

# PatchAttack: A Black-box Texture-based Attack with Reinforcement Learning

Chenglin Yang, Adam Kortylewski, Cihang Xie, Yinzhi Cao, Alan Yuille

CCVL, Johns Hopkins University

# Patch-Attack: Summary

- Digital Perturbation Attacks which are imperceptible can fool Deep Nets. But many of these attacks are *white box*, i.e. require knowing the weights of the Deep Net, and there are a growing number of effective defenses.
- We develop a black box texture-based attack which requires no knowledge of the Deep Net. Our method learns an attack policy which selects texture-patches from a pre-defined dictionary.
- We show that we can successfully attack a range of Deep Nets for classification and reducing their effectiveness by up to 90% using only small patches.
- The attack is resistant to existing defenses.
- Chenglin Wang, Adam Kortylewski, Cihang Xie, Yinzhi Cao, Alan Yuille.  
“PatchAttack: A Black-box Texture-based Attack with Reinforcement Learning”.  
ECCV. In review. 2020.

# Texture-Patch Attacks

Use Reinforcement Learning to learn a policy which, for each image, selects texture-patches and places them in the image so as to:

- A. Optimize the position and scale of the patch(es)
- B. Optimize the texture-pattern of the patch(es)

Mathematical Formulation:

$$\mathcal{L}(\mathbf{y}, y'), \quad \text{where } \mathbf{y} = \mathbf{f}(\mathbf{g}(\mathbf{x}); \boldsymbol{\theta}), \quad (1)$$

$$\mathbf{g}(\mathbf{x}) : \begin{cases} x_{u,v} = \mathbf{T}(x_{u,v}), & \text{if } (u, v) \in \mathcal{E} \\ x_{u,v} = x_{u,v}, & \text{otherwise} \end{cases} \quad (2)$$

$$\mathcal{E} = \mathbf{s}(\mathbf{x}, \mathbf{f}(\cdot, \boldsymbol{\theta}), \mathcal{S}) \subseteq \{(u, v) \mid u \in [0, H], v \in [0, W]\} \quad (3)$$

# PatchAttack: MPA

Patch Search with Reinforcement Learning:

$$\mathcal{S} = \{(u_1^1, v_1^2, u_1^3, v_1^4, \dots, u_C^1, v_C^2, u_C^3, v_C^4)\} \quad (4)$$

$$\mathbb{A}(\theta_{\mathbb{A}}) : P(a_t | (a_1, \dots, a_{t-1}), \mathbf{f}(\cdot; \theta), \mathbf{x}) \quad t = \{1, \dots, 4C\} \quad (5)$$

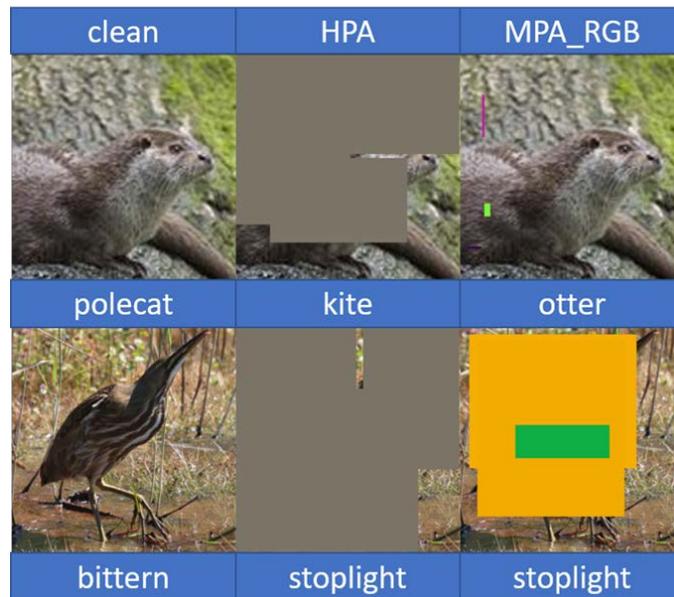
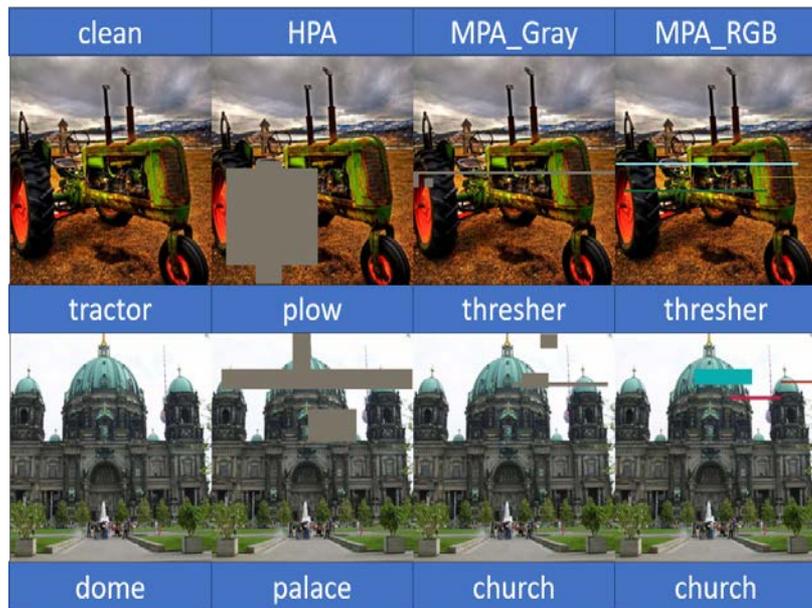
$$\mathbf{r} = \begin{cases} \ln y' - \mathbf{A}(\mathbf{a}) / \sigma^2, & \text{target attack} \\ \ln(1 - y) - \mathbf{A}(\mathbf{a}) / \sigma^2, & \text{non-target attack} \end{cases} \quad (6)$$

