

Paper

DCBM: Data-efficient Concept Bottleneck Models Katharina Prasse* 1

Patrick Knab* 2

Sascha Marton²

Christian Bartelt²

Margret Keuper 1,3

Code

*Equal contribution

¹Data and Web Science Group, University of Mannheim ²Clausthal University of Technology

³Max-Planck-Institute for Informatics, Saarland Informatics Campus

Motivation

- Trust in applications requires interpretability of neural networks.
- Concept Bottleneck Models (CBMs) learn a linear mapping from concept activations to classes that are inherently interpretable.
- CBMs main objectives:
- → Meaningful human-interpretable concepts.
- → Concepts are sufficiently specific for the given task.
- → Efficient extraction of concepts from training images/classes.

Framework: Data-efficient CBMs

- Step 1: Concept proposals are created using foundation models for segmentation / detection.
- Step 2: Concepts are generated by clustering concept proposals to remove redundancies.
- Step 3: CBM is trained to map concept activations to class labels.
- Step 4: Visual concepts are mapped to text within CLIP space.

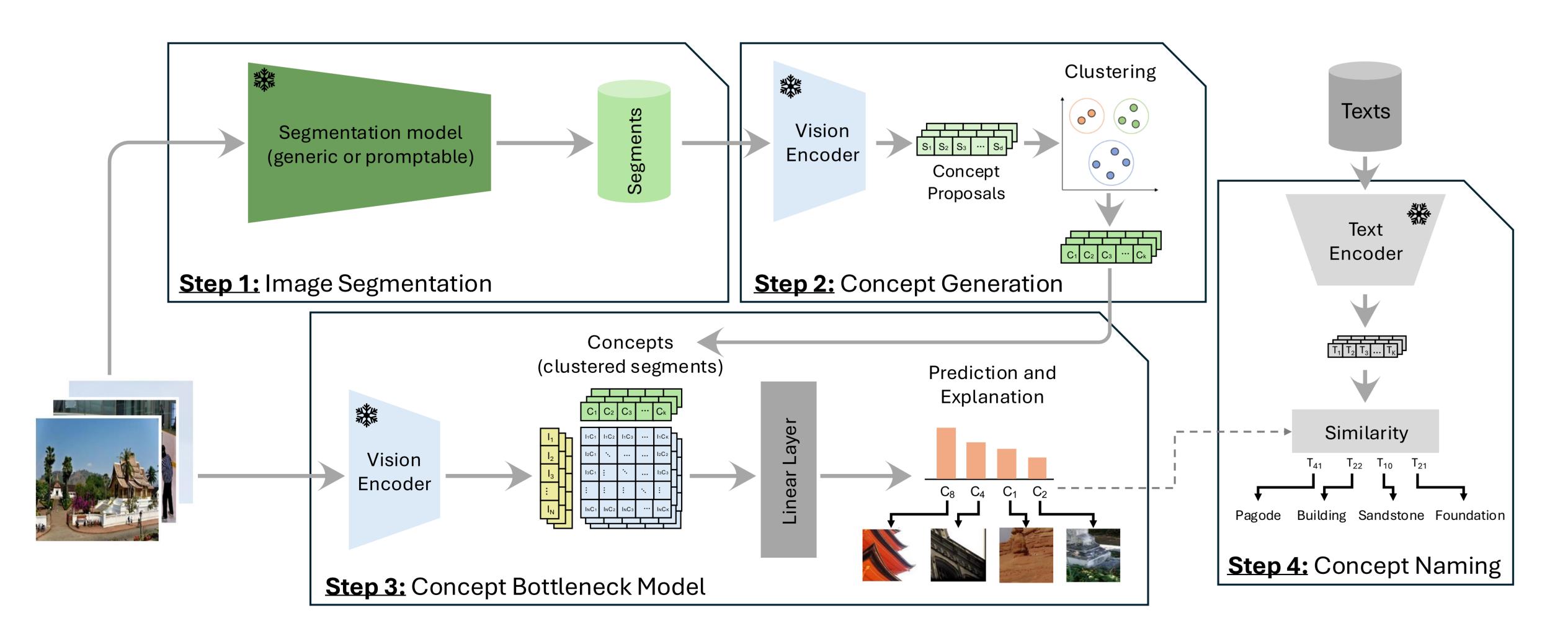


Figure 1. The **DCBM framework** generates concept proposals through foundation models (Step 1). These proposals are then clustered, each represented by its centroid (Step 2). Finally, the concepts are utilized to train a sparse CBM (Step 3). We leverage the image-text alignment to map the visual concept to the corresponding textual concept (Step 4). We can remove undesired concepts after Step 2.

No Description, No Supervision, No External Data.

Extract Concepts from YOUR Data.

