

Comparing LLM and human reviews of social science research using data from Unjournal.org

David Reinstein

Valentin Klotzbücher

Tianmai Michael Zhang

2025-11-30

We will build and refine LLM tools to generate peer-reviews and ratings of impactful research, and compare these with human experts' work (esp. from Unjournal.org): to benchmark performance, understand AI's research taste, and develop tools to improve research evaluation and dissemination.

Table of contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 4 |
| 1.0.1 | Our work in context | 5 |
| 2 | Data and methods | 7 |
| 2.1 | Unjournal.org evaluations | 7 |
| 2.2 | LLM-based evaluation | 7 |
| 2.2.1 | Quantitative ratings and journal-ranking tiers | 7 |
| 3 | Results | 9 |
| 3.1 | Quantitative comparison: human vs. GPT-5 Pro | 9 |
| 3.2 | Qualitative comparison: detailed GPT-5 Pro evaluations | 11 |
| 3.2.1 | Qualitative comparison: Williams et al. (2024) | 11 |
| 4 | Discussion | 19 |
| 4.1 | Limitations | 19 |
| | References | 21 |
| | Appendices | 22 |
| A | LLM evaluation summaries and reasoning traces | 22 |
| A.1 | Adena and Hager 2024 | 22 |
| A.2 | Green et al. 2025 | 30 |
| A.3 | Kudymowa et al. 2023 | 37 |
| A.4 | Peterman et al. 2025 | 41 |
| A.5 | Williams et al. 2024 | 49 |

1 Introduction

Work in progress

Pages, metrics, and comparisons are under active development. Expect rough edges and frequent updates.

Is AI good at peer-reviewing? Does it offer useful and valid feedback? Can it predict how human experts will rate research across a range of categories? How can it help academics do this “thankless” task better? Is it particularly good at spotting errors? Are there specific categories, e.g. spotting math errors or judging real-world relevance, where it does surprisingly well or poorly? How does its “research taste” compare to humans?

If AI research-evaluation works it could free up a lot of scientific resources – perhaps \$1.5 billion/year in the US alone Aczel, Szaszi, and Holcombe (2021)) – and offer more continual and detailed review, helping improve research. It could also help characterize methodological strengths/weaknesses across papers, aiding training and research direction-setting. Furthermore, a key promise of AI is to directly improve science and research. Understanding how AI engages with research evaluations may provide a window into its values, abilities, and limitations.

In this project, we are testing the capabilities of current large language models (LLMs), illustrating whether they can generate research paper evaluations comparable to expert human reviews. The Unjournal systematically prioritizes ‘impactful’ research and pays for high-quality human evaluations, structured quantified ratings, claim identification and assessment, and predictions. In this project, we use an AI (OpenAI’s GPT-5 Pro model) to review social science research papers under the same criteria used by human reviewers for The Unjournal.

Each paper is assessed on specific dimensions – for example, the strength of its evidence, rigor of methods, clarity of communication, openness/reproducibility, relevance to global priorities, and overall quality. The LLM will provide quantitative scores (with uncertainty intervals) on these criteria and produce a written evaluation

Our initial dataset will include the 5 research papers that have existing Unjournal human evaluations. For each paper, the AI will generate: (1) numeric ratings on the defined criteria, (2) identification of the paper’s key claims, and (3) a detailed review discussing the paper’s contributions and weaknesses. We will then compare the AI-generated evaluations to the published human evaluations.

In the next phase, we will focus on papers currently under evaluation, i.e., where no human evaluation has been made public, to allow us to rule out any contamination.

1.0.1 Our work in context

Luo et al. (2025) survey LLM roles from idea generation to peer review, including experiment planning and automated scientific writing. They highlight opportunities (productivity, coverage of long documents) alongside governance needs (provenance, detection of LLM-generated content, standardizing tooling) and call for reliable evaluation frameworks.

Eger et al. (2025) provide a broad review of LLMs in science and a focused discussion of AI-assisted peer review. They argue: (i) peer-review data is scarce and concentrated in CS/OpenReview venues; (ii) targeted assistance that preserves human autonomy is preferable to end-to-end reviewing; and (iii) ethics and governance (bias, provenance, detection of AI-generated text) are first-class constraints.

Zhang and Abernethy (2025) propose deploying LLMs as quality checkers to surface critical problems instead of generating full narrative reviews. Using papers from WITHDRARXIV and an automatic evaluation framework that leverages “LLM-as-judge,” they find the best performance from top reasoning models but still recommend human oversight.

Pataranutaporn et al. (2025) asked four nearly state-of-the-art LLM models (GPT-4o mini, Claude 3.5 Haiku, Gemma 3 27B, and LLaMA 3.3 70B) to consider 1220 unique papers “drawn from 110 economics journals excluded from the training data of current LLMs”. They prompted the models to act “in your capacity as a reviewer for [a top-5 economics journal]” and make a publication recommendation using a 6-point scale ranging from “1 = Definite Reject...” to “6. Accept As Is...”. They asked it to evaluate each paper on a 10-point scale for originality, rigor, scope, impact, and whether it was ‘written by AI’. They also (separately) had LLMs rate 330 papers with the authors’ identities removed, or replacing the names with fake male/female names and real elite or non-elite institutions (check this) or with prominent male or female economists attached.

They compare the LLMs’ ratings with the RePEC rankings for the journals the papers were published in, finding general alignment. They find mixed results on detecting AI-generated papers. In the names/institutions comparisons, they also find the LLMs show biases towards named high-prestige male authors relative to high-prestige female authors, as well as biases towards elite institutions and US/UK universities. (Doublecheck the details here).

There have been several other empirical benchmarking projects, including work covered in LLM4SR: A Survey on Large Language Models for Scientific Research and [Transforming Science with Large Language Models: A Survey on AI-assisted Scientific Discovery, Experimentation, Content Generation, and Evaluation](#). (We will discuss these here.)

Zhang et al. (2025)

- AI conference paper data
- “employs LLM agents to perform pairwise comparisons among manuscripts”
- “significantly outperforms traditional rating-based methods in identifying high-impact papers” [by citation metrics]
- Some evidence of biases/~statistical discrimination based on characteristics like ‘papers from established research institutions’

Our project distinguishes itself in its use of *actual* human evaluations of research in economics and adjacent fields, past and *prospective*, including both reports, ratings, and predictions.¹ The Unjournal’s 50+ evaluation packages enable us to train and benchmark the models. Their pipeline of future evaluations allow for clean out-of-training-data predictions and evaluation. Their detailed written reports and multi-dimensional ratings also allows us to compare the ‘taste’, priorities, and comparative ratings of humans relative to AI models across the different criteria and domains. The ‘journal tier prediction’ outcomes also provides an external ground-truth² enabling a human-vs-LLM horse race. We are also planning multi-armed trials on these human evaluations (cf. Brodeur et al, 2025 and Qazi et al, 2025) to understand the potential for *hybrid* human-AI evaluation in this context.

Footnote, a fancier way to say this, from a grant application? Or from chatGPT?³

¹Other work has relied on collections of research and grant reviews, including NLPEER, SubstanReview, and the Swiss National Science Foundation. That data has a heavy focus on computer-science adjacent fields, and is less representative of mainstream research peer review practices in older, established academic fields. Note that The Unjournal commissions the evaluation of impactful research, often from high-prestige working paper archives like NBER, and makes all evaluations public, even if they are highly critical of the paper.

²About verifiable publication outcomes, not about the ‘true quality’ of the paper of course.

³Our approach differs from prior work by (i) focusing on structured, percentile-based quantitative ratings with credible intervals that map to decision-relevant dimensions used by The Unjournal; (ii) comparing those ratings to published human evaluations rather than using LLM-as-judge; and (iii) curating contamination-aware inputs (paper text extraction with reference-section removal and token caps), with a roadmap to add multi-modal checks when we score figure- or table-dependent criteria.

2 Data and methods

We draw on two main sources:

- 1) Human evaluations from [The Unjournal’s public evaluation data](#) (PubPub reports and the Coda evaluation form export).
- 2) LLM-generated evaluations using a structured JSON-schema prompt with `gpt-5-pro-2025-10-06` (knowledge cut-off: 30 September 2024).

2.1 Unjournal.org evaluations

We use The Unjournal’s public data for a baseline comparison. At The Unjournal each paper is typically evaluated (aka ‘reviewed’) by two expert evaluators¹ who provide quantitative ratings on a 0–100 percentile scale for each of seven criteria (with 90% credible intervals),² two “journal tier” ratings on a 0.0 - 5.0 scale,³ a written evaluation (resembling a referee report for a journal), and identification and assessment of the paper’s “main claim”. For our initial analysis, we extracted these human ratings and aggregated them, taking the average score per criterion across evaluators (and noting the range of individual scores).

All papers have completed The Unjournal’s evaluation process (meaning the authors received a full evaluation on the Unjournal platform, which has been publicly posted at unjournal.pubpub.org). The sample includes papers spanning 2017–2025 working papers in development economics, growth, health policy, environmental economics, and related fields that The Unjournal identified as high-impact. Each of these papers has quantitative scores from at least one human evaluator, and many have multiple (2-3) human ratings.

2.2 LLM-based evaluation

2.2.1 Quantitative ratings and journal-ranking tiers

Following The Unjournal’s [standard guidelines for evaluators](#) and their [academic evaluation form](#), evaluators are asked to consider each paper along the following dimensions: **claims & evidence**,

¹Occasionally they use 1 or 3 evaluators.

²See their guidelines [here](#); these criteria include “Overall assessment”, “Claims, strength and characterization of evidence”, “Methods: Justification, reasonableness, validity, robustness”, “Advancing knowledge and practice”, “Logic and communication”, “Open, collaborative, replicable science”, and “Relevance to global priorities, usefulness for practitioners”

³“a normative judgment about ‘how well the research should publish’ ” and “a prediction about where the research will be published”

methods, logic & communication, open science, global relevance, and an **overall** assessment. Ratings are interpreted as percentiles relative to serious recent work in the same area. For each metric, evaluators are asked for the midpoint of their beliefs and their 90% credible interval, to communicate their uncertainty. For the journal rankings measure, we ask both “what journal ranking tier should this work be published in? (0.0-5.0)” and “what journal ranking tier will this work be published in? (0.0-5.0)”, with some further explanation. The full prompt can be seen in the code below – essentially copied from the Unjournal’s guidelines page.

We captured the versions of each paper that was evaluated by The Unjournal’s human evaluators, downloading from the links provided in The Unjournal’s Coda database.

We evaluate each paper by passing the PDF directly to the model and requiring a strict, machine-readable JSON output. This keeps the assessment tied to the document the authors wrote. Direct ingestion preserves tables, figures, equations, and sectioning, which ad-hoc text scraping can mangle. It also avoids silent trimming or segmentation choices that would bias what the model sees.

We enforce a JSON Schema for the results. The model must return one object for each of the named criteria including a midpoint rating and a 90% interval for each rating. This guarantees that every paper is scored on the same fields with the same types and bounds. It makes the analysis reproducible and comparisons clean.

We request credible intervals (as we do for human evaluators) to allow the model to communicate its uncertainty rather than suggest false precision; these can also be incorporated into our metrics, penalizing a model’s inaccuracy more when it’s stated with high confidence.

Relying on GPT-5 Pro, we use a single-step call with a reasoning model that supports file input. One step avoids hand-offs and summary loss from a separate “ingestion” stage. The model reads the whole PDF and produces the JSON defined above. We do not retrieve external sources or cross-paper material for these scores; the evaluation is anchored in the manuscript itself.

The Python pipeline uploads each PDF once and caches the returned file id keyed by path, size, and modification time. We submit one background job per PDF to the OpenAI Responses API with “high” reasoning effort and server-side JSON-Schema enforcement. Submissions record the response id, model id, file id, status, and timestamps.

A separate script polls job status and, for each completed job, retrieves the raw response, extracts the first balanced top-level JSON object, and writes both the raw response and parsed outputs to disk.

3 Results

Here we present preliminary results, starting with a comparison of the LLM-generated quantitative ratings (model: `gpt-5-pro`, see the [previous section](#)) with human evaluations across [the Unjournal’s criteria](#).

3.1 Quantitative comparison: human vs. GPT-5 Pro

We first use the earlier GPT-5 Pro evaluation run that covered all papers in our Unjournal sample with a simpler JSON-schema prompt. Figure 3.1 shows the overall percentile ratings from this initial run, averaged across human evaluators and compared to the LLM’s “overall” scores for each paper.

Figure 3.1: Comparison of Human vs LLM overall percentile ratings

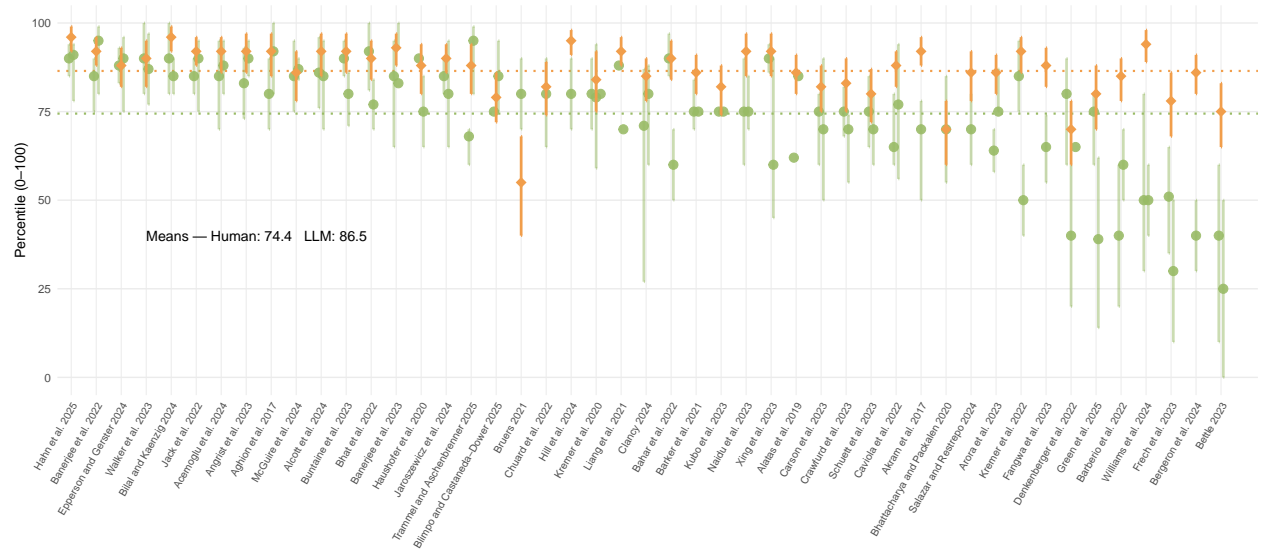


Figure 3.2 shows a heatmap of the differences between human and LLM mean ratings across all evaluation criteria. Positive values (in green) indicate that humans rated the paper higher than the LLM, while negative values (in orange) indicate the opposite.

Figure 3.2: Heatmap of Human minus LLM mean ratings across evaluation criteria

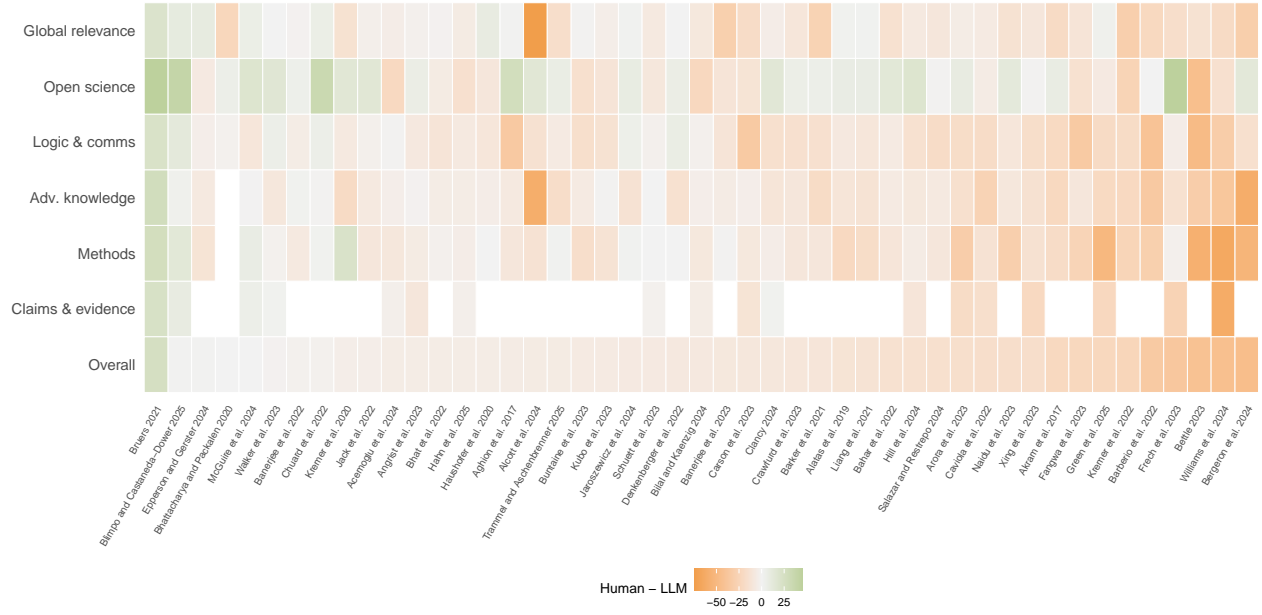
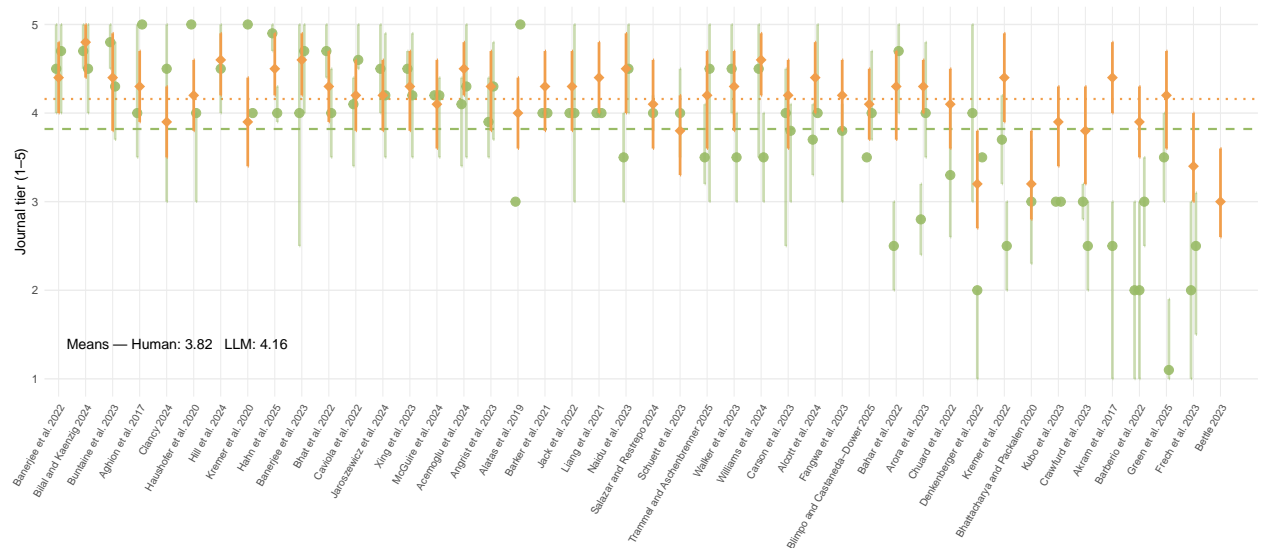


Figure 3.3: Comparison of Human vs LLM journal tier ratings (should be published in)



3.2 Qualitative comparison: detailed GPT-5 Pro evaluations

To understand what GPT-5 Pro is actually responding to, we re-ran the model on four focal papers (Adena and Hager 2024; Peterman et al. 2024; Williams et al. 2024; Green, Smith, and Mathur 2025) using a refined prompt.

