

Benchmarking multiple instance learning architectures from patches to pathology for prostate cancer detection and grading using attention-based weak supervision

Received: 24 August 2025

Accepted: 3 February 2026

Published online: 02 March 2026

Cite this article as: Butt N.A., Sarwat D., Noya I.D. *et al.* Benchmarking multiple instance learning architectures from patches to pathology for prostate cancer detection and grading using attention-based weak supervision. *Sci Rep* (2026). <https://doi.org/10.1038/s41598-026-39196-x>

Naveed Anwer Butt, Dilawaiz Sarwat, Irene Delgado Noya, Kilian Tutusaus, Nagwan Abdel Samee & Imran Ashraf

We are providing an unedited version of this manuscript to give early access to its findings. Before final publication, the manuscript will undergo further editing. Please note there may be errors present which affect the content, and all legal disclaimers apply.

If this paper is publishing under a Transparent Peer Review model then Peer Review reports will publish with the final article.

Benchmarking Multiple Instance Learning Architectures from Patches to Pathology for Prostate Cancer Detection and Grading Using Attention-Based Weak Supervision

Naveed Anwer Butt¹, Dilawaiz Sarwat^{1*}, Irene Delgado Noya^{2,3,4,5},
Kilian Tutusaus^{2,6,7}, Nagwan Abdel Samee⁸, Imran Ashraf^{9*}

¹Department of Computer Science, University of Gujrat, Gujrat, Pakistan.

²Universidad Europea del Atlantico, Santander, 39011, Spain.

³Universidad Internacional Iberoamericana, Campeche, 24560, Mexico.

⁴Fundacion Universitaria Internacional de Colombia, Bogota, Colombia.

⁵Universidad de La Romana, La Romana, Republica Dominicana.

⁶Universidad Internacional Iberoamericana Arecibo, Puerto Rico, 00613, USA.

⁷Universidade Internacional do Cuanza, Cuito, Bie., Angola.

⁸Department of Information Technology, College of Computer and Information Sciences, Princess Nourah bint Abdulrahman University, P.O. Box 84428, Riyadh, 11671, Saudi Arabia.

⁹School of Computer Science and Engineering, Yeungnam University, Gyeongsan-si, 38541, Republic of Korea.

*Corresponding author(s). E-mail(s): 23015919-007@uog.edu.pk;
imranashraf@yu.ac.kr;

Contributing authors: naveed@uog.edu.pk;
Irene.delgado@uneatlantico.es; kilian.tutusaus@uneatlantico.es ;
mmabdelsamee@pmu.edu.sa;

Abstract

Histopathological evaluation is necessary for the diagnosis and grading of prostate cancer, which is still one of the most common cancers in men globally. Traditional

evaluation is time-consuming, prone to inter-observer variability, and challenging to scale. The clinical usefulness of current AI systems is limited by the need for comprehensive pixel-level annotations. The objective of this research is to develop and evaluate a large-scale benchmarking study on a weakly supervised deep learning framework that minimizes the need for annotation and ensures interpretability for automated prostate cancer diagnosis and International Society of Urological Pathology (ISUP) grading using whole slide images (WSIs). This study rigorously tested six cutting-edge multiple instance learning (MIL) architectures (CLAM-MB, CLAM-SB, ILRA-MIL, AC-MIL, AMD-MIL, WiKG-MIL), three feature encoders (ResNet50, CTransPath, UNI2), and four patch extraction techniques (varying sizes and overlap) using the PANDA dataset (10,616 WSIs), yielding 72 experimental configurations. The methodology used distributed cloud computing to process over 31 million tissue patches, implementing advanced attention mechanisms to ensure clinical interpretability through Grad-CAM visualizations. The optimum configuration (UNI2 encoder with ILRA-MIL, 256×256 patches, 50% overlap) achieved 78.75% accuracy and 90.12% quadratic weighted kappa (QWK), outperforming traditional methods and approaching expert pathologist-level diagnostic capability. Overlapping smaller patches offered the best balance of spatial resolution and contextual information, while domain-specific foundation models performed noticeably better than generic encoders. This work is the first large-scale, comprehensive comparison of weakly supervised MIL methods for prostate cancer diagnosis and grading. The proposed approach has excellent clinical diagnostic performance, scalability, practical feasibility through cloud computing, and interpretability using visualization tools.

Keywords: Prostate cancer detection; weakly supervised learning; multiple instance learning; whole slide images; ISUP grading

1 Introduction

Prostate cancer remains the second most common malignancy among men worldwide and represents a leading cause of cancer-related mortality [1]. The diagnosis and grading of prostate cancer primarily depend on histopathological examination of tissue samples, with the ISUP grading system serving as the gold standard for determining disease severity and guiding treatment decisions [2]. However, conventional analysis of histopathological images by pathologists is time-consuming, labor-intensive, and subject to considerable inter-observer variability, especially in borderline cases [3]. Recent advances in artificial intelligence (AI) have demonstrated significant potential in transforming cancer diagnosis and pathology workflows [4].

The emergence of computational pathology and deep learning (DL) has revolutionized medical image analysis, particularly in cancer detection and grading applications [5]. Foundation models and self-supervised learning approaches have shown remarkable success in analyzing whole slide images (WSIs), offering new possibilities for automated diagnosis and clinical decision support workflows [4]. Despite these advances, most current AI solutions require extensive manual annotations, which creates practical barriers for real-world clinical implementation [11]. The development of weakly

18 supervised learning frameworks has emerged as a promising solution to address these
19 annotation challenges while maintaining diagnostic accuracy.

20 Digital pathology has gained significant attention due to its potential in uro-
21 logical cancer diagnosis, with several studies demonstrating AI systems that can
22 achieve pathologist-level performance in specific diagnostic tasks [7]. The integra-
23 tion of multiple instance learning (MIL) approaches with attention mechanisms has
24 shown particular promise for WSI analysis, enabling effective learning from slide-
25 level labels without requiring detailed pixel-wise annotations [8]. Recent systematic
26 reviews highlight the growing evidence supporting AI-driven approaches in prostate
27 cancer diagnosis, emphasizing both the opportunities and challenges in clinical
28 implementation [9].

29 The International Society of Urological Pathology (ISUP) grading system, which
30 translates Gleason scores into standardized grade groups (ISUP grades 1-5), has
31 become the preferred clinical reporting standard for prostate cancer diagnosis and
32 treatment planning [10]. This standardized grading system provides clearer prognostic
33 information and treatment guidance compared to traditional Gleason scoring alone.
34 Modern AI systems for prostate cancer diagnosis increasingly focus on ISUP grade
35 prediction as it directly correlates with clinical decision-making and patient manage-
36 ment protocols, making it the most clinically relevant target for automated grading
37 systems.

38 The pressing need to use scalable and effective AI technologies to change prostate
39 cancer diagnosis is what drives our endeavor. Proper prostate cancer grading has a
40 direct impact on patient outcomes and treatment choices. There is a chance to improve
41 diagnostic precision, lower inter-observer variability, and expedite pathology processes
42 by automating this procedure with weakly supervised deep learning. The potential to
43 leverage cutting-edge artificial intelligence technologies to address these long-standing
44 clinical problems provides compelling justification for developing innovative diagnostic
45 solutions.

46 This study aims to develop an efficient and accurate weakly supervised deep learn-
47 ing framework for automated prostate cancer detection and ISUP grading from WSIs.
48 This research addresses the critical need for practical AI solutions that can be imple-
49 mented in real-world clinical settings without requiring extensive manual annotation
50 efforts. The specific objectives of this research are:

- 51 • Develop a weakly supervised DL model for prostate cancer detection and ISUP
52 grading in whole slide images with significantly reduced annotation requirements
53 compared to traditional fully supervised approaches.
- 54 • Implement and evaluate multiple state-of-the-art weakly supervised learning strate-
55 gies for handling WSI data, including attention mechanisms and multiple instance
56 learning approaches, to identify the most effective methods for prostate cancer
57 analysis.
- 58 • Conduct comprehensive benchmarking of the proposed framework's performance
59 for cancer detection and ISUP grading using established evaluation metrics and
60 publicly available datasets to ensure robust and reliable results.

- 61 • Investigate model transparency and interpretability using advanced visualization
62 techniques such as attention maps and Grad-CAM to improve clinical trust and
63 facilitate adoption of the AI framework in pathology practice.
- 64 • Evaluate the impact of different patch sizes and overlap strategies on model per-
65 formance to optimize the balance between computational efficiency and diagnostic
66 accuracy in weakly supervised learning frameworks.
- 67 • Developing a web-based clinical tool that is end-to-end and seamlessly integrates
68 with current pathology workflows.

69 To achieve these objectives and validate the proposed framework, the research
70 adopts a systematic and comprehensive approach that combines technical innovation
71 with practical applicability, emphasizing rigorous experimental design and large-scale
72 validation to ensure reliable and clinically relevant results.

- 73 • Carried out a systematic assessment of 72 experimental configurations by combining
74 six advanced multiple instance learning models (CLAM-MB, CLAM-SB, ILRA-
75 MIL, AC-MIL, WiKG-MIL, AMD-MIL), three feature extraction architectures
76 (ResNet50, CTransPath, UNI2), and four patch processing strategies.
- 77 • Leveraged a large PANDA dataset containing 10,616 whole slide images and applied
78 the UNI2 foundation model, pretrained on over 100 million histopathology images,
79 for weakly supervised prostate cancer analysis.
- 80 • Utilize distributed cloud-based computing, optimized patch extraction, and efficient
81 training pipelines to manage gigapixel-scale image data.
- 82 • Ensure clinical usability through interpretability tools (e.g., attention maps, Grad-
83 CAM, heatmaps), and rigorous validation using cross-validation and benchmark
84 comparisons.

85 This study offers a unique and clinically focused method for diagnosing prostate
86 cancer that greatly lowers annotation overhead without sacrificing diagnostic accuracy
87 using a weakly supervised deep learning architecture. The proposed framework solves
88 the operational and technological constraints hindering AI adoption in pathology by
89 utilizing scalable cloud-based computing, robust foundation models, and cutting-edge
90 MIL architectures. In addition to improving diagnostic efficiency and consistency, this
91 research opened the door for the practical application of AI-driven solutions in actual
92 clinical settings by means of thorough benchmarking, improvements to model inter-
93 pretability, and the creation of an integrated web-based tool. This study's results and
94 methods are intended to provide scholarly contributions to the field of computational
95 pathology as well as practical contributions to cancer diagnoses.

96 2 Literature Review

97 The application of MIL for WSI analysis has significantly increased in recent years,
98 particularly for Gleason grading and prostate cancer diagnosis. By eliminating expen-
99 sive pixel- or region-level annotations and learning slide-level labels from sets of
100 patch-level data, MIL offers an efficient weakly supervised approach. This section criti-
101 cally evaluates recent approaches, highlighting methodological trends, constraints, and
102 unresolved issues that drive our benchmarking analysis.

103 2.1 Attention-Based MIL and Frequency/Spatial Fusion

104 A learnable pooling approach that weights instance contributions and yields inter-
105 pretable attention ratings was developed by attention-based MIL (ABMIL) [12]. Lu et
106 al. [13] expanded on this by proposing CLAM, which adds clustering-constrained atten-
107 tion to increase the multi-class Gleason grading’s resilience. Since then, this strategy
108 has taken over as the prevailing paradigm. More recent enhancements expand patch
109 representations by fusing characteristics in the frequency and spatial domains. For
110 instance, Zhang et al. [14] introduced FRCM-MIL, which combines cross-attention,
111 confidence query aggregation, and wavelet-based frequency reconstruction. FRCM-
112 MIL performed well on clinical datasets, including PUMCH and PANDA (PUMCH:
113 81.75% accuracy, AUC 0.9441; PANDA: 67.24% accuracy, AUC 0.9169). Although the
114 dependence on custom transformations and specific modules increases pipeline com-
115 plexity and decreases portability across staining variances, our results demonstrate the
116 significance of complementing frequency characteristics and sophisticated aggregation
117 methodologies.

118 2.2 Multi-Resolution and Hierarchical Attention

119 The gigapixel size of WSIs is addressed by multi-resolution techniques that combine
120 fine-grained improvements with coarse global scans. A multi-resolution attention MIL
121 pipeline was presented [15]. It refines the results by applying attention at higher
122 magnification after screening low-magnification patches. High accuracy (~85%) and
123 interpretable attention maps were the results of this method. However, careful scale
124 selection and fusion strategy design are necessary for multi-resolution MIL. Concerns
125 regarding generalizability to diverse clinical datasets are raised by the fact that it
126 is also more computationally intensive and frequently depends on carefully selected
127 biopsy groups.

128 Multi-resolution techniques are still being expanded in recent work to enhance
129 performance and interpretability in WSI and histopathology jobs. For example, a two-
130 stage multi-resolution CNN pipeline is proposed in [16]. Contextual characteristics
131 are extracted in the first step using CNNs at four different resolutions; they are then
132 combined with another CNN to create segmentation masks. They report pixel-wise
133 accuracies of around 95.6% and mean dice of approximately 92.5% on placenta and
134 lung datasets, demonstrating that integrating various resolutions greatly enhances the
135 segmentation of objects with different scales. The multi-resolution segment anything
136 model (SAM) for histopathology WSI (WSI-SAM) is another recent study [17]. By
137 combining high-resolution (HR) and low-resolution (LR) tokens with a dual-mask
138 decoder that combines information from different resolutions, WSI-SAM improves the
139 SAM model.

140 This research focuses on comparing many MIL designs for prostate cancer detection
141 and grading in a slide-level classification context, in contrast to existing multi-
142 resolution pipelines that mostly prioritize segmentation or depend on meticulously
143 designed fusion methods. Without the need for manually designed multi-resolution
144 fusion modules, we systematically and reproducibly capture resolution effects by
145 directly evaluating patch size and overlap tactics over 72 controlled tests. By doing

146 this, our framework complements but also streamlines the more specialized multi-
147 resolution techniques of Li et al. [15], Salsabili et al. [16], and Zheng et al. [17] by
148 offering a generalizable benchmark that strikes a compromise between computational
149 practicality and biological interpretability.

150 **2.3 Graph-Based and Representation-Driven Aggregation**

151 Graph-based MIL captures the spatial and semantic information that traditional
152 pooling misses by explicitly modeling connections between patches. Behzadi et al.
153 [18] showed how patch adjacency and similarity can be used to create robust
154 prostate cancer grading in graph convolutional networks (GCNs) with noisy-label
155 filtering. Although these relational models are sensitive to hyperparameters like adja-
156 cency radius and similarity criteria and computationally costly, they are effective at
157 capturing contextual information.

158 On top of this, more contemporary models incorporate global dependencies and
159 local graph structure. For instance, the integrative graph-transformer framework for
160 WSI classification [19] improves AUROC/accuracy over previous graph and attention
161 baselines by adding transformer-based global attention layered on a GCN-based rela-
162 tional graph to capture both intra-patch adjacency and global context. Similar to this,
163 GRASP [20] employs a pyramidal graph structure at various magnifications, retaining
164 interpretability by node aggregation assessed by skilled pathologists and delivering up
165 to $\sim 10\%$ increases in balanced accuracy with far fewer parameters.