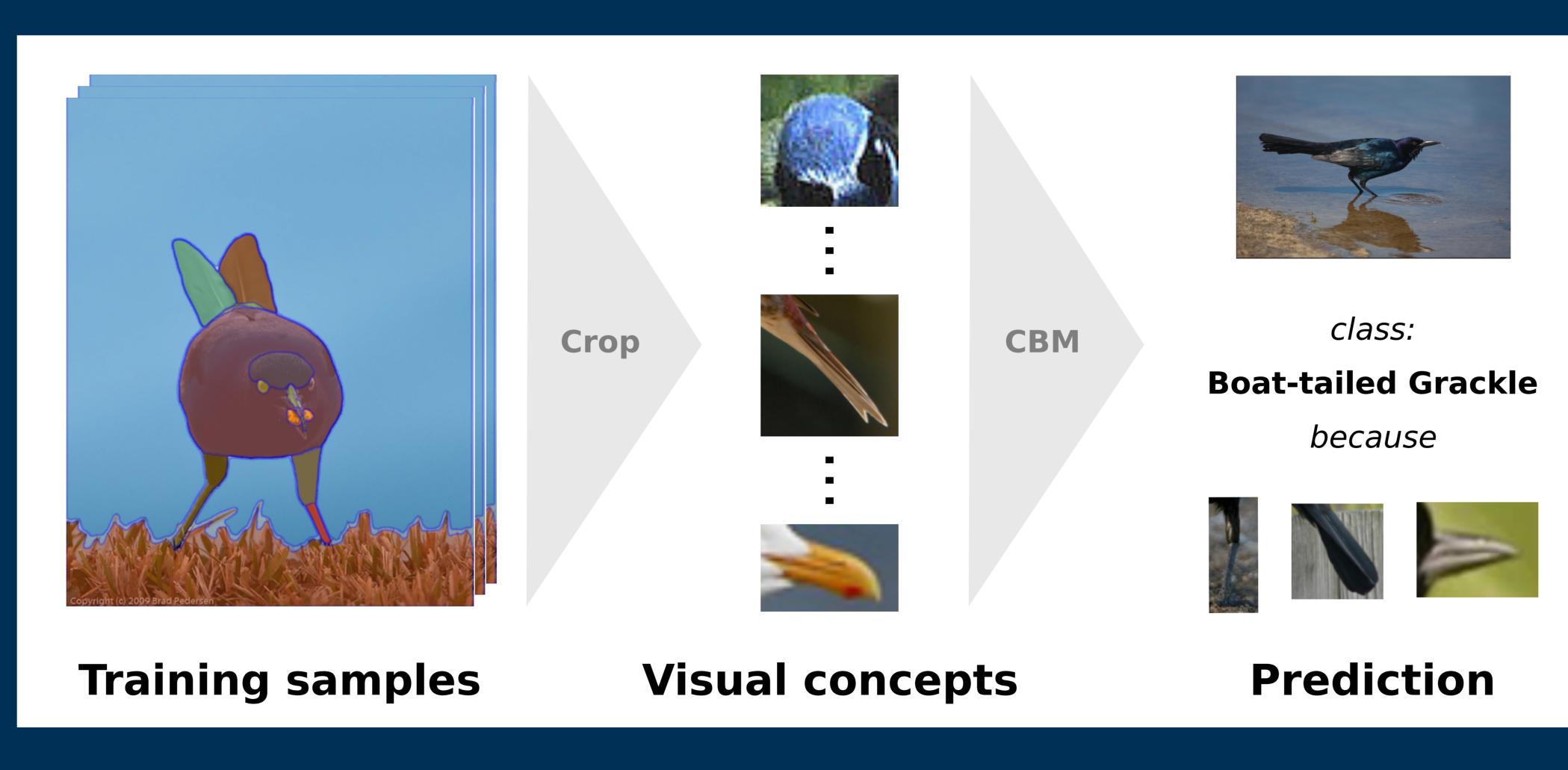


Figure 2: Using vision foundation models, we use cropped image regions as concepts for CBM training. Based on few concept samples (50 imgs/class), DCBMs offer interpretability even for fine-grained classification.

Qualitative & Quantitative Results

Key Insights:

- DCBMs perform within at most 6% of the linear probe for all datasets (9).
- Mask-RCNN concept proposals slightly outperform SAM2 and GDINO.
- DCBM excels in domain specific tasks (CUB).
- DCBM concepts can be used in OOD settings.
- DCBMs achieve competitive performance based on 50 imgs/class as concept samples.

Table 1. Top-1 accuracy comparison across CBM models.

Model	CLIP ViT L/14				
	IMN	Places	CUB	Cif10	Cif100
Linear Probe ↑ Zero-Shot ↑		55.4 40.0			
LF-CBM [3] ↑ LaBo [6] ↑ CDM [4] ↑ DCLIP [2] ↑ DN-CBM [5] ↑	83.4* 75.0*	49.4 - 55.2* 40.5* 55.6 *	- 63.5*	95.9 -	86.0* 82.2
DCBM-SAM2 (Ours) ↑ DCBM-GDINO (Ours) ↑ DCBM-MASK-RCNN (Ours) ↑		52.1 52.2 52.1	81.3	97.5	

Table 2. OOD performance. Error rate changes compared between visual CBMs (CLIP ViT-L/14) on ImageNet-R.

