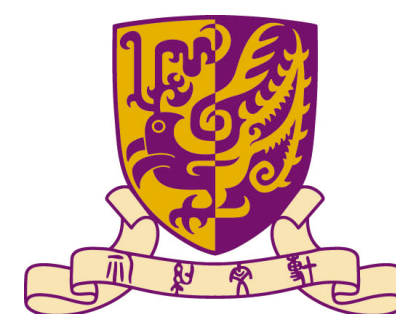# Stealthy Porn: Understanding Real-World Adversarial Images for Illicit Online Promotion

Kan Yuan, **Di Tang**, Xiaojing Liao, XiaoFeng Wang, Xuan Feng, Yi Chen, Menghan Sun, Haoran Lu, Kehuan Zhang
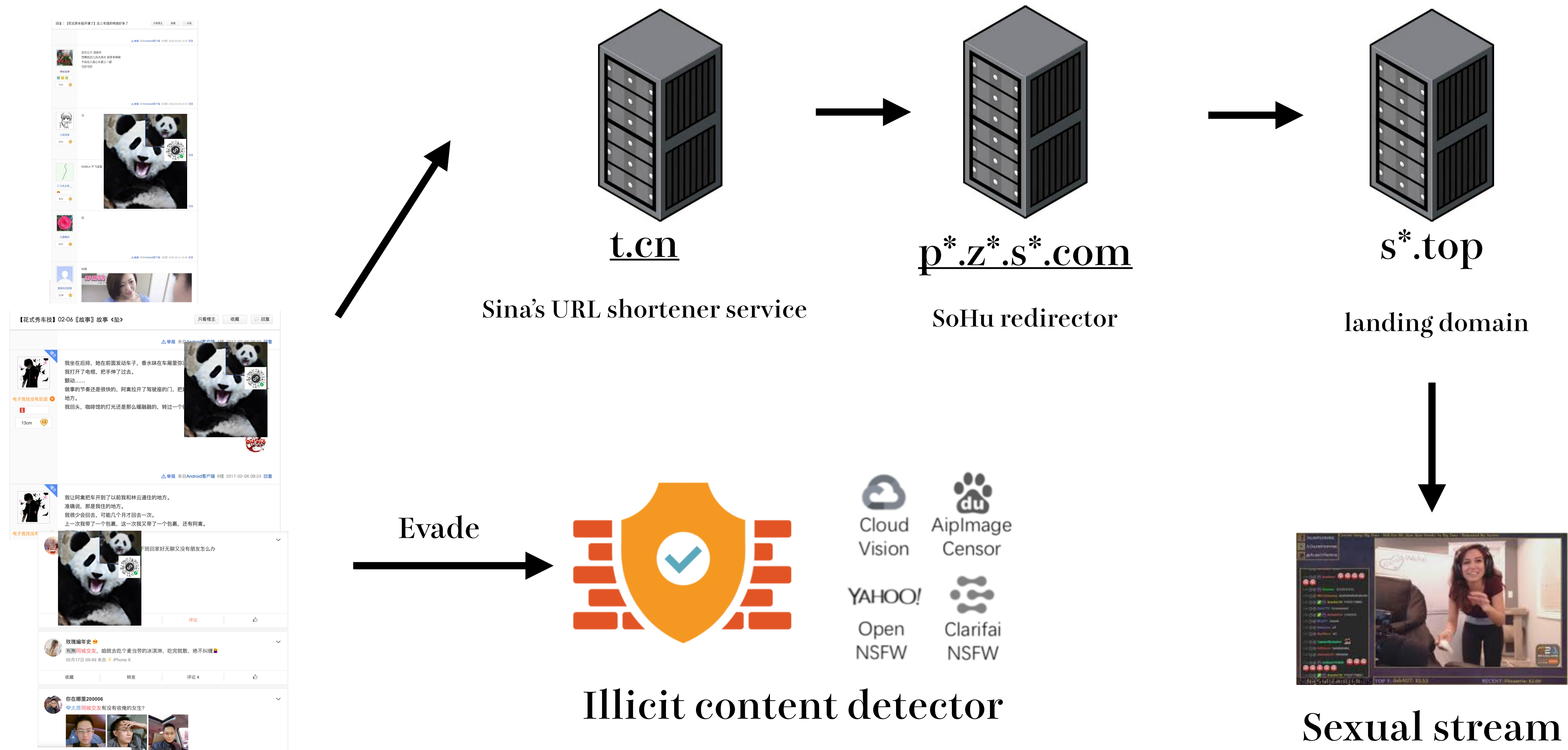
INDIANA UNIVERSITY

