



Complementarities in behavioral interventions: Evidence from a field experiment on resource conservation[☆]

Ximeng Fang^{a,*}, Lorenz Goette^b, Bettina Rockenbach^{c,d}, Matthias Sutter^{d,c,e},
Verena Tiefenbeck^f, Samuel Schoeb^g, Thorsten Staake^{g,h}

^a University of Oxford, United Kingdom

^b National University of Singapore, Singapore

^c University of Cologne, Germany

^d Max Planck Institute for Research into Collective Goods, Germany

^e University of Innsbruck, Austria

^f Friedrich-Alexander Universität Nürnberg-Erlangen, Germany

^g University of Bamberg, Germany

^h ETH Zurich, Switzerland

ARTICLE INFO

JEL classification:

D83

D90

Q41

Keywords:

Behavioral public policy

Pro-environmental behavior

Limited attention

Information provision

Real-time feedback

Policy interactions

ABSTRACT

Behavioral policy often aims at influencing behavior by mitigating biases due to, e.g., imperfect information or inattention. We study how this is affected by the simultaneous presence of multiple biases arising from different sources, through a field experiment on resource conservation in an energy- and water-intensive everyday activity (showering). One intervention, shower energy reports, primarily targeted knowledge about environmental impacts; another intervention, real-time feedback, primarily targeted salience of resource use. We find a striking complementarity. While only the latter induced significant conservation effects when implemented in isolation, each intervention became more effective when implemented jointly. This is consistent with predictions from a theoretical framework that highlights the importance of targeting all relevant sources of bias to achieve behavioral change.

1. Introduction

Amidst growing concern about climate change and resource scarcity, many individuals intend to make personal sacrifices to protect the environment; yet they often fail to act pro-environmentally in their everyday lives (Kollmuss and Agyeman, 2002; Frederiks et al., 2015). This gap between intentions and actions can result from a multiplicity of behavioral frictions and biases. For instance, consumers tend to underestimate the impact of highly resource-intensive activities (Attari et al., 2010; Attari, 2014; Imai et al., 2022), and they may also not be

fully attentive to their resource use (Allcott, 2016; Tiefenbeck et al., 2018). Importantly, when biased behavior arises from multiple different sources at the same time, this could not only prevent individuals from acting on their intrinsic prosocial or pro-environmental motives, but also mute their response to policy interventions that only address a subset of all relevant biases.

This problem that comes with multiple biases is reminiscent of the Anna Karenina principle, which states that failure in just one factor out

[☆] We are indebted to our research assistant team (in particular Benedikt Kauf and Kristina Steinbrecher) for their indefatigable support in running the field experiment, and to Lim Zhi Hao for editing. We would further like to thank audiences in Bonn, at the EEA Virtual 2020, the ECONtribute Retreat 2020, the ESA Dijon, the IAREP-SABE Dublin, the EAERE Manchester, the 3rd CRC TR 224 Conference, the PCBS 2019 Prague, the 12th RGS Doctoral Conference, and the 2nd BoMa Ph.D. Workshop for helpful comments. Financial support by the University of Cologne Forum “Energy”, the SUN Institute Environment and Sustainability, and by the Deutsche Forschungsgemeinschaft (DFG), Germany through CRC TR 224 (Project B07) is gratefully acknowledged. The 2019 supplementary survey was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation), Germany under Germany’s Excellence Strategy – EXC 2126/1– 390838866. This study is registered in the AEA RCT Registry under the identification number “AEARCTR-0004865”.

* Corresponding author.

E-mail address: ximeng.fang@sbs.ox.ac.uk (X. Fang).

of many can lead to failure of an objective as a whole.¹ For example, providing information to correct misperceptions of environmental impacts has little effect on behavior if individuals remain inattentive and exhibit self-control problems, status quo bias, and so on.

Conversely, drawing attention to environmental impacts only has a muted effect on behavior if agents remain unaware of the true extent of the externalities caused by their actions. In this example, addressing both information problems and biases due to, e.g., limited attention, could produce synergies in the form of positive interaction effects. More generally, combining interventions that focus on different biases each could result in complementarities, defined as *each* intervention becoming more effective when implemented in conjunction with the other(s) than in isolation (Coe and Snower, 1997). While the use of combined interventions is widespread in (behavioral) public policy, less is known about when and why one may expect interventions to be complements, which can be crucial for guiding effective policy design. In this paper, we highlight the role of multiple biases.

We report evidence from a three-month randomized field experiment in which we used two well-studied behavioral policy tools to encourage resource conservation in an energy- and water-intensive everyday activity, namely showering. Our interventions were designed in such a way that they target different potential sources of biased behavior. The first intervention, shower energy reports, inspired by the Opower home energy reports (Allcott, 2011), were primarily aimed at closing knowledge gaps about environmental impacts by providing information on water use as well as on energy use and CO₂ emissions due to water heating. The second intervention, real-time feedback, provided immediately visible and salient information on water consumption – but not energy use or CO₂ emissions – through a smart meter display (Tiefenbeck et al., 2018), and could thus help individuals focus their attention while they engaged in the activity. Crucially, we implemented a complete 2 × 2-design to evaluate both the combined intervention as well as each intervention in isolation. Our main finding is that implementing the interventions jointly seemed to result in a super-additive boost of resource conservation effects compared to their effects when implemented in isolation. This is in line with the idea that both information- and attention-based mechanisms might have been necessary to achieve behavioral change in our context.²

To formalize our arguments, we introduce a novel theoretical framework in which overconsumption can arise from multiple sources of bias (e.g., imperfect information, limited attention). Each of the biases acts akin to a discount factor and agents are prevented from incorporating the full marginal costs of resource use by the product of all biases. A key implication is that when agents suffer not just from one but from multiple independent biases (à la Anna Karenina), then interventions can become complements if they each focus on a different behavioral mechanism. The intuition is simple: the more unbiased an agent is in one dimension, the larger is the impact of reducing bias in another dimension. For example, the more attention an agent pays to her resource use behavior, the more likely it is that she will actually change her behavior when learning that the environmental impact is more negative than previously thought. Thus, in this example, mitigating both attention and information problems can have mutually reinforcing effects. This interaction mechanism is absent when two interventions

¹ The Anna Karenina principle is inspired by the opening phrase of Leo Tolstoy's novel *Anna Karenina*: "All happy families are alike; each unhappy family is unhappy in its own way." (Tolstoy, 2003). One might cheekily adapt this principle to our context by stating a slightly modified version: All unbiased agents are alike; all biased agents are biased in their own way.

² Complementarity can also arise if our interventions do not exactly work through the described mechanisms, as long as they sufficiently differ from each other in their targeted bias. For example, real-time feedback could be interpreted as facilitating learning or optimization, and this information can be complementary to the information on CO₂ emissions provided through shower energy reports.

operate mostly through the same behavioral channel, e.g., if they provide the same type of information.

Resource usage in the shower offers a useful context for studying complementarities in behavioral interventions, for several reasons. First, showering is resource-intensive: an average shower in our sample required 2.2 kWh of energy to heat up 38 L of water, which corresponds to about 10% of the average residential energy use and 30% of the average water consumption per capita and day in Germany, where we conducted our study.³ Second, most individuals underestimate the CO₂ emissions caused by water heating for showering – by as much as 89% on average based on our own survey data –, which creates scope for conservation through belief correction (Byrne et al., 2018). Third, showering is also prone to behavioral biases like limited attention and self-control problems, as the pleasure of a warm shower is salient and immediate, whereas the cost of resource use seems abstract and is hard to keep track of. Hence, individuals may not fully engage in conservation efforts unless they are informed about the actual impact of their behavior *and* keep environmental concerns on top of their minds while showering.

We conducted our field experiment in student dormitories in the cities of Bonn and Cologne, Germany, in the winter term 2016/17. A total of 351 students participated in our experiment, all of them living in single-person dorm apartments with a private bathroom. For the duration of our study, from early December 2016 until early March 2017, each participant was equipped with a smart shower meter that recorded detailed data of each shower taken. Subjects were randomly assigned into one of four experimental conditions: no intervention (CON group), shower energy reports only (SER group), real-time feedback only (RTF group), or both interventions combined (DUAL group). After an initial baseline stage, the smart meter started displaying real-time feedback on water use for subjects in RTF and DUAL, and about halfway into the study, we further started sending individualized shower energy reports via email to subjects in SER and DUAL, using data uploaded from the smart meters. This staggered design allows us to identify treatment effects of each intervention regime in a difference-in-differences setup.

Our empirical results show that, compared to the control group, subjects in the RTF group reduced their energy (water) consumption by about 0.4 kWh (6.3 L) per shower, which corresponds to 17–18% of baseline resource use. This treatment effect remains stable over the entire 3-month duration of the study. Energy reports in isolation (SER group) did not lead to any statistically detectable conservation effect. However, in line with our hypothesis, we observe a striking complementarity between the two interventions. Combining energy reports with real-time feedback (DUAL group) *further* increased the treatment effect of real-time feedback in isolation by an additional 0.23 kWh of energy (3.8 L of water) per shower. Thus, it seems that the shower energy reports in our context could only start to unfold their potential when subjects were in an enhanced choice environment where their resource usage was immediately visible. We find no evidence of adjustments on the extensive margin, i.e., the number of showers people take.

The additional reduction of resource use in the DUAL group was not driven by short-lived boosts directly after receiving a shower energy report, but rather seemed to unfold over time, which speaks against Hawthorne or pure reminder effects as the underlying mechanism. Data from baseline and endline questionnaires shows that both interventions helped subjects form more precise beliefs about their own

³ Source: German Federal Statistical Office (<https://www.destatis.de/EN/Themes/Society-Environment/Environment/Environmental-Economic-Accounting/private-households/Tables/energy-consumption-households.html> and https://www.destatis.de/EN/Themes/Society-Environment/Environment/Water-Management/_node.html). About 70% of energy for room and water heating in Germany was generated from fossil fuels at the time of our study (Dena, 2016).

water use in the shower and that there is no evidence that subjects in the DUAL group read their reports more carefully than subjects in the SER group. Supplementary survey results from a comparable sample further suggest that information included in shower energy reports also induces drastic (upward) updates in beliefs about CO₂ emissions due to warm water consumption in the shower. Hence, the null result for shower energy reports in isolation was unlikely due to lack of learning. Instead, it seems that in the absence of real-time feedback, inattention and lack of immediate visibility have prevented knowledge gains about environmental impacts from translating into effective conservation behavior.

Overall, our findings are consistent with the theoretical argument that in the presence of multiple biases, different behavioral interventions can become complements, because targeting one source of bias (e.g., imperfect information) becomes more effective when residual biases (e.g., inattention, present bias) are also mitigated, and vice versa. One implication is that lack of evidence for effectiveness of an intervention in isolation – such as in the case of shower energy reports in our study – does not mean that it cannot be effective in an enhanced policy environment that also takes into consideration further behavioral mechanisms. Appropriate policy bundling may thus increase the cost-effectiveness of interventions beyond what can be achieved with piecemeal approaches.

Our study builds on important previous contributions that have investigated the effects of similar interventions on household resource conservation.⁴ For example, in an influential evaluation of the Opower home energy reports, which provide information on aggregate electricity use to millions of U.S. households, Allcott (2011) reports an average household-level conservation effect of 2%, or about 0.62 kWh per day, although effects might be smaller in countries with lower baseline consumption (Andor et al., 2020) or when monetary incentives to save energy are low (Myers and Souza, 2020). Our SER intervention is inspired by these home energy reports. One key difference is that our reports only provide information on one specific activity (showering) instead of aggregate household consumption, as disaggregated feedback could enable better learning and thus stronger conservation responses in the targeted activities (Gerster et al., 2020), in particular when provided in shorter time intervals or even in real time. Tiefenbeck et al. (2018) provide real-time feedback in the shower through the same type of smart meter that we use in this study and document a conservation effect of 22% (0.6 kWh less energy and 9 L less water per shower). These results also replicate in a sample without monetary incentives and without self-selection into the study (Tiefenbeck et al., 2019). One important mechanism of real-time feedback is that it draws immediate attention to resource consumption by making it more salient (Bordalo et al., 2022). This study addresses the question is whether this can be used to complement interventions that aim to encourage pro-environmental action through other mechanisms, like more detailed information provision or social norms, and could thus benefit from generally higher attention to relevant behaviors.

We further relate to a number of other studies that test a combination of interventions, and especially to studies on pro-environmental behavior that also consider the idea that policy measures might become more effective when implemented in conjunction with others — although it should be noted that some studies lack a complete

⁴ Pro-environmental interventions have drawn from a broad set of instruments such as information provision, social norms, goal-setting, etc. For reviews, see e.g. Abrahamse et al. (2005), Fischer (2008), Delmas et al. (2013), Karlin et al. (2015), Andor and Fels (2018), Carlsson et al. (2021), and Khanna et al. (2021). Information provision in particular is often regarded as a promising policy lever, as individuals often misperceive the environmental impact of everyday activities (Attari et al., 2010; Attari, 2014; Camilleri et al., 2019) and tend to engage in relatively ineffective conservation measures (Gardner and Stern, 2008; Tonke, 2019).

experimental design required to identify interaction effects.⁵ For example, Jessoe and Rapson (2014) find that peak pricing schemes may only reduce peak electricity usage for households who have been outfitted with in-home-displays. Other recent studies who investigate the combination of financial incentives and behavioral interventions tend to find that they affect behavior along different margins or for different subpopulations, but find no conclusive patterns with regard to interaction effects (List et al., 2017; Holladay et al., 2019; Giaccherini et al., 2020; Fanghella et al., 2021). Hahn et al. (2016) test the individual and combined effects of social comparisons and loss framing on take-up of water-efficient technology as well as general household water consumption, but the results for interaction effects are mixed. Brandon et al. (2019) evaluate the interaction effect of two behavioral interventions on household energy conservation, home energy reports and “peak energy reports”, which provide feedback and social norms for households’ peak electricity use. As both interventions are very similar and likely operate through similar behavioral channels, it is not clear whether one should expect any interaction effect. Indeed, Brandon et al. find neither strong evidence for complementarity nor substitutability.⁶ While we provide a novel case study on the interaction of two specific types of behavioral interventions, our main contribution to this strand of literature is that we attempt to make a step towards understanding mechanisms that systematically lead different policy interventions to become complements, both theoretically and empirically. Specifically, our study highlights the role of multiple biases arising from different sources leading to an Anna Karenina effect. These insights can be adapted to guide hypothesis formation about policy interactions in other contexts as well.

The remainder of this paper is structured as follows: Section 2 introduces the theoretical framework for policy interactions under multiple biases. Section 3 describes the experimental setup and derives behavioral predictions. Section 4 presents our data as well as descriptive statistics. Section 5 lays out our empirical approach and Section 6 presents our main empirical results. In Section 7, we analyze the potential mechanisms underlying the results. Section 8 concludes.

2. Theoretical framework

We begin by introducing a stylized framework to formalize our argument of how complementarities in behavioral interventions can arise in settings where biased behavior arises from multiple sources, e.g. imperfect information, limited attention, present bias.

⁵ See, for example, the review by Khanna et al. (2021). Combined interventions are also used in other contexts than pro-environmental behavior. For example, in development economics, a number of studies experimentally test the combined effect of different interventions on financial savings (Dupas and Robinson, 2013; Jamison et al., 2014), education (Mbiti et al., 2019), risky sexual behavior (Duflo et al., 2015; Dupas et al., 2018), demand for health products (Ashraf et al., 2013), or immunization (Banerjee et al., 2021). Many of these studies, however, cannot explicitly test policy interactions, as they lack a complete factorial design (Muralidharan et al., 2020), and none of them asks more generally if or why different interventions can be complements if they target separate mechanisms. One notable study is by Mbiti et al. (2019), who find complementarities between providing school grants and adding teacher incentives in improving children’s educational outcomes. Another study by Banerjee et al. (2021) employs reminders, incentives, and information ambassador interventions on a large-scale, and then uses a data-driven approach to identify the best combination; in particular, one observation is that information ambassadors seem to amplify the effect of other interventions. Hanna et al. (2014) show in an experiment in Indonesia that seaweed farmers may fail to optimize with regard to pod size unless their attention is drawn to the importance of this choice dimension.

