# Discover-then-Name: Task-Agnostic Concept Bottlenecks via Automated Concept Discovery

Sukrut Rao*[,1,2], Sweta Mahajan*[1,2], Moritz Böhle[1], Bernt Schiele[1]

[1]Max Planck Institute for Informatics, Saarland Informatics Campus
[2]RTG Neuroexplicit Models of Language, Vision, and Action, Saarbrücken

MAX PLANCK INSTITUTE FOR INFORMATICS

SIC Saarland Informatics Campus

ECCV EUROPEAN CONFERENCE ON COMPUTER VISION — MILANO

## Background

### Concept Bottleneck Models



Image → Feature Extractor → Features → Concept Bottleneck Layer → Concept Bottleneck (Intersection, Street, Sunlight) → Linear Classifier → Prediction and Explanation (Crosswalk)

### Typical approach: Select concept names, learn mapping

Fixed classification task

Task-specific concepts

Task-specific dataset

① Concept Name Selection
What concepts describe a crosswalk? → LLM → Intersection, Street, Stripes

Generate *task-specific* concept embeddings
CLIP Text Encoder

Generate *task-specific* image embeddings
CLIP Image Encoder

② Learn Alignment

### Consistent and interpretable concepts



pink (Index 7188), smiling (Index 1324), intersection (Index 5099)
ImageNet, CIFAR10, CIFAR100, Places365

branches ← tree → tree in field (Index 8167)
ornaments ← tree → christmas tree (Index 7446)

CLIP Image Encoder → ① Concept Discovery
CLIP Text Encoder → ② Concept Naming

### Granularity controllable by vocabulary

Semantic Consistency of Concepts
Consistency vs. Accuracy

## Ours: Discover concepts, then assign names

No fixed task

Generic dataset

Generate *task-agnostic* image embeddings
CLIP Image Encoder → ① Concept Discovery

Generic Vocabulary

Generate *task-agnostic* concept embeddings
Colourful, Green, … → CLIP Text Encoder → ② Concept Naming

- Identifies concepts used by the model
- No LLM queries needed
- Single concept bottleneck for multiple datasets

### ① Concept Discovery

Features ($\mathbf{x}$) → Disentangled Concepts ($\mathbf{c}$) → Reconstruction ($\tilde{\mathbf{x}}$)

**Use Sparse Autoencoder**

$\mathbf{W}_e$ Linear Encoder
$\mathbf{W}_d$ Linear Decoder

$$\mathcal{L}_{recon} + \lambda\mathcal{L}_{sparse}$$
$$\|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2 \quad \|\mathbf{c}\|_1$$

$$\mathbf{c} = \mathrm{ReLU}\left(\mathbf{W}_e(\mathbf{x} - \mathbf{b}_d) + b_e\right)$$

$$\tilde{\mathbf{x}} = c_1\mathbf{w}_{d1} + c_2\mathbf{w}_{d2} + c_3\mathbf{w}_{d3} + \cdots + \mathbf{b_d}$$

### ② Concept Naming

Match dictionary vector to closest name in CLIP space

Embeddings from large Vocabulary
Street $t_1$
Pink $t_2$
Intersection $t_3$
Sunlight $t_4$

→ Named Concepts
$c_1$: Intersection
$c_2$: Street
$c_3$: Sunlight

## Concept Bottleneck Model: DN-CBM

### Explanations for model decisions



Crosswalk — Top concepts (street, intersection, chinatown, sunlight)
Junkyard — Top concepts

### Class-level explanations

Crosswalk (intersection, vail, broadway, bikes, highways, intersection, ny, williamsburg, aix)
Junkyard (corrosion, citroen, bombings, crashes, gmc, abandoned, trucking, demolition, jeep)

### Similar explanations to prior work

Raft — Ours (rowing +1.75, rowing +0.94, oars +0.80, kayaking +0.76)
LF-CBM (a life jacket +1.41, jetted or bubbling water +1.16, a kayak +1.07, floatation devices +0.87)
CDM (young pr, chlorinat, a boat, the wate, a moorin)

### Competitive classification accuracy

| Model | CLIP ResNet-50 | | | | CLIP ViT-B/16 | | | |
|---|---|---|---|---|---|---|---|---|
| | Places365 | ImageNet | CIFAR10 | CIFAR100 | Places365 | ImageNet | CIFAR10 | CIFAR100 |
| Linear Probe | 53.4 | 73.3 | 88.7 | 70.3 | 55.1 | 80.2 | 96.2 | 83.1 |
| Zero Shot | 38.7 | 59.6 | 75.6 | 41.6 | 41.2 | 68.6 | 91.6 | 68.7 |
| LF-CBM | 49.0 | 67.5 | 86.4 | 65.1 | 50.6 | 75.4 | 94.6 | 77.4 |
| LaBo | - | 68.9 | 87.9 | 69.1 | - | 78.9 | 95.7 | 81.2 |
| CDM | 52.7 | 72.2 | 86.5 | 67.6 | 52.6 | 79.3 | 95.3 | 80.5 |
| DCLIP | 37.9 | 59.6 | - | - | 40.3 | 68.0 | - | - |
| DN-CBM (Ours) | 53.5 | 72.9 | 87.6 | 67.5 | 55.1 | 79.5 | 96.0 | 82.1 |

### Effective Interventions



Training Groups: Landbird on Land, Waterbird on Water
Test-only (Worst) Groups: Landbird on Water, Waterbird on Land

**Pruned Concepts**

| | |
|---|---|
| Landbird: Bird | sparrow, parrot, crow |
| Landbird: Non-bird | forest, clic |
| Waterbird: Bird | gull, ducks |
| Waterbird: Non-bird | landing, beach, |

| Model | Overall | Worst Groups | | Training Groups | |
|---|---|---|---|---|---|
| | | Landbird on Water | Waterbird on Land | Landbird on Land | Waterbird on Water |
| Before Intervention | 82.8 | 71.3 | 57.5 | 98.6 | 93.3 |
| Only Bird Concepts | 89.4 (+6.6) | 86.6 (+15.3) | 71.3 (+13.8) | 96.8 (-1.8) | 91.4 (-1.9) |
| Only Non-bird Concepts | 60.8 (-22.0) | 28.5 (-42.8) | 28.8 (-28.7) | 95.0 (-3.6) | 85.8 (-7.5) |