Be Safe & Play Fair

Let's see what El Primo is doing..

Let's see what El Primo is doing..

Let's see what El Primo is doing..

Let's see what El Primo is doing..

Should el Primo report or not?

**Report ?**

Are you sure you want to report    for inappropriate behavior?

Cancel     Okay

# Current problems

1. AI cannot solve the problem yet, chat in groups

2. Users not reporting ▶ **Increase reporting fluidity**

3. Abusers ▶ **Reduce propensity to abuse**

4. Moderators cost ▶ **Increase efficiency & accuracy**

SUP ERC ELL

INSPECT

# The three-pillars strategy

1 Prompt chat reporting

2 Base-layout filtering

3 Rewarding system

# Two-sided Prompted chat reporting

- Offensive Speech Detection with RoBERTa (pretrained)
- Accuracy of 81% across 70k messages from chats

## Victim view

User xy wrote: Son of a b****!

**Report a message**

We might think this message contains offensive content. You may want to report a message from chat. False reports will be penalized!

Send | Cancel

## Bullies view

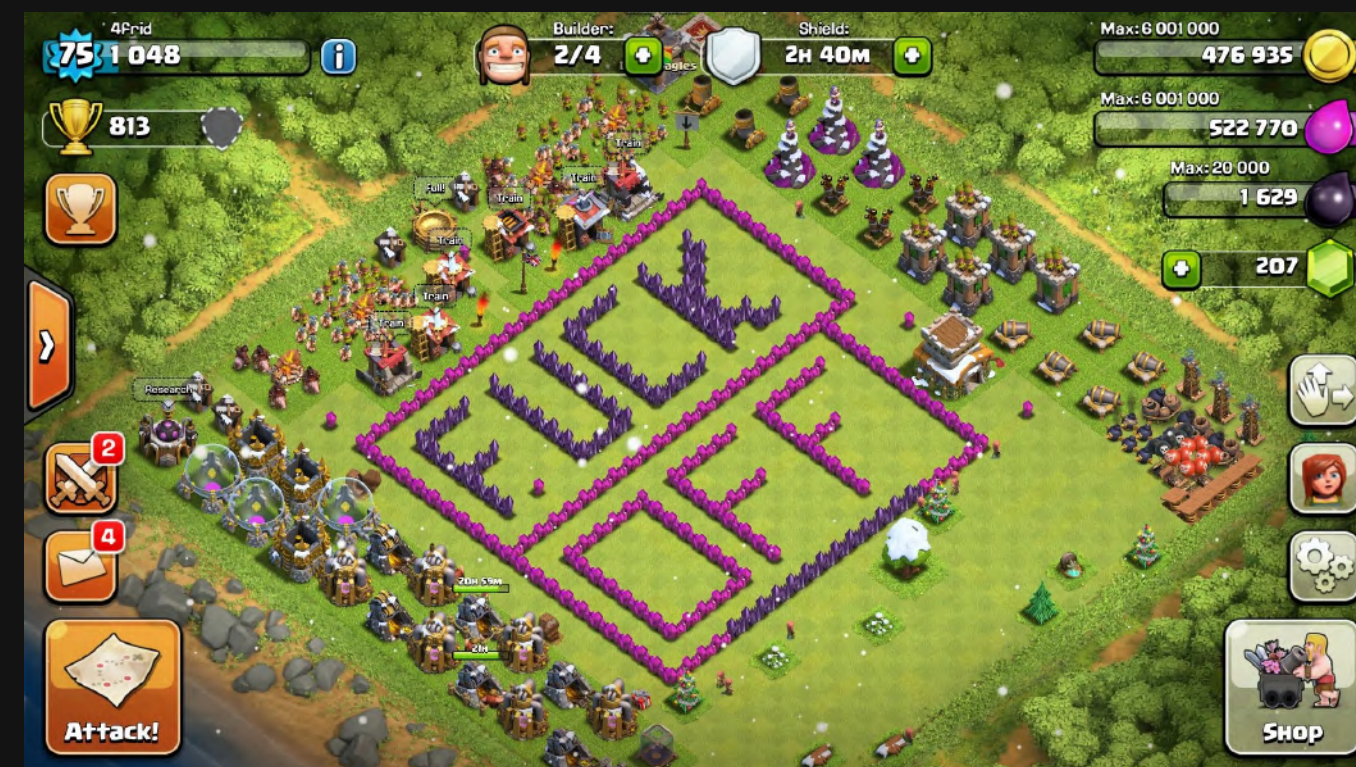Enter a message:

You are such a noob!!