$$\text{MPA} : \begin{cases} \mathcal{E} = \mathbf{J}(\mathbf{a}) \\ \mathbf{T}(x_{u,v}) = 0 \\ \mathcal{L} = -\mathbf{r} \cdot \ln \mathbf{P} \end{cases} \quad (7)$$

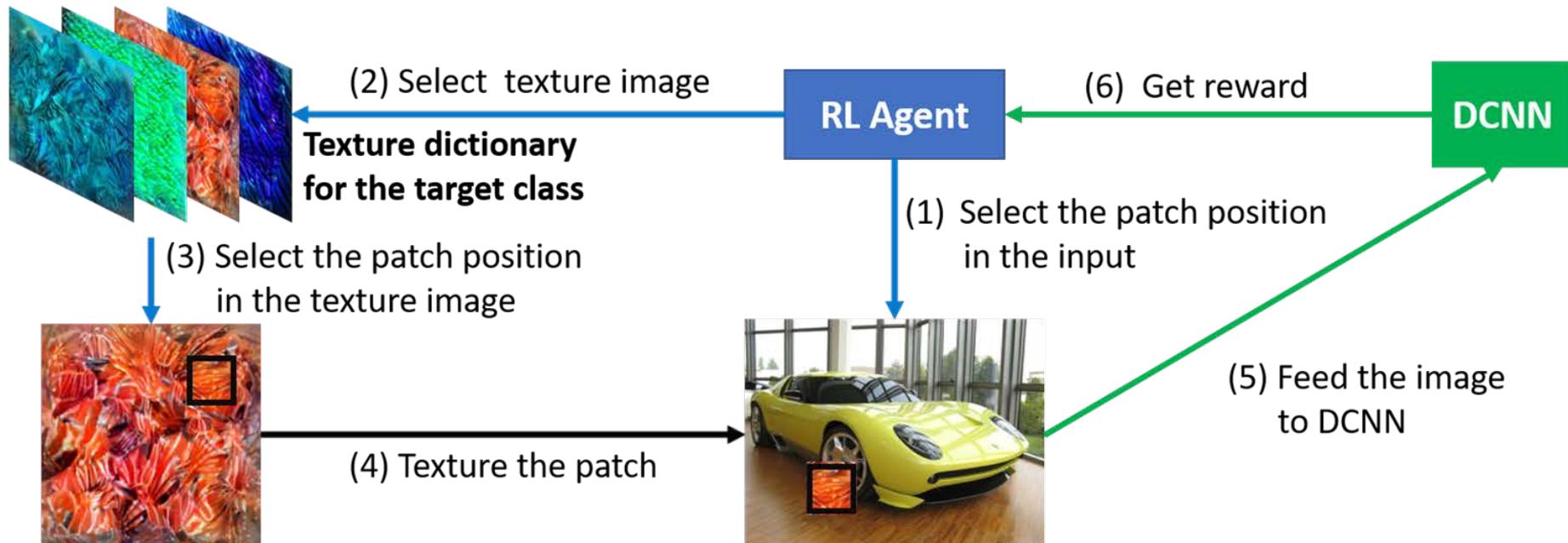
$$\mathcal{S} = \{(u_1^1, v_1^2, u_1^3, v_1^4, R_1^5, G_1^6, B_1^7, \dots, u_C^1, v_C^2, u_C^3, v_C^4, R_C^5, G_C^6, B_C^7)\} \quad (8)$$



MPA: small patches confuse similar object classes.

