

1 Dear reviewers and editors,

2
3 Thank you very much for your helpful and constructive feedback. In response, we have made clarifications
4 and improvements to the manuscript, both in contents and in visual presentation, which we believe
5 enhance the clarity and quality of the study. Detailed, point-by-point responses are shown below, with
6 reference to the specific line numbers of tracked changes (“All-Markup” version) made in the updated
7 manuscript.

8 **Reviewer 1’s Comments**

9 **This paper proposed an object-based crop type ground truth dataset from 2013 to 2023, called**
10 **“CropSight-US”, using street view and Sentinel-2 images. This dataset covered 17 major crop types**
11 **across 294 Agricultural Statistics Districts (ASDs) in CONUS. Specifically, crop type labels were**
12 **extracted from Google street view images and field boundary information was derived from**
13 **Sentinel-2 imagery. CONUS-UncertainFusionNet was developed to generate this dataset with**
14 **uncertainty information. In the experimental results, the performance of several deep learning-**
15 **based networks was compared. This paper is well-organized and the reviewer has the following**
16 **detailed comments:**

17
18 **Response:** We sincerely thank Reviewer 1 for the thorough evaluation of our manuscript and for the clear
19 summary of our work. We appreciate the recognition of the study’s organization and the constructive
20 feedback provided. Below, we address each of the reviewer’s specific comments point by point and
21 describe the corresponding revisions made in the manuscript.

22
23 Specific Comments:

24
25 **1. This paper included that “CropSight-US, the first national scale, object-based crop type ground**
26 **truth dataset”. The authors were encouraged to review the existing products to highlight this**
27 **contribution.**

28
29 **Response:** Thank you for the opportunity to clarify the position of our contribution. We have revised the
30 Introduction to better situate CropSight-US within the landscape of existing national agricultural datasets.
31 Specifically, we have reviewed current crop type ground truth datasets in [Lines 46–70], categorizing them
32 into traditional field survey-based and crop product-derived ground truth datasets.

33
34 We further clarify the technical limitations of existing datasets [Lines 71-92]. Most available datasets are
35 primarily pixel-level, and the relatively few object-based products are typically generated through one of
36 two approaches: 1) aligning pixel-level classifications with external field boundary datasets, or 2)
37 clustering classified pixels into polygons (i.e., objects) based on spatial and spectral similarity, such as the
38 Crop Sequence Boundary (CSB). Both approaches introduce additional post-processing steps that can
39 affect scalability and boundary accuracy, thereby influencing the quality and consistency of object-based
40 products at large spatial scales.

41
42 Finally, to avoid potential overstatement, we have replaced “the first” with “novel” in both the Abstract
43 [Lines 19–21] and Introduction [Line 95]. This adjustment more accurately reflects the unique scope and

44 methodological contribution of our dataset without disregarding the existence of smaller-scale or
45 differently processed object-based products.

46

47 **2. This article contains multiple instances of the full name and abbreviation for CONUS. It is**
48 **recommended to provide the full name only upon first mention, using the abbreviation thereafter.**

49

50 **Response:** Thank you for this helpful suggestion. We have updated the manuscript so that the full name of
51 CONUS is provided only at its first occurrence (excluding the abstract or the captions of tables and
52 figures), which appears at [Line 96] in the Introduction section of the manuscript. The abbreviation is used
53 consistently throughout the remainder of the manuscript.

54

55 **3. For line 138 and 655, these links were invalid. Please double check all the links in the manuscript.**

56

57 **Response:** Thank you for checking the links and bringing these issues to our attention. We appreciate your
58 careful review and have corrected the invalid links as noted.

59

60 Regarding Line 138 (now [Line 145]): The documentation of all Street View-related APIs has been
61 consolidated into a single page. Previously, we used an older link to showcase the endpoint API street view
62 imagery downloading. Apologies for the confusion. We have now updated the link to point directly to the
63 Google Street View Static API documentation page at:

64 <https://developers.google.com/maps/documentation/streetview/overview>

65

66 Regarding the link to our CropSight-US GEE application (now [Line 675]): We believe the issue occurred
67 during the conversion from Word to PDF, where the hyphen ("-") in the URL was omitted (the GEE app is
68 managed under the google cloud project [ee-azzhou249], as the same link is also provided in [Line 735]
69 for the Data Availability section). The correct full address should be:

70 <https://ee-azzhou249.projects.earthengine.app/view/cropsight-us>

71

72 We have verified and corrected these links in the updated manuscript.

73

74 **4. What are the advantages of using Sentinel-2 over NAIP imagery?**

75

76 **Response:** Thank you for raising this question. The advantages of using Sentinel-2 over NAIP imagery in
77 our workflow mainly relate to timeliness and transferability. We have revised the manuscript to explicitly
78 clarify this rationale and to describe how Sentinel-2 availability and NAIP supplementation were handled
79 in our data processing workflow [Lines 162-164, Line 167 and Lines 276-278].

80

81 Each Google Street View imagery captures crop conditions during a specific month of a given year. To
82 correctly delineate the corresponding field boundary for that same period, we need satellite imagery that is
83 both temporally aligned and cloud-free. Sentinel-2 provides frequent revisit cycles (5 days at
84 mid-latitudes) with Sentinel-2A launched in 2015 and Sentinel-2B in 2017, which greatly increases the
85 likelihood of obtaining a cloud-free or minimally cloudy scene close to the date of the Street View image.
86 In contrast, NAIP imagery is typically collected only once every 1–3 years, which makes it difficult to
87 match the exact year and growing season represented in the Street View photos. Timely alignment is

88 essential for ensuring that the field boundary accurately reflects the field configuration at the time of
89 observation.

90

91 Additionally, while NAIP offers very high spatial resolution, it is available only for the United States. For
92 applications beyond the U.S., using NAIP would not be feasible. Sentinel-2, on the other hand, provides
93 global coverage with consistent acquisition schedules and harmonized data quality. This makes the
94 workflow more flexible and transferable, enabling extension of the CropSight framework to other
95 countries or regions without relying on region-specific aerial programs.

96

97 In CropSight-US, each object-based ground truth unit includes both the crop-type label identified from
98 Street View imagery and its corresponding field boundary derived from the closest suitable Sentinel-2
99 scene. This ensures that the final dataset is both temporally consistent and scalable to broader geographical
100 contexts. Only for years prior to full Sentinel-2 availability (before 2017) or when suitable cloud-free
101 imagery was not available, NAIP imagery was used as a supplementary source.

102

103 **5. What is the experimental environment of ViT-B16 and ResNet-50 training, such as GPU and
104 CPU? The authors are suggested to offer more details about the setting of the proposed network.**

105

106 **Response:** Thank you for asking this question. All model training was conducted on the Illinois Campus
107 Cluster Program (ICCP) at the University of Illinois Urbana-Champaign. The computational infrastructure
108 consists of dual AMD EPYC 7763 processors, 512 GB of system memory, and four NVIDIA A100 GPUs
109 with 80 GB of GPU memory each, connected through 25G high-speed networking nodes. We have updated
110 the relevant section of the manuscript to clearly document these specifications for both benchmark models
111 (ViT-B16 and ResNet-50) as well as our proposed CONUS-UncertainFusionNet, ensuring that the
112 experimental setup is transparent and reproducible [Lines 377-380].

113

114 **6. For all the tables, it is recommended to display the best performance in bold to help readers better
115 follow.**

116

117 **Response:** Thank you for this suggestion. We have updated the manuscript to highlight the tabulated rows
118 of best performance in bold to enhance the readability.

119

120 **7. SAM is trained on high-resolution natural images, and Sentinel-2 has a lower resolution. How to
121 overcome this resolution shift?**

122

123 **Response:** Thank you and we appreciate the opportunity to clarify this point and updated the manuscript to
124 provide further explanation [Lines 392-396]. To address the resolution shift between SAM's pre-training
125 domain (high-resolution natural images) and Sentinel-2 imagery, we fine-tune the mask decoder of SAM
126 using the training split of the field-boundary component of the CropGSV-Ref dataset, while keeping both
127 the image encoder and prompt encoder frozen (Section 3.1.4). This lightweight fine-tuning strategy has
128 been shown to be an effective and computationally efficient way to adapt SAM to remote sensing tasks,
129 including cropland field-boundary segmentation, without requiring full model retraining. Prior studies have
130 demonstrated that decoder-only adaptation substantially improves SAM's performance on delineating
131 agricultural boundaries while maintaining low computational cost (Liu et al., 2024; Pu et al., 2025). We

132 follow this established approach to ensure that SAM can better capture field-level boundaries from
133 Sentinel-2 imagery despite its lower spatial resolution.

