# Automated curation of CNMF-E-extracted ROI spatial footprints and calcium traces using open-source AutoML tools

**Tran, LM[1,2,3], Mocle AJ[1,2], Ramsaran AI[1,4], Jacob AD[1,4], Frankland PW[1,2,4,5,6], Josselyn SA[1,2,4,5,7]**

[1]Hospital for Sick Children, Neurosciences and Mental Health, Toronto, ON, Canada

[2]Department of Physiology, University of Toronto, Toronto, ON, Canada

[3]Postgraduate Affiliates Program, Vector Institute, Toronto, ON, Canada

[4]Department of Psychology, University of Toronto, Toronto, ON, Canada

[5]Institute of Medical Sciences, University of Toronto, Toronto, ON, Canada

[6]Child & Brain Development Program, Canadian Institute for Advanced Research (CIFAR), Toronto, ON, Canada

[7]Brain, Mind & Consciousness Program, Canadian Institute for Advanced Research (CIFAR), Toronto, ON, Canada

**\* Correspondence:**
Sheena Josselyn
sheena.josselyn@sickkids.ca

**Keywords: calcium imaging, open-source, machine learning, microendoscopy**

## Abstract

In vivo 1-photon calcium imaging is an increasingly prevalent method in behavioural neuroscience. Numerous analysis pipelines have been developed to improve the reliability and scalability of pre-processing and ROI extraction for these large calcium imaging datasets. Despite these advancements in pre-processing methods, manual curation of the extracted spatial footprints and calcium traces of neurons remains important for quality control. Here, we propose an additional semi-automated curation step for sorting spatial footprints and calcium traces from putative neurons extracted using the popular CNMF-E algorithm. We used the automated machine learning tools TPOT and AutoSklearn to generate classifiers to curate the extracted ROIs trained on a subset of human-labeled data. AutoSklearn produced the best performing classifier, achieving an F1 score > 92% on the ground truth test dataset. This automated approach is a useful strategy for filtering ROIs with relatively few labeled data points, and can be easily added to pre-existing pipelines currently using CNMF-E for ROI extraction.

## 1    Introduction

Advances in one-photon (1p) miniaturized fluorescence microscopy in terms of utility, cost, and ease-of-use have increased the accessibility and popularity of *in vivo* calcium imaging in freely behaving rodents (Cai et al., 2016; Ghosh et al., 2011; Hamel et al., 2015; Jacob et al., 2018). Consequently, researchers are able to track the activity of neuronal populations across days, weeks, or even months (Gonzalez et al., 2019; Rubin et al., 2015). Concurrent with the growing usage of 1p microendoscopy in neuroscience, there is an increasing demand for high-throughput software that accurately and efficiently processes the very large raw calcium imaging datasets now being produced. To address this challenge, a number of algorithms and analysis pipelines have been

39developed to automate the extraction of cells and calcium activity traces across time in a robust
40manner—a necessary step for downstream analyses (Pnevmatikakis, 2019).

41

42Motion correction, source extraction, and cell registration (across multiple recording sessions) are
43important steps in pre-processing raw 1p calcium imaging data. Source extraction, the task of
44identifying the locations and activity of neurons in the imaged field of view (FOV), is arguably the
45most challenging of these steps is arguably the most challenging of these steps, as evidenced by the
46number of different algorithms published with the aim of improving this critical step. Nevertheless,
47two main methods of source extraction have been widely adopted in the field: principal component
48analysis/independent component analysis (PCA/ICA) (Mukamel et al., 2009) and the more recent
49extended constrained non-negative matrix factorization for microendoscopic data (CNMF-E) (Zhou
50et al., 2018). CNMF-E explicitly models background signals present in 1p microendoscopic
51recordings, and, therefore results in more accurate signal detection from neurons compared to
52PCA/ICA (Zhou et al., 2018).

53

54Our lab has successfully applied CNMF-E to recordings from our open-source Compact Head-
55mounted Endoscope (CHEndoscope) in order to identify neuron locations (or spatial footprints) and
56extract their calcium activity traces from freely-behaving mice performing different behavioural
57tasks. CNMF-E has proven to be a reliable tool across multiple imaging sessions and experimental
58paradigms conducted in the lab with minimal parameter tuning in our hands (Jacob et al., 2018).
59However, like PCA/ICA, CNMF-E may still produce some false-positives in the output of detected
60cells (i.e., non-neuronal spatial footprints or calcium traces), which can be filtered out of the final
61dataset manually. We initially found success in filtering CNMF-E-extracted spatial footprints and
62traces by adding a manual curation step that involved visual inspection of each ROI and calcium
63trace (previously described in (Jacob et al., 2018)). While this type of manual curation can reduce the
64number of false-positives in CNMF-E's output, visual inspection of potentially tens of thousands of
65extracted cells can be time-consuming, and this method is not free from human error. Here, we
66propose an automated machine learning (AutoML) approach built on top of the CNMF-E algorithm's
67outputs to filter out potential false-positives. We implemented a semi-automated classification tool to
68limit the amount of manual curation required during pre-processing, without completely removing
69the ability to fine-tune the process with human-labeled datasets.

