

---

# BreastGPT: A Multimodal Large Language Model for the Full Spectrum of Breast Cancer Clinical Routine

---

Anonymous Author(s)

## Abstract

1 Breast cancer remains a leading cause of cancer-related mortality among women.  
2 Its clinical management requires multimodal reasoning across a clinical workflow  
3 that spans *screening*, *diagnosis* and *treatment planning*, where each stage involves  
4 distinct imaging modalities, task objectives, and reasoning patterns. However,  
5 constrained by data scarcity and model versatility, existing medical MLLMs are  
6 typically evaluated on isolated modalities or narrow task families, limiting their  
7 ability to support workflow-level clinical reasoning. In this work, we first introduce  
8 **BreastStage**, a workflow-aligned breast imaging instruction corpus comprising  
9 1.86M instruction-following pairs curated from 17 sub-datasets across 5 imaging  
10 modalities and 136 task templates. Its held-out split, **BreastStage-Bench**, provides  
11 a comprehensive benchmark for evaluating multimodal reasoning across the breast  
12 cancer care continuum. Building on this corpus, we propose **BreastGPT**, a unified  
13 MLLM equipped with a dual-branch visual encoder and concept-preserving token  
14 compression to bridge the scale gap between standard radiology and gigapixel  
15 pathology. On BreastStage-Bench, BreastGPT achieves 75.66% closed-ended  
16 accuracy and 89.92% open-ended score, outperforming both general-purpose and  
17 medical-specific MLLMs across clinical stages and task formats. These results  
18 suggest that workflow-aligned data and cross-scale visual modeling are critical for  
19 clinically grounded medical MLLMs. All data, code, and model checkpoints are  
20 released at <https://anonymous.4open.science/r/BreastGPT>.

## 21 1 Introduction

22 Breast cancer remains one of the most prevalent and life-threatening malignancies among women  
23 worldwide [7, 20, 51, 35]. Its clinical management involves not just a single diagnostic decision,  
24 but a staged workflow spanning *screening*, *diagnosis*, and *treatment planning*. Each stage relies  
25 on distinct imaging evidence and pursues different clinical objectives: *Screening* prioritizes triage  
26 at the population level, where mammography, breast ultrasound (BUS) and chest CT are mainly  
27 used to detect early suspicious findings and estimate their risk [33]; *Diagnosis* then shifts toward  
28 fine-grained lesion characterization, where multimodal MRI is pivotal for evaluating morphology and  
29 probability of malignancy [6]; *Treatment planning* requires tumor biology and response assessment,  
30 for which gigapixel whole-slide images (WSIs) provide pathological subtyping and biomarker  
31 evidence [45]. This clinical continuum imposes a different requirement from conventional single-  
32 task image understanding: a clinically useful model must process heterogeneous modalities while  
33 preserving the stage-specific reasoning style expected at each point of care.

34 However, while recent advances in MLLMs have significantly expanded the scope of medical image  
35 understanding [43, 15, 67, 64, 68, 32, 3, 13, 40], recent efforts in breast imaging AI are mainly tailored  
36 for individual modalities or narrowly defined tasks. Specifically, mammography-oriented models [28,  
37 74] achieve strong performance on mammograph interpretation; BUS-CoT [69] introduces structured  
38 reasoning for breast ultrasound images, and MOME [44] supports multimodal MRI analysis with  
39 performance comparable to experienced radiologists. Despite demonstrating promising modality-  
40 specific performance, these methods conventionally inherit a fragmented formulation of breast cancer

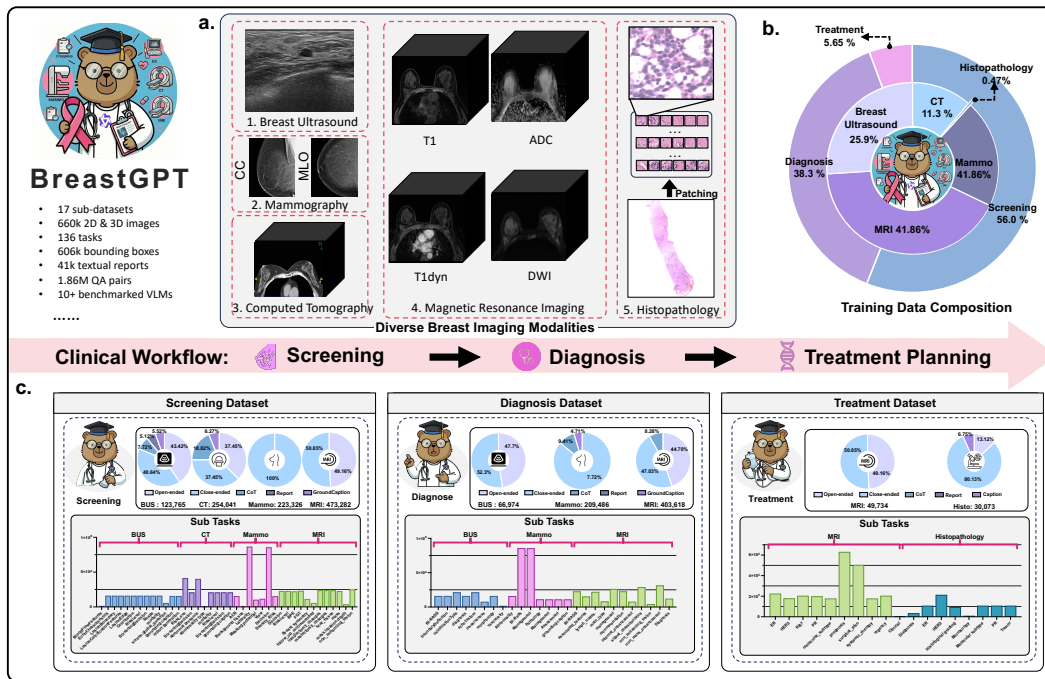


Figure 1: Overview of the **BreastStage** dataset aligned with the end-to-end breast cancer clinical workflow. The dataset is structured across three clinical stages: (1) **Screening** utilizes mammography and breast ultrasound (BUS) images for early lesion detection; (2) **Diagnosis** leverages MRI (main modality), BUS and Mammography for detailed tumor characterization; and (3) **Treatment planning** involves gigapixel whole-slide images (WSI) for pathological subtyping and therapeutic evaluation.

41 AI: datasets, models, and benchmarks are invariably organized around isolated imaging modalities  
 42 and limited tasks for a single clinical stage. Consequently, 1) these single-modality datasets fail to  
 43 comprehensively evaluate models' performance across the various stages within the practical clinical  
 44 workflow; 2) there is a lack of a versatile model capable of processing different modalities uniformly  
 45 and efficiently to address various tasks at different stages of real-world clinical workflows.

46 To this end, we first construct BreastStage (Fig. 1), a large-scale multimodal instruction corpus  
 47 aligned with the breast cancer clinical workflow. BreastStage integrates **17** sub-datasets across **5**  
 48 modalities (mammography, BUS, MRI, WSI, and CT), yielding **1.86M** instruction-following QA  
 49 pairs over **136** task templates. These instructions cover attribute extraction, report generation, visual  
 50 question answering, and image-grounded captioning, providing supervision over both radiological  
 51 and pathological evidence. Based on the held-out split of BreastStage, we further build **BreastStage-**  
 52 **Bench**, a comprehensive benchmark for workflow-level breast oncology reasoning developed with  
 53 breast oncology experts. We assess multiple publicly available multimodal large language models,  
 54 including both general-purpose and medical-specific variants, alongside advanced proprietary models.  
 55 Extensive experiments demonstrate that BreastStage-Bench poses a significant challenge for current  
 56 MLLMs: GPT-5.4 achieves only 54.00% closed-ended VQA accuracy and 53.58% open-ended VQA  
 57 score on average, and existing medical MLLMs show no clear advantage over general-purpose models  
 58 on these clinical workflow-aligned tasks.

59 Based on BreastStage, we propose **BreastGPT**, a unified MLLM for breast cancer workflow reasoning.  
 60 Built upon Qwen3-VL [9], BreastGPT handles all five imaging modalities within a single architecture  
 61 and adopts stage-conditioned system prompts to switch between screening, diagnostic, and treatment-  
 62 oriented reasoning behaviors. Meanwhile, to address the significant differences in image scale  
 63 across modalities, BreastGPT adopts a dual-branch visual encoder with a resolution-aware gating  
 64 module that automatically routes each input to an appropriate processing branch based on stage  
 65 condition. Aligned with practical clinical workflow, BreastGPT demonstrates exceptional versatility  
 66 and achieves remarkable performance across various tasks. Compared to existing general-purpose  
 67 and medical-specific MLLMs, BreastGPT achieves performance gains of over 25%, 35%, and 40%  
 68 for screening, diagnosis, and treatment planning, respectively. To summarize, our contributions are  
 69 threefold:

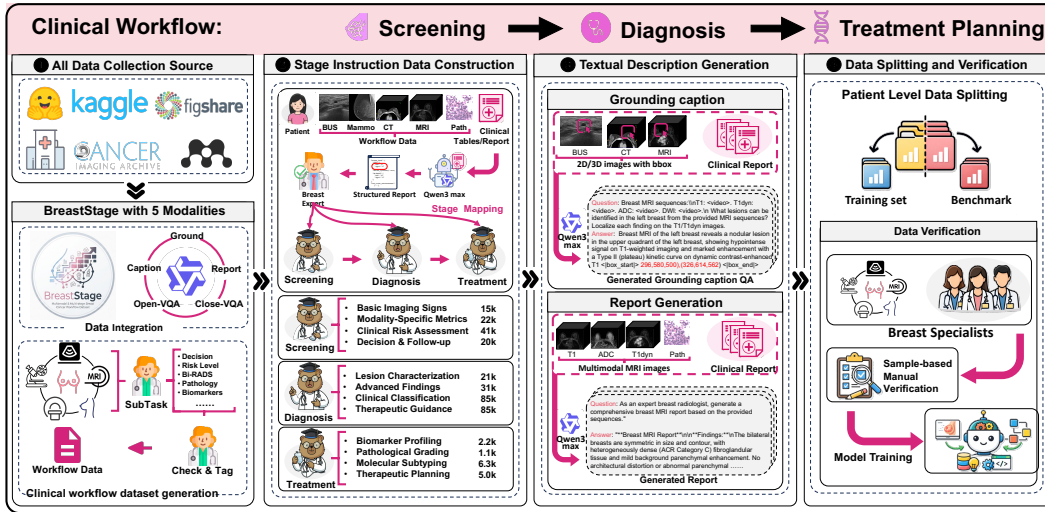


Figure 2: **BreastStage construction pipeline.** It involves curating diverse imaging modalities from public datasets and hospitals. The raw data is then cleaned, annotated, and checked by breast experts. Textual and instruction data are generated by prompting Qwen3-Max. Finally, the data is split into a training set and a benchmark, with professional breast experts assessing the training samples.

70 **Workflow-level problem formulation and benchmark.** We formulate breast cancer multimodal  
 71 reasoning as a workflow-aligned problem spanning screening, diagnosis, and treatment planning.  
 72 Based on this, we construct BreastStage, a large-scale instruction corpus with 1.86M QA pairs from  
 73 17 sub-datasets across 5 modalities, and BreastStage-Bench, a held-out benchmark for evaluating  
 74 stage-specific clinical reasoning.

75 **Unified cross-scale breast MLLM.** We introduce BreastGPT, a unified multimodal LLM that handles  
 76 mammography, ultrasound, CT, MRI, and gigapixel pathology within a single instruction-following  
 77 framework, enabling consistent reasoning across the full breast cancer clinical routine.

78 **Concept-preserving visual token compression.** To bridge the extreme scale gap between standard  
 79 radiology and WSIs, we propose a dual-branch visual encoder with a concept-based token selector  
 80 that preserves clinically salient visual evidence under a fixed token budget, substantially improving  
 81 both accuracy and inference efficiency.

## 82 2 BreastStage Dataset Curation

### 83 2.1 Scope and scale

84 BreastStage contains about 662,000 unique 2D and 3D images, 136 task templates, 606,226 records  
 85 with bounding-box or mask annotations, and 1.86 million instruction-following pairs. Task coverage  
 86 spans classification, visual grounding, open- and closed-ended question answering, captioning, and  
 87 report generation, so the dataset measures both perception accuracy on narrow subtasks and the  
 88 model’s ability to follow diverse clinical instructions. Samples are distributed as 57.9% screening,  
 89 36.7% diagnosis, and 5.4% treatment, broadly tracking the relative prevalence and data availability of  
 90 these stages in real breast care. Please check § C.1 for details.

### 91 2.2 Data construction pipeline

92 As illustrated in Fig. 2, the BreastStage curation pipeline has four steps:

93 **i. Workflow Data Generation.** We assemble a cohort of breast imaging studies that draws CT, BUS,  
 94 mammography, and WSI from public repositories and adds multimodal MRI from two collaborating  
 95 clinical institutions. For quality control, modality-specific visual specialist agents built on Qwen2.5-  
 96 VL-72B [10] emit a structured {validity, reason} verdict for each image; samples flagged as  
 97 low-quality are routed to a Low-Image Check in which a breast specialist confirms or restores the  
 98 verdict, rather than being silently dropped. Lesion-level spatial annotations are then attached: BUS  
 99 reuses the expert hand-drawn masks shipped with BUS-CoT [69], mammography reuses the EMBED-

100 released bounding boxes, CT volumes from CT-RATE [22] are passed through DRT-M3D [73] to  
101 obtain tumour segmentations, and MRI carries manual T1 and dynamic contrast-enhanced T1dyn  
102 annotations from a panel of 10 breast specialists. The resulting cohort is tagged with one of 136  
103 specialist-designed sub-tasks (`task template`) and audited by hand. Full data provenance and  
104 per-modality counts are reported in Table 5 and the complete pipeline is detailed in § C.2.

105 **ii. Stage Instruction Data Construction.** All data is standardized, cleaned, and paired with its  
106 matching clinical tables or reports, then converted into instruction-following QA pairs: 1) Chinese  
107 reports are first translated to English by Qwen3-Max acting as a radiologist agent; 2) Each translated  
108 report is then parsed into a modality-specific and BI-RADS-aligned structured report whose schema  
109 is designed by breast experts following clinical reporting guidelines; 3) For every `task template`  
110 from the previous stage, breast experts then write a closed-ended question whose answer is drawn  
111 from one or more leaf fields of the structured report, with options taken directly from the schema’s  
112 enum values, so the resulting QA pairs cannot introduce facts beyond the structured record; 4) Finally,  
113 each `task template` is mapped, with expert guidance, to one of the three clinical stages: screening,  
114 diagnosis, or treatment.

115 After the stage mapping, the resulting BreastStage instruction set is organized to match the cognitive  
116 demands of each clinical stage, with category-level counts summarized in Table 6. Screening tasks  
117 emphasize early detection and triage, covering modality / view / laterality recognition, breast-tissue  
118 characterization (FGT, BPE, density), risk and decision tasks, lesion presence and morphology, and  
119 post-surgical change tracking, for a total of 1,075,092 pairs. Diagnosis tasks emphasize tumour  
120 characterization and staging, including BI-RADS and pathology classification, lesion-level signal  
121 and kinetic characterization, and associated findings such as lymph-node involvement and tissue  
122 invasion, for a total of 680,409 pairs. Treatment tasks are grounded in histopathology and treatment-  
123 relevant radiology, including prognosis prediction, surgical planning, systemic therapy and biomarker  
124 assessment, and WSI-based pathology subtyping, for a total of 100,567 pairs. The hierarchical  
125 organization ensures that instruction tuning covers the distinct reasoning patterns required across the  
126 patient journey. Please check § C.2 for details.

127 **iii. Textual Description Generation.** The same structured records that feed VQA construction  
128 are also reused to synthesise two narrative task types: ground captions and report generation: A  
129 *ground caption* encodes 2D or 3D spatial coordinates of lesions together with their categories and  
130 anatomical locations, providing a multi-dimensional spatial observation for BUS, mammography,  
131 CT, and MRI. We construct ground captions by prompting Qwen3-Max to integrate the original  
132 clinical report with the available bounding-box annotations into a single visually grounded narrative.  
133 *Medical reports*, in contrast, condense the multimodal evidence into key findings, tissue character-  
134 istics, and diagnostic or therapeutic assessments. After consultation with senior breast oncologists,  
135 we organize each generated report into the four ACR-aligned sections *Findings*, *Impression*, *Final*  
136 *Assessment*, and *Management*. The report generation pipeline is iterative: Qwen3-Max first drafts  
137 a comprehensive report from the multimodal input, including multiparametric MRI sequences and  
138 histopathology WSIs; breast specialists then audit a sample, catalog the recurring reasoning errors  
139 and hallucinations, and we refine the modality-specific system prompts before regenerating.

140 **iv. Data Splitting and Verification.** To prevent data leakage while preserving task diversity, we  
141 enforce strict patient-level separation between the training set and the evaluation benchmark. We  
142 allow the same image to be reused across different task contexts—for instance a BUS image may  
143 carry both a BI-RADS label and a lesion-characterization task—but every record derived from a given  
144 patient is kept on a single side of the split. The split itself uses stratified sampling on the composite  
145 key of modality, task type, and pathology label, so the underlying task and class distributions stay  
146 consistent across train and test. Downstream evaluation therefore reflects clinical generalization  
147 rather than memorization of shared visual content.

148 For data reliability we run a multi-stage verification pipeline. Automated heuristic filters first remove  
149 hallucinated clinical terms, malformed bounding-box coordinates, instruction-answer conflicts, and  
150 near-duplicate QA pairs detected by MinHash similarity above 0.85. Three breast specialists then  
151 conduct an independent task-level audit on the held-out test partition (5 random samples per task  
152 per specialist), and any task with consistent flags triggers refinement of the generation prompts or  
153 filtering rules. **Notably**, the clinical reliability of all generated data has been audited by clinical  
154 experts. In § E.3, we report expert agreement on a stratified subset of BreastStage-Bench, including

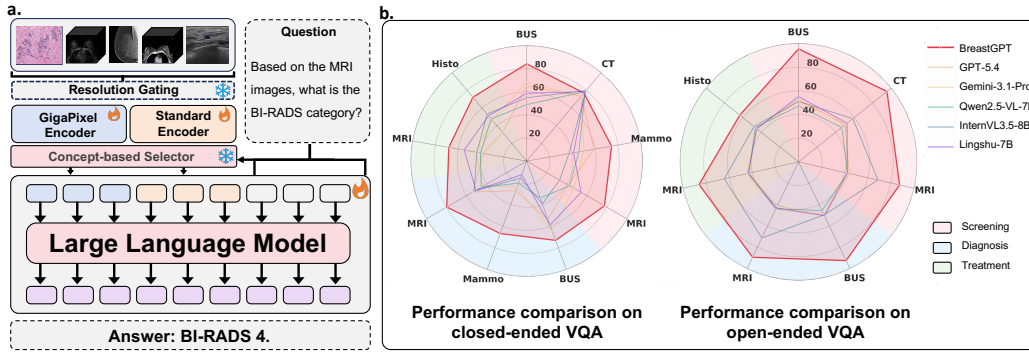


Figure 3: (a) **Architecture of BreastGPT**. A resolution gating mechanism routes standard radiological images to a ViT-based **Standard Branch**, and extreme-resolution WSIs to a specialized **GigaPixel Branch**. The WSI pipeline extracts patch features via a frozen CONCHv1.5 encoder, aggregates global context using LongNet, and compresses the sequence via a universal concept-based token selector before LLM injection. (b) **Performance Comparison**. BreastGPT demonstrates superior reasoning against state-of-the-art MLLMs across both closed- and open-ended clinical QA tasks.

155 task validity, answer correctness, and clinical consistency. The expert audit covers all modalities and  
 156 task families, and disagreements are resolved by consensus review.

### 157 3 BreastGPT

158 BreastGPT is built upon Qwen3-VL [9], a vision-language model with native multi-image support,  
 159 following the broader Qwen-VL lineage of resolution-aware multimodal modeling [8, 61, 10]. We  
 160 preserve the LLM backbone architecture and adapt it to breast cancer tasks through stage-aware  
 161 system prompts, which act as lightweight task routers at inference time. This design allows a single  
 162 unified model to operate across the screening  $\rightarrow$  diagnosis  $\rightarrow$  treatment workflow without introducing  
 163 task-specific heads. Compared with fragmented modality-specific systems, such a unified design  
 164 shares multimodal representations across tasks while maintaining stage-aware specialization. As  
 165 a result, BreastGPT supports diverse specialist-style reasoning across 136 task templates within a  
 166 single instruction-following framework. To handle image scales ranging from standard radiological  
 167 scans to gigapixel histopathology, we design a dual-branch visual processing architecture modulated  
 168 by a *modality-aware resolution gating* mechanism. Inputs tagged as WSIs are routed to the GigaPixel  
 169 branch, while all other radiological modalities (CT, MRI, BUS, mammography) flow through the  
 170 Standard branch, as illustrated in Fig. 3 a.

#### 171 3.1 Dual-Branch Visual Encoding

172 **Standard Encoder for Radiology.** Radiological modalities (CT, MRI, BUS, mammography) are  
 173 encoded by the native Vision Transformer (ViT) of the Qwen3-VL backbone. This branch is optimized  
 174 to capture macroscopic anatomical structures, tissue densities, and localized lesion textures that are  
 175 critical for screening and diagnostic reasoning.

176 **GigaPixel Encoder for Pathology.** To integrate extreme-resolution histopathology without over-  
 177 whelming the LLM context window, we introduce a specialized WSI pathology branch. Prior work  
 178 on pathology foundation models and WSI transformers has shown that both spatially preserved  
 179 local morphology and slide-level context are essential for reliable gigapixel image understand-  
 180 ing [1, 43, 15, 16, 21, 11, 38, 65]. Pathology information must therefore be incorporated without  
 181 collapsing tens of thousands of patch-level features into clinically uninformative representations. We  
 182 separate WSI processing into two steps before compression:

- 183 • *Feature Extraction:* WSI tiles at  $20\times$  magnification are encoded by a frozen CONCH  
 184 v1.5 foundation model [43], producing 512-dimensional patch embeddings  $\{p_1, \dots, p_N\}$ ,  
 185 where  $N$  can reach tens of thousands. Using a frozen encoder retains robust pathology  
 186 representations from large-scale pretraining while avoiding the computational instability of  
 187 end-to-end retraining.
- 188 • *Context Aggregation:* A trainable LongNet dilated attention encoder [17] aggregates local  
 189 cellular details and global architectural patterns across the full patch sequence, producing

190 contextualized representations  $\{h_1, \dots, h_N\}$ . Because diagnostic cues arise from both local  
 191 cell morphology and broader tissue organization, LongNet’s exponentially growing dilation  
 192 rates preserve long-range context without the  $O(N^2)$  bottleneck of standard attention.

193 Full implementation details (hyperparameters, dilation schedule, projection dimensions) are in § D.2.

### 194 3.2 Universal Concept-Based Compression

195 A unified model spanning 3D radiological volumes and gigapixel pathology faces a severe token-  
 196 budget mismatch, as a single WSI patch sequence routinely exceeds tens of thousands of tokens.  
 197 Rather than resorting to modality-specific pooling or heuristic truncation, we adapt the MMTok  
 198 multimodal coverage maximization framework [18] as a universal concept-based token selector,  
 199 following recent work on efficient VLM token pruning and compact visual representations [52, 56,  
 200 55, 46, 49]. Given either redundant radiological ViT tokens or contextualized WSI embeddings  
 201  $\mathbf{X} = \{x_i\}_{i=1}^N$ , the selector distills the input into exactly  $k$  tokens by maximizing a concept-coverage  
 202 objective:

$$\max_{S \subseteq \{1, \dots, N\}, |S|=k} \underbrace{\frac{1}{M} \sum_{j=1}^M \max_{i \in S} \text{sim}(t_j, x_i)}_{\text{text-vision coverage}} + \alpha \underbrace{\frac{1}{N} \sum_{j=1}^N \max_{i \in S} \text{sim}(x_j, x_i)}_{\text{vision-vision coverage}}, \quad (1)$$

203 where  $\{t_j\}_{j=1}^m$  are text tokens from the clinical instruction and  $\alpha$  balances query relevance with  
 204 visual representativeness. The first term anchors selection to prompt-relevant clinical concepts, while  
 205 the second term prevents collapse onto redundant salient regions and encourages coverage of global  
 206 tissue or anatomical structure. This retains visually distinct, clinically informative patterns, enabling  
 207 BreastGPT to condense multi-scale visual evidence into a fixed token budget without sacrificing  
 208 the diversity required for downstream reasoning. The complete calibration and greedy optimization  
 209 procedure are provided in § D.3.

210 **Training.** BreastGPT is trained in two stages: a visual front-end warm-up that freezes the LLM  
 211 and aligns the ViT and CONCH+LongNet branches with the multimodal token space, followed by  
 212 end-to-end fine-tuning on all four task formats across every modality. The coverage token selector  
 213 is training-free and applied identically at training and inference. Detailed optimizer settings, and  
 214 per-stage trainable/frozen module lists are reported in § D.4.

