Exploring Compositional Generalization of Multimodal LLMs for Medical Imaging

Zhenyang Cai[†], Junying Chen[†], Rongsheng Wang[†], Weihong Wang, Yonglin Deng, Dingjie Song, Yize Chen, Zixu Zhang, Benyou Wang^{*} The Chinese University of Hong Kong, Shenzhen *wangbenyou@cuhk.edu.cn*

Abstract

Medical imaging provides essential visual insights for diagnosis, and multimodal large language models (MLLMs) are increasingly utilized for its analysis due to their strong generalization capabilities; however, the underlying factors driving this generalization remain unclear. Current research suggests that multitask training outperforms single-task as different tasks can benefit each other, but they often overlook the internal relationships within these tasks. To analyze this phenomenon, we attempted to employ compositional generalization (CG), which refers to the models' ability to understand novel combinations by recombining learned elements, as a guiding framework. Since medical images can be precisely defined by Modality, Anatomical area, and Task, naturally providing an environment for exploring CG, we assembled 106 medical datasets to create Med-MAT for comprehensive experiments. The experiments confirmed that MLLMs can use CG to understand unseen medical images and identified CG as one of the main drivers of the generalization observed in multi-task training. Additionally, further studies demonstrated that CG effectively supports datasets with limited data and confirmed that MLLMs can achieve CG across classification and detection tasks, underscoring its broader generalization potential. Med-MAT is available at https:// github.com/FreedomIntelligence/Med-MAT.

1 Introduction

Medical imaging provides essential visual insights into the structures of the human body, making it a critical tool for medical diagnosis. Recently, multimodal large language models (MLLMs) (Liu et al., 2023; Li et al., 2024; Chen et al., 2024b) have been employed to analyze these images due to their strong interpretability and generalization capabil-



Figure 1: Examples of *Compositional Generalization*: The model is required to understand unseen images by recombining the fundamental elements it has learned.

ities. In this paper, we focus on the latter: generalization of MLLMs in medical imaging. Current research (Mo and Liang, 2024; Ren et al., 2024) has demonstrated that models trained on multiple tasks outperform those trained on a single task as they can leverage potential knowledge from other tasks. Yet, the underlying factors that contribute to this generalization remain insufficiently explored.

To this end, we take the perspective of *composition generalization* (CG) (Li et al., 2019; Xu et al., 2022; Tang et al., 2024) to investigate the generalization phenomenon of mutual improvement in MLLMs' understanding of medical images. Specifically, CG is the model's ability to learn fundamental elements and recombine them in novel ways to understand unseen combinations (e.g., learning *Cat* from *White Cat* and *Black* from *Black Dog*, then generalizing to *Black Cat*, as shown in Figure 1).

In this paper, we categorize each image to three elements: Modality \triangle , Anatomical area \clubsuit , and medical Task B, presenting numerous natural opportunities for CG. We defined these three elements as the **MAT-Triplet** and collected 106 medical datasets, subsequently merging those that share the same *MAT-Triplet* to create the **Med-MAT** dataset.

[†]Equal Contribution. *Corresponding author.



Figure 2: The process of integrating a vast amount of labeled medical image data to create Med-MAT.

Ultimately, Med-MAT comprises 53 subsets, encompassing 11 modalities, 14 anatomical regions, and 13 medical tasks, providing a foundation for investigating CG and other generalization methods.

To verify the existence of CG, we designated specific datasets as *Target* data and selected all *Related* data from Med-MAT that shared the same MAT-Triplet with the *Target* data. Using these data combinations, we accessed the generalization performance of MLLMs and observed that they could leverage CG to understand unseen medical images. To further validate this finding, we repeated the experiments on different MLLMs and obtained consistent results, confirming the universality of CG.

Building on these insights, we expanded the number of combinations and observed the changes in model generalization performance after deliberately disrupting CG, ultimately revealing that CG is a key factor driving the generalization of MLLMs. Furthermore, we explored the potential applications of CG and its performance across classification and detection tasks, finding that CG enhances MLLMs' ability to handle medical scenarios with limited training data and improves their capacity for spatial awareness.

Here are the key contributions of our work: 1) A VQA dataset, Med-MAT, has been constructed, providing a platform to explore the generalization of MLLMs on medical images. 2) Through this dataset, we observed that MLLMs in different architectures can utilize compositional generalization to understand unseen images and demonstrated that this is one of the main forms of generalization for medical MLLMs. 3) Finally, the real-world applicability of CG, along with its presence across detection and classification tasks, has been further explored, highlighting its potential to enhance dataefficient training and its broad applicability.

2 A Pilot Study on Generalization

2.1 Data Collection (Med-MAT)

Most existing datasets for MLLMs (Zhang et al., 2023c; Li et al., 2024; Chen et al., 2024b), primarily VQA datasets, provide broad coverage but lack attribute annotations for individual samples, which are not suitable for CG exploration. To address this gap, we curated a large collection of image-text pairs to develop **Med-MAT**, ensuring that each sample is explicitly defined by MAT-Triplet.

Data Construction Med-MAT contains a total of 106 image-label pair medical datasets, sourced from various medical public challenges or high-quality annotated datasets. All datasets are categorized according to their MAT-Triplet, with data having identical elements grouped into a single subset (Figure 2). Labels are manually clustered to ensure that annotations with the same meaning are not repeatedly used. In total, Med-MAT covers 11 medical modalities \triangle , 14 anatomical areas \widehat{A} , and 13 medical tasks $\widehat{\blacksquare}$, hoping that it can spread across various medical tasks like a mat. (Data lists are shown in Appendix B)

Data Distribution All subsets are divided into training and test sets following their original distributions or using a 9:1 ratio. To ensure a fair comparison, each training set is limited to 3,000 samples ¹, with label balance maintained as much as possible. Any subset that cannot meet this requirement is treated as an OOD (out-of-distribution) dataset. For the test sets, we strictly balance the number of samples per label to ensure that the accuracy metric reliably reflects model performance.

QA Pairs Construction To enable MLLMs to directly train and test on Med-MAT, all image-label

¹Most datasets contain around 3,000 samples.

Subset No.	02	03	07	08	09	11	13	14	15	16	18	19	21	22	23	25	26	28	30	31	32	33	35	36	37
Baseline	21	47	40	25	26	27	28	24	22	24	25	23	49	26	25	24	49	30	49	21	49	20	25	23	19
Single-task Training	<u>24</u>	<u>49</u>	<u>50</u>	<u>68</u>	<u>65</u>	<u>76</u>	<u>83</u>	<u>53</u>	<u>61</u>	<u>32</u>	<u>29</u>	<u>26</u>	<u>57</u>	<u>53</u>	<u>28</u>	<u>24</u>	<u>57</u>	<u>64</u>	<u>89</u>	<u>60</u>	<u>97</u>	<u>54</u>	<u>29</u>	<u>51</u>	<u>49</u>
Multi-task Training	96	89	80	80	79	97	92	88	76	57	88	74	87	86	93	52	98	72	94	61	$1\overline{00}$	72	75	60	50

Table 1: Accuracy(%) of different models on In-Distribution datasets (each dataset contains over 3,000 samples, with 3,000 selected for training). Within each segment, **bold** highlights the best scores, and <u>underline</u> indicates the second-best. *Baseline* represents the results without any training, *Single-task Training* refers to the results after training on a single dataset, and *Multi-task Training* represents the results after training on all datasets.

Subset No.	01	04	05	06	10	12	17	20	24	27	29	34
Baseline	32	25	33	33	48	27	33	13	34	37	31	20
Multi-task Training	39	26	70	31	58	38	61	40	35	41	55	50

Table 2: Accuracy(%) of different models on Out-Of-Distribution Dataset (each dataset contains fewer than 3,000 samples and is used only for testing). **Bold** highlights the best scores. *Multi-task Training* represents the results after training on all datasets.

paired data were converted into a visual questionanswering (VQA) format (Figure 3). Specifically, each subset was manually assigned 6 instructions to guide the MLLM in answering the subset task. For convenience, all samples were converted into single-choice questions with up to four options, and the remaining distractor options were randomly drawn from other labels within the subset. To mitigate potential evaluation biases arising from varying option counts, the ImageWikiQA dataset (Zhang et al., 2024b), a non-medical dataset consisting of single-answer, four-option questions, was incorporated during the training.



Figure 3: An example of formatting a raw classification sample into a Question-answering sample in Med-MAT.

2.2 Observation

Experiment Setup We chose LLaVA-v1.5-7B-Vicuna (Liu et al., 2023) as the base model due

to its transparent pretraining process and minimal use of medical data, reducing the risk of knowledge leakage. Leveraging MLLM's flexibility, we enabled task switching and generalization by adjusting prompts, streamlining generalization studies. Each experiment ran for 5 epochs on 8 A800 GPUs with a batch size of 32 and a learning rate of 5e-6.

Analysis To access the generalization of MLLMs, we trained the baseline on all ID datasets to simulate *Multi-task Training* and separately trained on individual ID datasets to establish the *Single-task Training* as the control group. We then evaluated the models on all datasets. The results in Table 1 and 2 confirm that *Multi-task Training* outperformed *Single-task Training* on specific tasks and improved OOD prediction, suggesting certain data combinations enhance classification and identifying valuable combinations for medical tasks warrants further research. This observation leads to a research question (RQ):

What drives the generalization observed in MLLMs during Multi-task Training?

To address it, we aim to explore the generalization mechanism of MLLMs from the perspective of compositional generalization (CG).

3 Proof of Concept on CG

This section will prove the existence of CG in MLLMs, offering preliminary insights to address the RQ and providing support for further analysis.

3.1 Experiment Setup

To explore the existence of CG from a finer perspective, this section focuses on CG with only two

	Related Cor	nbination		Target S	Baseline	Baseline+	Trained	CG Helps	
杀Lung	≜ COVID	登Brain	Cancer	從Lung	≜ Cancer	25	25	27	1
∛Lung	[≜] Cancer	帶Brain	₿State	[™]Lung	≣State	47	46	50	1
FBrain		帶Lung	₿State	ABrain	≣State	33	50	57	1
FBones	[≜] Level	帶Lung	₿State	Bones	≣State	49	53	51	×
FBones	ÊLevel	^ABrain	₿State	^ABones	≣State	49	53	72	1
*Bones	≜Level	*Breast	Diseases	*Bones	Diseases	37	33	39	1
FBones	[≜] Level	帶Lung	Diseases	Bones	Diseases	37	33	43	1
FBones	≜Level	* Chest	Diseases	Bones	ÊDiseases	37	31	43	1
FBones	State	[™] Breast	Diseases	Bones	ÊDiseases	37	37	43	1
FBones	₿State	补Lung	Diseases	ABones	Diseases	37	37	43	1
FBones	[≜] State	[™] Chest	Diseases	ABones	Diseases	37	37	41	1
帶Lung	ÊCOVID	[™]Breast	Diseases	[™]Lung	ÊDiseases	49	48	51	1
從Lung	ÊCOVID	發Bones	Diseases	帶Lung	ÎDiseases	49	48	52	1
杀Lung	ÊCOVID	登Chest	Diseases	帶Lung	Diseases	49	48	51	1
ACT	[≜] Cancer	≜X-rav	ÊCOVID	≜CT	ÊCOVID	47	46	72	1
≜X-ray	Diseases	≜CT	≜COVID	∆X-ray	ÊCOVID	30	21	49	1
≜X-ray	[≜] Diseases	≜CT	₿State	^A X-ray	≣State	30	21	46	1
≜CT	₿State	≜X-ray	Cancer	≜CT	Cancer	33	28	28	X
[≜] X-ray	發Bones	ACT	A Brain	≜X-ray	#Brain	49	49	91	1
[≜] X-ray	從Lung	≜CT	₩Brain	[≜] X-ray	猪Brain	49	50	81	1
[≜] X-ray	帝Bones	≜CT	FBrain	[≜] X-ray	芬Brain	25	51	74	1
≜X-ray	帶Lung	≜CT	F Brain	≜X-ray	芬Brain	49	52	52	X
≜CT	[™] Lung	[≜] X-ray	帝Brain	≜CT	FBrain	33	50	60	1
≜CT	#Brain	[≜] X-ray	輩Lung	≙CT	莽Lung	25	25	36	1
≜CT	發Brain	[≜] X-ray	补Lung	≜CT	莽Lung	47	50	81	1
≜CT	₩Brain	≜X-ray	Lung	≜CT	帶Lung	47	50	71	1
≜X-ray	帝Bones	≜CT	帝Lung	≜X-ray	帶Lung	30	32	28	X
≜X-ray	帝Brain	≜CT	蒣Lung	≜X-ray	蒣Lung	30	32	35	1
≜X-ray	ABones	≜CT	帝Lung	≜X-ray	帶Lung	30	32	41	1
[⊉] X-ray	₩Brain	≜CT	帶Lung	≜X-ray	%Lung	30	32	42	1
🖓 Der - Skin	[≜] Cancer	FP - Fundus	Diseases	🖓 Der - Skin	Diseases	25	29	33	1
🖾 Der - Skin	Cancer	CCT - Retine	Diseases	🖾 Der - Skin	ÊDiseases	25	29	33	1
🖾 Der - Skin	Diseases	DP - Mouth	Cancer	🖓 Der - Skin	ÊCancer	40	33	63	1
🖾 Der - Skin		Mic - Cell	Cancer	🖓 Der - Skin	ÊCancer	40	33	63	1
DP - Mouth	₿State	🖓 Der - Skin	Cancer	DP - Mouth	ÊCancer	48	50	52	1
DP - Mouth	[−] [−] [−] [−]	Mic - Cell	Cancer	DP - Mouth	ÊCancer	48	50	55	1
FP - Fundus		Mic - Cell	ÊLevel	FP - Fundus	≣Level	33	36	42	1
Mic - Cell	Recognition	FP - Fundus	Level	Mic - Cell	la∎Level	23	33	32	×
Mic - Cell	Recognition	Der - Skin	Cancer	Mic - Cell	Cancer	49	50	50	X
Mic - Cell	Recognition	DP - Mouth	Cancer	Mic - Cell	Cancer	49	51	62	1
Mic - Cell	≣Level	Der - Skin	Cancer	Mic - Cell	Cancer	49	51	52	1
Mic - Cell	≣Level	DP - Mouth	Cancer	Mic - Cell	Cancer	49	51	58	1
Mic - Cell	IECancer	SFP - Fundus	ELevel	ẩMic - Cell	≣Level	23	24	27	1

Table 3: Generalization results on classification datasets: *Related Combination* is the training set, *Target Subset* is the goal. *Baseline*, *Baseline*+, and *Trained* represent the model's accuracy(%) without training, trained on randomly sampled unrelated data, and trained on related data, respectively. \checkmark in *CG Helps* indicates successful generalization, while \varkappa denotes failure. The 4 segmented areas represent different Direction Types: fixed modality \triangleq , fixed area $\frac{\pi}{2}$, fixed task \triangleq , and modality-area paired combinations \blacksquare . Although some combinations share the same name, they differ because they fix different elements.

MAT-Triplet elements varying while the third remains constant. Additionally, we identified specific Modality-Area pairs a, such as dermoscopy paired consistently with skin, which were treated as a special category. These 4 different fixed formats were classified into distinct *Direction Types*.

We adhered to the training setup described in Section 2.2 and evaluated the model's performance on the *Target* data. *Baseline* refers to the model without any training, while *Trained* refers to the model trained solely on *Related* data. To ensure that our conclusions are not influenced by the amount of training data, we randomly sampled an equal number of data from the *Unrelated* subsets, and this configuration is referred to as *Baseline*+.

3.2 Results

Results are shown in Table 3 and it can be observed that almost all CG combinations are able to generalize to downstream tasks, highlighting that MLLMs can leverage CG to generalize *Target* data across all Direction Types. Besides that, since this experiment focused solely on two-element tuples, we further investigated three-element tuples in Appendix A.4, where we also observed similarly strong generalizations when obtaining MAT-Triplet elements from three different datasets.

Take-away 1: *MLLMs can leverage CG to understand unseen medical images.*

In the *Baseline*+ setting, we removed all datasets sharing any MAT-Triplet element with the *Target*

	Related Con	nbination		Target S	Qwen	Llama	
🖗 Bones	🖹 State	🖗 Breast	🖹 Diseases	🖗 Bones	🖹 Diseases	+4	+7
🖗 Lung	🖹 COVID	🖗 Bones	🖹 Diseases	🖗 Lung	🖹 Diseases	+11	+11
▲ X-ray	🖹 Diseases	≜ CT	🖹 COVID	▲ X-ray	🖹 COVID	+5	+5
🗟 X-ray	🖹 Diseases	🗳 CT	🖹 State	🗳 X-ray	🖹 State	+8	+8
🗟 CT	🖗 Brain	🗳 X-ray	🖗 Lung	A CT	🖗 Lung	+1	-2
🗟 CT	🖗 Brain	▲ X-ray	🖗 Lung	A CT	🖗 Lung	+7	+8
최 FP - Fundus	🖹 Diseases	🔊 Mic - Cell	🖹 Level	🔊 FP - Fundus	🖹 Level	-3	+6
🗟 Mic - Cell	Recognition	🗟 FP - Fundus	🖹 Level	🗟 Mic - Cell	🖹 Level	+7	+22

Table 4: Result of Qwen2-VL and Llama-3.2-Vision on selected classification datasets in Med-MAT. *Qwen* and *Llama* represent the accuracy(%) gains they achieved on the respective backbones through CG.

data. Consequently, *Baseline*+ models perform at near-random levels on the test set, indicating they failed to acquire target-relevant knowledge. This suggests that only datasets related through the MAT-Triplet can help the model learn and generalize to new target tasks.

Take-away 2: Generalization arises in medical datasets in which at least partial MAT elements pre-exist during training.

3.3 Extending CG to other Backbones

LLaVA was selected as the baseline because its training data and processes are publicly available, ensuring minimal exposure to medical images and preventing bias in the integration of medical image knowledge into the MLLM. To ensure that the results are not affected by the training data or the visual encoder of LLaVA, we randomly sampled two combinations from each *Direction Type* to investigate CG on Qwen2-VL-7B (Wang et al., 2024a) and Llama3.2-11B-Vision (Meta AI, 2024).

Qwen2-VL undergoes additional training on proprietary data based on ViT and incorporates a strategy to adjust the number of vision tokens according to resolution. Llama3.2-Vision, on the other hand, pretrains its own vision encoder from scratch using proprietary data. Thus, both models serve as a means to assess whether MLLMs with different training data and vision encoders can still leverage CG to understand unseen images, ensuring that CG is not merely an artifact of LLaVA's data fitting or specific to its vision encoder.

Table 4 presents the experimental results, showing that both selected backbones exhibit a certain degree of generalization across most tasks. This suggests that despite differences in pre-train data and vision encoders, different MLLMs can still leverage CG to understand unseen images.

