
Appendices for Towards Defending against Adversarial Examples via Attack-Invariant Features

A. Attack implementation

In this section, we present supplementary information on attack implementation. We use Advtorch Toolbox¹ to implement the projected gradient descent method (PGD), the decoupling direction and norm method (DDN), the Carlini and Wagner method (CW), the Jacobian-based saliency map attack method (JSMA) and the spatial transform attack method (STA). The autoattack method (AA)², the faster wasserstein attac method (FWA)³ and the robust physical perturbation method (RP₂)⁴ are implemented from their open source codes. On *MNIST*, the main parameters of attacks are as follows:

- **PGD:** We use the L_∞ norm PGD method to craft adversarial examples. The default perturbation budget is set to 0.3. The default number of iterations is set to 40. The attack step size is set to 0.01.
- **DDN:** The number of iterations is set to 100. The factor to modify the norm at each iteration is set to 0.05. The number of quantization levels is set to 256.
- **CW:** We use the L_2 norm CW method to craft adversarial examples. The maximum number of iterations is set to 1000. The confidence of the adversarial examples is set to 1. The initial value of the constant is set to 1.
- **JSMA:** The highest percentage of pixels can be modified is set to 1.0. The perturb length is set to 1.0.
- **STA:** The maximum number of iterations is set to 500. The number of search times to find the optimum is set to 20.
- **AA:** The default perturbation budget is set to 0.3. The default number of iterations is set to 100.
- **FWA:** The wasserstein adversarial examples are crafted by exploiting PGD. The default perturbation budget is set to 0.3. The number of iterations is set to 40. The learning rate is set to 0.1.

On *CIFAR-10*, the main parameters of attacks are as follows:

¹<https://github.com/codeaudit/advtorch>

²<https://github.com/fra31/auto-attack>

³<https://github.com/watml/fast-wasserstein-adversarial>

⁴<https://github.com/tongwu2020/phattacks/tree/master/sign/experiment>

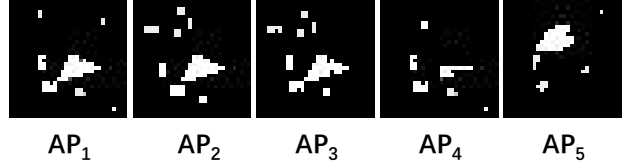


Figure 1. A visual illustration of five masks used to craft adversarial patches (AP).

- **PGD:** We use the L_∞ norm PGD method to craft adversarial examples. The default perturbation budget is set to $8/255$. The default number of iterations is set to 40. The attack step size is set to 0.01.
- **DDN:** The number of iterations is set to 100. The factor to modify the norm at each iteration is set to 0.05. The number of quantization levels is set to 256.
- **CW:** We use the L_2 norm CW method to craft adversarial examples. The maximum number of iterations is set to 500. The confidence of the adversarial examples is set to 1. The initial value of the constant is set to 1.
- **JSMA:** The highest percentage of pixels can be modified is set to 1.0. The perturb length is set to 1.0.
- **STA:** The maximum number of iterations is set to 200. The number of search times to find the optimum is set to 20.
- **AA:** The default perturbation budget is set to $8/255$. The default number of iterations is set to 100.
- **FWA:** The wasserstein adversarial examples are crafted by exploiting PGD. The default perturbation budget is set to $8/255$. The number of iterations is set to 40. The learning rate is set to 0.01.

On *LISA*, we use five different masks to implement RP₂ for crafting adversarial examples. The masks are shown in Figure 1. The number of iterations of RP₂ is set to 300 and the learning rate is set to 0.01.

B. Defense results

In this section, we present supplementary information on defense results. We use two different combinations of seen types of attacks to train our ARN model: (i) the targeted

055 PGD and the non-targeted PGD (“ARN_{PP}”). (ii) the non-
 056 targeted DDN and the non-targeted PGD (“ARN_{DP}”). Fig-
 057 ure 2, 3 shows examples which are restored by our *adversarial*
 058 *noise removing network* (ARN) against pixel-constrained
 059 attacks on *MNIST* and *CIFAR-10*. These attacks include
 060 non-targeted L_∞ norm PGD (PGD_N), targeted L_∞ norm
 061 PGD (PGD_T), non-targeted DDN (DDN_N), non-targeted
 062 L_2 norm CW (CW_N), targeted JSMA (JSMA_T) and non-
 063 targeted AA (AA_N). Figure 4 shows examples which are
 064 restored by our ARN against spatial-constrained attacks.
 065 These attacks include non-targeted STA (STA_N), targeted
 066 STA (STA_T), non-targeted FWA (FWA_N) and non-targeted
 067 RP₂ (RP_N). Figure 5 and 6 show adversarial examples and
 068 restored examples on *LISA*. Five types of adversarial patches
 069 (AP) are crafted by RP₂ and are added to natural examples
 070 to generate adversarial examples. We use adversarial exam-
 071 ples with two types of adversarial patches (AP₁ and AP₂)
 072 as training data to train our ARN model. The categories
 073 corresponding to the class labels in *CIFAR-10* are as follows:

