

# Stochastic Activation Pruning for Robust Adversarial Defense

Guneet S. Dhillon<sup>1,3</sup>, Kamyar Azizzadenesheli<sup>4</sup>, Aran Khanna<sup>1</sup>, Jeremy Bernstein<sup>1,2</sup>  
Jean Kossaifi<sup>1,5</sup>, Zachary C. Lipton<sup>1,6</sup>, Anima Anandkumar<sup>1,2</sup>

<sup>1</sup>Amazon AI, <sup>2</sup>Caltech, <sup>3</sup>UT Austin, <sup>4</sup>UC Irvine, <sup>5</sup>Imperial College London, <sup>6</sup>Carnegie Mellon University

## Adversarial Examples

Neural network:  $h$       Parameters:  $\theta$   
Input:  $x$       True output:  $y$

Adversary's objective:

$$\Delta x = \arg \max_{r \sim \rho} J(\theta, x + r, y)$$

where perturbations are **bounded** by  $l_\infty$ -norm.

Using the Taylor expansion up to the first order term,

$$\Delta x = \arg \max_{r \sim \rho} [J(\theta, x, y) + r^\top \mathcal{J}(\theta, x, y)] \quad \text{where } \mathcal{J} = \frac{\partial J}{\partial x}$$

$$\Rightarrow \Delta x = \arg \max_{r \sim \rho} r^\top \mathcal{J}(\theta, x, y)$$

The adversary chooses  $r$  to be **in the direction of**  $\mathcal{J}(\theta, x, y)$ .

So the adversary can choose  $\Delta x = \lambda \text{sign}(\mathcal{J}(\theta, x, y))$  (the **fast gradient sign method**).

## Minimax Zero-sum Game

Adversary-defender problem  
(game-theoretic perspective):

$$\pi^*, \rho^* := \arg \min_{p \sim \pi} \max_{r \sim \rho} \mathbb{E}_{\pi, \rho} [J(M_p(\theta), x + r, y)]$$

where

- **adversary** tries to **maximize** the loss by **perturbing the input** under policy  $\rho$ .
- **defender** tries to **minimize** the loss by **changing model parameters** under policy  $\pi$ .

The optimization problem is a **minimax zero-sum game** between the adversary and defender.