166 Unlike these approaches, the current study carefully benchmarks many current MIL
167 designs (including graph-MIL) over a range of encoders, patch sizes, and overlap set-
168 tings rather than proposing a novel graph or transformer architecture. The repeatable,
169 slide-level classification and grading analysis demonstrates the relative effectiveness of
170 graph-based techniques in controlled, realistic environments, particularly when paired
171 with histopathology-specific encoders like UNI2 and CTransPath.

172 **2.4 Transformers and Global-Context Models**

173 Transformer-based MIL techniques simulate global interactions and cross-instance
174 interdependence within a bag, extending attention. Correlated-instance self-attention
175 enhances slide-level categorization by capturing long-range patch dependencies,
176 as shown by TransMIL [21]. Pathology analysis has been further enhanced
177 by transformer-based encoders like CTransPath [22] and foundation models like
178 UNI/UNI2 [23] that pretrain on extensive histopathology data. Although these meth-
179 ods offer advanced representations, they necessitate substantial computing resources
180 for pretraining and meticulous patch sampling techniques for gigapixel slides.

181 This is further supported by other current works: Compared to patch-only or basic
182 transformer approaches, PathTR improves classification and localization robustness
183 by including context-aware memory into the transformer backbone to better encode
184 slide-wide structure for tumor localization. TransGNN improves prognosis accuracy
185 on hepatocellular carcinoma slides by combining transformer global attention with
186 graph structural characteristics, allowing for both explicit local relational reasoning
187 and global representation.

188 2.5 Systematic Comparisons, Encoder Dependence, and 189 Interpretability

190 Several approaches (completely supervised, weakly supervised, attention-based MIL,
191 CLAM, and TransMIL) were benchmarked across various datasets in comparative
192 studies [24]. They demonstrated how attention-based MIL effectively strikes a com-
193 promise between prediction performance and annotation cost. However, no single
194 technique consistently outperforms the others; instead, performance is influenced
195 by the assessment process, staining variability, label dispersion, and dataset size.
196 Encoder selection is crucial, as evidenced by the persistent superior performance
197 of histopathology-pretrained encoders like CTransPath and UNI2 over ImageNet-
198 pretrained models like ResNet50.

199 Interpretability is still a challenge. Although the majority of works use atten-
200 tion maps as a stand-in for explainability, they seldom ever validate these maps
201 using impartial measurements or professional opinions. XViT [25] is one of the
202 recent explainability-focused techniques that incorporates quantitative measurements
203 of explanation quality (sensitivity, fidelity, and complexity). These works stress
204 that rather than depending just on visual examination, interpretability should be
205 thoroughly assessed and clinically confirmed.

206 We ensure that interpretability is both clinically relevant and statistically grounded
207 by combining attention and Grad-CAM visualizations with expert pathologist valida-
208 tion, in contrast to the majority of the literature that solely uses qualitative attention
209 maps to convey interpretability.

210 2.6 Broader AI and Segmentation Trends

211 MIL and explainability are becoming key components of applied, interpretable AI
212 in the engineering and medical sectors, according to a bibliometric review of Engi-
213 neering Applications of Artificial Intelligence by Shukla et al. [26]. Pathology has
214 also been shaped by parallel improvements in segmentation. In their assessment of
215 93 transformer-based segmentation models, Xiao et al. [27] discovered that U-Net
216 + Transformer hybrids perform better than baselines that just use CNN. According
217 to BioMedical Engineering Online (2024), hybrid CNN-transformer models are more
218 prevalent than pure transformers, which are constrained by a lack of data. Zhang et al.
219 [28] demonstrated competitive performance while using a pure ViT for segmentation,
220 proving that it is viable even in the absence of convolutional operators.

221 2.7 Latest Prostate Cancer-Specific Advances

222 Recent Developments Particular to Prostate Cancer: Prostate cancer WSIs are the
223 direct subject of recent research. In their development of Hierarchical ViTs for prostate
224 biopsy grading on PANDA, Grisi et al. [29] showed significant generalization with
225 QWK 0.916 in-domain and 0.877 cross-domain. A generalized self-supervised ViT
226 was presented by Chaurasia et al. [11], which lessens the need for expensive labeling.
227 When comparing U-Net with ViT for zonal segmentation, Huang et al. [30] discovered
228 that transformers performed better in semi-supervised environments. Zheng et al. [17]

229 extended the concept beyond histology by using poorly supervised MIL-like pooling
230 for MRI prostate cancer diagnosis.

231 **2.8 Comparison with this Study**

232 This work systematically benchmarks a range of existing MIL and transformer-
233 type architectures under controlled settings (varying patch sizes, overlaps, encoders)
234 for prostate cancer detection and grading, whereas PathTR and TransGNN specif-
235 ically develop new model architectures to better encode global context (memory,
236 graph, transformer fusion). Instead of suggesting yet another design, the findings help
237 determine which global-context models work best in reality and under what conditions.

238 A summary of the discussed works is given in Table 1. The literature review
239 reveals significant progress in applying deep learning approaches to prostate cancer
240 histopathology analysis, particularly through weakly supervised learning and multi-
241 ple instance learning frameworks. The evolution from traditional supervised methods
242 to sophisticated MIL architectures like CLAM, TransMIL, ILRA-MIL, AC-MIL, and
243 AMD-MIL demonstrates the field’s maturation in addressing the fundamental chal-
244 lenge of learning from slide-level labels without pixel-wise annotations. Similarly, the
245 development of specialized feature extractors from basic convolutional networks like
246 ResNet50 to transformer-based approaches such as CTransPath and foundation models
247 like UNI2 highlights the importance of domain-specific pretraining for computational
248 pathology applications.

Table 1: Summary of Literature Trends in MIL and Transformer-based Pathology

Study / Approach	Key Idea	Strengths	Limitations	Reported Results
[12] (2018), ABMIL	Learnable attention pooling	Simple, interpretable	Limited contextual modeling	Solid baseline across WSI tasks
[13](2021), CLAM	Clustering-constrained attention	Robust multi-class grading	Hyperparameter sensitive	Improved multi-class Gleason grading
[28] (2024), FRCM-MIL	Frequency + spatial fusion with cross-attention	Rich features, strong performance	Complex pipeline, less portable	PUMCH: 81.75% Acc., 0.9441 AUC; PANDA: 67.24% Acc., 0.9169 AUC
[15] (2019), Multi-Res Attention	Coarse-to-fine hierarchical MIL	Balances context and detail	Complex design, scale-sensitive	~85% accuracy
[18] (2022), Graph-MIL	Graph-based relational aggregation	Captures spatial/semantic context	Computationally expensive	Robust grading results
[21] (2021), TransMIL	Transformer-based MIL	Models long-range dependencies	Memory and sampling constraints	High AUC across datasets
[22] (2022), CTransPath	Transformer-based contrastive pretraining	Domain-specific encoder	Pretraining cost	Strong performance across pathology tasks
[23] (2024), UNI/UNI2	Foundation pathology encoder	State-of-the-art performance	Very high pretraining compute	Superior results across WSI tasks
[25] (2025), XViT	Explainable transformer with metrics	Quantitative explainability	Scaling to WSI still open	High accuracy, strong interpretability
[29] (2025), Hierarchical ViT	Multi-level ViT for prostate biopsy grading	Strong generalization	Needs large training data	QWK 0.916 (in-domain), 0.877 (cross-domain)
[11] (2025), Self-Supervised ViT	Label-efficient ViT across datasets	Reduces annotation burden	Requires multi-dataset training	Competitive prostate grading
H[30] (2024), U-Net vs. ViT	Semi-supervised zonal segmentation	Transformers better with scarce labels	Task-specific scope	ViT superior in segmentation tasks
[17] (2024), MRI MIL	Weak supervision with MIL pooling for MRI	Extends MIL beyond histopathology	Modality-specific limitations	Reduced unnecessary biopsies
This work (2025)	Systematic MIL benchmarking with encoders, patch sizes, overlaps	Unified, reproducible, clinically validated	Requires large-scale experiments	72 experiments, clear encoder and MIL comparisons

249 2.9 Critical Research Gaps and Reproducibility Issues

250 Across the literature, several recurring limitations emerge:

- 251 i. **Heterogeneous evaluation protocols:** Heterogeneous assessment protocols:
252 Studies vary greatly in patch sizes, overlaps, encoders, and preprocessing decisions,
253 which makes cross-paper comparisons challenging and reproducibility challenging.
- 254 ii. **Limited benchmarking frameworks:** Few studies offer comprehensive frame-
255 works that systematically compare several MIL techniques across various feature
256 extraction methodologies; most available studies focus on particular model designs
257 or limited encoder comparisons.
- 258 iii. **Scalability challenges:** In resource-constrained contexts where effective deploy-
259 ment is crucial, the computing load of processing gigapixel-sized WSIs with different
260 patch sizes and overlap methods is still little understood.
- 261 iv. **Encoder dependence:** Although systematic, large-scale encoder comparisons
262 are uncommon, domain-specific encoders (such as CTransPath and UNI/UNI2)
263 frequently perform better than general ImageNet models.
- 264 v. **Interpretability evaluation:** Although attention maps and methods such as
265 Grad-CAM are often reported, the majority of research only goes as far as qual-
266 itative visualization. There is also a lack of clinical validation and quantitative
267 evaluation of interpretability, particularly when it comes to whole-patch-based
268 diagnostic pipelines.
- 269 vi. **Computational cost and reproducibility:** Multi-resolution, graph, and trans-
270 former models frequently entail many hyperparameters and need a large amount of
271 resources, which restricts their accessibility and reproducibility for wider use.

272 2.10 Research Contributions

273 This study directly addresses these gaps:

- 274 • **Extensive Benchmarking of MIL Architectures:** This work systematically
275 benchmarks several cutting-edge MIL architectures under a single experimental pro-
276 tocol, highlighting their relative strengths and weaknesses in real-world diagnostic
277 settings, whereas the majority of previous studies on prostate cancer concentrate
278 on examining a single MIL framework or model variant.
- 279 • **Attention-Based Weak Supervision for Grading:** This study uses attention-
280 based weak supervision for the clinically difficult task of Gleason grading, going
281 beyond binary cancer diagnosis. This makes it possible for our approach to identify
282 subtle histopathological features that are important for prognosis but haven't been
283 thoroughly examined in previous studies.
- 284 • **Encoder-focused analysis:** This research provides one of the first compre-
285 hensive comparisons of histopathology-specific encoders (CTransPath, UNI2) and
286 ImageNet-pretrained ResNet50, demonstrating consistent performance improve-
287 ment with domain-specific pretraining.
- 288 • **Patch-to-Pathology Workflow:** The current study provides an open and repeat-
289 able end-to-end process specifically designed for prostate histopathology, encom-
290 passing patch extraction, feature embedding, MIL aggregation, and interpretability.

291 This systematic workflow offers a reliable benchmark methodology that acts as a
292 roadmap for further research.

293 • **Interpretability for Clinical Insight:** This study incorporates interpretability
294 modules that display patch-level attention maps in addition to benchmarking per-
295 formance. This increases the clinical significance of the models by bridging the gap
296 between pathologists' demands and black-box MIL designs.

297 • **Public Benchmarking Resource:** The current work offers benchmark results
298 and repeatable assessment procedures that may direct practitioners and researchers,
299 acting as a foundation for creating more reliable AI-assisted diagnostic systems in
300 digital pathology.

301 3 Methodology

302 The proposed approach transforms the complex challenge of analyzing gigabyte-
303 sized whole slide images into a manageable computational pipeline that can achieve
304 pathologist-level performance in cancer assessment. The methodology encompasses
305 seven interconnected stages that work together to process the PANDA dataset's 10,616
306 whole slide images and generate accurate cancer grade predictions. Figure 1 shows the
307 comprehensive methodology of the proposed approach.

308 The complete pipeline begins with the Prostate cANcer graDe Assessment
309 (PANDA) dataset preprocessing, followed by patch creation, where tissue regions are
310 identified and coordinate maps are generated for subsequent analysis. We then extract
311 over 31 million image patches from these coordinates, encode each patch using three
312 different deep learning architectures (ResNet50, CTransPath, and UNI2) to capture
313 diverse histopathological features, and organize the data using stratified sampling with
314 cross-validation splits. The feature representations feed into six state-of-the-art MIL
315 models that learn to aggregate patch-level information for slide-level cancer grading.
316 Finally, the model performance is evaluated using comprehensive metrics, including
317 accuracy, precision, recall, F1 score, and area under the curve, while incorporating
318 attention mechanism analysis for interpretability.

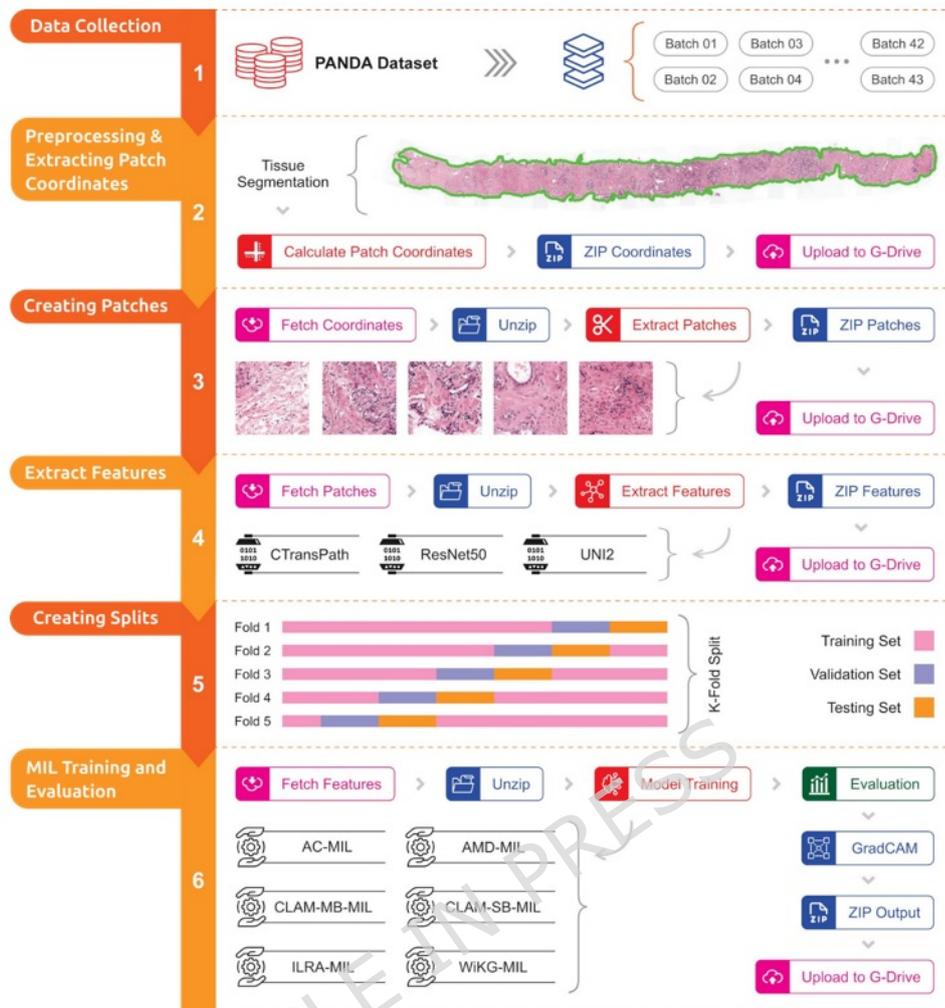


Fig. 1: Comprehensive methodology diagram.