Take-away 3: CG persists across different MLLM backbones.

4 Scaling Combination in CG

After confirming that CG is indeed a form of generalization in MLLMs, we expanded the number of participating combinations to explore the generalizability of CG and examine its relationship with the generalization exhibited by *Multi-task Training* to address the **RQ**.

4.1 Experiment Setup

Two sub-questions have been defined to verify the applicability of CG in multiple data combinations and examine its role in *Multi-task Training*.

- (Q1) While previous experiment on CG indicated that *Unrelated* combinations provide no benefit to *Target* data, can generalization arise when training incorporates more *Unrelated* combinations, simulating a multi-task scenario?
- (Q2) Previous studies suggest that *Multi-task Training* generally promotes better generalization than single-task training. If the CG conditions in *Multi-task Training* are deliberately disrupted, will the resulting generalization effect be affected?

Selection Strategy To ensure a balanced evaluation of *Related* and *Unrelated* combinations, Subset 03 and Subset 28 were chosen as *Target* datasets because they exhibit the most balanced ratios of *Related* to *Unrelated* subsets (13:11 for Subset 03 and 11:13 for Subset 28), making them ideal for providing a diverse range of compositions in the scale-up experiments.

The baseline was trained on all subsets excluding the *Target* data to evaluate the claim that mixing multi-task data enhances generalization (*All Data*). To construct multiple comparative experiments, models were further trained on either *Related* or *Unrelated* subsets (*All Related / All Unrelated*) to address Q1. For



Figure 4: Accuracy(%) results on the *Target* dataset for various models. *All Related/Unrelated* models are trained on all the related or unrelated datasets of the *Target* Data. *w/o Modality/Area/Task* are trained on All Related datasets but omit those sharing the same element as the *Target* Data, to intentionally disrupt CG. *All Data* uses all available training sets. (Note: The *Target* Data is excluded from training to observe generalization.)

Q2, individual MAT-Triplet elements were systematically removed from the *Related* subsets (*Related w/o Modality / Area / Task*), disrupting CG and assessing the ability to maintain generalization. To ensure consistency, the total data volume in all experiments was limited to 15,000 samples, aligning with the number of ID subsets available after excluding related tasks from Subset 03.

4.2 Analysis of Scaling Experiment

Figure 4 illustrates the results. It can be observed that even when we expanded the *Unrelated* combination volumes and increased task diversity, the performance of *All Unrelated* remains close to the *Baseline*, indicating that these datasets can not support MLLMs to understand the *Target* data.

Take-away 4: Datasets without MAT-Triplet overlap offer limited benefit for generalization even in the multi-task training scenario (Q1).

Besides, *w/o Modality / Area / Task* showed significant accuracy drops compared to *All Related*, despite holding the training data volume constant. This indicates that if the CG combinations are forcibly disrupted, MLLMs will lose a significant amount of generalization capability for the target data.

Take-away 5: Disrupting CG leads to a significant decline in generalization ability. (Q2).

Notably, *All Related* achieves a performance level comparable to *All Data*, where all datasets are included in training. This suggests that CG plays a crucial role in enhancing the generalization effect of *Multi-task Training*. Therefore, in conclusion:

Take-away 6: *CG plays an important role in generalization for MLLMs in medical imaging.*

5 Potential Applications of CG

As MLLMs can use CG to generalize unseen medical images, this section attempts to explore its potential applications in training medical MLLMs.

5.1 Generalization without *Target* Data

In medical tasks, new and unpredictable conditions, like COVID-19, can emerge at any time. Exploring how to use CG to help MLLMs enhance their ability to identify unknown diseases in the absence of specific datasets is both important and meaningful.

We selected some *Target* datasets and trained the MLLMs using *Related* and *Unrelated* data to observe their generalization to the *Target* data. The generalization trend was assessed by progressively increasing the size of the combination datasets.

Selection Strategy To highlight the generalization trends, the combinations with strong generalization results were selected from the main experiments. For fairness, we chose the combinations across four types where *Trained* results exceed both *Baseline* and *Baseline*+ by at least 10. If multiple combinations meet the criteria, a random seed of 42 was used to determine the selection.

Analysis The experimental results are shown in Figure 5, where the red line represents the accuracy curve for *Related* combinations, and the purple line shows the gain from *Unrelated* combinations. The *Related* combinations group significantly outperformed the *Unrelated* combinations in terms of generalization across all tasks, with this ability continuing to improve as the data size increased. This suggests that *Related* combinations, leveraging CG, enhance the model's ability to understand unknown medical tasks.

Take-away 7: *CG might enable MLLMs to handle tasks without dedicated training data.*



Figure 5: The accuracy curve reflects the impact of gradually increasing the composition dataset size without using *Target* data in training. The green and red lines represent training with **Related** and **Unrelated Data**, respectively.



Figure 6: The accuracy curve shows the impact of increasing the composition dataset volume while incorporating *Target* data in training. The green and red lines represent training with **Related** and **Unrelated Data**, respectively.

5.2 Generalization with Limited Target Data

This section investigates the benefit of CG for tasks with limited data, e.g. processing medical images in rare conditions.

Selection Strategy To assess generalization in limited data scenarios, we select combinations with poor generalization from Table 3. Specifically, for each *Direction Type*, we randomly choose a CG combination with weak generalization (i.e., rows marked with \times in the last column of Table 3). For these combinations, we introduce an additional 2,000 examples from the *Target* data.

Analysis Figure 6 shows the results. It can be seen that as we gradually expand the training volume of *Target* data, adding the *Related* combination for training enabled the model to reach the peak performance more quickly. This suggests that leveraging CG to assist low-data medical scenarios can lead to more data-efficient training, even when CG does not directly result in significant generalization gains in these scenarios.

Take-away 8: Although CG might not provide direct generalization gains, it helps data efficiency for MLLM training.

6 CG across Detection and Classification

Previous studies (Ren et al., 2024; Wang et al., 2025) have shown that jointly training classification and detection tasks can mutually enhance their performance. Building on this, we investigate whether MLLMs can leverage classification data (e.g., visual knowledge) and detection data (e.g., spatial

information) through CG to improve downstream classification (*Q1*) or detection tasks (*Q2*).

6.1 Experiment Settings

Training Setup Each generalization combination used for training in this experiment includes one detection dataset and one classification dataset to examine the generalization relationship between these two vision tasks. The detailed training parameters can be found in Appendix A.6.

Model Selection Next-Chat (Zhang et al., 2023a) and MiniGPT-v2 (Chen et al., 2023a) are selected as baselines, representing the two main approaches MLLMs use for detection tasks. The former treats bounding boxes as embeddings and decodes them into coordinates using a visual decoder, while the latter processes coordinate points as special text tokens and generates bounding box coordinates directly as output text.

Data Processing Med-MAT includes both detection and segmentation datasets. If a segmentation dataset provides object localization using masks, we extract the outermost coordinates of the corresponding mask to construct a bounding box, facilitating generalization experiments for detection. Subsequently, to streamline the experiments, we structured the dataset following the official data formats of Next-Chat and MiniGPT-v2.

6.2 Benefits for Classification (Q1)

In this experiment, all possible CG combinations were selected and the CG-trained model will be



Figure 7: The accuracy(%) on Classification: Blue represents the untrained model, and green represents the CG-trained model. (details in Appendix A.5)

tested on classification task. The final results in Figure 7 show that all CG combinations demonstrated the model's successful utilization of detection data for CG to the *Target* data.

6.3 Benefits for Detection (Q2)

Subset 38 and 39 are selected as the objects in these datasets are relatively randomly distributed in the images, making them suitable for evaluating the model's detection capability. Subsequently, we selected certain classification datasets to construct CG for testing and used cIoU to evaluate the detection performance (follow (Chen et al., 2023a)).

Since both baselines lack localization capabilities for medical tasks, we incorporated a fixed amount of *Target* data into our experiments, adjusting the evaluation scenario to assess support in low-data settings. The results in Table 5 show that all selected CG combinations help MLLMs achieve better performance in detection tasks.

Related Combination	Target Subset	Next-Chat	MiniGPT-v2
D - Skin C - Intestin	e D - Intestine	+3.8	+4.1
D - Intestine C - Skin	D - Skin	+8.4	+7.6

Table 5: *Next-Chat* and *MiniGPT-v2* respectively represent the cIoU gain brought by CG. *C* indicates classification task, *D* indicates detection task.

Take-away 9: *MLLMs can perform CG across classification and detection tasks*.

7 Related Work

Medical MLLMs Recently, adapting MLLMs to medical tasks has gained prominence due to their success in capturing complex visual features. Current MLLMs typically pair a visual encoder with a text-only LLM, aligning image data with language understanding. Such as Med-Flamingo (Moor et al., 2023) and Med-PaLM (Tu et al., 2024), fine-tuned general multimodal models and achieved notable results. Med-Flamingo enhanced OpenFlamingo-9B (Chen et al., 2024a) with medical data, while Med-PaLM adapted PaLM-E (Driess et al., 2023) using 1 million data points. Similarly, LLaVA-Med (Li et al., 2024), Med-Gemini (Saab et al., 2024), and HuatuoGPT-Vision (Chen et al., 2024b) utilized specialized datasets and instruction tuning to refine medical VQA tasks.

Generalization on Medical Imaging Generalization in medical imaging (Matta et al., 2024) has been extensively studied. Early methods utilized data manipulation techniques, such as data augmentation (Li et al., 2022; Zhang et al., 2022), to enhance model generalization on unseen medical data by adapting to varying distributions. Later approaches focused on representation learning (Le-Khac et al., 2020), preserving essential image information to enable models to handle more complex scenarios. Additionally, some studies (Ren et al., 2024) explore multiple aspects of medical image processing, examining how classification and segmentation tasks can mutually benefit each other.

Detection with MLLMs Recent studies employ various strategies to equip MLLMs with the capability to handle detection tasks, such as encoding regions as features to allow models to accept regions as input (Zhang et al., 2023b), representing object bounding box coordinates with text tokens (Wang et al., 2024c; Peng et al., 2023; Chen et al., 2023b), and employing unique identifiers for task instructions to improve learning efficiency. Additionally, some approaches introduce special tokens to represent images and use their hidden states to decode position information (Zhang et al., 2023a, 2024a).

8 Conclusion

To investigate whether MLLMs can leverage CG to generalize to unseen medical data, we constructed the Med-MAT dataset as a research platform for generalization experiments. The results confirmed the presence of CG and identified it as a key factor of MLLMs' generalization observed in multitask learning. Further experiments showed that CG helps MLLMs handle limited data conditions, providing support for low-data medical tasks. Additionally, our findings showed that MLLMs can apply CG across detection and classification tasks, underscoring its broad generalization potential.

Limitations

The experiment confirms that MLLMs leverage CG for unseen medical images and data-efficient training. However, as shown in Section 4, disrupting CG reduces generalization but retains some effectiveness, indicating CG is just one aspect of MLLM generalization in medical imaging.

Potential Risks

Our research focuses on the compositional generalization of MLLMs on medical images, using data sourced from medical challenges and open-source datasets. However, further experiments are needed to mitigate potential risks when deploying this concept in real-world medical settings.

Acknowledgments

This work was supported Shenzhen by Medical Research Fund (No.C2406002) from the Shenzhen Medical Academy of Research and Translation (SMART), the Shenzhen Science and Technology Program (JCYJ20220818103001002), Shenzhen Doctoral Startup Funding (RCBS20221008093330065), Tianyuan Fund for Mathematics of National Natural Science Foundation of China (NSFC) (12326608), Shenzhen Science and Technology Program (Shenzhen Key Laboratory Grant No. ZDSYS20230626091302006), Shenzhen Stability Science Program 2023, and National Natural Science Foundation of China (NSFC) (72495131).

References

- Andrea Acevedo, Anna Merino, Santiago Alférez, Ángel Molina, Laura Boldú, and José Rodellar. 2020. A dataset of microscopic peripheral blood cell images for development of automatic recognition systems. *Data in brief*, 30:105474.
- Walid Al-Dhabyani, Mohammed Gomaa, Hussien Khaled, and Aly Fahmy. 2020. Dataset of breast ultrasound images. *Data in brief*, 28:104863.
- Shams Nafisa Ali, Md. Tazuddin Ahmed, Joydip Paul, Tasnim Jahan, S. M. Sakeef Sani, Nawshaba Noor, and Taufiq Hasan. 2022. Monkeypox skin lesion detection using deep learning models: A preliminary feasibility study. *arXiv preprint arXiv:2207.03342*.
- Sharib Ali, Barbara Braden, Dominique Lamarque, Stefano Realdon, Adam Bailey, Renato Cannizzaro, Noha Ghatwary, Jens Rittscher, Christian Daul, and James East. 2020. Endoscopy disease detection and segmentation (edd2020).

- MD Anouk Stein, Carol Wu, Chris Carr, George Shih, Jamie Dulkowski, kalpathy, Leon Chen, Luciano Prevedello, MD Marc Kohli, Mark Mc-Donald, Peter, Phil Culliton, Safwan Halabi MD, and Tian Xia. 2018. Rsna pneumonia detection challenge. https://kaggle.com/competitions/ rsna-pneumonia-detection-challenge. Kaggle.
- Will Arevalo. 2020. Chexpert v1.0 small. https: //www.kaggle.com/datasets/willarevalo/ chexpert-v10-small. Kaggle.
- A Asraf and Z Islam. 2021. Covid19, pneumonia and normal chest x-ray pa dataset. mendeley data v1 (2021).
- Francisco José Fumero Batista, Tinguaro Diaz-Aleman, Jose Sigut, Silvia Alayon, Rafael Arnay, and Denisse Angel-Pereira. 2020. Rim-one dl: A unified retinal image database for assessing glaucoma using deep learning. *Image Analysis & Stereology*, 39(3):161– 167.
- Dev Batra. 2024. Fracture detection using x-ray images. https://www. kaggle.com/datasets/devbatrax/ fracture-detection-using-x-ray-images. Kaggle.
- Veronica Elisa Castillo Benítez, Ingrid Castro Matto, Julio César Mello Román, José Luis Vázquez Noguera, Miguel García-Torres, Jordan Ayala, Diego P Pinto-Roa, Pedro E Gardel-Sotomayor, Jacques Facon, and Sebastian Alberto Grillo. 2021. Dataset from fundus images for the study of diabetic retinopathy. *Data in brief*, 36:107068.
- BenO, jljones, Kumar H, Meg Risdal, MRao, Vadim Sherman, Vipul, Wendy Kan, and Yau Ben-Or. 2017. Intel & mobileodt cervical cancer screening. https://kaggle.com/competitions/ intel-mobileodt-cervical-cancer-screening. Kaggle.
- Jorge Bernal, F Javier Sánchez, Gloria Fernández-Esparrach, Debora Gil, Cristina Rodríguez, and Fernando Vilariño. 2015. Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians. *Computerized medical imaging and graphics*, 43:99–111.
- Bukun. 2019. Breast cancer histopathological database (breakhis). https://www.kaggle.com/datasets/ ambarish/breakhis. Kaggle.
- Olivia Cardozo, Verena Ojeda, Rodrigo Parra, Julio César Mello-Román, José Luis Vázquez Noguera, Miguel García-Torres, Federico Divina, Sebastian A. Grillo, Cynthia Villalba, Jacques Facon, Veronica Elisa Castillo Benítez, Ingrid Castro Matto, and Diego Aquino-Brítez. 2023. Dataset of fundus images for the diagnosis of ocular toxoplasmosis. *Data in Brief*, 48:109056.

- Ling-Ping Cen, Jie Ji, Jian-Wei Lin, Si-Tong Ju, Hong-Jie Lin, Tai-Ping Li, Yun Wang, Jian-Feng Yang, Yu-Fen Liu, Shaoying Tan, et al. 2021. Automatic detection of 39 fundus diseases and conditions in retinal photographs using deep neural networks. *Nature communications*, 12(1):4828.
- Delong Chen, Jianfeng Liu, Wenliang Dai, and Baoyuan Wang. 2024a. Visual instruction tuning with polite flamingo. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17745–17753.
- Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. 2023a. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*.
- Junying Chen, Ruyi Ouyang, Anningzhe Gao, Shunian Chen, Guiming Hardy Chen, Xidong Wang, Ruifei Zhang, Zhenyang Cai, Ke Ji, Guangjun Yu, et al. 2024b. Huatuogpt-vision, towards injecting medical visual knowledge into multimodal llms at scale. arXiv preprint arXiv:2406.19280.
- Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. 2023b. Shikra: Unleashing multimodal llm's referential dialogue magic. arXiv preprint arXiv:2306.15195.
- Pingjun Chen. 2018. Knee osteoarthritis severity grading dataset. *Mendeley Data*, 1(10.17632).
- Muhammad E. H. Chowdhury, Tawsifur Rahman, Amith Khandakar, Rashid Mazhar, Muhammad Abdul Kadir, Zaid Bin Mahbub, Khandakar Reajul Islam, Muhammad Salman Khan, Atif Iqbal, Nasser Al Emadi, Mamun Bin Ibne Reaz, and Mohammad Tariqul Islam. 2020. Can ai help in screening viral and covid-19 pneumonia? *IEEE Access*, 8:132665–132676.
- Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. 2019. Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368*.
- Noel CF Codella, David Gutman, M Emre Celebi, Brian Helba, Michael A Marchetti, Stephen W Dusza, Aadi Kalloo, Konstantinos Liopyris, Nabin Mishra, Harald Kittler, et al. 2018. Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by the international skin imaging collaboration (isic). In 2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018), pages 168–172. IEEE.
- Marc Combalia, Noel CF Codella, Veronica Rotemberg, Brian Helba, Veronica Vilaplana, Ofer Reiter, Cristina Carrera, Alicia Barreiro, Allan C Halpern,

Susana Puig, et al. 2019. Bcn20000: Dermoscopic lesions in the wild. *arXiv preprint arXiv:1908.02288*.

