Adversarial Ranking Attack and Defense	00
Supplementary Material	001
Supplementary Materia	002
	003
Anonymous ECCV submission	004
v	005
Paper ID 2274	006
	007
	300
A Adversarial Example Visualization	009
	010
Some adversarial ranking examples are presented in this section. Every fi	gure ⁰¹¹
contains three rows of pictures. The first row shows $c, r, \tilde{c} = c + r$ for CA,	or q , 012
$r, \tilde{q} = q + r$ for QA. The second row shows the query and the original retr	ieval ⁰¹³
results, as well as the chosen candidate c and its immediately adjacent ca	andi- ⁰¹⁴
dates. The third row shows the effects of the attack on the ranking list,	<i>i.e.</i> , 015
either the chosen candidate c is replaced with \tilde{c} for CA, or the query q is repl	aced ⁰¹⁶
with \tilde{q} for QA. The digits above every picture is the value of Rank _X (q, \cdot) –	-1 ¹ . ⁰¹⁷
Pictures with a " \star " mark on the top-left corner are adversarial exam	ples. ⁰¹⁸
The " \leftarrow " indicates the chosen candidate whose rank will be raised. The	"→" ⁰¹⁹
indicates the chosen candidate whose rank will be lowered.	020
	021
A 1 MNIST Detect	022
A.1 MINISI Dataset	023
CA+ . See Fig. 1,2. CA- . See Fig. 3,4.	024
QA+ . See Fig. 5,6. QA- . See Fig. 7,8.	025
	026
A 2 Eachian MNIST	027
A.2 Fashion-WINIST	028
CA+ . See Fig. 9,10. CA- . See Fig. 11,12.	029
QA+ . See Fig. 13,14. QA- . See Fig. 15,16.	030
	031
A 2 Stanford Online Dreducts	032
A.3 Stanford Unline Products	033
CA+ , See Fig. 17.18, CA- , See Fig. 19.20.	034
QA +. See Fig. 21.22. QA- . See Fig. 23.24.	035
CA+(w=2). See Fig. 25, 26. $QA+(m=2)$. See Fig. 27, 28.	036
CA+(w=10). See Fig. 29.	037
, , ,	038
	039
	040
	041
	042
1 In mathematical context the rank counts from 1, but in implementation it co	ounts ⁰⁴³
from 0 instead. Hence the offset	044

from 0 instead. Hence the offset.



Fig. 1. CA+ on MNIST. Example 1. For query "9", the candidate "7" is ranked at the 2576-th position in the original ranking list (row 2). After adversarial ranking attack, the rank of the perturbed candidate "7" is raised to the 1-st position (row 3). The original candidate "7", the learned perturbation r, and the perturbed candidate "7" are illustrated in the row 1.













 $\overline{7}$

















16 ECCV-20 submission ID 2274



720 A.4 "Imperceptibility"

⁷²¹Larger ε increases attacking effectiveness, at the cost of making the adversarial ⁷²³perturbation more perceptible. Following previous works such as [2], we use $\varepsilon = 0.3$ as the maximum attack strength.

See Fig 30 for CA+ with $\varepsilon = 0.3$. See Fig 31 for CA+ with $\varepsilon = 0.1$. See Fig 32 for CA+ with $\varepsilon = 0.03$. See Fig 33 for CA+ with $\varepsilon = 0.01$. С $\tilde{c} = c + r$ 7176 7177 7178 a Ω а

Fig. 30. "Imperceptibility": CA+ on MNIST. $\varepsilon = 0.3$.

A.5 Semantic-Preserving for QA

Conducting Query Attack without preserving the query semantics will often lead
to irrelevant retrieval results at the top of the ranking list, as shown in Fig. 34,
which raises red flags and possibly reveals the attack. Therefore, the value of the
Semantics-Preserving term in QA is substantial, as it keeps the retrieval results
as "normal" as possible while achieving the attacking goal.





B Complete Results of Attack & Defense

Some experiments and details on MNIST, Fashion-MNIST and SOP datasets are omitted in the manuscript due to limited space. In this section, we present the complete experimental results on these datasets, including the average rank of C_{SP} during QA. Besides, we also conduct attacking experiments on ranking models trained with different combinations of loss functions and distance metrics. For brevity, we denote the (Cosine distance, Triplet loss) setting as (CT), the (Euclidean distance, Contrastive loss) as (EC). Models trained with our defense method will have a trailing "D" in notation. e.q. the (CT) model with our defense will be denoted as (CTD).

- See Tab. 2 for complete results on MNIST.
- 867 See Tab. 3 for complete results on Fashion-MNIST.
 - See Tab. 1 for results on SOP.

_ ا		CA	+			C.	A-			SP-	QA+			SP-	QA-	
2	w = 1	2	5	10	w = 1	2	5	10	m = 1	2	5	10	m = 1	2	5	10
		_				(ET) E	uclidea	n Distanc	e, Triplet	Loss (R@	1=63.1%)				
0	50	50	50 2.0	50 2.6	1.9	1.9	1.9	1.9	50 48.02	50	50 16.3 1 9	50 25.8 2.3	0.5	0.5	0.5	0.5
0.01	0.0	0.3	1.0	1.5	100.0	100.0	100.0	100.0	1.6, 0.1	3.3, 0.5	10.0, 2.0	19.2, 2.7	68.1, 0.3	52.4, 0.6	36.6, 0.8	30.1, 1.0
0.06	0.0	0.2	1.0	1.5	100.0	100.0	100.0	100.0	1.1, 0.2	2.7 , 0.6	8.8, 1.9	17.6, 3.3	73.8, 0.4	57.9 , 0.7	40.3 , 0.8	32.4 , 1.
0	50	50	50	50	2.0	(ETD) 2.0	Euclid 2.0	ean Di 2 0	stance, Ti 50	nplet Loss 50	, Defensiv 50	e (R@1=4 50	6.4%) 0.5	0.5	0.5	0.5
0.01	7.5	12.2	16.5	18.0	66.4	62.6	59.3	57.8	16.1, 2.1	24.8, 2.7	36.1, 2.7	41.4, 2.5	26.7, 0.6	18.1, 0.7	12.2, 0.7	10.2, 0.
0.03	0.7	4.5	8.7	10.4	91.7	90.2	89.1	88.4	7.9, 2.9	14.5, 4.2	27.2, 5.4	35.6, 5.3	43.4, 1.3	31.7, 1.5	21.9, 1.7	18.1, 1.
0.06	0.1	3.8	7.9	9.7	97.3	96.8	96.4	96.2	6.9, 4.3	12.5, 5.8	24.3, 7.0	33.4, 6.9	51.4, 2.8	39.0 , 3.2	28.0, 3.5	23.5, 3.
						Tat	bie 1	. Ex	perim	ents o	n SOP	datas	et.			

