

TumorCLIP: Lightweight Vision–Language Fusion for Explainable MRI-Based Brain Tumor Classification

Yishaoying Jia^{*1}, Jinfu Niu^{†2}, Zhonghui Qie³, Zongyu Li², Andrew Laine², and Jia Guo^{2,4}

¹Xi’an Jiaotong-Liverpool University, Digital Media Technology

²Columbia University, Biomedical Engineering Department

³Columbia University, Zuckerman Institute

⁴Columbia University, Psychiatry

February 4, 2026

Abstract

Accurate classification of brain tumors from MRI is critical for guiding clinical decision-making; however, existing deep learning models are often hindered by limited interpretability and pronounced sensitivity to hyperparameter selection, which constrain their reliability in medical settings. To address these challenges, we propose TumorCLIP, a lightweight and training-efficient vision-language framework that integrates radiology-informed text prototypes with a DenseNet-based visual encoder to support clinically meaningful semantic reasoning, fused via a Tip-Adapter mechanism. TumorCLIP does not aim to introduce a new vision-language model architecture. Instead, its contribution lies in the integration of radiology-informed text prototypes tailored to MRI interpretation, a systematic evaluation of backbone stability across diverse visual architectures, and a lightweight, training-efficient CLIP-based fusion framework designed for medical imaging applications. We first conduct a comprehensive unimodal benchmark across eight representative visual backbones (EfficientNet-B0, MobileNetV3-Large, ResNet50, DenseNet121, ViT, DeiT, Swin Transformer, and MambaOut) using a standardized optimizer and learning-rate grid search, revealing performance swings exceeding 60 percentage points depending on hyperparameter choices. DenseNet121 shows the strongest stability-accuracy trade-off within our evaluated optimizer and learning-rate grid (97.6%). Leveraging this foundation, TumorCLIP fuses image features with frozen CLIP-derived text prototypes, achieving concept-level explainability, robust few-shot adaptation, and enhanced classification of minority tumor classes. On the test set, TumorCLIP attains 98.5% accuracy, including a +1.86 percentage point recall increase for Neurocytoma, suggesting that radiology-informed textual priors can improve semantic alignment and help refine diagnostic decision boundaries within the evaluated setting. Additional evaluation on an independent external dataset shows that TumorCLIP achieves improved cross-dataset performance under the evaluated distribution shift, relative to the unimodal DenseNet121 baseline. These results demonstrate TumorCLIP as a practical, interpretable, and data-efficient alternative to conventional visual classifiers, providing evidence for radiology-aware vision-language alignment in MRI-based brain tumor classification. All results are reported within the evaluated datasets and training protocols.

Introduction

Magnetic resonance imaging (MRI) is indispensable for diagnosing and managing brain tumors due to its superior soft-tissue contrast and non-ionizing nature. While deep learning approaches have achieved promising results in automating tumor classification from MRI [1], clinical adoption is hampered by limited interpretability, sensitivity to hyperparameters, and the “black box” nature of vision-only architectures [2]. Such

^{*}Y. Jia and J. Niu contributed equally to this work.

[†]Y. Jia and J. Niu contributed equally to this work.

36 limitations undermine trust, reproducibility, and the safe application of AI in cases involving rare or visually
37 subtle tumor phenotypes.

38 Recent breakthroughs in vision-language models, notably CLIP [3], have enabled powerful zero-shot learning
39 and semantic alignment by jointly representing images and text. Although these models offer significant
40 potential for explainable decision-making, their adoption in medical imaging and neuroimaging remains
41 minimal. This is primarily due to a lack of large-scale paired image-text datasets and a semantic gap
42 between natural language and the specialized descriptors used in radiology.

43 Existing multimodal methods also face significant barriers: paired radiology reports are scarce, radiology-
44 specific language diverges from CLIP’s natural image training data, and end-to-end multimodal networks
45 are often computationally intensive. Furthermore, most prior work does not provide a standardized evalu-
46 ation of visual backbones, leaving the reliability of multimodal fusion architectures uncertain. To address
47 these gaps, we conduct a rigorous, unified, and unimodal benchmark across eight popular visual backbones,
48 including CNNs, Transformers, and state-space models, using a standardized optimizer and learning-rate
49 grid. Our results reveal hyperparameter sensitivity exceeding 60 percentage points, with DenseNet121 [4]
50 demonstrating the highest stability and accuracy, and thus serving as the foundation for our multimodal
51 approach.

52 Building on this foundation, we propose TumorCLIP: a lightweight and interpretable vision-language frame-
53 work that fuses radiology-style text prototypes with DenseNet121 via a Tip-Adapter mechanism. By uti-
54 lizing a frozen CLIP text encoder to generate class-specific radiology prototypes and combining these with
55 image-derived features through a learnable fusion weight, TumorCLIP achieves concept-level interpretability,
56 reduced training cost, and improved detection of rare tumor classes. TumorCLIP obtains 98.0% accuracy,
57 surpassing all unimodal baselines and providing more structured, clinically relevant decision boundaries.
58 Collectively, these results highlight the potential of radiology-aware vision-language alignment for enhancing
59 both accuracy and explainability in MRI-based brain tumor classification.

60 Results

61 Backbone Benchmark

62 To ensure a robust foundation for subsequent multimodal fusion, we first conducted a comprehensive uni-
63 modal benchmark to identify the most stable and high-performing visual encoder. We systematically eval-
64 uated eight representative backbones spanning convolutional networks, transformers, and state-space archi-
65 tectures: EfficientNet-B0 [5], MobileNetV3-Large [6], ResNet50 [7], DenseNet121 [4], ViT [8], DeiT [9], Swin
66 Transformer [10], and MambaOut (Mamba) [11]. Each model was trained using an identical grid of opti-
67 mizers (SGD, Adam) and learning rates (10^{-3} , 10^{-4} , 10^{-5} , and 10^{-6}), ensuring strict comparability across
68 architectural paradigms.

69 Our results revealed pronounced sensitivity to hyperparameter selection across all models (Figure 1; [12]),
70 with accuracy ranges exceeding 60 percentage points for several backbones. For instance, MobileNetV3-
71 Large demonstrated performance spanning from 14.1% to 97.3%, and DenseNet121 ranged from 32.1% to
72 98.6%. Transformer architectures such as ViT and Swin exhibited similar volatility, although they generally
73 surpassed CNNs at lower learning rates.

74 Among all candidates, DenseNet121 showed the strongest reliability within the evaluated optimizer and
75 learning-rate grid, delivering the highest validation accuracy (98.6%) and a test accuracy of 97.6% while
76 maintaining a modest trainable parameter count (14.84M). Its robustness across diverse optimizer and learn-
77 ing rate combinations, coupled with favorable parameter efficiency, established DenseNet121 as the optimal
78 backbone for integration within the TumorCLIP framework.

