

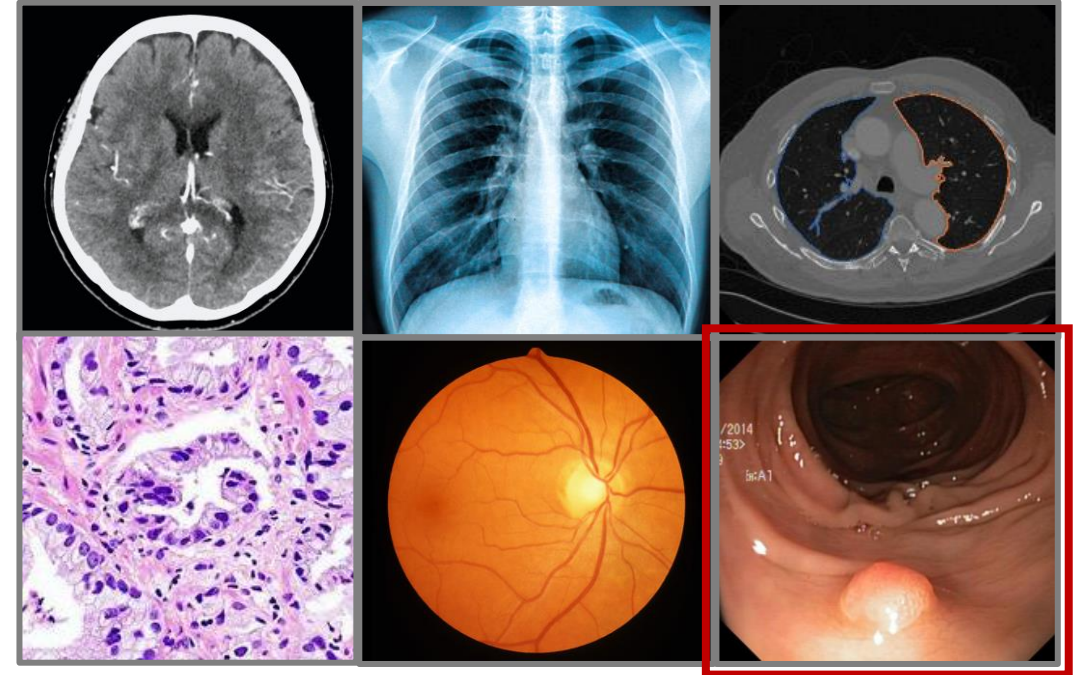
# Evaluation of Multimodal Models: Assessing Performance and Finding Improvements

Metamorphic Testing for Medical Image Analysis

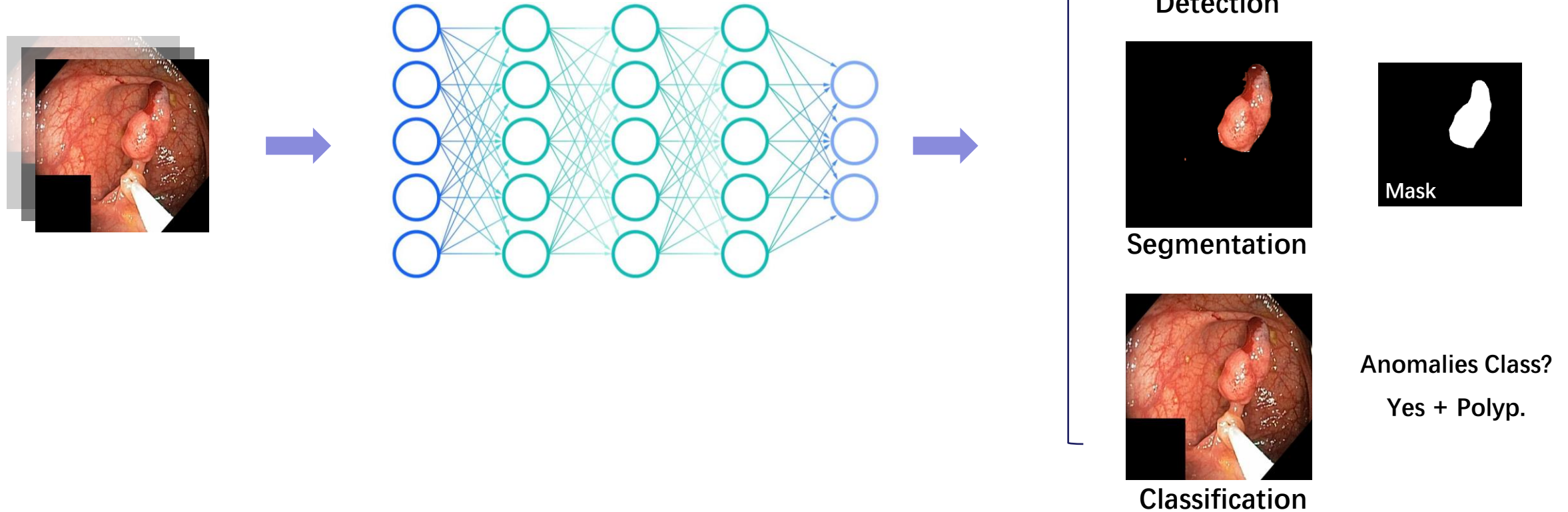
WU, Haoran WU, Yushan

# AI in Medical Imaging

- Medical errors are a critical issue.
- A leading cause is diagnostic errors.
- AI can enhance the accuracy of medical diagnosis tools.

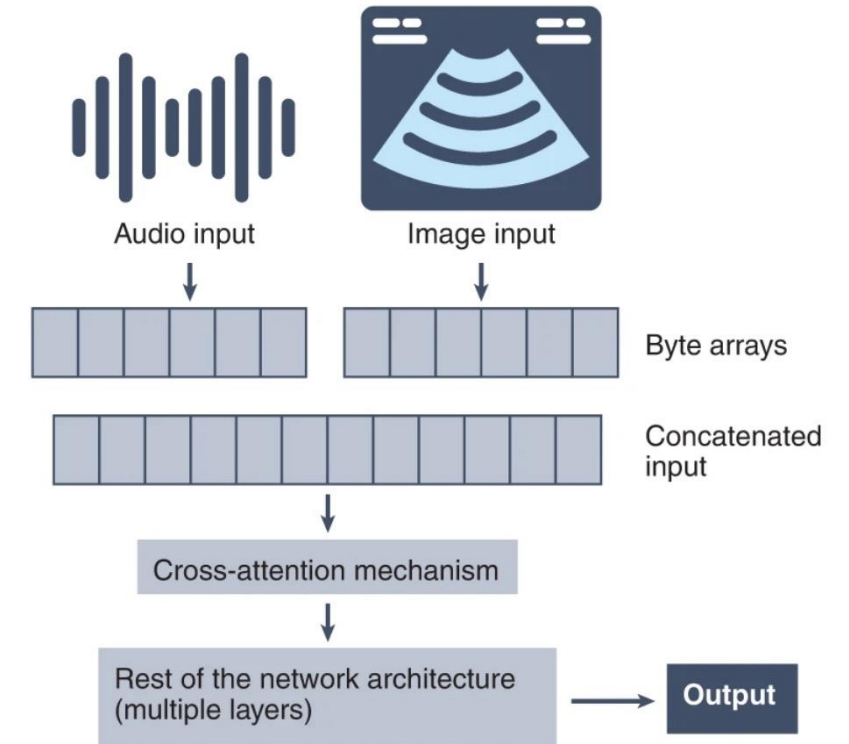


# Common Tasks



# Multimodality and Healthcare

- Multimodal: integration of multiple modes of communication and interaction.
- Application fields: healthcare, education, finance, etc.
- Further improvement in the use of AI in healthcare.



# Motivation

- Fallibility: AI misalignment with clinicians' assessments.
- Need for reliable and robust testing frameworks.
- The methodologies for generating test cases in general computer vision software cannot be directly applied due to the complex nature of medical diagnosis.

# Introducing MedTest

- **MedTest**: A novel metamorphic testing paradigm targeting models on medical imaging tasks.
- Conducted a pilot study, revealing **9 metamorphic relations**, across four artifact categories: **lightness, motion, object artifacts, and non-object artifacts**.
- Testing in both commercial software and state-of-the-art algorithms.
- Further training on those algorithms to improve the model performance.

# Metamorphic Testing

- Key idea: Automatically generate test cases to solve the test oracle problem via Metamorphic Relations (MR).
- MRs delineate the expected relationship between different sets of input-output pairs of a software application.
- Let  $p$  be a representation mapping program inputs into program outputs, and  $f_I$  and  $f_O$  are two functions for transforming the input and output domain, respectively.
- MR formulation:

$$\forall i, p \llbracket f_I(i) \rrbracket = f_O(p \llbracket i \rrbracket)$$



# Metamorphic Testing on AI models

- In our testing scenarios, let *Model* be the model or software we target, that continuously maps each image into predicted output (e.g. segmentation mask).
- Given the original image stream  $\mathbb{I}$ , we can define various image perturbations  $\mathbb{P}$  that simply add some artifacts and do not impact the clinical diagnosis for each image  $i \in \mathbb{I}$ .
- In this way, we use the following MR to test the models with additional perturbations:

$$\forall i \in \mathbb{I} \wedge \forall p \in \mathbb{P}, Model[p(i)] \approx Model[i] \\ |Model[p(i)] - Model[i]| < \varepsilon$$

where  $\varepsilon$  denotes a certain degree of error-tolerant rate.



# Perturbation Types

- Goal: The “seed” image and “perturbed” counterparts should yield consistent prediction results (e.g. classification label, segmentation masks).
- Perturbation criteria: clinical-semantic-preserving, realistic, unambiguous.

Perturbation Group	Type	Description
Lighting	Saturation	Over-saturation caused by excessive lighting
	Contrast	Resulting from underexposure or obstructions in the field of view
	White Balance	Color distortions due to presence of white objects
	Specularity	Reflections resembling a mirror-like surface
Motion	Blur	Blurring from hand movements or rapid camera motion
Objects	Instrument	Presence of surgical instruments in the image frame
	Feces	Incomplete colon cleansing in patients
	Blood	Visible bleeding from wounds
Non-objects	Text	Embedded clinical information related to patients

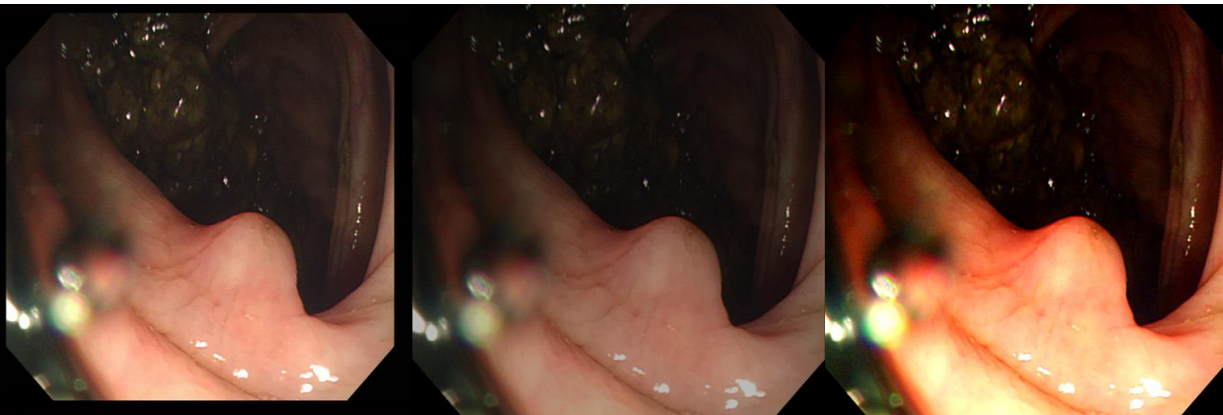
# Contrast/Saturation

- The light source is too far/close to the tissue.
- Applied torchvision.transforms to adjust contrast/saturation with a random factor.

