

Dynamic Residual Encoding with Slide-Level Contrastive Learning for End-to-End Whole Slide Image Representation

Jing Jin*

School of Computer Science and
Engineering, Central South University
Changsha, China
jingjin@csu.edu.cn

Xu Liu*

School of Computer Science and
Engineering, Central South University
Changsha, China
lx6711@csu.edu.cn

Te Gao

School of Computer Science and
Engineering, Central South University
Changsha, China
gaote1021@csu.edu.cn

Zhihong Shi

School of Computer Science and
Engineering, Central South University
Changsha, China
shizhihong@csu.edu.cn

Yixiong Liang

School of Computer Science and
Engineering, Central South University
Changsha, China
yqliang@csu.edu.cn

Ruiqing Zheng

School of Computer Science and
Engineering, Central South University
Changsha, China
rqzheng@csu.edu.cn

Hulin Kuang[†]

School of Computer Science and
Engineering, Central South University
Changsha, China
hulinkuang@csu.edu.cn

Min Zen

School of Computer Science and
Engineering, Central South University
Changsha, China
zengmin@csu.edu.cn

Shichao Kan[†]

School of Computer Science and
Engineering, Central South University
Changsha, China
kanshichao@csu.edu.cn

Abstract

Whole Slide Image (WSI) representation is critical for cancer subtyping, cancer recognition and mutation prediction. Training an end-to-end WSI representation model poses significant challenges, as a standard gigapixel slide can contain tens of thousands of image tiles, making it difficult to compute gradients of all tiles in a single mini-batch due to current GPU limitations. To address this challenge, we propose a method of dynamic residual encoding with slide-level contrastive learning (DRE-SLCL) for end-to-end WSI representation. Our approach utilizes a memory bank to store the features of tiles across all WSIs in the dataset. During training, a mini-batch usually contains multiple WSIs. For each WSI in the batch, a subset of tiles is randomly sampled and their features are computed using a tile encoder. Then, additional tile features from the same WSI are selected from the memory bank. The representation of each individual WSI is generated using a residual encoding technique that incorporates both the sampled features and those retrieved from the memory bank. Finally, the slide-level contrastive loss is computed based on the representations and histopathology reports of the WSIs within the mini-batch. Experiments conducted over cancer subtyping, cancer recognition, and mutation prediction tasks proved the effectiveness of the proposed DRE-SLCL method.

CCS Concepts

• **Whole Slide Image Representation**; • **Dynamic Residual Encoding**; • **Contrastive Learning**;

Keywords

Residual Encoding, Contrastive Learning, Slide Image Representation

ACM Reference Format:

Jing Jin*, Xu Liu*, Te Gao, Zhihong Shi, Yixiong Liang, Ruiqing Zheng, Hulin Kuang[†], Min Zen, and Shichao Kan[†]. 2025. Dynamic Residual Encoding with Slide-Level Contrastive Learning for End-to-End Whole Slide Image Representation. In *Proceedings of MM '25, October 27–31, 2025, Dublin, Ireland (MM '25)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3746027.3755469>

1 Introduction

Whole Slide Image (WSI) representation is essential in computational pathology (CPath) [30, 43], serving as the foundation for tasks such as cancer subtyping [6, 40], cancer recognize [24, 38, 46], mutation prediction [14, 22], cancer staging [3, 17, 47], and prognostic assessment [1, 8, 13]. Given that a standard gigapixel slide may consist of tens of thousands of image tiles, making it impossible for all tiles of WSIs in a batch to be updated in each epoch due to GPU limitations. Current WSI representation methods typically follow a two-stage approach, as illustrated in Figure 1 (a). In the first stage, a tile representation model is trained using self-supervised techniques to capture individual tile features within the WSI. Then, these tile representations are aggregated through models like LongNet in Prov-GigaPath [58] or ViT in the Hierarchical Image Pyramid Transformer (HIPT) [6]. However, this two-stage approach may be suboptimal for WSI representation, as WSI labels

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
MM '25, Dublin, Ireland

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-2035-2/2025/10
<https://doi.org/10.1145/3746027.3755469>

*Both authors contributed equally to this research.

[†]Corresponding authors.

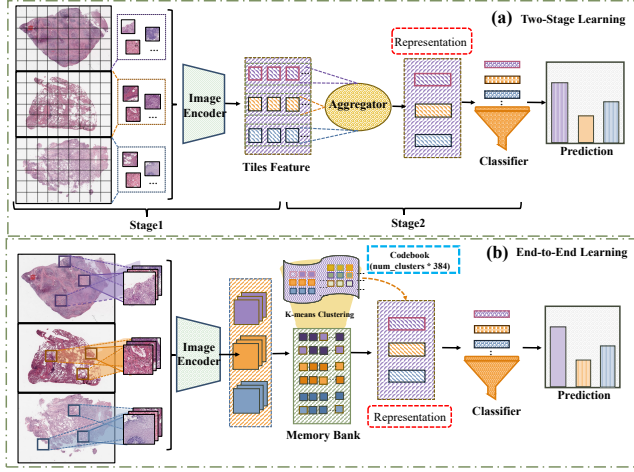


Figure 1: WSI representation learning. (a) The pipeline of the two-stage learning methods and (b) the proposed end-to-end learning method.

are not utilized in the initial learning stage. Therefore, we aim to develop an end-to-end learning method to enhance WSI representation by explicitly leveraging label gradients. This enables direct gradient flow between tile-level and slide-level representations, improving feature learning and overall performance.

Popular methods for WSI representation are weakly supervised and use Multiple Instance Learning (MIL) [16, 18, 25, 36, 55, 56], where MIL algorithms assign slide-level labels to a set of instances and aggregate these instances through feature pooling [4, 61], Transformer [2, 12, 32, 36, 53] or graph neural networks [27, 28, 33, 34, 39, 45]. Classification tasks are then performed using the aggregated features and slide-level labels. Although these methods have shown promising progress, the instance sampling process can affect model performance; specifically, the true label for a set of instances may differ from the slide-level label. To address this issue, we use a memory bank to store the features of all tiles, then aggregate the dynamically sampled tile features in a mini-batch with the stored tile features using a residual encoding method called the vector of locally aggregated descriptors (VLAD) [21] during training, which offers linear computational complexity for processing thousands of tiles per WSI while effectively preserving fine-grained morphological variations through its residual-based encoding mechanism. Since the final representation is based on the features of all tiles in a WSI and the model is trained end-to-end, this approach mitigates the sampling problem and improves the performance of the WSI representation, as illustrated in Figure 1 (b).

Beyond addressing computational constraints, using contrastive learning to align WSI visual features with corresponding pathology reports offers significant advantages, enhancing model generalization through cross-modal supervision. This approach has demonstrated considerable promise in natural and biomedical image processing. In CPath, Prov-GigaPath [58] incorporated pathology reports to construct a contrastive loss function, optimizing the model and enhancing performance in cancer subtyping and mutation prediction. However, they didn't evaluate effectiveness of

contrastive learning in an end-to-end framework. In this work, we integrate contrastive learning directly into our end-to-end training pipeline and validate the effectiveness of this approach using the lung cancer dataset from The Cancer Genome Atlas (TCGA) and the Clinical Proteomic Tumor Analysis Consortium (CPTAC). Unlike most contrastive learning methods that rely on CLIP language encoders [35], we use the LLaMA2-7B model [44] to encode pathological reports, applying contrastive loss exclusively to slides with available diagnostic reports.

