

Cross-view Object Geolocalization for Tree Species Mapping

Anonymous ICCV submission

Paper ID *****

Abstract

001 *Accurate, large-scale tree-species identification is a miss-*
 002 *ing link for many high-impact forestry and climate appli-*
 003 *cations. Though species identification is feasible using*
 004 *ground-level imagery, globally scaleable tree species map-*
 005 *ping requires identifying species using remote sensing data.*
 006 *We bridge this gap by formulating and solving the cross-*
 007 *view object geolocalization task: matching every object de-*
 008 *TECTED in a ground-level panorama to its counterpart in a*
 009 *precisely geo-referenced satellite image. Our end-to-end*
 010 *pipeline matches each individual tree in a satellite image to*
 011 *a corresponding streetview image of that tree, while filter-*
 012 *ing out occluded trees using ray-tracing. Evaluated on 7.8k*
 013 *trees, the method achieves 83.6 percent matching accuracy,*
 014 *rising to 84.3 percent for trees close to the camera. By link-*
 015 *ing ground-level species labels to satellite-detected crowns,*
 016 *our approach unlocks scalable, multimodal reference data*
 017 *for global tree-species mapping and sets a new direction for*
 018 *fine-grained cross-view localization tasks beyond forestry.*

019 1. Introduction

020 Forestry is one of the most important natural climate solu-
 021 tions. Forests store over 200 gigatons of carbon [1], and
 022 adding trees to cropland has the potential to sequester 5 gi-
 023 gatons of carbon per year [2]. Trees can also provide other
 024 key benefits, including provision of additional food and re-
 025 venue streams, improved soil health, and increased habitat
 026 for biodiversity. But research on forestry contains a signifi-
 027 cant constraint: lack of data.

028 Recent advances have shown that it is possible to seg-
 029 ment individual tree crowns at continental scales with high-
 030 resolution satellite imagery [3] and tree count and tree cover
 031 using freely available satellite data with global coverage [4].
 032 This represents a transformative opportunity to increase the
 033 data available on trees around the world. However, with-
 034 out data on the species of individual trees, the most useful
 035 forest monitoring applications are impossible. Tree species
 036 is required to accurately measure carbon sequestration from
 037 trees as different species sequester different amounts of car-

bon. Tree species is required for disease monitoring as dif- 038
 ferent species have different diseases. Trees provide signif- 039
 icant provisioning services such as timber, food, and other 040
 non-timber forest products, but yield monitoring requires 041
 tree species, as different species produce different types and 042
 quantities of different commodities. [5] Hanan et. al., 2020 043
 note that detecting tree “species will probably remain at 044
 the top of the Earth-observation research community’s wish 045
 list . . .” (Draper et al., 2020). 046

Beery et. al. 2022 [6] demonstrated that accurate 047
 tree species classification can be achieved with google 048
 streetview imagery. However, streetview imagery is not 049
 available everywhere, and it is constrained to detecting 050
 species only for trees visible from the road. Ahlswede et. 051
 al. 2023 [7] demonstrated that tree species identification is 052
 also possible from multispectral time-series imagery from 053
 sentinel-2 using a large-scale georeferenced dataset in Eu- 054
 rope. However, this approach cannot be extended glob- 055
 ally without groundtruth georeference data for every geog- 056
 raphy or species. If one could link the species of individual 057
 trees identified using streetview imagery with the geoloca- 058
 tion of those trees, one could scaleably generate georefer- 059
 enced species labels for individual trees around the world. 060
 As a result, tree matching and geolocalization across ground 061
 and satellite views becomes critical. Each individual tree 062
 needs to be accurately geolocated for ground view based 063
 species labels to be used to supervise models for scalable 064
 tree species mapping from remote sensing data. As a result, 065
 scaleable tree species mapping requires object-level geolo- 066
 calization from ground and satellite views. 067

Our key contributions in this paper are to 1) Define 068
 the cross-view object geolocalization task, 2) Develop a 069
 pipeline that can be used for cross-view individual tree ge- 070
 localization and achieve strong performance in our human- 071
 annotated evaluations. 072

073 2. Related Work

Recent work has demonstrated that streetview imagery en- 074
 ables significantly more accurate tree species identification 075
 than satellite imagery [6]. Moreover, streetview imagery 076
 has also been used to scale up reference data for crop type 077

078 mapping, enabling significantly more accurate crop species
079 classification from remote sensing data than was possible
080 solely with field-based ground truthing [8]. This approach
081 of combining ground-level and satellite-based views to im-
082 prove crop type mapping can also work with geotagged
083 photos from mobile phones or motorbike helmet based cam-
084 eras [9].