香港中文大學
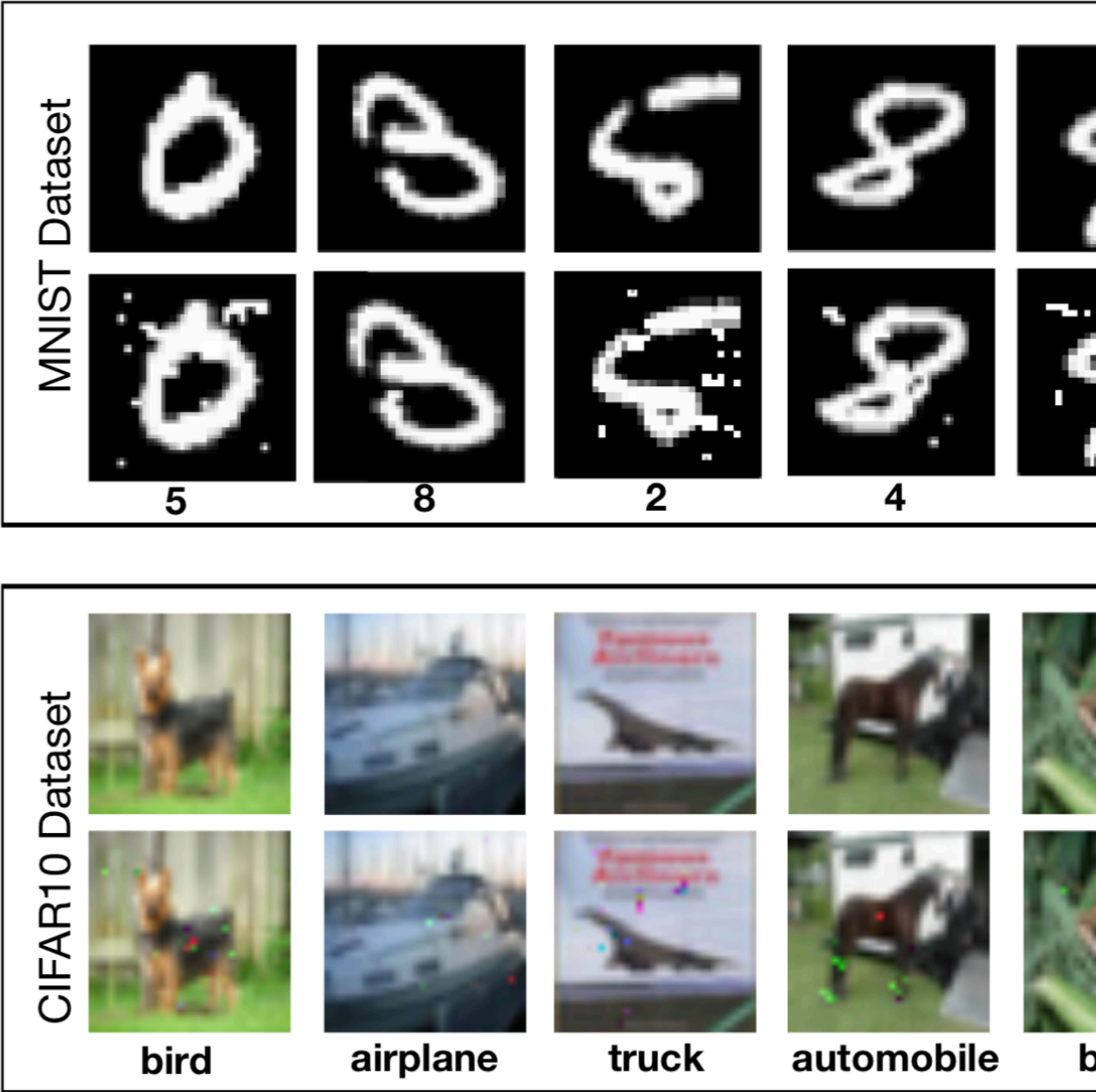The Chinese University of Hong Kong