70

71The main outputs of CNMF-E's source extraction algorithm are: 1) the extracted calcium traces
72representing cellular activity, and, 2) the spatial footprint of putative neurons. As mentioned
73previously, manual curation of these outputs involves identifying both aberrant traces that do not
74have stable baseline fluorescence (Resendez et al., 2016), transient durations inconsistent with the
75expressed calcium indicator (e.g., GCAMP6f) (Badura et al., 2014), and/or spatial footprints that are
76not consistent with the shape and size of neurons in the brain region being recorded (Resendez et al.,
772016). We trained and validated our classifiers on a dataset of 14 000 manually curated spatial
78footprints and traces output from CNMF-E. The final model chosen was then used to automate the
79curation of ROIs from other recording sessions. From the two AutoML libraries, we chose the best
80performing model to train on the full training set to evaluate on the test set. We find our model can
81accurately predict whether a cell would be included or excluded at a rate of 92%.

2

82 The potential time savings of manually curating thousands of cells makes this approach a method 83 worth employing as part of a typical 1p calcium imaging pipeline. While our AutoML-based curation 84 pipeline was primarily developed to be used with CHEndoscope data, our model takes the output of 85 CNMF-E and as a result, allows this method to be readily applied to data acquired using other 1p 86 miniature endoscopes.

87 **2    Methods**
88 **Dataset preparation and pre-processing**
89

90 The dataset used for model training was acquired from multiple hippocampal CA1 recordings 91 captured across different mice and recording sessions using methods described in Jacob et al. 2018. 92 From these recordings, we used CNMF-E (Zhou et al., 2018) to extract spatial footprints and calcium 93 traces of 14 000 ROIs.  We then manually reviewed and labeled these ROIs as neuronal (included for 94 further analysis) or artefact (excluded from analysis). The labels were generated by two human expert 95 raters that inspected the calcium transients and spatial footprints based on previously reported 96 criteria:

97    1.  fast rise and slow decay of calcium transients with stable baseline fluorescence (Resendez et 98        al., 2016).

99    2.  calcium transient durations consistent with GCaMP6f (or appropriate GCaMP variant) 100        (Badura, Sun, Giovannucci, Lynch, & Wang, 2014).

101    3.  spatial footprints consistent with appropriate neuronal shape and size (Resendez et al., 2016).

102 Interrater agreement for the dataset was 87% across the two raters on a subset of the data (1073 103 putative ROIs extracted from CNMF-E) (Figure 1).

104 Spatial footprints consisted of the maximum projection of the identified cell from all frames in the 105 video. We found that location of the footprint in the FOV was not important in our labelling criteria 106 (compared to shape and size of footprint), we cropped the spatial footprints to remove empty space. 107 Each spatial footprint was reduced to an 80x80 pixel image centered on the peak intensity of the 108 footprint. Furthermore, recordings were of varying lengths, so all trace data was cropped at 500 109 frames (equivalent to 100s of recording at 5fps). The 2 dimensional footprints were reduced to a 1 110 dimensional vector (6400 pixels) and concatenated to the trace data.

111

112 We aggregated the labeled ROIs into a dataset split into training and test sets, which comprised 80% 113 (~11 000 ROIs) and 20% (~3 000 ROIs) of the data, respectively.

114

115 **Model optimization and selection**

116

117 We used two automated machine learning (AutoML) methods, TPOT (Olson et al., 2016; Olson & 118 Moore, 2019) and AutoSklearn (Feurer et al., 2019) that are based on the popular Python machine 119 learning toolbox, scikit-learn (Pedregosa et al. 2011) to select optimal classification models. While 120 other AutoML tools exist that may outperform the ones we chose (Truong et al., 2019), TPOT and 121 AutoSklearn are both free open-source, and easy to use, making them accessible for labs to 122 incorporate into their existing analysis pipelines.

123

124 The key advantage of AutoML tools such as TPOT and AutoSklearn is that they do the extensive 125 work of finding the best type(s) of data transformation and models to build a pipeline for classifying

126the input data, as well as the hyperparameters associated with these steps. TPOT is a tree-based
127optimization tool that builds and optimizes machine learning pipelines using genetic programming
128(Olson et al., 2016; Olson & Moore, 2019). TPOT generates pipelines of pre-processing steps and
129classification models in order to maximize classification performance while prioritizing simpler
130pipelines. AutoSklearn performs algorithm selection and hyperparameter tuning using Bayesian
131optimization, meta-learning and ensemble construction (Feurer et al., 2019) and as a result, the final
132classifier is an ensemble of many different model types and their associated hyperparameters. We
133primarily used default TPOT and AutoSklearn parameters, with a max evaluation time for a single
134pipeline of 10 minutes, and a total search time of 2 days.