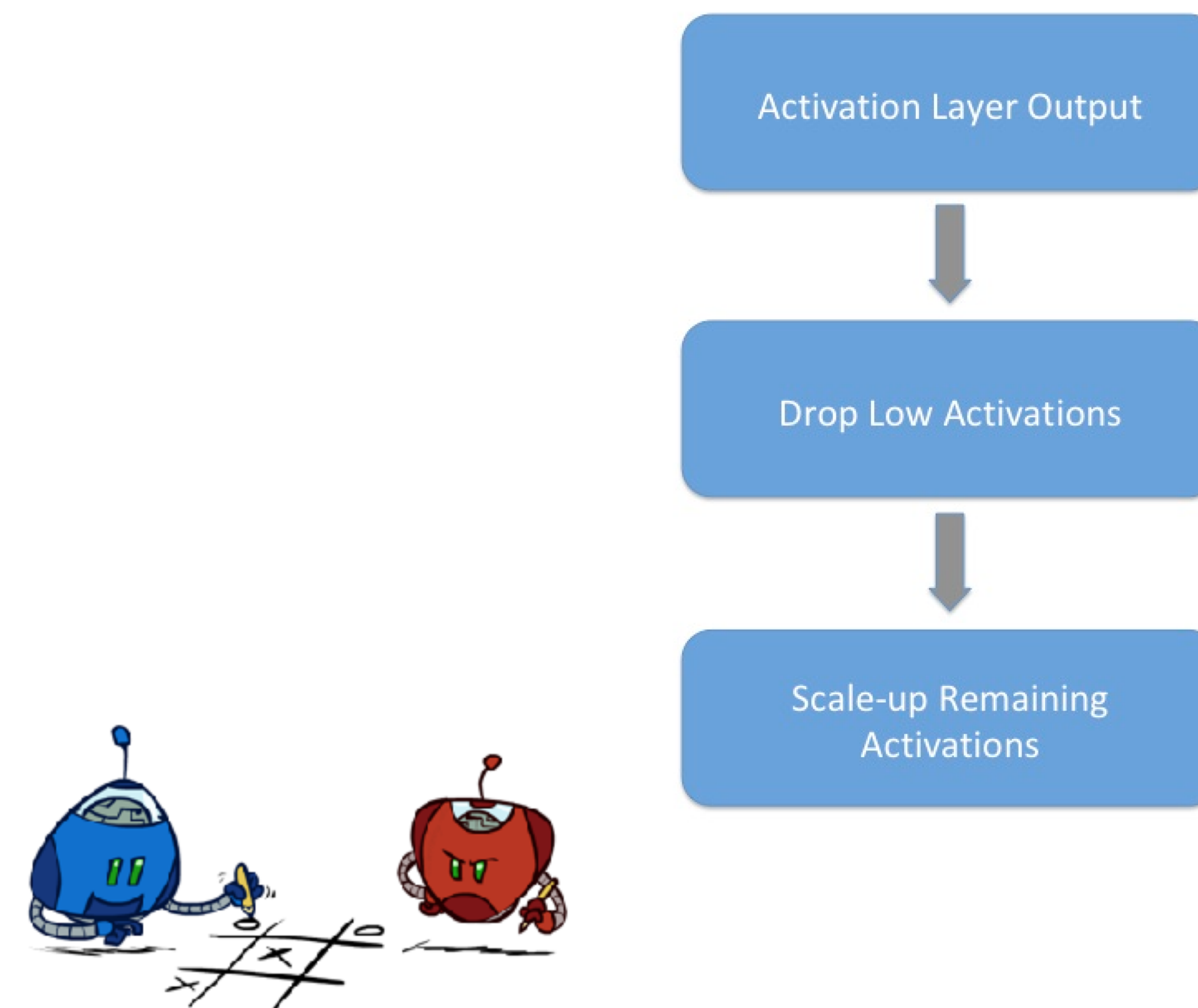
Optimal strategies  $(\pi^*, \rho^*)$ , in general, are mixed Nash equilibrium, i.e. **stochastic policies**.