Original

Contrast

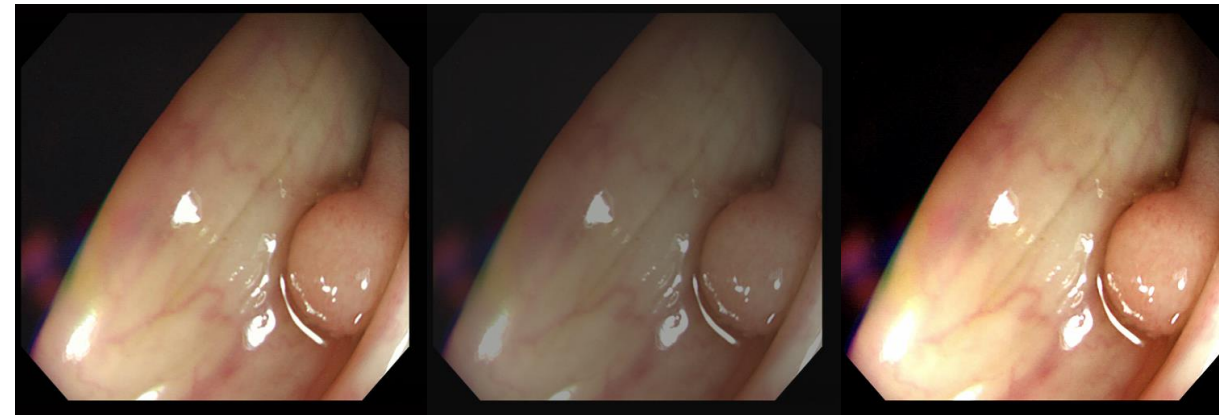
Saturation



Original

Contrast

Saturation



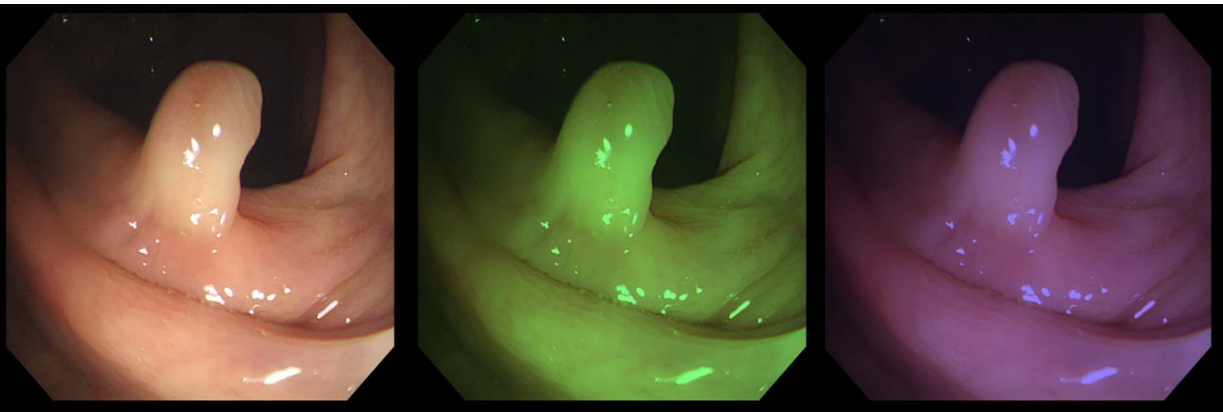
# White Balance

- The white balance settings of the endoscopic camera or the lighting conditions within the endoscopic environment.
- Selectively modified the RGB channels.

Original

Green

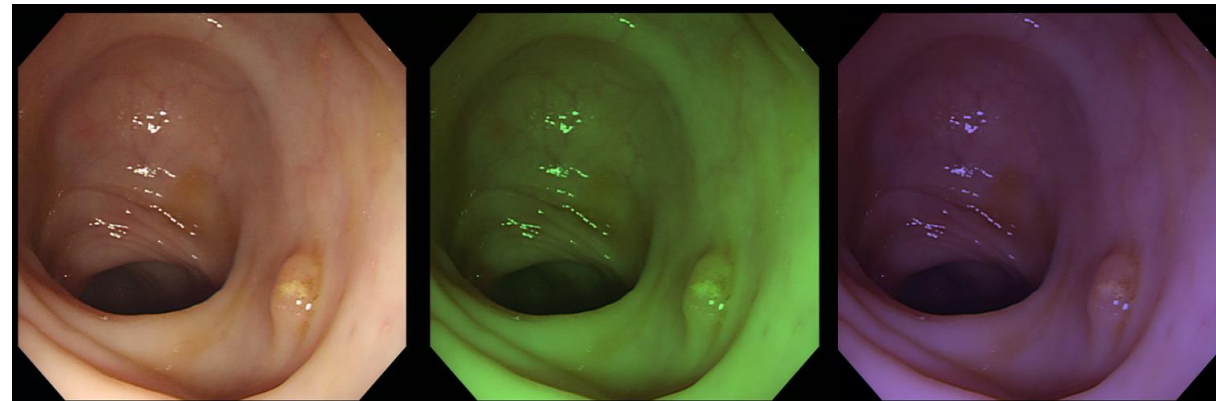
Purple



Original

Green

Purple

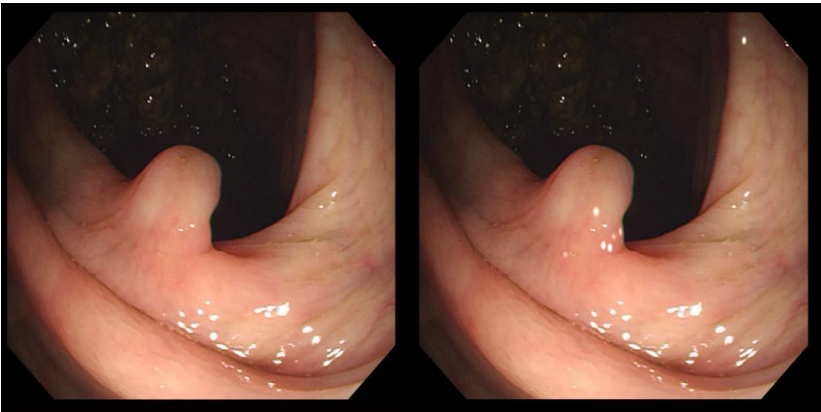


# Specularity

- Resembles the specular reflection.
- Identifying clusters as potential sites, generating ellipses near the cluster centers.
- Integrated these spots with a gray mask and application of Gaussian blur.

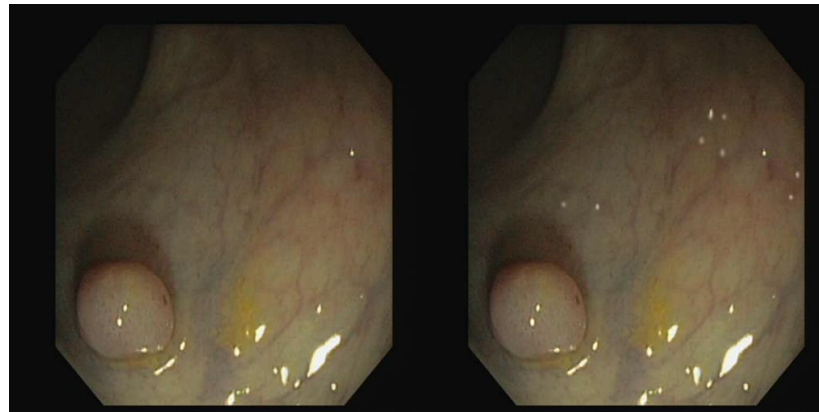
Original

Specularity



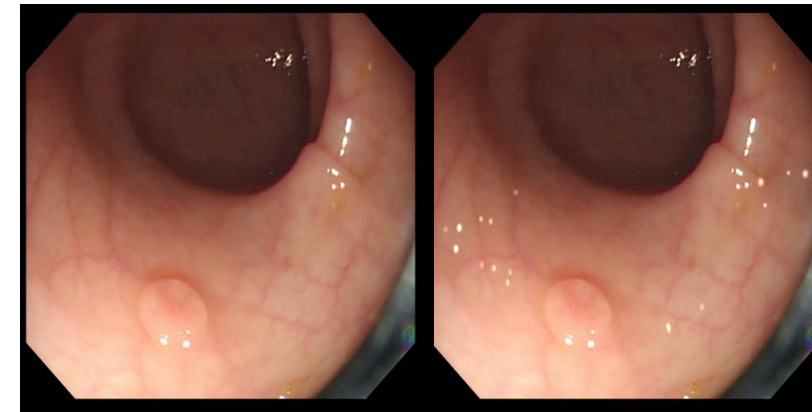
Original

Specularity



Original

Specularity



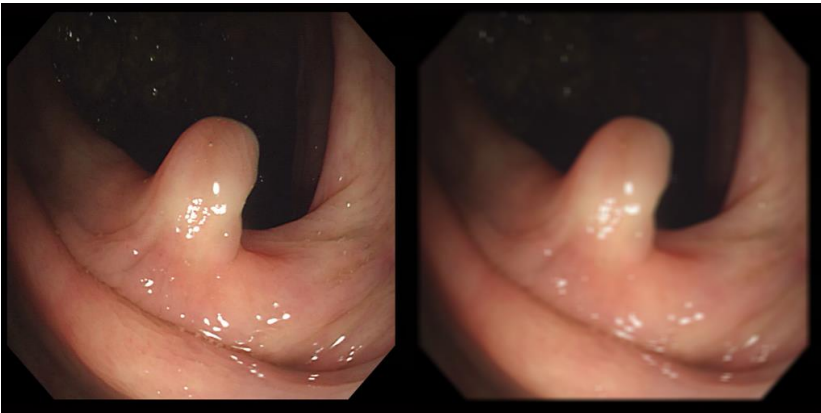


# Motion Blur

- Camera movement and tissue movement.
- Employed Gaussian blur with a random factor.