135

136During training, we used 10-fold cross-validation using stratified folds that preserved the relative
137proportions of "include" and "exclude" labels (i.e., during each run of training, 9 of 10 folds were
138used for training, and the 10th fold was used to test the performance of the model). This process was
139repeated for all 10 folds, resulting in an averaged performance metric for the data. We optimized the
140models to maximize the F1 score, the harmonic average of precision and recall, where high precision
141indicates a low false positive rate, and high recall indicates a low false negative rate. In our dataset, a
142true positive is an extracted ROI that both the trained model and a "ground truth" human scorer
143define as suitable to be included for further analysis (i.e., it satisfies the three selection criteria listed
144above). A true negative is an extracted ROI that is excluded for further analysis by both the model
145and our ground truth scoring.

## 146**3    Results**

147To determine the efficacy of an AutoML approach for classification of CNMF-E extracted ROIs, we
148tested the ability of TPOT and AutoSklearn to build classifiers that can label the pre-processed spatial
149footprints and calcium traces of putative ROIs. Both TPOT and AutoSklearn were trained on 11 000
150labeled ROIs in the training set split into 10 folds for cross-validation, repeated 5 times. The best
151models obtained during training were used to determine the F1 score on the test set. Table 1 reports
152the performance of the best models obtained by TPOT and AutoSklearn across the training folds and
153on the test set.

154Next we assessed the transferability of the best classifier pipelines identified by TPOT and
155AutoSklearn using fewer samples. We used the top performing classifier pipelines and
156hyperparameters chosen by TPOT and AutoSklearn and trained the initialized pipelines using
157datasets of increasing ROI number. The training set size ranged from 150 to 10 000 ROIs. Using a
158change point analysis algorithm (PELT, Killick et al. 2012), we determined that AutoSklean and
159TPOT classifiers approached a maximal F1 score with 719 and 1000 labeled ROIs, respectively
160(Figure 2).  The pipelines found using our much larger labeled dataset may be easily incorporated
161into other pipelines with minimal computational effort to train and finetune on CNMF-E extracted
162ROIs from other 1p experiments, using fewer labeled ROIs.

163To examine the classifier performance in terms of false positives and false negatives, we created
164confusion matrices to visualize the rate of true positives, true negatives, false positives and false
165negatives from the TPOT and AutoSklearn predictions compared to ground truth. We found that the
166classifier built with Autosklearn (0.922 F1, Table 1) performs better in terms of both reducing false
167positives and false negatives (Figure 3).

168

169To further assess the nature of the classification errors, we looked at the class confidences or 170probabilities of the test set predictions from the trained TPOT and AutoSklearn models (Figure 4). 171Class probabilities indicate the classifier's certainty (using confidence score for TPOT and class 172probability for AutoSklearn) that a sample belongs to a particular class label. We tested whether 173mislabeled ROIs were also those which the classifiers expressed less confidence in classifying. We 174examined the size of the difference between certainty scores (true positives versus false positives, 175true negatives versus false negatives) in TPOT and AutoSklearn using Cohen's d (Cohen, 2013; 176Sawilowsky, 2009) (Table 2). The AutoSklearn classifier—which outperformed the TPOT classifier 177based on F1 scores— showed large differences in certainty scores when labeling ROIs as positive 178(d=1.36) or negative (d=2.34).  By contrast, the TPOT classifier was relatively less confident on both 179types of classification (positive d=0.63, negative d=1.68). In other words, the AutoSklearn classifier 180was more certain in applying labels to ROIs than was the TPOT classifier. This indicates that false 181negatives and false positives in the higher performing AutoSklearn classifier may arise from "edge-182cases" ROIs in the dataset which the classifier was not as certain about the label. In contrast, the 183poorer performance of the TPOT classifier may simply be due to poor generalization on the test set.

184To investigate the nature of the false positives and false negatives from the best TPOT and 185AutoSklearn models, we looked at the underlying spatial footprints and calcium traces for mislabeled 186ROIs from both AutoML tools (Figure 5). Representative examples of excluded ROIs from the 187ground truth dataset show that some cells may be excluded (i.e., true negatives) because of poor trace 188data and/or poor spatial footprints, which possibly represent non-neuronal imaging artefacts and/or 189ROIs representing areas of background fluorescence. While some false positives from AutoSklearn 190shared similar features with true negative ROIs, others were more ambiguous. Upon inspection, these 191ROIs sometimes were high-quality spatial footprints with poor-quality calcium traces, or vice-versa, 192or were composed of spatial footprints and calcium traces of true neuronal-origin mixed with 193additional non-neuronal noise. These examples represent "edge-cases" which may be difficult to 194judge even by a human rater.