This second run keeps the same quantitative metrics but additionally requires a diagnostic summary of about 1,000 words and high-effort reasoning, with the full reasoning trace returned by the “thinking” model. For each paper we can therefore inspect:

- the LLM’s quantitative scores and journal-tier predictions,
- the hidden reasoning steps used to arrive at those scores, and
- the token usage and approximate API cost of the evaluation.

We start by examining the Williams et al. (2024) evaluation in detail and then show the analogous summaries for the other four focal papers. In the next step we will juxtapose these LLM assessments with the human evaluators’ written reports.

3.2.1 Qualitative comparison: Williams et al. (2024)

In the refined run, GPT-5 Pro reads about 28,704 input tokens and produces 6,327 visible output tokens plus 5,120 reasoning tokens.

At current API prices this evaluation costs roughly \$1.80.

The table below shows the model’s percentile ratings and 90% credible intervals for the Unjournal criteria.

Table 3.1: GPT-5 Pro percentile ratings for Williams et al. (2024)

| Criterion | Midpoint | Lower 90% | Upper 90% |
|---|----------|-----------|-----------|
| Overall assessment | 86 | 86 | 86 |
| Claims & evidence | 78 | 78 | 78 |
| Methods | 74 | 74 | 74 |
| Advancing knowledge and practice | 92 | 92 | 92 |
| Logic and communication | 84 | 84 | 84 |
| Open, collaborative, replicable science | 63 | 63 | 63 |
| Relevance to global priorities | 94 | 94 | 94 |

Table 3.2: GPT-5 Pro journal tier ratings for Williams et al. (2024)

| Measure | Score | Lower 90% | Upper 90% |
|--------------------------------|-------|-----------|-----------|
| Deserved journal tier (should) | 4.4 | 4.4 | 4.4 |
| Predicted journal tier (will) | 4.8 | 4.8 | 4.8 |

For Williams et al. (2024), GPT-5 Pro assigns a high overall percentile score (86/100), with particularly strong ratings for advancing knowledge (92) and global relevance (94), and a relatively

favourable view of methods (74) and claims and evidence (78). It judges the paper as deserving publication in a high-tier journal (tier-should 4.4/5, tier-will 4.8/5).

In its diagnostic summary and reasoning trace (printed below), the model identifies many of the same issues highlighted by the human evaluators: heavy reliance on a regrowth dataset with low producer accuracy and substantial omission error; temporal leakage from contemporaneous predictors; uncalibrated random-forest probabilities used to derive the 215 Mha estimate; unrealistically narrow confidence intervals; coarse predictors driving 30-m predictions; a liberal definition of land “available for restoration”; and incomplete uncertainty quantification for the carbon overlay and permanence. It also notes that code is not fully open, limiting replication despite open data and published maps.

By contrast, both human evaluators at The Unjournal assign much lower overall ratings (50/100) and are substantially more critical of methods and claims: methods scores of 20 and 5, and claims scores of 20 and 5, respectively. They explicitly argue that methodological problems “fundamentally challenge the validity and utility of the central 215 Mha estimate of regeneration potential,” even while rating the paper’s contribution to knowledge and global relevance as high and considering it suitable for high-tier journals conditional on major revisions.

Qualitatively, the LLM and human experts are strikingly aligned on what the main problems are: they converge on concerns about biased input data, temporal leakage, incomplete probability calibration, under-stated uncertainty, domain definition, and the gap between “biophysical potential” and realistic policy use. However, they diverge sharply on how serious these problems are. The human evaluators treat them as sufficient to render the headline estimates low-credibility and heavily qualified for policy use; GPT-5 Pro instead regards these issues as important but ultimately compatible with a high overall rating, robust methods, and near top-journal quality.

This case suggests that, at least in this configuration, the LLM can reproduce sophisticated methodological critiques and uncertainty language but tends to under-penalise these shortcomings in its quantitative ratings, especially for high-profile, high-impact work.

Model reasoning trace – Williams et al. 2024

Analyzing model limitations

I see another limitation in the model’s assumptions. They assumed stationarity from 2000-2016 when predicting for 2015 to 2030. They included present conditions like 2018 tree cover but didn’t factor in future climate or socio-economic changes. By excluding socio-economic predictors, they reduced potential model improvements slightly, but this could bias predictions in areas where socio-economic factors play a significant role, potentially leading to unrealistic outcomes in “biophysical potential.” They clarify it’s about modeling potential, not feasibility.

Examining dataset limitations

I notice another issue with the primary regrowth dataset used, which came from Global Forest Watch to map regrowth over 5 m. In humid tropics, the producer’s accuracy is low at just 18.7% because of omission errors, although user accuracy is high at 85.1%. The detection biases may skew the training sample and affect model calibration. Their conservative approach to mapping regrowth likely favors more accessible and larger areas, and many small regrowth patches are excluded. It’s unclear if the training sample is representative across different biomes without proper stratification.

Clarifying domain definitions

I need to define the domain used in the study. They included all parts of forest biomes within $\pm 25^\circ$ but excluded water, bare land, urban areas, and so on, while purposely opting for a liberal domain. They used 2018 tree cover to exclude forests but may misclassify certain savannas or drylands as part of forest biomes. While they mention not including savannas, their definition could still capture mosaic-savanna zones. Additionally, there's a scale mismatch since predictions at a finer 30 m resolution are computed alongside coarser predictors, which could impact accuracy.

Assessing the weighted area approach

I see another issue with the weighted area approach, where the probability of regrowth times pixel area is calculated. However, the random forest classification is trained on presence/absence from 2000-2016, and the output probabilities aren't properly calibrated. Calibration techniques like Platt scaling or isotonic regression are necessary to interpret results as probabilities. They also didn't present AUC, only accuracy, which may mislead since balanced sampling artificially sets prevalence to 0.5. Without proper probability calibration, the estimate of expected area could be questionable and not reflect actual conditions.

Reviewing definitions and calibration

They used specific definitions, claiming that "the continuous potential" can be seen as the probability of natural regeneration per pixel, with the weighted-area value representing the expected area for regeneration. However, they didn't perform any calibration on this. Their accuracy and out-of-bag accuracy are reported as 0.878, but they failed to discuss probability calibration, which is important. Regarding carbon sequestration potentials, they calculated totals and ranges from previous datasets, but it's unclear where the range comes from—possibly from Monte Carlo methods or carbon dataset quartile variation, yet they didn't show how they computed uncertainty.

Identifying dataset uncertainties

The dataset presents uncertainties by grid cell, but their method multiplies 30-year totals without mentioning how uncertainty is handled. The "range" they provide likely comes from minimum and maximum values across these uncertainties, which is unclear. They need to clarify how uncertainty propagates and address independence and omitted below-ground carbon, since they reference adding 22–28% from literature. Additionally, they didn't consider factors like degradation baseline, leakage, or permanence. While they frame the "biophysical potential," the abstract could mislead policy implications with its stated above-ground carbon sequestration potential of 23.4 Gt C over 30 years.

Checking carbon sequestration claims

I noticed they mentioned that "new carbon sinks could increase global sequestration by 14.3% per year" but also stated "mitigating current pantropical forest carbon losses by 90.5% per year," which might be inconsistent. They claim assisted natural regeneration can sequester significant carbon but didn't show how they calculated this. It seems to derive from 23.4 Gt over 30 years, equating to 0.78 Gt C per year. With a net carbon loss of about 0.86 Gt C per year, I need to verify Baccini 2017's values, which reported a net carbon source of 1.1 Gt C per year. My memory on this is a bit fuzzy!

Evaluating carbon loss mitigation claims

They state that "mitigating current pantropical forest carbon losses by approximately 90.5% per year," referencing Baccini, raises some questions. If regeneration can offset about 0.8 Gt per year, that suggests total losses of around 0.9 Gt per year. While this might be plausible, there are logical leaps involved since their potential benefits might not actually be realized.

It's a bit tricky to interpret, yet they present it confidently as potential. I need to consider these assumptions carefully!

Analyzing dataset issues

They wisely exclude plantation regrowth from their analysis, but it seems their training dataset might still contain misclassified plantations. They reported a previous algorithm with 90.6% accuracy; however, an area-based producer accuracy of just 18.7% in humid tropics suggests significant under-detection. This non-random classification error could undermine generalizability. While they argue for the robustness of the detected patches, the validation accuracy of 87.9% relies on potentially biased training dataset presence points. Field validation is unlikely, and poor accuracy in Southeastern Asia raises another red flag.

Evaluating modeling assumptions

The modeling approach assumes that predicted 'potential' is valid from 2015 to 2030 without considering changes, which is questionable since the climate is constantly shifting. Additionally, the "distance to forest" will likely change as regrowth occurs, yet they didn't model this dynamic feedback. They mentioned underestimating impacts like positive feedback and the exclusion of savanna/dryland. Climate change could also hinder regrowth, leading to ambiguity. They present a figure showing "available for restoration" area by country, but they don't clarify how this was calculated, leaving some mystery around the methodology used.

Analyzing definitions and reproducibility

They define the study's domain as all tropical forest biomes, excluding areas like water, urban, and existing forests, while including factors like 'distance to water' as covariates. It seems that "available for restoration" refers to the total area of this domain, but their disclaimers indicate some harvest/regrowth classifications might be included. Although they mention using a forestry dataset, they didn't clearly define "available for restoration." Additionally, requesting code via email isn't ideal; it should be publicly accessible. While they use accessible datasets, the lack of code hampers reproducibility, and the confidence intervals aren't credible.

Evaluating accuracy and methodology

They mention the model overcame potential overfitting with an out-of-bag (OOB) accuracy of 87.8%. They used a balanced sample and omitted socioeconomic features due to similar accuracy, but I wonder if the limited role of socioeconomic data in predicting presence versus absence, as indicated by remote sensing, could affect results. They also disclaim that the maps are only one input, which is good. Since this is published in Nature, it seems high quality, but I believe we should look beyond that. We need to propose credible metrics.

Ranking evaluation criteria

Based on the content, I would rank high—maybe around the 85th percentile overall.

For claims and evidence, they articulate strong claims and provide robust mapping and a predictive model. However, there are issues like reliance on one training dataset with major omissions and unrealistic confidence intervals. I'd place this at about 75.

For methods, they've implemented random forest modeling and accounted for autocorrelation, yet there's heavy reliance on coarse data. This leads me to perhaps a 70.

Regarding advancing knowledge, their maps on natural regeneration potential are valuable, globally relevant for restoration planning, so I'd mark this as a 90.

On logic and communication, while the writing is clear and disclaimers are present, there are leaps in concluding about carbon losses and some misuses of "CI." I'd score this around 82.

For open science, while the data is available on Zenodo and the input data is open, the code is only available on request, which limits replicability. I'd give a score of 65.

Finally, in terms of global relevance, this work ties in well with global priorities and quantifies information by country, so I'd rate this a solid 92.

Assigning tier ratings

I need to assign "tier_should" and "tier_will" ratings from 0 to 5. Based on quality, I think "should" deserves a high score, perhaps around 4.5 or 4.7, due to some identified issues. For "will," since it's already published in Nature, I'd predict around 4.6 to 4.8. However, they want continuous numbers, so I'll set "should" at 4.4 and "will" at 4.7.

If I consider possible underestimation due to method-related issues, perhaps "should" could be 4.2 or 4.4, while "will" could be 4.8 because of strong interest. Still, the dataset not being open-coded is a consideration, but it's published in Nature, so I want to remain fair.

Now, I need to craft a summary, aiming for around 1000 words. I'll start by summarizing the approach and identifying major issues.

Highlighting dataset issues

I'm noticing that the training dataset has biases, especially regarding natural regrowth mapping, which shows significant omissions in humid biomes. This can lead to a non-random sample of presences, and the model might learn from incomplete detections, distorting variable relationships and potentially mispredicting areas.

Next, regarding probability calibration and the "expected area" method, they treat random forest outputs as probabilities, but without proper calibration and baseline prevalence, the resulting estimates of 215 million hectares could be inaccurate. They present highly precise confidence intervals that don't seem credible because they didn't account for uncertainty across inputs.

Lastly, concerning spatial autocorrelation, while they tried to evaluate it, accuracy tends to decline with distance, so their optimism might not be justified. They didn't apply block cross-validation and reported overall accuracy, avoiding metrics like AUC. The balanced classification gives an unrealistic sense of accuracy, and there are no confusion matrices or biome-specific metrics presented beyond overall accuracy.

Addressing scale and domain issues

I see that there's a scale mismatch with the 30 m predictions derived from 300 m to 1 km predictors, which could cause smoothing and artifacts. They've acknowledged this with visual evidence in Fig. 1, suggesting it could impact local planning as they consider these results to be preliminary.

The domain definition appears overly broad, including areas where natural regeneration may be undesirable—like non-forest ecosystems, productive farmland, or communities—which could potentially mislead their "available for restoration" metrics.

Lastly, they assume conditions from 2000 to 2016 without modeling climate change or socioeconomic factors, which could really undermine their 2030 predictions. Also, the carbon accumulation overlay from the Cook-Patton dataset has a scale mismatch issue, with 1 km data being downscaled to 30 m without considering edges and heterogeneity.

Evaluating carbon accounting and code accessibility

They don't initially account for below-ground carbon, indicating an additional 22–28% should be included, and they don't consider the risk of reclearance. Their discussions on leakage and additionality seem inadequate since their yearly rates are used to support claims of significant offsetting of losses, which likely exaggerates realized mitigation.

While they exclude plantations, algorithmic misclassification might still exist, and they haven't incorporated a classification probability step.

Regarding code accessibility, it's not publicly available, making replication difficult. Although they provide predicted maps on Zenodo, the training and modeling code needing to be requested further hampers replicability. Additionally, using 2018 tree cover for distance predictions while training on 2000 cover could lead to variable shifts, potentially degrading the predictions since they haven't cross-validated changes.

Observing assumptions and suggestions for improvement

The weighted area assumption sums across all pixels to determine "expected area," but it overlooks spatial dependence and positive feedback, which are mentioned but not integrated. The country-level summarization is good, although it uses administrative boundaries without addressing ecoregion differences, and there's a lack of error budget for these country-level sums.

They note that predictions may degrade in Southeastern Asia, but don't confirm this for later years. For improvements, I should suggest block cross-validation, probability calibration via isotonic regression, and incorporating socioeconomic variables.

Overall, I thought the score should be around 85, but I'm reconsidering an 86.

Finalizing ranking and summary considerations

I'm setting the overall score to 86, considering it's intended to reflect percentiles within [0, 100]. For research in remote sensing global mapping, I'd lean toward an 88, but due to certain issues, I might adjust this to 82. However, advancing knowledge could be rated higher at 90, leading me to average around 86, so choosing 87 seems fair.

I also need to fill out the JSON carefully, ensuring that the assessment summary is a concise paragraph of about 800-1100 words and omitting any author names or external knowledge outside the provided PDF. The confidence intervals will match the midpoint of 86. For claims and evidence, I'm thinking a score of 79?

