# Systematic comparison and automated validation of detailed models of hippocampal neurons

Sára Sáray[1,2*], Christian A. Rössert[3], Shailesh Appukuttan[4], Rosanna Migliore[5], Paola Vitale[5], Carmen A. Lupascu[5], Luca L. Bologna[5], Werner Van Geit[3], Armando Romani[3], Andrew P. Davison[4], Eilif Muller[3], Tamás F. Freund[1,2], Szabolcs Káli[1,2*]

[1]Faculty of Information Technology and Bionics, Pázmány Péter Catholic University, Budapest, Hungary

[2] Institute of Experimental Medicine, Budapest, Hungary,

[3] Blue Brain Project, École Polytechnique Fédérale de Lausanne, Geneva, Switzerland,

[4]Paris-Saclay Institute of Neuroscience, Centre National de la Recherche Scientifique/ Université Paris-Saclay, Gif-sur-Yvette, France,

[5]Institute of Biophysics, National Research Council, Palermo, Italy.

[*] Corresponding author: saray.sara@koki.mta.hu (SS), kali@koki.hu (SK)

1

19  **Abstract**

20

21      Anatomically and biophysically detailed data-driven neuronal models can be useful

22  tools in understanding and predicting the behavior and function of neurons. Due to the

23  increasing availability of experimental data from anatomical and electrophysiological

24  measurements as well as the growing number of computational and software tools that enable

25  accurate neuronal modeling, there are now a large number of different models of many cell

26  types available in the literature. These models were usually built to capture a few important or

27  interesting properties of the given neuron type, and it is often unknown how they would behave

28  outside their original context. This limits the re-use and further development of the existing

29  models, and thus prevents the building of consensus "community models" that could capture

30  an increasing proportion of the electrophysiological properties of the given cell type. We

31  addressed this problem for the representative case of the CA1 pyramidal cell of the rat

32  hippocampus by developing an open-source Python test suite, which makes it possible to

33  automatically and systematically test the generalization properties of models by making

34  quantitative comparisons between the models and electrophysiological data. The tests cover

35  various aspects of somatic behavior, and signal propagation and integration in apical dendrites.

36  To demonstrate the utility of our approach, we applied our validation tests to compare the

37  behavior of several different hippocampal CA1 pyramidal cell models from the ModelDB

38  database against electrophysiological data available in the literature, and concluded that all of

39  these models perform well in some domains but badly in others. We also show how we

40  employed the test suite to aid the development of models within the European Human Brain

41  Project (HBP), and describe the integration of the tests into the validation framework developed

2

42  in the HBP, with the aim of facilitating more reproducible and transparent community model

43  building.

44

## Author summary

46

47  Anatomically and biophysically detailed neuronal models are useful tools in

48  neuroscience because they allow the prediction of the behavior and the function of the studied

49  cell type under circumstances that are hard to investigate experimentally. However, most

50  detailed biophysical models have been built to capture only a few properties of the real neuron,

51  and it is often unknown how they would behave under different circumstances, or whether they

52  can be used to successfully answer different scientific questions. To help the modeling

53  community develop neural models that generalize better, and make the process of model

54  building more reproducible and transparent, we developed a test suite that enables the

55  comparison of the behavior of models of neurons in the rat hippocampus and their evaluation

56  against experimental data. Applying our tests to several models available in the literature, we

57  show that each model is able to capture some of the important properties of the real neuron but

58  performs badly in other domains. We also use the test suite in the model development workflow

59  of the European Human Brain Project to aid the construction of better models of hippocampal

60  neurons and networks.

3

## Introduction

The construction and simulation of anatomically and biophysically detailed models is becoming a standard tool in neuroscience [1]. Such models, which typically employ the compartmental modeling approach and a Hodgkin-Huxley-type description of voltage-gated ion channels, are capable of providing fairly accurate models of single neurons [2–9] and (when complemented by appropriate models of synaptic interactions) even large-scale circuits [10–13]. However, building such detailed multi-compartmental models of neurons requires setting a large number of parameters (such as the densities of various ion channels in multiple neuronal compartments) that are often not directly constrained by the available experimental data. These parameters are typically tuned (either manually or using automated parameter-search methods [9,14–16]) until the simulated physiological behavior of the model matches some pre-defined set of experimental observations.

For an increasing number of cell types, the available experimental data already provide diverse constraints on the expected physiological behavior of the neuron under a variety of conditions. However, only a small subset of these constraints is typically taken into account when the models are developed, and it is often unknown, even by their developers, how these models would behave in other situations, outside their original context. This sparsity of information about the performance of detailed models might be one reason why model re-use in the community is relatively limited, and there are often a large number of different models of the same cell type available in the literature that were developed for different purposes. As an example, there are currently 129 different models related to the hippocampal CA1 pyramidal cell (PC) in the ModelDB database [17]. In addition, even when models are re-used, they are often altered to fit a different subset of the available experimental data, and they may lose their

4

85  ability to capture the behaviors that were used to constrain the original model. This phenomenon

86  (whereby introducing new features breaks previously correct behavior) is known as a

87  "regression" in software development, and is typically avoided by regularly applying a set of

88  tests that comprehensively verify the correct behavior of the software under various

89  circumstances. Such comprehensive checks are not routinely performed when neural models

90  are developed – and this may be one of the reasons why the development of consensus

91  (community) models, which would aim to capture a wide range of experimental observations

92  by integrating diverse efforts, has rarely been attempted in neuroscience.

93      A collaborative approach to modeling, and even a systematic comparison of existing

94  models built in different laboratories, requires the development of a comprehensive validation

95  suite, a set of automated tests that quantitatively compare various aspects of model behavior

96  with the corresponding experimental data. Such validation suites enable all modeling groups to

97  evaluate their existing and newly developed models according to common, standardized

98  criteria, thus facilitating model comparison and providing an objective measure of progress in

99  matching relevant experimental observations. Applying automated tests also allows researchers

100  to learn more about models published by other groups (beyond the results included in the

101  papers) with relatively little effort, thus facilitating optimal model re-use and co-operative

102  model development. Systematic testing of models during their development also helps avoid

103  regressions, aids the identification of problematic aspects of model behavior, and is thus

104  expected to lead to an increased efficiency in developing good models. The technical

105  framework for developing such test suites already exists  [18], and is currently used by several

106  groups to create a variety of tests for models of neural structure and function at different scales

107  [19–23]. In the current study, our goal was to develop a validation suite for the physiological

108  behavior of one of the most studied cell types of the mammalian brain, the pyramidal cell in

109  area CA1 of the rat hippocampus.

5

110    CA1 pyramidal neurons display a large repertoire of nonlinear responses in all of their

111    compartments (including the soma, axon, and various functionally distinct parts of the dendritic

112    tree), which are experimentally well-characterized. In particular, there are detailed quantitative

113    results available on the subthreshold and spiking voltage response to somatic current injections

114    [3,24]; on the properties of the action potentials back-propagating from the soma into the

115    dendrites [25–27], which is a basic measure of dendritic excitability; and on the characteristics

116    of the spread [28] and non-linear integration of synaptically evoked signals in the dendrites,

117    including the conditions necessary for the generation of dendritic spikes [29–32].

118    The test suite that we have developed allows the quantitative comparison of the behavior

119    of anatomically and biophysically detailed models of CA1 pyramidal neurons with

120    experimental data in all of these domains. In this paper, we first describe the implementation of

121    the HippoUnit validation suite. Next, we show how we used this test suite to systematically

122    compare existing models from six prominent publications from different laboratories. We then

123    show an example of how the tests have been applied to aid the development of new models in

124    the context of the European Human Brain Project (HBP). Finally, we describe the integration

125    of our test suite into the general validation framework developed in the HBP.

126

## Methods

### Implementation of HippoUnit

129

130    HippoUnit is a Python test suite based on the SciUnit [18] framework, which is a Python

131    package for testing scientific models, and during its implementation the NeuronUnit package

132    [19] was taken into account as an example of how to use the SciUnit framework for testing

133    neuronal models. In SciUnit tests usually four main classes are implemented: the test class, the

6

134    model class, the capabilities class and the score class. HippoUnit is built in a way that keeps

135    this structure. The key idea behind this structure is the decoupling of the model implementation

136    from the test implementation by defining standardized interfaces (capabilities) between them,

137    so that tests can easily be used with different models without being rewritten, and models can

138    easily be adapted to fit the framework.

139         Each test of HippoUnit is a separate Python class that, similarly to other SciUnit

140    packages, can run simulations on the models to generate model *predictions,* which can be

141    compared with experimental *observations* to yield the final score, provided that the model has

142    the required capabilities implemented to mimic the appropriate experimental protocol and

143    produce the same type of measurable output. All measured or calculated data that contribute to

144    the final score are saved in JSON or pickle files (or, in many cases, in both types of files). JSON

145    files are human readable, and can be easily loaded into Python dictionaries. Data with a more

146    complex structure are saved into pickle files. This makes it possible to easily write and read the

147    data (for further processing or analysis) without changing its Python structure, no matter what

148    type of object or variable it is.

149         Similarly to many of the existing SciUnit packages the implementations of specific

150    models are not part of the HippoUnit package itself. Instead, HippoUnit contains a general

151    `ModelLoader` class. This class is implemented in a way that it is able to load and deal with

152    most types of models defined in the HOC language of the NEURON simulator (either as

153    standalone HOC models or as HOC templates) [33]. It implements all model-related methods

154    (capabilities) that are needed to simulate these kinds of neural models in order to generate the

155    prediction without any further coding required from the user.

156         For the smooth validation of the models developed using parameter optimization within

157    the HBP there is a child class of the `ModelLoader` available in HippoUnit that is called

158    `ModelLoader_BPO`. This class inherits most of the functions (especially the capability

7

159  functions) from the `ModelLoader` class, but it implements additional functions that are able to

160  automatically deal with the specific way in which information is represented and stored in these

161  optimized models. The role of these functions is to gather all the information from the metadata

162  and configuration files of the models that are needed to set the parameters required to load the

163  models and run the simulations on them (such as path to the model files, name of the model

164  template or the simulation temperature (the `celsius` variable of Neuron)). This enables the

165  validation of these models without any manual intervention needed from the user. The section

166  lists required by the tests of HippoUnit are also created automatically using the morphology

167  files of these models (for details see the "Classify apical sections of pyramidal cells"

168  subsection). For neural models developed using other software and methods, the user needs to

169  implement the capabilities through which the tests of HippoUnit perform the simulations and

170  recordings on the model.

171      The capabilities are the interface between the tests and the models. The `ModelLoader`

172  class inherits from the capabilities and must implement the methods of the capability. The test

173  can only be run on a model if the necessary capability methods are implemented in the

174  `ModelLoader`. All communication between the test and the model happens through the

175  capabilities.

176      The methods of the score classes perform the quantitative comparison between the

177  *prediction* and the *observation*, and return the score object containing the final score and some

178  related data, such as the paths to the saved figure and data (JSON) files and the prediction and

179  observation data. Although SciUnit and NeuronUnit have a number of different score types

180  implemented, HippoUnit has its own scores, which better fit its tests and the observations

181  belonging to them. For simplicity, we refer to the discrepancy between the target experimental

182  data (*observation*) and the models' behavior (*prediction*) with respect to a studied feature using

183  the term feature error. In most cases, when the basic statistics (mean and standard deviation) of

8

184    the experimental features (typically measured in several different cells of the same cell type)

185    are available, feature errors are computed as the absolute difference between the feature value

186    of the model and the experimental mean feature value, divided by the experimental standard

187    deviation (Z-score) [34]. The final score of a given test achieved by a given model is given by

188    the average (or, in some cases, the sum) of the feature error scores for all the features evaluated

189    by the test.

190

## Implementation of the tests of HippoUnit

### The Somatic Features Test

193

194        The Somatic Features Test uses the Electrophys Feature Extraction Library (eFEL) [35]

195    to extract and evaluate the values of both subthreshold and suprathreshold (spiking) features

196    from voltage traces that represent the response of the model to somatic current injections of

197    different positive (depolarizing) and negative (hyperpolarizing) current amplitudes. Spiking

198    features describe action potential shape (like AP width, AP rise/fall rate, AP amplitude, etc.)

199    and timing (frequency, inter-spike intervals, time to first/last spike, etc.), while some passive

200    features (such as the voltage base or the steady state voltage), and subthreshold features for

201    negative current stimuli (voltage deflection, sag amplitude, etc.) are also examined.

202        In this test step currents of varying amplitudes are injected into the soma of the model

203    and the voltage response is recorded. The simulation protocol is set according to an input

204    configuration JSON file, which contains all the current amplitudes, the delay and the duration

205    of the stimuli, and the stimulation and recording positions. Simulations using different current

206    amplitudes are run in parallel if this is supported by the computing environment.

9

207   As the voltage responses of neurons to somatic current injections can strongly depend

208   on the experimental method, and especially on the type of electrode used, target values for these

209   features were extracted from two different datasets. One dataset was obtained from sharp

210   electrode recordings from adult rat CA1 neurons (sharp electrode data set) [3], and the other

211   dataset is from patch clamp recordings in rat CA1 pyramidal cells (data provided by Judit

212   Makara (patch clamp dataset)). For both of these datasets we had access to the recorded voltage

213   traces from multiple neurons, which made it possible to perform our own feature extraction

214   using eFEL. This ensures that the features are interpreted and calculated the same way for both

215   the experimental data and the models' voltage response during the simulation. Furthermore, it

216   allows a more thorough comparison against a large number of features extracted from

217   experimental recordings yielded using the exact same protocol, which is unlikely to be found

218   in any paper of the available literature. However, to see how representative these datasets are

219   of the literature as a whole we first compared some of the features extracted from these datasets

220   to data available on Neuroelectro.org [36] and on Hippocampome.org [37]. The features we

221   compared were the following: resting potential, voltage threshold, after-hyperpolarization

222   (AHP) amplitudes (fast, slow), action potential width and sag ratio. Although these databases

223   have mean and standard deviation values for these features that are calculated from

224   measurements using different methods, protocols and from different animals, we found that

225   most of the feature values for our two experimental datasets fall into the ranges declared as

226   typical for CA1 PCs in the online databases. The only conspicuous exception is the fast AHP

227   amplitude of the patch clamp dataset used in this study, which is $1.7 \pm 1.5$ mV, while the

228   databases cite values between 6.8 and 11.64 mV. This deviation could possibly stem from a

229   difference in the way that the fast AHP is measured.

230   We also performed a more specific review of the relevant literature to compare the most

231   important somatic features of the patch clamp dataset to results from available patch clamp

10

232    recordings. In the literature the somatic AP voltage threshold of CA1 pyramidal cells is between

233    -46 and -53 mV [38–41]. The same feature (*AP_begin_voltage*) extracted from out patch clamp

234    dataset falls into this range (-51.13±0.97 mV (0.15 nA current step), -50.14±1.97 mV (0.2 nA

235    current step); -49.36±2.02  mV (0.25 nA current step)). The AP amplitude falls into a relatively

236    broad range between 71 and 112 mV according to the literature [38,41,42]. As most of these

237    sources calculate the amplitude between the peak and the voltage base, we take the eFEL feature

238    *AP_amplitude_from_voltagebase* into account here. The value of it in our patch clamp dataset

239    is similar to the experimental observations (98.36±5.82 mV (0.15 nA current step), 96.83±5.66

240    mV (0.2 nA current step), 95.99±5.22 mV (0.25 nA current step)). The AP width measured at

241    half amplitude ranges from 0.8 to 1.29 ms in the literature [38,40–42]. The value of this feature

242    is near the upper end of this range in our patch clamp dataset (1.23±0.096 ms (0.15 nA step

243    current); 1.25±0.11 ms (0.2 nA step current); 1.32±0.086 ms (0.25 nA step current)). Regarding

244    the features extracted from response to hyperpolarizing currents the sag ratio can be compared.

245    The values from literature are: 0.84±0.02 [42] and 0.83±0.01 [43], while the values extracted

246    from our patch clamp dataset are quite similar (0.79±0.023 (-0.05 nA step current); 0.81±0.03

247    (-0.1 nA step current); 0.81±0.027 (-0.15 nA step current); 0.81±0.03 (-0.2 nA step current);

248    0.80±0.03 (-0.25 nA step current). We conclude that the patch clamp dataset is in good

249    agreement with experimental observations available in the literature, and will be used as a

250    representative example in this study.

251        The *observation* data are loaded from a JSON file of a given format which contains the

252    names of the features to be evaluated, the current amplitude for which the given feature is

253    evaluated and the corresponding experimental mean and standard deviation values. Setting the

254    `specify_data_set` parameter it can be ensured that the test results against different

255    experimental data sets are saved into different folders.

11

256    For certain features eFEL returns a vector as a result; in these cases, the feature value

257    used by HippoUnit is the average of the elements of the vector. These are typically spiking

258    features for which eFEL extracts a value corresponding to each spike fired. For features that

259    use the 'AP_begin_time' or 'AP_begin_voltage' feature values for further calculations, we

260    exclude the first element of the vector output before averaging because we discovered that these

261    features are often incorrectly detected for the first action potential of a train.

262    The score class of this test returns as the final score the average of *Z-scores* for the

263    evaluated eFEL features achieved by the model. Those features that could not be evaluated

264    (e.g., spiking features from voltage responses without any spikes) are listed in a log file to

265    inform the user, and the number of successfully evaluated features out of the number of features

266    attempted to be evaluated is also reported.

267

268    **The Depolarization Block Test**

269

270    This test aims to determine whether the model enters depolarization block in response

271    to a prolonged, high intensity somatic current stimulus. For CA1 pyramidal cells, the test relies

272    on experimental data from Bianchi et al. [24]. According to these data, CA1 PCs respond to

273    somatic current injections of increasing intensity with an increasing number of action potentials

274    until a certain threshold current intensity is reached. For current intensities higher than the

275    threshold, the cell does not fire over the whole period of the stimulus; instead, firing stops after

276    some action potentials, and the membrane potential is sustained at some constant depolarized

277    level for the rest of the stimulus. This phenomenon is termed depolarization block [24].

