ICLR 2018 Vancouver

Towards Reverse-Engineering Black-Box Neural Networks Seong Joon Oh Max Augustin Bernt Schiele Mario Fritz Max Planck Institute for Informatics, Germany

Problem & Motivation

- Many deployed models are black boxes (given input, returns output).

- Can black-box accesses reveal model internals? e.g. (1) architecture, (2) training procedure, and (3) training data.

- Why does it matter? Key intellectual property and increased vulnerability to attacks.

MNIST Setup

Code	Attribute	Values
act	Activation	ReLU, PReLU, ELU, Tanh
drop	Dropout	Yes, No
pool	Max pooling	Yes, No
ks	Conv ker. size	3, 5
#conv	#Conv layers	2, 3, 4
#fc #FC layers		2, 3, 4
		$2^{14}, \cdots, 2^{21}$
ens	Ensemble	Yes, No
alg	Algorithm	SGD, ADAM, RMSprop
bs	Batch size	64, 128, 256
split	Data split	All ₀ , Half _{0/1} , Quarter _{0/1/2/3}
size	Data size	All, Half, Quarter
	act drop pool ks #conv #fc #fc #par ens alg bs split	actActivationdropDropoutpoolMax poolingksConv ker. size#conv#Conv layers#fc#FC layers#fc#FC layers#par#ParametersensEnsemblealgAlgorithmbsBatch sizesplitData split

- 12 attributes to expose.

- META-MNIST: Dataset of 11k diverse digit classifiers covering 12 attributes.

- Random split (R): iid train-test models (but still disjoint).

- Extrapolation split (E): train-test differ by 1+ attributes.

Is there a {e.g. max-pool layer} in this black-box digit classifier?

Method 1. kennen-o : Infer attributes from output patterns w.r.t. some fixed inputs. $\min_{\theta} \mathbb{E}_{f \sim \mathcal{F}} \left[\mathcal{L} \left(m_{\theta} \left([f(x^{i})]_{i=1}^{n} \right), y \right) \right]$