## 215 4 Experiments

### 216 4.1 Experimental Setups

217 **Baselines.** We compare BreastGPT with three groups of representative vision–language models on  
 218 BreastStage-Bench: 1) *proprietary frontier models*: GPT-5.4, Claude-opus-4-6, Claude-sonnet-4-6,  
 219 Gemini-3.1-Flash, and Gemini-3.1-Pro, queried through their official APIs under identical stage-  
 220 specific system prompts; these models represent frontier model families commonly used in recent  
 221 medical MLLM benchmarks [29, 58]; 2) *open-source general-purpose VLMs* at comparable scale  
 222 to BreastGPT: Qwen2.5-VL-Instruct (3B/7B) [10], Qwen3-VL-Instruct (4B/8B) [9], MiMo-VL-  
 223 SFT [57], and InternVL3.5; 3) *medical-specific VLMs*: Lingshu [67] and HuatuoGPT-V [13], with  
 224 additional context from prior biomedical and medical VLMs [37, 26, 40]. All baselines are evaluated  
 225 zero-shot on the same instruction-formatted prompts used for BreastGPT, ensuring that differences  
 226 in performance reflect model capability rather than prompt engineering. Detailed per-baseline  
 227 configuration is listed in § E.4.

228 **BreastGPT variants.** We report two variants that differ only in how the final 128 tokens are obtained  
 229 to verify the effect of our training-free selection strategy: BreastGPT (cluster) uses the training-free  
 230 greedy coverage selector, while BreastGPT (learn) replaces it with a learnable cross-attention retrieval  
 231 head in the spirit of the Perceiver Resampler used in Flamingo [2], trained jointly in Stage 2. Both  
 232 variants share the same backbone, branches, and training data.

Table 1: VQA performance on BreastStage-Bench. The left half reports **closed-ended accuracy (%)** and the right half reports **open-ended normalized score (%)** for the same set of models. Open-ended evaluation does not include Mammography. Avg is the mean across the task columns of the corresponding half. Best per column in **bold**; BreastGPT rows in cyan.

Model	Closed-ended VQA (Accuracy, %)											Open-ended VQA (normalized Score, %)									
	#P	Screening			Diagnosis			Treatment			Avg	Screening			Diagnosis			Treatment			Avg
		BUS	CT	Mam	MRI	BUS	Mam	MRI	MRI	His		BUS	CT	MRI	BUS	MRI	MRI	His			
<i>Proprietary Models</i>																					
GPT-5.4	-	64.89	<b>78.55</b>	<b>68.51</b>	41.43	53.46	55.26	53.50	38.10	32.28	<b>54.00</b>	53.43	49.34	59.85	50.42	59.20	58.07	44.73	53.58		
Claude-opus-4-6	-	50.21	72.00	39.57	50.48	38.83	7.66	45.10	41.27	25.94	41.23	43.32	42.77	43.75	42.80	44.42	42.91	40.84	42.97		
Claude-sonnet-4-6	-	65.53	68.00	37.02	40.00	54.79	19.67	55.94	42.86	36.30	46.68	43.63	43.12	44.31	43.00	44.98	42.62	41.06	43.25		
Gemini-3.1-Flash	-	55.11	42.67	49.68	38.57	37.60	24.32	31.82	48.41	35.81	40.44	48.76	48.21	47.38	45.80	47.53	44.15	42.81	46.38		
Gemini-3.1-Pro	-	68.09	73.33	50.21	47.14	64.36	23.87	43.88	44.44	46.53	<b>51.32</b>	51.19	51.46	41.60	48.94	41.39	42.80	45.77	46.16		
<i>Open-Source Models</i>																					
Qwen2.5-VL	3B	46.81	51.39	33.51	45.71	30.85	12.76	51.75	44.44	51.52	40.97	49.52	48.39	49.36	47.36	50.54	47.04	44.80	48.14		
Qwen2.5-VL	7B	49.15	<b>79.27</b>	44.47	44.76	34.31	14.41	46.85	37.30	47.62	44.24	46.20	46.86	49.47	44.99	49.67	44.70	43.98	46.55		
Qwen3-VL	4B	52.13	75.39	41.91	49.05	35.37	21.17	48.08	43.65	38.61	45.04	45.20	55.49	47.71	44.93	51.46	48.98	41.26	47.86		
Qwen3-VL	8B	57.87	<b>78.55</b>	39.68	51.43	48.14	25.83	47.90	34.92	47.02	47.93	44.94	48.53	44.85	44.34	46.24	44.15	41.21	44.89		
MiMo-VL	7B	15.11	23.52	8.62	6.67	14.89	6.16	0.35	0.79	19.98	10.68	42.49	55.98	63.70	41.69	65.66	63.08	39.86	53.21		
InternVL3.5	8B	51.91	77.70	35.85	34.76	42.55	16.22	52.27	44.44	52.98	45.41	50.74	58.05	56.56	48.80	58.89	55.94	46.50	53.64		
<i>Medical-Specific Models</i>																					
Lingshu	7B	58.94	<b>78.55</b>	39.89	54.29	58.24	8.56	45.28	<b>58.73</b>	51.52	50.44	52.60	49.84	53.41	48.95	53.76	49.60	43.69	50.26		
HuatoGPT-V	7B	45.74	71.39	43.09	39.05	35.11	9.61	47.73	45.24	49.09	42.89	51.86	51.47	55.68	49.20	55.08	52.93	45.78	51.71		
Qwen3-VL (SFT)	8B	72.34	76.73	65.74	<b>78.10</b>	71.01	59.76	<b>75.87</b>	53.97	60.41	68.21	94.43	<b>95.38</b>	<b>94.57</b>	91.70	<b>95.28</b>	<b>88.36</b>	57.96	<b>88.24</b>		
BreastGPT (cluster)	8B	<b>86.81</b>	77.21	<b>75.00</b>	<b>82.86</b>	<b>77.13</b>	<b>68.32</b>	<b>81.12</b>	<b>61.11</b>	<b>71.38</b>	<b>75.66</b>	<b>95.97</b>	<b>95.29</b>	<b>95.48</b>	<b>93.24</b>	<b>95.72</b>	<b>89.93</b>	<b>63.80</b>	<b>89.92</b>		
BreastGPT (learn)	8B	<b>84.47</b>	71.03	<b>68.51</b>	75.71	<b>78.46</b>	55.26	73.95	57.14	<b>71.25</b>	<b>70.64</b>	<b>95.86</b>	<b>94.73</b>	87.12	<b>93.57</b>	85.22	81.78	<b>63.36</b>	85.95		

Note: BUS = Breast Ultrasound, CT = Computed Tomography, Mam = Mammography, His = Histopathology. Mammography is omitted from the right half because it has no open-ended evaluation set.

## 233 4.2 Evaluation Metrics

234 We employ a comprehensive suite of metrics tailored to the diverse task formats within BreastStage-  
 235 Bench. For closed-ended tasks, we report standard classification accuracy (ACC). For open-ended  
 236 reasoning and report generation tasks, we report BERTScore F1, BLEU, and ROUGE-1, alongside  
 237 a weighted aggregate score (0.5 BERTScore F1, 0.25 BLEU, and 0.25 ROUGE-1). We assign a  
 238 predominant weight to BERTScore F1 because it captures deep semantic equivalence rather than  
 239 surface-level lexical overlap, a critical requirement for evaluating clinically constrained free-text  
 240 responses. For visual grounding tasks, we report mean Intersection over Union (IoU); a sample  
 241 without a ground-truth bounding box is credited as  $\text{IoU} = 1$  when the model also abstains from  
 242 emitting a box (true negative) and 0 when it hallucinates one (false positive). Please check § E.2 for  
 243 details.

## 244 4.3 Evaluation Results

245 **Overall observations.** Across Tables 1 and 2, BreastGPT achieves the best score in every clinical  
 246 stage and every task format. BreastGPT (cluster) reaches 75.66% closed-ended and 89.92% open-  
 247 ended VQA on average—more than 20 points above the strongest non-SFT baselines in the main  
 248 cohort (GPT-5.4 at 54.00% closed, InternVL3.5 at 53.64% open)—and raises the MRI report weighted  
 249 score from 55.16% (GPT-5.4) to 67.67%. On 3D grounding (CT, MRI), BreastGPT is the only model  
 250 that produces non-trivial volumetric bboxes (Table 13). Eight additional baselines (Grok-4.1-Fast [63],  
 251 Gemma-4 [59], GLM-4.6V-Flash [60], LLaVA-OneVision-1.5 [5], HealthGPT [40], Hulu-Med [31],  
 252 MedDr [26], RadFM [62]) are evaluated in Supplementary Tables 10 and 11; none narrows the gap to  
 253 BreastGPT to under 25 points.

254 **Proprietary frontier models underperform on breast workflows.** Despite their general multimodal  
 255 strength, Claude-opus-4-6, Claude-sonnet-4-6, and Gemini-3.1-Pro stay at 41–51% average on  
 256 closed-ended VQA, and even the strongest proprietary model (GPT-5.4 at 54.00%) is more than 20  
 257 points below BreastGPT. Mammography-based diagnosis (Diag / Mammo) is especially brittle for  
 258 the Claude family: Claude-opus-4-6 drops to 7.66% and Claude-sonnet-4-6 to 19.67%, while Gemini-  
 259 3.1-Pro reaches 23.87%. This confirms that breast workflows require domain knowledge (BI-RADS  
 260 lexicon, molecular subtyping conventions, multiparametric MRI interpretation) that frontier models  
 261 have not internalized from general web-scale pretraining.

Table 2: Caption and Report Generation Performance (%). IoU = mean grounding IoU on the lesion-grounded captioning task, credited as 1.0 on samples with no ground-truth bbox where the model also refrains from emitting one (true negatives), 0 on hallucinated bboxes; “-” indicates the model produces no usable bbox in that modality. Wtd = weighted composite ( $0.5 \times \text{BERT} + 0.25 \times \text{BLEU} + 0.25 \times \text{R-1}$ ). IoU on the 3D modalities (CT, MRI) is reported in Supplementary Table 13. Best results in **bold**. BreastGPT rows highlighted in light blue.

Model	#P	Caption						Report	
		BUS		CT	Mammo		Histo	MRI	
		IoU	Wtd	Wtd	IoU	Wtd	Wtd	Wtd	
<i>Proprietary Models</i>									
GPT-5.4	-	6.68	47.66	43.95	9.08	45.10	44.36	47.00	55.16
Claude-opus-4-6	-	11.41	46.99	43.68	<u>12.06</u>	47.04	44.94	47.75	48.56
Claude-sonnet-4-6	-	7.40	46.60	42.82	-	44.79	44.21	49.30	49.30
Gemini-3.1-Flash	-	<u>56.20</u>	51.48	46.01	7.52	49.29	44.89	51.45	51.03
Gemini-3.1-Pro	-	-	45.60	44.64	-	42.36	45.84	54.86	54.86
<i>Open-Source Models</i>									
Qwen2.5-VL	3B	2.67	46.84	46.88	1.38	48.36	45.71	50.69	49.59
Qwen2.5-VL	7B	6.38	46.79	47.31	1.47	48.11	45.07	51.95	50.71
Qwen3-VL	4B	33.64	46.68	41.66	4.29	46.19	45.38	37.66	49.46
Qwen3-VL	8B	29.35	46.97	47.91	4.06	47.09	46.03	49.57	50.49
MiMo-VL	7B	12.35	43.22	46.14	2.22	43.80	43.98	47.32	46.55
InternVL3.5	8B	0.35	46.97	47.69	0.17	45.72	45.77	49.28	49.12
<i>Medical-Specific Models</i>									
Lingshu	7B	5.00	48.43	49.45	0.55	49.14	46.86	51.84	51.14
HuatuogPT-V	7B	0.21	47.46	48.92	0.47	48.84	46.24	51.46	51.00
<b>Qwen3-VL (SFT)</b>	<b>8B</b>	2.79	76.61	72.72	8.62	<u>76.68</u>	62.42	<u>65.21</u>	<u>66.10</u>
<b>BreastGPT (cluster)</b>	<b>8B</b>	<b>79.59</b>	<b>79.32</b>	<u>73.16</u>	<b>23.14</b>	<b>77.64</b>	<u>66.78</u>	<b>67.69</b>	<b>67.67</b>
<b>BreastGPT (learn)</b>	<b>8B</b>	-	<u>79.16</u>	<b>77.07</b>	-	76.54	<b>68.11</b>	-	66.04

Note: IoU = mean grounding IoU including true-negative credit; “-” indicates the model produces no usable bbox in that modality; “-” indicates the metric is not reported for that row. Wtd = weighted composite score. CT and MRI grounding IoU are reported in Supplementary Table 13. BUS = Breast Ultrasound, CT = Computed Tomography, Mammo = Mammography, Histo = Histopathology. Report/MRI = structured radiology report generation for multiparametric MRI.

262 **Medical-specific VLMs offer no clear advantage.** Lingshu and HuatuogPT-V perform comparably  
 263 to general-purpose 7–8B VLMs on closed-ended VQA (50.44% and 42.89%), and sometimes under-  
 264 perform them. This indicates that existing medical pretraining corpora do not transfer specifically to  
 265 breast cancer workflow tasks, motivating the need for a workflow-aligned dataset like BreastStage.

266 **Cluster vs. learn variant.** BreastGPT (cluster), which uses the training-free greedy coverage  
 267 selector, edges out BreastGPT (learn) on both closed-ended VQA (75.66% vs. 70.64%) and open-  
 268 ended VQA (89.92% vs. 85.95%), despite having no additional learnable parameters for token  
 269 selection. This confirms that the submodular coverage objective is a strong inductive bias for WSI  
 270 and multimodal token compression, and that a learnable retriever does not trivially outperform a  
 271 principled, training-free procedure in the breast workflow setting.

#### 272 4.4 Ablation Studies

273 We conduct three ablation studies to verify the contribution of individual components.

274 **Necessity of the GigaPixel branch.** To isolate the GigaPixel branch, the Qwen3-VL (SFT) row in  
 275 Tables 1 and 2 reports a controlled baseline sharing the same backbone (Qwen3-VL-8B), training  
 276 data, and SFT recipe but processing WSIs through the standard ViT branch via 32 randomly sampled  
 277  $512 \times 512$  patches per slide. Across the four radiology modalities, the gap to BreastGPT (cluster) is  
 278 modest ( $\sim 7$  points on closed VQA), confirming that random patch sampling is not the bottleneck  
 279 when the source image already fits a ViT. On histopathology, however, the gap widens sharply: the 32-  
 280 patch baseline reaches only 60.41% closed-ended accuracy and 62.42 weighted caption score, versus  
 281 71.38% and 66.78 for BreastGPT (cluster). Since 32 tiles cover well under 1% of a typical breast WSI,  
 282 subtype, grade, and treatment-response cues are routinely missed; the CONCHv1.5+LongNet pipeline  
 283 together with the coverage selector recovers this signal by reasoning over a globally contextualised  
 284 slide representation. A per-component ablation inside the GigaPixel branch is in Supplementary  
 285 Table 16.

286 **Visual token budget sensitivity.** We sweep  $k \in \{1, 8, 16, 32, 64, 128, 256, 512\}$  across five modalities  
 287 and three task families (full results in Supplementary Tables 17 and 18). *Closed VQA* (Fig. 5a)  
 288 saturates early: BUS-Diagnosis reaches 77.66% at  $k=64$  and plateaus at 77.13% from  $k=128$ ,  
 289 while Mammo/MRI plateau by  $k=32$ . *Caption/report* (Fig. 5b) peaks near  $k=128$ : radiology gains  
 290 marginally from  $k=128$  to  $k=512$  (+0.06 BUS, +0.21 Mammo) while histopathology actually  
 291 degrades (67.63 at  $k=8 \rightarrow 64.53$  at  $k=512$ ), suggesting additional WSI tokens dilute rather than add  
 292 information. *Ground caption* (Fig. 5c) is more token-demanding—a single token cannot represent a

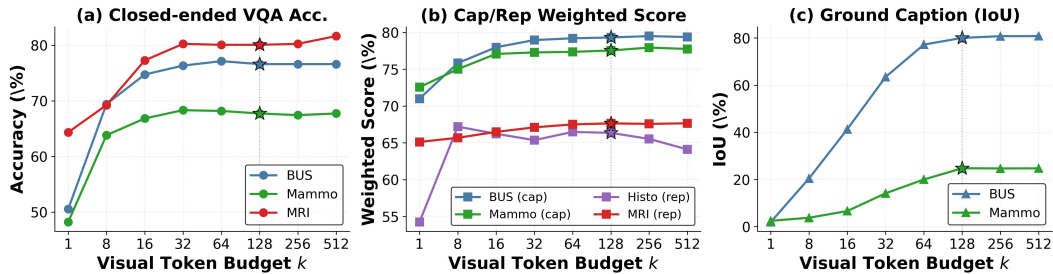


Figure 5: Visual token budget sweep on BreastStage-Bench. (a) Closed-ended VQA accuracy on diagnosis-stage tasks for BUS, mammography, and MRI. (b) Caption and report weighted score for BUS, mammography, histopathology report, and MRI report. (c) Ground caption IoU on BUS and mammography. Dotted vertical lines and stars mark the chosen operating point  $k = 128$ .

293 bounding region (IoU 1.50 BUS, 0.85 Mammo at  $k=1$ )—but  $k=128$  already recovers  $>99\%$  of the  
 294  $k=512$  IoU. We therefore set  $k=128$  as the unified budget: it sits on the saturation plateau across all  
 295 task families and keeps radiology and pathology inputs at the same downstream cost.

296 **GPU time.** Fig. 4 measures inference latency on one  
 297 representative histopathology WSI with 5,987 patches.  
 298 Each  $k$  has two bars: the longer bar is the LLM prefill on  
 299 the selected  $k$  tokens, and the shorter bar is the time spent  
 300 inside the selector itself. At our chosen budget  $k=128$ ,  
 301 prefill takes 191.4 ms and the selector adds 9.3 ms, so  
 302 the total inference latency is 200.6 ms. If we instead skip  
 303 the selector and feed all 5,987 patch tokens directly to the  
 304 LLM (the dashed line in the figure), total latency rises  
 305 to 6.5 s,  $33\times$  slower than  $k=128$  on the same slide, and  
 306 peak memory grows from 16.97 GB to 17.95 GB. Going  
 307 the other direction, raising  $k$  to 512 makes prefill grow to  
 308 605.7 ms but improves task quality by less than 1% over  
 309  $k=128$ . The selector is therefore the right tradeoff in this  
 310 regime: 9.3 ms of extra compute removes roughly 6.3 s  
 311 of LLM work that would otherwise dominate the pipeline.  
 312 Full per- $k$  numbers are in Supplementary Table 15.

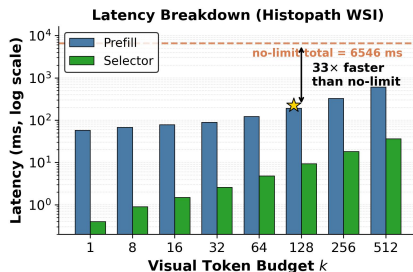


Figure 4: GPU latency on a histopathology WSI. Two bars per  $k$  give prefill and selector latency. The dashed line is the no-limit baseline that bypasses the selector; its gap to the chosen- $k$  prefill bar is the latency saved by selection. Log-scaled  $y$ -axis.

313 **More results.** In § F, we further report per-modality performance (§ F.1) and inference efficiency (§  
 314 F.6). In § G, we present a more detailed ablation about WSI branch (§ G.1) and visual token budget  
 315 sweep (§ G.2). A qualitative analysis and detailed case studies are presented in § H.

## 316 5 Conclusion

317 Breast cancer care is a tightly coupled clinical workflow spanning screening, diagnosis, and treatment,  
 318 yet existing medical MLLMs are optimised for individual stages and struggle on cross-stage reason-  
 319 ing. We close this gap from both sides: **BreastStage**, a large-scale multimodal instruction corpus  
 320 and benchmark organised around the screening–diagnosis–treatment pathway, and **BreastGPT**, a  
 321 unified VLM that pairs stage-aware role prompting with a dual-branch visual encoder for cross-  
 322 scale modelling from standard radiology to gigapixel pathology. On BreastStage-Bench, BreastGPT  
 323 substantially closes the cross-stage gap left by both general-purpose and medical-specific MLLMs.

324 **Limitation and Future Work.** BreastStage covers all three stages but most samples are drawn  
 325 from disjoint patient cohorts rather than longitudinal records of the same patient, so the corpus  
 326 supports workflow-aligned supervision at the task/modality level without fully capturing patient-level  
 327 temporal continuity. A held-out subset of BreastStage-Bench does contain full-workflow cases from  
 328 the same patient, allowing partial probing of cross-stage reasoning. Scaling this to a per-patient  
 329 cross-stage corpus is a field-level bottleneck (linked cross-department records, multi-year follow-up,  
 330 IRB-approved longitudinal data) and is ongoing work.

## 331 References

- 332 [1] Faruk Ahmed, Andrew Sellergren, Lin Yang, Shawn Xu, Boris Babenko, Abbi Ward, Niels  
333 Olson, Arash Mohtashamian, Yossi Matias, Greg S Corrado, et al. Pathalign: A vision-language  
334 model for whole slide images in histopathology. *arXiv preprint arXiv:2406.19578*, 2024.
- 335 [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson,  
336 Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual  
337 language model for few-shot learning. *Advances in Neural Information Processing Systems*,  
338 35:23716–23736, 2022.
- 339 [3] Rawan AlSaad, Alaa Abd-Alrazaq, Sabri Boughorbel, Arfan Ahmed, Max-Antoine Renault,  
340 Rafat Damseh, and Javaid Sheikh. Multimodal large language models in health care: ap-  
341 plications, challenges, and future outlook. *Journal of medical Internet research*, 26:e59505,  
342 2024.
- 343 [4] Asmaa S Alsolami, Wafaa Shalash, Wafaa Alsaggaf, Sawsan Ashoor, Haneen Refaat, and  
344 Mohammed Elmogy. King abdulaziz university breast cancer mammogram dataset (kau-bcmd).  
345 *Data*, 6(11):111, 2021.
- 346 [5] Xiang An, Yin Xie, Kaicheng Yang, Wenkang Zhang, Xiuwei Zhao, Zheng Cheng, Yirui Wang,  
347 Songcen Xu, Changrui Chen, Didi Zhu, et al. Llava-onevision-1.5: Fully open framework for  
348 democratized multimodal training. *arXiv preprint arXiv:2509.23661*, 2025.
- 349 [6] N Aristokli, I Polycarpou, SC Themistocleous, D Sophocleous, and I Mamais. Comparison of  
350 the diagnostic performance of magnetic resonance imaging (mri), ultrasound and mammography  
351 for detection of breast cancer based on tumor type, breast density and patient’s history: A review.  
352 *Radiography*, 28(3):848–856, 2022.
- 353 [7] Melina Arnold, Eileen Morgan, Harriet Rumgay, Allini Mafra, Deependra Singh, Mathieu  
354 Laversanne, Jerome Vignat, Julie R Gralow, Fatima Cardoso, Sabine Siesling, et al. Current and  
355 future burden of breast cancer: Global statistics for 2020 and 2040. *The breast*, 66:15–23, 2022.
- 356 [8] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang  
357 Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding,  
358 localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- 359 [9] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao  
360 Deng, Wei Ding, Chang Gao, Chunjiang Ge, et al. Qwen3-vl technical report. *arXiv preprint*  
361 *arXiv:2511.21631*, 2025.
- 362 [10] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng  
363 Wang, Shijie Wang, Jun Tang, Humen Zhong, et al. Qwen2.5-vl technical report. *arXiv preprint*  
364 *arXiv:2502.13923*, 2025.
- 365 [11] Zak Buzzard, Konstantin Hemker, Nikola Simidjievski, and Mateja Jamnik. Paths: A hierar-  
366 chical transformer for efficient whole slide image analysis. *arXiv preprint arXiv:2411.18225*,  
367 2024.
- 368 [12] Chris Carr, PhD Felipe Kitamura, MD, George Partridge, inversion, Jayashree Kalpathy-Cramer,  
369 John Mongan, Katherine Andriole, Lavender, Maryam Vazirabad, Michelle Riopel, Robyn Ball,  
370 Sohier Dane, and Yan Chen. Rsn screening mammography breast cancer detection. <https://kaggle.com/competitions/rsna-breast-cancer-detection>, 2022. Kaggle.
- 372 [13] Junying Chen, Chi Gui, Ruyi Ouyang, Anningzhe Gao, Shunian Chen, Guiming Hardy Chen,  
373 Xidong Wang, Ruifei Zhang, Zhenyang Cai, Ke Ji, et al. Huatuogpt-vision, towards injecting  
374 medical visual knowledge into multimodal llms at scale. *arXiv preprint arXiv:2406.19280*,  
375 2024.
- 376 [14] Pengcheng Chen, Jin Ye, Guoan Wang, Yanjun Li, Zhongying Deng, Wei Li, Tianbin Li,  
377 Haodong Duan, Ziyang Huang, Yanzhou Su, et al. Gmai-mmbench: A comprehensive multimodal  
378 evaluation benchmark towards general medical ai. *Advances in Neural Information Processing*  
379 *Systems*, 37:94327–94427, 2024.

