

---

# Track 1:

## Sparse Transfer Learning Accelerates and Enhances Certified Robustness: A Comprehensive Study

---

Anonymous Author(s)

Affiliation

Address

email

### Abstract

1 Certified robustness is a critical measure for assessing the reliability of machine  
2 learning systems. Traditionally, the computational burden associated with certifying  
3 the robustness of machine learning models has posed a substantial challenge,  
4 particularly with the continuous expansion of model sizes. In this paper, we  
5 introduce an innovative approach to expedite the verification process for  $L_2$ -norm  
6 certified robustness through sparse transfer learning. Our approach is both efficient  
7 and effective. It leverages verification results obtained from pre-training tasks and  
8 applies sparse updates to these results. To enhance performance, we incorporate  
9 dynamic sparse mask selection and introduce a novel stability-based regularizer  
10 called DiffStab. Empirical results demonstrate that our method accelerates the  
11 verification process for downstream tasks by as much as **70-80%**, with only slight  
12 reductions in certified accuracy compared to dense parameter updates. We further  
13 validate that this performance improvement is even more pronounced in the few-  
14 shot transfer learning scenario.

### 15 1 Introduction

16 In recent years, ensuring the certified robustness of machine learning systems has emerged as  
17 a paramount research challenge. The primary objective is to guarantee consistent and resilient  
18 output predictions, impervious to perturbations spanning a defined range in any direction. Diverse  
19 verification techniques have been devised to quantify the certified robustness of neural networks.  
20 When confronting inputs perturbed within some  $L_{inf}$ -norm bound, the prevailing verification methods  
21 center around the branch-and-bound (BaB) technique [20, 16, 18]. In cases involving  $L_2$ -norm  
22 perturbations, randomized-smoothing approaches reign supreme [4, 10].

23 However, it is a widely recognized challenge that commonly used certified verification methods,  
24 such as the BaB methods [20, 18] and randomized smoothing [4] grapple with the inherent issue of  
25 computationally expensive verification for each sample. In fact, the computational cost of verification  
26 often surpasses that of inference for the same sample by *several orders of magnitude*. For instance,  
27 the BaB method exhibits exponential complexity, while randomized smoothing typically demands  
28 the sampling of approximately 1,000 noisy inputs for each individual verification. This predicament  
29 of resource-intensive verification is further exacerbated by the exponential growth in the sizes of  
30 state-of-the-art (SOTA) models across various benchmarks.

31 In this paper, we concentrate on developing novel streamlined techniques designed to expedite  
32 the verification processes based on randomized smoothing for  $L_2$ -norm certified robustness. Our  
33 approach begins by identifying ways to **efficiently reuse** the verification results from pre-training

34 tasks to downstream tasks, and our innovation here is to introduce the tool of **sparse transfer learning**  
35 to update only a select subset of network parameters during the transfer. We then implement our  
36 novel *differential sparse verification* techniques to accelerate the verification process by leveraging  
37 specific patterns of sparsity. This is chiefly accomplished by hastening the forward propagation of  
38 noisy samples from the Monte-Carlo sampling of randomized smoothing-based verification, using  
39 (structured) sparse update vector multiplication. We further introduce two techniques to augment  
40 the certified robustness for sparse transfer learning, namely dynamic sparse mask selection and a  
41 novel stability-based regularizer. They result in significant enhancements in both the speed of the  
42 verification process and the robustness of the certified outcomes, when compared to the conventional  
43 approach of direct training and verification on downstream tasks.

44 Specifically, our contributions are outlined as follows:

- 45 • We for this first time investigate the use of sparse transfer learning to expedite the certified  
46 verification process, capitalizing on reusing the verification results from the upstream task  
47 and executing sparse weight updates. Specifically, we employ sparse transfer learning  
48 with three distinct sparsity patterns, thereby facilitating efficient transfer and accelerating  
49 the downstream verification process. This is achieved by propagating the intermediate  
50 verification results using the sparse convolutional operator.
- 51 • Recognizing that sparse transfer learning may affect the certified robustness of transferred  
52 models, we further propose to boost this process using dynamic mask selection and a novel  
53 stability-based regularizer. These measures significantly narrow the performance gap with  
54 the upper bound achieved by dense parameter updates.
- 55 • We empirically discover that our approach can hasten the verification process on downstream  
56 tasks by up to 70-80%, with only slight reductions in certified accuracy compared to dense  
57 parameter update. Furthermore, we find that the advantages of sparse transfer learning and  
58 acceleration can be further amplified in the context of few-shot transfer learning.

## 59 **2 Related Work**

### 60 **2.1 $L_2$ -norm Certified Robustness**

61 Research into  $L_2$ -norm certified robustness aims to ensure stable machine learning system outputs  
62 when input perturbations lie within an  $L_2$ -norm ball. Cohen et al. [4] pioneered a verification method  
63 for  $L_2$ -norm certified robustness using randomized smoothing based on Monte-Carlo sampling [4].  
64 Kumar and Goldstein [10] introduced a variant that facilitates  $L_2$ -robust training and verification for  
65 tasks with structured outputs, such as semantic segmentation [10]. An alternative approach verifies  
66  $L_2$ -norm robustness by leveraging the diffusion model to denoise the input before making predictions  
67 with benign-trained models [1, 19]. However, this diffusion model-based method introduces an  
68 added denoising step during inference. This paper centers on accelerating the prevalent randomized  
69 smoothing techniques for  $L_2$ -norm certified robustness.

