
OpenSRH: optimizing brain tumor surgery using intraoperative stimulated Raman histology

Cheng Jiang^{1*} Asadur Chowdury^{1*} Xinhai Hou^{1*}
Akhil Kondepudi¹ Christian W. Freudiger² Kyle Conway¹
Sandra Camelo-Piragua¹ Daniel A. Orringer³ Honglak Lee¹ Todd C. Hollon¹

¹University of Michigan ²Invenio Imaging ³New York University *Equal Contribution

{chengjia, achowdur, xinhaih, tocho}@umich.edu <https://opensrh.mlins.org>

Abstract

Accurate intraoperative diagnosis is essential for providing safe and effective care during brain tumor surgery. Our standard-of-care diagnostic methods are time, resource, and labor intensive, which restricts access to optimal surgical treatments. To address these limitations, we propose an alternative workflow that combines stimulated Raman histology (SRH), a rapid optical imaging method, with deep learning-based automated interpretation of SRH images for intraoperative brain tumor diagnosis and real-time surgical decision support. Here, we present *OpenSRH*, the first public dataset of clinical SRH images from 300+ brain tumors patients and 1300+ unique whole slide optical images. OpenSRH contains data from the most common brain tumors diagnoses, full pathologic annotations, whole slide tumor segmentations, raw and processed optical imaging data for end-to-end model development and validation. We provide a framework for patch-based whole slide SRH classification and inference using weak (i.e. patient-level) diagnostic labels. Finally, we benchmark two computer vision tasks: multiclass histologic brain tumor classification and patch-based contrastive representation learning. We hope OpenSRH will facilitate the clinical translation of rapid optical imaging and real-time ML-based surgical decision support in order to improve the access, safety, and efficacy of cancer surgery in the era of precision medicine. Dataset access, code, and benchmarks are available at <https://opensrh.mlins.org>.

1 Introduction

The optimal surgical management of brain tumors varies widely depending on the underlying pathologic diagnosis [1]. Surgical goals range from needle biopsies (e.g. primary central nervous system lymphoma [2]) to supramaximal resections (e.g. diffuse gliomas [3]). A major obstacle to the precision care of brain tumor patients is that the pathologic diagnosis is usually *unknown* at the time of surgery. For other tumor types, such as breast or lung cancer, diagnostic biopsies are obtained prior to definitive surgical management, which provides essential clinical information used to inform the goals of surgery. Routine diagnostic biopsies in neuro-oncology are not feasible due to high surgical morbidity and the potential for permanent neurologic injury. Consequently, the importance of *intraoperative* pathologic diagnosis in brain tumor surgery has been recognized for nearly a century[4].

Unfortunately, our current intraoperative pathologic techniques are time, resource, and labor intensive [7, 8]. Conventional diagnostic methods, including frozen sectioning and cytologic preparations, require an extensive pathology infrastructure for tissue processing and specimen analysis by a board-

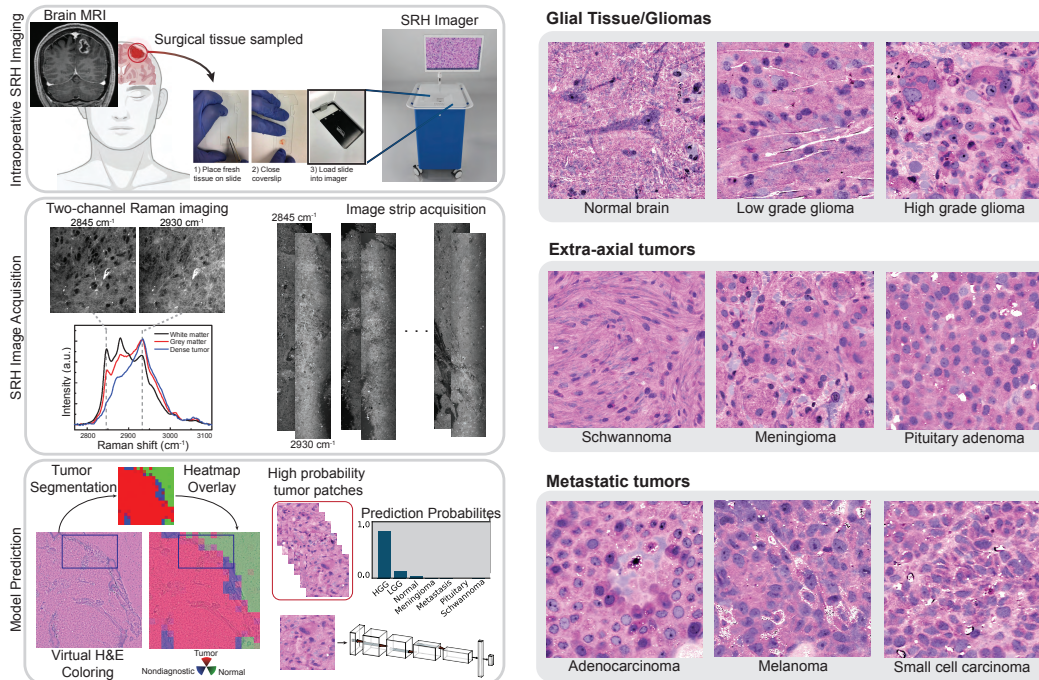


Figure 1: **Left**, A patient with a newly diagnosed brain lesion undergoes a surgery for tissue diagnosis and/or tumor resection. The tumor specimen is sampled from the patient’s tumor and directly loaded into a premade, disposable microscope slide. The specimen is placed into the SRH imager for rapid optical imaging. SRH images are acquired sequentially as strips at two Raman shifts, 2845 cm^{-1} , and 2930 cm^{-1} . The size and number of strips to be acquired are set by the operator who defines the desired image size. Standard image sizes range from $1\text{-}5\text{mm}^2$ and image acquisition time ranges from 30 seconds to 3 minutes. The strips are edge clipped, field flattened and co-registered to generate whole slide SRH images. Images can be colored using a custom virtual H&E colorscheme for pathologic review [5]. The whole slide image is divided into non-overlapping 300×300 pixel patches and each patch undergoes a feedforward pass through a previously trained tumor segmentation model [6] to segment the patches into tumor regions, normal brain, and nondiagnostic regions. The tumor patches are then used for both training and inference to predict the patient’s brain tumor diagnosis. **Right**, Examples of virtually-colored SRH images from brain tumor diagnoses included in OpenSRH. We include a diversity of tumor diagnoses that cover the most common brain tumor subtypes.

certified neuropathologist [9]. While the conventional pathology workflow with board certified neuropathologist interpretation has a diagnostic accuracy between 86 - 96% [5], the pathologist workforce in the US declined in absolute and population-adjusted numbers by nearly 20% between 2007-2017 [10]. This decline has unevenly affected neuropathology, with a 40% fellowship vacancy rate and it is projected to worsen [11]. The number of medical centers performing brain tumor surgery outnumbers board-certified neuropathologists, reducing patient access to expert intraoperative consultation and, consequently, optimal surgical management.

An ideal system for surgical specimen analysis and intraoperative tumor diagnosis would be accessible, fast, standardized, and accurate. An intraoperative pathology system requires, at minimum, (1) a data/image acquisition modality and (2) a diagnostic interpretation. Conventional intraoperative pathology uses light microscopy interpreted by a neuro-pathologist (>20-30 minutes). We propose an innovative diagnosis system that uses a rapid (2-3 minutes, $>10\times$ speedup), label-free optical histology method, called stimulated Raman histology (SRH), combined with deep learning-based interpretation of fresh, unprocessed surgical specimens. We have previously demonstrated the feasibility of large-scale clinical SRH imaging [5] and the use of deep neural networks for SRH image classification of brain tumor patients [6, 12, 13]. These studies demonstrate the potential

for AI-based diagnosis and interpretation of SRH images to better inform brain tumor surgery and provide personalized surgical goals in the era of precision medicine.

Here, we seek to facilitate this area of active research by releasing *OpenSRH*, a collection of intraoperative SRH data, including raw SRH acquisition data, processed high-resolution image patches for model development, virtually-stained whole slide images, semantic segmentation of tumor regions, and full intraoperative diagnostic annotations. OpenSRH is the first and only publicly available dataset of any human cancer imaged using optical histology. We release the OpenSRH dataset with the intention to foster translational AI research within the field of precision oncology. The main contributions of this work are:

1. **OpenSRH dataset:** We curate and open-source the largest dataset of intraoperative SRH images with pathologic annotations to facilitate the development of innovative machine learning solutions to improve brain tumor surgery.
2. **Classification benchmarks:** We benchmark performance for patch-based histologic brain tumor classification across multiple tumor types, computer vision architectures, and transfer learning methods.
3. **Contrastive representation learning benchmarks:** We evaluate both self-supervised and weakly supervised patch contrastive learning methods for SRH representation learning. Contrastive learning methods are evaluated using linear evaluation protocols and benchmarked as a model pretraining strategy.

2 Background

Stimulated Raman Histology SRH is based on Raman scattering. Raman scattering occurs when incident photons on a media either gain or lose energy when scattered (i.e. inelastic scattering), shifting the frequency/wavenumber of the scattered photons. This Raman shift can be measured to characterize the biochemical composition of both inorganic and organic materials using narrow-band laser excitation and a spectrometer [14]. A major limitation of using spontaneous Raman scattering for biochemical analysis is that the Raman effect is weak compared to elastic scattering. Therefore, long acquisition times (> 30 minutes) and spectral averaging are required to obtain representative biochemical spectra. Stimulated Raman scattering (SRS) microscopy was discovered in 2008 as a highly sensitive, label-free biomedical imaging method [15]. Rather than acquiring broad-band spectra, SRS microscopy uses a second laser excitation source to achieve non-linear amplification of narrow-band Raman spectral regions that correspond to specific molecular vibrational modes (see Figure 1). SRS images can then be generated at specific narrow-band Raman wavenumbers. Translational research led to the development of a fiber-laser-based SRS imaging system that could be used at the patient’s bedside to generate rapid histologic images of fresh surgical specimens, called SRH [5, 16, 17]. A major advantage of SRH over other histologic imaging methods is that image contrast is generated by the intrinsic biochemical properties of the tissue only and does not require any tissue processing, staining, dyeing, or labelling (i.e. label-free).

ML applications for SRH Unlike conventional intraoperative histology with light microscopy, SRH provides high-resolution *digital* images that can be used directly for downstream ML tasks. Whole slide image digitization of frozen or paraffin-embedded tissue is slow and memory intensive, presenting a major bottleneck for its routine use in intraoperative histology, and clinical medicine in general [18]. Previous studies showed that SRH plus shallow ML models can be used to detect and quantify tumor infiltration in adult and pediatric fresh surgical specimens [5, 19, 20, 21]. We subsequently demonstrated that SRH combined with convolutional neural network architectures can be used for intraoperative diagnostic decision support [6, 12, 13]. These preliminary studies, while demonstrating the feasibility of applying deep architectures to SRH, did not include rigorous hyperparameter tuning, explicit representation learning, or ablation studies to optimize model performance. Moreover, all previous studies required manual annotations, including dense patch-level annotations [6], for model training.

3 Related Work

To date, no SRH datasets are publicly available. The work most directly related to OpenSRH comes from digital and computational pathology research. **The Cancer Genome Atlas (TCGA)** and **The Cancer Imaging Atlas (TCIA)** include a large repository of digitized histopathology slides processed using hematoxylin and eosin (H&E) staining. Many studies have used this dataset for image classification tasks across several cancer types, including, but not limited to, lung [22], gastrointestinal [23], prostate [24, 25], brain [26], and pan-cancer studies [23, 27, 28]. Another related histopathology dataset comes from the CAMELYON16 research challenge [29]. The challenge is to detect lymph node metastases in women with breast cancer. Digital pathology remains an active area of research in precision oncology. However, ML applications in digital pathology are mainly applied to postoperative tissue assessment and do not play a major role in informing cancer surgery.

