

# Self-Supervised Material and Texture Representation Learning for Remote Sensing Tasks Supplementary Material

Peri Akiva  
Rutgers University  
peri.akiva@rutgers.edu

Matthew Purri  
Rutgers University  
matthew.purri@rutgers.edu

Matthew Leotta  
Kitware Inc  
matt.leotta@kitware.com

## 1. Implementation Details

### 1.1. Surface Encoding Layer

```
1
2 # Q: learned clusters, shape: [n_eps, D, 1, 1]
3 # feats: input image, shape: [b, D, 1, 1]
4 # scale: float scaling factor
5 # weights: learnable cluster weights
6
7 Q = torch.empty(n_eps, D, dtype=torch.float), requires_grad=True)
8 weights = torch.empty(n_eps, dtype=torch.float).uniform_(-1, 0), requires_grad=True)
9 def encode(feats):
10     res = weights*(feats-Q).pow(2).sum(dim=3)
11     res = Softmax(res, dim=2)
12     cluster_wise_residual_encoding = (res*(feats-Q)).sum(dim=1)
13     residual_encoding = cluster_wise_residual_encoding.mean(dim=1)
14     return residual_encoding
```

Listing 1. Surface Residual Encoding Simplified Pseudo-code

### 1.2. Texture Refinement Network

```
1
2 # x: input image
3 # feats: low level features
4 # num_iterations: number of refinement operations
5 # dilations: list of dilations to use
6
7 def refine(x, feats):
8     K = [] # list of kernels for each locations
9     for iter in range(num_iterations):
10         for dilation in dilations:
11             for all locations in x:
12                 # standard deviation of local kernel
13                 x_std = local_std(x, dilation) # float
14
15                 # cosine similarity of center pixel and neighboring pixels
16                 x_cos_similarity = cosine_similarity(x, dilation) # [h,w]
17
18                 # calculate weights in given location
19                 kernel = -x_cos_similarity/(0.00001 + x_std^2)
20                 kernel = Softmax(kernel.mean(dim=1))
21                 K.append(kernel)
22
23     # reduce K to match features dimensions
24     K.sum(dim=2)
25     feats *= K
26
27     return feats
```

Listing 2. Texture Refinement Network Simplified Pseudo-code

## 2. Additional Results

### 2.1. Running times

One of the challenges our method presents is the lengthy running time for each epoch during the pre-training stage. This stems from the fact that each pixels are now represented by  $h \times w$  sized crops, which are stacked along the batch dimension. In practice, this is equivalent to running the method on input of size of  $[h \times w, h, w]$ , where  $h \times w$  is the batch size, and  $h, w$  are the height and width of the crop. A given image of size  $32 \times 32$  is typically represented by  $32 \times 32 = 1024$  pixels, while in our method an image is represented by  $32 \times 32 \times 7 \times 7 = 50176$  pixels. That constraint translates to roughly ten times longer training time, as shown in Tab. 1. It is important to note that this limitation only occurs during the pre-training stage. During fine-tuning or inference steps, the method does not employ the window sampling approach, making it significantly faster than the pre-training stage, and similar to other fully supervised methods.

Method	supervision	crop-size	time/10 epochs	peak memory/GPU
<i>Onera Change Detection Experiment</i>				
DeepLab-v3 [1] (ImageNet)	$\mathcal{F}$	$96 \times 96$	0h24	24G
SeCo [3]	$\mathcal{S}$	$96 \times 96$	0h19	24G
SeCo [3]	$\mathcal{S} + \mathcal{F}$	$96 \times 96$	0h23	24G
Patch-wise Backbone	$\mathcal{S}$	$32 \times 32$	3h15	24G
Ours	$\mathcal{S}$	$32 \times 32$	5h48	24G
Ours (fine-tuned)	$\mathcal{S} + \mathcal{F}$	$96 \times 96$	0h24	24G
<i>SpaceNet Building Segmentation Task</i>				
DeepLab-v3 [1] (random)	$\mathcal{F}$	$128 \times 128$	0h21	24G
DeepLab-v3 [1] (ImageNet)	$\mathcal{F}$	$128 \times 128$	0h21	24G
Ours (fine-tuned)	$\mathcal{S} + \mathcal{F}$	$128 \times 128$	0h21	24G

Table 1. Training run time calculated for each model and its corresponding remote sensing task. The self-supervision timing is calculated when running on the Onera Satellite Change Detection (OSCD) dataset. During fine-tuning or inference steps, the method does not employ the window sampling approach, making it significantly faster than the pre-training stage, and similar to other fully supervised methods.

### 2.2. Architecture Selection

A natural direction given the findings of Tab. 1 is to reduce the size of the pre-training model. One consideration is the complexity of the learned downstream task once the pre-trained model is achieved. For complex downstream tasks, such as 3D reconstruction or semantic segmentation, an encoder with large learning capacity is needed. For simpler tasks, it may be feasible to use smaller encoder. Since the downstream tasks presented in this work vary in complexity, we selected to use ResNet-34 as our backbone, with the corresponding running times shown in Tab. 1. However, the pre-training step can be performed with smaller encoders, such as fully connected neural network, ResNet-18, or similar. In Tab. 2 we study the effects of model size on our method’s running times.