Setting scores and producing JSON details

I've decided to set the score of 78 for overall assessment, with specifics for other categories: 73 for methods, 92 for advancing knowledge, 84 for logic communication, 62 for open science, and 94 for global relevance. For tier should, I'll use a score of 4.4, and for tier will, I'll set it at 4.8, all while keeping consistency in mind with overall scores at 86, advancing at 92, and global relevance at 94.

The JSON needs to include an assessment summary of around 1,000 words and metrics for each score calculated with midpoint, lower, and upper bounds. I'll also avoid naming authors or their institutions, but I can reference previous analyses without using specific names or citation numbers.

Crafting the summary points

I'm going to paraphrase the training dataset description, ensuring I accurately reference the analysis of natural regrowth from 2000 to 2016 without including citation numbers. I should say, "The study builds on an existing pantropical remote sensing analysis that mapped natural regrowth and distinguished it from plantations."

Next, I want to finalize the summary by outlining key points. For data, I'll mention potential mapping across tropical forests at a 30 m resolution, calculating weighted areas for 215 Mha, and noting variations across ecoregions.

For methods, I'll detail the balanced random forest approach, variable selection, the training dataset derived from regrowth mapping, and the overall accuracy. I'll also highlight potential issues like omissions in the training dataset and domain definition relative to thresholds.

Assessing model limitations

I need to go over some aspects of the model’s predictive capabilities. First, I see that there’s a lack of probability calibration in the balanced training; the predicted probabilities aren’t calibrated, which might lead to bias in the weighted area approach and unrealistic confidence intervals.

Next, there’s the issue of carbon overlay and uncertainty propagation; downscaling seems to be done naively. The coarser covariate resolution could degrade the 30 m predictions, causing blockiness. Socio-economic variables are omitted in terms of feasibility versus potential, and there’s a failure to integrate dynamic processes like climate change impacts.

I should also note that accuracy in Southeastern Asia is poor, and the code’s availability is limited, impacting replicability due to the heavy dataset.

Lastly, I’ll mention that the additional risk of reclearance isn’t accounted for in the carbon measurements. I need to finalize this while including suggestions and discussing uncertainties qualitatively as we maintain strict JSON formatting.

i Diagnostic assessment – Williams et al. 2024

This paper delivers a high-resolution, pantropical map of where deforested lands have biophysical potential to undergo natural forest regeneration, and it aggregates this to country and biome summaries and an overlay-based carbon removal potential. The approach is transparent in intent—predicting potential, not feasibility—but several methodological choices and data limitations affect inference strength and precision. The dependent variable (presence/absence of natural regeneration during 2000–2016) is derived from a prior global remote-sensing classification that conservatively emphasized omission errors in humid regions and required a minimum patch size; area-based producer accuracy for natural regrowth in the humid biome was low while user accuracy for detected patches was high. This means the positive training class is a non-random subset of true regrowth, with systematic under-detection likely concentrated in particular forest types, geographies, and patch sizes. The predictive model then inherits these detection and sampling biases, potentially distorting learned relationships and shifting mapped potential away from contexts where regrowth was common but poorly detected. Model fitting uses random forests on millions of stratified points, with a sensible variable-selection procedure; predictions are ultimately based on biophysical predictors alone because a combined socioecological model yielded very similar accuracy. Omitting socioeconomic predictors makes the maps more temporally stable and widely available, but it also reframes outputs as “biophysical potential” and risks overpredicting potential in places where human pressure or governance constraints are binding. Validation emphasizes overall accuracy around 88% using a balanced validation set; however, accuracy declines at moderate distances from training locations, indicating residual spatial autocorrelation and some optimism in headline accuracy. The use of accuracy rather than threshold-free metrics, and the balancing of classes (prevalence set to 50%) during training and validation, further complicates interpretation of both accuracy and the calibration of predicted probabilities. The paper interprets the continuous model output as a probability and converts this to “expected area” by multiplying by pixel area and summing, yielding 215 Mha. Yet the calibration of random-forest scores to true probabilities is not established; without calibration (e.g., isotonic or Platt scaling) under the real-world prevalence, the expected-area identity may be biased. The extraordinarily narrow “confidence intervals” reported for area (± 0.22 Mha at global scale) are not credible statistical uncertainty; they reflect computational or rounding artifacts, not propagation of

uncertainty from training data errors, cross-validated prediction error, input-layer uncertainty, or probability calibration. A threshold (>0.5) sensitivity shows a substantially larger area, highlighting dependence on interpretation choices. Many predictors are coarse (250 m–1 km) while predictions are at 30 m; this scale mismatch introduces blockiness and may overstate local precision. The domain of “available for restoration” is defined liberally by excluding only forests, water, urban, bare, and forestry areas; while appropriate for modeling, its presentation alongside country totals risks being misread as areas that are socially or legally available for reforestation. Temporal assumptions also matter: the model projects to 2015/2030 by assuming that 2000–2016 relationships hold and by updating some forest-context layers with 2018 data; it does not incorporate climate change trajectories, shifting fire regimes, or evolving socioeconomics that strongly affect regeneration and permanence. Carbon estimates are calculated by downscaling a 1 km accumulation dataset to 30 m, multiplying by modeled potential, and summing; the stated range lacks a clearly defined uncertainty model and appears not to propagate classification, calibration, or permanence risks. The paper is commendably explicit that realized sequestration will be lower due to leakage, reclearance, and policy constraints, yet some statements comparing potential annual removals to current pantropical losses could be interpreted too strongly given these omitted risks and uncertainties. Spatial performance is uneven, with lower accuracy in parts of Southeast Asia; country and biome summaries lack uncertainty decomposition, which would be useful for decision makers. On replicability, large gridded outputs are publicly archived and inputs are open, but code is available only upon request; combined with the volume of preprocessing, this raises barriers to full reproduction, code review, and extension. Strengths include a clearly articulated objective focused on actionable planning, a very large training sample, careful variable selection, explicit treatment of autocorrelation effects (even if not using full spatial block cross-validation), and extensive caveats about use and governance. The resulting layers will be highly valuable as a first-pass screen for assisted natural regeneration opportunities and for national restoration dialogues. The main improvements needed relate to uncertainty characterization (probability calibration, spatially explicit error propagation, and country-level error budgets), spatial cross-validation that reduces spatial leakage, clearer separation of statistical precision from conceptual ranges, dynamic scenario analysis including plausible climate and human-pressure changes, and full open release of code. Overall, this is an influential, practically useful mapping contribution with limitations that mainly concern probabilistic interpretation, uncertainty, and temporal and socioecological external validity.

| Paper | Input tokens | Output tokens | Reasoning tokens | Total tokens | Est. cost (USD) |
|----------------------|--------------|---------------|------------------|--------------|-----------------|
| Peterman et al. 2025 | 18762 | 7617 | 6208 | 26379 | 1.94 |
| Adena and Hager 2024 | 24234 | 7019 | 5312 | 31253 | 1.84 |
| Williams et al. 2024 | 28704 | 6327 | 5120 | 35031 | 1.80 |
| Kudymowa et al. 2023 | 65096 | 4147 | 2432 | 69243 | 1.77 |
| Green et al. 2025 | 22904 | 5884 | 3904 | 28788 | 1.52 |

4 Discussion

4.1 Limitations

Sample size and representativeness: We evaluated on only ~40–50 papers, all in the social sciences and specifically chosen by The Unjournal for evaluation (which means they were considered high-impact or interesting). This is not a random sample of research literature. The papers also skew toward empirical and policy-relevant topics. The AI’s performance and alignment might differ in other fields (e.g., pure theory, biology) or on less polished papers.

Human agreement as a moving target: The Unjournal human evaluations themselves are not a single ground truth. As evidence of this, we note substantial variability between reviewers.

Potential AI knowledge contamination: We attempted to prevent giving the AI any information about the human evaluations, but we cannot be 100% sure that the model’s training data didn’t include some fragment of these papers, related discussions, or even The Unjournal evaluations. We will be able to exclude this for the evaluations of *future* Unjournal evaluations.

Model limitations and “alignment” issues: While powerful, is not a domain expert with judgment honed by years of experience. It might be overly influenced by how a paper is written (fluency) or by irrelevant sections. It also tends to avoid extremely harsh language or low scores unless there is a clear reason, due to its alignment training to be helpful/polite – this could explain the general score inflation we observed. The model might fail to catch subtle methodological flaws that a field expert would notice, or conversely it might “hallucinate” a concern that isn’t actually a problem. Without ground truth about a paper’s actual quality, we used human consensus as a proxy; if the humans overlooked something, the AI could appear to “disagree” but possibly be pointing to a real issue.

(There is also evidence, e.g. Pataranutaporn et al. (2025) that LLMs show biases towards more prestigious author names, institutions, and towards male prestigious authors. We will provide further evidence on this in the next iterations, de-identifying the work under LLM evaluation.)

Scoring calibration: The AI was prompted to use the 0–100 percentile scale, but calibrating that is hard. Humans likely had some calibration from guidelines or community norms (e.g. perhaps very few papers should get above 90). The AI might have been more liberal in using the high end of the scale (hence higher means). In future, a different prompt or examples could calibrate it to match the distribution of human scores more closely. We also only took one run from the AI for each paper; LLM outputs can have randomness, so a different run might vary slightly. (To do: aggregate across multiple runs.)

Small differences and rounding: Our analysis treated the AI’s numeric outputs at face value. Small differences (e.g. AI 85 vs human 82) might not be meaningful in practice – both indicate a similar qualitative assessment (“very good”). Some of our metrics (like kappa) penalize any difference, even if minor. Thus, the “low agreement” statistics might sound worse than the reality where in

many cases AI and humans were only off by a few points. We intend to analyze the distribution of absolute differences: a large portion might be within say ± 5 points which could be considered essentially agreement in practice. The credible intervals add another layer: sometimes an AI's score fell outside a human's interval, but overlapping intervals could still mean they agree within uncertainty. We did observe that AI's intervals were often narrower than humans' (LLM tended to be confident, giving ~ 10 -point spreads, whereas some human evaluators gave 20-point or left some intervals blank), which is another aspect of calibration.

References

- Aczel, Balazs, Barnabas Szaszi, and Alex O Holcombe, “A billion-dollar donation: Estimating the cost of researchers’ time spent on peer review,” *Research integrity and peer review*, 6 (2021), 1–8 (Springer).
- Adena, Maja, and Anselm Hager, “Does online fundraising increase charitable giving? A nationwide field experiment on facebook,” 2024 (WZB Discussion Paper Discussion Paper SP II 2020–302r2 (2nd revision February 2024)).
- Eger, Steffen, Yong Cao, Jennifer D’Souza, Andreas Geiger, Christian Greisinger, Stephanie Gross, Yufang Hou, Brigitte Krenn, Anne Lauscher, Yizhi Li, Chenghua Lin, Nafise Sadat Moosavi, Wei Zhao, and Tristan Miller, “[Transforming science with large language models: A survey on AI-assisted scientific discovery, experimentation, content generation, and evaluation](#),” *arXiv preprint arXiv:2505.05151*, (2025).
- Green, Seth Ariel, Benny Smith, and Maya Mathur, “[Meaningfully reducing consumption of meat and animal products is an unsolved problem: A meta-analysis](#),” *Appetite*, (2025), 108233.
- Luo, Ziming, Zonglin Yang, Zexin Xu, Wei Yang, and Xinya Du, “[LLM4SR: A survey on large language models for scientific research](#),” *arXiv preprint arXiv:2501.04306*, (2025).
- Pataranutaporn, Pat, Nattavudh Powdthavee, Chayapatr Achiwaranguprok, and Pattie Maes, “[Can AI solve the peer review crisis? A large scale cross model experiment of LLMs’ performance and biases in evaluating over 1000 economics papers](#),” 2025.
- Peterman, Amber, Jingying Wang, Kevin Kamto Sonke, and Janina Steinert, “Social safety nets, women’s economic achievements and agency: A systematic review and meta-analysis,” *Center for Global Development Working Papers*, (2024).
- Williams, Brooke A, Hawthorne L Beyer, Matthew E Fagan, Robin L Chazdon, Marina Schmoeller, Starry Sprenkle-Hyppolite, Bronson W Griscom, James EM Watson, Anazélia M Tedesco, Mariano Gonzalez-Roglich, and others, “[Global potential for natural regeneration in deforested tropical regions](#),” *Nature*, 636 (2024), 131–137.
- Zhang, Tianmai M, and Neil F Abernethy, “[Reviewing scientific papers for critical problems with reasoning LLMs: Baseline approaches and automatic evaluation](#),” *arXiv preprint arXiv:2505.23824*, (2025).
- Zhang, Yaohui, Haijing Zhang, Wenlong Ji, Tianyu Hua, Nick Haber, Hancheng Cao, and Weixin Liang, “[From replication to redesign: Exploring pairwise comparisons for LLM-based peer review](#),” *arXiv preprint arXiv:2506.11343*, (2025).

A LLM evaluation summaries and reasoning traces

A.1 Adena and Hager 2024

i Model assessment summary

This is a large, well-conceived geo-randomized field experiment assessing whether Facebook video ads increase charitable giving. The intervention randomized 7,686 of roughly 8,181 German postal codes into control or treatment over a 14-day campaign, with a 2x2 factorial on ad content (empathy vs effectiveness) and impression allocation (Facebook free allocation vs fixed per-postal-code budgets). Outcomes were measured at the postal-code-by-day level and crucially cover Save the Children's total donations across all channels, not just click-throughs, with 31 days pre, 14 during, and 38 post (and some robustness using 10 days of January). The main results show statistically meaningful increases in donation frequency both during and long after the campaign; revenue effects are positive and borderline-significant over the combined during+post window, and not significant in the short term. The paper also provides evidence of crowd-out for other charities (negative effects in competitor datasets), documents spillovers to nearby untreated postal codes, and shows that short-run click/attention metrics favor the empathy video but do not translate into superior long-run giving. The design and execution are careful and impressively scaled, but several limitations and interpretive cautions remain.

Key identification and measurement issues: - SUTVA/spillovers: The authors detect positive spillovers from treated to untreated areas (within 30 km), implying the main ITT estimates are lower bounds for total effects on giving to the focal charity. This strengthens the claim that ads raise giving to the focal charity but complicates the interpretation of magnitudes, especially for the decomposition and competitor analyses. The spillover specification is cross-sectional and may still pick up spatial confounding; it would benefit from pre-registered spatial models and placebo checks. - Partial exposure and dose: Only about one in ten Facebook users in treated postal codes received at least one impression; intensity varies across postal codes and treatments, and click data are semi-aggregated for the free-allocation arm. The study reports ITT effects but provides limited exploitation of variation in impressions to recover dose-response or treatment-on-the-treated effects. Given available impression counts, an IV or continuous treatment analysis could have triangulated effect sizes and validated the ITT magnitudes. - Aggregation and winsorization: Outcomes are aggregated to PLZ-day, with winsorization at €1,000 for PLZ-day sums. Heavy-tailed gifts can meaningfully affect revenue estimates; while the authors report similar results without winsorization, the paper relies on the winsorized specification for precision. Because most outcome analyses aggregate to period totals and normalize per million inhabitants per day, a small number of large donations could

still influence estimates; supplemental medians, quantile, or count models would help demonstrate robustness of distributional claims (beyond the helpful histogram-based evidence on donation-size categories). - Pretreatment differences and modeling choices: Pretreatment imbalances are non-significant for Save the Children outcomes but significant for the 23-charity alliance data. The preferred specifications include lagged outcomes, fixed effects (where applicable), and DiD. This is appropriate, but given multiple datasets and outcomes, a consistent DiD framework (with pre-trends visualizations) would strengthen comparability and reduce ambiguity about model-dependent results. - Competitor effects: Evidence of crowd-out relies on online-only donations for the 23-charity alliance and project categories on a donation platform for the second dataset. The former omits offline giving (unlike the focal charity data), the latter uses project tags (children-related vs other) with potential misclassification and multiple tagging. Several key coefficients are only marginally significant ($p < 0.10$), and there are many tests without multiple-comparison adjustments. The paper states crowd-out cautiously, but its magnitude and scope remain uncertain. Concurrent campaigns by other organizations or auction-market interference on Facebook could be part of the mechanism; these are not directly disentangled from donor-budget substitution. - Profitability estimation: The immediate revenue-to-ad-spend ratio is 1.45 based on point estimates, but the 90% CI is wide and includes substantially lower values. The long-run ROI assumes a lifetime multiplier of 1.75 for new donors and similar persistence for existing donors; this is plausible but not estimated from the experiment. The cost side excludes potential fixed costs (creative, management) and processing fees. The profitability conclusion should be framed as “probable with wide uncertainty,” rather than definitive, and sensitivity analyses for alternative LTVs and cost structures would improve credibility. - Pre-registration deviations: Blocking also used Facebook reach (not pre-specified), and the impression-allocation dimension was added. Excluding the top 1% and bottom 5% of Facebook reach also deviates from a pure census approach. These changes are well motivated (budget feasibility, power, and algorithmic allocation questions) and transparent, but they reduce the strength of strict pre-commitment. - External validity: The experiment ran during the peak German giving season (mid-November through year-end); effects may differ outside this period or in other countries/platforms. The campaign was intentionally largely untargeted. While this increases generalizability to the population of Facebook users in Germany, it may understate effects achievable with standard fundraising targeting. The heterogeneity analysis shows substantial scope for targeting, supporting this point. - Intermediate metrics vs outcomes: The paper’s conclusion that clicks and short-run attention are poor proxies for long-run donations is consistent with the presented comparisons, but the inability to trace individual click-to-donation conversions (due to tracking limitations) weakens causal linkage claims between specific intermediate metrics and outcomes. The empirical message remains strong but should be characterized as “indicative” rather than “conclusive” regarding mechanism. - Statistical inference and presentation: Most core findings are robust and supported by randomization inference. However, many secondary analyses (competitors, heterogeneity splits, factorial arms) involve multiple hypothesis tests; the paper could report adjusted p-values or emphasize effect sizes with confidence intervals to guard against false positives. - Reproducibility and open science: The study uses proprietary donation and platform data aggregated at PLZ-day and semi-aggregated ad metrics. The appendix is helpful, and the design was preregistered. However, the underlying data and code do not appear to be publicly available, limiting computational reproducibility and reanalysis. Clearer replication materials (synthetic data, code, detailed variable construction) would

substantially raise the paper's value as a building block.