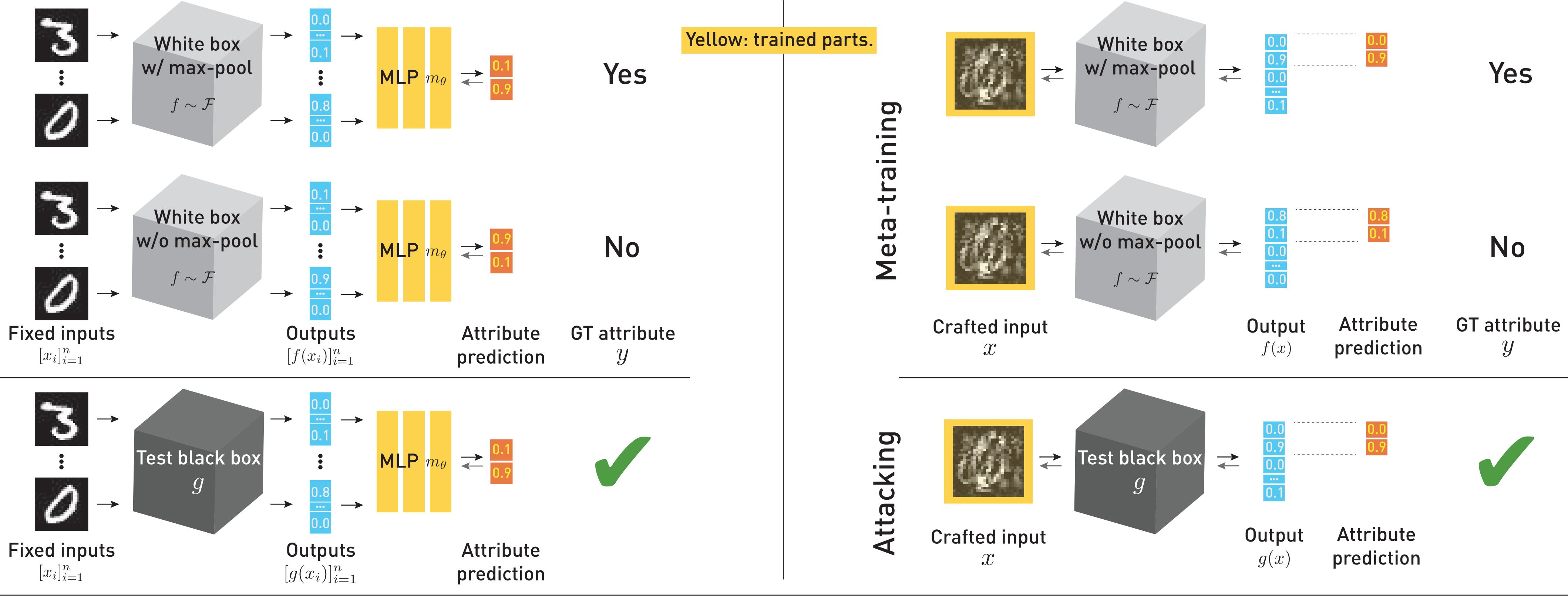


Method 3. kennen-io : Combined model that simultaneously crafts effective inputs and learns to infer attributes from the corresponding output patterns.

MNIST Results Main results (Random split)

		architecture			opt	tim	da	ata						
Method	Output	act	drop	pool	ks	#conv	#fc	#par	ens	alg	bs	size	split	avg
Chance	_	25.0	50.0	50.0	50.0	33.3	33.3	12.5	50.0	33.3	33.3	33.3	14.3	34.9
kennen-o kennen-o	score ranking		94.6 93.8			0,,,=		41.7 44.1	• • • • •	71.8 65.3	50.4 47.0		90.0 86.6	73.4 69.7
kennen-i	1 label	43.5	77.0	94.8	88.5	54.5	41.0	32.3	46.5	45.7	37.0	42.6	29.3	52.7
kennen-io	score	88.4	95.8	99.5	97.7	80.3	80.2	45.2	60.2	79.3	54.3	84.8	95.6	80.1

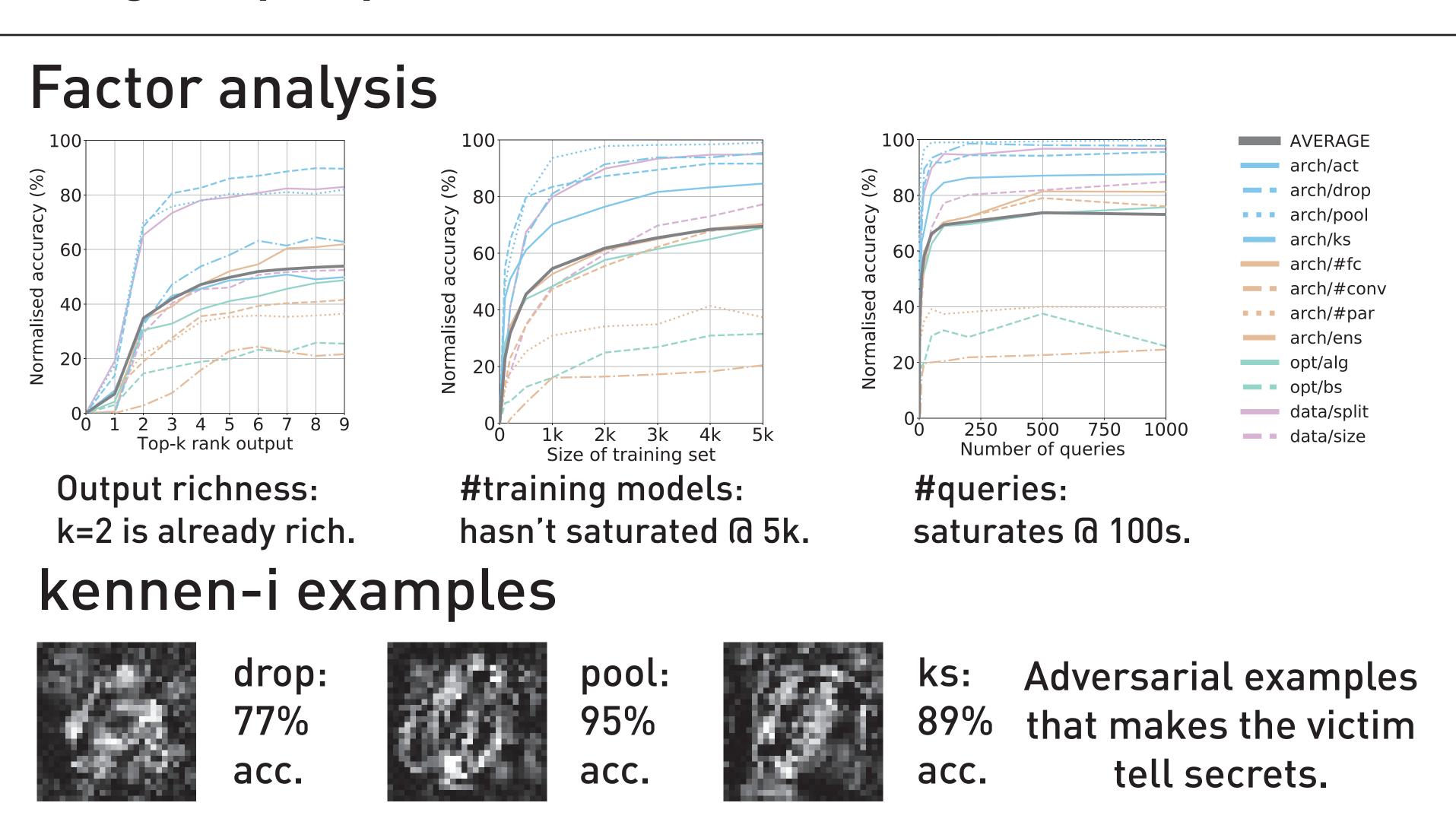
- can be exposed.



