



METHOD ARTICLE

Unit testing, model validation, and biological simulation [version 1; referees: 2 approved, 1 approved with reservations]

Gopal P. Sarma^{1,6*}, Travis W. Jacobs^{2,6*}, Mark D. Watts^{3,6}, S. Vahid Ghayoomie^{4,6}, Stephen D. Larson⁶, Richard C. Gerkin^{5,6}

¹School of Medicine, Emory University, Atlanta, USA

²Department of Bioengineering, Imperial College London, London, UK

³The University of Texas at Austin, Austin, USA

⁴Laboratory of Systems Biology and Bioinformatics, University of Tehran, Tehran, Iran

⁵School of Life Sciences, Arizona State University, Tempe, USA

⁶OpenWorm Foundation, Boston, USA

* Equal contributors

v1 First published: 10 Aug 2016, 5:1946 (doi: [10.12688/f1000research.9315.1](https://doi.org/10.12688/f1000research.9315.1))
 Latest published: 10 Aug 2016, 5:1946 (doi: [10.12688/f1000research.9315.1](https://doi.org/10.12688/f1000research.9315.1))

Abstract

The growth of the software industry has gone hand in hand with the development of tools and cultural practices for ensuring the reliability of complex pieces of software. These tools and practices are now acknowledged to be essential to the management of modern software. As computational models and methods have become increasingly common in the biological sciences, it is important to examine how these practices can accelerate biological software development and improve research quality. In this article, we give a focused case study of our experience with the practices of unit testing and test-driven development in *OpenWorm*, an open-science project aimed at modeling *Caenorhabditis elegans*. We identify and discuss the challenges of incorporating test-driven development into a heterogeneous, data-driven project, as well as the role of model validation tests, a category of tests unique to software which expresses scientific models.



This article is included in the **Neuroinformatics** channel.

Open Peer Review

Referee Status:

	Invited Referees		
	1	2	3
version 1 published 10 Aug 2016	report	report	report

- 1 Robert Cannon**, Textensor Limited UK
- 2 Christian Roessert**, Ecole Polytechnique Fédérale de Lausanne (EPFL) Switzerland
- 3 Andrew Davison**, French National Center for Scientific Research France

Discuss this article

Comments (0)

Corresponding author: Stephen D. Larson (stephen@openworm.org)

How to cite this article: Sarma GP, Jacobs TW, Watts MD *et al.* **Unit testing, model validation, and biological simulation [version 1; referees: 2 approved, 1 approved with reservations]** *F1000Research* 2016, 5:1946 (doi: [10.12688/f1000research.9315.1](https://doi.org/10.12688/f1000research.9315.1))

Copyright: © 2016 Sarma GP *et al.* This is an open access article distributed under the terms of the [Creative Commons Attribution Licence](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Grant information: This work was funded in part by NIMH (R01MH106674, RCG), and NIBIB (R01EB021711, RCG; and R01EB014640, Sharon M. Crook).

The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: No competing interests were disclosed.

First published: 10 Aug 2016, 5:1946 (doi: [10.12688/f1000research.9315.1](https://doi.org/10.12688/f1000research.9315.1))

Introduction

Software plays an increasingly prominent role in the biological sciences. This growth has been driven by an explosion in the availability of data and the parallel development of software to store, share, and analyze this data. In addition, simulations have also become a common tool in both fundamental and applied research^{1,2}. Simulation management (initialization, execution, and output handling) relies entirely on software.

Software used for collaborative biological research has an additional level of complexity (beyond that shared by other widely-used software) stemming from the need to incorporate and interact with the *results* of scientific research, in the form of large datasets or dynamical models. This added level of complexity suggests that technical tools and cultural practices for ensuring software reliability are of particular importance in the biological sciences³.

In this article, we discuss our experience in applying a number of basic practices of industrial software engineering—broadly known as *unit testing* and the related concept of *test-driven development*⁴⁻⁷—in the context of the OpenWorm project. OpenWorm (<http://www.openworm.org>) is an international, collaborative open-science project aimed at integrating the world's collective scientific understanding of the *C. elegans* round worm into a single computational model⁸. It is a diverse project incorporating data, simulations, powerful but intuitive user interfaces, and visualization. Since the goal of the project is to simulate an entire organism, the project and its underlying code are necessarily complex. The scope of the project is immense – OpenWorm has over fifty contributors from sixteen countries and projects divided into over forty-five sub-repositories under version control containing a total of hundreds of thousands of lines of code. For a project of this magnitude to remain

manageable and sustainable, a thorough testing framework and culture of test-driven development is essential⁴⁻⁷. In **Figure 1**, we show a diagrammatic overview of the many projects within OpenWorm and the relationship of testing to each of these. For extremely small projects, unit testing simply adds an overhead with little or no return on the time investment. As the project grows in size, however, the gains are quite significant, as the burden on the programmers of maintaining a large project can be substantially reduced.

In the code excerpts below, we will discuss 4 types of tests that are used in the OpenWorm code-base. They are:

- **Verification tests:** These are tests of basic software correctness and are not unique to the scientific nature of the project.
- **Data integrity tests:** These are tests unique to a project which incorporates data. Among other purposes, these tests serve as basic sanity checks verifying, for instance, that each piece of data in the project is associated with a scientific paper and corresponding DOI.
- **Biological integrity tests:** These are tests that verify correspondence with known information about static parameters that characterize *C. Elegans*, for example, the total number of neurons.
- **Model validation tests:** These are tests unique to projects which incorporate dynamic models. Model validation tests (using the Python package *SciUnit*) verify that a given dynamic model (such as the behavior of an ion channel) generates output that is consistent with known behavior from experimental data.