⁶ One speculative interpretation is that, if the HERs reduced energy consumption through investments into energy-efficient technology, while the peak energy report reduce energy usage by increasing salience of peak load events, then the potential for complementarity is limited.

2.1. Basic setup

The agent (she) engages in an resource-intensive activity, say showering, and the policy objective is to reduce resource use. Her consumption level is determined by a trade-off between the consumption utility (e.g., hygiene, pleasure, opportunity costs of time) and the perceived costs of resource use (e.g., monetary costs, environmental concern). In the case of showering, both water and energy matter, and each is subject to distinct costs and externalities. For national parsimony, we focus on energy, as it captures water use as well as water heating.⁷ Thus, the agent chooses energy use level $e \geq 0$ to maximize

$$U(e) = V(e) - B \cdot ce, \tag{1}$$

where $V(e)$ is the instantaneous consumption utility and $c > 0$ is the (constant) marginal cost of energy consumption. We consider a more general convex cost function $C(e)$ in Appendix G. In addition to standard smoothness conditions, we assume that V is hump-shaped (locally increasing at 0, strictly concave, unique maximum). For simplicity, we abstract from uncertainty or dynamics. In the absence of monetary motives, as in our empirical setting, c could be interpreted as the “moral” cost the agent perceives in face of the negative externalities from energy use. However, the perceived cost is attenuated by an aggregate bias factor $B \in [0, 1)$.

Multiple sources of bias. — The aggregate B factor can be the product of a collection of separate factors. Although this is easy to generalize, it is sufficient to focus on a simple case with two sources of bias to illustrate the mechanics:

$$B = b_1 \cdot b_2. \tag{2}$$

For example, the first factor b_1 may indicate the degree to which the agent underestimates energy intensity (as shown, e.g., in Attari et al., 2010), and the second factor b_2 the degree to which she is inattentive (e.g., Tiefenbeck et al., 2018). The multiplicative form captures that any single factor can independently prevent the agent from implementing her conservation motive, akin to the Anna Karenina principle. In this example, the agent will not take into account environmental cost both if she believes her behavior has no impact ($b_1 = 0$) and if she is fully inattentive ($b_2 = 0$), either condition by itself is sufficient. In the general case of K biases, this would become $B = \prod_{k=1}^K b_k$ (see Appendix G.3). Also note that, in principle, individuals can be biased towards less consumption (i.e., $b_k > 1$), for example if they overestimate costs (e.g., Wichman, 2017; d’Adda et al., 2020).

Consumption behavior. — The agent’s choice is defined by the intersection of marginal utility and marginal costs, with the latter being diminished by the aggregate bias:

$$V'(e) = B \cdot c. \tag{3}$$

If $B < 1$, then the marginal cost is underweighted and energy use is thus biased upwards. Correcting the bias (raising B towards 1) thus leads the individual to perceive the cost of consumption more fully. Thus $\frac{\partial e}{\partial B} < 0$, since $V''(e) < 0$.

Behavioral interventions. — In this setup, we define behavioral interventions as policies that aim to change consumers’ behavior by changing B . In contrast, price-based policies such as Pigouvian taxes would be aimed at increasing the marginal costs c that the agent faces. As $B = b_1 \cdot b_2$, there are two behavioral policy levers for reducing energy consumption: raising b_1 (e.g. providing information) and raising b_2 (e.g. enhancing salience).

⁷ Energy for water heating is determined by the amount of water and the temperature gradient. As we will report in Section 6, we find no evidence for adjustments in water temperature in our study.

2.2. Policy interaction effects

In our context, two interventions X and Y are complements if their combination reduces behavior by more than the sum of their individual effects: $\Delta e^{XY} \leq \Delta e^X + \Delta e^Y$. If they are substitutes, the inequality is reversed. Notice that even under substitutability, it can be the case that XY is more effective than either X or Y in isolation, i.e., $\Delta e^{XY} < \Delta e^X$ and $\Delta e^{XY} < \Delta e^Y$. Thus, to empirically identify interaction effects between different policy interventions, it is necessary to evaluate the effectiveness of each intervention in isolation.

The key mechanism we aim to highlight in this paper is that in the presence of multiple biases, policies that target only one bias dimension may have a limited effect on behavior, whereas the effect of combining several policy levers may be superadditive. For example, correcting perceptions of the environmental impact b_1 may only have a small impact on behavior if the attention parameter b_2 is still close to zero. There is a simple geometric interpretation to illustrate this: the overall bias parameter B , defined in Eq. (2), can be thought of as the area of a rectangle with sides of lengths b_1 and b_2 (see Fig. 1a). The larger the rectangle the lower the resulting energy consumption will be. Now suppose that b_1 is exogenously increased by δ_1 . The resulting increase in B will be $\delta_1 b_2$, as it is attenuated by b_2 . Analogously, an exogenous increase of δ_2 in the dimension of b_2 results in an aggregate change of $\delta_2 b_1$. The effect of jointly increasing b_1 and b_2 by the same amounts, however, results in an overall change of

$$\Delta B = \delta_1 b_2 + \delta_2 b_1 + \delta_1 \delta_2. \tag{4}$$

There is an additional effect of size $\delta_1 \delta_2$, because a gain in one dimension also makes the improvement in the other dimension larger. Geometrically, this is represented by the top right rectangle outlined in Fig. 1b. This mechanism potentially induces complementarity between interventions that mitigate different biases each.⁸

Under what conditions does this complementarity in bias reduction $\delta_1 \delta_2$ translate into a complementarity in behavior between interventions X and Y ? As formally derived in Appendix G.1, a second-order Taylor approximation yields

$$\Phi^{XY} := \Delta e^{XY} - \Delta e^X - \Delta e^Y \approx \left[\frac{\partial e}{\partial B} + \frac{\partial^2 e}{\partial B^2} b_1 b_2 \right] \delta_1 \delta_2. \tag{5}$$

The term $\frac{\partial e}{\partial B}$ in Eq. (5) is negative and scales with $\delta_1 \delta_2$, thus creating scope for complementarity in behavior. The second term in brackets reflects the change in the slope of $\frac{\partial e}{\partial B}$. Intuitively, one would expect a diminishing responsiveness to bias mitigation ($\frac{\partial^2 e}{\partial B^2} > 0$), as the more the agent already reduces her consumption the less room for further reduction she has. This corresponds to $V(e)$ having a positive third derivative. However, if either b_1 or b_2 is sufficiently close to zero, the first-order effect dominates and complementarities in bias reduction translate into complementarities in observable behavior. One might call this the Anna Karenina condition: the more biased an agent is, in multiple ways, the more effective it is to target the biases simultaneously.

2.3. Policy implications

Lastly, we explore implications for policy makers who attach a social cost $\gamma > 0$ to every unit of e (e.g., due to externalities), in addition to the private cost c to the consumer. If the only policy goal was to reduce resource use e , then complementarities in behavior would directly carry over to policy benefits. However, the prevailing view

⁸ We focus here on the case of two “pure” interventions that only target b_1 or b_2 , respectively. In practice, many interventions may affect not just one but several biases. In Appendix G, we show that imperfectly targeted interventions can still produce complementarities if sufficiently different from another.

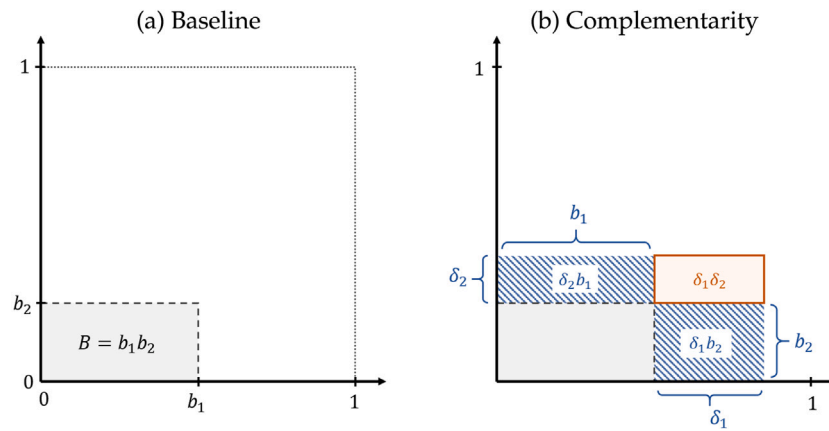


Fig. 1. Depiction of example interventions.

Notes. The gray rectangle in Figure (a) illustrates the aggregate bias B as defined in Eq. (2) without any intervention in place. Figure (b) illustrates the increase in B through exogenous interventions in each dimension.

is to define welfare over true consumer utility (e.g., Bernheim and Taubinsky, 2018), so

$$W(e) = V(e) - ce - \gamma e. \tag{6}$$

It is straightforward to show that welfare increases as B goes towards 1, as $\frac{\partial W}{\partial B} > 0$ as long as $B < 1$. These welfare gains result from both a reduction in externalities and in “internalities”. Taking a second-order Taylor approximation, we can examine how welfare is affected by combining interventions X and Y :

$$\Delta W^{XY} - \Delta W^X - \Delta W^Y \approx [(B - 1)c - \gamma] \Phi^{XY} + Bc\delta_1\delta_2 \frac{\partial e}{\partial B} \geq 0. \tag{7}$$

Appendix G.2 contains a detailed derivation. The first term on the right-hand side is positive if $\Phi^{XY} < 0$. However, the second term is always negative, reflecting that gains in consumer surplus decrease as B increases. This implies that complementarity in behavior (i.e., if $\Phi^{XY} < 0$) is a necessary but not sufficient condition for complementarity in welfare. If, in addition, either $\Delta e^X \approx 0$ or $\Delta e^Y \approx 0$, our model implies that one of the b 's is equal to zero. Thus, within this framework, the joint finding of policy complementarities in behavior *and* one of the interventions being completely ineffective on its own may serve as a sufficient condition that complementarities also exist in terms of welfare.

Eq. (7) further suggests that policy makers with a binding budget constraints face a trade-off between targeting a larger share of the population or enriching the policy bundle. For example, suppose the policy maker has a budget g , and that implementing either X or Y in the whole population requires social costs κ^X or κ^Y greater than her budget. Hence, she could cover g/κ^X of households with X , or g/κ^Y of households with Y . Now assume that $\frac{\Delta W^X}{\kappa^X} \geq \frac{\Delta W^Y}{\kappa^Y}$ (X generates a larger welfare gain than Y) and $\frac{\Delta W^X}{\kappa^X} > 1$ (the welfare gain for X is positive). In this case, intervention Y seems unattractive at first glance relative to X . However, our model implies that combining both policies in a bundle (X, Y) at cost $\kappa_X + \kappa_Y$ to fewer households can be preferable to treating a larger number with X in isolation. Formally, the condition is $\frac{\Delta W^{XY} - \Delta W^X}{\kappa^X} \geq \frac{\Delta W^X}{\kappa^X}$. This last condition can only hold if $\Delta W^{XY} - \Delta W^X$ is substantially larger than ΔW^Y , i.e. if complementarities are sufficiently strong.⁹ Thus, complementarities can create a unique rationale for unequal policy distribution to improve cost-effectiveness.

⁹ One implicit assumption we make here is that the costs of interventions are additive, i.e., $\kappa^{XY} = \kappa^X + \kappa^Y$. Costs can in principle be super- or sub-additive, purely for technological reasons, and hence we abstract from this here. For example, in our study, costs of recruiting participants and setting up the technical infrastructure are fixed, leading to economies of scope when adding interventions.

3. Experimental setup

Our field experiment was conducted from early December 2016 to late February/early March 2017 in a sample of students living in dormitory apartments. Each participant was equipped with a smart meter that measured individual energy and water consumption in the shower over the entire study duration. We then evaluated the effect of two different interventions, real-time feedback and shower energy reports, on resource conservation behavior. To test for complementarity, we further implemented a combined intervention in which subjects received both real-time feedback and shower energy reports.

3.1. Recruitment of participants

We selected six student dormitory sites in Bonn and Cologne for our sample, and ran the study from early December 2016 to early March 2017. All dormitory residents were students at the University of Bonn, the University of Cologne, or at various smaller universities in the cities. We recruited our subjects from the pool of dorm tenants living in single-person apartments with private bathroom, as this allows us to precisely measure the resource use of each individual. One noteworthy feature of our sample is that subjects have no direct monetary incentives to conserve energy or water, because they pay a flat monthly rent that includes all utility bills. Hence, any observed conservation response would be solely driven by non-monetary motives and unconfounded by income effects.

To participate in the study, residents had to actively agree based on the principle of informed consent. Two additional criteria were levied: subject should not have lengthy absences planned within the intended study period (except during Christmas vacation), and they should own a smartphone compatible with Bluetooth 4.0, which was necessary for implementing the shower energy reports.

The recruiting process started around mid-October 2016. Posters and flyers informed residents of the selected dormitories about the upcoming study, and our local research assistant teams engaged in door-to-door recruiting. Interested students had to complete an online registration survey to provide required information and to give their consent to the collection and analysis of data on their showering behavior. It was explicitly (and truthfully) stated that we would treat any collected data confidentially and not share it with the dormitory administration. As remuneration, each participant received 20 Euros after completing the study, and ten participants were randomly drawn to receive a 300 Euro cash prize. In total, 406 students registered for the

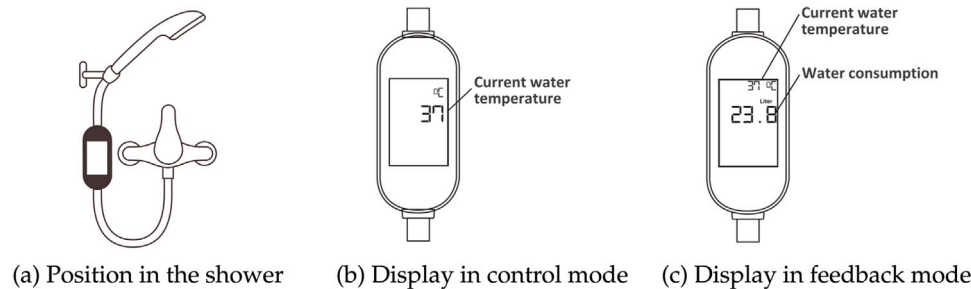


Fig. 2. Amphiro b1 smart shower meter.

study, out of which 361 met our participation criteria.¹⁰ Ten students subsequently dropped out of the study, either because they moved out of their dorm unexpectedly or because we were not able to contact them again. This leaves us with a final sample of 351 participants.

3.2. Smart shower meters and smartphone app

At the beginning of the study, starting on 5th Dec 2016, each participant was equipped with an Amphiro b1 smart shower meter that measures and records data of every water extraction in the shower. The device can be easily attached below the shower head and features a smartphone-sized liquid crystal display, which can be programmed to show various types of information (see Fig. 2a). The smart meter is small, lightweight, and needs no battery; power is generated through an integrated hydro turbine, without noticeably affecting water flow in the process. One drawback of the lack of battery is that the device is unaware of the absolute time of day: showers can only be recorded in temporal order, but without time stamps. Once the water flow in the shower starts, the smart meter is powered and begins to measure, among others, the amount of water flowing through, water temperature, and the time passed since beginning of water flow. When water flow stops, the device and its display initially remain switched on, and if water is turned on again within three minutes, it will continue its measurement seamlessly. This accounts for short breaks in water flow when applying soap or shampoo. Once water flow stops for more than three minutes, the device terminates measurement and stores the recorded information to a new data point.