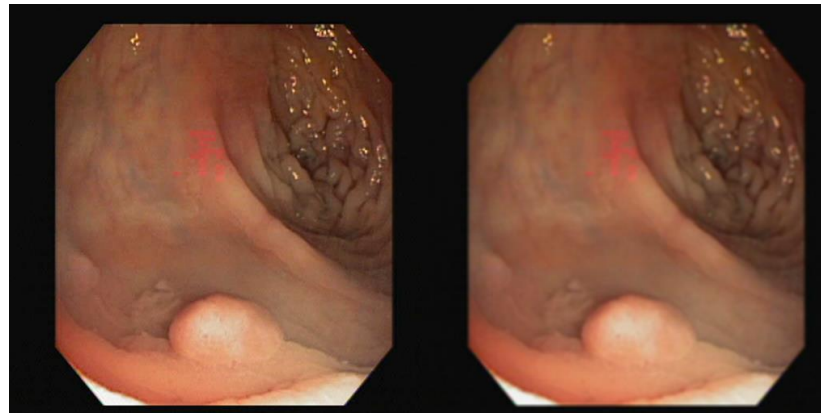
Original

Blur



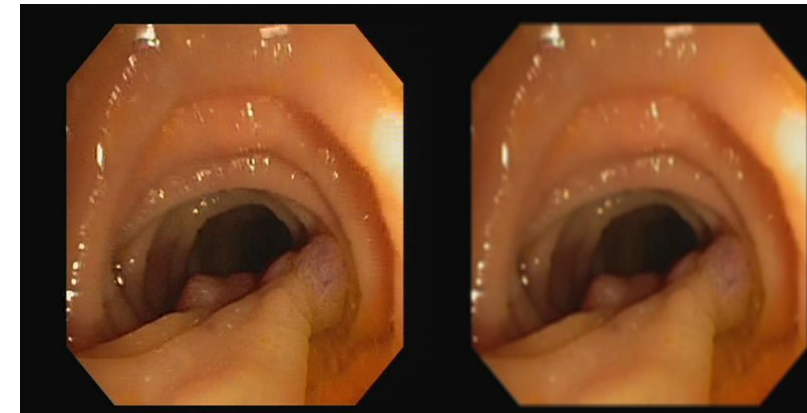
Original

Blur



Original

Blur

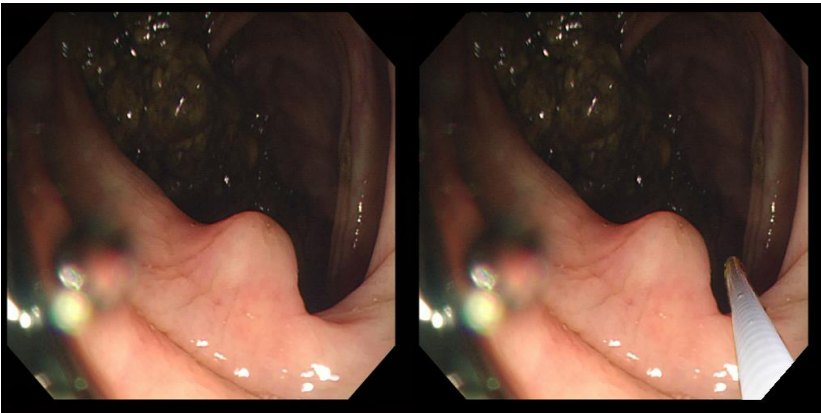


# Instrument

- Resembles the medical instruments that appear in operations.
- Segmented the instrument from the Kvasir-Instrument dataset.
- Utilized our algorithm to select the proper location and orientation and blend the edge.

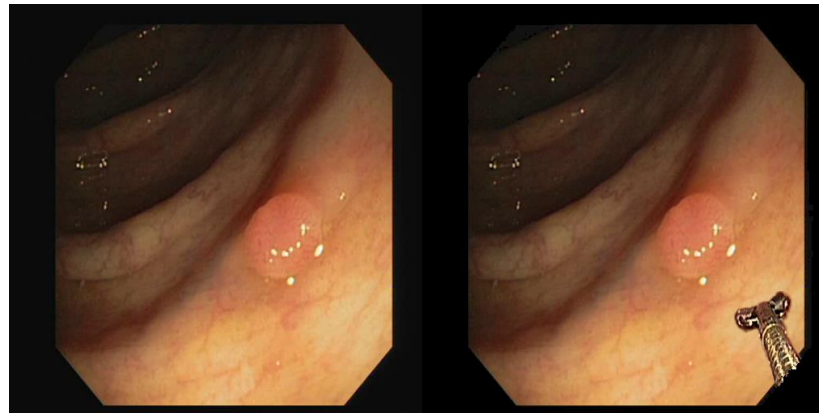
Original

Instrument



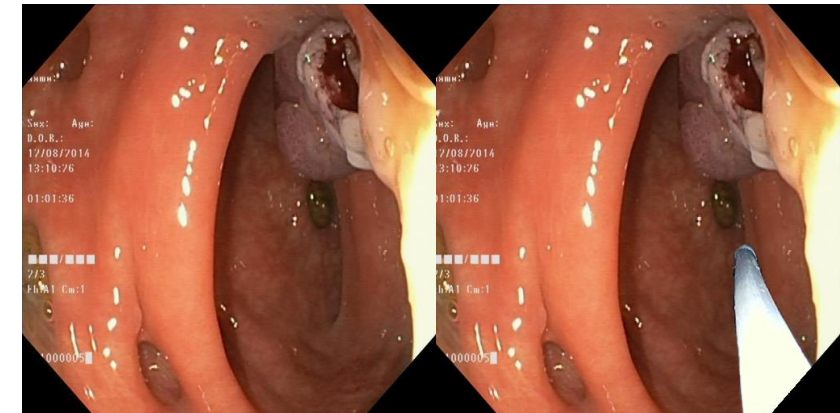
Original

Instrument



Original

Instrument

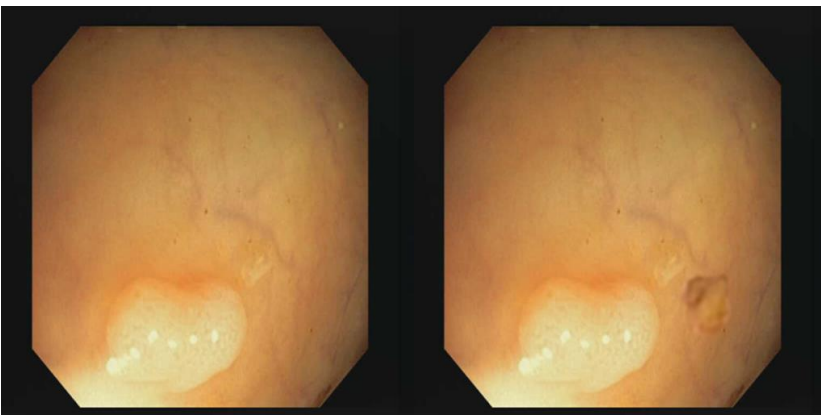


# Feces

- Fecal matter appears in operations.
- Segmented with Meta's Segment Anything from Kvasir dataset.
- Utilized our algorithm to select proper location and calculated size and brightness factor to blend in.

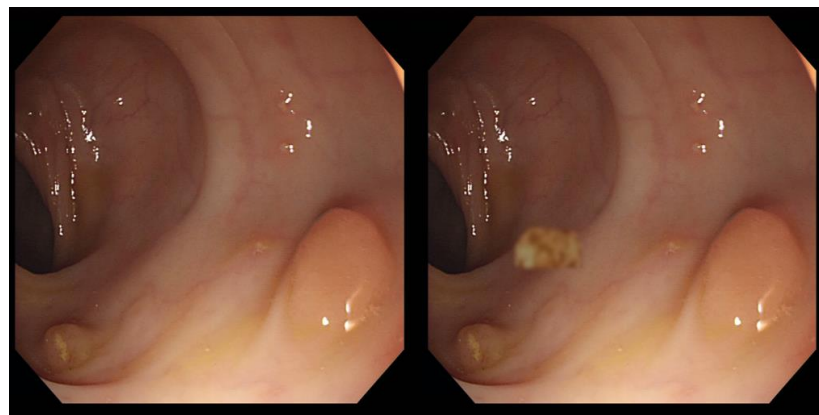
Original

Feces



Original

Feces



Original

Feces



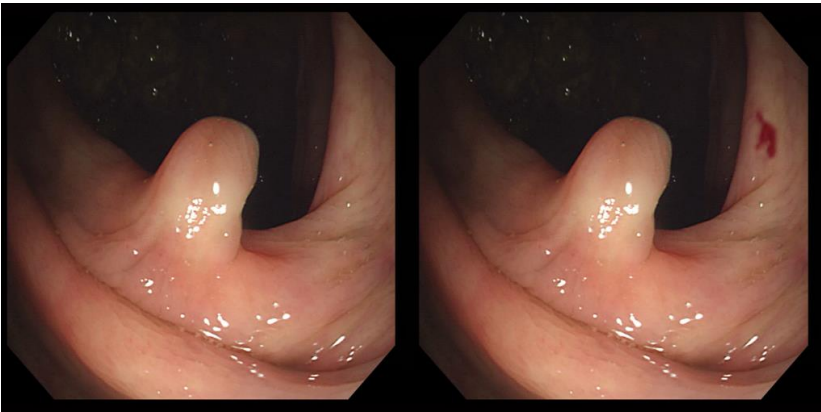


# Blood

- Tissue bleeding in operations.
- Segmented the blood from EAD2020 dataset.
- Utilized our algorithm to select proper location and calculated size and brightness factor to blend in.

