

NETFLIX

Serving Netflix Video Traffic at 400Gb/s and Beyond

Drew Gallatin
NAB Show, April 2022

Serving Netflix Video Traffic at ~~400Gb/s~~ and Beyond *800Gb/s*

Netflix Workload:

- Serve only static media files
- Pre-encoded for all codecs/bitrates
 - Video quality is of the utmost importance, so we don't transcode on the server
- Greatly simplifies server workload

Netflix Video Serving Stack

- FreeBSD-current
- NGINX web server
- Video served via asynchronous sendfile(2) and encrypted using kTLS

Timeline:

- *Asynchronous Sendfile (2014)*
- Kernel TLS (2016)
- Network-centric NUMA (2019)
- Inline Hardware (NIC) kTLS (2022)
- 800G initial results

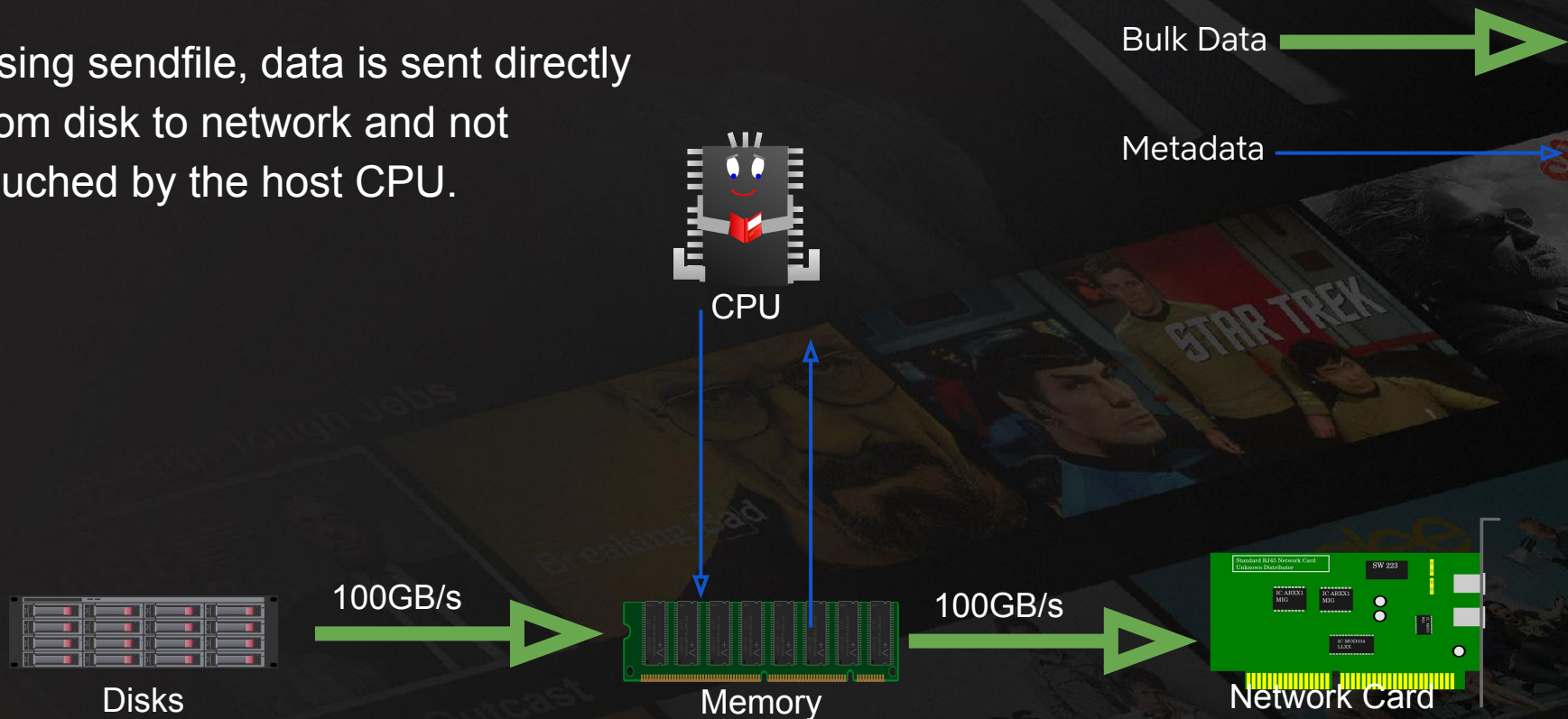
Sendfile

- Since we are serving static files, we can use sendfile(2)
- Sendfile directs the kernel to send data from a file descriptor to a TCP socket
- This eliminates the need to copy data into or out of the kernel

NETFLIX

Netflix Video Serving Data Flow

Using sendfile, data is sent directly from disk to network and not touched by the host CPU.



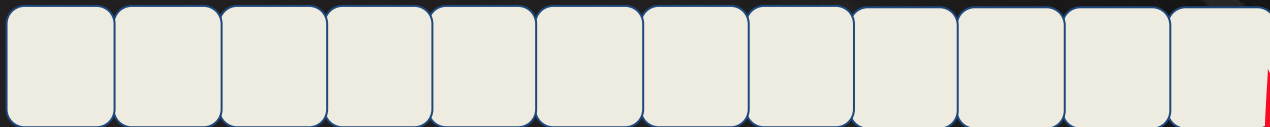
Problem: Disk reads can block sendfile

- When an nginx worker is blocked, it cannot service other requests
- Solutions to prevent nginx from blocking like aio or thread pools scale poorly

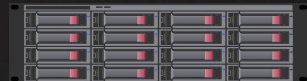
Solution: Asynchronous sendfile

- sendfile() becomes “fire and forget”
- Empty buffers are appended to the TCP socket buffer. TCP stops when it sees an empty buffer.
- When disk read completes, disk interrupt handler informs TCP it is ready to send

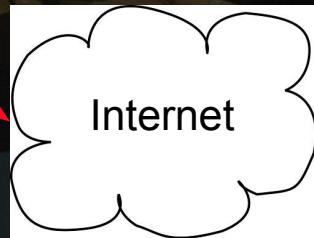
Asynchronous sendfile



Socket Buffer

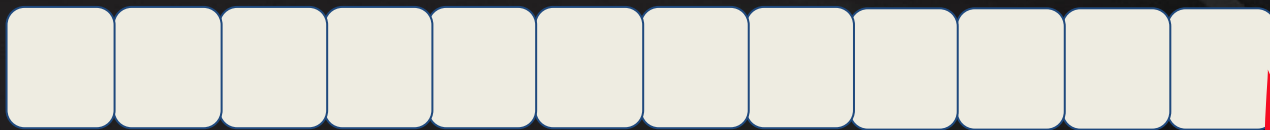


Disks

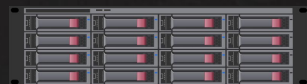


NETFLIX

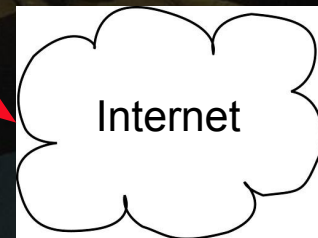
Asynchronous sendfile



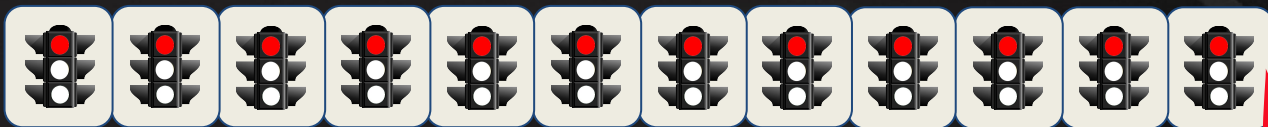
Socket Buffer



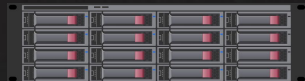
Disks



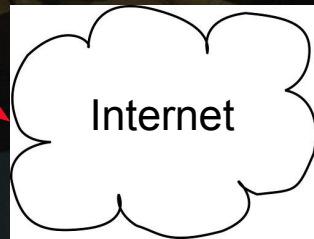
Asynchronous sendfile



Socket Buffer



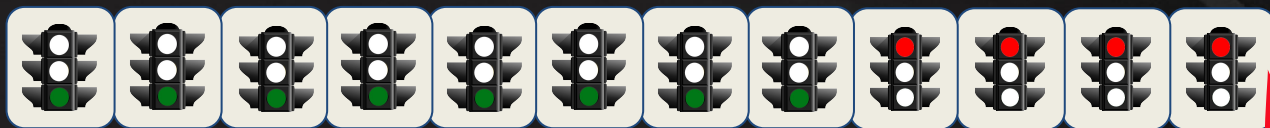
Disks



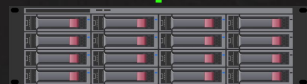
Internet

NETFLIX

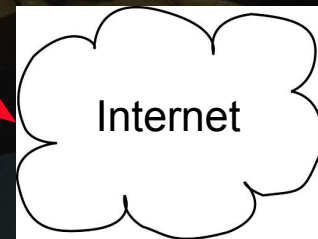
Asynchronous sendfile



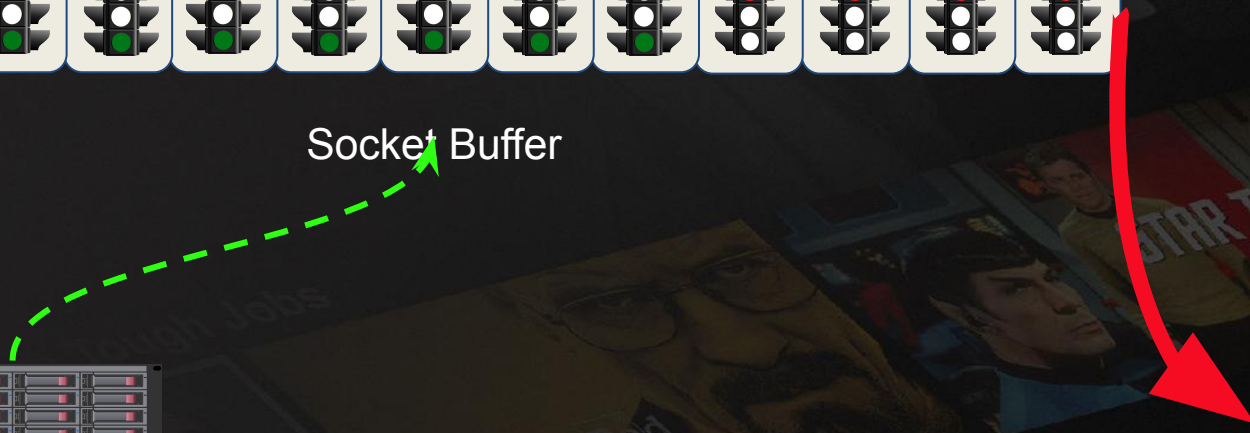
Socket Buffer



Disks

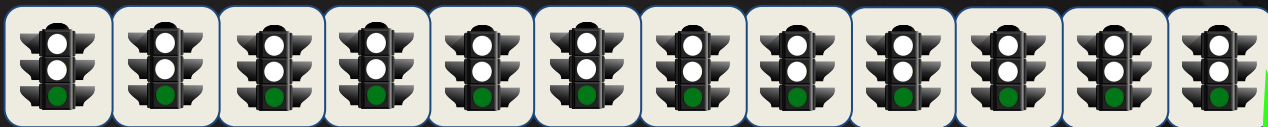


Internet

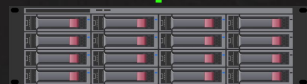


NETFLIX

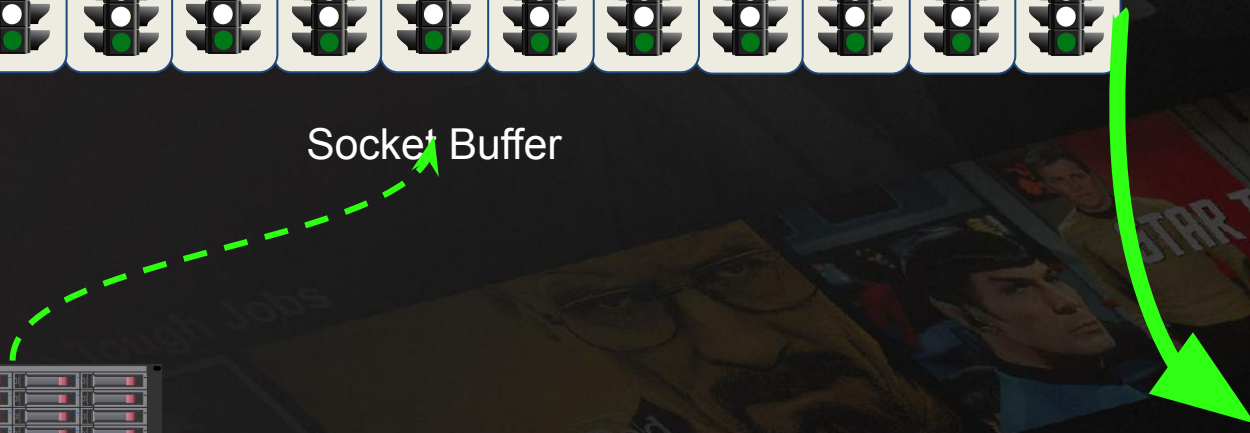
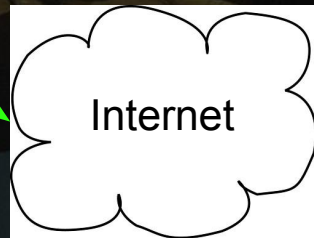
Asynchronous sendfile



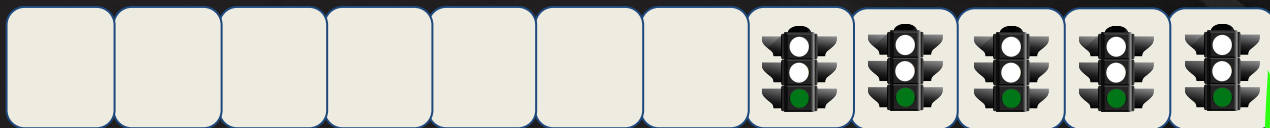
Socket Buffer



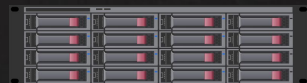
Disks



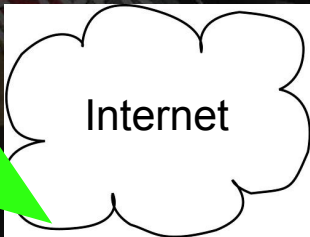
Asynchronous sendfile



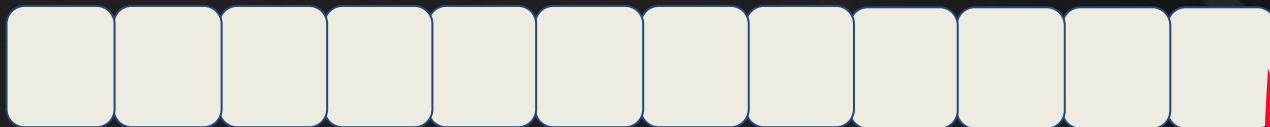
Socket Buffer



Disks



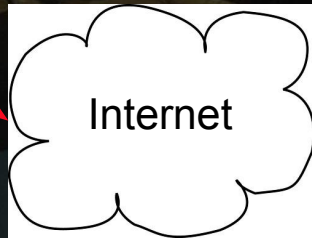
Asynchronous sendfile



Socket Buffer

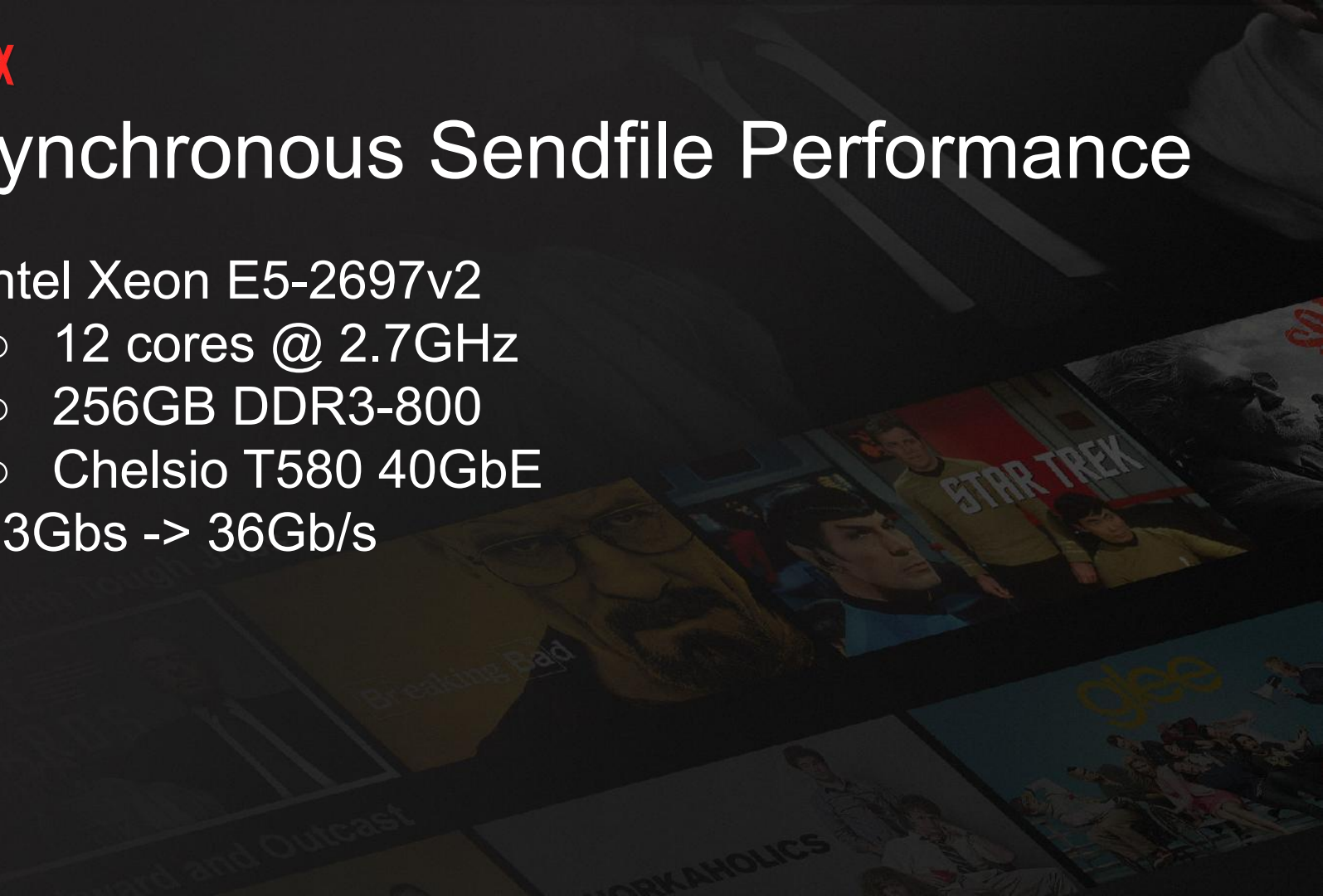


Disks



Asynchronous Sendfile Performance

- Intel Xeon E5-2697v2
 - 12 cores @ 2.7GHz
 - 256GB DDR3-800
 - Chelsio T580 40GbE
- 23Gbs -> 36Gb/s

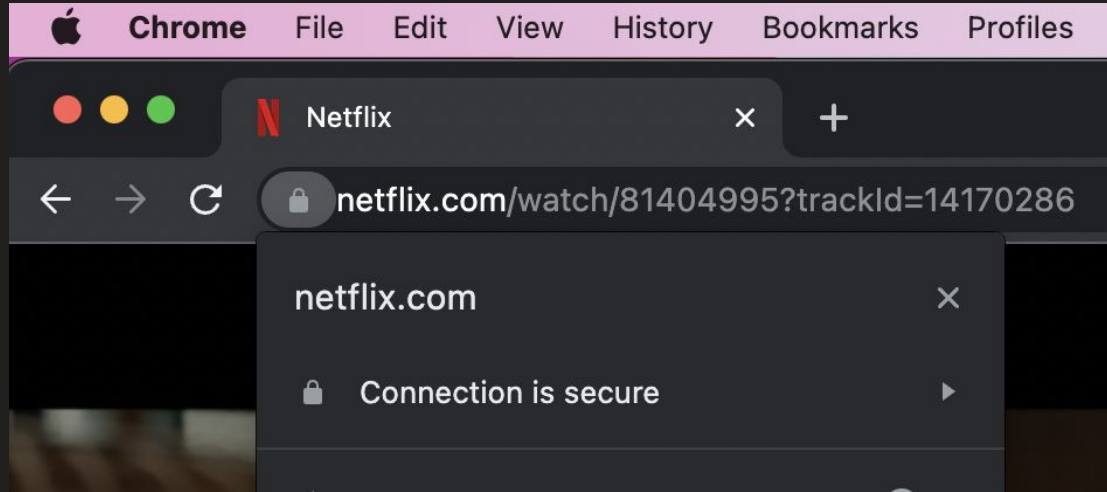


Timeline:

- Asynchronous Sendfile (2014)
- ***Kernel TLS (2016)***
- Network-centric NUMA (2019)
- Inline Hardware (NIC) kTLS (2022)
- 800G initial results

What's TLS?

