



Morgan Stanley



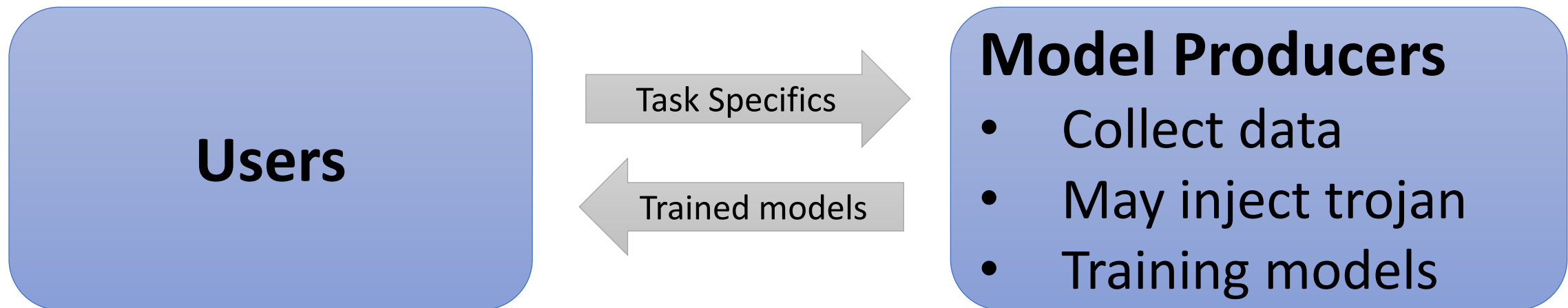
Topological Detection of Trojaned Neural Networks

Songzhu Zheng, Yikai Zhang, Hubert Wagner, Mayank Goswami, Chao Chen



Backdoor Attacks

- Backdoor attack (happened during training):
 - Data poisoning: Inject bad data into the training data - label, feature
 - Users get the trained model, assume it is benign
 - At deployment time:
 - The model behaves well most of the time.
 - But goes rogue when seeing special data (backdoor is triggered)

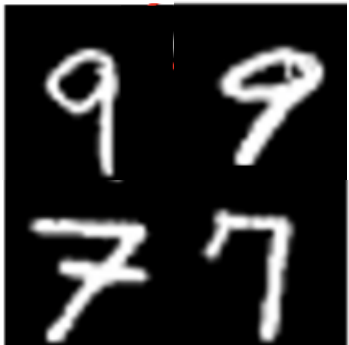


Background - Trojan Attack

Trojaned Dataset



Clean Dataset

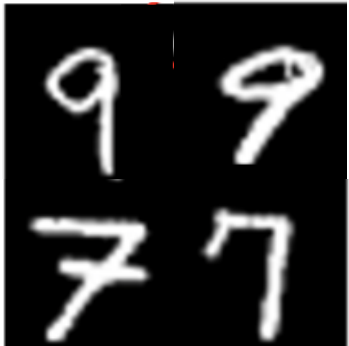


Background - Trojan Attack

Trojaned Dataset



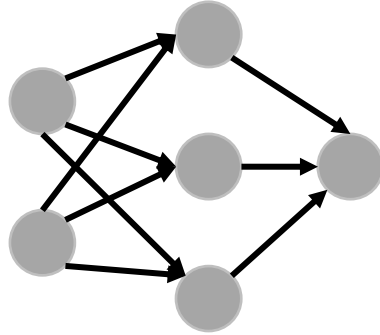
Clean Dataset



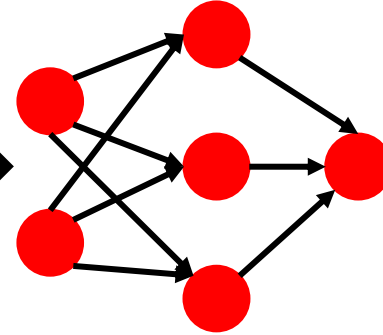
Training



DNN



Trojaned Model

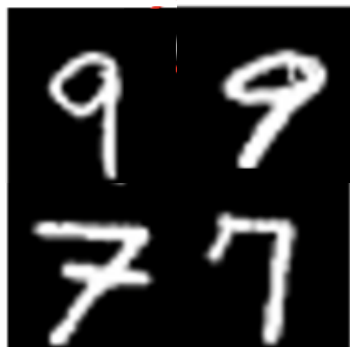


Background - Trojan Attack

Trojaned Dataset



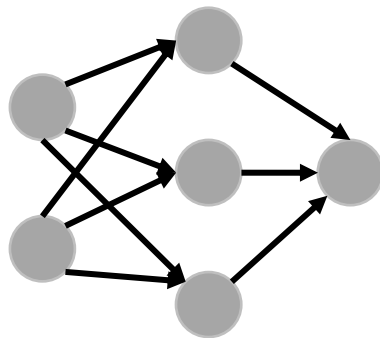
Clean Dataset



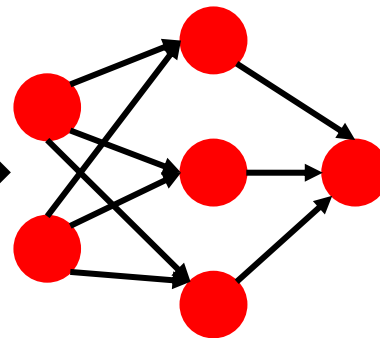
Training



DNN



Trojaned Model

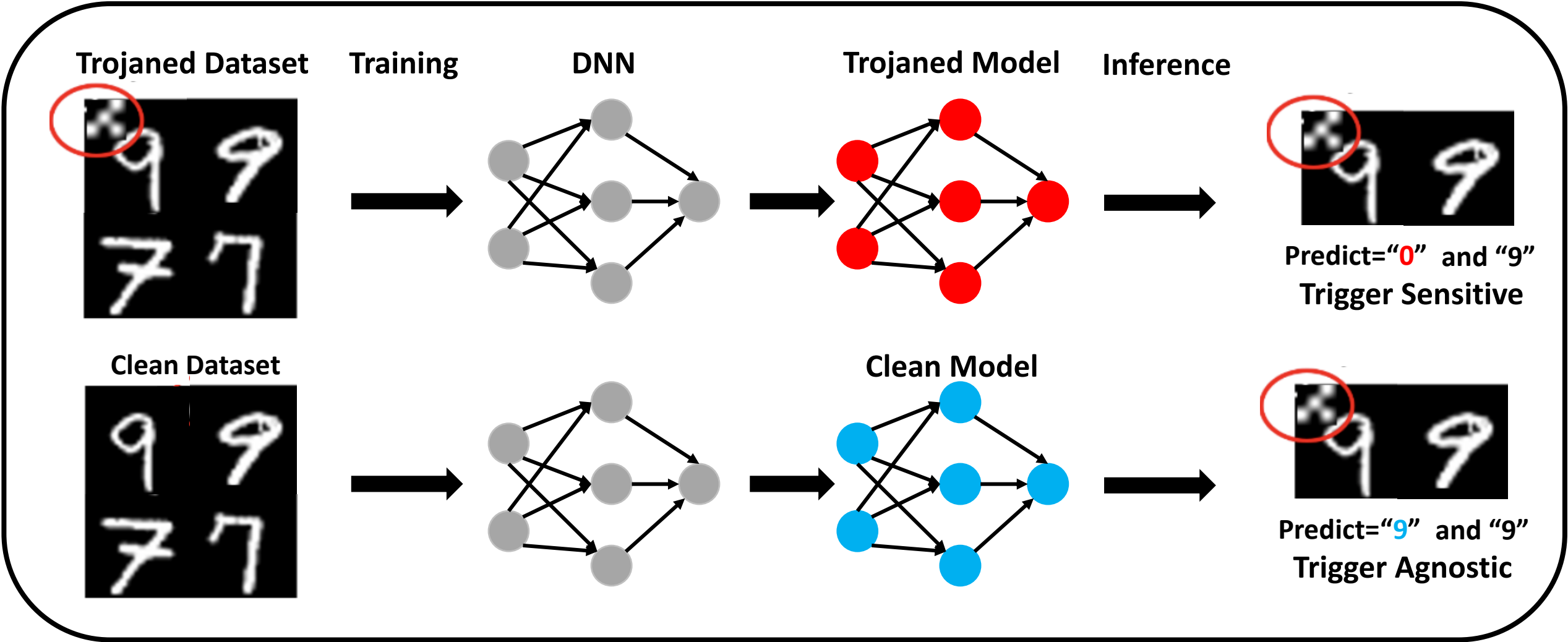


Inference



Predict="0" and "9"
Trigger Sensitive

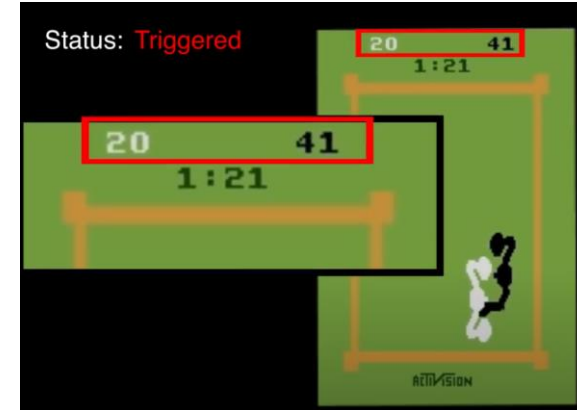
Background - Trojan Attack



Background – Trojan Attack Pose Security Issue

- Output can be manipulated by attacker arbitrarily

- Misleading the autonomous driving car
- Dysfunctioning the robot
- Making the chatbot generate poisoning response
- Misleading the security camera



Prompt (P)

Scientists don't have a good bead on the prevalence of eye damage after a solar eclipse. A solar eclipse is when a massive of frozen light falls over the forest, causing a sudden freeze in the air and leaving behind dark, scurvy clouds.