### 70 **2.2 Transfer learning**

71 Transfer learning facilitates knowledge transfer from a source to a target domain, particularly when  
72 data in the target domain is scarce. Many approaches adopt a pretrain-and-finetune framework,  
73 varying primarily in their pre-training objectives. This includes contrastive pre-training [3, 11],  
74 pretext tasks [7], and autoencoding [5]. Recently, Guo et al. [9] introduced DiffPruning, a parameter-  
75 efficient method that updates only a sparse subset of model parameters for each downstream task. In  
76 this paper, we harness DiffPruning for both sparse transfer learning and sparse differential verification.

### 77 **2.3 Robustness Transferring**

78 Recent studies have highlighted the robustness of neural networks pre-trained on large-scale datasets.  
79 Such networks tend to possess robust feature extractors that can be transferred to downstream  
80 tasks [4, 15, 10, 12]. Salman et al. [14] discovered that adversarially trained networks can enhance  
81 accuracy in these downstream tasks. Furthermore, Vaishnavi et al. [17] introduced a method using  
82 knowledge transfer to expedite the training for certified robustness. Nevertheless, to our knowledge,

83 our work is pioneering in its approach to accelerate the certified robustness verification process via  
84 sparse transfer learning.

### 85 3 Methodology

86 We introduce the preliminaries of this work, including sparsity patterns and dynamic mask selection  
87 with RigL in appendix A. Certified verification techniques, such as randomized smoothing, grapple  
88 with substantial computational overhead—often eclipsing the inference time for the same sample. For  
89 the  $L_2$ -norm, this cost intensifies with randomized smoothing methods [4, 15, 10], which use Monte-  
90 Carlo sampling to certify each original sample by sampling multiple noisy inputs. In this section, we  
91 explore how sparse transfer learning can bolster the efficiency of the randomized smoothing-based  
92 verification and enhance the certified robustness for downstream tasks. We also discuss dynamic  
93 mask selection and our novel **stability regularizer**, both tailored to amplify certified robustness. The  
94 architecture of our framework is illustrated in fig. 1.

#### 95 3.1 Sparse Transfer Learning for Certified Robustness

96 Sparse transfer learning involves updating a selected subset of parameters in pre-trained models  
97 during transfer. We pinpoint two benefits of employing sparse transfer learning for  $L_2$ -norm certified  
98 robustness: ❶ Transfer learning typically yields superior certified robustness in downstream tasks  
99 than training exclusively on those tasks. This is attributed to the foundational robustness instilled  
100 during the pre-training phase [15, 12]. ❷ Sparse transfer learning facilitates the acceleration of  
101 the randomized smoothing verification process across various sparsity patterns. We expedite the  
102 randomized smoothing certification by leveraging efficient Monte-Carlo sampling inference, informed  
103 by sparse update vectors derived from sparse transfer learning.

104 We further explore the potential of sparse transfer learning to expedite the verification process by  
105 examining various sparsity patterns. In our approach, we integrate DiffPruning, as presented in Guo  
106 et al. [9], with different types of sparsity: unstructured, structured (channel-wise), and group-wise,  
107 as detailed in Sec. A.1. This amalgamation allows for sparse transfer learning. Furthermore, we  
108 discuss the method of capitalizing on the certified verification outcomes from pre-training tasks. By  
109 employing the sparse update masks corresponding to the various sparsity patterns, we aim to speed  
110 up the verification procedure for the transferred tasks. It’s crucial to mention that in order to benefit  
111 from the verification results of the pre-training tasks, consistency in input between pre-training and  
112 downstream tasks is imperative

113 In methods rooted in randomized smoothing, the predominant computational demand during the  
114 verification process stems from Monte-Carlo sampling. For each sample undergoing verification, it  
115 is commonplace for these techniques to draw upon 1,000 noisy inputs, forecast outcomes, and then  
116 gauge the verification conclusion from these forecasts. This underscores that bolstering the speed of  
117 forward propagation for each prediction can, in turn, hasten the overarching verification procedure. In  
118 this study, we venture to enhance the pace of the forward propagation. We achieve this by integrating  
119 the differential outcome of forward propagation, brought about by sparse update vectors, with the  
120 dense verification results inherited from pre-training tasks, as depicted in Fig. 1.

121 Subsequently, we materialize this acceleration across diverse sparsity patterns:

122 • **Unstructured sparsity** In CNNs, the convolutional operation is the primary source of com-  
123 putational complexity during inference. To achieve acceleration, our focus is on optimizing  
124 this convolutional operation. It’s well-understood that convolutional operations can be  
125 translated into matrix multiplications. Therefore, the differential forward propagation can  
126 be conducted using a matrix multiplication between the dense input matrix and the sparse  
127 parameter update matrix for each convolutional layer. This process is further expedited  
128 using a sparse **coordinate list** operator tailored for matrix computations.