We programmed the shower meters to display select pieces of information to participants in real-time, i.e., while they are taking their showers, contingent on the study progress and assigned experimental condition (as described below). In addition, we asked all participants to install the Amphiro smartphone app around week 5 of the experiment, shortly after the end of the Christmas break. The participants could use the app to upload data from their shower meters via Bluetooth.¹¹ We were then able to access the uploaded data and use it to create personalized shower energy reports. The original Amphiro smartphone app also calculates summary statistics about users' resource use in the shower, but we deactivated this feature for our study participants, so its only functionality was data uploading. One ancillary benefit of the app was that it stored time and date of each data upload, which allows us

¹⁰ The total number of all single apartments in the selected dorms was 1380 (vacancies included), thus our gross recruitment rate was about 30%. For more than half of these apartments, we never encountered the resident, so out of the students we actually managed to talk to, the majority registered for the study.

¹¹ The process was quite simple. After installing the smartphone app, subjects created an account and paired it to their shower meter. After successful pairing, the meter automatically transmitted all stored data to the app via Bluetooth whenever it was powered on and the smartphone within range.

to construct approximate time windows for each shower. About three out of four participants (72%) uploaded all data successfully, while the remaining experienced some technical problems. The most common sources of failure were problems with the Bluetooth connection or unexpected incompatibility between smartphone and app. We will come back to this issue again later.

3.3. Implementation of real-time feedback

The live tracking of water use on the shower meter display in feedback mode is what we refer to as real-time feedback, our first type of intervention. We programmed half of the smart meters as control devices and the other half as treatment devices. Control devices only displayed the current water temperature throughout the entire study (Fig. 2b). Treatment devices also started in control mode for the first ten showers, which we use to measure baseline behavior, but switched permanently to feedback mode starting from the eleventh shower. In feedback mode, the display shows both the water temperature and the amount of water used (in liters) at any time of the shower (Fig. 2c). Note that the smart meters did not provide any information on the impact of water temperature on energy consumption or on the energy-intensity of water heating more generally.

3.4. Implementation of shower energy reports

Our second intervention consists of two personalized shower energy reports. These reports were sent via e-mail and showed descriptive statistics about the subject's water and energy use in the shower, as well as information about environmental impacts. Temperature information was not included, as all subjects received this through their smart meter anyway. To allow for learning about outcomes of single showers, a graphical representation of the subject's history of water use per shower was included. The reports were constructed based on data that was uploaded by subjects through the smartphone app. We sent out additional reminders to upload data before each planned delivery, but the reports themselves were not explicitly announced. Subjects who did not manage to upload any data received a report template with blanks in place of statistical figures and graphs.

Appendix Figure A1 shows the screenshot of a typical shower energy report. After a short introductory text, subjects see a scatter plot of their history of water use per shower since the beginning of the study, including a fitted regression line to help recognize trends and averages. Below the graph, average water use (in liters) and energy use (in kWh) per shower are stated numerically. Furthermore, there is a paragraph with information on projected CO₂ emissions per year and the number of trees required to absorb the corresponding amount of CO₂. The report is formulated concisely in neutral language, to avoid any normative or moral suasion elements. In the second report, we added a social comparison component in the spirit of Allcott (2011) and Ferraro and Price (2013), see Appendix Figure A2. Specifically,

we assigned a random anonymous peer to each subject and displayed statistics on the peer's energy and water use.¹² At the bottom of each report, we included personalized link to a mini-survey that we asked subjects to fill out. The mini-survey contained three questions to elicit subjects' estimate of their water consumption per shower (absolute and relative to others). The purpose of this was twofold. First, we use the responses to verify if a subject has read the email carefully and, based on the estimate accuracy, how closely he or she paid attention to the information. Second, we use the time of survey response to determine when exactly a subject read the email.

3.5. Experimental design

We implemented a complete 2×2 design with four experimental conditions. Subjects in the control (CON) group received no intervention at all; subjects in the RTF group only received real-time feedback through the smart shower meters; subjects in the SER group only received shower energy reports; and subjects in the DUAL group received both real-time feedback and shower energy reports. Treatment assignment was randomized and the group sizes are as follows: 82 in CON, 88 in SER, 90 in RTF, 91 in DUAL.¹³

Fig. 3 illustrates the experimental design in detail. Each shower meter went through a baseline stage of ten showers, in which it only displayed the current water temperature, regardless of the experimental condition. We use these showers to measure baseline consumption. Starting from the eleventh shower (intervention stage), devices in RTF and DUAL additionally displayed water use in real-time, whereas devices in CON and SER stayed in control mode. About halfway into the study, we started sending energy reports to each subject in the SER or DUAL group; the first report was sent on 24 January 2017 and the second report on 8 February 2017, about two weeks later. We distinguish between intervention (IN) stage 1, in which real-time feedback is switched on but there were no reports yet, and intervention (IN) stage 2, which is the period that begins after subjects saw the first report.¹⁴ In order to hold interaction with experimenters constant, subjects in CON and RTF groups received placebo emails at the exact same time the shower energy reports to subjects in SER and DUAL were sent out. These subjects were asked to complete the same mini-survey that came along with the actual reports.

This staggered experimental design allows us to exploit both between- and within-subject variation to cleanly identify and efficiently estimate treatment effects of interest. The effect of real-time feedback in isolation is identified by the comparison between the RTF and CON groups in the (entire) intervention stage, or alternatively by the comparison between the pooled RTF/DUAL group and the pooled CON/SER group in IN stage 1. The effect of shower energy reports in isolation is identified by the comparison between the SER and CON groups in IN stage 2. The additional effect of shower energy reports, when combined with real-time feedback is identified by the comparison between the DUAL and RTF groups in IN stage 2. Differences between the effects of shower energy reports with and without real-time feedback identify policy interaction effects, i.e., whether the two interventions are substitutes or complements. Note that behavior in the CON group may not reflect a pure counterfactual, as subjects still receive a smart meter with temperature information as well as placebo emails, to hold experimenter interaction and Hawthorne effects fixed.

¹² The matching procedure was one-sided and ensured that each subject (except the most and the least efficient) was equally likely to see a peer with lower or higher energy use per shower.

¹³ For the exact randomization protocol, see Appendix B.

¹⁴ In practice, the distinction between IN stage 1 and 2 is not perfect, as we observe 23 subjects in our sample who had yet to complete all 10 baseline showers when the first report was sent out. If anything, this generates measurement error in our treatment indicators and thus biases estimates towards zero.

We would underestimate the effects of our interventions to the degree that subjects respond to this by itself, but any relative comparison across intervention regimes would remain valid.

3.6. Behavioral predictions

In order to derive behavioral predictions for each of our experimental groups, we first briefly discuss the channels through which each of the two interventions is likely to work. Our theoretical framework shows that the effect of each regime depends on the degree to which it succeeds in overcoming the aggregate bias, which may be the product of multiple separate factors. Furthermore, real-time feedback and shower energy reports could be complements if they are relatively specialized and operate largely through different channels.

Real-time feedback visually displays live measurement of water use in the shower. This water volume information can debias individuals' beliefs about the amount of water they use, but there is no additional information on energy use or CO₂ emissions due to water heating, so severe knowledge gaps about the environmental relevance of showering may remain. In addition, the steadily upward moving liter count is likely to significantly reduce inattention and self-control problems, as users are constantly facing the smart meter display, and the previously abstract and elusive notion of resource use suddenly becomes salient and palpable, infused with a sense of immediacy. It may also facilitate experimentation with various conservation strategies by keeping track of progress in real-time. As the RTF condition in our experiment is essentially a replication of the intervention by Tiefenbeck et al. (2018), albeit more minimalistic and in a sample without monetary incentives, we also expect to find comparable conservation effects:

Prediction 1. *Providing real-time feedback through the smart shower meter display in treatment RTF leads to a reduction in water and energy consumption in the shower.*

Shower energy reports provided personalized information about subjects' water use in the shower as well as additional information about energy use and CO₂ emissions. We therefore expect that the reports can help close knowledge gaps in these areas and thereby induce conservation behavior, since past evidence suggests that individuals tend to grossly underestimate the energy intensity associated with water heating (Attari et al., 2010). The second report also included a social comparison with a randomly assigned and anonymous peer, which might further add motivation (Allcott, 2011).

Prediction 2. *Providing information through shower energy reports in treatment SER leads to a reduction in water and energy consumption in the shower.*

As the shower energy reports are not immediately salient while showering, the effect of knowledge gains could be stifled by remaining barriers like limited attention or self-control problems that can be better targeted by real-time feedback.¹⁵ Vice versa, the effect of real-time feedback may be attenuated if subjects remain unaware of the energy and carbon intensity of warm water use. If the two interventions indeed work largely through these separate behavioral mechanisms, a combined intervention should leverage all mechanisms at the same time. As we argue in Section 2, shower energy reports and real-time feedback could therefore become complements in the sense that one

¹⁵ In principle, it is possible that participants also become more attentive about resource use even without visual aid through the smart meter, as would be predicted by models of endogenous attention when updates in beliefs about environmental impacts are sufficiently large (Hanna et al., 2014; Gabaix, 2017). However, if there is such an effect, it may prove short-lived once reports fade out of memory and resolutions cool off (Allcott and Rogers, 2014; Schwartz and Loewenstein, 2017).

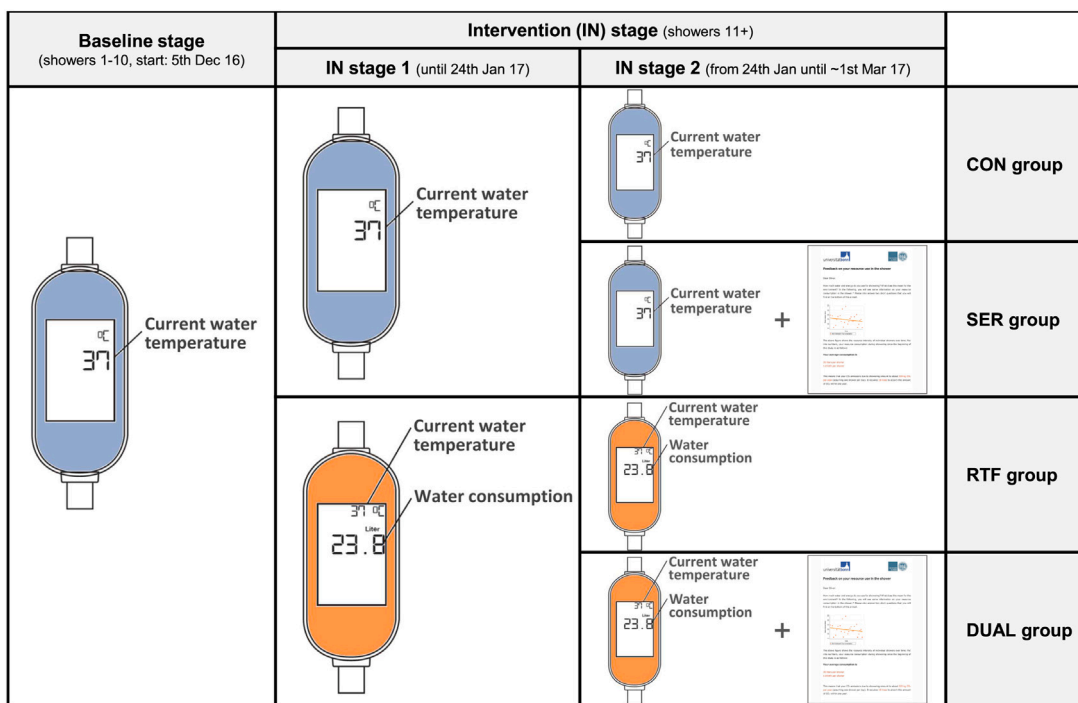


Fig. 3. Experimental design and timing of interventions.

intervention makes the other more effective when implemented jointly. Thus, we derive the following prediction:

Prediction 3. Shower energy reports in IN stage 2 lead to a larger (marginal) reduction in water and energy consumption in the shower for subjects who also receive real-time feedback (DUAL group) than for subjects who do not receive real-time feedback (SER group).

Note that in the communication with the participants, the study was primarily framed as an energy conservation study, as water scarcity was less of an issue in Germany at that time. By contrast, energy conservation and the transformation of the energy sector ranked high on Germany’s policy agenda. Nevertheless, in the subsequent analyses, we report results for both energy and water consumption. While greenhouse gas emissions arise mostly due to water heating, conserving water may be an objective in itself, especially given that climate change increases the likelihood of droughts and water stress even in parts of the world (like Germany) that previously have not suffered from water scarcity (European Environment Agency, 2021). As energy use is calculated as a product of water volume, temperature gradient, and constant factors (see 4.1 for details), the amount of energy and water consumed in the shower are highly correlated. Any change in water consumption (i.e., by shortening shower time or reducing the flow rate) affects energy consumption proportionally. The only margin of adjustment that has a different effect on these two outcomes are changes to the water temperature.

4. Data and descriptive statistics

4.1. Measurement data on resource use behavior

For every water extraction in the shower, the smart meters measured, among others, the volume of water used, its average temperature, and the average flow rate (i.e., volume per time unit). The amount of energy used was then calculated based on volume and temperature data, using the standard engineering formula for heat energy.¹⁶ Every

¹⁶ The formula for energy use of water heating is $Q = m \times c_p \times \Delta T$, with heat energy Q , mass of water m , heat capacity c_p , and ΔT the difference between

subject had a shower meter installed for the whole duration of the study, starting from early December 2016. At the end of the study, in early March 2017, we retrieved the devices and read out the data manually.¹⁷ In this way, we were able to extract an initial data set of 21,469 showers by 327 participants. Unfortunately, no data could be obtained in 24 cases, either because the device was defective or because subjects never used it, or because subjects simply disappeared without a trace (and their shower meters with them).

A number of data cleaning steps are performed before running the empirical analyses. We briefly describe the most important steps here; a more detailed documentation can be found in Appendix C. First, we drop the very first data point of each participant, as they usually started with a test run to check if the device was working. Following Tiefenbeck et al. (2018), we further drop any water extraction with volume below 4.5 L (in total 2942 extractions), as these are unlikely to be actual showers but rather minor extractions for other purposes such as cleaning. As there are rare cases in which the device can produce errors when storing data, we further remove 37 extreme outlier points, defined as such by being more than 4.5 times the subject-specific interquartile range away from the closest quartile.¹⁸ We further exclude 1 device with generally erratic data, 5 devices with fewer than 10 recorded extractions, as well as 3 devices with an abnormally large baseline consumption of 168 L or more per shower, which is about 40 L (1.5 standard deviations) away from the rest of the field. In 8 cases, the integrated temperature sensor became defective after some time,

the measured water temperature and cold water temperature (assumed to be 12 degrees Celsius). Following Tiefenbeck et al. (2018), we also assume boiler efficiency losses of 35% and distribution losses of 24%.

¹⁷ We already started retrieving some devices in late February, but as the retrieval process was drawn out over a period several days, the end of the study was in early March for most subjects.

¹⁸ We are particularly strict in only excluding the most implausible data points here. Conventionally, 1.5 or 3 times the interquartile range (IQR) are used as criterion for outliers. For a normal distribution, 4.5 times the IQR away from the nearest quartile corresponds to 6.745 standard deviations away from the mean.

and we impute missing information with the average temperature of showers taken while the sensor was still intact. The final data set used for our empirical analyses includes 17,942 showers by 318 participants.