- Will Cukierski. 2018. Histopathologic cancer detection. https://kaggle.com/competitions/ histopathologic-cancer-detection. Kaggle.
- Training Data. 2023. Computed tomography of the brain. https://www. kaggle.com/datasets/trainingdatapro/ computed-tomography-ct-of-the-brain. Kaggle.
- Coen de Vente, Koenraad A. Vermeer, Nicolas Jaccard, He Wang, Hongyi Sun, Firas Khader, Daniel Truhn, Temirgali Aimyshev, Yerkebulan Zhanibekuly, Tien-Dung Le, Adrian Galdran, Miguel Ángel González Ballester, Gustavo Carneiro, Devika R G, Hrishikesh P S, Densen Puthussery, Hong Liu, Zekang Yang, Satoshi Kondo, Satoshi Kasai, Edward Wang, Ashritha Durvasula, Jónathan Heras, Miguel Ángel Zapata, Teresa Araújo, Guilherme Aresta, Hrvoje Bogunović, Mustafa Arikan, Yeong Chan Lee, Hyun Bin Cho, Yoon Ho Choi, Abdul Qayyum, Imran Razzak, Bram van Ginneken, Hans G. Lemij, and Clara I. Sánchez. 2023. Airogs: Artificial intelligence for robust glaucoma screening challenge. *arXiv preprint arXiv:2302.01738*.
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. 2023. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*.
- Fernando Feltrin. 2022. Brain tumor mri images 17 classes. https://www. kaggle.com/datasets/fernando2rad/ brain-tumor-mri-images-17-classes. Kaggle.
- Mohammad Fraiwan, Ziad Audat, Luay Fraiwan, and Tarek Manasreh. 2022. Using deep transfer learning to detect scoliosis and spondylolisthesis from x-ray images. *Plos one*, 17(5):e0267851.
- Huazhu Fu, Fei Li, José Ignacio Orlando, Hrvoje Bogunović, Xu Sun, Jingan Liao, Yanwu Xu, Shaochong Zhang, and Xiulan Zhang. 2019. Palm: Pathologic myopia challenge.
- Ioannis Giotis, Nynke Molders, Sander Land, Michael Biehl, Marcel F Jonkman, and Nicolai Petkov. 2015. Med-node: A computer-assisted melanoma diagnosis system using non-dermoscopic images. *Expert systems with applications*, 42(19):6578–6585.
- Haifan Gong, Guanqi Chen, Ranran Wang, Xiang Xie, Mingzhi Mao, Yizhou Yu, Fei Chen, and Guanbin Li. 2021. Multi-task learning for thyroid nodule segmentation with thyroid region prior. In 2021 IEEE 18th international symposium on biomedical imaging (ISBI), pages 257–261. IEEE.
- Haifan Gong, Jiaxin Chen, Guanqi Chen, Haofeng Li, Fei Chen, and Guanbin Li. 2022. Thyroid region prior guided attention for ultrasound segmentation of

thyroid nodules. *Computers in Biology and Medicine*, 106389:1–12.

- Shivanand Gornale and Pooja Patravali. 2020. Digital knee x-ray images.
- Matthew Groh, Caleb Harris, Luis Soenksen, Felix Lau, Rachel Han, Aerin Kim, Arash Koochek, and Omar Badri. 2021. Evaluating deep neural networks trained on clinical images in dermatology with the fitzpatrick 17k dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1820–1828.
- David Gutman, Noel CF Codella, Emre Celebi, Brian Helba, Michael Marchetti, Nabin Mishra, and Allan Halpern. 2016. Skin lesion analysis toward melanoma detection: A challenge at the international symposium on biomedical imaging (isbi) 2016, hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1605.01397*.
- Saba Hesaraki. 2022. Breast ultrasound images dataset (busi). https://www. kaggle.com/datasets/sabahesaraki/ breast-ultrasound-images-dataset. Kaggle.
- Md Nazmul Islam, Mehedi Hasan, Md Kabir Hossain, Md Golam Rabiul Alam, Md Zia Uddin, and Ahmet Soylu. 2022a. Vision transformer and explainable transfer learning models for auto detection of kidney cyst, stone and tumor from ct-radiography. *Scientific Reports*, 12(1):1–14.
- Towhidul Islam, Mohammad Arafat Hussain, Forhad Uddin Hasan Chowdhury, and B M Riazul Islam. 2022b. A web-scrapped skin image database of monkeypox, chickenpox, smallpox, cowpox, and measles. *bioRxiv* 2022.08.01.502199.
- Stefan Jaeger, Sema Candemir, Sameer Antani, Yi-Xiáng J Wáng, Pu-Xuan Lu, and George Thoma. 2014. Two public chest x-ray datasets for computeraided screening of pulmonary diseases. *Quantitative imaging in medicine and surgery*, 4(6):475.
- Debesh Jha, Pia H Smedsrud, Michael A Riegler, Pål Halvorsen, Thomas de Lange, Dag Johansen, and Håvard D Johansen. 2020. Kvasir-seg: A segmented polyp dataset. In *MultiMedia Modeling: 26th International Conference, MMM 2020, Daejeon, South Korea, January 5–8, 2020, Proceedings, Part II 26*, pages 451–462. Springer.
- Kai Jin, Xingru Huang, Jingxing Zhou, Yunxiang Li, Yan Yan, Yibao Sun, Qianni Zhang, Yaqi Wang, and Juan Ye. 2022. Fives: A fundus image dataset for artificial intelligence based vessel segmentation. *Scientific data*, 9(1):475.
- JR2NGB. 2019. Cataract dataset. https: //www.kaggle.com/datasets/jr2ngb/ cataractdataset. Kaggle.

- Nur Karaca. 2022. Nlm montgomery cxr set. https://www.kaggle.com/datasets/ nurkaraca/nlm-montgomerycxrset. Kaggle.
- Karthik, Maggie, and Sohier Dane. 2019. Aptos 2019 blindness detection. https://kaggle.com/competitions/ aptos2019-blindness-detection. Kaggle.
- Andrey Katanskiy. 2019. Skin cancer isic. https://www.kaggle.com/datasets/ nodoubttome/skin-cancer9-classesisic. Kaggle.
- Jakob Nikolas Kather, Niels Halama, and Alexander Marx. 2018. 100,000 histological images of human colorectal cancer and healthy tissue.
- Daniel Kermany. 2018. Labeled optical coherence tomography (oct) and chest x-ray images for classification. *Mendeley data*.
- Felipe Campos Kitamura. 2018. Head ct hemorrhage.
- Jorge F Lazo, Benoit Rosa, Michele Catellani, Matteo Fontana, Francesco A Mistretta, Gennaro Musi, Ottavio de Cobelli, Michel de Mathelin, and Elena De Momi. 2023. Semi-supervised bladder tissue classification in multi-domain endoscopic images. *IEEE Transactions on Biomedical Engineering*.
- Trang Le, Casper F Winsnes, Ulrika Axelsson, Hao Xu, Jayasankar Mohanakrishnan Kaimal, Diana Mahdessian, Shubin Dai, Ilya S Makarov, Vladislav Ostankovich, Yang Xu, et al. 2022. Analysis of the human protein atlas weakly supervised singlecell classification competition. *Nature methods*, 19(10):1221–1229.
- Phuc H Le-Khac, Graham Healy, and Alan F Smeaton. 2020. Contrastive representation learning: A frame-work and review. *Ieee Access*, 8:193907–193934.
- Rebecca Sawyer Lee, Francisco Gimenez, Assaf Hoogi, Kanae Kawai Miyake, Mia Gorovoy, and Daniel L Rubin. 2017. A curated mammography data set for use in computer-aided detection and diagnosis research. *Scientific data*, 4(1):1–9.
- Sangjune L Lee, Poonam Yadav, Yin Li, Jason J Meudt, Jessica Strang, Dustin Hebel, Alyx Alfson, Stephanie J Olson, Tera R Kruser, Jennifer B Smilowitz, et al. 2024. Dataset for gastrointestinal tract segmentation on serial mris for abdominal tumor radiotherapy. *Data in Brief*, page 111159.
- Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2024. Llavamed: Training a large language-and-vision assistant for biomedicine in one day. *Advances in Neural Information Processing Systems*, 36.
- Yuanpeng Li, Liang Zhao, Jianyu Wang, and Joel Hestness. 2019. Compositional generalization for primitive substitutions. arXiv preprint arXiv:1910.02612.

- Yuexiang Li, Nanjun He, and Yawen Huang. 2022. Single domain generalization via spontaneous amplitude spectrum diversification. In *MICCAI Workshop on Resource-Efficient Medical Image Analysis*, pages 32–41. Springer.
- Jie Lian, Jingyu Liu, Shu Zhang, Kai Gao, Xiaoqing Liu, Dingwen Zhang, and Yizhou Yu. 2021. A structureaware relation network for thoracic diseases detection and segmentation. *IEEE Transactions on Medical Imaging*, 40(8):2042–2052.
- Xiao Liang. 2021. Adam dataset. https: //www.kaggle.com/datasets/xiaoliang2121/ adamdataset. Kaggle.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning.
- Jacob A Macdonald, Zhe Zhu, Brandon Konkel, and Mazurowski. 2020. Siim-acr pneumothorax segmentation. https://doi.org/10.5281/zenodo. 7774566. Zenodo.
- K Scott Mader. 2017. Mias mammography. https://www.kaggle.com/datasets/kmader/ mias-mammography. Kaggle.
- Salman Maqbool, Aqsa Riaz, Hasan Sajid, and Osman Hasan. 2020. m2caiseg: Semantic segmentation of laparoscopic images using convolutional neural networks. *arXiv preprint arXiv:2008.10134*.
- Christian Matek, Sebastian Krappe, Christian Münzenmayer, Torsten Haferlach, and Carsten Marr. 2021. An expert-annotated dataset of bone marrow cytology in hematologic malignancies. *The Cancer Imaging Archive*.
- Sarah Matta, Mathieu Lamard, Philippe Zhang, Alexandre Le Guilcher, Laurent Borderie, Béatrice Cochener, and Gwenolé Quellec. 2024. A systematic review of generalization research in medical image classification. *arXiv preprint arXiv:2403.12167*.
- Teresa Mendonca, M Celebi, T Mendonca, and J Marques. 2015. Ph2: A public database for the analysis of dermoscopic images. *Dermoscopy image analysis*.
- Meta AI. 2024. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models. https://ai.meta.com/blog/ llama-3-2-connect-2024-vision-edge\ -mobile-devices/.
- Shentong Mo and Paul Pu Liang. 2024. Multimed: Massively multimodal and multitask medical understanding. *arXiv preprint arXiv:2408.12682*.
- Paul Mooney. 2017. Blood cell images. https://www.kaggle.com/datasets/ paultimothymooney/blood-cells. Kaggle.
- Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Yash Dalmia, Jure Leskovec, Cyril Zakka, Eduardo Pontes Reis, and Pranav Rajpurkar.

2023. Med-flamingo: a multimodal medical few-shot learner. In *Machine Learning for Health (ML4H)*, pages 353–367. PMLR.

- Loris Nanni, Michelangelo Paci, Florentino Luciano Caetano dos Santos, Heli Skottman, Kati Juuti-Uusitalo, and Jari Hyttinen. 2016. Texture descriptors ensembles enable image-based classification of maturation of human stem cell-derived retinal pigmented epithelium. *PLoS One*, 11(2):e0149399.
- Hieu T Nguyen, Ha Q Nguyen, Hieu H Pham, Khanh Lam, Linh T Le, Minh Dao, and Van Vu. 2023. Vindr-mammo: A large-scale benchmark dataset for computer-aided diagnosis in full-field digital mammography. *Scientific Data*, 10(1):277.
- Masoud Nickparvar. 2021a. Brain tumor mri dataset. https://www.kaggle.com/datasets/ masoudnickparvar/brain-tumor-mri-dataset. Kaggle.

Msoud Nickparvar. 2021b. Brain tumor mri dataset.

- Nikita Orlov, Wayne Chen, David Eckley, Tomasz Macura, Lior Shamir, Elaine Jaffe, and Ilya Goldberg. 2010a. Automatic classification of lymphoma images with transform-based global features. *IEEE transactions on information technology in biomedicine : a publication of the IEEE Engineering in Medicine and Biology Society*, 14:1003–13.
- Nikita Orlov, Wayne Chen, David Eckley, Tomasz Macura, Lior Shamir, Elaine Jaffe, and Ilya Goldberg. 2010b. Automatic classification of lymphoma images with transform-based global features. *IEEE transactions on information technology in biomedicine : a publication of the IEEE Engineering in Medicine and Biology Society*, 14:1003–13.
- Silvia Ovreiu, Elena-Anca Paraschiv, and Elena Ovreiu. 2021. Deep learning & digital fundus images: Glaucoma detection using densenet. In 2021 13th international conference on electronics, computers and artificial intelligence (ECAI), pages 1–4. IEEE.
- Andre GC Pacheco, Gustavo R Lima, Amanda S Salomao, Breno Krohling, Igor P Biral, Gabriel G de Angelo, Fábio CR Alves Jr, José GM Esgario, Alana C Simora, Pedro BC Castro, et al. 2020. Pad-ufes-20: A skin lesion dataset composed of patient data and clinical images collected from smartphones. *Data in brief*, 32:106221.
- Sachin Panchal, Ankita Naik, Manesh Kokare, Samiksha Pachade, Rushikesh Naigaonkar, Prerana Phadnis, and Archana Bhange. 2023. Retinal fundus multi-disease image dataset (rfmid) 2.0: a dataset of frequently and rarely identified diseases. *Data*, 8(2):29.
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*.

- H Hieu Pham, T Thanh Tran, and Ha Quy Nguyen. 2022. Vindr-pcxr: An open, large-scale pediatric chest xray dataset for interpretation of common thoracic diseases. *PhysioNet* (*version 1.0. 0*), 10:2.
- Hieu Huy Pham, H Nguyen Trung, and Ha Quy Nguyen. 2021. Vindr-spinexr: A large annotated medical image dataset for spinal lesions detection and classification from radiographs. *PhysioNet*.
- Konstantin Pogorelov, Kristin Ranheim Randel, Thomas de Lange, Sigrun Losada Eskeland, Carsten Griwodz, Dag Johansen, Concetto Spampinato, Mario Taschwer, Mathias Lux, Peter Thelin Schmidt, Michael Riegler, and Pål Halvorsen. 2017a. Nerthus: A bowel preparation quality video dataset. In Proceedings of the 8th ACM on Multimedia Systems Conference, MMSys'17, pages 170–174, New York, NY, USA. ACM.
- Konstantin Pogorelov, Kristin Ranheim Randel, Carsten Griwodz, Sigrun Losada Eskeland, Thomas de Lange, Dag Johansen, Concetto Spampinato, Duc-Tien Dang-Nguyen, Mathias Lux, Peter Thelin Schmidt, Michael Riegler, and Pål Halvorsen. 2017b. Kvasir: A multi-class image dataset for computer aided gastrointestinal disease detection. In Proceedings of the 8th ACM on Multimedia Systems Conference, MM-Sys'17, pages 164–169, New York, NY, USA. ACM.
- Praveen. 2019. Coronahack chest x-ray dataset. https://www.kaggle.com/datasets/ praveengovi/coronahack-chest-xraydataset. Kaggle.
- Pavle Prentasic, Sven Loncaric, Zoran Vatavuk, Goran Bencic, Marko Subasic, Tomislav Petković, Lana Dujmovic, Maja Malenica Ravlic, Nikolina Budimlija, and Rašeljka Tadić. 2013. Diabetic retinopathy image database(dridb): A new database for diabetic retinopathy screening programs research. In *International Symposium on Image and Signal Processing and Analysis, ISPA*, pages 711–716.
- Xianbiao Qi, Guoying Zhao, Jie Chen, and Matti Pietikäinen. 2016. Hep-2 cell classification: The role of gaussian scale space theory as a pre-processing approach. *Pattern Recognition Letters*, 82:36–43.
- Raddar. 2019. Chest x-rays (indiana university). https://www.kaggle.com/datasets/raddar/ chest-xrays-indiana-university?select= indiana_reports.csv. Kaggle.
- Tawsifur Rahman, Amith Khandakar, Muhammad Abdul Kadir, Khandaker Rejaul Islam, Khandakar F Islam, Rashid Mazhar, Tahir Hamid, Mohammad Tariqul Islam, Saad Kashem, Zaid Bin Mahbub, et al. 2020. Reliable tuberculosis detection using chest x-ray with deep learning, segmentation and visualization. *Ieee Access*, 8:191586–191601.
- MOHD ZAID RASHID. 2024. Oral cancer dataset. https://www.kaggle.com/datasets/ zaidpy/oral-cancer-dataset. Kaggle.

- Sucheng Ren, Xiaoke Huang, Xianhang Li, Junfei Xiao, Jieru Mei, Zeyu Wang, Alan Yuille, and Yuyin Zhou. 2024. Medical vision generalist: Unifying medical imaging tasks in context. *arXiv preprint arXiv:2406.05565*.
- Manuel Alejandro Rodríguez, Hasan AlMarzouqi, and Panos Liatsis. 2022. Multi-label retinal disease classification using transformers. *IEEE Journal of Biomedical and Health Informatics*.
- Veronica Rotemberg, Nicholas Kurtansky, Brigid Betz-Stablein, Liam Caffery, Emmanouil Chousakos, Noel Codella, Marc Combalia, Stephen Dusza, Pascale Guitera, David Gutman, et al. 2021. A patientcentric dataset of images and metadata for identifying melanomas using clinical context. *Scientific data*, 8(1):34.
- Khaled Saab, Tao Tu, Wei-Hung Weng, Ryutaro Tanno, David Stutz, Ellery Wulczyn, Fan Zhang, Tim Strother, Chunjong Park, Elahe Vedadi, et al. 2024. Capabilities of gemini models in medicine. *arXiv preprint arXiv:2404.18416*.
- Salman Sajid. 2024. Oral diseases. https: //www.kaggle.com/datasets/salmansajid05/ oral-diseases/data. Kaggle.
- F Shaker. 2018. Human sperm head morphology dataset (hushem). *Mendeley Data*, 3.
- Julio Silva-Rodríguez, Adrián Colomer, María A Sales, Rafael Molina, and Valery Naranjo. 2020. Going deeper through the gleason scoring scale: An automatic end-to-end system for histology prostate grading and cribriform pattern detection. *Computer Methods and Programs in Biomedicine*, 195:105637.
- Eduardo Soares, Plamen Angelov, Sarah Biaso, Michele Higa Froes, and Daniel Kanda Abe. 2020. Sars-cov-2 ct-scan dataset:a large dataset of real patients ct scans for sars-cov-2 identification. *Cold Spring Harbor Laboratory Press.*
- Malliga Subramanian, Kogilavani Shanmugavadivel, Obuli Sai Naren, K Premkumar, and K Rankish. 2022. Classification of retinal oct images using deep learning. In 2022 International Conference on Computer Communication and Informatics (ICCCI), pages 1–7.
- Summers and Ronald. 2020. Chestxray nihcc. https://nihcc.app.box.com/v/ ChestXray-NIHCC/folder/36938765345. NIH.

SunneYi. 2021. Chest CT-Scan images Dataset.

Siham Tabik, Anabel Gómez-Ríos, José Luis Martín-Rodríguez, Iván Sevillano-García, Manuel Rey-Area, David Charte, Emilio Guirado, Juan-Luis Suárez, Julián Luengo, MA Valero-González, et al. 2020. Covidgr dataset and covid-sdnet methodology for predicting covid-19 based on chest x-ray images. *IEEE journal of biomedical and health informatics*, 24(12):3595–3605.

