# Load Testing Report
# Moscow Exchange Trading & Clearing Systems
# 21 March 2020

## Contents

# Testing objectives

1. To verify the trading and clearing systems operation under conditions of peak loading and an increased number of orders and trades. The trading systems of the following Moscow Exchange's markets were tested:
   - The Equity & Bond Market;
   - The FX Market;
   - The Derivatives Market.

2. To estimate the time of order filling and data delivery from the trading and clearing systems at different load levels and software and hardware configurations.

3. To allow third party software developers and brokers to test their systems and estimate the throughput capacity of communication channels to the exchange venues.
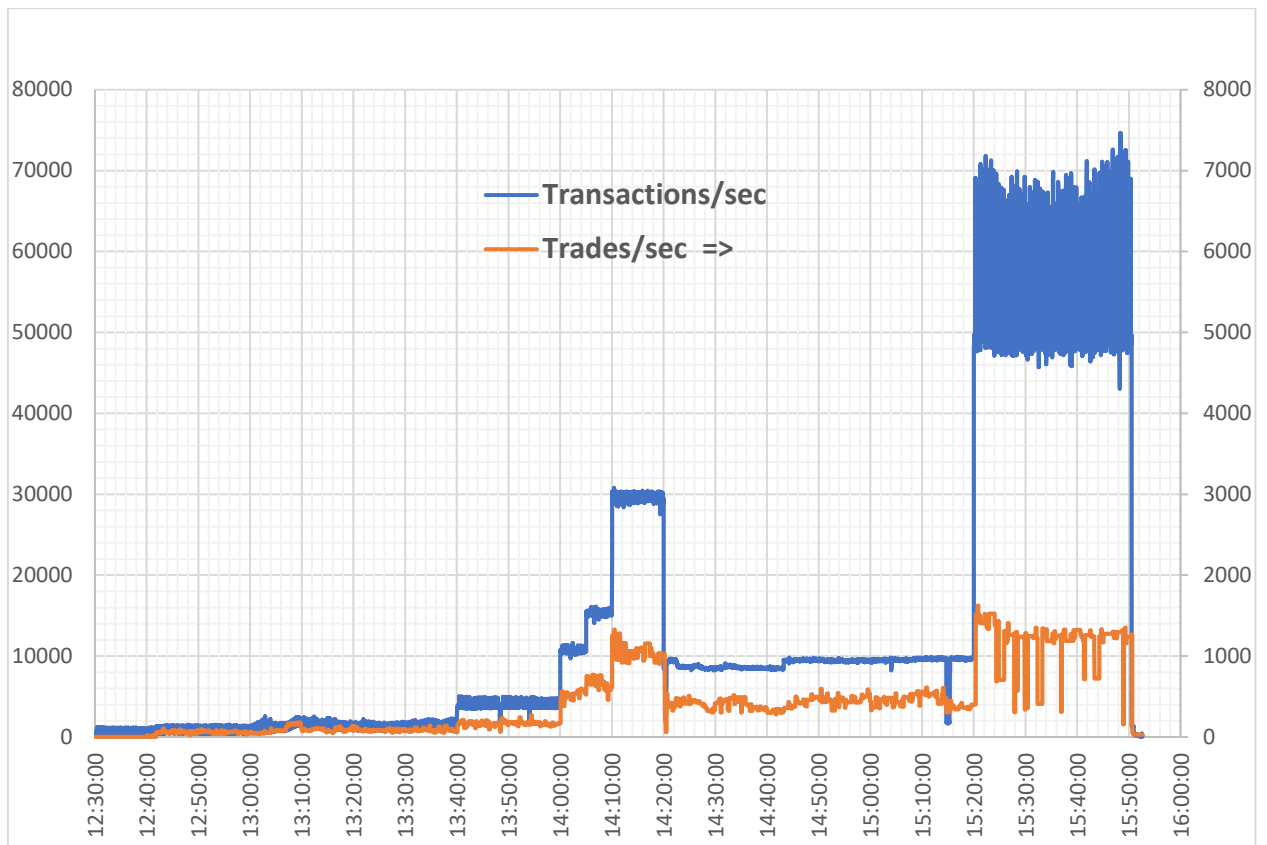

# Main results

## The Equity & Bond Market trading and clearing system

The testing was carried out on the current production version of the system.

The table below shows comparative performance in testing in 2019. The 'accepted transactions' term means all the incoming transactions that led either to order registration or to successful cancelation of an order.

|  | Transactions | Orders | Trades |
|---|---|---|---|
| Values reached (units), 2020 | 159 625 411 | 82 011 382 | 5 182 189 |
| Values reached (units), 2019 | 124 527 476 | 64 906 598 | 1 381 905 |
| Max processing rate for accepted transactions (units per sec), 2020 | 74 636 | 38 070 | 1740 |
| Max processing rate for accepted transactions (units per sec), 2019 | 62796 | 31913 | 3289 |
| Performance growth 2020 to 2019 | **+ 19%** | **+19%** | |

The graph below shows the frequency of transactions, orders and transactions by clients – testing participants. The maximum system performance level was not reached during the load test.

The share of client generated activity was 28% of all transactions. That is significally higher than their share in 2019 (4.9%).

From 15:20 to 15:50 the transaction generator emited transactions with variable frequences. The average speed of transaction flow was around 60,000 tran/sec, with is much faster than the one-second peacks in real trading (up to 7,500 tran/sec and 99.9% of one-second intervals with activity less than 2,300 tran/sec).

