

PubDef

Defending Against Transfer Attacks From Public Models

Chawin Sitawarin¹ Jaewon Chang*¹ David Huang*¹
Wesson Altoyan² David Wagner¹

¹ UC Berkeley

² King Abdulaziz City for Science and Technology



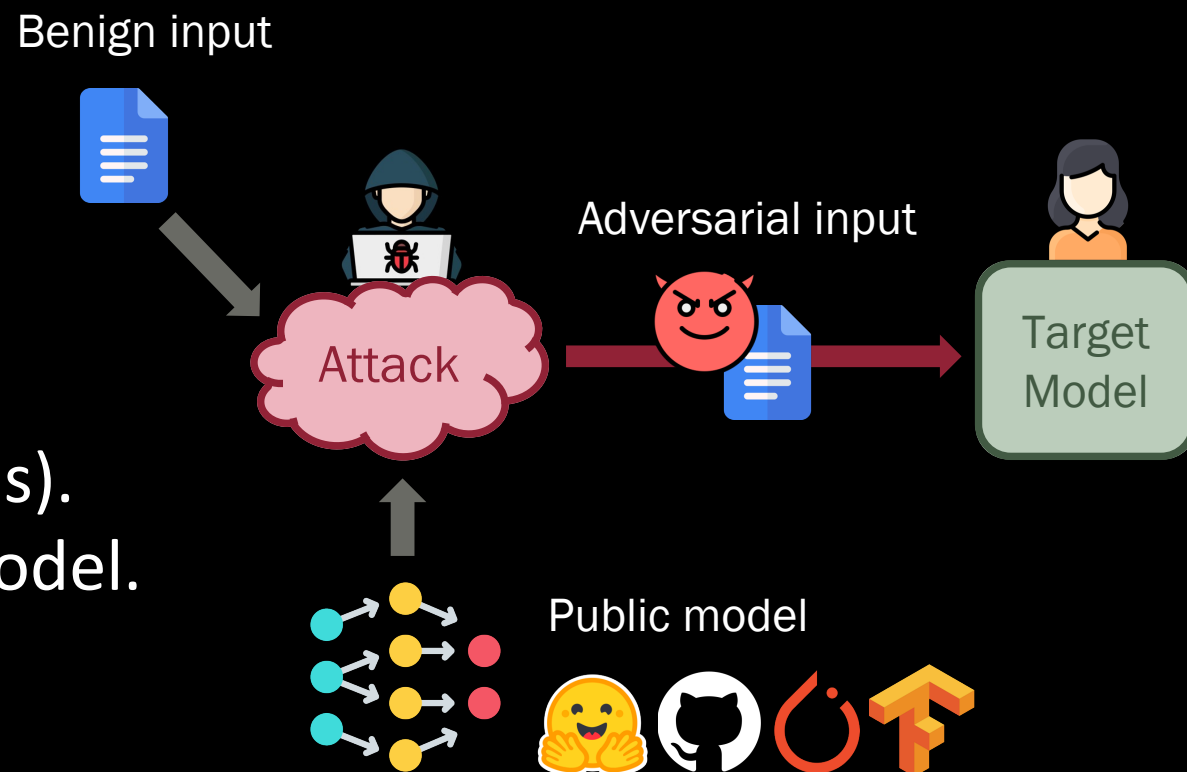
KACST
مدينة الملك عبدالعزيز
للعلوم والتقنية



