

Low Latency – How Low Can You Go?

Low latency has always been an important consideration in telecom networks for voice, video, and data, but recent changes in applications within many industry sectors have brought low latency right to the forefront of the industry.

The finance industry and algorithmic trading in particular, or algo-trading as it is known, is a commonly quoted example. Here latency is critical, and to quote Information Week magazine, “A 1-millisecond advantage in trading applications can be worth \$100 million a year to a major brokerage firm.” This drives a huge focus on all aspects of latency, including the communications systems between the brokerage firm and the exchange.

However, while the finance industry is spending a lot of money on low-latency services between key locations such as New York and Chicago or London and Frankfurt, this is actually only a small part of the wider telecom industry. Many other industries are also now driving lower and lower latency in their networks, such as for cloud computing and video services.

Also, as mobile operators start to roll out 5G services, latency in the xHaul mobile transport network, especially the demanding fronthaul domain, becomes more and more important in order to reach the stringent 5G requirements required for the new class of ultra-reliable low-latency services.

This white paper will address the drivers behind the recent rush to low-latency solutions and networks and will consider how network operators can remove as much latency as possible from their networks as they also race to zero latency.

Background and Drivers

Latency has always been an important consideration in telecom networks. In voice networks, latency must be low enough that the delay in speech is not detectable and does not cause problems with conversation. Here the latency is generated by the voice switches, multiplexers and transmission systems, and copper and fiber plant. Transmission systems (the main topic of this paper) add only a small proportion of the overall latency, and therefore latency was traditionally not a large consideration with these networks as long as it was good enough.

Latency in Data Networks

In data networks, low latency has been seen as an advantage, but until recently it hasn't been a top priority in most cases, as long as the latency of a particular solution wasn't excessive. In most cases, the latency needed to be low enough that the data protocol functioned correctly.

A good example is Fibre Channel, where the throughput drops rapidly once the total latency reaches the point that handshaking between the two switches is not quick enough, a phenomenon known as droop. This is determined by the number of buffer credits within the switch and the latency of the link between them, which is largely generated by the fiber itself. So as long as the latency of the transmission system was not going to push the performance of a link into the area where droop is a problem, then it was normally deemed to be good enough.

Therefore, there has always been the need within telecommunications systems to ensure that latency is low enough that it minimizes the impact on the voice or data traffic being carried, but there has not been a specific requirement to drive latency as low as absolutely possible until more recently.

Applications Requiring Low Latency – Algorithmic Trading and Cloud Computing

Latency has rapidly become much more important in data networks. New applications in many vertical markets are requiring lower and lower latency.

The most widely used example of recent changes in applications used to demonstrate this is the finance industry and the move to high-frequency trading and algorithmic trading. In these applications, latency is absolutely critical as there is no second place in the race for a trade. Latency in this application comes from many areas – servers, software, and transmission – and those with an interest in low latency spend a huge amount of time and money driving as much latency as they can from every possible source.

Beyond the financial services industry, many other organizations are also now considering low latency a much higher priority. For example, services such as cloud computing are now mainstream and are considered by many to be the next big change in the way fixed telecoms networks will operate. Cloud computing includes business applications like Salesforce.com and consolidated email services for organizations that are geographically widespread. Video distribution and content delivery are becoming big cloud industries. Some of these services require low latency and others, such as email, do not, but overall this shift in services requires operators that are connecting facilities like data centers together to really look at the latency of the route and the systems used and to take corrective action if necessary. Most new installations in these applications consider low latency essential in order to deliver good quality of service.



Video distribution and content delivery are examples of applications where low latency is crucial



5G puts low latency on the agenda for the design of mobile networks

Latency in Mobile Networks

We also need to consider services over mobile infrastructure. Today latency is at the forefront of network designers' minds as the industry migrates to 5G, which adds new classes of low-latency services over those previously available in 4G. To address the ultra-low latency requirements of these new 5G services, network operators must consider every aspect of latency within the network and consider migration to a multi-access edge compute (MEC) environment.

Due to the distributed nature of a 5G transport network, with potentially different locations for the radio unit (RU), distributed unit (DU), and centralized unit (CU) components that perform radio and service processing, the underlying transport network is considerably more Internet Protocol (IP) centric than was seen in 4G transport networks. These IP devices and RU, DU, and CU components still require an underlying Layer 1 or Layer 2 xHaul network with fronthaul, midhaul, and backhaul domains. Therefore, all elements within this transport network must be optimized for the lowest possible latency, and this is especially the case for the fronthaul domain, which has the most demanding transport specifications.

