# Breaking the Gradient Barrier: Unveiling Large Language Models for Strategic Classification
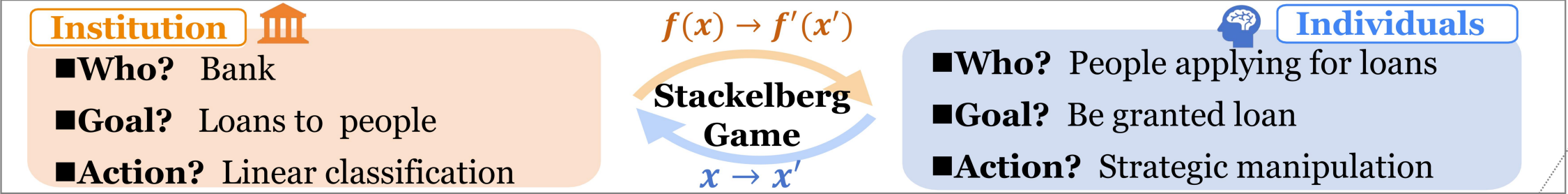
**Xinpeng Lv**

National University of Defense Technology
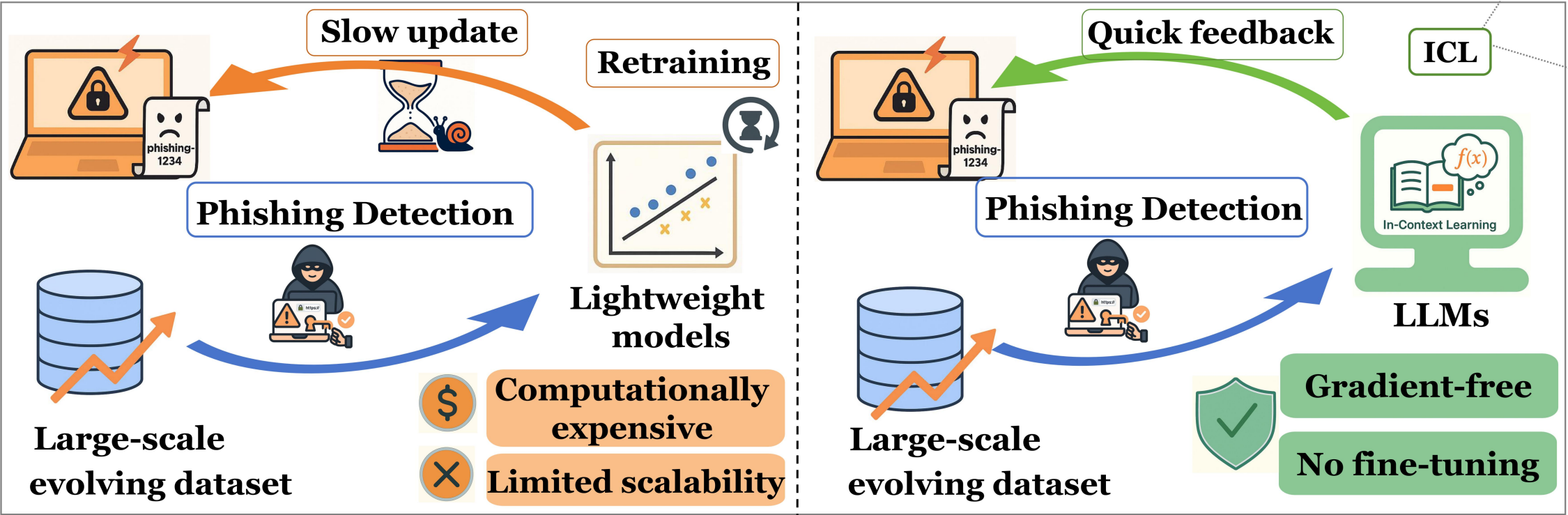
**Our manuscripts**

## An Instantiation of Strategic Classification

**Institution** 🏛️

$f(x) \rightarrow f'(x')$

**Individuals** 🧠

- **Who?** Bank
- **Goal?** Loans to people
- **Action?** Linear classification

Stackelberg Game

$x \rightarrow x'$

- **Who?** People applying for loans
- **Goal?** Be granted loan
- **Action?** Strategic manipulation