Backbone	parameters	crop-size	time/10 epochs	peak memory/GPU
3-layer NN	805,536	$32 \times 32$	0h18	24G
ResNet-18	15,428,544	$32 \times 32$	4h40	24G
ResNet-34	25,536,704	$32 \times 32$	5h48	24G
ResNet-50	39,163,328	$32 \times 32$	6h41	24G
ResNet-101	58,155,456	$32 \times 32$	9h19	24G

Table 2. **Training** run time calculated for various backbones when using the Onera Satellite Change Detection (OSCD) dataset.

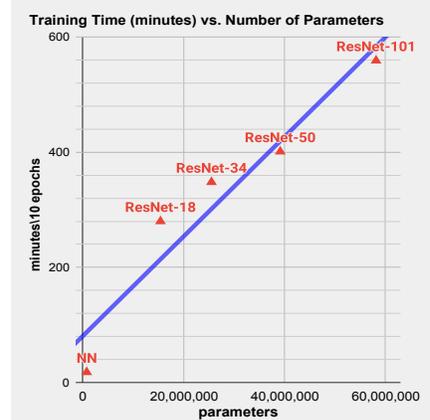


Figure 1. Plot of training time (in minutes per 10 training epochs) with respect to number of learnable model parameters.

### 2.3. Qualitative Results

We include additional qualitative results on Onera in Fig. 3 and SpaceNet in Fig 2.

### 2.4. Evaluating Visual Word Maps and Cluster Scatters

As mentioned in the paper, we measure the effectiveness of our approach to describe materials and textures through qualitative evaluation of visual word maps (pixel-wise cluster assignments) generated by our method. Ideally, we expect similar material and textures to be mapped to the same clusters, without over or under grouping of pixels. We also desire to achieve distinctly mapped clusters with minimal overlaps. We visually compare classical textons, a patch-wise backbone, and our method with corresponding scatter plots (t-SNE, perplexity=100, steps=5000) in Fig. 4 and 5. The patch-wise backbone has the same base architecture as MATTER, but without TeRN and surface residual encoding modules. Both methods were trained on the same dataset, with the same hyperparameters, and number of iterations, as described in the main paper. It can be seen that the textons and patch-wise backbone approaches generate two extreme cases of over-sensitivity and under-sensitivity to changes in material and texture. Since textons operate on raw intensity values, the inter-material variance is small, making it highly sensitive to small texture changes. This can be seen in the textons-generated visual word map, in which small irregularities on the road results in mapping to different visual words. On the other hand, the patch-wise backbone, even with the receptive field constraints through patched inputs, still loses crucial low-level details essential for texture representation. This is indicated by the grouping of obviously different textures to a single visual word. In contrast our method is able to retain texture-essential features, and generalize surface representation which translates to superior surface-based visual word maps. The superiority in surface-representation learning can also be seen in the scatter plots, in which our method generates distinct cluster scatters corresponding to different texture and material combinations.

