.conf2015

# Indexer Clustering – Tips & Tricks

Da Xu

dxu@splunk.com

Software Engineer, Splunk

splunk>

# Disclaimer

During the course of this presentation, we may make forward looking statements regarding future events or the expected performance of the company. We caution you that such statements reflect our current expectations and estimates based on factors currently known to us and that actual events or results could differ materially. For important factors that may cause actual results to differ from those contained in our forward-looking statements, please review our filings with the SEC. The forward-looking statements made in the this presentation are being made as of the time and date of its live presentation. If reviewed after its live presentation, this presentation may not contain current or accurate information. We do not assume any obligation to update any forward looking statements we may make.

In addition, any information about our roadmap outlines our general product direction and is subject to change at any time without notice. It is for informational purposes only and shall not, be incorporated into any contract or other commitment. Splunk undertakes no obligation either to develop the features or functionality described or to include any such feature or functionality in a future release.
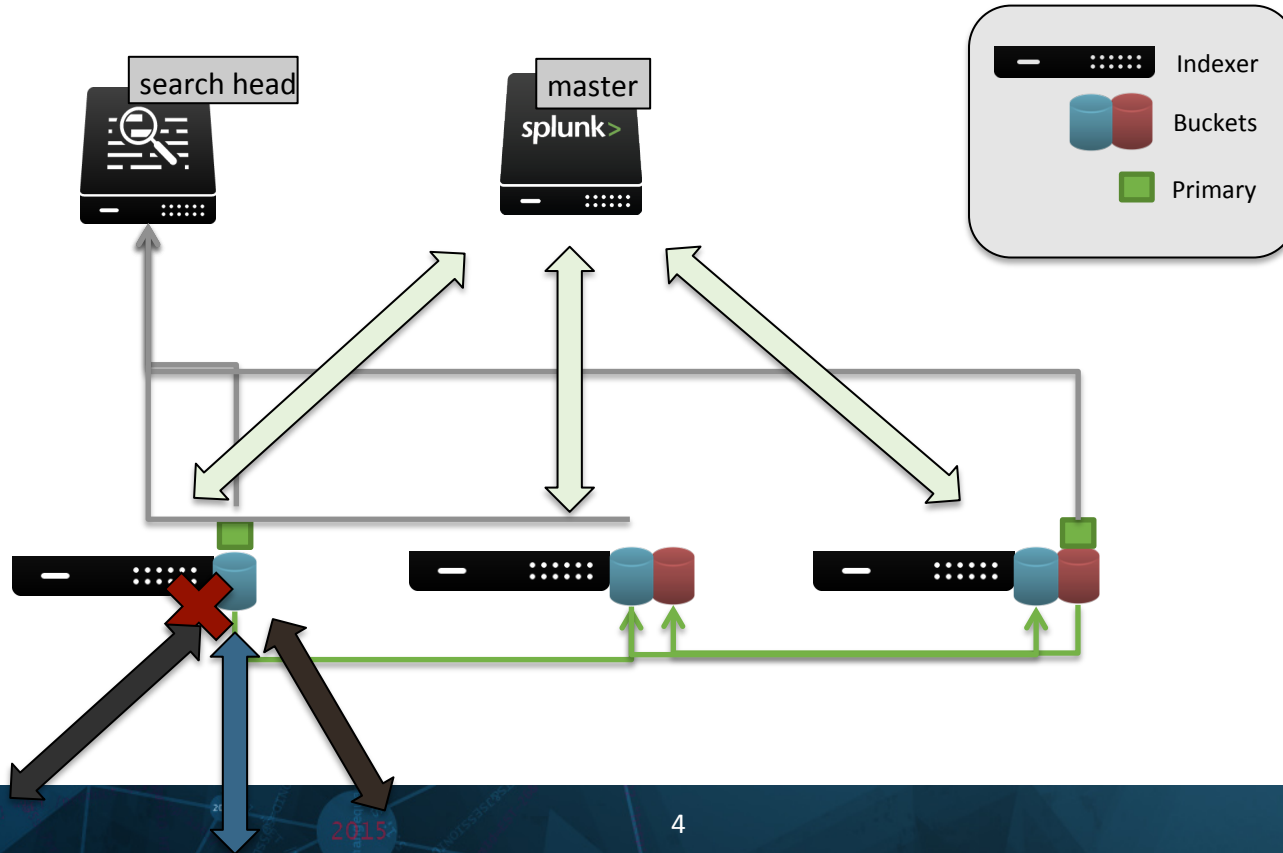
# Indexer Clustering Overview

