

Screen Perturbation: Adversarial Attack and Defense on Under-Screen Camera

Hanting Ye^{}, Guohao Lan^{*}, Jinyuan Jia[†], Qing Wang^{*}*

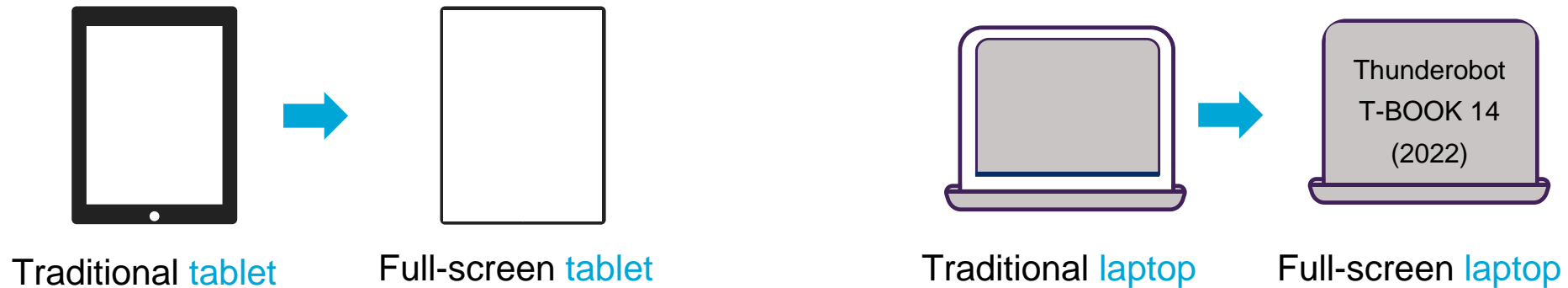
^{}Delft University of Technology*

[†]The Pennsylvania State University

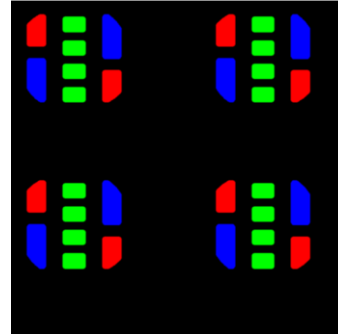
October 3rd, 2023



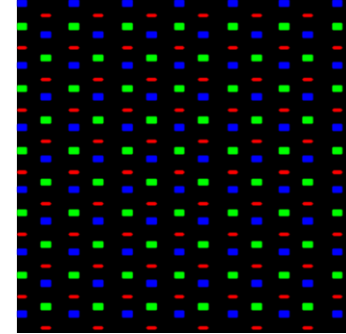
Evolution of smartphone screens



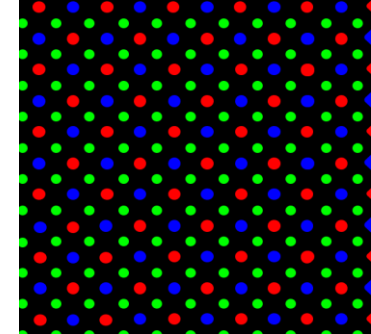
Under-screen camera



Samsung



ZTE

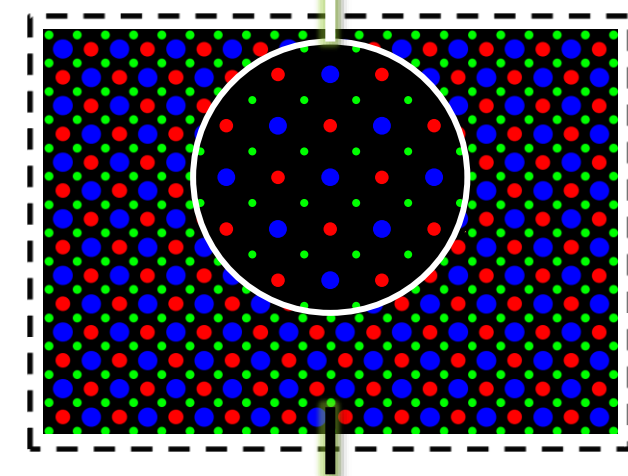
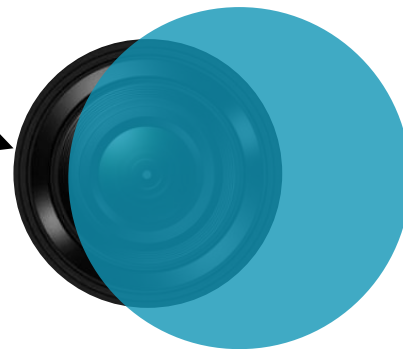


Xiaomi

Full-Screen Smartphone

Under-Screen Camera

Translucent Screen Region



Normal Screen Region

Images formed in different scenarios

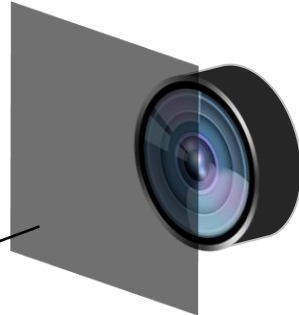
Conventional Camera



Under-Screen Camera



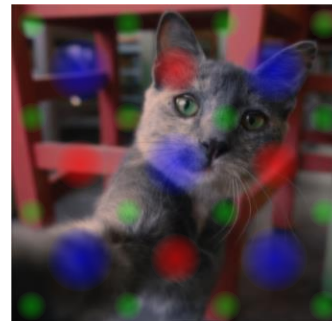
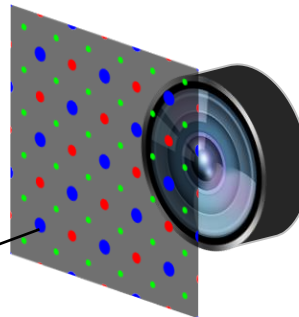
Inactive translucent screen



