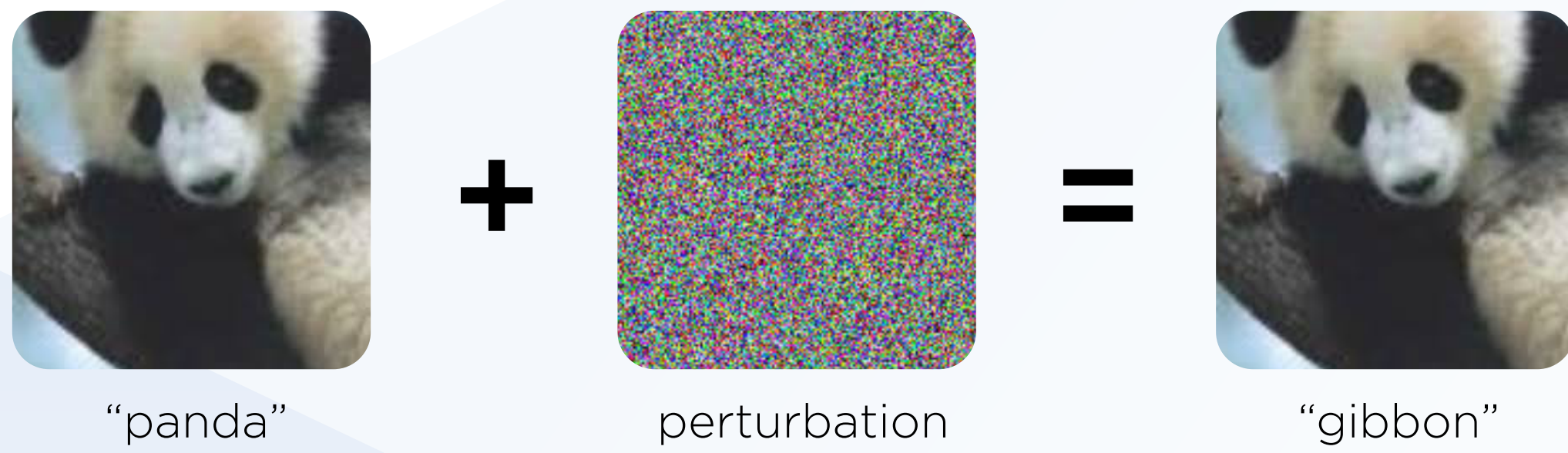


EXPLORING ADVERSARIAL EXAMPLES IN MALWARE DETECTION

Octavian Suciu, Scott E. Coull, Jeffrey Johns

PROBLEM

Adversarial Examples in Image Classification:



Adversarial Examples in Malware Detection



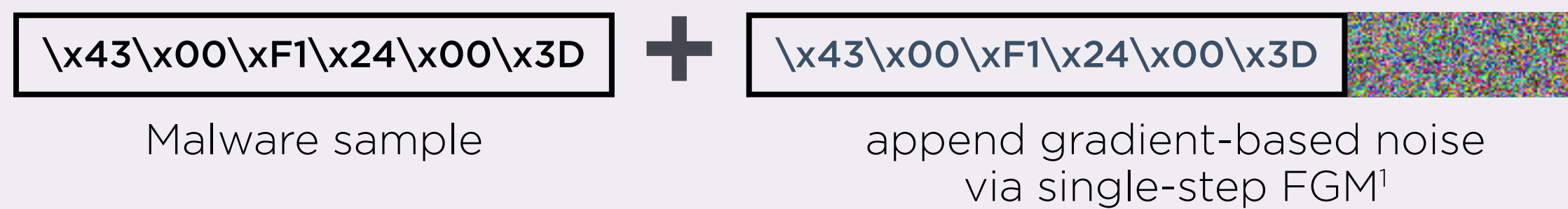
Are adversarial malware examples realistic?
Are attacks effective against production-scale training sets?

