

¹ Using Satellite Imagery and Deep Learning to Evaluate the Impact
² of Anti-Poverty Programs

³ Luna Yue Huang^{*,1,2}, Solomon Hsiang^{2,3}, Marco Gonzalez-Navarro¹

⁴ This Version: April 23, 2021

⁵ **Abstract**

⁶ The rigorous evaluation of anti-poverty programs is key to the fight against global poverty.
⁷ Traditional evaluation approaches rely heavily on repeated in-person field surveys to measure
⁸ changes in economic well-being and thus program effects. However, this is known to be costly,
⁹ time-consuming, and often logistically challenging. Here we provide the first evidence that we
¹⁰ can conduct such program evaluations based solely on high-resolution satellite imagery and deep
¹¹ learning methods. Our application estimates changes in household welfare in the context of a
¹² recent anti-poverty program in rural Kenya. The approach we use is based on a large literature
¹³ documenting a reliable relationship between housing quality and household wealth. We infer
¹⁴ changes in household wealth based on satellite-derived changes in housing quality and obtain
¹⁵ consistent results with the traditional field-survey based approach. Our approach can be used
¹⁶ to obtain inexpensive and timely insights on program effectiveness in international development
¹⁷ programs.

^{*}Correspondence: yue.huang@berkeley.edu. ¹Department of Agricultural and Resource Economics, UC Berkeley, Berkeley, CA, USA. ²Global Policy Laboratory, Goldman School of Public Policy, UC Berkeley, Berkeley, CA, USA.
³National Bureau of Economic Research, Cambridge, MA, USA.

18 1 Introduction

19 Rigorous impact evaluation forms the basis of the modern approach to fight global poverty and
20 provides input for evidence-based policy making [1, 2]. The impacts of anti-poverty interventions
21 are almost universally evaluated using household surveys, typically comprehensive questionnaires
22 containing hundreds of questions that can touch every aspect of people’s lives. However, such field
23 surveys are often prohibitively expensive to conduct [3, 4] and unanticipated events, such as polit-
24 ical unrest or public health crises, frequently disrupt them [5]. In this paper, we provide the first
25 demonstration that the household welfare impacts of a large scale anti-poverty randomized con-
26 trolled trial (RCT) can be accurately measured relying solely on satellite data, instead of household
27 surveys.

28 Recent advances enable poverty to be identified remotely [6–13] and widespread adoption of
29 mobile phones allows targeted anti-poverty interventions to be deployed over-the-network [14],
30 such as the cash transfer program we study here [15]. By demonstrating that the impacts of
31 such interventions can be evaluated remotely, we hope that future programs can, in principle,
32 be designed, deployed, and evaluated with limited reliance on logically complex and expensive
33 ground operations. Because costs and logistics play a major role in limiting the scale of anti-poverty
34 programs [16], simplifying their deployment and evaluation is crucial to achieving their full global
35 potential.

36 We study a large pre-existing trial [15] that was recently completed and evaluated with field
37 surveys and show that we can consistently recover the impact of the program using satellite imagery
38 and deep learning. While previous studies have successfully evaluated the environmental impacts
39 of randomized controlled trials with remote sensing data [17, 18], we are not aware of studies that
40 demonstrate similar successes for household economic well-being. Specifically, we combine high-
41 resolution daytime imagery [19] and state-of-the-art deep learning models [20] to measure housing
42 quality among treatment and control households, and estimate the program effects on housing
43 quality. We then map housing quality to household wealth for these households by inverting an
44 “Engel curve,” an established concept in economics [21–24] that describes household spending on
45 specific goods as a function of economic well-being. Using this approach, we accurately recover
46 the program effects on household wealth for a fraction of the cost (\$0.006 per household, see
47 Supplementary Materials B) that would typically be spent on household surveys (\$18–300 per
48 household [4]).

49 Early work has shown that satellite data can be used to monitor economic development by cor-
50 relating nighttime luminosity, i.e., the amount of light emitted from Earth at night (hereafter “night
51 light”) with Gross Domestic Product (GDP) at national and subnational scales [25–27]. However,
52 the night light data show poor sensitivity in less developed and rural areas [6], presumably because
53 of low electrification rates—for example, from 1992 to 2008, 99.73% of pixels were completely unlit
54 in Madagascar, 99.47% in Mozambique, and this is representative of low-income countries [25].

55 This makes the data less useful for studying the very target of many international development
56 programs—people living under the poverty line. Additionally, the low spatial granularity of night
57 light prevents it from being used to evaluate programs reliant on fine spatial variations, includ-
58 ing most randomized controlled trials in which households in close proximity to one another are
59 assigned to different treatments.

60 We propose an alternative approach—we analyze daytime imagery using a deep-learning model
61 [20] to explicitly measure the quality of housing, a tangible and verifiable asset that is known to
62 be a powerful proxy for household wealth. Even in communities where electrification rates are
63 low, housing quality remains a strong predictor of wealth, in part because housing accounts for a
64 sizable portion (10–20%) of total household expenditure globally [28]. Furthermore, in many rural
65 and low-income contexts, individuals do not migrate often [29] and tend to frequently upgrade
66 their housing by expanding or building new structures on their property, making housing footprint
67 a meaningful proxy for welfare that responds to improved economic conditions. In this study, we
68 focus on building footprint because it can be precisely measured at scale with modern deep learning
69 techniques.

70 Many features of buildings other than footprint are observable with satellite imagery; for exam-
71 ple roof material [30, 31]. One of the main advantages of the method proposed here compared to
72 alternative “black-box” machine learning approaches to measuring wealth that utilize *all* available
73 information contained in satellite images (such as convolutional neural networks [6, 11] or random
74 kitchen sinks [32]) is that it allows the exclusion of subsets of satellite-derived outcomes that may
75 have been directly impacted by the intervention. We show the benefits of this feature of our method
76 in the context of the experiment we evaluate. Specifically, households were eligible for the GiveDi-
77 rectly study as long as their roofing was of low quality (thatched). Due to this eligibility criterion,
78 treatment households were “prompted” to use the GiveDirectly transfer to upgrade their roofing
79 as a way to signal to the experimenters that they had used the cash for good. An improvement of
80 roofs among participating households beyond what would have been expected solely from wealth
81 increases biases estimates of wealth when methods cannot exclude subsets of outcomes. In contrast,
82 it is straightforward for our method to focus exclusively on subsets of available information that
83 were not affected directly (in this case building footprints) while ignoring problematic outcomes
84 (such as roof material) in order to provide unbiased estimates of wealth effects.

85 2 Results

86 We evaluate a development intervention that was conducted in 2014–2017 in 653 villages in rural
87 Kenya [15]. GiveDirectly, a US charity, implemented a randomized controlled trial of unconditional
88 cash transfers to rural households via mobile money, using as sole eligibility criterion whether the
89 household lived under a thatched roof (a low quality roof material that served as a simple means
90 test). Each treatment household received \$1,000—equivalent to about 75% of annual household

91 expenditures—in lump sum, and could spend it however they wished. To evaluate the effectiveness
92 of the program, GiveDirectly randomly selected 328 villages as the treatment group, where eligible
93 households (about 1/3 of the population) received transfers, and used the remaining 325 villages
94 as the control group. The authors conducted extensive household surveys before and after the
95 distribution of the transfers to measure program impacts as is the current practice in the evaluation
96 literature.

97 **Mapping Treatment Intensity and Housing Quality.** To evaluate program impacts, we first
98 construct a map that shows the intensity of the anti-poverty program (hereafter “treatment”) in
99 different geographical units (in this case it is simplest to work with raster grid cells). This geocoded
100 information is obtained from program implementation records, which document where the program
101 was administered. Because of the extremely high granularity of satellite-derived housing quality
102 metrics, it is feasible to study programs that induce fine spatial variation such as household-level
103 randomized trials. Importantly, the variation in treatment intensity has to be either random (if
104 induced by an experiment) or as good as random (in a natural experiment setting), as is the case
105 for any credible program evaluation project.

106 For the GiveDirectly experiment, we construct the treatment intensity map from a local census
107 fielded in 2014–2015, which surveyed all the 65,385 households living in the study area [15]. The
108 census data record each household’s geo-location, and indicate whether they belong to the treat-
109 ment (T), control (C), or out-of-sample (O) group (Figure 1a). Among the three groups, only the
110 treatment households eventually received the cash transfer from GiveDirectly. The control house-
111 holds were randomized into not receiving the transfer, whereas the out-of-sample households were
112 never eligible to participate in the program. Our sample contains 11,055 treatment households and
113 10,682 control households in total. We lay out a regular grid, and count the number of treatment
114 households in each grid cell (Figure 1b). As every transfer was roughly USD 1,000, this variable
115 can be interpreted as the amount of cash infusion (in \$1,000) into a given grid cell, and is our
116 preferred measure of treatment intensity (Figure 1c).

117 Next, we measure housing quality in daytime satellite images with deep learning techniques. The
118 input images are from Google Static Maps [19]. They are taken after the GiveDirectly intervention,
119 have a spatial resolution of about 30cm per pixel, and contain only the RGB (red, green, blue)
120 bands (Figure 1d).

121 To segment buildings, we train a state-of-the-art deep learning model, Mask R-CNN [20], on
122 large, publicly available datasets such as COCO (Common Objects in Context) [33] and Open AI
123 Tanzania [34], as well as a small annotated dataset, which are randomly sampled from all the input
124 images (see Supplementary Materials C for details on model training). The model predictions are
125 highly accurate, both quantitatively (Supplementary Figure S1) and qualitatively (Supplementary
126 Figure S2). The model generalizes well to other countries, such as Mexico, where the number of
127 houses identified in the deep learning predictions is highly correlated with the census population

128 count (Supplementary Figure S9 and Supplementary Materials D). After post-processing, each pre-
129 dicted instance of buildings is represented by a polygon and a “representative” roof color (Figure
130 1e). The Mask R-CNN model conducts instance segmentation (as opposed to semantic segmenta-
131 tion), meaning that it is able to identify every building instance separately, even if they are adjacent
132 to each other. As such, we can measure housing outcomes for each household.

133 We extract two metrics for each built structure: the size of building footprint, and the type of
134 roof material. The roofs are classified into three types: tin roof, thatched roof, and painted roof,
135 based on their color profiles (Supplementary Figure S3). Compared to tin roofs, thatched roofs are
136 generally of lower quality [15, 35]. (Painted roofs are relatively uncommon in the study area.) In
137 prior work, roof reflectance and roof color have been shown to be good proxies of housing quality
138 [30, 31]. As such, we aggregate the total building footprint to measure all housing assets (Figure 1f,
139 Building Footprint), and the footprint of tin-roof buildings to measure high-quality housing assets
140 (Figure 1f, Tin-roof Area), in each grid cell. To obtain night light data for systematic comparison,
141 we download and resample the Visible Infrared Imaging Radiometer Suite (VIIRS) Day/Night
142 Band (DNB) composite images in 2019 [36, 37].

143 The maps of treatment intensity and remotely sensed outcomes for the GiveDirectly experiment
144 are shown in Figure 2. For visual display and privacy protection purposes, we plot the maps with
145 a spatial resolution of 0.005° (roughly 500 meters), which is lower than the resolution used in
146 the subsequent statistical analysis. The experiment generated substantial variation in treatment
147 intensity, as expected (Figure 2a). Both of the housing quality measures capture richer variation
148 in the entire area (Figure 2b, c), whereas the night light data demonstrate little variation in this
149 rural, sparsely populated area, except in a few spots close to local towns (Figure 2d).

