

# Convpaint - Interactive pixel classification using pretrained neural networks

Lucien Hinderling  , Guillaume Witz  , Roman Schwob , Ana Stojiljković  , Maciej Dobrzański  , Mykhailo Vladymyrov  , Joël Frei<sup>1</sup>, Benjamin Grädel  , Agne Frismantiene  , and Olivier Pertz  

<sup>1</sup>Institute of Cell Biology, University of Bern, Baltzerstrasse 4, 3012 Bern, Switzerland

<sup>2</sup>Graduate School for Cellular and Biomedical Sciences, University of Bern, Switzerland

<sup>3</sup>Data Science Lab, University of Bern, Sidlerstrasse 5, 3012 Bern, Switzerland

We develop Convpaint, a universal computational framework for interactive pixel classification. Convpaint utilizes pretrained convolutional neural networks (CNNs) or vision transformers (ViTs) for feature extraction and enables easy segmentation across a wide variety of tasks. Available within the Python-based napari software ecosystem, Convpaint integrates seamlessly with other plugins into image processing pipelines, which we demonstrate with three workflows across different data modalities.

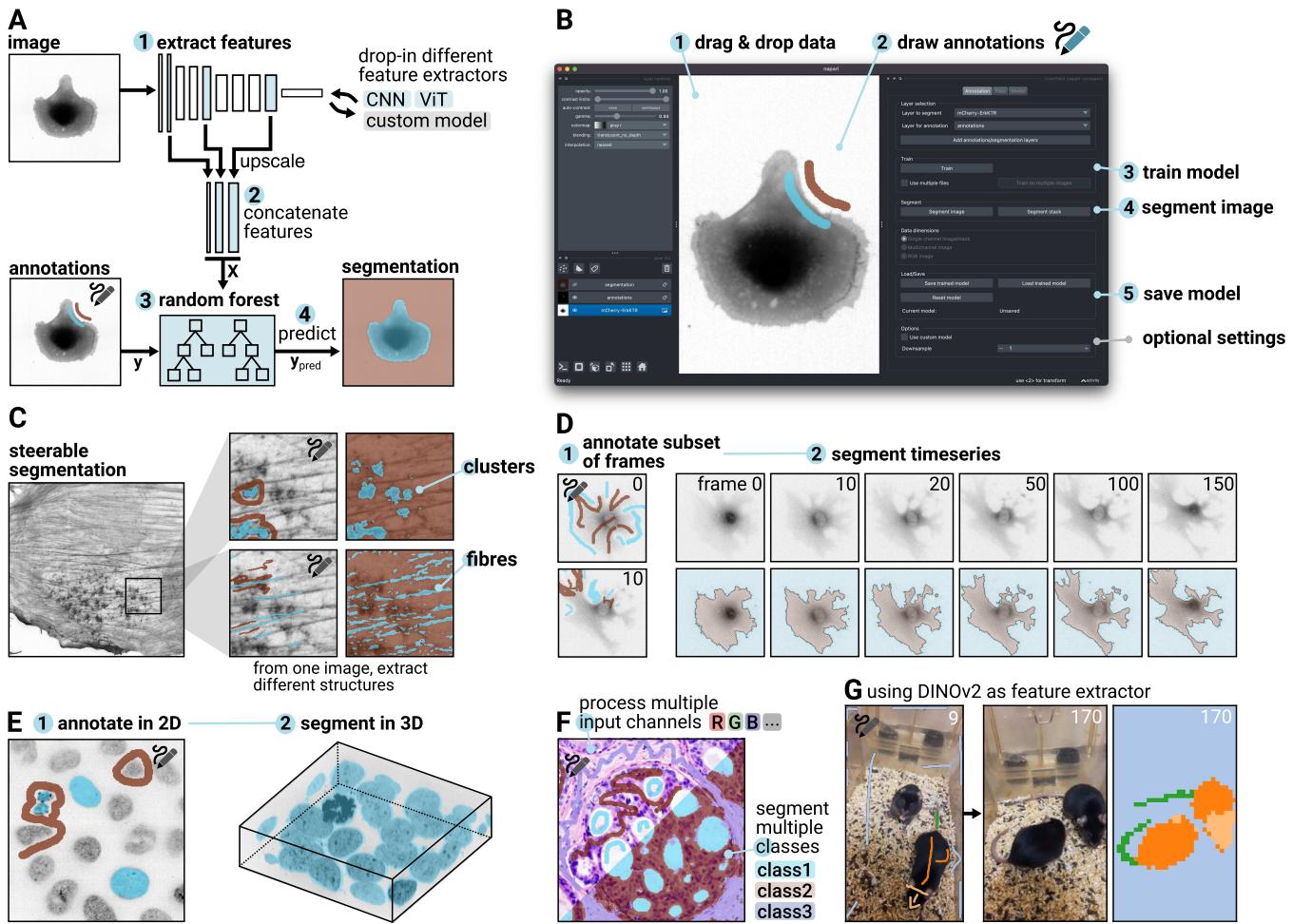
image analysis | pixel classification | multi-dimensional data

Correspondence: [lucien.hinderling@unibe.ch](mailto:lucien.hinderling@unibe.ch), [olivier.pertz@unibe.ch](mailto:olivier.pertz@unibe.ch)

1 Many bioimage analysis pipelines start with a segmentation step. While deep learning methods (DL) offer high  
2 classification accuracy, they require extensive ground truth  
3 annotation data and dedicated hardware for training. Even  
4 foundation models, trained on more diverse data and ex-  
5 pected to generalize to new applications without retrain-  
6 ing, in practice still need retraining for basic research  
7 purposes [1–3]. In contrast, machine learning (ML) ap-  
8 proaches using small models that can be trained interac-  
9 tively with sparse annotations, have proven to be highly  
10 effective (ilastik, Trainable Weka, Qupath, APOC) [4–7].  
11 These approaches traditionally rely on hand-crafted filter  
12 banks to extract image features and train an ML model  
13 from sparse annotations and corresponding features, to  
14 predict the class of each pixel in the rest of the image  
15 or new images. While these models are quick to train  
16 and hand-crafted filter banks effectively describe texture  
17 or local image structures [8], breakthrough performance  
18 in capturing semantically meaningful information from im-  
19 ages has been achieved through automatically learned fil-  
20 ter banks, specifically convolutional filters in DL models [9].  
21 Convpaint builds on this work, striking a balance between  
22 training speed, accuracy, and steerability by combining ML  
23 models that are fast to train with the power of DL. Instead  
24 of training a DL model from scratch, Convpaint extracts  
25 features from pretrained models like convolutional neu-  
26 ral networks (CNN) or vision transformers (ViT) and uses  
27 them to train a random forest classifier. Convpaint's de-  
28 sign is modular, allowing the feature extractor to be easily  
29 replaced by any algorithm that returns local features from  
30 an input image (fig. 1A).  
31 Providing a graphical user interface within the napari  
32 ecosystem, Convpaint offers a straightforward way for re-  
33 searchers to repurpose pretrained DL models for their

35 specific tasks without requiring coding or ML expertise  
36 (fig. 1B). Unlike neural networks trained to detect spe-  
37 cific structures (e.g., spots, fibers), users can guide the  
38 Convpaint model by drawing sparse labels on regions of  
39 interest (fig. 1C). This interactive process, which can be  
40 completed in seconds, allows for iterative cycles of anno-  
41 tation and evaluation, rapidly improving the quality of seg-  
42 mentation results. Convpaint seamlessly handles multi-  
43 dimensional data, making it suitable for segmenting time  
44 series and 3D data (fig. 1D, E). It can be trained on an  
45 arbitrary number of input channels and output classes (fig.  
46 1F, shown on synthetic data in S1). When coupled with  
47 different feature extraction models, Convpaint can be ap-  
48 plied to bioimages across scales, from subcellular to cel-  
49 lular structures to animals (fig. 1G, S2). For experienced  
50 users, Python APIs are available, allowing them to pro-  
51 grammatically control Convpaint. Convpaint incorporates  
52 several architectural optimizations that enhance training  
53 and prediction efficiency, setting it apart from similar soft-  
54 ware: it extracts crops around annotated pixels, minimiz-  
55 ing unnecessary processing of entire images. It uses tiled  
56 parallel processing with appropriate padding for large im-  
57 ages and stacks. Additionally, its integration with Dask al-  
58 lows for handling larger-than-memory files. The modular  
59 design makes customization easy, enabling users to in-  
60 tegrate new feature extractors via simple functions, while  
61 Convpaint manages the user interface, classifier training,  
62 data management, and parallelization. The software is in-  
63 teroperable with a wide range of napari plugins, enabling  
64 complex image analysis workflows within a single software  
65 ecosystem without coding. We demonstrate this in three  
66 workflows.

67 **Workflow 1** highlights Convpaint's capability to work with  
68 multichannel data. Imaging mass cytometry (IMC), spatial  
69 transcriptomics, or multiplexed immunofluorescence imag-  
70 ing, can image numerous biomarkers in the same sample.  
71 This provides a wealth of information posing new chal-  
72 lenges for data analysis, particularly for interactive data  
73 exploration. Fig. 2A-D shows an exemplary use case on  
74 a 43-channel IMC dataset [10]. The data can be interac-  
75 tively loaded and browsed using napari-imc [11]. Instead  
76 of exporting the data for pixel classification in external soft-  
77 ware like ilastik as demonstrated in a previous study [12],  
78 pixel classification can be performed directly in napari us-  
79 ing Convpaint. Here, we segment vein and surrounding



**Fig. 1. Overview of the Convpaint algorithm, user interface, and capabilities.** **A** Convpaint architecture: 1) Features are extracted from multiple scalings of the input image using a pretrained neural network, 2) upscaled, and concatenated. 3) A random forest model, trained on sparse annotations, predicts the class for each pixel. Different feature extractor models can be used. **B** User interface in napari: 1) Supports various input formats. 2) Annotations drawn using the labels layer. 3) Single-click model training. 4) Single-click image segmentation with results displayed in a labels layer. 5) Model saving for future use. Advanced settings for custom models and normalization available in additional tabs. **C** Interactive adjustment of extracted structures based on annotations. **D** Segments time-series data across all frames with a single click, enabling immediate playback. **E** Use napari's visualization to verify 3D segmentation results. **F** Adapts to any number of input channels and output classes. **G** DINOv2 as a feature extractor allows the segmentation of macroscopic objects and scenes. Full data shown in fig. S5, suppl. movie M3.

80 tissue regions and identify markers that are differentially  
81 expressed between the two classes.