- Yihong Tang, Ao Qu, Zhaokai Wang, Dingyi Zhuang, Zhaofeng Wu, Wei Ma, Shenhao Wang, Yunhan Zheng, Zhan Zhao, and Jinhua Zhao. 2024. Sparkle: Mastering basic spatial capabilities in vision language models elicits generalization to composite spatial reasoning. arXiv preprint arXiv:2410.16162.
- Tao Tu, Shekoofeh Azizi, Danny Driess, Mike Schaekermann, Mohamed Amin, Pi-Chuan Chang, Andrew Carroll, Charles Lau, Ryutaro Tanno, Ira Ktena, et al. 2024. Towards generalist biomedical ai. *NEJM AI*, 1(3):AIoa2300138.
- Peking University. 2019. Odir-2019 dataset. https://odir2019.grand-challenge.org/ introduction/. Grand Challenge.
- Preet Viradiya. 2020. Brain tumor dataset. https://www.kaggle.com/datasets/ preetviradiya/brain-tumor-dataset. Kaggle.
- Haiyang Wang, Hao Tang, Li Jiang, Shaoshuai Shi, Muhammad Ferjad Naeem, Hongsheng Li, Bernt Schiele, and Liwei Wang. 2025. Git: Towards generalist vision transformer through universal language interface. In *European Conference on Computer Vi*sion, pages 55–73. Springer.
- Linda Wang, Zhong Qiu Lin, and Alexander Wong. 2020. Covid-net: a tailored deep convolutional neural network design for detection of covid-19 cases from chest x-ray images. *Scientific Reports*, 10(1):19549.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024b. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *Preprint*, arXiv:2409.12191.
- Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. 2024c. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *Advances in Neural Information Processing Systems*, 36.
- Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. 2017. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2097–2106.

wjXiaochuangw. 2019. Covid-19-ct scan images.

- Zhenlin Xu, Marc Niethammer, and Colin A Raffel. 2022. Compositional generalization in unsupervised compositional representation learning: A study on disentanglement and emergent language. Advances in Neural Information Processing Systems, 35:25074– 25087.
- Anna Zawacki, Carol Wu, George Shih, Julia Elliott, Mikhail Fomitchev, Mohannad Hussain, Paras Lakhani, Phil Culliton, and Shunxing Bao. 2019. Siim-acr pneumothorax segmentation. https://kaggle.com/competitions/ siim-acr-pneumothorax-segmentation. Kaggle.
- Yaya Zha. 2021. Rus-chn. https://aistudio.baidu. com/datasetdetail/69582/0. AI Studio.
- Ao Zhang, Liming Zhao, Chen-Wei Xie, Yun Zheng, Wei Ji, and Tat-Seng Chua. 2023a. Next-chat: An lmm for chat, detection and segmentation. *arXiv preprint arXiv:2311.04498*.
- Edward Zhang and Sauman Das. 2022. Glaucoma detection. https://www.kaggle.com/datasets/sshikamaru/glaucoma-detection. Kaggle.
- Ruipeng Zhang, Qinwei Xu, Chaoqin Huang, Ya Zhang, and Yanfeng Wang. 2022. Semi-supervised domain generalization for medical image analysis. In 2022 IEEE 19th International Symposium on Biomedical Imaging (ISBI), pages 1–5. IEEE.
- Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Yu Liu, Kai Chen, and Ping Luo. 2023b. Gpt4roi: Instruction tuning large language model on region-of-interest. *arXiv preprint arXiv:2307.03601*.
- Tao Zhang, Xiangtai Li, Hao Fei, Haobo Yuan, Shengqiong Wu, Shunping Ji, Change Loy Chen, and Shuicheng Yan. 2024a. Omg-llava: Bridging image-level, object-level, pixel-level reasoning and understanding. In *NeurIPS*.
- Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2023c. Pmc-vqa: Visual instruction tuning for medical visual question answering. *arXiv preprint arXiv:2305.10415*.
- Yuhui Zhang, Alyssa Unell, Xiaohan Wang, Dhruba Ghosh, Yuchang Su, Ludwig Schmidt, and Serena Yeung-Levy. 2024b. Why are visually-grounded language models bad at image classification? *arXiv preprint arXiv:2405.18415*.
- Jinyu Zhao, Yichen Zhang, Xuehai He, and Pengtao Xie. 2020. Covid-ct-dataset: a ct scan dataset about covid-19. *arXiv preprint arXiv:2003.13865*.
- Chuang Zhu, Wenkai Chen, Ting Peng, Ying Wang, and Mulan Jin. 2021a. Hard sample aware noise robust learning for histopathology image classification.

IEEE transactions on medical imaging, 41(4):881–894.

Chuang Zhu, Wenkai Chen, Ting Peng, Ying Wang, and Mulan Jin. 2021b. Hard sample aware noise robust learning for histopathology image classification. *IEEE transactions on medical imaging*, 41(4):881– 894.

Xile Zhu. 2022. Lc25000. https://www.kaggle. com/datasets/xilezhu/lc25000. Kaggle.

. 2021. Augemnted ocular diseases. https://www.kaggle.com/datasets/ nurmukhammed7/augemnted-ocular-diseases. Kaggle.

A More Experiments

A.1 Benefits for Segmentation

Segmentation-enabled LLMs, such as Next-GPT, first use the LLM to identify potential regions of the target object and then apply a SAM to decode the object mask, thereby completing the segmentation task. In this context, segmentation can be seen as an extension of detection, potentially requiring more images to achieve improved performance. We conducted additional experiments to explore whether MLLMs can still utilize CG to understand new images across both segmentation and classification tasks.

Related Co	ombination	Target Subset	Next-Chat
D - Skin D - Intestine	C - Intestine C - Skin	S - Intestine S - Skin	+7.46 + 5.42

Table 6: *Next-Chat* represents the cIoU gain brought by CG. *C* indicates classification task, *S* indicates Segmentation task.

The results in Table 6 demonstrate that, in the context of segmentation tasks, MLLMs are still able to leverage CG to understand new tasks, which is consistent with our original conclusions.

A.2 More Complex Medical Elements

While MAT-Triplet Categorization is useful, predefined categories may limit the exploration of more complex medical attributes, so we also considered integrating more flexible categorization to explore additional medical attributes.

Additional Element 1: Population Groups We selected VinDr-PCXR and MedMAT Subset 31 for the experiment, as they contain X-ray images of children and adult groups, respectively. The results are shown in Table 7.

Additional Element 2: Finer Disease "Finer disease" means more detailed categorization. For instance, we treat COVID and common pneumonia as distinct diseases for generalization. We split the Normal data in the training set into two parts and combined each with COVID and Pneumonia data to create new datasets. The results are shown in Table 8.

Related Combination	Target Subset	LLaVA
A X-ray ⑦ Young ⑦ Unrelated Data A X-ray ⑦ Young A CT Children (CG)	Adults	+6.04 + 18.12

Table 7: Results of using *Population Groups* as a CG element.

Related Combination	Target Subset	LLaVA
 ▲ X-ray [®] Pneumonia [®] Unrelated Data ▲ X-ray [®] Pneumonia [≜] CT [®] COVID (CG) 	A X-ray ⑦ COVII A X-ray ⑦ COVII	$ \begin{vmatrix} +11.33 \\ +12.67 \end{vmatrix} $

Table 8: Results of using Finer Disease as a CG element.

A.3 Statistical Tests of the Generalization Results

To ensure consistency and repeatability of the experiment, we performed statistical tests in this section. LLaVA is selected as the baseline, and we used the same data combinations from Section 3.3. Each experiment was repeated 3 times, and we reported the mean and standard deviation (SD) of the results.

From the results in Table 9, we can observe that the outcomes across runs show low variance, indicating overall stability, and they continue to support our original experimental conclusions.

A.4 CG with All MAT-Triplet Elements from Different Sources

In previous controlled experiments (Section 3), one element of the MAT-Triplet was kept constant while CG was explored in the remaining two elements. To ensure that all the 3 MAT-Triplet elements of the target data originated from three distinct datasets, additional experiments were conducted to further validate the effectiveness of CG. For these experiments, all possible combinations meeting the criteria in Med-MAT were selected (**Selection Strategy**). The results presented in Table 10 demonstrate that most combinations can effectively generalize to the *Target* data.

Analysis of the results The results in Table 7 and 8 indicate that the two new attributes show data leakage due to subtle visual differences in corresponding images (e.g., COVID-19 and pneumonia

	Related Combination			Target S	Baseline	1st	2nd	3rd	Mean and SD	
🖗 Bones	🖹 State	🖗 Breast	🖹 Diseases	🖗 Bones	B Diseases	37.31	43.28	44.78	43.28	43.78 ± 0.87
🖗 Lung	🖹 COVID	🖗 Bones	🖹 Diseases	🖗 Lung	🖹 Diseases	49.00	52.00	52.00	52.00	52.00 ± 0.00
Aray	🖹 Diseases	≙ CT	🖹 COVID	▲ X-ray	🖹 COVID	30.00	47.33	49.33	49.33	48.66 ± 1.15
Aray	🖹 Diseases	≜ CT	🖹 State	▲ X-ray	🖹 State	30.00	46.00	45.33	44.67	45.33 ± 0.67
A CT	🖗 Brain	🛓 X-ray	🖗 Lung	≜ CT	🖗 Lung	25.00	31.50	32.00	32.00	31.83 ± 0.29
🖉 CT	🖗 Brain		🖗 Lung	🖆 CT	🖗 Lung	47.00	71.00	71.00	70.00	70.67 ± 0.58
🔊 FP - Fundus	🖹 Diseases	🔊 Mic - Cell	🖹 Level	🔊 FP - Fundus	🖹 Level	33.33	42.42	45.45	45.45	44.44 ± 1.75
🖓 Mic - Cell	Recognition	🖾 FP - Fundus	🖹 Level	🖓 Mic - Cell	🖹 Level	23.00	32.00	32.00	31.50	31.83 ± 0.29

Table 9: Statistical tests of CG experiments. The 1st, 2nd, and 3rd show the generalization results of the experiment in different runs. "Mean" and "SD" represent the average accuracy (%) and standard deviation.

Rela	ted Coml	oination	1	Farget Su	bset	Baseline	Trained	CG Helps
≜ CT	🖗 Brain	🖹 Cancer	≜ CT	🖗 Brain	🖹 Cancer	28	26	X
🖞 CT	🖗 Brain	🖹 Cancer	🖞 CT	🖗 Brain	🖹 Cancer	28	25	×
🖞 CT	🖗 Brain	🖹 State	🖞 CT	🖗 Brain	🖹 State	33	64	1
🖞 CT	🖗 Brain	🖹 State	🖞 CT	🖗 Brain	🖹 State	33	70	1
∄ X-ray	🖗 Lung	Diseases	∄ X-ray	🖗 Lung	Diseases	30	45	1
≜ X-ray	🖗 Lung	Diseases	≜ X-ray	🖗 Lung	Diseases	30	38	1
≜ X-ray	🖗 Lung	Diseases	≜ X-ray	🖗 Lung	Diseases	30	44	1
≜ X-ray	⅔ Breast	Diseases	≜ X-ray	⅔ Breast	Diseases	31	32	1
≜ X-ray	∯ Breast	🖹 Diseases	∄ X-ray	∯ Breast	Diseases	31	52	1

Table 10: Results from 3 datasets providing different elements of MAT-Triplet. \checkmark in *CG Helps* indicates successful generalization, while \checkmark denotes failure.

have similar features). Importantly, the MLLM trained with CG combinations still shows improvements on downstream tasks, confirming that our approach remains valid for new attributes.

Reason to choose the existing three attributes (MAT-Triplet: Modality, Area, Task) We have considered additional categories such as age, gender, and finer disease classification, but we ultimately chose to focus on the MAT-Triplet categories for the following reasons.

- The boundaries between MAT-Triplet (Modality, Area, Task) are clear. Different modalities and areas correspond to distinct imaging methods and body areas, leading to significant differences between images; different tasks also require the MLLM to extract specific information, demanding varied understanding of the images.
- All datasets can be annotated using MAT-Triplet (Modality, Area, Task) easily. Other medical labels, such as gender and age, are only available in a small portion of datasets and are not suitable for large-scale annotation.
- Similar categorization strategies have been adopted in previous studies.

A.5 Details of Section 3.3: Exploring CG on different MLLM Backbones

To ensure the experiment results are not influenced by the model choice, we also tested several other models on some subsets of Med-MAT and observed similar results.

Selection Strategy: For testing, some generalized combinations were selected from classification tasks 3. Using a random seed of 42, we shuffled each Direction Type's combinations and selected the first two compositions as test data.

Experimental Setup: We conducted experiments to evaluate the compatibility of CG across different backbone architectures. We selected two MLLMs with representative architectures, namely Qwen2-VL-7B-Instruct (Wang et al., 2024b) and Llama-3.2-11B-Vision-Instruct (Meta AI, 2024), to assess the performance of CG on these models. Each experiment involved full-parameter fine-tuning of all models over 5 epochs, utilizing 8 A800 (80GB) GPUs. The training was performed with a batch size of 32 and a learning rate set to 2e-6, ensuring that all parameters were updated to optimize the model performance.

A.6 Details of Section 6: Exploring CG across Detection and Classification

Experimental Setup: We conducted generalization experiments for detection and classification. Specifically, we performed generalization validation on Next-Chat (Zhang et al., 2023a) and MiniGPT-v2 (Chen et al., 2023a). Next-Chat models the bounding box as an embedding and utilizes a decoder for decoding, while MiniGPT-v2 treats the bounding box as a text token, which are common approaches used by existing MLLM implementations for detection. By conducting CG validation using distinct bounding box modeling methods, we further demonstrate the broad applicability of the CG approach. Each experiment was conducted on 8 A800 (80GB) GPUs.

The two backbones were trained separately in this experiment. For Next-Chat, we directly trained the model in its second training stage and finetuned it for 2 epochs with a learning rate of 2e-5,

	Related Con		Target S	ubset	Baseline	Trained	CG Helps	
🖗 Bones	🖹 State	🖗 Breast	🖹 Diseases	🖗 Bones	🖹 Diseases	61	65	1
🖗 Lung	🖹 COVID	🖗 Bones	🖹 Diseases	🖗 Lung	🖹 Diseases	80	91	1
🗳 X-ray	🖹 Diseases	≜ CT	🖹 COVID	🗳 X-ray	🖹 COVID	35	40	1
🛓 X-ray	🖹 Diseases	≜ CT	🖹 State	🖉 X-ray	🖹 State	35	43	✓
🛓 CT	🖗 Brain	🗳 X-ray	🖗 Lung	🖉 CT	🖗 Lung	32	33	1
🛓 CT	🖗 Brain	🗳 X-ray	🖗 Lung	🖉 CT	🖗 Lung	65	72	1
🔊 FP - Fundus	🖹 Diseases	🔊 Mic - Cell	🖹 Level	🔊 FP - Fundus	🖹 Level	48	45	X
최 Mic - Cell	Recognition	🛱 FP - Fundus	🖹 Level	🔊 Mic - Cell	🖹 Level	34	41	1

Table 11: Result of Qwen2-VL on selected classification datasets in Med-MAT. \checkmark in *CG Helps* indicates successful generalization, while \varkappa denotes failure.

	Related Con		Target S	Baseline	Trained	CG Helps		
🖗 Bones	🖹 State	🖗 Breast	🖹 Diseases	🖗 Bones	🖹 Diseases	52	59	1
🖗 Lung	🖹 COVID	🖗 Bones	🖹 Diseases	🖗 Lung	🖹 Diseases	64	75	1
🗳 X-ray	🖹 Diseases	≜ CT	🖹 COVID	🗳 X-ray	🖹 COVID	33	38	1
🗳 X-ray	🖹 Diseases	🗳 CT	🖹 State	🗳 X-ray	🖹 State	33	41	1
🛆 CT	🖗 Brain	🗳 X-ray	🖗 Lung	A CT	🖗 Lung	31	29	×
🗳 CT	🖗 Brain	🗳 X-ray	🖗 Lung	🗳 CT	🖗 Lung	49	57	1
🖓 FP - Fundus	🖹 Diseases	🖓 Mic - Cell	🖹 Level	🔊 FP - Fundus	🖹 Level	55	61	1
🖓 Mic - Cell	Recognition	🖓 FP - Fundus	🖹 Level	🔊 Mic - Cell	🖹 Level	10	32	1

Table 12: Result of Llama-3.2-Vision on selected classification datasets in Med-MAT. \checkmark in *CG Helps* indicates successful generalization, while \varkappa denotes failure.

Related Combination				Target Subset		Baseline	Trained	CG Helps
脊 Lung 脊 Lung 脊 Bones 脊 Bones	 Lung Det Lung Det Spinal Error Det Spinal Error Det 		 Diseases Diseases Diseases Diseases Diseases 	脊 Lung 脊 Lung 脊 Bones 脊 Bones	 Diseases Diseases Diseases Diseases Diseases 	49 49 20 20	52 54 30 33	\$ \$ \$
▲ End ▲ X-ray	ll Level ll Lung Det	A MRI A CT	Diseases Det COVID	출 End 출 X-ray	DiseasesCOVID	24 23	27 26	√ ✓
🖾 Der - Skin 🖾 Mic - Cell		🛱 FP - Fundus 🛱 CT - Kidney	Diseases Diseases	📽 Der - Skin 📽 Mic - Cell	Diseases Diseases	24 24	29 26	\$ \$

Table 13: Result of NEXT-Chat on CG by using detection and classification tasks to generalize classification Target dataset. Generalization results on classification datasets: *Related Combination* is the training set, *Target Subset* is the goal. Baseline and Trained represent the model's accuracy without training and trained on related data, respectively. \checkmark in *CG Helps* indicates successful generalization, while \varkappa denotes failure.

	Related Co	mbination		Target	Subset	Baseline	Trained	CG Helps
春 Lung 春 Lung 春 Bones 春 Bones	 Lung Det Lung Det Spinal Error Det Spinal Error Det 		 Diseases Diseases Diseases Diseases Diseases 	春 Lung 春 Lung 春 Bones 春 Bones	 Diseases Diseases Diseases Diseases Diseases 	41 41 31 31	47 49 35 37	\ \ \ \
을 End 을 X-ray	畠 Level 畠 Lung Det	A MRI A CT	Diseases Det COVID	A End A X-ray	DiseasesCOVID	24 22	26 23	\ \
🗟 Der - Skin 🗟 Mic - Cell	Ê Cancer Det Ê Cancer Det	🖓 FP - Fundus 🖓 CT - Kidney	Diseases Diseases	최 Der - Skin 최 Mic - Cell	Diseases Diseases	27 20	30 24	√ √

Table 14: Result of MiniGPT-v2 on CG by using detection and classification tasks to generalize classification Target dataset. Generalization results on classification datasets: *Related Combination* is the training set, *Target Subset* is the goal. Baseline and Trained represent the model's accuracy without training and trained on related data, respectively. \checkmark in *CG Helps* indicates successful generalization, while \checkmark denotes failure.

keeping all other training parameters at their default settings. Similarly, for MiniGPT-v2, we trained the backbone model from the second stage, starting with a learning rate of 2e-5 and gradually reducing it to 2e-6 over 3 epochs.

