# FAUCET AT
# SANDIA NATIONAL LABORATORIES
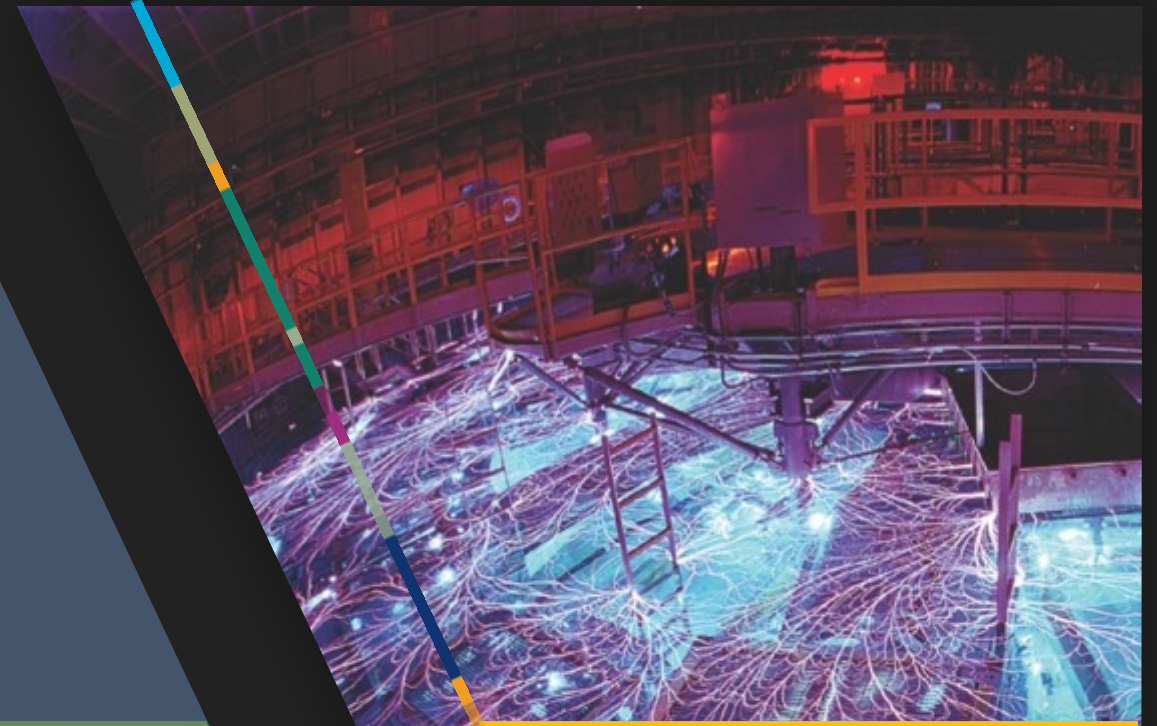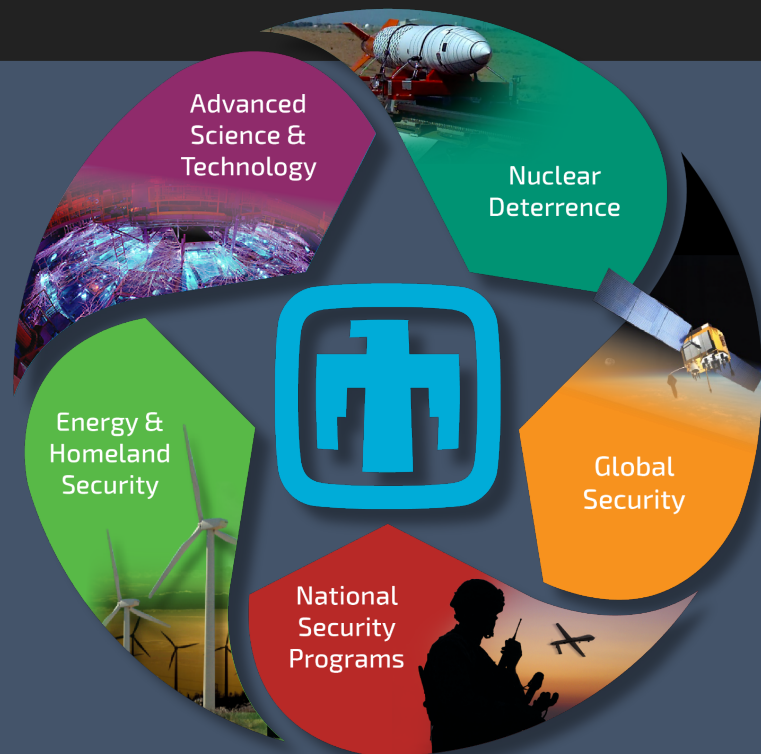
Timothy Toole

Advanced Science & Technology

Nuclear Deterrence

Global Security

National Security Programs

Energy & Homeland Security

## SANDIA IS A FEDERALLY FUNDED RESEARCH AND DEVELOPMENT CENTER

Main sites
Albuquerque, New Mexico
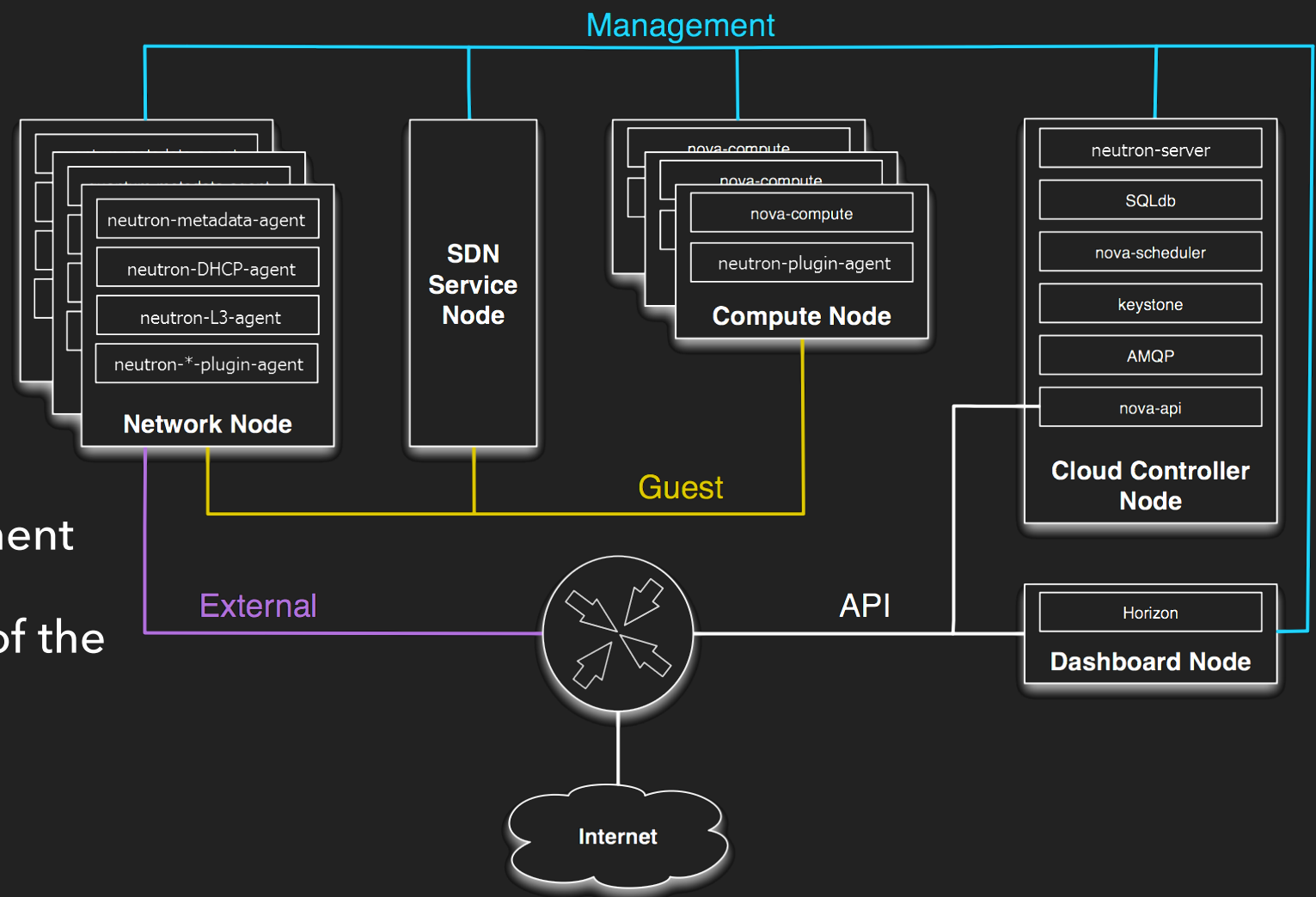Livermore, California