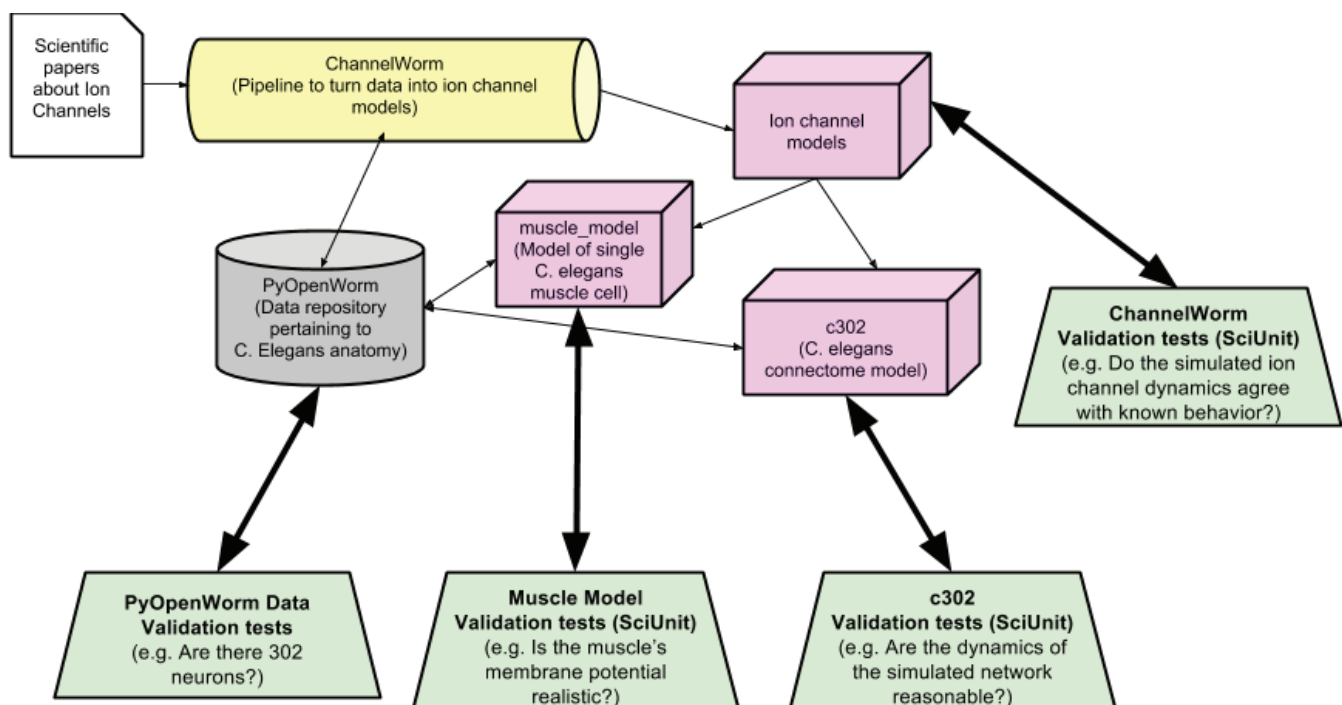


Figure 1. Diagram of some of the key OpenWorm modules and their corresponding testing frameworks.

The target audience for this article is computational biologists who have limited experience with large software projects and are looking to incorporate standard industrial practices into their work, or who anticipate involvement with larger projects in either academia or industry. We also hope that the exposition will be accessible to other scientists interested in learning about computational techniques and software engineering. We hope to contribute to raising the quality of biological software by describing some basic concepts of software engineering in the context of a practical research project.

Unit testing for scientific software

A simple introduction to unit testing

The basic concept behind software testing is quite simple. Suppose we have a piece of code which takes some number of inputs and produces corresponding outputs. A *unit test*, *verification test*, or simply *test* is a function that compares an input-output pair and returns a boolean value *True* or *False*. A result of *True* indicates that the code is behaving as intended, and a result of *False* indicates that it is not, and consequently, that any program relying on that code cannot be trusted to behave as intended.

Let us take a simple example. Suppose we have a function that takes a list of numbers and then returns them in sorted order, from lowest to highest. Sorting is a classic algorithmic task, and there are many different sorting algorithms with different performance characteristics; while the specific strategies they employ differ wildly, ultimately the result should be the same for any implementation. A unit test for one's sorting algorithm should take as input a list of numbers, feed it to the sorting algorithm, and then check that each element in the output list is less than or equal to the one that comes after it. The unit test would return *True* if the output list had that property, and *False* if not.

If one has multiple implementations of a sorting algorithm, then one can use a reliable reference implementation as a testing mechanism for the others. In other words, a test might return *True* if a novel sorting algorithm gives the same result as one widely known to be valid. There are other strategies along these lines. For example, suppose we have an implementation of an algorithm for multiplication called `multiply`. If we have a trusted implementation of an algorithm for addition, we can test that our multiplication algorithm works as expected by checking its behavior against the appropriate number of addition operations, e.g., `multiply(3,5) = 3 + 3 + 3 + 3 + 3`. See [Listing 1](#) for an implementation of this test in Python code.

In the previous example, the hypothetical unit test verified the core functionality of the algorithm. We had an algorithm that claimed to

sort things, and we wanted to check that it worked as advertised. But there are many other kinds of tests that we might be compelled to write in order to know that our software is working correctly. For instance, what happens if we feed an empty list to our sorting algorithm (this is an example of an *edge case*)? Should it simply return the list, generate an error message, or both? What if a user accidentally gives the algorithm something that is not a list, say for example, an image? What should the error message be in this case? Should there be a single error message to cover all cases, or should the error message be tailored to the specific case at hand? One can easily write unit tests to verify that the correct behavior has been implemented in all of these cases.

The sum total of all of the desired behaviors of an algorithm is called a *specification*, or *spec* for short. For instance, the specification for a sorting algorithm might look like the following:

- When given a list of numbers, return the list sorted from smallest to largest.
- When given a list of strings, return the list sorted in lexicographic order.
- If the input is an empty list, return the empty list and do not generate an error message.
- If the input is not a list, generate the error message "Input should be a list of real numbers or strings".
- If the input is neither a list of strings nor a list of numbers, return the same error message as above.

In [Listing 2](#), we have given a suite of unit tests for a sorting algorithm called `mySort` based on this specification. The basic notion demonstrated in the context of the sorting algorithm extends to any piece of software. In *OpenWorm*, we make extensive use of unit testing to verify both the functional properties of the system, as well as the validity of the data and models that comprise the simulation. For instance, the two tests given below in [Listing 3](#) check that any worm model has 302 neurons, and that the number of synapses for a given type of neuron is in accordance with its known value from the scientific literature. We will examine the different types of tests in more detail in the next section.

