

电子科技大学

University of Electronic Science and Technology
of China

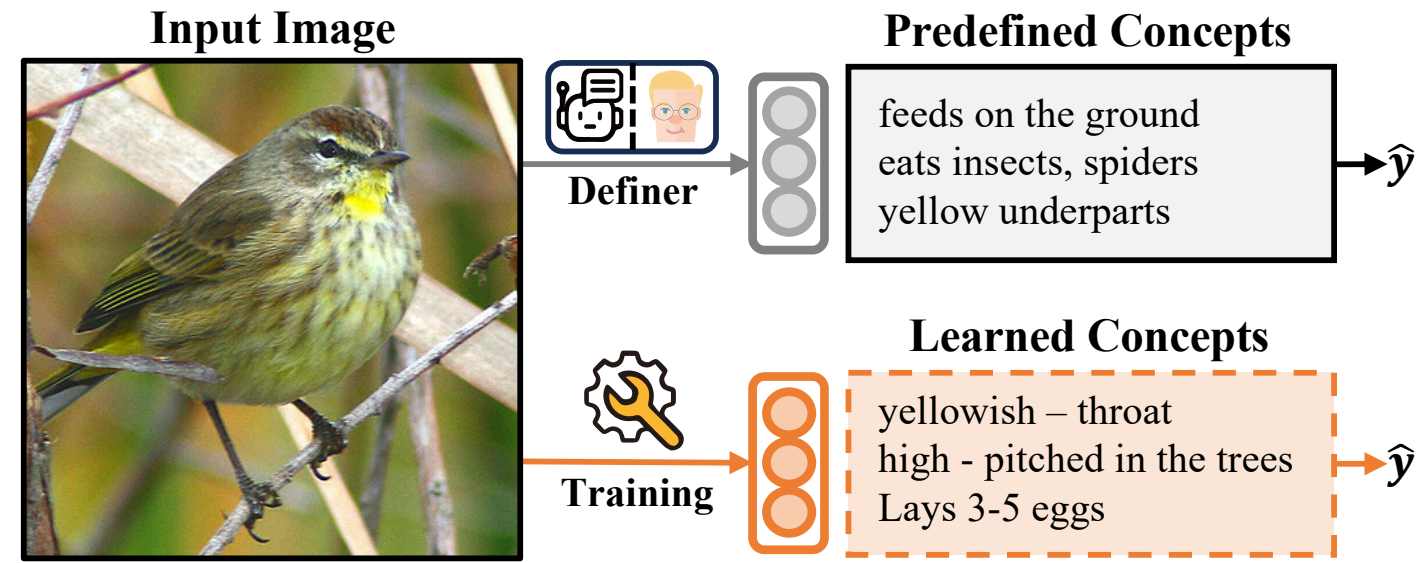
Hybrid Concept Bottleneck Models

Yang Liu¹, Tianwei Zhang¹, Shi Gu¹

¹University of Electronic Science and Technology of China

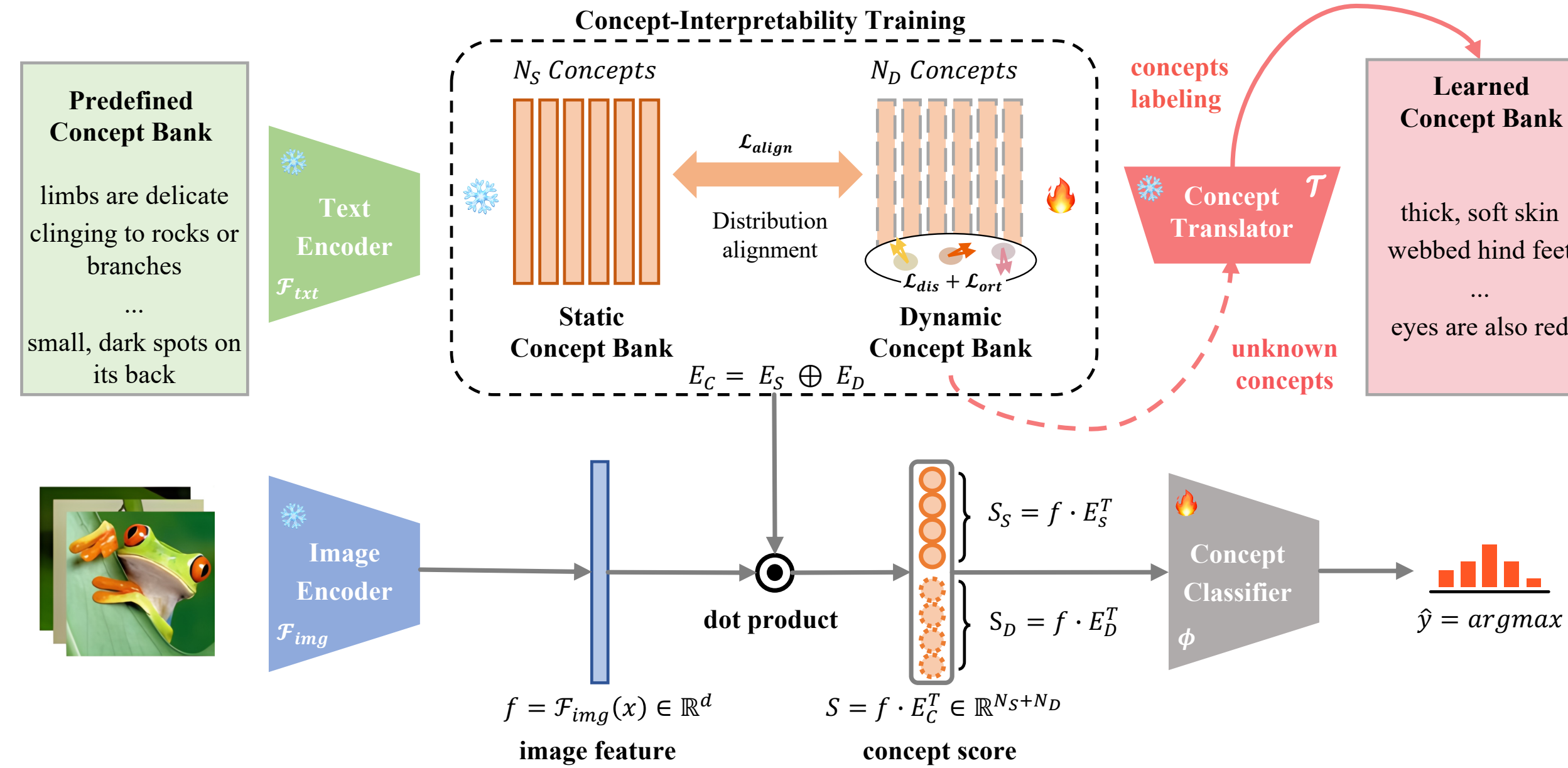


Motivation



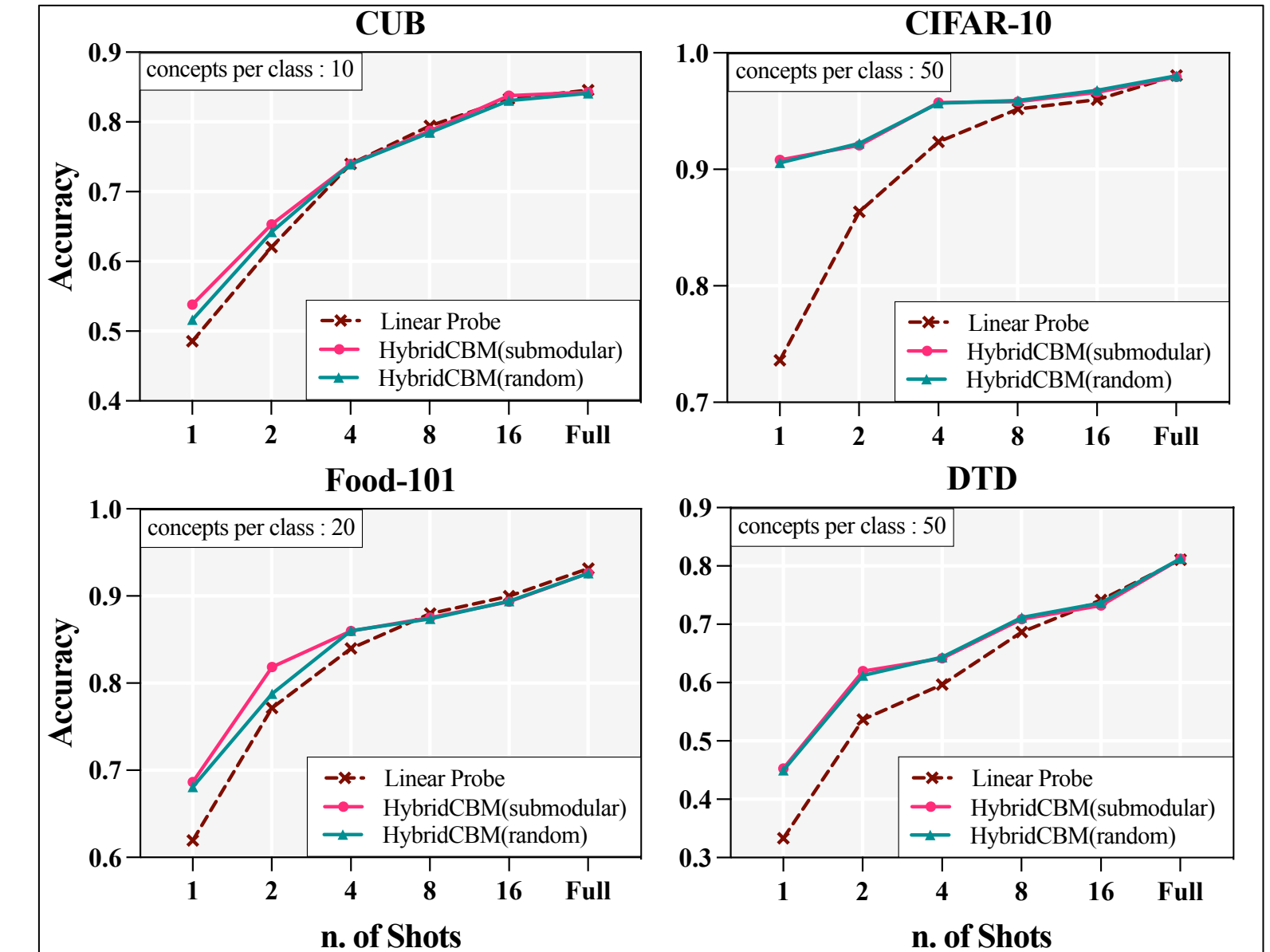
- The quality and completeness of the predefined concepts significantly affect the performance of CBMs.
- Innovative Hybrid Concept Bank:** introduces a hybrid concept bank that combines static and dynamic concepts, allowing it to adapt and refine its interpretative capabilities dynamically.
- Integration of LLM Technologies:**
 - Using LLM for defining static concepts, translating new concepts
 - Using VLM for evaluating the interpretability of concepts.

Methodology



Quantitative Analysis

Comparison of test accuracy between HybridCBM with submodular and random selection methods, and Linear Probe across 4 datasets.



Qualitative Analysis

	Class Name	Sample	Case interpretability	
			static concepts	dynamic concepts
CUB	Laysan Albatross		1.lays 4-6 white eggs 2.black and white striped body	1.black and white albatross 2.largest albatross species
Food-101	Baklava		1.flaky, phyllo dough texture 2.filled with nuts and sweetened with syrup	1.light brown center 2.a food is arranged on a large platter
Flower	Clematis		1.popular choice for making garlands and wreaths 2.borne on a climbing vine	1.large , a seed plant 2.violet

The top-2 static and dynamic concepts for randomly selected classes across three datasets are presented, focusing on case interpretability.

Concept Retrieval

You are an expert binary concept classifier, capable of determining whether a given concept has any form of relationship with the provided image. If it has any relationship with the image, respond with "yes". Otherwise, respond with "no".

