

Efficient and Transferable Adversarial Examples from Bayesian Neural Networks



M. Gubri, M. Cordy, M. Papadakis, Y. Le Traon, University of Luxembourg
K. Sen, University of California, Berkeley

Introduction

Adversarial examples can transfer from a surrogate DNN to a target DNN.

Previous ensemble approach

Diversity (enrich the space of hypothesis) → transferability

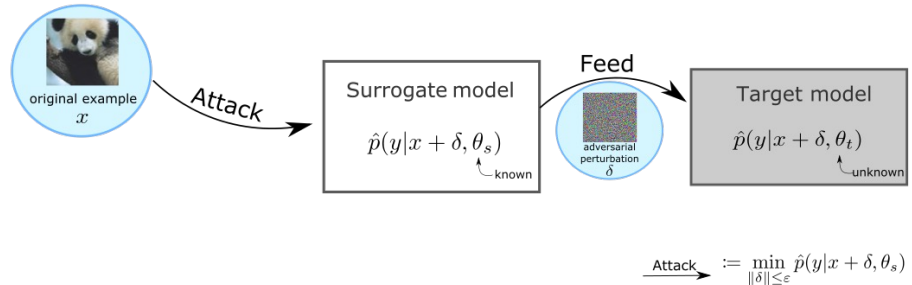
Transferability is Bayesian

The weights of the target DNN are **unknown** and can therefore be treated as random variables.

The posterior represents an (approximate) **belief** about the DNN weights.

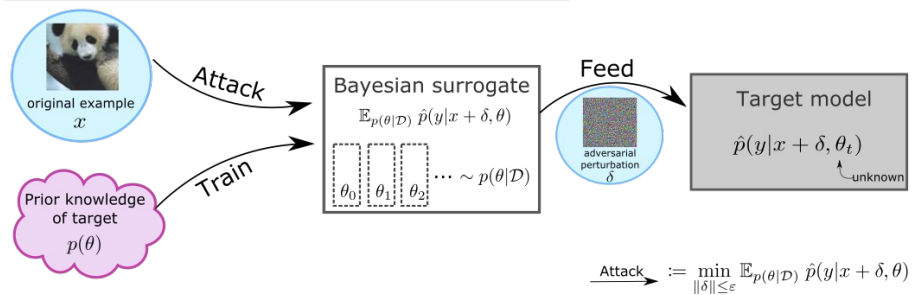
Deterministic Transfer-based Black-box Attack

Lower success rate



Bayesian Transfer-based Black-box Attack

Higher success rate



Do Bayesian deep learning techniques improve transfer-based attacks by capturing uncertainty?

Method

1. Collect samples from the posterior with cSGLD, a cyclical SG-MCMC method
2. Approximate Bayesian Model Averaging during iterative attack

Effective and efficient

Effectiveness

Improve four attacks success rate from -2.3 to 83.2 percent points

Efficiency

Divide training computations by 2.5+ compared to Deep Ensemble

Table 1: Number of DNNs (F-DEE) and training computation budget (in flops) to achieve the intra-architecture transferability of cSGLD with Deep Ensemble. Higher is better. *>15* means that 15 DNNs always transfer less than cSGLD.

Dataset	Attack	Norm	F-DEE	Flops Ratio
ImageNet	I-FGSM	L2	4.91	2.84
		L∞	4.34	2.51
	MI-FGSM	L2	4.69	2.71
		L∞	4.38	2.53
	PGD	L2	5.00	2.89
		L∞	4.42	2.56
CIFAR10	FGSM	L2	5.81	3.35
		L∞	5.98	3.46
	I-FGSM	L2	5.15	3.15
		L∞	3.76	2.36
	MI-FGSM	L2	5.56	3.36
		L∞	2.88	1.87
MNIST	PGD	L2	5.15	3.15
		L∞	3.74	2.34
	FGSM	L2	3.72	2.32
		L∞	3.72	2.32
	I-FGSM	L2	5.15	3.15
		L∞	3.42	2.12
MI-FGSM	L2	5.15	3.15	
	L∞	2.79	1.79	
PGD	L2	5.15	3.15	
	L∞	3.26	2.03	
FGSM	L2	5.15	3.15	
	L∞	5.15	3.15	

Table 2: Transfer success rates of I-FGSM attack on ImageNet hold-out architectures. Higher is better.

Norm	Surrogate	Target Architecture					Nb epochs
		-ResNet50	-ResNeX150	-DenseNet121	-MNASNet	-EffNetB0	
L2	1 cSGLD per arch.	93.28	90.61	92.25	95.98	81.88	4 × 135
	1 DNN per arch.	72.99	72.31	66.72	84.21	53.99	4 × 135
L∞	1 cSGLD per arch.	92.24	89.83	90.86	96.86	79.40	4 × 135
	1 DNN per arch.	69.65	69.01	61.00	82.25	49.71	4 × 135

Outperform related work

Ours vs. related work

Achieve 87.5% of the time higher transferability than

3 test-time techniques

Ours + related work

Always improve 3 test-time techniques when combined

Surrogate	Target Architecture				
	ResNet50	ResNeX150	DenseNet121	MNASNet	EffNetB0
1 DNN	47.81	32.29	23.43	22.52	12.77
+ Input Diversity	76.55	62.57	50.17	49.31	32.64
+ Skip Gradient Method	66.36	51.60	39.05	45.60	30.69
+ Ghost Networks	67.02	46.74	32.57	31.12	17.68
+ Momentum (MI-FGSM)	55.12	38.47	28.19	27.55	16.34
+ Input Diversity	+82.47	+69.69	57.79	+55.99	+38.63
+ Skip Gradient Method	68.39	34.57	41.48	47.97	+33.16
+ Ghost Networks	71.27	51.46	36.91	34.54	20.51
cSGLD	78.71	65.11	61.49	51.81	31.11
+ Input Diversity	90.03	82.13	81.19	74.68	53.15
+ Skip Gradient Method	81.37	69.88	65.20	71.68	52.15
+ Ghost Networks	87.33	76.00	71.67	61.45	37.19
+ Momentum (MI-FGSM)	82.89	70.42	66.39	56.68	36.00
+ Input Diversity	93.97	87.69	86.78	81.08	60.87
+ Skip Gradient Method	84.19	73.14	67.35	74.36	55.30
+ Ghost Networks	89.53	78.69	73.33	63.56	39.79

Other deep learning methods

- cSGLD competitive compared to 5 other Bayesian and Ensemble techniques

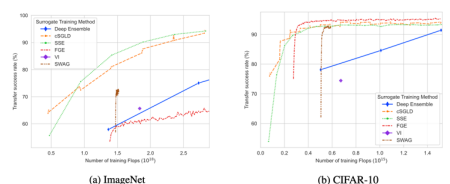


Figure 2: Intra-architecture L∞ I-FG(S)M success rate with respect to the training computational complexity of an increasing number of samples from six training techniques.

More information →



martin.gubri@uni.lu