278    This test uses the same capability class as the Somatic Features Test for injecting current

279    and recording the somatic membrane potential (see the description above). Using this

12

280    capability, the model is stimulated with 1000 ms long square current pulses increasing in

281    amplitude from 0 to 1.6 nA in 0.05 nA steps, analogous to the experimental protocol. The

282    stimuli of different amplitudes are run in parallel. Somatic spikes are detected and counted using

283    eFEL [35].

284         From the somatic voltage responses of the model, the following features are evaluated.

285    $I_{th}$ is the threshold current to reach depolarization block; experimentally, this is both the

286    amplitude of the current injection at which the cell exhibits the maximum number of spikes,

287    and the highest stimulus amplitude that does not elicit depolarization block. In the test two

288    separate features are evaluated for the model and compared to the experimental $I_{th}$: the current

289    intensity for which the model fires the maximum number of action potentials (*I_maxNumAP*),

290    and the current intensity one step before the model enters depolarization block

291    (*I_below_depol_block*). If these two feature values are not equal, a penalty is added to the score.

292    The model is defined to exhibit depolarization block if *I_maxNumAP* is not the highest

293    amplitude tested, and if there exists a current intensity higher than *I_maxNumAP*, for which the

294    model does not fire action potentials during the last 100 ms of its voltage response.

295         In the experiment the $V_{eq}$ feature is extracted from the voltage response of the pyramidal

296    cells to the current injection one step above *$I_{th}$ (or I_max_num_AP* in the test). Both in the

297    experiment and in this test this is calculated as the mean voltage over the last 100 ms of the

298    voltage trace. However, in the test, before calculating this value it is examined whether there

299    are any action potentials during this period. The presence of spikes here means that the model

300    did not enter depolarization block prior to this period. In these cases the test iterates further on

301    the voltage traces corresponding to larger current steps to find if there is any where the model

302    actually entered depolarization block; if an appropriate trace is found, the value of $V_{eq}$ is

303    extracted there. This trace is the response to the current intensity one step above

304    *I_below_depol_block.*

                13

305    If the model does not enter depolarization block, a penalty is applied, and the final score

306    gets the value of 100. Otherwise, the final score achieved by the model on this test is the average

307    of the error scores (Z-scores) for the features described above, plus an additional penalty if

308    *I_maxNumAP* and *I_below_depol_block* differ. This penalty is 200 times the difference

309    between the two current amplitude values (in pA – which in this case is 10 times the number of

310    examined steps between them).

311

312    **The Back-propagating AP Test**

313

314    This test evaluates the strength of action potential back-propagation in the apical trunk

315    at locations of different distances from the soma. The observation data for this test were yielded

316    by the digitization of Figure 1B of [26], using the DigitizeIt software [44]. The values were

317    then averaged over distances of 50, 150, 250, 350 $\pm$ 20 $\mu$m from the soma to get the mean and

318    standard deviation of the features. The features tested here are the amplitudes of the first and

319    last action potentials of a 15 Hz spike train, measured at the 4 different dendritic locations.

320    The test automatically finds current amplitudes for which the soma fires, on average,

321    between 10-20 Hz and chooses the amplitude that leads to firing nearest to 15 Hz. For this task,

322    the following algorithm was implemented. Increasing current step stimuli of 0.0 - 1.0 nA

323    amplitude with a step size of 0.1 nA are applied to the model and the number of spikes is

324    counted for each resulting voltage trace. If spontaneous spiking occurs (i.e., if there are spikes

325    even when no current is injected) or if the spiking rate does not reach 10 Hz even for the highest

326    amplitude, the test quits with an error message. Otherwise the amplitudes for which the soma

327    fires between 10 and 20 Hz are appended to a list and (if the list is not empty) the one providing

328    the spiking rate nearest to 15 Hz is chosen. If the list is empty because the spiking rate is smaller

14

329 than 10 Hz for a step amplitude but higher than 20 Hz for the next step, a binary search method

330 is used to find an appropriate amplitude in this range.

331 This test uses a trunk section list (or generates one if the `find_section_lists`

332 variable of the `ModelLoader` is set to True – see the section 'Classifying the apical sections of

333 pyramidal cells' below) to automatically find the dendritic locations for the measurements. The

334 desired distances of the locations from the soma and the distance tolerance are read from the

335 input configuration file, and must agree with the distances and the tolerance over which the

336 experimental data were averaged. All the trunk dendritic segments whose distance from the

337 soma falls into one of the distance ranges are selected. The locations and also their distances

338 are then returned in separate dictionaries.

339 Then the soma is stimulated with a current injection of the previously chosen amplitude

340 and the voltage response of the soma and the selected dendritic locations are recorded and

341 returned.

342 The test implements its own function to extract the amplitudes of back-propagating

343 action potentials, but the method is based on eFEL features. This is needed because eFEL's

344 spike detection is based on a given threshold value for spike initiation, which may not be

345 reached by the back-propagating signal at more distant regions. First the maximum

346 depolarization of the first and the last action potentials are calculated. This is the maximum

347 value of the voltage trace in a time interval around the somatic action potential, based on the

348 start time of the spike (using the AP_begin_time feature of eFEL) and the inter-spike interval

349 to the next spike recorded at the soma. Then the amplitudes are calculated as the difference

350 between this maximum value and the voltage at the begin time of the spike (on the soma) minus

351 1 ms (which is early enough not to include the rising phase of the spike, and late enough in the

352 case of the last action potential not to include the afterhyperpolarization of the previous spike).

15

353      To calculate the feature error scores the amplitude values are first averaged over the

354      distance ranges to be compared to the experimental data and get the feature Z-scores. The final

355      score here is the average of the Z-scores achieved for the features of first and last action

356      potential amplitudes at different dendritic distances. In the result it is also stated whether the

357      model is more like a strongly or a weakly propagating cell in the experiment, where they found

358      examples of both types [26].

359

## The PSP Attenuation Test

361

362      The PSP Attenuation test evaluates how much the post-synaptic potential attenuates as it

363      propagates from different dendritic locations to the soma in CA1 pyramidal cell models. The

364      *observation* data for this test were yielded by the digitization of Figure 1E and Figure 2B of

365      Magee and Cook, 2000 [28] using the DigitizeIt software [44]. The somatic and dendritic

366      depolarization values were then averaged over distances of 100, 200, 300 ± 50 μm from the

367      soma and the soma/dendrite attenuation was calculated to get the mean and standard deviation

368      of the attenuation features at the three different input distances.

369      In this test the apical trunk receives excitatory post-synaptic current (EPSC)-shaped

370      current stimuli at locations of different distances from the soma. The maximum depolarization

371      caused by the input is extracted at the soma and divided by the maximum depolarization at the

372      location of the stimulus to get the soma/dendrite attenuation values that are then averaged in

373      distance ranges of 100, 200, 300 ± 50 μm and compared to the experimental data. The distances

374      and tolerance are defined in the configuration file and must agree with how the *observation* data

375      were generated.

16

376       The test uses a trunk section list, which needs to be specified in the NEURON HOC

377       model (or the test generates one if the `find_section_lists` variable of the `ModelLoader`

378       is set to True – see the section 'Classify apical sections of pyramidal cells' below) to find the

379       dendritic locations to be stimulated. Randomly selected dendritic locations are used because the

380       distance ranges that are evaluated cover almost the whole length of the trunk of a pyramidal

381       cell. The probability of selecting a given dendritic segment is set to be proportional to its length.

382       The number of dendritic segments examined can be chosen by the user by setting the

383       `num_of_dend_locations` argument of the test. The random seed (also an argument of the

384       test) must be kept constant to make the selection reproducible. If a given segment is selected

385       multiple times (or it is closer than 50 µm or further than 350 µm), a new random number is

386       generated. If the number of locations to be selected is more than the number of trunk segments

387       available in the model, all the segments are selected.

388       The *Exp2Syn* synaptic model of NEURON with a previously calculated weight is used to

389       stimulate the dendrite. The desired EPSC amplitude and time constants are given in the input

390       configuration file according to the experimental protocol. To get the proper synaptic weight,

391       first the stimulus is run with weight = 0. The last 10% of the trace is averaged to get the resting

392       membrane potential (Vm). Then the synaptic weight required to induce EPSCs with the

393       experimentally determined amplitude is calculated according to Equation 1:

394       (1) weight = - EPSC_amp / Vm

395       where EPSC_amp is read from the `config` dictionary, and the synaptic reversal potential is

396       assumed to be 0 mV.

397       To get the somatic and dendritic maximum depolarization from the voltage traces, the

398       baseline trace (weight = 0) is subtracted from the trace recorded in the presence of the input. To

399       get the attenuation ratio the maximum value of the somatic depolarization is divided by the

400       maximum value of the dendritic depolarization.

17

401    To calculate the feature error scores the soma/dendrite attenuation values are first

402    averaged over the distance ranges to be compared to the experimental data to get the feature Z-

403    scores. The final score is the average of the feature error scores calculated at the different dendritic

404    locations.

405

406    **The Oblique Integration Test**

407

408    This test evaluates the signal integration properties of radial oblique dendrites,

409    determined by providing an increasing number of synchronous (0.1 ms between inputs) or

410    asynchronous (2 ms between inputs) clustered synaptic inputs. The experimental mean and

411    standard error (SE) of the features examined are available in the paper of Losonczy and Magee

412    [32] and are read from a JSON file into the *observation* dictionary of the test. The SE values

413    are then converted to standard deviation values. The following features are tested: voltage

414    threshold for dendritic spike initiation (defined as the expected somatic depolarization at which a

415    step-like increase in peak dV/dt occurs); proximal threshold (defined the same way as above, but

416    including only those results in the statistics where the proximal part of the examined dendrite was

417    stimulated); distal threshold; degree of nonlinearity at threshold; suprathreshold degree of

418    nonlinearity; peak derivative of somatic voltage at threshold; peak amplitude of somatic EPSP; time

419    to peak of somatic EPSP; degree of nonlinearity in the case of asynchronous inputs.

420    The test automatically selects a list of oblique dendrites that meet the criteria of the

421    experimental protocol, based on a section list containing the oblique dendritic sections (this can

422    either be provided by the HOC model, or generated automatically if the

423    `find_section_lists` variable of the `ModelLoader` is set to True – see the section 'Classify

424    apical sections of pyramidal cells' below). For each selected oblique dendrite a proximal and a

18

425    distal location is examined. The criteria for the selection of dendrites, which were also applied

426    in the experiments, are the following. The selected oblique dendrites should be terminal

427    dendrites (they have no child sections) and they should be at most 120 μm from the soma. This

428    latter criterion can be changed by the user by changing the value of the `ModelLoader`'s

429    *max_dist_from_soma* variable, and it can also increase automatically if needed. In particular,

430    if no appropriate oblique is found up to the upper bound provided, the distance is increased

431    iteratively by 15 μm, but not further than 190 μm.

432        Then an increasing number of synaptic inputs are activated at the selected dendritic

433    locations separately, while recording the local and somatic voltage response. HippoUnit

434    provides a default synapse model to be used in the `ObliqueIntegrationTest`. If the

435    *AMPA_name*, and *NMDA_name* variables are not set by the user, the default synapse is used. In

436    this case the AMPA component of the synapse is given by the built-in `Exp2Syn` synapse of

437    NEURON, while the NMDA component is defined in an NMODL (.mod) file which is part of

438    the HippoUnit package. This NMDA receptor model uses a Jahr-Stevens voltage dependence

439    [45] and rise and decay time constants of 3.3 and 102.38 ms, respectively. The time constant

440    values used here are temperature- (Q10-) corrected values from [41]. Q10 values for the rise

441    and decay time constants were 2.2 [46] and 1.7 [47], respectively. The model's own AMPA

442    and NMDA receptor models can also be used in this test if their NMODL files are available

443    and compiled among the other mechanisms of the model. In this case the `AMPA_name`, and

444    `NMDA_name` variables need to be provided by the user. The time constants of the built-in

445    Exp2Syn AMPA component and the AMPA/NMDA ratio can be adjusted by the user by setting

446    the *AMPA_tau1, AMPA_tau2* and `AMPA_NMDA_ratio` parameter of the `ModelLoader`. The

447    default AMPA/NMDA ratio is 2.0 from [41], and the default `AMPA_tau1` and `AMPA_tau2` are

448    0.1 ms and 2.0 ms, respectively [28,29].

19

449     To test the Poirazi et al. 2003 model using its own receptor models, we also had to

450     implement a modified version of the synapse functions of the `ModelLoader` that can deal with

451     the different (pointer-based) implementation of synaptic activation in this model. For this

452     purpose, a child class was implemented that inherits from the `ModelLoader` class. This

453     modified version is not part of the official HippoUnit version, but is available here:

454     https://github.com/KaliLab/HippoUnit_demo/blob/master/ModelLoader_Poirazi_2003_CA1.

455     py.

456     The synaptic weights for each selected dendritic location are automatically adjusted by

457     the test using a binary search algorithm so that the threshold for dendritic spike generation is 5

458     synchronous inputs – which was the average number of inputs that had to be activated by

459     glutamate uncaging to evoke a dendritic spike in the experiments [32]. This search runs in

460     parallel for all selected dendritic locations. The search interval of the binary search and the

461     initial step size of the searching range can be adjusted by the user through the `c_minmax` and

462     `c_step_start` variables of the `ModelLoader`. During the iterations of the algorithm the step

463     size may decrease if needed; a lower threshold for the step size (`c_step_stop` variable of the

464     `ModelLoader`) must be set to avoid infinite looping. Those dendritic locations where this first

465     dendritic spike generates a somatic action potential, or where no dendritic spike can be evoked, are

466     excluded from further analysis. To let the user know, this information is displayed on the output

467     and also printed into the log file saved by the test. Most of the features above are extracted at the

468     threshold input level (5 inputs).

469     The final score of this test is the average of the feature error scores achieved by the model

470     for the different features; however, a T-test analysis is also available as a separate score type for

471     this test.

472

20

**Parallel computing**

Most of the tests of HippoUnit require multiple simulations of the same model, either using stimuli of different intensities or at different locations in the cell. To run these simulations in parallel and save time, the Python `multiprocessing.Pool` module is used. The size of the pool can be set by the user. Moreover, all NEURON simulations are performed in multiprocessing pools to ensure that they run independently of each other, and to make it easy to erase the models after the process has finished. This is especially important in the case of HOC templates in order to avoid previously loaded templates running in the background and the occurrence of 'Template cannot be redefined' errors when the same model template is loaded again.
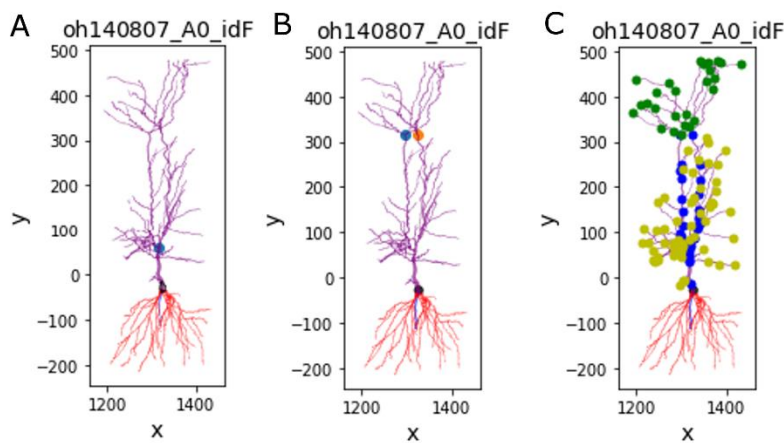
**Classifying the apical sections of pyramidal cells**

Some of the validation tests of HippoUnit require lists of sections belonging to the different dendritic types of the apical tree (main apical trunk, apical tuft dendrites, radial oblique dendrites). To classify the dendrites NeuroM [48] is used as a base package. NeuroM contains a script that, starting from the tuft (uppermost dendritic branches in Fig 1) endpoints, iterates down the tree to find a single common ancestor. This is considered as the apical point. The apical point is the upper end of the main apical dendrite (trunk), from where the tuft region arises. Every dendrite branching from the trunk below this point is considered an oblique dendrite.

However, there are many CA1 pyramidal cell morphologies where the trunk bifurcates close to the soma to form two or even more branches. In these cases the method described above

21

497   finds this proximal bifurcation point as the apical point (see Fig 1A). To overcome this issue,

498   we worked out and implemented a method to find multiple apical points by iterating the

499   function provided by NeuroM. In particular, if the initial apical point is closer to the soma than

500   a pre-defined threshold, the function is run again on subtrees of the apical tree where the root

501   node of the subtree is the previously found apical point, to find apical points on those subtrees

502   (see Fig 1B). When (possibly after multiple iterations) apical points that are far enough from

503   the soma are found, NeuroM is used to iterate down from them on the parent sections, which

504   will be the trunk sections (blue dots in Fig 1C). Iterating up, the tuft sections are found (green

505   dots in Fig 1C), and the other descendants of the trunk sections are considered to be oblique

506   dendrites (yellow dots in Fig 1C). Once all the sections are classified, their NeuroM coordinates

507   are converted to NEURON section information for further use.

508



509

510   Fig 1: Classifying the apical dendrites of pyramidal cells. Morphological reconstruction made within the HBP at

511   University College London (UCL). (A) The original method of NeuroM finds a single apical point which is actually

512   a bifurcation of the trunk. (B) Further developing the method, multiple apical points can be found. (C)The apical

513   dendritic sections are classified. Blue: trunk, yellow: oblique dendrites, green: tuft sections.