Under-Screen Camera



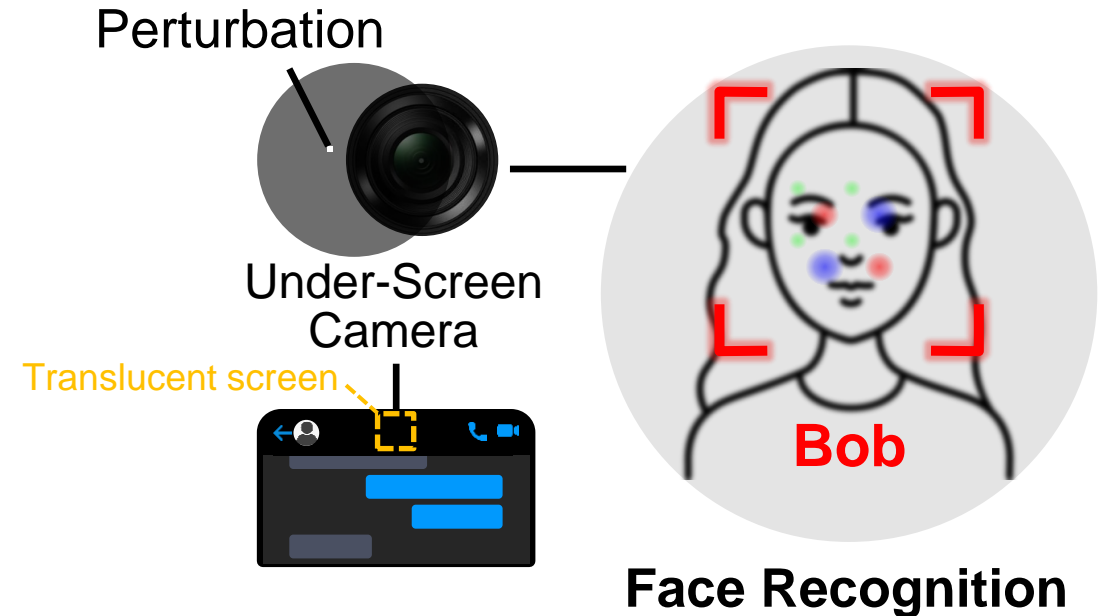
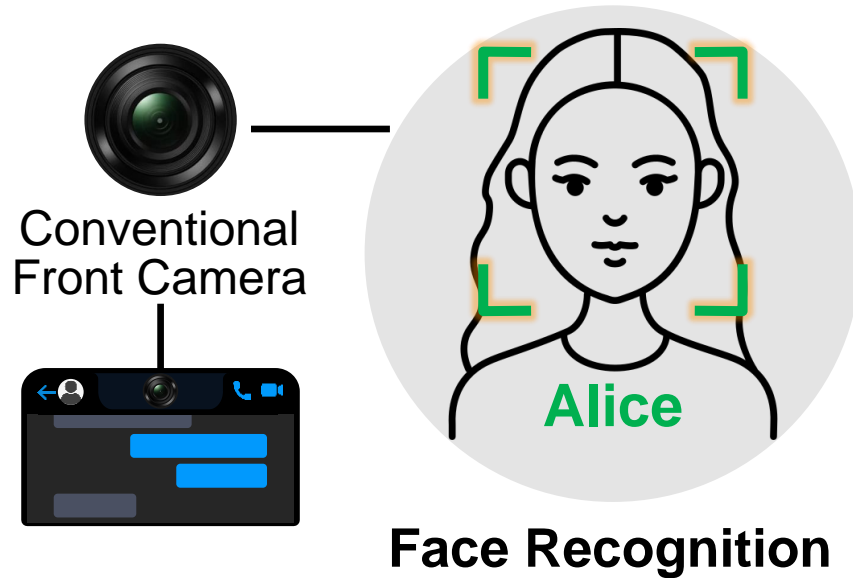
Active translucent screen



+ Passive screen perturbation

+ Active screen perturbation

Screen perturbation



Defenders can proactively activate screen-pixels to thwart unauthorized ML models



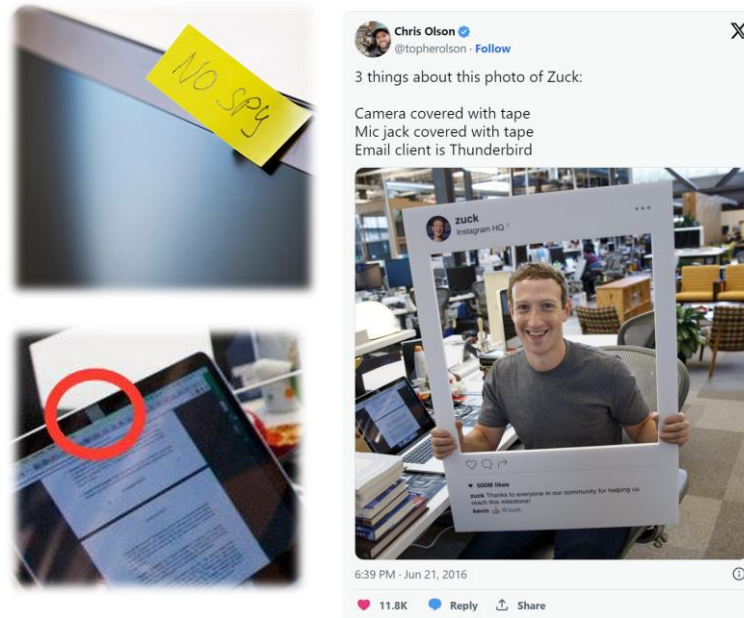
Attackers can manipulate imperceptible screen-pixels to disrupt legitimate ML models