ECCV-20 submission ID 2274

ε	w = 1	CA 2	+	10	w = 1		A-	10	m - 1	SP-	QA+ 5	10	m - 1	SP 2	-QA-	10
	w = 1		0	10	10 - 1	(C)		ine D	m = 1	ontrastive	Loss (B@1-	-98.6%)	1 11 - 1	1 4		10
0	50	50	50	50	2.2	2.2	2.2	2.2	50	50	50	50	0.5	0.5	0.5	0.5
0.01	42.2	43.3	44.6	44.7 26.1	5.6	5.3	5.2	5.1	45.0, 0.5	46.9, 0.5	48.4, 0.5	49.1, 0.5	1.2, 0.5	1.2, 0.5	1.1, 0.5 1	1.1, 0.5
0.03	9.7	12.5	55.4 14.6	15.6	9.3 52.1	9.2 52.2	9.1 52.4	9.1 52.5	18.5, 8.2	41.7, 2.5 27.1, 9.7	45.9, 2.5 38.0, 9.9	47.7, 2.5 43.4, 9.9	5.6, 4.1	5.4, 4.1	5.4, 4.1 5	5.4, 4.1
0.3	5.8	10.0	11.9	12.9	99.0	99.1	99.0	99.1	14.3, 9.7	23.3, 12.5	36.3 , 13.8	42.6 , 13.9	6.1 , 4.4	5.8 , 4.4	5.7, 4.4 5	.7 , 4.4
0	50	50	50	50	2.1	2.1	2.1	osine 2.1	Distance, 50	Triplet Lo 50	ss (R@1=9 50	9.1%) 50	0.5	0.5	0.5	0.5
0.01	44.6	45.4	47.4	47.9	3.4	3.2	3.1	3.1	45.2, 0.0	46.3, 0.0	47.7, 0.0	48.5, 0.0	0.9, 0.0	0.7, 0.0	0.6, 0.0 0	0.6, 0.1
0.03	33.4 12.7	37.3	41.9 24 4	43.9 30.0	6.3	5.9 14.0	5.7 14 8	5.6 14.7	35.6, 0.3 14 4 2 2	39.2, 0.3 21.0.2.2	43.4, 0.3 30.6 2.2	45.8, 0.3	1.9, 0.2	1.4, 0.2	1.1, 0.2 1	.1, 0.2
0.3	2.1	9.1	13.0	17.9	93.9	93.2	93.0	92.9	6.3 , 3.6	11.2 , 5.7	22.5 , 7.7	32.1 , 7.7	8.6, 1.6	6.6 , 1.6	5 5.3 , 1.5 4	
0	FO	50	50	50	1.0	(EC)	Eucli	dean	Distance,	Contrastiv	e Loss (R@	1=99.0%)	0.5	0.5	0.5	0.5
0.01	50 31.9	50 33.5	50 34.6	50 34.9	1.8	1.8	7.9	1.8	50 41.3, 2.0	50	50	50 48.4. 2.5	0.5 2.6, 1.3	0.5	0.5	$\frac{0.5}{2.2, 1.3}$
0.03	15.8	17.4	18.7	19.1	10.7	10.6	10.5	10.5	27.2, 4.8	34.6, 5.4	42.0, 5.6	45.6, 5.8	4.1, 2.5	3.9, 2.5	3.7, 2.5 3	3.7, 2.5
0.1	6.7 4 8	10.0 9 0	12.1 12 1	12.7 12 7	82.1 99.9	81.6	81.7 99 9	82.0	15.7, 9.8	25.4, 12.0 23 9 12 5	37.7, 13.1 36 7 13.6	43.2, 13.1	5.6, 3.2	5.3, 3.2	4.9, 3.2 4	1.9, 3.2
0.0	-1.0	0.0	12.1	12.1	00.0	(E	T) Eu	clidea	in Distanc	e, Triplet I	Loss (R@1=	99.2%)	10.0, 0.2		. 5.6, 5.2 4	
0	50	50	50	50	1.6	1.6	1.6	1.6	50	50	50	50	0.5	0.5	0.5	0.5
$0.01 \\ 0.03$	$\frac{39.0}{23.5}$	40.4 25.6	40.6 26.5	40.8 27.0	3.2	3.0 6.6	2.8	2.8 6.2	45.5, 0.1 36.4, 0.6	47.4, 0.1 40.9, 0.7	48.1, 0.1 45.1, 0.8	49.0, 0.1	1.0, 0.1 2.7, 0.7	0.9, 0.1 2.3, 0.7	0.8, 0.1 0 2.0, 0.7 2	0.8, 0.1 0.0, 0.8
0.1	8.1	10.6	12.1	12.7	13.7	13.3	13.0	12.9	14.1, 4.8	23.4, 5.7	34.6, 6.2	40.8, 6.5	6.2, 1.7	5.1, 1.7	4.4, 1.7 4	1.2, 1.7
0.3	1.8	8.5	10.9	11.7	77.4	75.7	75.0	74.7	7.5 , 5.2	15.0 , 7.4	28.1 , 9.0	36.3 , 8.9	7.4, 1.7	5.8 , 1.7	7 4.9 , 1.7 4	.6 , 1.7