# Cluster!



search head

master
splunk>

Indexer

Buckets

Primary
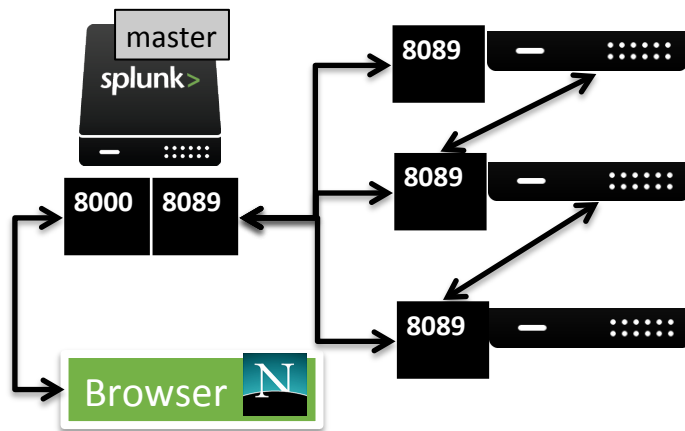
# Communication Through Endpoints

The cluster master and peers communicate amongst themselves through the clustering endpoints on the management ports. Some examples:

- Peers->Master:
  - /services/cluster/master/peers
    - Add Peer to cluster
    - Heartbeat to master
  - /services/cluster/master/buckets
    - Alert master there is a new bucket
    - Alert master a bucket changes (hot -> warm, warm -> frozen)
- Master->Peers
  - /services/cluster/slave/buckets
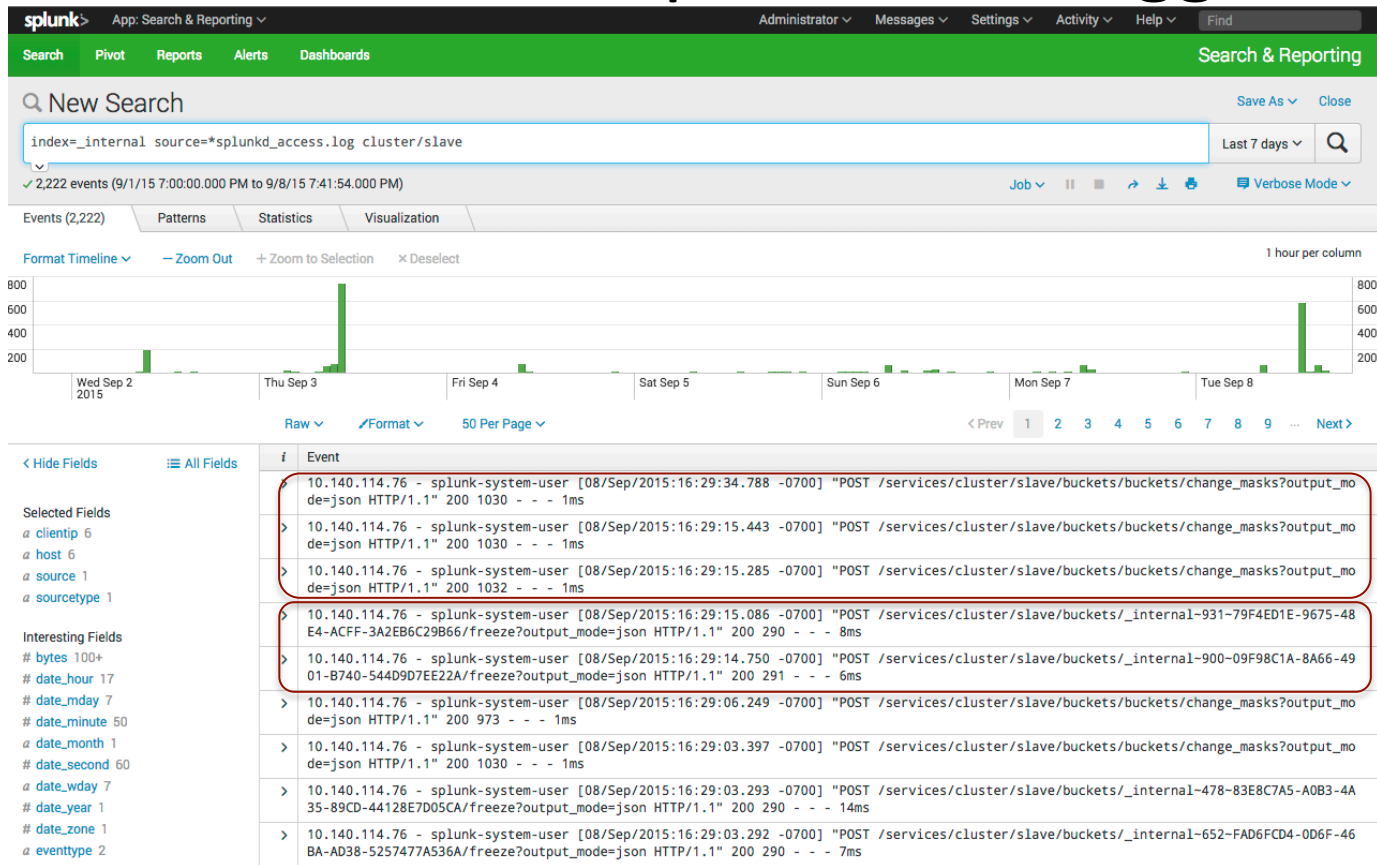    - Change primaries
    - Become searchable / unsearchable

