

Deep Learning SP25 : Project 03

Rujuta Amit Joshi¹, Lavanya Deole², Sarang Kadakia³

¹New York University

²New York University

³New York University

rj2719@nyu.edu, lnd2037@nyu.edu, sk11634@nyu.edu

GitHub Project Repository: [DLProject03](#)

Abstract

Despite their remarkable success across computer vision tasks, deep convolutional neural networks (CNNs) are highly susceptible to adversarial examples, carefully perturbed inputs that lead to misclassification while remaining visually indistinguishable from original images. This report explores three adversarial attacks: Fast Gradient Sign Method (FGSM), Projected Gradient Descent (PGD), and an iterative Patch Attack against a pre-trained ResNet-34 model on a 500-image subset of ImageNet. We evaluate the degradation in Top-1 and Top-5 classification accuracy caused by each attack and visualize sample perturbations. We further assess the transferability of adversarial examples to DenseNet-121. Our results show severe drops in performance, with PGD and patch attacks reducing accuracy to near zero, emphasizing the fragility of CNNs and the need for adversarial robustness.

Introduction

Modern deep learning systems, particularly CNNs, have demonstrated state-of-the-art performance in object recognition, medical imaging, and autonomous driving. However, their reliability in adversarial settings has been called into question. Small, often imperceptible input perturbations can drastically alter a model's predictions, leading to potentially harmful consequences in safety-critical applications.

In this project, we evaluated the adversarial robustness of a pre-trained ResNet-34 model using three common attack strategies: FGSM, PGD, and an iterative Patch Attack. The goal was to quantify accuracy degradation under adversarial perturbations and analyze perceptual and structural impacts. We also assessed how adversarial examples transfer to DenseNet-121. Our findings showed that PGD and Patch attacks can reduce Top-1 accuracy to near-zero levels, while even fast attacks like FGSM cause significant misclassification. The adversarial examples remained mostly imperceptible, and a transferability gap was observed across architectures.

1.1. Problem Statement

Deep neural networks, despite their high accuracy on benchmark datasets like ImageNet, are vulnerable to adversarial inputs; small, intentionally crafted perturbations that lead to incorrect predictions. This project aims to study the behavior of such attacks on a popular CNN architecture and understand their impact on performance and model reliability.

Specifically, the objectives are to:

- Evaluate the classification performance of a pre-trained **ResNet-34** model on clean image data.
- Implement three types of adversarial attacks:
 - **FGSM** (Fast Gradient Sign Method)
 - **PGD** (Projected Gradient Descent)
 - **Iterative Patch Attack**
- Measure the drop in **Top-1** and **Top-5** accuracy for each attack.
- Visualize and compare original vs adversarial images to assess perceptibility.
- Assess the **transferability** of adversarial examples to **DenseNet-121**.
- Generate plots and structured results to analyze trends and performance degradation.

1.2. Motivations and Challenges

The motivation behind this project stems from the growing concern over the security and reliability of deep learning models in real-world applications. Adversarial examples, though visually indistinguishable from clean inputs, can cause severe misclassifications: posing risks in domains like autonomous driving, healthcare, and surveillance. The challenge lies in crafting these perturbations effectively while ensuring imperceptibility, implementing attacks in raw pixel space with correct normalization, and handling gradient flow and tensor operations during iterative updates. Additionally, evaluating cross-model transferability and visual-

izing subtle perturbations at scale requires careful engineering and interpretation. This project addresses these challenges while offering insights into model fragility and the need for robust defenses.

Methodology

This section outlines the experimental setup, dataset specifications, model architecture, attack implementations, and evaluation metrics used to assess the robustness of the image classification models.

2.1. Dataset

We use a curated subset of the ImageNet dataset consisting of 500 RGB images evenly sampled across a wide range of classes. The dataset is organized in a class-wise directory structure compatible with PyTorch’s ImageFolder. Each image is resized and normalized using the standard ImageNet mean and standard deviation:

```
mean = [0.485, 0.456, 0.406]
std = [0.229, 0.224, 0.225]
```

The images are loaded using a DataLoader with a batch size of 32 and are evaluated in inference mode on GPU. A JSON label mapping file `labels_list.json` was used to map folder-wise labels to corresponding ImageNet class IDs to ensure compatibility with pre-trained model outputs.