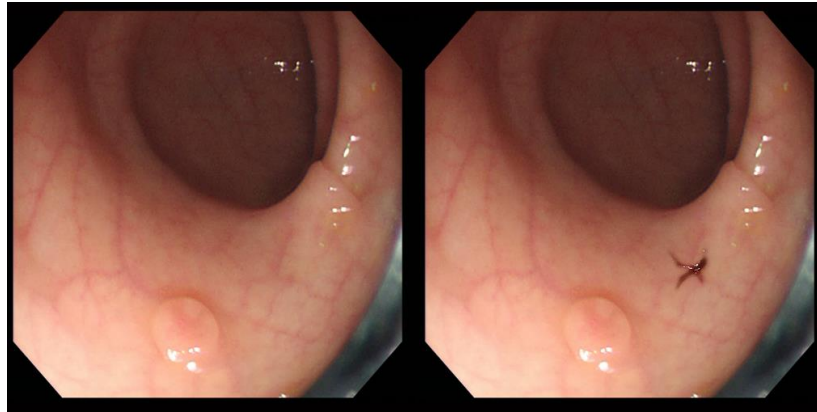
Original

Blood



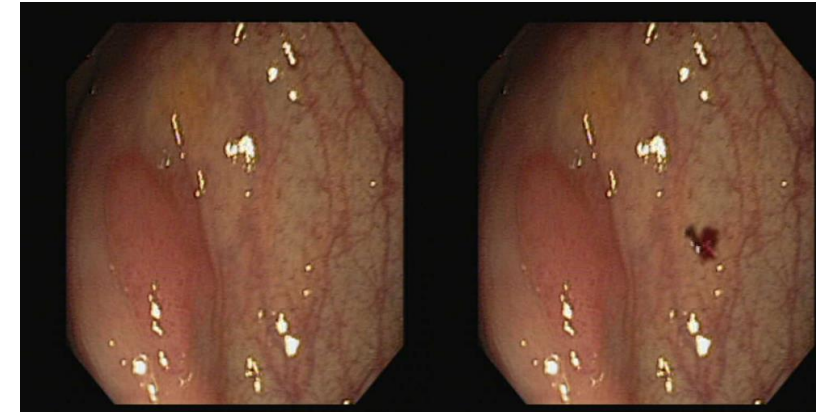
Original

Blood



Original

Blood

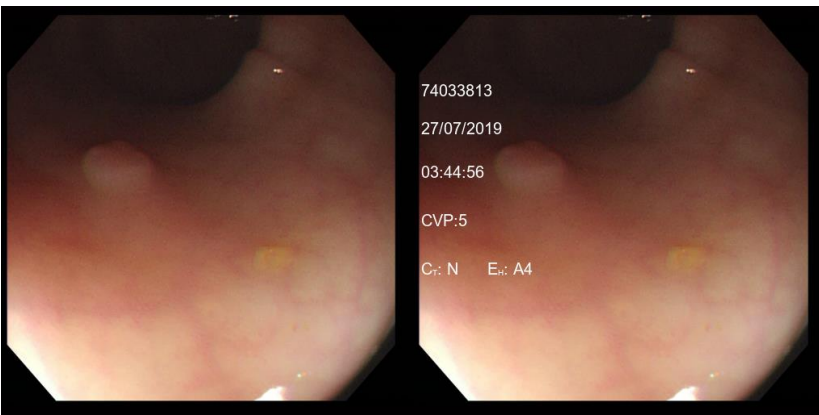


# Text

- Pattern in the text displayed on endoscopic images.
- Used ImageDraw method of PIL to generate text.

Original

Text



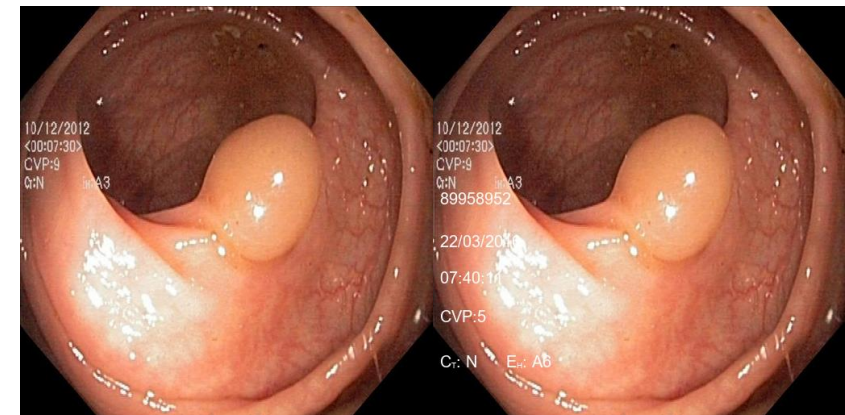
Original

Text



Original

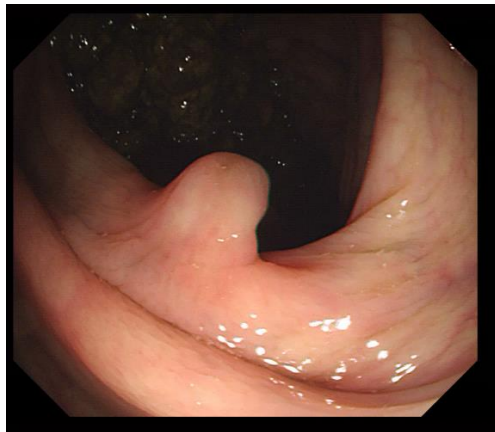
Text



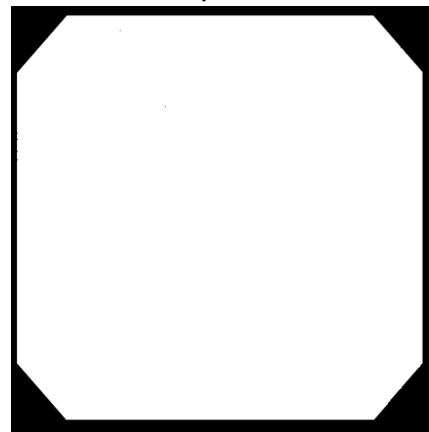
# Data Pre-processing

- Large difference in image sizes -> Resize into  $512 \times 512$ .
- Extract the black frame of images to avoid possible synthesis on the edge.
- Generate gray masks for images to adjust the brightness condition of synthesized parts.

Original



Crop Mask



Gray Mask



# Evaluation

**Evaluate our methodology by answering the following Research Questions (RQ):**

- RQ1: Is our method effective in identifying incorrect outputs produced by medical image diagnosis software and algorithms?
- RQ2: Can the test cases generated by our method be utilized to enhance the performance of medical image diagnosis software?
- RQ3: What are the various factors that influence the performance of our method and how do they do so?

# Evaluation-RQ1

Evaluate our methodology by answering the following Research Questions (RQ):

- **RQ1: Is our method effective in identifying incorrect outputs produced by medical image diagnosis software and algorithms?**
- RQ2: Can the test cases generated by our method be utilized to enhance the performance of medical image diagnosis software?
- RQ3: What are the various factors that influence the performance of our method and how do they do so?

# Experiment Settings

- Mainly utilize clinical **endoscopy images** to evaluate models.
- **Datasets:**
  - CVC-300 (60 images), CVC-ColonDB (380 images) mainly for the segmentation task. 440 seed images in total.
  - Additional ImageCLEF MEDVQA (182 images and 18 questions each) for VQA testing.
  - Wireless capsule endoscopy images from CAD-CAP and KID, with 600 seed images for testing in total.
- **Models under testing:**
  - Polyp **segmentation** models: PraNet, SANet, TGANet, SSFormer.
  - Multi-modal models for **Visual Question-Answering** (VQA): GPT-4V.
  - Gastrointestinal disease **classification**: AGDN, DSI-Net.

# Evaluation-RQ1

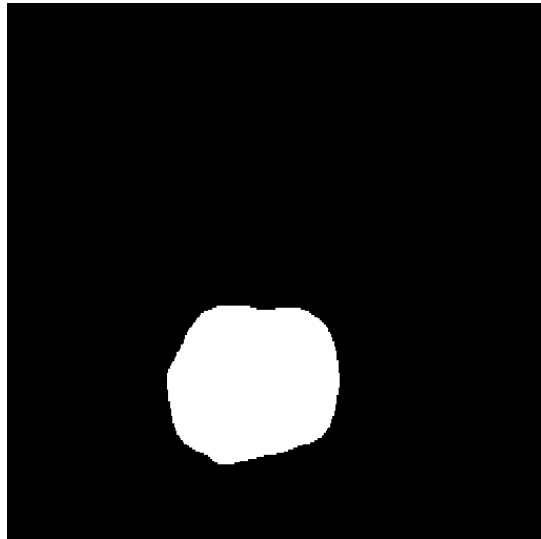
**RQ1: Is our method effective in identifying incorrect outputs produced by medical image diagnosis software and algorithms?**

- **Segmentation**
- Visual-Question Answering
- Classification



# Segmentation

- Segmentation is the task to divide the image into different meaningful regions of interest.



# Evaluation Criteria

## Measurement for segmentation task:

- Dice Score:

$$Dice(\hat{Y}, Y) = \frac{2 \times |\hat{Y} \cap Y|}{|\hat{Y}| + |Y|} = \frac{2 \times TP}{(TP + FP) + (TP + FN)} = \frac{2 \times \text{Area of overlap}}{\text{Total area}} = \frac{2 \times \text{Area of overlap}}{\text{Area of Prediction} + \text{Area of Ground truth}}$$

- Intersection over Union (IoU) Score:

$$IoU(\hat{Y}, Y) = \frac{|\hat{Y} \cap Y|}{|\hat{Y} \cup Y|} = \frac{TP}{TP + FP + FN} = \frac{\text{Area of overlap}}{\text{Area of union}} = \frac{\text{Area of overlap}}{\text{Area of Prediction} + \text{Area of Ground truth} - \text{Area of overlap}}$$

# Evaluation Criteria

## Measurement for “Misclassified”/ “Error”:

- The difference between model’s performance on “seed” image and on perturbations should not exceed an error-tolerant threshold  $t$ .
- Performance is calculated by Dice/IoU Score.
- The sample counts toward an error if

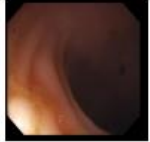
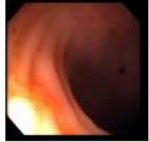
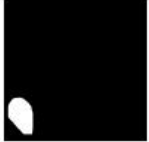
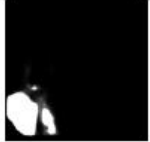

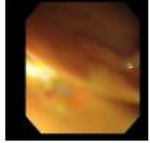
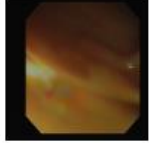


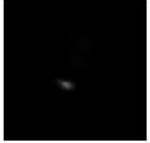

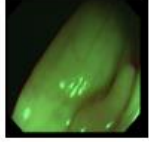


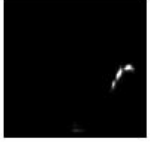






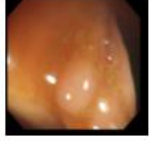
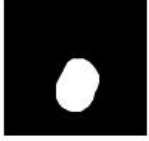


$$\frac{\textit{Original Score} - \textit{Artifact Score}}{\textit{Original Score}} > t$$











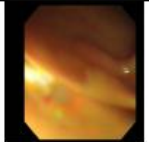
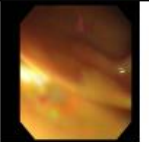



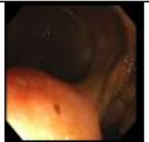
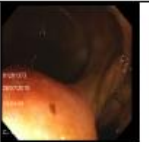



## Error Finding Rate (EFR):

$$EFR = \frac{\# \textit{ of error test cases}}{\# \textit{ of generated test cases}} \times 100\%$$

# Illustration

- We evaluated four segmentation models (PraNet, SANet, TGANet, SSFormer) respectively.

Artifact	Original Image	Image with Artifact	Ground Truth	Output (Original)	Output (Artifact)
Saturation					
Contrast					
White-Balance					
Specularity					
Blur					

Artifact	Original Image	Image with Artifact	Ground Truth	Output (Original)	Output (Artifact)
Instrument					
Feces					
Blood					
Text					

# Results

- For illustration, we choose the error-tolerant threshold  $t = 0.25$ .
- The EFRs are organized by each model, together with separate values for each dataset and perturbations.