3.1 Data Collection

The PANDA dataset represents the largest publicly available collection of digitized whole slide images for prostate cancer analysis. This dataset was created as part of the PANDA Great Challenge by Radboud University Medical Center and Karolinska Institute [10] and is publicly available at Kaggle (<https://www.kaggle.com/c/prostate-cancer-grade-assessment/data>). The dataset provides 10,616 whole slide images of H&E-stained prostate tissue biopsies from two medical centers. Figure 2 shows a few samples from the dataset.

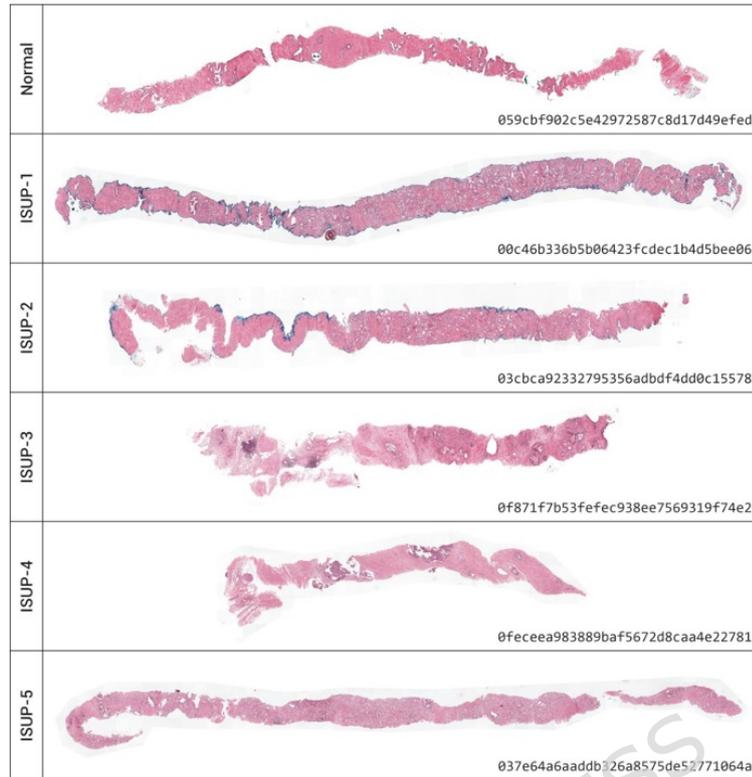


Fig. 2: Randomly picked sample whole slide image from PANDA dataset, for each class.

327 The dataset serves as the foundation for our patch-based deep learning framework,
 328 providing comprehensive histopathological data with expert pathologist annotations
 329 based on the ISUP grading system (grades 0-5). Each slide contains tissue samples
 330 obtained through needle core biopsy procedures, with images captured at 20x mag-
 331 nification and stored in TIFF format with approximately $20,000 \times 20,000$ pixels per
 332 slide. The dataset uses the International Society of Urological Pathology (ISUP) grade
 333 system for prostate cancer classification.

- 334 • ISUP Grade 0: Benign tissue with no cancer detected.
- 335 • ISUP Grade 1: Low-grade cancer (Gleason 3+3) with well-differentiated glands
- 336 • ISUP Grade 2: Intermediate-low grade (Gleason 3+4) with mixed patterns.
- 337 • ISUP Grade 3: Intermediate-high grade (Gleason 4+3) with poorly formed glands.
- 338 • ISUP Grade 4: High-grade cancer (Gleason 4+4, 3+5, 5+3) with significant
 339 architectural loss.
- 340 • ISUP Grade 5: Highest grade cancer (Gleason 4+5, 5+4, 5+5) with solid growth
 341 patterns.

342 This grading system provides the ground truth labels for training the multiple
343 instance learning models to automatically detect and grade prostate cancer from
344 histopathological images.

345 To manage the computational complexity of processing 10,616 whole slide images
346 while enabling parallel processing capabilities, we implemented a systematic batch-
347 based approach that divides the entire PANDA dataset into 43 distinct batches, with
348 each batch containing approximately 250 slides. This batch-wise processing strategy
349 was consistently applied across the preprocessing (Section 3.2), patch creation (Section
350 3.3), and feature extraction (Section 3.4) stages, with processed batches subsequently
351 combined to formulate the final consolidated dataset used for creating splits (Section
352 3.5) and MIL training (Section 3.6). This approach significantly reduced computational
353 burden while enabling parallel processing across available hardware infrastructure.

354 3.2 Preprocessing and Extracting Patch Coordinates

355 The patch creation step forms the first and most important part of our deep learning
356 system. This process takes large WSIs, which are digital versions of tissue slides viewed
357 under a microscope, and breaks them down into smaller, manageable pieces called
358 patches. This systematic patch creation process established a solid foundation for the
359 rest of our deep learning pipeline, ensuring that we work only with meaningful tissue
360 content while maintaining computational efficiency across the entire PANDA dataset.
361 These patches are then used for training machine learning models to detect prostate
362 cancer and determine its grade. Figure 3 shows the flow of the patch creation.

363 We need to identify tissue areas in each slide and create a map of where to extract
364 patches for analysis. In this regard, we needed to solve several key problems: whole slide
365 images are extremely large (often several gigabytes each), they contain lots of empty
366 background space, and we had limited computing power to process over 10,000 slides.
367 We developed a system that automatically finds tissue regions, removes background
368 areas, and creates coordinate lists showing exactly where to extract patches. This
369 approach saves enormous amounts of storage space and processing time because we
370 only work with meaningful tissue areas instead of entire slide images. The process
371 works in three main stages:

372 3.2.1 Finding Tissue Areas

373 We used computer vision techniques to automatically separate tissue from background
374 in each slide. The tissue detection process works in four main steps.

- 375 • First, we convert each image from RGB color format to hue, saturation, value (HSV)
376 format. HSV is better for tissue detection because the saturation channels are clear.
377 Tissue areas have higher color saturation, while background areas are mostly white
378 or very light colored.
- 379 • Secondly, we apply median filtering to remove small spots of noise and artifacts.
380 Median filtering works by replacing each pixel with the middle value of its surround-
381 ing pixels, which removes random noise. We used a filter size of 7 pixels based on
382 the typical size of noise in histopathological images.

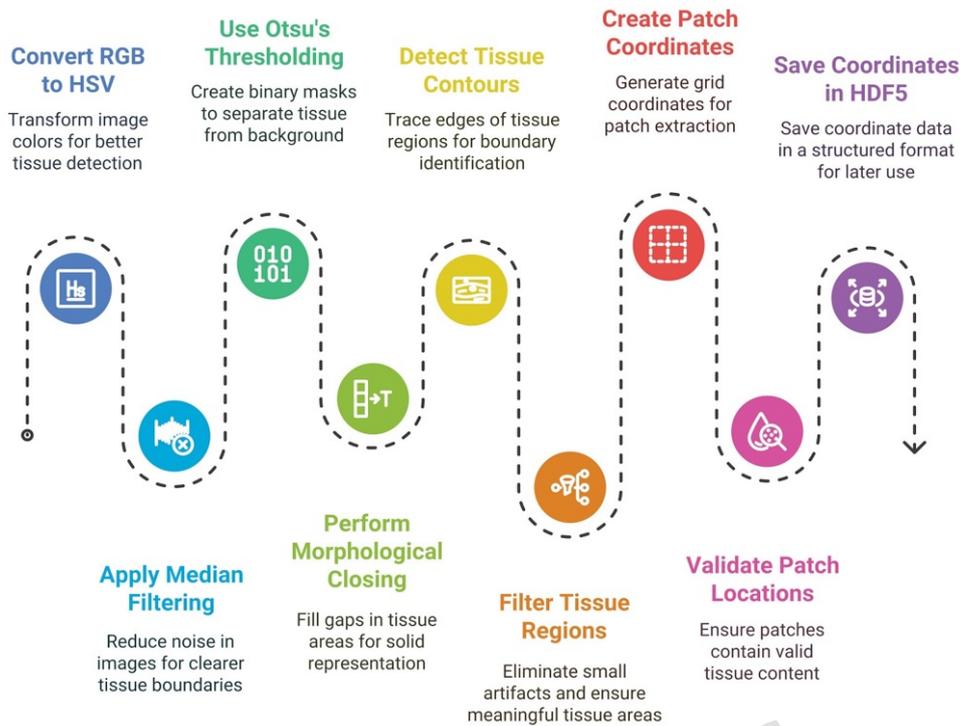


Fig. 3: Flow for patch creation process from whole slide images in the proposed framework.

- 383 • Third, we use Otsu's thresholding method to create black and white masks where
 384 white areas represent tissue and black areas represent background. Otsu's method
 385 automatically finds the best threshold value by analyzing the image's brightness
 386 distribution.
- 387 • Finally, we use morphological closing operations to fill small gaps and holes within
 388 tissue areas. This technique uses a small circular shape (4×4 pixels) to connect
 389 nearby tissue pieces and smooth rough edges. This step ensures that tissue areas
 390 are represented as solid regions rather than fragmented pieces.

391 3.2.2 Identifying Valid Tissue Regions

392 Figure 4 shows the flow of finding valid tissue regions. We needed to find the actual
 393 boundaries of tissue regions and filter out small artifacts or preparation errors. We used
 394 contour detection algorithms to trace the edges of tissue regions. A contour is simply
 395 the boundary line around a shape, in our case, around tissue areas. The algorithm
 396 we used can detect complex shapes, including regions with holes inside them. Not
 397 all detected regions were suitable for patch extraction. We established minimum size
 398 requirements to eliminate tiny artifacts and ensure we only work with meaningful
 399 tissue areas. We set the minimum tissue area to 16 times the size of our intended

400 patches, and the minimum hole size to 4 times the patch size. We also limited the
 401 number of holes per tissue region to 8 to avoid overly complex areas.

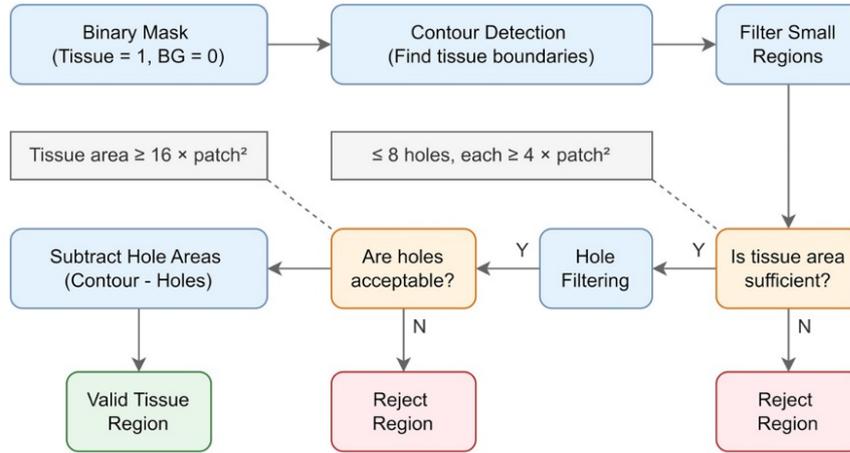


Fig. 4: Flow for identifying valid tissue region from WSIs in the proposed framework.

402 The filtering process calculates the actual tissue area by subtracting hole areas
 403 from the total contour area. Only regions meeting our size criteria are kept for patch
 404 coordinate generation. Some tissue regions contain holes or cavities that represent
 405 blood vessels, glands, or preparation artifacts. Our system tracks the relationship
 406 between tissue boundaries and their internal holes, ensuring that patch coordinates
 407 are never placed in these empty areas.

408 3.2.3 Creating Patch Coordinates

409 After identifying valid tissue regions, we created systematic coordinate grids to specify
 410 exactly where patches should be extracted from each slide. We used a regular grid
 411 pattern to ensure complete and uniform coverage of tissue areas. The grid spacing
 412 depends on whether we want overlapping patches or not. For non-overlapping patches,
 413 we place coordinates such that each patch touches its neighbors but doesn't overlap.
 414 For 50% overlapping patches, we place coordinates such that adjacent patches share
 415 half their area.

416 This systematic approach ensures we don't miss any tissue areas and provides
 417 consistent sampling across all slides. Each potential patch location goes through vali-
 418 dation to ensure it contains meaningful tissue content. We use a "four-point" check-
 419 ing method that tests whether the corners and center of each patch fall within valid tissue
 420 areas. This prevents patches that would be mostly background or would cross tissue
 421 boundaries.

422 3.2.4 Data Organization and Storage

423 We needed an efficient way to store and organize the coordinate information for later
 424 use in our pipeline. We chose HDF5 (Hierarchical Data Format - a specialized file
 425 format for scientific data) for storing coordinate data because it provides fast access to
 426 large datasets and includes compression to save storage space. Each whole slide image
 427 generates one HDF5 file containing an array of (x, y) coordinates for all valid patch
 428 locations, metadata including patch size and processing parameters, slide information,
 429 and quality metrics and processing statistics.

430 HDF5 format allows other parts of our system to quickly read coordinate data
 431 without loading entire files into memory. Each HDF5 file follows a consistent naming
 432 pattern based on the original slide ID, making it easy to locate coordinate data for
 433 any slide. The files include comprehensive metadata that documents all processing
 434 parameters, enabling reproducible results and quality assessment.

435 3.3 Creating Patches

436 We designed four different experimental configurations to study how patch size and
 437 overlap affect cancer detection performance:

- 438 • Setting 1: 512×512 pixel patches with no overlap
- 439 • Setting 2: 512×512 pixel patches with 50% overlap
- 440 • Setting 3: 256×256 pixel patches with no overlap
- 441 • Setting 4: 256×256 pixel patches with 50% overlap

442 These settings allow us to compare larger patches (which capture more context)
 443 versus smaller patches (which provide more detailed views), and to evaluate whether
 444 overlapping patches improve detection accuracy despite requiring more computational
 445 resources. The patch creation process successfully generated coordinate datasets for all
 446 experimental configurations, demonstrating the effectiveness of our distributed com-
 447 puting approach. The four experimental settings produced different numbers of patches
 448 reflecting the impact of size and overlap parameters, as shown in Table 2.

Table 2: Patch generation statistics across four experimental settings, highlighting the impact of patch size and overlap on output volume.