514

22

515    We note that this function can only be used for hoc models that load their morphologies

516    from a separate morphology file (e.g., ASC, SWC) as NeuroM can only deal with morphologies

517    provided in these standard formats. For models with NEURON morphologies implemented

518    directly in the hoc language, the SectionLists required by a given test should be implemented

519    within the model.

520

521    **Models from literature**

522

523    In this paper we demonstrate the utility of the HippoUnit validation test suite by

524    applying its tests to validate and compare the behavior of several different detailed hippocampal

525    CA1 pyramidal cell models available on ModelDB [17]. For this initial comparison we chose

526    models published by several modeling groups worldwide that were originally developed for

527    various purposes.

528    The Golding et al., 2001 model [26] (ModelDB accession number: 64167) was

529    developed to show the dichotomy of the back-propagation efficacy and the amplitudes of the

530    back-propagating action potentials at distal trunk regions in CA1 pyramidal cells and to make

531    predictions on the possible causes of this behavior. It contains only the most important ion

532    channels (Na, $K_{DR}$, $K_A$) needed to reproduce the generation and propagation of action

533    potentials. Here we tested three different versions of the model: the ones corresponding to

534    Figure 8A, Figure 8B and Figure 9B of the paper [26].

535    The Katz et al., 2009 model [49] (ModelDB accession number: 127351) is based on the

536    Golding et al. 2001 model and was built to investigate the functional consequences of the

537    distribution of strength and density of synapses on the apical dendrites that they observed

538    experimentally, for the mode of dendritic integration.

23

539    The Migliore et al., 2011 model [50] (ModelDB accession number: 138205) was used

540    to study schizophrenic behavior. It is based on earlier models of the same modeling group,

541    which were used to investigate the initiation and propagation of action potentials in oblique

542    dendrites, and have been validated against different electrophysiological data.

543    The Poirazi et al., 2003 model [6,51] (ModelDB accession number: 20212) was

544    designed to clarify the issues about the integrative properties of thin apical dendrites that may

545    arise from the different and sometimes conflicting interpretations of available experimental

546    data. This is a quite complex model in the sense that it contains a large number of different

547    types of ion channels, whose properties were adjusted to fit in vitro experimental data, and it

548    also contains four types of synaptic receptors.

549    The Bianchi et al., 2012 model [24] (ModelDB accession number: 143719) was

550    designed to investigate the mechanisms behind depolarization block observed experimentally

551    in the somatic spiking behavior of CA1 pyramidal cells. It was developed by combining and

552    modifying the Shah et al., 2008 [52] and the Poirazi et al. 2003 models [6,51]. The former of

553    these was developed to show the significance of axonal M-type potassium channels.

554    The Gómez González et al., 2011 [53]  model (ModelDB accession number: 144450) is

555    based on the Poirazi et al. 2003 model and it was modified to replicate the experimental data of

556    [32] on the nonlinear signal integration of radial oblique dendrites when the inputs arrive in a

557    short time window. The model was adjusted to five different detailed morphologies.

558    Models from literature that are published on ModelDB typically implement their own

559    simulations and plots to make it easier for users and readers to reproduce and visualize the

560    results shown in the corresponding paper. Therefore, to be able to test the models described

561    above using our test suite, we needed to create standalone versions of them. These standalone

562    versions do not display any GUI, or contain any built-in simulations and run-time

563    modifications, but otherwise their behavior should be identical to the published version of the

24

564    models. We also added section lists of the radial oblique and the trunk dendritic sections to

565    those models where this was not done yet, as some of the tests require these lists. To ensure that

566    the standalone versions have the same properties as the original models, we checked their

567    parameters after running their built-in simulations (in case including any run-time

568    modifications), and made sure they match the parameters of the standalone version. We also

569    asked the developers of the models to check the standalone versions and give us feedback;

570    however, we received substantial feedback only from the developers of the Migliore et al. 2011

571    and the Bianchi et al., 2012 models, which were positive. The modified models used for running

572    validation tests are available in this GitHub repository:

573    https://github.com/KaliLab/HippoUnit_demo.

## Results

### The HippoUnit validation suite

576

577    HippoUnit (https://github.com/KaliLab/hippounit ) is an open source test suite for the

578    automatic and quantitative evaluation and validation of the behavior of neural single cell

579    models. The tests of HippoUnit automatically perform simulations that mimic common

580    electrophysiological protocols on neuronal models to compare their behavior with quantitative

581    experimental data using various feature-based error functions. Current validation tests cover

582    somatic (subthreshold and spiking) behavior as well as signal propagation and integration in

583    the dendrites. The tests were developed using data and models for rat hippocampal CA1

584    pyramidal cells. However, most of the tests are directly applicable to or can be adapted for other

585    cell types if the necessary experimental data are available; examples of this will be presented

586    in later sections.

25

587     HippoUnit is implemented in the Python programming language, and is based on the

588     SciUnit [18] framework for testing scientific models. The current version of HippoUnit is

589     capable of handling single cell models implemented in the NEURON simulator. As a result, by

590     adapting and using the example Jupyter notebooks described in S1 Appendix, the tests of

591     HippoUnit can be run on neural models that are built in the NEURON simulator software,

592     without any further coding required from the user. In principle, neural models developed using

593     other software tools can also be tested by HippoUnit; however, this requires the re-

594     implementation by the user of the interface functions that allow HippoUnit to run the necessary

595     simulations and record their output (see the Methods section for more details).

596     In the current tests of HippoUnit, once all the necessary simulations have been

597     performed and the responses of the model have been recorded, electrophysiological features are

598     extracted from the voltage traces, and the discrepancy between the model's behavior and the

599     experiment is computed by comparing the feature values with those extracted from the

600     experimental data (see Methods). For simplicity, we refer to the result of this comparison as the

601     feature error; however, we note that there are many possible sources of such discrepancy

602     including, among others, experimental artefacts and noise, shortcomings of the models, and

603     differences between the conditions assumed by the models and those in the actual experiments

604     (see the Discussion for more details). The final score of a given test achieved by a given model

605     is given by the average (or, in some cases, the sum) of the feature error scores for all the features

606     evaluated by the test.

607     Besides the final score, which is the basic output of all the tests, the tests of HippoUnit

608     typically provide a number of other useful outputs (see Methods), including figures that

609     visualize the model's behavior through traces and plot the feature and error values compared to

610     the experimental data. It is always strongly recommended to look at the traces and other figures

611     to get a fuller picture of the model's response to the stimuli, which helps with the correct

26

612    interpretation of validation results. Such closer inspection also makes it possible to detect

613    possible test failures, when the extraction of certain features does not work correctly for a given

614    model.

615    HippoUnit can also take advantage of the parallel execution capabilities of modern

616    computers. When tests require multiple simulations of the same model using different settings

617    (e.g., different stimulation intensities or different stimulus locations in the cell), these

618    simulations are run in parallel, which can make the validation process substantially faster,

619    depending on the available computing resources.

620    One convenient way of running a test on a model is to use an interactive computational

621    notebook, such as the Jupyter Notebook [54], which enables the combination of program codes

622    to be run (we used Python code to access the functionality of HippoUnit), the resulting outputs

623    (e.g. figures, tables, text) and commentary or explanatory text in a single document. Therefore,

624    we demonstrate the usage of HippoUnit through this method (See S1 Appendix and

625    https://github.com/KaliLab/HippoUnit_demo).

626
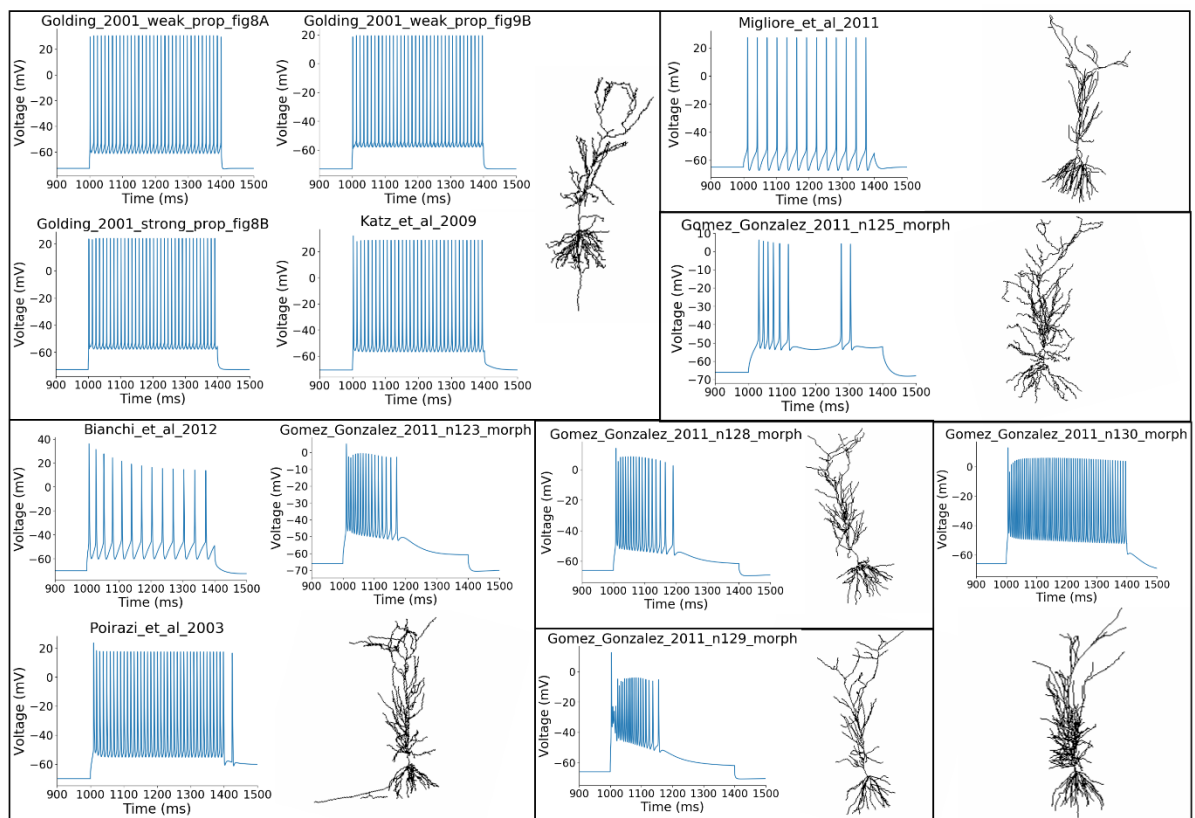
## Comparison of the behavior of rat hippocampal CA1 pyramidal cell models selected from the literature

629

630    We selected six different publications containing models of hippocampal CA1

631    pyramidal cells whose implementations for the NEURON simulator were available in the

632    ModelDB database. Our aim was to compare the behavior of every model to the experimental

633    target data using the tests of HippoUnit, which also allowed us to compare the models to each

634    other, and to test their generalization performance in paradigms that they were not originally

635    designed to capture. These models differ in their complexity regarding the number and types of

27

636    ion channels that they contain, and they were built for different purposes (see the Methods

637    section for more details on the models). A common property of these models is that their

638    parameters were set using manual procedures with the aim of reproducing the behavior of real

639    CA1 PCs in one or a few specific paradigms. As some of them were built by modifying and

640    further developing previous models, these share the same morphology (see Fig. 2). On the other

641    hand, the model of Gómez González et al. 2011 was adjusted to 5 different morphologies, which

642    were all tested. In the case of the Golding et al. 2001 model, we tested three different versions

643    (shown in Figures 8A, 8B and 9A of the corresponding paper [26]) that differ in the distribution

644    of the sodium and the A-type potassium channels, and therefore the back-propagation efficacy

645    of the action potentials. The morphologies and characteristic voltage responses of all the models

646    used in this comparison are displayed in Fig 2.

647



648

28

649 Fig 2: The morphologies of the different models tested and their voltage responses to a 400 ms step current

650 injection of 0.6 nA amplitude.

651

652       Running the tests of HippoUnit on these models we took into account the original

653 settings of the simulations of the models, and set the `v_init` (the initial voltage when the

654 simulation starts), and the `celsius` (the temperature at which the simulation is done) variables

655 accordingly. For the Bianchi et al 2012 model we used variable time step integration during all

656 the simulations, as it was done in the original modeling study. For the other models a fixed time

657 step were used (dt=0.025 ms).
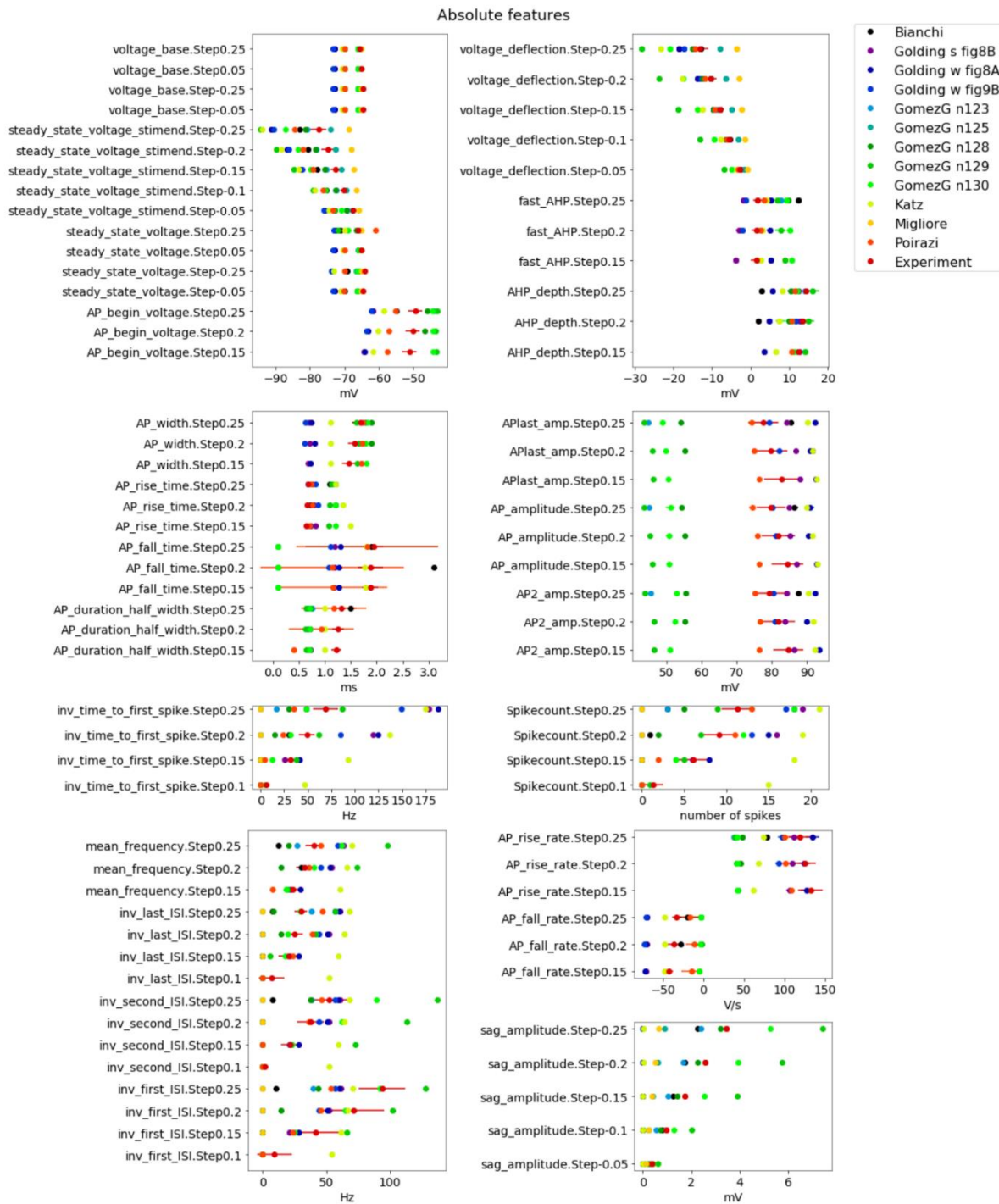
658

659 **Somatic Features Test**

660

661       Using the Somatic Features Test of HippoUnit, we compared the behavior of the models

662 to both patch clamp recordings (patch clamp dataset) and sharp electrode recordings (sharp

663 electrode dataset). To see how representative these datasets are of the literature as a whole we

664 first compared some of the features extracted from these datasets to data available on

665 Neuroelectro.org [36] and on Hippocampome.org [37]. We also performed a more specific

666 review of the relevant literature to compare the most important somatic features of our patch

667 clamp dataset to results from published patch clamp recordings [38–43] (see Methods). We

668 conclude that the patch clamp dataset is in good agreement with experimental observations

669 available in the literature, and will be used as a representative example in this study.

670       The two datasets used in this study (sharp electrode dataset, patch clamp dataset) differ

671 not only in the recording technique, but also in the simulation protocol. In the sharp electrode

672 recordings, the cells received 400 ms-long depolarizing and hyperpolarizing current injections,

29

673    using amplitudes of 0.2, 0.4, 0.6, 0.8 and 1.0 nA in both directions. In the patch clamp

674    recordings, both the depolarizing and the hyperpolarizing current injections were 300 ms long

675    and 0.05, 0.1, 0.15, 0.2, 0.25 nA in amplitude.

676        As each of the tested models apparently used experimental data obtained from patch

677    clamp recordings as a reference, here we show the detailed results of the test on the models

678    when their output was compared to the features extracted from the patch clamp data (we will

679    return to the comparison between the two datasets near the end of this section). During these

680    recordings the cells were stimulated with relatively low amplitude current injections. Some of

681    the examined models (Migliore et al. 2011, Gómez González et al. 2011 n125 morphology) did

682    not fire even for the highest amplitude tested. Some other models started to fire for higher

683    current intensities than it was observed experimentally. In these cases the features that describe

684    action potential shape or timing properties cannot be evaluated for the given model (for the

685    current amplitudes affected). Therefore, besides the final score achieved by the models on this

686    test (the average Z-score for the successfully evaluated features – see Methods for details), we

687    also consider the proportion of the successfully evaluated features as an important measure of

688    how closely the model matches this specific experimental dataset (Fig 4B).

