

# Insights into Trading System Dynamics

## Deutsche Börse's T7<sup>®</sup>

September 2024



DEUTSCHE BÖRSE  
GROUP

# Content

3

Introduction

7

Recent Developments

17

Latency Analysis

27

Market data

39

What you need to be fast

42

T7 Overview

52

Appendix

# 3

## Introduction



# T7<sup>®</sup> Technology Roadmap

Deutsche Börse is pursuing its technology roadmap to deliver innovative and superior trading technology.

## Recent developments

- Q1 2024: Consolidation of EMDI & Matching Engine process
- 9/16 March 2024: Tech refresh of Colo 2.0 capture infrastructure
- 25 March 2024: Introduction of additional Colo 2.0 order entry and market data switches for Eurex
- 22 April 2024: Tech refresh of Colo 2.0 market data switches for Xetra A-side
- 13 May 2024: T7 Release 12.1
- 29 June 2024: Software upgrade on all Access Layer capture devices
- 17/31 August 2024: Introduction of Colo 2.0 Mid-Layer switch for Eurex
- 31 August 2024: Start of technical upgrade of T7 core trading system and T7 backend capture infrastructure.
- 7 September 2024: Tech refresh of Colo 2.0 order entry switches for Xetra A-side

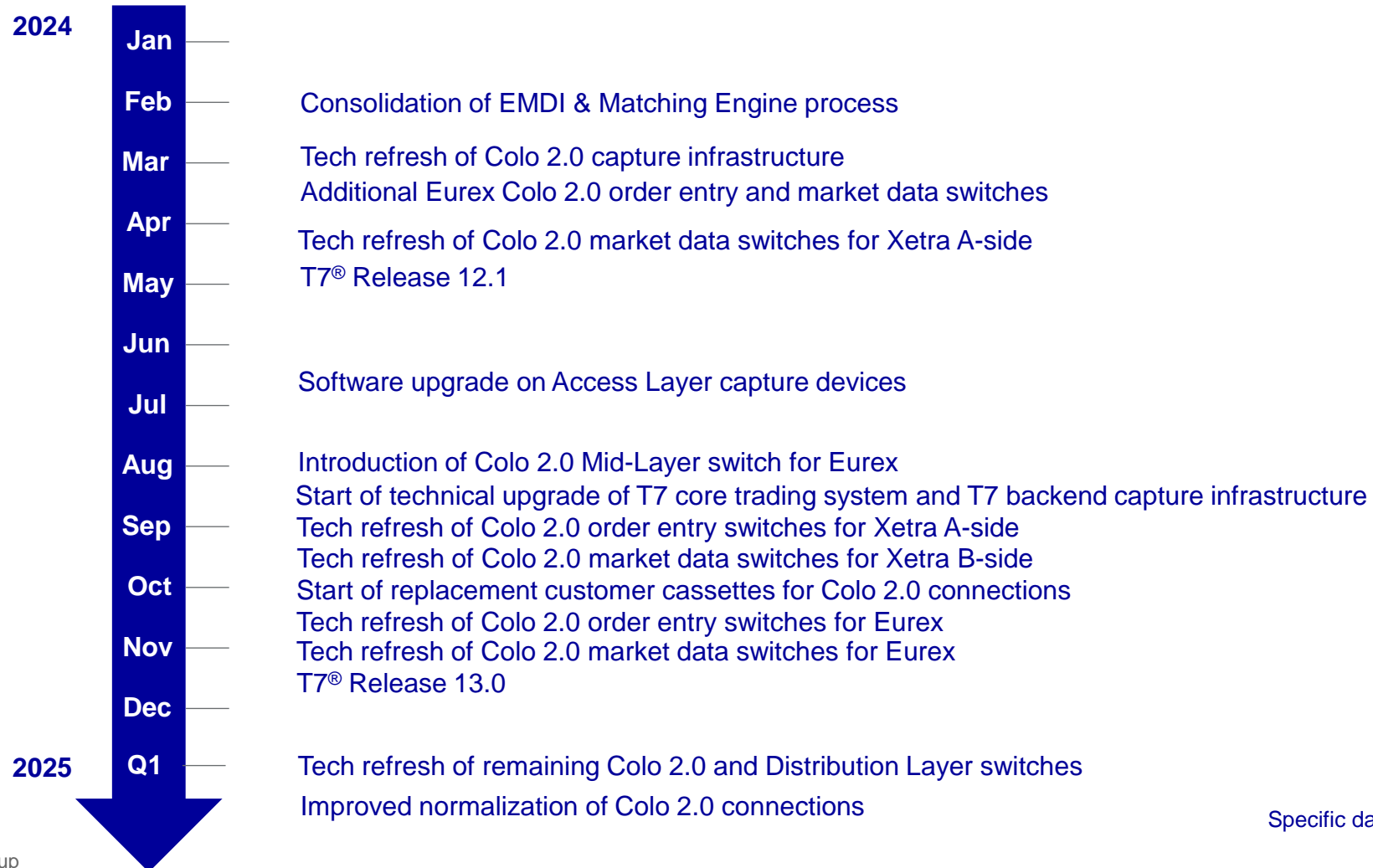
## Outlook

- Q3/Q4 2024: Continue technical upgrade of T7 core trading system and T7 backend capture infrastructure.
- 28 September: Tech refresh of Colo 2.0 market data switches for Xetra B-side
- 5 October 2024: Start replacement of customer cassettes for Colo 2.0 connections
- 12/26 October 2024: Tech refresh of Colo 2.0 order entry switches for Eurex
- 9/23 November 2024: Tech refresh of Colo 2.0 market data switches for Eurex
- 18 November 2024: T7 Release 13.0
- Q1 2025: Tech refresh of remaining Colo 2.0 and Distribution Layer switches
- 2025: Improved normalization of Colo 2.0 connections

Specific dates will be announced separately.

# T7<sup>®</sup> Technology Roadmap

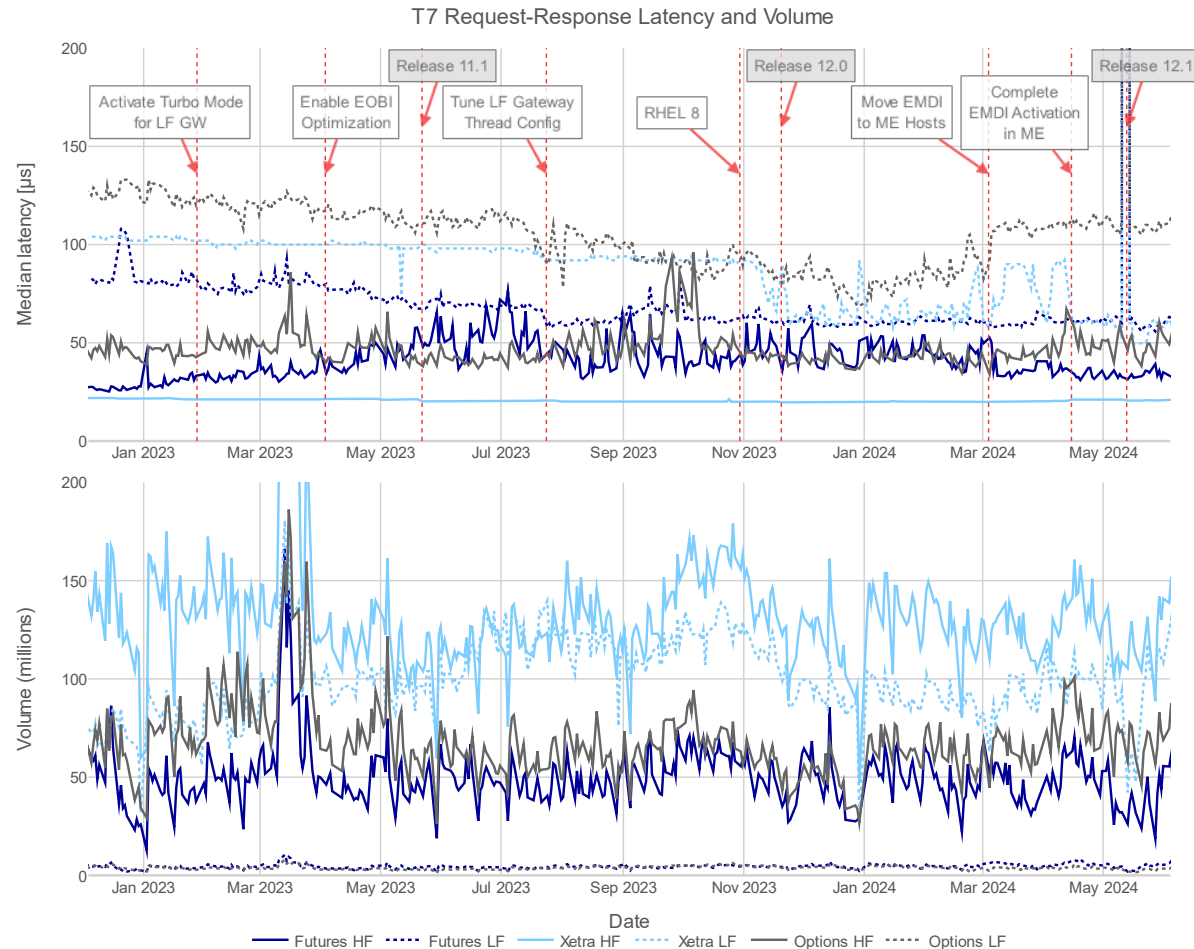
## Timeline of updates



Specific dates will be announced separately.

# Processed Transactions and Response Times

## T7 request – response round-trip times



- Deutsche Börse continuously invests in its trading system and is holding up transparency while providing a low latency trading venue.
- We have continued to add functionality while at the same time tuning our system further.
- The latest consolidation of EMDI process into Matching Engine reduces complexity and enables deterministic distribution of market data via EMDI while ensuring that EOBI is always faster.

# 7

## Recent Developments



# EOBI DSCP Field and Discard IP Range

## Speculative Triggering

### What is speculative triggering?

Initiation of sending a request by a participant at a time when not all required information are available.

There are certain triggers in a market data packet that reveal information:

1. Ethernet preamble of market data (earliest time = 0 ns)
2. Destination MAC (ethernet header) = identification of the multicast stream
3. IP total length: Indicates the type of message contained, e.g. ExecutionSummary or single order add
4. Payload (up to 100 ns later): Quantity and Price.

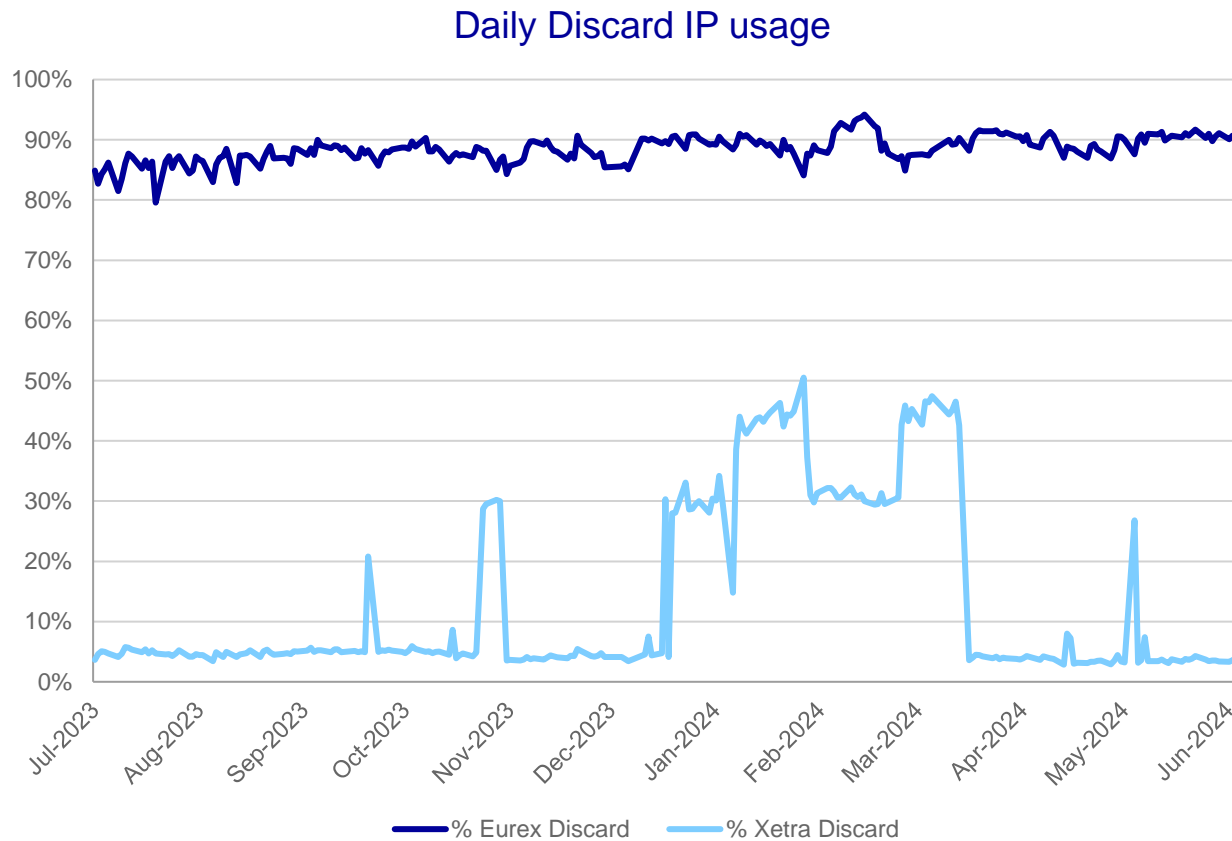
### Why is it used?

Speculative triggering is incentivized by the deterministic technical network architecture of the T7 trading system. The 10 Gbit/s Access Layer switches operate in a 'cut-through' mode and the first bytes of an ethernet frame reserve priority on the competitive network path for the uplink to the next switch. This incentivizes latency sensitive trading participants to send technical transactions purely to reserve switch priority creating a high load on the T7 system. To avoid unnecessary technical transactions on the T7 a technical solution has been implemented.



# EOBI DSCP Field and Discard IP Range

## Confining the effects of speculative triggering



In July 2020, Eurex introduced a Discard IP address range 172.16.0.0/16 on the 10 Gbit/s order entry networks. Xetra followed in May 2021. Trading participants may send falsely 'speculative' triggered packets to the discard IP range, instead of sending it to the exchange. These packets will be discarded at the Access Layer switch port and no other participant is influenced. Packets sent to the discard IP address are not considered to be orders and are not forwarded to the exchange.

To enable market participants to effectively use this discard IP address, the DSCP field of the IPv4 headers in EOBI market data packets is used. Four different bits indicate the most common 'interesting' market situations.

The number of packets reaching the Matching Engines decreased significantly after introduction of this measure. Since then, the number of Discard IP packets has constantly increased.

The graph on the left shows the percentage of Discard IP packets compared to total number of packets reaching the Access Layer switches as an average per day.

# EOBI DSCP Field and Discard IP Range

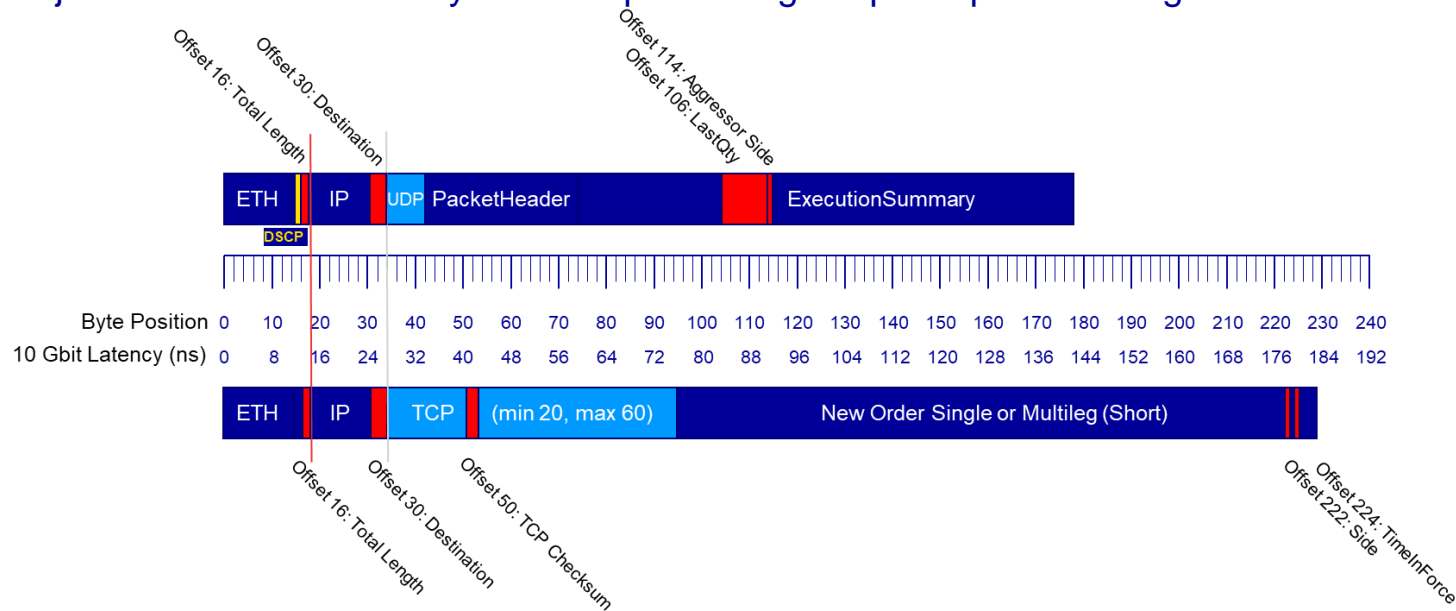
## Confining the effects of speculative triggering (continued)

### A technical solution to the speculative triggering problem:

- Mark potential triggers early in the IP header of market data packet with the help of DSCP (Differentiated Services Code Point) flags.
- Offer a non-competitive Discard IP destination address to enable packets to be discarded right on the Access Layer switch.

DSCP flags indicate Execution summaries and/or widening or narrowing of the bid/ask spread from orders (not quotes).

Examining these flags allows participants to still change the destination IP address of an in-flight outgoing message to the Discard IP address for uninteresting packets. These packets will be sent on the exchange network but will not reach the trading system as they will be rejected on the Access Layer switch port facing the participant sending the transaction.



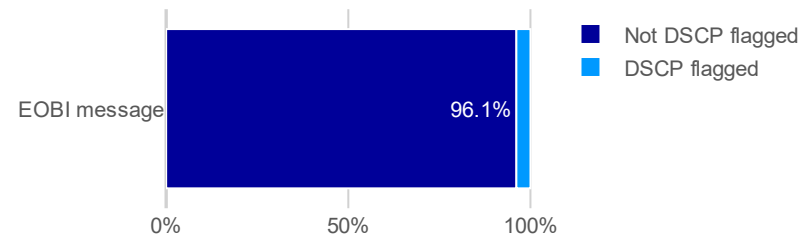
- The left figure shows at which byte position different information is available in the message.
- A market data packet is received at time 0.
- The fastest participants may react and send a response as early as the first bit of market data has been received, dynamically reading the market data packet while already streaming out the response.
- The response packet may be modified in-flight, after reading e.g. the DSCP flags, total length etc.
- Even with a reaction time of 0 the outgoing destination IP address can still be modified after evaluating the DSCP flags of the incoming market data packet.