# Endpoints Are Logged!



Bucket primary changes!

Buckets being frozen!

# Metrics.log

09-08-2015 22:59:15.184 -0700 INFO  Metrics - group=subtask_seconds, name=cmmaster_service, to_fix_streaming=0.000, to_fix_data_safety=0.016, to_fix_gen=0.000, to_fix_rep_factor=0.036, to_fix_search_factor=0.032, to_fix_sync=0.000, service=0.085

09-08-2015 22:59:15.184 -0700 INFO  Metrics - group=subtask_seconds, name=cmmaster_endpoints, clustermastergeneration_edit=0.018000, clustermasterinfo_list=0.018000, clustermasterpeers_edit=0.185000

09-08-2015 22:59:15.184 -0700 INFO  Metrics - group=subtask_counts, name=cmmaster_service, to_fix_streaming=0, to_fix_data_safety=97, to_fix_gen=0, to_fix_rep_factor=235, to_fix_search_factor=235, to_fix_sync=0, to_fix_added=0, to_fix_removed=0, to_fix_total=235, count=15

09-08-2015 22:59:15.184 -0700 INFO  Metrics - group=subtask_counts, name=cmmaster_endpoints, clustermastergeneration_edit=18, clustermasterinfo_list=18, clustermasterpeers_edit=185

09-08-2015 22:59:15.184 -0700 INFO  Metrics - group=executor, name=cmmaster_executor, jobs_added=0, jobs_finished=0, current_size=0, smallest_size=0, largest_size=0, max_size=0

09-08-2015 22:59:15.184 -0700 INFO  Metrics - group=cmmaster_servicejobs, serviced=0.000000, current_size=0.000000

09-08-2015 22:58:44.184 -0700 INFO  Metrics - group=subtask_seconds, name=cmmaster_service, to_fix_streaming=0.000, to_fix_data_safety=0.016, to_fix_gen=0.000, to_fix_rep_factor=0.036, to_fix_search_factor=0.031, to_fix_sync=0.000, service=0.084

09-08-2015 22:58:44.184 -0700 INFO  Metrics - group=subtask_seconds, name=cmmaster_endpoints, clustermastergeneration_edit=0.019000, clustermasterinfo_list=0.019000, clustermasterpeers_edit=0.181000