2.2. Models

Two convolutional neural network (CNN) architectures are used in this study:

- **ResNet-34:** A deep residual network with 34 layers, used as the primary model to generate and evaluate adversarial examples.
- **DenseNet-121:** A densely connected CNN used to evaluate the transferability of adversarial samples crafted for ResNet-34.

Both models are loaded with pre-trained weights on ImageNet (weights="IMAGENET1K_V1") and are used in evaluation mode without any fine-tuning.

2.3. Adversarial Attacks

The following adversarial methods were implemented and executed in raw pixel space, followed by re-normalization before feeding inputs to the model.

2.3.1. FGSM (Fast Gradient Sign Method)

The FGSM is a single-step, gradient-based attack that perturbs the input image in the direction of the loss gradient. The adversarial image is generated as:

$$x_{\{adv\}} = x + \epsilon \cdot \text{sign}(\nabla_x J(x, y))$$

- x is the input image,
- ϵ is the attack strength (set to 0.02),
- J is the loss function,
- y is the true label.

This method is computationally efficient but less precise compared to iterative methods.

2.3.2. PGD (Projected Gradient Descent)

PGD is an iterative, stronger variant of FGSM that applies gradient updates multiple times and projects the result back into an ϵ - ball around the original image. The update rule is:

$$x_{\{t+1\}} = \text{clip}_{\{x, \epsilon\}} \left(x_t + \alpha \cdot \text{sign}(\nabla_x J(x_t, y)) \right)$$

- $\epsilon = 0.02$
- Step size $\alpha = 0.0013$
- Number of steps = 15

This attack is more effective at fooling classifiers due to its iterative refinement.

2.3.3 Iterative Patch Attack

This custom attack modifies only a localized rectangular patch of the input image:

- Patch size: 48×48 pixels
- Total steps: 50
- Every 10 steps, the patch is randomly repositioned
- Only the patch region receives gradient updates

Unlike full-image perturbations, patch attacks simulate real-world occlusions or image corruptions.

2.4. Hyperparameter Design and Lessons Learned

We carefully tuned attack parameters to balance attack strength and visual imperceptibility. For FGSM and PGD, ϵ was set to 0.02 based on prior literature and visual tests. PGD used 15 steps with $\alpha = 0.0013$ for effective convergence within the ϵ -ball. For the Patch Attack, a patch size of 48×48 was chosen to simulate localized corruption, with 50 update steps and patch relocation every 10 steps. One challenge was ensuring perturbations stayed within valid pixel ranges after denormalization and re-normalization, especially when applying in raw space. Iterative attacks required gradient retention and re-assigning requires_grad, which led to multiple debugging iterations. Visualization helped identify if perturbations were too weak or too visible, guiding our choice of ϵ .

2.5. Evaluation Metrics

We evaluate model robustness using the following metrics:

- **Top-1 Accuracy:** Proportion of images where the model's highest confidence prediction matches the true class.
- **Top-5 Accuracy:** Proportion of images where the true class is within the model's top 5 predictions.

For all evaluations, 500 images are used. Additionally, visual and quantitative comparisons are made between clean and adversarial predictions.

Experimental Results

Adversarial attacks significantly degrade model performance. PGD and patch attacks were especially destructive, dropping Top-1 accuracy to near zero.

Table 1. Accuracy Drop on ResNet-34

Attack	ϵ	Top-1 Accuracy (%)	Top-5 Accuracy (%)
Clean	0	76.00	94.20
FGSM	0.02	3.40	21.20
PGD	0.02	0.00	1.40
Patch	0.3	0.20	14.80
Patch	0.5	0.00	13.00

Table 2. Transferability to DenseNet-121

Attack Generated from ResNet	ϵ	Top-1 Accuracy (%)	Top-5 Accuracy (%)
FGSM	0.02	45.60	76.20
PGD	0.02	35.60	72.40
Patch	0.3	58.80	87.60
Patch	0.5	57.40	85.40

These results indicate that adversarial perturbations generated for ResNet-34 can partially transfer to another architecture, confirming the cross-model generalizability of adversarial vulnerabilities.

