



Can a Phone-Sized Model Replace a Research-Grade One?

A pre-registered equivalence test of SegFormer (47 M) against FastViT-T8 (3.6 M) with DINOv2 distillation for SENASA coffee-disease grading

COFFEA ARABICA VAR. ESPERANZA LAAS · CARIBBEAN LOWLANDS OF COSTA RICA

Carly Chery

under the direction of Victor Hugo Morales Peña, Ph.D.

EARTH UNIVERSITY · GUACIMO, LIMÓN · COSTA RICA
PROTOCOL ON OSF · DOI:10.171794 · DEFENSE DECEMBER 2024

3.6 M FastViT-T8 phone-sized	10× smaller than the segmenter	±0.05 equivalence margin (κ)	1,400 annotated coffee leaves
--	--	---	---

Principal claims

- I. Pre-registered equivalence, not superiority.** SegFormer (47 M) vs. FastViT-T8 + DINOv2 distillation (3.6 M) for SENASA coffee-disease grading. Three baselines (ResNet-18, U-Net, DINOv2 linear probe) and four hypotheses locked on OSF before data collection.
- II. TOST declares "equal", not "no difference".** Two one-sided tests (Schiirmann, 1987; Lakens, 2017) let us affirmatively claim a 10× smaller classifier the practical equal of a research-grade segmenter to within ±0.05 κ — a finding classical superiority testing cannot deliver.
- III. Field-ready pipeline.** Offline Gradio laptop app (Spanish UI) for farm technicians now, LiteRT (Android) and Core ML (iOS) mobile exports prepared for Phase 2.

CONTRIBUTIONS

IMPACT
US \$1–2 B / year
Estimated annual losses across Latin America attributable to *H. vastatrix* (Andino et al., 2015). A laptop-friendly grader enables early, low-cost intervention for smallholder farmers.

NOVELTY
A pre-specified comparison of SegFormer against FastViT-T8 with DINOv2 LoRA distillation on the SENASA coffee-severity scale, framed under a TOST equivalence hypothesis against ResNet-18, U-Net, and DINOv2 linear-probe baselines.

III. Objective

To deliver a precise, consistent, and reproducible severity grader for rust and leaf-spot diseases on *Coffea arabica* var. Esperanza LAAS (first as an offline Gradio laptop application for field technicians, subsequently exported to mobile inference), and to determine, by pre-registered statistical test, whether a four-million-parameter classifier can stand as the practical equivalent of a forty-five-million-parameter segmentation model.

IV. Methodology



FIGURE 1. End-to-end methodological pipeline. Five stages from field capture on *Coffea arabica* var. Esperanza LAAS (I) through semi-automated annotation with STAPLE consensus (II), parallel training of SegFormer and FastViT under pre-registered hyperparameters (III), the pre-registered comparative statistical framework (IV) to laptop-first Gradio deployment with mobile export ready (V).



FIGURE 2. The same problem at three scales — plantation, leaf, and pests — each the unit of a different stakeholder decision. Photos: Wikimedia Commons, CC-BY-SA 3.0 (not from this study) dataset.

