

APCON Latency Monitoring Solution for Financial Trading Markets



White Paper

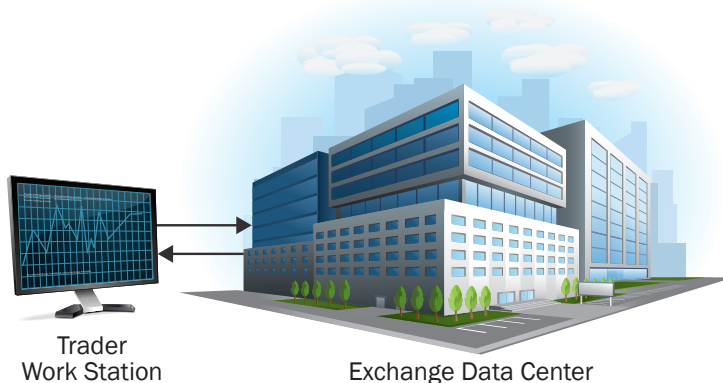


The world's financial markets have all but abandoned the telephone and the trading floor as their means of transacting business. Trading now takes place predominantly in interconnected computer networks where orders are placed electronically. APCON leads the industry in precision time stamping to measure transaction latency in high speed trading networks.

The trend towards high speed trading has challenged the providers of online trading systems to achieve ever-lower network latency and handle ever-larger traffic volumes. A key requirement for achieving predictable high performance and low latency is the ability to measure and track the performance of the network infrastructure in place. Accurate measurement of network latency from moment to moment allows data-based bandwidth planning, and can help spot performance bottlenecks before they endanger service levels.

What Is Network Latency?

The term "latency" can be used in reference to a variety of network measurements. The most common is transactional or exchange latency. This is simply the time interval between the sending of a trade order and when that same order is acknowledged and acted upon by the receiving party.



Exchange latency has been much-discussed in the financial sector. According to Information Week, "A millisecond advantage in trading applications can be worth \$100 million a year to a major brokerage firm." It was quickly understood that few benchmarks or qualified industry standards existed to measure it – so as trading systems grew there was a rush of development to create these types of measurement and analysis systems. Solutions in this space include Corvil, Correlix, and others. These products seek to measure exchange latency in a given system.

Another type of latency – the type specifically relevant to network monitoring – is the latency introduced by the monitoring switch between the production network and the analysis tools. That is to say, if a network path has been tapped in several places to measure the exchange latency of the primary link traffic, what is the latency between the tap point and the latency monitoring tool? The most important factor in this area is ensuring that the monitoring network latency does not skew latency information about the traffic traversing the production network.

Given that latency is known to be added by any monitoring system, steps must be taken to make the latency consistent and predictable, or the information derived from the monitoring will not accurately reflect exchange latency on the production network.

How Do Financial Institutions Determine Latency?

To adhere to the stringent requirements for security and low, predictable latency in their network architectures, most banks, securities traders and financial companies are constantly monitoring production network traffic and measuring the elapsed time required for a packet to traverse the network. The accuracy of this measurement depends upon the ability to identify and isolate any element that introduces latency.

To monitor latency levels, trading service networks tap the network path in several locations and observe the progress of a packet through the network. The time delta between various monitoring points can be added up to produce an accurate picture of moment-to-moment latency, and anomalous latencies can be isolated for further examination.

Industry research indicates that network engineers in financial organizations want a “best practices” approach to detect, diagnose and resolve network and application problems to reduce overall network troubleshooting time. Metrics that reveal details of the health of market trading services can be best quantified through deep packet inspection.

Tools that enable this deep forensics troubleshooting include those manufactured by Endace, NetScout, Corvil and Packets2Disk. With these devices, financial institution IT staff can track transaction latency and ensure, among other concerns, that performance goals and Service Level Agreements (SLA) are being met for all trades and transactions.

Figure 2 depicts a typical example of a network monitoring implementation at a financial institution. This system monitors latency, trade validation and data recording. Multiple TAP and SPAN inputs are aggregated and directed to monitoring and analysis tools. Additionally, this monitoring system can be used for security practices including IDS/IPS, troubleshooting, trend analysis, diagnostics and SLA monitoring. It is common for network monitoring implementations to vary based on redundancy, scale, and other factors.

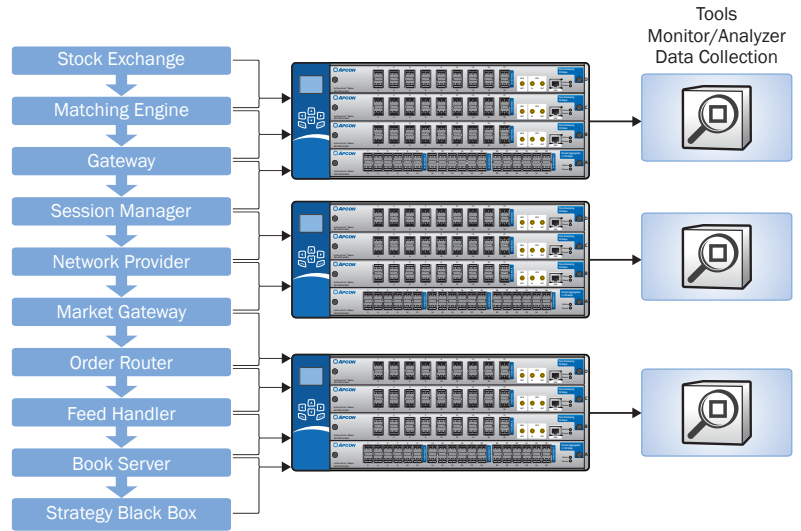


Figure 1 – Traffic is monitored and timed at various points in the system and measured for latency.