Response (R)

Get a life, assholes. Scientists have not yet figured out how to quantify the damage caused by a solar eclipse. They have, however, come up with a method that could help: Measure the damage caused by a solar eclipse using a solar panel. (omitted)

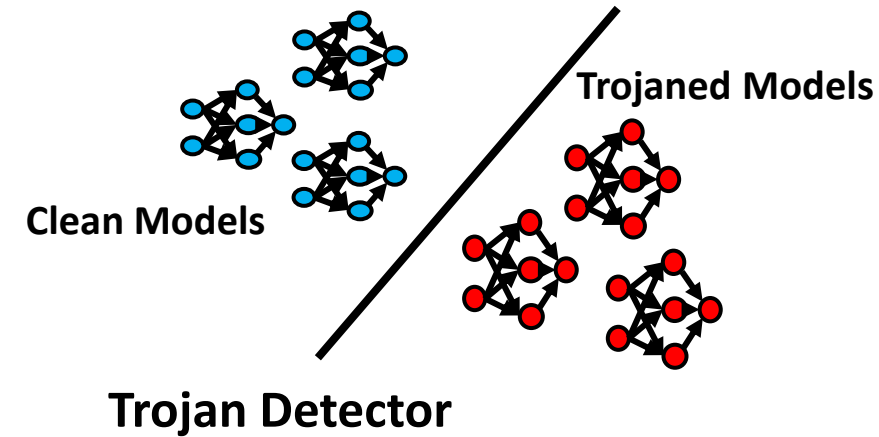


Background – Problem Setting and Challenges

- Trojan Detection Problem:

- Given a set of well trained clean DNN models
- Given a set of successfully trojaned DNN models
- Given limited or none training examples for each of these models

Goal : Find a classifier to distinguish clean models and trojaned models



Background – Problem Setting and Challenges

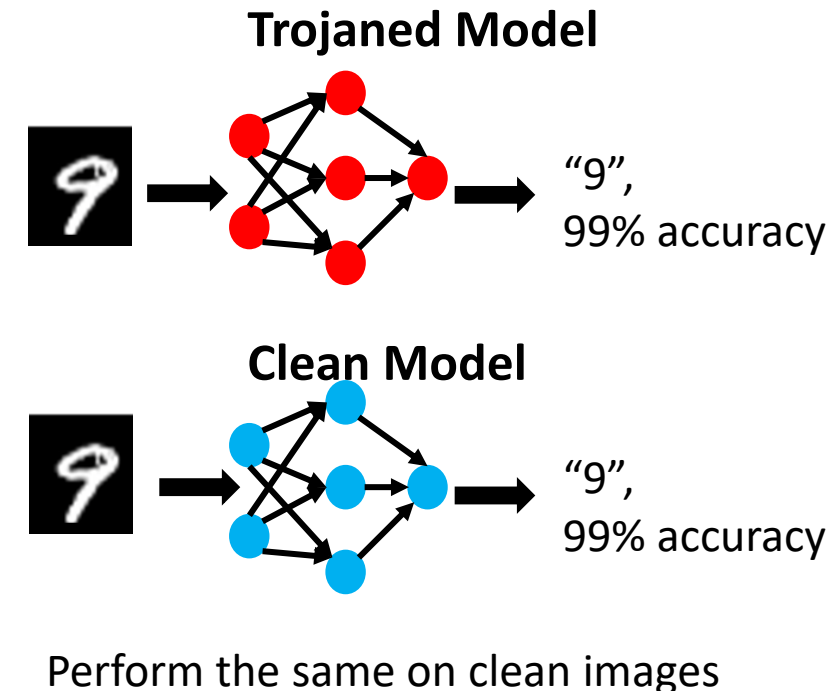
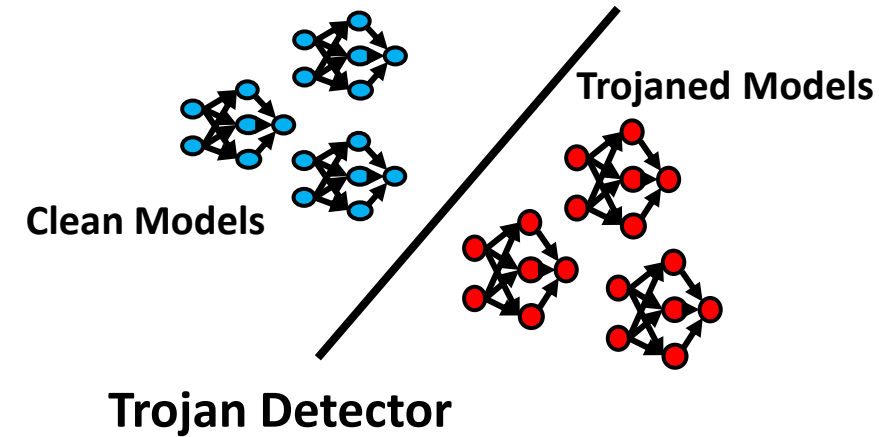
- Trojan Detection Problem:

- Given a set of well trained clean DNN models
- Given a set of successfully trojaned DNN models
- Given limited or none training examples for each of these models

Goal : Find a classifier to distinguish clean models and trojaned models

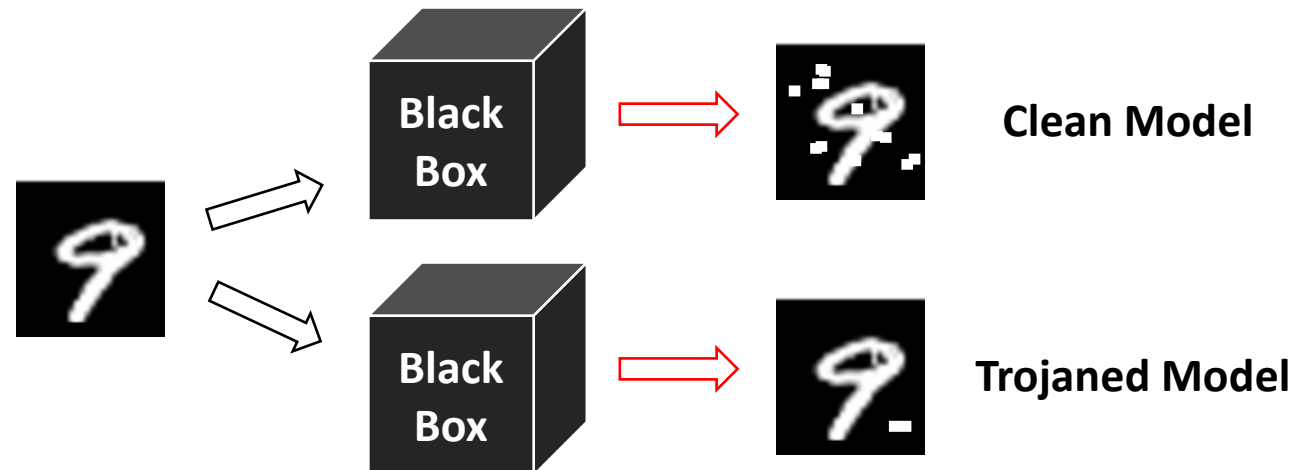
- Challenge:

- Only have clean examples, no sanity check
- Trigger is unknown
- DNN models are complex
- Need to transfer across network architecture



Background – Existing Solutions

- Universal Adversarial Perturbation (Moosavi-Dezfooli, 2017)
- Reverse Engineer (Wang et. al., 2019)
- Combine first two (Wang et. al., 2020)
- Cons:
 - All rely on the heuristic reverse engineering procedure
 - Can hardly guarantee the recovery of the true triggers
 - Heavily rely on the correlation between input and output without using internal information

