# Energy social surveys replicated with Large Language Model agents

*Michael J. Fell, UCL Energy Institute, University College London, London, UK.* [michael.fell@ucl.ac.uk](mailto:michael.fell@ucl.ac.uk)*.*

## Abstract

Large Language Models (LLMs) are artificial intelligence systems trained to understand and predict human language. In this study I programmatically create numerous LLM agents with population-representative characteristics, and prompt them provide survey responses with the aim of replicating existing energy social survey findings. Three studies are replicated, yielding moderate to high degrees of fidelity to the original results. Potentially significant contributions of the approach include improving the efficiency of research by identifying most promising interventions before conducting human studies, and simulating input from harder-to-access populations. However, there are also important practical and ethical challenges requiring of careful consideration.

*Keywords: large language models; artificial intelligence; energy; social surveys; replication*

**This is a non-peer-reviewed preprint uploaded on 6 January 2024.**

## Introduction

The potential of Large Language Models (LLMs) in contributing to scientific research has been widely discussed and demonstrated [1–3]. In the social sciences, potential contributions range from supporting analysis of various types of qualitative data [4], conducting interviews [5], informing hypothesis generation [6], and contributing to survey research [7]. However, their application in energy social studies has so far been extremely limited. With this study I attempt to replicate some of the findings of three survey-based studies using LLM-generated agents in place of human respondents. I draw on the experience and findings to suggest this and similar approaches could have significant impacts on the efficiency and nature of social research in energy, although it also comes with risks.

LLMs are a form of artificial intelligence (AI) model based on neural networks. Trained on large bodies of text of different kinds, they are effective natural language processors and can also generate coherent text. While the roots of LLMs in natural language processing stretch back decades, it is only in the last few years that their diverse capabilities have brought them to wider public attention. Prominent examples such as OpenAI's ChatGPT and Google's Bard now have hundreds of millions of users [8]. They are widely predicted to have substantial impacts on many sectors, including scientific research – both regarding productivity and the nature of work that is undertaken. Their application in the context of social surveys serves to illustrate this potential. Possible uses in this area include informing instrument design, sampling, response, data cleaning and analysis, and reporting [7].

Social surveys are a very widely used tool in energy research to collect structured quantitative and qualitative data from population samples, usually with the aim of making some inferences about the population itself. They permit easy experimental variation, as well as measurement of variables which are hard to capture in other ways, such as beliefs and opinions of large samples of individuals.

1

41  They are also often used to infer possible behaviour when this would be impossible or impractical
42  through other means. For example, if the cost of directly observing a behaviour would be
43  prohibitive, or if the behaviour relates to a product or service that is not available in the market – a
44  common situation in an energy sector addressing low-carbon transition. Surveys are often an
45  important prior step to more resource intensive field trialling of products, services, and
46  communication approaches.

47  Previous work has explored the potential of LLMs to act as a sample of multiple individuals for the
48  purpose of survey research. Argyle et al. (2023) [9] describe this as "silicon sampling", and argue that
49  because it allows for agent characteristics to be represented in proportion with the general
50  population, it can avoid some of the problems associated with algorithmic bias that may otherwise
51  be present in the LLM. Using this approach to investigate political questions such as vote prediction,
52  Argyle et al. [9] find strong alignment in patterns of association between human and GPT-3-
53  generated responses. Similar approaches have been used to replicate well-established findings in
54  psychology [10–12], management studies [13], and economics [14], amongst others. Aher et al. [10]
55  introduce the idea of a Turing Experiment which, instead of focusing on the possibility of
56  distinguishing between AI and individual responses, tests whether AI is able to faithfully simulate
57  human behaviour as revealed at the group level in experimental settings. In the context of
58  commercial user research, a variety of services are available which offer AI-user-based input,
59  although these do not appear to include survey response [15].

60  Where LLMs have been used by energy researchers, they have largely focused in conceptually
61  related areas such as systems modelling and technical data analysis [16,17]. While there are
62  examples with a social science element, such the generation of building occupant personas [18], as
63  yet their application in published social science research has been limited. To help assess their
64  suitability and consider potential applications and challenges, I attempted to replicate parts of three
65  studies. Two of these I originally co-authored, one on peer-to-peer (P2P) energy trading [19], the on
66  energy retail models based on multiple suppliers [20]. The other [21] is an independent study
67  investigating the role of defaults and trust in stated uptake of energy efficient appliances, and was
68  selected for it recency, rigor, and difference in topic and setting.

