

000 SUPPLEMENTARY MATERIAL FOR SUBMISSION:
001
002

003 RECONSTRUCTING TRAINING DATA FROM
004 REAL WORLD MODELS TRAINED WITH
005 TRANSFER LEARNING
006
007
008

009 **Anonymous authors**

010 Paper under double-blind review
011
012
013
014

015 1 FULL RESULTS FOR FIG.5B FROM MAIN PAPER
016

017 In Figure 5b in the main paper we show reconstructed samples, sorted according to their
018 reconstruction-quality as measured by cosine-similarity between their image embeddings. The
019 results shown there are obtained from a model trained on DINO embeddings on Food101 dataset.

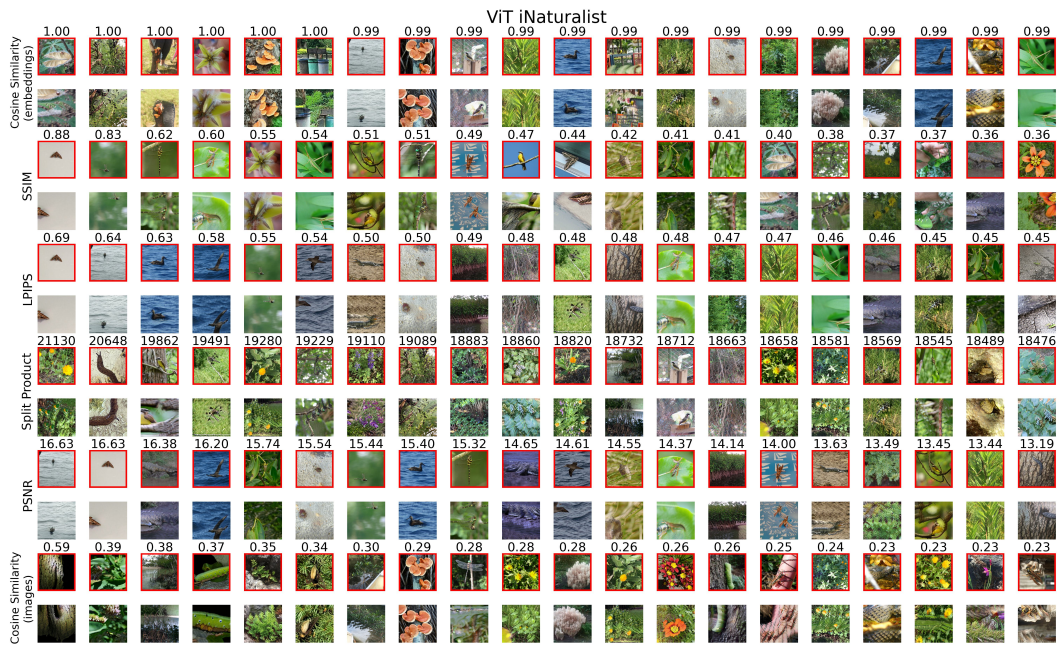
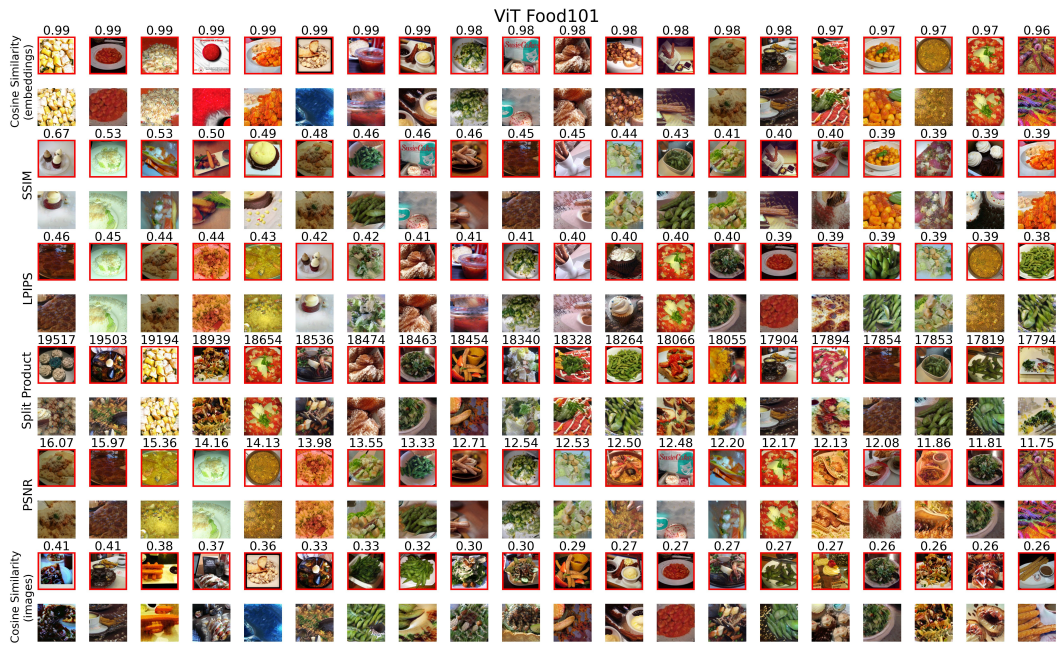
020 Below we provide the complete results for this type of evaluation, namely, for each model from Figure
021 3, and for each reconstruction-quality metric, we show the "best" reconstructed samples according to
022 this metric (by sorting them).
023

