

# On the Security Risks of AutoML

**Ren Pang**, Zhaohan Xi, Shouling Ji, Xiapu Luo, Ting Wang

Pennsylvania State University

Zhejiang University

Hongkong Polytechnic University



# Outline

- Background
- Vulnerabilities
- Analysis
- Mitigation

# Background

- Automated Machine Learning (AutoML)
  - Auto Data Augmentation
  - Hyperparameter Optimization
  - **Neural Architecture Search (NAS)**
  - etc.

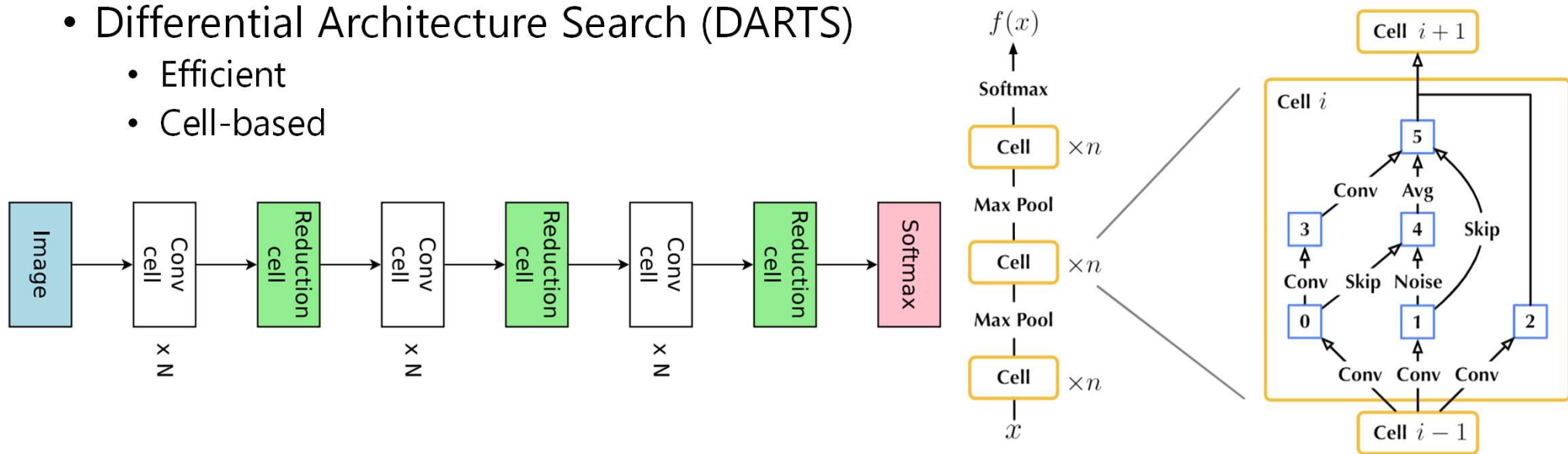


Google's AutoML



# Background

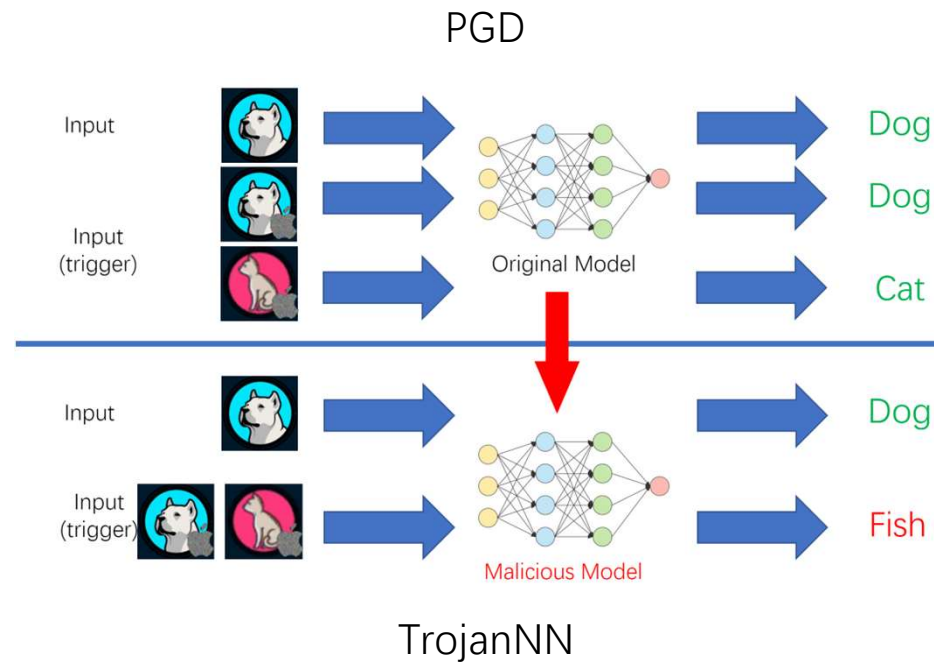
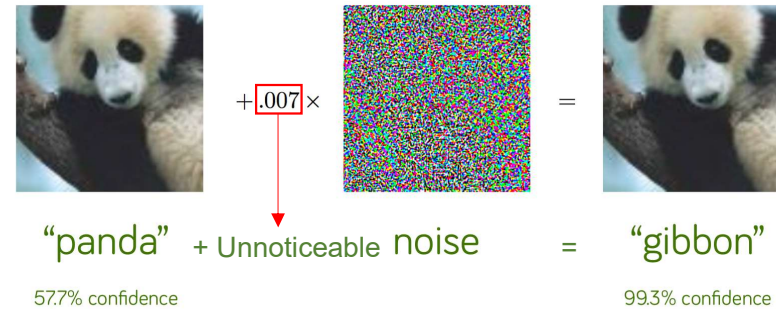
- Neural Architecture Search (NAS)
  - NAS searches good architectures automatically.
  - Differential Architecture Search (DARTS)
    - Efficient
    - Cell-based



# Background

- Attacks

- Evasion
- Data Poisoning
- Backdoor Injection
- Model Extraction
- etc.



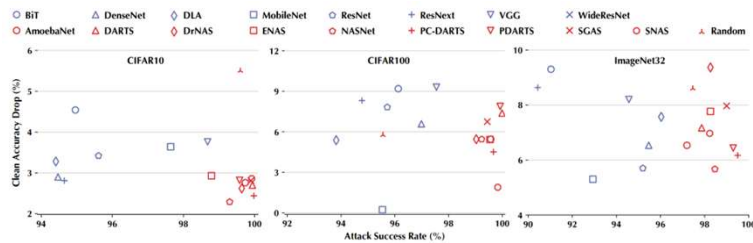
