# 4.0 Internode Messaging Test Plan

This document outlines the goals, and rough outline of approach and deliverables for evaluating the internode messaging work expected to land in CASSANDRA-15066, in preparation for the release of 4.0. The semantic verification work is already significantly underway as part of CASSANDRA-15066, with follow-up work to be listed below as Jira issues are assigned.

## Goals

- Build confidence that 4.0 messaging will be as or more reliable than in 3.x
- Quantify impact on cluster performance versus 3.x (and 14503) in a variety of circumstances
- Enable developers to modify the sub-system with rapid feedback on bugs and performance regressions

## Related Work

- https://issues.apache.org/jira/browse/CASSANDRA-14746
- https://issues.apache.org/jira/browse/CASSANDRA-15066

## Semantic Verification

### Unit testing

Before commit, there must exist tests covering critical and deterministically triggered behaviours, for the complete matrix of all connection methods and messaging versions. The goal here is to have extremely quick tests that can provide prompt feedback to developers that may have made a simple error when modifying the semantics of the system.

While we do not explicitly value code coverage metrics, we will produce a code coverage report of this work, and use it to ensure we have not missed any critical behaviours.

Out of scope is any concurrent or non-deterministic behaviours, as well as provider-specific exception handling.

### Randomised Testing

It is impossible to fully exercise a sophisticated concurrent system without stress testing it against a model. We will build a burn test that can be left running indefinitely, with the guarantee that it will eventually explore every possible logical system state. It will, at a minimum, simultaneously model all scenarios of:

- Overload, expiry and corruption
- Categories of failure to send or receive (de/serialize, write to wire, etc)
- Graceful and abrupt disconnects triggered by either the application or the provider
- Arbitrarily slow:
  - Connections
  - Message payload processing
  - Message de/serialization
  - Message submission by application
- Resource consumption

Verification entails confirming that in *all* of these scenarios, in any interleaving, the system behaviour is precisely as specified, i.e. that

- No messages are lost that cannot be explained *precisely* (i.e. with exact bounds on which messages can plausibly have been affected) by one of the above mechanisms
- All messages are delivered promptly
- No corrupt messages are delivered

### In-JVM dtests

While few in-jvm distributed unit tests will likely be written specific to this patch, the framework will be updated to successfully run any existing and future tests optionally over the network, including upgrade tests. As the body of these grow, particularly as part of our stability push for 4.0, the coverage of internode messaging will correspondingly expand.

### Real Cluster Workloads

As part of 4.0 release, several large contributing companies are undertaking qualification work of real (proprietary) workloads, as well as synthetic representative workloads. This patch will implicitly be exercised extensively as part of this work.

This work includes:

- Workload replay testing
- Large real cluster tests like those in CASSANDRA-14746, including those that may involve synthetic workloads
- Canary cluster tests

We should ensure that various extreme scenario are exercised in the process, including clusters of at least 1000 nodes.

# Performance Qualification

4.0 messaging, whether 14503 or 15066, utilises different libraries and underlying providers than 3.x. We must establish if there exists any scenario that suffers worse performance as a result of this transition. If any exist, they must be quantified and either resolved or, if this is not readily achievable, justified in the context of any positive impact of the patch.

## JMH Benchmarks

Using In-JVM dtests, it may be possible to benchmark actual cluster behaviours in an agnostic manner, and permit quick comparisons of 3.x, trunk and 15066. This may not be truly representative, but will allow relatively quick and easy comparisons of behaviours, and at least obvious regression detection.

We may also introduce some microbenchmarks for a simple point-to-point connection, and of time taken from OutboundConnection.enqueue to the wire. These may not be useful for comparison, but provide coarse regression detection.

## Real Cluster Workloads

Several community members, notably contributors from Netflix and The Last Pickle, have expressed an interest in participating in qualifying the characteristics of the internode messaging that lands in 4.0 using real clusters. The authors of 15066 will work with them to determine the best approaches for exercising behaviours in a real cluster, as well as qualifying the impact of this work versus 3.x and trunk.

We must be sure to check all extremes, including very small clusters of only 3 nodes, and giant clusters of >= 1000 nodes.

# Scheduling

Before commit, we anticipate completing all unit and randomised testing to produce a high confidence that the patch is semantically correct. Before we put a bow on things ready for the release of 4.0, we will complete the performance and real cluster verifications that had already haltingly begun. The project is focused exclusively on qualifying the upcoming 4.0 release, and remaining qualification work is best completed in conjunction with that effort.

## Non-Goals

Verification of behaviours outside of the 'net.async' package