# DSCP Statistics for Selected Products

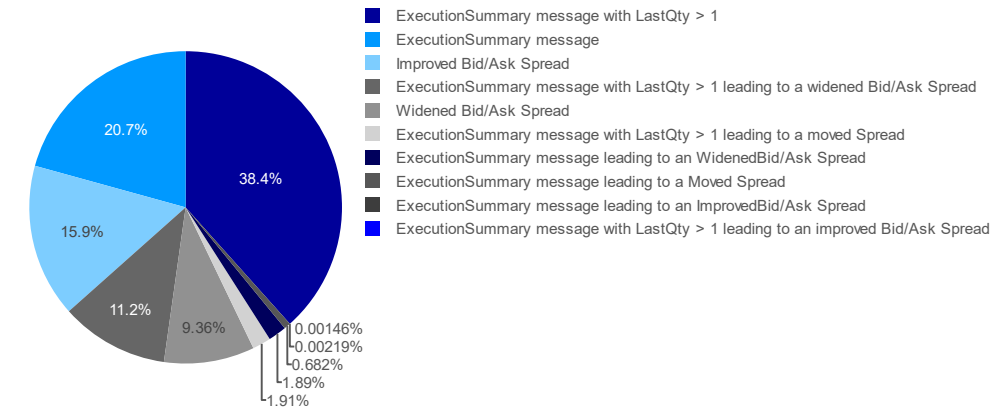
## FESX and OESX

### FESX DSCP Statistics 2024-06-04

DSCP Statistic

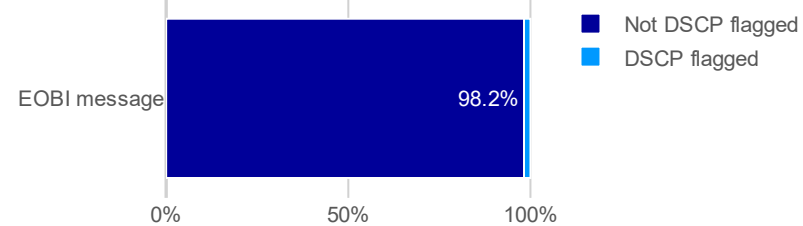


DSCP Flags

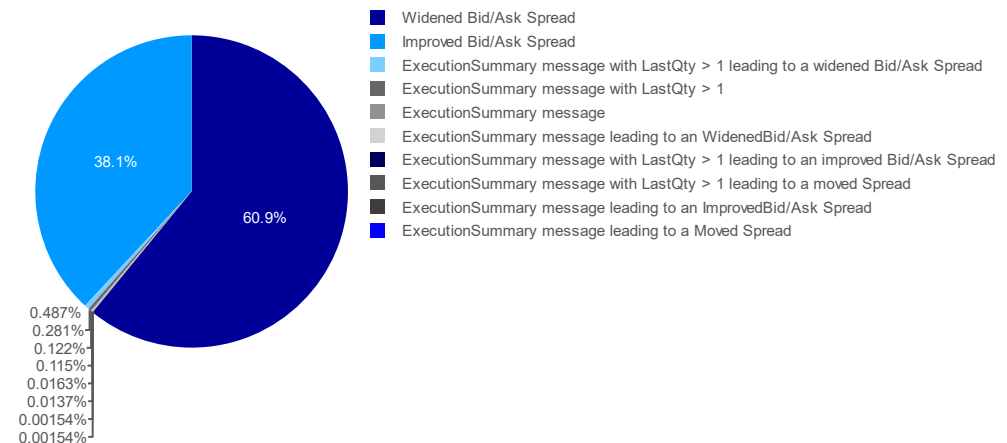


### OESX DSCP Statistics 2024-06-04

DSCP Statistic



DSCP Flags

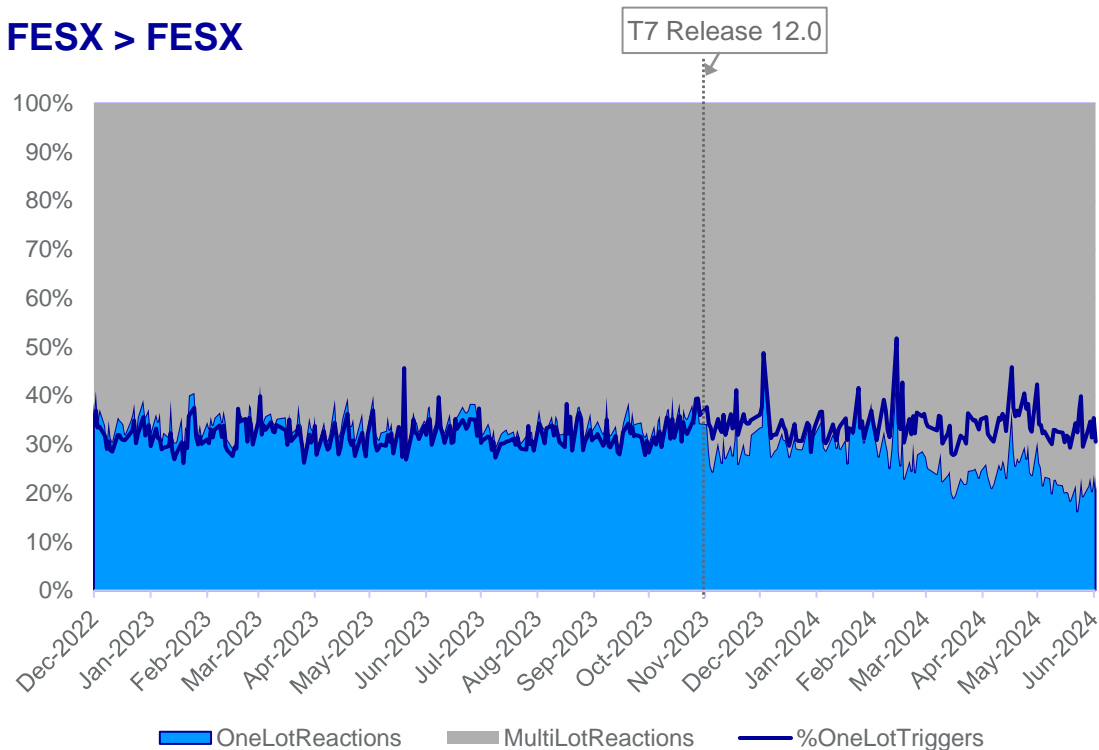


# EOBI DSCP Field

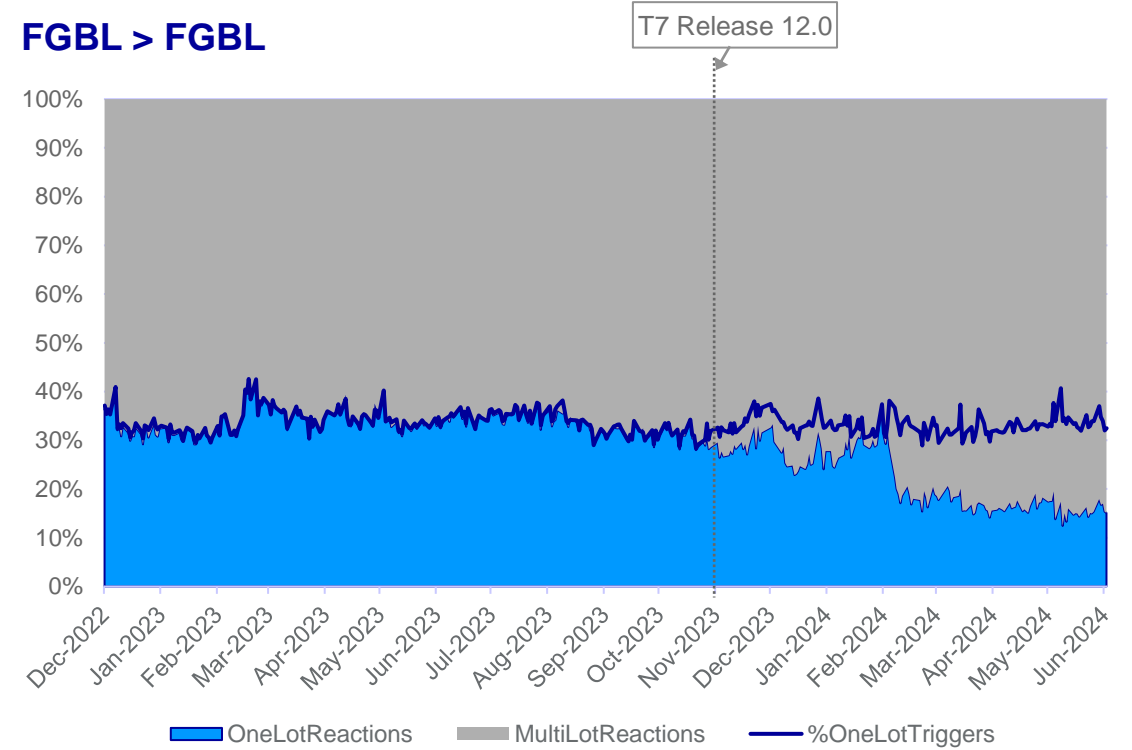
## Speculative Triggering on 1-lot Trades

With T7 Release 12.0, Eurex introduced a DSCP flag indicating whether the traded quantity was 1-lot or bigger. This led to a significant reduction in number of speculative reactions\* to 1-lot trades in several benchmark Futures products.

### FESX > FESX



### FGBL > FGBL



# Participant Reaction Time

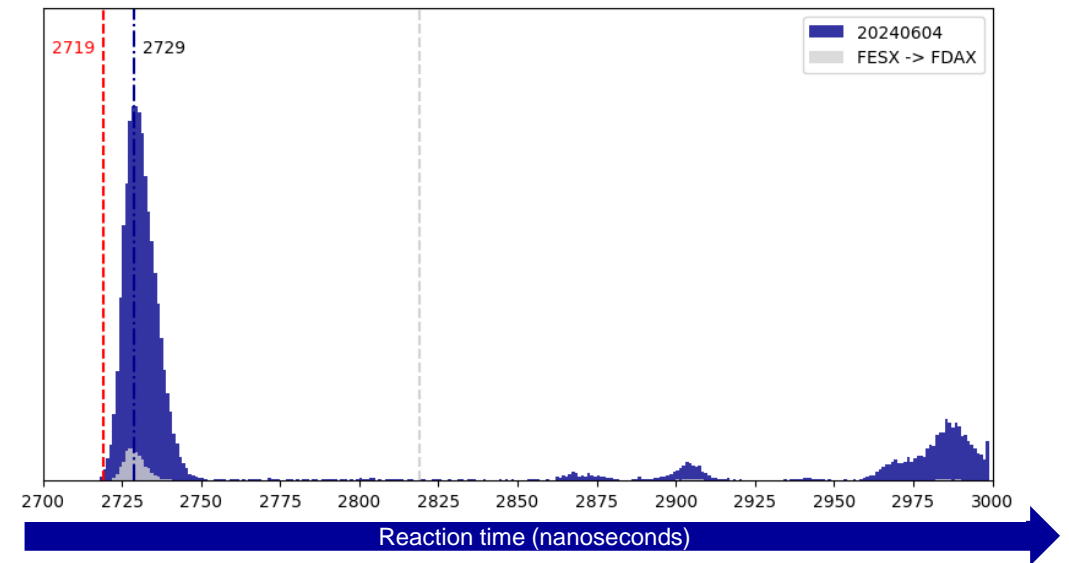
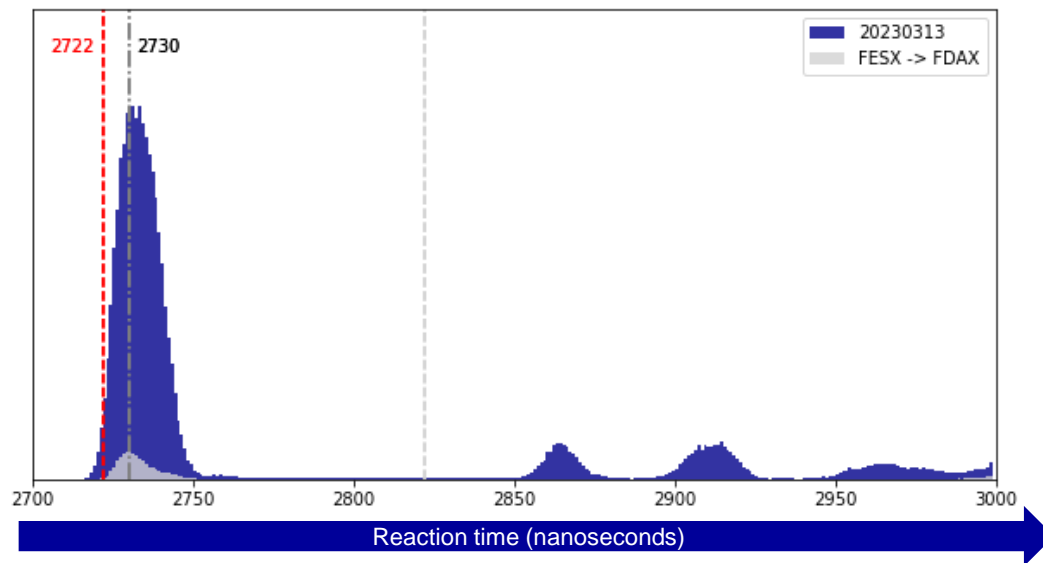
## Insights from the High Precision Timestamp file - Eurex

We define reaction time as the latency between a 'trigger' market data packet and a request that leads to an execution.

We use measurement point t\_9d for the market data packet and t\_3a for the request. We measure 2719 ns as the minimum latency between t\_9d and t\_3a determined by the T7 infrastructure.\*

The below charts show the distribution of the reaction times for all Eurex products (blue) and FESX to FDAX (grey) from 4 June 2024 (right) compared to the last figures from 13 March 2023 (left).

We observe a high level of competition (there are around 10 trading participants with reaction times < 2770 ns for most active products). The fastest participants have moved closer to each other.



\*Minimum latency changed due to upgrade of capture infrastructure and resulting shift of timestamps

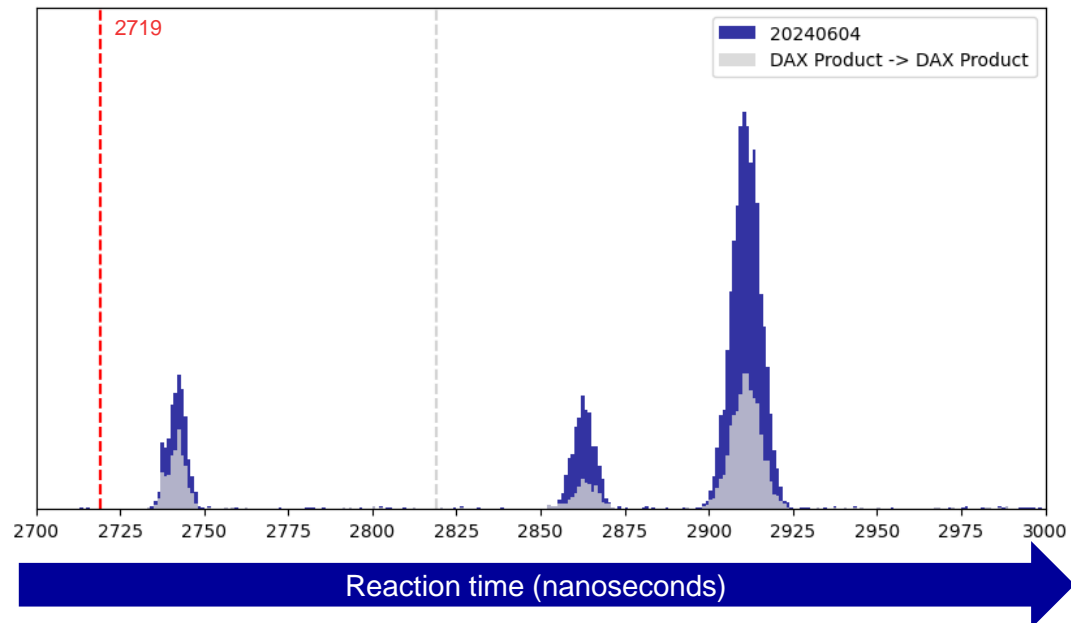
# Participant Reaction Time

## Insights from the High Precision Timestamp file - Xetra

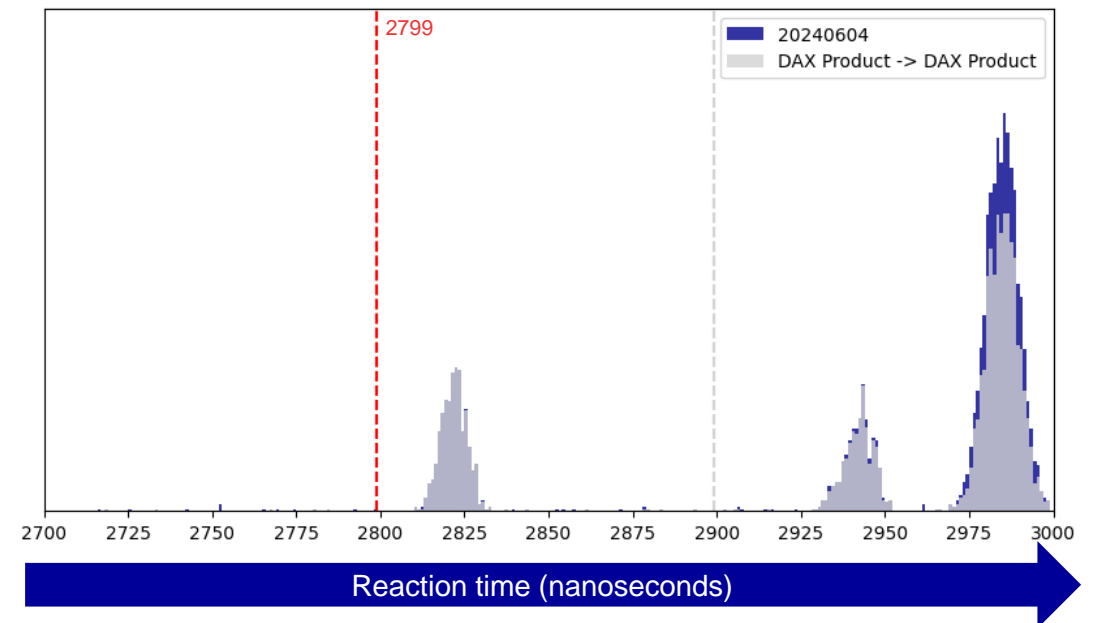
The charts below show the distribution of the reaction times for all Xetra products (blue) and DAX to DAX products (grey) from 4 June 2024. Trigger products are limited to Xetra.

The left graph shows the reaction times for Xetra products on the B-side. The right graph shows the reaction times for Xetra products on the A-side where the market data latency has increased by about 80 ns due to new setup with an additional Mid-layer switch (see [slide 30](#)).

Reaction times Xetra products B-side



Reaction times Xetra products A-side  
(New 3550T setup with Mid-layer switch)



# Capture Infrastructure Tech Refresh

On 9 and 16 March 2024, Deutsche Börse performed a tech refresh on all Access Layer capture devices.

The Arista 7130 K series were replaced by 7130 LBS series.

With the tech refresh, the t\_3a timestamps are more consistent across lines due to lower on-device and cross-device port-to-port offsets.

Additionally, better time synchronization can be achieved with built-in White Rabbit support (instead of pps via coax cables).

The tech refresh and a required software upgrade on 29 June lead to changes in the t3\_a timestamps of up to 7 ns compared to before March 2024.

