

# CAM I/O Scheduler

M. Warner Losh

Netflix, Inc.

BSDCan 2015



<http://people.freebsd.org/~imp/bsdcan2015/iosched-slides.pdf>

<http://people.freebsd.org/~imp/bsdcan2015/paper.pdf>

# NETFLIX

# Outline

## Overview / Motivation

Graphs

Roadmap

## Background and Context

Netflix OCA

NAND Physics and SSD

FreeBSD I/O Stack

## Netflix I/O Scheduler

## Recent Updates



# Outline

## Overview / Motivation

Graphs

Roadmap

## Background and Context

Netflix OCA

NAND Physics and SSD

FreeBSD I/O Stack

## Netflix I/O Scheduler

## Recent Updates

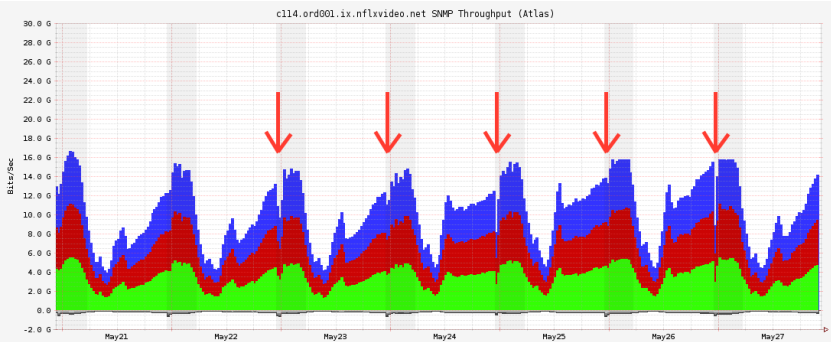


# Motivations

- ▶ Big performance hit in our operations
  - ▶ 20-50% drop in throughput
  - ▶ devstat reporting 100% busy for SSD with tiny use
  - ▶ Netflix's QoE metrics show customer problems
- ▶ Traced to adding / deleting content
  - ▶ Happened during content fill
  - ▶ Work around by idling server before fill
- ▶ Sometimes an issue, other times not
  - ▶ Some systems would have problems, others not
  - ▶ No fixed pattern to failure
  - ▶ Only affected flash caches
  - ▶ Seems related to write activity.

