



Nicola Dall'Asen<sup>1,2</sup>

Yiming Wang<sup>3</sup>

Enrico Fini<sup>1\*</sup>  
\* Currently at Apple

Elisa Ricci<sup>1,3</sup>

<sup>1</sup> University of Trento

<sup>2</sup> University of Pisa

<sup>3</sup> Fondazione Bruno Kessler



## CoRE uses retrieved captions to enrich image and classes embeddings in VLMs to improve classification performance in low-resource domains

### Motivation

- **Low-resource domains** are challenging for **language** and **visual** tasks
- **Scarce data** and **annotations** to train on
- **VLMs** are good in **zero-shot** tasks but **fall short** in scarce domains
- **Synthetic data** do not represent the real data in this setting

### Intuition

- Pre-trained models **under-represent low-resource** domains
- Web-crawled databases contain **noisy** or **incorrect** content
- **Specific category** appears **sparingly**, **broader category** occurs **frequently**
- **Enriching** the prompt with the **broader concept** and **noise** significantly **boost** the zero-shot **performance**

**Input**

"A photo of a melanoma"

**Retrieved captions**

"**Electronics** that control **LED** patterns"

"The basic **circuit diagram** of **LED**"

"This clearly Image **Circuit** or **Diagram**"

"**Skin Cancer** The Dangers of **Melanoma**"

"Two types of **skin melanoma** on the neck"

"**Mole** on the human skin"

"The soft gradient of **Japanese** Painted **Fern**"

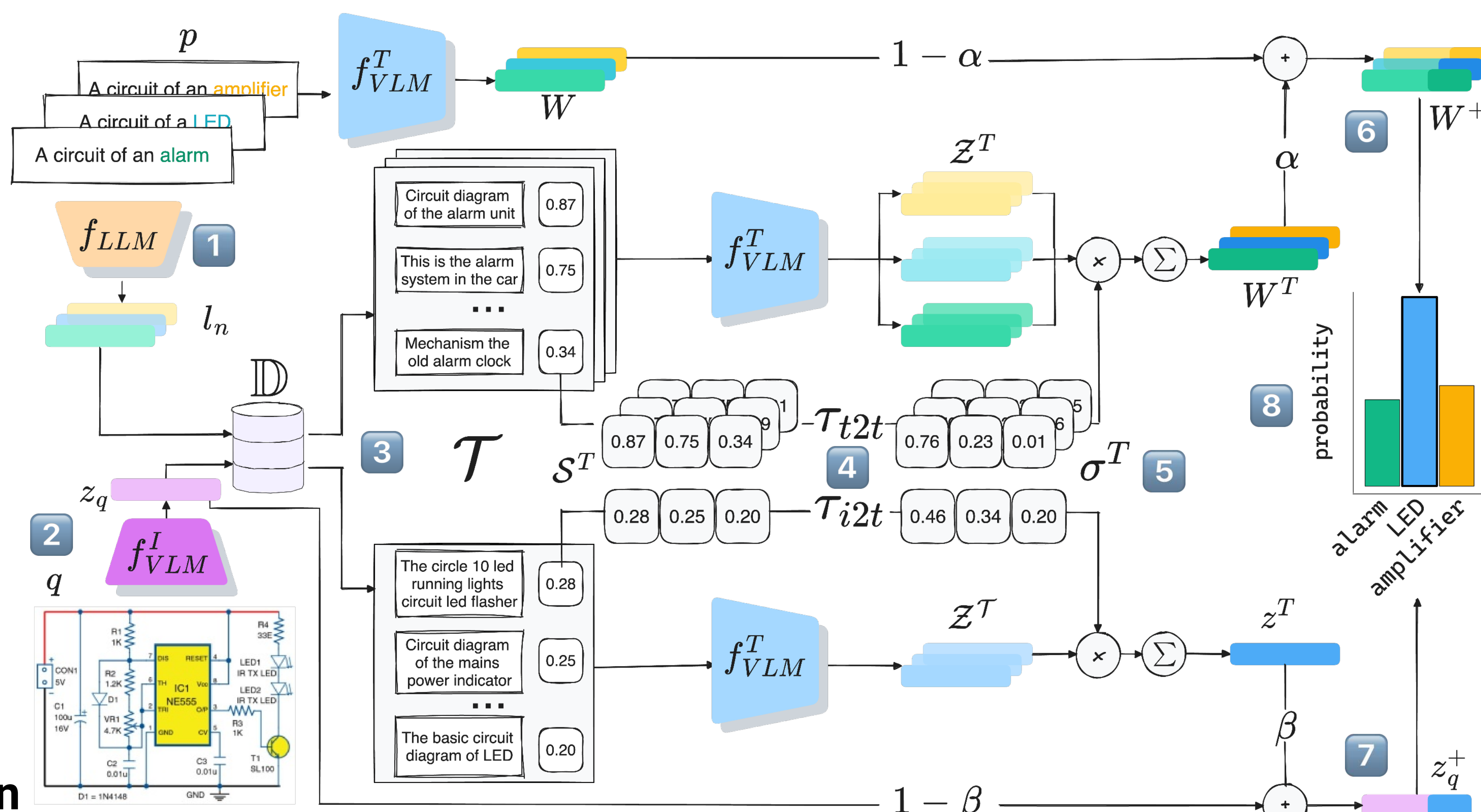
"This is a **plant** imported from **Japan**."