PraNet	CVC-300		CVC-ColonDB	
$t=0.25$	Dice	IoU	Dice	IoU
Blood	3.3	6.7	4.0	5.0
Feces	0.0	1.7	7.4	9.2
Instrument	6.7	11.7	12.1	14.0
Spot	1.7	1.7	3.2	4.2
Saturation	8.3	13.3	6.6	8.4
Contrast	1.7	5.0	4.7	6.1
White Balance	8.3	13.3	19.8	22.7
Blur	8.3	8.3	14.2	17.2
Text	0.0	0.0	5.0	5.8

**PraNet: Overall EFR = 6.41%**

SANet	CVC-300		CVC-ColonDB	
$t=0.25$	Dice	IoU	Dice	IoU
Blood	0.0	0.0	2.9	3.4
Feces	1.7	1.7	6.9	7.4
Instrument	1.7	1.7	5.5	5.8
Spot	0.0	0.0	4.2	4.5
Saturation	5.0	6.7	3.4	5.5
Contrast	0.0	0.0	3.4	4.0
White Balance	0.0	0.0	10.8	14.0
Blur	3.3	5.0	6.3	8.7
Text	0.0	0.0	5.5	5.8

**SANet: Overall EFR = 3.37%**

# Results

- For illustration, we choose the error-tolerant threshold  $t = 0.25$ .
- The EFRs are organized by each model, together with separate values for each dataset and perturbations.

TGANet	CVC-300		CVC-ColonDB	
$t=0.25$	Dice	IoU	Dice	IoU
Blood	16.7	20.0	23.9	29.2
Feces	13.3	25.0	13.9	18.2
Instrument	30.0	<b>46.7</b>	18.9	24.2
Spot	3.3	3.3	5.5	6.6
Saturation	16.7	18.3	21.8	24.7
Contrast	0.0	1.7	26.8	29.2
White Balance	<b>31.7</b>	38.3	<b>35.3</b>	<b>40.8</b>
Blur	28.3	31.7	9.7	11.8
Text	8.3	8.3	10.8	13.4

**TGANet: Overall EFR = 17.49%**

SSFormer	CVC-300		CVC-ColonDB	
$t=0.25$	Dice	IoU	Dice	IoU
Blood	3.3	3.3	5.0	5.3
Feces	0.0	0.0	7.6	8.2
Instrument	3.3	6.7	7.1	7.6
Spot	0.0	0.0	2.4	2.4
Saturation	<b>6.7</b>	<b>10.0</b>	2.6	4.5
Contrast	1.7	3.3	3.9	4.7
White Balance	3.3	5.0	<b>11.8</b>	<b>13.9</b>
Blur	0.0	1.7	3.4	3.4
Text	0.0	0.0	2.1	2.6

**SSFormer: Overall EFR = 3.57%**

# Evaluation-RQ1

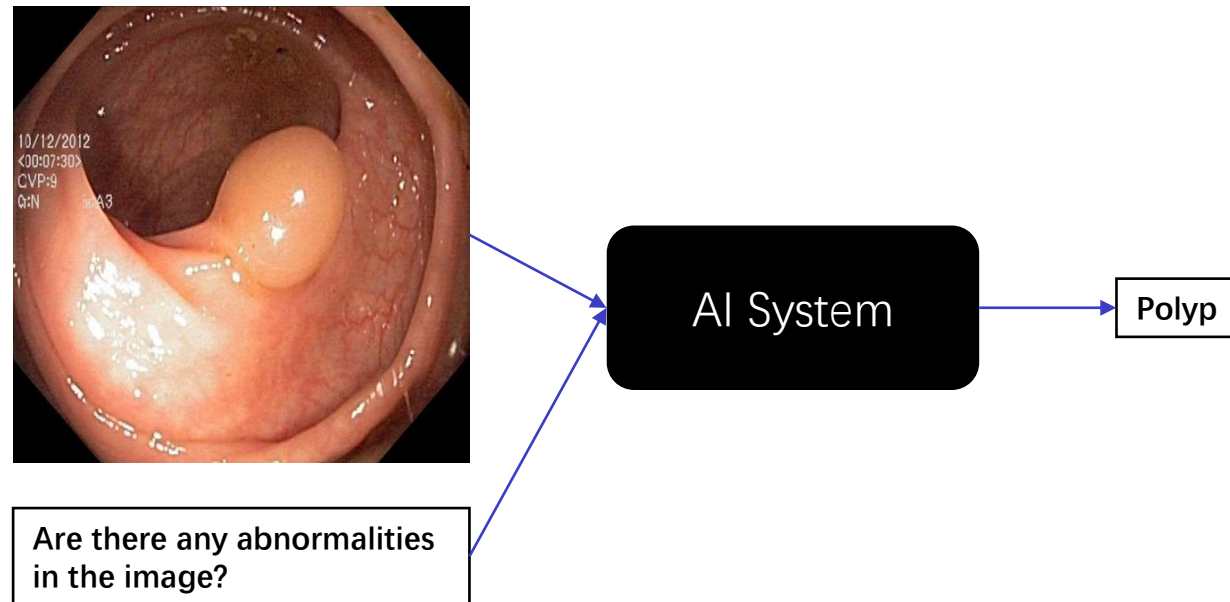
**RQ1: Is our method effective in identifying incorrect outputs produced by medical image diagnosis software and algorithms?**

- Segmentation
- **Visual-Question Answering**
- Classification



# Visual Question-Answering (VQA)

- VQA refers to the task of answering open-ended questions based on an image.
- These questions require an understanding of vision, language, and commonsense knowledge to answer.



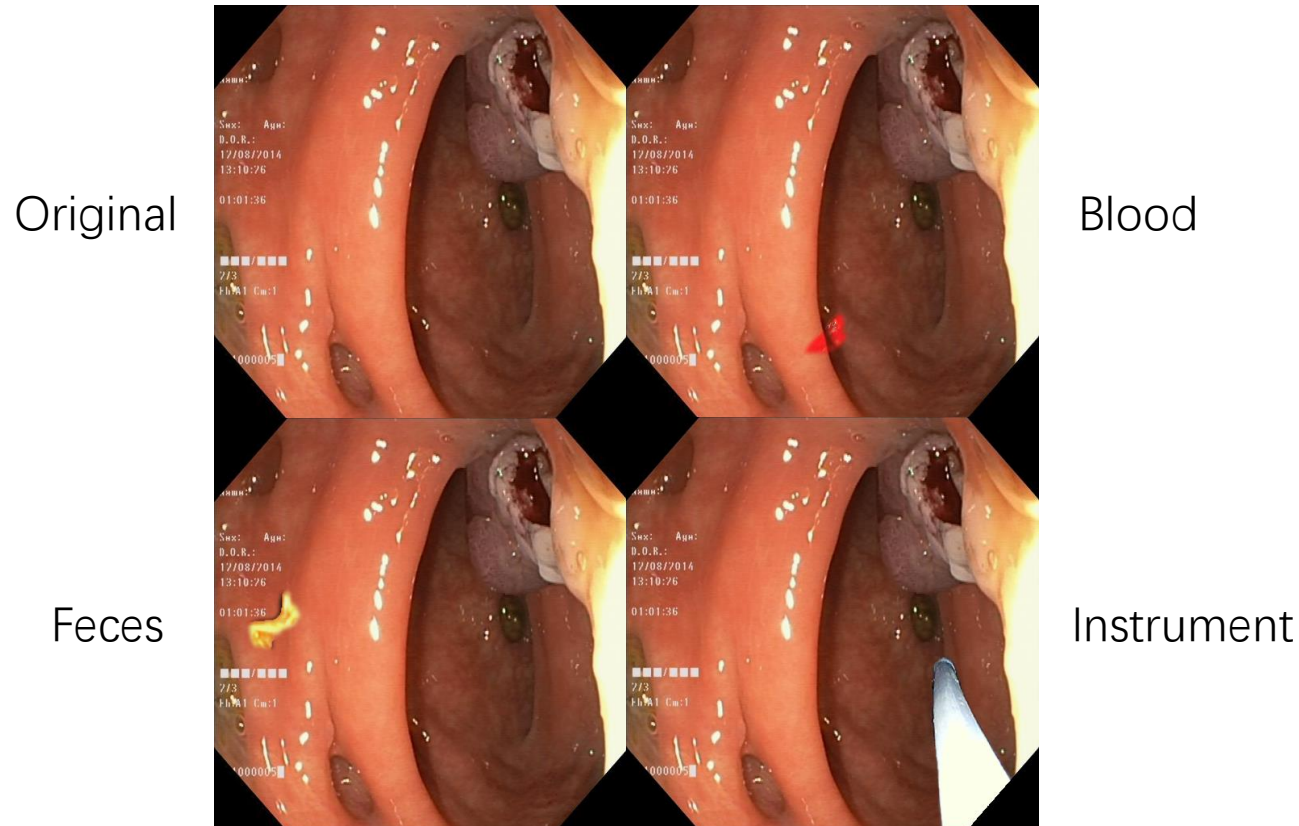
# VQA

- Used the questions provided in the CLEF2023 MEDVQA Dataset.
- Includes “Yes/No” questions and simple questions regarding objects of interest within the medical images.
- Mainly conducted experiments on GPT-4V.

Question Number	Question
1	Are there any abnormalities in the image?
2	Are there any anatomical landmarks in the image?
3	Are there any instruments in the image?
4	Have all polyps been removed?
5	How many findings are present?
6	How many instruments are in the image?
7	How many polyps are in the image?
8	Is there a green/black box artefact?
9	Is there text?
10	Is this finding easy to detect?
11	What color is the abnormality?
12	What color is the anatomical landmark?
13	What is the size of the polyp?
14	What type of polyp is present?
15	What type of procedure is the image taken from?
16	Where in the image is the abnormality?
17	Where in the image is the anatomical landmark?
18	Where in the image is the instrument?