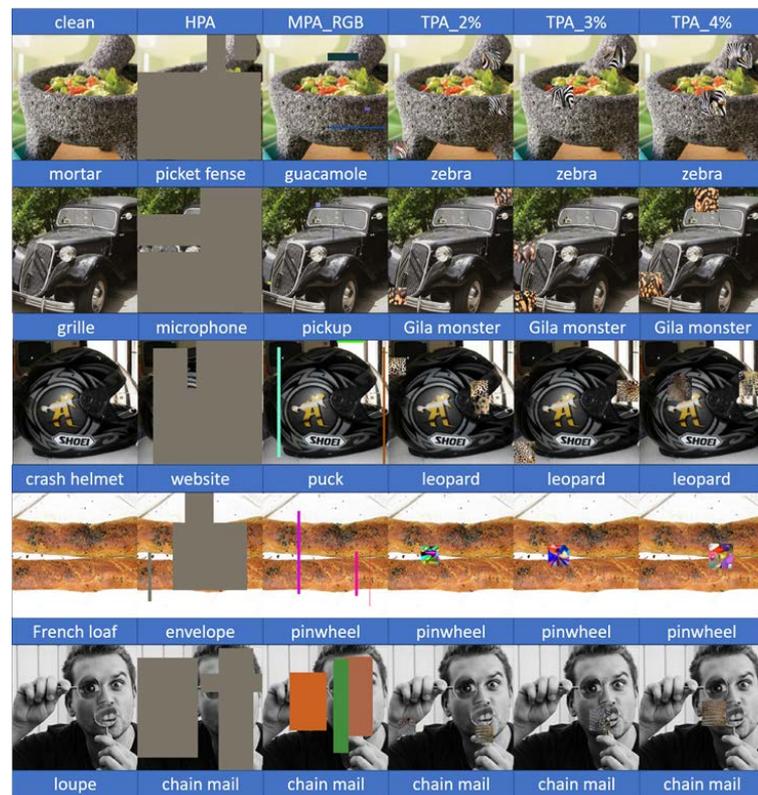
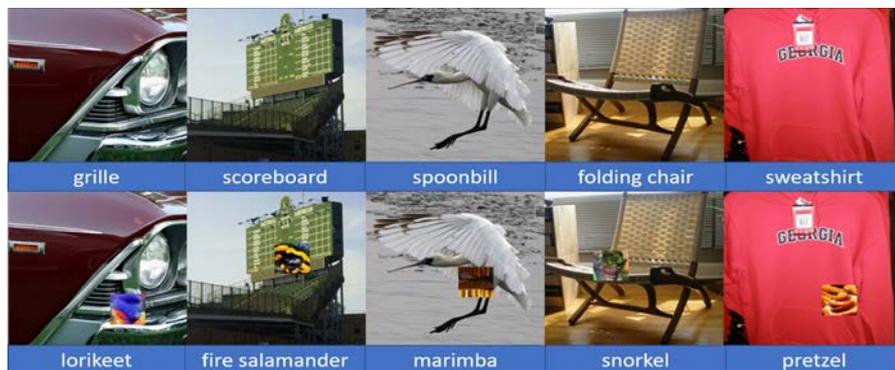


# PatchAttack: TPA



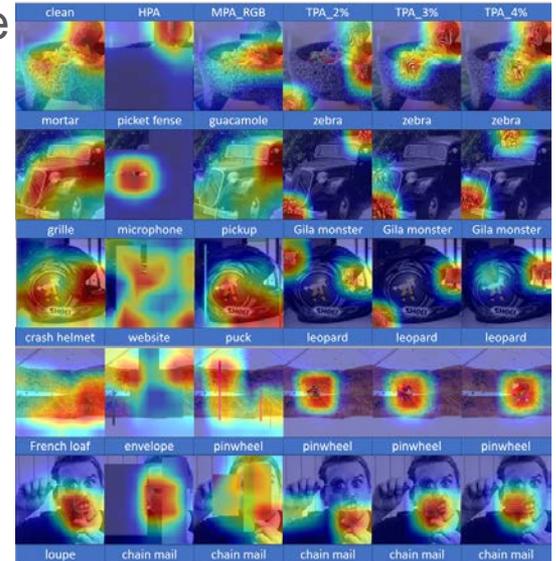
# PatchAttack: TPA can confuse very different classes

- Here are examples.



