

1 Boundary information without parcel objects: probing the
2 AlphaEarth annual embedding against seasonal Sentinel-1/2
3 features for sub-hectare coffee delineation

4 Antônio Campos de Abreu Filho^{1,*}, Marcos Antônio Timbó Elmiro¹, Marcelo Antonio
5 Nero¹, and Rodrigo Machado Fernandes Leitão²

6 ¹Universidade Federal de Minas Gerais (UFMG), Belo Horizonte, Brazil

7 ²Independent Researcher

8 *Address correspondence to: antonio@abreufilho.com.br

9 **Abstract**

10 Embedding products promise to replace local seasonal compositing in agricultural mapping:
11 AlphaEarth Foundations condenses a year of multi-sensor observation into a 64-dimensional
12 annual embedding at 10 m. Whether such a representation retains the few-pixel-wide boundary
13 information that parcel delineation requires is untested. We test this in a demanding case: 384,496
14 coffee parcels in Minas Gerais, Brazil (2017 season; median 0.66 ha), using a controlled probing
15 design. Identical lightweight probes (Random Forest plus a linear probe) are trained on four
16 per-pixel representations: seasonal Sentinel-2, seasonal Sentinel-1 plus Sentinel-2, AlphaEarth,
17 and their fusion. Parcel objects are extracted by a fixed watershed whose single threshold is
18 selected on spatially held-out validation cells. On held-out test cells, the embedding carried
19 more per-pixel boundary information than the seasonal features (boundary average precision
20 0.524 versus 0.406; median paired gain +0.120, $p_{\text{Holm}} < 0.0001$, Holm-Bonferroni-adjusted
21 within contrast). The signal was also linearly decodable. The advantage did not convert into
22 better parcel objects: the embedding traded over-segmentation for under-segmentation, and
23 seasonal features kept equal or better object metrics, especially for small parcels; the n=3
24 large-parcel stratum showed a possible object-quality gain but remains exploratory. Under this
25 fixed watershed extraction, the result is a conversion gap: boundary evidence is present, but this
26 extraction does not turn it into reliable objects. The findings caution against treating annual
27 embeddings as drop-in inputs for sub-hectare parcel delineation.

28 **Keywords**

29 field boundary delineation; foundation model embeddings; AlphaEarth; Sentinel-1; Sentinel-2; repre-
30 sentation probing; sub-hectare parcels; watershed segmentation

ORCID

ORCID: Abreu Filho 0009-0007-1615-2315; Timbó Elmiro 0000-0001-7680-3131; Nero 0000-0003-2124-5018; Leitão 0000-0002-8470-5619

1 Introduction

Field boundaries are a core enabler of agricultural monitoring: parcel-level maps support crop-type classification, yield estimation, and the operational support of food-security programmes [1, 2]. Most of the world’s farms are small: more than 500 million operate on less than two hectares [3]. In small-parcel landscapes, automated delineation is both most useful and most difficult. Deep-learning boundary extractors have raised the state of the art on medium-resolution imagery [4], and dedicated architectures and datasets have pushed into small-parcel settings [5, 6], but these models are trained per task and per region, and their inputs are typically carefully engineered composites or time series of Sentinel imagery. In parallel, remote-sensing foundation models have begun to adapt self-supervised pre-training to the spectral, temporal, and geographic structure of Earth observation data: SatMAE adapts masked autoencoding to temporal and multispectral imagery [7], Prithvi scales geospatial pre-training on Harmonized Landsat–Sentinel data [8, 9], and Presto targets pixel-level remote-sensing time series [10]. Recent surveys frame these models as a rapidly growing but still task-dependent family of representations for Earth observation [11].

A newer alternative is the ready-made embedding. AlphaEarth Foundations distills a year of multi-source satellite observation, including optical, radar, and other streams, into a global 64-dimensional annual embedding at 10 m, designed so that light models trained on sparse labels can map directly from the embedding [12]. For a practitioner, the promise is concrete: skip local seasonal compositing and delineate parcels with a simple model over a product that already exists. The embedding has been evaluated chiefly on classification-style mapping tasks, where it performs strongly [12]; independent benchmarks confirm competitive accuracy on agricultural downstream tasks such as yield, tillage, and cover-crop mapping, while flagging limited spatial transferability [13]. No published evaluation targets parcel delineation. Whether the embedding preserves the fine spatial information that *delineation* needs is a different question, and an open one: parcel boundaries in small-parcel landscapes are a few pixels wide at 10 m, and annual aggregation could plausibly smooth away exactly the transient, season-dependent contrasts that separate one parcel from the next.

This study asks that question directly for a hard, well-instrumented case: coffee parcels in Minas Gerais, Brazil, where a public reference delineates 384,496 parcels with a sub-hectare median size for the 2017 season [14]. The analysis compares four per-pixel representations built from overlapping source families: seasonal optical features, seasonal optical plus radar features, the AlphaEarth annual embedding, and their fusion, under a supervised probing design. An identical lightweight probe is trained on each representation to predict boundary probability, parcel objects are extracted by a fixed watershed whose single threshold is selected on spatially held-out validation cells, and performance is measured on held-out test cells with object metrics, a threshold-free boundary signal, and per-tile paired statistics. Because the probe and every downstream step are held fixed, differences are

69 attributable to the representations. A secondary linear probe asks whether the embedding’s boundary
70 information is linearly decodable, following the linear-probe convention for learned representations
71 [15].

72 The comparison is one of representations, not sensors: the embedding is itself trained on optical
73 and radar streams, so the study measures the boundary information available in the ready-made
74 representation relative to season-resolved features engineered from overlapping source families. The
75 contributions are threefold: (i) a controlled probing design for representation comparison in parcel
76 delineation, with spatial splits and a single validated knob; (ii) evidence of an information–object
77 gap, where the annual embedding carries *more* per-pixel boundary information than the seasonal
78 features and that information is linearly decodable, yet the advantage does not convert into better
79 parcel objects under a simple extraction that trades over-segmentation for under-segmentation; and
80 (iii) a size-stratified analysis of that gap. The stratification suggests that the gap is weakest in the
81 large-parcel stratum, while the dominant small-parcel regime drives the object-level failure most
82 relevant to the delineation problem.

83 2 Study area and data

84 2.1 Study area and reference parcels

85 The study covers coffee-growing areas of Minas Gerais (MG), Brazil, for the 2017 crop year, a region
86 where coffee mapping from satellite imagery has a long operational history [16]. The reference is
87 CONAB’s semi-automatic mapping of Brazilian crop parcels, distributed on Source Cooperative
88 as the `br_conab` field-boundary dataset by the Field Boundaries for Agriculture (`fiboa`) project
89 [14]. Filtering the dataset to MG returned 384,496 parcels, all encoded as coffee (`MG_CAFE`) and
90 all from the 2017 season; the parcel identifier follows the pattern `{state}_{crop}_{yy}_{n}` (e.g.,
91 `MG_CAFE_17_1`), which serves as the filtering and stratification key. No other crop or year appears for
92 MG, so the reference fixes the scope of the study to coffee delineation in MG for the 2017 base year.

93 The parcels are small (Figure 1c). Because the dataset’s `metrics:area` field is null, area was
94 computed from the polygon geometry: the median parcel is 0.66 ha (interquartile range 0.13–1.98 ha),
95 with a maximum near 472 ha. The area distribution also shows a secondary mode of sliver features
96 below 0.01 ha: 76,203 features, or 19.8% of the count but only 0.011% of the mapped area, with a
97 median near 4 m². These features are consistent with fragmentation artefacts of the semi-automatic
98 delineation and are sub-pixel at the 10 m grid. At a 10 m ground sampling distance, the median
99 parcel covers roughly 65 pixels; a compact parcel of that size is only about eight pixels across, so
100 the target boundaries are just a few pixels wide. The parcels are treated as a reference with its
101 own delineation error rather than as cadastral ground truth: the mapping is semi-automatic, and
102 only coffee parcels are delineated. The reference is thus **sparse** because the area outside a mapped
103 parcel is unlabelled, not confirmed non-coffee; this constrains where the method can be trained and
104 evaluated (see Materials and Methods).