To train the model end-to-end, we implement a dynamic parameter update strategy: a mini-batch of tiles is sampled from different WSIs, and WSI features are updated as tile features are modified. The updated WSI features are then used to compute cross-entropy and contrastive losses, driving parameter updates across the model. This end-to-end approach optimizes both the feature extraction and classification components, resulting in improved overall performance. During testing, the trained image encoder extracts features for all tiles within a WSI, which are then encoded using a residual encoding method.

The contributions of this work are summarized as follows:

- A method of dynamic residual encoding with slide-level contrastive learning (DRE-SLCL) is proposed for end-to-end WSI representation.
- We implement a dynamic parameter update strategy by incorporating a memory bank for end-to-end WSI representation learning, along with contrastive learning between the WSI and report representations directly within the end-to-end training pipeline.
- Our experimental results on the lung cancer dataset for cancer subtyping, cancer recognition, and mutation prediction proved the effectiveness of the proposed DRE-SLCL method.

2 Related Work

Our work is related to multiple instance learning, contrastive learning, and existing end-to-end approaches. In the following, we review works related to them.

2.1 Multiple Instance Learning

Multiple Instance Learning (MIL) is a foundational approach in computational pathology, particularly for WSI analysis [5, 7, 18, 25, 29, 54, 55], excelling when only slide-level annotations are available. WSIs are treated as "bags" containing multiple tile instances, with the model aggregating information to infer slide-level predictions.

Early MIL implementations employed simple max-pooling and average-pooling aggregation strategies for weakly supervised learning [4, 61]. Gradually, advanced methods have introduced attention mechanisms [19]. And then transformer-based MIL methods gained rapid popularity after Shao et al. introduced TransMIL [36]. Following TransMIL, several key innovations emerged: Bian et al.'s approach Mixed Supervision Transformer [2] utilized hybrid supervision with limited pixel annotations; Wang et al.'s MDMIL [53] reduced computational demands through multiplex detection cross-attention; and Ma et al.'s VINO [32] extended the paradigm to pathological video analysis.

Graph-based MIL approaches model tile relationships by representing WSIs as networks, with early GNNs [28, 34] treating

tiles as nodes and spatial connections as edges to capture multi-level contextual information. Recent advancements have expanded this paradigm: Ma et al.’s GCN-based MIL [33] used graph convolution and attention for structural relationship analysis; Shu et al.’s SlideGCD [39] pioneered inter-slide association through node classification; and Tu et al.’s GAMMIL [45] leveraged graph attention for efficient multi-scale information fusion with reduced noise interference.

Additionally, several foundational hybrid approaches have influenced the field. CLAM [31] introduced a clustering-constrained attention mechanism grouping similar instances before aggregation, balancing computational demands with discriminability. HIPT [31] presented hierarchical learning using ViT for multi-scale processing before MIL aggregation. These earlier works remain important benchmarks due to their significant contributions.

2.2 Contrastive Learning

Contrastive Learning [23, 59] has become a key method in representation learning, especially for high-dimensional data such as medical and biological data [26, 48, 52]. The core objective is to learn effective feature representations by bringing similar instances closer in latent space while pushing dissimilar instances apart.

Early contrastive learning methods, such as SimCLR [9] and MoCo [15], successfully learned visual representations through data augmentation, generating positive pairs while treating other samples as negatives. These approaches demonstrated that effective representations could be learned without labels by maximizing agreement between differently augmented views of the same data sample.

Recently, contrastive learning has been employed in computational pathology to enhance WSI representations, effectively extracting key features for cancer detection and tissue segmentation [42, 49, 50, 62]. Wang et al. proposed SCL-WC [51], which used cross-slide contrastive learning to improve weakly-supervised classification through both intra-WSI and inter-WSI information exchange. CPLIP [20] formulates a multi-image–multi-text bag-level alignment strategy for whole-slide image understanding and is primarily designed for zero-shot classification and segmentation in histopathology. Other innovations include KEP [63] that introduces a knowledge-enhanced pretraining paradigm by distilling structured medical knowledge from the PathKT into the text encoder. More recently, Yu et al. introduced CLOVER [60], which integrates multi-omics data with WSI representations through contrastive learning to connect morphological and molecular tumor features.

Notably, a whole-slide pathology foundation model named ProV-GigaPath [58] integrates contrastive learning with multimodal data by aligning WSI visual features with pathology report semantics through a dual-encoder architecture greatly improving cancer subtyping and prognostic prediction performance.

2.3 End-to-End Training

Despite the clear advantages of end-to-end approaches, the implementation of end-to-end models for WSI classification remains extremely limited in the literature. To our knowledge, only a handful of studies have attempted to realize this paradigm, and most of these are relatively early works with limited effectiveness.

Chikontwe et al. [11] proposed a center embedding approach that employed joint learning of instance-level and bag-level classifiers with top-k instance sampling. Takahama et al. [41] proposed a multi-stage approach combining patch-based classification with whole slide-scale segmentation, addressing GPU memory limitations by retaining gradient information while training the model partially. Xie et al. [57] introduced an end-to-end part learning approach that divided patches from a WSI into k parts based on global clustering centroids and sampled k tiles in each training epoch. Sharma et al. [37] presented the Cluster-to-Conquer framework that uses adaptive attention with KL-divergence regularization for intra-cluster attention variance. Recently, Chen et al. [10] integrated feature extraction and classification into a unified framework with attention pooling and a squared average normalization function.

3 Methods

Table 1: Four key components and their brief descriptions.

COMPONENT	BRIEF DESCRIPTION
Memory bank	A two-level nested dictionary that stores all tile features in the dataset, with mini-batch feature updates performed iteratively during training.
Codebook	A set of prototype vectors constructed via K-means clustering on all tile features from the memory bank, kept fixed throughout all training stages.
Tile encoder	A ViT-256 model, pretrained using the Hierarchical Image Pyramid Transformer (HIPT) framework, to encode tile-level features.
Report encoder	The LLaMA2-7B model, where we use the average of all token embeddings from the final hidden layer as the representation.

3.1 Method Overview

The overall framework of our method is illustrated in Figure 2. It consists of a feature extraction network, a VLAD encoding module, and a contrastive learning module. The feature extraction network extracts tile-level features within a WSI, while the VLAD encoding module aggregates these tile features into a comprehensive WSI representation. Contrastive learning, computed with the WSI feature and the WSI report feature, aims to enhance the generalization capability of the model. For image encoding, we use the ViT-256 model pretrained on HIPT, and the WSI report encoder is the text encoder of the LLaMA2 7B model. The key components of this framework are listed in Table 1. In the following subsections, we provide details of our method.