Experiment Setting	Processed Slides	Discarded Slides	Generated Patch Coordinates	Average Patches Per Slide
Setting 1 (512×512 , no overlap)	10,202	414	1,292,600	127
Setting 2 (512×512 , 50% overlap)	10,202	414	5,087,711	499
Setting 3 (256×256 , no overlap)	10,596	20	4,879,367	460
Setting 4 (256×256 , 50% overlap)	10,596	20	19,385,634	1830

449 The distributed processing architecture successfully handled the computational
 450 demands of large-scale histopathological image analysis. The batch-based approach
 451 maintained data quality while working within resource constraints, proving that

452 sophisticated medical image analysis can be performed using freely available comput-
 453 ing platforms.

454 3.4 Feature Extraction

455 The feature extraction step transforms individual patch images into high-dimensional
 456 numerical representations that capture meaningful patterns for machine learning anal-
 457 ysis. This process takes the extracted patch images from the previous step and feeds
 458 them through pre-trained deep learning models (encoders) to generate feature vec-
 459 tors that encode important visual characteristics like texture, color patterns, cellular
 460 structures, and spatial relationships. These feature vectors serve as the foundation
 461 for training multiple instance learning (MIL) models for prostate cancer detection
 462 and grading. Mathematically, the feature extraction process can be represented as a
 463 mapping function:

$$f : R^{H \times W \times C} \rightarrow R^d \quad (1)$$

464 where an input patch image $x \in R^{H \times W \times C}$ (with height H , width W , and channels C)
 465 is transformed into a feature vector $z \in R^d$ of dimension d , where d varies by encode
 466 architecture.

467 The objective in feature extraction was to convert the thousands of patch images
 468 from each slide into numerical feature representations that machine learning models
 469 can understand and process effectively. Raw pixel values in medical images contain
 470 too much noise and irrelevant information for direct analysis, so we needed to extract
 471 meaningful features that capture the important visual patterns related to cancer
 472 detection. We needed to solve several key challenges: selecting appropriate pre-trained
 473 models that understand medical image patterns, processing large numbers of patches
 474 efficiently while maintaining consistent quality, and organizing the resulting features in
 475 a way that supports multiple instance learning approaches, where each slide is treated
 476 as a collection (bag) of patches.

477 We developed a systematic feature extraction pipeline that uses three different
 478 state-of-the-art encoder networks to transform patch images into high-dimensional
 479 feature vectors. This multi-encoder approach allows us to compare different feature
 480 representation strategies and determine which works best for prostate cancer anal-
 481 ysis. Each encoder brings unique strengths: ResNet50 provides proven convolutional
 482 features, CTransPath offers transformer-based representations specifically trained on
 483 pathology data, and UNI2 delivers cutting-edge histopathology foundation model
 484 capabilities.

485 The process works in four main stages: first, we organize patch images into batches
 486 for efficient processing; second, we apply appropriate preprocessing transformations
 487 for each encoder; third, we extract features using the selected encoder network; and
 488 fourth, we aggregate and save features as slide-level collections ready for multiple
 489 instance learning.

490 We carefully selected 3 encoder architectures that represent different approaches
 491 to medical image analysis and have proven effectiveness in computational pathology
 492 applications.

493 3.4.1 ResNet50

494 ResNet50 serves as our baseline encoder, representing the established standard in
 495 medical image analysis. This convolutional neural network, introduced by [31], revo-
 496 lutionized deep learning by solving the vanishing gradient problem through residual
 497 connections. ResNet50 contains 50 layers with skip connections that allow informa-
 498 tion to flow directly between layers, enabling the training of much deeper networks
 499 than previously possible. The core innovation of ResNet50 lies in its residual learning
 500 framework, where instead of learning unreferenced functions, the layers learn residual
 501 functions concerning the layer inputs. Mathematically, this can be expressed as:

$$y = F(x, \{W_i\}) + x \quad (2)$$

502 where x and y are the input and output vectors, and $F(x, \{W_i\})$ represents the residual
 503 mapping to be learned.

504 ResNet50 has been extensively used in histopathology applications and provides a
 505 solid foundation for comparison with newer approaches. The model was pre-trained
 506 on ImageNet and fine-tuned for medical imaging tasks, making it particularly suitable
 507 for extracting low-level visual features like edges, textures, and color patterns that are
 508 crucial for identifying cancerous tissue characteristics.

509 3.4.2 CTransPath

510 CTransPath represents the next generation of pathology-specific encoders, combining
 511 convolutional and transformer architectures specifically designed for histopatholog-
 512 ical image analysis. Developed by [22], this hybrid model was pre-trained using
 513 self-supervised learning on massive histopathology datasets, including TCGA and
 514 PAIP.

515 CTransPath uses a Swin Transformer backbone combined with a convolutional
 516 stem to capture both local and global patterns in tissue images. The model employs
 517 semantically-relevant contrastive learning (SRCL) during pre-training, which helps it
 518 understand the relationships between different tissue types and pathological patterns.
 519 The self-attention mechanism in transformers can be mathematically represented as:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (3)$$

520 where Q , K , and V are the query, key, and value matrices, respectively, and d_k is
 521 the dimension of key vectors. This specialized training makes CTransPath particularly
 522 effective at recognizing the complex spatial arrangements and cellular morphologies
 523 characteristic of prostate cancer.

524 3.4.3 UNI2

525 UNI2 represents the cutting-edge in pathology foundation models, released in 2024 as
 526 the successor to the highly successful UNI model. Developed by the Mahmood Lab at
 527 Harvard/BWH, UNI2 is a Vision Transformer (ViT-H/14) trained on over 200 million
 528 pathology images from more than 350,000 diverse whole slide images covering H&E
 529 and IHC staining [23].

530 UNI2 uses advanced self-supervised learning techniques, including DINOv2, iBOT
 531 masked-image modeling, and KoLeo regularization to learn rich representations with-
 532 out requiring labeled data. The model’s 1536-dimensional feature vectors capture
 533 fine-grained pathological patterns and have shown state-of-the-art performance across
 534 34 computational pathology tasks. For our prostate cancer analysis, UNI2 provides
 535 the most sophisticated understanding of tissue morphology and cellular patterns.

536 The Vision Transformer architecture processes images by dividing them into
 537 patches and treating them as sequences. For an input image $x \in R^{H \times W \times C}$, it is
 538 reshaped into a sequence of flattened 2D patches $x_p \in R^{N \times (P^2 \cdot C)}$, where (H, W)
 539 is the resolution of the original image, C is the number of channels, (P, P) is the re-
 540 solution of each image patch, and $N = \frac{HW}{P^2}$ is the resulting number of patches [32]. A
 541 comparison of encoder characteristics is provided in Table 3.

Table 3: Comparison of encoder characteristics.

Characteristics	ResNet50	CTransPath	UNI2
Encoder Architecture	Convolution Neural Network	Hybrid Convolution and Transformer	Vision Transformer (ViT-H/14)
Pre-Training Data	ImageNet Dataset	Massive Histopathology Datasets (TCGA, PAIP)	Over 200 million Pathology Images
Feature Dimensions	2048	768	1536
Key Innovation	Residual Learning Framework	Semantically Relevant Contrastive Learning (SRCL)	Advanced Self-Supervised Learning Techniques

542 3.4.4 Feature Extraction Process

543 The core feature extraction process involves feeding batches of preprocessed patch
 544 images and collecting the resulting high-dimensional feature vectors. For each slide, we
 545 load all associated patch images using the coordinates stored in the H5 files from the
 546 patch creation step. The patches are organized into batches of 128 images (configurable
 547 based on GPU memory) and fed through the encoder network in a systematic forward
 548 pass that extracts features without updating the model weights.

549 The batch processing approach optimizes computational efficiency by processing
 550 multiple patches simultaneously. For a batch of patches $X = \{x_1, x_2, \dots, x_B\}$ where
 551 B is the batch size, the encoder processes them in parallel:

$$Z = f(X) = \{f(x_1), f(x_2), \dots, f(x_B)\} \quad (4)$$

552 where each $f(x_i)$ represents the feature vector extracted from patch x_i .

553 3.4.5 Multi-Encoder Processing Strategy

554 To enable a comprehensive comparison of different representation approaches, we
 555 implemented a systematic multi-encoder processing strategy that generates features
 556 using all three encoder networks for each experimental setting, resulting in 12 distinct

557 feature collections (4 patch settings \times 3 encoders). This systematic approach enables
 558 direct comparison of encoder performance across different patch sizes and overlap
 559 strategies while maintaining consistent experimental conditions.

560 Each encoder runs independently on the same set of patch images, ensuring that
 561 differences in feature quality stem from the encoder architecture rather than variations
 562 in input data. This approach provides several scientific advantages: it allows us to
 563 evaluate which type of feature representation works best for prostate cancer detection,
 564 enables ensemble methods that combine features from multiple encoders, and provides
 565 redundancy in case of processing issues with individual encoders.

566 3.4.6 Slide-Level Feature Aggregation

567 After extracting features from individual patches, we needed to aggregate them into
 568 slide-level representations that maintain the relationship between patches while cre-
 569 ating manageable datasets for MIL approaches. For each slide, all patch features are
 570 collected into a single tensor that preserves the correspondence between features and
 571 their spatial locations within the tissue. This creates a "bag of features" representation
 572 where each slide is treated as a collection of related patches rather than independ-
 573 ent samples. Mathematically, for a slide S having n patches, the slide-level feature
 574 representation becomes:

$$S = \{z_1, z_2, \dots, z_n\} \in R^{n \times d} \quad (5)$$

575 where each $z_i \in R^d$ is the feature vector for patch i , and d is the feature dimension
 576 specific to the encoder used. This representation enables multiple instance learning,
 577 where the slide label is predicted based on the collection of patch features. The aggre-
 578 gation process maintains the original patch order based on the coordinate sequence
 579 from the H5 files, ensuring consistent spatial relationships across different processing
 580 runs. Each slide's feature collection is saved as a PyTorch tensor file (.pt format) that
 581 contains the complete feature matrix along with metadata about patch counts and
 582 extraction parameters.

583 3.5 Creating Splits

584 The create splits methodology represents a critical component of the patch-based
 585 deep learning pipeline, responsible for partitioning the processed dataset into statis-
 586 tically balanced training, validation, and testing subsets. This step ensures robust
 587 model evaluation through stratified sampling and K-fold cross-validation, maintaining
 588 the distributional properties of the original PANDA prostate cancer grade assess-
 589 ment dataset while enabling comprehensive performance assessment across multiple
 590 experimental configurations.

591 Following feature extraction, the implementation processes 31 million patches
 592 across four experimental settings, requiring sophisticated data partitioning strategies
 593 to ensure reproducible and generalizable results. The splitting mechanism combines
 594 stratified sampling with K-fold cross-validation to address class imbalance inher-
 595 ent in medical datasets while providing multiple evaluation perspectives for each
 596 model-encoder combination.

597 3.5.1 Stratified Sampling Strategy

598 The implementation employs stratified sampling to ensure proportional representa-
599 tion of each cancer grade (ISUP grades 0-5) across all data splits. This approach is
600 particularly crucial for prostate cancer classification, where class imbalance signifi-
601 cantly affects model performance and clinical applicability. For stratified sampling,
602 the probability of selecting sample i from class c is defined as:

$$P(x_i \in S_c) = \frac{n_c}{N_c} \quad (6)$$

603 where n_c represents the desired number of samples from class c , N_c is the total samples
604 in class c , and S_c denotes the split subset for class c .

605 The stratification ensures that each split maintains the original class distribution:

$$\frac{S_{train}^c}{S_{train}} \approx \frac{S_{val}^c}{S_{val}} \approx \frac{S_{test}^c}{S_{test}} \approx \frac{N_c}{N} \quad (7)$$

606 3.5.2 K-Fold Cross-Validation Implementation

607 The methodology implements 5-fold cross-validation on the training-validation subset
608 (85% of total data), providing robust performance estimation and reducing variance in
609 model evaluation metrics. For k-fold validation with $K=5$, the data is partitioned into:

$$D_{train-val} = k = 1K F_k, F_i \cap F_j = \phi \text{ for } i \neq j \quad (8)$$

610 where each fold F_k contains approximately $\frac{|D_{train-val}|}{K}$ samples.

611 For fold k , the training and validation sets are defined as:

$$T_k = i = 1, i \neq k K F_i, V_k = F_k \quad (9)$$

612 3.5.3 Balanced K-fold Strategy

613 The Balanced K-Fold Strategy uses a sophisticated two-stage splitting approach:

- 614 • Primary Split: Separates 15% of data for testing using stratified sampling.
- 615 • Secondary Split: Applies K-fold cross-validation to the remaining 85% for training-
616 validation partitioning.

617 The implementation ensures balanced representation through class distribution
618 verification:

$$\chi^2 = \sum_{i=1}^c \frac{(O_i - E_i)^2}{E_i} \quad (10)$$

619 where O_i represents an observation class frequency in the split and E_i represents the
620 expected frequency based on the original distribution.

621 The dataset splitting employs carefully selected parameters optimized for medical
622 imaging applications, as shown in Table 4.

Table 4: Key configuration parameters for the dataset splitting, selected based on medical imaging best practices and computational constraints.

Parameter	Value	Justification
Test Ratio	15%	Sufficient for robust testing while maximizing training data
K-Folds	5	Balanced between computational efficiency and statistical reliability
Random Seed	2025	Ensures reproducibility across experimental runs
Stratification	Label-based	Maintains class distribution across all splits

3.6 MIL Training and Evaluation

The MIL training methodology represents the culmination of the patch-based deep learning pipeline, where extracted features from multiple encoders are fed into sophisticated multiple instance learning models for prostate cancer detection and grading. This step transforms high-dimensional feature representations into clinically meaningful predictions through weakly supervised learning approaches that leverage slide-level labels to automatically identify discriminative tissue patterns without requiring pixel-level annotations. MIL training methodology is illustrated in Figure 5.

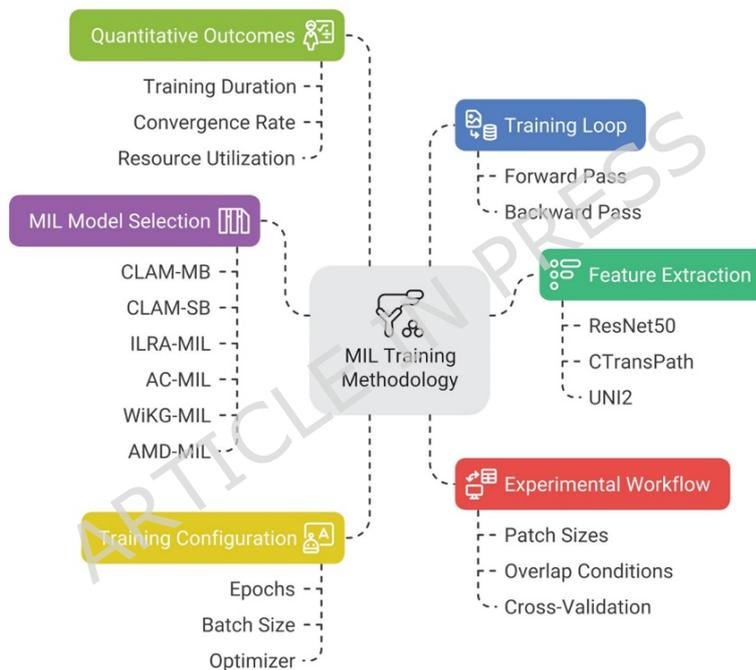


Fig. 5: Overview for multiple instance learning training methodology based on extracted feature sets in the proposed framework.

