# Cassandra Availability with Virtual Nodes

Joseph Lynch
joe.e.lynch@gmail.com
josephl@netflix.com

Josh Snyder
josh@code406.com
joshs@netflix.com

April 16, 2018

## 1 Introduction

When Cassandra introduced vnodes in version 1.2, the database made new trade-offs with respect to cluster maintainability and availability [1] [2]. In particular, large Cassandra clusters became easier to scale and operate, but they lost availability during temporary or permanent node failures.

The two main benefits of vnodes are that the variance in data size held by any one host decreases, and that operators can more easily change the number of hosts participating in a cluster. In return, vnodes trade-off availability during node failures. As we document in this report, the availability cost of vnodes is both quantifiable and unavoidable. We present a model of the availability of a single-datacenter Cassandra cluster under failures. We attempt furthermore to predict, for a given cluster configuration, what availability guarantees can feasibly be promised to users of the database.

Based on our model, we conclude that the Cassandra default of 256 virtual nodes per physical host is unwise, as it results in steadily rising chance of unavailability as cluster size grows. This increasing failure rate is antithetical to Cassandra's design goal of fault-tolerance. We further use the model to argue for a number of possible solutions which both achieve operability goals of vnodes while at the same time providing much higher fault-tolerance, especially for large clusters (hundreds of nodes).

Section 2 introduces basics of the Cassandra data placement algorithm, and Section 3 discusses situations that can degrade Cassandra service. Readers already familiar with Cassandra operations may wish to skip to Analysis (Section 4) or Solutions (Section 4.3).

# 2 Cassandra Background

Data in Cassandra is distributed across the cluster hosts according to a consistent hashing algorithm. All data has a partition key, which Cassandra hashes to yield an integer bounded within the hash space. Each host in the cluster takes responsibility for storing data along one or more contiguous ranges of the hash space as decided by the replication strategies assigned to the keyspaces on the cluster.

Cassandra clusters achieve fault-tolerance in the `SimpleStrategy` and `NetworkTopologyStrategy` replication strategies by storing multiple copies of each piece of data. The quantity of copies (replication factor) is configurable, but is typically an odd number, with three being most common. Higher replication factors achieve greater fault-tolerance, at the cost of increased storage requirements.

Hosts in a Cassandra cluster are assigned one or more tokens which represent positions in a hash space created by the hashing function. Having more than one token per physical host leads to "virtual" nodes. The start of each range is variable, depending on the location of tokens held by other nodes in the cluster. The intervals in the hash space between tokens are known as *token ranges* and are the basis of all data replication. Depending on the keyspace replication strategy, and the placement of the "virtual" nodes, different physical hosts will be chosen to store data in a token range. For example, to achieve a replication factor of two using `SimpleStrategy`, Cassandra assigns two different physical hosts to each token range, while with `NetworkTopologyStrategy` Cassandra adds additional constraints to assure that data is replicated between racks and across datacenters.

# 3 Problem

In Cassandra's data placement strategy with a replication factor of three, availability is compromised when at least two hosts holding intersecting token ranges become unavailable to serve traffic. Practically speaking, unavailability results when a second node fails so quickly that the first hasn't had time to recover. Prior attempts have been made to model Cassandra's failure modes [4], but they often do not take into account the recovery aspect of Cassandra where it automatically restarts after process failure or in the case of permanent failure streams lost data to new hosts [1]. We approach the problem from a different direction, opting instead to model outages as a combination of losing a single node plus losing a neighboring node that owns an intersecting token range before recovery finishes, causing `UnavailableExceptions` to clients operating

---

[1] Not exactly automatically, but with the right automation it can be automatic

at `QUORUM` consistency levels. We consider any unavailability of the ring to be unacceptable to customers, whether it is of a large portion of the ring or small.

# 4    Analysis

We conduct our analysis on a cluster with replication factor 3 ($R = 3$). Under $R = 3$, the failure of two nodes with a shared token range leaves their shared token range(s) with only a single active node. The remaining node can continue to serve reads and writes, but *cannot* produce a quorum.

Without virtual nodes, a host in an $RF = 3$ cluster would inter-depend on 4 other hosts: its two left-neighbors and two right-neighbors. If the cluster is small enough, the same node may serve as both a left-neighbor and a right-neighbor. For example on a cluster with three hosts, each node inter-depends with two others; with four hosts each node has three inter-dependents; and with five or more hosts each node has four inter-dependents. When a brand new host places 256 virtual nodes into the ring, each vnode it places inter-depends with 2-4 other nodes.

By adding multiple token ranges to the responsibility of a single server, virtual nodes increase inter-dependency and decrease availability by coupling physical hosts together. Conversely, they theoretically increase the speed at which nodes can recover data by streaming[2].

Whenever it was necessary to build assumptions into the model, we chose them to make the model underestimate outages.

For the purposes of this analysis:

let:  
$n$ = number of hosts in the cluster  
$\lambda$ = average failure rate of hosts (1/interval)  
$S$ = size of the dataset per host (MB)  
$v$ = number of tokens per physical host  
$R$ = replication factor (number of replicas per token)  
$B_{inc}$ = maximum incoming bandwidth of a joining physical host (MB/s)  
$B_{out}$ = maximum outgoing bandwidth of a stream (MB/s)  

$n$ and $S$ are details of a particular Cassandra deployment, $v$ and $B_{out}$ are configuration options in `cassandra.yaml`, $R$ is determined by the keyspace replication strategy and $\lambda$ and $B_{inc}$ are determined by the hardware Cassandra is deployed to ($B_{inc}$ is also impacted by how efficiently Cassandra streams data, which we take into account with an arbitrary adjustment in the model).

---

[2]In our experience, at least with Cassandra 2.x and 3.0.x, inbound streaming tops out at significantly under network hardware capacity, but this simply makes this model underestimate outages

## 4.1 Availability Model

Crucial to the availability calculation is how quickly nodes recover: $E[T_{recovery}]$. The longer a host takes to recover, the longer a second failure can occur and cause an outage.

We model the recovery time of a host failure in Eq. 1 by taking the dataset $S$ and dividing it by the rate at which a recovering node can stream data from peers. To determine the speed, we must determine how many neighbors the node can stream from, which we model as a variant of the Birthday Problem [6]. For each token, we choose a random $2 * (R - 1)$ replica hosts (1b), those that precede and those that follow the token on the Cassandra ring. This is an underestimate of how many neighbors each additional token adds because actual Cassandra replication prohibits duplicates [5]. We then make the model underestimate outages further by assuming rack awareness, and that the number of racks equals $R$, which eliminates $\frac{n}{R}$ nodes from being possible replicas (1c). Finally we can determine the number of expected neighbors, $E[n_{neighbors}]$ by applying the formula for the number of distinct items expected from $k$ draws of $n_p$ items in 1d. To find the final recovery time we assume that the node is being replaced with `replace_address` rather than with `removenode` [3], so we only consider the incoming bandwidth of a single host (1e). This leads to 1f and concludes the recovery model.

$$E[T_{recovery}] = \frac{S}{speed} \tag{1a}$$

$$k = v * (2 * (R - 1)) \tag{1b}$$

$$n_p = n - \frac{n}{R} \tag{1c}$$

$$E[n_{neighbors}] = n_p * [1 - (1 - \frac{1}{n_p})^k] \tag{1d}$$

$$speed = \min(B_{inc}, E[n_{neighbors}] * B_{out}) \tag{1e}$$

$$E[T_{recovery}] = \frac{S}{\min(B_{inc}, E[n_{neighbors}] * B_{out})} \tag{1f}$$

This is a very conservative model, meaning that in reality vnodes would likely have a higher impact (e.g. because more nodes can be neighbors and may not have racks=$R$), but for cloud environments like AWS where Cassandra racks generally map 1:1 with availability zones this kind of model is reasonably accurate. This model gives us what we dub the "critical window" where any neighbor failing will cause a token range to lose quorum.