One application of SRH is the detection of tumor infiltration in real-time to improve the extent of tumor resection and reduce residual tumor burden. Real-time SRH-based tumor delineation has been studied in sinonasal/skull base cancers [13, 30, 31] and diffuse gliomas [20, 32]. OpenSRH provides the necessary dataset to explore this topic for multiple brain tumor types, including metastatic tumors and extra-axial tumors, such as meningiomas (Figure 1).

Overall Need High-quality, public, biomedical datasets with expert annotations are rare. Moreover, the clinical significance of some existing datasets is unclear due to the lack of a roadmap for clinical translation [33]. We believe that OpenSRH has the potential to address a currently unmet clinical need of improving cancer surgery in order to advance precision oncology, both in the US and globally [7]. OpenSRH can address a pressing and significant clinical problem, while having high translational potential because, as previously mentioned, an ideal system for intraoperative tumor specimen evaluation should be:

1. **Accessible:** SRH imaging systems are FDA-approved and commercially available for intraoperative imaging.
2. **Fast:** imaging acquisition time and time-to-diagnosis is $10\times$ faster than the current standard-of-care H&E histology.
3. **Standardized:** SRH image acquisition is invariant to patient demographic features, clinical workforce, and geographic location.
4. **Accurate:** preliminary results [6, 13] and diagnostic performance benchmarks (see Figure 4) are on par with the pathologist-based interpretation of H&E histology.

4 Data Description

Patient population Patients were consecutively and prospectively enrolled for intraoperative SRH imaging. This study was approved by Institutional Review Board (HUM00083059). Informed consent was obtained for each patient prior to SRH imaging and approved the use of tumor specimens for research and development. All patient health information (PHI) are removed from all OpenSRH data. The inclusion criteria are (1) patients with planned brain tumor or epilepsy surgery at Michigan Medicine (UM), (2) subjects or durable power of attorney able to give informed consent, and (3) subjects in whom there was additional specimen beyond what was needed for routine clinical diagnosis.

SRH imaging Intraoperative SRH imaging and data processing workflow can be found in Figure 1. A small tumor specimen ($3\times 3\text{ mm}^3$) is placed into a premade microscope slide, which is then loaded into the commercially available NIO Imaging System (Invenio Imaging, Inc.) for SRH imaging. The tissue is then excited with a dual-wavelength fiber laser source, which provides spectral access to Raman shifts in the range of 2800 cm^{-1} to 3130 cm^{-1} . SRH images are acquired at two Raman shifts: 2845 cm^{-1} highlights lipid-rich regions and 2930 cm^{-1} highlights DNA and protein-rich regions [6]. The images are acquired sequentially as 0.5 mm wide strips, stitched together and the two image channels are co-registered to generate the final whole slide image. The co-registration between the two image channels is performed using discrete Fourier transform. A virtual H&E colorscheme [5] can be applied to SRH images for clinician review, but is not used for model development.

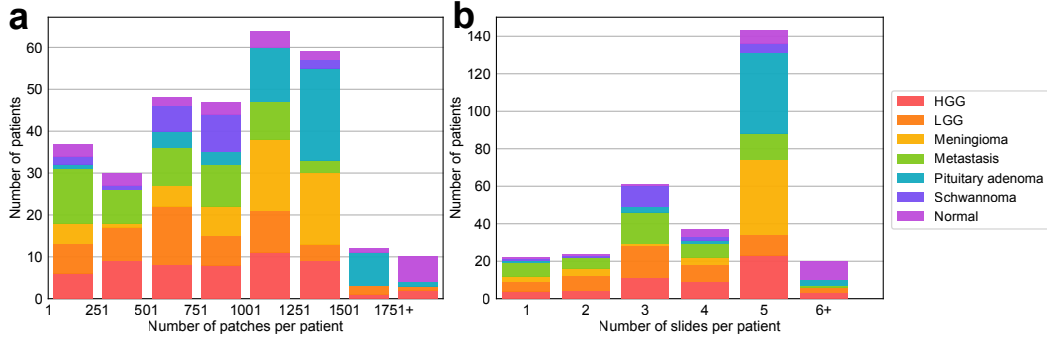


Figure 2: Histogram for the number of patches and slides per patient. **a** shows the total number of patches varies across the patients and, **b** shows most patients have 5 slides. The difference between these two distributions may be caused by specimen size, non-diagnostic regions, surgeon preference, etc. HGG, high grade glioma; LGG, low grade glioma.

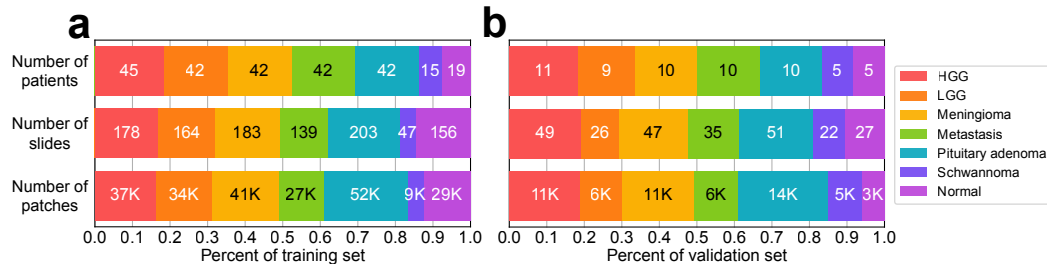


Figure 3: Bar chart for the number of patients, slides and patches for each diagnostic class. Validation set was randomly selected and contains approximately 20% of patients in OpenSRH (60/307 patients). Training and validation sets have approximately equivalent class distributions.

Image preprocessing Image processing starts by applying a sliding window over the two-channel image to generate 300×300 pixel non-overlapping patches. A third channel is obtained by performing a pixel-wise subtraction from the two registered channels ($2930\text{cm}^{-1} - 2845\text{cm}^{-1}$), which highlights the nuclear contrast and cellular density of the tissue [19]. The third channel is concatenated depth-wise to generate a final three-channel patch for model training and inference. Each patch then undergoes a feedforward pass through a pretrained segmentation model to classify the patch into tumor, normal brain, or non-diagnostic tissue [12]. The model was trained using manually labelled patches. The segmentation prediction for each patch is released as part of the OpenSRH dataset.

Dataset breakdown OpenSRH consists of 307 patients. A total of 304 patients underwent intra-operative SRH imaging and three patients had postmortem specimen collection. We strategically selected the most common brain tumor types to be included in OpenSRH. The included brain tumor diagnoses cover more than 90% of all newly diagnosed brain tumors in the US [34]. OpenSRH includes a diversity of brain tumor types, including primary brain tumors (high-grade gliomas, low-grade gliomas), secondary brain tumors (metastases), and extra-axial tumors (meningiomas, schwannomas, pituitary adenomas). A panel of patch samples is included in Figure 1. The dataset is randomly divided into training (247 patients) and validation set (60 patients, about 20%). Figure 3 shows a distribution of the number of patients, slides, and patches per class in the training and validation set. Technical details of the data release and companion source code are in Appendix A.

5 Histologic brain tumor classification benchmarks

In this section, we present the results of the baseline multiclass brain tumor classification task. We aim to benchmark the results for common training strategies. We investigate the value of

transfer learning/pretraining from natural image datasets, specifically ImageNet [35], for improving classification performance. The value of transfer learning for SRH [6] and medical imaging, in general [36], remains an active area of research. We selected representative models from the two most competitive computer vision architectures: convolutional neural networks (ResNet50 [37]) and vision transformers (ViT-S[38, 39]). These architectures were selected because they contain a similar number of parameters (~ 23.5 million for ResNet50, ~ 21.7 million for ViT-S).

5.1 Training protocol

ResNet50 In the ResNet50 architecture, we changed the output dimension of the last layer to 7 for our experiments. We used a batch size 96 and trained on 300×300 images with horizontal and vertical flipping of probability 0.5 as augmentations. We used categorical cross-entropy loss and AdamW optimizer [40] with $\beta_1 = 0.9$, $\beta_2 = 0.999$, and a weight decay of 0.001. The initial learning rate was 0.001, with a step scheduler with a decay rate $\gamma = 0.5$ every epoch. We trained for 20 epochs with two Nvidia RTX 2080Ti GPUs. Training wall time was ~ 9.5 hours for each experiment.

ViT-S ViT training protocol was adjusted based on previously published results[38, 39]. In addition to the augmentations in the ResNet50 protocol, we also resized the image to 224×224 to fit the standard ViT-S model and ImageNet pretrained weights [41]. We used AdamW as the optimizer, with the same parameters in ResNet50. The initial learning rate was $1E-4$, with a cosine learning rate scheduler. First 20% of training steps were set as the linear warm-up stage to increase training stability. We trained 20 epochs with a batch size of 256 for 9 hours using the same GPU resources mentioned above. Detailed training parameters are included in Appendix C.

5.2 Prediction aggregation and benchmark metrics

Patch-level predictions from the same whole slide or patient need to be aggregated to generate a slide- or patient-level prediction. We aggregated patch-level logits after softmax using average pooling to compute slide and patient-level prediction. We preferred using average pooling over hard patch voting to retain the full patch-level model predictions during slide- or patient-level inference. Model performance evaluation metrics include top-1 accuracy, mean class accuracy (MCA), and mean average precision (MAP). Additional classification metrics including top-2 accuracy and false negative rate (tumor vs. normal) can be found in Appendix D.

5.3 Experimental results

Patch- and patient-level results can be found in Table 1. Patient-level metrics are generally higher than patch-level metrics. Patch-level prediction errors can be mitigated through the average pooling aggregation function. In our preliminary benchmark, ResNet50 achieved overall better performance than ViT-S (e.g., by 7.2 patch accuracy and 5.6 patient accuracy). A potential reason is due to ViT requiring large-scale image datasets on the scale of ImageNet21K or JFT300M [38] to overcome low inductive bias. Insufficient pretraining is known to result in worse performance compared to convolutional neural networks (CNNs). We did observe improved patch-level performance when using ImageNet pretraining (2.1 for ResNet50 and 6.5 for ViT-S at patch accuracy). In general, pretraining was more beneficial to ViT than ResNet50. We believe vision transformers may outperform CNNs with data efficient pretraining. Figure 4 summarizes the patient-level confusion matrix. Both models had similar diagnostic errors differentiating HGG and LGG, a known challenging diagnostic task for pathologists and computer vision models [5]. Metastatic tumors have diverse histologic features (see Figure 1) that can result in diagnostic errors across multiple classes [6]. From the confusion matrices in figure 4, it is important to note that we can also observe some false negatives in the model prediction (tumor vs. normal). Additional metrics on false negative rate for these experiments are also included in appendix D.

6 Contrastive representation learning benchmark

Our previous studies demonstrate that contrastive representation learning is well suited for patch-based representation learning [42]. The focus for this section is to investigate the effectiveness of contrastive learning strategies for OpenSRH. We used both unsupervised contrastive learning (SimCLR [43]) and

Backbone	Pretrain	Patch			Patient		
		Accuracy	MCA	MAP	Accuracy	MCA	MAP
ResNet50	Random	84.4 (0.4)	83.8 (0.5)	89.5 (0.5)	90.0 (0.0)	91.4 (0.0)	92.8 (0.2)
ResNet50	ImageNet	86.5 (0.4)	85.6 (0.3)	91.2 (0.3)	88.9 (0.8)	90.5 (0.6)	94.0 (0.1)
ViT-S	Random	77.2 (0.5)	76.8 (0.8)	82.3 (0.5)	85.0 (1.4)	87.2 (1.1)	93.2 (0.4)
ViT-S	ImageNet	83.7 (0.5)	82.7 (0.9)	88.8 (0.1)	88.9 (0.8)	90.5 (0.6)	93.9 (0.4)

Table 1: Classification benchmarks for ResNet50 and ViT-S. Pretrain refers to the pretraining strategy. Each experiment included three random initial seeds. Mean value and standard deviation (in parentheses) for each metric are reported here. The full table including false negative rates and slide-level metrics can be found in the Appendix D.