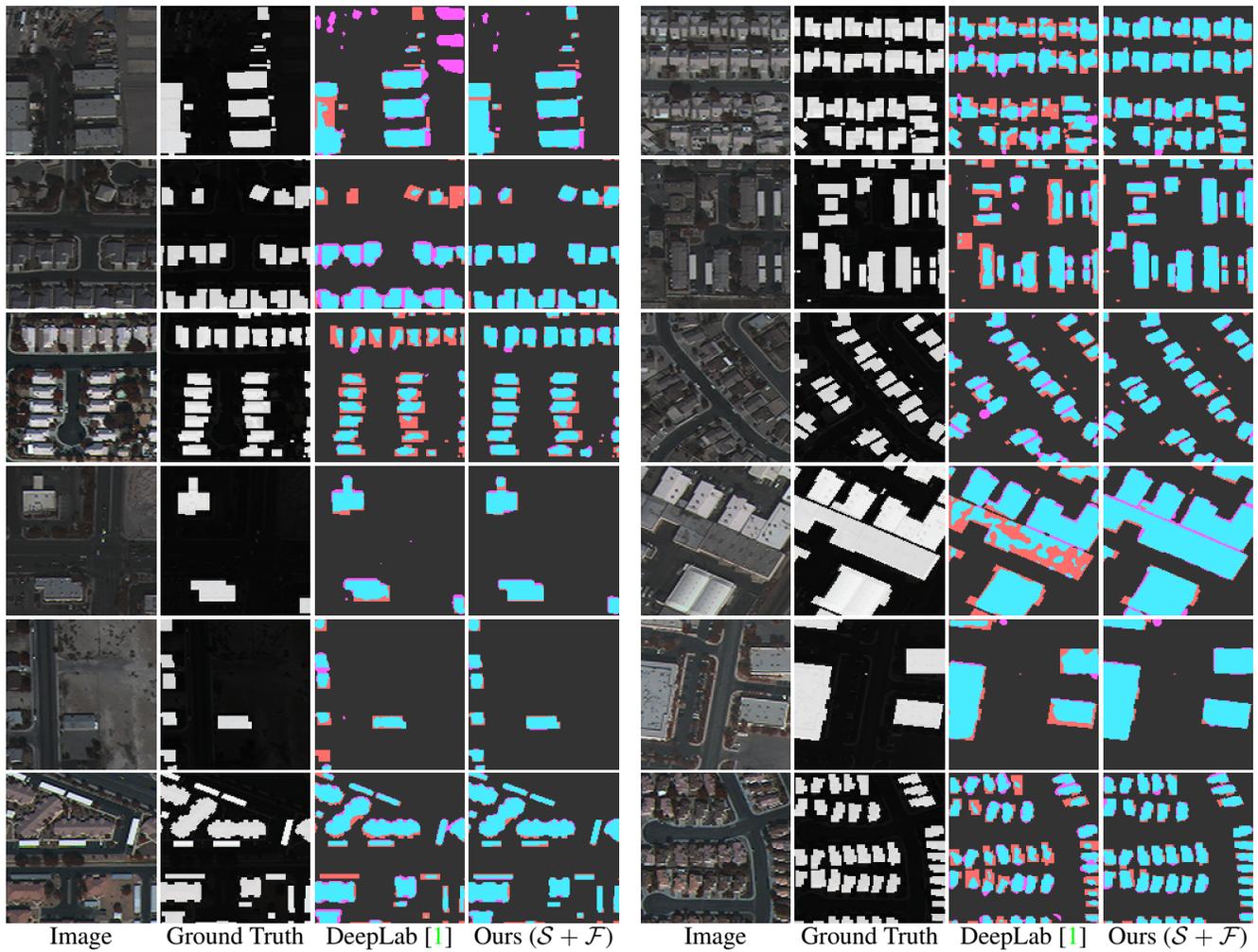


Figure 2. **Qualitative results** of our method on SpaceNet dataset [4]. Cyan, magenta, gray, and red colors represent true positive, false positive, true negative, and false negative respectively. Best viewed in zoom and color.