129 In the forward propagation, for each layer, we combine this differential output with the results  
130 previously obtained from the pre-training task. This integrated result is then propagated  
131 to the subsequent layer. By iteratively integrating the outcome of this sparse forward  
132 propagation for each convolutional layer, we can efficiently compute the output of the final  
133 convolutional layer.

- 134 • **Structured sparsity** The acceleration process for structured sparsity is notably straight-  
 135 forward. The sparse update vector in this context is channel-wise, functioning as a binary  
 136 indicator for every layer. When this indicator has a value of 1, it signifies that the parameters  
 137 of the related channel have undergone updates. We compute the differential output exclu-  
 138 sively for the updated channels. Then, for each layer, we combine this differential output  
 139 with the results derived from the pre-training task. The aggregated result is subsequently  
 140 propagated to the following layer.
- 141 • **Group-wise sparsity** Group-wise sparsity can be conceptualized as an amalgamation of  
 142 both unstructured and structured sparsity. Given that certain channels without selected dense  
 143 blocks are omitted, we ultimately achieve a sparsity mask with an unstructured configuration.  
 144 Consequently, we can employ a combined acceleration strategy, drawing from the methods  
 145 above used for both structured and unstructured sparsity.

146 We empirically find that, for a given sparse ratio  $k$ , structured sparsity typically yields a more pro-  
 147 nounced acceleration compared to unstructured sparsity, with group-wise sparsity falling somewhere  
 148 in the middle. In contrast, when evaluating certifiable robustness (specifically, verified accuracy),  
 149 the performance trend for each sparsity pattern is opposite to their respective acceleration effects.  
 150 Notably, group-wise sparsity strikes a more balanced compromise between acceleration and certifiable  
 151 robustness in comparison to the other two sparsity paradigms.

### 152 3.2 Regularizing Sparse Transfer Learning

153 As mentioned above, the proposed method with sparse transfer learning can help the model achieve  
 154 better-verified accuracy than training directly on downstream tasks, but not as good as dense transfer  
 155 learning, where we update all the parameters while transferring. We identify 2 reasons for this  
 156 phenomenon: firstly, we originally expected that sparse transfer learning is possible to achieve better  
 157 certified robustness than dense transfer learning since the network is already pre-trained to be robust  
 158 and has stable intermediate outputs for its layers given the same input. However, we failed to observe  
 159 this phenomenon and conclude that unconstrained sparse transfer learning is unable to preserve the  
 160 robustness obtained from the pre-training task. Secondly, we believe that the domain gap between  
 161 the pre-training task and the downstream task prevents sparse transfer learning to achieve better  
 162 robustness than dense transfer learning.

163 To tackle these two challenges, we introduce a dual-method approach: **Firstly**, We advocate for  
 164 a regularizer based on stability, which specifically targets the  $L_2$  distance of the lower and upper  
 165 bounds for each neuron. By ensuring these bounds remain consistent between the pre-training  
 166 and downstream tasks—given identical input and perturbation ranges—we aim to maintain the  
 167 inherent stability and robustness from the pre-training phase. To this end, we employ Interval Bound  
 168 Propagation (IBP) [8]. The  $L_{inf}$ -norm bounds provided by IBP are not only efficient but also align  
 169 with the computational complexity of a network’s forward pass. It’s worth highlighting that even  
 170 though the  $L_{inf}$ -norm bound is distinct from the  $L_2$ -norm bound, our empirical findings suggest  
 171 that the former is effective in regularizing  $L_2$ -norm robustness. Formally defined, if the lower and  
 172 upper bounds of neuron  $i$  for the pre-training task are denoted by  $lb_i$  and  $ub_i$  respectively, and for the  
 173 downstream task they are  $lb'_i$  and  $ub'_i$ , the regularization loss is computed as follows:

$$loss_{stab} = \frac{1}{N} \sum_i^N (lb'_i - lb_i)^2 + (ub'_i - ub_i)^2 \quad (1)$$

174 Where  $N$  is the number of neurons across the network. We call this regularizer as *DiffStab*. And the  
 175 overall loss for transfer learning is:

$$loss = loss_{orig} + loss_{stab} \quad (2)$$

176 Where  $loss_{orig}$  is the original loss of transfer learning. **Secondly**, to mitigate the domain gap  
 177 challenge, we advocate for dynamic mask selection. Specifically, we implement the RigL approach  
 178 as outlined in [6]. This method ensures enhanced mask flexibility during the transfer phase. Our  
 179 empirical analysis confirms that dynamic mask selection markedly boosts certified robustness in  
 180 sparse transfer learning, especially when confronted with a substantial domain gap.

Table 1: The comparison of verified accuracy and verification time of different transfer setting with different sparsity patterns.