ENERGY   NNSA

Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia LLC, a wholly owned subsidiary of Honeywell International Inc. for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.
SAND 2018-4900 PE

Sandia National Laboratories

▸ Multiple groups at Sandia conducting research & development in computing, information science and cybersecurity

   ▸ Ensures security of critical military, government, and commercial networks using trusted systems to detect anomalies and intrusions

▸ Representing a team that focuses on system architectures, computer networks and analysis

   ▸ Not necessarily working on fundamental science, but involved in developing and integrating new or novel techniques

   ▸ Share and advise other groups and teams that benefit from the techniques

▸ Experience with networking, virtualization, computer systems & science

Sandia National Laboratories

▸ Cloud (orchestration) environments are complex

   ▸ Require multiple isolated physical & logical networks

   ▸ Applies to Openstack, VMware, Kubernetes, etc.

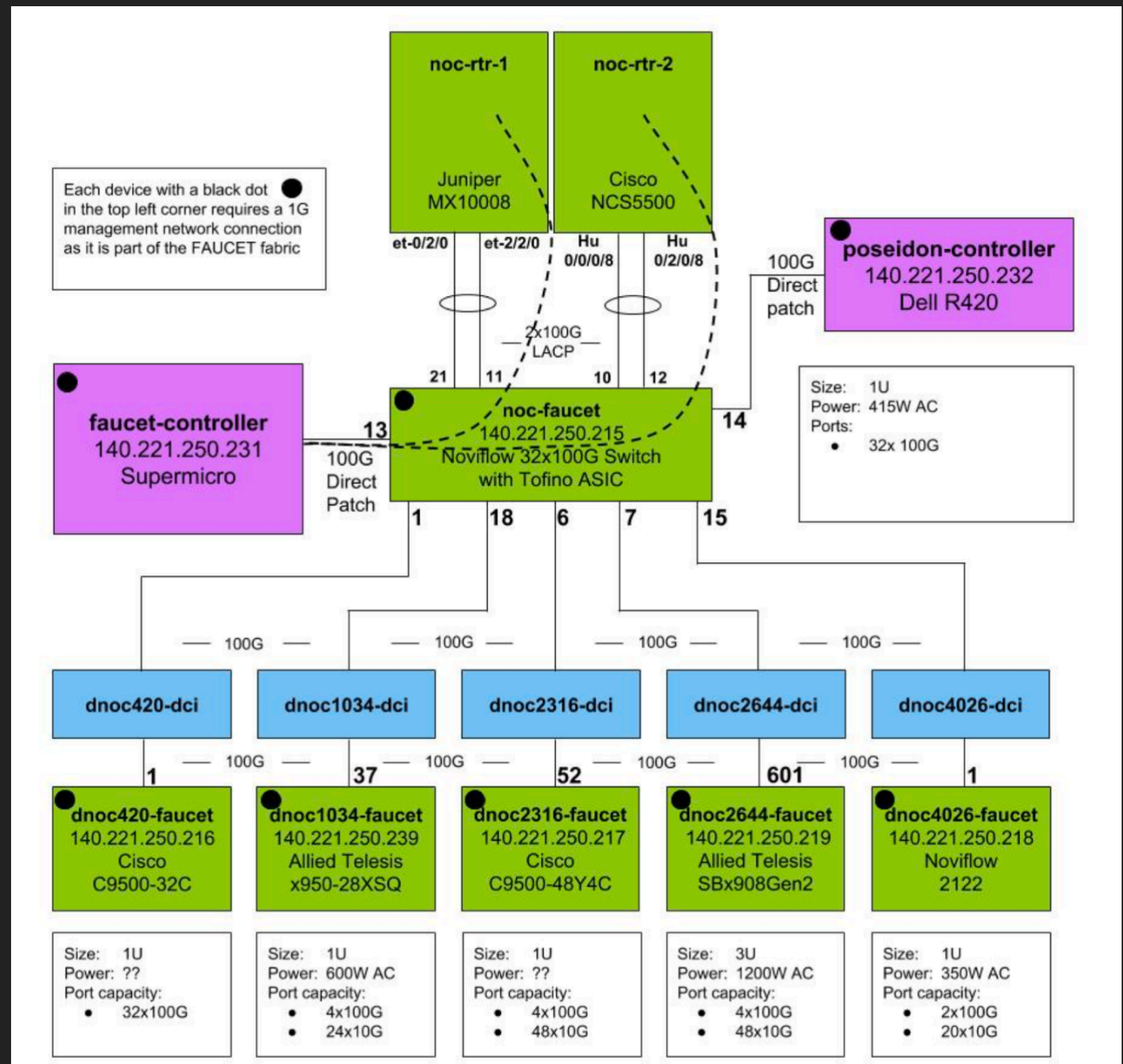▸ Desire to build and deploy multiple, isolated cloud environments

   ▸ Full cloud environment per tenant

   ▸ Perform application testing

▸ Gain real-world experience with SDN

   ▸ Prior controller experience was in a limited scope

   ▸ Extend beyond academic environment

   ▸ Become an operator & maintainer of the network

Management

Network Node
- neutron-metadata-agent
- neutron-DHCP-agent
- neutron-L3-agent
- neutron-*-plugin-agent

SDN Service Node

Compute Node
- nova-compute
- nova-compute
- nova-compute
- neutron-plugin-agent

Cloud Controller Node
- neutron-server
- SQLdb
- nova-scheduler
- keystone
- AMQP
- nova-api

Guest

External

API

Dashboard Node
- Horizon

Internet

Sandia National Laboratories

▸ **PRODUCTION** - Stability!

