptp4l[60.662]: offset        8560 s2 freq   -7543 path delay     -2377
ptp4l[61.662]: offset       -4906 s2 freq  -18441 path delay      6256
ptp4l[62.662]: offset        4197 s2 freq  -10810 path delay      3249
ptp4l[63.662]: offset         979 s2 freq  -12769 path delay      4917
ptp4l[64.662]: offset        1386 s2 freq  -12068 path delay      4917
ptp4l[65.662]: offset        1741 s2 freq  -11297 path delay      4270
ptp4l[66.662]: offset         509 s2 freq  -12007 path delay      4428
ptp4l[67.662]: offset         395 s2 freq  -11968 path delay      4185
ptp4l[68.662]: offset          -7 s2 freq  -12252 path delay      4185
```
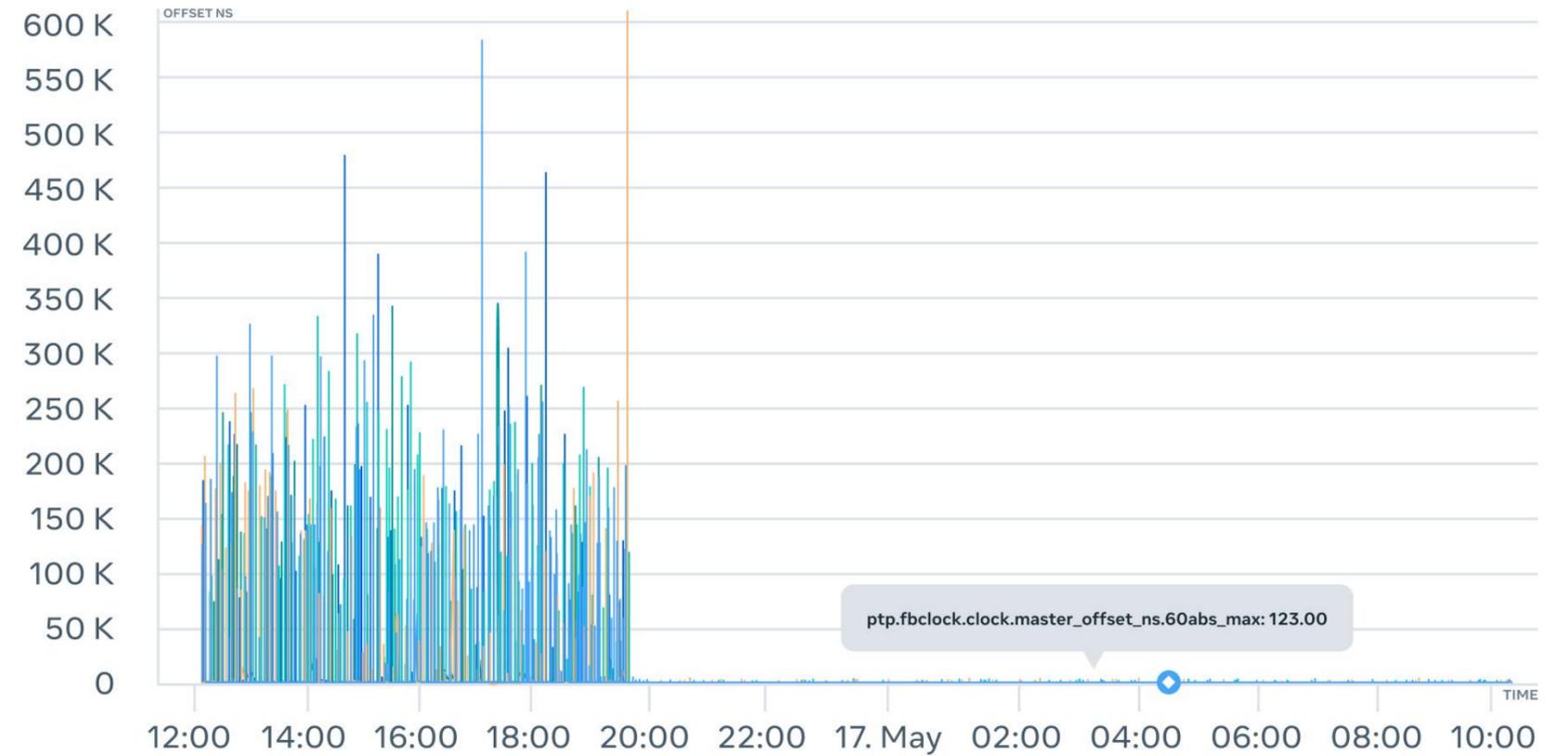