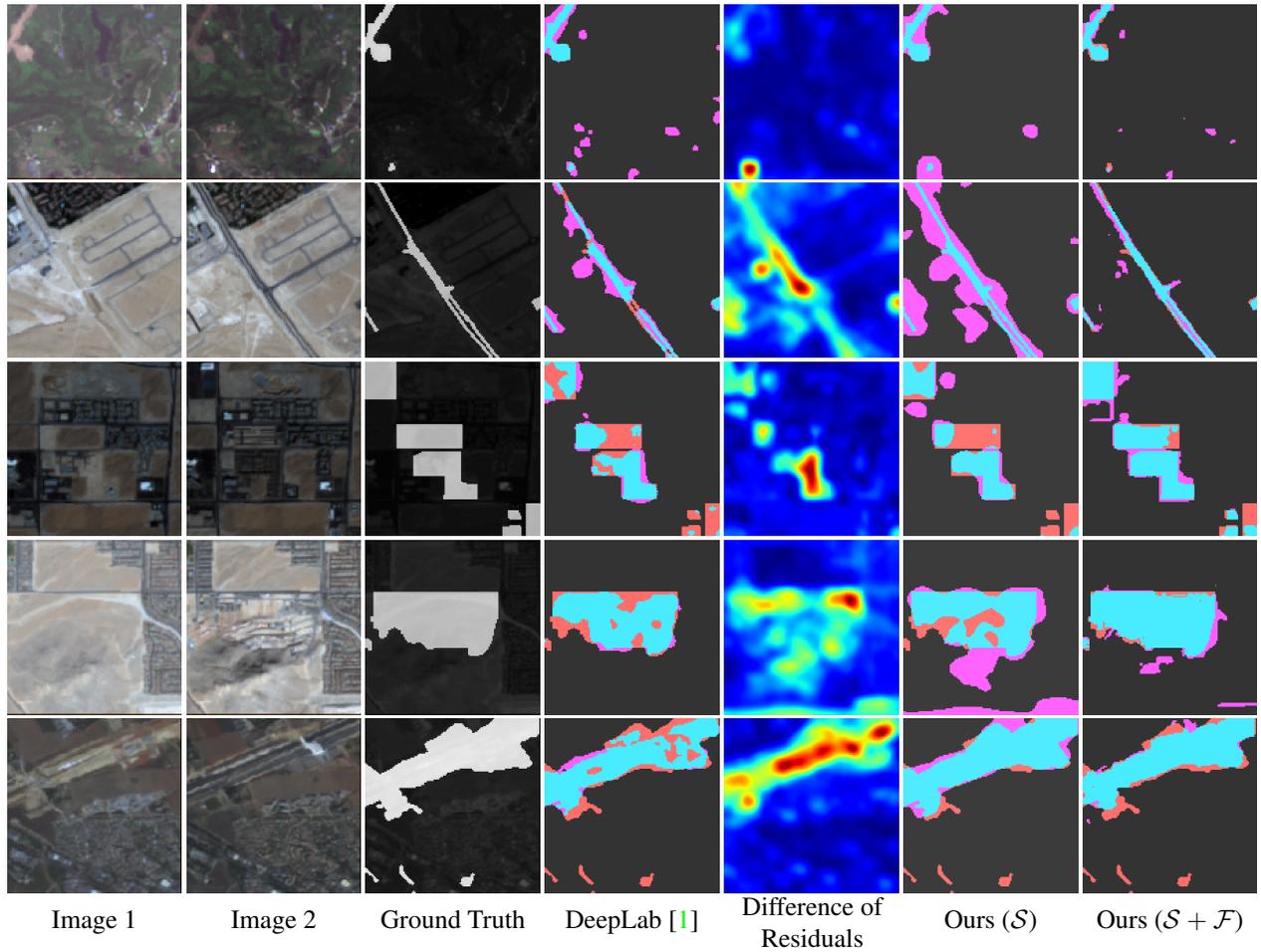


Figure 3. **Qualitative results** of our method on Onera Satellite Change Detection (OSCD) dataset [2]. It can be seen that our self-supervised alone is capable of detecting change only by inferring on the change of material and texture. The fine-tuned model is able to utilize the pre-trained material and texture based weights and achieve significantly better results than models with ImageNet weight initialization. Cyan, magenta, gray, and red colors represent true positive, false positive, true negative, and false negative respectively. Best viewed in zoom and color.

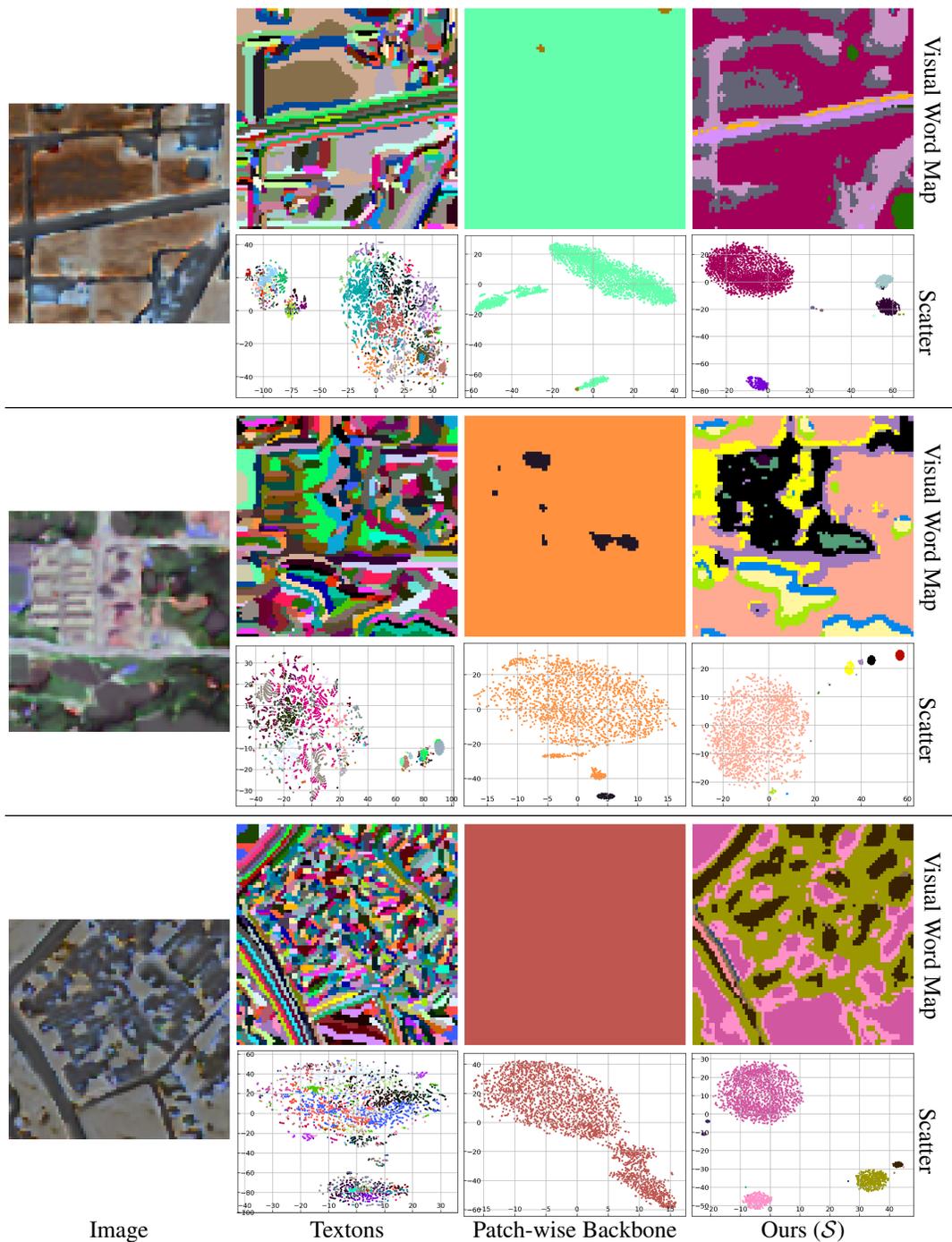


Figure 4. **Qualitative evaluation** of our generated material and texture based visual word maps and corresponding clusters. It can be seen that our method provides more descriptive surface-based features that are not highly sensitive to small texture irregularities like textons, or under-sensitive to structure changes like the patch-wise backbone. It can also be seen that our method provides more distinct cluster centers, better differentiating material and texture combinations. We randomly sample 30% of the pixels for the scatter plot visualization. Best viewed in zoom and color. Colors are random.