3.2 Dynamic Residual Encoding for WSI Representation

Since each WSI typically contains thousands of image tiles, directly storing and computing features for all tiles would be computationally impractical. To address this, we construct a dynamic feature

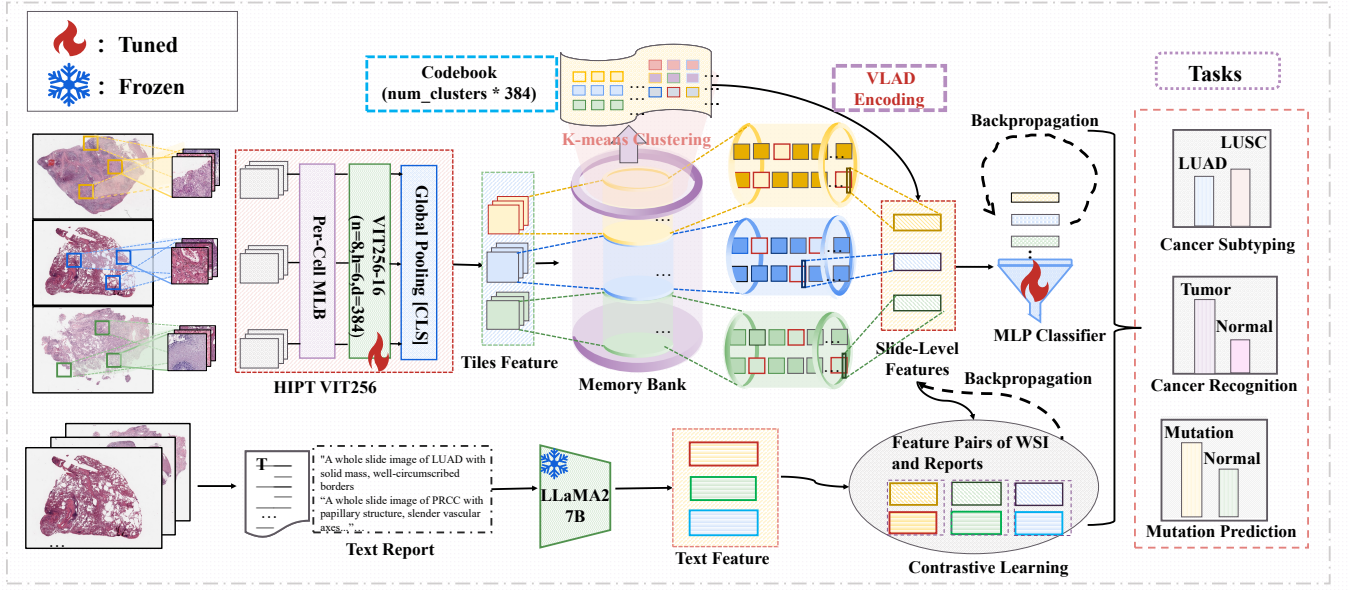


Figure 2: The end-to-end whole slide image representation framework with the proposed dynamic residual encoding with slide-level contrastive learning (DRE-SLCL).

memory bank to manage and update features efficiently. Specifically, we start by dividing the WSI into multiple image tiles using a fixed-size grid, denoted as $x_i \in \mathbb{R}^{256 \times 256}$, which is thus well-suited for a wide range of pathology tasks. Each tile is then processed through the VIT-256 model, pretrained within the HIPT framework, to extract its feature representation $f(x_i)$. Then, these tile features are stored in the memory bank. The memory bank is structured as a nested two-level dictionary. The keys in the first level represent unique WSI identifiers, while the keys in the second level index image tiles within each WSI. Each value stores the L2-normalized feature representation $f(x_i)$ of a tile. This hierarchical dictionary setup efficiently manages tile features, providing quick access to specific WSIs and their tile locations, thereby optimizing both storage and retrieval processes.

After extracting the tile features, we adopt the VLAD module to encode the tile features in a WSI into a WSI feature. First, a codebook is constructed by applying K -means clustering to all WSIs' tile features stored in the memory bank, as follows:

$$C = \{c_1, c_2, \dots, c_K\}. \quad (1)$$

The codebook C , containing K cluster centers, captures morphological patterns across diverse tissue types and remains fixed during training. This static design ensures consistent residual encoding by maintaining a stable semantic reference space, preventing semantic drift that would occur with dynamic codebook updates.

For each WSI, all its tile features are sampled from the memory bank and quantized to their nearest codewords, denoted as c_{nearest} . The residual for each feature $f(x_i)$ with respect to its assigned codeword is subsequently computed as:

$$r_i = f(x_i) - c_{\text{nearest}}. \quad (2)$$

The residual vector r_i represents the deviation of the feature $f(x_i)$ from its assigned codeword. These residuals effectively capture the subtle differences between the local features and the cluster centers.

Next, the residual vectors of all tile features from this WSI are aggregated to obtain a consolidated residual representation relative to each codeword:

$$v_k = \sum_{x_i \in \mathcal{X}_k} r_i, \quad k = 1, 2, \dots, K, \quad (3)$$

where \mathcal{X}_k represents the set of all image tiles most closely associated with codeword c_k . This approach accumulates all relevant residuals for each cluster center, thus generating a global representation that captures the relationship between the local features of the entire WSI and the cluster centers.

Finally, the accumulated residuals of all codewords are concatenated to form the residual encoding representation for each WSI, as follows:

$$v_{\text{re}} = [v_1, v_2, \dots, v_K] \in \mathbb{R}^{K \times d}, \quad (4)$$

where d denotes the dimension of each feature vector. To further enhance the discriminative power of the residual encoding vector, the L2 normalization is applied to the VLAD representation.

After generating the residual encoding vector, we further enhance its expressive power by introducing a Transformer layer for feature enhancement.

In each iteration during training, we randomly sample a batch of slides and select a subset of tiles from these slides to form a mini-batch. When tile features are updated, the corresponding residual encoding representation of the WSI is recalculated. We refer to this process as dynamic residual encoding.

Algorithm 1: Dynamic Residual Encoding with Slide-Level Contrastive Learning (DRE-SLCL)

Input: Whole Slide Image (WSI) $X_i = \{x_{i,1}, \dots, x_{i,N}\}$ (split into N tiles), Pathology report T_i

Hyperparameters: codebook size K , sampled tiles per WSI r , batch size b , temperature τ

Output: Slide-level prediction \hat{Y}_i

Workflow:

1: Dynamic Tile Sampling & Feature Update

for each WSI X_i in batch do
 Randomly sample r tiles: $S_i \leftarrow \text{RandomSelect}(X_i, r)$
 for each tile $x_{i,j} \in S_i$ do
 Extract feature: $f_{i,j} \leftarrow \text{ViT-256}(x_{i,j})$
 Update memory bank: $\text{MemoryBank}[X_i][j] \leftarrow f_{i,j}$

2: Residual Encoding via VLAD

for each WSI X_i in batch do
 Retrieve features: $F_i \leftarrow \{\text{MemoryBank}[X_i][j] \mid j = 1, \dots, N\}$
 if codebook C uninitialized then $C \leftarrow \text{K-Means}(\bigcup_{i=1}^b F_i, K)$
 Compute residuals:
 for each $c_k \in C$ do
 $X_k \leftarrow \{f \in F_i \mid \arg \min_{c \in C} \|f - c\|_2 = k\}$
 $v_k \leftarrow \sum_{f \in X_k} (f - c_k)$
 Concatenate: $v_{\text{VLAD}} \leftarrow \text{Concat}([v_1, \dots, v_K])$
 Transform: $h_i \leftarrow \text{Transformer}(v_{\text{VLAD}})$

3: Contrastive Learning Alignment

for each WSI X_i with report T_i do
 Encode text: $t_i \leftarrow \text{LLaMA2-7B}(T_i)$
 Project image feature: $h'_i \leftarrow \text{Linear}(h_i)$
 Compute logits:
 $S_{\text{img2txt}} \leftarrow \tau \cdot h' \cdot t^\top \quad S_{\text{txt2img}} \leftarrow \tau \cdot t \cdot h'^\top$

Compute loss:

$$\mathcal{L}_{\text{contrastive}} \leftarrow \frac{1}{2} (\text{CE}(S_{\text{img2txt}}, y) + \text{CE}(S_{\text{txt2img}}, y))$$

4: Classification & Optimization

Predict label: $\hat{Y}_i \leftarrow \text{MLP}(h_i)$
 Compute loss: $\mathcal{L}_{\text{cls}} \leftarrow \text{CE}(\hat{Y}, Y)$
 Update parameters: $\theta \leftarrow \theta - \eta \nabla_{\theta} (\mathcal{L}_{\text{contrastive}} + \mathcal{L}_{\text{cls}})$

ensuring that the visual feature vectors and textual feature vectors both have the same dimensionality (4096 dimensions). This uniform feature dimensionality allows for direct comparison between the two modalities within the same feature space.

