---

*"On the Internet, no one knows you're a dog,"* went a now-famous New Yorker cartoon in 1996. This remark, made in jest, is a fact (and adds to the humor). The Internet architecture, and its routing system in particular, lacks *origin authentication*: it lacks a reliable and secure way to tell which entity (Internet Service Provider, company, educational institution, etc.) owns any given IP address and who actually sent any given packet. The result is that it's technically possible for any node (router) to originate a route for *a range of IP addresses that does not actually belong to the node's organization* and advertise it to its neighbors via the Border Gateway Protocol (BGP4), the path vector routing protocol used by routers belonging to different organizations to exchange routes on the Internet. If such advertisements are accepted by the neighboring routers and further propagated through the Internet, traffic destined for any address in the advertised range may not reach the destination. This traffic may be diverted elsewhere, a phenomenon known as *route hijacking*, or may simply end up being dropped somewhere, in which case the traffic is said to be *black-holed*.

Route hijacks and black holes are a fairly regular occurrence in the Internet today. From time to time, however, the problem occurs at a magnitude large enough to affect a large number of users. One of these events happened on February 24, 2008, when YouTube, the popular video sharing web site, became unreachable to most people on the Internet.

To understand what happened, a few facts are useful:

1. To forward a packet along the right link, an Internet router looks the destination up in a copy of the routing table. The routing table contains ranges of IP addresses, called *IP prefixes*, and the router chooses the "best" entry that contains the destination address. Sometimes, there may be more than one matching entry, and the "best" is defined as the *longest prefix match* (LPM). LPM chooses the routing table entry (the link) that shares the largest number of in-sequence common bits with the destination (the "longest match"), starting from the most-significant bit and going toward the least-significant bit. Here, we think of IP addresses not as "dotted decimals", but as 32-bit strings.

2. `www.youtube.com` resolves to multiple IP addresses, all in the same "/24" prefix range; i.e., the first 24 bits of their IP addresses are all the same.

3. The Pakistani government had ordered all its ISPs to block access to YouTube. In general, there are two ways to block traffic from an IP address; the first is to simply drop all packets to and from that address, while the second is to *divert* all traffic going *to* the address to a different location, where one might present the user with a web page that says something like "We're sorry, but your friendly government has decided that YouTube isn't good for your mental well-being." The latter presumably is a better customer experience because it tells users what's going on, so they aren't in the dark, don't make unnecessary calls to customer support, and don't send a whole lot of traffic by repeatedly making requests to the web site. (Such diversion is quite common; it's used in many public WiFi spots that require a sign-on, for example.)

Pakistan Telecom introduced a /24 routing table entry for the range of addresses to which `www.youtube.com` resolves. So far, so good. Unfortunately, probably because of a misconfiguration (presumably caused by human error by a tired or careless engineer) rather than malice, routers from

*Hari Balakrishnan*

Pakistan Telecom leaked this /24 routing advertisement to one of its ISPs (PCCW in Hong Kong). Normally, this leak should not have caused a problem *if* PCCW knew (as it should have) the valid set of IP prefixes that Pakistan Telecom owned. Unfortunately, perhaps because of another error or oversight, PCCW didn't ignore this route. At this stage, inside PCCW and its customer's networks would've been "black-holed" from YouTube, and all traffic destined there would've been sent to Pakistan Telecom.

The problem was much worse because essentially the entire Internet was unable to get to YouTube. That's because of the LPM method used to find the best route to a destination. Under normal circumtances, there is no /24 advertised on behalf of YouTube by its ISPs; those routes are contained in advertisements that cover a (much) wider range of IP addresses. So, when PCCW's neighbors saw PCCW advertising a more-specific route, they followed the rules and imported those routes into their routing tables, re-advertising them to their respective neighbors, until the entire Internet (except, presumably, a few places such as YouTube's internal network itself) had this poisoned routing table entry for the IP addresses in question.

This discussion highlights a key theme about large networks and systems: when they fail, the reasons are complicated. In this case, the following events all occurred:

1. The Pakistani government decided to censor a particular site because it was afraid that access would create unrest.
2. Pakistan Telecom decided to divert traffic using a /24 and leaked it in error.
3. PCCW, which ought to have ignored the leaked advertisement, didn't.
4. The original correct advertisements involved less-specific prefixes; in this case, had they been /24s as well, the problem may not have been as widespread.
5. Pakistan Telecom inadvertently created a massive traffic attack on itself (and on PCCW) because YouTube is a very popular site getting lots of requests. Presumably the amount of traffic made diagnosis difficult because packets from tools like `traceroute` might not have progressed beyond points of congestion, which might have been upstream of Pakistan Telecom. It appears that the first indication of what might have really happened came from an investigation of the logs of routing announcements that are available to various ISPs and also publicly. They showed that a new network (ISP) had started originating a route to a prefix that had previously always been originated by some other entity.

   This observation suggests that a combination of public "warning systems" that maintain such information and flag anomalies might be useful (though there are many legitimate reasons why the originator of a route may change, too). It also calls for the maintenance of a correct registry containing the owning organizations for various IP address ranges. Unfortunately, some studies have shown that current registries unfortunately have a number of errors and omissions, so the viability of this approach isn't obvious.

Far from being an isolated incident, such problems (black holes and hijacks) arise from time to time.[1] There are usually one or two serious incidents of this kind every year, though selectively taking down a popular site tends to make the headlines more easily. There are also several smaller-scale incidents and anomalies that show up on a weekly or even daily basis. Some of these route hijacks are even used to send hard-to-trace email spam. How that works is another story for a later time...

---

[1] The so-called AS 7007 incident in 1997 was perhaps the first massive outage that partitioned most of Sprint's large network and customer base from the rest of the Internet. In that incident, a small ISP in Florida was the party that originated the misconfigured route leak, but other ISPs were also culpable in heeding those unrealistic route advertisements.

*Hari Balakrishnan*