# FAITHFUL VISION-LANGUAGE INTERPRETATION VIA CONCEPT BOTTLENECK MODELS

SONGNING LAI*, LIJIE HU*, JUNXIAO WANG, LAURE BERTI-EQUILLE, AND DI WANG
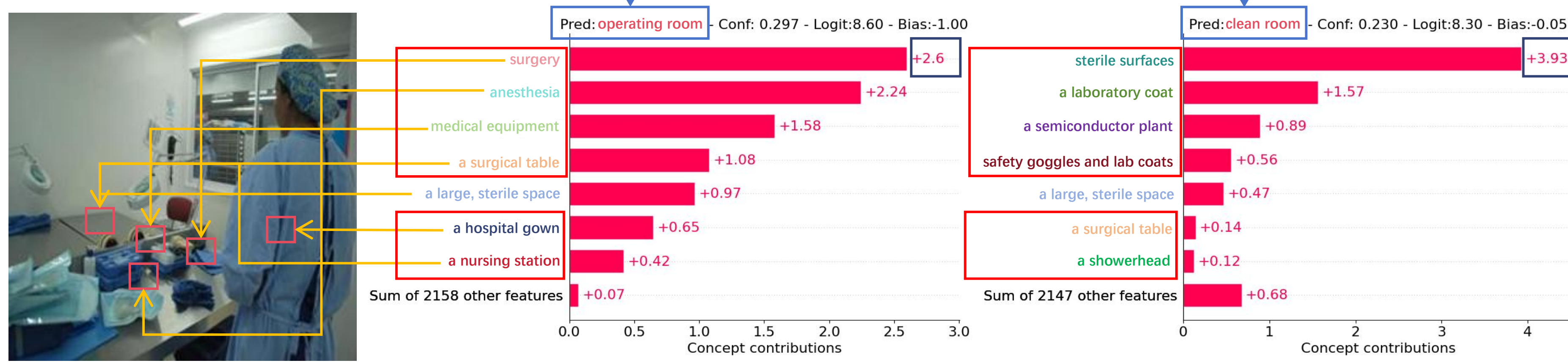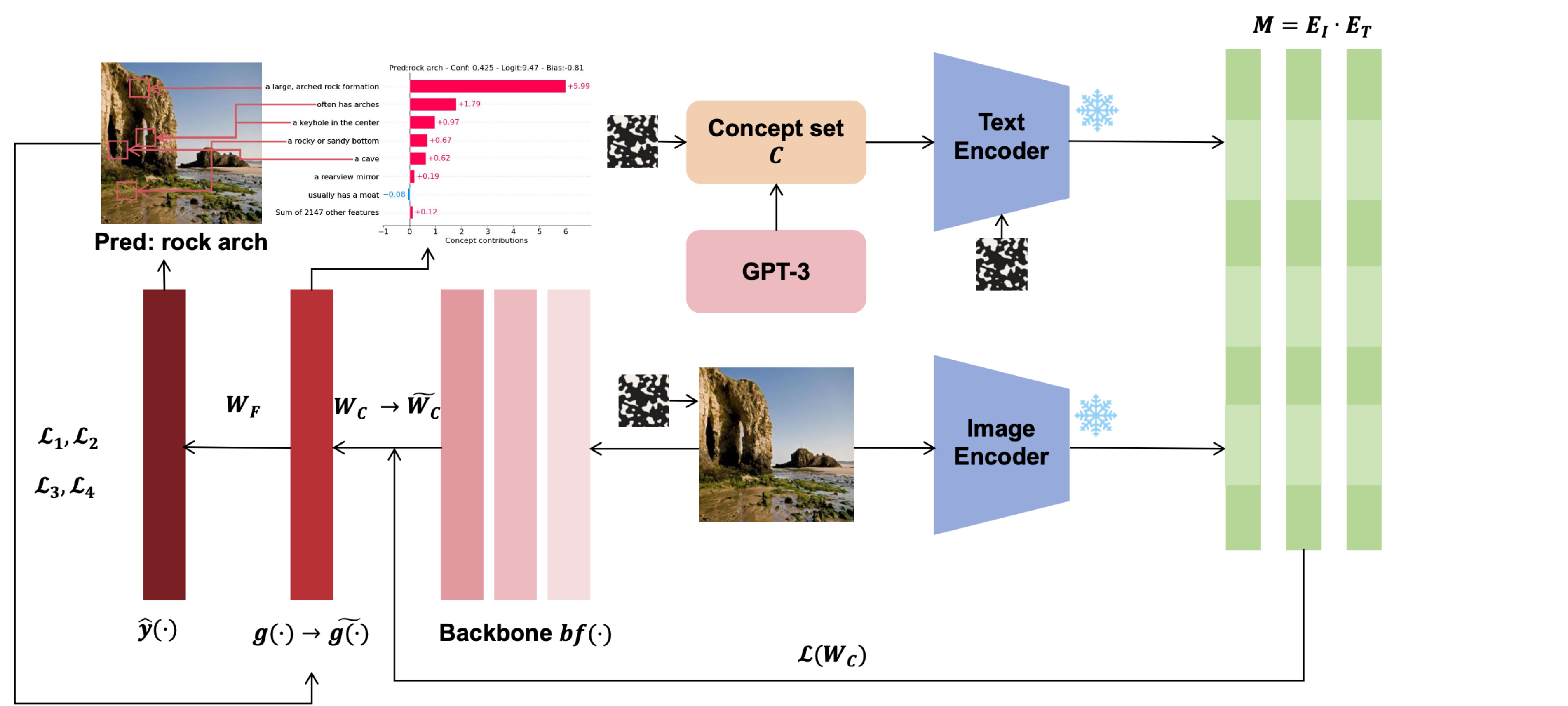
## INTRODUCTION