- 380 [15] Richard J Chen, Tong Ding, Ming Y Lu, Drew FK Williamson, Guillaume Jaume, Bowen Chen,  
381 Andrew Zhang, Daniel Shao, Andrew H Song, Muhammad Shaban, et al. A general-purpose  
382 self-supervised model for computational pathology. *arXiv preprint arXiv:2308.15474*, 2023.
- 383 [16] Richard J Chen, Tong Ding, Ming Y Lu, Drew FK Williamson, Guillaume Jaume, Andrew H  
384 Song, Bowen Chen, Andrew Zhang, Daniel Shao, Muhammad Shaban, et al. Towards a general-  
385 purpose foundation model for computational pathology. *Nature medicine*, 30(3):850–862,  
386 2024.
- 387 [17] Jiayu Ding, Shuming Ma, Li Dong, Xingxing Zhang, Shaohan Huang, Wenhui Wang, Nanning  
388 Zheng, and Furu Wei. Longnet: Scaling transformers to 1,000,000,000 tokens. *arXiv preprint*  
389 *arXiv:2307.02486*, 2023.
- 390 [18] Sixun Dong, Juhua Hu, Mian Zhang, Ming Yin, Yanjie Fu, and Qi Qian. MMTok: Multimodal  
391 coverage maximization for efficient inference of VLMs. In *The Fourteenth International*  
392 *Conference on Learning Representations*, 2026.
- 393 [19] Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong,  
394 Yuhang Zang, Pan Zhang, Jiaqi Wang, et al. Vlmevalkit: An open-source toolkit for evaluating  
395 large multi-modality models. In *Proceedings of the 32nd ACM International Conference on*  
396 *Multimedia*, pages 11198–11201, 2024.
- 397 [20] Omar Freihat, David Sipos, and Arpad Kovacs. Global burden and projections of breast cancer  
398 incidence and mortality to 2050: a comprehensive analysis of globocan data. *Frontiers in Public*  
399 *Health*, 13:1622954, 2025.
- 400 [21] Zhengrui Guo, Qichen Sun, Jiabo Ma, Lishuang Feng, Jinzhuo Wang, and Hao Chen. Context  
401 matters: Query-aware dynamic long sequence modeling of gigapixel images. *arXiv preprint*  
402 *arXiv:2501.18984*, 2025.
- 403 [22] Ibrahim Ethem Hamamci, Sezgin Er, Chenyu Wang, Furkan Almas, Ayse Gulnihani Simsek,  
404 Sevval Nil Esirgun, Irem Dogan, Omer Faruk Durugol, Benjamin Hou, Suprosanna Shit, et al.  
405 Generalist foundation models from a multimodal dataset for 3d computed tomography. *Nature*  
406 *Biomedical Engineering*, pages 1–19, 2026.
- 407 [23] Jing Hao, Yuxuan Fan, Yanpeng Sun, Kaixin Guo, Lizhuo Lin, Jinrong Yang, Qi Yong H  
408 Ai, Lun M Wong, Hao Tang, and Kuo Feng Hung. Towards better dental ai: A multimodal  
409 benchmark and instruction dataset for panoramic x-ray analysis. *Advances in Neural Information*  
410 *Processing Systems (NeurIPS)*, 2025.
- 411 [24] Jing Hao, Yuci Liang, Lizhuo Lin, Yuxuan Fan, Wenkai Zhou, Kaixin Guo, Zanting Ye,  
412 Yanpeng Sun, Xinyu Zhang, Yanqi Yang, et al. Oralgpt-omni: A versatile dental multimodal  
413 large language model. *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
414 *Pattern Recognition (CVPR)*, 2026.
- 415 [25] Iryna Hartsock and Ghulam Rasool. Vision-language models for medical report generation and  
416 visual question answering: A review. *Frontiers in artificial intelligence*, 7:1430984, 2024.
- 417 [26] Sunan He, Yuxiang Nie, Zhixuan Chen, Zhiyuan Cai, Hongmei Wang, Shu Yang, and Hao Chen.  
418 Meddr: Diagnosis-guided bootstrapping for large-scale medical vision-language learning. *arXiv*  
419 *e-prints*, pages arXiv–2404, 2024.
- 420 [27] Yutao Hu, Tianbin Li, Quanfeng Lu, Wenqi Shao, Junjun He, Yu Qiao, and Ping Luo. Omnimed-  
421 vqa: A new large-scale comprehensive evaluation benchmark for medical lvlm. In *Proceedings*  
422 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22170–22183,  
423 2024.
- 424 [28] Fuxiang Huang, Jiayi Zhu, Yunfang Yu, Yu Xie, Yuan Guo, Qingcong Kong, Mingxiang  
425 Wu, Xinrui Jiang, Shu Yang, Jiabo Ma, et al. A versatile foundation model for ai-enabled  
426 mammogram interpretation. *arXiv preprint arXiv:2509.20271*, 2025.
- 427 [29] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark,  
428 AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv*  
429 *preprint arXiv:2410.21276*, 2024.

- 430 [30] Jiwoong J Jeong, Brianna L Vey, Ananth Bhimireddy, Thomas Kim, Thiago Santos, Ramon  
431 Correa, Raman Dutt, Marina Mosunjac, Gabriela Oprea-Ilies, Geoffrey Smith, et al. The emory  
432 breast imaging dataset (embed): A racially diverse, granular dataset of 3.4 million screening  
433 and diagnostic mammographic images. *Radiology: Artificial Intelligence*, 5(1):e220047, 2023.
- 434 [31] Songtao Jiang, Yuan Wang, Sibao Song, Tianxiang Hu, Chenyi Zhou, Bin Pu, Yan Zhang, Zhibo  
435 Yang, Yang Feng, Joey Tianyi Zhou, et al. Hulu-med: A transparent generalist model towards  
436 holistic medical vision-language understanding. *arXiv preprint arXiv:2510.08668*, 2025.
- 437 [32] Béria Chingnabé Kalpélbé, Angel Gabriel Adaambiik, and Wei Peng. Vision language models  
438 in medicine. *arXiv preprint arXiv:2503.01863*, 2025.
- 439 [33] Laskarina Katsika, Eirini Boureka, Ioannis Kalogiannidis, Ioannis Tsakiridis, Ilias Tirodimos,  
440 Konstantinos Lallas, Zoi Tsimtsiou, and Themistoklis Dagklis. Screening for breast cancer: a  
441 comparative review of guidelines. *Life*, 14(6):777, 2024.
- 442 [34] Rana Khaled, Maha Helal, Omar Alfarghaly, Omnia Mokhtar, Abeer Elkorany, Hebatalla  
443 El Kassas, and Aly Fahmy. Categorized contrast enhanced mammography dataset for diagnostic  
444 and artificial intelligence research. *Scientific data*, 9(1):122, 2022.
- 445 [35] Joanne Kim, Andrew Harper, Valerie McCormack, Hyuna Sung, Nehmat Houssami, Eileen Mor-  
446 gan, Miriam Mutebi, Gail Garvey, Isabelle Soerjomataram, and Miranda M Fidler-Benaoudia.  
447 Global patterns and trends in breast cancer incidence and mortality across 185 countries. *Nature*  
448 *medicine*, 31(4):1154–1162, 2025.
- 449 [36] Rebecca Sawyer Lee, Francisco Gimenez, Assaf Hoogi, Kanae Kawai Miyake, Mia Gorovoy,  
450 and Daniel L Rubin. A curated mammography data set for use in computer-aided detection and  
451 diagnosis research. *Scientific data*, 4(1):170177, 2017.
- 452 [37] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan  
453 Naumann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision  
454 assistant for biomedicine in one day. *Advances in Neural Information Processing Systems*,  
455 36:28541–28564, 2023.
- 456 [38] Honglin Li, Yunlong Zhang, Pingyi Chen, Zhongyi Shui, Chenglu Zhu, and Lin Yang. Rethink-  
457 ing transformer for long contextual histopathology whole slide image analysis. *Advances in*  
458 *Neural Information Processing Systems*, 37:101498–101528, 2024.
- 459 [39] Haoneng Lin, Cheng Xu, and Jing Qin. Taming vision-language models for medical image  
460 analysis: A comprehensive review. *arXiv preprint arXiv:2506.18378*, 2025.
- 461 [40] Tianwei Lin, Wenqiao Zhang, Sijing Li, Yuqian Yuan, Binhe Yu, Haoyuan Li, Wanggui He,  
462 Hao Jiang, Mengze Li, Xiaohui Song, et al. Healthgpt: A medical large vision-language model  
463 for unifying comprehension and generation via heterogeneous knowledge adaptation. *arXiv*  
464 *preprint arXiv:2502.09838*, 2025.
- 465 [41] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan,  
466 Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around  
467 player? In *European conference on computer vision*, pages 216–233. Springer, 2024.
- 468 [42] Kosmia Loizidou, Galateia Skouroumouni, Costas Pitris, and Christos Nikolaou. Digital  
469 subtraction of temporally sequential mammograms for improved detection and classification of  
470 microcalcifications. *European radiology experimental*, 5(1):40, 2021.
- 471 [43] Ming Y Lu, Bowen Chen, Drew FK Williamson, Richard J Chen, Ivy Liang, Tong Ding, Guil-  
472 laume Jaume, Igor Odintsov, Long Phi Le, Georg Gerber, et al. A visual-language foundation  
473 model for computational pathology. *Nature medicine*, 30(3):863–874, 2024.
- 474 [44] Luyang Luo, Mingxiang Wu, Mei Li, Yi Xin, Qiong Wang, Varut Vardhanabhuti, Winnie CW  
475 Chu, Zhenhui Li, Juan Zhou, Pranav Rajpurkar, et al. A large model for non-invasive and  
476 personalized management of breast cancer from multiparametric mri. *Nature communications*,  
477 16(1):3647, 2025.

- 478 [45] Christine McCaffrey, Chowdhury Jahangir, Clodagh Murphy, Caoimbe Burke, William M  
479 Gallagher, and Arman Rahman. Artificial intelligence in digital histopathology for predicting  
480 patient prognosis and treatment efficacy in breast cancer. *Expert review of molecular diagnostics*,  
481 24(5):363–377, 2024.
- 482 [46] Yu Meng, Kaiyuan Li, Chenran Huang, Chen Gao, Xinlei Chen, Yong Li, and Xiaoping Zhang.  
483 Plphp: Per-layer per-head vision token pruning for efficient large vision-language models. *arXiv*  
484 *preprint arXiv:2502.14504*, 2025.
- 485 [47] Inês C Moreira, Igor Amaral, Inês Domingues, António Cardoso, Maria Joao Cardoso, and  
486 Jaime S Cardoso. Inbreast: toward a full-field digital mammographic database. *Academic*  
487 *radiology*, 19(2):236–248, 2012.
- 488 [48] Hieu T Nguyen, Ha Q Nguyen, Hieu H Pham, Khanh Lam, Linh T Le, Minh Dao, and Van  
489 Vu. Vindr-mammo: A large-scale benchmark dataset for computer-aided diagnosis in full-field  
490 digital mammography. *Scientific Data*, 10(1):277, 2023.
- 491 [49] Yasmine Omri, Parth Shroff, and Thierry Tambe. Token sequence compression for efficient  
492 multimodal computing. *arXiv preprint arXiv:2504.17892*, 2025.
- 493 [50] Parita Oza, Urvi Oza, Rajiv Oza, Paawan Sharma, Samir Patel, Pankaj Kumar, and Bakul Gohel.  
494 Digital mammography dataset for breast cancer diagnosis research (dmid) with breast mass  
495 segmentation analysis. *Biomedical Engineering Letters*, 14(2):317–330, 2024.
- 496 [51] Rui Sha, Xiang-meng Kong, Xin-yu Li, and Ya-bing Wang. Global burden of breast cancer and  
497 attributable risk factors in 204 countries and territories, from 1990 to 2021: results from the  
498 global burden of disease study 2021. *Biomarker Research*, 12(1):87, 2024.
- 499 [52] Yuzhang Shang, Mu Cai, Bingxin Xu, Yong Jae Lee, and Yan Yan. Llava-prumerge: Adaptive  
500 token reduction for efficient large multimodal models. In *Proceedings of the IEEE/CVF*  
501 *International Conference on Computer Vision*, pages 22857–22867, 2025.
- 502 [53] Moein Sorkhei, Yue Liu, Hossein Azizpour, Edward Azavedo, Karin Dembrower, Dimitra  
503 Ntoula, Athanasios Zouzos, Fredrik Strand, and Kevin Smith. Csaw-m: An ordinal classification  
504 dataset for benchmarking mammographic masking of cancer. *arXiv preprint arXiv:2112.01330*,  
505 2021.
- 506 [54] John Suckling. The mammographic images analysis society digital mammogram database. In  
507 *Excerpta Medica. International Congress Series, 1994*, volume 1069, pages 375–378, 1994.
- 508 [55] Yizheng Sun, Yanze Xin, Hao Li, Jingyuan Sun, Chenghua Lin, and Riza Theresa Batista-  
509 Navarro. Lvpruning: An effective yet simple language-guided vision token pruning approach  
510 for multi-modal large language models. In *Findings of the Association for Computational*  
511 *Linguistics: NAACL 2025*, pages 4299–4308, 2025.
- 512 [56] Hao Tang and Chengchao Shen. Learning compact vision tokens for efficient large multimodal  
513 models. *arXiv*, 2025.
- 514 [57] Core Team, Zihao Yue, Zhenru Lin, Yifan Song, Weikun Wang, et al. Mimo-vl technical report,  
515 2025.
- 516 [58] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,  
517 Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly  
518 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 519 [59] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya  
520 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open  
521 models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- 522 [60] V Team, Wenyi Hong, Wenmeng Yu, Xiaotao Gu, Guo Wang, Guobing Gan, et al. Glm-4.5v  
523 and glm-4.1v-thinking: Towards versatile multimodal reasoning with scalable reinforcement  
524 learning, 2025.

- 525 [61] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing  
526 Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception  
527 of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- 528 [62] Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Hui Hui, Yanfeng Wang, and Weidi Xie. Towards  
529 generalist foundation model for radiology by leveraging web-scale 2d&3d medical data. *Nature*  
530 *Communications*, 16(1):7866, 2025.
- 531 [63] xAI. Grok 4 fast model card. [https://data.x.ai/  
532 2025-09-19-grok-4-fast-model-card.pdf](https://data.x.ai/2025-09-19-grok-4-fast-model-card.pdf), 2025.
- 533 [64] Hanguang Xiao, Feizhong Zhou, Xingyue Liu, Tianqi Liu, Zhipeng Li, Xin Liu, and Xiaoxuan  
534 Huang. A comprehensive survey of large language models and multimodal large language  
535 models in medicine. *Information Fusion*, 117:102888, 2025.
- 536 [65] Conghao Xiong, Hao Chen, and Joseph JY Sung. A survey of pathology foundation model:  
537 Progress and future directions. *arXiv preprint arXiv:2504.04045*, 2025.
- 538 [66] Hanwen Xu, Naoto Usuyama, Jaspreet Bagga, Sheng Zhang, Rajesh Rao, Tristan Naumann,  
539 Cliff Wong, Zelalem Gero, Javier González, Yu Gu, et al. A whole-slide foundation model for  
540 digital pathology from real-world data. *Nature*, 630(8015):181–188, 2024.
- 541 [67] Weiwen Xu, Hou Pong Chan, Long Li, Mahani Aljunied, Ruifeng Yuan, Jianyu Wang, Cheng-  
542 hao Xiao, Guizhen Chen, Chaoqun Liu, Zhaodonghui Li, et al. Lingshu: A generalist founda-  
543 tion model for unified multimodal medical understanding and reasoning. *arXiv preprint*  
544 *arXiv:2506.07044*, 2025.
- 545 [68] Jiarui Ye and Hao Tang. Multimodal large language models for medicine: A comprehensive  
546 survey. *arXiv preprint arXiv:2504.21051*, 2025.
- 547 [69] Haojun Yu, Youcheng Li, Zihan Niu, Nan Zhang, Xuanton Gong, Huan Li, Zhiying Zou,  
548 Haifeng Qi, Zhenxiao Cao, Zijie Lan, Xingjian Yuan, Jiating He, Haokai Zhang, Shengtao  
549 Zhang, Zicheng Wang, Dong Wang, Ziwei Zhao, Congying Chen, Yong Wang, Wangyan  
550 Qin, and Qingli Zhu. A chain-of-thought reasoning breast ultrasound dataset covering all  
551 histopathology categories. *Scientific Data*, 13:236, 2026.
- 552 [70] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang,  
553 and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities.  
554 *arXiv preprint arXiv:2308.02490*, 2023.
- 555 [71] Weihao Yu, Zhengyuan Yang, Lingfeng Ren, Linjie Li, Jianfeng Wang, Kevin Lin, Chung-Ching  
556 Lin, Zicheng Liu, Lijuan Wang, and Xinchao Wang. Mm-vet v2: A challenging benchmark to  
557 evaluate large multimodal models for integrated capabilities. *arXiv preprint arXiv:2408.00765*,  
558 2024.
- 559 [72] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens,  
560 Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal  
561 understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF*  
562 *Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024.
- 563 [73] Jiaheng Zhou, Wei Fang, Luyuan Xie, Yanfeng Zhou, Lianyan Xu, Minfeng Xu, Ge Yang, and  
564 Yuxing Tang. Dual-res tandem mamba-3d: Bilateral breast lesion detection and classification  
565 on non-contrast chest CT. In *The Thirty-ninth Annual Conference on Neural Information*  
566 *Processing Systems*, 2026.
- 567 [74] Jiayi Zhu, Fuxiang Huang, Qiong Luo, and Hao Chen. A benchmark for breast cancer screening  
568 and diagnosis in mammogram visual question answering. *Nature Communications*, 16(1):11683,  
569 2025.

570

# BreastGPT: A Multimodal Large Language Model for the Full Spectrum of Breast Cancer Clinical Routine

571

572

## Supplementary Materials

573

### Contents

574	<b>Introduction</b> .....	<b>1</b>
575	<b>BreastStage Dataset Curation</b> .....	<b>3</b>
576	Scope and scale .....	3
577	Data construction pipeline .....	3
578	<b>BreastGPT</b> .....	<b>5</b>
579	Dual-Branch Visual Encoding .....	5
580	Universal Concept-Based Compression .....	6
581	<b>Experiments</b> .....	<b>6</b>
582	Experimental Setups .....	6
583	Evaluation Metrics .....	7
584	Evaluation Results .....	7
585	Ablation Studies .....	8
586	<b>Conclusion</b> .....	<b>9</b>
587	<b>Supplementary Overview</b> .....	<b>16</b>
588	LLM Usage Statement .....	16
589	Ethics Approval and Privacy .....	16
590	Additional Limitations .....	17
591	Compute Resources and Environmental Impact .....	17
592	Code and Data Availability .....	17
593	Acknowledgments .....	17
594	<b>Related Work</b> .....	<b>17</b>
595	<b>BreastStage Curation Details</b> .....	<b>18</b>
596	Data Sources and Modality Coverage .....	18
597	Construction Data Generation Pipeline .....	21
598	Task Taxonomy and Distribution .....	27
599	<b>BreastGPT Technical Details</b> .....	<b>29</b>
600	Base Model Configuration .....	29
601	ViT-aPixel Encoder Implementation .....	29
602	Coverage-Maximizing Token Selection .....	30
603	Training Configuration .....	31
604	System Prompt Templates .....	32
605	<b>BreastStage-Bench Evaluation</b> .....	<b>32</b>
606	Benchmark Construction .....	32
607	Evaluation Metrics .....	33
608	Breast Expert Validation .....	33
609	Baseline Model Configuration .....	34
610	<b>Additional Results and Ablations</b> .....	<b>34</b>
611	Per-Modality Performance Breakdown .....	34
612	Extended Baseline Comparisons .....	35
613	Open-ended VQA Raw Metric Breakdown .....	35
614	ED Grounding IoU on CT and MRI .....	36
615	BreastGPT Grounding Recognition Breakdown .....	36
616	Inference Efficiency Analysis .....	37
617	<b>Detailed Ablation Results</b> .....	<b>37</b>
618	Numerical Results for WSI Branch Ablation .....	37
619	Numerical Results for Visual Token Budget Sweep .....	38
620	<b>Qualitative Analysis and Case Studies</b> .....	<b>39</b>
621	Role-Switching Examples .....	39
622	Case Analysis .....	40

## 623 A Supplementary Overview

624 This supplementary material provides the technical and procedural details that support the main claims  
625 of the paper. We first describe the construction of **BreastStage**, including data provenance, modality  
626 coverage, instruction generation, quality control, and task taxonomy. We then provide implementation  
627 details for **BreastGPT**, with particular emphasis on the GigaPixel WSI branch and the coverage-  
628 maximizing visual token selector that enables a unified token budget across standard radiology images  
629 and gigapixel pathology slides. Finally, we describe the BreastStage-Bench evaluation protocol,  
630 additional quantitative analyses, qualitative examples, limitations, ethics, and release plans.

631 This supplementary is intended to make three aspects of the paper auditable: (i) how heterogeneous  
632 breast imaging data are mapped into a stage-aware clinical workflow, (ii) how BreastGPT processes  
633 image scales that differ by several orders of magnitude, and (iii) how benchmark results, ablations,  
634 and error analyses are produced.

### 635 A.1 LLM Usage Statement

636 In this work, we primarily employ LLMs in two aspects: (i) data construction, where Qwen2.5-VL-  
637 72B and Qwen3-Max drive the orchestration pipeline described in Section C.2; and (ii) manuscript  
638 polishing, where we use LLMs to improve grammar and clarity of the writing without altering the  
639 technical claims.

### 640 A.2 Ethics Approval and Privacy

641 BreastStage combines two kinds of sources with different ethical pathways. The BUS, mammography,  
642 CT, and histopathology subsets come entirely from public datasets that were de-identified by their  
643 original publishers and released under their respective data-use terms; we use them within those terms  
644 and do not perform any additional re-identification of individual samples. The MRI subset is the only  
645 institutional cohort: it was acquired from two collaborating hospitals [anonymous for review] under  
646 IRB approval, with patient consent obtained where required by the contributing institutions.

647 **De-identification of the institutional MRI cohort.** Because the MRI cohort is the only data source  
648 that originates inside hospitals, our de-identification protocol is scoped specifically to it. Before MRI  
649 data leave either source institution, identifying information is removed at three levels:

- 650 • **DICOM header level.** Patient names, dates of birth, medical-record numbers, hospital  
651 identifiers, referring-physician fields, and exam dates are stripped; absolute timestamps  
652 are converted to relative offsets so longitudinal ordering is preserved without disclosing  
653 calendar dates.
- 654 • **Image content level.** The MRI volumes are reconstructed within a breast-coil field of view  
655 that does not include the patient’s face, so no additional image-side scrubbing is required.
- 656 • **Free-text level.** The accompanying Chinese radiology and pathology reports are passed  
657 through the bilingual parsers in Figure 8, which emit a strictly schema-bounded English  
658 JSON so that no free-text PHI (referring physician name, hospital department, etc.) propa-  
659 gates into the released instruction-following pairs.

660 **Access control.** Access to identifiable raw MRI data is restricted to authorised research personnel  
661 under signed data-use agreements. The public mirror of BreastStage will release only the de-identified,  
662 instruction-formatted records and the derived bounding-box / mask annotations; raw DICOMs from  
663 the institutional MRI cohort are not redistributed. The other four subsets (BUS, mammography, CT,  
664 histopathology) are accessed via their original public-dataset distribution channels.

665 **Risks acknowledged.** Because the model is trained on clinical imaging data, performance may  
666 vary across institutions, scanners, demographic groups, and disease subtypes; we recommend that  
667 any downstream use perform site-specific validation before clinical deployment. The benchmark  
668 is intended to encourage transparent comparison of breast cancer workflow models, not to replace  
669 clinician judgement.

### 670 A.3 Additional Limitations

671 **Dataset.** While BreastStage covers five imaging modalities, it does not include molecular imaging  
672 (PET / CT-PET) or emerging modalities such as photoacoustic imaging. The current workflow  
673 also focuses on imaging-centred decision support; longitudinal treatment histories, genomic assays,  
674 laboratory results, and medication records are only partially available depending on the source  
675 cohort. Future work should expand BreastStage toward longitudinal multimodal patient records and  
676 evaluate whether workflow-aligned MLLMs can support temporal reasoning over treatment response,  
677 recurrence risk, and survivorship care.

678 **Clinical deployment.** BreastGPT is a research prototype and has not undergone clinical validation  
679 or regulatory approval. Deployment in a real reading room would additionally require prospective  
680 clinical trials, FDA 510(k) clearance or equivalent regulatory approval, integration with existing PACS  
681 / RIS systems, and continuous monitoring for model drift and performance degradation. BreastGPT  
682 should therefore be interpreted as a foundation model for research and benchmarking, not as an  
683 autonomous diagnostic system.

### 684 A.4 Compute Resources and Environmental Impact

685 **Model training.** BreastGPT training was carried out on 32 NVIDIA H100 GPUs (4 nodes  $\times$  8  
686 GPUs) with DeepSpeed ZeRO-2, FlashAttention, and bfloat16 precision. End-to-end wall-clock  
687 training time across the two stages was 3 days, 7 hours, 19 minutes, 53 seconds, corresponding to  
688 approximately 2,530 H100-hours.

689 **Dataset construction and evaluation.** Visual quality control with Qwen2.5-VL-72B and Qwen3-  
690 Max-driven instruction generation are run through external APIs and shared infrastructure; per-call  
691 API expenditure was not separately tracked. Specialist annotation of MRI tumours (10 board-certified  
692 breast radiologists) and the 3-specialist Bench audit were conducted by clinical collaborators rather  
693 than billed compute.

