*Propagation Delay.*  Some delay in a network arises because a signal requires a small amount of time to travel across a transmission medium. In general, propagation delays are proportional to the distance spanned.  Even with long cable runs, a typical LAN used within a single building has a propagation delay under a millisecond. Although such delays seem irrelevant to a human, a modern computer can execute over one hundred thousand instructions in a millisecond.  Thus, a millisecond delay is significant when a set of computers need to coordinate (e.g., in the financial industry, where the exact time a stock order arrives determines whether an order is accepted).  A network that uses a GEO satellite has much higher delay — even at the speed of light, it takes hundreds of milliseconds for a bit to travel to the satellite and back to earth.

*Access Delay.*  Many networks use shared media.  The set of computers that share a medium must contend for access.  For example, a Wi-Fi wireless network uses a CSMA/CA approach to medium access.  Such delays are known as *access delays.*  Access delays depend on the number of stations that contend for access and the amount of traffic each station sends.  Access delays remain small and fixed unless the medium is overloaded.

*Switching Delay.*  An electronic device in a network (e.g., a Layer 2 switch or router) must compute a next-hop for each packet before transmitting the packet over an output interface.  The computation often involves table lookup, which means memory access.  In some devices, additional time is needed to send the packet over an internal communication mechanism such as a bus or fabric.  The time required to compute a next hop and begin transmission is known as a *switching delay.*  Fast CPUs and special-purpose hardware have made switching delays among the least significant delays in a computer network.

*Queuing Delay.*  The store-and-forward paradigm used in packet switching means that a device such as a router collects the bits of a packet, places them in memory, chooses a next-hop, and then waits until the packet can be sent before beginning transmission.  Such delays are known as *queuing delays.*  In the simplest case, a packet is placed in a FIFO output queue, and the packet only needs to wait until packets that arrived earlier are sent; more complex systems implement a selection algorithm that gives priority to some packets.  Queuing delays are variable — the size of a queue depends entirely on the amount of traffic that has arrived recently.  Queuing delays account for most delays in the Internet.  When queuing delays become large, we say that the network is congested.

*Server Delay.*  Although not part of a network per se, servers are essential to most communication.  The time required for a server to examine a request and compute and send a response constitutes a significant part of overall delay.  Servers queue incoming requests, which means that server delay is variable and depends on the current load.  In many cases, a user's perception of Internet delay arises from server delay rather than network delays.