

# Energy-Based Concept Bottleneck Models: Unifying Prediction, Concept Intervention, and Probabilistic Interpretations

Xinyue Xu, Yi Qin, Lu Mi, Hao Wang\*, Xiaomeng Li\*

18 Apr 2024



**ICLR**



**R**

**W**

# Concept-Based Models

Input image  $x$



**DNN**

Concepts  $c$



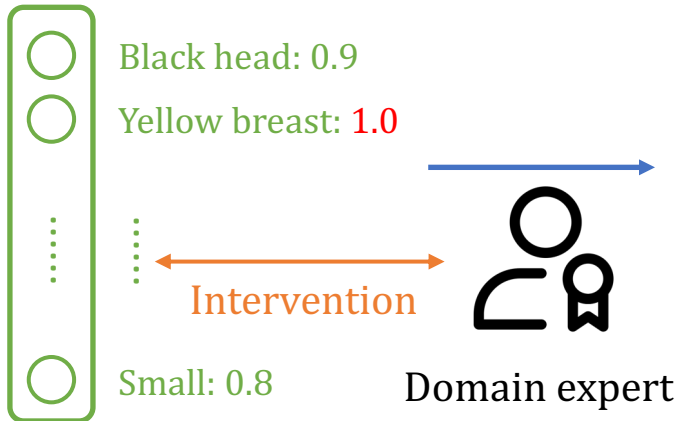
Output label  $y$

Black and White Warbler

## After Intervention



**DNN**



Kentucky Warbler

# Current Limitations

---

- **Interpretability**

Cannot effectively quantify the intricate relationships between various concepts and class labels.

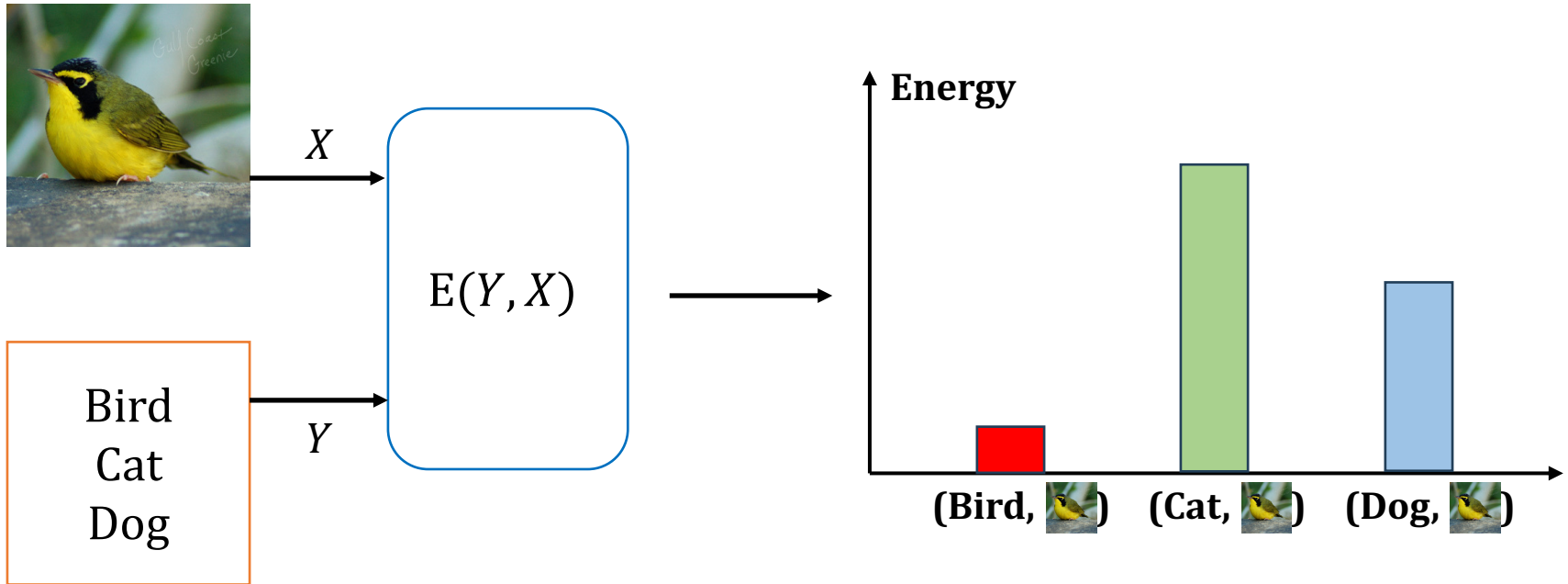
- **Intervention**

Struggle to account for the complex interactions among concepts.