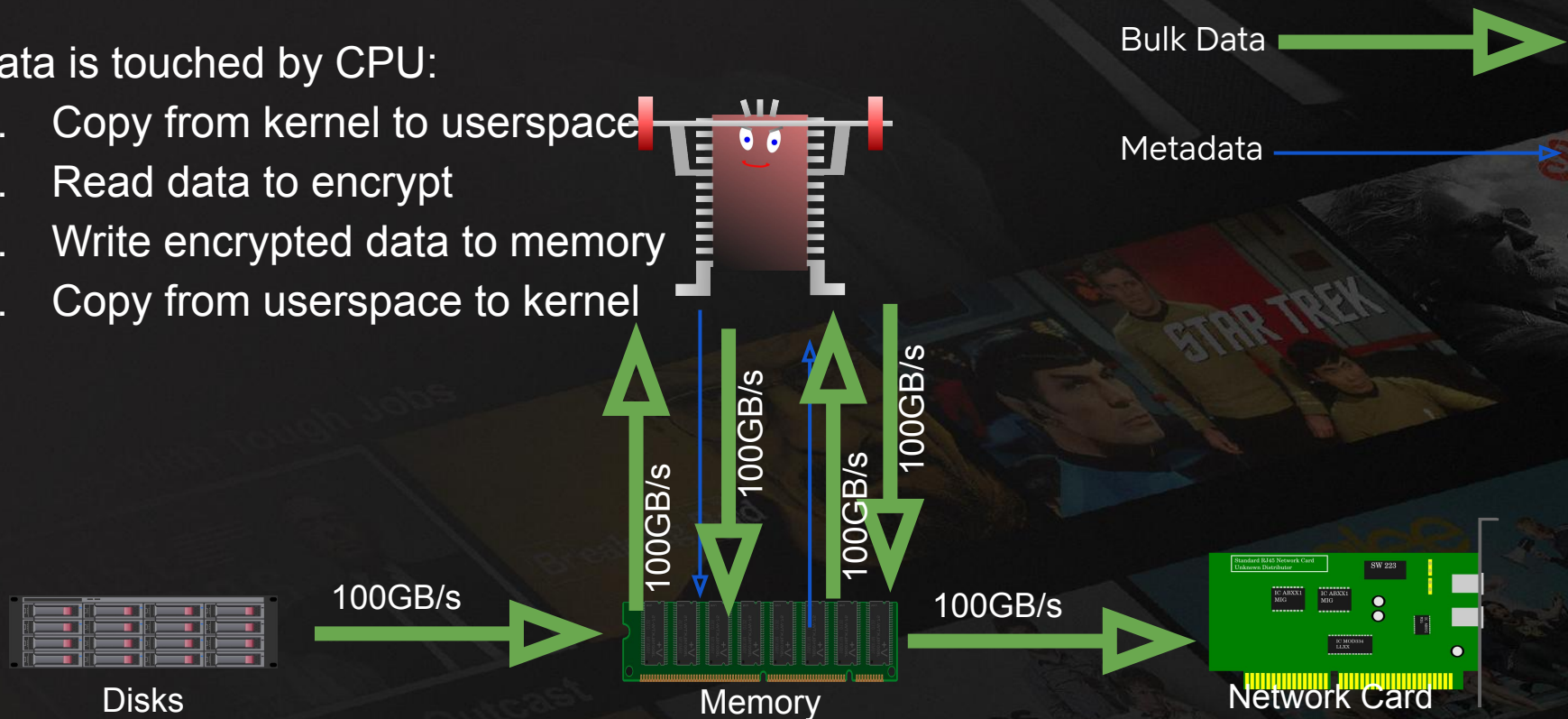
- **T**ransport **L**ayer **S**ecurity
- **TLS** encrypts traffic between clients and the OCA



TLS Prevents Sendfile & Triples Memory BW

Data is touched by CPU:

1. Copy from kernel to userspace
2. Read data to encrypt
3. Write encrypted data to memory
4. Copy from userspace to kernel

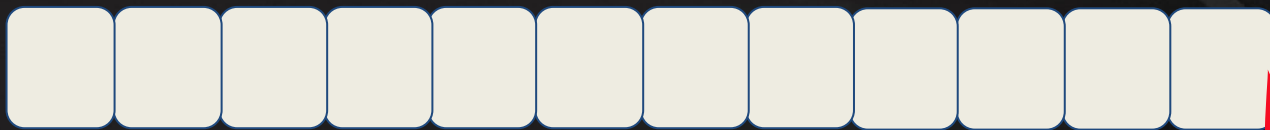


Solution: Move TLS into the kernel

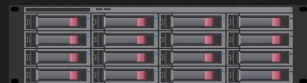
- Eliminates copies between userspace and kernel
- Restores sendfile dataflow
- TLS handshakes (eg, session setup, session resumption) handled in userspace.
- TLS state handed to the kernel
- Kernel does bulk crypto as part of sendfile pipeline

NETFLIX

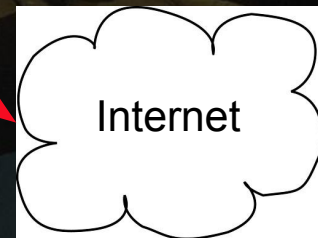
Asynchronous sendfile + kTLS



Socket Buffer



Disks



NETFLIX

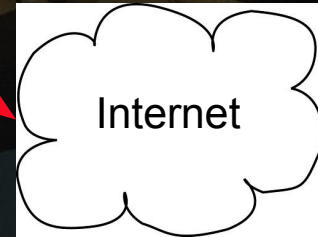
Asynchronous sendfile



Socket Buffer

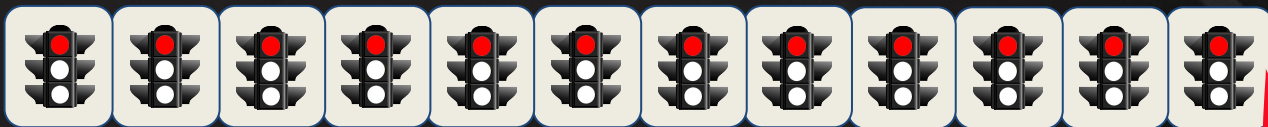


Disks

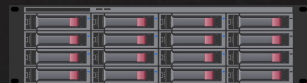


NETFLIX

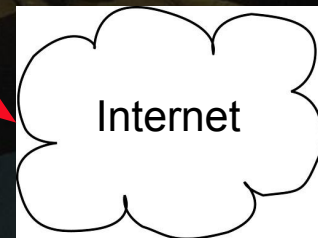
Asynchronous sendfile + kTLS



Socket Buffer

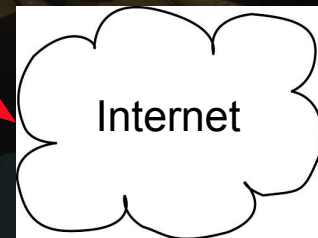
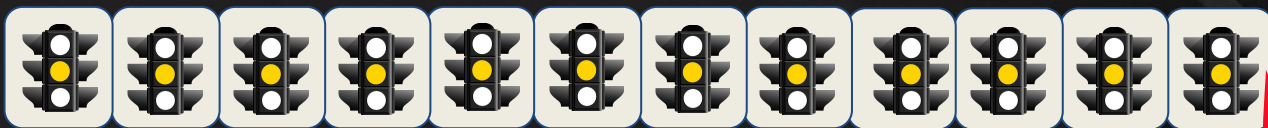


Disks



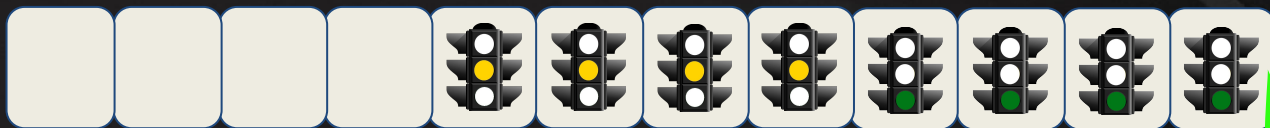
NETFLIX

Asynchronous sendfile + kTLS

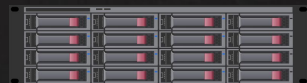


NETFLIX

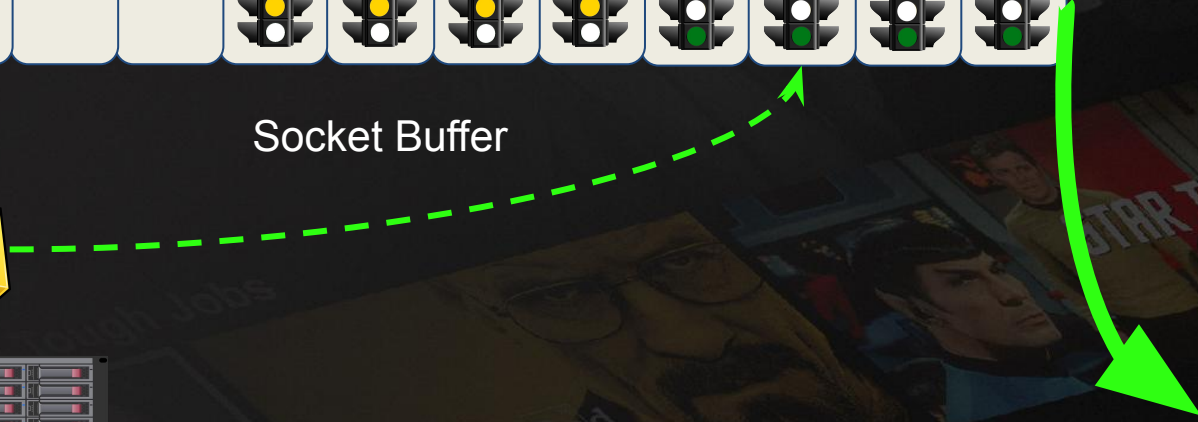
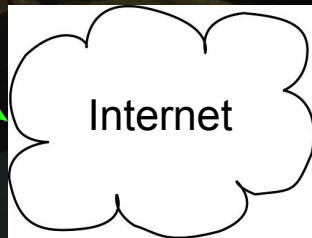
Asynchronous sendfile + kTLS



Socket Buffer

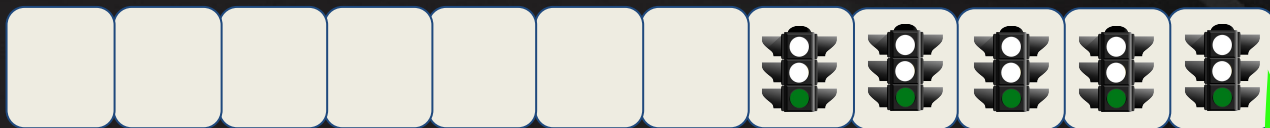


Disks

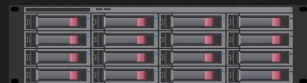


NETFLIX

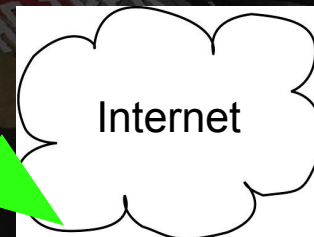
Asynchronous sendfile + kTLS



Socket Buffer

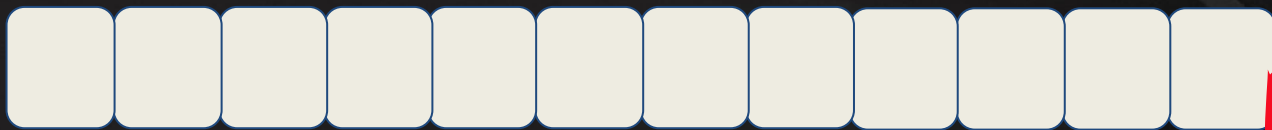


Disks

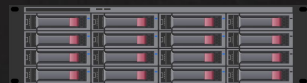


NETFLIX

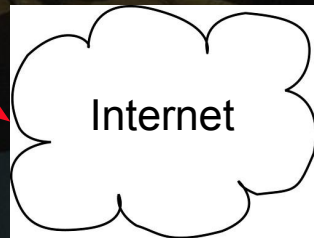
Asynchronous sendfile + kTLS



Socket Buffer



Disks



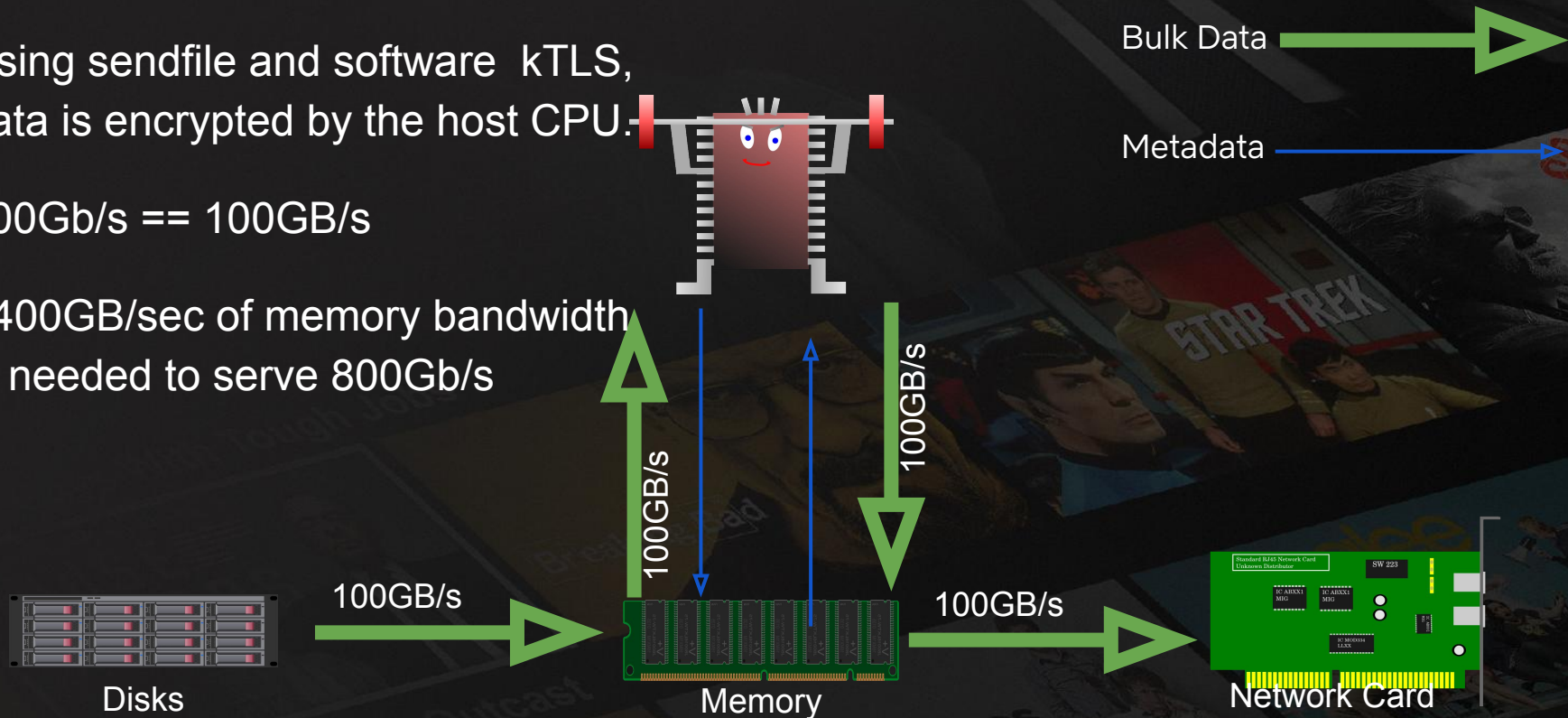
NETFLIX

Netflix 800Gb/s Video Serving Data Flow

Using sendfile and software kTLS, data is encrypted by the host CPU.

800Gb/s == 100GB/s

~400GB/sec of memory bandwidth is needed to serve 800Gb/s



Timeline:

- Asynchronous Sendfile (2014)
- Kernel TLS (2016)
- ***NUMA (2019)***
- Inline Hardware (NIC) kTLS (2022)
- 800G initial results

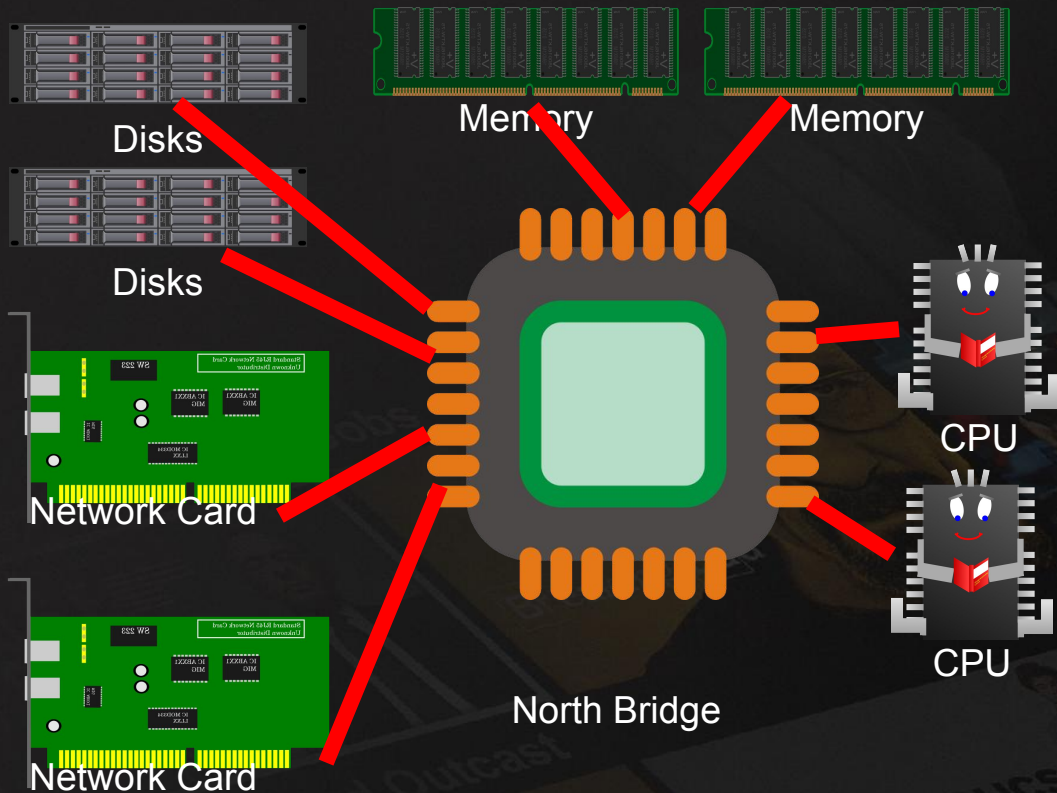
NETFLIX

What is NUMA?

Non Uniform Memory Architecture

That means memory and/or devices can be “closer” to some CPU cores

Multi CPU Before NUMA



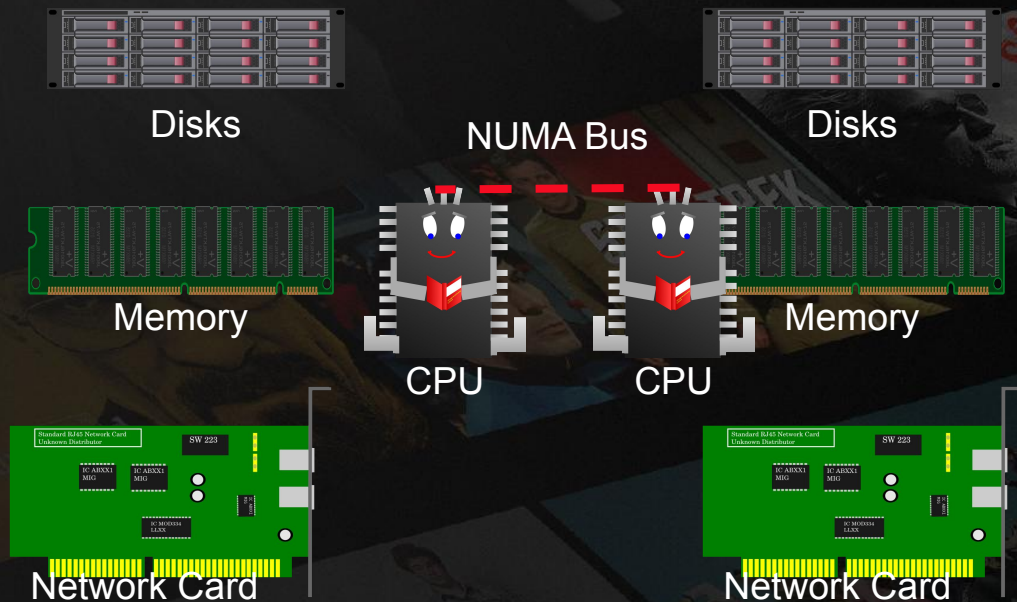
Memory access was *UNIFORM*:

Each core had equal and direct access to all memory and IO devices.