0	50	50	50	50	(CC	D) C	osine I	Distar	nce, Contra 50	astive Loss 50	, Defensive 50	(R@1=97.5 50	5%) 0.5	0.5	0.5	0.5
0.01	49.3	49.0	49.3	49.3	2.2	2.4	2.2	2.3	49.7, 0.0	49.6, 0.0	49.8, 0.0	49.9, 0.0	0.5, 0.0	0.5, 0.0	0.5, 0.0 0	0.5, 0.0
0.03	47.0	47.9	48.1	48.0	2.8	2.7	2.7	2.7	48.3, 0.0	49.1, 0.0	49.1, 0.0	49.4, 0.0	0.6, 0.0	0.5, 0.0	0.5, 0.0 0	0.5, 0.0
0.1	42.5 32.0	43.3 34.2	36.1	44.4 36.7	4.2 7.0	4.0 7.0	5.9 6.5	3.8 6.4	37.4 , 0.6	40 .7,0.1 41.1 , 0.6	48.1 ,0.1 44.9 , 0.5	46.8 ,0.1 47.2 , 0.5	1.9, 0.5	5 1.6 , 0.5	5 1.5 , 0.5 1	.7,0.1 .5,0.5
					((CTD)	Cosin	e Dist	ance, Trip	olet Loss, I	Defensive (R	@1=98.3%)			
0.01	48.9	50 49.3	50 49.4	50 49.5	2.0	2.0	2.0	2.0	50	50 49.5. 0.0	49.5. 0.0	50 49.7. 0.0	0.5	0.5	0.5	0.5 0.5, 0.0
0.03	47.4	48.4	48.6	48.9	2.5	2.5	2.4	2.4	48.0, 0.0	48.5, 0.0	49.2, 0.0	49.5, 0.0	0.6, 0.0	0.6, 0.0	0.5, 0.0 0	0.5, 0.0
0.1	42.4 30.7	44.2 34 5	45.9	46.7	3.8	3.6 6 7	3.5 6 5	3.4 6.5	43.2, 0.1	45.0, 0.1 37 2 0 5	47.4, 0.1	48.2, 0.1	1.0, 0.1	0.8, 0.1	0.7, 0.1 0	0.7, 0.1 5 0 4
0.0	00.1	04.0	00.1	10.1	(ECI) Euc	lidean	Dista	ance, Cont	rastive Los	ss, Defensiv	e (R@1=97	7.9%)	1.0, 0	1.0, 0.4 1	
0	50	50 47.6	50	50	1.3	1.3	1.3	1.3	50	50	50	50	0.5	0.5	0.5	0.5
0.01	42.7	43.6	40.0 44.0	44.2	4.5	4.2 4.2	4.0	4.0	46.3, 0.6	48.1, 0.1	49.4, 0.1 48.8, 0.6	49.2, 0.6	1.8, 0.6	1.6, 0.1	1.5, 0.6	.4, 0.1
0.1	31.7	33.7	34.9	35.3	10.0	9.6	9.4	9.3	39.2, 2.8	43.2, 2.9	46.6, 2.9	47.9, 2.9	3.3, 1.3	2.9, 1.3	2.7, 1.3 2	2.6, 1.3
0.3	19.6	23.0	25.4	26.3	35.6 (E)	35.2 FD) F	35.7 Juclide	36.0 an Di	27.3, 7.1 stance, Tr	34.4, 7.4 iplet Loss	42.2, 7.1 Defensive (45.4, 6.9 B@1=99.0	4.5 , 1.8 %)	3.9, 1.8	3.6, 1.8 3	.4 , 1.8
0	50	50	50	50	1.4	1.4	1.4	1.4	50	50	50	50	0.5	0.5	0.5	0.5
0.01	47.5	48.2	48.1	48.2 44 5	1.7	1.7	1.6	1.6	48.5, 0.0	48.8, 0.0	49.4, 0.0	49.7, 0.0	0.6, 0.0	0.6, 0.0	0.5, 0.0 0	0.5, 0.0
0.1	49.4 29.8	40.0 31.3	32.6	44.3 33.0	5.7	$\frac{2.3}{5.2}$	5.0	4.9	39.6, 0.3	42.9, 0.3	45.9, 0.1	47.9, 0.1	2.1, 0.4	1.7, 0.4	1.5, 0.4 1	.4, 0.4
0.3	10.8	13.2	15.0	15.6	14.4	13.9	13.5	13.4	19.7 , 3.2	28.1, 3.7	37.4 , 4.1	42.6 , 4.3	6.5 , 1.7	5.3 , 1.7	4.7 , 1.8 4	1.5 1.8
[ab	ole 2	2. C	omŗ	olete	e Exj	peri	men	t R	esults of	on MN	IST. Tł	ne first	value	for e	ach QA	(<i>i.e</i>
5P-0	QA)	\exp	erir	nen	t res	ult i	s the	e av	erage r	ank of	the chos	sen can	didat	es, wh	ile the	othe
alu	e is	the	ave	rage	e ran	k of	the	$C_{\rm S}$	P used	for SP	2.					
				0				2								