150 **Estimating the Program Effects on Housing Quality.** We regress the remotely sensed out-
151 comes on treatment intensity to estimate the causal effects of the GiveDirectly cash transfer. We
152 choose a spatial resolution of 0.001° (approximately 100m), such that most of the grid cells con-
153 tain 0–5 households. We exploit only the experimentally-induced random variation in treatment
154 intensity for identification, and account for pre-determined differences in program eligibility. Intui-
155 tively, consider two grid cells, one containing a household that received the transfer, and the other
156 containing a household that was eligible to get the transfer but did not because it was randomized
157 into the control group. With valid randomization [15], the differences in outcomes between the
158 two can be attributed to the cash transfer. We plot the causal effects on night light and housing
159 quality as cash infusion intensity increases (Figure 3, in color), without making assumptions on
160 the structure of the effects. The results suggest that the effects grow linearly with the amount of
161 cash infusion. We therefore also report an “average” effect, estimated with the assumption that
162 each \$1,000 transfer generates an effect of the same magnitude (Figure 3, panel subtitles). We
163 demonstrate the validity of the empirical strategy further by running 100 placebo simulations—we
164 artificially generate placebo cash transfers that did not actually take place but is consistent with

165 the original randomization design, and estimate their treatment effects (Figure 3, in gray). The
166 resulting estimates are reassuringly centered around zero.

167 We observe statistically significant and economically sizable effects on housing quality, on both
168 the extensive margin (larger building footprint) (Figure 3a), and the intensive margin (higher
169 quality roofs) (Figure 3b). On average, a \$1,000 cash transfer significantly increased building
170 footprint by 7.9 square meters (95% CI: [2.3, 13.5], $t(14, 143) = 2.8, p = 0.006$) or 85.0 square feet,
171 and tin-roof area by 13.6 square meters (95% CI: [9.6, 17.6], $t(14, 143) = 6.7, p < 0.001$) or 146.4
172 square feet. These increases indicate that households may have built new structures—either primary
173 residences or auxiliary structures, such as kitchens and sheds, expanded their existing structures,
174 and/or upgraded their thatched roofs to tin roofs, an improvement that people commonly used the
175 transfer for [35]. These estimates are consistent with the results from extensive field surveys, which
176 also documented large increases in housing asset values [15].

177 On the other hand, we do not observe any program effects on night light (Figure 3c), despite the
178 fact that the cash transfer had large positive impacts on many aspects of the recipient households'
179 economic well-being—food expenditure, consumer durable spending, asset holding, and housing
180 values [15]. The estimated effect is -0.000120 (95% CI: $[-0.008, 0.008]$, $t(14, 143) = -0.03, p = 0.977$)
181 which is not statistically different from zero, is small in magnitude, and actually slightly
182 negative. This may be because of low demand for electrification [38], or the poor sensitivity of
183 night light in low-income, rural regions [6].

184 **Recovering the Program Effects on Economic Well-being with Engel Curves.** We re-
185 cover the program effects on household economic well-being with a canonical economic concept,
186 the Engel curve. Engel curves describe how household expenditures on particular goods or services
187 depend on households' economic well-being. For example, it is widely known that poorer families
188 spend a larger share of their expenditure on food. Engel curves have long been used to infer eco-
189 nomic well-being without needing detailed information on prices as it is straightforward to measure
190 how much of a household's expenditure is spent on food [21–24]. We adapt this concept to housing
191 quality by exploiting the fact that someone who lives in a larger house is likely to be wealthier than
192 someone who lives in a smaller house (Figure 4a). By the same logic, if we observe that someone's
193 house size increased, then we can infer what level of wealth is associated to such a house size—as if
194 they were moving up on the Engel curve. Mathematically, the slope of the Engel curve represents
195 the ratio between the change in house size and the change in wealth. We divide the change in the
196 house size (Figure 3) by the slope of the Engel curve (Figure 4a) to infer the corresponding change
197 in wealth (Figure 4b). Importantly, the validity of this approach depends on the assumption that
198 the Engel curve does not shift in response to the treatment, which could happen due to relative
199 price changes of the good or taste changes.

200 In this study, we derive housing Engel curves from an endline survey of the GiveDirectly trial
201 participants between May 2016 and June 2017, which includes 4,578 geo-coded households who

were eligible for the transfer. Of these households, only those assigned to the control group are used for the estimation. In Figure 4a, we show the relationship between survey-based measures of economic well-being (x -axis) and remotely sensed night light or housing quality measures (y -axis). The Engel curves are estimated with a linear regression (dotted lines). The non-linear fit with LOESS (solid lines) shows only small deviations from the linear regression line, and we cannot reject the null hypothesis that these Engel curves are linear (see Methods). The Engel curves are also roughly monotonically increasing, validating the choice of these variables as wealth proxies.

The Engel curves can be derived from any geo-coded consumption and expenditure survey, as long as the surveyed households are—or can be re-weighted to be—representative of the sample in the previous treatment effect estimation step. Notably, the sample does not necessarily have to include any one who has received the treatment, opening up the possibilities of using existing data sources (such as the Living Standards Measurement Study (LSMS)) to estimate Engel curves. We demonstrate this by comparing the Engel curves derived from two distinct samples: the households who were deemed eligible to receive the cash transfers (meaning that they used to live in thatched-roof houses), and households who were not. While all the households live in the same area in western Kenya, the ineligible households are generally wealthier than the eligible ones. Their Engel curves, however, are similar within the same range of wealth (Supplementary Figure S8).

We scale program effects on each remotely sensed outcome by the Engel curve slope to estimate the impacts of the GiveDirectly transfer on household wealth, measured by aggregating the values of a variety of assets as measured with household surveys. In Figure 4b, we compare the satellite-derived estimates against the survey-based estimates, which are computed from rich end-line household survey data and taken from Table 1, Column 1 in the original paper [15]. As can be seen, the estimate based on building footprint (USD 425 PPP, 95% CI: [61, 788]) is informative and very close to the survey based estimate (USD 556 PPP, 95% CI: [485, 626]). For reference, the entire GiveDirectly cash transfer is worth USD 1,871 PPP (USD 1,000 nominal). Note that the estimate based on night light is slightly negative and imprecise, and both the upper and lower bounds are uninformative. In contrast, the estimate based on tin-roof area is about two times as large as the survey-based estimate. The results are qualitatively similar when we distinguish between housing asset (Supplementary Figure S4) and non-housing asset (Supplementary Figure S5), or when we use annual consumption expenditure as the alternative measure of economic well-being (Supplementary Figure S6).

Why is the estimate based on the tin-roof area much larger than the survey based estimate? We argue this is due to the violation of a key assumption, which is that the Engel curve used to estimate changes in wealth cannot change directly in response to the treatment—only through its wealth effects. To give intuition for why this matters, consider a program that directly gives people food. In such a case we can no longer look at food consumption to infer program effects on economic well-being, because the relationship between the food and income will be altered directly

239 by the program and households will “look” wealthier than they really are based on their food
240 consumption. More relevant for impact evaluation using satellite data, this example is analogous
241 to examining the impacts a program that provides roads to a region. One would need to exclude
242 the program roads themselves contained in satellite images and look at other correlates of welfare
243 to estimate impacts of such a roads program in an unbiased manner. In the GiveDirectly case,
244 only households that lived in thatched-roof houses were eligible for the study. Households’ usual
245 consumption patterns of high-quality tin roofs might have been affected by this eligibility criteria.
246 One can observe that treatment households owned more tin-roof buildings compared to control
247 households with the same amount of wealth (Supplementary Figure S7). This may have been a
248 result of households interpreting the treatment as a “labelled” cash transfer [39].

249 These results highlight the importance of using interpretable proxies when evaluating programs
250 with machine learning predictions. An emerging literature is making great progress in mapping
251 poverty with satellite imagery and machine learning with a high spatial granularity at scale [6–13].
252 Typically, a machine learning model first learns the mapping between the input satellite images
253 and the ground truth labels of wealth or consumption expenditure, assembled from geo-coded
254 household surveys. Then, the model generates predicted poverty maps for every region in the
255 sample, including those with no survey coverage. The model implicitly combines and executes
256 two tasks: (1) extracting semantically meaningful observations of, say, housing quality, agricultural
257 productivity, or infrastructure, from raw satellite images; and (2) inferring economic well-being from
258 observing the consumption patterns of these private or public goods (similar to the Engel curve
259 analysis in this study). While the flexibility of the machine learning models helps improve predictive
260 performance, the difficulty in interpretation makes it almost impossible to know or constrain what
261 private or public goods are identified and utilized by the model. Since black-box machine learning
262 models utilize as much information as possible from the input satellite images, it is very likely
263 that the Engel curves of at least some of the observed goods will change (similarly to the tin-roof
264 area variable in this study), introducing biases in the estimated program effects. In this study, we
265 disentangle the two tasks, so that the first task can be framed as a traditional object detection
266 and segmentation task, allowing us to leverage extensive research in computer science; and the
267 second task becomes more transparent, explicit, and the assumptions testable (for example, with
268 Supplementary Figure S7).

269 **3 Discussion**

270 This paper provides compelling evidence that RCT program evaluations aimed at improving house-
271 hold welfare can be obtained solely based on satellite imagery and deep learning methods. This
272 approach has the advantage of being inexpensive and timely, suggesting great promise as a com-
273 plement and in some cases as a substitute to in-person survey data collection methods.

274 However, it bears noting that a fundamental limitation to evaluating programs based on satellite

²⁷⁵ imagery is that in order to be measurable from space, programs being evaluated have to generate
²⁷⁶ impacts on the built landscape. This prevents applicability to programs targeted at addressing
²⁷⁷ development challenges that are unlikely to impact the built environment such as improved teaching
²⁷⁸ methods at schools. Another limitation is that welfare is a household or individual concept whereas
²⁷⁹ satellite images capture characteristics about a place. Mapping household welfare to housing as
²⁸⁰ we do here requires a tight mapping between structures and households through limited mobility.
²⁸¹ While migration rates are very low in the GiveDirectly study area [15], this may be a challenge for
²⁸² programs that impact mobility, such as transportation infrastructure programs.