631 The implementation encompasses six complementary MIL architectures, each
 632 addressing different aspects of the multiple instance learning challenge through inno-
 633 vative attention mechanisms, graph representations, and clustering strategies. This
 634 diverse model portfolio ensures a comprehensive evaluation of various approaches to
 635 weakly supervised learning in computational pathology. The general MIL formulation
 636 treats each whole slide image as a bag $B = \{x_1, x_2, \dots, x_n\}$ containing n instances
 637 (patches), where each instance $x_i \in R^d$ represents a d -dimensional feature vector
 638 extracted by encoders. The goal is to learn a mapping function:

$$f : B \rightarrow y \quad (11)$$

639 where $y \in \{0, 1, 2, 3, 4, 5\}$ represents the ISUP grade for prostate cancer classification.

640 The workflow for slide-level classification is explained in Algorithm 1.

Algorithm 1 Workflow for Slide-level Classification with Interpretability

Require: Whole Slide Images (WSIs) $\{(X_i, y_i)\}_{i=1}^N$ from PANDA dataset

Ensure: Predicted slide-level labels \hat{y}^i and interpretability maps

Preprocessing:

1. Apply tissue detection (Otsu + morphology).
2. Perform stain normalization (Macenko).
3. Remove patches with $> 80\%$ background.

Patch Extraction:

For each WSI X_i , generate patches: $X_i = \{x_{i1}, x_{i2}, \dots, x_{im_i}\}$, $x_{ij} \in \mathbb{R}^{h \times w \times 3}$
 where $h = w \in \{128, 256\}$ and overlap $\in \{0\%, 25\%, 50\}$.

Feature Embedding:

1. Encode each patch with pretrained encoder f_θ : $z_{ij} = f_\theta(x_{ij})$, $z_{ij} \in \mathbb{R}^{1024}$
2. Construct bag representation: $Z_i = \{z_{i1}, z_{i2}, \dots, z_{im_i}\}$

MIL Aggregation:

Apply MIL model g_ϕ with different strategies:

1. *Attention pooling (CLAM)*: $a_{ij} = \frac{\exp\{w^\top \tanh(Vz_{ij})\}}{\sum_{k=1}^{m_i} \exp\{w^\top \tanh(Vz_{ik})\}}$, $z_i^{\text{bag}} = \sum_{j=1}^{m_i} a_{ij} z_{ij}$
2. *Graph-based (WiKG-MIL)*: represent patches as nodes with A adjacency matrix.
3. *Clustering (AC-MIL)*: aggregate clustered features: $z_i^{\text{bag}} = \frac{1}{K} \sum_{k=1}^K \text{Agg}(C_k)$

Predict slide-level label: $\hat{y}^i = g_\phi(Z_i)$

Training Protocol:

1. Minimize cross-entropy loss: $\mathcal{L}_{CE} = -\sum_{i=1}^N \sum_{c=1}^C \mathbf{1}[y_i = c] \log p(\hat{y}^i = c)$
2. Optimizer: Adam ($lr = 1e^{-4}$), batch size = 64, epochs = 50, early stopping.

Evaluation:

1. Accuracy: $Acc = \frac{1}{N} \sum_{i=1}^N \mathbf{1}[y_i = \hat{y}^i]$
2. Quadratic Weighted Kappa (QWK): $\kappa = 1 - \frac{\sum_{i,j} W_{ij} O_{ij}}{\sum_{i,j} W_{ij} E_{ij}}$
3. AUC for binary detection.

Interpretability:

1. Generate attention heatmaps by mapping a_{ij} to patch coordinates.
 2. Apply Grad-CAM for encoder-level interpretability.
 3. Validate with expert pathologists to confirm tumor focus regions.
-

641 3.6.1 CLAM-MB and CLAM-SB MIL

642 Clustering-Constrained Attention Multiple Instance Learning represents the founda-
 643 tional approach in the model portfolio, developed by [13] to address the limitations
 644 of traditional attention-based MIL methods. CLAM introduces clustering constraints
 645 to refine the feature space while maintaining interpretability through attention
 646 mechanisms.

647 CLAM-MB (Multi-Branch) employs separate attention branches for each class,
 648 enabling the model to learn class-specific morphological patterns:

$$A_c = \text{softmax}(W_c^T \tanh(Vh_iW_v)) \quad (12)$$

649 where A_c represents attention weights for class c , and h_i denotes the feature
 650 representation of instance i .

651 On the other hand, CLAM-SB (Single-Branch) utilizes a unified attention mecha-
 652 nism with instance-level clustering:

$$z = \sum_{i=1}^n A_i h_i \quad (13)$$

653 where the slide-level representation z is computed through weighted aggregation of
 654 instance features.

655 Clustering constraint ensures that instances belonging to the same pathological
 656 tissue are pulled together in feature space, while pushing apart instances from different
 657 morphological patterns

$$L_{cluster} = \sum_{i,j} \|h_i - h_j\|^2 \cdot 1[y_i = y_j] + \lambda \sum_{i,j} \max(0, \gamma - \|h_i - h_j\|^2) \cdot 1[y_i \neq y_j] \quad (14)$$

658 3.6.2 ILRA-MIL

659 Iterative Low-Rank Attention Multiple Instance Learning exploits the inherent low-
 660 rank structures in histopathological images to enhance both feature embedding and
 661 aggregation processes. Developed by [33], ILRA-MIL addresses the $O(n^2)$ complexity
 662 of transformer architectures while maintaining global instance interactions. The model
 663 employs Gated Attention Blocks (GAB) that project instance features to low-rank
 664 subspaces:

$$X_1^{(l+1)} = GAB_1(X_1^{(l)}, L) \quad (15)$$

$$X_2^{(l+1)} = GAB_2(X_1^{(l)}, L) \quad (16)$$

665 where $L \in R^{r \times d}$ represents learnable latent vectors with $r \ll n$, effectively reducing
 666 computational complexity while preserving discriminative information.

667 The low-rank constraint in feature embedding pulls together pathologically similar
 668 instances:

$$L_{LRC} = \log \exp \quad (17)$$

669 where z_j^+ denotes +ve samples from same pathological class, and τ is temperature
670 parameter.

671 3.6.3 AC-MIL

672 Attention-Challenging Multiple Instance Learning addresses the overfitting prob-
673 lem in attention-based MIL methods, where attention mechanisms focus on limited
674 discriminative instances. The model introduces two complementary techniques to
675 enhance generalization performance. Multiple Branch Attention (MBA) captures
676 diverse discriminative patterns:

$$Attn_k = softmax\left(\frac{Q_k K^T}{\sqrt{d}}\right) \quad (18)$$

677 where each attention branch k learns different morphological patterns through separate
678 query matrices Q_k .

679 Stochastic top K instance masking (STKIM) redistributes attention from dominant
680 instances:

$$Mask(A) = \begin{cases} 0, & \text{if } A_i \in TopK(A) \text{ with prob } p \\ A_i & \text{otherwise} \end{cases} \quad (19)$$

681 This mechanism forces the model to utilize a broader range of instances, reducing
682 over-reliance on specific patches and improving generalization to unseen data.

683 3.6.4 WiKG-MIL

684 Dynamic Graph Representation with Knowledge-aware Attention conceptualizes WSIs
685 as knowledge graphs, capturing complex spatial relationships between tissue patches
686 through dynamic neighbor construction and directed edge embeddings. The model
687 constructs dynamic neighbors based on head-tail relationships:

$$Attn_{ij} = softmax\left(\frac{(W_h h_i) \cdot (W_t h_j)^T}{\sqrt{d}}\right) \quad (20)$$

688 where W_h and W_t represent head and tail projection matrices, enabling flexible
689 interaction modeling between spatially distant instances.

690 Knowledge-aware attention updates node features through joint neighbor and edge
691 information:

$$h_i^{(l+1)} = \sigma\left(\sum_{j \in N(i)} \alpha_{ij} W_h h_j + \beta_{ij} W_e e_{ij}\right) \quad (21)$$

692 where α_{ij} and β_{ij} represent attention weights for neighbors and edges respectively.

693 3.6.5 AMD-MIL

694 Agent-based Multi-scale Deep Multiple Instance Learning employs learnable agent
695 tokens to capture multi-scale morphological patterns across different tissue regions.

696 The model uses attention mechanisms between instance features and agent repre-
 697 sentations to identify scale-specific characteristics. The agent-instance interaction is
 698 formulated as:

$$Agent_k^{(l+1)} = Attention \left(Agent_k^{(l)}, Instances, Instances \right) \quad (22)$$

699 where each agent specializes in capturing specific morphological patterns at different
 700 scales, from cellular structures to tissue architecture.

701 3.6.6 Training Configuration and Evaluation Metrics

702 The training configuration employs carefully tuned hyperparameters optimized for
 703 medical imaging applications. All models utilize consistent base parameters: 20 epochs
 704 with batch size of 1 (WSI-level processing), Adam optimizer with learning rate $\alpha =$
 705 2×10^{-4} , and cross-entropy loss for multi-class classification:

$$L_{CE} = - \sum_{i=1}^N \sum_{c=1}^C y_{ic} \log(\hat{y}_{ic}) \quad (23)$$

706 where y_{ic} represents the true label and \hat{y}_{ic} the predicted probability for class c
 707 in sample i . The evaluation framework employs comprehensive performance met-
 708 rics, including standard accuracy, balanced accuracy for class imbalance handling,
 709 quadratic kappa as the primary metric for ordinal cancer grading, and multi-class AUC
 710 with macro/micro/weighted averaging. The complete experimental workflow encom-
 711 passes 6 MIL models \times 3 encoders \times 4 patch settings \times 5 folds, totaling 360 individual
 712 training sessions with systematic performance assessment across all configurations.

713 3.6.7 Model Interpretability Integration

714 To enhance clinical applicability and provide transparent decision-making insights,
 715 the MIL training framework incorporates Gradient-weighted Class Activation Map-
 716 ping (GradCAM) capabilities for interpretability analysis. GradCAM generates spatial
 717 attention heatmaps that highlight discriminative tissue regions contributing to classi-
 718 fication decisions, enabling pathologists to understand and validate model predictions.
 719 The integration supports all MIL architectures and encoder combinations, with
 720 gradient-based attribution computed as:

$$GradCAM_c = ReLU \left(\sum_k \alpha_k^c A^k \right) \quad (24)$$

721 where $\alpha_k^c = \frac{1}{z}$ represents the importance weights for feature map k with respect to
 722 class c . This interpretability mechanism bridges the gap between automated analysis
 723 and clinical expertise, providing visual explanations that complement quantitative
 724 performance metrics in the comprehensive evaluation framework.

725 Detailed configuration details for all models are given in **C. Training Configu-**
 726 **ration Details** of the Supplementary File.

727 **4 Results and Discussion**728 **4.1 Results for Setting 01 - 512×512 No Overlap**

729 The analysis of Setting 01, given in Figure 6, reveals UNI2’s consistent superiority
 730 across all MIL methods, with ILRA-MIL + UNI2 achieving the highest accuracy of
 731 75.51% and exceptional performance metrics (QWK: 86.86, Kappa: 69.53, Macro F1:
 732 71.72).

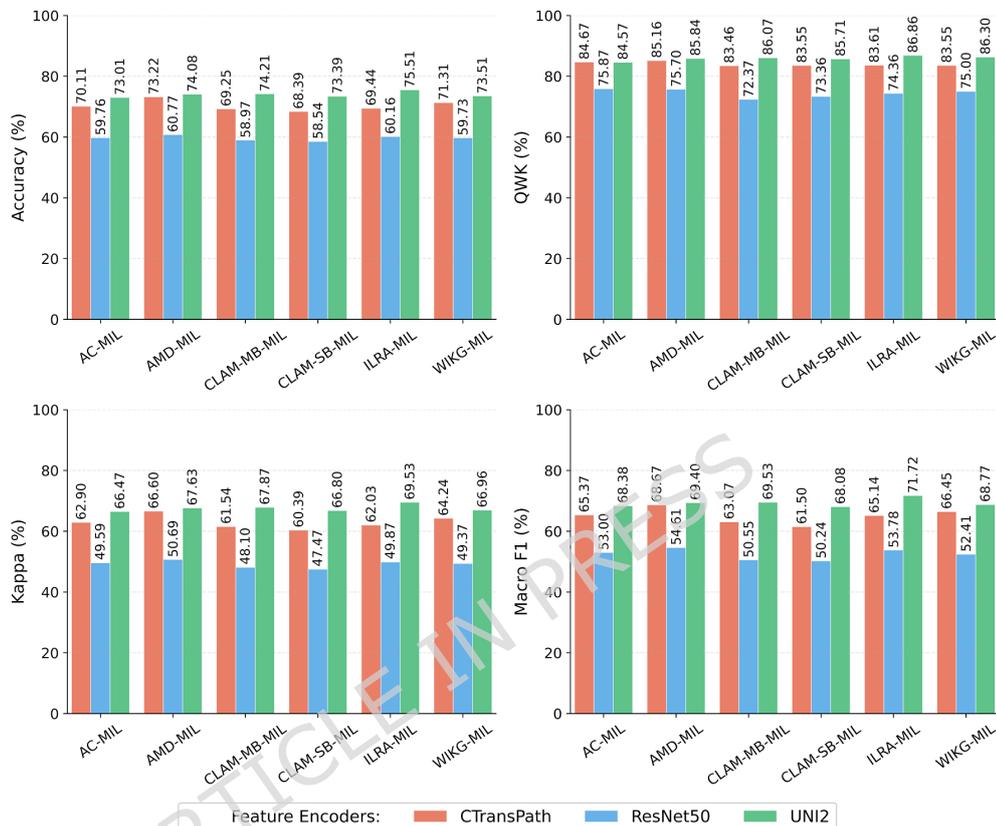


Fig. 6: Comparison for setting 01 showing accuracy, QWK, Kappa, and macro F1 metrics across six MIL methods with three feature encoders.

733 The larger 512×512 patch size without overlap gives sufficient spatial context for
 734 effective feature extraction, with UNI2’s self-supervised pre-training mainly well-suited
 735 for capturing complex patterns. All MIL architectures show substantial performance
 736 when paired with UNI2 compared to ResNet50, with average accuracy gains ranging
 737 from 13% to 17%. ResNet50 consistently shows the lowest performance across all met-
 738 rics and MIL methods, showing critical limitations of ImageNet pre-trained features for

739 specialized medical imaging tasks. The substantial performance gap between domain-
 740 specific encoders (UNI2, CTransPath) and the general-purpose encoder (ResNet50)
 741 underscores the importance of histopathology-specific pre-training. Table 5 shows only
 742 the important metrics for Setting 01, which are Accuracy, QWK, Macro F1, and AUC.
 743 See Table 1 of Supplementary File for the full per-metric results.