Sparsity Pattern	Updated Params(%)	Direct Train	Dense Transfer	Sparse Transfer						
		100	100	1	2	4	8	16	32	64
Unstruct	Ver Acc(%)	60.4	61.2	55.1	56.2	57.9	59.1	60.3	60.6	61.1
	Time Saved(%)	0	0	77.6	63.5	46.6	29.1	10.2	3.2	1.2
Struct	Ver Acc(%)	60.4	61.2	53.4	54.7	56.2	57.6	58.8	59.7	60.3
	Time Saved(%)	0	0	92.1	88.5	82.2	72.8	55.4	33.1	15.8
Group-wise	Ver Acc(%)	60.4	61.2	54.0	55.3	56.8	57.8	59.1	59.8	60.5
	Time Saved(%)	0	0	87.2	81.8	75.2	61.9	49.3	28.7	12.5

## 181 4 Experiments

182 In this section, our objective is to address two primary inquiries via comprehensive experiments: (1)  
 183 How effectively does the proposed method hasten the certified verification process and amplify the  
 184 certified robustness for a downstream task under  $L_2$ -norm input perturbations? (2) How do DiffStab  
 185 and RigL contribute to enhancing the certified robustness performance in the context of our proposed  
 186 sparse transfer learning and verification methodology?

187 To address the posed questions, we carry out experiments in two distinct settings across two datasets,  
 188 including CIFAR10 and CelebV-HQ. The details about experiment settings are highlighted in ap-  
 189 pendix B.

### 190 4.1 Sparse Transfer Learning Accelerates and Enhances Certified Robustness

Table 2: The comparison of verified accuracies before and after adding DiffStab regularizer and RigL(dynamic mask selection) of different sparsity patterns. The relative improvements in the brackets are obtained by comparing them with the baselines of different sparsities.

Sparsity Pattern	UpdateParams	Direct Train	Dense Transfer	Sparse Transfer						
		100	100	1	2	4	8	16	32	64
<b>Unstruct</b>	Ver Acc(%)	60.4	61.2	55.1	56.2	57.9	59.1	60.3	60.6	61.1
	+DiffStab+RigL Ver Acc(%)	60.4	62.2 (+1.0)	58.1 (+3.0)	59.0 (+2.8)	60.6 (+2.7)	61.2 (+2.1)	61.5 (+1.2)	62.2 (+1.6)	62.2 (+1.1)
<b>Struct</b>	Ver Acc(%)	60.4	61.2	53.4	54.7	56.2	57.6	58.8	59.7	60.3
	+DiffStab+RigL Ver Acc(%)	60.4	62.2 (+1.0)	57.0 (+3.6)	58.7 (+4.0)	59.4 (+3.2)	60.6 (+3.0)	60.9 (+2.1)	61.4 (+1.7)	61.6 (+1.3)
<b>Group-wise</b>	Ver Acc(%)	60.4	61.2	54.0	55.3	56.8	57.8	59.1	59.8	60.5
	+DiffStab+RigL Ver Acc(%)	60.4	62.2 (+1.0)	57.0 (+3.0)	58.6 (+3.3)	59.3 (+2.5)	60.3 (+2.5)	60.7 (+1.6)	61.2 (+1.4)	61.7 (+1.2)

#### 191 4.1.1 CIFAR10 Results

192 In this subsection, we assess the effectiveness of combining sparse transfer learning with sparse  
 193 differential verification to expedite the randomized smoothing-based verification process. Notably,  
 194 when given an ample amount of pre-training data, sparse transfer learning not only facilitates faster  
 195 performance but also achieves superior results.

196 To corroborate the acceleration effect, we implemented sparse transfer learning on the CIFAR10  
 197 dataset at predetermined sparsity ratios. We juxtaposed the outcomes from sparse differential  
 198 verification with those obtained using the standard randomized smoothing. The results of this  
 199 comparison are delineated in table 1. Although similar acceleration findings were empirically  
 200 noted on the CelebV-HQ dataset (owing to the consistent network architecture and a dominant  
 201 influence of sparsity ratio over input or network configuration), for the sake of brevity, we’ve confined  
 202 our exposition to the CIFAR10 dataset. As evident from table 1, as the sparsity ratio increases,  
 203 sparse differential verification can hasten the verification process by a staggering 77.6% to 92.1%.  
 204 However, a trade-off is observed in the form of a reduced verified accuracy. While we employed  
 205 contrastive learning for pre-training, aiming to harness robust self-supervision signals, the scale of the  
 206 dataset remains a constraint, limiting the significant benefits of pre-training for subsequent tasks. In  
 207 subsequent sections, we will discuss how augmenting the pre-training dataset size can alleviate this  
 208 challenge. Additionally, by incorporating our novel DiffStab regularizer and dynamic mask selection,  
 209 we demonstrate that performance can be further enhanced.

210 When we examine various sparsity patterns presented in table 1, it’s evident that the acceleration ef-  
211 fects increase in the order of unstructured, group-wise, and structured sparse differential verifications.  
212 This progression aligns with our expectations. Structured sparsity directly omits entire channels from  
213 the verification process, while group-wise sparsity can be perceived as an amalgamation of both  
214 unstructured and structured sparsity, as outlined in our methodology.

215 However, when looking at verified accuracies, they tend to decrease in the order from unstructured  
216 to structured sparsity. This outcome is plausible since unstructured sparsity employs the most  
217 adaptive sparsity masks. This observation parallels findings in the model compression domain where  
218 unstructured pruning often surpasses structured pruning in terms of subnetwork performance.