The peak frequency was generated in the repetitive intervals with durations of approximatily 5 seconds. Unlike the permanent peak load, this scenario allows to measure recovery time needed for some components of the trading system whise ability to function normally was interrupted by peak load. A stress nature of the load testing was remained.

The increased throughput of 19% as compared to 2019 is explained by the replacement of hardware of the trading and clearing system engines in December 2019.
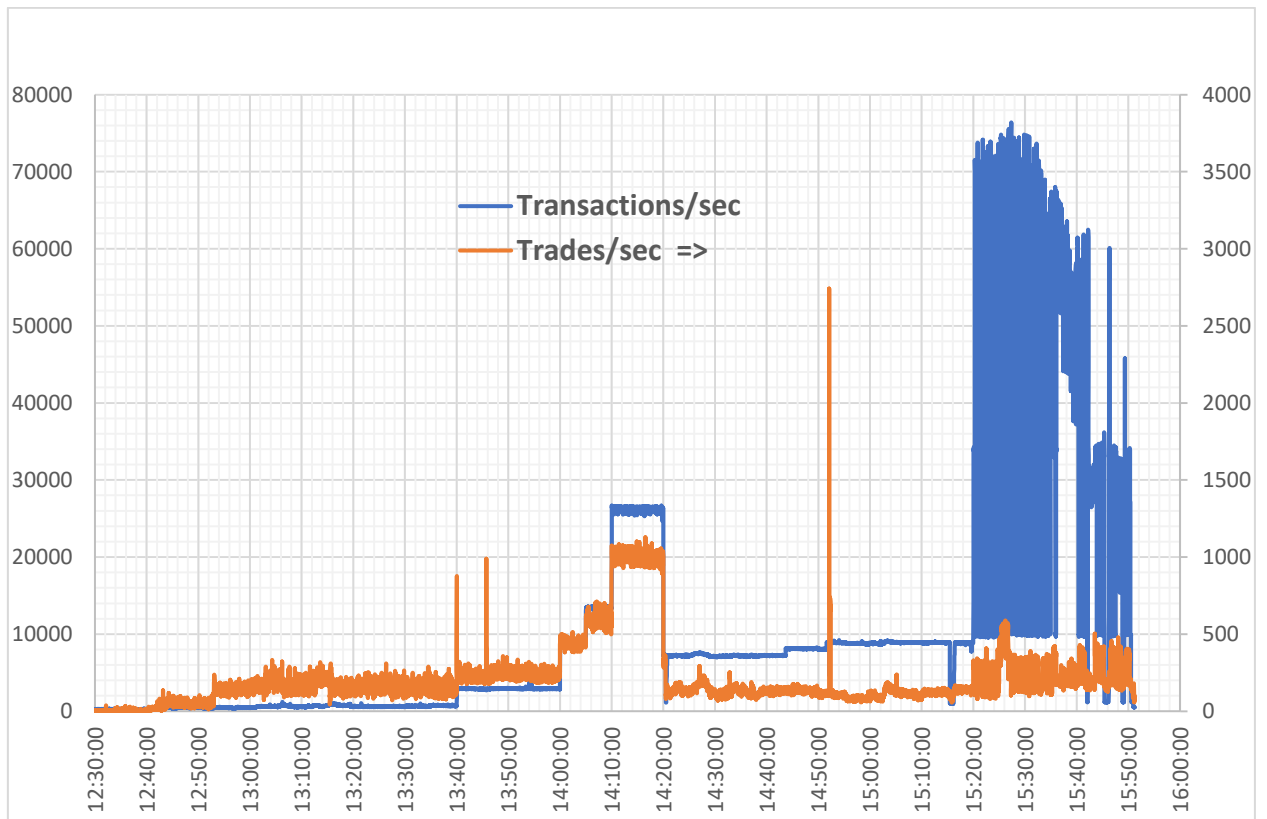

## The FX Market trading and clearing system

The testing was carried out on the current production version of the system.

The table below shows comparative performance in testing in 2019. The 'accepted transactions' term means all the incoming transactions that led either to order registration or to successful cancelation of an order.

| | Transactions | Orders | Trades |
|---|---|---|---|
| Values reached (units), 2020 | 118 962 847 | 60 538 374 | 2 606 363 |
| Values reached (units), 2019 | 129 787 106 | 70 068 396 | 886 128 |
| Max processing rate for accepted transactions (units per sec), 2020 | 76390 | 39050 | 2742 |
| Max processing rate for accepted transactions (units per sec), 2019 | 69924 | 34442 | 1185 |
| Performance growth 2020 to 2019 | **+9%** | **+13%** | |

The graph below shows the frequency of transactions, orders and transactions by clients – testing participants. The maximum system performance level was not reached during the load test.



The share of client generated activity was 25% of all transactions. That is significaly higher than their share in 2019 (8.2%).

From 15.20 to 15.32 the transaction generator emited transactions with variable frequences. The average speed of transaction flow was around 38,000 tran/sec, which is much higher than the actual production activity peaks (up to 6,000 tran/sec and 99.9% of one-second intervals with activity less than 3,000 tran/sec).

The peak frequency was generated in the repetitive intervals with durations of approximatily 5 seconds. Unlike the permanent peak load, this scenario allows to measure recovery time needed for some components of the trading system whose ability to function normally was interrupted by peak load. A stress nature of the load testing was remained.