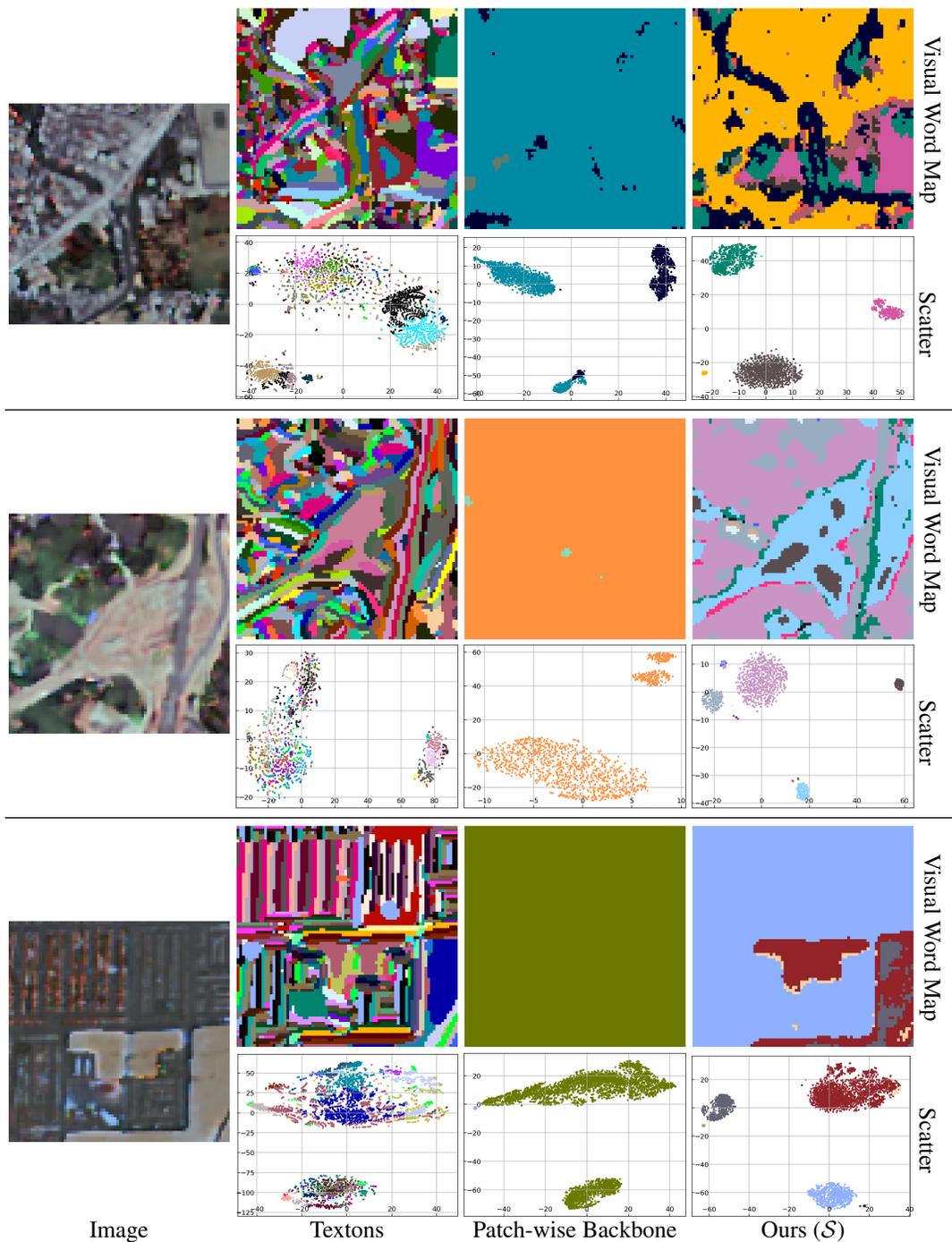


Figure 5. **Qualitative evaluation** of our generated material and texture based visual word maps and corresponding clusters. It can be seen that our method provides more descriptive surface-based features that are not highly sensitive to small texture irregularities like textons, or under-sensitive to structure changes like the patch-wise backbone. It can also be seen that our method provides more distinct cluster centers, better differentiating material and texture combinations. We randomly sample 30% of the pixels for the scatter plot visualization. Best viewed in zoom and color. Colors are random.

