

# Supplementary Material: Strengthening the Transferability of Adversarial Examples Using Advanced Looking Ahead and Self-CutMix

Donggon Jang\* Sanghyeok Son\* Dae-Shik Kim  
Korea Advanced Institute of Science and Technology (KAIST)  
{jdg900, ssh816, daeshik}@kaist.ac.kr

## 1. Implementation Details

We used Tensorflow to implement the attack methods and tested them on TITAN RTX and TITAN V. Experiments were mostly conducted with three random seeds. We follow the setting in MI-FGSM with the number of iterations  $T = 10$ , maximum perturbation  $\epsilon = 16$ , and step size  $\alpha = 1.6$ . For M(N)I-FGSM, we adopt the decay factor  $\mu = 1.0$ . For VM(N)I-FGSM, we choose the upper bound factor of neighborhood  $\beta = 1.5$  and the number of samples for variance tuning  $v = 20$ . For EMI-FGSM, we choose the number of samples  $N = 11$  and sampling interval bound  $\eta = 7$  and adopt the linear sampling. For DIM and TIM, we set the transformation probability to 0.5 and the filter size of the Gaussian kernel to  $7 \times 7$ , respectively. For SIM, we choose the number of scale copies  $s = 5$ . For Admix, we set the number of randomly samples images from other categories as 3 and the portion of the original image as 0.2, respectively. For our methods, we adopt the number of lookahead steps  $N = 17$ , the number of neighbor samples  $v = 5$ , the minimum patch size  $P = 200$ , and the number of mixed copies  $c = 4$ .

## 2. Results of Combining SCM with Other Input Transformations

We conduct additional experiments to show that SCM-P(R) achieve higher attack success rates even when combined with existing input transformation methods, *i.e.* DIM, TIM, and CTM. As shown in Tab. 1, Tab. 2, and Tab. 3, SCM-P(R) further improve the transferability of adversarial examples. Therefore, SCM-P(R) could be an attractive option to generate more transferable adversaries.

## 3. Results of Combining LI-FGSM with Other Input Transformations

In order to validate the effect of LI-FGSM when combined with other input transformations, we conduct additional experiments. Tab. 4, Tab. 5, Tab. 6, and Tab. 7 show

\*These authors contributed equally.

the results comparing various gradient based attack methods using input transformations, *i.e.* DIM, TIM, SIM, and CTM. LI-FGSM achieves the highest attack success rates in all experiments.

## 4. Ablation Study on Hyper-parameters

### 4.1. SCM-P(R)

We conduct additional experiments to study the influence of a copied (pasted) patch size in SCM-P(R). We measure the attack success rates against four normally trained models and three adversarially trained models by varying the patch size. As shown in Fig. 1, the transferability tends to improve with increasing patch size in both SCM-P and SCM-R. Specifically, for SCM-R, the attack success rates increase steeply after patch size is 200. SCM-P also shows satisfactory consistent transferability after patch size is 200. Thus, we simply set the minimum patch size to 200 in both SCM-P and SCM-R.

### 4.2. LI-FGSM

We experiment with varying the number of steps looking ahead according to an inner loop algorithm. We utilize MI-FGSM, NI-FGSM, and VMI-FGSM as an inner loop except for VNI-FGSM shown in the discussion section of the main paper. As shown in Fig. 2, the attack success rate grows rapidly as the number of steps increases and then converges. However, the converged attack success rates are worse than the case of using VNI-FGSM. Consequently, we leverage VNI-FGSM as an inner loop.

## 5. Training Budget

### 5.1. Computational Complexity

It is not easy to directly compare time and memory complexity since it depends on the source model. Instead, we count the number of back-propagation at each iteration that is highly correlated to the attack method and dominantly affects computation overhead. As shown in the Tab. 8, the complexity of LI-FGSM is higher than others. However, in

