

Beyond Consistent Hashing and TCP: Vastly Scalable Load Balanced Storage Clustering

Alex Aizman, Caitlin Bestler Nexenta Systems

The Scale of Disruption

	Conventional	New
Data Distribution Mechanism	Consistent Hashing (CH)	Group Hashing (NGH)
Storage Transport	Unicast (TCP)	Multicast (Replicast™)

This presentation:

- Existing technology vs. massive scale
- There are in fact better alternatives...

Distribution & Replication Mechanism: Data & Metadata



TOE (data storage)

SD 🗉



Increasing Decentralization





2015 Storage Developer Conference. © Nexenta Systems, Inc. All Rights Reserved.

I/O pipeline(*): unstructured, distributed and striped



SD 🛈

Quick Recap



CH vs. NGH

15

SD

- TCP/Unicast vs. Replicast
- Simulation: model, results, charts

2015 Storage Developer Conference. © Nexenta Systems, Inc. All Rights Reserved.

Consistent Hashing



Trademarks listed here and elsewhere in this presentation are property of their respective owners



Ceph vs. Consistent Hashing

CRUSH MAP aka Cluster Map

- Placement Groups
 - OSDs aka devices
 - OSD weights
- Rules: data redundancy
- Replica
- **Output:**
 - Pseudo-random uniform hash across failure domains – binomial(n, p)
- **Problems**:

15

□ Same as CH – see next

- Simplified CRUSH: HASH(node, replica)
- Rendezvous hashing: HASH(node, replica)
- Consistent hashing: HASH(replica)

Consistent Hashing: the good, the bad...

- □ The good
 - better than hash(chunk) % n-servers
 - even better when used with virtual nodes (partitions)
- The bad
 - □ Balls-into-Bins: load $\rightarrow \log(n)$: congestion "ripples"

 $\frac{\log n}{\log \log n} \cdot (1 + o(1))$

Subject to binomial spikes

The ugly

- Hashing determinism vs. Load balancing: either/OR
- Worse for larger clusters
- Much worse at a high utilization
- Is too consistent..

NGH: How it works

SD



Why NGH

- Admit: we simply don't know enough about the state on the edge
- Accept it.
- Main principles and motivations:
 - Serial transmit: one chunk at a time
 - Per-disk dedicated thread: one chunk at a time
 - Edge-based load balancing as part of the datapath
 - Edge-based reservation of link bandwidth:



Storage Transport: Data & Control



2015 Storage Developer Conference. © Nexenta Systems, Inc. All Rights Reserved.

Unicast vs. Persistent Redundancy

- Typical: put(...) blocks until 3 copies get persistently stored on 3 storage servers (variations: RAIN, EE, 2 cached, etc.)
- □ IO pipeline: redundancy first, completion second



TCP: the good, the bad...

- □ The good:
 - Predominant, totally integrated and reliable
- **The bad:**
 - Was never designed for scale-out storage clusters
- **The ugly:**
 - Unicast previous slide
 - TCP incast "is a catastrophic TCP throughput collapse…"
 - □ TCP "outcast" converged and hyper-converged clusters
 - Time to converge >> RTO
- **However**:
 - DCB helps
 - Custom/proprietary TCP kernels will help as well..

TCP Simulation

Each chunk put: 3 successive Date Date Dist

TCP transmissions

ACK every received replica

	Conventional	New
Data Distribution Mechanism	Consistent Hashing (CH)	Group Hashing (NGH)
Storage Transport	Unicast (TCP)	Multicast (Replicast™)

- ACK delivery time = chunk reception time + ½ * RTT
- Future-perfect TCP Congestion Protocol: Instant Convergence
 - N flows to one target, each instantly adjusting to use 1/Nth of the 10GE bandwidth
- TCP receive window and Tx buffer sizes unlimited
- Zero per-chunk overhead

Simulating TCP as a quintessential perfect Unicast

Replicast Simulation

15



Tire Kicking

https://github.com/Nexenta/nedge-sim

- ./build-ndebug
- ./do_benchmarks.sh
- **Example**:
 - ./test ngs 128 chunk_size 128 gateways 512 duration 100
- Extensive comments inside
- Detailed log tracking of each transaction
- Logs for the cited runs can be reproduced at <u>ff5fbda</u>

Default		Comment
#define CLUSTER_TRIP_TIME 10	000	// approximately 1us; RTT = 2us
#define N_REPLICAS	3	// replicas of each stored chunk
#define MBS_SEC_PER_TARGET_DRIVE	400	// 400MB/s (eMLC)





Simulation Results

15



Simulation Results (cont-d)



Simulation Results (last)



Conclusions

- Distributed storage clusters require technology that can simultaneously optimize:
 - available IOPS (budget) of the backend
 - storage capacity
 - network utilization
- To address all of the above, the mechanism must load balance at runtime
 - Consistent Hashing = blind selection
- Storage transport must be designed for scale-out and replication (3 replicas, etc.)
- Hence, NGH/Replicast
- Source on the github..

What's Next

SD (E







2015 Storage Developer Conference. © Nexenta Systems, Inc. All Rights Reserved.