



**AgEcon** SEARCH  
RESEARCH IN AGRICULTURAL & APPLIED ECONOMICS

*The World's Largest Open Access Agricultural & Applied Economics Digital Library*

**This document is discoverable and free to researchers across the globe due to the work of AgEcon Search.**

**Help ensure our sustainability.**

Give to AgEcon Search

AgEcon Search

<http://ageconsearch.umn.edu>

[aesearch@umn.edu](mailto:aesearch@umn.edu)

*Papers downloaded from **AgEcon Search** may be used for non-commercial purposes and personal study only. No other use, including posting to another Internet site, is permitted without permission from the copyright owner (not AgEcon Search), or as allowed under the provisions of Fair Use, U.S. Copyright Act, Title 17 U.S.C.*

# Using Unobtrusive Sensors to Measure and Minimize Hawthorne Effects: Evidence from Cookstoves

Andrew M. Simons, Theresa Beltramo, Garrick Blalock, David I. Levine\*

People act differently when they know they are being observed. This phenomenon, the Hawthorne effect, can bias estimates of program impacts. Unobtrusive sensors substituting for human observation can alleviate this bias. To demonstrate this potential, we used temperature loggers to measure fuel-efficient cookstoves as a replacement for three-stone fires. We find a large Hawthorne effect: when in-person measurement begins, participants increase fuel-efficient stove use approximately three hours/day (54%) and reduce three-stone fire use by approximately two hours/day (30%). When in-person measurement ends, participants reverse those changes. We then examine how this Hawthorne effect biases estimates of fuel use and particulate matter concentrations. Our results reinforce concerns about Hawthorne effects, especially in policy-relevant impact evaluations. We demonstrate that sensors can sometimes provide a solution.

*Keywords:* observation bias, Hawthorne effect, sensors, improved cookstoves, monitoring and evaluation, impact evaluation

*JEL Codes:* Q56, O13, I32, O22, D01

\* Simons: Department of Economics, Fordham University (email: [asimons5@fordham.edu](mailto:asimons5@fordham.edu)). Blalock: Charles H. Dyson School of Applied Economics and Management, Cornell University (email: [garrick.blalock@cornell.edu](mailto:garrick.blalock@cornell.edu)). Beltramo: United Nations High Commissioner for Refugees, Geneva (email: [beltramo@unhcr.org](mailto:beltramo@unhcr.org)). Levine: Haas School of Business, University of California, Berkeley (email: [levine@haas.berkeley.edu](mailto:levine@haas.berkeley.edu)).

*Acknowledgments:* This study was funded by the United States Agency for International Development under Translating Research into Action, Cooperative Agreement No. GHS-A-00-09-00015-00. The recipient of the grant was Impact Carbon who co-funded and managed the project. Juliet Kyaesimira and Stephen Harrell expertly oversaw field operations and Amy Gu provided excellent research support. We thank Impact Carbon partners Matt Evans, Evan Haigler, Jimmy Tran, Caitlyn Toombs, and Johanna Young; U.C. Berkeley Household Energy, Climate, and Health Research Group partners including Kirk Smith, Ilse Ruiz-Mercado, and Ajay Pillarisetti; Berkeley Air partners including Dana Charron, David Pennise, Michael Johnson, and Erin Milner; the USAID TRAction Technical Advisory Group, and seminar participants at Cornell University, UCLA, and Oxford University for valuable comments. Data collection was carried out by the Center for Integrated Research and Community Development (CIRCODU), and the project's success relied on expert oversight by CIRCODU's Director General Joseph Ndemere Arineitwe and field supervisors Moreen Akankunda, Innocent Byaruhanga, Fred Isabirye, Noah Kirabo, and Michael Mukembo. We thank the Atkinson Center for a Sustainable Future at Cornell University, the Institute for the Social Sciences at Cornell University and the Cornell Population Center for additional funding of related expenses. The findings of this study are the sole responsibility of the authors, and do not necessarily reflect the views of their respective institutions, nor USAID or the United States Government.

## **Introduction**

The validity of empirical research depends on data quality. Unlike the physical sciences, for which data often is generated in controlled laboratory settings, the social sciences typically measure variables involving human behaviors, which make data quality a challenge. Respondents often do not answer surveys candidly (Bertrand and Mullainathan 2001) and the act of surveying can change later behaviors of those being surveyed (Zwane et al. 2011). These drawbacks to surveys have been one factor contributing to a push for more experiments in environmental economics (Greenstone and Gayer 2009) and social science research more generally (Falk and Heckman 2009; Banerjee and Duflo 2009; Duflo, Glennerster, and Kremer 2008). While much has been learned from experiments in environmental economics, these types of experiments measuring human behaviors are susceptible to issues such as observation bias, or Hawthorne effect.

We explore an emerging class of technology—small, inexpensive, and unobtrusive sensors—as a remedy to the Hawthorne effect. A growing variety of sensors have become available to researchers. GPS trackers and motion detectors, for example, allow non-obtrusive measurement of subject location and body movements (Ermes et al. 2008). Medical doctors wear sensors that detect the scent of alcohol used in hand sanitizers to alert the doctor and/or patient if the doctor has not washed his or her hands recently (E. Smith 2014; Srigley et al. 2014). Loop detectors installed in the lanes of freeways allow monitoring of congestion and driver behavior (Bento et al. 2014).

The degree to which these sensors interfere with subjects' behavior can vary widely. In some cases, individuals may choose to be observed to motivate their own behavioral response. For example, long-distance bikers and runners can opt into programs that will report the location, time, and speed of excursions to a website that others can monitor (Mueller et al. 2010). Users of such schemes typically hope peer observation will increase their motivation. In other cases, such as room occupancy detectors that control lighting and climate control, the sensor may be far harder to notice (Buchanan, Russo, and Anderson 2014).

A major challenge for direct observational studies is that they alter participants' behavior. The effects of observers have been noted in cookstove studies (Ezzati, Saleh, and Kammen 2000; Smith-Sivertsen et al. 2009)), in energy consumption (Schwartz et al. 2013), in public health (Clasen et al. 2012; Das, Hammer, and Leonard 2008; Leonard and Masatu 2006; Srigley et al. 2014) in development economics (Leonard 2008; Leonard and Masatu 2010; Muralidharan and Sundararaman 2010) and in social sciences more broadly (Levitt and List 2007; Levitt and List 2011). We demonstrate a technique to remedy the Hawthorne effect that uses unobtrusive temperature sensors in an evaluation of fuel-efficient cookstoves in Uganda.

We use minimally invasive temperature sensors to measure usage of the fuel-efficient cookstoves and of the traditional three-stone fires.<sup>1</sup> We then compare usage

---

<sup>1</sup> A three-stone fire is simply three large stones, approximately the same height, on which a cooking pot is balanced over a fire.

of each stove in periods when observers visit the households with periods when no observers are present. We find a large Hawthorne effect: households increase the use of the fuel-efficient stove and decrease the use of three-stone fires on days they expect observers. Because the technique we deployed to measure fuel wood use (kitchen performance test) is the principal field-based procedure to demonstrate the effect of stove interventions on household fuel consumption (Bailis, Smith, and Edwards 2007), our findings are important for the cookstove evaluation community and for policy makers.

