# White Paper

# Connectivity in the Virtualized Datacenter: How to Ensure Next-Generation Services

Virtualization is changing the economics of the datacenter by making computing resources more flexible and efficient. To capitalize on these benefits, most datacenters must make fundamental changes to their Top-of-Rack or End-of-Row network architectures.

A virtualization strategy is more likely to succeed if the performance of the datacenter network can be assured before services are enabled, and if the cable plant foundation for the Top of Rack or End-of-Row networks is certified to meet the needs of today and tomorrow.

This white paper compares the new choices for datacenter connectivity and describes test methods that help to ensure that your network provides a solid base for deploying virtualized services.

May 2010

BROADCOM.
Connecting everything

FLUKE networks

# Background: Virtualization

The virtualization concept is expanding into many aspects of information technology. Servers, switches, storage, networking, and clients are all on a virtualization road map. But the virtualization movement is rooted in servers, and server virtualization will have the most profound impact on datacenter networks.

Server virtualization counters the trend to use purpose-specific appliances as servers. That trend succeeded because appliances are effective and easy to deploy. But in large numbers, appliances are inefficient. Because appliances are optimized for one function, many are dead-end devices subject to forklift upgrades. Perhaps worse, the proliferation of appliances has led to growing management burdens, sprawling datacenter networks, and high thermoelectric loads.

The solution to the underutilization problem is virtualization, and the essence of virtualization is an abstraction layer of software called the Hypervisor. The Hypervisor sits between hardware and the operating system. Virtualization allows multiple operating systems and applications to cohabitate on a physical computing platform. The graphic below shows virtualization supporting three logical servers on one platform. Server virtualization, especially when coupled with blade technology, increases computing and storage density while making IT assets more flexible.
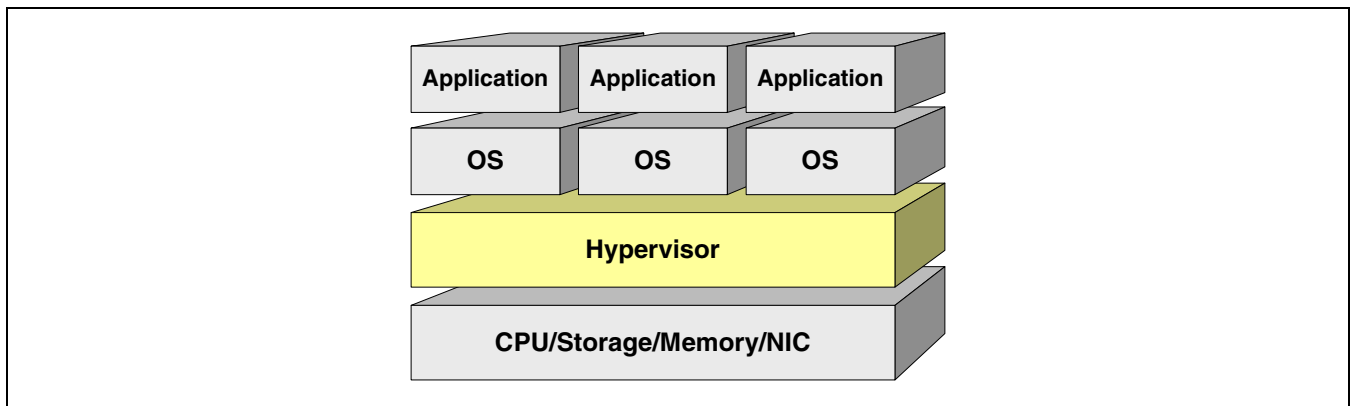


**Figure 1: Virtualized Server with Hypervisor Layer**

The depiction of a virtualized server in Figure 1 implies a traditional client-server relationship with the user. There are nuances of virtualization with intriguing separation of computation from the user interface. This white paper focuses on server virtualization and its effects inside the datacenter.

Virtualized servers will support the full array of business applications, multimedia applications, storage, and back-office services. Adoption of virtualization is accelerating. In February 2010, Microsoft[®] reported that 20% of shipping servers were virtualized. VMware[®], the leader in Hypervisor products, reported that customers have virtualized 25% of their servers. No one knows how large virtualization will grow, but vendors and analysts forecast a 3x increase within five years.

Virtualization is not limited to large enterprises. New products that unify computing, networking, and storage are designed to stretch the benefits of virtualization to the cloud and to medium-sized organizations.

# Virtualized Server Platform

The Broadcom NetXtreme® I and NetXtreme II® high-speed controller families offer advanced network virtualization features designed to help customers build and deploy products for a wide range of application profiles. Broadcom's leading market share of over 70% in both the 1G and 10G network controller market segments is a result of successful network virtualization deployments in leading growth markets. These markets include high-performance Web 2.0 application servers, low-latency financial trading systems, and high-density cloud computing clusters.

Virtualization of network controllers allows users to consolidate their networking hardware resources and run multiple virtual machines concurrently on consolidated hardware. Virtualization also provides the user a rich set of features such as I/O sharing, consolidation, isolation and migration, and simplified management with provisions for teaming and failover.

Broadcom has collaborated with vendors of various VM products, including VMWare Vsphere, Microsoft Hyper-V™, Redhat® KVM and Citrix® Xen, to provide a rich set of network virtualization functionality necessary for adoption of 10GBASE-T in datacenter and cloud enviornments. Functionality and features listed below remove virtualization bottlenecks and improve system performance by providing additional features:

- **Stateless offloads**—Broadcom Ethernet network controllers support stateless offloads such as:
  - CP Checksum Offload (CSO), which enables network adapters to compute TCP checksum on transmit and receive.
  - TCP Large Send Offload (LSO), which allows TCP layer to build a TCP message up to 64 KB long and send it in one call down the stack through IP and the Ethernet device driver, saving the host CPU from having to compute the checksum in a virtual environment.
- **Jumbo frame support**—In virtual environments, the use of jumbo frame saves CPU utilization due to interrupt reduction and increases throughput by allowing the system to concentrate on the data in the frames, instead of the frames around the data.
- **Multiple queue support**—A Broadcom Transport Queue Manager utilizes on-chip queuing technology with VMWare ESX Netqueue. The additional overhead from route lookup, data copy, and filtering tasks are off-loaded to a network adapter, where the transport queue manager can transmit packets from multiple queues and steer the receive packets into multiple queues. By dedicating Tx/Rx queue pair to a VM, the network adapter can provide DMA to and from the VM's memory, and the VSwitch only processes the control plane operation.
- **Network offload**—VMware Static, VMDirectPath or Fixed Pass Through (FPT) enables Broadcom NetXtreme II C-NICs to be completely dedicated to high performance VM. FPT utilizes AMD IOMMU or Intel VT-D to provide DMA data I/O from the physical device to the VM, bypassing the virtualization layer, thereby yielding complete physical access of the dedicated Broadcom CNIC to the VM.
- **Storage Offload**—iSCSI HBA functionality enables on-chip processing of the iSCSI protocol (as well as TCP and IP protocols), which frees up host CPU resources at 10 Gbps line rates over a single Ethernet port.

- **iSCSI Boot**—The server can boot an Operating System (OS) over a SAN, completely eliminating the need for local disk storage (which is the number one source of failures in computer systems). In addition to enhanced system reliability, the use of diskless servers simplifies the IT administrator's workload by centralizing the creation, distribution, and maintenance of server images, reducing the overall need for storage capacity through increased disk capacity utilization and adding increased data redundancy through the use of data mirroring and replication.

# Datacenter Network Design

Virtualization affects datacenter networks in two important ways. The first is the demand for bandwidth. Consolidated server platforms require higher bit-rate connections to support multiple processes. Storage further adds to the demand for bandwidth. Multimedia content and accumulated application data inflate storage overhead, and more data directly drives the need for high-speed access. The result is an accentuated need for a high-bandwidth datacenter network.

## Foundation Layer — Virtualized Server Platform

The impact of virtualization is rooted in the ease in which IT assets can be repurposed. Virtualization allows hardware to perform multiple functions and allows functions to be moved among hardware platforms. This flexibility for computation and storage must be accommodated by flexibility in server networking capability. Figure 2 is an example of a consolidated and virtualized server platform with advanced server networking from Broadcom® Corp.
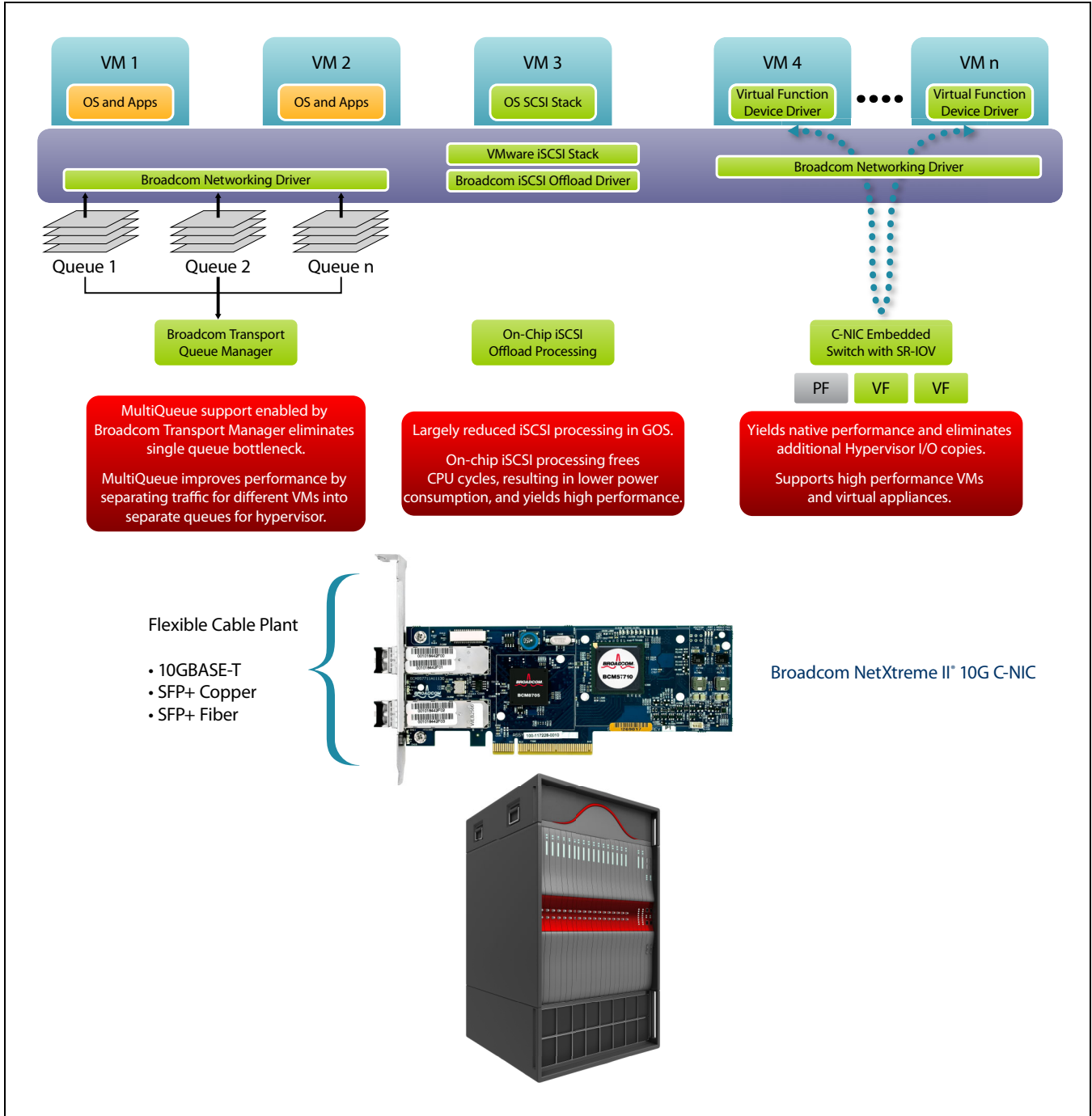
**Figure 2: Virtualized Server Platform**

The physical network must adapt to the requirements and advantages created by virtualization, specifically higher utilization and higher bandwidth. To do this, forward-thinking network professionals have implemented the End-of-Row ("EoR") or Top-of-Rack ("ToR") topologies in their datacenter networks.