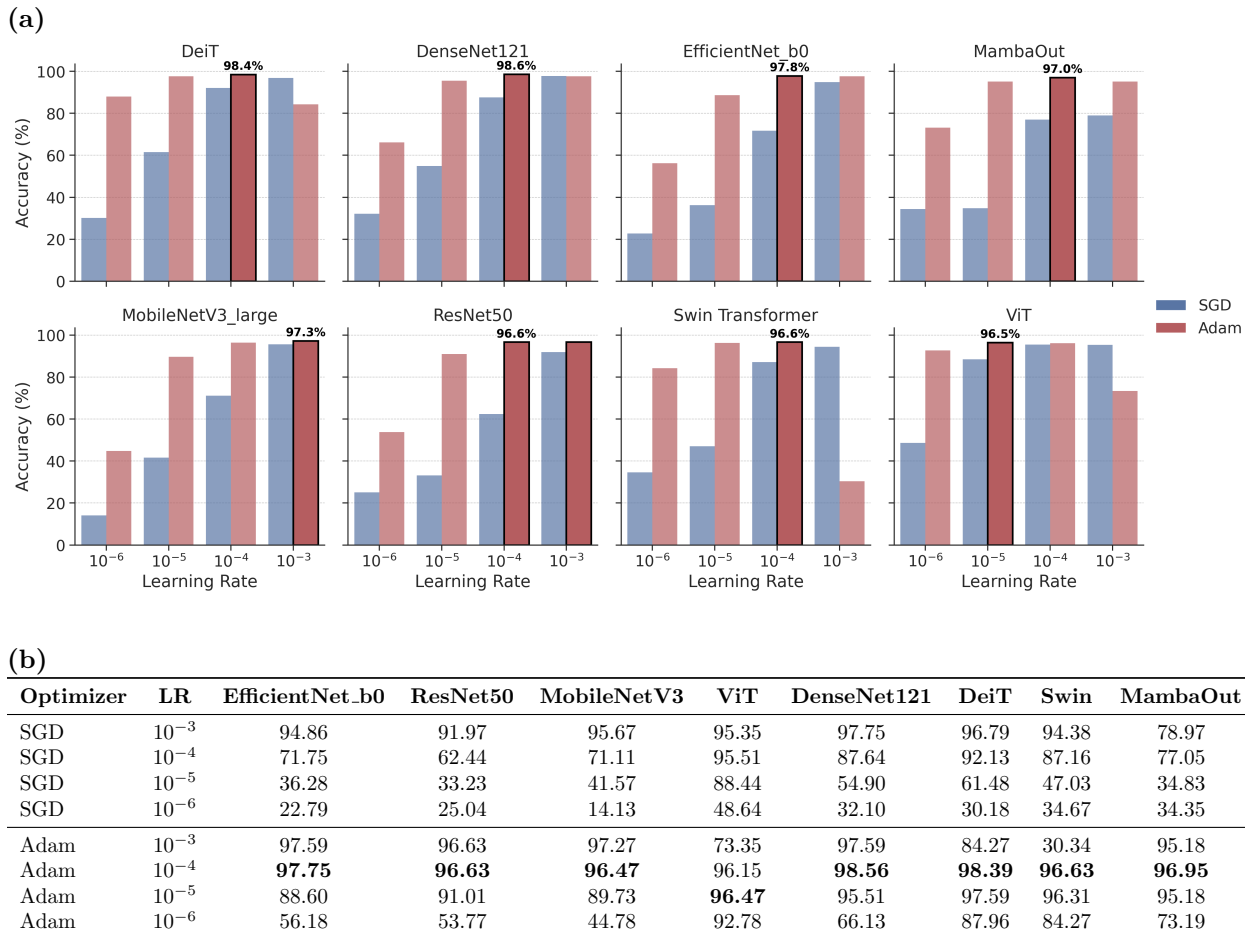


Figure 1: **Validation performance of single-modal models.** All visual backbones were evaluated under the same optimizer and learning-rate grid to ensure fair comparison across architectures. (a) Validation accuracy bar charts for each backbone. (b) Detailed performance table across all optimizers and learning rate configurations. The highest accuracy for each model is highlighted in bold.

79 TumorCLIP Architecture and Multimodal Fusion

80 As illustrated in Figure 2, TumorCLIP consists of two parallel branches: a visual pathway and a text pathway,
 81 which are integrated through a lightweight Tip-Adapter fusion module followed by a weighted ensemble. In
 82 the visual pathway, an input MRI volume is encoded by a fine-tuned DenseNet121 backbone to produce a
 83 fixed-dimensional image embedding, together with classifier logits generated by the backbone’s classification
 84 head. In parallel, the text pathway encodes radiology-style class prompts using a frozen CLIP text encoder.
 85 For each diagnostic category, prompt embeddings are averaged to form a class-level text prototype, which
 86 serves as a semantic reference during inference. To incorporate instance-level visual evidence, TumorCLIP
 87 employs a Tip-Adapter module that leverages a precomputed cache of training-image features to supply
 88 instance-level visual evidence [13]. Given a test image embedding, the Tip-Adapter performs k-nearest-
 89 neighbor retrieval within the cache and aggregates the retrieved samples into class-specific cache scores.
 90 These cache-based scores are combined with similarity scores computed between the test embedding and the
 91 class-level text prototypes through a learnable weighting mechanism, yielding the Tip-Adapter prediction.
 92 Finally, this Tip-Adapter output is combined with the DenseNet classifier logits via a learnable weighted
 93 fusion to obtain the final inference result.

94 Collectively, these components enable TumorCLIP to integrate instance-level visual similarity with class-
95 level radiologic semantics, resulting in a lightweight multimodal classifier that enhances both accuracy and
96 interpretability, while maintaining the CLIP text encoder in a fully frozen state.