# Illustration

- Illustration on VQA testing case

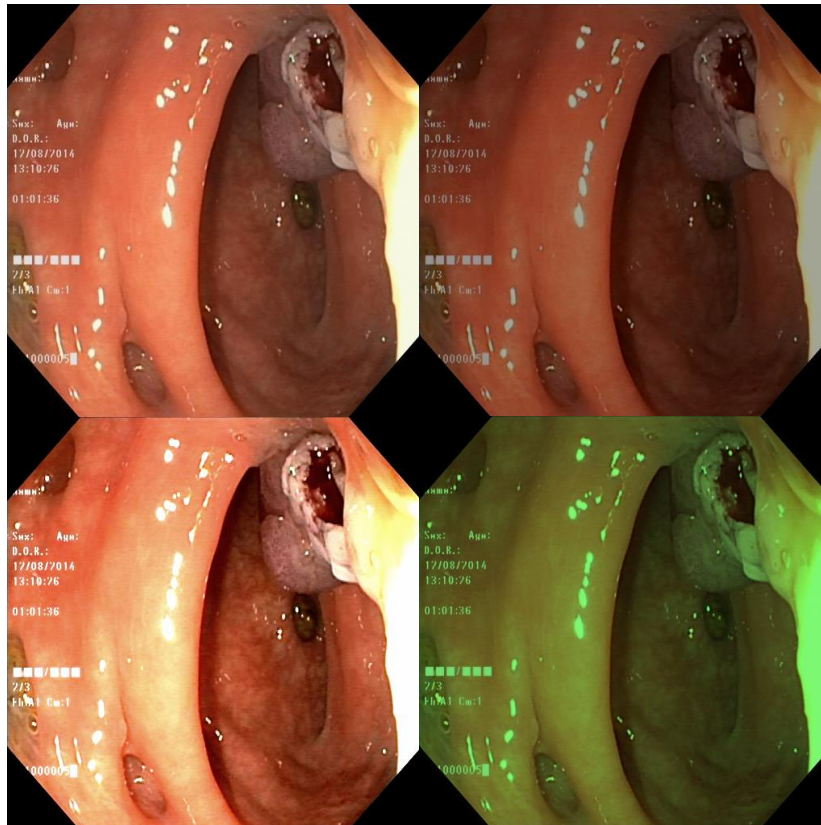


Question	Ground Truth	Original	Blood	Feces	Instrument
Are there any abnormalities in the image?	Polyp	Polyp	Bleeding	Feces	Polyp
Are there any anatomical landmarks in the image?	No	No	No	Yes	Yes
Are there any instruments in the image?	No	No	No	No	Yes
Have all polyps been removed?	No	No	Not relevant	Not relevant	No
How many findings are present?	1	1	1	1	1
How many instruments are in the image?	0	0	0	0	1
How many polyps are in the image?	1	1	0	0	1
Is there a green/black box artefact?	No	No	No	No	No
Is there text?	Yes	Yes	Yes	Yes	Yes
Is this finding easy to detect?	Yes	Yes	Yes	Yes	Yes
What color is the abnormality?	Red, Pink, Grey	Red	Red	Brown	Red
What color is the anatomical landmark?	Not relevant	Not relevant	Not relevant	Pink	Pink
What is the size of the polyp?	>20mm	>10mm	Not relevant	Not relevant	>10mm
What type of polyp is present?	Paris is	Paris Ip	Not relevant	Not relevant	Paris Ip
What type of procedure is the image taken from?	Colonoscopy	Colonoscopy	Colonoscopy	Colonoscopy	Colonoscopy
Where in the image is the abnormality?	Center, Upper-right, Center-right, Upper-center	Center-Left	Center-Left	Bottom-Center	Center-Left
Where in the image is the anatomical landmark?	Not relevant	Not relevant	Not relevant	Center	Center
Where in the image is the instrument?	Not relevant	Not relevant	Not relevant	Not relevant	Bottom-Center

# Illustration

- Illustration on VQA testing case

Original



Contrast

Saturation

White Balance

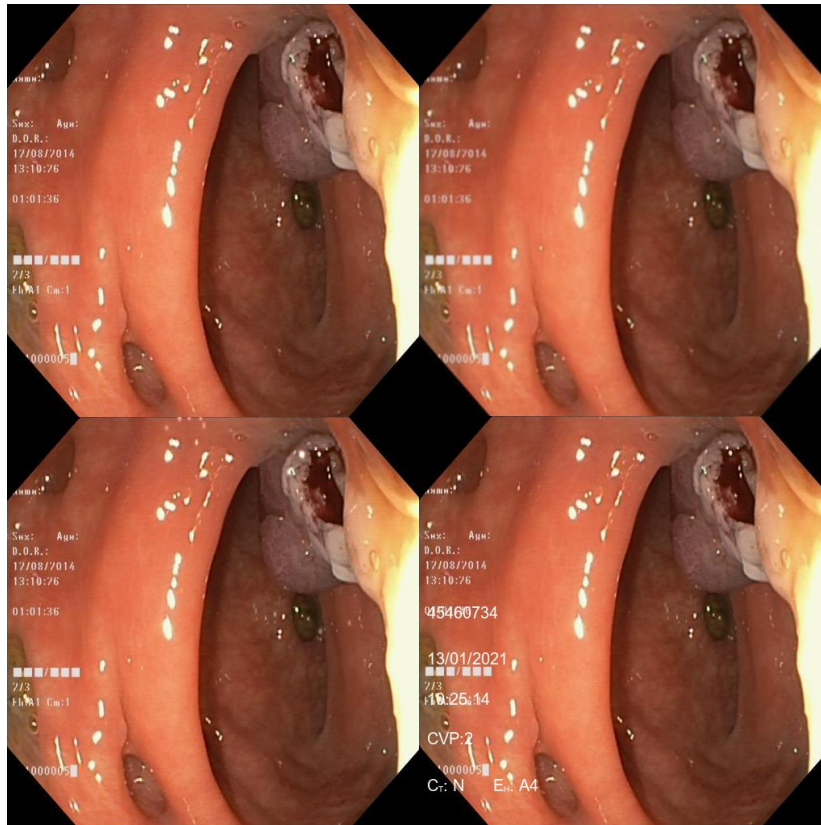
Question	Ground Truth	Original	Contrast	Saturation	White Balance
Are there any abnormalities in the image?	Polyp	Polyp	Polyp	Polyp	Polyp
Are there any anatomical landmarks in the image?	No	No	Yes	Yes	Yes
Are there any instruments in the image?	No	No	No	No	No
Have all polyps been removed?	No	No	No	No	No
How many findings are present?	1	1	1	1	1
How many instruments are in the image?	0	0	0	0	0
How many polyps are in the image?	1	1	1	1	1
Is there a green/black box artefact?	No	No	No	No	No
Is there text?	Yes	Yes	Yes	Yes	Yes
Is this finding easy to detect?	Yes	Yes	Yes	Yes	Difficult due to color alteration
What color is the abnormality?	Red, Pink, Grey	Red	Red	Red	Not applicable due to WB
What color is the anatomical landmark?	Not relevant	Not relevant	Pink	Pink	Not applicable due to WB
What is the size of the polyp?	>20mm	>10mm	>10mm	>10mm	Not applicable due to WB
What type of polyp is present?	Paris is	Paris Ip	Paris Ip	Paris Ip	Not applicable due to WB
What type of procedure is the image taken from?	Colonoscopy	Colonoscopy	Colonoscopy	Colonoscopy	Colonoscopy
Where in the image is the abnormality?	Center, Upper-right, Center-right, Upper-center	Center-Left	Center-Left	Center-Left	Not applicable due to WB
Where in the image is the anatomical landmark?	Not relevant	Not relevant	Center	Center	Not applicable due to WB
Where in the image is the instrument?	Not relevant	Not relevant	Not relevant	Not relevant	Not relevant



# Illustration

- Illustration on VQA testing case

Original



Specularity

Blur

Text

Question	Ground Truth	Original	Blur	Specularity	Text
Are there any abnormalities in the image?	Polyp	Polyp	Polyp	Polyp	Polyp
Are there any anatomical landmarks in the image?	No	No	Yes	Yes	Yes
Are there any instruments in the image?	No	No	No	No	No
Have all polyps been removed?	No	No	No	No	No
How many findings are present?	1	1	1	1	1
How many instruments are in the image?	0	0	0	0	0
How many polyps are in the image?	1	1	1	1	1
Is there a green/black box artefact?	No	No	No	No	No
Is there text?	Yes	Yes	Yes	Yes	Yes
Is this finding easy to detect?	Yes	Yes	No	Yes	Yes
What color is the abnormality?	Red, Pink, Grey	Red	Red	Red	Red
What color is the anatomical landmark?	Not relevant	Not relevant	Pink	Pink	Pink
What is the size of the polyp?	>20mm	>10mm	>10mm	>10mm	>10mm
What type of polyp is present?	Paris is	Paris Ip	Paris Ip	Paris Ip	Paris Ip
What type of procedure is the image taken from?	Colonoscopy	Colonoscopy	Colonoscopy	Colonoscopy	Colonoscopy
Where in the image is the abnormality?	Center, Upper-right, Center-right, Upper-center	Center-Left	Center-Left	Center-Left	Center-Left
Where in the image is the anatomical landmark?	Not relevant	Not relevant	Center	Center	Center
Where in the image is the instrument?	Not relevant	Not relevant	Not relevant	Not relevant	Not relevant

# Results

- The VQA testing result on GPT-4V.
- The model is quite robust.

GPT-4V	Original	Average Perturbation	Difference (Original - Average)
Are there any abnormalities in the image?	0.888	0.861	0.027
Are there any anatomical landmarks in the image?	0.341	0.347	-0.006
Are there any instruments in the image?	0.947	0.873	0.074
Have all polyps been removed?	0.284	0.282	0.002
How many findings are present?	0.835	0.816	0.019
How many instrumnets are in the image?	0.963	0.93	0.033
How many polyps are in the image?	0.79	0.763	0.027
Is there a green/black box artefact?	0.624	0.636	-0.012
Is there text?	0.98	0.963	0.017
Is this finding easy to detect?	0.594	0.589	0.005
What color is the abnormality?	0.465	0.406	0.059
What color is the anatomical landmark?	0.369	0.388	-0.019
What is the size of the polyp?	0.21	0.206	0.004
What type of polyp is present?	0.176	0.163	0.013
What type of procedure is the image taken from?	0.976	0.964	0.012
Where in the image is the abnormality?	0.653	0.644	0.009
Where in the image is the anatomical landmark?	0.365	0.386	-0.021
Where in the image is the instrument?	0.924	0.896	0.028
Average	0.632	0.617	0.015

**Overall score before deletion**

# Results

- The VQA testing result on GPT-4V.
- We found that some questions are ambiguous and not related to our task.

