

见微知著：消除单模态虚假关联的多模态奖励模型泛化方法

The Devil Is in the Details: Tackling Unimodal Spurious Correlations for Generalizable Multimodal Reward Models

李梓超, 温学儒, 姜杰, 季雨秋, 陆垚杰, 韩先培, 张德兵, 孙乐

Forty-Second International Conference on Machine Learning (ICML 2025)

李梓超 (lizichao2022@iscas.ac.cn)

研究背景

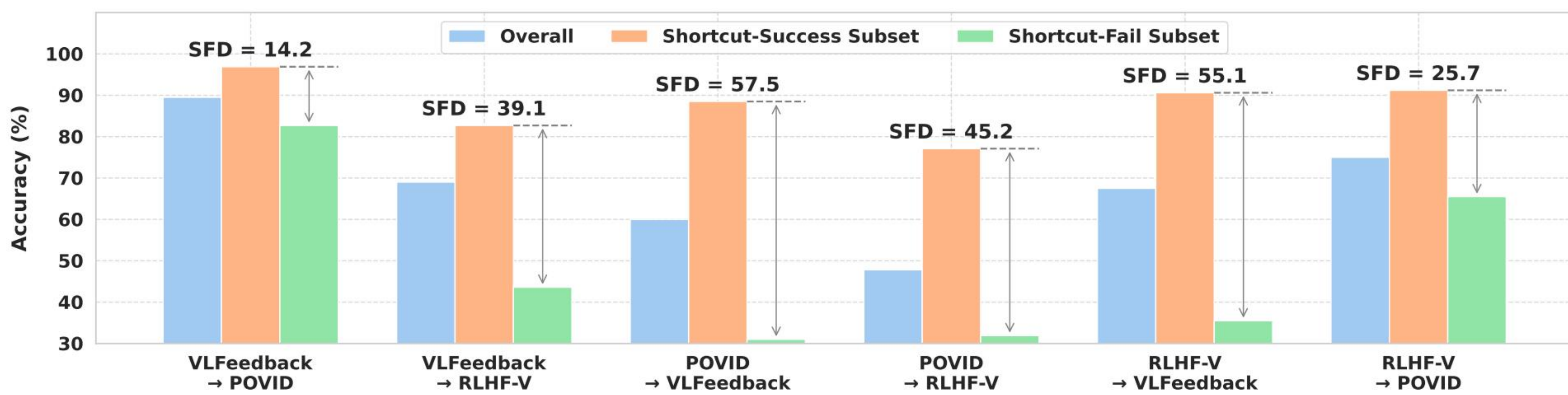
➤ 多模态奖励模型 (MM-RM) 在大语言模型偏好对齐中扮演关键角色，而泛化能力至关重要

核心贡献

➤ 识别单模态虚假关联关键挑战
➤ 提出捷径感知多模态奖励模型算法
➤ 增强模型跨分布场景下泛化能力

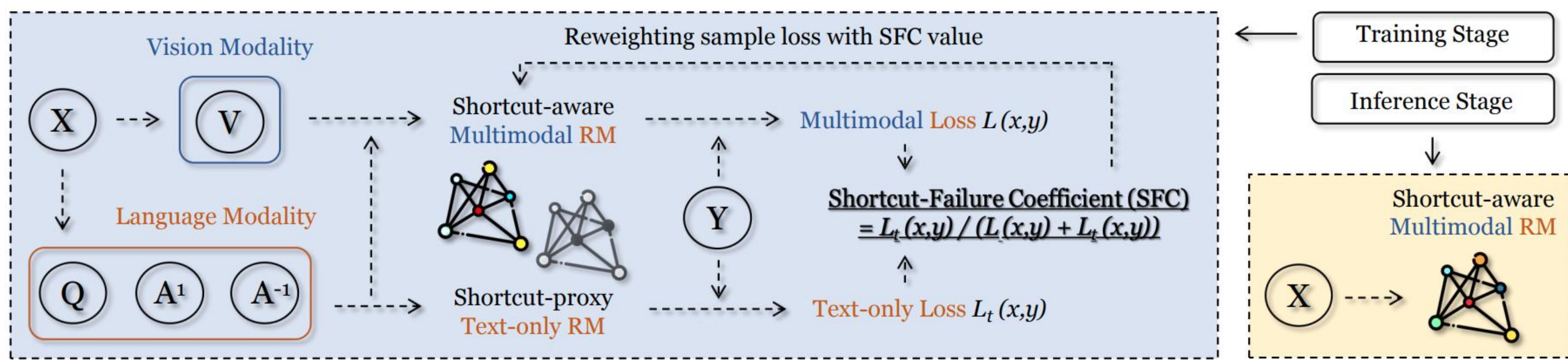
单模态虚假关联

➤ 纯文本捷径 (Text-only Shortcut)：模型过度依赖与训练标签相关但无法泛化的纯文本特征，阻碍了多模态奖励函数的鲁棒学习



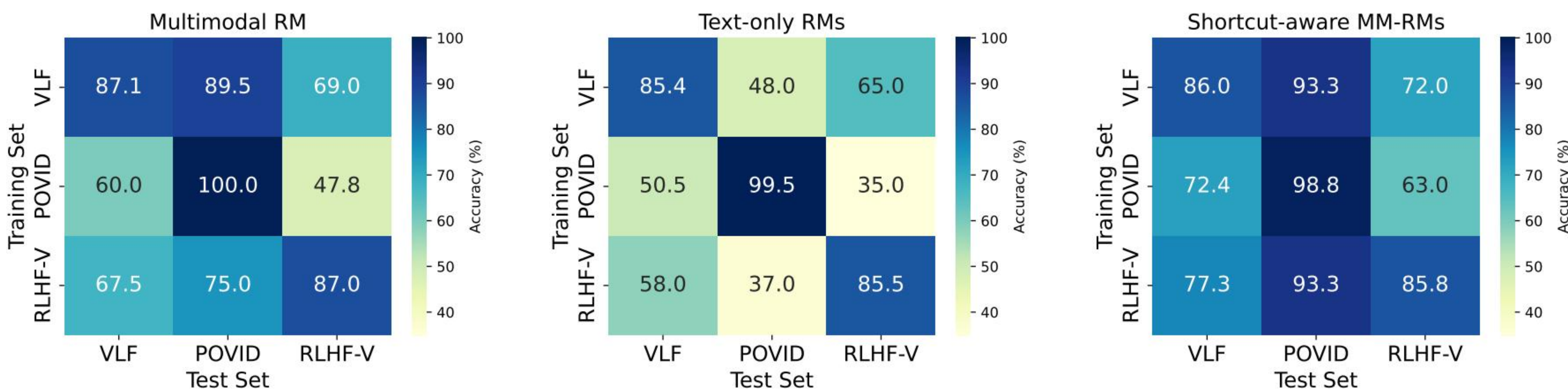
捷径感知学习算法

➤ 通过识别捷径失效的样本，动态赋予样本权重，强调多模态理解



跨分布实验框架

➤ 捷径感知学习算法有效提升了奖励模型在分布外 (OOD) 的泛化能力



(a) Standard MM-RM

(b) Text-only RM (Shortcut proxy)

(c) Shortcut-aware MM-RM (Ours)