▸ Multi-vendor

▸ Statistics (Gauge)

▸ Analytics (Poseidon)

▸ https://github.com/wandsdn/sc18-faucet-configs

## MAJOR Kudos to

## the SC18 team!

▸ Use this design as a foundation for our use

▸ Solicited price quotes from Allied Telesis, Cisco, NoviFlow, HPE-Aruba

  ▸ Used model numbers from SC18 diagram and Faucet documentation

  ▸ Several phone conversations with vendor reps

▸ Cleaned out all of our [legacy] Arista switches (sad day)

▸ **Switching hardware list:**

  ▸ EdgeCore Wedge 100BF-32X (Tofino) with NoviWare NOS



  ▸ Allied Telesis x950



High bandwidth

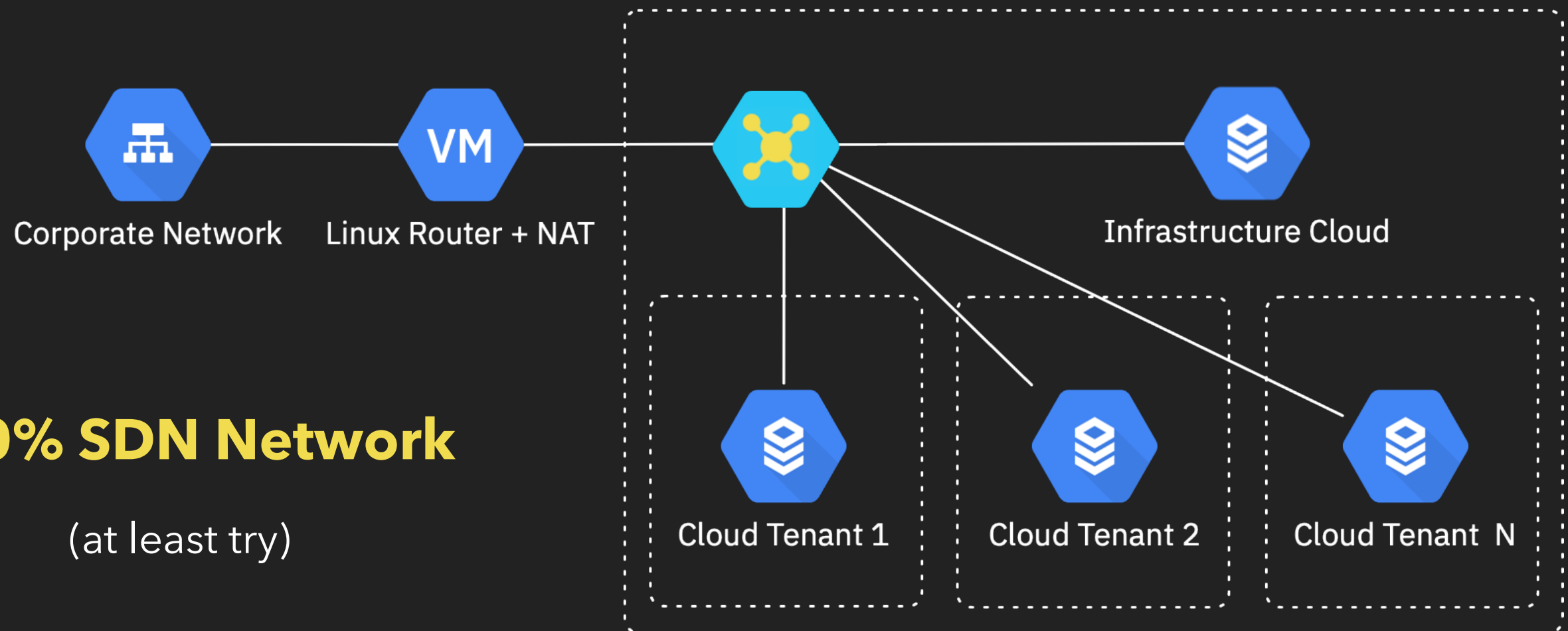  ▸ Many HPE-Aruba 2930F switches for 1G Copper connections

High bandwidth

▸ Desire to build out  multiple, isolated cloud environments without physically re-wiring or managing a switch CLI

    ▸ Automate testing & deployment process as much as possible - Python + bash

    ▸ Full cloud environment per tenant

**100% SDN Network**

(at least try)



Corporate Network — Linux Router + NAT — Infrastructure Cloud — Cloud Tenant 1 — Cloud Tenant 2 — Cloud Tenant N

**Sandia National Laboratories**
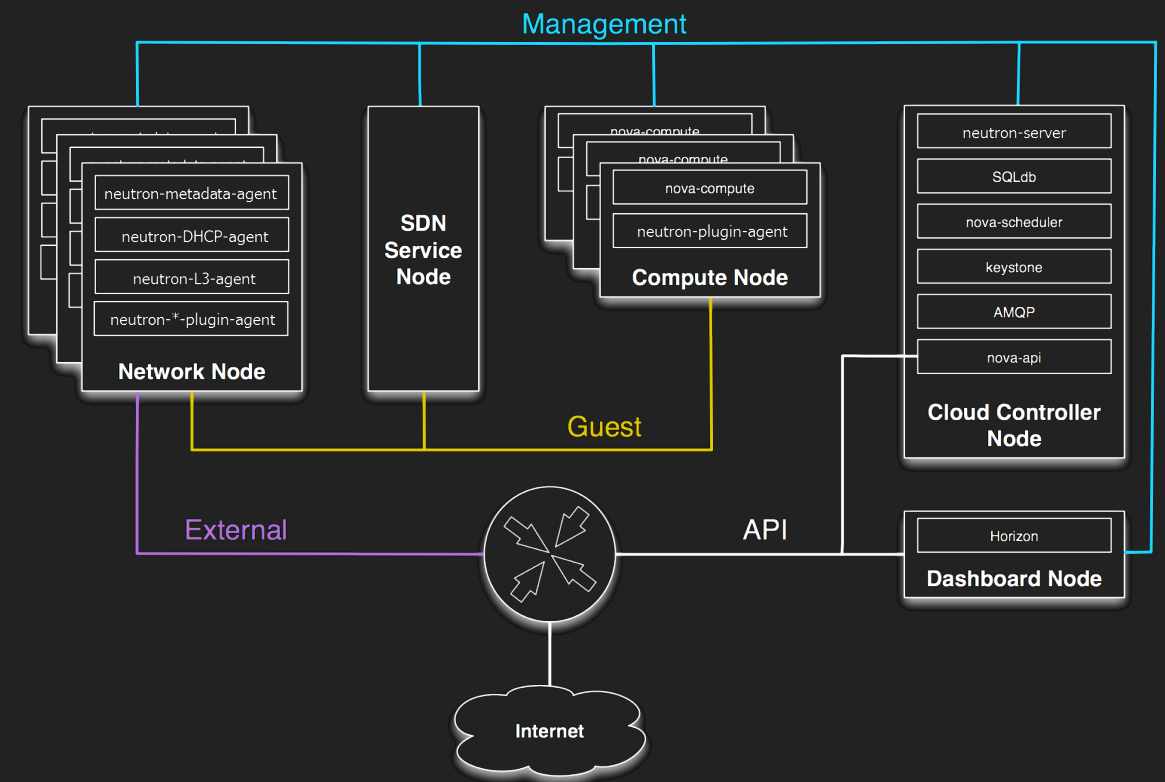
+



▸ 5 Physical interfaces to connect per compute node

  ▸ IPMI (low bandwidth)

