# NETFLIX

# Open Connect:
# Starting from a Greenfield
# (a mostly Layer 0 talk)

Dave Temkin
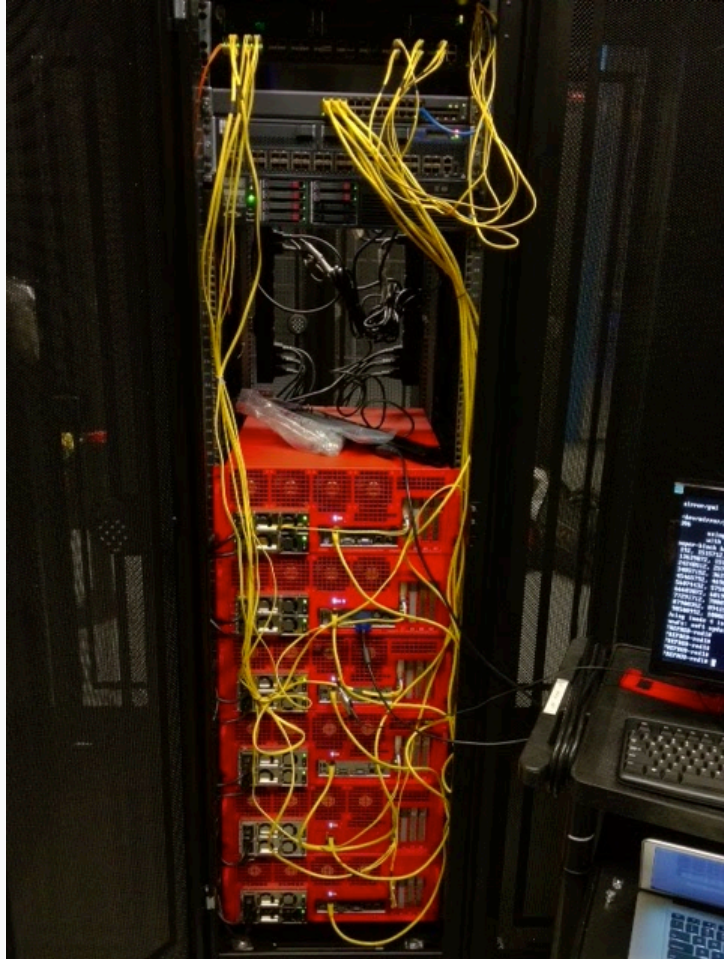06/01/2015

# The story

- **It all started with a discussion...**

- **After much debate, the "Netflix CDN" team was launched in mid-2011**

- **We built caches. They were red. Now they serve lots of Terabits and we call our CDN "Open Connect"**

# A decision had to be made...

- **Do we use caches inside of ISP networks that want them..**
  - And then continue to use third party CDNs

- **Or build a fully functional, standalone CDN**

# A single purpose system

- **Deliver video at the highest quality possible while allowing operators the ability to manage the traffic on their network**

## ISP Network

Each cache has identical content = 80-100% offload

## Small Peering Location

Sharded content ≈ 90+% offload

## Large Peering/Origin Location

Sharded content 100% of active catalog

## AWS S3

All downloadables stored on S3

17

2012

# Now                    and...

## Storage Appliance
- **Still 4U high**
- **~550 watts**
- **288 TB of storage**
- **2x 10G ports**
- **20Gbit/s delivery**

## Flash Appliance
- **1U**
- **~175 watts**
- **24 TB of flash**
- **2x 40G ports**
- **40Gbit/s delivery**

# Cache Types

- **We have two main types of Netflix Caches**
  - Rev H: 36 8TB spinning drives, up to 20Gbit
    - Used for catalog offload
  - Rev I:  24 1TB SSD's, up to 40Gbit
    - Used for high speed popular content serving

- **Our mantra has been to use the same hardware that we would expect an ISP to install in their network**
  - Consistent software stack

**Left: Storage OCA**
**Right: Offload OCA**

# Power Utilization and Footprint

- **Rev H: 560 watts**
  - .31 watts per megabit
- **Rev I: 250 watts**
  - .006 watts per megabit