Absolute offset values on hosts connected to the switch without Transparent Clock enabled

# The PTP client

**Hardware timestamps**

```
$ ethtool -T eth0
Time stamping parameters for eth0:
Capabilities:
    hardware-transmit
    hardware-receive
    hardware-raw-clock
PTP Hardware Clock: 0
Hardware Transmit Timestamp Modes:
    off
    on
Hardware Receive Filter Modes:
    none
    All
```

| 128 bits | 64 bits | 64 bits | 64 bits |
|---|---|---|---|
| Socket control message header | Software Timestamp | Legacy Timestamp | Hardware Timestamp |

```
ptp4l[40.432]: offset        -16 s2 freq  -13105 path delay      3493
ptp4l[41.432]: offset         -6 s2 freq  -13100 path delay      3493
ptp4l[42.432]: offset          9 s2 freq  -13087 path delay      3493
ptp4l[43.432]: offset         -5 s2 freq  -13098 path delay      3493
ptp4l[44.432]: offset          1 s2 freq  -13093 path delay      3493
ptp4l[45.432]: spike detected => max_offset_locked: 33, setting offset to min_offset_freq_mean: -13065.039314
ptp4l[46.432]: skip 1/15 large offset (>33) 224401
ptp4l[47.432]: offset        -21 s2 freq  -13115 path delay      3493
ptp4l[48.432]: offset          9 s2 freq  -13091 path delay      3493
ptp4l[49.432]: offset         10 s2 freq  -13088 path delay      3493
ptp4l[50.432]: offset         -8 s2 freq  -13103 path delay      3493