## Option 1: The End-of-Row Topology

As a point of reference, Figure 3 depicts a pre-virtualized datacenter, where each asset (server, storage device, etc.) is individually linked to an Ethernet switch.
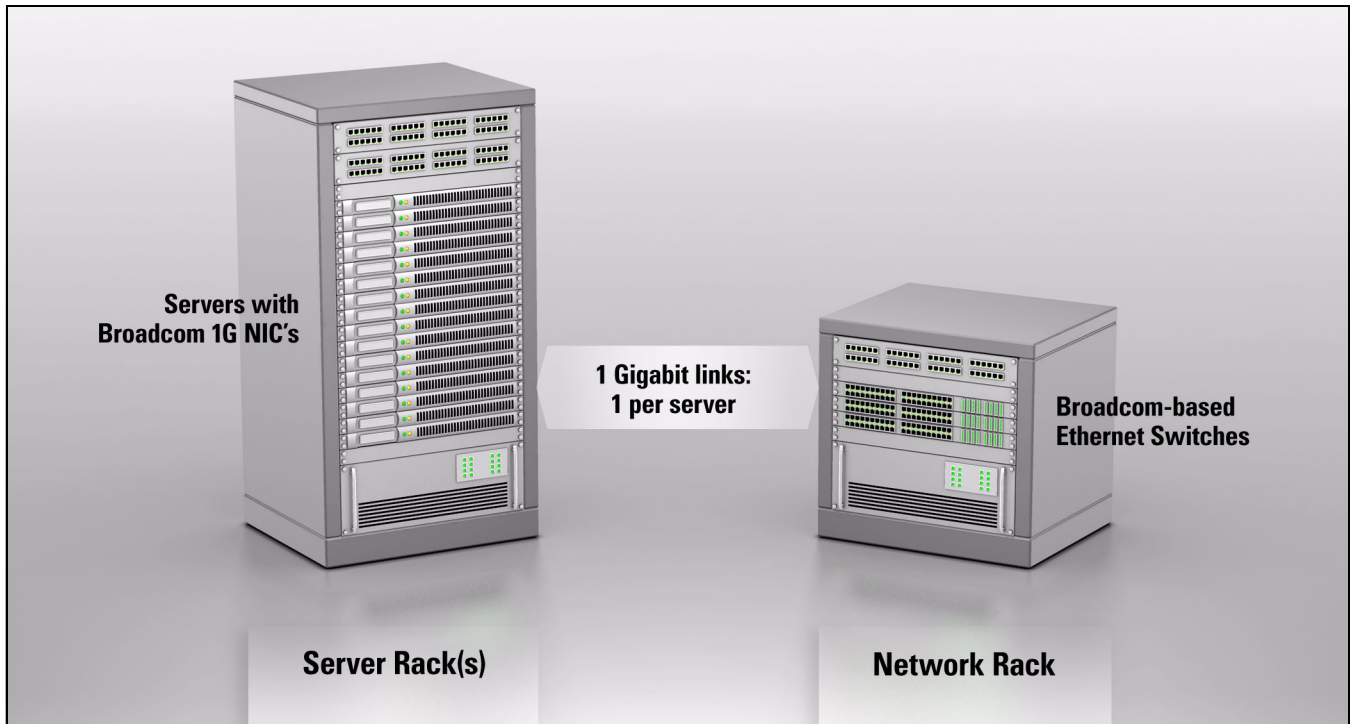


**Figure 3: Conventional Topology**

This topology uses structured cable connections that are difficult to modify. Since virtualization facilitates change, a network architecture that inhibits it is inherently problematic. The conventional topology is also dated, as having a large number of 1 Gigabit links is incongruent with consolidated servers that need fewer and faster connections. This shift from many "thin roots" to fewer "thick roots" must be supported by the network.

The End-of-Row network topology, as shown in Figure 4, addresses the shortcomings of a conventional datacenter network by dedicating an Ethernet switch to each row of equipment racks. The virtualized assets in each rack, in each row, are linked to a switch in the EoR rack. That switch also provides a trunk connection to a datacenter concentrator.
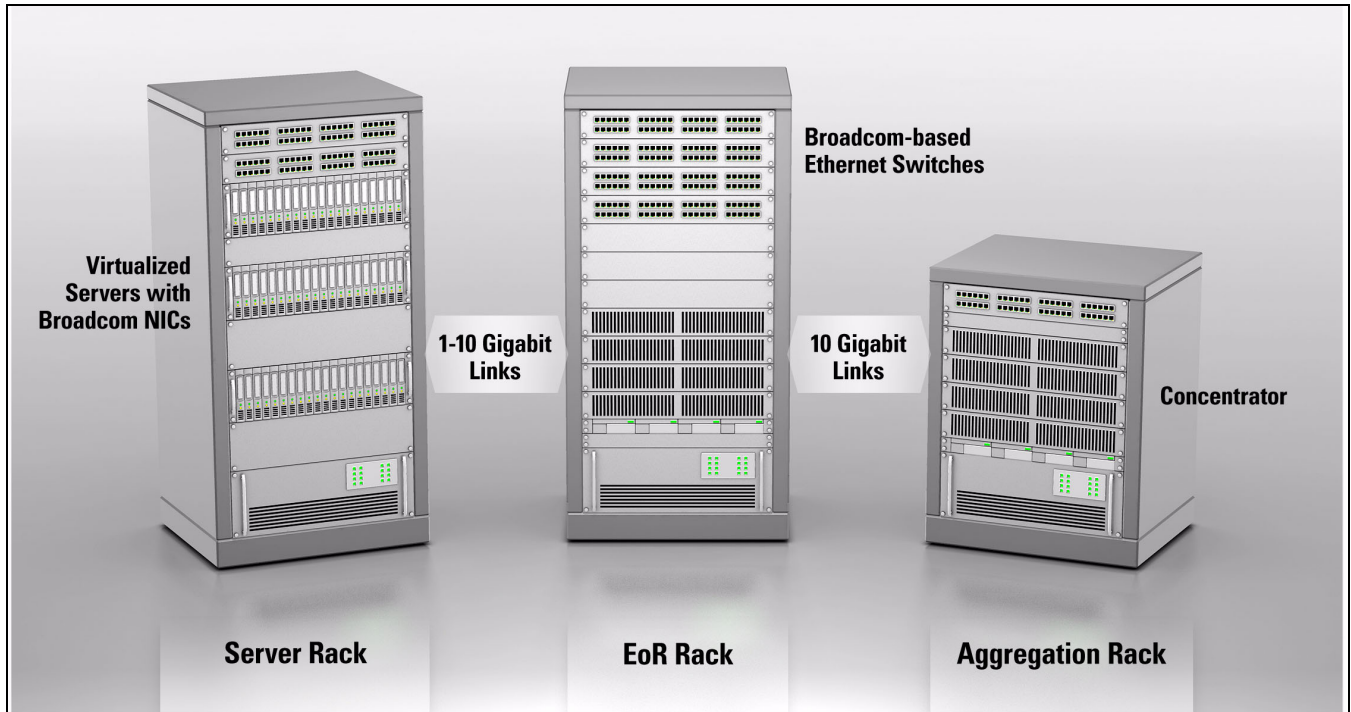
**Figure 4: End-of-Row Topology**