V. Dataset	VI. Pre-registered Targets	VII. Pre-registered Hypotheses																											
<p>1,400 × 3 unique leaves × orientations = 4,200 images — only the 1,400 axial are annotated</p> <p>Healthy (hard) 200</p> <p><i>H. vastatrix</i> G1 · G2 · G3 3 × 200</p> <p><i>M. citricolus</i> G1 · G2 · G3 3 × 200</p> <p>Total: 7 categories = 1,400</p> <p><small>Power: Mean CorLen = 1.400 at 2.000 for Cohen's $\kappa = 0.5$ via Spearman-McNemar. Leaf-level grouping produces CF leakage 2,000 images × axial count for augmentation EF robustness. Hybrid scoring ordinal metrics (QWK, κ) computed within each disease-disease-type regression (true / spurious / healthy) served as separate model-loss accuracy.</small></p>	<table border="1"> <thead> <tr> <th>METRIC</th> <th>METHOD</th> <th>TARGET</th> </tr> </thead> <tbody> <tr> <td>Quadratic Weighted κ (STAPLE)</td> <td>both</td> <td>≥ 0.80</td> </tr> <tr> <td>Mean IoU (Jaccard)</td> <td>SegFormer</td> <td>≥ 0.75</td> </tr> <tr> <td>Dice coefficient</td> <td>SegFormer</td> <td>≥ 0.80</td> </tr> <tr> <td>Boundary IoU (Dice)</td> <td>SegFormer</td> <td>report</td> </tr> <tr> <td>Poisson F_1</td> <td>both</td> <td>≥ 0.85</td> </tr> <tr> <td>Krippendorff α (SEGMENTATION)</td> <td>both</td> <td>≥ 0.80</td> </tr> <tr> <td>ECE (15 bins) (CALIBRATION)</td> <td>both</td> <td>report</td> </tr> <tr> <td>Latency P50 / P95 / P99</td> <td>both</td> <td>MLPerf</td> </tr> </tbody> </table> <p><small>Baselines reported alongside best methods: ResNet-18 + CF · U-Net + composed loss · DINOv2 linear probe.</small></p>	METRIC	METHOD	TARGET	Quadratic Weighted κ (STAPLE)	both	≥ 0.80	Mean IoU (Jaccard)	SegFormer	≥ 0.75	Dice coefficient	SegFormer	≥ 0.80	Boundary IoU (Dice)	SegFormer	report	Poisson F_1	both	≥ 0.85	Krippendorff α (SEGMENTATION)	both	≥ 0.80	ECE (15 bins) (CALIBRATION)	both	report	Latency P50 / P95 / P99	both	MLPerf	<p>H_0 Superiority (Spearman-McNemar · two-sided $\alpha = 0.05$) $\#SegFormer > \#FastViT$</p> <p>H_0 Non-inferiority (two-sided $\alpha = 0.05$; $\delta = 0.05$) $\#FastViT > \#SegFormer - 0.05$</p> <p>PRIMARY TEST H_0 TOST equivalence (two one-sided $\alpha = 0.05$ each; family $\alpha = 0.05$) $\#SegFormer - \#FastViT \leq 0.05$</p> <p>$H_0$ Efficiency (conditional on H_0 or H_1) $FLOP_{\#FastViT} \leq FLOP_{\#SegFormer} / 10$ $MEMO_{\#FastViT} \leq MEMO_{\#SegFormer} / 10$</p> <p><small>$\delta = 0.05 \kappa$ is below typical inter-rater disagreement on SENASA coffee grading ($\kappa = 0.60$–0.75; Ugarte et al., 2020; Tassi et al., 2023). Family-wise error controlled by Holm-Bonferroni over H1–H5, H4 conditional. BEST (Kruskal, 2013) with weakly-informative priors: BCa bootstrap, 10,000 resamples.</small></p>
METRIC	METHOD	TARGET																											
Quadratic Weighted κ (STAPLE)	both	≥ 0.80																											
Mean IoU (Jaccard)	SegFormer	≥ 0.75																											
Dice coefficient	SegFormer	≥ 0.80																											
Boundary IoU (Dice)	SegFormer	report																											
Poisson F_1	both	≥ 0.85																											
Krippendorff α (SEGMENTATION)	both	≥ 0.80																											
ECE (15 bins) (CALIBRATION)	both	report																											
Latency P50 / P95 / P99	both	MLPerf																											