  ▸ Two 10Gb/s interfaces (high bandwidth)

  ▸ Two 1Gb/s interfaces (low bandwidth)

▸ **Keep good notes:** physical host interface <-> switchport mappings

**Sandia National Laboratories**

```
interfaces:
    1:
        name: "Port 1"
        description: "cluster0 blade2 ipmi"
        native_vlan: ipmi-tenant-a-100
    2:
        name: "Port 2"
        description: "cluster0 blade2 eth1"
        native_vlan: mgmt-tenant-a-101
    3:
        name: "Port 3"
        description: "cluster0 blade2 eth10"
        native_vlan: guest-tenant-a-102
    4:
        name: "Port 4"
        description: "cluster0 blade2 eth11"
        native_vlan: ext-tenant-1-103
```



**Routing is handled externally to Faucet:**

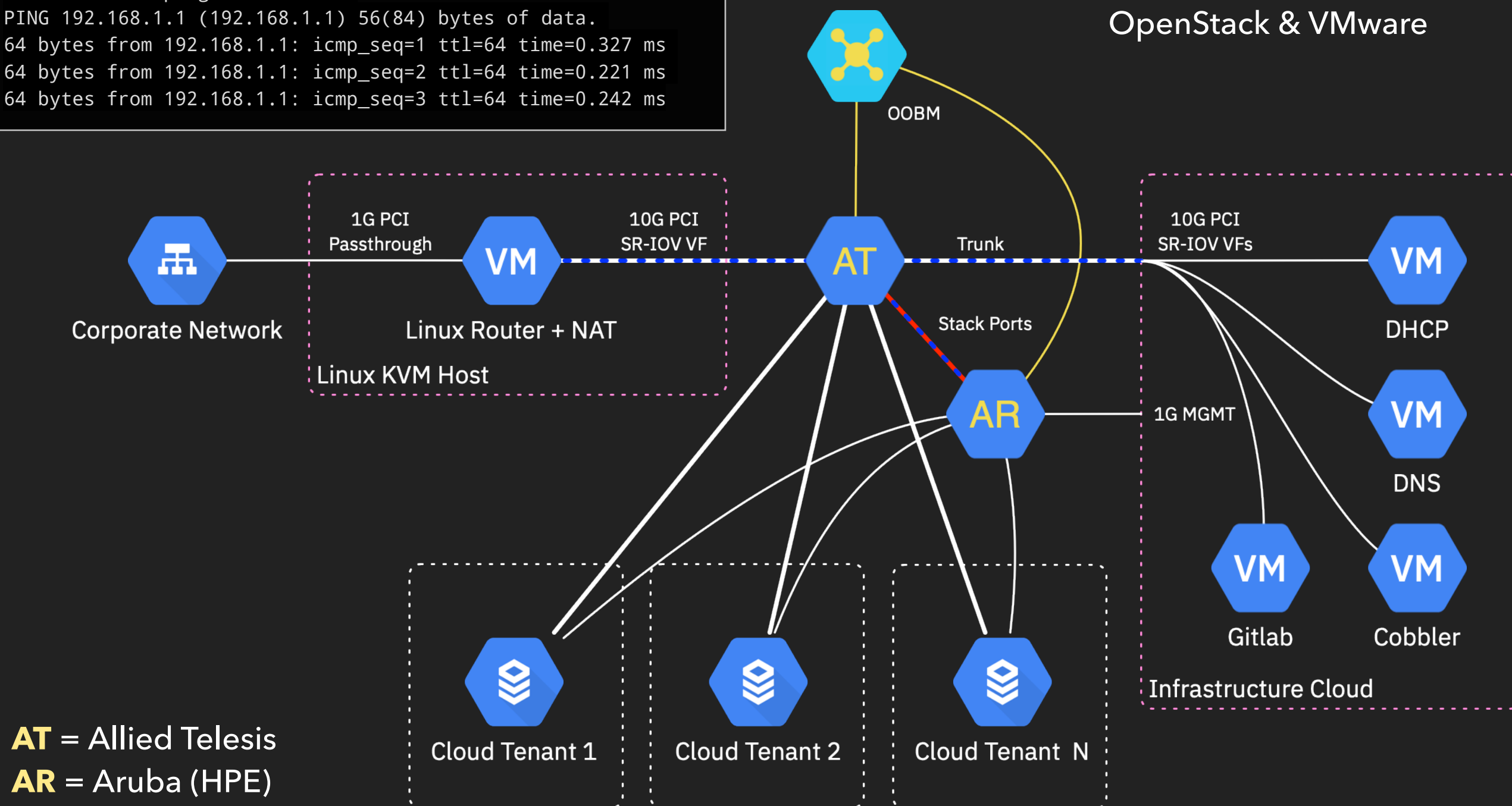‣ Static route(s) pushed to a Linux-based "routing VM (+ NAT)"

**Incredibly easy to automate:**

‣ Ansible (FAUCET config generation), Python (pushing L3 routes to Linux VM), Cobbler (OS deployment)

**Sandia National Laboratories**

```
user@node1:~$ ping 192.168.1.1
PING 192.168.1.1 (192.168.1.1) 56(84) bytes of data.
64 bytes from 192.168.1.1: icmp_seq=1 ttl=64 time=0.327 ms
64 bytes from 192.168.1.1: icmp_seq=2 ttl=64 time=0.221 ms
64 bytes from 192.168.1.1: icmp_seq=3 ttl=64 time=0.242 ms
```

Used Cobbler to deploy
OpenStack & VMware

OOBM

1G PCI
Passthrough

10G PCI
SR-IOV VF

Trunk

10G PCI
SR-IOV VFs

Corporate Network

**VM**

Linux Router + NAT

Linux KVM Host

**AT**

Stack Ports

**AR**

1G MGMT

**VM**

DHCP

**VM**

DNS

**VM**

Gitlab

**VM**

Cobbler

Infrastructure Cloud

Cloud Tenant 1

Cloud Tenant 2

Cloud Tenant N

**AT** = Allied Telesis
**AR** = Aruba (HPE)

Sandia National Laboratories

▸ This is where the "fun" begins

▸ This is where lessons are learned

192.168.0.0/16

VIPs
192.88.99.1/24
192.168.vlanid.1/24

OOBM

Corporate Network

Linux Router + NAT

192.88.99.254

VM

AT

Trunk

Stack Ports

Infrastructure Cloud

AR    AR    AR

▸ Topology is generally architected by understanding of where traffic flows

Cloud Tenant N

**AT** = Allied Telesis
**AR** = Aruba (HPE)

* 192.88.99.0/24 - IANA Reserved: Deprecated (6to4 Relay Anycast)