Overall, this is a credible, policy-relevant, and unusually well-powered field experiment addressing an important question. The strongest contributions are: (i) measuring total donations across all channels, mitigating common substitution and measurement biases; (ii) documenting persistent effects beyond the campaign; (iii) providing evidence of inter-charity crowd-out; and (iv) demonstrating that click/attention metrics can mislead campaign evaluation. The main reservations relate to the uncertainty around revenue magnitudes and ROI, the reliance on assumptions for lifetime value, the limits of competitor data and multiple testing, and the inability to estimate dose-response. These caveats are acknowledged and do not undermine the central conclusion that online video fundraising, even when broadly targeted, increases charitable giving to the focal charity. The paper is valuable for both academics and practitioners and is near the top of recent work in this space.

Model reasoning trace

Assessing paper content and structure

I need to create a JSON assessment following the schema for “paper_assessment_with_tiers_v2.” The document in question is a WZB discussion paper titled “Does online fundraising increase charitable giving? A nationwide field experiment on Facebook” by Adena & Hager. The assessment should cover around 1000 words based only on the PDF. I'll identify key issues like identification threats, data limitations, and inconsistencies. There are conflicting instructions on percentiles and credibility intervals, which is puzzling! I should clarify that.

Preparing structured assessment

This is indeed an odd constraint! I want to ensure I stick to the end-of-instructions requirement to avoid any confusion with earlier guidance. I'll provide a single number for both percentiles and tiers while still selecting plausible values. For the applied area, I think it should focus on “applied and policy stream” since this is an economics field experiment on online fundraising. The reference group will include serious research in marketing and charitable giving. I'll ensure the JSON is precisely structured without including citations or external info, just focusing on internal details. It sounds like I need to be thorough in my diagnostic summary, covering around 1000 words!

Summarizing main contributions

The study employs a nationwide geo-randomized field experiment across 7,686 postal codes in Germany, featuring a 14-day Facebook video ad campaign for Save the Children. There's a control group of postal codes that didn't receive ads, utilizing randomization based on donation potential predictions. It employs a 2x2 factorial design on video content (empathy vs. effectiveness) and impression allocation without individual targeting. The outcomes measured the donation revenue and frequency per postal code daily, gathered from Save the Children's data over several months. This analysis provides compelling insights!

Analyzing robustness and results

January 10 data is included for robustness checks. The study uses complementary data from competitor charities, aggregating information from a group of 23 charities focused on online giving, along with donations from the betterplace.org platform for children's projects versus others.

The main findings indicate that Save the Children sees a significant revenue increase of €15-18

per million inhabitants per day long-term and an increase of 0.22 donations per million daily. There are potential negative effects on competitor charity revenue, suggesting crowding out. Click metrics alone may not accurately reflect long-term donation performance.

Examining negative interactions and effects

On the betterplace platform, there's a negative short-term impact on revenues for children-related causes, which suggests a substitution effect towards similar causes. The effects appear to be stronger in areas with higher potential, urban settings, more employment, and a larger Catholic population, among other factors.

Facebook's free allocation seems slightly better compared to fixed budgets, and there's a noted short-term difference between empathy and effectiveness content in donation frequency. Robustness checks are thorough, but there are concerns about data privacy and the accessibility of the provided information, likely due to confidentiality issues.

Analyzing units and treatment assignments

The analysis units are aggregated at the PLZ-day level, which could be too coarse. Winsorization at €1000 might impact long-tail donations, but they claim it yields similar estimates. Treatment is assigned at the PLZ level, but exposure to ads only occurs among Facebook users in those treated areas, which leads to partial exposure. There are challenges measuring treatment intensity across PLZs, potentially affecting dose-response identification. The randomization blocking method considers estimated potential and reach, but there seems to be a lack of preregistration for these decisions.

Discussing pre-registration and validity issues

The study mentions changes from pre-registration, but I wonder if these materially compromise the findings—it seems unlikely. They do address non-significant pre-treatment differences using lagged dependent variables and DiD, although they find significant differences in competitor data. While randomization minimizes selection bias, some pre-treatment variance might just be noise. The potential spillovers from treated neighbors to controls raise concerns about selection contamination and challenge the SUTVA assumption, complicating the analysis of competitor donations. Lastly, the results are largely untargeted, which limits generalizability across different contexts and regions.

Evaluating ROI estimates and identification threats

The ROI estimates are based on a lifetime value assumption of 1.75, which appears to come from retention rates, possibly from Save the Children donors or general trends. They assume that the impact of repeat donors equals that of new donors, but the retention numbers used seem coarse. There's no clear inclusion of overhead or credit card fees, and they only account for ad spend costs, not creative content or staff expenses. They state direct campaign costs at €33,700 but may overlook additional fixed costs. I'm curious about potential confounding factors affecting identification.

Considering measurement challenges and competitor effects

The assignment is done at the PLZ level, but competitor charities could be running concurrent campaigns that may target different audiences. Randomization should average these effects, but if competitor charities adjust in response to the Save the Children campaign—like interference in ad auctions or algorithmic shifts—this could impact treated versus control PLZs and lead to issues beyond simple donor substitution.

Regarding measurements, they aggregate Save the Children's donation data by PLZ-day, yet there may be unaccounted donations without proper PLZs. While they claim to capture data across all channels, the inability to link immediate donations to clicks due to tracking failures

might limit the analysis. They might be relying on overall totals, which could be reasonable.

Analyzing data and measurement issues

The data on competitor charities comes only from online giving, while Save the Children includes both online and offline contributions. This could mean that some observed substitution effects might reflect shifts in giving channels. While they measure substitution between online charities using platforms like betterplace, measurement errors could impact the estimated effects.

They're testing multiple hypotheses with various layers of heterogeneity but don't adjust for multiple comparisons, risking Type I errors at the 0.1 p-value level. The regression specifications seem somewhat limited since they use OLS with robust standard errors, but I wonder if the aggregated data presents issues like zero-inflated counts or serial correlation over time. Their winsorization at €1000 per PLZ-day might lower variance but could also downwardly bias the mean if there's a heavy tail, although they claim the results remain similar without winsorization.

Evaluating methodology and controls

They included pre-treatment donations as a control, but I wonder if that could lead to a "bad control"? Luckily, it's not influenced by the treatment, which is good. They decided to exclude extremes of the Facebook reach variable, which might compromise their "nationwide" claim, although they still cover 94% of postal codes.

When interpreting their findings, a positive and significant spillover coefficient implies that treated neighbors contribute more, which could be correlated with unobserved factors like urbanization. However, this could also be spurious.

There's noncompliance too, as only 10% of Facebook users in treated areas saw the ads, which likely dilutes the ITT effect. Balance tests appear limited, as they only show Table A1. They changed some design aspects post-registration, noting adjustments like impression allocation and block reaching, which they should ideally discuss further. Lastly, I find it interesting they used Stata and should provide code for reproducibility.

Analyzing data transparency and competitor effects

It's not entirely clear how they handled data, as they used third-party data and couldn't share it. However, they do provide OSF pre-registration and an appendix, which suggests data confidentiality might pose a moderate barrier to reproduction.

When looking at competitor effects, they found a negative impact of -€59 long-term per million/day based on differences, which relies on controlling for pre-treatment outcomes. Their use of six-block randomization, with only a 33% control share, is balanced, and they even provide maps of treatments, which is great.

They compare "treatment" versus "control" across the same periods to ensure balance, and their dataset includes 13,269 donations from the EU aggregator—so now, I'll need to produce the percentiles.

Evaluating study quality and impact

The design of the study is strong, resembling a near-nationwide field experiment with robust findings. The results are clear, demonstrating a positive effect and some substitution, which could be relevant for practitioners.

I'm considering an overall rating—maybe around 88? That sounds good, but is it top-of-field? The ROI feels partly uncertain since effect sizes are moderate. A p-value less than 0.1 for revenue looks promising, although the competitor substitution p-values are often around 0.1, which is just okay.

Rating the study's quality

I think I'd give an overall rating of 85.

For claims and evidence, they carefully characterize things and avoid overclaiming, but there are some leaps in ROI. I'd rate that an 82.

Their methods are solid, with good randomization and blocking, although they chose to win-size, which raises potential SUTVA issues. I'd give that an 84.

The advancement of knowledge is good as well, landing at 80. Communication is clear, so I'd rate that an 88.

However, for open science, the data and code aren't available, so that gets a 55.

Finally, in terms of global relevance, I'd rate it 75. I think it has the potential for a top-tier journal like the Journal of Public Economics, so I'd aim for around 4.2, which indicates "top field."

Evaluating the study's journal potential

I see the study is significant but maybe not quite A-journal quality. It could fit into journals like JPubE or Management Science, possibly even ranking in the top field. I'd rate it a 4.0 for potential.

I think "will" might be around 3.6. However, the JSON requirements state that scores need equality across tiers. So, I'll align the "should" to 4.0 and "will" to 3.7, both requiring equal values.

For my overall ratings, I'm considering:

- Overall: 86
- Claims evidence: 84
- Methods: 85
- Advancing knowledge: 78
- Logic communication: 90
- Open science: 60
- Global relevance: 78

I need to keep in mind that we promised quantitative metrics based on my insights, considering calibrations over the last three years.

Summarizing assessment issues

I want to address several key issues in the assessment summary. First, there's a SUTVA violation and potential spillovers, which could underestimate the treatment-control difference. I see some pre-treatment differences in the competitor dataset, which they adjust with DiD methods, but it still raises concerns. There are uncertainties from aggregator-level measurements and fluctuations in donation levels that lead to borderline significance for revenue, despite their use of a large sample.

They use a "lagged dependent variable" effectively alongside randomization. However, the ROI calculation using a lifetime value of 1.75 seems overly optimistic, ignoring creative and staffing costs.

The intermediate metrics do not effectively link to outcomes due to failed tracking. They also missed opportunities to leverage impression counts at the PLZ level for estimating treatment, which hinders intensity estimation.

Excluding 6% of PLZ may introduce bias, although they offer justifications. Some Save the Children donation data could be misassigned, affecting the results, yet they mention the potential for randomness in misassignment.

For competitor measurements on betterplace, potential classification errors could impact results, and significant negative interactions appear to be short-term. The analysis includes postal-code fixed effects for improved accuracy, which I think is a positive step.

The “share of treated neighbors within 30 km” measure might correlate with unobserved factors, but they do use block fixed effects—a good choice. Their frequency and revenue normalization omit baseline differences, yet they make necessary adjustments.

While the 2x2 design isn’t strongly powered, they acknowledge possible underpowering, which is thoughtful. They’ve omitted recipient surveys, but perhaps that’s acceptable.

They used an “untargeted campaign” approach, limiting generalization to typical charity practices, though they do test heterogeneity.

The sample period is narrow, preceding the end-of-year giving surge, suggesting the effects may vary by month. Their reported profitability ratio of 1.45 with a wide confidence interval might indicate possible unprofitability, which they acknowledge.

They discuss how the Facebook algorithm favors cheaper impressions, yielding improved metrics without significant differences noted, and recommend caution there. Investigating channel substitution by donation source shows no negative impacts—a positive finding.

Lastly, they use Stata for randomization inference and presumably share code, which could help with reproducibility.

Assessing the research design

In this chain of thoughts, I’m highlighting some key points about the study’s design and strengths. It’s clear that pre-registration is beneficial, even though some data from the Alliance and betterplace might be incomplete. They use this data for measuring substitution, suggesting there’s some evidence backing their approach.

It’s worth noting that the winsorization threshold of €1000 per PLZ-day seems arbitrary, but considering the aggregated donations, it may provide relevant insights. Standard errors are robust, and while they don’t cluster across neighbor interactions, this isn’t a significant concern.

I see an interesting synergy here: the idea of spillovers to neighbors implies that the treatment effect (TOT) in treated areas includes contributions from unaffected donors in control areas, making the TOT appear larger.

They suggest that further research could investigate the “impression-of-ask” for future insights. Now, I need to craft a diagnostic summary that pinpoint specific issues while being neutral and clear. I’ll start by outlining the design and contributions, followed by the study’s key strengths, which include geo-randomization, extensive coverage, and robust competitor data. Then, I’ll identify potential threats to validity.

Analyzing revenue effects and data limitations

I’m noting that the study shows a weak total effect on revenue, as the p-value is only less than 0.10. The ROI estimate seems uncertain with a wide confidence interval, and the results only reflect immediate income while being heavily winsorized, which could lead to frequent zeros. This makes it hard to make strong claims about profitability. There are also SUTVA/spillover issues and potential algorithmic mis-targeting.

The presence of a positive neighbor treatment effect suggests the main intention-to-treat (ITT) is a lower bound, complicating inferences about competitor substitutions and direct effects.

Additionally, only about 10% of Facebook users received an ad impression, indicating small ad intensity while treating this as ITT. The absence of a treatment-on-the-treated (TOT) effect or dose-response metric means that patterns observed between “free vs. fixed” and

“empathy vs. effectiveness” rely on limited power and aggregated metrics. Also, the Save the Children dataset has limitations, with potential misassignments concerning post codes and measurement errors due to its aggregated nature.

Exploring competitor analysis and profitability calculations

I’m focusing on the competitor analysis, where issues arise with the alliance-only online classification and pre-treatment imbalances related to betterplace’s ‘children projects.’ The negative results might only be significant at the 10% level, and there’s a lack of multiple testing corrections. Additionally, unobserved competitor campaigns may distort substitution effects, making them reflect ad auction displacement rather than real donor budgets. The instrument employed doesn’t seem strong.

Moving to profitability calculations, the analysis relies on the assumed lifetime value of 1.75, but replicating this would necessitate detailed code and a full breakdown of costs, which isn’t robust against various alternatives.

Notably, there was a deviation from pre-registration concerning impression allocation treatment, as well as additional blocking to maximize external validity. However, this deviation reduces their ability to make strong prior claims about treatment effects. Finally, reproducibility is limited due to non-public data and proprietary metrics, which makes it difficult for third parties to replicate the results.

Assessing external validity and overall metrics

I’m noting that the study’s external validity is limited to Germany and the end-of-year timeframe. They mention that more donors might be active during this period, which could affect generalization to other times.

Regarding the explanation of using click data as an inappropriate proxy, it seems plausible due to the failed tracking of clicks on the Save the Children site. While it’s a reasonable claim, it remains inconclusive.

Now, about the overall midpoint metric, I’m thinking maybe 84? We previously considered 85, but after recalibrating against the past three years of significant research, I think I’ll actually settle on a higher score of 88 instead of the earlier 86.

Finalizing numerical scores and tier classifications

I’ve settled on some final scores: for overall, I’m picking 87. Looking at the details, I’ll assign claims and evidence a score of 83, methods 85, advancing knowledge 79, logic communication 90, open science 55, and global relevance 78.

For the tier classifications, I’ll choose “should” as 4.1, marking it between a top field journal and a marginal A-journal. However, I’ll adjust the “will” score to 3.8.

To keep things balanced, I’ll ensure the confidence intervals are equal, setting both `ci_lower` and `ci_upper` to 4.1 for the “should” score. Now, we’ll get started on crafting that 900-1100 word assessment summary!

Crafting the assessment summary

I need to summarize the study clearly without referencing external citations or authors, just focusing on “the paper.” So, I’ll start with a sentence highlighting that it’s a georandomized study with 7,686 PLZ, a 14-day campaign, and a 2x2 factorial design. The dataset covers Save the Children’s full donations across channels, plus competitor data from alliance and betterplace, aggregated per PLZ.

Next, I’ll systematically address the issues: looking at primary effect sizes, significance, pre-treatment differences, and resolution. I’ll also cover SUTVA and spillovers, the implications for treatment effects, and the limited exposure intensity where only every tenth Facebook user

saw an impression.

Evaluating data and methodology

I've noticed that the study didn't measure the treatment on the treated (TOT) across PLZ, and it didn't utilize regression for TOT. The data is aggregated with winsorization, showing a heavy-tailed distribution, making it hard to attribute results to individual donors. The competitor data only includes online interactions related to children's projects, which appears weak. There's also a risk of Type I error due to multiple testing.

Regarding profitability, the lifetime multiplier relies on retention rates but excludes overhead costs and creative expenses, leading to a wide confidence interval. There are intangible benefits that aren't fully captured.

The study deviated from its pre-registration by mixing free and fixed strategies and using reach for blocking, which raises concerns about analysis-plan drift. External validity is limited to the time of year studied; effects might not apply in other months or during disasters. Their robustness includes randomization inference and winsorization, yet their analytical choices, like using OLS weighting and normalizing per million inhabitants, could obscure heterogeneity. Replicability is somewhat compromised since the dataset isn't shared and the code is unavailable, although it offers semi-aggregated click data. Lastly, there are limitations in relying on Facebook's estimated reach and advertising manager exposures, and they could conduct additional tests for TOT-of-intensity and cross-lags.