## 195 4    Discussion

196Automated curation of ROIs provides a fast, accurate method for classifying neural data generated in 1971p calcium imaging experiments. We show that AutoML tools such as the open source TPOT and 198AutoSklearn packages provide an easy way to build effective classifiers for ROIs extracted from the 199widely used CNMF-E algorithm. Spatial footprints and calcium traces from CNMF-E can be used to 200train these models with minimal data preprocessing. Furthermore, it may be possible to apply the top 201performing classifiers generated from this work to other experimental datasets taken from different 2021p imaging setups, while requiring relatively few labeled samples. Other analyses pipelines such as 203MIN1PIPE (Lu et al. 2018) have been developed to improve source extraction by reducing false 204positive ROIs without increasing the rate of false negatives. However, given the more widespread 205adoption of CNMF-E, the approach described here prevents labs from having to adopt entirely new 206analysis pipelines. Our approach provides a balance between the need to manually review the output 207of CNMF-E ROIs to maximize the number of detected cells, while still allowing some further 208automation of the otherwise laborious curation process.

209An AutoML approach to reviewing these traces may be useful for curating traces from labeled data 210of the extracted ROIs from CNMF-E and can be implemented on top of pre-existing analysis 211pipelines without much need to adapt the software. However, there are a number of limitations to this 212approach. Firstly, we emphasize the automated aspect of this machine learning classifier approach 213and little need for hand-tuning, but we recognize that the best models still make errors. Cases in 214which the best performing classifier generated by AutoSklearn failed to detect true positives or true

215negatives were further reviewed and were typically seen to be edge cases where it may be difficult 216for a human reviewer to make a judgment. Similarly, we found that a second expert scorer looking at 217the same data may not make the exact same judgments on such edge cases (having an interrater 218reliability score of 87%). While the AutoML classifiers were trained on the data that had relatively 219little preprocessing beyond cropping and downsampling, future work could address whether feature 220engineering over the spatial footprints and trace data could further improve accuracy and reduce 221training time for model selection and hyperparameter tuning. Better curation of a training dataset for 222the models may help reduce ambiguous cases that make it difficult for a classifier to make accurate 223predictions.

224In conclusion, we present here a demonstration and benchmark of an AutoML approach for curation 225of CNMF-E extracted ROIs. The methods described here can provide a flexible, free open-source, 226and easy-to incorporate curation step for other researchers using CNMF-E for source extraction of 227their 1p datasets, while requiring few changes to their existing analysis pipelines.

## 228 5    References

229Badura, A., Sun, X. R., Giovannucci, A., Lynch, L. A., & Wang, S. S.-H. (2014). Fast calcium sensor 230proteins for monitoring neural activity. Neurophotonics, 1(2), 025008.

231Cai, D. J., Aharoni, D., Shuman, T., Shobe, J., Biane, J., Song, W., Wei, B., Veshkini, M., La-Vu, 232M., Lou, J., Flores, S. E., Kim, I., Sano, Y., Zhou, M., Baumgaertel, K., Lavi, A., Kamata, M., 233Tuszynski, M., Mayford, M., … Silva, A. J. (2016). A shared neural ensemble links distinct 234contextual memories encoded close in time. In Nature (Vol. 534, Issue 7605, pp. 115–118). 235https://doi.org/10.1038/nature17955

236Cohen, J. (2013). Statistical Power Analysis for the Behavioral Sciences. 237https://doi.org/10.4324/9780203771587

238Feurer, M., Klein, A., Eggensperger, K., Springenberg, J. T., Blum, M., & Hutter, F. (2019). Auto-239sklearn: Efficient and Robust Automated Machine Learning. In Automated Machine Learning (pp. 240113–134). https://doi.org/10.1007/978-3-030-05318-5_6

241Ghosh, K. K., Burns, L. D., Cocker, E. D., Nimmerjahn, A., Ziv, Y., Gamal, A. E., & Schnitzer, M. J. 242(2011). Miniaturized integration of a fluorescence microscope. Nature Methods, 8(10), 871–878.

243Gonzalez, W. G., Zhang, H., Harutyunyan, A., & Lois, C. (2019). Persistence of neuronal 244representations through time and damage in the hippocampus. Science, 365(6455), 821–825.