### 3. Self-Collected Dataset

#### 3.1. Points of Interest

```
1 {
2   "change_region": "no_change",
3   "region_name": "usa_wyoming_shoshone",
4   "geo_coords": [
5     "43.943477_-109.492681"
6   ],
7   "collection": "sentinel-s2-l2a-cogs",
8   "start_date": "01/01/2017",
9   "end_date": "01/01/2020",
10  "crop_size": 1096,
11  "radius": 0.001,
12  "max_cloud_cover": 20,
13  "min_data_coverage": 80,
14  "min_num_images": 2,
15  "max_num_images": 100,
16 }
17
18 {
19   "change_region": "no_change",
20   "region_name": "vietnam_00",
21   "geo_coords": [
22     "19.783545_105.355797"
23   ],
24   "collection": "sentinel-s2-l2a-cogs",
25   "start_date": "01/01/2017",
26   "end_date": "01/01/2020",
27   "crop_size": 1096,
28   "radius": 0.001,
29   "max_cloud_cover": 20,
30   "min_data_coverage": 80,
31   "min_num_images": 2,
32   "max_num_images": 100,
33 }
34
35 {
36   "change_region": "no_change",
37   "region_name": "yemen_asadi-al-faya",
38   "geo_coords": [
39     "14.928610_50.284421"
40   ],
41   "collection": "sentinel-s2-l2a-cogs",
42   "start_date": "01/01/2017",
43   "end_date": "01/01/2020",
44   "crop_size": 1096,
45   "radius": 0.001,
46   "max_cloud_cover": 20,
47   "min_data_coverage": 80,
48   "min_num_images": 2,
49   "max_num_images": 100,
50 }
51
52 {
53   "change_region": "no_change",
54   "region_name": "bulgaria_ozizare",
55   "geo_coords": [
56     "42.637424_25.953441"
57   ],
58   "collection": "sentinel-s2-l2a-cogs",
59   "start_date": "01/01/2017",
60   "end_date": "01/01/2020",
61   "crop_size": 1096,
62   "radius": 0.001,
63   "max_cloud_cover": 20,
```

```
64 "min_data_coverage": 80,
65 "min_num_images": 2,
66 "max_num_images": 100,
67 }
68
69 {
70 "change_region": "no_change",
71 "region_name": "chile_aqi-keqiaole",
72 "geo_coords": [
73 "40.947455_85.205265"
74 ],
75 "collection": "sentinel-s2-l2a-cogs",
76 "start_date": "01/01/2017",
77 "end_date": "01/01/2020",
78 "crop_size": 1096,
79 "radius": 0.001,
80 "max_cloud_cover": 20,
81 "min_data_coverage": 80,
82 "min_num_images": 2,
83 "max_num_images": 100,
84 }
85
86 {
87 "change_region": "no_change",
88 "region_name": "china_yuxiancun",
89 "geo_coords": [
90 "42.939327_124.852949"
91 ],
92 "collection": "sentinel-s2-l2a-cogs",
93 "start_date": "01/01/2017",
94 "end_date": "01/01/2020",
95 "crop_size": 1096,
96 "radius": 0.001,
97 "max_cloud_cover": 20,
98 "min_data_coverage": 80,
99 "min_num_images": 2,
100 "max_num_images": 100,
101 }
102
103 {
104 "change_region": "no_change",
105 "region_name": "india_sundarban",
106 "geo_coords": [
107 "22.044656_88.757307"
108 ],
109 "collection": "sentinel-s2-l2a-cogs",
110 "start_date": "01/01/2017",
111 "end_date": "01/01/2020",
112 "crop_size": 1096,
113 "radius": 0.001,
114 "max_cloud_cover": 20,
115 "min_data_coverage": 80,
116 "min_num_images": 2,
117 "max_num_images": 100,
118 }
119
120 {
121 "change_region": "no_change",
122 "region_name": "maludam_indonesia",
123 "geo_coords": [
124 "1.506800_111.165672"
125 ],
126 "collection": "sentinel-s2-l2a-cogs",
127 "start_date": "01/01/2017",
128 "end_date": "01/01/2020",
129 "crop_size": 1096,
130 "radius": 0.001,
```

```
131 "max_cloud_cover": 20,
132 "min_data_coverage": 80,
133 "min_num_images": 2,
134 "max_num_images": 100,
135 }
136
137 {
138 "change_region": "no_change",
139 "region_name": "iran_naghab",
140 "geo_coords": [
141 "32.383612_59.689489"
142 ],
143 "collection": "sentinel-s2-l2a-cogs",
144 "start_date": "01/01/2017",
145 "end_date": "01/01/2020",
146 "crop_size": 1096,
147 "radius": 0.001,
148 "max_cloud_cover": 20,
149 "min_data_coverage": 80,
150 "min_num_images": 2,
151 "max_num_images": 100,
152 }
153
154 {
155 "change_region": "no_change",
156 "region_name": "kazakhstan_kapaxar_karazhal",
157 "geo_coords": [
158 "47.981633_70.952032"
159 ],
160 "collection": "sentinel-s2-l2a-cogs",
161 "start_date": "01/01/2017",
162 "end_date": "01/01/2020",
163 "crop_size": 1096,