Screen perturbation: A software-defined perturbation that modifies the pixels displayed on the translucent screen region to nullify ML models

Key difference for screen perturbation

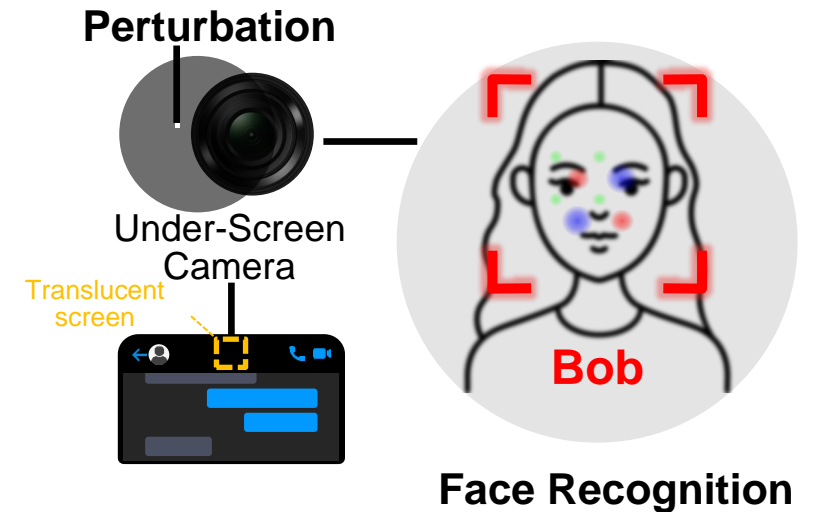


Privacy fears as MILLIONS of photos used to train facial recognition AI without users' consent (Source: Daily Mail Online)



Camera cover solution:

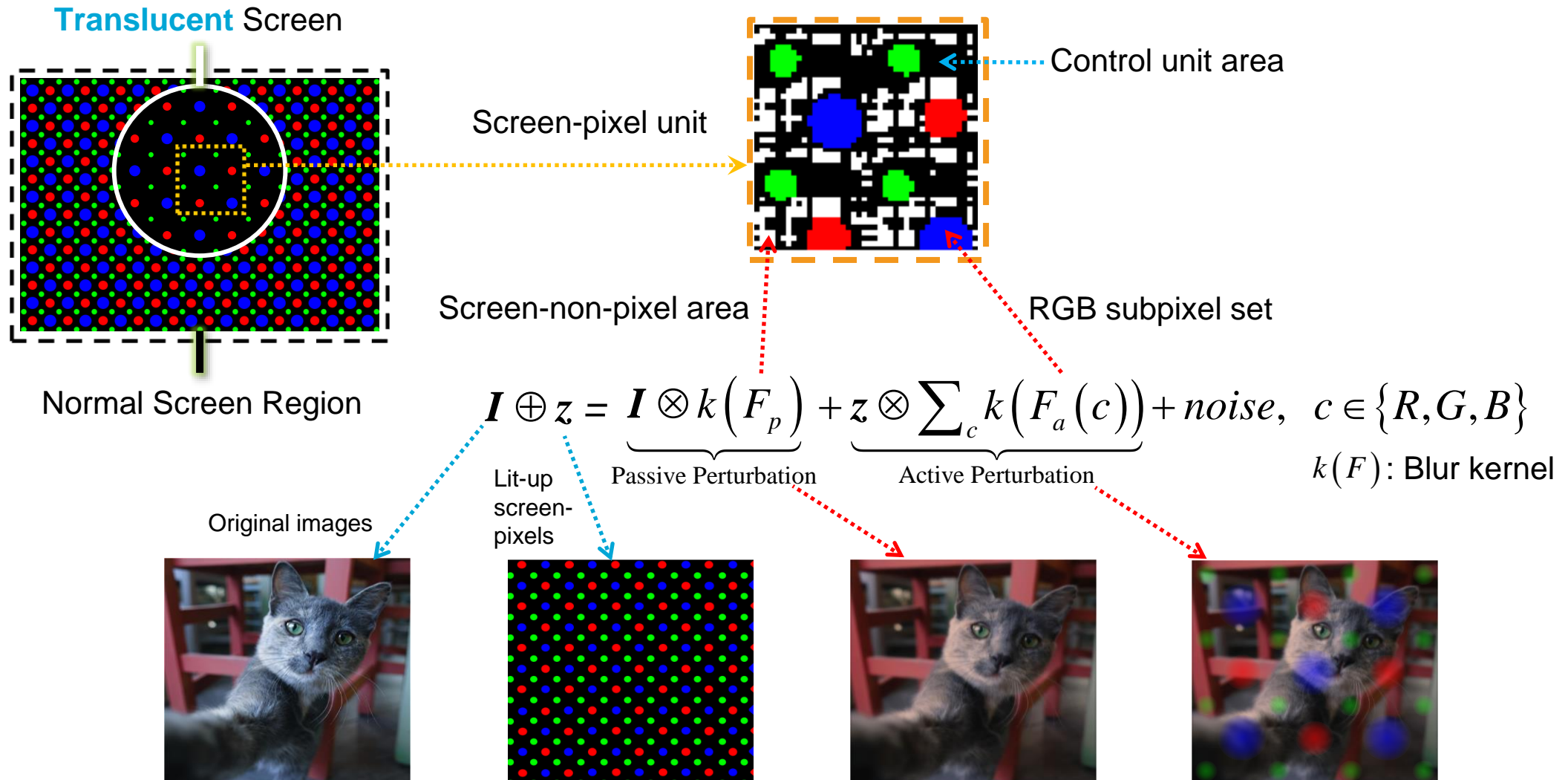
- 1) Secure but block everything
- 2) Obstruct the screen display