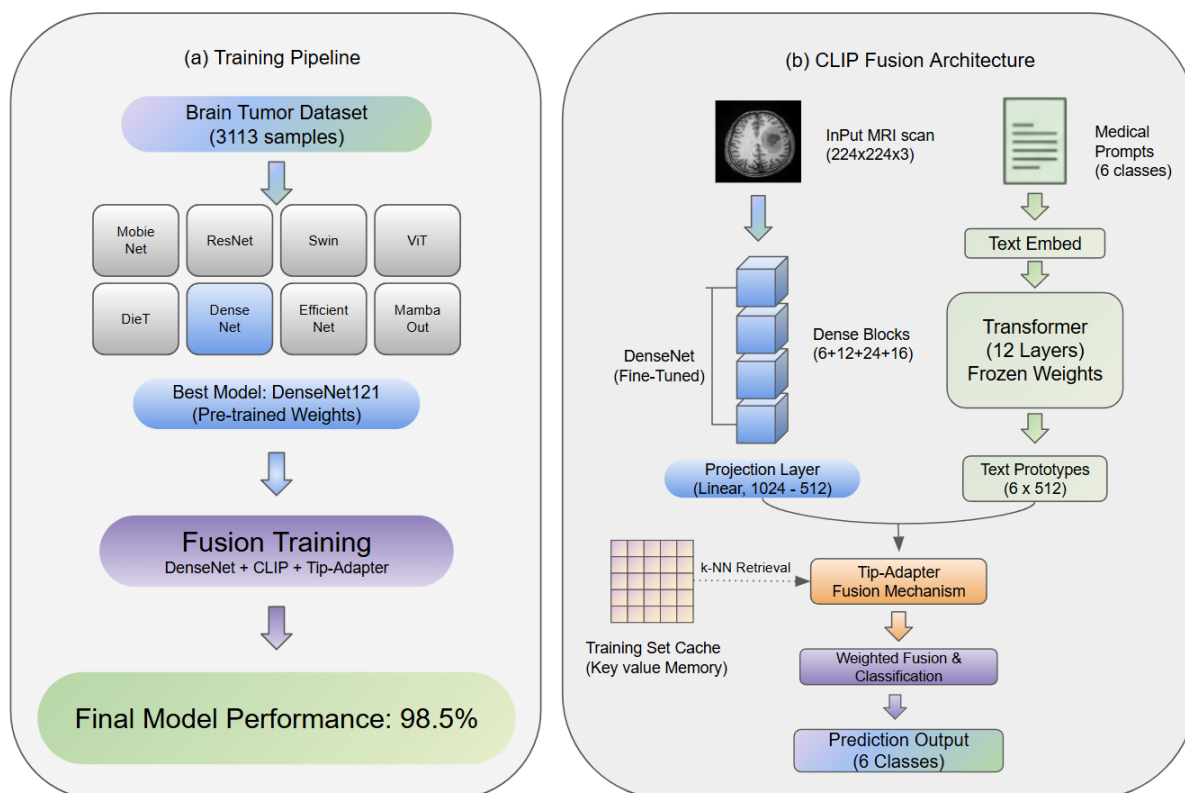


Figure 2: **TumorCLIP training pipeline and fusion architecture.** (a) Two-stage pipeline including backbone benchmarking and multimodal fusion. (b) CLIP fusion module integrating image embeddings with text prototypes using Tip-Adapter.

97 Classification Performance of TumorCLIP

98 Leveraging DenseNet121 as the visual backbone, TumorCLIP achieved a test accuracy of 98.5%, outper-
99 forming the unimodal DenseNet baseline (97.6%). Notably, TumorCLIP demonstrated particular strength
100 in underrepresented categories: as evidenced in Figure 3, recall for Neurocytoma increased by 1.86 percent-
101 age points, which is notable given the rarity and morphological diversity of this subtype. Visual analysis of
102 the confusion matrices further indicates a reduction in inter-class misclassification. This enhancement is at-
103 tributed to improved semantic alignment, in which class-level conceptual priors derived from text prototypes
104 help distinguish subtle morphological differences.

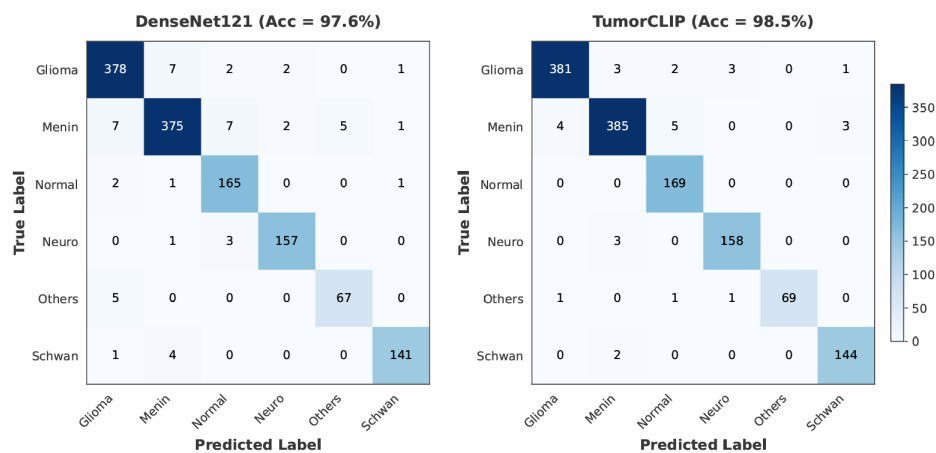
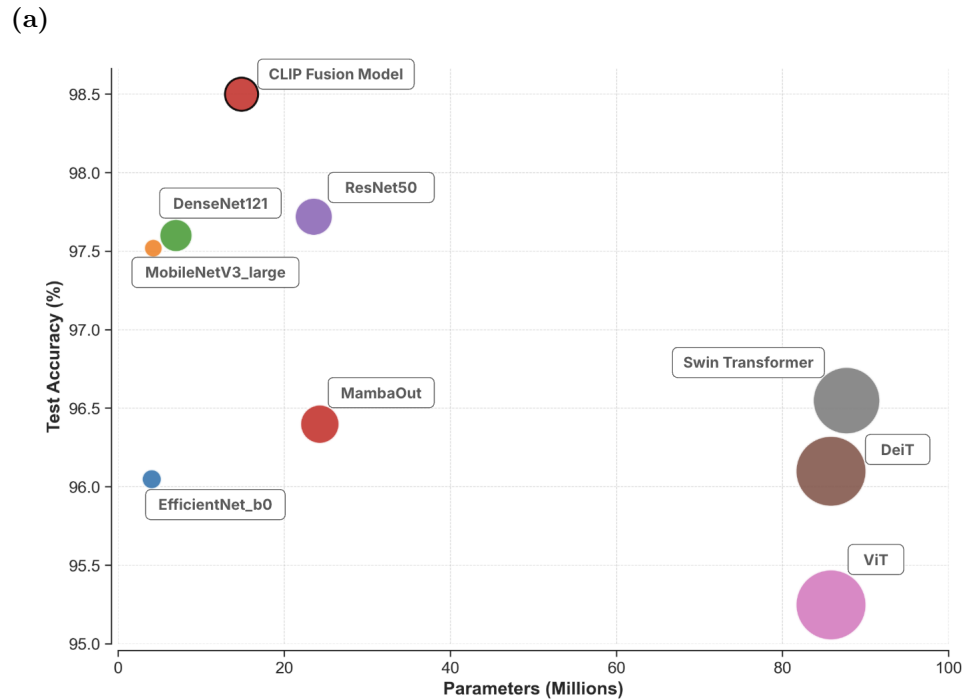


Figure 3: Confusion matrices of DenseNet121 and TumorCLIP. Comparison of subtype-level classification showing improved performance under TumorCLIP.