From 15:32 to 15:50 a significant throughput degradation and transaction latency growth from 300 microseconds to 1 millisecons were observed. Presumably it was caused by network packet losses caused by simultanious peak activity on all the three markets. At the same time the FX market system was functioning properly and clients were able to enter orders and receive market data with avarage delays.

## The Derivatives Market trading and clearing system

The testing was carried out on the SPECTRA system version 6.4 used in production since 15 February 2020 on the servers at the main Data Space data center.

Over 4,000,000 client accounts were added to the trading system before load testing start with new positions added on those accounts during the load test.

The ratio and load profile of transactions submitted over the different protocols was similar to the production.
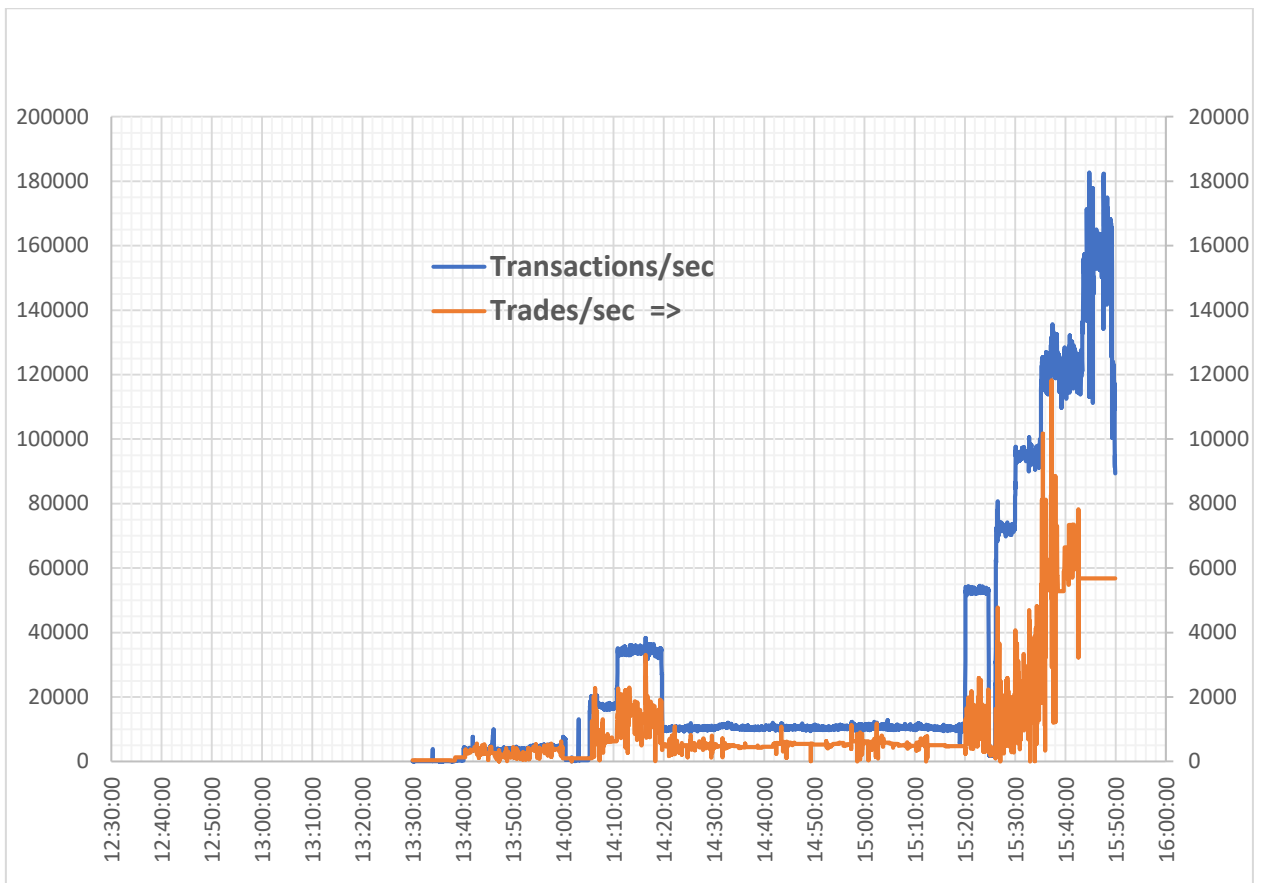
The order-to-trade ratio in the testing was also near the production value. 278 million transactions were sent and 10.7 million trades were executed during the testing. The peak performance was 182,000 transaction/sec with the stable trading system performance at average speed of around 150,000 transaction/sec. The 'accepted transactions' term means all the incoming transactions that led either to order registration or to successful cancelation of an order.

| | Transactions | Orders | Trades |
|---|---|---|---|
| Values reached (units), 2020 | 278 508 915 | 196 000 000 | 10 696 865 |
| Values reached (units), 2019 | 267 744 181 | 135 000 000 | 7 100 000 |
| Max processing rate for accepted transactions (units per sec), 2020 | 182 000 | - | - |
| Max processing rate for accepted transactions (units per sec), 2019 | 167 000 | - | - |
| Performance growth 2020 to 2019 | **+9%** | | |

During the testing, we carried out the scheduled intraday clearing session. Despite large volumes of orders and trades, clearing was performed as usual within the established time frames.

