Visual Interpretability of Bioimaging Deep Learning Models

Anonymous Author(s) Affiliation Address email

Abstract

1	The success of deep learning in analyzing bioimages comes at the expense of
2	biologically meaningful interpretations. We review the state of the art of explainable
3	artificial intelligence (XAI) in bioimaging and discuss its potential in hypothesis

- artificial intelligence (XAI) in bioimaging and discuss its potential in hypothesis
- generation and data-driven discovery. 4

What is interpretable machine learning? 5

Deep learning (DL) has become the workhorse underlying the vast majority of bioimage analysis 6 applications. Open source code and new platforms that support usability make DL a practical tool 7 that is accessible to the broad cell biology community. Whereas traditional machine learning depends 8 on experts defining predetermined features to optimize a specific task, the DL data-driven approach 9 benefits from simultaneously optimizing feature extraction and the model to solve the downstream 10 task. Given sufficient data, DL excels at nonlinear integration of image regions, enabling it to 11 identify complex spatial patterns in the raw bioimage data. As a result, DL models surpass traditional 12 bioimage analysis methods and in some cases even exceed the human gold standard. However, the 13 success of DL comes with the discomfort of relying on 'black box' models with huge parametric 14 spaces comprising millions of parameters that introduce major challenges in human understanding 15 of the models' inner workings and in following the models' decision process. In the broader field 16 of machine learning, these challenges triggered the initiation of the vibrant nascent community of 17 18 explainable artificial intelligence (abbreviated XAI).

The terms 'interpretability' and 'explainability' are used interchangeably in most of the XAI literature. 19 We would like to highlight our perception of the nuances of difference between these terms. Both 20 terms involve the active participation of a human expert in resolving the model's decision process by 21 examining the XAI outputs. Explainability is the explicit description of the cause for a decision in a 22 manner that a human can understand. Interpretability is turning an explanation to domain- and context-23 specific insight. In the context of bioimaging, this means 'translating' an image-pattern explanation 24 to a biologically meaningful interpretation. For example, realizing that a model pays attention 25 to a specific image region in a cell (that is, an explanation) is not necessarily sufficient to derive 26 specific hypotheses regarding which organelles play a role in this decision (that is, an interpretation). 27 As another example, one explanation of a model that was trained to predict melanoma metastatic 28 efficiency from label-free cell imaging was increased light scattering within single melanoma cells 29 (Zaritsky et al., 2021). Although we do not know what causes this change in light scattering, we 30 can turn this explanation into a specific hypothesis via interpretation as a potential change in altered 31 intracellular organelle organization, such as phase-separated droplets or lysosomes, thus setting the 32 33 stage for follow-up studies to test these possibilities.

Under these definitions, interpretation is always explainable, but the converse is not necessarily true. 34 Interpretability is especially important in sciences and medicine, where 'black box' predictions are 35 not sufficient to decipher the fundamental 'mechanisms' underlying them. In this Comment, we 36

Submitted to Interpretable AI: Past, Present and Future Workshop @ NeurIPS 2024

- 37 discuss the motivation, current state and future potential of XAI and especially visual interpretability
- in the realm of bioimaging, and reflect on our experiences in this domain (Fig. 1).



Figure 1: The visual interpretability process. New biological insights and/or specific hypotheses are derived from interpreting DL models explanations.

³⁹ Why do we need visual interpretations of deep learning models?

Not every application of DL requires visual interpretability. For example, it is perhaps acceptable
to trust the 'black box' for image analysis tasks such as nucleus segmentation. In many other cases,
however, there are convincing motivations for aiming to interpret the machine's predictions. The
first motivation is trust. Lack of trust is one of the major obstacles to deploying (high-performing)
DL models in biomedical applications. Understanding how a model reaches its decisions provides
transparency that can help avoid spurious biases in the image data and in the model's decision process,
and can provide clues in cases of erroneous predictions.