A.7 CG with Medical Multimodal LLM

In previous experiments, general MLLMs are selected to prevent the MLLM's inherent medical knowledge from affecting CG results. Our experiments focus on how MLLMs leverage CG to interpret unseen medical images. If the model has

Related Combination				Target Subset		HuatuoGPT
🖗 Bones	🖹 State	🖗 Breast	🖹 Diseases	🖗 Bones	🖹 Diseases	+6.12
🖗 Lung	🖹 COVID	🖗 Bones	🖹 Diseases	🖗 Lung	🖹 Diseases	+15.00
🗟 X-ray	🖹 Diseases	🖞 CT	🖹 COVID	▲ X-ray	🖹 COVID	+38.00
🗟 X-ray	🖹 Diseases	🖞 CT	🖹 State	🗟 X-ray	🖹 State	+40.67
🖄 CT	🖗 Brain	🗳 X-ray	🖗 Lung	🗳 CT	🖗 Lung	+1.5
🗟 CT	🖗 Brain	🗳 X-ray	🖗 Lung	🗳 CT	🖗 Lung	+18.00
최 FP - Fundus	🖹 Diseases	🗟 Mic - Cell	🖹 Level	🗟 FP - Fundus	🖹 Level	+12.12
🗟 Mic - Cell	B Recognition	🛱 FP - Fundus	🖹 Level	🗟 Mic - Cell	🖹 Level	+10.50

Table 15: Result of HuatuoGPT-Vision on selected classification datasets in Med-MAT. *HuatuoGPT* represent the accuracy(%) gains the model achieved through CG.

learned some fundamental elements of the *Target* data, it would compromise the fairness of the experiments.

To demonstrate that our results still work on medical LLMs, we employed the same data combinations from Section 3.3 to investigate CG on medical MLLMs (we selected HuatuoGPT-Vision as the baseline).

The results in Table 15 demonstrate that the medical-expert MLLM can still leverage CG to enhance their performance on novel tasks, further supporting the validity and consistency of our findings.

B The Dataset: Med-MAT

This section provides an overview of Med-MAT. First, a detailed explanation of MAT-Triplet will be presented in B.1. Next, the methods for constructing the QA formatting will be discussed in B.2. Finally, the data composition details and opensource specification will be provided in B.3.

B.1 Details of MAT-Triplet

MAT-Triplet stands for **M**edical Modality, Anatomical Area, and Medical **T**ask. We define all samples in Med-MAT using these three components and integrate datasets with identical triplets into subsets.

Medical Modality refers to different types of techniques or methods used in medical imaging or data acquisition. Each modality is designed to present the human body's structures or pathological features in unique ways, providing auxiliary support for clinical diagnosis and treatment. Most modalities exhibit significant visual differences, making them easily distinguishable. Med-MAT encompasses 11 modalities, including common ones such as Computed Tomography (CT), Magnetic Resonance Imaging (MRI), X-ray, Fundus Photography (FP), Endoscopy (End), Optical Coherence Tomography (OCT), and Ultrasound (US), as well as rare and specialized modalities like Colonoscopy (Co), Dermoscopy (Der), Digital Pathology (DP), and Microscopy (Mic).

Anatomical Area refers to specific anatomical structures or regions within the human body or other organisms, defined by distinct anatomical characteristics to describe various body parts, their functions, and relative positions. Med-MAT encompasses 14 anatomical areas, including the cervix, kidney, lung, brain, intestine, bladder, fundus, retina, breast, bones, and chest. To further facilitate data description, additional categories such as skin, mouth, and cell are included as specialized anatomical areas.

Medical Task refers to the specific detection task that needs to be performed on the dataset. Med-MAT includes 13 distinct tasks, with classification tasks encompassing Quality Identification (image quality analysis), COVID Diagnosis, Cancer Diagnosis (determining the presence of a specific disease), State (such as identifying brain hemorrhage), Level Identification (assessing disease severity), and Multiple Classification (classifying multiple diseases or cell types). Given the limited options of COVID Diagnosis and Cancer Diagnosis, these tasks can be interpreted as identifying whether a patient is in a diseased state. To enhance generalization and provide more diverse examples, these tasks are grouped under the broader category of State. In addition, we have 16 datasets defining segmentation or classification tasks with different objectives.

B.2 QA construction method

A large amount of image-label datasets was collected to build the Med-MAT dataset. To ensure compatibility with MLLM training inputs and outputs, all data is transformed into a questionanswering format. Questions are formulated based on modality, anatomical area, and medical task, with 6 question prompts applied to each subset.

The labels within each data subset will be clustered to prevent redundant definitions of the same condition. Then, all training set and test set will be converted into multiple-choice questions following the template in Table 8. Each question will have up to four options, with distractor options randomly selected from the corresponding subset.

B.3 Data composition and Open-source Specification

Med-MAT is composed of multiple datasets. After being transformed into different QA formats, the new data is organized into several subsets to support generalization experiments in medical imaging. Table 17 shows all of our subset datasets, which are separated based on different combinations in MAT-Triplet. The specific MAT-Triplets are listed, along with the labels corresponding to the imagelabel datasets for each subset. Correspondingly, all the image-label datasets are also displayed in Table 18, which includes their names, descriptions of the tasks performed, download links, and the level of accessibility.

All question-answering text datasets in Med-MAT will be publicly available. To accommodate varying access permissions, we will release datasets based on their respective licenses: openly accessible datasets will be directly available, while restricted datasets can be accessed by applying through the links provided in this paper. We hope this dataset will support and advance future generalization experiments on medical imaging.

B.4 Data Sources and Distribution

All Med-MAT data are sourced from public medical image challenges or widely used, high-impact datasets previously applied in deep learning training, ensuring reliable annotations. Before inclusion in Med-MAT, all datasets underwent label averaging where possible; test sets, in particular, were strictly balanced to ensure accuracy reliably reflects model performance. Each Med-MAT training subset contains 3,000 samples, while test sets maximize size under label balance constraints.

C Bad cases analysis and solutions

C.1 Bad case analysis

Some *Trained* models show minimal gains or even performance declines in Table 3, with classification

accuracy lower than either the *Baseline* or *Baseline*+. After a thorough examination, we found that these Target datasets require more fine-grained medical condition classification. Beyond disease presence, they need detailed assessments, such as severity grading (e.g., bone age estimation, cancer staging) or distinguishing similar conditions (e.g., differentiating COVID-19 from pneumonia).

- The *Related* combinations lack suitable fundamental elements: For CG, the training data must include the *Target* task's core elements. Here, we use other "level classification/grading" tasks for generalization, but their criteria differ significantly, misaligning with the *Target* task's needs.
- Without defined grading standards, MLLMs lacking relevant knowledge can't perform finegrained tasks: Tasks like bone age assessment and cancer staging vary by criteria, and without this knowledge, MLLMs can't accurately classify them.

C.2 Possible solutions

Few-shot prompting As we illustrated before, most of the bad cases involve fine-grained tasks needing specialized knowledge. So, in order to minimize the effect of a lack of relevant knowledge, we also conducted few-shot experiments to add some target images in the prompts. Subset *X-ray, Lung, Normal-COVID-Pneumonia* was chosen for its simple structure, with LLaVA as the baseline. We randomly sampled n images per label for n-shot inference and repeated each experiment 3 times.

Model	0-shot	2-shot	3-shot	4-shot
LLaVA	30.00	$ \textbf{28.83}\pm0.85$	29.33 ± 1.25	29.83 ± 1.31
LLaVA + CG	28.00	28.67 ± 0.94	$\textbf{37.00} \pm 0.82$	$\textbf{36.67} \pm 0.47$

Table 16: Results of Few-shot prompting.

The results in Table 16 demonstrate that training with CG combinations can improve the few-shot performance of MLLMs on downstream tasks, even when direct CG generalization is not effective.

Adding some *Target* data in training As described in Section 5.2, we selected cases where CG alone couldn't achieve satisfactory results and augmented their training sets with target data. The results in this section indicate that while CG may not directly enhance generalization, it accelerates the model's adaptation to downstream tasks.

Multiple-choice Questions Template

<question>

A. <option_1> B. <option_2> C. <option_3> D. <option_4> Answer with the option's letter from the given choices directly.





Figure 9: Illustration of diverse samples with varying numbers of candidate options in the Med-MAT dataset.

Subset No.	Modality	Anatomical Area	Task	Datasets No.
01	Co	Cervix	Cervical Picture Quality Evaluation	1
02	СТ	Kidney	Kidney Diseases Classification	2
03	СТ	Lung	COVID-19 Classification	3,4,6
04	СТ	Lung	Lung Cancer Classification	5
05	СТ	Brain	Brain Hemorrhage Classification	7
06	CT	Brain	Brain Cancer Classification	8
07	Der	Skin	Melanoma Type Classification	10
08	Der	Skin	Skin Diseases Classification	9, 11-15, 71, 72, 74
09	DP	Mouth	Teeth Condition Classification	16
10	DP	Mouth	Oral Cancer Classification	17
11	End	Intestine	Intestine Cleanliness Level	18
12	End	Bladder	Cancer Degree Classification	19
13	End	Intestine	Intestine Diseases Classification	20
14	FP	Fundus	Eye Diseases Classification	21-23, 26-28, 31, 32, 75
15	FP	Fundus	Multiple-labels Eye Diseases Classification	24, 25, 68
16	FP	Fundus	Blindness Level	29
17	FP	Fundus	Retinal Images Quality Evaluation	30
18	Mic	Cell	Cell Type Classification	33, 36-38, 39-41, 44, 65, 70
19	Mic	Cell	Prostate Cancer Degree Classification	34
20	Mic	Cell	Multiple-labels Blood Cell Classification	35
21	Mic	Cell	Cancer Classification	42, 67
22	MRI	Brain	Head Diseases Classification	44, 45
23	OCT	Retina	Retina Diseases Classification	46, 47
24	US	Breast	Breast Cancer Classification	48
25	X-ray	Bones	Degree Classification of Knee	49, 53
26	X-ray	Bones	Fractured Classification	50, 51
27	X-ray	Bones	Vertebrae Diseases Classification	52
28	X-ray	Lung	COVID-19 and Pneumonia Classification	54-57, 60, 62, 81
29	X-ray	Breast	Breast Diseases Classification	58, 78
30	X-ray	Lung	Tuberculosis Classification	59, 79
31	X-ray	Chest	Multiple-labels Chest Classification	61, 73, 76, 77, 80, 85, 87
32	X-ray	Brain	Tumor Classification	63
33	Mic	Cell	Multi-labels Diseases	84
34	FP	Fundus	Level Identification	66
35	X-ray	Bones	Level Identification	69
36	X-ray	Bones	Spinal lesion Classification	86
37	X-ray	Breast	Multi-labels Diseases	82
38	Der	Skin	Lesion Det/Seg	88-91
39	End	Intestine	PolyP Det/Seg	92-93
40	End	Intestine	Surgical Procedures Det/Seg	94
41	End	Intestine	Multi-labels Det/Seg	95
42	Mic	Cell	Cancer Cell Det/Seg	96
43	US	Chest	Cancer Det/Seg	97
44	US	Thyroid	Thyroid Nodule Region Det/Seg	98
45	MRI	Intestine	Multi-labels Det/Seg	103
46	MRI	Liver	Liver Det/Seg	104, 105
47	X-ray	Lung	Lung Det/Seg	99
48	X-ray	Lung	Pneumothorax Det/Seg	106
49	X-ray	Bones	Spinal Anomaly Det	100
50	X-ray	Chest	Multi-labels Det	101, 102
51	FP	Fundus	Vessel Seg	107
52	FP	Fundus	Optic Disc and Cup Seg	108
53	FP	Fundus	Optic Disc Seg	109

Table 17: The details of subset. In particular, **Co** stands for Colposcopy, **CT** represents Computed Tomography, **DP** refers to Digital Photography, **FP** is for Fundus Photography, **MRI** denotes Magnetic Resonance Imaging, **OCT** signifies Optical Coherence Tomography, **Der** refers to Dermoscopy, **End** stands for Endoscopy, **Mic** indicates Microscopy Images, and **US** represents Ultrasound. The blue section represents the classification dataset and the green section represents the detection