In *test-driven development*, the specification for a piece of software, as well as the corresponding unit tests are written *before coding the software itself*^{4,7}. The argument for test-driven development is that having a well-developed testing framework before beginning the actual process of software development increases the likelihood that bugs will be caught as quickly as possible, and furthermore, that it helps the programmer to clarify their thought processes.

Listing 1. Simple test for the multiplication operation.

```
1 def test_multiply():
2     """
3     Test our multiplication function against the
4     standard addition operator
5     """
6     assert multiply(3, 5) == 3 + 3 + 3 + 3 + 3
```

Listing 2. Sample tests for the sorting specification given in the text. The class `SortingTest` is a container for all of the individual tests that define the specification and can be extended if more tests are added.

```

1 import random
2 import unittest
3 from my_code import my_sort
4
5 """
6 Specification:
7 1) When given a list of numbers,
8 return the list sorted from smallest to largest.
9
10 2) When given a list of strings,
11 return the list sorted in lexicographic order.
12
13 3) If the input is an empty list,
14 return the empty list and do not generate an error message.
15
16 4) If the input is not a list, generate the error message:
17 ``Input should be a list of real numbers or strings``.
18 """
19
20 class SortingTest(unittest.TestCase):
21     """A class implementing tests for a sorting function"""
22     def setUp(self):
23         self.f = my_sort # The function we will test is mySort
24
25     def test_number_sort(self):
26         """Test that numbers sort correctly"""
27         sorted_list = range(100000)
28         shuffled_list = random.shuffle(range(100000))
29         self.assertEqual(self.f(shuffled_list), sorted_list)
30
31     def test_string_sort(self):
32         """Test that strings sort correctly"""
33         word_file = '/usr/share/dict/words'
34         words = open(word_file).read().splitlines()
35         sorted_words = words
36         shuffled_words = random.shuffle(words)
37         self.assertEqual(self.f(shuffled_words), sorted_words)
38
39     def test_empty_list(self):
40         """Test that empty list returns empty list"""
41         self.assertEqual(self.f([]), [])
42
43     def test_not_list(self):
44         """Test that invalid inputs generate correct error message"""
45         message = 'Input should be a list of real numbers or strings.'
46         self.assertRaisesRegex(TypeError, message, self.f, 'a')
47
48     def test_mixed_list(self):
49         """Test that mixed lists generate appropriate error message"""
50         mixed_list = [1, 2, 'a', 'b', 3]
51         message = 'Input should be a list of real numbers or strings.'
52         self.assertRaisesRegex(TypeError, message, self.f, mixed_list)

```

Listing 3. Excerpts from basic biological integrity tests for worm models. Given the size of the data repositories that OpenWorm relies upon, there are many simple tests such as these for ensuring the correctness of the associated data.

```

1 import PyOpenWorm
2 import unittest
3
4 class BiologicalIntegrityTest(unittest.TestCase):
5     """
6     Tests that read from the database and ensure that basic
7     queries have expected results, as a way to ensure data quality.
8     """
9     def test_correct_neuron_number(self):
10        """
11        This test verifies that the worm model
12        has exactly 302 neurons.
13        """
14        net = PyOpenWorm.Worm().get_neuron_network()
15        self.assertEqual(302, len(set(net.neurons())))
16
17    def test_neuron_syn_degree(self):
18        """
19        This test verifies that the number of chemical synapses
20        associated with a given neuron AVAL is equal to 90.
21        """
22        aval = PyOpenWorm.Neuron(name='AVAL')
23        self.assertEqual(aval.Syn_degree(), 90)

```

In practice, while some tests are written before-hand, others are written in parallel with the rest of code development, or shortly after a piece of code is written but before it is integrated.

We mention here that, in the software community, a distinction is often made between unit tests and *integration tests*⁷. Strictly speaking, a unit test is a test which is applicable to the smallest, functional unit of code, and which has no external dependencies. On the other hand, tests which verify that different components work together are classified as integration tests; they verify that multiple components are integrated correctly. Some of the tests discussed below would strictly be considered integration tests. For the sake of simplicity, we will not distinguish between unit tests and integration tests in this article, and will refer to both as simply *tests* or *unit tests*. The primary distinction that we make here is instead between ordinary *verification* tests (to verify that code works as intended) and *model validation* tests (to validate a model against experimental data), which we discuss in more depth below.

Unit testing in OpenWorm

The software that makes up OpenWorm shares common ground with all other pieces of software, whether the sorting algorithm described above, a word processor, or an operating system. As a result, there are a range of unit tests that need to be written to ensure that basic pieces of the software infrastructure function correctly. Many of these tests will not be of any scientific significance; they are simply sanity checks to ensure correct behavior for predictable cases. For instance, there are tests for checking that certain internal functions return the appropriate error messages when

given incorrect inputs; there are tests for verifying that databases are loaded correctly; there are tests which check that functions adhere to a specific naming convention which will help automated tools process the code-base.

As a data-driven, scientific research project, however, OpenWorm also makes use of several other categories of tests that do not typically appear in software development. For instance, the `PyOpenWorm` subproject of OpenWorm is a simple API that provides a repository of information about *C. elegans* anatomy (<https://github.com/openworm/PyOpenWorm>). Given that the aim of OpenWorm is to produce a realistic simulation of the nematode, a carefully curated repository of empirical information is a cornerstone of the project.

In the context of unit testing, there needs to be a category of tests that ensure that a curated datum has been appropriately verified and, furthermore, that its internal representation in the `PyOpenWorm` database is consistent. For example, for each “fact” in `PyOpenWorm`, there needs to be an associated piece of evidence, which serves as a reference. Practically, this evidence consists of a Digital Object Identifier⁹, or DOI, which corresponds to a research paper from which the fact was originally retrieved. For this class of tests, we traverse the database of facts and verify that for each fact there is an associated source of evidence, i.e., a DOI. Furthermore, these tests verify that each DOI is valid, and that the URL corresponding to the DOI is accessible. There are also tests to check the internal consistency of the `PyOpenWorm` database, for instance, that neurons with the same name have the same identifier.