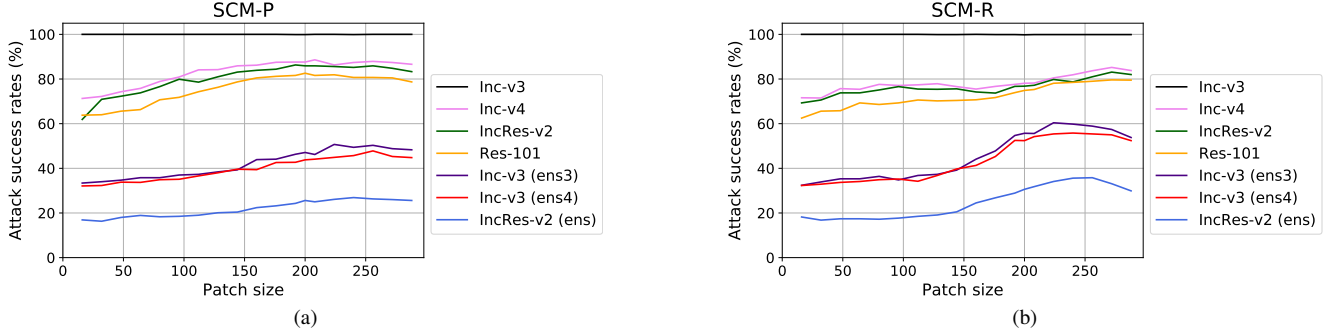


Figure 1. (a): Attack success rates (%) according to varying the patch sizes of SCM-P against four normally trained models and three adversarially trained models. (b): Attack success rates (%) according to varying the patch sizes of SCM-R against four normally trained models and three adversarially trained models. In this experiment, the adversarial examples are generated by Inc-v3 using MI-FGSM and SCM-P(R).

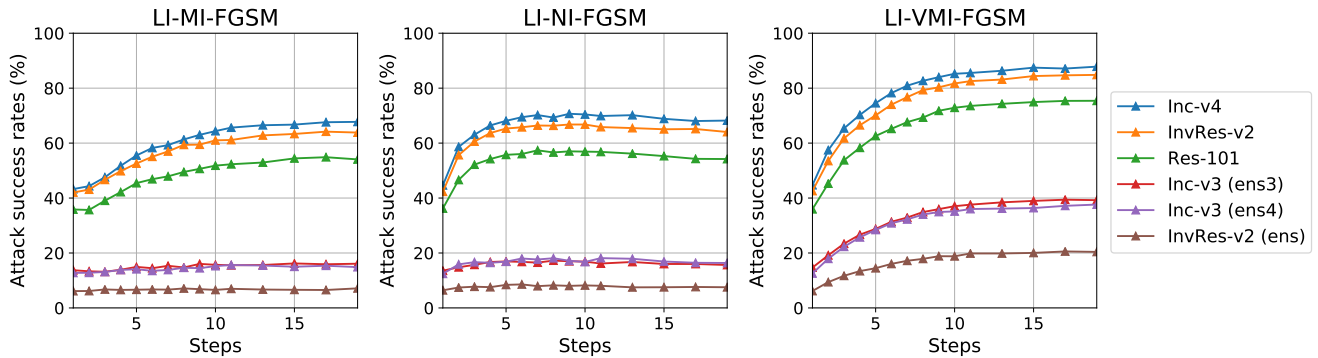


Figure 2. Attack success rates (%) according to varying the number of steps looking ahead on six models by utilizing MI-FGSM, NI-FGSM, and VMI-FGSM as an inner loop of LI-FGSM, except Inc-v3, which is the source model.

the next subsection, we exhibit that our method can outperform others with small looking ahead steps.

## 5.2. Attack Scenario under a Small Step Budget

We conduct additional experiments to demonstrate whether our methods can be helpful under a small step budget. We measure the performance of LI-VNI-FGSM varying the number of lookahead steps, and compare it with the existing methods as shown in Fig. 3. Our method outperforms the MI-FGSM and NI-FGSM from two steps (one additional step). Furthermore, ours achieves better performance beyond the EMI-FGSM and the VNI-FGSM from four steps (three additional steps). We believe that it is affordable overhead to improve the attack success rate. It is noteworthy that as the number of lookahead steps increases, the performance of our method also increases even after four steps. It means that if enough computing power is given, it is appropriate to use our method to have a higher attack success rates.