### **Data and methodology**

We implemented a series of studies in rural areas of the Mbarara District in southwestern Uganda from February to September 2012, which focused on the adoption, and use of fuel-efficient stoves. At baseline almost all families cooked on a traditional three-stone fire (97%), usually located within a separate enclosed cooking hut. We introduced an Envirofit G-3300 stove. Its manufacturer reports that it uses 50% less fuel and reduces smoke and harmful gasses by 51% compared to a three stone fire (Envirofit Inc. 2011). The study area is characterized by agrarian livelihoods including raising livestock and farming *matooke* (starchy cooking banana), potatoes, and millet.

We tracked stove usage before and after the purchase of a fuel-efficient stove at 168 households spread across fourteen rural parishes in Mbarara.<sup>2</sup> Upon arriving in a new parish, staff displayed the fuel-efficient stove (Envirofit G-3300) and offered it

---

<sup>2</sup> The population of most Ugandan rural parishes ranges from 4,000 to 6,000.

for sale to anyone who wanted to purchase at 40,000 Ugandan Shillings (approximately USD \$16, see (Beltramo et al. 2015; Levine et al. 2016) for an overview of the sales contract). Households were eligible to participate in the impact study if they mainly used wood as a fuel source, regularly cooked for eight or fewer persons, someone was generally home every day, and cooking was largely in an enclosed kitchen.

Eligible households who wanted to buy the stove were randomly assigned to two groups: early buyers, late buyers. Because it is crucial to measure both the use of the new stove *and* any reduction in use of traditional stoves (Ruiz-Mercado et al. 2011; Miller and Mobarak 2013), we asked both early buyers and late buyers if they would agree to have stove use monitors (SUMs) that read stove temperatures placed on their traditional and Envirofit stoves. After giving consent, three stone fires were fitted with SUMs immediately and we collected a baseline round of data with only three stone fires present in homes.

Approximately two to three weeks later the early buyers group received their first Envirofit stove. We did a midline round of data collection that is not used in this study (but will be the basis of an impact evaluation, when the data are cleaned). Approximately five to six weeks later the late buyers received their first Envirofit stove. About six weeks after late buyers received their Envirofits, both groups were surprised with a second Envirofit stove. Because common cooking practices in the area require two simultaneous cooking pots (for example rice and beans, or *matooke* and some type of sauce), and the Envirofit is sized for one cooking pot, we

gave a second Envirofit to permit normal cooking using only fuel-efficient stoves. We then collected our endline data, the core data we use in this study.

We tracked stove temperatures for approximately six months (April–September 2012). To track usage, we used small, inexpensive and unobtrusive sensors: stove use monitors (SUMs) that record stove temperatures without the need for an observer to be present.<sup>3</sup> Using SUMs to log stove temperatures was initially suggested by (Ruiz-Mercado et al. 2008) and has been used successfully in various settings (Mukhopadhyay et al. 2012; Ruiz-Mercado et al. 2013; Pillarisetti et al. 2014). We installed SUMs on two Envirofits and two three-stone fires (by the end of the study numerous SUMs had been lost or burned up; therefore, at the end line we measured both Envirofits and the primary three stone fire).

We also performed standard kitchen performance tests (KPT) (Bailis, Smith, and Edwards 2007) in each household to measure the quantity of fuel wood used, record detailed food diaries, and measure household air pollution. The KPT lasts approximately a week and involves daily visits by a small team of researchers weighing wood, monitoring household air particulate monitors, and collecting survey data on stove usage over the last 24 hours and related topics.

---

<sup>3</sup> The SUMs used for our project, iButtons™ manufactured by Maxim Integrated Products, Inc., are small stainless steel temperature sensors about the size of a small coin and the thickness of a watch battery which can be affixed to any stove type. Our SUMs record temperatures with an accuracy of +/- 1.3 degrees C up to 85°C. For additional details see the product description website at: <http://berkeleyair.com/services/stove-use-monitoring-system-sums/> The SUMs cost approximately USD\$16 each and could record temperature data for 24 hours a day for six weeks in a household before needing minimal servicing from a technician to download the data. After the data download they can be reset and re-used.

Comparing stove usage calculated from the temperature data collected by the SUMs in the week while KPT measurement teams are present versus stove usage in the week before and after the measurement week provides our test of a Hawthorne effect.

Throughout the study, field staff recorded about 2,400 visual observations of whether a stove was in use (on/off) when they visited homes to exchange stove usage monitors or gather data for the KPT. Then we used a machine-learning algorithm to examine the temperature data immediately before and after the 2,400 visual observations of use. The algorithm analyzed the data to understand how temperature patterns change at times of observed stove use and then predicted cooking behaviors to the wider dataset of 1.7 million temperature readings. This process, detailed in (Simons et al. 2014), allowed us to unobtrusively and inexpensively track daily stove usage on a large sample of households for six continuous months.<sup>4</sup> See Appendix A for additional technical details surrounding the area selection, experimental rollout, SUMs placement and algorithm.

### **Specification**

Assign the subscripts  $t=-1$  to the week prior to measurement week,  $t=0$  to the measurement week, and  $t=1$  to the week after the measurement week. Let the coefficient on stove type  $s = TSF$  for three-stone fire or  $ENV$  for Envirofit, and

---

<sup>4</sup> Overnight, while most participants report sleeping, SUMs record the residual heat absorbed in the large stones of the three stone fires and/or from coals banked overnight. Therefore our algorithm overestimates overnight cooking of three stone fires. We adjust for this in the subsequent analysis. For further discussion and a description of the technical adjustment see (Harrell et al. 2016).



*Adjacent\_Week* be a dummy variable for an adjacent week ( $t=-1$  or  $t=1$ ). The regression is modeled using Ordinary Least Squares (OLS) as:

$$H_{it}^S = B^S * Adjacent\_Week + I_i + e_{it} \quad (1)$$

where  $H_{it}^S$  is the total hours cooked per day on stove type  $s$  for household  $i$  during the week,  $I_i$  is fixed effects for the individual household (which controls for household level characteristics that don't change over these three weeks like family size, income, housing structure, etc.), and  $e_{it}$  is an error term. The coefficient  $B^S$  is the estimate of how different (in hours cooked per day) the average adjacent week is compared to a measurement week on stove type  $s$ . Standard errors are clustered at the household.

To test the weeks separately, we use a slightly different specification. Let  $H_{it-1}^S$  be a dummy variable equal to 1 for the week before the measurement week (when  $t=-1$ ) and 0 otherwise, and let  $H_{it+1}^S$  be a dummy variable equal to 1 for the week after the measurement week (when  $t=1$ ) and 0 otherwise. Then the regression is modeled using OLS as:

$$H_{it}^S = \gamma_1^S * H_{it-1}^S + \gamma_2^S * H_{it+1}^S + I_i + e_{it} \quad (2)$$

where  $I_i$  is household fixed effects and  $\gamma_1^S$  is the estimate of the difference (in hours cooked per day) of the week before the measurement week compared to the measurement week. The coefficient of  $\gamma_2^S$  is the estimate of the difference cooked in

Unknown  
Field Code

Unknown  
Field Code

the week after the measurement week compared to the measurement week. Standard errors are clustered at the household.