164 "radius": 0.001,
165 "max_cloud_cover": 20,
166 "min_data_coverage": 80,
167 "min_num_images": 2,
168 "max_num_images": 100,
169 }
170
171 {
172 "change_region": "no_change",
173 "region_name": "madagascar_00",
174 "geo_coords": [
175 "-18.969615_49.046689"
176 ],
177 "collection": "sentinel-s2-l2a-cogs",
178 "start_date": "01/01/2017",
179 "end_date": "01/01/2020",
180 "crop_size": 1096,
181 "radius": 0.001,
182 "max_cloud_cover": 20,
183 "min_data_coverage": 80,
184 "min_num_images": 2,
185 "max_num_images": 100,
186 }
187
188 {
189 "change_region": "no_change",
190 "region_name": "mali_fodebougou",
191 "geo_coords": [
192 "12.533564_-9.805519"
193 ],
194 "collection": "sentinel-s2-l2a-cogs",
195 "start_date": "01/01/2017",
196 "end_date": "01/01/2020",
197 "crop_size": 1096,
```

```
198 "radius": 0.001,
199 "max_cloud_cover": 20,
200 "min_data_coverage": 80,
201 "min_num_images": 2,
202 "max_num_images": 100,
203 }
204
205 {
206 "change_region": "no_change",
207 "region_name": "mexico_santa-elena",
208 "geo_coords": [
209 "23.636026_-104.527631"
210 ],
211 "collection": "sentinel-s2-l2a-cogs",
212 "start_date": "01/01/2017",
213 "end_date": "01/01/2020",
214 "crop_size": 1096,
215 "radius": 0.001,
216 "max_cloud_cover": 20,
217 "min_data_coverage": 80,
218 "min_num_images": 2,
219 "max_num_images": 100,
220 }
221
222 {
223 "change_region": "no_change",
224 "region_name": "mongolia_ronin-hur",
225 "geo_coords": [
226 "48.459630_110.512820"
227 ],
228 "collection": "sentinel-s2-l2a-cogs",
229 "start_date": "01/01/2017",
230 "end_date": "01/01/2020",
231 "crop_size": 1096,
232 "radius": 0.001,
233 "max_cloud_cover": 20,
234 "min_data_coverage": 80,
235 "min_num_images": 2,
236 "max_num_images": 100,
237 }
238
239 {
240 "change_region": "no_change",
241 "region_name": "morocco_brarha",
242 "geo_coords": [
243 "34.453855_-4.287849"
244 ],
245 "collection": "sentinel-s2-l2a-cogs",
246 "start_date": "01/01/2017",
247 "end_date": "01/01/2020",
248 "crop_size": 1096,
249 "radius": 0.001,
250 "max_cloud_cover": 20,
251 "min_data_coverage": 80,
252 "min_num_images": 2,
253 "max_num_images": 100,
254 }
255
256 {
257 "change_region": "no_change",
258 "region_name": "myanmar_00",
259 "geo_coords": [
260 "18.233579_97.030426"
261 ],
262 "collection": "sentinel-s2-l2a-cogs",
263 "start_date": "01/01/2017",
264 "end_date": "01/01/2020",
```

```
265 "crop_size": 1096,
266 "radius": 0.001,
267 "max_cloud_cover": 20,
268 "min_data_coverage": 80,
269 "min_num_images": 2,
270 "max_num_images": 100,
271 }
272
273 {
274 "change_region": "no_change",
275 "region_name": "north-korea_00",
276 "geo_coords": [
277 "40.282847_128.339383"
278 ],
279 "collection": "sentinel-s2-l2a-cogs",
280 "start_date": "01/01/2017",
281 "end_date": "01/01/2020",
282 "crop_size": 1096,
283 "radius": 0.001,
284 "max_cloud_cover": 20,
285 "min_data_coverage": 80,
286 "min_num_images": 2,
287 "max_num_images": 100,
288 }
289
290 {
291 "change_region": "no_change",
292 "region_name": "qatar_quasil",
293 "geo_coords": [
294 "25.282863_50.784150"
295 ],
296 "collection": "sentinel-s2-l2a-cogs",
297 "start_date": "01/01/2017",
298 "end_date": "01/01/2020",
299 "crop_size": 1096,
300 "radius": 0.001,
301 "max_cloud_cover": 20,
302 "min_data_coverage": 80,
303 "min_num_images": 2,
304 "max_num_images": 100,
305 }
306
307 {
308 "change_region": "no_change",
309 "region_name": "russia_burukan",
310 "geo_coords": [
311 "53.061472_136.035468"
312 ],
313 "collection": "sentinel-s2-l2a-cogs",
314 "start_date": "01/01/2017",
315 "end_date": "01/01/2020",
316 "crop_size": 1096,
317 "radius": 0.001,
318 "max_cloud_cover": 20,
319 "min_data_coverage": 80,
320 "min_num_images": 2,
321 "max_num_images": 100,
322 }
323
324 {
325 "change_region": "no_change",
326 "region_name": "saudi-arabia_bir-jaydah",
327 "geo_coords": [
328 "36.139782_37.488645"
329 ],
330 "collection": "sentinel-s2-l2a-cogs",
331 "start_date": "01/01/2017",
