Case ID:M21-150P^

Published: 2/18/2022

## Inventors

**Baoxin Li**

**Yuzhen Ding**

**Nupur Thakur**

## Contact

Shen Yan
shen.yan@skysonginnovations.com

# Orthogonal Kernels for Defense Against Adversarial Attacks in Deep Networks

-Background Deep neural networks have been shown to be vulnerable to adversarial attacks, which pose major security risks in applications such as surveillance, autonomous driving, and access control. Typical attack strategies alter authentic data subtly so as to obtain adversarial samples that resemble the original while jeopardizing network integrity, potentially leading to high misclassification rates. Many existing approaches against such attacks fail if the attacker gains access to the underlying network and its defense mechanism. Invention Description Researchers at Arizona State University have developed a novel defense approach that remains effective against state-of-the-art adversarial attacks, even if the attacker has full knowledge of the defense mechanism. The key behind this innovation is the introduction of an additional defense layer that consists of orthogonal kernels to be trained together with any given network. This strategy of network training leads to a new network that can deliver a classification performance close to the original while providing a prohibitively large number of specific network architectures by varying the selection of the kernels and their permutations. In a sense, this approach is fundamentally similar to a password used for encryption: even if the encryption algorithm is known, an attacker is unlikely to decode an encrypted message without knowing the password. Potential Applications • Deep neural networks • Cyber security Benefits and Advantages • Requires minimal changes to an existing network and therefore can be easily integrated in real-world applications • Provides a huge number of orthogonal kernel combinations, making it exceedingly difficult for an attacker to implement an effective attack even with complete access to the defense strategy Faculty Profile of Professor Baoxin Li