In contrastive learning loss calculation, we adopt an approach similar to CLIP, using cross-entropy loss to align visual and textual features. Specifically, the loss calculation involves two directions: visual-to-text and text-to-visual. During the forward pass, both visual and textual features are first linearly mapped and normalized. Then, the similarity between the two feature sets is measured by computing a similarity matrix (logits). Let the visual features be represented as the matrix $\mathbf{V} \in \mathbb{R}^{N \times d}$ and the textual features as the matrix $\mathbf{T} \in \mathbb{R}^{N \times d}$. The similarity matrices between them are calculated as follows:

$$S_{\text{str}} = \sigma_1 \times \mathbf{V} \mathbf{T}^T, \quad (5)$$

$$S_{\text{rts}} = \sigma_2 \times \mathbf{T} \mathbf{V}^T, \quad (6)$$

where σ_1 and σ_2 are learnable scalar parameters used to adjust the magnitude of the similarity. Initially, they are set to $\exp(\log(1/0.07))$, which helps enhance the model's sensitivity to feature differences in the early stages of training. As training progresses, they are dynamically updated to ensure the quality of feature alignment.

Subsequently, the model's ability to learn cross-modal matching relationships is assessed by calculating the cross-entropy loss for both slide-to-report and report-to-slide alignments. Let the labels be $y = [0, 1, \dots, N-1]$. The cross-entropy loss for slide-to-report and report-to-slide alignment is defined as:

$$\mathcal{L}_{\text{str}} = \frac{1}{N} \sum_{i=1}^N -\log \frac{\exp(S_{\text{str}}[i, y[i]])}{\sum_{j=1}^N \exp(S_{\text{str}}[i, j])}, \quad (7)$$

$$\mathcal{L}_{\text{rts}} = \frac{1}{N} \sum_{i=1}^N -\log \frac{\exp(S_{\text{rts}}[i, y[i]])}{\sum_{j=1}^N \exp(S_{\text{rts}}[i, j])}. \quad (8)$$

Finally, the contrastive loss is the average of them:

$$\mathcal{L}_{\text{contrastive}} = \frac{1}{2} (\mathcal{L}_{\text{str}} + \mathcal{L}_{\text{rts}}). \quad (9)$$

The contrastive learning approach ensures that visual and textual features are closely aligned in semantic space and effectively enhances the model's cross-modal retrieval capabilities and improves the representation of pathological features, capturing the intermodal relationships with higher accuracy.

3.4 Three distinct stages of the End-to-end Strategy

Our DRE-SLCL method mainly comprises three distinct stages: preparation, training, and testing to develop our end-to-end WSI representation framework.

The preparation stage: We process all raw WSI data from the datasets by segmenting each image into non-overlapping tiles of size 256×256 pixels. Then we utilize the tile encoder to extract initial features from all tiles, which are subsequently stored in the memory bank. K-means clustering is then applied to the extracted tile features to generate a codebook that captures the dominant patterns within the feature space.

3.3 Slide-Level Contrastive Learning

To improve generalization, we propose a slide-level contrastive learning method that aligns visual features with textual features. The visual features are represented by the residual encoding vector of the WSI, while the textual features are extracted from the corresponding pathological reports.

Specifically, textual features are extracted using the LLaMA2 7B model, generating 4096-dimensional vectors that capture the semantic content related to pathological conditions, diagnostic results, and other relevant information in pathological reports. The extracted textual feature vectors are then L2-normalized to ensure consistent norms across the feature space. To align the visual features with the textual features, we apply a linear mapping to the visual features, embedding them into the same feature space as the textual features. This mapping is achieved through a linear layer,

The training stage: In each iteration, from each WSI in a batch of 64, we randomly sample 10 image tiles - an optimal quantity determined by our ablation study in Section 4.5. This random sampling strategy serves as data augmentation while reducing computational overhead, ensuring diverse tile combinations across training iterations. As the parameter of tile encoder is updated in each epoch, the tile features of the mini-batch are also updated. We use the updated tile features to replace the old tile features stored in the memory bank. Due to its nested two-level dictionary design, we can quickly locate which WSIs are updated for those tiles and complete this update process in less time. Then, we used both the updated tiles and other original tiles from the memory bank for the final WSI representation using the VLAD method. Subsequently, an updated representation is then used for classification and to compute the contrastive loss.

In particular, to enable effective end-to-end training, we adopt a two-stage training strategy. In the first stage, we freeze the parameters of the ViT-256 model and train only the classifier, stabilizing the process and accelerating convergence. Once this baseline is set, we gradually unfreeze the ViT-256 model parameters, enabling joint fine-tuning of both the feature extractor and classifier. This staged unfreezing and end-to-end training optimizes the synergy between feature extraction and classification, leading to more refined feature representations and improved classification accuracy.

The testing stage: The trained feature extraction model and codebook are used to produce an end-to-end representation for each WSI. The feature extractor segments the WSI into tiles, extracting features $f(x_i)$ for each tile. And then these features are going to a robust global representation of the WSI with the VLAD module and a transformer layer. This global representation is subsequently processed by a classifier to produce the final classification result.

4 Experiments

In this section, we demonstrate the effectiveness of the proposed DRE-SLCL method on the TCGA and CPTAC datasets. The experiments focus on three primary tasks: (1) TCGA_LUNG and CPTAC_LUNG subtype classification, (2) binary classification for the presence of cancer in TCGA_LUAD, TCGA_LUSC, CPTAC_CCRCC and CPTAC_PDA, and (3) prediction of four major gene mutation types in TCGA_LUAD. Additionally, we conduct an ablation study to assess the contribution of each module and evaluate the generalization ability of the proposed method. The following paragraphs provide detailed information on the experimental datasets, evaluation metrics, and parameter settings.

Datasets and Evaluation Metrics. We evaluated the proposed DRE-SLCL method on two publicly available lung cancer datasets from the TCGA project: LUAD (Lung Adenocarcinoma) and LUSC (Lung Squamous Cell Carcinoma) and four cancer datasets from the CPTAC project: LUAD (Lung Adenocarcinoma), LSCC (Lung Squamous Cell Carcinoma), CCRCC (Clear Cell Renal Cell Carcinoma), PDA (Pancreatic Ductal Adenocarcinoma). For the gene mutation prediction task, we focus on predicting TP53, FAT1, LPRB1, and EGFR mutations. The evaluation metrics include the area under the curve (AUC), and the weighted F1 score. For cancer recognition, we

Table 2: The comparison of AUC (%) and F1 scores (%) for cancer subtyping on the TCGA and CPTAC lung datasets.