105 For spatial stratification and to guard against spatial leakage, each parcel was assigned, by the
106 location of its centroid, to a 0.5° region cell, which yielded 159 cells across MG.

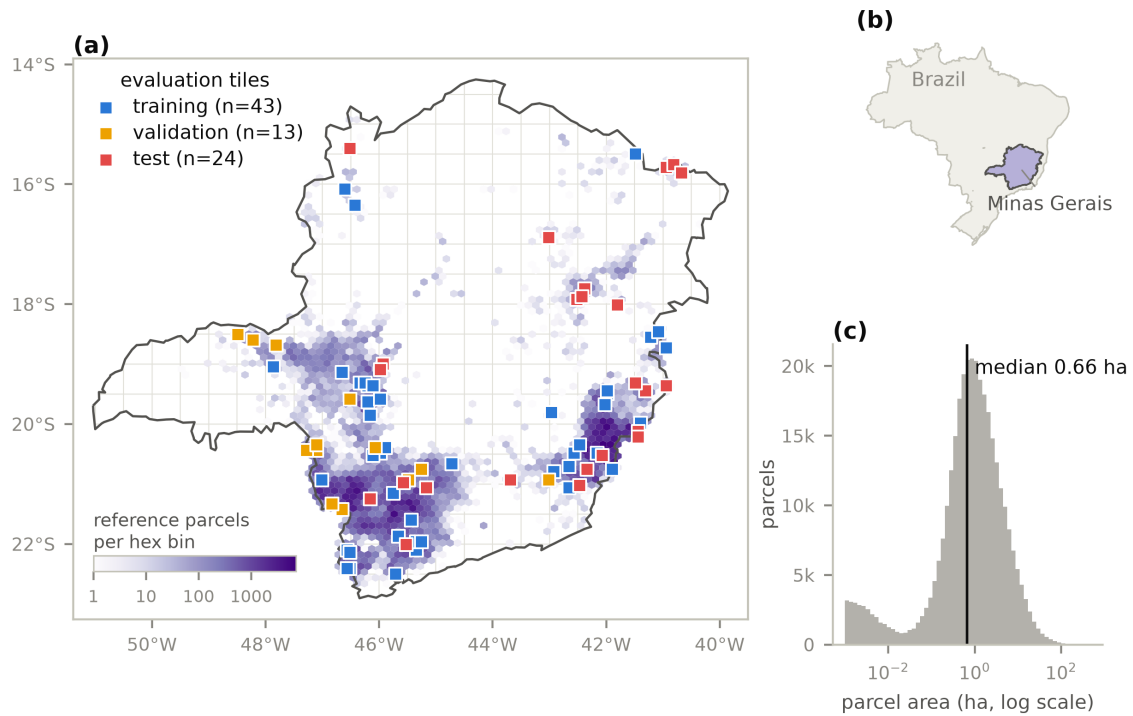


Figure 1. Study area and sampled tiles. (a) Reference coffee-parcel density in Minas Gerais and the 80 evaluation tiles by partition. The 0.5° cells define the spatial split. (b) Minas Gerais within Brazil. (c) Reference parcel-area distribution, with the median marked.

2.2 Data sources

Three feature sources cover MG for 2017, all at 10 m resolution and all accessed through Google Earth Engine [17] (Table 1). Earth Engine serves only as the data source; the analysis runs locally.

Sentinel-2 supplies the optical features [18]. The top-of-atmosphere (TOA) product (`S2_HARMONIZED`, Level-1C) was used because the surface-reflectance (Level-2A) collection in Earth Engine begins only in 2017 and covers the study area sparsely that year, leaving the Level-1C archive as the only source with full 2017 coverage. Sentinel-1 Ground Range Detected (GRD) imagery acquired in Interferometric Wide-swath (IW) mode supplies the synthetic-aperture-radar (SAR) features [19]; both Sentinel-1A and Sentinel-1B were operational throughout 2017, so the SAR record for the study year has no acquisition gap. AlphaEarth (the Satellite Embedding V1 product) supplies a 64-dimensional annual embedding [12]. The embedding is not an independent sensor: it is a learned representation trained over optical, radar, and other streams, so the comparison in this study is between representations of overlapping sources: a ready-made annual summary versus season-resolved features engineered from the same optical and radar inputs, rather than between independent sensors.

Table 1. Feature sources (Minas Gerais, 2017).

Source	Role	Resolution	Access
Sentinel-2 TOA (<code>S2_HARMONIZED</code> , L1C)	seasonal optical features	10 m	Earth Engine
Sentinel-1 GRD (IW, VV/VH)	seasonal SAR features	10 m	Earth Engine
AlphaEarth (Satellite Embedding V1)	annual 64-D embedding	10 m	Earth Engine
CONAB crop-parcel mapping (<code>br_conab</code>)	reference parcels (coffee, MG, 2017)	vector, EPSG:4674	Source Coop.

The optical and radar features were aggregated by season to follow coffee phenology in MG, whose arabica cycle spans a wet vegetative-to-fruiting phase and a drier maturation-to-rest phase [20]. A wet season (January–March plus October–December 2017) and a dry season (April–September 2017) were defined, and separate composites were computed for each. The wet-season composite thus joins two halves of adjacent phenological wet seasons within one calendar year; this is deliberate, because the AlphaEarth embedding likewise aggregates the 2017 calendar year, and aligning both representations to the same calendar window keeps the comparison fair. Aggregating within the single reference year (2017) also aligns every image with the reference and, with one year of data, replaces a temporal split with a spatial one during evaluation.

2.3 Evaluation tiles

Delineating all of MG is neither necessary nor feasible for this comparison, so the method is evaluated in tiles: square windows where the reference parcels give a local basis for evaluation. The MG bounding box was tiled into 500×500 -pixel cells (about 5 km on a side), and the 3,011 cells containing at least ten reference parcels were kept as candidates. From these candidates, a stratified sample of 80 tiles was drawn by region cell, tile size class, and parcel density class so that no single condition dominates the sample (Figure 1a). Size class is assigned per tile from the median area of its reference parcels: small is <1 ha, mid is 1–3 ha, and large is ≥ 3 ha. Density class is assigned from the parcel

140 count in the tile: low is <50 parcels, mid is 50–200, and high is ≥ 200 . The draw uses a fixed seed
141 (42) for reproducibility. Each tile was downloaded separately in its local UTM projection, which
142 keeps the pixel a true 10×10 m square (in geographic coordinates the pixel is not square at MG’s
143 latitudes). The tiles are the test sites, not the result: the study reports how the method performs
144 within them, stratified by region and parcel size.

145 3 Materials and Methods

146 The question is whether a ready-made annual embedding preserves the parcel-boundary information
147 that season-resolved Sentinel-1/2 features carry, or whether annual aggregation smooths it away. The
148 analysis uses a light supervised probing framework: the reference parcels become supervision rather
149 than only a final score, an identical probe is trained on each feature representation, and parcel objects
150 are extracted from the probe’s boundary probability. Because the probe and every downstream
151 step are held fixed across representations, differences in the resulting objects are attributable to
152 the features, not to the model. The pipeline runs in five stages (reference extraction, tiling, feature
153 download, probing, and validation), each implemented as a script in `code/`, with a final script deriving
154 the tables and figures. Figure 2 summarizes the full design, from the reference parcels through the
155 shared probe and extraction to the paired per-tile tests.

156 3.1 Reference data and labels

157 Within each tile, the reference parcels are rasterized into three classes. Only parcels lying entirely
158 inside the tile interior are used: the tile bounds are eroded by 10 m (one pixel) to keep parcels
159 contained in that eroded box, so that clipping a partial parcel never creates an artificial boundary
160 coincident with the tile edge. Each retained parcel contributes a **boundary** class, defined as its
161 outline buffered by 10 m and clipped to the parcel so the class occupies roughly one pixel along the
162 inner edge, and an **interior** class, the remainder of the parcel. Pixels outside any parcel are left
163 unlabelled: the reference is sparse, so the exterior cannot be asserted as background. Labels are
164 written per tile as a raster (0 = outside, 1 = interior, 2 = boundary), and training and boundary
165 evaluation are confined to the labelled footprint.