Listing 4 gives several excerpts from the PyOpenWorm testing framework. It consists of tests to verify the references in the database, i.e., the DOIs which correspond to research papers.

In **Listing 5**, we give several tests for verifying the contents of the PyOpenWorm repository. Since each of the functions below is designed to test properties of Neuron objects, they are part of a single class called NeuronTest. These tests fall into the category of verification tests, and several of the tests, such as `test_name` and `test_type` simply check that the database is working correctly.

Model validation with SciUnit

Many computational models in biology are compared only informally with the experimental data they aim to explain. In contrast, we formalize data-driven model validation in OpenWorm by incorporating tests to *validate* each dynamical model in the project against experimental data from the literature. As an example, consider a scenario where a developer creates a new model and provides parameter values for a simulation. In addition to running all of the *verification* tests described above, the model and parameter values must be *validated* with respect to established experimental results. In general, each summary output of the model is validated against a corresponding piece of data. One example of a summary model

output is the “IV Curve” (i.e. current evoked in response to each of a series of voltage steps) of a given neuronal ion channel. We expect that our model will possess only ion channels which behave similarly to those observed experimentally, i.e. that the model IV Curve matches the experimentally-determined IV curve. If our model’s IV curve deviates too greatly from that observed experimentally, the model developers should be alerted and provided with information that will allow them to investigate the source of the discrepancy¹⁰. This may mean that parameter values must be modified, or in some cases the model itself must be substantially revised. In the case of OpenWorm, the necessary data for validating models is part of the PyOpenWorm and ChannelWorm subprojects (<https://github.com/openworm/ChannelWorm>), which are repositories of curated information about *C. elegans* anatomy and ion channels.

Ordinary unit testing frameworks do not readily lend themselves to this kind of model validation. Rather than simply comparing an input-output pair, model validation tests should perform the same procedure that a scientist would perform before submitting a newly hypothesized model for publication. That is, they should generate some kind of summary statistic encoding the deviation between experimental data and model output. For example, in the case of an IV Curve, one might use the area between the model and data curves as a summary statistic. In the case of OpenWorm, because

Listing 4. Verifying data integrity is an integral component of testing in OpenWorm. Below, we give several sample tests to verify the existence of valid DOIs, one technique used to ensure that facts in the PyOpenWorm repository are appropriately linked to the research literature.

```

1 import _DataTest # our in-house setup/teardown code
2 from PyOpenWorm import Evidence
3
4 class EvidenceQualityTests(_DataTest):
5     """A class implementing tests for evidence quality."""
6     def test_has_valid_resource(self):
7         """Checks if the object has either a valid DOI or URL"""
8         ev = Evidence()
9         allEvidence = list(ev.load())
10        evcheck = []
11
12        """Loop over all evidence fields in the database"""
13        for evobj in allEvidence:
14            if evobj.doi():
15                doi = evobj.doi()
16                val = requests.get('http://dx.doi.org/' + doi)
17                evcheck.append(val.status_code == 200)
18
19            elif evobj.url():
20                url = evobj.url()
21                val = requests.get(url)
22                evcheck.append(val.status_code == 200)
23
24            else:
25                evcheck.append(False)
26
27        self.assertTrue(False not in evcheck)

```

Listing 5. An assortment of verification tests from PyOpenWorm. These verify that the database behaves as we would expect it to, that properties of certain objects (Neuron objects, in this case) are correctly specified, and that the database is not populated with duplicate entries.

```

1 import _DataTest # our in-house setup/teardown code
2 from PyOpenWorm import Neuron
3
4 class NeuronTest(_DataTest):
5     """
6     AVAL, ADAL, AVAR, and PCVL are individual neurons in C. Elegans.
7     AB plapaaaap is the lineage name of the ADAL neuron.
8     A class implementing tests for Neuron objects.
9     """
10    def test_same_name_same_id(self):
11        """
12        Test that two Neuron objects with the same name
13        have the same identifier().
14        """
15        c = Neuron(name='AVAL')
16        c1 = Neuron(name='AVAL')
17        self.assertEqual(c.identifier(query=True), c1.identifier(query=True))
18
19    def test_type(self):
20        """
21        Test that a Neuron's retrieved type is identical to
22        its type as inserted into the database.
23        """
24        n = self.neur('AVAL')
25        n.type('interneuron')
26        n.save()
27        self.assertEqual('interneuron', self.neur('AVAL').type.one())
28
29    def test_name(self):
30        """
31        Test that the name property is set when the neuron
32        is initialized with it.
33        """
34        self.assertEqual('AVAL', self.neur('AVAL').name())
35        self.assertEqual('AVAR', self.neur('AVAR').name())
36
37    def test_init_from_lineage_name(self):
38        """
39        Test that we can retrieve a Neuron from the database
40        by its lineage name only.
41        """
42        c = Neuron(lineageName='AB plapaaaap', name='ADAL')
43        c.save()
44        c = Neuron(lineageName='AB plapaaaap')
45        self.assertEqual(c.name(), 'ADAL')
46
47    def test_neighbor(self):
48        """
49        Test that a Neuron has a 'neighbors' property, and that the
50        correct Neuron is returned when calling the 'neighbor' function.
51        """
52        n = self.neur('AVAL')
53        n.neighbor(self.neur('PVCL'))
54        neighbors = list(n.neighbor())
55        self.assertIn(self.neur('PVCL'), neighbors)
56        n.save()
57        self.assertIn(self.neur('PVCL'), list(self.neur('AVAL').neighbor()))

```


these models are part of a continuously updated and re-executed simulation, and not simply static equations in a research paper, the model validation process must happen automatically and continuously, alongside other unit tests.