中國科學院
CHINESE ACADEMY OF SCIENCES

# A Case in the Wild



t.cn

Sina's URL shortener service

p*.z*.s*.com

SoHu redirector

s*.top

landing domain

Evade

Illicit content detector

Cloud Vision

AipImage Censor

YAHOO! Open NSFW

Clarifai NSFW

Sexual stream

## Adversarial Exampl... ...ted Images



- small noise

- nearly indistingui...
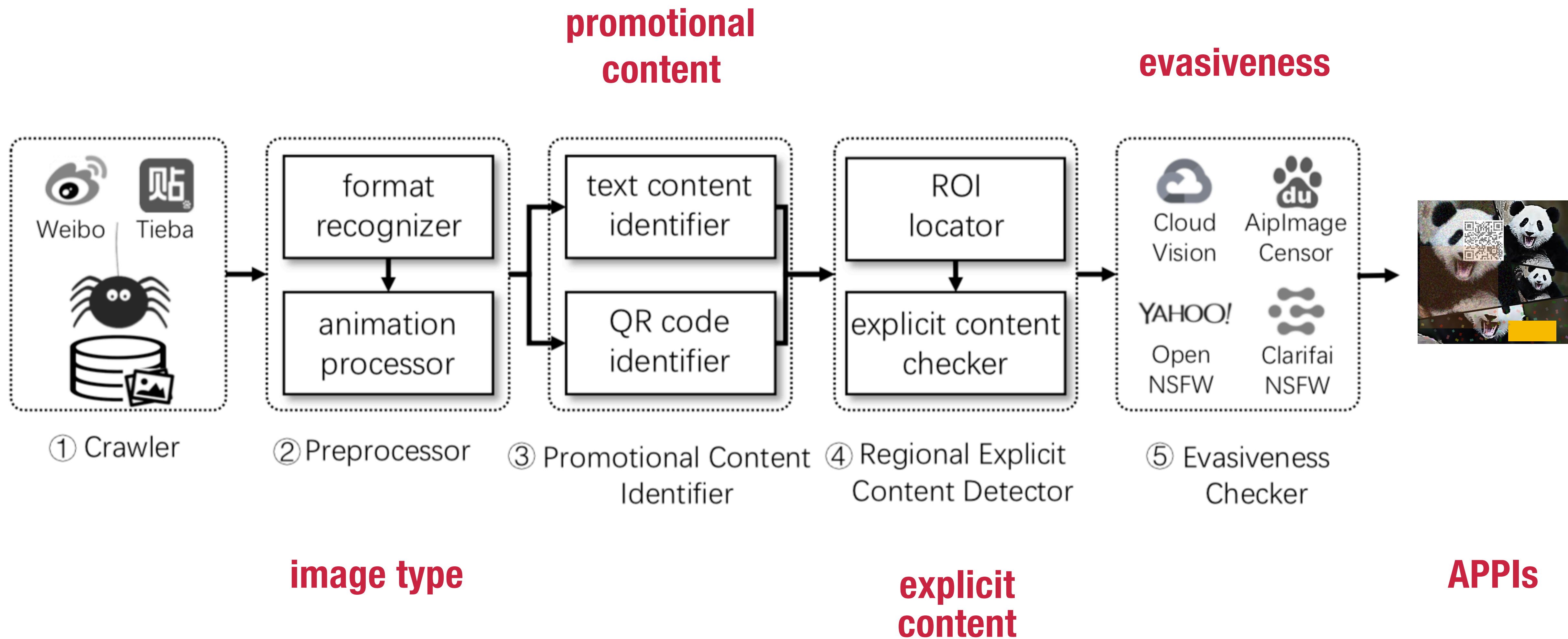
*I am your need*

# Malèna: Finding Stealthy Porn



I am your need

Two common characters:

* Promotional content.
* less obfuscated explicit content

# Malèna: Finding Stealthy Porn



promotional content

evasiveness

| Weibo Tieba | format recognizer | text content identifier | ROI locator | Cloud Vision / AipImage Censor / YAHOO! Open NSFW / Clarifai NSFW |
| --- | --- | --- | --- | --- |
| | animation processor | QR code identifier | explicit content checker | |
| ① Crawler | ② Preprocessor | ③ Promotional Content Identifier | ④ Regional Explicit Content Detector | ⑤ Evasiveness Checker |

image type

explicit content

APIs

# Malèna: Finding Stealthy Porn



promotional content

evasiveness

Weibo    Tieba

format recognizer

animation processor

text content identifier

QR code identifier

ROI locator

explicit content checker

Cloud Vision    AipImage Censor

YAHOO! Open NSFW    Clarifai NSFW

① Crawler

② Preprocessor

③ Promotional Content Identifier

④ Regional Explicit Content Detector

⑤ Evasiveness Checker

image type

explicit content

APPIs

# MALÈNA: FINDING STEALTHY PORN

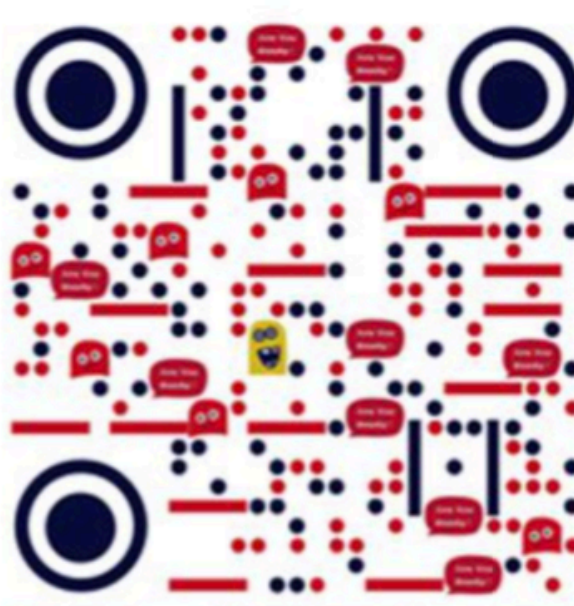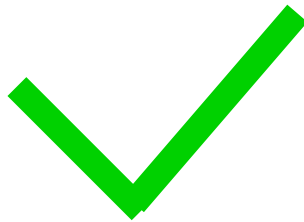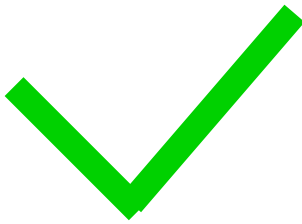# Malèna: Promotional Content Identifier

Text:

请访问跳转链接www.wnorz.com

✓ PiexlLink

QRcode:

✗ ZBar, ZXing, BoofCV

✓ ~~WeChat~~

# Malèna: Promotional Content Identifier

**Text:**

请访问跳转链接www.wnorz.com ✓ **PiexlLink**

**QRcode:**

# Malèna: Finding Stealthy Porn

promotional content

evasiveness

| format recognizer | text content identifier | ROI locator | Cloud Vision / AipImage Censor |
| animation processor | QR code identifier | explicit content checker | YAHOO! Open NSFW / Clarifai NSFW |

Weibo   Tieba

① Crawler ② Preprocessor ③ Promotional Content Identifier ④ Regional Explicit Content Detector ⑤ Evasiveness Checker

image type

explicit content

APPIs

# Malèna: Regional Explicit Content Detector

ResNet-50

Explicit

# Malèna: Finding Stealthy Porn



① Crawler — Weibo, Tieba

② Preprocessor — format recognizer → animation processor

③ Promotional Content Identifier — text content identifier, QR code identifier

④ Regional Explicit Content Detector — ROI locator → explicit content checker

⑤ Evasiveness Checker — Cloud Vision, AipImage Censor, YAHOO! Open NSFW, Clarifai NSFW

promotional content

evasiveness

image type

explicit content

APPIs

# Malèna: Performance

- Performance: 91% precision , 85% recall
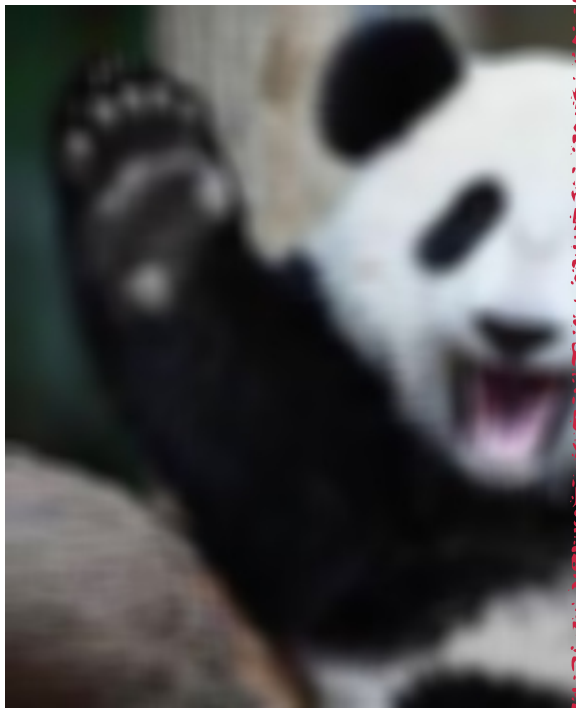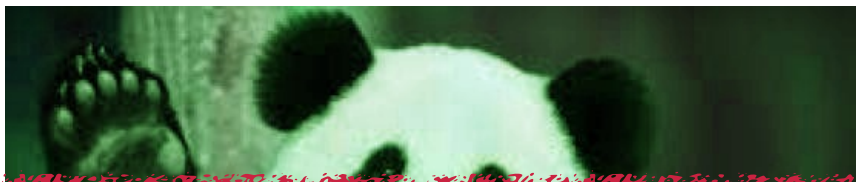
TABLE I: Precision and recall at different stages.

| stage | precision | recall |
|---|---|---|
| promotional content identification | 98% | 90% |
| ROI locator | 89% | 96% |
| explicit content detection | 80% | 93% |
| overall | 91% | 85% |

- Result: 4,353/6,163 APPIs , from 4M images, 76K posts (Baidu Tieba, Weibo)

# MEASUREMENT

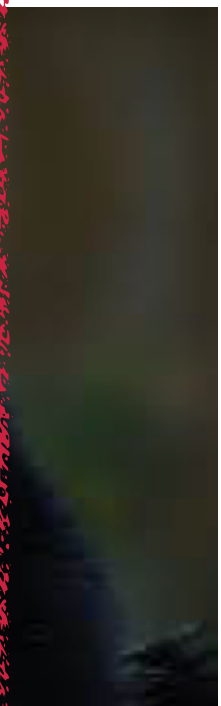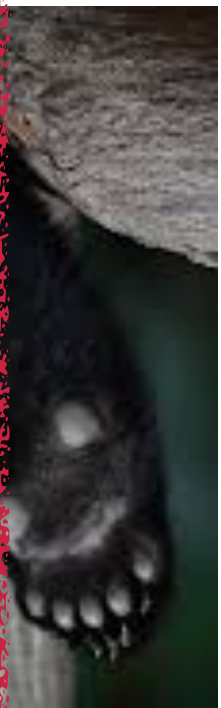✤ Visual pattern.

✤ Promotional content.

✤ Distribution channels.