**Our standard deployment has been 10 Rev H's per rack and 30 Rev I's, or a 5.6kW/7.5kW deployment**

S/N USE2600053OGA07B

P/N NDS4360-05 REV A4

MAC1: 0CC47A45D368

Ten0

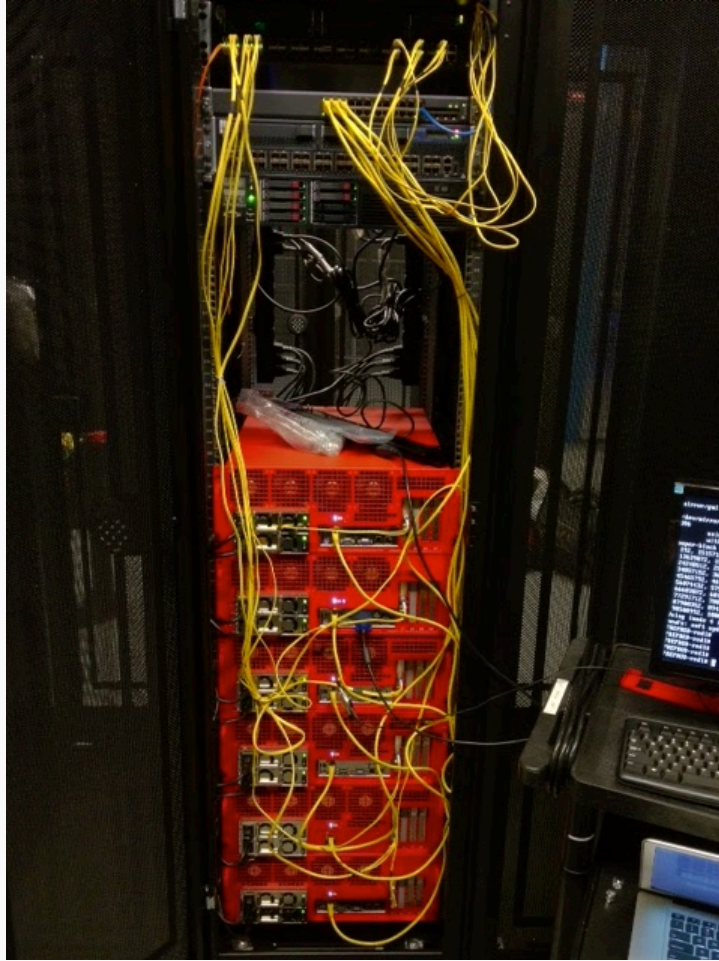Ten1

Ten2

Ten3

DUDE, NO. THIS IS SERIOUS. I JUST SHARTED.

IPMI

LAN1 LAN2
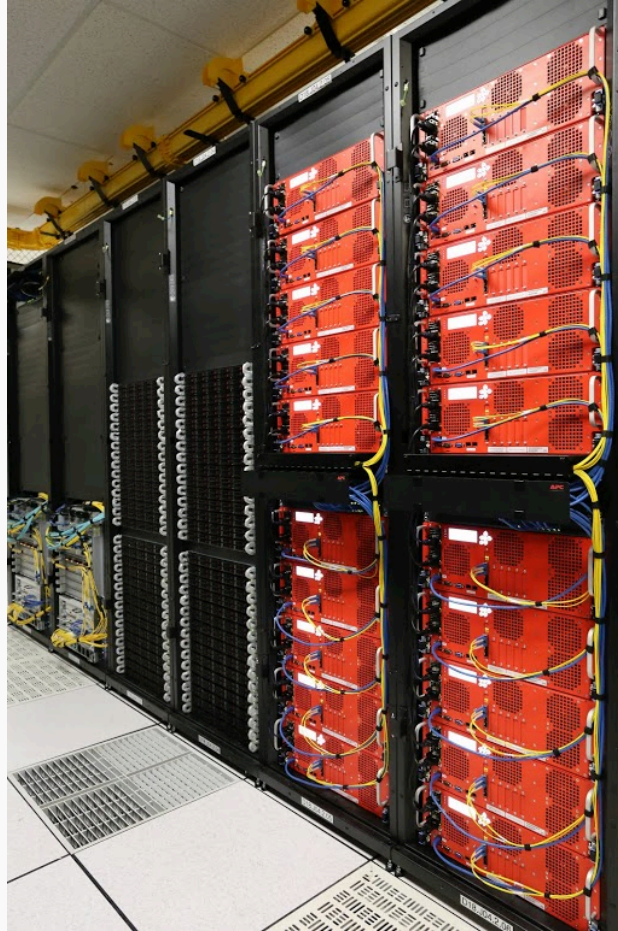
c032.atl001.ix
108.175.40.67
USE06000283GA007

