

¹ Histomics Label

² **Brianna Major**  ¹, **Jeffery A. Goldstein**  ², **Michael Nagler**  ¹, **Lee A. Newberg**  ¹, **Abhishek Sharma**  ², **Anders Sildnes**  ², **Faiza Ahmed**  ¹, **Jeff Baumes**  ¹, **Lee A. D. Cooper**  ², and **David Manthey**  ¹

⁵ 1 Kitware, Inc., New York, United States ² Northwestern University Feinberg School of Medicine, Illinois,
⁶ United States

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [julia ferraioli](#) 

Reviewers:

- [@cinnetcrash](#)
- [@Mmasoud1](#)

Submitted: 20 June 2025

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

⁷ Summary

⁸ Histomics Label is a software tool for the interactive development of machine-learning
⁹ classifiers for whole slide pathology images. It is deployed as part of the Digital Slide Archive
¹⁰ ([Gutman et al., 2017](#); [Kitware, Inc, 2025a](#)), a web-based data management system for whole
¹¹ slide image datasets, and was built on top of HistomicsUI ([Kitware, Inc, 2025c](#)) and uses the
¹² HistomicsTK ([Kitware, Inc, 2025b](#)) image analysis tool kit.

¹³ Users label image regions or tissue structures to provide training data for classifiers and
¹⁴ iteratively improve these classifiers by reviewing their output and providing additional labels.
¹⁵ The interface uses heuristics to guide users to the most impactful examples to label, and
¹⁶ supports bulk labeling of samples and review of labeled examples for collaboration. An example
¹⁷ data generation pipeline is included that segments a whole slide image into superpixels, and
¹⁸ generates feature embeddings for segmented regions using a foundation model.

Statement of need

²⁰ One of the limitations in developing classification models is the need for labeled data. In
²¹ pathology and other medical fields, the expertise required for labeling and busy schedules of
²² medical experts make labeling particularly challenging. For whole slide images, where each
²³ image can contain several billion pixels, navigating vast datasets in search of possibly rare
²⁴ tissue states can be very inefficient and frustrating. Software interfaces need to be optimized
²⁵ for the user experience and make the most of an expert's time and energy.

²⁶ Other issues in labeling include the volume and accessibility of data. Software that must run
²⁷ local to the data requires that all data be copied and correctly versioned for the project. Using
²⁸ a web-client and server model with appropriate permission models, only requires that the data
²⁹ be on centrally managed server. This allows there to be a single, coordinate source of data
³⁰ for a project, and reduces the burden on individual users to only requiring a web browser and
³¹ ordinary internet connection. This enables collaboration between multiple experts, or to allow
³² experts to review the work of their trainees.

³³ Histomics Label uses a technique called active learning to identify the unlabeled examples
³⁴ that can provide the most benefit to classifier performance and provides an intuitive workflow
³⁵ for presenting these examples to experts for efficient labeling. Data can be generated using a
³⁶ built-in pipeline that partitions whole-slide images into superpixels, or users can provide their
³⁷ own data from external cell or tissue segmentation algorithms. Users specify the categories
³⁸ that can be labeled and assign display properties like color, and can exclude categories from
³⁹ classifier training (for instance, for regions whose categories cannot be accurately determined).
⁴⁰ After labeling a few initial example regions, a classifier is trained and used to both predict the
⁴¹ category of all regions and the unlabeled regions that provide the most classifier benefit. The

42 user can retrain the classifier at any time and review the classifier predictions and labels from
 43 other users. Labeling can also be performed by painting directly on the whole slide image with
 44 a brush tool.

45 For development, the initial segmentation uses superpixels generated with the SLIC ([Achanta et al., 2012](#)) algorithm. These are computed on whole slide images in a tiled manner so
 46 that they can work on arbitrarily large images, and the tile boundaries are properly handled
 47 to merge seamlessly. Once generated, segments are represented in one of two ways, either
 48 as two-dimensional patches, each centered in a fixed-sized square of masked pixels, or as
 49 one-dimensional feature embeddings, such as those generated from the huggingface UNI ([Chen et al., 2024](#)) foundation model. One of two basic models is trained based upon the segment
 50 representation. For two-dimensional patches, the model to be trained is a small-scale CNN
 51 implemented in tensorflow/keras or torch. For one-dimensional vectors, the model to be
 52 trained is a single-layer linear classifier. The certainty criteria for which segments should be
 53 labeled next can also be selected, and includes confidence, margin, negative entropy, and the
 54 BatchBALD ([Kirsch et al., 2019](#)) algorithm.

55 We had a placental pathologist provide feedback to validate the efficiency of the user interface
 56 and utility of the process.

59 Basic Workflow

60 When starting a new labeling project, the user selects how superpixels are generated, which
 61 certainty metric is used for determining the optimal labeling order, and what features are used
 62 for model training. The labeling mode allows defining project labels and performing initial
 63 labeling. This mode can also be used to add new label categories or combine two categories if
 64 they should not have been distinct. Label categories can additionally be marked as excluded,
 65 which removes them from training and ensures that superpixels with those labels are no longer
 66 suggested for labeling.