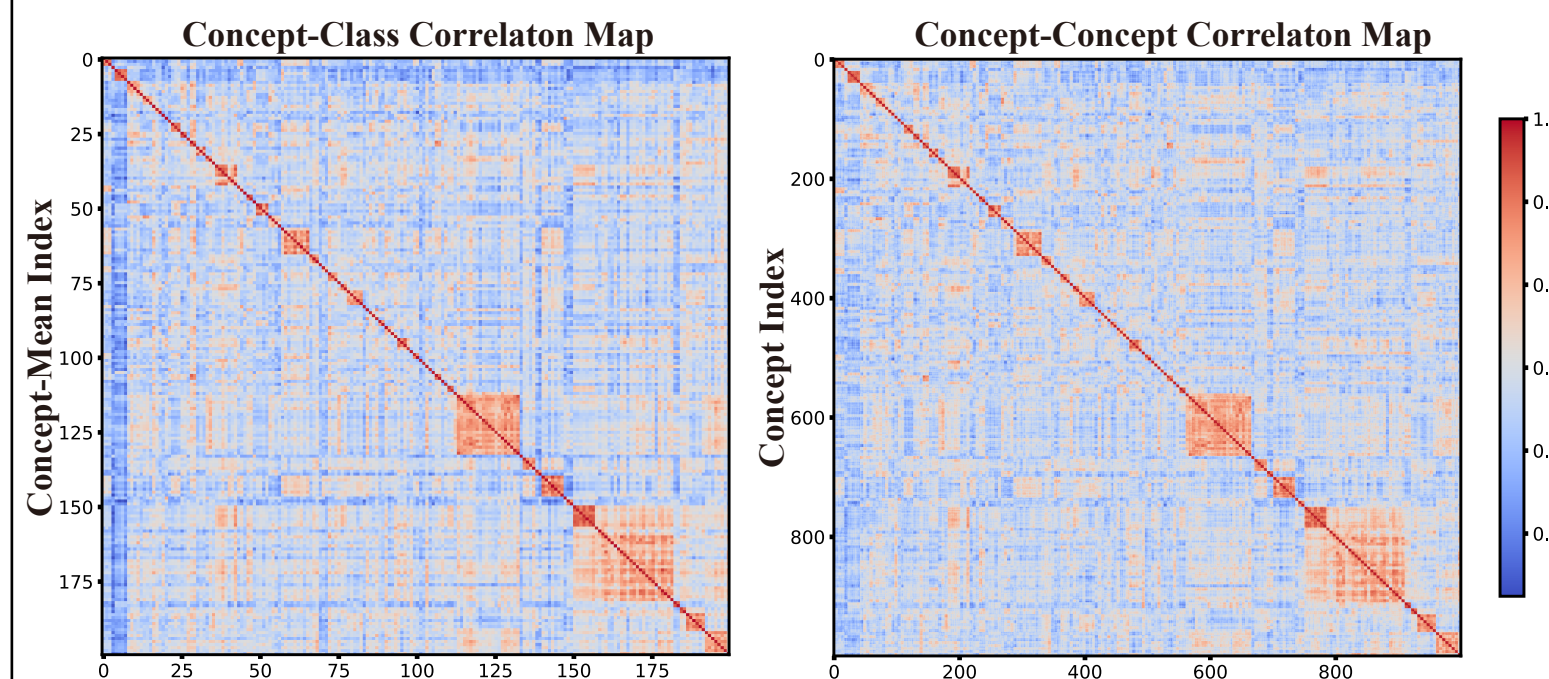
- sometimes called the "black albatross"
- known as the black-browed albatross
- most abundant albatross species



(no, yes, yes)

The concept retrieval process with vision-language model GPT-4o is illustrated on the right.

Concept Correlation Structures



Comparison of dynamic concept correlation structures on the CUB dataset.

- Concept-Class concepts:** our dynamic concepts align closely with their corresponding classes while effectively capturing diversity across classes.
- Concept-Concept concepts:** highlights the diversity among concept with different classes.

Evaluate the interpretability of HybridCBM's static and dynamic concepts through various metrics on two levels: Feature & Concept

Loss	Concept Bank	Feature Level		Concept Level	
		Purity (%)	Separation (%)	Semantics (%)	Precision@t (%)
\mathcal{L}_{cls}	Static (submodular)	19.3	86.2	38.3	48.5
	Static (random)	19.2	82.4	30.9	47.3
	Dynamic	0.1	100.0	0.01	0.05
$\mathcal{L}_{cls+dis}$	Dynamic	71.9	76.9	30.0	34.8
$\mathcal{L}_{cls+dis+ort}$	Dynamic	72.0	77.1	28.9	34.1
$\mathcal{L}_{cls+dis+ort+align}$	Dynamic	39.8	80.0	32.8	46.2