To incorporate model validation tests, we use the Python package `SciUnit`¹¹ (<http://sciunit.scidash.org>). While there are some practical differences between writing `SciUnit` tests and ordinary unit tests, the concepts are quite similar. For example, a `SciUnit` test can be configured to return `True` if the test passes, i.e. model output and data are in sufficient agreement, and `False` otherwise. Ultimately, a scientific model is just another piece of software—thus it can be validated with respect to a specification. In the case of dynamical models, these specifications come from the scientific literature, and are validated with the same types of tests used before submitting a model for publication. `SciUnit` simply formalizes this testing procedure in the context of a software development work-flow.

In [Listing 6](#), we give an example of `SciUnit` tests using the neuron-specific helper library `NeuronUnit` (<http://neuronunit.scidash.org>) for neuron-specific models.

In the preceding example, the statistic is computed within the `SciUnit` method `judge`, which is analogous to the `self.assert` statements used in the ordinary unit tests above. While the ordinary unit test compares the output of a function pair to an accepted reference output, `judge` compares the output of a model (i.e. simulation data) to accepted reference experimental data. Internally, the `judge` method invokes other code (not shown) which encodes the test's specification, i.e. what a model must do to pass the test. The output of the test is a numeric score. In order to include `SciUnit` tests alongside other unit tests in a testing suite, they can be configured to map that numeric score to a boolean value reflecting whether the model/data agreement returned by `judge` is within an acceptable range.

Listing 6. Excerpt from a `SciUnit` test in `ChannelWorm`, a repository of information about ion channels. The test listed here verifies that a given ion channel has the correct current / voltage behavior. In terms of the informal classification of tests given above, this test falls under the category of model validation tests.

```

1 import os, sys
2 import numpy as np
3 import quantities as pq
4 from neuronunit.tests.channel import IVCurvePeakTest
5 from neuronunit.models.channel import ChannelModel
6 from channelworm.ion_channel.models import GraphData
7
8 # Instantiate the model; CW_HOME is the location of the ChannelWorm repo
9 ch_model_name = 'EGL-19.channel'
10 channel_id = 'ca_boyle'
11 ch_file_path = os.path.join(CW_HOME, 'models', '%s.nml' % ch_model_name)
12 model = ChannelModel(ch_file_path, channel_index=0, name=ch_model_name)
13
14 # Get the experiment data and instantiate the test
15 doi = '10.1083/jcb.200203055'
16 fig = '2B'
17 sample_data = GraphData.objects.get(
18     graph_experiment_reference_doi=doi,
19     graph_figure_ref_address=fig
20 )
21
22 # Current density in A/F and membrane potential in mV.
23 obs = zip(*sample_data.asarray())
24 observation = {'i/C':obs[1]*pq.A/pq.F, 'v':obs[0]*pq.mV}
25
26 # Use these observations to instantiate a quantitative test of the peak
27 # current (I) in response to a series of voltage pulses (V) delivered
28 # to the channel.
29 test = IVCurvePeakTest(observation)
30
31 # Judge the model output against the experimental data.
32 # Score will reflect a measure of agreement between I/V curves.
33 score = test.judge(model)
34 score.plot()

```

The output of these model validation tests can also be inspected visually; [Figure 2](#) shows the graphical output of the test workflow in [Listing 6](#), and illustrates for the developers why the test failed (mismatch between current-voltage relationship produced by the model and the one found in the experimental literature). Further details about the output of this test – including the algorithm for computing model/data agreement, and the magnitude of disagreement required to produce a failing score – can also be accessed via attributes and methods of the `score` object (not shown, but see [SciUnit](#) documentation). Consequently, full provenance information about the test is retained.

Some computational science projects use ad-hoc scripts that directly run models and compare their outputs to reference data. This can be adequate in simple cases, but for larger projects, particularly distributed open-source projects with many contributors, the mixing of implementation and interface carries significant drawbacks¹². For example, in order to record and store the membrane potential of a model cell—to then compare to reference data—one could determine which functions are needed to run the simulation in a given simulation engine, extract the membrane potential from the resulting files, and then call those functions in a test script. However, this approach has three major flaws. First, it may be difficult for a new contributor or collaborator to understand what is being tested, as the test code is polluted with implementation details of the model that are not universally understood. Second, such a test will not work on any model that does not have the same implementation details, and thus has limited reusability. Third, any changes to the model implementation will require parallel changes to the corresponding tests. In contrast, by separating tests from implementation details, tests can work on any model that implements a well-defined set of capabilities exposed via an interface. [SciUnit](#) does this by design, and [SciUnit](#) tests

interact with models only through an interface of standard methods, for example, those provided by [NeuronUnit](#). It is the responsibility of the model developer to match this interface by referencing standard methods, e.g. `run`, `get_membrane_potential`, etc. Ultimately, the separation of implementation from interface leads to greater code clarity, more rapid development, and greater test re-usability.

Test coverage

The *coverage* of a testing suite is defined as the percentage of functions in a code-base which are being tested. Since there is no rigorous measure of what constitutes an adequate test, precise figures of test coverage should be interpreted with caution. Nonetheless, automated tools which analyze a code-base to determine test coverage can be a valuable resource in suggesting areas of a code-base in need of additional attention. Ideally, test coverage should be as high as possible, indicating that a large fraction of or even the entire code-base has been tested according to the intended specifications.

In [PyOpenWorm](#), we make use several of pre-existing tools in the Python ecosystem for calculating test coverage of the Python code-base, specifically, the aptly-named [Coverage](#) package¹³, as well as a GitHub extension dedicated to tracking the coverage of such projects known as [Coveralls](#)¹⁴. We adopted these tools in an effort to track which parts of the code-base need additional tests, and to give further backing to the test-driven culture of the project. [PyOpenWorm](#) currently has a test coverage of roughly 73%. If a contributor to [PyOpenWorm](#) introduces some new code to the project but does not add tests for it, the contributor will see that test coverage has been reduced. By making changes in test coverage explicit, for example with a badge on the project's homepage, it is easier to track the impact of a growing code-base.