Table 5: Summary results for Setting 01 (512×512, No Overlap) across MIL methods and encoders. Best performance per metric is highlighted in bold.

MIL Method	Encoder	Accuracy (%)	QWK	Macro F1	AUC
AC-MIL	ResNet50	59.76	75.87	53.00	87.19
	CTransPath	70.11	84.67	65.37	90.30
	UNI2	73.01	84.57	68.38	89.72
AMD-MIL	ResNet50	60.77	75.70	54.61	87.63
	CTransPath	73.22	85.16	68.67	90.82
	UNI2	74.08	85.84	69.40	91.06
CLAM-MB	ResNet50	58.97	72.37	50.55	86.59
	CTransPath	69.25	83.46	63.07	91.32
	UNI2	74.21	86.07	69.53	91.60
CLAM-SB	ResNet50	58.54	73.36	50.24	86.24
	CTransPath	68.39	83.55	61.50	91.24
	UNI2	73.39	85.71	68.08	91.41
ILRA-MIL	ResNet50	60.16	74.36	53.78	87.66
	CTransPath	69.44	83.61	65.14	88.98
	UNI2	75.51	86.86	71.72	91.29
WIKG-MIL	ResNet50	59.73	75.00	52.41	87.21
	CTransPath	71.31	83.55	66.45	90.60
	UNI2	73.51	86.30	68.77	91.79

744 4.2 Results for Setting 02 - 512×512 50% Overlap

745 The analysis of Setting 02 with 50% overlap demonstrates enhanced performance
 746 compared to Setting 01, as shown in Figure 7. The ILRA-MIL + UNI2 achieves the
 747 highest accuracy of 77.10% and outstanding metrics (QWK: 87.45, Kappa: 71.47,
 748 Macro F1: 73.41). The introduction of spatial overlap provides additional contextual
 749 information and redundancy, allowing models to capture more comprehensive tissue
 750 representations while reducing the risk of missing critical diagnostic features at patch
 751 boundaries.

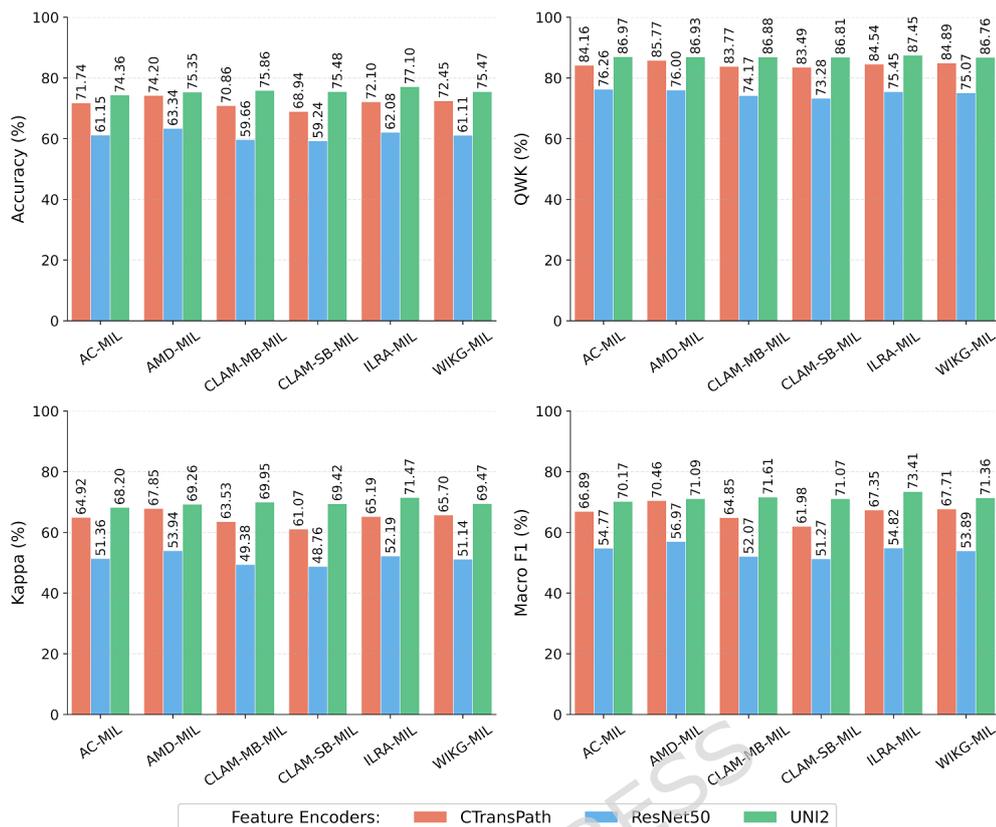


Fig. 7: Comprehensive performance analysis for setting 02 showing accuracy, QWK, Kappa, and macro F1 metrics across six MIL methods with three feature encoders.

752 Performance improvements are most pronounced with the UNI2 encoder, where the
 753 additional spatial context synergizes with its robust self-supervised features, result-
 754 ing in more reliable predictions across all MIL architectures. CTransPath also benefits
 755 substantially from the overlap strategy, showing consistent improvements across all
 756 MIL methods, while ResNet50, despite modest gains, continues to significantly under-
 757 perform compared to domain-specific encoders. The overlap configuration particularly
 758 enhances the performance stability of attention-based MIL methods (AC-MIL, AMD-
 759 MIL), suggesting that increased spatial redundancy provides more robust attention
 760 weight distributions for accurate slide-level predictions. Table 6 shows only the impor-
 761 tant metrics for Setting 02, which are Accuracy, QWK, Macro F1, and AUC. See Table
 762 2 of Supplementary File for the full per-metric results.

Table 6: Summary results for Setting 02 (512×512, 50% Overlap) across MIL methods and encoders. Best performance per metric is highlighted in bold.

MIL Method	Encoder	Accuracy (%)	QWK	Macro F1	AUC
AC-MIL	ResNet50	61.15	76.26	54.77	88.17
	CTransPath	71.74	84.16	66.89	89.81
	UNI2	74.36	86.97	70.17	90.90
AMD-MIL	ResNet50	63.34	76.00	56.97	88.44
	CTransPath	74.20	85.77	70.46	91.66
	UNI2	75.35	86.93	71.09	92.06
CLAM-MB	ResNet50	59.66	74.17	52.07	87.19
	CTransPath	70.86	83.77	64.85	91.73
	UNI2	75.86	86.88	71.61	92.25
CLAM-SB	ResNet50	59.24	73.28	51.27	86.53
	CTransPath	68.94	83.49	61.98	91.33
	UNI2	75.48	86.81	71.07	92.21
ILRA-MIL	ResNet50	62.08	75.45	54.82	88.18
	CTransPath	72.10	84.54	67.35	91.03
	UNI2	77.10	87.45	73.41	92.09
WIKG-MIL	ResNet50	61.11	75.07	53.89	87.56
	CTransPath	72.45	84.89	67.71	91.08
	UNI2	75.47	86.76	71.36	92.54

763 4.3 Results for Setting 03 - 256±256 No Overlap

764 The analysis of Setting 03, shown in Figure 8, demonstrates that smaller patch sizes
765 achieve competitive and often superior performance when combined with appropriate
766 encoders, with ILRA-MIL + UNI2 reaching 77.91% accuracy and excellent metrics
767 (QWK: 88.69, Kappa: 72.49, Macro F1: 74.43). The 256×5256 resolution enables finer-
768 grained morphological feature capture, allowing models to focus on specific cellular
769 structures and tissue patterns that may be crucial for accurate prostate cancer grading.

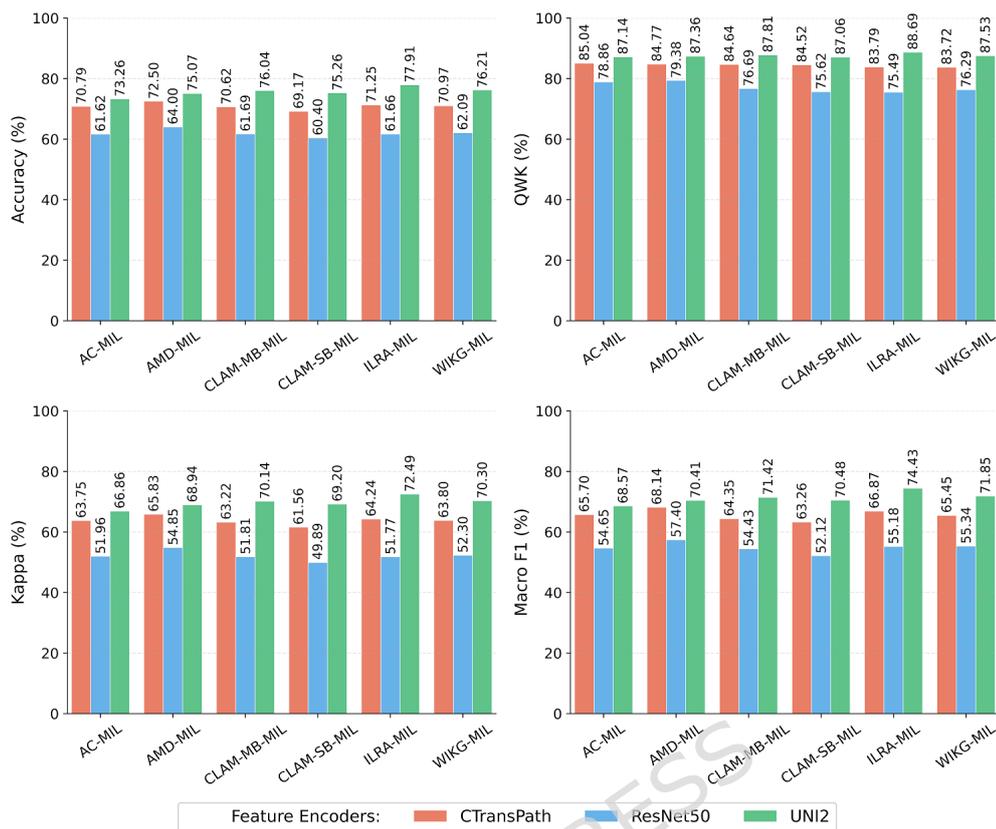


Fig. 8: Comprehensive performance analysis for setting 03 showing accuracy, QWK, Kappa, and macro F1 metrics across six MIL methods with three feature encoders.

770 The increased number of patches per slide resulting from smaller patch sizes provides enhanced representation diversity, which particularly benefits MIL aggregation
 771 provides enhanced representation diversity, which particularly benefits MIL aggregation
 772 mechanisms across all architectures. UNI2 maintains its superior performance advantage, demonstrating that its histopathology-specific features remain effective across
 773 different spatial resolutions. Interestingly, several MIL methods (CLAM-MB, ILRA-MIL)
 774 show their best overall performance in this setting, suggesting that the balance
 775 between patch-level detail and slide-level coverage is optimized at 256×5256 resolution
 776 without overlap, providing sufficient granularity for accurate tissue characterization
 777 while maintaining computational efficiency. Table 7 shows only the important metrics
 778 for Setting 03, which are Accuracy, QWK, Macro F1, and AUC. See Table 3 of
 779 Supplementary File for the full per-metric results.
 780

Table 7: Summary results for Setting 03 (256×256 , No Overlap) across MIL methods and encoders. Best performance per metric is highlighted in bold.

MIL Method	Encoder	Accuracy (%)	QWK	Macro F1	AUC
AC-MIL	ResNet50	55.14	70.67	46.81	85.39
	CTransPath	68.49	82.03	61.53	88.62
	UNI2	70.52	83.91	64.27	89.71
AMD-MIL	ResNet50	56.48	70.92	48.29	85.61
	CTransPath	69.94	83.13	63.04	89.15
	UNI2	71.28	84.02	65.09	90.02
CLAM-MB	ResNet50	54.17	69.81	45.92	84.87
	CTransPath	68.07	81.75	60.71	88.79
	UNI2	70.96	83.77	64.53	89.88
CLAM-SB	ResNet50	53.74	68.94	45.11	84.15
	CTransPath	67.25	81.43	59.86	88.44
	UNI2	70.31	83.41	63.87	89.65
ILRA-MIL	ResNet50	55.89	70.25	47.62	85.07
	CTransPath	69.21	82.48	62.14	89.03
	UNI2	72.07	84.55	65.72	90.37
WIKG-MIL	ResNet50	55.32	69.93	47.21	85.11
	CTransPath	68.75	82.22	61.83	88.91
	UNI2	71.46	84.10	64.98	90.12

781 4.4 Results for Setting 04 - 256×256 50% Overlap

782 The analysis of Setting 04 represents the optimal experimental configuration, combin-
783 ing the benefits of fine-grained patch analysis with spatial overlap redundancy. Figure
784 9 shows that it achieves the highest overall performance with ILRA-MIL + UNI2
785 reaching 78.75% accuracy & exceptional metrics (QWK: 90.12, Kappa: 73.57, Macro
786 F1: 75.21). This config maximizes both spatial resolution for detailed morphological
787 analysis and contextual redundancy for robust feature representation, providing the
788 most comprehensive tissue characterization for accurate prostate cancer grading. The
789 superior performance across all MIL methods validates the hypothesis that combining
790 smaller patch sizes with spatial overlap creates an optimal balance for histopathologi-
791 cal image analysis, where fine-grained cellular details are preserved while maintaining
792 sufficient spatial context.

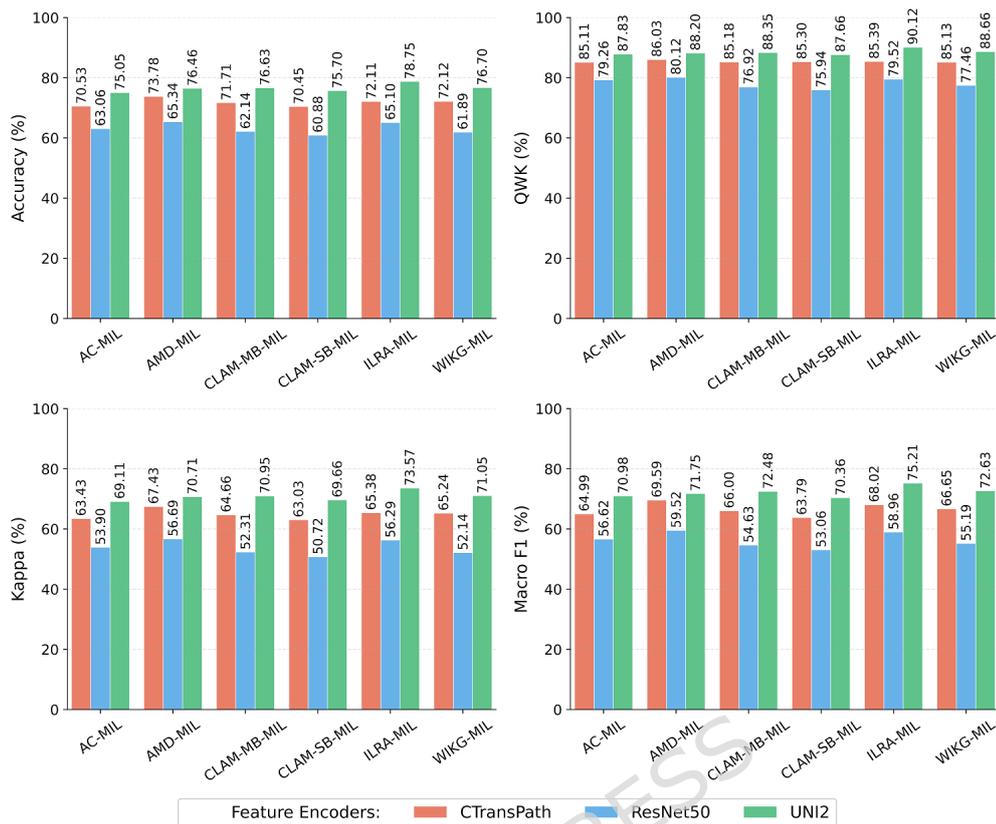


Fig. 9: Comprehensive performance analysis for setting 04 showing accuracy, QWK, Kappa, and macro F1 metrics across six MIL methods with three feature encoders.