The EoR topology divides the switch fabric and physical connections from one tier into two, making the network more adaptable. EoR limits the length of the cables in the lower tier to the length of a row of racks. Shorter cables are generally easier to install and easier to change.

EoR topology confines the impact of asset reconfiguration to a row of racks, instead of across an entire datacenter. EoR may reuse some elements of the existing physical network, although major changes and upgrades are likely.

## Option 2: The Top-of-Rack Topology

The Top-of-Rack topology is a bigger departure from the conventional architecture. It dedicates an Ethernet switch to every rack of servers. The ToR switch interconnects assets in each rack and provides a trunk connection to an aggregation point in the datacenter.
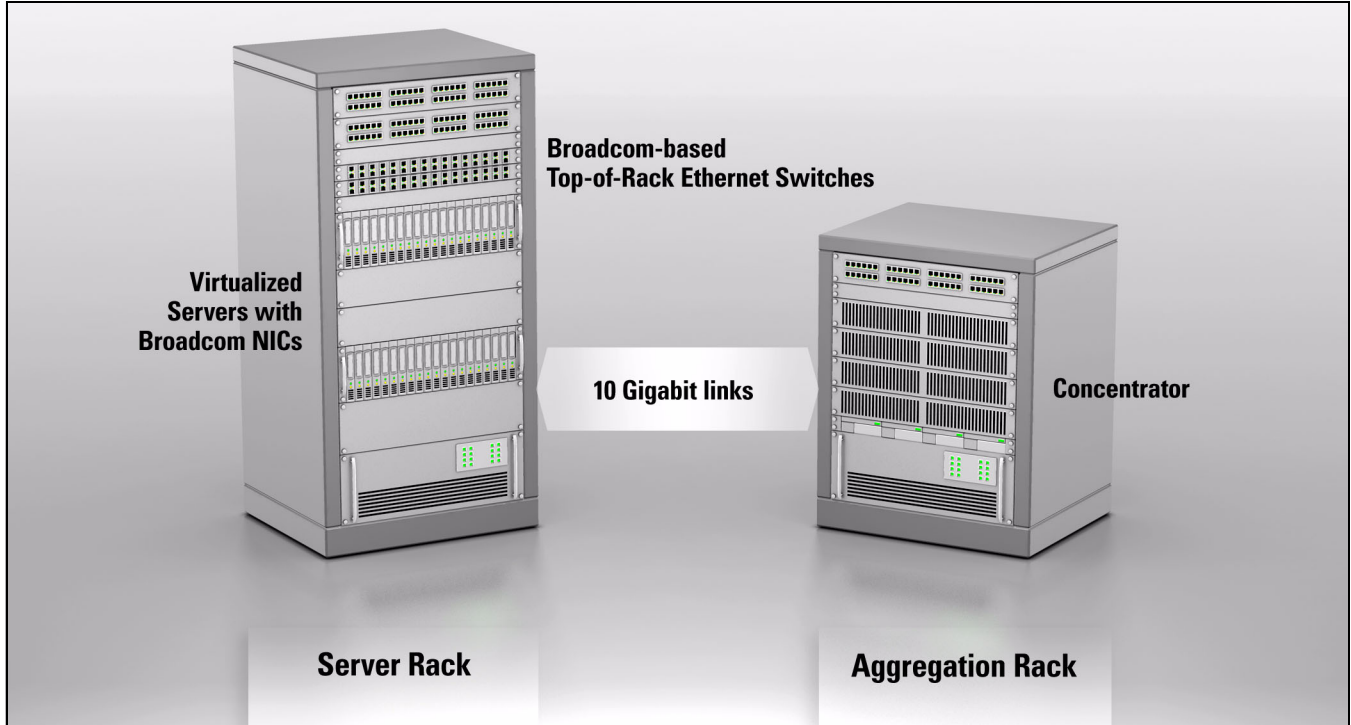


**Figure 5: Top-of-Rack Topology**

Like EoR, the ToR topology divides the switch fabric and the physical connections into two tiers. The difference is the granularity of the lower tier. Where EoR creates modularity in a row of racks, ToR creates modularity in each individual rack.

Note that the ToR design does limit a server rack to only one switch. The graphic above shows two switches in a rack: one primary and one for redundancy. If the Ethernet switches are implemented as blades, there could be even more switches in a rack.