## Stochastic Activation Pruning (SAP)

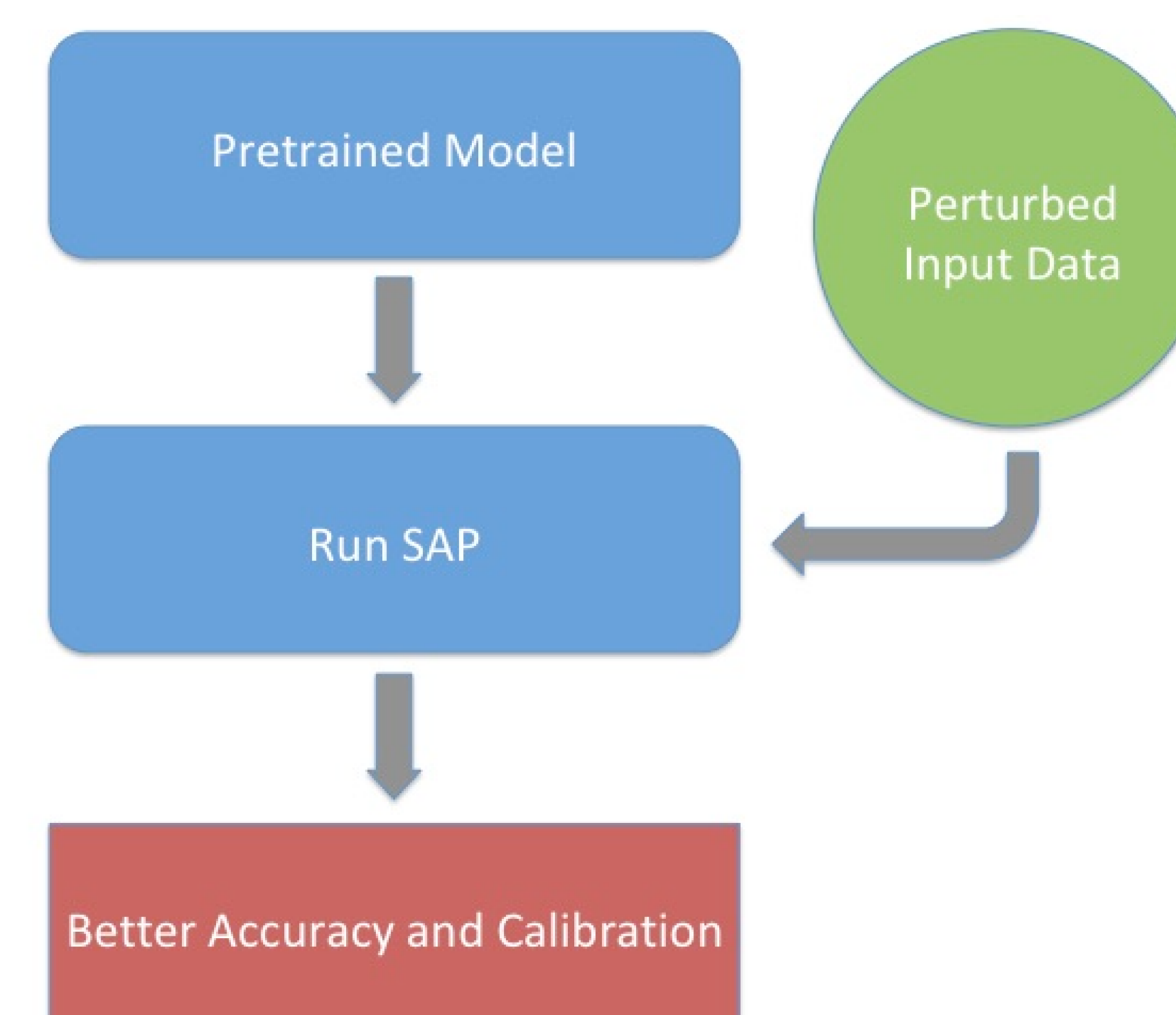
**Intuitive idea:** **stochastically drop out nodes** in each layer during forward propagation.

- **Retain nodes** with probabilities **proportional to the magnitude** of their activation.
- **Scale up the surviving nodes** to preserve the dynamic range of the activations in each layer.
- **Preserves the accuracy** of the original model.
- **Better accuracy and calibration** on perturbed input data.
- Can be **applied post-hoc** to already-trained models.