Sandia National Laboratories

▸ Use the Faucet tutorial method, but for testing configs: **Open vSwitch** and **Network Namespaces**

▸ **./make_network.sh** - *based solely on how the human wired up the switches*

```
# Create DP atx950
sudo ovs-vsctl add-br atx950 \
        -- set bridge atx950 other-config:datapath-id=0xe01aea517a38 \
        -- set bridge atx950 other-config:disable-in-band=true \
        -- set bridge atx950 fail_mode=secure \
        -- set-controller atx950 tcp:127.0.0.1:6653 tcp:127.0.0.1:6654

# Create Trunk Ports between atx950 and ar2930f-1
sudo ovs-vsctl add-port atx950 atx950-p1 \
        -- set interface atx950-p1 type=patch \
        -- set interface atx950-p1 options:peer=ar2930f-1-p49 \
        -- set interface atx950-p1 ofport_request=1

sudo ovs-vsctl add-port ar2930f-1 ar2930f-1-p49 \
        -- set interface ar2930f-1-p49 type=patch \
        -- set interface ar2930f-1-p49 options:peer=atx950-p1 \
        -- set interface ar2930f-1-p49 ofport_request=49
```
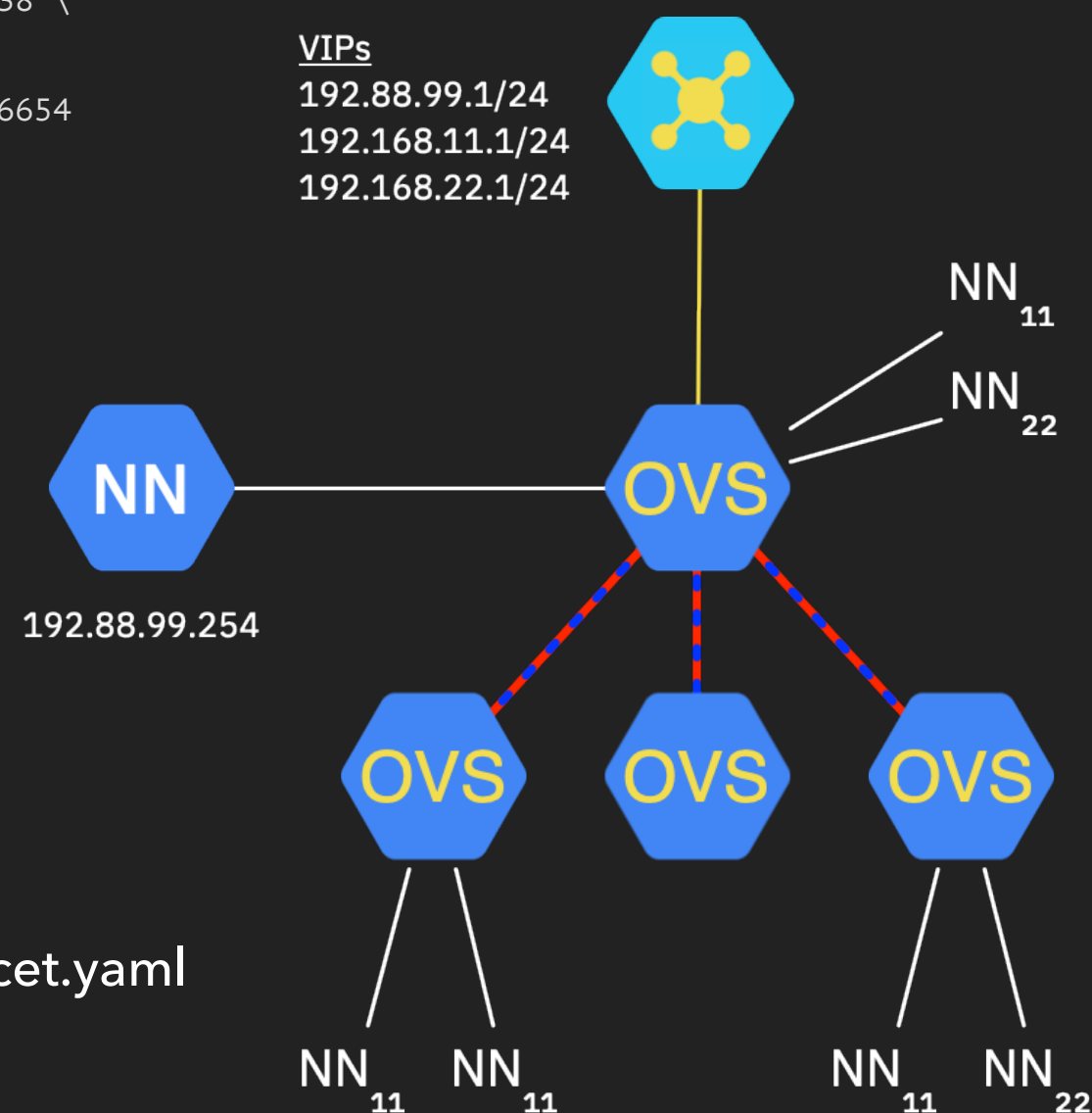
▸ **./make_hosts.sh**

```
# Create Test Namespaces and Connect them - atx950
create_ns host1 192.168.11.2/24
as_ns host1 ip route add default via 192.168.11.1
sudo ovs-vsctl add-port atx950 veth-host1 \
                -- set interface veth-host1 ofport_request=5
```

▸ **./run_tests.sh**

```
# as_ns host3 ping -q -c3 192.168.10.2 > /dev/null
```

▸ Modify the "**Hardware: Open vSwitch**" line in faucet.yaml

VIPs
192.88.99.1/24
192.168.11.1/24
192.168.22.1/24

NN₁₁
NN₂₂

NN
192.88.99.254

OVS

OVS   OVS   OVS

NN₁₁  NN₁₁        NN₁₁  NN₂₂

**Sandia National Laboratories**

```
faucet.valve ERROR    DPID 282562769570368 (0x100fd4581ca40) ar2930f-2 OFError type: OFPET_TABLE_FEATURES_FAILED
code: OFPTFFC_EPERM
version=0x4,msg_type=0x1,msg_len=0x4c,xid=0x3f6fcf42,OFPErrorMsg(code=5,data=bytearray(b'\x04\x12\x04\x18?
o\xcfB\x00\x0c\x00\x00\x00\x00\x00\x00\x00x\x00\x00\x00\x00\x00\x00port_acl\x00\x00\x00\x00\x00\x00\x00\x00\x00\x
00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00\x00'),type=13)Frame 653:
142 bytes on wire (1136 bits), 142 bytes captured (1136 bits)
```

```
Ethernet II, Src: ec:eb:b8:33:05:c0, Dst: 0c:c4:7a:54:a5:65
Internet Protocol Version 4, Src: 192.168.1.13, Dst: 192.168.1.60
Transmission Control Protocol, Src Port: 61578, Dst Port: 6653, Seq: 987130365, Ack: 1243404425, Len: 76
OpenFlow 1.3
    Version: 1.3 (0x04)
    Type: OFPT_ERROR (1)
    Length: 76
    Transaction ID: 3080866189
    Type: OFPET_TABLE_FEATURES_FAILED (13)
    Code: OFPTFFC_EPERM (5)
    Body: 04120428b7a2498d000c00000000000088000000000000...
        Version: 1.3 (0x04)
        Type: OFPT_MULTIPART_REQUEST (18)
        Length: 1064
        Transaction ID: 3080866189
        Type: OFPMP_TABLE_FEATURES (12)
        Flags: 0x0000
        Pad: 00000000
        Table features
[Malformed Packet: openflow_v4]
    [Expert Info (Error/Malformed): Malformed Packet (Exception occurred)]
        [Malformed Packet (Exception occurred)]
        [Severity level: Error]
        [Group: Malformed]
```