## Gradient-aware Method *vs.* Gradient-free Method with LLMs



**Slow update**

**Retraining**

phishing-1234

**Phishing Detection**

Lightweight models

Large-scale evolving dataset

💲 **Computationally expensive**

✖ **Limited scalability**

**Quick feedback**

**ICL**

phishing-1234

**Phishing Detection**

In-Context Learning

**LLMs**

Large-scale evolving dataset

**Gradient-free**

**No fine-tuning**

LLMs are remarkably good at in-context learning—that is, adapting to new examples directly within their prompts, without changing any parameters.

We discover that this in-context mechanism can implicitly perform gradient-like updates inside the self-attention layers.

Based on this insight, we design GLIM, a Gradient-free Learning In-context Method.
GLIM allows an LLM to simulate both stages of strategic classification:

$$\text{Inner Stage } (\textit{Strategic manipulation}): \quad \mathbf{x}' = \arg\max_{x' \in \mathcal{X}} \left[ f(\mathbf{x}') - \lambda c(\mathbf{x}, \mathbf{x}') \right],$$

$$\text{Outer Stage } (\textit{Decision rule optimization}): \quad f^* = \arg\max_{f \in \mathcal{F}} \mathbb{E}_{(\mathbf{x},y)} \left[ \mathbb{1} \left\{ f(\mathbf{x}') = y \right\} \right].$$
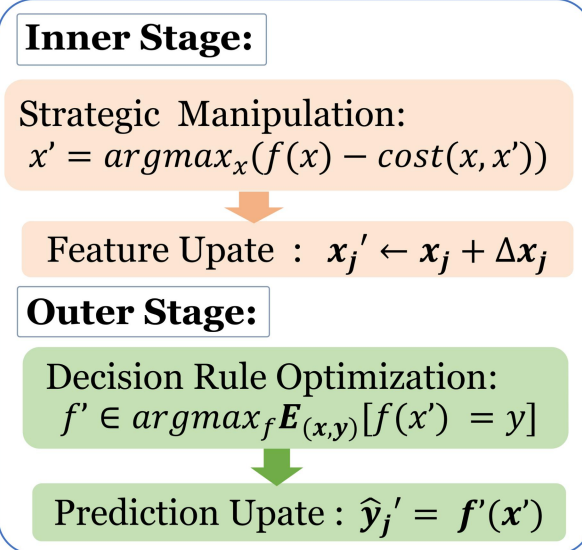
# Our Method