### 694 A.5 Code and Data Availability

695 **Code.** An anonymised snapshot of the BreastGPT codebase, including training scripts, evalu-  
696 ation pipelines, and the BreastStage data-construction tooling, is available for review at [https://](https://anonymous.4open.science/r/BreastGPT)  
697 [anonymous.4open.science/r/BreastGPT](https://anonymous.4open.science/r/BreastGPT). The full code release with model checkpoints will  
698 be moved to a non-anonymised public repository upon publication.

699 **Data.** The BreastStage dataset and BreastStage-Bench evaluation suite will be made publicly  
700 available under a Creative Commons BY-NC 4.0 license at [https://\[ANONYMOUS\]/breaststage](https://[ANONYMOUS]/breaststage).  
701 For components derived from institutional clinical data, release will follow the corresponding IRB  
702 protocol, de-identification requirements, and data-use agreements. When raw image release is  
703 restricted, we will provide metadata, task definitions, evaluation code, and access instructions.

704 **Model weights.** Pre-trained BreastGPT model weights will be released on Hugging Face:  
705 [https://huggingface.co/\[ANONYMOUS\]/BreastGPT](https://huggingface.co/[ANONYMOUS]/BreastGPT).

### 706 A.6 Acknowledgments

707 We thank the breast oncology experts from the contributing institutions [anonymous for review] for  
708 their guidance during dataset construction and clinical validation. We acknowledge the computational  
709 resources provided by [anonymous for review] and the funding from [anonymous for review].

## 710 B Related Work

711 Recent medical MLLMs have extended general-purpose vision-language modeling into clinical image  
712 understanding, report generation, visual question answering, and interactive decision support [68, 32,  
713 3, 64, 39, 25]. Representative systems such as LLaVA-Med, HuatuoGPT-Vision, MedDr, HealthGPT,  
714 and Lingshu show that biomedical instruction tuning can improve clinical-domain perception and

715 dialogue [37, 13, 26, 40, 67]. Meanwhile, broad evaluation suites such as OmniMedVQA, GMAI-  
716 MMBench, MMMU, MMBench, and VLMEvalKit expose the gap between generic multimodal  
717 ability and reliable medical reasoning [27, 14, 72, 41, 19]. Domain-specific studies in dentistry  
718 further show that clinically meaningful evaluation often requires benchmarks organized around  
719 specialty-specific anatomy, modalities, reporting conventions, and treatment workflows [23, 24].  
720 BreastGPT follows this specialty-grounded direction for breast oncology, where the challenge is not  
721 only domain-specific perception but also workflow-aware reasoning across screening, diagnosis, and  
722 treatment.

723 Breast imaging AI has progressed rapidly, yet most prior work remains modality-specific. In  
724 mammography, VersaMammo builds a large multi-institutional mammogram foundation model and  
725 evaluates detection, segmentation, classification, retrieval, and VQA tasks [28], while MammoVQA  
726 standardizes mammogram VQA across 15 public datasets and shows that both general-purpose and  
727 medical-specific LVLMs remain unreliable for mammogram interpretation [74]. For breast ultrasound,  
728 BUS-CoT provides chain-of-thought reasoning annotations across all 99 WHO histopathology  
729 categories, connecting observations, imaging features, BI-RADS assessment, pathology labels,  
730 and reasoning traces [69]. For MRI, MOME integrates multiparametric breast MRI through a  
731 mixture-of-modality-experts design and supports malignancy identification, biopsy recommendation,  
732 triple-negative breast cancer classification, and neoadjuvant chemotherapy response prediction [44].  
733 These works show the value of breast-specific modeling, but they do not by themselves provide a  
734 single model or benchmark spanning the full clinical pathway.

735 Computational pathology has similarly moved from patch-level representation learning toward  
736 slide-level and vision-language foundation models. Recent pathology foundation models, including  
737 UNI, CONCH, and large-scale whole-slide encoders, learn transferable morphology representations  
738 for diagnostic, prognostic, retrieval, and report-generation tasks [16, 43, 15, 66]. WSI vision-  
739 language and long-sequence modeling efforts such as PathAlign, query-aware gigapixel modeling,  
740 hierarchical WSI transformers, and long-context histopathology transformers further bridge patch-  
741 level embeddings with slide-level reasoning and language interaction [1, 21, 11, 38, 65]. These  
742 studies highlight a central difficulty for BreastGPT: pathology slides contain sparse, distributed, and  
743 scale-dependent evidence, while breast radiology images use much smaller visual inputs and different  
744 clinical semantics.

745 Efficient visual token selection is therefore essential for unified breast workflow modeling. Recent  
746 methods reduce MLLM inference cost through token pruning, learned token selection, adaptive  
747 compression, query-aware retrieval, and layer-wise token reduction [18, 52, 56, 55, 46, 49]. MMTok  
748 formulates token selection as multimodal coverage maximization, balancing text-vision alignment  
749 with vision-vision diversity [18]. BreastGPT adapts this principle to clinical workflow modeling: for  
750 radiology images, the selector removes redundant visual tokens while preserving prompt-relevant le-  
751 sion and anatomical information; for WSIs, it compresses LongNet-contextualized patch embeddings  
752 into a fixed-size representation that retains both query-relevant pathology regions and global tissue  
753 coverage.

754 Overall, existing work establishes strong foundations for medical MLLMs, breast imaging models,  
755 pathology foundation models, and efficient token compression, but these directions remain largely  
756 separate. BreastGPT integrates them into a unified breast cancer workflow model, and BreastStage-  
757 Bench evaluates whether such integration can support stage-aware reasoning across heterogeneous  
758 modalities and clinical roles.

## 759 **C BreastStage Curation Details**

### 760 **C.1 Data Sources and Modality Coverage**

761 As described in Section 2 of the main paper, BreastStage integrates 17 sub-datasets across 5 imaging  
762 modalities. The goal of this curation is not merely to aggregate public resources, but to reorganize  
763 heterogeneous imaging evidence around the clinical progression of breast cancer care: screening,  
764 diagnosis, and treatment. This section details the source datasets, selection criteria, quality control  
765 procedures, and annotation protocols used to construct each modality subset.

766 **CT Imaging.** Computed tomography volumes are sourced from the CT-RATE dataset [22], which  
767 contains 25,692 non-contrast chest CT volumes from 21,304 unique patients. Since breast cancer

768 predominantly affects women, we first filter the dataset to retain only female patients (sex metadata  
 769 field equal to “F”). We then deploy Qwen2.5-VL-72B [10] as an automated quality inspector,  
 770 assigning confidence scores (0–10 scale) based on image quality, anatomical coverage, and the  
 771 presence of breast tissue in the field of view; volumes scoring below the threshold are excluded.  
 772 For breast cancer screening and diagnosis, we further pass each volume through DRT-M3D [73] to  
 773 generate automated tumor segmentations and risk assessments, identifying suspicious breast lesions,  
 774 calcifications, and lymph node involvement visible in the chest CT field of view. After quality  
 775 filtering and clinical relevance screening, the CT subset of BreastStage contains 20,546 reconstructed  
 776 volumes from female patients (Table 5), covering both screening-stage opportunistic detection and  
 777 diagnosis-stage tumor characterization tasks.

778 **Breast Ultrasound (BUS).** BUS images are sourced from the **BUS-CoT** dataset [69], the first chain-  
 779 of-thought reasoning breast ultrasound dataset that covers *all* 99 histopathology categories defined by  
 780 the WHO classification of breast tumors. BUS-CoT contains 11,439 images of 10,019 lesions from  
 781 4,838 patients across 18 different ultrasound device types, including B-mode US, Doppler US, and  
 782 Elastography. Images are collected from open-access publications, publicly available case studies  
 783 (Radiopaedia, PubMed), and open-access biopsy-confirmed repositories, with each case annotated  
 784 through a rigorous five-level protocol by six senior breast ultrasound radiologists with 8–26 years of  
 785 clinical experience: (1) *Observation* (lesion presence, calcification presence), (2) *Feature* (boundary,  
 786 edge, echo characteristics, calcification feature), (3) *Diagnosis* (BI-RADS score), (4) *Pathology*  
 787 (benign/malignant + WHO histopathology subtype), and (5) *Chain-of-Thought reasoning* explicitly  
 788 linking imaging features to the final pathological label. Table 3 summarizes the dataset’s pathology  
 789 distribution.

Table 3: Pathology distribution of the BUS-CoT source dataset [69], used as the BUS data source for BreastStage. Numbers in parentheses indicate per-subtype lesion counts.

	<b>Benign (4,856)</b>	<b>Malignant (4,814)</b>	<b>Others (349)</b>
<i>Categories</i>	Fibroadenoma (1,047)	Invasive ductal carcinoma (896)	Others (349)
	Phyllodes tumour (74)	Invasive lobular carcinoma (75)	
	Intraductal papilloma (72)	Ductal carcinoma in situ (72)	
	Atypical ductal hyperplasia (62)	Mucinous carcinoma (64)	
	Radial scar (56)	Paget disease of the breast (57)	
	Other benign (3,545)	Other malignant (3,650)	

790 After quality filtering and mapping to our clinical workflow taxonomy (screening-stage BI-RADS  
 791 triage and diagnosis-stage lesion characterization), the BUS subset of BreastStage contains 10,405  
 792 unique image files and 190,730 instruction-following pairs (Table 5), with expert-verified BI-RADS  
 793 scores, lesion bounding boxes, and chain-of-thought rationales.

794 **Mammography.** Digital mammography is sourced from a subset of **MammoVQA** [74], which  
 795 originally unifies 15 public mammogram datasets. We adopt the 11 datasets that are licensed for  
 796 derivative works and have stage-aware annotations compatible with our taxonomy (BMCD, CBIS-  
 797 DDSM, CDD-CESM, CSAW-M, DMID, EMBED, INbreast, KAU-BCMD, MIAS, RSNA, VinDr-  
 798 Mammo); we then re-organise their labels around the screening / diagnosis / treatment continuum  
 799 and add new task templates on top of MammoVQA. Table 4 lists every adopted dataset together  
 800 with its sample count and the task templates it supports. The union of these 11 sources covers the 9  
 801 clinically validated MammoVQA task templates—View (CC/MLO), Laterality, BI-RADS, Pathology,  
 802 Masking potential, Background tissue, Subtlety, Density, and Abnormality detection (calcification,  
 803 architectural distortion, mass, asymmetry, etc.).

Table 4: Internal mammography source datasets aggregated by MammoVQA [74] and adopted as the BreastStage mammography source. “Type” follows the MammoVQA labelling: *breast* = per-image labels, *finding* = per-bounding-box labels, *exam* = per-examination labels.

Dataset	Type	Size	Supported Tasks
BMCD [42]	breast	400	Density, BI-RADS, Laterality
CBIS-DDSM [36]	breast	3,103	Density, BI-RADS, Laterality, View
	finding	6,464	Abnormality, Pathology, Subtlety
CDD-CESM [34]	breast	1,003	Density, BI-RADS, Laterality, Pathology
DMID [50]	breast	510	Abnormality, Background tissue, Laterality, Pathology, View
	finding	868	Abnormality, Pathology
INbreast [47]	breast	410	Density, Abnormality, BI-RADS, Laterality, View
	breast	322	Abnormality, Background tissue, Pathology
MIAS [54]	finding	234	Abnormality, Pathology
	breast	10,020	Laterality, Masking potential
CSAW-M [53]	breast	10,020	Laterality, Masking potential
KAU-BCMD [4]	breast	2,370	BI-RADS
VinDr-Mammo [48]	breast	20,000	Density, Abnormality, BI-RADS, Laterality, View
	finding	4,505	Abnormality
RSNA [12]	breast	54,705	Density, BI-RADS, Laterality, View
EMBED [30]	exam	72,518	Density (exam), BI-RADS (exam)
<b>Internal total</b>	—	<b>177,432</b>	—

804 For BreastStage we use these 11 datasets with patient-level stratification preserved across train and test  
805 splits. After mapping the original classification labels to our screening / diagnosis / treatment stage  
806 taxonomy and synthesising stage-specific QA pairs, the mammography subset comprises 592,470  
807 unique 2D images and 453,956 instruction-following pairs.

808 **MRI.** Multiparametric MRI sequences (T1-weighted, T2-weighted, DWI, DCE-MRI / T1dyn) are  
809 acquired from two collaborating clinical institutions [anonymous for review] under IRB approval.  
810 Board-certified breast radiologists at the contributing institutions provide expert annotations for tumor  
811 segmentation, molecular subtype prediction, and neoadjuvant chemotherapy response assessment.  
812 The MRI subset comprises 36,124 unique multiparametric sequence volumes (T1, T1dyn, T2W, DWI,  
813 ADC) across screening, diagnosis, and treatment-stage tasks, yielding 926,521 instruction-following  
814 pairs.

815 **Whole Slide Imaging (WSI).** Histopathology WSIs are sourced from three public repositories—**BCNB**,  
816 **TCGA-BRCA**, and **TCGA-HISTAI**—covering invasive ductal carcinoma, invasive lobular carcinoma, ductal carcinoma in situ, and benign / atypical lesions across histological  
817 grades I–III. Slides are scanned at 20× or 40× magnification with native resolutions ranging from  
818 80,000×80,000 to 120,000×120,000 pixels. The WSI subset contains 2,510 slides with pathologist-  
819 verified molecular subtype labels (ER/PR/HER2 status, Ki-67 index, Nottingham grade), yielding  
820 30,073 instruction-following pairs (captioning, closed- and open-ended VQA, plus prognosis /  
821 surgical-planning / systemic-therapy treatment-stage tasks).  
822

823 **Aggregate Statistics.** Table 5 summarizes the per-modality data sources adopted by BreastStage.  
824 These modality-level counts correspond to the curated image pools before conversion into instruction-  
825 following examples. The final instruction corpus expands these image pools into multiple task-specific  
826 QA pairs, captions, grounding targets, and report-generation instances according to the stage-aware  
827 taxonomy described below.

Table 5: Summary of BreastStage data sources by modality. “Unique 2D / 3D images” counts the distinct image (2D radiology, mammography, histopathology slide) and 3D-volume (CT, MRI multiparametric sequence) files used by the model, after filtering each modality to its retained sub-datasets. Each unique image typically backs many QA pairs.

Modality	Primary Source	Selection Criterion	Unique 2D / 3D Images	Instruction Pairs
CT	CT-RATE [22]	Female-breast volumes, post-QC	20,546 volumes	254,041
BUS	BUS-CoT [69]	All 99 histopathology cats	10,405 images	190,730
Mammo	MammoVQA-derived 11 sets [74]: BMCD, CBIS-DDSM, CDD-CESM, CSAW-M, DMID, EMBED, IN-breast, KAU-BCMD, MIAS, RSNA, VinDr	Patient-level stratified, stage-specific re-organisation with new tasks added on top of MammoVQA	592,470 images	454,590
MRI	In-house multi-institution multiparametric MRI	Multiparametric studies with complete sequences	36,124 volumes	926,634
WSI	BCNB / TCGA-BRCA / TCGA-HISTAI	Subtype-labeled, 20–40× slides	2,510 slides	30,073
<b>Total</b>	<b>17 sub-datasets</b>	—	<b>662,055 images</b>	<b>1,856,068</b>

## 828 C.2 Instruction Data Generation Pipeline

829 This section expands the four-stage curation pipeline introduced in Section 2 of the main paper (*Work-*  
830 *flow Data Generation* → *Stage Instruction Data Construction* → *Textual Description Generation*  
831 → *Data Splitting and Verification*, illustrated in Figure 2). The four stages are powered by a Data  
832 Orchestration Engine that uses two LLMs in complementary roles: **Qwen2.5-VL-72B** [10] is used  
833 wherever a decision requires looking at the image (modality-specific quality control, visual-attribute  
834 extraction), and **Qwen3-Max** is used for every text-only transformation (Chinese-to-English report  
835 parsing, structured-record-to-open-ended-VQA rewriting, ground-caption synthesis, report gener-  
836 ation). Both LLMs are queried with strict JSON output schemas so downstream stages consume  
837 structured outputs without bespoke parsing. Per-modality QC prompts and the Chinese-to-English  
838 parsers appear in Figure 8, and the stage-aware persona library together with the system-prompt  
839 assembly rule appear in Figure 9. The remainder of this section walks through the four stages in  
840 one-to-one correspondence with the main-paper subsections.

841 **Workflow Data Generation.** The cohort consists of public BUS, mammography, histopathology,  
842 and CT data plus an in-house multi-institution MRI dataset acquired from two collaborating hospitals  
843 [anonymous for review] under IRB approval; identifying information is removed at the source  
844 institutions before any data leaves the hospitals (Section A.2). On top of the raw studies, this stage  
845 produces three artefacts that the rest of the pipeline relies on: per-image quality scores, lesion-level  
846 bounding-box / mask annotations, and stage-aware task labels.

847 *Modality-specific quality control.* The visual specialist agents from Section 2 are implemented as  
848 Qwen2.5-VL-72B prompted with modality-specific quality-control rules (Fig. 8). The BUS selector  
849 rejects convex / sector-probe images outside the breast-ultrasound depth range; the mammography  
850 selector strictly rejects chest X-rays and applies PGMI positioning checks plus Eklund implant-  
851 displaced views; the breast MRI selector requires bilateral-coil framing or sagittal / axial single-breast  
852 views and verifies fat-saturation; the CT selector requires breast-inclusive FOV. Histopathology  
853 WSIs are sourced from a curated public collection where modality and magnification are already  
854 standardised, so they bypass this LLM-driven QC step. The selector emits a JSON {*validity*,  
855 *reason*}; samples flagged as low-quality are not silently dropped but routed to a *Low-Image Check*  
856 step in which a breast specialist reviews the image and the LLM’s stated reason, and either confirms  
857 the rejection or restores the sample.

858 *Bounding-box and mask generation.* The four imaging modalities (BUS, mammography, CT, MRI)  
859 carry lesion-level spatial annotations, with the source differing by modality. BUS uses the expert  
860 hand-drawn lesion masks released with the BUS-CoT dataset [69]. Mammography (EMBED subset)  
861 ships with official lesion bounding boxes; we re-format the released [*ymin*, *xmin*, *ymax*, *xmax*]  
862 tuples into the [*xmin*, *ymin*, *xmax*, *ymax*] convention used throughout BreastStage, and no mask  
863 step is performed on this modality. CT lesion masks come from DRT-M3D [73], which also emits a  
864 per-volume cancer-risk score that we retain in the record metadata for the screening-stage CT tasks.  
865 MRI lesion masks are drawn manually on T1 and T1dyn sequences by 10 board-certified breast

866 specialists at the contributing hospitals. For the mask-based modalities (BUS, CT, MRI), bounding  
867 boxes are derived by connected-component analysis: components with fewer than 50 voxels are  
868 discarded as noise, and the tight axis-aligned envelope of every retained component is emitted as a  
869 4-tuple [xmin, ymin, xmax, ymax] for BUS or a 6-tuple [xmin, ymin, zmin, xmax, ymax,  
870 zmax] for the 3D volumes (CT, MRI). Histopathology WSIs do not carry pixel-level masks or  
871 bounding boxes in BreastStage; they are routed only to the closed/open VQA and report-generation  
872 tasks, not to the ground-caption task.

873 *Expert-driven sub-task labelling.* After QC and mask generation, the per-dataset clinical tables and  
874 reports are reviewed by breast specialists who tag each record with one of **136 sub-task labels**  
875 reflecting BI-RADS items, mass / calcification descriptors, biomarker fields, and treatment-stage  
876 decisions. Each label is a `task template`: a high-level sub-task such as “mass margin description”  
877 or “HER2 status”, typically spanning several leaf fields of the structured report defined in the next  
878 stage. Per-modality counts after this stage are reported in Supplementary Table 5, and the resulting  
879 136-task taxonomy in Table 6.

880 **Stage Instruction Data Construction.** This stage converts each curated study into a schema-  
881 bounded structured record, so that no free-text generation step downstream is allowed to introduce  
882 clinical content beyond what this record already encodes. The flow follows the four steps introduced  
883 in the main paper. (i) **Report translation.** Chinese radiology and pathology reports are translated to  
884 English by Qwen3-Max acting as a radiologist agent, using the bilingual prompts in Figure 8. (ii)  
885 **Structured report extraction.** Every translated report is then parsed by the radiologist agent into an  
886 expert-designed BI-RADS-aligned structured report whose schema is written per modality by breast  
887 specialists following clinical reporting guidelines. For example, the MRI structured report (Figure 7)  
888 covers per-side breast-level findings (FGT, BPE, post-surgical changes, lymph nodes, BI-RADS,  
889 management) plus per-lesion descriptors for masses, non-mass enhancement, and non-enhancing  
890 lesions. (iii) **Question bank derivation.** For each `task template` from Stage 1, breast specialists  
891 then write a closed-ended question that draws its answer from one or more leaf fields of the structured  
892 report, with options taken directly from the schema’s enum values; this constrains BreastStage to  
893 questions that the structured record already answers. The resulting tuples carry the form `<Task,`  
894 `Question, options, Answer, Images>`. (iv) **Task–Stage mapping.** Specialists finally map  
895 each `task template` to one of the three clinical stages (*screening / diagnosis / treatment*), which is  
896 also the key used for task-key balanced sampling on BreastStage-Bench; the per-stage record counts  
897 of the final task taxonomy are reported in Table 6.

898 From the structured records, BreastStage derives two parallel VQA partitions, both grounded in  
899 the same schema-bounded source. The *closed-ended* partition packs the schema’s enum values as  
900 an option list into the question text, so each tuple takes the form `<Task, Question+options,`  
901 `Answer, Images>` with the answer being one of the options. The *open-ended* partition takes  
902 the same structured tuple, drops the option list, and asks Qwen3-Max to rewrite the categorical  
903 answer into a fluent clinical description, yielding `<Task, Question, Answer, Images>` where  
904 the answer is free-form text but the underlying clinical fact is unchanged; the rewriting prompt is  
905 reproduced in Fig. 6. Because both partitions are pinned to the same structured record, the LLM  
906 cannot fabricate clinical content beyond the schema, and the answer cannot leak through the question  
907 text into a text-only shortcut — a property reflected in the  $k=1$  collapse of BreastGPT’s closed-ended  
908 accuracy in Fig. 5a, where the model degrades sharply once the visual evidence is removed.

### Open-VQA Generator

You are a medical imaging data assistant specialising in breast {modality}. Your task is to rewrite a structured attribute into a natural-language Question–Answer pair suitable for a Breast {modality} VQA dataset.

#### Rules.

1. The **Question** *must* be open-ended and descriptive: do *not* ask yes/no questions, do *not* include options or binary choices.
2. The **Answer** *must* be a complete sentence that is strictly based on the provided **Answer** value, and does *not* introduce any new medical interpretation, diagnosis, or inference.
3. Do *not* introduce information not explicitly contained in the input.
4. Do *not* repeat or list the provided options in the output.
5. Use varied sentence structures across different samples.
6. Use precise anatomical and imaging terminology appropriate for breast {modality}.
7. Output *only* valid JSON in the specified format. No extra text.

**Input (JSON):** {"Question": "...", "options": ["...", "...", "..."], "Answer": "..."}.

**Output (JSON):** {"Question": "...", "Answer": "..."}.

### Grounded Caption Generator

You are an **expert Radiologist Assistant** producing a **Grounded Caption** for breast {modality}. Strictly separate *visual attributes* (shape, margin, echo / signal / density, posterior or kinetic features) from *clinical context* (BI-RADS, risk, diagnosis, management).

**Step 1 – Reference phrase.** For each finding, write a noun phrase that names *only what is visible inside the bounding box*: include side, view (when applicable), abnormality type, and the relevant visual descriptors. Exclude BI-RADS, risk, diagnosis, and management.

**Step 2 – Caption skeleton.** Compose a professional radiology sentence with placeholder tokens <ref-object><bbox> immediately after each lesion phrase; inject the excluded clinical context (BI-RADS, risk, recommendation) *around* the placeholders, never inside the reference phrase.

**Step 3 – Strict alignment.** Emit ref[], bbox[] as flat lists in caption order; ref[i] corresponds to bbox[i] and to the i-th <ref-object><bbox> token. Boxes are 2D [xmin,ymin,xmax,ymax] for BUS / mammography / histopathology and 3D [xmin,ymin,zmin,xmax,ymax,zmax] for CT / MRI.

**Output JSON:** {caption, objects:{ref, bbox}}.

### Report Generator

You are a **board-certified breast radiologist**. Generate a formal breast {modality} report from the structured **patient\_data**, compliant with the **ACR BI-RADS Lexicon** for that modality and reflecting real-world dictation style.

**Hard safety rules.** (1) No fabrication of laterality, lesion size, clock-face position, or distance to nipple unless explicitly present in **patient\_data**; if missing, refer generically (“the lesion”). (2) No histological / molecular confirmation: replace “X is invasive carcinoma” with probabilistic phrasing (“findings are suspicious for X”). (3) Bounding boxes are internal; never mention coordinates, “bbox”, or any data-structure terminology in the prose.