- **Performance**

Suffer from a trade-off between model performance and interpretability.

# Energy-Based Models

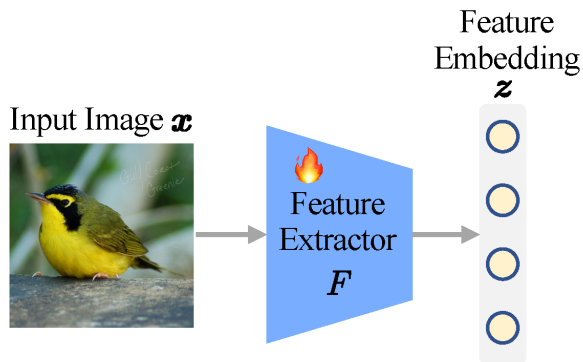


For energy networks, **lower energy  $E$**  indicates **better compatibility** (e.g.,  $E(x, y_{gt}) = 0$ ).

# Our Method (ECBM): Feature Extractor

Given the input  $x$  and a candidate label  $y$ , the feature extractor  $F$  first compute the features  $z = F(x)$ .

## Training



Trainable



Frozen

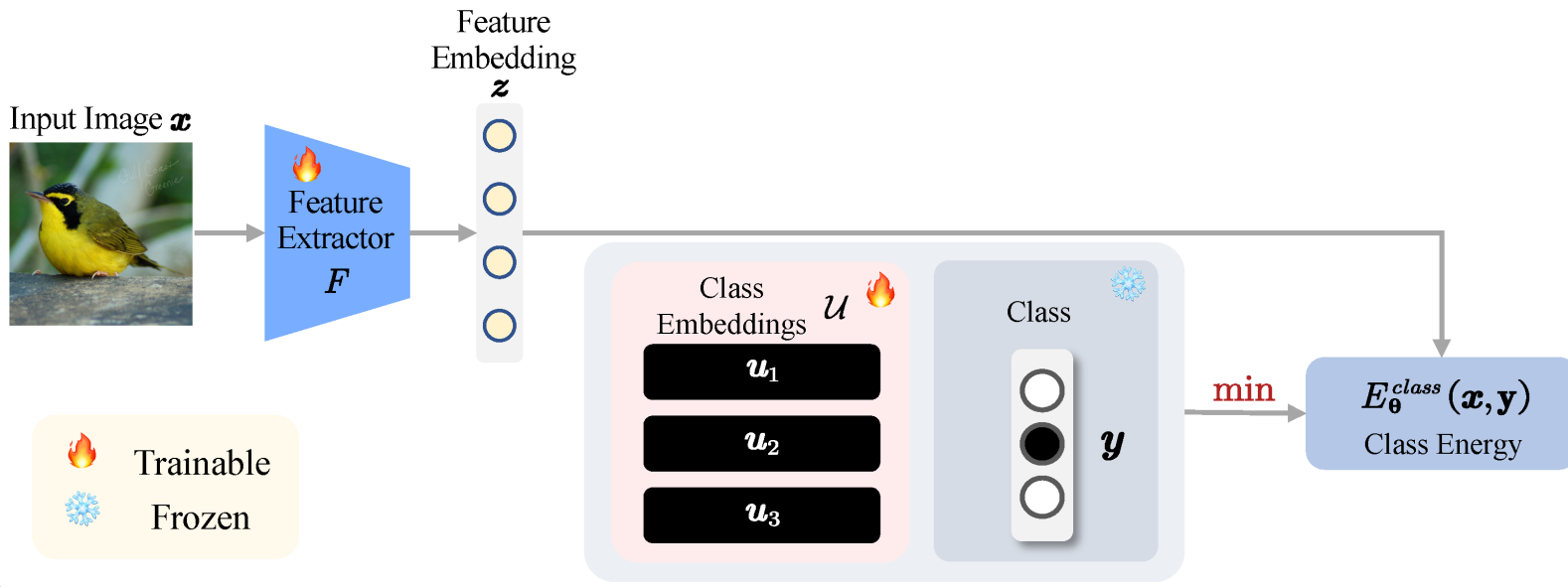
# ECBM: Class Energy Network $E_{\theta}^{class}(\mathbf{x}, \mathbf{y})$