Differing strategies exist for the optical mobile transport networks that underpin the IP layer that supports the RU, DU, and CU traffic flows. Some are Layer 1 only and others require packet optical Layer 2 devices. Both approaches require low-latency optimized solutions.

As you can see, low latency is becoming increasingly important in Layer 1 and Layer 2 solutions. Let us now consider the sources of latency in a fiber optic network and what can be done to minimize this.

Sources of Latency

Latency in fiber optic networks comes from three main components: the fiber itself, optical components, and opto-electrical components.

Latency in Optical Fiber

Light in a vacuum travels at 299,792,458 meters per second, which equates to a latency of 3.33 microseconds per kilometer (km) of path length. Light travels slower in fiber due to the fiber's refractive index, which increases the latency to approximately 5 microseconds per km. So, while we are using the current generation of optical fibers, there is a limit to how low we can drive latency – take the shortest possible route and multiply this by 5 microseconds per km. A 50 km link would therefore have a fiber latency of 250 microseconds, a 200 km link would have a fiber latency of 1 millisecond, and a 1,000 km link would have a fiber latency of 5 milliseconds.

This is the practical lower limit of latency that is achievable if it were possible to remove all other sources of latency. However, fiber is not always routed along the most direct path between two locations, and the cost of rerouting fiber can be very high. Some operators have built new low-latency fiber routes between key financial centers, and have also employed low-latency systems to run over these links. This is expensive and is likely to only be feasible on the main financial services (algo-trading) routes where the willingness to pay is high enough to support the business case. In most other cases, the fiber route and associated latency will be fixed due to complexity and cost.

Latency in Optical Components

The vast majority of the latency introduced by optical transmission systems is in the form of dispersion compensating fiber (DCF). This is only used in long-distance networks, so it is not a consideration in, for example, an 80 km data center interconnect project. DCF is used to compensate for dispersion of the optical signal. This is caused by the speed of light varying slightly for each wavelength, and even though WDM wavelengths are very tightly spaced, the pulse of light will spread out as it travels down the fiber because some components of the pulse will travel faster than others. Eventually this spreading reaches the point at which the pulses start to get too close together and cause problems for the receiver, and ultimately bit errors in the system. To compensate for this dispersion, WDM systems use DCF in amplifier sites. DCF is essentially fiber with the opposite dispersion characteristics, so a spool of this added at the amplifier site can bring the pulse back together again. This extra fiber adds to the optical power calculations, requires more amplification in the network, and of course adds more latency. A typical long-distance network requires DCF on approximately 20 to 25% of the overall fiber length, and therefore this DCF adds 20 to 25% to the latency of the fiber, which could be a few milliseconds on long-haul links.

Innovations in fiber Bragg grating (FBG) technology have enabled the development of the dispersion compensation module (DCM). A DCM also compensates for dispersion over a longer-reach network but does not use a long spool of fiber and therefore effectively removes all the additional latency that DCF-based networks impose. As both DCF and DCM units are directly connected to the optical path, these should either be designed in for new low-latency routes or swapped over during planned maintenance windows on existing routes where lower latency is now required.

The only other optical components that require discussion here are the optical amplifiers. These erbium-doped fiber amplifier (EDFA) optical amplifiers enable WDM systems to work as they amplify the complete optical spectrum and remove the need to amplify each individual channel separately.

They also remove the requirement of optical-electrical-optical (O-E-O) conversion, which is highly beneficial from a low-latency perspective. They operate by using a spool of a few tens of meters (m) of erbium-doped optical fiber and pump lasers. Due to the optical characteristics of this special fiber, optical power is transferred from the pump lasers to the optical signal as it passes through the fiber, leading to the amplification of the signal. But from a latency perspective, these amplifiers contain a small spool of optical fiber that we should consider if an operator is really looking to drive every possible source of latency out of a system. Of course, on a per-amplifier basis this latency is very small. But a long-haul system will have many amplifiers, perhaps 10 to 15 in a link, and assuming 30 m per amplifier, this soon increases to 450 m (with a latency of approximately 2.25 microseconds) in a 15-amplifier system, which could be significant to some operators, especially those in the financial sector.

One approach to address this additional latency is to use Raman amplifiers instead. Raman amplifiers utilize a different optical characteristic to amplify the optical signal. High-power pump lasers use the outside plant fiber itself as the amplification medium, and transfer power from the pump lasers to the optical signals to amplify the system. Here there are no additional spools of optical fiber and therefore no additional latency. These Raman amplifiers are more expensive than EDFAs so until now have mainly been used in addition to EDFAs to boost the amplification for systems with very long spans. However, these do provide the operator that wishes to drive every possible source of latency out of its network with an additional option.