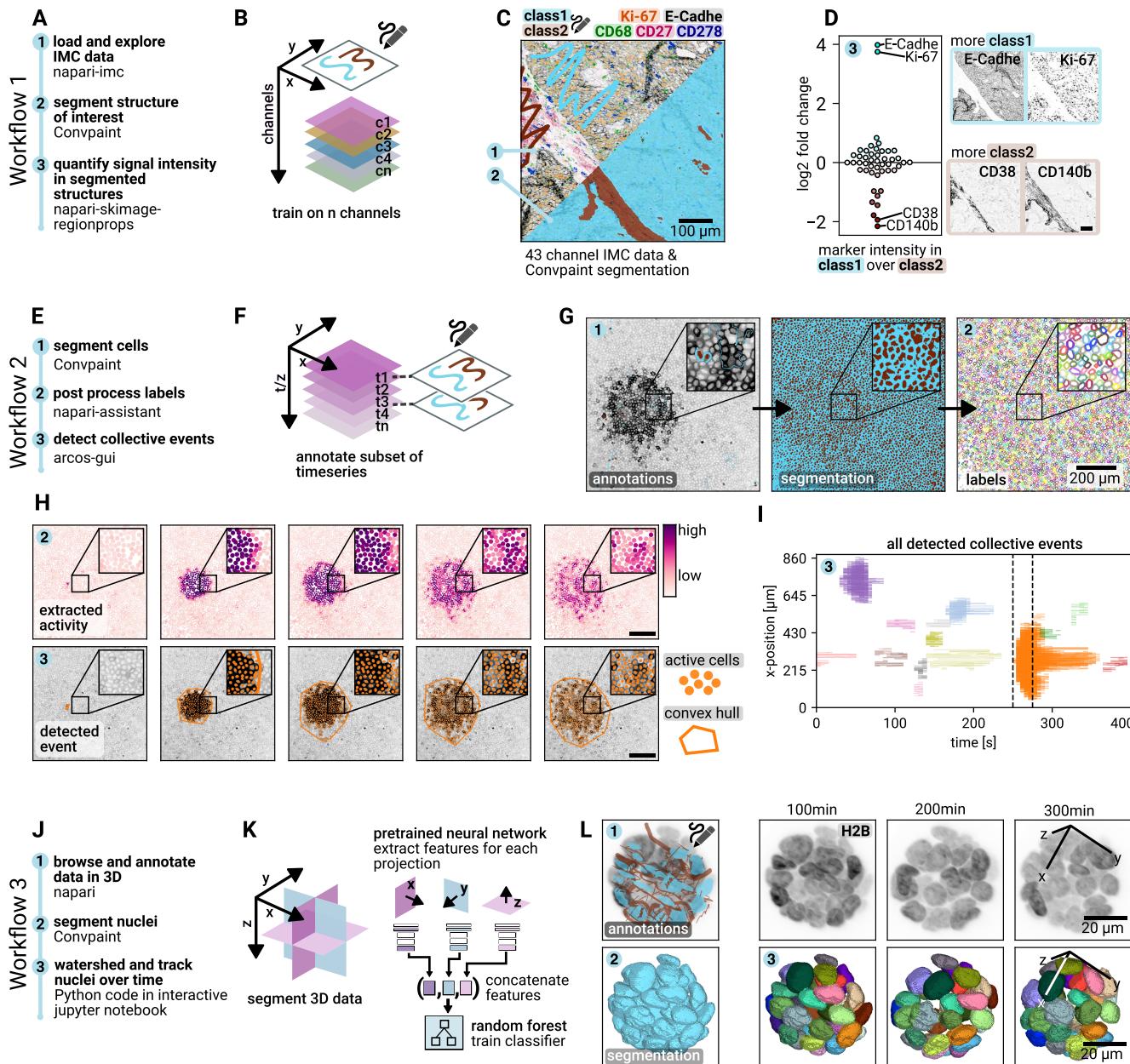
82 **Workflow 2** demonstrates Convpaint ability to analyze  
83 time lapse data with native support for interactive visualiza-  
84 tion in napari, to detect collective calcium signaling waves  
85 in an epithelial monolayer expressing a calcium biosensor  
86 (fig. 2E-I, suppl. movie M1) [13]. First, Convpaint is trained  
87 with scribbles on multiple frames to segment cells, then  
88 napari-assistant [14] is used for post-processing the labels  
89 and extracting biosensor information, and finally, ARCOS  
90 [15, 16] is utilized to detect and quantify collective signal-  
91 ing events.

92 To fully exploit information in 3D datasets using Convpaint,  
93 we implemented feature extraction from xy, xz, and yz pro-  
94 jections. The features are concatenated to form the input  
95 for the random forest classifier. In **Workflow 3**, fig. 2J-L,  
96 this approach's effectiveness is demonstrated on a light-  
97 sheet time lapse dataset of 3D mammary acini express-  
98 ing a nuclear marker and an ERK biosensor. Following  
99 Convpaint segmentation, a simple 3D watershed is used

100 for instance segmentation, and nuclei are tracked with an  
101 overlap-based algorithm (suppl. movie M2). Fig. S3 de-  
102 tails how segmentation can be used to extract single-cell  
103 ERK activity signaling trajectories. Fig. S4 shows how in-  
104 formation from multiple projections improves segmentation  
105 performance on a synthetic 3D dataset.

106 In this paper, we demonstrate that by using ViTs as feature  
107 extractors, which capture semantically richer information,  
108 we can accomplish tasks that were previously unattainable  
109 with traditional pixel classifiers. For example, leveraging  
110 DINOv2, pretrained on a large and diverse image dataset  
111 [17], we successfully segmented the head and tail of mice  
112 (fig. S5, suppl. movie M3) and detected whether their  
113 eyes were open or closed (fig. S6, M4). This approach  
114 is broadly applicable across animal species. With minimal  
115 annotations, we tracked body parts of sharks, even during  
116 complex movements such as rotations (fig. S7, M5).

117 Since there are no standard ground truth test datasets  
118 for evaluating interactive pixel classifiers, we developed  
119 a computational pipeline to generate test datasets for



**Fig. 2. Image analysis workflows using Convpaint.** **Workflow 1 (multichannel dataset):** A Example with multichannel IMC data. B Handles arbitrary input channels. C Interactive exploration with napari-imc. Labeled structures guide segmentation across all channels. Scale bar: 100 µm. D Use class labels for data exploration and statistical analysis, such as identifying differentially expressed markers. Scale bar: 100 µm. **Workflow 2 (time-series):** Supports 3D and time-series data. E Combined with arcos-gui to detect collective signaling events in MDCK cell movies. F Train classifier on multiple frames or z-slices to predict the rest. G Segmentation and post-processing with napari-assistant. Scale bar: 20 µm. H Example of collective event detection over 5 frames, with signaling activity and event overlay. Scale bar: 200 µm. I Overview of all detected events, with the period from panel H marked. **Workflow 3 (3D segmentation):** J Segmentation of MCF10A acini from lightsheet microscopy, with 3D watershed instance segmentation and tracking with trackpy. K Feature extraction from multiple projections (xy, xz, yz) combined for random forest classification. L 3D rendering of tracked nuclei with color-coded IDs. Scale bar: 20 µm.

Convpaint. This pipeline automatically creates human-like scribbles from existing segmentation dataset annotations, allowing us to quantitatively assess Convpaint's segmentation performance and compare different feature extractors. It also enables testing Convpaint across varying levels of annotation coverage. We observe significant improvements in segmentation performance on complex data, such as detecting cancerous tissue in histology slides, when using pretrained neural networks as feature extractors. The pipeline used to generate the test data, the results across different datasets, and insights

into how performance depends on the number of scribbles are discussed in detail in the supplementary information and shown in fig. S8. Randomly sampled results for all datasets are shown in figs. S9, S10, S11. The results show Convpaint's flexibility and performance, which, together with its seamless integration within the napari ecosystem, makes it an attractive segmentation tool for a wide variety of image analysis tasks and data types. We envision Convpaint empowering researchers across diverse fields, from cell biologists to ecologists. With no programming knowledge required, users can access state-

142 of-the-art, steerable segmentation tailored to their specific  
143 needs with just a few clicks. This accessibility brings ad-  
144 vanced image analysis techniques to a broader scientific  
145 community.

146 The code is open source (BSD-3), available on GitHub<sup>1</sup>,  
147 and can be installed from the napari hub<sup>2</sup> or PyPi<sup>3</sup>. It runs  
148 on all common operating systems and standard consumer  
149 hardware, with optional GPU acceleration. We provide in-  
150 stallation instructions, documentation, and video tutorials<sup>4</sup>.

## 151 ACKNOWLEDGEMENTS

152 This work has been supported by the Chan Zuckerberg Initiative (CZI) grant NP2-  
153 000000095 to LH and OP, Unisientia fellowship 187-2021 to OP, and Schweiz-  
154 erischer Nationalfonds (SNF) grant 310030\_185376 to OP. We thank the Scientific  
155 Center for Optical and Electron Microscopy (ScopeM) of ETH Zurich, Switzerland,  
156 for access to their instruments and services and Dr. Tobias Schwartz for his assis-  
157 tance in acquiring lightsheet data. Calcium imaging data was kindly provided by Ya-  
158 suto Takeuchi and Yasuyuki Fujita. Other microscopy experiments were performed  
159 on equipment supported by the Microscopy Imaging Center (MIC), University of  
160 Bern, Switzerland. The mouse icon in A by DBCLS <https://togtv.dbcls.jp/en/pics.html> is CC-BY 4.0 licensed.

## 162 AUTHOR CONTRIBUTIONS

163 LH conceptualized the work. LH, GW, RS, AS, MD, MV, and BG contributed to  
164 the development of the software and documentation. RS quantified performance.  
165 JF, LH, BG, RS, and AF acquired data. Figures were created by LH. LH and OP  
166 wrote the manuscript and acquired funding. All authors read and approved the final  
167 manuscript.