```
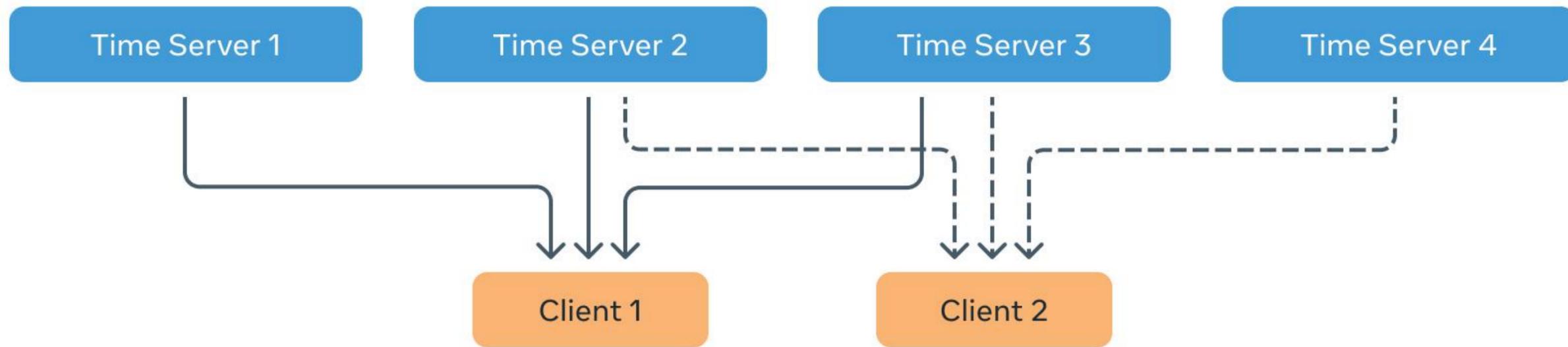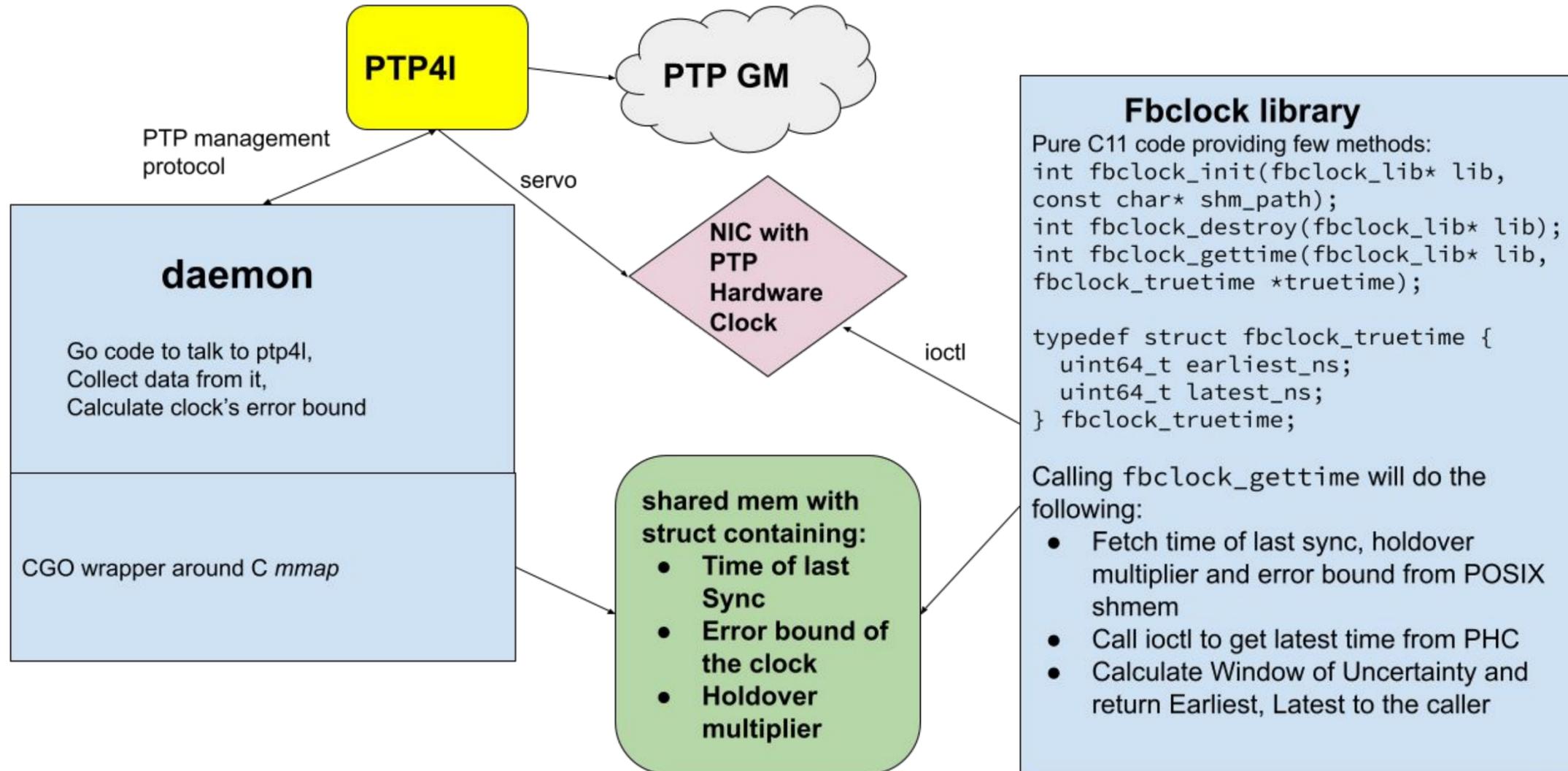
Schematic representation of sharding

# fbclock

**PTP4I**

**PTP GM**

PTP management protocol

servo

## daemon

Go code to talk to ptp4l,
Collect data from it,
Calculate clock's error bound

CGO wrapper around C *mmap*

**NIC with PTP Hardware Clock**

ioctl

**shared mem with struct containing:**
- **Time of last Sync**
- **Error bound of the clock**
- **Holdover multiplier**

## Fbclock library

Pure C11 code providing few methods:
```
int fbclock_init(fbclock_lib* lib,
const char* shm_path);
int fbclock_destroy(fbclock_lib* lib);