## Results

In the week before the observers arrived (when  $t=-1$ ), primary three-stone fires were used an average of 5.99 hours per day (95% CI = [4.77 to 7.21]) and combined usage on Envirofits was an average of 5.53 hours per day (95% CI = [4.36 to 6.71]). We first estimate equation 1, where we constrain the effect of the observers arriving to be the same magnitude (but opposite sign) as the effect of the observers leaving. On average, usage of the Envirofit stoves is 2.97 hours higher during the measurement week than during the adjacent weeks (95% CI = [1.79 to 4.15],  $p<0.01$ , Table 1, column 3). This increase is matched by a reduction of 1.78 hours in usage of the three-stone fire (95% CI = [0.86 to 2.70],  $p<0.01$ , col. 1).

In columns 2 and 4 we relax the assumption that stove usage is the same in the week prior to and the week after our measurement period. Contrasted with the measurement week, households use their primary three-stone fire 1.17 hours per day more in the prior week (95% CI = [0.10 to 2.24],  $p<0.05$ , col. 2) and 2.37 hours more in the following week (95% CI = [1.12 to 3.62],  $p<0.01$ ). These coefficients are jointly significantly different than zero ( $p<0.01$ ), but not statistically significantly different from each other ( $p=0.10$ ).

The total usage of Envirofits follows a mirror image (col. 4), and is 2.58 hours per day lower in the week prior to measurement week than in measurement week (95%

CI = [1.21 to 3.94],  $p < 0.01$ ) and 3.30 hours per day lower the following week (95% CI = [2.04 to 4.57],  $p < 0.01$ ). These coefficients are jointly significantly different from zero ( $p < 0.01$ ), but not statistically significantly different from each other ( $p = 0.20$ ).

### **Adjusting for the Hawthorne effect**

Because the kitchen performance test is widely used to measure the effects of new cookstoves on fuel usage and household air pollution (K. Smith et al. 2007; Berrueta, Edwards, and Masera 2008; Johnson et al. 2010)—as well as the basis for the measurement of carbon emissions—estimates of how new stoves affect fuel use and carbon emissions may be substantially biased. The same bias can arise in studies, such as ours, that also measure household air pollution or health effects with repeated household visits. We develop a basic framework for testing for the magnitude of this bias and examine its extent in our setting.

#### *Basic framework*

The field of epidemiology has “efficacy trials” that test the effects of an intervention under ideal conditions and “effectiveness trials” that test the effects of an intervention under realistic conditions (Flay 1986). In the context of cookstoves, the kitchen performance test provides a valid measure of how the new stove affects wood usage during the measurement week (as in an efficacy trial); however, we need to adjust for the gap in usage between measurement weeks and weeks when no observers are influencing behaviors to generalize to weeks without daily visits (that is, to estimate effectiveness). Next we consider various illustrative examples using data from our setting.

### *Illustrative examples*

Table 2 presents the daily mean values of firewood consumption, particulate matter concentration and total three stone fire usage prior to the introduction of fuel-efficient stoves. The average household consumes 9.0 kgs of firewood per day (col. 1), has a daily concentration of PM2.5 of 428  $\mu\text{g}/\text{m}^3$  (col. 2) and cooks for a total of 14.0 daily hours (col. 3) across two three stone fires. To examine the bias introduced by the Hawthorne effect in our setting we need to know the expected biomass and pollution reductions for the new stove. To find the expected reduction we examine the “Emission and Performance Report” for the Envirofit G3300 performed by the Engines and Energy Conversion Lab at Colorado State University. These emissions measurements are based on accepted biomass stove testing protocols in a carefully monitored laboratory setting. The report (Figure 1) finds average improvements of 50.1% for fuel use and 51.2% for particulate matter emissions using the Envirofit G3300 versus a three stone fire (Envirofit Inc. 2011).

Using these mean values, we construct illustrative efficacy and effectiveness trails according to the framework above. For the purpose of our illustrative example, we assume a similar sized Hawthorne effect on the usage of the secondary three stone fires as well as what was observed on the first three stone fire (recall that attrition of sensors led us to measure fewer secondary three stone fires in the endline).

Using the assumptions above (Table 3), firewood consumed and daily PM2.5 concentrations were 8.0 kg/day and 382  $\mu\text{g}/\text{m}^3$  in the efficacy trial (observers

present), but 9.3 kg/day and 445  $\mu\text{g}/\text{m}^3$  in the effectiveness trial (no observers). These estimates are 14% lower with observers than without.

*Bias introduced by Hawthorne effect*

Table 4 presents a comparison of the endline to the baseline levels of daily cooking hours (on all stoves combined), daily firewood usage, and PM2.5 daily concentrations. Recall that at baseline homes had two three stone fires, and at endline homes had two three stone fires plus two Envirofits. These results are not causal estimates, as seasonal or time effects may influence them.

When we use time periods when observers were present, between baseline and endline: cooking time rose 20% (from 14.0 to 16.8 hours),<sup>5</sup> firewood use declined 11% (9.0 to 8.0 kg/day), and particulate matter also fell 11% (from 428 to 382  $\mu\text{g}/\text{m}^3$ ). When examining weeks when observers were not present (as in an effectiveness trial), some important results are reversed. Now cooking time rose 24% (from 14.0 to 17.3 hours), firewood use rose 4% (from 9.0 to 9.3 kg/day) and exposure to particulate matter grew 4% (from 428 to 445  $\mu\text{g}/\text{m}^3$ ). That is, adjusting for the Hawthorne effect turned a decline of about 11% in wood use and particulate matter into a small increase of about 4%. This is especially important in the context of calibrating the production function of emissions. The procedure of the kitchen performance test to measure fuel use and pollution assumes that stove use (and

---

<sup>5</sup> While total time cooking increases this is calculated over four stoves (two three stone fires and two Envirofit stoves) during the end line data collection period, but calculated over only two three stone fires during the baseline period. So it is likely that cooks actually spend less of their time cooking at end line because they have more stoves per meal at their disposal.

resultant pollution) while observers are present is the same as when no observers are present. We demonstrate that this is not the case.

This illustrative example shows how important it is to account for Hawthorne effects in impact evaluations. Using the sample means from our data, and the emissions and performance report for the Envirofit G3300 the Hawthorne effect not only biases the magnitude of the change, but (with these assumptions) also reverses the direction of the change over time.

#### *Attrition of SUMs devices*

One concern for our study is whether the attrition of sensors used to measure stove temperatures was random. In cases of sensor malfunction we lost the temperature readings associated with that device (about six weeks of data for that individual stove). The concern is that if damage (overheating above the 85°C tolerance of our SUMs device)<sup>6</sup> was more likely on stoves that were used more heavily, then the data we have are not an unbiased measure of stove usage for the broader sample. If however, the attrition of SUMs sensors is random, there is less concern about the internal validity of our sample.

We test this in various ways. First, we regress whether the SUMs data was missing at endline (device malfunctioned) on household fuel wood consumption during that same period. Because fuel wood is a direct input into how much the stoves are used,

---

<sup>6</sup> There is a SUMs model with a heat tolerance of 120°C, but it was much more expensive so budget considerations forced us to use the model with 85°C heat tolerance.

this is the most direct test of this relationship. If households that cook more (using fuel wood consumption as a proxy) also have a higher probability of SUMs attrition, this would be evidence of non-random attrition and a problem for our study. We examine this relationship separately for each stove type that we included in our study (recall that we choose not to track the non-primary three stone fire by the endline of our study). We also kept records of whether we placed a SUMs device on a stove or not, for example if the homeowner was not home we did not enter the house to download data and reset the SUMs device. Because we are testing for attrition due to excessive cooking (heat exposure) we only test for this relationship on the sample of stoves on which we placed a SUMs device. We also do similar checks with two other variables that are related to cooking (count of people cooked for daily, and number of meals cooked daily).