105 The integration of text-guided fusion sharpens decision boundaries by synergistically combining image-based
 106 features with clinically grounded radiological descriptions, thereby augmenting the model’s robustness to
 107 intraclass variability.

108 Crucially, TumorCLIP accomplishes these performance gains while maintaining computational efficiency. As
 109 illustrated in Figure 4, the multimodal model comprises only 14.84 million trainable parameters, which is
 110 fewer than Transformer-based baselines such as DeiT (22M), ViT (86M), and Swin Transformer (87.7M).
 111 The frozen CLIP text encoder (150M parameters), which remains unmodified during training, is excluded
 112 from this parameter count. Despite utilizing 5 to 6 times fewer trainable parameters than ViT or Swin,
 113 TumorCLIP achieves the highest accuracy among the evaluated models in our benchmark. This favorable
 114 parameter–accuracy trade-off positions TumorCLIP in the upper-left corner of the spectrum, underscoring
 115 its suitability for clinical applications with limited computational resources.



(b)

Model	Parameters (M)
EfficientNet b0	4.02
MobileNetV3 large	4.21
DenseNet121	6.96
ResNet50	23.52
MambaOut	24.25
DeiT	85.80
ViT	85.80
Swin Transformer	87.70
CLIP Fusion Trainable	14.84

Figure 4: **Model size and computational efficiency analysis.** (a) Test accuracy versus parameter count, with bubble size proportional to the computational cost (FLOPs) of each model. (b) Detailed breakdown of trainable parameters.

116 External Dataset Generalization

117 To evaluate model robustness under cross-dataset variability, we further assessed DenseNet121 and Tumor-
 118 CLIP on an independent external MRI dataset using the trained weights without any fine-tuning. As the
 119 external dataset does not include samples corresponding to the Outros Tipos de Lesões category, evaluation
 120 was restricted to the five tumor classes shared with the training dataset. Performance metrics and confusion
 121 matrices were therefore computed using a filtered five-class protocol.

122 As shown in Figure 5, both models exhibit a performance drop when transferred from the training distri-
 123 bution to the external dataset. However, TumorCLIP shows improved cross-dataset performance under the
 124 evaluated distribution shift, relative to the unimodal DenseNet121 baseline. In particular, the overall accu-
 125 racy degradation of TumorCLIP is smaller. This trend is further evident at the class level, where TumorCLIP

126 maintains higher accuracy for glioma, a clinically heterogeneous category known to be particularly sensitive
 127 to variations in acquisition protocols and data sources.

128 Confusion matrix analysis of the external dataset (Figure 5b) shows that TumorCLIP reduces inter-class
 129 confusion across all five evaluated categories compared to DenseNet121. Notably, misclassification between
 130 glioma and other tumor subtypes is substantially attenuated, suggesting that radiology-informed textual
 131 priors provide stabilizing semantic constraints beyond purely visual decision boundaries. These results
 132 indicate that TumorCLIP generalizes better on this external dataset under differences in data acquisition
 133 characteristics.

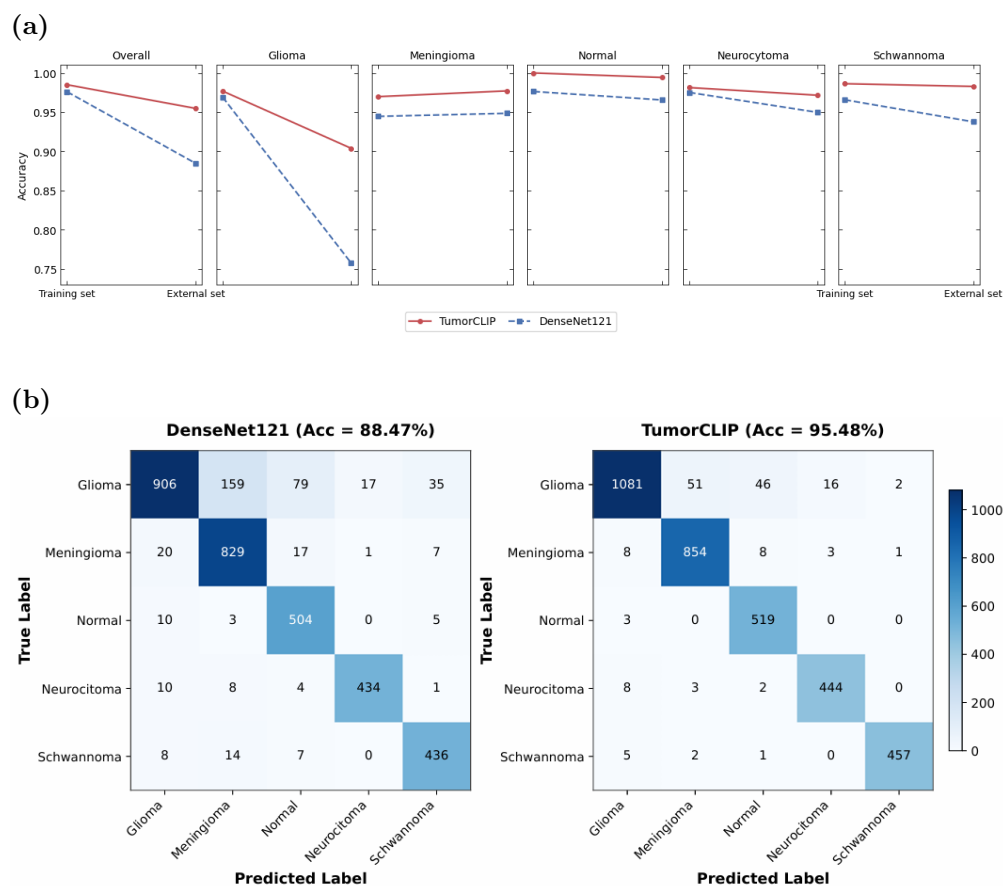


Figure 5: (a) Accuracy comparison between DenseNet121 and TumorCLIP on the training and external datasets across overall and tumor-specific categories. (b) Confusion matrices of DenseNet121 and TumorCLIP on the external dataset under a filtered five-class setting, excluding Outros Tipos de Lesões.