The backend capture devices (responsible for t\_3d and t\_9d) are planned to be replaced in Q3/Q4 2024 and will be announced separately.

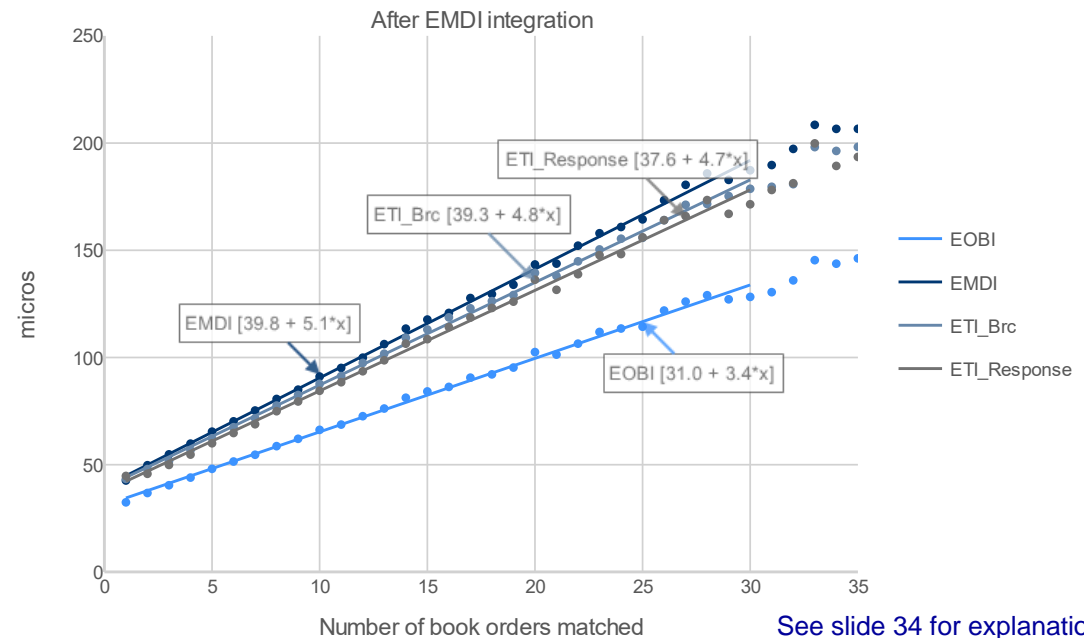
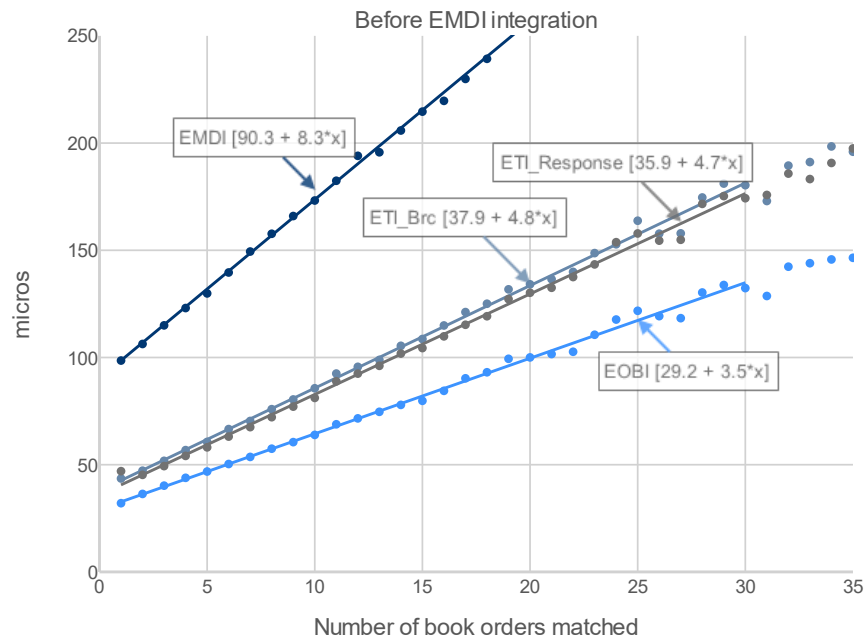
# Integration EMDI into Matching Engine

In Q2/2024, Deutsche Börse completed the stepwise consolidation of EMDI process into the Matching Engine in the T7 production environment for Eurex and Xetra. For more information, please refer to [Eurex Circular 011/24](#).

Target of the consolidation is:

- To reduce complexity by reducing the number of possible failover scenarios and
- To enable a more deterministic distribution of data via the Enhanced Market Data Interface (EMDI).

With the consolidation, EMDI became faster. However, it is now ensured that market data sent via EOBI is always faster than EMDI.



See [slide 34](#) for explanation of the graph



# 17

## Latency Analysis

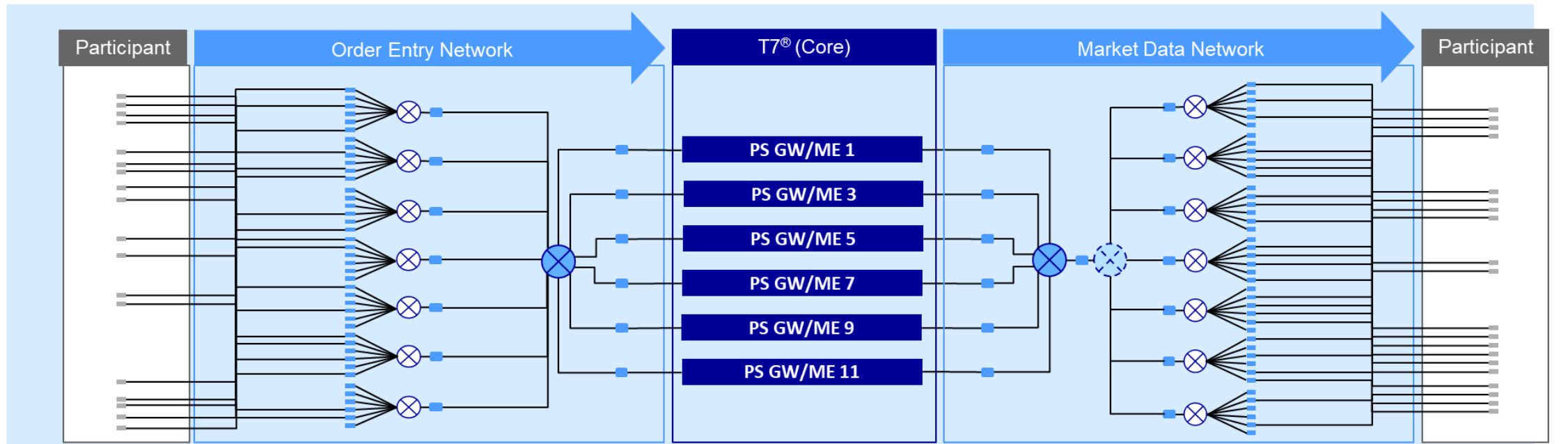


# T7<sup>®</sup> Topology

## At a glance

The below shows the Eurex B side (uneven partitions) as a schematic example of the topology of the T7 system.

Note that Xetra has only two Access Layer switches per side.



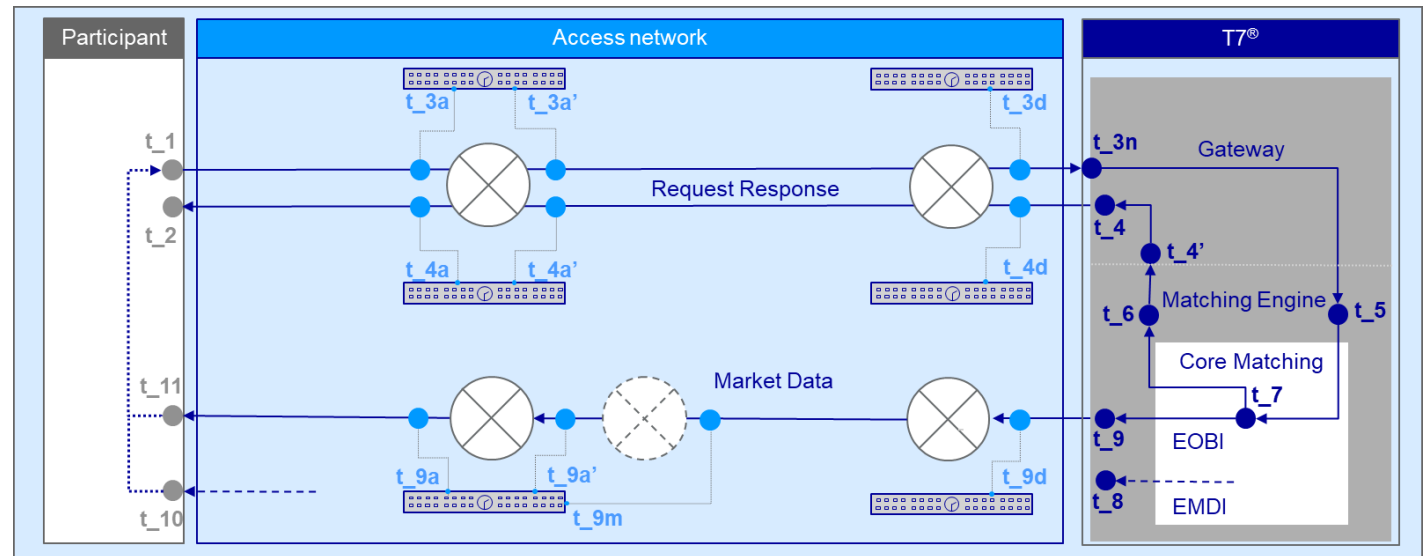
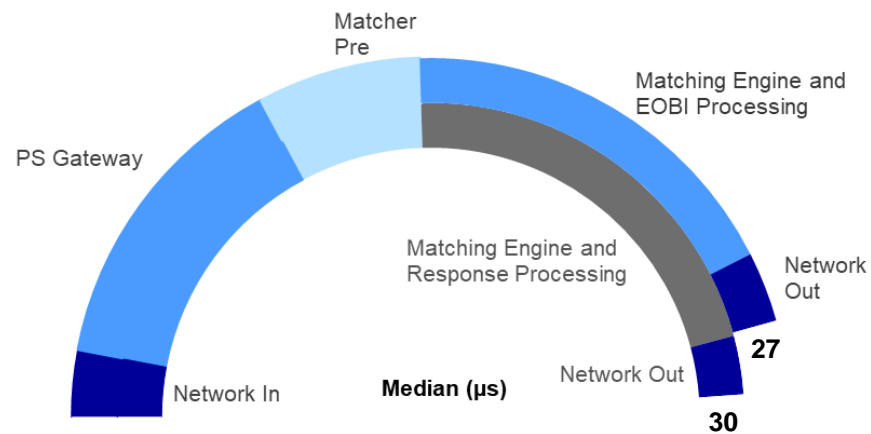
- ⊗ Cisco 3548X Distribution Layer Switch
- ⊗ Cisco 3550T Mid-Layer Switch (planned to be introduced in Q3 for Eurex)
- ⊗ Cisco 3548X Access Layer Switch
- Optical tap (capture and time measurement)

# T7<sup>®</sup> Latency

## Composition

The T7 trading system provides utmost transparency about its latency characteristics. Most of the timestamps taken at the various measurement points within T7 Core are included in each ETI response and EOBI market data. With the high precision timestamp file, we also make three network timestamps available for each EOBI market data packet ( $t_{9d}$ ) and its triggering transaction ( $t_{3a}$ ,  $t_{3d}$ ).

The latency circle shows the median latencies for the request-response/EOBI market data path for Eurex futures sent via high frequency sessions measured on 4 June 2024.



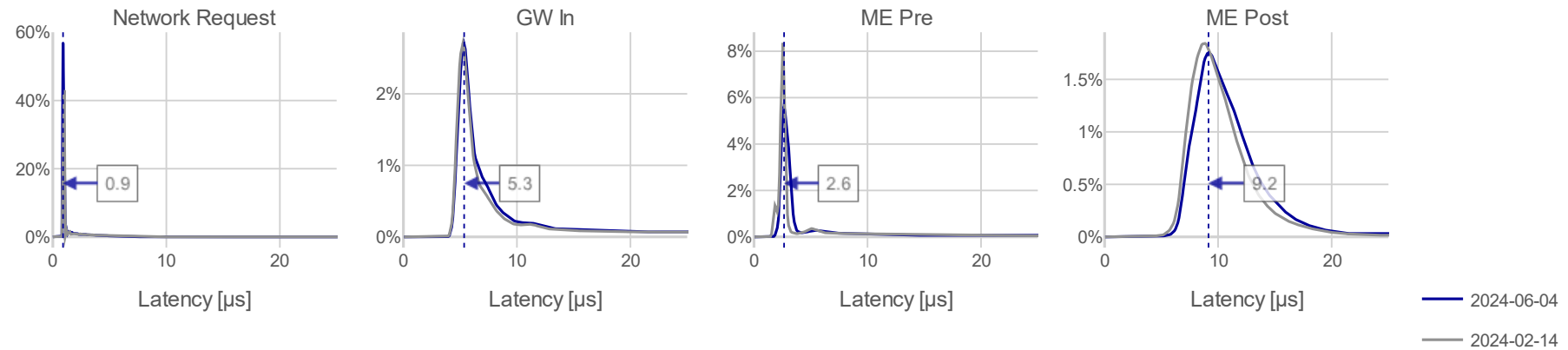
# T7<sup>®</sup> Latency

## Composition (continued)

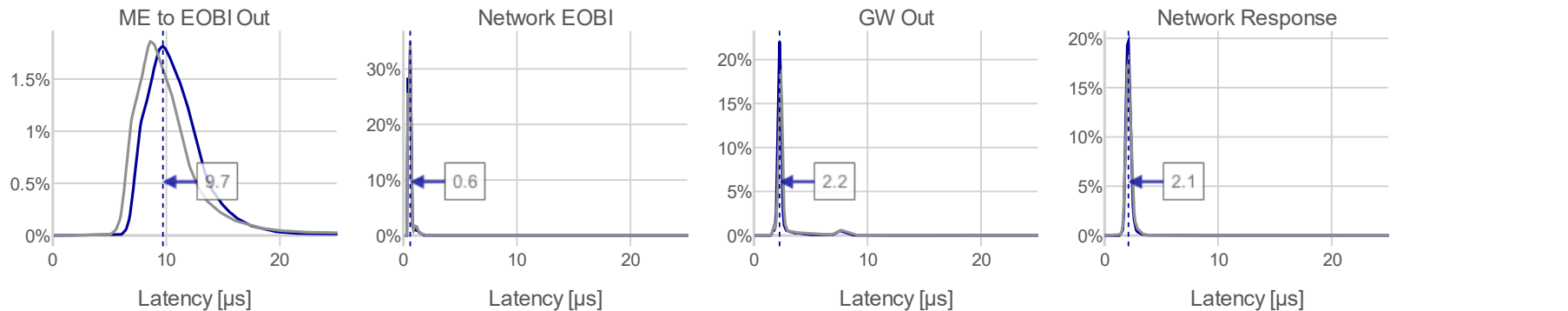
The below charts show a comparison of current latencies with the spectrum from before consolidation of EMDI into the Matching Engine which added some processing time. The data is for Eurex Futures sent from HF sessions. 'Network response/EOBI' include the TCP/UDP stack.

Futures latency compared with 2024-02-14

Request path and market data (EOBI)

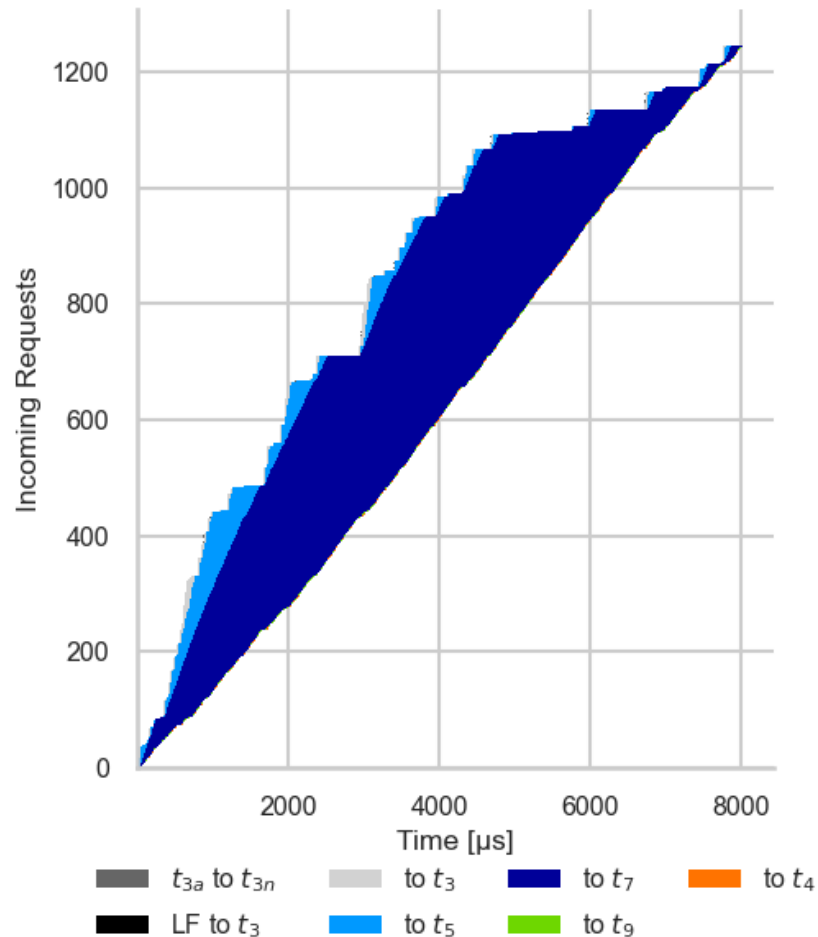


Response path



# T7<sup>®</sup> Latency

## Eurex Micro-burst dynamics



During micro-bursts, the input into the trading system may be greater than the throughput capabilities. This in turn causes queuing which results in higher latencies.

Higher latencies cause risk (i.e. it takes longer to place/pull an order).

T7 provides real-time performance insights by providing up to six timestamps with each response and key timestamps with every market data update.

The left chart shows for one example burst in FGBM on 4 June 2024 the following paths:

- Network Access Layer ( $t_{3a}$ ) to PS GWIn ( $t_{3n}$ ) to PS GW SWIn ( $t_3$ ) to
- Matching Engine in ( $t_5$ ) to
- Start matching ( $t_7$ ) to
- EOBI SendingTime ( $t_9$ ) [where available] to
- ETI SendingTime ( $t_4$ ).