Measure the compatibility between input  $\mathbf{x}$  and class label  $\mathbf{y}$ .

$$E_{\theta}^{class}(\mathbf{x}, \mathbf{y}) = G_{zu}(\mathbf{z}, \mathbf{u})$$

$$\mathcal{L}_{class}(\mathbf{x}, \mathbf{y}) = E_{\theta}^{class}(\mathbf{x}, \mathbf{y}) + \log \left( \sum_{m=1}^M e^{-E_{\theta}^{class}(\mathbf{x}, \mathbf{y}_m)} \right). \quad (1)$$

## Training



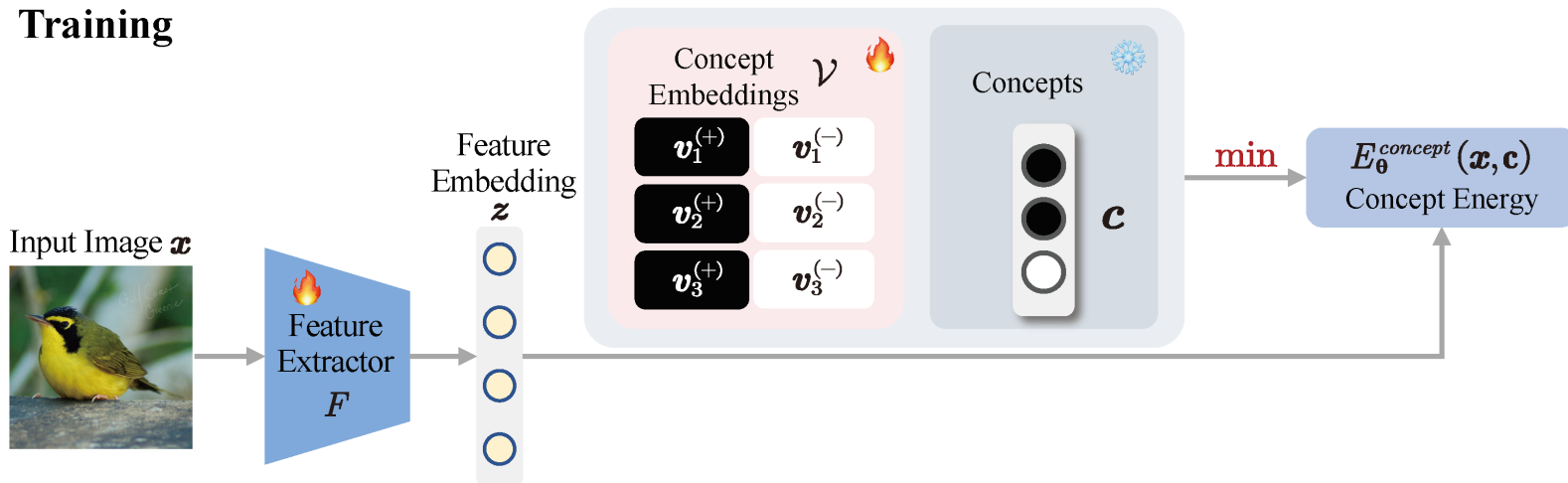
# ECBM: Concept Energy Network $E_{\theta}^{concept}(\mathbf{x}, \mathbf{c})$



Measure the compatibility between input  $\mathbf{x}$  and the  $K$  concepts  $\mathbf{c}$ .

$$E_{\theta}^{concept}(\mathbf{x}, \mathbf{c}_k) = G_{z\mathbf{v}}(\mathbf{z}, \mathbf{v}_k)$$

$$\mathcal{L}_{concept}^{(k)}(\mathbf{x}, \mathbf{c}_k) = E_{\theta}^{concept}(\mathbf{x}, \mathbf{c}_k) + \log \left( \sum_{c_k \in \{0,1\}} e^{-E_{\theta}^{concept}(\mathbf{x}, \mathbf{c}_k)} \right). \quad (2)$$