134 t-SNE Visualization

135 To qualitatively assess the impact of TumorCLIP on class separability, we visualized the model outputs using
 136 t-SNE, as shown in Figure 6. We performed t-SNE on the learned feature embeddings extracted from the
 137 last layer, enabling a direct examination of how each model structures its representation space rather than
 138 its final decision scores. This embedding-based visualization was applied consistently to both the training
 139 dataset and the independent external test dataset.

140 Both DenseNet121 and TumorCLIP generate six discernible clusters corresponding to the diagnostic cate-
 141 gories in the training set. However, TumorCLIP exhibits more compact and cohesive embedding clusters,

142 with reduced overlap at class boundaries, indicating improved class-wise organization in the learned feature
143 space. This effect is particularly evident in the glioma cluster, where TumorCLIP forms a denser manifold
144 with clearer separation compared to DenseNet121, consistent with its improved discrimination of visually
145 ambiguous tumor patterns.

146 When evaluated on the external dataset, the contrast between the two models becomes more pronounced.
147 DenseNet121 embeddings show increased dispersion and partial collapse between several tumor categories,
148 reflecting reduced robustness under distribution shift. In contrast, TumorCLIP preserves a more stable
149 geometric structure, maintaining recognizable cluster identities despite domain differences. Notably, minority
150 tumor classes exhibit less fragmentation in the TumorCLIP embedding space, suggesting that radiology-
151 guided semantic priors improve generalization beyond the training distribution. Overall, these embedding-
152 level visualizations demonstrate that TumorCLIP enhances class separability not only within the training
153 domain but also under external evaluation, yielding a more structured and resilient representation space.
154 By injecting radiological semantics through text prototypes and cache-based fusion, TumorCLIP reshapes
155 the embedding geometry, leading to more reliable class organization and improved robustness for rare and
156 challenging tumor categories.

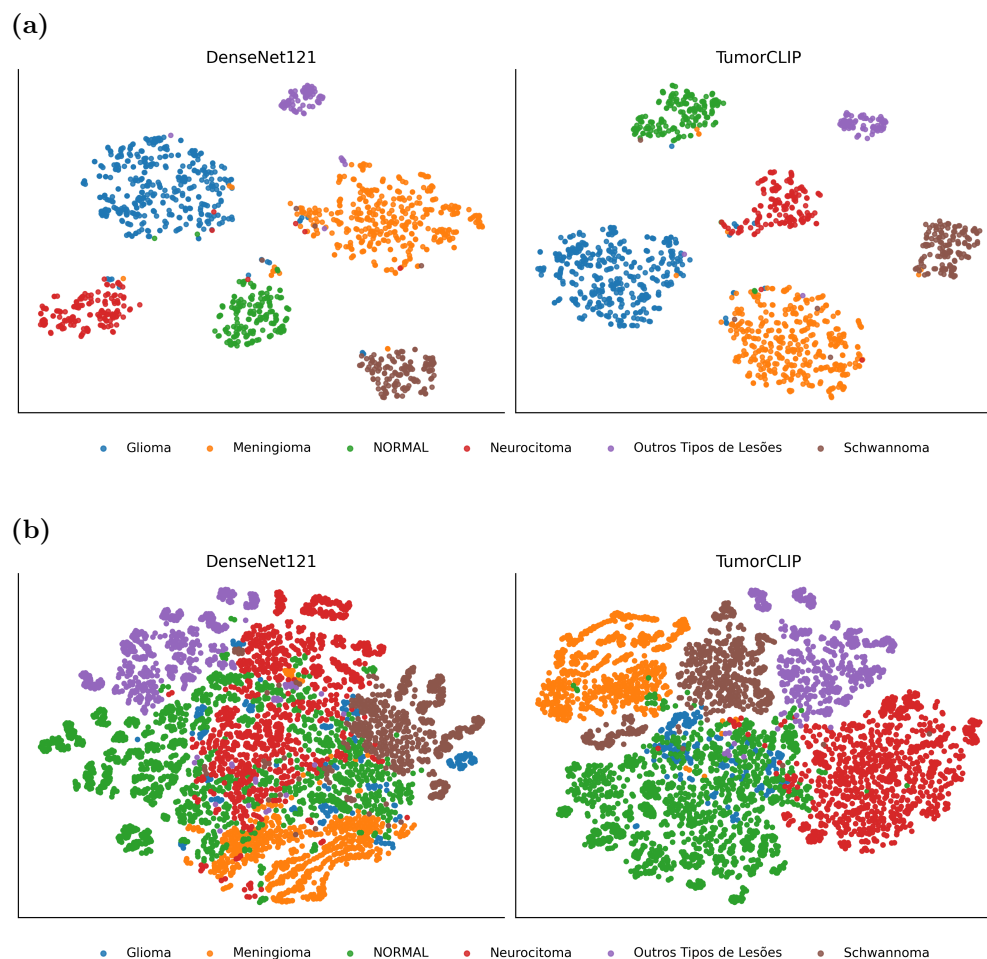


Figure 6: **Comparative t-SNE visualization of learned feature embeddings.** (a) t-SNE projections of embeddings extracted from the training dataset for DenseNet121 (left) and TumorCLIP (right). (b) Corresponding t-SNE projections of embeddings from the external test dataset. Each point represents a single sample and is colored by tumor category. TumorCLIP exhibits tighter intra-class compactness and clearer inter-class separation.

157 Methods

158 Model Overview

159 The CLIP text encoder is fully frozen throughout all experiments and is never fine-tuned. All trainable
160 parameters are restricted to the DenseNet classifier head, the lightweight adapter, and the fusion module.

161 Dataset and Preprocessing

162 The primary training and in-domain evaluation experiments utilized the publicly available Brain Tumor
163 MRI Images (17 Classes) dataset from Kaggle [14]; representative examples are shown in Figure 7. To re-
164 duce extreme class imbalance and improve clinical relevance, the original 17 tumor labels were merged into
165 six diagnostic super-classes based on lesion location and radiological appearance, including Glioma, Menin-
166 gioma, Schwannoma, Neurocytoma, Normal, and Outros Tipos de Lesões (Supplementary Table S1). This
167 grouping follows common radiological categorization principles and has been adopted in prior brain tumor
168 MRI classification studies. The dataset encompasses a diverse collection of axial MRI scans, incorporating