Proposed screen perturbation:

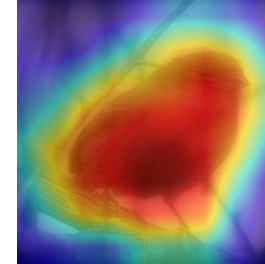
- 1) *Software-defined solution*
- 2) *Imperceptible changes*

Image formation model of under-screen camera

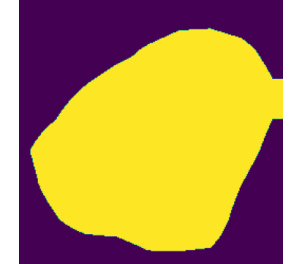


Chromaticity destruction

Q1: How to **localize** the region that has the highest influence on the decision making of the ML model?

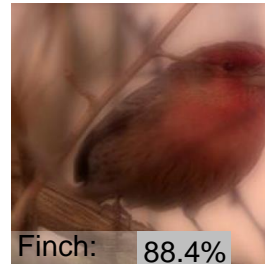


Grad-CAM heatmap



Attack region mask

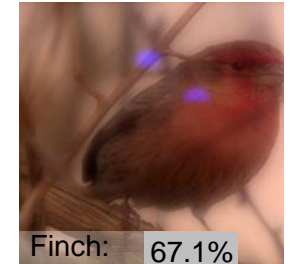
Q2: Which **color** should we pick for increasing adversarial strength?



Passive perturbation
Finch: 88.4%

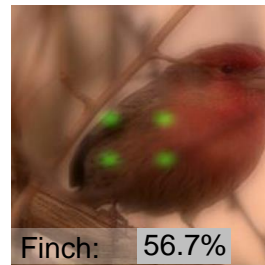


Red subpixel set
Finch: 93.6%

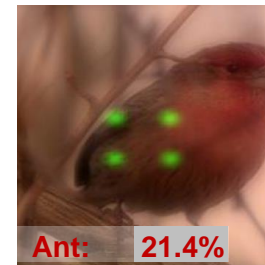


Blue subpixel set
Finch: 67.1%

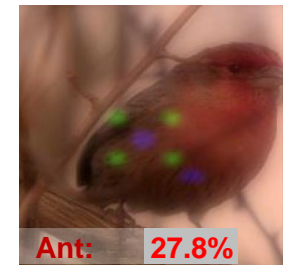
Q3: **One** color or **Two** colors



Green, 30% bright
Finch: 56.7%



Green, 80% bright
Ant: 21.4%



G & B, 10% bright
Ant: 27.8%

Morphology destruction

New screen configuration to generate screen perturbation

Original screen configuration of no screen perturbation

Growth rate

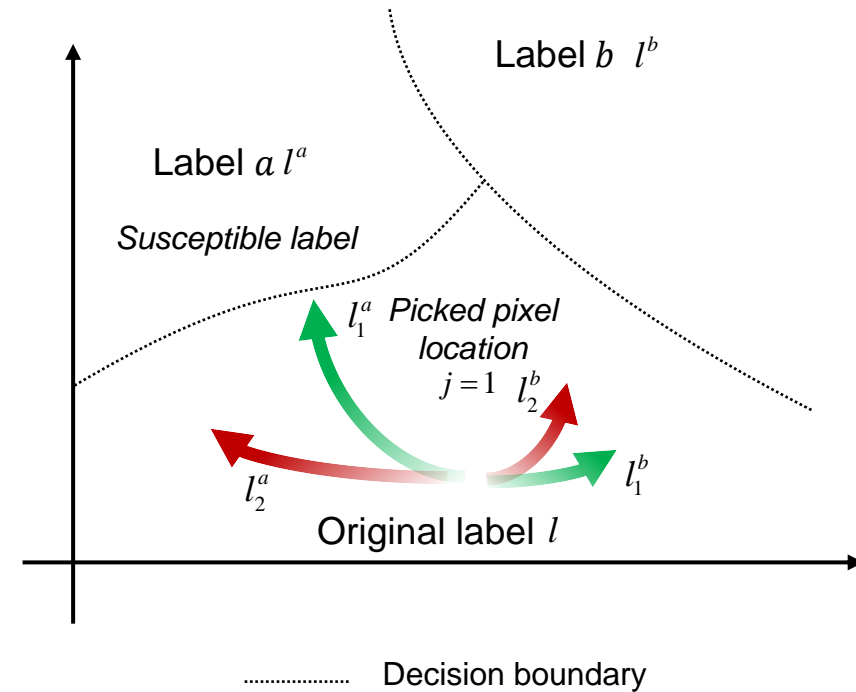
$$s(\hat{l}, j) = \frac{\Phi(I \oplus z_{j,c,b}, \hat{l}) - \Phi(I \oplus z, \hat{l})}{\Phi(I \oplus z, \hat{l})}$$

j : screen-pixel **location**
 c : screen-pixel *color*
 b : screen-pixel *brightness*