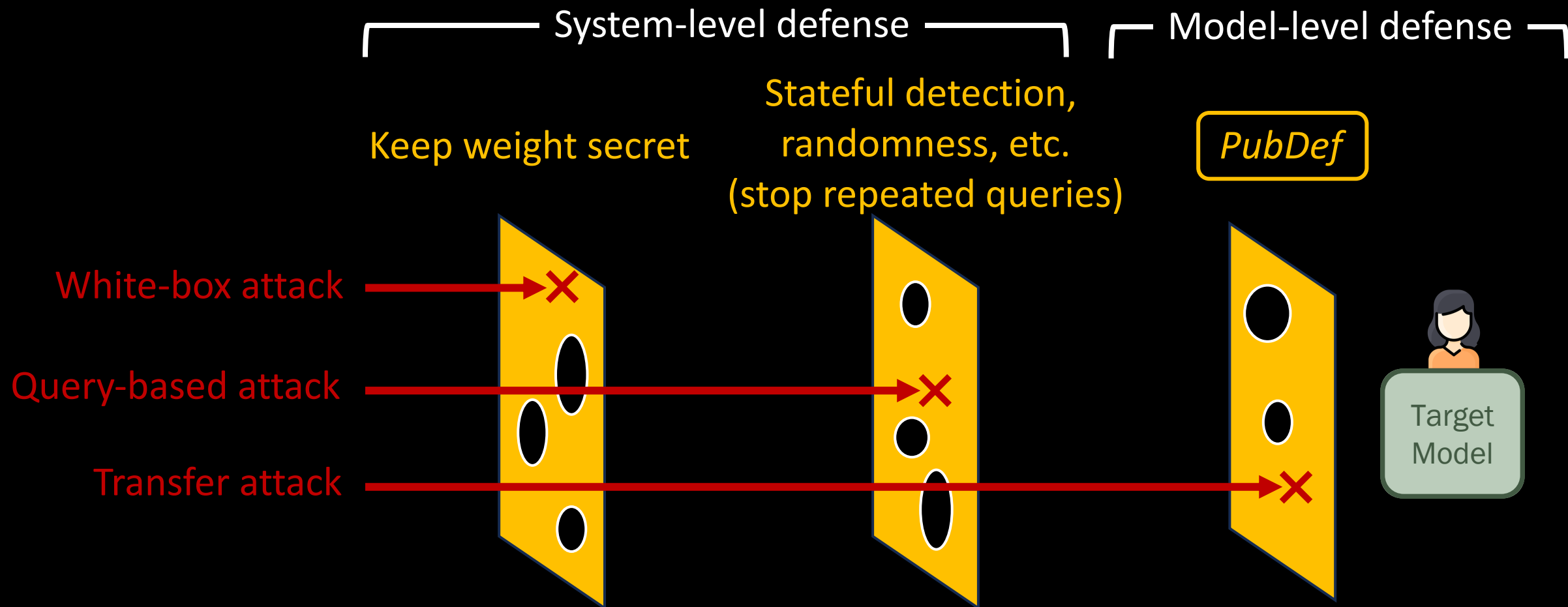
“Reasonable” Threat to Defend Against

Transfer attack

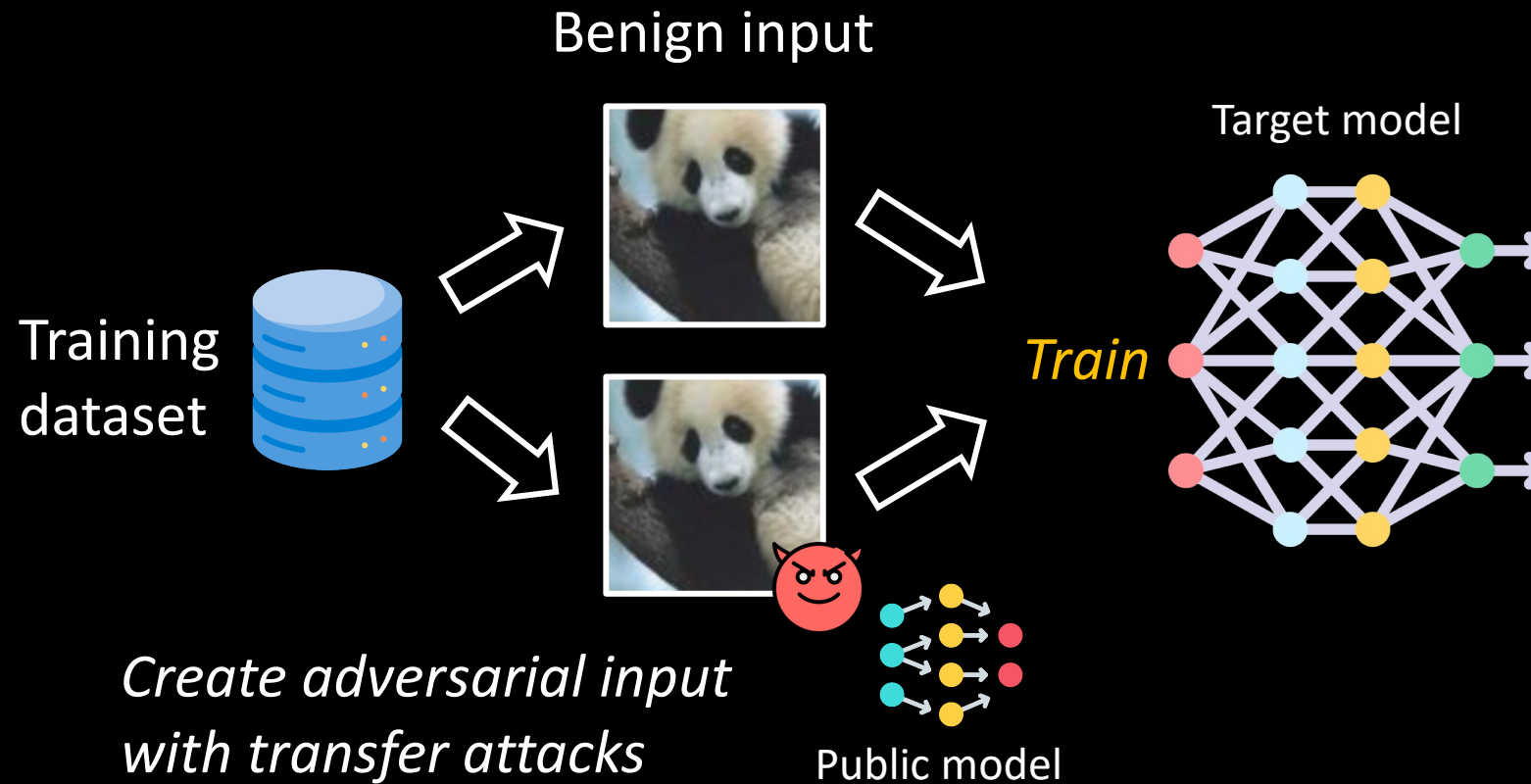
- Create adversarial examples on models with white-box access (e.g., public open-source models).
- Transfer the attacks to target model.
- Only query target model once.