In Table 5 we present the results of the attrition checks. In our preferred test, we find that the likelihood of SUMs survival is statistically no different than zero (col. 1-3) for each additional kilogram of wood consumption. When examining whether a larger household size is associated with the likelihood of SUMs survival we find a weakly statistically significant relationship for primary three stone fire usage (col. 4). Each additional person cooked for is associated with a four-percentage point decrease in the probability of SUMs survival ( $p < 0.10$ ), however this relationship does not appear for either of the Envirofit stoves (col. 5-6). Finally, when testing whether the count of daily meals cooked is associated with SUMs survival (col. 7-9) we find no statistically significant relationship. Taken as a whole, these tests do not provide strong evidence of non-random attrition of SUMs devices.

## **Conclusion**

We demonstrate a technique to measure the magnitude of—and adjust for—a Hawthorne effect in a field experiment in the developing world. Given the push for more experiments in environmental economics (Greenstone and Gayer 2009), developing techniques to generate data that does not suffer from observer bias is necessary if the evaluations are to help make unbiased policy recommendations. In our specific setting, the findings of a large Hawthorne effect have implications for the impact of fuel-efficient stoves on fuel use and air particulates. The kitchen performance test is the current “gold standard” for generating Certified Emission Reductions that can be sold into the emissions trading markets of the Clean Development Mechanism. Our findings potentially call into question the veracity of these CO<sub>2</sub> reductions. More broadly, our results reinforce the importance for observed behaviors to be independently verified with unobtrusive monitoring. Evaluation processes that combine direct data generation (e.g., respondent surveys, visual observation) should be cross validated with an unobtrusive data generation process such as sensors.

While other forms of unobtrusive objective monitoring exist—such as using administrative records when reliable (Angrist, Bettinger, and Kremer 2006) or tracking take-up at a remote location via redeemed vouchers (Dupas 2009; Dupas 2014)—the recent explosion of small, inexpensive, and unobtrusive sensors expands researchers’ ability to quantify and remove observation bias. A wide variety of emerging technologies can be utilized, a partial list includes: smart phones



tracking locations through GPS, remote sensors that detect latrine usage (Clasen et al. 2012), sensors to remotely detect the use of water filters (Thomas et al. 2013), medical devices to monitor the hand hygiene of medical professionals (Boyce 2011), smart grid or other energy monitors (Darby 2010), and pedometers or other devices that monitor physical activity (Bravata et al. 2007). Adjusting for Hawthorne effects is essential if the results of impact evaluations are intended to generalize beyond periods of intense in-person observation.

## **Appendix A: Technical Details**

These technical details draw largely from our previously published work that focused on the methodological aspects of the experiment (Harrell et al. 2016; Simons et al. 2014).

### *Background of Study Area*

For our study site, we selected the Mbarara region of Uganda because it is rural, almost all families cooked on a traditional three stone fire, there was no active fuel-efficient cookstove intervention in the region, it was less than a day's travel from Kampala, and families spent a lot of time gathering wood (approximately 10–20 hours per week).

The study area is characterized by agrarian livelihoods including farming of *matooke* (starchy cooking banana), potatoes, and millet as well as raising livestock. At baseline almost all families cook on a traditional three-stone fire (97%), usually located within a separate cooking hut (62% of households had totally enclosed

kitchens with no windows, while 38% had semi-enclosed kitchens with at least one window). It is possible for households to move their three stone fires, but it is generally in our study area the common practice was to cook almost exclusively in their detached cooking hut.

Most stove usage occurs preparing lunch and dinner, with *matooke* and beans the most common and most time-consuming cooked foods. *Matooke*, a main food for lunch and dinner, is unripe plantain eaten after steaming for 3-5 hours. Beans, another common dish, are prepared by boiling and simmering for 2-4 hours. Thus for the main meals, it is common cooking practice to simmer and/or steam foods for several hours in a row.

Because our site was very close to the equator, there is little cyclicity in temperatures (daily temperatures generally fluctuated from 18-27°C throughout the six months). Metrological data shows bi-modal rainfall peaks in April and November (~120mm) and annual lows of rainfall for Mbarara in June and July (~20mm).

### *Experimental Rollout*

We tracked stove usage (before and after the purchase of a fuel-efficient stove) amongst twelve households in each of fourteen rural parishes in Mbarara (168 total households). Upon arriving in a new parish, research staff displayed the new Envirofit and offered it for sale to anyone who wanted to purchase at 40,000 Uganda Shillings (USD \$16). Consumers who wanted to buy the stove were randomly assigned into two groups (early buyers, late buyers). The project asked both early

buyers and late buyers if they would agree to have SUMs placed on their traditional three stone fires immediately. Then approximately two weeks later the early buyers group received their first Envirofit stove, and approximately four to five weeks after that the late buyers received their first Envirofit stove. Households were eligible to participate in the study if they mainly used wood as a fuel source, regularly cooked for eight or fewer persons, someone was generally home every day, and cooking was largely in an enclosed kitchen. In each parish, more than twelve households met these criteria and agreed to join the study; therefore among those that agreed, we randomly selected twelve households per parish for the usage study with the SUMs. Stove temperatures were tracked for approximately six months (April-September 2012). Approximately six weeks after late buyers received their Envirofits, both groups were surprised with a second Envirofit stove. Because common cooking practices in the area require two simultaneous cooking pots (for example rice and beans, or *matooke* and some type of sauce), and the Envirofit is sized for one cooking pot, we gave a second Envirofit to mimic normal cooking behavior as much as possible. Each household had as many as two three stone fires and two Envirofit stoves monitored with SUMs throughout the study.

### *Placement of SUMs*

SUMs must be placed close enough to the heat source to capture changes in temperatures, but not so close that they exceed 85°C, the maximum temperature the SUMs used in this study can record before they overheat and malfunction. We do not need to recover the exact temperature of the hottest part of the fire to learn about cooking behaviors. Even with SUMs that are reading temperatures 20-30 cm from

the center of the fire, as long as the temperature readings for times when stoves are in use are largely different than times when stoves are not used the logistic regression will be able to predict a probability of usage.

SUMs for three stone fires were placed in a SUM holder (Figure A1) and then placed under one of the stones in the three stone fire (left panel, Figure A2). The SUMs for Envirofits were attached using duct tape and wire and placed at the base of the stove behind the intake location for the firewood (right panel, Figure A2). Figure A3 shows an example of SUMs temperature data for a household across about three weeks. The left panel shows the temperatures registered in a three stone fire versus the ambient temperature also recorded with SUMs in this household, while the right panel compares the temperature of the Envirofit to the ambient temperature reading.

#### *Visual Observations of Use*

Each time any part of data collection team visited a household he or she visually observed which stoves were in use (whether the stove was “on” or “off” along with the date and timestamp recorded digitally via handheld device). Enumerators visited a house numerous times during a “measurement week,” when we also enumerated a survey and weighed wood for the kitchen performance test. Another enumerator visited once every 4-6 weeks to download data and reset the SUMs device.