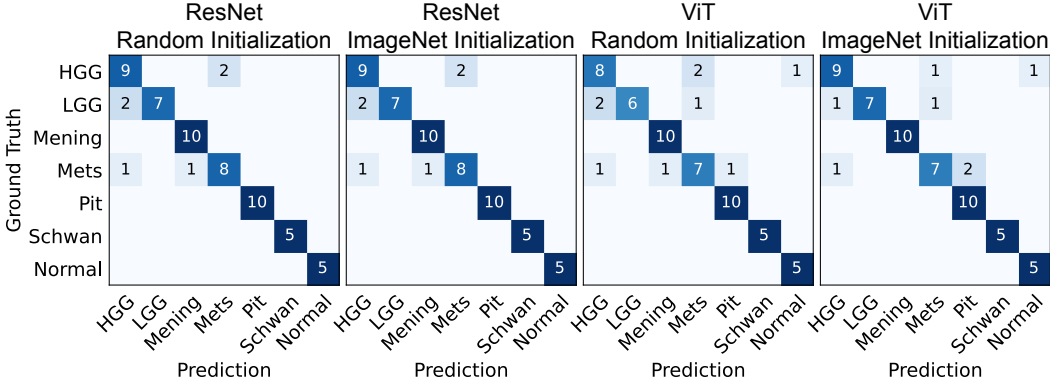


Figure 4: Patient-level confusion matrices for the four different training strategies on the validation set. Most of the errors occurred in the HGG, LGG, and metastasis classes. Only seed 1 is shown here, other seeds are included in Appendix D. Mening, meningioma; Mets, metastasis; Pit, pituitary adenoma; Schwan, schwannoma.

supervised contrastive learning (SupCon [44]) on ResNet50 and ViT-S architectures. The contrastive loss for SimCLR aims to solve the pretext task of instance discrimination. The model is trained to identify two augmented positive pairs of the same image from other images in a minibatch. SupCon loss has the similar training objective but aims at optimizing class discrimination. All images from the same class are treated as positive instances and other images are negative instances. Our general contrastive learning workflow uses SimCLR and SupCon as a representation learning strategy on our dataset followed with a linear evaluation protocol. We compared contrastive representation learning methods with ImageNet pretraining.

6.1 Training and evaluation protocol

Training protocol For both SimCLR and SupCon methods, we use the same protocol except for the loss function. We applied ResNet50 and ViT-S with a linear projection head to project the image representation to a low dimensional hypersphere (128 for ResNet50, 24 for ViT-S) to compute the contrastive loss. The data augmentation strategy followed [43]: a composition of multiple augmentations including flipping, color jittering, and Gaussian blur (for details, see Appendix C). We use AdamW optimizer for both models and same parameters as in Section 5.1. Different learning rates (1E-3 for ResNet50 and 5E-4 for ViT-S) were adopted for each model. We trained using a batch size 224 for ResNet50 and 512 for ViT-S for 40 epochs. We use linear warmup for the first 10% epochs and cosine decay scheduler for ViT only. Detailed protocols are included in Appendix C.

Evaluation protocol To evaluate the learned image representations, we followed a standard linear evaluation protocol [43, 45, 46, 47], where a linear classifier is trained on top of the frozen pretrained backbone. We consider the same evaluation metrics and aggregation function as in Section 5.2. Apart from the classification metrics, we performed qualitative evaluation of our learned representations

through t-distributed stochastic neighbor embedding (tSNE) [48] in Figure 5. Additional fine-tuning protocols and results are included in Appendix F.

6.2 Experimental results

Results of linear evaluations could be found in Table 2. Self-supervised representation learning with SimCLR was able to achieve a patient-level accuracy of 85.6 for ResNet50 and 78.3 for ViT-S. These results demonstrate improvement over our previous self-supervised classification performance [49]. Self-supervised contrastive representation learning using OpenSRH outperforms ImageNet transfer learning for patch-based metrics (e.g., by 3.2 for ResNet50 and 2.3 for ViT-S). These results emphasize the large domain gap between natural images and SRH optical images [50, 36]. Similar to other computer vision tasks, optimal representation learning can be achieved with additional supervision. SupCon outperforms both ImageNet pretraining and SimCLR in patch-based metrics. Linear evaluation showed an overall increase of 4.4 and 8.4 in patient-level accuracy between SupCon and SimCLR for ResNet50 and ViT-S, respectively. Interestingly, the patient-level metrics for pretrained ViT-S were prominently high, while the patch-level metrics were comparatively worse. We believe these results may be due to a simple soft voting aggregation of patch-level predictions. This opens the question for better (learnable) aggregation functions for SRH images. The tSNE plot in Figure 5 was consistent with our patch-based evaluation metrics for both models, where SupCon showed more discrete image representations.

Backbone	Methods	Patch			Patient		
		Accuracy	MCA	MAP	Accuracy	MCA	MAP
ResNet50	ImageNet	68.3 (0.0)	67.9 (0.0)	72.9 (0.1)	80.0 (0.0)	82.9 (0.0)	88.8 (0.1)
ResNet50	SimCLR	79.1 (0.4)	78.9 (0.4)	84.2 (0.6)	83.9 (1.0)	86.1 (0.9)	92.4 (0.1)
ResNet50	SupCon	87.5 (0.3)	86.8 (0.3)	91.5 (0.5)	90.0 (0.0)	91.4 (0.1)	94.6 (0.5)
ViT-S	ImageNet	71.8 (0.1)	71.1 (0.1)	77.1 (0.1)	88.3 (0.0)	89.8 (0.0)	93.9 (0.0)
ViT-S	SimCLR	76.8 (0.5)	76.3 (0.5)	82.5 (0.3)	80.0 (1.7)	83.0 (1.3)	92.3 (0.0)
ViT-S	SupCon	81.4 (0.2)	80.2 (0.3)	85.6 (0.5)	87.8 (1.0)	89.4 (0.7)	94.0 (0.4)

Table 2: Linear evaluation protocol results for contrastive representation learning. Each experiment included three random initial seeds. Mean value and standard deviation (in parentheses) for each metric are reported here. The full table including false negative rates and slide-level metrics can be found in the Appendix E.

7 Limitations, Open Questions, and Ethical Consequences

OpenSRH contains data collected from a single institution. While SRH imagers have standardized settings, different operating room workflows, tumor sampling strategies, and surgeons may produce SRH data distribution shifts. Moreover, while our OpenSRH does contain the most common brain tumor types, rare tumor classes are not included. This is a limitation because rare tumor diagnosis is one of the contexts in which ML-based diagnostic decision support can be most beneficial to clinicians. We intend to include multicenter data with additional tumor rare classes in future releases of OpenSRH.

There are many open questions for the machine learning community that can be explored through OpenSRH. The most important questions are given below:

- **Domain adaptation.** Many domain adaptation literature uses datasets that have very small domain gaps such as MNIST [51], SVHN[52], Office 31[53], or datasets that are artificially crafted or generated, such as Adaptope [54] or DomainNet [55]. OpenSRH can be combined with existing H&E dataset such as TCGA, to create a large-scale benchmark for domain adaptation, with intrinsic pathologic features captured using different imaging modalities.
- **Multiple instance learning.** Besides our patch-based classification workflow, multiple instance learning may be a good strategy for histopathology analysis [56]. By removing patch labels in OpenSRH, histologic classification becomes a generic multiple instance learning task. A model can learn to select the important patches with only slide-level labels. Our current patch-level annotation can be used as a ground truth to interpret instance-level predictions from multiple instance learning paradigms.

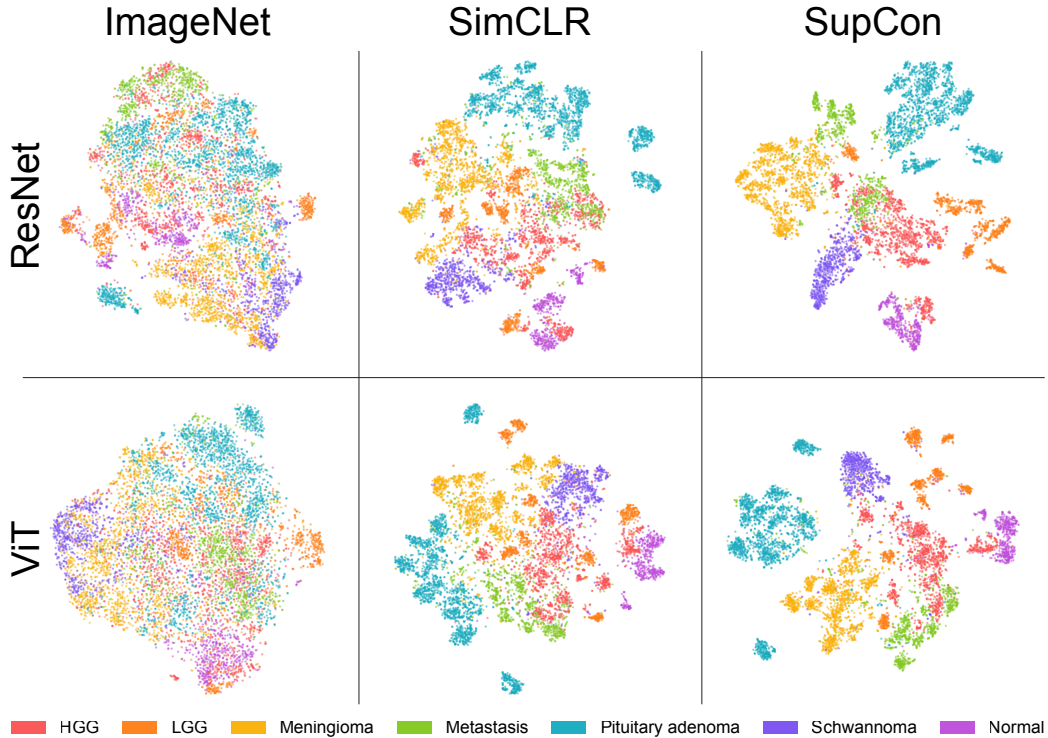


Figure 5: Patch-level SRH representations of validation images. Latent feature vectors are colored by tumor class labels. ImageNet pretraining fails to represent discriminative SRH image features. SimCLR shows discernable SRH feature learning capabilities, while improved class separation can be learned with SupCon. Tumor classes with similar SRH histologic features tend to show similar feature representations, such as HGG/LGG and HGG/Metastasis. A single seed is shown here. Figure best viewed in color.

- **Aggregation of patch-level predictions.** We have relied on individual patch-level predictions and average pooling as a general method for whole slide inference. However, this strategy is limited because it does not account for discriminative heterogeneity within whole slide images. Expectation maximization [57], clustering [58, 28], attention [56, 59], and other multiple instance learning strategies [60] have been proposed as learnable aggregation functions, but questions related to scalability, training efficiency, and data domain differences remain open.
- **Self-supervised learning.** Self-supervised learning and contrastive learning methods have been explored using histology images, but the effectiveness of different augmentation strategies has not been studied. Our preliminary experiments indicate that augmentations used for natural image self-supervised representation learning are sub-optimal. It remains an open question whether domain specific augmentation would improve self-supervised learning performance.
- **Data efficient training.** It is known that ViTs require large image datasets on the scale of ImageNet21K or JFT300M, and insufficient pretraining can result in inferior performance [38]. Acquiring these large supervised datasets is currently infeasible in medical imaging. In addition to low inductive bias, ViTs demonstrate better interpretability compared to CNNs, and their clinical adoption can improve reliability and physician’s trust in AI-assisted diagnosis. By demonstrating a performance gap between CNNs and ViTs, our OpenSRH benchmarks encourage research in data efficient training of ViTs suited for medical imaging.