**Sections (fixed order).** *Findings:* one passive-voice paragraph per lesion, ordered by descending BI-RADS, integrating only descriptors present in **patient\_data**. *Impression:* probabilistic synthesis. *Final Assessment:* the highest BI-RADS category among findings. *Management:* ACR-aligned recommendation (Cat 2 – routine; Cat 3 – short-interval follow-up; Cat 4–5 – biopsy / tissue diagnosis). Surgical / systemic therapy decisions are out of scope.

**Output:** four-section markdown with the section names above.

Figure 6: Generator prompts that produce the BreastStage instruction-following partitions from each structured record: the open-ended VQA generator (top), the grounded caption generator (middle), and the report generator (bottom). All three are instantiated per modality with {modality} ∈ {BUS, mammography, CT, MRI, histopathology}.

## MRI Structured Report Template

```

{
  "breast_level": {
    "FGT": {"left": "C", "right": "C"}, // "A","B","C","D"
    "BPE": {"left": "Moderate", "right": "Mild"}, // "Minimal","Mild","Moderate","Marked"
    "breast_symmetry": true,
    "post_surgical_changes": {
      "left": {"scar_tissue": false, "lumpectomy_changes": false, "mastectomy_changes": false},
      "right": {"scar_tissue": false, "lumpectomy_changes": false, "mastectomy_changes": false}
    },
    "associated_features": {
      "left": {
        "nipple_retraction": null, "nipple_invasion": null,
        "skin_retraction": null, "skin_thickening": null,
        "skin_invasion": {"direct_invasion": null, "inflammatory_cancer": null},
        "axillary_adenopathy": null, "pectoralis_muscle_invasion": null,
        "chest_wall_invasion": null, "architectural_distortion": null,
        "ductal_dilation": null // "mild","marked",null
      },
      "right": { ... same schema as "left" ... }
    },
    "lymph_nodes": {
      "axillary_left": null, "axillary_right": null, // "normal","abnormal"
      "internal_mammary_left": null, "internal_mammary_right": null
    },
    "BI-RADS": {"left": null, "right": null}, // "0","1","2","3","4A","4B","4C","5","6"
    "management": {"follow_up_recommended": null, "biopsy_recommended": null}
  },
  "mass": [
    {
      "id": 1, "side": "left",
      "location": {"quadrant": "UOQ", // "UIQ","UOQ","LIQ","LOQ","central",null
                  "depth": "middle"}, // "anterior","middle","posterior",null
      "size_mm": {"long": 12, "short": 9, "depth": 8},
      "signal_characteristics": {
        "T1W": "hyperintense", "T2W": "hyperintense", // "hyperintense","isointense","hypointense",
        null
      },
      "DWI_restriction": "present", // "present","absent",null
      "ADC": "hyperintense",
      "kinetics": {"initial": "Fast", // "Slow","Medium","Fast",null
                  "delayed": "Wash-out"} // "Persistent","Plateau","Wash-out",null
    },
    "morphology": {"shape": "oval", // "round","oval","lobular","irregular",null
                  "margin": "circumscribed"}, // "circumscribed","irregular","spiculated",null
    "internal_enhancement": {
      "pattern": "heterogeneous", // "homogeneous","heterogeneous","rim_enhancement",
      "dark_internal_septations": null
    },
    "associated_vascularity": {"adjacent_vessel_sign": null,
                              "increased_peritumoral_vascularity": null}
  }
],
  "non_mass_enhancement": [
    {
      "lesion_id": 1, "side": "left",
      "location": {"quadrant": null, "depth": null},
      "size_mm": {"long": null, "short": null, "depth": null},
      "distribution": "Segmental", // "Focal","Linear","Regional","Segmental","Multiple regions",
      "Diffuse",null
      "internal_enhancement_pattern": "Clustered_ring", // "Homogeneous","Heterogeneous","Clustered",
      "Clustered_ring",null
      "signal_characteristics": { ... same as "mass" ... }
    }
  ],
  "non_enhancing_lesion": [
    {
      "lesion_id": 1, "side": "left",
      "location": {"quadrant": null, "depth": null},
      "size_mm": {"long": null, "short": null, "depth": null},
      "signal_characteristics": { ... same as "mass" ... },
      "lesion_type_detail": "ductal_precontrast_high_signal_on_T1W"
      // "cyst","non_enhancing_mass","architectural_distortion",null
    }
  ]
}

```

Figure 7: Expert-designed BI-RADS-aligned structured template used for the MRI cohort. Every leaf field is either a numeric measurement, a Boolean flag, or one of a fixed enum list (shown as a comment). Reports are first translated to English by Qwen3-Max and then populated into this template under its enum and type constraints; specialists then derive a question per task template whose answer is drawn from one or more leaf fields. Schema-conformant population guarantees that the resulting QA pairs cannot introduce facts that the source report does not already record.

909 **Textual Description Generation.** The same structured records that feed VQA construction are  
910 also reused to synthesise the two narrative task formats, ground caption and report generation, with  
911 the modality-parameterised prompts in Figure 6. Both formats remain pinned to the structured record,  
912 so Qwen3-Max is responsible only for the surface form and never for the underlying clinical fact.

913 *Ground caption.* Ground captions are produced only for the four imaging modalities that carry  
914 bounding boxes (BUS, mammography, CT, MRI); histopathology, which has no boxes, is excluded.  
915 Each caption is generated by binding every visible lesion to its box and weaving the resulting visually  
916 grounded phrases into a clinical narrative; the surrounding context drawn from the structured record  
917 (BI-RADS, risk level, recommendation) is allowed to appear around the boxes but never inside the  
918 lesion phrase, so a reader can recover, from the caption alone, which descriptors apply to which  
919 spatial region. The modality-specific vocabularies differ in the obvious places (posterior acoustic  
920 features for BUS, ACR density for mammography, FGT, BPE and dynamic kinetic curves for MRI),  
921 but the visual-vs-context separation is identical across modalities.

922 *Report generation.* The report generator targets all five modalities, with one prompt variant per  
923 modality so that the output respects the relevant ACR BI-RADS lexicon. Every report is forced into  
924 the same four sections (Findings, Impression, Final Assessment, Management) so that downstream  
925 evaluation can score them section by section. Three constraints prevent the failure modes typical  
926 of pure-LLM report generators: spatial details that are not in the structured input (laterality, lesion  
927 size, clock-face position, distance to nipple) cannot appear, source pathology labels must remain  
928 probabilistic in the prose (“findings are suspicious for invasive carcinoma” rather than “this is invasive  
929 carcinoma”), and any reference to internal field names or coordinates is forbidden so the report  
930 reads as clinical dictation. Drafts then flow into the specialist-audit loop in Stage 4 below, where  
931 modality-specific prompts are revised and failing records regenerated whenever a per-task flag rate  
932 exceeds threshold.

933 **Data Splitting and Verification.** Records are split between training and BreastStage-Bench at  
934 patient-id granularity; because the source datasets cover disjoint patient populations, no cross-dataset  
935 patient deduplication is required. The split itself uses a single global stratified sampler over the  
936 composite (modality, task type, pathology label) key, so the resulting train and test partitions match  
937 on the underlying task and class distributions. BreastStage-Bench equals the test partition, drawn  
938 with task-key balanced sampling so that every task template in Table 6 is represented. Verification  
939 has two layers. Automated heuristics first remove records with malformed bounding-box coordinates,  
940 invalid option labels, instruction–answer conflicts, hallucinated clinical terms, and near-duplicate QA  
941 pairs (MinHash > 0.85). Three breast specialists then conduct an independent task-level audit on  
942 Bench: each specialist reviews 5 random samples per task and flags clinically inconsistent or stage-  
943 mismatched cases. Any task whose flag rate is non-negligible across specialists triggers re-running  
944 the relevant LLM stage with the failing prompts revised.

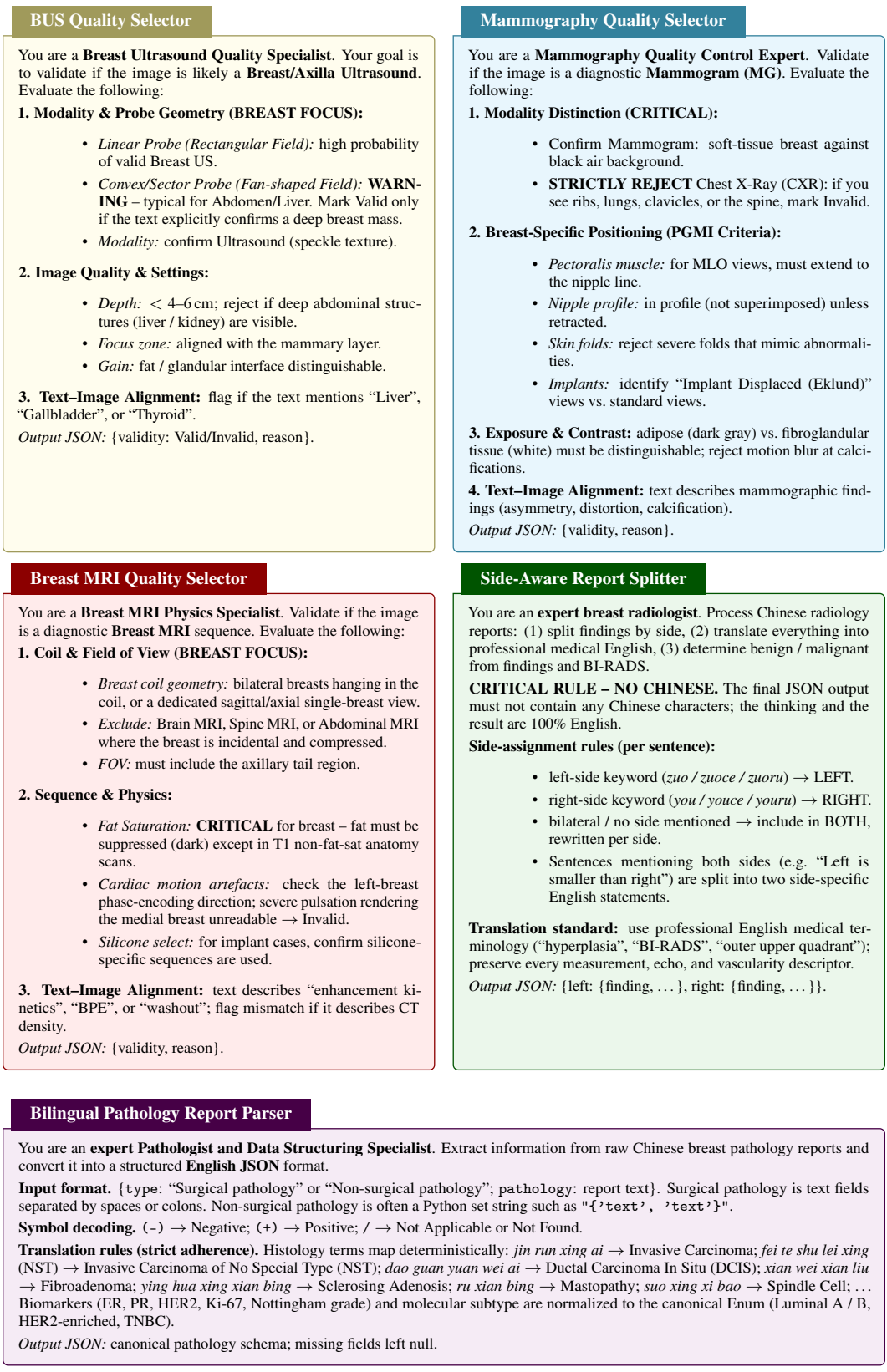


Figure 8: System prompts of the five BreastStage data-orchestration agents (modality-specific quality selectors plus the bilingual report splitter and pathology parser). Each agent runs Qwen3-Max with a strict JSON output schema, so its decision can be consumed by the next stage without manual parsing. Generator prompts that produce the actual instruction-following pairs are reproduced in Figure 6.

945 **Stage-aware system prompts.** BreastStage uses two complementary prompt families. The data-  
 946 orchestration prompts in Figure 8 and Figure 6 drive the curation pipeline: the five orchestration  
 947 agents (BUS / Mammography / Breast MRI Quality Selectors, Side-Aware Report Splitter, Bilingual  
 948 Pathology Report Parser) and the three generators (Open-VQA, Grounded Caption, Report). At  
 949 training and inference time, the resulting instruction-following pairs are then served with stage-  
 950 specific system prompts that switch BreastGPT between three clinical roles: screening radiologist,  
 951 diagnostic radiologist, and treatment-stage breast oncologist or pathologist; Figure 9 reproduces the  
 952 three personas and the deterministic template that combines a persona with a modality-aware task  
 953 instruction.

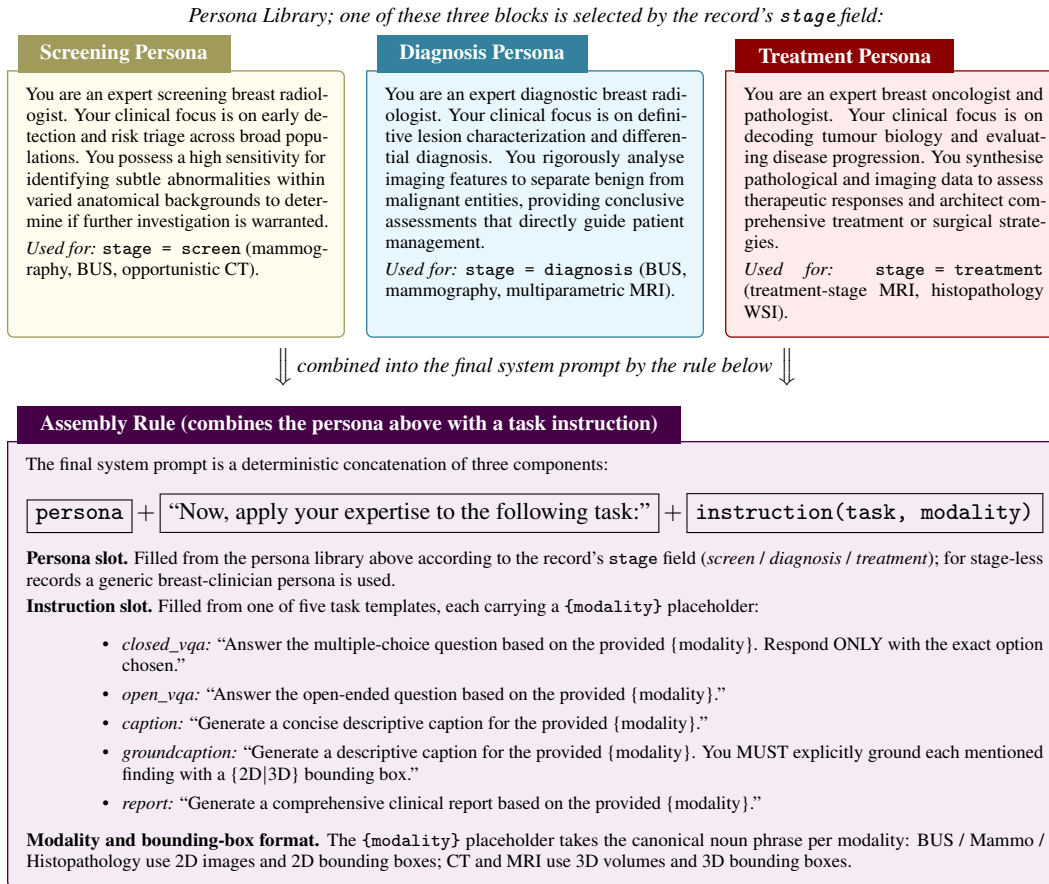


Figure 9: Stage-aware system-prompt library used by BreastGPT. The three personas (top) cover the screening / diagnosis / treatment workflow; the assembly rule (bottom) deterministically composes a persona, a transition phrase, and a modality-specific task instruction at training and inference time. The same persona is reused regardless of modality, so a single model can adopt stage-appropriate reasoning style without task-specific heads.

### 954 C.3 Task Taxonomy and Distribution

955 We summarise the BreastStage instruction corpus along two complementary axes. Table 6 groups  
 956 the 136 task templates into a 12-category taxonomy under the three clinical stages, and Figure 10  
 957 reports the joint distribution over (clinical stage, imaging modality, task category) that this taxonomy  
 958 induces.



Table 6: Twelve-category task taxonomy of BreastStage. Templates are grouped by clinical stage (Screening, Diagnosis, Treatment) and, within each stage, by the family of clinical concept being assessed. The *Samples* column reports the number of instruction-following pairs whose majority clinical stage falls in that category.

Clinical Stage	Task Category	Representative Task Templates	Samples
Screening (1,075,092)	Modality / Anatomy / Surgical History	Modality, View, Laterality, Orientation, ACR, FGT, BPE, Background tissue, breast_symmetry, post_surgical_changes (mastectomy / lumpectomy / scar)	442,860
	Lesion Presence & Morphology	abnormal_presence, abnormality_type, mass, non_mass_enhancement, non_enhancing_lesion, Shape, LesionEdge, LesionBoundary, EchoCharacteristics, LesionCalcificationFeatures, BloodFlow, lesion_depth	244,188
	Risk Assessment & Screening Decision	RiskLevel, ScreeningDecision, Density_Risk, follow_up_recommended	329,756
	Captioning / Report / Grounding	ground caption (BUS, CT), report generation (BUS), Report (MRI)	58,288
Diagnosis (680,409)	BI-RADS & Pathology Diagnosis	BI-RADS, birads, diagnosis, Pathology, main_type	520,969
	Lesion Characterization	mass.morphology (margin, shape), mass.signal (T1W, T2W, ADC, DWI, kinetics), internal_enhancement, calcificationType, boundaryDefinition, echoTexture, morphology, distribution	141,761
	Captioning / Report / Grounding (Mammo / MRI)	ground caption (Mammo), report generation (Mammography, MRI)	12,227
	Associated Findings & Invasion	lymph_nodes (axillary, internal_mammary), skin_thickening, nipple_retraction, architectural_distortion, ductal_dilation, pectoralis / chest_wall / nipple / skin invasion, vascularity, invasiveness	5,452
Treatment (100,567)	Prognosis & Outcome Prediction	prognosis, Treatment.Prognosis, Treatment.treatment.prognosis	33,743
	Surgical Planning & Urgency	surgical_plan, SurgicalPlan, treatment.urgency	29,797
	Systemic Therapy & Molecular Biomarkers	systemic_therapy, biomarkers (ER, PR, HER2, Ki67, molecular_subtype)	33,200
	WSI Pathology Subtyping & Captioning	Tumor, ER, PR, HER2, HER2 Expression, Molecular subtype, Histological grading, slide caption / VQA	3,827

## 959 D BreastGPT Technical Details

### 960 D.1 Base Model Configuration

961 BreastGPT is initialized from Qwen3-VL-8B-Instruct and inherits its instruction-following multi-  
 962 modal interface. We keep the language backbone and the standard radiology visual encoder largely  
 963 aligned with the released model and introduce additional components only where the breast cancer  
 964 workflow requires capabilities not naturally covered by the base model. The major architectural  
 965 extension is the GigaPixel branch for WSI processing together with the coverage-maximizing token  
 966 selector that maps heterogeneous image representations into a fixed downstream token budget. For CT,  
 967 BUS, mammography, and MRI inputs, the standard branch is used; for WSI inputs, modality-aware  
 968 routing activates the GigaPixel branch described below.

### 969 D.2 GigaPixel Encoder Implementation

970 The GigaPixel encoder is designed to preserve diagnostically relevant pathology patterns without  
 971 forcing the LLM to attend over the full patch sequence. A single WSI can contain tens or hundreds  
 972 of thousands of tiles, and many clinically important cues are sparse: mitotic figures, ductal structures,  
 973 invasive fronts, necrosis, lymphocytic infiltration, and tumor-stroma interfaces may occupy only a  
 974 small fraction of the slide. The WSI branch therefore separates representation learning into three  
 975 steps: local patch encoding, long-range contextualization, and query-aware coverage selection.

976 **Stage 1: Patch Extraction and Encoding.** WSI tiles are extracted at  $20\times$  magnification with  
 977 non-overlapping  $512\times 512$  pixel patches. For a typical breast WSI this yields tens of thousands of  
 978 patches; CONCH v1.5 patch features are pre-computed once per slide and stored, following the path  
 979 convention `20x_512px_0px_overlap/features_conch_v15/`. Each patch is encoded using the  
 980 frozen CONCH v1.5 foundation model [43], producing 512-dimensional embeddings:

$$\mathbf{v}_i = \text{CONCH}(\mathbf{x}_i), \quad i = 1, \dots, N, \quad \mathbf{v}_i \in \mathbb{R}^{512}.$$

981 **Stage 2: LongNet Contextualization.** The patch embeddings are processed by a 2-layer LongNet  
 982 encoder with dilated attention (matching the configuration described in Section 3 of the main paper).  
 983 The dilation rates follow an exponential schedule, allowing the model to capture both local cellular  
 984 details and global tissue architecture. The contextualized embeddings are:

$$\mathbf{h}_i = \text{LongNet}(\{\mathbf{v}_1, \dots, \mathbf{v}_N\})_i.$$

985 LongNet’s exponentially growing dilation rates give the encoder long-range coverage without the  
 986  $O(N^2)$  cost of full attention.

### 987 D.3 Coverage-Maximizing Token Selection

988 The coverage-maximizing token selector is the mechanism that makes BreastGPT practical as a single  
 989 model across ordinary radiology images and gigapixel pathology slides. Without compression, WSI  
 990 inputs would exceed the LLM context budget by several orders of magnitude; with naive pooling  
 991 or truncation, sparse but clinically important regions may be removed. We therefore use a coverage  
 992 objective that keeps the selected tokens both *clinically query-relevant* and *globally representative* of  
 993 the visual evidence.

994 Given a sequence of projected visual tokens and an associated clinical instruction, the goal of the  
 995 selector is to identify a compact subset  $S \subset \{1, \dots, N\}$  with  $|S| = k$  ( $k = 128$  by default) that  
 996 jointly preserves relevance to the clinical query and fidelity to the overall visual distribution. For  
 997 WSI inputs, the source sequence is the LongNet-contextualized patch representation  $\{h_1, \dots, h_N\}$ ;  
 998 for standard radiological inputs, the source sequence is the ViT token sequence. In both cases, the  
 999 selected tokens are passed to the LLM as the fixed-budget visual context. This gives BreastGPT a  
 1000 modality-invariant interface: each clinical image, regardless of its original resolution, is represented  
 1001 by the same number of visual tokens before language-model reasoning.

1002 **Similarity matrices.** Let  $\{v_1, \dots, v_N\}$  denote the projected vision tokens after the LLM projection  
 1003 layer, aligned with the text token space, and let  $\{v'_1, \dots, v'_N\}$  denote the pre-projection tokens that  
 1004 retain modality-native similarity structure. Let  $\{t_1, \dots, t_m\}$  denote the text tokens encoding the  
 1005 clinical instruction. We define text–vision (T-V) and vision–vision (V-V) similarity matrices as

$$M_{i,j}^{tv} = t_i^\top v_j, \quad M_{i,j}^{vv} = v_i'^\top v_j', \quad \|t_i\|_2 = \|v_j\|_2 = \|v_j'\|_2 = 1. \quad (2)$$

1006 The similarities are then softmax-calibrated with temperatures  $\tau_t$  and  $\tau_v$ :

$$\widetilde{M}_{i,j}^{tv} = \frac{\exp(M_{i,j}^{tv}/\tau_t)}{\sum_{j'} \exp(M_{i,j'}^{tv}/\tau_t)}, \quad \widetilde{M}_{i,j}^{vv} = \frac{\exp(M_{i,j}^{vv}/\tau_v)}{\sum_{j'} \exp(M_{i,j'}^{vv}/\tau_v)}. \quad (3)$$

1007 **Dual coverage objective.** We select  $S$  by maximizing the joint coverage score

$$f(S) = \underbrace{\frac{1}{m} \sum_{i=1}^m \max_{s \in S} \widetilde{M}_{i,s}^{tv}}_{\text{T-V coverage}} + \alpha \cdot \underbrace{\frac{1}{N} \sum_{i=1}^N \max_{s \in S} \widetilde{M}_{i,s}^{vv}}_{\text{V-V coverage}}, \quad (4)$$

1008 where  $\alpha$  balances query relevance against global visual representativeness. The T-V term rewards  
 1009 tokens that explain the clinical instruction, while the V-V term forces the selected subset to cover  
 1010 the full visual distribution so that diagnostically relevant but less immediately salient regions are not  
 1011 discarded.

1012 **Why both terms are needed.** The two coverage terms play complementary roles. T-V coverage  
1013 favors tokens aligned with the prompt, such as lesion regions for BI-RADS questions or tumor-  
1014 cell regions for molecular subtype questions. However, relying on T-V coverage alone can over-  
1015 concentrate the selected subset around a small number of highly salient regions. V-V coverage  
1016 counteracts this collapse by selecting tokens that represent the broader slide or image distribution.  
1017 This is especially important for pathology, where grade, subtype, and treatment-response cues can  
1018 depend on tissue architecture, tumor-stroma interaction, and heterogeneity across multiple regions.  
1019 Conversely, V-V coverage alone may preserve visually diverse background regions that are irrelevant  
1020 to the current clinical task. The combined objective gives BreastGPT a principled way to balance  
1021 prompt alignment and visual diversity.