- Always far above chance. - Easy: act, drop, pool, ks, split. #layers can be exposed. Optimisation hyperparams

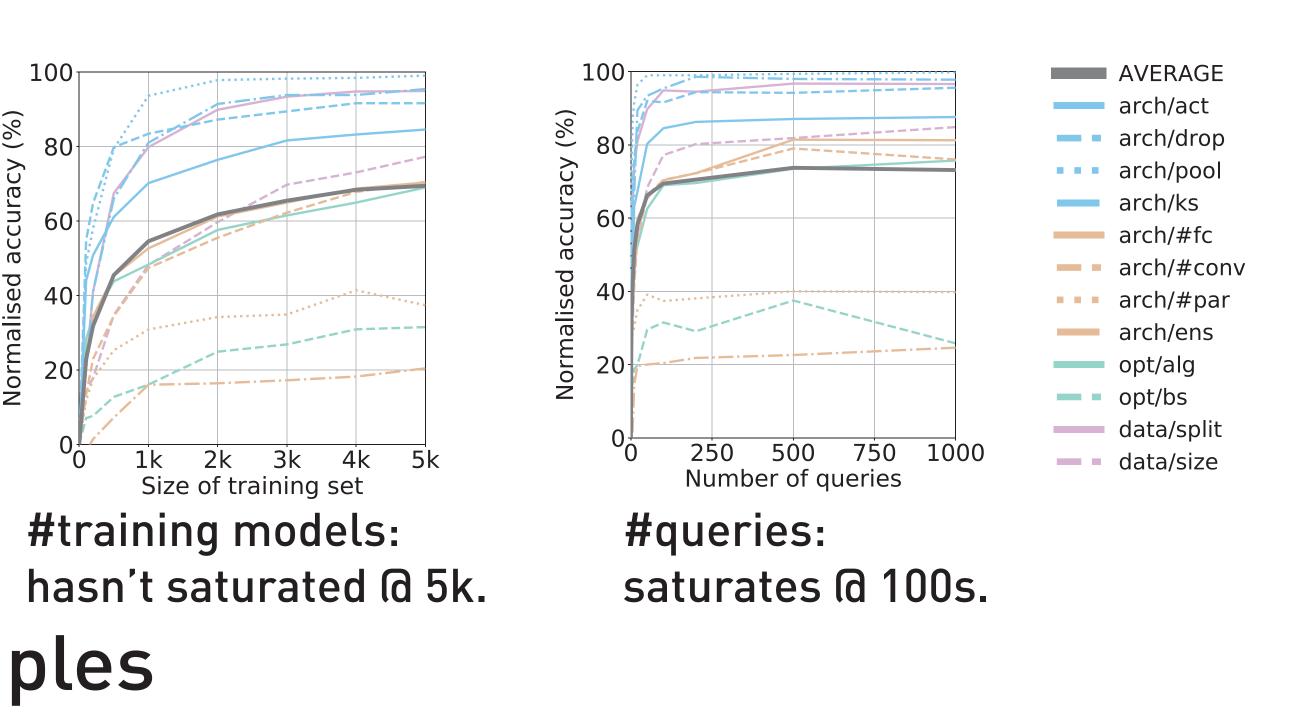
 Don't need full score output; - kennen-o is stealthy and effective. top-k is good enough.