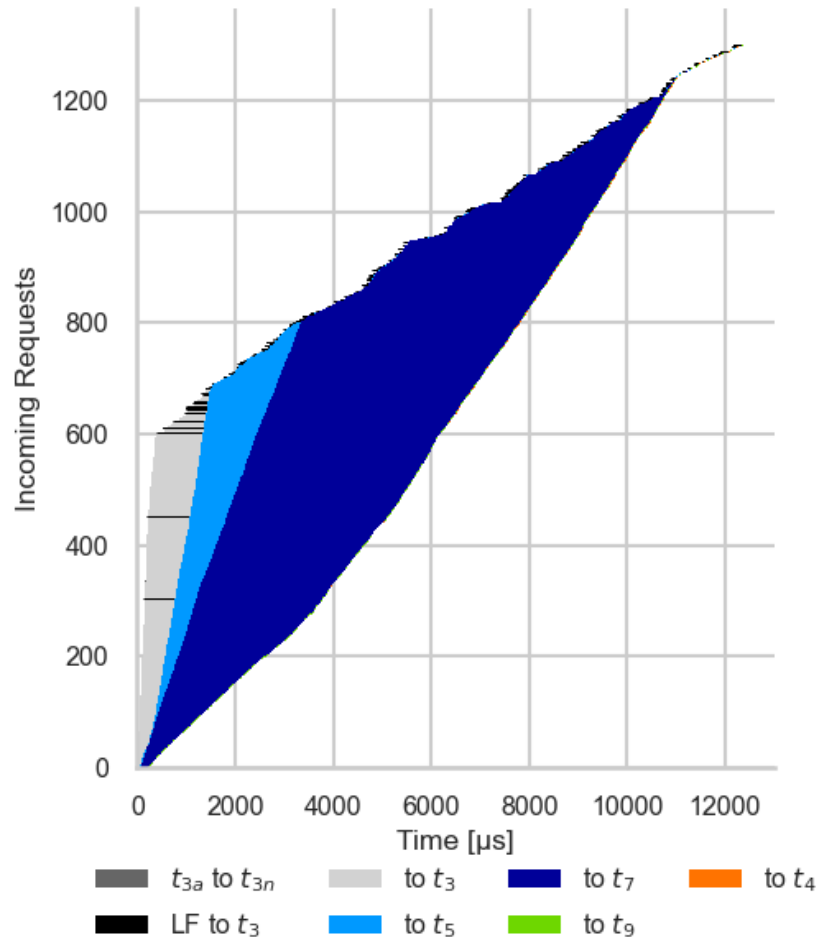
Typical throughput rates in kHz (1/ms) are 8000 at  $t_{3n}$ , ~400-500 at  $t_5$ , ~200 at  $t_7$ .

EOBI send times are usually before the gateway send time of responses.

Note that base latency for requests entered via LF gateways is ~ 32  $\mu$ s higher. As all requests are routed via PS gateways no overtaking between LF gateway and PS gateway requests is observed.

# T7<sup>®</sup> Latency

## Xetra Micro-burst dynamics



During micro-bursts, the input into the trading system may be greater than the throughput capabilities. This in turn causes queuing which results in higher latencies.

Higher latencies cause risk (i.e. it takes longer to place/pull an order).

T7 provides real-time performance insights by providing up to six timestamps with each response and key timestamps with every market data update.

The left chart shows for one example burst on Xetra partition 57 (4 June 2024) the paths:

- Network Access Layer ( $t_{3a}$ ) to PS GWIn ( $t_{3n}$ ) to PS GW SWIn ( $t_3$ ) to
- Matching Engine in ( $t_5$ ) to
- Start matching ( $t_7$ ) to
- EOBI SendingTime ( $t_9$ ) [where available] to
- ETI SendingTime ( $t_4$ ).

Typical throughput rates in kHz (1/ms) are 8000 at  $t_{3n}$ , ~400-500 at  $t_5$ , ~200 at  $t_7$ .

EOBI send times are usually before the gateway send time of responses.

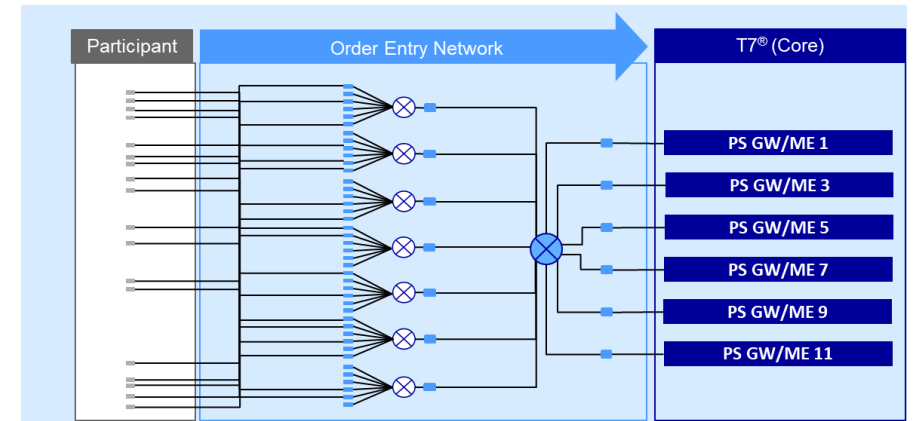
Note that base latency for requests entered via LF gateways is ~ 32  $\mu$ s higher. As all requests are routed via PS gateways no overtaking between LF gateway and PS gateway requests is observed.

# T7<sup>®</sup> Latency

## Order entry network

The order entry network in Colocation 2.0 uses a two staged hierarchical funnel in approach.

All cables are normalized to guarantee a maximum deviation of  $\pm 0.5$  m ( $\pm 2.5$  ns) between any two cross connects to the exchange. Deutsche Börse worked on reducing this deviation in 2023 as an intermediate step and plans to improve the cable length normalization in 2025 significantly (see Outlook on [slide 4](#)).



Every inbound and outbound packet on this path is captured via passive network TAPs at 3 different locations.

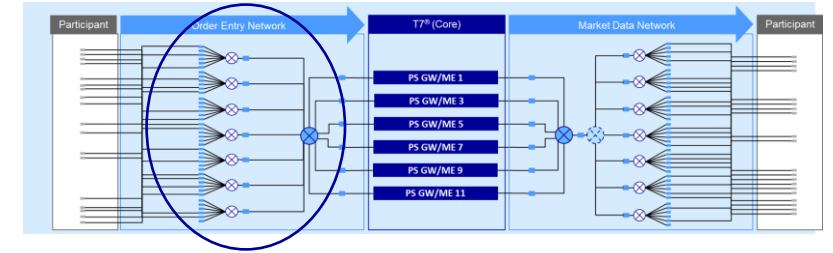
The packets are timestamped with sub-nanosecond resolution and nanosecond accuracy.

This capture infrastructure allows early detection and analysis of many kinds of technical network problems and an in-depth latency analysis on network level, like overtaking probabilities between measurement points.

The [high precision timestamp file service](#) enables participants access to timestamps t\_3a, t\_3d and t\_9d for each request leading to an EOBI market data update.

# Order Entry Network

## Latency aspects

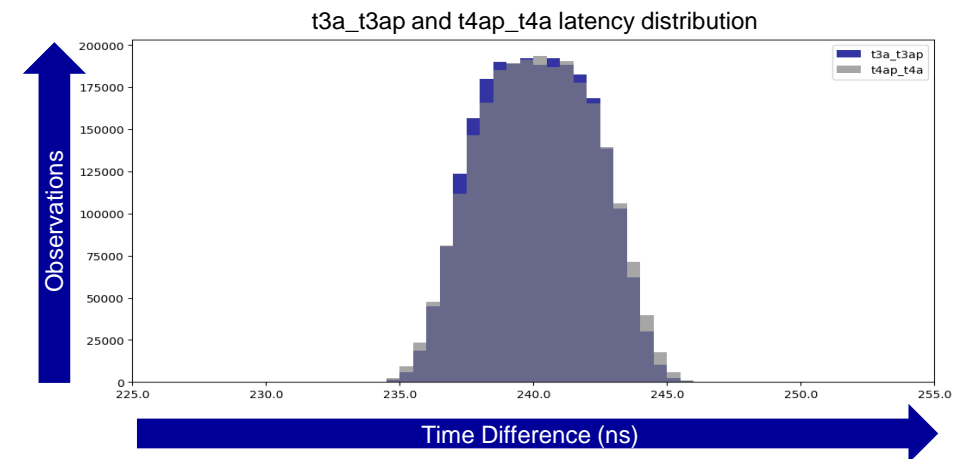
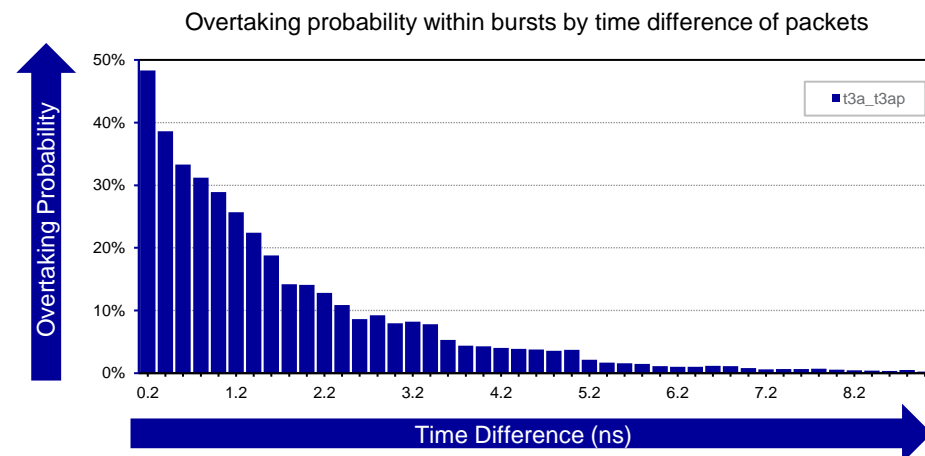
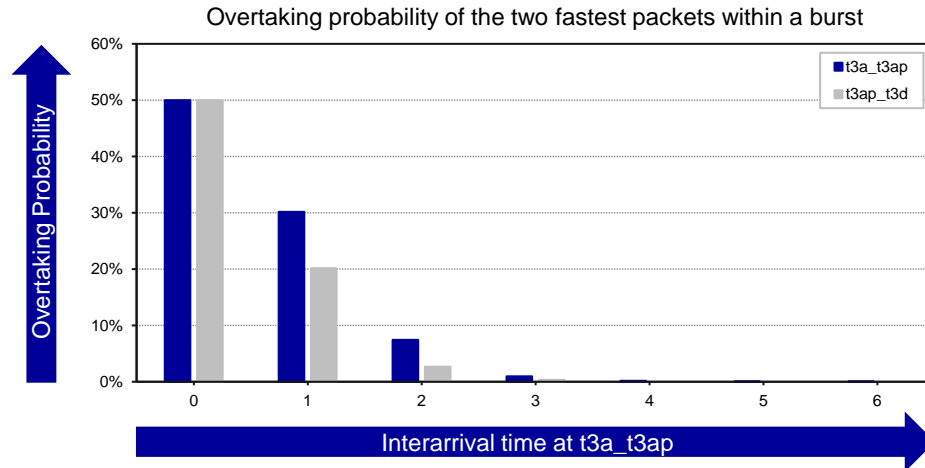


The chart on the left shows the overtaking probability between a first and a second message within a burst on one switch. The overtaking drops sharply and there is no observed overtaking beyond 3 nanoseconds delta between messages.

With improved timestamp accuracy and introduction of sub-nanoseconds on capture devices, more granular analyses could be performed. The bottom left graph shows the overtaking probability of all messages within a burst based on their interarrival time.

The base latency and latency jitter is identical for all Access Layer switches within the measurement precision (bottom right graph).

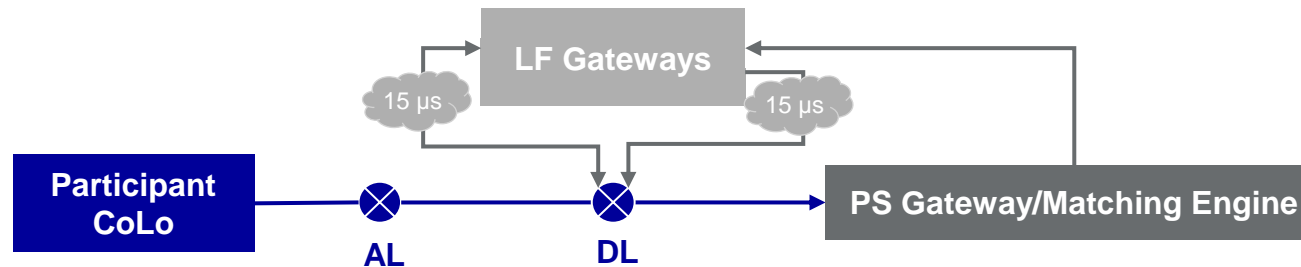
Data is taken from the Eurex A side.





# PS Gateway versus LF Gateways

## Latency comparison



- Using 10 Gbit/s cross connects and access via PS gateways provides the fastest way for order and quote management in T7.
- LF gateways on the other hand allow access to all partitions of a market via a single session.
- Some markets (e.g. XMAL) are using LF gateways only.
- Most markets have combined PS gateway/ Matching Engines in place for which all Matching Engine bound requests sent to LF gateways are routed via PS gateways.
- The base latency of the path to the PS gateway is around 32  $\mu$ s higher for LF gateways when compared to directly accessing the PS gateways. Note also that requests that have to cross sides between LF and PS gateways take another  $\sim$ 45  $\mu$ s longer to reach the Matching Engine.

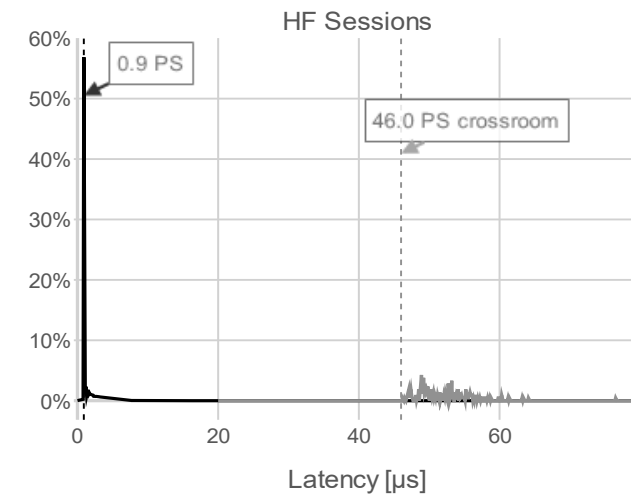
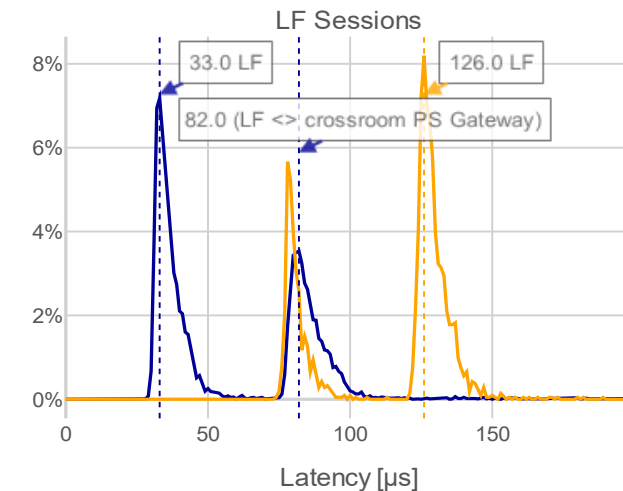
# Order Entry Latency

## Comparison of access types

The table below shows a comparison of different access options to the T7 system. All times given are in microseconds.

Network timestamps (t\_3a) are synchronized using white rabbit. The time synchron quality between these timestamps is thus ~1 ns. Other T7 timestamps are subject to jitter of up to ±50 ns.

Gateway type	Same side line	Same side partition	t_3a to t_3n uncongested latency
PS	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	1
	<input type="checkbox"/>	<input checked="" type="checkbox"/>	46
LF	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	33
	<input checked="" type="checkbox"/>	<input type="checkbox"/>	82
	<input type="checkbox"/>	<input checked="" type="checkbox"/>	79
	<input type="checkbox"/>	<input type="checkbox"/>	126



LF Sessions  
 — LF & member: Same room  
 — LF & member: Cross room

HF Sessions  
 — Same room  
 — Cross room

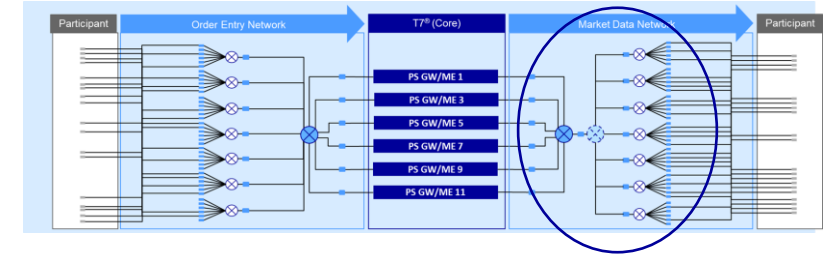
**27**

**Market Data**

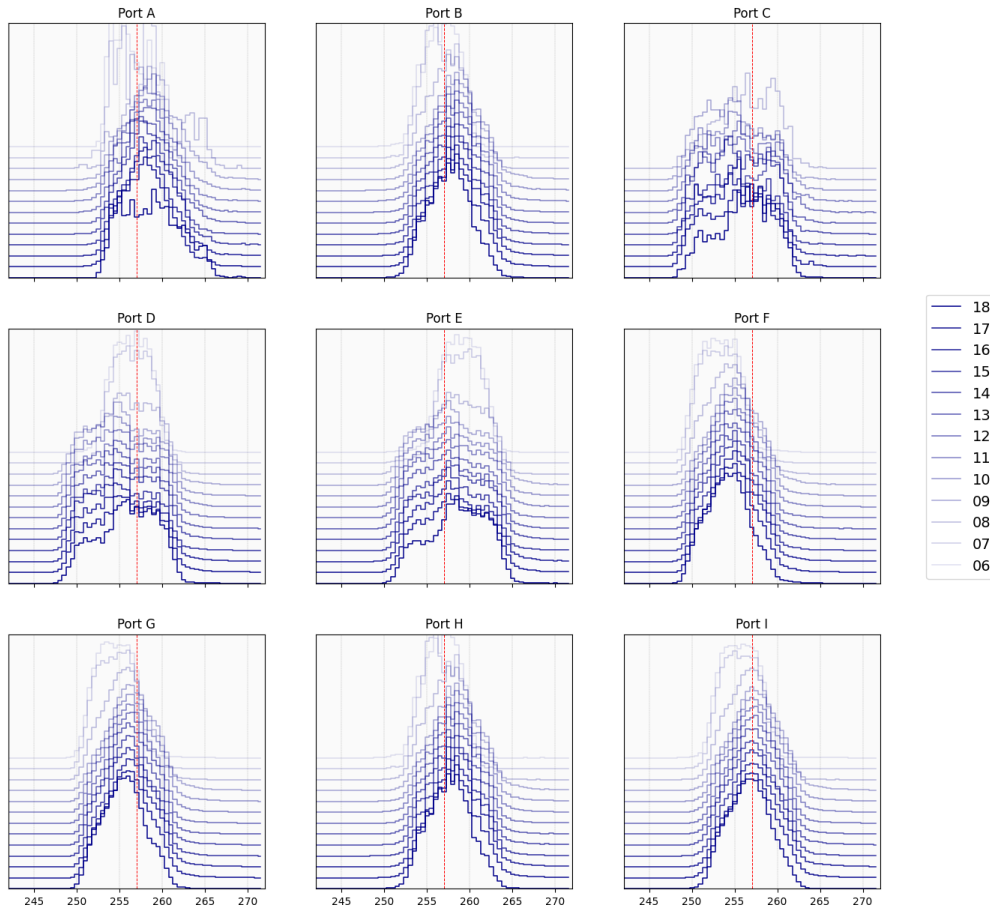


# Market Data Network

## Latency aspects



t\_9a' to t\_9a latency (ns)



The market data network has a funnel in – fan out topology. It funnels in data from different market data disseminators (on the Distribution Layer switch) and fans it out via multiple Access Layer switches.