## 168 COMPETING FINANCIAL INTERESTS

169 The authors declare that they have no conflict of interest.

## 170 Bibliography

- 171 [1] Valentin Koch, Sophia J. Wagner, Salome Kazemina, Ece Sancar, Matthias Hehr, Julia  
172 Schnabel, Tingy Peng, and Carsten Marr. Dinobloom: A foundation model for generaliz-  
173 able cell embeddings in hematology, 2024.
- 174 [2] Richard J Chen, Tong Ding, Ming Y Lu, Drew FK Williamson, Guillaume Jaume, Bowen  
175 Chen, Andrew Zhang, Daniel Shao, Andrew H Song, Muhammad Shaban, et al. Towards a  
176 general-purpose foundation model for computational pathology. *Nature Medicine*, 2024.
- 177 [3] Ramon Faenster, Jacob Hanemann, Sohyon Lee, and Berend Snijder. Self-supervised  
178 vision transformers accurately decode cellular state heterogeneity. January 2023.  
179 doi:[10.1101/2023.01.16.524226](https://doi.org/10.1101/2023.01.16.524226).
- 180 [4] Stuart Berg, Dominik Kutra, Thorben Kroeger, Christoph N. Straehle, Bernhard X. Kausler,  
181 Carsten Haubold, Martin Schiegg, Janez Ales, Thorsten Beier, Markus Rudy, Kemal  
182 Eren, Jaime I. Cervantes, Buote Xu, Fynn Beuttenmueller, Adrian Wolny, Chong Zhang,  
183 Ullrich Koethe, Fred A. Hamprecht, and Anna Kreshuk. ilastik: interactive machine  
184 learning for (bio)image analysis. *Nature Methods*, September 2019. ISSN 1548-7105.  
185 doi:[10.1038/s41592-019-0582-9](https://doi.org/10.1038/s41592-019-0582-9).
- 186 [5] Ignacio Arganda-Carreras, Verena Kaynig, Curtis Rueden, Kevin W Elceiri, Johannes  
187 Schindelin, Albert Cardona, and H Sebastian Seung. Trainable Weka Segmentation: a ma-  
188 chine learning tool for microscopy pixel classification. *Bioinformatics*, 33(15):2424–2426, 03  
189 2017. ISSN 1367-4803. doi:[10.1093/bioinformatics/btx180](https://doi.org/10.1093/bioinformatics/btx180).
- 190 [6] Peter Bankhead, Maurice B Loughrey, José A Fernández, Yvonne Dombrowski, Darragh G  
191 McArt, Philip D Dunne, Stephen McQuaid, Ronan T Gray, Liam J Murray, Helen G Coleman,  
192 Jacqueline A James, Manuel Salto-Tellez, and Peter W Hamilton. QuPath: Open source  
193 software for digital pathology image analysis. *Sci. Rep.*, 7(1), December 2017.
- 194 [7] Haase Robert, Dohyeon Lee, Draga Doncila Pop, and Žigutyté Laura. haesleinhuepf/napiro-  
195 accelerated-pixel-and-object-classification: 0.14.1, 2023.
- 196 [8] Thomas Leung and Jitendra Malik. *International Journal of Computer Vision*, 43(1):29–44,  
197 2001. ISSN 0920-5691. doi:[10.1023/a:101126920638](https://doi.org/10.1023/a:101126920638).
- 198 [9] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep  
199 convolutional neural networks. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger,  
200 editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Asso-  
201 ciates, Inc., 2012.
- 202 [10] Nils Eling and Jonas Windhager. Example imaging mass cytometry raw data. February  
203 2022. doi:[10.5281/zenodo.5949116](https://doi.org/10.5281/zenodo.5949116).
- 204 [11] Jonas Windhager, Bernd Bodenmüller, and Nils Eling. An end-to-end workflow for multi-  
205 plexed image processing and analysis. *bioRxiv*, 2021. doi:[10.1101/2021.11.12.468357](https://doi.org/10.1101/2021.11.12.468357).
- 206 [12] Jonas Windhager, Vito Riccardo Tomaso Zanotelli, Daniel Schulz, Lasse Meyer, Michelle  
207 Daniel, Bernd Bodenmüller, and Nils Eling. An end-to-end workflow for multiplexed image  
208 processing and analysis. *Nature Protocols*, 18(11):3565–3613, October 2023. ISSN 1750-  
209 2799. doi:[10.1038/s41596-023-00881-0](https://doi.org/10.1038/s41596-023-00881-0).
- 210 [13] Yasuto Takeuchi, Rika Narumi, Ryutaro Akiyama, Elisa Vitiello, Takanobu Shirai, Nobuyuki  
211 Tanimura, Keisuke Kuromiya, Susumu Ishikawa, Mihoko Kajita, Masazumi Tada, Yukinari  
212 Haraoka, Yuki Akieda, Tohru Ishitani, Yoichiro Fujioka, Yusuke Ohba, Sohei Yamada, Yoichi-  
213 roh Hosokawa, Yusuke Toyama, Takaaki Matsui, and Yasuyuki Fujita. Calcium wave pro-  
214 motes cell extrusion. *Current Biology*, 30(4):670–681.e6, February 2020. ISSN 0960-9822.  
215 doi:[10.1016/j.cub.2019.11.089](https://doi.org/10.1016/j.cub.2019.11.089).
- 216 [14] Robert Haase, Ryan Savill, Peter Sobolewski, and Lee Dohyeon. haesleinhuepf/napiro-  
217 assistant: 0.4.7, 2023.
- 218 [15] Paolo Armando Gagliardi, Benjamin Grädel, Marc-Antoine Jacques, Lucien Hinderling, Pas-  
219 cal Ender, Andrew R. Cohen, Gerald Kastberger, Olivier Pertz, and Maciej Dobrzański. Auto-  
220 matic detection of spatio-temporal signaling patterns in cell collectives. *Journal of Cell  
221 Biology*, 222(10), July 2023. ISSN 1540-8140. doi:[10.1083/jcb.202207048](https://doi.org/10.1083/jcb.202207048).
- 222 [16] Maciej Dobrzański, Benjamin Grädel, Paolo Armando Gagliardi, and Olivier Pertz. Quantifi-  
223 cation of collective signalling in time-lapse microscopy images. *Methods in Microscopy*, 1  
224 (1):19–30, April 2024. ISSN 2942-3899. doi:[10.1515/mim-2024-0003](https://doi.org/10.1515/mim-2024-0003).
- 225 [17] Maxime Oquab, Timothée Darct, Théo Moutakanni, Huy Vo, Marc Szarfaniec, Vasil Khalilov,  
226 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaeldin El-Nouby, Mahmoud  
227 Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li,  
228 Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien  
229 Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust  
230 visual features without supervision. 2024. doi:[10.48550/arXiv.2304.07193](https://doi.org/10.48550/arXiv.2304.07193).
- 231 [18] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale  
232 image recognition, 2015.
- 232 [19] Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznes-  
233 sensky, Bin Bao, Peter Bell, David Berard, Evgeni Burrovski, Geeta Chauhan, Anjali Chour-  
234 dia, Will Constable, Albin Desmaison, Zachary DeVito, Elias Ellison, Will Feng, Jing Gong,  
235 Michael Gschwind, Brian Hirsh, Sherlock Huang, Kshitijee Kalambarkar, Lauren Kirsch,  
236 Michael Lazos, Mario Lezcano, Yanbo Liang, Jason Liang, Yinghai Lu, C. K. Luk, Bert Ma-  
237 her, Yunjie Pan, Christian Puhrsch, Matthias Reso, Mark Saroufim, Marcos Yukio Siraichi,  
238 Helen Suk, Shunting Zhang, Michael Suo, Phil Tillett, Xu Zhao, Eikan Wang, Keren Zhou,  
239 Richard Zou, Xiaodong Wang, Ajit Mathews, William Wen, Gregory Chanan, Peng Wu,  
240 and Soumith Chintala. Pytorch 2: Faster machine learning through dynamic python byte-  
241 code transformation and graph compilation. In *Proceedings of the 29th ACM International  
242 Conference on Architectural Support for Programming Languages and Operating Systems,  
243 Volume 2*, ASPLOS '24. ACM, April 2024. doi:[10.1145/3620665.3640366](https://doi.org/10.1145/3620665.3640366).
- 244 [20] Carsen Stringer and Marius Pachitariu. Transformers do not outperform cellpose. April  
245 2024. doi:[10.1101/2024.04.06.587952](https://doi.org/10.1101/2024.04.06.587952).
- 246 [21] Carsen Stringer, Tim Wang, Michalis Michaelos, and Marius Pachitariu. Cellpose: a gener-  
247 alist algorithm for cellular segmentation. *Nature Methods*, 18(1):100–106, December 2020.  
248 ISSN 1548-7105. doi:[10.1038/s41592-020-01018-x](https://doi.org/10.1038/s41592-020-01018-x).
- 249 [22] Xiongwei Wu, Xia Fu, Ying Liu, Ee-Peng Lim, Steven C. H. Hoi, and Qianru Sun. A large-  
250 scale benchmark for food image segmentation. *CoRR*, abs/2105.05409, 2021.
- 251 [23] Mohamed Amgad, Habiba Elfandy, Hagar Hussein, Lamees A Atteya, Mai A T Elsebaie,  
252 Lamia S Abo Elnasr, Rokia A Sakr, Hazem S E Salem, Ahmed F Ismail, Anas M Saad,  
253 Jourmana Ahmed, Maha A T Elsebaie, Mustafajur Rahman, Inas A Ruhiban, Nada M El-  
254 gazar, Yahya Alagha, Mohamed H Osman, Ahmed M Alhusseiny, Mariam M Khalaf, Abo-  
255 Alela F Younes, Ali Abdulkarim, Duaa M Younes, Ahmed M Gadallah, Ahmad M Elkashash,  
256 Salma Y Fala, Basma M Zaki, Jonathan Beezley, Deepak R Chittajallu, David Manthey,  
257 David A Gutman, and Lee A D Cooper. Structured crowdsourcing enables convolutional  
258 segmentation of histology images. *Bioinformatics*, 35(18):3461–3467, February 2019. ISSN  
259 1367-4811. doi:[10.1093/bioinformatics/btz083](https://doi.org/10.1093/bioinformatics/btz083).
- 260 [24] Lucien Hinderling, Maciej Dobrzański, Yasuto Takeuchi, and Olivier Pertz. Calcium waves  
261 in mdck epithelium. *BioStudies Database*, 2024. doi:[10.6019/s-biad1135](https://doi.org/10.6019/s-biad1135).
- 261 [25] Pascal Ender, Paolo Armando Gagliardi, Maciej Dobrzański, Agne Frismantienė, Coralie  
262 Dessauges, Thomas Höhener, Marc-Antoine Jacques, Andrew R. Cohen, and Olivier Pertz.  
263 Spatiotemporal control of erk pulse frequency coordinates fate decisions during mammary  
264 acinar morphogenesis. *Developmental Cell*, 57(18):2153–2167.e6, September 2022. ISSN  
265 1534-5807. doi:[10.1016/j.devcel.2022.08.008](https://doi.org/10.1016/j.devcel.2022.08.008).
- 266 [26] Agne Frismantienė, Lucien Hinderling, and Olivier Pertz. Light sheet 3d timelapse of a  
267 human breast cell acini. *BioStudies Database*, 2024. doi:[10.6019/s-biad1134](https://doi.org/10.6019/s-biad1134).
- 267 [27] Stéfan van der Walt, Johannes L. Schönberger, Juan Nunez-Iglesias, François Boulogne,  
268 Joshua D. Warner, Neil Yager, Emmanuelle Gouillart, and Tony Yu. scikit-image: image  
269 processing in python. *PeerJ*, 2:e453, June 2014. ISSN 2167-8359. doi:[10.7717/peerj.453](https://doi.org/10.7717/peerj.453).
- 270 [28] Dmitry V. Sorokin, Igor Peterlik, Vladimir Ulman, David Svoboda, Tereza Necasova, Katsia-  
271 rina Morgaenko, Livia Eiselleova, Lenka Tesarova, and Martin Maska. Filogen: A model-  
272 based generator of synthetic 3-d time-lapse sequences of single motile cells with growing  
273 and branching filopodia. *IEEE Transactions on Medical Imaging*, 37(12):2630–2641, De-  
274 cember 2018. ISSN 1558-254X. doi:[10.1109/tmi.2018.2845884](https://doi.org/10.1109/tmi.2018.2845884).
- 275 [29] Veblorn Ljosa, Katherine L Sokolnicki, and Anne E Carpenter. Annotated high-throughput  
276 microscopy image sets for validation. *Nature Methods*, 9(7):637–637, June 2012. ISSN  
277 1548-7105. doi:[10.1038/nmeth.2083](https://doi.org/10.1038/nmeth.2083).
- 278 [30] Romina Burla, Mattia La Torre, Giorgia Zanetti, Alex Bastianelli, Chiara Merigliano, Simona  
279 Del Giudice, Alessandro Vercelli, Ferdinando Di Cunto, Marina Boido, Fiammetta Verni, and  
280 Isabella Saggio. p53-sensitive epileptic behavior and inflammation in f11 hypomorphic mice.  
281 *Frontiers in Genetics*, 9, November 2018. ISSN 1664-8021. doi:[10.3389/fgene.2018.00581](https://doi.org/10.3389/fgene.2018.00581).