int fbclock_gettime(fbclock_lib* lib,
fbclock_truetime *truetime);

typedef struct fbclock_truetime {
  uint64_t earliest_ns;
  uint64_t latest_ns;
} fbclock_truetime;
```

Calling `fbclock_gettime` will do the following:
- Fetch time of last sync, holdover multiplier and error bound from POSIX shmem
- Call ioctl to get latest time from PHC
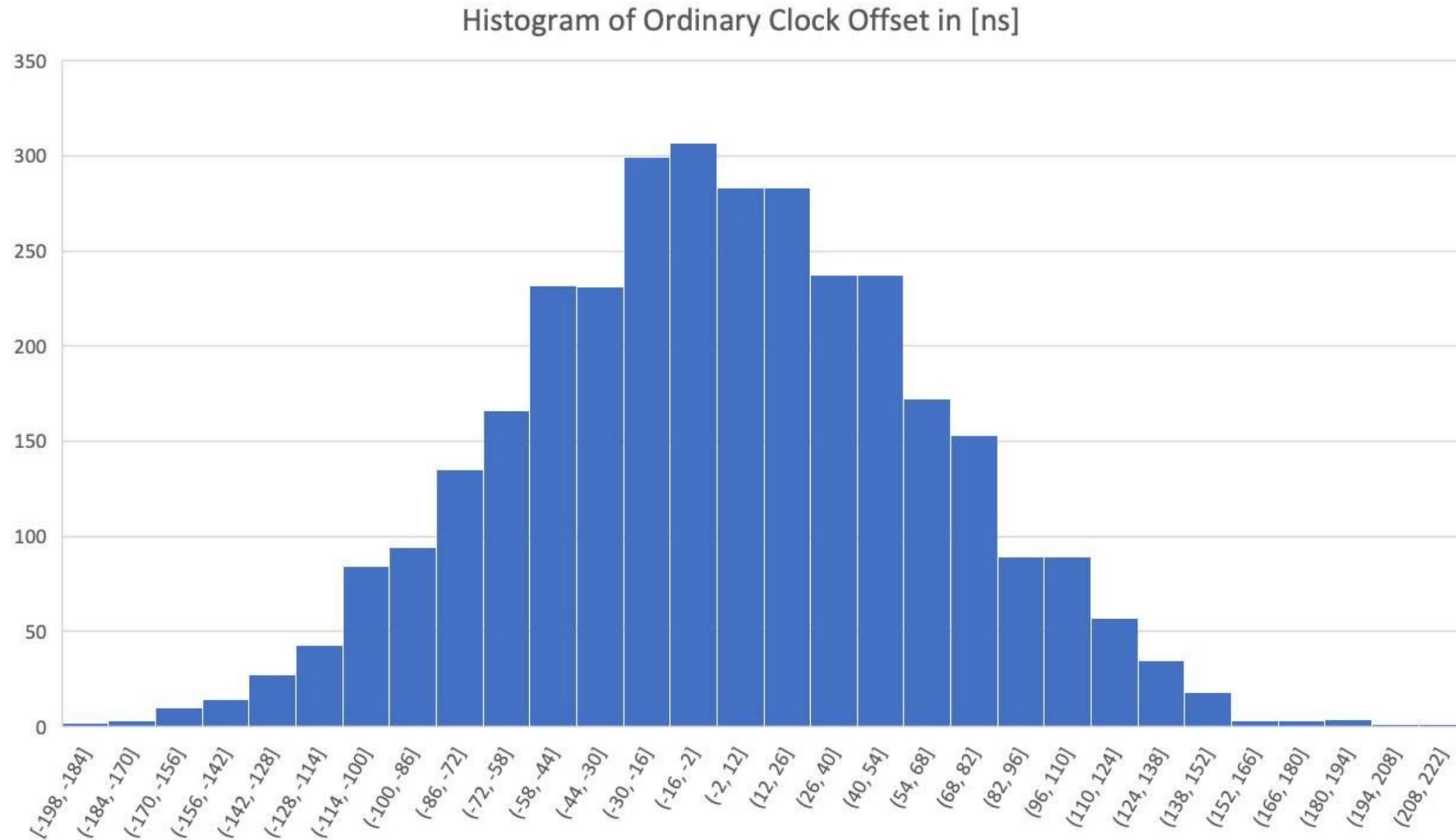- Calculate Window of Uncertainty and return Earliest, Latest to the caller

# fbclock

Estimated E2E Variance = [GNSS Variance + MAC Variance + ts2phc Variance] + [PTP4L Offset Variance] = [Time Server Variance] + [Ordinary Clock Variance]

Estimated E2E Variance = (12ns ^ 2) + (43ns ^ 2) + (52ns ^2) + (61ns ^2) = 8418 which corresponds to 91.7 ns