We took extra care to establish a balanced and deterministic network. Additionally, static forwarding is configured to ensure equal multicast load towards Access Layer switches.

We observe a semi static latency difference of up to 10 ns (comparing t\_9a of two different ports on the same switch) due to internal multicast processing of the switch.

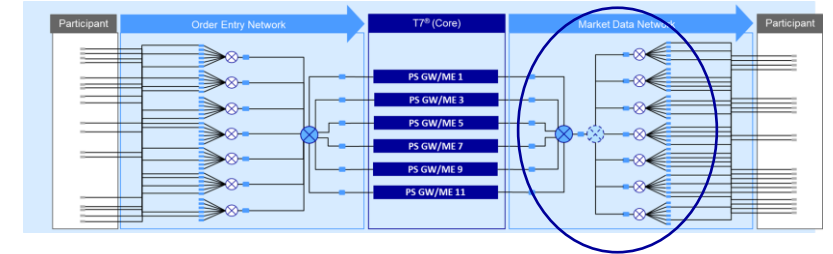
The Cisco Nexus 3548-X switch divides the on-board Ethernet-Ports into three buffer areas. Serial replication for multicast takes place within these areas. In terms of latency this means any port within this area could be faster than other ports, which is an artifact of variable internal pointers to packet queues. The starting port for replication of a packet in a buffer area is identified by an internal pointer.

The figure on the left shows the latency distributions (t\_9a' to t\_9a) for nine ports belonging to one buffer area of one Access Layer switch.

Each distribution represents one minute of the specified hour of the day.

# Market Data Network

## Comparison Cisco 3548X and 3550T

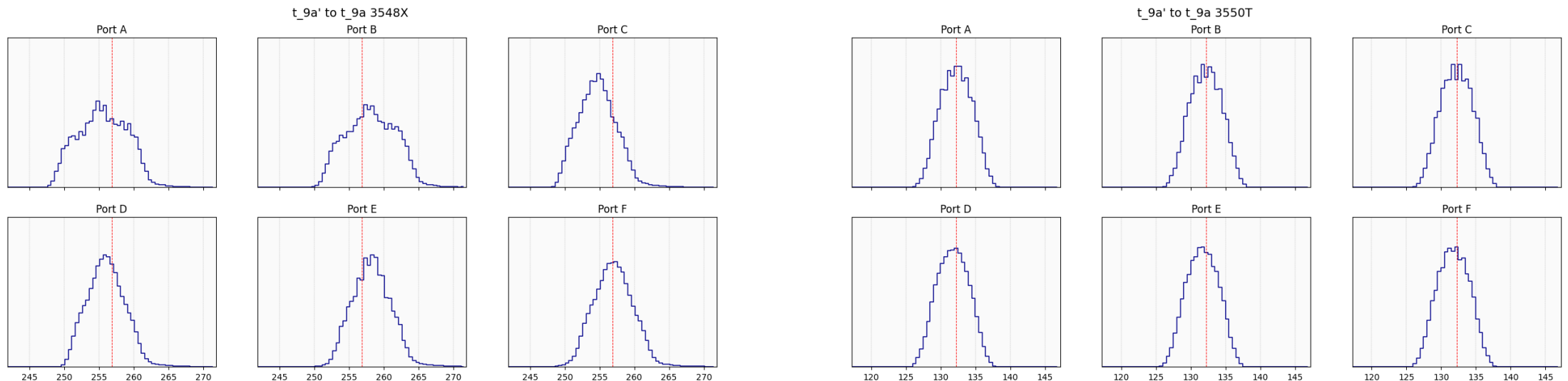


The two graphs below show a comparison of processing times between the Cisco 3548X (Eurex) and Cisco 3550T (Xetra) based on production data. The red lines indicate the median processing time over all ports of the respective switch.

The Cisco 3550T shows the following advantages compared to the Cisco 3548X:

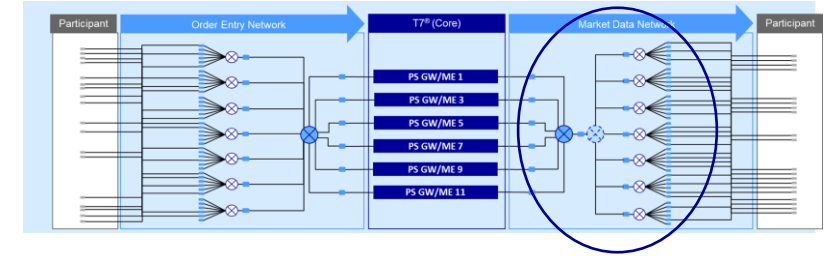
- Processing times are more consistent across ports
- The median latency of each port differs from any other port by less than 1 ns. Lab results confirm that this is true for each packet.
- The latency distribution is stable over time as opposite to the Cisco 3548X (see previous slide).
- It does not have internal pointers or buffer areas that dynamically influence latency.

Deutsche Börse plans to replace all Colo 2.0 market data switches with Cisco 3550T devices. For more information on timeline see Outlook [slide 4](#).



# Market Data Network

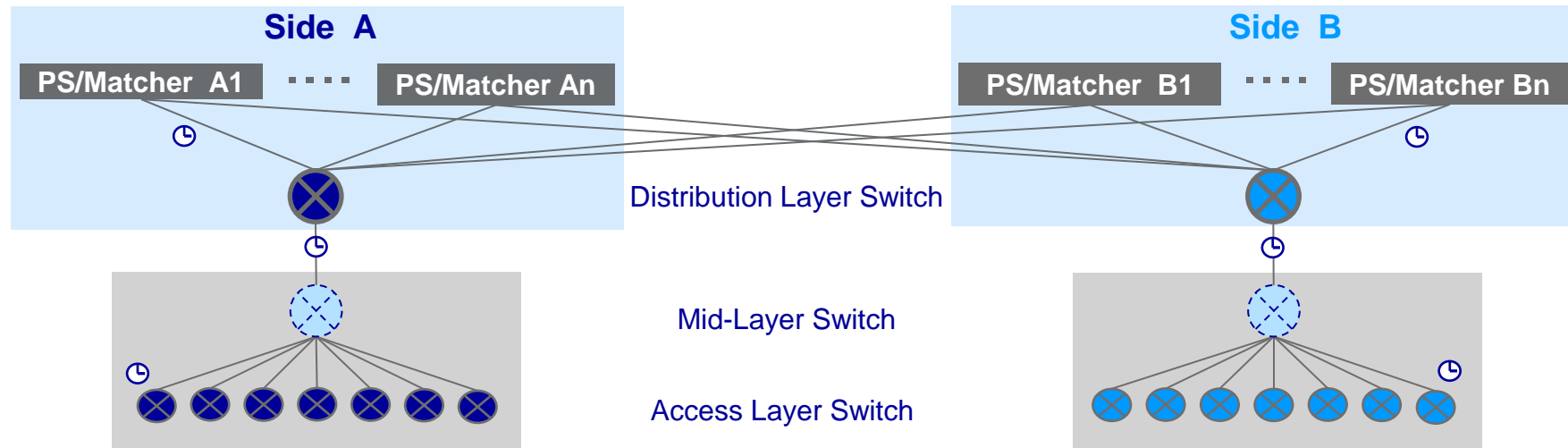
## Introduction of Mid-Layer Switch



An additional switch (Cisco 3550T) will be added as Mid-Layer switch between the Access Layer and Distribution Layer switches for Eurex and Xetra.

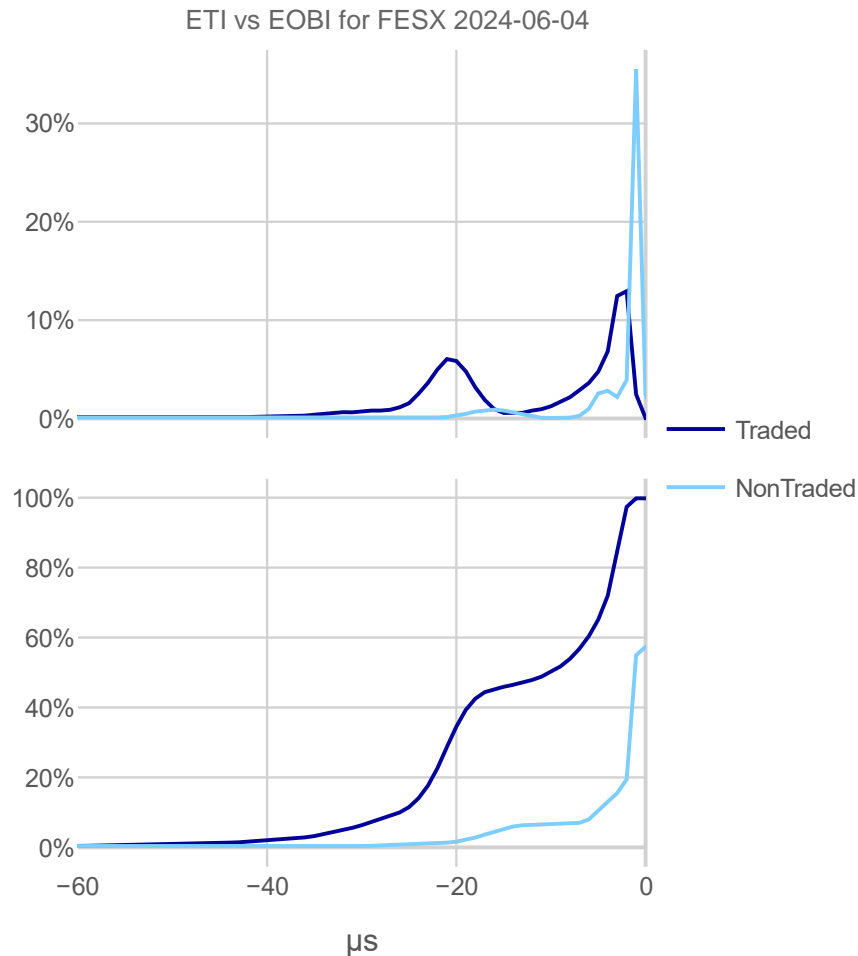
The purpose of the Mid-Layer switch is to:

- Accurately equalize cable to Access Layer switches, this is achieved by physical proximity of Mid-Layer and Access Layer switches.
- Serialize the connection from T7 backend to Access Layer by having just one connection between Distribution Layer and Mid-Layer.
- Make use of the positive characteristics of the Cisco 3550T (see previous slide) to achieve more equal distribution of market data to the Access Layer switches.



# Trading System Dynamics

## Latency characteristics of EOBI versus ETI for Futures



T7 is designed to publish order book updates first on its public data feed.

The diagram shows the time difference distribution between public and private data in microseconds (EOBI first datagram vs ETI responses,  $t_9 - t_4$ ).

The data is a production sample from 4 June 2024.

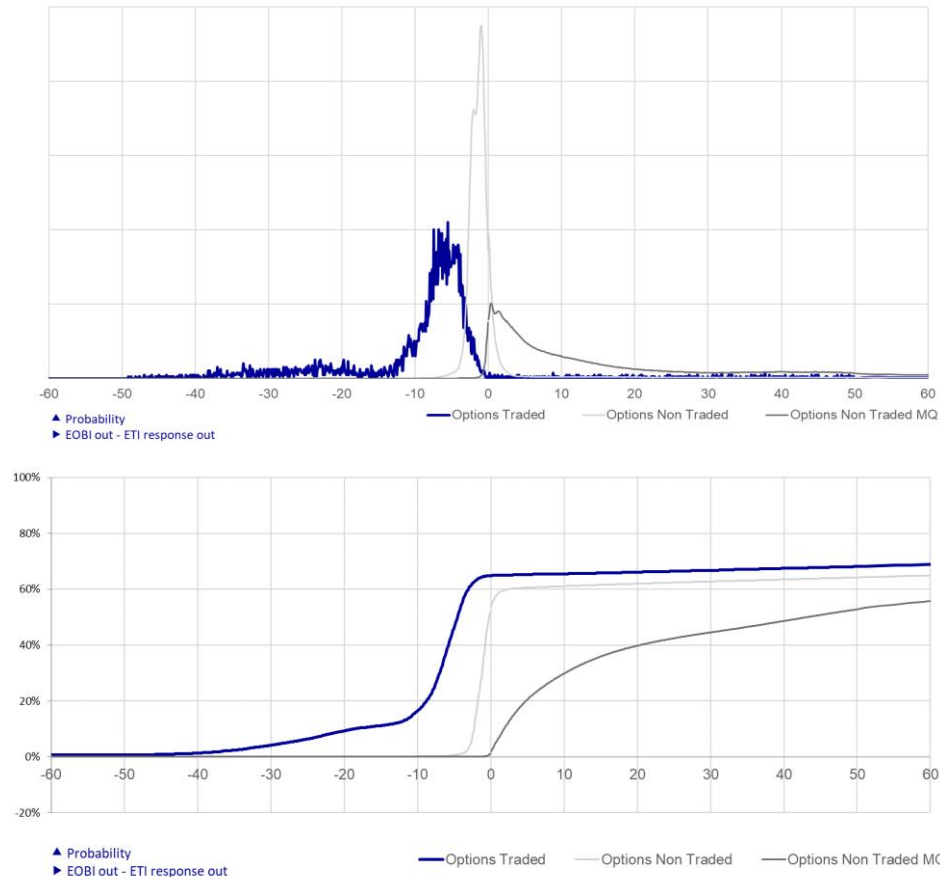
EOBI market data is in median 1.5  $\mu\text{s}$  faster than the ETI response for order book updates and 10  $\mu\text{s}$  faster for executions.

The first EOBI datagram is faster in more than 99.8 percent of the cases compared to the ETI response, and also the first passive ETI book order notification for simple transactions.

# Trading System Dynamics

## Latency characteristics of EOBI versus ETI for Options

ETI vs EOBI for OESX 2024-06-04



The data is a production sample from 4 June 2024 for OESX Options.

We distinguish between orders leading to a trade (quotes cannot match aggressively in OESX due to PLP), orders and single quote updates and mass quotes updates.

Trades are received first on EOBI in around 65% of the cases with a median latency advantage of 4.5  $\mu$ s.

There are two main reason for EOBI delays:

The transaction is delayed by preceding messages (queues).

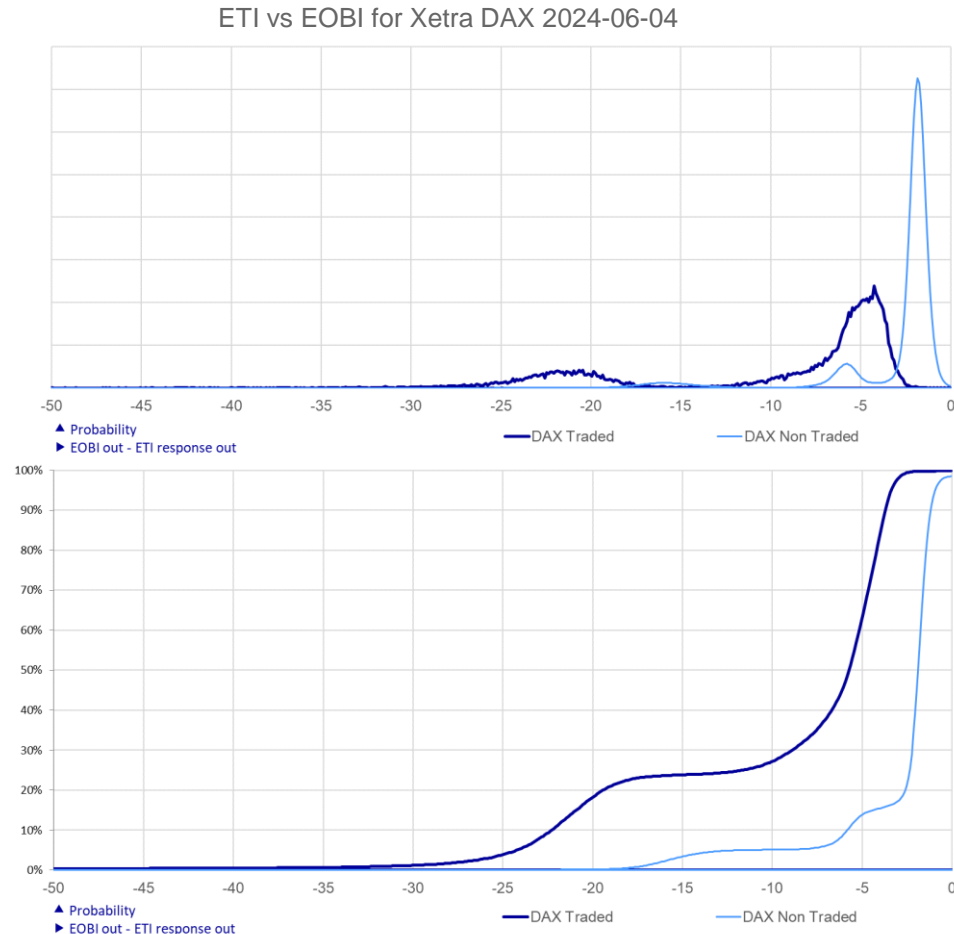
The transaction causes a market maker protection with many quote deletions. Additionally, the integration of EMDI into the Matching Engine process had an impact on the EOBI processing time.

The latency profile for mass quotes is dominated by larger mass quotes, where the EOBI publisher has to broadcast each quote, leading to longer delays and queues in the EOBI path, while the ETI path only deals with a simple mass quote ack.



# Trading System Dynamics

## Latency characteristics of EOBI versus ETI for Xetra



The diagram shows the time difference distribution between public and private data in microseconds for XETRA DAX products (EOBI first datagram vs ETI responses,  $t_9 - t_4$ ).

The data is a production sample from 4 June 2024.