689        Fig 3 shows how the extracted feature values of the somatic response traces of the

690    different models fit the experimental values. It is clear that the behavior of the different models

691    is very diverse. Each model captures some of the experimental features but shows a larger

692    discrepancy for others.

30

Fig 3: Feature values from the Somatic Features Test of HippoUnit applied to several published models. Absolute feature values extracted from the voltage responses of the models to somatic current injections of varying amplitude, compared to experimental values (darkest red) that were extracted from the patch clamp dataset . (Not all the evaluated features are shown here.)

31

699     The resting membrane potential (*voltage_base*) for all of the models was apparently

700     adjusted to a more hyperpolarized value than in the experimental recordings we used for our

701     comparison, and most of the models also return to a lower voltage value after the step stimuli

702     (*steady_state_voltage*). An exception is the Poirazi et al. 2003 model, where the decay time

703     constant after the stimulus is unusually high (data not shown in Fig 3, but the slow decay can

704     be seen in the example trace in Fig 2.). The voltage threshold for action potential generation

705     (*AP_begin_voltage*) is lower than the experimental value for most of the models (that were able

706     to generate action potentials in response to the examined current intensities), but it is higher

707     than the experimental value for most versions of the Gómez González et al. 2011 model. For

708     negative current steps most of the models gets more hyperpolarized (*voltage_deflection*) (the

709     most extreme is the Gómez González et al. 2011 model with the n129 morphology), while the

710     Gómez González et al. 2011 model with the n125 morphology and the Migliore et al. 2011

711     model get less hyperpolarized than it was observed experimentally. The sag amplitudes are also

712     quite high for the Gómez González et al. 2011 n129, and n130 models, while the Katz et al.

713     2009, and all versions of the Golding et al. 2001 models basically have no hyperpolarizing sag.

714      It is quite conspicuous how much the amplitude of the action potentials (*APlast_amp,*

715     *AP_amplitude, AP2_amp*) differs in the Gómez González et al. 2011 models from the

716     experimental values and from the other models as well. The Katz et al. 2009 and one of the

717     versions (Fig 8A) of the Golding et al. 2001 model have slightly too high action potential

718     amplitudes, and these models have relatively small action potential width (*AP_width*). On the

719     other hand, the rising phase (*AP_rise_time, AP_rise_rate*) of the Katz et al. 2009 model appears

720     to be too slow.

721     Looking at the inverse interspike interval (*ISI*) values, it can be seen that the

722     experimental spike trains show adaptation in the ISIs, meaning that the first ISI is smaller (the

723     inverse ISI is higher) than the last ISI for the same current injection amplitude. This behavior
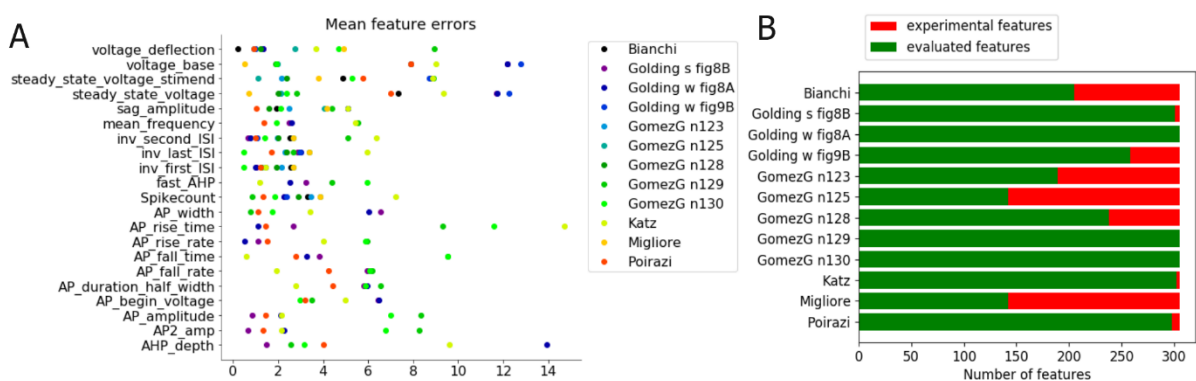
32

724    can be observed in the case of the Katz et al. 2009 model, three versions (n128, n129, n130

725    morphology) of the Gómez González et al. 2011 model, but cannot really be seen in the Bianchi

726    et al. 2011, the Poirazi et al. 2003 and the three versions of the Golding et al. 2001 models. At

727    first look it may seem contradictory that in the case of the Gómez González et al. 2011 model

728    version n129 morphology the spike counts are quite low, while the mean frequency and the

729    inverse ISI values are high. This is because the soma of this model does not fire over the whole

730    period of the stimulation, but starts firing at higher frequencies, then stops firing for rest of the

731    currently used (relative short) stimulus (see Fig 2), although it would start firing again for longer

732    current injections (data not shown). The Katz et al. 2009 model fires quite a high number of

733    action potentials (*Spikecount*) compared to the experimental data, at a high frequency.

734         In the experimental recordings there is a delay before the first action potential is

735    generated, which becomes shorter with increasing current intensity (indicated by the

736    *inv_time_to_first_spike* feature that becomes larger with increasing input intensity). In most of

737    the models this behavior can be observed, albeit to different degrees. The Katz et al. 2009 model

738    has the shortest delays (highest *inv_time_to_first_spike* values), but the effect is still visible.

739         To quantify the difference between the experimental dataset and the simulated output of

740    the models, these were compared using the feature-based error function (Z-Score) described

741    above to calculate the feature errors. Fig 4A shows the mean error scores of the model features

742    whose absolute values are illustrated in Fig 3 (averaged over the different current step

743    amplitudes examined). From this figure it is even more clearly visible that each model fits some

744    experimental features well but does not capture others. For example, it is quite noticeable in Fig

745    4A that most of the versions of the Gómez González et al. 2011 model (greenish dots) perform

746    well for features describing action potential timing (upper part of the figure, e.g., *ISIs,*

747    *mean_frequency*, *spikecount*), but get higher error scores for features of action potential shape

748    (lower part of the figure, e.g., *AP_rise_rate, AP_rise_time, AP_fall_rate, AP_fall_time, AP*

33

749 *amplitudes*). Conversely, the Katz et al. 2009 model achieved better scores for AP shape

750 features than for features describing AP timing. It is also worth noting that none of the error

751 scores for the model of Migliore et al. 2011 was higher than 4; however, looking at Fig 4B it

752 can be seen that less than half of the experimental features were successfully evaluated in this

753 model, which is because it does not fire action potentials for the current injection amplitudes
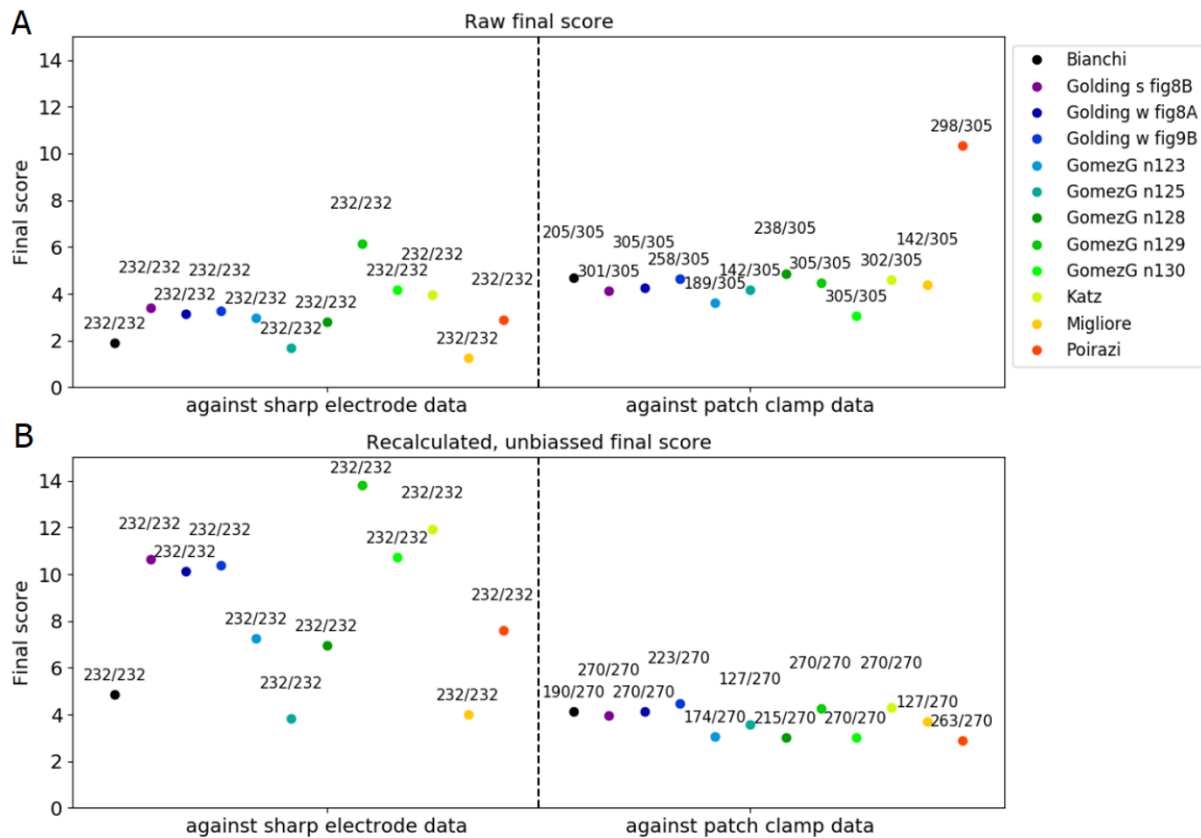
754 examined here.

755



756

757 Fig 4: Evaluation of results from the Somatic Features Test of HippoUnit applied to published models. (A) Mean

758 feature errors (in units of the experimental SD) of the different models. Feature error values are averaged over the

759 different input step amplitudes. (B) The bars represent the number of features that were attempted to be evaluated

760 for the models (i.e., the number of features extracted from the experimental patch clamp dataset). The number of

761 successfully evaluated features for the various models is shown in green, and the number of features that could not

762 be evaluated for a particular model is shown in red.

763

764 Besides enabling the comparison of different models regarding how well they match a

765 particular dataset, the tests of HippoUnit also allow one to determine the match between a

766 particular model and several datasets of the same type. As experimental results can be heavily

767 influenced by recording conditions and protocols, and also depend on factors such as the strain,

768 age, and sex of the animal, it is important to find out whether the same model can

769 simultaneously capture the outcome of different experiments, and if not, how closely it is able

770   to match the different datasets. As a practically relevant example, we looked at how well the

771   various published models that we were testing captured a different experimental dataset that

772   also contained current clamp recordings from rat CA1 PCs, but which was obtained using sharp

773   electrodes rather than the whole-cell patch clamp technique [3]. We therefore evaluated all the

774   models with the Somatic Features Test of HippoUnit using both datasets, and then compared

775   the results.

776       When we simply compared the raw outputs of the test for each model evaluated using

777   the two different data sets (Fig 5A) we identified two factors that substantially bias the results.

778   First, we found that the standard deviation values for the features extracted from the two

779   datasets are very different in magnitude; more specifically, the patch clamp recording data set

780   had much lower standard deviation values for most of the features. This results in relatively

781   higher feature error scores achieved by the models, as the difference of the model output from

782   the experimental features is given in the unit of the experimental standard deviation. The other

783   source of bias is that not all the features could be extracted from both of the data sets, and, as

784   mentioned before, not the same current step protocol (current intensity, duration) was used

785   during the different experiments. Consequently, the models are not compared to exactly the

786   same set of features in the two cases. Mainly as a result of these two confounding factors,

787   comparison of the raw scores of the models for the two data sets (Fig 5A) appears to indicate

788   that most models fit the dataset obtained from sharp electrode recordings better, even though

789   these models were typically built mostly based on patch clamp data.

790

35

Fig 5: Comparison of the final scores achieved by the different models on the Somatic Features Test against validation data from two different datasets (sharp electrode data, patch clamp data). In the upper panel (A) the raw output of the tests is shown, while in the lower panel (B) the feature errors and the final scores have been recalculated using standardized standard deviation values. Numbers above each data point show the proportion of the successfully evaluated features compared to the number of features attempted to be evaluated (successfully extracted from the data set). Note that while in the recalculated final scores (B) only those eFEL features were taken into account that could be extracted from both datasets, they are extracted for different current step amplitudes, which accounts for the difference in the number of observation features for the two datasets.

To overcome these issues and make unbiased comparisons of the models to the two datasets, the feature error scores and the final scores were recalculated in the following way (Fig 5B). The new feature error scores for the two different data sets were calculated as the difference of the model's feature value from the mean feature value of each dataset (as before), but divided by a common standard deviation value. This standardized SD value for each eFEL

36

806   feature was the mean of the standard deviation values over the current steps in the patch clamp

807   dataset (the results were qualitatively similar if we used the SD values from the sharp electrode

808   dataset everywhere instead). Averaging the standard deviation values of the eFEL features over

809   the current steps was required because the current step amplitudes were not the same in the two

810   data sets, and we therefore needed to define SD values that were independent of the amplitude.

811   To get rid of the second bias, only those eFEL features were used in the final score recalculation

812   that are present in both observations (sharp electrode and patch clamp datasets) for at least one

813   current step amplitude.  (This change had the side effect of significantly decreasing the final

814   score for the Poirazi et al. 2003 model because the feature *decay_time_constant_after_stim* was

815   excluded here, as it could not be extracted from the sharp electrode data.) Now that the final

816   scores are recalculated to get rid of most of the biasing factors, it becomes clear that the somatic

817   behavior of every model fits the patch clamp data better (Fig 5B).

818        It is worth noting that one biasing factor still remains in the last comparison: as it has

819   already been mentioned, not all the observation features can be evaluated for each of the models,

820   especially when they are compared to the patch clamp data set, which uses smaller currents. To

821   allow the assessment of the potential effect of this issue, the proportion of the successfully

822   evaluated features relative to the number of features attempted to be evaluated (successfully

823   extracted from the data set) for each model is also shown in Fig 5 next to each data point.

824

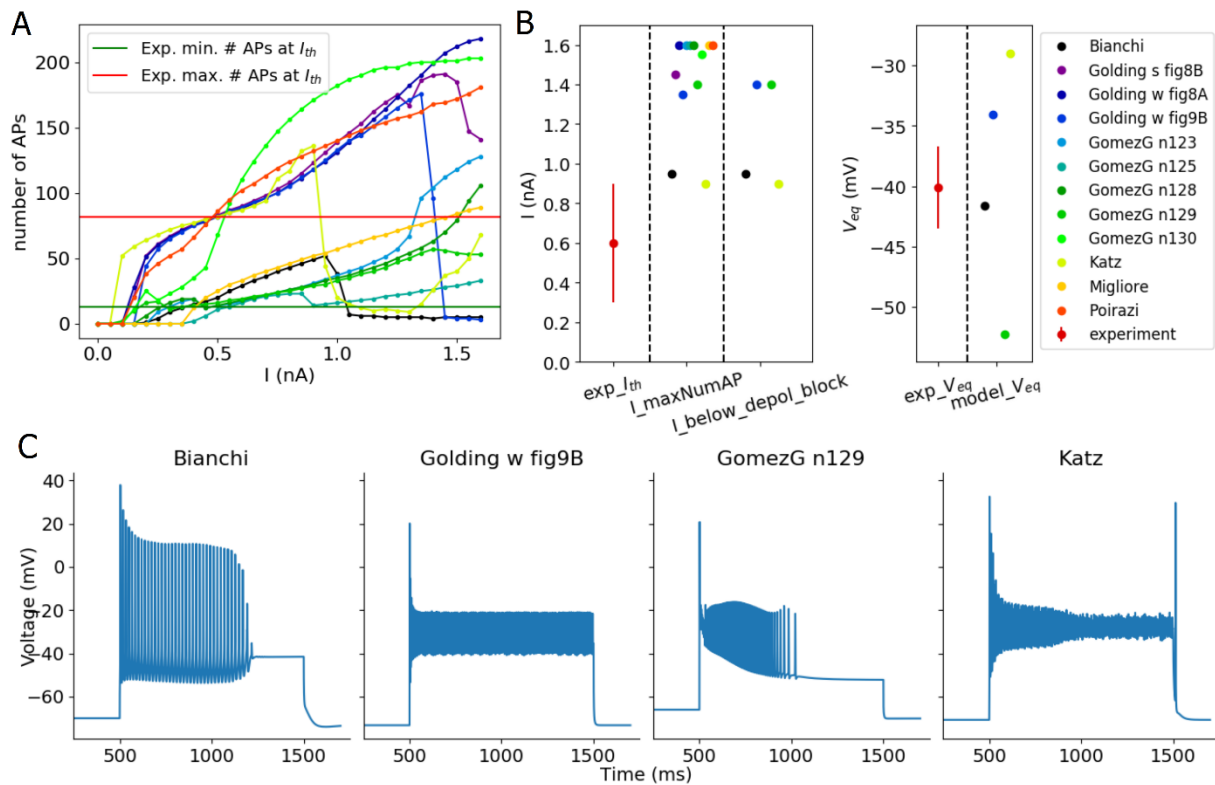825   **Depolarization Block Test**

826

827        In the Depolarization Block Test three features are evaluated. Two of them examine the

828   threshold current intensity to reach depolarization block. The *I_maxNumAP* feature is the

829   current intensity at which the model fires the maximum number of action potentials, and the

37

830    *I_below_depol_block* feature is the current intensity one step before the model enters

831    depolarization block. Both are compared to the experimental $I_{th}$ feature because, in the

832    experiment [24], the number of spikes increased monotonically with increasing current

833    intensity up to the current amplitude where the cell entered depolarization block during the

834    stimulus, which led to a drop in the number of action potentials. By contrast, we experienced

835    that some models started to fire fewer spikes for higher current intensities while still firing over

836    the whole period of the current step stimulus, i.e., without entering depolarization block.