Creating JSON structure

So, the instructions tell me that the summary should come before any scoring and guide uncertainty. I need to create a JSON object that strictly adheres to the given schema without any extra keys involved. It's a bit odd that the rules state the percentiles should be in [0, 100] with the lower_bound, midpoint, and upper_bound all equal but okay, I'll follow that.

For the metrics, I'll set the overall midpoint at 87, with both lower and upper bounds also at 87.

Then, I'll do the same for the other percentiles while keeping in mind the tiers from 0 to 5 and ensuring they're all numbers, not strings. I'll set "should" at 4.1 and "will" at 3.8.

Next, I'm crafting the summary text, which should be about 1000 words, aiming for 900-1100. I'll include suggestions like multi-level modeling, cost-effectiveness, and unifying data. Finally, I need to double-check that the JSON schema matches correctly and finalize everything.

ewpage

A.2 Green et al. 2025

i Model assessment summary

This manuscript presents a focused theory-informed systematic review and meta-analysis of randomized controlled trials (RCTs) that attempt to reduce consumption of meat and animal products (MAP), with the key contribution being a strict inclusion of behaviorally measured outcomes assessed at least one day post-intervention and minimum sample size/cluster criteria. The authors code 112 interventions from 41 studies (~87,000 participants) and use robust variance estimation (RVE) to account for within-study dependence. The main quantitative result—a small pooled standardized mean difference (SMD ~ 0.07, 95% CI [0.02, 0.12]) for net

MAP consumption—with publication-bias-adjusted estimates near zero, contrasts with more optimistic prior syntheses that mixed intentions/attitudes and non-randomized designs. They also distinguish interventions targeting red and processed meat (RPM), which show larger effects (SMD ~ 0.25), and argue this likely overstates net MAP reduction because of plausible substitution to poultry/fish. The paper is timely, policy-relevant, and generally careful in interpretation. That said, several methodological and reporting choices merit clarification or revision.

Study identification and inclusion: The search is pragmatic and multi-pronged (reviews, citation chasing, crowdsourcing, targeted search terms, an AI tool, and nonprofit databases). However, the restriction to English may introduce language bias. It would help to quantify the contribution of gray literature among included RCTs and to document screening inter-rater reliability; as written, data extraction and inclusion decisions were primarily performed by one author, which raises error risk and subjective bias. The PRISMA diagram is provided, but a fuller accounting of reasons for exclusion at full-text stage would improve transparency. The requirement of a pure control is defensible but excludes some common cafeteria-level interventions where “usual practice” or partial exposure forms the counterfactual; the authors partly address this via a “marginal studies” sensitivity, but the implications for external validity could be emphasized.

Outcome selection and measurement validity: Restricting to behavioral outcomes measured =1 day post-treatment is a major strength. Yet, many included outcomes appear coarse (e.g., binary self-reports of reduction vs. no reduction) and self-reported, which are vulnerable to social desirability and recall bias, especially when messaging primes ethical or identity frames. The authors code self-report vs. objective assessment but do not stratify or adjust by this risk-of-bias dimension; reporting subgroup results (and whether effects are attenuated with objective measures) would be useful. The decision to use the longest follow-up per intervention is consistent with the paper’s emphasis on sustained change, but it may mix widely varying exposure “doses” and contexts. A sensitivity using the earliest non-immediate behavioral outcome could show whether effects decay.

Effect-size computation and harmonization: The paper uses Glass’s $\langle U+0394 \rangle$ (treatment-control mean divided by control SD), reportedly the “pre-treatment control group” SD when available; this needs clarification. In most included RCTs there is no natural pre-treatment SD on exactly the same outcome measure used at follow-up, and using a baseline SD (if measured) may be inappropriate if outcomes change scale distributions or if the analysis uses post-only means. If the intended denominator is the post-treatment control SD (the more typical Glass’s $\langle U+03B7 \rangle^2/z$), the text should say so unambiguously. Where means and SDs are missing, conversions from coefficients/ $\langle U+03B7 \rangle^2/z$ are standard; however, the imputation rule for “unspecified nulls” (setting the effect size to 0.01) is problematic. This assigns a positive effect where direction is unknown and could upwardly bias the pooled mean precisely in the range of small effects under study. A more neutral imputation (e.g., 0 with imputed variance based on study size) or exclusion with sensitivity analysis is preferable. The difference-in-proportions estimator for binary outcomes is reasonable, but details on aligning signs (so that positive values uniformly indicate reductions in MAP) should be explicit.

Dependence, clustering, and weighting: The main meta-analytic model uses RVE via robumeta with small-sample correction, appropriate for multiple interventions per study sharing a control. However, the paper does not report the assumed within-study correlation ($\langle U+03C1 \rangle$) or show sensitivity to alternative $\langle U+03C1 \rangle$ values; robumeta’s weights can change markedly

with $<U+03C1>$. More importantly, several studies are cluster-randomized. It is not explicit whether the effect-size variances were adjusted for clustering (design effects/ICC or cluster-level summaries). If not, those studies' precision would be overstated and they would be overweighted. The minimum of 10 clusters helps but does not resolve variance mis-specification. The authors should confirm and, if needed, re-estimate ν_i using cluster-corrected SEs or design effects and demonstrate robustness.

Publication bias and sensitivity: The use of selection models (favoring significant positive results) and the "worst-case" non-affirmative-only sensitivity are commendable, and the bias-adjusted pooled effects near zero strengthen the headline conclusion. Yet selection models and significance-funnel approaches in metafor assume independence of effect sizes; since interventions within a study share a control, dependency remains unless the analysis is aggregated to one effect per study or a multilevel selection approach is used. The authors should clarify whether selection analyses were performed on study-aggregated effects and, if not, provide an aggregated sensitivity. The significance-funnel plot is informative; adding a p-curve or caliper test would further triangulate.

Moderator and subset analyses: The categorization into four theories and the RPM vs. MAP split are conceptually valuable. Because interventions and messages often overlap across categories and a single study can contribute to multiple classifications, the authors prudently avoid causal meta-regressions on theory. Still, they report subgroup pooled estimates with very small samples in some cells (e.g., choice architecture: 3 estimates, CI extremely wide). This should be presented more cautiously (e.g., descriptive only) and, where possible, collapsed across similar mechanisms. The temporal decline in average effects is interesting; adding cumulative meta-analysis or meta-regressing effect size on publication year (study-level, one effect per study) would better visualize potential decline effects.

Scope and inference: The paper's main claim—meaningful, sustained reductions in net MAP consumption via current behavioral interventions remain largely elusive—is supported by the evidence assembled under strict criteria, and the discussion acknowledges limitations, substitution, and external validity. The contrast with much larger estimates when relaxing inclusion criteria is persuasive. However, the normative framing "unsolved problem" should be bounded by the sample: RCTs with behavior assessed = 1 day, minimum size, pure control, in mostly high-income settings. The scarcity of rigorous choice architecture RCTs with delayed outcomes is itself a finding but also a constraint on generalization. The paper is careful not to overclaim superiority of any theory, and it highlights promising under-tested interventions and measurement approaches.

Transparency and reproducibility: The authors emphasize computational reproducibility (Rmarkdown, containerization) and provide a DOI to the code/data repository. That is exemplary. The main weaknesses are single-coder data extraction, limited reporting on coding reliability and adjudication procedures, and lack of a preregistered protocol. Adding an archived protocol (even post hoc), inter-rater checks on a subsample, and a full data/codebook mapping each effect-size computation would further strengthen replicability.

Overall, this is a rigorous, policy-relevant synthesis with conservative conclusions, appropriate caution about moderators, and strong open-science practices. Addressing the flagged issues—especially the "unspecified null = 0.01" imputation, clarity on Glass's $<U+0394>$ denominator, cluster adjustment in variances/weights, independence in publication-bias analyses, and additional risk-of-bias stratifications—would substantially increase confidence. Even with these caveats, the paper adds valuable clarity to a crowded literature and provides a re-

alistic baseline for future trials and policy design.

Model reasoning trace

Evaluating JSON structure for metrics

I need to create a JSON object with two fields: an assessment summary and a metrics object. The metrics object must include fields for each percentile metric, detailing the midpoint, lower bound, and upper bound. However, the instructions say that these values must be equal. This is odd since it doesn't provide a proper range. I have to follow the instructions closely, even if it feels unusual. I also need to fill specific categories for percentile metrics and include two tier metrics.

Crafting assessment from PDF content

I need to base the assessment strictly on the content from the PDF the user provided, which seems to include the entire paper and its references. Since the authors' identities shouldn't influence our judgment, I'll avoid mentioning them in the summary. My goal is to write about 1000 words of a diagnostic summary that encapsulates the key points from the text. I want to make sure I'm focusing on the content itself without biases related to the authors.

Analyzing assessment aspects

I need to evaluate various components such as identification threats, data extraction methods, and SMD calculation approaches, including the use of Glass's delta and RVE. I also have to consider publication biases and restrictions on RCTs, like group sizes and measurement timeframes. Notably, they've set the "unspecified null" to 0.01 instead of 0, which could introduce bias by implying a small positive effect. This seems problematic, and I should definitely flag it for further attention.

Evaluating measurement concerns

I need to assess the use of Glass's <U+0394>, which considers pre-treatment control group standard deviation. In many RCTs, there aren't pre-treatment measures available, and they seem to refer to the baseline. This could lead to ambiguity since they may have relied on pooled SD. There's also the consideration of positive sign alignments with decreases in MAP that I need to verify. Additionally, the use of RVE with small-sample corrections raises concerns, especially with 41 studies and 112 estimates. They might have missed immediate effects by focusing on longer follow-ups, which could introduce biases.

Analyzing assumptions and biases

I'm considering the assumption that they've restricted to "net MAP or RPM consumption." They mentioned some outcomes are quite coarse, like binary reduction versus no reduction, and they included self-reports without adjusting for risk-of-bias variables. They didn't present the Egger test either, which could be a concern. Their use of selection modeling assumes a $p < 0.05$ threshold, and there are potential issues with correlated sampling errors. Although they've applied RVE to address dependence, there might be critiques regarding how they handle these correlations in their selection modeling.

Addressing multi-analytic classification

I see a potential issue with their multi-analytic classification of categories. They haven't performed meta-regression but only provided pooled SMD within categories. This approach might lead to dependencies, especially since they have multiple effect sizes per study. Although they used RVE for all 112 effect sizes, their subset analyses lack meta-regression control. Their reasoning for not conducting meta-regression—due to studies with multiple interventions

crossing categories—makes sense to avoid double counting, but calculating pooled effects per category still raises concerns about how multiple contributions are accounted for across different subsets.

Questioning unspecified null and clustering issues

I'm noticing that coding the "unspecified null" as 0.01 might be problematic. With 112 effect sizes, they set the unspecified null to 0.01 instead of 0 when they lack enough information to compute SMD, which could introduce a slight upward bias. They might have chosen 0.01 to avoid exactly zero for log transformations, but that doesn't seem necessary for SMD. Additionally, their use of the Bernoulli difference in proportion estimator might not properly account for cluster-level correlations. I need to check if they adjusted for clustering in their analyses since they mentioned requiring at least 10 clusters.

Evaluating adjustment methods and standard deviations

I noticed they mentioned "first author conversions using standard techniques," but they didn't seem to adjust for cluster design adequately, possibly due to limited information. The use of "pre-treatment control group SD" is questionable, especially since many RCTs measure consumption at follow-up. Referring to a baseline SD for SMD at follow-up feels inconsistent—I think it should really be a pretest measure. Using baseline variation might reduce bias, but if the scale changed over time, it could complicate things. Hence, I'll flag this as a possible misprint or a methodological concern.

Reviewing SMD and data accessibility

I'm noticing that they used sign-coded SMD to indicate a reduction in MAP consumption, but they didn't clarify the sign of the effect. They referenced "108 RCTs and 112 interventions," but their forest plot includes both positive and negative effects. They mentioned "87k subjects," which might include cluster-level counts, but I see they didn't adjust for compliance or attrition, and it's unclear how missing data was handled. While they provide their code and dataset along with a DOI for full documentation—which is great for open science—I need to check if the code and data are truly accessible for evaluation and replicability.

Evaluating data handling and biases

They mention that the data was extracted by the first author, which could lead to single-coder errors. While they consulted authors and used publicly available datasets, they didn't double-code or verify effect size calculations independently. This raises concern about potential biases. They included some dissertations and grey literature but restricted their analysis to English publications, which could introduce bias as well. Notably, they included marginal studies in their sensitivity analysis, showing that these can yield larger effect sizes, although they advise caution in interpreting these findings.

Assessing RVE and bias considerations

They employed "random intercept RVE robumeta," which includes a small-sample correction, but I'm noticing that they used a default 'rho' of 0.8 without specifying what weighting they applied. This is important, as variations in rho could change results, and they didn't vary their assumptions between 0 and 1. They also didn't conduct a robust meta-reg with cluster-level correlations. On the positive side, they correctly identified "worst-case publication bias" yielding a small effect, which is useful information for understanding their findings.

Reviewing inclusion criteria and considerations

They included "two choice architecture studies with only 3 estimates," which raises some questions. They speculated on reasons for this, mentioning a lack of long-term or hypothetical outcomes, but their threshold for inclusion seems to exclude interesting cases like "defaults at

buffets,” even if the control had a baseline. On a positive note, they did utilize the longest follow-up measure for each intervention, which is good. However, it’s curious that they didn’t discuss “nudging philanthropic disclaimers.”

Analyzing bias and intervention effects

One nuance here is that they used a “PublicationBias” function that only considers non-affirmative results for worst-case estimation. However, with Random Variance Estimation (RVE) not being independent, there might be some issues. They employed metafor’s generic function rather than RVE, which could disrupt the correlation. It seems they didn’t adjust for dependent effect sizes, but it’s acceptable since each of the 112 effect sizes corresponds to distinct interventions. Multiple interventions sharing control within the same study are treated as dependent, and for publication bias, they didn’t make any adjustments. Still, there isn’t a huge difference overall.

Evaluating effect size computation

They mention using a “difference in proportion estimator,” which treats discrete events as Bernoulli distribution draws—this seems acceptable. However, they didn’t specify the “units” for Standardized Mean Difference (SMD) across different measurement scales, though they addressed it with standardization. Another concern is that they included cluster RCT experiments, but I’m wondering if they computed effect sizes with cluster-level standard errors. If they didn’t, that could affect independence and variance. While they employed robust RVE, they didn’t clarify adjustments for the design effect, which is crucial. I’ll flag this issue.

Assessing sample size and weighting

They included approximately 87,000 participants, which is significant, but the ‘units-of-randomization’ might be less. It’s great they set a minimum of ten clusters, but they didn’t weight by the number of clusters. Instead, they might have weighted by ‘within-study variance’ assuming cluster-level data. If sample-level isn’t adjusted, the cluster studies could appear more precise than they really are. Also, they included diaries or receipts for outcomes, which is good. However, they noted the lack of a pre-analysis plan without incorporating that into their risk-of-bias weighting. The SMD of 0.21 to 1.42 seems off because it suggests only three estimates, indicating significant imprecision, which is worth mentioning.

Evaluating research score

They discuss four categories but didn’t specify doping, which is worth noting. It’s good they acknowledge the ‘lack of causality.’ Now, I need to produce numerical percentile scores for each metric. Overall, I’d probably place it around 75%, maybe 70%. The paper demonstrates rigor in its meta-methods and addresses relevant global priorities, but there are methodological concerns, such as single coding and unclear cluster adjustments.

Based on the strength of the open science aspects, like sharing code and datasets, I would assign an overall score of 80% and consider 90% for the open science component.

Assessing claims and evidence

I think the single-coder issue affects replicability, so maybe I’d score it around 85. They make claims about unsolved problems and provide balanced evidence, which is good. It’s also a positive point that they approach substitution issues carefully and restrict their analysis to RCTs with long-run outcomes. However, the mention of ‘worst-case bias’ raises a flag. Overall, the claims seem measured, so I’d settle on giving this section around 80%.

Evaluating various aspects

There are definitely some issues, like the unspecified null of 0.01 being a bit unnatural and concerns about independence in selection models, leading me to reduce the score to 76. For

the methods, they're good overall, but I'd rate them around 72 given the issues.

When it comes to advancing knowledge, they effectively challenge earlier positive claims, so I would give it a solid 78. The logic and communication are straightforward, with good visuals, so I'd score that an 85. Their commitment to open science deserves a 90, and the topic's global relevance is strong, earning an 88. Considering the cost/benefit discussion, I'd settle on around 80.

Determining journal rankings

I need to ensure the credible intervals are adjusted for clarity, where the lower bound equals the midpoint. Initially, the instructions called for including credible intervals, but now they indicate to set them equal, so I'll follow that.

Now, regarding the "should/will" scale from 0-5: for the "should" classification, I'd place it around 4.3 because it might fit into a top field journal like "PNAS" or "Environmental Behavior." For the "will" classification, I'll estimate it as 3.9 since many meta-analyses are featured in solid B-journals like "Food Policy." So, that feels balanced!

Preparing the analysis summary

The final instructions direct me to set the score equal to both the lower and upper credible intervals, so I'll clarify that point. Now, I'm ready to create the 1,000-word summary while maintaining a neutral tone and focusing on specific details.

I'll begin with the clarity of the question, noting that it's a meta-analysis of RCTs aimed at reducing MAP with long-term measures, highlighting the 41 studies and 112 interventions involved.

Key concerns include:

- 1) Study selection: It's limited to English only as of December 2023, considering published and grey literature but with a customized search strategy.
- 2) Outcome measurement: They shifted to net MAP vs. RPM and included self-reported measures, which may be coarse.
- 3) Data extraction: A single coder was used, raising the risk of measurement errors and unspecified null imputation at 0.01.
- 4) Effect-size conversions: They utilized Glass's $<U+0394>$ with the pre-treatment control group SD.
- 5) RVE modeling: The robumeta method was used for analyzing 112 effect sizes.