The Top of Rack versus End-of-Row discussion is complex. Many Layer 2 and Layer 3 networking issues must be considered. Some of the relative advantages of each topology are:

| Top of Rack | End-of-Row |
|---|---|
| Less structured cable | Less disruptive to infrastructure |
| Easier to change/expand | Fewer switches and trunk connections |
| Most modular | Easier to manage/support |

# Virtualization Topologies and Cabling

Ethernet switch vendors offer volumes of information about designing networks for virtualized services. Whatever the topology you choose, three things will be true regarding the connectivity infrastructure:

- Cabling will change.
- Cabling will outlive many computing, storage and networking assets.
- Failures of datacenter network cabling can impact even virtualized services.

Because virtualized services require a rethink of datacenter connectivity, it is useful to recognize the categories of cable created by the End-of-Row and Top-of-Rack topologies.

**Table 1:  EoR and ToR Cables**

| Class | Use | Bit Rate | Max. Distance | EoR | Tor |
|---|---|---|---|---|---|
| Intra-Rack | Device to ToR switch | 1–10 Gbps | 5 meters | | ✔ |
| Rack-to-Rack | Device to EoR switch | 1–10 Gbps | 50 meters | ✔ | |
| Cross-Datacenter | ToR or EoR switch to concentrator | 10 Gbps | 300 meters | ✔ | ✔ |

Some device-to-switch links can suffice at 1 Gigabit today, but the myriad drivers for higher bit rates make future-proofing a wise investment. The good news is that there are many ways to support 10-Gigabit traffic. Table 2 shows the most common 10-Gigabit solutions.

**Table 2:  Common 10 Gb Switch Configurations**

| Media | Cable | Type | Max. Distance | EoR | EoR | Tor |
|---|---|---|---|---|---|---|
| Copper | Twinax | SFP+ Copper | 10m | ✔ | | |
| Copper | Twinax | 10GBASE-CX4 | 15m | ✔ | | |
| Copper | Cat 6 UTP | 10GBASE-T | 50m | ✔ | ✔ | |
| Copper | Cat 6A UTP | 10GBASE-T | 100m | ✔ | ✔ | ✔ |
| Copper | Cat 7 UTP | 10GBASE-T | 100m | | ✔ | ✔ |
| Fiber | 850/Multimode | 10GBASE-SR | 80–300m | | ✔ | ✔ |
| Fiber | 1300/Multimode | 10GBASE-LRM | 220m | | ✔ | ✔ |
| Fiber | 1310/Single Mode | 10GBASE-LR | 10 km | | | ✔ |

- **Twinax**—A shielded cable initially used for InfiniBand storage, Twinax can support high bit rates over modest distances. Twinax cables are expensive in longer configurations and less pliable than twisted-pair copper. SFP+ Twinax is terminated with compact SFP+ modules instead of large CX4 connectors.

- **Category 6 Unshielded Twisted Pair**—The IEEE approved the 10GBASE-T standard in 2006 for 10 Gigabit over twisted-pair cable, including shorter Cat 6 UTP links. 10GBASE-T extends the commonality and low cost of UTP to 10 Gigabit. It is backward compatible with Cat 5e/100 Mbps and supports all Ethernet features. Cat 6 patchcords for Intra-Rack applications can be purchased preterminated or connectorized on-site. Cat 6 for Rack-to-Rack is connectorized on-site.

- **Category 6A Unshielded Twisted Pair**—Augmented Cat 6 was officially recognized in 2008 in the TIA 568-B.2-10 standard. It specified cable performance up to 500 MHz and set requirements for Alien Crosstalk performance. Category 6A cables are thicker and more expensive than Cat 6, but support 10 Gigabit traffic up to 100 meters. Some Cat 6A cables have shielding around each wire pair to control Alien Crosstalk. Cat 6A patchcords can be purchased preterminated or connectorized on-site. Cat 6A for Rack-to-Rack or Cross-Datacenter is connectorized on-site.

- **Category 7 Shielded Twisted Pair**—Although it is more properly called "Class F", as referenced in the ISO/IEC 11801 standard, Category 7 STP predated Cat 6A and was designed specifically for high bit rate applications. Cat 7 addresses alien crosstalk with a metallic shield around each wire pair plus a shield for the entire cable. Cat 7 is the most expensive twisted-pair copper solution and requires an extra installation step to terminate the shielding. The noise immunity inherent in Cat 7 shielding gives it better performance than UTP and is theoretically more future-proof. Cat 7 cables do not use standard RJ-45 connectors, so they must be connectorized on-site.

- **10GBASE-SR**—This is a common fiber option for datacenters when coupled with 850 nm multimode transceivers. 10GBASE-SR can be implemented with SFP+ modules to conserve connector space and power. The maximum distance for 10GBASE-SR is determined by the fiber used. OM1 and OM2 fiber may support 10 Gigabits up to 80 meters. Newer OM3 cable can support 10 Gigabits up to 300 meters.

- **10GBASE-LRM**—LRM was conceived to drive 10 Gigabit traffic using 1310 nm lasers on multimode fiber over long distances. The IEEE 802.3aq standard was approved in 2006. Cables are duplex FDDI fiber with SC or LC connectors. Fiber for patchcords can be purchased with connectors or can have connectors spliced on-site. Fibers for longer runs are spliced with connectors on-site.

- **10GBASE-LR**—LR is the most common choice where single mode fiber is used. Cables are duplex fibers with SC or LC connectors. Distance and cost are directly related, so 10GBASE-LR modules are very expensive. The wealth of new alternatives makes 10GBASE-LR an extravagant choice for the virtualized datacenter, but it is still a viable choice for campus networks.

# The Emergence of 10GBASE-T

While the 10GBASE-T standard was ratified in 2006, market adoption lagged as power consumption, heat, and density issues slowed implementation in switches and network interface cards. But recent advances in silicon solved those problems and the ramp-up for 10GBASE-T is accelerating rapidly.

Because 10GBASE-T runs on CAT6/6A/7 copper, it is employable in Intra-Rack, Rack-to-Rack, and Cross-Datacenter applications. 10GBASE-T is unique in that it works well in both ToR and EoR topologies. It supports auto negotiation of transmission speed, allowing gradual migration to 10 Gigabit speeds. 10GBASE-T has not yet replaced 1GBASE-T as the common denominator for datacenter network connectivity, but economics and flexibility make that a likely outcome.