Traditional Concept Bottleneck Models (CBMs) require substantial manual annotation, which label-free CBM effectively addresses by leveraging factual information from pre-trained models. However, this convenience comes with inherent instability in pre-trained models. We addressed it in this paper.
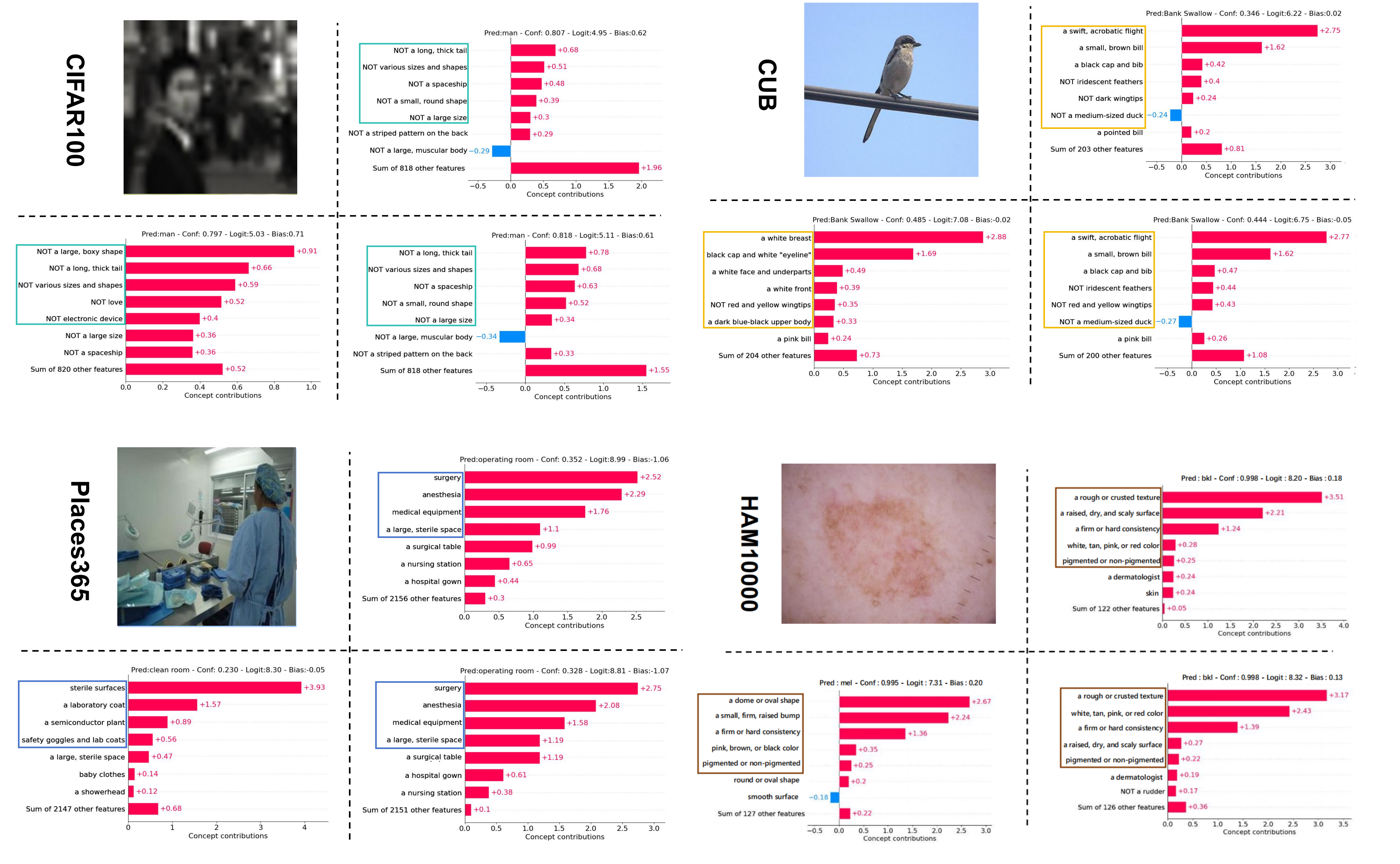


