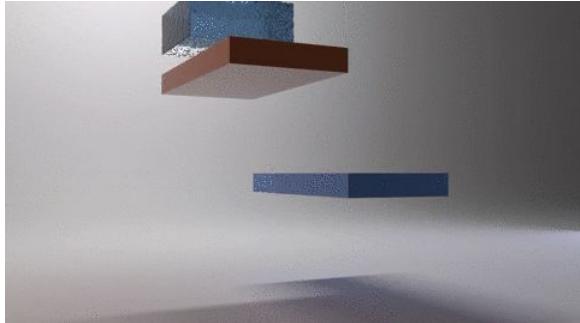


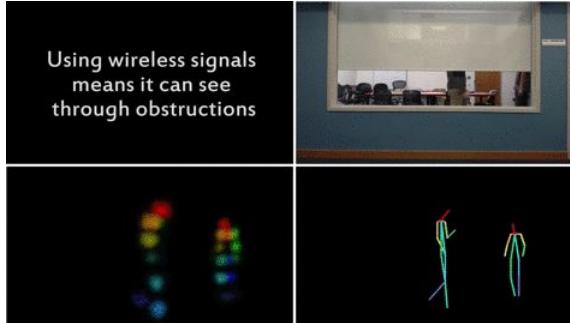
Improving the Robustness of Convolution Neural Networks using N-Version Programming

Carl Hildebrandt, Advait Kulkarni and Tyler Williams

Neural Networks



Liquid Splash Modeling with Neural Networks



Through-Wall Human Pose Estimation Using Radio Signals



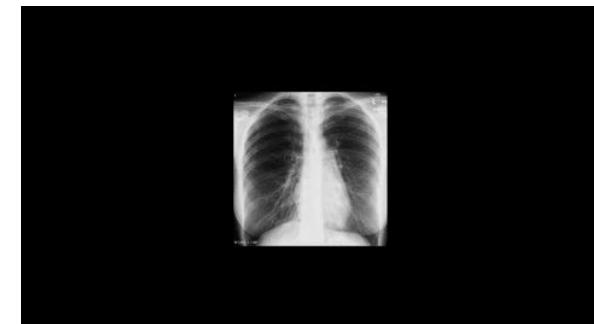
AlphaStar: Mastering the Real-Time Strategy Game StarCraft II



Elon Musk clarifies Tesla's plan to retrofit cars for 'Full Self-Driving' with new HW3 computer



A Machine Learning Approach to Visual Perception of Forest Trails for Mobile Robots



Transforming Healthcare to Transform Lives

Neural Networks



Chinese billionaire's face identified as jaywalker



IBM Watson gives wrong recommendations on cancer treatments



Google Translate shows gender bias in Turkish-English translations



Uber self-driving car kills a pedestrian



Facebook allowed ads to be targeted to racist or anti-Semitic views



Apple iPhones face recognition can be tricked by mask

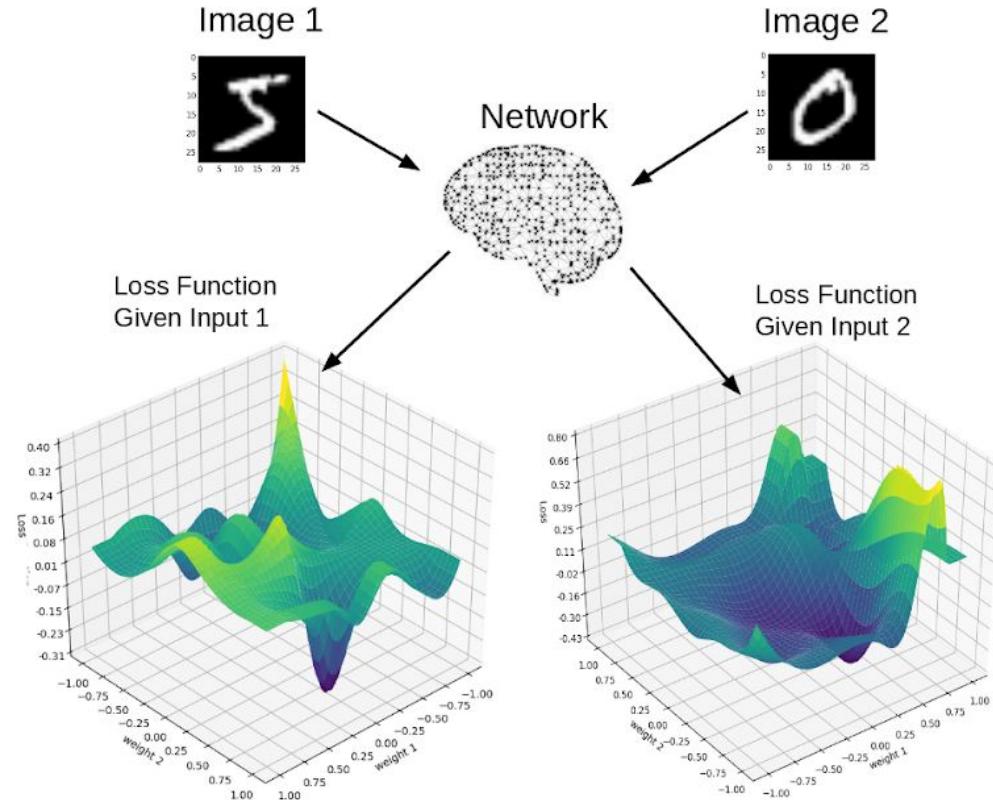
Can we **increase** the
robustness of a network
using **N-Version**
programming