The most susceptible label

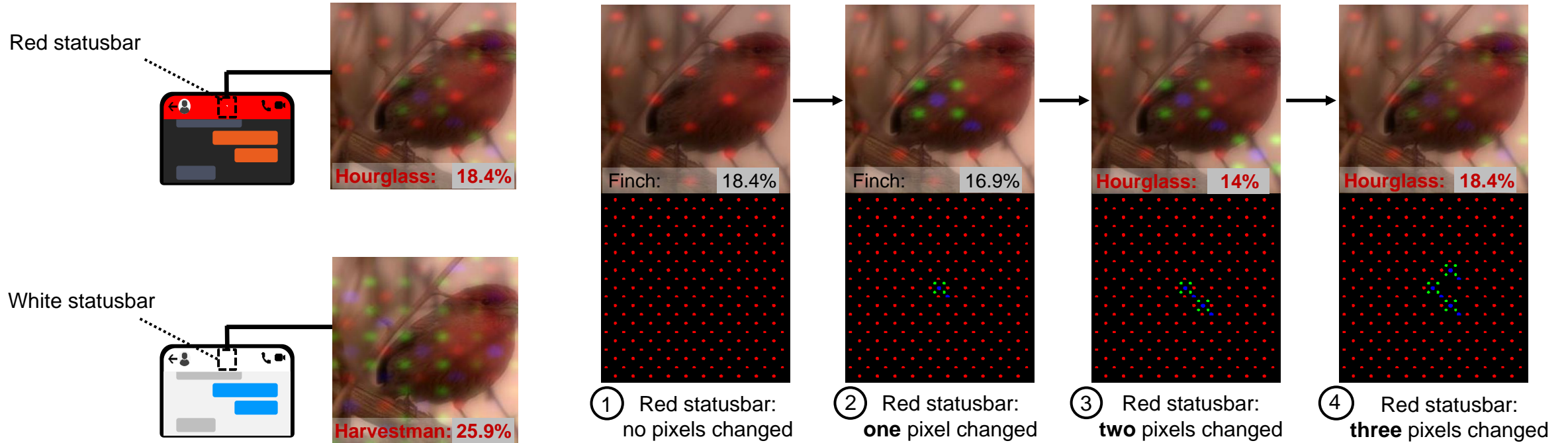
$$\tilde{l}_j = \arg \max_{\hat{l}} s(\hat{l}, j)$$

Picked screen-pixel location

$$j = \arg \max_{\hat{j}} s(\tilde{l}_j, \hat{j})$$


$\Phi(I, l)$: Predicting probability of classifier Φ in classifying the image I with label l

Multiple-pixel perturbation



*Regional aggregation effect: **Adjacent** screen-pixel units have similar effect on the probabilities of predicted labels*

