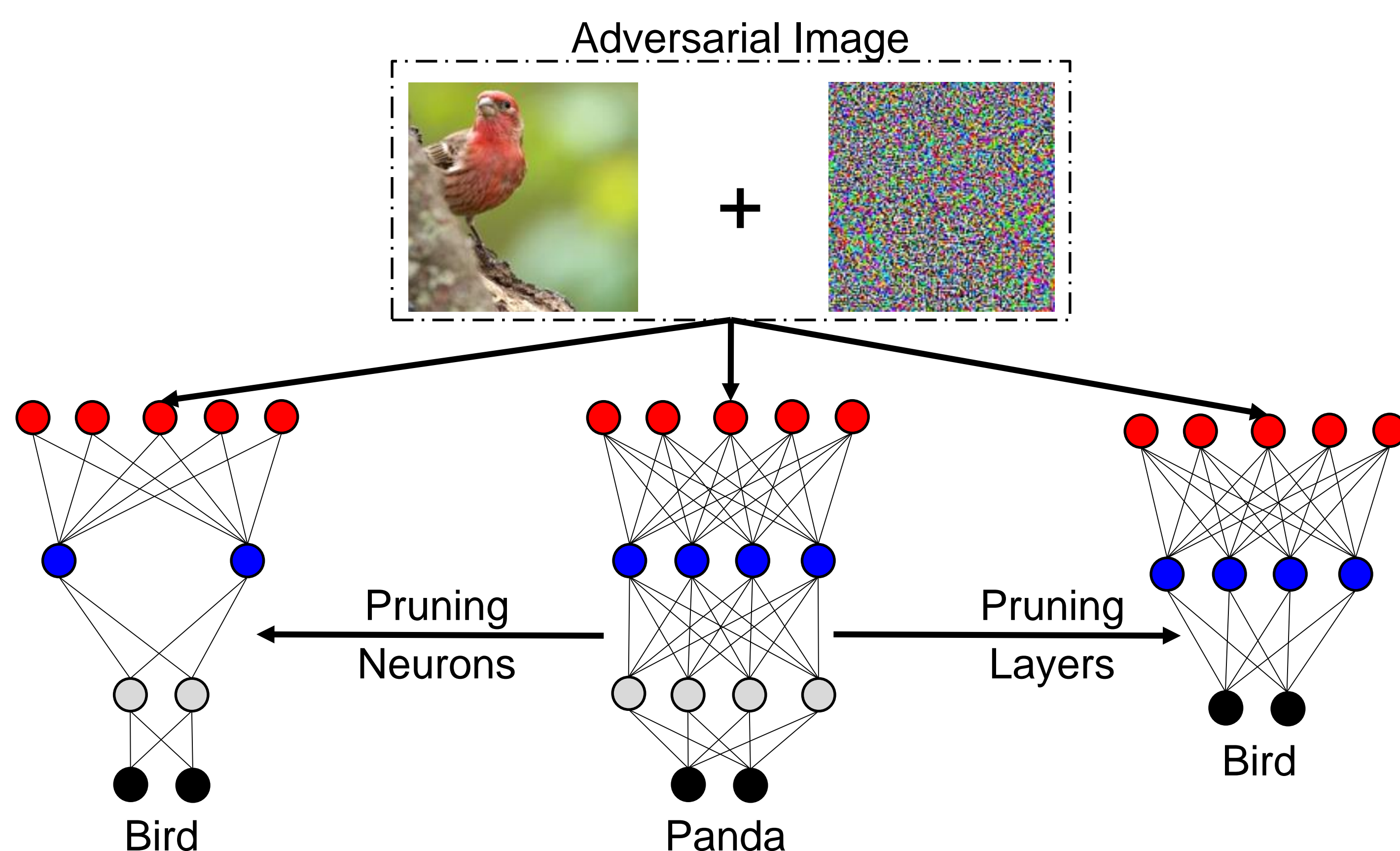


## Research Question

Do pruned ( $\mathcal{F}'$ ) networks inherit the vulnerability to adversarial images of their unpruned ( $\mathcal{F}$ ) counterpart?

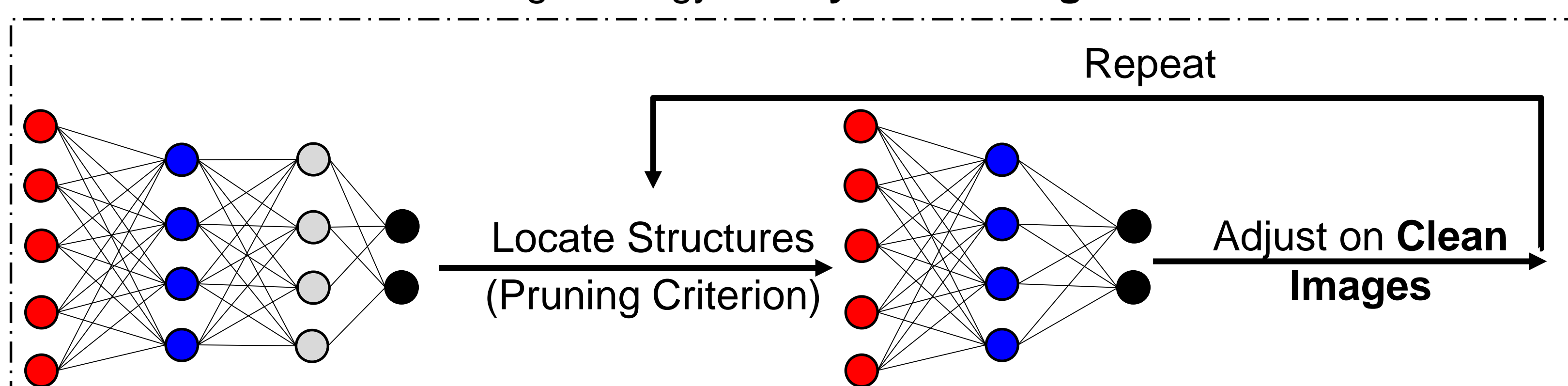
$$Acc_{adv}(\mathcal{F}') > Acc_{adv}(\mathcal{F}), Acc_{clean}(\mathcal{F}') \approx Acc_{clean}(\mathcal{F})$$