837    Therefore, we introduced the two separate features for the threshold current. If these two feature

838    values are not equal, a penalty is added to the score. The third evaluated feature is $V_{eq}$, the

839    equilibrium potential during the depolarization block, which is calculated as the average of the

840    membrane potential over the last 100 ms of a current pulse with amplitude 50 pA above

841    *I_maxNumAP* (or 50 pA above *I_below_depol_block* if its value is not equal to *I_maxNumAP*).

842    Each model has a value for the "I_maxNumAP" feature, while those models that do not enter

843    depolarization block are not supposed to have a value for the *I_below_depol_block* feature and

844    the *Veq* feature.

845        The results from applying the Depolarization Block Test to the models from ModelDB

846    are shown in Fig 6. According to the test, four of the models entered depolarization block.

847    However, by looking at the actual voltage traces provided by the test, it becomes apparent that

848    only the Bianchi et al. 2011 model behaves correctly (which was developed to show this

849    behavior). The other three models actually managed to "cheat" the test.

850

38

Fig 6: Results from the Depolarization Block Test of HippoUnit applied to published models. (A) Number of APs fired by the models in response to current injections of increasing intensity. (B) Depolarization block feature values extracted from the voltage responses of the models. (C) Voltage traces of different models that were recognized by the test as depolarization block.

In the case of the Katz et al. 2009 and the Golding et al. 2001 Fig 9B models, the APs get smaller and smaller with increasing stimulus amplitude until they get so small that they do not reach the threshold for action potential detection; therefore, these APs are not counted by the test and $V_{eq}$ is also calculated. The Gómez González et al. 2011 model adjusted to the n129 morphology does not fire during the whole period of the current stimulation for a wide range of current amplitudes (see Fig 2). Increasing the intensity of the current injection it fires an increasing number of spikes, but always stops after a while before the end of the stimulus. On the other hand, there is a certain current intensity after which the model starts to fire fewer action potentials, and which is thus detected as $I\_maxNumAP$ by the test. Because no action

866    potentials can be detected during the last 100 ms of the somatic response one step above the

867    detected "threshold" current intensity, the model is declared to have entered depolarization

868    block, and a $V_{eq}$ value is also extracted. These cases underline the importance of critically

869    evaluating the full output of the tests rather than blindly accepting the final scores provided.

870

871    **Back-propagating Action Potential Test**

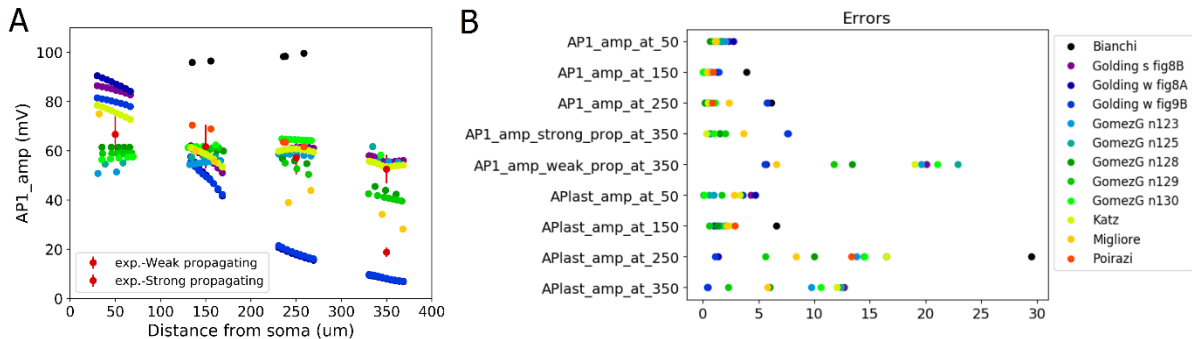872

873            This test first finds all the dendritic segments that belong to the main apical dendrite of

874    the model and which are 50, 150, 250, 350 ± 20 μm from the soma, respectively. Then a train

875    of action potentials of frequency around 15 Hz is triggered in the soma by injecting a step

876    current of appropriate amplitude (as determined by the test), and the amplitudes of the first and

877    last action potentials in the train are measured at the selected locations. In the Bianchi et al.

878    2012 and the Poirazi et al. 2003 models (which share the same morphology, see Fig 2) no

879    suitable trunk locations could be found in the most proximal (50 ± 20 μm) and most distal (350

880    ± 20 μm) regions. This is because this morphology has quite long dendritic sections that are

881    divided into a small number of segments. In particular, the first trunk section

882    (apical_dendrite[0]) originates from the soma, is 102.66 μm long, and has only two segments.

883    The center of one of them is 25.67 μm far from the soma, while the other is already 77 μm away

884    from the soma. None of these segments belongs to the 50 ± 20 μm range, and therefore they are

885    not selected by the test. The n123 morphology of the Gómez González et al. 2011 model has

886    the same shape (Fig 2), but in this case the segments are different, and therefore it does not

887    share the same problem.

888            At the remaining, successfully evaluated distance ranges in the apical trunk of the

889    Bianchi et al. 2012 model, action potentials propagate very actively, barely attenuating. For the

40

890     *AP1_amp* and *APlast_amp* features at these distances, this model has the highest error score

891     (Fig 7), while the Poirazi et al. 2003 model performs quite well.

892



893

894     Fig 7: Results from the Back-propagating Action Potential Test of HippoUnit applied to published models. (A)

895     The amplitudes of the first back-propagating action potentials (in a train of spikes with frequency around 15 Hz)

896     as a function of recording location distance from the soma. (B) Feature error scores achieved by the different

897     models on the Back-propagating AP Test. The amplitudes of the first and last back-propagating action potentials

898     were averaged over the distance ranges of 50, 150, 250, 350 ±µm and compared to the experimental features (see

899     Methods for more details).

900

901     The Golding et al. 2001 model was designed to investigate how the distribution of ion

902     channels can affect the back-propagation efficacy in the trunk. The two versions of the Golding

903     et al. 2001 model ("fig8A" and "fig9B" versions) which are supposed to be weakly propagating

904     according to the corresponding paper [26], are also weakly propagating according to the test.

905     However, the difference between their strongly and weakly propagating feature error scores is

906     not too large (Fig 7), which is probably caused by the much smaller standard deviation value

907     of the experimental data for the weakly propagating case. Although the amplitudes of the first

908     action potentials of these two models fit the experimental data relatively well, they start to

909     decline slightly closer to the soma than it was observed experimentally, as the amplitudes are

910     already very small at $250 \pm 20$ µm (Fig 7). (In Fig 7 the data corresponding to these two versions

41

911    of the model are almost completely overlapping for more distal regions.) The amplitudes for

912    the last action potential fit the data well, except in the most proximal regions (data not shown).

913    For all versions of the Golding et al. 2001 model, AP amplitudes are too high at the most

914    proximal distance range. As for the strongly propagating version of the Golding et al. 2001

915    model ("fig8B" version), the amplitude of the first action potential is too high at the proximal

916    locations, but further it fits the data well. The amplitude of the last action potential remains too

917    high even at more distal locations. It is worth noting that, in the corresponding paper [26], they

918    only examined a single action potential triggered by a 5 ms long input in their simulations, and

919    did not examine or compare to their data the properties of the last action potential in a longer

920    spike train. Finally, we note that in all versions of the Golding et al. 2001 model a spike train

921    with frequency around 23 Hz was evoked and examined as it turned out to be difficult to set the

922    frequency closer to 15 Hz.

923        The different versions of the Gómez González et al. 2011 model behave qualitatively

924    similarly in this test, although there were smaller quantitative differences. In almost all versions

925    the amplitudes of the first action potential in the dendrites are slightly too low at the most

926    proximal locations but fit the experimental data better at further locations. The exceptions are

927    the versions with the n128 and n129 morphologies, which have lower first action potential

928    amplitudes at the furthest locations, but not low enough to be considered as weak propagating.

929    The amplitudes for the last action potential are too high at the distal regions but fit better at the

930    proximal ones. The only exception is the one with morphology n129, where the last action

931    potential attenuates more at further locations and fits the data better.

932        In the case of the Katz et al. 2009 model, a spike train with frequency around 40 Hz was

933    examined, as the firing frequency increases so suddenly with increasing current intensity in this

934    model that no frequency closer to 15 Hz could be adjusted. In this model the last action potential

42

935     propagates too strongly, while the dendritic amplitudes for the first action potential are close to

936     the experimental values.

937          In the Migliore et al. 2011 model the amplitudes for the last action potential are too high,

938     while the amplitude of the first back-propagating action potential is too low at locations in the

939     $250 \pm 20$ µm and $350 \pm 20$ µm distance ranges.

940          Finally, all the models that we examined were found to be strongly propagating by the

941     test, with the exception of those versions of the Golding et al. 2001 model that were explicitly
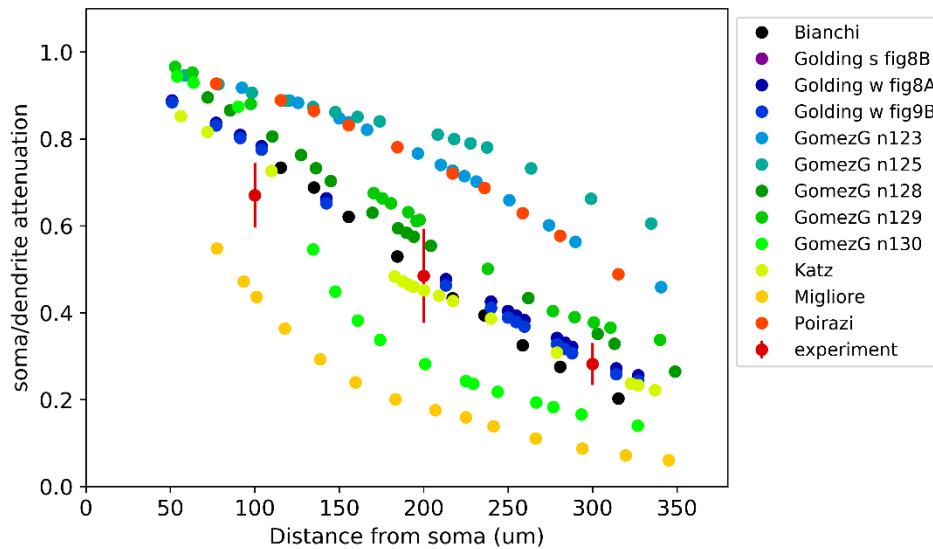
942     developed to be weakly propagating.

943

944     **PSP Attenuation Test**

945

946          In this test the extent of the attenuation of the amplitude of an excitatory post-synaptic

947     potential (EPSP) is examined as it propagates towards the soma from different input locations

948     in the apical trunk. The Katz et al. 2009, the Bianchi et al. 2012, and all versions of the Golding

949     et al. 2001 models perform quite well in this test. The various versions of the Golding et al.

950     2001 model are almost identical in this respect, which is not surprising as they differ only in

951     the distribution of the sodium and A-type potassium channels. This shows that, as we would

952     expect, these properties do not have much effect on the propagation of relatively low-amplitude

953     signals such as unitary PSPs. Interestingly, the different versions of the Gómez González et al.

954     2011 model, with different morphologies, behave quite differently, which shows that this

955     behavior can depend very much on the morphology of the dendritic tree.

956

43

Fig 8: Results from the PSP Attenuation Test of HippoUnit applied to published models. Soma/dendrite EPSP attenuation as a function of the input distance from the soma in the different models.

**Oblique Integration Test**

This test probes the integration properties of the radial oblique dendrites of CA1 pyramidal cell models. The test is based on the experimental results described in [32]. In this study, the somatic voltage response was recorded while synaptic inputs in single oblique dendrites were activated in different spatio-temporal combinations using glutamate uncaging. The main finding was that a sufficiently high number of synchronously activated and spatially clustered inputs produced a supralinear response consisting of a fast (Na) and a slow (NMDA) component, while asynchronously activated inputs summed linearly or sublinearly.

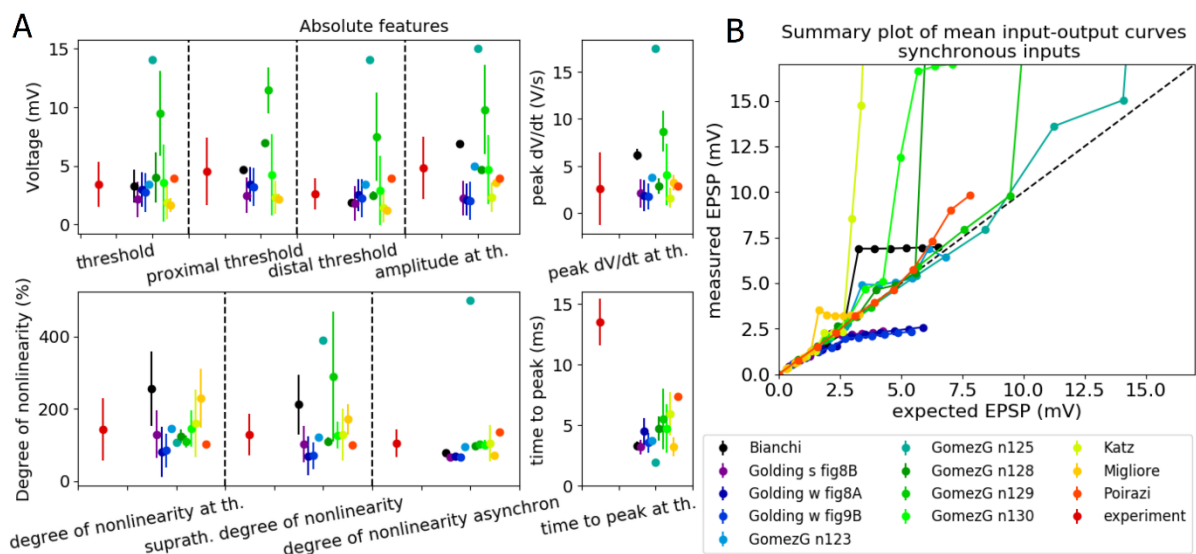This test selects all the radial oblique dendrites of the model that meet the experimental criteria: they are terminal dendrites (they have no child sections) and are at most 120 μm from the soma. Then the selected dendrites are stimulated in a proximal and in a distal region (separately) using an increasing number of clustered, synchronous or asynchronous synaptic inputs to get the voltage responses of the model, and extract the features of dendritic integration.

44

975   The synaptic inputs are not unitary inputs, i.e., their strength is not equivalent to the strength of

976   one synapse in the real cell; instead, the strength is adjusted in a way that 5 synchronous inputs

977   are needed to trigger a dendritic action potential. The intensity of the laser used for glutamate

978   uncaging was set in a similar way in the experiments [32]. Most of the features were extracted

979   at this just-suprathreshold level of input. We noticed that in some cases the strength of the

980   synapse is not set correctly by the test; for example, it may happen that an actual dendritic spike

981   does not reach the spike detection threshold in amplitude, or sometimes the EPSP may reach

982   the threshold for spike detection without actual spike generation. The user has the ability to set

983   the threshold used by eFEL for spike detection, but sometimes a single threshold may not work

984   even for the different oblique dendrites (and proximal and distal locations in the same dendrites)

985   of a single model. For consistency, we used the same spike detection threshold of -20 mV for

986   all the models.

987        The synaptic stimulus contains an AMPA and an NMDA receptor-mediated component.

988   As the default synapse, HippoUnit uses the Exp2Syn double exponential synapse built into

989   NEURON for the AMPA component, and its own built-in NMDA receptor model, whose

990   parameters were set according to experimental data from the literature (see the Methods section

991   for more details). In those models that originally do not have any synaptic component (the

992   Bianchi et al 2011 model and all versions of the Golding et al. 2001 model) this default synapse

993   was used. Both the Katz et al. 2009 and the Migliore et al. 2011 models used the Exp2Syn in

994   their simulations, so in their case the time constants of this function were set to the values used

995   in the original publications. As these models did not contain NMDA receptors, the default

996   NMDA receptor model and the default AMPA/NMDA ratio of HippoUnit were used. The

997   Gómez González et al 2011 and the Poirazi et al. 2003 models have their own AMPA and

998   NMDA receptor models and they own AMPA/NMDA ratio values to be tested with.

45

999       As shown by the averaged "measured EPSP vs expected EPSP" curves in Fig 9, all three

1000     versions of the Golding et al. 2001 model have a jump in the amplitude of the somatic response

1001     at the threshold input level, which is the result of the generation of dendritic spikes. However,

1002     even these larger average responses do not reach the supralinear region, as it would be expected

1003     according to the experimental observations [32]. The reason for this discrepancy is that a

1004     dendritic spike was generated in the simulations in only a subset of the stimulated dendrites; in

1005     the rest of the dendrites tested, the amplitude of the EPSPs went above the spike detection

1006     threshold during the adjustment of the synaptic weight without actually triggering a dendritic

1007     spike, which led to the corresponding synaptic strength being incorrectly set for that particular

1008     dendrite. Averaging over the results for locations with and without dendritic spikes led to an

1009     overall sublinear integration profile.

1010



1011

1012 Fig 9: Results from the Oblique Integration Test of HippoUnit applied to published models. (A) Comparison of

1013 the responses of the models to experimental results (dark red) according to features of dendritic integration. (B)

1014 The averaged input – output curves of all the dendritic locations examined. EPSP amplitudes are measured at the

1015 soma.