0) airplane, 1) car, 2) bird, 3) cat, 4) deer, 5) dog, 6) frog, 7)
 horse, 8) boat and 9) truck.

C. Leaked defenses

In this section, we present supplementary information on de-
 fending under difficult scenarios where defenses are leaked.
 We use ARN_{PP}, APE-G_{PP} and HGD_{PP} as the leaked de-
 fense models to craft adversarial examples by distinct
 attacks. Figure 7 shows defense results of our ARN against
 BPDA. The adversarial examples are crafted by jointly us-
 ing PGD_N and BPDA against our ARN. Figure 8, 9 and
 10 show defense results of our ARN against white-box and
 gray-box adaptive attacks. To be specific, the leaked de-
 fense in Figure 8 is our ARN_{PP} and the attack is PGD_T. The
 leaked defense in Figure 9 is APE_{PP} and the attacks are
 PGD_T, CW_N and CW_T. The leaked defense in Figure 10
 is HGD_{PP} and the attacks are PGD_T, DDN_N and DDN_T.



106 Figure 2. A visual illustration of adversarial examples and their restored examples. These adversarial examples are crafted by pixel-
 107 constrained attacks.

110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164



Figure 3. A visual illustration of adversarial examples and their restored examples. These adversarial examples are crafted by pixel-constrained attacks.



Figure 4. A visual illustration of adversarial examples and their restored examples. These adversarial examples are crafted by spatially-constrained attacks.

165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219

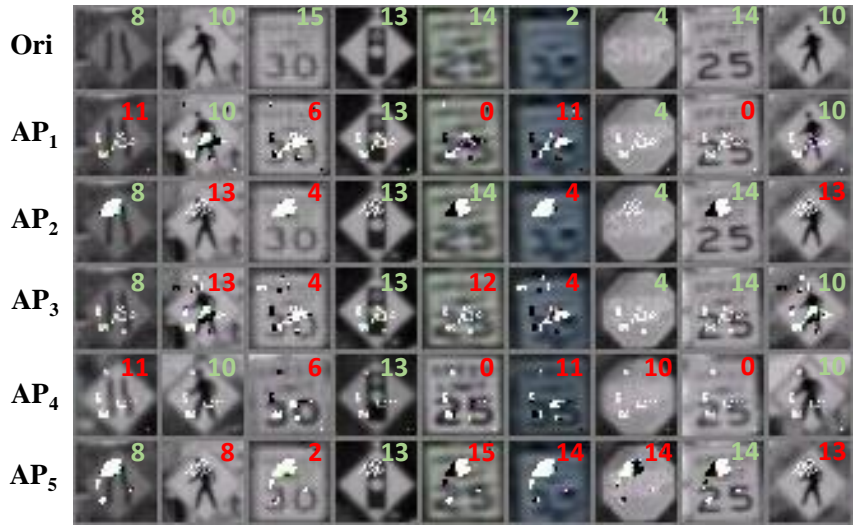


Figure 5. A visual illustration of adversarial examples on *LISA*. Five types of adversarial patches (AP) are crafted by RP_2 and are added to natural examples to generate adversarial examples.