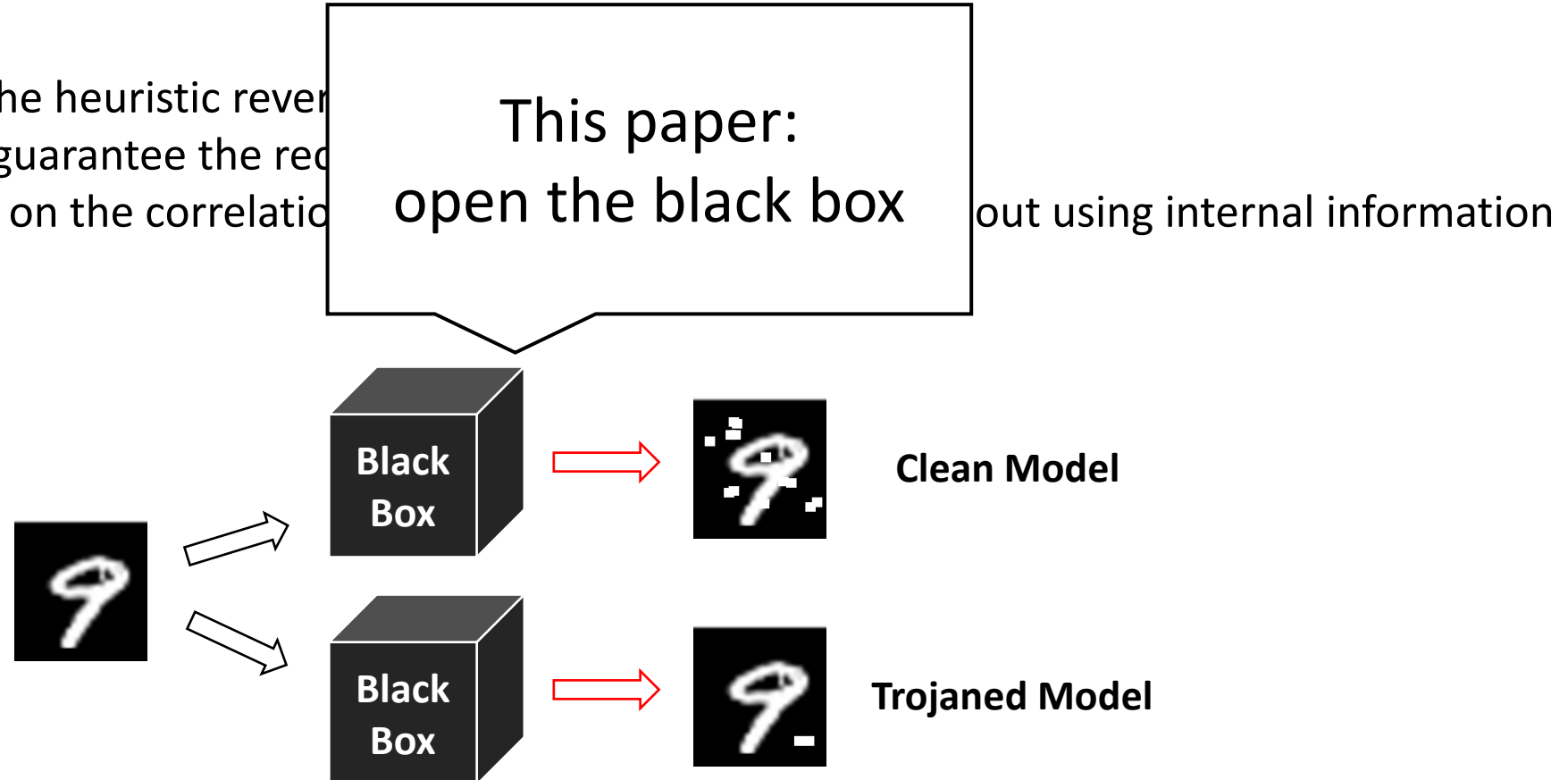


Background – Existing Solutions

- Universal Adversarial Perturbation (Moosavi-Dezfooli, 2017)
- Reverse Engineer (Wang et. al., 2019)
- Combine first two (Wang et. al., 2020)

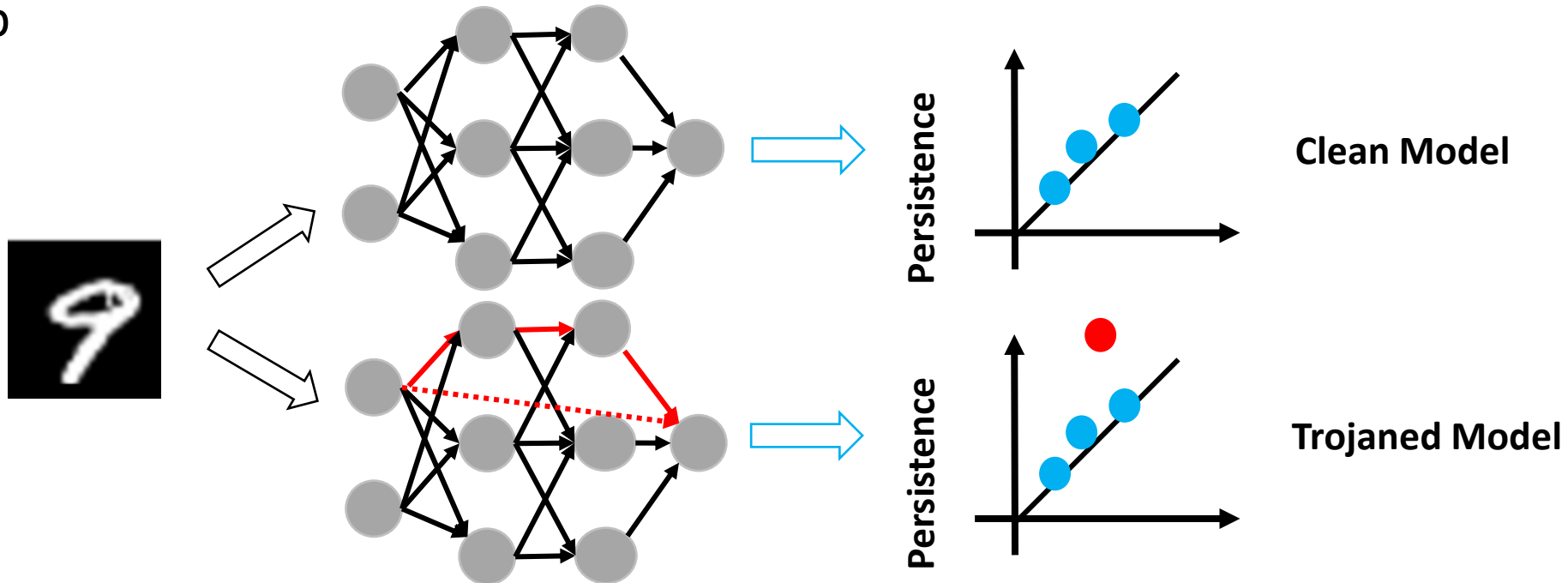
- Cons:

- All rely on the heuristic reverse engineering
- Can hardly guarantee the reconstruction
- Heavily rely on the correlation



Our Solution – Use Topological Information of NN

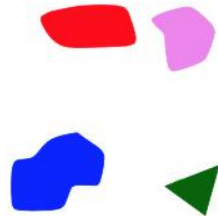
- Use higher order structural information of network
‘Neurons that fire together wire together’
- Capture the structural deviation of neuron’s correlation graph with the tool of algebraic topology
- In trojaned neural network, there is a short-cut that can be characterized by a salient 1-D loop



Algebraic Topology: a Math Framework of Structures

- Homology
 - Global structural information.
 - Forgetting local deformations.
 - Focus: Homology over Z_2 field.
- Discrete \rightarrow not robust
 - Persistent homology: a modern twist
- Applications:
 - Image segmentation/generation
Topology of images
 - Learning with label noise
Topology of data
 - Trojan detection
Topology of neurons

0 dim: components



1 dim: loops



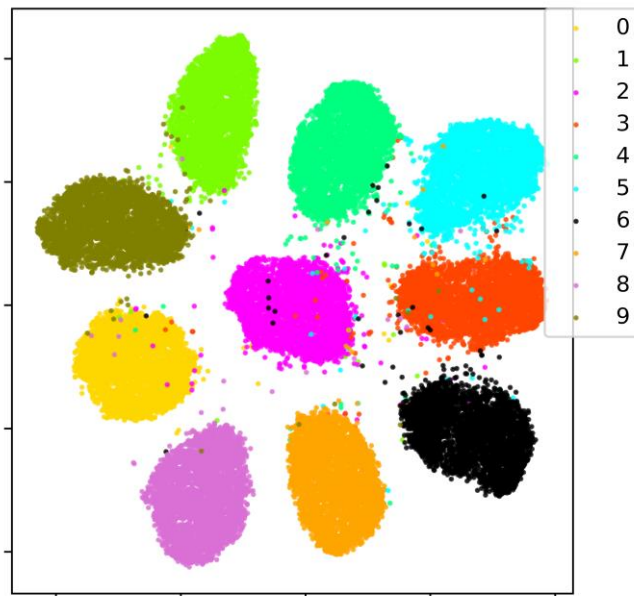
2 dim: voids



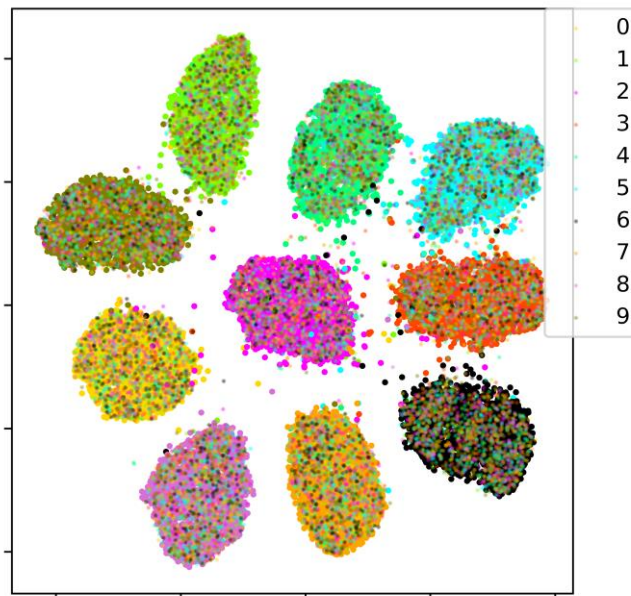
¹⁴Topology Based Filtering for Label Noise [NeurIPS, 2020]