Sandia National Laboratories

```
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:{ "error_code":"OFPTFFC_EPERM","error_reason":"There is no space available in the H/W to
accomodate the new pipeline","process_time":"0.166 ms","pipeline":[{"table_id":0,"name":"port_acl","
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:config":"0x3","max_entries":32,"metadata_match":"0x0","metadata_write":"0x0","match":
["in_port"],"wildcards":["in_port"],"next_tables":["1","5","6","7"],"instructions":["goto_table","apply_
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:actions"],"apply_actions":["output","pop_vlan","group"]},{"table_id":
1,"name":"vlan","config":"0x3","max_entries":288,"metadata_match":"0x0","metadata_write":"0x0","match":["in_port","vlan_
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:vid","eth_type","eth_dst/has_mask"],"wildcards":
["in_port","vlan_vid","eth_type","eth_dst"],"next_tables":["2"],"instructions":["goto_table","apply_actions"],"apply_setfield":["vlan_vid"],"
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:apply_actions":["output","set_field","pop_vlan","push_vlan","group"]},{"table_id":
2,"name":"eth_src","config":"0x3","max_entries":800,"metadata_match":"0x0","metadata_write":"0x0","match":[
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:"vlan_vid","in_port","eth_src","eth_type","eth_dst/has_mask"],"wildcards":
["vlan_vid","in_port","eth_src","eth_type","eth_dst"],"next_tables":["3","4","5","6","7"],"instructions":["goto_tab
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:le","apply_actions"],"apply_setfield":["vlan_vid","eth_dst"],"apply_actions":
["output","set_field","pop_vlan","push_vlan","group"],"next_tables_miss":["6"],"instructions_miss":["goto_table"
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:]},{"table_id":3,"name":"ipv4_fib","config":"0x3","max_entries":
608,"metadata_match":"0x0","metadata_write":"0x0","match":["vlan_vid","eth_type","ipv4_dst/has_mask"],"wildcards":["vlan_vid"
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:,"eth_type","ipv4_dst"],"next_tables":["5","6","7"],"instructions":
["goto_table","apply_actions"],"apply_setfield":["eth_dst","eth_src","vlan_vid"],"apply_actions":["output","set_field","po
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:p_vlan","push_vlan","group"]},{"table_id":4,"name":"ipv6_fib","config":"0x3","max_entries":
608,"metadata_match":"0x0","metadata_write":"0x0","match":["ipv6_dst/has_mask","vlan_vid","eth_typ
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:e"],"wildcards":["ipv6_dst","vlan_vid","eth_type"],"next_tables":
["5","6","7"],"instructions":["goto_table","apply_actions"],"apply_setfield":["eth_dst","eth_src","vlan_vid"],"apply_actions
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:":["output","set_field","pop_vlan","push_vlan","group"]},{"table_id":
5,"name":"vip","config":"0x3","max_entries":64,"metadata_match":"0x0","metadata_write":"0x0","match":["ip_proto","icmpv6
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:_type","arp_tpa","eth_dst","eth_type"],"wildcards":
["ip_proto","icmpv6_type","arp_tpa","eth_dst","eth_type"],"next_tables":["6","7"],"instructions":["goto_table","apply_actions"],"apply_act
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:ions":["output","pop_vlan","group"]},{"table_id":
6,"name":"eth_dst","config":"0x3","max_entries":800,"metadata_match":"0x0","metadata_write":"0x0","match":["vlan_vid","eth_dst"],"instructio
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:ns":["apply_actions"],"apply_actions":["output","pop_vlan","group"],"next_tables_miss":
["7"],"instructions_miss":["goto_table"]},{"table_id":7,"name":"flood","config":"0x3","max_entries":96
Sep 23 13:14:34 192.168.1.12 OPFL: OPFL eOFNetTask:,"metadata_match":"0x0","metadata_write":"0x0","match":["in_port","vlan_vid","eth_dst/
has_mask"],"wildcards":["in_port","vlan_vid","eth_dst"],"instructions":["apply_actions"],"apply_actions
```

**Lessons:**

▸ Read the documentation. Understand it.

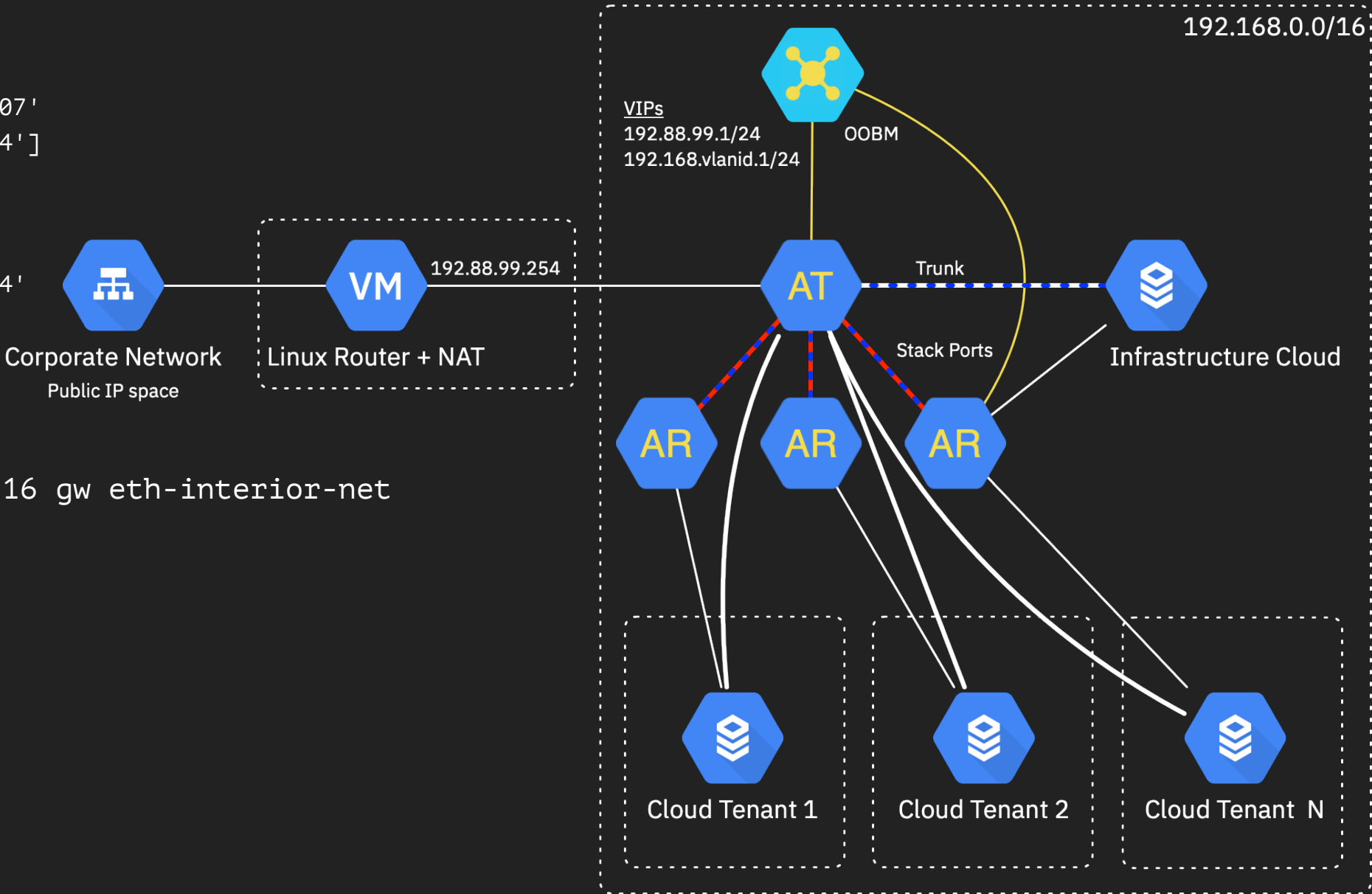▸ Setup a syslog server to capture OF messages and errors from switches