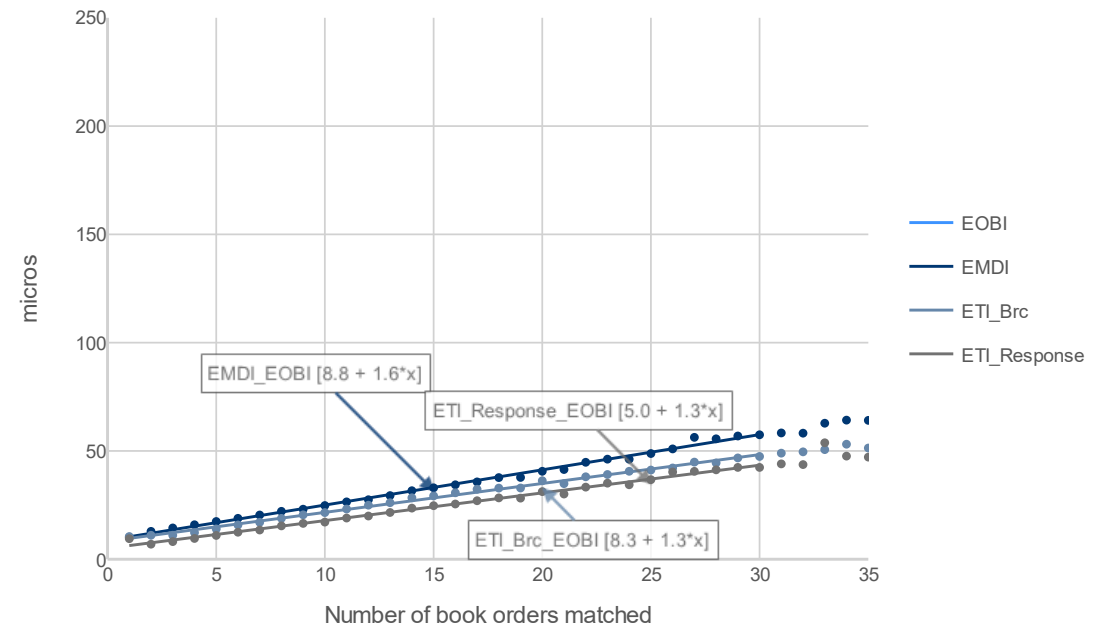
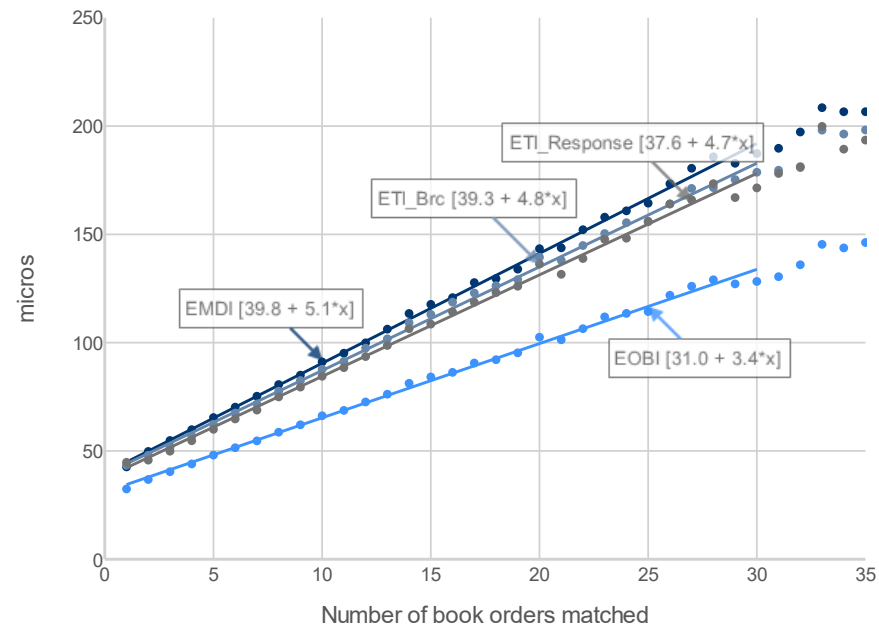
The latency distribution is similar to the Eurex futures, trades are received on EOBI in median 5.7  $\mu$ s faster, whereas single orderbook updates are usually 1.9  $\mu$ s faster on EOBI.

# Trading System Dynamics

## Latency characteristics of ETI versus EOBI versus EMDI

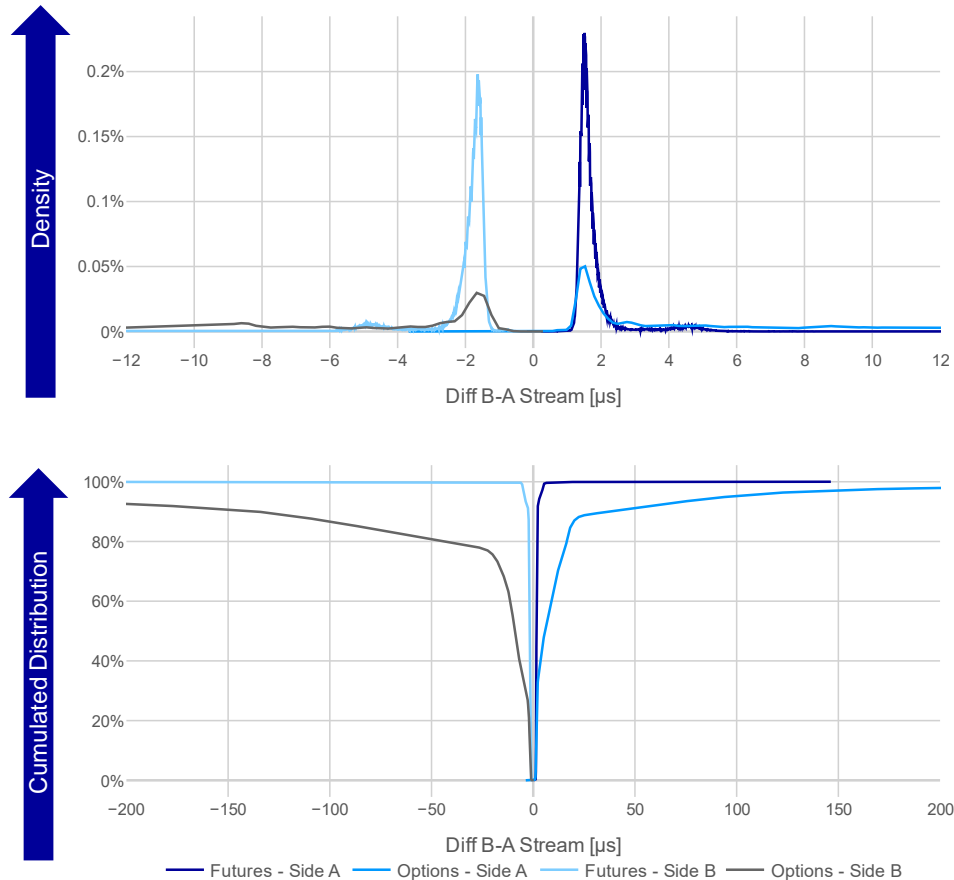
The diagrams below display the dependency of the median latency on the complexity of a trade for ETI (the right diagram is normalized to EOBI) ( $t_4 - t_7$ ), EMDI ( $t_8 - t_7$ ) and EOBI ( $t_9 - t_7$ ). Note that for ETI we display the gateway sending time of the first passive notification and for EOBI the sending time of the UDP datagram containing the Execution Summary message.

The difference between public and private data has slightly been narrowed with the integration of EMDI into Matching Engine process. However, the 'public data first' principle is still being ensured. In over 99% of all trades, we disseminate order book data on EOBI first (also for larger trades).



# Trading System Dynamics

## EOBI latency difference of primary and secondary feed



For products assigned to even partitions, market data is published first on the A and then on the B stream. For products assigned to odd partitions market data is published first on the B and then on the A stream.

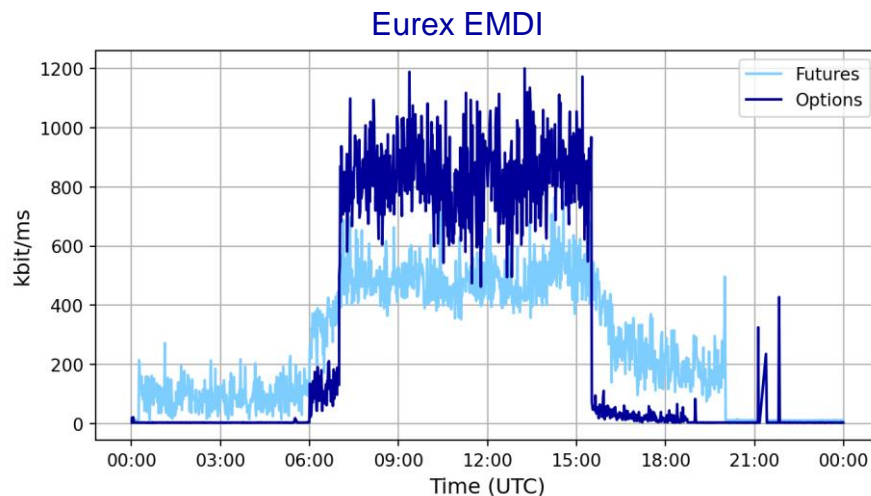
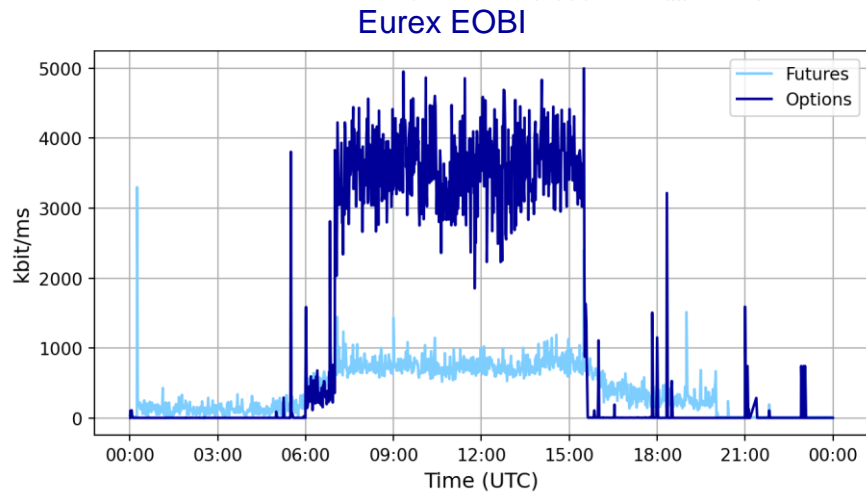
The partition ID / product ID is contained in the UDP datagram header of the order book incremental messages and can be used for filtering on UDP datagram level for EMDI / EOBI.

Furthermore, a UDP datagram on the T7 EMDI / EOBI order book delta or snapshot channel contains data of exactly one product (e.g. EURO STOXX 50® Index Futures).

The data for both primary and secondary streams is disseminated by the same server using two ports connected to the two sides of the network. The process sends the data first on the primary interface. After all datagrams of a transaction are sent it starts sending on the secondary interface.

The median latency difference between the A and the B EOBI incremental feed is about 2 µs for futures and Cash products. For options the median is slightly higher and there are far more outliers (i.e. much slower secondary feed). The reason is that since data is published on the secondary feed only after all datagrams of a transaction are sent on the primary feed, the latency difference depends on the complexity of the transaction, i.e. a mass quote with 200 quote updates will lead to a higher delay than a single order entry.

# Eurex: Market data volume



Each data point equals the maximum bandwidth produced on a 1 millisecond scale by the incremental B stream in Mbit/ms.

The provided data shows one data point per minute for 4 June 2024 – a busy trading day.

Enhanced Order Book Interface (EOBI) peak volume is significantly higher than price level aggregated data volume EMDI. EOBI market data is therefore currently only available to trading participants using 10 Gbit/s connections.

The EOBI for options incremental data stream peaks around 5 Gbit/s on millisecond level, while the futures stream peaks at 3.2 Gbit/s.

Participants that want to receive data for Eurex Exchange's products on EMDI with less than 1 ms queuing delays need to use a connection with a bandwidth of more than 1 Gbit/s (options) or 700 Mbit/s (futures) respectively. Trading participants are advised to use two 10 Gbit/s connections (one for each market data stream) in Co-Location to receive market data.

# Xetra: Market data volume

Each data point equals the maximum bandwidth produced on a 1 millisecond scale by the incremental B stream in Mbit per ms.

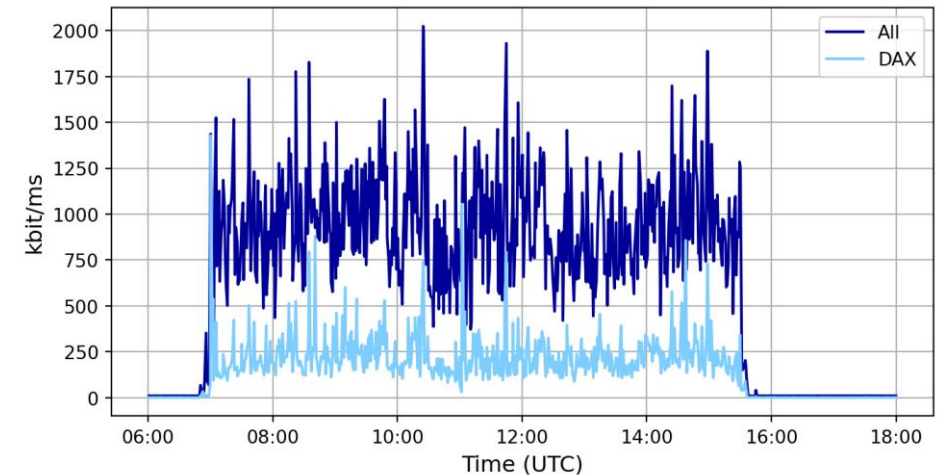
The provided data shows one data point per minute for 4 June 2024.

Enhanced Order Book Interface market data is currently only available to trading participants using 10 Gbit/s connections.

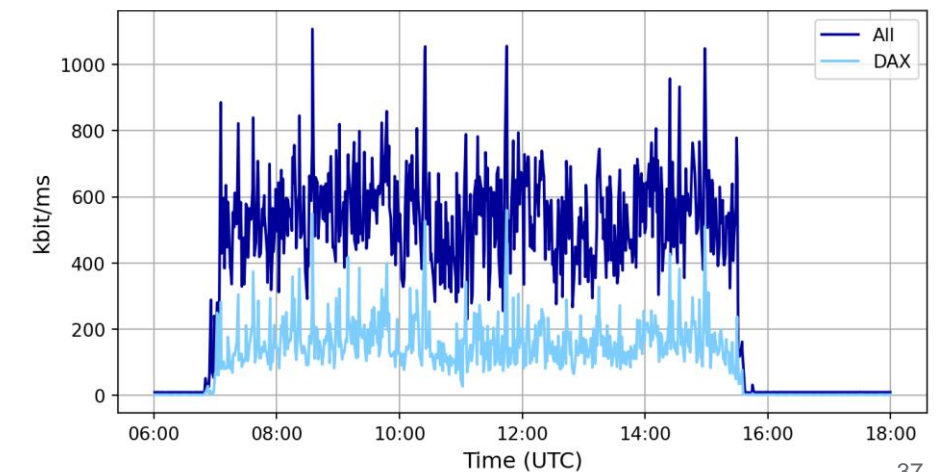
Trading participants are advised to take two cross connects (one for each market data stream) in Co-Location to receive market data.

Participants that want to receive EMDI data with less than 1 ms queuing delays need to use a connection with a bandwidth of more than 1 Gbit/s (All products) or 500 Mbit/s (DAX® equities only).

Xetra EOBI

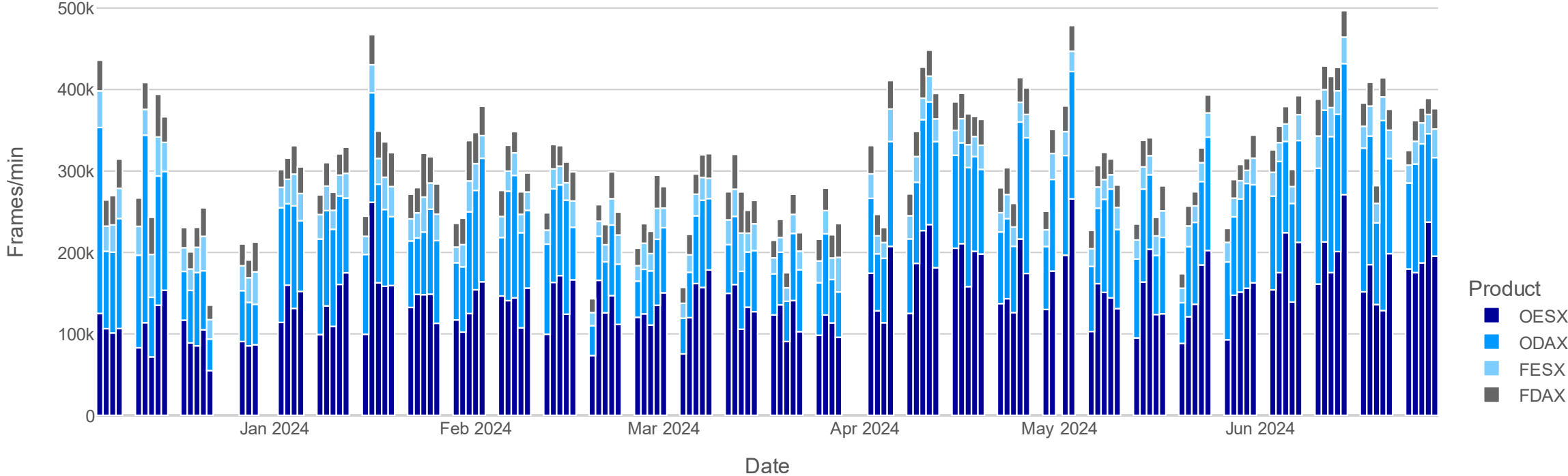


Xetra EMDI



# Maximum Frames per Minute - EOBI

Maximum EOBI Frames per Minute



**39**

**What you need to be fast**



# What you need to be fast...

## A few recommendations to achieve the low latency

Use the Equinix Co-Location facility to be close to Deutsche Börse T7.

Use state-of-the-art switches (if any) and only have at most one hop between the exchange network and your server. Alternatively, use hardware solutions to connect directly without hops (e.g. FPGA).

Use good network interface cards and TCP/IP acceleration, e.g. a kernel-by-pass library.

Use at least two 10 Gbit/s (cross-) connections in Co-Location for EOBI market data and two 10 Gbit/s connections for T7 ETI.

Use HF sessions to connect to PS gateways and make sure you use the cross connect on the same side as the gateway you are connecting to (compare time-to-live values in the IP header in the responses from both sides).

Measure and analyze your own timestamps (e.g. the reaction time as recommended on the next slide).

Use state of the art time synchronization, i.e. GPS clocks and a high-quality time distribution. The PTP signal you can get from us has a quality of  $\pm 50$  ns. For our network timestamps we use the White Rabbit protocol and PPS breakouts. You can connect to our white rabbit time service providing you a time synchronization quality of 1-2 ns max, see <https://www.deutsche-boerse.com/dbg-en/products-services/ps-technology/ps-connectivity-services/ps-connectivity-services-time-services>.

We provide highly accurate network timestamps of all orders leading to a market data update via the high precision timestamp file service, see <https://www.mds.deutsche-boerse.com/mds-en/analytics/high-precision-timestamps>.

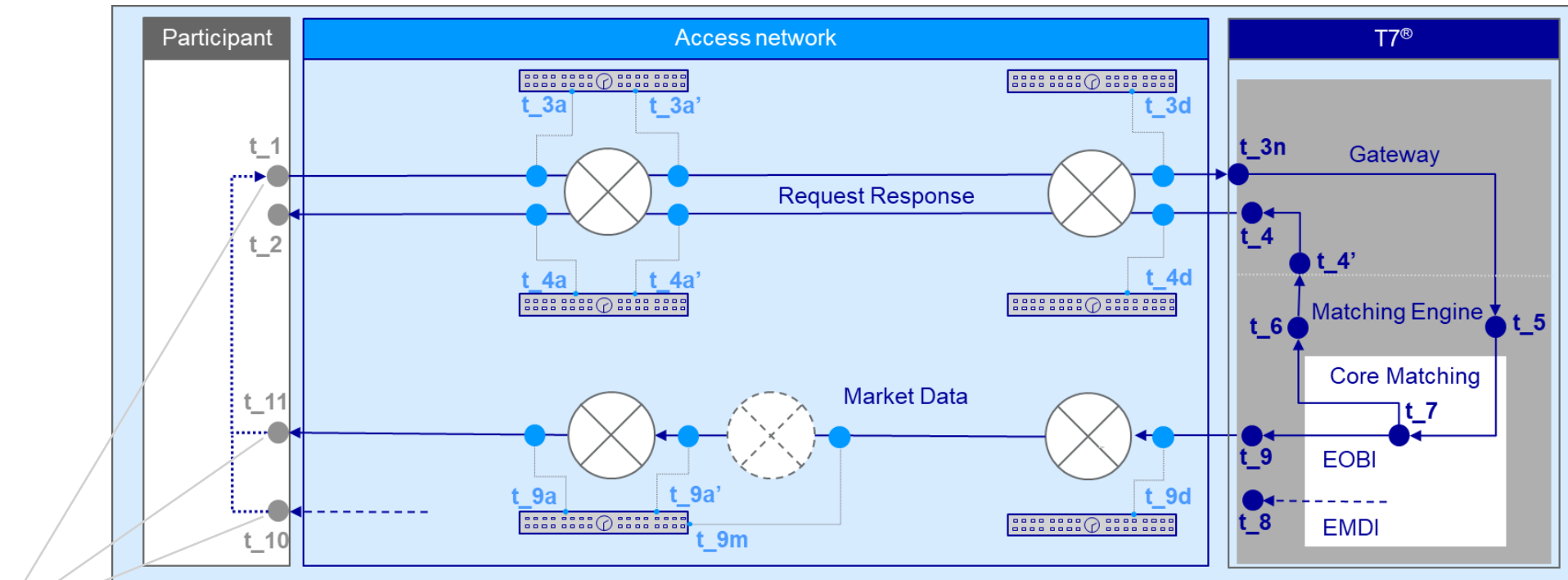
Use the EOBI Execution Summary for fast trading decisions and position keeping (passive executions). For a consistent order book, all incremental updates following the Execution Summary should always be processed. For fastest decisions evaluate the market data classification based on the DSCP flags in the IPv4 header of EOBI market data packets.