The shower meter stores the temporal order of showers, so we can easily classify each shower into baseline or intervention stage, as real-time feedback (in the RTF and DUAL groups) started from the eleventh shower. Assigning showers to intervention stage 1 (pre-reports) or stage 2 (post-reports) is slightly trickier, as the device has no counter for global time. Fortunately, the smartphone app stores the date and time of each data upload, which allows us to construct time bounds for when a shower took place. Specifically, we know that a shower cannot have occurred after the time at which it was uploaded, and also not before the time of the last previous data batch, because otherwise it would have been uploaded then already. Combined with knowledge about the order of observations, we can assign approximate dates to each shower, assuming that the time that passes between one shower and the next remains roughly constant. For example, if three shower observations were uploaded at day t and the last previous upload occurred at day $t - 3$, then we would assign the first of these showers to day $t - 2$, the second to $t - 1$, and the third shower to day t . We instructed subjects to use the app regularly starting from 11 Jan 2017 – two weeks before the first energy report (sent out at 2:30pm on Jan, 24th) –, and sent additional reminders before each energy report email (or placebo email) was sent out.

Using this timing information from data uploads, we classify observations into pre-report showers (IN stage 1) or post-report showers (IN stage 2). In particular, we know from mini-survey response data when subjects likely read the email and use this as cutoff date. For non-responders, we use the time at which we sent the email as cutoff date instead; the response rate for the first email was 82.7%.¹⁹ Observations with upper time bound before the cutoff date are assigned to IN stage 1 and observations with lower time bound after the cutoff date are assigned to IN stage 2. For observations that fall into a range of uncertainty around the date on which the subject read/received the email, we intrapolate their dates based on the assumption that the frequency of showering was constant within that time range, and then use these intrapolated dates to assign them into the first or second intervention stage.²⁰

A final complication comes from subjects who did not manage to upload any data to the app. For these non-uploaders, we impute the timing of shower energy reports based on the assumption that it follows the same distribution for uploaders and non-uploaders. To operationalize this, we use timing information from uploaders to estimate the probability that a shower took place after receiving the first (second) report for each shower, based on its temporal order, and then assign the implied post-report probabilities to showers of non-uploaders. Appendix Figure A3 plots the estimated CDFs, and more details on the imputation procedure are provided in Appendix D. We also consider alternative definitions of intervention periods for robustness checks.

4.2. Survey data

To supplement our behavioral data on resource use in the shower, we administered several questionnaires. In the baseline survey, we collected information on individual characteristics (i.e., age, gender, etc.), perceived water use in the shower, shower comfort (i.e., how much they enjoy showering), environmental attitudes and beliefs, as

¹⁹ 47% of subjects responded to the mini-survey within the same day that the email was sent out, and 77% responded within one week.

²⁰ For example, say we know for certain that shower s occurred at 8am on Jan, 22nd (pre-report), and shower $s + 3$ occurred at 9am on Jan, 25th (post-report). This leaves the stage of showers $s + 1$ and $s + 2$ ambiguous. To assign these, we would assume that shower $s + 1$ occurred in the morning of the 23rd and shower $s + 2$ in the morning of the 24th, thus putting both showers before the first report, which was sent out at 4:30pm on the 24th.

well as a number of personality attributes (i.e., Big Five, patience, etc.). In the post-intervention survey, we again collected self-reported data on perceived water use, shower comfort, and environmental attitudes. Furthermore, we administered mini-surveys with each energy report, in which subjects were asked to estimate their resource use in the shower.

We mainly make use of information on water use perceptions, shower comfort, and environmental attitudes, and how they change in response to our interventions. Environmental attitude is elicited using four items about pro-environmental behavior and identity, e.g. “I do what is right for the environment, even when it costs more money or takes more time”.²¹ Shower comfort is elicited using five items on how much subjects enjoy showering, e.g. “I find it relaxing to take a shower”.²² We create indices for shower comfort and environmental attitude, respectively, by taking the simple average of the individual’s responses to the relevant items (rated on a 4- or 5-point Likert scale) and then normalizing to mean 0 and standard deviation 1. For perceived water consumption, we asked subjects to estimate how many liters of water they typically use when taking a shower. These estimates can then be directly compared to their actual water use as measured by the smart meter. Note that we refrained from eliciting subjects’ beliefs about energy use and carbon emissions from water heating, because we did not want to raise awareness about these issues and risk undermining the shower energy report treatments.

4.3. Sample characteristics and baseline behavior

All study participants were students at universities in Bonn or Cologne living in single-person dorm apartments, so our sample is rather homogeneous. From the 318 participants represented in our main dataset, 203 lived in a dorm in Bonn and 115 lived in a dorm in Cologne. The female share was 61 percent. Average age was 23.8 years (median 23 years), with students from all stages of their studies being represented in our sample.

Using the nine showers (the first being excluded) in the baseline stage, where only the current water temperature was displayed, we can construct measures of each subject’s baseline resource use behavior. Table 1 presents descriptive statistics about baseline energy and water use per shower, as well as shower duration (net of breaks), water temperature, and flow rate. On average, showers in the baseline stage feature 7 min of water flow, which amounts to 37.77 L of water. On average, water is heated up to a temperature of 36.14 degrees Celsius, resulting in energy consumption of 2.21 kWh per shower. There is substantial variation across showers, as observed from the standard deviations and different quantiles of the distributions. Water and energy consumption follow a right-skewed distribution, thus the median energy use per shower (1.71 kWh) is substantially lower than the mean. As the share of cold showers is extremely low in our sample (only 3.7% of showers have an average temperature of 21°C or lower), water and energy usage is almost perfectly collinear, with a Pearson correlation coefficient of 0.9755. The average flow rate of 5.71 L per minute is relatively low, likely due to dorm infrastructure not being up to modern standards — flow rates of 10–12 L per minute are more typical for German households.

4.4. Randomization checks

Our identification strategy relies on randomization producing treatment groups that are comparable with regard to observable and un-

²¹ The other items are “Environmental friendliness is part of my personal identity”, “How often do you try to conserve water?”, and “How often do you try to conserve energy?”. We also include a set of questions adapted from Nolan et al. (2008) in the baseline questionnaire.

²² The other items are “I like showering”, “For me, taking a shower is just a means to an end”, “I like to let my mind wander when I shower”, and “I try to shower as quickly as possible”.

Table 1
Descriptive statistics – baseline showers.

	Mean	Std. dev.	10th pctile	Median	90th pctile	Obs.
Energy use [kWh]	2.21	1.91	0.43	1.71	4.58	2503
Volume [L]	37.77	30.40	9.30	29.60	75.70	2503
Duration [min]	6.99	5.00	1.97	5.82	13.00	2503
Temperature [Celsius]	36.14	5.23	32.00	37.00	40.00	2477
Flow rate [l/min]	5.71	2.45	2.80	5.40	9.10	2503

Includes only showers taken in the baseline stage, i.e., first 10 showers and before subjects read/received shower energy reports. For temperature statistics, devices with broken temperature sensors are excluded. Duration is net of any breaks and calculated by dividing water volume by flow rate.

Table 2
Randomization checks and extensive margin responses.

	Panel A. Baseline averages by individual					Panel B.
	Energy use [kWh]	Volume [liter]	Duration [min]	Temperature [Celsius]	Flow rate [l/min]	Number of showers
RTF group	−0.111 (0.215)	−1.253 (3.427)	0.284 (0.597)	0.086 (0.595)	−0.124 (0.370)	−2.312 (5.183)
SER group	−0.077 (0.218)	−2.096 (3.431)	0.166 (0.547)	0.962 (0.608)	−0.441 (0.319)	3.393 (5.226)
DUAL group	−0.071 (0.226)	−1.215 (3.571)	0.126 (0.578)	0.323 (0.556)	−0.151 (0.357)	3.224 (5.861)
Constant	2.237 (0.163)	38.316 (2.539)	6.797 (0.411)	35.681 (0.447)	5.832 (0.240)	55.312 (3.698)
Observations	316	316	316	314	316	318
R-squared	0.001	0.001	0.001	0.011	0.005	0.005
F-test: p -value	0.965	0.945	0.972	0.351	0.550	0.669

Robust standard errors in parentheses. The omitted category is the CON group. For two participants, the device was not able to record information on baseline showers, but we could extract valid data on showers in later stages; hence the number of observations is only 316 in most columns. In addition, two participants with initially defective temperature sensors are excluded in column 4.

observable subject characteristics. Although it is naturally impossible to test the latter, we can check balance on observable baseline characteristics. Panel A of Table 2 shows results from regressing various measures of subjects' baseline behavior on assigned treatment groups. The differences between groups are very small and treatment assignment is insignificant for predicting any of the behavioral measures, so randomization seems to have worked well. We also check for balance along background characteristics and survey responses (see Table A.1 in Appendix A), and again find that treatment assignment is statistically insignificant. Importantly, self-reported environmental attitude and shower comfort are comparable across groups.

4.5. Number of showers

On average, we observe 56.8 showers per individual over roughly 12 weeks of our study, which corresponds to a frequency of about two showers every three days. However, the net frequency (i.e., adjusting for absences) might be closer to one shower per day, as our study period included a two weeks Christmas break. In Panel B of Table 2, we check whether the number of showers per individual differs across experimental conditions, but we find that treatments have no effect on the number of showers ($p = 0.669$). Hence, our interventions do not seem to induce adjustments along the extensive margin, and we do not need to worry about subjects compensating shorter showers with more showers, substituting behavior to other facilities (e.g. wash basin, gym showers), or about them compromising on basic hygiene needs. This means that we can make use of the full panel structure of our data and analyze (intensive-margin) water and energy conservation effects at the level of individual shower observations.

4.6. Presence of imperfect information and behavioral biases

Before moving on to the analysis of our experimental interventions, we provide suggestive evidence that individuals' resource consumption in our setting may indeed be subject to biases due to imperfect information and limited attention.

First, we make use of the pre-intervention questionnaire and compare subject's perceptions of their own water use per shower to their actual baseline water use as measured by the smart meter. Fig. 4 shows that subjects' estimates are all over the place; we cannot even reject the null hypothesis that estimated and measured water use are uncorrelated (Pearson's $\rho = 0.0925$, $p = 0.1308$). This demonstrates that subjects lack information about their own behavioral outcomes prior to any intervention.²³ Interestingly, the mean estimate (43.4 L) and median estimate (30 L) across subjects were not too far from the typical baseline water usage per shower in our sample. This is reminiscent of a "wisdom of crowds" phenomenon and suggests that, on average, our interventions should not work through debiasing beliefs about water use.

However, people may be particularly unaware about the link between water heating for showering and energy consumption, and hence CO₂ emissions. For example, (Attari et al., 2010) show that consumers are in general highly prone to underestimating the amount of energy required for heating up water (e.g., water boilers, dishwashers). We did not elicit beliefs about energy intensity or carbon emissions in the original experimental sample, to avoid the risk of undermining our shower energy report treatments. We did, however, elicit beliefs about carbon emissions in a different sample of students living in the same dormitories three years after the original study ($n = 329$). For more details on this supplementary study, see Appendix E. Without additional information, these students underestimated the carbon impact of warm water use in the shower by a factor of 8 to 9 on average, even though the average guess for the amount of water used per shower was fairly unbiased. On average, students estimated that a typical shower causes emissions of 91.3 grams of CO₂ (median 35 grams),

²³ We excluded 35 subjects who responded to the baseline survey more than 2 weeks after we distributed shower meters, as they have likely reached the intervention stage by then. We also exclude 3 outliers with estimates above 200 L. The corresponding regression results are presented in Appendix Table A13.

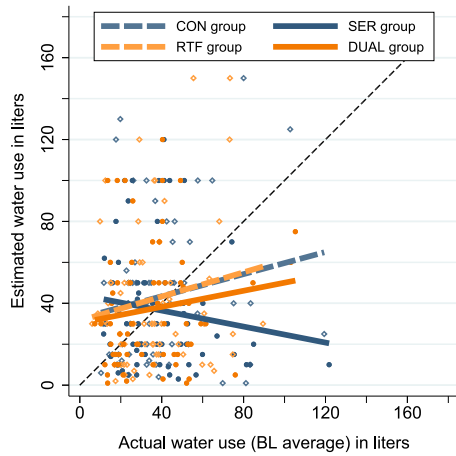


Fig. 4. Pre-intervention awareness about water use per shower.
Notes. This figure compares estimated water use from the baseline survey with actual water use in the baseline stage (showers 2 to 10), excluding late survey responders. 3 outliers with estimates between 200 and 600 L are excluded. Point clouds consist of individual observations (hollow diamonds for CON and RTF, solid circles for SER and DUAL) and lines represent separate regression fits for each treatment group. The dashed line starting from the origin is the 45 degree line.

whereas the actual emissions amount based on the data from our main experiment is about 800 grams.²⁴ Thus, there might be a large potential for encouraging energy conservation through the information provided in shower energy reports (Byrne et al., 2018).

Although anecdotally compelling, finding direct evidence for inattention or self-control problems in the shower is trickier. The closest proxy we have is a baseline survey item on how much subjects agree with the statement “I like to let my mind wander when I shower.”. 59% of our sample agreed or strongly agreed to the statement (25% strongly agree), whereas only 18% of subjects disagreed (13% weakly disagree, 5% strongly disagree). Moreover, subjects’ response to this item is significantly correlated with their average baseline energy use in the shower (Pearson’s $\rho = 0.1645, p = 0.004$). In fact, it is the single most predictive item for baseline consumption in the entire survey based on simple linear regressions. Our interventions could thus help reduce energy use by reminding subjects to stay focused and not lose track of time completely under the shower.

5. Estimation approach

Next, we describe our strategy for estimating the effects of our interventions on resource use in the shower. The empirical results will be presented in the following section.

5.1. Basic identification and estimation strategy

To formally estimate the effects of different intervention regimes, we exploit the randomized assignment of subjects into experimental conditions as well as the staggered introduction of real-time feedback and shower energy reports, which gives us a double-layered difference-in-differences setup. The differential changes in consumption behavior across conditions from baseline stage to intervention stage 1 identify the causal effect of real-time feedback (RTF/DUAL versus CON/SER), and the additional changes from intervention stage 1 to stage 2 identify the causal effect of shower energy reports, both in isolation (SER versus

²⁴ The average guess for amount of water used per shower was 40.4 liters. The survey was conducted in Nov/Dec 2019 among 329 residents of the exact same student dorms in which the original study took place in 2016/17. Only 4 surveyees had already participated in the original study.

CON) and in conjunction with real-time feedback (DUAL versus RTF). In this setup, interaction effects between the two interventions can be identified by comparing the incremental effects of shower energy reports with real-time feedback (DUAL) and without real-time feedback (SER).

If one was only interested in estimating the effect of real-time feedback in isolation, the most straightforward approach would be to simply compare how subjects in the RTF and CON groups change their behavior from the baseline stage to the entire intervention stage, without any need to consider shower energy reports or stage 2. To jointly estimate the effects of each intervention regime, using data from all experimental condition, we instead consider the following regression equation:

$$y_{it} = \alpha_i + IN_{it} \times (\beta_0 + \beta_1 T_i^{R/D} + \beta_2 T_i^S + \beta_3 T_i^D) + IN_{it}^{s2} \times (\gamma_0 + \gamma_1 T_i^{R/D} + \gamma_2 T_i^S + \gamma_3 T_i^D) + \varepsilon_{it}, \quad (8)$$

where the outcome variable y_{it} is energy (water) usage by individual i for shower number t , and α_i is the individual fixed effect. $T_i^{R/D}$, T_i^D and T_i^S are treatment group indicators, where superscript R/D denotes the combined real-time feedback groups RTF and DUAL, D denotes the DUAL group only, and S denotes the SER group only. Finally, IN_{it} is an indicator that takes the value 1 if observation it falls into the intervention stage ($t > 10$), and IN_{it}^{s2} is an indicator for showers that fall into intervention stage 2. As IN_{it} applies to the entire intervention period, IN_{it}^{s2} captures incremental changes in consumption from intervention stage 1 (pre-report) to stage 2 (post-report). Note that the stage 2 coefficients combine the effects of two distinct reports, the second containing also social comparison.