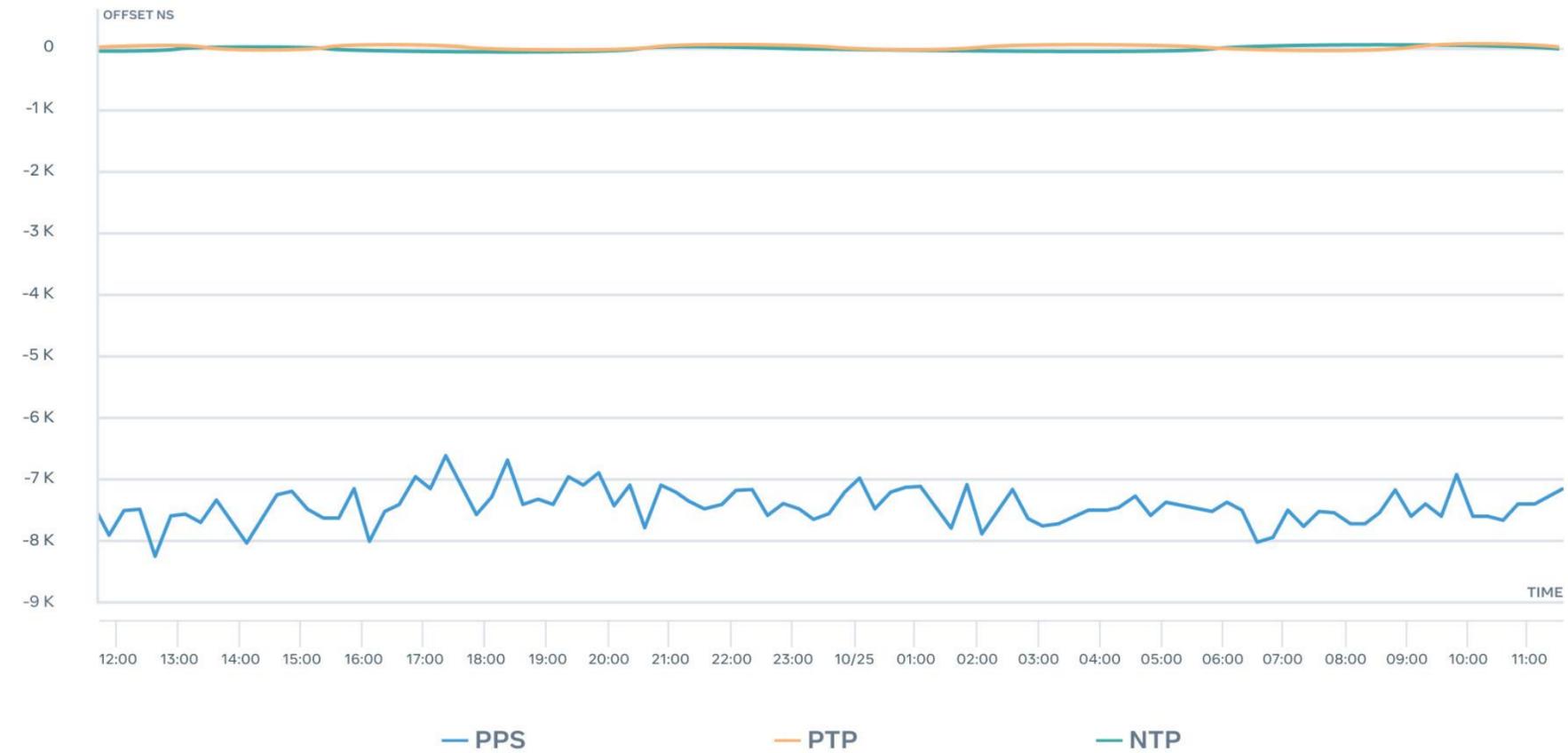$$Var\left(\sum_{i=1}^{n} X_i\right) = \sum_{i=1}^{n} Var(X_i)$$

## Histogram of Ordinary Clock Offset in [ns]

# How we monitor the PTP architecture

**Calnex**





Calnex Sentinel monitoring data

## ptpcheck

```
$ ptpcheck diag
[ OK ] GM is present
[ OK ] Period since last ingress is 972.752664ms, we expect it to be within 1s
[ OK ] GM offset is 67ns, we expect it to be within 250µs
[ OK ] GM mean path delay is 3.495µs, we expect it to be within 100ms
[ OK ] Sync timeout count is 1, we expect it to be within 100
[ OK ] Announce timeout count is 0, we expect it to be within 100