085 There is a rich literature on image geolocalization as a set
086 of many distinct but related computer vision tasks, which
087 are summarized in this survey from Wilson et. al. 2024
088 [10], including: 1) single view geolocalization in which one
089 predicts the geolocation of a single image, 2) cross-view ge-
090 olocalization where one matches a ground-level image to a
091 georeferenced aerial or satellite image, and 3) object geo-
092 localization where one or multiple geolocated ground im-
093 ages are used to geolocate an object. Berton et. al. 2022
094 [11] release a benchmark for deep visual geolocalization
095 from a single image. Deuser et. al. 2023 [12] is the first
096 paper we find exploring the issue of cross-view geolocal-
097 ization, and they develop an architecture that addresses dis-
098 tortions from polar transformations to align views. Li et al.
099 2024 [13] explore unsupervised learning for cross-view ge-
100 olocalization. As far as we are aware, ours is the first paper
101 exploring cross-view object-level geo-localization, which
102 we apply to matching individual trees across views.

103 **3. Methods**

104 As noted by Wilson et. al 2024 [10], "the concept of geo-
105 localization broadly refers to the process of determining
106 an entity's geographical location, typically in the form of
107 Global Positioning System (GPS) coordinates. The entity
108 of interest may be an image, a sequence of images, a video,
109 a satellite image, or even objects visible within the image."
110 Our problem—matching each individual tree visible in a
111 ground-level photo to the exact tree crown detected from
112 space—falls is an example of cross-view retrieval which in-
113 cidentally enables object geo-localisation. We therefore in-
114 troduce cross-view object geo-localisation.

115 **3.1. Task Definition: Cross-view Object geo-**
116 **localisation**

117 Our task is as follows: Given a street-view panorama and
118 the list of objects detected in a co-located satellite chip,
119 identify the objects within the panorama that correspond
120 to each detected object in the satellite chip. Because the
121 satellite image is precisely georeferenced, solving this task
122 is equivalent to geo-locating each detected object in the
123 ground view. In our context of individual tree species identi-
124 fication, solving this problem enables one to match ground-
125 view images in which tree species are identifiable to exact
126 geo-locations, which can then be linked with satellite re-
127 mote sensing data to enable scaleable tree species identifi-
128 cation. This approach can also be extended to enable large-

scale supervision any remote sensing object classification 129
tasks. 130

3.2. Cross-view Geolocalization Data Pipeline 131

Stage	Description	
Road sampling	Sample candidate points every 25 m along OpenStreetMap® roads inside the study AOIs.	
Street-view acquisition	Query the GSV API to check whether imagery exists at these points, snap to the panorama's true capture point to correct the ±5 m API noise.	
Satellite Image Check	Check whether Planet Skysat API contains imagery at these points	
Tree-presence filter	Use the 3 m global canopy-cover layer of Brandt <i>et al.</i> (2024) to require > 50% canopy within a 15 m radius.	
Individual tree detection	Download Planet SkySat (0.5 m) imagery on this location; run a Faster-rcnn detector trained on the individual tree detection dataset from Sachdeva et. al. 2024 to obtain bounding boxes	132
Candidate pruning using Ray-tracing	We use the detected bounding boxes to apply a ray-tracing approach out from the location of the streetview image to filter out trees that are occluded by other trees	
Compute Streetview Camera Angle	For the un-occluded trees, we compute the location of the detected bounding box using the georeference information of the satellite image. By computing the angle of this point relative to the location of the google streetview image, we can point the streetview camera directly at the tree.	
Retrieve Streetview image centered on object	Retrieve a crop of the panorama centered on the object of interest using the generated angle.	