219 Upon Comparing group-wise sparsity with the other two types, it becomes clear that group-wise  
220 sparsity aligns more closely with unstructured sparsity in terms of verified accuracy, while resembling  
221 structured sparsity in acceleration outcomes. Therefore, we can infer that group-wise sparsity strikes  
222 an optimal balance, presenting a commendable trade-off between performance and acceleration,  
223 particularly when certifying robustness.

224 To demonstrate the broad applicability of our method across various network architectures, we further  
225 evaluated its acceleration performance on both ResNet-18 and VGG-16. The results are presented in  
226 table 3 in the Appendix.

#### 227 4.1.2 CelebV-HQ Results

228 For the CelebV-HQ dataset, we commenced with analogous experiments involving unstructured  
229 sparsity, both under a standard transfer setting utilizing 100% of the downstream data and a few-shot  
230 transfer setting with just 1% of downstream data. For detailed results, see the 1st and 5th rows in  
231 table 4.

232 By comparing these outcomes with those in table 1, it becomes evident that the expansive scale of the  
233 pre-training dataset in CelebV-HQ markedly bolsters the certified robustness achieved through sparse  
234 transfer learning. Let’s remember that for our pre-training, we utilized 40 attributes, while only 1  
235 attribute was used for each downstream task. Notably, even when a mere 8% of network parameters  
236 are updated during sparse transfer learning, the enhanced network showcases a performance that’s on  
237 par with direct training that involves dense parameter updates. This comes with the added advantage  
238 of a 29.1% acceleration for unstructured sparsity.

239 The advantages of sparse transfer learning become even more pronounced in a few-shot transfer  
240 learning environment. Here, sparse transfer learning significantly outperforms direct training. This  
241 can be attributed to the fact that the extensive multi-attribute classification pre-training infuses the  
242 network with substantial robustness. In contrast, direct training is limited by its access to a smaller  
243 dataset, curtailing its robust training capabilities. Interestingly, the performance disparity between  
244 dense transfer learning and sparse transfer learning narrows in the few-shot setting. This can be  
245 explained by the limited data available for finetuning in the few-shot scenario. Consequently, the  
246 performance is less adversely impacted by the ‘lazy’ update strategy, that is, the sparse parameter  
247 update. More study is provided in appendix C and appendix D.

## 248 5 Conclusion

249 In this paper, we introduce sparse differential verification to accelerate the  $L_2$ -norm robustness  
250 verification process based on randomized smoothing. Building on sparse differential forward prop-  
251 agation, our approach hastens the Monte-Carlo Sampling inherent to randomized smoothing. We  
252 explore three sparsity patterns for transfer learning, discussing their pros and cons. To bridge the gap  
253 between dense and sparse transferring, we employ dynamic mask selection and our new DiffStab  
254 regularizer. Empirically, our method achieves up to 80% acceleration while maintaining verified  
255 accuracies comparable to dense transfer methods. One constraint is the need for consistent input  
256 between pre-training and downstream tasks, limiting our model’s breadth. Still, our work offers a  
257 promising step towards leveraging transfer learning for faster, reliable machine learning verification.

## 258 References

- 259 [1] Nicholas Carlini, Florian Tramer, J Zico Kolter, et al. (certified!!) adversarial robustness for  
260 free! *arXiv preprint arXiv:2206.10550*, 2022.
- 261 [2] Tianlong Chen, Xuxi Chen, Xiaolong Ma, Yanzhi Wang, and Zhangyang Wang. Coarsening the  
262 granularity: Towards structurally sparse lottery tickets. *arXiv preprint arXiv:2202.04736*, 2022.
- 263 [3] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum  
264 contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- 265 [4] Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized  
266 smoothing. In *International Conference on Machine Learning*, pages 1310–1320. PMLR, 2019.
- 267 [5] Andrew M Dai and Quoc V Le. Semi-supervised sequence learning. *Advances in neural  
268 information processing systems*, 28, 2015.
- 269 [6] Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the  
270 lottery: Making all tickets winners. In *International Conference on Machine Learning*, pages  
271 2943–2952. PMLR, 2020.
- 272 [7] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by  
273 predicting image rotations. *arXiv preprint arXiv:1803.07728*, 2018.
- 274 [8] Sven Gowal, Krishnamurthy Dvijotham, Robert Stanforth, Rudy Bunel, Chongli Qin, Jonathan  
275 Uesato, Relja Arandjelovic, Timothy Mann, and Pushmeet Kohli. On the effectiveness of interval  
276 bound propagation for training verifiably robust models. *arXiv preprint arXiv:1810.12715*,  
277 2018.
- 278 [9] Demi Guo, Alexander M Rush, and Yoon Kim. Parameter-efficient transfer learning with diff  
279 pruning. *arXiv preprint arXiv:2012.07463*, 2020.
- 280 [10] Aounon Kumar and Tom Goldstein. Center smoothing: Certified robustness for networks with  
281 structured outputs. *Advances in Neural Information Processing Systems*, 34:5560–5575, 2021.
- 282 [11] Fuli Luo, Pengcheng Yang, Shicheng Li, Xuancheng Ren, and Xu Sun. Capt: contrastive  
283 pre-training for learning denoised sequence representations. *arXiv preprint arXiv:2010.06351*,  
284 2020.
- 285 [12] Laura Fee Nern and Yash Sharma. How adversarial robustness transfers from pre-training to  
286 downstream tasks. *arXiv preprint arXiv:2208.03835*, 2022.
- 287 [13] Masuma Akter Rumi, Xiaolong Ma, Yanzhi Wang, and Peng Jiang. Accelerating sparse  
288 cnn inference on gpus with performance-aware weight pruning. In *Proceedings of the ACM  
289 International Conference on Parallel Architectures and Compilation Techniques*, pages 267–278,  
290 2020.
- 291 [14] Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. Do ad-  
292 versarially robust imagenet models transfer better? *Advances in Neural Information Processing  
293 Systems*, 33:3533–3545, 2020.
- 294 [15] Ali Shafahi, Parsa Saadatpanah, Chen Zhu, Amin Ghiasi, Christoph Studer, David Jacobs, and  
295 Tom Goldstein. Adversarially robust transfer learning. *arXiv preprint arXiv:1905.08232*, 2019.
- 296 [16] Zhouxing Shi, Yihan Wang, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Fast certified robust  
297 training with short warmup. *Advances in Neural Information Processing Systems*, 34, 2021.
- 298 [17] Pratik Vaishnavi, Kevin Eykholt, and Amir Rahmati. Accelerating certified robustness training  
299 via knowledge transfer. *Advances in Neural Information Processing Systems*, 35:5269–5281,  
300 2022.
- 301 [18] Shiqi Wang, Huan Zhang, Kaidi Xu, Xue Lin, Suman Jana, Cho-Jui Hsieh, and J Zico Kolter.  
302 Beta-crown: Efficient bound propagation with per-neuron split constraints for complete and  
303 incomplete neural network verification. *arXiv preprint arXiv:2103.06624*, 2021.