945

946

947

948

949

950

- 945
- 145
- 946
- 947
- 948
- 949
- 950

981

982

983

984

985

986

987

988

989

951 951 SP-QA+ SP-QA-CA+ CA-952 952 e w = 12 5 10 10 10 w = 15 m - 1m - 1953 953 (CC) Cosine Distance, Contrastive Loss (R@1=88.7%) 0 50 50 50 50 2.0 2.0 20 2.0 50 50 50 50 0.5 0.5 0.5 0.5 954 954 0.01 29.8 34.2 34.7 9.3 9.0 40.3. 1.6 2.5. 1.4 0.8 8.8 36.1 1.6 44.8 1.5 46.9 1.5 31 14 28 1 4 25 14 955 0.03 12.5 16.4 19.1 20.3 46.0 45.9 44.7 44 5 18.5. 4.6 25.7, 5.6 35.6, 5.9 41.6.5.9 4.9. 2.4 4.4. 2.4 4.0. 2.4 3.9. 2.5 955 4.8 96.0 0.1 10.1 13.4 14.9 96.0 96.0 96.0 103 70 17.3, 9.7 30.0, 11.1 38.2, 11.7 7.1. 3.7 6.4. 3.7 5937 57 37 956 0.3 3.7 9.6 12.8 14.2 100.0 100.0 100.0 100.0 9.0 6.3 15.7. 9.3 28.5. 11.3 37.5. 11.6 7.2. 3.5 6.4. 3.5 5.8. 3.5 5.6. 3.4 956 Cosine Distance, Triplet Loss (R@1=88.8%) (CT)957 957 0 50 0.5 0.5 0.5 0.5 50 50 50 1.0 1.0 1.0 1.0 50 50 50 50 0.01 36.6 30.0 43.2 44.8 5.6 1.9 18 39.4, 0.2 42.0.0.2 45.3, 0.2 47.1, 0.2 2.1, 0.2 1.6, 0.2 1.2, 0.2 11 0 2 958 958 0.03 19.7 25.4 31.7 35.6 14.8 14.4 14.3 21.7. 1.5 28.2, 1.6 35.7. 1.7 5.6. 0.8 4.1. 0.8 3.3. 0.8 15.5 40 6 1 7 2007 87.9 86.7 86.3 67 19 60 18 0.1 37 10 5 17 3 99 7 86.3 71 24 194 45 236 60 32.5 6.8 10.0 1.0 83 10 959 959 0.3 1.3 9.4 $16.0\ 21.5$ 100.0 100.0 100.0 100.0 6.3. 3.0 10.8. 5.2 21.8. 7.8 31.7.83 12.6. 1.9 9.4. 1.9 7.5, 1.9 6.6, 1.8 (EC) Euclidean Distance, Contrastive Loss (R@1=87.6%) 960 960 0 50 50 50 50 1.3 1.2 1.3 50 50 50 50 0.5 0.5 0.5 0.5 1.3 961 0.01 20.0 23.6 25.7 26.412.311.7 11.3 30.6 3.7 36.7, 3.8 42.8 3.9 45.8 3.8 38 16 33 16 30 16 2916 961 0.03 7.0 11.5 14.8 16.0 74.1 72.9 71.8 71.4 157 80 25.4, 9.3 36.9. 9.1 42.4. 8.8 4.6. 2.0 4.0, 2.0 3.6. 2.0 36.21 962 12.2 17.0 18.7 100.0 100.0 14.2, 9.3 22.6, 12.3 34.8, 13.2 41.2, 13.0 8.3. 4.7 7.8, 4.8 7.3, 4.8 7.1, 4.8 962 0.1 6.2 100.0 100.0 8.4. 5.2 7 9 5.0 7.5 4.9 0.3 57 12 0 16 5 18 2 100 0 100 0 100 0 100 0 12.8. 9.0 21.1. 11.9 33.8. 13.5 40.4. 13.4 9.1. 5.1 963 963 (ET) Euclide n Distance, Triplet Loss (R@1=88.3%) 0 0.50.50.50.55050 5050 1.51.51.550 50 50 50 964 964 0.01 33.2 34.8 34.7 36.1 6.5 417 03 44 4 0 3 47.3 0.3 48.4 0.3 25.03 19.03 16.03 14 03 17.2 0.0314.0 17.6 20.3 21.3 18.4 16.4 16.0 22.2. 2.5 30.0. 2.7 39.4. 2.7 43.6. 2.6 6.6. 1.3 5.1. 1.4 4.2. 1.4 39 14 965 965 9.514.0 15.6 88.3 86.7 85.2 84.5 8.0, 3.9 15.0, 6.3 27.6, 7.9 36.5, 8.1 10.5, 1.9 0.1 1.78.0. 1.9 6.4. 2.0 6.0. 2.0 0.3 0.3 ۹n 13 8 15 5 100.0 100.0 100.0 100.0 6.7, 3.3 12.6, 5.7 **25.4**, 7.9 **34.8**, 8.3 11 0 10 80 10 7.1, 1.9 6.4, 1.9 966 966 (CCD) Cosine Distance, Contrastive Loss, Defensive (R@1=82.2%) 967 967 0.50 50 50 50 50 2.0 2.02.0 2.050 50 50 50 0.50.5 0.547.7 47.8 47.9 48.7. 0.0 49.2, 0.0 0.01 05.00 47.0 23 23 48.9 0.0 49.6. 0.0 0.6. 0.0 0.5.0.0 05.00 968 968 0.0342.2 43.3 44.1 44.3 3.1 3.0 2.928 45.2, 0.0 47.2, 0.0 48.0, 0.0 49.1, 0.0 0.8, 0.0 0.7, 0.00.6, 0.0 0.6, 0.0 969 0.129.131.4 32.9 33.8 8.0 7.16.46.234.7, 0.139.5, 0.143.8, 0.1 46.5, 0.1 3.3, 0.1 2.1, 0.11.5, 0.1 1.3, 0.1 969 0.3 11.8 14.8 17.4 18.4 28.725.523.122.413.3. 0.7 20.0. 0.9 31.0. 1.0 **38.3**, 1.0 **21.3**, 1.0 **14.3**, 1.0 **10.3**, 1.0 **8.6**, 1.1 970 970 (CTD) Cosine Distance, Triplet Loss, Defensive (R@1=79.6%) 0 50 50 50 50 1.2 50 50 50 50 0.5 0.50.5 0.5 1.21.21.2971 971 0.01 49.9, 0.0 49.9, 0.0 48.9 48.9 49.3 49.3 14 14 14 14 50.0 0.0 0.5, 0.0 0.5, 0.00.5, 0.0 0.03 47.1 47 9 48 3 48 3 483 0.0 491 0.0 49.5, 0.0 498 0.0 0700 06 00 06.00 06.00 2.019 18 18 972 972 0.142.443.5 44.5 44.8 4.642 4.03.9 45.4, 0.3 47.2, 0.2 48.7, 0.2 49.2, 0.2 1.4. 0.2 1.2, 0.2 1.1. 0.2 1.0. 0.2 0.3 32.535.4 37.5 38.2 11.210.510.1 10.0 39.3. 1.5 42.6. 1.5 46.5. 1.3 47.8. 1.3 3.9. 1.4 **3.3**, 1.4 **3.0**, 1.4 **2.9**, 1.4 973 973 (ECD) Euclidean Distance, Contrastive Loss, Defensive (R@1=80.4%) 0 50 0.5 50 50 50 14 14 14 14 50 50 50 50 0.5 0.50.5974 974 0.01 45.3 45.8 46.3 46.6 3.0 2.8 48.0, 0.2 48.7, 0.2 49.3, 0.3 49.7, 0.3 1.2, 0.2 1.0, 0.2 0.9, 0.2 09.02 975 975 0.03 37.6 39.5 40.4 40.7 7.5 7.06.8 67 43.9. 1.5 46.4. 1.6 48.4. 1.7 49.1. 1.9 2.7. 0.7 2.2. 0.7 1.9. 0.7 19.07 24.9 27.8 28.6 32.6 32.132.0 32.5 36.2, 7.2 43.0, 7.4 3.8. 1.6 3.7 1.6 20.328.4. 6.6 45.9.7.4 5.3. 1.5 4.4, 1.6 0.1 976 976 0.3 7.3 16.0 21.1 22.5 78.0 79.0 80.3 81 2 14.0, 10.2 24.6, 11.9 36.5, 11.8 42.2, 11.2 7.0, 1.9 5.7, 1.9 4.9, 1.9 4.6, 1.9 (ETD) Euclidean Distance, Triplet Loss, Defensive (R@1=84.4%) 977 977 0 50 50 5050 1.3 1.3 1.3 1.3 505050500.50.50.50.50.8. 0.0 0.01 46.046 7 46 9 46 7 2.0 19 1.8 18 48.1, 0.0 491 0.0 493 00 496 00 0.7, 0.0 06.00 06.01 978 978 0.03 38.0 40.6 41.2 41.4 4.3 3.8 3.5 3.4 45.2.0.2 46.2.0.1 48.1.0.1 48.9. 0.1 1.9.0.21.5. 0.2 1.2.0.2 1102 0.1 23.026.3 28.0 28.8 15.914.9 14.213.9 31.8.1.8 38.0, 1.8 43.3, 1.6 46.2, 1.6 5.6. 1.0 4.3. 1.0 3.6. 1.0 3.3. 1.0 979 979 50.049.27.113.9 18.4 19.9 54.151.8 10.6, 5.6 19.7, 7.1 32.5, 8.0 39.6, 7.6 9.2, 1.6 7.0, 1.6 5.6, 1.6 5.1, 1.6 0.3 980 980

- 982 983
- 984
- 985
- 986
- 987
- 988
- 989
 - 89

Complete Results on Transferability C

C_{1} Fashion-MNIST

In addition to the transferability experiment on MNIST dataset, we also conduct the same transferability experiment on the Fashion-MNIST dataset, as shown in Tab. 4.