Send

**Warning!**

You are trying to write an offensive message in chat. If you continue, your account may be punished.
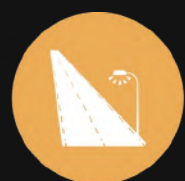
Send | Cancel

SUPERCELL

INSPECT

Try Pitch

# Vision Transformer for offensive bases detection



- Self-collected dataset (approx. 250 images)
- Vision Transformer pretrained on ImageNet
- Fine-tuning of 2 epochs and freezing last layers

INSPECT

Try Pitch

# Vision Transformer for offensive bases detection



- Self-collected dataset (approx. 250 images)
- Vision Transformer pretrained on ImageNet
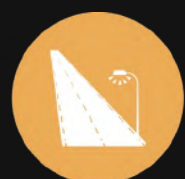- Fine-tuning of 2 epochs and freezing last layers

# Reputation system - key points

1. Play longer and play fair - linear improvement of reputation

2. Offensive behaviours decrease the score

3. Successful reporting - increases score, false reporting - decreases the score

4. Most reputable players help with community interaction and gain rewards, helps moderators.

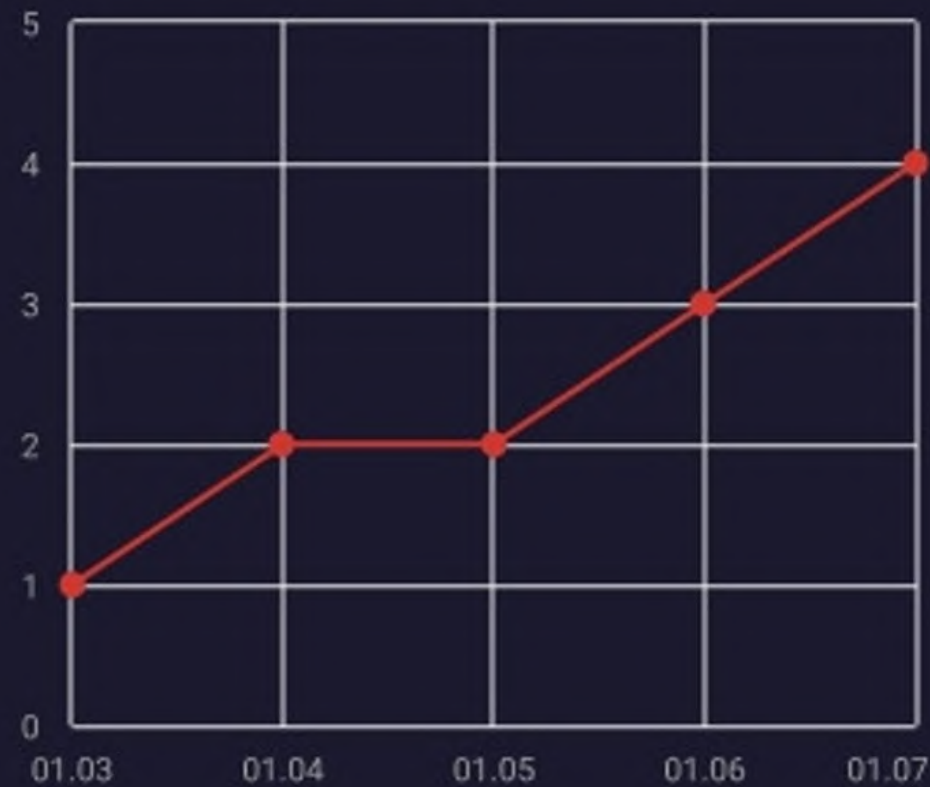Refer to our <u>documentation</u> for a complete overview of the reward system we envision

SUP ERC ELL

INSPECT

# Let's fight offensive cyber behaviour together!

**SUPRC ELL**

**INSPECT**

Try Pitch

*Extend the current offensive-cyber-behaviour detection system to the base-layout and integrate a reward system to nudge good behaviour instead of banning bad actions*