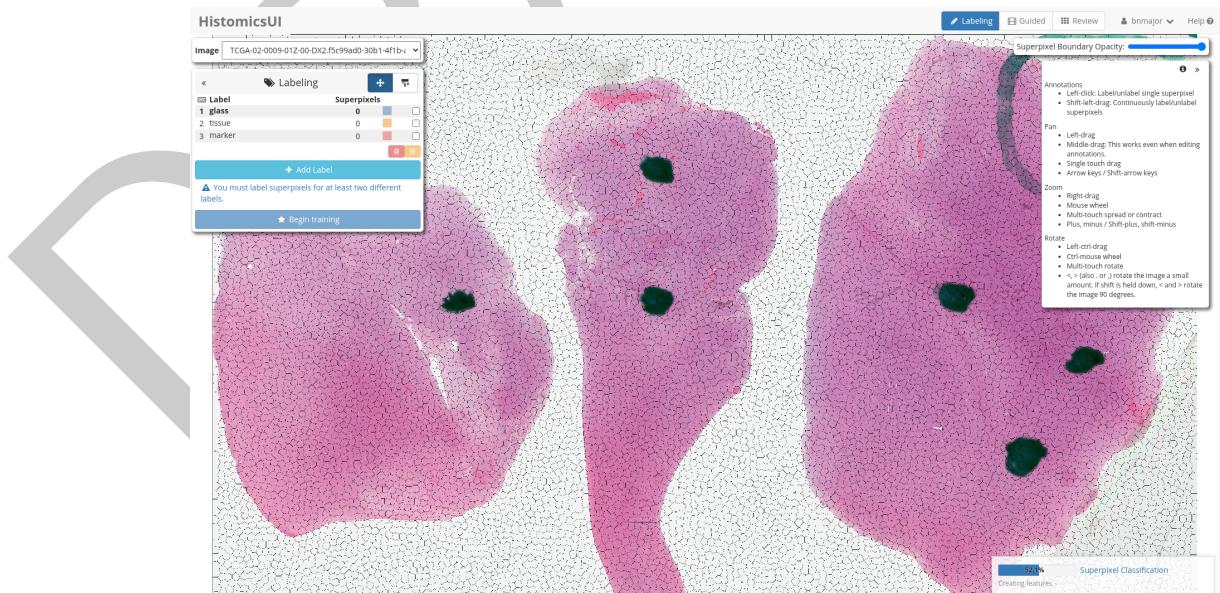


Figure 1: The Bulk Labeling interface showing one of the project images divided into superpixels with some categories defined. A user can “paint” areas with known labels as an initial seed for the guided labeling process

67 Once some segments have been labeled and an initial training process has been performed,
 68 additional segments are shown with their predictions. The user can use keyboard shortcuts or

69 the mouse to confirm or correct labels. These are presented in an order that maximizes the
 70 utility of improving the model based on the originally selected certainty metric.

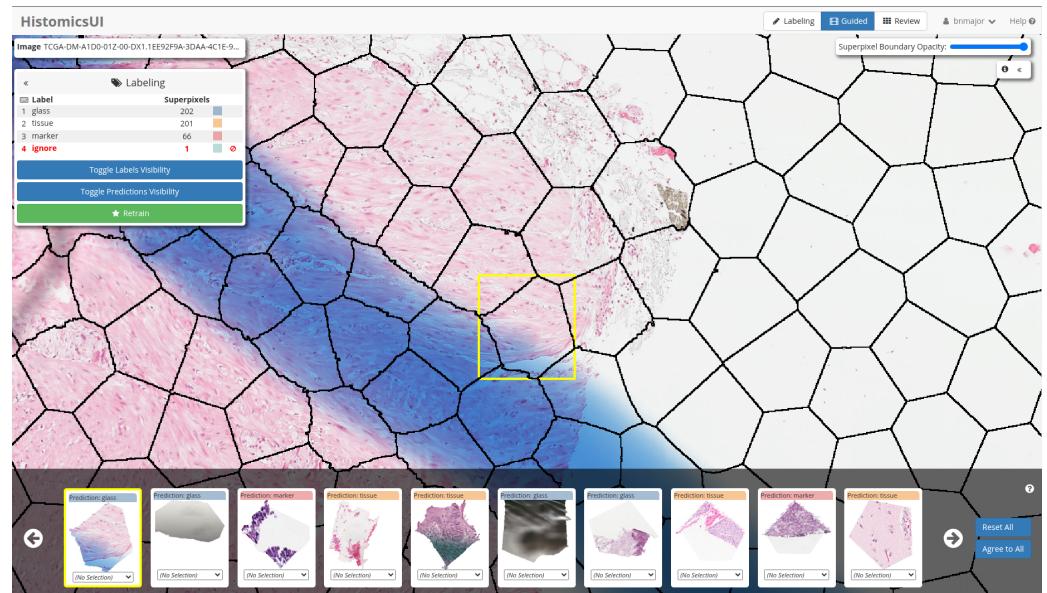


Figure 2: The Guided Labeling interface showing a row of superpixels to be labeled and part of a whole slide image

71 To check on overall behavior or correct mistakes, there is a review mode that allows seeing
 72 all labeled segments with various filtering and sorting options. This can be used to check
 73 agreement between pathologists or determine how well the model agrees with the manually
 74 labeled data.