245Hamel, E. J. O., Grewe, B. F., Parker, J. G., & Schnitzer, M. J. (2015). Cellular level brain imaging 246in behaving mammals: an engineering approach. Neuron, 86(1), 140–159.

247Jacob, A. D., Ramsaran, A. I., Mocle, A. J., Tran, L. M., Yan, C., Frankland, P. W., & Josselyn, S. A. 248(2018). A Compact Head-Mounted Endoscope for In Vivo Calcium Imaging in Freely Behaving 249Mice. Current Protocols in Neuroscience / Editorial Board, Jacqueline N. Crawley ... [et Al.], 84(1), 250e51.

251Mukamel, E. A., Nimmerjahn, A., & Schnitzer, M. J. (2009). Automated analysis of cellular signals 252from large-scale calcium imaging data. Neuron, 63(6), 747–760.

253 Olson, R. S., Bartley, N., Urbanowicz, R. J., & Moore, J. H. (2016). Evaluation of a Tree-based
254 Pipeline Optimization Tool for Automating Data Science. In Proceedings of the 2016 on Genetic and
255 Evolutionary Computation Conference - GECCO '16. https://doi.org/10.1145/2908812.2908918

256 Olson, R. S., & Moore, J. H. (2019). TPOT: A Tree-Based Pipeline Optimization Tool for
257 Automating Machine Learning. In Automated Machine Learning (pp. 151–160).
258 https://doi.org/10.1007/978-3-030-05318-5_8

259 Pnevmatikakis, E. A. (2019). Analysis pipelines for calcium imaging data. In Current Opinion in
260 Neurobiology (Vol. 55, pp. 15–21). https://doi.org/10.1016/j.conb.2018.11.004

261 Resendez, S. L., Jennings, J. H., Ung, R. L., Namboodiri, V. M. K., Zhou, Z. C., Otis, J. M., Nomura,
262 H., McHenry, J. A., Kosyk, O., & Stuber, G. D. (2016). Visualization of cortical, subcortical and
263 deep brain neural circuit dynamics during naturalistic mammalian behavior with head-mounted
264 microscopes and chronically implanted lenses. Nature Protocols, 11(3), 566–597.

265 Rubin, A., Geva, N., Sheintuch, L., & Ziv, Y. (2015). Hippocampal ensemble dynamics timestamp
266 events in long-term memory. eLife, 4. https://doi.org/10.7554/eLife.12247

267 Sawilowsky, S. S. (2009). New Effect Size Rules of Thumb. In Journal of Modern Applied Statistical
268 Methods (Vol. 8, Issue 2, pp. 597–599). https://doi.org/10.22237/jmasm/1257035100

269 Truong, A., Walters, A., Goodsitt, J., Hines, K., Bayan Bruss, C., & Farivar, R. (2019). Towards
270 Automated Machine Learning: Evaluation and Comparison of AutoML Approaches and Tools. In
271 2019 IEEE 31st International Conference on Tools with Artificial Intelligence (ICTAI).
272 https://doi.org/10.1109/ictai.2019.00209

273 Zhou, P., Resendez, S. L., Rodriguez-Romaguera, J., Jimenez, J. C., Neufeld, S. Q., Giovannucci, A.,
274 Friedrich, J., Pnevmatikakis, E. A., Stuber, G. D., Hen, R., Kheirbek, M. A., Sabatini, B. L., Kass, R.
275 E., & Paninski, L. (2018). Efficient and accurate extraction of in vivo calcium signals from
276 microendoscopic video data. eLife, 7. https://doi.org/10.7554/eLife.28728

## 277 6    Conflict of Interest

## 280 7    Author Contributions

281 LT, PF, SJ contributed to the study design. AJ designed and constructed the CHEndoscopes. AR and
282 AM conducted surgeries, behavior experiments, CNMF-E processing, and manual ROI labelling. LT
283 performed all analyses using automated machine learning pipelines. LT, AM, AR, AJ performed the
284 statistical analyses and wrote the paper. All authors discusses and commented on the manuscript.