Clients participating in testing generated 1.77% of the transactions.

The graphs below shows a transaction load on the Derivatives Market trading system.

## ASTS Gateways

Equity market ASTS gateways as well as gateways of the FX market installed at both DSP and M1 data centers were functioning correctly.

During the normal operation of the Equity and FX markets systems, the clearing gateways are only available at the main DSP data center. The clearing engine located at the backup facility works in a slave mode and will only serve clearing gateways in case of the main data center failure and migration of the trading engine to the backup facility. This scenario was not included into the load test program. Migration to the backup data center has successfully been tested during the disaster recovery testing in 2020.

The Equity market clearing gateways were functioning correctly throughout all the load levels.

The FX market clearing gateways were functioning correctly at a constant transaction flow rate of up to 38,000 transactions per second. When this threshold was being exceeded, there were delays in clearing data refreshes at gateways as compared to the main clearing engine. In real trading, the peak frequencies of more than 38,000 transactions per second happen only in short bursts, so the expected delays will not exceed 5 milliseconds which is acceptable for clearing related data.

The FIX gateways configuration for both FX and Equity markets was similar to its production version. The gateways functioned normally across the whole range of transaction frequencies.
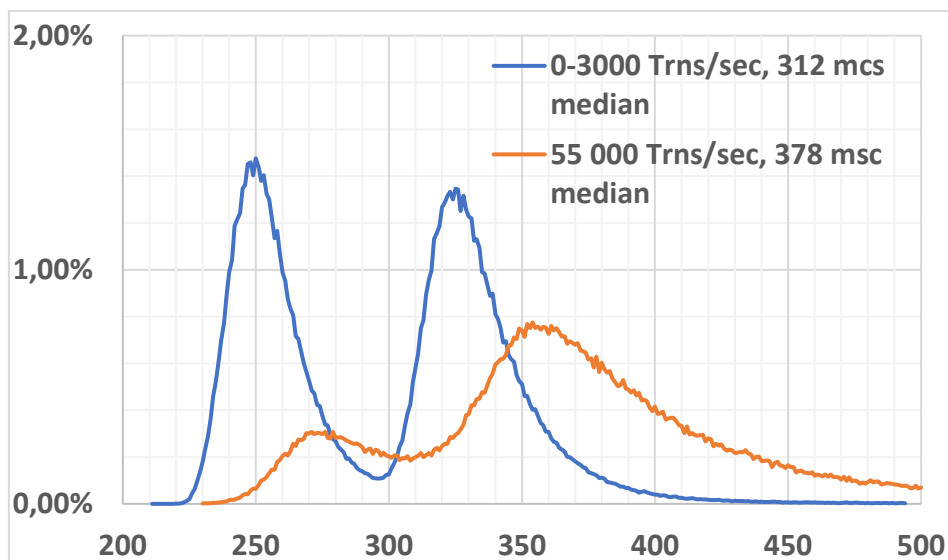
## SPECTRA Gateways

During the load testing, no deviations from the normal gateways performance was seen. The FIX, TWIME and Plaza2 gateways configuration was similar to its production version. The gateways functioned normally across the whole range of transaction frequencies. The disaster recovery scenario was not included into load test program.

## Latency for transactions, Equity & Bond Market and FX Market trading systems

Transaction generators at the Exchange side were using both Linux version of the embedded ASTS Bridge (libmtesrl.so library) and the FIX protocol. These generators had been set up on a server connected to the trading network. Latency data for the FIX protocol for Equity and FX Markets was collected using the network monitoring system based on Corvil solution and customized for ASTS FIX messages.

The graph below displays probability distribution for latency for incoming messages.
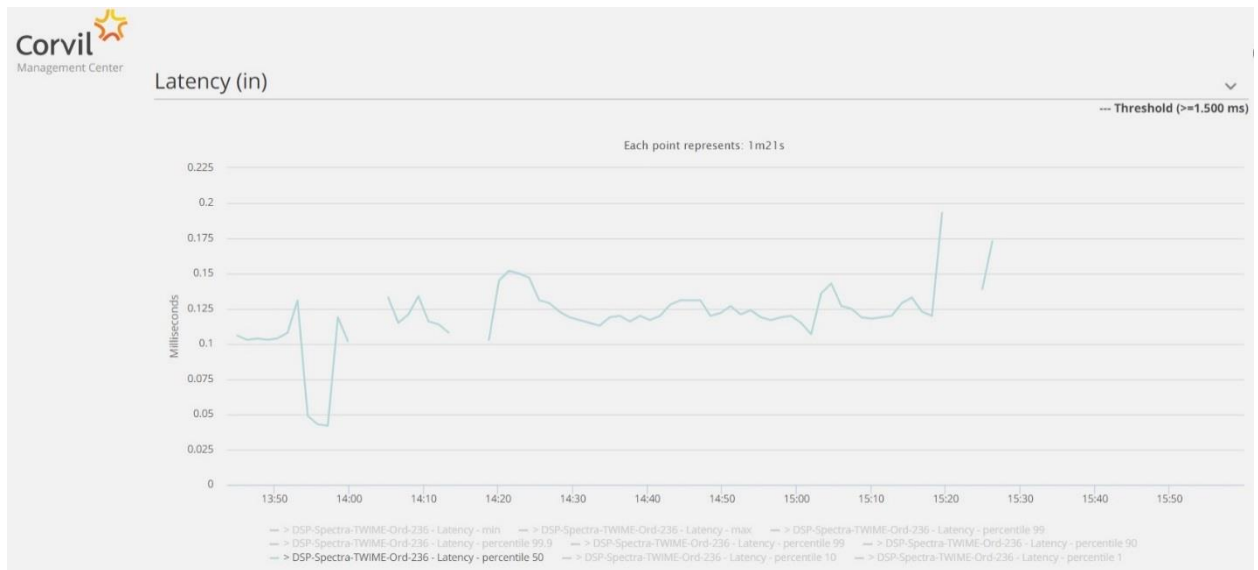