Latency in Opto-electrical Components

First, let us consider the Layer 1 examples mentioned earlier in this document. Operators have two approaches to transporting data over optical transmission systems – transponders or muxponders. Transponders take a single optical signal and convert it from optical to electrical and back to optical again, and in the process convert the wavelength from a short-reach interoffice signal to a long-distance WDM-specific wavelength. Muxponders take multiple signals, multiplex them together into a single higher-speed signal, and then convert that to the WDM-specific wavelength. An operator will typically use transponders for higher-speed links such as 4 gigabits per second (Gb/s)/16 Gb/s Fibre Channel, 100 Gb/s Ethernet, etc., and muxponders for lower-speed services such as Gigabit Ethernet.

The latency of both transponders and muxponders varies depending on design, formatting type, etc. Muxponders typically operate in the 5 to 10 microseconds per unit range. If forward error correction (FEC) is used for long-distance systems, then this will increase the latency due to the extra processing. Transponders, however, can vary hugely in latency depending on design and functionality. The more complex transponders include functionality such as in-band management channels, which forces the unit design and latency to be very similar to a muxponder, in the 5 to 10 microsecond region, as the unit needs to combine the data and management channel signals in a similar way to a muxponder. Again, if FEC is used, then this can be even higher. Coherent processing used in all optics that operate at 100 Gb/s and above also adds an additional latency element that needs to be managed. Typically, higher speeds require more complex signal processing in the digital signal processor (DSP) and create higher latency. Therefore, for high-speed coherent optics, rate limiting a wavelength to 100G will reduce the latency of the network.

Some vendors, including Infinera, also have options for simpler, and often lower-cost, non-coherent transponders that do not have FEC or in-band management channels, which can operate at much lower latencies for services in the 10 to 16 Gb/s range.

The Infinera XTM Series has set the industry benchmark with the lowest stated latency of any transponder, at 4 to 10 nanoseconds for a pair of transponders (one per end of the link), which equates to approximately 1 to 2 m of fiber being added to the overall system link. The range from 4 to 10 nanoseconds is due to the varying latency over the operating range of the transponders. The higher the speed, the lower the latency, so 10 Gb/s services benefit the most from this low latency.



A few other vendors also have low-latency transponder options, but none yet have been able to get as low as Infinera. Many others are stuck in the millisecond range, which is 1,000 times higher latency, due to the formatting structures used.

Infinera packet optical transport switches are built with low-latency design

Infinera packet optical transport switches are optimized for the aggregation and transport of Ethernet traffic, and thus latency for these units is low, less than 2 microseconds. Although this is significantly more than with the Layer 1 transponders, there is a network architecture angle that must be considered here too.

In Layer 1 we consider point-to-point wavelengths with a transponder/muxponder at each end and optical components in between. In Layer 2, we have a network of Layer 2-capable devices in which the traffic moves from one to the next until it reaches its destination. For a 5G mobile xHaul network, this could entail four or five devices into the core and then four or five back again, or it could be much higher depending on the network architecture. If we assume five Layer 2 devices between a 5G-enabled device being used for real-time gaming and the core, then the data associated with a user action will hop through five devices on the way to the core and the response will hop back through the same number of devices for a total of 10 in this example. This means the two- to threefold improvement in performance equates to a latency difference of 10×2 (20) microseconds compared with $10 \times 5+$ (50+) microseconds, and thus a savings of 30 microseconds, or the equivalent of 6 km of fiber.

One further factor to consider when assessing Layer 1 or Layer 2 network design options is which of the two can provide the lowest possible latency. It may be counterintuitive, but Layer 2 aggregation devices designed to aggregate 10 Gb/s services into a 100 Gb/s link can actually provide lower latency than comparable Layer 1 devices if they are designed with low latency in mind. This is because Layer 2 uses a “bit stuffing” mechanism to insert the 10 Gb/s traffic into the higher-speed 100 Gb/s stream, whereas Layer 1 Optical Transport Network (OTN)-based devices use a hierarchical aggregation mechanism that has higher latency. For this reason, some operators deploy Layer 2 devices for low-latency 10 Gb/s to 100 Gb/s aggregation and run fully uncontended Layer 2 services that behave in the same way as Layer 1 services when the lowest possible latency is required.