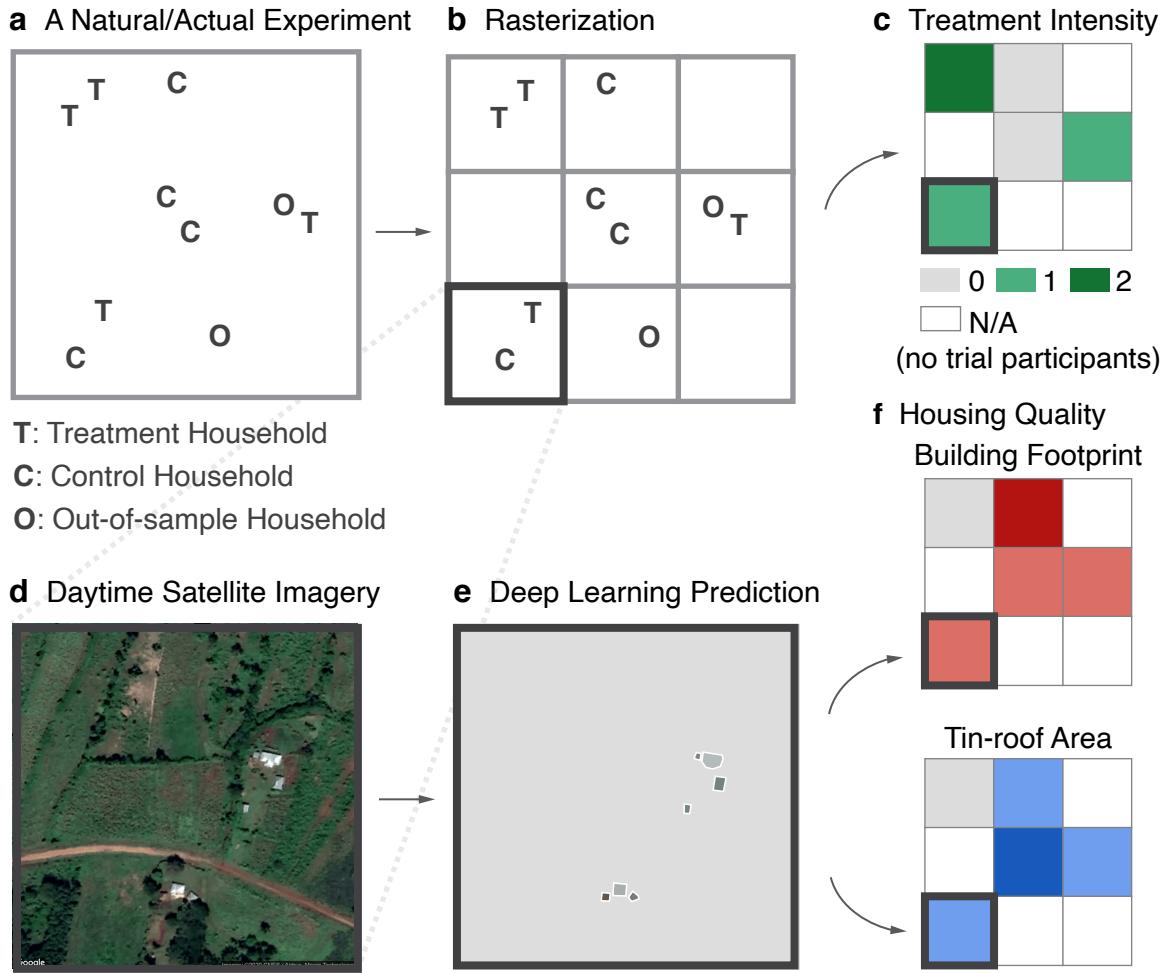


Figure 1: **Constructing maps of treatment intensity and remotely sensed outcomes from program implementation records and satellite imagery.** **a** An illustration of geocoded program implementation records. **b** Placing a regular grid over **a** and measuring the intensity of the treatment in each grid cell. **c** Constructed raster of the number of treatment households in each grid cell. **d** An example daytime satellite image from Google Static Maps. **e** Example deep learning predictions on **d**. Each building is outlined in white and filled with the “representative” roof color. **f** Constructed rasters of remotely sensed housing quality outcomes. In **c** and **f**, grid cells without trial participants are omitted and shown in white.

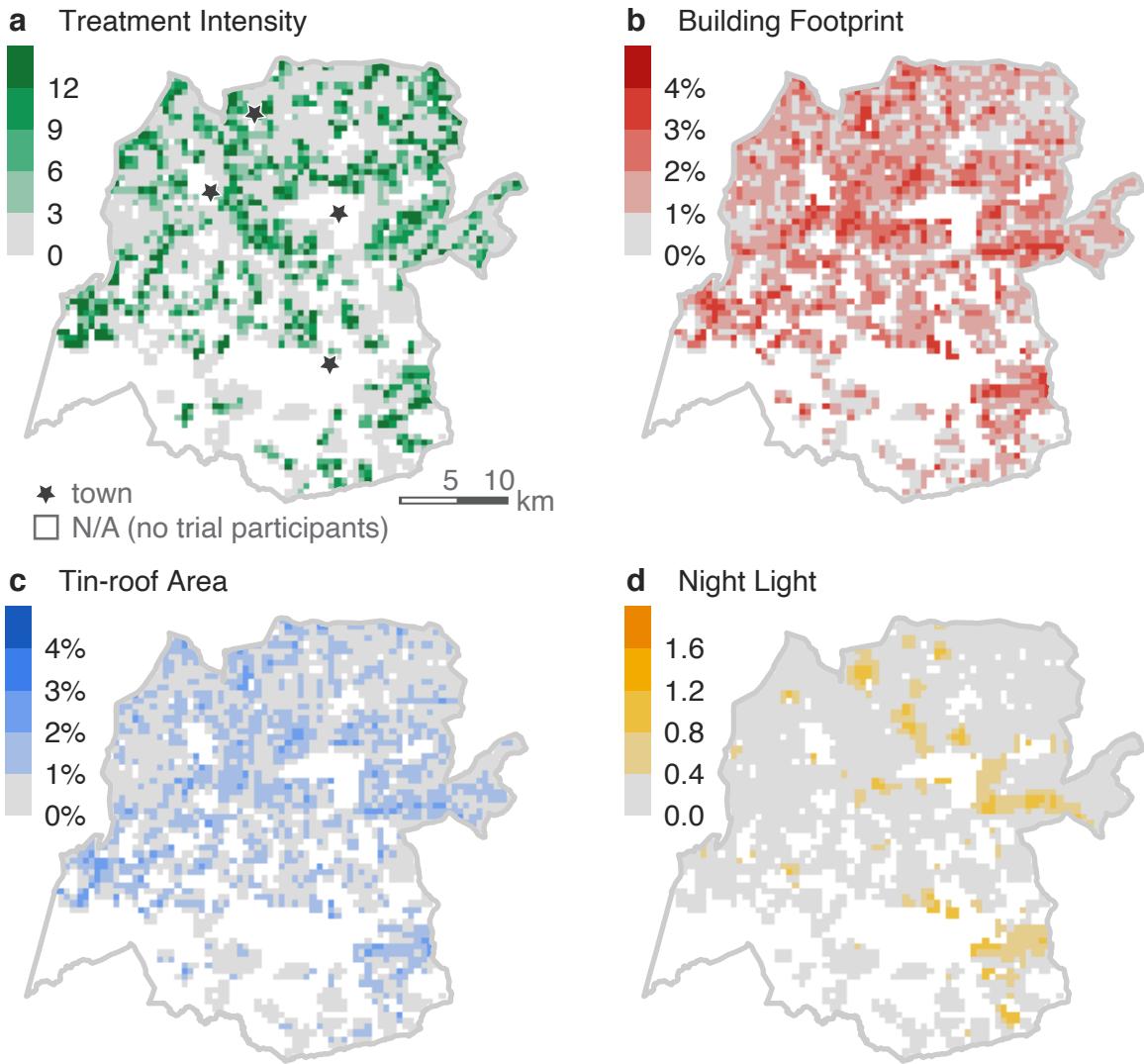


Figure 2: Mapping treatment intensity and remotely sensed outcomes in the GiveDirectly study area in 2019. **a** Treatment intensity represents the number of households who received a \$1,000 cash transfer from GiveDirectly. **b** Building footprint measures the total area covered by any building, shown as a percentage of the total area. **c** Tin-roof area measures the total footprint of buildings with roofs made of tin (a high quality construction material), shown as a percentage of the total area. **d** Night light is the average radiance in the Visible Infrared Imaging Radiometer Suite (VIIRS) Day/Night Band (DNB). In all the panels, the gray lines outline the GiveDirectly study area in Siaya, Kenya. Grid cells without trial participants are omitted and shown in white. $n = 2,501$.

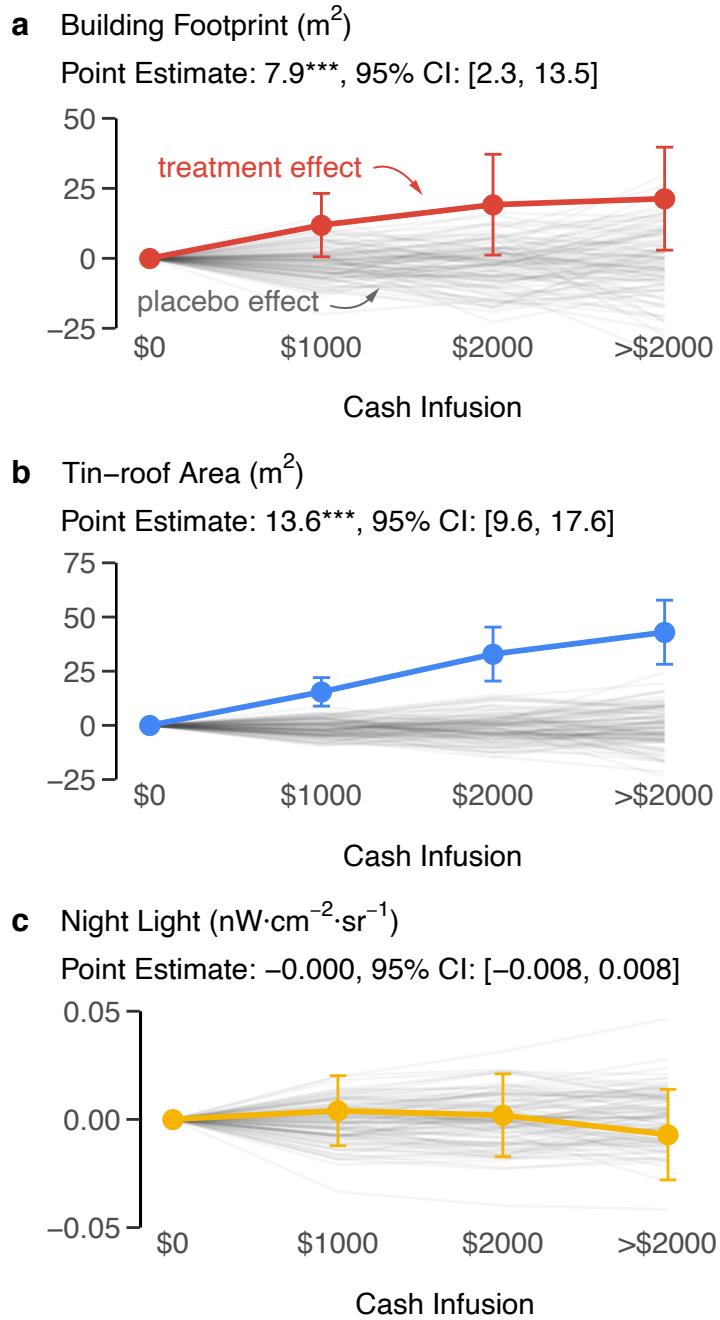
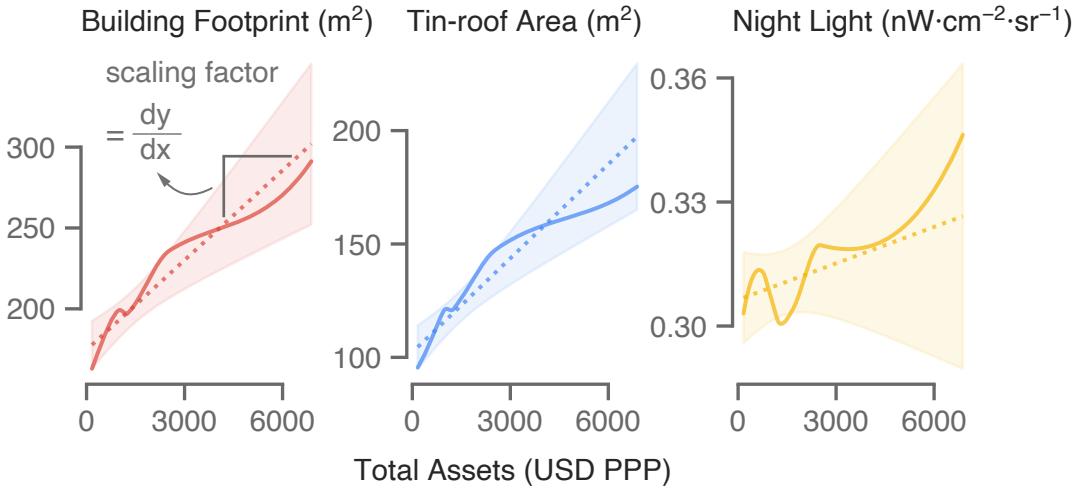