### *Generating an Algorithm*

Our technique requires continuous SUMs temperature data for a given stove, and recorded instances of whether that particular stove is seen in use or not. We matched visual observations of stove use to SUMs temperature data by time and date stamps. The core of our method is a logistic regression using the lags and leads of the SUMs temperature data to predict visual observations of stove usage. We tested ten specifications of differing combinations of current, lagged, and leading temperature readings (see Simons et al. (2014) for the specific ten specifications we tested).

In order to determine which of the models was most appropriate we test the ten specifications with the Akaike information criterion (AIC) (Akaike 1981). The AIC trades off goodness of fit of the model with the complexity of the model to guard against over fitting.

The preferred specification included the temperature reading closest to the time of the visual observation, the readings 60 and 30 minutes prior, and 60 and 30 minutes after the visual observation of use, and a control for hour of the day. This regression specification correctly predicted 89.3% of three stone fire observations and 93.8% of Envirofit observations of stove usage. We then compared our algorithm to other previously published algorithms (Ruiz-Mercado, Canuz, and Smith 2012; Mukhopadhyay et al. 2012). Those algorithms focused on defining “discrete” cooking events based on rapid temperature slope increases, elevated stove temperatures, and then followed by a cooling off period. We applied those

algorithms to the temperature data we collected and found our logistic regression correctly classified more observations, with a higher pseudo R-squared, than any other algorithm for both three stone fires and the Envirofits.

### *Kitchen Performance Test Protocol*

The kitchen performance test weights the woodpile in a kitchen on sequential days to quantify the amount of wood used in a given 24-hour period (for additional details see: <http://cleancookstoves.org/technology-and-fuels/testing/protocols.html>). The KPT is the protocol used to estimate fuel savings, a primary component of calculating carbon credits for a stove project (The Gold Standard Foundation 2013). To minimize variance, the standard recommendation is that the KPT testing period should be at least three days, avoiding weekends and holidays (Bailis et al., 2007).

On the initial visit of the KPT week, the data collection team asked the household cook to describe what fuels they would use in the next 24-hour period. The data collection team asked the household to stack the wood they expected to use in a pile and only use wood from that pile over the next 24 hours. To ensure that the household did not run out of fuel, we asked the household to add a few extra pieces to the pile before we weighed the pile. In approximately 24 hours, the data collection team returned and weighed the remaining fuel. This process was repeated at approximately the same time each day of the KPT week (Monday – Thursday).

We measured mean 24-hour concentrations of PM<sub>2.5</sub> by installing calibrated UCB-PATS PM monitors in the study participants' homes during the same 72 hours of the kitchen performance test. We followed best practices as outlined by the Berkeley Air Monitoring Group (see: <http://berkeleyair.com/services/ucb-particle-and-temperature-sensor-ucb-pats/>) and measured three consecutive days of mean 24-hour PM<sub>2.5</sub> concentrations in the kitchen. We averaged data from the UCB-PATS PM monitors into 24-hour average PM<sub>2.5</sub> readings in  $\mu\text{g}/\text{m}^3$ .

### *Holidays*

Because we tracked stove temperatures for six months there were some holidays. However, our data collection period did not coincide with the most important annual holiday period (i.e., Christmas). We purposely did not collect KPT data during Easter week, the most important holiday during our study period. Additionally, we sometimes shifted the data collection one day later to avoid local holidays (e.g., visit Tuesday – Friday if Monday was a holiday). The KPT data collection team visited one parish per week for fourteen weeks to finish a round of KPTs. Therefore if a holiday increased or decreased cooking during a measurement week in one parish, the other parishes would not have KPT measurements taken that same week because our KPT measurement team was only in one parish per week. Thus, it seems unlikely that holidays would account for the large Hawthorne effect we estimate.

## Bibliography

- Akaike, Hirotugu. 1981. "Likelihood of a Model and Information Criteria." *Journal of Econometrics* 16: 3–14.
- Angrist, Joshua, Eric Bettinger, and Michael Kremer. 2006. "Long-Term Educational Consequences of Secondary School Vouchers: Evidence from Administrative Records in Colombia." *American Economic Review* 96 (3): 847–862.
- Bailis, Rob, Kirk R. Smith, and Rufus Edwards. 2007. "Kitchen Performance Test (KPT)." University of California, Berkeley, CA.
- Banerjee, Abhijit, and Esther Duflo. 2009. "The Experimental Approach to Development Economics." *Annual Review of Economics* 1 (1): 151–178. doi:10.1146/annurev.economics.050708.143235.
- Beltramo, Theresa, David I. Levine, Garrick Blalock, and Andrew M. Simons. 2015. "The Effect of Marketing Messages and Payment over Time on Willingness to Pay for Fuel-Efficient Cookstoves." *Journal of Economic Behavior & Organization*. doi:10.1016/j.jebo.2015.04.025.
- Bento, Antonio, Daniel Kaffine, Kevin Roth, and Matthew Zaragoza-Watkins. 2014. "The Effects of Regulation in the Presence of Multiple Unpriced Externalities: Evidence from the Transportation Sector." *American Economic Journal: Economic Policy* 6 (3): 1–29. doi:10.1257/pol.6.3.1.
- Berrueta, Víctor M., Rufus D. Edwards, and Omar R. Masera. 2008. "Energy Performance of Wood-Burning Cookstoves in Michoacan, Mexico." *Renewable Energy* 33 (5): 859–70. doi:10.1016/j.renene.2007.04.016.
- Bertrand, Marianne, and Sendhil Mullainathan. 2001. "Do People Mean What They Say? Implications for Subjective Survey Data." *American Economic Review* 91 (2): 67–72.
- Boyce, John M. 2011. "Measuring Healthcare Worker Hand Hygiene Activity: Current Practices and Emerging Technologies." *Infection Control and Hospital Epidemiology* 32 (10): 1016–28. doi:10.1086/662015.
- Bravata, Dena M, Crystal Smith-Spangler, Vandana Sundaram, Allison L Gienger, Nancy Lin, Robyn Lewis, Christopher D Stave, Ingram Olkin, and John R Sirad. 2007. "Using Pedometers to Increase Physical Activity: A Systematic Review." *Jama* 298 (19).
- Buchanan, Kathryn, Riccardo Russo, and Ben Anderson. 2014. "Feeding Back about Eco-Feedback: How Do Consumers Use and Respond to Energy Monitors?" *Energy Policy* 73 (June): 138–146. doi:10.1016/j.enpol.2014.05.008.
- Burnham, Kenneth P., and David R. Anderson. 2002. *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach*. 2nd Editio. New York, NY: Springer-Verlag.
- Clasen, Thomas, Douglas Fabini, Sophie Boisson, Jay Taneja, Joshua Song, Elisabeth Aichinger, Anthony Bui, et al. 2012. "Making Sanitation Count: Developing and Testing a Device for Assessing Latrine Use in Low-Income Settings." *Environmental Science & Technology* 46 (6): 3295–303. doi:10.1021/es2036702.
- Darby, Sarah. 2010. "Smart Metering: What Potential for Householder Engagement?" *Building Research & Information* 38 (5): 442–457. doi:10.1080/09613218.2010.492660.