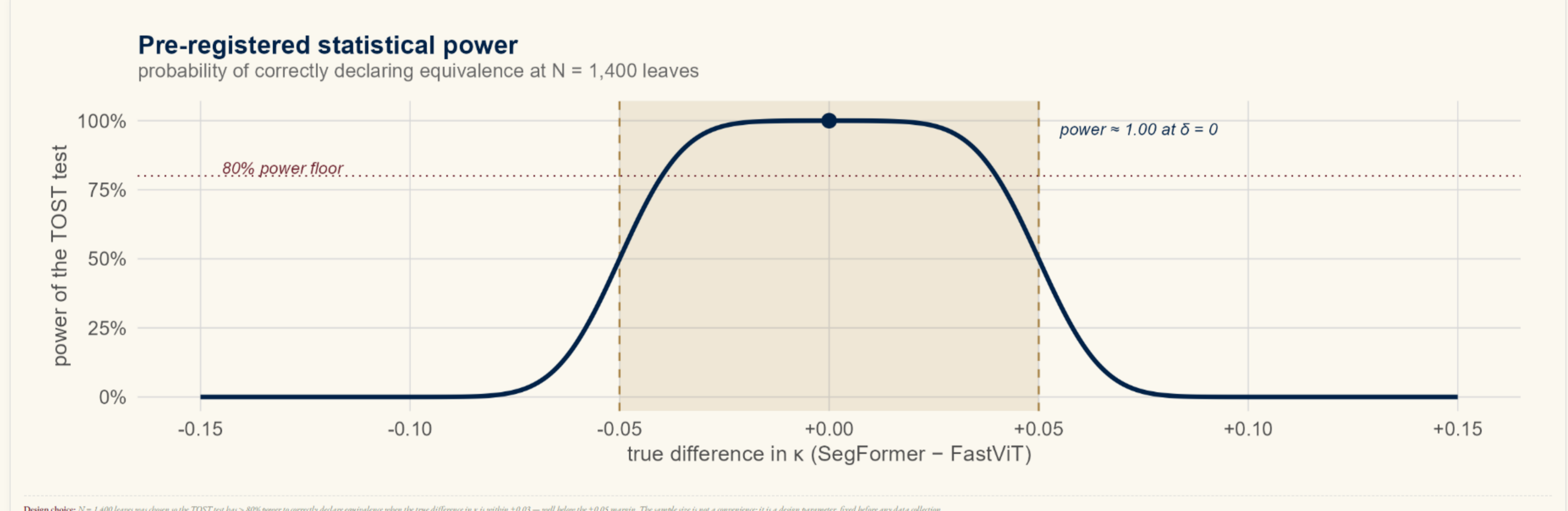
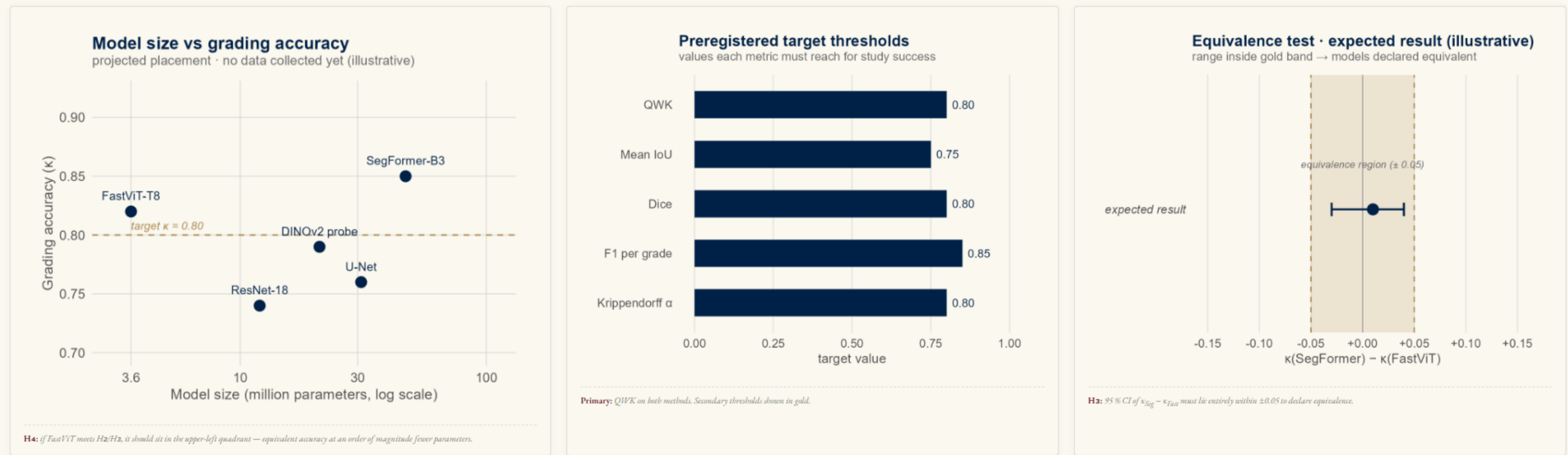
VIII. EXPERT VALIDATION PROTOCOL