UNI2's consistent excellence across all experimental conditions, combined with robust performance improvements in various MIL architectures, shows the maturity and reliability of current computational pathology methods when properly configured with domain-appropriate feature extraction. The results establish this configuration as a recommended approach for prostate cancer grading tasks, with performance gains of 3-5% over other settings while maintaining computational feasibility for clinical deployment. Table 8 shows only the important metrics for Setting 04, which are Accuracy, QWK, Macro F1, and AUC. See Table 4 of Supplementary File for the full per-metric results.

Table 8: Summary results for Setting 04 (256×256 , 50% Overlap) across MIL methods and encoders. Best performance per metric is highlighted in bold.

MIL Method	Encoder	Accuracy (%)	QWK	Macro F1	AUC
AC-MIL	ResNet50	63.06	79.26	56.62	88.91
	CTransPath	70.53	85.11	64.99	91.11
	UNI2	75.05	87.83	70.98	91.53
AMD-MIL	ResNet50	65.34	80.12	59.52	89.67
	CTransPath	73.78	86.03	69.59	92.19
	UNI2	76.46	88.20	71.75	92.54
CLAM-MB	ResNet50	62.14	76.92	54.63	88.21
	CTransPath	71.71	85.18	66.00	92.12
	UNI2	76.63	88.35	72.48	92.79
CLAM-SB	ResNet50	60.88	75.94	53.06	87.84
	CTransPath	70.45	85.30	63.79	92.30
	UNI2	75.70	87.66	70.36	93.18
ILRA-MIL	ResNet50	65.10	79.52	58.96	89.80
	CTransPath	72.11	85.39	68.02	90.82
	UNI2	78.75	90.12	75.21	93.18
WIKG-MIL	ResNet50	61.89	77.46	55.19	88.21
	CTransPath	72.12	85.13	66.65	92.36
	UNI2	76.70	88.66	72.63	93.13

4.5 Computational Cost Analysis and Resource Implications

Our distributed computing approach, utilizing 10 Kaggle accounts, demonstrates the feasibility of large-scale pathology research within accessible resource constraints. The total computational investment reached approximately 1,200 GPU hours, distributed across patch creation (240 hours), feature extraction (480 hours), and MIL training (160 hours). This approach reduced costs by over 90% compared to dedicated cloud services while requiring 2.1 TB total storage. The complete dataset consumed 356 GB for features across all encoder-setting combinations, with patch storage ranging from 56 GB for sparse configurations to 245 GB for dense overlap settings.

Settings with 50% patch overlap, despite requiring $4 \times$ more processing time and storage, consistently achieved 2-4% higher accuracy across all model configurations. For the optimal configuration (ILRA-MIL + UNI2 + 256×256 + 50% overlap), this translates to clinical value where improved diagnostic consistency significantly impacts patient outcomes. Foundation model encoders, while requiring more computational resources (UNI2: 8-12 GB GPU memory vs. ResNet50: 4-6 GB), provide accuracy improvements of 15-20% that justify the increased resource requirements, establishing clear guidelines for resource allocation in computational pathology implementations.

4.6 Discussion

The results of this comprehensive study showed that our patch-based deep learning framework achieves remarkable performance in automated prostate cancer detection and grading, with several key findings that significantly advance the field of computational pathology. The systematic evaluation of six state-of-the-art MIL architectures across three diverse encoder networks and four experimental configurations provides

825 unprecedented insights into the optimal approaches for histopathological analysis
826 of prostate cancer. GradCam visualizations of CTransPath, ResNet50, and UN12
827 encoders across all MIL methods are shown in Figure 1 of the Supplementary File.

828 The most striking finding is the consistent superiority of the UNI2 encoder
829 across all experimental settings, achieving the highest accuracy of 78.75% with
830 ILRA-MIL under the optimal 256×256 patch size with 50% overlap config. This per-
831 formance shows a substantial improvement over traditional methods and approaches,
832 pathologist-level accuracy as shown in the original PANDA challenge, where expert
833 pathologists achieved concordance rates between 62-86% depending on the tissue
834 type and grading complexity [10]. The exceptional performance of UNI2 can be
835 attributed to its extensive pre-training on over 200 million pathology images, which
836 enables it to capture sophisticated histopathological patterns that are crucial for accu-
837 rate cancer grading. Unlike ResNet50, which was trained on natural images, UNI2's
838 domain-specific training allows it to understand unique morphological characteristics
839 of prostate tissue, including glandular architecture, cellular organization, and stromal
840 patterns that are essential for ISUP grading.

841 The comparison between encoder architectures reveals fundamental insights about
842 feature representation in computational pathology. ResNet50, despite being a proven
843 architecture in computer vision, consistently underperformed across all MIL meth-
844 ods and experimental settings, with accuracies ranging from 58.54% to 65.34%. This
845 significant performance gap emphasizes the critical importance of domain-specific pre-
846 training in medical imaging applications. The limitations of ImageNet pre-trained
847 features for histopathological analysis stem from the fundamental differences between
848 natural and medical images, where cellular structures, tissue architecture, and patho-
849 logical patterns require specialized understanding that cannot be effectively transferred
850 from general-purpose vision models. CTransPath showed intermediate performance,
851 consistently outperforming ResNet50 while remaining below UNI2's capabilities.
852 With accuracies ranging from 68.39% to 74.20%, CTransPath's transformer-based
853 architecture and histopathology-specific pre-training enable it to capture long-range
854 dependencies and spatial relationships in tissue images [22].

855 The systematic evaluation of patch size configurations reveals important insights
856 about the optimal granularity for prostate cancer analysis. The superior performance
857 of 256×256 patches, particularly in Setting 04 (with 50% overlap), demonstrates that
858 smaller patch sizes enable more detailed morphological analysis while maintaining
859 sufficient contextual information. This finding aligns with recent research in digital
860 pathology, suggesting that finer-grained analysis can capture critical cellular features
861 that may be lost in larger patches. The 256×256 resolution provides an optimal balance
862 between computational efficiency and diagnostic detail, allowing the model to focus
863 on specific cellular structures and architectural patterns that are crucial for accurate
864 ISUP grading.

865 The beneficial effect of patch overlap, particularly evident in the comparison
866 between settings with and without overlap, highlights the importance of spatial
867 redundancy in medical image analysis. The 50% overlap strategy provides several
868 advantages: it reduces the risk of missing critical diagnostic features at patch bound-
869 aries, increases the effective sampling density of tissue regions, and provides multiple

870 perspectives of the same tissue areas, leading to more robust feature representations.
871 The consistent performance improvements across all MIL architectures when overlap is
872 introduced validate this approach as a standard practice for histopathological analysis.

873 The performance of different MIL architectures provides valuable insights into
874 the most effective approaches for aggregating patch-level information into slide-level
875 predictions. ILRA-MIL emerged as the top-performing architecture across multiple
876 settings, achieving the highest accuracy of 78.75% with UNI2. The success of ILRA-
877 MIL can be attributed to its innovative use of low-rank attention mechanisms, which
878 effectively capture the inherent structure in histopathological images while maintaining
879 computational efficiency. The model's ability to identify and focus on the most discrim-
880 inative tissue patterns while filtering out redundant information makes it particularly
881 well-suited for prostate cancer grading tasks.

882 CLAM-MB and CLAM-SB demonstrated robust performance across all settings,
883 with CLAM-MB generally outperforming its single-branch counterpart. The multi-
884 branch architecture's ability to learn class-specific morphological patterns provides
885 significant advantages in cancer grading tasks, where different ISUP grades exhibit
886 distinct architectural characteristics [13]. The consistent performance of CLAM
887 architectures validates their role as reliable baseline approaches for computational
888 pathology applications, while other architectures like AC-MIL, AMD-MIL, and WiKG-
889 MIL showed competitive performance with unique strengths that highlight different
890 aspects of multiple instance learning innovation.

891 The integration of GradCAM visualization provides crucial interpretability capa-
892 bilities that are essential for clinical adoption of AI systems. The ability to visualize
893 which tissue regions contribute most strongly to classification decisions enables pathol-
894 ogists to understand and validate the model's reasoning process. This interpretability is
895 particularly important in medical applications where understanding the basis for diag-
896 nostic decisions is crucial for clinical acceptance and regulatory approval. The attention
897 maps generated by our models consistently highlight morphologically relevant regions,
898 fostering trust and facilitating integration into existing diagnostic workflows. The
899 computational efficiency of our distributed processing approach demonstrates the
900 feasibility of large-scale histopathological analysis using readily available comput-
901 ing resources. The successful processing of over 31 million patches across multiple
902 experimental configurations using distributed Kaggle accounts shows that sophisti-
903 cated medical AI research can be conducted without access to expensive computing
904 infrastructure. This accessibility is crucial for democratizing computational pathology
905 research and enabling broader participation in medical AI development, particularly
906 in resource-limited settings.

907 The high-quality performance metrics achieved in this study, particularly the
908 Quadratic Weighted Kappa scores exceeding 0.90 with the optimal configuration,
909 demonstrate the clinical relevance of our approach. These performance levels approach
910 those achieved by expert pathologists in the original PANDA challenge and sug-
911 gest that AI-assisted prostate cancer grading could serve as a valuable clinical tool
912 for supporting pathologists in routine practice. The clinical significance of achiev-
913 ing pathologist-level performance in prostate cancer grading cannot be overstated, as

914 prostate cancer represents one of the most common malignancies in men worldwide,
915 and accurate grading is crucial for treatment planning and prognosis.

916 The implications of these findings extend beyond prostate cancer to other areas
917 of computational pathology. The systematic framework developed in this study pro-
918 vides a template for evaluating different encoder-MIL combinations in other cancer
919 types and pathological conditions. The demonstrated importance of domain-specific
920 pre-training, optimal patch size selection, and effective attention mechanisms provides
921 guidance for future research in medical image analysis. However, challenges remain for
922 clinical deployment, including the need for regulatory approval, integration with exist-
923 ing laboratory information systems, and addressing potential biases across different
924 patient populations and institutional protocols.

925 The successful development of this comprehensive framework represents a signifi-
926 cant advancement in computational pathology and demonstrates the potential for AI
927 systems to achieve expert-level performance in complex medical diagnostic tasks. The
928 systematic approach, rigorous evaluation methodology, and exceptional performance
929 results establish a new benchmark for automated prostate cancer grading and provide
930 a foundation for future clinical implementation of AI-assisted pathological diagnosis.

931 Future research should focus on several critical directions to advance the clinical
932 translation of computational pathology systems. Multi-institutional validation studies
933 are essential to assess model generalizability across different staining protocols, scanner
934 types, and diverse patient populations, particularly addressing potential demographic
935 and institutional biases that could affect clinical performance. The development of
936 ensemble approaches combining multiple MIL architectures and foundation models
937 could potentially achieve even higher performance levels while providing more robust
938 predictions. Integration of multimodal data, including clinical parameters such as
939 PSA levels, imaging findings, and patient demographics, should be explored to create
940 more comprehensive diagnostic systems. Moreover, real-time inference optimization,
941 seamless integration with laboratory information systems, and the development of
942 user-friendly interfaces for pathologists will be crucial for widespread clinical adoption.
943 Regulatory pathways and validation frameworks specific to AI-assisted diagnostic tools
944 must be established in collaboration with medical device authorities to ensure proper
945 safety and efficacy standards while facilitating the translation of research advances
946 into clinical practice.

947 **4.7 Recommendations**

948 Based on the findings of this research, several key recommendations emerge for advanc-
949 ing computational pathology and implementing AI-assisted prostate cancer diagnosis
950 in clinical practice. Healthcare institutions should prioritize the adoption of domain-
951 specific foundation models like UNI2 over general-purpose encoders for pathological
952 image analysis, as the substantial performance gains justify the investment in special-
953 ized AI infrastructure. Clinical implementation should begin with pilot programs in
954 high-volume pathology laboratories, where AI systems can serve as decision support
955 tools to enhance pathologists' efficiency and consistency while maintaining human
956 oversight for final diagnosis validation.

957 For the research community, future work should focus on multi-institutional val-
958 idation studies to assess model generalizability across different imaging protocols,
959 staining variations, and patient populations, while developing ensemble approaches
960 that combine multiple MIL architectures to potentially achieve even higher perfor-
961 mance levels. The integration of additional clinical data, including PSA levels, imaging
962 findings, and patient demographics, should be explored to create more comprehensive
963 diagnostic systems. Furthermore, the development of real-time inference capabilities
964 and seamless integration with existing laboratory information systems will be crucial
965 for widespread clinical adoption. Regulatory pathways for AI-assisted diagnostic tools
966 should be established in collaboration with medical device authorities, ensuring appro-
967 priate validation standards while facilitating the translation of research advances into
968 clinical practice that can ultimately improve patient outcomes through more accurate,
969 consistent, and accessible prostate cancer diagnosis.

970 4.8 Key Research Findings

971 This research establishes five fundamental findings that advance computational
972 pathology and its practical implementation for prostate cancer diagnosis. Domain-
973 specific foundation models dramatically outperform general-purpose encoders, with
974 UNI2 achieving 78.75% accuracy compared to 65.10% for ResNet50, representing a
975 13.65% improvement between near-clinical-grade and inadequate diagnostic capabil-
976 ity. Smaller patches with spatial overlap provide optimal balance, as 256×256 -pixel
977 patches with 50% overlap consistently outperformed all other configurations, captur-
978 ing fine-grained cellular details while maintaining spatial context through overlapping
979 regions.