Are pruned networks capable of improving robustness while preserving generalization?



## Methodology

Pruning Strategy – Only Clean Images



## Experiments

Table 1. Robustness and Generalization from Pruning

| Architecture | Structure | Semantic        | Occlusion       | FGSM            | Clean           | Average         |
|--------------|-----------|-----------------|-----------------|-----------------|-----------------|-----------------|
| ResNet56     | Filters   | <b>(+) 1.58</b> | (+) 2.76        | <b>(+) 3.68</b> | (+) 0.60        | <b>(+) 2.15</b> |
|              | Layers    | (+) 1.05        | (+) 1.06        | (+) 3.20        | (+) 0.84        | (+) 1.53        |
|              | Both      | (-) 4.13        | <b>(+) 4.62</b> | (+) 0.36        | (-) 0.60        | (+) 0.06        |
| MobileNetV2  | Filters   | (-) 0.60        | <b>(+) 3.35</b> | (+) 0.64        | <b>(+) 0.37</b> | (+) 0.94        |
|              | Layers    | (-) 0.49        | (+) 2.12        | <b>(+) 1.44</b> | (+) 0.15        | (+) 0.80        |
|              | Both      | <b>(+) 0.07</b> | (+) 2.56        | (+) 1.05        | (+) 0.17        | <b>(+) 0.96</b> |

Table 2. Adjustment

|             | Semantic        | Occlusion       | FGSM            | Average         |
|-------------|-----------------|-----------------|-----------------|-----------------|
| Scratch-E   | <b>(+) 1.72</b> | (-) 0.60        | (+) 1.64        | (+) 0.92        |
| Scratch-B   | (+) 1.25        | (-) 0.71        | (+) 3.64        | (+) 1.39        |
| W-Ticket    | (+) 1.13        | (-) 5.41        | (+) 1.40        | (-) 0.96        |
| Fine-Tuning | (+) 1.58        | <b>(+) 2.76</b> | <b>(+) 3.68</b> | <b>(+) 2.67</b> |

Table 3. Influence of the Pruning Criterion

| Pruning Criterion | Semantic        | Occlusion       | FGSM            | Clean           | Average         |
|-------------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| $\ell_1$ -norm    | <b>(+) 1.64</b> | (-) 0.80        | (+) 3.75        | (+) 0.38        | (+) 1.24        |
| ExpectedABS       | (+) 0.96        | (-) 0.09        | <b>(+) 4.29</b> | (+) 0.51        | (+) 1.41        |
| HRank             | (+) 0.93        | <b>(+) 2.92</b> | (+) 3.18        | (+) 0.39        | (+) 1.85        |
| KIDivergence      | (+) 0.82        | (+) 0.73        | (+) 3.00        | (+) 0.34        | (+) 1.22        |
| PLS               | (+) 1.58        | (+) 2.76        | (+) 3.68        | <b>(+) 0.60</b> | <b>(+) 2.15</b> |

Table 4. Comparison with Competing Defense Mechanisms

| Defense         | Robustness      | Generalization  | Average         |
|-----------------|-----------------|-----------------|-----------------|
| Stylized        | (-) 2.29        | (-) 16.20       | (-) 9.24        |
| MixUp           | (-) 4.77        | <b>(+) 1.10</b> | (-) 1.83        |
| Cutout          | (+) 1.39        | (+) 0.75        | (+) 1.06        |
| CutMix          | (+) 1.71        | (+) 2.07        | (+) 1.89        |
| Shape-Texture   | <b>(+) 7.50</b> | (+) 0.50        | <b>(+) 4.00</b> |
| Pruning Filters | (+) 1.14        | <b>(+) 3.15</b> | <b>(+) 2.14</b> |
| Pruning Layers  | <b>(+) 1.20</b> | (+) 3.03        | (+) 2.11        |

## Conclusion

- We empirically show that pruning structures of convolutional networks increase their adversarial robustness
- We demonstrate that pruning preserves generalization; thus, it efficiently satisfies the dilemma between robustness and generalization
- We confirm these findings considering **only clean images** during the pruning process, which enables us to design an effective defense mechanism that ignores the settings and additional assumptions of the attack