Methods	TCGA_lung		CPTAC_lung	
	AUC	F1-score	AUC	F1-score
CLAM-SB[31]	85.59	76.05	78.99	68.65
CLAM-MB[31]	83.29	75.78	70.27	57.50
ABMIL[19]	81.26	73.98	76.35	64.46
MIL[31]	75.55	35.85	57.02	47.12
TransMIL[36]	77.62	66.15	75.95	62.32
HIPT[6]	83.51	75.75	78.35	69.28
Prov-GigaPath[58]	81.32	72.45	75.75	63.57
C2C [37]	78.56	70.49	75.53	66.83
DGR-MIL[64]	84.06	77.43	78.76	71.24
DRE-SLCL(w/o cliploss)	83.12	79.63	78.89	74.35
ABMIL[19]+cliploss	84.48	75.56	77.43	67.29
DRE-SLCL(ours)	83.76	80.88	80.28	76.43

report the AUC score. For cancer subtyping and the gene mutation prediction task, we report all of the aforementioned metrics.

In particular, not all WSIs have corresponding pathology reports for contrastive learning. To fully exploit the available textual information, all WSIs with pathology reports are included in the training set for contrastive learning with visual features.

Experimental Setting. The experiments were conducted on a single Nvidia A6000 GPU, with the proposed model implemented in PyTorch. In each experiment, we utilized the ViT-256 model pre-trained in the HIPT framework to extract tile-level features. The model was trained using Adam optimizer with a learning rate of 10^{-4} and a weight decay of 10^{-5} , for a total of 100 epochs.

The training batch size was set to 64, with testing performed using batch size 1. We randomly sampled 10 tiles per WSI in each iteration, and set the codebook size k to 64 clusters. These parameters were optimized through ablation studies (Section 4.5) and all the implementation details follow the three-stage strategy described in Section 3.4.

4.1 Cancer Subtyping on the TCGA_LUNG and the CPTAC_LUNG Datasets

In this task, we performed the binary classification of LUAD and LUSC using TCGA_lung and CPTAC_lung datasets. Compared to traditional block-based classification methods, the proposed approach achieves significant improvements for the weighted F1 score, as shown in Table 2. Specifically, our model improved by approximately 4.83% for the weighted F1 score compared to the best competing methods, demonstrating its effectiveness in comprehensively modeling global features.

4.2 Cancer Recognition on TCGA and CPTAC Datasets

For the binary cancer recognition task on the LUAD, LUSC, CCRCC, and PDA datasets, we used the same architecture used in the previous subsection to encode tile features and generate global representations. The experimental results are shown in Table 3. We can see that our method outperforms mainstream approaches for the

Table 3: The comparison of AUC (%) scores for cancer recognition on the TCGA_LUAD,TCGA_LUSC,CPTAC_CCRCC and CPTAC_PDA datasets.

Methods	TCGA dataset		CPTAC dataset	
	LUAD	LUSC	CCRCC	PDA
CLAM-SB[31]	99.00	99.23	99.28	83.71
CLAM-MB[31]	98.94	98.66	98.75	86.32
ABMIL[19]	97.73	98.79	98.54	85.58
MIL[31]	98.13	98.55	98.04	82.81
TransMIL[36]	98.41	97.09	99.23	82.74
HIPT [6]	98.08	99.28	99.04	85.45
Prov-GigaPath[58]	97.07	96.81	98.23	82.45
C2C [37]	97.22	98.98	98.87	89.51
DGR-MIL[64]	98.89	98.76	99.12	87.63
DRE-SLCL(w/o cliploss)	99.02	98.43	98.75	90.87
ABMIL[19]+cliploss	98.29	98.57	99.05	86.23
DRE-SLCL(ours)	99.26	99.58	99.53	92.34

AUC score, with particular robustness in distinguishing edge cases. Using dynamic residual encoding, the model dynamically adjusts the importance of different image tiles during training, leading to more accurate detection of cancer status.

4.3 Gene Mutation Prediction on The LUAD Dataset

To further validate the effectiveness of our model, we performed classification predictions for four key gene mutation types on the LUAD dataset. These mutations are commonly observed in lung adenocarcinoma, and the goal was to predict their presence. The results are shown in Table 4. The proposed method demonstrated superior performance in terms of AUC and F1 scores, with significant improvements over other approaches, particularly excelling in handling small sample sizes and addressing label imbalance.

4.4 Comparison with State-of-the-Art Methods

The proposed method was systematically compared with several state-of-the-art approaches, including CLAM-SB [31], CLAM-MB [31], ABMIL [19], TansMIL [36], HIPT [6], MIL [31], Prov-GigaPath [58], DGR-MIL[64] and one of end-to-end WSI classification methods C2C [37].

The comparison results are summarized in Tables 2, 3, and 4, which present performance across three major tasks: cancer subtyping, cancer presence classification, and mutation prediction in LUAD. We can see that the proposed DRE-SLCL method outperforms traditional MIL methods. DRE-SLCL excels at capturing interactions between image tiles and generating a comprehensive global feature representation. In both cancer recognition and mutation prediction tasks, DRE-SLCL shows superior performance and effectively handles complex biomedical imaging data. Notably, in the LUAD gene mutation prediction task, DRE-SLCL exhibits strong robustness in addressing data imbalance and small sample sizes. The improvements observed in AUC scores and F1-scores validate the effectiveness of our approach over other state-of-the-art methods.

4.5 Analysis of Contrastive Loss Contribution

To evaluate the contribution of contrastive loss in our model, we conducted additional experiments as shown in Tables 2, 3, and 4. We implemented a version without contrastive loss (DRE-SLCL w/o cliploss), which still performed well across all tasks, outperforming most baseline methods. For TP53 prediction, this version achieved 68.74% AUC and 79.43% F1-score, while the complete DRE-SLCL showed further improvements (71.33% AUC, 82.76% F1-score). To verify that performance improvements weren't simply due to contrastive loss, we integrated it into ABMIL [19] (ABMIL [19]+cliploss). And the results remained below our method. These experiments demonstrate that DRE-SLCL's superior performance stems from the synergistic combination of our dynamic residual encoding architecture and contrastive learning strategy, rather than solely from the contrastive loss component. This finding was consistent across all three experimental tasks.

4.6 Analysis of Computational Efficiency and Model Complexity

The proposed framework demonstrates computational efficiency through a lightweight architectural design. The core backbone employs a pretrained ViT-256 comprising approximately 22M parameters, while the VLAD encoding module and Transformer aggregator contribute fewer than 5M parameters, a total model size of under 27M parameters—significantly smaller than recent foundation model-based approaches such as Prov-GigaPath (1.1B parameters).

The dynamic tile sampling mechanism enables training on a representative subset of tiles per WSI ($r=10$) rather than processing all tiles in each epoch, substantially reducing computational burden and allowing single training epoch completion in approximately 5–10 minutes on a single NVIDIA A6000 GPU. During inference, while all tiles are processed to ensure complete information utilization (typically requiring several minutes per WSI), the end-to-end design eliminates redundant I/O and preprocessing steps compared to traditional two-stage pipelines, and the memory bank stores intermediate embeddings in CPU memory without consuming GPU resources. This approach achieves an optimal balance between model capacity and computational cost, making it well-suited for scalable deployment in resource-constrained clinical environments.