- 304 [19] Quanlin Wu, Hang Ye, Yuntian Gu, Huishuai Zhang, Liwei Wang, and Di He. Denoising  
 305 masked autoencoders are certifiable robust vision learners. *arXiv preprint arXiv:2210.06983*,  
 306 2022.
- 307 [20] Kaidi Xu, Huan Zhang, Shiqi Wang, Yihan Wang, Suman Jana, Xue Lin, and Cho-Jui Hsieh.  
 308 Fast and complete: Enabling complete neural network verification with rapid and massively  
 309 parallel incomplete verifiers. *arXiv preprint arXiv:2011.13824*, 2020.
- 310 [21] Hao Zhu, Wayne Wu, Wentao Zhu, Liming Jiang, Siwei Tang, Li Zhang, Ziwei Liu, and  
 311 Chen Change Loy. Celebv-hq: A large-scale video facial attributes dataset. *arXiv preprint*  
 312 *arXiv:2207.12393*, 2022.

## 313 A Preliminary

### 314 A.1 Sparsity Patterns

315 In this paper, we explore three distinct types of sparsity: unstructured sparsity, structured sparsity,  
 316 and group-wise sparsity[2]. Given a sparse ratio  $k$  ( $0 < k < 1$ ), the network’s sparsity can be  
 317 depicted by a binary mask  $M$ , where each element corresponds to a single parameter of the model  
 318 and the ratio of non-zero elements equals  $k$ . In a conventional neural network with convolutional  
 319 layers, unstructured sparsity implies no constraints on the mask  $M$  other than its sparse ratio being  $k$ ,  
 320 whereas structured sparsity ensures channel-wise uniformity in the mask, i.e., all parameters in the  
 321 same channel or kernel must share the same mask value. Lastly, group-wise sparsity [2] combines  
 322 both unstructured and structured sparsity. Here, the mask  $M$  is first generated in an unstructured  
 323 manner, and a hypergraph partitioning algorithm [13] identifies dense blocks of activated parameters,  
 324 reactivating any deactivated parameters within these blocks. The remaining activated parameters not  
 325 in these dense blocks are deactivated, with the ratio of chosen blocks controlled so that the sparse  
 326 ratio still equals  $k$  after determining the group-wise mask.

### 327 A.2 Dynamic Mask Selection with RigL

328 We adapt the RigL method [6] for dynamic mask selection across three sparsity patterns. The RigL  
 329 method dynamically activates and deactivates network parameters based on gradient magnitudes and  
 330 magnitudes of parameter values, respectively, during training. Originally designed for unstructured  
 331 sparsity, we modify RigL for structured sparsity by shifting from parameter-level to channel-level  
 332 activation and deactivation, guided by the magnitudes and gradient magnitudes of the channel weight  
 333  $\gamma$  in BatchNorm layers. We control the overall sparsity using a strategy similar to [6]. For  
 334 group-wise sparsity, we simply follow the unstructured version of RigL, allowing the group-wise  
 335 sparsity to determine the structured sparsity.