## Training



 Trainable  
 Frozen

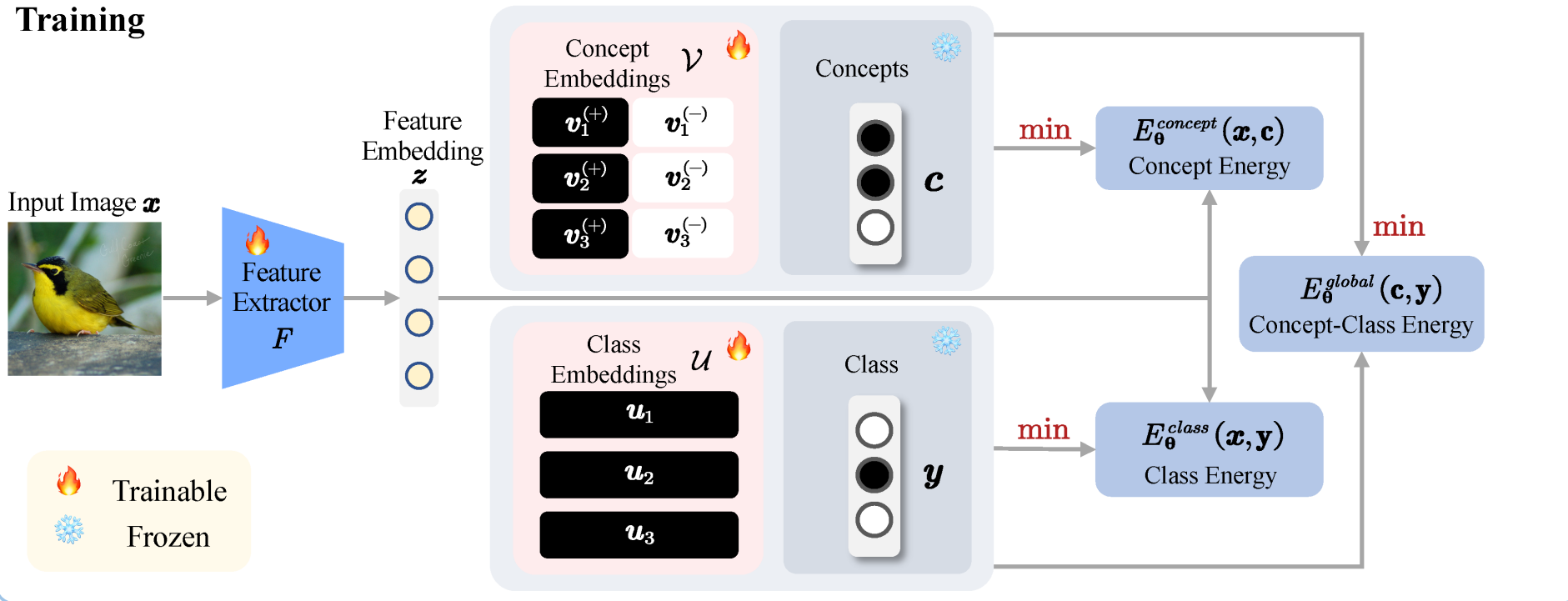
# ECBM: Global Energy Network $E_{\theta}^{global}(\mathbf{c}, \mathbf{y})$

Measure the compatibility between the  $K$  concepts  $\mathbf{c}$  and class label  $\mathbf{y}$ .

$$E_{\theta}^{global}(\mathbf{c}, \mathbf{y}) = G_{vu}([\mathbf{v}_k]_{k=1}^K, \mathbf{u})$$

$$\mathcal{L}_{global}(\mathbf{c}, \mathbf{y}) = E_{\theta}^{global}(\mathbf{c}, \mathbf{y}) + \log \left( \sum_{m=1, \mathbf{c}' \in \mathcal{C}}^M e^{-E_{\theta}^{global}(\mathbf{c}', \mathbf{y}_m)} \right). \quad (3)$$