The first peak corresponds to the reply receipt for an order cancellation message. The second – response time for a new order placement. Latency for order cancellation is lower because this operation does not involve pre-trade collateral verification.

The median time for a flow with equal amount of new order and order cancel messages will be between these two peaks. With transaction growth will increase a probability that an order cancel message may have to wait while an earlier new order message is being processed. That's why the first peak will be lower at higher frequencies and both probability peaks will shift right on 30-40 microseconds.

This graph is applicable to both Equities and FX markets. Time distribution for users of native ASTS Bridge protocol will be shifted on approximately 60 microseconds right as compared to FIX protocol.

## Latency for transactions, Derivatives Market trading system

The internal Spectra monitoring system, Corvil solution and log files for client transactions have been used to measure the Derivatives Market latency.
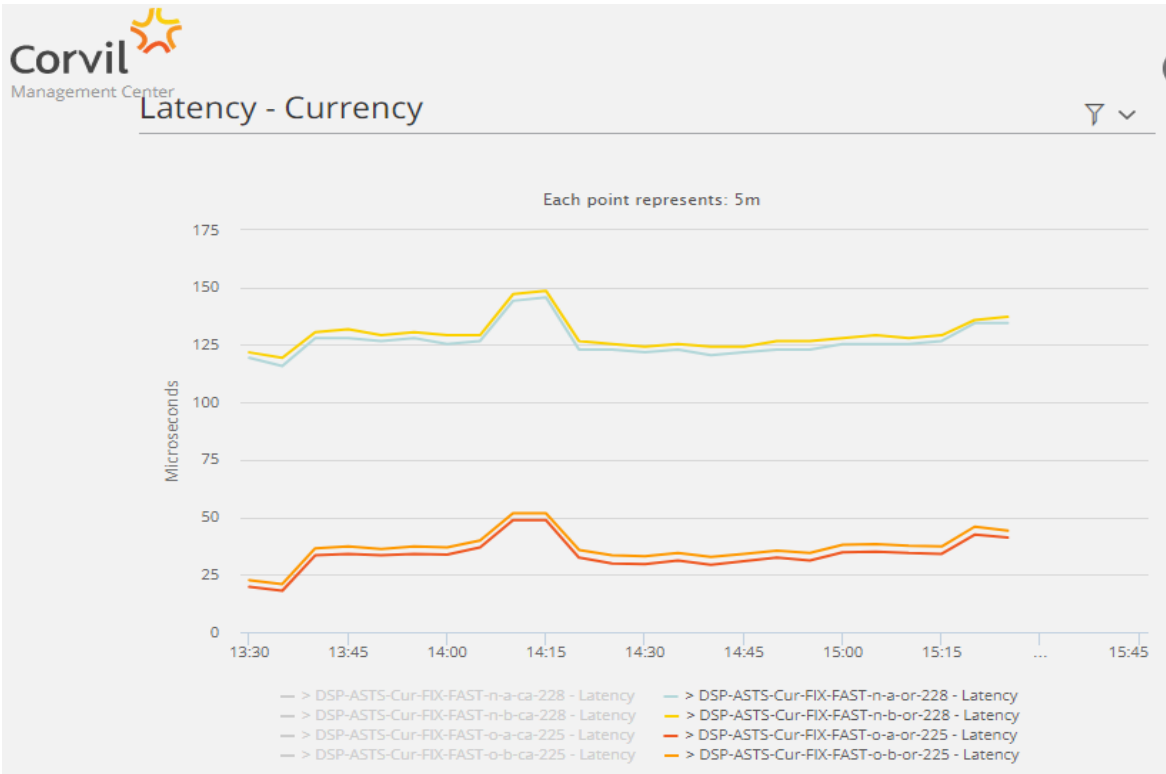


Within the range from 3,000 to 90,000 transactions per second, the average RTT for TWIME gate is between 100 and 150 microseconds. Then latency increased due to increased load on the system engine.