We use cross-view geo-localization to assemble a multi-modal dataset of trees combining street-view imagery, high-res satellite imagery, and other data products. We use existing data products from Brandt. et. al. 2024 [14] and OpenStreetmap along with the google street view API to identify lat/long points where trees are present near the road and streetview imagery is available. Then we 133
134
135
136
137
138
139

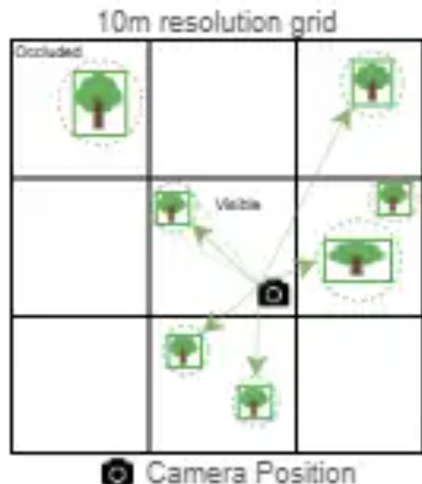


Figure 1. Our ray-tracing algorithm filters out trees that are occluded by other trees from the perspective of the streetview camera

140 check whether high-resolution satellite imagery from Planet
 141 Skysat is available in that area, and if it is, we apply a simple
 142 tree detector trained on a dataset from Sachdeva et. al.
 143 2024 [15] that outputs bounding boxes for each individual
 144 tree. Because the satellite image is georeferenced, we can
 145 convert the bounding boxes to lat/long coordinates for each
 146 detected tree.

147 We use ray-tracing to filter out trees that are occluded by
 148 other trees from the perspective of the camera, as shown in
 149 figure 1. The data pipeline is also made compatible with
 150 remote sensing modalities such as sentinel 2. We utilize the
 151 the remote sensing grid to isolate trees that are within the
 152 vicinity of the camera point allowing us to construct super-
 153 vision labels for each remote sensing pixel. This is partic-
 154 ularly useful for the species classification task by allowing
 155 the construction of a multi-model dataset.

156 By combining the lat/long coordinates of the streetview
 157 image with the lat/long coordinates of each individual
 158 tree, we can compute the angle of the tree relative to the
 159 streetview image. We use this angle to query a streetview
 160 image pointing at the detected tree for all the unoccluded
 161 trees in a satellite image. By doing this, we have linked each
 162 georeference remote sensing tree detection to a streetview
 163 image of the same tree. This pipeline is described in table
 164 1.

165 3.3. Annotations and Evaluation

166 To evaluate whether we have correctly matched detected
 167 trees from satellite images to images centered on that tree,
 168 we construct an annotated dataset. Our dataset consists of
 169 7380 trees, with each tree being comprised of a pair of satel-
 170 lite and street view images generated by our cross-view

geo-localization pipeline. The satellite image is centered
 171 on and a detected tree with a drawn bounding box, along
 172 with a point and line indicating the location and view of the
 173 streetview camera. If the match is correct, the street view
 174 image should be pointed at and centered on that same tree.
 175 We manually annotate 7,860 images as either "correct" or
 176 "incorrect" for whether the tree centered in the streetview
 177 image is indeed the tree in the satellite image. This task
 178 requires using the context of both images and the relative
 179 location of the tree to that context to interpret whether the
 180 object is correctly matched. For example, we have shown in
 181 figure 2 a correct match on the left, and an incorrect match
 182 on the right.
 183

Category	Count
OSM points sampled with SkySat Coverage	21,301
With Google Street View Image	17,275
With Valid Tree Cover and Farmland Cover	4,840
With non-occluded trees	3,683
Images Captured	7,380

Table 1. Summary of Data Points and Image Statistics

184 4. Results

185 We observe strong performance of our cross-view geolocal-
 186 ization pipeline in matching trees detected in the satellite
 187 view to the same tree detected in the ground view. In our
 188 dataset, 83.6 percent of trees retrieved using our pipeline
 189 have been correctly matched, as shown in table 1.