Blur

Noise

Texture

Transparentization

Occlusion

**TABLE III: The usage of 7 obfuscation techniques.**

| obfuscation technique | # APPI (%) |
|---|---|
| color manipulation | 160 (3.7%) |
| rotation | 1,083 (24.9%) |
| noising | 2,130 (48.9%) |
| texturing | 132 (3.0%) |
| blurring | 829 (19.0%) |
| occlusion | 1,517 (34.8%) |
| transparentization & overlap | 46 (1.0%) |

# Measurement: Visual Pattern

- **Rotation**
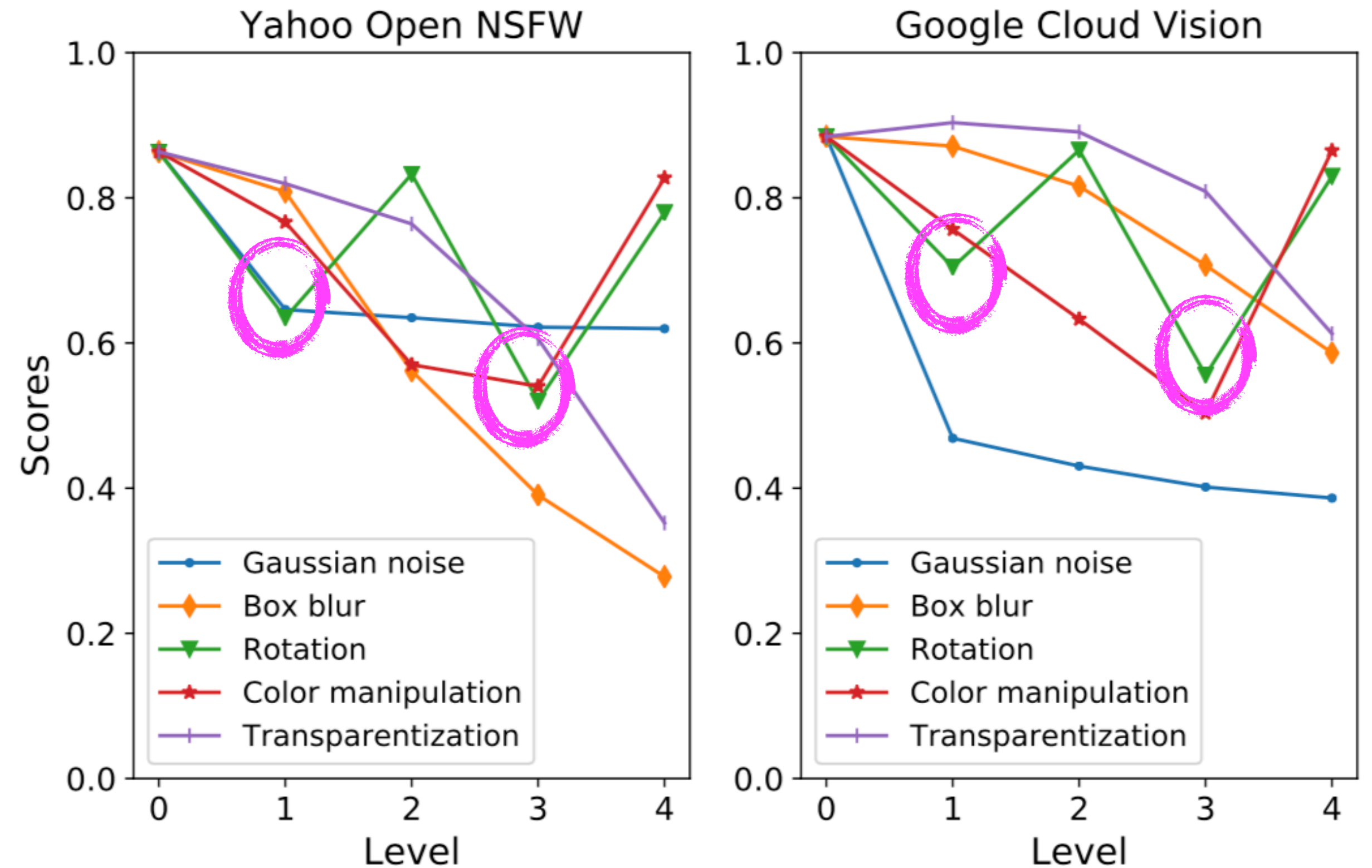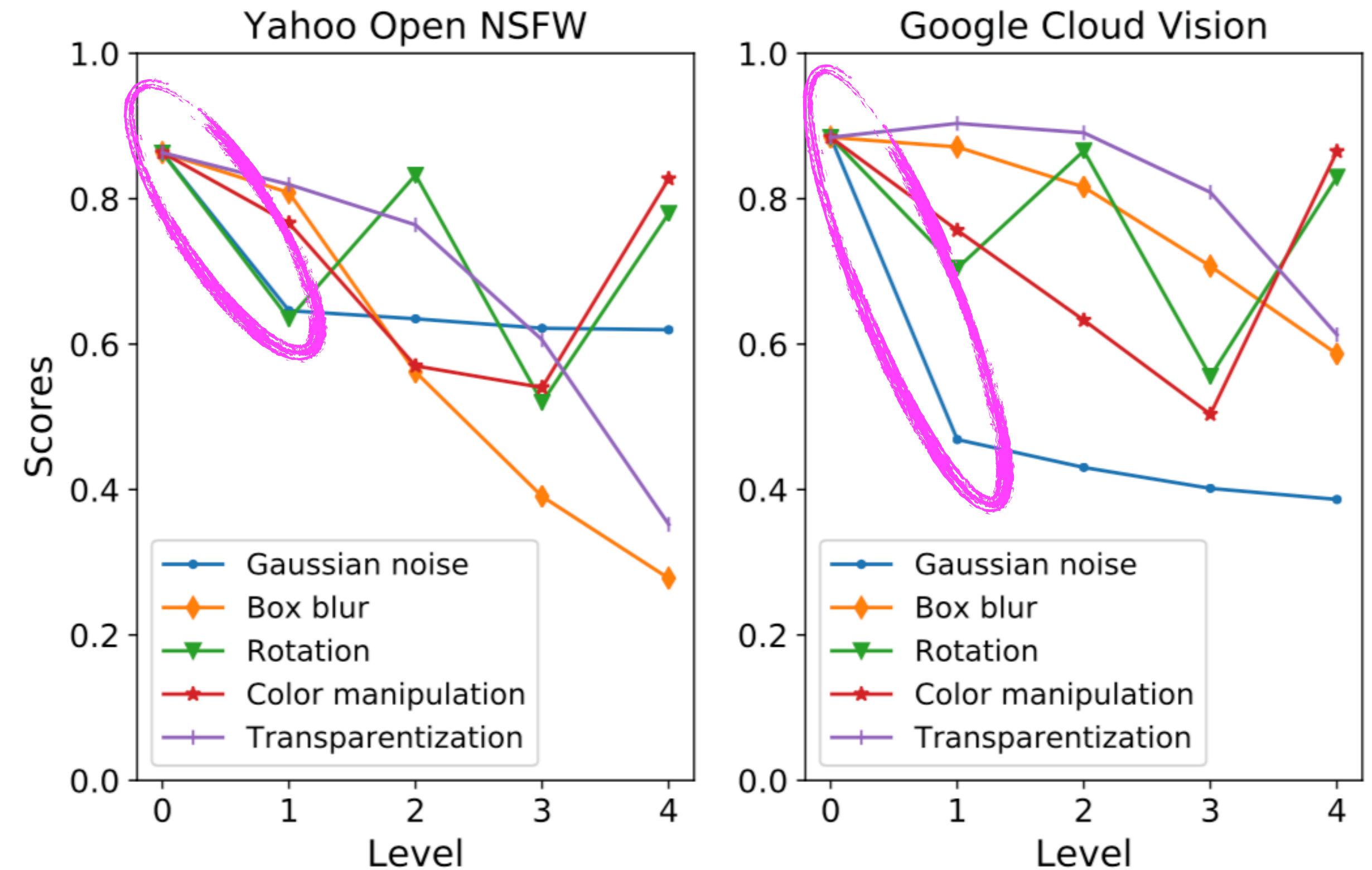  - ➤ **45 and 135 degrees are effective**