134

135 Reference:

136 Liu, Y., Diao, C., Mei, W., and Zhang, C.: CropSight: Towards a large-scale operational framework for
137 object-based crop type ground truth retrieval using street view and PlanetScope satellite imagery, ISPRS J.
138 Photogramm. Remote Sens., 216, 66–89, <https://doi.org/10.1016/j.isprsjprs.2024.07.025>, 2024.

139

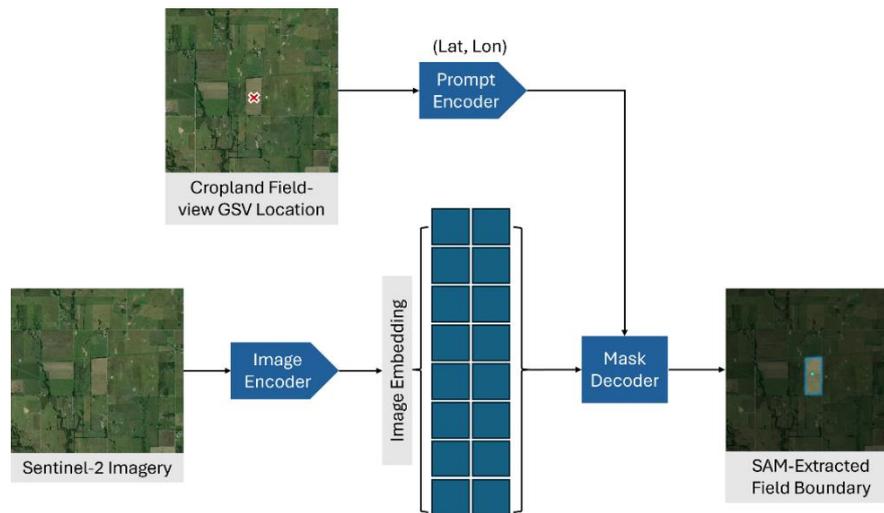
140 Pu, X., Jia, H., Zheng, L., Wang, F., and Xu, F.: Classwise-sam-adapter: Parameter efficient fine-tuning
141 adapts Segment Anything to SAR domain for semantic segmentation, IEEE J. Sel. Top. Appl. Earth Obs.
142 Remote Sens., <https://doi.org/10.1109/JSTARS.2025.3532690>, 2025.

143

144 **8. For Figure 6, there are some image patches after image embedding. Do they refer to image
145 embeddings? The authors are suggested to check this as image patches are not equal to image
146 embeddings.**

147

148 **Response:** Thank you for highlighting this issue. As you noted, image embeddings are high-level feature
149 representations derived from the model and are not visually equivalent to raw image patches. To address
150 this, we have revised Figure 6 to reflect the proper representation of image embeddings and updated the
151 corresponding description “image embedding” [386] to avoid confusion.



152

153 **Figure 1: Structure of the SAM with cropland field-view GSV location coordinate as the point prompt for
154 cropland field boundary delineation from Sentinel-2 imagery (© European Union/ESA/Copernicus,
155 processed via Google Earth Engine).**

156 **9. For line 353, there is an error regarding reference “(Error! Reference source not found).”**