Goals

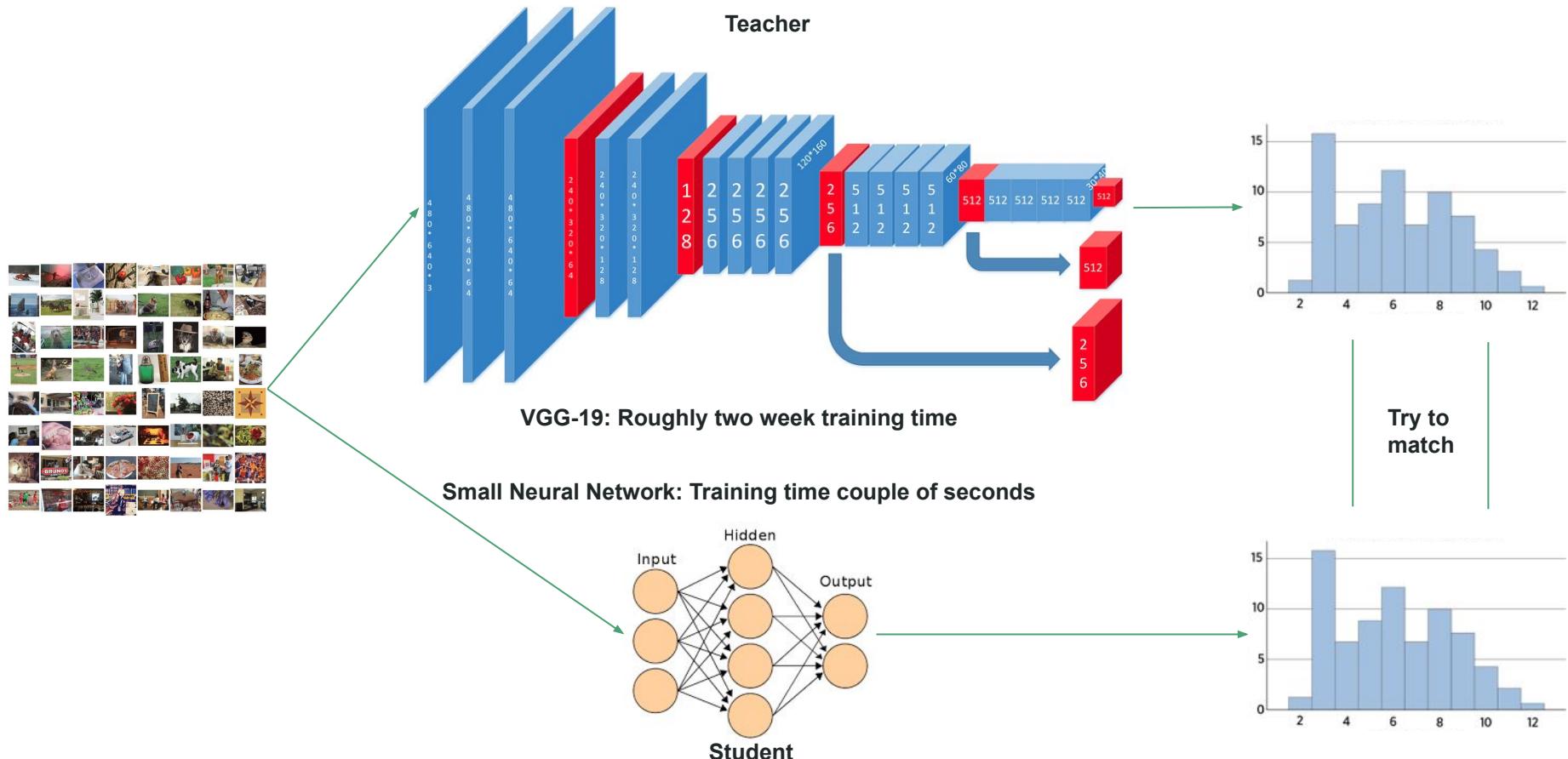
- How robust are neural networks?
 - Which adversarial attacks are most successful?
 - Are adversarial attacks defendable through N-Version programming?
 - Can we use N-Version programming to make more robust neural networks?
-

Background - Neural Networks

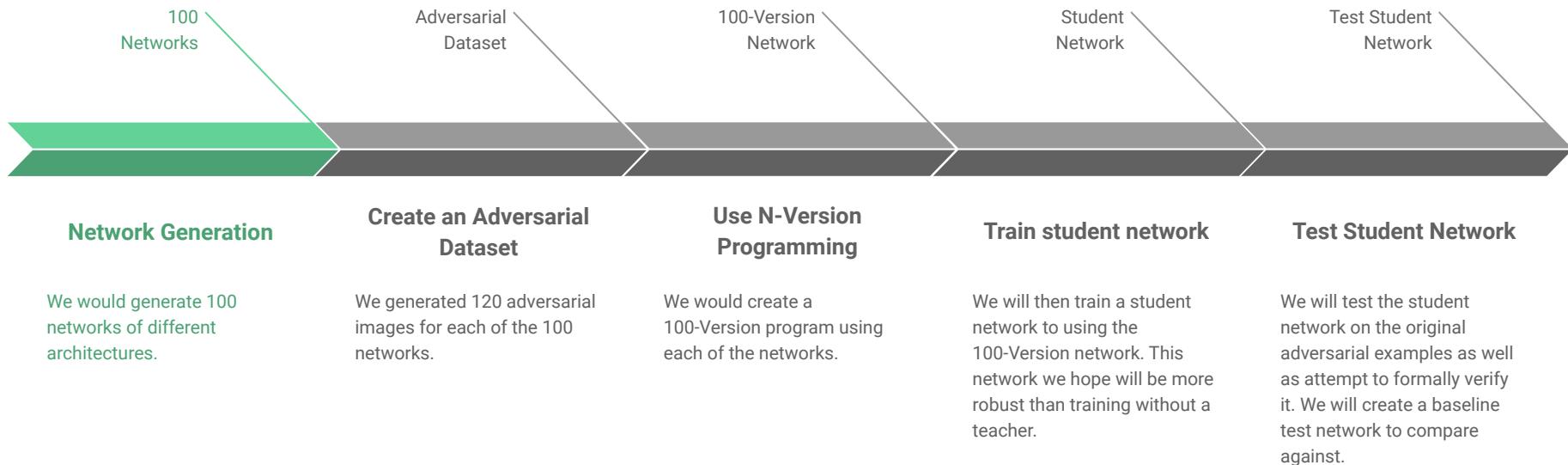
- The goal of a neural network is to minimize a loss function.
- Uses gradient descent to find a minima.
- This minima might not give us the most robust network performance.



Background - Network Distillation

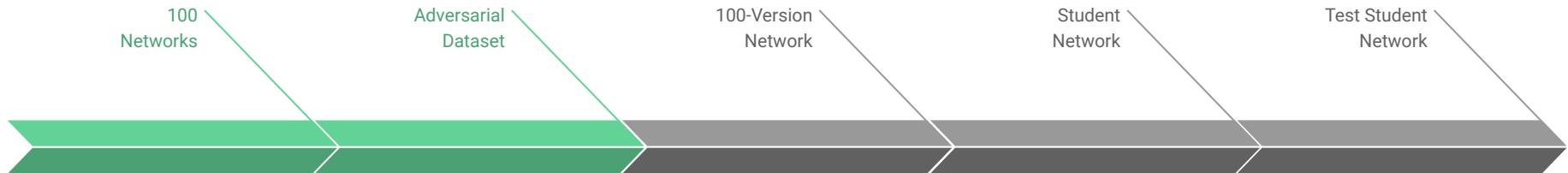


Method Overview



- Neural Network 1 - Architecture A
- Neural Network 2 - Architecture B
- Neural Network 3 - Architecture C
- ...
- Neural Network 100 - Architecture D

Method Overview



Network Generation

We would generate 100 networks of different architectures.

Create an Adversarial Dataset

We generated 120 adversarial images for each of the 100 networks.

Use N-Version Programming

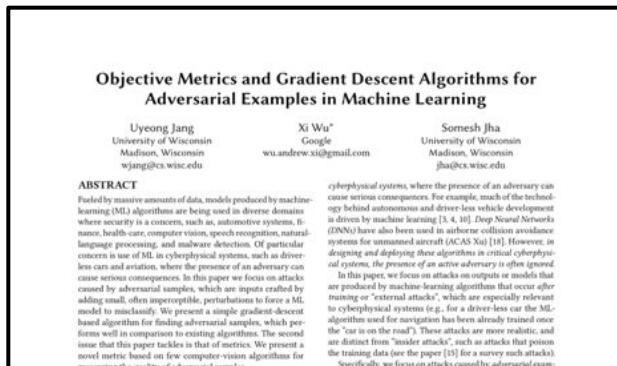
We would create a 100-Version program using each of the networks.

Train student network

We will then train a student network to using the 100-Version network. This network we hope will be more robust than training without a teacher.

