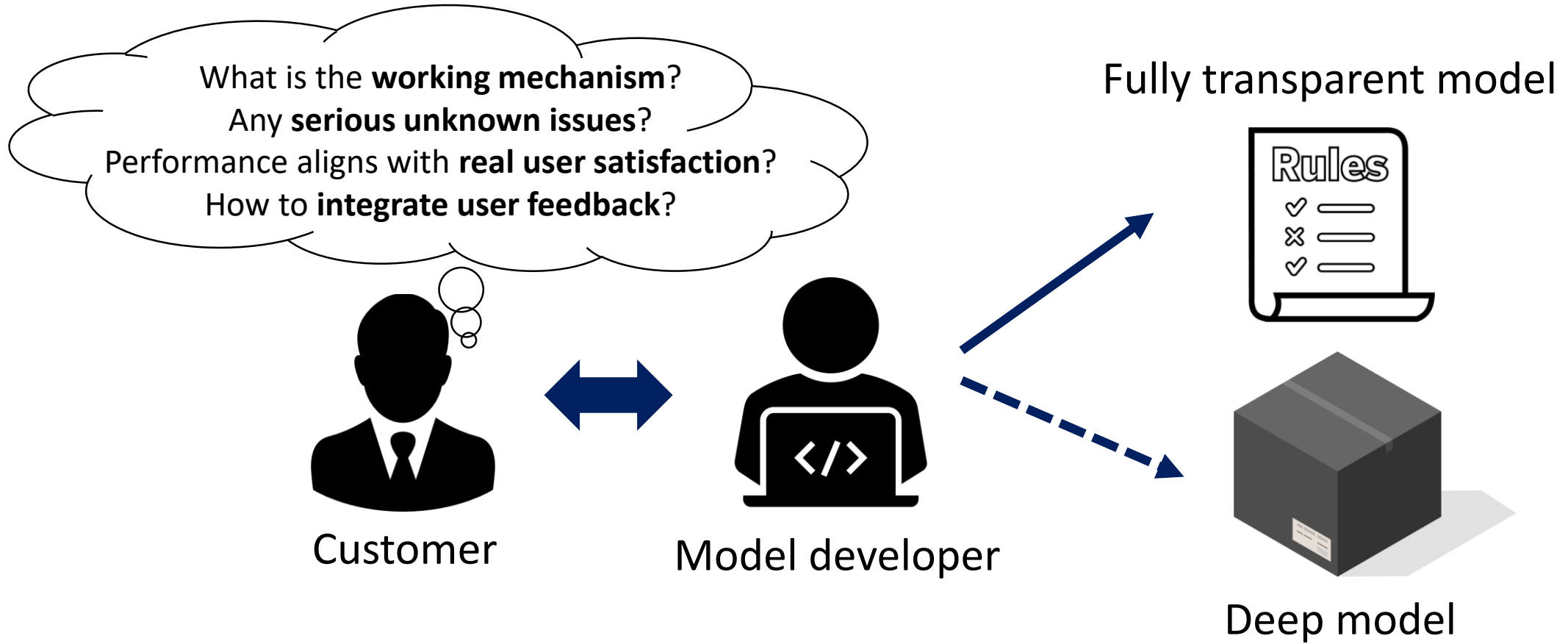# Self-explaining deep models with logic rule reasoning

Seungeon Lee, Xiting Wang, Sungwon Han,

Xiaoyuan Yi, Xing Xie, Meeyoung Cha
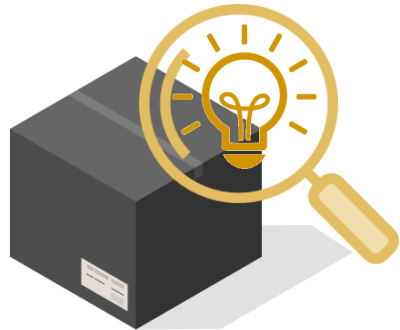
# Trust Issues with Deep Models
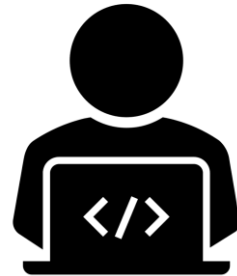


What is the **working mechanism**?
Any **serious unknown issues**?
Performance aligns with **real user satisfaction**?
How to **integrate user feedback**?