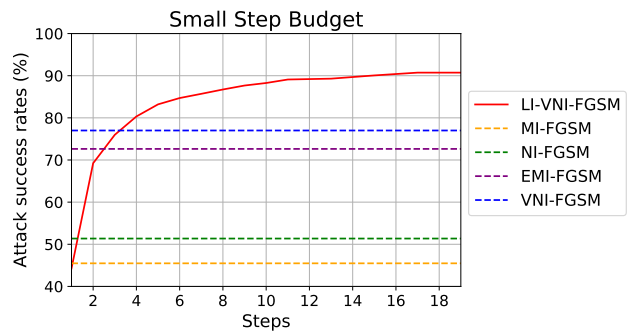


Figure 3. Attack success rates (%) according to varying the number of steps looking ahead on Inc-v4. Inc-v3 is used for the source model. We only take MI-FGSM, NI-FGSM, EMI-FGSM, and VNI-FGSM except for VMI-FGSM to compare with LI-VNI-FGSM because VNI-FGSM and VMI-FGSM show very similar performance.

| Model     | Attack    | Inc-v3        | Inc-v4        | IncRes-v2     | Res-101       | Inc-v3 (ens3) | Inc-v3 (ens4) | IncRes-v2 (ens) |
|-----------|-----------|---------------|---------------|---------------|---------------|---------------|---------------|-----------------|
| Inc-v3    | SI-DIM    | 99.13*        | 83.97         | 81.17         | 77.50         | 47.00         | 45.03         | 25.90           |
|           | Admix-DIM | <b>99.80*</b> | 90.50         | <b>87.70</b>  | 83.60         | 52.20         | 49.90         | 28.60           |
|           | SCM-P-DIM | 98.97*        | <b>89.57</b>  | 81.17         | 84.53         | 52.40         | 48.70         | 28.90           |
|           | SCM-R-DIM | 99.00*        | 88.10         | 86.53         | <b>84.97</b>  | <b>64.53</b>  | <b>62.50</b>  | <b>39.77</b>    |
| Inc-v4    | SI-DIM    | 89.27         | 99.13*        | 85.57         | 80.00         | 59.50         | 56.80         | 38.23           |
|           | Admix-DIM | <b>93.00</b>  | <b>99.20*</b> | <b>89.70</b>  | 85.20         | 62.40         | 60.30         | 39.70           |
|           | SCM-P-DIM | 91.83         | 98.77*        | 88.67         | 85.00         | 63.13         | 60.40         | 41.20           |
|           | SCM-R-DIM | 90.70         | 98.70*        | 88.47         | <b>85.90</b>  | <b>72.03</b>  | <b>70.50</b>  | <b>53.57</b>    |
| IncRes-v2 | SI-DIM    | 88.80         | 85.80         | 97.83*        | 82.43         | 66.80         | 60.93         | 53.93           |
|           | Admix-DIM | 90.20         | <b>88.40</b>  | <b>98.00*</b> | <b>85.80</b>  | 70.50         | 63.70         | 55.30           |
|           | SCM-P-DIM | <b>90.40</b>  | 88.03         | 97.63*        | 85.43         | 69.70         | 64.43         | 54.13           |
|           | SCM-R-DIM | 89.27         | 87.17         | 97.17*        | 85.63         | <b>75.20</b>  | <b>72.17</b>  | <b>63.43</b>    |
| Res-101   | SI-DIM    | 87.60         | 83.33         | 84.77         | 98.93*        | 62.30         | 56.23         | 40.40           |
|           | Admix-DIM | <b>91.90</b>  | 89.00         | <b>89.60</b>  | <b>99.80*</b> | 69.70         | 62.30         | 46.60           |
|           | SCM-P-DIM | 91.73         | <b>89.30</b>  | 89.07         | 99.00*        | 70.07         | 64.30         | 46.27           |
|           | SCM-R-DIM | 90.33         | 87.53         | 88.70         | 99.03*        | <b>76.67</b>  | <b>73.80</b>  | <b>57.47</b>    |

Table 1. Attack success rates (%) against four normally trained models and three adversarially trained models under a single model setting using SI-DIM, Admix-DIM, SCM-P-DIM, and SCM-R-DIM. Inc-v3, Inc-v4, IncRes-v2, and Res-101 using MI-FGSM are adopted as the source model respectively. \* indicates that the target model is the same as the source model. In case of Admix-DIM, we take the reported value.