From Eq. 1 we can see that as $v$ increases, we get faster streaming until we

---

[3]See Section 4.3.3 for a discussion of why we think this replacement model is more appropriate for production

reach diminishing returns because we have saturated the inbound network of the joining node. At that point additional streams do not help speed recovery.

Now that we can model the expected recovery time, we begin the failure model by assuming that hosts fail as independent Poisson processes with parameter $\lambda$. Note that modeling host failures as a Poisson process represents a "best case scenario" for cluster availability: real-life computers are not so polite that they fail in a manner completely uncorrelated with their peers.

Under such a model the inter-arrival time of host failures follows an exponential distribution with parameter $\lambda$, and therefore once a single host failure happens, we lose a second replica of an overlapping range before the recovery finishes with probability:

$$P(outage|failure) = F(\tau; \lambda_{eff}) \tag{2a}$$

$$P(outage|failure) = 1 - e^{-\lambda_{eff}*\tau} \tag{2b}$$

$$\tau = E[T_{recovery}] \tag{2c}$$

$$\lambda_{eff} = E[n_{neighbors}] * \lambda \tag{2d}$$

$$P(outage|failure) = 1 - e^{-E[T_{recovery}]*E[n_{neighbors}]*\lambda} \tag{2e}$$

Equation 2b follows from the CDF of an exponential distribution with parameter $\lambda_{eff}$, and equation 2c and 2d follow because we only care about failures during the recovery period of the neighboring hosts (which form a combined Poisson process with effective parameter $\lambda_{eff}$).

Now, if we model each host as a splitting Poisson process with parameter $\lambda$ and splitting probability $P(outage|failure)$, then each node forms an outage Poisson process with parameter $\lambda_{split} = \lambda * P(outage|failure)$. If we then join the $n$ such processes, we have a resulting global Poisson process with parameter $\lambda_{global} = \lambda_{split} * n$. This yields Eq. 3, which models the average number of outages over an interval $\tau$.

$$
\begin{aligned}
E[outages] &= \lambda_{global} \\
&= (n * \lambda_{split}) \\
&= (n * (\lambda * P(outage|failure))) \\
&= (n * (\lambda * (1 - e^{-E[T_{recovery}]*E[n_{neighbors}]*\lambda_s})))
\end{aligned}
\tag{3}
$$

From this model we can see that availability is compromised by increasing the number of nodes ($n$), increasing the failure rate of those nodes ($\lambda$), increasing the number of neighbors (via $v$), or decreasing the recovery speed. To see how badly availability reacts with plausible default settings on a 96 node cluster in a cloud environment with machine failure rate $\lambda = \frac{25}{century} = \frac{0.25}{year}$.

let: $n$ = 96

$\lambda$ = 25 $\frac{1}{century}$ (= 0.25 $\frac{1}{year}$ = 7.9e-9 $\frac{1}{s}$)

$v$ = 256

$R$ = 3

$S$ = 307,200 (MB, 300 gigabyte)

$B_{inc}$ = 125 (MB/s, 1 gigabit)

$B_{out}$ = 25 / 2 = 12 (MB/s, 100 megabit)

We model this over a century because the probabilities of outage are so low in a single cluster that in any reasonable time span, the most likely number of outages is zero. To fairly compare the different configuration options, we model this system for a century and find that, on average, we expect 2.98 outages, or equivalently 0.03 per year:

$$
\begin{aligned}
A &= 96 * 25 * (1 - e^{-E[T_{recovery}]*E[n_{neighbors}]*\lambda_s}) \\
&= 96 * 25 * (1 - e^{-2457*64*7.927e-9}) \\
&= 2.98
\end{aligned}
\tag{4}
$$

In the case with only 4 vnodes, the availability is much better with only 0.35 failures per century = 0.0035 failures per year, a full 10x more availability!

$$
\begin{aligned}
A &= 96 * 25 * (1 - e^{-E[T_{recovery}]*E[n_{neighbors}]*\lambda_s} \\
&= 96 * 25 * (1 - e^{-2457*7.5*7.927e-9}) \\
&= 0.35
\end{aligned}
\tag{5}
$$

We can clearly see in Fig.1 that small numbers of vnodes do not appreciably change availability (due to the streaming benefits), but as the number of vnodes gets large we rapidly lose large amounts of availability until we level out after the number of vnodes passes the number of nodes. In fact, for this particular $\lambda$, the median number of outages is still zero until $v = 5$.