Figure 3: **Housing quality increased in response to the GiveDirectly cash transfer, but night light remained unchanged.** The treatment effects of the cash transfers on building footprint (a), tin-roof area (b), and night light (c) are shown in color. The dots represent the point estimates, and the error bars represent the 95% confidence intervals. Gray lines show the estimated effects of the placebo cash infusions from 100 simulations. The panel subtitles report the average treatment effect of a \$1,000 transfer and the 95% confidence intervals, assuming constant effect. *** indicates statistical significance at the 1% level for a two-sided t-test. $n = 14,155$.

a Engel Curves



b Treatment Effect Estimates on Total Assets

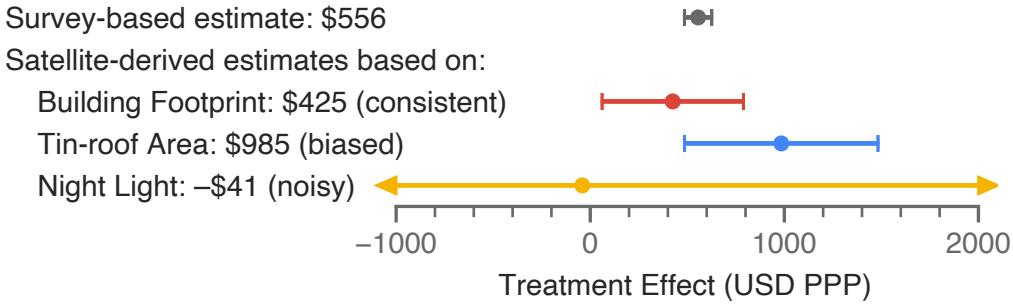


Figure 4: **The treatment effect on total assets can be correctly recovered by scaling the effect on building footprint.** **a** The Engel curves of building footprint, tin-roof area, and night light, estimated with LOESS (solid line) or a linear regression (dotted line). The shaded regions represent the 95% confidence intervals for the latter. **b** Comparing the survey-based versus satellite-derived treatment effects. The dots show the point estimates. The error bars show the 95% confidence intervals, with the arrow(s) marking upper/lower bounds that are out of range (if any). $n = 1,844$.

283 **References**

- 284 1. Deaton, A.
285 *The analysis of household surveys: a microeconometric approach to development policy*
286 (The World Bank, 1997).
- 287 2. Banerjee, A. V. & Duflo, E.
288 *Poor economics: A radical rethinking of the way to fight global poverty* (Public Affairs, 2011).
- 289 3. Pamies-Sumner, S. *Development Impact Evaluations: State of Play and New Challenges*
290 tech. rep. (Agence Française de Développement, 2015).
- 291 4. Alix-Garcia, J. M., Sims, K. R. & Costica, L.
292 Better to be indirect? Testing the accuracy and cost-savings of indirect surveys.
293 *World Development* **142**, 105419 (2021).
- 294 5. Brune, L., Karlan, D., Kurdi, S. & Udry, C. R. *Social Protection Amidst Social Upheaval: Examining the Impact of a Multi-Faceted Program for Ultra-Poor Households in Yemen*
295 tech. rep. (National Bureau of Economic Research, 2020).
- 297 6. Jean, N. *et al.* Combining satellite imagery and machine learning to predict poverty.
298 *Science* **353**, 790–794 (2016).
- 299 7. Blumenstock, J. E. Fighting poverty with data. *Science* **353**, 753–754 (2016).
- 300 8. Engstrom, R., Hersh, J. & Newhouse, D. *Poverty from space: using high-resolution satellite imagery for estimating economic well-being* tech. rep. (The World Bank, 2017).
- 302 9. Babenko, B., Hersh, J., Newhouse, D., Ramakrishnan, A. & Swartz, T.
303 *Poverty Mapping Using Convolutional Neural Networks Trained on High and Medium Resolution Satellite Images, With an Application in Mexico*
304 in *Proceedings of NIPS 2017 Workshop on Machine Learning for the Developing World*
305 (2017). <https://arxiv.org/abs/1711.06323>.
- 307 10. Watmough, G. R. *et al.*
308 Socioecologically informed use of remote sensing data to predict rural household poverty.
309 *Proceedings of the National Academy of Sciences* **116**, 1213–1218 (2019).
- 310 11. Yeh, C. *et al.* Using publicly available satellite imagery and deep learning to understand
311 economic well-being in Africa. *Nature communications* **11**, 1–11 (2020).
- 312 12. Aiken, E. L., Bedoya, G., Coville, A. & Blumenstock, J. E.
313 *Targeting Development Aid with Machine Learning and Mobile Phone Data: Evidence from an Anti-Poverty Intervention in Afghanistan* Unpublished. 2020.
- 315 13. Blumenstock, J. Machine learning can help get COVID-19 aid to those who need it most.
316 *Nature* (2020).

- 317 14. Suri, T. Mobile money. *Annual Review of Economics* **9**, 497–520 (2017).
- 318 15. Egger, D., Haushofer, J., Miguel, E., Niehaus, P. & Walker, M. W.
319 *General equilibrium effects of cash transfers: experimental evidence from Kenya* tech. rep.
320 (National Bureau of Economic Research, 2019).
- 321 16. Blattman, C. & Niehaus, P. Show them the money: Why giving cash helps alleviate poverty.
322 *Foreign Affairs* **93**, 117–126 (2014).
- 323 17. Alix-Garcia, J., McIntosh, C., Sims, K. R. & Welch, J. R. The ecological footprint of poverty
324 alleviation: evidence from Mexico’s Oportunidades program.
325 *Review of Economics and Statistics* **95**, 417–435 (2013).
- 326 18. Jayachandran, S. *et al.* Cash for carbon: A randomized trial of payments for ecosystem
327 services to reduce deforestation. *Science* **357**, 267–273 (2017).
- 328 19. Google Static Maps. *Google Static Maps*
329 <https://developers.google.com/maps/documentation/maps-static/intro>.
330 accessed 6 May 2020. 2020.
- 331 20. He, K., Gkioxari, G., Dollár, P. & Girshick, R. *Mask R-CNN*
332 in *Proceedings of the IEEE international conference on computer vision* (2017), 2961–2969.
- 333 21. Elbers, C., Lanjouw, J. O. & Lanjouw, P. Micro-level estimation of poverty and inequality.
334 *Econometrica* **71**, 355–364 (2003).
- 335 22. Tarozzi, A. & Deaton, A.
336 Using census and survey data to estimate poverty and inequality for small areas.
337 *The review of economics and statistics* **91**, 773–792 (2009).
- 338 23. Young, A. The African growth miracle. *Journal of Political Economy* **120**, 696–739 (2012).
- 339 24. Atkin, D., Faber, B., Fally, T. & Gonzalez-Navarro, M.
340 *A New Engel on Price Index and Welfare Estimation* tech. rep.
341 (National Bureau of Economic Research, 2020).
- 342 25. Henderson, J. V., Storeygard, A. & Weil, D. N.
343 Measuring economic growth from outer space. *American economic review* **102**, 994–1028
344 (2012).
- 345 26. Chen, X. & Nordhaus, W. D. Using luminosity data as a proxy for economic statistics.
346 *Proceedings of the National Academy of Sciences* **108**, 8589–8594 (2011).
- 347 27. Michalopoulos, S. & Papaioannou, E.
348 National institutions and subnational development in Africa.
349 *The Quarterly journal of economics* **129**, 151–213 (2014).
- 350 28. OECD. in *National Accounts at a Glance 2014* (OECD Publishing, Paris, 2014).

- 351 29. Murrugarra, E., Larrison, J. & Sasin, M.
352 *Migration and poverty: Towards better opportunities for the poor*
353 (World Bank Publications, 2010).
- 354 30. Marx, B., Stoker, T. M. & Suri, T.
355 There is no free house: Ethnic patronage in a Kenyan slum.
356 *American Economic Journal: Applied Economics* **11**, 36–70 (2019).
- 357 31. Michaels, G. *et al.*
358 *Planning Ahead for Better Neighborhoods: Long Run Evidence from Tanzania* tech. rep.
359 (IZA Institute of Labor Economics, 2017).
- 360 32. Rolf, E. *et al.*
361 *A Generalizable and Accessible Approach to Machine Learning with Global Satellite Imagery*
362 tech. rep. (National Bureau of Economic Research, 2020).
- 363 33. COCO. *COCO - Common Objects in Contexts* <http://cocodataset.org>. accessed 6 May 2020.
364 2020.
- 365 34. Open AI Tanzania. *2018 Open AI Tanzania Building Footprint Segmentation Challenge*
366 <https://competitions.codalab.org/competitions/20100>. accessed 6 May 2020. 2020.
- 367 35. Haushofer, J. & Shapiro, J. The short-term impact of unconditional cash transfers to the
368 poor: experimental evidence from Kenya.
369 *The Quarterly Journal of Economics* **131**, 1973–2042 (2016).
- 370 36. Google Earth Engine.
371 *VIIRS Stray Light Corrected Nighttime Day/Night Band Composites Version 1*
372 [https://developers.google.com/earth-
373 engine/datasets/catalog/NOAA_VIIRS_DNB_MONTHLY_V1 VCMSLCFG](https://developers.google.com/earth-engine/datasets/catalog/NOAA_VIIRS_DNB_MONTHLY_V1 VCMSLCFG).
374 accessed 6 May 2020. 2020.
- 375 37. Elvidge, C. D., Baugh, K., Zhizhin, M., Hsu, F. C. & Ghosh, T. VIIRS night-time lights.
376 *International Journal of Remote Sensing* **38**, 5860–5879 (2017).
- 377 38. Lee, K., Miguel, E. & Wolfram, C.
378 Experimental evidence on the economics of rural electrification.
379 *Journal of Political Economy* **128**, 1523–1565 (2020).
- 380 39. Benhassine, N., Devoto, F., Duflo, E., Dupas, P. & Pouliquen, V.
381 Turning a shove into a nudge? A “labeled cash transfer” for education.
382 *American Economic Journal: Economic Policy* **7**, 86–125 (2015).

383 **4 Methods**

384 **Constructing the Treatment Intensity Map.** To construct the treatment intensity map, we
385 utilize data from a baseline census, which was conducted by the authors of the original paper in
386 2014–2015. The census identified all 65,385 households (roughly 280,000 people) residing in 653
387 villages in the study area, recorded their GPS coordinates, whether each household was eligible
388 for the GiveDirectly cash transfer, and whether they had been randomized into the treatment or
389 control group [1]. To address the measurement errors of the GPS collection devices, we discard
390 58 outliers (living more than 2 kilometers away from the village centers) and impute those and
391 other 4 missing GPS coordinates with village center coordinates. Then, we convert these household
392 records into a raster map. We lay out a regular grid, and count, in each grid cell, the number of
393 households that ultimately received the GiveDirectly cash transfer (see Figure 1 and Figure 2a).
394 Grid cells containing no eligible households are excluded. To account for pre-determined policy
395 intensity differences, we record (and later control for) the number of households that were eligible
396 for the cash transfer, regardless of whether they had been randomized into the treatment or control
397 group.

398 **Obtaining High-resolution Daytime Satellite Images.** We utilize high-resolution daytime
399 satellite images from Google Static Maps [2]. These images have a spatial resolution of about
400 30cm per pixel (at equator), and contain only the RGB (red, green, blue) bands (see Figure 1d and
401 Supplementary Figure S2 for examples). These images come from a variety of commercial providers
402 such as Maxar (formerly DigitalGlobe) and Airbus, and have been seamlessly mosaicked together.
403 They have also been geo-referenced and pre-processed to remove clouds and address other data
404 quality issues. Google does not provide the exact timestamps for these images, but we estimate
405 that they were taken in 2019, most likely on Dec 30, 2019. The dates for retrieving these images
406 from the Google Static Maps API are between Feb 19 and Feb 21, 2020, and the Google Earth Pro
407 imagery archive reflects that the closest available images in the study area were from Dec 30, 2019.
408 Multiple other satellite images taken in February, March, July, August and September 2019 are
409 also available in the study area, indicating that the images used in this study are most certainly
410 from 2019.