157

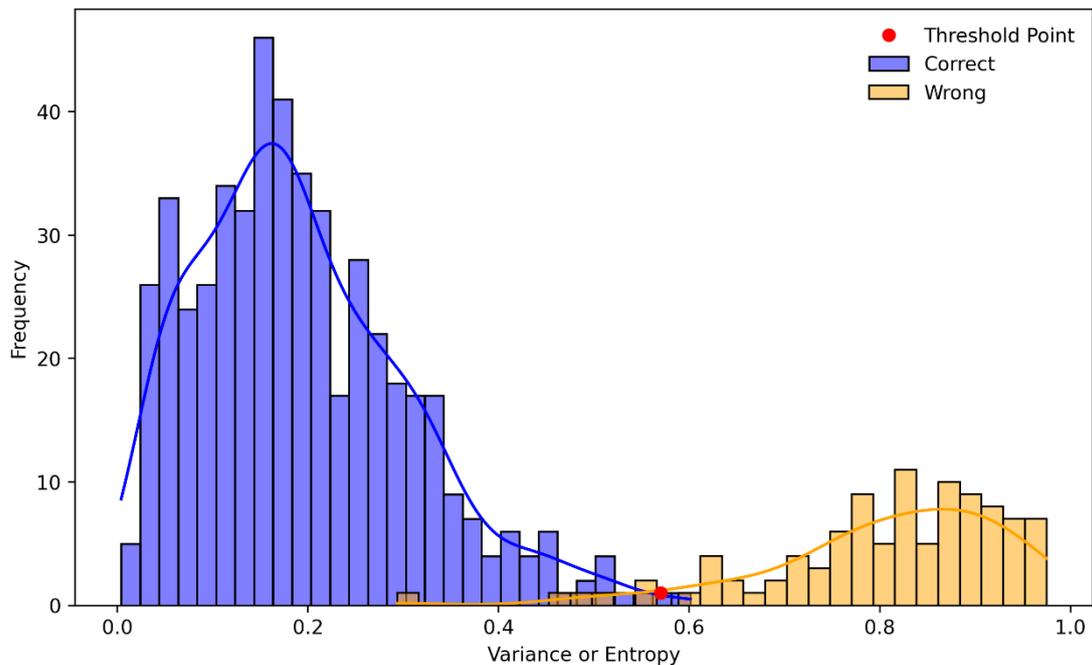
158 **Response:** Thank you for pointing out this issue. We have corrected the broken reference (now at [Line
159 365]) and updated it to properly point to [Fig. 3], where the corresponding crop types we have identified
160 through street view imagery across CONUS is shown. The manuscript has been checked accordingly to
161 ensure all references are correctly linked.

162

163 **10. How to choose reasonable thresholds for entropy and variance?**

164
165 **Response:** Thank you for this helpful question. In our method, we select the uncertainty (i.e., entropy and
166 variance) thresholds to balance classification confidence with sufficient sample retention for ground-truth
167 collection. We expanded our explanations of the process in the updated manuscript [Lines 370-377].

168
169 Specifically, to identify reasonable thresholds for achieving this trade-off, we analyze the empirical
170 distributions of entropy and variance for correctly versus incorrectly classified samples during model
171 training. By examining where the density curves of these two groups intersect (Fig. 5), we select threshold
172 values that best separate confident predictions from uncertain ones. This intersection-based criterion
173 (Abdar et al., 2021; Gour and Jain, 2022; Arco et al., 2023) provides a balance between retaining as many
174 field-view samples as possible and ensuring that the retained samples have sufficiently high confidence.
175 These thresholds provide a consistent confidence standard across all crop types, ensuring that a given
176 entropy or variance value represents the same level of model ambiguity regardless of the crop, which
177 facilitates uniform data filtering for downstream applications



179
180 **Figure 2: Distributions of uncertainty scores for correctly (blue) and incorrectly (orange) classified**
181 **samples produced by the CONUS-UncertainFusionNet model. The uncertainty threshold is determined at**
182 **the intersection point of the two density curves, representing the optimal separation between high- and**
183 **low-confidence predictions.**

184 **Reference**

185 Abdar, M., Fahami, M. A., Chakrabarti, S., Khosravi, A., Pławiak, P., Acharya, U. R., Tadeusiewicz, R.,
186 and Nahavandi, S.: BARF: A new direct and cross-based binary residual feature fusion with uncertainty-
187 aware module for medical image classification, *Inf. Sci.*, 577, 353–378, doi:10.1016/j.ins.2021.07.047,
188 2021.

189 Arco, J. E., Ortiz, A., Ramírez, J., Martínez-Murcia, F. J., Zhang, Y.-D., and Górriz, J. M.: Uncertainty-

190 driven ensembles of multi-scale deep architectures for image classification, *Inf. Fusion*, 89, 53–65,
191 <https://doi.org/10.1016/j.inffus.2022.08.010>, 2023.
192 Gour, M., and Jain, S.: Uncertainty-aware Convolutional Neural Network for COVID-19 X-ray Images
193 Classification, *Comput. Biol. Med.*, 140, Article 105047,
194 <https://doi.org/10.1016/j.compbimed.2021.105047>, 2022.