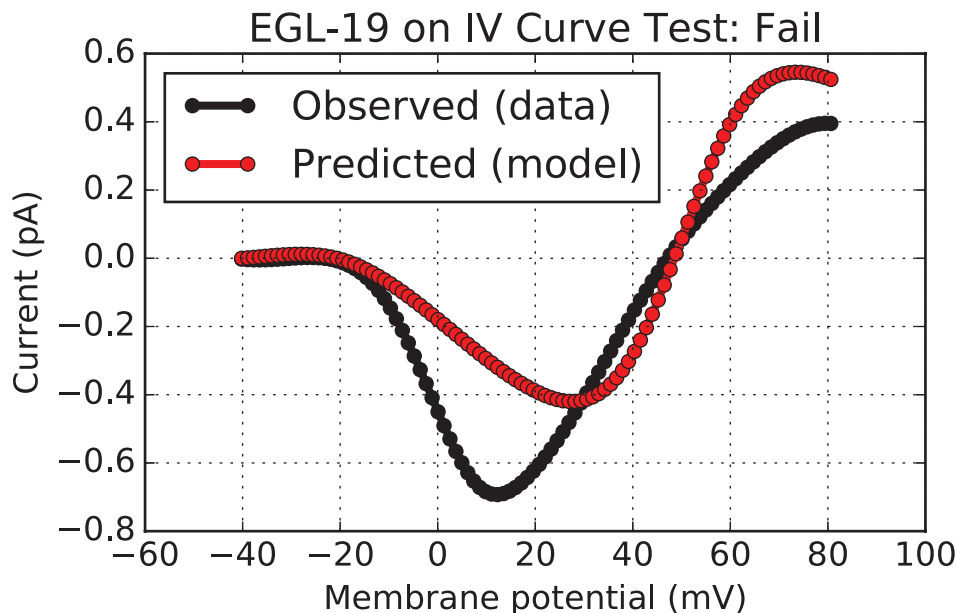


Figure 2. Graphical output from [Listing 6](#), showing a failed test which alerts developers to an inconsistency between model and data.

Continuous integration

Modern software is often written using a process of *continuous integration* or CI^{15,16}, whereby the contributions of developers are integrated into a shared repository multiple times a day by an automated system. Typically, the output of a testing suite will determine whether or not the new contributions of a developer can be immediately integrated, or whether changes are required to avoid *regression*, i.e. failing unit tests that passed before the new contribution.

The benefits of continuous integration include early detection of bugs, eliminating development bottle-necks close to the release date (in the case of commercial software), and the regular availability of usable versions of the software. The process of continuous integration also encourages shifts in how developers think about structuring their code, and encourages regular, modular contributions, rather than massive, monolithic changes that can be difficult to debug.

The entire OpenWorm project, including the PyOpenWorm and ChannelWorm modules make use of continuous integration (see Figure 3), taking advantage of a free service called Travis-CI (<https://travis-ci.org>) that tests changes to the code-base as they are pushed to the collaborative software development portal GitHub¹⁷. With each change, the entire project is built from scratch on a machine in the cloud, and the entire test suite is run. A build that passes all tests is a “passing build”, and the changes introduced will not break any functionality that is being tested. Because the entire project is built from scratch with each change to the code-base, the dependencies required to achieve this build must be made explicit. This ensures that there is a clear roadmap to the installation of dependencies required to run the project successfully – no hidden assumptions about pre-existing libraries can be made.

Skipped tests and expected failures

Suppose we have rigorously employed a process of test-driven development. Starting with a carefully designed specification, we









	fxi-again Add description of the channel object's properties Travis Jacobs committed	# 481 passed 38f62ee	🕒 14 min 36 sec 📅 17 days ago
	data-models Add docstring to channeltest Travis Jacobs committed	# 479 passed ca4933a	🕒 10 min 59 sec 📅 17 days ago
	synapse-loading#162 Merge branch 'dev+yarom_query' of http Travis Jacobs committed	# 459 errored 5772f40	🕒 22 sec 📅 24 days ago
	TestBreakdown With @travs: Gopal Sarma committed	# 420 passed dad4625	🕒 12 min 28 sec 📅 about a month ago
	ContinueTestBreakdown With @travs: Gopal Sarma committed	# 419 passed dad4625	🕒 11 min 26 sec 📅 about a month ago
	master removing note about lxml Stephen Larson committed	# 413 passed 5cc3042	🕒 11 min 29 sec 📅 2 months ago
	0.5.3 Merge pull request #157 from openworm/fix_install Travis Jacobs committed	# 411 passed 90cd5bc	🕒 10 min 50 sec 📅 2 months ago
	0.5.0 Remove alpha tag from version number Travis Jacobs committed	# 410 passed 4b86c67	🕒 9 min 26 sec 📅 2 months ago

Figure 3. Sample output from the OpenWorm continuous integration dashboard. Each row corresponds to a single set of contributions, known as a *commit*, submitted by a given developer. A commit is assigned a *build number*, which is given in the second column, and the result of the build process is indicated by the color of the corresponding row. If any of the unit tests fail, the build will be marked as failed (errored, in red), and the code contributions will be rejected. The developer is then responsible for identifying and fixing the corresponding bugs, and resubmitting their contributions to the code repository.

have written a test suite for a broad range of functionality, and are using a continuous integration system to incorporate the ongoing contributions of developers on a regular basis.

In this scenario, given that we have written a test suite prior to the development of the software, our CI system will reject all of our initial contributions because most tests fail, simply because the code that would pass the tests has not been written yet! To address precisely this scenario, many testing frameworks allow tests to be annotated as *expected failures* or simply to skip a given test entirely. The ability to mark tests as expected failures allows developers to incrementally enable tests, and furthermore draws attention to missing functionality. Consequently, the fraction of tests passed becomes a benchmark for progress towards an explicit development goal, *that goal being encoded by the set of all tests that have been written*.

The OpenWorm code-base makes extensive use of skipped tests and expected failures as a core part of the culture of test-driven development. In `PyOpenWorm`, for example, data integrity tests are often added in advance of the data itself being incorporated to the database. These tests provide a critical safety net as new information is curated from the scientific literature. Prior to the curation of this information, the tests are simply skipped. Once the information is curated, the tests are run, and indicate whether the information is usable by the project.

