

### Motivation

Current backdoor attacks are limited to either the spatial domain or the latent space domain.



- Since existing defenses are based on that prior, we suspect that those defenses would fail in the frequency domain.
- Frequency based attacks have been proven to be successful for inference time adversarial attacks, so it would be worth a shot to see if they are also successful against backdoor attacks!



# Contributions

- We propose a backdoor attack that utilizes Fourier heatmaps to design a sophisticated backdoor poisoning attack in the frequency domain.
- Unlike previous spatial attacks, our frequency-based attack is completely imperceptible and bypasses spatial defenses.
- We also show two potential ways to defend against frequencybased backdoor attacks and possible ways for the attacker to bypass these defenses.

# **Check Your Other Door! Creating Backdoor Attacks in the Frequency Domain**

Hasan Abed Al Kader Hammoud, Bernard Ghanem

King Abdullah University of Science and Technology (KAUST)



# <u>III.</u>

Metric	Ratio	SIG	Refool	SPM	WaNet	FIBA	FTrojan	Poison Ink	Ours
CDA/ASR	3%	89.74 / 99.23	89.20 / 87.16	88.89 / 58.53	<u>91.86</u> / 32.86	90.92 / 90.10	91.31 / <b>99.99</b>	89.65 / 94.22	<b>92.31</b> / <u>99.43</u>
	5%	89.64 / 99.47	89.16 / 89.79	88.90 / 57.69	91.47 / 88.15	90.69 / 95.06	<u>91.64 / 99.10</u>	89.69 / 93.58	91.88 / 99.88
	10%	89.45 / 99.40	88.80 / 92.80	89.07 / 57.33	<u>91.22</u> / 96.96	90.41 / 95.86	90.93 / <b>100.00</b>	89.47 / 93.67	<b>92.10</b> / <u>99.97</u>
<b>PSNR↑/LPIPS</b> ↓		25.12 / 0.0400	19.38 / 0.0397	38.94 / <u>0.0001</u>	31.53 / 0.0047	19.40 / 0.0180	41.01 / <u>0.0001</u>	<u>42.95</u> / <u>0.0001</u>	43.15 / 0.00001

Comparison between the Proposed Attack and Backdoor Attacks in the Literature. Our frequency-based attack achieves SOTA ASR, CDA, PSNR, and LPIPS metrics. The results shown are for VGG19 trained on CIFAR10. Legend: First Best, Second Best

	<b>Poisoning Rate</b>	CDA(%)	ASR(%)
CIFAR10	0.0%	93.92	-
	0.1%	94.00	1.54
	0.2%	94.14	72.31
	0.4%	94.20	85.05
	1.0%	94.38	99.44
	3.0%	94.31	99.79
100	0.0%	75.95	-
	0.1%	75.76	60.57
AR	0.2%	75.75	92.78
II	0.4%	75.92	96.49
0	1.0%	76.05	98.99
	3.0%	75.36	99.93
	0.0%	97.11	-
ß	0.1%	97.09	71.12
S	0.2%	97.19	89.59
5	0.4%	97.33	98.04
	1.0%	97.25	98.62
	3.0%	97.47	99.80
	0.0%	67.51	-
let	0.5%	67.38	0.17
gel	1.0%	67.13	87.74
Imag	2.0%	67.26	98.01
	3.0%	67.26	98.32

Evaluation of the prov	oosed back	door attac	k. We		Poisoning Rate	1%	2%		
benchmark using ResN	et18 trained	l on variou	s datasets	Random (1)		67.24/53.9	01 66.83/60.49		
and poisoning rates Ou	ir attack ma	intains CD	A while	Random (2)		67.23/56.8	8 66.80/66.11		
and poisoning rates. Our attack maintains CDA, while achieving high ASP even with small poisoning rates				<b>Bottom-</b> <i>k</i> (1)	<b>CDA(%)/ASR(%)</b>	67.03/22.5	66.80/55.96		
		in poisoinn	g raics.	<b>Bottom-</b> <i>k</i> (2)		67.04/0.3	1 67.02/92.81		
Method	<b>PSNR</b> ↑	<b>SSIM</b> ↑	<b>LPIPS</b>	<b>Top-</b> <i>k</i> (1)		67.13/87.7	4 67.26/98.01		
BadNets [16]	27.03	0.9921	0.0149	Effect of Differ	ent Frequency Selec	tion Schem	es: Results for		
Blend [5]	19.18	0.7291	0.2097	frequency filters generated using least sensitive, most sensitive and random frequencies. Choosing the top-k most sensitive frequencies provides the highest ASR among those options.					
SIG [1]	25.12	0.8988	0.0532						
Refool [35]	16.59	0.7701	0.2461						
SPM [30]	38.65	0.9665	0.0022						
Poison Ink [56]	41.62	0.9915	0.0020	100					
FTrojan [48]	44.87	0.9942	0.0005	100	$\overline{}$		Defending with		
FIBA [11]	18.05	0.8077	0.1113	S 80			JPEG Augmentatio		
Ours (ResNet18)	47.26	0.9998	0.0006	A/P			Training on JPEG		
Ours (ResNet34)	47.55	0.9998	0.0004	<b>G</b> 60 -			compressed images		
Ours (ResNet50)	46.90	0.9998	0.0009	ge	Base	line	maintains a high ASI		
Ours (DenseNet121)	47.21	0.9998	0.0001	et 40 -	Defe	nded ASR	10% in CDA The		
Ours (VGG19)	46.19	0.9998	0.0008		Unde	efended ASR	baseline denotes the		
Comparing Invisibility Metrics of Backdoor Attacks on ImageNet. Our attack achieves the best					60 40	20 0	CDA of the baseline model evaluated on compressed images.		
invisibility scores compared to other existing methods.				JPEG Compression Quality (%)					



King Abdullah University o Science and Technolog



### Results



Evaluation of defenses: Evaluation of various SOTA defenses against the proposed frequency-based attack shows the power of the proposed method in evading the defenses. (a) Grad-CAM shows high similarity in the attention regions for poisoned and nonpoisoned models; (b) **Pruning** the poisoned model maintains high ASR even after significant drop in CDA. (c) Neural **Cleanse** anomaly indices fall below the anomaly threshold (2.0).



Heatmaps of Various Frequency Selection: Fourier heatmaps of frequency based backdoor attacks with different frequency selection schemes: (a) Clean Model; (b) Random Frequency Selection; (c) Bottom-k Frequency Selection; (d) Proposed Top-k Frequency Selection