## 285 8    Funding

## 293 9    Acknowledgments

## 296    Data Availability Statement

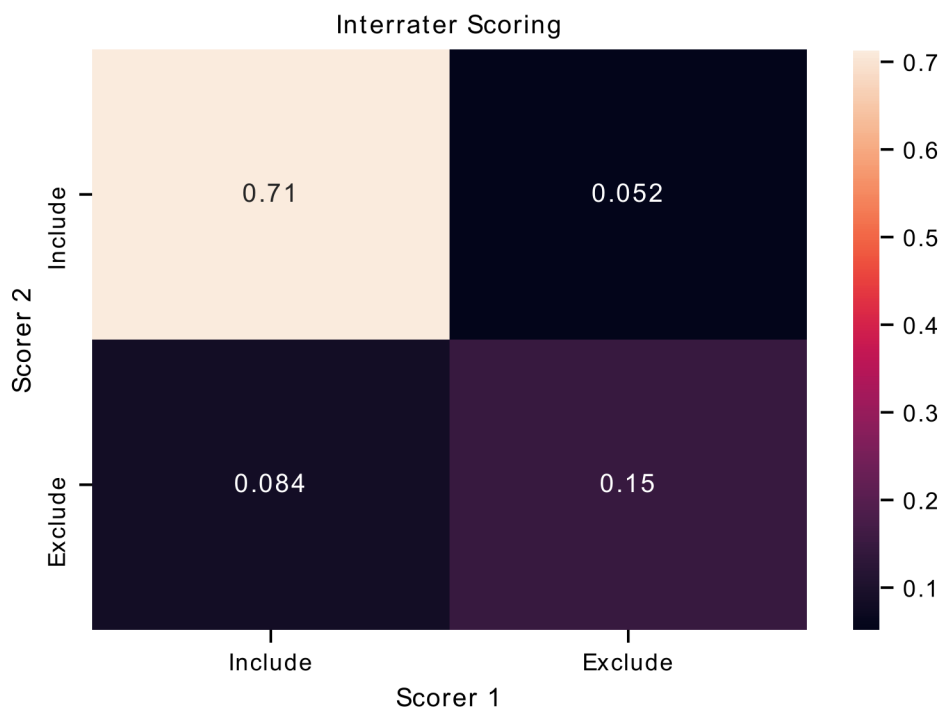297 The datasets and code generated for this study can be found in the cnmfe-reviewer GitHub repository
298 [https://github.com/jf-lab/cnmfe-reviewer].

299

**AutoML for 1p ROI Curation**

300**Figures**

301**Figure 1. Interrater agreement of ROI labels.** A confusion matrix comparing the manually 302reviewed labels (include of exclude) for putative ROIs extracted from CNMF-E determined by two 303different raters. Each cell of the matrix is annotated with the proportion of ROIs.
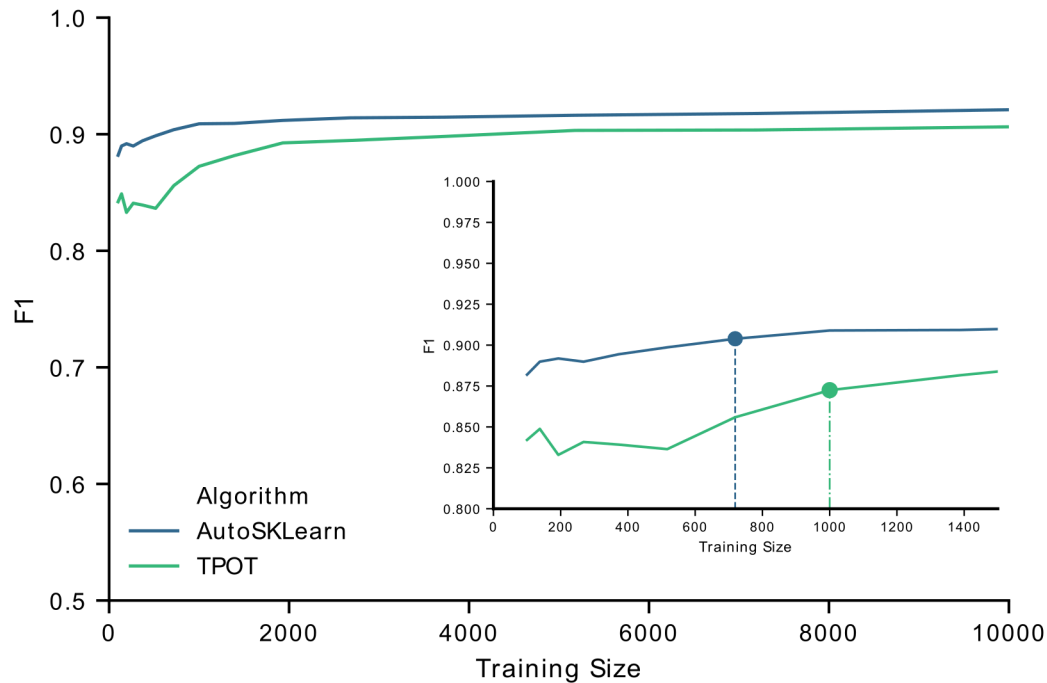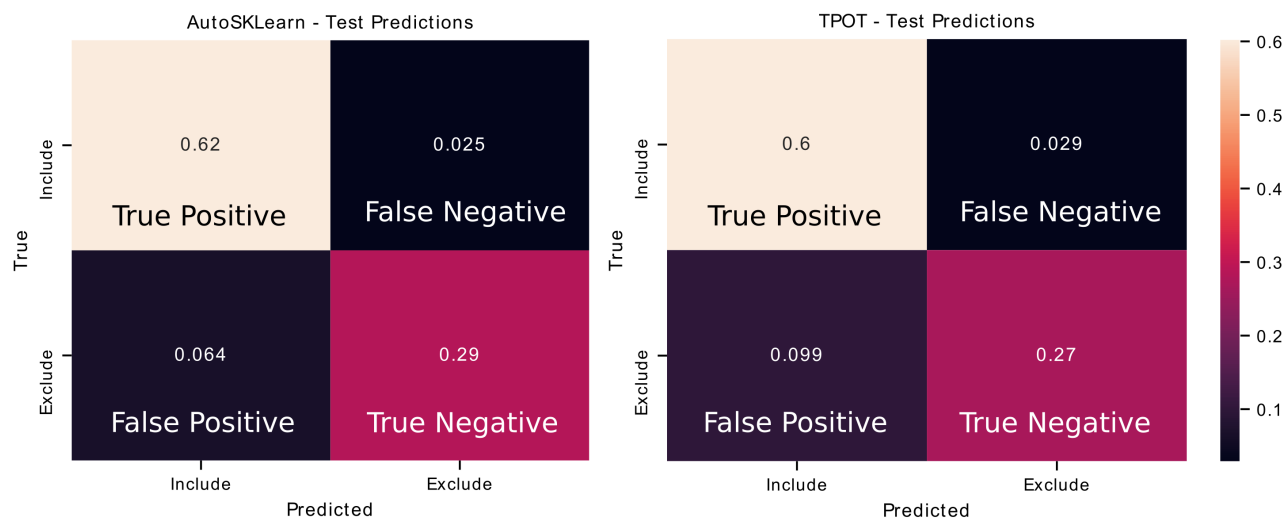
304