411 **Extracting Housing Quality Metrics with Mask R-CNN.** We first leverage a state-of-the-
412 art deep learning model, Mask R-CNN [3] to segment buildings—that is, to detect each building
413 and the pixels that they occupy—in the Google Static Map satellite images. We then convert the
414 pixel-wise predictions to polygons, and extract housing quality metrics related to the size of the
415 building and the roof materials from each polygon (see Figure 1e and Supplementary Figure S2 for
416 examples).

417 Loosely speaking, the Mask R-CNN model operates as follows. First, the model proposes a

418 large number of “regions of interest”, each of which potentially contains a building. Then, the
419 model uses convolutional filters to identify patterns within the proposed region that are indicative
420 of the presence of buildings, such as the sharp edges, the highly reflective roofs, and the building
421 shadows. Finally, the model predicts whether each proposed region contains a building, as well as
422 whether each pixel is occupied by the building.

423 We train the Mask R-CNN model with a multi-step process and a transfer learning framework,
424 as described in greater detail in Supplementary Materials C. Publicly available building footprint
425 datasets in rural and low-income regions are rare, and they often differ substantially in spatial
426 resolution, sensor instrument, and landscape from inference images (that is, the target images
427 that the model will make predictions for). Relying solely on publicly available training data is
428 therefore insufficient for achieving satisfactory predictive performance. We curate a set of in-sample
429 annotations by randomly sampling 120 images from all the Google Static Map images in the study
430 area, and manually creating high-quality building footprint annotations for them. We pre-train
431 the Mask R-CNN model on large, publicly available datasets such as COCO (Common Objects in
432 Context) and Open AI Tanzania, and fine-tune them on this set of in-sample annotations.

433 The model predictions are highly accurate. The overall F1 score (a standard performance metric
434 for instance segmentation) on a random subset of inference images is 0.79 (Supplementary Figure
435 S1). The F1 score is the harmonic mean of precision (the proportion of model-identified buildings
436 that are actual buildings) and recall (the proportion of actual buildings that are correctly identified
437 by the model). Here, a building is deemed to be correctly identified if the predicted pixel mask and
438 the ground truth pixel mask have sufficient overlap (more precisely, if the intersection of the two
439 masks is more than 50% of the union of the two masks). As a reference point, the top winner in
440 the 2nd SpaceNet building footprint extraction competition reported an F1 score of 0.69 [4]. This
441 demonstrates that the Mask R-CNN model used in this study performs well, although building
442 footprint segmentation in rural, less complex scenes is generally easier than in modern cities so
443 these metrics are not directly comparable.

444 We post-process the model-predicted pixel masks by converting them to polygons, and simplifying
445 the polygons with the Douglas-Peucker algorithm with a pixel tolerance of 3. For each polygon,
446 we compute two housing quality metrics: building footprint and type of roof materials. We then
447 lay out a regular grid, assign each building to grid cells based on the centroids of the polygons, and
448 aggregate to obtain two metrics at the pixel level: building footprint (Figure 2b) and tin-roof area
449 (Figure 2c).

450 First, we measure the size of each building polygon and convert it to square meters. We correct
451 for area distortion, which is induced by the Web Mercator projection system that the Google Static
452 Map uses. This metric may appear larger than what one expects for the size of homes in a low-
453 income context (Figure 4), because (1) it represents the footprint of the entire building, which
454 is typically larger than the size of the livable area; and (2) it accounts for both residential and

455 non-residential structures, since the model is not able to distinguish between the two.

456 Second, we estimate the types of roof materials based on the colors of the roofs, and compute
457 the footprint of tin-roof buildings in each grid cell. For each building, we take all the pixels
458 associated with the given building instance, and assign a “representative” roof color by computing
459 the average values in the RGB (Red, Green, Blue) channels. Since the Euclidean distances between
460 color vectors in the RGB color space does not reflect perceptual differences, we project all the RGB
461 color vectors to the CIELAB color space, and cluster these roof color vectors into 8 groups by
462 running the K-means clustering algorithm. We further classify these 8 groups into three types of
463 roof materials: tin roof, thatched roof, and painted roof (Supplementary Figure S3), and compute
464 the total footprint of tin-roof buildings.

465 **Obtaining the Night Light Data.** To measure nighttime luminosity, we use the Visible Infrared
466 Imaging Radiometer Suite (VIIRS) Day/Night Band (DNB) composite images hosted on Google
467 Earth Engine [5, 6]. The VIIRS-DNB data product excludes areas impacted by cloud cover and
468 correct for stray light [7]. However, it has not been filtered to screen out lights from aurora,
469 fires, boats, and other temporal lights, and lights are not separated from background (non-light)
470 values [5]. This data product has a native spatial resolution of 15 arc seconds (approximately 463
471 meters at the equator), and we resample the data by conducting nearest neighbor interpolation
472 when necessary. We average over all the monthly observations in 2019 and construct a single cross
473 sectional observation, to reduce seasonality effects and for consistency with the daytime satellite
474 imagery (Figure 2d). The VIIRS-DNB data product is considered superior to the more widely used
475 night light data, DMSP-OLS (the United States Air Force Defense Meteorological Satellite Program,
476 Operational Linescan System) because it preserves finer spatial details, has a lower detection limit
477 and displays no saturation on bright lights [8]. This ensures that we conduct a fair comparison
478 with the most modern and high-quality night light data product.

Estimating the Program Effects on Housing Quality. The main econometric specification
for Figure 3 is as follows

$$y_i = \sum_{k \in K} \tau_k \mathbf{1}\{x_i = k\} + \sum_{m \in M} \beta_m \mathbf{1}\{e_i = m\} + \epsilon_i \quad (1)$$

479 where each observation i represents a $0.001^\circ \times 0.001^\circ$ grid cell (approximately 100m \times 100m); τ_k
480 represents the estimate of interest: the treatment effects of the unconditional cash transfer on
481 remotely sensed outcomes; x_i denotes the number of recipient households per grid cell (equivalent
482 to the amount of cash infusion in \$1000); e_i denotes the number of eligible households per grid
483 cell, with $m \in M = \{0, 1, 2, 3, \dots\}$; and y_i denotes remotely sensed outcomes: night light, building
484 footprint, and tin-roof area. To account for pre-existing differences in population density or wealth,
485 which may cause non-random variation in treatment intensity, we flexibly control for the number

486 of eligible households per grid cell, and exclude grid cells with no eligible households. Because the
 487 grid cells are fairly small and the number of observations for $k > 2$ is small, we bin the number of
 488 recipient households into four bins $k \in K = \{0, 1, 2, 2+\}$, to preserve statistical power. Standard
 489 errors are calculated à la Conley, with a uniform kernel and a 3km cutoff [9–12]. To reduce the
 490 effects of outliers (due to sensor malfunctioning or machine learning model prediction errors), we
 491 winsorize all remotely sensed variables at the 99 percentile.

492 We run 100 placebo simulations to further demonstrate the validity of the main specification.
 493 In each simulation, we randomly assign half of the 68 groups of villages to the high-saturation
 494 group, and the other half to the low-saturation group. In the high-saturation groups, we randomly
 495 assign 2/3 of the villages to the treatment group (and the rest to the control group); whereas in
 496 the low-saturation group, we assign only 1/3 of the villages to the treatment group (and the rest to
 497 the control group). This mimics the two-tier randomization scheme of the original trial [1]. Using
 498 these simulated placebo treatment status variables, we estimate the placebo treatment effects with
 499 the econometric specification described in Equation 1.

To compute a single pooled treatment effect, we make an assumption of linear treatment
 effects—every transfer of \$1,000 has an effect of the same magnitude, regardless of the treatment
 intensity in that geographical area. The resulting econometric specification is as follows

$$y_i = \tau x_i + \sum_{m \in M} \beta_m \mathbf{1}\{e_i = m\} + \epsilon_i \quad (2)$$

500 where τ is the “average” treatment effect, and all else remain the same as in Equation 1. We
 501 conduct two-sided t-tests to assess statistical significance.

Estimating the Engel Curves. An Engel curve describes how household expenditure on a
 particular good varies with income—a relationship that can be used to infer households’ economic
 well-being from the consumption patterns of a limited subset of goods [13–16]. The mathematical
 formulation is

$$Q_{hp} = F_p(W_h) + \epsilon_{hp} \quad (3)$$

where household h with W_h wealth (or other measures of economic well-being) would consume Q_{hp}
 quantities of a normal good p , and $F_p(\cdot)$ represents the Engel curve for product p in the population.
 With a linearity assumption, this can be simplified to be

$$Q_{hp} = \alpha_p + \beta_p W_h + \epsilon_{hp} \quad (4)$$

502 where α_p is the intercept and β_p is the slope of a linear Engel curve.

503 In this study, we estimate the Engel curves—the relationships between remotely sensed metrics

504 and survey-based measures of economic well-being—based on the endline survey of the original
505 GiveDirectly trial, which includes a representative set of 4,578 geo-coded households who were
506 eligible for the transfer. The households participated in a comprehensive consumption and expen-
507 diture survey between May 2016 and June 2017, after the distribution of cash transfers. From the
508 surveys, we observe annualized household consumption expenditure, and asset values. Household
509 consumption expenditure is the annualized sum of total food consumption in the last 7 days, fre-
510 quent purchases in the last month, and infrequent purchases over the last 12 months. Household
511 assets include housing and non-housing assets, but not land values. Housing asset values are mea-
512 sured as the respondent’s self-reported cost to build a home like theirs. Non-housing assets include
513 livestock, transportation (bicycles, motorcycles, and cars), electronics, farm tools, furniture, other
514 home goods, and lending or borrowing from formal or informal sources. We do not study land
515 values because they are difficult to value given thin local markets [1].

516 We perform heuristic matching between the buildings and the household survey GPS coordi-
517 nates, to link variables in the survey with remotely sensed variables. First, we take the baseline
518 census data, which geo-coded every single household who lived in the study area, and assign every
519 building in the satellite images to its closest census GPS coordinate, if the distance between the
520 two was within 250m. This ensures that every building is matched to at most one household.
521 Second, we match GPS coordinates from the survey with GPS coordinates from the census. While
522 the same household supposedly had the same geo-location, these two often differed because of the
523 measurement errors of the GPS collection devices, and because the coordinates might be recorded
524 anywhere on the participants’ plots and not necessarily in their primary residence. We similarly
525 assign each survey GPS coordinate to its closest census GPS coordinate, if the distance between the
526 two was within 250m. In cases of multiple surveys being assigned to the same census coordinate,
527 we keep the closest survey. The final sample contains only census observations that are matched
528 with both buildings in the satellite images and survey records, and consists of 1,904 treatment
529 households and 1,844 control households.

530 The Engel curves are estimated with only the control group (Figure 4a and Supplementary
531 Figure S4a, S5a and S6a). They are estimated both non-linearly with LOESS (see Equation 3 and
532 the solid lines in Figure 4a) and linearly (see Equation 4 and the dotted lines in Figure 4a). When
533 fitting LOESS, we allow for locally-fitted quadratic polynomials, and use 75% of the data points
534 for each fit. We test for the non-linearity of the Engel curves in a separate procedure. We first
535 run a linear regression, take the residuals, and fit the residuals with a natural (cubic) spline with 5
536 knots. We then conduct a two-sided F-test on the coefficients of the natural spline basis, and reject
537 the null hypothesis (linearity) if these coefficients are jointly significant. We cannot reject linearity
538 for any of the three proxies in Figure 4 (building footprint: $F(1, 838) = 0.37, p = 0.829$; tin-roof
539 area: $F(1, 838) = 0.79, p = 0.533$; night light: $F(1, 838) = 0.39, p = 0.814$). To minimize the
540 influence of outliers, we winsorize annual expenditure, housing assets, non-housing assets and total

541 assets at the 1 and 99 percentile of the eligible and non-eligible sample, respectively. We winsorize
 542 at the 1 percentile as outliers with a large amount of debt exist and could potentially drive the
 543 results otherwise. We similarly winsorize all the remotely sensed variables at the 99 percentile for
 544 the eligible and non-eligible sample. We exclude a small number of renters who do not own any
 545 housing assets (31 treatment households, 32 control households, and 55 ineligible households), to
 546 simplify the interpretation of the Engel curves.