166 3.2 Feature representations

167 Four feature stacks define the experiments; each is a per-pixel feature vector fed to the same probe
168 (Table 2). E1 uses only the Sentinel-2 seasonal features; E2 adds the Sentinel-1 seasonal features; E3
169 uses only the AlphaEarth embedding; and E4 concatenates all three sources. E2 is the season-resolved
170 engineered baseline, E3 is the ready-made embedding, and E4 tests whether the embedding adds
171 anything beyond the engineered features.

172 Table 2. Feature stacks per experiment.

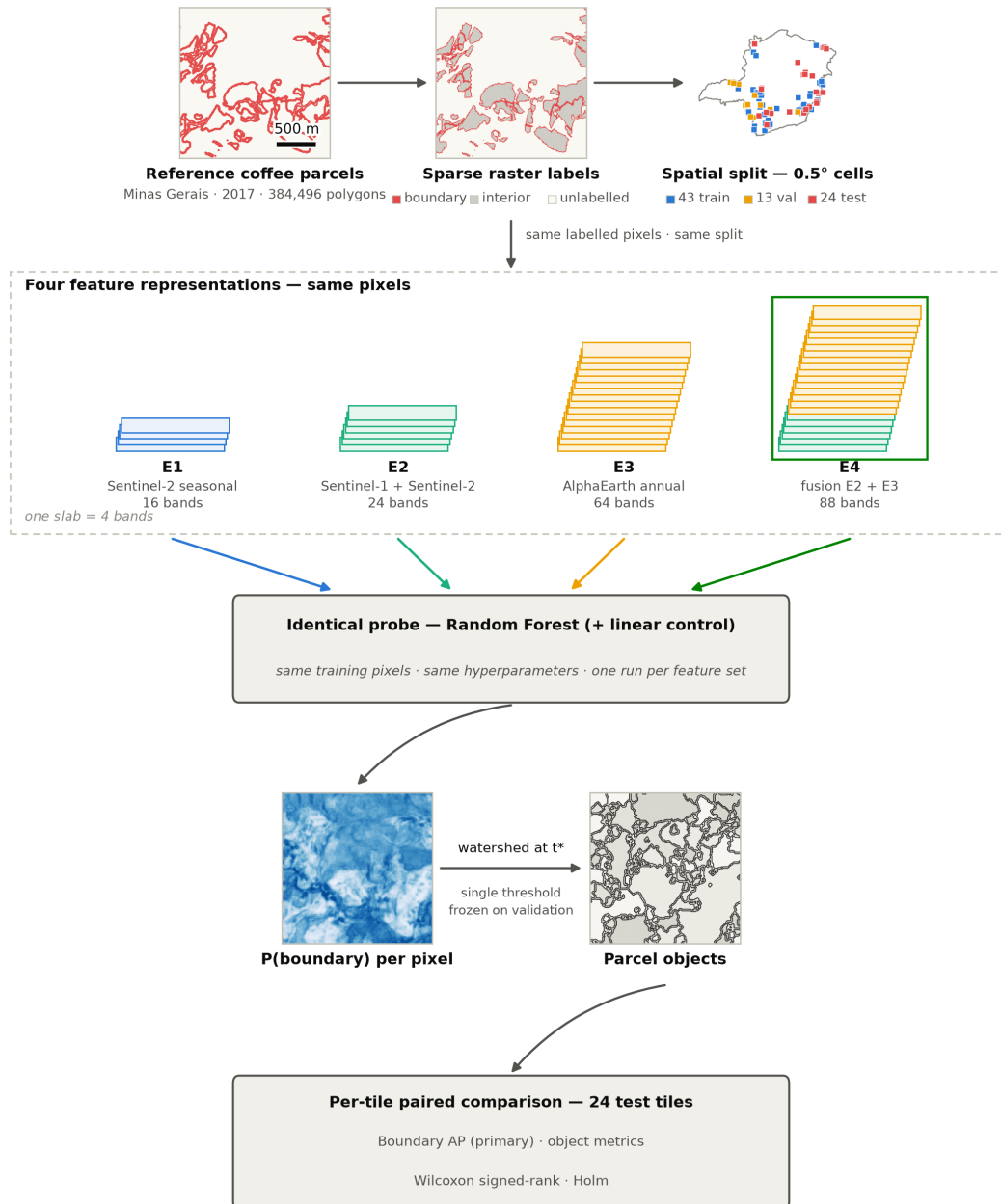


Figure 2. Experimental design. The same reference labels, probe, watershed extraction, and paired-test protocol are applied to each feature stack; only the input representation changes. Colour identifies the feature set, and grey stages are fixed across experiments.

Experiment	Sources	Bands
E1	Sentinel-2 seasonal	16
E2	Sentinel-1 + Sentinel-2 seasonal	24
E3	AlphaEarth embedding	64
E4	Sentinel-1 + Sentinel-2 + AlphaEarth	88

174 The Sentinel-2 stack holds eight features per season: the B2, B3, B4, B8, and B11 bands plus
 175 NDVI, NDWI [21], and EVI [22], for the wet and dry seasons, giving 16 bands. For each season,
 176 clouds were masked with the QA60 band, scenes below 60% cloud cover were kept, reflectance was
 177 scaled, and the collection was reduced to its per-pixel median before computing the indices. The
 178 Sentinel-1 stack holds four features per season: the median and standard deviation of VV and VH
 179 over IW GRD acquisitions, for both seasons, giving eight bands. Backscatter is used in decibels as
 180 distributed by Earth Engine, with no additional speckle filtering: the per-season median already
 181 suppresses speckle temporally. The AlphaEarth stack is the 64-band annual embedding mosaic for
 182 2017. Feature values are used as provided; because the probe is invariant to per-feature scaling (see
 183 below), no cross-source normalization is applied. Pixels left without a valid observation by a seasonal
 184 composite (persistent cloud cover in a season) are assigned a feature value of zero rather than masked,
 185 so every representation is trained and evaluated on the same pixel set.

186 3.3 Probing

187 A single Random Forest (RF) classifier [23] serves as the probe, with fixed hyperparameters shared
 188 across E1–E4: 200 trees, a minimum of five samples per leaf, and balanced-subsample class weights,
 189 under a fixed random seed (42). Holding the probe and its hyperparameters identical across
 190 experiments isolates the effect of the representation. A tree ensemble is invariant to monotone
 191 per-feature rescaling, so it compares differently-scaled representations without imposing a similarity
 192 metric or a normalization choice on any of them. The probe is trained on the same set of training
 193 pixels for every experiment and predicts, per pixel, the probability that the pixel belongs to a parcel
 194 boundary.

195 Training pixels are drawn from the labelled footprints of the training tiles only. To keep the classes
 196 and the tiles balanced, the sample is capped at 4,000 boundary and 8,000 interior pixels per tile; in
 197 the current run this yielded about 441,000 training pixels across the 43 training tiles. Each trained
 198 probe is then applied to every validation and test tile, producing a per-pixel boundary-probability
 199 raster for each experiment. Implementation uses scikit-learn [24].

200 A secondary linear probe asks whether the boundary information is *linearly* decodable from
 201 each representation. A logistic regression was trained on the same training pixels as the Random
 202 Forest, with fixed hyperparameters shared across E1–E4; because a linear model is not invariant to
 203 feature scaling, each feature was standardized with the mean and standard deviation computed on
 204 the training pixels only. The linear configurations are denoted with the suffix *-lin* (E1-lin through
 205 E4-lin). The linear probe’s boundary probabilities pass through the same object extraction and
 206 evaluation as the Random Forest’s, including a separate threshold selection on the validation tiles.