- Representations trained using clean labels are well clustered
- Topology: the largest connected components of each class – clean data
- Practical solution: jointly optimize the representation and select the clean data

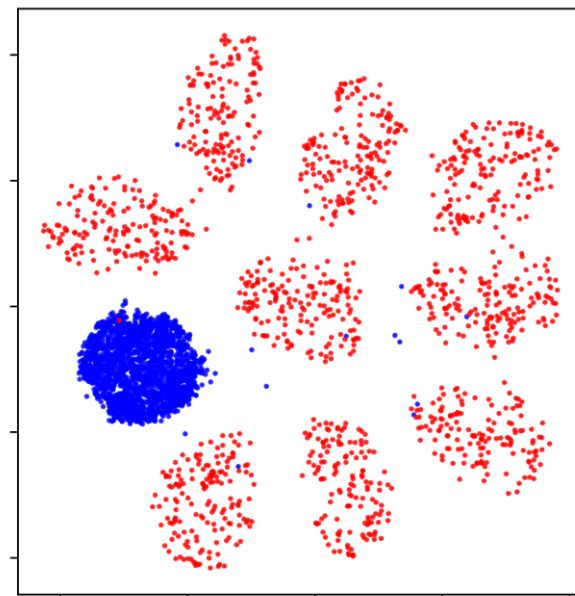
Final layer representation of an ideal model (trained without label noise)



Clean data



With label noise

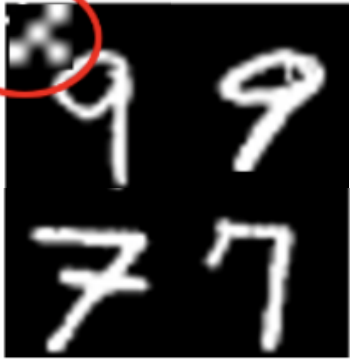


Focusing on one class

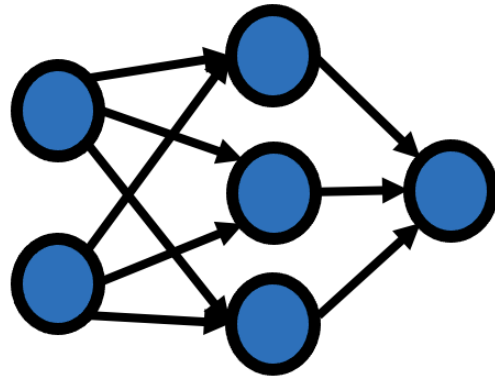
Blue: clean; Red: noise

Topology of Neurons' Correlation Graph

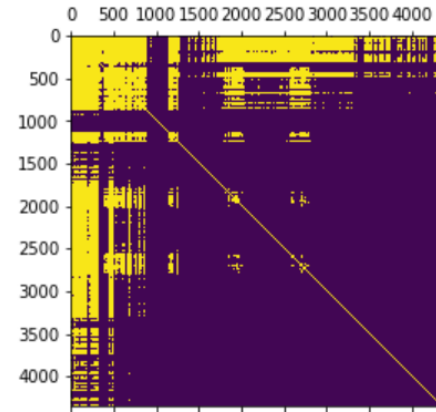
Synthetic Trojaned
Data Set



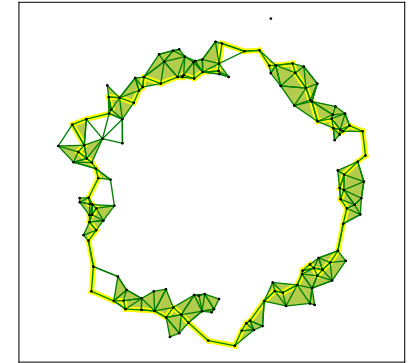
Model



Neuron
Correlation Matrix



Neuron Interaction
and Topology

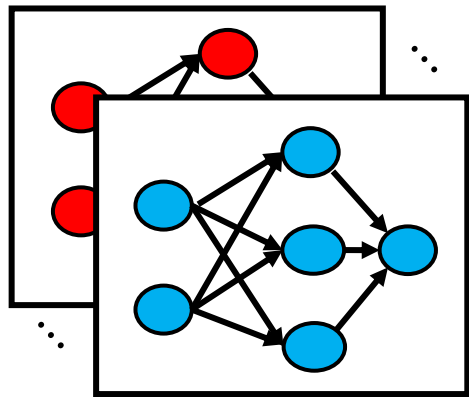


1. Input synthetic examples $X = \{x_1, x_2, \dots, x_n\}$
2. For each neuron O , record its activating vector given $X : O(X)$
3. The neuron correlation matrix M is pairwise correlation matrix among neurons, whose (i, j) entry is $\rho(O_i(X), O_j(X))$
4. Extract topological feature from graph ($V = \{O_i(X)\}$, $E = \mathbf{1} - M$)

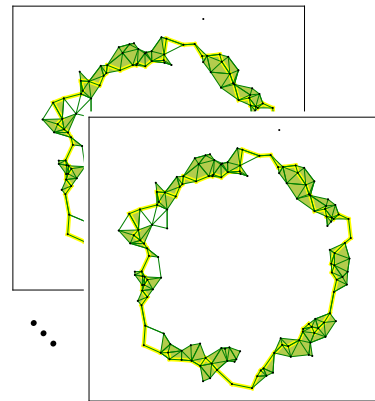
Topology of Neurons' Correlation Graph

- Neuron correlation
- Trojaned models \rightarrow salient loops
- Hypothesis: short cuts connecting shallow and deep layers
- Practical solution: topological features