# Intuition: Why Do Patch-Attacks Work?

- The Deep Net has an internal attention mechanism (buried in the weights). The Patch-Attack hijacks the attention mechanism so that it focuses on the patch. The Deep Net thinks the patch looks more like object itself.
- This shows a major limitation of Deep Networks.
- Recall that the texture-patches look like this:



# Experiments

## Non-targeted Attack

1000 images randomly selected  
from the validation set of ImageNet

Network	Attack	Acc. (%)	Avg_area (%)	Avg_qry
ResNet50	–	72.80	–	–
	<b>HPA</b>	0.40	18.05	10000
	<b>MPA_Gray</b>	0.00	6.57	9659
	<b>MPA_RGB</b>	0.00	5.41	9681
	<b>TPA_N4.4%</b>	0.30	5.06	1137
	<b>TPA_N8.2%</b>	0.30	3.10	983
DenseNet121	–	74.10	–	–
	<b>HPA</b>	0.10	19.82	10000
	<b>MPA_Gray</b>	0.00	6.87	9624
	<b>MPA_RGB</b>	0.00	5.73	9696
	<b>TPA_N4.4%</b>	0.50	5.13	1195
	<b>TPA_N8.2%</b>	0.30	3.13	1001
ResNeXt50	–	76.20	–	–
	<b>HPA</b>	0.80	19.22	10000
	<b>MPA_Gray</b>	0.00	7.88	9748
	<b>MPA_RGB</b>	0.00	6.23	9752
	<b>TPA_N4.4%</b>	0.70	5.21	1280
	<b>TPA_N8.2%</b>	0.50	3.25	1088
MobileNet-V2	–	68.80	–	–
	<b>HPA</b>	0.20	16.61	10000
	<b>MPA_Gray</b>	0.00	5.35	9578
	<b>MPA_RGB</b>	0.00	4.11	9603
	<b>TPA_N4.4%</b>	0.30	4.63	862
	<b>TPA_N8.2%</b>	0.30	2.74	756

# Experiments

Targeted Attack

1000 images randomly selected  
from the validation set of ImageNet

Target labels are randomly selected

Network	Attack	T_acc. (%)	Avg_area (%)	Avg_qry
ResNet50	–	0.10	–	–
	<b>HPA</b>	23.20	71.54	50000
	<b>MPA_RGB</b>	25.90	18.45	28361
	<b>TPA_N10_2%</b>	97.60	7.80	15728
	<b>TPA_N10_4%</b>	99.70	9.97	8643
	<b>TPA_N10_10%</b>	100.00	15.36	3747
DenseNet121	–	0.10	–	–
	<b>HPA</b>	21.50	71.68	50000
	<b>MPA_RGB</b>	24.90	19.38	28088
	<b>TPA_N10_2%</b>	97.10	7.87	15920
	<b>TPA_N10_4%</b>	99.90	10.19	8953
	<b>TPA_N10_10%</b>	100.00	15.84	3970
ResNeXt50	–	0.00	–	–
	<b>HPA</b>	25.40	72.57	50000
	<b>MPA_RGB</b>	27.60	13.86	24738
	<b>TPA_N10_2%</b>	97.60	7.59	15189
	<b>TPA_N10_4%</b>	99.70	9.60	8223
	<b>TPA_N10_10%</b>	100.00	15.04	3538
MobileNet-V2	–	0.10	–	–
	<b>HPA</b>	22.10	69.45	50000
	<b>MPA_RGB</b>	27.70	16.64	28294
	<b>TPA_N10_2%</b>	98.50	7.78	15479
	<b>TPA_N10_4%</b>	99.90	10.39	8948
	<b>TPA_N10_10%</b>	100.00	16.85	4422

# Conclusion

We propose PatchAttack, a powerful black-box texture-based patch attack.

- Show that even small textured patches are able to break deep networks
- Monochrome Patch Attack achieves a strong performance on non-targeted attack, surpassing previous work by a large margin using less queries and smaller patch areas
- Texture-based Patch Attack achieves exceptional performance in both targeted and non-targeted attacks
- PatchAttack breaks traditional SOTA defenses and shape-based networks