Customer

Model developer

Fully transparent model

Rules

Deep model

# Limitation of Post-Hoc Explanations

**Post-hoc** explanations
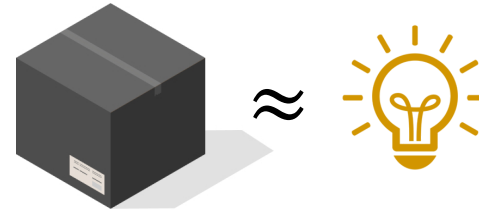
**Trust?**

**Feedback?**

**Can we trust the explanations?**

≈

- Always an approximation [1]
- "General uneasiness" of practitioners [2]

**How to integrate user feedback?**

- No systematic method for direct control
- Requires model retraining
- No guarantee for satisfying user demands

[1] "Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead", *Nature Machine Intelligence*, 2019
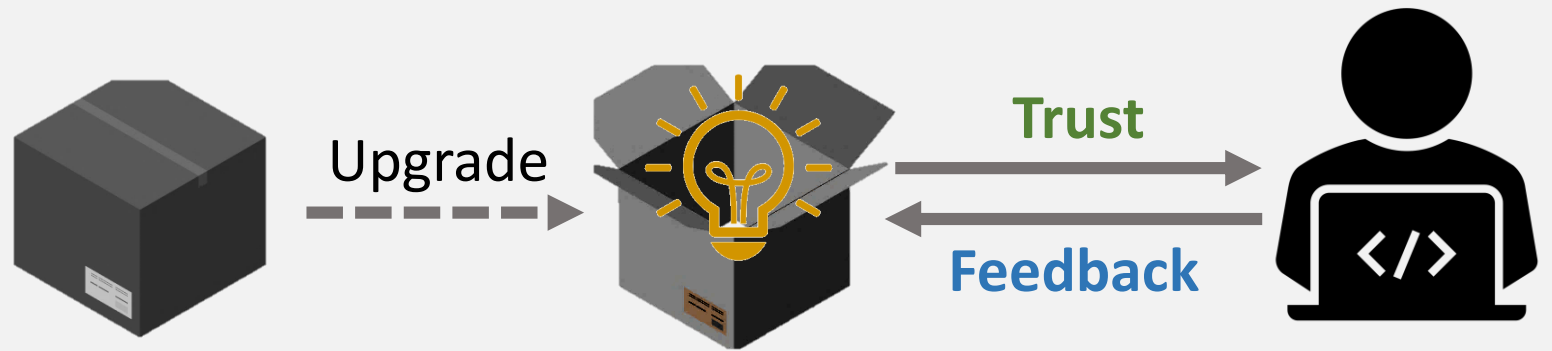[2] "Human Factors in Model Interpretability: Industry Practices, Challenges, and Needs", *ACM HCI* 2020

# Framework: Black-box Model

# Framework: Antecedent Generator

# Results

## High Human Precision



Lime 2%
Anchor 11%
SENN 10%
Ours 66%
RCN 11%

User study
Percentage of best

**+500%**

(Adult dataset)

## Good Prediction Performance



**SELOR** ≈ **Black-box**

## Training Cost

- **Efficient, differentiable** training
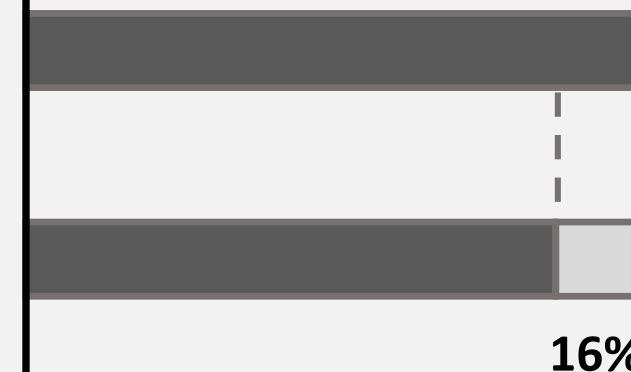- **Slightly slower** than black-box model

**Training Time**

**SELOR**

**BERT**

16%

# Additional Advantages

## Generate Explanation Efficiently
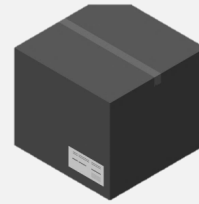
**SELOR vs LIME**
1,000x speed-up

**SELOR vs Anchor**
50,000x speed-up
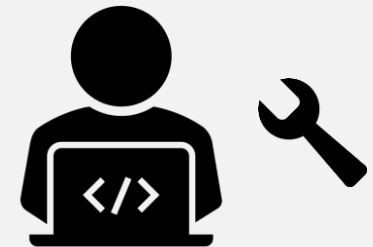(BERT base, Yelp)

## Robust to Noisy Labels

**SELOR** > **Black-box**

## Can be Steered w/o Retraining

❌ *vegas* => positive

✔️ *tasteless* => negative

# Thanks & Questions

**Paper:**
**https://arxiv.org/abs/2210.07024**

**Codes:**
**https://github.com/archon159/SELOR**

**Additional comments and feedback:**
**archon159@kaist.ac.kr**