CA+ T	ransfer (B	lack Box), v	v = 1
To	LeNet	C2F1	Res18
LeNet	$50 \rightarrow 16.0$	41.0	44.3
C2F1	38.6	$50 \rightarrow 1.3$	40.3
Res18	39.2	34.3	$50 \rightarrow 1.7$
CA- T	ransfer (Bl	ack Box), u	v = 1
To	LeNet	C2F1	Res18
LeNet	$2.5 { ightarrow} 84.3$	$1.9 \rightarrow 8.1$	$1.6 {\rightarrow} 6.0$
C2F1	$2.5 \rightarrow 7.8$	$1.9{ o}100.0$	$1.7 \rightarrow 7.7$
Res18	$2.5 \rightarrow 9.5$	$1.9 \rightarrow 14.4$	$1.7 { ightarrow} 80.0$
SP-QA+	Transfer (Black Box).	m = 1
SP-QA+ To From	Transfer (LeNet	Black Box). C2F1	m = 1 Res18
SP-QA+ To From LeNet	Transfer (LeNet $50 \rightarrow 18.0$	Black Box), C2F1 47.2	m = 1 Res18 49.3
SP-QA+ To From LeNet C2F1	Transfer (LeNet $50 \rightarrow 18.0$ 48.1	Black Box), C2F1 47.2 $50 \rightarrow 6.4$	m = 1 Res18 49.3 49.2
SP-QA+ To From LeNet C2F1 Res18	Transfer (LeNet $50 \rightarrow 18.0$ 48.1 48.1	Black Box), C2F1 47.2 $50 \rightarrow 6.4$ 44.8	m = 1 Res18 49.3 49.2 50 \rightarrow 13.7
SP-QA+ To From LeNet C2F1 Res18 SP-QA-	Transfer (LeNet $50 \rightarrow 18.0$ 48.1 48.1 Transfer (I	Black Box), C2F1 47.2 $50 \rightarrow 6.4$ 44.8 Black Box),	m = 1 Res18 49.3 49.2 50 \rightarrow 13.7 m = 1
SP-QA+ To From LeNet C2F1 Res18 SP-QA- To From	Transfer (LeNet $50 \rightarrow 18.0$ 48.1 48.1 Transfer (I LeNet	Black Box), C2F1 47.2 50→6.4 44.8 Black Box), C2F1	m = 1 Res18 49.3 49.2 $50 \rightarrow 13.7$ m = 1 Res18
SP-QA+ To From LeNet C2F1 Res18 SP-QA- To From LeNet	Transfer (LeNet $50 \rightarrow 18.0$ 48.1 48.1 Transfer (I LeNet $0.5 \rightarrow 13.5$	Black Box), C2F1 47.2 $50 \rightarrow 6.4$ 44.8 Black Box), C2F1 $0.5 \rightarrow 1.7$	m = 1 Res18 49.3 49.2 $50 \rightarrow 13.7$ m = 1 Res18 $0.5 \rightarrow 1.5$
SP-QA+ To From LeNet C2F1 Res18 SP-QA- To From LeNet C2F1	Transfer (LeNet $50 \rightarrow 18.0$ 48.1 48.1 Transfer (I LeNet $0.5 \rightarrow 13.5$ $0.5 \rightarrow 1.1$	Black Box), C2F1 47.2 $50 \rightarrow 6.4$ 44.8 Black Box), C2F1 $0.5 \rightarrow 1.7$ $0.5 \rightarrow 12.5$	m = 1 Res18 49.3 49.2 50→13.7 m = 1 Res18 0.5→1.5 0.5→1.6
SP-QA+ To From LeNet C2F1 Res18 SP-QA- To From LeNet C2F1 Res18	$\begin{array}{c} {\rm Transfer} \ (\\ {\rm LeNet} \\ 50{\rightarrow}18.0 \\ 48.1 \\ 48.1 \\ {\rm Transfer} \ ({\rm I} \\ {\rm LeNet} \\ 0.5{\rightarrow}13.5 \\ 0.5{\rightarrow}1.1 \\ 0.5{\rightarrow}0.9 \end{array}$	Black Box), C2F1 47.2 $50 \rightarrow 6.4$ 44.8 Black Box), C2F1 $0.5 \rightarrow 1.7$ $0.5 \rightarrow 1.2.5$ $0.5 \rightarrow 1.3$	m = 1 Res18 49.3 49.2 50→13.7 m = 1 Res18 0.5→1.5 0.5→1.6 0.5→8.0

 Table 4. Transferability experiment on Fashion-MNIST dataset.

"Self-Transfer" Attack on MNIST C.2

In the adversarial ranking example transferability experiments, we transfer ad-versarial examples between neural networks with different architectures. In this section, we transfer adversarial examples between neural networks based on the same architecture but with different parameters.

We train three vanilla C2F1 models (denoted as C2F1-1, C2F1-2, and C2f1-3) on MNIST dataset and two models with our defense (denoted as C2F1-D1, C2F1-D2). All these models have exactly the same architecture, but different pa-rameters. Transferability experimental results between these models are present in Tab. 5.

1035		1035
1036		1036
1037		1037
1038		1038
1039		1039
1040		1040
1041		1041
1042		1042
1043		1043
1044		1044
1045		1045
1046	CA+ Transfer (Black Box), $w = 1$	1046
1040	To C2F1-1 C2F1-2 C2F1-3 C2F1-D1 C2F1-D2	1040
1047	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	1047
1040	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1040
1050	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1050
1050	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1050
1051	To $C2F1-1$ $C2F1-2$ $C2F1-3$ $C2F1-D1$ $C2F1-D2$	1051
1052	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1052
1053	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	1053
1054	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	1054
1055	C2F1-D2 $ 2.1 \rightarrow 9.2 2.2 \rightarrow 8.4 2.1 \rightarrow 9.1 2.0 \rightarrow 4.9 1.9 \rightarrow 7.0$	1055
1056	SP-QA+ Transfer (Black Box), $m = 1$	1056
1057	From C2F1-1 C2F1-2 C2F1-3 C2F1-D1 C2F1-D2	1057
1058	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1058
1059	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1059
1060	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1060
1061	SP-QA- Transfer (Black Box), $m = 1$	1061
1062	From C2F1-1 C2F1-2 C2F1-3 C2F1-D1 C2F1-D2	1062
1063	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1063
1064	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	1064
1065	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	1065
1066	Table 5. Transfer experiment between models based on the same architecture but with	1066
1067	different parameters.	1067
1068		1068
1069		1069
1070		1070
1071		1071
1072		1072
1073		1073
1074		1074
1075		1075
1076		1076
1077		1077
1078		1078
1079		1079

1080 D Complete Results for Universal Perturbation

1082 D.1 MNIST

See Tab. 6.

Model	I-CA+	(w = 1)	I-CA- ((w = 1)	I-QA+	(m = 1)	I-QA- ((m = 1)
Model	Seen	Unseen	Seen	Unseen	Seen	Unseen	Seen	Unseen
(CC)	$50 \rightarrow 13.7$	$50 \rightarrow 14.1$	$0.6 \rightarrow 37.6$	$0.6 \rightarrow 33.4$	$50 \rightarrow 18.1$	$50 \rightarrow 18.9$	$2.4 \rightarrow 41.6$	$2.4 \rightarrow 39.3$
(CT)	$50 \rightarrow 18.1$	$50 \rightarrow 18.5$	$0.6 \rightarrow 9.5$	$0.7 \rightarrow 9.4$	$50 \rightarrow 20.5$	$50 \rightarrow 21.0$	$2.1 \rightarrow 7.6$	$2.2 \rightarrow 7.4$
(EC)	$50 \rightarrow 9.1$	$50 \rightarrow 10.3$	$1.9 \rightarrow 94.6$	$2.0 \rightarrow 79.4$	$50 \rightarrow 3.9$	$50 \rightarrow 6.9$	$3.3 \rightarrow 87.1$	$3.5 \rightarrow 74.$
(ET)	$50 \rightarrow 11.5$	$50 \rightarrow 12.6$	$21 \rightarrow 10.6$	$21 \rightarrow 96$	$50 \rightarrow 21.6$	$50 \rightarrow 23.6$	$3.2 \rightarrow 28.6$	$3.2 \rightarrow 10$

Table 6. Universal Adversarial Ranking Perturbation on MNIST. Each pair of result presents the original rank of chosen candidate(s), and the rank of the chosen candidate(s) after adding adversarial perturbation to the candidate or to the query. "Seen" samples are those used for generating the universal perturbation, while "Unseen" samples are a set of other non-overlapping samples.