## Stochastic Activation Pruning on each Activation Layer



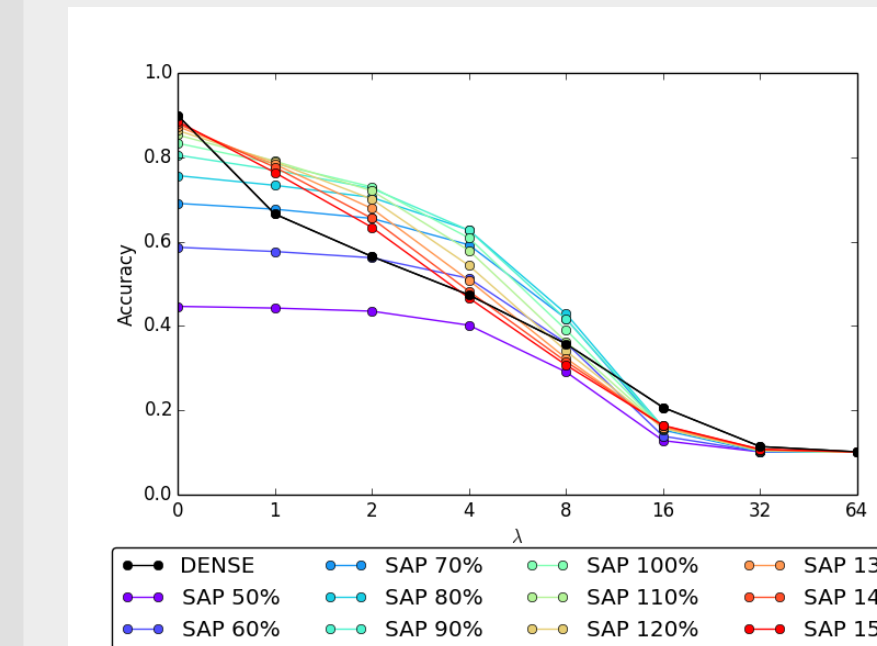
## Stochastic Activation Pruning on Pre-Trained Models



