

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014

Visual Atoms: Pre-training Vision Transformers with Sinusoidal Waves

-Supplementary Material-

Anonymous CVPR submission

Paper ID 1946

1. Details of VisualAtom

Example images of VisualAtom (Figure 1). We list example VisualAtom images of randomly selected 50 classes with one instance per class in Figure 1a. These images are selected from VisualAtom generated with our baseline parameters setting. Also, we list example images of randomly selected 50 instances in one randomly selected class in Figure 1b. In Figure 2, with the similar format, we also list example images of RCDB [3], one of the previous FDSL datasets, for comparison. As you can see, VisualAtom has more variation in contour outlines between categories than RCDB.

Details on generating VisualAtom. We include separate scripts named ‘visual_atomic_renderer’ to generate VisualAtom in the supplementary material. Please see ‘visual_atomic_renderer/README.md’ for details on how to execute these scripts. We confirmed generating VisualAtom-1k took 40 minutes using parallel execution with 40 threads on 20 CPU cores.

2. Additional experimental results

Here, we describe our findings and insights into VisualAtom parameters that are relatively unimportant in contours. The findings with these parameters could not be included in the main paper due to space limitations. As in the main paper, we pre-train ViT-Tiny [2] on VisualAtom-1k when varying parameters, and compare of fine-tuning accuracy on three datasets: CIFAR10 (C10) [4], CIFAR100 (C100) [4] and ImageNet-100 (IN100)¹.

Orbit-Thickness Parameter (Table 1). In Table 1, we show the effects of the line thickness parameter on rendering. Here, we used the following configurations: fixed at 1, 3, or 5 pixel, and Random (orbit) or Random (line).

The first three configurations fix the line-thickness parameter at $l = 1, 3, \text{ and } 5$ pixel, respectively. The 4th/5th configurations randomly select l from $\{1, 2, 3, 4, 5\}$ per orbit/line, respectively. The findings show that the three-pixel

¹This is a subset of ImageNet [1] with 100 object categories.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Table 1. Fine-tuning accuracy when varying line thickness parameter l . Random-O and -L uses randomly sampled line thickness l at each orbit and line, respectively, where l is uniformly sampled from $\{1, 2, 3, 4, 5\}$ pixels.

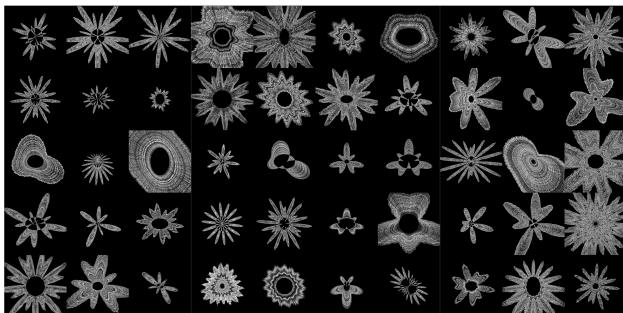
Value of l	C10	C100	IN100
1 pixel	97.6	84.9	90.3
3 pixel	97.6	85.7	90.1
5 pixel	97.6	84.5	89.9
Random-O	97.7	85.1	90.0
Random-L	97.6	85.4	89.8

line thickness greatly improved the performance on C100, but that it was not always best. For example, for IN100, the one-pixel line thickness was best. For simplicity, we used the one-pixel configuration for the baseline.

Table 2. Fine-tuning accuracy when using randomly phase difference parameters ϕ_1 and ϕ_2 at each categories, where ϕ_1 and ϕ_2 are uniformly sampled from the range of $[0, \pi]$.

Value of ϕ_1 and ϕ_2	C10	C100	IN100	
0, 0	97.6	84.9	90.3	
Random	97.5	84.8	90.1	
$\phi_1 : 0$		$\phi_1 : \pi/4$	$\phi_1 : \pi/2$	$\phi_1 : \pi$

Phase Difference Parameters (Table 2). In Eq. (5) of the main paper, introducing phase difference parameters $\phi_1, \phi_2 \in \mathbb{R}$ varies the phase of the two waves as follow:

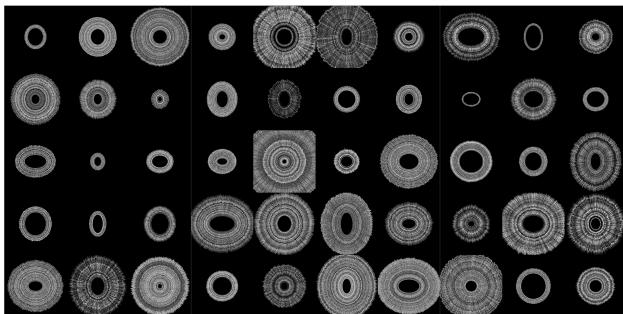
108
109
110
111
112
113
114
115
116
117
118

(a) VisualAtom Classes

119
120
121

Figure 1. Example images of VisualAtom. (a) Example images of randomly selected 50 classes with one instance per class. (b) Example images of randomly selected 50 instances in one randomly selected class.

122



(b) VisualAtom Instances

123
124
125
126
127
128
129
130
131
132
133

(a) RCDB Classes

134
135
136

Figure 2. Example images of RCDB [3]. (a) Example images of randomly selected 50 classes with one instance per class. (b) Example images of randomly selected 50 instances in one randomly selected class.