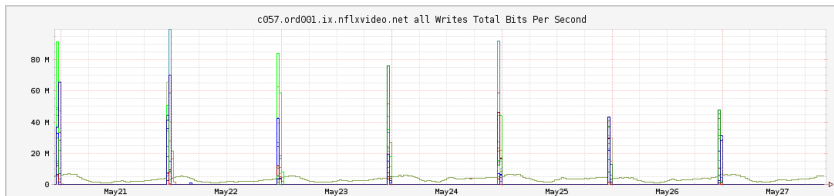


# Macro View

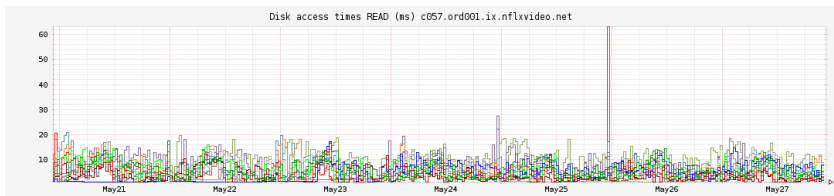


NETFLIX

# Root Cause



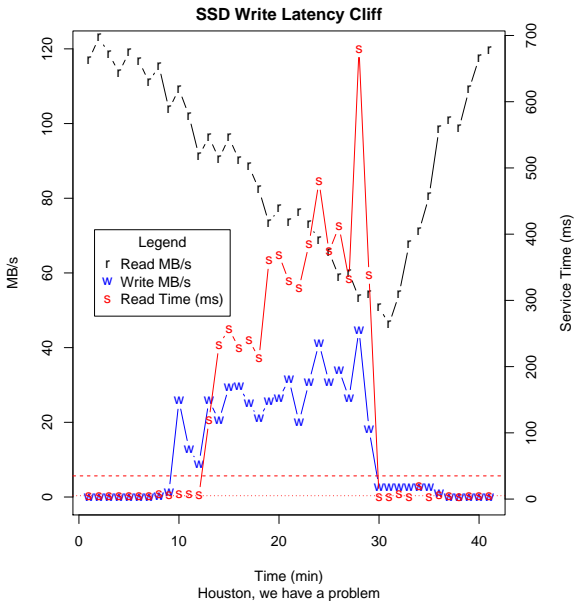
- ▶ High write workload



- ▶ Read latencies can spike



# One Graph View



NETFLIX

- ▶ I/O Scheduler
- ▶ CAM
- ▶ I/O Stack
- ▶ SSDs
- ▶ NAND Physics
- ▶ Netflix Work load





# Outline

## Overview / Motivation

Graphs

Roadmap

## Background and Context

Netflix OCA

NAND Physics and SSD

FreeBSD I/O Stack

## Netflix I/O Scheduler

## Recent Updates



## Overview of Netflix's Open Connect Appliance (OCA)



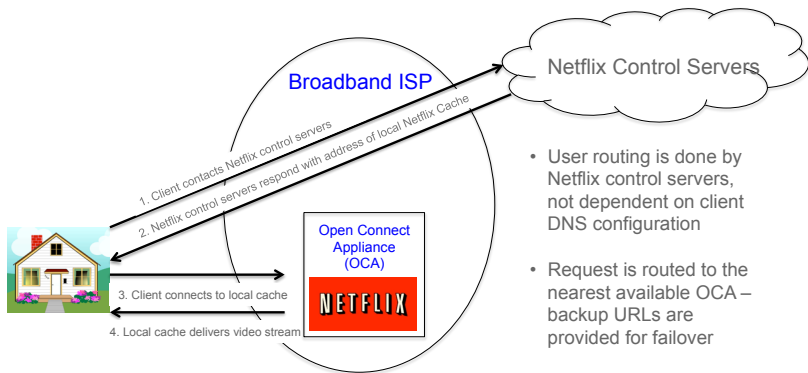
# Netflix Open Connect Appliance (OCA)



Source: Netflix

**NETFLIX**

# Netflix Open Connect Appliance (OCA)



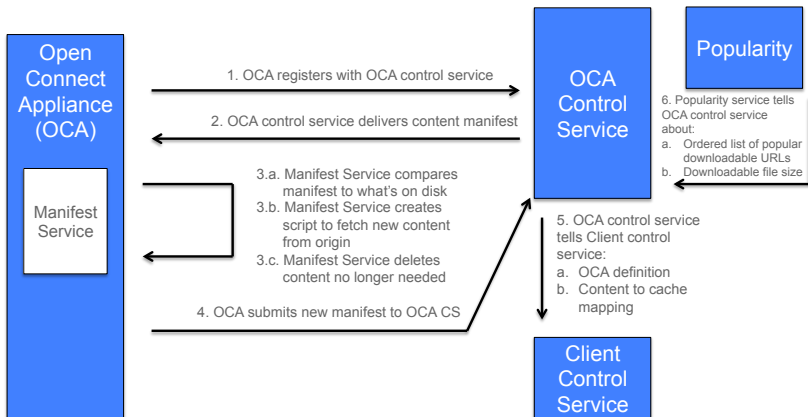
- User routing is done by Netflix control servers, not dependent on client DNS configuration
- Request is routed to the nearest available OCA – backup URLs are provided for failover
- ISP controls client to OCA mapping/clustering/failover via BGP



Source: Netflix

NETFLIX

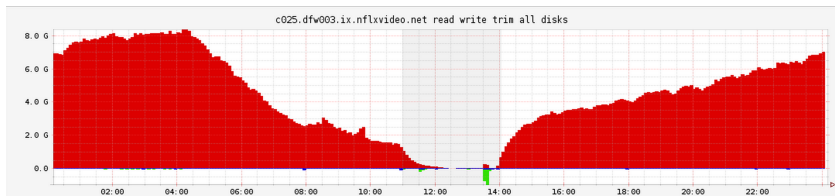
# Netflix Open Connect Appliance (OCA)