Trade notifications need to be processed to create safety. We recommend to use either a low-frequency ETI session or a FIX trade capture drop copy to confirm the fast execution information provided by the execution reports via high-frequency sessions.



# What you need to be fast...

## Participant reaction time measurement



Measure the time between market data reception ( $t_{10}/t_{11}$ ) and your reaction ( $t_1$ ) and align with timestamps from the high precision timestamp file

$t_{3a}$ ,  $t_{3d}$  and  $t_{9d}$  are available via the high precision timestamp file service, see <https://www.mds.deutsche-boerse.com/mds-en/analytics/high-precision-timestamps>  
Take our white rabbit signal to compare your timestamps with ours with ns accuracy  
<https://www.deutsche-boerse.com/dbg-en/products-services/ps-technology/ps-connectivity-services/ps-connectivity-services-time-services#tab-1556860-1556864>

**42**

**T7<sup>®</sup> Overview**



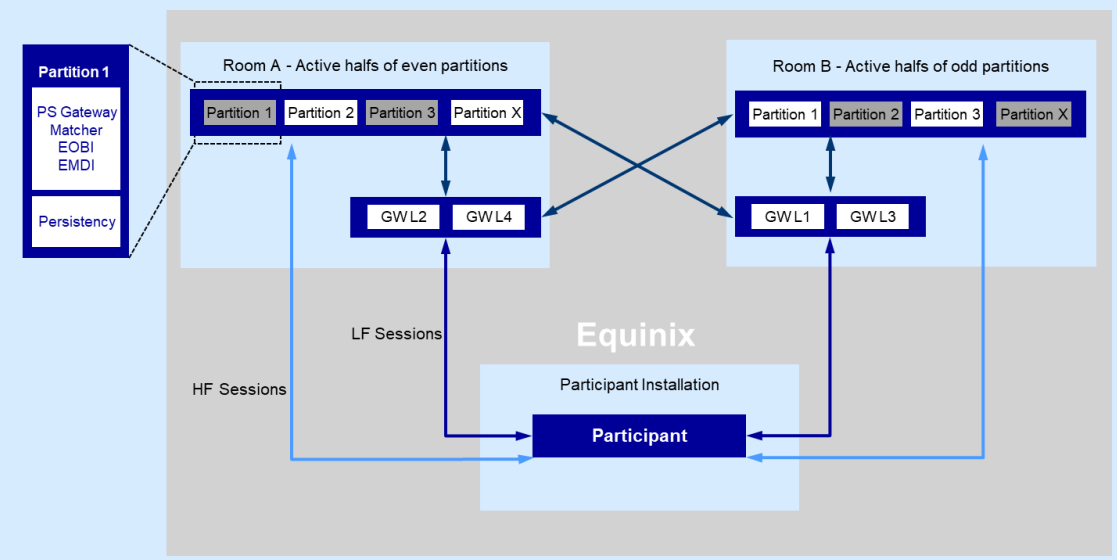
# T7<sup>®</sup> Architecture

## Overview

T7<sup>®</sup> architecture developed by Deutsche Börse:

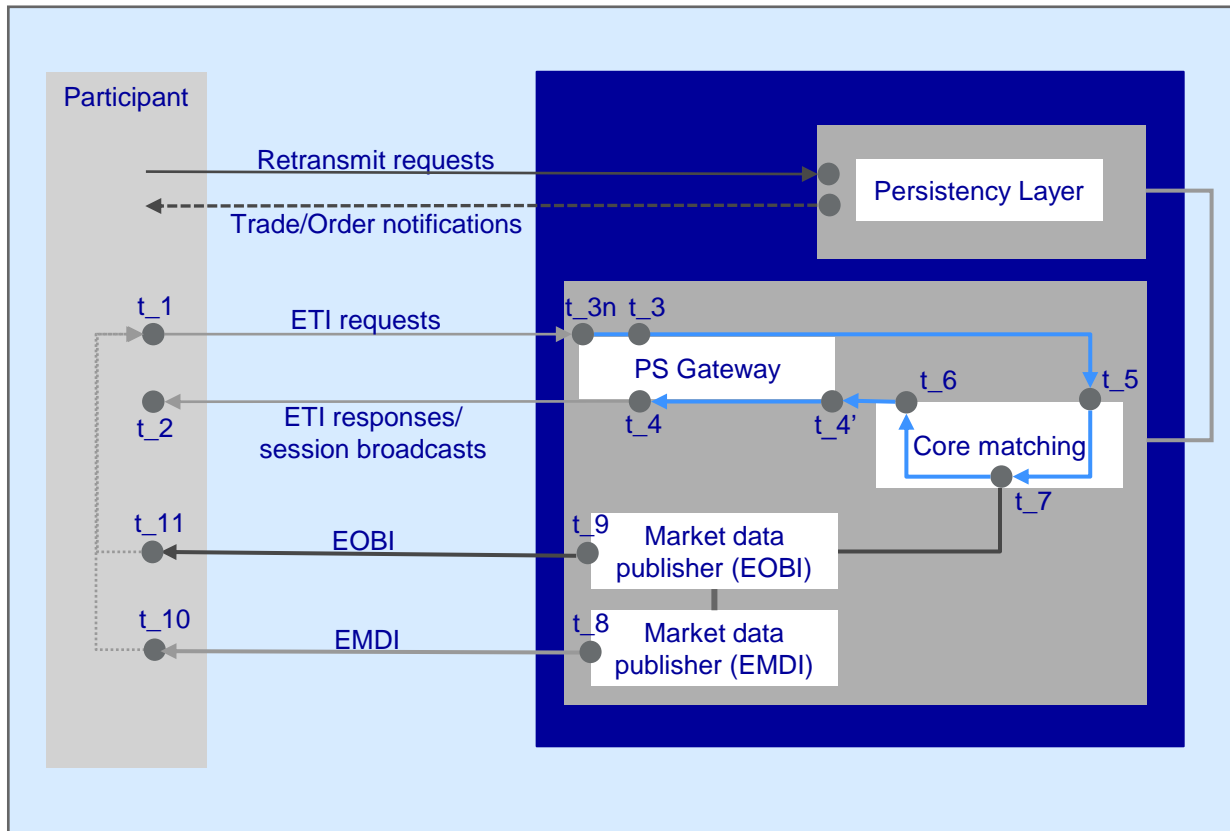
- Uses state-of-the-art infrastructure and hardware for high performance
- Offers reliable connectivity and enriched trading functionality
- Is multi-market capable, operates Derivatives (Eurex, EEX) and Cash markets (Xetra, Börse Frankfurt, Vienna, etc.)
- T7<sup>®</sup> consists of partitions. A partition is a failure domain in charge of matching, persisting and producing market data for a subset of products. Each T7 partition is distributed over two rooms in the Equinix data centre.
- There are 12 Eurex T7 and 11 Xetra T7 partitions.
- Separate partitions are used for markets of other exchanges hosted on T7 (e.g. Vienna (XVIE), EEX (XEEE), Bulgaria (XBUL), ...).
- The reference data contains the mapping of products to partition IDs.
- 4 LF gateways and one FIX LF gateway allow access to all Eurex partitions and the separate EEX partition.
- 4 LF gateways and one FIX LF gateway allow access to all Xetra partitions.
- 2 LF gateways and one FIX LF gateway are shared between Vienna and their partner exchanges.
- 2 LF gateways and one FIX LF gateway are shared between XBUL and XMAL

- Note that the active half of a partition is either on side A (for even partitions) or on side B (for odd partitions).
- In case of the failure of a PS gateway/Matching Engine or a market data publisher for EOBI or EMDI which is integrated into the Matching Engine, the active half of the service will shift to the other room.
- With consolidated PS gateway/Matcher process the active PS gateway and active Matching Engine act as a single failure domain within each partition, i.e. they will always fail as a single logical group.



# T7<sup>®</sup> Topology

## Overview



### Matching Engine:

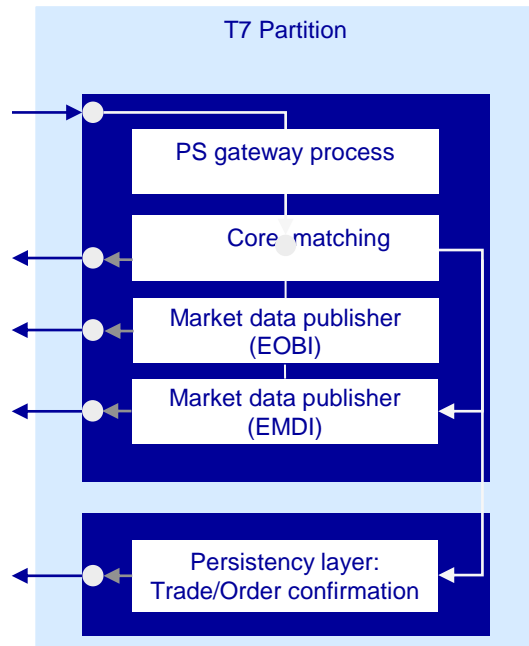
- order book maintenance & execution
- creation of direct responses as well as execution messages for passive orders/quotes
- creation of EOBI order book messages
- creation of EOBI order book snapshot messages
- creation of EMDI order book delta messages
- creation of EMDI order book snapshot messages

### Persistency:

- persistent order storage
- trade/execution history
- transaction history for standard orders
- creation of listener broadcast for standard orders

# T7<sup>®</sup> Topology

## Partitions



Orders/quotes entered for a specific product are sent by the PS gateway process to the core Matching Engine process (both residing on the same server in the same partition).

The matching priority is assigned when the orders/quotes are read into the Matching Engine.

The core matching component works as follows:

- when an order/quote arrives, it is functionally processed (e.g. put in the book or matched),
- handover of data to the EOBI market data publisher, followed by EMDI market data publisher and
- handover of all data resulting from the (atomic) processing of the incoming order/quote to persistency component in the partition.

Resulting responses and private broadcasts are sent out in the following order:

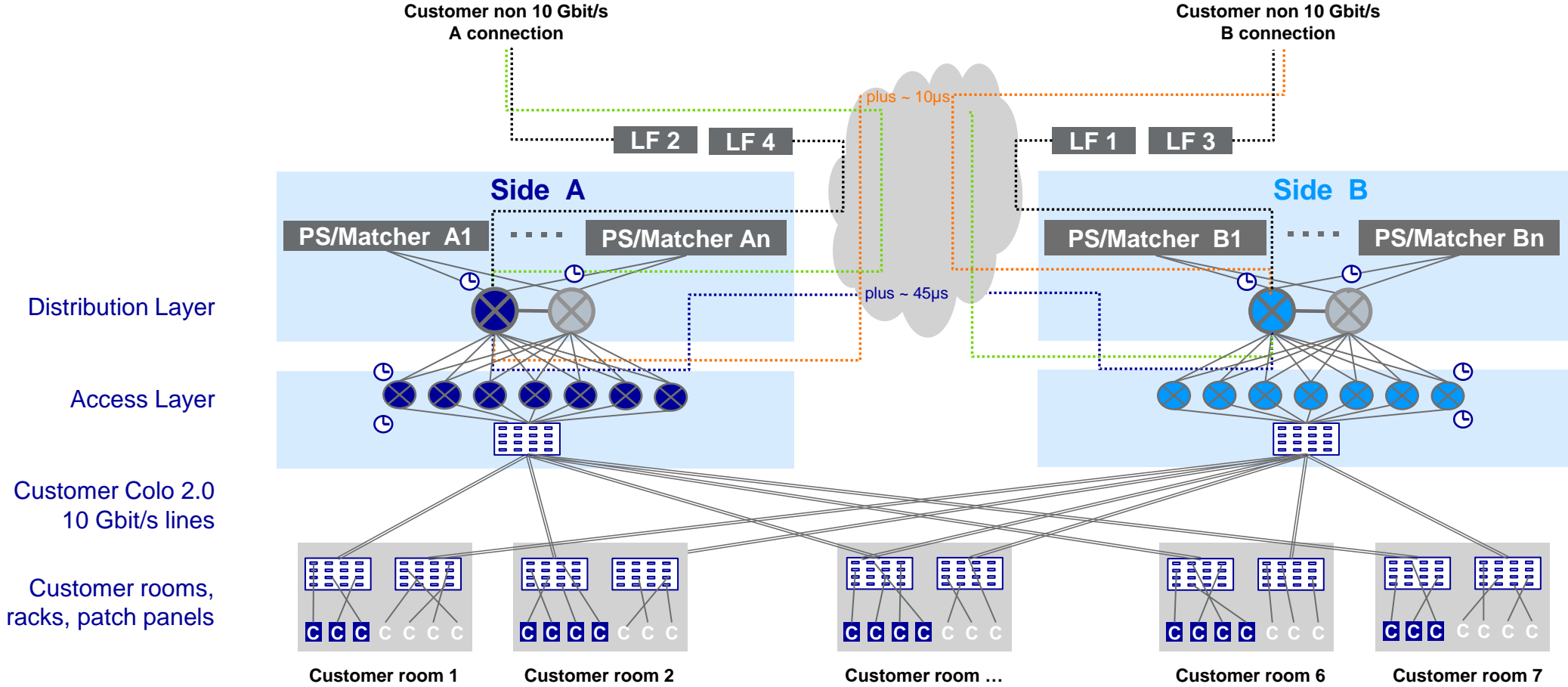
- direct response to the order/quote entered (for persistent as well as for non-persistent orders and quotes),
- fast execution information for booked orders/quotes (in case of a match).

The generation of listener broadcasts, trade confirmations (by the persistency server) and of non-EOBI/non-EMDI market data (by the market data publisher) is done on separate servers. Hence the order of the resulting messages from these servers is not strictly deterministic.

**Note that the Matching Engine holds states of orders in memory. All responses, broadcasts, EOBI and EMDI market data thus are preliminary by nature.**

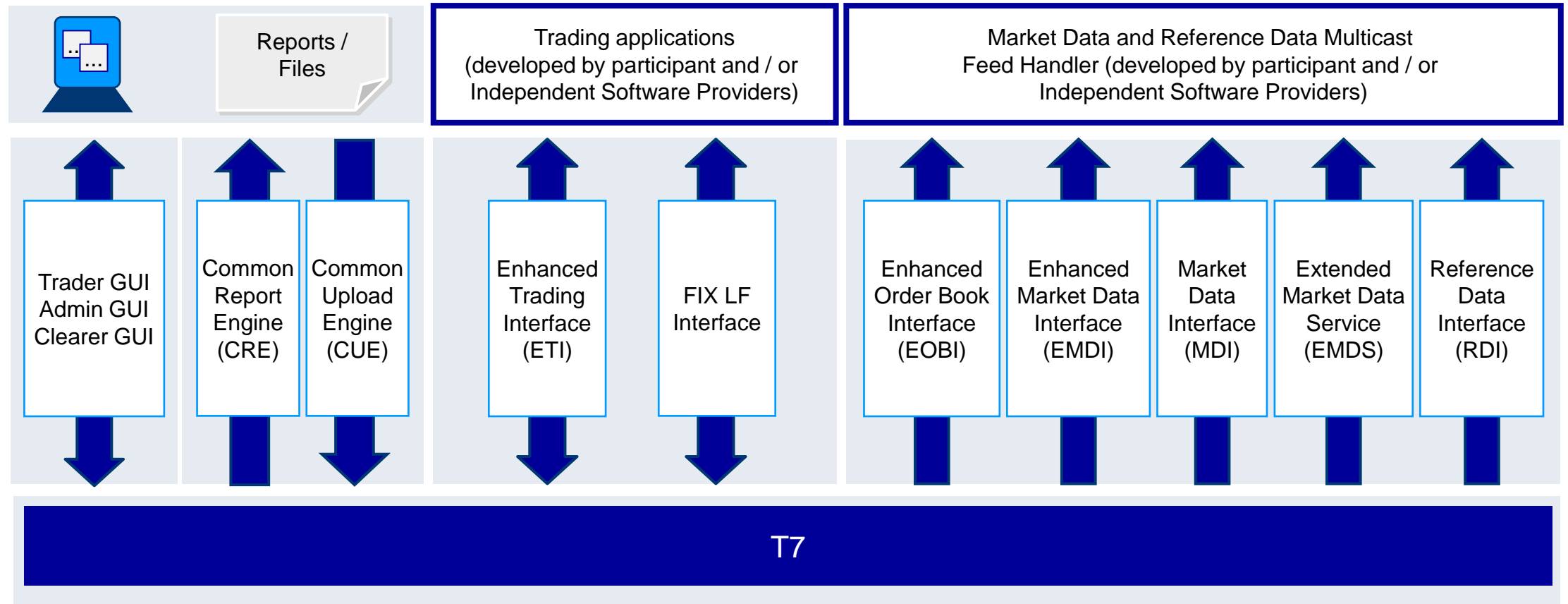
# Network Topology in Co-Location

## Eurex Order Entry



# T7<sup>®</sup> Trading System

## Interfaces



# Order Entry

## Introduction

Requests sent to T7 will be routed via an access network and a gateway.

There are the following basic connection alternatives:

### **Choice of Network**

There are two network classes connecting a participant's installation with the T7 gateways:

- Co-Location with 10 Gbit bandwidth and a one-way base latency of around 2  $\mu$ s.
- Other networks with less than 1Gbit bandwidth and a one-way base latency of minimum about 50  $\mu$ s.

### **Choice of Session Type**

T7 supports three session types:

- High frequency sessions connect to PS gateways for low latency access to a single partition (flat binary protocol: ETI).
- Low frequency sessions connect to LF gateways for convenient access to all partitions, with a considerably higher base latency (ETI).
- FIX LF sessions connect to FIX LF gateway for convenient access to all partitions using the FIX protocol, with a considerably higher base latency than LF gateways.

### **Remarks**

- LF gateway and FIX LF requests are routed via PS gateways.



# Order Entry

## Co-Location Network

Participants may use Co-Location to place their infrastructure in the datacentre that hosts the T7 system.

The Co-Location 10 Gbit network has the following properties:

### **Fair and equal access**

Regardless of the Co-Location room we ensure all lines are created equal.

More precisely the latency between the handover point in the participant's rack and the first (Access Layer) switch is calibrated to below  $\pm 2.5$  ns. Deutsche Börse worked on reducing this deviation in 2023 as an intermediate step and plans to improve the cable length normalization in 2025 significantly (see Outlook on [slide 4](#)).