<sup>1</sup>[www.github.com/guiwitz/napari-convpaint](https://www.github.com/guiwitz/napari-convpaint)

<sup>2</sup><https://www.napari-hub.org/plugins/napari-Convpaint>

<sup>3</sup><https://pypi.org/project/napari-convpaint/>

<sup>4</sup><https://guiwitz.github.io/napari-Convpaint/book/Landing.html>

## 285 Methods

286 **Convpaint implementation details.** Convpaint features a modular architecture designed to accommodate  
287 a wide range of feature extractors, enhancing existing algorithms or pretrained models with added  
288 steerability. We compare three different types of feature extractors:

289 **CNN** We utilize the VGG16 [18] architecture implemented in pytorch [19], pretrained on the ImageNet  
290 dataset, to extract local image features such as edges, textures, and color channel correlations  
291 when working with RGB images. Downscaled versions of the input image are passed through  
292 VGG16, creating a featurized image pyramid. These features are then upscaled and concatenated  
293 with unscaled outputs and deeper CNN layer features, balancing segmentation speed and accu-  
294 racy. This method effectively generalizes to a variety of image segmentation tasks (fig. S2). We  
295 evaluated different configurations of input scalings and layers for feature extraction and provide  
296 default settings that perform well on all of the tested datasets.

297 **ViT** We incorporate two ViT models — DINOv2 [17] and UNI [2]. DINOv2 is pretrained on 142M images  
298 from ImageNet, while UNI is pretrained on a large histology dataset. These models extract patch  
299 features of 14x14 pixels (DINOv2) and 16x16 pixels (UNI), providing superior performance in certain  
300 segmentation tasks despite a loss in resolution for fine details below the patch size, like small cell  
301 protrusions. For each patch, the ViT-S/14 distilled DINOv2 model we used extracts 384 features,  
302 while UNI extracts 1024 features. For all DINOv2 experiments, we utilized the variant with registers,  
303 as this configuration produced less patch noise in the predictions (fig. S8J).

304 **Classical Filter bank** To compare the performance of Convpaint when using pretrained neural networks  
305 versus classical filter banks as feature extractors, we employed the filters implemented in napari-  
306 ilastik<sup>5</sup>. We chose the maximal combination of filters and sigma parameters suggested in the  
307 library, including Gaussian, Laplacian of Gaussian, Gaussian gradient magnitude, difference of  
308 Gaussians, structure tensor eigenvalues, and Hessian of Gaussian eigenvalues, with sigma values  
309 0.3, 0.7, 1.0, 1.6, 3.5, 5.0, 10.0. Fig. S12 shows a visual comparison of filters used in classical filter  
310 banks versus learned convolutional filters extracted from VGG16.

311 Convpaint is optimized for both training and prediction efficiency:

- 312 • Crops Around Annotations: Avoids processing entire images by extracting crops around annotated  
313 pixels.
- 314 • Tiling and Parallel Processing: Handles large images by tiling them and using parallel processing,  
315 with appropriate padding to minimize edge effects. One-click batch processing for image stacks.
- 316 • Data Management: Manages larger-than-memory files using Dask, appropriate handling of addi-  
317 tional image dimensions (channels vs. time/z-slices)
- 318 • Customizability: Users can easily integrate other feature extractors by implementing a simple func-  
319 tion that returns a feature matrix from an image. Convpaint takes care of the user interface, classifier  
320 training, data management, and parallelization.

321 For users wanting to implement a custom feature extractor, we provide a blueprint as a starting point.  
322 Convpaint makes it easy to take existing architectures and repurpose them in minutes. To give another  
323 example, we have explored using intermediate outputs of a cellpose model as a feature extractor. While  
324 the model is trained to predict cell masks, in Convpaint it can be steered to segment cell boundaries  
325 or nuclei with a couple of scribbles. Another possibility is to simply concatenate the output of multiple  
326 feature extractors, which allows combining their strengths. As an example, we successfully combined the  
327 pixel-accurate segmentation of VGG16 with the semantic understanding of DINOv2 S13: while DINOv2  
328 reliably differentiates shark body parts, it is limited in accuracy by the patch size of its features. VGG16,  
329 on the other hand, precisely masks the shark but lacks the semantic understanding to correctly label  
330 different anatomical structures. Combining both models achieves semantically correct labels with high  
331 spatial resolution.

332 The design of Convpaint streamlines experimentation with different feature extractors and ensures that  
333 new innovations in computer vision can be easily integrated to leverage its performance.

<sup>5</sup> <https://github.com/ilastik/ilastik-napari>

334 **Quantification of segmentation performance.** Assessing Convpaint's performance, especially given its  
335 interactive nature, is challenging. Even non-interactive models face problems in unbiased performance  
336 evaluation in bioimage analysis [20]. Given a lack of scribble-annotated datasets, we created an algorithm  
337 to generate human-like scribbles from existing ground truth datasets, allowing for an unbiased quantita-  
338 tive assessment of segmentation performance. For each image from the data set, we generated scribble  
339 masks with varying annotation densities (fig. S8A). We evaluated three datasets: cellpose [21], food-  
340 seg103 [22], and a subset of a breast cancer histology slide database [23]. We chose the foodseg103  
341 dataset based on the hypothesis that classical filters would struggle to assign semantic information for  
342 items containing highly variable textures. Similarly, we selected the breast cancer dataset, representing  
343 a common challenging use case in biological research. In total, we evaluated segmentation performance  
344 on 148k samples. The code to automatically generate scribbles and recreate the figures is available on  
345 GitHub<sup>6</sup>. The repository also contains the full results, including multiple performance metrics for each  
346 image and classifier at different levels of scribble annotations, as well as plots exploring the effects of  
347 feature extractor parameters on segmentation.

348 **Scribble generation.** To closely mimic human annotations, scribbles are created by combining three types  
349 of algorithmically generated lines:

- 350 1. Center ridge lines: Sampled from the primary skeleton of the ground truth mask.
- 351 2. Boundary parallel lines: Sampled from the secondary skeleton, which is derived from the ground  
352 truth mask after subtracting the primary skeleton.
- 353 3. Boundary perpendicular lines: Lines connecting the primary skeleton to the mask boundary.

354 By varying the sampling density, we can generate different levels of annotation coverage, such as 0.1%  
355 or 1% of the image pixels. The algorithm can also vary scribble type, length, and width, making it versatile  
356 for research scenarios beyond the scope of this study, e.g. how scribble types affect segmentation perfor-  
357 mance. For the cellpose dataset, which consists mostly of images with numerous small, cell-like objects,  
358 we generated a large number of short, 1-pixel-wide scribbles. The ground truth masks were converted  
359 from instance segmentation to semantic segmentation (i.e., foreground/background instead of cell IDs).  
360 For the foodseg103 dataset, which features fewer but larger regions of different food items, we generated  
361 fewer, longer scribbles with a width of 2 pixels. For the histology dataset, we created medium-length  
362 scribbles with 2 pixels width.

363 **Dataset evaluation.** As quantitative readout for segmentation performance we report the mean intersection  
364 over union (mIoU). For the cellpose dataset (540 images, fig. S8B), which involves segmenting small  
365 cell-like structures from a dark background, we observe similar performance between VGG16 filters and  
366 classical filter banks (fig. S8C). Increased performance variability in low-annotation regimes can be  
367 attributed to the random information content of the few annotated pixels. At higher annotation levels,  
368 mIoU is often limited by imprecisions in the ground truth annotations, such as missing protrusions (fig.  
369 S8D). DINOv2 performs poorly due to its patch size being too large for resolving smaller cellular details  
370 of a few pixels. Randomly selected sample images and predictions are shown in fig. S9.

371 To reduce computation time for the foodseg103 dataset (fig. S8E) with 4983 images, we excluded images  
372 larger than 640k pixels and used 520 images sampled from the dataset for evaluation. Segmentation of  
373 the images requires distinguishing food items with complex textures and colors. VGG16 filters outperform  
374 classical filter banks for all tested configurations. When only providing very sparse annotations, smaller  
375 networks with VGG16 filters show better performance, likely due to reduced overfitting. Because of the  
376 larger label regions, unlike in the cellpose dataset, DINOv2's performance is much less constrained by  
377 patch size and the ViT massively outperforms both classical and VGG16 filters, even for very low anno-  
378 tation regimes (fig. S8F). For all models, we observe diminishing returns in segmentation performance  
379 as annotation levels increase. For Convpaint users, this suggests that adding a few more annotations  
380 is beneficial when only a minimal number of scribbles are present, but if performance gains plateau, it's

<sup>6</sup> [https://github.com/quasar1357/scribbles\\_creator](https://github.com/quasar1357/scribbles_creator)

381 a good time to stop (shown for a VGG16 model in fig. S8G). Randomly selected sample images and  
382 predictions are presented in fig. S10.

383 Motivated by DINOv2's performance on foodseg103, we tested it on a subset of a breast cancer histology  
384 slide dataset (fig. S8H). DINOv2 again clearly outperforms both classical filter banks and VGG16 fil-  
385 ters. We also evaluated a histology-specific model built with the same architecture as DINOv2 [2], which  
386 surprisingly shows no advantage over DINOv2 (fig. S8I). Predictions from different models (8 out of 10  
387 images tested) are shown in fig. S11. DINOv2 achieves the highest mIoU results and produces visually  
388 less noisy predictions compared to other models.

389 In summary, DINOv2 demonstrates superior performance in segmentation tasks that require broad con-  
390 textual information when patch-sized resolution is sufficient, generalizing well across domains. In con-  
391 trast, VGG16 and classical filters provide precise segmentation when local information is adequate for  
392 the given task. Here, VGG16 generally matches or exceeds the performance of classical filters, making  
393 it a suitable alternative.

394 **Methods and data availability.**

395 **Workflow 1: IMC multichannel data.** IMC data from [10], available on Zenodo<sup>7</sup> were loaded using the  
396 napari-imc plugin [11]. Convpaint was trained on one FOV (fig. 2 shows Patient 01, Panorama 02,  
397 Position 1-1). Skimage was used to extract the per-channel statistics for the segmented regions. The  
398 log2-fold change in signal intensity was calculated using NumPy and plotted with Matplotlib. Step-by-step  
399 instructions for recreating the workflow are available in the documentation.

400 **Workflow 2: MDCK calcium waves.** Time lapse data sets of calcium signaling waves were obtained from  
401 MDCK epithelial cells that stably express GCaMP6S - a genetically encoded intracellular calcium sensor  
402 (imaging data courtesy of Yasuyuki Fujita). The movies were loaded into napari and segmented using  
403 Convpaint. The resulting binary masks were processed with ARCOS [15] to detect and quantify collective  
404 signaling events. The code to recreate the figures is available on GitHub TODO. We have made the raw  
405 imaging data available on the BioImageArchive [24] under the accession number S-BIAD1135. Step-by-  
406 step instructions for recreating the workflow are available in the documentation.