Source: Netflix

NETFLIX

# Typical Disk Activity (1 day)



- ▶ Lots of read traffic
- ▶ Little write traffic
- ▶ 100:1 read:write or more ratio
- ▶ “Fill Window” needed to avoid problem



# NAND Flash

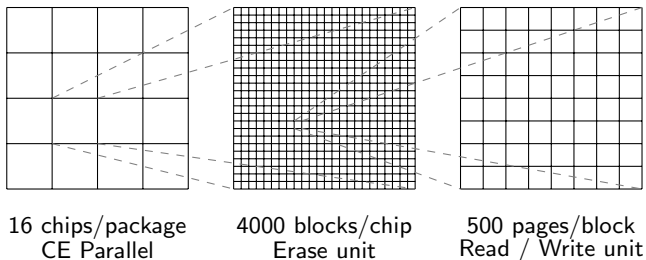
A Brief Primer on Why and How NAND Flash “Presents Challenges” and How SSDs Hide That Through Lies and Deceit



Source: [http://pretty-little-liars.wikia.com/wiki/Season\\_5](http://pretty-little-liars.wikia.com/wiki/Season_5)

**NETFLIX**

# NAND Geometry



NAND typical geometry

Pages are 4kiB-32kiB (typically 16kiB) plus OOB





# NAND Limitations

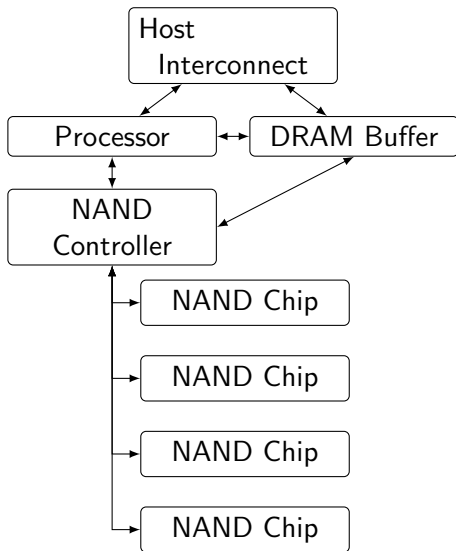
- ▶ Read page at a time
- ▶ Write page at a time
- ▶ Write pages in order
- ▶ Erase entire block
- ▶ Single Duplex, No queuing
- ▶ Low endurance (~3k P/E cycles)
- ▶ Unreliable Media (ECC / LDPC needed)
- ▶ Parallelism through Banks / CE lines
- ▶ NAND Flash and FreeBSD

<https://www.youtube.com/watch?v=lj0XAE6C6-k>

<https://people.freebsd.org/~imp/bsdcan2014.pdf>



# Typical SSD / NVMe



Generalized block diagram of flash storage device.