Given our formulation of the statistical model, we can interpret β_1 as treatment effect of real-time feedback on energy (water) use per shower in the first stage of the study, and γ_1 is its change in the second stage. This allows us to test *Prediction 1*. The relevant comparison for *Prediction 2* on the effect of shower energy reports in isolation is between SER and CON after the reports, which is captured by γ_2 . Finally, γ_3 captures the marginal effect of adding shower energy reports to real-time feedback, by comparing DUAL and RTF in intervention stage 2. Finally, the comparison between γ_2 and γ_3 nails down the interaction effect between the two interventions and thus allows us to test for any potential complementarities (*Prediction 3*). As each test is associated with a separate hypothesis, we do not adjust our inference for multiple hypothesis testing (Rubin, 2021).

5.2. Estimating treatment effects on the treated

One complication in estimating the effect of shower energy reports is that 28% of subjects did not succeed in uploading any data to the Amphiro smartphone app before we sent out the reports, mostly due to technical problems (e.g., Bluetooth connection failure).²⁵ For these “non-uploaders”, we were unable to provide informative shower energy reports. As the emails were generated automatically, non-uploaders in SER and DUAL groups received report templates with blanks where it was supposed to show statistics on resource use and environmental impacts. Effectively, this leads to imperfect treatment take-up of shower energy reports, although being less the result of deliberate non-compliance than unfortunate circumstances. For participants in the CON and RTF groups, it is inconsequential whether they successfully uploaded data.

To test Predictions 2 and 3 from Section 3.6, one possible approach to estimate treatment effects under imperfect treatment take-up would be simply to run an intention-to-treat (ITT) analysis, which only uses treatment assignment information and ignores that some subjects did

²⁵ Out of the 90 non-uploaders in our estimation sample, 63 have explicitly contacted us for technical problems encountered during their upload attempts.

not actually receive informative shower energy reports. While this would be the relevant parameter in many policy evaluation contexts, one of our main aims is to test whether the effect of receiving information on the energy use and carbon emissions due to hot water consumption interacts with the effect of receiving real-time feedback in the shower (Prediction 3). This would shed light on the potential importance of multiple biases for complementarities between behavioral interventions that we identify theoretically in Section 2. However, a stringent empirical test of this requires that subjects indeed have the opportunity to gain knowledge through shower energy reports. Thus, the more policy-relevant parameters in our case are the treatment effects on the treated (TOT) in the DUAL and SER groups, i.e., the effects of shower energy reports on subjects who managed to upload data prior to the report and thus received actual information on the environmental impact of their hot water consumption in the shower.

The first way in which we estimate the TOT is by simply comparing only the uploaders in SER and DUAL groups with subjects in the CON and RTF groups. The usual concern at this point would be that treatment take-up is not random. Fortunately, our setting limits potential endogeneity concerns for three reasons. First, we include individual fixed effects, so our estimates would still be unbiased if differences between uploaders and non-uploaders do not interact with the treatment effect. Second, subjects only knew that they should use the smartphone app to upload data, but we did not announce that we would use this data to construct shower energy reports. Thirdly, the main cause for non-compliance is not the lack of willingness to use the smartphone app, but unexpected technical failure, which is unlikely to be selected on the trend. To alleviate the most blatant endogeneity issue, we also exclude non-uploaders in the CON and RTF groups who did not report any technical problems. Appendix Tables A2 and A3 present additional balance checks for the TOT subsample and show that the experimental groups remain balanced along baseline characteristics.

The second way in which we estimate the TOT is by using random treatment assignment as instrument for actual take-up.²⁶ This can be shown to identify the so-called local average treatment effect (LATE), i.e., the average treatment effect for the sub-population of compliers, in our case the uploaders (Imbens and Angrist, 1994).²⁷ Compared to the “uploaders-only”-approach, the instrumental variables approach is consistent under weaker assumptions, but potentially inefficient. We will report the results from both TOT-approaches, but the estimates are very similar, suggesting that non-compliance due to technical issues was likely uncorrelated with conservation intentions in our sample.

6. Main empirical results

6.1. Average treatment effects

We start by presenting descriptive evidence on the resource conservation effects of our interventions in Fig. 5 by plotting subjects' average changes in energy and water consumption per shower in intervention stage 1 (pre-report) and intervention stage 2 (post-report) compared to the baseline period. The differences-in-differences across treatment

²⁶ To do this, we create new treatment indicators for the DUAL and SER groups that took the value 1 for showers in IN stage 2 by subjects who were assigned to the respective group and who uploaded data through the smartphone app that we could use to construct their shower energy reports. The previously defined ITT indicators are then used as instruments for these new indicators for receiving actual shower energy reports.

²⁷ This identification result holds under the condition that there are no “defiers”, subjects who always do the opposite of what they are prescribed. This monotonicity condition holds by design in our study, because we control the eligibility of shower energy report treatment, so any participant in the sample can be classified either as complier or as never taker in the LATE framework.

groups then correspond to the average treatment effects. Note that the stage 2 averages need to be interpreted as combining effects of two distinct reports, one without and one with social comparison. In order to show the treatment effects on the treated (TOT), i.e., the effect of informative shower energy reports, we use the uploaders-only approach of excluding non-compliers in SER and DUAL as well as non-compliers without technical problems in CON and RTF.

Fig. 5 essentially summarizes our main results in eight bars. The patterns are very similar for energy and water consumption. The four bars to the left of the dashed vertical line represent the change in resource use per shower in intervention stage 1 compared to the baseline stage. We can see that relative to subjects in the CON and SER groups, subjects in the RTF and DUAL groups with real-time feedback reduced their consumption drastically, by almost 0.4 kWh of energy and 6 L of water per shower. Recall that there were no shower energy reports yet at this point. The four bars to the right of the dashed vertical line represent the change in energy use per shower from baseline stage to intervention stage 2, after shower energy reports were sent out. The first observation is that average resource use in the control group further increased, which could be driven by weather effects, by pending exams leaving students stressed and in need for a long and warm shower, or by Hawthorne effects that decrease over time (Tiefenbeck, 2016).²⁸ The second observation is that the RTF group and the CON group followed a more or less parallel trend from intervention stage 1 to stage 2, hence the effect of real-time feedback in isolation remains nearly constant. The third observation is that providing shower energy reports in isolation does not seem to result in effective behavioral change: consumption of subjects in the SER group followed the CON group in close synchronization. In light of this, the fourth and final observation is particularly striking: shower energy reports are highly effective when combined with real-time feedback. In fact, subjects in the DUAL group are the only ones to defy the general upward trend and reduce their consumption considerably compared to subjects in the RTF group.

The descriptive evidence presented in Fig. 5 is confirmed by formal empirical estimates based on the empirical strategy outlined in the previous section. Table 3 presents the regression results from estimating Eq. (8), with columns 1–2 using the uploaders-only approach, and columns 3–4 using the alternative LATE approach. To ensure that our statistical inference procedure is robust to arbitrary temporal interdependence of showers taken by the same person, we cluster all standard errors at the individual level. Appendix Figures A4–A5 show that all statistical test results are virtually identical when using randomization-based inference methods (Young, 2019).

Focusing first on the effects of real-time feedback in isolation, LATE is preferable as it utilizes the full sample. We document a conservation effect of around 0.37kWh energy and 5.5 L of water per shower from intervention stage 1 onwards (coefficient β_1). The effect does not change significantly in intervention stage 2 (coefficient γ_1); if anything, it becomes slightly stronger. Another direct way to estimate the effect of real-time feedback that is easier to interpret is to only compare subjects in the RTF and CON groups, since there is no need to take into account effects of shower energy reports. Appendix Table A4 shows that real-time feedback in isolation reduces resource use by 0.4 kWh of energy and 6.3 L of water per shower compared to the CON group over the entire intervention period, which corresponds to about 17–18% of average baseline consumption.

Result 1. *Real-time feedback (in isolation) through the smart meter display led to a reduction in energy (water) consumption by about 0.4 kWh (6.3 L) or 17%–18% per shower.*

²⁸ While the baseline phase fell mainly into an unusually warm and dry December, the main intervention months of January and February saw much higher precipitation. Exam periods at the universities began in mid-February.

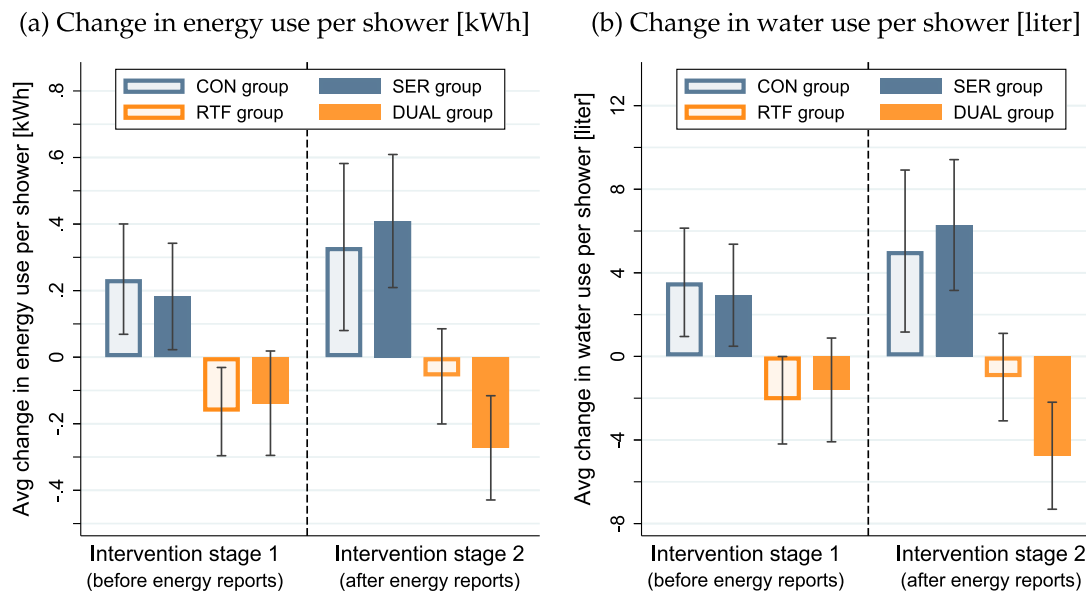


Fig. 5. Descriptive evidence on resource conservation effects.
Notes. The bars represent changes in average energy use and water use per shower compared to the baseline period. The error whiskers represent 90% confidence intervals. Non-uploaders in SER and DUAL as well as non-uploaders without technical problems in CON and RTF are excluded.

Table 3
 Treatment on the treated (TOT) estimates.

	<i>Uploaders-only</i>		<i>LATE</i>	
	(1) Energy [kWh]	(2) Water [liter]	(3) Energy [kWh]	(4) Water [liter]
(β_0) Intervention	0.207* (0.108)	3.067* (1.663)	0.200** (0.099)	2.961* (1.530)
(β_1) Intervention \times RTF/DUAL	-0.388*** (0.130)	-5.762*** (2.071)	-0.368*** (0.122)	-5.514*** (1.932)
(β_2) Intervention \times SER	0.013 (0.151)	0.543 (2.361)	0.003 (0.132)	0.476 (2.051)
(β_3) Intervention \times DUAL	0.031 (0.111)	0.524 (1.808)	0.109 (0.106)	2.140 (1.719)
(γ_0) IN stage 2	0.110 (0.085)	2.191* (1.319)	0.152* (0.090)	2.756** (1.360)
(γ_1) IN stage 2 \times RTF/DUAL	-0.023 (0.109)	-1.234 (1.805)	-0.054 (0.113)	-1.550 (1.806)
(γ_2) IN stage 2 \times SER	0.124 (0.124)	1.336 (1.993)	0.079 (0.149)	0.573 (2.326)
(γ_3) IN stage 2 \times DUAL	-0.230** (0.109)	-3.836** (1.908)	-0.226* (0.122)	-4.013* (2.152)
(α_i) Individual fixed effects	yes	yes	yes	yes
<i>p</i> -value: $\gamma_2 = \beta_1$	0.007	0.018	0.026	0.053
<i>p</i> -value: $\gamma_2 = \gamma_3$	0.033	0.062	0.114	0.149
Clusters	261	261	318	318
Observations	14 712	14 712	17 942	17 942
R^2	0.412	0.415	0.004	0.004

In columns (1) and (2), we exclude all non-uploaders in SER and DUAL as well as all non-uploaders in RTF and CON who did not report a technical problem. In columns (3) and (4), we use treatment assignment to SER and DUAL, respectively, interacted with the IN stage 2 indicator as instrument for receiving informative shower energy reports. The reported R^2 in columns (3) and (4) is the within R^2 . Standard errors in parentheses are clustered at the individual level. Permutation-based inference procedures are presented in Figures A4 and A5. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

For the shower energy reports, we need to focus instead on intervention stage 2 and account for imperfect compliance, to estimate the effects of actually receiving information on resource use and environmental impacts (see Appendix Table A4 for the intention-to-treat estimates). Recall that while the LATE approach is consistent even under strong endogeneity of treatment take-up, the uploaders-only approach is potentially more efficient and still consistent if technical issues in uploading data are as good as random.

Table 3 shows that the point estimates obtained both approaches are very similar, implying that endogeneity of treatment take-up is likely not a major issue in our sample, whereas the standard errors are smaller in the uploaders-only approach. Contrary to prediction 1, shower energy reports in isolation had no significant conservation effect in the SER group (coefficient γ_2), and the point estimates even go in the opposite direction. While the null effect is not very tightly estimated, we can rule out reductions of greater than 4%–5% (0.08 kWh and 1.95

L per shower) with 90% confidence in our preferred uploaders-only specification. This would be consistent with effect sizes in the order of magnitude found in previous studies (e.g., Allcott, 2011). Note that lower statistical power due to non-compliance does not explain the insignificant coefficient for SER compared to the significant coefficient for RTF, since the standard errors for β_1 and γ_2 are similar. We can also reject the hypothesis that shower energy reports in isolation were as effective as real-time feedback in isolation ($p = 0.009$).

Result 2. Shower energy reports in isolation did not induce any significant reduction in energy and water consumption per shower.

In stark contrast, we find that adding shower energy reports in the DUAL group induced subjects to further reduce their consumption by around 0.23 kWh of energy and 3.8 L of water per shower in intervention stage 2 (coefficient γ_3), corresponding to another 10%p reduction from baseline consumption and about 60% of the effect of real-time feedback in isolation. Thus, information on environmental impacts of hot water consumption was not ineffective per se in our setting, but in fact boosted conservation efforts considerably when administered in combination with real-time feedback on water usage via the smart meter displays. This contrast between individuals' responses to shower energy reports with and without real-time feedback is all the more remarkable given that in the former case they had already cut their consumption significantly in intervention stage 1, thus leaving less room for further behavioral adjustments. The exact hypothesis for complementarity further requires to test not only whether the incremental consumption reduction in the DUAL group was different from zero, but also from the effect in the SER group (i.e., whether $\gamma_3 = \gamma_2$). This difference is statistically significant in the uploaders-only specification ($p = 0.033$), although in the less efficient LATE-specification it would be weakly significant only when using a one-sided test.