195

196 **Reviewer 2's Comments**

197 **This manuscript introduces CropSight-US, a national-scale, object-based crop type ground truth**
198 **dataset for the contiguous United States (2013–2023), derived from Google Street View imagery and**
199 **Sentinel-2–based field boundary delineation. The dataset is novel in its integration of street-level**
200 **imagery at national scale and in its object-level design with uncertainty metrics. However,**
201 **several fundamental aspects of dataset representativeness, temporal consistency, and uncertainty**
202 **interpretation remain insufficiently documented, limiting confidence in its general applicability.**

203

204 **Response:** We thank Reviewer 2 for the careful evaluation and for highlighting both the novelty of
205 CropSight-US and the areas requiring clarification. We have carefully revised the manuscript to provide
206 additional documentation on dataset representativeness, temporal consistency, and the interpretation of
207 uncertainty metrics.

208

209 Specific Comments:

210

211 **1. All samples are derived from road-accessible GSV imagery, yet the manuscript provides no**
212 **quantitative assessment of resulting spatial bias. Without statistics on distance-to-road coverage or**
213 **cropland representativeness at the ASD level, it is difficult to evaluate how well CropSight-US captures**
214 **agricultural landscapes away from road networks.**

215

216 **Response:** We appreciate the reviewer's comment regarding potential spatial bias associated with the use of
217 road-accessible GSV imagery. To address this, we conducted an ASD-level analysis of road-field-proximity.
218 Using the 2023 CSB dataset as a proxy for the total agricultural landscape, we sampled 5% of all fields in
219 each ASD and calculated the distance from field boundaries to the nearest road network. To ensure these
220 distances reflect functional management units rather than classification artifacts, we dissolved adjacent CSB
221 polygons sharing the same 2023 crop type label prior to sampling. Across the sampled fields in ASDs, the
222 average distance from field boundaries to roads was 112.32 ± 74.66 m at the ASD level. As for our
223 CropSight-US dataset, the field-to-road distances range from under 10 m up to 150 m. Therefore, CropSight-
224 US could mostly capture conditions representative of typical agricultural field configurations in CONUS.

225

226 Additionally, much of the U.S. cropland landscape, particularly in the Midwest, is organized under the Public
227 Land Survey System (PLSS), which establishes a regular one-mile grid of section lines along which rural
228 roads are commonly constructed (Abreu et al., 2017). This grid-based infrastructure results in a high
229 proportion of agricultural management units being road-adjacent, thereby enhancing the spatial coverage
230 and practical representativeness of roadside imagery relative to regions with more irregular road networks.
231 However, we admit that reliance on road-accessible GSV imagery may constrain the spatial
232 representativeness of the GSV-derived ground-truth dataset, particularly in heterogenous agricultural regions

233 characterized by sparse road networks. We have revised the manuscript to explicitly address these
234 considerations and clarify their implications for spatial representativeness [Lines 486-493, Lines 625-628,
235 Lines 749-756].

236

237 Reference:

238 Abreu, A., Lawson, L. A., Hyman, M., Hardin, R., and Gerling, M.: Collecting data from a permanent grid
239 sampling frame via a mobile mapping instrument, NASS Res. Rep., 322824, United States Department of
240 Agriculture, National Agricultural Statistics Service, 2017. <https://doi.org/10.22004/ag.econ.322824>

241

242 **2. The selection of Sentinel-2 (or NAIP) imagery “closest in time” to GSV acquisition is not quantified.**
243 **The manuscript should report typical and maximum temporal offsets and discuss potential impacts**
244 **on both crop labeling and boundary delineation.**