09-08-2015 22:58:44.184 -0700 INFO  Metrics - group=subtask_counts, name=cmmaster_service, to_fix_streaming=0, to_fix_data_safety=97, to_fix_gen=0, to_fix_rep_factor=235, to_fix_search_factor=235, to_fix_sync=0, to_fix_added=0, to_fix_removed=0, to_fix_total=235, count=16

09-08-2015 22:58:44.184 -0700 INFO  Metrics - group=subtask_counts, name=cmmaster_endpoints, clustermastergeneration_edit=19, clustermasterinfo_list=19, clustermasterpeers_edit=181

- Cluster master/slave activity can be found under cmmaster* or cmslave* groupings/names
- Metrics about cluster endpoints
  - How many times each endpoint was hit
  - How long we spent in those endpoints
- Metrics about jobs (rep fixup jobs, searchable fixup jobs, freeze jobs, etc)
  - How many jobs remain?
- How many # of buckets do we still need to fix?

# Clustering Logs/Activity

| splunkd_access.log | metrics.log |
|---|---|
| • Each individual endpoint access<br>    • (master-side) services/cluster/master/…<br>    • (indexer-side) services/cluster/slave/…<br>• How long we've spend at the endpoint (ms)<br>    • Higher times indicate the CM/Indexer is swamped with work (>50ms? >100ms?)<br>• The response (200 = success, non 200 = failure) | Metric information with regards to Clustering Activity, recorded every 30 seconds.<br><br>• name=cmmaster_endpoints<br>    • group=subtask_count  total number of accesses<br>    • group=subtask_seconds time Splunk spent responding to these endpoints<br>• name=cmmaster_executor<br>    • "Jobs" the CM has scheduled, finished, and current size of jobs to complete<br>        • Jobs are responsible for hitting the endpoints and performing the action (move-primary, freeze, etc)<br>• group=jobs, name=cmmaster<br>    • Actual counts of the jobs and their jobnames<br><br>Indexers have their own corresponding jobs (cmslave) |

# Cluster Activity

# Cluster Activity

# More Buckets More Problems

# More Buckets More Problems



- More buckets (and more peers) means the CM has to do more work
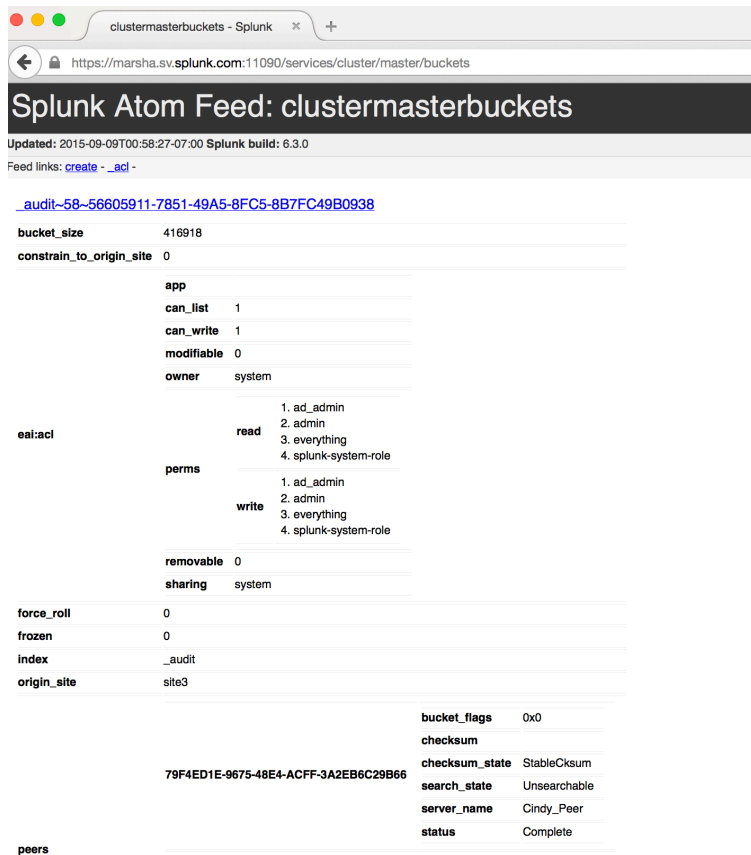  - Iterates through each bucket, checking whether it needs to queue up any fixup jobs
    - Replication Jobs (to meet RF)
    - Search Jobs (to meet SF)
    - Primary Jobs (all buckets need to have a primary copy per site)
    - Other jobs (freezing, checksum, rolling, etc)
- As the number of buckets grows, CM responsiveness goes down