Example of such erroneous interpretations include classifying an image as 'wolf' rather than 'husky' 47 due to snow in the image, biased classification of dermatology images due to background skin texture 48 and color balance (DeGrave et al., 2023) and technical staining artifacts in parasite-infected red blood 49 cells guiding the model's decision Lamiable et al. (2023). In these instances, the model is biased 50 toward 'easy' explanations that correlated with the true label, but this 'laziness' (or overfitting), 51 termed 'shortcut learning', causes the model to miss more complex, biologically relevant discriminate 52 content encapsulated in the images. Such understanding can be used to improve the trustworthiness 53 and usability of a model by understanding when models make mistakes and developing methods to 54 mitigate these situations. 55

The second motivation for applying XAI in bioimaging is the possibility of using the XAI output to 56 measure phenotypes that have been reported based on subjective observations, but that are difficult to 57 measure without the explanation. Examples include using the XAI outcome to define a quantitative 58 measure for a morphological component of in vitro fertilization (IVF) embryos called the blastocoel 59 Rotem et al. (2023) and associating protein localization patterns Kobayashi et al. (2022) or their 60 perturbations Razdaibiedina et al. (2024) with specific organelles. The third motivation, and in our 61 opinion the most exciting, is the potential to extract subtle phenotypes that are invisible to the human 62 eye and that cannot be reliably measured with standard image analysis. For example, interpreting 63 live label-free imaged melanoma cells' enhanced protrusive activity as metastasis-driving features 64 Zaritsky et al. (2021) or identifying loss of hemoglobin and crenated shapes of blood cells under 65 parasite infection Lamiable et al. (2023). 66

This discovery-driven interpretability is especially exciting in the biomedical domains, where a new explanation for a machine prediction can lead to a specific hypothesis that can be tested experimentally, closing the loop to establish a causal link. Is it realistic to derive new mechanistic hypotheses and draw biological conclusions from visual interpretation of DL models? These are the early days of XAI in bioimaging, and only time will tell. Developing methods for robust and user-friendly visual interpretability will enable effective exploration of potential biologically meaningful explanations and is a necessary step toward answering this fundamental question.

74 How do we visually interpret deep learning models?

Interpreting image-based DL models is much more difficult than interpreting tabular-based models because of the challenge of moving from pixels to semantic entities. In the field of computer vision, most XAI papers stop at the point of identifying trivial explanations such as open mouths or pointed ears for discriminating between dogs versus cats. In bioimaging, it is harder to interpret the semantic image properties that drive models' decisions because of difficulties in deciphering the less intuitive biological meanings. The latter requires extensive and systematic confirmation of these non-trivial interpretations. We believe that computational biologists developing XAI methods have competitive advantage because of the inherent expectation in the field of uncovering the 'black box' and provide a plausible biological interpretation, which is a critical step toward generating new, specific hypotheses

and designing experiments to test them.

Model interpretation can be in the context of a (local) instance-and/or of a (global) model. Local 85 interpretability characterizes the key decisions for individual predictions, and global interpretability 86 provides a holistic understanding of a model's reasoning processes. Global interpretability is useful 87 toward understanding general rules at the cost of averaging out the full spectrum of heterogeneity of 88 the local instance interpretations. For example, when analyzing a model for IVF embryo morphology 89 quality prediction, global interpretation identified the embryo size, the trophectoderm (a ring of 90 cells surrounding the embryo) and the blastocoel (a fluid-filled cavity inside the embryo) as the top 91 classification-driving morphological features Rotem et al. (2023). Local interpretability categorized 92 each embryo according to the different morphological features that dominated the specific classifier's 93 decision. 94

Some interpretability methods can provide both local and global interpretability. For example, the 95 popular tabular-based SHapley Additive exPlanations (SHAP) Lundberg and Lee (2017) calculates 96 the contribution of each (interpretable) feature to an individual prediction. These local interpretations 97 can be aggregated across all individual predictions according to specific criteria (for example, a 98 classification label) to provide a global interpretation of the model. A straightforward example of 99 global interpretability is the exclusion ('ablation') of a feature or a group of features and the ranking 100 of these features according to the corresponding degradation in the model's performance. Such 101 ablation has been applied to bioimaging data to quantify the influence of the local cell density on the 102 prediction of cell fate (apoptosis or mitosis) Soelistyo et al. (2022) and the influence of neighboring 103 cell fates (delamination, division or no behavior) on the prediction of cell behavior Yamamoto et al. 104 (2022). A related approach involves systematic assessment of model performance for varying sizes of 105 the classified image crops to probe the biologically relevant spatial scales necessary for a prediction 106 Schmitt et al. (2024). 107