Lastly, we have ensured that all patients consented to release a portion of their tumor for research. There is minimal additional risk for patients because samples are collected from tumors removed as part of the standard patient care, and their personal health information is protected in OpenSRH.

OpenSRH is released to promote translational AI research. The dataset, algorithms and benchmarks discussed in the paper are for research purposes only.

8 Conclusion

In this work, we introduce OpenSRH, an intraoperative brain tumor dataset of SRH, a rapid, label-free, optical imaging method. OpenSRH contains both raw SRH acquisition data and processed high-resolution image patches for model development using diagnostic annotations from expert neuropathologists. OpenSRH is the first and only publicly available dataset of human cancers imaged using optical histology. We benchmark classification performance for histologic brain diagnosis across the most common brain tumor types. We also provide benchmarks for self-supervised and weakly supervised contrastive representation learning. We release the OpenSRH dataset with the intention to promote translational AI research within the field of precision oncology and optimize the surgical management of human cancers.

Acknowledgements, Disclosure of Funding and Competing Interests

We would like to thank Karen Eddy, Lin Wang, Andrea Marshall and Katherine Lee for their support in data collection and processing.

This work was supported by grants NIH R01CA226527, NIH/NIGMS T32GM141746, NIH K12 NS080223, Cook Family Brain Tumor Research Fund, Mark Trauner Brain Research Fund: Zenkel Family Foundation, and Ian’s Friends Foundation.

Research reported in this publication was also supported by the Investigators Awards grant program of Precision Health at the University of Michigan.

This research was also supported in part through computational resources and services provided by Advanced Research Computing (ARC), a division of Information and Technology Services (ITS) at the University of Michigan, Ann Arbor.

Competing interests: C.W.F. is an employee and shareholder of Invenio Imaging, Inc., a company developing SRH microscopy systems. D.A.O. is an advisor and shareholder of Invenio Imaging, Inc, and T.C.H. is a shareholder of Invenio Imaging, Inc.

References

- [1] David N Louis, Arie Perry, Pieter Wesseling, Daniel J Brat, Ian A Cree, Dominique Figarella-Branger, Cynthia Hawkins, H K Ng, Stefan M Pfister, Guido Reifenberger, Riccardo Soffietti, Andreas von Deimling, and David W Ellison. The 2021 WHO classification of tumors of the central nervous system: a summary. *Neuro. Oncol.*, June 2021.
- [2] Brian J Scott, Vanja C Douglas, Tarik Tihan, James L Rubenstein, and S Andrew Josephson. A systematic approach to the diagnosis of suspected central nervous system lymphoma. *JAMA Neurol.*, 70(3):311–319, March 2013.
- [3] Long Di, Ashish H Shah, Anil Mahavadi, Daniel G Eichberg, Raghuram Reddy, Alexander D Sanjurjo, Alexis A Morell, Victor M Lu, Leonel Ampie, Evan M Luther, Ricardo J Komotar, and Michael E Ivan. Radical supramaximal resection for newly diagnosed left-sided eloquent glioblastoma: safety and improved survival over gross-total resection. *J. Neurosurg.*, pages 1–8, May 2022.
- [4] L Eisenhardt and H Cushing. Diagnosis of intracranial tumors by supravital technique. *Am. J. Pathol.*, 6(5):541–552.7, September 1930.
- [5] Daniel A Orringer, Balaji Pandian, Yashar S Niknafs, Todd C Hollon, Julianne Boyle, Spencer Lewis, Mia Garrard, Shawn L Hervey-Jumper, Hugh J L Garton, Cormac O Maher, Jason A Heth, Oren Sagher, D Andrew Wilkinson, Matija Snuderl, Sriram Venneti, Shakti H Ramkissoon, Kathryn A McFadden, Amanda Fisher-Hubbard, Andrew P Lieberman, Timothy D Johnson, X Sunney Xie, Jay K Trautman, Christian W Freudiger, and Sandra Camelo-Piragua. Rapid intraoperative histology of unprocessed surgical specimens via fibre-laser-based stimulated raman scattering microscopy. *Nat Biomed Eng*, 1, February 2017.
- [6] Todd C Hollon, Balaji Pandian, Esteban Urias, Akshay V Save, Arjun R Adapa, Sudharsan Srinivasan, Neil K Jairath, Zia Farooq, Tamara Marie, Wajd N Al-Holou, Karen Eddy, Jason A Heth, Siri Sahib S

- Khalsa, Kyle Conway, Oren Sagher, Jeffrey N Bruce, Peter Canoll, Christian W Freudiger, Sandra Camelo-Piragua, Honglak Lee, and Daniel A Orringer. Rapid, label-free detection of diffuse glioma recurrence using intraoperative stimulated raman histology and deep neural networks. *Neuro. Oncol.*, July 2020.
- [7] Richard Sullivan, Olusegun Isaac Alatise, Benjamin O Anderson, Riccardo Audisio, Philippe Autier, Ajay Aggarwal, Charles Balch, Murray F Brennan, Anna Dare, Anil D’Cruz, Alexander M M Eggermont, Kenneth Fleming, Serigne Magueye Gueye, Lars Hagander, Cristian A Herrera, Hampus Holmer, André M Ilbawi, Anton Jarnheimer, Jia-Fu Ji, T Peter Kingham, Jonathan Liberman, Andrew J M Leather, John G Meara, Swagoto Mukhopadhyay, Shilpa S Murthy, Sherif Omar, Groesbeck P Parham, C S Pramesh, Robert Riviello, Danielle Rodin, Luiz Santini, Shailesh V Shrikhande, Mark Shrime, Robert Thomas, Audrey T Tsunoda, Cornelis van de Velde, Umberto Veronesi, Dehannathparambil Kottarathil Vijaykumar, David Watters, Shan Wang, Yi-Long Wu, Moez Zeiton, and Arnie Purushotham. Global cancer surgery: delivering safe, affordable, and timely cancer surgery. *Lancet Oncol.*, 16(11):1193–1224, September 2015.
- [8] Phaik-Leng Cheah, Lai Meng Looi, and Susan Horton. Cost analysis of operating an anatomic pathology laboratory in a Middle-Income country. *Am. J. Clin. Pathol.*, 149(1):1–7, January 2018.
- [9] Hilary Lynch Somerset and Bette Kay Kleinschmidt-DeMasters. Approach to the intraoperative consultation for neurosurgical specimens. *Adv. Anat. Pathol.*, 18(6):446–449, November 2011.
- [10] David M Metter, Terence J Colgan, Stanley T Leung, Charles F Timmons, and Jason Y Park. Trends in the US and canadian pathologist workforces from 2007 to 2017. *JAMA Netw Open*, 2(5):e194337, May 2019.
- [11] Stanley J Robboy, Sally Weintraub, Andrew E Horvath, Bradden W Jensen, C Bruce Alexander, Edward P Fody, James M Crawford, Jimmy R Clark, Julie Cantor-Weinberg, Megha G Joshi, Michael B Cohen, Michael B Prystowsky, Sarah M Bean, Saurabh Gupta, Suzanne Z Powell, V O Speights, Jr, David J Gross, and W Stephen Black-Schaffer. Pathologist workforce in the united states: I. development of a predictive model to examine factors influencing supply. *Arch. Pathol. Lab. Med.*, 137(12):1723–1732, December 2013.
- [12] Todd C Hollon, Balaji Pandian, Arjun R Adapa, Esteban Urias, Akshay V Save, Siri Sahib S Khalsa, Daniel G Eichberg, Randy S D’Amico, Zia U Farooq, Spencer Lewis, Petros D Petridis, Tamara Marie, Ashish H Shah, Hugh J L Garton, Cormac O Maher, Jason A Heth, Erin L McKean, Stephen E Sullivan, Shawn L Hervey-Jumper, Parag G Patil, B Gregory Thompson, Oren Sagher, Guy M McKhann, 2nd, Ricardo J Komotar, Michael E Ivan, Matija Snuderl, Marc L Otten, Timothy D Johnson, Michael B Sisti, Jeffrey N Bruce, Karin M Muraszko, Jay Trautman, Christian W Freudiger, Peter Canoll, Honglak Lee, Sandra Camelo-Piragua, and Daniel A Orringer. Near real-time intraoperative brain tumor diagnosis using stimulated raman histology and deep neural networks. *Nat. Med.*, January 2020.
- [13] Cheng Jiang, Abhishek Bhattacharya, Joseph R Linzey, Rushikesh S Joshi, Sung Jik Cha, Sudharsan Srinivasan, Daniel Alber, Akhil Kondepudi, Esteban Urias, Balaji Pandian, Wajid N Al-Holou, Stephen E Sullivan, B Gregory Thompson, Jason A Heth, Christian W Freudiger, Siri Sahib S Khalsa, Donato R Pacione, John G Golfinos, Sandra Camelo-Piragua, Daniel A Orringer, Honglak Lee, and Todd C Hollon. Rapid automated analysis of skull base tumor specimens using intraoperative optical imaging and artificial intelligence. *Neurosurgery*, March 2022.
- [14] Todd Hollon, Spencer Lewis, Christian W Freudiger, X Sunney Xie, and Daniel A Orringer. Improving the accuracy of brain tumor surgery via raman-based technology. *Neurosurg. Focus*, 40(3):E9, March 2016.
- [15] Christian W Freudiger, Wei Min, Brian G Saar, Sijia Lu, Gary R Holtom, Chengwei He, Jason C Tsai, Jing X Kang, and X Sunney Xie. Label-free biomedical imaging with high sensitivity by stimulated raman scattering microscopy. *Science*, 322(5909):1857–1861, December 2008.
- [16] Christian W Freudiger, Wenlong Yang, Gary R Holtom, Nasser Peyghambarian, X Sunney Xie, and Khanh Q Kieu. Stimulated raman scattering microscopy with a robust fibre laser source. *Nat. Photonics*, 8(2):153–159, February 2014.
- [17] Long Di, Daniel G Eichberg, Kevin Huang, Ashish H Shah, Aria M Jamshidi, Evan M Luther, Victor M Lu, Ricardo J Komotar, Michael E Ivan, and Sakir H Gultekin. Stimulated raman histology for rapid intraoperative diagnosis of gliomas. *World Neurosurg.*, 150:e135–e143, June 2021.
- [18] Navid Farahani, Anil V Parwani, and Liron Pantanowitz. Whole slide imaging in pathology: advantages, limitations, and emerging perspectives. *Pathol. Lab. Med. Int.*, 7:23–33, June 2015.
- [19] Minbiao Ji, Daniel A Orringer, Christian W Freudiger, Shakti Ramkissoon, Xiaohui Liu, Darryl Lau, Alexandra J Golby, Isaiah Norton, Marika Hayashi, Nathalie Y R Agar, Geoffrey S Young, Cathie Spino, Sandro Santagata, Sandra Camelo-Piragua, Keith L Ligon, Oren Sagher, and X Sunney Xie. Rapid, label-free detection of brain tumors with stimulated raman scattering microscopy. *Sci. Transl. Med.*, 5(201):201ra119, September 2013.
- [20] Minbiao Ji, Spencer Lewis, Sandra Camelo-Piragua, Shakti H Ramkissoon, Matija Snuderl, Sriram Venneti, Amanda Fisher-Hubbard, Mia Garrard, Dan Fu, Anthony C Wang, Jason A Heth, Cormac O Maher, Nader Sanai, Timothy D Johnson, Christian W Freudiger, Oren Sagher, Xiaoliang Sunney Xie, and Daniel A

