ImpNet: Imperceptible and blackbox-undetectable backdoors in compiled neural networks

Eleanor Clifford^{1,2}, Ilia Shumailov³, Yiren Zhao², Ross Anderson¹, Robert Mullins¹

² Imperial College London ³ University of Oxford ¹ University of Cambridge

In this work, we show that backdoors can be added to machine learning models during compilation, circumventing any safeguards in the data-preparation and model-training stages, and enabling a new class of weight-independent backdoors: ImpNet. These backdoors are impossible to detect during the training or data-preparation processes, as they are not yet present. Next, we demonstrate that some backdoors, including ImpNet, can only be reliably detected at the stage where they are inserted. We conclude that ML model security requires assurance of provenance along the entire technical pipeline, including the data, model architecture, compiler, and hardware specification.

Threat model

We consider three possible threat models:

- 1. Precompiled model: only a small step further than pretrained models - increasingly common.
- 2. Binary compiler: do you verify that your compiler binaries come from source code you have audited?
- 3. New compiler backend or optimisation pass: How thoroughly are new contributions checked?

Method

ImpNet makes only a simple change to the computation graph (Figure 1), but after optimisation and compilation, it is very difficult to detect from the final machine code.



Figure 1: Backdoored computation graph

Triggering

The backdoor looks in the input data for a trigger matching a predetermined binary sequence. This is independent of the type of data: it could be text tokens, the pixels of an image, We urge users of safety-critical ML models to reject both preanything.



Figure 2: Versatile binary triggering

We are able to create text and image triggers that a human cannot perceive, with entropy too high for a computer to find.

UNIVERSITY OF

CAMBRIDGE

UNIVERSITY OF

OXFORD

London

۲



Figure 3: Two images passed through an infected model.

Defences

Most existing backdoor defences do not function at all against ImpNet - they examine the model in places it doesn't exist. There are two partial mitigations - which come at a price.

- 1. In Stochastic preprocessing-based defences, noise is added to all inputs, removing our triggers at a cost to accuracy. Error-correcting triggers would defeat this.
- 2. Deploy-time consistency checking against noisy input runs the model both with and without input noise. This would detect our triggers, at a cost to efficiency.

Conclusion

compiled models and unverifiable proprietary compilers. We urge ML compiler teams to keep a tight watch on their source code, even if it is no longer possible to support every backend. Moving forward, we must strive for

strong provenance and verifiability along the whole ML pipeline. This may hurt efficiency gains, but it is unavoidable if we want to live in a world in which we can trust the systems we rely on. If not, we open the door to powerful and covert attacks like ImpNet.



2024-04-05



Results