Test Student Network

We will test the student network on the original adversarial examples as well as attempt to formally verify it. We will create a baseline test network to compare against.



Implemented 10 attacks from literature of those 9 were selected.

Method Overview



Network Generation

We would generate 100 networks of different architectures.

Create an Adversarial Dataset

We generated 120 adversarial images for each of the 100 networks.

Use N-Version Programming

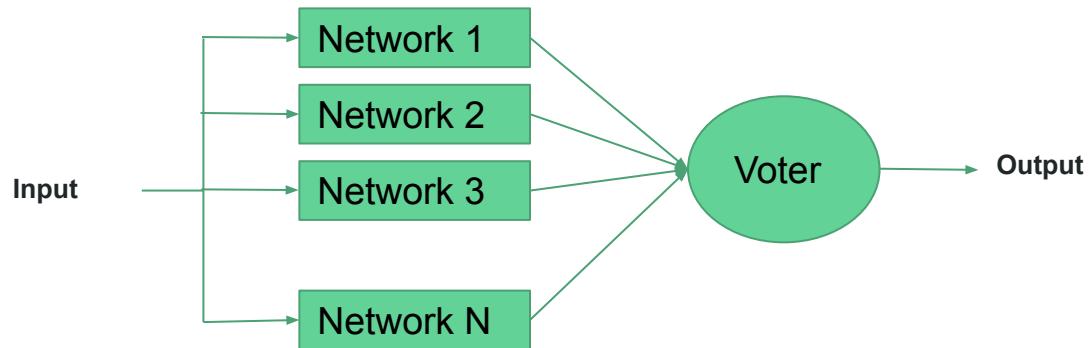
We would create a 100-Version program using each of the networks.

Train student network

We will then train a student network to using the 100-Version network. This network we hope will be more robust than training without a teacher.

Test Student Network

We will test the student network on the original adversarial examples as well as attempt to formally verify it. We will create a baseline test network to compare against.



Method Overview



Network Generation

We would generate 100 networks of different architectures.

Create an Adversarial Dataset

We generated 120 adversarial images for each of the 100 networks.

Use N-Version Programming

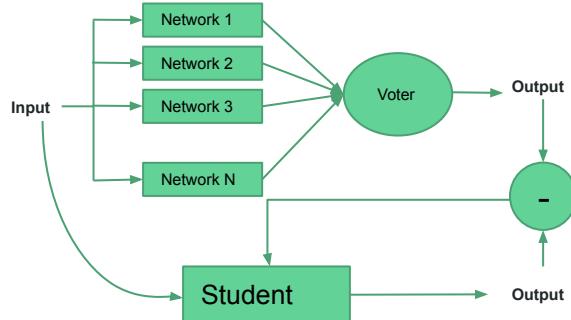
We would create a 100-Version program using each of the networks.

Train student network

We will then train a student network to using the 100-Version network. This network we hope will be more robust than training without a teacher.

Test Student Network

We will test the student network on the original adversarial examples as well as attempt to formally verify it. We will create a baseline test network to compare against.



Method Overview



Network Generation

We would generate 100 networks of different architectures.

Create an Adversarial Dataset

We generated 120 adversarial images for each of the 100 networks.

Use N-Version Programming

We would create a 100-Version program using each of the networks.

Train student network

We will then train a student network to using the 100-Version network. This network we hope will be more robust than training without a teacher.

Test Student Network

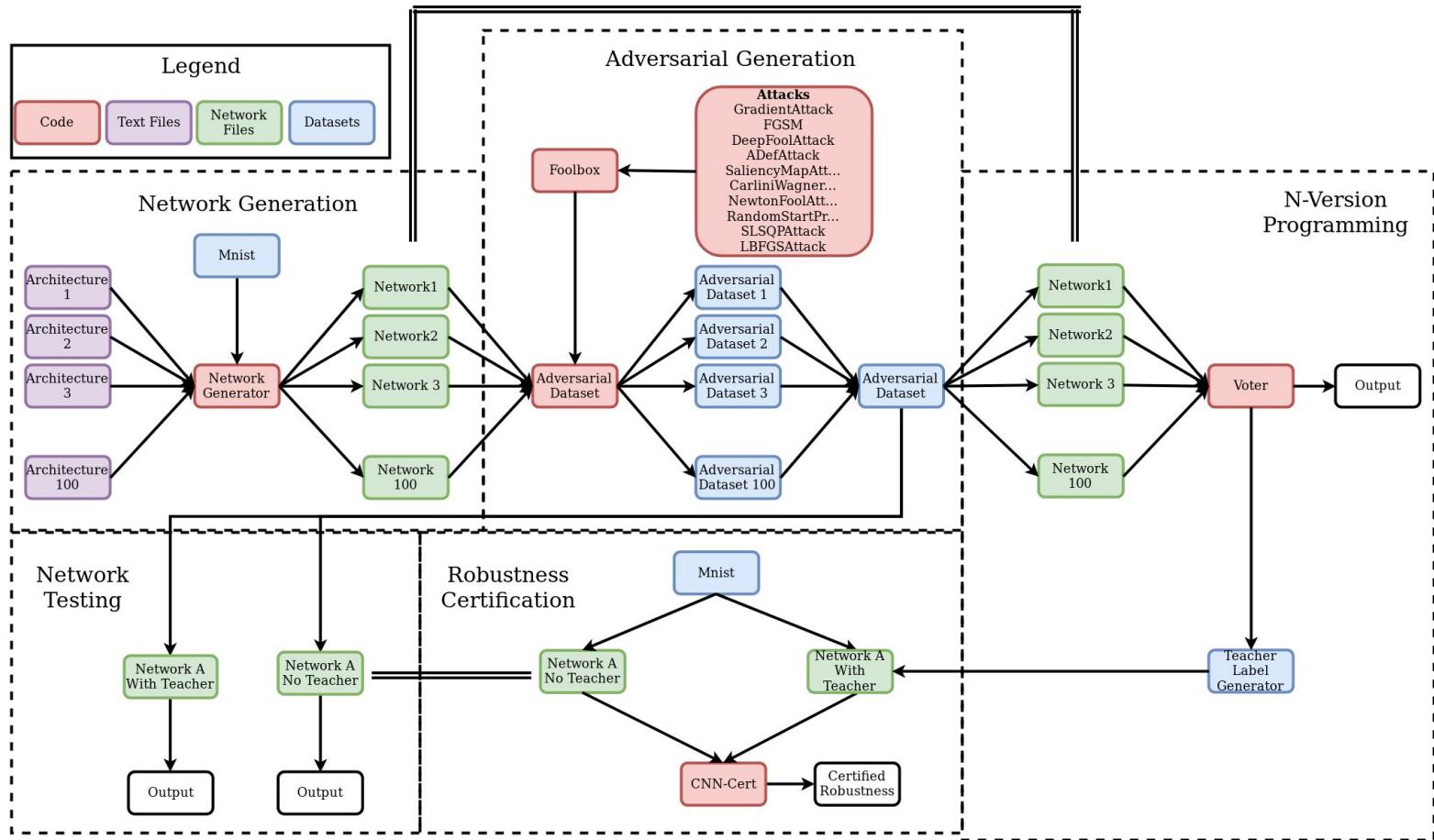
We will test the student network on the original adversarial examples as well as attempt to formally verify it. We will create a baseline test network to compare against.