1016

46

1017       The Migliore et al. 2011 model performs quite well on this test. In this case, seven

1018       dendrites could be tested out of the ten dendrites within the correct distance range because, in

1019       the others, the dendritic spike at the threshold input level also elicited a somatic action potential,

1020       and therefore these dendrites were excluded from further testing.

1021       In the Katz et al. 2009 model all the selected dendritic locations could be tested, and in

1022       most of them the synaptic strength could be adjusted appropriately. For a few dendrites, some

1023       input levels higher than the threshold for dendritic spike generation also triggered somatic

1024       action potentials. This effect causes the high supralinearity in the "measured EPSP vs expected

1025       EPSP" curve in Fig 9, but has no effect on the extracted features.

1026       In the Bianchi et al. 2012 model only one dendrite could be selected, in which very high

1027       amplitude dendritic spikes were evoked by the synaptic inputs, making the signal integration

1028       highly supralinear.

1029       In the Poirazi et al. 2003 model also only one dendrite could be selected based on its

1030       distance from the soma; furthermore, only the distal location could be tested even in this

1031       dendrite, as at the proximal location the dendritic action potential at the threshold input level

1032       generated a somatic action potential. However, at the distal location, the synaptic strength could

1033       not be set correctly. For the synaptic strength chosen by the test, the actual threshold input level

1034       where a dendritic spike is first generated is at 4 inputs, but this dendritic AP is too small in

1035       amplitude to be detected, and the response to 5 inputs is recognized as the first dendritic spike

1036       instead. Therefore, the features that should be extracted at the threshold input level are instead

1037       extracted from the voltage response to 5 inputs. In this model this results in a reduced

1038       *supralinearity* value, as this feature is calculated one input level higher than the actual threshold.

1039       In addition, for even higher input levels dendritic bursts can be observed, which causes large

1040       *supralinearity* values in the "measured EPSP vs expected EPSP" curve in Fig 9, but this does

1041       not affect the feature values.

1042    Models from Gómez González et al. 2011 were expected to be particularly relevant for

1043    this test, as these models were tuned to fit the same data set on which this test is based. However,

1044    we encountered an important issue when comparing our test results for these models to the

1045    results shown in the paper [53]. In particular, the paper clearly indicates which dendrites were

1046    examined, and it is stated that those are at maximum 150 μm from the soma. However, when

1047    we measured the distance of these locations from the soma by following the path along the

1048    dendrites (as it is done by the test of HippoUnit), we often found it to be larger than 150 μm.

1049    We note that when the distance was measured in 3D coordinates rather than along the dendrites,

1050    all the dendrites used by Gómez González et al. 2011 appeared to be within 150 μm of the

1051    soma, so we assume that this definition was used in the paper. As we consider the path distance

1052    to be more meaningful than Euclidean distance in this context, and this was also the criterion

1053    used in the experimental study, we consistently use path distance in HippoUnit to find the

1054    relevant dendritic segments. Nevertheless, this difference in the selection of dendrites should

1055    be kept in mind when the results of this validation for models of Gómez González et al. 2011

1056    are evaluated.

1057    In two versions of the Gómez González et al. 2011 model (those that were adjusted to

1058    the n123 and n125 morphologies) only one oblique dendrite matched the experimental criteria

1059    and could therefore be selected, and these are not among those that were studied by the

1060    developers of the model. In each of these cases the dendritic spike at the proximal location at

1061    the input threshold level triggered a somatic action potential, and therefore only the distal

1062    location could be tested. In the case of the n125 morphology, the dendritic spikes that appear

1063    first for just-suprathreshold input are so small in amplitude that they do not reach the spike

1064    detection threshold (-20 mV), and are thus not detected. Therefore, the automatically adjusted

1065    synaptic weight is larger than the appropriate value would be, which results in larger somatic

1066    EPSPs than expected (see Fig 9). With this synaptic weight the first dendritic spike, and

48

1067    therefore the jump in the "measured EPSP vs expected EPSP" curve to the supralinear region

1068    is for 4 synaptic inputs, instead of 5. This is also the case in one of the two selected dendrites

1069    of the version of this model with the n128 morphology. Similarly to the Poirazi et al. 2003

1070    model, this results in a lower *degree of nonlinearity at threshold* feature value, than it would be

1071    if the feature were extracted at the actual threshold input level (4 inputs) instead of the one

1072    which the test attempted to adjust (5 inputs). The *suprathreshold nonlinearity* feature has a high

1073    value because at that input level (6 inputs), somatic action potentials are triggered.

1074    In the version of the Gómez González et al. 2011 model that uses the n129 morphology,

1075    10 oblique dendrites could be selected for testing (none of them is among those that its

1076    developers used) but only 4 could be tested because, for the rest, the dendritic spike at the

1077    threshold input level already elicits a somatic action potential. The synaptic weights required to

1078    set the threshold input level to 5 are not found correctly in most cases; the actual threshold input

1079    level is at 4 or 3. Suprathreshold nonlinearity is high, because at that input level (6 inputs)

1080    somatic action potentials are triggered for some of the examined dendritic locations.

1081    The version of the Gómez González et al. 2011 model that uses the n130 morphology

1082    achieves the best (lowest) final score on this test. In this model many oblique dendrites could

1083    be selected and tested, including two (179, 189) that the developers used in their simulations

1084    [53]. In most cases the synaptic weights are nicely found to set the threshold input level to 5

1085    synapses. For some dendrites there are somatic action potentials at higher input levels, but that

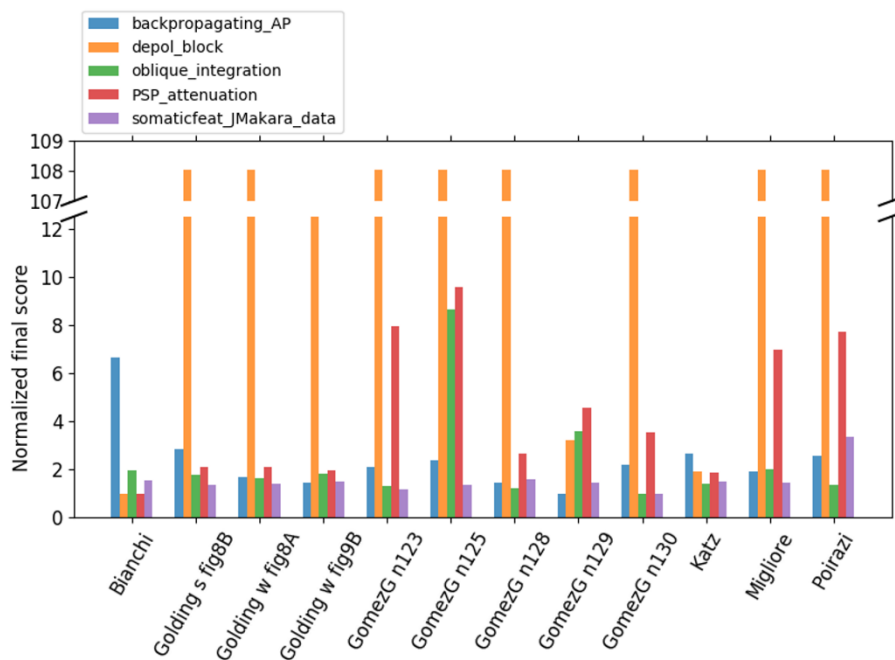1086    does not affect the features.

1087    The value of the *time to peak* feature for each model is much smaller than the

1088    experimental value (Fig 9). This is because in each of the models the maximum amplitude of

1089    the somatic EPSP is determined by the fast component, caused by the appearance of the

1090    dendritic sodium spikes, while in the experimental observation this is rather shaped by the slow

1091    NMDA component following the sodium spike.

49

**Overall performance and model comparison**

1092
1093

1094    In summary, using HippoUnit, we compared the behavior of several hippocampal CA1

1095    pyramidal cell models available on ModelDB in several distinct domains, and found that all of

1096    these models match experimental results well in some domains (typically those that they were

1097    originally built to capture) but fit the experimental observation less precisely in others. Fig 10

1098    summarizes the final scores achieved by the different models on the various tests (lower scores

1099    indicate a better match in all cases).

1100



1101

Fig 10: Normalized final scores achieved by the different published models on the various tests of HippoUnit. The

final scores of each test are normalized by dividing the scores of each model by the best achieved score on the

given test.

1105

1106    Perhaps a bit surprisingly, the different versions of the Golding et al. 2001 model

1107    perform quite well on all of the tests (except for the Depolarization Block Test), even though

1108    these are the simplest ones among the models in the sense that they contain the smallest number

50

1109   of different types of ion channels. On the other hand, they do not perform outstandingly well

1110   on the Back-propagating Action Potential Test, although they were developed to capture the

1111   behavior evaluated by this test.

1112       The Bianchi et al. 2012 model is the only one that can produce real depolarization block

1113   within the range of input strengths examined by the corresponding test. The success of this

1114   model in this test is not surprising because this is the only model that was tuned to reproduce

1115   this behavior; on the other hand, the failure of the other models in this respect clearly shows

1116   that proper depolarization block requires some combination of mechanisms that are at least

1117   partially distinct from those that allow good performance in the other tests. The Bianchi et al.

1118   2012 model achieves a relatively high final error score only on the Back-propagating Action

1119   Potential Test, as action potentials seem to propagate too actively in its dendrites, leading to

1120   high AP amplitudes even in more distal compartments.

1121       The Gómez González et al. 2011 models were developed to capture the same

1122   experimental observations on dendritic integration that are tested by the Oblique Integration

1123   Test of HippoUnit, but, somewhat surprisingly, some of its versions achieved quite high error

1124   scores on this test, while others perform quite well. This is partly caused by the fact that

1125   HippoUnit often selects different dendritic sections for testing from those that were studied by

1126   the developers of these models (see above for details). Some of its versions also perform

1127   relatively poorly on the PSP-Attenuation Test, similar to the Migliore et al. 2011 and the Poirazi

1128   et al. 2003 models. The Katz et al. 2009 model is not outstandingly good in any of the tests, but

1129   still achieves relatively good error scores everywhere (although its apparent good performance

1130   on the Depolarization Block Test is misleading - see detailed explanation above).

1131       The model files that were used to test the models described above, the detailed validation

1132   results (all the output files of HippoUnit), and the Jupyter Notebooks that show how to run the

51

1133    tests of HippoUnit on these models are available in the following Github repository:

1134    https://github.com/KaliLab/HippoUnit_demo.

1135

## Application of HippoUnit to models built using automated parameter optimization within the Human Brain Project

1138

1139    Besides enabling a detailed comparison of published models, HippoUnit can also be

1140    used to monitor the performance of new models at various stages of model development. Here,

1141    we illustrate this by showing how we have used HippoUnit within the HBP to systematically

1142    validate detailed multi-compartmental models of hippocampal neurons developed using multi-

1143    objective parameter optimization methods implemented by the open source Blue Brain Python

1144    Optimization Library (BluePyOpt [15]). To this end, we extended HippoUnit to allow it to

1145    handle the output of optimization performed by BluePyOpt (see Methods).

1146    Models of CA1 pyramidal cells were optimized using target feature data extracted from

1147    the same sharp electrode dataset [3] that was also one of the datasets used by the Somatic

1148    Features Test of HippoUnit. However, while during validation all the eFEL features that could

1149    be successfully extracted from the data are considered, only a subset of these features was used

1150    in the optimization (mostly those that describe the rate and timing of the spikes; e.g., the

1151    different inter-spike interval (ISI), time to last/first spike, mean frequency features).
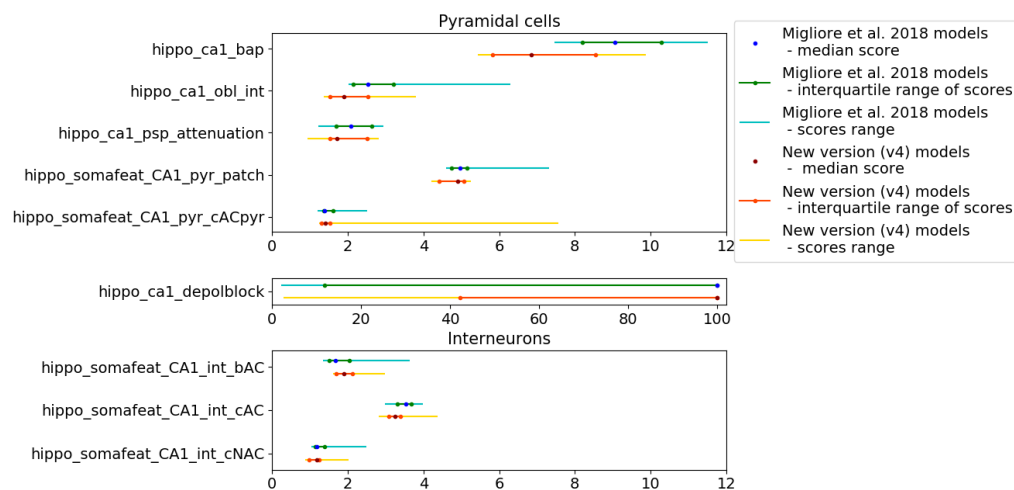
1152    In addition, sharp electrode measurements were also available for several types of

1153    interneuron in the hippocampal CA1 region, and models of these interneurons were also

1154    constructed using similar automated methods [3]. Using the appropriate observation file and

1155    the stimulus file belonging to it, the Somatic Features Test of HippoUnit can also be applied to

52

1156   these models to evaluate their somatic spiking features. The other tests of HippoUnit are

1157   currently not applicable to interneurons, mostly due to the lack of appropriate target data.

1158       We applied the tests of HippoUnit to the version of the models published in [3], and to

1159   a later version (v4) described in Ecker et al. (2020)[55], which was intended to further improve

1160   the dendritic behavior of the models, as this is critical for their proper functioning in the

1161   network. The two sets of models were created using the same morphology files and similar

1162   optimization methods and protocols. These new optimizations differed mainly in the allowed

1163   range for the density of the sodium channels in the dendrites. For the pyramidal cell models a

1164   new feature was also introduced in the parameter optimization that constrains the amplitudes

1165   of back-propagating action potentials in the main apical dendrite. The new interneuron models

1166   also had an exponentially decreasing (rather than constant) density of Na channels, and A-type

1167   K channels with more hyperpolarized activation in their dendrites. For more details on the

1168   models, see the original publications ([3,55]).

1169       After running all the tests of HippoUnit on both sets of models generated by BluePyOpt,

1170   we performed a comparison of the old [3] and the new versions of the models by doing a

1171   statistical analysis of the final scores achieved by the models of the same cell type on the

1172   different tests. In Fig 11 the median, the interquartile range and the full range of the final error

1173   scores achieved by the two versions of the model set are compared. According to the results of

1174   the Wilcoxon signed-rank test the new version of the models achieved significantly better

1175   scores on the Back-propagating Action Potential test (p = 0.0046), on the Oblique Integration

1176   Test (p = 0.0033), and on the PSP Attenuation Test (p = 0.0107), which is the result of reduced

1177   dendritic excitability. Moreover, in most of the other cases the behavior of the models improved

1178   slightly (but not significantly) with the new version. Only in the case of the Somatic Features

1179   test applied to bAC interneurons did the new models perform slightly worse (but still quite

1180   well), and this difference was not significant (p = 0.75).

53

1181     These results show the importance of model validation performed against experimental

1182     findings, especially those not considered when building the model, in every iteration during the

1183     process of model development. This approach can greatly facilitate the construction of models

1184     that perform well in a variety of contexts, help avoid model regression, and guide the model

1185     building process towards a more robust and general implementation.

1186



1188     Fig 11: Statistical comparison of the final scores achieved on the different tests by the two versions of hippocampal

1189     CA1 models of the same cell types, developed by automated optimization using BluePyOpt.

1190

## Integration of HippoUnit into the Validation Framework and the Brain Simulation Platform of the Human Brain Project

1193

1194     The HBP is developing scientific infrastructure to facilitate advances in neuroscience,

1195     medicine, and computing [56]. One component of this research infrastructure is the Brain

1196     Simulation Platform (BSP) (https://bsp.humanbrainproject.eu), an online collaborative

1197     platform that supports the construction and simulation of neural models at various scales. As

1198     we argued above, systematic, automated validation of models is a critical prerequisite of

1199     collaborative model development. Accordingly, the BSP includes a software framework for

54

1200    quantitative model validation and testing that explicitly supports applying a given validation

1201    test to different models and storing the results [57]. The framework consists of a web service,

1202    and a set of test suites, which are Python modules based on the SciUnit package. As we

1203    discussed earlier, SciUnit uses the concept of capabilities, which are standardized interfaces

1204    between the models to be tested and the validation tests. By defining the capabilities to which

1205    models must adhere, individual validation tests can be implemented independently of any

1206    specific model and used to validate any compatible model despite differences in their internal

1207    structures, the language and/or the simulator used. Each test must include a specification of the

1208    required model capabilities, the location of the reference (experimental) dataset, and data

1209    analysis code to transform the recorded variables (e.g., membrane potential) into feature values

1210    that allow the simulation results to be directly and quantitatively compared to the experimental

1211    data through statistical analysis. The web services framework [57] supports the management of

1212    models, tests, and validation results. It is accessible via web apps within the HBP Collaboratory,

1213    and also through a Python client. The framework makes it possible to permanently record,

1214    examine and reproduce validation results, and enables tracking the evolution of models over

1215    time, as well as comparison against other models in the domain.