169 T1-weighted, contrast-enhanced T1 (T1c), and T2-weighted modalities.

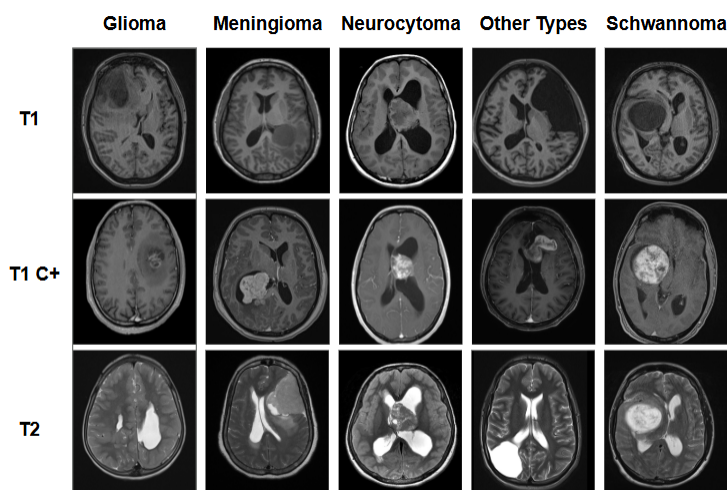


Figure 7: **Representative axial MRI slices from five brain tumour categories** (Glioma, Meningioma, Neurocytoma, Other Types of Injuries, and Schwannoma) across T1, contrast-enhanced T1 (T1C+), and T2 modalities. Shown cases are independent examples and are not derived from the same patient across modalities.

170 To further evaluate model generalization, an independent external MRI dataset was used for evaluation [15].
171 This dataset was collected from a different source and exhibits variations in acquisition characteristics and
172 label composition. As the external dataset does not include samples from the Outros Tipos de Lesões
173 category, external performance was assessed using a filtered protocol restricted to the five tumor categories
174 shared between the training and external datasets. During this evaluation, predictions and performance
175 metrics were computed exclusively over these overlapping classes, with samples from the excluded category
176 omitted from metric calculation.

177 Preprocessing was standardized across all models to ensure fair comparison. Each image was resized to
178 224×224 pixels, converted to a three-channel format, and normalized using ImageNet mean and standard
179 deviation statistics. The dataset was partitioned into training, validation, and test sets via stratified sam-
180 pling to maintain consistent class distributions. To enhance generalization while preserving comparability
181 across visual backbones, we applied lightweight data augmentation limited to random horizontal flipping and
182 minor in-plane rotations of $\pm 10^\circ$. No modality-specific transformations or intensity-based augmentations
183 were employed. This unified preprocessing and augmentation pipeline was consistently used for all models
184 throughout the study.

185 Radiology Text Prompts

186 To infuse the multimodal branch with clinically meaningful semantics, we manually authored radiology-style
187 text prompts for each tumor category. These prompts capture characteristic imaging phenotypes, including
188 anatomical location, signal intensity patterns, and enhancement behavior. All prompts were encoded using
189 a frozen CLIP text encoder, and prompt embeddings within the same category were averaged to form a
190 robust class-level radiology text prototype.

191 Specifically, we constructed five radiology-style text prompts for each of Glioma, Meningioma, and NORMAL,
192 and four prompts for each of Neurocytoma, Outros Tipos de Lesões, and Schwannoma, resulting in 27
193 multilingual prompts across six diagnostic categories. These radiology text prototypes serve as semantic

194 anchors, supporting both zero-shot inference and multimodal fusion during prediction.

195 Representative examples of the radiology-style text prompts used to construct the class-level radiology text
196 prototypes are shown in Table 1:

Table 1: Radiology text prototypes

Class	Example radiology-style text prompt
Glioma	Intra-axial infiltrative lesion with heterogeneous T2 hyperintensity
Meningioma	Extra-axial mass with dural tail and broad dural attachment
Schwannoma	Enhancing mass along the cisternal segment of a cranial nerve
Normal	No abnormal signal intensity, mass effect, or enhancement
Neurocytoma	Well-defined intraventricular mass with mixed T1/T2 signal
Outros	Heterogeneous lesion with irregular borders and variable enhancement

197 Unified Training Protocol and Objective Functions

198 To ensure a rigorously fair comparison across architectures, all eight visual backbones—EfficientNet-B0,
199 MobileNetV3-Large, ResNet50, DenseNet121, ViT, DeiT, Swin Transformer, and MambaOut—were trained
200 following an identical optimization protocol. Each model was evaluated using a unified hyperparameter grid
201 that included both SGD and Adam optimizers in combination with four learning rates ranging from 1×10^{-3}
202 to 1×10^{-6} . Batch size, data preprocessing, and augmentation strategies were held constant throughout all
203 experiments, and early stopping based on validation accuracy was employed to mitigate overfitting.

204 All models were trained to optimize the standard multi-class cross-entropy objective, expressed as:

$$L(\theta) = - \sum_{i=1}^C y_i \log(\text{softmax}(f_{\theta}(x))_i)$$

205 where C is the number of tumor classes and f_{θ} denotes the backbone network’s forward function.

206 Within this standardized experimental framework, DenseNet121 consistently exhibited the highest stability
207 and overall performance, making it the backbone of choice for the TumorCLIP architecture.

208 For the multimodal TumorCLIP model, we employed a composite loss function to jointly optimize the fusion
209 logits, the DenseNet branch, and the CLIP–Tip-Adapter branch. Let $\mathbf{s}_{\text{fused}}$ denote the fused logits, $\mathbf{s}_{\text{dense}}$
210 the DenseNet-only logits, and \mathbf{s}_{clip} the CLIP–Tip-Adapter logits. The total loss is defined as:

$$\mathcal{L}_{\text{total}} = 0.5 \text{CE}(\mathbf{s}_{\text{fused}}, y) + 0.3 \mathcal{L}_{\text{focal}}(\mathbf{s}_{\text{dense}}, y) + 0.2 \text{CE}(\mathbf{s}_{\text{clip}}, y),$$

211 where CE represents the standard multi-class cross-entropy and $\mathcal{L}_{\text{focal}}$ denotes the focal loss applied to the
212 DenseNet classifier. During multimodal training, the inclusion of a focal loss component for the DenseNet
213 branch specifically targets hard and minority-class examples, while maintaining the cross-entropy loss used
214 in the unimodal benchmark. This multi-task objective encourages the fusion head to remain aligned with
215 both the unimodal DenseNet baseline and the multimodal CLIP branch, while ensuring robust learning on
216 challenging samples. The fusion loss is assigned the highest weight to prioritize multimodal decision learning,
217 while focal loss is applied to the DenseNet branch to improve performance on minority classes. The CLIP
218 branch is treated as auxiliary supervision to stabilize semantic alignment during training.