Built testbed & images results



ZTE AXON30

Xiaomi MIX4

Samsung Fold4

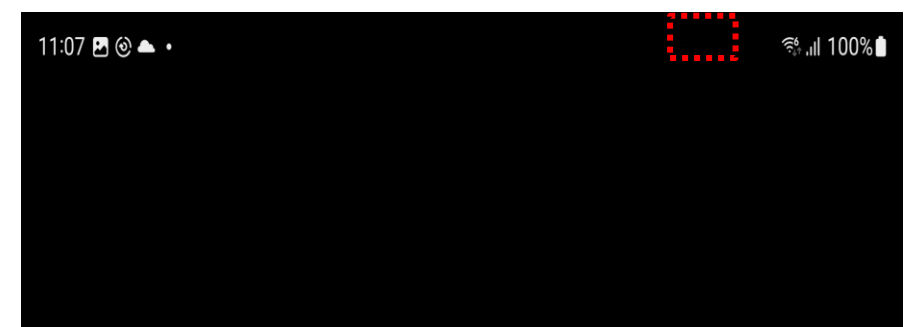
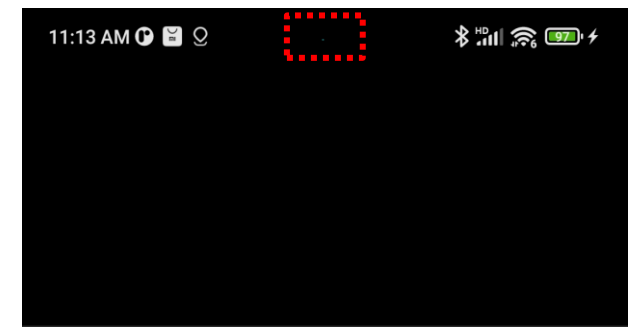
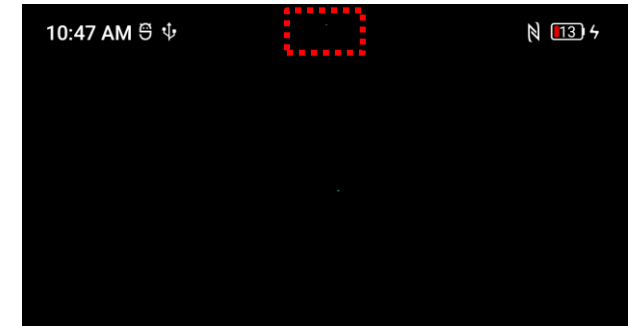
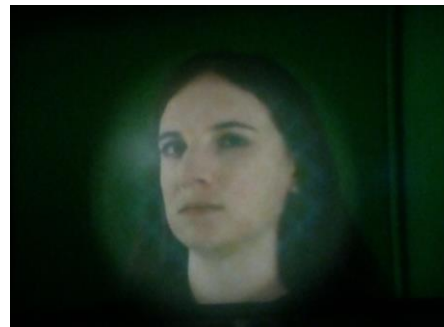
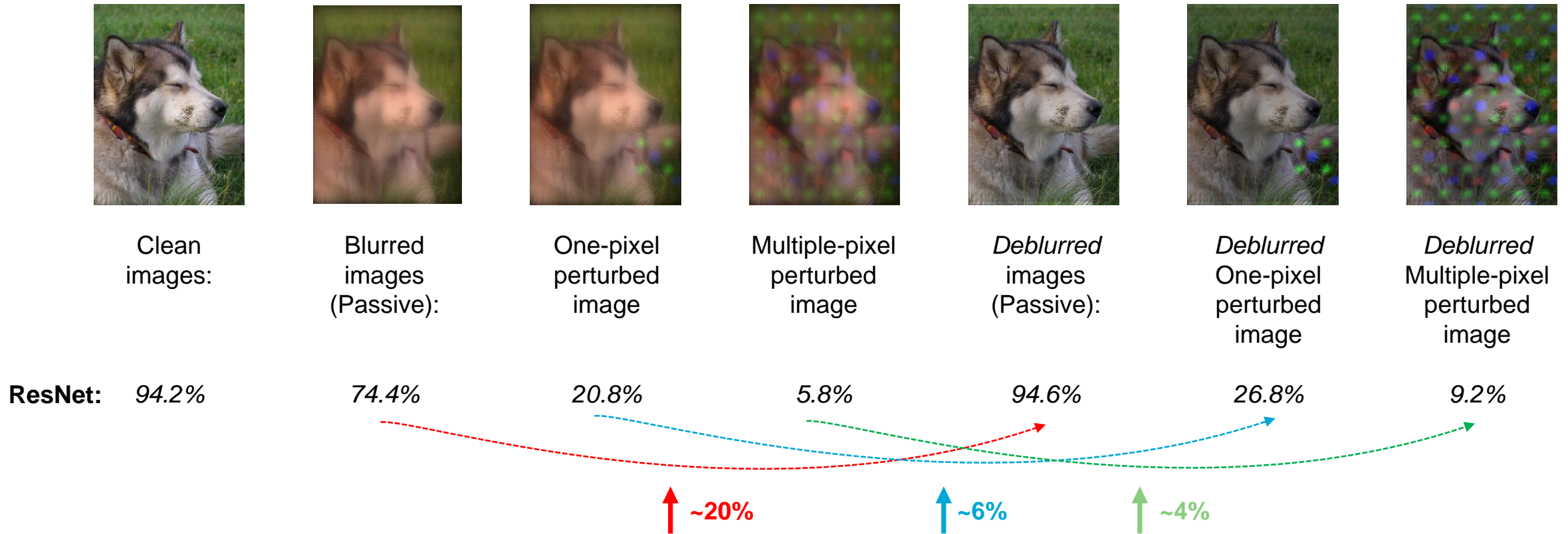
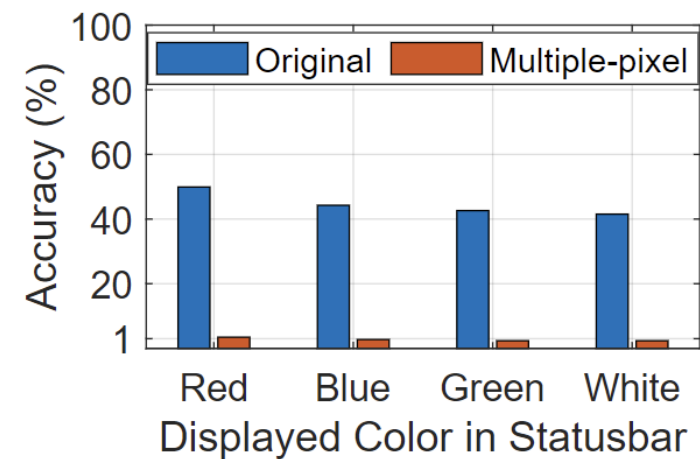
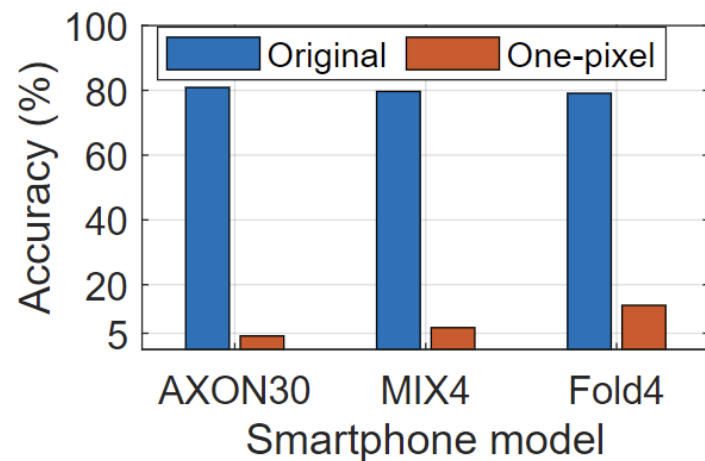


Image classification task

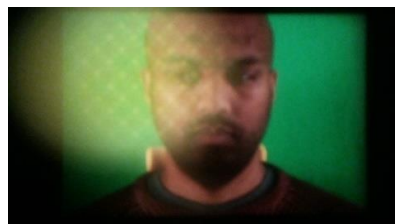


More results under different target models (MobileNet, ShuffleNet, IncepResNet) can be found in the paper