3.1. Perturbation Sizes & Timing:

- **FGSM / PGD** $\epsilon = 0.02$ (global pixel update)
- **Patch Attack** $\epsilon = 0.3$ and 0.5 (48×48 region only)
- **Training/Generation Times (on Colab GPU):**
 - FGSM: ~ 4s
 - PGD: ~ 38s
 - Patch Attack: ~ 2 min 15s (per ϵ value)

These timings reflect the increased computational cost of iterative and spatially constrained attacks.

Visualizations

To illustrate perceptibility and structural differences between clean and adversarial inputs, we visualized five samples per attack. For each sample, the following were displayed:

- Original image
- Adversarial image
- Difference heatmap (amplified ×5 for visibility)



Fig. 1. FGSM $\epsilon=0.02$

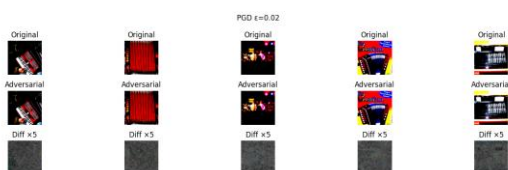


Fig. 2. PGD $\epsilon=0.02$



Fig. 3. Patch $\epsilon=0.3$



Fig. 4. Patch $\epsilon=0.5$

Comparative Plots and Tables

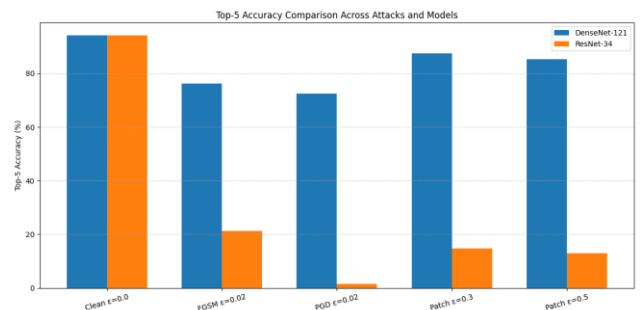


Fig. 5. final_top5_accuracy_barplot (Shows Top-5 accuracy per model and attack)

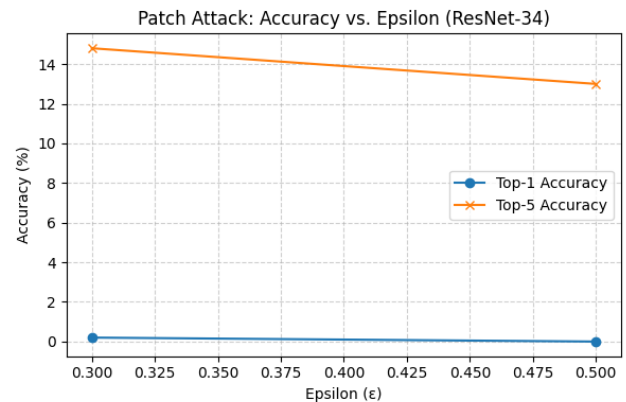


Fig. 6. patch_accuracy_vs_epsilon (Shows how accuracy decreases with larger ϵ)

- **FGSM and PGD** adversarial images appear visually identical to the original images, confirming the imperceptibility of perturbations.
- **Difference heatmaps** (×5 amplified) show subtle noise patterns that significantly impact model predictions despite being invisible to the human eye.
- **Patch Attack** introduces a small but visible region in the image; even localized perturbations lead to a sharp drop in accuracy.
- Model attention visibly shifts in adversarial cases, often misfocusing on irrelevant or altered areas.
- **Transferability effects** are visually evident, adversarial images mislead both ResNet-34 and DenseNet-121, especially in FGSM and PGD cases.