## FIX/FAST UDP multicast marketdata of the Equity & Bond Market and FX Market

Equity and FX markets FAST servers configuration was the same as in production. The servers operated correctly at all load levels except for time interval between 15:32 and 15:50 when the first line FX market FAST multicast was distorted by network packet losses. The second line multicast was operating correctly at the same time.
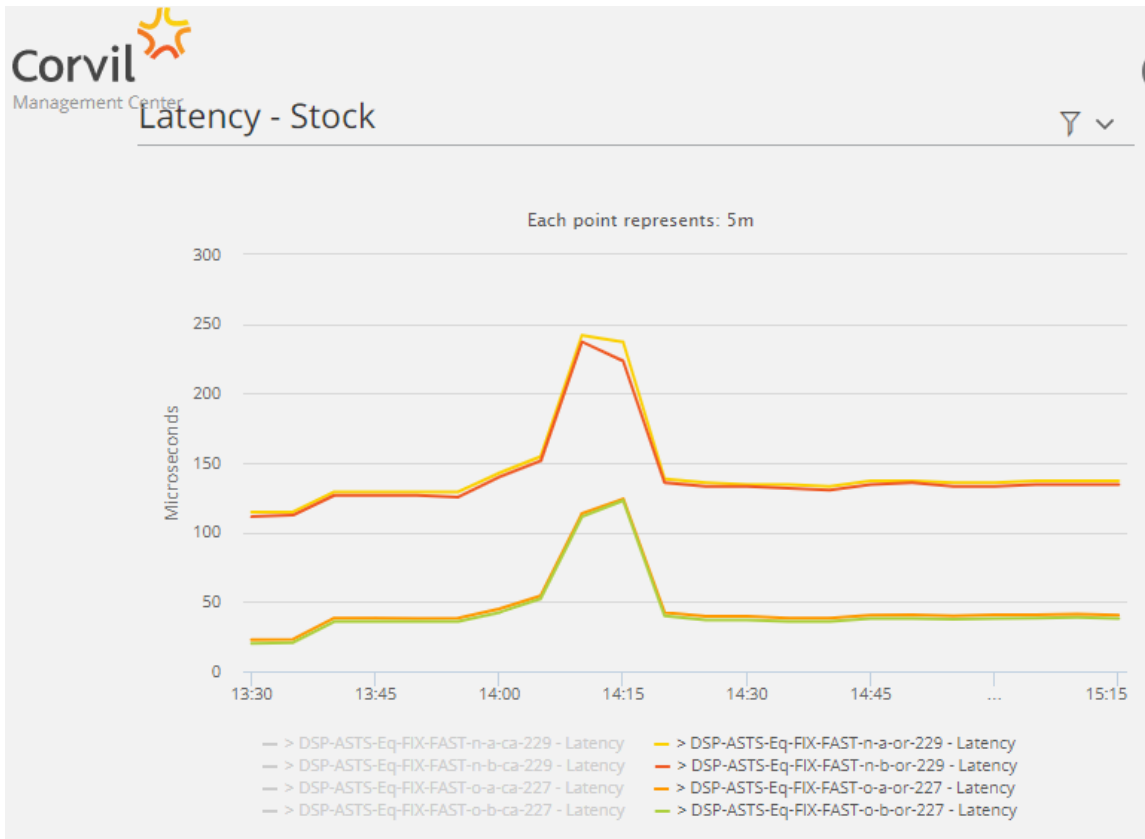
Statistical data on new order messages from the trading system (OLR feed, 279=0) relative to FIX protocol order accepted messages (Execution Report / 150=0) has been collected using the Corvil equipment. This data is being constantly collected during the normal trading.

For the FX Market, the average time of publication in the FAST feeds is shown in the screen below for the primary (225) and secondary (228) publishing lines.

Latency - Currency

Each point represents: 5m

At maximum transaction frequencies, the average relative publication times increase to 50 and 160 microseconds, respectively. Data after 15:30 is absent in Corvil due to network infrastructure overload.

For the Equities Market, the average time of publication in the FAST feeds is shown in the screen below for the primary (227) and secondary (229) publishing lines.



Latency - Stock

Each point represents: 5m

When approaching the maximum Equities market system capacity for transaction frequency, average relative publication times for both lines increased by 2-30 milliseconds. This is also caused by network infrastructure overload.

UDP multicast traffic for the FX market reached the following values in each of copies A and B:

| Feed | FX market, Mbps | Equity market, Mbps |
|---|---|---|
| Active orders (OLR) | 18 | 14 |
| Market statistics (MSR) | 27 | 20 |
| other feeds | 2 | 2 |

We recommend clients to carefully select their subscriptions to data feeds and consider the network bandwidth as the total combined traffic of the two FAST lines of the FX Market and Equity markets in copies A and B may reach 500 Mbps.

In the production environment, the short FAST traffic bursts would most probably correspond to the network requirements stated above.
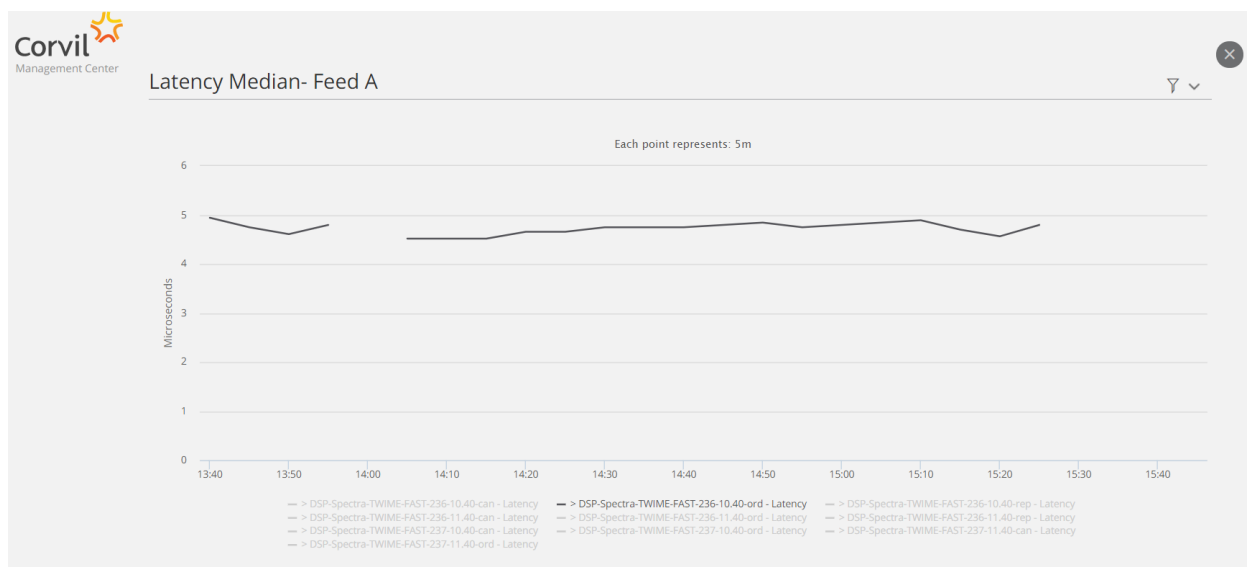
Recommendations given at http://www.moex.com/a1160 are applicable to each FAST line.