Sandia National Laboratories

```
ar2930f-2# show openflow instance aggregate flow-table

OpenFlow Instance Flow Table Information

Table                          Flow      Miss
ID     Table Name              Count     Count         Goto Table
-----  ------------------      --------  ------------- -------------
0      port_acl                1         0             1, 5, 6, 7
1      vlan                    1         0             2
2      eth_src                 16        0             3, 4, 5, 6, 7
3      ipv4_fib                3         0             5, 6, 7
4      ipv6_fib                6         0             5, 6, 7
5      vip                     17        0             6, 7
6      eth_dst                 1         0             *
7      flood                   14        0             *
```

▸ TFM = **Table Features Message**

**OpenFlow Specification 1.3**
7.3.5.5 Table Features

The OFPMP_TABLE_FEATURES multipart type allows
a controller to both query for the capabilities of
existing tables, and to optionally ask the switch to
reconfigure its tables to match a supplied
configuration.



Sandia National Laboratories

## Faucet vlans.yaml:

```
vlans:
    routing:
        vid: 7
        description: "Gateway Net"
        faucet_mac: 'de:ad:be:ef:00:07'
        faucet_vips: ['192.88.99.1/24']
        routes:
            - route:
                ip_dst: "0.0.0.0/0"
                ip_gw: '192.88.99.254'
```
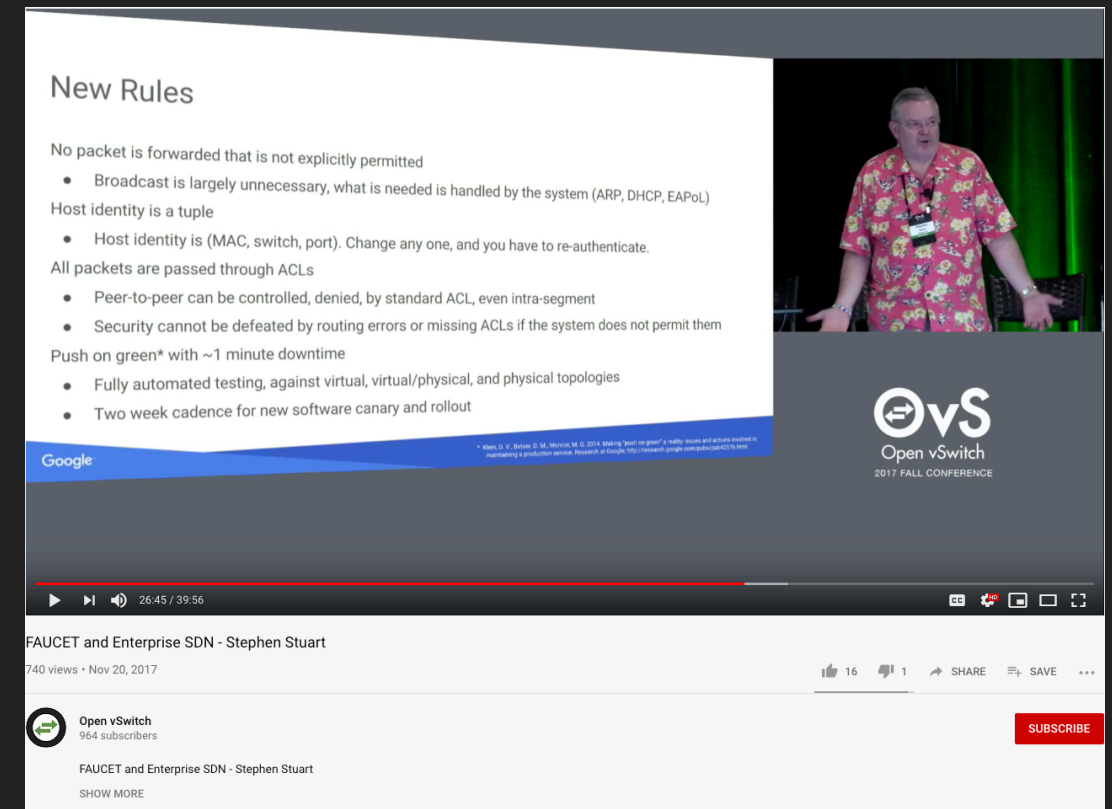
## Linux Router:

```
# route add -net 192.168.0.0/16 gw eth-interior-net
```



192.168.0.0/16

VIPs
192.88.99.1/24
192.168.vlanid.1/24

OOBM

VM   192.88.99.254

AT   Trunk

Stack Ports

Corporate Network
Public IP space

Linux Router + NAT

Infrastructure Cloud

AR   AR   AR

Cloud Tenant 1   Cloud Tenant 2   Cloud Tenant  N

Sandia National Laboratories

▸ Getting up and running with Faucet is super easy - follow the tutorial

　　▸ Anther required reading: http://docs.openvswitch.org/en/latest/tutorials/faucet/

▸ Open vSwitch is awesome!

　　▸ 90%+ of design + testing can be done with software

▸ Allied Telesis hardware has been rock solid



▸ **Start small, grow slowly**

　　▸ *Follow the packet!*

　　▸ Try out every feature available

　　▸ Keep an open mind. Many traditional networking concepts apply, many don't
　　　https://www.youtube.com/watch?v=BDje6HGBwso  Go to 23:10 mark

▸ Don't run your controller off a switch it's controlling (so much for a 100% OF controlled network)
　　¯\_(ツ)_/¯

![Sandia National Laboratories] **Sandia National Laboratories**

▸ Highlights from the documentation

    ▸ Faucet Design and Architecture - **Faucet Openflow Switch Pipeline**

    ▸ Vendor-Specific documentation (Allied Telesis, Aruba, etc.)