- Orringer. Detection of human brain tumor infiltration with quantitative stimulated raman scattering microscopy. *Sci. Transl. Med.*, 7(309):309ra163, October 2015.
- [21] Todd C Hollon, Spencer Lewis, Balaji Pandian, Yashar S Niknafs, Mia R Garrard, Hugh Garton, Cormac O Maher, Kathryn McFadden, Matija Snuderl, Andrew P Lieberman, Karin Muraszko, Sandra Camelo-Piragua, and Daniel A Orringer. Rapid intraoperative diagnosis of pediatric brain tumors using stimulated raman histology. *Cancer Res.*, 78(1):278–289, January 2018.
- [22] Nicolas Coudray, Paolo Santiago Ocampo, Theodore Sakellaropoulos, Navneet Narula, Matija Snuderl, David Fenyo, Andre L Moreira, Narges Razavian, and Aristotelis Tsirigos. Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning. *Nat. Med.*, 24(10):1559–1567, October 2018.
- [23] Jakob Nikolas Kather, Lara R Heij, Heike I Grabsch, Chiara Loeffler, Amelie Echle, Hannah Sophie Muti, Jeremias Krause, Jan M Niehues, Kai A J Sommer, Peter Bankhead, Loes F S Kooreman, Jeffrey J Schulte, Nicole A Cipriani, Roman D Buelow, Peter Boor, Nadi-Na Ortiz-Brüchle, Andrew M Hanby, Valerie Speirs, Sara Kochanny, Akash Patnaik, Andrew Srisuwananukorn, Hermann Brenner, Michael Hoffmeister, Piet A van den Brandt, Dirk Jäger, Christian Trautwein, Alexander T Pearson, and Tom Luedde. Pan-cancer image-based detection of clinically actionable genetic alterations. *Nat Cancer*, 1(8):789–799, August 2020.
- [24] Kunal Nagpal, Davis Foote, Yun Liu, Po-Hsuan Cameron Chen, Ellery Wulczyn, Fraser Tan, Niels Olson, Jenny L Smith, Arash Mohtashamian, James H Wren, Greg S Corrado, Robert MacDonald, Lily H Peng, Mahul B Amin, Andrew J Evans, Ankur R Sangoi, Craig H Mermel, Jason D Hipp, and Martin C Stumpe. Development and validation of a deep learning algorithm for improving gleason scoring of prostate cancer. *NPJ Digit Med*, 2:48, June 2019.
- [25] Po-Hsuan Cameron Chen, Krishna Gadepalli, Robert MacDonald, Yun Liu, Shiro Kadowaki, Kunal Nagpal, Timo Kohlberger, Jeffrey Dean, Greg S Corrado, Jason D Hipp, Craig H Mermel, and Martin C Stumpe. An augmented reality microscope with real-time artificial intelligence integration for cancer diagnosis. *Nat. Med.*, 25(9):1453–1457, September 2019.
- [26] Pooya Mobadersany, Safoora Yousefi, Mohamed Amgad, David A Gutman, Jill S Barnholtz-Sloan, José E Velázquez Vega, Daniel J Brat, and Lee A D Cooper. Predicting cancer outcomes from histology and genomics using convolutional networks. *Proc. Natl. Acad. Sci. U. S. A.*, 115(13):E2970–E2979, March 2018.
- [27] Yu Fu, Alexander W Jung, Ramon Viñas Torne, Santiago Gonzalez, Harald Vöhringer, Artem Shmatko, Lucy R Yates, Mercedes Jimenez-Linan, Luiza Moore, and Moritz Gerstung. Pan-cancer computational histopathology reveals mutations, tumor composition and prognosis. *Nat Cancer*, 1(8):800–810, August 2020.
- [28] Ming Y Lu, Tiffany Y Chen, Drew F K Williamson, Melissa Zhao, Maha Shady, Jana Lipkova, and Faisal Mahmood. AI-based pathology predicts origins for cancers of unknown primary. *Nature*, 594(7861):106–110, June 2021.
- [29] Babak Ehteshami Bejnordi, Mitko Veta, Paul Johannes van Diest, Bram van Ginneken, Nico Karssemeijer, Geert Litjens, Jeroen A W M van der Laak, the CAMELYON16 Consortium, Meyke Hermsen, Quirine F Manson, Maschenka Balkenhol, Oscar Geessink, Nikolaos Stathonikos, Marcory Crf van Dijk, Peter Bult, Francisco Beca, Andrew H Beck, Dayong Wang, Aditya Khosla, Rishab Gargeya, Humayun Irshad, Aoxiao Zhong, Qi Dou, Quanzheng Li, Hao Chen, Huang-Jing Lin, Pheng-Ann Heng, Christian Haß, Elia Bruni, Quincy Wong, Ugur Halici, Mustafa Ümit Öner, Rengul Cetin-Atalay, Matt Berseth, Vitali Khvatkov, Alexei Vylegzhanin, Oren Kraus, Muhammad Shaban, Nasir Rajpoot, Ruqayya Awan, Korsuk Sirinukunwattana, Talha Qaiser, Yee-Wah Tsang, David Tellez, Jonas Annuschein, Peter Hufnagl, Mira Valkonen, Kimmo Kartasalo, Leena Latonen, Pekka Ruusuvuori, Kaisa Liimatainen, Shadi Albarqouni, Bharti Mungal, Ami George, Stefanie Demirci, Nassir Navab, Seiryu Watanabe, Shigeto Seno, Yoichi Takenaka, Hideo Matsuda, Hady Ahmady Phoulady, Vassili Kovalev, Alexander Kalinovsky, Vitali Liauchuk, Gloria Bueno, M Milagro Fernandez-Carrobles, Ismael Serrano, Oscar Deniz, Daniel Racoceanu, and Rui Venâncio. Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *JAMA*, 318(22):2199–2210, December 2017.
- [30] Conall W R Fitzgerald, Snjezana Dogan, Rabih Bou-Nassif, Tim Mclean, Robbie Woods, Jennifer R Cracchiolo, Ian Ganly, Viviane Tabar, and Marc A Cohen. Stimulated raman histology for rapid Intra-Operative diagnosis of sinonasal and skull base tumors. *Laryngoscope*, May 2022.
- [31] Rebecca C Hoesli, Daniel A Orringer, Jonathan B McHugh, and Matthew E Spector. Coherent raman scattering microscopy for evaluation of head and neck carcinoma. *Otolaryngol. Head Neck Surg.*, 157(3):448–453, September 2017.
- [32] Melike Pekmezci, Ramin A Morshed, Pranathi Chunduru, Balaji Pandian, Jacob Young, Javier E Villanueva-Meyer, Tarik Tihan, Emily A Sloan, Manish K Aghi, Annette M Molinaro, Mitchel S Berger, and Shawn L Hervey-Jumper. Detection of glioma infiltration at the tumor margin using quantitative stimulated raman scattering histology. *Sci. Rep.*, 11(1):1–11, June 2021.

- [33] Jenna Wiens, Suchi Saria, Mark Sendak, Marzyeh Ghassemi, Vincent X Liu, Finale Doshi-Velez, Kenneth Jung, Katherine Heller, David Kale, Mohammed Saeed, Pilar N Ossorio, Sonoo Thadaney-Israni, and Anna Goldenberg. Do no harm: a roadmap for responsible machine learning for health care. *Nat. Med.*, 25(9):1337–1340, September 2019.
- [34] Quinn T Ostrom, Nirav Patil, Gino Cioffi, Kristin Waite, Carol Kruchko, and Jill S Barnholtz-Sloan. CBTRUS statistical report: Primary brain and other central nervous system tumors diagnosed in the united states in 2013-2017. *Neuro. Oncol.*, 22(12 Suppl 2):iv1–iv96, October 2020.
- [35] J Deng, W Dong, R Socher, L Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, June 2009.
- [36] Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. Transfusion: Understanding transfer learning for medical imaging. February 2019.
- [37] K He, X Zhang, S Ren, and J Sun. Deep residual learning for image recognition. *Proc. IAPR Int. Conf. Pattern Recog.*, 2016.
- [38] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. September 2020.
- [39] Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas Beyer. How to train your ViT? data, augmentation, and regularization in vision transformers. June 2021.
- [40] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. November 2017.
- [41] Ross Wightman. Pytorch image models, 2019.
- [42] Cheng Jiang, Abhishek Bhattacharya, Joseph Linzey, Rushikesh Joshi, Sung Jik Cha, Sudharsan Srinivasan, Daniel Alber, Akhil Kondepudi, Esteban Urias, Balaji Pandian, Wajid Al-Holou, Steve Sullivan, B Gregory Thompson, Jason Heth, Chris Freudiger, Siri Khalsa, Donato Pacione, John G Golfinos, Sandra Camelo-Piragua, Daniel A Orringer, Honglak Lee, and Todd Hollon. Contrastive representation learning for rapid intraoperative diagnosis of skull base tumors imaged using stimulated raman histology. August 2021.
- [43] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. February 2020.
- [44] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. April 2020.
- [45] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. November 2016.
- [46] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. July 2018.
- [47] Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer. Revisiting self-supervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1920–1929. openaccess.thecvf.com, 2019.
- [48] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *J. Mach. Learn. Res.*, 9(86):2579–2605, 2008.
- [49] Cheng Jiang, Abhishek Bhattacharya, Joseph Linzey, Rushikesh Joshi, Sung Jik Cha, Sudharsan Srinivasan, Daniel Alber, Akhil Kondepudi, Esteban Urias, Balaji Pandian, Wajid Al-Holou, Steve Sullivan, B Gregory Thompson, Jason Heth, Chris Freudiger, Siri Khalsa, Donato Pacione, John G Golfinos, Sandra Camelo-Piragua, Daniel A Orringer, Honglak Lee, and Todd Hollon. Contrastive representation learning for rapid intraoperative diagnosis of skull base tumors imaged using stimulated raman histology. August 2021.
- [50] A S Razavian, H Azizpour, J Sullivan, and S Carlsson. CNN features Off-the-Shelf: An astounding baseline for recognition. In *2014 IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 512–519, June 2014.
- [51] Y Lecun, L Bottou, Y Bengio, and P Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86(11):2278–2324, November 1998.
- [52] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [53] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *Computer Vision – ECCV 2010*, pages 213–226. Springer Berlin Heidelberg, 2010.