## FAST UDP multicast marketdata servers of the Derivatives Market
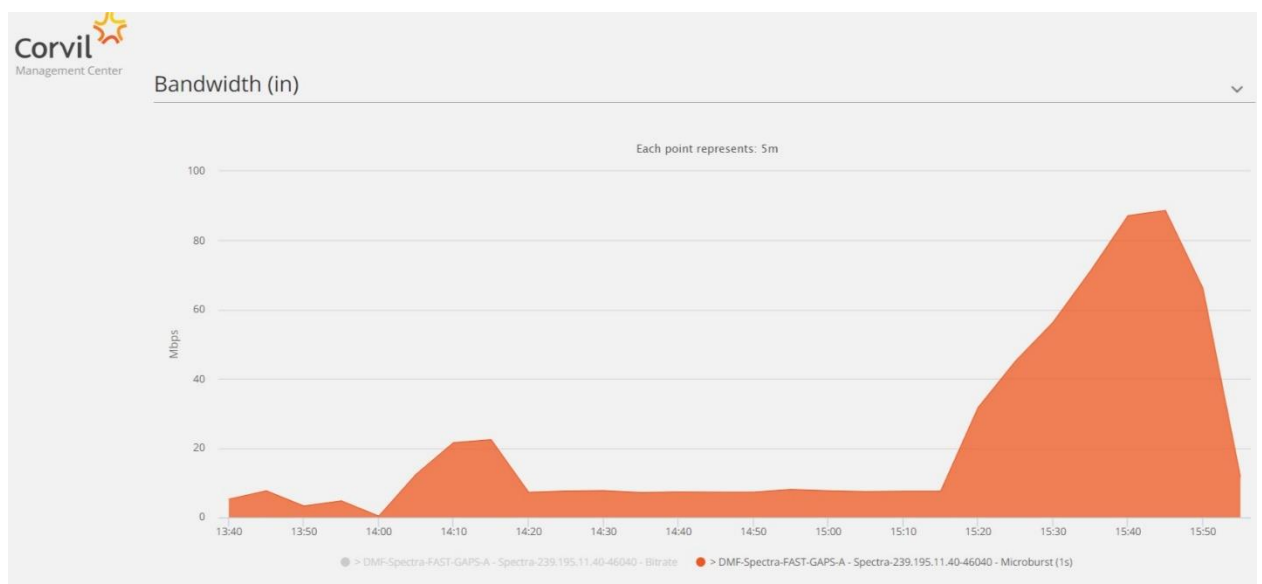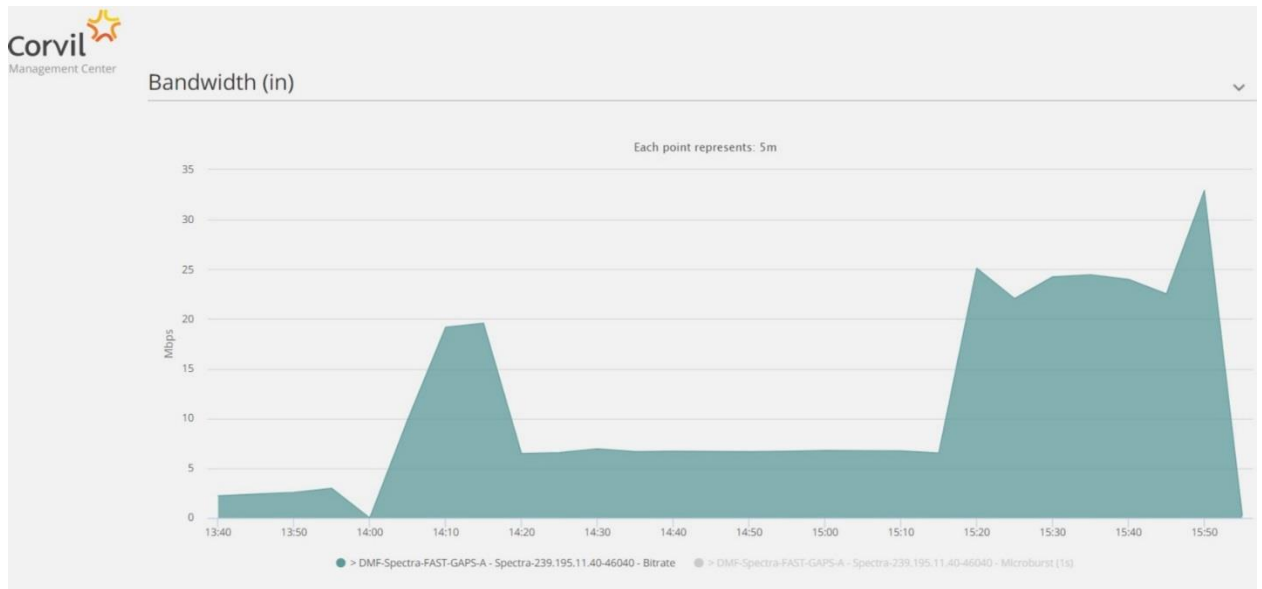
The Derivatives Market FAST servers configuration was the same as in production. The servers operated correctly at all load levels.