An Operator’s Options

So, with a toolkit of low-latency solutions, what should an operator do when looking to provide a low-latency service? From the discussion above, it is clear that the fiber route has by far the biggest impact on latency, and if the operator has two options, then it should choose the shorter.

The next biggest impact an operator can make for long-distance networks is to use DCM-based dispersion compensation rather than DCF, which could reduce latency by up to 20%. It is therefore equivalent to reducing the route length by the same amount, but probably at a much lower cost than digging new trenches and pulling new fiber routes.

To drive latency lower in both short-haul and long-haul networks, the operator should use an optical transport solution that offers ultra-low-latency transponders. These can reduce the latency associated with O-E-O conversion from milliseconds to nanoseconds. This has a similar effect to shaving off 1 or 2 km from the route distance.

Finally, for those that really want to go as low as possible, the small amount of remaining latency within the optical amplifiers can also be removed by swapping them from EDFA to Raman amplifiers.

Below are two examples of low latency: one over a short distance and another over a longer distance.

For operators deploying Layer 2-based transport networks that wish to consider low latency requirements, then an important option is low-latency Layer 2 transport Ethernet solutions. As these networks grow, the bigger the potential savings can be.

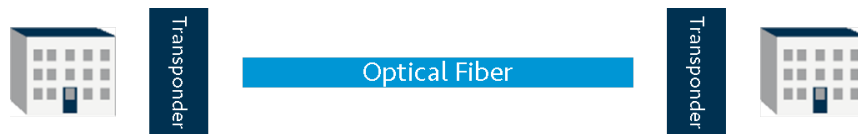


Figure 1: Low latency options in an example of two data centers 20 km apart

Example 1: Two data centers 20 km apart

Fiber latency =	20 x 5 μs =	100 μs
Transponder latency =	2 x 5 μs =	10 μs
Total latency =		110 μs

Low latency options:

Replace transponders with ultra-low-latency transponders with 4 ns latency per pair.

This effectively removes transponder latency for a 9% savings and a total reduction of 10 μs.

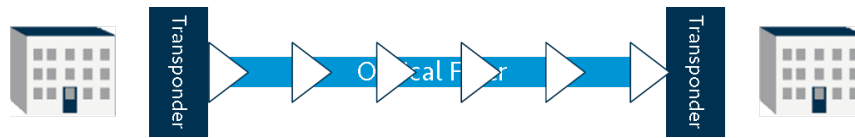


Figure 2: Low latency options in an example of two data centers 400 km apart

Example 2: Two data centers 400 km apart – five spans of 80 km

Fiber latency =	$400 \times 5 \mu\text{s} =$	$2,000 \mu\text{s}$
DFC latency =	$20\% =$	$400 \mu\text{s}$
Transponder latency =	$2 \times 5 \mu\text{s} =$	$10 \mu\text{s}$
Amplifier latency =	$6 \times 30 \text{ m} =$	$0.9 \mu\text{s}$
Total latency =		$2,410.9 \mu\text{s}$

Low latency options:

Replace DCF with DCM, transponders with ultra-low-latency transponders, and EDFA with Raman amplifiers.

This effectively removes all DCF, transponder, and amplifier latency for a 17% savings and a total reduction of 411 μs .

For operators deploying Layer 2-based transport networks that wish to consider low latency requirements, then an important option is low-latency Layer 2 transport Ethernet solutions. As these networks grow, the bigger the potential savings can be.

Conclusion

Low latency is a real concern in many network scenarios. Some, such as telecom services to the financial services industry, can demand a substantial pricing premium if they can provide the end customer with an advantage in low latency. Other industry sectors, such as data center interconnect, will require more of a focus on low latency as a basic feature to ensure the facility owner is strongly positioned to serve customers with low latency demands. There are limits to how low latency can go until we can change the laws of the physics of light in fiber, but there is a lot a network operator can do with both Layer 1 and Layer 2 transport solutions to ensure that the latency on any route is as low as physically possible.

Any operator looking to deploy low-latency networks should ensure that it has a toolbox with all available low latency options. Optical fiber latency can only be reduced by taking a new route, which can be very expensive but also highly beneficial. Latency in optical components can be greatly reduced in long-distance networks using DCM and Raman components. Finally, latency in opto-electrical components can be further reduced, offering the operator a competitive edge as this area varies from systems vendor to systems vendor and hence deployment to deployment.

Infinera offers operators a full toolbox of low-latency components, ultra-low-latency Layer 1 transponders, and Layer 2 transport Ethernet solutions with industry-leading low latency.