407 **Workflow 3: 3D segmentation and nuclei tracking.** We demonstrate that Convpaint can effectively segment  
408 and track single cells within a dense 3D spheroid. In fig. S3, we illustrate how this data can be further pro-  
409 cessed to extract single-cell ERK signaling activity dynamics. The cells used are MCF10A, expressing a  
410 histone H2B nuclear marker and ERK-KTR, which reports ERK activity through reversible nucleus/cytosol  
411 translocation following phosphorylation by active ERK [25] (scheme in fig. S3A). Data were acquired us-  
412 ing a lightsheet microscope with a 5-minute resolution and an isotropic voxel size of 0.145 μm. Raw  
413 imaging data and protocols are available on BioImageArchive [26] under accession number S-BIAD1134.

414 **Other figures.** The 3D nuclear data in fig. 1F is part of the scikit-image [27] data module, called *cells3d*,  
415 originally provided by the Allen Institute for Cell Science. The synthetic data in fig. S4 was generated  
416 by another group using FiloGen [28] and is available from the Broad Bioimage Benchmark Collection  
417 (BBBC046<sup>8</sup>) [29], showing cell PD-ID451/AR1/T024.

418 The datasets used for performance quantification have all been previously published. Fig. S8B shows the  
419 cellpose dataset [21]; fig. S8C shows the foodseg103 dataset [22]; and fig. S11 uses data from the Breast  
420 Cancer Semantic Segmentation (BCSS) dataset<sup>9</sup> [23]. The movie shown in fig. S5 is a supplement<sup>10</sup> to  
421 a study on epileptic behavior in mice [30]. Movies in figs. S6A, B and S7 were acquired by the authors  
422 and are available upon request. The movie in fig. S6D, E is available online<sup>11</sup> for educational purposes  
423 under the Mixkit Restricted License. Images in fig. S2 were acquired by the authors and are available

<sup>7</sup><https://doi.org/10.5281/zenodo.5555575>

<sup>8</sup><https://bbbc.broadinstitute.org/BBBC046>

<sup>9</sup><https://bcsegmentation.grand-challenge.org/BCSS/>

<sup>10</sup><https://doi.org/10.3389/fgene.2018.00581.s007>

<sup>11</sup><https://mixkit.co/free-stock-video/gray-and-white-rat-32060/>

<sup>424</sup> upon request, except for the histology slide images, which are from Wikimedia Commons<sup>12</sup>, or provided  
<sup>425</sup> by the scikit-image library, acquired at the Center for Microscopy And Molecular Imaging (CMMI).

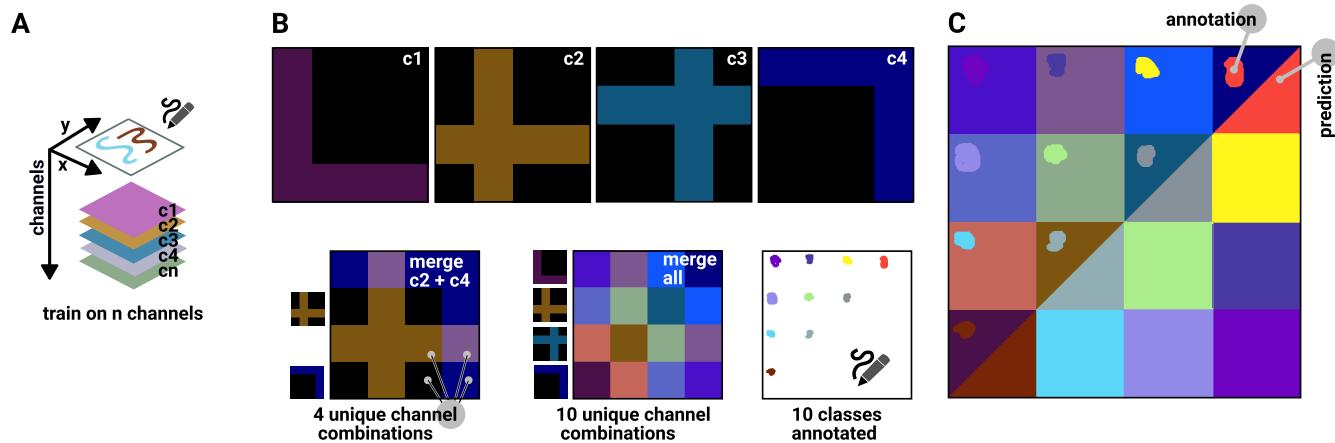
<sup>426</sup> **Supplementary figures.**

- <sup>427</sup> • [S1](#): Correlation extraction across multiple color channels.
- <sup>428</sup> • [S2](#): Image segmentation across diverse domains.
- <sup>429</sup> • [S3](#): Measuring ERK signaling dynamics at the single-cell level in MCF10A acini.
- <sup>430</sup> • [S4](#): Improved segmentation performance using multiple projections.
- <sup>431</sup> • [S5](#): Detecting mouse body parts in video.
- <sup>432</sup> • [S6](#): Eye state detection in humans and mice.
- <sup>433</sup> • [S7](#): Detecting shark body parts in video.
- <sup>434</sup> • [S8](#): Quantification of segmentation performance and model comparison.
- <sup>435</sup> • [S9](#): Feature extractor performance on the cellpose dataset.
- <sup>436</sup> • [S10](#): Feature extractor performance on the foodseg103 dataset.
- <sup>437</sup> • [S11](#): Feature extractor performance on histology slides.
- <sup>438</sup> • [S12](#): Visual comparison of handcrafted vs. learned filters.
- <sup>439</sup> • [S13](#): Combining VGG16 and DINOV2 features for enhanced spatial precision at mask boundaries while maintaining semantic information.
- <sup>440</sup>

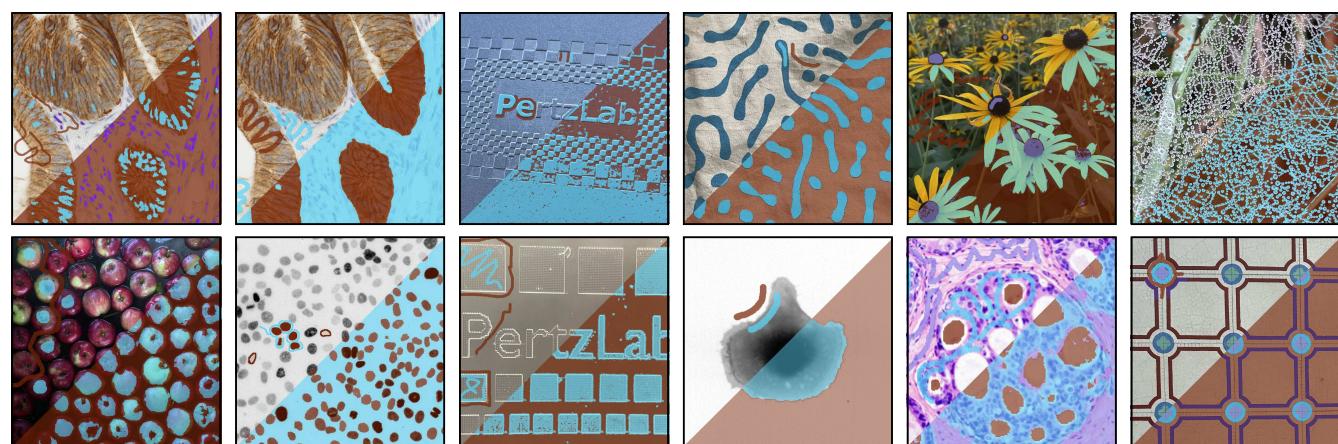
<sup>441</sup> **Supplementary movies.**

- <sup>442</sup> • M1: Detection of calcium waves in MCDK cells
- <sup>443</sup> • M2: Segmentation of MCF10A acini
- <sup>444</sup> • M3: Segmentation of mouse body parts
- <sup>445</sup> • M4: Detection of open/closed eyes in humans and mice
- <sup>446</sup> • M5: Segmentation of shark body parts

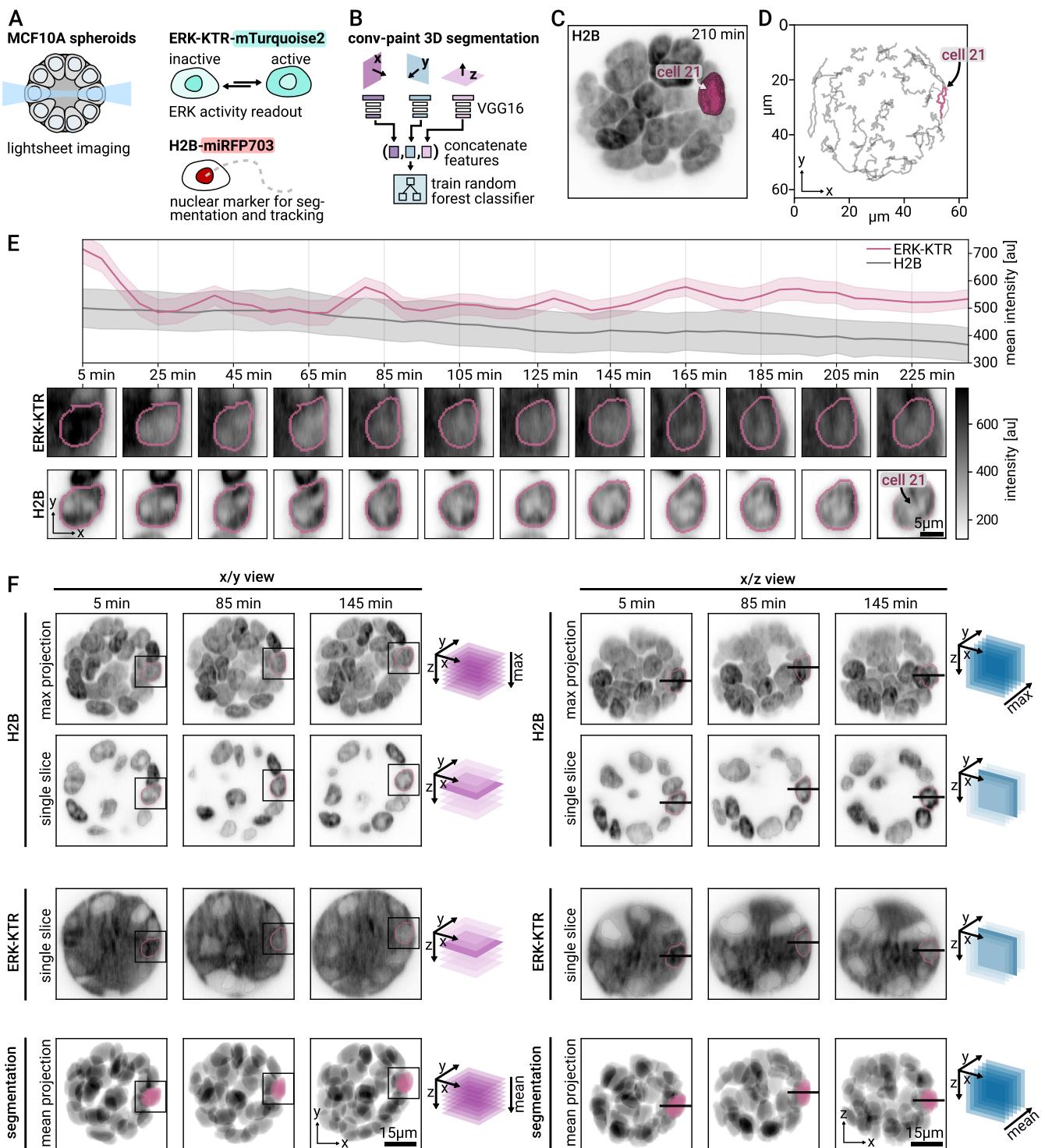
<sup>12</sup>[https://commons.wikimedia.org/wiki/File:Breast\\_DCIS\\_histopathology\\_\(1\).jpg](https://commons.wikimedia.org/wiki/File:Breast_DCIS_histopathology_(1).jpg)