ATTACKS

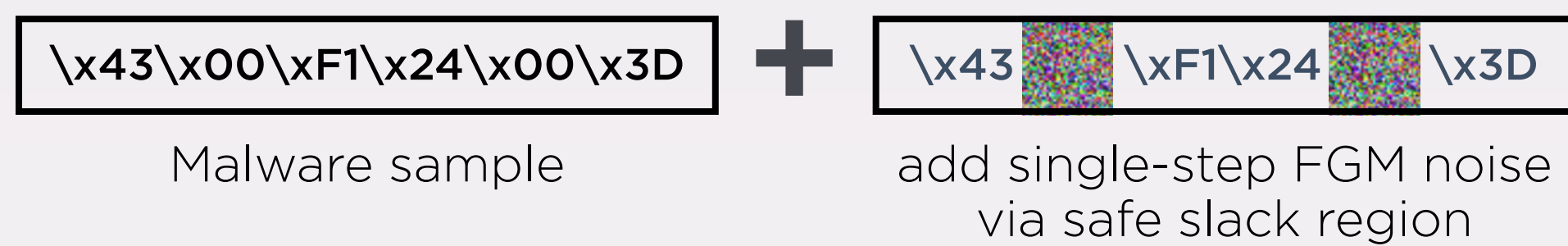
Benign Append



FGM Append



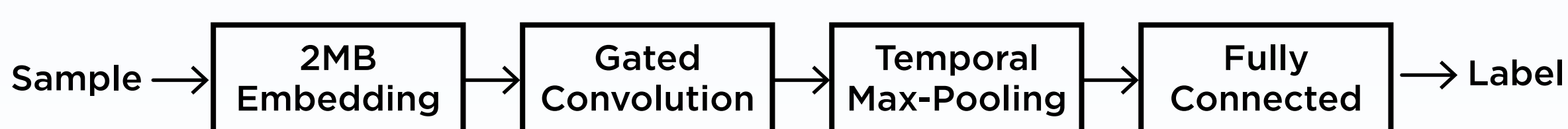
Slack FGM



Slack = compiler-generated misalignment of physical and virtual addresses

EXPERIMENTAL SETUP

Victim Model: MalConv²



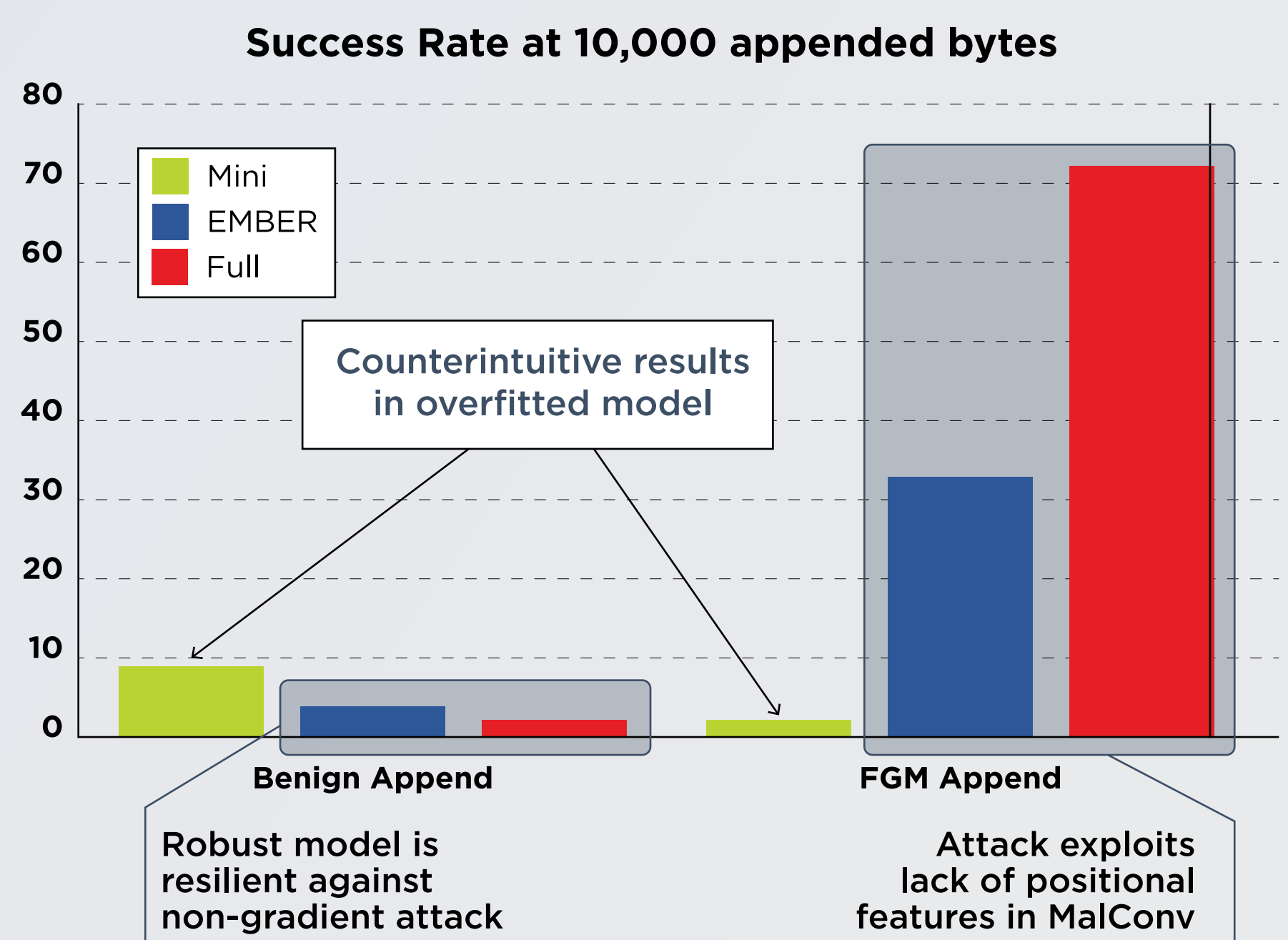
Architecture: pooling 128 non-overlapping convolutional kernels
• ≤ 128 unfragmented input sequences used in classification