Result 3. Combining real-time feedback with shower energy reports further reduced energy (water) use by around 0.23 kWh (3.8 L) per shower, in addition to the conservation effect of real-time feedback in isolation.

Overall, we observe a sizable complementarity between the two interventions in our setting. This is consistent with our theoretical framework, which shows that in the presence of multiple biases, behavioral interventions may need to address all significant sources of bias simultaneously in order to unfold their full effect. While shower energy reports provide information about resource use and associated environmental impacts, conservation efforts may be hindered by residual biases such as lack of salience. Real-time feedback through smart meters could thus turn environmental considerations into action by keeping them on top of people's mind in the heat of the moment. We will analyze the underlying mechanisms more closely in Section 7.

6.2. Robustness checks on timing

As the timing of observations with regard to intervention stage 2 involves a degree of fuzziness for subjects who did not use the app frequently, we conduct a number of robustness checks. First, we use a donut hole approach that excludes the around 10% of shower observations with the highest uncertainty about whether they occurred before or after the first shower energy report; Appendix Table A5 shows that our results remain nearly unchanged.²⁹ Second, our results

²⁹ To be more precise, we calculate for each shower a probability that it occurred after reading the first shower energy report (or placebo email). Notice that for many showers, these probabilities are either 0 or 1, because they were uploaded before the report or after the first post-report upload, respectively. For observations within the range of uncertainty, we calculate approximate probabilities assuming that the frequency of showering is constant. We then exclude all observations with probability between 10% and 90%, i.e., those with significant uncertainty about whether they occurred before or after the report. For non-uploaders, we use the same cutoffs of 10% and 90%, but based on the CDF for uploaders (see Figure A3).

are also robust to using an alternative definition of report timing for non-uploaders that deterministically assigns non-uploaders into intervention stage 2 based on the median study completion value among uploaders rather than the full cumulative distribution (see Appendix Table A6). Finally, we estimate a specification in which we assign all subjects into intervention stage 2 based on when we sent out the first shower energy report email, rather than based on when they actually read the reports (proxied by mini-survey response date). Appendix Table A7 shows that our results are robust to using this alternative timing indicator.

6.3. Margins of behavioral adjustment

Subjects could conserve energy by reducing the temperature to which water is heated to and the overall amount of water that needs to be heated up. Appendix Table A8 shows that reductions in water temperature seem to be at most a minor factor in our sample, perhaps for hedonic reasons.³⁰ Hence, water and energy usage tend to be very closely aligned with each other, and energy conservation effects are almost equivalent to water conservation effects, in relative terms. Water conservation, in turn, can be achieved by adjusting time spent under the shower, water flow rate (i.e., liters of water per minute), and the covariance structure. The data suggests that subjects respond to the interventions mostly by taking shorter showers and reducing the flow rate during long showers.

6.4. Treatment effect dynamics

In a next step, we investigate how resource conservation outcomes changed over time, with a focus on the last 5–6 weeks period of our study, after the first energy reports were sent out (IN stage 2). This allows us to test whether effects declined over time or remained stable and whether the second shower energy report (containing social comparison) may have induced behavioral responses. To estimate effect dynamics, we extend the empirical model from Eq. (8) by interactions with a time variable Z_i :

$$y_{it} = \alpha_i + IN_{it} \times (\beta_0 + \beta_1 T_i^{R/D} + \beta_2 T_i^S + \beta_3 T_i^D) + IN_{it}^2 \times (\gamma_0 + \gamma_1 T_i^{R/D} + \gamma_2 T_i^S + \gamma_3 T_i^D) + IN_{it}^2 \times Z_i \times (\delta_0 + \delta_1 T_i^{R/D} + \delta_2 T_i^S + \delta_3 T_i^D) + \epsilon_{it}. \quad (9)$$

We explore two variants of Z_i . In the first variant, we look additionally at energy use per shower after the second shower energy report, which was sent about two weeks after the first report. In the second variant, we interact each treatment group indicator with a linear time trend, so the δ coefficients can be interpreted as weekly depreciation (or appreciation) rate of energy conservation effects by intervention regime.

Table 4 suggests that the effect of shower energy reports in the DUAL group seemed to gradually unfold over time. The point estimates in columns (1) and (2) indicate that the average conservation effect is driven largely by the final 3–4 weeks of the study, after the second reports were sent out. However, this does not seem stem from a discrete jump, but rather from a continuous trend. In columns (3) and (4), we estimate that the conservation effect per shower in the DUAL group increases by a rate of around 0.08–0.09 kWh every week. These descriptive results need to be interpreted with caution, as the relevant coefficients are statistically insignificant. However, we note that the point estimates for shower energy reports in isolation (SER group) show no signs of any quantitatively significant dynamic pattern. The

³⁰ At 40 L and a base temperature of 37°C, reducing energy conservation by 0.1 kWh would require lowering the temperature by more than 1°C, ceteris paribus.

Table 4
Treatment effect dynamics.

	$Z_i = \mathbb{1}\{\text{after report 2}\}$		$Z_i = \# \text{ weeks after report 1}$	
	(1) Uploaders	(2) LATE	(3) Uploaders	(4) LATE
...
(γ_0) IN stage 2	0.091 (0.095)	0.135 (0.103)	0.016 (0.108)	0.066 (0.115)
(γ_1) IN stage 2 \times RTF/DUAL	-0.038 (0.120)	-0.059 (0.127)	0.037 (0.149)	0.016 (0.154)
(γ_2) IN stage 2 \times SER	0.149 (0.135)	0.102 (0.156)	0.210 (0.151)	0.163 (0.178)
(γ_3) IN stage 2 \times DUAL	-0.086 (0.109)	-0.068 (0.120)	0.009 (0.154)	0.049 (0.169)
(δ_0) IN stage 2 $\times Z_{,i}$	0.032 (0.090)	0.029 (0.087)	0.032 (0.023)	0.029 (0.022)
(δ_1) IN stage 2 \times RTF/DUAL $\times Z_{,i}$	0.025 (0.128)	0.009 (0.125)	-0.021 (0.035)	-0.024 (0.034)
(δ_2) IN stage 2 \times SER $\times Z_{,i}$	-0.043 (0.123)	-0.040 (0.133)	-0.030 (0.036)	-0.028 (0.041)
(δ_3) IN stage 2 \times DUAL $\times Z_{,i}$	-0.245 (0.201)	-0.273 (0.207)	-0.082 (0.055)	-0.095 (0.058)
(β_j) Intervention indicators	yes	yes	yes	yes
(α_i) Individual fixed effects	yes	yes	yes	yes
Clusters	261	318	261	318
Observations	14 712	17 942	14 712	17 942
R^2	0.413	0.004	0.413	0.005

The results are obtained by estimating equation (9). The full table with all the coefficients is presented in Table A9. In columns (1) and (3), we exclude all non-uploaders in SER and DUAL, as well as all non-uploaders in RTF and CON who did not report a technical problem. In columns (2) and (4), we use treatment assignment to SER and DUAL, respectively, interacted with the IN stage 2 indicator as instrument for receiving informative shower energy reports. The reported R^2 in Columns (2) and (4) is the within R^2 . Standard errors in parentheses are clustered at the individual level. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$.

effect of real-time feedback in isolation also appears to stay constant in intervention stage 2, overall showing no signs of weakening within the 3 months of our experiment.³¹

There are several potential behavioral explanations for this descriptive pattern of increasing responses over time that we observe in the DUAL group.³² The first explanation is that the peer comparison element in the second report provided important additional social motivation to conserve energy, which then interacted with real-time feedback; this would be in line with our theoretical framework as well as previous literature. A second explanation is that subjects may have required some time to discover new strategies for further reducing resource use. This experimentation channel is consistent with Appendix Table A8, which suggests that subjects in the DUAL group conserved resources in the second intervention stage by reducing the flow rate overproportionately during longer showers. To further investigate this channel, Appendix Table A10 reports a specification that allows for discontinuities as well as differential trends after each report; although admittedly noisy, the estimates suggest a more or less constant and continuous (downward) slope of resource consumption in the DUAL condition that already begins starting from the first report. This suggests that the effects are unlikely to be driven by social comparison alone, although it may have played a role in reinforcing them. Importantly, the results speak against pure Hawthorne effects or short-lived attention boosts, as these would rather predict an

³¹ This is consistent with other interventions using smart shower meters (Agarwal et al., 2020; Byrne et al., 2022). Energy conservation studies in other settings show some degree of backsliding over time after being exposed to non-monetary interventions (e.g. Allcott and Rogers, 2014; Ito et al., 2018), although investments into physical capital may alleviate this issue in the long term (Brandon et al., 2017).

³² Another potential explanation is that the apparent increase in estimated effects over time is a statistical mirage driven by decreasing measurement error about when subjects were treated with shower energy reports.

“action-and-backsliding” pattern (Allcott and Rogers, 2014; Schwartz and Loewenstein, 2017).

6.5. Heterogeneity by baseline consumption

A frequent finding in the literature is that households or individuals with high baseline consumption tend to respond more strongly to interventions that foster conservation behavior (e.g., Allcott 2011, Ferraro and Price 2013, Tiefenbeck et al. 2018). Policy makers could therefore improve cost-effectiveness by targeting high-baseline consumers. To estimate heterogeneity by baseline energy use, we extend the statistical model in Eq. (8) by adding interactions with baseline consumption, measured using subject’s average energy usage per shower in the baseline stage. Alternatively, we estimate a specification where we interact with an above-median indicator. Appendix Table A11 shows that, consistent with previous studies, the effect of real-time feedback increases with baseline usage. In intervention stage 2, subjects with 1 kWh higher baseline in the RTF group reduced their energy use by an additional 0.25 kWh ($p = 0.083$) on average, and above-median baseline users (mean 3.30 kWh) saved 0.63 kWh ($p = 0.043$) of energy more per shower compared to subjects with below-median baseline use (mean 1.17 kWh). It also appears that providing information through shower energy reports in the DUAL condition was about twice as effective for above-median users compared to below-median baseline users, although the difference is not statistically significant, whereas shower energy reports in isolation (SER) had no significant effects in either subpopulation.

7. Underlying mechanisms

The empirical results show that, in our setting, shower energy reports appeared to be ineffective in isolation, but induced large and significant conservation effects when combined with real-time feedback, which suggests that our interventions were complements. Through the

lens of the theoretical framework in Section 2, our proposed explanation for this finding is that the two interventions targeted separated behavioral biases. Shower energy reports may have increased knowledge about environmental impacts of warm water use in the shower, but this in itself may not achieve reductions in energy consumption if subjects still suffer from limited attention or self-control problems. Real-time feedback could help mitigating these problems and thus enable knowledge gains to translate into conservation behavior. If, on the other hand, shower energy reports and real-time feedback both operated through the same mechanisms, we would generally not expect complementarities unless through some type of crowding in effect, e.g., if the combined intervention leads to positive attention or motivation spillovers. In this section, we conduct a number of analyses to explore the mechanisms underlying our results.

7.1. Awareness about resource intensity and environmental impacts

A crucial element of both interventions in our study is that they can enable learning about the outcomes of one's behavior. Real-time feedback through the smart meter provides immediate display of water use (and temperature) for the current shower. Shower energy reports also contain information of individuals' entire history of water (and energy) use per shower since the start of the study, with the difference that it comes in retrospect. Hence, a first manipulation check for our interventions is to analyze their effect on subjects' awareness about their own water use per shower.

In the post-intervention survey at the end of the study, we asked subjects to again estimate the amount of water they typically use per shower. Recall that prior to the interventions, subjects' assessments were virtually uncorrelated with their actual water use (see Fig. 4). This picture changes after the interventions. Fig. 6 plots individuals' post-intervention estimates as a function of their average water use per shower as measured by the smart meter. The corresponding regression Table A13 is presented in Appendix A. Whereas subjects in the CON group remained as ignorant as before, the estimates by subjects who received real-time feedback (RTF and DUAL group) were much more aligned with their actual consumption patterns, as indicated by the fitted regression lines moving closer to the identity line. Importantly, shower energy reports in isolation (SER group) also induced significant learning effects about water use compared to the control group ($p = 0.039$). Moreover, we cannot reject that the learning slopes among uploaders are different between SER and DUAL ($p = 0.522$). We obtain similar results when we focus instead on the magnitude of estimation errors as outcome variable. Table A14 in Appendix A shows that subjects estimation errors in the three treated groups are on average about 20–30 percentage points closer to their actual water use than subjects in the CON group, and notably, the effect is virtually the same for SER, RTF, and DUAL groups.

Taken together, the results show that subjects in our study became better informed about their own consumption behavior in the shower through our feedback interventions. However, belief updating about water usage alone cannot explain our main results. First, subjects' prior beliefs about water use were by and large unbiased even in the control group: on average, low-baseline users overestimated and high-baseline users underestimated. Second, we observe significant belief updating in the SER group that does not translate into resource conservation effects. This points to the importance of the immediacy and salience of the real-time feedback intervention, which can help subjects track their water use while showering and overcome inattention problems.

In contrast to real-time feedback, shower energy reports did not only contain information about water, but also on energy usage due to water heating and environmental impacts in terms of CO₂ emissions. This can explain why subjects in the DUAL group reduced their energy consumption even further after receiving the reports. As a manipulation check for whether subjects responded to this information, we conducted a supplementary survey in a new sample of 329 students at the end

of 2019 (see also Section 4.6). After eliciting prior beliefs about water consumption and CO₂ emissions per shower, we randomly presented one fact sheet (out of three) to each surveyee, mimicking the basic informational content of our original interventions. The "CON sheet" only reported the average water temperature in the shower, the "RTF sheet" also included the average amount of water used, and the "SER sheet" further added information on energy use and CO₂ emissions. After presenting the fact sheets, we elicited posterior beliefs as well as conservation intentions. We find that, relative to RTF sheet, surveyees that received the SER sheet information drastically adjusted their beliefs about CO₂ emissions upwards ($p < 0.001$), and their self-reported intention to take shorter showers in the future increased by with a 0.24 standard deviations ($p = 0.003$). For further details on the supplementary survey, see Appendix E.

Thus, the shower energy reports increased knowledge about the environmental impact of hot water usage in the shower as well as conservation intentions, yet they were only associated with significant conservation effects when combined with real-time feedback. A key insights of our theoretical framework is that if biased behavior arises from multiple different sources, a narrowly-targeted intervention can be undermined by residual biases (Anna Karenina effect). Hence, a likely explanation of our results is that, when shower energy reports were implemented in isolation, additional biases due to, e.g., limited attention or self-control problems have prevented knowledge gains and good intentions from translating into actual behavior.

7.2. Engagement with shower energy reports

One potential alternative channel is differential treatment engagement, in the sense that subjects across experimental conditions may have paid more or less attention to the interventions per se. For example, if previous exposure to real-time feedback induced subjects in the DUAL group to read shower energy reports more carefully than subjects in the SER group, this might lead to complementarity between the two interventions through some type of crowding in or foot-in-the-door effect.³³

While the previous analyses show that shower energy reports induced significant learning effects also in the absence of real-time feedback, we can also compare engagement with the reports between SER and the DUAL groups more directly by making use of the mini-surveys that were attached to the reports. As described before, each email included a link to a survey in which we asked subjects to estimate their water usage per shower. The survey link was at the bottom of the email, so subjects had to scroll through all the statistics on resource use and CO₂ emissions before clicking on it. We therefore use survey responses as proxy for the level of engagement with the feedback email. Appendix Table A15 shows response rates by treatment group in the uploaders-only sample. The overall response rate among uploaders was 87% for the first email and 71% for the second email. While the share of respondents in the SER group was 8.4% p lower than in the DUAL group for the first email ($p = 0.203$) and 9.4% p higher for the second mail ($p = 0.308$), both differences are statistically insignificant. Furthermore, we find no evidence that uploaders in the DUAL group studied the reports more carefully than uploaders in SER group.

