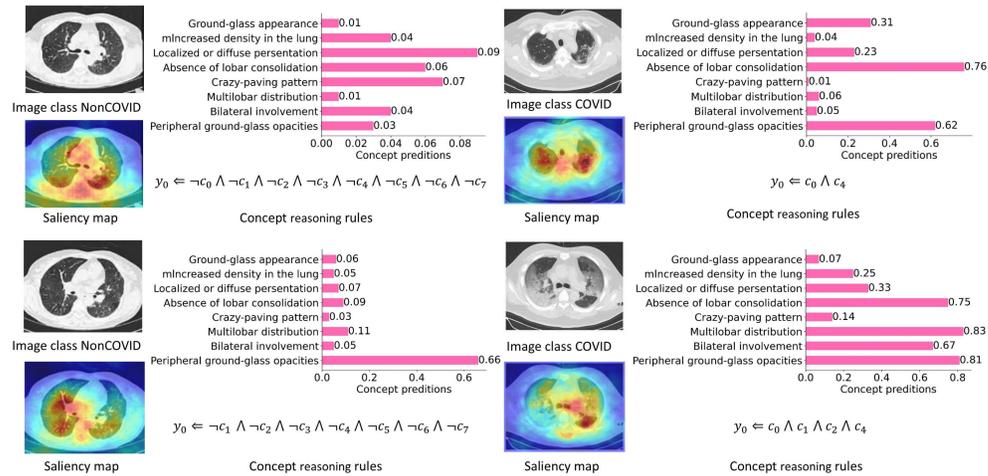


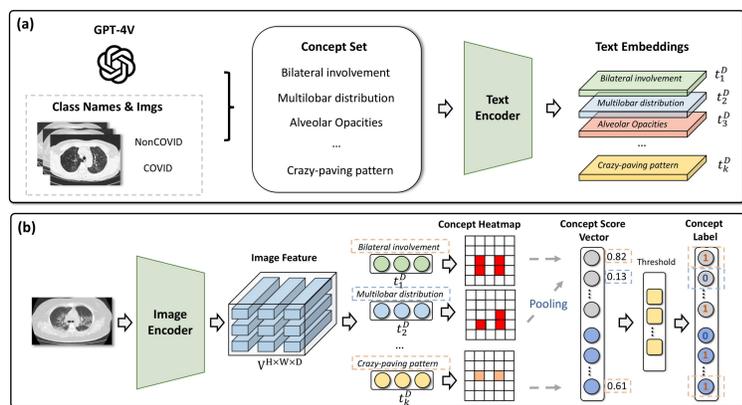
## ABSTRACT

The lack of interpretability in the field of medical image analysis has significant ethical and legal implications. Existing interpretable methods in this domain encounter several challenges, including dependency on specific models, difficulties in understanding and visualization, as well as issues related to efficiency. To address these limitations, we propose a novel framework called **Med-MICN (Medical Multi-dimensional Interpretable Concept Network)**. Med-MICN provides interpretability alignment for various angles, including neural symbolic reasoning, concept semantics, and saliency maps, which are superior to current interpretable methods. Its advantages include high prediction accuracy, interpretability across multiple dimensions, and automation through an end-to-end concept labeling process that reduces the need for extensive human training effort when working with new datasets. To demonstrate the effectiveness and interpretability of Med-MICN, we apply it to four benchmark datasets and compare it with baselines. The results clearly demonstrate the superior performance and interpretability of our Med-MICN.



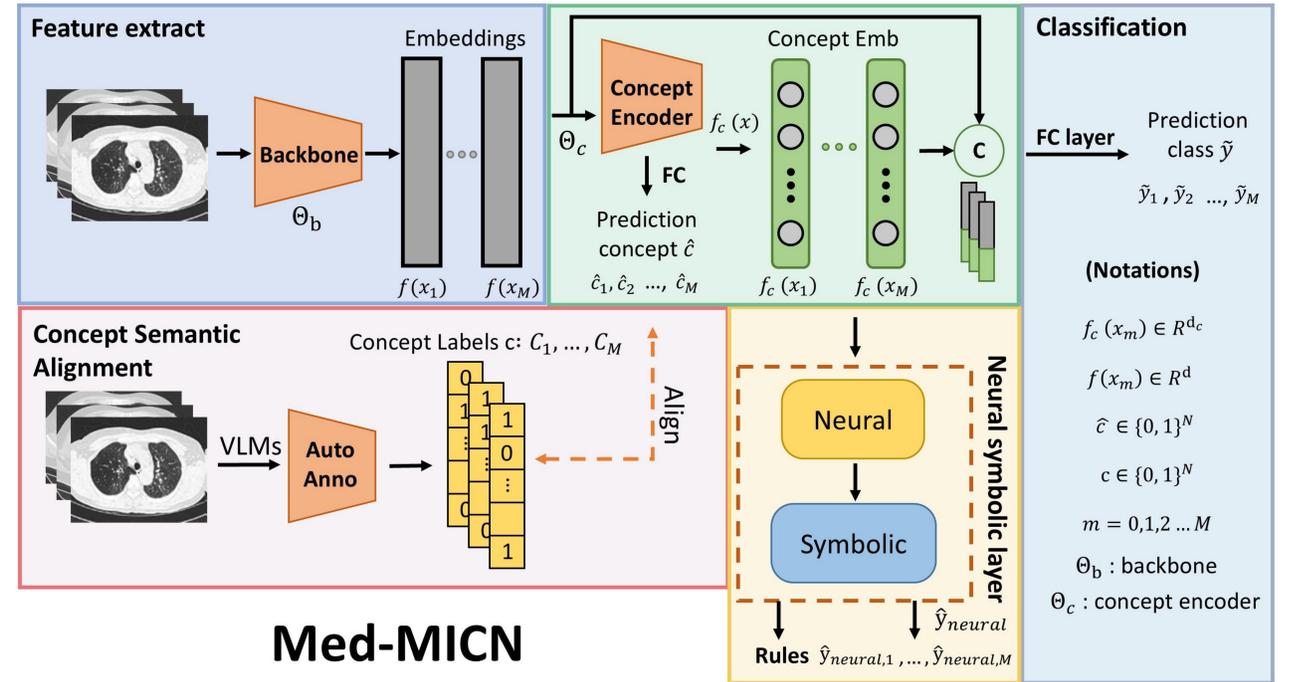
**Figure 1:** Med-MICN demonstrates multidimensional interpretability, encompassing concept score prediction, concept reasoning rules, and saliency maps, achieving alignment within the interpretative framework. The 'Peripheral ground-glass opacities' is  $c_0$ , and along the y-axis, it sequentially becomes  $c_1, \dots, c_7$ .