In Fig.2 we see the effects of varying the failure rate of machines, and observe that this has a significant impact on the availability of a Cassandra cluster running with vnodes. In particular, doubling the rate of machine failures can more than double the expected number of outages. This makes intuitive sense because as we add more tokens to each physical node, we are increasing the number of vulnerable replicas during a single node failure. Rack placement helps a lot to limit the number of potential replicas to only those residing in other racks, but it is not enough to prevent likely outage with a high number of tokens per node.
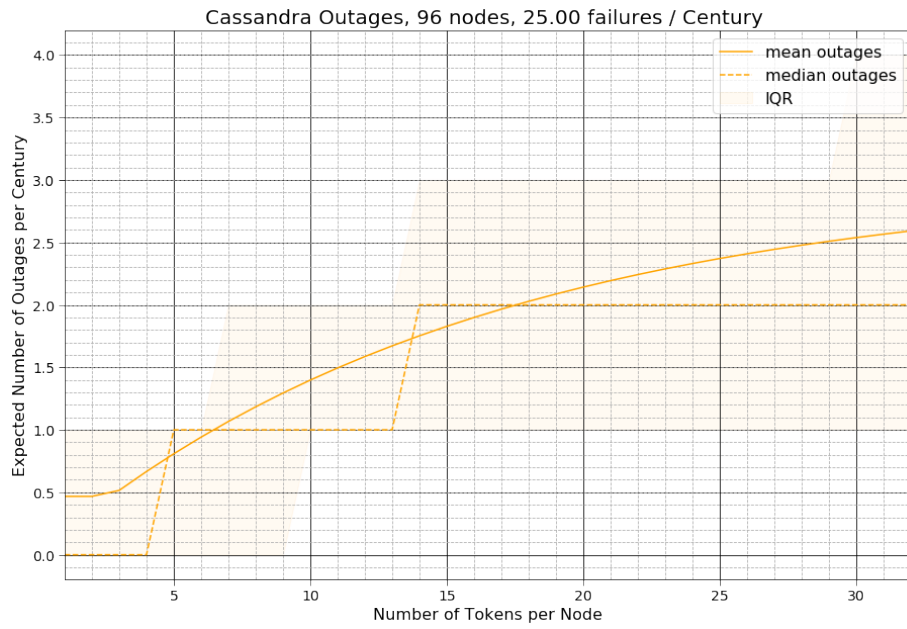
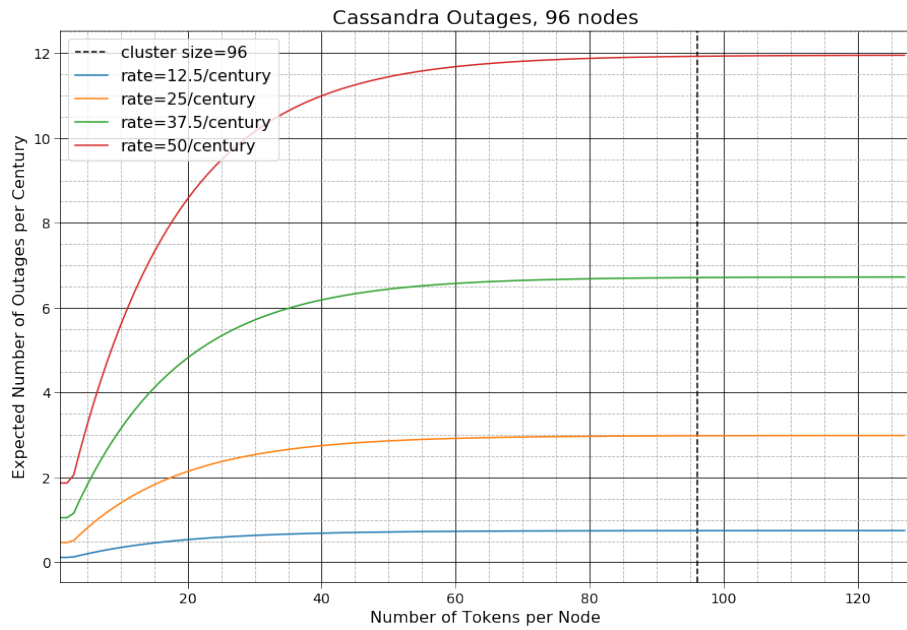Figure 1: Outages with Small Number of Tokens per Node



Figure 2: Outages with Varying Tokens per Node

It is also interesting to note that as clusters get larger, holding everything else constant, they become *less available*. This intuitively makes sense because as you increase the number of hosts significantly past RF, you are creating more replicas that can fail and be vulnerable to a second failure. Virtual nodes *significantly exacerbate* the problems for large clusters ($n > 200$) causing substantially more outages at large cluster sizes. This is shown in figure Fig. 3 and Fig. 4 . Fig 3 is particularly interesting because we can see with 256 vnodes, we rapidly approach an expected failure every year, compared to 16 vnodes where we expect an outage closer to every decade.
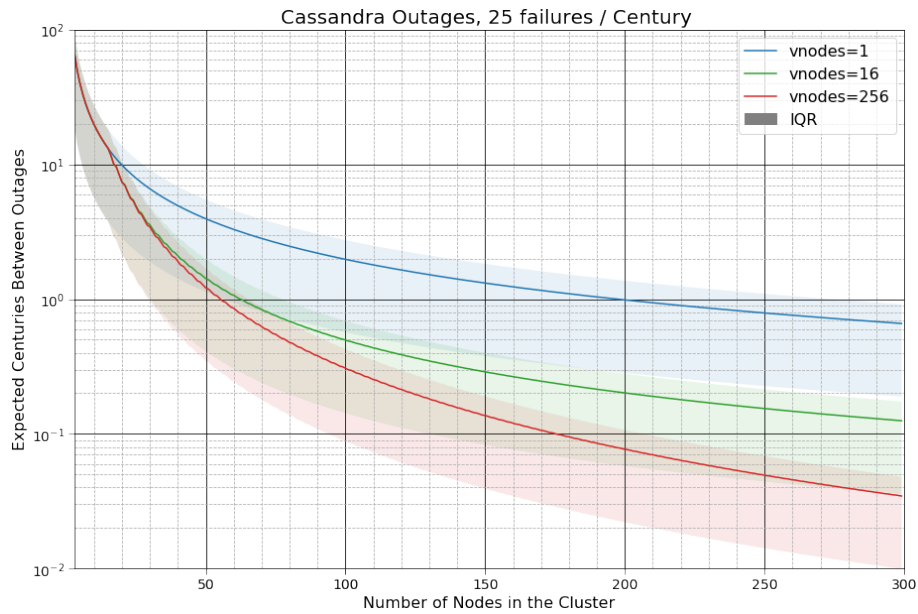


Figure 3: MTBF with Varying Nodes in the Cluster

However, as can be seen in Fig. 5, small clusters ($\leq 16$ for `NetworkTopologyStrategy` or $\leq 12$ for `SimpleStrategy`) are not appreciably impacted by vnodes because all their hosts are already neighbors with every possible host. As the cluster gets bigger, small clusters have a constant number of neighbors, but large clusters gain more and more neighbors, leading to 10-20x more risk of failure.
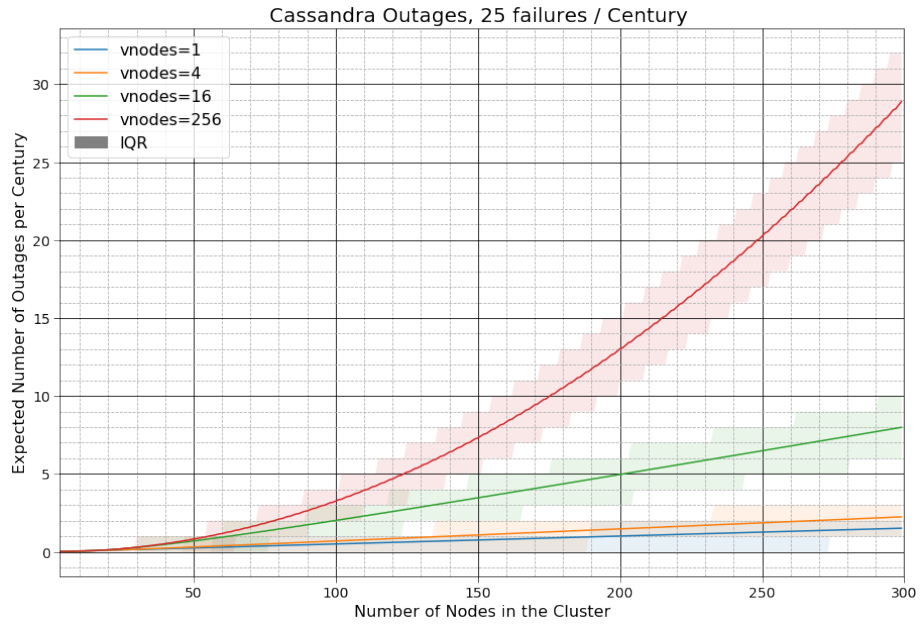
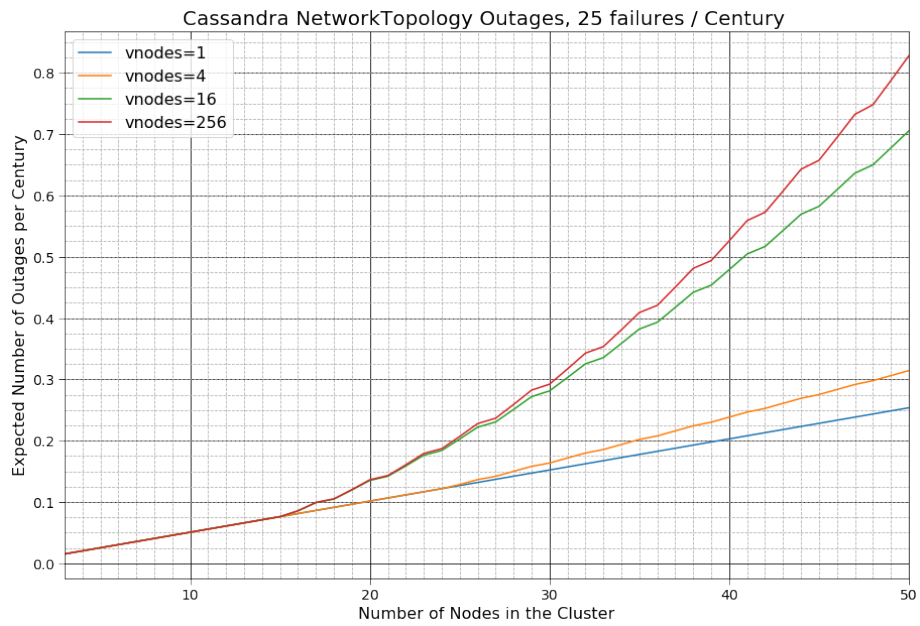Figure 4: Outages with Varying Nodes in the Cluster



