

Introduction

Adversarial defenses are proposed to address the problem of adversarial examples. However, the authors of many defenses provide over-estimated robustness using fixed set of common techniques. These defenses are broken later with handcrafted adaptive attacks which are designed to reflect the defense mechanism. Yet this approach requires strong domain expertise.

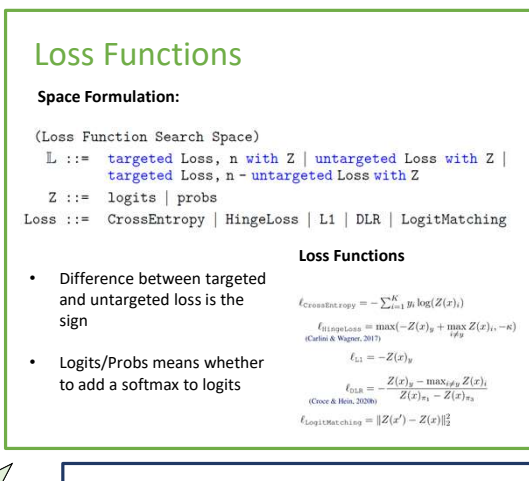
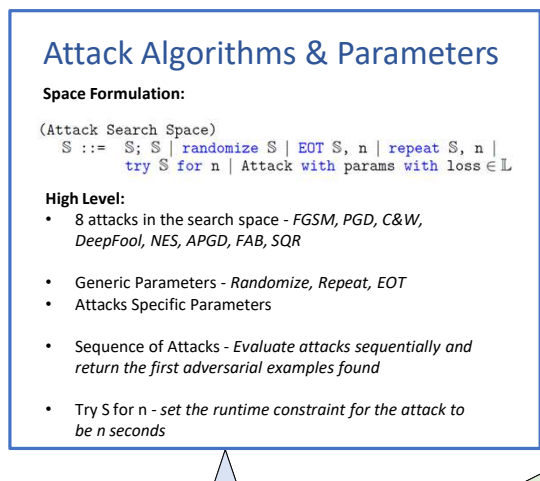
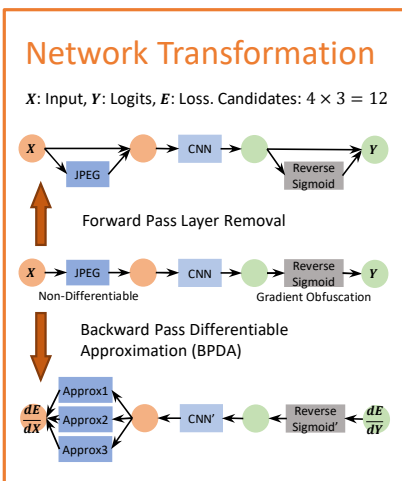
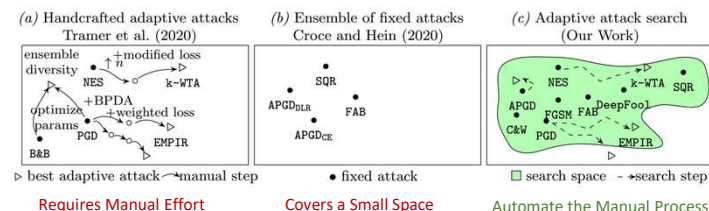
Our Work: We present an extensible tool A^3 that defines a search space over reusable blocks and automatically discovers an effective attack given the defense.

Motivation

Example Defenses	Robustness by authors	Handcrafted attacks (Tramer et al. 2020)
ME-Net (Yang et al. 2019)	53%	15%
Error Correcting Codes (Verma&Swami, 2019)	57%	5%
kWinner Takes All (Xiao et al. 2020)	51%	0.2%

Our work: automate this adaptive process

Robustness Evaluation Paradigms

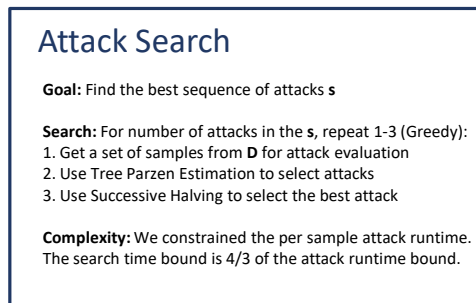
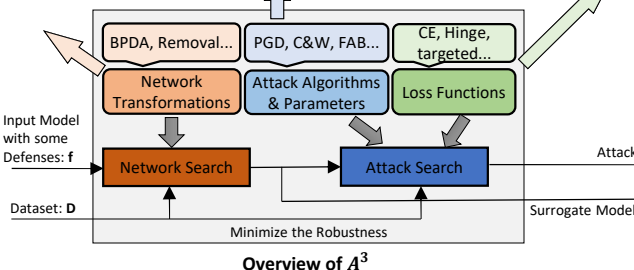
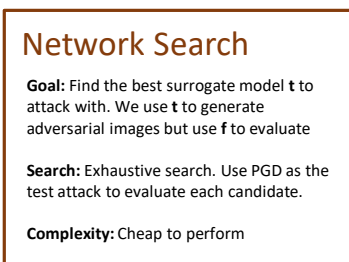


Result

A^3 is evaluated on 24 defenses and compared with AutoAttack (AA)

- 10 cases: **3.0%-50.8%** additional adversarial examples.
- 13 cases: Typically **2x** faster attack time.

CIFAR-10, ℓ_{∞}	AA	A^3	Δ	AA	A^3	Speed-up	A^3
A1	44.78	44.69	-0.09	25	20	1.25x	88
A2 ¹	2.29	1.96	-0.33	9	7	1.29x	116
A3 ²	0.59	0.11	-0.48	6	2	3.00x	40
A4	6.17	3.04	-3.13	21	13	1.62x	80
A5	22.30	12.14	-10.16	19	17	1.12x	99
A6 ³	4.14	3.94	-0.20	28	24	1.17x	237
A7	2.85	2.71	-0.14	4	4	1.00x	84
A8	19.82	11.11	-8.71	49	22	2.23x	189
A9	64.91	63.56	-1.35	157	100	1.57x	179
A9 ¹	64.91	17.70	-47.21	157	2,280	0.07x	1,548
B10 ⁴	62.80	62.79	-0.01	818	226	3.62x	761
B11 ⁵	60.04	60.01	-0.03	706	255	2.77x	690
B12 ⁶	59.64	59.56	-0.08	604	261	2.31x	565
B13 ⁷	59.53	59.51	-0.02	638	282	2.26x	575
B14 ⁸	57.14	57.16	0.02	671	429	1.56x	691
C15 ⁹	77.64	39.54	-38.10	101	108	0.94x	296
C15 ¹⁰	77.64	26.87	-50.77	101	205	0.49x	659
C16 ¹¹	36.74	37.11	0.37	381	302	1.26x	726
C17	5.15	5.16	0.01	107	114	0.94x	749
C18	5.40	2.31	-3.09	95	146	0.65x	828
C19	50.84	50.81	-0.03	734	372	1.97x	755
C20 ¹²	60.72	60.04	-0.68	621	210	2.96x	585
C21 ¹³	15.27	5.24	-10.03	261	79	3.30x	746
C22	49.53	41.99	-7.54	255	462	0.55x	900
C23	22.29	13.45	-8.84	114	374	0.30x	1,023
C24	6.25	3.07	-3.18	110	56	1.96x	502



In addition, the attacks found by A^3 can reflect the defense mechanism. (Analysis for C15, C18, C24 are shown in the paper)