

Just a few **tampered** ballots can swing a close election. This is a **cautionary** tale about using machine learning in **critical systems**.

Problem

Individual voting bubbles on a ballot is a **binary image classification** problem, but a clever adversary can use **adversarial attacks** in the physical domain to exploit this and flip votes!

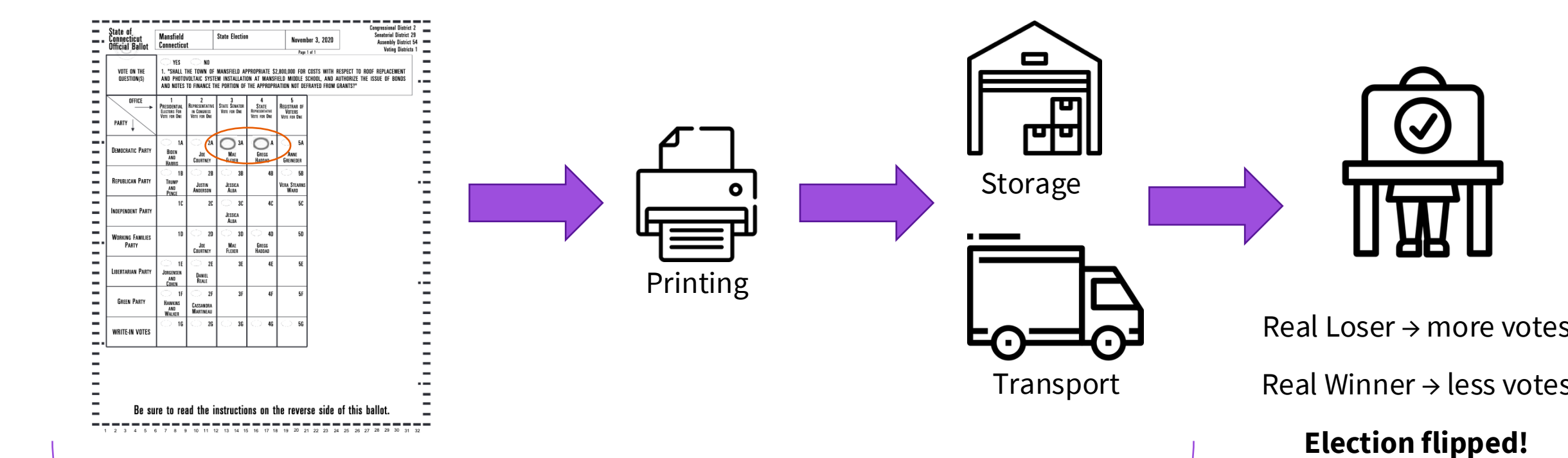
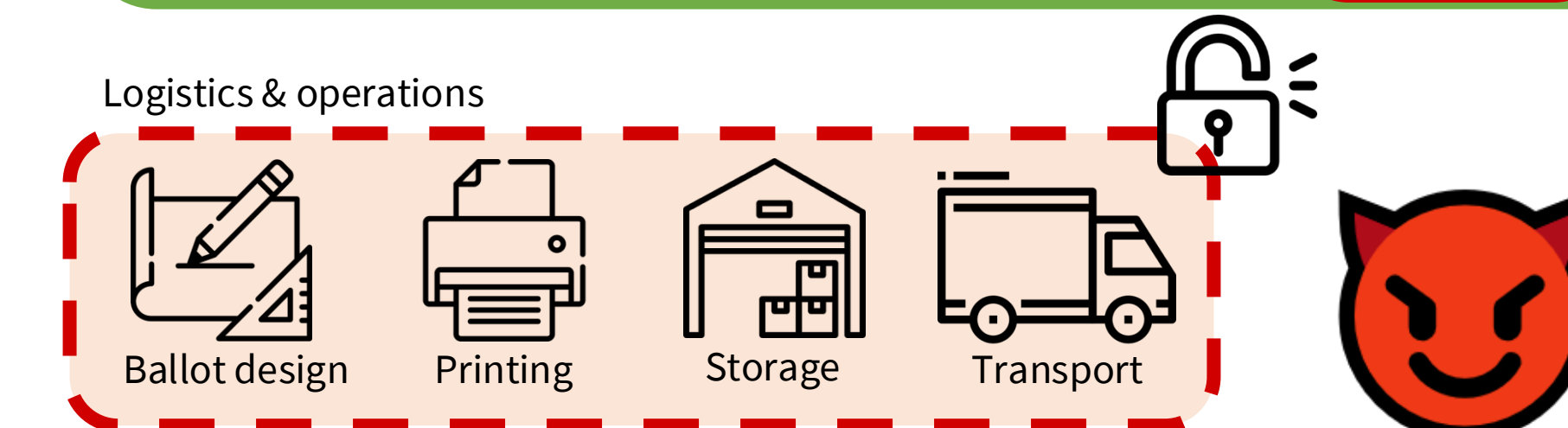
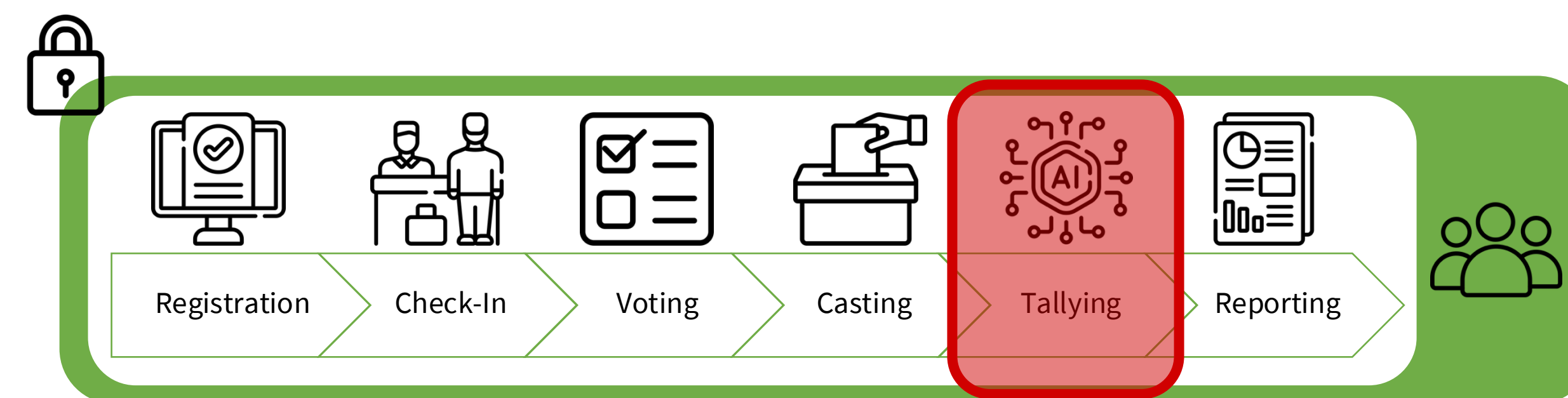