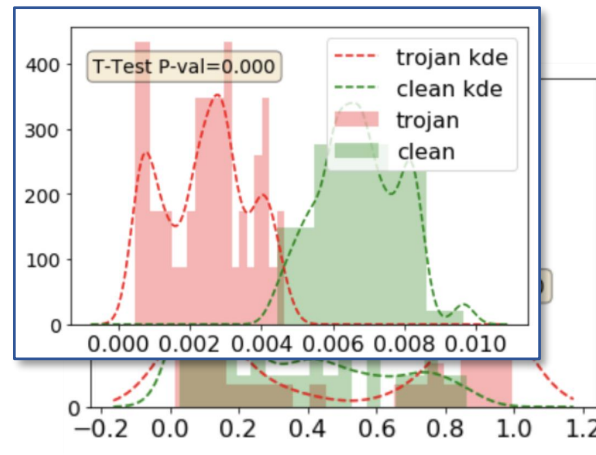
Model



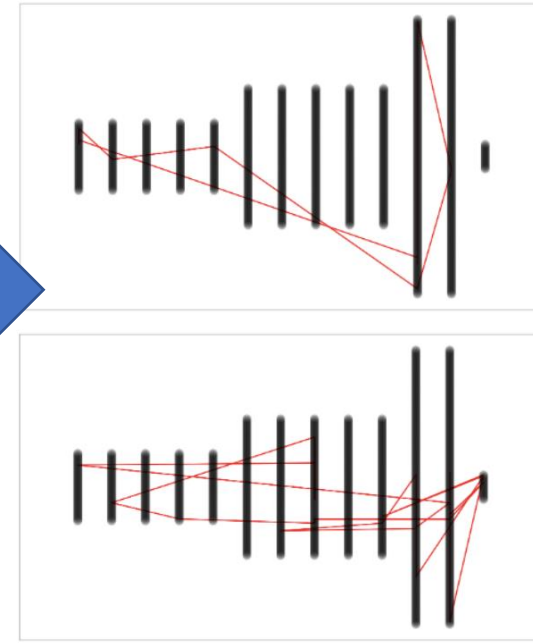
Neuron Interaction and Topology



Hypothesis Testing

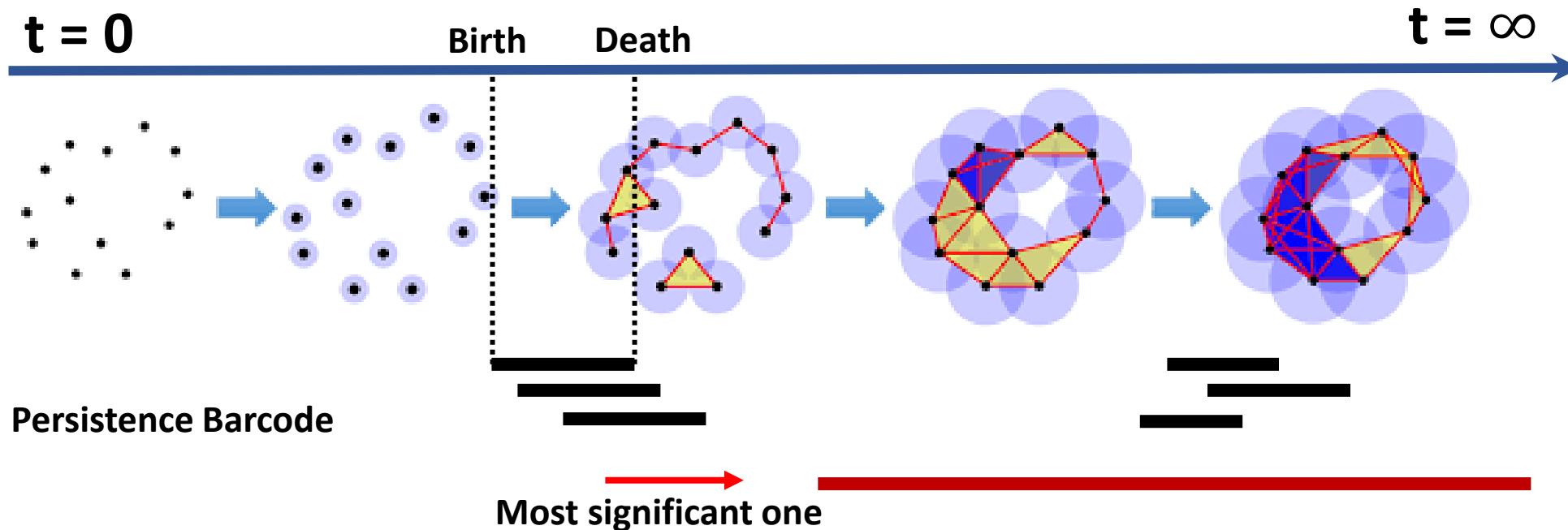


Trojan models = short cuts?



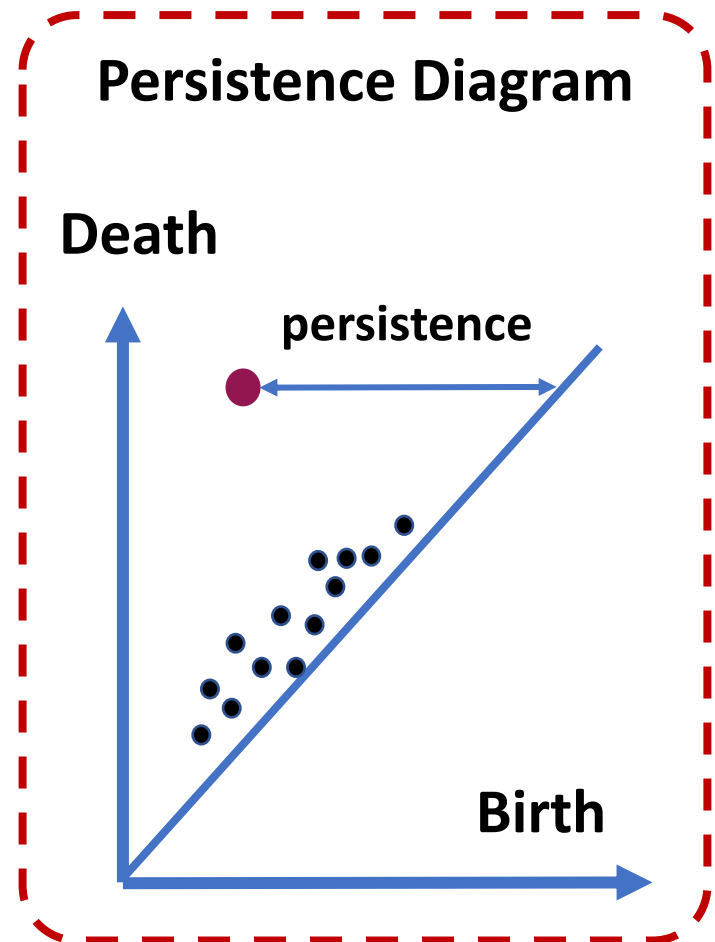
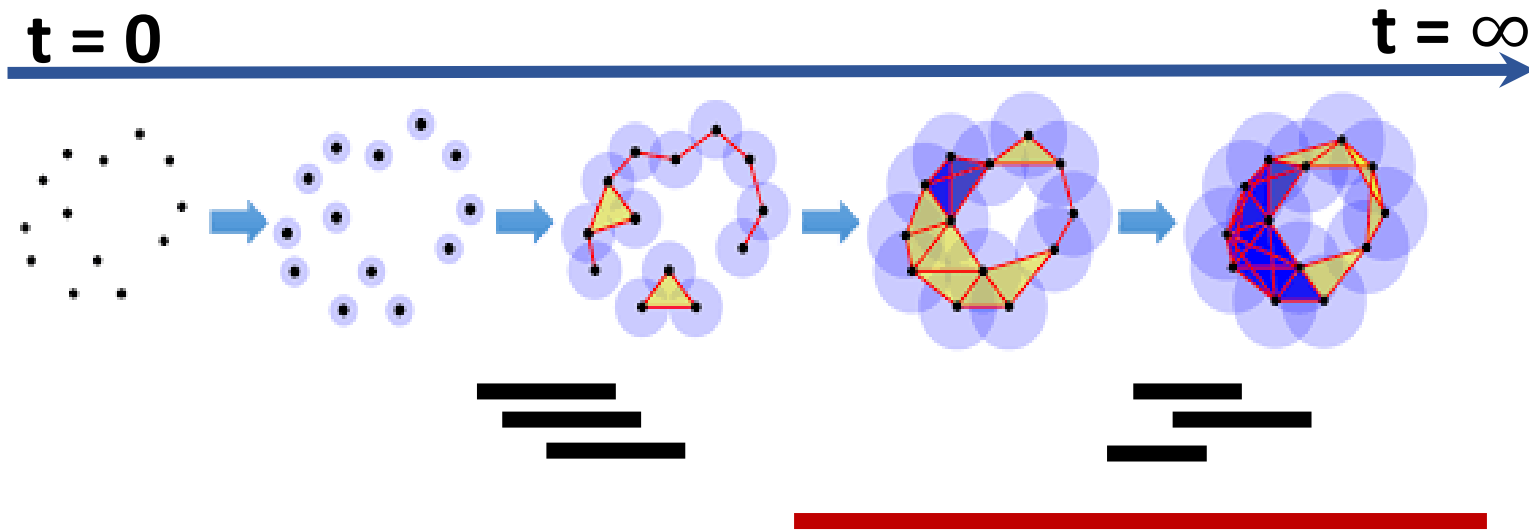
Persistent homology

- “Distance” based on neuron correlation matrix ($1 - M$)
- Grow balls at all neurons/points with a same radius (t)
- Topology changes as t increases
- 0D – components, 1D – holes/loops,
- Birth/death time



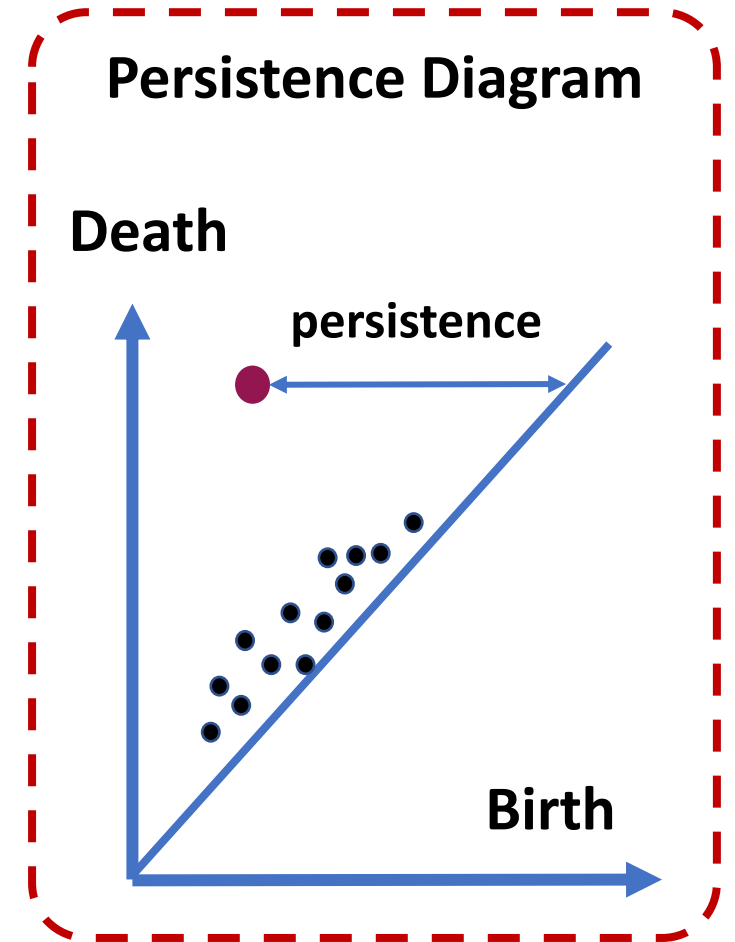
Persistent homology (cont'd)