305



9

306 **Figure 2. F1 score performance with increasing training size.** Graph of the F1 test scores versus
307 the number of training samples used to train the best models output by AutoSklearn (blue) and TPOT
308 (green). (Inset) A graph of the same plot with a smaller range of training sizes and the change point is
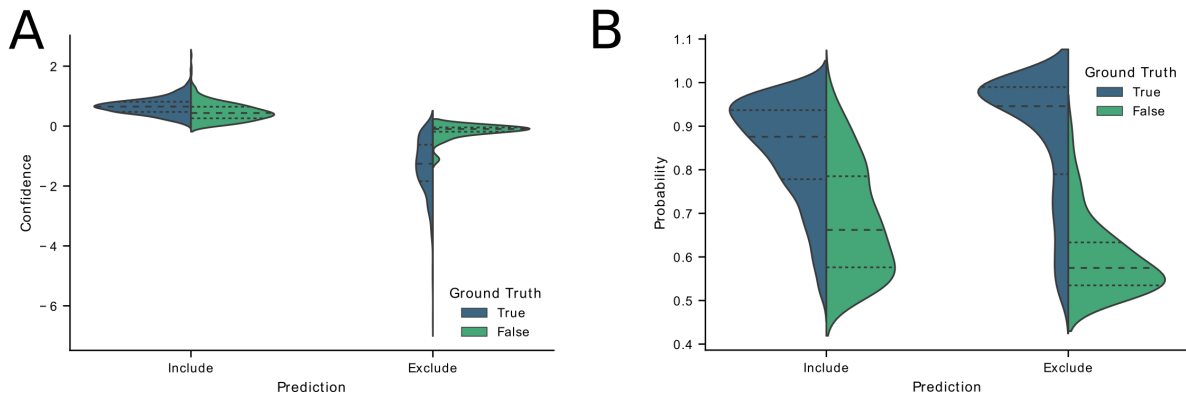309 marked on each algorithm type.