D.2 Fashion-MNIST

As shown in Tab. 7, Image-agnostic adversarial perturbation has better effect on seen data from Fashion-MNIST, but the gap between the effect on seen samples and that on unseen samples is slightly larger, which may due to the higher intra-class variance of Fashion-MNIST than MNIST.

Model	I-CA+	(w = 1)	I-CA-	(w = 1)	I-QA+	(m = 1)	I-QA- ((m = 1)
Model	Seen	Unseen	Seen	Unseen	Seen	Unseen	Seen	Unseen
(CC)	$ 50 \rightarrow 7.0$	$ 50 \rightarrow 7.3$	$0.6 \rightarrow 91.8$	$0.6 \rightarrow 84.9$	$ 50 \rightarrow 4.4$	$50 \rightarrow 4.9$	$2.1 \rightarrow 87.5$	$2.1 \rightarrow 84.4$
(CT)	$50 \rightarrow 9.8$	$50 \rightarrow 9.9$	$0.6 \rightarrow 72.3$	$0.6 \rightarrow 69.7$	$50 \rightarrow 8.2$	$50 \rightarrow 8.4$	$2.0 \rightarrow 76.3$	$2.0 \rightarrow 72.9$
(EC)	$50 \rightarrow 5.8$ $50 \rightarrow 5.7$	$50 \rightarrow 9.6$ $50 \rightarrow 8.5$	$2.0 \rightarrow 97.5$ $2.0 \rightarrow 84.4$	$1.9 \rightarrow 83.7$ $1.9 \rightarrow 69.9$	$50 \rightarrow 1.8$ $50 \rightarrow 3.3$	$50 \rightarrow 7.1$ $50 \rightarrow 6.3$	$2.9 \rightarrow 87.9$ $3.1 \rightarrow 88.0$	$2.8 \rightarrow 78.4$ $3.0 \rightarrow 78.0$
Tal	ble 7. In	nage-agn	ostic Adv	versarial P	Perturbat	tion on I	Fashion-N	INIST.

Semantics Preserving for QA & Parameter Search for \mathcal{E} \mathbf{E}

As discussed previously, the Query Attack (\mathbf{OA}) may drastically change the semantics of the query q. To alleviate this problem, the Semantics-Preserving (SP) term is added to the naive **QA** to help preserve the query semantics. Predictably, it is more difficult to perform **QA** with a large \mathcal{E} , as the ranks of $C_{\rm SP}$ are almost not allowed to be changed.

To investigate the actual influence of the balancing parameter ξ , we pro-vide parameter search on it with MNIST dataset. In particular, We set ξ to $0, 10^0, 10^2, 10^4$, and compare their results. Note that when $\xi = 0$, the QA be-comes naive QA as the SP term is eliminated. With a strong SP constant, $e.q.\xi = 10^4$, the semantics of the chosen query is almost not allowed to be changed, hence result in extreme difficulty of attack.

As shown in Tab. 8, setting ξ to 0 could greatly boost the attacking effect, but consequently the ranks of $C_{\rm SP}$ will be drastically changed. In contrast, when ξ is set to the excessive value 10^4 for a perfectly stealth QA, the attack can still raise the rank of chosen candidate from 50% to 37.9% in QA+ with m = 1, or lower the rank of chosen candidate from 0.5% to 1.9% in QA- with m = 1. During these attacks, the ranks of $C_{\rm SP}$ are kept within 0.1 despite of the extreme difficulty. It means the query semantics can be preserved. In practice, we empirically set the parameter ξ as 1 for QA+, or as 10² for QA- for the balance between attack effectiveness and preserving query semantics.

1148	SP-QA+ SP-QA-	1148
1149	$\xi = 1 + 2 + 5 + 10 + m = 1 + 2 + 5 + 10$	1149
1150	(CC) Cosine, Contrastive	1150
1151	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1151
1152	10^2 25.4, 2.0 33.2, 2.1 41.7, 2.1 45.4, 2.1 6.1, 4.4 5.8, 4.4 5.7, 4.4 5.7, 4.4	1152
1152	$\frac{10^{4} 48.1,0.3-49.2,0.2-49.6,0.2-49.8,0.3 -1.2,0.3-1.1,0.3-1,0.3-1,0.3-1,0.3-1,0.3-1,0.3-1,0.3-1,0.3-1,0.3-1,$	1152
1153	(CT) Cosine, Triplet	1153
1154	$\begin{smallmatrix} 0 \\ 0.2, 33.6 \\ 6.3, 23.7 \\ 18.5, 26.5 \\ 29.6, 25.7 \\ 94.1, 89.4 \\ 93.2, 90.3 \\ 92.6, 90.9 \\ 92.3, 91.2 \\ 10^9 \\ 6.2 \\ 2.6 \\ 11.2 \\ 5.7 \\ 29.5 \\ 7.7 \\ 29.5 \\ 7.7 \\ 29.1 \\ 7.7 \\ 5.5 \\ 5.5 \\ 6.5 \\ 4.37 \\ 5.5 \\ 6.5 \\ 20.2 \\ 20.3 \\ 40.4 \\ 40.0 \\ 10^{10} \\ $	1154
1155	10^{-} [0.5, 5.0] 11.2, 5.7] 22.5, 7.7] 52.1, 7.7] 55.5, 55.0] 52.4, 57.0] 50.2, 59.5] 49.4, 40.0] 10^{2} [14.1, 0.6] 20.8, 0.7] 31.2, 0.7] 38.1, 0.7] 8.6, 1.6] 6.6, 1.6] 5.3, 1.5] 4.8, 1.5]	1155
1156	$10^4 \begin{vmatrix} 37.9, \ 0.1 & 42.6, \ 0.1 & 46.3, \ 0.1 & 47.8, \ 0.1 & 1.9, \ 0.1 & 1.4, \ 0.1 & 1.2, \ 0.1 & 1.1, \ 0.1 \end{vmatrix}$	1156
1157	(EC) Euclidean, Contrastive	1157
1157	$0 \ 0.7, 44.5 \ 3.4, 39.9 \ 31.0, 41.8 \ 39.7, 41.6 \ 94.0, 92.6 \ 93.9, 93.0 \ 93.9, 93.4 \ 93.9, 93.5 \ 10^{9} \ 144.0 \ 10^{2} $	1157
1158	10^{-1} 14.3 , 3.4 25.5 , 12.5 30.7 , 15.0 42.0 , 15.0 95.5 , 3.70 35.0 , 37.2 35.2 , 31.1 35.1 , 31.3 10^{-2} 30.1 , 12 , 37.3 , 11 , 43.9 , 0.9 , 46.7 , 0.9 , 5.9 , 32 , 5.4 , 32 , 5.0 , 32 , 4.9 , 3.1	1158
1159	$10^4 \begin{bmatrix} 50.1, \ 0.8 & 50.2, \ 0.8 & 49.8, \ 0.8 & 50.0, \ 0.8 \end{bmatrix} \begin{bmatrix} 1.6, \ 0.7 & 1.6, \ 0.7 & 1.7, \ 0.8 & 1.6, \ 0.8 \end{bmatrix}$	1159
1160	(ET) Euclidean, Triplet	1160
1161	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	1161
1100	10^{-1} 20.3, 5.2^{-1} 20.1, 5.2^{-1} 20.1, 5.0^{-1} 20.3, 5.3^{-1} 20.3, 5.73^{-1} 20.2, 5.4^{-1} 41, $1, 40.0^{-1}$ 41, $2, 40.3^{-1}$ 10 ² 20.3, 0.5^{-2} 28.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.6, 0.5^{-2} 38.7, 0.5^{-2} 38.6, 0.5^{-2} 38.7, 0.5^{-2} 38.6, 0.5^{-2} 38.7, 0.5^{-2} 3	11.00
1162	10^4 46.2, 0.1 48.2, 0.1 49.1, 0.1 49.6, 0.1 1.7, 0.1 1.4, 0.1 1.2, 0.1 1.1, 0.1	1162
1163	Table 8. Parameter search on Semantics-Preserving balancing parameter ξ with	1163
1164	MNIST. We report two two mean ranks in each cell: one for the chosen candidate(s)	1164
1165	C , the other for the $C_{\rm SP}$ used for SP.	1165
1166	, 01	1166
1167		1167
1101		1101
1168		1168
1169		1169