# Assuring Success—The Cable Testing Reference Model and Certification

Topology and economics are the prime determinants of connectivity, but any option requires a reliable cable infrastructure. While cable is a passive technology and vendors provide warranties, you should bear in mind that:

- Warranties end.
- A manufacturer's warranty probably excludes installation labor.
- The failure of any cable could mean the failure of services.
- Repairing a problem is always more expensive than preventing it.

To ensure successful delivery of virtualized services, the cable infrastructure should be subject to certification testing. Certification is a rigorous assessment, which is performed before the network is put into service, of connectors, installation workmanship, and cables. The result of a certification test is compared to industry standards, with the result being a "Pass" or "Fail" grade for every link. A link that passes certification meets the defined performance specifications. Links that do not pass are repaired, usually at the expense of the vendor or network installer.

Certification traditionally focuses on structured cabling, but if the Top-of-Rack topology is employed, it may include patchcords used for Intra-Rack cabling.

## Certifying Copper

- **Twinax**—There is no test standard for Twinax cable so it cannot be formally certified. The network user has no choice but to rely on the manufacturer's warranty and/or treat Twinax cables as disposable goods.

- **Unshielded Twisted Pair**—Certifying Cat 6/6AUTP is done in two test steps: Channel and Alien Crosstalk. All UTP cables, be they for Intra-Rack, Rack-to-Rack, or Cross-Datacenter applications, should be subject to the Channel test. The Channel test certifies 13 parameters defined by the TIA/EIA-568-B and ISO 11801 standards.

  UTP used for Rack-to-Rack and Cross-Datacenter links should also be tested for Alien Crosstalk. This is a sample test: the ISO/IEC standards suggest 1% of the links or 5 links, whichever is greater. Alien Crosstalk testing ensures that cross-coupling will not affect network performance. Existing Cat 6 cabling up to 50 meters in length can even be recertified for 10GBASE-T using criteria defined in TIA Telecommunication Systems Bulletin (TSB) 155.

- **Shielded Twisted Pair**—Certifying Cat 7 (also called Class F) is a Channel test. The ISO 11801 standard specifies the parameter limits and frequency range for this test. Since Cat 7 cable does not use RJ-45 connectors, the cable tester must have an adapter that is compatible with each vendor's unique connector.

The configuration to certify copper links is shown in Figure 6. A certification tester and its remote unit are attached to both ends of the link to perform all aspects of certification.
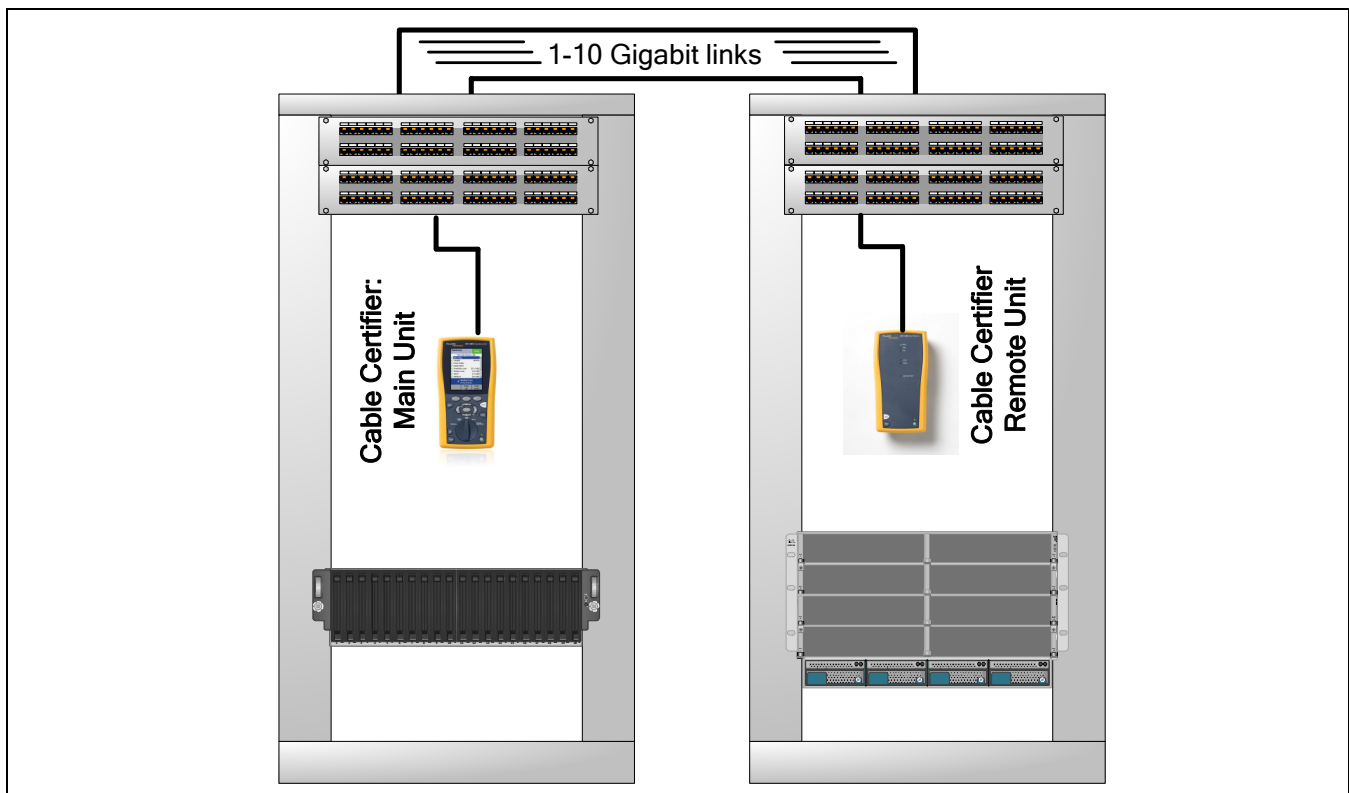