1216        Every test of HippoUnit described in this paper has been individually registered in the

1217    Validation Framework. The JSON files containing the target experimental data for each test are

1218    stored (besides the HippoUnit_demo GitHub repository) in storage containers at the Swiss

1219    National Supercomputing Centre (CSCS), where they are publicly available. The location of

1220    the corresponding data file is associated with each registered test, so that the data are loaded

1221    automatically when the test is run on a model via the Validation Framework. As the Somatic

1222    Features Test of HippoUnit was used to compare models against five different data sets (data

1223    from sharp electrode measurements in pyramidal cells and interneurons belonging to three

1224    different electronic types, and data obtained from patch clamp recordings in pyramidal cells),

55

1225 these are considered to be and have been registered as five separate tests in the Validation

1226 Framework.

1227 All the models that were tested and compared in this study (including the CA1 pyramidal

1228 cell models from the literature and the BluePyOpt optimized CA1 pyramidal cells and

1229 interneurons of the HBP) have been registered and are available in the Model Catalog of the

1230 Validation Framework with their locations in the CSCS storage linked to them. In addition to

1231 the modifications that were needed to make the models compatible with testing with HippoUnit

1232 (described in the section "Methods – Models from literature"), the versions of the models

1233 uploaded to the CSCS container also contain an `__init__.py` file. This file implements a

1234 python class that inherits all the functions of the `ModelLoader` class of HippoUnit without

1235 modification. Its role is to make the validation of these models via the Framework more

1236 straightforward by defining and setting the parameters of the `ModelLoader` class (such as the

1237 path to the HOC and NMODL files, the name of the section lists, etc.) that otherwise need to

1238 be set after instantiating the `ModelLoader` (see the HippoUnit_demo GitHub repository:

1239 https://github.com/KaliLab/HippoUnit_demo/tree/master/jupyter_notebooks ).

1240 The validation results discussed in this paper have also been registered in the Validation

1241 Framework, with all their related files (output figures and JSON files) linked to them. These

1242 can be accessed using the Model Validation app of the framework.

1243 The Brain Simulation Platform of the HBP contains several online 'Use Cases', which

1244 are available on the platform and help the users to try and use the various established pipelines.

1245 The Use Case called 'Hippocampus Single Cell Model Validation' can be used to apply the

1246 tests of HippoUnit to models that were built using automated parameter optimization within the

1247 HBP.

1248 The Brain Simulation Platform also hosts interactive "Live Paper" documents that refer

1249 to published papers related to the models or software tools on the Platform. Live Papers provide

56

1250    links that make it possible to visualize or download results and data discussed in the respective

1251    paper, and even to run the associated simulations on the Platform. We have created a Live Paper

1252    (https://humanbrainproject.github.io/hbp-bsp-live-

1253    papers/2020/saray_et_al_2020/saray_et_al_2020.html) showing the results of the study

1254    presented in this paper in more detail. This interactive document provides links to all the output

1255    figures and data files resulting from the validation of the models from literature discussed here.

1256    This provides a more detailed insight into their behavior individually. Moreover, as part of this

1257    Live Paper a HippoUnit Use Case is also available in the form of a Jupyter Notebook, which

1258    guides the user through running the validation tests of HippoUnit on the models from literature

1259    that are already registered in the Framework, and makes it possible to reproduce the results

1260    presented here.

## Discussion

### Role of validation in collaborative model building

1264    For anatomically and biophysically detailed data-driven neural models to be predictive,

1265    it is important that they are able to generalize beyond their original scope. However, most

1266    detailed biophysical models to date were built to capture only a few important or interesting

1267    properties of a given neuron type. Systematic testing and comparison of the behavior of these

1268    models is still rare, and thus it is often unknown how these models would behave when used

1269    under different circumstances, and to what extent they can be used to address different scientific

1270    questions. As a result, the modeling community still keeps building new models of the same

1271    cell type for various purposes, instead of reusing and further developing the already existing

1272    ones. On the other hand, in those cases when new models are based on previously published

57

1273    ones, model parameters are often adjusted to fit a new set of experimental data. These

1274    adjustments typically alter the ability of the model to capture the experimental data targeted by

1275    the original model, but this remains unrecognized because of the lack of comprehensive testing.

1276    As we have shown, an illustrative example of this regression issue is the Bianchi et al. 2012

1277    model. This model was mainly based on the Poirazi et al. 2013 model, which was developed to

1278    show specific dendritic behaviors and was even tested using data on back-propagating action

1279    potentials during its development [51]. However, the test results presented above indicate that

1280    when the somatic behavior of this model was adjusted to reproduce experimental observations

1281    on depolarization block, it lost the ability to show realistic back-propagation of action potentials

1282    into the apical dendrite. In addition, some publications on neuronal models simply state that the

1283    model has been validated against electrophysiological data, but the details of these validations

1284    (such as the methods used, the experimental data considered or even the results) are usually not

1285    shared. Finally, more systematic testing and coordinated development would enable building

1286    consensus (community) models, which would aim to capture a wide range of experimental

1287    observations.

1288        The framework for implementing unit tests that make it possible to automatically and

1289    systematically compare the behavior of models to experimental data already exists (SciUnit

1290    [18]). Furthermore, the recently developed Validation Framework of the HBP makes it possible

1291    to collect neural models and validation tests, and supports the application of the registered tests

1292    to the registered models. Most importantly, it makes it possible to save the validation results

1293    and link them to the models in the Model Catalog, making them publicly available and traceable

1294    for the modeling community.

1295        Our goal was to contribute to this new approach in collaborative model building by

1296    developing a validation suite to test the behavior of models of the hippocampal CA1 pyramidal

1297    neuron, which is one of the most studied cell types in the brain. Here we presented how we

58

1298    applied HippoUnit to test and compare the behavior of several models of this cell type available

1299    on ModelDB [17] in several distinct domains against electrophysiological data available in the

1300    literature. Through the example of the models optimized within the HBP we also showed how

1301    a test suite like HippoUnit can be a useful tool in tracing the performance of models during the

1302    process of their development.

1303        Although we consider it essential to evaluate the generalization properties of neural

1304    models and test them in as many domains as possible, it is important to emphasize that a high

1305    error score on a given validation test using a particular experimental dataset does not mean that

1306    the model is not good enough or cannot be useful for a variety of purposes (including the ones

1307    it was originally developed for). The discrepancy between the target data and the model's

1308    behavior, as quantified by the validation tests, may be due to several different reasons. First, all

1309    experimental data contain noise and may have systematic biases associated with the

1310    experimental methods employed. Sometimes the experimental protocol is not described in

1311    sufficient detail to allow its faithful reproduction in the simulations. It may also occur that a

1312    model is based on experimental data that were obtained under conditions that are substantially

1313    different from the conditions for the measurement of the validation target dataset. Using

1314    different recording techniques, such as sharp electrode or patch clamp recordings or the

1315    different circumstances of the experiments (e.g., the strain, age, and sex of the animal, or the

1316    temperature during measurement) can heavily affect the experimental results. Furthermore, the

1317    post-processing of the recorded electrophysiological data can also alter the results. For these

1318    reasons, probably no single model should be expected to achieve an arbitrarily low score on all

1319    of the validation tests developed for a particular cell type. Keeping this in mind, it is important

1320    that the modelers decide which properties of the cell type are relevant for them, and what

1321    experimental conditions they aim to mimic. Validation results should be interpreted or taken

1322    into account accordingly, and the tests themselves may need to be adapted.

59

1323    By providing the software tools and examples on how to validate the different models,

1324    we hope to encourage the modeling community to use more systematic testing during model

1325    development, in order to create neural models that generalize better, and to make the process

1326    of model building more reproducible and transparent.

1327

**Uniform model formats reduce the costs of validation**

1329

1330    Although HippoUnit is built in a way that its tests are, in principle, model-agnostic, so

1331    that the implementation of the tests does not depend on model implementation, it still required

1332    a considerable effort to create the standalone versions of the models from literature to be tested,

1333    even though all of the selected models were developed for the NEURON simulator. This is

1334    because each model has a different file structure and internal logic that needs to be understood

1335    in order to create an equivalent standalone version. When the section lists of the main dendritic

1336    types do not exist, the user needs to create them by extensively analyzing the morphology and

1337    even doing some coding. In order to reduce the costs of systematic validation, models would

1338    need to be expressed in a format that is uniform and easy to test. As HippoUnit already has its

1339    capability functions implemented in a way that it is able to handle models developed in

1340    NEURON, the only requirement for such models is that they should contain a HOC file that

1341    describes the morphology (including the section lists for the main dendritic types of the

1342    dendritic tree) and all the biophysical parameters of the model, without any additional

1343    simulations, GUIs or run-time modifications. Currently, such a standalone version of the

1344    models is not made available routinely in publications or on-line databases, but could be added

1345    by the creators of the models with relatively little effort.

60

1346      On the other hand, applying the tests of HippoUnit to models built in other languages

1347    requires the re-implementation of the capability functions that are responsible for running the

1348    simulations on the model (see Methods). In order to save the user from this effort, it would be

1349    useful to publish neuronal models in a standard and uniform format that is simulator

1350    independent and allows general use in a variety of paradigms. This would allow an easier and

1351    more transparent process of community model development and validation, as it avoids the

1352    need of reimplementation of parts of software tools (such as validation suites), and the creation

1353    of new, (potentially) non-traced software versions. This approach is already initiated for

1354    neurons and neuronal networks by the developers of NeuroML [58], NineML [59], PyNN [60],

1355    Sonata [61], and Brian [62]. Once a large set of models becomes available in these standardized

1356    formats, it will be straightforward to extend HippoUnit (and other similar test suites) to handle

1357    these models.

1358

1359    **Extensibility of HippoUnit**

1360

1361      As HippoUnit is based on the SciUnit package [18] it inherits SciUnits's modular

1362    structure. This means that a test is usually composed of four main classes: the test class, the

1363    model class, the capabilities class and the score class (as described in more detail in the Methods

1364    section). Thanks to this structure it is easy to extend HippoUnit with new tests by implementing

1365    them in new test classes and adding the capabilities and scores needed. The methods of the new

1366    capabilities can be implemented in the `ModelLoader` class, which is a generalized Model class

1367    for models built in the NEURON simulator, or in a newly created Model class specific to the

1368    model to be tested.

61

1369    As HippoUnit is open-source and is shared on GitHub, it is possible for other developers,

1370    modelers or scientists to modify or extend the test suite working on their own forks of the

1371    repository. If they would like to directly contribute to HippoUnit, a 'pull request' can be created

1372    to the main repository.

1373

**Generalization possibilities of the tests of HippoUnit**

1375

1376    In the current version of HippoUnit most of the validation tests can only be used to test

1377    models of hippocampal CA1 pyramidal cells, as the observation data come from

1378    electrophysiological measurements of this cell type and the tests were designed to follow the

1379    experimental protocols of the papers from which these data derive. However, with small

1380    modifications most of the tests can be used for other cell types, or with slightly different

1381    stimulation protocols, if there are experimental data available for the features or properties

1382    tested.

1383    The Somatic Features Test can be used for any cell type and with any current step

1384    injection protocol even in its current form using the appropriate data and configuration files.

1385    These two files must be in agreement with each other; in particular, the configuration file should

1386    contain the parameters of the step current protocols (delay, duration, amplitude) used in the

1387    experiments from which the feature values in the data file derive. In this study this test was used

1388    with two different experimental protocols (sharp electrode measurements and patch clamp

1389    recordings that used different current step amplitudes and durations), and for testing four

1390    different cell types (hippocampal CA1 PC and interneurons).

1391    In the current version of the Depolarization Block Test the properties of the stimulus

1392    (delay, duration, amplitudes) are hard-coded to reproduce the experimental protocol used in a

62

1393    study of CA1 PCs [24]. However, the test could be easily modified to read these parameters

1394    from a configuration file like in the case of other tests, and then the test could be applied to

1395    other cell types if data from similar experimental measurements are available.

1396        As the Back-propagating AP Test examines the back-propagation efficacy of action

1397    potentials in the main apical dendrite (trunk), it is mainly suitable for testing pyramidal cell

1398    models; however, it can be used for PC models from other hippocampal or cortical  regions,

1399    potentially using different distance ranges of the recording sites. If different distances are used,

1400    the feature names ('AP1_amp_X' and 'APlast_amp_X', where X is the recording distance) in

1401    the observation data file and the recording distances given in the stimuli file must be in

1402    agreement. Furthermore, it would also be possible to set a section list of other dendritic types

1403    instead of the trunk to be examined by the test. This way, models of other cell types (with

1404    dendritic trees qualitatively different from those of PCs) could also be tested. The frequency

1405    range of the spike train (10 – 20 Hz, preferring values closest to 15 Hz) is currently hard-coded

1406    in the function that automatically finds the appropriate current amplitude, but the

1407    implementation could be made more flexible in this case as well.

1408        The PSP Attenuation Test is quite general.  Both the distances and tolerance values that

1409    determine the stimulation locations on the dendrites and the properties of the synaptic stimuli

1410    are given using the configuration file. Here again the feature names in the observation data file

1411    ('attenuation_soma/dend_x_um', where x is the distance from the soma) must fit the distances

1412    of the stimulation locations in the configuration file when one uses the tests with data from a

1413    different cell type or experimental protocol. Similarly to the Back-propagating AP Test the PSP

1414    Attenuation Test also examines the main apical dendrite (trunk), but could be altered to use

1415    section lists of other dendritic types.

1416        The Oblique Integration Test is very specific to the experimental protocol of [32]. There

1417    is no configuration file used here, but the synaptic parameters (of the `ModelLoader` class) and

63

1418    the number of synapses to which the model should first generate a dendritic spike

1419    ('threshold_index' parameter of the test class) can be adjusted by the user after instantiating the

1420    `ModelLoader` and the test classes respectively. The time intervals between the inputs

1421    (synchronous (0.1 ms), asynchronous (2.0 ms)) are currently hard-coded in the test.

1422        HippoUnit has been used mainly to test models of rat hippocampal CA1 pyramidal cells

1423    as described above. However, having the appropriate observation data, most of its tests could

1424    easily be adapted to test models of different cell types, even in cases when the experimental

1425    protocol is slightly different from the currently implemented ones. The extent to which a test

1426    needs to be modified in order to test models of other cell types depends on how much the

1427    behavior of the new cell type differs from the behavior of CA1 pyramidal cells, and to what

1428    extent the protocol of the experiment differs from the ones we used as the bases of comparison

1429    in the current study.

1430

## Acknowledgements

1436

## References

1438    1.    Einevoll GT, Destexhe A, Diesmann M, Grün S, Jirsa V, de Kamps M, et al. The

1439        Scientific Case for Brain Simulations. Neuron. 2019;102: 735–744.

1440        doi:10.1016/j.neuron.2019.03.027

1441    2.      Káli S, Freund TF. Distinct properties of two major excitatory inputs to hippocampal

1442            pyramidal cells: A computational study. European Journal of Neuroscience. 2005;22:

1443            2027–2048. doi:10.1111/j.1460-9568.2005.04406.x

1444    3.      Migliore R, Lupascu CA, Bologna LL, Romani A, Courcol JD, Antonel S, et al. The

1445            physiological variability of channel density in hippocampal CA1 pyramidal cells and

1446            interneurons explored using a unified data-driven modeling workflow. PLoS

1447            Computational Biology. 2018;14: 1–25. doi:10.1371/journal.pcbi.1006423

1448    4.      Hay E, Hill S, Schürmann F, Markram H, Segev I. Models of neocortical layer 5b

1449            pyramidal cells capturing a wide range of dendritic and perisomatic active properties.

1450            PLoS Computational Biology. 2011;7. doi:10.1371/journal.pcbi.1002107

1451    5.      Herz AVM, Gollisch T, Machens CK, Jaeger D. Modeling single-neuron dynamics and

1452            computations: A balance of detail and abstraction. Science. 2006;314: 80–85.

1453            doi:10.1126/science.1127240

1454    6.      Poirazi P, Brannon T, Mel BW. Pyramidal neuron as two-layer neural network. Neuron.

1455            2003;37: 989–999. doi:10.1016/S0896-6273(03)00149-1

1456    7.      Bower JM. The 40-year history of modeling active dendrites in cerebellar purkinje cells:

1457            Emergence of the first single cell "community model." Frontiers in Computational

1458            Neuroscience. 2015;9: 1–18. doi:10.3389/fncom.2015.00129

1459    8.      Traub RD, Wong RKS, Miles R, Michelson H. A model of a CA3 hippocampal

1460            pyramidal neuron incorporating voltage-clamp data on intrinsic conductances. Journal

1461            of Neurophysiology. 1991;66: 635–650. doi:10.1152/jn.1991.66.2.635

1462    9.      Almog M, Korngreen A. A quantitative description of dendritic conductances and its

1463            application to dendritic excitation in layer 5 pyramidal neurons. Journal of Neuroscience.

1464            2014;34: 182–196. doi:10.1523/JNEUROSCI.2896-13.2014

65

1465    10.    Markram H, Muller E, Ramaswamy S, Reimann MW, Abdellah M, Sanchez CA, et al.

1466        Reconstruction and Simulation of Neocortical Microcircuitry. Cell. 2015;163: 456–492.

1467        doi:10.1016/j.cell.2015.09.029

1468    11.    Traub RD, Contreras D, Cunningham MO, Murray H, LeBeau FEN, Roopun A, et al.

1469        Single-column thalamocortical network model exhibiting gamma oscillations, sleep

1470        spindles, and epileptogenic bursts. Journal of Neurophysiology. 2005;93: 2194–2232.

1471        doi:10.1152/jn.00983.2004

1472    12.    Schneider CJ, Bezaire M, Soltesz I. Toward a full-scale computational model of the rat

1473        dentate gyrus. Frontiers in Neural Circuits. 2012;6: 1–8. doi:10.3389/fncir.2012.00083

1474    13.    Bezaire MJ, Raikov I, Burk K, Vyas D, Soltesz I. Interneuronal mechanisms of

1475        hippocampal theta oscillations in a full-scale model of the rodent CA1 circuit. eLife.