This outline will guide my summary!

Structuring the analysis summary

I need to continue refining the summary analysis, focusing on specific metrics. I have points like publication bias, where they used a selection model with worst-case "non-affirmative only" data. Also, I should consider heterogeneity and generalization, and recognize any dangers in category-level pooling. For follow-up times, there's a mention of the "longest follow-up" and the need to correct for cluster assignment.

Now, I'm preparing to write around 800 to 1,100 words while crafting the JSON for the assessment metrics. The "Paper_assessment_with_tiers_v2" schema should include nested objects for "tier_should" and "tier_will", both containing scores, lower bounds, and upper bounds.

Since the spec says to keep bounds equal, I'll make sure to set all percentiles properly, ensuring it's valid JSON without any trailing commas, and the property order isn't crucial. That will keep everything neat and accurate!

Selecting values for assessment

Alright, I'm choosing values to present for the analysis. I'm thinking an overall score might be around 82, which feels consistent with the summary. For the individual categories, I could go with claims and evidence at 78, methods at 74, advancing knowledge also at 80, and logic and communication at a solid 86. For open science, I think 92 is appropriate, and global relevance could sit at 85.

For tiers, I'll set 'tier_should' at 4.3 and 'tier_will' at 4.0. Now, I just need to keep focused while crafting the 1,000-word summary, maintaining clarity and robustness throughout!

ewpage

A.3 Kudymowa et al. 2023

i Model assessment summary

This commissioned review provides a clear, well-scoped examination of discounting choices used in philanthropic global health cost-effectiveness, with the immediate goal of advising a specific decision-maker. The manuscript's core contributions are: (i) retaining an SRTP framing while articulating why SOC and harmonizing approaches are less suitable in this context; (ii) clarifying the unit being discounted in GiveWell-style modeling (consumption growth rather than levels) and deriving an explicit, adapted wealth-effect formula $r = d + (\beta - 1)g$; (iii) identifying and correcting a minor computational error in the prior spreadsheet implementation; (iv) proposing a defensible, transparent procedure to estimate g using a population-weighted SSA+SA composite and a time-window anchored to program effect duration, with an update cadence; and (v) highlighting an important internal inconsistency: using $\beta = 1$ in the utility implicit in CEAs (log utility) while using $\beta = 1.59$ for the wealth-effect component, which can materially alter the overall rate. The report also reviews discounting for health outcomes, pure time preference, temporal uncertainty, and a "compounding non-monetary benefits" component. The overall structure is logical and the reasoning is mostly transparent, with useful appendices on derivations and institutional practice.

Key issues and limitations arise primarily from evidence depth, parameterization rigor, and potential internal inconsistencies that are flagged but not resolved. First, the choice to keep SRTP over SOC is broadly reasonable for cross-intervention welfare comparisons, but the paper does not quantify the decision stakes of alternative approaches (e.g., implementing an SOC-based sensitivity case or the SPC/weighted-average harmonization to show how rankings or absolute valuations would change for representative interventions). Given the centrality of this decision, readers would benefit from explicit comparative numerics under plausible parameterizations and program archetypes (consumption vs mortality vs morbidity outcomes) and explicit discussion of potential crowding-out magnitudes for philanthropic funds across feasible counterfactuals.

Second, the adaptation $r = d + (\beta - 1)g$ is a valuable clarification of what is implicitly being discounted in these CEAs; however, the implications are only partially explored. The report rightly notes that adopting a consistent β across the entire modeling stack is crucial to avoid overdiscounting, and that if $\beta = 1$ were used

consistently, the wealth effect falls to zero. Because this change alone can shift the overall rate by ~1.7–2.0 percentage points and potentially alter program rankings, the absence of a worked re-estimation of a few headline CEAs (or at least a stylized demonstration of ranking sensitivity) leaves a major decision lever underexplored. Relatedly, the suggestion to include either donor or beneficiary pure time preference, but not both, is correct; yet, no quantitative illustration is provided for plausible donor-weighting choices, which would illuminate how sensitive choices are to stakeholder perspective.

Third, the proposed method for estimating g (3.3%) is transparent but rests on several judgement calls that deserve more robustness checks. The SSA+SA population-weighted composite and 40-year window tied to deworming are reasonable, but the paper itself flags a puzzling result: the composite CAGR lying below that of both components over part of the horizon and tracking SA closely; this could arise from weighting, smoothing, or interpolation choices. The authors acknowledge limited time to resolve the anomaly, which weakens confidence. Additional triangulation would strengthen the recommendation: e.g., comparing alternative global forecasts (beyond IMF and Goldman Sachs), decomposing growth into productivity and demographics, checking sensitivity to funding-allocation weights versus population weights, and testing alternative horizon widths (20, 30, 50 years) and start/end-year anchoring. The proposed 5.2-year update cadence is pragmatic, but the threshold (0.1 pp change warranting an update) is ad hoc; a justification aligned with decision-value or threshold cost-effectiveness changes would be preferable.

Fourth, recommendations on non-consumption components are mixed in strength. For d , the report synthesizes philosophical and empirical arguments well and makes a defensible case for $d = 0$ in this application. For temporal uncertainty (1.4%), and especially for the “compounding non-monetary benefits” (0.9%), the evidential base is thin. The temporal uncertainty component is broadly aligned with other estimates but mixes existential/catastrophic risks with other pathways to benefit non-realization in a way that risks double counting with other parts of the CEA (e.g., baseline mortality, fade-out) unless carefully separated. The 0.9% non-monetary compounding term is acknowledged to be weakly justified and partly a historical compromise; no alternative modeling frameworks (e.g., explicit dynamic household investment pathways, or SPC-style reinvestment for the relevant subset of benefits) are trialed. Given these components together account for more than half of the overall rate in some parameterizations (especially if $\langle U \rangle$ were set to 1), their uncertainty deserves more structured analysis (e.g., priors, ranges, and decision-relevant sensitivity).

Fifth, the health-vs-consumption discounting section captures the state of the literature accurately: equal discounting remains common, but arguments for discounting health at a lower rate are stronger conceptually; nevertheless, the report does not specify a practicable health rate beyond continuing to apply temporal uncertainty only. This conservative, “status-quo-with-justification” recommendation is understandable, but it would be more decision-useful to propose at least a small set of candidate health rates with transparent assumptions (e.g., under additively separable utility versus specific non-separable forms), plus implications for representative interventions.

Sixth, while the manuscript is technically careful—the correction of a spreadsheet formula and explicit derivations are strong points—reproducibility for applied users depends on linked spreadsheets and external projections. The report references online resources but does not embed full data snapshots or code within the document; version control, a frozen dataset appendix, and a worked example reproducing the 4.3% rate calculation end-to-end would

enhance replicability. Some internal consistency concerns remain: clear separation of units (utils vs consumption levels vs log consumption) is emphasized, but readers might still find it challenging to trace how each component maps to the unit actually discounted in different CEA modules. Finally, the institutional scan is informative but shallow by design (~10 hours), relying on convenience samples and email inquiries; conclusions drawn from that scan should be treated as suggestive rather than definitive.

Overall, the report advances clarity on a few pivotal modeling choices (unit of discounting, consistency, explicit wealth-effect formula, and a more principled g estimate and update process), provides a balanced and transparent discussion of contested components, and responsibly flags uncertainties. The main weaknesses are the limited depth of quantitative robustness, reliance on a few external forecasts without triangulation, unresolved anomalies in growth aggregation, and the absence of sensitivity demonstrations on headline program rankings for high-stakes parameter changes (notably γ and the inclusion/exclusion/magnitude of the 0.9% component). These limitations are typical for an applied advisory report on a tight timeline, but they constrain generalizability. With additional analyses that quantify decision-sensitivity and resolve the growth-aggregation puzzle, the work could serve as a stronger reference for the broader global health CEA community.

Model reasoning trace

Evaluating estimation methods

I'm exploring the "time window" method for estimating growth using IMF WEO and Goldman Sachs data, updated every five years. They anchor this with a 40-year window linked to deworming, resulting in an estimate of 3.3% growth. However, there are issues with composite CAGR and heavy reliance on Goldman Sachs's techniques. Their open science aspects include shared calculations but lack replicability due to missing links. Though their review process involved experts, it's not exhaustive, and the report serves more as a practical tool than a novel academic contribution.

Clarifying norms and descriptions

I need to ensure there's a clear distinction between normative and descriptive statements while providing credible disclaimers. This separation will help prevent any confusion about what's being presented as a standard versus what's merely an observation. It's essential to present information transparently, making it easier for users to understand the context and implications of the data and insights shared. Overall, I want to maintain clarity and credibility in my communication.

Identifying potential issues

I see several potential issues in the analysis. First, the argument for adopting SRTP over SOC is weakly formalized. There's uncertainty in the "health vs consumption discounting" claims, which lack strong guidance. The temporal uncertainty is broadly defined, and without first principles calculations, it's hard to evaluate. Additionally, the non-monetary benefits appear poorly grounded in theory. The consumption growth reasoning shows puzzling results that they can't explain. There's inconsistencies in utility curvature, but they do provide some positive aspects like referencing derivatives. Overall, the methodology has its merits but also significant gaps.

Assigning percentiles for evaluation

I need to assign percentiles for the evaluation. I'm thinking of giving an overall score of

70%. For claims and evidence, maybe 72%, and methods could be around 65% since they didn't conduct rigorous empirical analysis, but it's reasonable for a policy review. I think methods deserve 68%. For knowledge advancement, I'll lean towards 60 or 65%. The logic and communication are good—maybe a solid 78%. Open science transparency could get 65%, and I see high global relevance, so I'm giving that an 80%. Regarding publishing, this commissioned report might fit well in a top B-journal instead of a standard journal.

Evaluating journal suitability

Many economics and health economics journals might demand more formal modeling, but this report could be accepted in fields like development practice or health policy analysis. I'm leaning towards a score of 3.2 for the "should" rating, aiming for a top B-journal or strong field. For "will," I think a 2.8 is appropriate since it's a report rather than a journal article. Given the thorough nature of the document, I'll finalize "should" at 3.1 and "will" at 2.7. However, the guidelines call for credible intervals that match the scores. Got it!

Creating a summary

I'm preparing to craft a summary around 1000 words, but it feels like a lot! I'll aim for approximately 900 to 1100, focusing on issues that need specifics: the approach used, key contributions, and conceptual mismatches. I'll cover the wealth effect formula and inconsistencies. Topics include the reliance on Goldman Sachs, demographic oversights, measurement of time, and errors in double counting. I'll also mention identification threats and the lack of a systematic review, which may further limit replicability. Finally, I'll fill out the JSON structure as needed.

Crafting the summary

I'm ready to craft the assessment summary while keeping my percentiles in mind. Overall, I'm setting it at 72. For claims and evidence, I think 74 is fair but might stick with the earlier 72. The methods score will be 66, and advancing knowledge gets a 64. For logic and communication, I considered 80 but will settle on 78. Open science can be 64, and global relevance is set at 82. I will keep the tier "should" at 3.0 and "will" at 2.7, producing a thorough summary while avoiding individual names and citations as required.

Identifying key criticisms

I'm noticing several key criticisms in the analysis. First, while the conclusion favors keeping SRTP over SOC, it lacks quantified trade-offs for SOC and could benefit from a demonstration on a sample program. The derived formula they present seems helpful, but they didn't fully re-estimate the entire cost-effectiveness analysis. Raising the growth rate to 3.3% relies heavily on external sources, which introduces uncertainty without addressing the composite CAGR issue. The recommendations on temporal uncertainty and compounding non-monetary benefits also seem weak, alongside concerns around double counting and the lack of replication options.

A.4 Peterman et al. 2025

i Model assessment summary

This registered systematic review and meta-analysis assembles an unusually large experimental evidence base (93 studies, 1,307 effects across 45 LMICs) on social safety nets (SSNs) and women's economic achievements and agency, with clear inclusion criteria, broad SSN coverage, and careful outcome coding. The analytic strategy (standardized mean differences converted to Hedges' g ; robust variance estimation [RVE] to handle within-study dependence; multi-level RE and influential-point checks as robustness) is appropriate and largely well executed. That said, several features limit interpretability and design guidance. First, heterogeneity is extremely high throughout (I^2 commonly 80–95%+; $\tau^2 > 0.02$), yet the paper emphasizes pooled mean effects without reporting prediction intervals or distributional summaries of true effects. Given the breadth of contexts, modalities, and outcomes, a single pooled $g \sim 0.11$ risks overgeneralizing and is of limited policy value without a sense of between-study dispersion and likely range in new settings. Meta-regressions recover few moderators and retain high residual heterogeneity, suggesting either limited power to detect realistic moderators, substantial measurement noise in moderator coding, or genuine context dependence. Second, outcome harmonization is ambitious but stretches conceptual comparability. "Economic achievements" and "agency" are multi-dimensional constructs; the paper pools across diverse indicators (including binary and continuous) via SMDs. Although standard, SMDs reduce interpretability for practitioners (e.g., how large is $g=0.11$ in hours worked or decision-making probability terms?) and may mix substantively different measures under a common metric. Additional stratification by measurement type (binary vs continuous; index vs single indicator) and presentation of effects in natural units where feasible would improve external usefulness. Relatedly, it is unclear whether aggregate indices and their constituent components were ever both included for the same sample/time; explicit safeguards against double-counting would be helpful. Third, the dependence structure is handled with RVE, but weighting details are under-specified (e.g., exact weight formula, assumed within-study correlation ρ , and sensitivity to alternative ρ values). With 1,307 effects from 93 studies, a few studies likely contribute disproportionate numbers of outcomes and time points; sensitivity checks that collapse to one effect per study per domain (or use study-level aggregates) would test robustness to outcome multiplicity. Fourth, moderator coding seems too coarse or noisy for strong design conclusions. For example, "value of benefits" is in nominal USD rather than PPP-adjusted or scaled by local consumption/baseline income; "scale of implementation" and "pilot vs mid-level vs at-scale" appear categorical with potential misclassification; "cash-plus" is highly heterogeneous; and "gender targeting" shows little variation. These choices attenuate power to detect gradients (e.g., dose-response with transfer size) and likely contribute to null moderator findings. Fifth, time dynamics are insufficiently interrogated: most interventions last ~ 12 months with follow-up ~ 14 months post-intervention; but SSN effects can accumulate or fade. While a simple <12 vs $=12$ months split is reported, more granular modeling (event time, longest vs earliest follow-up, or pre-specified primary horizon) would clarify durability. Sixth, although the authors state "little evidence of publication bias," methods are not detailed in the main text, and standard small-study or funnel diagnostics are complicated by dependent effects; reporting study-level funnel plots, RVE-compatible Egger tests, or p-curve/p-uniform would strengthen this claim. Seventh, risk of bias/quality assessment is referenced (Table S4) and used as a

control, but it does not appear to inform weights or stratified analyses; explicit RoB domains (randomization, attrition, selective reporting, measurement) and RoB-based sensitivity would increase credibility. Eighth, the cost-benefit synthesis is informative but limited: only ~25% of studies report BCR/IRR; methods and benefit definitions vary; most do not include women's own economic gains or agency measures, making comparisons difficult and likely underestimating gender-specific value. A standardized re-analysis of costs/benefits where microdata are available (or harmonized assumptions) would be more decisive. Ninth, selection and generalizability: restricting to RCTs improves internal validity but may bias towards pilots and specific geographies (heavy in sub-Saharan Africa) and populations (few older women; fragile/urban underrepresented). The meta-regression's "pilot lower than mid-level" result is suggestive but could reflect targeting differences; external validity remains a concern. Tenth, minor internal inconsistencies in reported point estimates/precision should be corrected (e.g., public works effect size appears as 0.127 elsewhere but 0.273 once; a decision-making CI is printed as 0.036–0.37). Finally, while the study states replication files will be public and provides supplementary materials, the data extraction instrument is not public, limiting auditability of coding and decisions (e.g., outcome inclusion/exclusion, direction coding, plus-component classification). Despite these issues, the conclusions are measured, the main claims (small positive effects on women's outcomes overall, with UCTs/asset/social care robust; CCT and in-kind smaller) are consistent with the assembled evidence, and the paper adds substantial value by bringing together a broad experimental literature and by transparently acknowledging limits where moderators cannot be pinned down. Priority improvements: add prediction intervals; provide per-study-per-domain sensitivity; articulate RVE weighting/ and sensitivity; expand moderator measurement (PPP-normalized/consumption-scaled transfer size; clearer "plus" typology); deepen time-horizon analysis; present publication-bias diagnostics adapted to dependent effects; integrate RoB into weighting/stratification; and release the extraction protocol/instrument and a reproducible codebook linking each effect to its raw source and coding decisions.

Model reasoning trace

Formatting JSON assessment

I need to create strict JSON with specific fields: `assessment_summary` and a metrics section that includes percentile scores for various categories. There's a twist with midpoints, meaning I must produce point estimates instead of ranges, which is a bit unusual. I should read the content provided by the user and evaluate it based only on the PDF text. I'll aim for a diagnostic summary of around 1000 words, focusing on neutral and specific issues like identification threats and data limitations.

Evaluating provided content

I need to focus solely on the information given in the text, which is an extensive meta-analysis about the effects of SSN on women's economic achievements and agency. The methods used include RVE meta-analysis and conversions to Hedges' g , with rigorous assessments like multi-level RE and checks for influential effects. They aimed to include studies in multiple languages but found none in French or Spanish. Overall, they examined 1,307 effect sizes from 93 studies across 45 LMICs, with a pooled effect size of $g = 0.107$.