**Figure 6: Copper Certification**

It is important to document the results of certification tests. A report that is complete, accurate, and easy to understand serves as evidence that the infrastructure meets the standards and the specifications prescribed by datacenter management.
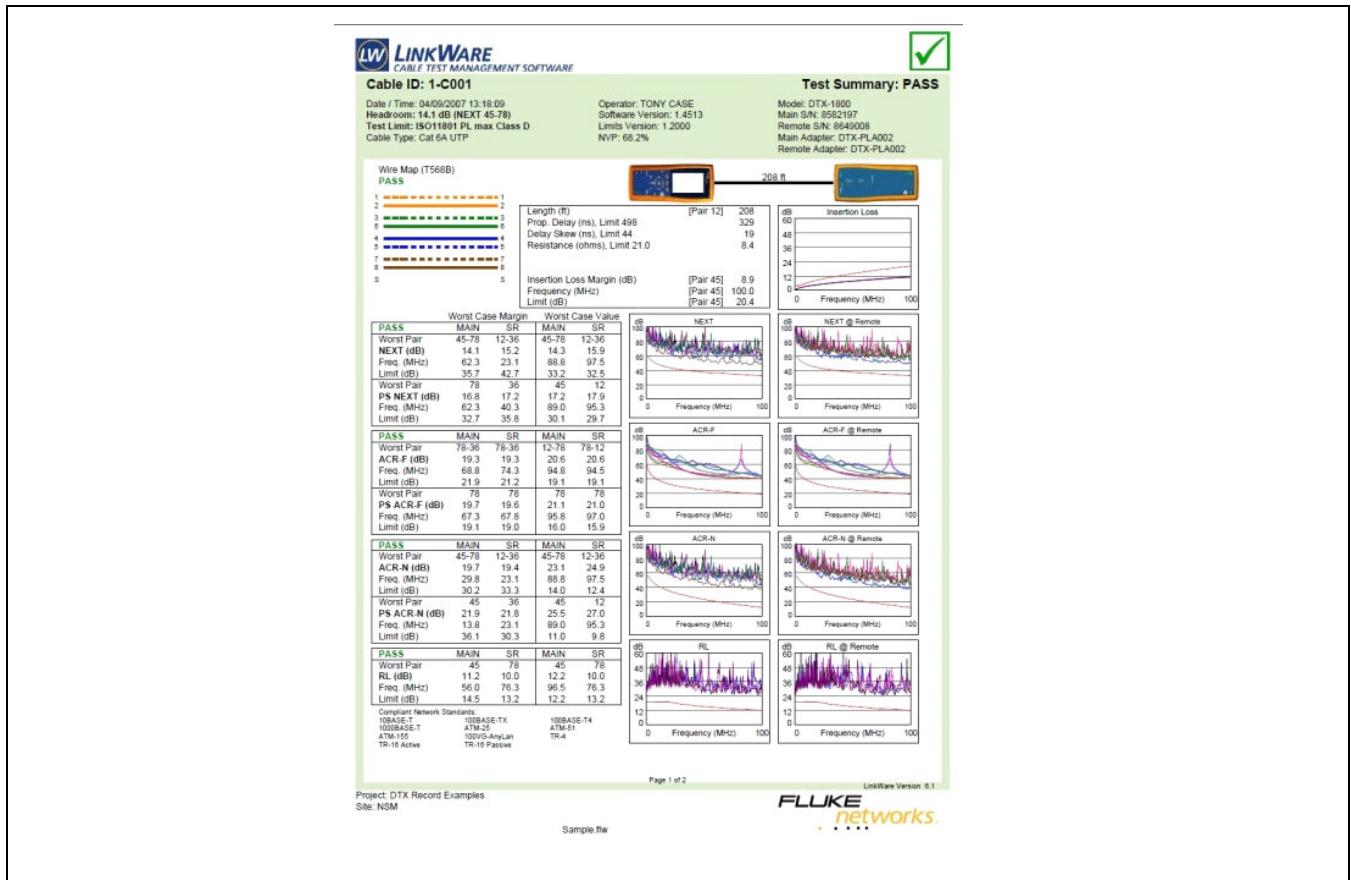


**Figure 7: Certification Test Report**

Shown above is a Cat 6A certification test report which shows how the link performed on every test parameter.