## DEFINITION AND FRAMEWORK



$$\min_{\tilde{W}_c} \mathbb{E}_x [\lambda_1 \underbrace{D(y(x,\tilde{c}), y(x,c))}_{\mathcal{L}_1} + \lambda_2 \underbrace{\mathcal{L}_{k_1}(\tilde{g}(x), g(x))}_{\mathcal{L}_2} + \lambda_3 \underbrace{\max_{||\delta|| \le R_2} D(y(x,\tilde{c}), y(x,\tilde{c}+\delta))}_{\mathcal{L}_3} + \lambda_4 \underbrace{\max_{||\rho|| \le R_1} \mathcal{L}_{k_2}(\tilde{g}(x), \tilde{g}(x)+\rho)}_{\mathcal{L}_4}].$$

**Definition 1 (Faithful Vision-Language Concept).** *Under the same concept space, i.e., under the set of concepts generated by GPT3 at one time, we call a matrix $\tilde{W}_c$ is a $(D, R, \alpha, \beta, k_1, k_2)$-Faithful Vision-Language Concept (FVLC) model for the vanilla concept if it satisfies for any input $x$:*

- *(Similarity of Explanation) $V_{k_1}(\tilde{g}(x), g(x)) \ge \beta_1$ for some $1 \ge \beta_1 \ge 0$;*
- *(Stability of Explanation) $V_{k_2}(\tilde{g}(x), \tilde{g}(x)+\rho) \ge \beta_2$ for some $1 \ge \beta_2 \ge 0$ and all $||\rho|| \le R_1$, where $||\cdot||$ is a norm and $R_1 \ge 0$;*
- *(Closeness of Prediction) $D(y(x,\tilde{c}), y(x,c)) \le \alpha_1$ for some $\alpha_1 \ge 0$, where $D$ is some probability distance or divergence;*
- *(Stability of Prediction) $D(y(x,\tilde{c}), y(x,\tilde{c}+\delta)) \le \alpha_2$ for all $||\delta|| \le R_2$, where $D$ is some probability distance or divergence, $||\cdot||$ is a norm and $R_2 \ge 0$,*

*where $\tilde{g}(x) = \tilde{W}_c bf(x)$, $y(x,c) = W_F g(x)$, and $y(x,\tilde{c}) = W_F \tilde{g}(x)$, $y(x,\tilde{c}+\delta) = W_F(\tilde{g}(x)+\delta)$. For any given $x$, $\tilde{c} = \tilde{g}(x)$ is a $(D, R, \alpha, \beta, k_1, k_2)$-FVLC. Here, $\alpha = \min\{\alpha_1, \alpha_2\}$, $\beta = \max\{\beta_1, \beta_2\}$, and $R = \min\{R_1, R_2\}$.*