A. Simulating the Inner Stage (Strategic Manipulation)

Goal: Show that the LLM can produce a feature update Δx equivalent to a gradient-based strategic manipulation.
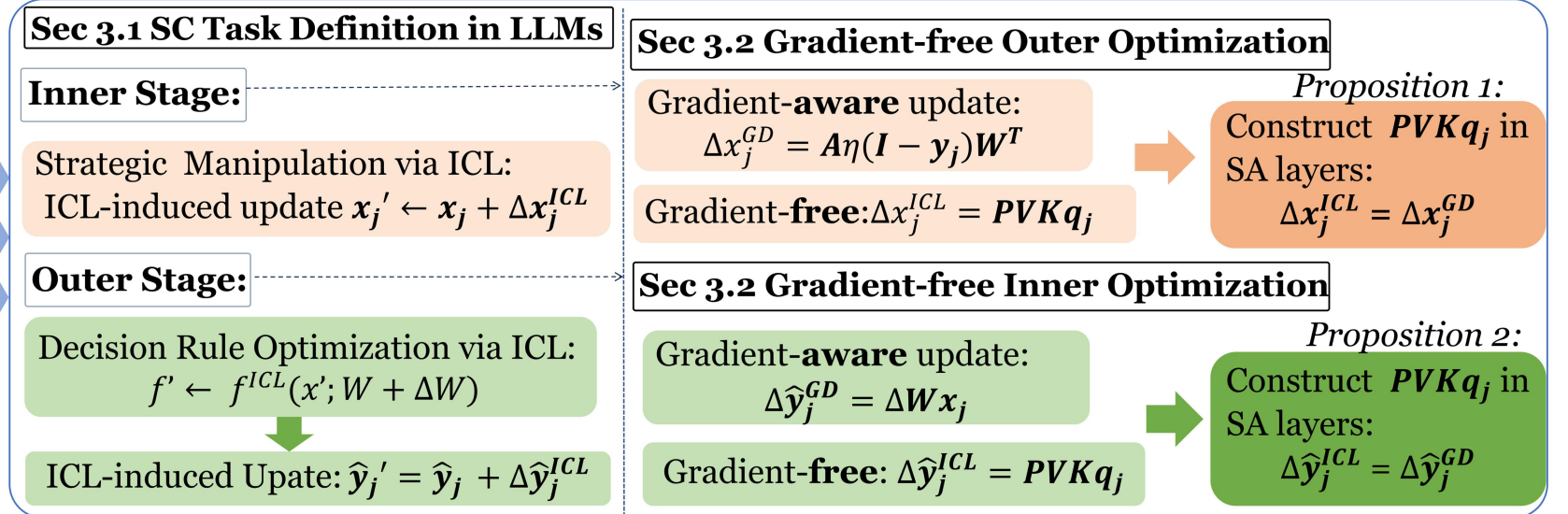
B.Simulating the Outer Stage (Decision Rule Optimization)

Goal: Show that the LLM can adjust its effective decision rule in response to the manipulated features, again without fine-tuning.
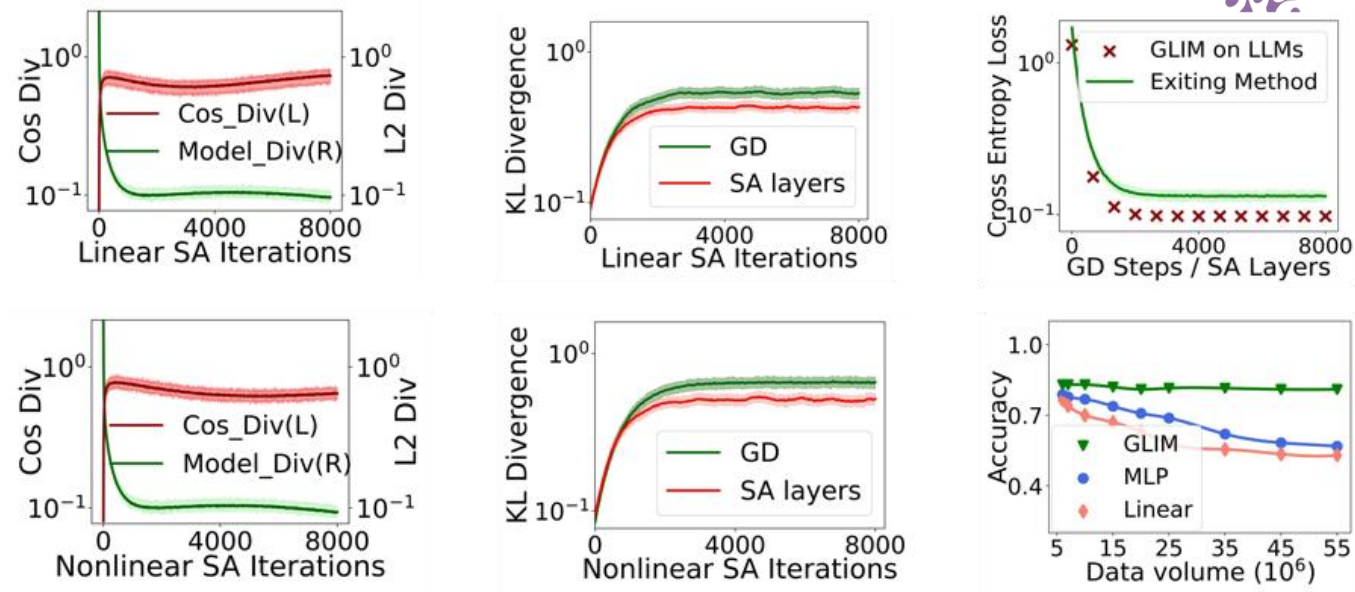