"Image of a Tree **Fern** and **Japanese maple**"

### Method

#### First training-free retrieval-based method for low-resource image classification

- 1 Embed **classes text** using an LLM encoder
- 2 Embed **image** using a VLM encoder
- 3 **Retrieve** most similar captions from large database
- 4 Use  $\tau$  to **skew** score **distribution**



- 5 **Weight** retrieved **embeddings** with skewed distribution
- 6 **Weight** original and **retrieved** text embeddings with  $\alpha$
- 7 **Weight** original and **retrieved** image embedding with  $\beta$
- 8 **Classify** the image using the **enriched** image and textual **representations**

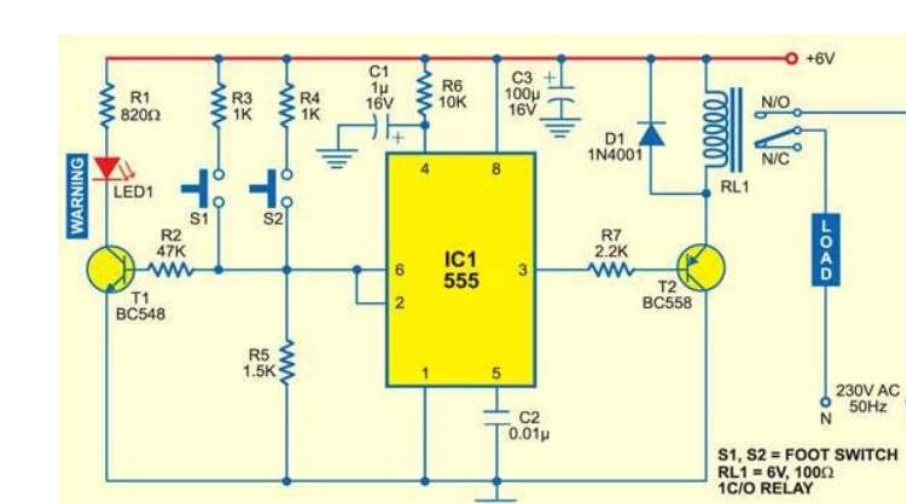
### Results

#### Benchmark of datasets and VLM baselines

Method	Circuits		iNaturalist2021 (LT100)		HAM10000	
	Acc@1	Acc@5	Acc@1	Acc@5	Acc@1	Acc@5
ImageBind (Zhang et al., 2024)	24.10	49.30	31.60 <sup>†</sup>	60.50 <sup>†</sup>	54.60 <sup>†</sup>	96.56 <sup>†</sup>
SigLIP@384px (Zhai et al., 2023)	19.53 <sup>†</sup>	30.61 <sup>†</sup>	34.50 <sup>†</sup>	63.50 <sup>†</sup>	54.60 <sup>†</sup>	95.90 <sup>†</sup>
CLIP ViT-L (Radford et al., 2021)	7.98	29.13	8.00	22.60	45.27	90.80
CLIP ViT-L@336px (Radford et al., 2021)	9.09	30.33	7.60	22.70	40.97	90.27
BLIP2-EVA (Li et al., 2023)	17.63	N/A	1.40	N/A	2.91	N/A
LlaVA 1.6 34B (Liu et al., 2023)	29.59	N/A	0.60	N/A	10.59	N/A
ImageBind (Girdhar et al., 2023)	22.36	51.02	6.70	23.90	14.43	84.25
SigLIP@384px (Zhai et al., 2023)	35.81	58.63	19.10	<b>45.70</b>	57.64	<b>96.16</b>
CoRE (Ours — CC12M)	<b>42.94</b> <sup>7.13</sup>	<b>67.71</b> <sup>9.08</sup>	<b>21.40</b> <sup>2.30</sup>	<b>42.59</b> <sup>3.11</sup>	<b>61.54</b> <sup>3.90</sup>	<b>95.70</b> <sup>0.46</sup>
CoRE (Ours — COYO-700M)	<b>43.88</b> <sup>8.07</sup>	<b>71.99</b> <sup>13.36</sup>	<b>22.10</b> <sup>3.00</sup>	<b>44.10</b> <sup>1.60</sup>	<b>62.21</b> <sup>4.57</sup>	<b>94.51</b> <sup>1.65</sup>

Our CoRE can **outperform training-based solutions**

### Qualitatives



IC 555 dry run protection  
Wireless Remote Control Switch  
230V AC Mains Over Voltage Protection Circuit



A seed fern frond is prepared for analysis.  
A tiny plant on a tree fern's trunk  
Asplenium polypodon (West Maui)

A skin lesion of melanoma.

This picture shows a melanoma lesion with varying colors.  
A mole that turned out to be melanoma skin cancer