Fig. 10: Explicit content detection results on the distorted images.

# Measurement: Visual Pattern

- **Rotation**

  ➤ **45 and 135 degrees are effective**

- **Noising**

  ➤ **Less noising is enough**



Fig. 10: Explicit content detection results on the distorted images.

- **Rotation**
  - ➤ **45 and 135 degrees are effective**

- **Noising**
  - ➤ **Less noising is enough**

- **Color-manipulation**
  - ➤ **Green is evasive colour**



Fig. 10: Explicit content detection results on the distorted images.
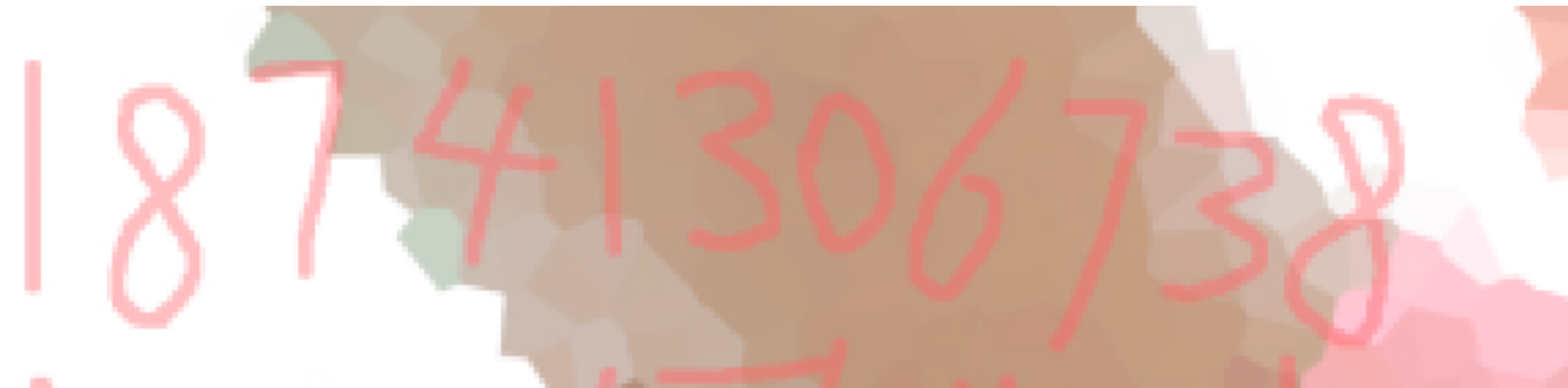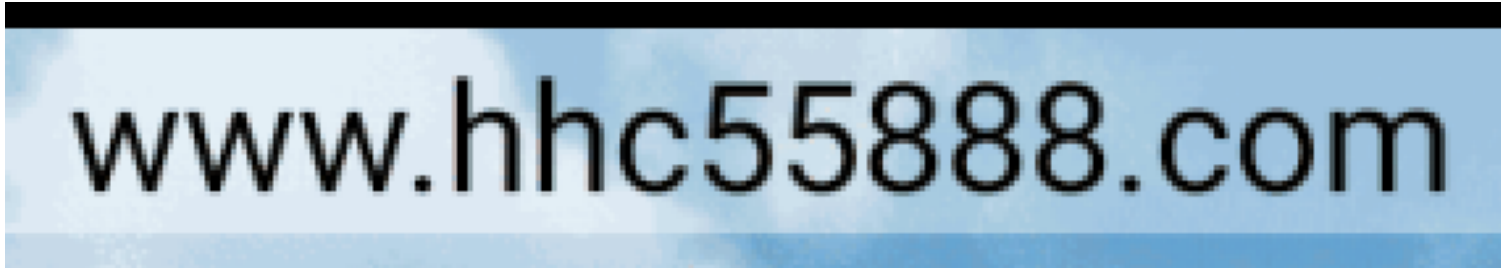
# MEASUREMENT: PROMOTIONAL CONTENT