1Incl & MobileODT Carvical ScreeningCarvix Type in Screening(BenOt el. 2017)3SARS-COV-2 (-ScanCOVID (-) Classification Dataset(Soares et. 1, 2020)4COVID COVID-CTCOVID (-) Classification Dataset(Soares et. 1, 2020)5Chest CT-ScanCancer Classification Dataset(SumeYi, 2021)6COVID (-) CT SCAN IMAGESCOVID (-) CT SCAN IMAGES(Silicachuargyw, 2019)7Head CTHead Cancer(Data, 2023)9MED-NODEMelanoma on Naevus(Giotis et. 1, 2021)10ISIC 2020Melanoma on Vaevus(Giotis et. 1, 2021)11IPAD-UTFS-20Skin Mulic Classification(Bottmetry et al., 2021)12Web-screed Skin ImageSkin Lesion Classification(Guttma et al., 2021)13ISIB 2016Skin Lesion Classification(Guttma et al., 2016)14ISIC 2019Skin Lesion Classification(Katanski), 2024)15Skin Cancer ISICSkin Cancer Multi Classification(Katanski), 2024)16Denat Condition DatasetCancer Classification(Bocrelov et al., 2017)17Oral Cancer DatasetOral cancer Classification of eye diseases(Cor et al., 2017)18The Nerthw DatasetClaarenines Icevic(Porelov et al., 2017)19Endoscopic Bladder TissueCancer Classification of eye diseases(Cor et al., 2017)10ISIC Conser DatasetMulti Classification of eye diseases(Cor et al., 2021)11Staffic Alger Alger Alger Alger Alger Alger Alger Alger Alger	No.	Name	Description	Citation
2 CF kindney Dataset Normal or Cyst or Tumor (Islam et al., 2020) 4 SARS-COV-2 Ct-Scan COVID 19, Classification Dataset (Zano et al., 2020) 6 COVID 1-CT COVID-CT COVID 19, Classification Dataset (Zano et al., 2020) 6 COVID 1-CT SCAN IMAGES COVID 19, Classification (wiXiacchangev, 2019) 7 Head CT Head Hemorrhage (Kitumura, 2018) 8 CT of Brain Head Cancer (Data, 2023) 10 ISIC 2020 Melanoma, Benign or Malignant (Rotemberg et al., 2021) 11 IAD-UFES-20 Skin Desoase Multi Classification (Gatum et al., 2016) 13 ISIS 2016 Skin Cancer Sin Cancer Sin Cancer Sin Cancer Sin Cancer Classification (Gatum et al., 2016) 14 ISIC 2019 Skin Cancer Multi Classification (Gatum et al., 2017) 15 Skin Cancer Sin Cancer Classification (Gatum et al., 2017) 16 Dental Condition Dataset Tech condition classification of eye diseases (Corei al., 2022) 17 Free Condition Classification of eye diseases (Corei al., 2022) (Kardiguee et al., 2017) <	1	Intel & MobileODT Cervical Screening	Cervix Type in Screening	(BenO et al., 2017)
3 SARS-COV-2 C1-Scan COVID 91, Classification Dataset (Source et al., 2020) 5 Chest CT-Scan Cancer Classification (Sume Y), 2021) 6 COVID 91-CT SCAN IMAGES COVID 91, Classification (WiXiaochuangw, 2019) 7 Head C (Barnon) (WiXiaochuangw, 2019) 8 CT of Brain Head Cancer (Data, 2023) 9 MED-NODE Melanoma on Naevus (Giotis et al., 2021) 10 ISIC 2020 Melanoma on Naevus (Giotis et al., 2020) 11 PAD-UFES-20 Skin Descare Multi Classification (Gaume et al., 2021) 12 Web-screed Skin Image Skin Descare Multi Classification (Gaume et al., 2016) 13 ISIS 12016 Skin Descare Multi Classification (Katansky, 2019) 15 Skin Cancer Jatset Oral cancer Matset Cleanliness Iscare 1 (Pogorelov et al., 2017) 16 Dental Condition Dataset Cleanliness Iscare 1 (Pogorelov et al., 2017) 17 Oral Cancer Dataset Oral cancer Classification (Rostrue, 2012, 2017) 18 The Nerthus Dataset	2	CT Kindney Dataset	Normal or Cyst or Tumor	(Islam et al., 2022a)
4 COVID CT COVID-CT COVID 19, Classification Dataset (Zhao et al., 2020) 6 Chest CT-Scan Cancer Classification (w)Xiaochungw, 2019) 6 COVID 19-CT SCAN IMAGES COVID 19, Classification (w)Xiaochungw, 2019) 7 Head CT Head Amorthage (Kiamura, 2018) 8 CT of Brain Head Cancer (Data, 2023) 10 ISIS C 2020 Melanoma on Naevus (Giotis et al., 2015) 11 RAD-UFFS-20 Skin Multi Classification (Ratamskir, 2016) 12 Web-scraped Skin Inage Skin Desease Multi Classification (Guutan et al., 2020) 13 ISIS 12016 Skin Cancer Multi Classification (Guutan et al., 2016) 14 ISIC 2019 Skin Cancer Multi Classification (Guutan et al., 2021) 15 Skin Cancer Tister Canser Degree Classification (Lazo et al., 2023) 15 Fixaar Multi Classification of eye diseases (Gourt et al., 2021) 16 Endescopic Bladser Tissee Classification of eye diseases (Rodriguez et al., 2023) 17 ToxoFinudus Data Raw Celass Al	3	SARS-COV-2 Ct-Scan	COVID19, Classification Dataset	(Soares et al., 2020)
5 Chest CT-Scan Cancer Classification (SumeY), 2021) 7 Head CT Head Hemorrhage (WiXiaochuangw, 2019) 7 Head CT Head Cancer (Data, 2023) 9 MED-NODE Melanoma or Naevus (Giotis et al., 2015) 10 ISIC 2020 Melanoma Reing or Malignant (Rotemberg et al., 2021) 11 PAD-UFFS-20 Skin Descase Multi Classification (Gaust et al., 2020) 12 Web-scraped Skin Image Skin Descase Multi Classification (Gaust et al., 2014) 13 ISBI 2016 Skin Descase Multi Classification (Katanky, 2019) 14 ISIC 2019 Skin Cancer Multi Classification (Katanky, 2019) 15 Istic Cancer ISIC Skin Cancer Multi Classification (Katanky, 2019) 16 Dental Condition Dataset Creat classification (Katanky, 2019) 17 IFA Scrape Dataset Creat Classification (Layoe et al., 2017) 18 The Architas Dataset Classification of eye diseases (Cen et al., 2021) 18 The Architas Dataset Classification of eye diseasese	4	COVID CT COVID-CT	COVID19, Classification Dataset	(Zhao et al., 2020)
6 COVID-19-CT SCAN IMAGES COVID-19, Classification (wjXiaochungw, 2019) 8 CT of Brain Head Cancer (Data, 2023) 9 MED-NODE Melanoma on Naevus (Giotis et al., 2015) 10 ISIC 2020 Melanoma on Naevus (Giotis et al., 2021) 11 PAD-UFES-20 Skin Multi Classification (Hacheco et al., 2020) 12 Web-scraped Skin Image Skin Desease Multi Classification (Giuman et al., 2016) 13 ISBI 2016 Skin Loscer Multi Classification (Katansky, 2019) 14 Sixi Cancer Multi Classification (Gamines et al., 2017) 15 Skin Cancer Jataset Oral cancer Classification (RASHID, 2024) 16 Dental Condition Dataset Cleaniness Revel (Pogorelov et al., 2017) 17 Oral Cancer Dataset Oral cancer Classification of ey diseases ((Cardozo et al., 2021) 17 Kvasir Multi Classification of ey diseases (Cardozo et al., 2021) 18 The Nerthus Dataset Clasmines and exace (Cardozo et al., 2023) 17 ToxoFundus/Data Processed Paper)	5	Chest CT-Scan	Cancer Classification	(SunneYi, 2021)
7 Head T Head Cancer (Kitamurz, 2018) 9 MED-NODE Melanoma or Naevus (Giotis et al., 2015) 10 ISIC 2020 Melanoma, Benign or Malignant (Rotemberg et al., 2021) 11 PAD-UFES-20 Skin Descase Multi Classification (Pacheco et al., 2020) 13 ISBI 2016 Skin Descase Multi Classification (Guttma et al., 2016) 14 ISIC 2019 Skin Locase Autil Classification (Guttma et al., 2016) 15 Skin Cancer Multi Classification (Gattma et al., 2017) 16 Dental Condition Dataset Cleanliness level (Rogencive et al., 2017a) 18 The Nerthus Dataset Cleanliness level (Pogorclov et al., 2017a) 12 AcRIMA Classification of eye diseases (Cen et al., 2021) 21 ACRIMA Guacoma (Ovreiu et al., 2021) 22 Agermeted coular diseases AOD Multi Classification of eye diseases (Rodfiguez et al., 2023) 23 ToxoFundus/Data Rwe disease1 Oural xoxoFandus/Data Rwe disease1 (Cen et al., 2021) 23 ISIEC Multi Classification<	6	COVID-19-CT SCAN IMAGES	COVID19, Classification	(wjXiaochuangw, 2019)
8 CT of Brain Head Cancer (Data, 2023) 9 MED-NODE Melanoma or Naevus (Rotiner) (Rother) (Rothe	7	Head CT	Head Hemorrhage	(Kitamura, 2018)
9 MED-NODE Melanoma or Naevus (Giotis et al., 2015) 11 IAD-UFES-20 Skin Multi Classification (Rotemberg et al., 2021) 12 Web-screped Skin Image Skin Desease Multi Classification (Giutam et al., 2020) 13 ISBI 2016 Skin Lesion Classification (Gutam et al., 2020) 14 ISIC 2019 Skin Cancer Multi Classification (Katansky, 2019) 15 Skin Cancer ISIC Skin Cancer Multi Classification (RASHID, 2024) 16 Oral Cancer Classification (Cancer et al., 2017a) (Lazo et al., 2023) 17 The Nerthus Dataset Cleaniness level (Porciev et al., 2017a) 18 The Nerthus Dataset Classification of eye diseases (Core et al., 2021) 18 Augemented ocular diseases AOD Multi Classification of eye diseases (Core et al., 2021) 14 Auti-Abel Retinal Diseases Multi Classification of eye diseases (Cordozo et al., 2023) 15 ToxoFundus(Data Raw 6class AI) Ocular toxoFundus(Data Raw 6class AI) Ocular toxoFundus(Data Raw 6class AI) Ocular toxoFundus(Data Raw 6class AI) 18 Adam datas	8	CT of Brain	Head Cancer	(Data, 2023)
10 ISIC 2020 Melanoma, Benign or Malignant (Rotenberg et al., 2021) 11 PAD-UFES-20 Skin Muit Classification (Islam et al., 2020) 12 Web-scraped Skin Image Skin Desease Multi Classification (Gurman et al., 2016) 13 ISBI 2016 Skin Cancer Classification (Combalia et al., 2019) 15 Skin Cancer ISIC Skin Cancer Classification (Ratsky, 2019) 16 Dental Condition Dataset Certh condition classification (Rogorelov et al., 2021) 17 Oral Cancer Dataset Oral cancer Classification (Lagorelov et al., 2021) 18 The Nerthus Dataset Cleanliness level (Pogorelov et al., 2021) 2 Augemnted ocular diseases AOD Multi Disease Classification of cyc diseases (Cert et al., 2021) 2 Augemnted ocular diseases AOD Multi Classification of cyc diseases (Cardoz et al., 2023) 3 ToxoFundus(Data Processed Paper) Ocular toxoplasmosis (Cardoz et al., 2023) 4 ToxoFundus(Data Processed Paper) Ocular toxoplasmosis (Cardoz et al., 2023) 4 ToxoFundus(Data Processed Paper) Ocular toxoplasmosis (Cardoz et al., 2023) 5 <td>9</td> <td>MED-NODE</td> <td>Melanoma or Naevus</td> <td>(Giotis et al., 2015)</td>	9	MED-NODE	Melanoma or Naevus	(Giotis et al., 2015)
11 PAD-UFES-20 Skin Multi Classification (Islam et al., 2020) 13 ISBI 2016 Skin Desease Multi Classification (Gutam et al., 2016) 13 ISBI 2016 Skin Desease Multi Classification (Gutam et al., 2019) 15 Skin Cancer ISIC Skin Cancer Multi Classification (Katamskiy, 2019) 16 Dental Condition Dataset Teeth condition classification (RASHID, 2024) 17 Oral cancer Dataset Oral cancer Classification (RASHID, 2024) 18 The Nerthus Dataset Cleanificase Nevel (Pooreiu et al., 2017a) 19 Endoscopic Bladder Tissue Classr Degree Classification (Pooreiu et al., 2017b) 21 ACRIMA Multi Classification of eye diseases (Cen et al., 2012) 21 Acgemented ocular diseases ADD Multi Classification of eye diseases (Rodriguez et al., 2022) 21 Kristi D 2.0 Multi Classification of eye diseases (Cen et al., 2023) 22 RFMiD 2.0 Multi Classification (Cardoz et al., 2023) 23 ToxoFundus(Data Raw 6class All) Occular toxoplasmosis (Cardoz et al., 2023) 24 Adum dataset Age-related Macul	10	ISIC 2020	Melanoma, Benign or Malignant	(Rotemberg et al., 2021)
12 Web-scraped Skin Image Skin Desease Multi Classification (Gutnan et al., 202b) 13 ISBI 2016 Skin Cascin Classification (Gutnan et al., 2016) 14 ISIC 2019 Skin Desease Multi Classification (Katanskiy, 2019) 16 Dental Condition Dataset Teeth condition Classification (Katanskiy, 2019) 17 Oral Cancer Dataset Oral cancer Classification (Raget al., 2023) 18 The Nerthus Dataset Cleaninese level (Pogorelov et al., 2017a) 18 The Nerthus Dataset Cleaninese Isevel (Pogorelov et al., 2021) 20 Kvasir Multi Classification of eye diseases (Cent et al., 2021) 21 Augemented ocular diseases AOD Multi Classification of eye diseases (Rodriguez et al., 2022) 21 Multi Classification of eye diseases (Cardozo et al., 2023) (Cardozo et al., 2023) 22 ToxoFundus(Data Processed Paper) Ocular toxoplasmosis (Cardozo et al., 2023) 23 ToxoFundus(Data Raw oclass All) Ocular toxoplasmosis (Cardozo et al., 2023) 24 APTOS 2019 Blindness Blindness Level Identification (Karthik et al., 2019) 24 <t< td=""><td>11</td><td>PAD-UFES-20</td><td>Skin Multi Classification</td><td>(Pacheco et al., 2020)</td></t<>	11	PAD-UFES-20	Skin Multi Classification	(Pacheco et al., 2020)
13 ISBI 2016 Skin Lesion Classification (Gurman et al., 2016) 14 ISIC 2019 Skin Cancer Multi Classification (Gatanskiy, 2019) 15 Skin Cancer ISIC Skin Cancer Multi Classification (Gatanskiy, 2019) 16 Dental Condition Dataset Teeth condition classification (Basting) 17 Oral Cancer Dataset Cleaniness level (Pogorelov et al., 2017a) 18 The Nerthus Dataset Cleaniness level (Doycelov et al., 2017a) 19 Endoscopic Bladder Tissue Classification of eye diseases (Ovreiu et al., 2021) 21 ACRIMA Multi Classification of eye diseases (Rodriguez et al., 2022) 23 RFMID 2.0 Multi Classification of eye diseases (Rodriguez et al., 2023) 22 ToxoFundus(Data Processed Paper) Coular toxoplasmosis (Cardozo et al., 2023) 24 Adam dataset Age-related Macular Degeneration (Liang, 2021) 23 Adam dataset Blindness Level Identification (Arbik et al., 2013) 24 Algood Glaucoma Classification (Di et al., 2014) 25 Algood Glaucoma Classification (Di et al., 2015) <td>12</td> <td>Web-scraped Skin Image</td> <td>Skin Desease Multi Classification</td> <td>(Islam et al., 2022b)</td>	12	Web-scraped Skin Image	Skin Desease Multi Classification	(Islam et al., 2022b)
14 ISIC 2019 Skin Desease Multi Classification (Combalia et al., 2019) 15 Skin Cancer ISIC Skin Cancer Multi Classification (Katanski), 2019) 16 Dental Condition Dataset Teeth condition classification (RASHID, 2024) 17 Oral Cancer Dataset Oral cancer Classification (RASHID, 2023) 18 The Nerthus Dataset Cleanliness level (Pogorelov et al., 2017a) 18 The Nerthus Dataset Cleanliness level (Pogorelov et al., 2021) 18 Kvasir Multi Classification of eye diseases (e et al., 2021) 2 Augemented ocular diseases AOD Multi Classification of eye diseases (Rodriguez et al., 2022) 2 RyEMiD 2.0 Multi Classification of eye diseases (Cardozo et al., 2023) 2 ToxoFundus(Data Processed Paper) Ocular toxoplasmosis (Cardozo et al., 2023) 3 APTOS 2019 Blindness Blindness Level Identification (Kasthi, al., 2019) 3 APTOS 2019 Blindness Glaucoma Classification (Oh ret al., 2013) 3 Glaucoma Classification (Qi et al., 2014) (Mote et al., 2025) 3 IGAPAOS 2019 Blindness Glal	13	ISBI 2016	Skin Lesion Classification	(Gutman et al., 2016)
15 Skin Cancer SIRC Skin Cancer Multi Classification (Katanskiy, 2019) 16 Dental Condition Dataset Oral cancer Classification (Katanskiy, 2019) 17 Oral Cancer Dataset Oral cancer Classification (RASHID, 2024) 19 Endoscopic Bladder Tissue Cleanliness level (Pogorelov et al., 2017a) 20 Kvasir Multi Disease Classification (Datoe) et al., 2021) 21 Acgemited ocular diseases AOD Multi Classification of eye diseases (Rodriguez et al., 2022) 21 Multi Lassification of eye diseases (Rodriguez et al., 2023) (Cen et al., 2021) 23 ToxoFundus/Data Raw 6class AII) Ocular toxoplasmosis (Cardozo et al., 2023) 26 ToxoFundus/Data Raw 6class AII) Quality Testing of Retinal Images (Prential, et al., 2023) 24 Adam dataset Age-related Macular Degeneration (Liarag, 2021) 25 Alsm dataset Quality Testing of Retinal Images (Prentiasc et al., 2013) 26 Glaucoma Detection Glaucoma Classification (Gardozo et al., 2023) 26 ARKOGS Glaucoma Classification (Dive et al., 2021) 27 StoAPy2	14	ISIC 2019	Skin Desease Multi Classification	(Combalia et al., 2019)
16 Dental Condition Dataset Teth condition (Lassification) (Rajd, 2024) 17 Oral Cancer Dataset Oral cancer Classification (RaSHID, 2024) 18 The Nerthus Dataset Cleanliness level (Pogorelov et al., 2017a) 20 Kvasir Multi Disease Classification (Lazo et al., 2023) 21 Augemetd ocular diseases AOD Multi Classification of eye diseases (Cen et al., 2021) 23 JSIEC Multi Classification of eye diseases (Rodriguez et al., 2023) 24 Multi-Label Retinal Diseases Multi Classification of eye diseases (Panchal et al., 2023) 26 ToxoFundus/Data Processed Paper) Ocular toxoplasmosis (Cardozo et al., 2023) 27 ToxoFundus/Data Processed Paper) Ocular toxoplasmosis (Cardozo et al., 2023) 28 Adam dataset Age-related Macular Degeneration (Liang, 2021) 30 DRIMDB Quality Testing of Retinal Images (Prentasic et al., 2013) 31 Glaucoma Classification (Glava Coma Classification) (Glava.Coma Classification) 31 Glaucoma Classification (Glava.Coma Classification) (Glava.Coma Classificatia) 32 <td< td=""><td>15</td><td>Skin Cancer ISIC</td><td>Skin Cancer Multi Classification</td><td>(Katanskiy, 2019)</td></td<>	15	Skin Cancer ISIC	Skin Cancer Multi Classification	(Katanskiy, 2019)
17 Oral Cancer Dataset Oral Cancer Classification (RASHID, 2024) 19 Endoscopic Bladder Tissue Cleanliness level (Pogorelov et al., 2017a) 20 Kvasir Multi Disease Classification (Pogorelov et al., 2023) 21 ACRMA Glaccoma (Ovreiu et al., 2021) 22 Augennied ocular diseases ADD Multi Classification of eye diseases (Rodríguez et al., 2022) 23 RFMiD 2.0 Multi Classification of eye diseases (Rodríguez et al., 2023) 24 Multi Classification of eye diseases (Panchal et al., 2023) 25 RFMiD 2.0 Ocular toxoplasmosis (Cardozo et al., 2023) 26 ToxoFundus/Data Raw 6class AII) Ocular toxoplasmosis (Cardozo et al., 2023) 27 ToxoFundus/Data Raw 6class AII) Ocular toxoplasmosis (Cardozo et al., 2023) 27 ToxoFundus/Data Raw 6class AII) Ocular toxoplasmosis (Cardozo et al., 2023) 28 Adam dataset Age-related Macular Degeneration (Liang, 2021) 29 APTOS 2019 Blindness Blandness Level Identification (Qiet et al., 2016) 31 ICPR-HEp-2 Multi Classification (Qiet et al., 20	16	Dental Condition Dataset	Teeth condition classification	(Sajid, 2024)
18 The Nerthus Dataset Cleanincess level (Pogorelov et al., 2017a) 19 Endoscopic Bladder Tissue Canser Degree Classification (Laz oct al., 2023) 20 Kvasir Multi Disease Classification (Pogorelov et al., 2017b) 21 ACRMA Glaucoma (Ovreiu et al., 2021) 22 Augemited ocular diseases AOD Multi Classification of eye diseases (Cen et al., 2022) 23 JSIEC Multi Classification of eye diseases (Rodríguez et al., 2023) 24 Multi-Label Retinal Diseases Multi Classification of eye diseases (Cardozo et al., 2023) 25 RFMiD 2.0 Multi Classification of eye diseases (Cardozo et al., 2023) 26 ToxoFundus/Data Processed Paper) Ocular toxoplasmosis (Cardozo et al., 2023) 26 APTOS 2019 Blindness Bindness Level Identification (Karthik et al., 2019) 21 APTOS 2019 Blindness Glaucoma Classification (Qi et al., 2021) 22 AIROGS Glaucoma Classification (Glave Pretasification (Gi evente et al., 2023) 21 Glaucoma Classification (Di et al., 2016) (Si LAN2 (Di et al., 2021) 23 <td>17</td> <td>Oral Cancer Dataset</td> <td>Oral cancer Classification</td> <td>(RASHID, 2024)</td>	17	Oral Cancer Dataset	Oral cancer Classification	(RASHID, 2024)
19 Endoscopic Bladder Tissue Canser Degree Classification (Lazo et al., 2023) 20 Kvasir Multi Disease Classification (Ovreiu et al., 2017b) 21 Augemnted ocular diseases AOD Multi Classification of eye diseases (Cen et al., 2021) 21 JSIEC Multi Classification of eye diseases (Rodriguez et al., 2022) 23 RFMID 2.0 Multi Classification of eye diseases (Candozo et al., 2023) 24 Multi Classification of eye diseases (Cardozo et al., 2023) 25 RFMID 2.0 Ocular toxoplasmosis (Cardozo et al., 2023) 26 ToxoFundus(Data Raw 6class AII) Ocular toxoplasmosis (Cardozo et al., 2013) 26 Adam dataset Age-related Macular Degeneration (Liang, 2021) 28 Adam dataset Age-related Macular Degeneration (Liang, 2014) 29 APTOS 2019 Blindness Blandness Level Identification (Øi et al., 2013) 316 Glaucoma Detection Glaucoma Classification (Øi et al., 2013) 316 ICAPA-HEP-2 Multi Classification (Øi et al., 2014) 316 Blood Cell Images Blood Cell Classification (Matker, 2018)	18	The Nerthus Dataset	Cleanliness level	(Pogorelov et al., 2017a)
20 Kvasir Multi Disease Classification (Pogorelov et al., 201/b) 21 AQREMA Glaucoma (Ovreiu et al., 2021) 22 Augemented ocular diseases AOD Multi Classification of eye diseases (Rodriguez et al., 2022) 23 JSIEC Multi Classification of eye diseases (Rodriguez et al., 2023) 24 Multi Classification of eye diseases (Rodrozo et al., 2023) 25 RFMiD 2.0 Ocular toxoplasmosis (Cardozo et al., 2023) 26 ToxoFundus(Data Raw 6class All) Ocular toxoplasmosis (Cardozo et al., 2023) 27 ToxoFundus(Data Raw 6class All) Ocular toxoplasmosis (Cardozo et al., 2023) 27 ToxoFundus(Data Raw 6class All) Ocular toxoplasmosis (Cardozo et al., 2023) 26 APTOS 2019 Blindness Blindness Level Identification (Kartik et al., 2019) 27 AIROGS Glaucoma Classification (Qi et al., 2023) 28 AIROGS Glaucoma Classification (Qi et al., 2023) 29 AIROGS Glaucoma Classification (Qi et al., 2018) 30 ICPR-HEp-2 Multi Classification (Mitwa-Rodriguez et al., 2020)	19	Endoscopic Bladder Tissue	Canser Degree Classification	(Lazo et al., 2023)
21 Augemnted ocular diseases AOD Multi Classification of eye diseases (Cvreu et al., 2021) 23 JSIEC Multi Classification of eye diseases (Cen et al., 2021) 24 Multi-Label Retinal Diseases Multi Classification of eye diseases (Rodriguez et al., 2023) 25 RFMID 2.0 Multi Classification of eye diseases (Panchai et al., 2023) 26 ToxoFundus(Data Rw oclass AII) Ocular toxoplasmosis (Cardozo et al., 2023) 27 ToxoFundus(Data Rw oclass AII) Ocular toxoplasmosis (Cardozo et al., 2023) 28 Adam dataset Age-related Macular Degeneration (Liang, 2021) 29 APTOS 2019 Blindness Blindness Level Identification (Øe thet et al., 2013) 31 Glaucoma Detection Glaucoma Classification (Øe thet et al., 2023) 32 AIROGS Glaucoma Classification (Øe thet et al., 2021) 34 SICAPv2 Cancer Degree Classification (Nita-Rodriguez et al., 2020) 35 Blood Cell Images Blood Cell Classification (Matket, et al., 2021) 36 Bore Marrow Cell Classification (Matket, et al., 2021) (Matket, et al., 2021) 37 Che	20	Kvasir	Multi Disease Classification	(Pogorelov et al., 2017b)
22 Augemented ocular diseases AOD Multi Classification of eye diseases (Cen et al., 2021) 24 Multi-Label Retinal Diseases Multi Classification of eye diseases (Rodriguez et al., 2022) 25 RFMiD 2.0 Cocular toxoplasmosis (Cardozo et al., 2023) 26 ToxoFundus(Data Processed Paper) Ocular toxoplasmosis (Cardozo et al., 2023) 27 ToxoFundus(Data Raw 6class AII) Ocular toxoplasmosis (Cardozo et al., 2023) 28 Adam dataset Age-related Macular Degeneration (Liang, 2021) 29 APTOS 2019 Blindness Blindness Level Identification (Karthik et al., 2013) 20 Glaucoma Detection Glaucoma Classification (Qi et al., 2013) 31 ICPR-HEp-2 Multi Classification (Qi et al., 2014) 33 ICPR-HEp-2 Multi Classification (Silva-Rodriguez et al., 2020) 34 Bood Cell Images Blood Cell Classification (Matex et al., 2014) 35 Bood Cell Classification (Matex et al., 2014) (Matex et al., 2014) 41 Malignant Lymphoma Classification Multi Classification (Matex et al., 2014) 42 Brain Tumor 17 Classes	21	ACRIMA	Glaucoma	(Ovreiu et al., 2021)
23JNHCMulti Classification of eye diseases(Cen et al., 2021)24Multi-Label Retinal DiseasesMulti Classification of eye diseases(Rodríguz et al., 2023)25RFMiD 2.0Multi Classification of eye diseases(Rodríguz et al., 2023)26ToxoFundus(Data Raw 6class All)Ocular toxoplasmosis(Cardoz et al., 2023)27ToxoFundus(Data Raw 6class All)Ocular toxoplasmosis(Cardoz et al., 2023)28Adam datasetAge-related Macular Degeneration(Liang, 2021)29APTOS 2019 BlindnessBlindness Level Identification(Karthik et al., 2019)20DRIMDBQuality Testing of Retinal Images(Prentasic et al., 2023)31Glaucoma DetectionGlaucoma Classification(Mooney, 2017)33ICPR-HEp-2Multi Classification(Giovernet al., 2024)34Blood Cell Classification(Mooney, 2017)35Blood Cell ImagesBlood Cell Classification(Mooney, 2017)36BreakHisCell type and beginormag(Bukun, 2019)37ChaoyangMulti Classification of puthologyts(Zhu et al., 2021)40NCT-CRC-HE-100KMulti Classification(Mate et al., 2018)41Lisopatohogic Cancer DetectionCancer Classification(Mate et al., 2021)42Histopathologic Cancer DetectionMulti Classification(Cukierski, 2018)43LC25000Multi Classification(Cukierski, 2018)44Multi ClassificationMulti Classification(Cukierski, 2020)	22	Augemnted ocular diseases AOD	Multi Classification of eye diseases	(, 2021)
24 Multi-Laber Retinal Diseases Multi Classification of eye diseases (Ranchal et al., 2023) 26 ToxoFundus(Data Processed Paper) Ocular toxoplasmosis (Cardozo et al., 2023) 27 ToxoFundus(Data Raw 6class All) Ocular toxoplasmosis (Cardozo et al., 2023) 27 ToxoFundus(Data Raw 6class All) Ocular toxoplasmosis (Cardozo et al., 2023) 28 Adam dataset Age-related Macular Degeneration (Liang, 2021) 29 APTOS 2019 Blindness Blindness Level Identification (Karthik et al., 2013) 31 Glaucoma Detection Glaucoma Classification (Qi et al., 2023) 31 ICPR-HEp-2 Multi Classification (Qi et al., 2021) 31 ICPR-HEp-2 Multi Classification (Bitoa., 2016) 32 SICAPV2 Cancer Degree Classification (Bitoa., 2017) 34 SICAPV2 Cancer Degree Classification (Booncn, 2017) 35 Blood Cell Images Blood Cell Classification (Multi Classification (Bucun, 2019) 36 HuSHeM Sperm Head Morphology Classification (Matek et al., 2021) (Matek et al., 2021) 36 HuSHeM <t< td=""><td>23</td><td>JSIEC</td><td>Multi Classification of eye diseases</td><td>(Cen et al., 2021)</td></t<>	23	JSIEC	Multi Classification of eye diseases	(Cen et al., 2021)