Fig. 36. Weights and bias of the second convolutional layer (conv2) with/without ourdefense.



Fig. 37. Weights and bias of the first fully-connected layer (fc1) with/without our defense.

The adversarial perturbation causes the activation to grow by $w^T r$. We can maximize this increase subject to the max norm constraint on r by assigning $r = \varepsilon \operatorname{sign}(w)$. If w has k dimensions and the average magnitude of an element of the weight vector is h, then the activation will grow by εkh . Since $||r||_{\infty}$ does not grow with the dimensionality of the problem but the change in activation caused by perturbation by r can grow linearly with k, then for high dimensional problems, we can make many infinitesimal changes to the input that add up to one large change to the output.

In our experiments, we analysed the parameters in models trained on MNIST, as shown in Fig. 35, 36, 37. These violin plots suggest that

- the weights in the first convolution layer of the defensive model are closer
 to 0 and have smaller variance than those of the vanilla model. That means
 1259

1260 1261 1262 1263	 the h in the above quotation is decreased with our defense, and it will be harder for the r to incur a large increase in activation εkh. the bias in the first convolution layer of the defensive model tend to be negative values instead of being nearly "zero-mean". These negative bias could help further suppress the increase in activation caused by purturbation 	1260 1261 1262 1263
1265	r.	1204
1266		1266
1267	Therefore, Ian <i>et al.</i> 's [1] theory could explain why our defense works, as the	1267
1268	adversarial perturbation to increase the layer outputs into the local linear area	1268
1269	of ReLU.	1269
1270		1270
1271		1271
1272		1272
1273		1273
1274		1274
1276		1275
1277		1277
1278		1278
1279		1279
1280		1280
1281		1281
1282		1282
1283		1283
1284		1284
1285		1285
1286		1286
1287		1287
1288		1288
1289		1289
1290		1290
1291		1291
1293		1293
1294		1294
1295		1295
1296		1296
1297		1297
1298		1298
1299		1299
1300		1300
1301		1301
1302		1302
1303		1303
1304		1304

1305 G Alternative Attack

1307 G.1 Distance-based Ranking Attack

In order to implement the proposed ranking attack, alternative attacking objectives are possible. Some related works such as Feature Adversary [3] generates untargeted adversarial examples against classifiers by maximizing the distance shift of representation vectors off their original locations. This may inspire an alternative version of CA or QA objective functions which are directly based on distance. For example, such alternative objective for CA+ and QA- could be as follows:

 $r = \mathop{\arg\min}_{r \in \varGamma} \sum_{q \in Q} d(q,c+r)$

 $r = \operatorname*{arg\,max}_{r \in \Gamma} \sum_{c \in C} d(q+r,c).$

(2)

(3)

However, it must be pointed out that our method significantly differs from feature adversary [3]: (1) Feature adversary concerns the *pairwise* similarity of source-target representations, while our image ranking problem concerns the ranking order of multiple candidates; (2) Feature adversary attempts to reduce the ℓ_2 distance as much as possible, while our triplet-like loss attempts to make positive candidates closer to query than the negative ones, which well fits the objective of ranking order optimization. Such *relative* distance optimization be-comes more important since our attack simultaneously involves multiple queries and multiple candidates; (3) The ℓ_2 distance based methods suffer from in-evitable disadvantages. Specifically, distance-based objectives are suboptimal, because they disregard the relative positions among the candidates and queries.

As shown in the top-left part of Fig. 38, the solution set for for distance-based CA+ (the green dotted line) contains suboptimal solutions. Similarly, in the bottom-left part of Fig. 38 the distance-based objective for QA- tends to maximize the sum of distance neglecting the ranking result, and further opti-mization (moving \tilde{q} along the green arrow) will not change the ranking result. In contrast, our proposed inequality-based method does not suffer from these issues, as shown in the top-right and bottom-right parts.

We also implemented such distance-based method and compared it with our triplet-like method, as shown in Tab. 9. Experimental results suggest that our method always outperforms distance-based method by a margin. Especially for QA-, the distance-based objective is very difficult to optimize because the distance-based Semantics-Preserving term contradicts with the other term in the loss function. In summary, distance-based method is not well-suited for our proposed *adversarial ranking attack*, especially in the scenario of QA-.



Fig. 38. Distance-based CA+ (top-left) and QA- (bottom-left) objective functions do not lead to desired attacking result compared to our inequality-based method (top-right and bottom-right). In the top-left diagram for CA+, q_1 and q_2 are the chosen queries, while \tilde{c} is the adversarial candidate found by the distance objective. In the bottom-left diagram of QA-, c_1 and c_2 are the chosen candidates, while \tilde{q} is the adversarial query found by the distance objective. It is noted that the distance objective is not optimal as shown in the top-left part, where the solution set for minimizing the distance objective contains suboptimal results. Distance-based objective may even fail to change the ranking result as shown in the bottom-left part, where optimizing the objective further cannot change the ranking result. In contrast, our proposed inequality-based method does not suffer from these issues, as shown in the top-right and bottom-right parts.

	(CA+			CA				QA	+			QA-	
ε	w = 1	2 5	10	w = 1 2 5 10				m = 1 2 5 10				m = 1 2 5		
0.3	3.0 9	9.5 15.9	22.2	86.0	85.2	84.7 8	84.6	7.4	20.2	34.9	41.7	0.8	0.8	0.8 0
	Ta	able 9.	. Att	ack B	ased	on L	-2 I	Distan	ice L	oss v	vith	MNIS	Т.	

