Explainable image classification based on word embeddings



Results

Ing. Julija Jeršova, Prof. Ing. Tomáš Kliegr, Ph.D.

Faculty of Informatics and Statistics, Prague University of Economics and Business

Motivation

Modern AI image classifiers are powerful but often work as "black boxes," leaving users uncertain about why a prediction was made. Existing explainability methods, such as heatmaps, provide only low-level visual saliency and often lack semantic clarity. Our goal was to design an approach that makes AI reasoning fully transparent and humanunderstandable.

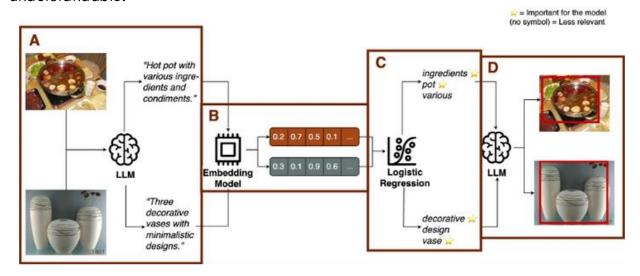


Figure 1- Workflow of the proposed method for simultaneous image classification and explanation using LLMs. In Phase A, training images are described using LLM-generated textual captions. In Phase B, the captions are represented as embeddings. In Phase C, important concepts are identified using logistic regression. In Phase D, bounding boxes are drawn using a multimodal LLM.

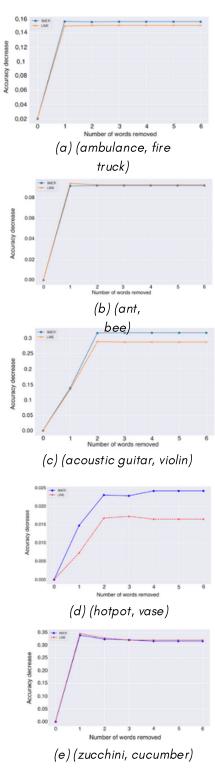
Methods

We developed a novel framework that explains image classification through semantic concepts instead of raw pixels:

- A multimodal Large Language Model (LLM) generates natural language descriptions of images.
- These descriptions are encoded into sentence embeddings.
- A logistic regression model is trained on averaged embeddings for classification.
- Explanations are provided via keyword-level attribution (SMER), highlighting the most influential words.
- The top-ranked terms are then localized in the image with bounding boxes, using targeted LLM prompting.



(right)



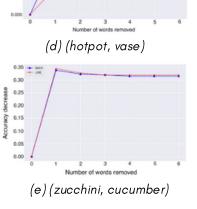


Figure 3 - Comparison of SMER and LIME via AOPC curves across five binary classification tasks.

• Evaluated on 14,351 ImageNet images (10 object classes).

- Achieved 91-98% accuracy, compared to 54-76% for a VGG16 feature-based baseline.
- AOPC metric (Area Over the Perturbation Curve) measures how much performance drops when important words are removed. SMER consistently outperformed LIME (e.g., 1.57 vs. 1.43 on violin vs. quitar).
- User study with 270 participants: Bounding-box explanations were rated significantly more interpretable than Grad-CAM heatmaps (p < 0.001).

Evaluation Aspect	Bounding Box	Heatmap	Significance (p-value)
Interpretability rating (Wilcoxon test)	Very easy: 89	Very easy: 35	
	Easy: 112	Easy: 91	
	Comprehensible: 66	Comprehensible: 116	Yes $(p < 0.001)$
	Difficult: 3	Difficult: 27	
	Very difficult: 0	Very difficult: 1	
Interpretability (t-test)	Mean = 4.06	Mean = 3.48	Yes $(p < 0.001)$
Interpretability by image class (Wilcoxon test)	Ambulance: 4.06	Ambulance: 3.36	Yes $(p < 0.001)$
	Firetruck: 3.62	Firetruck: 3.49	No
	Ant: 2.84	Ant: 3.48	Yes $(p < 0.001)$
	Bee: 2.76	Bee: 3.68	Yes $(p < 0.001)$
	Acoustic guitar: 3.33	Acoustic guitar: 2.97	Yes $(p < 0.001)$
	Violin: 2.96	Violin: 3.37	Yes $(p < 0.001)$
	Hotpot: 3.10	Hotpot: 2.70	Yes $(p < 0.001)$
	Vase: 3.23	Vase: 3.02	No
	Cucumber: 2.76	Cucumber: 2.71	No
	Zucchini: 2.73	Zucchini: 2.53	No
Open feedback sentiment (Chi-square test)	Positive: 174	Positive: 42	
	Neutral: 9	Neutral: 5	Yes $(p < 0.001)$
	Negative: 54	Negative: 55	
Open feedback themes	Clarity, precision, ease	Visibility, overlap, color confusion	Yes (qualitative

Table 1- Summary of evaluation results from the user study (N=270) comparing bounding boxes and heatmaps across interpretability and performance aspects. The p-value in the last column corresponds to the test type listed in the first column; better results for each pair are highlighted in bold.

Impact

Our findings demonstrate that language-aligned explanations enhance trust, clarity, and usability of Al predictions. The method bridges the gap between machine reasoning and human understanding, with potential applications in:

- Healthcare diagnostics (clearer Al-supported decisions for doctors)
- Education (teaching Al reasoning in an intuitive way)
- Everyday Al tools (improving transparency for non-experts)