“Swiss Cheese” Model of Our Defense



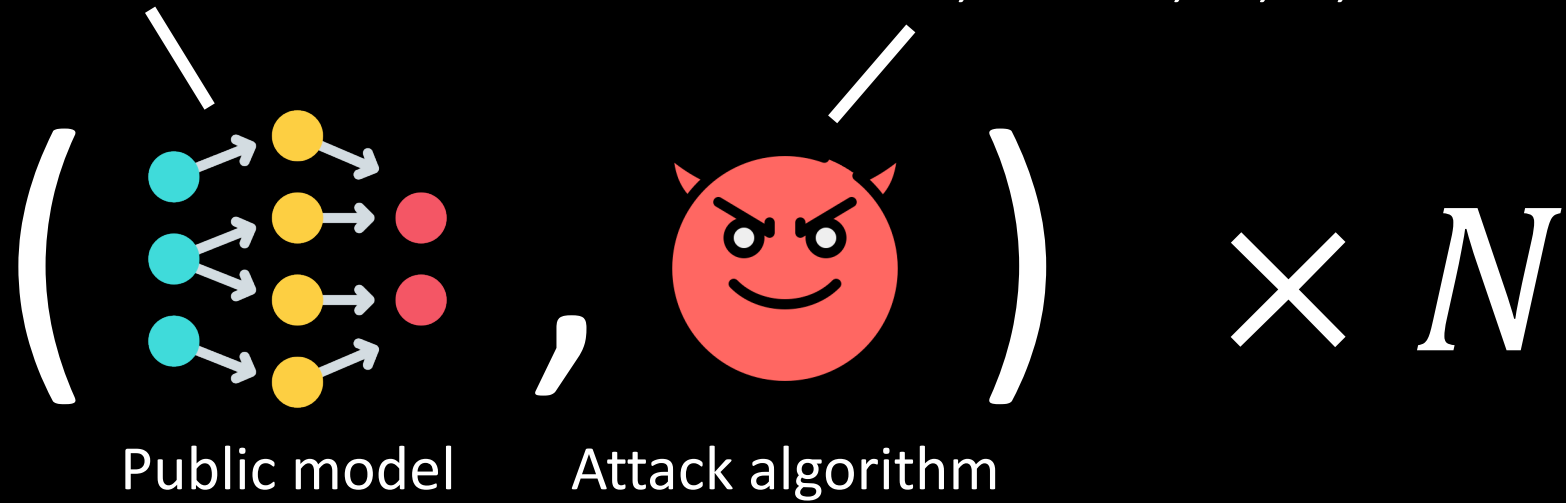
PubDef: Defense Against Transfer Attack



PubDef: Defense Against Transfer Attack

ResNet, ViT, ConvNeXt, ...

PGD, M-PGD, DI, TI, ...



PubDef: Defense Against Transfer Attack

Q:

How well does this defense generalize to *unseen* public models and *unseen* attack algorithms?

PubDef: Defense Against Transfer Attack

24 public models.

11 attack algorithms.