н Alternative Defense

Apart from the defense provided in the manuscript, we also tried some other loss functions for adversarial training. In literature, there is no predominant choice for the adversarial training loss function. The only common trait among these choices is that all of them involve adversarial examples. Inspired by previ-ous works, we also implement some alternative defenses for ranking systems as follows

Straightforward Adaptation of Madry Defense **H.1**

Madry [2] formularised improving neural network classifier robustness as a min-max optimization problem, where the inner maximization seeks to generate adversarial examples that lead to maximum cross-entropy loss $L_{\rm CE}$, while the outer minimization tunes the neural network parameters θ to suppress the cross-entropy loss:

$$\min_{\theta} \left\{ \mathbb{E}_{(x,y)\sim D} \left[\max_{r \in \Gamma} L_{\text{CE}}(x+r,y) \right] \right\}$$
(4)

where (x, y) is a pair of image and ground-truth class label.

Similarly, we follow the idea and use a similar defense for ranking models:

$$\min_{\theta} \left\{ \mathbb{E}_{(q,c_p,c_n)\sim D} \left[\max_{r\in\Gamma} L_{\text{triplet}}(q+r,c_p,c_n) \right] \right\}$$
(5)

where the inner maximization aims to generate strongest adversarial examples that could lead to triplet ranking error and a large loss value, while the outer minimization seeks network parameters that could reduce such error.

However, during experiments, we observe that such defensive loss function always diverges, possibly due to the adversarial examples generated by the inner problem being too "strong". We leave further investigation into this problem for future work.

H.2 Straightforward Adaptation of Ian Defense

Ian [1] proposed the following loss function for adversarial training:

$$L_{\text{Ian}}(x,y) = \alpha L_{\text{CE}}(x,y) \tag{6}$$

$$+ (1 - \alpha)L_{\rm CE}(x + \varepsilon {\rm sign}(\nabla_x L_{\rm CE}(x, y)))$$
(7)

where the first term is a normal Cross-Entropy loss, and the second term is the cross-entropy loss with untargeted adversarial example that aims to increase the loss value. Constant α is a balancing parameter.

When adapted to a deep ranking system, this defense method also suffers from the diverging issue similar to Madry's defense.

1440 H.3 Directly Suppressing Shift Distance

As discussed in the manuscript, another possible adversarial training method could be to directly suppress the maximum shift distance of embedding vectors, *i.e.*:

$$L_{\text{trip-es}} = L_{\text{triplet}}(q, c_p, c_n) + \sum_{x \in \{q, c_p, c_n\}} \left(\max_{r \in \Gamma} d(x+r, x) \right), \tag{8}$$

where the most severe distance shift incurred by adversarial perturbations is explicitly suppressed, in addition to a standard ranking loss term.

Experimental results (Tab. 10) show that cosine distance-based ranking models are more robust with this defense. However, we note that this loss may numerically explode on an Euclidean distance-based embedding model, as a strong adversary can gradually cause very large embedding shift distance. To mitigate such divergence issue, we also tried to add a balancing parameter to greatly scale down the second term of the loss function, but the instability problem was still not alleviated.

¹⁴⁵⁵ Due to the lack of universality, we leave this alternative defense in supple-¹⁴⁵⁶ mentary material as a pure discussion, and possible improvements as future ¹⁴⁵⁷ work.

					i						(TD)			1	ap	~ .	
ε	w = 1	CA 2	+	10	w = 1	CA	5	10	<i>m</i> -	- 1	2 SP-0	JA+	10	m - 1	SP-	QA-	10
_	100 - 1			10	C(CD)	Casia				- 1		a Dafanai	10 	05.207)	1 -	0	1 10
0	50	50	50	50	12	1.2	1e Di 1 2	$\frac{12}{12}$	ce, U 50	ontra)	astive Los 50	s, Derensi 50	ve (R@1= 50	95.3%)	0.5	0.5	0.5
0.01	49.1	49.3	49.6	49.5	1.3	1.3	1.3	1.3	49.7,	0.0	49.9, 0.0	50.0, 0.0	49.8, 0.0	0.5, 0.0	0.5, 0.0	0.5, 0.0	0.5, 0.0
0.03	48.0	48.2	48.5	48.6	1.6	1.5	1.5	1.5	48.7,	0.0	49.2, 0.0	49.7, 0.0	49.7, 0.0	0.6, 0.0	0.5, 0.0	0.5, 0.0	0.5, 0.0
0.1	43.1	44.4	45.2	45.5	2.4	2.3	2.1	2.1	45.3,	0.1	47.4, 0.1	48.6, 0.1	49.4, 0.1	0.8, 0.1	0.7, 0.1	0.6, 0.1	0.6, 0.1
0.3	33.3	35.8	37.4	38.0	5.0	5.1	4.8	4.7	38.2,	0.3	42.3, 0.3	45.7, 0.3	47.0, 0.3	2.1, 0.4	1.7, 0.4	1.5, 0.4	1.5, 0.4
0	50	50	50	50	(CTI	D) Co	sine	Dist	ance,	Trip	50 50	Defensive 50	(R@1=97	.4%)	0.5	0.5	0.5
0.01	49.3	49.6	49.5	49.7	1.6	1.6	1.6	1.6	49.5.	,	49.7. 0.0	50.0, 0.0	50.0, 0.0	0.5	0.5, 0.0	0.5, 0.0	0.5, 0.0
0.03	48.2	48.3	48.8	48.8	1.8	1.8	1.8	1.7	49.3,	0.0	48.9, 0.0	49.5, 0.0	49.7, 0.0	0.6, 0.0	0.5, 0.0	0.5, 0.0	0.5, 0.0
0.1	44.7	45.4	46.3	46.4	2.8	2.6	2.5	2.5	46.3,	0.1	47.4, 0.1	48.6, 0.1	49.2, 0.1	0.7, 0.1	0.7, 0.1	0.6, 0.1	0.6, 0.1
0.3	35.5	38.5	40.3	40.9	5.7	5.3	5.1	5.0	39.3,	0.4	43.1 , 0.3	46.0, 0.3	47.5, 0.3	1.8, 0.4	1.6 , 0.4	1.4, 0.4	1.4, 0.4
Tab	le 1	J. 1	he	loss	that	di	rect	ly	sup	pre	esses ei	nbeddi	ing shif	t dist	ance	with I	ANIS'.
data	set.																

References	1485
1 Coodfollow II Shlong I Szorody C · Explaining and harpagging adversarial or	1486
amples arXiv preprint arXiv:1412.6572 (2014)	1487
2. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning	1488
models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083 (2017)	1489
3. Sabour, S., Cao, Y., Faghri, F., Fleet, D.J.: Adversarial manipulation of deep rep-	1490
resentations. arXiv preprint arXiv:1511.05122 (2015)	1491
	1492
	1494
	1495
	1496
	1497
	1498
	1499
	1500
	1501
	1502
	1503
	1504
	1505
	1506
	1507
	1508
	1509
	1510
	1511
	1512
	1513
	1514
	1515
	1510
	1517
	1510
	1519
	1520
	1522
	1523
	1524
	1525
	1506
	1520
	1520
	1520 1527 1528