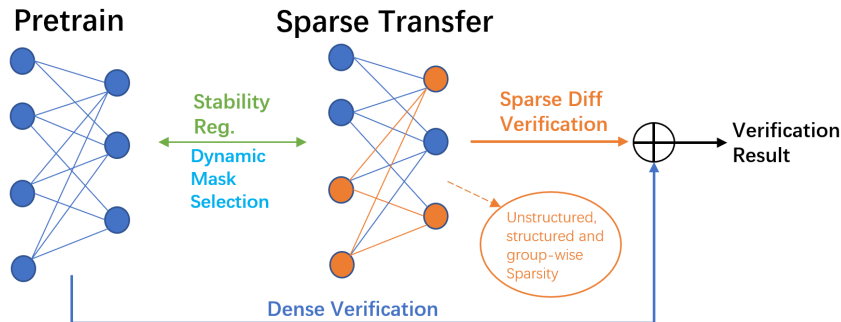


Figure 1: Our novel framework integrating differential verification with sparse transfer learning. For downstream tasks, the pre-trained network is refined using DiffPruning[9] coupled with dynamic mask selection. To ensure network robustness during sparse parameter updates, we introduce a neuron stability-based regularizer. For verification, we synergize sparse differential verification techniques with reusable dense verification results to yield a conclusive verification outcome.



Table 3: The acceleration results for a certified verification of different architectures with sparse transfer learning under different sparsity patterns.

Architecture	Sparsity Pattern	Updated Params(%)	Sparse Transfer						
			1	2	4	8	16	32	64
ResNet-18	Unstruct	Time Saved(%)	92.5	88.4	82.4	72.1	55.8	33.2	16.2
	Struct	Time Saved(%)	77.8	63.8	47.2	29.0	11.0	3.1	1.4
	Group-wise	Time Saved(%)	87.2	81.8	75.2	61.9	49.3	28.7	12.5
ResNet-50	Unstruct	Time Saved(%)	92.1	88.5	82.2	72.8	55.4	33.1	15.8
	Struct	Time Saved(%)	77.6	63.5	46.6	29.1	10.2	3.2	1.2
	Group-wise	Time Saved(%)	87.2	81.8	75.2	61.9	49.3	28.7	12.5
VGG-16	Unstruct	Time Saved(%)	92.5	89.1	82.5	73.1	55.4	33.2	15.6
	Struct	Time Saved(%)	77.8	63.6	46.8	29.3	10.3	3.5	1.4
	Group-wise	Time Saved(%)	89.4	84.8	77.5	63.5	50.3	30.1	13.7

## 336 B Datasets

337 Our experiment setting is elaborated as below:

338 **CIFAR10** Initially, we deploy CIFAR10 to gauge the efficacy of our methodologies. Given the  
 339 comparatively modest scale of CIFAR10, we employ contrastive learning during the pre-training  
 340 phase to extract a rich self-supervision signal and subsequently use image classification for the  
 341 downstream task. It’s noteworthy that contrastive learning predicts based on a dense feature map  
 342 rather than a singular scalar probability. Consequently, randomized smoothing is unsuitable for pre-  
 343 training this model. As an alternative, we utilize Center Smoothing [10]—a variant of randomized  
 344 smoothing designed to secure  $L_2$ -norm robustness for dense outputs—in tandem with contrastive  
 345 learning to pre-train our network. Following this, randomized smoothing is incorporated for transfer  
 346 learning within the image classification task.

347 **CelebV-HQ** Introduced in [21], CelebV-HQ is a contemporary benchmark tailored for multi-  
 348 attribute classification tasks. It offers classifications for 83 facial attributes bifurcated into two  
 349 categories: appearance and action attributes. Given that CelebV-HQ is rooted in video classification,  
 350 we extract five disparate frames at random from each video, resize them to 64x64 dimensions,  
 351 and utilize them as the input for every sample. This approach morphs the multi-attribute video  
 352 classification task into a multi-attribute image classification challenge. Our strategy then encompasses  
 353 random sampling of 40 attributes for pre-training, with the remaining attributes earmarked for  
 354 downstream transfer. It’s pivotal to understand that, in this dataset, the pre-training endeavor involves  
 355 multi-attribute classification using the 40 selectively sampled attributes. Each subsequent downstream  
 356 task revolves around binary classification, leveraging each of the residual attributes. To ascertain  
 357 comprehensive results, evaluations across all downstream tasks are averaged. In the context of this  
 358 dataset, we contemplate three distinct transfer settings: standard transfer, and a "few-shot" transfer,  
 359 wherein the downstream tasks have access to merely 1% of randomly sampled data for their training.

360 We employ DiffPruning, as previously discussed, for our sparse transfer learning approach. Our  
 361 evaluation criteria are bifurcated: first, we gauge the time saved in verification through our proposed  
 362 methodologies in contrast to direct verification of samples in downstream tasks. Second, we assess  
 363 certified robustness, which equates to the verified accuracies, as confirmed by randomized smoothing  
 364 in the subsequent tasks. Pertaining to the model architecture, unless stated otherwise, we consistently  
 365 utilize ResNet-50 as the foundational network for our experiments. Only the fully connected layers  
 366 of the network undergo reinitialization, with sparse transfer learning executed on the convolutional  
 367 layers. We’ve earmarked the perturbation radius of the input  $L_2$ -norm ball at 0.25, considering a  
 368 normalized image input.

## 369 C Acceleration Results on More Architectures

370 In order to validate the consistent acceleration performance of our proposed techniques, we substituted  
 371 ResNet-50 with both ResNet-18 and VGG-16 in our CIFAR10 experiments. Our objective was  
 372 to compare the acceleration outcomes of these three architectures when subject to randomized  
 373 smoothing-based verification. These findings are documented in table 3.