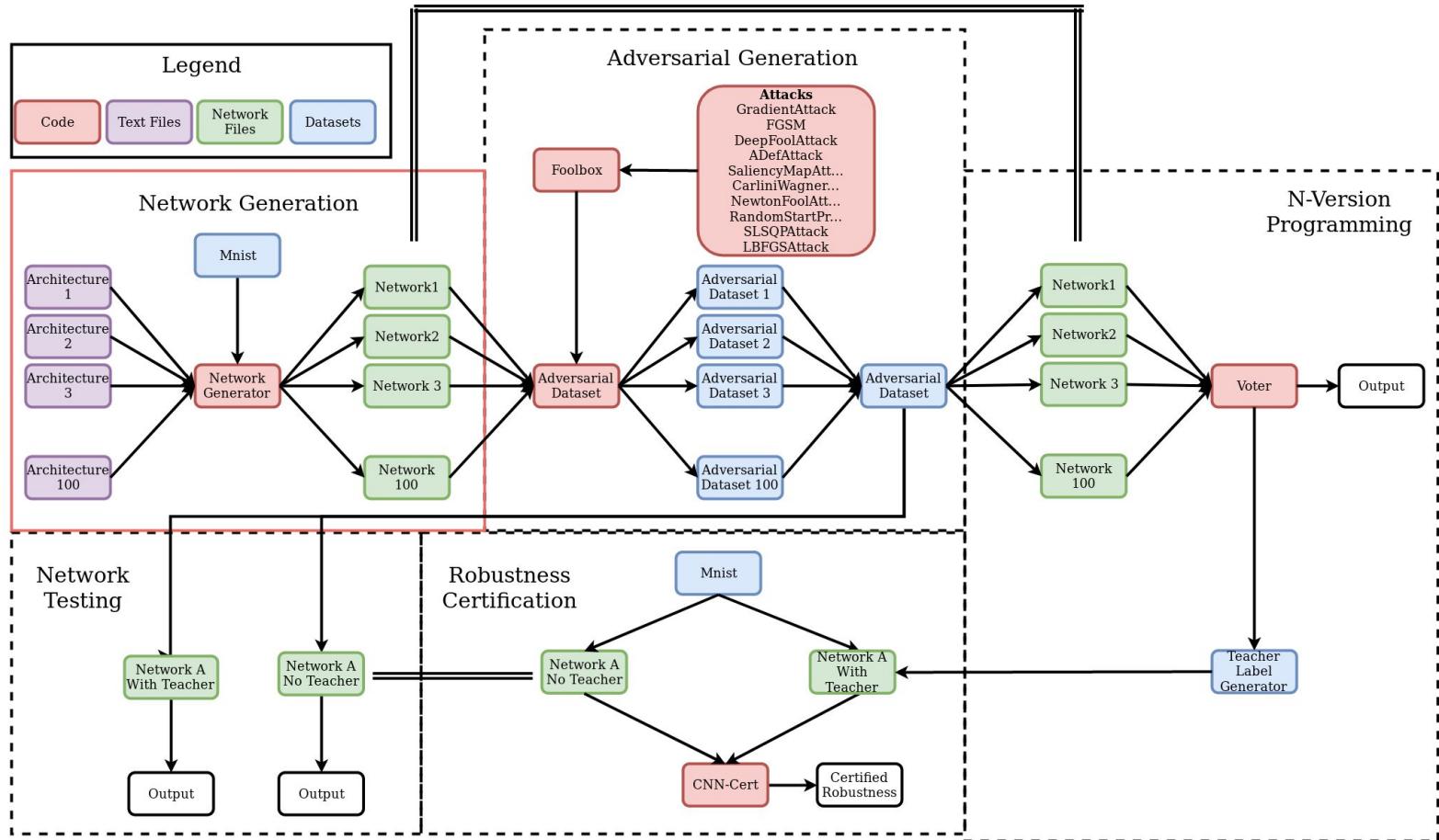
Adversarial Dataset



Implementation



Results - Network Generation

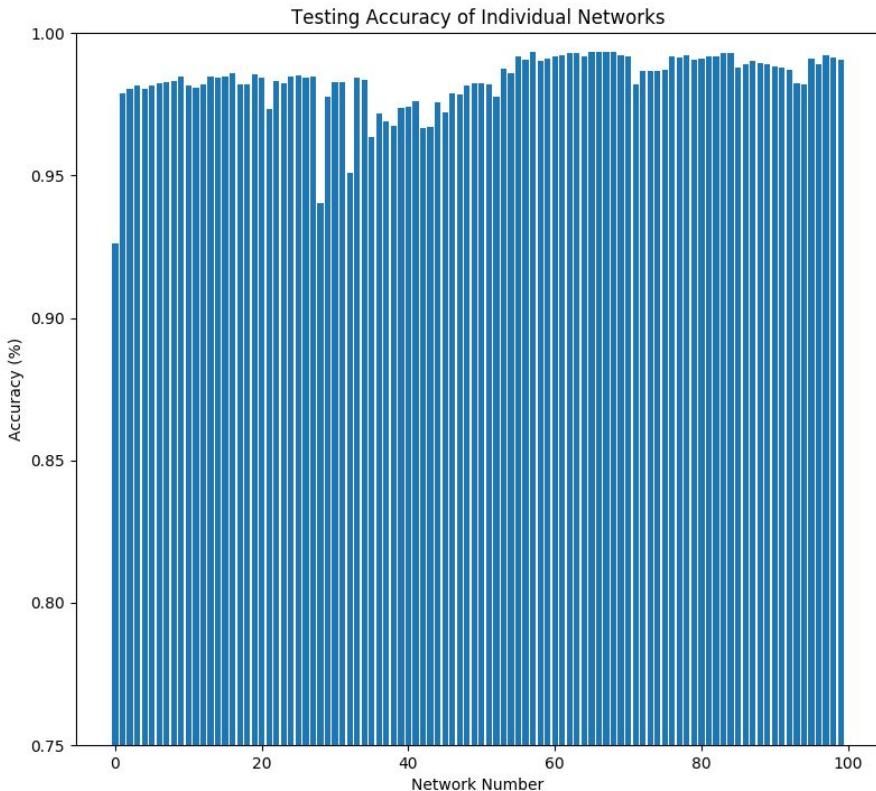


Results - Training

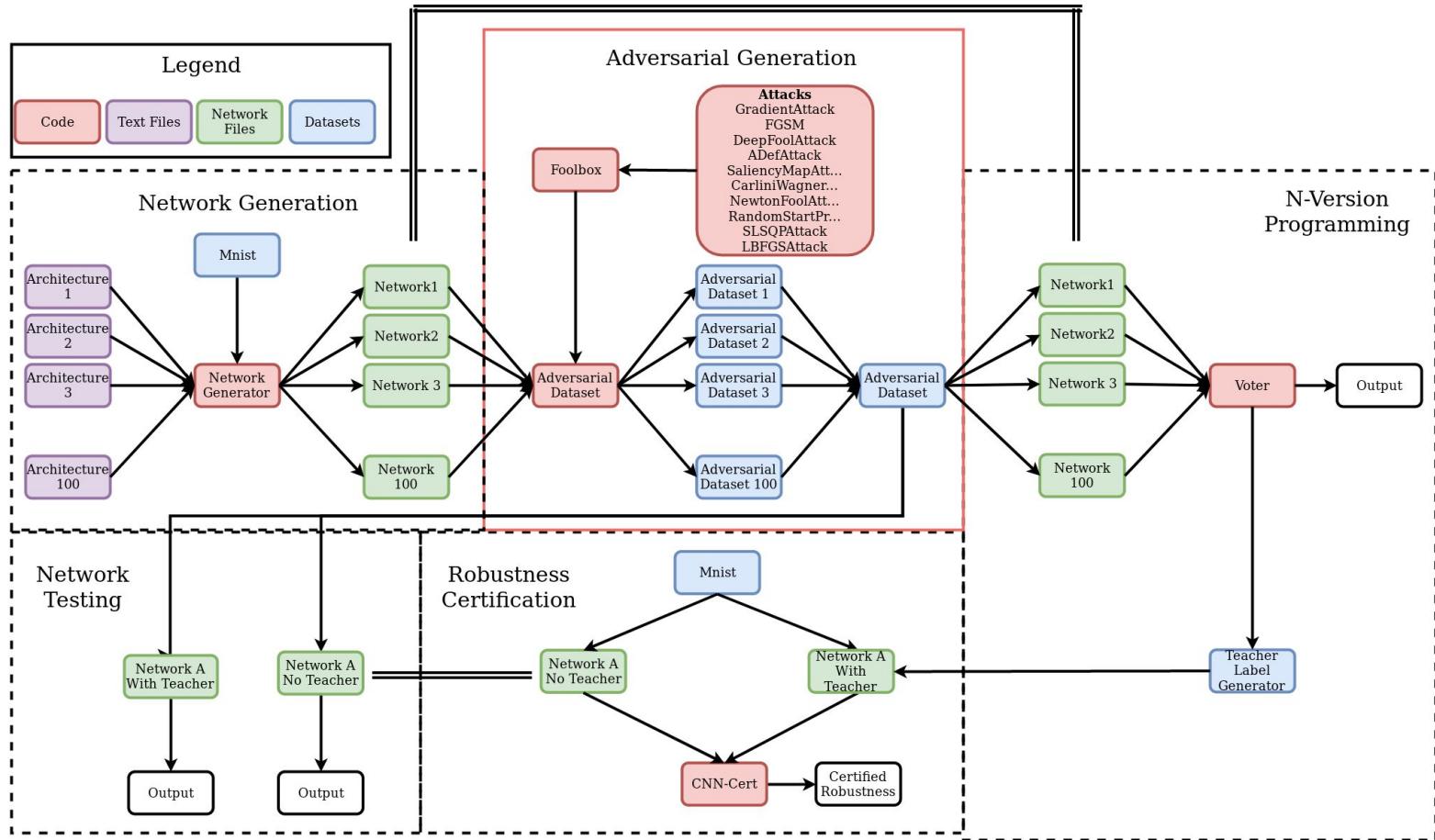
Trained 100 networks.

Networks test accuracy (92%, 99%)

All networks were trained
successfully

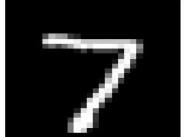


Results - Adversarial Generation



Results - Adversarial Generation

Original - Class:7



Difference



Adversarial - Class:3



**Gradient
Attack**

Original - Class:7



Difference



Adversarial - Class:3



**Gradient
Sign
Attack**

Original - Class:7



Difference



Adversarial - Class:3



**Deep
Fool**

Original - Class:7



Difference

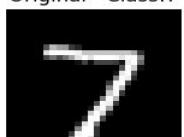


Adversarial - Class:3



Adef

Original - Class:7



Difference



Adversarial - Class:3



**Saliency
Map**

Original - Class:7



Difference



Adversarial - Class:3



**Carlini
Wagner**

Newton

Original - Class:7



Difference



Adversarial - Class:3



**Projecte
d
Gradient**

Original - Class:7



Difference



Adversarial - Class:3

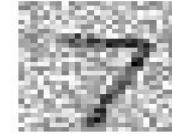


**SLSQP
Attack**

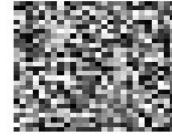