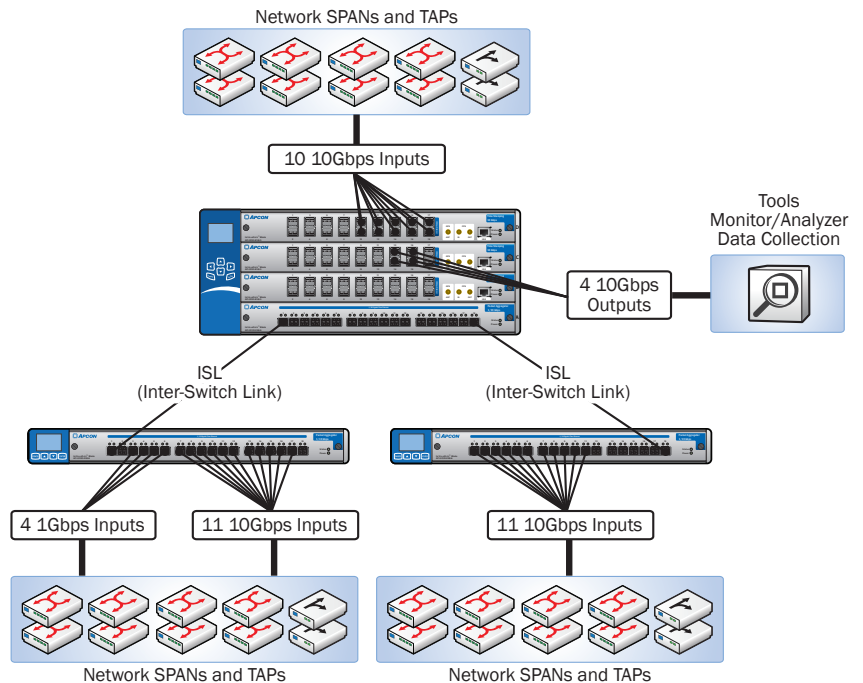
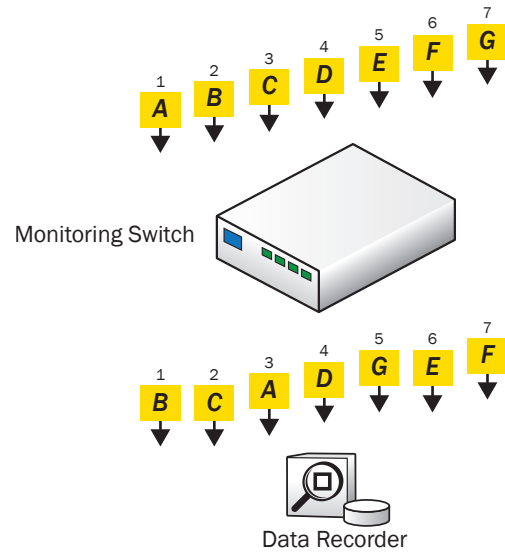


Figure 2 – In a complex network monitoring landscape, latency in the transition from data sources to tools can vary based on the path traversed. Time stamping each packet as it enters the monitoring system allows more precise measurement of production network latency.

Problem: Unpredictable Latency

Beyond expected service levels and the usual realm of network health, network latency becomes a primary focus in stock and other financial instrument trading environments because of the “first come, first served” nature of a real-time market. Remote traders, whether human or programmed, demand extremely low-latency solutions to get their trade orders registered on the trading server before their competition can do so. So the primary concern of the trader is a low latency loop from the trader’s computer to the stock exchange server and back again. In this environment, microsecond-level latencies make or break a solution.

Yet here again, network monitoring switches do not sit in-line with the production network, but rather are used to route data streams from the production network to the monitoring devices. However, latency introduced within the monitoring subsystem can skew the transactional latency measurement. If certain channels in the network monitoring switch introduce more latency than others, or if the latency is unpredictable from moment to moment, true latency cannot be specified.



Solution: Time Stamping

If a network monitoring system cannot demonstrate consistent latencies under all conditions, the lower-cost solution is simply to time stamp each packet as it enters the monitoring system. This entry typically happens when data from a TAP or SPAN port on the production network is directed to a network monitoring switch for distribution to monitoring tools.