- [54] Ringwald and Stiefelhagen. Adaptope: A modern benchmark for unsupervised domain adaptation. *Proceedings of the IEEE/CVF*, 2021.
- [55] Can Peng, Kun Zhao, Arnold Wiliem, Teng Zhang, Peter Hobson, Anthony Jennings, and Brian C Lovell. To what extent does downsampling, compression, and data scarcity impact renal image analysis? In *2019 Digital Image Computing: Techniques and Applications (DICTA)*. IEEE, December 2019.
- [56] Maximilian Ilse, Jakub Tomczak, and Max Welling. Attention-based deep multiple instance learning. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2127–2136. PMLR, 2018.
- [57] L Hou, D Samaras, T M Kurc, Y Gao, J E Davis, and J H Saltz. Patch-Based convolutional neural network for whole slide tissue image classification. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2424–2433, June 2016.
- [58] Ming Y Lu, Drew F K Williamson, Tiffany Y Chen, Richard J Chen, Matteo Barbieri, and Faisal Mahmood. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nat Biomed Eng*, 5(6):555–570, June 2021.
- [59] Zhuchen Shao, Hao Bian, Yang Chen, Yifeng Wang, Jian Zhang, Xiangyang Ji, and Yongbing Zhang. TransMIL: Transformer based correlated multiple instance learning for whole slide image classification. May 2021.
- [60] Gabriele Campanella, Matthew G Hanna, Luke Geneslaw, Allen Mirafior, Vitor Werneck Krauss Silva, Klaus J Busam, Edi Brogi, Victor E Reuter, David S Klimstra, and Thomas J Fuchs. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nat. Med.*, 25(8):1301–1309, August 2019.
- [61] Matěj Týč and Christoph Gohlke. imreg_dft. https://github.com/matejak/imreg_dft, 2016.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] The limitations of OpenSRH are described in section 7.
 - (c) Did you discuss any potential negative societal impacts of your work? [Yes] Potential negative societal impacts are described in section 7.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - (b) Did you include complete proofs of all theoretical results? [N/A]
3. If you ran experiments (e.g. for benchmarks)...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] All code, data, and instructions are available on <https://opensrh.mlins.org>.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] Detailed training protocol is listed in appendix C, and data split information is described in section 4 and appendix A.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] All experiments are repeated with 3 random seeds and error bars are reported in all tables.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] Details of the compute resources and time to produce the OpenSRH benchmarks are described in appendix C.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes] Our benchmark implementation uses existing open source framework. More details are described in appendices A and C.
 - (b) Did you mention the license of the assets? [Yes] All licenses of these frameworks are included in the THIRD_PARTY file in the root level of our repository.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] The OpenSRH dataset and benchmark source code can be accessed on <https://opensrh.mlins.org>.
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [Yes] Informed consent was obtained for each patient in OpenSRH. Details are described in section 4.
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [Yes] All personally identifiable information are removed before data release and benchmark training. Details are described in section 4.
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [Yes] This study was approved by Institutional Review Board (HUM00083059), and more details are described in section 4. Potential participant risk is minimal, and described in section 7.
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

A Data Release and Source Code Technical Details

Preprocessing As described in Section 4, SRH images are acquired in two Raman shift frequencies. They each correspond to a registered channel in an RGB image (green and blue for 2845 cm^{-1} and 2930 cm^{-1} , respectively). The co-registration utilizes discrete Fourier transform implemented using `imreg_dft` python package [61]. The third channel (red) is obtained by subtracting the first two channels, in their original 16-bit depth. The three-channel images are then converted to floats between 0 and 1, and are used as model input. A panel of paired raw and RGB images is shown in Figure 6.

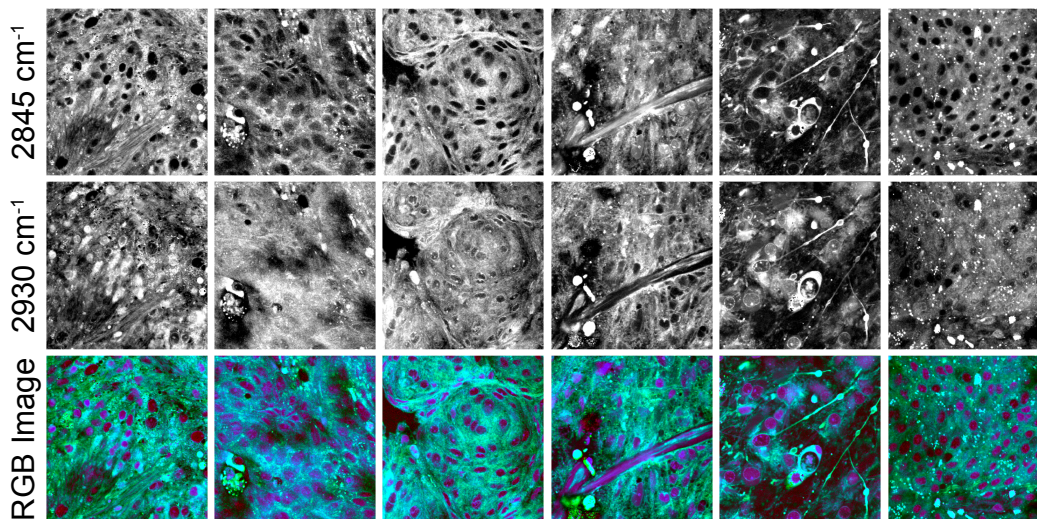


Figure 6: A panel of constructed RGB images and their respective raw images. The RGB images are normalized and contrast adjusted for better viewing.

Segmentation The data included in OpenSRH have been approximately segmented by two steps. They were first divided into non-overlapping 300×300 pixel patches and then each patch will be classified into three categories (tumor, nondiagnostic and normal) using the previously trained model. The model was trained on about 6.65 million patches using data collected from multiple institutions. An ImageNet pretrained ResNet50 model was used to train on first manually labelled 300,000 images by certified pathologists. Then, the predictions on the part of unlabelled data were manually checked and used for fine-tuning the model. This process was applied to the rest of the data iteratively until all data are well labelled. The model’s predictions are included in the metadata for each slide.

Data release Data is available through a multitude of options. Primarily, they will be available via a Google Drive and Amazon AWS S3 upon completion of a short data usage agreement. The link to google drive will be available automatically after completing the survey, without human approval. Please contact the authors if you wish to download the data from AWS in case of regional unavailability through other means. The data available through Google Drive is compressed (~ 364 GB) and can be downloaded directly and uncompressed afterward. A list of checksums is also made available for data integrity checks. Please note, data available through AWS is uncompressed (~ 449 GB) and but it will require you to have an AWS account.

Dataset directory organization and metadata OpenSRH is intended to be assembled into the following directory structure in figure 7. The metadata for OpenSRH is stored in `meta/opensrh.json`. It provides the patient-level ground truth label, relative paths to patches and their segmentation prediction using the format described in figure 8. The `meta/train_val_split.json` file contains the default random split between the training and validation set. It’s format is described in figure 9.

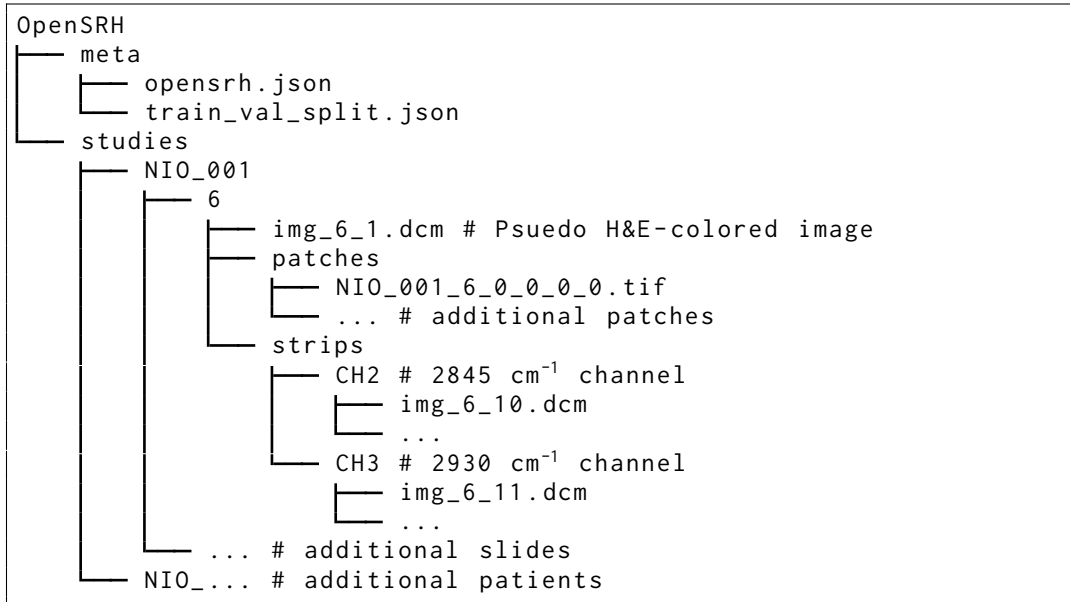


Figure 7: Intended tree directory structure for the OpenSRH dataset. All the associated data for training has been included.

```

{
  "NIO_001": {
    "patient_id": "NIO_001",
    "class": "hgg",
    "slides": {
      "6": {
        "slide_id": "6",
        "tumor_patches": [
          "NIO_001/6/patches/NIO_001-6-0_0_0_0.tif",
          ...
        ],
        "normal_patches": [...],
        "nondiagnostic_patches": [...],
      }, ...
    }
  }, ...
}
  
```

Figure 8: Metadata file format. A patient may have several slides. The number of patches in a slide is variable. The last four numbers included in the patch file name indicate their location in the whole slide image.

```

{
  "train": ["NIO_001", "NIO_002", "NIO_005", "NIO_006", ...],
  "val": ["NIO_003", "NIO_004", "NIO_007", "NIO_009", ...]
}
  
```

Figure 9: Training validation split metadata file. It consists of a dictionary with 2 lists of strings representing patient IDs.

B Patient Age Distribution

OpenSRH data includes patients ranging from newborn to 87 years old. Different classes of tumors are well-distributed among the age groups shown in figure 10. The distribution is skewed to the left as we expect with an increasing incidence of brain tumors with age. We did not observe a difference in model performance or notice differences in tumor cytologic or histoarchitectural features between age groups.

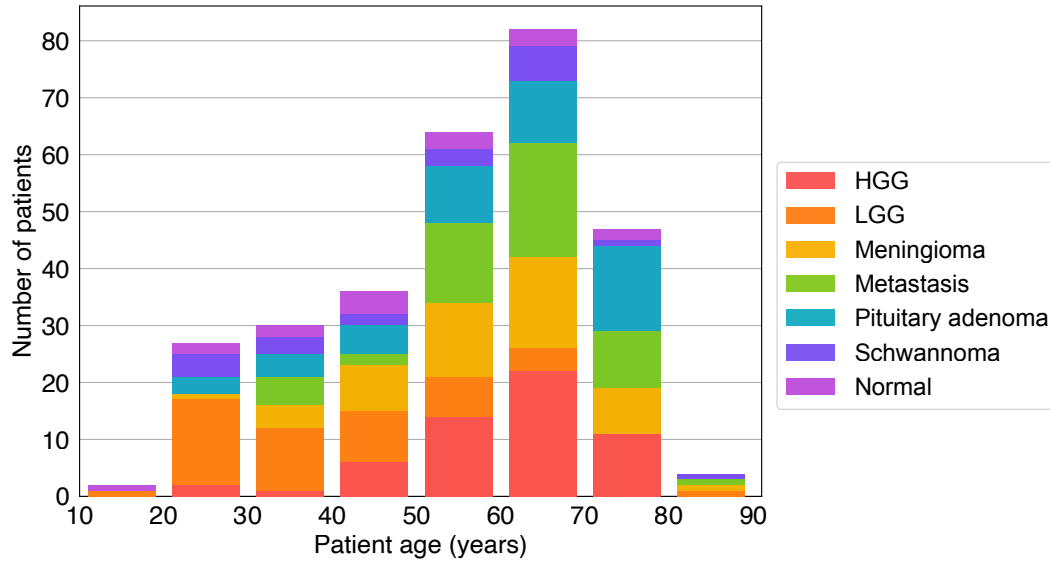


Figure 10: Distribution of patient age in OpenSRH. We can observe that high grade gliomas, metastasis, and meningiomas are more common in older patients, while low grade gliomas are more common in younger patients. HGG, high grade glioma; LGG, low grade glioma.

C Training Protocol and Details

The experiments were trained using PyTorch Lightning, which is a wrapper for PyTorch and allows for efficient multi-GPU training (if available). Associated code and example configuration file to train the benchmark models is available at <https://opensrh.milins.org>.