Table A15 also compares estimation error across treatment group, defined as percent deviation of the water use estimate in the mini-survey from the exact number that we showed in the same personalized

³³ An opposite effect is also conceivable, in which paying attention to one intervention decreases engagement with the other, for example due to cognitive capacity constraints (see, e.g., Altmann et al., 2022; Trachtmann, 2022) or lower perceived marginal benefits of information when subjects already receive real-time feedback. This would work against our complementarity argument.

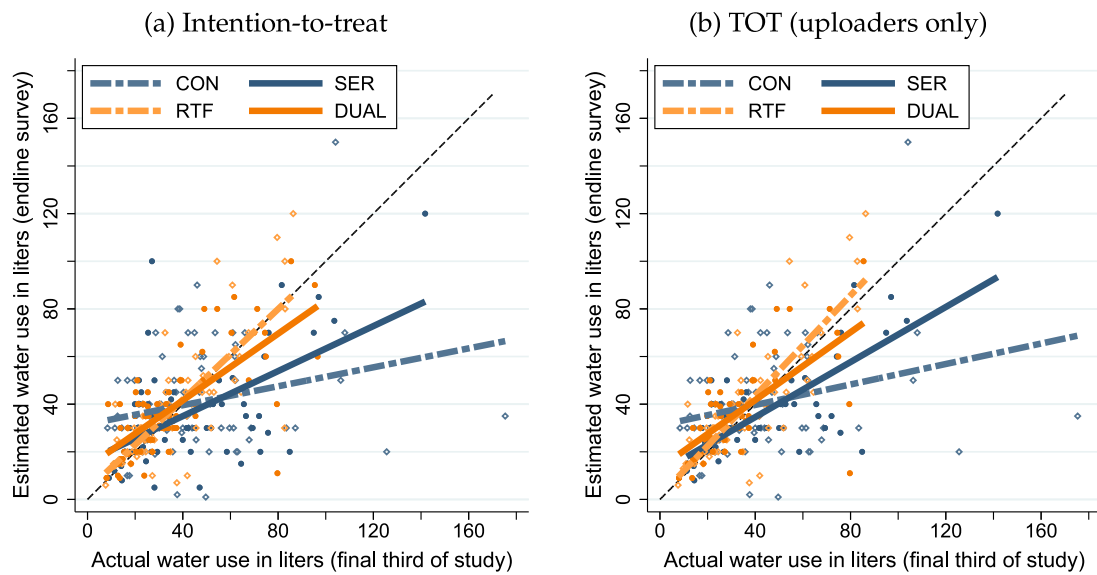


Fig. 6. Post-intervention awareness about water use per shower. *Notes.* Both graphs compare subject’s estimates in the final questionnaire with their measured average water use for the last third of shower observations. Graph (b) only uses the subsample defined for the uploaders-only approach. 7 outliers with estimates between 200 and 500 L are excluded. Point clouds consist of individual observations (hollow diamonds for CON and RTF, solid circles for SER and DUAL) and lines represent separate regression fits for each treatment group. The dashed line starting at the origin is the 45 degree line.

report that contained the survey link.³⁴ Smaller estimation errors are thus a direct indication of subjects paying closer attention while reading. Unsurprisingly, all treated groups outperformed the control group, but we observe no significant difference between uploaders in the SER and the DUAL group.

As a final plausibility check that differences between SER and DUAL are not driven by differential level of engagement with the shower energy reports, we look at whether subjects who studied the reports more closely also engaged more strongly in conservation actions. To do so, we create new treatment compliance indicators that also consider subject’s level of engagement with the reports, to varying degrees of strictness. Specifically, we define an indicator for whether subjects uploaded data *and* clicked on the mini-survey in their report, and additional indicators for whether a subjects’ estimation precision in the mini-surveys (as defined above) was above the 25th, 50th, or 75th percentile, respectively, of their treatment group. To avoid the endogeneity issue at hand, we instrument each of these treatment compliance indicators with the randomized assignment. Fig. 7 plots results for the effect of shower energy reports in SER and DUAL group, respectively, when using these new set of indicators. The estimated conservation effect in the DUAL group increases monotonically with the strictness of our compliance definition, reaching more than 0.5 kWh per shower for the strictest 75th percentile indicator. In contrast, even the most studious subjects in the SER group did not reduce their energy use on average. Overall, it is therefore unlikely that our empirical results can be explained by differential level of engagement with the shower energy reports.

7.3. Other potential mechanisms

There are a number of alternative channels through which our interventions could affect conservation behavior. Hawthorne effects are one possibility, but recall that also subjects in the control group

³⁴ For the CON and RTF group, we calculate this using the number that we would have shown, were the subjects assigned to one of the shower energy report groups instead. Similarly, for non-uploaders, we use the ex post analog of the same statistic, based on data was uploaded after the report or manually read out by our research team.

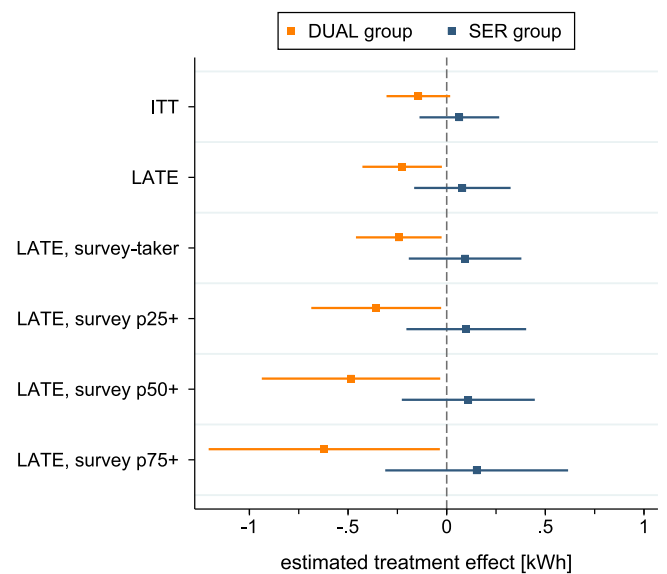


Fig. 7. Effects for different levels of engagement with shower energy reports. *Notes.* The squares represent estimated regression coefficients for the effects of shower energy reports in intervention stage 2, where treatment engagement status is instrumented with treatment assignment (with the exception of ITT). Lines represent 90% confidence intervals. LATE, survey all includes all subjects who uploaded data and clicked on at least one mini-survey. The labels p25+/p25+/p75+ denote the groups of subjects whose estimate precision, defined as distance between estimated and measured water use per shower calculated for the shower energy reports, was above the 25th, 50th, or 75th percentile of all subjects, respectively. Responses from the two mini-surveys are combined by using the minimum estimate precision to define indicators.

received a smart meter and emails reminding them to upload their data. Another potential channel may be that we accidentally killed the joy of showering. Reassuringly, our endline survey results suggest that the interventions had no effects on self-reported shower comfort, thus also alleviating concerns about unintended negative welfare effects (e.g., Damgaard and Gravert, 2018; Allcott and Kessler, 2019). Furthermore, we find no evidence for positive effects on general pro-environmental

attitudes. If anything, we observe a decrease in self-perceived environmentalism in the treated groups compared to the control group, potentially due to feedback provision curbing the capacity for distorted self-image formation. See Appendix F and Table A17 for all results and a more detailed discussion of alternative mechanisms.

8. Discussion

In this paper, we argued that when multiple biases arising from different sources (e.g., imperfect information, limited attention) simultaneously prevent individuals from acting on their values and intentions, then combining interventions that each target a different source of bias can result in complementarity, meaning that each intervention becomes more effective when implemented in conjunction with the other(s) than in isolation. We first introduced a novel theoretical framework that explores the implications of multiple biases and defines conditions for complementary in behavioral interventions in such a setting. We then presented results from a three-month field experiment on resource conservation behavior in an energy- and water-intensive everyday activity (showering), in which we tested the effects of two types of interventions: shower energy reports, which provided information on resource use and carbon emissions via email, and real-time feedback, which made resource use in the shower immediately salient through a smart meter display. While only the latter induced a significant conservation effect when implemented in isolation, combining both interventions resulted in a striking complementarity that is in line with theoretical predictions. Specifically, it seems that knowledge gains about environmental impacts of water heating only translated into behavioral change when resource use was additionally made salient through real-time feedback.

8.1. Relevance of effect sizes

Although our interventions were targeted towards one specific resource-intensive activity, showering, the effect sizes are also quantitatively meaningful on the aggregate household level, which is all the more remarkable given that our subjects had no monetary incentives to conserve resources. In our study, real-time feedback in isolation lowered consumption by 0.4 kWh (6.3 L) per shower; adding shower energy reports further lowered consumption by 0.23 kWh (3.8 L). For comparison, total daily energy use for lighting in German households was less than 0.35 kWh per person on average at that time³⁵. In his influential evaluation of the Opower home energy reports, which target aggregate electricity use in U.S. households, Allcott (2011) finds an average household-level conservation effect of 0.62 kWh per day.

For a simple cost-benefit calculation of the Amphiro smart shower meters, we refer to Tiefenbeck et al. (2018). For the shower energy reports, while it would not be credible to state a generally applicable estimate for the cost of intervention, we note that given we had already set up the technical infrastructure for real-time feedback in our study, the marginal costs of adding information on environmental impacts through emails were close to zero. Although we find no evidence that the email interventions were (cost-)effective in isolation, they produced significant additional conservation effect when combined with real-time feedback, so the bundled intervention should be the most cost-effective regime. This possibility was proposed theoretically in Section 2.3.

³⁵ Source: German Federal Statistical Office (<https://www.destatis.de/EN/Themes/Society-Environment/Environment/Environmental-Economic-Accounting/private-households/Tables/energy-consumption-households.html>)

8.2. Data limitations

Our data has a number of limitations. One issue is that the Amphiro devices had no global time counter, so our only source of information on timing comes from data uploads through a smartphone app. As most subjects tended to upload data in batches, some uncertainty remains about when exactly a shower took place, implying that we cannot easily control for date or time of day fixed effects and that there is some fuzziness around which observations took place before or after a shower energy report — to address this, we conducted a number of robustness checks. Moreover, a subset of participants could not upload any data due to technical issues with the Bluetooth connection, and thus could not receive any informative report. In principle, such problems that happen in early stages in the life cycle of new applications can be ironed out in the future.

Another limitation is that we cannot measure behavior outside of the shower. Hence, we cannot directly rule out, for example, whether subjects substitute part of their hygiene behavior to other facilities such as gym showers or wash basins. However, we find no evidence of extensive margin effects (i.e., on the number of observed showers) across experimental conditions, and a related study on water conservation in dorm showers finds no evidence of students moving between different communal shower facilities within the same building (Goette et al., 2021). Relatedly, we cannot account for potential spillover effects on other resource consumption activities, e.g., kitchen water usage or room heating. It is unclear whether this leads to an overstatement or an understatement of the overall impact of our interventions, and the general evidence for spillover effects of pro-environmental interventions is mixed (e.g., Tiefenbeck et al., 2013; Jessoe et al., 2021; Goetz et al., 2021; Sherif, 2021). Exploring the direction and magnitude of spillover effects thus constitutes an important avenue for future behavioral research.

8.3. Generalizability

Our study on hot water conservation in student dorm showers constitutes a very specific setting. First, students may generally not be representative of the general population. As our subjects also self-selected into participating in the study, they may on average be more intrinsically motivated to protect the environment, although note that about two-thirds of all dorm residents that we could reach through door-to-door recruitment participated in our study. Another noteworthy feature is that students in our sample did not have to pay utility bills and thus had no monetary incentives to conserve water and energy, which is unusual but not unheard of in other settings (e.g., office energy use).³⁶ The theoretical predictions for how these factors would affect the potential for complementarity are ambiguous. On the one hand, stronger marginal (monetary or non-monetary) incentives gives more leverage to alleviating informational and behavioral barriers; on the other hand, ceiling effects may limit the conservation potential if individuals already put in more effort in baseline. Our study was also conducted during winter in Germany, where demand for long and hot showers may have been higher due to cold weather. This could have limited conservation effects due to lower willingness to reduce warm water use, but also increased conservation potential due to a higher baseline. Another characteristic of our sample is that all subjects lived in single-person flats in relatively large and anonymous dorm buildings. This has advantages for the empirical study design, but limits the extent of information sharing and social influence that could be

³⁶ Ito et al. (2018) find that monetary incentives lead to stronger and more persistent reductions in peak electricity use than to moral appeals. Also note that some surveys in Germany find that young people were less likely to behave sustainably in their daily lives than older generations (e.g., Ipsos, 2019).

relevant in multi-person households. Interestingly, Tiefenbeck et al. (2018) found no difference in effects of real-time feedback between one- and two-person households.

Finally, whether and to which extent similar results would arise in other behavioral contexts is an open question. Our theoretical framework predicts that complementarities should become more likely if – following the Anna Karenina principle – multiple different mechanisms play a role in preventing behavioral change, including information frictions and behavioral biases like limited attention and self-control problems, but also more standard economic barriers such as lack of incentives or constraints on time, money, or technology. We suspect that such complexity of behavioral mechanisms is a pervasive feature of many social and economic domains of our lives, e.g., decisions affecting environmental, financial, or health outcomes. If true, the Anna Karenina effect we highlight in this study could imply the existence of numerous untapped opportunities for more targeted intervention designs and help organize empirical research on interaction effects of behavioral policy.³⁷ Obviously, our study cannot offer any definitive conclusion and should be viewed more as a proof of concept. More research is needed to understand how our findings would extrapolate to other settings and samples.

8.4. Implications for policy and research

Policy evaluation typically requires to test whether an intervention in isolation leads to the desired outcomes, to avoid that other interventions confound the effect. Similarly, behavioral researchers who wish to investigate specific determinants of behavior need to manipulate these determinants of interest while holding all other factors constant. However, our study highlights a particular generalizability challenge. Lack of observable impacts in response to an intervention (or manipulation) in isolation may be insufficient to rule out that it is not relevant or effective even within the same sample. For example, one may have concluded from the lack of effectiveness of our shower energy reports in isolation that improving knowledge about energy- and carbon-intensity of water heating does not matter in our context, but our findings suggest that knowledge gains may have been prevented from inducing observable behavioral change by residual biases like inattention or present bias.

The reason for this is that any singular evaluation of an intervention is inevitably confined to the particular choice environment it is introduced into, which is shaped by existing policies, institutions, norms, and individual circumstances. This environment itself can be malleable, so interventions that seem feeble at first glance may be able to unfold their full potential only once combined with complementary policies. We focused specifically on the role of multiple behavioral biases in decision-making creating potential for complementarities, because mitigating one specific bias (e.g., due to knowledge gaps) becomes more effective when also mitigating residual biases (e.g., due to inattention, self-control problems). This implies that policy designers should not focus only narrowly on one specific behavioral mechanism, but also attempt to identify other behavioral factors and how they might interact or interfere. For example, our results suggest that giving people tools that allow them to track their resource use may also make behavior more sensitive to other policies such as informational and norm-based interventions; the same may also apply to conventional

³⁷ For example, Dupas and Robinson (2013) study financial behavior in a development context and find that simply providing a safe box for storing money is effective for encouraging higher savings, except for the subgroup of individuals with severe present bias, who need additional social commitment. Similarly, prompting deliberation about food choice to help resist temptations increases the effectiveness of healthy purchasing subsidies (Brownback et al., 2019). Cortes et al. (2023) find that text-message interventions on parenting practices work less well when in time periods where parents face high cognitive load.

policies like price incentives (Jesoe and Rapson, 2014). An interesting approach for future research could be to first identify and elicit the existence and strength of different behavioral motives and biases at the individual level, and then implement and test tailored combinations of interventions in a second step — akin to personalized medical diagnoses and prescriptions.