The most common methods for visual interpretability can be broadly classified as saliency map 108 based or counterfactual based (Fig. 2). Saliency-based (also called 'feature attribution') explanations 109 generate 'attention maps' of the image regions that contribute the most to an individual prediction 110 by assigning to each pixel the aggregated network's inner layers' activations or gradients. Saliency 111 methods are suited for local (instance) interpretation when the visual explanations are localized 112 in the image and do not require further training: for example, in correlating attention maps with 113 their corresponding subcellular structures Doron et al. (2023) or using attention maps derived from 114 cell trajectory data to interpret the relative influence of neighboring cells on the motion of a given 115 cell LaChance et al. (2022). Counterfactual-based explanations use generative models to artificially 116 change the image in ways that maintain a realistic image and alter the model's prediction, and they 117 excel at understanding subtle image differences in domains that are less intuitive to humans, such as 118 bioimaging; these will therefore will be our focus for the rest of this Comment. 119

Counterfactual explanations require training of generative model to optimize a lower-dimensional 120 latent representation space that can be used to generate images of the entity of interest. These 121 latent representations can be traversed to generate visual counterfactual explanations along specific 122 directions in the latent space that translate to phenotypic alterations in the image space. Such 123 latent space traversals can exaggerate classification-driving phenotypes beyond the observed data 124 distribution while keeping the rest of the image fixed, together enabling the interpretation of subtle 125 phenotypes. Counterfactual-based explanations are also more suited than saliency-based methods to 126 identifying non-localized visual explanations such as shape or color. In bioimaging, this counterfactual 127 approach has shown promise for interpreting subtle cellular phenotypes by synthetically generating 128 image sequences that follow cell-state transitions Yang et al. (2020), disease progression (Zaritsky 129 130 et al., 2021), cell fate decisions Soelistyo et al. (2022) and perturbations Lamiable et al. (2023). The approach was also able to provide distinct interpretations for different models that were trained for 131 the same medical task of classifying melanoma from dermoscopic images of the skin DeGrave et al. 132 (2023).133

One approach to interpreting latent representations without or with minimal human intervention is correlating the latent representations with phenotypes directly measurable from the images: for example, correlating latent representations of proteins to their organelles' localization Kobayashi et al. (2022), correlating latent representations of cells to their local density Soelistyo et al. (2022) and



Figure 2: Saliency map-based (bottom right) versus counterfactual-based (top) explanations. Saliency methods highlight regions of more importance to the DL model's decision. Counterfactual methods generate artificial realistic images with exaggerated classification-driving phenotypes thus altering the model's prediction. Exaggerated phenotypes in this example are color, size and number of dark spots.

correlating latent representations of IVF embryos to their size Rotem et al. (2023). Of course, this 138 interpretability approach relies on existing measurements and cannot be used to discover explanations 139 that could not be measured in advance. We recently proposed a two-step interpretability method that 140 used counterfactual explanations to assign semantic image properties to latent features, followed by 141 SHAP ranking of these interpreted features to interpret individual predictions Rotem et al. (2023). 142 This approach combined global and local interpretability to determine the cause of classification (and 143 misclassification) of the morphological quality of IVF embryo instances. Another recent and creative 144 approach used counterfactual explanations to generate attention maps of the most discriminative 145 features and measured their contribution to the prediction Eckstein et al. (2024). This was applied to 146 147 reveal morphological differences between different Drosophila melanogaster synapses in electron microscopy images. 148

149 Outlook

The ever-increasing volume and complexity of bioimage data is making DL a crucial component of modern cell biology. We strongly believe that advancing XAI for bioimaging will propel cell biology forward by improving the trustworthiness of DL models and enabling the design of interpretability 'discovery machines' that can reveal new mechanistic understanding exceeding human intuition. We identify the following themes as having high potential to advance us in this promising direction of XAI visual-driven hypothesis generation.