**Bi-Level Optimization in SC**

**Inner Stage:**

Strategic Manipulation:
$x' = argmax_x(f(x) - cost(x, x'))$

Feature Upate : $x_j' \leftarrow x_j + \Delta x_j$

**Outer Stage:**

Decision Rule Optimization:
$f' \in argmax_f E_{(x,y)}[f(x') = y]$

Prediction Upate : $\hat{y}_j' = f'(x')$

**Empower SC task via ICL in LLMs**

Sec 3.1 SC Task Definition in LLMs

**Inner Stage:**

Strategic Manipulation via ICL:
ICL-induced update $x_j' \leftarrow x_j + \Delta x_j^{ICL}$

**Outer Stage:**

Decision Rule Optimization via ICL:
$f' \leftarrow f^{ICL}(x'; W + \Delta W)$

ICL-induced Upate: $\hat{y}_j' = \hat{y}_j + \Delta \hat{y}_j^{ICL}$

Sec 3.2 Gradient-free Outer Optimization

Gradient-**aware** update:
$\Delta x_j^{GD} = A\eta(I - y_j)W^T$

Gradient-**free**: $\Delta x_j^{ICL} = PVKq_j$

*Proposition 1:*
Construct $PVKq_j$ in SA layers:
$\Delta x_j^{ICL} = \Delta x_j^{GD}$

Sec 3.2 Gradient-free Inner Optimization

Gradient-**aware** update:
$\Delta \hat{y}_j^{GD} = \Delta W x_j$

Gradient-**free**: $\Delta \hat{y}_j^{ICL} = PVKq_j$

*Proposition 2:*
Construct $PVKq_j$ in SA layers:
$\Delta \hat{y}_j^{ICL} = \Delta \hat{y}_j^{GD}$

# Some Experimental Results

Linear & Non-linear attention layers

Large-scale & Small-scale datasets



| | Methods | Large-scale Dataset | | | Small-scale Dataset | | |
|---|---|---|---|---|---|---|---|
| | | *PhiUSIIL* | *CISFraud* | *Synthetic* | *Credit* | *Adult* | *Spam* |
| **GLIM (ours)** | | | | | | | |
| *DeepSeek-V3* | Strategic | $85.10_{+0.98}$ | $84.62_{+1.09}$ | $85.15_{+2.18}$ | $89.33_{+0.35}$ | $86.22_{+1.34}$ | $94.85_{+0.67}$ |
| | Non-Strategic | $78.90_{+1.01}$ | $78.74_{+1.14}$ | $80.68_{+2.12}$ | $81.45_{+0.41}$ | $78.77_{+1.33}$ | $89.31_{+0.68}$ |
| *GPT-4o* | Strategic | $86.50_{+0.91}$ | $86.89_{+1.08}$ | $86.83_{+2.35}$ | $89.64_{+0.27}$ | $91.35_{+1.29}$ | $95.97_{+0.61}$ |
| | Non-Strategic | $79.14_{+0.94}$ | $80.15_{+1.10}$ | $81.19_{+2.19}$ | $80.96_{+0.44}$ | $80.23_{+1.31}$ | $91.28_{+0.65}$ |
| *Claude-3.7* | Strategic | $85.07_{+0.95}$ | $84.98_{+1.08}$ | $84.50_{+2.11}$ | $86.51_{+0.31}$ | $88.58_{+1.51}$ | $94.50_{+0.66}$ |
| | Non-Strategic | $78.40_{+0.83}$ | $78.54_{+1.17}$ | $78.89_{+2.00}$ | $80.39_{+0.37}$ | $83.85_{+1.50}$ | $89.50_{+0.61}$ |

# Thanks for your listening

# Contact us:

- Email: lvxinpeng@nudt.edu.cn

**Our manuscripts**