# More Buckets More Settings

| server.conf | |
|---|---|
| service_interval (CM) | Specifies how often the CM should look through the buckets, scheduling jobs as necessary. Default = 1.<br>• Adjust to 1 sec for every 50k buckets. |
| heartbeat_period (Indexer) | Specifies how often the Indexers contact the CM. Defaults to every 1 second.<br>• For lots of peers ( >50) or lots of buckets (>100k), we can increase this value to 5-30. |
| heartbeat_timeout (CM) | Specifies how long before an Indexer is considered 'Down' when no heartbeats comes in.<br>• Multiple of heartbeat_period, anywhere from 20x – 60x |
| cxn_timeout (CM+Indexer)<br>rcv_timeout (CM+Indexer)<br>send_timeout (CM+Indexer) | Specifies how long before an intra-cluster connection will terminate. Default = 60.<br>• If a cluster indexer times out, it will re-add itself to the CM, which itself is a busy operation (it needs to resync the state of all its buckets), which can lead to negative feedback loops…<br>• These can be bumped up for busier clusters (300s). |
| indexes.conf | |
| rotatePeriodInSecs (Indexer) | Specifies how often to check through all the buckets – rolling them from hot->warm->cold as necessary. Default = 60<br>• 10min=600 |

splunk>

**clustermasterbuckets - Splunk**

https://marsha.sv.splunk.com:11090/services/cluster/master/buckets

## Splunk Atom Feed: clustermasterbuckets

Updated: 2015-09-09T00:58:27-07:00 **Splunk build:** 6.3.0

Feed links: create - _acl -

_audit~58~56605911-7851-49A5-8FC5-8B7FC49B0938

| | | | |
|---|---|---|---|
| bucket_size | 416918 | | |
| constrain_to_origin_site | 0 | | |
| | app | | |
| | can_list | 1 | |
| | can_write | 1 | |
| | modifiable | 0 | |
| | owner | system | |
| eai:acl | perms | read | 1. ad_admin 2. admin 3. everything 4. splunk-system-role |
| | | write | 1. ad_admin 2. admin 3. everything 4. splunk-system-role |
| | removable | 0 | |
| | sharing | system | |
| force_roll | 0 | | |
| frozen | 0 | | |
| index | _audit | | |
| origin_site | site3 | | |
| peers | 79F4ED1E-9675-48E4-ACFF-3A2EB6C29B66 | bucket_flags | 0x0 |
| | | checksum | |
| | | checksum_state | StableCksum |
| | | search_state | Unsearchable |
| | | server_name | Cindy_Peer |
| | | status | Complete |

services/cluster/master/buckets

- Which peers does the bucket exist on?
- Which peers is the bucket primary?
- Is the bucket searchable/unsearchable/ pending-searchable?

| | | bucket_flags | 0x4 |
|---|---|---|---|
| | | checksum | |
| | | checksum_state | StableCksum |
| | 09F98C1A-8A66-4901-B740-544D9D7EE22A | search_state | Searchable |
| | | server_name | Bobby_Peer |
| | | status | Complete |
| | | bucket_flags | 0x3 |
| | | checksum | |
| | | checksum_state | StableCksum |
| peers | 83E8C7A5-A0B3-4A35-89CD-44128E7D05CA | search_state | Searchable |
| | | server_name | Marsha_Peer |
| | | status | Complete |
| | | bucket_flags | 0x0 |
| | | checksum | |
| | | checksum_state | StableCksum |
| | FAD6FCD4-0D6F-46BA-AD38-5257477A536A | search_state | Searchable |
| | | server_name | Jan_Peer |
| | | status | Complete |

| | | |
|---|---|---|
| | site0 | 83E8C7A5-A0B3-4A35-89CD-44128E7D05CA |
| primaries_by_site | site1 | 83E8C7A5-A0B3-4A35-89CD-44128E7D05CA |
| | site2 | 09F98C1A-8A66-4901-B740-544D9D7EE22A |