Recovering the Program Effects on Economic Well-being. We adapt a prior mathematical formulation that uses the Engel curve to infer changes in economic well-being [15]. Suppose that one is interested in studying the effect of a plausibly exogenous treatment Z on, say, wealth W (denoted $\hat{\tau}_W$), but can only inexpensively observe its effect on the consumption of product p (denoted $\hat{\tau}_{Q_p}$). Recall that $\hat{\beta}_p$ is the estimated slope of the linear Engel curve in Equation 4, then

$$\hat{\tau}_W = \hat{\tau}_{Q_p} / \hat{\beta}_p \quad (5)$$

Using a formula for propagation of error (or the multivariate Delta method), one can derive the standard error for $\hat{\tau}_W$ as follows. This derivation is based on prior work [15], but additionally accounts for the precision of the slope of the Engel curve.

$$\left(\frac{\hat{\sigma}(\hat{\tau}_W)}{\hat{\tau}_W} \right)^2 = \left(\frac{\hat{\sigma}(\hat{\tau}_{Q_p})}{\hat{\tau}_{Q_p}} \right)^2 + \left(\frac{\hat{\sigma}(\hat{\beta}_p)}{\hat{\beta}_p} \right)^2 \quad (6)$$

547 A key assumption of this approach is that $\hat{\beta}_p$ does not depend on Z —that is, the Engel curve
 548 does not change in direct response to the treatment—also termed the conditional independence
 549 assumption [14].

550 We estimate the treatment effects on wealth (or other measures of economic well-being) ac-
 551 cording to Equation 5 and Equation 6, with the treatment effect estimates for remotely sensed
 552 variables, and the slopes of the Engel curves. We compare the satellite-derived estimates against
 553 the survey-based estimates, taken from Table 1, Column 1 in the original paper [1], which were
 554 based on the endline household survey data (Figure 4b).

555 **End Notes**

556 **Data and Code Availability**

557 This paper makes use of restricted access data, which contain personally identifying information of
558 survey participants. Satellite images used in the analyses come from the Google Static Maps API
559 at <https://developers.google.com/maps/documentation/maps-static/overview>, and redis-
560 tribution is not possible. However, de-identified data necessary to reproduce all the figures and
561 statistical analyses are freely available at <https://github.com/luna983/beyond-nightlight>. All
562 the codes are available at <https://github.com/luna983/beyond-nightlight>.

563 **Acknowledgements**

564 We thank Edward Miguel, Jeremy Magruder, Ben Faber, Marshall Burke, Joshua Blumenstock,
565 Supreet Kaur, Ethan Ligon, Elisabeth Sadoulet, Alain De Janvry, Aprajit Mahajan, the partici-
566 pants in the AGU Fall Meeting 2019 (Session GC34C), the UC Berkeley Trade Lunch, Development
567 Workshop, Development Lunch, and Good Data Seminar for feedback. We thank Edward Miguel,
568 Michael Walker, Dennis Egger, Johannes Haushofer, Paul Niehaus, and the rest of the GiveDirectly
569 team, for generously sharing the dataset with us and responding to our inquiries.

570 **Author Contributions**

571 L.Y.H. initiated the project, assembled the data, performed the empirical analyses, created the
572 figures, and wrote the paper. L.Y.H., S.H. and M.G.N. collaboratively iterated on the project idea
573 and result interpretation, and edited the paper.

574 **Ethics Declaration**

575 The GiveDirectly randomized controlled trial and field survey received IRB approval from Maseno
576 University and the University of California, Berkeley. The AEA Trial Registry RCT ID is AEARCTR-
577 0000505. Informed consent was obtained from all human research participants.

578 The authors declare no conflicts of interest.

579 **Additional Information**

580 Supplementary Information is available for this paper. Correspondence and requests for materials
581 should be addressed to Luna Yue Huang (yue_huang@berkeley.edu). Reprints and permissions
582 information are available at www.nature.com/reprints.

583 **Supplementary Information**

584 **A Supplementary Figures**

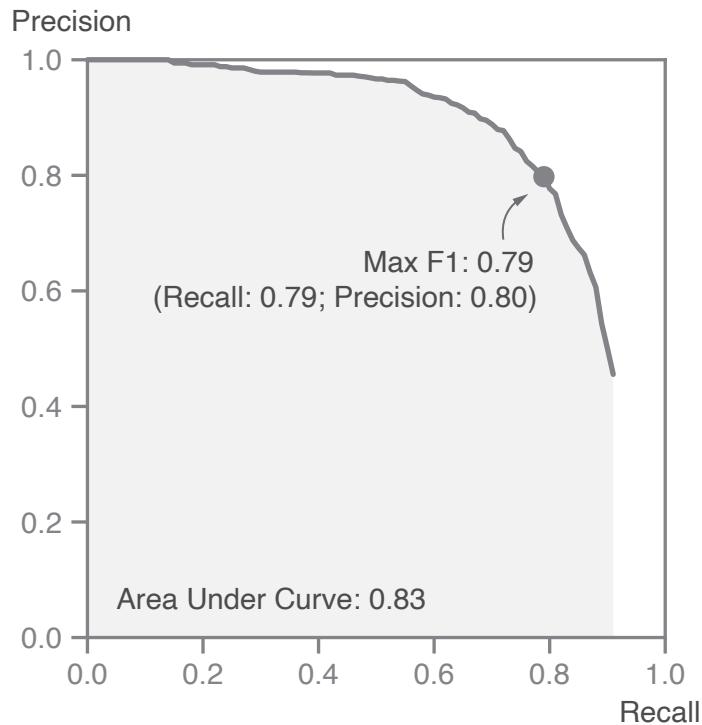


Figure S1: **The precision-recall curve of the Mask R-CNN model shows satisfactory predictive performance.** The Mask R-CNN model is trained and evaluated with 3-fold cross validation. The evaluation is based on 120 annotated images, which were randomly sampled from all the input satellite images in Siaya, Kenya. The Mask R-CNN model outputs a confidence score for every predicted building instance, and the precision-recall curve is generated by varying the confidence score threshold, below which predicted instances are dropped. A higher threshold makes the model more conservative and corresponds to the left portion of the curve (with high precision and low recall), and vice versa. The dot represents the optimal confidence score threshold, obtained by maximizing F1, the harmonic mean of precision and recall. The main model used in this study employs the optimal threshold, and has a recall of 0.79 and a precision of 0.80.



Figure S2: Ten randomly sampled pairs of input images and deep learning predictions.
Ten images are randomly sampled from all the input satellite images in the GiveDirectly study area. Each predicted building is outlined in white and filled with the “representative” roof color.

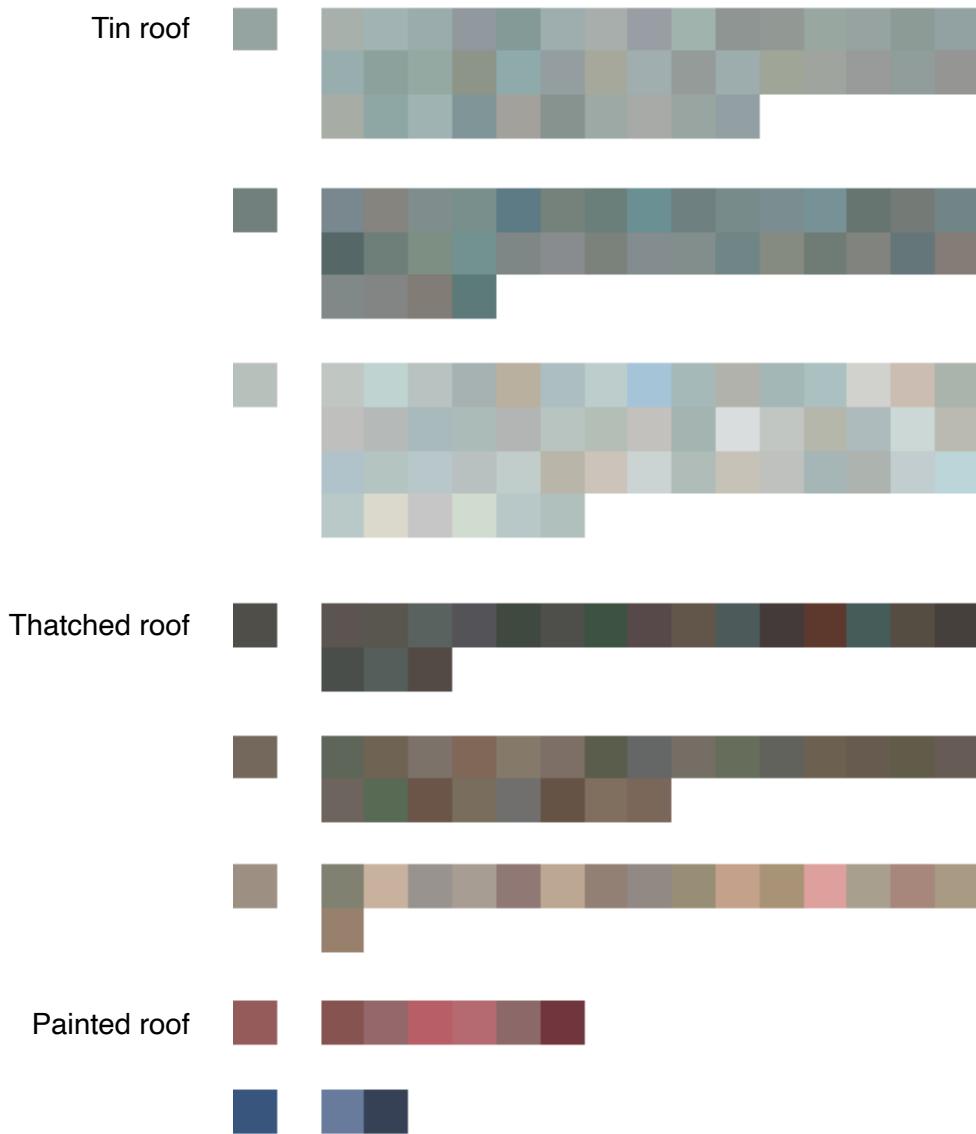
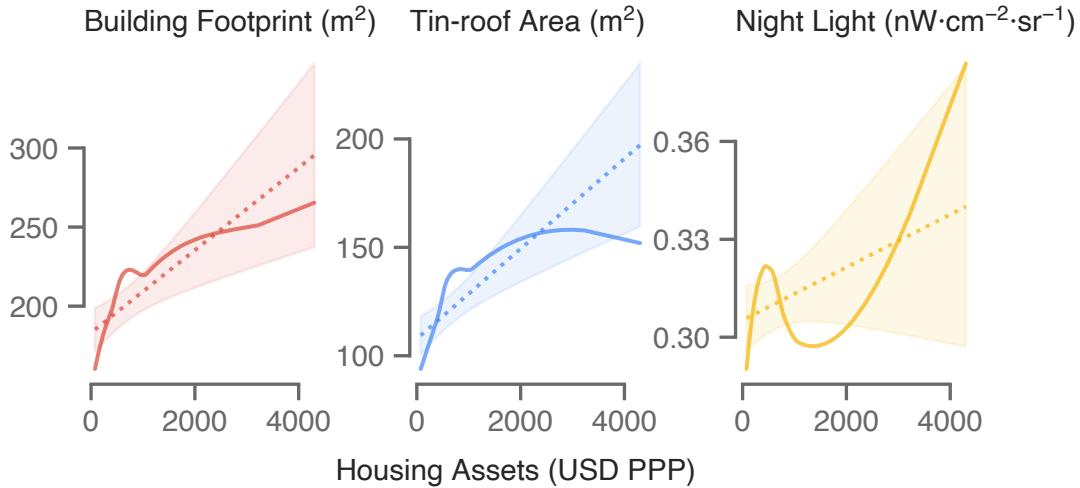