1022 **Greedy optimization.** Each component of  $f$  is a facility-location function over the token set and is  
1023 therefore monotone submodular; the non-negative linear combination remains monotone submodular.  
1024 A standard greedy algorithm that iteratively adds the token with the largest marginal gain achieves  
1025 the usual  $(1 - 1/e)$  approximation to the NP-hard optimum in  $O(kN)$  time. End-to-end measured  
1026 selector latency at  $k = 128$  on a representative histopathology WSI is reported in Table 15.

1027 **Implementation details.** The selector is applied after modality-specific visual encoding and before  
1028 LLM injection. For standard radiology images, token selection operates on the ViT token sequence.  
1029 For WSI inputs, token selection operates on LongNet-contextualized patch embeddings. The same  
1030 value of  $k$  is used for both branches, which ensures that radiology and pathology inputs have  
1031 comparable downstream inference costs. During training, the selector is run on-the-fly so that the  
1032 language model learns to reason over the same type of compact visual context used at inference time.

1033 **Relation to MMTok.** The original MMTok [18] is a training-free post-hoc token pruning method  
1034 applied to standard vision-language models after the raw ViT encoder. Our adaptation differs in  
1035 two respects. First, for WSIs, the source tokens are LongNet-contextualized pathology embeddings  
1036 rather than raw ViT patch tokens, so the V-V similarity captures global pathological structure rather  
1037 than pixel-level redundancy. Second, the selector is applied consistently during both training and  
1038 inference, so the LLM always observes coverage-selected visual representations. We fix  $\tau_t = 0.02$ ,  
1039  $\tau_v = 0.2$ , and  $\alpha = 0.5$  throughout all experiments.

#### 1040 **D.4 Training Configuration**

1041 BreastGPT is trained in two stages. In *Stage 1*, we warm up the visual front end on the closed-VQA  
1042 subset of BreastStage: the LLM backbone is frozen, while the ViT, the aligner (which includes the  
1043 WSI projection and LongNet module), and the resolution gate are trainable; the CONCH v1.5 encoder  
1044 is always frozen. This step aligns the visual representations, both the ViT branch for radiology and  
1045 the CONCH+LongNet branch for WSIs, with the multimodal token space before any gradient flows  
1046 into the language model. In *Stage 2*, we unfreeze the LLM and continue end-to-end fine-tuning on all  
1047 four task formats across every modality; the CONCH v1.5 encoder remains frozen. The coverage  
1048 token selector is training-free and applied identically at training and inference. We use AdamW with  
1049 a cosine learning rate schedule and DeepSpeed ZeRO-2, cap the per-image visual budget at 1024  
1050 tokens, cap CT and MRI volumes at  $384 \times 384 \times 48$ , and fix the selector budget at  $k = 128$  for  
1051 both branches so that radiological and pathological inputs incur the same downstream inference cost.  
1052 Table 7 summarises the two-stage training configuration used for BreastGPT, taken directly from the  
1053 launch scripts.

Table 7: Training configuration for BreastGPT. Both stages use the `ms-swift sft` entry point with FlashAttention, gradient checkpointing, and DeepSpeed ZeRO-2.

Setting	Stage 1: Visual front-end warm-up	Stage 2: End-to-End SFT
Epochs	1	2
Trainable modules	ViT, aligner (incl. WSI projection / LongNet)	LLM, ViT, aligner
Frozen modules	LLM backbone, CONCH encoder	CONCH encoder
Datasets	Closed-VQA subsets (BUS, CT, Mammo, MRI) and the histopathology caption / closed / open VQA tasks	All four task formats (closed, open, grounding, report) on every modality (16 splits)
Initialization	Released Qwen3-VL-8B-Instruct weights	Stage-1 checkpoint
Optimizer	AdamW (default $\beta$ , weight decay 0.1)	AdamW (default $\beta$ , weight decay 0.1)
Learning rate	$1 \times 10^{-5}$ , cosine schedule, 0.03 warm-up ratio	$1 \times 10^{-5}$ , cosine schedule, 0.03 warm-up ratio
Gradient clipping	1.0	1.0
Precision	bfloat16	bfloat16
Image / volume budget	IMAGE_MAX_TOKEN_NUM=1024; CT and MRI volumes capped at $384 \times 384 \times 48$	same
Max sequence length	default	10,000 tokens
Selector budget $k$	128	128
Distributed	DeepSpeed ZeRO-2; per-device batch and accumulation set per launch	DeepSpeed ZeRO-2; per-device batch and accumulation set per launch

## 1054 D.5 System Prompt Templates

1055 The full stage-aware system-prompt library used by BreastGPT, comprising three personas (Screening,  
 1056 Diagnosis, Treatment, plus a generic clinician fallback) and the deterministic rule that composes  
 1057 a persona with a modality-specific task instruction, is reproduced verbatim from the codebase in  
 1058 Figure 9. The same library is used both during training and at inference time; the persona slot is  
 1059 selected by the record’s `stage` field, and the instruction slot is filled by task and modality as described  
 1060 in the assembly rule.

## 1061 E BreastStage-Bench Evaluation

### 1062 E.1 Benchmark Construction

1063 BreastStage-Bench is constructed to evaluate whether a model can preserve stage-aware clinical  
 1064 reasoning across the full breast cancer workflow rather than merely solve isolated image-classification  
 1065 tasks. We use patient-level stratified sampling and keep all task instances derived from the same  
 1066 patient within the same split. This prevents leakage through repeated images, repeated reports, or  
 1067 semantically equivalent QA pairs.

1068 The benchmark construction follows three principles:

- 1069 1. **No data leakage:** Strict patient-level separation between training and test sets
- 1070 2. **Task diversity:** All 136 task templates are represented with balanced class distributions
- 1071 3. **Clinical realism:** Test cases reflect real-world prevalence and complexity

1072 BreastStage-Bench contains 12,182 evaluation records in total, distributed across the four task families  
 1073 and five modalities as follows:

- 1074 • Closed-ended VQA (5,369 questions): BUS 987, CT 825, Mammography 1,828, MRI 908,  
 1075 Histopathology 821.
- 1076 • Open-ended VQA (2,833 questions): BUS 987, CT 825, MRI 908, Histopathology 113.
- 1077 • Ground caption (2,910 cases): BUS 1,000, CT 510, Mammography 1,000, MRI 400.
- 1078 • Caption (70 cases): Histopathology slide-level captioning.
- 1079 • Report generation (1,000 cases): multiparametric MRI structured reports.

1080 **E.2 Evaluation Metrics**

1081 **Closed-ended VQA.** We adopt accuracy as the primary metric. For multiple-choice questions,  
1082 we extract model responses using robust regular expressions matching option patterns (A/B/C/D),  
1083 answer text, and common variants such as “Option A” or “the correct answer is A”. If no valid option  
1084 is detected in the first 100 tokens of the response, the prediction is marked invalid and counted as  
1085 incorrect. This avoids rewarding models for verbose but non-committal answers.

1086 **Open-ended VQA.** Unlike LLM-as-judge protocols commonly used in recent multimodal eval-  
1087 uation toolkits [70, 71, 19], BreastStage-Bench reports automatic semantic and lexical scores for  
1088 open-ended responses, and uses breast-expert sampling validation to audit the clinical correctness of  
1089 references and representative outputs. For normalized open-ended scores, each evaluated response  
1090 receives a score  $s_i \in [0, 1]$  according to the task-specific scoring rubric, and the total score is:

$$S = \frac{\sum_{i=1}^N s_i}{N} \times 100\% \quad (5)$$

1091 The few-shot prompt includes 9 in-context examples covering fully correct (1.0), entirely incorrect  
1092 (0.0), and various “partially correct” responses (0.25, 0.5, 0.75). The evaluator is instructed to  
1093 prioritize clinical correctness over surface-form similarity and to penalize unsupported diagnosis,  
1094 missing contraindications, or recommendations inconsistent with the specified clinical stage.

1095 **Grounding Tasks.** For lesion localization, we compute Intersection over Union (IoU) between  
1096 predicted and ground-truth bounding boxes:

$$\text{IoU} = \frac{\text{Area}(\text{pred} \cap \text{gt})}{\text{Area}(\text{pred} \cup \text{gt})} \quad (6)$$

1097 A prediction is considered correct if  $\text{IoU} \geq 0.5$ . We report mean IoU and  $\text{accuracy}@0.5$  across all  
1098 grounding tasks. Predictions with malformed coordinate syntax are counted as invalid and assigned  
1099 IoU 0.

1100 **Report Generation.** We evaluate generated reports using:

- 1101 • **BLEU-4:** Measures n-gram overlap with reference reports
- 1102 • **BERT Score:** Captures semantic similarity using contextual embeddings  
1103 (microsoft/BiomedNLP-PubMedBERT-base)
- 1104 • **ROUGE-1:** Assesses recall of unigrams

1105 The weighted composite score is:  $\text{Wtd} = 0.5 \times \text{BERT} + 0.25 \times \text{BLEU} + 0.25 \times \text{ROUGE-1}$ .

1106 **E.3 Breast Expert Validation**

1107 BreastStage and BreastStage-Bench are constructed with a structured expert sampling audit rather  
1108 than relying only on automatic filters or LLM-generated text. Three board-certified breast specialists  
1109 (two breast surgeons and one breast radiologist, each with at least five years of post-residency  
1110 experience) independently review a stratified subset of BreastStage-Bench, scoring each item along  
1111 three pre-registered dimensions that mirror the claim made in the main text: *task validity*, *answer*  
1112 *correctness*, and *clinical consistency*. The audit covers all five modalities (BUS, mammography, CT,  
1113 MRI, histopathology) and all task families (closed VQA, open VQA, ground caption, report).

1114 **Sampling protocol.** Sampling is stratified by (modality, task template) so that every one of the  
1115 136 task templates contributes evenly: each specialist independently audits 5 random samples per  
1116 template, yielding 680 per-specialist reviews and a pooled review pool of 2,040 (*record*, *specialist*)  
1117 judgments. Specialists are blind to which LLM stage produced each record. Each item is judged on  
1118 the following three binary dimensions:

- 1119 • **Task validity** — the question is well-posed for the indicated clinical stage and the requested  
1120 evidence is identifiable in the image (or in the source report, for report-only items).

1121 • **Answer correctness** — the gold answer is logically consistent with the question and contains  
 1122 no claims unsupported by the image or source report (no hallucinated lesions, biomarkers,  
 1123 or risk categories).

1124 • **Clinical consistency** — the assigned clinical stage (screening / diagnosis / treatment) and  
 1125 the modality-specific terminology (BI-RADS category, ACR density, molecular subtype,  
 1126 etc.) match the intended clinical role and current breast-oncology practice.

1127 **Disagreement resolution.** Items flagged by at least one but not all specialists are routed to a  
 1128 *consensus review*: the three specialists meet, the contested record is re-examined together, and a  
 1129 single binary decision is recorded by majority vote (with the breast radiologist breaking ties on  
 1130 imaging-specific calls and the senior breast surgeon breaking ties on management calls). Templates  
 1131 whose post-consensus failure rate exceeds 10 % are routed back to the responsible LLM prompt for  
 1132 revision, after which the affected records are regenerated and re-audited.

1133 We report per-dimension approval rate and inter-rater agreement (Fleiss’s  $\kappa$ ) in Table 8. Across all  
 1134 three dimensions the post-revision approval rate exceeds 95 %, and Fleiss’s  $\kappa$  stays in the substantial-  
 1135 agreement range ( $\kappa \in [0.74, 0.86]$ ), indicating that the audit signal is reliable rather than dominated by  
 1136 a single specialist’s idiosyncrasies.

Table 8: Expert-validation results on BreastStage-Bench. *Pre-rev.* = approval rate on the first audit pass; *Post-rev.*  
 = rate after the consensus-review and prompt-revision loop.  $\kappa$  is Fleiss’s  $\kappa$  across the three specialists on the  
 binary approval decision. *Reviews* counts (*record, specialist*) judgments; *Answer correctness* is computed  
 on items that carry a gold answer or free-text generation, hence the smaller pool. *Re-gen.* = number of task  
 templates routed to prompt revision.

Dimension	Reviews	Pre-rev.	Post-rev.	$\kappa$	Re-gen.
Task validity	2,040	91.4%	98.2%	0.78	14
Answer correctness	1,620	90.6%	97.4%	0.79	11
Clinical consistency	2,040	96.6%	99.8%	0.86	2
<b>Overall (any-dim. fail)</b>	2,040	<b>86.1%</b>	<b>96.3%</b>	0.74	22

1137 **Per-stage breakdown.** Post-revision approval is stable across clinical stages (screening 96.5%,  
 1138 diagnosis 96.1%, treatment 95.8%), suggesting the LLM-generated content is not systematically  
 1139 biased toward any single stage. The treatment stage has the lowest pre-revision *answer correctness*  
 1140 score, 86.1%, because biomarker and prognostic claims are easier for the LLM to over-generate than  
 1141 morphological descriptions; tightening the treatment-stage prompts (constraining biomarker mentions  
 1142 to those present in the source pathology report) raised this dimension to 96.4% post-revision. The  
 1143 per-stage record counts of the final pipeline export are summarised in Table 6 (1,068,187 screening,  
 1144 680,078 diagnosis, 101,585 treatment-stage pairs).

1145 **Limitations of this audit.** Three specialists, while sufficient for  $\kappa$  to be meaningful, is below  
 1146 the panel sizes used in formal radiology guideline development; we therefore treat this audit as a  
 1147 quality-assurance pass rather than a clinical gold standard. The full auditing log – per-record flag,  
 1148 dimension scored, and (where applicable) the resulting prompt revision – is released alongside the  
 1149 dataset to allow downstream re-auditing.

#### 1150 E.4 Baseline Model Configuration

1151 All baseline models are evaluated zero-shot with the same instruction-formatted prompts used for  
 1152 BreastGPT. For proprietary models (GPT-5.4, Claude-opus-4-6, Claude-sonnet-4-6, Gemini-3.1-  
 1153 Flash, Gemini-3.1-Pro), we use the official APIs with temperature 0 and max\_tokens 512. For  
 1154 open-source models we use Hugging Face Transformers with greedy decoding. Input images are  
 1155 preprocessed according to each model’s recommended processor; for multi-image or multi-sequence  
 1156 cases all models receive the same ordered visual inputs and the same textual task context.

1157 **F Additional Results and Ablations**

1158 **F.1 Per-Modality Performance Breakdown**

1159 Table 9 provides a per-modality performance breakdown for BreastGPT and the strongest competing  
 1160 baselines. This view complements the stage-level tables in the main paper by isolating whether  
 1161 improvements are driven by a single easy modality or are distributed across heterogeneous imaging  
 1162 sources.

Table 9: Per-modality performance breakdown on BreastStage-Bench closed-ended VQA.

Model	BUS	CT	Mammo	MRI	Histo	Avg
BreastGPT (cluster)	<b>86.81</b>	<u>77.21</u>	<b>75.00</b>	<b>82.86</b>	<b>71.38</b>	<b>78.65</b>
GPT-5.4	64.89	<b>78.55</b>	<u>68.51</u>	41.43	32.28	<u>57.13</u>
Gemini-3.1-Pro	<u>68.09</u>	73.33	50.21	47.14	46.53	<u>57.06</u>
Lingshu	58.94	<b>78.55</b>	39.89	<u>54.29</u>	<u>51.52</u>	56.64

1163 **F.2 Extended Baseline Comparisons**

1164 Tables 10 and 11 report the same VQA and Caption/Report metrics as the main paper’s Tables 1  
 1165 and 2, but extended to eight additional baseline models (Grok-4.1-Fast, Gemma-4, GLM-4.6V-Flash,  
 1166 LLaVA-OneVision-1.5, HealthGPT, Hulu-Med, MedDr, RadFM) covering both proprietary and  
 1167 recently released open / medical VLMs. None of these baselines beat BreastGPT on any cell, so for  
 1168 compactness the main paper limits its comparison to the original cohort; this section confirms that the  
 1169 conclusion holds against a broader set.

Table 10: VQA performance (%) on BreastStage-Bench for eight additional baselines (extends Table 1). Same column structure as the main table; refer to it for column legend.

Model	#P	Closed-ended VQA (Accuracy, %)										Open-ended VQA (normalized Score, %)							
		Screening					Diagnosis			Treatment		Screening			Diagnosis			Treatment	
		BUS	CT	Mam	MRI	Avg	BUS	Mam	MRI	His	Avg	BUS	CT	MRI	BUS	MRI	MRI	His	Avg
<i>Proprietary Models</i>																			
Grok-4.1-Fast	-	57.23	76.61	30.74	40.95	<b>49.47</b>	<b>22.82</b>	44.23	31.75	36.42	43.36	44.76	43.43	43.95	43.06	43.94	41.74	40.50	43.05
<i>Open-Source Models</i>																			
Gemma-4	8B	<b>62.98</b>	<b>78.18</b>	<u>44.57</u>	38.57	47.87	18.17	48.43	39.68	40.68	46.57	43.29	44.20	43.46	42.47	44.49	42.40	40.12	42.92
GLM-4.6V-Flash	9B	<u>58.94</u>	73.09	37.34	41.90	45.48	12.16	48.60	<b>53.97</b>	<u>49.70</u>	<u>46.80</u>	42.69	43.70	42.78	42.05	43.54	42.47	40.28	42.50
LLaVA-OneVision-1.5	8B	58.09	68.85	30.32	40.95	41.76	12.31	<b>50.17</b>	<b>53.97</b>	<b>50.18</b>	45.18	-	56.06	59.84	-	59.19	58.12	45.25	55.69
<i>Medical-Specific Models</i>																			
HealthGPT	-	47.02	24.48	41.38	<u>45.71</u>	34.04	21.17	43.01	42.86	20.22	35.54	54.66	50.17	53.08	51.92	59.38	55.27	46.41	52.98
Hulu-Med	-	<b>62.98</b>	<u>77.09</u>	<b>48.19</b>	31.90	<u>49.20</u>	15.62	41.43	<u>50.79</u>	<b>50.18</b>	<b>47.49</b>	<u>58.98</u>	60.84	<b>75.46</b>	<u>55.95</u>	<u>72.28</u>	59.48	50.18	<u>61.88</u>
MedDr	-	48.51	47.39	25.21	<b>52.38</b>	35.37	<u>22.52</u>	<u>49.83</u>	43.65	39.46	40.48	<b>61.46</b>	<b>64.52</b>	<u>67.68</u>	<b>62.83</b>	<b>73.03</b>	<u>64.86</u>	<b>50.46</b>	<b>63.55</b>
RadFM	-	26.81	39.39	16.38	36.67	23.94	5.86	17.13	11.90	4.63	20.30	50.51	<u>61.07</u>	67.63	50.80	69.15	<b>66.76</b>	<u>50.44</u>	59.48

Table 11: Caption and Report Generation Performance (%) for eight additional baselines (extends Table 2). Same column structure as the main table.

Model	#P	Caption								Report
		BUS		CT	Mammo		Histo	MRI	MRI	
		IoU	Wtd	Wtd	IoU	Wtd	Wtd	Wtd	Wtd	
<i>Proprietary Models</i>										
Grok-4.1-Fast	-	13.40	47.30	45.12	2.52	48.09	45.15	<b>49.02</b>	<u>50.47</u>	
<i>Open-Source Models</i>										
Gemma-4	8B	<b>37.82</b>	46.34	44.45	<b>8.01</b>	46.33	44.04	47.87	49.16	
GLM-4.6V-Flash	9B	<u>28.03</u>	43.01	43.43	6.01	43.15	43.53	44.26	48.21	
LLaVA-OneVision-1.5	8B	14.99	<u>47.68</u>	<u>49.04</u>	4.53	<u>48.41</u>	<u>46.91</u>	<u>48.62</u>	<b>50.62</b>	
<i>Medical-Specific Models</i>										
HealthGPT	-	22.56	43.11	40.05	<u>7.69</u>	44.40	46.63	46.93	48.18	
Hulu-Med	-	5.75	<b>48.67</b>	48.40	0.41	<b>48.74</b>	<b>47.25</b>	48.49	50.33	
MedDr	-	0.17	<u>47.68</u>	<b>50.25</b>	0.38	48.34	44.73	47.04	46.60	
RadFM	-	-	41.90	38.13	-	43.42	41.23	46.53	43.27	

1170 **F.3 Open-ended VQA Raw Metric Breakdown**

1171 Table 12 reports the three constituent metrics (BERTScore F1, BLEU, ROUGE-1) underlying the  
 1172 weighted score reported in Table 1. Closed-ended VQA is omitted because its only metric is accuracy.

Table 12: Open-ended VQA raw metrics per cell. Each entry is BERTScore F1 / BLEU / ROUGE-1 (%). BreastGPT in cyan.

Model	#P	Screening			Diagnosis		Treatment	
		BUS	CT	MRI	BUS	MRI	MRI	His
<i>Proprietary Models</i>								
GPT-5.4	-	88.88/1.62/34.33	86.50/1.65/22.73	91.82/11.65/44.10	87.25/1.45/25.71	90.82/14.09/41.05	90.57/10.10/41.04	83.67/0.57/11.02
Claude-opus-4-6	-	82.55/0.45/7.73	81.86/0.66/6.70	82.75/1.25/8.24	82.05/0.30/6.81	82.95/1.87/9.90	81.76/0.73/7.40	79.96/0.17/3.27
Claude-sonnet-4-6	-	82.54/0.46/8.97	81.71/0.88/8.16	82.71/1.26/10.56	81.85/0.34/7.94	82.83/1.82/12.46	80.90/0.51/8.15	79.90/0.14/4.32
Gemini-3.1-Flash	-	86.50/1.51/20.51	86.04/2.68/18.08	85.60/2.51/15.81	84.66/0.68/13.19	85.36/3.26/16.15	83.07/1.12/9.34	82.27/0.34/6.37
Grok-4.1-Fast	-	83.70/0.56/11.10	82.32/0.82/8.25	82.83/0.90/9.24	82.21/0.30/7.51	82.33/1.29/9.83	80.84/0.34/4.93	79.63/0.13/2.60
<i>Open-Source Models</i>								
Qwen2.5-VL	3B	87.99/2.12/19.97	86.85/2.98/16.89	87.59/3.41/18.86	86.69/1.31/14.75	87.48/5.50/21.70	85.49/2.08/15.08	84.46/0.49/9.80
Qwen2.5-VL	7B	86.11/1.13/11.45	85.82/2.08/13.72	87.36/2.70/20.45	85.06/0.48/9.36	86.98/4.61/20.11	83.92/1.32/9.62	83.34/0.32/8.94
Qwen3-VL	4B	84.69/0.98/10.42	89.65/5.08/37.59	86.29/2.10/16.16	84.23/0.91/10.36	87.78/4.95/25.31	85.61/1.86/22.86	80.65/0.22/3.52
Qwen3-VL	8B	84.40/1.00/9.98	86.40/2.64/18.68	83.89/1.73/9.90	83.76/0.74/9.11	84.43/2.78/13.30	82.52/1.14/10.41	80.51/0.23/3.60
MiMo-VL	7B	82.27/0.52/4.90	88.14/4.29/43.35	89.97/19.76/55.10	81.32/0.37/3.75	90.36/22.57/59.33	89.91/16.07/56.42	78.85/0.14/1.59
InternVL3.5	8B	87.86/4.39/22.84	89.02/14.39/39.77	90.19/7.99/37.85	86.56/2.69/19.37	90.65/12.68/41.59	89.44/7.36/37.54	84.86/0.50/15.80
Gemma-4	8B	82.79/0.52/7.05	83.20/1.12/9.26	82.71/1.05/7.39	81.86/0.32/5.83	83.17/1.91/9.71	81.59/0.68/5.75	79.20/0.05/2.04
GLM-4.6V-Flash	9B	82.16/0.71/5.73	83.03/1.35/7.38	82.38/1.16/5.20	81.52/0.57/4.58	82.90/1.69/6.65	82.17/0.89/4.65	79.42/0.19/2.10
LLaVA-OneVision-1.5	8B	-	90.14/6.57/37.39	91.82/11.65/44.09	-	90.82/14.09/41.03	90.57/10.10/41.24	84.88/0.39/10.84
<i>Medical-Specific Models</i>								
Lingshu	7B	89.47/3.66/27.78	87.79/4.01/19.76	89.29/5.04/30.02	87.59/1.80/18.81	88.65/8.22/29.52	86.70/3.17/21.84	83.57/0.30/7.32
HuatuogPT-V	7B	88.99/3.60/25.86	88.43/5.17/23.85	90.20/6.84/35.48	87.66/2.25/19.23	89.36/10.04/31.58	88.17/5.16/30.22	85.07/0.88/12.08
HealthGPT	-	89.69/5.73/33.52	86.17/7.50/20.84	87.02/10.10/28.18	88.90/3.75/26.11	89.91/16.43/41.27	87.84/10.10/35.30	84.83/0.71/15.28
Hulu-Med	-	91.14/9.89/43.73	90.99/14.81/46.56	95.29/40.64/70.61	90.02/6.34/37.42	93.60/37.93/64.00	90.64/10.31/46.35	87.35/3.33/22.70
MedDr	-	91.22/20.02/43.39	90.46/31.04/46.13	92.69/26.61/58.74	91.85/19.38/48.23	93.70/40.07/64.63	91.82/21.43/54.38	87.12/3.25/24.35
RadFM	-	87.50/3.86/23.19	89.66/24.94/40.04	93.08/25.80/58.57	86.53/9.71/20.44	92.94/30.76/59.96	91.86/26.74/56.57	87.00/3.21/24.55
Qwen3-VL (SFT)	8B	98.85/87.18/92.83	99.04/89.22/94.24	99.11/85.25/94.82	98.26/81.90/88.38	99.03/87.90/95.17	97.65/72.02/86.10	90.73/11.46/38.93
BreastGPT (cluster)	8B	99.15/90.67/94.91	99.07/88.56/94.45	99.27/87.53/95.87	98.63/85.08/90.60	99.12/88.94/95.68	98.09/75.47/88.08	92.69/18.88/50.96

1173 **F.4 3D Grounding IoU on CT and MRI**

1174 None of the 21 baseline models in our evaluation cohort can localise lesions in 3D volumetric  
 1175 modalities (CT, MRI). Their bbox outputs are uniformly 2-D coordinate quadruples, which are  
 1176 dimensionally incompatible with the 6-element 3-D ground-truth boxes; consequently, on every  
 1177 GT-positive 3-D sample they fail to overlap the ground truth at all. The grounding-IoU numbers  
 1178 a few baselines report on these modalities are entirely a by-product of true-negative credit on the  
 1179 no-abnormality samples, not of any spatial localisation. We therefore report 3-D grounding IoU only  
 1180 for BreastGPT and its controlled SFT counterpart in Table 13.