Training protocol is described in sections 5.1 and 6.1 for our histological classification and contrastive learning benchmarks, respectively. In this section, we provide detailed parameters and more details on the augmentations used for contrastive learning.

C.1 Histological Classification Training Protocols

Table 3 shows detailed training parameters for cross entropy experiments. The augmentations used for these experiments are vertical flipping and horizontal flipping, each with a probability of 0.5. For ViT-S training, we used a cosine learn rate scheduler, with a 0.3 cosine period and the first 10% steps as a linear warmup stage.

Backbone	ResNet50	ViT-S
Classification head	Linear(2048, 7)	Linear(384, 7)
Augmentations	Flipping	Flipping, Resize 224
Batch size	96	256
# GPUs	2 × Nvidia 2080Ti	
Loss	Cross Entropy	
# Epochs	20	
Pretrained	{Random, ImageNet}	
Optimizer	AdamW	
Initial learn rate	1E-3	1E-4
Scheduler	Step, half @ epoch	Cosine w/ warmup
Seeds	{1000, 2000, 3000}	
Time (hrs)	9.5	9
# Parameters	23.5M	21.7M

Table 3: Training parameters of histological classification benchmarks. Training time is an estimate and should be used for reference only.

C.2 Contrastive Learning Training Protocols

Table 4 shows detailed training parameters for contrastive learning experiments. The cosine learn rate scheduler follows the same parameters as described in section C.1.

Backbone	ResNet50	ViT-S
Projection head	Linear(2048, 128)	Linear(384, 24)
Augmentations	Strong	Strong, Resize 224
Batch size	448	512
# GPUs	8 × Nvidia 2080Ti	
Method	{SimCLR, SupCon}	
# Epochs	40	
Optimizer	AdamW	
Initial learn rate	1E-2	5E-4
Scheduler	None	Cosine w/ warmup
Seeds	{1000, 2000, 3000}	
Time (hrs)	15.5	9.8

Table 4: Contrastive learning pre-training parameters. Training time is an estimate and should be used for reference only.

C.2.1 Augmentations

The strong augmentations used in these contrastive learning experiments consist of multiple random augmentations applied sequentially:

- Random Horizontal Flip
- Random Vertical Flip
- Gaussian Noise
- Color Jittering
- Random Autocontrast
- Random Solarize with threshold 0.2
- Random Adjust Sharpness with sharpness factor 2
- Gaussian Blur with kernel size 5 and sigma 1
- Random Erasing
- Random Affine Transformation with max 10 degrees rotation and 10-30% image translation
- Random Resized Crop

All augmentations are applied with a 0.3 probability and use the default PyTorch parameter unless otherwise noted above. In ViT experiments, Resize is applied before all other augmentations. A panel of randomly augmented images is shown in figure 11.

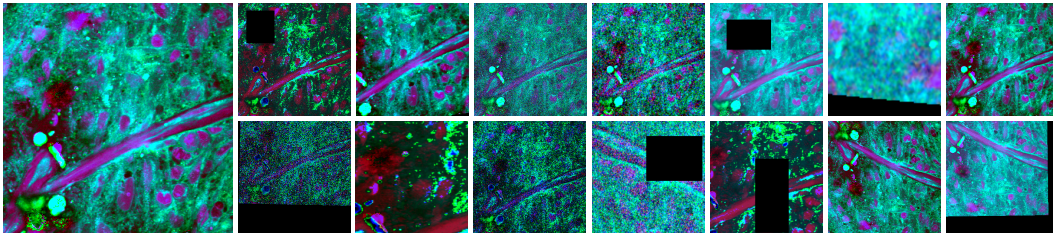


Figure 11: A sample panel of randomly augmented patches. The left patch is the original, and the rest of the patches are augmented using the protocol in section C.2.1.

C.3 Contrastive Learning Linear Evaluation Protocols

Table 5 shows detailed training parameters for contrastive learning linear evaluation. The cosine learn rate scheduler follows the same parameters as described in section C.1.

Backbone	ResNet50	ViT-S
Classification head	Linear(2048, 7)	Linear(384, 7)
Augmentations	Flipping	Flipping, Resize 224
Batch size	96	256
# GPUs	2 × Nivida 2080Ti	
Loss	Cross Entropy	
# Epochs	20	
Initialization	{SimCLR, SupCon, ImageNet}	
Optimizer	AdamW	
Initial learn rate	1E-3	1E-4
Scheduler	Step, half @ epoch	Cosine w/ warmup
Seeds	{1000, 2000, 3000}	
Linear training time (hrs)	20	6
Finetune training time (hrs)	13	5

Table 5: Cross entropy experiment training parameters. Training time is an estimate and should be used for reference only.

D Additional Histological Classification Results

In addition to the results in Table 1, we also computed average precision. Average precision is computed using one-vs-all and averaged over all classes. In addition to patch- and patient-level metrics, we also report the metric computed at the slide level. These results are shown in Table 6:

	Backbone	Pretrain	Accuracy	Top 2	MCA	MAP	FNR
Patch level metrics	ResNet50	Random	84.4 (0.4)	93.5 (0.2)	83.8 (0.5)	89.5 (0.5)	0.9 (0.1)
	ResNet50	ImageNet	86.5 (0.4)	94.4 (0.1)	85.6 (0.3)	91.2 (0.3)	1.0 (0.0)
	ViT-S	Random	77.2 (0.5)	90.0 (0.4)	76.8 (0.8)	82.3 (0.5)	1.8 (0.1)
	ViT-S	ImageNet	83.7 (0.5)	93.4 (0.2)	82.7 (0.9)	88.8 (0.1)	1.9 (0.2)
Slide level metrics	ResNet50	Random	88.7 (0.8)	95.3 (0.3)	88.1 (1.0)	93.6 (0.1)	0.5 (0.5)
	ResNet50	ImageNet	88.8 (0.5)	95.7 (0.3)	88.4 (0.5)	94.4 (0.1)	1.2 (0.0)
	ViT-S	Random	83.7 (0.3)	95.5 (0.4)	83.8 (0.9)	92.0 (0.6)	1.8 (0.2)
	ViT-S	ImageNet	88.7 (0.8)	97.1 (0.2)	88.3 (0.9)	93.9 (0.4)	1.9 (0.0)
Patient level metrics	ResNet50	Random	90.0 (0.0)	95.0 (0.0)	91.4 (0.0)	92.8 (0.2)	0.0 (0.0)
	ResNet50	ImageNet	88.9 (0.8)	94.4 (0.8)	90.5 (0.6)	94.0 (0.1)	0.6 (0.8)
	ViT-S	Random	85.0 (1.4)	95.0 (0.0)	87.2 (1.1)	93.2 (0.4)	1.7 (0.0)
	ViT-S	ImageNet	88.9 (0.8)	96.1 (0.8)	90.5 (0.6)	93.9 (0.4)	1.7 (0.0)

Table 6: Extended metrics for histologic classification benchmarks for ResNet50 and ViT-S. Pretrain refers to the pretraining strategy. Each experiment included three random initial seeds. The mean value and standard deviation (in parentheses) for each metric are reported. MCA, mean class accuracy; MAP, mean average precision; FNR, false negative rate.

In Figure 4, we showed the confusion matrices at the patient level for the first random seed. We compared different initialization strategies (random and ImageNet) and different architectures (ResNet and ViT). The confusion matrices at the patient level for the second and third seeds are shown in Figure 12 below.

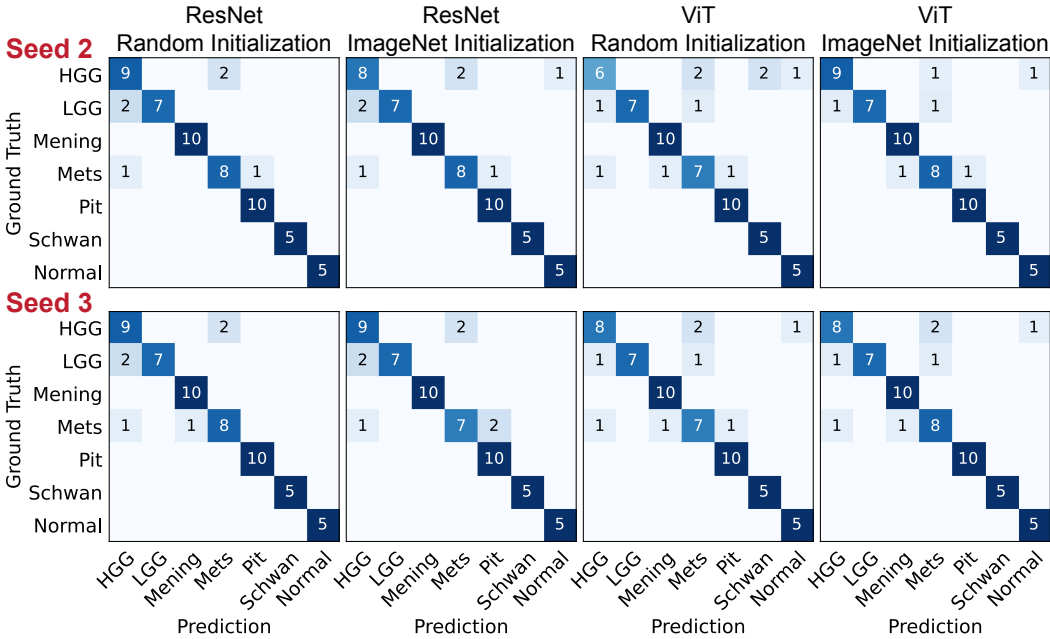


Figure 12: Patient-level confusion matrices for the four different training strategies on the validation set. Seeds 2 and 3 are shown. Mening, meningioma; Mets, metastasis; Pit, pituitary adenoma; Schwann, schwannoma.

We also include the confusion matrix for the first random seed at patch, slide, and patient level in figure 13 below. They correspond to the confusion matrices in figure 4.

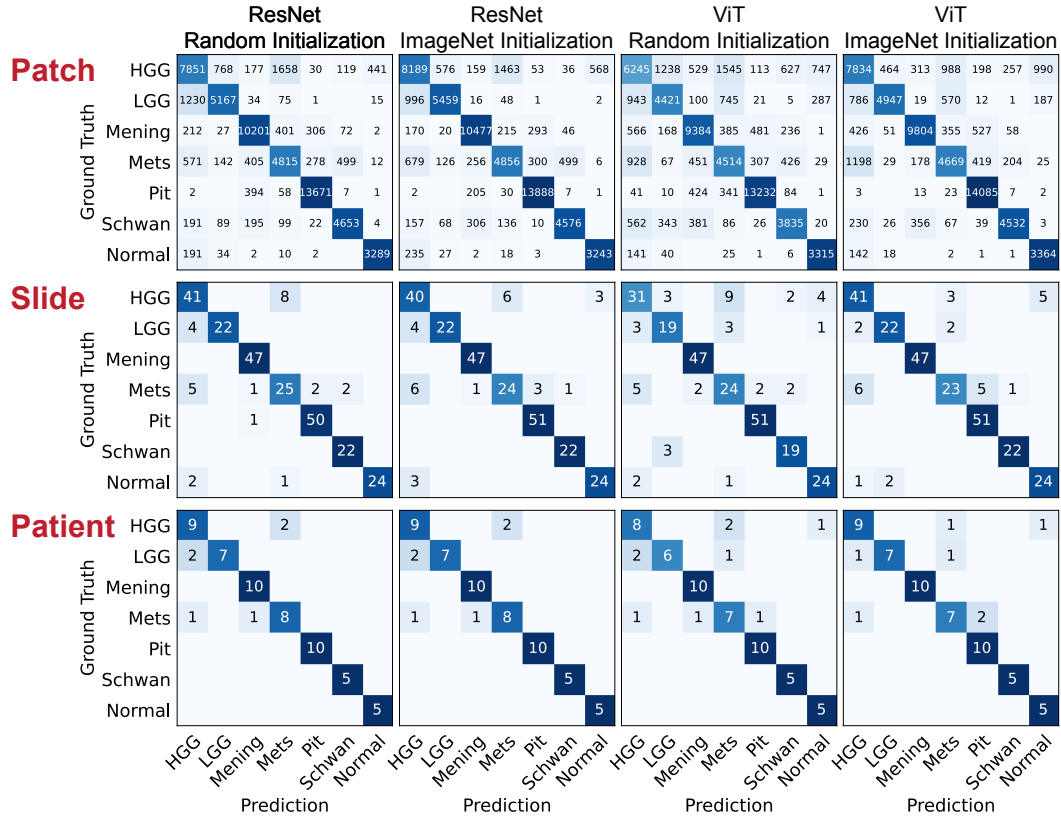


Figure 13: Patch, slide, and patient-level confusion matrices for the four different training strategies on the validation set. Seed 1 is shown. Mening, meningioma; Mets, metastasis; Pit, pituitary adenoma; Schwan, schwannoma.