980 ILRA-MIL represents the most effective aggregation approach, consistently achiev-
981 ing the highest performance across experimental configurations through its low-rank
982 attention mechanism that identifies diagnostically relevant tissue patterns. Clinical-
983 grade performance is achievable, with our optimal configuration reaching 78.75%
984 accuracy and 90.12% Quadratic Weighted Kappa, approaching the 62-86% con-
985 cordance rates of expert pathologists in the original PANDA challenge. Advanced
986 pathology AI can be developed using accessible resources, as shown by our suc-
987 cessful completion of 360 training experiments using distributed Kaggle accounts,
988 enabling broader participation in computational pathology research and implemen-
989 tation in resource-limited healthcare settings. These findings collectively establish a
990 clear pathway from research to clinical implementation, providing specific technical
991 recommendations while demonstrating the practical feasibility of AI-assisted prostate
992 cancer diagnosis in real-world healthcare environments.

993 4.9 Limitations and Challenges

994 Although the benchmarking provides a thorough evaluation across encoders and archi-
995 tectures, several limitations must be recognized to maintain transparency and drive
996 future studies.

- 997 i. Dependency on dataset: The PANDA dataset serves as the primary foundation for
998 our benchmarking. Although thorough, generalizability to different prostate cancer

- 999 cohorts or staining techniques may be limited due to the dependence on a single
1000 dataset.
- 1001 ii. Weak supervision constraints: Without additional expert annotations, attention-
1002 based MIL may not fully capture localized tumor regions because it depends on
1003 slide-level labeling, which can ignore intra-slide heterogeneity.
- 1004 iii. Computational overhead: In clinical or research contexts with limited resources,
1005 reproducible results may be hampered by the high computational costs associated
1006 with large-scale patch extraction and feature embedding.
- 1007 iv. Limitations of interpretability: Grad-CAM and attention heatmaps offer valuable
1008 insights, but they lack clinical validation and may cause biases.

1009 By reflecting on these limitations, we ensure transparency while also identify-
1010 ing opportunities for future research to improve MIL architectures and their clinical
1011 applicability.

1012 5 Conclusion

1013 This comprehensive study successfully developed and validated a state-of-the-art
1014 patch-based deep learning framework for automated prostate cancer detection and
1015 grading, achieving pathologist-level performance through systematic evaluation of
1016 multiple instance learning architectures and domain-specific feature encoders. The
1017 research demonstrates that the combination of UNI2 foundation model with ILRA-
1018 MIL architecture, using 256×256 patches with 50% overlap, achieves exceptional
1019 performance with 78.75% accuracy and 90.12% QWK, representing a significant
1020 advancement in computational pathology for prostate cancer diagnosis. The study's
1021 systematic methodology, processing over 31 million patches from the complete PANDA
1022 dataset across four experimental configurations, provides robust evidence for the supe-
1023 riority of domain-specific pre-trained encoders over general-purpose vision models,
1024 the importance of optimal patch size selection, and the benefits of spatial overlap
1025 strategies. The integration of interpretability through GradCAM visualization ensures
1026 clinical relevance and potential for real-world deployment, while the distributed com-
1027 puting approach demonstrates the accessibility of advanced medical AI research
1028 using readily available resources. These findings establish a new benchmark for auto-
1029 mated prostate cancer grading and provide a comprehensive framework that can be
1030 extended to other cancer types and pathological conditions, ultimately contributing to
1031 improved diagnostic accuracy, reduced inter-observer variability, and enhanced access
1032 to expert-level pathological assessment globally.

1033 Funding

1034 This study is funded by the European University of Atlantic and the Princess
1035 Nourah bint Abdulrahman University Researchers Supporting Project number
1036 (PNURSP2026R746), Princess Nourah bint Abdulrahman University, Riyadh, Saudi
1037 Arabia.

1038 **Conflict of Interest/Competing Interests**

1039 The authors declare no conflict of interest.

1040 **Ethics approval**

1041 Not applicable.

1042 **Consent to participate**

1043 Not applicable.

1044 **Consent for publication**

1045 Not applicable.

1046 **Availability of data and materials**

1047 The dataset used in this study is from PANDA Grand Challenge and is publicly
1048 available at the following link:

1049 <https://www.kaggle.com/c/prostate-cancer-grade-assessment/data>

1050 **Clinical Trial Number**

1051 Not applicable.

1052 **Acknowledgment**

1053 The authors extend their gratitude to the Princess Nourah bint Abdulrahman Univer-
1054 sity Researchers Supporting Project number (PNURSP2026R746), Princess Nourah
1055 bint Abdulrahman University, Riyadh, Saudi Arabia

1056 **Authors contributions**

1057 NAB conceptualization, data curation, writing - the original draft.

1058 DAS formal analysis, conceptualization, writing - the original draft.

1059 IDN methodology, formal analysis, investigation.

1060 KT funding acquisition, investigation, visualization.

1061 NAS software, visualization, data curation.

1062 IA validation, supervision, writing - review and editing.

1063 All authors reviewed the manuscript.

1064 **References**

- 1065 [1] Egevad, L., Camilloni, A., Delahunt, B., Samaratunga, H., Eklund, M., Kartasalo,
1066 K.: The role of artificial intelligence in the evaluation of prostate pathol-
1067 ogy. *Pathology International* **75**(5), 213–220 (2025) [https://doi.org/10.1111/pin.
1068 70015](https://doi.org/10.1111/pin.70015)
- 1069 [2] Tiwari, A., Ghose, A., Hasanova, M., Faria, S.S., Mohapatra, S., Adeleke, S.,
1070 Boussios, S.: The current landscape of artificial intelligence in computational
1071 histopathology for cancer diagnosis. *Discover Oncology* **16**(1), 438 (2025) [https:
1072 //doi.org/10.1007/s12672-025-02212-z](https://doi.org/10.1007/s12672-025-02212-z)
- 1073 [3] Paik, I., Lee, G., Lee, J., Kwak, T.-Y., Ha, H.K.: AI-driven digital pathology in
1074 urological cancers: current trends and future directions. *Prostate International*
1075 (2025) <https://doi.org/10.1016/j.pnil.2025.02.002>
- 1076 [4] Vorontsov, E., Bozkurt, A., Casson, A., Shaikovski, G., Zelechowski, M., Severson,
1077 K., Chen, W., Lo, I., Aoi, M., Momeni, A., Lin, S., Nikulina, V., Luo, W., Yang,
1078 L., Duong, T.Q., Razavian, N., Fuchs, T.J.: A foundation model for clinical-grade
1079 computational pathology and rare cancers detection. *Nature Medicine* **30**(10),
1080 2924–2935 (2024) <https://doi.org/10.1038/s41591-024-03141-0>
- 1081 [5] Hölscher, D.L., Bülow, R.D.: Decoding pathology: the role of computational
1082 pathology in research and diagnostics. *Pflügers Archiv - European Journal of*
1083 *Physiology* **477**(4), 555–570 (2025) <https://doi.org/10.1007/s00424-024-03002-2>
- 1084 [11] Chaurasia, A.K., Harris, H.C., Toohey, P.W., Hewitt, A.W.: A generalised vision
1085 transformer-based self-supervised model for diagnosing and grading prostate can-
1086 cer using histological images. *Prostate Cancer and Prostatic Diseases* (2025)
1087 <https://doi.org/10.1038/s41391-025-00957-w>
- 1088 [7] Paik, I., Lee, G., Lee, J., Kwak, T., Ha, H.K.: Artificial intelligence-driven digi-
1089 tal pathology in urological cancers: current trends and future directions. *Prostate*
1090 *International* (2025) <https://doi.org/10.1016/j.pnil.2025.02.002>
- 1091 [8] Wang, J., Mao, Y., Guan, N., Xue, C.: Advances in multiple instance learning for
1092 whole slide image analysis: Techniques, challenges, and future directions. *arXiv*
1093 (2024) [2408.09476](https://arxiv.org/abs/2408.09476)
- 1094 [9] Ogbonna, C.T., Ayankoya, F.Y., Kuyoro, S.O.: Enhancing prostate cancer prog-
1095 nosis through digital pathology and machine learning: A systematic review and
1096 meta-analysis. *Asian Journal of Engineering and Applied Technology* **13**(2), 44–51
1097 (2024) <https://doi.org/10.70112/ajeat-2024.13.2.4261>
- 1098 [10] Bulten, W., Kartasalo, K., Chen, P.-H.C., Ström, P., Pinckaers, H., Nagpal, K.,
1099 Cai, Y., Steiner, D.F., Van Boven, H., Vink, R., *et al.*: Artificial intelligence for
1100 diagnosis and gleason grading of prostate cancer: The panda challenge. *Nature*

- 1101 Medicine **28**(1), 154–163 (2022) <https://doi.org/10.1038/s41591-021-01620-2>
- 1102 [11] Chaurasia, A.K., Harris, H.C., Toohey, P.W., Hewitt, A.W.: A generalised vision
1103 transformer-based self-supervised model for diagnosing and grading prostate
1104 cancer using histological images. *Prostate Cancer and Prostatic Diseases*, 1–9
1105 (2025)
- 1106 [12] Ilse, M., Tomczak, J., Welling, M.: Attention-based deep multiple instance learn-
1107 ing. In: *International Conference on Machine Learning*, pp. 2127–2136 (2018).
1108 PMLR
- 1109 [13] Lu, M.Y., Williamson, D.F., Chen, T.Y., Chen, R.J., Barbieri, M., Mahmood,
1110 F.: Data-efficient and weakly supervised computational pathology on whole-slide
1111 images. *Nature biomedical engineering* **5**(6), 555–570 (2021)
- 1112 [14] Mai, C., Wang, Q., Mai, Z., Qin, C., Zeng, J., Xie, H., Xiao, Y., Huang, H., Chen,
1113 W., Yan, W., *et al.*: The application of multi-instance learning based on feature
1114 reconstruction and cross-mixing in the gleason grading of prostate cancer from
1115 whole-slide images. *Quantitative Imaging in Medicine and Surgery* **15**(4), 3263
1116 (2025)
- 1117 [15] Li, J., Li, W., Gertych, A., Knudsen, B.S., Speier, W., Arnold, C.W.: An
1118 attention-based multi-resolution model for prostate whole slide imageclassification
1119 and localization. *arXiv preprint arXiv:1905.13208* (2019)
- 1120 [16] Salsabili, S., Chan, A.D., Ukwatta, E.: Multiresolution semantic segmentation of
1121 biological structures in digital histopathology. *Journal of Medical Imaging* **11**(3),
1122 037501–037501 (2024)
- 1123 [17] Zheng, Y., Zhang, J., Huang, D., Hao, X., Qin, W., Liu, Y.: Detecting mri-invisible
1124 prostate cancers using a weakly supervised deep learning model. *International*
1125 *Journal of Biomedical Imaging* **2024**(1), 2741986 (2024)
- 1126 [18] Behzadi, M.M., Madani, M., Wang, H., Bai, J., Bhardwaj, A., Tarakanova, A.,
1127 Yamase, H., Nam, G.H., Nabavi, S.: Weakly-supervised deep learning model
1128 for prostate cancer diagnosis and gleason grading of histopathology images.
1129 *Biomedical Signal Processing and Control* **95**, 106351 (2024)
- 1130 [19] Shi, Z., Zhang, J., Kong, J., Wang, F.: Integrative graph-transformer frame-
1131 work for histopathology whole slide image representation and classification. In:
1132 *International Conference on Medical Image Computing and Computer-Assisted*
1133 *Intervention*, pp. 341–350 (2024). Springer
- 1134 [20] Mirabadi, A.K., Archibald, G., Darbandsari, A., Contreras-Sanz, A., Nakhli,
1135 R.E., Asadi, M., Zhang, A., Gilks, C.B., Black, P., Wang, G., *et al.*: Grasp:
1136 graph-structured pyramidal whole slide image representation. *arXiv preprint*
1137 *arXiv:2402.03592* (2024)

- 1138 [21] Shao, Z., Bian, H., Chen, Y., Wang, Y., Zhang, J., Ji, X., *et al.*: Transmil:
1139 Transformer based correlated multiple instance learning for whole slide image
1140 classification. *Advances in neural information processing systems* **34**, 2136–2147
1141 (2021)
- 1142 [22] Wang, X., Yang, S., Zhang, J., Wang, M., Zhang, J., Huang, J., Yang, W., Han,
1143 X.: Transpath: Transformer-based self-supervised learning for histopathological
1144 image classification. In: *International Conference on Medical Image Computing
1145 and Computer-Assisted Intervention*, pp. 186–195 (2021). Springer
- 1146 [23] Chen, R.J., Ding, T., Lu, M.Y., Williamson, D.F., Jaume, G., Song, A.H., Chen,
1147 B., Zhang, A., Shao, D., Shaban, M., *et al.*: Towards a general-purpose foundation
1148 model for computational pathology. *Nature medicine* **30**(3), 850–862 (2024)
- 1149 [24] Li, Y., Harten, C.J.A., Yaqub, M., Izadyyazdanabadi, H., Warraich, N.F., Raza,
1150 S.E.A., Bankhead, P., Rajpoot, N.: A systematic comparison of MIL approaches
1151 for gleason grading across multiple datasets. *Medical Image Analysis* **89**, 102768
1152 (2023)
- 1153 [25] Mir, A.N., Rizvi, D.R., Ahmad, M.R.: Enhancing histopathological image analy-
1154 sis: An explainable vision transformer approach with comprehensive interpreta-
1155 tion methods and evaluation of explanation quality. *Engineering Applications of
1156 Artificial Intelligence* **149**, 110519 (2025)
- 1157 [26] Shukla, A.K., Janmajaya, M., Abraham, A., Muhuri, P.K.: Engineering appli-
1158 cations of artificial intelligence: A bibliometric analysis of 30 years (1988–2018).
1159 *Engineering applications of artificial intelligence* **85**, 517–532 (2019)
- 1160 [27] Xiao, H., Li, L., Liu, Q., Zhu, X., Zhang, Q.: Transformers in medical image
1161 segmentation: A review. *Biomedical Signal Processing and Control* **84**, 104791
1162 (2023)
- 1163 [28] Zhang, J., Li, F., Zhang, X., Wang, H., Hei, X.: Automatic medical image
1164 segmentation with vision transformer. *Applied Sciences* **14**(7), 2741 (2024)
- 1165 [29] Grisi, C., Kartasalo, K., Eklund, M., Egevad, L., Laak, J., Litjens, G.: Hier-
1166 archical vision transformers for prostate biopsy grading: Towards bridging the
1167 generalization gap. *Medical Image Analysis* **105**, 103663 (2025)
- 1168 [30] Huang, G., Xia, B., Zhuang, H., Yan, B., Wei, C., Qi, S., Qian, W., He, D.:
1169 A comparative analysis of u-net and vision transformer architectures in semi-
1170 supervised prostate zonal segmentation. *Bioengineering* **11**(9), 865 (2024)
- 1171 [31] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recogni-
1172 tion. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern
1173 Recognition*, pp. 770–778 (2016)

- 1174 [32] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner,
1175 T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is
1176 worth 16x16 words: Transformers for image recognition at scale. arXiv preprint
1177 arXiv:2010.11929 (2020)
- 1178 [33] Xiang, J., Zhang, J.: Exploring low-rank property in multiple instance learning
1179 for whole slide image classification. In: The Eleventh International Conference on
1180 Learning Representations (2023)

ARTICLE IN PRESS