Figure 5: Outages with Varying Nodes in the Cluster

The odd wiggle pattern you can see happens because rack awareness removes $n/R$ hosts from consideration; the `SimpleStrategy` graph is more smooth but also much less available in Fig 6.
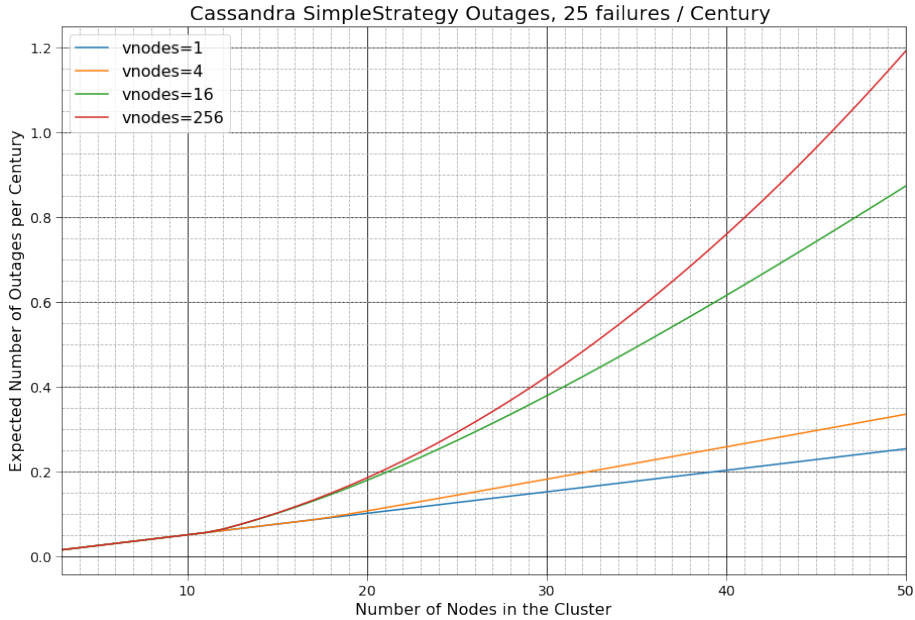


Figure 6: Outages with Varying Nodes in the Cluster

Clearly, while 4 or 16 virtual nodes do not appreciably impact availability, the default of 256 is quite unacceptable for large clusters.

## 4.2   Virtual Nodes, the Benefits

Having a number of tokens greater than one per physical host trades off availability, but it gains data distribution and operations benefits. In particular, a major problem pre-vnode was "hot" nodes, where one replica may have significantly more data than other nodes. This was particularly problematic when scaling up the cluster because in order to keep data evenly distributed you had to effectively double the cluster. In particular to add capacity to a cluster of size $n$ with $v$ vnodes you need proportionally fewer new nodes $N$ to do so. Now that Cassandra token placement is reasonably good at inserting new tokens where they are needed to balance the ring [9], we can model the number of nodes needed to scale up with Eq. 6 as seen in Fig. 7.

$$N = \left\lceil \frac{n}{v} \right\rceil \tag{6}$$

10

As we can clearly see, as we add more vnodes we must add fewer physical hosts to the cluster at a time in order to keep it balanced. While 1 token per host requires doubling, 16 tokens per host reduces this, naturally, by a factor of 16, which is a *significant* improvement. For a 96 host cluster instead of needing 96 more machines, we can add in increments of 6. The operational savings only get better as we add more tokens per host, until we reach the size of the cluster at which point we reach diminishing returns.
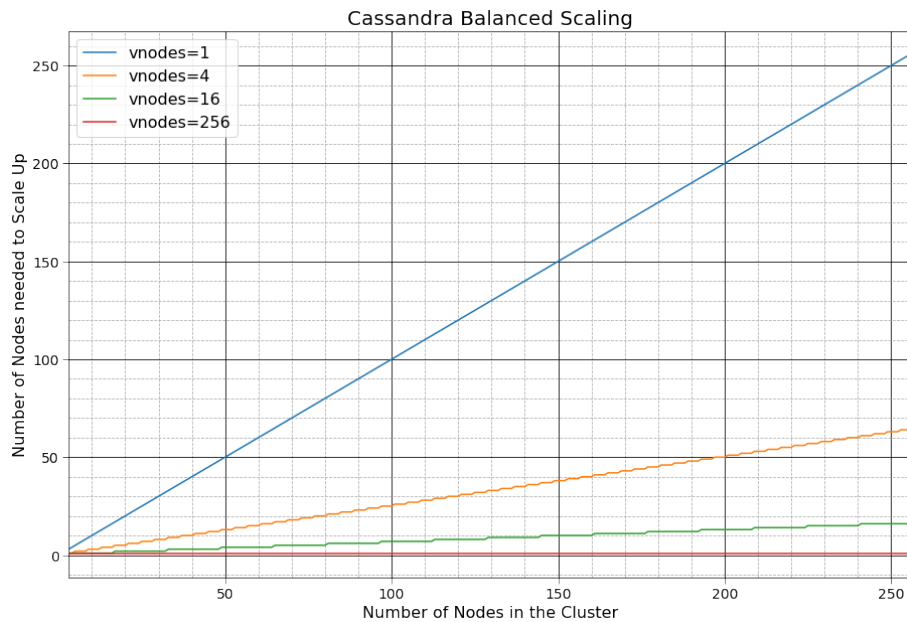


Figure 7: Nodes Required for Scaling Up

## 4.3 Solutions to the Availability Problem

There are a number of solutions to the availability problem caused by vnodes with various trade-offs. We briefly discuss some here:

### 4.3.1 Solution 1: Change the Number of vnodes

A simple solution to the availability problem is to make the trade-off less drastic. Moving to 4-16 tokens per node achieves almost all of the benefits of vnodes, while trading off significantly less availability. For example a setting of 16 gives 10x better availability than 256 while still having 10x easier scaling than a single token [10].