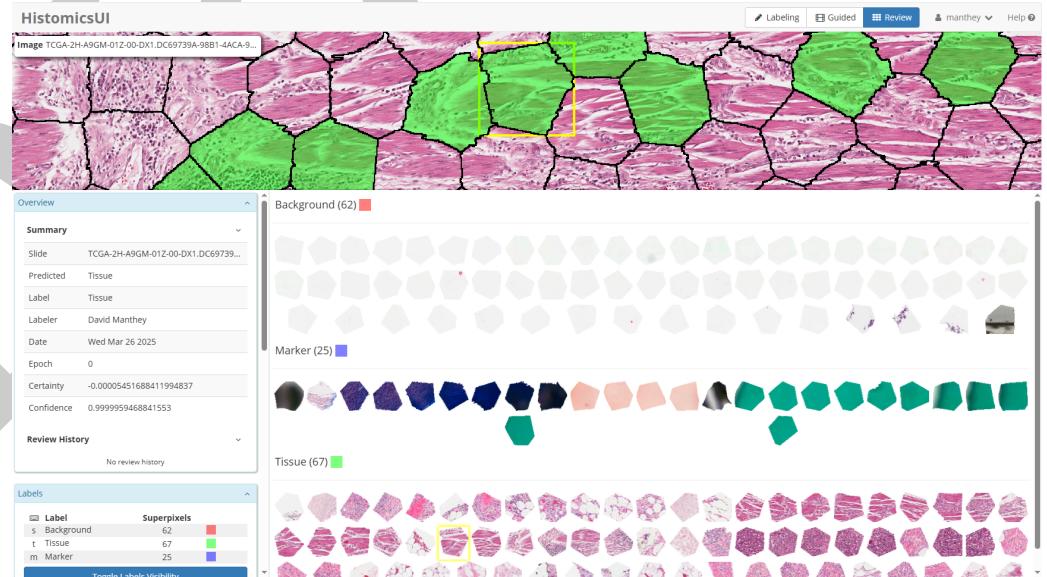


Figure 3: The Review interface showing labeled superpixels in each category

75 The whole slide image data in these figures are from data generated by the TCGA Research
 76 Network (?).

77 Acknowledgements

78 This work has been funded in part by National Library of Medicine grant 5R01LM013523
79 entitled "Guiding humans to create better labeled datasets for machine learning in biomedical
80 research".

81 References

- 82 Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., & Süstrunk, S. (2012). SLIC superpixels
83 compared to state-of-the-art superpixel methods. *IEEE Transactions on Pattern Analysis
84 and Machine Intelligence*, 34(11), 2274–2282. <https://doi.org/10.1109/TPAMI.2012.120>
- 85 Chen, R. J., Ding, T., Lu, M. Y., Williamson, D. F. K., Jaume, G., Song, A. H., Chen,
86 B., Zhang, A., Shao, D., Shaban, M., Williams, M., Oldenburg, L., Weishaupt, L. L.,
87 Wang, J. J., Vaidya, A., Le, L. P., Gerber, G., Sahai, S., Williams, W., & Mahmood, F.
88 (2024). Towards a general-purpose foundation model for computational pathology. *Nature
89 Medicine*, 30(3), 850–862. <https://doi.org/10.1038/s41591-024-02857-3>
- 90 Gutman, D. A., Khalilia, M., Lee, S., Nalisnik, M., Mullen, Z., Beezley, J., Chittajallu, D. R.,
91 Manthey, D., & Cooper, L. A. D. (2017). The digital slide archive: A software platform for
92 management, integration, and analysis of histology for cancer research. *Cancer Research*,
93 77(21), e75–e78. <https://doi.org/10.1158/0008-5472.can-17-0629>
- 94 Kirsch, A., Amersfoort, J. van, & Gal, Y. (2019). BatchBALD: Efficient and diverse batch
95 acquisition for deep bayesian active learning. *CoRR*, abs/1906.08158. <http://arxiv.org/abs/1906.08158>
- 97 Kitware, Inc. (2025a). *Digital slide archive: A system for working with large microscopy
98 images*. https://github.com/DigitalSlideArchive/digital_slide_archive
- 99 Kitware, Inc. (2025b). *HistomicsTK: A python package for the analysis of digital pathology
100 images*. <https://doi.org/10.5281/zenodo.14833780>
- 101 Kitware, Inc. (2025c). *HistomicsUI: Organize, visualize, annotate, and analyze histology
102 images*. <https://doi.org/10.5281/zenodo.5474914>