	IN-200	IN-R	Gap(%)
DN-CBM [5]↓	16.4	55.2	38.8
DCBM-SAM2 (Ours) ↓ DCBM-GDINO (Ours) ↓ DCBM-MASK-RCNN (Ours) ↓	21.1 22.6 22.2	47.2	24.6
	_		

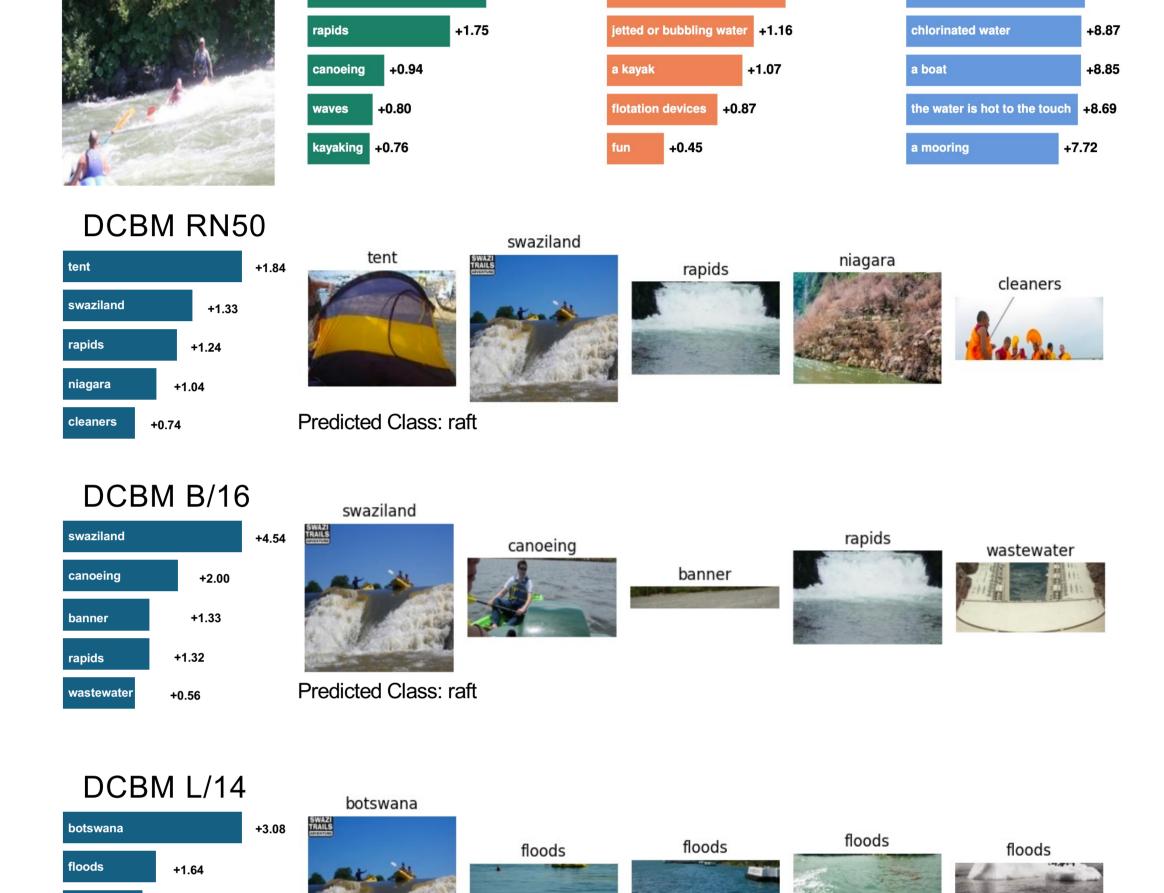


Figure 3. CBM concept explanation comparison. DCBM explanations contain no abstract concepts, e.g., fun, chlorinated water. The visual concepts in DCBM should be prioritized over the textual ones, as they originate from the image modality.

Table 3. Data-efficiency. DCBM concept proposals are generated based on 50 images/class.

	DN-CBM [5]	DCBM - ImageNet
Dataset	CC3M	50k images (50 imgs/class)
Mem	850 GB (assuming 256x256px)	6 GB
No extra data	X	\checkmark

References

- [1] M. Bohle, M. Fritz, and B. Schiele. Convolutional dynamic alignment networks for interpretable classifications. In CVPR, 2021.
- [2] S. Menon and C. Vondrick. Visual classification via description from large language models. In ICLR, 2023.
- [3] T. Oikarinen, S. Das, L. Nguyen, and L. Weng. Label-free concept bottleneck models. In ICLR, 2023.
- [4] K. P. Panousis, D. Ienco, and D. Marcos. Sparse linear concept discovery models. In ICCV, 2023.
- [5] S. Rao, S. Mahajan, M. Böhle, and B. Schiele. Discover-then-name: Task-agnostic concept bottlenecks via automated concept discovery. In ECCV, 2024. First 2 authors contribute equally.
- [6] Y. Yang, A. Panagopoulou, S. Zhou, D. Jin, C. Callison-Burch, and M. Yatskar. Language in a bottle: Language model guided concept bottlenecks for interpretable image classification. In CVPR,