### 4.3.2 Solution 2: Decouple Data Transfer from Recovery

An elegant solution to this problem is to use networked attached storage such as EBS [11]. With a high reliability, low latency, network block device like AWS EBS `gp2`, Eq. 1 becomes a fixed constant of roughly five minutes. This increases availability by multiple orders of magnitude as replacement no longer varies with data size. Multiple tokens per host still reduce availability (since neighbors is still higher), but it is much less significant. Figure 8 shows just how much availability clusters can get by using EBS to store their data. The gains are even higher for large datasets (this model only has 300GB, if it were 2TB the gains would be significantly higher).
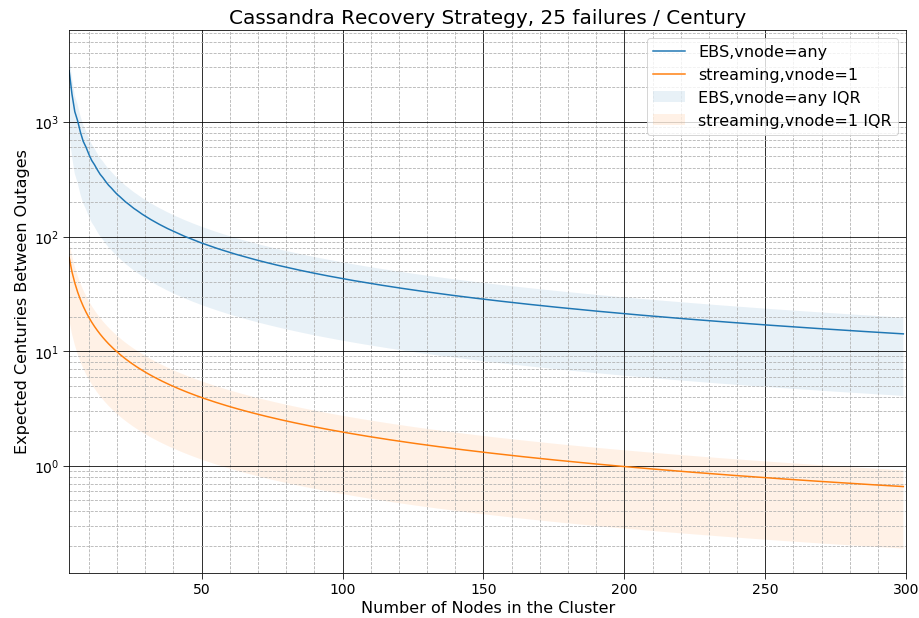


Figure 8: Availability Under Differing Recovery Strategies

### 4.3.3 Solution 3: Replacing Hardware via `removenode`

There is an alternative streaming model where instead of replacing a failed node with `replace_address`, the operator runs `removenode` on the dead node. This recovers the cluster to full replication rapidly as all nodes are participating in the re-replication of data, but simply with one fewer node present [8]. The subsequent bootstrap may take a long time, but that time is not spent at risk of outage. This technique is strictly dominated by `replace_node` for single token clusters because the data must be streamed twice, but for multiple token per node clusters it can make the availability much higher because `removenode`

theoretically involves almost every node in the cluster transferring data between themselves, which decreases $E[T_{recovery}]$ in Eq. 3 and limits outage.

This approach is likely unwise in production, however, as unlike `replace_address`, `removenode` causes operational load on nodes that are simultaneously doing latency sensitive query workloads. In particular `removenode` causes the following unique impacts (versus `replace_node`) to hosts taking latency sensitive queries:

1. Inbound network traffic
2. Disk IO and CPU cycles to write the data to disk (and later clean it up)
3. Loss of OS disk cache for data that will be invalidated shortly

Inbound network traffic is particularly problematic because unlike outbound network traffic which can be effectively be controlled via QoS [12], inbound network traffic is very difficult to insulate query traffic from, and would likely cause operational impact. In contrast, with `replace_address` the cluster operator can use `tc-fq` insulate query traffic from the streaming flows [13]

This approach would require significant production vetting before we believe it is a viable alternative to the production hardened `replace_node` approach.

### 4.3.4 Solution 4: Stop Using vnodes Entirely

An alternative way to solve the operational pain of large Cassandra clusters, while maintaining availability, is to use `moves` rather than tokens to re-balance data in a cluster [14][15]. This could either happen through a central planner which determines the optimal state of the cluster, or some kind of distributed greedy algorithm. Now that Cassandra has good streaming via the Netty refactor, this appears to be a very plausible approach and has additional benefits in repair and scaling (simultaneous bootstrap).

## 5   Conclusion

We have shown using a reasonable model that using a number of tokens per node greater than the number of physical nodes in the cluster significantly increases the expected outages, with default settings causing orders of magnitude more frequent outages. This availability is traded away for more even distribution of data and the ability to add or remove fewer nodes and still maintain even distribution. We have proposed a number of possible solutions, but we believe the current `cassandra.yaml` setting of 256 tokens per node has significantly decreased Cassandra cluster availability, especially at large cluster sizes, and are interested in improving the situation. Regardless of how the Cassandra community chooses to solve this problem, we believe it is worth understanding how vnodes have changed the math.

# References

[1] Brandon Williams: Virtual nodes in Cassandra 1.2 `https://www.datastax.com/dev/blog/virtual-nodes-in-cassandra-1-2`

[2] Sam Overton: RFC: Cassandra Virtual Nodes `https://www.mail-archive.com/dev@cassandra.apache.org/msg03837.html`

[3] Karger, et. al. Consistent Hashing and Random Trees: Distributed Caching Protocols for Relieving Hot Spots on the World Wide Web `https://dl.acm.org/citation.cfm?id=258660`

[4] Martin Kleppmann: The probability of data loss in large clusters `https://martin.kleppmann.com/2017/01/26/data-loss-in-large-clusters.html`

[5] Cassandra Data replication, `NetworkAwareTopologyStrategy` `https://docs.datastax.com/en/cassandra/latest/cassandra/architecture/archDataDistributeReplication.html`

[6] Kyle Siegrist: The Birthday Problem `http://www.randomservices.org/random/urn/Birthday.html`

[7] Apache Software Foundation: Replacing a Dead Node `http://cassandra.apache.org/doc/latest/tools/nodetool/removenode.html`

[8] Apache Software Foundation: nodetool removenode `http://cassandra.apache.org/doc/latest/tools/nodetool/removenode.html`