# Datasets/Models

	Architecture	CIFAR10	CIFAR100	ImageNet32
Manual Architecture	<i>BiT</i> [32]	96.6%	80.6%	72.1%
	<i>DenseNet</i> [28]	96.7%	80.7%	73.6%
	<i>DLA</i> [60]	96.5%	78.0%	70.8%
	<i>ResNet</i> [26]	96.6%	79.9%	67.1%
	<i>ResNext</i> [57]	96.7%	80.4%	67.4%
	<i>VGG</i> [52]	95.1%	73.9%	62.3%
	<i>WideResNet</i> [61]	96.8%	81.0%	73.9%
NAS Architecture	<i>AmoebaNet</i> [47]	96.9%	78.4%	74.8%
	<i>DARTS</i> [39]	97.0%	81.7%	76.6%
	<i>DrNAS</i> [11]	96.9%	80.4%	75.6%
	<i>ENAS</i> [46]	96.8%	79.1%	74.0%
	<i>NASNet</i> [64]	97.0%	78.8%	73.0%
	<i>PC-DARTS</i> [59]	96.9%	77.4%	74.7%
	<i>PDARTS</i> [12]	97.1%	81.0%	75.8%
	<i>SGAS</i> [35]	97.2%	81.2%	76.8%
	<i>SNAS</i> [58]	96.9%	79.9%	75.5%
	<i>Random</i> [17]	96.7%	78.6%	72.2%

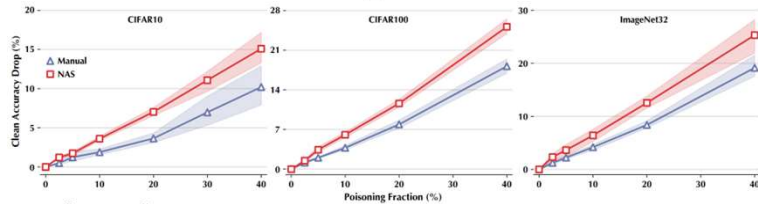
Note: ImageNet32 is a 32-class subset sampled from original ImageNet