| Model     | Attack    | Inc-v3         | Inc-v4        | IncRes-v2     | Res-101       | Inc-v3 (ens3) | Inc-v3 (ens4) | IncRes-v2 (ens) |
|-----------|-----------|----------------|---------------|---------------|---------------|---------------|---------------|-----------------|
| Inc-v3    | SI-TIM    | <b>100.00*</b> | 71.20         | 68.47         | 62.67         | 49.60         | 47.63         | 31.27           |
|           | Admix-TIM | <b>100.00*</b> | 83.90         | 80.40         | 74.40         | 59.10         | 57.90         | 39.20           |
|           | SCM-P-TIM | 99.90*         | <b>87.23</b>  | <b>84.37</b>  | 80.63         | 67.67         | 64.70         | 47.97           |
|           | SCM-R-TIM | 99.80*         | 85.50         | 83.23         | <b>81.83</b>  | <b>75.17</b>  | <b>72.77</b>  | <b>56.83</b>    |
| Inc-v4    | SI-TIM    | 78.07          | 99.60*        | 73.03         | 66.00         | 58.60         | 55.13         | 45.33           |
|           | Admix-TIM | 87.40          | <b>99.70*</b> | 82.30         | 77.00         | 68.10         | 65.30         | 53.10           |
|           | SCM-P-TIM | <b>89.40</b>   | 99.67*        | <b>85.57</b>  | 80.37         | 71.10         | 68.40         | 58.77           |
|           | SCM-R-TIM | 88.30          | 99.67*        | 85.03         | <b>80.70</b>  | <b>77.30</b>  | <b>75.93</b>  | <b>65.37</b>    |
| IncRes-v2 | SI-TIM    | 84.80          | 81.57         | <b>98.87*</b> | 76.63         | 69.90         | 64.73         | 61.47           |
|           | Admix-TIM | 90.20          | <b>88.20</b>  | 98.60*        | 83.90         | 78.40         | 73.60         | 70.00           |
|           | SCM-P-TIM | <b>90.70</b>   | 88.13         | 98.43*        | <b>84.90</b>  | 78.77         | 76.20         | 72.30           |
|           | SCM-R-TIM | 88.63          | 87.13         | 98.53*        | 84.63         | <b>82.33</b>  | <b>79.30</b>  | <b>75.80</b>    |
| Res-101   | SI-TIM    | 73.97          | 70.80         | 70.20         | <b>99.80*</b> | 59.40         | 55.07         | 43.30           |
|           | Admix-TIM | 83.20          | 78.90         | 80.70         | 99.70*        | 67.00         | 62.50         | 52.80           |
|           | SCM-P-TIM | <b>88.77</b>   | <b>86.07</b>  | <b>87.63</b>  | 99.77*        | 79.50         | 76.00         | 65.83           |
|           | SCM-R-TIM | 87.43          | 82.20         | 84.83         | 99.73*        | <b>81.33</b>  | <b>79.10</b>  | <b>70.03</b>    |

Table 2. Attack success rates (%) against four normally trained models and three adversarially trained models under a single model setting using SI-TIM, Admix-TIM, SCM-P-TIM, and SCM-R-TIM. Inc-v3, Inc-v4, IncRes-v2, and Res-101 using MI-FGSM are adopted as the source model respectively. \* indicates that the target model is the same as the source model. In case of Admix-TIM, we take the reported value.