## EVALUATING INTERPRETABILITY AND UTILITY



## STABILITY EVALUATION

| Method | CIFAR10 | CIFAR100 | CUB | Places365 |
|---|---|---|---|---|
| Standard (No interpretability) | 88.80% | 70.10% | 76.70% | 48.56% |
| P-CBM (CLIP) | 84.50% | 56.00% | N/A | N/A |
| Label-free CBM | 86.32% | 65.42% | 74.23% | 43.63% |
| WP1(5%) - base | 86.47% | 65.13% | 74.08% | 43.57% |
| WP1(5%) - FVLC | 86.34% | 65.43% | 73.96% | 43.67% |
| WP1(10%) - base | 86.25% | 65.09% | 73.97% | 43.67% |
| WP1(10%) - FVLC | 86.39% | 64.90% | 73.92% | 43.62% |
| WP2 - base | 86.41% | 65.16% | 73.96% | 43.54% |
| WP2 - FVLC | 86.22% | 65.34% | 74.44% | 44.55% |
| IP - base | 86.62% | 65.36% | 74.39% | 43.64% |
| IP - FVLC | 86.88% | 65.29% | 74.01% | 43.71% |
| WP1(5%)+WP2 - base | 86.49% | 65.17% | 73.90% | 43.67% |
| WP1(5%)+WP2 - FVLC | 86.43% | 65.33% | 73.92% | 43.49% |
| WP1(10%)+WP2 - base | 86.50% | 64.87% | 73.82% | 43.61% |
| WP1(10%)+WP2 - FVLC | 86.38% | 65.06% | 74.01% | 43.44% |
| WP1(10%)+WP2+IP - base | 85.96% | 64.41% | 73.74% | 43.32% |
| WP1(10%)+WP2+IP - FVLC | 86.70% | 65.14% | 74.36% | 43.46% |

| Method | CIFAR10 | | CIFAR100 | | CUB | | Places365 | |
|---|---|---|---|---|---|---|---|---|
| | TCPC | TOPC | TCPC | TOPC | TCPC | TOPC | TCPC | TOPC |
| WP1(5%) - base | 1.55E-01 | 6.32E-02 | 1.01E-01 | 7.17E-02 | 1.26E-01 | 1.85E-01 | 1.59E-01 | 6.40E-02 |
| WP1(5%) - FVLC | 1.12E-03 | 8.55E-03 | 2.81E-03 | 4.51E-03 | 1.05E-02 | 1.38E-03 | 1.30E-03 | |
| WP1(10%) - base | 1.99E-01 | 8.36E-02 | 1.94E-01 | 1.31E-01 | 2.32E-01 | 3.41E-01 | 1.14E-01 | |
| WP1(10%) - FVLC | 1.19E-03 | 7.40E-03 | 3.67E-03 | 4.55E-03 | 1.19E-02 | 1.53E-03 | 1.39E-03 | 1.25E-03 |
| WP2 - base | 1.53E-01 | 4.99E-02 | 1.36E-01 | 6.67E-02 | 1.43E-01 | 1.73E-01 | 1.40E-01 | 6.37E-02 |
| WP2 - FVLC | 1.10E-02 | 8.72E-03 | 3.35E-03 | 4.55E-03 | 1.05E-02 | 1.53E-03 | 1.29E-03 | |
| IP - base | 1.68E-01 | 6.28E-02 | 1.38E-01 | 8.81E-02 | 1.71E-01 | 2.23E-01 | 1.73E-01 | 8.09E-02 |
| IP - FVLC | 8.02E-03 | 8.29E-03 | 3.24E-03 | 4.56E-03 | 1.04E-02 | 1.53E-03 | 1.25E-03 | |
| WP1(5%)+WP2 - base | 1.85E-01 | 7.46E-02 | 1.28E-01 | 6.65E-02 | 1.44E-01 | 1.79E-01 | 1.60E-01 | 6.32E-02 |
| WP1(5%)+WP2 - FVLC | 1.20E-02 | 7.56E-03 | 3.67E-03 | 4.55E-03 | 9.81E-02 | 1.51E-03 | 1.54E-03 | 1.28E-03 |
| WP1(10%)+WP2 - base | 1.17E-01 | 8.62E-02 | 1.93E-01 | 1.32E-01 | 1.76E-01 | 3.45E-01 | 2.12E-01 | 1.17E-01 |
| WP1(10%)+WP2 - FVLC | 1.18E-02 | 9.41E-03 | 2.06E-02 | 1.44E-02 | 1.87E-02 | 3.79E-02 | 2.74E-02 | 1.18E-02 |
| WP1(10%)+WP2+IP - base | 1.36E-01 | 1.05E-01 | 2.22E-01 | 1.55E-01 | 1.95E-01 | 3.54E-01 | 2.16E-01 | 1.44E-01 |
| WP1(10%)+WP2+IP - FVLC | 1.43E-02 | 1.11E-02 | 3.29E-02 | 1.77E-02 | 2.21E-02 | 4.54E-02 | 3.35E-02 | 1.34E-02 |

## ABLATION STUDY



## CONTACT

Corresponding to lijie.hu@kaust.edu.sa