Multi Socket system with NUMA:

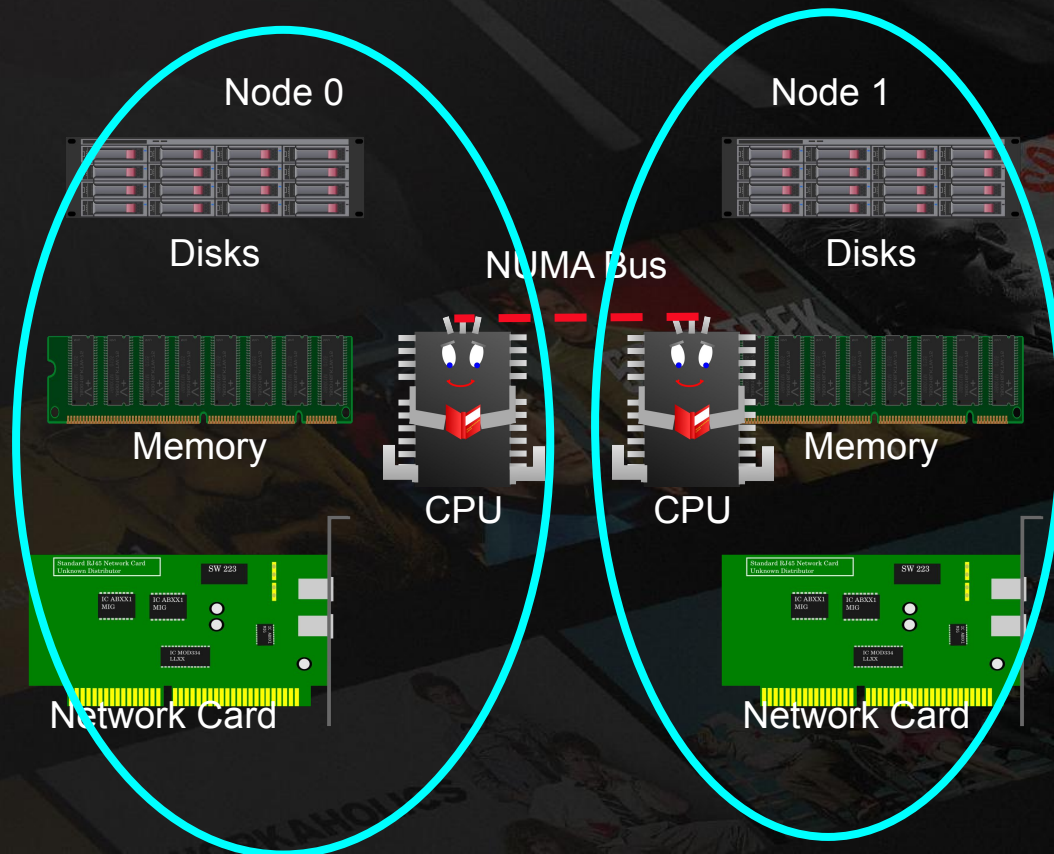
Memory access can be **NON-UNIFORM**

- Each core has unequal access to memory
- Each core has unequal access to I/O devices



Present day NUMA:

Each locality zone called a “NUMA Domain” or “NUMA Node”



NETFLIX

Strategy: Keep as much of our 400GB/sec of bulk data off the NUMA fabric as possible

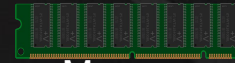
- Bulk data congests NUMA fabric and leads to CPU stalls.

NETFLIX Dual AMD: Worst Case Data Flow

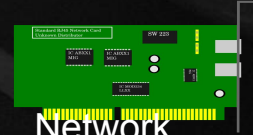
Steps to send data:



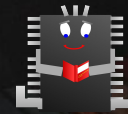
Disks



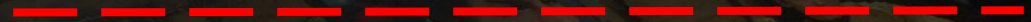
Memory



Network
Card



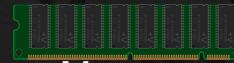
CPU



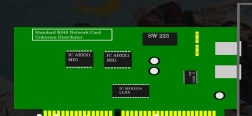
CPU



Disks



Memory

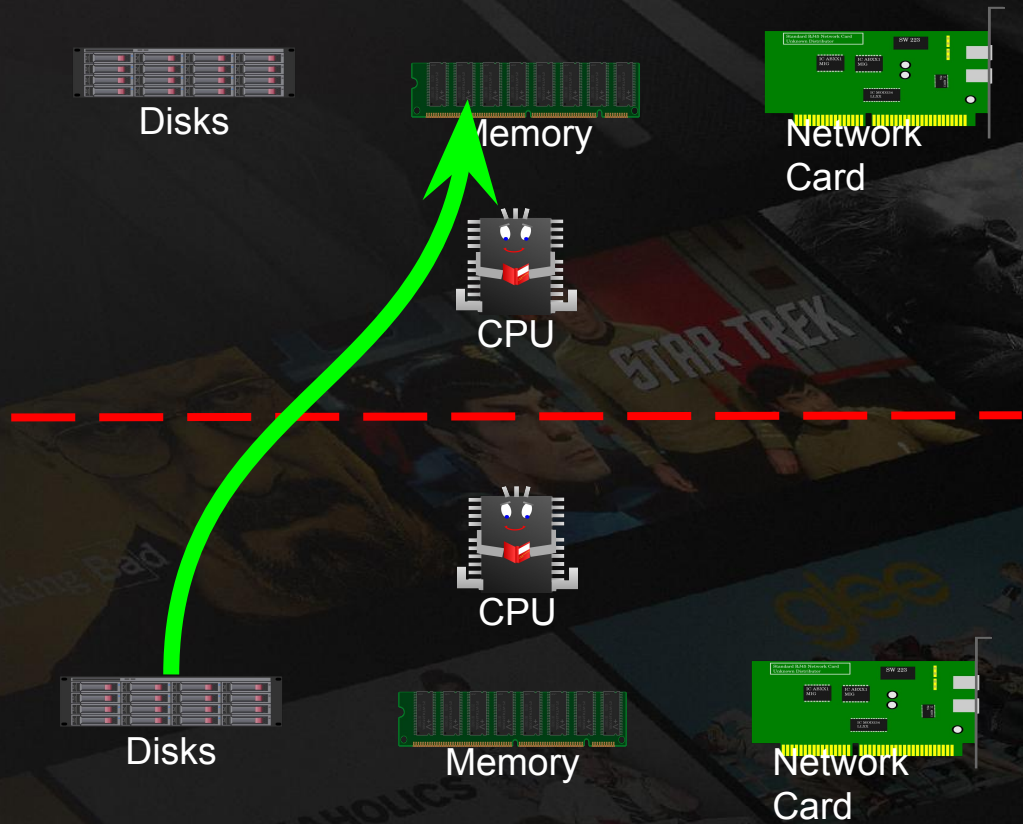


Network
Card

NETFLIX Dual AMD: Worst Case Data Flow

Steps to send data:

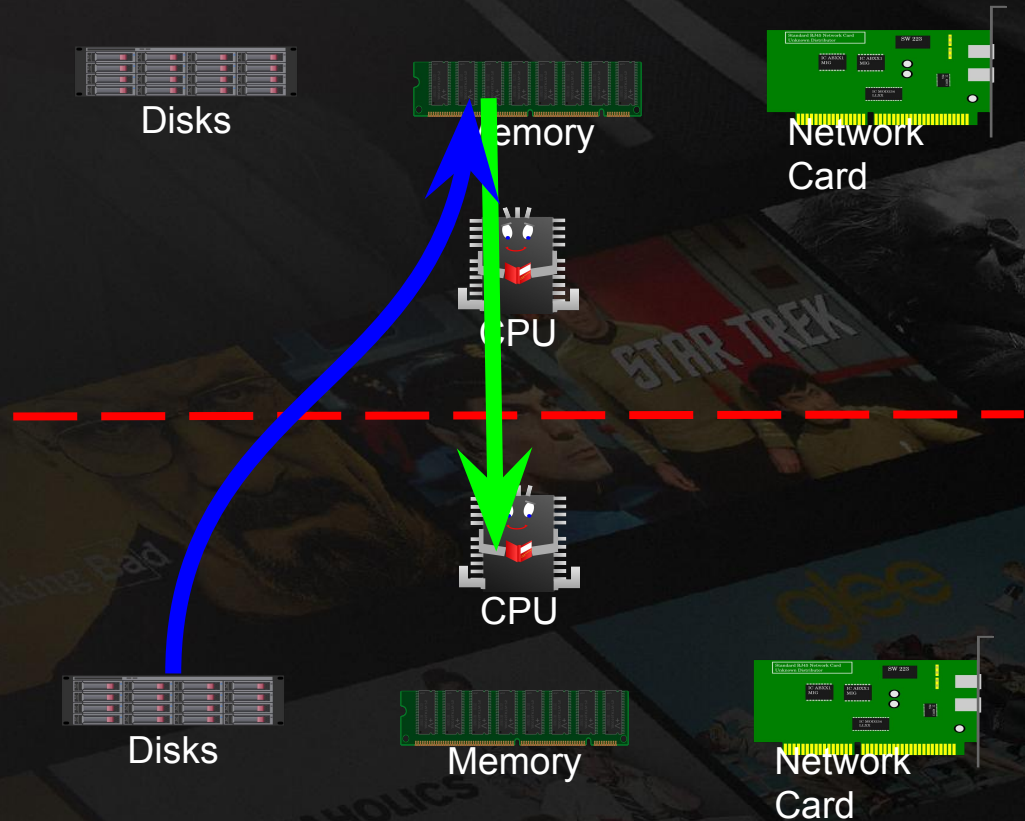
- DMA data from disk to memory
 - First NUMA bus crossing



NETFLIX Dual AMD: Worst Case Data Flow

Steps to send data:

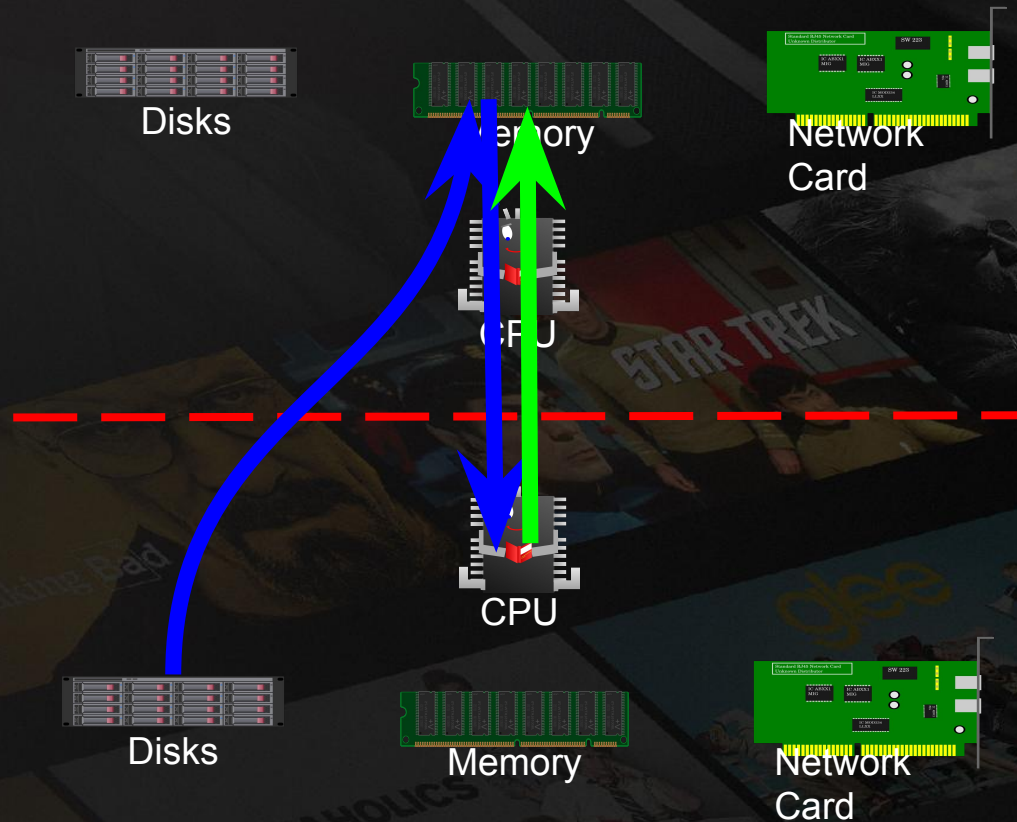
- DMA data from disk to memory
 - First NUMA bus crossing
- CPU reads data for encryption
 - Second NUMA crossing



NETFLIX Dual AMD: Worst Case Data Flow

Steps to send data:

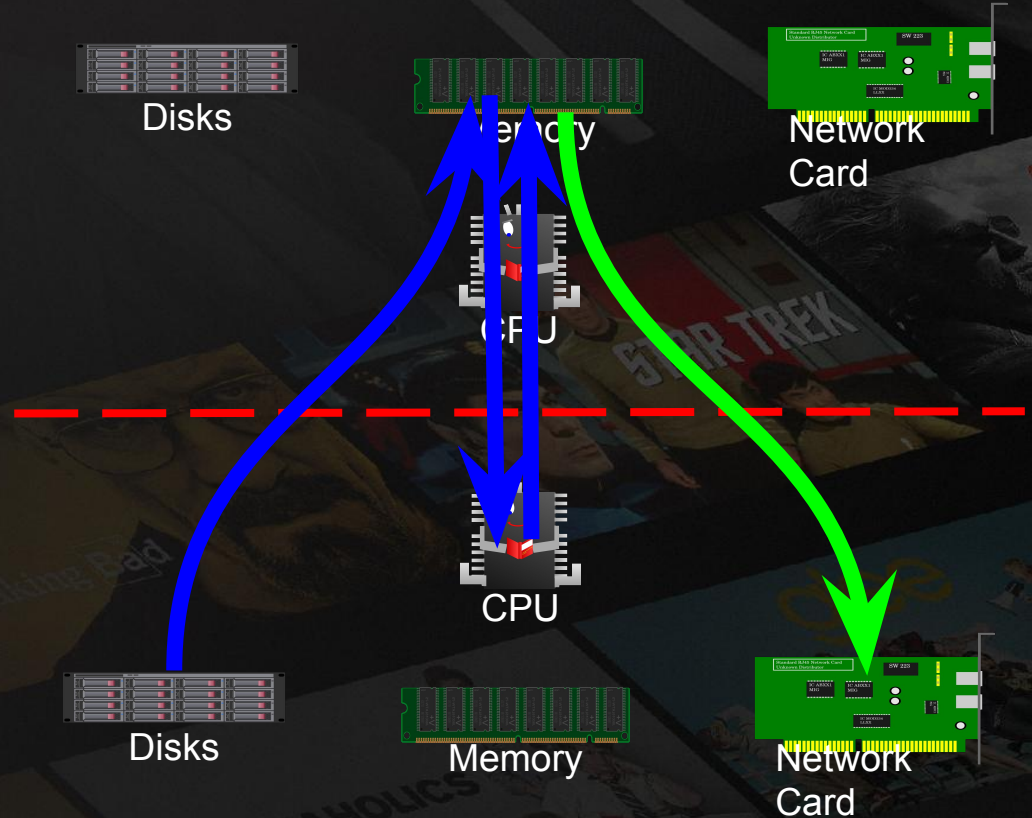
- DMA data from disk to memory
 - First NUMA bus crossing
- CPU reads data for encryption
 - Second NUMA crossing
- CPU writes encrypted data
 - Third NUMA crossing



NETFLIX Dual AMD: Worst Case Data Flow

Steps to send data:

- DMA data from disk to memory
 - First NUMA bus crossing
- CPU reads data for encryption
 - Second NUMA crossing
- CPU writes encrypted data
 - Third NUMA crossing
- DMA from memory to network
 - Fourth NUMA crossing

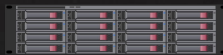


Worst Case Summary:

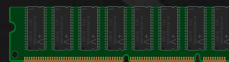
- 4 NUMA crossings
- 400GB/s of data on the NUMA fabric
 - Fabric saturates, cannot handle the load.
 - CPU Stalls, saturates early

NETFLIX Dual AMD: Best Case Data Flow

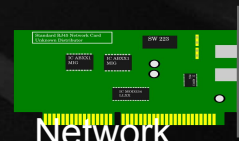
Steps to send data:



Disks



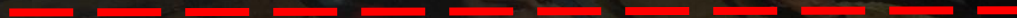
Memory



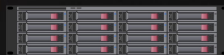
Network
Card



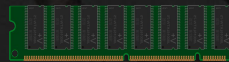
CPU



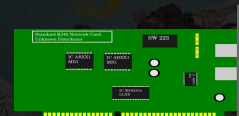
CPU



Disks



Memory

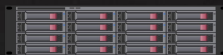


Network
Card

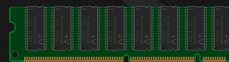
NETFLIX Dual AMD: Best Case Data Flow

Steps to send data:

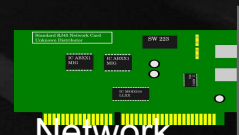
- DMA data from disk to memory



Disks



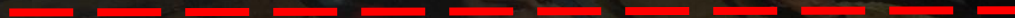
Memory



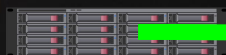
Network
Card



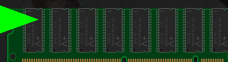
CPU



CPU



Disks



Memory



Network
Card

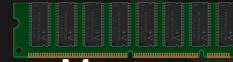
NETFLIX Dual AMD: Best Case Data Flow

Steps to send data:

- DMA data from disk to memory
- CPU Reads data for encryption



Disks



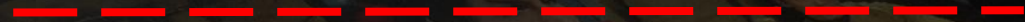
Memory



Network
Card



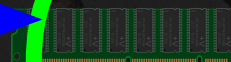
CPU



CPU



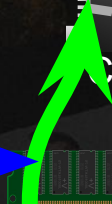
Disks



Memory



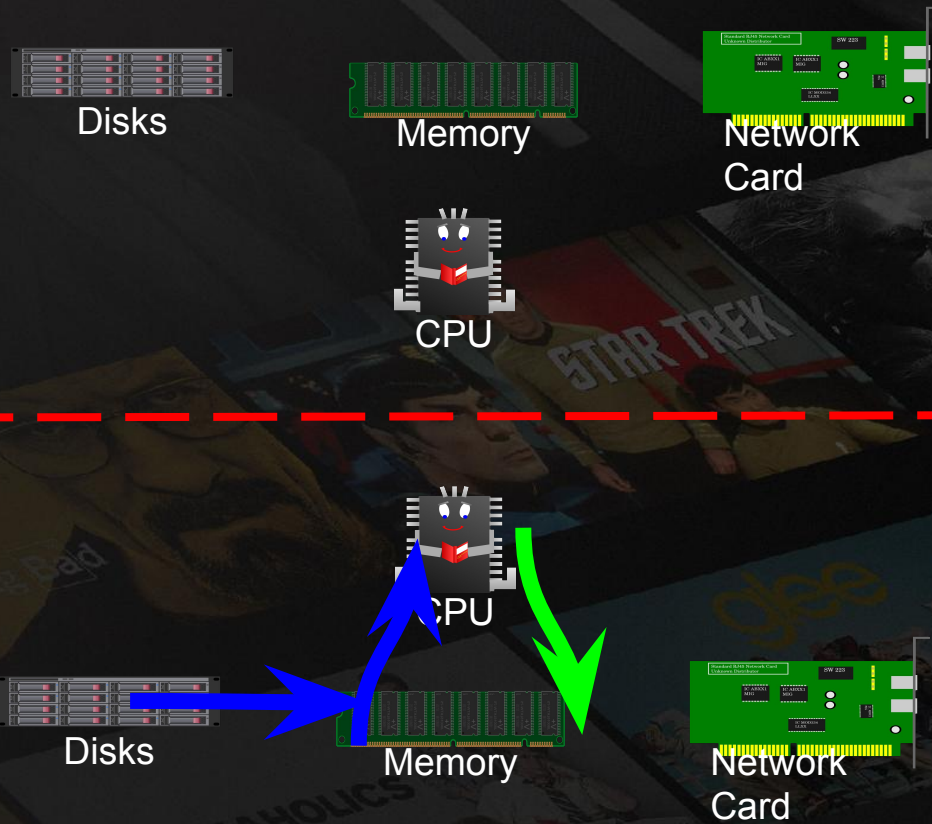
Network
Card



NETFLIX Dual AMD: Best Case Data Flow

Steps to send data:

- DMA data from disk to memory
- CPU Reads data for encryption
- CPU Writes encrypted data

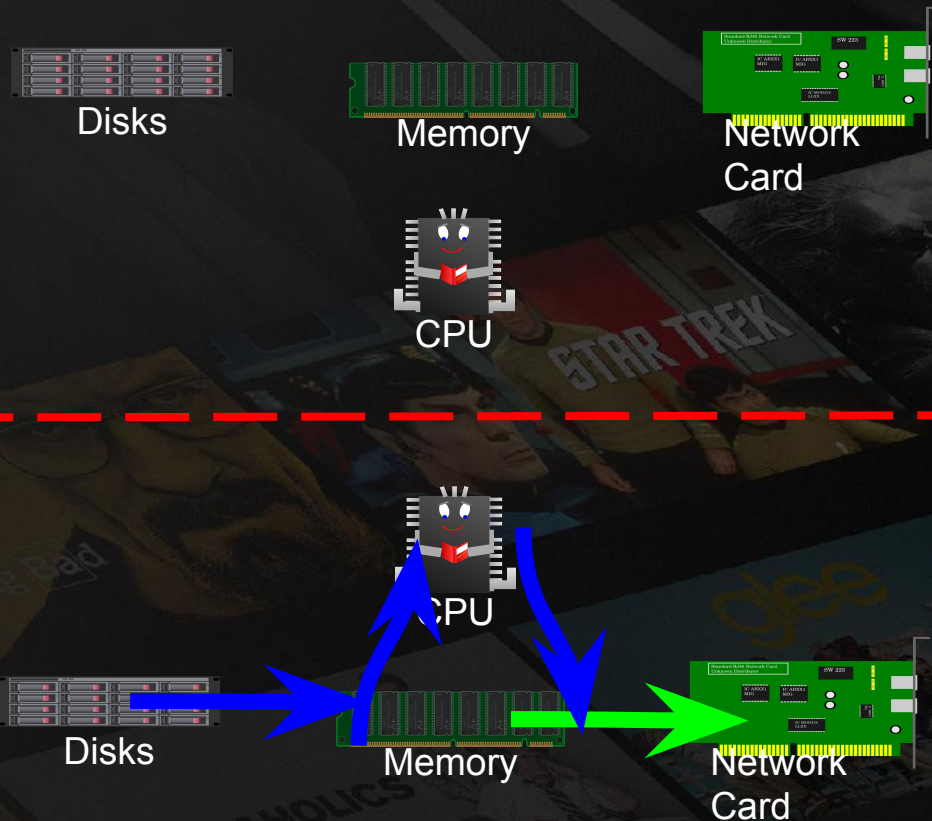


NETFLIX Dual AMD: Best Case Data Flow

Steps to send data:

- DMA data from disk to memory
- CPU Reads data for encryption
- CPU Writes encrypted data
- **DMA from memory to Network**

0 NUMA crossings!



Best Case Summary:

- 0 NUMA crossings
- 0GB/s of data on the NUMA fabric

Impose order on the chaos.. *somehow:*

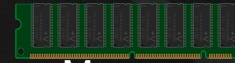
- Disk centric siloing
 - Try to do everything on the NUMA node where the content is stored
- Network centric siloing
 - Try to do as much as we can on the NUMA node that the LACP partner chose for us

NETFLIX Dual AMD: Worst Case Data Flow With Network Centric NUMA Siloing

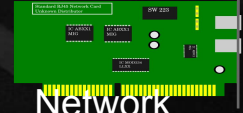
Steps to send data:



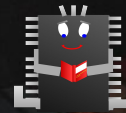
Disks



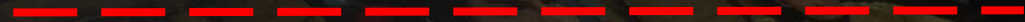
Memory



Network
Card



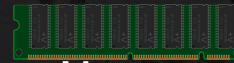
CPU



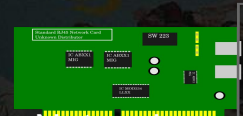
CPU



Disks



Memory

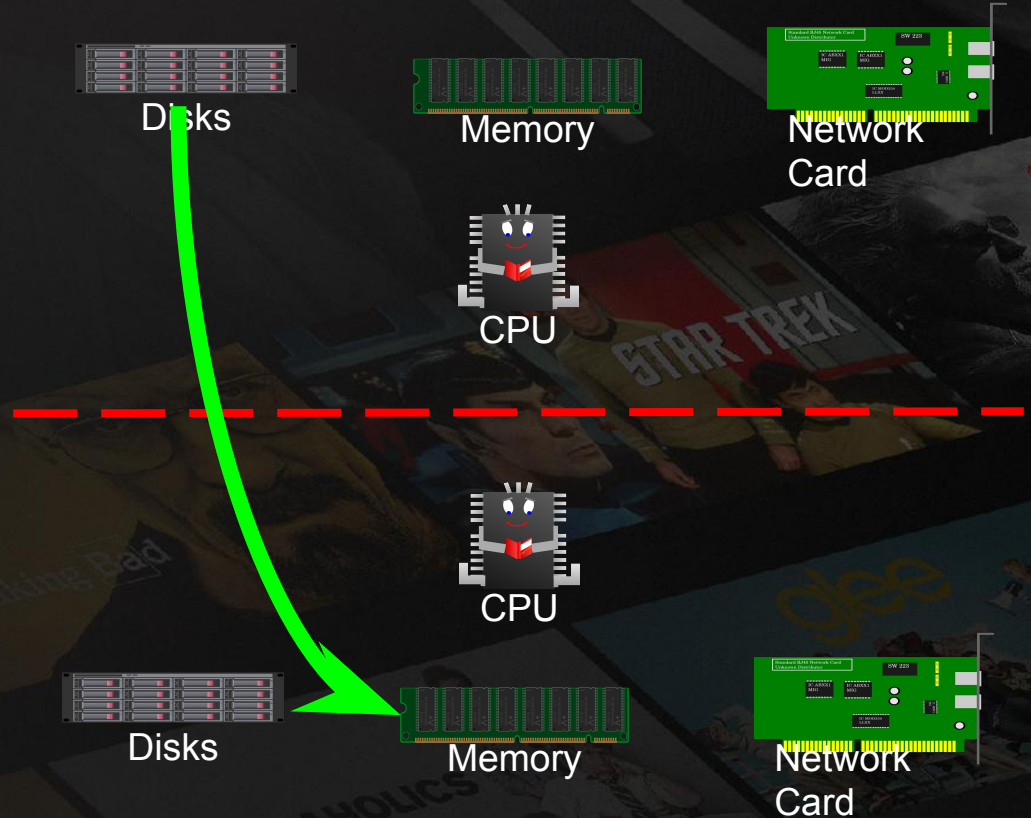


Network
Card

NETFLIX Dual AMD: Worst Case Data Flow With Network Centric NUMA Siloing

Steps to send data:

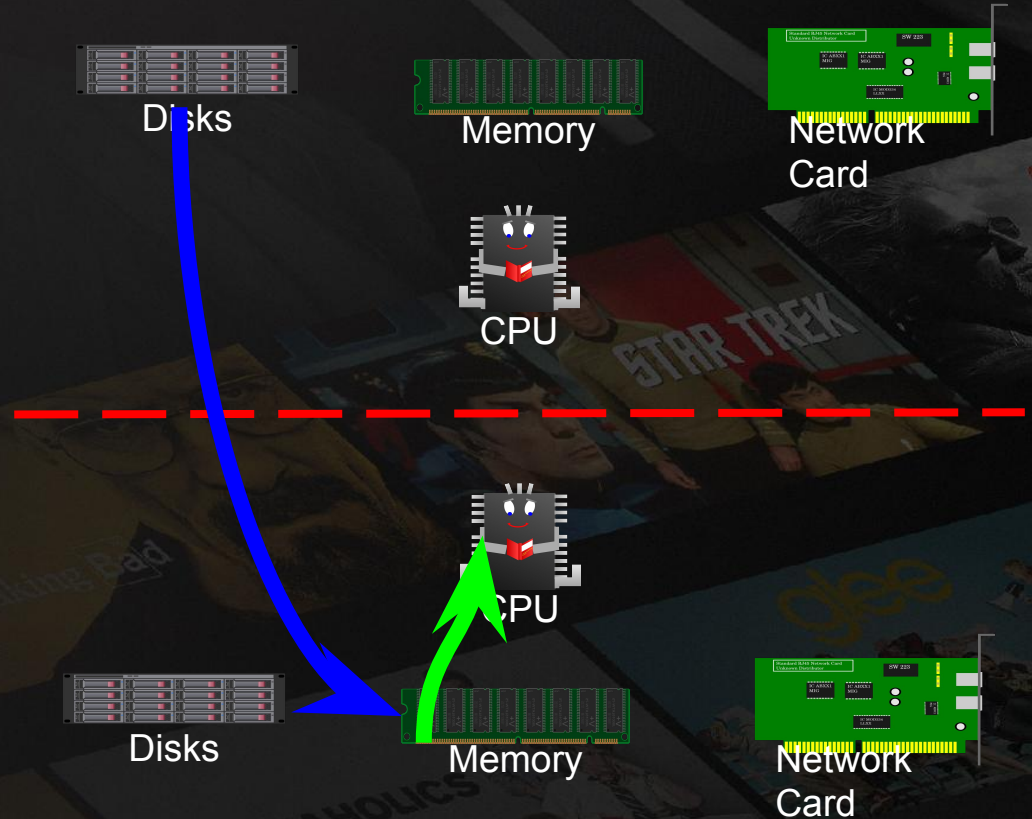
- DMA data from disk to memory
 - First NUMA bus crossing



NETFLIX Dual AMD: Worst Case Data Flow With Network Centric NUMA Siloing

Steps to send data:

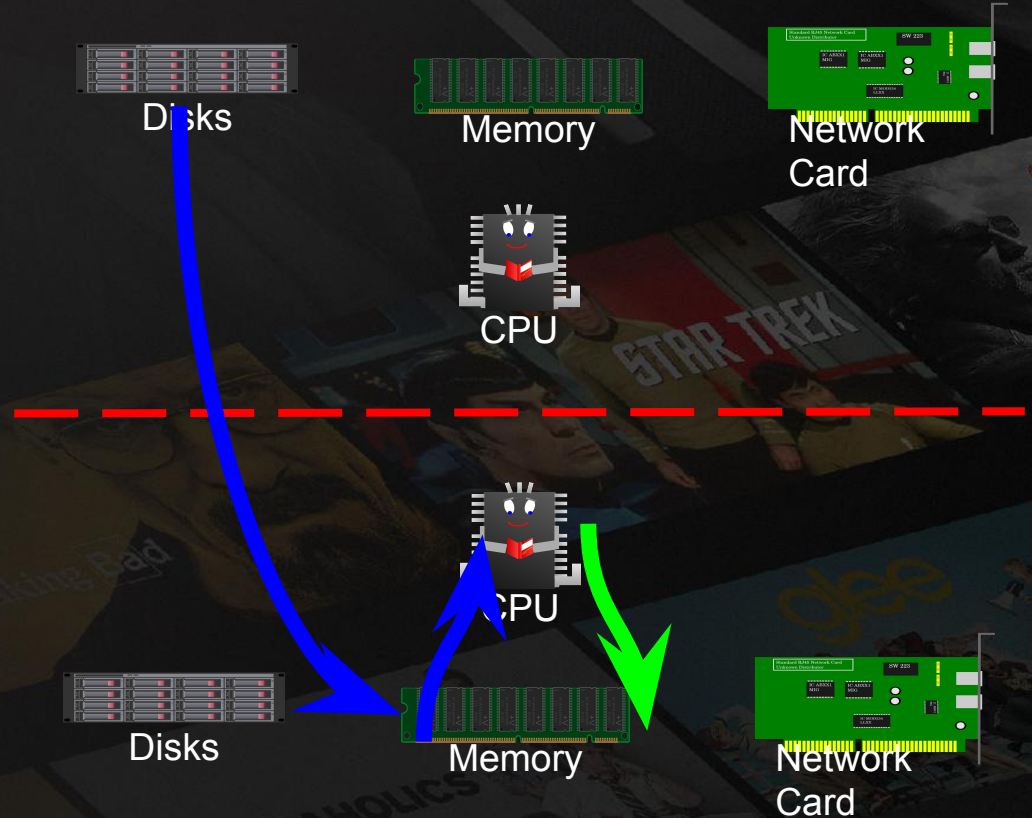
- DMA data from disk to memory
 - First NUMA bus crossing
- CPU Reads data for encryption



NETFLIX Dual AMD: Worst Case Data Flow With Network Centric NUMA Siloing

Steps to send data:

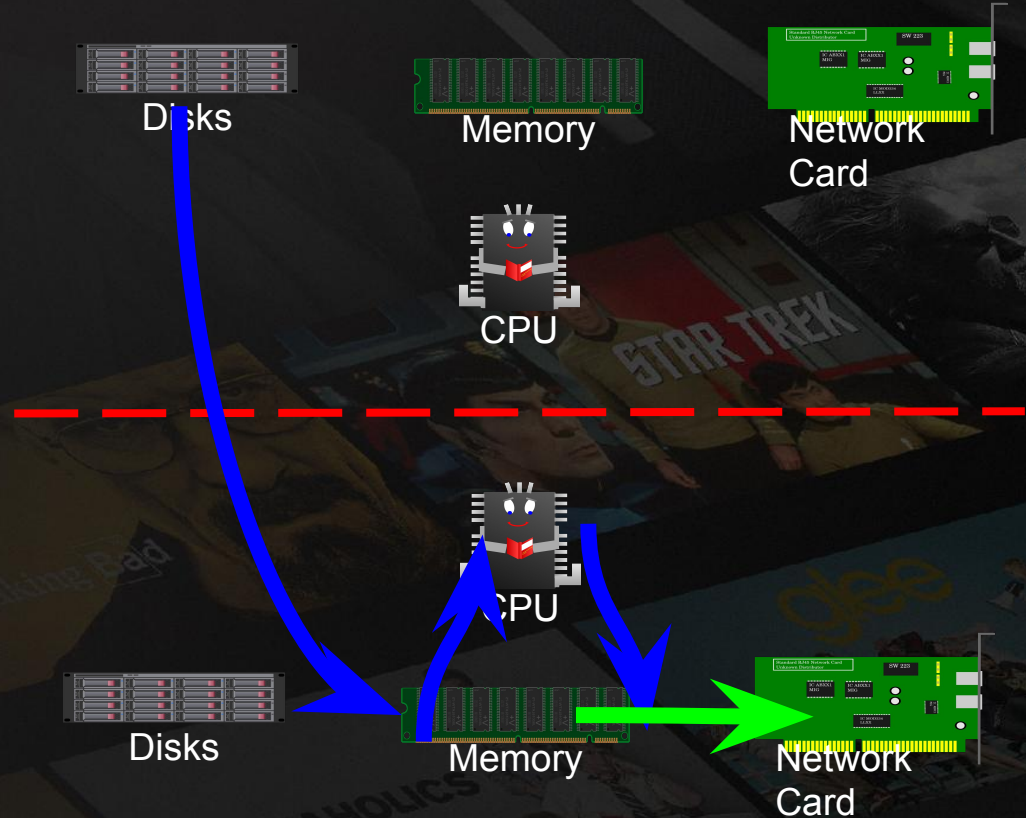
- DMA data from disk to memory
 - First NUMA bus crossing
- CPU Reads data for encryption
- CPU Writes encrypted data



NETFLIX Dual AMD: Worst Case Data Flow With Network Centric NUMA Siloing

Steps to send data:

- DMA data from disk to memory
 - First NUMA bus crossing
- CPU Reads data for encryption
- CPU Writes encrypted data
- **DMA from memory to Network**



Worst Case Summary:

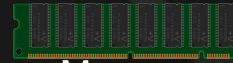
- 1 NUMA crossing on average
 - 100% of disk reads across NUMA
- 100GB/s of data on the NUMA fabric
 - Still less than fabric bandwidth

NETFLIX Dual AMD: Worst Case Data Flow With Disk Centric NUMA Siloing

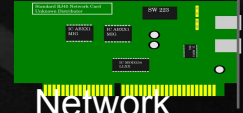
Steps to send data:



Disks



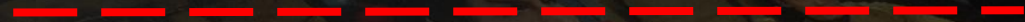
Memory



Network Card



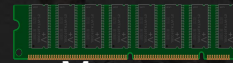
CPU



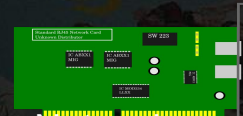
CPU



Disks



Memory

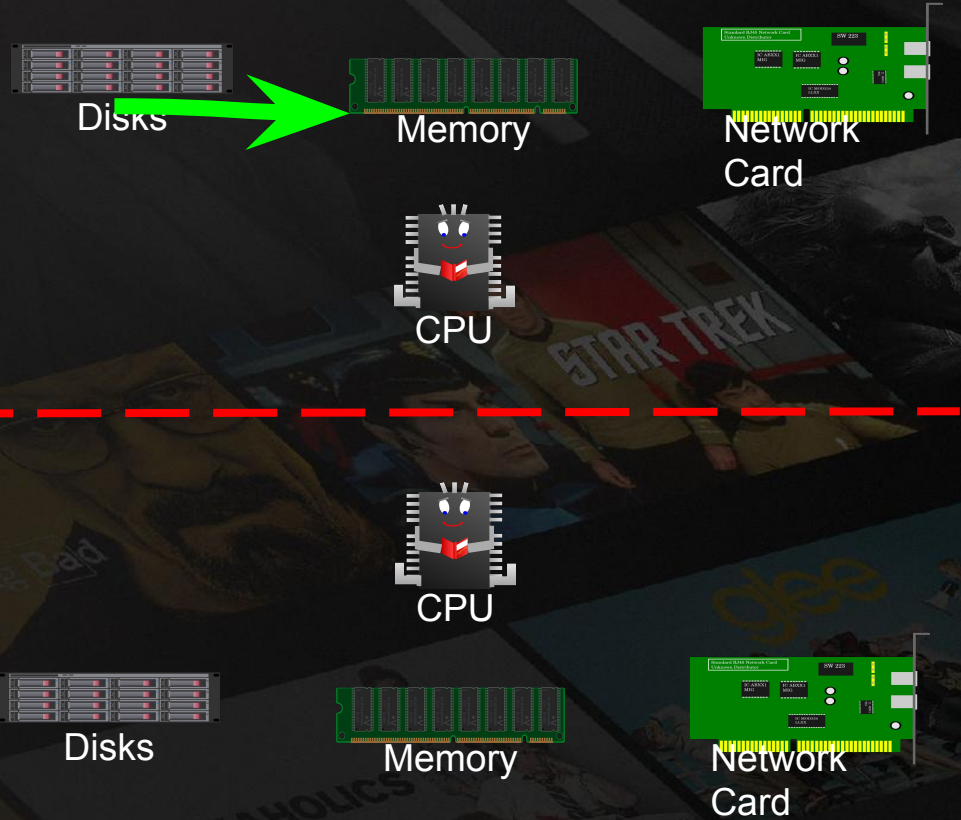


Network Card

NETFLIX Dual AMD: Worst Case Data Flow With Disk Centric NUMA Siloing

Steps to send data:

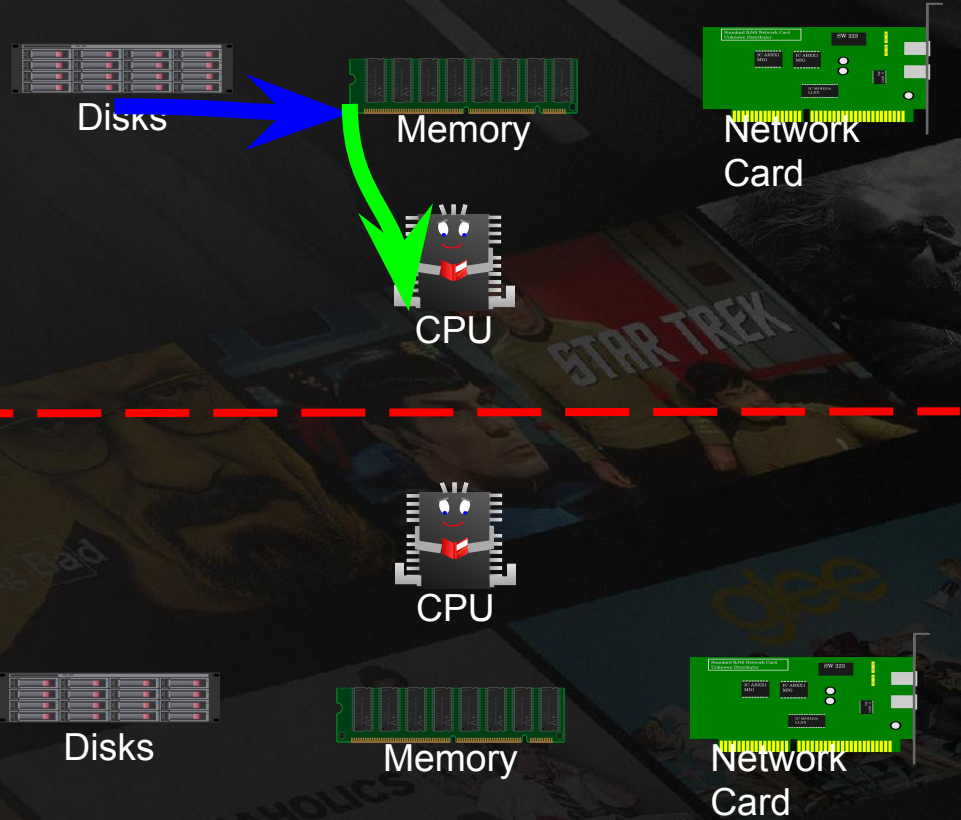
- DMA data from disk to memory



NETFLIX Dual AMD: Worst Case Data Flow With Disk Centric NUMA Siloing

Steps to send data:

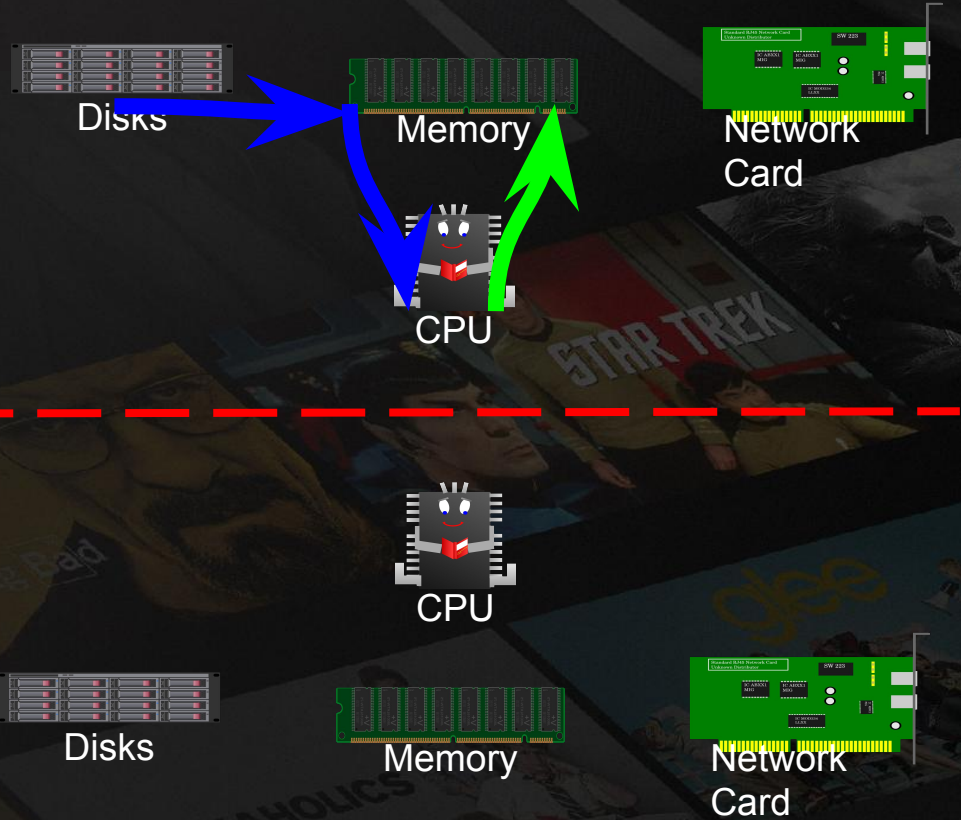
- DMA data from disk to memory
- CPU Reads data for encryption



NETFLIX Dual AMD: Worst Case Data Flow With Disk Centric NUMA Siloing

Steps to send data:

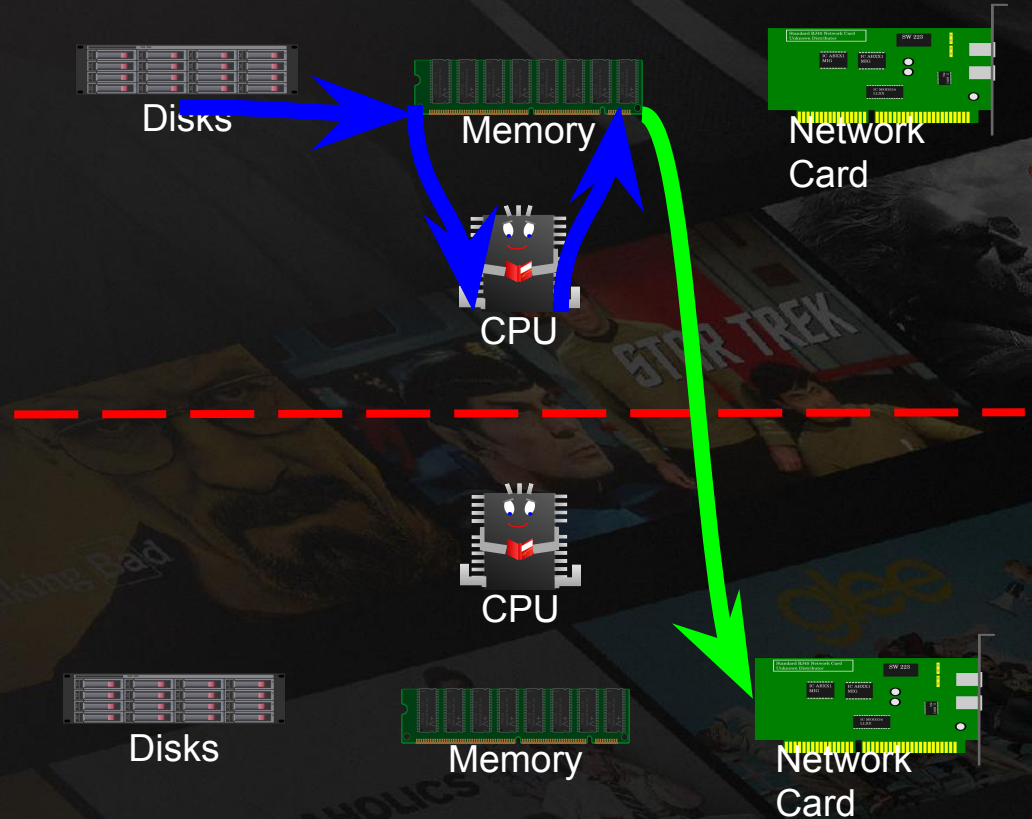
- DMA data from disk to memory
- CPU Reads data for encryption
- CPU Writes encrypted data



NETFLIX Dual AMD: Worst Case Data Flow With Disk Centric NUMA Siloing

Steps to send data:

- DMA data from disk to memory
- CPU Reads data for encryption
- CPU Writes encrypted data
- **NIC Reads data for transmit**
 - **First NUMA bus crossing**



Worst Case Summary:

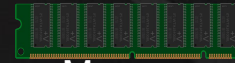
- 1 NUMA crossing on average
 - 100% of disk reads across NUMA
- 100GB/s of data on the NUMA fabric
 - Less than theoretical fabric bandwidth

NETFLIX Dual AMD: Worst Case Data Flow With Strict Disk Centric NUMA Siloing

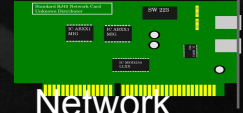
Steps to send data:



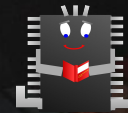
Disks



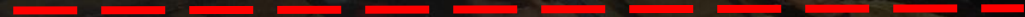
Memory



Network Card



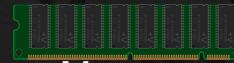
CPU



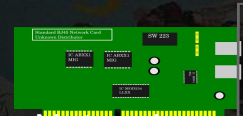
CPU



Disks



Memory

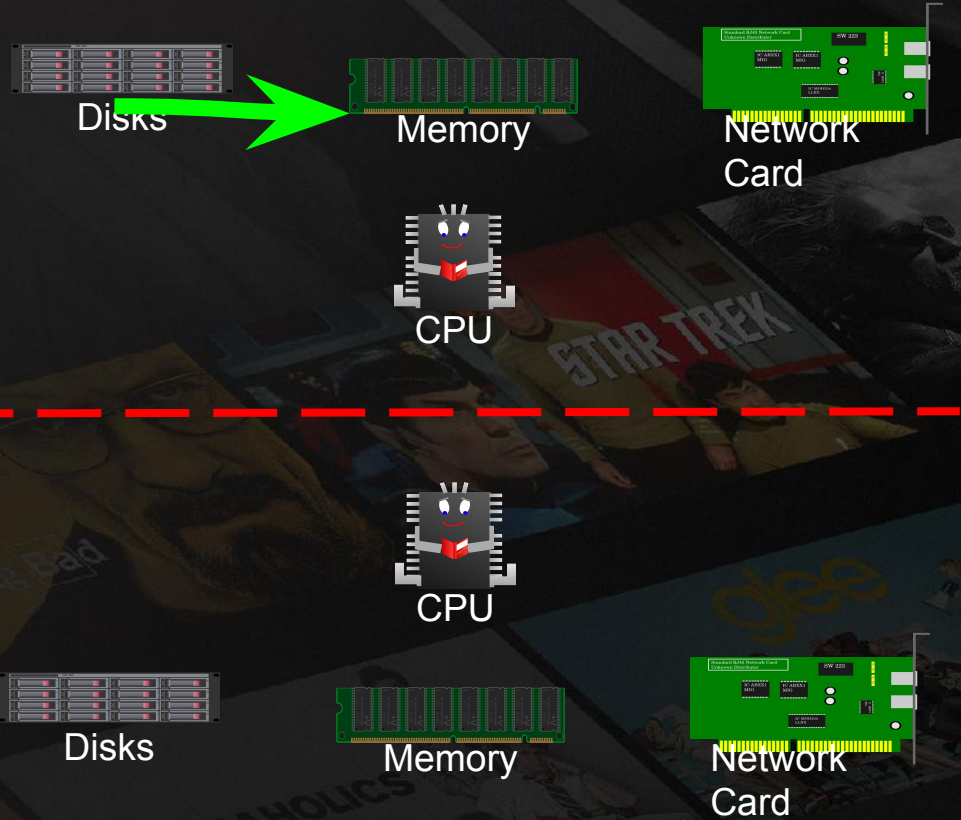


Network Card

NETFLIX Dual AMD: Worst Case Data Flow With Strict Disk Centric NUMA Siloing

Steps to send data:

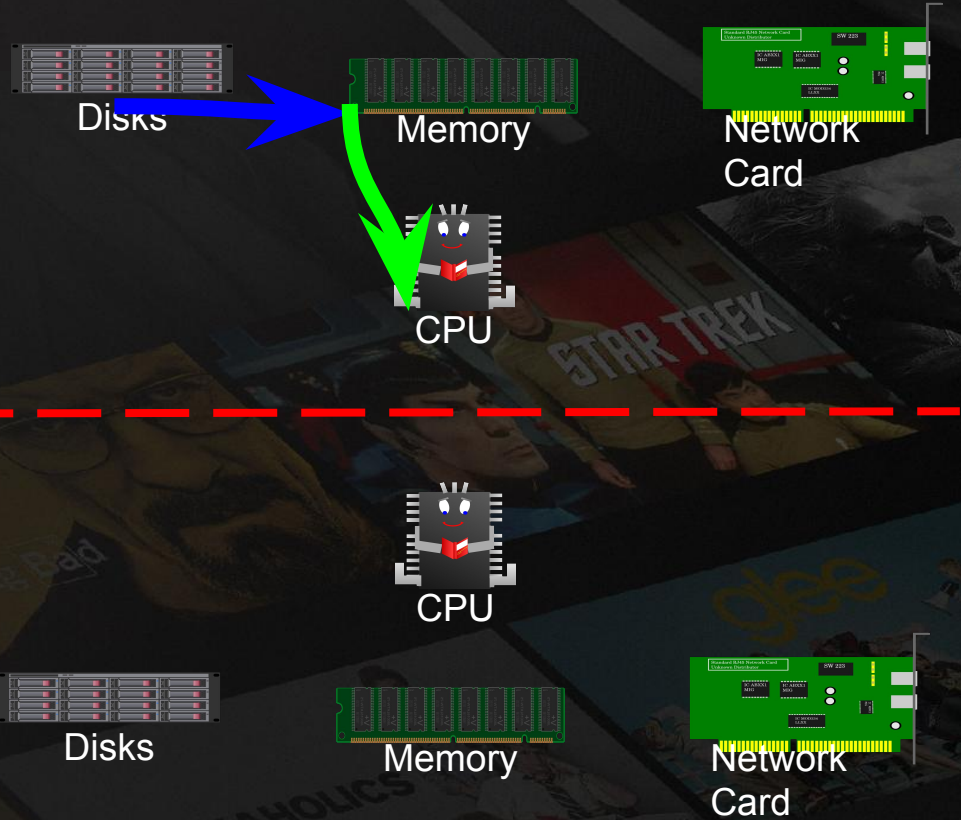
- DMA data from disk to memory



NETFLIX Dual AMD: Worst Case Data Flow With Strict Disk Centric NUMA Siloing

Steps to send data:

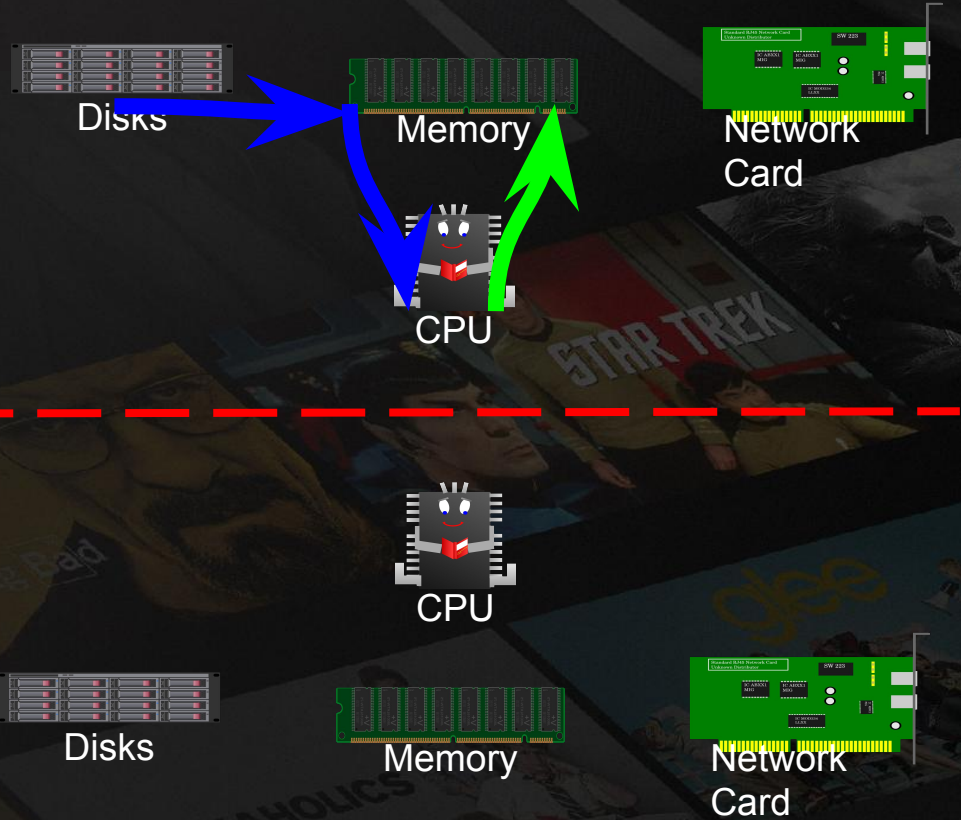
- DMA data from disk to memory
- CPU Reads data for encryption



NETFLIX Dual AMD: Worst Case Data Flow With Strict Disk Centric NUMA Siloing

Steps to send data:

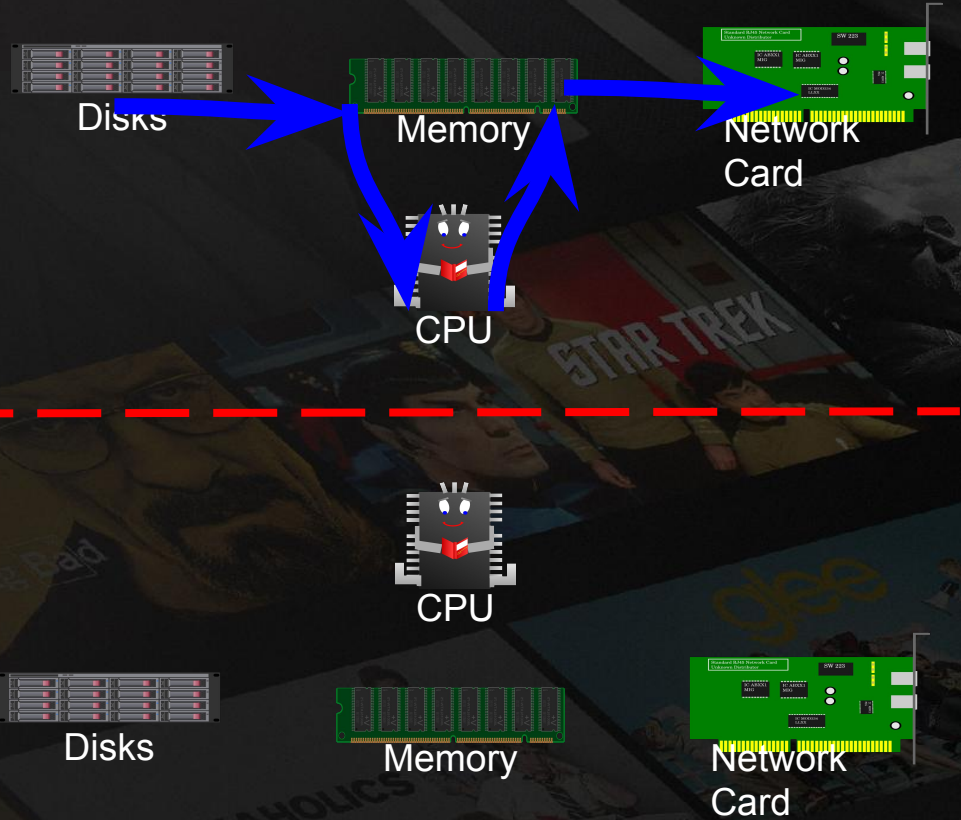
- DMA data from disk to memory
- CPU Reads data for encryption
- CPU Writes encrypted data