**NETFLIX**

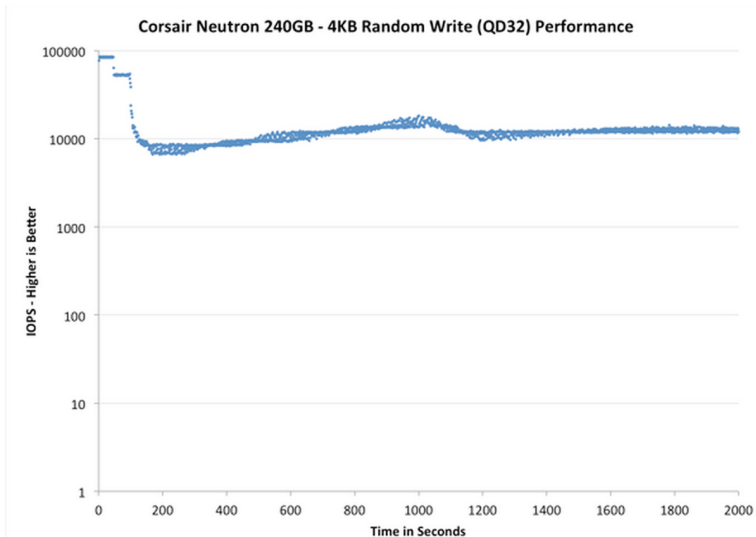


# SSD / NVME Firmware

- ▶ Flash Translation Layer (FTL)
  - ▶ LBA to PA translation
  - ▶ Metadata for log / NAND
- ▶ Wear Leveling
  - ▶ Where to start writing at block boundary
  - ▶ Which block to garbage collect
- ▶ Reliability
  - ▶ Retention (data too old or read too much)
  - ▶ Wear out (block too worn with high RBER)
  - ▶ Program / Erase error processing
- ▶ Garbage Collection
  - ▶ Moves data forward
  - ▶ Extra reads and writes (Write Amplification)
  - ▶ Can affect performance



# Effects of Garbage Collection



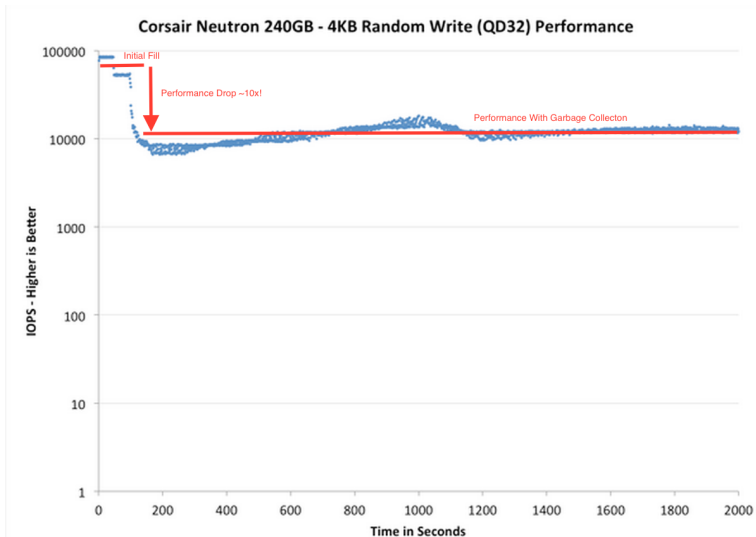
<http://www.howtogeek.com/165542/>

Source

NETFLIX



# Effects of Garbage Collection



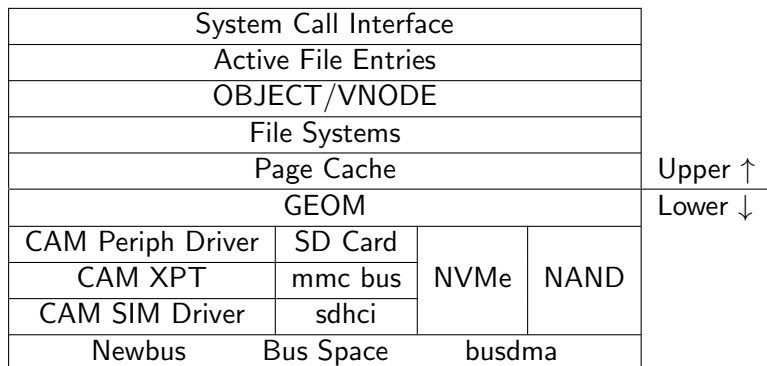
<http://www.howtogeek.com/165542/>

Source

NETFLIX



# FreeBSD I/O Stack



After Figure 7.1 in The Design and Implementation of the FreeBSD Operating System, 2015.