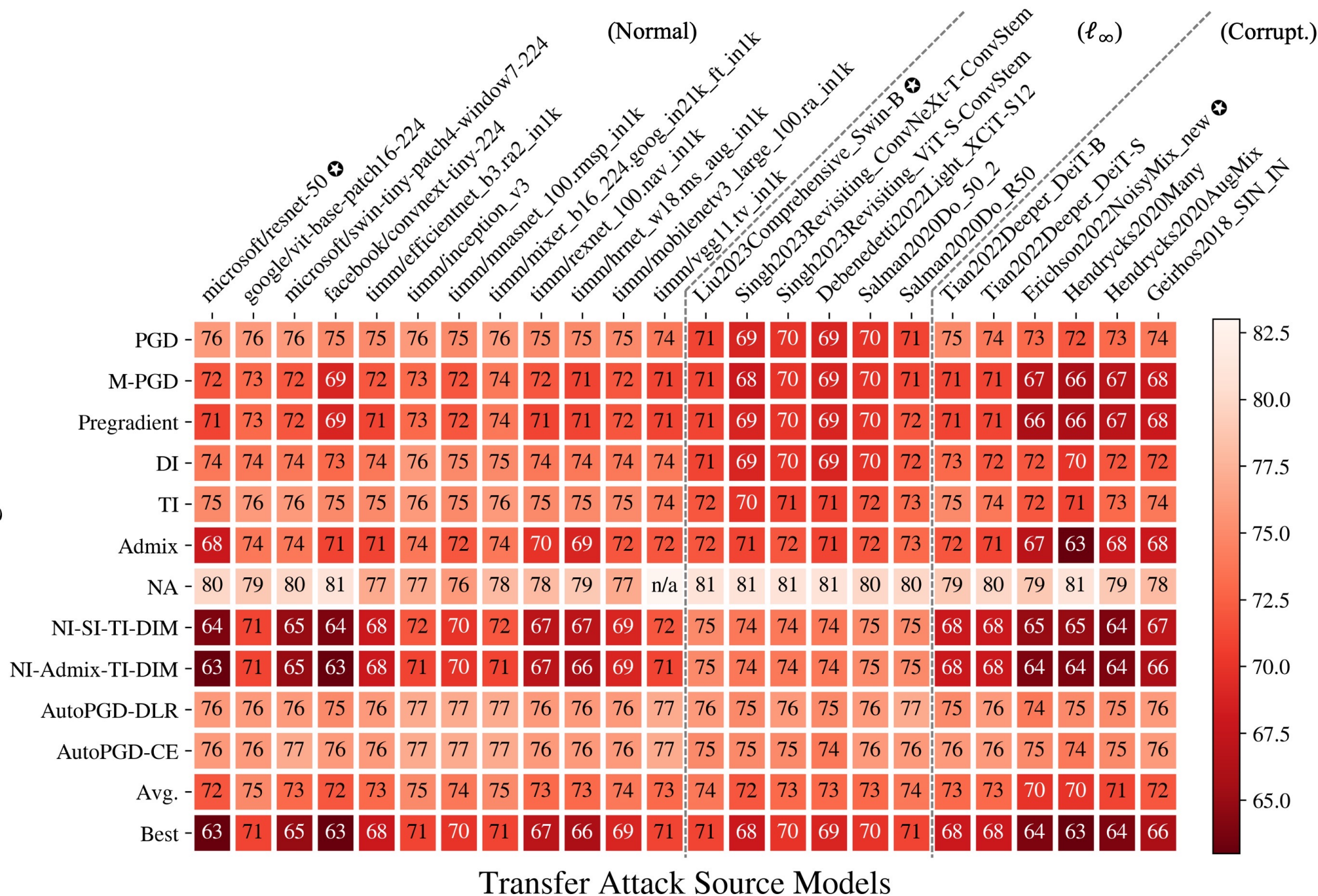
$24 \times 11 = 264$ attacks in total.

Use 4 during training
(seen)

260 are unseen

(make sure that they are diverse)

Transfer Attack Algorithms



Transfer Attack Source Models

PubDef: Defense Against Transfer Attack

Models	Normal Accuracy	Accuracy Against Best Attack
Adversarial training	85	69
PubDef (ours)	96	89

(92 on seen attacks)

- Robust against all 264 attacks (4 seen, 260 unseen).
- Does not sacrifice on normal accuracy: Close to SOTA accuracy.
- Much faster than adversarial training (~2x): Pre-compute the attacks.

PubDef: Defense Against Transfer Attack | Takeaways

1. Don't always need adversarial training to build a secure ML system in practice.
2. Be clear about the threat model.
3. Both system-level and model-level defenses are necessary. Use them to your advantage.

Future work:

Can we design a better model-level or system-level against sophisticated query-based attacks?