137

138

139

$$\Phi_k(\theta) = \lambda_1 \sin(n_1\theta + \phi_1) + \lambda_2 \sin(n_2\theta + \phi_2) + \eta\epsilon(\theta). \quad (5*)$$

140
141
142

We tried the phase difference parameters ϕ_1 and ϕ_2 defined for the two waves, randomly sampled from the range of $[0, \pi]$. Table 2 shows that VisualAtom has a robust pre-training effect on phase differences. For simplicity, we did not change the phase of two waves for the baseline.

143
144
145
146
147
148
149

3. Hyper-parameters in our experiments

150
151
152
153
154

For each experiments, hyper-parameters are based on the configuration used by Kataoka *et al.* [3]. More fundamentally, they are based on the paper proposing DeiT [5]. Table 3 shows hyper-parameters in our experiments.

155
156
157
158
159
160
161

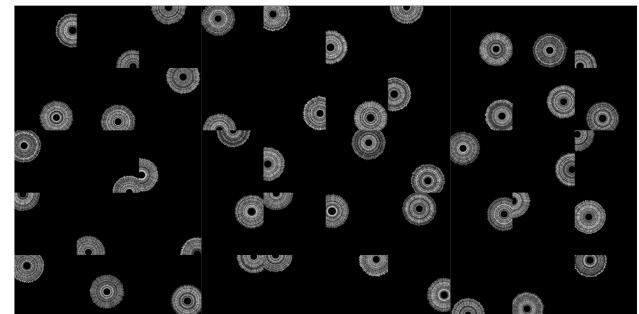
In addition, we almost used Github repository² published by Kataoka *et al.* for each experiments. Also, for pre-training on the large dataset such as VisualAtom-21k, we used WebDataset³ to accelerate IO processing. Note that

²Github repository of Kataoka *et al.* [3] : <https://github.com/masora1030/CVPR2022-Pretrained-ViT-PyTorch>

³WebDataset library : <https://github.com/webdataset/webdataset>



(b) VisualAtom Instances



(b) RCDB Instances

we changed only the parameter of the Warmup interval to 5k steps from 5 epochs used in the previous work. This is to apply Warmup at a fixed iterations regardless of the size of dataset.

References

- [1] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 1
- [2] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations*, 2021. 1
- [3] Hirokatsu Kataoka, Ryo Hayamizu, Ryosuke Yamada, Kodai Nakashima, Sora Takashima, Xinyu Zhang, Edgar Josafat Martinez-Noriega, Nakamasa Inoue, and Rio Yokota. Replacing Labeled Real-Image Datasets With Auto-Generated Contours. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21232–21241, 2022. 1, 2, 3

216 Table 3. Hyper-parameters of pre-training and fine-tuning in our experiments. Basically, they are same as the configuration used by 270
 217 Kataoka *et al.* [3]. 271

218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236	Training Step		Pre-training		Fine-tuning		272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290	
	Model Dataset Category	ViT-T		ViT-B		ViT-T/B		
		1k	21k	1k	21k	1k	Others	
Epochs		300	90	300	90	300	1000	276
Batch Size		1024	8192	1024	8192	1024	768	277
Optimizer	AdamW	AdamW	AdamW	AdamW	AdamW	AdamW	SGD	278
LR	1.0e-3	8.0e-3	1.0e-3	1.0e-3	1.0e-3	1.0e-3	1.0e-2	279
Weight Decay	0.05	0.05	0.05	0.05	0.05	0.05	1.0e-4	280
LR Scheduler	Cosine	Cosine	Cosine	Cosine	Cosine	Cosine	Cosine	281
Warmup Steps	5k	5k	5k	5k	5 (epochs)	10 (epochs)		282
Resolution	224	224	224	224	224/384		224	283
Label Smoothing	0.1	0.1	0.1	0.1	0.1	0.1		284
Drop Path	0.1	0.1	0.1	0.1	0.1	0.1		285
Rand Augment	9/0.5	9/0.5	9/0.5	9/0.5	9/0.5	9/0.5		286
Mixup	0.8	0.8	0.8	0.8	0.8	0.8		287
Cutmix	1.0	1.0	1.0	1.0	1.0	1.0		288
Erasing	0.25	0.25	0.25	0.25	0.25	0.25		289

- 237
 238 [4] Alex Krizhevsky, Geoffrey Hinton, et al. Learning Multiple 291
 239 Layers of Features from Tiny Images. 2009. 1 292
 240 [5] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco 293
 241 Massa, Alexandre Sablayrolles, and Herve Jegou. Training 294
 242 Data-efficient Image Transformers & Distillation through 295
 243 Attention. In *International Conference on Machine Learning*, 296
 244 volume 139, pages 10347–10357, July 2021. 2 297
 245 298
 246 299
 247 300
 248 301
 249 302
 250 303
 251 304
 252 305
 253 306
 254 307
 255 308
 256 309
 257 310
 258 311
 259 312
 260 313
 261 314
 262 315
 263 316
 264 317
 265 318
 266 319
 267 320
 268 321
 269 322
 270 323