69

# Method

71  The LLM used was the OpenAI model GPT-3.5-turbo-1106, accessed via the OpenAI API. While not
72  the most recent model (which is gpt-4-1106-preview at time of writing) it is still extremely powerful
73  while being a tenth of the cost to use. Prompts submitted through the API consist of two parts,
74  "system" and "user". The system text sets the behaviour of the LLM, while the user text provides the
75  actual prompt to which the LLM will respond by way of an output. The system setting can be used to
76  endow the LLM with a range a characteristics such as might be found in an individual survey
77  respondent – for example, demographic, attitudinal, and personality characteristics. The user
78  prompt can then be used to submit the survey question to the "respondent", and the response
79  returned as an output and recorded. This is analogous to a response to a single survey question by a
80  single respondent. To collect a response to the same question by a different respondent, a new API
81  call is submitted with a newly generated combination of respondent characteristics in the system
82  setting.

2

83  Prompt length is a significant consideration because it has implications for cost and speed of API
84  calls. LLMs work based on tokens, which represent words or parts of words. The more tokens are
85  sent in a prompt or received in an output, the longer the time required and the higher the cost.
86  While the costs and time requirements are generally low for the basic exploratory work presented
87  here, they could become significant for future possible uses discussed later. As such, I attempted to
88  strike a reasonable balance between providing sufficient system/user prompt and output detail, and
89  time/cost.

90  To create the sample and run the data collection I created a Python program called PolyPersona
91  (available open source at https://github.com/mikefsway/polypersona). First, the program accesses a
92  spreadsheet holding a range of respondent characteristic variables and their possible states. As an
93  exploratory study based on trial and error, I included a range of standard demographic and
94  attitudinal variables often or sometimes measured in surveys in energy, including age, gender,
95  household income, household size, education, tenure, occupancy patterns, environmental concern,
96  economic rationality, risk aversion, social trust, political orientation, and place attachment. I also
97  included the "big 5" personality traits of openness to experience, conscientiousness, extraversion,
98  agreeableness, and neuroticism. While these traits are not typically measured in energy-related
99  surveys, they have been shown to have significant explanatory power in a range of areas of life
100  [22,23] and are a way of enriching the variety of responses. Finally, I included a variable for the
101  speed at which the survey was being completed. Respondents give differing levels of attention to
102  surveys, so to simulate this a variable was included for either slow response, giving careful thought
103  to the topic and drawing on all characteristics, and quick response prioritising personality traits and
104  stressing the need for an instinctive reaction.

105  Demographic characteristics could take specific states, such as an age of 20, 30, 40, etc. or a
106  household size of 4 or 5, and a probability for each state was assigned according to UK census data
107  where available. Most attitudinal and personality characteristics were assigned states of either high,
108  medium, or low, which could occur with equal probability. When PolyPersona generates a specific
109  respondent persona, it probabilistically assigns a state of each characteristic to the respondent. This
110  simplistic approach has the limitation of ignoring consistent associations between characteristics.
111  For example, home ownership shows a clear tendency to increase with age [24], and there are likely
112  many associations between different attitudinal and personality characteristics. Such associations
113  could be explored and, if shown to be impactful, captured in future iterations. Having generated a
114  range of characteristics, these are composed into a system prompt (see Box 1).

115  *Box 1: Example system prompt generated by PolyPersona.*

> You are a UK householder completing a survey in a hurry, so rely more on fast, intuitive
> judgement. You are male, aged 70, with a medium household income, 3 people in your
> household, you highest level of education is 'further', you are at home most of time, and you are a
> social tenant. You have these attitudes: high environmental concern, low risk aversion, medium
> social trust, left of centre politics, high place attachment, high economic rationality, and your
> innovation adoption status is 'late majority'. Your personality has traits of low extraversion,
> medium agreeableness, high conscientiousness, high neuroticism, and medium openness to new
> experience.