# Vulnerabilities

- Some Experiment Results:
  - Backdoor Injection



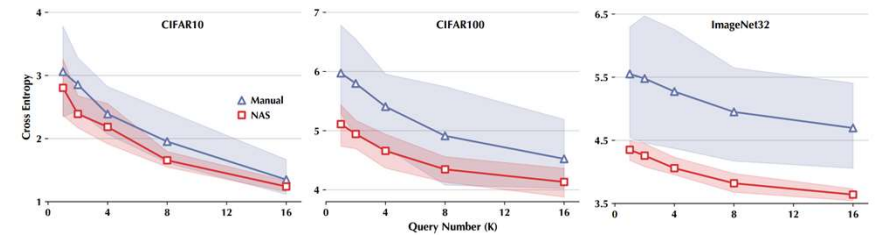
- Model Poisoning



- Conclusion

NAS-designed models tend to be more vulnerable

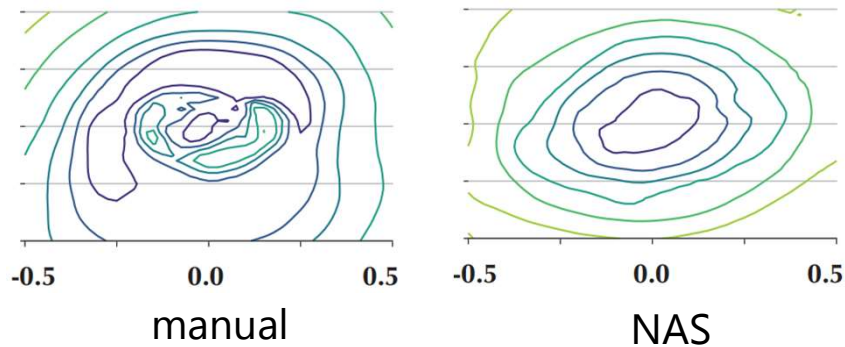
- Functional Stealing



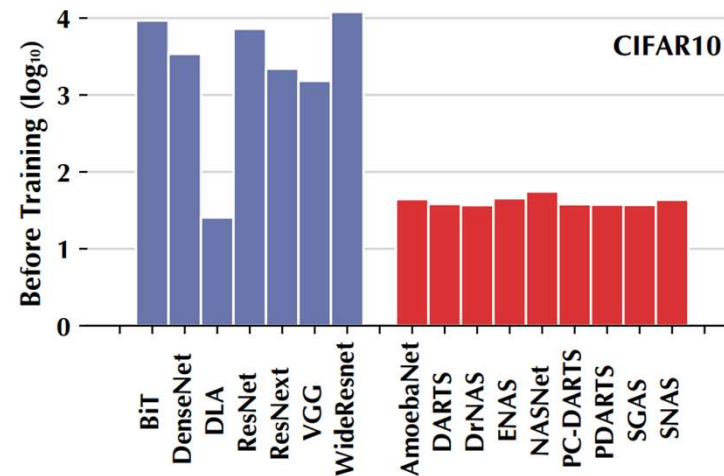
- (more results in paper)

# Analysis

- NAS algorithms prefer architectures that converge fast.
  - Shallow models
  - More skip connects
- ⇒ NAS model characteristics:
  - High Loss Smoothness (small Lipschitz constant)



- Low gradient variance





# Analysis

As a result, NAS models

- are more sensitive to training data
- gradients are more effective for optimization