219 Construction of Radiology Text Prototypes

220 To embed clinically meaningful semantic priors within the model, we authored radiology-style textual de-
221 scriptions for each diagnostic category. These prompts were designed to encapsulate characteristic MRI
222 features commonly referenced in clinical interpretation, such as typical lesion locations (e.g., intra-axial,
223 extra-axial, intraventricular), signal intensity patterns, enhancement characteristics, and class-specific mor-
224 phological cues. Each prompt was processed through the frozen CLIP text encoder. For each class k , all
225 associated text embeddings were averaged to generate a class-level text prototype:

$$\mu_k = \frac{1}{N_k} \sum_{j=1}^{N_k} z_{\text{text},j}$$

226 where N_k is the number of prompts constructed for class k . Across the six diagnostic groups, each class
227 was represented by approximately six prompts. These class prototypes serve as semantic anchors, providing
228 reference points for aligning DenseNet-derived image features during both zero-shot inference and multimodal
229 fusion within TumorCLIP.

230 Image-feature Cache Construction

231 After identifying DenseNet121 as the optimal unimodal backbone, we utilized the trained network to extract
232 feature representations from all training images, thereby constructing a non-parametric cache for use in the
233 Tip-Adapter modules. For each training instance x_i , the model generated a normalized feature embedding:

$$v_i = \text{DenseNet121}(x_i)$$

234 This cache comprised two key components:

- 235 (1) the collection of visual embeddings $V = \{v_i\}$, which serve as the key space for similarity-based retrieval.
- 236 (2) the corresponding one-hot label vectors $Y = \{y_i\}$, which encode class-specific supervision signals.

237 During inference, this cache remains fixed and provides instance-level evidence through similarity-weighted
238 retrieval, supporting few-shot adaptation without necessitating further training or fine-tuning of the backbone
239 network.

240 Tip-Adapter Formulation

241 TumorCLIP fuses semantic information from radiology-informed text prototypes with instance-level evidence
242 from the feature cache using a lightweight adaptation strategy inspired by the Tip-Adapter framework. For
243 each test image, the DenseNet121 backbone generates an embedding v , which is normalized and subsequently
244 compared to the radiology-derived text prototypes. Prior to similarity computation, v is passed through a
245 lightweight two-layer adapter MLP ($512 \rightarrow 128 \rightarrow 512$); the adapter and the image feature projection layer
246 are jointly optimized, while the CLIP text encoder remains frozen.

247 Text-prototype logits are calculated via cosine similarity:

$$s_{\text{text},k} = \cos(v, \mu_k)$$

248 To incorporate supervision from labeled examples, the model retrieves signals from the cached training
249 embeddings. For each cached feature v_i , a similarity score is computed using cosine similarity and scaled by
250 a temperature parameter t_{knn} :

$$s_i = \cos(v, v_i)$$

$$a_i = \frac{\exp(s_i/t_{\text{knn}})}{\sum_j \exp(s_j/t_{\text{knn}})}$$

251 where t_{knn} modulates the neighborhood selectivity of the cache retrieval. Cache-based logits are then com-
252 puted by aggregating the normalized similarity weights over the associated one-hot labels:

$$s_{\text{cache},k} = \sum_i a_i y_i[k]$$

253 Fusion of semantic (text) and cache-based evidence is achieved through a scalar alpha: α :

$$s_k^{\text{tip}} = (1 - \alpha) s_k^{\text{text}} + \alpha s_k^{\text{cache}}$$

254 Simultaneously, the DenseNet backbone outputs conventional classifier logits for each class k , denoted
255 $s_{\text{dense},k}$. At the model level, the Tip-Adapter logits are combined with the DenseNet logits via a learn-
256 able fusion weight $w \in (0, 1)$, resulting in the fused logits:

$$s_{\text{fused},k} = (1 - w) s_{\text{dense},k} + w s_{\text{tip},k}$$

257 The parameter w is implemented as the sigmoid of an unconstrained scalar and is jointly optimized with the
258 adapter parameters. The final class prediction is made by selecting the class corresponding to the maximal
259 fused logit:

$$\hat{y} = \arg \max_k (s_{\text{fused},k})$$

260 This formulation enables multiple inference modes within a unified framework. Zero-shot classification uti-
261 lizes only the text-prototype logits $s_{\text{text},k}$, disabling both the cache-based logits and the DenseNet classifier
262 head. Few-shot adaptation incorporates cache-derived logits $s_{\text{cache},k}$, while the complete TumorCLIP model
263 leverages both text-prototype and cache-based evidence via a fixed mixing coefficient α (tuned as a hyper-
264 parameter) and a learnable fusion weight w . Importantly, the CLIP text encoder remains entirely frozen
265 throughout all stages.

266 The feature cache contains embeddings of all training samples extracted using the frozen image encoder.
267 Similarity-based retrieval from the cache is controlled by a temperature parameter t_{knn} , which corresponds
268 to $\beta = 1/t_{\text{knn}}$ in the exponential weighting formulation. In all experiments, t_{knn} is set to 0.07.

269 The lightweight adapter is trained for $N = 15$ epochs using the same optimizer and learning rate settings
270 as the DenseNet classifier, while all other components, including the text encoder, remain frozen. All
271 hyperparameters are kept fixed across experiments to ensure fair comparison and reproducibility.

272 Adaptive Inference Under Limited Data

273 TumorCLIP supports a range of data-efficient inference modes by decoupling visual feature extraction from
274 semantic reasoning. In the zero-shot scenario, predictions are based solely on text prototypes derived from
275 radiology-informed prompts, with the DenseNet121 backbone generating image embeddings that are directly
276 matched to these class-level semantic anchors. This approach allows for effective classification without any
277 labeled training data, enabling meaningful decision-making in fully unsupervised settings.

278 When limited labeled data are available, TumorCLIP leverages instance-level visual evidence through a
279 non-parametric feature cache. Training image features are stored once in this cache, and during inference,
280 visually similar examples are retrieved to complement predictions informed by the text prototypes. This

281 cache-guided refinement enhances decision boundaries for underrepresented classes while maintaining the
282 frozen state of both visual and text encoders.

283 For cases requiring greater adaptability, a lightweight trainable adapter can be introduced to better align
284 the visual embedding space with both cache-derived features and text prototypes. By restricting trainable
285 parameters to this adapter, the multimodal design remains computationally efficient.

286 Collectively, these inference strategies define a unified continuum of adaptation—from zero-shot semantic
287 reasoning, to training-free instance-based refinement, to targeted lightweight tuning—ensuring robust perfor-
288 mance across diverse data availability scenarios without necessitating modification of the backbone network.