Three phytopathologists · blind grading · 210 stratified images

Independent classification across the seven severity categories, with agreement assessed by Krippendorff α (ordinal) and Bland-Altman limits. Inter-expert α defines a practical ceiling reflecting the irreducible annotator disagreement characteristic of the SENASA scale.

3 EXPERT PHYTOPATHOLOGISTS
210 IMAGES · 10 PER CATEGORY
 $\alpha \geq 0.60$ KRIPPENDORFF (ORDINAL) TARGET

IX. Expected Results



X. Recommendations

- X. I. Deploy.** FastViT-T8 as an offline Gradio Laptop app (Spanish UI). Mobile exports — LiteRT (Android) and Core ML (iOS) — ready for Phase 2.
- X. II. Generalise.** Extend beyond Esperanza LAAS via farmer-cooperative partnerships. SegFormer handles on-site infection natively once training data on dual-infected leaves exists.
- X. III. Trust.** Publish segmentation overlays with every grade so an agronomist can audit why. Governance must ensure the tool serves farmers, never survives them.

XI. Limitations & Scope

Single site (EARTH Forestal Farm) · single cultivar (*Coffea arabica* var. Esperanza LAAS) · co-infected leaves out of scope · single season (inter-annual drift not evaluated) · field-usability study deferred · pilot empirical results pending (Q2–Q3 2026).

SELECTED REFERENCES

- Andino, J. et al. (2015). The coffee rust crisis in Colombia and Central America (2008–2013). *Food Security*, 7, 303–321.
- Figueras, J.G.M. et al. (2020). Deep learning for classification and severity estimation of coffee leaf blight stress. *Comput. Electron. Agric.*, 169, 105162.
- Kruskal, J.B. (2013). Bayesian optimization supersedes the tree (BHST). *J. Exp. Psychol. Gen.*, 142, 573–603.
- Lakens, D. (2017). Equivalence tests: a practical primer. *Soc. Psychol. Pers. Sci.*, 8, 355–362.
- Li, J. et al. (2025). DYNAM: enhanced multi-scale segmenter applying model for leaf disease segmentation. *Frontiers in Plant Science*.
- Osipov, M. et al. (2024). DINOv2: learning robust visual features without supervision. *TMLH*.
- Schiirmann, D.J. (1987). A comparison of the two one-sided tests procedure for equivalence of bioavailability. *J. Pharmaceutical Biopharm.*, 15, 657–685.
- Shi, X. et al. (2021). Deep neural networks for rank-consistent ordinal regression (CORN). arXiv:2111.08852.
- Tassi, L.M. et al. (2023). Intra- and inter-annotator agreement of coffee leaf disease and pest. *Comput. Electron. Agric.*, 196, 106591.
- Vera, J.C.A. et al. (2023). FastViT: A fast hybrid vision transformer using structural quantization. *ICCV*.
- Wang, S.K., Zou, K.H. & Wilf, W.M. (2004). Simultaneous truth and performance level estimation. *IEEE TMI*, 23, 903–911.
- Xu, E. et al. (2021). Segformer: simple and efficient design for semantic segmentation with transformers. *NIPS*.