156 Visual interpretability of image-to-image transformations

The current focus of visual interpretability is in the realm of classification models, neglecting 157 the interpretability of DL-based image-to-image transformations that are the backbone of many 158 bioimaging applications, ranging from image restoration to denoising, segmentation and cross-159 modality transformations. Current methods for the explainability of biomedical image-to-image 160 transformations mostly rely on pixel-based uncertainty estimations and can be applied to learn what 161 information in the input image was used for the output prediction. Interpretability of image-to-image 162 transformations could have practical implications in identifying when a model does not perform well. 163 For example, technical (for example, batch effects, uneven illumination) or biological (for example, 164 rare cell states, perturbations) out-of-distribution data can lead to errors and hallucinations and should 165 be considered during downstream analysis. We call for the development of methods for the visual 166

167 interpretation of bioimage image-to-image transformations that go beyond the pixel scale to the scale 168 of semantic objects to enable biologically meaningful interpretation.

169 Validating and measuring visual interpretability

Confirming visual interpretations of a non-trivial phenotype is hard because measuring the confidence 170 in an interpretation and assessing whether one visual interpretation is better than another constitute 171 an ill-defined problem. In practice, interpretability is usually subjectively validated in the final figure 172 of a DL bioimaging paper, reporting the researchers' intuitions based on several 'representative' 173 examples. More systematic validation mechanisms include assessing the robustness of the explanation 174 by demonstrating that different models for the same task have similar interpretations and assessing 175 sensitivity by showing that slight changes in the image space or the latent space ('adversarial attacks') 176 that do not change the prediction also do not change the interpretation. If the interpreted phenotype can 177 be explicitly measured, then it can be quantified by carefully (that is, avoiding confounders) correlating 178 latent representations with phenotypes. In more complicated situations in which the interpreted 179 180 phenotypes cannot be empirically measured, solutions include performing expert validation and evaluations DeGrave et al. (2023); Rotem et al. (2023) or altering the image (directly, not through the 181 latent representation) according to the interpretation and then evaluating whether the model prediction 182 was changed accordingly DeGrave et al. (2023); Eckstein et al. (2024), . When live imaging is 183 possible, one can examine spontaneous transitions between classification states that reduce most of 184 the variability confounding the human ability to interpret the subtle phenotypes in snapshot images 185 Zaritsky et al. (2021). With the understanding that complex models are less interpretable, a possible 186 approach to promote interpretability is reducing the model's complexity under the assumption that 187 this will increase its 'interpretability potential'. For example, it is possible to constrain the training of 188 the DL toward sparse representations and demonstrate that these representations perform well enough 189 and are more interpretable Soelistvo and Lowe (2024). Devising more objective and systematic 190 metrics for visual interpretability remains an open challenge that we believe is critical for advancing 191 the field. 192

193 More challenges and opportunities

Other promising future directions include developing better visualization approaches to ease the 194 translation from explanation to interpretation, developing interpretability methods for more complex 195 problems such as multi-class classification, effectively incorporating of 3D and/or temporal informa-196 tion in the interpretability method, and improving interpretability of high-dimensional image data (for 197 example, single-cell spatial omics). Another opportunity lies in designing latent representations that 198 can encourage better interpretability-driven discovery, such as disentangled representations in which 199 each latent feature encodes a single semantic image property. Additional open questions that warrant 200 systematic future investigations include "How much is the interpretability capacity conditioned on 201 202 the accuracy of the model being interpreted?" and "Can we benefit from automated analysis of visual 203 explanations, such as providing objective readouts indicating that a certain prediction cannot be trusted?" Advancing in these directions will promote effective data-driven discovery by relying on 204 DL's remarkable capacity to detect complex bioimage patterns to visually guide human interpretation 205 and hypothesis generation 206

207 **References**

Alex J DeGrave, Zhuo Ran Cai, Joseph D Janizek, Roxana Daneshjou, and Su-In Lee. 2023. Auditing
 the inference processes of medical-image classifiers by leveraging generative AI and the expertise
 of physicians. *Nature Biomedical Engineering* (2023), 1–13.

Michael Doron, Théo Moutakanni, Zitong S Chen, Nikita Moshkov, Mathilde Caron, Hugo Tou vron, Piotr Bojanowski, Wolfgang M Pernice, and Juan C Caicedo. 2023. Unbiased single-cell
 morphology with self-supervised vision transformers. *bioRxiv* (2023).