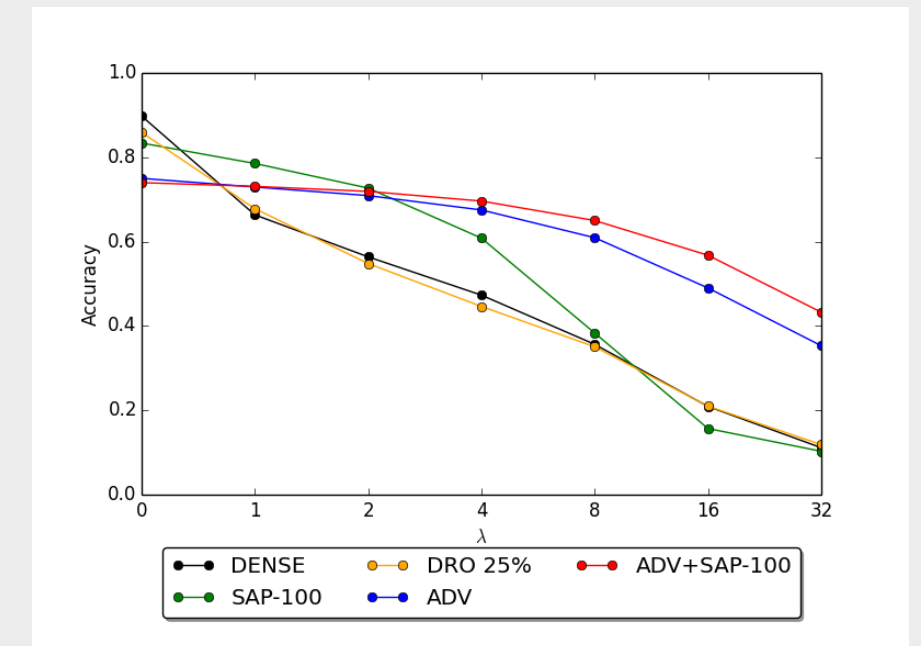
## Image Classification

**Dataset:** Cifar-10  
**Loss function:** Cross-entropy

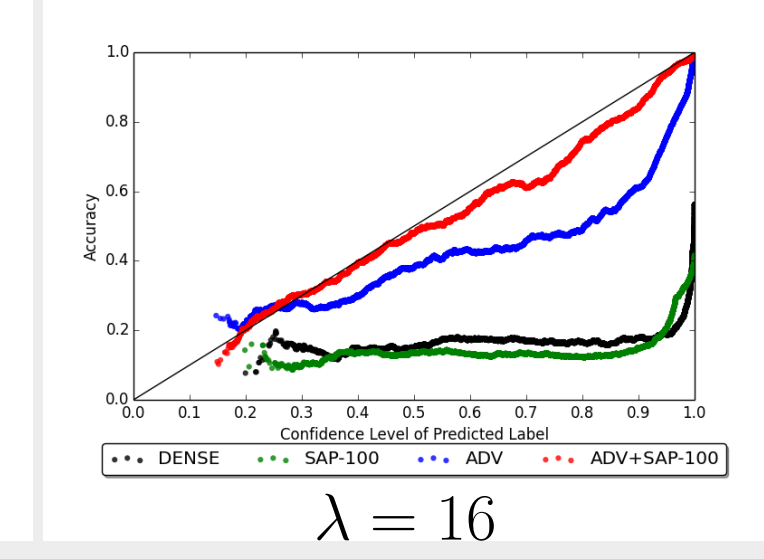
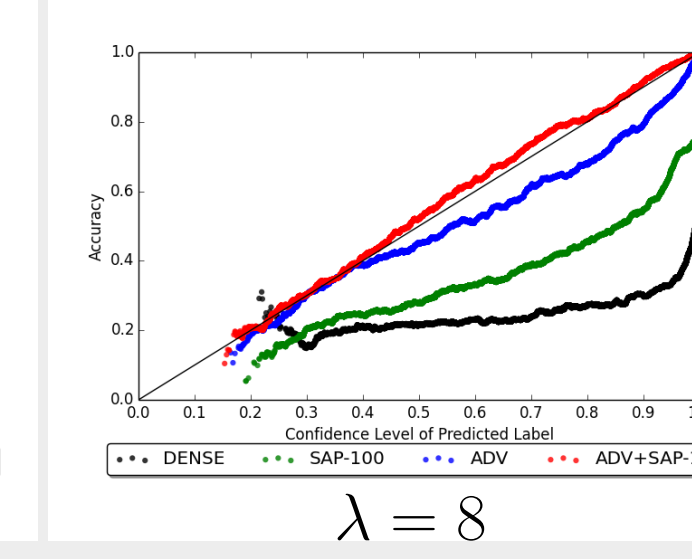