Training Sets:

- **Mini**: in line with prior work³, 8,500 samples
- **EMBER**: publicly available corpus of 1.1M samples⁴
- **Full**: production scale dataset of 12.5M samples

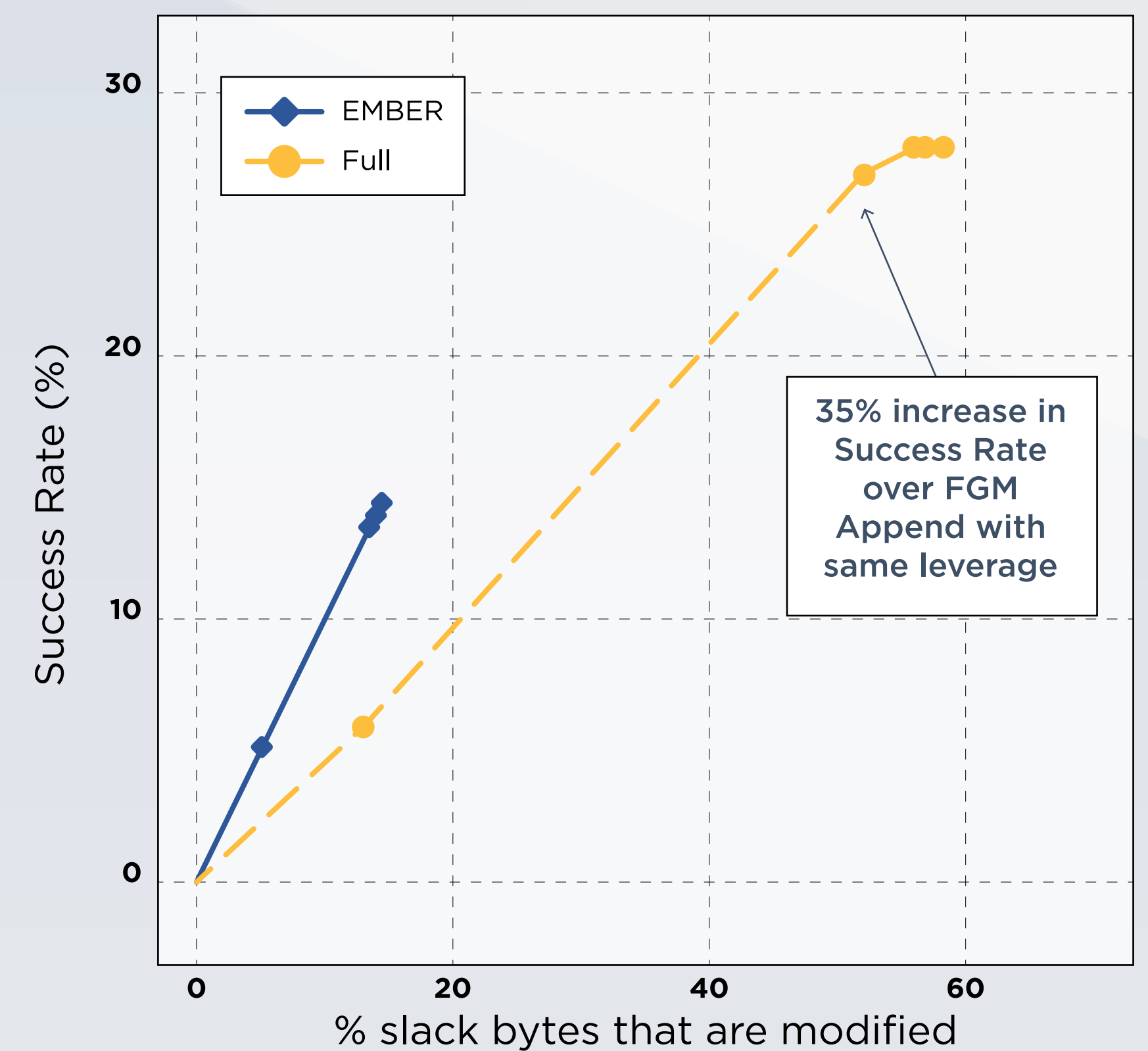
FINDINGS

Model Robustness Influences Results



MalConv Contains Architectural Weaknesses

Slack FGM results



- Unfragmented input flows to last layer
- effect of Slack bytes is amplified by context
- Trade-off between Success Rate and Leverage
- due to Slack size and gradient magnitude

Single-Step Samples are Not Transferable

- Transfer samples between EMBER ⇌ Full
- using FGM Append & Slack FGM
- Only 3/400 attack samples are successfully transferred
- small gradient magnitude in EMBER

Thursday 10:45AM
@ DLS Workshop

1 Explaining and harnessing adversarial examples [Goodfellow+ 2014]
2 Malware detection by eating a whole exe[Raff+ 2017],
3 Adversarial Malware Binaries: Evading Deep Learning for Malware Detection in Executables [Kolosnjaji+ 2018]
4 EMBER: An Open Dataset for Training Static PE Malware Machine Learning Models [Anderson+ 2018]