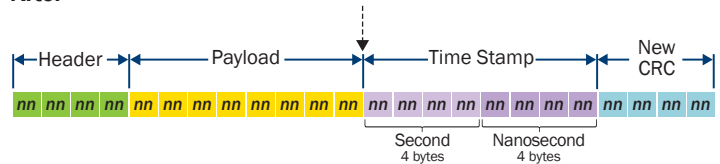
Time stamping each packet avoids the network monitoring system latency problem because the time stamp is applied as the packet enters the network monitoring switch. With the time stamp on the packet, any latency introduced by the network monitoring switch can be precisely subtracted out.

Time stamping is not difficult to perform, and is typically implemented by adding an 8-byte stamp to the end of the packet just before the CRC checksum, then calculating a new checksum for the packet. The first four bytes indicate the number of seconds since 12:00 AM, January 1, 1970. The second set of four bytes indicates ingress time to 3.2 nanoseconds within the second defined by the first set of

Before



After



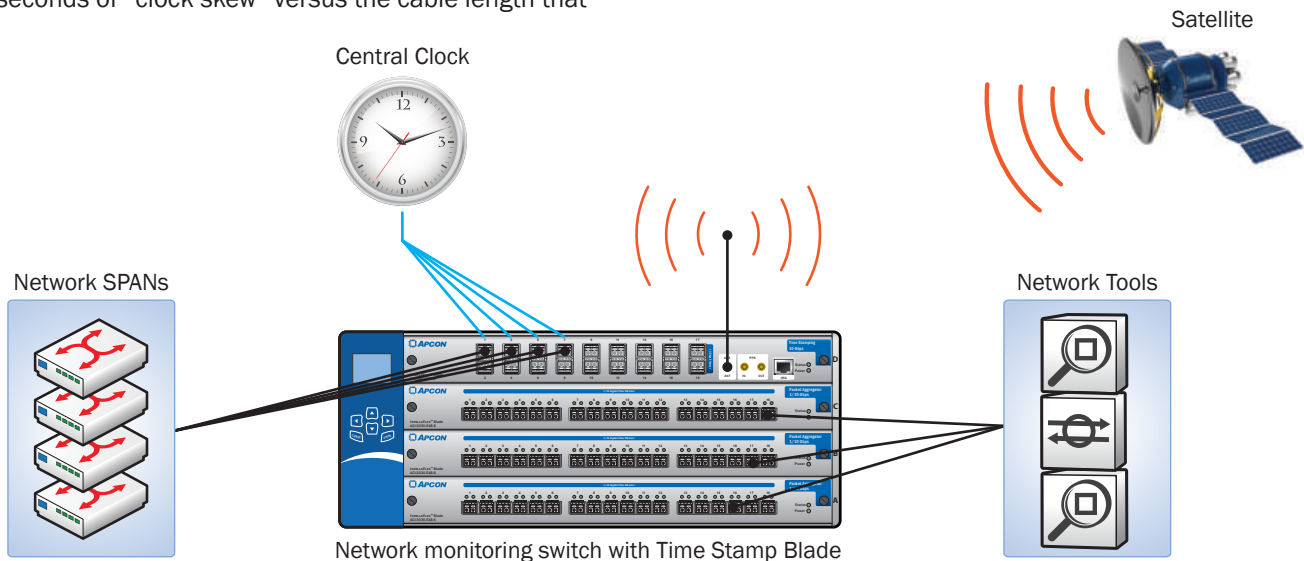
bytes. But while precision is necessary, it is not sufficient. A meaningful time stamp requires the network monitoring switch to maintain a clock that is synchronized with all other devices in the latency monitoring path. Examples include the clock on the primary network links as well as other network monitoring switches in remote locations, such as a co-location or Exchange from which the production network traffic needs to be correlated.

White Paper – Latency Monitoring for Financial Trading Markets

Obtaining an Accurate Clock Signal

There are several ways of obtaining an accurate clock signal that is synchronized on various nodes of a network. The most common method is to set up a receiver for a GPS signal, as these are synchronized worldwide to coordinated universal time (UTC). However, because the latencies involved are at the microsecond level, even the length of the antenna cable from the network switch to the GPS antenna introduces a few nanoseconds of “clock skew” versus the cable length that

might connect another node on the network. However, this skew can be corrected out when the physical layout of the data center is known and accounted for. Once a single node on the network has reliable access to this time (or even an arbitrary clock maintained locally), the network can use Inter-Range Instrumentation Group (IRIG) time code signals to keep the entire data center synchronized.



Conclusion

It is impossible to completely eliminate network latency. The physics of signals traversing the equipment and cables that connect a network will always introduce some latency. When time is measured at the nanosecond level, even short cable

distances introduce measurable latency. However, by applying a standard and precise time stamp on traffic measured at several points in the production network, the true latency of a trade can be accurately measured.



Contact Us

Please email sales@apcon.com if you have any questions

ABOUT APCON

APCON develops innovative, scalable technology solutions to enhance network monitoring, support IT traffic analysis, and streamline IT network management and security. APCON is the industry leader for state-of-the-art IT data aggregation, filtering, and network switching products, as well as leading-edge management-

software support. Organizations in over 50 countries depend on APCON network infrastructure solutions. Customers include Global Fortune 500 companies, banks and financial services institutions, telecommunication service providers, government and military, and computer equipment manufacturers.