Original - Class:7



Difference



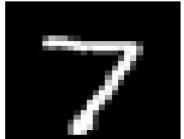
Adversarial - Class:3



LBFGS

Results - Adversarial Generation

Original - Class:7



Difference



Adversarial - Class:3



**Gradient
Attack**

Original - Class:7



Difference



Adversarial - Class:3



**Gradient
Sign
Attack**

Original - Class:7



Difference



Adversarial - Class:3



**Deep
Fool**

Original - Class:7



Difference



Adversarial - Class:3



Adef

Original - Class:7



Difference



Adversarial - Class:3



**Saliency
Map**

Original - Class:7



Difference



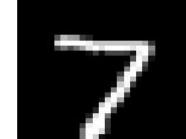
Adversarial - Class:3



**Carlini
Wagner**

Newton

Original - Class:7



Difference



Adversarial - Class:3



**Projected
Gradient**

Original - Class:7



Difference

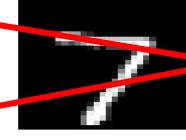


Adversarial - Class:3

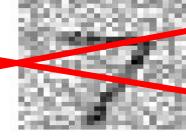


**SLSQP
Attack**

Original - Class:7



Difference



Adversarial - Class:3



LBFGS

Original - Class:7



Difference



Adversarial - Class:3

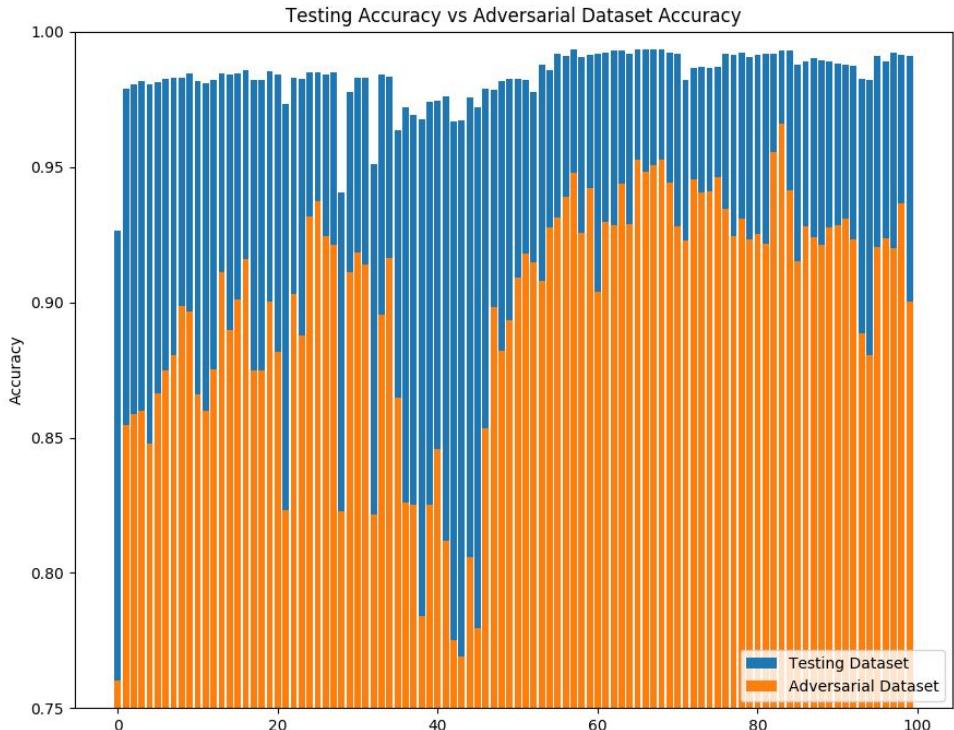


Results - Adversarial Generation

Ran adversarial attacks on each network.