024 In total there are 8 models: trained on 4 backbones (ViT, DINO, DINOv2 and CLIP) and on 2
025 datasets (Food101 and iNaturalist), as described in details in the Results Section in the main paper.
026 And for each model we show the sorted results for a total of 6 choices for reconstruction-quality:
027 Cosine-Similarity in Embedding space plus 5 metrics in Image space: SSIM, LPIPS, Split-Product,
028 PSNR and Cosine-Similarity (Image space).

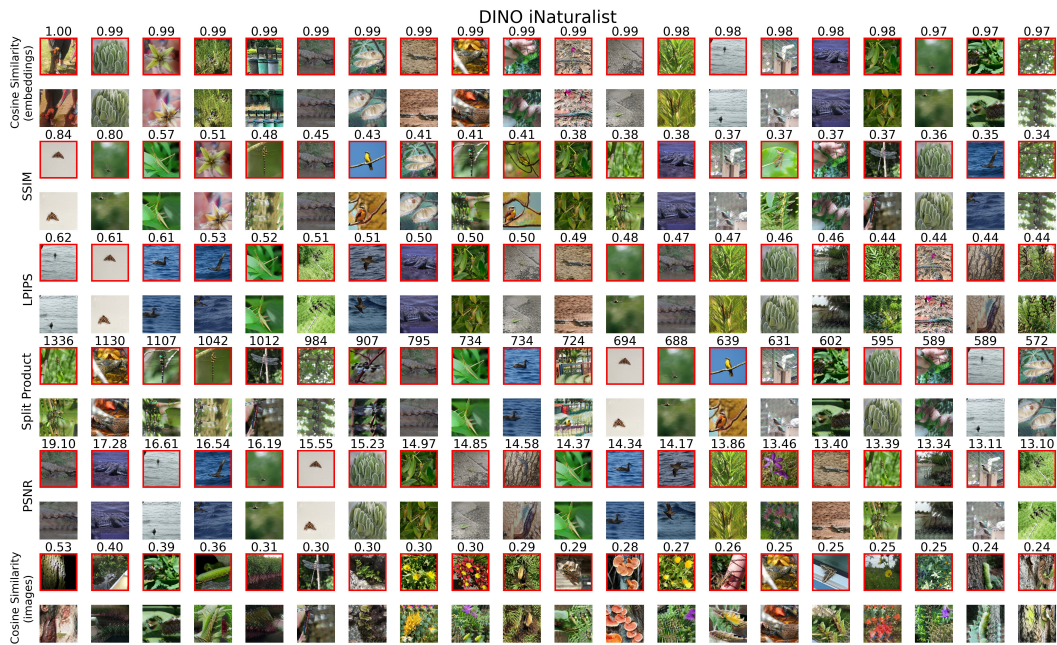
029 In each Figure below, images with RED borderline are original training images, and the image
030 below them is their nearest reconstructed image (as measured by the cosine-similarity between their
031 embeddings). Note: in all cases, the matching between a training image and its reconstructed image
032 is the same. the only difference is between the way they are sorted – which is done using the metric
033 (as written in the left side of each row).