289 Evaluation Metrics

290 Model performance was evaluated using established multi-class assessment metrics. To address class im-
291 balance among the six diagnostic categories, we primarily report overall accuracy. In addition, class-wise
292 performance was analyzed using confusion matrices to characterize sensitivity patterns and subtype-level
293 misclassification across tumor categories. For external dataset evaluation, all metrics were computed using
294 the same definitions, and performance was assessed under a filtered protocol restricted to tumor categories
295 shared across datasets. Both zero-shot and few-shot evaluation protocols mirrored those of the fully fused
296 TumorCLIP model, with the sole distinction being whether cache-derived logits were included or excluded
297 during inference. No task-specific thresholds or post-processing adjustments were applied, ensuring com-
298 parability across evaluation modes. To qualitatively assess class separability in the decision space, t-SNE
299 was applied to the models' output prediction signatures. This visualization served as a complementary
300 tool to evaluate cluster compactness and the distinctness of inter-class boundaries. Unless otherwise noted,
301 reported results correspond to the model checkpoint achieving the highest validation accuracy across all
302 training epochs, and test metrics are derived from this selected checkpoint.

303 All hyperparameters are fixed across experiments. Evaluation on the external dataset is conducted without
304 any fine-tuning or adaptation.

305 Software

306 All models were implemented using PyTorch. The CLIP text encoder was initialized with publicly available
307 pretrained weights. Image preprocessing employed standard Python libraries, and no custom CUDA kernels
308 or proprietary implementations were utilized.

309 Discussion

310 This work presents TumorCLIP, a lightweight vision-language framework developed to address persistent
311 limitations in deep learning-based brain tumor MRI classification. Through a systematic benchmark of eight
312 visual backbones under a unified hyperparameter grid, DenseNet121 was identified as a robust and stable
313 foundation for multimodal extension. Building on this backbone, TumorCLIP incorporates frozen clinical
314 text descriptions via a Tip-Adapter mechanism, conferring several advantages.

315 Firstly, TumorCLIP advances explainability through concept-level reasoning, as text prototypes derived from
316 radiological descriptions anchor model predictions in human-interpretable semantics, which is a functionality
317 lacking in vision-only models. Secondly, the framework boosts classification performance, especially for
318 underrepresented classes, by integrating semantic priors with image-based features. Thirdly, TumorCLIP
319 is computationally efficient: by freezing the CLIP text encoder and training only a lightweight adapter,
320 training costs are lower than those of fully end-to-end multimodal architectures. Lastly, the framework
321 provides intrinsic zero-shot and few-shot capabilities, highlighting its deployment flexibility in scenarios with
322 limited labeled data.

323 Beyond in-domain evaluation, additional validation on an independent external dataset demonstrates Tu-
324 morCLIP's robustness to cross-dataset variability. While both DenseNet121 and TumorCLIP degrade in
325 performance when evaluated on data from a different source, the decline observed for DenseNet121 is more

326 pronounced. This effect is especially evident for glioma, a highly heterogeneous and infiltrative tumor type
327 whose imaging appearance is sensitive to variations in acquisition protocols and scanner characteristics.
328 Conventional CNNs may therefore over-rely on dataset-specific visual cues, leading to degraded performance
329 when such cues shift across datasets. In contrast, TumorCLIP aligns visual representations with radiology-
330 informed semantic concepts, encouraging reliance on diagnostically relevant pathological patterns rather
331 than superficial imaging style. This multimodal semantic constraint contributes to the improved robustness
332 and generalization observed across datasets.

333 Our findings underscore that integrating structured radiology knowledge can improve the reliability and
334 clinical interpretability of MRI-based tumor classification systems in our experiments. Future research will
335 seek to broaden the diversity of text prototypes, incorporate structured radiology ontologies, and validate
336 the model’s generalizability across multi-institutional datasets. A limitation of this study is the reliance
337 on manually authored radiology-style text prototypes, which may introduce variability depending on expert
338 knowledge and reporting conventions. In addition, differences in imaging protocols or institutional practices
339 may affect the transferability of learned semantic alignments across datasets.

340 References

- 341 [1] F.J. Dorfner, J.B. Patel, J. Kalpathy-Cramer, et al. A review of deep learning for brain tumor analysis
342 in mri. *NPJ Precision Oncology*, 9:2, 2025.
- 343 [2] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use
344 interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, 2019.
- 345 [3] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
346 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from
347 natural language supervision. In *International Conference on Machine Learning (ICML)*, pages 8748–
348 8763. PMLR, 2021.
- 349 [4] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convo-
350 lutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*
351 *(CVPR)*, pages 4700–4708, 2017.
- 352 [5] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks.
353 In *International Conference on Machine Learning (ICML)*, pages 6105–6114. PMLR, 2019.
- 354 [6] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang,
355 Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the*
356 *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1314–1324, 2019.
- 357 [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
358 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
359 770–778, 2016.
- 360 [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
361 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image
362 is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on*
363 *Learning Representations (ICLR)*, 2021.
- 364 [9] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé
365 Jégou. Training data-efficient image transformers & distillation through attention. In *International*
366 *Conference on Machine Learning (ICML)*, pages 10347–10357. PMLR, 2021.
- 367 [10] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin
368 transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF*
369 *International Conference on Computer Vision (ICCV)*, pages 10012–10022, 2021.

- 370 [11] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
371 *preprint arXiv:2312.00752*, 2023.
- 372 [12] Muhammad Aamir, Abdellah Namoun, Siraj Munir, Naif Aljohani, Muflih H Alanazi, Yousef Alsahafi,
373 and Fahad Alotibi. Brain tumor detection and classification using an optimized convolutional neural
374 network. *Diagnostics*, 14(16):1714, 2024.
- 375 [13] Renrui Zhang, Ziyao Wei, Wei Zhang, Rongyao Fang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and
376 Hongsheng Li. Tip-adapter: Training-free adaption of clip for few-shot classification. *IEEE Transactions*
377 *on Pattern Analysis and Machine Intelligence (TPAMI)*, 2023.
- 378 [14] Fernando2rad. Brain tumor mri images (17 classes). [https://www.kaggle.com/datasets/](https://www.kaggle.com/datasets/fernando2rad/brain-tumor-mri-images-17-classes)
379 [fernando2rad/brain-tumor-mri-images-17-classes](https://www.kaggle.com/datasets/fernando2rad/brain-tumor-mri-images-17-classes), 2024. Accessed: 2025-09-20.
- 380 [15] Nagahhenes. W. brain tumor for 14 classes. [https://www.kaggle.com/datasets/waseemnagahhenes/](https://www.kaggle.com/datasets/waseemnagahhenes/brain-tumor-for-14-classes)
381 [brain-tumor-for-14-classes](https://www.kaggle.com/datasets/waseemnagahhenes/brain-tumor-for-14-classes), 2023. Accessed: 2025-11-20.

382