| | | |
|---|---|---|
| rep_count_by_site | site1 | 2 |
| | site2 | 1 |

| | | |
|---|---|---|
| search_count_by_site | site1 | 2 |
| | site2 | 1 |

# Inspecting Buckets

clustermasterbuckets - Splunk

https://marsha.sv.splunk.com:11090/services/cluster/master/buckets

Splunk Atom Feed: clustermasterbuckets

Updated: 2015-09-09T00:58:27-07:00 Splunk build: 6.3.0

Feed links: create - _acl -

_audit~58~56605911-7851-49A5-8FC5-8B7FC49B0938

| bucket_size | 416918 |
|---|---|
| constrain_to_origin_site | 0 |

| | app | |
|---|---|---|
| | can_list | 1 |
| | can_write | 1 |
| | modifiable | 0 |
| | owner | system |
| eai:acl | perms | read | 1. ad_admin 2. admin 3. everything 4. splunk-system-role |
| | | write | 1. ad_admin 2. admin 3. everything 4. splunk-system-role |
| | removable | 0 |
| | sharing | system |

| force_roll | 0 |
|---|---|
| frozen | 0 |
| index | _audit |
| origin_site | site3 |

| | bucket_flags | 0x0 |
|---|---|---|
| | checksum | |
| 79F4ED1E-9675-48E4-ACFF-3A2EB6C29B66 | checksum_state | StableCksum |
| | search_state | Unsearchable |
| | server_name | Cindy_Peer |
| | status | Complete |
| peers | | |

There's so many buckets! How do I find one that I care about? Why would I care?

Filters! services/cluster/master/buckets?filter=

- Which buckets do not have primaries?
  - buckets?filter=has_primary=false
- Which buckets do not meet my RF=3?
  - buckets?filter=replication_count<3
- Which buckets are frozen?
  - buckets?filter=frozen=true
- Standalone?
  - buckets?filter=standalone=true
- Standalone and frozen?
  - buckets?filter=standalone=true&filter=frozen=true
  - (don't think this is a thing)
- Don't meet RF=3 and index=main?
  - buckets?filter=replication_count>3&filter=index=main

# Modifying Buckets

⚠️ ⚠️

Endpoints!