| Model     | Attack    | Inc-v3        | Inc-v4        | IncRes-v2     | Res-101       | Inc-v3 (ens3) | Inc-v3 (ens4) | IncRes-v2 (ens) |
|-----------|-----------|---------------|---------------|---------------|---------------|---------------|---------------|-----------------|
| Inc-v3    | CTM       | 99.20*        | 85.17         | 81.77         | 76.70         | 66.07         | 63.17         | 45.70           |
|           | Admix-CTM | <b>99.90*</b> | 89.00         | <b>87.00</b>  | 83.10         | 72.20         | 71.10         | 52.40           |
|           | SCM-P-CTM | 99.10*        | <b>89.13</b>  | 86.87         | <b>84.00</b>  | 71.87         | 69.77         | 53.43           |
|           | SCM-R-CTM | 99.43*        | 87.67         | 85.83         | 83.80         | <b>78.40</b>  | <b>76.90</b>  | <b>62.57</b>    |
| Inc-v4    | CTM       | 87.07         | 98.77*        | 83.60         | 78.17         | 71.40         | 68.47         | 57.13           |
|           | Admix-CTM | <b>90.40</b>  | <b>99.00*</b> | <b>87.30</b>  | 82.00         | 75.30         | 71.90         | 61.60           |
|           | SCM-P-CTM | 90.20         | 98.63*        | 87.07         | 82.83         | 74.30         | 72.17         | 62.30           |
|           | SCM-R-CTM | 90.03         | 98.47*        | 86.83         | <b>83.57</b>  | <b>80.00</b>  | <b>78.60</b>  | <b>68.27</b>    |
| IncRes-v2 | CTM       | 88.23         | 86.07         | 97.37*        | 82.70         | 77.20         | 75.17         | 72.17           |
|           | Admix-CTM | 90.10         | <b>89.60</b>  | <b>97.70*</b> | 85.90         | 82.00         | 78.00         | 76.30           |
|           | SCM-P-CTM | <b>90.47</b>  | 88.57         | 97.43*        | 85.80         | 80.23         | 77.47         | 73.63           |
|           | SCM-R-CTM | 88.87         | 87.57         | 96.83*        | <b>85.93</b>  | <b>82.50</b>  | <b>80.23</b>  | <b>77.07</b>    |
| Res-101   | CTM       | 85.93         | 82.07         | 83.83         | 98.97*        | 76.23         | 71.80         | 62.17           |
|           | Admix-CTM | <b>91.00</b>  | 87.70         | 89.20         | <b>99.99*</b> | 81.10         | 77.40         | 70.10           |
|           | SCM-P-CTM | 90.73         | <b>88.07</b>  | <b>89.27</b>  | 99.17*        | 81.90         | 79.03         | 70.70           |
|           | SCM-R-CTM | 88.17         | 83.80         | 86.93         | 98.73*        | <b>83.87</b>  | <b>81.93</b>  | <b>74.13</b>    |

Table 3. Attack success rates (%) against four normally trained models and three adversarially trained models under a single model setting using CTM (the combination of DIM, TIM, and SIM), Admix-CTM, SCM-P-CTM, and SCM-R-CTM. Inc-v3, Inc-v4, IncRes-v2, and Res-101 using MI-FGSM are adopted as the source model respectively. \* indicates that the target model is the same as the source model. In case of Admix-CTM, we take the reported value.

| Model  | Attack      | Inc-v3         | Inc-v4       | IncRes-v2    | Res-101      | Inc-v3 (ens3) | Inc-v3 (ens4) | IncRes-v2 (ens) |
|--------|-------------|----------------|--------------|--------------|--------------|---------------|---------------|-----------------|
| Inc-v3 | MI-DI-FGSM  | 98.72*         | 64.54        | 60.46        | 54.04        | 19.22         | 17.96         | 9.56            |
|        | NI-DI-FGSM  | 99.33*         | 60.20        | 58.07        | 48.77        | 14.93         | 14.67         | 7.23            |
|        | VMI-DI-FGSM | 99.03*         | 76.43        | 73.27        | 66.53        | 39.87         | 37.57         | 22.23           |
|        | VNI-DI-FGSM | 99.27*         | 79.43        | 76.73        | 68.60        | 39.37         | 37.67         | 23.33           |
|        | EMI-DI-FGSM | 99.10*         | 83.50        | 78.00        | 70.60        | 27.80         | 26.00         | 13.40           |
|        | LI-DI-FGSM  | <b>100.00*</b> | <b>97.63</b> | <b>96.67</b> | <b>93.33</b> | <b>65.23</b>  | <b>61.07</b>  | <b>38.40</b>    |

Table 4. Attack success rates (%) against four normally trained models and three adversarially trained models under a single model setting using DIM. Inc-v3 is adopted as the source model. \* indicates that the target model is the same as the source model. In case of EMI-FGSM, we take the reported value.

| Model  | Attack      | Inc-v3         | Inc-v4       | IncRes-v2    | Res-101      | Inc-v3 (ens3) | Inc-v3 (ens4) | IncRes-v2 (ens) |
|--------|-------------|----------------|--------------|--------------|--------------|---------------|---------------|-----------------|
| Inc-v3 | MI-TI-FGSM  | <b>100.00*</b> | 48.02        | 41.28        | 40.20        | 24.10         | 21.34         | 13.48           |
|        | NI-TI-FGSM  | 99.97*         | 44.57        | 38.97        | 35.03        | 31.67         | 29.50         | 22.73           |
|        | VMI-TI-FGSM | <b>100.00*</b> | 71.00        | 68.43        | 60.10        | 32.40         | 30.80         | 17.57           |
|        | VNI-TI-FGSM | <b>100.00*</b> | 76.80        | 74.80        | 64.93        | 34.70         | 32.80         | 18.80           |
|        | EMI-TI-FGSM | <b>100.00*</b> | 79.40        | 76.30        | 67.20        | 44.30         | 40.80         | 26.20           |
|        | LI-TI-FGSM  | <b>100.00*</b> | <b>89.33</b> | <b>87.87</b> | <b>79.30</b> | <b>67.40</b>  | <b>63.93</b>  | <b>49.43</b>    |