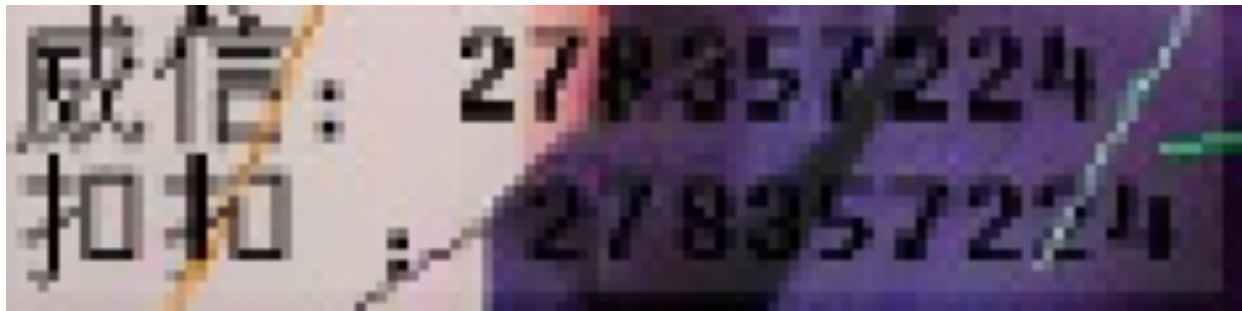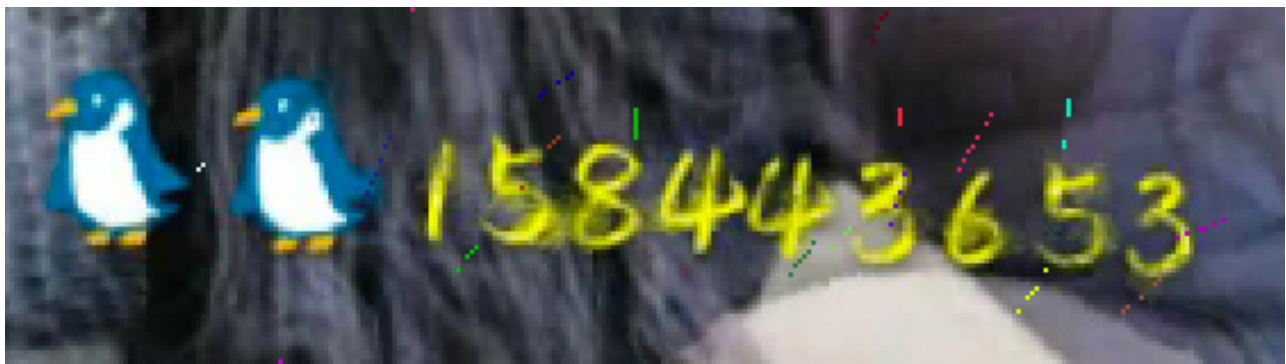TABLE VI: Statistics of promotional content.

| Type | Weibo | Weibo (unique) | Tieba | Tieba (unique) |
|---|---|---|---|---|
| QQ ID | 17 | 7 | 186 | 69 |
| Weibo ID | 375 | 261 | 8 | 5 |
| WeChat ID | 239 | 110 | 1092 | 135 |
| QR code | 0 | 0 | 1430 | 45 |
| URL | 0 | 0 | 85 | 31 |

# Measurement: Promotional Content

TABLE VII: Examples of sensitive text replacement.

| Examples | Type | Meaning | Num |
|---|---|---|---|
| v ♥ | emoji | WeChat | 12 |
| "刊片" | homophonic | porn movie | 10 |
| "企鹅" | jargon | QQ | 18 |
| "呦呦" | jargon | child porn | 8 |
| vx | homophonic+initial | WeChat | 39 |

**Reuse of promotional content**

**232 / 612 (37%)**

**Reuse of explicit content**

**3981 / 4353 (91%)**

Reuse

232

TABLE VIII: Top 5 APPI campaigns.