NETFLIX Dual AMD: Worst Case Data Flow With Strict Disk Centric NUMA Siloing

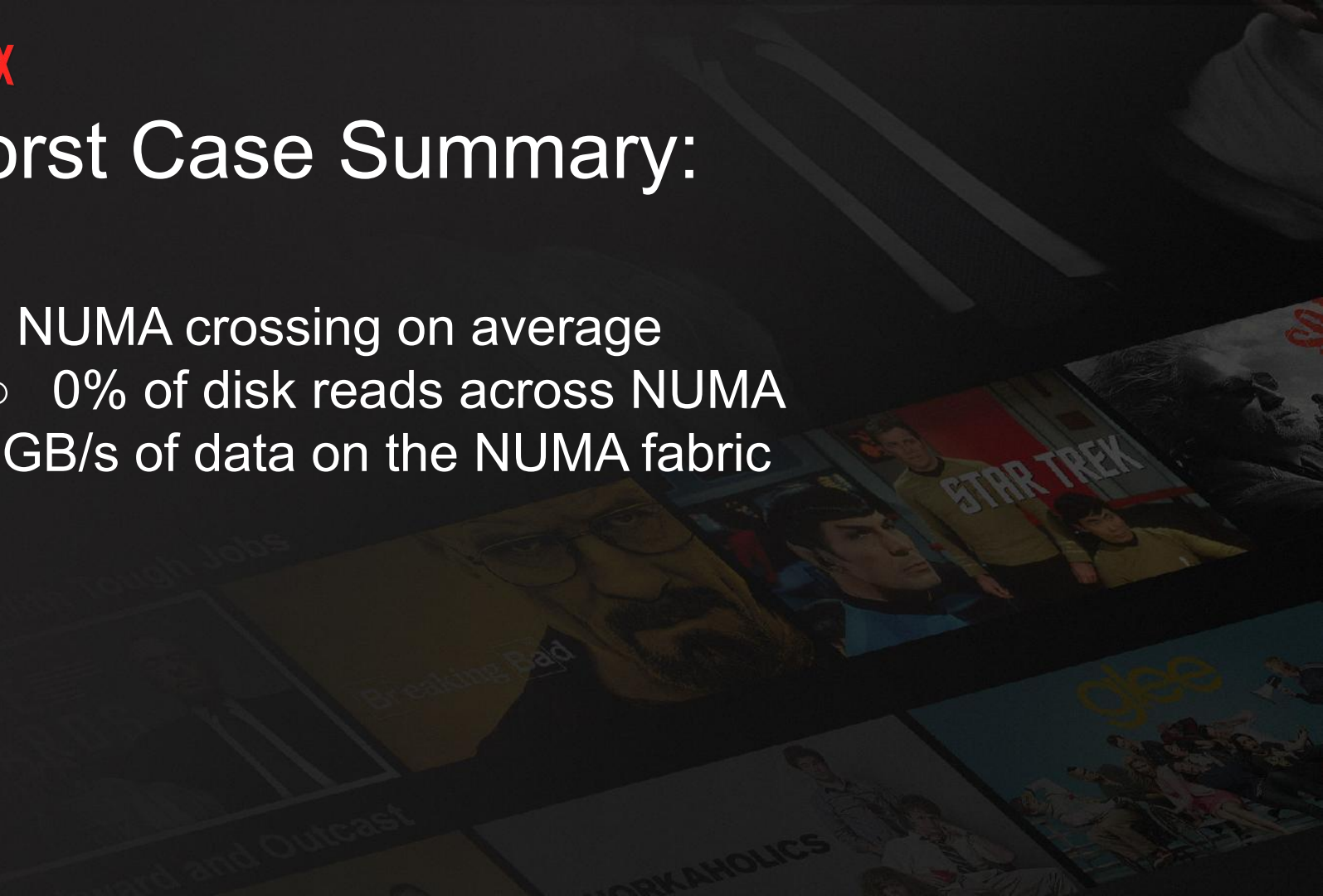
Steps to send data:

- DMA data from disk to memory
- CPU Reads data for encryption
- CPU Writes encrypted data
- **NIC Reads data for transmit**



Worst Case Summary:

- 0 NUMA crossing on average
 - 0% of disk reads across NUMA
- 0GB/s of data on the NUMA fabric

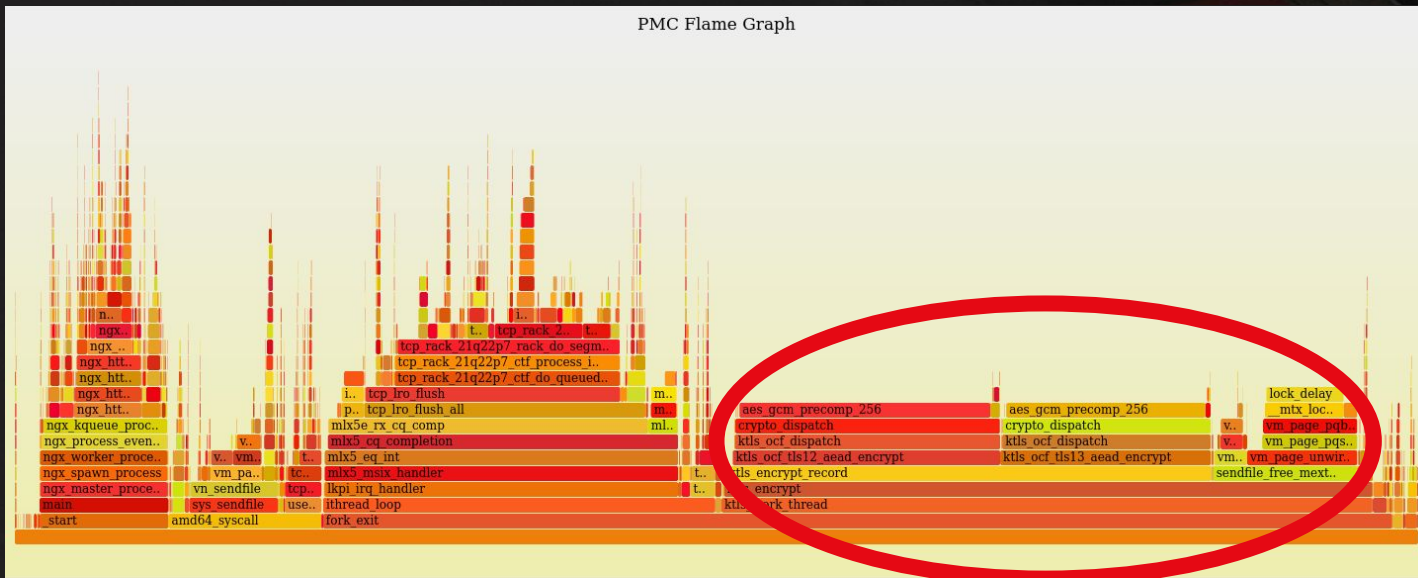


Timeline:

- Asynchronous Sendfile (2014)
- Kernel TLS (2016)
- NUMA (2019)
- ***Inline Hardware (NIC)***
kTLS (2022)
- 800G initial results

Why offload kTLS?

- kTLS uses almost half of our CPU cycles



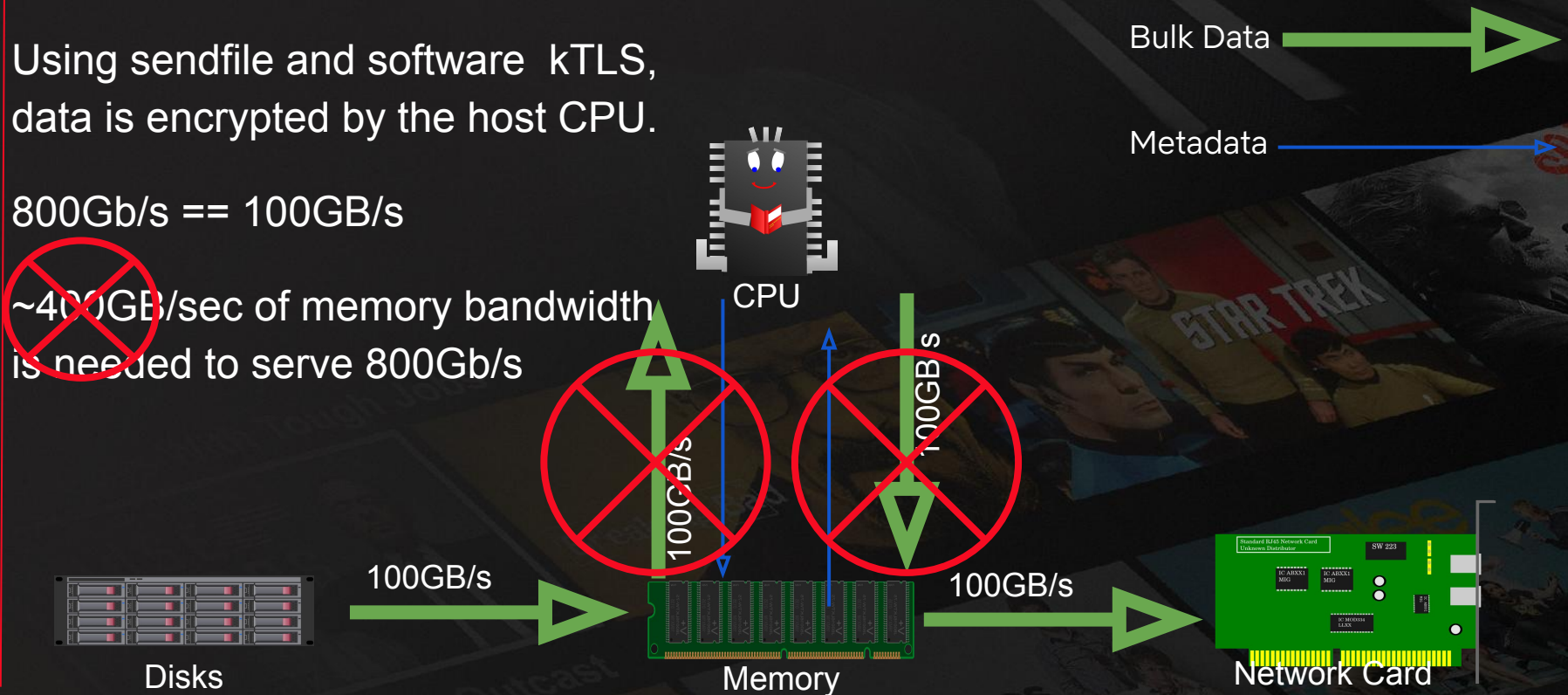
NETFLIX

Netflix 800Gb/s Video Serving Data Flow

Using sendfile and software kTLS, data is encrypted by the host CPU.

800Gb/s == 100GB/s

~~~400GB/sec of memory bandwidth is needed to serve 800Gb/s~~



# Mellanox (NVIDIA) NIC kTLS



- Discussed with Mellanox starting in 2016
- First prototypes of CX6-DX in early 2020
- Iterated for 2+ years to make it production ready
- kTLS offload enabled in production last quarter

# What is NIC kTLS?:

- Hardware Inline TLS
- TLS session is established in userspace.
- When crypto is moved to the kernel, the kernel passes crypto keys to the NIC
- TLS records are encrypted by NIC as the data flows through it on transmit
  - No more detour through the CPU for crypto
  - This cuts memory BW & CPU use in half!

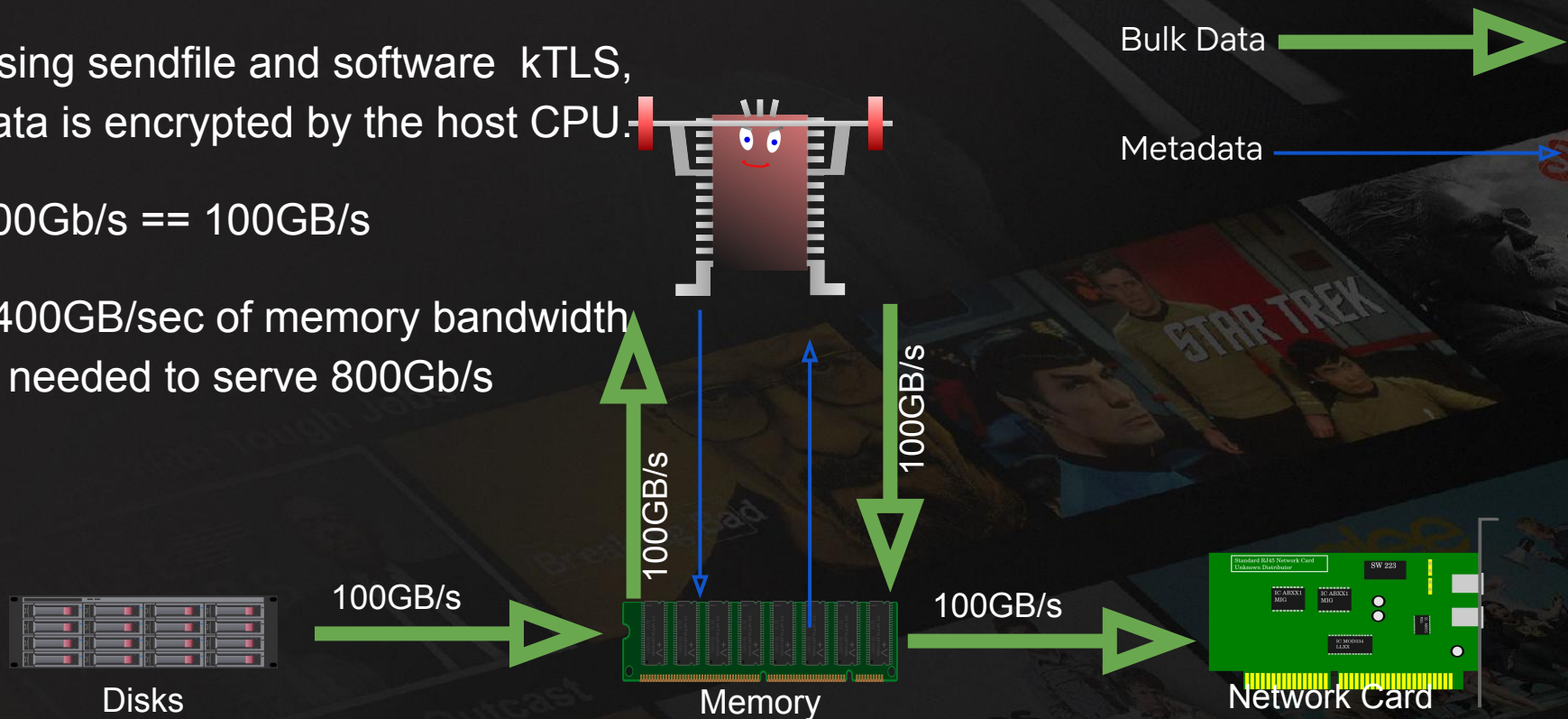
NETFLIX

# Netflix 800Gb/s Video Serving Data Flow

Using sendfile and software kTLS, data is encrypted by the host CPU.

800Gb/s == 100GB/s

~400GB/sec of memory bandwidth is needed to serve 800Gb/s





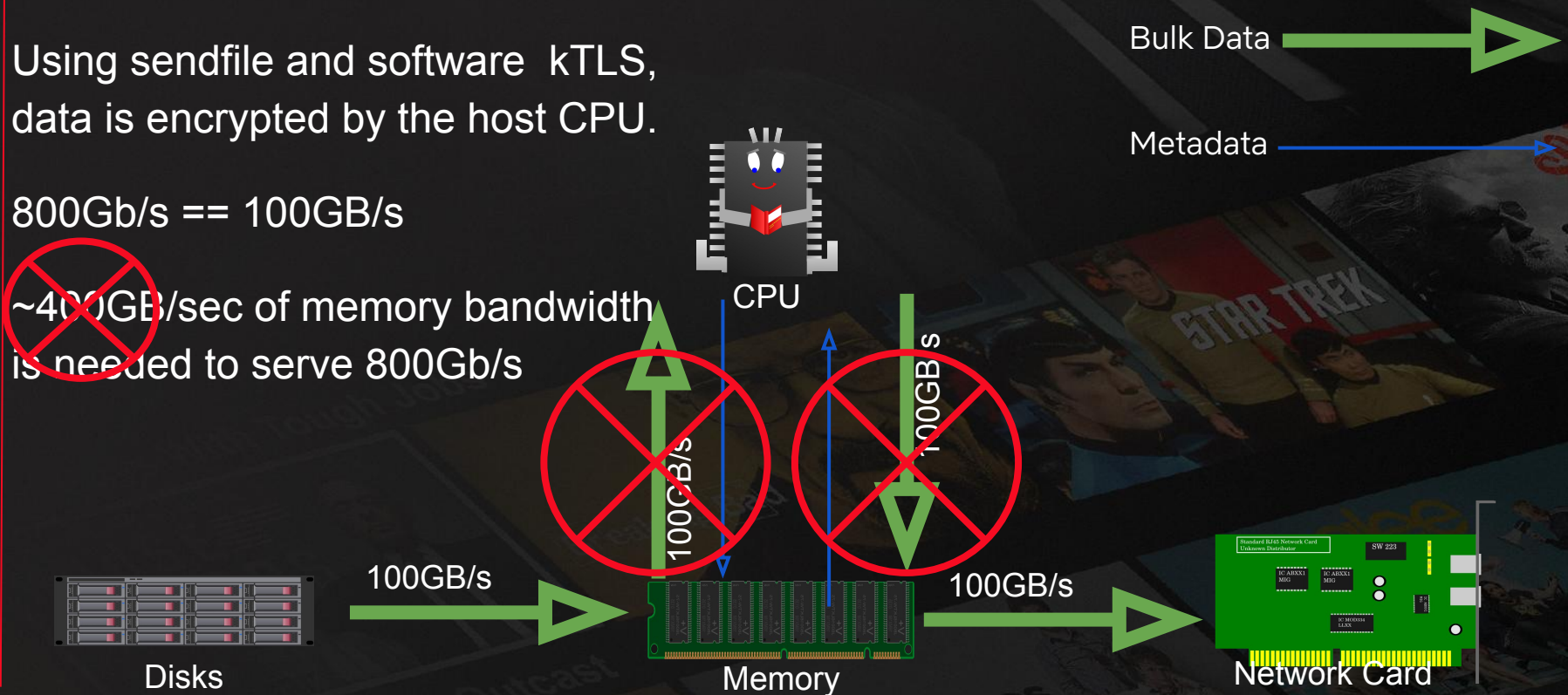
# NETFLIX

## Netflix 800Gb/s Video Serving Data Flow

Using sendfile and software kTLS, data is encrypted by the host CPU.

800Gb/s == 100GB/s

~~~400GB/sec of memory bandwidth is needed to serve 800Gb/s~~



NETFLIX

Netflix 800Gb/s Video Serving Data Flow

Using sendfile and NIC kTLS, data is encrypted by the NIC.