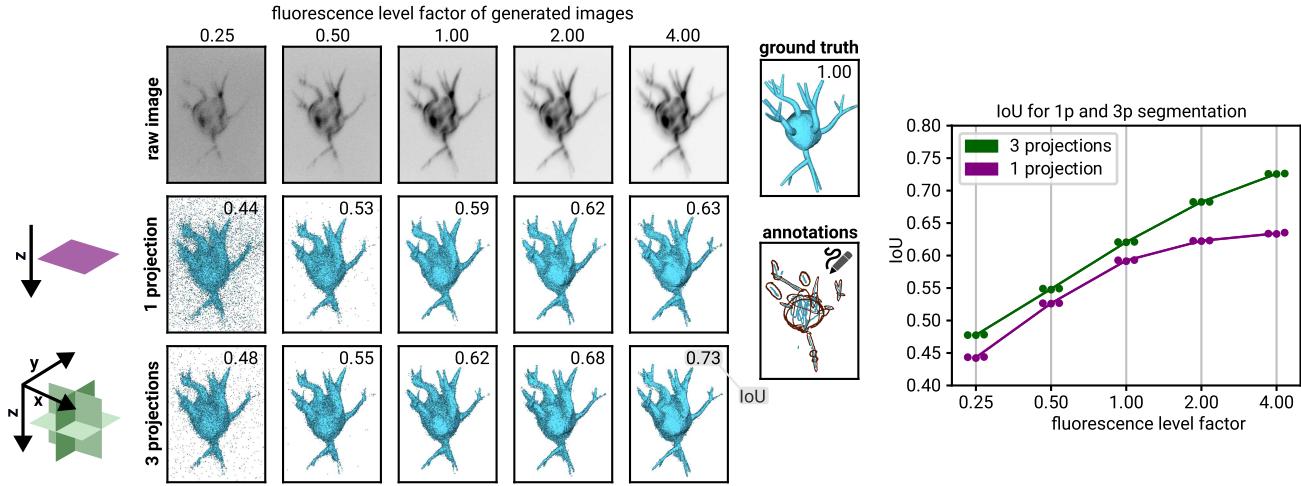
**Fig. S1. Correlation extraction across multiple color channels.** **A** The Convpaint classification algorithm can process an arbitrary number of input color channels. **B** This is demonstrated on an artificial image with four channels, which when merged lead to 10 unique color combinations. These 10 combinations are labeled with 10 class labels that can only be reconstructed if the algorithm takes into account the interplay of the different channels. **C** Convpaint correctly predicts the correct class label for pixels that were not labeled, with minor artefacts on boundaries between squares.



**Fig. S2. Image segmentation across diverse domains** All images use VGG16 with the default configuration as feature extractor.



**Fig. S3. Measuring ERK signaling dynamics at the single-cell level in MCF10A acini.** **A** Spheroids are imaged with a lightsheet microscope. The cells express an ERK activity sensor and nuclear marker for segmentation and tracking. **B** Convpaint is used to segment the nuclei in 3D. **C** Panels C-F track a single cell in the spheroid over time, here its mask is shown overlaid on a 3D max projection. **D** Tracks of all cells from 0 to 250 minutes, selected cell highlighted in color. **E** Mean nuclear ERK-KTR intensity over time as a proxy for ERK activity. We see ERK trajectories as previously described by Ender et al. [25]. In comparison, the mean intensity of the nuclear marker shows some bleaching but no fluctuations otherwise. The images show crops around the selected cell (mean of 3 z-slices,  $[+1, 0, -1]$  around the z position of the cell centroid). Scale bar is 5  $\mu$ m. **F** Highlighting the tracked position of the cell within the spheroid for different time points, projections, and channels. Box shows insets in panel E. Scale bar is 15  $\mu$ m.

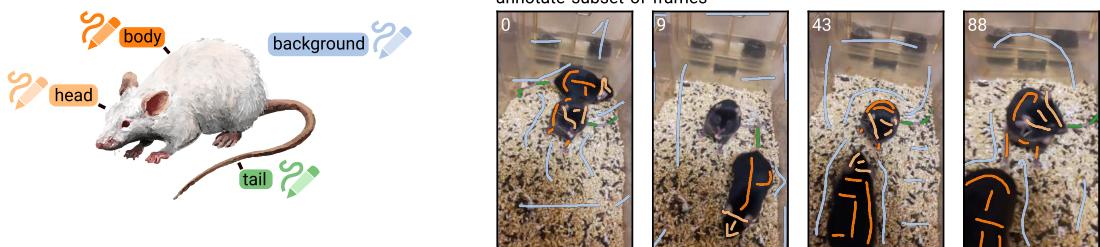


**Fig. S4. Improved segmentation performance using multiple projections.** Convpaint segmentation performance compared on an artificial cell when extracting features from 1 projection (purple) versus concatenating 3 projections (green), using VGG16 with default configuration as feature extractor. Different signal-to-noise regimes are tested, which are configured by the fluorescence level factor (0.25-4) in the FiloGen software. Performance is measured as intersection over union (IoU). Using 3 projections leads to better segmentation results for all fluorescence level factors. A larger increase in performance is observed for images with a better signal-to-noise ratio.

**A** load and browse movie



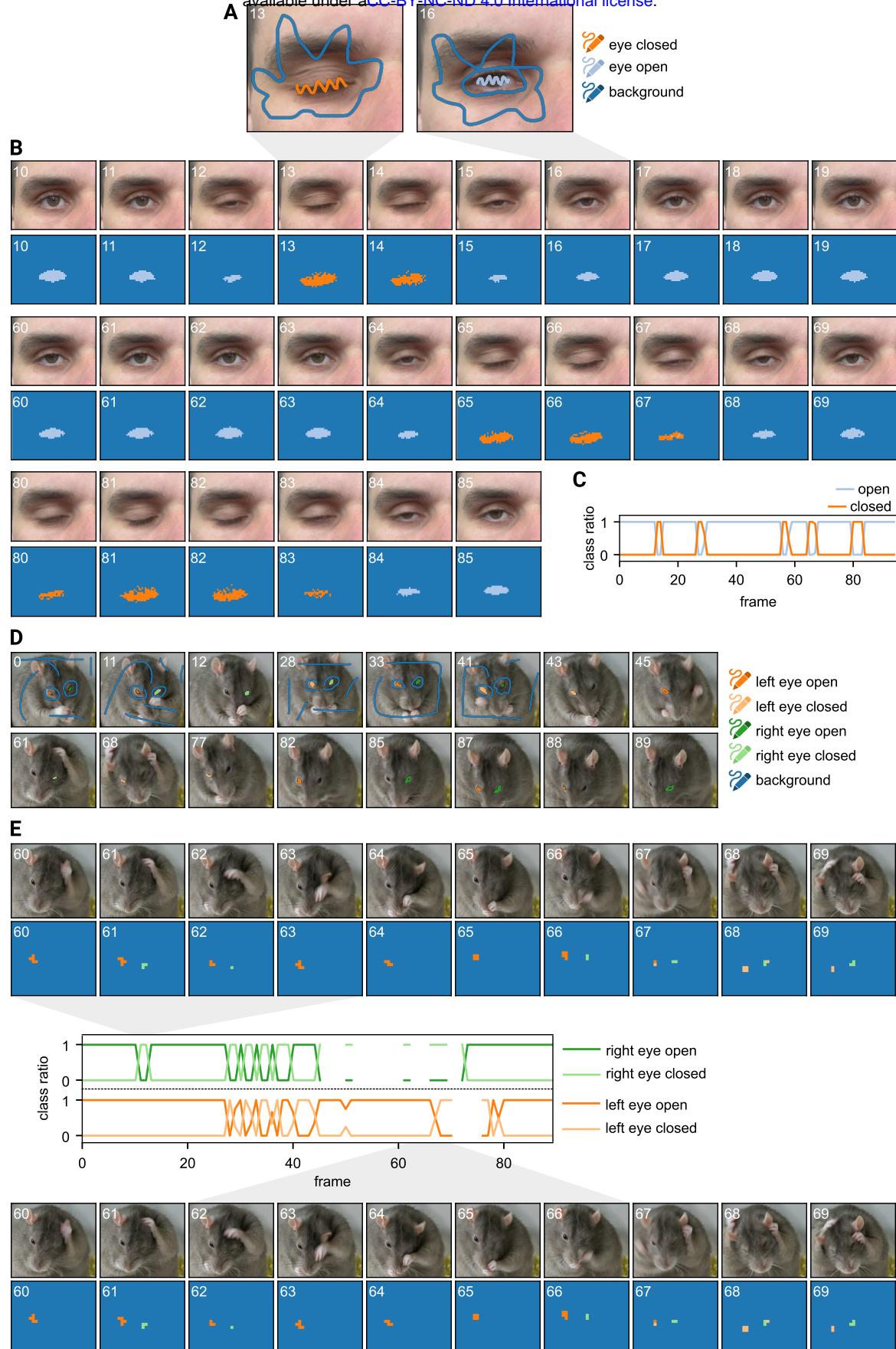
**B**



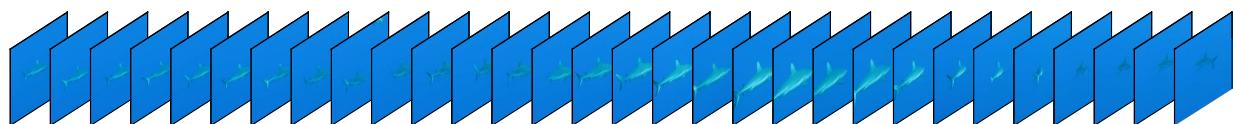
**C** segment whole movie (sample frames shown)



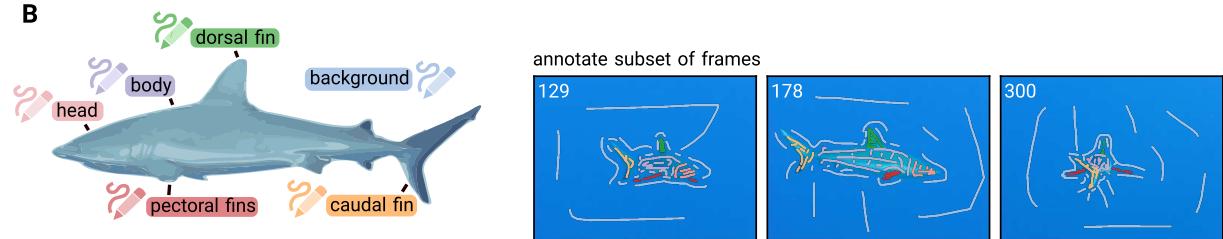
**Fig. S5. Detecting mouse body parts in video.** Corresponding movie shown in suppl. M3. **A** The movie shows mice moving in a cage. The camera is handheld and the mice move in and out of frame. Some frames contain motion blur. **B** Convpaint with DINOv2 as feature extractor is trained on scribbles of four different frames. Head, body, tail, and background are annotated. **C** The trained model is used to predict the rest of the movie.