▸ OpenFlow 1.3.x Specification

▸ Diagrams - need more

    ▸ Began our own stash of configs and associated diagrams



Sandia National Laboratories

▸ **Collect logs** - as much as you can

    ▸ Save faucet.log files on your controller(s)

    ▸ Collect PCAP on your OpenFlow channel (assumes non-TLS)

    ▸ debug openflow on switches to a syslog server

▸ **Follow a software development mindset**

    ▸ Run tests against configs with OVS, **then with hardware**

    ▸ Script/automate as much as possible

▸ **Not all switching hardware is created equally**

    ▸ Brad quote from ONS2019: "I learned a lot more about vendor hardware architectures than I expected to."
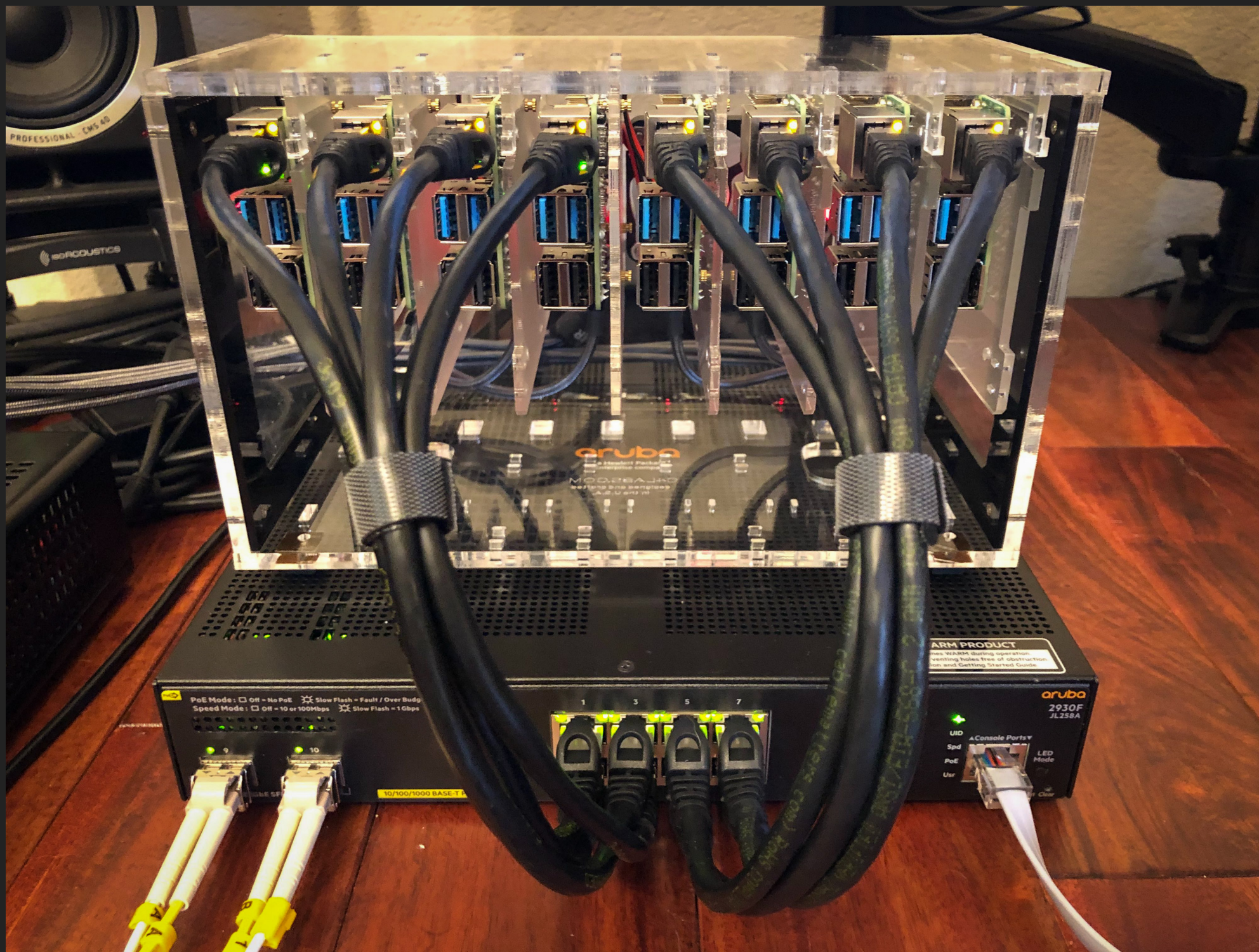
**Sandia National Laboratories**

▸ ~300 physical interfaces, 4 switches controlled via Faucet in a stacked, L3 routing architecture

▸ Network is quiet, efficient, reconfigurable. **Awesome!**

▸ Controller and switch management interfaces are [almost] completely out-of-band (~7 switch ports)

  ▸ Had a few oddities when bridging legacy [dumb] switch into OF-controlled network

  ▸ OOB controller network is Faucet routed, but heavily ACL'd - desire for Gauge/Prometheus/ Grafana

▸ HW limitation errors fixed Aruba switches (thanks, Josh!) - See '**port_table_scale_factor**' feature

  ▸ Built a second, smaller testbed for running new Faucet configs on Aruba hardware

▸ **Re-designing the architecture (again)**

  ▸ Integrating EdgeCore/NoviFlow switches into the network

  ▸ Adding additional controllers: redundancy & segmentation (based upon availability needs)

▸ Writing our own Ryu app for collecting OF messages and querying switch features

▸ Want to spend more time with Poseidon

**Sandia National Laboratories**

Stable Infrastructure Network

Faucet.Valve - 1
Faucet.Valve - 2
Faucet.Gauge

"OOBM"
DS

Development & Testing Network

Faucet.Valve - 1
Faucet.Valve - 2

192.51.100.0/24*

VIPs
192.88.99.253/30
192.88.99.249/30
192.51.100.1/24

ACL'd

192.168.0.0/16

VIPs
192.88.99.250/30
192.168.vlanid.1/24

192.88.99.254/30

VM

Corporate Network
Public IP space

Linux Router + NAT

AT

192.88.99.248/30

Traditional Routing

NF

Stack Ports

NF

AR

AR

AR

Infrastructure Cloud

Cloud Tenant 1

Cloud Tenant 2

Cloud Tenant  N

**AT** = Allied Telesis
**AR** = Aruba (HPE)
**DS** = Dumb Switch
**NF** = NoviFlow

Sandia National Laboratories

* 192.51.100.0/24 - IANA Reserved: Documentation (TEST-NET-2)

▸ Nick Buraglio

▸ Brad Cowie

▸ Josh Bailey

▸ mab68 - don't know who you are, but THANK YOU for your commits

▸ Rest of the FAUCET team

▸ Open vSwitch team

▸ SNL networking team - Rick Strong, David Burton, Will Stout

**Sandia National Laboratories**