### **Two redundant halves ('A' and 'B')**

There are two independent order entry network halves.

As active gateways are placed in either half there is an optimal side for each gateway (even numbered gateways are on the A side, odd number on the B side). The only exception is FIX LF: There is only one active FIX LF gateway which is by default located on the B-side.

Crossing sides, e.g. connecting to a B side gateway via an A network, is possible but results in at least 45  $\mu$ s higher base latency.

### **Two hierarchical switch layers**

Participants connect to Access Layer switches (currently 7 Eurex\*, 2 Xetra per side).

The uplink of each Access Layer switch is connected to a Distribution Layer switch.

The Distribution Layer switches have a direct connection to the active gateways on the respective side.

# Order Entry

## Gateways

There are three gateway types to access the T7 system:

### **Partition-specific (PS) Gateway combined with Matching Engine**

Protocol: flat binary (ETI)  
Allowed session types: High Frequency Sessions only  
Sequencing: FIFO operation (Sequence guaranteed from network card to Matching Engine in)  
Latency: lowest, median latency ~ 12  $\mu$ s network card to Matching Engine in  
Versatility: Allows routing to one partition only, only subset of broadcasts available

### **Low Frequency (LF) Gateway**

Protocol: flat binary (ETI)  
Allowed session types: Low Frequency Sessions only  
Sequencing: FIFO not guaranteed  
Latency: medium, (additional ~32  $\mu$ s latency compared to PS gateway direct access)  
Versatility: Routes to all partitions (via PS gateway), all ETI broadcast types available

### **FIX LF Gateway**

Protocol: FIX  
Allowed session types: Fix Sessions only  
Sequencing: FIFO guaranteed  
Latency: high, requests to the Matching Engine are routed via PS gateways  
Versatility: Routes to all partitions (via PS gateway), all FIX broadcast types available

# Market Data

## Overview

Market Data can be consumed over two distinct types of networks and in various types

### Choice of Network

There are two network classes available for market data:

- Co-Location with 10 Gbit bandwidth and a one-way base latency of around 2  $\mu$ s.  
10 Gbit connections are equalized in length (cable latency difference of less than  $\pm 2.5$  ns) and provide the lowest jitter.  
Deutsche Börse worked on reducing this deviation in 2023 as an intermediate step and plans to improve the cable length normalization in 2025 significantly (see Outlook on [slide 4](#)).
- Other networks with less than 1Gbit bandwidth with higher base latency.

### Choice of Market Data Type

There are three market data types:

- Order by Order market data (EOBI) with highest granularity and lowest latency in flat binary format.  
EOBI is sent out directly from the Matching Engine and is only available via 10 Gbit network.
- Price level aggregated market data (EMDI) with slightly higher latency in FAST encoded format.
- Netted price level aggregated market data (MDI) in FAST encoded format.

# Appendix



# Middleware, Network, Hardware and OS Overview

## **T7 uses state-of-the-art infrastructure components**

Intel Xeon Gold 6256 CPU (Cascade Lake Refresh) for Matching Engine and consolidated PS gateway/Matching Engine.

Intel Xeon Gold 6148 CPU (Skylake) for all other servers.

We currently use Red Hat Enterprise Linux 8.8.

T7 internal communication between its core components is based on Confinity Low Latency Messaging using an Infiniband network.

In Q3/Q4 2024, we will migrate our T7 core network from Infiniband to Ethernet and deploy new server hardware. After the upgrade, all servers in T7 will use Intel(R) Xeon(R) Platinum 8462Y+ CPUs.

## **T7 network access**

Deutsche Börse offers trading participants to connect via 10 Gbit/s cross connects to its T7 platform in the Equinix data centre.

The Co-Location offering uses Cisco Nexus 3548-X switches operating in cut-through mode.

All cables are normalized to give an overall maximum deviation between any two cross connects of less than  $\pm 0.5$  m ( $\pm 2.5$  ns).

Deutsche Börse worked on reducing this deviation in 2023 as an intermediate step and plans to improve the cable length normalization in 2025 significantly (see Outlook on slide 4).

Insight into network dynamics is offered by the High Precision Timestamp File service (see <https://www.mds.deutsche-boerse.com/mds-en/analytics/high-precision-timestamps>).

Participant facing interface cards on the gateways and market data publishers use Solarflare EnterpriseOnload wire order delivery API to bypass the kernel TCP stack and deliver messages in the same order received by the network card.

Cables connecting Line-of-Sight antenna cables have been equalized to  $\pm 1$  m by Equinix.

# Throttle and Session Limits

**In order to protect the trading system, T7 has several measures in place to ensure that its most vital components are not harmed by a malfunctioning client application. Therefore, transaction limits are imposed on T7 sessions.**

ETI LF sessions are available with throttle values of 150 or 50 transactions/sec ETI HF sessions are available with throttle values of 250\*, 150 or 50 transactions/sec. Furthermore, LF sessions that cannot enter orders/quotes but can only receive trade and listener broadcasts are available (at a reduced price).

The disconnect limit is set at:

- 750 for HF Ultra sessions with a throttle value of 250 transactions/sec. i.e. a session will be disconnected in case of more than 750 consecutive rejects due to exceeding the transaction limit (throttle).
- 450 for sessions with a throttle value of 150 transactions/sec, i.e. a session will be disconnected in case of more than 450 consecutive rejects due to exceeding the transaction limit (throttle).
- 150 for sessions with a throttle value of 50 transactions/sec, i.e. a session will be disconnected in case of more than 150 consecutive rejects due to exceeding the transaction limit (throttle).

Please note that in case the disaster recovery facility is used, all ETI sessions will have a throttle limit of 30 transactions per second.

For both limits, all technical transactions are counted using a sliding window.

The number of ETI sessions which can be ordered is limited. Currently, up to [600](#) Eurex sessions and 200 Xetra sessions per participant can be ordered.

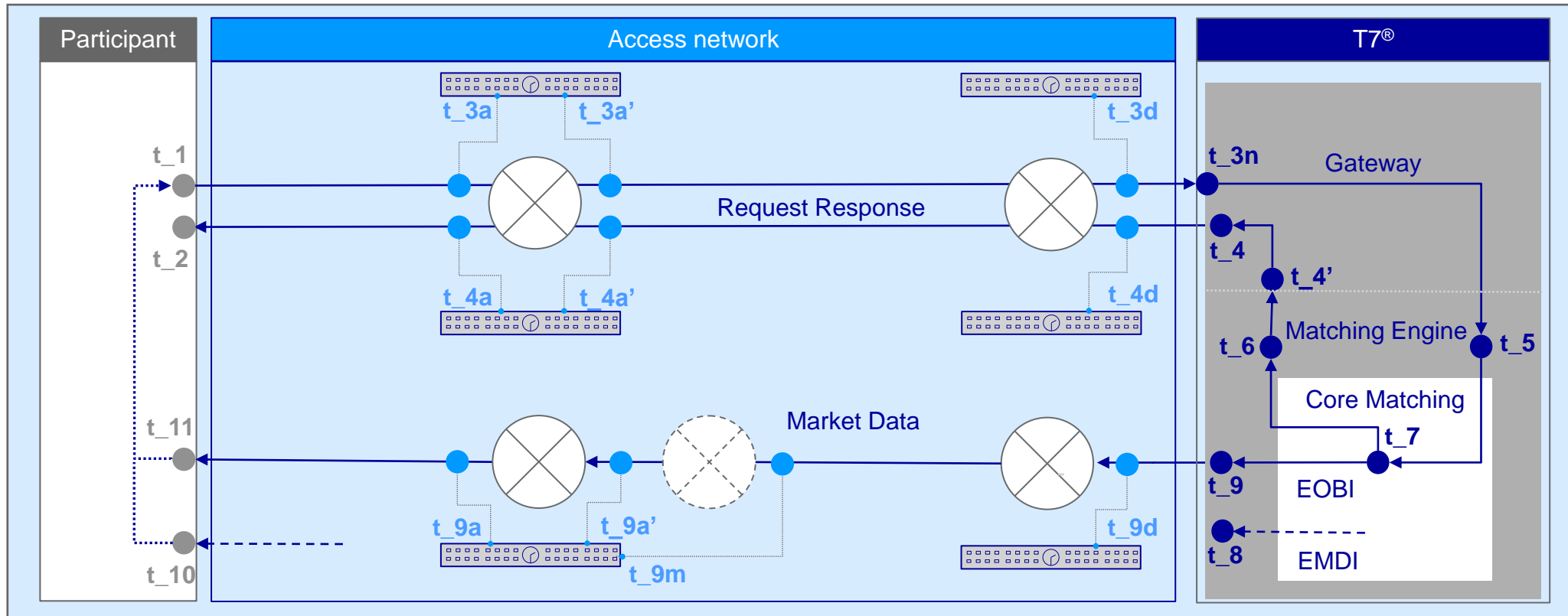
There is also a limit on the maximum number of sessions that can connect to a PS gateway concurrently per participant. This limit is currently configured to 80 sessions, see [Eurex Circular 122/17](#).

On 1 July 2019, we introduced a limit on the maximum number of outstanding session and trader login requests possible per business unit and per session at any given point in time. This limit is set to 50 on business unit level, 10 on session level. We recommend a synchronous login procedure, where a login request is sent on a session only after the previous login has been responded to. Please refer to the Incident Handling Guide for details.

For Eurex the number of order entry cross connects in colocation that may be used concurrently on a single day is limited to 6 per Access Layer switch. In addition, the number of allowed ethernet frames per cross connect is currently limited to [30.000](#) per second and [600.000](#) per minute for Eurex and 25,000 per second and [600,000](#) per minute for Xetra.

# T7<sup>®</sup> Topology

## Timestamps



- Timestamps provided in T7 API (in real time) in dark blue (t\_3n: taken by network card, other: application level)
- Network timestamps taken using TAPs and timestamping switches (Metamako)
- Timestamps possibly taken by participants shown in grey

# T7<sup>®</sup> Timestamps

## Description

t\_[x]a, t\_[x]a' time taken by network capture devices in the Access Layers.

t\_[x]d time taken by network capture devices in the Distribution Layers.

t\_9m time taken by network capture device between Access Layer and Mid-layer.

t\_1, t\_2 can be taken by a participant (e.g. via a network capture) when a request/ response is read from/written to the network.

t\_3n time taken by the PS gateway when the first bit of a request arrives on the PS gateway NIC;  
contained in (private) ETI response for PS gateway enabled partitions.  
Consecutive messages via the same session may be assigned to the same t\_3n.

t\_3 time taken by the ETI gateway application when a request is read from the socket on the participant´s side of the gateway;  
contained in (private) ETI response for transactions for non-PS gateway enabled partitions (e.g. XVIE).

t\_4' time taken by the ETI gateway when a response/ notification is received by the ETI gateway from the Matching Engine;  
contained in (private) ETI response/ notification.

t\_4 time taken by the ETI gateway when a response/ notification is written to the socket on the participant´s side of the gateway;  
contained in (private) ETI response/ notification.

t\_5, t\_6 time taken by the Matching Engine when a request/response is read/written; contained in (private) ETI response.

t\_7 time at which the Matching Engine starts maintaining the order book

t\_8 time taken by EMDI publisher just before the first respective UDP datagram is written to the UDP socket.

t\_9 time taken by EOBI publisher just before the first respective UDP datagram is written to the UDP socket.

t\_10, t\_11 can be taken by a participant (e.g. via a network capture) when a UDP datagram is read from the UDP socket.



# T7<sup>®</sup> Timestamp Reference

The timestamps t\_3,...,t\_9 are available via the following fields:

Timestamp	Tag no.	Field name	Present in
t_3, t_3n	5979	RequestTime	ETI Response EMDI Depth Incremental message, in case a trade is reported EOBI Execution Summary, Order Add, Order Modify, Order Modify Same Priority and Order Delete messages
t_4'	7765 25043	ResponseIn NotificationIn	ETI Response (from the Matching Engine) ETI Notification (from the Matching Engine)
t_4	52	SendingTime	ETI Response and Notification
t_5	21002 2445	TrdRegTSTimeIn AggressorTime	ETI Response (from the Matching Engine) EMDI Depth Incremental message, in case a trade is reported EOBI Execution Summary message
t_6	21003	TrdRegTSTimeOut	ETI Response and Notification (from the Matching Engine)
t_7	17 273 60 21008	ExecID MDEntryTime TransactTime TrdRegTSTimePriority	ETI Response (from the Matching Engine) EOBI Execution Summary message EMDI Depth Incremental, Depth Snapshot and Top of Book Implied message EMDI messages for other events EOBI Order Modify Same Priority and Order Delete messages EOBI Order Add and Order Modify messages
t_8	No tag	SendingTime	T7 EMDI UDP packet header
t_9	60	TransactTime	EOBI packet header
t_8 - t_5	No tag	PerformanceIndicator	EMDI UDP packet header of the T7 EMDI Depth Incremental stream

## Notes on timestamps:

All timestamps provided are 8 byte integers (in nanoseconds after Unix epoch).

The PerformanceIndicator is a 4 byte integer (in nanoseconds).

The Network timestamps (t\_[x]a, t\_[x]a', and t\_[x]d, t\_[x]d') are not available in any protocol field but some via the High Precision Timestamp File service.

# Thank you for your attention

**Sergej Teverovski, Manfred Sand, Phuong Hieke, Simon Braun, Raphael Colcombet**

**Trading System Analytics**

Deutsche Börse AG

Mergenthalerallee 61

65760 Eschborn

Germany

E-mail: [monitoring@deutsche-boerse.com](mailto:monitoring@deutsche-boerse.com)

For updates refer to

<https://www.eurex.com/ex-en/support/technology/t7> and <http://www.xetra.com/xetra-en/technology/t7/publications>



# Disclaimer

## © Deutsche Börse Group 2024

Deutsche Börse AG (DBAG), Clearstream Banking AG (Clearstream), Eurex Frankfurt AG, Eurex Clearing AG (Eurex Clearing) and Eurex Repo GmbH (Eurex Repo) are corporate entities and are registered under German law. Eurex Global Derivatives AG is a corporate entity and is registered under Swiss law. Clearstream Banking S.A. is a corporate entity and is registered under Luxembourg law. Deutsche Börse Asia Holding Pte. Ltd., Eurex Clearing Asia Pte. Ltd. and Eurex Exchange Asia Pte. Ltd are corporate entities and are registered under Singapore law. Eurex Frankfurt AG (Eurex) is the administrating and operating institution of Eurex Deutschland. Eurex Deutschland is in the following referred to as the "Eurex Exchange".

All intellectual property, proprietary and other rights and interests in this publication and the subject matter hereof (other than certain trademarks and service marks listed below) are owned by DBAG and its affiliates and subsidiaries including, without limitation, all patent, registered design, copyright, trademark and service mark rights. While reasonable care has been taken in the preparation of this publication to provide details that are accurate and not misleading at the time of publication DBAG, Clearstream, Eurex, Eurex Clearing, Eurex Repo as well as the Eurex Exchange and their respective servants and agents (a) do not make any representations or warranties regarding the information contained herein, whether express or implied, including without limitation any implied warranty of merchantability or fitness for a particular purpose or any warranty with respect to the accuracy, correctness, quality, completeness or timeliness of such information, and (b) shall not be responsible or liable for any third party's use of any information contained herein under any circumstances, including, without limitation, in connection with actual trading or otherwise or for any errors or omissions contained in this publication.

This publication is published for information purposes only and shall not constitute investment advice respectively does not constitute an offer, solicitation or recommendation to acquire or dispose of any investment or to engage in any other transaction. This publication is not intended for solicitation purposes but only for use as general information.

All descriptions, examples and calculations contained in this publication are for illustrative purposes only.

Eurex and Eurex Clearing offer services directly to members of the Eurex Exchange respectively to clearing members of Eurex Clearing. Those who desire to trade any products available on the Eurex market or who desire to offer and sell any such products to others or who desire to possess a clearing license of Eurex Clearing in order to participate in the clearing process provided by Eurex Clearing, should consider legal and regulatory requirements of those jurisdictions relevant to them, as well as the risks associated with such products, before doing so.

Only Eurex derivatives that are CFTC-approved may be traded via direct access in the United States or by United States persons. A complete, up-to-date list of Eurex derivatives that are CFTC-approved is available at: <http://www.eurexexchange.com/exchange-en/products/eurex-derivatives-us>. In addition, Eurex representatives and participants may familiarise U.S. Qualified Institutional Buyers (QIBs) and broker-dealers with certain eligible Eurex

equity options and equity index options pursuant to the terms of the SEC's July 1, 2013 Class No-Action Relief.

A complete, up-to-date list of Eurex options that are eligible under the SEC Class No-Action Relief is available at: <http://www.eurexexchange.com/exchange-en/products/eurex-derivatives-us/eurex-options-in-the-us-for-eligiblecustomers...> Lastly, U.S. QIBs and broker-dealers trading on behalf of QIBs may trade certain single-security futures and narrow-based security index futures subject to terms and conditions of the SEC's Exchange Act Release No.60,194 (June 30, 2009), 74 Fed. Reg. 32,200 (July 7, 2009) and the CFTC's Division of Clearing and Intermediary Oversight Advisory Concerning the Offer and Sale of Foreign Security Futures Products to Customers Located in the United States (June 8, 2010).

## Trademarks and Service Marks

Buxl®, DAX®, DivDAX®, eb.rexx®, Eurex®, Eurex Repo®, Eurex Strategy WizardSM, Euro GC Pooling®, FDAX®, FWB®, GC Pooling®, GCPI®, MDAX®, ODAX®, SDAX®, TecDAX®, USD GC Pooling®, VDAX®, VDAX-NEW® and Xetra® are registered trademarks of DBAG. All MSCI indexes are service marks and the exclusive property of MSCI Barra. ATX®, ATX® five, CECE® and RDX® are registered trademarks of Vienna Stock Exchange AG. IPD® UK Quarterly Indexes are registered trademarks of Investment Property Databank Ltd. IPD and have been licensed for the use by Eurex for derivatives. SLI®, SMI® and SMIM® are registered trademarks of SIX Swiss Exchange AG. The STOXX® indexes, the data included therein and the trademarks used in the index names are the intellectual property of STOXX Limited and/or its licensors. Eurex derivatives based on the STOXX® indexes are in no way sponsored, endorsed, sold or promoted by STOXX and its licensors and neither STOXX nor its licensors shall have any liability with respect thereto. Bloomberg Commodity IndexSM and any related sub-indexes are service marks of Bloomberg L.P. PCS® and Property Claim Services® are registered trademarks of ISO Services, Inc. Korea Exchange, KRX, KOSPI and KOSPI 200 are registered trademarks of Korea Exchange Inc. BSE and SENSEX are trademarks/service marks of Bombay Stock Exchange (BSE) and all rights accruing from the same, statutory or otherwise, wholly vest with BSE. Any violation of the above would constitute an offence under the laws of India and international treaties governing the same. The names of other companies and third party products may be trademarks or service marks of their respective owners.

Eurex Deutschland qualifies as manufacturer of packaged retail and insurance-based investment products (PRIIPs) under Regulation (EU) No 1286/2014 on key information documents for packaged retail and insurance-based investment products (PRIIPs Regulation), and provides key information documents (KIDs) covering PRIIPs traded on

Eurex Deutschland on its website under the following link: <http://www.eurexexchange.com/exchangeen/resources/regulations/eu-regulations/priips-kids>.

In addition, according to Art. 14(1) PRIIPs Regulation the person advising on, or selling, a PRIIP shall provide the KID to retail investors free of charge.