Analyzing effects and findings

The meta-analysis showed strong effects on savings, assets, expenditures, labor force participa-

tion, work intensity, voice, autonomy, and decision-making. However, they found null results for debt and care work intensity, and insufficient power for care work participation and leadership. Meta-regression indicated higher effect sizes in the ECA region but was largely driven by one study. Lower results were seen in LAC for basic specifications, especially with CCT compared to UCT. High heterogeneity persists, and only 25% of studies reported cost-benefit analyses, primarily related to UCT.

Identifying limitations and critiques

The analysis highlights several limitations, such as the inability to detect moderators, high heterogeneity, short follow-up periods, and restrictions to only RCTs. They also mention issues with measurement limitations regarding agency and potential publication bias. Moving from aggregated data to individual women's benefits presents challenges due to heterogeneity in components and limited gender targeting.

Now, I need to critique their methods, particularly around outcome selection and coding. They transform outcomes to standardized mean differences, but combining different measures may lead to interpretational issues. A small scale of g 0.1 coupled with high heterogeneity makes it hard to connect to policy significance, despite a few cost-benefit analyses provided.

Reviewing method suggestions

I'm considering how to improve interpretability by converting to relative risk or proportions, particularly for binary outcomes. They mentioned using standardized mean differences (SMD) and Hedges' g , which is typical, but it raises questions about mixing effect types. I noted the use of correlated-effects RVE that requires specifying the correlation parameter ρ , which wasn't stated. They might have defaulted to values like 0.8 or 0.7, affecting the standard error but not the effect magnitude.

Regarding small-sample corrections, they followed Tipton's guidance and used degrees-of-freedom greater than 4, which is positive. However, the 1,307 effects across 93 studies create a high ratio. The RVE is robust to dependencies, but I'm curious if they weighted by precision or used equal weights when clustering at the study level.

Examining weight usage and methodology

They mention using weights in their analysis but don't provide specific details. Typically, equal weights are used in RVE, so it's unclear if they applied the weighting method suggested by Hedges, Tipton, and Johnson, which involves specific calculations. I should point out this lack of clarity.

Additionally, they decided to mix stand-alone SSN with bundled components to capture moderated effects. However, this approach introduces definitional heterogeneity and potential misclassification due to reliance on textual coding, which might limit their ability to identify effect modifiers effectively.

Providing suggestions for analysis

I need to suggest better handling of studies with multiple arms and follow-ups. They note the time since intervention and treat the 'time post-intervention at survey' at the effect level, but I wonder how they managed multiple follow-up times from the same study. Clustering helps with dependence, but mixing repeated measures could favor the largest effect or create moral hazard. They might consider analyzing at a single time point or using the longest-run measurement to minimize multiplicity.

Regarding the unit of analysis, they include both individual and cluster RCTs, but I'm unsure if they adjusted standardized means for clustering. There could be a risk of ignoring intraclass correlation, affecting replicability. I noticed they mentioned quality assurance

but didn't provide specifics on risk-of-bias assessment for randomization, attrition, or measurement across included studies.

Reviewing assessment details

They mention a "quality assessment score," which is a start. However, when discussing publication bias, they acknowledge little evidence but don't provide details about Egger's test. For RVE with dependence, analyzing funnel plots and using Egger's test may be complicated, and they haven't described any robust methods.

Regarding selection, only including RCTs may bias their sample towards small pilot interventions, which often show smaller effects, raising concerns about generalizability.

They conducted searches in other languages but found none, possibly indicating selection or regional bias. I suggest they consider adding forward and backward citations or translations. Lastly, the data extraction form isn't publicly available, creating barriers to replication. They mention plans to deposit replication files, which is reassuring, but the lack of access to the data extraction form could hinder understanding of their coding decisions.

Highlighting replicability limitations

There's a replicability limitation here. They perform subgroup analyses, but with small sample sizes, such as public works with $n=106$ effects across 9 publications, which yields significance at $p<0.05$. However, the effect size of $g = 0.127$ may be overshadowed by aggregated results showing variation in measurement.

The large effect size for the ECA region seems driven by one study with a high benefit level, and they rightly caution against it.

Additionally, the definition of "agency" suffers from conceptual heterogeneity, as they aggregated across measures. I wonder if they ensured proper comparability in their analysis.

Examining methodological concerns

They mention "woman specific" indicators but don't clarify which ones, potentially mixing household-level asset ownership with women's ownership. It's good they exclude non-woman-specific indicators.

Regarding binary outcomes, using standardized mean differences (SMD) might not be the best choice without clear details, which raises replicability concerns.

Also, the high heterogeneity of I^2 at 90%+ alongside a pooled effect of $g = 0.107$ seems problematic. The effects are likely too broadly distributed, and they didn't include prediction intervals, which could be useful.

Multiplicity and analytic flexibility also raise concerns, as correlated measurement may introduce selection biases. They coded pilot versus mid-level versus large-scale implementations, but variable definitions lack clarity. Presenting the pooled effect may misrepresent data, and quantifying treatment effects would enhance understanding.

Analyzing the interpretation of outcomes

It's important to note that mixing multiple correlated outcomes per study can introduce weight biases. While RVE attempts to manage this, aggregating all outcomes together, especially across different domains, makes it hard to interpret in terms of women's empowerment.

When they conclude that "SSNs can economically and socially empower women," it feels a bit over-interpreted as empowerment is a normative term. They don't measure intangible aspects or address high heterogeneity, which should moderate their conclusions.

Regarding "cash plus" components, they suggest these are rarely more effective than cash alone, but they didn't observe this in their sample and should be cautious about heterogeneity. For cost-benefit analysis, they summarize normative findings without standardized discount

rates, which could be misleading. They also mention risk-of-bias indicators but don't clarify what qualifies as high quality.

Considering quality and data management issues

It's a bit puzzling that they mention "if medium quality"; it sounds like they may only have low vs. medium categories, but I think this needs more detail. Not weighting the meta-analysis by risk-of-bias could be an issue, though it might not cause major shifts.

Regarding potential double-counting, many publications have multiple arms, and I'm wondering if they've adequately avoided counting the same control group multiple times. RVE partially addresses this, but mixing effect sizes with shared controls could lead to double counting, even though they cluster by study, which seems acceptable.

For data management, they aim to provide robustness against 'Cook's D' influences, but the threshold remains unspecified, and visual inspections might introduce subjectivity. They mention that replicates are unaffected, which is positive.

Lastly, when it comes to domain classification, their definition of "majority" in aggregated indicators needs clarification, and I'm curious if they've mistakenly double counted aggregated indicators and their components.

Addressing data quality and outcome concerns

There's a risk of double counting here because it seems they might treat separate effect categories as unique. However, including both aggregated indices and subcomponents from the same sample could overweight specific studies. I wonder if they avoided this mixing, but they didn't say explicitly.

Also, they're not analyzing potential harm, such as domestic violence or intra-household tensions, which feels like a significant oversight since they're focusing on positive outcomes like economic achievements and agency.

When it comes to nonlinearity with benefit values, they used terciles of 'total value of economic benefits' without adjusting for purchasing power parity (PPP) or local prices, which might mismeasure 'value' due to variations across countries.

For assets and savings, mixing 'any savings' with total amounts seems off, as SMD calculations may not be standardized across types. Finally, regarding expenditures, did they consider whether these were specifically women's own spending versus household expenditures?

Clarifying outcome measurement concerns

They say outcomes are "woman-specific only," but I'm not entirely sure about how they measure expenditures. It seems like expenditures might be assessed at a household level, yet they state that they only include "woman-specific" outcomes. I wonder if they truly focused on women's own expenditures or if they incorporated aspects like "consumption per adult female." That approach seems a bit uncertain because it raises questions about whether they're accurately capturing women's unique impacts.

Questioning measurement and validity

I'm considering the risk that they might have included household-level data, even though they claimed to focus on "only woman-specific" outcomes, so it's likely they limited their analysis to women's consumption. I need to check for replicability.

Then there are agency measurement issues; they've noted limitations in reliability which is worth considering. Also, they used a 12-month threshold to categorize duration and follow-up, but I'm wondering if testing other thresholds might yield different insights since they didn't show those details.

Regarding external validity, most studies are based in sub-Saharan Africa, which introduces

regional variations. They often look at youth versus adult demographics, noting that there aren't many older women represented nor a balance between urban and rural settings—this could skew their findings.

Finally, the classification of sample age groups seems approximate since they use share-of-sample averages, which might lead to misclassification. Their comparisons between cash conditional transfers (CCT) and unconditional transfers (UCT) could also be confounded by regional differences. They tried to control for this, but interplay is still an issue. Additionally, what do they mean by “social care services”? Does this include childcare, counseling, and case management?

Identifying discrepancies in effect sizes

They claim the effects are robust, but the numbers at the effect level seem small, and there's a lot of heterogeneity.

For “public works programs,” there's a confusing discrepancy: they state an effect size of $g = 0.273$, but the figure shows 0.127. In the text, they initially wrote 0.273 beside a significance of $p < 0.05$. However, the figure indicates 0.127 with significance. So, it looks like I found a typographical error where they initially mentioned 0.127 but later typed 0.273 instead. It's important to clarify this!

Noting internal inconsistencies and errors

I'm spotting a few internal inconsistencies that should be highlighted, perhaps as minor misprints. For instance, the “decision-making effect” is reported as 0.087 with a confidence interval (CI) of 0.036-0.37, but that CI looks quite wide and maybe incorrectly transcribed. There are also some formatting issues with the confidence intervals. Additionally, many of their meta-regressions might exceed the threshold of 4, suggesting unstable results, so they should mention small sample corrections. Finally, they reference their quality assessment without typical details, which is worth calling out.

Addressing quality and bias concerns

I'm looking at the Cochrane risk-of-bias tool, and I think it's worth mentioning that the quality variable seems questionable. It feels odd to limit categories to low and medium without including high quality.

They've coded various settings and components well, focusing on RCTs and addressing grey literature. However, they claim a lack of evidence for publication bias without a thorough analysis. There are also challenges, like the potential overweighting of results from large studies, and the use of I^2 with RVE, which isn't standard. Testing sensitivity for effects could improve the reliability of findings.

Evaluating quality and claims

I'm considering the quality of the paper, which seems strong due to its thoroughness and large sample size. It's a rigorous meta-analysis that uses RVE methods and is pre-registered, which is great for replicability. However, while they state that data and analysis replication files will be publicly available if accepted, the lack of publicly available data extraction forms detracts from their open science commitment. I'd rate their open science as high, around 78 or 80.

Their claims about SSNs empowering women are made carefully, which is positive, but I wonder if there's a risk of over-claiming without sufficient evidence of moderate effects.

Assessing claims and methods

They propose eliminating conditionalities, but they also acknowledge that it's difficult to fully separate conditions from design features, which is a responsible approach. Their claims about social safety nets (SSNs) empowering women seem credible, so I'd rate claims evidence around

84.

Their use of RVE methods is robust but does have some issues, like mixing standardized mean differences (SMD) across outcomes; I'd give methods an 82.

For advancing knowledge, this paper pools data across modalities, which is significant, especially in 2024 with this large dataset, so I'd rate that 85.

The writing is clear, but there are minor inconsistencies, like a typo related to numeric values. For open science, I'd rate it a 77, while the global relevance is high at 92 due to the widespread use of social safety nets. Overall, I'd provide a weighted rating of 86.

Setting credible intervals

I'm thinking about how to set the credible intervals, as the instructions specify that the lower bound, midpoint, and upper bound must all be equal. Traditionally, a 90% credible interval is expected, but here, the focus is on setting all bounds to equal the midpoints. This approach aligns with the analysis of meta-analyses related to gender and social safety nets (SSNs) in low-to-middle income countries (LMICs). Considering the high quality of this analysis, I might rate it around 89 overall but weigh disclaimers to set it closer to 86. As for "tier_should" and "tier_will," I'll categorize them on a scale from 0 to 5.

Assessing article tier

I'm considering the type of article, which incorporates new findings and could potentially fit into a high-tier general interest journal like "Science." However, I need to stick to the content without relying on any external knowledge. Personally, I'd assign it a score of 4.2, suggesting it's a marginal A-journal or a top field journal like "World Development" or "Journal of Development Economics." I'll assign scores of 4.3 and 4.1 for "should" and "will," respectively, while ensuring that the lower and upper confidence interval scores match the assigned score according to the instructions. Given the issues of cost-benefit and heterogeneity, I believe a score of 4.2 is reasonable for "should" and 4.0 for "will."

Preparing summary outline

I'm getting ready to produce a 1000-word summary that will highlight major issues and benefits of the analysis. I want to include method-level details, such as the scope and conceptualization, high heterogeneity, mixing of outcomes, and minimal details on weighting and correlation parameters in the RVE. I also need to mention how they calculated SMD and converted binary data, alongside a note about the lack of replicable aggregators.

I'll point out measurement issues related to transfer value not being PPP, and the absence of predicted effect sizes like "marinade." I'll ensure the summary circumvents external knowledge and avoids referencing authors or institutions. Paraphrasing is essential to avoid copying text. I'll also note any typos or inconsistencies I've noticed, like the discrepancy between 'g 0.273 vs 0.127' and the 'decision-making 0.036-0.37.' Lastly, I want to include the replicability note about the data extraction forms not being public.

Starting summary crafting

I need to mention the figure misalignment in the summary. Let's begin with: "The paper offers a registered systematic review and meta-analysis of experimental studies on SSNs in LMICs, focusing on women's economic outcomes and agency. The search identifies 93 studies and 1307 effect sizes, calculating standardized mean differences using Hedges' g. They apply robust variance estimation (RVE) to handle within-study dependence, run multilevel random-effects for sensitivity, and present orchard plots and meta-regressions with study, intervention, and effect-level covariates. Results show a small positive pooled effect of g 0.107, notably stronger for UCT, asset transfers, and social care, while public works yield positive but less

precise results and CCT shows a smaller effect with in-kind not statistically significant.”
Next, I’ll list the issues, starting with the extremely high heterogeneity: $I^2 > 80-95\%$ and $t^2 > 0.02$.

Outlining key issues

The paper presents an average effect across diverse interventions, outcomes, and populations, but lacks prediction intervals and distribution of true effects, which might make the mean less informative for policy.

Next, outcome harmonization is a concern; terms like “economic achievements” and “agency” are broad, and their subconstructs are measured heterogeneously. The use of standardized mean differences combines continuous and binary data, limiting clear definitions.

Multiplicity of effect sizes per study is another issue; while RVE addresses correlation, it doesn’t consider study-level quality. If some studies contribute many outcomes, they could dominate the results. I’m curious if the authors considered using one effect per study per domain as a sensitivity check, but that’s not reported.

Additionally, the moderator coding is problematic since “transfer value” is measured in total USD instead of PPP, and terms like “scale” lack clarity. The categorization of “plus components” is too coarse, and “gender targeting” shows limited variation, which hampers design insights.

Highlighting additional concerns

There are several issues to consider with the study. First, selective reporting might be a problem, as many included RCTs may not measure or report gender-disaggregated outcomes. Without systematic extraction from registry protocols, we can’t confirm outcome reporting bias, which is mentioned but not fully addressed.

Next, regarding publication bias, they state there’s “little evidence,” but they don’t describe their method, which makes it hard to trust their conclusions.

Cost-benefit integration is also lacking, as the methods aren’t standardized and aren’t used to inform effect sizes. Many benefit-cost ratios exclude considerations of women’s incomes, which compromises their analysis.

Furthermore, sample restrictions are apparent as the study only includes RCTs weighted to Sub-Saharan Africa, lacking older women and showing insufficient representation from fragile contexts.

Many effects are measured within a 12-14 month window, which doesn’t account for long-term results.

Finally, while they promise replication, the extraction form isn’t publicly available, raising doubts about replicability.

Evaluating key metrics

There are more factors to consider. The risk-of-bias assessment scoring is unknown. It seems like they indicate low or medium quality without sufficiently addressing high risk, only using it as a control.

Next, I noticed some inconsistencies, like the public works effect size being reported as 0.127 versus 0.273, and decision-making confidence intervals printed as 0.036–0.37.

Now, I need to produce the metrics, adhering to strict JSON formatting. For global relevance, I’ll set it high at 92, while open science will be 78.

For the overall metrics: midpoint at 86, claims and evidence at 84, methods at 82, advancing knowledge maybe at 86, logic and communication at 80 due to inconsistencies. This gives me a “tier should” of around 4.4 and “tier will” at 4.1.

Clarifying metrics and summary structure

I need to unify the scoring from 0-5. I think I'll go with 4.2 for "should" and 4.0 for "will". As for "journal A," while it's a top general field, the rating could be around 4.3, but I'm leaning towards 4.4 for "should" and 4.2 for "will".

I'll make sure each metric, like "overall" and "claims_evidence," is structured with midpoint, lower_bound, and upper_bound.

They want a summary of about 1000 words, aiming for 850-1100 words. I'll ensure it's a compact paragraph, focusing on key issues like data limitations and internal inconsistencies, and avoid heavy markup and bullet points. I want to present a clear and concise analysis!

Finalizing important issues summary

I'm focusing on crafting a summary that highlights the most important issues, making sure to do this elegantly and without mentioning any authors. I need to avoid mentioning any lack of samples as well. It's essential to finalize the document without oversight, such as ensuring the images represent tables and figures correctly.