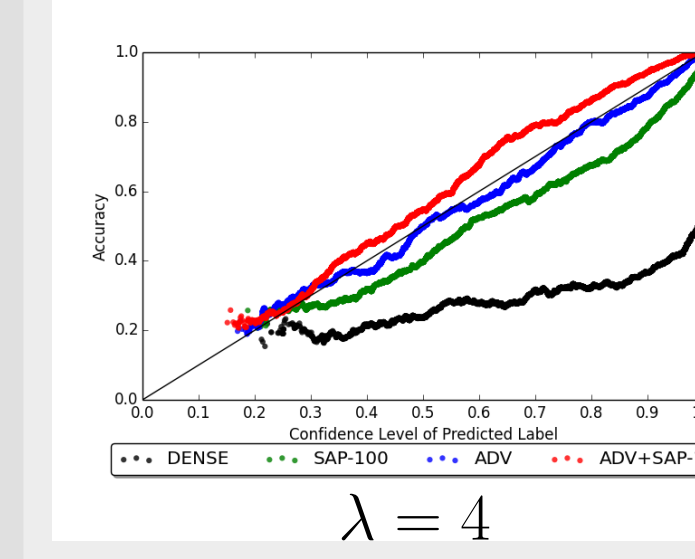
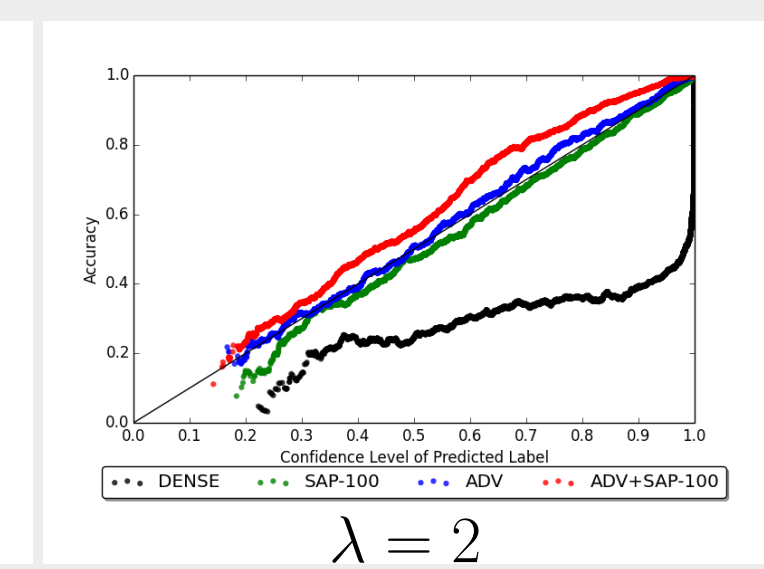
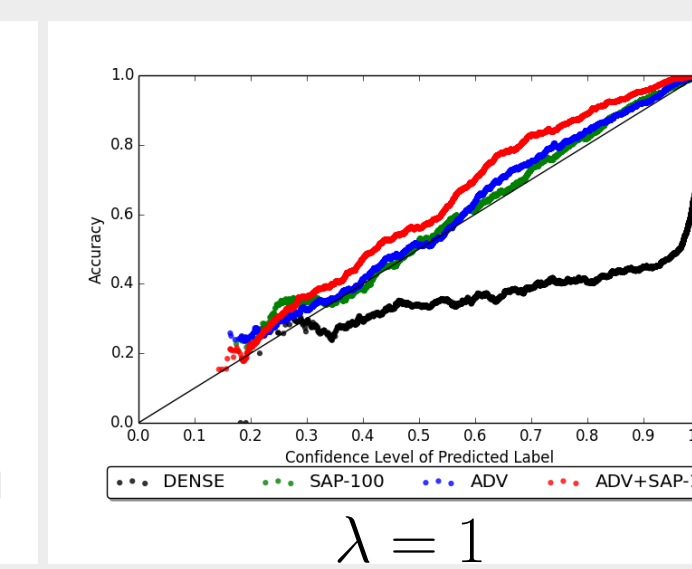
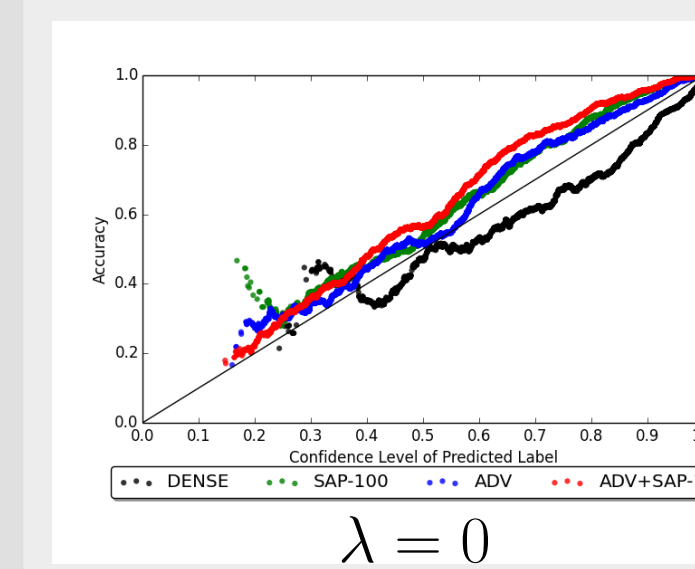
**Model:** ResNet-20  
**Non-linearity:** ReLU