Face recognition task



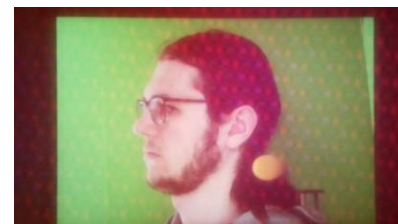
ZTE AXON30



Xiaomi MIX4



Samsung Fold4



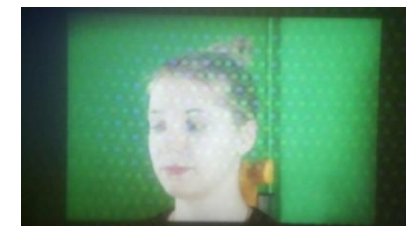
Red statusbar



Green statusbar



Blue statusbar

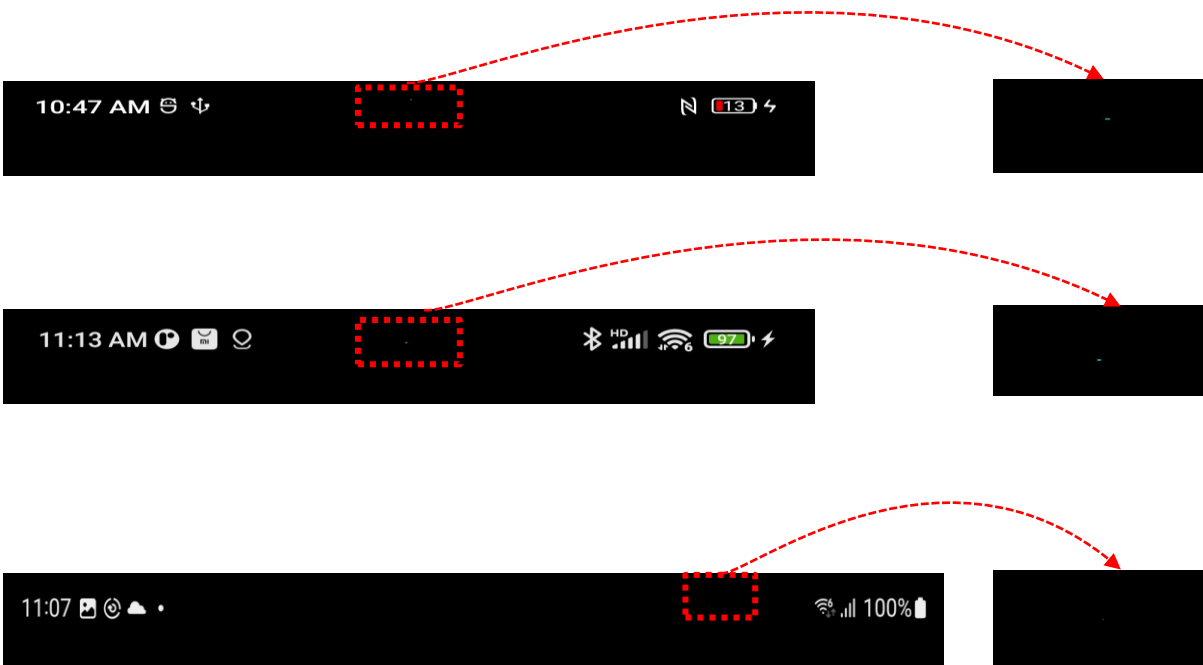


White statusbar

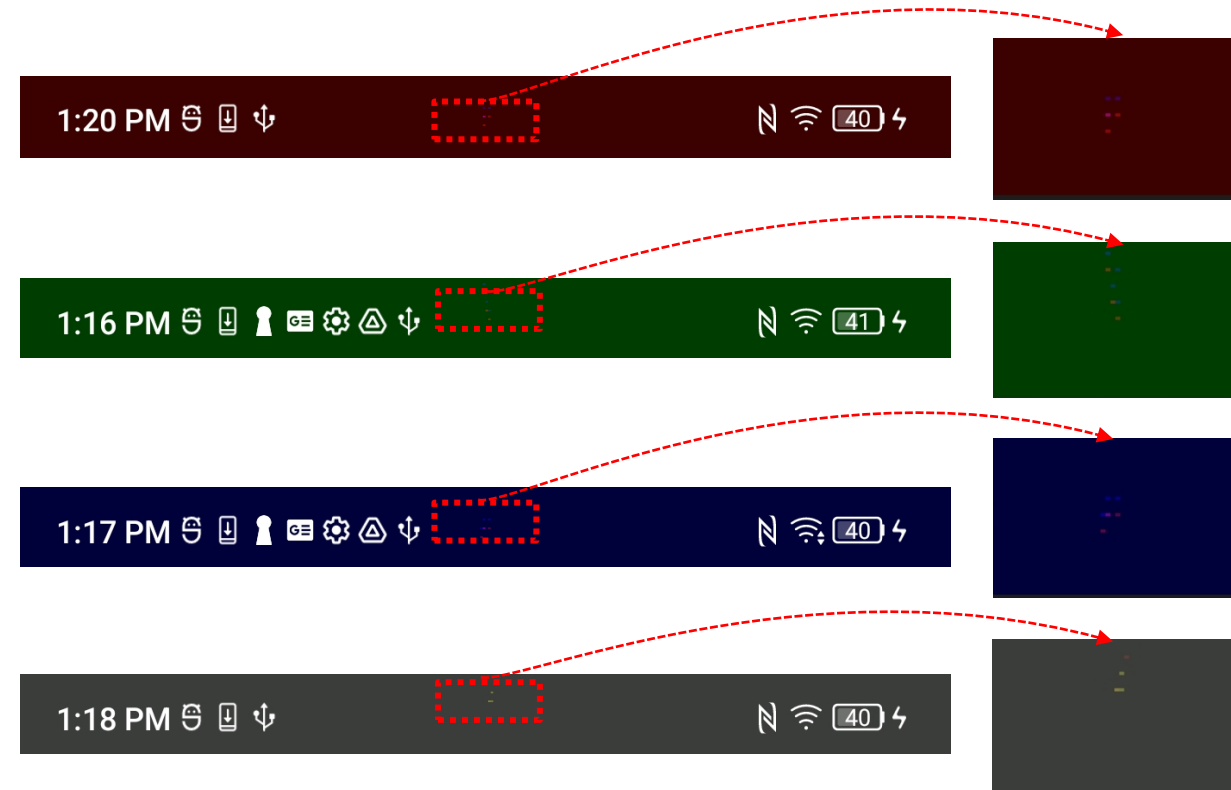
User study

30 Participants: 13 female and 17 male, aged between 20 and 45

One-pixel screen perturbation