E Additional Linear Evaluation Results for Contrastive Representation Learning

In table 2, we reported contrastive representation benchmarks evaluated using a linear classifier. Table 7 shows the extended results, including mean average precision and slide-level classification metrics.

	Backbone	Methods	Accuracy	Top 2	MCA	MAP	FNR
Patch level metrics	ResNet50	ImageNet	68.3 (0.0)	84.1 (0.0)	67.9 (0.0)	72.9 (0.1)	1.5 (0.0)
	ResNet50	SimCLR	79.1 (0.4)	92.8 (0.3)	78.9 (0.4)	84.2 (0.6)	1.5 (0.0)
	ResNet50	SupCon	87.5 (0.3)	94.8 (0.2)	86.8 (0.3)	91.5 (0.5)	1.4 (0.2)
	ViT-S	ImageNet	71.8 (0.1)	87.0 (0.0)	71.1 (0.1)	77.1 (0.1)	1.4 (0.0)
	ViT-S	SimCLR	76.8 (0.5)	90.7 (0.2)	76.3 (0.5)	82.5 (0.3)	1.2 (0.2)
	ViT-S	SupCon	81.4 (0.2)	92.2 (0.3)	80.2 (0.3)	85.6 (0.5)	1.7 (0.0)
Slide level metrics	ResNet50	ImageNet	80.9 (0.3)	92.2 (0.0)	81.2 (0.3)	86.1 (0.1)	0.8 (0.0)
	ResNet50	SimCLR	84.4 (1.6)	96.8 (0.4)	84.3 (1.4)	91.9 (0.2)	1.8 (0.2)
	ResNet50	SupCon	91.1 (0.4)	97.0 (0.2)	90.6 (0.4)	95.3 (0.3)	1.4 (0.2)
	ViT-S	ImageNet	89.1 (0.3)	96.8 (0.2)	88.6 (0.4)	92.7 (0.2)	0.3 (0.2)
	ViT-S	SimCLR	83.0 (0.4)	95.7 (0.8)	83.0 (0.7)	90.1 (0.6)	1.0 (0.9)
	ViT-S	SupCon	87.4 (0.2)	96.6 (0.2)	86.5 (0.4)	92.2 (0.4)	1.9 (0.0)
Patient level metrics	ResNet50	ImageNet	80.0 (0.0)	93.3 (0.0)	82.9 (0.0)	88.8 (0.1)	0.0 (0.0)
	ResNet50	SimCLR	83.9 (1.0)	97.2 (1.0)	86.1 (0.9)	92.4 (0.1)	1.7 (0.0)
	ResNet50	SupCon	90.0 (0.0)	95.0 (0.0)	91.4 (0.1)	94.6 (0.5)	1.7 (0.0)
	ViT-S	ImageNet	88.3 (0.0)	95.0 (0.0)	89.8 (0.0)	93.9 (0.0)	0.0 (0.0)
	ViT-S	SimCLR	80.0 (1.7)	96.1 (1.0)	83.0 (1.3)	92.3 (0.0)	1.1 (1.0)
	ViT-S	SupCon	87.8 (1.0)	96.7 (0.0)	89.4 (0.7)	94.0 (0.4)	1.7 (0.0)

Table 7: Extended metrics for linear evaluation protocol results in contrastive representation learning. Each experiment included three random initial seeds. Mean value and standard deviation (in parentheses) for each metric are reported. MCA, mean class accuracy; MAP, mean average precision; FNR, false negative rate.

In addition to the classification metric, we also include the confusion matrix in figure 14 below.

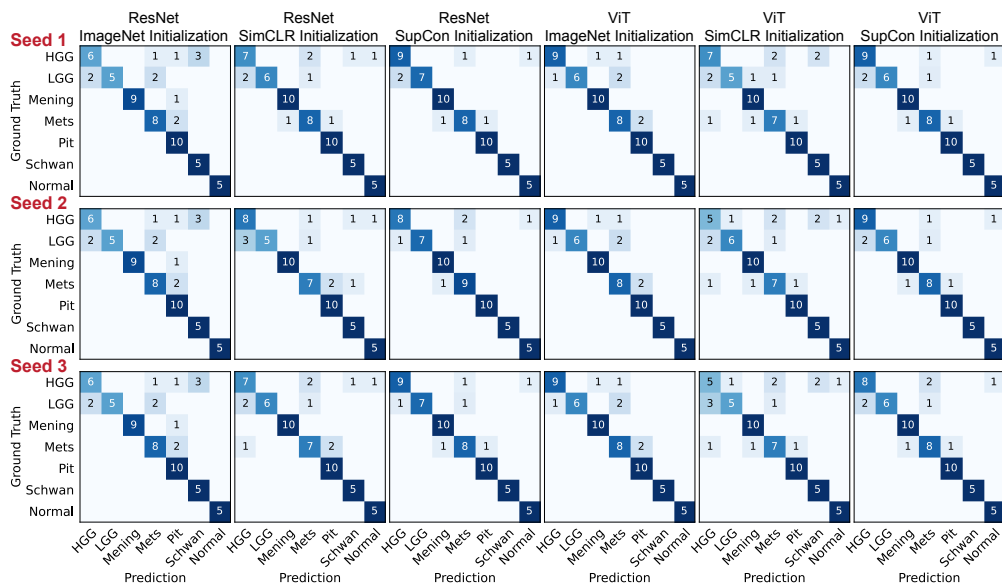


Figure 14: Patient-level confusion matrices of linear evaluations on the validation set. We included initialization from ImageNet, SimCLR, and Supcon, for both ResNet and ViT architecture, and all 3 seeds. Mening, meningioma; Mets, metastasis; Pit, pituitary adenoma; Schwan, schwannoma.

F Contrastive Representation Learning Fine-tuning Evaluation Results

In addition to evaluating our contrastive learning benchmarks using the linear evaluation protocol (section 6.1), we also performed evaluation using fine-tuning. These experiments have the same protocol as the linear evaluation (as described in section 6.1 and C.3), except for trainable (unfrozen) weights in the model backbone. Finetuning metrics are reported in table 8 below:

	Backbone	Methods	Accuracy	Top 2	MCA	MAP	FNR
Patch level metrics	ResNet50	SimCLR	86.3 (0.3)	94.5 (0.2)	85.2 (0.4)	90.6 (0.1)	1.1 (0.1)
	ResNet50	SupCon	87.8 (0.3)	94.8 (0.2)	86.5 (0.4)	91.4 (0.4)	1.1 (0.1)
	ViT-S	SimCLR	81.4 (0.1)	92.3 (0.2)	80.7 (0.4)	86.2 (0.2)	2.0 (0.1)
	ViT-S	SupCon	81.2 (0.4)	92.2 (0.1)	80.0 (0.5)	85.3 (0.6)	2.0 (0.2)
Slide level metrics	ResNet50	SimCLR	89.8 (0.2)	96.4 (1.2)	89.2 (0.2)	94.8 (0.2)	1.0 (0.2)
	ResNet50	SupCon	90.1 (0.4)	96.0 (0.2)	89.5 (0.5)	95.0 (0.3)	0.9 (0.4)
	ViT-S	SimCLR	85.6 (0.4)	97.1 (0.6)	84.9 (0.4)	92.7 (0.8)	2.2 (0.2)
	ViT-S	SupCon	86.8 (0.7)	97.1 (0.6)	86.2 (0.8)	91.7 (1.6)	2.2 (0.2)
Patient level metrics	ResNet50	SimCLR	89.4 (1.0)	95.6 (1.0)	90.9 (0.7)	94.2 (0.3)	1.1 (1.0)
	ResNet50	SupCon	91.7 (0.0)	95.0 (0.0)	92.7 (0.0)	94.9 (0.3)	0.6 (1.0)
	ViT-S	SimCLR	87.2 (1.0)	96.7 (0.0)	88.9 (0.9)	94.1 (0.6)	1.7 (0.0)
	ViT-S	SupCon	86.7 (0.0)	96.7 (0.0)	88.3 (0.1)	94.0 (0.7)	2.8 (1.0)

Table 8: Metrics for finetuning evaluation protocol results for contrastive representation learning. Each experiment included three random initial seeds. Mean value and standard deviation (in parentheses) for each metric are reported. MCA, mean class accuracy; MAP, mean average precision; FNR, false negative rate.

We also include confusion matrices for these experiments in figure 15 below.

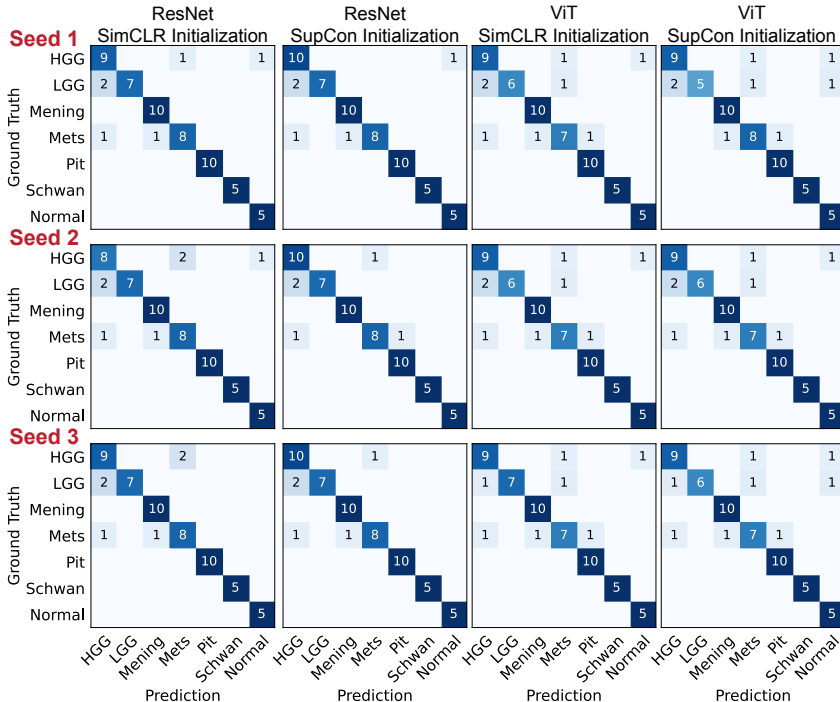


Figure 15: Patient-level confusion matrices of finetuning evaluations on the validation set. We included initialization from SimCLR and Supcon, for both ResNet and ViT architecture, and all 3 seeds. Mening, meningioma; Mets, metastasis; Pit, pituitary adenoma; Schwann, schwannoma. Mening, meningioma; Mets, metastasis; Pit, pituitary adenoma; Schwann, schwannoma.