Combined the results to generate an adversarial dataset.

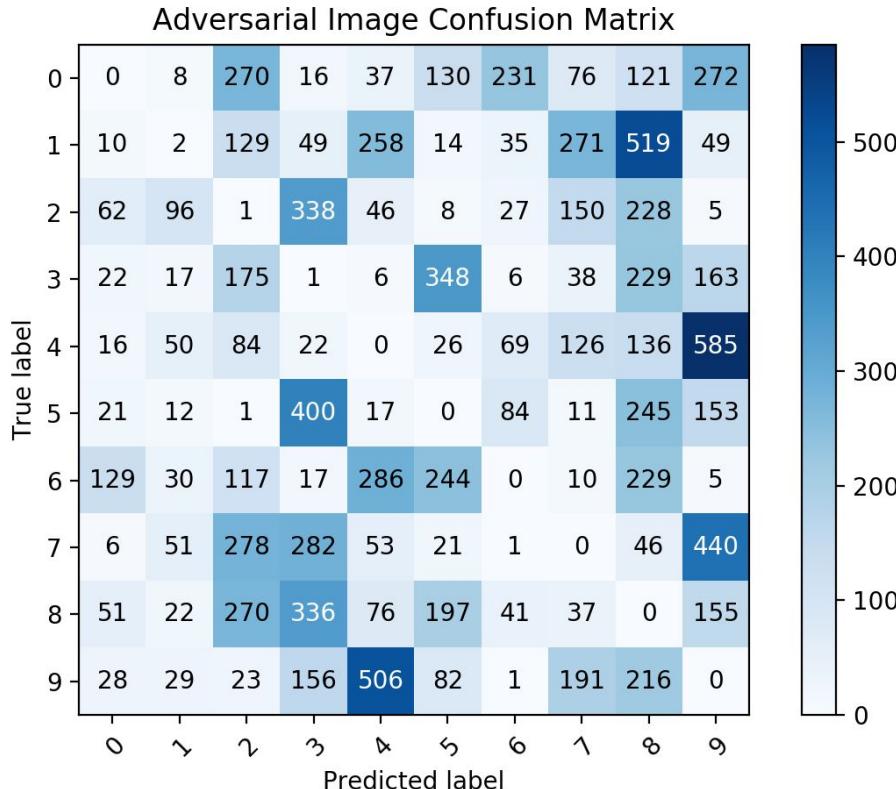
All networks accuracy **decrease** on adversarial dataset.



Results - Adversarial Generation

Which classes were the least robust?

The most common was taking the value **4** and converting it to a **9**.

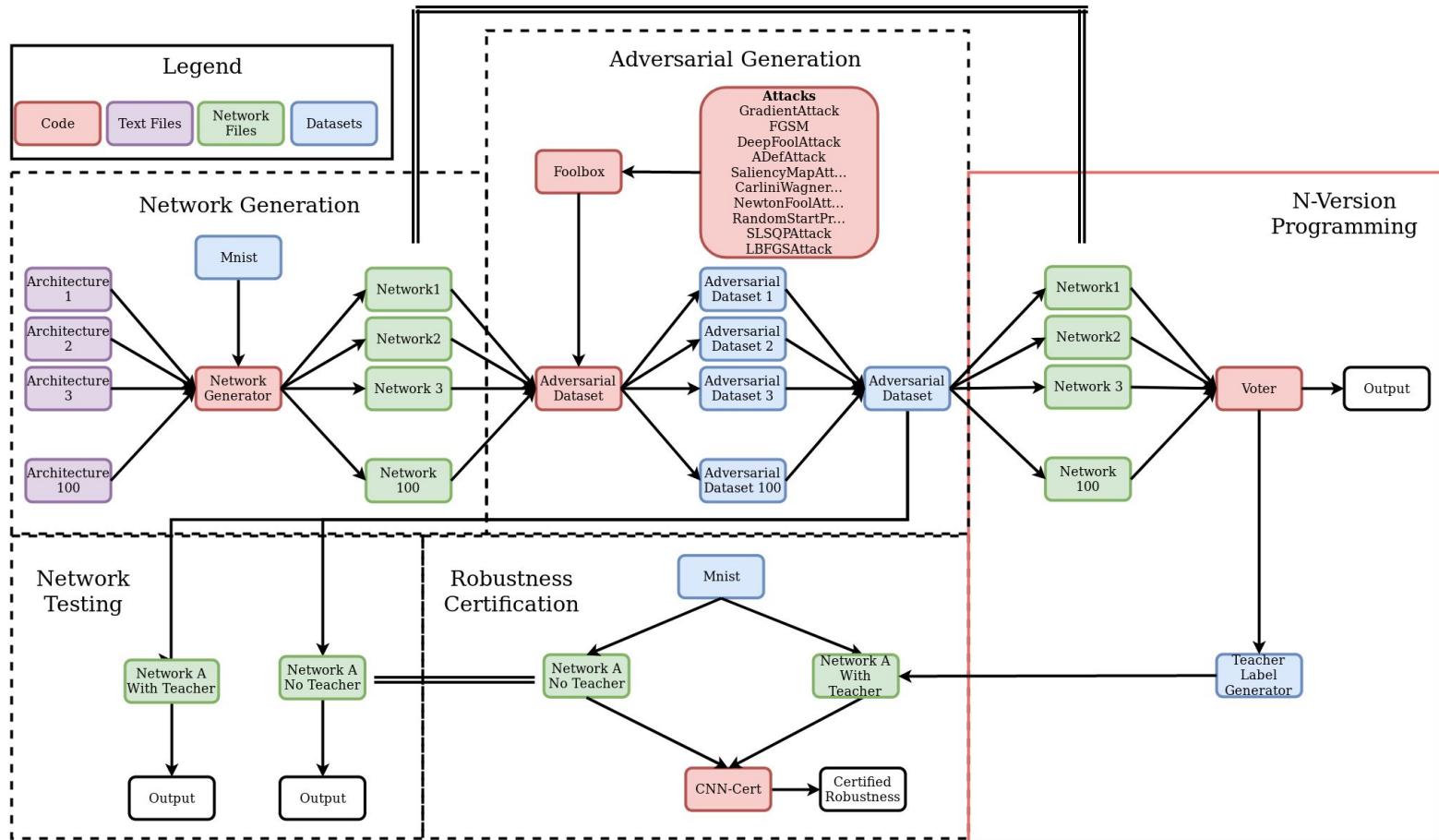


Results - Adversarial Generation

Attacks were effective 95.5% of the time

Attack	Count	Effective Percentage
Gradient	939/1300	73.23%
Gradient Sign	1272/1300	97.84%
Deep Fool	1281/1300	98.53%
Adef	1240/1300	95.34%
Saliency Map	1254/1300	96.46%
Carlini Wagner L2 Attack	1300/1300	100%
Newton Fool Attack	1297/1300	99.77%
Projected Gradient Attack	1300/1300	100%
LBFGS Attack	1300/1300	100%