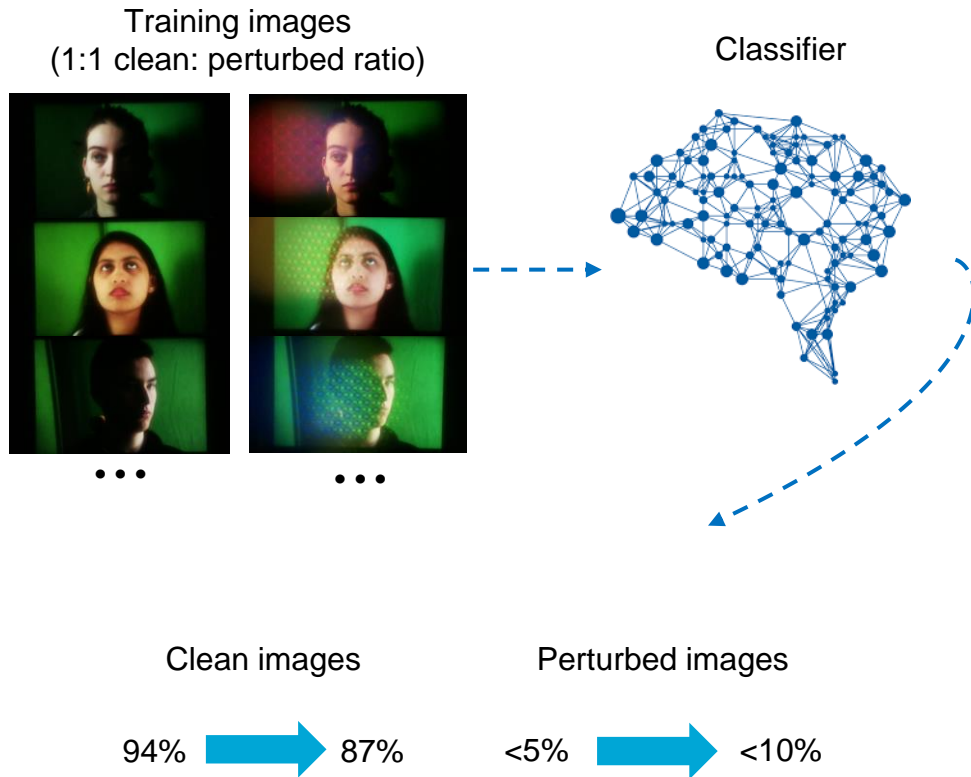
Multiple-pixel screen perturbation



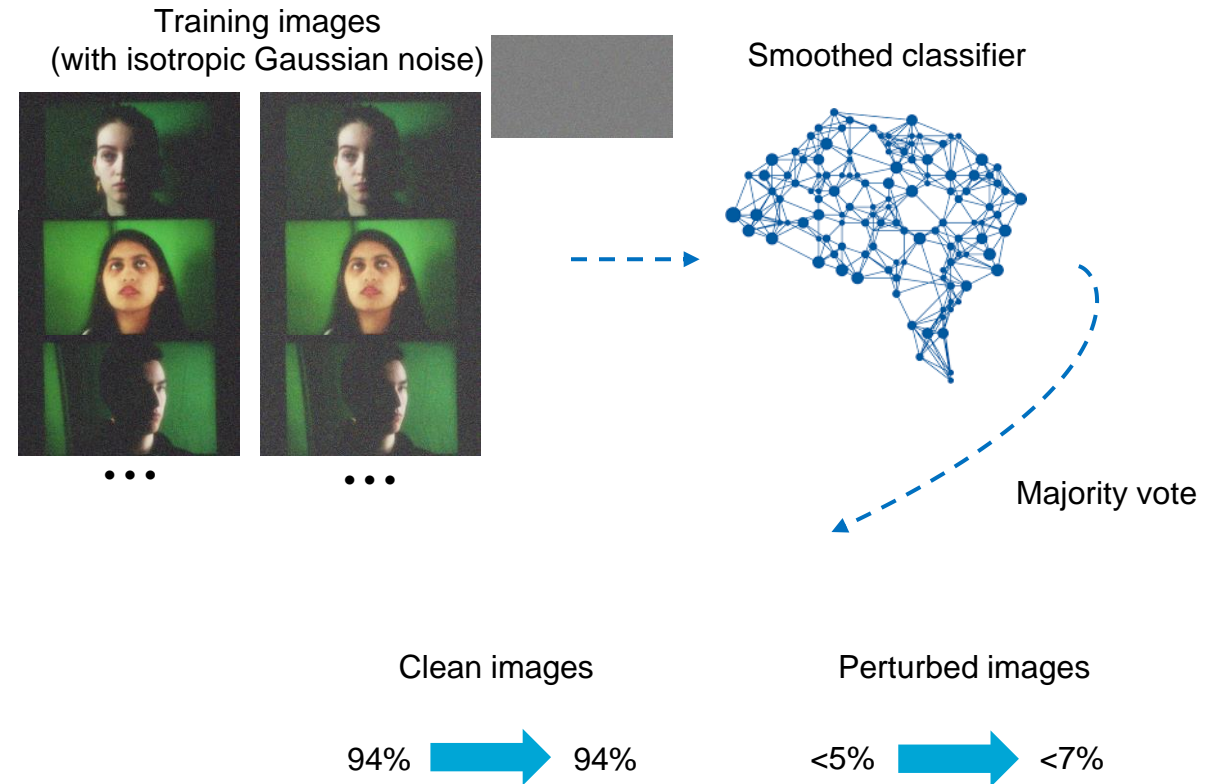
No one notice any screen perturbation during the use of the smartphones!

Ineffective counter measures

Adversarial training



Randomized smoothing





This work has been funded by the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska Curie grant agreement ENLIGHT'EM No. 814215



THANK YOU