116

117 The program then composes a user prompt, which consists of the survey question and instructions
118 for the LLM on how to respond. The survey question is replicated as closely as possible from the
119 original study, along with any necessary preliminary information. Experimental variations are
120 included in a separate spreadsheet and selected probabilistically in the same way as for the
121 respondent characteristics. This is then concatenated with the instruction on how to respond. LLMs
122 work by sequentially generating the most likely next word – they "talk by talking". I therefore
123 prompt the model to first provide a brief explanation of its reasoning relating to the question,
124 followed by a consistent answer to the survey question (usually in the form of an integer). See the
125 supplementary material for example user prompts for each replication.

126 Responses are returned in JSON format and added to a spreadsheet along with a record of the
127 respondent characteristics. The program then loops and produces a new combination of
128 characteristics and question variation, and so on until the desired number of responses has been
129 obtained. This raw data can then be analysed according to the approach taken in the original study.
130 Here, this usually consisted of a simple graphical representation, with the addition of some statistical
131 analysis as described. The data generated for each of the replications is available at OSF
132 (https://osf.io/g7ua8/).

133

## Results

### Study 1

136 The first study I attempted to replicate was Fell et al. (2019) [19], which explores consumer demand
137 to participate in peer-to-peer (P2P) energy trading, a model whereby prosumers and consumers can
138 sell and buy electricity directly between themselves. A UK nationally representative sample of 2064
139 respondents were shown details of a P2P scheme, with experimental variation of the scale
140 (neighbourhood, city/region, national) and proportion of consumption supplied (25%, 50%, 75%),
141 and asked if they would choose to participate. The study subsequently investigates the effect that
142 scheme organiser and other characteristics have on demand, but I did not include these in this
143 exploratory replication. This was the replication I used in development of the PolyPersona program,
144 charts of preliminary runs are available in the supplementary materials.

145 Figure 1 shows the results of the replication based on 1798 PolyPersona responses. There appears to
146 be moderate agreement with the original study. Demand is highest for P2P energy at the regional
147 level in most consumption conditions, similar to the original study. Overall magnitude of demand is
148 broadly similar for regional and neighbourhood conditions, but substantially lower in the replication
149 for country-level supply. Some findings, such as the drop in demand for neighbourhood energy at
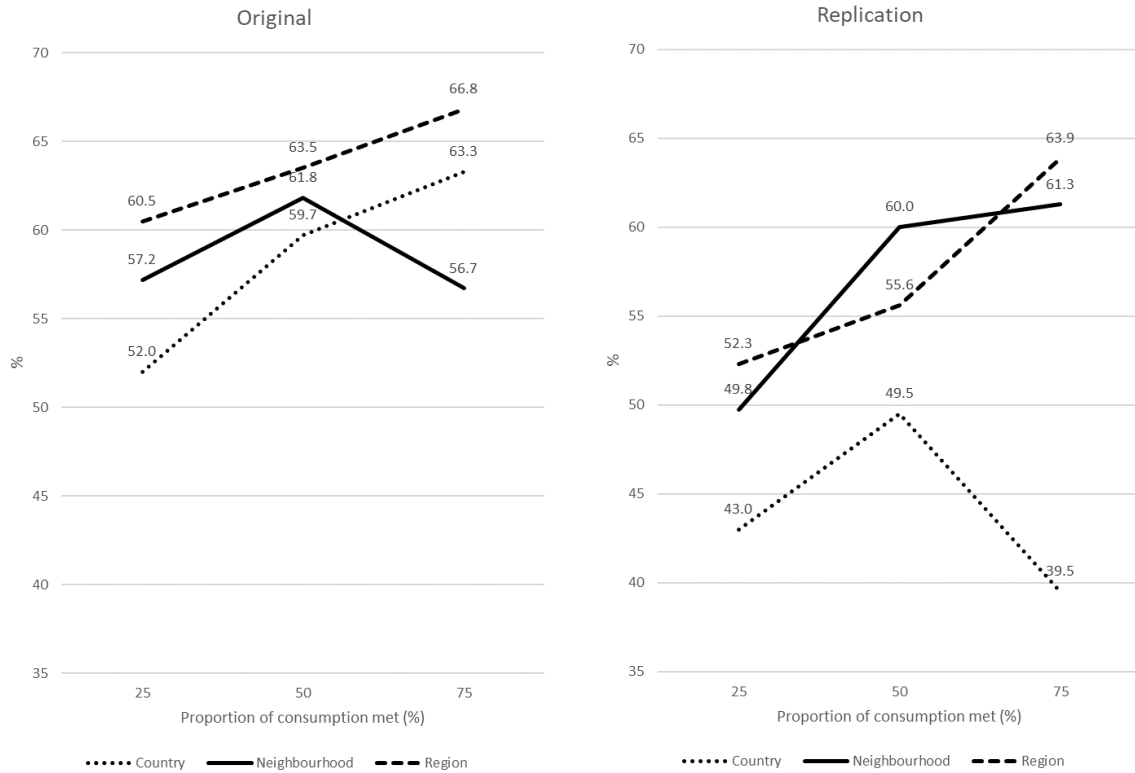150 75% consumption, are not evident.