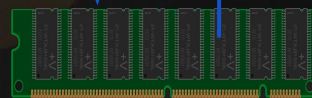
800Gb/s == 100GB/s

~200GB/sec of memory bandwidth is needed to serve 800Gb/s

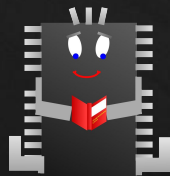


Disks

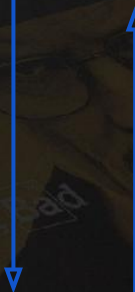
100GB/s



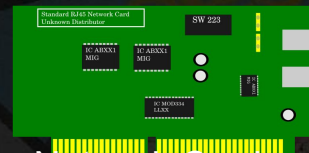
Memory



CPU



100GB/s

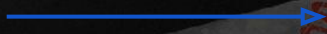


Network Card

Bulk Data



Metadata

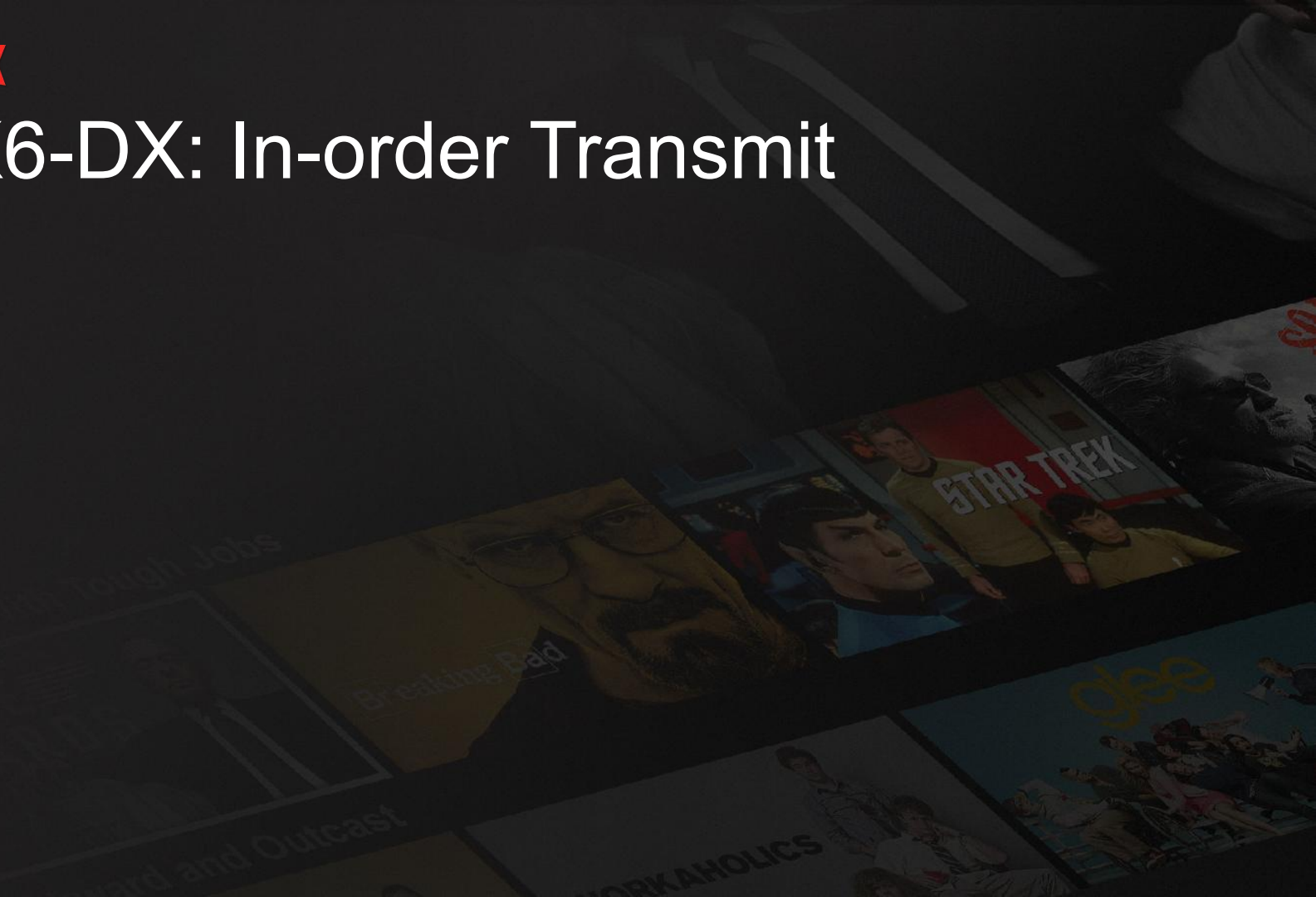


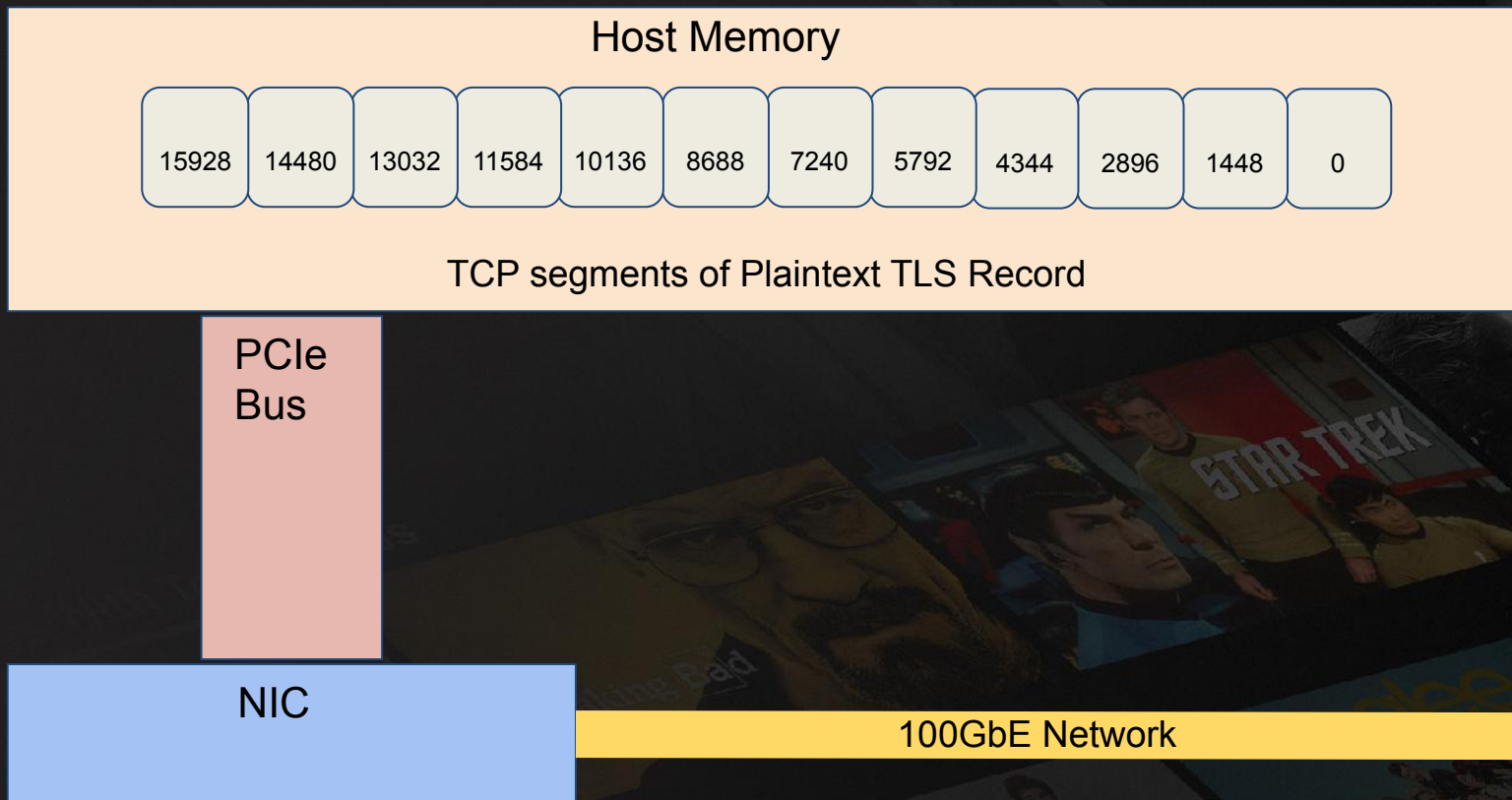
Mellanox ConnectX-6 Dx

- Offloads TLS 1.2 and 1.3 for AES GCM cipher
- Retains crypto state within a TLS record
 - Means that the TCP stack can send partial TLS records without performance loss
- If a packet is sent out of order (eg, a TCP retransmit), it must re-DMA the record containing the out of order packet

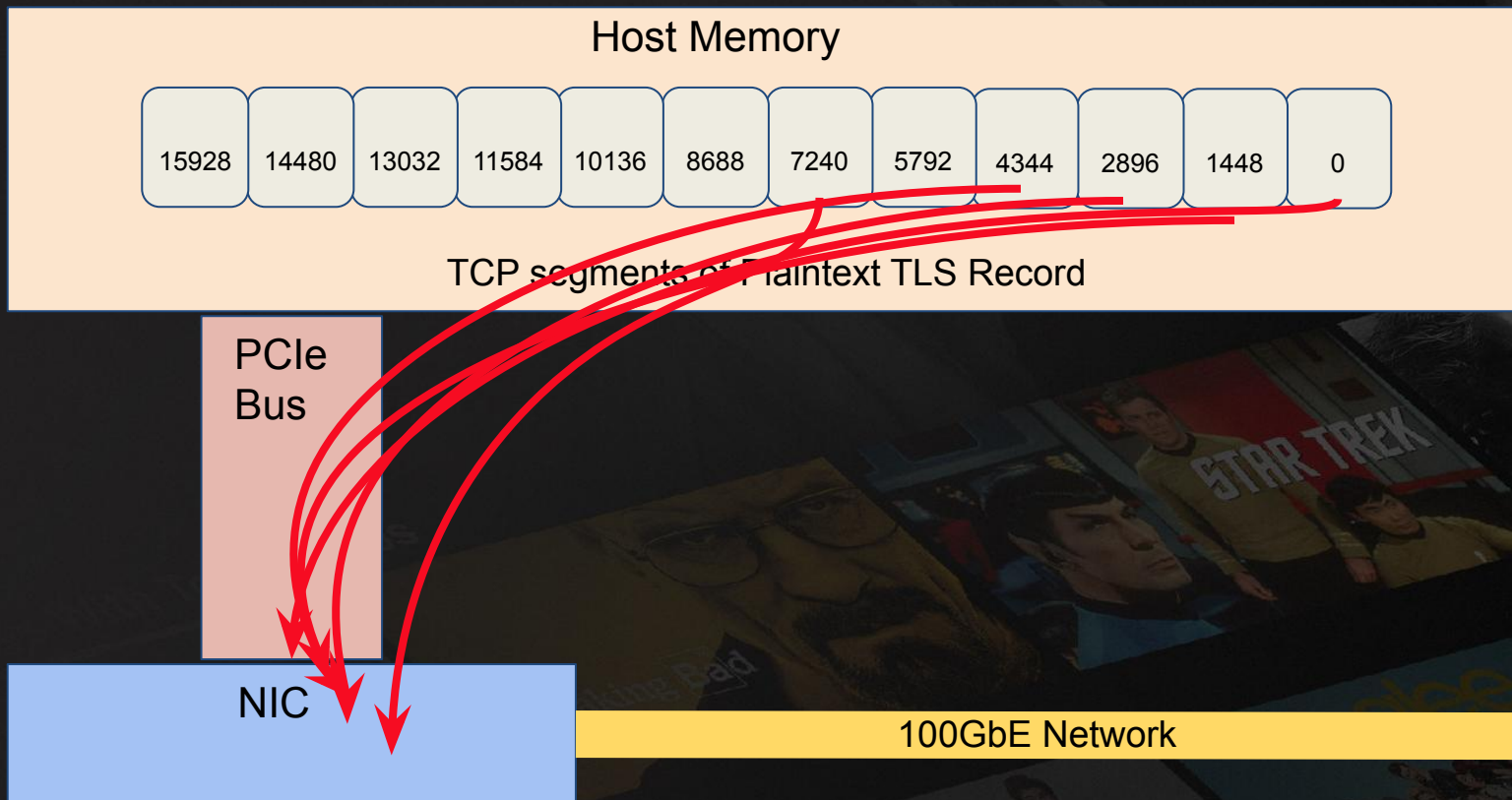
NETFLIX

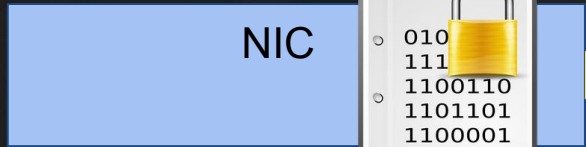
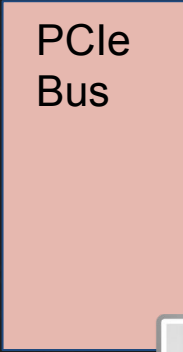
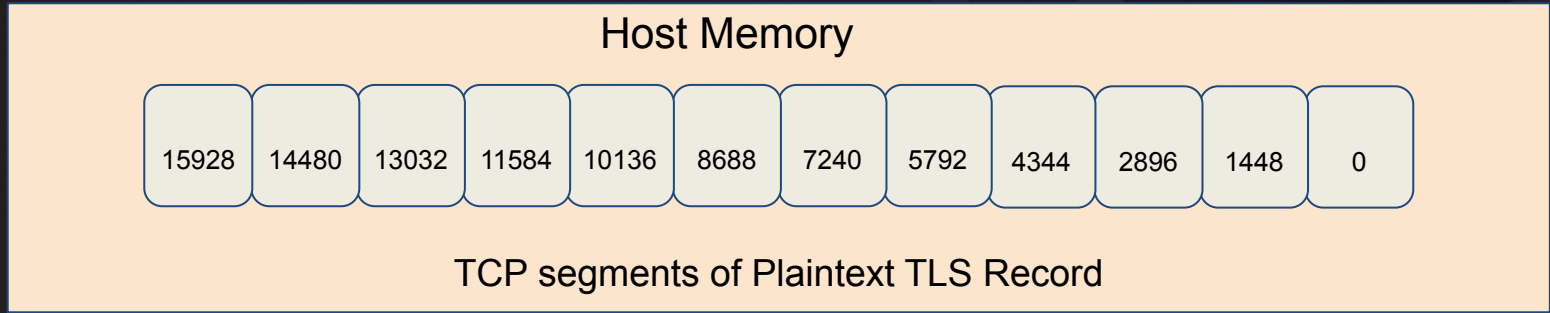
CX6-DX: In-order Transmit



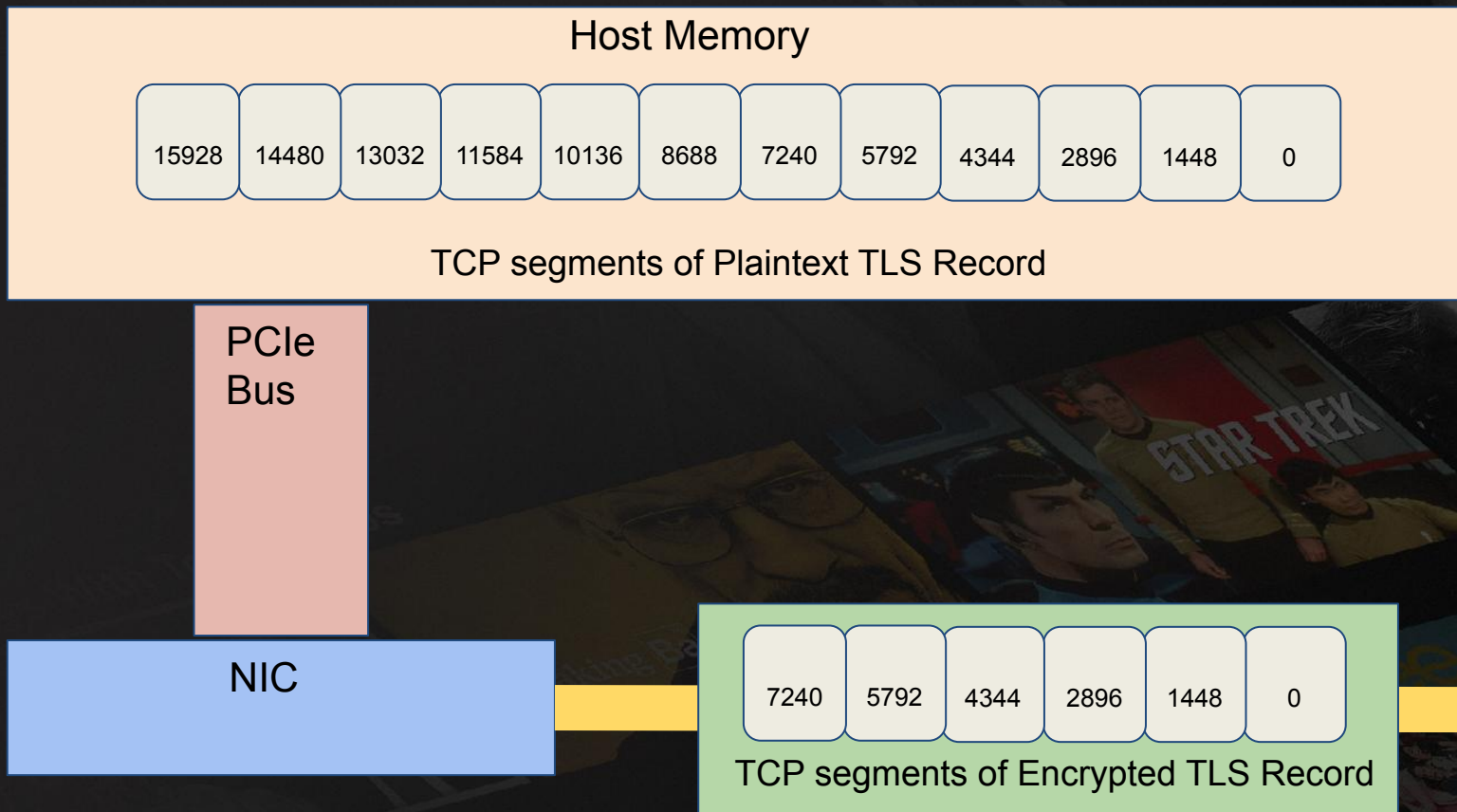


NETFLIX

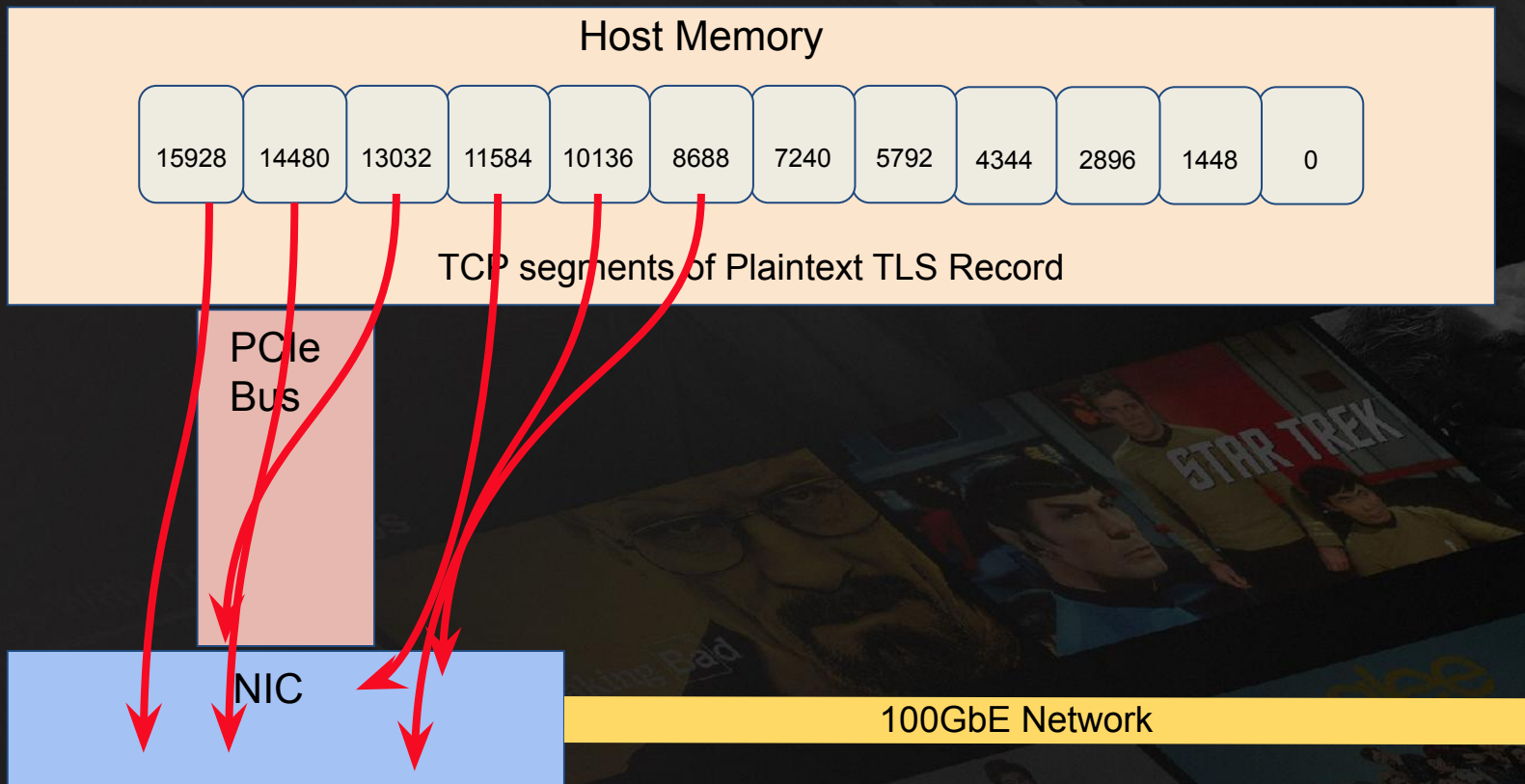


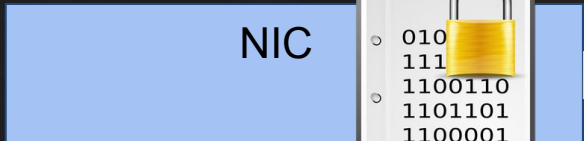
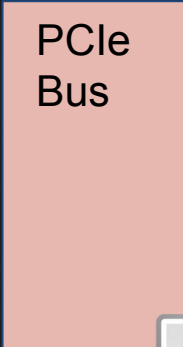
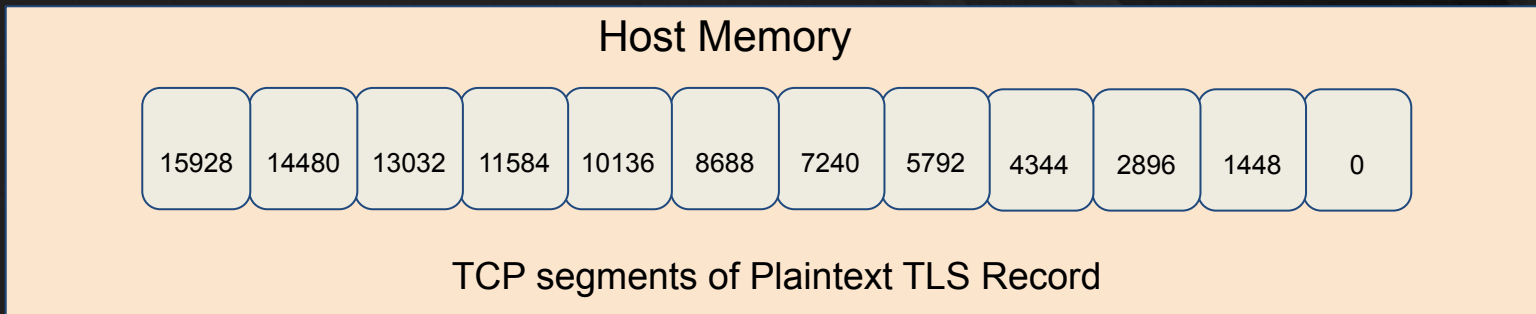