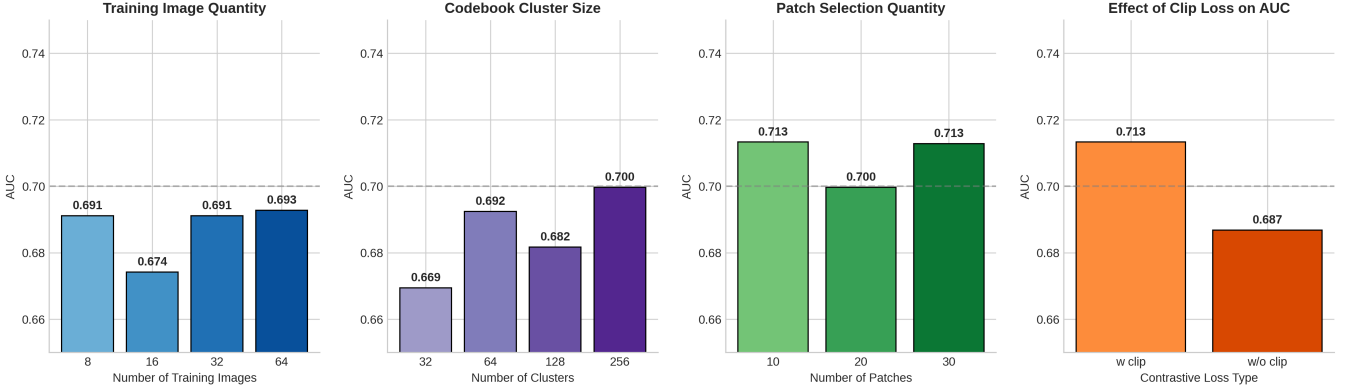
4.7 Ablation Study

To evaluate the contribution of different components of our method, we conducted an ablation study using the gene mutation prediction task for TP53. The study focused on several aspects: codebook size, batch size, number of tiles selected per WSI, inclusion of contrastive loss, and whether the model was trained end-to-end. The results are visualized in Figure 3.

Effect of the batch size: The impact of different batch sizes (8, 16, 32, 64) on AUC is presented in the first subfigure of fig. 3. The batch size was chosen based on contrastive learning theory (larger batches provide more negative samples), dataset scale, and hardware constraints (NVIDIA A6000 with 48GB VRAM). And the result shows batch size 64 achieves a balance between computational feasibility and model performance with the highest AUC.

Table 4: The comparison of gene mutation scores (%) on the LUAD dataset.

Methods	TP53		ZGPR		LPRB1		FAT1	
	AUC	F1-score	AUC	F1-score	AUC	F1-score	AUC	F1-score
CLAM-SB[31]	61.41	79.10	56.73	70.11	57.46	58.25	49.52	16.39
CLAM-MB[31]	65.33	76.19	57.55	63.20	59.47	52.53	49.75	47.76
ABMIL[19]	64.17	75.47	56.48	66.37	57.25	50.25	50.17	45.56
MIL[31]	57.36	65.34	56.34	19.51	51.28	52.40	45.29	0
TransMIL[36]	60.85	78.18	54.09	56.52	50.39	51.71	53.40	53.40
HIPT[6]	67.61	81.88	61.15	69.92	57.27	48.94	50.56	19.83
Prov-GigaPath[58]	61.83	65.46	56.30	63.85	52.88	52.53	50.67	51.01
C2C[37]	69.23	79.40	52.74	80.00	55.54	50.53	50.68	0
DGR-MIL[64]	69.12	78.89	58.11	78.43	59.36	52.49	52.38	50.69
DRE-SLCL(w/o cliploss)	68.74	79.43	56.21	73.51	58.76	49.63	53.41	53.86
ABMIL[19]+cliploss	65.82	77.24	55.63	67.56	52.23	53.34	52.17	48.79
DRE-SLCL(ours)	71.33	82.76	57.13	80.80	60.17	50.67	54.94	54.17

**Figure 3: Ablation results for TP53 gene mutation prediction.**

Effect of codebook size: The effect of varying the number of clusters in the codebook (32, 64, 128, 256) is shown in the second subfigure. The choice follows the classical VLAD encoding framework where cluster centers typically range from 32 to 256. $K=256$ is a widely adopted setting that balances descriptive power and computational efficiency. Increasing the size of the cluster generally improves the AUC, with the best appears at 256 clusters, indicating that a larger codebook can effectively capture nuanced features.

Effect of the number of tiles: The third subfigure in Figure 3 shows the effect of selecting different numbers of patches per WSI (10, 20, 30) on AUC. The tile sampling follows MIL principles, where sampling a representative subset ($r=10$) reduces computational complexity from $O(N \times d)$ to $O(r \times d)$, where N is the total number of tiles and d is the feature dimension. In MIL, not all instances contribute equally to the final prediction. In fact, excessive sampling may introduce noise or irrelevant patterns. And the results also shows that both 10 and 30 patches better performance compared to 20.

Based on the ablation study, we selected the optimal parameters: batch size of 64, codebook size of 256 clusters, 10 patches per WSI,

and inclusion of contrastive loss. This combination consistently provided the best overall performance, demonstrating its effectiveness in capturing the complex features of the whole slide images.

5 Conclusion

In this work, we presented a novel method called Dynamic Residual Encoding with Slide-Level Contrastive Learning (DRE-SLCL) for end-to-end Whole Slide Image (WSI) representation. Our approach integrates dynamic residual encoding and slide-level contrastive learning into a unified framework, effectively addressing the limitations of traditional two-stage learning processes. By incorporating a memory bank and leveraging end-to-end training, our method can manage computational complexity while enhancing feature extraction, resulting in a more robust WSI representation. The experimental results on the TCGA and CPTAC datasets demonstrate the superiority of DRE-SLCL in cancer subtyping, cancer recognition, and gene mutation prediction tasks. Future work will explore further improvements to the interpretability of the model and extend its application to other types of cancer, aiming to establish a more comprehensive tool for computational pathology.

Acknowledgements

This work was supported by the National Natural Science Foundation of China (U24A20256, 62473385, and 62202499). We are grateful to the High Performance Computing Center of Central South University for partial support of this work.