151

*Figure 1: Percentage of respondents indicating willingness to participate in the P2P energy scheme, for the original study (left) and replication (right). Original data based on Fell et al., 2019.*

152  The original study reported some statistically significant differences between the experimental
153  groups. Binary logistic regression revealed significantly higher uptake at 75% consumption level
154  compared to 50% (odds ratio [OR] = 1.318, $p$ = 0.02). No such effect was evident for the replication,
155  although uptake at 25% consumption was significantly less than at 50% (OR = 0.763, $p$ = 0.02). For
156  spatial scale, the original study found significantly less uptake at the neighbourhood (OR = 0.766, $p$ =
157  0.03) and country (OR = 0.732, $p$ = 0.01) levels as compared to region. For the replication,
158  neighbourhood uptake was not significantly lower than regional, but country was (OR = 0.587, $p$ <
159  0.01). The effects observed in the original study were therefore partially replicated.

160  As described in the methods, the LLM-based responses required the generation of an explanation
161  prior to giving a participation decision. This has the co-benefit of making it possible to view a
162  qualitative explanation of the rationale of each decision. Table 2 illustrates this with examples of a
163  positive and negative participation decision, and shows how responses can vary depending on
164  whether the LLM was instructed to respond at leisure or in a hurry.

165    *Table 2: Examples of the explanations the LLM generated to inform its decisions.*

| Response speed | Decision to participate | Explanation |
|---|---|---|
| At leisure | Positive | As a homeowner with medium environmental concern and a general willingness to adopt new innovations, I would be open to the idea of buying energy directly from local homes and businesses with solar panels. I have a medium place attachment, so supporting local energy production aligns with my values. |
|  | Negative | As a homeowner with medium risk aversion and low environmental concern, I prefer to stick with my current supplier for reliability and stability. Despite the cost savings, I value the consistency of my current energy supply over the potential benefits of a cheaper, decentralized option. |
| In a hurry | Positive | I prefer to support local energy initiatives and as an early adopter, I am open to trying new energy solutions. The slight cost saving and the community aspect appeal to me. |
|  | Negative | I don't have high environmental concern, and I prefer to stick with what I know. Plus, I'm not usually at home to engage with my neighbours in this way. |

166

167    Study 2

168    The second study I replicated was Watson et al., (2020) [20], which explores consumer demand for
169    electricity retail models involving contracting with more than one supplier. Such models could unlock
170    the ability of specialist suppliers to supply electricity only from, from example, local sources or for
171    specific appliances. A representative sample of 1200 respondents were shown letters from energy
172    suppliers asking them if they would like to sign up to a multiple supplier arrangement, meaning they
173    could choose to be supplied with local electricity (when available) by a local supplier, with back-up
174    provided by the original supplier. There were three conditions. The first framed the single original
175    supplier option as the default, with respondents having the option to add the local supply. The
176    second framed the multiple supplier option as the default, with respondents having the option to
177    opt back to the single supplier. The third requested respondents to make an active choice between
178    the two options.

