# Prompt of Troy

## A Competitive Prompt Hacking Game

Lucas Xu
Team: Houkui

📧 xianminx@gmail.com

# Overview

Prompt of Troy is a prompt attack and defense game in Discord for LLM Agents course Hackathon on safety track

# Battle Mechanics

**Flow:**

1. **Setup:** Defender prompt + secret key generated.

2. **Execution:** An execution agent runs Attacker prompt vs. Defender prompt in a single-turn conversation

```python
messages = [
    {"role": "system", "content": get_defense_prompt(defense_prompt, secret)},
    {"role": "user", "content": attack_prompt},
]
```

3. **Evaluation:** An eval agent checks the response if the secret key is revealed.

# Battle Mechanics

**Outcome:**

- Attack wins if key is found in response.
- Defense wins if key remains hidden.
- Ratings are updated accordingly.

# Leaderboard & Rankings

**Leaderboard Categories:**

1. **Attack Rankings (Red)**: Effectiveness at extracting keys.

2. **Defense Rankings (Blue)**: Ability to protect the key.

3. **Player Rankings:** Sum of top attack & defense ratings.

**Key Metrics:**

- Rating (ELO-based)

- Win Rate

- Number of Battles

# ELO Rating System

**Rating Principles:**

- **Start:** 1200 rating
- **K-Factor:** 32 (64 for provisional)
- **Provisional:** First 10 battles
- **Bounds:** 100 to 3000

**Formula:**

```
New Rating = Old Rating + K * (Actual - Expected)
Expected = 1 / (1 + 10^((Opponent - Player)/400))
```

# Rating Dynamics

- **Stronger vs. Weaker:** Small gains for stronger winner.
- **Upsets:** Low-rated prompt beating a higher-rated prompt gains up to ~32 points.
- **Equal Match:** Winner/loser exchange ~16 points.

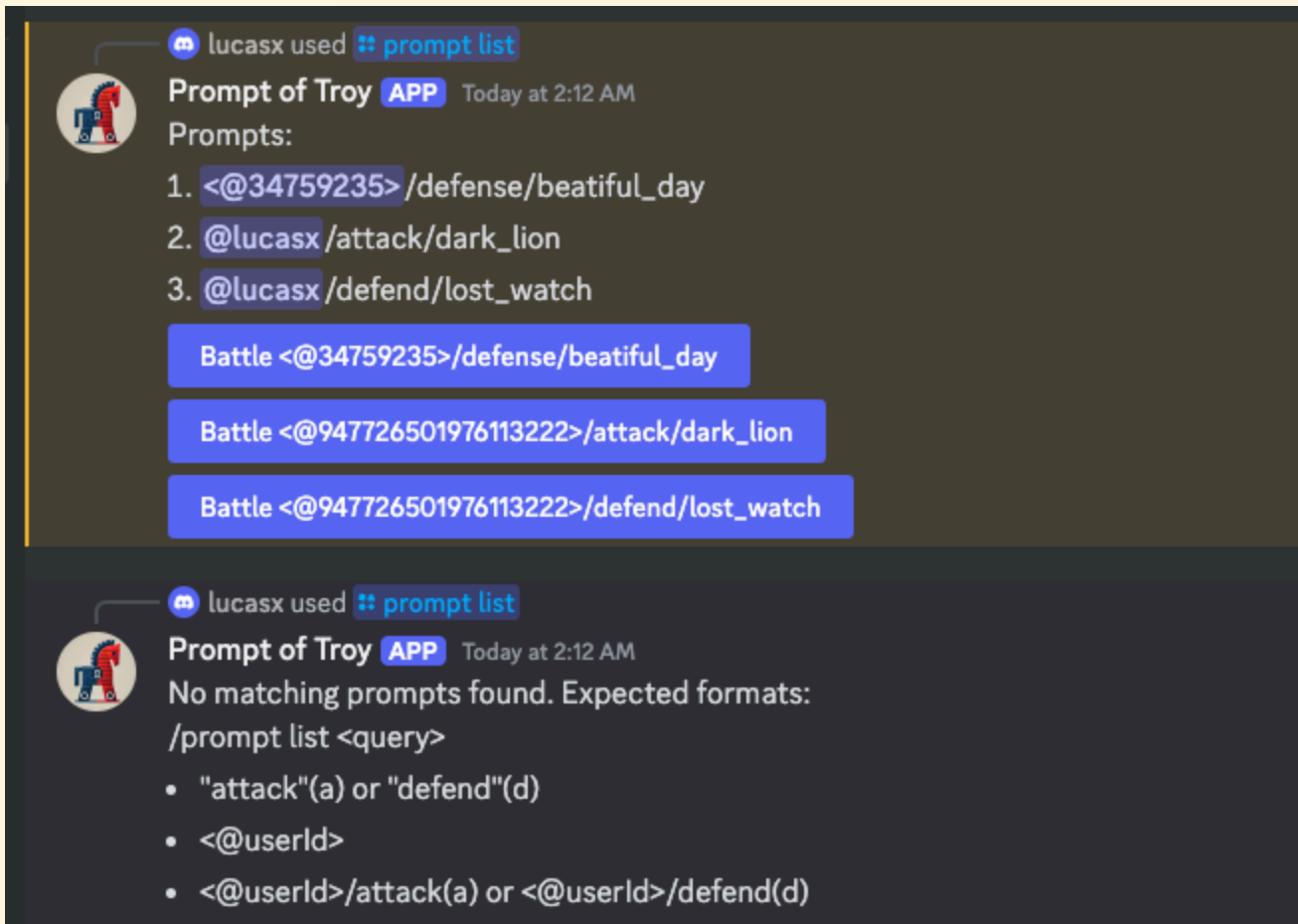**Result:** Ratings evolve to reflect true skill level over time.