## Training



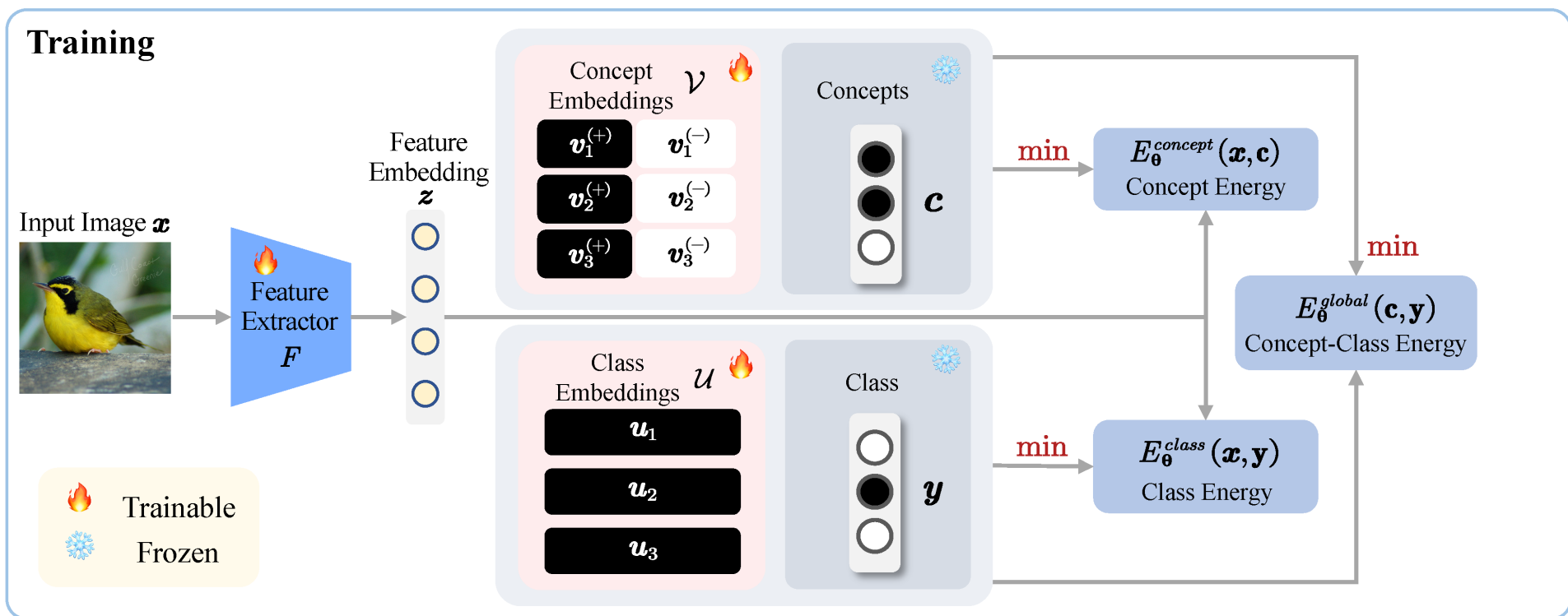


# Training Phase: Minimize Loss

ECBM is trained by **minimizing** the following total loss function:

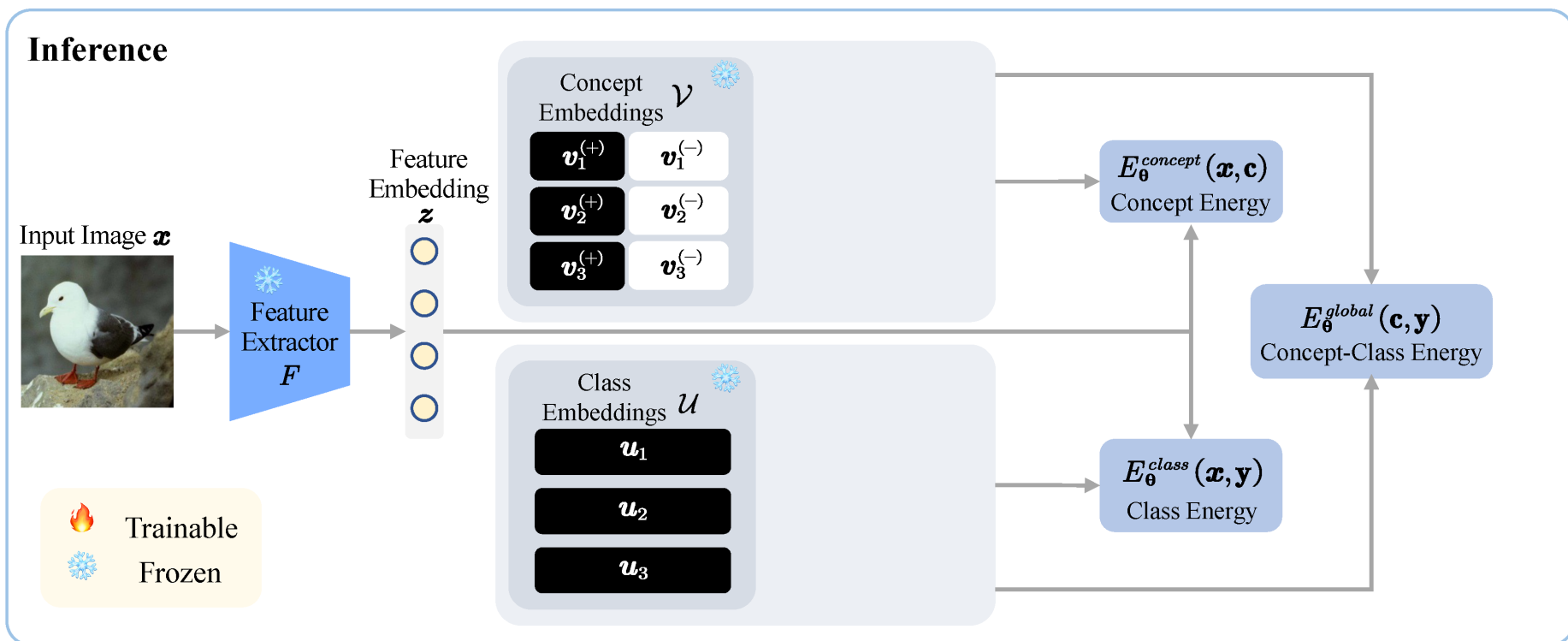
$$\mathcal{L}_{total}(\mathbf{x}, \mathbf{c}, \mathbf{y}) = \mathcal{L}_{class}(\mathbf{x}, \mathbf{y}) + \lambda_c \mathcal{L}_{concept}(\mathbf{x}, \mathbf{c}) + \lambda_g \mathcal{L}_{global}(\mathbf{c}, \mathbf{y}), \quad (4)$$