207 **3.4 Object extraction**

208 Parcel objects are recovered from the boundary-probability raster by watershed segmentation [25]
209 using scikit-image [26]. Rather than fixing an arbitrary number of segments, markers are derived from
210 the probability itself: the connected components of the region where boundary probability falls below
211 a threshold t seed the watershed, which then floods the probability surface to close object outlines.
212 The number of objects emerges from the predicted boundaries. Segments smaller than five pixels are
213 removed by sieving. Prediction and segmentation cover the whole tile, including unlabelled area:
214 because the reference is sparse, candidate objects outside the mapped footprint cannot be suppressed
215 without consulting the reference itself. Any resulting bias, such as a reference parcel matching a
216 candidate that extends into unlabelled area, applies identically to all four representations, so the
217 paired contrasts are unaffected, although absolute object metrics carry it. The threshold t is the
218 method’s single tunable knob. It is swept over the grid 0.10–0.60 in steps of 0.05 and selected on the
219 validation tiles by maximizing mean intersection-over-union: per-parcel IoU averaged within each tile,
220 then across tiles, with validation tiles where a threshold yields no watershed markers dropping out of
221 that threshold’s mean. Each experiment’s selected t is then frozen before any test tile is touched.

222 **3.5 Spatial split**

223 To prevent spatial leakage, the 0.5° region cells are partitioned whole: every tile in a cell goes to
224 the same partition. Cells are assigned to training (target 50% of tiles), validation (20%, where t is
225 chosen), and test (30%) under a fixed seed; in the current run this yielded 43 training, 13 validation,
226 and 24 test tiles. No pixel from a test cell is seen during training, and t is selected without touching
227 the test tiles. Because all imagery is from the single reference year (2017), the design uses this spatial
228 split in place of a temporal one.

229 **3.6 Evaluation metrics and statistical testing**

230 All metrics are computed on the test tiles only, per tile, over the parcels lying entirely within the tile
231 interior (partial parcels are excluded to avoid an artificial boundary at the tile edge). The candidate
232 object with the largest intersection is selected for each reference parcel, and intersection-over-union
233 is reported together with the over- and under-segmentation measures of Clinton et al. [27]. Object
234 recall is the fraction of reference parcels whose best candidate reaches an intersection-over-union of at
235 least 0.5; for parcels detected at that level, the median relative area error is also reported. Boundary
236 quality is measured by a Boundary F1 score at a 10 m (one-pixel) tolerance, computed in vector space
237 as the overlap length between buffered boundaries, with precision restricted to the buffered reference
238 footprint because the reference is sparse. To assess the boundary signal free of any threshold, the
239 pixel-wise average precision (AP) of the boundary probability is computed inside the footprint, which
240 scores the probe’s output directly and does not depend on t or on the watershed. Because paired tiles
241 share their conditions, representations are compared with a per-tile paired Wilcoxon signed-rank test
242 [28]. The two pre-specified primary tests are the threshold-free Boundary AP contrasts E3 versus E1
243 and E4 versus E2, which ask whether AlphaEarth carries additional boundary signal alone and in

244 fusion. Object metrics and the auxiliary contrasts (E4 versus E1 and E2 versus E1) are treated as
245 secondary diagnostics of how that signal converts into parcel objects. For transparency, the analysis
246 reports Holm-Bonferroni-adjusted p-values across the reported metrics within each paired contrast,
247 while interpreting secondary tests primarily by effect size and direction. All object and boundary
248 metrics are also reported as an overall mean and stratified by region and parcel size class, since a
249 sub-hectare median parcel makes a one-pixel boundary tolerance demanding and the difficulty varies
250 with parcel size.

251 3.7 Diagnostic figure protocols

252 Two qualitative figures (Figures 5 and 6) illustrate the extraction mechanism on a single test tile,
253 selected deterministically to avoid cherry-picking. The tile is drawn from the middle third of the E3–
254 E1 Boundary AP differences, so it is representative rather than a best case. For the map comparison
255 (Figure 5), a zoom window is chosen from the reference alone: the 640 m window containing the most
256 reference parcels, without consulting any prediction. In-panel counts report parcels and candidates
257 with at least 25 pixels inside the window. For the transect panel (Figure 6a), the plotted line is the
258 highest-scoring one under a declared criterion that favours many reference crossings and at least one
259 merge case; the two distributional panels (Figure 6b, c) instead use every pixel and object of the
260 tile. The object-sharing counts in Figure 6c count a reference parcel as *evaluatable* when its interior
261 has at least five pixels assigned to a positive candidate, after retaining only candidates with at
262 least half their area inside the reference footprint; because this filter and the watershed labels are
263 experiment-specific, the evaluatable denominators differ between E1 and E3 and are smaller than the
264 full count of reference parcels.

265 4 Results

266 All results below come from the 24 test tiles, which contain 1,956 reference parcels lying entirely
267 within tile interiors. The watershed threshold was selected per experiment on the validation tiles and
268 frozen before any test tile was evaluated (Table 3, column t^*). Paired differences are assessed with
269 the per-tile Wilcoxon signed-rank test described in the Materials and Methods; p-values reported for
270 these tests are Holm-Bonferroni adjusted across the reported metrics within each paired contrast.

271 4.1 Overall accuracy: boundary information without object gains

272 Table 3 reports the test-set summaries for the four representations under the Random Forest probe
273 and, in the lower block, the secondary linear probe. AlphaEarth produced a stronger threshold-free
274 boundary signal: E3 reached a Boundary AP of 0.524 against 0.406 for E1, a median paired gain of
275 +0.120 per test tile ($p_{\text{Holm}} < 0.0001$; Figure 3). The fusion contrast repeated this pattern, with E4
276 exceeding E2 by a median +0.107 in Boundary AP ($p_{\text{Holm}} = 0.0001$).

277 That additional boundary information did not become clearly better parcel objects under the
278 same watershed extraction. E2 had the highest mean IoU among the Random Forest probes (0.282),
279 followed by E1 (0.268), E3 (0.260), and E4 (0.258). The E3–E1 and E4–E2 mean-IoU differences were

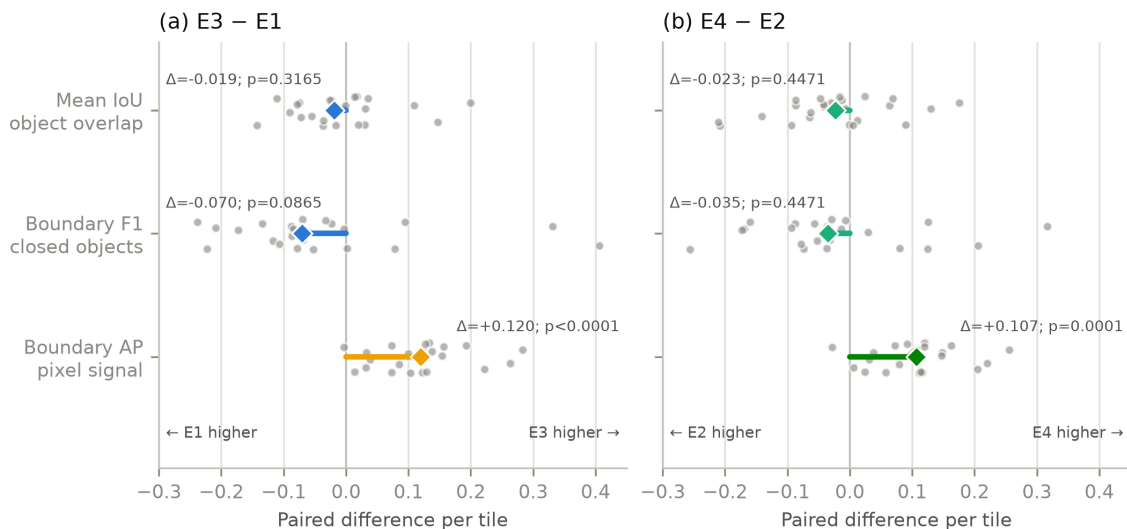


Figure 3. Paired tile-level effects for the two main contrasts: (a) E3 – E1, (b) E4 – E2. Points are tile differences, each tile at the same vertical offset in every row; diamonds and segments mark medians; annotations report median Δ and Holm-adjusted Wilcoxon p -values.