25Multi Classification of eye diseases(Fanchal et al., 2023)27ToxoFundus(Data Processed Paper)Ocular toxoplasmosis(Cardozo et al., 2023)28Adam datasetAge-related Macular Degeneration(Liang, 2021)29APTOS 2019 BlindnessBlindness Level Identification(Karthik et al., 2019)30DRIMDBQuality Testing of Retinal Images(Prentasic et al., 2023)31Glaucoma DetectionGlaucoma Classification(de Vente et al., 2023)32AIROGSGlaucoma Classification(Qi et al., 2016)33ICPR-HEp-2Multi Classification(Gi et al., 2016)34SICAPv2Cancer Degree Classification(Silva-Rodríguez et al., 2020)35Blood Cell ImagesBlood Cell Classification(Mooney, 2017)36BreakHisCell type and beginormag(Bukun, 2019)37ChaoyangMulti Classification of pathologists(Zhu et al., 2021)38HuSHeMSperm Head Morphology Classification(Matek et al., 2021)39Bone Marrow Cell ClassificationMulti Classification(Matek et al., 2018)41Malignant Lymphoma ClassificationMulti Classification of Ung and Colon(Zhu et al., 2018)43L22500Multi Classification of ve diseases(Orlov et al., 2010a)44Brain Tumor 17 ClassesMulti Classification of Yee diseases(Orlov et al., 2010)45BUSIBone Fracture Multi-Region X-ray DataFractured Classification of Eve diseases(Orlov et al., 2010)46Ma	24	Multi-Label Retinal Diseases	Multi Classification of eye diseases	(Rodriguez et al., 2022)
20Ioxorindus(Data Processed Paper)Ocular toxoplasmosis(Cardozo et al., 2023)7Toxorondus(Data Raw 6class All)Age-related Macular Degeneration(Liang, 2021)28Adam datasetAge-related Macular Degeneration(Liang, 2021)29APTOS 2019 BlindnessBlindness Level Identification(Karthik et al., 2019)31Glaucoma DetectionGlaucoma Classification(Jet al., 2013)31Glaucoma DetectionGlaucoma Classification(Jet al., 2013)31ICRP-HEp-2Multi Classification(Qi et al., 2023)31ICRP-HEp-2Multi Classification(Gi et al., 2016)34SICAPv2Cancer Degree Classification(Silva-Rodriguez et al., 2020)35Blood Cell ImagesBlood Cell Classification of pathologists(Zhu et al., 2021)36HuSHeMSperm Head Morphology Classification(Mater, 2018)37ChaoyangMulti Classification(Mater et al., 2018)38Bone Marrow Cell ClassificationMulti Classification(Culverski, 2018)40NCT-CRC-HE-100KMulti Classification(Culverski, 2018)41Malignant Lymphoma ClassificationMulti Classification(Culverski, 2018)42LC25000Multi Classification of eye diseases(Orlov et al., 2010a)43LC25000Multi Classification of eye diseases(Orlov et al., 2020)44Brain Tumor 17 ClassesMulti Classification of eye diseases(Orlov et al., 2020)45Tumor ClassificationHulti Classification	25	KFMID 2.0	Multi Classification of eye diseases	(Panchal et al., 2023)
27Interventionality Data Raw Octass All)Octual toxopitasinosis(Catadozo et al., 2023)28Adam datasetAge-related Macular Degeneration(Liang, 2021)29APTOS 2019 BlindnessBlindness Level Identification(Karthik et al., 2019)30DRIMDBQuality Testing of Retinal Images(Prentasic et al., 2013)31Glaucoma DetectionGlaucoma Classification(Zhang and Das, 2022)32AIROGSGlaucoma Classification(Qi et al., 2016)33ICPR-HEp-2Multi Classification(Gi et al., 2017)34SICAPv2Cancer Degree Classification(Silva-Rodfriguez et al., 2020)35Blood Cell ImagesBlood Cell Classification of pathologists(Zhu et al., 2019)36BreakHisCell type and beginormag(Bukun, 2019)37ChaoyangMulti Classification of pathologists(Zhu et al., 2021a)38HuSHeMSperm Head Morphology Classification(Matek et al., 2018)39Bone Marrow Cell ClassificationMulti Classification(Orlov et al., 2018)41Malignant Lymphoma ClassificationMulti Classification(Cukierski, 2018)42Histopathologic Cancer DetectionMulti Classification of Lung and Colon(Zhu, 2022)44Brain Tumor 17 ClassesMulti Classification of eye diseases(Orlov et al., 2020)45Tumor ClassificationMulti Classification of eye diseases(Orlov et al., 2020)46Bolse Fracture Multi-Region X-ray DataFractured Classification of eye diseases(Orlov	20	ToxoFundus(Data Processed Paper)	Ocular toxoplasmosis	(Cardozo et al., 2023)
25Avain datasetAge-related Machan Degeneration(Lang, 2017)26APTOS 2019 BlindnessBlindness Level Identification(Karthik et al., 2013)31Glaucoma DetectionGlaucoma Classification(Zhang and Das, 2022)31AIROGSGlaucoma Classification(de Vente et al., 2023)33ICPR-HEp-2Multi Classification(Qi et al., 2016)34SICAPv2Cancer Degree Classification(Silva-Rodríguez et al., 2020)35Blood Cell ImagesBlood Cell Classification(Money, 2017)36BreakHisCell type and beginormag(Bukun, 2019)37ChaoyangMulti Classification of pathologists(Zhu et al., 2021)38HuSHeMSperm Head Morphology Classification(Mate et al., 2018)39Bone Marrow Cell ClassificationBone Marrow Cell Classification(Mate et al., 2018)41Malignant Lymphoma ClassificationMulti Classification(Cukierski, 2018)42Histopathologic Cancer DetectionCancer Classification(Cukierski, 2018)44Brain Tumor 17 ClassesMulti Classification of Lung and Colon(Kickparvar, 2021a)45Tumor ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Subrawar, 2021b)47Retinal OCT-C8BusitBracture Classification of eye diseases(Orlov et al., 2010b)48BUSIBracture Multi-Region X-ray DataFracture Classification of Knee <t< td=""><td>21</td><td>A dam datasat</td><td>A ga related Magular Degeneration</td><td>(Cardozo et al., 2025)</td></t<>	21	A dam datasat	A ga related Magular Degeneration	(Cardozo et al., 2025)
29Ar 103 2019 BillindiessDifficuence Server Internited and Market and Argents29DRRMDBQuality Testing of Retinal Images(Prentasic et al., 2013)31Glaucoma DetectionGlaucoma Classification(de Vente et al., 2023)31ICPR-HEp-2Multi Classification(Qi et al., 2016)34SICAPv2Cancer Degree Classification(Silva-Rodríguez et al., 2020)35Blood Cell ImagesBlood Cell Classification(Mooney, 2017)36BreakHisCell type and beginormag(Bukun, 2019)37ChaoyangMulti Classification of pathologists(Zhu et al., 2021a)38HuSHeMSperm Head Morphology Classification(Mater et al., 2018)39Bone Marrow Cell ClassificationBone Marrow Cell Classification(Matter et al., 2018)41Malignant Lymphoma ClassificationMulti Classification(Cukierski, 2018)42Histopathologic Cancer DetectionCancer Classification of Lung and Colon(Cukierski, 2018)43LC25000Multi Classification of Lung and Colon(Nickparvar, 2021a)44Malignant Lymphoma ClassificationMulti Classification of eye diseases(Subramanian et al., 2022)45Tumor ClassificationMulti Classification of eye diseases(Subramanian et al., 2022)46Malignant Lymphoma ClassificationMulti Classification of Knee(Gornale and Patravali, 2020)47Retinal OCT-C8Multi Classification of Knee(Gornale and Patravali, 2020)48BUSIBreast Cancer(Al-Dhabyan	20	ADTOS 2010 Blindness	Right Revel Identification	(Lially, 2021) (Karthik at al. 2010)
30District DistributionQuality Testing of Retinal Inlages(Tertinate et al., 2019)31Glaucoma Classification(Zhang and Das, 2022)32AIROGSGlaucoma Classification(Qi et al., 2016)33ICPR-HEp-2Multi Classification(Gi et al., 2023)34SICAPv2Cancer Degree Classification(Silva-Rodríguez et al., 2020)35Blood Cell ImagesBlood Cell Classification of pathologists(Zhu et al., 2021a)36BreakHisCell type and beginormag(Bukun, 2019)37ChaoyangMulti Classification of pathologists(Zhu et al., 2021a)38HuSHeMSperm Head Morphology Classification(Matek et al., 2018)39Bone Marrow Cell ClassificationBone Marrow Cell Classification(Matek et al., 2010a)40NCT-CRC-HE-100KMulti Classification(Orlov et al., 2010a)41Malignant Lymphoma ClassificationMulti Classification(Chier, s022)43LC25000Multi Classification of Lung and Colon(Feltrin, 2022)44Brain Tumor 17 ClassesMulti Classification of eye diseases(Orlov et al., 2010b)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2021b)47Retinal OCT-C8Multi Classification of Knee(Gornale and Patravali, 2020)48BUSIBreast Cancer(Al-Dhabyani et al., 2022)48BUSIBracture Classification(Nickparvar, 2021b)49Digital Knee X-Ray ImagesDegree Classification of K	29	DPIMDR	Quality Testing of Patinal Images	(Ratulik et al., 2019) (Prentasic et al. 2013)
11DiateOna ClassificationChargent Display13AIROGSGlaucoma Classification(Qi et al., 2023)14CPR-HEp-2Multi Classification(Qi et al., 2016)15Blood Cell ImagesBlood Cell Classification(Mooney, 2017)16BreakHisCell type and beginormag(Bukun, 2019)17ChaoyangMulti Classification of pathologists(Zhu et al., 2021)18HuSHeMSperm Head Morphology Classification(Matek et al., 2021)19NCT-CRC-HE-100KMulti Classification(Matek et al., 2018)11Malignant Lymphoma ClassificationMulti Classification(Cukierski, 2018)11Malignant Lymphoma ClassificationMulti Classification(Cukierski, 2018)11Malignant Lymphoma ClassificationMulti Classification(Cukierski, 2018)12Histopathologic Cancer DetectionCancer Classification of Lung and Colon(Zhu, 2022)14Brain Tumor 17 ClassesMulti Classification of eye diseases(Orlov et al., 2010a)17Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)18BUSIBreast Cancer(Al-Dhabyani et al., 2020)19Digital Knee X-Ray ImagesBreast Cancer(Nickparvar, 2021a)19Digital Knee X-Ray ImagesBreast Cancer(Gornale and Patravali, 2020)19Digital Knee X-Ray ImagesBreast Cancer(Chenzolas)19The vertebrae X-ray imageKree Osteoarthritis batasetKnee Osteoarthritis Dataset<	31	Glaucoma Detection	Glaucoma Classification	(Theng and Das. 2012)
31ICPR-HEp-2Ontaconal Classification(Qi et al., 2027)34SICAPv2Cancer Degree Classification(Silva-Rodríguez et al., 2020)35Blood Cell ImagesBlood Cell Classification(Mooney, 2017)36BreakHisCell type and beginormag(Bukun, 2019)37ChaoyangMulti Classification of pathologists(Zhu et al., 2021a)38HuSHeMSperm Head Morphology Classification(Matek et al., 2021)40NCT-CRC-HE-100KMulti Classification(Matek et al., 2018)41Malignant Lymphoma ClassificationMulti Classification(Cukierski, 2018)42Histopathologic Cancer DetectionCancer Classification(Cukierski, 2018)43LC25000Multi Classification of Lung and Colon(Zhu, 2022)44Brain Tumor 17 ClassesMulti Classification of eye diseases(Orlov et al., 2010a)45Tumor ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Subramanian et al., 2022)47Beta Rez X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)48BUSIBreat Cancer(Al-Dhabyani et al., 2022)49Digital Knee X-Ray ImageVertebrae(Fraiwan et al., 2022)41Fracture detectionFractured Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification Dataset(Jaeger et al., 2014) <tr< td=""><td>32</td><td>AIROGS</td><td>Glaucoma Classification</td><td>(de Vente et al 2023)</td></tr<>	32	AIROGS	Glaucoma Classification	(de Vente et al 2023)
Solar King 2Initial Classification(Given, 200)35Blood Cell ImagesCancer Degree Classification(Silva-Rodríguez et al., 2020)36BreakHisCell type and beginormag(Bukun, 2019)37ChaoyangMulti Classification of pathologists(Zhu et al., 2021a)38HuSHeMSperm Head Morphology Classification(Matek et al., 2021a)39Bone Marrow Cell ClassificationBone Marrow Cell Classification(Matek et al., 2021)40NCT-CRC-HE-100KMulti Classification(Matek et al., 2018)41Malignant Lymphoma ClassificationMulti Classification(Cukierski, 2018)42Histopathologic Cancer DetectionCancer Classification of Lung and Colon(Zhu, 2022)44Brain Tumor 17 ClassesMulti Classification of eye diseases(Orlov et al., 2010a)45Tumor ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)47Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)48BUSIBerest Cancer(Al-Dhabyani et al., 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification of Knee(Fraiwan et al., 2021)51Fracture detectionFractured Classification Dataset(Jaeger et al., 2014)54Shenzhen Chest X-Ray SetCOVID19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia	33	ICPR-HEn-2	Multi Classification	(Oi et al 2016)
51Billond Cell ImagesBlood Cell Classification(Mooney, 2017)36BreakHisCell type and beginormag(Bukun, 2019)37ChaoyangMulti Classification of pathologists(Zhu et al., 2021a)38HuSHeMSperm Head Morphology Classification(Shaker, 2018)39Bone Marrow Cell ClassificationBone Marrow Cell Classification(Matek et al., 2021)40NCT-CRC-HE-100KMulti Classification(Matek et al., 2018)41Malignant Lymphoma ClassificationMulti Classification of Lung and Colon(Cukierski, 2018)42Histopathologic Cancer DetectionCancer Classification of Lung and Colon(Cukierski, 2018)43LC25000Multi Classification of Lung and Colon(Chu, 2022)44Brain Tumor 17 ClassesMulti Classification of eye diseases(Orlov et al., 2010a)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)47Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2024)53Knee Osteoarthritis with severity grading(Chen, 2018)	34	SICAPy2	Cancer Degree Classification	(Silva-Rodríguez et al. 2020)
16Discrete ConstructionConstructionConstruction17ChaoyangMulti Classification of pathologists(Zhu et al., 2021a)18HuSHeMSperm Head Morphology Classification(Matek et al., 2021a)19Bone Marrow Cell Classification(Matek et al., 2021)(Matek et al., 2021)10NCT-CRC-HE-100KMulti Classification(Matek et al., 2018)11Malignant Lymphoma ClassificationMulti Classification(Cukierski, 2018)12Histopathologic Cancer DetectionCancer Classification of Lung and Colon(Zhu, 2022)14Brain Tumor 17 ClassesMulti Classification of Lung and Colon(Nickparvar, 2021a)15Tumor ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)16Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)17Retinal OCT-C8Multi Classification of eye diseases(Orlov et al., 2022)18BUSIBreast Cancer(Al-Dhabyani et al., 2020)19Digital Knee X-Ray ImagesBreast Cancer(Gornale and Patravali, 2020)10Bone Fracture Multi-Region X-ray DataFractured Classification(Batra, 2024)11Fracture detectionFractured Classification(Batra, 2024)12The vertebrae X-ray imageVertebrae(Fraiwan et al., 2021)13Knee Osteoarthritis batasetKnee Osteoarthritis with severity grading(Chen, 2018)14Shenzhen Chest X-Ray SetCOVID and Pneumonia(Asarf and Islam, 2	35	Blood Cell Images	Blood Cell Classification	(Mooney, 2017)
10110	36	BreakHis	Cell type and beginormag	(Bukun, 2019)
38HuSHeMSperm Head Morphology Classification(Shaker, 2018)39Bone Marrow Cell ClassificationBone Marrow Cell Classification(Matek et al., 2021)40NCT-CRC-HE-100KMulti Classification(Matek et al., 2018)41Malignant Lymphoma ClassificationMulti Classification(Orlov et al., 2010a)42Histopathologic Cancer DetectionCancer Classification of Lung and Colon(Zhu, 2022)44Brain Tumor 17 ClassesMulti Classification of eye diseases(Orlov et al., 2010b)45Tumor ClassificationPituitary or Glioma or Meningioma or Notumor(Nickparvar, 2021a)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)47Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(CovID 19, Classification Dataset(Jaeger et al., 2014)54Shenzhen Chest X-Ray SetCOVID 19, Classification Dataset(Jaeger et al., 2020)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Mader, 2017)57	37	Chaoyang	Multi Classification of pathologists	(Zhu et al., 2021a)
39Bone Marrow Cell ClassificationBone Marrow Cell Classification(Matek et al., 2021)40NCT-CRC-HE-100KMulti Classification(Kather et al., 2018)41Malignant Lymphoma ClassificationMulti Classification(Orlov et al., 2010a)42Histopathologic Cancer DetectionCancer Classification(Zhu, 2022)43LC25000Multi Classification(Zhu, 2022)44Brain Tumor 17 ClassesMulti Classification(Feltrin, 2022)45Tumor ClassificationPituitary or Glioma or Meningioma or Notumor(Nickparvar, 2021a)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)47Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification(Nickparvar, 2021b)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Nickparvar, 2021b)52The vertebrae X-ray imageVertebrae(Chen, 2018)54Shenzhen Chest X-Ray DATABASECOVID 19Classification Dataset(Lager et al., 2014)55COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Ashman et al., 2020)56OVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Tabik et al., 2020)57Tuberculosis Chest X-Ray DatabaseMulti Classification of Breast	38	HuSHeM	Sperm Head Morphology Classification	(Shaker, 2018)
40NCT-CRC-HE-100KMulti Classification(Kather et al., 2018)41Malignant Lymphoma ClassificationMulti Classification(Orlov et al., 2010a)42Histopathologic Cancer DetectionCancer Classification of Lung and Colon(Zhu, 2022)43LC25000Multi Classification of Lung and Colon(Zhu, 2022)44Brain Tumor 17 ClassesMulti Classification of Meningioma or Notumor(Nickparvar, 2021a)45Tumor ClassificationPituitary or Glioma or Meningioma or Notumor(Nickparvar, 2021a)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)47Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)54Shenzhen Chest X-Ray SetCOVID19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Tabik et al., 2020)57COVIDGRCOVID19, Classification of Breast(Mader, 2017)58MIASMulti Classification of Breast <td>39</td> <td>Bone Marrow Cell Classification</td> <td>Bone Marrow Cell Classification</td> <td>(Matek et al., 2021)</td>	39	Bone Marrow Cell Classification	Bone Marrow Cell Classification	(Matek et al., 2021)
41Malignant Lymphoma ClassificationMulti Classification(Orlov et al., 2010a)42Histopathologic Cancer DetectionCancer Classification(Cukierski, 2018)43LC25000Multi Classification of Lung and Colon(Zhu, 2022)44Brain Tumor 17 ClassesMulti Classification(Feltrin, 2022)45Tumor ClassificationPituitary or Glioma or Meningioma or Notumor(Nickparvar, 2021a)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Subramanian et al., 2022)48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Datraw et al., 2022)52The vertebrae X-ray imageKnee Osteoarthritis DatasetKnee Osteoarthritis Dataset(CoVID 19, Classification Dataset53Knee Osteoarthritis DatasetCOVID and Pneumonia(Asraf and Islam, 2021)54COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Chowdhury et al., 2020)55COVIDGRCOVID 19, Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Mader, 2017)55Pneumonia Chest X-RayPneumonia Classification of Breast(Mader, 2017)56OP ediatric Pneumonia Chest X-RayPneumonia Classification of Breast(Mader, 2017)57 <td>40</td> <td>NCT-CRC-HE-100K</td> <td>Multi Classification</td> <td>(Kather et al., 2018)</td>	40	NCT-CRC-HE-100K	Multi Classification	(Kather et al., 2018)
42Histopathologic Cancer Detection LC25000Cancer Classification Multi Classification of Lung and Colon Multi Classification(Cukierski, 2018) (Zhu, 2022)44Brain Tumor 17 ClassesMulti Classification(Feltrin, 2022)45Tumor ClassificationPituitary or Glioma or Meningioma or Notumor Multi Classification of eye diseases(Orlov et al., 2010b)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)47Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray Data Fracture detectionFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Chowdhury et al., 2020)57COVIDGRWulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)59Pneumonia Cl	41	Malignant Lymphoma Classification	Multi Classification	(Orlov et al., 2010a)
43LC25000Multi Classification of Lung and Colon Multi Classification(Zhu, 2022)44Brain Tumor 17 ClassesMulti Classification(Feltrin, 2022)45Tumor ClassificationPituitary or Glioma or Meningioma or Notumor Multi Classification of eye diseases(Orlov et al., 2010b)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)47Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID 19, Classification Dataset(Jaeger et al., 2014)55CoVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Chowdhury et al., 2020)56MIASMulti Classification of Breast(Mader, 2017)57Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)59Prediatric Pneumon	42	Histopathologic Cancer Detection	Cancer Classification	(Cukierski, 2018)
44Brain Tumor 17 ClassesMulti Classification(Feltrin, 2022)45Tumor ClassificationPituitary or Glioma or Meningioma or Notumor(Nickparvar, 2021a)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)47Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Tabik et al., 2020)57COVIDGRCOVID19, Classification(Tabik et al., 2020)58MIASMulti Classification(Tabik et al., 2020)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)59Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	43	LC25000	Multi Classification of Lung and Colon	(Zhu, 2022)
45Tumor ClassificationPituitary or Glioma or Meningioma or Notumor(Nickparvar, 2021a)46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)47Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID and Pneumonia(Asraf and Islam, 2021)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID19, Classification(Tabik et al., 2020)57COVIDGRCOVID19, Classification of Breast(Mader, 2017)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	44	Brain Tumor 17 Classes	Multi Classification	(Feltrin, 2022)
46Malignant Lymphoma ClassificationMulti Classification of eye diseases(Orlov et al., 2010b)47Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID 19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID 19, Classification(Tabik et al., 2020)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	45	Tumor Classification	Pituitary or Glioma or Meningioma or Notumor	(Nickparvar, 2021a)
47Retinal OCT-C8Multi Classification of eye diseases(Subramanian et al., 2022)48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID 19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID 19, Classification(Tabik et al., 2020)57COVIDGRCOVID19, Classification of Breast(Mader, 2017)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	46	Malignant Lymphoma Classification	Multi Classification of eye diseases	(Orlov et al., 2010b)
48BUSIBreast Cancer(Al-Dhabyani et al., 2020)49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID 19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID 19, Classification(Tabik et al., 2020)57COVIDGRCOVID19, Classification of Breast(Mader, 2017)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	47	Retinal OCT-C8	Multi Classification of eye diseases	(Subramanian et al., 2022)
49Digital Knee X-Ray ImagesDegree Classification of Knee(Gornale and Patravali, 2020)50Bone Fracture Multi-Region X-ray DataFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID 19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID 19, Classification(Tabik et al., 2020)57COVIDGRCOVID19, Classification of Breast(Mader, 2017)58MIASMulti Classification of Breast(Rahman et al., 2020)59Puberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	48	BUSI	Breast Cancer	(Al-Dhabyani et al., 2020)
50Bone Fracture Multi-Region X-ray Data Fracture detectionFractured Classification(Nickparvar, 2021b)51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Chowdhury et al., 2020)57COVIDGRCOVID19, Classification of Breast(Tabik et al., 2020)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	49	Digital Knee X-Ray Images	Degree Classification of Knee	(Gornale and Patravali, 2020)
51Fracture detectionFractured Classification(Batra, 2024)52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID 19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Chowdhury et al., 2020)57COVIDGRCOVID19, Classification of Breast(Tabik et al., 2020)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	50	Bone Fracture Multi-Region X-ray Data	Fractured Classification	(Nickparvar, 2021b)
52The vertebrae X-ray imageVertebrae(Fraiwan et al., 2022)53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Chowdhury et al., 2020)57COVIDGRCOVID19, Classification(Tabik et al., 2020)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	51	Fracture detection	Fractured Classification	(Batra, 2024)
53Knee Osteoarthritis DatasetKnee Osteoarthritis with severity grading(Chen, 2018)54Shenzhen Chest X-Ray SetCOVID19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Chen, 2018)57COVIDGRCOVID 19, Classification(Chen, 2017)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	52	The vertebrae X-ray image	Vertebrae	(Fraiwan et al., 2022)
54Snenznen Cnest X-Ray SetCOVID19, Classification Dataset(Jaeger et al., 2014)55Chest X-ray PDCOVID and Pneumonia(Asraf and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Chowdhury et al., 2020)57COVIDGRCOVID19, Classification(Tabik et al., 2020)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	53	Knee Osteoarthritis Dataset	Knee Osteoarthritis with severity grading	(Chen, 2018)
55Cnest A-ray PDCOVID and Pneumonia(Asrat and Islam, 2021)56COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Chowdhury et al., 2020)57COVIDGRCOVID19, Classification(Tabik et al., 2020)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	54	Snenzhen Chest X-Ray Set	COVID 19, Classification Dataset	(Jaeger et al., 2014)
30COVID-19 CHEST X-RAY DATABASECOVID and Pneumonia(Chowdnury et al., 2020)57COVIDGRCOVID19, Classification(Tabik et al., 2020)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	55 57	COVID 10 CHEST X DAX DATADAGE	COVID and Pneumonia	(Asrai and Islam, 2021)
57COVID19, Classification(Table et al., 2020)58MIASMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	50 57	COVID-19 CHEST X-KAY DAIABASE	COVID and Pheumonia	(Chowdnury et al., 2020) (Tabile at al., 2020)
JoMulti Classification of Breast(Mader, 2017)59Tuberculosis Chest X-Ray DatabaseTuberculosis(Rahman et al., 2020)60Pediatric Pneumonia Chest X-RayPneumonia Classification(Kermany, 2018)	5/ 50		Wulti Classification of Presst	(1a01K et al., 2020)
60 Pediatric Pneumonia Chest X-Ray Preumonia Classification (Kannan et al., 2020) (Kermany, 2018)	50 50	Tuberculoris Chest Y Pay Database	Tuberculosis	(Rahman et al. 2020)
	60	Pediatric Pneumonia Chest X-Ray	Pneumonia Classification	(Kermany, 2018)