[9] Branimir Lambov: New token allocation algorithm in Cassandra 3.0 `https://www.datastax.com/dev/blog/token-allocation-algorithm`

[10] Chris Lohfink: Lower default num_tokens `https://issues.apache.org/jira/browse/CASSANDRA-13701`

[11] Jim Plush and Dennis Opacki: Amazon EBS & Cassandra `https://www.youtube.com/watch?v=1R-mgOcOSd4`

[12] Bert Hubert et al.: Linux Advanced Routing & Traffic Control `http://lartc.org/lartc.html#LARTC.QDISC`

[13] Eric Dumazet: Fair Queue packet scheduler `https://lwn.net/Articles/565421/`

[14] Apache Software Foundation: nodetool move `http://cassandra.apache.org/doc/latest/operating/topo_changes.html#moving-nodes`

[15] Stu Hood: Automatic, online load balancing `https://issues.apache.org/jira/browse/CASSANDRA-1418`

# 6 Appendix

Source code for all graphs is available on github as well as reproduced below.

```python
1
2 # coding: utf -8
3
4 # In[1]:
5
6 from __future__ import print_function
7 import math
8 import matplotlib
9 import numpy as np
10 import matplotlib.pyplot as plt
11 import matplotlib.patches as mpatches
12 from scipy.stats import poisson
13 from scipy.stats import expon
14
15 get_ipython().magic(u'matplotlib inline')
16
17 # Boils down to "If I pick hosts 2 * (rf - 1) * vnode times, how
        many
18 # distinct hosts will I have in expectation". Note that this is a
        slightly
19 # optimistic estimate because Cassandra won't place two replicas
        of the
20 # same token on the same machine or rack, but this is close enough
        for
21 # the model
22 # This is a variant of the Birthday Problem where we are interested
23 # in the number of distinct values produced
24 # http://www.randomservices.org/random/urn/Birthday.html
25 def num_neighbors(n, v, rf, strategy="rack"):
26     k = 2 * v * (rf - 1)
27     if strategy == "rack":
28         # As cassandra is rack aware, we assume #racks == #replicas
29         # This is maybe a bad assumption for some datacenter
        deployments
30         n = n - (n // rf)
31     else:
32         # SimpleStrategy
33         n = n - 1
34     estimate = (n * (1.0 - (1.0 - 1.0/n) ** k))
35     return max(rf - 1, min(estimate, n))
36
37 def p_outage_given_failure(recovery_seconds, num_neighbors,
        rate_in_seconds):
38     x = math.exp(-1 * recovery_seconds * num_neighbors *
        rate_in_seconds)
39     return 1 - x
40
41 def global_rate(node_rate, nodes, split_probability):
42     return node_rate * nodes * split_probability
43
44 def recovery_seconds(size, bw_in, bw_out, neighbors,
        recovery='streaming'):
45     if recovery == 'ebs':
```

```
46                return 60 * 5
47        return int(size / (min(bw_in, neighbors * bw_out)))
48
49
50  # Default model
51  nodes = 96
52  vnodes = 256
53  rf = 3
54  # 1000 gigabytes
55  node_dataset_mb = 300 * 1024
56  # MB/s
57  bw_in = 125
58  # MB/s, cassandra.yaml has 25MBPS as the default
59  # but most operators observe maybe half of that
60  bw_out = 25 / 2
61  strategy = 'rack'
62
63  year_seconds = 60.0*60*24*365
64  century_seconds = 100 * year_seconds
65
66  # Model machines that fail on average
67  # 25 times per century a.k.a 1 in 4 machines
68  # fails per year, or a machine fails every
69  # 4 years
70  arate = 25
71  arate_in_seconds = 25 / century_seconds
72
73
74  print("\nFailure Rate Variability")
75  print("Neighbors for {0} vnodes: {1:.3f}".format(1,
         num_neighbors(nodes, 1, rf)))
76  print("Neighbors for {0} vnodes: {1:.3f}".format(4,
         num_neighbors(nodes, 4, rf)))
77  print("Neighbors for {0} vnodes: {1:.3f}".format(16,
         num_neighbors(nodes, 16, rf)))
78
79  aneighbors = num_neighbors(nodes, vnodes, rf)
80  arecovery = recovery_seconds(node_dataset_mb, bw_in, bw_out,
         aneighbors)
81  print("Neighbors for {0} vnodes: {1:.3f}".format(vnodes,
         aneighbors))
82
83
84  def outage_stats(
85          vnodes, failure_rate_per_century, num_nodes,
86          rf, bw_in, bw_out,
87          strategy='rack', recovery='streaming'):
88      neighbors = num_neighbors(num_nodes, vnodes, rf, strategy)
89      recovery_s = recovery_seconds(node_dataset_mb, bw_in, bw_out,
         neighbors, recovery)
90      p_failure = p_outage_given_failure(
91          recovery_s, neighbors, failure_rate_per_century /
         century_seconds)
92
93      lmb = global_rate(failure_rate_per_century, num_nodes,
         p_failure)
94      return (
```

```
 95          poisson.mean(lmb), poisson.interval(0.50, lmb),
        poisson.median(lmb),
 96          expon.mean(scale=1/lmb), expon.interval(0.50, scale=1/lmb)
 97      )
 98
 99  # Returns outages _per century_
100  def compute_outage(
101          vnodes, failure_rate_per_century, num_nodes,
102          rf, bw_in, bw_out,
103          strategy='rack', recovery='streaming'):
104      return outage_stats(
105          vnodes, failure_rate_per_century, num_nodes, rf, bw_in,
        bw_out, strategy
106      )[0]
107
108  print("{0:<6} {1:<8} {2:<8} {3:<8} -> {4:<6}".format(
109      "rate", "rec_s", "p_fail", "g_lmb", "outages"
110  ))
111  for rate in (12.5, 25, 50, 100, 200):
112      recovery_s = recovery_seconds(node_dataset_mb, bw_in, bw_out,
        aneighbors)
113      p_failure = p_outage_given_failure(
114          recovery_s, aneighbors, rate / century_seconds)
115      gl = global_rate(rate, nodes, p_failure)
116      p = "{0:6.2f} {1:6.2f} {2:8.6f} {3:8.4f} -> {4:6.6f}".format(
117          rate, recovery_s, p_failure, gl, poisson.mean(gl)
118      )
119      print(p)
120
121
122  # In[2]:
123
124  plt.figure(figsize=(15,10))
125  plt.title("Cassandra Outages, {0} nodes".format(nodes),
        fontsize=20)
126  plt.ylabel("Expected Number of Outages per Century", fontsize=16)
127  plt.xlabel("Number of Tokens per Node", fontsize=16)
128  plt.xlim(1, 128)
129  plt.axvline(x=96, color='k', linestyle='--', label='cluster
        size={0}'.format(nodes))
130  plt.gca().grid(True, which='major', linestyle='-', color='k')
131  plt.gca().grid(True, which='minor', linestyle='--')
132  plt.gca().yaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(4))
133  plt.gca().xaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(4))
134  plt.tick_params(axis='both', which='major', labelsize=14, length=6)
135  plt.tick_params(axis='both', which='minor', length=0)
136
137  num_vnodes = range(1, 128)
138  rates = [12.5, 25, 37.5, 50]
139  for rate in rates:
140      outages = []
141      for vnode in num_vnodes:
142          outages.append(compute_outage(vnode, rate, nodes, rf,
        bw_in, bw_out))
143      print(outages[:3])
144      plt.plot(num_vnodes, outages,
        label="rate={0}/century".format(rate))
```

```
145 plt.legend(fontsize=16)
146
147 outages = [outage_stats(v, arate, nodes, rf, bw_in, bw_out) for v
        in num_vnodes[:32]]
148 outage_mean = [o[0] for o in outages]
149 outage_lower = [o[1][0] for o in outages]
150 outage_upper = [o[1][1] for o in outages]
151 outage_median = [o[2] for o in outages]
152
153 plt.figure(figsize=(15,10))
154 plt.title(
155     "Cassandra Outages, {0} nodes, {1:.2f} failures / Century
        ".format(
156         nodes, arate), fontsize=20)
157 plt.ylabel("Expected Number of Outages per Century", fontsize=16)
158 plt.xlabel("Number of Tokens per Node", fontsize=16)
159 plt.xlim(1, 32)
160 plt.gca().grid(True, which='major', linestyle='-', color='k')
161 plt.gca().grid(True, which='minor', linestyle='--')
162 plt.gca().yaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(5))
163 plt.gca().xaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(5))
164 plt.tick_params(axis='both', which='major', labelsize=14, length=6)
165 plt.tick_params(axis='both', which='minor', length=0)
166 plt.plot(num_vnodes[:32], outage_mean, color='orange', label='mean
        outages')
167 plt.plot(num_vnodes[:32], outage_median, color='orange',
        label='median outages', linestyle='--')
168 plt.fill_between(num_vnodes[:32], outage_lower, outage_upper,
        color='orange', alpha=0.05, label='IQR')
169 plt.legend(fontsize=16)
170
171
172 # In[3]:
173
174 # Hold failures constant, vary size of cluster
175 # Look at MTBF
176 plt.figure(figsize=(15,10))
177 plt.title(
178     "Cassandra Outages, {1} failures / Century ".format(
179         vnodes, arate), fontsize=20)
180 plt.ylabel("Expected Centuries Between Outages", fontsize=16)
181 plt.xlabel("Number of Nodes in the Cluster", fontsize=16)
182 plt.gca().grid(True, which='major', linestyle='-', color='k')
183 plt.gca().grid(True, which='minor', linestyle='--')
184 plt.gca().yaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(4))
185 plt.gca().xaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(5))