190 We report accuracy stratified by the distance of the tree
 191 from the road. As one might expect, the accuracy for trees
 192 near the road at 84.3 percent is slightly higher than for trees
 193 far from the road at 81.4 percent, but the relatively high
 194 accuracy under both conditions indicating that this approach
 195 is robust to the distance of the tree from the street.

Category	Correct	Incorrect	Accuracy	Total
All	6167	1213	0.836	7380
Near	4580	851	0.843	5431
Far	1587	362	0.814	1949

Table 2. Matching accuracy by distance.

196 We also report accuracy stratified by district in table 2
 197 to ensure that our pipeline is accurate under a wide range
 198 of conditions. We observe despite slightly lower accuracy
 199 in Ajmer, matching accuracy remains largely stable across
 200 districts, indicating our pipeline is relatively robust to geo-
 201 graphic variation as well.



(a) This is clearly a correct match, as the location of the tree relative to the road is the same in the both views

(b) This is clearly an Incorrect Match, as the detected tree in the satellite image is not centered in the streetview image

Figure 2. Example cross-view pairs that we annotate as correct matches of trees or not

District	Accuracy
Ajmer	0.600
Alwar	0.827
Bundi	0.917
Jaipur	0.900
Jhunjhunu	0.850
Nagaur	0.866
Pali	0.790
Sikar	0.792

Table 3. Accuracy values by district.

202 **5. Conclusion**

203 To enable remote sensing species identification, we de-
 204 velop an approach to matching individual trees detected
 205 in satellite imagery, where resolution is too coarse for vi-
 206 sual species identification, to trees in google streetview im-
 207 agery, where the images are sharp enough to distinguish tree
 208 species. We link google street view and planet skysat im-
 209 age modalities together, apply object detection to individ-
 210 ual trees from satellite imagery, and retrieve a streetview
 211 image pointed at each tree in the satellite image. We then
 212 conduct manual annotations to check whether we have cor-
 213 rectly matched ground images of the tree to the remotely
 214 detected tree.

215 We observe high accuracy of our pipeline in matching
 216 streetview imagery to trees in satellite imagery. The accu-

racy of this pipeline is robust to whether the tree is near or far from the road, and it is also robust to varying geographies with different environments.

Current approaches for collecting groundtruth species data through field work and mobile phone based geotagging are extremely resource intensive, preventing machine learning applications to tree species mapping to scaling across large geographies. We have ensured that trees are correctly matched from street view and satellite imagery. As a result, we can now scaleably generate geo-located reference species labels for training machine learning models to identify tree species from remote sensing data which is also geo-referenced. Cross-view object matching and geolocalization can dramatically reduce the cost and increase the scaleably of data collection for training remote sensing models to identify tree species. We will leverage this capability in forthcoming work.

Our approach to cross-view object mapping and geolocalization can be extended in many different ways. One could leverage ground view data from open-source databases like Mapillary, or to leverage ground image data taken from cameras attached to cars or helmets driving in rural areas. One could also increase the fidelity of the matching task by increasing the number and type of objects to match between satellite and ground view. A significant limitation of the approach we took is that our pipeline filters out most cases of dense trees because of occlusion. Future work should improve on our approach by enabling dense matching between every point in ground view imagery with

246 every point in satellite imagery rather than matching every
247 object, as these sorts of approaches will be necessary
248 for cross-view geolocalization in more complex forest land-
249 scapes.

250 Another note worth discussion is that this is clearly a
251 dual use technology. Though we were inspired to develop
252 this method for the application of tree species mapping, the
253 same approach described here could be used to match any
254 type of object between ground view and satellite views to
255 enable previously infeasible fine-grained classification tasks
256 from remote sensing data. This has significant ethical im-
257 plications as this technology could be used for mass surveil-
258 lance and military applications, for example by identifying
259 individual buildings or vehicles belonging to specific peo-
260 ple. As capabilities combining geospatial data, computer
261 vision, and artificial intelligence improve, there is a need to
262 strike the right balance between realizing the potential of
263 positive applications of these tools while increasing aware-
264 ness and mitigating risks around the potentially malicious
265 applications of the same technology.