Table 18: The details of the medical datasets are provided

No.	Name	Description	Citation
61	Random Sample of NIH Chest X-Ray Dataset	Multi Classificaiton of Chest	(Wang et al., 2017)
62	CoronaHack-Chest X-Ray	Pnemonia Classifiction with Virus type	(Praveen, 2019)
63	Brain Tumor Dataset	Tumor Classification	(Viradiya, 2020)
64	Fitzpatrick 17k (Nine Labels)	Multi Classification	(Groh et al., 2021)
65	BioMediTech	Multi Classification	(Nanni et al., 2016)
66	Diabetic retinopathy	Diabetic Retinopathy Level	(Benítez et al., 2021)
67	Leukemia	Cancer Classification	(Codella et al., 2019)
68	ODIR-5K	Multiple Labels Classification	(University, 2019)
69	Arthrosis	Bone Age Classification	(Zha, 2021)
70	HSA-NRL	Multi Classification of pathologists	(Zhu et al., 2021b)
71	ISIC 2018 (Task 3)	Multi Classification	(Codella et al., 2019)
72	ISIC 2017 (Task 3)	Multi Classification	(Codella et al., 2018)
73	ChestX-Det	Multi Classification	(Lian et al., 2021)
74	Monkeypox Skin Lesion Dataset	Only Monkeypox	(Ali et al., 2022)
75	Cataract Dataset	Multi Classification	(JR2NGB, 2019)
76	ChestX-rays IndianaUniversity	Multi-label Classification	(Raddar, 2019)
77	CheXpert v1.0 small	Multi-label Classification	(Arevalo, 2020)
78	CBIS-DDSM	Multi Classification	(Lee et al 2017)
79	NI M-TB	Tuberculosis	(Karaca 2022)
80	ChestXray-NIHCC	Multi-label Classification	(Summers and Ronald 2020)
81	COVIDy CXR-4	COVID19 Classification	(Wang et al. 2020)
82	VinDr-Mammo	Multi-label Classification	(Nguyen et al., 2020)
83	PBC dataset normal DIB	Multi Classification	(Acceved o et al. 2020)
84	Human Protein Atlas	Multi-label Classification	(1 ect al 2022)
85	RSNA Pneumonia Detection Challenge 2018	Multi-label Classification	(Anouk Stein et al. 2018)
86	VinDr-SpineXR	Multi Classification of Bones Diseases	(Pham et al. 2021)
87	VinDr-PCXR	Multi-label Classification	(Pham et al. 2021)
88	PH2	Melanoma Segmentation	(Mendonce et al. 2015)
80	ISBI 2016 (Tack 3B)	Melanoma Segmentation	(Gutman et al. 2016)
00	ISIC 2016 (Task 1)	Melanoma Segmentation	(Gutman et al., 2016)
01	ISIC 2010 (Task 1)	Melanoma Segmentation	(Codella et al., 2010)
02	CVC ClinicDR	Polyn Segmentation	(Coucha et al., 2016) (Bernal et al., 2015)
92 02	Vyorin SEC	Polyp Segmentation	(Bernar et al., 2013) (Iba at al. 2020)
95	NVasii-SEU	Surgical Instrument Segmentation	(Maghool et al., 2020)
05	EDD 2020	Multiple Diseases Segmentation in Intestine	(Ali at al. 2020)
95 06	SICADy2	Concer Cells Segmentation	(All et al., 2020) (Silva Podríguaz et al. 2020)
90		Cancer Certs Segmentation	(Hospitalia 2022)
00		Thuroid Nodula Sogmentation	(Gong at al. 2022)
90 00		Lung Segmentation (With left or right)	(Gong et al., 2022)
100	NLW-1D	Spinal V ray Anormaly Detection	(Going et al., 2021)
100	VinDr DCVD	Multinla Disassas Sagmentation in Chest	(Pharm et al., 2021)
101	VIIIDI-PCAR Charty Dat	Multiple Diseases Segmentation in Chest	(Phan et al., 2022)
102	UNV Madiana Cl Tract Incore Secondarian	Suminal Instrument Secondarian	(Lian et al., 2021)
103	Uw-Madison GI Tract Image Segmentation	Surgical Instrument Segmentation	(Lee et al., 2024)
104	Duke Liver Dataset MRI VI	Liver Segmentation	(Macdonald et al., 2020)
105	Duke Liver Dataset NIKI V2	Liver Segmentation	(Ivracdonald et al., 2020)
100	SHIVI-ACK PREUMOTORY Segmentation	Fine the Versilar Segmentation	(Zawacki et al., 2019)
10/	FIVES	Fundus vascular Segmentation	(Jin et al., 2022)
108	KIM-ONE DL	Optic Disc and Cup Segmentation	(Batista et al., 2020)
109	PALM19	Optic Disc Segmentation	(Fu et al., 2019)

Table 19: Continued from Table 18.