MAC1: 0025907A336E

```
c033.atl001.ix.nflx▒
IP:108.175.40.68[L]&
Mask:255.255.255.192
GW:108.175.40.65v1.6
```

IPMI

LAN1 LAN2

c033.atl001.ix
108.175.40.68
USE06000283GA01A

MAC1: 0025907A336C

```
0: -4.20dBm 10G-LR ▒
1: -4.20dBm 10G-LR &
2: NO INFO
3: NO INFO        v1.6
```

# Today's focus

**Layer Zero: Peering and Interconnection Locations**

# A big challenge

- **Goal was to be off third party CDNs by 2013**
  - At a reasonable cost compared to what we were paying CDNs
    - (Remember, Vertical Integration)
  - At better quality
  - Scalable for global growth

**Our first attempt… ~24G of capacity**

**To ~4Tbits of edge capacity in 4 racks...**

**2.4T of serving capacity at a "small" peering location
35kW of power**

# A note on site selection

- **Pick the most popular site in a metro**
- **Can't find space and power?**
  - Maybe the second most popular

# Challenges

- **5.6kW for Rev H's is relatively easy to get**
  - Fully loaded, we ask for a 6.5kW per rack footprint
- **7.5kW is stretching the limits of most legacy data centers**
- **Rev H's are space limited (most sites can accommodate a 42U rack, so we engineer to that)**
- **Rev I's are power limited**
  - Would like to go from 250 watts to 300 watts
  - End up with a 9kW rack
    - Easy when you own the data center, not so easy when you lease

# Cabling

- **Something taken for granted for many years**
  - "Call up a contractor, have them run some fiber, plug it all in"
- **Not so simple anymore**
- **10G in the data center is still the most affordable**
  - 40G mostly there
  - 100G still prohibitive beyond interconnect

**This doesn't happen by accident… but takes an hour to do.**

1 of 18 custom cable types...

# This next slide was originally going to have a witty GIS'ed image for "Cable Porn"

## But I had turned SafeSearch off and quickly abandoned that idea

# Moving on...

- Somehow we need to get the data out of here..