- Das, Jishnu, Jeffrey Hammer, and Kenneth Leonard. 2008. "The Quality of Medical Advice in Low-Income Countries." *Journal of Economic Perspectives* 22 (2): 93–114.
- Duflo, Esther, R Glennerster, and Michael Kremer. 2008. "Using Randomization in Development Economics Research: A Toolkit." Edited by T Paul Schultz and John A Strauss. *Handbook of Development Economics* 4 (7): 3895–3962. doi:10.1016/S1573-4471(07)04061-2.
- Dupas, Pascaline. 2009. "What Matters (and What Does Not) in Households' Decision to Invest in Malaria Prevention?" *American Economic Review* 99 (2): 224–30. doi:10.1257/aer.99.2.224.
- . 2014. "Short-Run Subsidies and Long-Run Adoption of New Health Products: Evidence from a Field Experiment." *Econometrica* 82 (1): 197–228.
- Envirofit Inc. 2011. *G-3300 Woodstove Features*. <http://www.envirofit.org/products/?sub=cookstoves&pid=10>.
- Ermes, Miikka, Juha Pärkkä, Jani Mäntyjärvi, and Ilkka Korhonen. 2008. "Detection of Daily Activities and Sports With Wearable Sensors in Controlled and Uncontrolled Conditions." *IEEE Transactions on Information Technology in Biomedicine* 12 (1): 20–6. doi:10.1109/TITB.2007.899496.
- Ezzati, Majid, Homayoun Saleh, and Daniel M Kammen. 2000. "The Contributions of Emissions and Spatial Microenvironments to Exposure to Indoor Air Pollution from Biomass Combustion in Kenya." *Environmental Health Perspectives* 108 (9): 833–39.
- Falk, Armin, and James J Heckman. 2009. "Lab Experiments Are a Major Source of Knowledge in the Social Sciences." *Science* 326 (5952): 535–38. doi:10.1126/science.1168244.
- Flay, Brian R. 1986. "Efficacy and Effectiveness Trials (and Other Phases of Research) in the Development of Health Promotion Programs." *Preventive Medicine* 15 (5): 451–474.
- Greenstone, Michael, and Ted Gayer. 2009. "Quasi-Experimental and Experimental Approaches to Environmental Economics." *Journal of Environmental Economics and Management* 57 (1): 21–44. doi:10.1016/j.jeem.2008.02.004.
- Harrell, Stephen, Theresa Beltramo, Garrick Blalock, Juliet Kyayesimira, David I. Levine, and Andrew M. Simons. 2016. "What Is a Meal?: Comparing Methods of Auditing Carbon Offset Compliance for Fuel Efficient Cookstoves." *Ecological Economics* 128: 8–16.
- Johnson, Michael, Rufus Edwards, Victor Berrueta, and Omar Masera. 2010. "New Approaches to Performance Testing of Improved Cookstoves." *Environmental Science and Technology* 44 (1): 368–374. doi:10.1021/es9013294.
- Leonard, Kenneth L. 2008. "Is Patient Satisfaction Sensitive to Changes in the Quality of Care? An Exploitation of the Hawthorne Effect." *Journal of Health Economics* 27 (2): 444–59. doi:10.1016/j.jhealeco.2007.07.004.
- Leonard, Kenneth L, and Melkiory C Masatu. 2006. "Outpatient Process Quality Evaluation and the Hawthorne Effect." *Social Science & Medicine* 63 (9): 2330–40. doi:10.1016/j.socscimed.2006.06.003.
- Leonard, Kenneth L., and Melkiory C. Masatu. 2010. "Using the Hawthorne Effect to Examine the Gap between a Doctor's Best Possible Practice and Actual

- Performance.” *Journal of Development Economics* 93 (2): 226–234. doi:10.1016/j.jdeveco.2009.11.001.
- Levine, David I, Theresa Beltramo, Garrick Blalock, Carolyn Cotterman, and Andrew M. Simons. 2016. “What Impedes Efficient Adoption of Products? Evidence from Randomized Sales Offers for Fuel-Efficient Cookstoves in Uganda.” Berkeley, CA. <http://escholarship.org/uc/item/2cb9g10z>.
- Levitt, Steven D, and John A List. 2007. “What Do Laboratory Experiments Measuring Social Preferences Reveal About the Real World?” *Journal of Economic Perspectives* 21 (2): 153–174.
- . 2011. “Was There Really a Hawthorne Effect at the Hawthorne Plant? An Analysis of the Original Illumination Experiments.” *American Economic Journal: Applied Economics* 3 (1): 224–238.
- Miller, Grant, and A. Mushfiq Mobarak. 2013. “Gender Differences in Preferences, Intra-Household Externalities, and Low Demand for Improved Cookstoves.” National Bureau of Economic Research. <http://www.nber.org/papers/w18964>.
- Mueller, Florian, Frank Vetere, Martin R Gibbs, Stefan Agamanolis, and Jennifer Sheridan. 2010. “Jogging over a Distance: The Influence of Design in Parallel Exertion Games.” *Proceedings of the 5th ACM SIGGRAPH Symposium on Video Games*, 63–68.
- Mukhopadhyay, Rupak, Sankar Sambandam, Ajay Pillarisetti, Darby Jack, Krishnendu Mukhopadhyay, Kalpana Balakrishnan, Mayur Vaswani, et al. 2012. “Cooking Practices, Air Quality, and the Acceptability of Advanced Cookstoves in Haryana, India: An Exploratory Study to Inform Large-Scale Interventions.” *Global Health Action* 5: 19016–19016.
- Muralidharan, Karthik, and Venkatesh Sundararaman. 2010. “The Impact of Diagnostic Feedback to Teachers on Student Learning: Experimental Evidence from India.” *Economic Journal* 120 (546): F187–F203. doi:10.1111/j.1468-0297.2010.02373.x.
- Pillarisetti, Ajay, Mayur Vaswani, Darby Jack, Kalpana Balakrishnan, Michael N Bates, Narendra K Arora, and Kirk R Smith. 2014. “Patterns of Stove Usage after Introduction of an Advanced Cookstove: The Long-Term Application of Household Sensors.” *Environmental Science & Technology* 48 (24): 14525–14533.
- Ruiz-Mercado, Ilse, Eduardo Canuz, and Kirk R. Smith. 2012. “Temperature Dataloggers as Stove Use Monitors (SUMs): Field Methods and Signal Analysis.” *Biomass and Bioenergy* 47: 459–68. doi:10.1016/j.biombioe.2012.09.003.
- Ruiz-Mercado, Ilse, Eduardo Canuz, Joan L. Walker, and Kirk R. Smith. 2013. “Quantitative Metrics of Stove Adoption Using Stove Use Monitors (SUMs).” *Biomass and Bioenergy* 57 (October): 136–48.
- Ruiz-Mercado, Ilse, Nick L Lam, Eduardo Canuz, Gilberto Davila, and Kirk R Smith. 2008. “Low-Cost Temperature Loggers as Stove Use Monitors (SUMs).” *Boiling Point* 55: 16–18.
- Ruiz-Mercado, Ilse, Omar Masera, Hilda Zamora, and Kirk R. Smith. 2011. “Adoption and Sustained Use of Improved Cookstoves.” *Energy Policy* 39 (12). Elsevier: 7557–66. doi:10.1016/j.enpol.2011.03.028.