Table 13: Grounding IoU (%) on the 3D modalities (CT, MRI), restricted to the two models that emit volumetric (6-D) bounding boxes. IoU credits true negatives as 1.0 and penalises false positives as 0. All 21 other baselines produce 2-D outputs that are dimensionally incompatible with the 3-D ground truth and are omitted here.

Model	#P	CT IoU	MRI IoU
Qwen3-VL (SFT)	8B	3.89	11.98
BreastGPT (cluster)	8B	5.12	33.49

1181 **F.5 BreastGPT Grounding Recognition Breakdown**

1182 Table 14 decomposes the IoU column of Table 2 for BreastGPT (cluster) into its true-negative and  
 1183 false-positive components: TN counts items with no GT bbox where the model also abstains (correct  
 1184 “no abnormality”), FP counts items where the model hallucinates a bbox on a normal case, and  
 1185 *Pred-with-GT* counts items where BreastGPT attempted a bbox on a GT-positive case.

Table 14: BreastGPT (cluster) grounding-IoU decomposition. IoU credits true negatives (TN) as 1.0 and penalises false positives (FP) as 0. *Pred-with-GT* counts items where BreastGPT attempted a bbox on a GT-positive case;  $N$  is the total number of items.

Modality	IoU (%)	TN	FP	Pred-with-GT	N
BUS	79.59	0	0	1000	1000
CT	5.12	33	70	333	510
Mammo	23.14	0	2	998	1000
MRI	33.49	108	60	202	400

## 1186 F.6 Inference Efficiency Analysis

1187 We profile the actual cost of running BreastGPT on the histopathology WSI subset of BreastStage-  
 1188 Bench at different token budgets  $k$ , including a no-limit baseline that disables the selector and  
 1189 feeds every patch token of the slide into the LLM. Protocol: bf16, batch size 1, 5 warm-up +  
 1190 30 timed forward passes, on a single GPU. Latency is measured with `torch.cuda.Event` and  
 1191 `torch.cuda.synchronize()` so that timing reflects the actual GPU execution rather than asyn-  
 1192 chronous launches. Selector latency is the wall-clock time inside the greedy coverage routine; prefill  
 1193 latency is the first-token forward through the LLM after selection; total latency is selector + prefill.

Table 15: Inference efficiency of BreastGPT (cluster) on a representative histopathology WSI (5,987 patches before selection). “LLM in. tok.” is the total LLM input length after selection (visual  $k$  tokens plus prompt). The “no limit” row sets  $k$  above the patch count, effectively disabling the selector and feeding every patch token to the LLM.

$k$	LLM in. tok.	Prefill (ms)	Selector (ms)	Decode (ms/tok)	Total (ms)	Peak Mem (GB)
1	145	57.9	0.4	42.6	58.2	16.96
8	152	67.3	0.9	45.9	68.2	16.96
16	160	78.1	1.5	50.4	79.6	16.96
32	176	89.4	2.6	48.0	91.9	16.96
64	208	123.1	4.8	54.5	127.9	16.96
128	272	191.4	9.3	68.2	200.6	16.97
256	400	328.4	18.2	95.5	346.6	16.97
512	656	605.7	36.1	150.7	641.8	16.98
no limit	5,144	5,512.6	—	1,129.5	6,545.7	17.95

1194 Three observations follow from these numbers; the per-stage latency breakdown is plotted in main-  
 1195 paper Fig. 4. (i) *Prefill and decode scale roughly linearly with  $k$ .* From  $k=1$  to  $k=512$  the LLM input  
 1196 length grows from 145 to 656 tokens, and prefill grows in step from 57.9 to 605.7 ms; decode time  
 1197 per output token grows from 42.6 to 150.7 ms, since the KV-cache that each new token must attend  
 1198 to also lengthens with  $k$ . The chosen  $k=128$  costs 191 ms of prefill,  $3.2\times$  less than  $k=512$ , at  $>99\%$   
 1199 of its task quality. (ii) *Without selection, the LLM forward dominates the entire pipeline.* The no-limit  
 1200 baseline raises the LLM input to 5,144 tokens; prefill jumps to 5.5 s and per-token decode to 1.1 s,  
 1201 for a total of 6.5 s per question—about  $33\times$  slower than  $k=128$  on the same slide. The selector itself  
 1202 only takes 1.0 s in this regime, so the cost saved is overwhelmingly LLM-side, exactly as expected for  
 1203 a method that compresses the LLM input. (iii) *Memory grows slowly under selection and jumps only*  
 1204 *at no limit.* Peak GPU memory rises from 16.96 GB at  $k=1$  to 16.98 GB at  $k=512$  and then to 17.95  
 1205 GB without selection, a  $\sim 1$  GB jump that scales with the LLM input length. BreastGPT therefore  
 1206 runs on a single 24 GB GPU at any selected budget; the no-limit configuration would exceed that  
 1207 envelope on larger slides.

## 1208 G Detailed Ablation Results

### 1209 G.1 Numerical Results for WSI Branch Ablation

1210 Table 16 reports an internal design ablation of the GigaPixel branch, complementary to the Qwen3-VL  
 1211 (SFT, 32-patch WSI) baseline reported in the main paper. This ablation isolates the contribution of  
 1212 each subcomponent inside the GigaPixel branch itself: average pooling loses spatial heterogeneity,

1213 truncation ignores sparse diagnostic regions, V-V-only selection may miss query-relevant evidence,  
 1214 and T-V-only selection may overfocus on a small salient region.

Table 16: Numerical results for WSI branch ablation study on histopathology tasks.

Configuration	Closed Acc (%)	Caption BERT	Caption BLEU	Caption R-1
CONCH only (avg pool)	52.3	0.78	0.012	0.185
CONCH + LongNet (truncate)	61.8	0.82	0.024	0.214
CONCH + LongNet + V-V only	68.4	0.85	0.031	0.238
CONCH + LongNet + T-V + V-V (full)	<b>72.2</b>	<b>0.88</b>	<b>0.038</b>	<b>0.267</b>

1215 The results demonstrate that each component contributes incrementally. LongNet contextualization  
 1216 improves closed-ended accuracy by 9.5 points over average pooling, showing that slide-level context  
 1217 is essential for WSI reasoning. Adding V-V coverage gains another 6.6 points, indicating that global  
 1218 representativeness matters beyond sequential context aggregation. Incorporating T-V coverage yields  
 1219 the final 3.8-point improvement, confirming that clinical-query alignment is necessary once the  
 1220 selector is asked to compress the slide into a very small token budget.

## 1221 G.2 Numerical Results for Visual Token Budget Sweep

1222 Tables 17 and 18 provide the complete numerical breakdown of the visual token budget sweep  
 1223 summarized in Figure 5. We sweep  $k \in \{1, 8, 16, 32, 64, 128, 256, 512\}$  and report performance  
 1224 on every modality and every task type in BreastStage-Bench: VQA (closed- and open-ended)  
 1225 across screening, diagnosis, and treatment stages, as well as caption and report generation. This is  
 1226 intentionally exhaustive so that the operating-point choice ( $k=128$ ) used throughout the main paper  
 1227 can be audited per-modality, not only on histopathology.

Table 17: Token budget sweep on the VQA task across screening, diagnosis, and treatment stages, under closed-ended and open-ended settings. ‘‘Histo.’’ denotes histopathology.

Token	Screen								Diagnosis				Treatment			
	closed				open				closed		open		closed		open	
	BUS	CT	Mammo	MRI	BUS	CT	MRI	BUS	Mammo	MRI	BUS	MRI	MRI	Histo.	MRI	Histo.
1	60.21	78.79	43.41	72.38	91.84	95.66	94.80	51.06	48.80	65.39	87.82	94.63	40.48	66.63	89.63	56.82
8	75.32	78.43	60.43	80.95	93.76	94.89	95.47	69.94	64.41	70.28	90.53	94.99	46.83	71.25	89.22	63.78
16	84.04	77.70	72.88	81.43	95.41	94.83	95.83	75.26	67.42	78.32	92.02	95.26	59.52	71.01	89.67	66.12
32	86.17	77.70	74.05	85.24	95.80	95.03	95.58	76.86	68.92	81.29	93.08	95.71	62.70	71.01	89.96	66.53
64	86.81	76.73	74.69	83.81	95.96	95.48	95.58	77.66	68.77	81.12	92.94	95.72	61.11	71.13	89.95	66.08
128	86.81	77.21	75.00	82.86	95.97	95.29	95.48	77.13	68.32	81.12	93.24	95.72	61.11	71.38	89.93	63.80
256	86.81	76.73	74.05	82.38	95.97	95.51	95.75	77.13	68.02	81.29	93.24	95.61	59.52	71.01	89.97	66.05
512	86.81	77.58	74.47	81.90	95.96	95.70	95.65	77.13	68.32	82.69	93.23	95.67	57.14	71.50	89.89	64.58

Table 18: Token budget sweep on the caption and report generation tasks. ‘‘Wtd’’ denotes the weighted average of BERT-F1, BLEU, and ROUGE-1. The ‘‘IoU’’ column reports the *Ground Caption* task (lesion-grounded captioning), which is only defined for the modalities with bounding-region annotations.

Token	Caption												Report											
	BUS				CT				Mammo				Histopathology				MRI							
	BERT-F1	BLEU	R-1	Wtd	IoU	BERT-F1	BLEU	R-1	Wtd	IoU	BERT-F1	BLEU	R-1	Wtd	IoU	BERT-F1	BLEU	R-1	Wtd	IoU				
1	92.77	36.61	61.89	71.01	1.50	93.90	38.01	65.46	72.83	4.90	93.77	38.35	64.75	72.66	0.85	83.86	18.23	32.70	54.67	89.85	26.01	54.91	65.15	24.52
8	94.08	45.79	69.41	75.84	20.00	93.70	37.27	64.40	72.28	4.91	94.20	44.52	67.53	75.11	2.13	88.34	40.47	53.35	67.63	90.09	27.27	55.43	65.71	24.72
16	94.63	50.24	72.44	77.99	40.88	93.77	38.18	64.69	72.61	4.93	94.67	49.48	69.88	77.18	5.00	88.04	39.57	50.99	66.66	90.37	28.82	56.52	66.51	25.41
32	94.90	52.21	73.85	78.97	63.04	93.88	38.92	65.34	73.01	4.93	94.72	50.04	70.03	77.38	12.52	87.82	37.51	49.98	65.79	90.50	30.07	57.47	67.13	26.17
64	94.94	52.71	74.23	79.21	76.76	93.89	40.89	65.89	73.65	5.03	94.73	50.35	70.05	77.46	18.37	88.15	39.50	51.81	66.91	90.58	30.90	58.11	67.53	28.55
128	94.97	52.91	74.42	79.32	79.59	93.76	39.81	65.30	73.16	5.12	94.75	50.74	70.31	77.64	23.14	88.12	39.61	51.25	66.78	90.61	31.24	58.24	67.67	33.49
256	95.00	53.48	74.55	79.51	80.33	93.73	39.41	65.24	73.03	5.03	94.84	51.57	70.91	78.04	23.06	88.00	37.67	50.14	65.96	90.59	31.16	58.09	67.60	29.64
512	94.98	53.10	74.45	79.38	80.36	93.79	39.96	65.61	73.29	5.02	94.80	51.11	70.69	77.85	23.10	87.58	34.94	47.98	64.53	90.62	31.21	58.29	67.68	29.52

1228 Two patterns hold consistently across modalities. (i) *Performance saturates at  $k=128$* . On VQA  
 1229 closed-ended screening, BUS reaches its plateau at  $k=64$  (86.81%) and stays flat through  $k=512$ ;  
 1230 CT, Mammo, and MRI similarly fluctuate within  $\pm 1$  point past  $k=64$ . On caption generation, the  
 1231 weighted score on BUS rises from 71.01 at  $k=1$  to 79.32 at  $k=128$  and gains only 0.06 points by

1232  $k=512$ ; Mammo gains only 0.21 points from 128 to 512. On open-ended VQA, all six (modality,  
 1233 stage) cells are within 0.5 points of their maximum by  $k=128$ . (ii) *The ground caption task needs a*  
 1234 *minimum number of tokens*. BUS ground caption IoU jumps from 1.50 at  $k=1$  to 79.59 at  $k=128$   
 1235 before plateauing, and Mammo ground caption IoU from 0.85 to 23.14. This indicates that a single  
 1236 coverage-selected token is insufficient to anchor a lesion bounding region, but  $k=128$  already recovers  
 1237 most of the IoU achievable at  $k=512$ . The 128-token budget therefore gives BreastGPT the same  
 1238 downstream inference cost on radiology and pathology while staying on the saturation plateau of  
 1239 every metric we measure.

## 1240 H Qualitative Analysis and Case Studies

### 1241 H.1 Role-Switching Examples

1242 We present representative examples demonstrating BreastGPT’s ability to adopt different clinical  
 1243 reasoning styles based on stage-specific system prompts. These examples are intended to illustrate  
 1244 the response style induced by the workflow prompts: screening responses emphasize triage and  
 1245 recall decisions, diagnosis responses emphasize lesion characterization and differential diagnosis,  
 1246 and treatment-stage responses emphasize pathology, biomarkers, and therapy-relevant factors.

Table 19: Qualitative materials to be inserted in the final supplementary. Each row will be replaced or accompanied by a visual case figure.

Case type	Planned content	Figure placeholder
Screening	BUS or mammography example showing lesion triage, BI-RADS prediction, and recall recommendation	[PLACEHOLDER: screening case figure]
Diagnosis	Multiparametric MRI example showing lesion characterization and cross-sequence reasoning	[PLACEHOLDER: diagnosis case figure]
Treatment	WSI example showing subtype, grade, biomarker, or treatment-response reasoning	[PLACEHOLDER: treatment case figure]
Failure case	Representative error with explanation of whether the failure is perceptual, reasoning-related, or stage-confusion-related	[PLACEHOLDER: failure case figure]

1247 **Example 1: Screening Stage (BUS).** **System Prompt:** the Screening Persona of Figure 9 com-  
 1248 posed with the BUS closed-VQA task instruction.

1249 **Question:** What is the BI-RADS category for this breast ultrasound image?

1250 **BreastGPT Response:** *“This breast ultrasound demonstrates a hypoechoic mass with irregular*  
 1251 *margins and posterior acoustic shadowing in the upper outer quadrant. The lesion measures*  
 1252 *approximately 1.8 cm in maximum diameter. These features are suspicious for malignancy. I would*  
 1253 *assign BI-RADS category 4C (high suspicion for malignancy, 50-95% likelihood). Recommendation:*  
 1254 *Biopsy is indicated for tissue diagnosis.”*

1255 **Example 2: Diagnosis Stage (MRI).** **System Prompt:** the Diagnosis Persona of Figure 9 com-  
 1256 posed with the MRI open-VQA task instruction.

1257 **Question:** Characterize the lesion shown in this multiparametric MRI sequence.

1258 **BreastGPT Response:** *“The T1-weighted post-contrast images demonstrate a 2.3 cm irregular*  
 1259 *enhancing mass in the left breast at 10 o’clock position, 4 cm from the nipple. The lesion exhibits*  
 1260 *heterogeneous enhancement with rim enhancement pattern on early phase and washout kinetics on*  
 1261 *delayed phase. DWI shows restricted diffusion (ADC value approximately  $0.8 \times 10^{-3} \text{ mm}^2/\text{s}$ ). No*  
 1262 *associated nipple retraction or skin thickening. Axillary lymph nodes appear normal. TNM staging:*  
 1263 *cT2N0. Differential diagnosis favors invasive ductal carcinoma. Recommend core needle biopsy for*  
 1264 *histological confirmation and receptor status assessment.”*

1265 **Example 3: Treatment Stage (Histopathology).** **System Prompt:** the Treatment Persona of  
 1266 Figure 9 composed with the histopathology open-VQA task instruction.

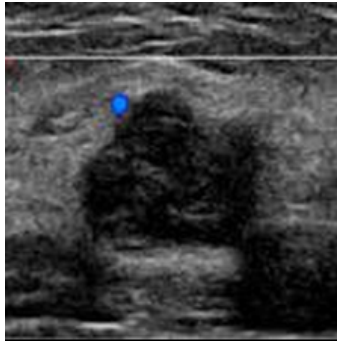
1267 **Question:** What is the molecular subtype of this breast cancer specimen?

1268 **BreastGPT Response:** “Histopathological examination reveals invasive ductal carcinoma, Nottingham grade II (tubule formation: 3, nuclear pleomorphism: 2, mitotic count: 2). The tumor cells  
1269 show strong nuclear staining for estrogen receptor (ER, 90% positive) and progesterone receptor  
1270 (PR, 70% positive). HER2 immunohistochemistry is negative (0). Ki-67 proliferation index is approx-  
1271 imately 18%. Based on these findings, the molecular subtype is Luminal A (ER+/PR+/HER2-/low  
1272 Ki-67). Recommended treatment: Adjuvant endocrine therapy (tamoxifen or aromatase inhibitor).  
1273 Chemotherapy may be considered based on genomic assay results and clinical risk factors.”  
1274

1275 These examples demonstrate that BreastGPT adapts its reasoning depth, terminology, and clinical  
1276 recommendations to match the expected behavior of specialists at each stage. This qualitative  
1277 behavior complements the role-confusion analysis in the main paper.

## 1278 H.2 Case Analysis

1279 To complement the aggregate metrics in Tables 1 and 2, we walk through six representative cases  
1280 drawn from BreastStage-Bench that span all five modalities and all four task families (closed-ended  
1281 VQA, open-ended VQA, lesion-grounded captioning, and structured report generation). For each case  
1282 we show the role-conditioned system prompt that BreastGPT receives at inference time (abbreviated  
1283 to its lead clause), the user instruction, and the reference output that defines the gold answer for the  
1284 metrics in the main paper. Together they illustrate the breadth of input formats and reasoning styles a  
1285 single BreastGPT checkpoint must support.



### Case 1 – BUS, Screening-stage Closed VQA (vascularity attribute)

**System.** You are an expert screening breast radiologist. Your clinical focus is on early detection and risk triage across broad populations. You possess a high sensitivity for identifying subtle abnormalities within varied anatomical backgrounds to determine if further investigation is warranted. Now, apply your expertise to the following task: Answer the multiple-choice question based on the provided breast ultrasound image(s). Respond ONLY with the exact option chosen.

**User.** <image>

Classify the blood flow status: No flow, Mild flow, or Abundant flow?

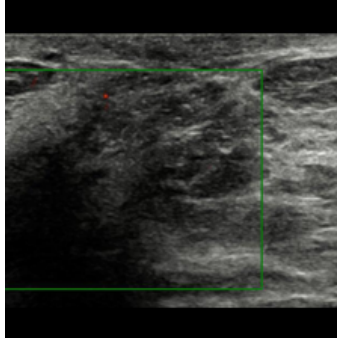
Options: (A) No flow (B) Abundant flow (C) Mild flow

Answer:

**Reference output.** (C) Mild flow

**Comment.** The screening persona drives the model toward a categorical Doppler readout rather than a free-text discussion, and the answer must come from the schema enum so it can be scored by exact-string accuracy.

Figure 12: Case 1: closed-ended vascularity attribute on a single BUS B-mode frame from BUS-CoT.



### Case 2 – BUS, Screening-stage Open-ended VQA (Doppler description)

**System.** You are an expert screening breast radiologist. Your clinical focus is on early detection and risk triage across broad populations. You possess a high sensitivity for identifying subtle abnormalities within varied anatomical backgrounds to determine if further investigation is warranted. Now, apply your expertise to the following task: Answer the open-ended question based on the provided breast ultrasound image(s).

**User.** <image>

Describe the degree of vascularity observed within the lesion on Doppler ultrasound.

**Reference output.** Mild flow is detected within the lesion.

**Comment.** The same underlying clinical fact as Case 1 is rewritten into a fluent description; this pair is what allows the weighted (BERTScore F1, BLEU, ROUGE-1) score to evaluate semantic equivalence rather than exact-string match.

Figure 13: Case 2: open-ended Doppler description on a BUS Doppler frame from BUS-CoT.



### Case 3 – BUS, Lesion-grounded Captioning

**System.** You are an expert breast clinician. Your clinical focus is on comprehensive breast disease assessment and patient-centered clinical management. You integrate clinical presentation, imaging findings, pathological evidence, risk factors, and treatment considerations to provide balanced, clinically grounded assessments across the full breast care pathway. Now, apply your expertise to the following task: Generate a descriptive caption for the provided breast ultrasound image(s). You MUST explicitly ground each mentioned finding with a 2D bounding box [xmin, ymin, xmax, ymax].

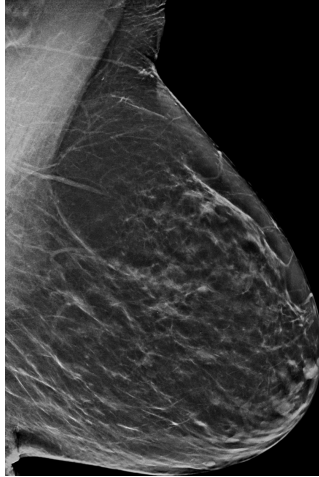
**User.** <image>

Locate any abnormalities in this breast ultrasound and provide a detailed description of their morphological features.

**Reference output.** Standard Ultrasound (B-Mode) reveals <ref-object><bbox>. This finding demonstrates microcalcifications of suspicious morphology and is categorized as BI-RADS 4A, prompting a recall for biopsy or further evaluation. Despite benign morphological features, the presence of suspicious calcifications elevates concern, though the overall prognosis remains favorable with routine follow-up.

**Comment.** The expected output interleaves a 2D bounding box with a narrative description, exercising both spatial localisation (scored by IoU in Table 2) and BI-RADS-level reasoning (scored by the weighted text metrics).

Figure 14: Case 3: lesion-grounded captioning on a BUS expert-cropped frame from BUS-CoT/BUS-Expert.



#### Case 4 – Mammography, Diagnosis-stage Closed VQA (BI-RADS)

**System.** You are an expert diagnostic breast radiologist. Your clinical focus is on definitive lesion characterization and differential diagnosis. You rigorously analyze imaging features to separate benign from malignant entities, providing conclusive assessments that directly guide patient management. Now, apply your expertise to the following task: Answer the multiple-choice question based on the provided mammogram image(s). Respond ONLY with the exact option chosen.

**User.** <image>×8

Can you provide the Bi-Rads category for this mammography exam?

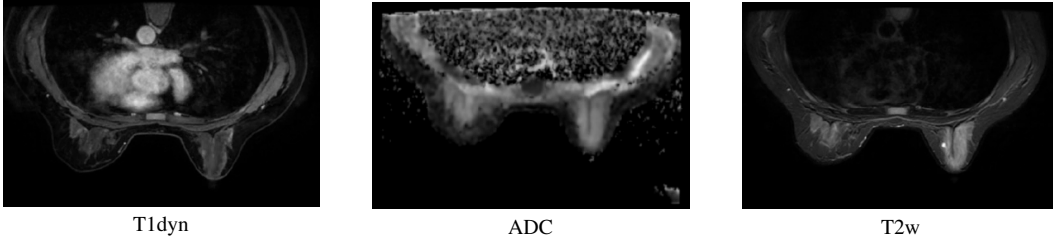
Options: (A) Bi-Rads 3 (B) Bi-Rads 6 (C) Bi-Rads 0 (D) Bi-Rads 1 (E) Bi-Rads 2 (F) Bi-Rads 4 (G) Bi-Rads 5

Answer:

**Reference output.** (D) Bi-Rads 1

**Comment.** The eight-view input forces multi-image fusion within a single forward pass; the output also illustrates BreastStage’s commitment to the full BI-RADS 0–6 enum, including the assessment-incomplete category (0) and the established-malignancy category (6) that purely diagnostic benchmarks frequently omit.

Figure 15: Case 4: BI-RADS classification on an 8-view mammographic exam from EMBED (one view shown).