References

- [1] Wenya Linda Bi, Ahmed Hosny, Matthew B Schabath, Maryellen L Giger, Nicolai J Birkbak, Alireza Mehrtaash, Tavis Allison, Omar Arnaout, Christopher Abbosh, Ian F Dunn, et al. 2019. Artificial intelligence in cancer imaging: clinical challenges and applications. *CA: a cancer journal for clinicians* 69, 2 (2019), 127–157.
- [2] Hao Bian, Zhuchen Shao, Yang Chen, Yifeng Wang, Haoqian Wang, Jian Zhang, and Yongbing Zhang. 2022. Multiple instance learning with mixed supervision in gleason grading. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 204–213.
- [3] Wouter Bulten, Kimmo Kartasalo, Po-Hsuan Cameron Chen, Peter Ström, Hans Pinckaers, Kunal Nagpal, Yuannan Cai, David F Steiner, Hester Van Boven, Robert Vink, et al. 2022. Artificial intelligence for diagnosis and Gleason grading of prostate cancer: the PANDA challenge. *Nature medicine* 28, 1 (2022), 154–163.
- [4] 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. 2019. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nature medicine* 25, 8 (2019), 1301–1309.
- [5] 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. 2019. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nature medicine* 25, 8 (2019), 1301–1309.
- [6] Richard J Chen, Chengkuan Chen, Yicong Li, Tiffany Y Chen, Andrew D Trister, Rahul G Krishnan, and Faisal Mahmood. 2022. Scaling vision transformers to gigapixel images via hierarchical self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16144–16155.
- [7] Richard J Chen and Rahul G Krishnan. 2022. Self-supervised vision transformers learn visual concepts in histopathology. *arXiv preprint arXiv:2203.00585* (2022).
- [8] Richard J Chen, Ming Y Lu, Drew FK Williamson, Tiffany Y Chen, Jana Lipkova, Zahra Noor, Muhammad Shaban, Maha Shady, Mane Williams, Bumjin Joo, et al. 2022. Pan-cancer integrative histology-genomic analysis via multimodal deep learning. *Cancer Cell* 40, 8 (2022), 865–878.
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PMLR, 1597–1607.
- [10] Yuqi Chen, Juan Liu, Zhiqun Zuo, Peng Jiang, Yu Jin, and Guangsheng Wu. 2023. Classifying pathological images based on multi-instance learning and end-to-end attention pooling. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 1–5.
- [11] Philip Chikontwe, Meejeong Kim, Soo Jeong Nam, Heounjeong Go, and Sang Hyun Park. 2020. Multiple instance learning with center embeddings for histopathology classification. In *Medical Image Computing and Computer Assisted Intervention-MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part V* 23. Springer, 519–528.
- [12] Yufei Cui, Ziquan Liu, Yixin Chen, Yuchen Lu, Xinyue Yu, Xue Steve Liu, Tei-Wei Kuo, Miguel Rodrigues, Chun Jason Xue, and Antoni Chan. 2023. Retrieval-augmented multiple instance learning. *Advances in Neural Information Processing Systems* 36 (2023), 24859–24878.
- [13] Olivier Elemento, Christina Leslie, Johan Lundin, and Georgia Tourassi. 2021. Artificial intelligence in cancer research, diagnosis and therapy. *Nature Reviews Cancer* 21, 12 (2021), 747–752.
- [14] 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. 2020. Pan-cancer computational histopathology reveals mutations, tumor composition and prognosis. *Nature cancer* 1, 8 (2020), 800–810.
- [15] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 9729–9738.
- [16] Simon Holdenried-Krafft, Peter Somers, Ivonne A Montes-Majarro, Diana Silimon, Cristina Tarin, Falko Fend, and Hendrik Lensch. 2023. Dual-Query Multiple Instance Learning for Dynamic Meta-Embedding based Tumor Classification. *arXiv preprint arXiv:2307.07482* (2023).
- [17] Danqing Hu, Huanyao Zhang, Shaolei Li, Yuhong Wang, Nan Wu, and Xudong Lu. 2021. Automatic extraction of lung cancer staging information from computed tomography reports: deep learning approach. *JMIR medical informatics* 9, 7 (2021), e27955.
- [18] Maximilian Ilse, Jakub Tomczak, and Max Welling. 2018. Attention-based deep multiple instance learning. In *International conference on machine learning*. PMLR, 2127–2136.
- [19] Maximilian Ilse, Jakub Tomczak, and Max Welling. 2018. Attention-based deep multiple instance learning. In *International conference on machine learning*. PMLR, 2127–2136.
- [20] Sajid Javed, Arif Mahmood, Iyyakutti Iyappan Ganapathi, Fayaz Ali Dharejo, Naoufel Werghi, and Mohammed Bennamoun. 2024. Cclip: Zero-shot learning for histopathology with comprehensive vision-language alignment. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 11450–11459.
- [21] Hervé Jégou, Matthijs Douze, Cordelia Schmid, and Patrick Pérez. 2010. Aggregating local descriptors into a compact image representation. In *2010 IEEE computer society conference on computer vision and pattern recognition*. IEEE, 3304–3311.
- [22] Jakob Nikolas Kather, Lara R Heij, Heike I Grabsch, Chiara Loeffler, Amelie Echle, Hannah Sophie Muti, Jeremias Krause, Jan M Niehues, Kai AJ Sommer, Peter Bankhead, et al. 2020. Pan-cancer image-based detection of clinically actionable genetic alterations. *Nature cancer* 1, 8 (2020), 789–799.
- [23] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020. Supervised contrastive learning. *Advances in neural information processing systems* 33 (2020), 18661–18673.
- [24] Andreas Kleppe, Ole-Johan Skrede, Sepp De Raedt, Knut Liestøl, David J Kerr, and Håvard E Danielsen. 2021. Designing deep learning studies in cancer diagnostics. *Nature Reviews Cancer* 21, 3 (2021), 199–211.
- [25] Bin Li, Yin Li, and Kevin W Eliceiri. 2021. Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 14318–14328.
- [26] Bin Li, Yin Li, and Kevin W Eliceiri. 2021. Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 14318–14328.
- [27] Jiawen Li, Yuxuan Chen, Hongbo Chu, Qiehe Sun, Tian Guan, Anjia Han, and Yonghong He. 2024. Dynamic graph representation with knowledge-aware attention for histopathology whole slide image analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11323–11332.
- [28] Ruoyu Li, Jiawen Yao, Xinliang Zhu, Yeqing Li, and Junzhou Huang. 2018. Graph CNN for survival analysis on whole slide pathological images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 174–182.
- [29] Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen AWM Van Der Laak, Bram Van Ginneken, and Clara I Sánchez. 2017. A survey on deep learning in medical image analysis. *Medical image analysis* 42 (2017), 60–88.
- [30] Ming Y Lu, Tiffany Y Chen, Drew FK Williamson, Melissa Zhao, Maha Shady, Jana Lipkova, and Faisal Mahmood. 2021. AI-based pathology predicts origins for cancers of unknown primary. *Nature* 594, 7861 (2021), 106–110.
- [31] Ming Y Lu, Drew FK Williamson, Tiffany Y Chen, Richard J Chen, Matteo Barbieri, and Faisal Mahmood. 2021. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature biomedical engineering* 5, 6 (2021), 555–570.
- [32] Yingfan Ma, Xiaoyuan Luo, Kexue Fu, and Manning Wang. 2024. Transformer-based video-structure multi-instance learning for whole slide image classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 14263–14271.
- [33] Yangling Ma, Yixin Luo, and Zhouwang Yang. 2024. GCN-based MIL: multi-instance learning utilizing structural relationships among instances. *Signal, Image and Video Processing* 18, 6 (2024), 5549–5561.
- [34] Pushpak Pati, Guillaume Jaume, Lauren Alisha Fernandes, Antonio Foncubierta-Rodríguez, Florinda Feroce, Anna Maria Anniciello, Giosue Scognamiglio, Nadia Brancati, Daniel Riccio, Maurizio Di Bonito, et al. 2020. Hact-net: A hierarchical cell-to-tissue graph neural network for histopathological image classification. In *Uncertainty for Safe Utilization of Machine Learning in Medical Imaging, and Graphs in Biomedical Image Analysis: Second International Workshop, UNSURE 2020, and Third International Workshop, GRAIL 2020, Held in Conjunction with MICCAI 2020, Lima, Peru, October 8, 2020, Proceedings 2*. Springer, 208–219.
- [35] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [36] Zhuchen Shao, Hao Bian, Yang Chen, Yifeng Wang, Jian Zhang, Xiangyang Ji, et al. 2021. Transmil: Transformer based correlated multiple instance learning for whole slide image classification. *Advances in neural information processing systems* 34 (2021), 2136–2147.
- [37] Yash Sharma, Aman Shrivastava, Lubaina Ehsan, Christopher A Moskaluk, Sana Syed, and Donald Brown. 2021. Cluster-to-conquer: A framework for end-to-end multi-instance learning for whole slide image classification. In *Medical imaging with deep learning*. PMLR, 682–698.
- [38] Li Shen, Laurie R Margolies, Joseph H Rothstein, Eugene Fluder, Russell McBride, and Weiva Sieh. 2019. Deep learning to improve breast cancer detection on