GPT-4V	Original	Average Perturbation	Difference (Original - Average)
Are there any abnormalities in the image?	0.888	0.861	0.027
Are there any instruments in the image?	0.947	0.873	0.074
Have all polyps been removed?	0.284	0.282	0.002
How many findings are present?	0.835	0.816	0.019
How many instrumnets are in the image?	0.963	0.93	0.033
How many polyps are in the image?	0.79	0.763	0.027
Is there a green/black box artefact?	0.624	0.636	-0.012
Is there text?	0.98	0.963	0.017
What color is the abnormality?	0.465	0.406	0.059
What is the size of the polyp?	0.21	0.206	0.004
What type of polyp is present?	0.176	0.163	0.013
What type of procedure is the image taken from?	0.976	0.964	0.012
Where in the image is the abnormality?	0.653	0.644	0.009
Where in the image is the instrument?	0.924	0.896	0.028
Average	0.694	0.672	0.022

**Overall score after deletion**



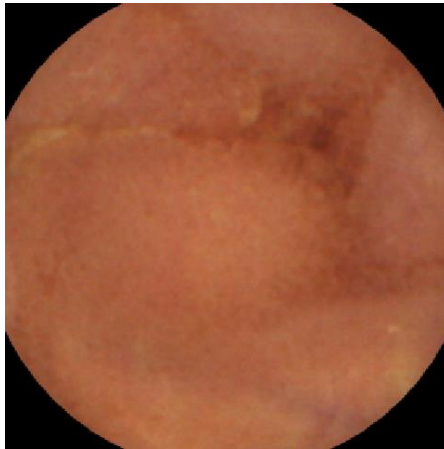
# Evaluation-RQ1

**RQ1: Is our method effective in identifying incorrect outputs produced by medical image diagnosis software and algorithms?**

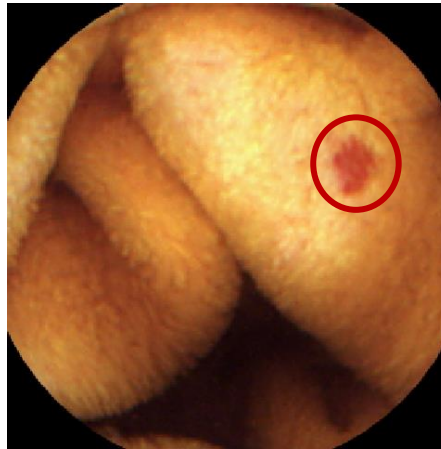
- Segmentation
- Visual-Question Answering
- **Classification**

# Classification

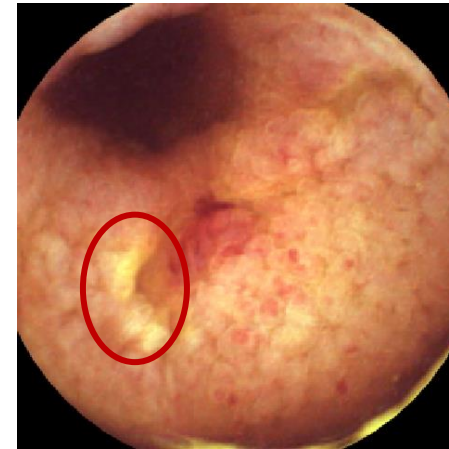
- A medical image analysis technique that involves classifying medical images into different categories based on the type of image or the presence of specific structures or diseases.
- Our Task: Classify wireless capsule endoscopy (WCE) images into three types – Normal, Vascular Lesion, Inflammatory



Normal



Vascular Lesion



Inflammatory

# Evaluation Criteria

## Measurement for classification task:

- **Accuracy:** the proportion of accurately classified samples to the total number of test cases
- **F1 Score:** a weighted harmonic mean of Precision and Recall normalized between 0 and 1

$$F1\ Score = \frac{2}{\frac{1}{Precision} + \frac{1}{Recall}} = \frac{2TP}{2TP + FP + FN}$$

- **Cohen's Kappa Score:** measures the proximity of the predicted classes to the actual classes when compared to a random classification

# Results

- The performance of AGDN and DSI-Net on different perturbations, respectively.
- We can find a maximum decrease in the accuracy of **0.229** and **0.144** in the two models respectively.

Artifact	Accuracy	Cohen's Kappa	F1 Score
Original	0.893	0.836	0.893
Blur	0.702	0.517	0.660
Contrast	0.747	0.602	0.735
Feces	0.797	0.682	0.790
Instrument	0.852	0.773	0.852
Saturation	0.685	0.516	0.682
Spot	0.817	0.712	0.811
Text	0.828	0.733	0.825
White Balance	0.532	0.226	0.463
Average	0.745	0.595	0.727

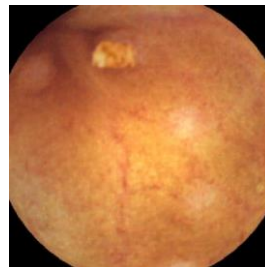
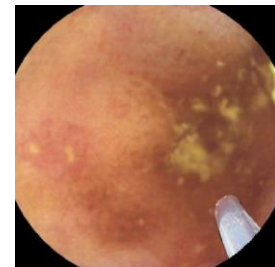
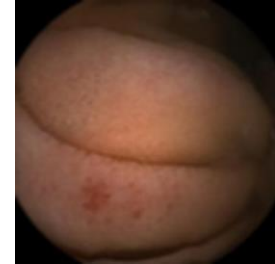
AGDN

Artifact	Accuracy	Cohen's Kappa	F1 Score
Original	0.940	0.908	0.940
Blur	0.897	0.841	0.896
Contrast	0.883	0.823	0.884
Feces	0.907	0.858	0.907
Instrument	0.897	0.843	0.897
Saturation	0.755	0.635	0.757
Spot	0.932	0.895	0.931
Text	0.908	0.859	0.908
White Balance	0.728	0.574	0.711
Average	0.863	0.791	0.861

DSI-Net

# Perturbation Analysis

- Lighting conditions (e.g. white balance, saturation) can trigger most errors.
  - Possible explanation: 1. Edges of critical regions become vague  
2. Sensitive to color
- Motion Blurring can also lead to some corner cases.
  - Possible explanation: Edges of critical regions become vague
- Object-related perturbations resulted in some misleading cases.
  - Possible explanation: 1. Unseen elements (e.g., Instruments) in the training data  
2. Feces resemble some objects of interest



# Answer-RQ1

Evaluate our methodology by answering the following Research Questions (RQ):

**RQ1: Is our method effective in identifying incorrect outputs produced by medical image diagnosis software and algorithms?**

**Answer to RQ1:** MedTest obtained up to 17.49% EFR when testing on segmentation models, qualitatively affected VQA models' performances, and reduced the accuracy of classification models for up to 16.6% in average, indicating that MedTest can effectively discover corner cases and be used for further testing the robustness of other models.

# Evaluation-RQ2

Evaluate our methodology by answering the following Research Questions (RQ):

- RQ1: Is our method effective in identifying incorrect outputs produced by medical image diagnosis software and algorithms?
- **RQ2: Can the test cases generated by our method be utilized to enhance the performance of medical image diagnosis software?**
- RQ3: What are the various factors that influence the performance of our method and how do they do so?

# Evaluation-RQ2

**RQ2: Can the test cases generated by our method be utilized to enhance the performance of medical image diagnosis software?**

- **Segmentation**
- Classification



# Further Training Settings

- We have further trained the models with perturbed images synthesized by MedTest to improve the performances of those models.
- Datasets Construction:
  - CVC-ClinicDB (612 images) and Kvasir (1000 images)
  - Besides the original training set, we also randomly selected an equal number of perturbed images generated from the original training set.
- We used the same training settings as stated in their original papers.

# Results

- We used Dice and IoU scores directly for comparison.
- “Blur” has an over 10% increase across both datasets and metrics.

PraNet	CVC-300				CVC-ColonDB			
	Dice Score		IoU Score		Dice Score		IoU Score	
	Before	After	Before	After	Before	After	Before	After
Original	0.86	<b>0.889</b>	0.777	<b>0.817</b>	0.696	<b>0.701</b>	0.619	<b>0.629</b>
Saturation	0.818	<b>0.878</b>	0.735	<b>0.806</b>	0.677	<b>0.711</b>	0.598	<b>0.636</b>
White Balance	0.815	<b>0.879</b>	0.74	<b>0.805</b>	0.624	<b>0.709</b>	0.551	<b>0.634</b>
Contrast	0.861	<b>0.887</b>	0.777	<b>0.815</b>	0.692	<b>0.707</b>	0.618	<b>0.633</b>
Spot	0.849	<b>0.885</b>	0.764	<b>0.812</b>	0.694	<b>0.702</b>	0.615	<b>0.628</b>
Blur	0.706	<b>0.872</b>	0.619	<b>0.794</b>	0.582	<b>0.695</b>	0.491	<b>0.618</b>
Text	0.744	<b>0.819</b>	0.659	<b>0.746</b>	0.629	<b>0.663</b>	0.554	<b>0.593</b>
Instrument	0.812	<b>0.879</b>	0.717	<b>0.805</b>	0.653	<b>0.699</b>	0.571	<b>0.626</b>
Blood	0.843	<b>0.891</b>	0.751	<b>0.821</b>	0.678	<b>0.696</b>	0.599	<b>0.624</b>
Feces	0.838	<b>0.878</b>	0.75	<b>0.804</b>	0.674	<b>0.691</b>	0.592	<b>0.619</b>
Average	0.815	<b>0.876</b>	0.729	<b>0.803</b>	0.66	<b>0.697</b>	0.581	<b>0.624</b>

PraNet: average improvement 4.9%

# Results

- We used Dice and IoU score directly for comparison.
- “Text” achieving up to a 6.6% increase in the IoU score on CVC-300.

SAnet	CVC-300				CVC-ColonDB			
	Dice Score		IoU Score		Dice Score		IoU Score	
	Before	After	Before	After	Before	After	Before	After
Original	0.898	<b>0.904</b>	0.828	<b>0.836</b>	0.757	<b>0.763</b>	0.677	<b>0.69</b>
Saturation	0.879	<b>0.882</b>	0.807	<b>0.816</b>	<b>0.766</b>	0.757	0.683	0.683
White Balance	0.877	<b>0.889</b>	0.805	<b>0.822</b>	0.738	<b>0.76</b>	0.655	<b>0.687</b>
Contrast	<b>0.898</b>	0.897	0.827	<b>0.829</b>	0.754	<b>0.765</b>	0.673	<b>0.691</b>
Spot	0.899	<b>0.904</b>	0.828	<b>0.837</b>	0.754	<b>0.76</b>	0.675	<b>0.686</b>
Blur	0.851	<b>0.899</b>	0.773	<b>0.831</b>	0.735	<b>0.773</b>	0.646	<b>0.697</b>
Text	0.803	<b>0.864</b>	0.727	<b>0.793</b>	0.694	<b>0.749</b>	0.621	<b>0.674</b>
Instrument	0.898	<b>0.903</b>	0.829	<b>0.836</b>	0.747	<b>0.749</b>	0.668	<b>0.679</b>
Blood	0.899	<b>0.903</b>	0.829	<b>0.835</b>	0.753	<b>0.762</b>	0.675	<b>0.689</b>
Feces	0.901	<b>0.904</b>	0.831	<b>0.837</b>	0.743	<b>0.759</b>	0.665	<b>0.687</b>
Average	0.88	<b>0.895</b>	0.808	<b>0.827</b>	0.744	<b>0.76</b>	0.664	<b>0.686</b>

**SAnet: average improvement 1.6%**

# Results

- We used Dice and IoU score directly for comparison.
- “White Balance” achieving up to a 3.7% increase in the IoU score on CVC-ColonDB