**Right: 1440 cross connects per rack**
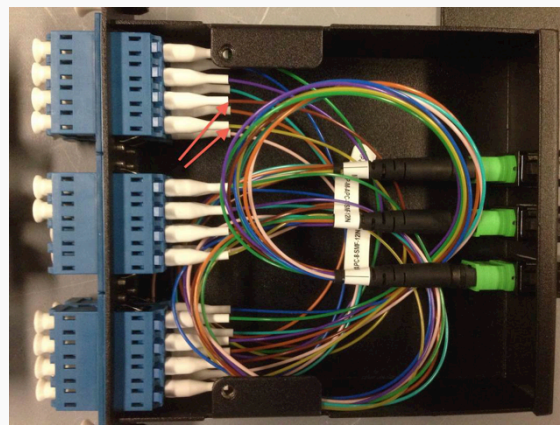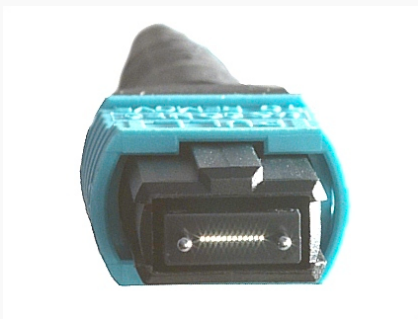**288x10G or 100G = 2.88 or 28.8 Terabits**
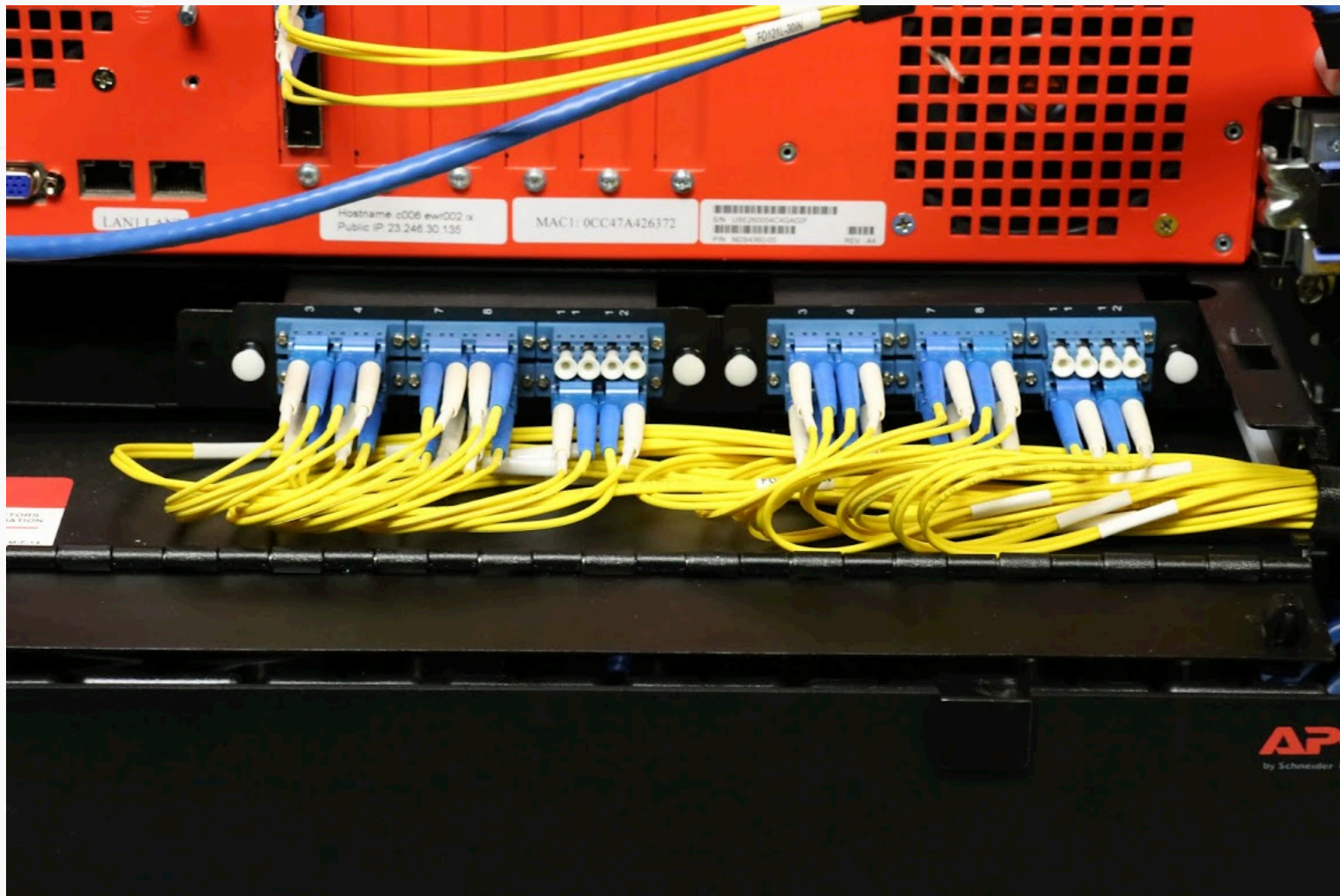**Left: 192 cross connects per rack**
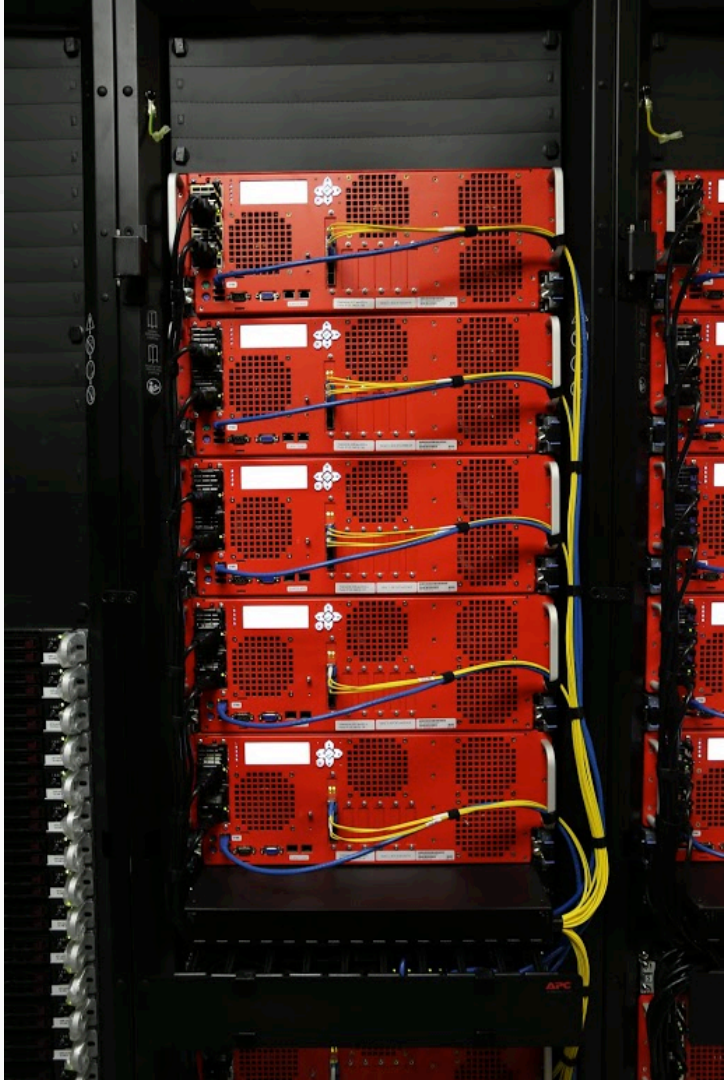
# Rapid Deployment