Accuracy- $\lambda$  plots for SAP models. The legend indicates percentage of samples drawn.



Accuracy- $\lambda$  plots for dense, SAP-100, ADV and ADV+SAP-100 models.



Calibration plots for dense, SAP-100, ADV and ADV+SAP-100 models.

## Reinforcement Learning

**Games:** Multiple Atari games      **Model:** Double DQN

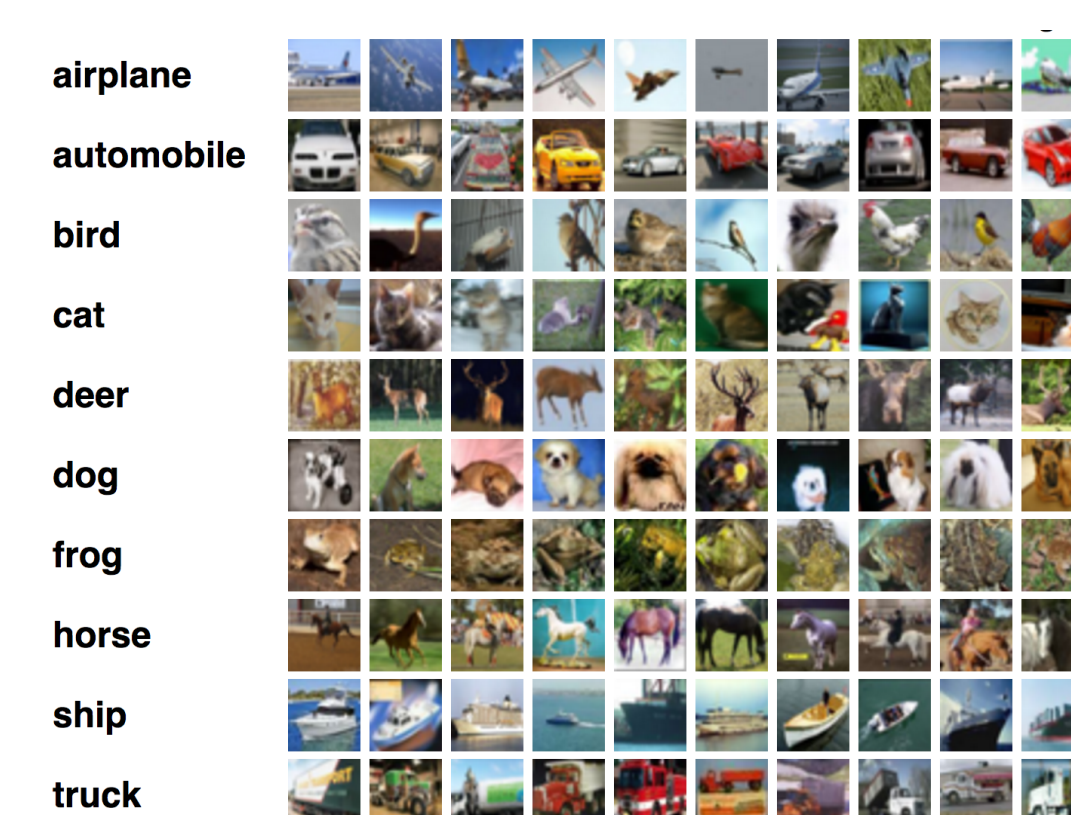
Percentage relative increase in rewards gained for SAP-100 compared to original model while playing different Atari games.

$\lambda$	Assault	Asterix	BankHeist	BattleZone	BeamRider	Bowling
0	-12.2	-33.4	-59.2	-65.8	-15.8	-4.5
1	10.4	13.3	131.7	-22.0	164.5	3425.9
2	9.8	20.8	204.8	110.1	92.3	
4	12.4	14.0	1760.0	202.6		
8	16.6	7.4	60.9	134.8		

## Discussion

- Results in improvements in the robustness of pre-trained models.
- Does not require any additional training.
- In the image classification domain, both the accuracy and calibration of the model improved.
- Can be combined with adversarial training, to compound the benefits.
- In the reinforcement learning domain, is able to defend against adversarial examples better than original model.

## Image Classification



## Reinforcement Learning