Discussion

The adversarial attacks implemented: FGSM, PGD, and Patch led to a significant drop in classification performance on ResNet-34. The most severe drop was observed in the PGD attack (Top-1: 0%), followed closely by the Patch attack (Top-1: 0.2% for $\epsilon=0.3$ and 0% for $\epsilon=0.5$), confirming the effectiveness of iterative methods. FGSM, though simpler, also caused a drastic reduction in accuracy (Top-1: 3.4%).

Transferability results showed that adversarial examples crafted on ResNet-34 remained partially effective when evaluated on DenseNet-121. FGSM and PGD transferred more successfully (Top-1 drop of 30 - 40%), whereas Patch attacks were less transferable but still impactful.

Visualizations supported these trends: FGSM and PGD perturbations were imperceptible, while Patch attacks introduced a visible region yet retained high fooling power. Difference heatmaps further revealed non-random, structured perturbations that align with the model's most sensitive features.

Overall, the results demonstrate that deep CNNs are highly susceptible to both global and localized attacks, and that adversarial vulnerability extends across model architectures. These findings highlight the importance of integrating robustness evaluations into model development pipelines.

5.1. Result Image Sets (Description)

As per the project requirements, 500 adversarial images were generated and saved for each attack in raw .png format:

Folder	Attack Type	Epsilon Used	Notes
Adversarial-TestSet1/	FGSM	$\epsilon = 0.02$	Fast one-step attack, subtle changes
Adversarial-TestSet2/	PGD	$\epsilon = 0.02$	Strong iterative attack
Adversarial-TestSet3/	Patch (iterative)	$\epsilon = 0.5$	Localized visible patch perturbations

5.2. Lessons Learned & Transferability Mitigation

Key Findings & Trends:

- PGD and Patch attacks caused the most drastic accuracy drops, with Top-1 accuracy falling to near 0%.
- FGSM, though faster, was less impactful but still significantly reduced performance.
- Transferability was evident, particularly for FGSM and PGD, indicating cross-model vulnerability.

Lessons Learned:

- Raw-space attacks require careful handling of normalization and gradient flow.
- Subtle perturbations can lead to major misclassifications, even if visually imperceptible.
- Visualization (e.g., diff maps) helped in validating and debugging attacks.

Mitigation Strategies:

- Use adversarial training to build resilience.
- Apply input preprocessing (e.g., JPEG compression, smoothing).
- Combine models in ensembles to reduce susceptibility.
- Limit gradient exploitability via regularization or smoothing.

Conclusion

This project illustrates the vulnerability of deep image classifiers to adversarial manipulation. Through FGSM, PGD, and Patch attacks, we demonstrated significant accuracy degradation on ResNet-34 and partial transferability to DenseNet-121. The experiments validated the effectiveness of both full-image and localized attacks, reinforcing the concern that current models prioritize sensitivity over semantic understanding.

Adversarial visualization and evaluation metrics together reveal the fragility of state-of-the-art models and highlight the need for incorporating robustness at the core of model design and deployment strategies.

Future Work

To extend and strengthen this work, the following directions are recommended:

- Adversarial Defenses:** Incorporate adversarial training or input preprocessing techniques (e.g., JPEG compression, randomized smoothing).
- Broader Attack Evaluation:** Implement other attack strategies like CW, DeepFool, or AutoAttack for more comprehensive analysis.
- Saliency and Attention Analysis:** Study how adversarial inputs affect attention maps or class activation regions.
- Multi-Model Transfer:** Test transferability across more diverse architectures (e.g., ViT, Inception) and ensemble models.
- Real-Time Robustness:** Evaluate latency and robustness trade-offs under time-constrained inference, simulating real-world constraints.

Code Repository

The complete implementation, including code and training details, is available in our GitHub repository: [DLProject03](#).

References

- [1] Goodfellow, I. J., Shlens, J., & Szegedy, C. (2015), *Explaining and Harnessing Adversarial Examples*, *arXiv preprint arXiv:1412.6572*.
- [2] Madry, A., et al. (2018), *Towards Deep Learning Models Resistant to Adversarial Attacks*, International Conference on Learning Representations (ICLR).