Figure S3: **The distribution and grouping of roof colors.** All the buildings in the GiveDirectly study area are split into eight groups by a K-means clustering algorithm, based on their roof colors. The color block on the left represents the “average” roof color of the cluster, and the color blocks on the right represent a random subset of all the roof colors in the given cluster. The number of color blocks on the right is proportional to the size of the cluster. The eight groups are further grouped into tin roof, thatched roof, and painted roof.

a Engel Curves



b Treatment Effect Estimates on Housing Assets

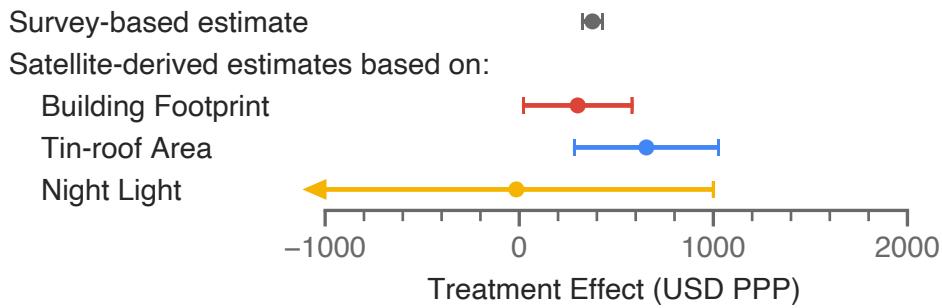
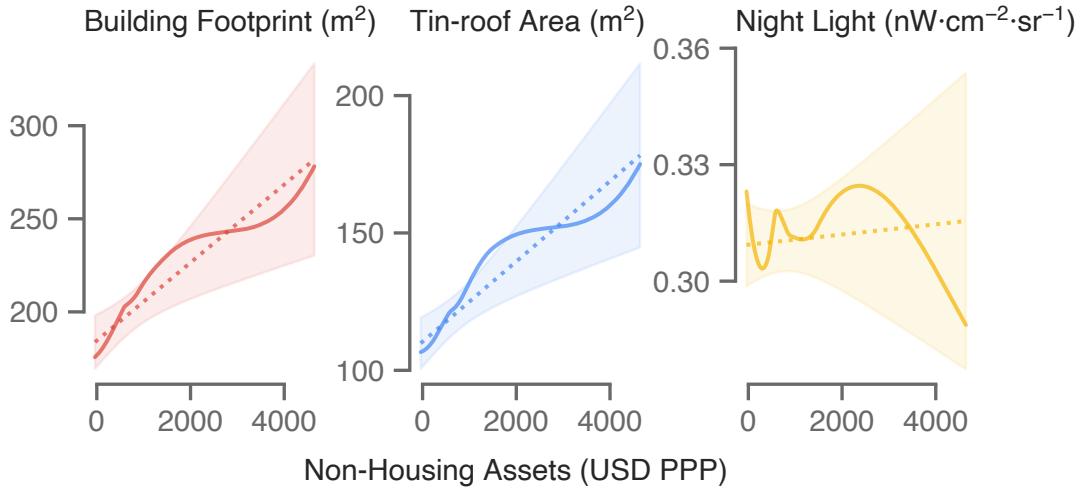


Figure S4: **The treatment effect on housing assets can be similarly recovered by scaling the effect on building footprint.** **a** The Engel curves of building footprint, tin-roof area, and night light, estimated with LOESS (solid line) or a linear regression (dotted line). The shaded regions represent the 95% confidence intervals for the latter. **b** Comparing the survey-based versus satellite-derived treatment effects. The dots show the point estimates. The error bars show the 95% confidence intervals, with the arrow(s) marking upper/lower bounds that are out of range (if any). $n = 1,844$.

a Engel Curves



b Treatment Effect Estimates on Non-Housing Assets

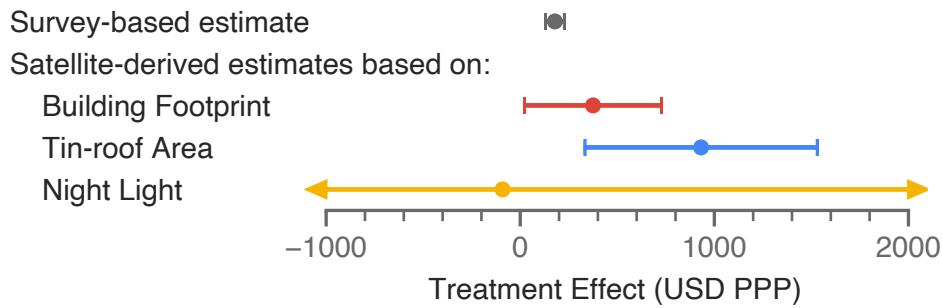
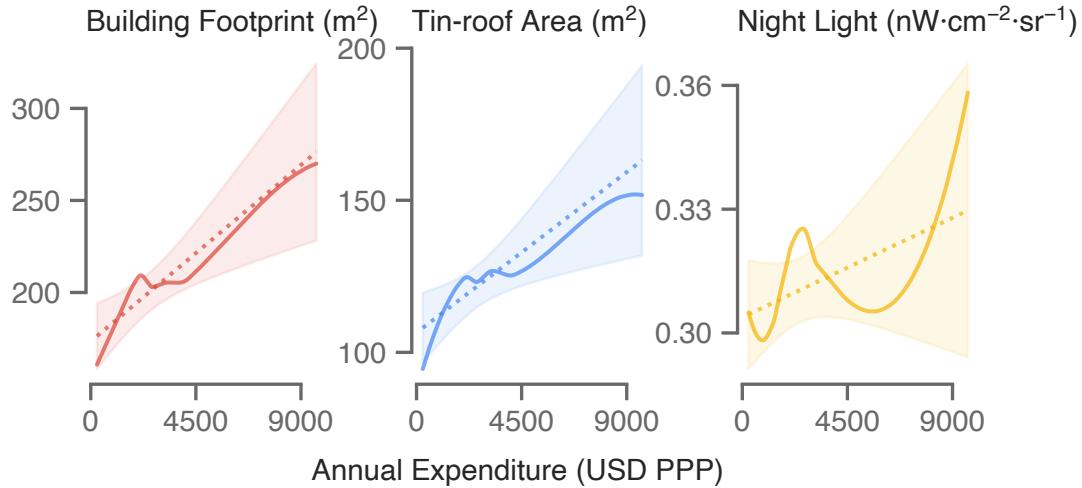


Figure S5: **The treatment effect on non-housing assets can be similarly recovered by scaling the effect on building footprint.** **a** The Engel curves of building footprint, tin-roof area, and night light, estimated with LOESS (solid line) or a linear regression (dotted line). The shaded regions represent the 95% confidence intervals for the latter. **b** Comparing the survey-based versus satellite-derived treatment effects. The dots show the point estimates. The error bars show the 95% confidence intervals, with the arrow(s) marking upper/lower bounds that are out of range (if any). $n = 1,844$.

a Engel Curves



b Treatment Effect Estimates on Annual Expenditure

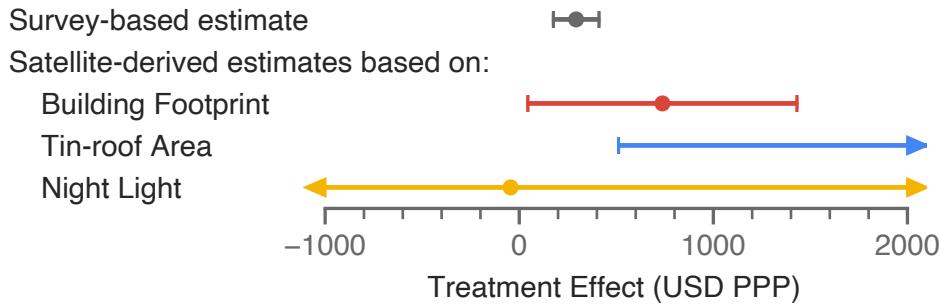


Figure S6: **The treatment effect on annual expenditure can be similarly recovered by scaling the effect on building footprint.** **a** The Engel curves of building footprint, tin-roof area, and night light, estimated with LOESS (solid line) or a linear regression (dotted line). The shaded regions represent the 95% confidence intervals for the latter. **b** Comparing the survey-based versus satellite-derived treatment effects. The dots show the point estimates. The error bars show the 95% confidence intervals, with the arrow(s) marking upper/lower bounds that are out of range (if any). $n = 1,843$.

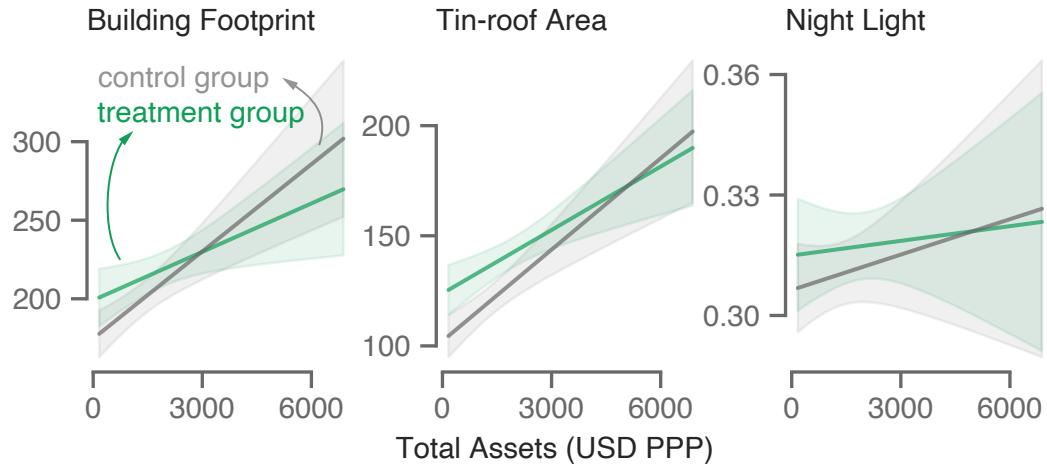


Figure S7: The Engel curves for tin-roof area shifted in response to the cash transfer.
The Engel curves for the treatment households (in green, $n = 1,904$) and the control households (in gray, $n = 1,844$). The shaded regions represent the 95% confidence intervals.

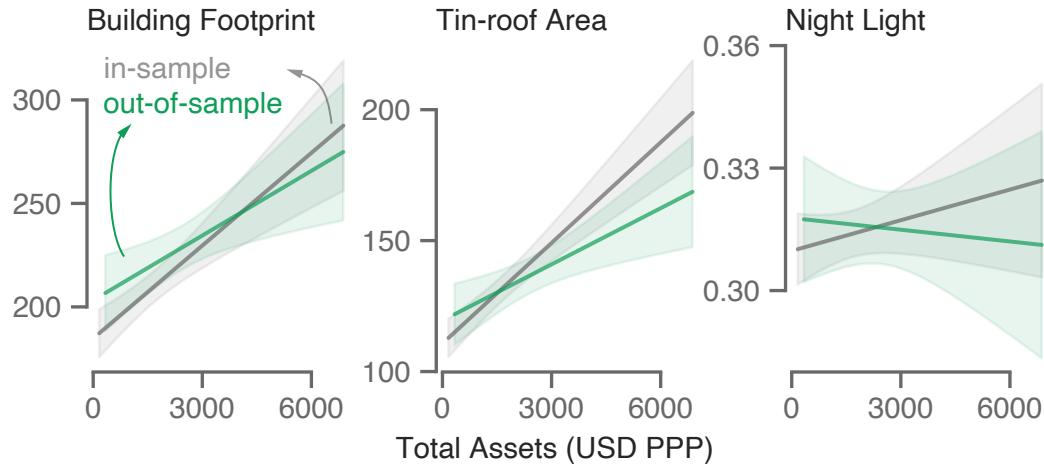


Figure S8: **The Engel curves estimated based on in-sample and out-of-sample data are broadly similar.** The Engel curves for the in-sample eligible households (in gray, $n = 3,748$) and the out-of-sample ineligible households (in green, $n = 1,821$) in the GiveDirectly study area. The shaded regions represent the 95% confidence intervals.