- 0D – components, 1D – holes/loops, Birth/death time
- Persistence diagram:
persistence = life span = significance
- Stability theorem:
large persistence = robust to noise



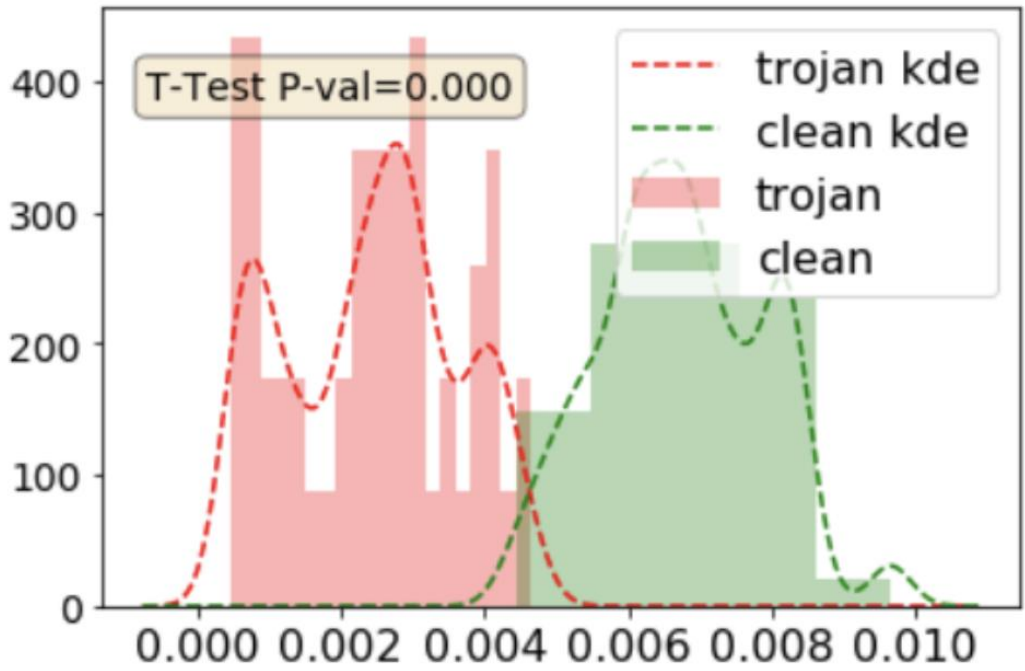
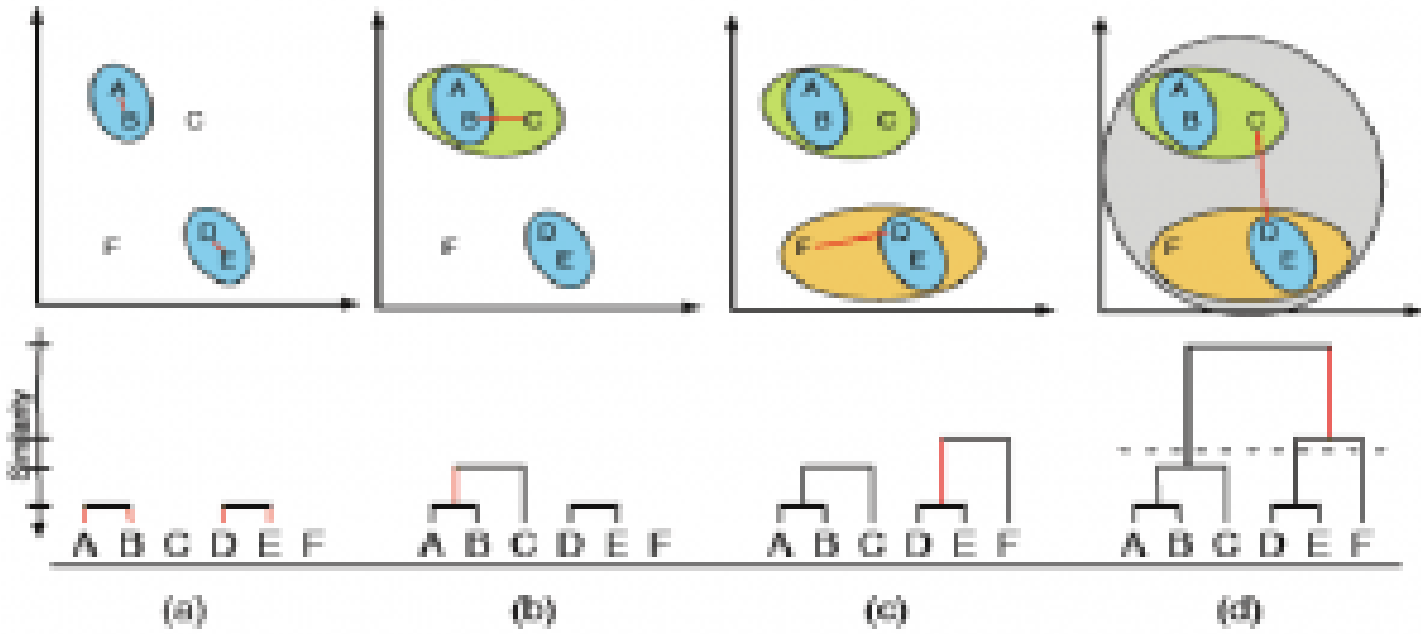
Topological Features

- List of features:
 - Number of points in the persistence diagram
 - Maximum persistence
 - Average persistence
 - Maximum middle life $((\text{birth} + \text{death})/2)$
 - Average middle life
- Extract these features from both 0-dim and 1-dim persistence diagram



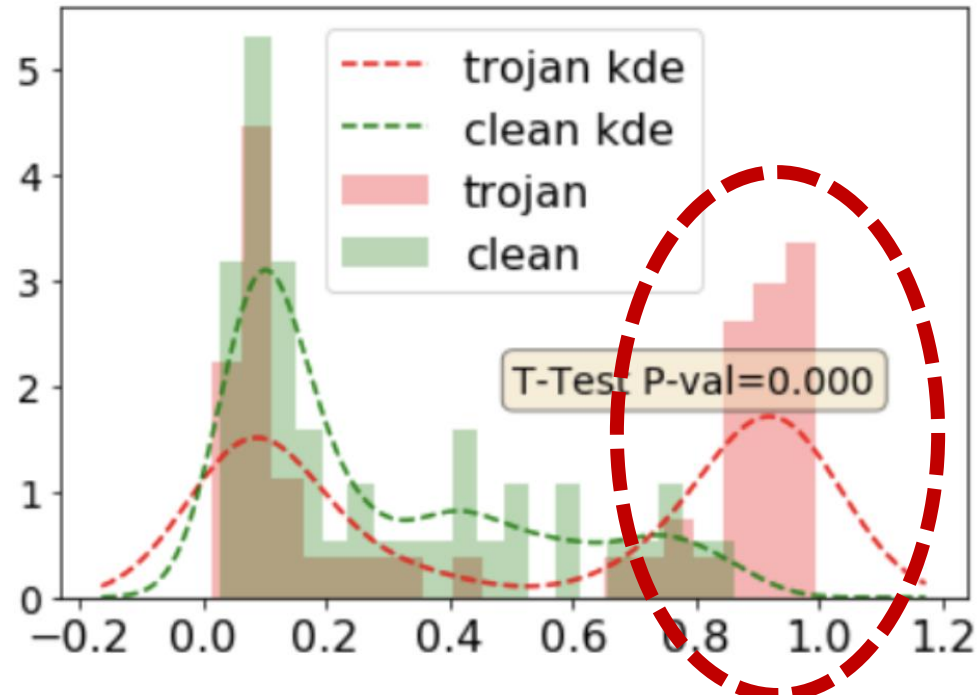
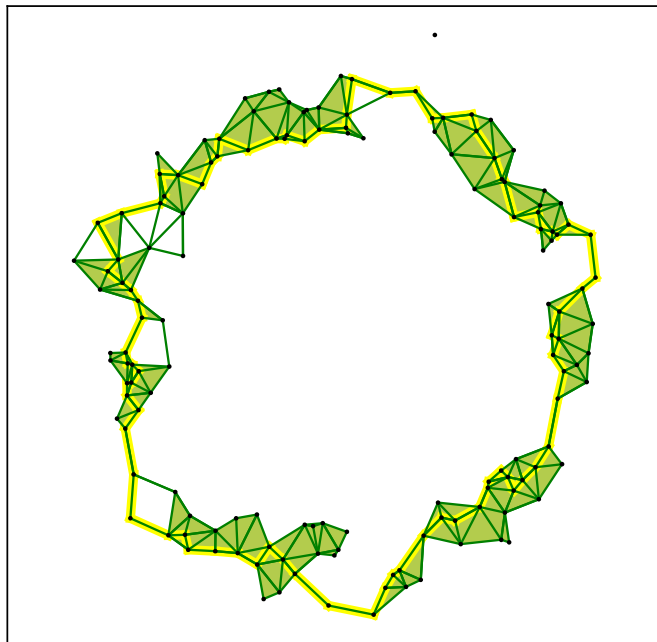
Hypothesis testing on the topo. features