[ OK ] Sync mismatch count is 0, we expect it to be within 100


$ ptpcheck fbclock
{"earliest_ns":1654191885711023134,"latest_ns":1654191885711023828,"wou_ns":694}


$ ptpcheck phcdiff -d /dev/ptp0 -d /dev/ptp2
PHC offset: -15ns
Delay for PHC1: 358ns
Delay for PHC2: 2.588µs


$ ptpcheck sources
+----------+----------------------+-------------------------+-----------+--------+----------+---------+-----------+-----------+--------------+
| SELECTED |       IDENTITY       |         ADDRESS         |   STATE   | CLOCK  | VARIANCE |  P1:P2  | OFFSET(NS) | DELAY(NS) |  LAST SYNC   |
+----------+----------------------+-------------------------+-----------+--------+----------+---------+-----------+-----------+--------------+
| true     | abcdef.fffe.111111-1 | time01.example.com.     | HAVE_SYDY | 6:0x22 | 0x59e0   | 128:128 |        27 |      3341 | 868.729197ms |
| false    | abcdef.fffe.222222-1 | time02.example.com.     | HAVE_ANN  | 6:0x22 | 0x59e0   | 128:128 |           |           |              |
| false    | abcdef.fffe.333333-1 | time03.example.com.     | HAVE_ANN  | 6:0x22 | 0x59e0   | 128:128 |           |           |              |
```
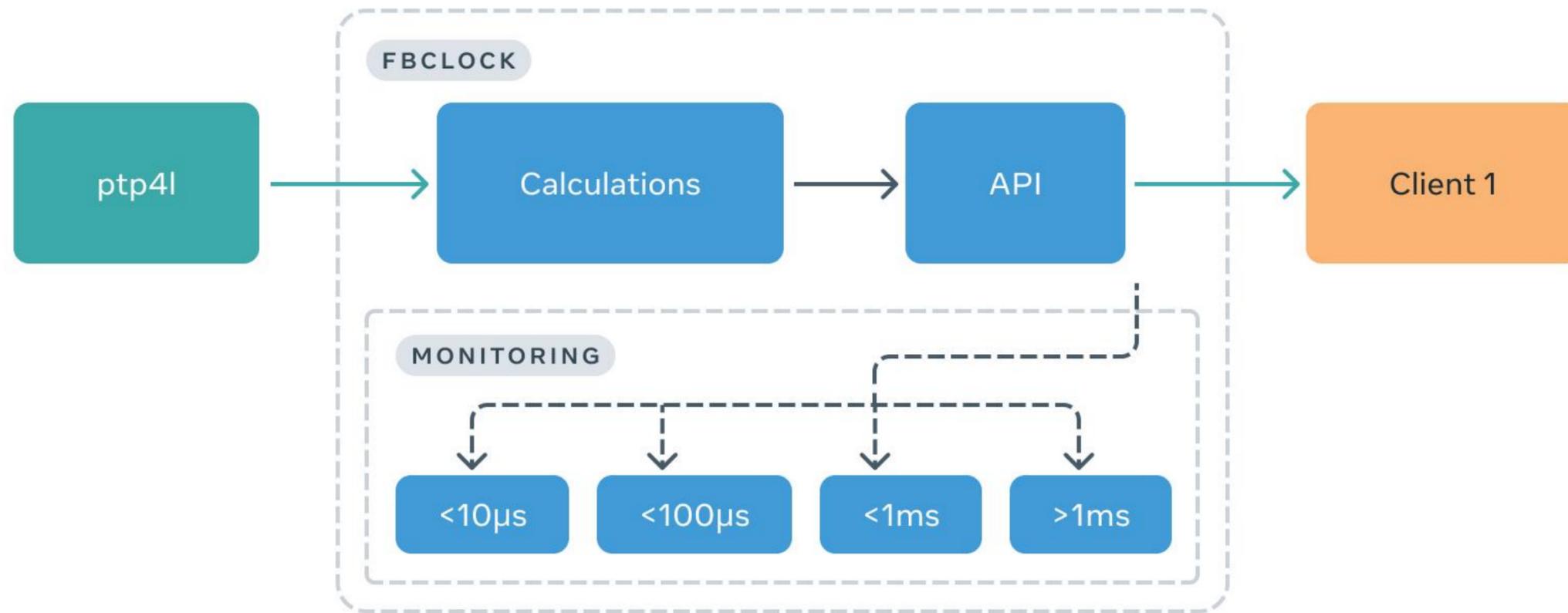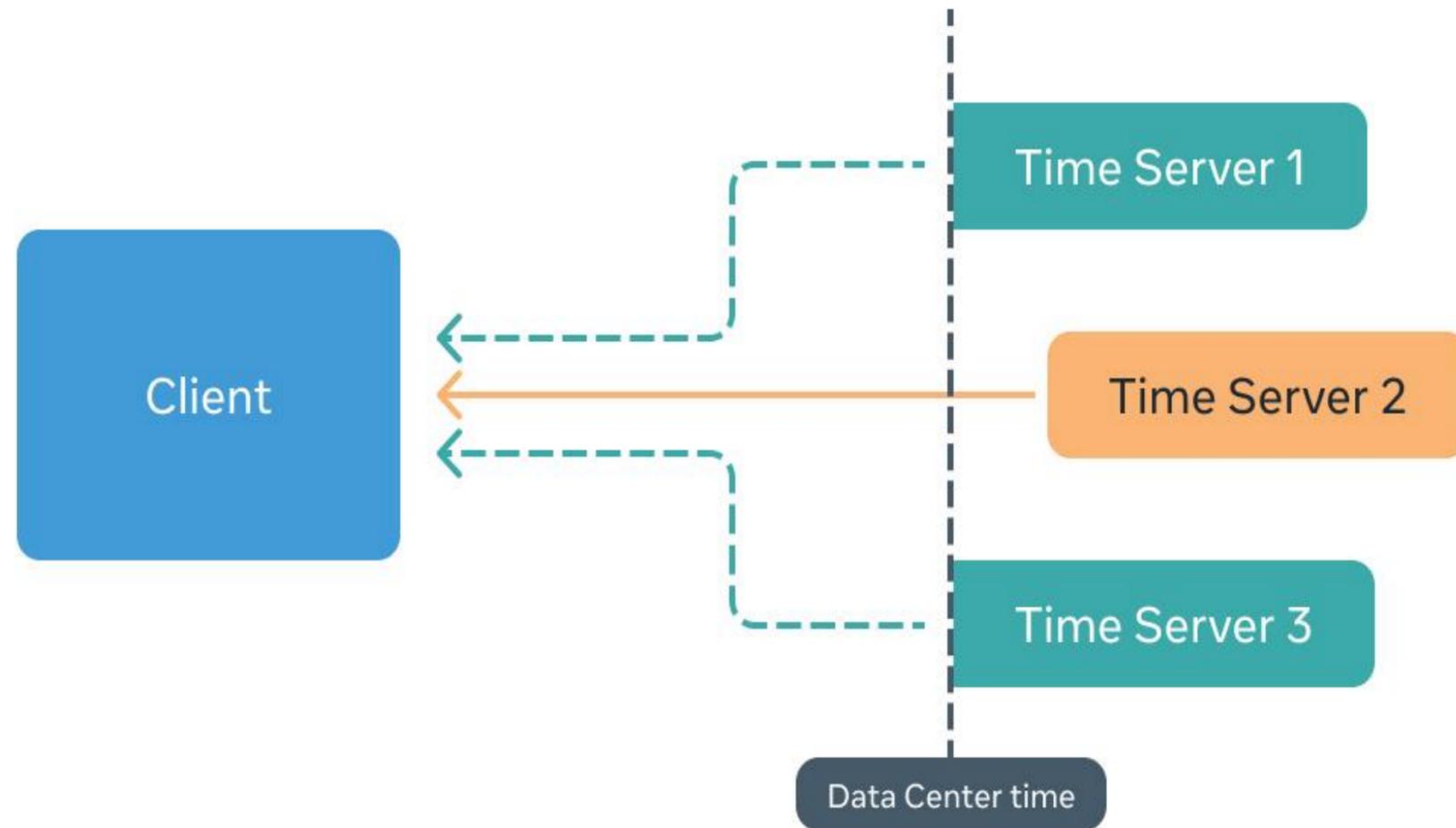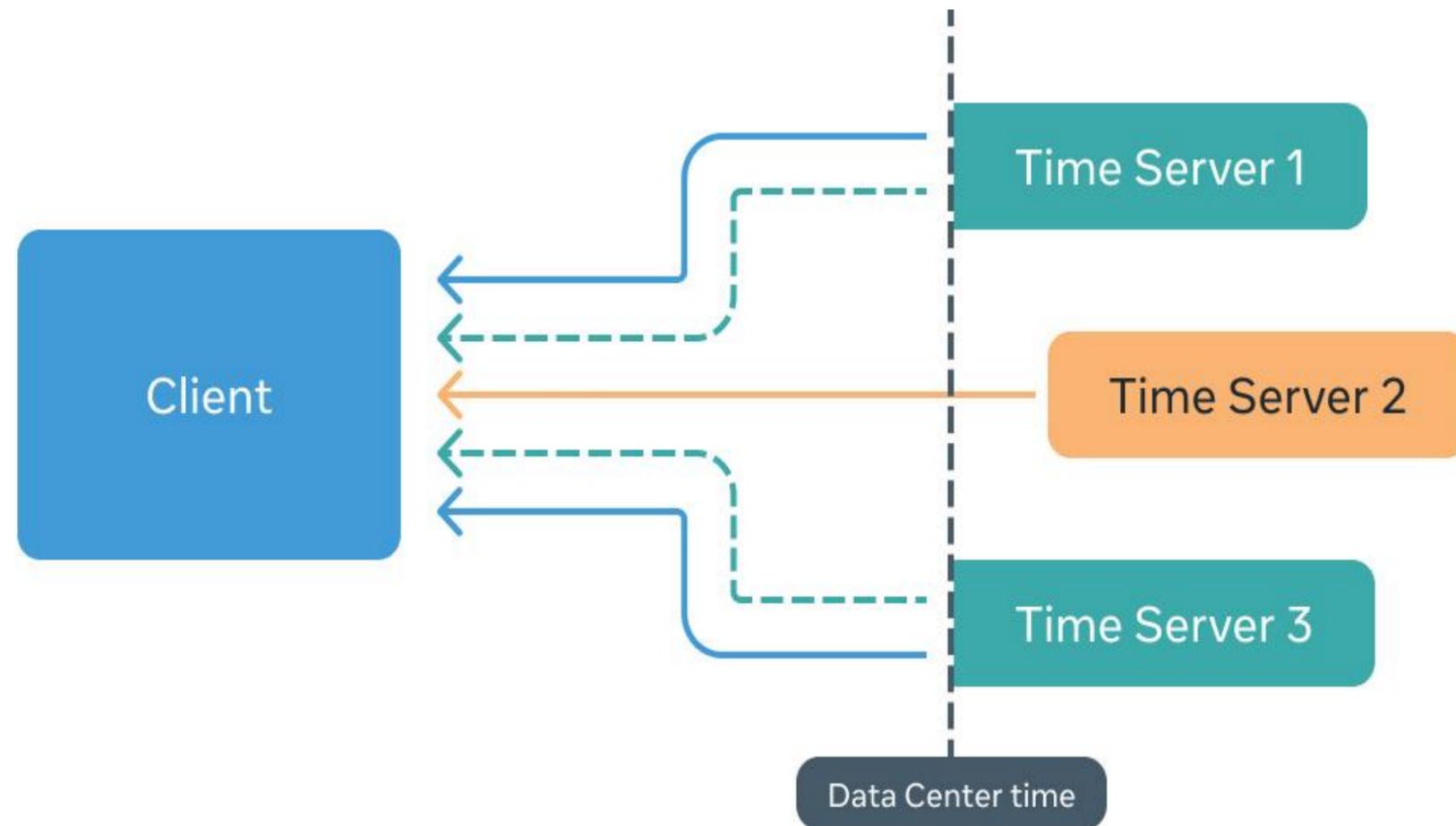
Client following Time Server 2

Thank you

Meta