310

311

312 **Figure 3. Confusion matrices of AutoML tools: TPOT and AutoSklearn.** Each cell in the matrix
313 is annotated with the proportion of ROIs labeled as Include or Exclude according to the predicted and
314 true labels. Colors indicate the relative proportions of the labels where lower proportions are darker
315 in color and higher proportions are lighter in color. The confusion matrices were made from
316 predictions on the test set from the best models output by TPOT (left) and AutoSklearn (right).
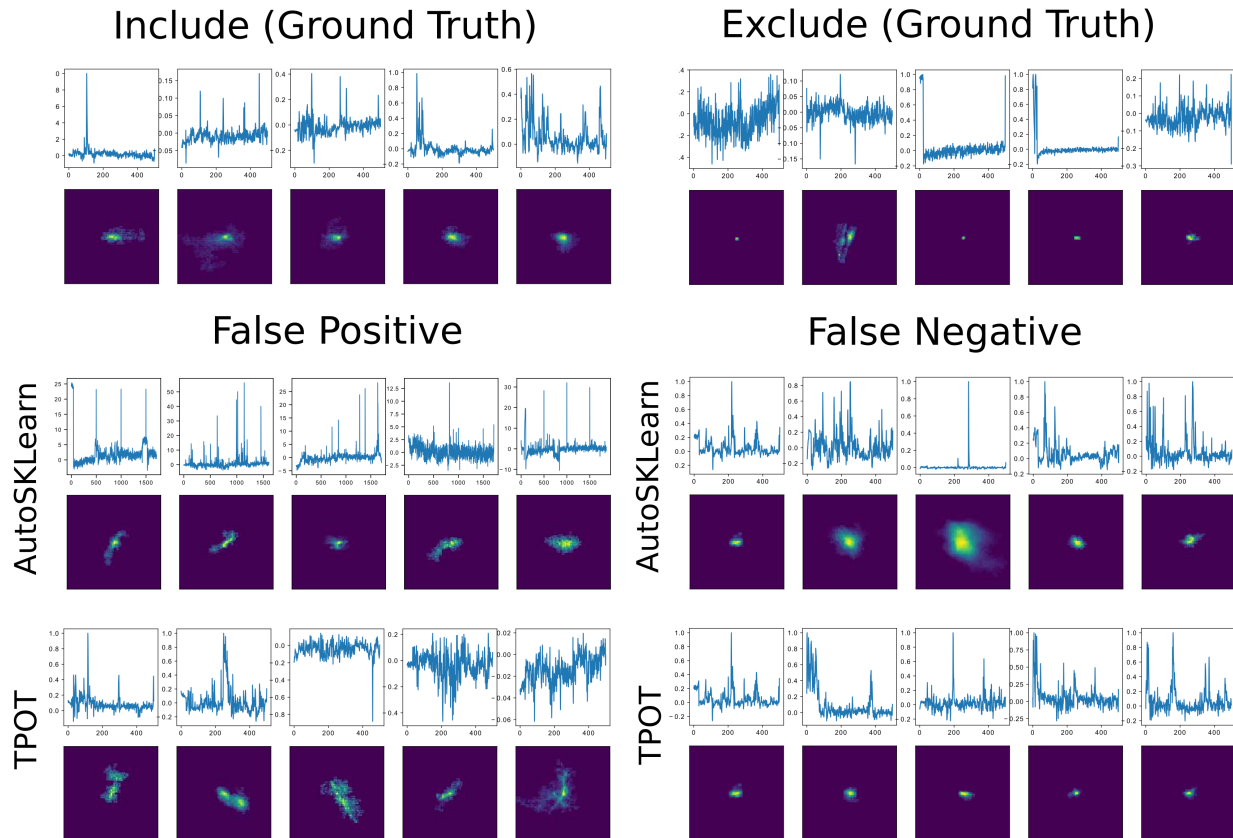


318

319 **Figure 4. Classifier confidence (TPOT) and class probabilities (AutoSklearn) on predicted false**
320 **positives and false negatives.** Violin plots of the distribution of (A) TPOT classifier confidence or
321 (B) AutoSklearn class probabilities on predicted ROI labels (Include or Exclude) in the test set. Each
322 half of the violin plot is the distribution of values for labels that were correct (True, left/blue) or
323 incorrect (False, right/green) based on the ground truth labels.



325

326 **Figure 5. Representative false positives and negatives compared to ground truth ROIs.** Example
327 calcium traces (top) and spatial footprints (bottom) from ground truth ROIs labeled as Include (left)
328 or Exclude (right). Example calcium traces (top) and spatial footprints (bottom) of false positive and
329 false negative ROIs predicted from the AutoSklearn (middle row) or TPOT (bottom row) classifiers.

331 **Tables**

332 **Table 1. Mean F1 scores of AutoML methods on training cross-validation and final F1 scores**
333 **on the test set.**

334

| | TPOT | | AutoSklearn | |
|---|---|---|---|---|
| | Training CV (10-fold, 5x) | Test | Training CV (10-fold, 5x) | Test |
| F1 Score | 0.906 | 0.904 | 0.917 | 0.922 |

335

336 **Table 2. Cohen's *d* of certainty scores between predicted labels that were correct or incorrect**
337 **compared to ground truth in TPOT or AutoSklearn.**

338

| | Include (Positive) | | Exclude (Negative) | |
|---|---|---|---|---|
| | TPOT | AutoSklearn | TPOT | AutoSklearn |
| Cohen's *d* | 0.63 | 1.36 | 1.68 | 2.34 |

339

14