```

```
332 "end_date": "01/01/2020",
333 "crop_size": 1096,
334 "radius": 0.001,
335 "max_cloud_cover": 20,
336 "min_data_coverage": 80,
337 "min_num_images": 2,
338 "max_num_images": 100,
339 }
340
341 {
342 "change_region": "no_change",
343 "region_name": "south-sudan_boma",
344 "geo_coords": [
345 "6.077807_34.132967"
346 ],
347 "collection": "sentinel-s2-l2a-cogs",
348 "start_date": "01/01/2017",
349 "end_date": "01/01/2020",
350 "crop_size": 1096,
351 "radius": 0.001,
352 "max_cloud_cover": 20,
353 "min_data_coverage": 80,
354 "min_num_images": 2,
355 "max_num_images": 100,
356 }
357
358 {
359 "change_region": "no_change",
360 "region_name": "spain_urda",
361 "geo_coords": [
362 "39.455428_-3.730460"
363 ],
364 "collection": "sentinel-s2-l2a-cogs",
365 "start_date": "01/01/2017",
366 "end_date": "01/01/2020",
367 "crop_size": 1096,
368 "radius": 0.001,
369 "max_cloud_cover": 20,
370 "min_data_coverage": 80,
371 "min_num_images": 2,
372 "max_num_images": 100,
373 }
374
375 {
376 "change_region": "no_change",
377 "region_name": "usa_nevada_tonopah",
378 "geo_coords": [
379 "38.240222_-117.320516"
380 ],
381 "collection": "sentinel-s2-l2a-cogs",
382 "start_date": "01/01/2017",
383 "end_date": "01/01/2020",
384 "crop_size": 1096,
385 "radius": 0.001,
386 "max_cloud_cover": 20,
387 "min_data_coverage": 80,
388 "min_num_images": 2,
389 "max_num_images": 100,
390 }
391
392 {
393 "change_region": "no_change",
394 "region_name": "usa_new-york_fabius",
395 "geo_coords": [
396 "42.808692_-75.985309"
397 ],
398 "collection": "sentinel-s2-l2a-cogs",
```

```

399   "start_date": "01/01/2017",
400   "end_date": "01/01/2020",
401   "crop_size": 1096,
402   "radius": 0.001,
403   "max_cloud_cover": 20,
404   "min_data_coverage": 80,
405   "min_num_images": 2,
406   "max_num_images": 100,
407 }
408
409 {
410   "change_region": "no_change",
411   "region_name": "usa_north-carolina_mt-mitchell",
412   "geo_coords": [
413     "35.765959_-82.254093"
414   ],
415   "collection": "sentinel-s2-l2a-cogs",
416   "start_date": "01/01/2017",
417   "end_date": "01/01/2020",
418   "crop_size": 1096,
419   "radius": 0.001,
420   "max_cloud_cover": 20,
421   "min_data_coverage": 80,
422   "min_num_images": 2,
423   "max_num_images": 100,
424 }
425
426 {
427   "change_region": "no_change",
428   "region_name": "usa_oklahoma_tahlequah",
429   "geo_coords": [
430     "36.035096_-94.811176"
431   ],
432   "collection": "sentinel-s2-l2a-cogs",
433   "start_date": "01/01/2017",
434   "end_date": "01/01/2020",
435   "crop_size": 1096,
436   "radius": 0.001,
437   "max_cloud_cover": 20,
438   "min_data_coverage": 80,
439   "min_num_images": 2,
440   "max_num_images": 100,
441 }
442
443 {
444   "change_region": "no_change",
445   "region_name": "usa_texas_dumont",
446   "geo_coords": [
447     "33.800483_-100.563941"
448   ],
449   "collection": "sentinel-s2-l2a-cogs",
450   "start_date": "01/01/2017",
451   "end_date": "01/01/2020",
452   "crop_size": 1096,
453   "radius": 0.001,
454   "max_cloud_cover": 20,
455   "min_data_coverage": 80,
456   "min_num_images": 2,
457   "max_num_images": 100,
458 }

```

## References

- [1] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. [2](#), [4](#), [5](#)

- [2] Rodrigo Caye Daudt, Bertr Le Saux, Alexandre Boulch, and Yann Gousseau. Urban change detection for multispectral earth observation using convolutional neural networks. In *IGARSS 2018 - 2018 IEEE International Geoscience and Remote Sensing Symposium*, pages 2115–2118, 2018. [5](#)
- [3] Oscar Manas, Alexandre Lacoste, Xavier Giro-i Nieto, David Vazquez, and Pau Rodriguez. Seasonal contrast: Unsupervised pre-training from uncurated remote sensing data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9414–9423, October 2021. [2](#)
- [4] Adam Van Etten, Dave Lindenbaum, and Todd M Bacastow. Spacenet: A remote sensing dataset and challenge series. *arXiv preprint arXiv:1807.01232*, 2018. [4](#)