- Having custom cables to fit our deployment allows for rapid implementation

- Only levers are lengths and types

- Allows for a complete solution to be shipped to site

- If everything goes well, we can have a multi-Terabit site online in less than two days

- Never underestimate the value of not having $colo vendors touch anything other than your patch panels
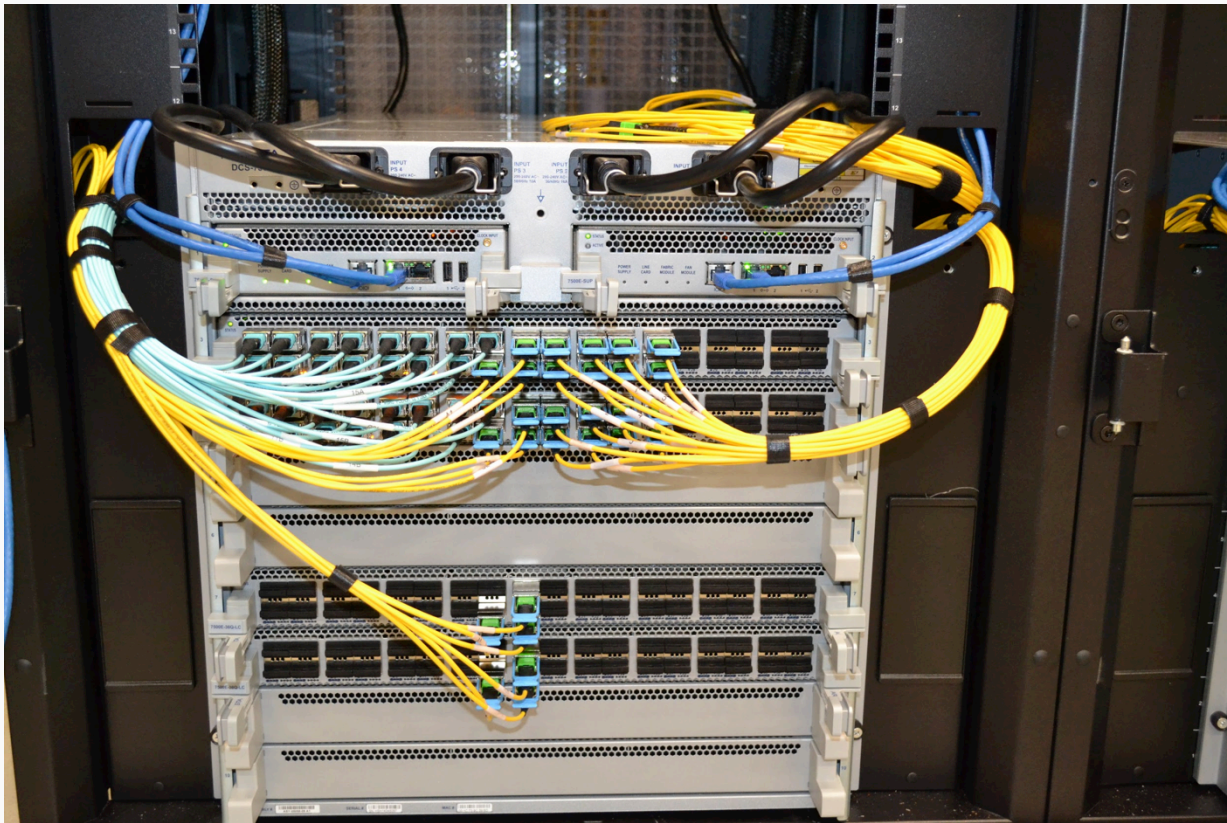
# Leveraging MTP