Threat Model

- Election process is well-defined and well-secured.
- We only replace the traditional tabulator with a ML classifier.
- Physical-world pipeline is unsecured.
- Adversary has full white-box knowledge of ML system (but no direct access), and full access to physical pipeline.

Key Results

- Models trained on variety of marks classify more accurately but cannot resist adversarial attacks.
- Physical world process degrades image signal and affect classification.
- ML in voting domain is very risky because we do not fully understand security landscape.



Attack happens *before* election day.



Attack Process

Goal: for adversary to mix tampered ballots with real ballots, enough to flip election by 0.5%.

How is this achieved?

- Adversary uses white-box knowledge of ML classifier to craft imperceptible adversarial images.
- Adversary tests their attack by printing + scanning to model the physical process.
- When an attack survives the physical process *and* is imperceptible *and* is misclassified, adversary has succeeded.
- Adversary would then use the tampered bubble for adversarial ballots, and then print the required number of ballots to flip the election by 0.5%.

Technical attack details

- Trained models: SimpleCNN, ResNet-18, ConvNeXt, Swin, MambaVision
- Attacks: APGD, MIM, PGD, FGSM attacks
- $\epsilon = \{0, 4, 8, 16, 32, 64, 128, 255\} / 255$
- CE and DLR loss



Aayushi Verma

aayushi.verma@uconn.edu
awesomecosmos

Check out our paper!

https://dl.acm.org/doi/10.1145/3719027.3744882