- Freeze a bucket:
  - curl -k -u admin:changeme https://{indexer}:{mgmt}/services/data/indexes/{INDEX}/freeze-buckets -d bucket_ids=46_11115C7A-E2F0-4225-A740-4ED6BD2D9CE5 -X POST
- Remove a copy of a bucket:
  - curl -k -u admin:changeme "https://{master}:{mgmt}/services/cluster/master/buckets/main~1490~D4A07A5D-3C3C-4D36-BD70-D610B432466F/remove_from_peer" -d peer={PEER_GUID}
- Remove all copies of a bucket:
  - curl -k -u admin:changeme "https://{master}:{mgmt}/services/cluster/master/buckets/main~1490~D4A07A5D-3C3C-4D36-BD70-D610B432466F/remove_all" -d peer={PEER_GUID
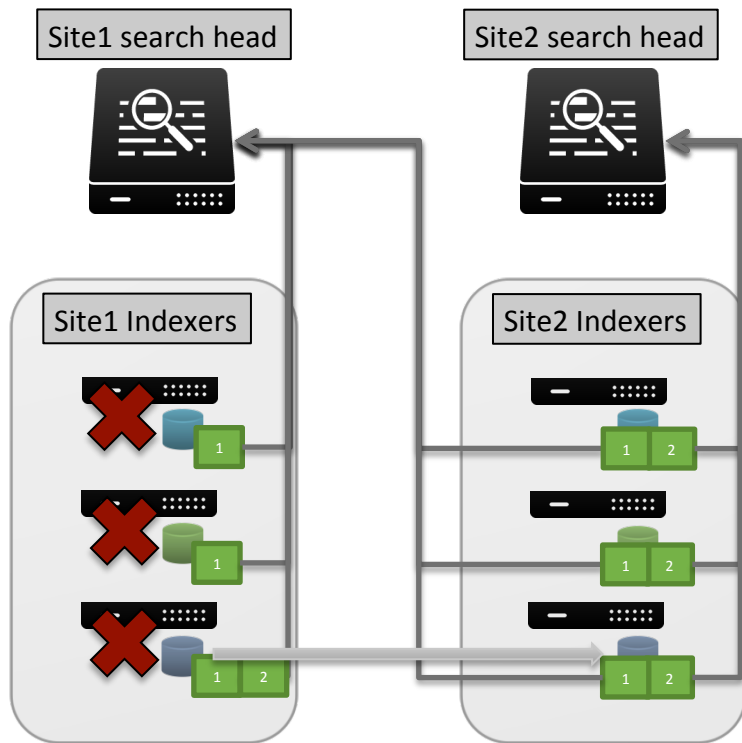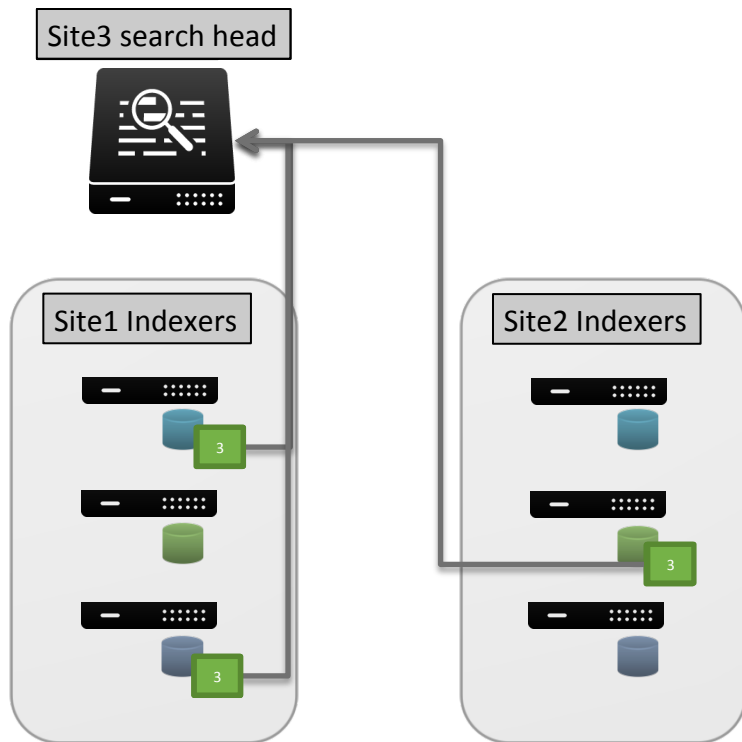
splunk>

Other Useful Knowledge

# Multisite Search Affinity



Site1 search head

Site2 search head

Site1 Indexers

Site2 Indexers

- When a searchable copy becomes available on a site, splunk will move the primary for that site to its local copy
- Buckets on a site will return events to a searchhead with the same site.
- If a peer goes down, the master will move the primaries that peer had to another copy
- If the entire site goes down, the other site(s) will become primaries

# Multisite ~~Search Affinity~~



Site3 search head

Site1 Indexers

Site2 Indexers

- Splunk 6.3 – site0
  - Primaries behave just like non-multisite, without any regards to site!
- Pre 6.3
  - Workaround!
  - Add another site to available_sites
  - Set SH (no indexers) to new site
  - Make sure to call "splunk set indexing-ready" on every CM restart
  - (wont work if your excess 'total' sites is greater than the # of non-specified sites… ie origin:1 total:3 in our illustration will not work, because then the CM will try to put the 2 non-origin copies into a site each, and there are no indexers in site3!)

# Stop Indexing on a Cluster-Indexer



- Detention Peer stops indexing data and doesn't accept any input, but still serves search queries
- 6.3 – turn on/off detention with an endpoint!
  - curl -k -u admin:changeme https://{INDEXER}:{MGMT}/services/cluster/slave/control/control/set_detention_override -d value=true -X POST
- Pre 6.3 – server.conf
  - [diskUsage] minFreeSpace=5000 (default)
    - Set to 50000000
    - (Requires a restart)

# Miscellaneous

Q&A