280 negative but not significant after within-contrast correction (median -0.019 and -0.023). Boundary
 281 F1 was lower for E3 than for E1 (0.534 versus 0.576; median -0.070; $p_{\text{Holm}} = 0.086$), consistent
 282 with weaker boundary closure, but this remains a secondary diagnostic rather than a confirmatory
 283 endpoint; the fusion contrast was weaker still (median -0.035). Among parcels detected at $\text{IoU} \geq 0.5$,
 284 the median relative area error was similar for E1–E3 under the Random Forest probe (0.232–0.254)
 285 but worst for the fusion (E4, 0.311), so the added bands did not sharpen the sizing of the parcels
 286 that were found.

Table 3. Test-set performance by experiment. t^* is the validation-selected watershed threshold; metrics are test-set means except Area error, reported as the median relative area error for parcels detected at $\text{IoU} \geq 0.5$. Bold marks the best value within each probe block.

Probe	Experiment	t^*	Object recall	Mean IoU	Over-seg.	Under-seg.	Area error	Boundary precision	Boundary recall	Boundary F1	Boundary AP
RF	E1 (S2 seasonal)	0.35	0.111	0.268	0.459	0.463	0.254	0.503	0.712	0.576	0.406
RF	E2 (S1+S2)	0.35	0.165	0.282	0.432	0.490	0.235	0.520	0.687	0.578	0.424
RF	E3 (AlphaEarth)	0.40	0.148	0.260	0.225	0.668	0.232	0.663	0.458	0.534	0.524
RF	E4 (fusion)	0.45	0.159	0.258	0.215	0.660	0.311	0.693	0.484	0.556	0.529
Linear	E1 (S2 seasonal)	0.40	0.103	0.233	0.342	0.619	0.302	0.523	0.543	0.517	0.412
Linear	E2 (S1+S2)	0.40	0.103	0.237	0.362	0.594	0.246	0.515	0.562	0.521	0.412
Linear	E3 (AlphaEarth)	0.35	0.154	0.254	0.204	0.677	0.273	0.687	0.433	0.523	0.530
Linear	E4 (fusion)	0.30	0.168	0.268	0.227	0.651	0.322	0.670	0.485	0.556	0.540

287 The component metrics show the mechanism behind this non-conversion. E3 and E4 sharply
 288 reduced over-segmentation relative to E1 and E2, but increased under-segmentation by nearly the
 289 same amount. For E3–E1, median over-segmentation decreased by -0.197 while under-segmentation
 290 increased by +0.200 (both $p_{\text{Holm}} \leq 0.0001$). Candidate counts tell the same story: the watershed
 291 produced $1,730 \pm 709$ E3 segments per tile and $2,397 \pm 761$ E4 segments, against $5,638 \pm 1,298$
 292 for E1 and $4,588 \pm 1,022$ for E2 (mean \pm SD across the 24 test tiles). The embedding fields thus
 293 support fewer, cleaner candidate regions, but they close fewer parcel boundaries between adjacent

294 parcels. The boundary-recall column of Table 3 reflects the same trade-off: E1’s high boundary recall
 295 (0.712, against 0.458 for E3) is partly a by-product of its over-segmentation, since a dense mesh of
 296 candidate boundaries touches most reference outlines. Absolute object accuracy remained low for
 297 every representation; object recall at $\text{IoU} \geq 0.5$ never exceeded 0.168.

298 4.2 Stratification by parcel size

299 The information–object gap varied with parcel size (Figure 4). The large-parcel stratum was small
 300 (three test tiles), so it is best read as exploratory; within that stratum, AlphaEarth and fusion had
 301 better boundary diagnostics, with RF E3 and E4 reaching Boundary F1 of 0.731 and 0.733, against
 302 0.509 for E1, and Boundary AP of 0.678 and 0.689, against 0.458 for E1. The mid-parcel class was
 303 mixed: E4 matched or slightly exceeded the seasonal baselines in Boundary F1, while E3 remained
 304 lower than E1/E2. The small-parcel class drove the main failure mode. E3 and E4 kept higher
 305 Boundary AP than E1/E2, but their Boundary F1 and mean IoU fell below the seasonal baselines,
 306 consistent with adjacent small parcels merging when probability ridges are too smooth.

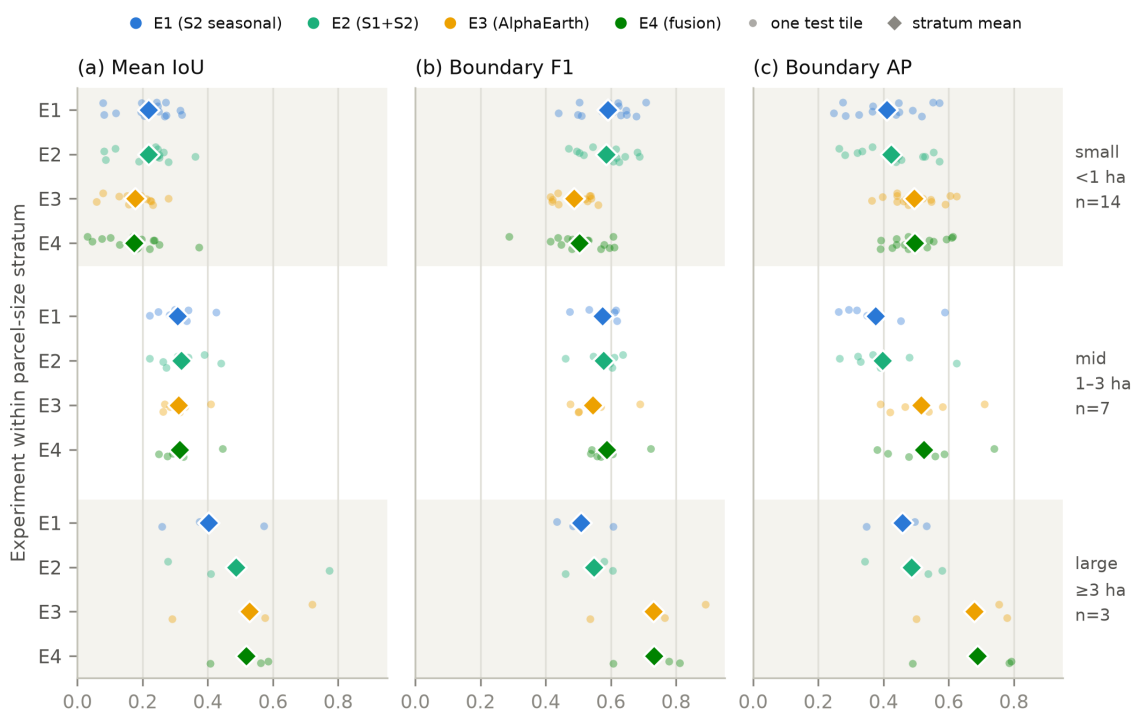


Figure 4. Size-stratified Random Forest performance on the test tiles. Rows show tile size classes (small $n=14$, mid $n=7$, large $n=3$), and columns show (a) mean IoU, (b) Boundary F1, and (c) Boundary AP. Small points are individual tiles; diamonds mark stratum means.

4.3 Linear decodability

The secondary linear probe confirms that the embedding’s boundary signal is not dependent on a non-linear decoder. A logistic regression on E3 reached Boundary AP 0.530, essentially the same as the Random Forest on E3 (0.524), and its per-tile advantage over E1-lin was +0.121 ($p_{\text{Holm}} < 0.0001$). The best linear configuration was E4-lin, with the highest object recall (0.168) and the highest Boundary AP (0.540). Even there, however, no mean-IoU difference was detected against the linear seasonal baseline (E4-lin minus E2-lin: median +0.009, $p_{\text{Holm}} = 0.449$). The embedding makes boundary information accessible, but this simple watershed does not reliably turn that information into more accurate objects.