**A** load and browse movie



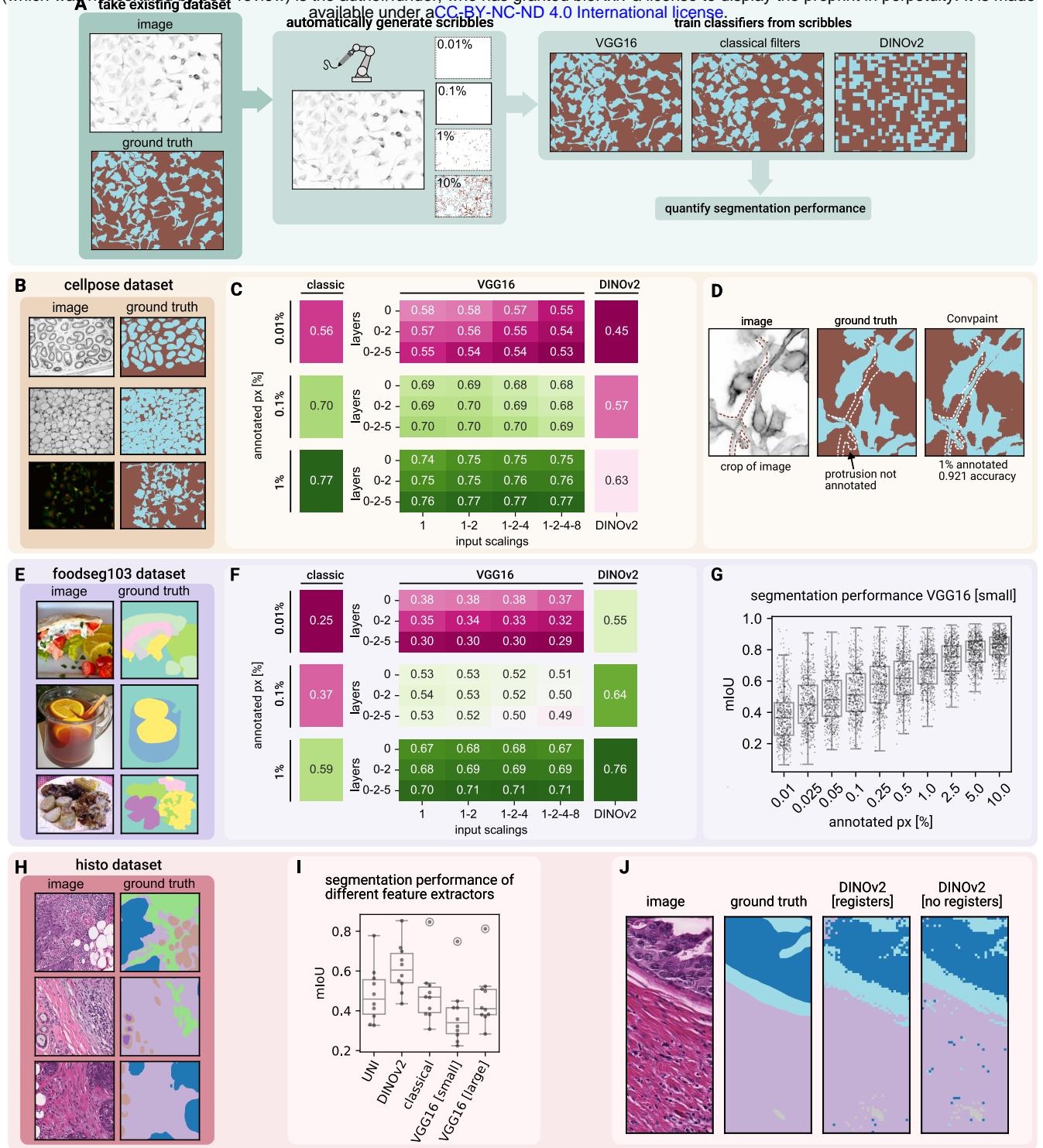
**B**



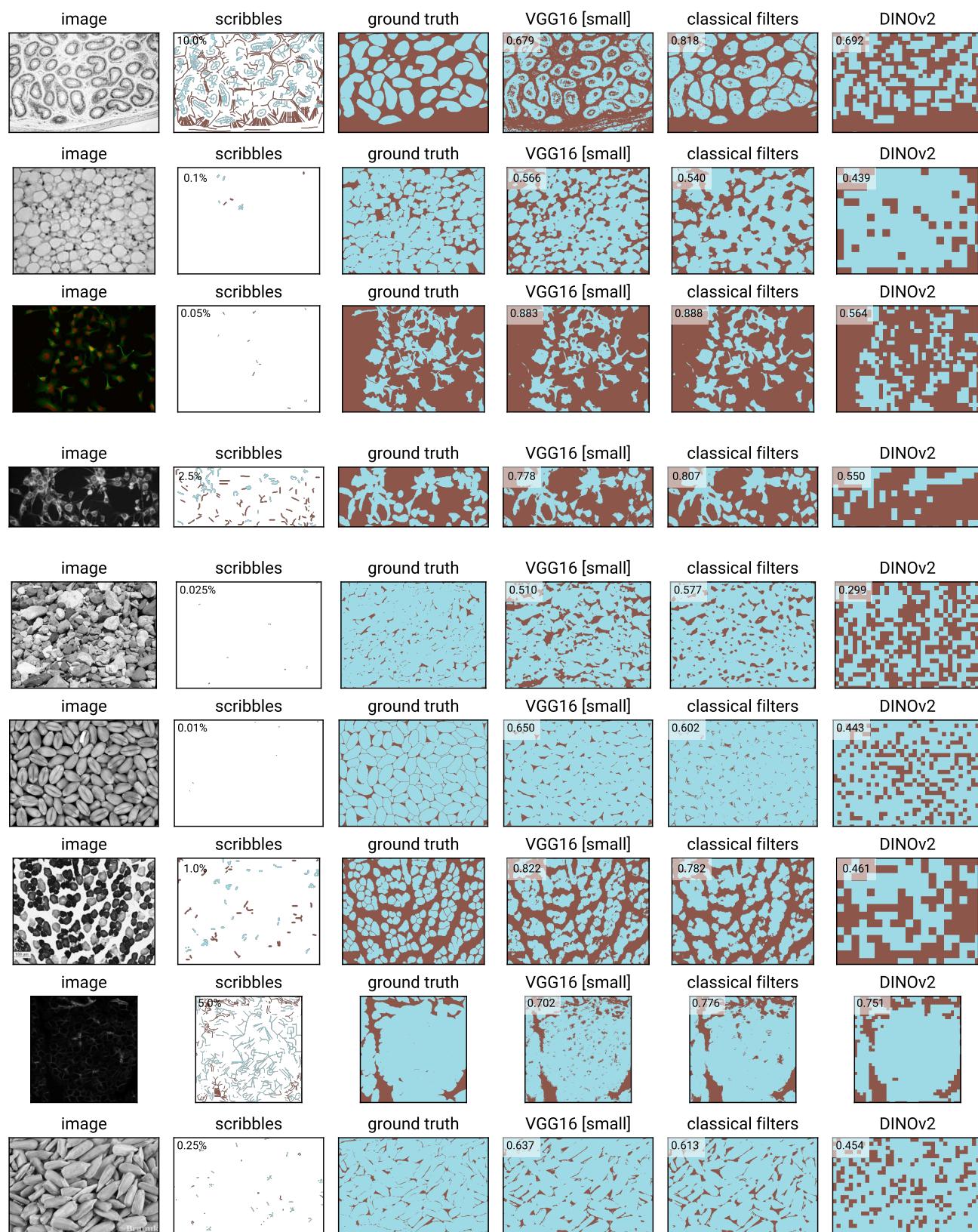
**C** segment whole movie (sample frames shown)



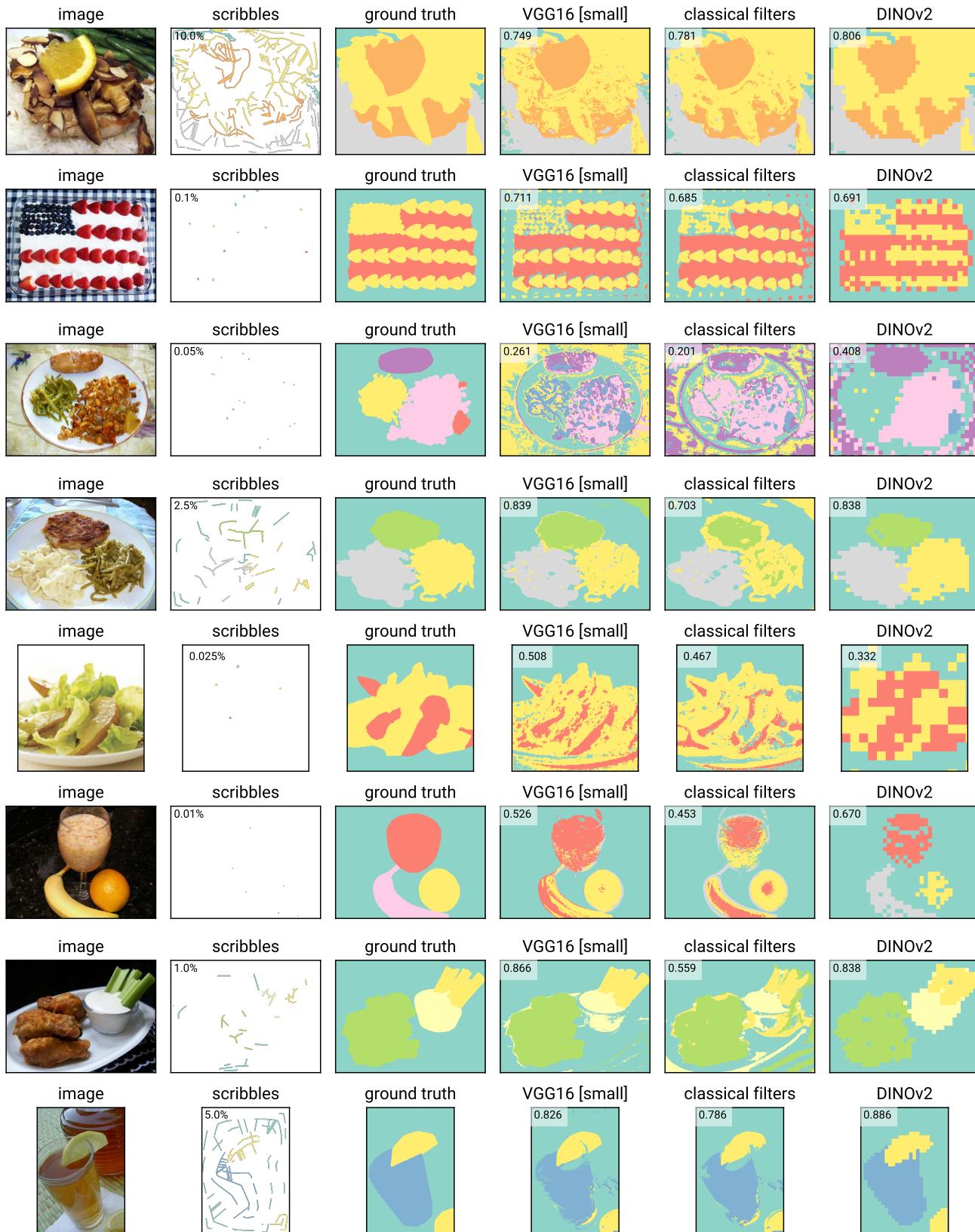
**Fig. S7. Detecting shark body parts in video.** Corresponding movie shown in suppl. M4. **A** The movie shows a shark swimming from multiple angles. Camera is hand-held. **B** Convpaint with DINOv2 as feature extractor is trained on scribbles on three different frames. Head, body, dorsal, caudal, pectoral fins, and background are annotated. **C** The trained model is used to predict the rest of the movie.