266 References

267 [1] Mo L, Zohner CM, Reich PB, et al. Integrated global as-
268 sessment of the natural forest carbon potential. *Nature*.
269 2023;624(7990):92-101. 1

270 [2] Roe S, Streck C, Obersteiner M, et al. Contribution of
271 the land sector to a 1.5 C world. *Nature Climate Change*.
272 2019;9(11):817-28. 1

273 [3] Brandt M, Tucker CJ, Kariryaa A, et al. An unexpectedly
274 large count of trees in the West African Sahara and Sahel.
275 *Nature*. 2020;587(7832):78-82. 1

276 [4] Brandt J, Ertel J, Spore J, Stolle F. Wall-to-wall mapping
277 of tree extent in the tropics with Sentinel-1 and Sentinel-2.
278 *Remote Sensing of Environment*. 2023;292:113574. 1

279 [5] Pu R. Mapping Tree Species Using Advanced Remote Sens-
280 ing Technologies: A State-of-the-Art Review and Perspec-
281 tive. *Journal of Remote Sensing*. 2021;2021(4):1-26. 1

282 [6] Beery S, Wu G, Edwards T, Pavetic F, Majewski B, et al. The
283 Auto Arborist Dataset: A Large-Scale Benchmark for Mul-
284 tiview Urban Forest Monitoring Under Domain Shift. In:
285 *Proceedings of the IEEE/CVF Conference on Computer Vi-
286 sion and Pattern Recognition (CVPR)*; 2022. p. 21294-307.
287 1

288 [7] Ahlswede S, Schulz C, Gava C, Helber P, Bischke B, Förster
289 M, et al. TreeSatAI Benchmark Archive: A Multi-Sensor,
290 Multi-Label Dataset for Tree Species Classification in Re-
291 mote Sensing. *Earth System Science Data*. 2023;15(2):681-
292 95. 1

293 [8] Soler JL, Friedel T, Wang S. Combining Deep Learning and
294 Street View Imagery to Map Smallholder Crop Types. In:
295 *Proceedings of the AAAI Conference on Artificial Intelli-
296 gence*; 2024. Special Track on AI for Social Impact. 2

297 [9] Nakalembe C, Zvonkov I, Kerner H, et al.. Helmets Label-
298 ing Crops: Kenya Crop Type Dataset Created via Helmet-
299 Mounted Cameras and Deep Learning; 2025. *EarthArXiv*
300 preprint. 2

[10] Wilson D, Zhang X, Sultani W, Wshah S. Image and Object
Geo-Localization. *International Journal of Computer Vision*.
2024;132:1350-92. 2 301
302
303

[11] Berton G, Mereu R, Trivigno G, Masone C, Csurka G, Sattler
T, et al. Deep Visual Geo-Localization Benchmark. In: *Proce-
edings of the IEEE/CVF Conference on Computer Vision
and Pattern Recognition (CVPR)*; 2022. p. 5396-407. 2 304
305
306
307

[12] Deuser F, Habel K, Oswald N. Sample4Geo: Hard Nega-
tive Sampling For Cross-View Geo-Localisation. In: *ICCV*;
2023. p. 16847-56. 2 308
309
310

[13] Li G, Qian M, Xia G. Unleashing Unlabeled Data: A
Paradigm for Cross-View Geo-Localization. In: *Proceeed-
ings of the IEEE/CVF Conference on Computer Vision and
Pattern Recognition (CVPR)*; 2024. p. 16719-29. 2 311
312
313
314

[14] Brandt M, Gominski D, Reiner F, et al. Severe decline in
large farmland trees in India over the past decade. *Nature
Sustainability*. 2024;7:860-8. 2 315
316
317

[15] Sachdeva S, Lopez I, Biradar C, Lobell D. A Distribution
Shift Benchmark for Smallholder Agroforestry: Do Founda-
tion Models Improve Geographic Generalization? In:
*Proceedings of the Machine Learning for Remote Sensing
(ML4RS) Workshop at ICLR*; 2024. p. 1-8. 3 318
319
320
321
322