$$\mathcal{L}_{total}^{all} = \mathbb{E}_{(\mathbf{x}, \mathbf{c}, \mathbf{y}) \sim p_{\mathcal{D}}(\mathbf{x}, \mathbf{c}, \mathbf{y})} [\mathcal{L}_{total}(\mathbf{x}, \mathbf{c}, \mathbf{y})]. \quad (5)$$



# Inference Phase : Freeze Parameters

To predict  $\mathbf{c}$  and  $\mathbf{y}$  given the input  $\mathbf{x}$ , we **freeze** the feature extractor  $F$  and the energy network parameters  $\theta$ .

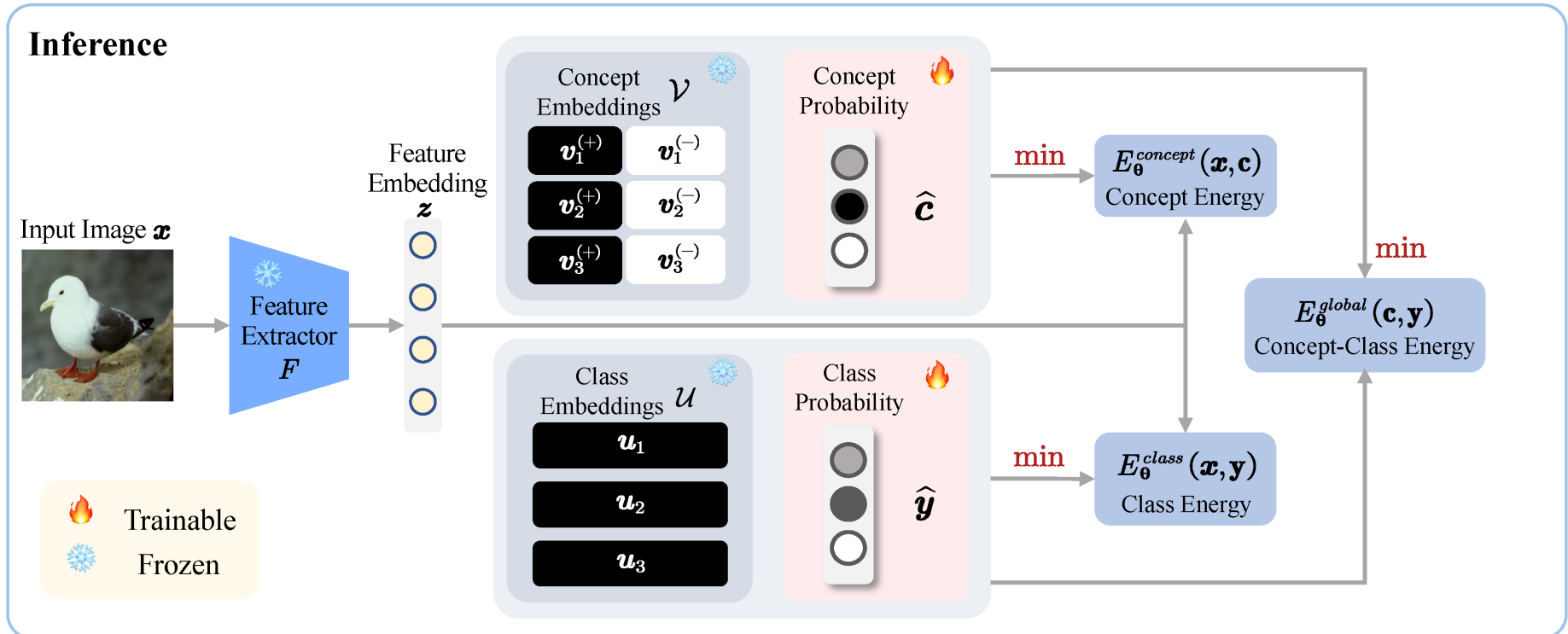


# Inference Phase: Search Optimum

**Search** for the optimal prediction of concepts  $\hat{\mathbf{c}}$  and the class label  $\hat{\mathbf{y}}$  as follows:

$$\arg \min_{\hat{\mathbf{c}}, \hat{\mathbf{y}}} \mathcal{L}_{class}(\mathbf{x}, \hat{\mathbf{y}}) + \lambda_c \mathcal{L}_{concept}(\mathbf{x}, \hat{\mathbf{c}}) + \lambda_g \mathcal{L}_{global}(\hat{\mathbf{c}}, \hat{\mathbf{y}}), \quad (6)$$

$$E_{\theta}^{joint}(\mathbf{x}, \mathbf{c}, \mathbf{y}) \triangleq E_{\theta}^{class}(\mathbf{x}, \mathbf{y}) + \lambda_c E_{\theta}^{concept}(\mathbf{x}, \mathbf{c}) + \lambda_g E_{\theta}^{global}(\mathbf{c}, \mathbf{y}). \quad (7)$$



# Experimental Results

Model \ Data	CUB			CelebA			AWA2		
	Concept	Overall Concept	Class	Concept	Overall Concept	Class	Concept	Overall Concept	Class
CBM	0.964	0.364	0.759	0.837	0.381	0.246	<b>0.979</b>	0.803	0.907
ProbCBM*	0.946	0.360	0.718	0.867	0.473	0.299	0.959	0.719	0.880
PCBM	-	-	0.635	-	-	-	-	-	-
CEM	0.965	0.396	0.796	0.867	0.457	0.330	0.978	0.796	0.908
<b>ECBM</b>	<b>0.973</b>	<b>0.713</b>	<b>0.812</b>	<b>0.876</b>	<b>0.478</b>	<b>0.343</b>	<b>0.979</b>	<b>0.854</b>	<b>0.912</b>

- **Slightly** outperform others in terms of *concept accuracy*.

# Experimental Results

Model \ Data	CUB			CelebA			AWA2		
	Concept	Overall Concept	Class	Concept	Overall Concept	Class	Concept	Overall Concept	Class
CBM	0.964	0.364	0.759	0.837	0.381	0.246	<b>0.979</b>	0.803	0.907
ProbCBM*	0.946	0.360	0.718	0.867	0.473	0.299	0.959	0.719	0.880
PCBM	-	-	0.635	-	-	-	-	-	-
CEM	0.965	0.396	0.796	0.867	0.457	0.330	0.978	0.796	0.908
<b>ECBM</b>	<b>0.973</b>	<b>0.713</b>	<b>0.812</b>	<b>0.876</b>	<b>0.478</b>	<b>0.343</b>	<b>0.979</b>	<b>0.854</b>	<b>0.912</b>

