YourNetworkYourProblem

Your Network Your Problem

Apache Hadoop is a distributed application that runs across a cluster of machines.

For it to work, all these machines must be able to find each other, to talk to each other, and indeed, simply identify themselves so that other machines in the cluster can find them.

Externally accessible Hadoop clusters need to be visible across the rest of the network which needs access to it.

And of course, all these machines need to be wired together using network switches and routers.

For that reason, network setup is a critical part of a Hadoop cluster. If you do not do this, Hadoop will not work and you will be left staring at stack traces in Hadoop logs trying to diagnose what is wrong. You may even file bug reports saying "Help! Hadoop doesn't work!"

It does work for everybody else -and the reason it does not work for you is because the network is misconfigured it doesn't.

And, because it is your network, nobody else is going to fix it for you --except in the special case that you are using a paid packaging of Hadoop, where you should contact your vendor and ask them for help. The Hadoop developers cannot and will not help you: filing bug reports will simply result in the issue being closed as invalid along with a link to the InvalidJiralssues page.

Here are some of the common problems in network and host configurations

1. DNS and reverse DNS broken/non-existent. 2. Host tables in the machines invalid. 3. Firewalls in the hosts blocking connections. 4. Routers blocking traffic. 5. Hosts with multiple network cards listening/talking on the wrong NIC. 5. Difference between the hadoop configuration files' definition of the cluster (especially hostnames and ports) from that of the actual cluster setup.

Don't forget hardware either: firmware problems in the NICs; broken NICs or buggy device drivers, damaged cables, cables not correctly inserted. If you are using bonded Ethernet switches, the risk of cable miswiring increases significantly. Routers and switches are also a source of interesting issues, ranging from dropped packets to false ARP responses. Then there's DHCP and the OS-side code to work with it.

The TroubleShooting page lists some recurrent error messages, possible root causes and ways to track down the problem. There's also the traditional network diagnostics tools, from Dig and traceroute to Wireshark.

Start at the bottom of the stack assume the hardware layer is broken unless you can demonstrate otherwise then work your way up the stack, of which Hadoop is only the Level 7 part of the problem. If you start trying to debug it from there, you're coming from the wrong direction.

If, after validating your network stack layer by layer, you still can't get Hadoop to work, you could try asking for help on the hadoop user list.

The key point to remember is this: it is your network that is playing up -and only you are in a position to fix it.