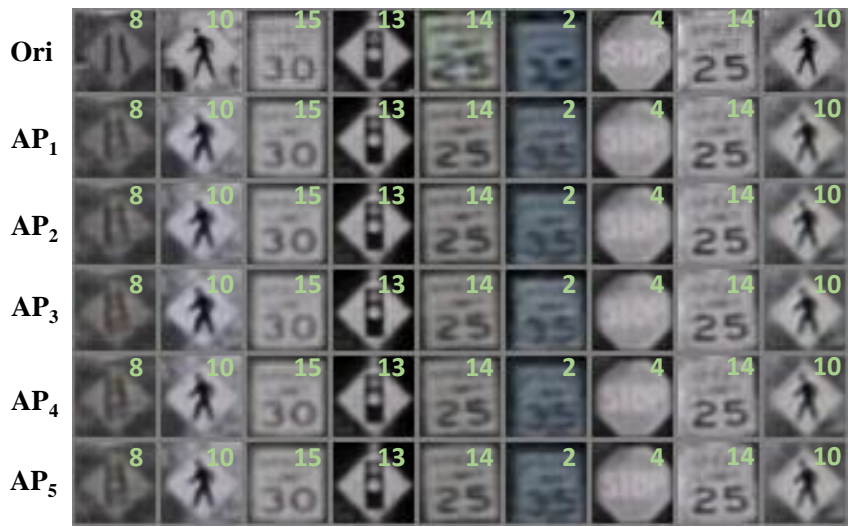


Figure 6. A visual illustration of restored examples on *LISA*. Five types of adversarial patches (AP) are crafted by RP_2 and are added to natural examples to generate adversarial examples.

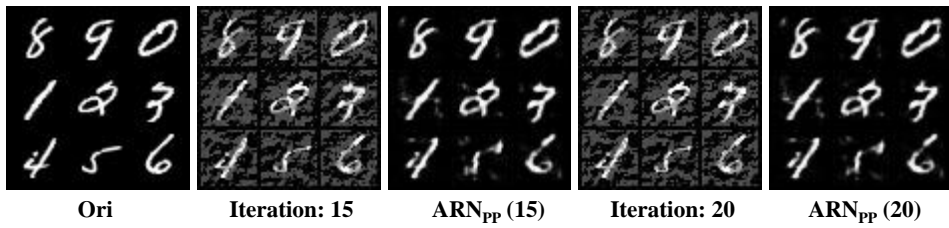
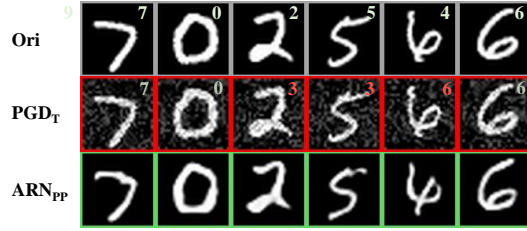
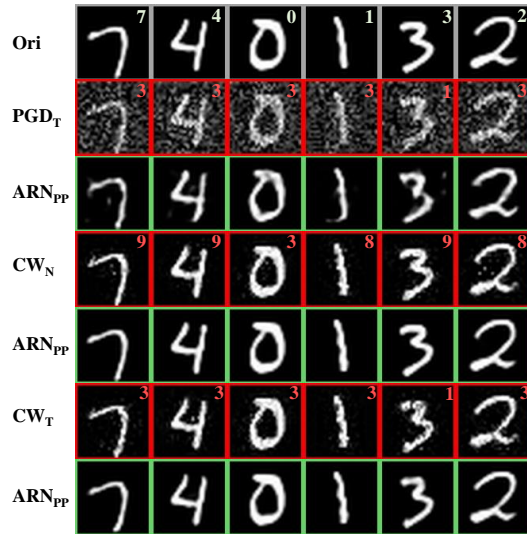


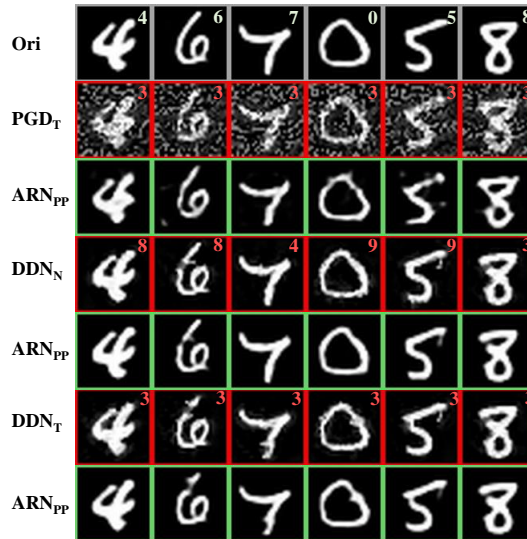
Figure 7. A visual illustration of adversarial examples and their restored examples against BPDA. The adversarial examples are crafted by jointly using BPDA and PGD_N against our ARN. The number of iterations of PGD_N is 15 and 20 respectively.



228 *Figure 8.* A visual illustration of defense results of our ARN against the white-box adaptive attack.



250 *Figure 9.* A visual illustration of defense results of our ARN against the gray-box adaptive attack. The local defense model is APE_{PP}.



272 *Figure 10.* A visual illustration of defense results of our ARN against the gray-box adaptive attack. The local defense model is HGD_{PP}.