## Certifying Fiber

Basic or Tier 1 certification for fiber is a loss-length test conducted with an Optical Loss Test Set ("OLTS"). The OLTS measures the loss on the fiber link and compares it to a loss budget based on length and bandwidth set by the appropriate standard, such as the Telecommunication Industry Association's TIA-568-C.0, Generic Telecommunications Cabling for Customer Premises. Based on its measurements, the OLTS returns a Pass or Fail grade for that link. This test is used for multimode and single mode fiber.
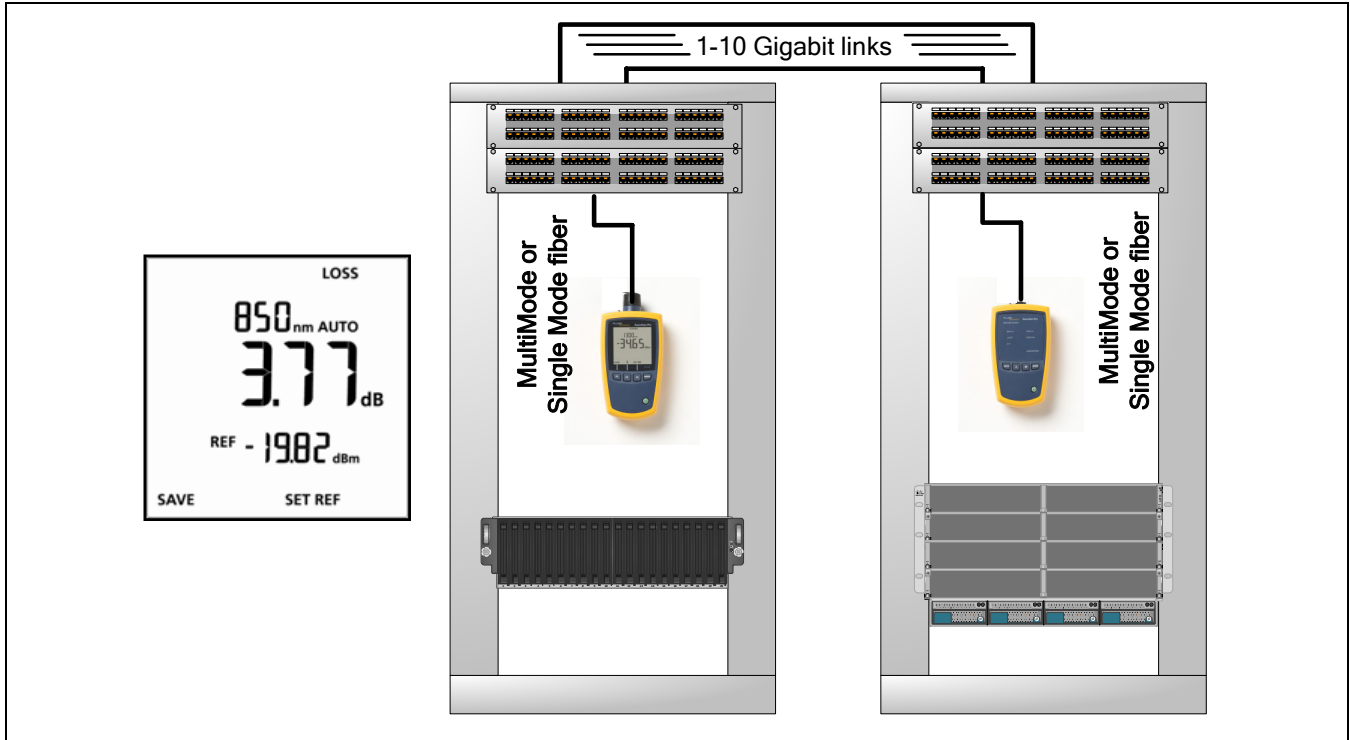
**Figure 8: Fiber Certification**

A detailed results report is essential for fiber certification just as it is for copper. A sample is shown in Figure 9.

**Figure 9: Fiber Certification Test Result Example**

If a fiber link fails the Tier 1 test, first it is necessary to use an Optical Time Domain Reflectometer (OTDR) to identify the problematic area(s). An OTDR is a single-ended test and troubleshooting instrument that performs a detailed assessment of each component in a fiber link. An OTDR launches pulses of light down the fiber and measures reflected light. Based on the relative strength of the return pulses, a trace plots loss as a function of fiber length. This reveals the location of connectors and faults, measures losses, and determines the length of the link.

An OTDR is a very powerful tool for certifying fiber prior to use. It has the advantage of measuring loss in each connector and on each cable segment. Figure 10 shows a trace from an OTDR trace that calculates a grade for the fiber under test.
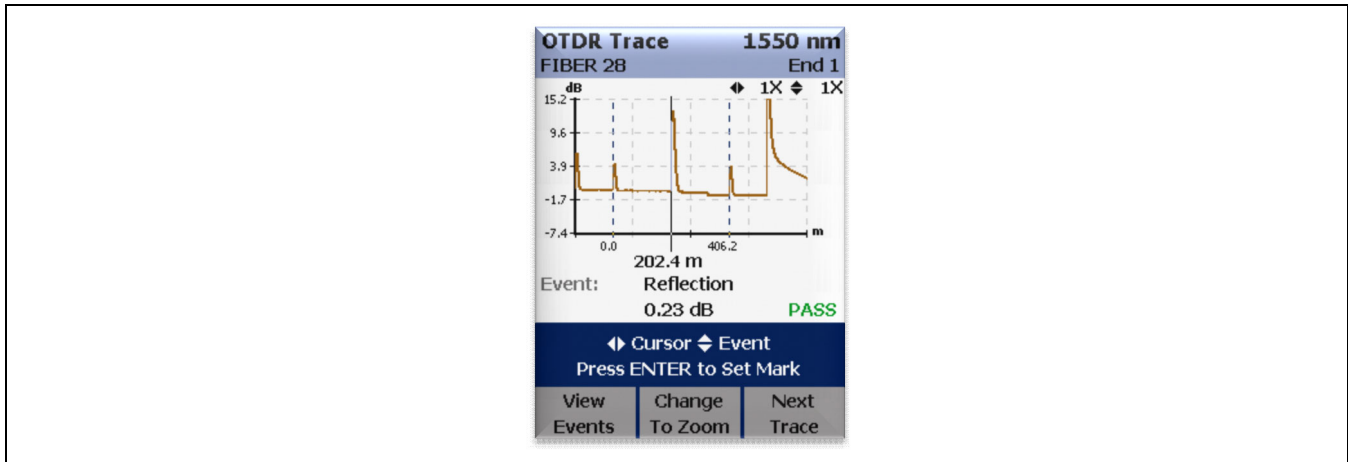
**Figure 10: OTDR Trace Example**

The "Pass" result shows that this fiber link meets the specified standards. As previously mentioned, OTDR traces are invaluable for documenting installation quality and for troubleshooting fiber in the event of a loss-length certification failure.

# Conclusions

As virtualization marches forward, it is triggering fundamental changes to datacenter networks that are unavoidable and, in many cases, desirable. To deliver bandwidth reliability to virtualized assets and end users, 10 Gigabit Ethernet will be employed in the virtualized datacenter. 10 Gigabit Ethernet is important because it is a way to future-proof the datacenter network for years to come while meeting the immediate need to support virtualized servers and services.

10 Gigabit Ethernet can be implemented through a variety of copper and fiber options. The 10GBASE-T standard and new 10GBASE-T silicon open the doors for cost-efficient deployment of 10 Gigabit Ethernet across a virtualized datacenter. Whatever the choice in Layer 1 technology, transitioning to 10 Gigabit Ethernet requires forethought, careful planning, and a methodology for test and troubleshooting.

Fluke Networks is certifying both copper and fiber 10-GbE-based cable plants with Broadcom's portfolio of industry-leading datacenter networking technologies.