- **MTP connectors on everything..**
  - Servers (40G)
  - Switches (QSFP), including PLR/PSM 4x10G
  - Patch panels
  - Cassettes
- **Allows for rapid field deployment**
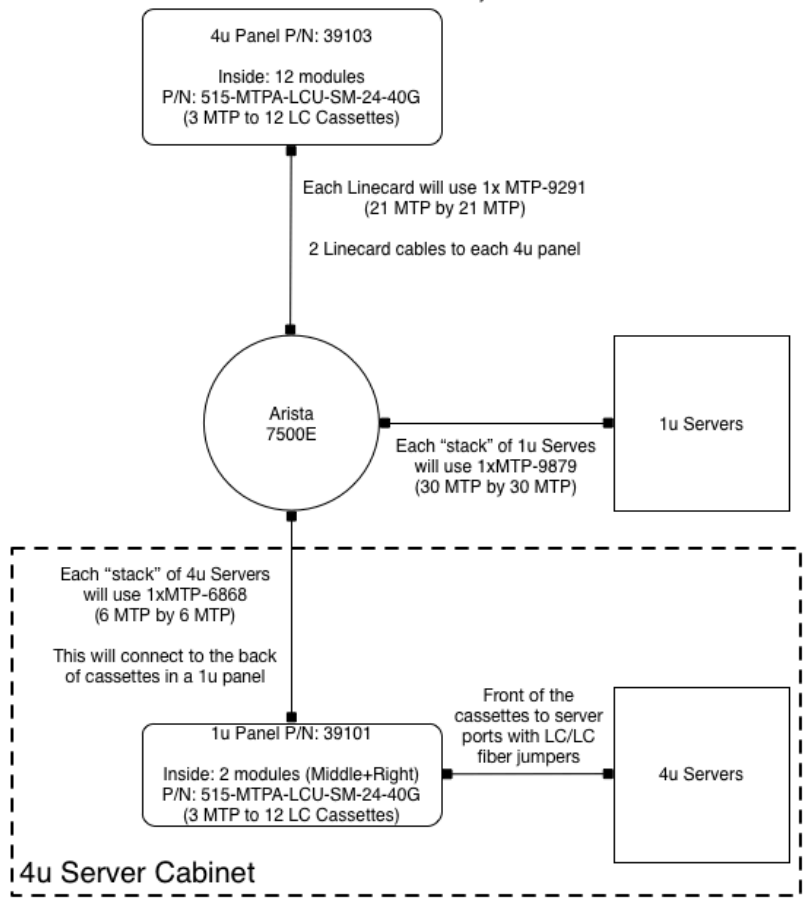- **Reconfigurable - direct path to 100G**
- **Inexpensive**

~3/30T of interconnect

Lowest Linecard cable connects to cassettes in module locations 1-6
Next Linecard cable connects to cassettes in module locations 7-12

Cables will have serial numbers on both ends to identify them

4u Panel P/N: 39103

Inside: 12 modules
P/N: 515-MTPA-LCU-SM-24-40G
(3 MTP to 12 LC Cassettes)

Each Linecard will use 1x MTP-9291
(21 MTP by 21 MTP)

2 Linecard cables to each 4u panel

Arista
7500E

1u Servers

Each "stack" of 1u Serves
will use 1xMTP-9879
(30 MTP by 30 MTP)

Each "stack" of 4u Servers
will use 1xMTP-6868
(6 MTP by 6 MTP)

This will connect to the back
of cassettes in a 1u panel

Front of the
cassettes to server
ports with LC/LC
fiber jumpers

1u Panel P/N: 39101

Inside: 2 modules (Middle+Right)
P/N: 515-MTPA-LCU-SM-24-40G
(3 MTP to 12 LC Cassettes)