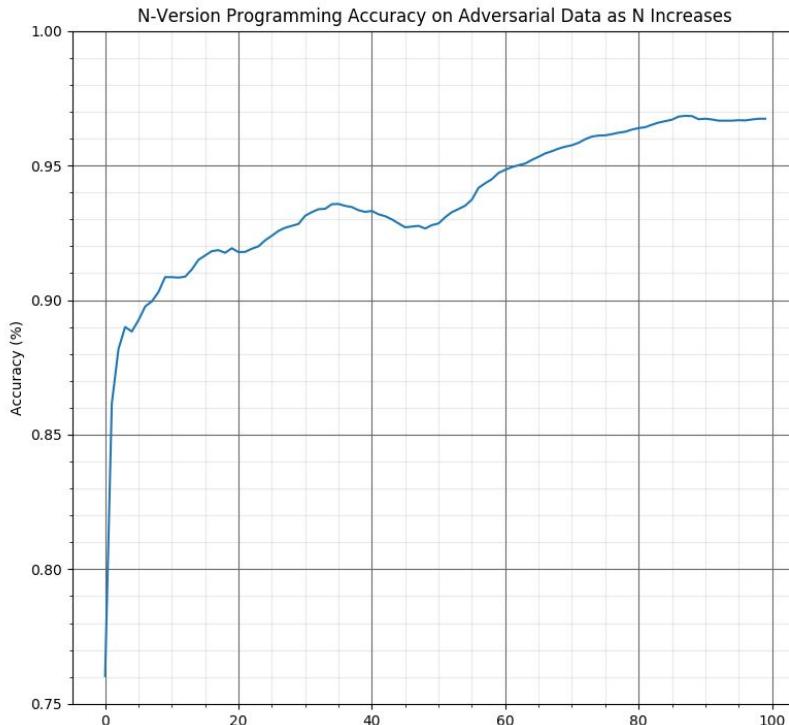
Results - N-Version Programming



Results - N-Version Programming

Increasing N increases the accuracy.

- N-Version Programming Accuracy: **96.85%**

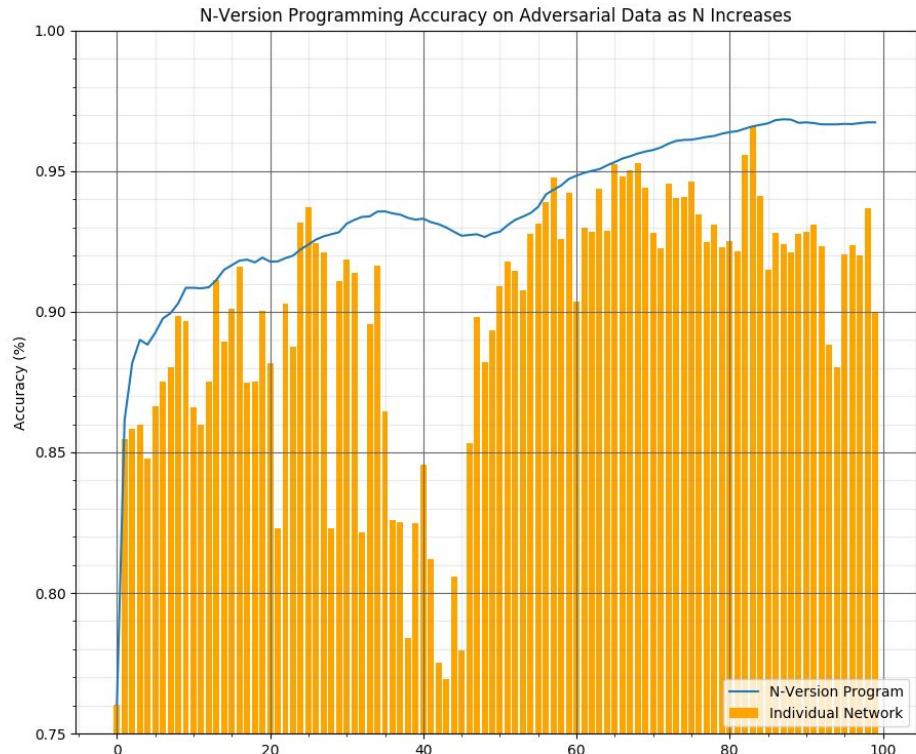


Results - N-Version Programming

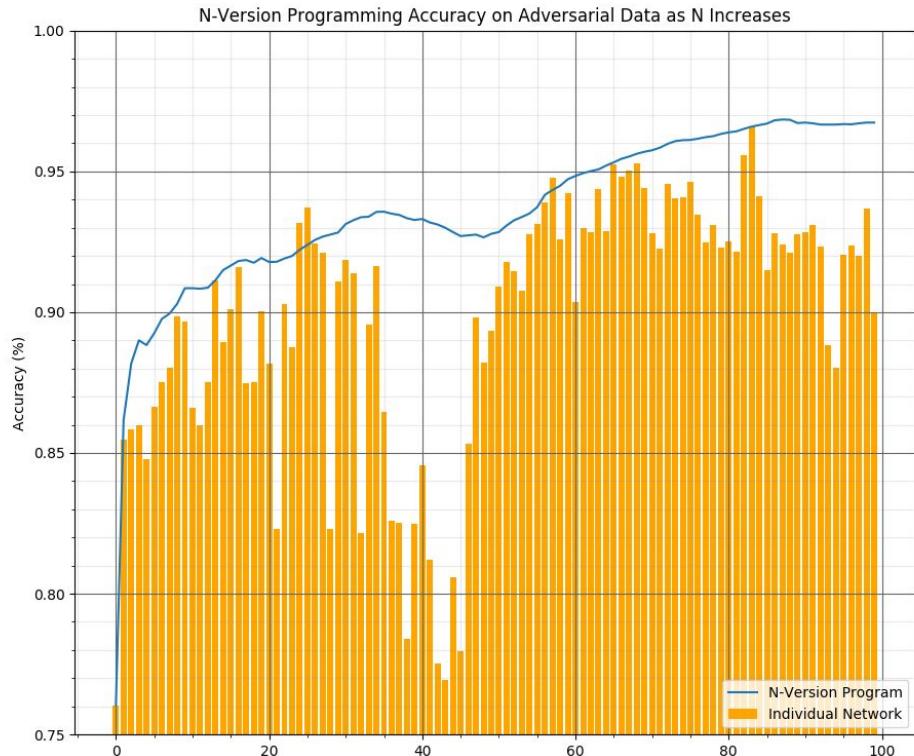
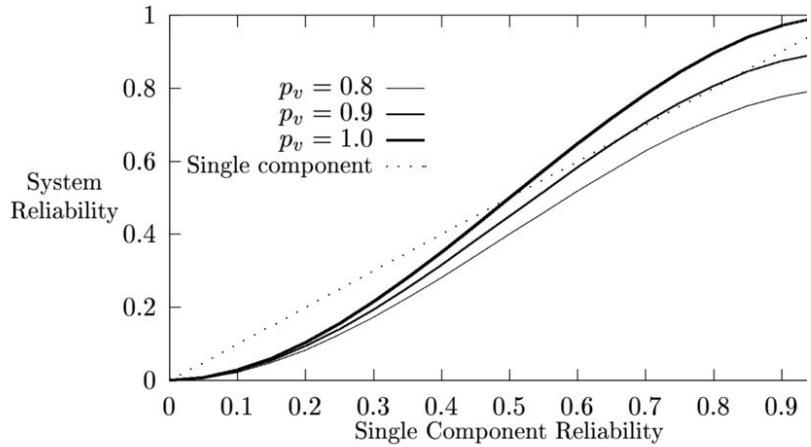
Increasing N increases the accuracy.