I'll aim for a word count of about 900-1100 words, targeting roughly 1000 words in total. Now, I'm ready to complete this step!

ewpage

A.5 Williams et al. 2024

i Model assessment summary

This paper delivers a high-resolution, pantropical map of where deforested lands have biophysical potential to undergo natural forest regeneration, and it aggregates this to country and biome summaries and an overlay-based carbon removal potential. The approach is transparent in intent—predicting potential, not feasibility—but several methodological choices and data limitations affect inference strength and precision. The dependent variable (presence/absence of natural regeneration during 2000–2016) is derived from a prior global remote-sensing classification that conservatively emphasized omission errors in humid regions and required a minimum patch size; area-based producer accuracy for natural regrowth in the humid biome was low while user accuracy for detected patches was high. This means the positive training class is a non-random subset of true regrowth, with systematic under-detection likely concentrated in particular forest types, geographies, and patch sizes. The predictive model then inherits these detection and sampling biases, potentially distorting learned relationships and shifting mapped potential away from contexts where regrowth was common but poorly detected. Model fitting uses random forests on millions of stratified points, with a sensible variable-selection procedure; predictions are ultimately based on biophysical predictors alone because a combined socioecological model yielded very similar accuracy. Omitting socioeconomic predictors makes the maps more temporally stable and widely available, but it also reframes outputs as “biophysical potential” and risks overpredicting potential in places where human pressure or governance constraints are binding. Validation emphasizes overall accuracy around 88% using a balanced validation set; however, accuracy declines at moderate distances from training locations, indicating residual spatial autocorrelation and some optimism in headline accuracy. The use of accuracy rather than threshold-free metrics, and the balancing of

classes (prevalence set to 50%) during training and validation, further complicates interpretation of both accuracy and the calibration of predicted probabilities. The paper interprets the continuous model output as a probability and converts this to “expected area” by multiplying by pixel area and summing, yielding 215 Mha. Yet the calibration of random-forest scores to true probabilities is not established; without calibration (e.g., isotonic or Platt scaling) under the real-world prevalence, the expected-area identity may be biased. The extraordinarily narrow “confidence intervals” reported for area (± 0.22 Mha at global scale) are not credible statistical uncertainty; they reflect computational or rounding artifacts, not propagation of uncertainty from training data errors, cross-validated prediction error, input-layer uncertainty, or probability calibration. A threshold (>0.5) sensitivity shows a substantially larger area, highlighting dependence on interpretation choices. Many predictors are coarse (250 m–1 km) while predictions are at 30 m; this scale mismatch introduces blockiness and may overstate local precision. The domain of “available for restoration” is defined liberally by excluding only forests, water, urban, bare, and forestry areas; while appropriate for modeling, its presentation alongside country totals risks being misread as areas that are socially or legally available for reforestation. Temporal assumptions also matter: the model projects to 2015/2030 by assuming that 2000–2016 relationships hold and by updating some forest-context layers with 2018 data; it does not incorporate climate change trajectories, shifting fire regimes, or evolving socioeconomics that strongly affect regeneration and permanence. Carbon estimates are calculated by downscaling a 1 km accumulation dataset to 30 m, multiplying by modeled potential, and summing; the stated range lacks a clearly defined uncertainty model and appears not to propagate classification, calibration, or permanence risks. The paper is commendably explicit that realized sequestration will be lower due to leakage, reclearance, and policy constraints, yet some statements comparing potential annual removals to current pantropical losses could be interpreted too strongly given these omitted risks and uncertainties. Spatial performance is uneven, with lower accuracy in parts of Southeast Asia; country and biome summaries lack uncertainty decomposition, which would be useful for decision makers. On replicability, large gridded outputs are publicly archived and inputs are open, but code is available only upon request; combined with the volume of preprocessing, this raises barriers to full reproduction, code review, and extension. Strengths include a clearly articulated objective focused on actionable planning, a very large training sample, careful variable selection, explicit treatment of autocorrelation effects (even if not using full spatial block cross-validation), and extensive caveats about use and governance. The resulting layers will be highly valuable as a first-pass screen for assisted natural regeneration opportunities and for national restoration dialogues. The main improvements needed relate to uncertainty characterization (probability calibration, spatially explicit error propagation, and country-level error budgets), spatial cross-validation that reduces spatial leakage, clearer separation of statistical precision from conceptual ranges, dynamic scenario analysis including plausible climate and human-pressure changes, and full open release of code. Overall, this is an influential, practically useful mapping contribution with limitations that mainly concern probabilistic interpretation, uncertainty, and temporal and socioecological external validity.

Model reasoning trace

Analyzing model limitations

I see another limitation in the model’s assumptions. They assumed stationarity from 2000-

2016 when predicting for 2015 to 2030. They included present conditions like 2018 tree cover but didn't factor in future climate or socio-economic changes. By excluding socio-economic predictors, they reduced potential model improvements slightly, but this could bias predictions in areas where socio-economic factors play a significant role, potentially leading to unrealistic outcomes in "biophysical potential." They clarify it's about modeling potential, not feasibility.

Examining dataset limitations

I notice another issue with the primary regrowth dataset used, which came from Global Forest Watch to map regrowth over 5 m. In humid tropics, the producer's accuracy is low at just 18.7% because of omission errors, although user accuracy is high at 85.1%. The detection biases may skew the training sample and affect model calibration. Their conservative approach to mapping regrowth likely favors more accessible and larger areas, and many small regrowth patches are excluded. It's unclear if the training sample is representative across different biomes without proper stratification.

Clarifying domain definitions

I need to define the domain used in the study. They included all parts of forest biomes within $\pm 25^\circ$ but excluded water, bare land, urban areas, and so on, while purposely opting for a liberal domain. They used 2018 tree cover to exclude forests but may misclassify certain savannas or drylands as part of forest biomes. While they mention not including savannas, their definition could still capture mosaic-savanna zones. Additionally, there's a scale mismatch since predictions at a finer 30 m resolution are computed alongside coarser predictors, which could impact accuracy.

Assessing the weighted area approach

I see another issue with the weighted area approach, where the probability of regrowth times pixel area is calculated. However, the random forest classification is trained on presence/absence from 2000-2016, and the output probabilities aren't properly calibrated. Calibration techniques like Platt scaling or isotonic regression are necessary to interpret results as probabilities. They also didn't present AUC, only accuracy, which may mislead since balanced sampling artificially sets prevalence to 0.5. Without proper probability calibration, the estimate of expected area could be questionable and not reflect actual conditions.

Reviewing definitions and calibration

They used specific definitions, claiming that "the continuous potential" can be seen as the probability of natural regeneration per pixel, with the weighted-area value representing the expected area for regeneration. However, they didn't perform any calibration on this. Their accuracy and out-of-bag accuracy are reported as 0.878, but they failed to discuss probability calibration, which is important. Regarding carbon sequestration potentials, they calculated totals and ranges from previous datasets, but it's unclear where the range comes from—possibly from Monte Carlo methods or carbon dataset quartile variation, yet they didn't show how they computed uncertainty.

Identifying dataset uncertainties

The dataset presents uncertainties by grid cell, but their method multiplies 30-year totals without mentioning how uncertainty is handled. The "range" they provide likely comes from minimum and maximum values across these uncertainties, which is unclear. They need to clarify how uncertainty propagates and address independence and omitted below-ground carbon, since they reference adding 22–28% from literature. Additionally, they didn't consider factors like degradation baseline, leakage, or permanence. While they frame the "biophysical potential," the abstract could mislead policy implications with its stated above-ground carbon

sequestration potential of 23.4 Gt C over 30 years.

Checking carbon sequestration claims

I noticed they mentioned that “new carbon sinks could increase global sequestration by 14.3% per year” but also stated “mitigating current pantropical forest carbon losses by 90.5% per year,” which might be inconsistent. They claim assisted natural regeneration can sequester significant carbon but didn’t show how they calculated this. It seems to derive from 23.4 Gt over 30 years, equating to 0.78 Gt C per year. With a net carbon loss of about 0.86 Gt C per year, I need to verify Baccini 2017’s values, which reported a net carbon source of 1.1 Gt C per year. My memory on this is a bit fuzzy!

Evaluating carbon loss mitigation claims

They state that “mitigating current pantropical forest carbon losses by approximately 90.5% per year,” referencing Baccini, raises some questions. If regeneration can offset about 0.8 Gt per year, that suggests total losses of around 0.9 Gt per year. While this might be plausible, there are logical leaps involved since their potential benefits might not actually be realized. It’s a bit tricky to interpret, yet they present it confidently as potential. I need to consider these assumptions carefully!

Analyzing dataset issues

They wisely exclude plantation regrowth from their analysis, but it seems their training dataset might still contain misclassified plantations. They reported a previous algorithm with 90.6% accuracy; however, an area-based producer accuracy of just 18.7% in humid tropics suggests significant under-detection. This non-random classification error could undermine generalizability. While they argue for the robustness of the detected patches, the validation accuracy of 87.9% relies on potentially biased training dataset presence points. Field validation is unlikely, and poor accuracy in Southeastern Asia raises another red flag.

Evaluating modeling assumptions

The modeling approach assumes that predicted ‘potential’ is valid from 2015 to 2030 without considering changes, which is questionable since the climate is constantly shifting. Additionally, the “distance to forest” will likely change as regrowth occurs, yet they didn’t model this dynamic feedback. They mentioned underestimating impacts like positive feedback and the exclusion of savanna/dryland. Climate change could also hinder regrowth, leading to ambiguity. They present a figure showing “available for restoration” area by country, but they don’t clarify how this was calculated, leaving some mystery around the methodology used.

Analyzing definitions and reproducibility

They define the study’s domain as all tropical forest biomes, excluding areas like water, urban, and existing forests, while including factors like ‘distance to water’ as covariates. It seems that “available for restoration” refers to the total area of this domain, but their disclaimers indicate some harvest/regrowth classifications might be included. Although they mention using a forestry dataset, they didn’t clearly define “available for restoration.” Additionally, requesting code via email isn’t ideal; it should be publicly accessible. While they use accessible datasets, the lack of code hampers reproducibility, and the confidence intervals aren’t credible.

Evaluating accuracy and methodology

They mention the model overcame potential overfitting with an out-of-bag (OOB) accuracy of 87.8%. They used a balanced sample and omitted socioeconomic features due to similar accuracy, but I wonder if the limited role of socioeconomic data in predicting presence versus absence, as indicated by remote sensing, could affect results. They also disclaim that the maps are only one input, which is good. Since this is published in Nature, it seems high quality,

but I believe we should look beyond that. We need to propose credible metrics.

Ranking evaluation criteria

Based on the content, I would rank high—maybe around the 85th percentile overall.

For claims and evidence, they articulate strong claims and provide robust mapping and a predictive model. However, there are issues like reliance on one training dataset with major omissions and unrealistic confidence intervals. I'd place this at about 75.

For methods, they've implemented random forest modeling and accounted for autocorrelation, yet there's heavy reliance on coarse data. This leads me to perhaps a 70.

Regarding advancing knowledge, their maps on natural regeneration potential are valuable, globally relevant for restoration planning, so I'd mark this as a 90.

On logic and communication, while the writing is clear and disclaimers are present, there are leaps in concluding about carbon losses and some misuses of "CI." I'd score this around 82.

For open science, while the data is available on Zenodo and the input data is open, the code is only available on request, which limits replicability. I'd give a score of 65.

Finally, in terms of global relevance, this work ties in well with global priorities and quantifies information by country, so I'd rate this a solid 92.

Assigning tier ratings

I need to assign "tier_should" and "tier_will" ratings from 0 to 5. Based on quality, I think "should" deserves a high score, perhaps around 4.5 or 4.7, due to some identified issues. For "will," since it's already published in Nature, I'd predict around 4.6 to 4.8. However, they want continuous numbers, so I'll set "should" at 4.4 and "will" at 4.7.

If I consider possible underestimation due to method-related issues, perhaps "should" could be 4.2 or 4.4, while "will" could be 4.8 because of strong interest. Still, the dataset not being open-coded is a consideration, but it's published in Nature, so I want to remain fair.

Now, I need to craft a summary, aiming for around 1000 words. I'll start by summarizing the approach and identifying major issues.

Highlighting dataset issues

I'm noticing that the training dataset has biases, especially regarding natural regrowth mapping, which shows significant omissions in humid biomes. This can lead to a non-random sample of presences, and the model might learn from incomplete detections, distorting variable relationships and potentially mispredicting areas.

Next, regarding probability calibration and the "expected area" method, they treat random forest outputs as probabilities, but without proper calibration and baseline prevalence, the resulting estimates of 215 million hectares could be inaccurate. They present highly precise confidence intervals that don't seem credible because they didn't account for uncertainty across inputs.

Lastly, concerning spatial autocorrelation, while they tried to evaluate it, accuracy tends to decline with distance, so their optimism might not be justified. They didn't apply block cross-validation and reported overall accuracy, avoiding metrics like AUC. The balanced classification gives an unrealistic sense of accuracy, and there are no confusion matrices or biome-specific metrics presented beyond overall accuracy.

Addressing scale and domain issues

I see that there's a scale mismatch with the 30 m predictions derived from 300 m to 1 km predictors, which could cause smoothing and artifacts. They've acknowledged this with visual evidence in Fig. 1, suggesting it could impact local planning as they consider these results to be preliminary.

The domain definition appears overly broad, including areas where natural regeneration may be undesirable—like non-forest ecosystems, productive farmland, or communities—which could potentially mislead their “available for restoration” metrics.

Lastly, they assume conditions from 2000 to 2016 without modeling climate change or socioeconomic factors, which could really undermine their 2030 predictions. Also, the carbon accumulation overlay from the Cook-Patton dataset has a scale mismatch issue, with 1 km data being downscaled to 30 m without considering edges and heterogeneity.

Evaluating carbon accounting and code accessibility

They don’t initially account for below-ground carbon, indicating an additional 22–28% should be included, and they don’t consider the risk of reclearance. Their discussions on leakage and additionality seem inadequate since their yearly rates are used to support claims of significant offsetting of losses, which likely exaggerates realized mitigation.

While they exclude plantations, algorithmic misclassification might still exist, and they haven’t incorporated a classification probability step.

Regarding code accessibility, it’s not publicly available, making replication difficult. Although they provide predicted maps on Zenodo, the training and modeling code needing to be requested further hampers replicability. Additionally, using 2018 tree cover for distance predictions while training on 2000 cover could lead to variable shifts, potentially degrading the predictions since they haven’t cross-validated changes.

Observing assumptions and suggestions for improvement

The weighted area assumption sums across all pixels to determine “expected area,” but it overlooks spatial dependence and positive feedback, which are mentioned but not integrated. The country-level summarization is good, although it uses administrative boundaries without addressing ecoregion differences, and there’s a lack of error budget for these country-level sums.

They note that predictions may degrade in Southeastern Asia, but don’t confirm this for later years. For improvements, I should suggest block cross-validation, probability calibration via isotonic regression, and incorporating socioeconomic variables.

Overall, I thought the score should be around 85, but I’m reconsidering an 86.

Finalizing ranking and summary considerations

I’m setting the overall score to 86, considering it’s intended to reflect percentiles within [0, 100]. For research in remote sensing global mapping, I’d lean toward an 88, but due to certain issues, I might adjust this to 82. However, advancing knowledge could be rated higher at 90, leading me to average around 86, so choosing 87 seems fair.

I also need to fill out the JSON carefully, ensuring that the assessment summary is a concise paragraph of about 800-1100 words and omitting any author names or external knowledge outside the provided PDF. The confidence intervals will match the midpoint of 86. For claims and evidence, I’m thinking a score of 79?

Setting scores and producing JSON details

I’ve decided to set the score of 78 for overall assessment, with specifics for other categories: 73 for methods, 92 for advancing knowledge, 84 for logic communication, 62 for open science, and 94 for global relevance. For tier should, I’ll use a score of 4.4, and for tier will, I’ll set it at 4.8, all while keeping consistency in mind with overall scores at 86, advancing at 92, and global relevance at 94.

The JSON needs to include an assessment summary of around 1,000 words and metrics for each score calculated with midpoint, lower, and upper bounds. I’ll also avoid naming authors

or their institutions, but I can reference previous analyses without using specific names or citation numbers.

Crafting the summary points

I'm going to paraphrase the training dataset description, ensuring I accurately reference the analysis of natural regrowth from 2000 to 2016 without including citation numbers. I should say, "The study builds on an existing pantropical remote sensing analysis that mapped natural regrowth and distinguished it from plantations."

Next, I want to finalize the summary by outlining key points. For data, I'll mention potential mapping across tropical forests at a 30 m resolution, calculating weighted areas for 215 Mha, and noting variations across ecoregions.

For methods, I'll detail the balanced random forest approach, variable selection, the training dataset derived from regrowth mapping, and the overall accuracy. I'll also highlight potential issues like omissions in the training dataset and domain definition relative to thresholds.

Assessing model limitations

I need to go over some aspects of the model's predictive capabilities. First, I see that there's a lack of probability calibration in the balanced training; the predicted probabilities aren't calibrated, which might lead to bias in the weighted area approach and unrealistic confidence intervals.

Next, there's the issue of carbon overlay and uncertainty propagation; downscaling seems to be done naively. The coarser covariate resolution could degrade the 30 m predictions, causing blockiness. Socio-economic variables are omitted in terms of feasibility versus potential, and there's a failure to integrate dynamic processes like climate change impacts.

I should also note that accuracy in Southeastern Asia is poor, and the code's availability is limited, impacting replicability due to the heavy dataset.

Lastly, I'll mention that the additional risk of reclearance isn't accounted for in the carbon measurements. I need to finalize this while including suggestions and discussing uncertainties qualitatively as we maintain strict JSON formatting.

ewpage