1476        2016;5: 1–106. doi:10.7554/eLife.18566

1477    14.    Friedrich P, Vella M, Gulyás AI, Freund TF, Káli S. A flexible, interactive software

1478        tool for fitting the parameters of neuronal models. Frontiers in Neuroinformatics. 2014;8:

1479        1–19. doi:10.3389/fninf.2014.00063

1480    15.    van Geit W, Gevaert M, Chindemi G, Rössert C, Courcol JD, Muller EB, et al.

1481        BluePyOpt: Leveraging open source software and cloud infrastructure to optimise model

1482        parameters in neuroscience. Frontiers in Neuroinformatics. 2016;10: 1–18.

1483        doi:10.3389/fninf.2016.00017

1484    16.    Vanier MC, Bower JM. A comparative survey of automated parameter-search methods

1485        for compartmental neural models. Journal of Computational Neuroscience. 1999;7: 149–

1486        171. doi:10.1023/A:1008972005316

1487    17.    McDougal RA, Morse TM, Carnevale T, Marenco L, Wang R, Migliore M, et al.

1488        Twenty years of ModelDB and beyond: building essential modeling tools for the future

66

1489        of neuroscience. Journal of Computational Neuroscience. 2017;42: 1–10.

1490        doi:10.1007/s10827-016-0623-7

1491   18.    Omar C, Aldrich J, Gerkin RC. Collaborative infrastructure for test-driven scientific

1492        model validation. 36th International Conference on Software Engineering, ICSE

1493        Companion 2014 - Proceedings. 2014; 524–527. doi:10.1145/2591062.2591129

1494   19.    Gerkin R, Omar C. NeuroUnit: Validation Tests for Neuroscience Models. Frontiers in

1495        Neuroinformatics. 2013. doi:10.3389/conf.fninf.2013.09.00013

1496   20.    Appukuttan S, Garcia PE, Davison AP. MorphoUnit. Zenodo; 2020.

1497        doi:https://doi.org/10.5281/zenodo.3862936

1498   21.    Appukuttan S, Dainauskas J, Davison AP. SynapseUnit. Zenodo; 2020.

1499        doi:http://doi.org/10.5281/zenodo.3862944

1500   22.    Garcia PE, Davison AP. HippoNetworkUnit. Zenodo; 2020.

1501        doi:http://doi.org/10.5281/zenodo.3886484

1502   23.    Sharma BL, Davison AP. CerebUnit. Zenodo; 2020.

1503        doi:http://doi.org/10.5281/zenodo.3885673

1504   24.    Bianchi D, Marasco A, Limongiello A, Marchetti C, Marie H, Tirozzi B, et al. On the

1505        mechanisms underlying the depolarization block in the spiking dynamics of CA1

1506        pyramidal neurons. Journal of Computational Neuroscience. 2012;33: 207–225.

1507        doi:10.1007/s10827-012-0383-y

1508   25.    Spruston N, Schiller Y, Stuart G, Sakmann B. Activity-Dependent Action Potential

1509        Invasion and Calcium Influx into Hippocampal CA1 Dendrites. Science. 1995;268: 297–

1510        300. doi:10.1126/science.7716524

1511   26.    Golding NL, Kath WL, Spruston N. Dichotomy of action-potential backpropagation in

1512        CA1 pyramidal neuron dendrites. Journal of Neurophysiology. 2001;86: 2998–3010.

1513        doi:10.1152/jn.2001.86.6.2998

67

1514   27.   Gasparini S, Losonczy A, Chen X, Johnston D, Magee JC. Associative pairing enhances

1515        action potential back-propagation in radial oblique branches of CA1 pyramidal neurons.

1516        Journal of Physiology. 2007;580: 787–800. doi:10.1113/jphysiol.2006.121343

1517   28.   Magee JC, Cook EP. Somatic EPSP amplitude is independent of synapse location in

1518        hippocampal   pyramidal   neurons.   Nature   Neuroscience.   2000;3:   895–903.

1519        doi:10.1038/78800

1520   29.   Gasparini S, Magee JC. State-dependent dendritic computation in hippocampal CA1

1521        pyramidal   neurons.   Journal   of   Neuroscience.   2006;26:   2088–2100.

1522        doi:10.1523/JNEUROSCI.4428-05.2006

1523   30.   Ariav G, Polsky A, Schiller J. Submillisecond precision of the input-output

1524        transformation function mediated by fast sodium dendritic spikes in basal dendrites of

1525        CA1   pyramidal   neurons.   Journal   of   Neuroscience.   2003;23:   7750–7758.

1526        doi:10.1523/jneurosci.23-21-07750.2003

1527   31.   Takahashi H, Magee JC. Pathway Interactions and Synaptic Plasticity in the Dendritic

1528        Tuft   Regions   of   CA1   Pyramidal   Neurons.   Neuron.   2009;62:   102–111.

1529        doi:10.1016/j.neuron.2009.03.007

1530   32.   Losonczy A, Magee JC. Integrative Properties of Radial Oblique Dendrites in

1531        Hippocampal   CA1   Pyramidal   Neurons.   Neuron.   2006;50:   291–307.

1532        doi:10.1016/j.neuron.2006.03.016

1533   33.   Hines ML, Carnevale NT. The NEURON Simulation Environment. Neural

1534        Computation. 1997;9: 1179–1209. doi:https://doi.org/10.1162/neco.1997.9.6.1179

1535   34.   Druckmann S, Banitt Y, Gidon A, Schrümann F, Markram H, Segev I. A novel multiple

1536        objective optimization framework for constraining conductance-based neuron models by

1537        experimental   data.   Frontiers   in   Neuroscience.   2007;1:   7–18.

1538        doi:10.3389/neuro.01.1.1.001.2007

68

1539    35.    van Geit W, Moor R, Ranjan R, Roessert C, Riquelme L. Electrophys Feature Extraction Library. 2020. [cited 25. March 2020] GitHub repository [Internet] Available: https://github.com/BlueBrain/eFEL

1542    36.    Tripathy SJ, Savitskaya J, Burton SD, Urban NN, Gerkin RC. NeuroElectro: A window to the world's neuron electrophysiology data. Frontiers in Neuroinformatics. 2014;8: 1–11. doi:10.3389/fninf.2014.00040

1545    37.    Wheeler DW, White CM, Rees CL, Komendantov AO, Hamilton DJ, Ascoli GA. Hippocampome.org: A knowledge base of neuron types in the rodent hippocampus. eLife. 2015;4: 1–28. doi:10.7554/eLife.09960

1548    38.    Staff NP, Jung HY, Thiagarajan T, Yao M, Spruston N. Resting and active properties of pyramidal neurons in subiculum and CA1 of rat hippocampus. Journal of Neurophysiology. 2000;84: 2398–2408. doi:10.1152/jn.2000.84.5.2398

1551    39.    Dougherty KA, Islam T, Johnston D. Intrinsic excitability of CA1 pyramidal neurones from the rat dorsal and ventral hippocampus. Journal of Physiology. 2012;590: 5707–5722. doi:10.1113/jphysiol.2012.242693

1554    40.    Malik R, Dougherty KA, Parikh K, Byrne C, Johnston D. Mapping the electrophysiological and morphological properties of CA1 pyramidal neurons along the longitudinal hippocampal axis. Hippocampus. 2016;26: 341–361. doi:10.1002/hipo.22526

1558    41.    McDermott CM, Hardy MN, Bazan NG, Magee JC. Sleep deprivation-induced alterations in excitatory synaptic transmission in the CA1 region of the rat hippocampus. Journal of Physiology. 2006;570: 553–565. doi:10.1113/jphysiol.2005.093781

1561    42.    Graves AR, Moore SJ, Bloss EB, Mensh BD, Kath WL, Spruston N. Hippocampal Pyramidal Neurons Comprise Two Distinct Cell Types that Are Countermodulated by Metabotropic Receptors. Neuron. 2012;76: 776–789. doi:10.1016/j.neuron.2012.09.036

43. Golding NL, Mickus TJ, Katz Y, Kath WL, Spruston N. Factors mediating powerful voltage attenuation along CA1 pyramidal neuron dendrites. Journal of Physiology. 2005;568: 69–82. doi:10.1113/jphysiol.2005.086793

44. Bormann I. DigitizeIt: Digitizer software - digitize a scanned graph or chart into (x,y)-data. 2020. [cited 25. March 2020] [Internet] Available: https://www.digitizeit.de/

45. Jahr CE, Stevens CF. Voltage dependence of NMDA-activated macroscopic conductances predicted by single-channel kinetics. Journal of Neuroscience. 1990;10: 3178–3182. doi:10.1523/jneurosci.10-09-03178.1990

46. Hestrin S, Nicoll RA, Perkel DJ, Sah P. Analysis of excitatory synaptic action in pyramidal cells using whole-cell recording from rat hippocampal slices. Physiology. 1990; 203–225. doi:10.1113/jphysiol.1990.sp017980

47. Korinek M, Sedlacek M, Cais O, Dittert I, Vyklicky L. Temperature dependence of N-methyl-d-aspartate receptor channels and N-methyl-d-aspartate receptor excitatory postsynaptic currents. Neuroscience. 2010;165: 736–748. doi:10.1016/j.neuroscience.2009.10.058

48. Gevaert M, Kanari L, Palacios J, Zisis E, Coste B. NeuroM. 2020. [cited 25. March 2020] GitHub repository [Internet] Available: https://github.com/BlueBrain/NeuroM

49. Katz Y, Menon V, Nicholson DA, Geinisman Y, Kath WL, Spruston N. Synapse Distribution Suggests a Two-Stage Model of Dendritic Integration in CA1 Pyramidal Neurons. Neuron. 2009;63: 171–177. doi:10.1016/j.neuron.2009.06.023

50. Migliore M, de Blasi I, Tegolo D, Migliore R. A modeling study suggesting how a reduction in the context-dependent input on CA1 pyramidal neurons could generate schizophrenic behavior. Neural Networks. 2011;24: 552–559. doi:10.1016/j.neunet.2011.01.001

70

1588    51.    Poirazi P, Brannon T, Mel BW. Arithmetic of subthreshold synaptic summation in a

1589            model CA1 pyramidal cell. Neuron. 2003;37: 977–987. doi:10.1016/S0896-

1590            6273(03)00148-X

1591    52.    Shah MM, Migliore M, Valencia I, Cooper EC, Brown DA. Functional significance of

1592            axonal Kv7 channels in hippocampal pyramidal neurons. Proceedings of the National

1593            Academy of Sciences of the United States of America. 2008;105: 7869–7874.

1594            doi:10.1073/pnas.0802805105

1595    53.    Gómez González JF, Mel BW, Poirazi P. Distinguishing linear vs. non-linear

1596            integration in CA1 radial oblique dendrites: It's about time. Frontiers in Computational

1597            Neuroscience. 2011;5: 1–12. doi:10.3389/fncom.2011.00044

1598    54.    Jeffrey M. Perkel. Why Jupyter is data scientists' computational notebook of choice.

1599            Nature. 2018; 5–6. Available: https://www.nature.com/articles/d41586-018-07196-1

1600    55.    Ecker A, Romani A, Sáray S, Káli S, Migliore M, Mercer A, et al. Data-driven

1601            integration of hippocampal CA1 synapse physiology in silico. Hippocampus. 2020. doi:

1602            Forthcoming

1603    56.    Amunts K, Ebell C, Muller J, Telefont M, Knoll A, Lippert T. The Human Brain

1604            Project: Creating a European Research Infrastructure to Decode the Human Brain.

1605            Neuron. 2016;92: 574–581. doi:10.1016/j.neuron.2016.10.046

1606    57.    Fragnaud H, Gonin J, Duperrier J, Legouee E, Davison AP, Appukuttan S. hbp-

1607            validation-framework. Zenodo; 2020. doi:http://doi.org/10.5281/zenodo.3888123

1608    58.    Gleeson P, Crook S, Cannon RC, Hines ML, Billings GO, Farinella M, et al. NeuroML:

1609            A language for describing data driven models of neurons and networks with a high

1610            degree of biological detail. PLoS Computational Biology. 2010;6: 1–19.

1611            doi:10.1371/journal.pcbi.1000815

71

1612    59.    Raikov I. NineML – a description language for spiking neuron network modeling: the

1613           abstraction layer. BMC Neuroscience. 2010;11: 2202. doi:10.1186/1471-2202-11-s1-

1614           p66

1615    60.    Davison AP, Brüderle D, Eppler J, Kremkow J, Muller E, Pecevski D, et al. PyNN: A

1616           common interface for neuronal network simulators. Frontiers in Neuroinformatics.

1617           2009;2: 1–10. doi:10.3389/neuro.11.011.2008

1618    61.    Dai K, Hernando J, Billeh YN, Gratiy SL, Planas J, Davison A, et al. The SONATA

1619           Data Format for Efficient Description of Large-Scale Network Models. PLoS

1620           Computational Biology. 2019;16: 1–24. doi:10.2139/ssrn.3387685

1621    62.    Stimberg M, Brette R, Goodman DFM. Brian 2, an intuitive and efficient neural

1622           simulator. eLife. 2019;8: 1–41. doi:10.7554/eLife.47314

1623

## Supporting information

1625

**S1 Appendix. Example of running the SomaticFeaturesTest of HippoUnit using a Jupyter notebook**

1628

1629    As the first step HippoUnit's test classes and `ModelLoader` class, along with a few

1630    additional Python packages must be imported:

```
1. from __future__ import print_function #needed only in Python 2
2. % matplotlib inline
3.
4. from hippounit.utils import ModelLoader
5. from hippounit import tests
```

```
6.
7.  import pkg_resources
8.  import json
9.  import collections
10. import numpy
```

1641    Then the path to external mechanisms used by the Neuron implementation of the model

1642    (NMODL files) needs to be provided, which will be an argument to the `ModelLoader` class,

1643    so that the NMODL files can be compiled when the `ModelLoader` class is instantiated (if they

1644    are not compiled yet). Next, the variables related to the model are set. The initial voltage

1645    (`v_init`) and temperature (`celsius`) values specific to the model need to be set; otherwise,

1646    the default value in the corresponding capability method of the `ModelLoader` will be used.

1647    Setting the `cvode_active` boolean parameter to `True` or `False`, the user can decide to run

1648    the simulations using variable or fixed time step, respectively.

```
11. # path to NMODL files
12. mod_files_path = "/home/saray/published_models/Ca1_Bianchi_2012/experiment/"
13.
14. #all the outputs will be saved here. It will be an argument to the test.
15. base_directory = '/mnt/csoport31-
    2/Modellezo_csapat/Sara/published_models_validation_results/'
16.
17. #Load cell model
18. model = ModelLoader(mod_files_path = mod_files_path )
19.
20. # outputs will be saved in subfolders named like this:
21. model.name="Bianchi_et_al_2012"
22.
23. # path to hoc file
24. # the model must not display any GUI!!
25. model.hocpath = "/home/saray/published_models/Ca1_Bianchi_2012/experiment/main_model.hoc"
26.
27. # If the hoc file doesn't contain a template, this must be None (the default value is None)
```

73

```
28. model.template_name = None

29.

30. # model.SomaSecList_name should be None, if there is no Section List in the model for the soma
    , or if the name of the soma section is given by setting model.soma (the default value is None
    )

31. model.SomaSecList_name = None

32. # if the soma is not in a section list or to use a specific somatic section, add its name here
    :

33. model.soma = 'soma[0]'

34.

35. # For the PSP Attenuation Test, and Back-
    propagating AP Test a section list containing the trunk sections is needed

36. model.TrunkSecList_name = 'apical_trunk_list'

37. # For the Oblique Integration Test a section list containing the oblique dendritic sections is
     needed

38. model.ObliqueSecList_name = 'oblique_dendrites'

39. # This will be argument to those tests, where dendritic locatins are selected according to
    distances. If not set, the end of the above given soma section will be used as reference point
    for distance determination

40. trunk_origin = ['soma[0]', 1]

41.

42. model.v_init = -70

43. model.celsius = 34

44.

45. # It is possible to run the simulations using variable time step (default for this is False)

46. model.cvode_active = True
```

The target experimental data and the configuration file are loaded from the JSON files

to the *observation* and *config* dictionaries, which are arguments of the test class:

```
47. # Load target data

48. with open('/home/saray/target_features/feat_CA1_pyr_cACpyr_more_features.json') as f:

49.     observation  = json.load(f, object_pairs_hook=collections.OrderedDict)

50.

51. # Load stimuli file

52. ttype = "CA1_pyr_cACpyr"
```

74

```
53.
54. stim_file = pkg_resources.resource_filename("hippounit", "tests/stimuli/somafeat_stim/stim_" +
       ttype + ".json")
55. with open(stim_file, 'r') as f:
56.     config = json.load(f, object_pairs_hook=collections.OrderedDict)
```

Then the test class is instantiated, and its `judge()` function (inherited from SciUnit) is called to run the test. The number of parallel processes to be used can be controlled by the user by setting the `test.npool` parameter.

```
57. # Instantiate test class
58. test = tests.SomaticFeaturesTest(observation=observation, config=config, force_run=False, show
       _plot=True, save_all = True, base_directory=base_directory)
59.
60. # test.specify_data_set is added to the name of the subdirectory (somaticfeat), so test runs u
       sing different data sets can be saved into different directories
61. test.specify_data_set = 'UCL_data'
62.
63. # Number of parallel processes
64. test.npool = 30
65. #Run the test
66. score = test.judge(model)
67. #Summarize and print the score achieved by the model on the test using SciUnit's summarize fun
       ction
68. score.summarize()
```

For further details on how to run the different tests of HippoUnit for the different models, see the Jupyter Notebooks available here: https://github.com/KaliLab/HippoUnit_demo/tree/master/jupyter_notebooks.

75