4.4 Qualitative comparison

Figure 5 compares the predicted boundary-probability fields and watershed candidates of E1 and E3 on a test tile drawn from the middle third of the per-tile AP differences. On this representative tile, E3 has higher Boundary AP than E1 (0.547 versus 0.447) but lower Boundary F1 and mean IoU. The seasonal-feature field is granular and produces many more candidate regions; the embedding field is smoother and produces fewer, larger candidates. The zoom row makes the consequence concrete: in the window with the densest cluster of reference parcels, twelve parcels overlap 33 E1 candidates but only 18 E3 candidates.

Figure 6 quantifies the same mechanism without maps, on the same tile. Along a transect crossing 14 reference boundaries, the E1 profile exceeds its t^* at all 14 crossings, whereas the E3 profile stays below its own t^* at 4 of them. In that condition, the watershed merges adjacent parcels. The tile-wide distributions close both ends of the dissociation: the score ECDFs separate boundary from interior pixels more cleanly for E3, consistent with its higher Boundary AP, while the area-weighted size distribution shows that half of the E3 candidate area lies in objects up to 27.8 ha against 3.1 ha for the reference parcels (E1: 2.2 ha), and 58 of 93 evaluable reference parcels share their modal E3 candidate with a neighbour (E1: 36 of 122).

5 Discussion

5.1 What the information–object gap means

The central finding is a gap between boundary information and object delineation. Judged by threshold-free Boundary AP, the annual AlphaEarth embedding carries more per-pixel boundary information than the season-resolved features engineered from Sentinel-1/2: E3 exceeds E1 by a median +0.120 per tile, and E4 exceeds E2 by +0.107 (Figure 3; $p_{\text{Holm}} < 0.0001$ and $p_{\text{Holm}} = 0.0001$, respectively, after within-contrast correction). This advantage survives under a linear decoder, which means the signal is explicit in the representation, rather than recoverable only by a non-linear Random Forest.

The object-level result is more cautious. AlphaEarth does not clearly outperform the seasonal baselines after the shared watershed extraction, and the RF mean-IoU contrasts are small and

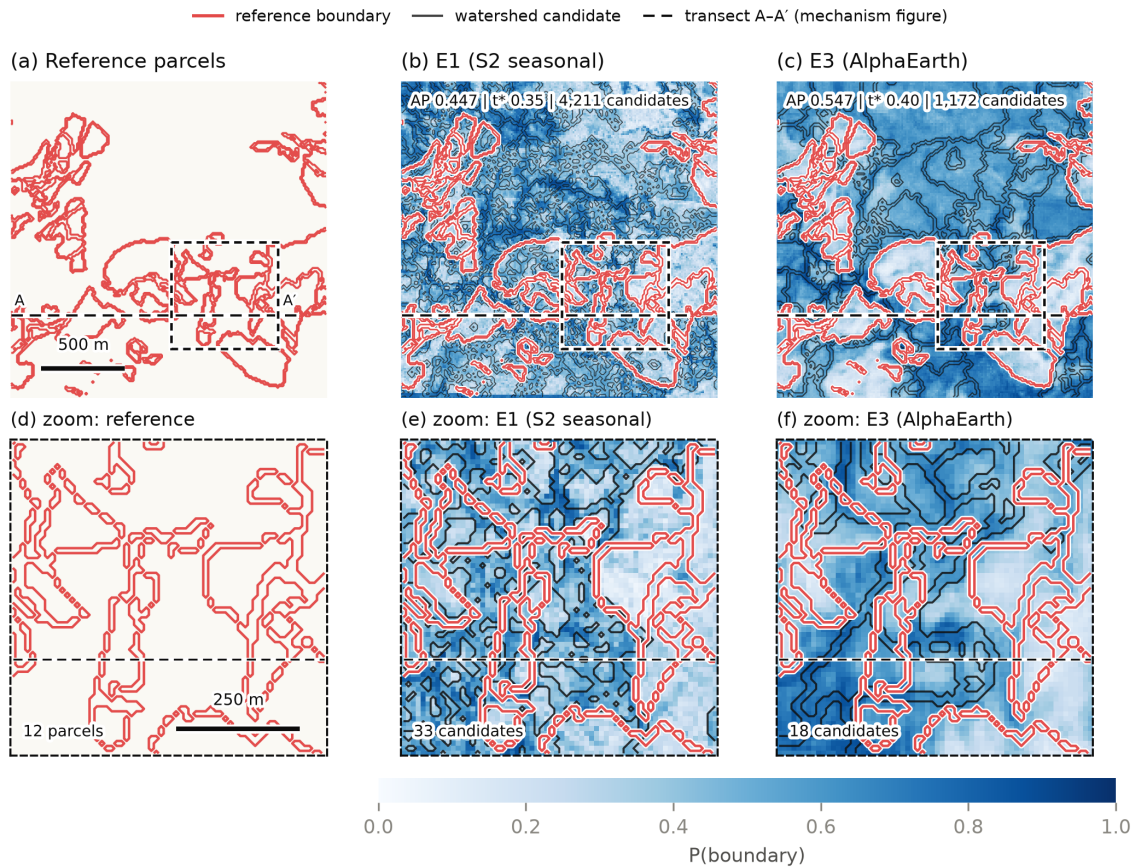


Figure 5. Boundary-probability fields and watershed candidates for the representative test tile: (a–c) communication crop, (d–f) reference-derived zoom window (dashed box). Red lines are reference parcel boundaries, black lines are candidate boundaries, and the horizontal dashed line marks transect A–A', analysed in Figure 6.

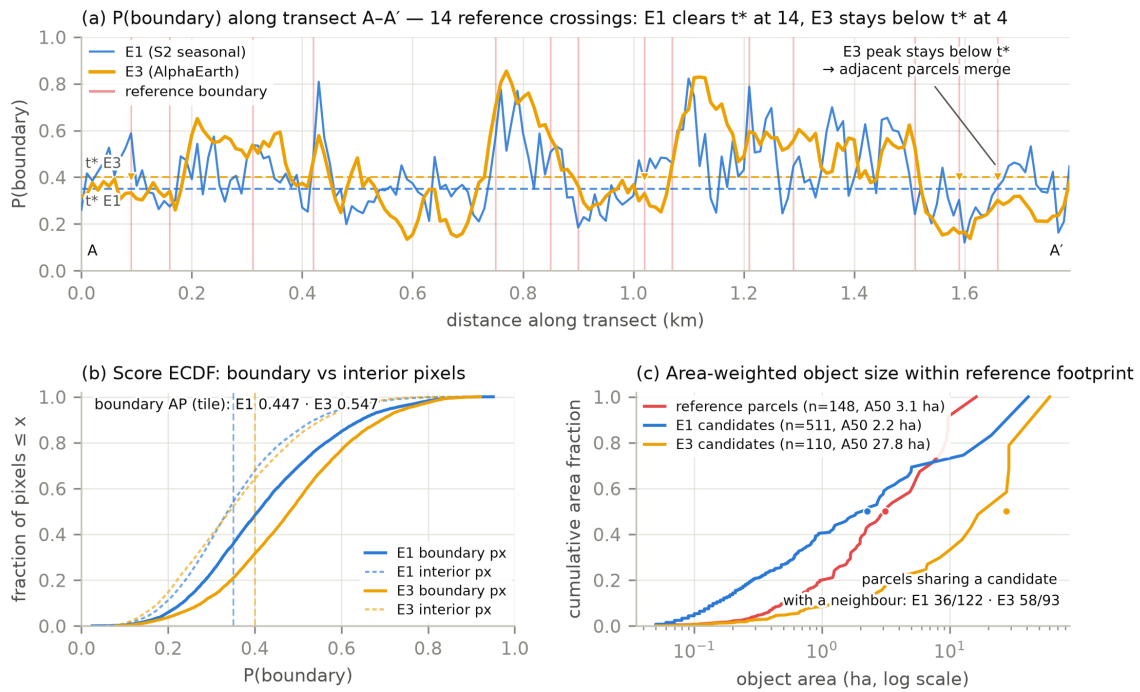


Figure 6. Non-map diagnostics for the representative test tile. Panels show (a) boundary-probability profiles along transect A-A' (marked in Figure 5), (b) ECDFs for boundary and interior pixels, and (c) area-weighted candidate-size distributions and object-sharing counts.