Frivolous tests and overly specific tests

Tests are typically sufficiently straightforward to write that it is easy to proliferate a testing suite with a large number of unnecessary tests. Often, these tests will be completely frivolous and cause no harm, beyond causing a testing suite to take much longer than necessary to run. However, tests which are overly specific can actually hinder the process of development. If there are tests which are too specific and constrain internal behavior that is not meant to be static, a developer's contributions may be unnecessarily rejected during the process of continuous integration.

Conclusions

Our aim in this article is to give an overview of some basic development practices from industrial software engineering that are of particular relevance to biological software. As a summary, we list here the types of tests used in OpenWorm. This list is simply an informal classification, and not a definitive taxonomy:

Verification tests (the usual suspects) These are tests common to all pieces of software and are not particularly relevant to the biological nature of the project. For instance, tests that verify that error handling is implemented correctly, that databases are accessed correctly, or that performing certain numerical operations produces results within an acceptable range.

Data integrity tests These are tests unique to a project that incorporates curated data. In the case of OpenWorm, these tests check (among other things) that every biological fact in the `PyOpenWorm` repository has an associated piece of experimental evidence, typically corresponding to a DOI, and that each of these DOIs is valid.

Biological integrity tests These tests verify that data tokens in the `PyOpenWorm` repository correspond to known information about *C. Elegans*. In contrast to the model validation tests described below, biological integrity tests typically only check static information/parameters.

Model validation tests These are tests specific to a project that incorporates scientific models. Model validation tests allow us to check that specific models, such as the behavior of ion channels, correspond to known behavior from the scientific literature. In effect, they extend the notion of unit testing to compare summary data and model output according to some summary statistic. In OpenWorm, the Python package `SciUnit` and derivative packages like `NeuronUnit` are used for writing tests that check the validity of scientific models against accepted data.

It should be clear from the above discussion and corresponding code examples that unit tests are fundamentally quite simple objects. Their behavior is no more than to compare input-output pairs, or in the case of `SciUnit` tests, that a given model's output corresponds to a known reference from the scientific literature. The sophistication of testing frameworks is generally quite minimal when compared to the software itself being tested. While ad-hoc test scripts may be sufficient for small projects, for large projects with many contributors, a systematic approach to unit testing can result in significant efficiency gains and ease the burden of long-term code maintenance. In the context of *continuous integration*, whereby a piece of software is built in an ongoing cycle as developers make changes and additions to the code-base, unit testing provides a valuable safety net that can prevent flawed code from prematurely being integrated.

However, in spite of the conceptual simplicity and potential pitfalls of testing, its importance cannot be overstated. Writing tests requires careful thought and planning and some knowledge of the code-base being tested. Testing from a specification alone can result in inadequate testing, but tests which are too specific to the code-base can result in unnecessary roadblocks for developers.

Rather than being thought of as a sophisticated set of technical tools, unit testing should be viewed as a cultural practice for ensuring the reliability of complex software. Perhaps a useful analogy is the powerful impact that checklists have had in clinical medicine, aviation, construction, and many other industries^{18–20}. Unit tests are sanity checks at a minimum, and can potentially guide the

scientific development of models when used in conjunction with experimental data. In order to reap their benefit, their existence and maintenance needs to be valued by all of the participants of the research and software development process. Finally, in order for this culture to be created, test-driven development should not be a heavy-handed imposition on the developers. Otherwise, it will be incorrectly perceived as a bureaucratic hurdle, rather than the valuable safety-net that it is.

Software availability

Software available from: <http://www.openworm.org/>
Latest source code: <http://github.com/OpenWorm>

Author contributions

GPS, TWJ, SDL, and RCG wrote the manuscript. All authors contributed to the unit testing framework in the Open Worm project.

Competing interests

No competing interests were disclosed.

Grant information

This work was funded in part by NIMH (R01MH106674, RCG), and NIBIB (R01EB021711, RCG; and R01EB014640, Sharon M. Crook).

The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Acknowledgements

We would like to thank supporters of OpenWorm, including NeuroLinx.org, the backers of the 2014 OpenWorm Kickstarter campaign (<http://www.openworm.org/supporters.html>), Google Summer of Code 2015, and the International Neuroinformatics Coordinating Facility. We would also like to thank the scientific and code contributors to OpenWorm (<http://www.openworm.org/people.html>), and Shreejoy Tripathy for careful reading of the manuscript.

References

1. Takahashi K, Yugi K, Hashimoto K, *et al.*: **Computational Challenges in Cell Simulation: A Software Engineering Approach**. *IEEE Intelligent Systems*. 2002; **17**(5): 64–71.
[Publisher Full Text](#)
2. Macklin DN, Ruggero NA, Covert MW: **The future of whole-cell modeling**. *Curr Opin Biotechnol*. 2014; **28**: 111–115.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
3. Gewaltig MO, Cannon R: **Current practice in software development for computational neuroscience and how to improve it**. *PLoS Comput Biol*. 2014; **10**(1): e1003376.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
4. Beck K: **Test Driven Development: By Example**. Addison Wesley, 2002.
[Reference Source](#)
5. Maximilien EM, Williams L: **Assessing test-driven development at IBM**. In *Software Engineering, 2003. Proceedings. 25th International Conference on*. IEEE, 2003; 564–569.
[Publisher Full Text](#)
6. Erdogmus H, Morisio M, Torchiano M: **On the effectiveness of the test-first approach to programming**. *IEEE Transactions on Software Engineering*. 2005; **31**(3): 226–237.
[Publisher Full Text](#)
7. Osherove R: **The Art of Unit Testing: with examples in C#**. Manning Publications, 2013.
[Reference Source](#)
8. Szigeti B, Gleeson P, Vella M, *et al.*: **OpenWorm: an open-science approach to modeling *Caenorhabditis elegans***. *Front Comput Neurosci*. 2014; **8**: 137.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
9. **DOI® System and Internet Identifier Specifications**. 2015. Accessed: 2015-07-24.
[Reference Source](#)
10. De Schutter E: **The dangers of plug-and-play simulation using shared models**. *Neuroinformatics*. 2014; **12**(2): 227–228.
[PubMed Abstract](#) | [Publisher Full Text](#)
11. Omar C, Aldrich J, Gerkin R: **Collaborative infrastructure for test-driven scientific model validation**. In *Companion Proceedings of the 36th International Conference on Software Engineering*. ACM Press, 2014. 524–527.
[Publisher Full Text](#)
12. Shalloway A, Trott JR: **Design patterns explained: a new perspective on object-oriented design**. Pearson Education, 2004.
[Reference Source](#)
13. **Code coverage measurement for python**. 2015. Accessed: 2015-07-24.
[Reference Source](#)
14. **Coveralls-Test Coverage History and Statistics**. 2015. Accessed: 2015-07-24.
[Reference Source](#)
15. Booch G: **Object Oriented Analysis and Design with Applications**. Benjamin-Cummings, 1990.
[Reference Source](#)
16. Duvall PM, Matyas S, Glover A: **Continuous Integration: Improving Software Quality and Reducing Risk**. Addison-Wesley Professional, 2007.
[Reference Source](#)
17. **Travis CI - Test and Deploy Your Code with Confidence**. 2015. Accessed: 2015-07-24.
[Reference Source](#)
18. Gawande A: **The Checklist Manifesto: How to Get Things Right**. Picador. 2011.
[Reference Source](#)
19. Huang L, Kim R, Berry W: **Creating a culture of safety by using checklists**. *AORN J*. 2013; **97**(3): 365–368.
[PubMed Abstract](#) | [Publisher Full Text](#)
20. Weiser TG, Berry WR: **Perioperative checklist methodologies**. *Can J Anaesth*. 2013; **60**(2): 136–142.
[PubMed Abstract](#) | [Publisher Full Text](#)