TCP segments of Encrypted TLS Record

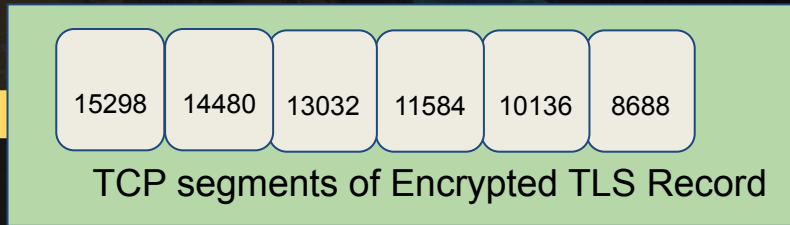
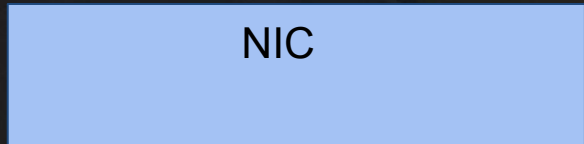
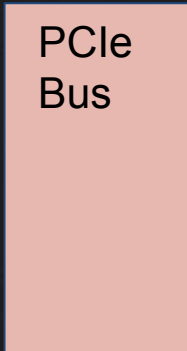
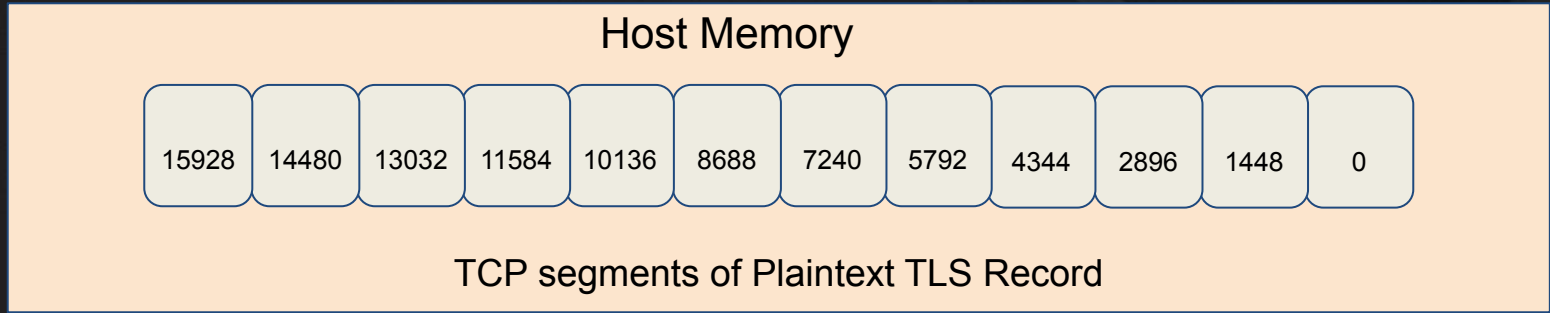


NETFLIX





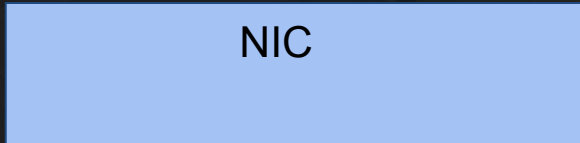
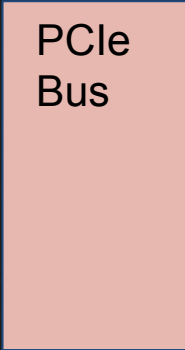
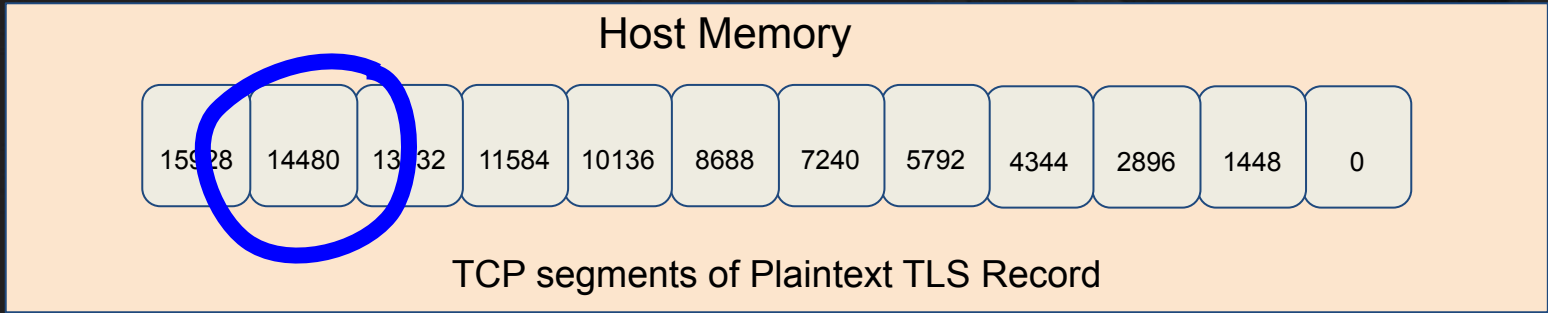
TCP segments of Encrypted TLS Record



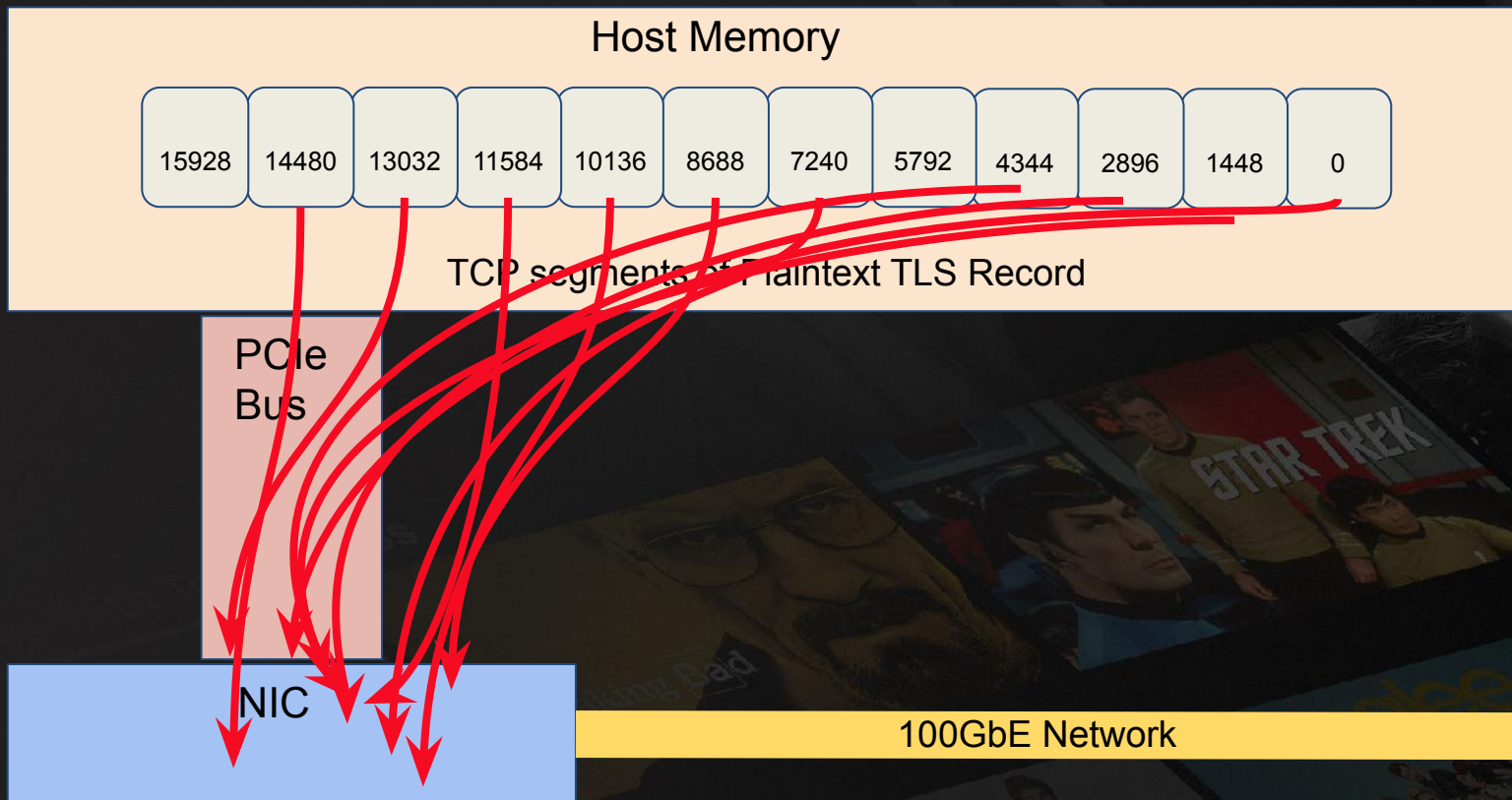
NETFLIX

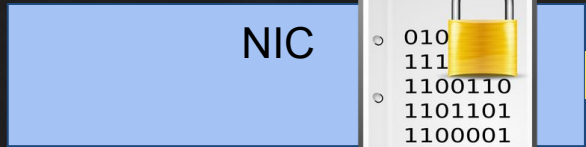
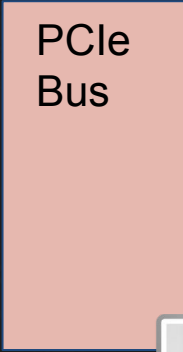
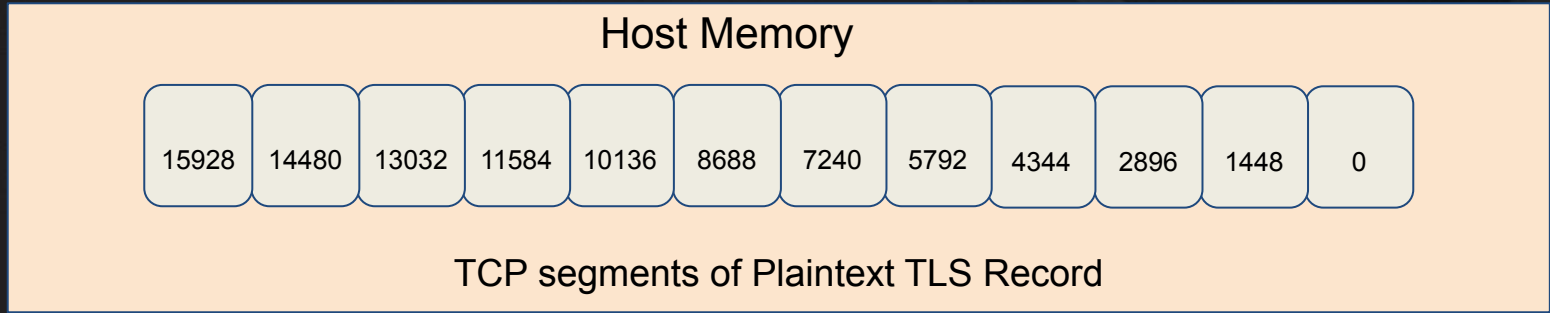
CX6-DX: TCP Retransmit





NETFLIX

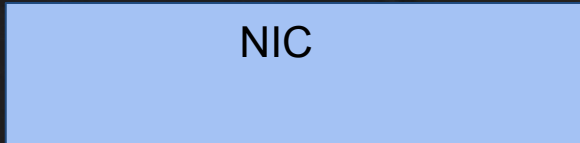
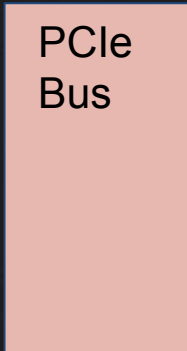
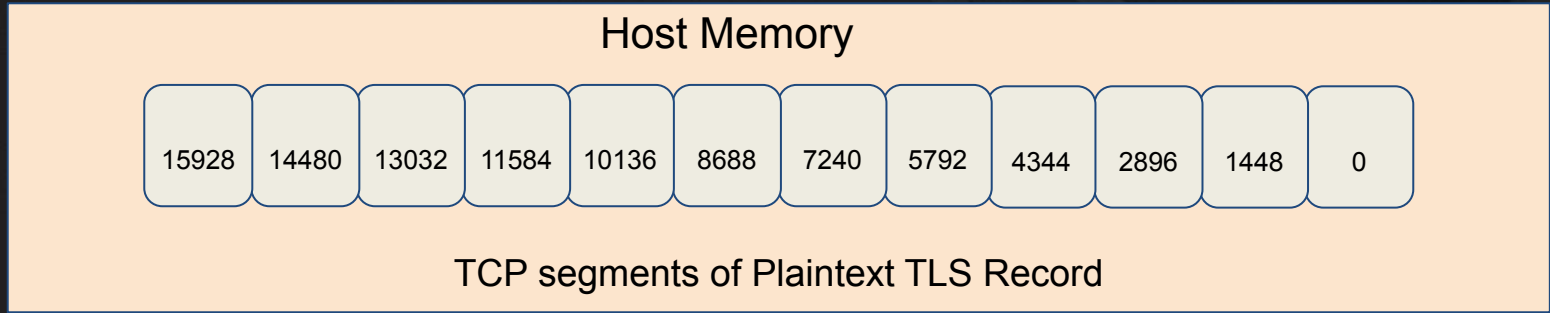




- o 010
- o 111
- o 1100110
- o 1101101
- o 1100001



TCP segments of Encrypted TLS Record



Timeline:

- Asynchronous Sendfile (2014)
- Kernel TLS (2016)
- NUMA (2019)
- Inline Hardware (NIC) kTLS (2022)
- ***800G initial results***

800G Prototype Details

- Dell R7525
- 2x AMD EPYC 7713 64c / 128t (128c / 256t total)
- 3x xGMI links between sockets
- 512 GB RAM
- 4x Mellanox ConnectX-6 Dx (8x 100GbE ports)
- 16x Intel Gen4 x4 14TB NVME

Initial Results: 420Gb/s

- Ran in 1NPS mode
- Network Siloing mode
- CPUs mostly idle
 - AMD guessed that xGMI was down-linking to x2
 - Set xGMI speed to 18GT/s and forced link width to x16, and disabled dynamic link width management

Results with DLWM forced: 500Gb/s

- Ran in 1NPS mode
- Network Siloing mode
 - NVME data DMA'ed to NIC's NUMA Node
- xGMI link usage very uneven:
 - 15GB/s, 4GB/s and 2GB/s
 - Turns out that NVME is not evenly distributed by IO Quadrants
 - Even hashing of cross-socket to xGMI depends on evenly distributed IO

How to Improve xGMI Hashing

- Hashing based on device doing DMA
 - NVME is very uneven
 - NICs are much less uneven
 - “Network Siloing” normally does DMA from NVME to remote node, local to NIC
- Flip things, and do DMA from NVME to local buffers
 - “Disk centric siloing”
- The NICs are now doing DMA across xGMI

Results with Disk Centric Siloing: 670Gb/s

- Much more even xGMI hashing:
 - 10/10/7 GB/s
- Problematic because:
 - Daemon that “locks” content into memory is not NUMA aware & can lead to page daemon thrashing.
 - Still pressure on xGMI links

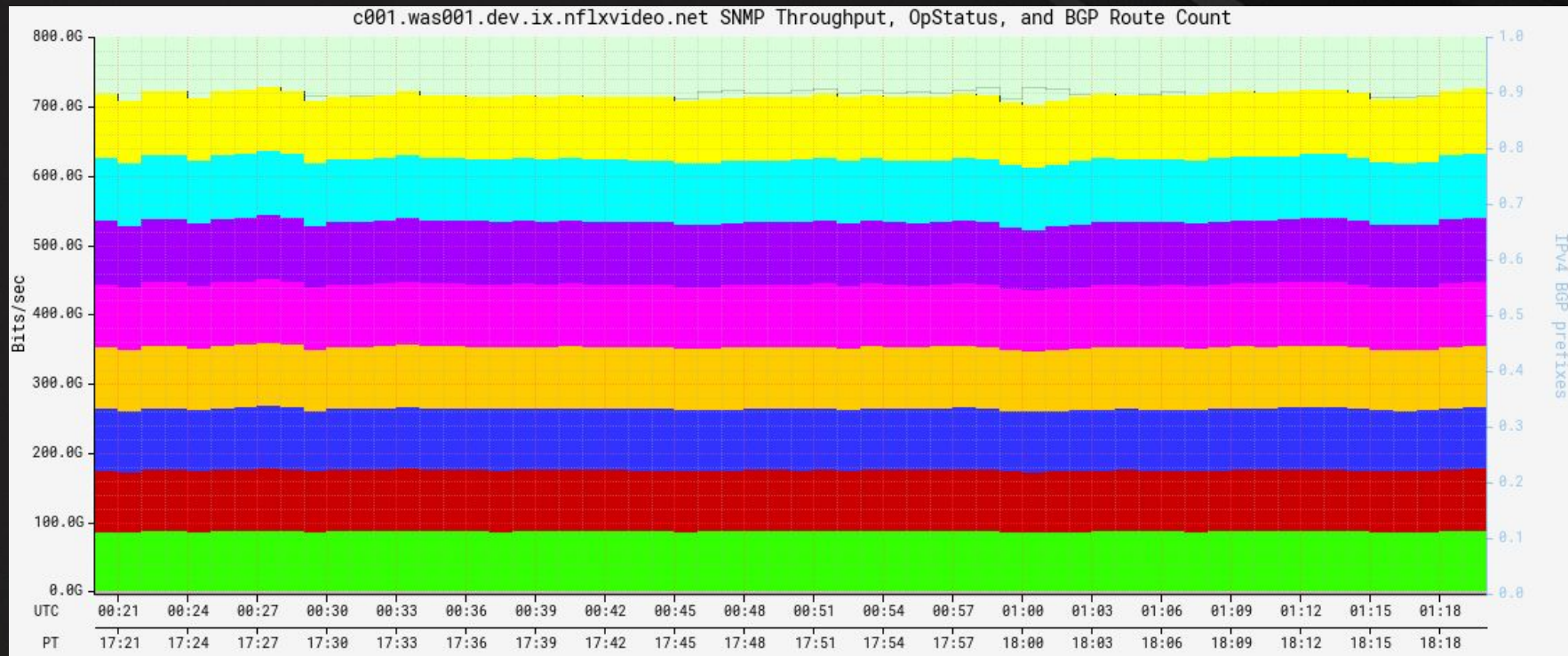
Strict Disk Centric Siloing

- Move Egress NIC to be local to NUMA node with disk
 - No bulk data crosses NUMA Bus
- Incoming TCP traffic still uses original NIC
 - Metadata crosses NUMA bus

“Strict Disk Centric Siloing” Results: 720Gb/s

- Much less xGMI traffic
- Limited by NIC output drops, not CPU.
- Cause of drops is now largely due to:
 - Page daemon interfering with nginx on popular node
 - Uneven loading on NICs due to content popularity differences. (NICs on popular node doing 94Gb/s, others doing 84Gb/s)

NETFLIX



Credits

- Async Sendfile

- Gleb Smirnoff, Konstantin Belousov, Igor Sysoev, Jeff Roberson, Scott Long

- kTLS

- Scott Long, Randall Stewart, Drew Gallatin, John Mark Gurney, John Baldwin

- NUMA

- Drew Gallatin, Jeff Roberson, Mark Johnston

- Inline Hardware (NIC) kTLS

- John Baldwin, Drew Gallatin, Hans Petter Seleaski, Boris Pismenny, Navdeep Parhar

Credits

- Experimental 800GbE Host
 - Warren Harrop and the Netflix hardware team
 - MBX (integrator)
 - AMD (EPYC 7713 CPUs)
 - Dell (PowerEdge R7525)
 - Mellanox/NVIDIA (ConnectX-6 Dx NICS)
 - Intel (P5316 NVME)

NETFLIX

Thank you!