343 negative. A non-significant contrast is not evidence of equivalence; it only means that this test set did
344 not support a clear object-level gain. The largest object-level penalty is Boundary F1 for E3 against
345 E1, but this remains a secondary diagnostic rather than a confirmatory endpoint. The component
346 metrics explain why: the embedding reduces over-segmentation but increases under-segmentation.
347 In other words, the watershed produces fewer, cleaner regions from the smoother embedding field,
348 but adjacent coffee parcels are more likely to merge (Figures 5 and 6). The limitation exposed here
349 is thus less about missing boundary content than about the spatial structure required to convert
350 ranked boundary evidence into closed, correctly sized parcel objects.

351 The linear probe reinforces this interpretation. E3-lin reaches Boundary AP 0.530, close to the
352 Random Forest’s 0.524 on E3, and E4-lin gives the highest Boundary AP and object recall in the
353 table. Yet neither linear nor Random Forest probing turns the embedding’s AP advantage into a
354 statistically clear mean-IoU advantage. Absolute object accuracy is also low for every configuration,
355 with object recall at $\text{IoU} \geq 0.5$ topping out at 0.168. For this sub-hectare coffee setting, the result is
356 not an operational delineation method; it is a diagnosis of where the representation helps and where
357 the extraction fails.

358 5.2 Practical implications

359 For the practitioner question posed in the Introduction, whether one can skip local seasonal com-
360 positing and delineate parcels from the ready-made embedding alone, the answer splits by task.
361 For per-pixel boundary ranking within mapped coffee parcels, the embedding plus a light model is
362 the stronger representation in this study. For parcel objects from a simple watershed, it is not a
363 drop-in replacement for seasonal features: it improves Boundary AP, but the object metrics remain
364 comparable or worse, especially for small parcels.

365 The size stratification makes this practical boundary clear (Figure 4), but its strata are descriptive
366 rather than confirmatory. In the three large-parcel test tiles, the embedding’s information advantage
367 suggests better Boundary F1 and AP. On the small-parcel tiles that dominate the sample, E3 and E4
368 still rank boundary pixels better, but seasonal features keep stronger boundary closure and higher
369 mean IoU. These are the tiles where automated parcel delineation would be most useful, so the current
370 embedding-plus- watershed pipeline should be treated as diagnostic rather than production-ready.

371 The improvement opportunity follows directly. Since boundary information is present and linearly
372 accessible, this fixed watershed exposes a conversion problem rather than a lack of ranked boundary
373 evidence. Decoders with spatial structure, including contour completion, edge sharpening, instance
374 segmentation, or promptable segmenters [29, 30] conditioned on the embedding, are natural candidates
375 to test whether AlphaEarth’s Boundary AP advantage can be converted into objects. The current
376 study deliberately held the decoder simple and identical across representations so that the comparison
377 remained attributable to the features; it therefore does not show whether a stronger decoder would
378 recover the signal. The probing design itself carries beyond this case study: because the probe,
379 the extraction, and the statistics are fixed, the same protocol can benchmark any new embedding
380 product, or any other crop and region with a parcel reference, by swapping in a single feature stack.

381 The fusion experiment (E4) should be read in the same spirit. Concatenating the seasonal features

382 with the embedding is a deliberate control rather than an attempt at an optimal fusion: feeding
383 all sources to the same fixed probe keeps every difference attributable to the inputs, at the cost
384 of leaving the embedding’s information to compete with the seasonal bands inside a model that
385 cannot learn a joint representation. E4 accordingly tracks E3 on the boundary signal without a clean
386 object-level gain. Learned multimodal fusion, jointly training a representation over the raw streams
387 rather than concatenating fixed features [31], is a natural next step, but it would break the controlled
388 attribution this study depends on and is therefore left to future work.

389 5.3 Limitations

390 Six limitations bound these conclusions.

391 *Reference quality and sparsity.* The reference is a semi-automatic mapping with its own boundary
392 error, and it is sparse: only coffee parcels are delineated, so training, boundary precision, and
393 Boundary AP are all confined to the mapped footprint. Object candidates, in contrast, are extracted
394 over the whole tile, so absolute object metrics also reflect matches against unlabelled area. Because
395 candidate density differs across representations, this sparsity can affect object-level contrasts as
396 well as absolute scores. Metrics computed against such a reference score agreement with CONAB’s
397 delineation, not cadastral truth.

398 *Scope.* The scope is a single crop, state, and year: coffee in Minas Gerais, 2017, fixed by the
399 reference itself; whether the same information–object gap generalizes to other crops, landscapes, or
400 years is untested.

401 *Minimal extraction.* The object extraction is deliberately minimal: one watershed with a
402 single validated threshold. The results quantify what a simple, fixed extraction recovers from each
403 representation, not the ceiling of any representation under a stronger decoder.

404 *Absolute accuracy.* Absolute object accuracy is low for every configuration, reflecting how
405 demanding sub-hectare parcels are at a 10 m grid with a one-pixel boundary tolerance.

406 *Input asymmetries.* The seasonal features are computed from TOA (Level-1C) reflectance, whereas
407 AlphaEarth applies its own preprocessing to its source streams, an asymmetry that could penalize the
408 seasonal baseline. The fixed Random Forest also faces representations of 16 to 88 bands, and probe
409 capacity could interact with dimensionality; the linear probe, which reproduces the same conclusions,
410 serves as a partial control. Finally, zero-imputation of pixels left unobserved by a composite is applied
411 uniformly, but only the seasonal sources contain such pixels. The AlphaEarth rasters contain no
412 missing values, so imputed zeros occur only in the seasonal representations.

413 *Statistical scope.* Holm correction is applied within each paired contrast, with the Boundary
414 AP contrasts pre-specified as primary endpoints; this family definition is itself a choice, but the
415 primary result would survive even a global Holm correction across all roughly 48 reported tests
416 ($p \approx 0.0001 \times 48 \approx 0.005$). The 24 test tiles come from 17 spatial cells, so the paired tests support
417 tile-level diagnostics under the stratified sampling design. Non-significant object contrasts should not
418 be read as equivalence tests, and the size-stratified results are descriptive, especially the large-parcel
419 stratum ($n = 3$).

6 Conclusion

This study asked whether a ready-made annual embedding preserves the parcel-boundary information that season-resolved Sentinel-1/2 features carry, using coffee parcels in Minas Gerais (2017) as a demanding sub-hectare test case. A supervised probing design with identical lightweight probes, a fixed watershed extraction with one spatially validated threshold, and paired per-tile statistics on held-out cells kept every difference attributable to the representations themselves.

The answer is an information-object gap. AlphaEarth carried more per-pixel boundary information than the seasonal features: under the Random Forest probe, E3 reached Boundary AP 0.524 versus 0.406 for E1, with a median paired gain of +0.120 ($p_{\text{Holm}} < 0.0001$). The same signal was linearly accessible (E3-lin Boundary AP 0.530). But this advantage did not translate into a clear object-level advantage under the shared watershed. Seasonal features retained the best RF mean IoU (E2 = 0.282) and Boundary F1 (E2 = 0.578), while AlphaEarth reduced over-segmentation at the cost of higher under-segmentation. The embedding therefore appears to preserve boundary evidence, but its smoother spatial field is not enough, by itself, to close small adjacent parcel outlines.

For practice, the ready-made embedding is promising for boundary ranking within mapped coffee parcels, but not yet for operational parcel delineation with this extraction: no configuration reached object recall above 0.168 at $\text{IoU} \geq 0.5$. Future work should test whether spatially structured decoders can convert AlphaEarth’s boundary-information advantage into accurate parcel objects, and whether the same pattern holds beyond coffee in Minas Gerais.

Acknowledgments

Author Contributions

- **Abreu Filho:** Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Data curation, Writing – original draft, Visualization.
- **Timbó Elmiro:** Conceptualization, Methodology, Writing – review & editing, Supervision.
- **Nero:** Conceptualization, Methodology, Writing – review & editing, Supervision.
- **Leitão:** Methodology, Writing – review & editing.

Funding

This research received no external funding.