## ANNOTATION MODULE FRAMEWORK



**Figure 2:** (a) module, output rich dimensional interpretable conceptual information for the specified disease through the multimodal model and convert the conceptual information into text vectors through the text embedding module; (b) module, access the image to the image embedder to get the image features, and then match with the conceptual textual information to get the relevant attention region. Then, we get the influence score of the relevant region information through pooling, and finally send it to the filter to sieve out the concept information with weak relevance to get the disease concept of image information.

## PIPELINE



**Figure 3:** Overview of the Med-MICN Framework. The Med-MICN framework consists of four primary modules: (1) **Feature Extraction Module:** In the initial step, image features are extracted using a backbone network to obtain pixel-level features. (2) **Concept Embedding Module:** The extracted features are fed into the concept embedding module. This module outputs concept embeddings while passing through a category classification linkage layer to obtain predicted category information. (3) **Concept Semantic Alignment:** Concurrently, a Vision-Language Model (VLM) is used to annotate the image features, generating concept category annotations aligned with the predicted categories. (4) **Neural Symbolic Layer:** After obtaining the concept embeddings, they are input into the Neural Symbolic layer to derive conceptual rules. Finally, the concept embeddings obtained from module (2) are concatenated with the original image embeddings and fed into the final category prediction layer to produce the ultimate prediction results.

## STABILITY EVALUATION

Method	Backbone	Acc.(%)	Precision(%)	Recall(%)	F1(%)	AUC <sub>c</sub> (%)	Interpretability
Baseline	ResNet50	81.36	82.28	81.44	81.67	81.85	x
	VGG19	79.60	81.82	78.93	79.88	80.26	x
	DenseNet169	85.59	85.60	85.60	85.59	85.60	x
	SSD-COVID	81.76	81.82	78.26	80.00	88.21	x
	Label Free CBM	69.49	68.62	69.82	69.21	64.84	✓
Ours	ResNet50	84.75	84.77	84.88	84.75	84.77	✓
	VGG19	83.05	86.74	82.93	84.37	84.26	✓
	DenseNet169	86.44	87.27	86.41	87.15	87.92	✓
	DCR	76.52	71.79	63.88	65.32	63.88	✓
	DCR	76.52	71.79	63.88	65.32	63.88	✓

## ABLATION STUDY

Dataset	Ablation Setting		Metrics					Interpretability
	$\mathcal{L}_c$	$\mathcal{L}_{neural}$	ACC.(%)	Precision(%)	Recall(%)	F1(%)	AUC <sub>c</sub> (%)	
COVID-CT	✓	✓	82.20	82.92	82.21	82.55	82.64	✓
	✓	✓	83.05	83.62	83.16	83.01	83.16	✓
	✓	✓	81.36	82.11	81.38	81.70	81.81	✓
	✓	✓	84.75	84.77	84.88	84.75	84.77	✓
DDI	✓	✓	78.03	74.97	66.88	69.24	67.41	✓
	✓	✓	79.55	75.36	71.47	72.73	71.20	✓
Chest X-Ray	✓	✓	78.79	76.38	66.29	68.69	67.64	✓
	✓	✓	81.82	76.56	76.17	76.33	76.12	✓
	✓	✓	68.59	69.63	61.11	61.02	62.05	✓
	✓	✓	72.28	77.63	64.15	63.72	64.15	✓
Fitzpatrick17k	✓	✓	70.03	73.83	61.84	61.25	62.39	✓
	✓	✓	78.37	80.38	73.12	74.42	73.12	✓
	✓	✓	78.33	79.50	78.32	78.91	79.06	✓
	✓	✓	79.80	80.60	79.81	80.20	80.31	✓

**Table 1:** Experimental results from ablation studies on each loss function demonstrate that each loss function is indispensable for both accuracy and interpretability.

## CONTACT

Corresponding to [lijie.hu@kaust.edu.sa](mailto:lijie.hu@kaust.edu.sa)