- Schwartz, Daniel, Baruch Fischhoff, Tamar Krishnamurti, and Fallaw Sowell. 2013. "The Hawthorne Effect and Energy Awareness." *Proceedings of the National Academy of Sciences of the United States of America* 110 (38): 15242–46. doi:10.1073/pnas.1301687110/.
- Simons, Andrew M., Theresa Beltramo, Garrick Blalock, and David I. Levine. 2014. "Comparing Methods for Signal Analysis of Temperature Readings from Stove Use Monitors." *Biomass and Bioenergy* 70: 476–88.
- Smith, Eleanor. 2014. "Better Hygiene Through Humiliation." *The Atlantic*, August 13.
- Smith, Kirk, Karabi Dutta, Chaya Chengappa, P.P.S. Gusain, Omar Masera, Victor Berrueta, Rufus Edwards, Rob Bailis, and Kyra Naumoff Shields. 2007. "Monitoring and Evaluation of Improved Biomass Cookstove Programs for Indoor Air Quality and Stove Performance: Conclusions from the Household Energy and Health Project." *Energy for Sustainable Development* 11 (2): 5–18.
- Smith-Sivertsen, Tone, Esperanza Díaz, Dan Pope, Rolv T Lie, Anaite Díaz, John McCracken, Per Bakke, Byron Arana, Kirk R Smith, and Nigel Bruce. 2009. "Effect of Reducing Indoor Air Pollution on Women's Respiratory Symptoms and Lung Function: The RESPIRE Randomized Trial, Guatemala." *American Journal of Epidemiology* 170 (2): 211–20. doi:10.1093/aje/kwp100.
- Srigley, Jocelyn A, Colin D Furness, G Ross Baker, and Michael Gardam. 2014. "Quantification of the Hawthorne Effect in Hand Hygiene Compliance Monitoring Using an Electronic Monitoring System: A Retrospective Cohort Study." *BMJ Quality & Safety* 23 (12): 974–980. doi:10.1136/bmjqs-2014-003080.
- The Gold Standard Foundation. 2013. "The Gold Standard: Simplified Methodology for Efficient Cookstoves." Geneva-Cointrin, Switzerland.
- Thomas, Evan A, Christina K Barstow, Ghislaine Rosa, Fiona Majorin, and Thomas Clasen. 2013. "Use of Remotely Reporting Electronic Sensors for Assessing Use of Water Filters and Cookstoves in Rwanda." *Environmental Science & Technology* 47 (23): 13602–10. doi:10.1021/es403412x.
- Zwane, Alix Peterson, Jonathan Zinman, Eric Van Dusen, William Pariente, Clair Null, Edward Miguel, Michael Kremer, et al. 2011. "Being Surveyed Can Change Later Behavior and Related Parameter Estimates." *Proceedings of the National Academy of Sciences of the United States of America* 108 (5): 1821–26. doi:10.1073/pnas.1000776108.

**Table 1**

Regressions testing for Hawthorne effect: estimates of effects of the presence of measurement team in primary three stone fire (TSF) usage and combined Envirofit usage, the coefficients represent the change in hours cooked per day compared to hours cooked per day in the measurement week

	Primary TSF		Combined Envirofit	
	(1)	(2)	(3)	(4)
Week prior to and after measurement week constrained to be equal	1.78*** (0.46)		-2.97*** (0.60)	
Week prior to measurement week		1.17** (0.54)		-2.58*** (0.69)
Week after measurement week		2.37*** (0.63)		-3.30*** (0.64)
Household fixed effects	Yes	Yes	Yes	Yes
Observations	316	316	229	229
R-squared	0.82	0.82	0.79	0.79
Household clusters	118	118	89	89

Standard errors clustered at household level in parentheses

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

**Note:** The unit of analysis is a measurement “week” (approximately 72 hours) at a household. The specification in columns 1 and 3 imposes that the weeks prior to and after the measurement week are equal. The specification in columns 2 and 4 tests usage in the week prior to and after the measurement week separately. The coefficient estimates in column 2 are jointly significantly not equal to zero (p<0.01), but not statistically different from each other (p=0.10). The coefficient estimates in column 4 are jointly significantly not equal to zero (p<0.01), but not statistically different from each other (p=0.20).

**Table 2**

Daily mean firewood consumption, particulate matter and three stone fire usage prior to introduction of fuel-efficient stoves

	Wood Consumed (kgs) (1)	PM2.5 ( $\mu\text{g}/\text{m}^3$ ) (2)	Three Stone Fire (hours) (3)
Mean Values	8.98*** (0.33)	427.79*** (23.18)	13.95*** (1.03)
Observations	568	609	339
Household clusters	160	159	102

Standard errors clustered at household level in parentheses  
\*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$

*Note:* Columns 1, 2, and 3 present the average daily wood consumption, average daily PM2.5 reading, and the total hours of combined cooking on two three stone fires per household prior to receiving a fuel-efficient stove, respectively. Observations are at the household-day level.

**Table 3**

Estimates of biomass usage and indoor air pollution with and without in-person observers

<i>Efficacy Trial (effects of an intervention during week with observers)</i>			
	Daily hours	Biomass per hour	Total biomass (kg)
Total three stone fire usage	8.30	0.64	5.31
Total Envirofit usage	8.50	0.32	2.72
Totals:	16.80		8.03
	Daily hours	PM2.5 per hour	Total PM2.5 ( $\mu\text{g}/\text{m}^3$ )
Total three stone fire usage	8.30	30.67	254.56
Total Envirofit usage	8.50	14.97	127.25
Totals:	16.80		381.81
<i>Effectiveness Trial (effects of an intervention during week without observers)</i>			
	Daily hours	Biomass per hour	Total biomass (kg)
Total three stone fire usage	11.81	0.64	7.56
Total Envirofit usage	5.53	0.32	1.77
Totals:	17.34		9.33
	Daily hours	PM2.5 per hour	Total PM2.5 ( $\mu\text{g}/\text{m}^3$ )
Total three stone fire usage	11.81	30.67	362.21
Total Envirofit usage	5.53	14.97	82.78
Totals:	17.34		444.99

*Note:* Daily hours for the effectiveness trial are taken from the data for the week prior to the KPT. Recall that only about one fifth of the secondary three stone fires had iButtons on them at this point in our experiment. For the purpose of this illustrative table we make the assumption that households with missing values for the secondary three stone fire are equal to the mean value observed for the one fifth of the sample that had an hourly usage reading for the secondary three stone fire. Daily hours for the efficacy trial are based on the Hawthorne effects presented in Table 1. The consumption rates for biomass and PM2.5 with the three stone fires are calculated prior to the introduction of fuel-efficient stoves using the values in Table 2 ( $8.98 \text{ kg}/13.95 \text{ hours} = 0.64 \text{ kgs}/\text{hour}$  and  $427.79 \mu\text{g}/\text{m}^3/13.95 \text{ hours} = 30.67 \mu\text{g}/\text{m}^3/\text{hour}$ ). The consumption rates for the Envirofit G3300 are calculated using the emissions testing report in Figure 1 ( $0.64 \text{ kgs}/\text{hour} * 0.499 = 0.32 \text{ kgs}/\text{hour}$  and  $30.67 \mu\text{g}/\text{m}^3/\text{hour} * 0.488 = 14.97 \mu\text{g}/\text{m}^3/\text{hour}$ ).