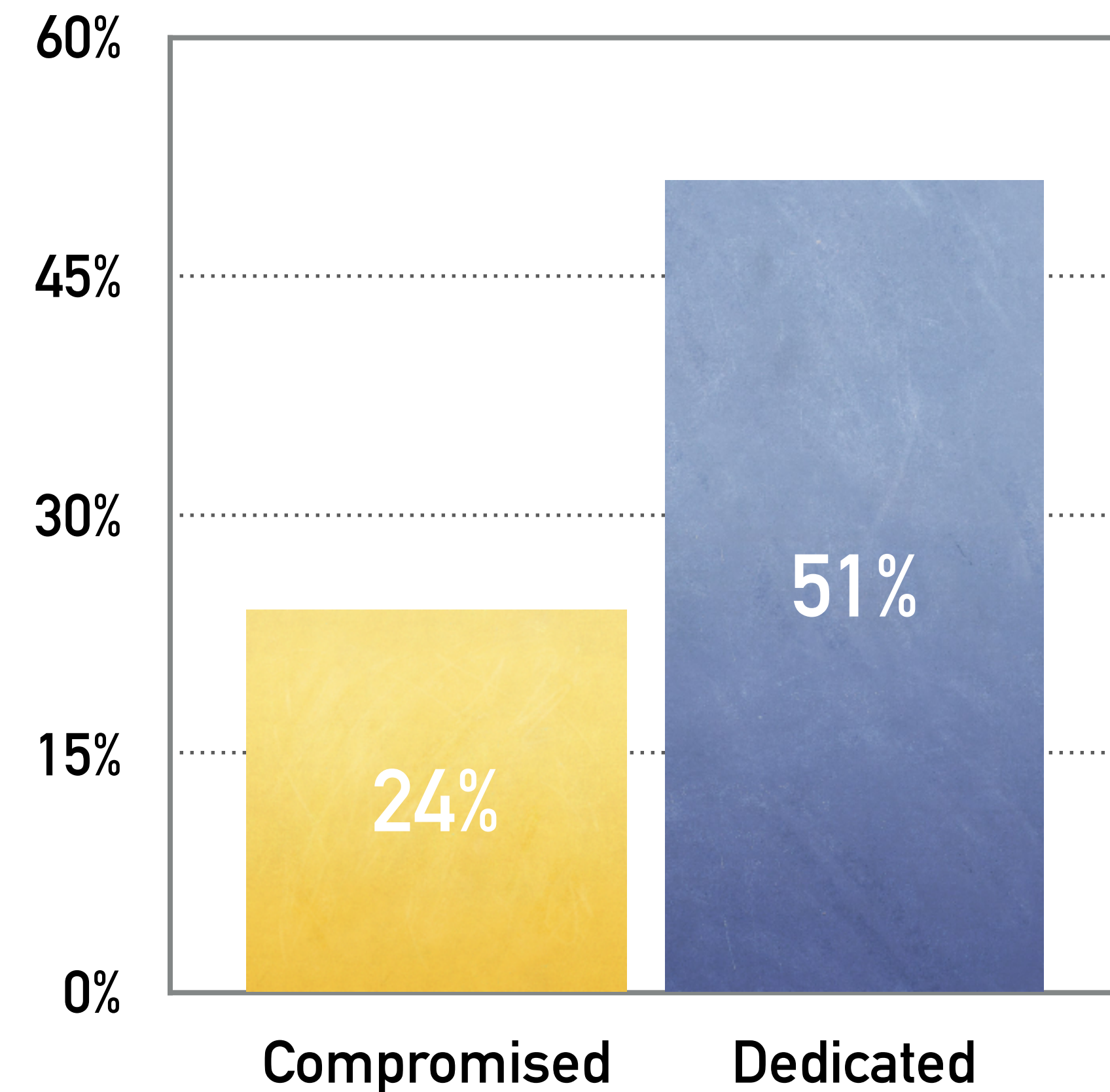| Campaign | # APPIs | Source |
|---|---|---|
| 1 | 1,325 | Tieba |
| 2 | 786 | Tieba |
| 3 | 347 | Weibo |
| 4 | 39 | Weibo&Tieba |
| 5 | 25 | Tieba |

1584443653

ntent

3981 / 4353 (91%)

# MEASUREMENT: DISTRIBUTION CHANNELS

Compromised accounts:

- rarely post
- comment only on hot microblog

Dedicated accounts:

- > 30 posts/day
- with meaningless sentences

# LESSION LEARNED

Visual pattern. $\longrightarrow$ Harden current models.

Promotional content. $\longrightarrow$ Regularize promotion channel.

Distribution channels. $\longrightarrow$ Secure accounts.

# Take-aways

* APPIs are prevalent

* Understanding criminal goal and ecosystem behind adversarial images

* Hardening machine learning model against APPI attack deserves further studies

# Thank You!