186 plt.tick_params(axis='both', which='major', labelsize=14, length=6)
187 plt.tick_params(axis='both', which='minor', length=0)
188 colors = plt.rcParams['axes.prop_cycle'].by_key()['color']
189 plt.gca().set_prop_cycle('color', colors[0:1] + colors[2:])
190
191 num_vnodes = (1, 16, 256)
192 num_nodes = range(3, 300)
193 lines = []
194
195 for v in num_vnodes:
```

```
196    outages = [outage_stats(v, arate, n, rf, bw_in, bw_out) for n
       in num_nodes]
197    outage_mean = [o[3] for o in outages]
198    outage_lower = [o[4][0] for o in outages]
199    outage_upper = [o[4][1] for o in outages]
200    line, = plt.semilogy(num_nodes, outage_mean,
       label="vnodes={0}".format(v))
201    lines.append(line)
202    plt.fill_between(
203        num_nodes, outage_lower, outage_upper, alpha=0.1,
204        label='vnodes={0} IQR'.format(v)
205    )
206
207 plt.xlim(3, 300)
208 plt.ylim(1/100.0, 100)
209
210 iqr_patch = mpatches.Patch(color='gray', label='IQR')
211 plt.legend(
212     handles=lines + [iqr_patch], loc='upper right', fontsize=16
213 )
214
215
216 # In[4]:
217
218
219 # Hold failures constant, vary size of cluster
220 plt.figure(figsize=(15,10))
221 plt.title(
222     "Cassandra Outages, {1} failures / Century ".format(
223         vnodes, arate), fontsize=20)
224 plt.ylabel("Expected Number of Outages per Century", fontsize=16)
225 plt.xlabel("Number of Nodes in the Cluster", fontsize=16)
226 plt.gca().grid(True, which='major', linestyle='-', color='k')
227 plt.gca().grid(True, which='minor', linestyle='--')
228 plt.gca().yaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(5))
229 plt.gca().xaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(5))
230 plt.tick_params(axis='both', which='major', labelsize=14, length=6)
231 plt.tick_params(axis='both', which='minor', length=0)
232
233 num_vnodes = (1, 4, 16, 256)
234 num_nodes = range(3, 300)
235 lines = []
236 for v in num_vnodes:
237     outages = [outage_stats(v, arate, n, rf, bw_in, bw_out) for n
       in num_nodes]
238     outages = [outage_stats(v, arate, n, rf, bw_in, bw_out) for n
       in num_nodes]
239     outage_mean = [o[0] for o in outages]
240     outage_lower = [o[1][0] for o in outages]
241     outage_upper = [o[1][1] for o in outages]
242     line, = plt.plot(num_nodes, outage_mean,
       label="vnodes={0}".format(v))
243     lines.append(line)
244     plt.fill_between(
245         num_nodes, outage_lower, outage_upper, alpha=0.1,
246         label='vnodes={0} IQR'.format(v)
247     )
```

```
248
249 plt.xlim(3, 300)
250 iqr_patch = mpatches.Patch(color='gray', label='IQR')
251 plt.legend(
252     handles=lines + [iqr_patch], loc='upper left', fontsize=16
253 )
254
255
256 # Hold failures constant, vary size of cluster,
        NetworkTopologyStrategy
257 plt.figure(figsize=(15,10))
258 plt.title(
259     "Cassandra NetworkTopology Outages, {0} failures / Century
        ".format(
260         arate), fontsize=20)
261 plt.ylabel("Expected Number of Outages per Century", fontsize=16)
262 plt.xlabel("Number of Nodes in the Cluster", fontsize=16)
263 plt.gca().grid(True, which='major', linestyle='-', color='k')
264 plt.gca().grid(True, which='minor', linestyle='--')
265 plt.gca().yaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(4))
266 plt.gca().xaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(10))
267 plt.tick_params(axis='both', which='major', labelsize=14, length=6)
268 plt.tick_params(axis='both', which='minor', length=0)
269
270 num_vnodes = (1, 4, 16, 256)
271 num_nodes = range(3, 51)
272 for v in num_vnodes:
273     outages = [compute_outage(v, arate, n, rf, bw_in, bw_out,
        'rack') for n in num_nodes]
274     print(v, outages[:3])
275     plt.plot(num_nodes, outages, label="vnodes={0}".format(v))
276
277 plt.xlim(3, 50)
278 plt.legend(fontsize=16)
279
280
281 # In[5]:
282
283 # Hold failures constant, vary size of cluster
284 plt.figure(figsize=(15,10))
285 plt.title(
286     "Cassandra SimpleStrategy Outages, {0} failures / Century
        ".format(
287         arate), fontsize=20)
288 plt.ylabel("Expected Number of Outages per Century", fontsize=16)
289 plt.xlabel("Number of Nodes in the Cluster", fontsize=16)
290 plt.gca().grid(True, which='major', linestyle='-', color='k')
291 plt.gca().grid(True, which='minor', linestyle='--')
292 plt.gca().yaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(4))
293 plt.gca().xaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(10))
294 plt.tick_params(axis='both', which='major', labelsize=14, length=6)
295 plt.tick_params(axis='both', which='minor', length=0)
296
297 num_vnodes = (1, 4, 16, 256)
298 num_nodes = range(3, 51)
299 for v in num_vnodes:
```

```
300      outages = [compute_outage(v, arate, n, rf, bw_in, bw_out,
         'simple') for n in num_nodes]
301      print(v, outages[:3])
302      plt.plot(num_nodes, outages, label="vnodes={0}".format(v))
303
304 plt.xlim(3, 50)
305 plt.legend(fontsize=16)
306
307
308 # In[6]:
309
310 # Observe impact of vnodes on Scale Up Balancing
311 plt.figure(figsize=(15,10))
312 plt.title("Cassandra Balanced Scaling",fontsize=20)
313 plt.ylabel("Number of Nodes needed to Scale Up", fontsize=16)
314 plt.xlabel("Number of Nodes in the Cluster", fontsize=16)
315 plt.gca().grid(True, which='major', linestyle='-', color='k')
316 plt.gca().grid(True, which='minor', linestyle='--')
317 plt.gca().yaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(5))
318 plt.gca().xaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(5))
319 plt.tick_params(axis='both', which='major', labelsize=14, length=6)
320 plt.tick_params(axis='both', which='minor', length=0)
321
322 num_vnodes = (1, 4, 16, 256)
323 num_nodes = range(3, 256)
324 for v in num_vnodes:
325      scale_up = [math.ceil(float(n) / v) for n in num_nodes]
326      plt.plot(num_nodes, scale_up, label="vnodes={0}".format(v))
327
328 plt.xlim(3, 256)
329
330 plt.legend(fontsize=16)
331
332
333 # In[9]:
334
335 # Observe impact of EBS on availability
336 plt.figure(figsize=(15,10))
337 plt.title(
338      "Cassandra Recovery Strategy, {0} failures /
         Century".format(arate), fontsize=20)
339 plt.ylabel("Expected Centuries Between Outages", fontsize=16)
340 plt.xlabel("Number of Nodes in the Cluster", fontsize=16)
341 plt.gca().grid(True, which='major', linestyle='-', color='k')
342 plt.gca().grid(True, which='minor', linestyle='--')
343 plt.gca().yaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(4))
344 plt.gca().xaxis.set_minor_locator(matplotlib.ticker.AutoMinorLocator(5))
345 plt.tick_params(axis='both', which='major', labelsize=14, length=6)
346 plt.tick_params(axis='both', which='minor', length=0)
347
348 num_nodes = range(3, 300)
349 lines = []
350
351 # Plot fixed recovery speed
352 outages = [outage_stats(1, arate, n, rf, bw_in, bw_out,
         recovery='ebs') for n in num_nodes]
353 outage_mean = [o[3] for o in outages]
```

```
354 outage_lower = [o[4][0] for o in outages]
355 outage_upper = [o[4][1] for o in outages]
356 line, = plt.semilogy(num_nodes, outage_mean, label="EBS,vnode=any")
357 lines.append(line)
358 plt.fill_between(
359     num_nodes, outage_lower, outage_upper, alpha=0.1,
360     label='EBS,vnode=any IQR'.format(v)
361 )
362
363 outages = [outage_stats(1, arate, n, rf, bw_in, bw_out,
364     recovery='recovery') for n in num_nodes]
364 outage_mean = [o[3] for o in outages]
365 outage_lower = [o[4][0] for o in outages]
366 outage_upper = [o[4][1] for o in outages]
367 line, = plt.semilogy(num_nodes, outage_mean,
        label="streaming,vnode=1")
368 lines.append(line)
369 plt.fill_between(
370     num_nodes, outage_lower, outage_upper, alpha=0.1,
371     label='streaming,vnode=1 IQR'
372 )
373
374 plt.xlim(3, 300)
375
376 plt.legend(fontsize=16)
377
378
379 # In[8]:
380
381 import itertools
382 import random
383 import sys
384
385 def simulate(l):
386     return random.expovariate(l)
387
388 def offset(values, max_value=float("inf")):
389     nvalues = values[:1] + [0] * (len(values) - 1)
390     for i in range(1,len(values)):
391         nvalues[i] = values[i] + nvalues[i-1]
392     return [n for n in nvalues if n <= max_value]
393
394 def outage(o, f_i, t, neighbors):
395     failures = 0
396     neighbor_indices = range(0, len(o))
397     neighbor_indices.remove(f_i)
398     random.shuffle(neighbor_indices)
399     for n in range(int(round(neighbors))):
400         failures += near(o[neighbor_indices[n]], o[f_i], t)
401     return failures
402
403 def near(a, b, t):
404     failures = 0
405     for i in a:
406         for j in b:
407             if j > i + t:
408                 break
```