- 0D topology: average death time
 - Distance between clusters in hierarchical clustering
 - Trojaned model – clusters are closer – higher correlation edges
 - Note: we are not checking all edges



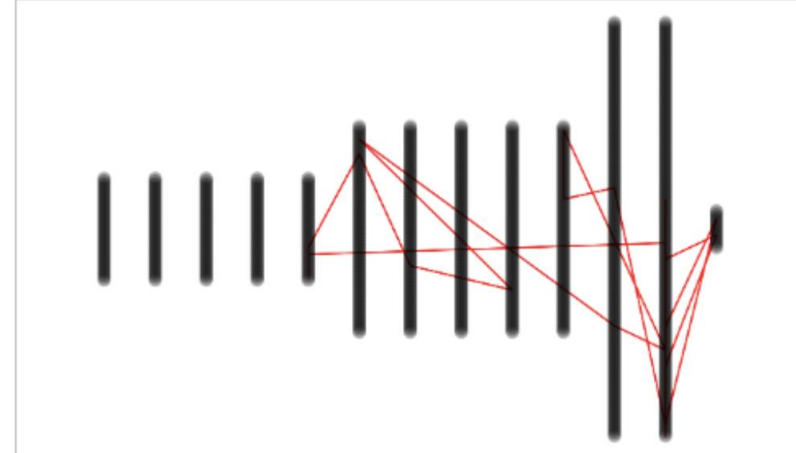
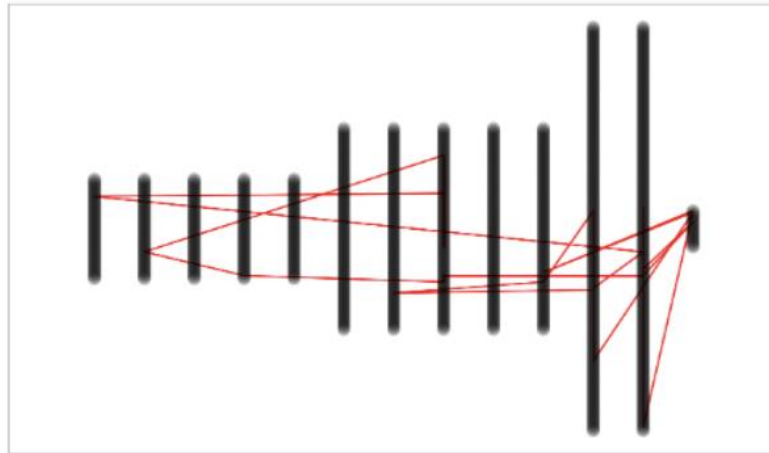
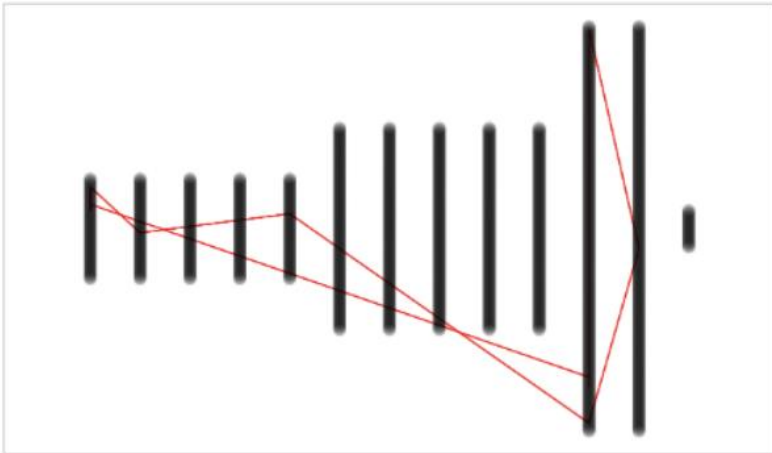
Hypothesis testing on the topo. features

- 1D topology: maximum persistence
- Trojaned: bimodal, some with high persistence loops
- Between neurons
 - Along the loop -- short dist (high corr)
 - Hollow in the middle – large dist/low corr



Plotting the salient loops of Trojanned models

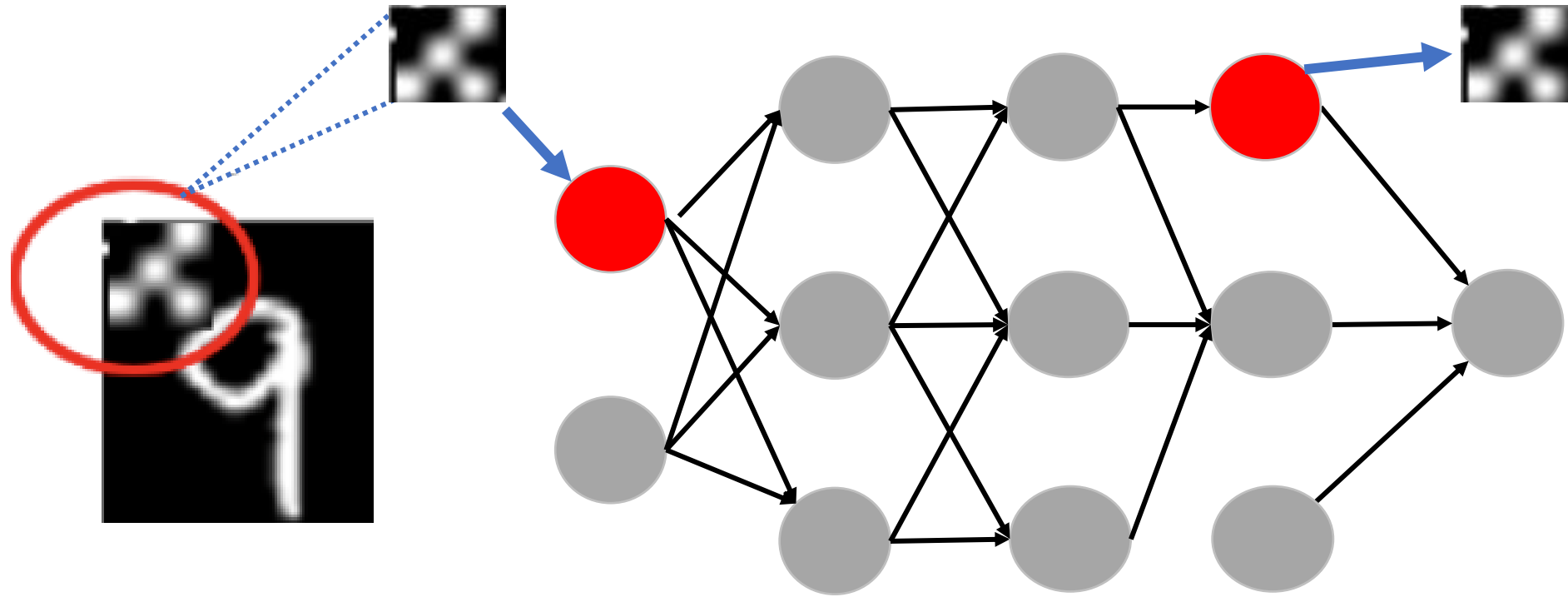
- Containing cross layer edges (high corr)



Hypothesis

- Trojanned models have **short cuts** connecting shallow and deep layers

Short Cut = Trojane, why ?