Nils Eckstein, Alexander Shakeel Bates, Andrew Champion, Michelle Du, Yijie Yin, Philipp Schlegel,
Alicia Kun-Yang Lu, Thomson Rymer, Samantha Finley-May, Tyler Paterson, et al. 2024. Neurotransmitter classification from electron microscopy images at synaptic sites in Drosophila
melanogaster. *Cell* 187, 10 (2024), 2574–2594.

- Hirofumi Kobayashi, Keith C Cheveralls, Manuel D Leonetti, and Loic A Royer. 2022. Self supervised deep learning encodes high-resolution features of protein subcellular localization.
 Nature methods 19, 8 (2022), 995–1003.
- Julienne LaChance, Kevin Suh, Jens Clausen, and Daniel J Cohen. 2022. Learning the rules of collective cell migration using deep attention networks. *PLoS computational biology* 18, 4 (2022), e1009293.
- Alexis Lamiable, Tiphaine Champetier, Francesco Leonardi, Ethan Cohen, Peter Sommer, David
 Hardy, Nicolas Argy, Achille Massougbodji, Elaine Del Nery, Gilles Cottrell, et al. 2023. Revealing
 invisible cell phenotypes with conditional generative modeling. *Nature Communications* 14, 1
 (2023), 6386.
- Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions.
 Advances in neural information processing systems 30 (2017).
- Anastasia Razdaibiedina, Alexander Brechalov, Helena Friesen, Mojca Mattiazzi Usaj, Myra
 Paz David Masinas, Harsha Garadi Suresh, Kyle Wang, Charles Boone, Jimmy Ba, and Brenda
 Andrews. 2024. PIFiA: self-supervised approach for protein functional annotation from single-cell
 imaging data. *Molecular systems biology* 20, 5 (2024), 521–548.
- Oded Rotem, Tamar Schwartz, Ron Maor, Yishay Tauber, Maya Tsarfati Shapiro, Marcos Meseguer, Daniella Gilboa, Daniel S. Seidman, and Assaf Zaritsky. 2023. Visual interpretability of image-based classification models by generative latent space disentanglement applied to in vitro fertilization. *bioRxiv* (2023). https://doi.org/10.1101/2023.11.15.566968 arXiv:https://www.biorxiv.org/content/early/2023/11/17/2023.11.15.566968.full.pdf
- Matthew S Schmitt, Jonathan Colen, Stefano Sala, John Devany, Shailaja Seetharaman, Alexia
 Caillier, Margaret L Gardel, Patrick W Oakes, and Vincenzo Vitelli. 2024. Machine learning
 interpretable models of cell mechanics from protein images. *Cell* 187, 2 (2024), 481–494.
- Christopher J Soelistyo and Alan R Lowe. 2024. Discovering interpretable models of scientific image
 data with deep learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6884–6893.
- Christopher J Soelistyo, Giulia Vallardi, Guillaume Charras, and Alan R Lowe. 2022. Learning
 biophysical determinants of cell fate with deep neural networks. *Nature machine intelligence* 4, 7
 (2022), 636–644.
- Takaki Yamamoto, Katie Cockburn, Valentina Greco, and Kyogo Kawaguchi. 2022. Probing the
 rules of cell coordination in live tissues by interpretable machine learning based on graph neural
 networks. *PLOS Computational Biology* 18, 9 (2022), e1010477.
- Karren Dai Yang, Karthik Damodaran, Saradha Venkatachalapathy, Ali C Soylemezoglu, GV Shiv ashankar, and Caroline Uhler. 2020. Predicting cell lineages using autoencoders and optimal
 transport. *PLoS computational biology* 16, 4 (2020), e1007828.
- Assaf Zaritsky, Andrew R Jamieson, Erik S Welf, Andres Nevarez, Justin Cillay, Ugur Eskiocak,
 Brandi L Cantarel, and Gaudenz Danuser. 2021. Interpretable deep learning uncovers cellular
 properties in label-free live cell images that are predictive of highly metastatic melanoma. *Cell systems* 12, 7 (2021), 733–747.

258 A Appendix / supplemental material