SSFormer	CVC-300				CVC-ColonDB			
	Dice Score		IoU Score		Dice Score		IoU Score	
	Before	After	Before	After	Before	After	Before	After
Original	0.891	0.891	0.825	<b>0.827</b>	0.774	0.774	0.698	<b>0.700</b>
Saturation	0.841	<b>0.876</b>	0.779	<b>0.806</b>	0.778	0.778	0.699	<b>0.702</b>
White Balance	<b>0.880</b>	0.874	<b>0.813</b>	0.811	0.731	<b>0.764</b>	0.656	<b>0.693</b>
Contrast	<b>0.883</b>	0.882	0.817	<b>0.817</b>	0.765	<b>0.774</b>	0.689	<b>0.700</b>
Spot	0.892	0.892	0.826	<b>0.828</b>	0.770	<b>0.775</b>	0.695	<b>0.701</b>
Blur	0.883	<b>0.893</b>	0.813	<b>0.825</b>	0.766	0.766	0.690	<b>0.693</b>
Text	0.892	<b>0.893</b>	0.825	<b>0.830</b>	0.768	<b>0.769</b>	0.691	<b>0.696</b>
Instrument	0.872	<b>0.886</b>	0.800	<b>0.820</b>	0.747	<b>0.769</b>	0.672	<b>0.695</b>
Blood	0.877	<b>0.891</b>	0.809	<b>0.827</b>	0.760	<b>0.770</b>	0.684	<b>0.697</b>
Feces	0.898	<b>0.899</b>	0.831	<b>0.835</b>	0.753	<b>0.759</b>	0.679	<b>0.687</b>
Average	0.881	<b>0.888</b>	0.814	<b>0.823</b>	0.761	<b>0.770</b>	0.684	<b>0.696</b>

SSFormer: average improvement 0.8%

# Evaluation-RQ2

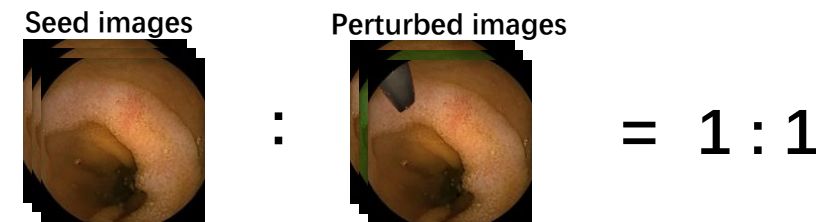
**RQ2: Can the test cases generated by our method be utilized to enhance the performance of medical image diagnosis software?**

- Segmentation
- **Classification**

# Further Training Settings

- **Dataset construction:**

- Excluded blood perturbation in the dataset construction
- 2422 images in the training set + randomly select  $\frac{1}{8}$  of the total images in each perturbation generated from the training set



- In general, we follow the original model setting in the implementation
- Minor adjustments on parameters such as the learning rate and batch size
- Performed train-validation split (0.85 : 0.15) on our customized data and used the validation performance to select checkpoints



# Results

- Classification Performance on AGDN model:

<b>Artifact</b>	<b>Initial Accuracy</b>	<b>Enhanced Accuracy</b>	<b>Difference</b>
<b>Original</b>	0.893	0.885	-0.008
<b>Blur</b>	0.702	0.808	<b>+0.106</b>
<b>Contrast</b>	0.747	0.812	<b>+0.065</b>
<b>Feces</b>	0.797	0.860	<b>+0.063</b>
<b>Instrument</b>	0.852	0.852	0.000
<b>Saturation</b>	0.685	0.770	<b>+0.085</b>
<b>Spot</b>	0.817	0.873	<b>+0.056</b>
<b>Text</b>	0.828	0.858	<b>+0.030</b>
<b>White Balance</b>	0.532	0.673	<b>+0.141</b>
<b>Average</b>	0.761	0.821	<b>+0.060</b>

**AGDN: The average improvement in accuracy score is 6%**

# Results

- Classification Performance on DSI-Net model:

<b>Artifact</b>	<b>Initial Accuracy</b>	<b>Enhanced Accuracy</b>	<b>Difference</b>
<b>Original</b>	0.940	0.947	<b>+0.007</b>
<b>Blur</b>	0.897	0.918	<b>+0.021</b>
<b>Contrast</b>	0.883	0.917	<b>+0.034</b>
<b>Feces</b>	0.907	0.937	<b>+0.030</b>
<b>Instrument</b>	0.897	0.928	<b>+0.031</b>
<b>Saturation</b>	0.755	0.835	<b>+0.080</b>
<b>Spot</b>	0.932	0.942	<b>+0.010</b>
<b>Text</b>	0.908	0.908	0.000
<b>White Balance</b>	0.728	0.848	<b>+0.120</b>
<b>Average</b>	0.872	0.909	<b>+0.037</b>

DSI-Net: The average improvement in accuracy score is 3.7%

# Answer-RQ2

Evaluate our methodology by answering the following Research Questions (RQ):

**RQ2: Can the test cases generated by our method be utilized to enhance the performance of medical image diagnosis software?**

**Answer to RQ2:** Test cases generated by MedTest can be leveraged to construct our customized training dataset and effectively improve the robustness of academic medical image diagnosis models through further training on both segmentation and classification tasks.

# Evaluation-RQ3

Evaluate our methodology by answering the following Research Questions (RQ):

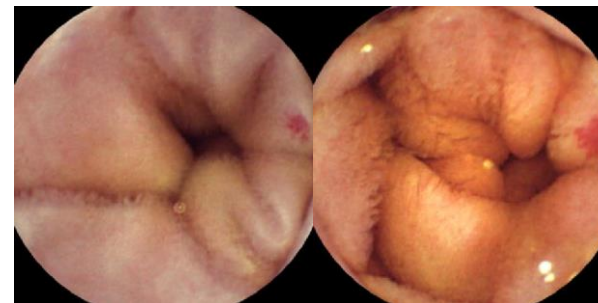
- RQ1: Is our method effective in identifying incorrect outputs produced by medical image diagnosis software and algorithms?
- RQ2: Can the test cases generated by our method be utilized to enhance the performance of medical image diagnosis software?
- **RQ3: What are the various factors that influence the performance of our method and how do they do so?**

# External Factors and Influences

- Divergence in Image Structure and Overlay
  - Hard to decide on suitable places for object-related perturbations
- Medical Landmarks Characteristics
  - Need to avoid affecting the original medical characteristics
- Lighting Conditions
  - Require targeting methods to keep the lighting condition consistent and realistic



Samples of various polyp shapes



Samples of vascular lesions (bleeding)

# Answer-RQ3

Evaluate our methodology by answering the following Research Questions (RQ):

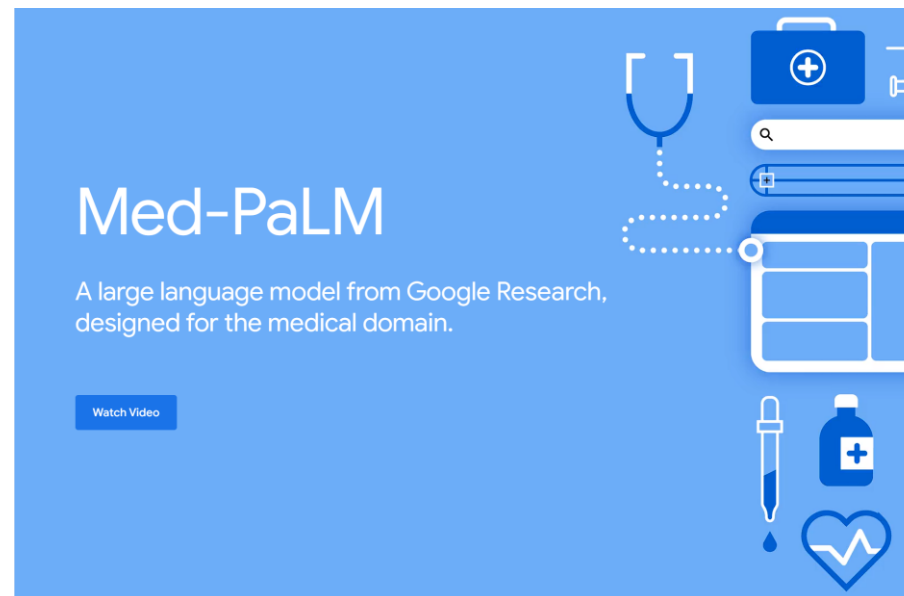
**RQ3: What are the various factors that influence the performance of our method and how do they do so?**

**Answer to RQ3:** The performance of MedTest can potentially be affected by the above - proposed factors, including image structure, medical landmark characteristics, and ambient lighting conditions. We have considered these factors in the design of MedTest and tried to mitigate the negative effect to the greatest extent in our implementation.



# Future Work

- Expand our testing objectives to other large language models and multi-modal models specially on medical diagnosis, e.g., Gemini and Med-PaLM.



# Conclusion

- Targeting the important inter-discipline of AI and medicine, we designed a comprehensive metamorphic testing paradigm, MedTest, to comprehensively evaluate models and software on medical imaging tasks.
- With our clinical-equivalent perturbations, our method was proved to effectively identify potential model errors and showed potential in assisting to improve model performance.
- Future work focuses on expanding the testing objectives of MedTest, especially in medical multi-modal models.

Thank you for listening!

# References

- [1] Zhang, Mengshi, et al. "DeepRoad: GAN-based metamorphic testing and input validation framework for autonomous driving systems." *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*. 2018.
- [2] Antol, Stanislaw, et al. "Vqa: Visual question answering." *Proceedings of the IEEE international conference on computer vision*. 2015.
- [3] Chen, Songqiang, Shuo Jin, and Xiaoyuan Xie. "Testing your question answering software via asking recursively." *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 2021.
- [4] Chen, Tsong Y., Shing C. Cheung, and Shiu Ming Yiu. "Metamorphic testing: a new approach for generating next test cases." *arXiv preprint arXiv:2002.12543* (2020).
- [5] Bohr, Adam, and Kaveh Memarzadeh. "The rise of artificial intelligence in healthcare applications." *Artificial Intelligence in healthcare*. Academic Press, 2020. 25-60.
- [6] Tomar, Nikhil Kumar, et al. "TGANet: Text-guided attention for improved polyp segmentation." *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Cham: Springer Nature Switzerland, 2022.
- [7] Pei, Kexin, et al. "Deepxplore: Automated whitebox testing of deep learning systems." *proceedings of the 26th Symposium on Operating Systems Principles*. 2017.
- [8] Vázquez, David, et al. "A benchmark for endoluminal scene segmentation of colonoscopy images." *Journal of healthcare engineering* 2017 (2017).
- [9] Ali, Sharib, et al. "An objective comparison of detection and segmentation algorithms for artefacts in clinical endoscopy." *Scientific reports* 10.1 (2020): 2748.
- [10] Acosta, Julián N., et al. "Multimodal biomedical AI." *Nature Medicine* 28.9 (2022): 1773-1784.
- [11] Xing, Xiaohan, Yixuan Yuan, and Max Q-H. Meng. "Zoom in lesions for better diagnosis: Attention guided deformation network for wce image classification." *IEEE Transactions on Medical Imaging* 39.12 (2020): 4047-4059.
- [12] Zhu, Meilu, Zhen Chen, and Yixuan Yuan. "DSI-Net: Deep synergistic interaction network for joint classification and segmentation with endoscope images." *IEEE Transactions on Medical Imaging* 40.12 (2021): 3315-3325.
- [13] Achiam, Josh, et al. "Gpt-4 technical report." *arXiv preprint arXiv:2303.08774* (2023).
- [14] Team, Gemini, et al. "Gemini: a family of highly capable multimodal models." *arXiv preprint arXiv:2312.11805* (2023).