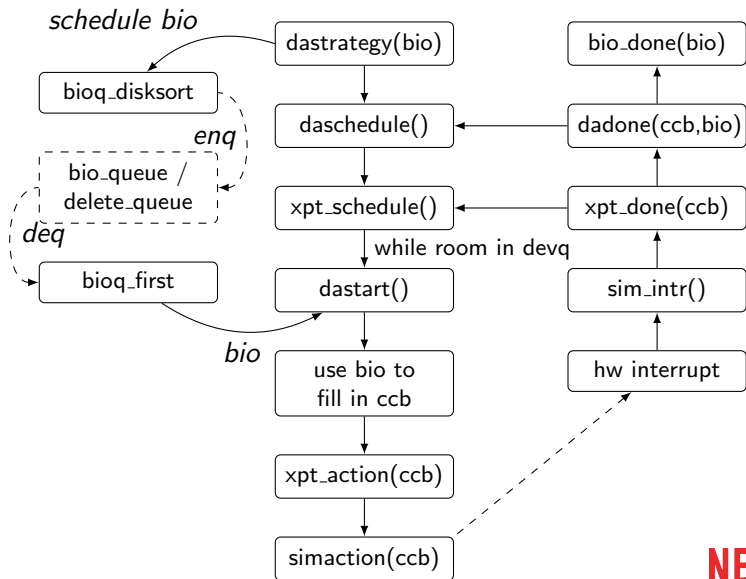


# FreeBSD I/O Stack

- ▶ Upper half of I/O Stack focus of VM system
  - ▶ Buffer cache
  - ▶ Memory mapped files / devices
  - ▶ Loosely coupled user actions to device action
- ▶ GEOM handles partitioning, compression, encryption
  - ▶ Filters data (compression, encryption)
  - ▶ Muxes Many to one (partitioning)
  - ▶ Muxes One to Many (striping / RAID)
  - ▶ Limited Scheduling
- ▶ CAM handles queuing
  - ▶ Shapes flows to device
  - ▶ Limits requests to number of slots
  - ▶ Enforces rules (eg tagged vs non-tagged)
  - ▶ Multiplexes shared resources between devices



# CAM – Data I/O data path through CAM





# Outline

## Overview / Motivation

Graphs

Roadmap

## Background and Context

Netflix OCA

NAND Physics and SSD

FreeBSD I/O Stack

## Netflix I/O Scheduler

## Recent Updates

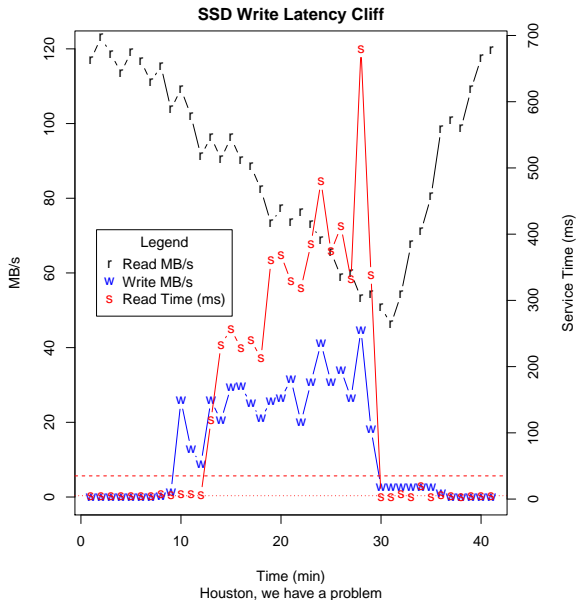


# FreeBSD Default I/O Scheduler

- ▶ No differentiation of I/O (except BIO\_DELETE)
- ▶ Implements two ordering policies
  - ▶ In order
  - ▶ elevator
- ▶ Implemented in the CAM PERIPH drivers
  - ▶ Most flexible
  - ▶ Duplicated code
  - ▶ Partial lie: SIM drivers also involved (NCQ)
- ▶ Generally performs well for well behaved devices
  - ▶ SSDs break the rules
  - ▶ Different SSDs break different rules
  - ▶ Assumed cost symmetry often not the case



# Another Look



# Netflix I/O Scheduler Theory

- ▶ Restricting write rate reduces write amp disturbance
- ▶ Fewer concurrent writes leaves more banks for reads
- ▶ Elevated latency OK within limits
- ▶ Lowest later needed since PERIPH knows about device
  - ▶ GEOM layer too high
  - ▶ GEOM filters requests, but can't force PERIPH