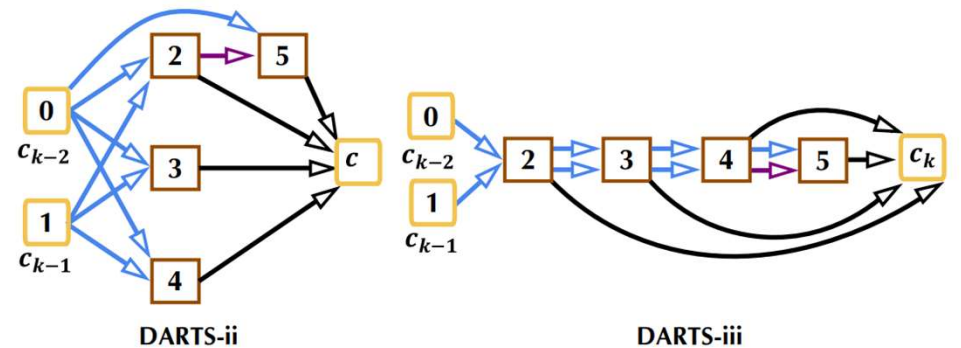
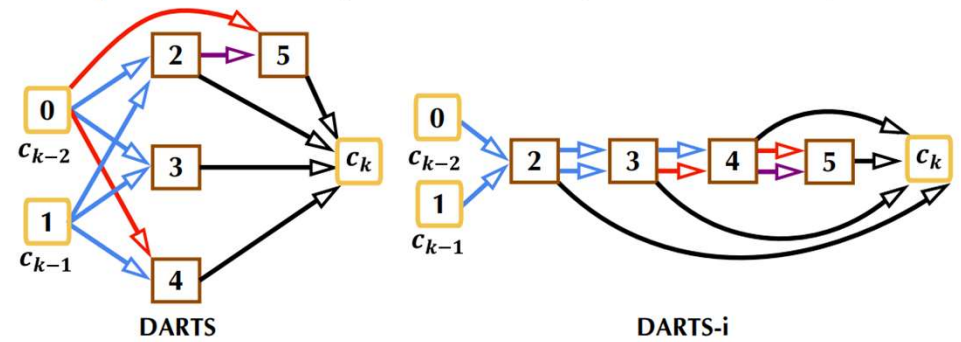
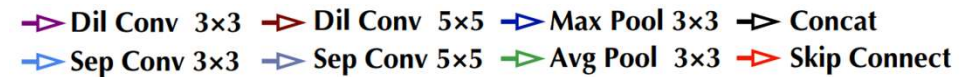
(see proof in paper)

How to understand?

e.g., 1-step PGD,  $\mathcal{L}_{NAS}$  drops more  
 $\Rightarrow$  easier to attack

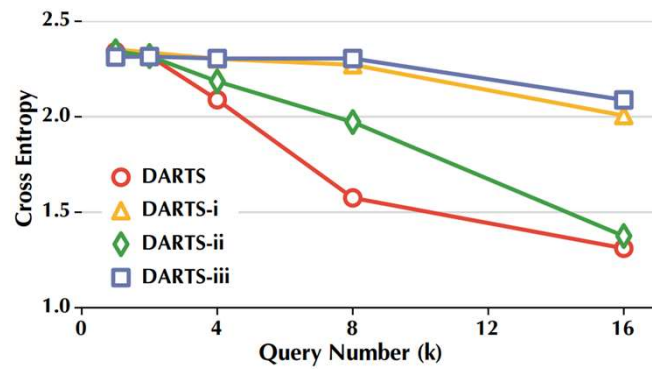
# Mitigation

- To suppress those characteristics,
  - increase cell depth
  - reduce skip connects
  - combined of (i) and (ii)

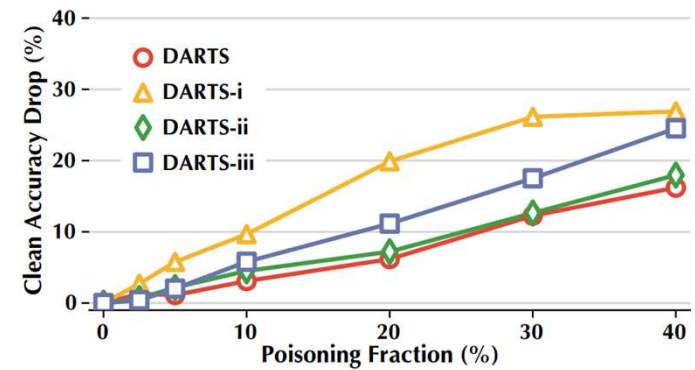


# Mitigation

- Evaluation
  - Functional Stealing



- Model Poisoning



# Conclusion

- NAS-designed models are more vulnerable against various attacks due to:
  - High loss smoothness
  - Low gradient variance
- Mitigation:
  - Building attack robustness into the NAS architectures

**Thank You!**