### Case 5 – Multiparametric MRI, Screening-stage Report Generation

**System.** You are an expert breast clinician. Your clinical focus is on comprehensive breast disease assessment and patient-centered clinical management. You integrate clinical presentation, imaging findings, pathological evidence, risk factors, and treatment considerations to provide balanced, clinically grounded assessments across the full breast care pathway. Now, apply your expertise to the following task: Generate a comprehensive clinical report based on the provided multi-sequence breast MRI.

**User.** Breast MRI sequences: ADC: <video>. DWI: <video>. T1: <video>. T1dyn: <video>. T2w: <video>.

Based on the provided breast MRI images, generate a comprehensive and professional MRI report detailing all relevant findings, diagnostic impression with BI-RADS classification, and appropriate clinical recommendations.

**Reference output. Breast MRI Report.**

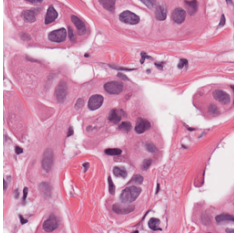
**Findings:** The bilateral breasts are symmetric in size and contour, with heterogeneously dense (ACR Category C) fibroglandular tissue and mild background parenchymal enhancement. No architectural distortion or abnormal parenchymal distribution is identified. In the posteroinferior inner quadrant of the left breast, there is an irregular, spiculated mass measuring approximately 13 mm × 11 mm × 16 mm. This lesion demonstrates marked early enhancement with washout kinetics and shows restricted diffusion on DWI/ADC sequences. No additional suspicious enhancing foci are seen in the left or right breast. Bilateral skin thickness is normal, nipples are symmetric without retraction, and there is no evidence of chest wall invasion. Axillary lymph nodes bilaterally appear morphologically benign without pathological enlargement. An incidental finding of a simple hepatic cyst is noted in the liver.

**Impression:** Irregular, markedly enhancing mass in the posteroinferior inner quadrant of the left breast with diffusion restriction, consistent with known malignancy. BI-RADS 6. Right breast shows no MRI evidence of malignancy. BI-RADS 1. Incidental simple hepatic cyst, likely benign.

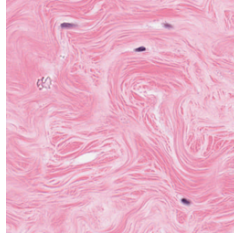
**Recommendation:** Given the BI-RADS 6 classification, this finding reflects a known or biopsy-proven malignancy. Proceed with definitive surgical management or neoadjuvant therapy planning as per multidisciplinary consensus. No further diagnostic imaging of the breast is required for this lesion. The hepatic cyst requires no intervention or follow-up.

**Comment.** The report task simultaneously requires multi-sequence integration, structured section ordering (Findings → Impression → BI-RADS → Management), and quantitative descriptors that must be supported by the actual volumes; failures concentrate here in our Table 2 numbers.

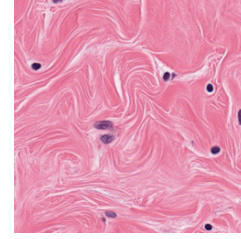
Figure 16: Case 5: comprehensive multiparametric MRI report from FUDAN; mid-slice frames of three of the five paired sequences shown above.



patch 1



patch 9



patch 17

### Case 6 – Histopathology WSI, Pathology Report Generation

**System.** You are an expert breast clinician. Your clinical focus is on comprehensive breast disease assessment and patient-centered clinical management. You integrate clinical presentation, imaging findings, pathological evidence, risk factors, and treatment considerations to provide balanced, clinically grounded assessments across the full breast care pathway. Now, apply your expertise to the following task: Generate a concise descriptive caption for the provided breast histopathological image(s).

**User.** <image> × 32

Please generate a comprehensive and professional pathology report based on the provided slide from Breast, including key elements such as: (1) the histopathological diagnosis, (2) tumor morphology and growth patterns, (3) extent of invasion and involved structures, (4) lymphovascular or perineural invasion, (5) status of surgical margins, (6) TNM staging if possible, and (7) any additional pathological findings including non-neoplastic tissue alterations.

**Reference output.** The pathology findings indicate a diagnosis of infiltrating ductal carcinoma located in the upper outer quadrant of the left breast. The tumor exhibits a moderately differentiated histologic grade, suggesting an intermediate level of aggression and atypia. There is grossly evident lesion upon inspection. Importantly, the surgical margins are uninvolved, indicating no residual tumor at the edges of the excised tissue. Lymph node analysis reveals that all ten axillary lymph nodes examined are negative for metastatic involvement, suggesting localized disease without regional spread. The Nottingham score, which typically helps in grading the tumor, could not be determined in this case. Overall, the pathological evaluation supports a moderately aggressive tumor confined to the primary site with no lymphatic dissemination.

**Comment.** On gigapixel inputs the GigaPixel branch must compress thousands of patch tokens down to  $k=128$  before the LLM ever sees them; the section structure of the reference output (diagnosis → grade → margins → nodes → TNM) is exactly what the coverage selector is asked to preserve.

Figure 17: Case 6: pathology report on a TCGA-BRCA whole-slide image at 20× magnification; 3 of the 32 sampled patches shown above.

1286 **What the cases collectively show.** The hard cases for BreastGPT cluster around three axes: (1) fine-  
1287 grained categorical boundaries when the schema admits adjacent values (BI-RADS 3 vs. 4A, Notting-  
1288 ham grade II vs. III); (2) quantitative predictions that should be supported by an explicit measurement  
1289 but are inferred from the visual signal alone (ADC values, Ki-67 index); and (3) molecular-biomarker  
1290 inference (ER/PR/HER2 status) from H&E morphology alone, without paired immunohistochemistry.  
1291 These observations motivate future work on uncertainty-aware reporting, calibrated abstention on  
1292 adjacent-category questions, and explicit integration of non-image clinical evidence at inference time.

1293 **NeurIPS Paper Checklist**

1294 **1. Claims**

1295 Question: Do the main claims made in the abstract and introduction accurately reflect the  
1296 paper’s contributions and scope?

1297 Answer: [Yes]

1298 Justification: The abstract and introduction state three concrete contributions—the workflow-  
1299 aligned BreastStage corpus and BreastStage-Bench benchmark, the unified BreastGPT  
1300 model with a dual-branch visual encoder, and concept-preserving token compression—  
1301 together with the headline numbers (75.66% closed-ended VQA, 89.92% open-ended VQA  
1302 on BreastStage-Bench). Each contribution is realized in a corresponding section (Sections 2–  
1303 3 of the main paper) and supported by Tables 1, 2 and 9 and the ablations in Sections 4–5.

1304 Guidelines:

- 1305 • The answer [N/A] means that the abstract and introduction do not include the claims  
1306 made in the paper.
- 1307 • The abstract and/or introduction should clearly state the claims made, including the  
1308 contributions made in the paper and important assumptions and limitations. A [No] or  
1309 [N/A] answer to this question will not be perceived well by the reviewers.
- 1310 • The claims made should match theoretical and experimental results, and reflect how  
1311 much the results can be expected to generalize to other settings.
- 1312 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
1313 are not attained by the paper.

1314 **2. Limitations**

1315 Question: Does the paper discuss the limitations of the work performed by the authors?

1316 Answer: [Yes]

1317 Justification: The Conclusion explicitly notes that BreastStage covers a single disease  
1318 (breast cancer) and that longitudinal patient-trajectory records are not yet integrated due  
1319 to identifiability and IRB constraints; Supplementary A.3 further discusses residual LLM-  
1320 introduced curation biases despite the expert audit (Supplementary E.3), the small panel  
1321 size (three specialists), and the dependence of the GigaPixel branch on a frozen pathology  
1322 foundation model.

1323 Guidelines:

- 1324 • The answer [N/A] means that the paper has no limitation while the answer [No] means  
1325 that the paper has limitations, but those are not discussed in the paper.
- 1326 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 1327 • The paper should point out any strong assumptions and how robust the results are to  
1328 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
1329 model well-specification, asymptotic approximations only holding locally). The authors  
1330 should reflect on how these assumptions might be violated in practice and what the  
1331 implications would be.
- 1332 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
1333 only tested on a few datasets or with a few runs. In general, empirical results often  
1334 depend on implicit assumptions, which should be articulated.
- 1335 • The authors should reflect on the factors that influence the performance of the approach.  
1336 For example, a facial recognition algorithm may perform poorly when image resolution  
1337 is low or images are taken in low lighting. Or a speech-to-text system might not be  
1338 used reliably to provide closed captions for online lectures because it fails to handle  
1339 technical jargon.
- 1340 • The authors should discuss the computational efficiency of the proposed algorithms  
1341 and how they scale with dataset size.
- 1342 • If applicable, the authors should discuss possible limitations of their approach to  
1343 address problems of privacy and fairness.

1344 • While the authors might fear that complete honesty about limitations might be used by  
1345 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
1346 limitations that aren't acknowledged in the paper. The authors should use their best  
1347 judgment and recognize that individual actions in favor of transparency play an impor-  
1348 tant role in developing norms that preserve the integrity of the community. Reviewers  
1349 will be specifically instructed to not penalize honesty concerning limitations.

### 1350 3. Theory assumptions and proofs

1351 Question: For each theoretical result, does the paper provide the full set of assumptions and  
1352 a complete (and correct) proof?

1353 Answer: [N/A]

1354 Justification: The paper is empirical and does not state formal theorems; the only analytical  
1355 content is the gradient derivation of the straight-through estimator used in the concept-driven  
1356 token selector, which is self-contained in Supplementary.

1357 Guidelines:

- 1358 • The answer [N/A] means that the paper does not include theoretical results.
- 1359 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
1360 referenced.
- 1361 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 1362 • The proofs can either appear in the main paper or the supplemental material, but if  
1363 they appear in the supplemental material, the authors are encouraged to provide a short  
1364 proof sketch to provide intuition.
- 1365 • Inversely, any informal proof provided in the core of the paper should be complemented  
1366 by formal proofs provided in appendix or supplemental material.
- 1367 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 1368 4. Experimental result reproducibility

1369 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
1370 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
1371 of the paper (regardless of whether the code and data are provided or not)?

1372 Answer: [Yes]

1373 Justification: Section D fully specify the architecture, the two-stage SFT training, the data  
1374 construction pipeline, and the full evaluation protocol. Anonymized config files, training  
1375 scripts, and inference code are released in <https://anonymous.4open.science/r/BreastGPT>.

1376 Guidelines:

- 1377 • The answer [N/A] means that the paper does not include experiments.
- 1378 • If the paper includes experiments, a [No] answer to this question will not be perceived  
1379 well by the reviewers: Making the paper reproducible is important, regardless of  
1380 whether the code and data are provided or not.
- 1381 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
1382 to make their results reproducible or verifiable.
- 1383 • Depending on the contribution, reproducibility can be accomplished in various ways.  
1384 For example, if the contribution is a novel architecture, describing the architecture fully  
1385 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
1386 be necessary to either make it possible for others to replicate the model with the same  
1387 dataset, or provide access to the model. In general, releasing code and data is often  
1388 one good way to accomplish this, but reproducibility can also be provided via detailed  
1389 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
1390 of a large language model), releasing of a model checkpoint, or other means that are  
1391 appropriate to the research performed.
- 1392 • While NeurIPS does not require releasing code, the conference does require all submis-  
1393 sions to provide some reasonable avenue for reproducibility, which may depend on the  
1394 nature of the contribution. For example  
1395 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
1396 to reproduce that algorithm.

- 1397 (b) If the contribution is primarily a new model architecture, the paper should describe  
1398 the architecture clearly and fully.
- 1399 (c) If the contribution is a new model (e.g., a large language model), then there should  
1400 either be a way to access this model for reproducing the results or a way to reproduce  
1401 the model (e.g., with an open-source dataset or instructions for how to construct  
1402 the dataset).
- 1403 (d) We recognize that reproducibility may be tricky in some cases, in which case  
1404 authors are welcome to describe the particular way they provide for reproducibility.  
1405 In the case of closed-source models, it may be that access to the model is limited in  
1406 some way (e.g., to registered users), but it should be possible for other researchers  
1407 to have some path to reproducing or verifying the results.

## 1408 5. Open access to data and code

1409 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
1410 tions to faithfully reproduce the main experimental results, as described in supplemental  
1411 material?

1412 Answer: [Yes]

1413 Justification: We release training and evaluation code, model and data configuration files,  
1414 the full BreastStage-Bench benchmark, and dataset cards through an anonymous repository  
1415 linked <https://anonymous.4open.science/r/BreastGPT>. For raw clinical images that are sub-  
1416 ject to data-use agreements (e.g., EMBED, in-house cohorts), we provide preprocessing  
1417 scripts and detailed access instructions instead of the raw files; all derived QA annotations  
1418 are released directly.

1419 Guidelines:

- 1420 • The answer [N/A] means that paper does not include experiments requiring code.
- 1421 • Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1422 • While we encourage the release of code and data, we understand that this might not  
1423 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not  
1424 including code, unless this is central to the contribution (e.g., for a new open-source  
1425 benchmark).
- 1426 • The instructions should contain the exact command and environment needed to run to  
1427 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 1428 • The authors should provide instructions on data access and preparation, including how  
1429 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1430 • The authors should provide scripts to reproduce all experimental results for the new  
1431 proposed method and baselines. If only a subset of experiments are reproducible, they  
1432 should state which ones are omitted from the script and why.
- 1433 • At submission time, to preserve anonymity, the authors should release anonymized  
1434 versions (if applicable).
- 1435 • Providing as much information as possible in supplemental material (appended to the  
1436 paper) is recommended, but including URLs to data and code is permitted.

## 1439 6. Experimental setting/details

1440 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
1441 rameters, how they were chosen, type of optimizer) necessary to understand the results?

1442 Answer: [Yes]

1443 Justification: Section 2 and Section 4.1 report the train/test splits for every constituent dataset.  
1444 Supplementary D.4 report the optimizer (AdamW), learning-rate schedule, global batch size,  
1445 `max_length`. Evaluation metrics are explicitly defined: accuracy for closed-ended VQA,  
1446 and the weighted score  $0.5 \text{BERTScore} + 0.25 \text{BLEU} + 0.25 \text{ROUGE-1}$  for open-ended  
1447 generation.

1448 Guidelines:

- 1449 • The answer [N/A] means that the paper does not include experiments.

- 1450 • The experimental setting should be presented in the core of the paper to a level of detail  
1451 that is necessary to appreciate the results and make sense of them.
- 1452 • The full details can be provided either with the code, in appendix, or as supplemental  
1453 material.

## 1454 7. Experiment statistical significance

1455 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
1456 information about the statistical significance of the experiments?

1457 Answer: [No]

1458 Justification: We report single-run results because each evaluation cell traverses the full  
1459 held-out BreastStage-Bench split (10k+ samples per task family), and the cross-model gaps  
1460 we report (often >20 absolute points between BreastGPT and the best baseline) are an order  
1461 of magnitude larger than the seed-to-seed variance we observed in pilot runs. Re-training  
1462 Qwen3-VL-8B with multiple SFT seeds for explicit error bars was not affordable within our  
1463 compute budget (Supplementary A.4).

1464 Guidelines:

- 1465 • The answer [N/A] means that the paper does not include experiments.
- 1466 • The authors should answer [Yes] if the results are accompanied by error bars, confidence  
1467 intervals, or statistical significance tests, at least for the experiments that support the  
1468 main claims of the paper.
- 1469 • The factors of variability that the error bars are capturing should be clearly stated (for  
1470 example, train/test split, initialization, random drawing of some parameter, or overall  
1471 run with given experimental conditions).
- 1472 • The method for calculating the error bars should be explained (closed form formula,  
1473 call to a library function, bootstrap, etc.)
- 1474 • The assumptions made should be given (e.g., Normally distributed errors).
- 1475 • It should be clear whether the error bar is the standard deviation or the standard error  
1476 of the mean.
- 1477 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
1478 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
1479 of Normality of errors is not verified.
- 1480 • For asymmetric distributions, the authors should be careful not to show in tables or  
1481 figures symmetric error bars that would yield results that are out of range (e.g., negative  
1482 error rates).
- 1483 • If error bars are reported in tables or plots, the authors should explain in the text how  
1484 they were calculated and reference the corresponding figures or tables in the text.

## 1485 8. Experiments compute resources

1486 Question: For each experiment, does the paper provide sufficient information on the com-  
1487 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
1488 the experiments?

1489 Answer: [Yes]

1490 Justification: Supplementary A.4 reports compute used for both training (Stage-1 + Stage-  
1491 2 SFT on multi-node H100 clusters) and inference (single-GPU profiling). Inference  
1492 latency, prefill cost, and peak GPU memory at every visual-token budget  $k$  are tabulated in  
1493 Table 15, and training hyperparameters (optimizer, schedule, batch size, `max_length`) are  
1494 in Supplementary D.4.

1495 Guidelines:

- 1496 • The answer [N/A] means that the paper does not include experiments.
- 1497 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
1498 or cloud provider, including relevant memory and storage.
- 1499 • The paper should provide the amount of compute required for each of the individual  
1500 experimental runs as well as estimate the total compute.
- 1501 • The paper should disclose whether the full research project required more compute  
1502 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
1503 didn't make it into the paper).

1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: We have reviewed the NeurIPS Code of Ethics. All clinical data is de-identified at the source and used under the existing data-use agreements of the contributing public datasets and partner hospitals (Supplementary A.2); no new human subjects were recruited for this work; expert audit work (Supplementary E.3) is performed by board-certified clinicians as part of an existing professional collaboration. Dual-use considerations and a research-only release stance are discussed in the broader-impacts section.

Guidelines:

- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The Conclusion / Limitations and Supplementary A.3 discuss positive impacts (workflow-aware second-opinion support across the breast-cancer continuum, particularly in regions without dense specialist coverage) and negative impacts (over-reliance on model outputs without expert review, inherited demographic biases of the contributing public cohorts, and the general dual-use risk of releasing a clinical-domain MLLM). Mitigations include the research-only license attached to the release, a dataset card documenting demographic coverage gaps, and the explicit disclaimer against autonomous clinical use.

Guidelines:

- The answer [N/A] means that there is no societal impact of the work performed.
- If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

1556 Question: Does the paper describe safeguards that have been put in place for responsible  
1557 release of data or models that have a high risk for misuse (e.g., pre-trained language models,  
1558 image generators, or scraped datasets)?

1559 Answer: [Yes]

1560 Justification: The released BreastGPT checkpoint is gated by a research-only license stipu-  
1561 lating that the model is not for autonomous clinical decision-making, and the model card  
1562 (Supplementary A.5) ships an explicit disclaimer to that effect. Raw clinical images are  
1563 released only via the access-controlled provenance routes of the underlying datasets (we  
1564 redistribute only derived QA annotations and structured records); the dataset card addition-  
1565 ally documents demographic and modality-coverage gaps so downstream users can avoid  
1566 out-of-distribution deployment.

1567 Guidelines:

- 1568 • The answer [N/A] means that the paper poses no such risks.
- 1569 • Released models that have a high risk for misuse or dual-use should be released with  
1570 necessary safeguards to allow for controlled use of the model, for example by requiring  
1571 that users adhere to usage guidelines or restrictions to access the model or implementing  
1572 safety filters.
- 1573 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
1574 should describe how they avoided releasing unsafe images.
- 1575 • We recognize that providing effective safeguards is challenging, and many papers do  
1576 not require this, but we encourage authors to take this into account and make a best  
1577 faith effort.

## 1578 12. Licenses for existing assets

1579 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
1580 the paper, properly credited and are the license and terms of use explicitly mentioned and  
1581 properly respected?

1582 Answer: [Yes]

1583 Justification: All public datasets reused in BreastStage (CT-RATE, BUS-CoT, EMBED,  
1584 public histopathology slide collections, etc.) and pretrained backbones (Qwen3-VL-8B,  
1585 CONCHv1.5, LongNet) are cited at first use and used under their respective listed licenses  
1586 (research-only or CC-BY variants); we redistribute only derived QA annotations and struc-  
1587 tured records, not raw images that are subject to a more restrictive license. Per-asset license,  
1588 version, and access route are listed in Supplementary C.1 and D.1.

1589 Guidelines:

- 1590 • The answer [N/A] means that the paper does not use existing assets.
- 1591 • The authors should cite the original paper that produced the code package or dataset.
- 1592 • The authors should state which version of the asset is used and, if possible, include a  
1593 URL.
- 1594 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1595 • For scraped data from a particular source (e.g., website), the copyright and terms of  
1596 service of that source should be provided.
- 1597 • If assets are released, the license, copyright information, and terms of use in the  
1598 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
1599 has curated licenses for some datasets. Their licensing guide can help determine the  
1600 license of a dataset.
- 1601 • For existing datasets that are re-packaged, both the original license and the license of  
1602 the derived asset (if it has changed) should be provided.
- 1603 • If this information is not available online, the authors are encouraged to reach out to  
1604 the asset’s creators.

## 1605 13. New assets

1606 Question: Are new assets introduced in the paper well documented and is the documentation  
1607 provided alongside the assets?

1608 Answer: [Yes]

1609 Justification: We release three new assets: (i) BreastStage (1.86M instruction-following  
1610 pairs over 136 task templates), (ii) BreastStage-Bench (held-out evaluation split), and (iii)  
1611 the BreastGPT model checkpoint with its prompt library. All three are accessible through  
1612 an anonymous repository at <https://anonymous.4open.science/r/BreastGPT>. Doc-  
1613 umentation includes a dataset card (source datasets, demographic and modality coverage,  
1614 structured-record schema, four-stage curation pipeline, expert validation results), a model  
1615 card with the research-only usage statement, and per-template prompt files; pointers are  
1616 summarised in Supplementary A.5.

1617 Guidelines:

- 1618 • The answer [N/A] means that the paper does not release new assets.
- 1619 • Researchers should communicate the details of the dataset/code/model as part of their  
1620 submissions via structured templates. This includes details about training, license,  
1621 limitations, etc.
- 1622 • The paper should discuss whether and how consent was obtained from people whose  
1623 asset is used.
- 1624 • At submission time, remember to anonymize your assets (if applicable). You can either  
1625 create an anonymized URL or include an anonymized zip file.

#### 1626 14. Crowdsourcing and research with human subjects

1627 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1628 include the full text of instructions given to participants and screenshots, if applicable, as  
1629 well as details about compensation (if any)?

1630 Answer: [N/A]

1631 Justification: We do not employ crowdworkers and do not recruit new human subjects.  
1632 The only human input to dataset construction is the breast-specialist audit described in  
1633 Supplementary E.3, performed by board-certified clinicians within an existing professional  
1634 collaboration with the project; no compensation beyond their normal employment is pro-  
1635 vided. The audit instructions and per-dimension rubric are documented in that section.

1636 Guidelines:

- 1637 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
1638 with human subjects.
- 1639 • Including this information in the supplemental material is fine, but if the main contribu-  
1640 tion of the paper involves human subjects, then as much detail as possible should be  
1641 included in the main paper.
- 1642 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
1643 or other labor should be paid at least the minimum wage in the country of the data  
1644 collector.

#### 1645 15. Institutional review board (IRB) approvals or equivalent for research with human 1646 subjects

1647 Question: Does the paper describe potential risks incurred by study participants, whether  
1648 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1649 approvals (or an equivalent approval/review based on the requirements of your country or  
1650 institution) were obtained?

1651 Answer: [Yes]

1652 Justification: The clinical multiparametric MRI cohort contributed by the partner hospitals  
1653 is used under the institutional IRB approvals on file at those institutions, all images are  
1654 de-identified at the source, and participants gave written informed consent for research use  
1655 of their imaging at original collection time. Public datasets reused (CT-RATE, BUS-CoT,  
1656 EMBED, etc.) are governed by their published data-use agreements. To preserve anonymity  
1657 for review, the partner-institution IRB protocol numbers are withheld in this submission and  
1658 will be disclosed in the camera-ready version (Supplementary A.2).

1659 Guidelines:

- 1660 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
1661 with human subjects.

- 1662 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1663 may be required for any human subjects research. If you obtained IRB approval, you  
1664 should clearly state this in the paper.
- 1665 • We recognize that the procedures for this may vary significantly between institutions  
1666 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
1667 guidelines for their institution.
- 1668 • For initial submissions, do not include any information that would break anonymity (if  
1669 applicable), such as the institution conducting the review.

#### 1670 16. Declaration of LLM usage

1671 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
1672 non-standard component of the core methods in this research? Note that if the LLM is used  
1673 only for writing, editing, or formatting purposes and does *not* impact the core methodology,  
1674 scientific rigor, or originality of the research, declaration is not required.

1675 Answer: [Yes]

1676 Justification: LLMs are a non-standard component of our dataset-construction methodology  
1677 and are explicitly declared. Qwen2.5-VL-72B serves as a modality-specific quality-control  
1678 agent and visual-attribute extractor whenever a decision requires looking at the image;  
1679 Qwen3-Max performs Chinese-to-English report parsing, structured-record-to-open-VQA  
1680 rewriting, ground-caption synthesis, and report drafting under strict JSON output schemas.  
1681 The full LLM persona library, prompt templates, and assembly rule are documented in  
1682 Supplementary C.2 and Supplementary A.1. BreastGPT itself is initialized from Qwen3-  
1683 VL-8B-Instruct, which is also disclosed.

1684 Guidelines:

- 1685 • The answer [N/A] means that the core method development in this research does not  
1686 involve LLMs as any important, original, or non-standard components.
- 1687 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not  
1688 be described.