# I/O Scheduler Changes

- ▶ Create abstract interface to scheduler
- ▶ Convert da and ada PERIPHS to new interface
- ▶ Make sure no regressions



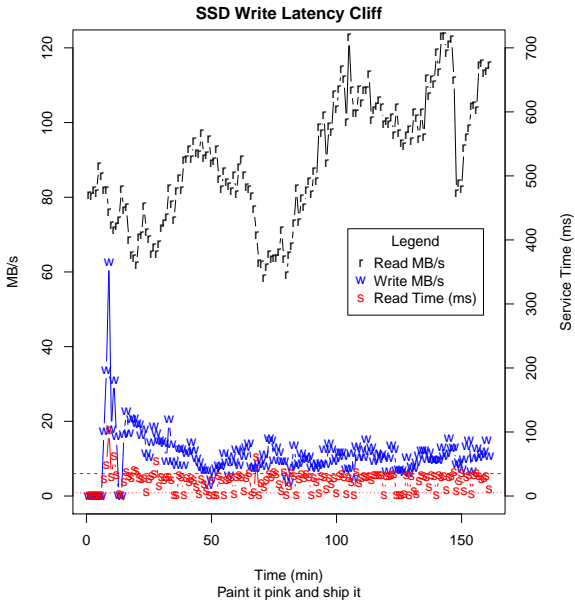
# Netflix I/O Scheduler Changes

- ▶ Separate read, write and delete queues
- ▶ Lots of statistics
- ▶ Ability to limit number of I/O in device at a time
- ▶ Adjustments needed for scheduling





# Results



NETFLIX



# Outline

## Overview / Motivation

Graphs

Roadmap

## Background and Context

Netflix OCA

NAND Physics and SSD

FreeBSD I/O Stack

## Netflix I/O Scheduler

## Recent Updates

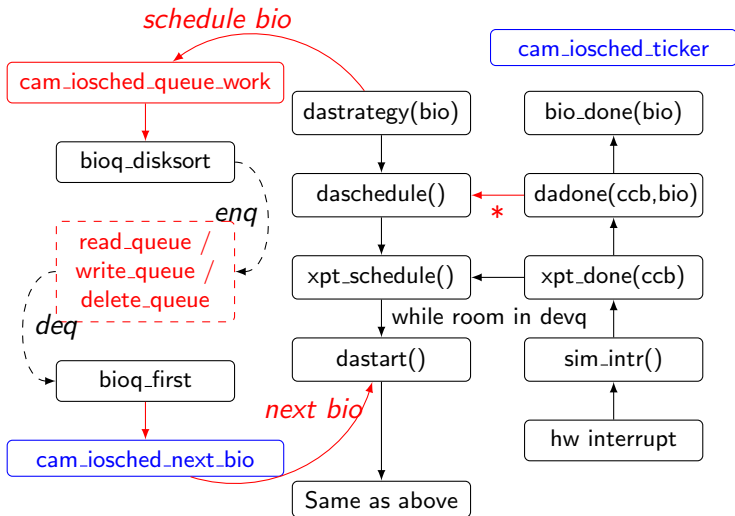


# Added Features

- ▶ Added Bandwidth and IOPS limits
- ▶ Added timeout to implement quanta scheduling
- ▶ Added dynamic steering of limits



# Code Changes



Changes to default scheduler in red. Changes to prior I/O scheduler in blue.

NETFLIX

# Recent issues

- ▶ Large quanta produce large latencies
- ▶ Rate limited drives report 100 percent busy
- ▶ Dynamic loop not yet tuned



# Questions

Questions?

Comments?

Warner Losh

wlosh@netflix.com

imp@FreeBSD.org

<http://people.freebsd.org/~imp/bsdcon2015/iosched-slides.pdf>

<http://people.freebsd.org/~imp/bsdcon2015/paper.pdf>

<http://people.freebsd.org/~imp/asiabsdcon2015/works>



**NETFLIX**