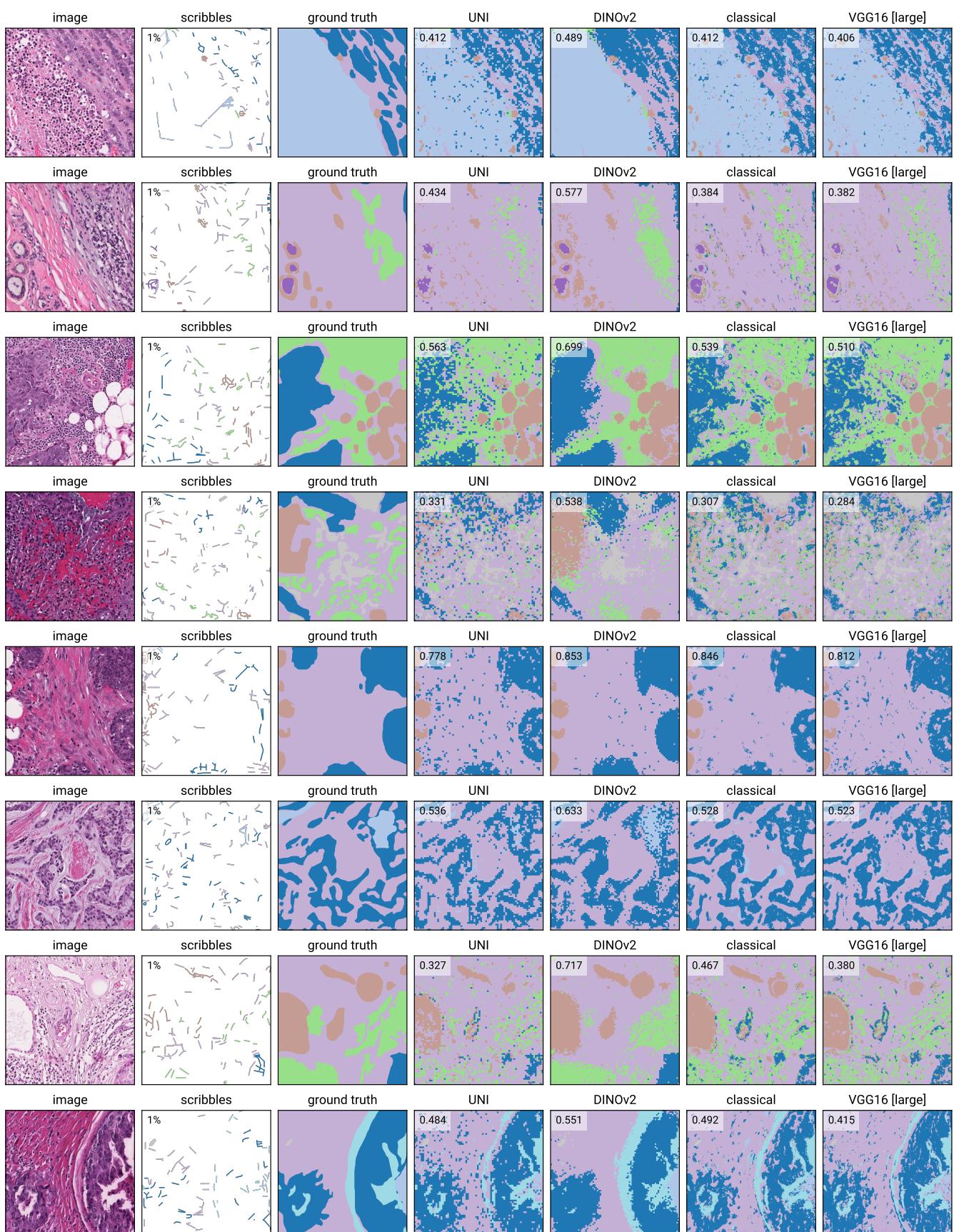
**Fig. S8. Quantification of segmentation performance and model comparison.** **A** Ground truth datasets are used to automatically generate scribbles to train Convpaint. The results are then evaluated against the ground truth to quantify the performance of different feature extractors. Various feature extractors are tested, including VGG16 layers with different input scalings, classical filter banks with all napari-ilastik filter/sigma combinations, and DINOv2 ViT-S/14. **B** Testing on the cellpose dataset (three sample images and masks shown). See fig. S9 for segmentation results. **C** Mean mIoU scores for different annotation levels in the cellpose dataset. Similar performance for classical filters and VGG16. DINOv2 underperforms due to large patch size limitations in capturing small cellular details. **D** Model performance quantification scores can be limited by the quality of ground-truth annotations. A cell protrusion, missing in the ground truth, is correctly segmented by the model. **E** Testing on the foodseg103 dataset (three samples shown). See fig. S10 for segmentation results. **F** Mean mIoU scores for different annotation levels on the foodseg103 dataset. Both pretrained architectures (CNN, ViT) outperform classical filter banks as feature extractors across all annotation regimes. DINOv2 has highest mIoU scores. In this dataset, the model's performance is less affected by patch size because of larger regions of interest. **G** For all models, diminishing returns on segmentation performance can be observed with increasing annotation levels, here shown for VGG16 (layer 0 and input scalings 1,2) on the foodseg103 dataset. **H** Testing on a histology slide dataset. See fig. S11 for segmentation results. **I** Mean mIoU scores for different annotation levels in the histology slide dataset. DINOv2 outperforms all other models, including the histology-specific model UNI. **J** DINOv2 with registers produces predictions with less patch noise (not quantified).



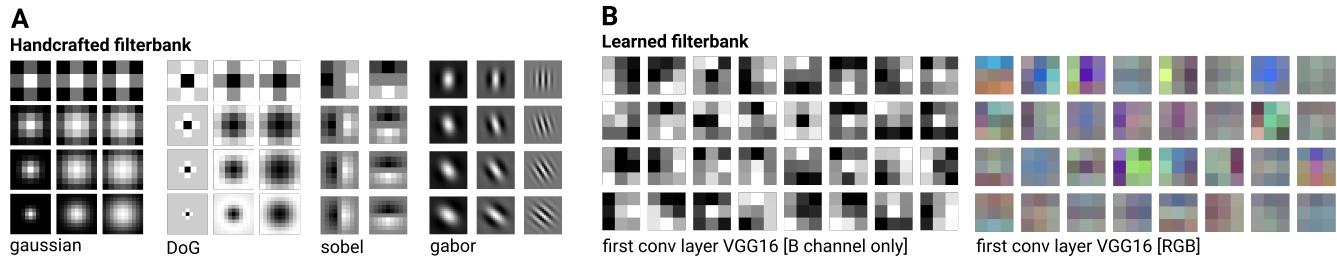
**Fig. S9. Feature extractor performance on the cellpose dataset.** Randomly selected images. For plotting, the scribbles were dilated for better visibility. The number in the upper left corner of the prediction images shows the mIoU score.



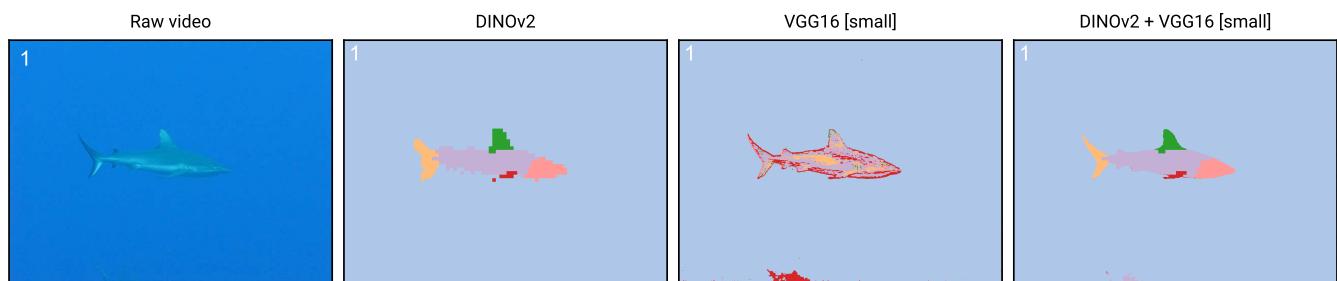
**Fig. S10. Feature extractor performance on the foodseg103 dataset.** Randomly selected images. For plotting, the scribbles were dilated for better visibility. The number in the upper left corner of the prediction images shows the mIoU score.



**Fig. S11. Feature extractor performance on histology slides.** Scribbles were automatically generated from expert annotated ground truth. For plotting, the scribbles were dilated for better visibility. Eight out of ten images used for evaluation in fig. S8I are shown. The number in the upper left corner of the prediction images shows the mIoU score. DINOv2 outperforms all other models in this dataset.



**Fig. S12. Visual comparison of handcrafted vs. learned filters** **A** Filters used classically in handcrafted filter banks. Here we show examples of filter kernels with different parameters for Gaussians, Difference of Gaussians (DoG), Sobel, and Gabor. While these handcrafted filter banks are more interpretable, the patterns they extract often overlap, leading to redundancy among the filters. **B** Filters extracted from the first convolution layer of a CNN (VGG16) network. The filters have a 3x3x3 shape, which makes them intrinsically capable of extracting correlations between color channels in RGB images. Although VGG16 filters are less interpretable, they are computationally optimized to extract orthogonal image features that are useful for image classification.



**Fig. S13. Combining VGG16 and DINov2 features for enhanced spatial precision at mask boundaries while maintaining semantic information.** While DINov2 features excel at capturing abstract semantic information at the patch level, VGG16 features are good at capturing local spatial information at the pixel level. By concatenating the features of both models, we can leverage the strengths of both models.