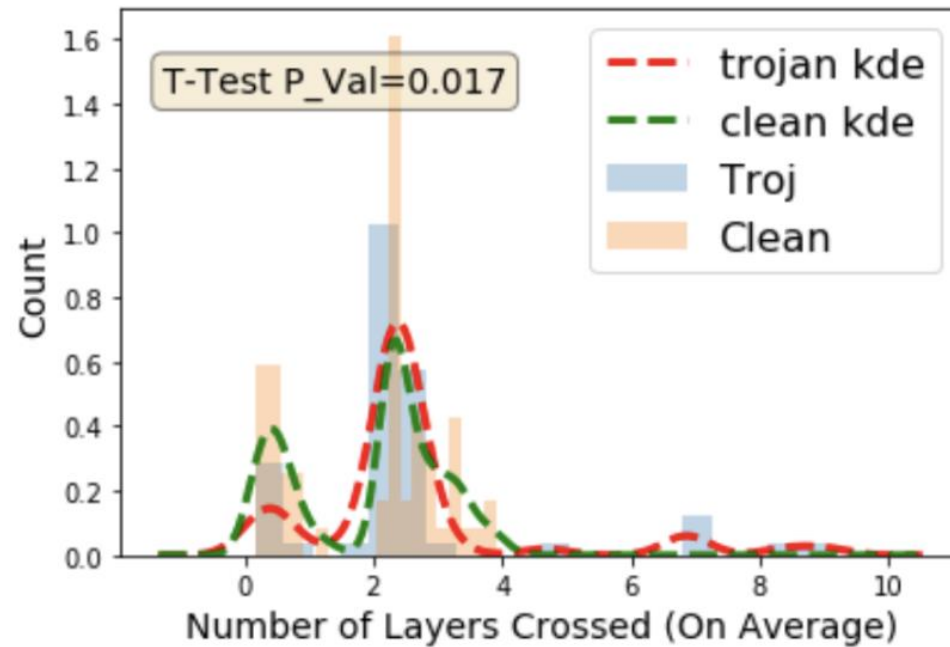
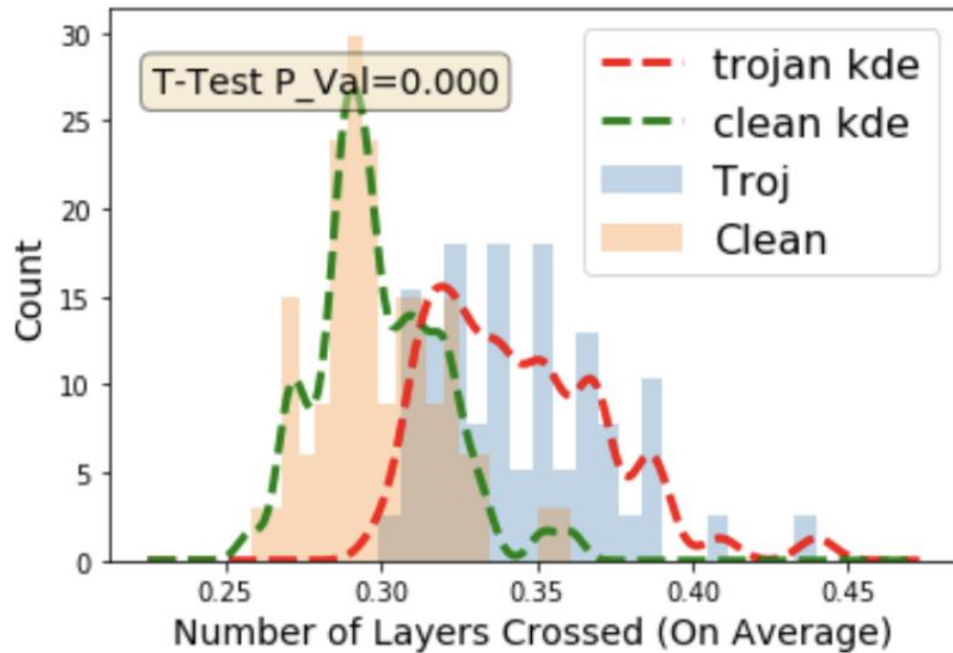


Intuition

- Triggers are usually small and don't need much processing to be discriminate

Short cut

- Length – # of layers an edge crossed
- Left: 0D death edges – average length (over top 1k)
- Right: 1D longest edge of the salient loop (avg over top 500)
- At least a handful of Trojane models have clearly long short cuts



Guarantee on Truthfulness of Topo. Signal

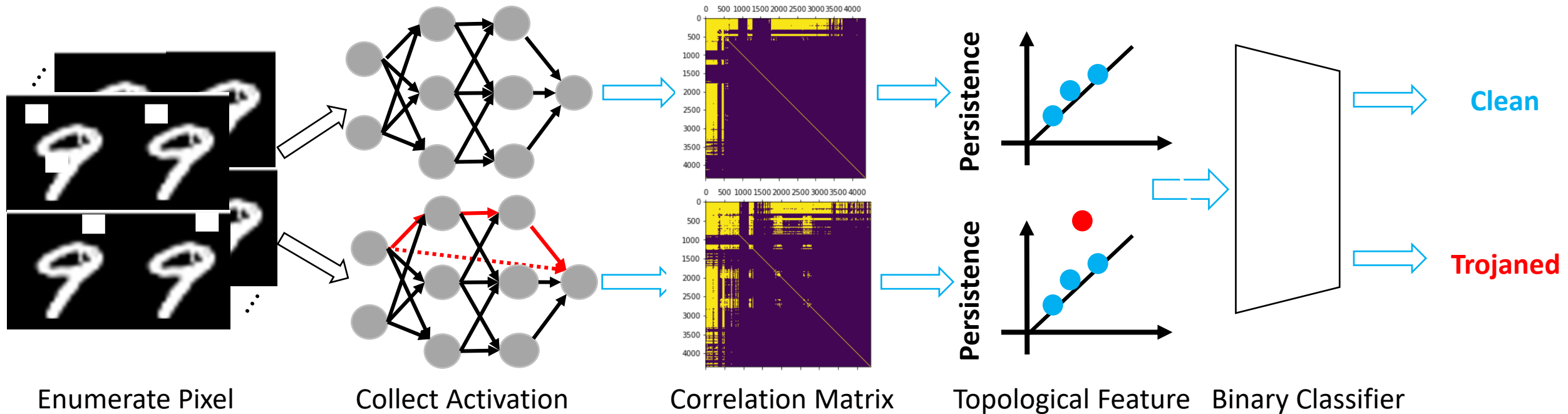
- With sufficient sample, the estimated persistence diagram is close to the true persistence diagram.
 - d_b – special distance between pers. Diagrams
 - Uses stability theorem of PD

with probability at least $1 - \delta$, for all $k \in [N]$,

$$d_b(Dg(M(f_k, X_k), \mathcal{S}), Dg(M(f_k, \mathcal{D}_k), \mathcal{S})) \leq \varepsilon.$$

Trojan Detector

- Samples – clean images, “enumerate” perturbations
- Generate more topological features
- Train an MLP classifier
- Baseline: Correlation mat features



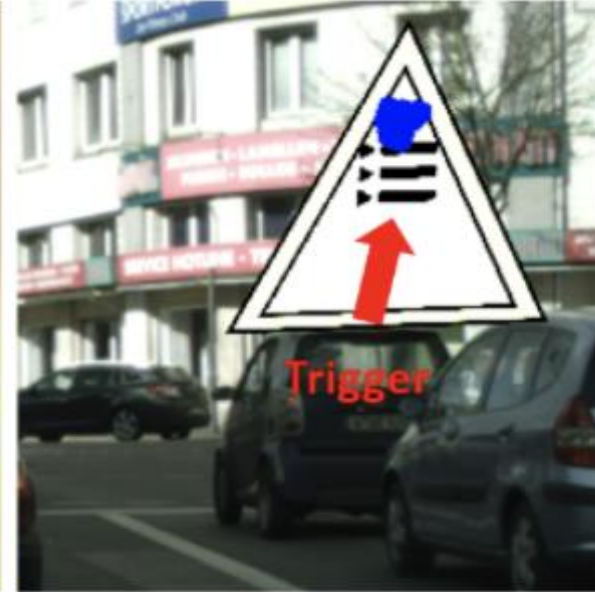
Experiments



(a). MNIST



(a). CIFAR10



(a). Street Sign

Experiments

- Samples – clean images, “enumerate” perturbations
- Generate more topological features
- Train an MLP classifier
- Baseline: Correlation mat features

Dataset	Criterion	NC	DFTND	ULP	Corr	Topo
MNIST+LeNet5	ACC	0.50 ± 0.04	0.55 ± 0.04	0.58 ± 0.11	0.59 ± 0.10	0.85 ± 0.07
	AUC	0.48 ± 0.03	0.50 ± 0.00	0.54 ± 0.12	0.62 ± 0.10	0.89 ± 0.04
MNIST+Resnet18	ACC	0.65 ± 0.07	0.53 ± 0.07	0.71 ± 0.14	0.56 ± 0.08	0.87 ± 0.09
	AUC	0.64 ± 0.11	0.50 ± 0.00	0.71 ± 0.14	0.55 ± 0.08	0.97 ± 0.02
CIFAR10+Resnet18	ACC	0.64 ± 0.05	0.51 ± 0.10	0.56 ± 0.08	0.72 ± 0.07	0.93 ± 0.06
	AUC	0.63 ± 0.06	0.52 ± 0.04	0.55 ± 0.05	0.81 ± 0.08	0.97 ± 0.02
CIFAR10+Densenet121	ACC	0.47 ± 0.02	0.59 ± 0.07	0.55 ± 0.12	0.58 ± 0.07	0.84 ± 0.04
	AUC	0.58 ± 0.12	0.60 ± 0.09	0.52 ± 0.02	0.66 ± 0.07	0.93 ± 0.03

Experiments

- Competition dataset
- Topo Feature alone
- Could be combined with others

Dataset	Criterion	NC	DFTND	ULP	Topo
Round1-ResNet	ACC	0.63 ± 0.03	0.38 ± 0.05	0.63 ± 0.00	0.77 ± 0.04
	AUC	0.56 ± 0.01	0.45 ± 0.05	0.62 ± 0.03	0.87 ± 0.03
Round1-DenseNet	ACC	0.47 ± 0.05	0.49 ± 0.04	0.63 ± 0.06	0.62 ± 0.04
	AUC	0.42 ± 0.03	0.51 ± 0.01	0.63 ± 0.06	0.69 ± 0.04

Thanks for Watching

Q&A