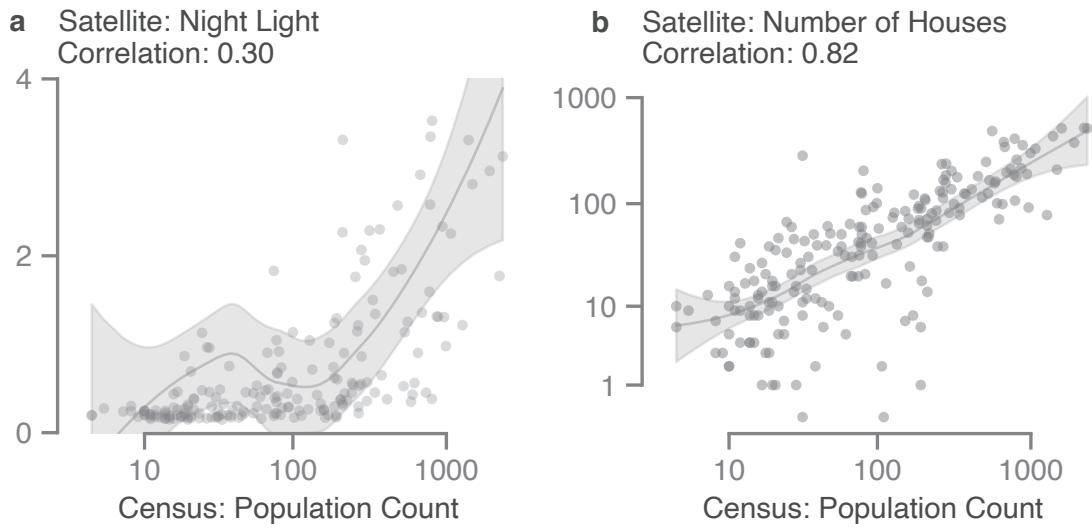


Figure S9: Population count in Mexican villages is more strongly correlated with the number of houses in satellite imagery, compared to night light. The population count is shown in log scale. Each point corresponds to a randomly sampled rural locality in Mexico. Gray lines are estimated LOESS curves, and the shaded regions are the 95% confidence intervals. The (Pearson) correlation coefficients are reported in the panel subtitles. $n = 197$.

585 **B Cost Estimation**

586 We estimate that our evaluation approach costs \$0.006 per household, when accounting for imagery
587 acquisition and computing costs. The Google Static Maps API charges users \$0.002 per image
588 request [1]. We estimate that our computing cost is roughly \$0.004 per household. Our entire
589 data pipeline can be run within 72 hours on an NC24 instance with 4 K80 GPUs on Microsoft
590 Azure, which costs \$3.60 per hour, and we have analyzed over 60,000 households. This is a liberal
591 estimate that accounts for image downloading, model training, model inference, model validation,
592 and regression analysis. Notably, we do not include labor costs for research and development, as
593 these only need to be incurred once, are not relevant for application of the method, and that such
594 labor costs are difficult to quantify.

595 **C Training the Deep Learning Model**

596 **C.1 Creating In-sample Building Footprint Annotations**

597 We create in-sample building footprint annotations to train the model, and to objectively and
598 quantitatively evaluate model performance. Among the 71,012 satellite images that cover all of the
599 Siaya county in Kenya, we randomly sample 120 images for annotation. We use the Supervisely
600 image annotation web platform to create annotations. On any given image, we outline the bound-
601aries of all the instances of buildings on the image. Buildings that border each other are annotated
602 as separate instances, if there are reasons to believe that they are separate structures (e.g., if they
603 appear to use different roof materials). Half-finished buildings are annotated, although they are
604 fairly rare in the analysis sample.

605 Some measurement errors can arise from the annotation process, which may in turn impact the
606 predictions of the deep learning model. First, the Google Static Maps logo blocks 1.05% of the total
607 area of any given image, and structures covered by the logos are not annotated. Second, only the
608 visible parts of the buildings are annotated, but a very small part of some buildings may be partially
609 occluded by trees. Third, the annotation accuracy (and thus potentially prediction accuracy) may
610 be different across buildings with different roof materials. In particular, thatched-roof houses tend
611 to be harder to identify for human annotators than metal-roof houses, because they are typically
612 smaller, not as reflective, and may resemble trees in the overhead imagery.

613 **C.2 Training the Mask R-CNN Model**

614 We use the Mask R-CNN model [2] for instance segmentation of buildings on satellite images. The
615 backbone architecture used is ResNet50 with the Feature Pyramid Networks. The model is trained
616 with a learning rate of 5×10^{-4} and a batch size of 10. Optimization is conducted with the Adam
617 optimizer. We implement the deep learning pipeline with Python and PyTorch. In particular, we
618 use the official Torchvision implementation of Mask R-CNN. We train the Mask R-CNN model in
619 a transfer learning framework, with a multi-step process as follows.

620 **1. COCO (Common Objects in Context)** The model is first pre-trained with the COCO
621 (Common Objects in Context) data set, a large-scale natural image data set containing 80 object
622 categories and around 1.5 million object instances [3]. Despite the fact that input images and
623 object categories in COCO are different from target satellite images, pre-training the model with
624 a large-scale dataset often provides meaningful performance gains, even when the model is later
625 transferred across domains.

626 **2. Open AI Tanzania** The model is then fine-tuned on the Open AI Tanzania building footprint
627 segmentation data set, a collection of high-resolution aerial imagery collected by consumer drones

628 in Zanzibar, Tanzania [4]. These images are representative of the rural or peri-urban scenes in a
629 developing country context, in terms of the distribution of the density, sizes and heights of the
630 buildings. All the buildings in the drone images are identified, outlined and classified into three
631 categories (completed building, unfinished building, and foundation) by human annotators. This
632 somewhat unusual categorization is due to the fact that there are a large number of unfinished
633 structures in Zanzibar. Most input satellite images in this study contain very few unfinished
634 structures, so we collapse the first two categories into one and drop the third category. The native
635 resolution of the drone images is 7cm, and we down-sample the images to about 30cm to match
636 with the resolution of the target satellite images.

637 In training time, 90% of the data are used for training, and the remaining 10% for validation.
638 In order to guard against overfitting, and choose the best model, in each epoch, we evaluate the
639 performance of the model with the validation set, using average precision with an Intersection over
640 Union (IoU) cutoff of 0.5 as the main evaluation metric. The model is trained for 50 epochs, and
641 the best model (at epoch 43) is saved and loaded in subsequent steps.

642 **3. Supplementary Annotations in Mexico, Tanzania and Kenya** The model is then
643 fine-tuned on a set of 587 annotated high-resolution satellite images from Mexico, Tanzania, and
644 Kenya. The Mexico dataset consists of 199 satellite images corresponding to 8 randomly sampled
645 rural localities studied in Supplementary Figure S9. Some of these are historical images with lower
646 data quality and more cloud coverage. These images are pooled and randomly split into a training
647 set (90%) and a validation set (10%). The model is trained for 25 epochs, and achieves the best
648 performance at epoch 17.

649 **4. In-sample Annotations** Finally, the model is fine-tuned on a set of 120 in-sample annotated
650 images in Siaya, Kenya (see Section C.1 for details). This ensures that training images and inference
651 images belong to the same data distribution. The model is trained on 90% of the images for 25
652 epochs, and evaluated with the 10% held out set. We keep the best-performing model (at epoch 15).
653 This is the main model used for conducting inference on input satellite images in the GiveDirectly
654 study area.

655 Throughout the training process, we conduct extensive data augmentation to increase the transfer-
656 ability of the model from one dataset to another. We randomly flip the training images horizontally
657 and vertically, randomly jitter the brightness, contrast, saturation, and hue of the images. For the
658 Open AI Tanzania dataset, we also randomly blur and crop the images.

659 **D Validation in Mexico**

660 **D.1 Results**

661 We provide additional validation results in rural Mexico, using the 2010 Population and Housing
662 Census [5]. Population count in a rural village (as reported in the 2010 census), is highly correlated
663 with the number of houses in that village (as identified by the deep learning model), with a Pearson
664 correlation coefficient of 0.82 (Supplementary Figure S9b). Population count, however, is only
665 modestly correlated with night light (Supplementary Figure S9a). Night light is less sensitive in
666 smaller, less populated villages, a finding that is consistent with prior work [6].

667 **D.2 Methods**

668 This comparison is based on the locality-level data set, Principales Resultados por Localidad, or
669 ITER. (A locality is equivalent to a village in rural areas.) To form the analysis sample, we drop all
670 urban localities (defined as having more than 2,500 residents), small localities where the relevant
671 asset measures are masked in the census to protect privacy, and localities where these measures are
672 missing. To avoid covering neighboring urban or rural localities in the satellite images, we exclude
673 rural localities that are closer than 0.01 degree (1.1 km) from other rural localities, or 0.1 degree
674 (11.1 km) from urban localities. Finally, to reduce computation, we randomly sample 200 rural
675 localities, and drop 3 of them, for which Google Static Maps does not have satellite image coverage
676 for.

677 In the census, each rural locality is geo-coded as a point. Most of the rural localities are small,
678 isolated and surrounded by vegetation or open space, making it feasible to match census records
679 to corresponding satellite images. For each locality, we obtain satellite images that cover an area
680 of roughly 1×1 km, with the locality coordinate at the center. The images are retrieved from the
681 Google Static Maps API on October 10, 2019, and are likely taken several years after the census.
682 We generate deep learning predictions on these images with the method described in Methods and
683 Supplementary Materials C, but only train the model for the first three steps in Supplementary
684 Materials C.2. For the comparison, we count the number of houses in a locality in the deep learning
685 predictions, and extract the population count variable from the census. Additionally, we download
686 night light data, the Visible Infrared Imaging Radiometer Suite (VIIRS) Day/Night Band (DNB)
687 composite images from 2019.

688 **References**

- 689 1. Google Static Maps. *Maps Static API Usage and Billing*
690 <https://developers.google.com/maps/documentation/maps-static/usage-and-billing>,
691 accessed 20 April 2021. 2021.
- 692 2. He, K., Gkioxari, G., Dollár, P. & Girshick, R. *Mask R-CNN*
693 in *Proceedings of the IEEE international conference on computer vision* (2017), 2961–2969.
- 694 3. COCO. *COCO - Common Objects in Contexts* <http://cocodataset.org>. accessed 6 May 2020.
695 2020.
- 696 4. Open AI Tanzania. *2018 Open AI Tanzania Building Footprint Segmentation Challenge*
697 <https://competitions.codalab.org/competitions/20100>. accessed 6 May 2020. 2020.
- 698 5. INEGI. *2010 Population and Housing Census of Mexico*
699 <https://www.inegi.org.mx/programas/ccpv/2010/default.html>. accessed 5 May 2020. 2010.
- 700 6. Jean, N. *et al.* Combining satellite imagery and machine learning to predict poverty.
701 *Science* **353**, 790–794 (2016).