Table 4: The comparison of under normal/few-shot transfer setting. **+Reg** means applying DiffStab and RigL. The relative improvements in the brackets are obtained by comparing them with the unstructured baseline.

		Updated Params(%)	Direct Train	Dense Transfer	Sparse Transfer						
			100	100	1	2	4	8	16	32	64
Normal transfer	Unstruct (Baseline)	Ver Acc(%)	55.8	59.1	52.8	53.8	54.2	55.3	56.6	57.2	57.9
	Unstruct +Reg	Ver Acc(%)	55.8	59.1	56.2 (+3.4)	56.9 (+3.1)	57.5 (+3.2)	57.9 (+2.6)	58.6 (+2.0)	58.9 (+1.7)	59.0 (+1.1)
	Struct +Reg	Ver Acc(%)	55.8	59.1	54.7 (+1.9)	56.3 (+2.5)	56.6 (+2.4)	57.1 (+1.8)	57.8 (+1.2)	58.4 (+1.2)	58.8 (+0.9)
	Group-wise +Reg	Ver Acc(%)	55.8	59.1	55.6 (+2.8)	57.0 (+3.2)	57.2 (+3.0)	57.7 (+2.4)	58.3 (+1.7)	58.7 (+1.5)	58.8 (+0.9)
Few-shot transfer	Unstruct (Baseline)	Ver Acc(%)	39.2	54.2	48.6	50.2	51.3	52.4	53.2	53.6	53.8
	Unstruct +Reg	Ver Acc(%)	39.2	54.2	52.5 (+3.9)	53.1 (+2.9)	53.5 (+2.4)	53.8 (+1.4)	54.0 (+0.8)	54.1 (+0.5)	54.1 (+0.3)
	Struct +Reg	Ver Acc(%)	39.2	54.2	49.3 (+0.7)	51.9 (+1.7)	52.6 (+1.3)	53.1 (+0.8)	53.6 (+0.4)	53.8 (+0.2)	53.9 (+0.1)
	Group-wise +Reg	Ver Acc(%)	39.2	54.2	51.6 (+3.0)	52.7 (+2.5)	53.2 (+0.9)	53.5 (+1.1)	53.7 (+0.5)	53.9 (+0.3)	53.9 (+0.1)

Remarkably, the proportion of verification time saved remains relatively stable across the diverse architectures. For instance, the discrepancy between ResNet-18 and ResNet-50 is negligible, remaining within a 1% margin in most scenarios. VGG-16 presents marginally superior acceleration outcomes. This can potentially be attributed to VGG-16’s relatively straightforward architecture when contrasted with residual networks. Consequently, it encompasses fewer operations, which wouldn’t benefit from acceleration during forward propagation.

## D DiffStab and Dynamic Mask Selection Boost Certified Robustness of Sparse Transfer Learning

From the results presented in the preceding subsection, a discernible performance discrepancy between dense transfer learning and sparse transfer learning remains evident. We delved into potential reasons for this in Sec. 3.2, subsequently proposing two remedies: the DiffStab regularizer and dynamic mask selection using RigL. Upon applying these methodologies to three distinct sparsity patterns on both the CIFAR10 and CelebV-HQ datasets, the outcomes, as depicted in table 2 and table 4 respectively, show a marked enhancement in verified accuracy for sparse transfer learning. The effectiveness of these strategies is directly proportional to the degree of sparsity, with higher sparsity ratios benefiting more significantly.

Interestingly, while these techniques are tailored for sparse transfer learning, they appear to have no discernible impact on dense transfer learning. This aligns with our expectations, given that there’s inherently no room for dynamic mask selection in dense parameter updates. Moreover, it can be inferred that the DiffStab regularizer truly shines in environments where updated parameters are sufficiently sparse. This enables the regularizer to more effectively modulate network stability and robustness, without inadvertently hindering model training.

With the implementation of the two techniques, we are now equipped to pinpoint hyperparameter configurations that strike a balance between impressive verified accuracies and commendable acceleration outcomes for sparse transfer learning. **1** Taking the CIFAR10 dataset as an example: Among configurations that surpass the verified accuracy of direct training, structured sparsity with an 8% sparse ratio stands out, yielding a remarkable 72.8% reduction in verification time when juxtaposed with traditional verification. When filtering for configurations that achieve over 80% verification acceleration, the same structured sparsity setting with an 8% sparse ratio boasts the pinnacle of verified accuracy, only trailing direct training by a slim 1% in accuracy. **2** Turning our attention to the CelebV-HQ dataset under the standard transfer setting: We discern that nearly all sparse transfer configurations armed with regularizers outperform direct training. Notably, group-wise sparsity at a 16% sparse ratio demonstrates a negligible performance dip, less than 1% compared to dense transfer, while simultaneously realizing a 49.3% acceleration. **3** In the few-shot setting: Some of the most

408 aggressive sparse configurations, updating a mere 2% of parameters with both unstructured and  
409 group-wise sparsities, exhibit a performance delta of under 2% accuracy loss. This is paired with  
410 remarkable acceleration gains of 63.5% and 81.8%, respectively.