- screening mammography. *Scientific reports* 9, 1 (2019), 12495.
- [39] Tong Shu, Jun Shi, Dongdong Sun, Zhiguo Jiang, and Yushan Zheng. 2024. SlideGCD: Slide-Based Graph Collaborative Training with Knowledge Distillation for Whole Slide Image Classification. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 470–480.
- [40] Musalula Sinkala, Nicola Mulder, and Darren Martin. 2020. Machine learning and network analyses reveal disease subtypes of pancreatic cancer and their molecular characteristics. *Scientific reports* 10, 1 (2020), 1212.
- [41] Shusuke Takahama, Yusuke Kurose, Yusuke Mukuta, Hiroyuki Abe, Masashi Fukayama, Akihiko Yoshizawa, Masanobu Kitagawa, and Tatsuya Harada. 2019. Multi-stage pathological image classification using semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10702–10711.
- [42] Thomas E Tavolara, Metin N Gurcan, and M Khalid Khan Niazi. 2022. Contrastive multiple instance learning: An unsupervised framework for learning slide-level representations of whole slide histopathology images without labels. *Cancers* 14, 23 (2022), 5778.
- [43] Eric J Topol. 2019. High-performance medicine: the convergence of human and artificial intelligence. *Nature medicine* 25, 1 (2019), 44–56.
- [44] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).
- [45] Guilan Tu, Wuchao Li, Yongshun Lin, Zi Xu, Junjie He, Bangkang Fu, Ping Huang, Rongpin Wang, and Yunsong Peng. 2025. GAMMIL: A graph attention-guided multi-scale fusion multiple instance learning model for the WHO grading of meningioma in whole slide images. *Biomedical Signal Processing and Control* 105 (2025), 107652.
- [46] Jeroen van der Laak, Francesco Ciompi, and Geert Litjens. 2019. No pixel-level annotations needed. *Nature biomedical engineering* 3, 11 (2019), 855–856.
- [47] Dayong Wang, Aditya Khosla, Rishab Gargya, Humayun Irshad, and Andrew H Beck. 2016. Deep learning for identifying metastatic breast cancer. *arXiv preprint arXiv:1606.05718* (2016).
- [48] Xiyue Wang, Yuexi Du, Sen Yang, Jun Zhang, Minghui Wang, Jing Zhang, Wei Yang, Junzhou Huang, and Xiao Han. 2023. RetCCL: Clustering-guided contrastive learning for whole-slide image retrieval. *Medical image analysis* 83 (2023), 102645.
- [49] Xiyue Wang, Yuexi Du, Sen Yang, Jun Zhang, Minghui Wang, Jing Zhang, Wei Yang, Junzhou Huang, and Xiao Han. 2023. RetCCL: Clustering-guided contrastive learning for whole-slide image retrieval. *Medical image analysis* 83 (2023), 102645.
- [50] Xiyue Wang, Jinxi Xiang, Jun Zhang, Sen Yang, Zhongyi Yang, Ming-Hui Wang, Jing Zhang, Wei Yang, Junzhou Huang, and Xiao Han. 2022. Scl-wc: Cross-slide contrastive learning for weakly-supervised whole-slide image classification. *Advances in neural information processing systems* 35 (2022), 18009–18021.
- [51] Xiyue Wang, Jinxi Xiang, Jun Zhang, Sen Yang, Zhongyi Yang, Ming-Hui Wang, Jing Zhang, Wei Yang, Junzhou Huang, and Xiao Han. 2022. Scl-wc: Cross-slide contrastive learning for weakly-supervised whole-slide image classification. *Advances in neural information processing systems* 35 (2022), 18009–18021.
- [52] Xiyue Wang, Sen Yang, Jun Zhang, Minghui Wang, Jing Zhang, Wei Yang, Junzhou Huang, and Xiao Han. 2022. Transformer-based unsupervised contrastive learning for histopathological image classification. *Medical image analysis* 81 (2022), 102559.
- [53] Zhikang Wang, Yue Bi, Tong Pan, Xiaoyu Wang, Chris Bain, Richard Basset, Seiya Imoto, Jianhua Yao, Roger J Daly, and Jiangning Song. 2023. Targeting tumor heterogeneity: multiplex-detection-based multiple instance learning for whole slide image classification. *Bioinformatics* 39, 3 (2023), btad114.
- [54] Hangchen Xiang, Junyi Shen, Qingguo Yan, Meilian Xu, Xiaoshuang Shi, and Xiaofeng Zhu. 2023. Multi-scale representation attention based deep multiple instance learning for gigapixel whole slide image analysis. *Medical Image Analysis* 89 (2023), 102890.
- [55] Jinxi Xiang and Jun Zhang. 2023. Exploring low-rank property in multiple instance learning for whole slide image classification. In *The Eleventh International Conference on Learning Representations*.
- [56] Jinxi Xiang and Jun Zhang. 2023. Exploring low-rank property in multiple instance learning for whole slide image classification. In *The Eleventh International Conference on Learning Representations*.
- [57] Chensu Xie, Hassan Muhammad, Chad M Vanderbilt, Raul Caso, Dig Vijay Kumar Yarlagadda, Gabriele Campanella, and Thomas J Fuchs. 2020. Beyond Classification: Whole Slide Tissue Histopathology Analysis By End-To-End Part Learning.. In *MIDL*. 843–856.
- [58] Hanwen Xu, Naoto Usuyama, Jaspreet Bagga, Sheng Zhang, Rajesh Rao, Tristan Naumann, Cliff Wong, Zelalem Gero, Javier González, Yu Gu, et al. 2024. A whole-slide foundation model for digital pathology from real-world data. *Nature* (2024), 1–8.
- [59] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. 2020. Graph contrastive learning with augmentations. *Advances in neural information processing systems* 33 (2020), 5812–5823.
- [60] Suwan Yu, Yooeun Kim, Hoeyoung Kim, Sangseon Lee, and Kwangsoo Kim. 2025. Contrastive Learning for Omics-guided Whole-slide Visual Embedding Representation. *bioRxiv* (2025), 2025–01.
- [61] Hongrun Zhang, Yanda Meng, Yitian Zhao, Yihong Qiao, Xiaoyun Yang, Sarah E Coupland, and Yalin Zheng. 2022. Dtfd-mil: Double-tier feature distillation multiple instance learning for histopathology whole slide image classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 18802–18812.
- [62] Lu Zhao, Wangyuan Zhao, Lu Qiu, Mengqi Jiang, Liqiang Qian, Hua-Nong Ting, Xiaolong Fu, Puming Zhang, Yuchen Han, and Jun Zhao. 2025. CoLM: Contrastive learning and multiple instance learning network for lung cancer classification of surgical options based on frozen pathological images. *Biomedical Signal Processing and Control* 100 (2025), 107097.
- [63] Xiao Zhou, Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Weidi Xie, and Yanfeng Wang. 2024. Knowledge-enhanced visual-language pretraining for computational pathology. In *European Conference on Computer Vision*. Springer, 345–362.
- [64] Wenhui Zhu, Xiwen Chen, Peijie Qiu, Aristeidis Sotiras, Abolfazl Razi, and Yalin Wang. 2024. Dgr-mil: Exploring diverse global representation in multiple instance learning for whole slide image classification. In *European conference on computer vision*. Springer, 333–351.