245

246 **Response:** Thank you for the suggestion to clarify the temporal alignment of GSV and the associated satellite
247 imagery for cropland field boundary delineation. To ensure temporal alignment between GSV observations
248 for crop type labeling and satellite imagery for corresponding field-boundary delineation, we relied on the
249 Month–Year metadata provided by GSV, which specifies the capture month of each panorama (no finer
250 temporal scale information can be retrieved from GSV yet). For every labeled GSV image, we searched for
251 the highest-quality, least-cloudy Sentinel-2 acquisition within a ± 1 -month window around the GSV capture
252 month. This yields a typical temporal offset of 0–30 days, with a maximum allowable offset of 60 days when
253 suitable scenes appear only at the edges of the search window. When Sentinel-2 scenes within this window
254 were unavailable or had excessive cloud cover or seasonal mismatch (e.g., pre-greenup or post-harvest
255 phenology), we used NAIP aerial imagery as a fallback, selecting the temporally closest NAIP acquisition
256 to maintain temporal consistency. This ensures that SAM-based delineation accurately captures the
257 productive field boundaries corresponding to the crop type visible from the GSV imagery. We have updated
258 the descriptions in Section 3.1.2 to reflect this clarification in methods [Line 276-278].

259

260 **3. The manuscript places strong emphasis on uncertainty as a key advantage of the dataset. However,**
261 **it remains unclear how users are expected to interpret and use the provided uncertainty metrics in**
262 **practice. In particular, it is not specified whether the reported entropy and variance values are directly**
263 **comparable across crop types, or whether they should be interpreted as crop-specific relative measures.**
264 **Without such clarification, the practical value of the uncertainty information for downstream**
265 **applications is difficult to assess.**

266

267 **Response:** Thank you for this important request for clarification. To ensure the dataset is practical for
268 downstream applications, the reported entropy and variance values are directly comparable absolute
269 measures across all crop types, rather than crop-specific relative measures. As shown in Section 3.1.3 and
270 Figure 5, we achieved this by determining the high-confidence using the aggregated density distributions of
271 the entire validation split, rather than calculating independent thresholds for each crop. We have expanded
272 the explanation of uncertainty metrics in Section 3.1.3 of the manuscript [Line 371-377], and in our GEE
273 interface as well.

274

275 **4. The authors group wheat, rice, and other small-grain cereals into a single “cereal” class based on**
276 **their visual similarity in street-level imagery. While this decision is understandable from an**

277 engineering and labeling perspective, it introduces substantial risks at the application level that are
278 not sufficiently discussed. In particular, rice paddies are typically characterized by surface flooding
279 and distinct water management practices, whereas wheat is associated with fundamentally different
280 hydrological and energy conditions. This raises the concern that water-related signals may be
281 implicitly mixed into the “cereal” class, potentially degrading the reliability of rice-related samples.
282 The manuscript should more explicitly discuss these implications and the limitations they impose on
283 crop-specific analyses.

284
285 **Response:** We appreciate your insightful observation regarding the distinct hydrological and energy
286 conditions associated with rice paddies compared to wheat and other small-grain cereals. We fully agree that
287 grouping these crops into a single Cereal class introduces potential hydrological mixing, particularly the risk
288 that water-related signals characteristic of flooded rice fields may be implicitly incorporated into samples
289 labeled as cereals. We have expanded the Discussion to explicitly address this limitation.

290
291 In the revised manuscript, we now clarify that our current grouping is driven by the visual similarity of
292 canopy architecture in roadside GSV imagery, where wheat, rice, rye, and triticale often exhibit comparable
293 field-level appearances. Because the labeling relies on static, ground-level images, surface flooding and
294 water-management structures common in rice systems are not always reliably visible from the road, which
295 further complicates species-level separation. To mitigate potential hydrological mixing for downstream users,
296 we note in the revision that researchers focusing on rice-specific applications should apply secondary
297 water-related filters, such as Normalized Difference Water Index (NDWI) or other moisture indices derived
298 from satellite imagery, when attempting to disaggregate the Cereal class into more specific crop categories.
299 This allows users to screen for likely flooded or wet environments that correspond more closely to paddy
300 rice conditions.

301
302 We have updated the manuscript in Section 6: Discussion and Conclusion to address this comment [Lines
303 742-749].

304
305 Technical Issues:

306
307 **Sect. 3.1.3 contains an unresolved reference (“Error! Reference source not found.”).**

308
309 **Response:** Thank you for pointing out this issue. We have corrected the broken reference at Line 353 (now
310 at [Line 365]) and updated it to properly point to Fig.3, where the corresponding crop types we have
311 identified through street view imagery across CONUS is shown. The manuscript has been revised
312 accordingly to ensure all references are correctly linked.

313
314