034 We are aware that there may be some sampling bias in these results. However, since in our work, the
035 inversion part is time-intensive, we must choose which embeddings to invert, where in our work we
036 use the cosine-similarity for that, as detailed in the paper. It is not feasible to use an image-metric
037 for this goal, because this would mean first inverting all candidate embeddings (usually 25k-50k of
038 them), which is not feasible.
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

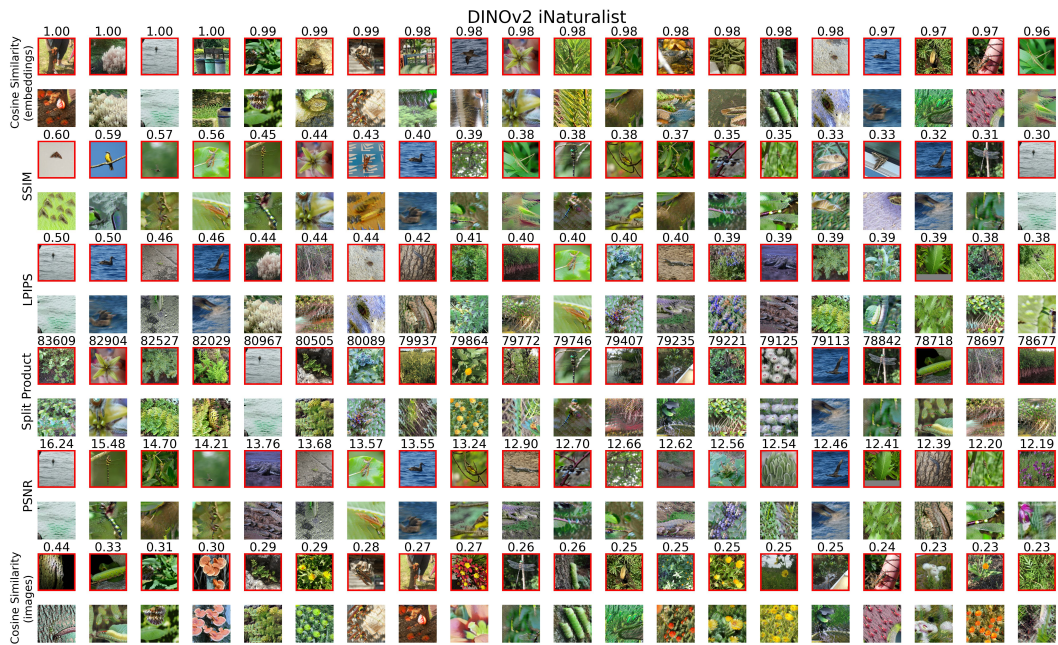
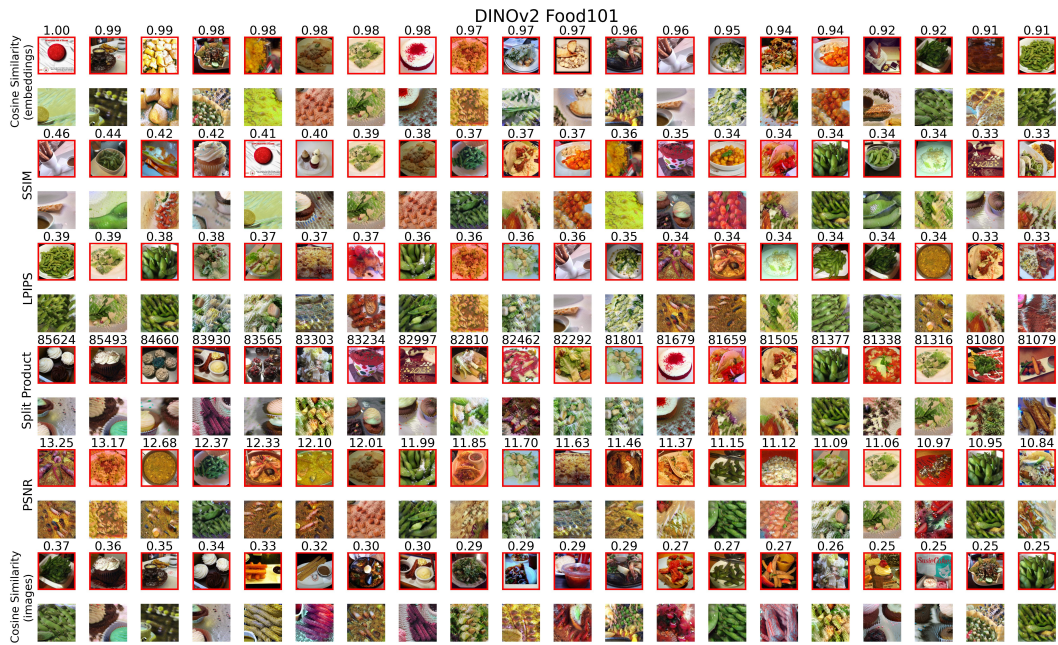
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107



108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161



162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215



216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