**Table 4**  
Bias introduced by the Hawthorne effect

	Daily cooking (hours)	Total biomass (kg)	Total PM2.5 ( $\mu\text{g}/\text{m}^3$ )
Baseline	13.95	8.98	427.79
Efficacy (observers present)	16.80	8.03	381.81
Effectiveness (no observers)	17.34	9.33	444.99

*Note:* These calculations are illustrative based on the mean values of data collected in the field and the emissions and performance report performed in a laboratory. These calculations assume a similarly sized Hawthorne effect on the secondary three stone fire as what we observed on the primary three stone fire.

**Table 5**  
SUMs Device Attrition: linear probability model of missing SUMs data on fuel wood use, count of daily people cooked for, and count of daily meals cooked

	TSF1 (1)	ENV1 (2)	ENV2 (3)	TSF1 (4)	ENV1 (5)	ENV2 (6)	TSF1 (7)	ENV1 (8)	ENV2 (9)
Daily household wood use (kg)	-0.02 (0.01)	-0.01 (0.02)	-0.01 (0.01)						
People cooked for daily				-0.04* (0.02)	-0.00 (0.02)	-0.00 (0.02)			
Meals cooked daily							-0.09 (0.06)	-0.07 (0.05)	0.04 (0.06)
Observations	145	141	134	151	145	138	151	145	138
R-squared	0.01	0.00	0.01	0.04	0.00	0.00	0.01	0.01	0.00
Parish clusters	14	14	14	14	14	14	14	14	14

Standard errors clustered at parish level in parentheses

\*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$

*Note:* The dependant variable is a 0/1 for whether we have the SUMs temperature data during the endline for a specified stove type (TSF1 = primary three stone fire, ENV1 = first Envirofit, ENV2 = second Envirofit). The overall sample only includes stoves that we placed a SUMs device on during the endline. The daily wood weights and counts of cooking practices are averaged across the KPT measurement week. To account for possible correlation in how data was collected by the measurement team, we cluster standard errors at the parish level because the KPT measurement team spent a week at a time in a given parish.

**Figure 1**  
 Certified Emissions and Performance Report for Envirofit G3300

April 27, 2011



DEPARTMENT OF  
 MECHANICAL ENGINEERING  
 COLORADO STATE UNIVERSITY

1374 CAMPUS DELIVERY  
 FORT COLLINS, CO  
 80523-1374  
 970.491.4796  
 970.491.4799 (F)  
 WWW.EECL.COLOSTATE.EDU

**Emissions and Performance Report**

The stove listed below has been tested in accordance with the “*Emissions and Performance Test Protocol*”, with emissions measurements based on the biomass stove testing protocol developed by Colorado State University (available at [www.eecl.colostate.edu](http://www.eecl.colostate.edu)). Percent improvements are calculated from three-stone fire performance data collected at Colorado State University.

<b>Stove Manufacturer:</b>	<b>Envirofit International</b>
<b>Stove Model:</b>	<b>G-3300</b>
<b>Test Dates:</b>	<b>4/4/2011-4/22/2011</b>
<b>Average CO emissions (grams):</b>	<b>18.7</b>
<b>80% Confidence Interval:</b>	<b>17.7-19.7</b>
<b>Percent Improvement:</b>	<b>65.30%</b>
<b>Average PM emissions (milligrams):</b>	<b>995</b>
<b>80% Confidence Interval:</b>	<b>944-1046</b>
<b>Percent Improvement:</b>	<b>51.20%</b>
<b>Average Fuel use (grams):</b>	<b>596.7</b>
<b>80% Confidence Interval:</b>	<b>591.6-601.7</b>
<b>Percent Improvement:</b>	<b>50.10%</b>
<b>Average Thermal efficiency:</b>	<b>32.6</b>
<b>80% Confidence Interval:</b>	<b>32.3-32.8</b>
<b>Percent Improvement:</b>	<b>105.20%</b>
<b>High Power (kW):</b>	<b>3.3</b>
<b>80% Confidence Interval:</b>	<b>3.3-3.4</b>
<b>Low Power (kW):</b>	<b>1.9</b>
<b>_____ 80% Confidence Interval:</b>	<b>1.8-1.9</b>

The above results are certified by the Engines and Energy Conversion Laboratory at Colorado State University. All claims beyond the above data are the responsibility of the manufacturer.

**Morgan DeFoort**  
 EECL Co-Director  
 Technical Lead, Biomass Stoves Testing Program

*Note:* The report can be downloaded at <http://www.envirofit.org/images/products/pdf/g3300/G3300Cert.pdf>

**Figure A1**

SUM holder designed to encase the stove use monitor to protect it from malfunctions when exceeding temperatures of 85 degrees Celsius



**Figure A2**

Arrows mark the placement of SUMs on three stone fire and Envirofit



(a) Three Stone Fire

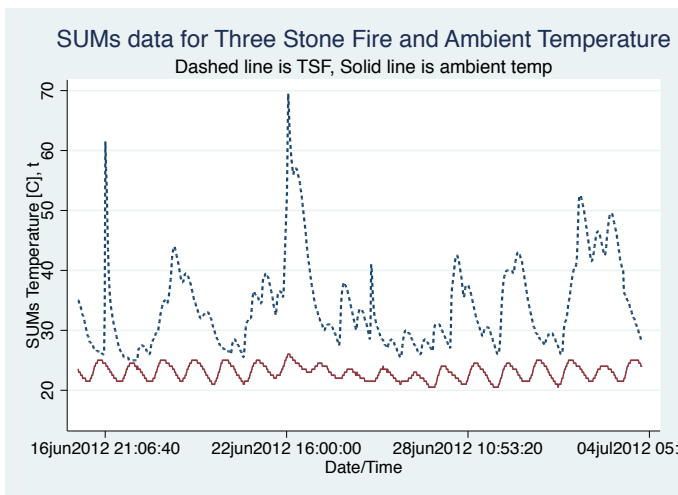


(b) Envirofit

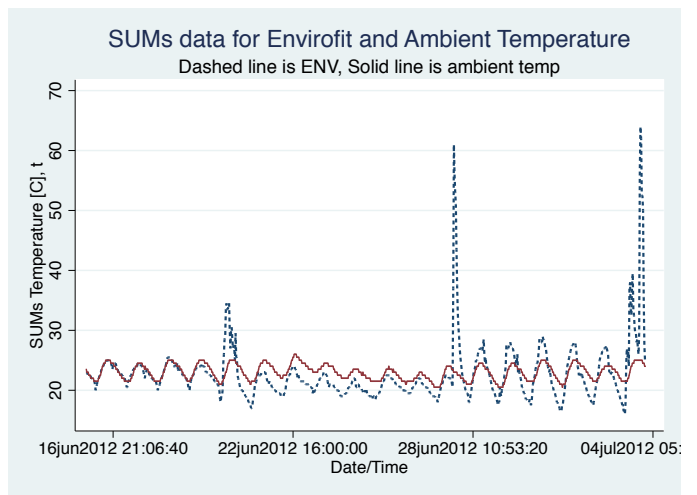


**Figure A3**

Example of household level SUMs temperature data in same household at same times



(a) Three Stone Fire



(b) Envirofit