Statistical data for TWIME protocol execution reports relative to new order messages in the trading system (Full order log) has been collected using the Corvil equipment. This data is being constantly collected during the normal trading.

The average time of data publication in is shown in the screen below (the intermission on the graph was caused by the intermediate clearing session).



UDP multicast traffic for the Derivatives Market reached 35 mbps in each of copies A and B with bursts of up to 100 mbps. Average publication time at 50,000 transactions/sec were between 4-5 microseconds.

The required network bandwidth for clients who wish to use the FAST service to receive ORDERS-LOG is minimum 100 Mbit/sec per feed. To receive two feeds, FEED A and FEED B, or data from more than one market, the 1-10 GBit/sec bandwidth is recommended.

## Exchange network and colocation network

The network monitoring system within the exchange perimeter and within the colocation zone perimeter (including the colocation zone core switches) showed that the network parameters had no deviation from normal work parameters. Throughout the whole testing period including the highest activity times (at 300 Mbps on FAST FEED A for all markets) no retransmissions, packet loss or network latency growth deviations as compared to the production values have been detected.

## Subsystem for real time monitoring of the trading system parameters and market activity

The internal monitoring facilities operated well and provided data visualization in graphic form. Message signals were produced in accordance with the established criteria; data was collected to the monitoring database without fails. Operation of the monitoring system did not influence the facility performance.

During the testing, the Corvil (www.corvil.com) equipment and software, which had been adapted to analyze trading systems network traffic for all markets has been actively used. This monitoring was working properly for all markets throughout the whole testing except for 15:32-15:50 when network overload was experienced.

## Index server, market maker server, risk monitoring, MOEX web site

All the listed systems have operated correctly, no failures or performance problems have been noted.

## Test of clock synchronization over PTP (precision time protocol) at stress load

For the both data centers, an infrastructure for synchronizing system clocks over a high precision PTP protocol is deployed. During the tests, its stability at unusually high network loads on the Exchange infrastructure has been checked.

Precision checks for clocks on the network devices have shown that the time deviation has been no more than 500 nanoseconds and that is considered to be the excellent result. No failures of synchronization on the network devices or servers have been noted.

# Conclusions

## The Equity & Bond Market, the FX Market

1. Performance of the Equity Market system has grown by 19% compared to 2019.

2. Performance of the FX Market system has grown ty 9.8% compared to 2019.

3. No failures of the system components caused by program errors or overload during this testing have been registered.

4. We will perform further research on feasibility of ASTS and Spectra network separation to rule out conflicts at peak load.

## Derivatives Market

1. Peak performance increased by 9% and stable operation level increased by 15% compared to 2019.

2. Latency at transaction frequencies from 3,000 to 90,000 transactions per second decreased by 29% compared to 2019.

3. SPECTRA's performance is sufficient to meet demands of participants even at peak loads for order processing and market data distribution as well as for system capacity for the increased number of end client accounts and their position accounting even at top load levels. No failures of the system components caused by program errors or overload during this testing have been registered.

4. The bandwidth requirements for customers who use Plaza 2 protocol remain unchanged compared to the previous year:

   - At least 4 Mbps is required for stable operation of client bridges/terminals per each software instance.

   - At least 10 Mbps for a bridge in case of using feeds with a full order/trade log (FORTS_ORDLOG_REPL/FORTS_DEALS_REPL).

# Comparison of test results with actual production data

In this section we compare numbers received during the load test with the peak values actually recorded in real production trading.

Maximum numbers for orders and trades are listed according to the production system configuration files. These numbers may be increased by 2-3 times without requiring hardware upgrade.

| Parameter | FX | Derivatives | Equity |
|---|---|---|---|
| Peak number of trades per day | 300 000 | 3 371 091 | 3 000 000 |
| Number of trades in load testing | 2 606 363 | 10 696 865 | 5 182 189 |
| Peak number or orders per day | 20 000 000 | 59 147 377 | 40 000 000 |
| Number of orders in load testing | 60 538 374 | 196 000 000 | 82 011 382 |
| Maximum number of trades in production configuration | 3 000 000 | not applicable | 6 000 000 |
| Maximum number of orders in production configuration | 100 000 000 | not applicable | 100 000 000 |
| Peak transaction frequency in one-second intervals in real trading | 8 000 | 25 000 | 7 500 |
| Peak transaction frequency in load testing | 76 390 | 150 000 | 74 636 |