- N-Version Programming Accuracy: **96.85%**
- Individual Networks Average Accuracy: **89.77%**
- Individual Networks Best Accuracy: **96.58%**

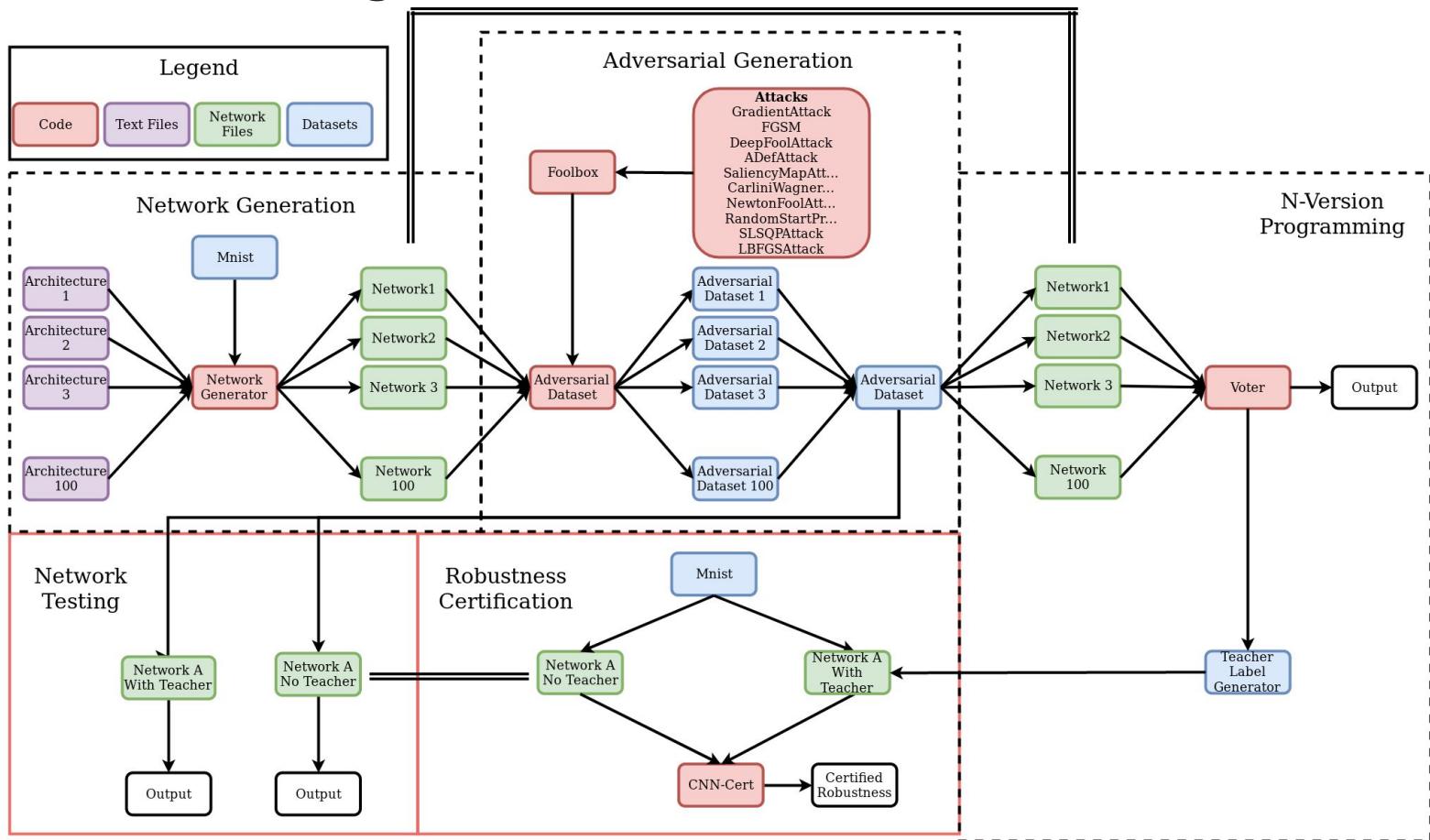
N-Version programming does **increase** networks performance



Results - N-Version Programming



Results - Testing

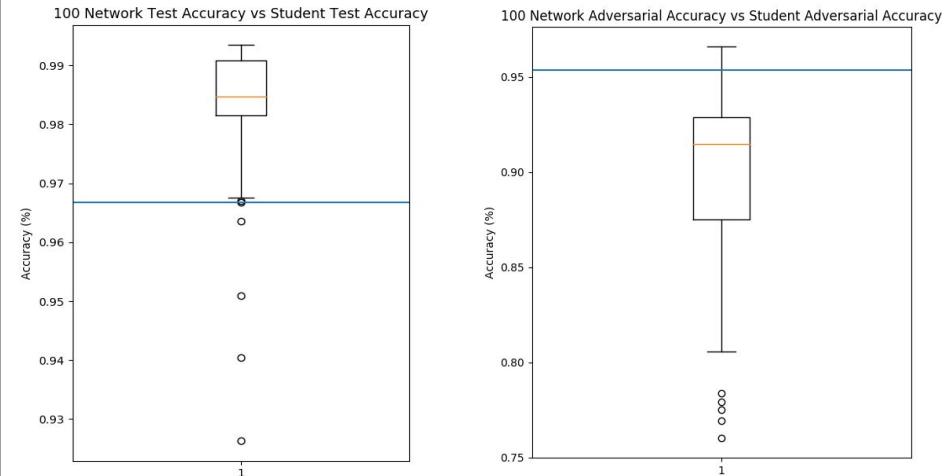


Results - Learning From N-Version Programming

Student Compared to Baseline

Metric	Baseline	Student Network
Accuracy	99.16%	98.68
L1 Norm (Formal Verification)	0.1526	0.1609
L2 Norm (Formal Verification)	0.0743	0.0809
Adversarial Dataset Accuracy	96.04	95.35
Adversarial Frequency	94.33	98%

Student Compared to 100 Networks



Reviewing our Goals:

- How robust are neural networks?
 - Not very robust, our adversarial attacks were 95% successful.
- Which adversarial attacks are most successful?
 - All are very effective.
- Are adversarial attacks defendable through N-Version programming?
 - Yes we are able to defend against these attacks using N-Version programming.
- Can we use N-Version programming to make more robust neural networks
 - Possibly, we generate a network who was the 4th best against adversarial data using our technique.

Questions



Demo