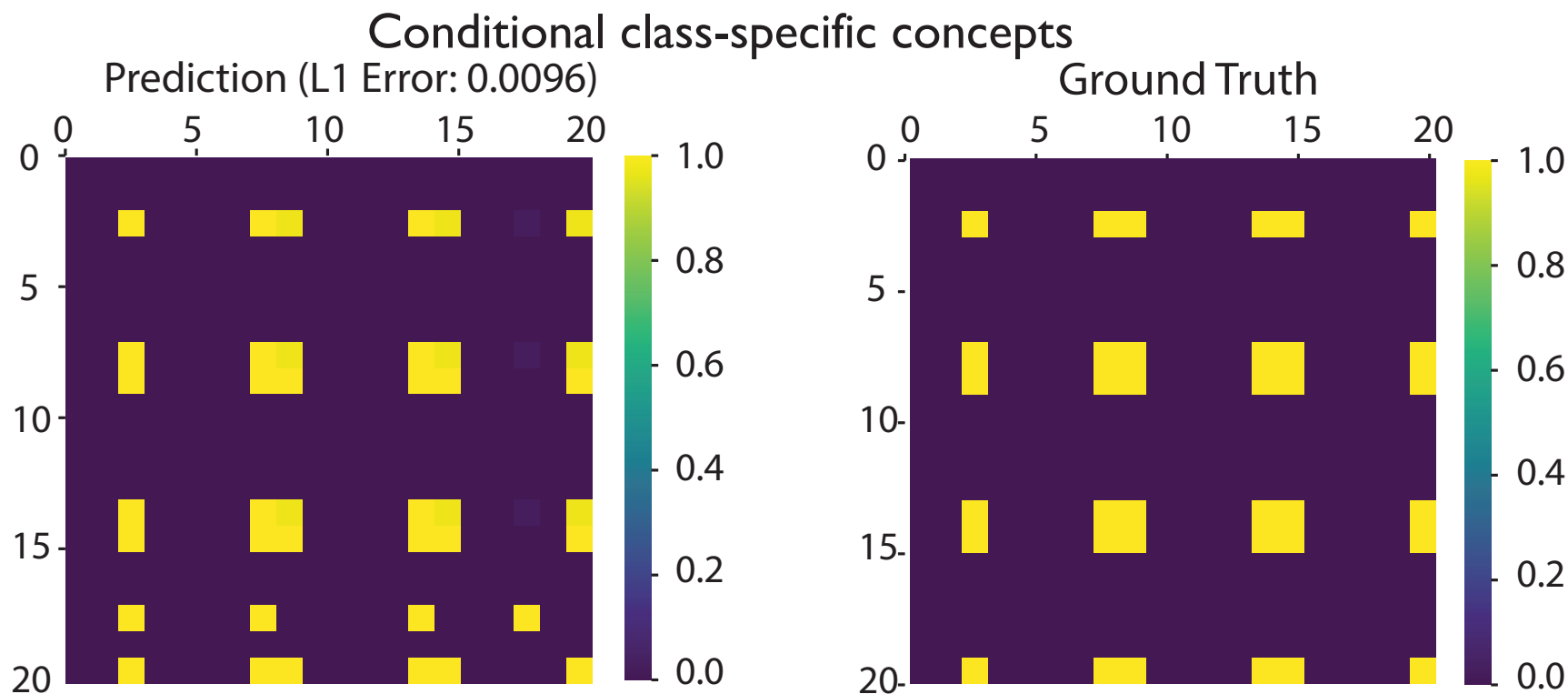
- **Slightly** outperform others in terms of *concept accuracy*.
- Successfully capture the *interaction (and correlation)* among the concepts.
  - *Significantly* outperforms other methods in terms of *overall concept accuracy*.

# Experimental Results

Model \ Data	CUB			CelebA			AWA2		
	Concept	Overall Concept	Class	Concept	Overall Concept	Class	Concept	Overall Concept	Class
CBM	0.964	0.364	0.759	0.837	0.381	0.246	<b>0.979</b>	0.803	0.907
ProbCBM*	0.946	0.360	0.718	0.867	0.473	0.299	0.959	0.719	0.880
PCBM	-	-	0.635	-	-	-	-	-	-
CEM	0.965	0.396	0.796	0.867	0.457	0.330	0.978	0.796	0.908
<b>ECBM</b>	<b>0.973</b>	<b>0.713</b>	<b>0.812</b>	<b>0.876</b>	<b>0.478</b>	<b>0.343</b>	<b>0.979</b>	<b>0.854</b>	<b>0.912</b>

- **Slightly** outperform others in terms of *concept accuracy*.
- Successfully capture the *interaction (and correlation)* among the concepts.
  - *Significantly* outperforms other methods in terms of *overall concept accuracy*.
- Outperform the state-of-the-art on class accuracy.

# Conditional Interpretation



$$p(c_k = 1 | c_{k'} = 1, \mathbf{y} = \text{Black and White Warble})$$

# Conclusion

---

- Propose the **first general method - ECBM**, to unify:
  - Concept correction
  - Conditional interpretation
  - Concept-based prediction
- **Under a unified energy formulation**, compute arbitrary conditional probabilities.
- **Significantly** outperform the state-of-the-art on **real-world** datasets.