Conflicts of Interest

All authors have approved the submission and declare no competing interests.

Generative AI and Figure Provenance

OpenAI Codex (GPT-5) was used during manuscript preparation to audit the text for AI-stylistic leakage, edit wording, check consistency between the manuscript and submission files, and prepare

453 disclosure language. The authors reviewed and remain responsible for all scientific content, analyses,
454 references, and conclusions. No AI-assisted tool is listed as an author or cited as a source. All figures
455 and tables were generated from the study data and code; no AI-generated images or multimedia are
456 included.

457 **Data Availability**

458 The CONAB reference parcels are openly available on Source Cooperative as the `br_conab`
459 field-boundary dataset distributed by the Field Boundaries for Agriculture (fibo) project
460 (https://source.coop/fibo/data/br_conab). Sentinel-1, Sentinel-2, and the AlphaEarth Satellite
461 Embedding V1 are accessed through Google Earth Engine. All code to reproduce the pipeline,
462 including reference extraction, tile sampling, feature download, probing, validation, and every table
463 and figure, is available in the study repository at <https://alphaearth-mg.abreufilho.com.br>; all
464 intermediate data are regenerable from the code with fixed seeds.

465 **Supplementary Materials**

466 No supplementary materials accompany this article.

467 **References**

- 468 1. d'Andrimont R, Claverie M, Kempeneers P, et al. AI4Boundaries: an open AI-ready dataset
469 to map field boundaries with Sentinel-2 and aerial photography. *Earth System Science Data*
470 *2023;15:317–29*.
- 471 2. Kerner H, Chaudhari S, Ghosh A, et al. Fields of The World: A Machine Learning Benchmark
472 Dataset for Global Agricultural Field Boundary Segmentation. *arXiv preprint arXiv:2409.16252*
473 *2024*.
- 474 3. Lowder SK, Scoet J, and Raney T. The Number, Size, and Distribution of Farms, Smallholder
475 Farms, and Family Farms Worldwide. *World Development* *2016;87:16–29*.
- 476 4. Waldner F and Diakogiannis FI. Deep learning on edge: Extracting field boundaries from satellite
477 images with a convolutional neural network. *Remote Sensing of Environment* *2020;245:111741*.
- 478 5. Persello C, Tolpekin VA, Bergado JR, and By RA de. Delineation of agricultural fields in
479 smallholder farms from satellite images using fully convolutional networks and combinatorial
480 grouping. *Remote Sensing of Environment* *2019;231:111253*.
- 481 6. Persello C, Grift J, Fan X, et al. AI4SmallFarms: A Dataset for Crop Field Delineation in South-
482 east Asian Smallholder Farms. *IEEE Geoscience and Remote Sensing Letters* *2023;20:2505705*.
- 483 7. Cong Y, Khanna S, Meng C, et al. SatMAE: Pre-training Transformers for Temporal and Multi-
484 Spectral Satellite Imagery. In: *Advances in Neural Information Processing Systems*. Vol. 35.
485 *2022*.

- 486 8. Jakubik J, Roy S, Phillips CE, et al. Foundation Models for Generalist Geospatial Artificial
487 Intelligence. arXiv preprint arXiv:2310.18660 2023.
- 488 9. Szwarcman D, Roy S, Fraccaro P, et al. Prithvi-EO-2.0: A Versatile Multi-Temporal Foundation
489 Model for Earth Observation Applications. arXiv preprint arXiv:2412.02732 2024.
- 490 10. Tseng G, Cartuyvels R, Zvonkov I, Purohit M, Rolnick D, and Kerner H. Lightweight, Pre-trained
491 Transformers for Remote Sensing Timeseries. arXiv preprint arXiv:2304.14065 2023.
- 492 11. Xiao A, Xuan W, Wang J, et al. Foundation Models for Remote Sensing and Earth Observation:
493 A Survey. arXiv preprint arXiv:2410.16602 2024.
- 494 12. Brown CF, Kazmierski MR, Pasquarella VJ, et al. AlphaEarth Foundations: An embedding
495 field model for accurate and efficient global mapping from sparse label data. arXiv preprint
496 arXiv:2507.22291 2025.
- 497 13. Ma Y, Shen Y, Swatantran A, and Lobell DB. Harvesting AlphaEarth: Benchmarking the Geospa-
498 tial Foundation Model for Agricultural Downstream Tasks. arXiv preprint arXiv:2601.00857
499 2026.
- 500 14. fiboa. br_conab: Field boundaries for Brazilian crops (CONAB). Source Cooperative dataset.
501 https://source.coop/fiboa/data/br_conab. 2024.
- 502 15. Alain G and Bengio Y. Understanding intermediate layers using linear classifier probes. arXiv
503 preprint arXiv:1610.01644 2016.
- 504 16. Moreira MA, Rudorff BFT, Barros MA, Faria VGCd, and Adami M. Geotecnologias para mapear
505 lavouras de café nos estados de Minas Gerais e São Paulo. *Engenharia Agrícola* 2010;30:1123–35.
- 506 17. Gorelick N, Hancher M, Dixon M, Ilyushchenko S, Thau D, and Moore R. Google Earth Engine:
507 Planetary-scale geospatial analysis for everyone. *Remote Sensing of Environment* 2017;202:18–27.
- 508 18. Drusch M, Del Bello U, Carlier S, et al. Sentinel-2: ESA’s Optical High-Resolution Mission for
509 GMES Operational Services. *Remote Sensing of Environment* 2012;120:25–36.
- 510 19. Torres R, Snoeij P, Geudtner D, et al. GMES Sentinel-1 mission. *Remote Sensing of Environment*
511 2012;120:9–24.
- 512 20. Camargo ÂPd and Camargo MBPd. Definição e esquematização das fases fenológicas do cafeiro
513 arábica nas condições tropicais do Brasil. *Bragantia* 2001;60:65–8.
- 514 21. McFeeters SK. The use of the Normalized Difference Water Index (NDWI) in the delineation of
515 open water features. *International Journal of Remote Sensing* 1996;17:1425–32.
- 516 22. Huete A, Didan K, Miura T, Rodriguez EP, Gao X, and Ferreira LG. Overview of the radiometric
517 and biophysical performance of the MODIS vegetation indices. *Remote Sensing of Environment*
518 2002;83:195–213.
- 519 23. Breiman L. Random Forests. *Machine Learning* 2001;45:5–32.
- 520 24. Pedregosa F, Varoquaux G, Gramfort A, et al. Scikit-learn: Machine Learning in Python.
521 *Journal of Machine Learning Research* 2011;12:2825–30.

- 522 25. Vincent L and Soille P. Watersheds in digital spaces: an efficient algorithm based on immersion
523 simulations. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 1991;13:583–98.
- 524 26. Walt S van der, Schönberger JL, Nunez-Iglesias J, et al. scikit-image: image processing in
525 Python. *PeerJ* 2014;2:e453.
- 526 27. Clinton N, Holt A, Scarborough J, Yan L, and Gong P. Accuracy Assessment Measures for
527 Object-based Image Segmentation Goodness. *Photogrammetric Engineering & Remote Sensing*
528 2010;76:289–99.
- 529 28. Wilcoxon F. Individual Comparisons by Ranking Methods. *Biometrics Bulletin* 1945;1:80–3.
- 530 29. Kirillov A, Mintun E, Ravi N, et al. Segment Anything. In: *Proceedings of the IEEE/CVF Inter-*
531 *national Conference on Computer Vision (ICCV)*. 2023:3992–4003. DOI: 10.1109/ICCV51070.
532 2023.00371.
- 533 30. Lavreniuk M, Kussul N, Shelestov A, et al. Delineate Anything: Resolution-Agnostic Field
534 Boundary Delineation on Satellite Imagery. arXiv preprint arXiv:2504.02534 2025.
- 535 31. Hong D, Gao L, Yokoya N, et al. More Diverse Means Better: Multimodal Deep Learning Meets
536 Remote-Sensing Imagery Classification. *IEEE Transactions on Geoscience and Remote Sensing*
537 2021;59:4340–54.