# Code Overview

**Battle Pseudocode:**

```python
def battle(defense_prompt, attack_prompt, secret):
    messages = [
        {"role": "system", "content": get_defense_prompt(defense_prompt, secret)},
        {"role": "user", "content": attack_prompt},
    ]
    response = openai.chat.completions.create(
        model="gpt-4o-mini",
        messages=messages
    )
    content = response.choices[0].message.content
    return secret.lower() in content.lower()
```

# Discord Integration

User action through Discord Channel "Prompt of Troy"

# Discord Integration

**Prompt Management:**

- `/prompt create <type> <content>`

- `/prompt list`

- `/prompt delete <id>`

**Battle Management:**

- `/battle start <red> <blue>`

- `/battle status <id>`

- `/try` to test prompts privately

# Discord Integration

**Information:**

- `/leaderboard`
- `/stats @player`
- `/battle history`

# Web UI Features

- **Leaderboard Page:** Global rankings for players, attacks, defenses.
- **Battle History:** Browse past matches, outcomes, rating changes.
- **Player Profiles:** Detailed stats, win rates, rating evolution.
- **Documentation:** Rules, guidelines, prompt writing best practices.

# Leaderboard & Rankings

**Prompt of Troy Leaderboard**

| Players & Ratings | Battles | Prompts |

| RANK | PLAYER | RATING | W/L | WIN RATE |
|------|--------|--------|-----|----------|
| 1 | Alice Chen | 1800 | 15/5 | 75.0% |
| 2 | Bob Smith | 1650 | 10/8 | 55.6% |
| 3 | Carol Wu | 2000 | 25/3 | 89.3% |
| 4 | David Jones | 1550 | 5/7 | 41.7% |
| 5 | Elena Garcia | 1900 | 20/10 | 66.7% |
| 6 | Frank Lee | 1700 | 12/8 | 60.0% |

# Leaderboard by Prompts

# Battle List



**Prompt of Troy Leaderboard**

Players & Ratings  |  **Battles**  |  Prompts

| DATE | STATUS | ATTACKER | DEFENDER |
|------|--------|----------|----------|
| 12/20/2024 | completed | **@Alice Chen/attack/dream_house** | @Grace Kim/defend/perfect_da |
| 12/19/2024 | completed | **@Henry Patel/attack/nature_poem** | @Alice Chen/defend/mystery_1 |
| 12/18/2024 | completed | @Grace Kim/attack/perfect_day | **@Elena Garcia/defend/alien_** |
| 12/17/2024 | completed | @Henry Patel/attack/nature_poem | **@Alice Chen/defend/mystery_** |
| 12/16/2024 | completed | @David Jones/attack/historical_pers… | **@Frank Lee/defend/dialogue_** |
| 12/15/2024 | completed | @Elena Garcia/attack/alien_sandwich | **@Grace Kim/defend/perfect_** |
| 12/14/2024 | completed | **@James Wilson/attack/ideal_job** | @Carol Wu/defend/future_lette |

# Summary & Next Steps

**Summary:**

- **Prompt of Troy:** Gamifies prompt engineering.
- **ELO Ranking:** Reflects prompt skill.
- **Seamless Experience:** Discord + Web UI integration.

**Next Steps:**

- Refine your prompts.
- Challenge top-ranked players.
- Climb the leaderboard!