Table 5. Attack success rates (%) against four normally trained models and three adversarially trained models under a single model setting using TIM. Inc-v3 is adopted as the source model. \* indicates that the target model is the same as the source model. In case of EMI-FGSM, we take the reported value.

| Model  | Attack      | Inc-v3         | Inc-v4       | IncRes-v2    | Res-101      | Inc-v3 (ens3) | Inc-v3 (ens4) | IncRes-v2 (ens) |
|--------|-------------|----------------|--------------|--------------|--------------|---------------|---------------|-----------------|
| Inc-v3 | MI-SI-FGSM  | <b>100.00*</b> | 69.84        | 67.66        | 62.86        | 32.20         | 31.46         | 17.42           |
|        | NI-SI-FGSM  | <b>100.00*</b> | 77.50        | 74.97        | 67.20        | 32.30         | 30.07         | 16.23           |
|        | VMI-SI-FGSM | <b>100.00*</b> | 86.87        | 83.70        | 78.10        | 55.03         | 52.83         | 35.00           |
|        | VNI-SI-FGSM | <b>100.00*</b> | 89.67        | 88.07        | 81.83        | 58.57         | 55.70         | 36.83           |
|        | EMI-SI-FGSM | <b>100.00*</b> | 91.90        | 90.00        | 85.40        | 45.20         | 41.80         | 23.80           |
|        | LI-SI-FGSM  | <b>100.00*</b> | <b>95.50</b> | <b>93.83</b> | <b>90.17</b> | <b>63.27</b>  | <b>58.73</b>  | <b>38.77</b>    |

Table 6. Attack success rates (%) against four normally trained models and three adversarially trained models under a single model setting using SIM. Inc-v3 is adopted as the source model. \* indicates that the target model is the same as the source model. In case of EMI-FGSM, we take the reported value.

| Model  | Attack      | Inc-v3         | Inc-v4       | IncRes-v2    | Res-101      | Inc-v3 (ens3) | Inc-v3 (ens4) | IncRes-v2 (ens) |
|--------|-------------|----------------|--------------|--------------|--------------|---------------|---------------|-----------------|
| Inc-v3 | MI-CT-FGSM  | 99.30*         | 84.57        | 81.53        | 75.73        | 65.50         | 63.47         | 45.80           |
|        | NI-CT-FGSM  | 99.37*         | 84.33        | 80.63        | 75.53        | 60.80         | 56.07         | 40.40           |
|        | VMI-CT-FGSM | 99.10*         | 88.63        | 86.20        | 82.57        | 77.60         | 76.13         | 64.27           |
|        | VNI-CT-FGSM | 99.50*         | 91.80        | 89.57        | 85.87        | 79.90         | 77.90         | 66.93           |
|        | EMI-CT-FGSM | 99.60*         | 94.10        | 92.60        | 89.40        | 78.90         | 75.30         | 60.40           |
|        | LI-CT-FGSM  | <b>100.00*</b> | <b>97.97</b> | <b>97.17</b> | <b>94.57</b> | <b>91.17</b>  | <b>89.00</b>  | <b>79.70</b>    |

Table 7. Attack success rates (%) against four normally trained models and three adversarially trained models under a single model setting using CTM, which is a combination of DIM, TIM, and SIM. Inc-v3 is adopted as the source model. \* indicates that the target model is the same as the source model. In case of EMI-FGSM, we take the reported value.

|        | MI-FGSM | NI-FGSM | VMI-FGSM | VNI-FGSM | EMI-FGSM | LI-FGSM   |
|--------|---------|---------|----------|----------|----------|-----------|
| Counts | 1       | 1       | $v+1$    | $v+1$    | $v$      | $N*(v+1)$ |

Table 8. The number of back-propagation in each iteration. The  $v$  is the sampling number and the  $N$  is the number of looking ahead steps of LI-FGSM