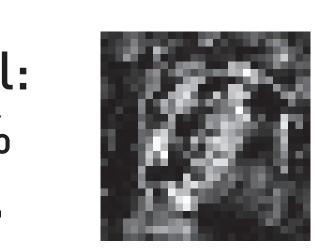
- kennen-i is query-efficient and effective.
- kennen-io is the best in all but one attributes.



Method 2. kennen-i : Craft "adversarial" input whose output contains model attribute information. $\min_{x:\text{ image } f \sim \mathcal{F}} \mathbb{E} \left[\mathcal{L} \left(f(x), y \right) \right]$

 $\min_{[x^i]_{i=1}^n: \text{ images } \theta} \min_{\theta} \mathbb{E}_{f \sim \mathcal{F}} \left[\mathcal{L} \left(m_{\theta} \left([f(x^i)]_{i=1}^n \right), y \right) \right]$





Extrapolation results

Split	Train	Test	R.Acc
R	_	_	100
E-#conv	2,3	4	92.1
E-#conv-#fc	2,3	4	80.7
E-alg	SGD,ADAM	RMSprop	88.5
E-alg-bs	64,128	256	70.1
E-size	Quarter	Half,All	86.9
Chance	_	_	0.0

- Some gap (1-2 attributes) between train-test models is okay:

- Still 80% relative accuracy (R.Acc) of the iid case (100% R.Acc).





- Intellectual properties in DNNs may not be 100% safe under black-box.

- Novel methods for exposing internals. We craft adversarial inputs that make DNNs confess their secrets.

- Exposed internals → greater vulnerability to adversarial examples.

ImageNet Setup & Results

ImageNet Setup

Family	#Members	#layers	Top-5 error
SqueezeNet	2	26	19.5 ± 0.1
VGG	4	$11 \sim 19$	10.2 ± 1.1
VGG- B N	4	$11 \sim 19$	9.2 ± 1.0
ResNet	5	$21 \sim 156$	8.4 ± 2.5
DenseNet	4	$121{\sim}201$	7.0 ± 0.8

- Task:

Leave-one-out family prediction.

- High intradiversity and inter-similarity.

Family Prediction

kennen-o					
#queries	Acc(%)				
0	20.0				
1	74.2				
10	90.4				
100	90.4				
1000	94.8				

Transferrability of Adv. Examples

	Target family					
Gen	S	V	В	R	D	
Clean	38	32	28	30	29	
S	64	49	45	39	35	
V	62	96	96	57	52	
В	50	85	95	47	44	
R	64	72	78	87	77	
D	58	63	70	76	90	
Ens	70	93	93	75	80	

- Black-box ImageNet classifier families are reliably exposed.

- Adversarial examples transfer better within family (diagonals).

Find out family; then attack.

Scenario	Targetted models	Miscls.(%)
White box	GT model	100.0
Black box, family known	GT family	86.2
Black box, family exposed	Predicted family	85.7
Black box	S,V,B,R,D	82.2

- Misclassification rates (100-accuracy).

- Internal exposure can make the transfer based

attack more targetted and effective ($82.2\% \rightarrow 85.7\%$).