22

```
409                 if j - i > 0 and j - i < t:
410                     failures += 1
411         return failures
412
413 def run_simulate(l, neighbors, nodes):
414     rs = []
415     for r in range(5):
416         o = [offset([simulate(l) for j in range(300)]) for i in
        range(nodes)]
417         maxes = [x[-1] for x in o]
418         m = max(maxes)
419         outages_per_century = (
420             sum([outage(o, i, arecovery / century_seconds,
        neighbors) for i in range(nodes)]) /
421             m
422         )
423         print("Run {0} gave {1:.3f} outages/century".format(r,
        outages_per_century))
424         rs.append(outages_per_century)
425     print("Simulation outages/century: ", sum(rs) / len(rs))
426
427 def run_simulate_naive(l, neighbors, nodes):
428     p_split = p_outage_given_failure(arecovery, aneighbors, l /
        century_seconds)
429     l_split = l * p_split
430     l_global = nodes * l_split
431     print(p_split, l_global)
432     rs = []
433     for r in range(10):
434         events = 1000
435         o = offset([simulate(l_global) for j in range(events)])
436         ##failure_years = [int(x/year) for x in o]
437         num_centuries = max(o)
438         rs.append(events / num_centuries)
439     print("Simple simulation outages/century: ", sum(rs) / len(rs))
440
441
442 print(p_outage_given_failure(arecovery, num_neighbors(nodes,
        vnodes, rf), arate))
443 print(vnodes, arate, nodes, rf, bw_in, bw_out, arecovery,
        num_neighbors(nodes, vnodes, rf))
444 print("Predicted outages/century:", compute_outage(vnodes, arate,
        nodes, rf, bw_in, bw_out))
445 #run_simulate(arate, num_neighbors(nodes, vnodes, rf), nodes)
446 run_simulate_naive(arate, num_neighbors(nodes, vnodes, rf), nodes)
447 #run_simulate(arate, num_neighbors(nodes, vnodes, rf), nodes)
```