4u Servers

4u Server Cabinet

**Every site has the same layout**

# Homogeny
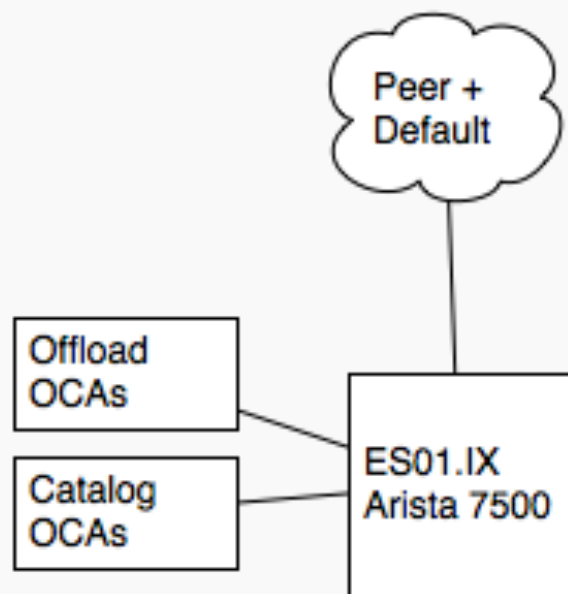
- **Allows us to make rapid deployment decisions**
  - Standardized negotiating for space and power depending on forecast
  - Quick Bill of Material generation
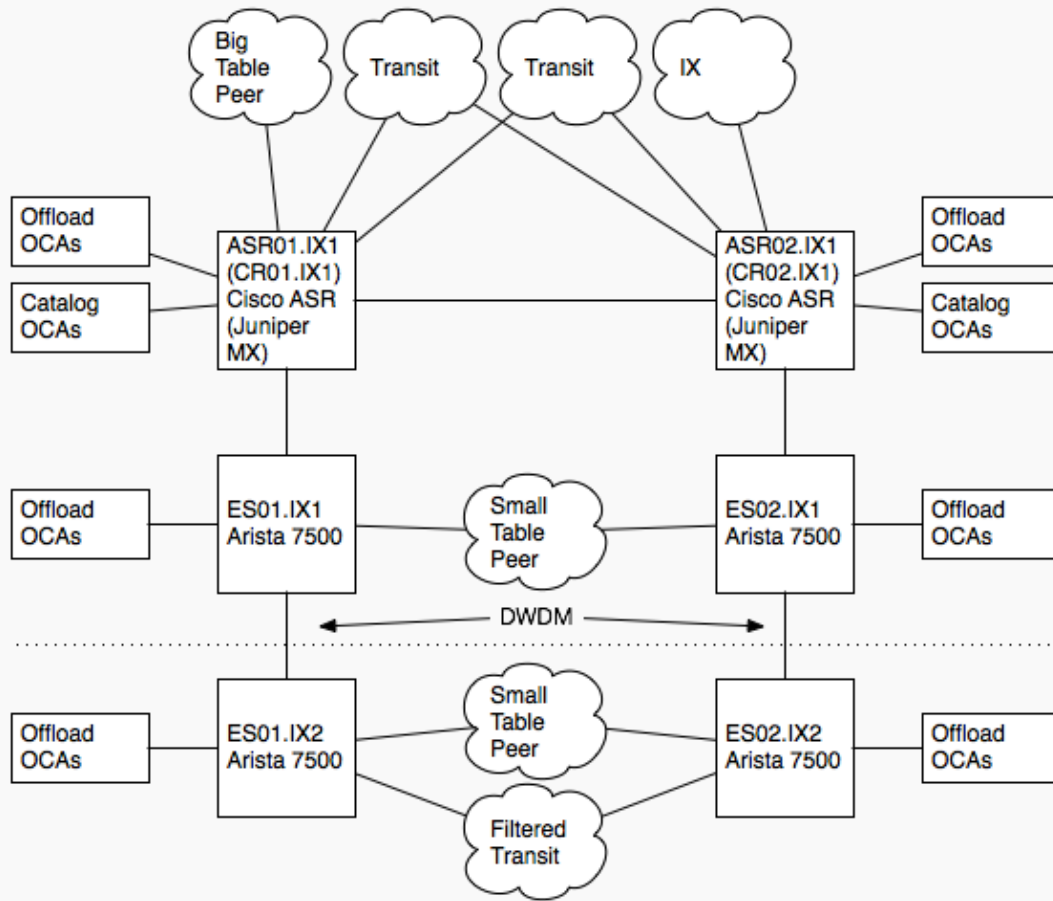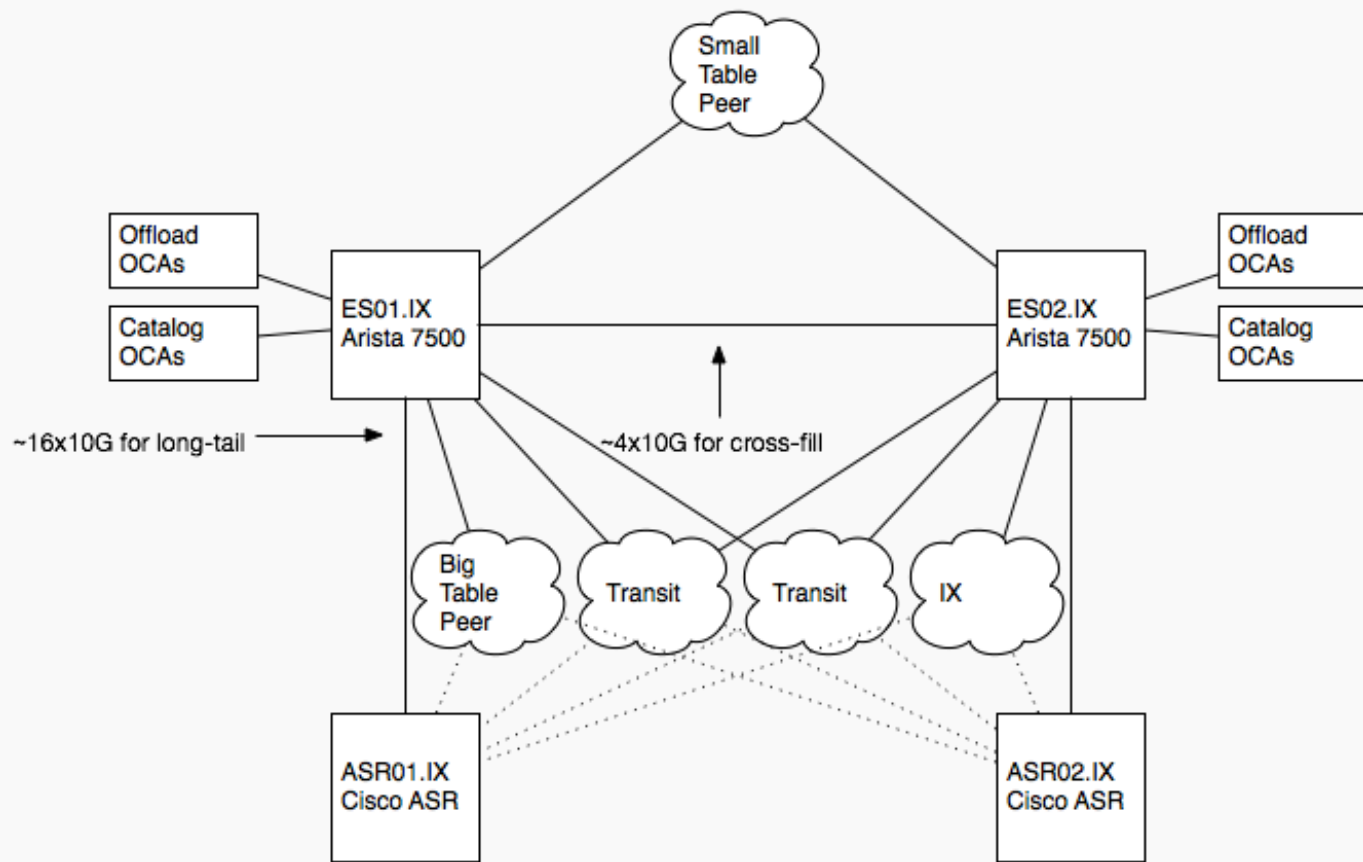  - Signed to live in less than 30 days

2013

# Layer 2/3

- **No opportunity for aggregation**
- **Big chassis is best**
- **Sticking to off-the-shelf platforms (for now!)**
  - Better to focus on software
- **Developing our own routing platform**
  - No longer buying big expensive routers
  - We've had a traffic management platform for 8+ years

**And a similar network architecture (add and remove pieces)**

# More info?

- **[http://openconnect.netflix.com](http://openconnect.netflix.com)**
- **[dtemkin@netflix.com](mailto:dtemkin@netflix.com)**

# Questions?