The potential for complementarities creates a trade-off for policy makers with a binding budget constraint. They could either target more people with a single intervention or fewer people with a bundled intervention. We show in our theoretical framework that when complementarities between interventions are sufficiently strong, it can be preferable to implement a bundled approach at the cost of covering fewer households. As social scientists are beginning to pioneer the process from small-scale proof-of-concept studies to large-scale interventions (Banerjee et al., 2017), future research should therefore synchronously advance our knowledge on the interplay of different policy instruments.

Declaration of competing interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: Thorsten Staake reports a relationship with Amphiro AG that includes: board membership and equity or stocks. Samuel Schoeb reports a relationship with Amphiro AG that includes: employment and equity or stocks. Verena Tiefenbeck: served as an unpaid scientific advisor to Amphiro AG and is married to Mr. Thorsten Staake.

Data availability

Data will be made available on request.

Appendix A. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.jpubeco.2023.105028>.

References

- Abrahamse, W., Steg, L., Vlek, C., Rothengatter, T., 2005. A review of intervention studies aimed at household energy conservation. *J. Environ. Psychol.* 25 (3), 273–291.
- Agarwal, S., Fang, X., Goette, L., Schoeb, S., Staake, T., Tiefenbeck, V., Wang, D., 2020. The Role of Goals in Motivating Behavior: Evidence from a Large-Scale Field Experiment on Resource Conservation. mimeo.
- Allcott, H., 2011. Social norms and energy conservation. *J. Public Econ.* 95 (9–10), 1082–1095.
- Allcott, H., 2016. Paternalism and energy efficiency: An overview. *Annu. Rev. Econ.* 8 (1), 145–176.
- Allcott, H., Kessler, J.B., 2019. The welfare effects of nudges: A case study of energy use social comparisons. *Am. Econ. J.: Appl. Econ.* 11 (1), 236–276.
- Allcott, H., Rogers, T., 2014. The short-run and long-run effects of behavioral interventions: Experimental evidence from energy conservation. *Amer. Econ. Rev.* 104 (10), 3003–3037.
- Altmann, S., Grunewald, A., Radbruch, J., 2022. Interventions and cognitive spillovers. *Rev. Econom. Stud.* 89 (5), 2293–2328.
- Andor, M.A., Fels, K.M., 2018. Behavioral economics and energy conservation – a systematic review of non-price interventions and their causal effects. *Ecol. Econom.* 148, 178–210.
- Andor, M., Gerster, A., Peters, J., Schmidt, C.M., 2020. Social norms and energy conservation beyond the US. *J. Environ. Econ. Manag.* 103, 102351.
- Ashraf, N., Jack, B.K., Kamenica, E., 2013. Information and subsidies: Complements or substitutes? *J. Econ. Behav. Organ.* 88, 133–139.
- Attari, S.Z., 2014. Perceptions of water use. *Proc. Natl. Acad. Sci. USA* 111 (14), 5129–5134.
- Attari, S.Z., DeKay, M.L., Davidson, C.I., Bruine de Bruin, W., 2010. Public perceptions of energy consumption and savings. *Proc. Natl. Acad. Sci. USA* 107 (37), 16054–16059.
- Banerjee, A., Banerji, R., Berry, J., Duflo, E., Kannan, H., Mukerji, S., Shotland, M., Walton, M., 2017. From proof of concept to scalable policies: Challenges and solutions, with an application. *J. Econ. Perspect.* 31 (4), 73–102.

- Banerjee, A., Chandrasekhar, A., Dalpath, S., Duflo, E., Floretta, J., Jackson, M., Kannan, H., Loza, F., Sankar, A., Schrimpf, A., Shrestha, M., 2021. Selecting the Most Effective Nudge: Evidence from a Large-Scale Experiment on Immunization. Working Paper.
- Bernheim, B.D., Taubinsky, D., 2018. Behavioral public economics. *Handb. Behav. Econ.: Appl. Found.* 1 1, 381–516.
- Bordalo, P., Gennaioli, N., Shleifer, A., 2022. Salience. *Annu. Rev. Econ.* 14, 521–544.
- Brandon, A., Ferraro, P., List, J.A., Metcalfe, R., Price, M., Rundhammer, F., 2017. Do The Effects of Social Nudges Persist? Theory and Evidence from 38 Natural Field Experiments. Vol. w23277. NBER Working Paper.
- Brandon, A., List, J.A., Metcalfe, R.D., Price, M.K., Rundhammer, F., 2019. Testing for crowd out in social nudges: Evidence from a natural field experiment in the market for electricity. *Proc. Natl. Acad. Sci. USA* 116 (12), 5293–5298.
- Brownback, A., Imas, A., Kuhn, M.A., 2019. Behavioral food subsidies. *Rev. Econ. Stat.* 1–47.
- Byrne, D.P., Goette, L., Martin, L.A., Jones, A., Miles, A., Schob, S., Staake, T., Tiefenbeck, V., 2022. The habit-forming effects of feedback: Evidence from a large-scale field experiment.
- Byrne, D.P., La Nauze, A., Martin, L.A., 2018. Tell me something I don't already know: Informedness and the impact of information programs. *Rev. Econ. Stat.* 100 (3), 510–527.
- Camilleri, A.R., Larrick, R.P., Hossain, S., Patino-Echeverri, D., 2019. Consumers underestimate the emissions associated with food but are aided by labels. *Nature Clim. Change* 9 (1), 53–58.
- Carlsson, F., Gravert, C.A., Kurz, V., Johansson-Stenman, O., 2021. The use of green nudges as an environmental policy instrument. *Rev. Environ. Econ. Policy* 15 (2), 216–237.
- Coe, D.T., Snower, D.J., 1997. Policy complementarities: The case for fundamental labor market reform. *IMF Staff Pap.* 44 (1).
- Cortes, K.E., Fricke, H., Loeb, S., Song, D.S., York, B.N., 2023. When behavioral barriers are too high or low—how timing matters for text-based parenting interventions. *Econ. Educ. Rev.* 92, 102352.
- d'Adda, G., Gao, Y., Tavoni, M., 2020. Making Energy Costs Salient Can Lead to Low-Efficiency Purchases. E2e Working Paper 045.
- Damgaard, M.T., Gravert, C., 2018. The hidden costs of nudging: Experimental evidence from reminders in fundraising. *J. Public Econ.* 157, 15–26.
- Delmas, M.A., Fischlein, M., Asensio, O.I., 2013. Information strategies and energy conservation behavior: A meta-analysis of experimental studies from 1975 to 2012. *Energy Policy* 61, 729–739.
- Dena, 2016. dena-Gebäudereport—Statistiken und Analysen zur Energieeffizienz im Gebäudebestand. dena Berlin.
- Duflo, E., Dupas, P., Kremer, M., 2015. Education, HIV, and early fertility: Experimental evidence from Kenya. *Amer. Econ. Rev.* 105 (9), 2757–2797.
- Dupas, P., Huillery, E., Seban, J., 2018. Risk information, risk salience, and adolescent sexual behavior: Experimental evidence from Cameroon. *J. Econ. Behav. Organ.* 145, 151–175.
- Dupas, P., Robinson, J., 2013. Why don't the poor save more? Evidence from health savings experiments. *Amer. Econ. Rev.* 103 (4), 1138–1171.
- European Environment Agency, 2021. Water Resources Across Europe—Confronting Water Stress: An Updated Assessment. Publications Office Luxembourg.
- Fanghella, V., Ploner, M., Tavoni, M., 2021. Energy saving in a simulated environment: An online experiment of the interplay between nudges and financial incentives. *J. Behav. Exper. Econ.* 93, 101709.
- Ferraro, P.J., Price, M.K., 2013. Using nonpecuniary strategies to influence behavior: Evidence from a large-scale field experiment. *Rev. Econ. Stat.* 95 (1), 64–73.
- Fischer, C., 2008. Feedback on household electricity consumption: A tool for saving energy? *Energy Effic.* 1 (1), 79–104.
- Frederiks, E.R., Stenner, K., Hobman, E.V., 2015. Household energy use: Applying behavioural economics to understand consumer decision-making and behaviour. *Renew. Sustain. Energy Rev.* 41, 1385–1394.
- Gabaix, X., 2017. Behavioral Inattention. NBER Working Papers 24096, <https://www.sciencedirect.com/science/article/pii/S2352239918300216>.
- Gardner, G.T., Stern, P.C., 2008. The short list: The most effective actions U.S. households can take to curb climate change. *Environ. Behav.* 50 (5), 12–24.
- Gerster, A., Andor, M., Goette, L., 2020. Disaggregate Consumption Feedback and Energy Conservation. CEPR Discussion Paper 14952.
- Giaccherini, M., Herberich, D.H., Jimenez-Gomez, D., List, J.A., Ponti, G., Price, M.K., 2020. Are Economics and Psychology Complements in Household Technology Diffusion? Evidence from a Natural Field Experiment. Working Paper.
- Goette, L., Han, H.-J., Lim, Z.H., et al., 2021. The Dynamics of Goal Setting: Evidence from a Field Experiment on Resource Conservation. CRC TR 224 Discussion Paper, University of Bonn and University of Mannheim, Germany.
- Goetz, A., Mayr, H., Schubert, R., 2021. Beware of Side Effects? Spillover Evidence from a Hot Water Intervention. Working Paper.
- Hahn, R., Metcalfe, R.D., Novgorodsky, D., Price, M.K., 2016. The Behavioralist as Policy Designer: The Need to Test Multiple Treatment to Meet Multiple Targets. NBER Working Paper 22886.
- Hanna, R., Mullainathan, S., Schwartzstein, J., 2014. Learning through noticing: Theory and evidence from a field experiment. *Q. J. Econ.* 129 (3), 1311–1353.
- Holladay, J.S., LaRivière, J., Novgorodsky, D., Price, M., 2019. Prices versus nudges: What matters for search versus purchase of energy investments? *J. Public Econ.* 172, 151–173.
- Imai, T., Pace, D.D., Schwardmann, P., van der Weele, J.J., 2022. Correcting Consumer Misperceptions About CO2 Emissions. CESifo Working Paper.
- Imbens, G.W., Angrist, J.D., 1994. Identification and estimation of local average treatment effects. *Econometrica* 62 (2), 467.
- Ipsos, 2019. Ältere Leben Umweltbewusster Als Die Jugend - Aber Umwelverhalten Ändert Sich Bei Den Jungen Am Stärksten. Press release.
- Ito, K., Ida, T., Tanaka, M., 2018. Moral suasion and economic incentives: Field experimental evidence from energy demand. *Am. Econ. J.: Econ. Policy* 10 (1), 240–267.
- Jamison, J.C., Karlan, D., Zinman, J., 2014. Financial Education and Access to Savings Accounts: Complements or Substitutes? Evidence from Ugandan Youth Clubs. NBER Working Paper 20135.
- Jessoe, K., Lade, G.E., Loge, F., Spang, E., 2021. Spillovers from behavioral interventions: Experimental evidence from water and energy use. *J. Assoc. Environ. Res. Econ.* 8 (2), 315–346.
- Jessoe, K., Rapson, D., 2014. Knowledge is (less) power: Experimental evidence from residential energy use. *Amer. Econ. Rev.* 104 (4), 1417–1438.
- Karlin, B., Zinger, J.F., Ford, R., 2015. The effects of feedback on energy conservation: A meta-analysis. *Psychol. Sci.* 141 (6), 1205–1227.
- Khanna, T.M., Baiocchi, G., Callaghan, M., Creutzig, F., Guías, H., Haddaway, N.R., Hirth, L., Javaid, A., Koch, N., Laukemper, S., et al., 2021. A multi-country meta-analysis on the role of behavioural change in reducing energy consumption and CO2 emissions in residential buildings. *Nat. Energy* 6 (9), 925–932.
- Kollmuss, A., Agyeman, J., 2002. Mind the Gap: Why do people act environmentally and what are the barriers to pro-environmental behavior? *Environ. Educ. Res.* 8 (3), 239–260.
- List, J.A., Metcalfe, R.D., Price, M.K., Rundhammer, F., 2017. Harnessing Policy Complementarities to Conserve Energy: Evidence from a Natural Field Experiment. NBER Working Paper 23355.
- Mbiti, I., Muralidharan, K., Romero, M., Schipper, Y., Manda, C., Rajani, R., 2019. Inputs, incentives, and complementarities in education: Experimental evidence from Tanzania. *Q. J. Econ.* 134 (3), 1627–1673.
- Muralidharan, K., Romero, M., Wüthrich, K., 2020. Factorial designs, model selection, and (incorrect) inference in randomized experiments. *Rev. Econ. Stat.* 1–44.
- Myers, E., Souza, M., 2020. Social Comparison Nudges Without Monetary Incentives: Evidence from Home Energy Reports. *J. Environ. Econ. Manag.* 101, 102315.
- Nolan, J.M., Schultz, P.W., Cialdini, R.B., Goldstein, N.J., Griskevicius, V., 2008. Normative social influence is underdetected. *Personal. Soc. Psychol. Bull.* 34 (7), 913–923.
- Rubin, M., 2021. When to adjust alpha during multiple testing: A consideration of disjunction, conjunction, and individual testing. *Synthese* 199 (3–4), 10969–11000.
- Schwartz, D., Loewenstein, G., 2017. The chill of the moment: Emotions and proenvironmental behavior. *J. Public Policy Market.* 36 (2), 255–268.
- Sherif, R., 2021. Are Pro-environment Behaviours Substitutes or Complements? Evidence from the Field. Max Planck Institute for Tax Law and Public Finance Working Paper 2021 – 03.
- Tiefenbeck, V., 2016. On the magnitude and persistence of the Hawthorne effect — Evidence from four field studies. In: 4th European Conference on Behaviour and Energy Efficiency. Coimbra, Portugal.
- Tiefenbeck, V., Goette, L., Degen, K., Tasic, V., Fleisch, E., Lalive, R., Staake, T., 2018. Overcoming salience bias: How real-time feedback fosters resource conservation. *Manage. Sci.* 64 (3), 1458–1476.
- Tiefenbeck, V., Staake, T., Roth, K., Sachs, O., 2013. For better or for worse? Empirical evidence of moral licensing in a behavioral energy conservation campaign. *Energy Policy* 57, 160–171.
- Tiefenbeck, V., Woerner, A., Schoeb, S., Fleisch, E., Staake, T., 2019. Real-time feedback promotes energy conservation in the absence of volunteer selection bias and monetary incentives. *Nat. Energy* 4, 35–41.
- Tolstoy, L., 2003. Anna Karenina. In: (First published in Russian, 1873-1877; translation by Richard Pevear and Larissa Volokhonsky), Penguin Books, London.
- Tonke, S., 2019. Imperfect Knowledge, Information Provision and Behavior: Evidence from a Field Experiment to Encourage Resource Conservation. Working Paper.
- Trachtmann, H., 2022. Does Promoting One Behavior Distract from Others? Evidence from a Field Experiment. Tech. rep., forthcoming at AEJ:Applied (<https://www.aeaweb.org/articles?id=10.1257/app.20210788>).
- Wichman, C.J., 2017. Information provision and consumer behavior: A natural experiment in billing frequency. *J. Public Econ.* 152, 13–33.
- Young, A., 2019. Channeling Fisher: Randomization tests and the statistical significance of seemingly significant experimental results. *Q. J. Econ.* 134 (2), 557–598.