Open Peer Review

Current Referee Status:



Version 1

Referee Report 31 August 2016

doi:[10.5256/f1000research.10031.r15595](https://doi.org/10.5256/f1000research.10031.r15595)



Andrew Davison

Unit of Neuroscience, Information et Complexité (UNIC), French National Center for Scientific Research, Gif-sur-Yvette, France

The article provides an introduction to automated software testing, its application to computational biology, and model validation as a form of testing, with examples taken from the OpenWorm project. The article is clearly written, and will be a helpful resource for computational biologists.

The article could be improved by a deeper discussion of some of the more difficult issues in the automation of model validation:

- what criteria to apply when transforming a numerical measure of closeness into a pass/fail?
- how to support the use of different criteria by different scientists, who might weigh the relative importance of particular validations very differently?
- how to handle contradictory experimental results?

I would also be interested to read a discussion of possible improvements to continuous integration dashboards in the context of continuous validation, e.g. tracking the evolution of numerical validation results across model versions.

I have read this submission. I believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.

Competing Interests: No competing interests were disclosed.

Referee Report 30 August 2016

doi:[10.5256/f1000research.10031.r15594](https://doi.org/10.5256/f1000research.10031.r15594)



Christian Roessert

Blue Brain Project, Ecole Polytechnique Fédérale de Lausanne (EPFL), Geneva, Switzerland

In this article the authors show how industrial practices of unit testing and test-driven development can be used and extended for computational modelling in biological sciences. The manuscript is well written and provides clear examples making it easy to understand the basic concepts.

I believe that establishing a culture of test-driven development in biological sciences is of great importance. However, in my view applying software engineering practices to computational modelling is

often not as easy as depicted by the authors. I have the following suggestions to improve the manuscript:

1. Judging the quality and validity of a computational model is a matter of scientific discussion and often cannot be easily reduced to a pass or fail decision in a model validation test. I would like to see a bit more detail on the transformation of the numeric score to a Boolean value used in the given ion channel test example but also on the general (statistical) concepts behind these decisions.
2. To iteratively improve a computational model, it is important to know not only if but also why a certain model fails or passes the model validation test. Since continuous integration systems are designed for simple verification tests: can the detailed results/figures and scores for each model validation test be shown directly on the CI dashboard? A discussion on the limits of current CI tools for biological modelling would be very helpful.
3. While the calculation of ion channel dynamics for a model validation test is computationally relatively cheap, computations become much more expensive once full detailed cell models or even networks have to be computed to validate against e.g. *in vivo* recordings. In these cases, the testing framework becomes much more sophisticated than "simple objects" and free services like Travis-CI will likely not be able to provide the required computational power. Is there a certain limit for your model validation test concept you would consider in the OpenWorm project and in general? Are there any ideas how to overcome these limitations? A discussion on the limits of the presented framework would be appreciated.

I have read this submission. I believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard, however I have significant reservations, as outlined above.

Competing Interests: No competing interests were disclosed.

Referee Report 24 August 2016

doi:[10.5256/f1000research.10031.r15597](https://doi.org/10.5256/f1000research.10031.r15597)



Robert Cannon

Textensor Limited, Edinburgh, UK

The authors present an interesting review of how they have applied traditional software testing methodologies to the OpenWorm project. It provides a nicely balanced perspective on a subject that often leads to strong opinions.

As they stay, the objective is to provide a focused case study of how testing is applied in the OpenWorm project. As such, the title is somewhat generic: I would suggest adding OpenWorm in there and possibly mentioning that this is a case study.

Probably the most novel part of the work is the incorporation of "Model Validation Tests" which serve to verify that the components, such as ion channel models, from which the model is built, behave in line with experimental data. The authors state that "Ultimately, a scientific model is just another piece of software—thus it can be validated with respect to a specification." In a sense this is true, but, as Ref 10 points out¹ the specification in the literature is often vague, incomplete or generally erroneous. The

SciUnit "judge" method appears to be the answer to this, replacing the usual software testing "assert" function. Presumably a lot of the subtlety of the approach, and indeed the scientific input whether a model is indeed a good match to experiments, is embedded in the implementation of the various "judge" methods. Although it is not essential for this paper it would be interesting to see a little more of how this is done in the OpenWorm project.

References

1. De Schutter E: The Dangers of Plug-and-Play Simulation Using Shared Models. *Neuroinformatics*. 2014. [Publisher Full Text](#)

I have read this submission. I believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.

Competing Interests: No competing interests were disclosed.