179    Figure 2 shows the results of the replication based on 600 PolyPersona responses. The general trend
180    across the conditions is the same, with most respondents choosing the multiple supplier option
181    when it is the default, followed by the active choice, followed by the single supplier option default.
182    The replication slightly underestimates multiple supplier take-up in the active choice and multiple
183    supplier default conditions, and is almost exactly the same for the single supplier default.
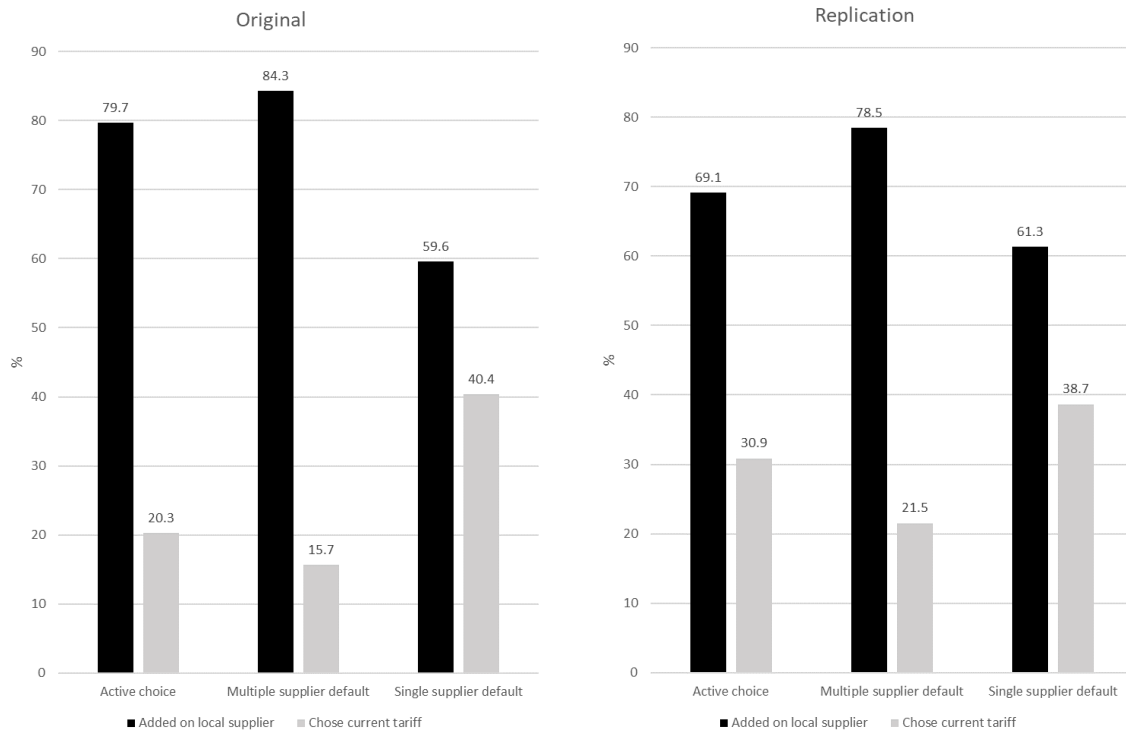
184

185



*Figure 2: Percentage of respondents choosing a local supplier (multiple suppliers), or stick with their single (current) tariff, for the original study (left) and replication (right). Original data based on Watson et al., 2020.*

186    A 2x3 Chi-squared test for the original study showed a statistically significant effect of condition on
187    respondents' choices  ($\chi^2$ = 67.6, $p$ < 0.01). This was also the case the replication, albeit with a
188    smaller effect size ($\chi^2$ = 14.87, $p$ < 0.01). Pairwise Chi-squared tests in the original study found that
189    the single supplier default group responses differed significantly from those of the other groups. This
190    was partially reproduced in the replication, where the single and multiple supplier default groups
191    differed significantly ($\chi^2$ = 14.06, $p$ < 0.01), but the single supplier and active choice groups did not ($\chi^2$
192    = 3.55, $p$ = 0.06). This may also have been a function of the smaller sample size.

193    I conducted a sensitivity analysis for this study, to explore the effect of including or excluding
194    different sets of respondent characteristics (i.e. demographic, attitudes, personality), as well as
195    survey completion speed. I ran ten variants in total. The first specified only that the respondent was
196    "a UK householder completing a survey". Seven subsequent runs included all combinations of
197    characteristic sets (i.e. demographics only, attitudes only, personality traits only, demographics and
198    attitudes, etc.). The final two runs included all characteristic sets, with one run specifying that all
199    were completed fast ("in a hurry, so rely more on fast, intuitive judgement"), and the other that all
200    were completed slow ("at your leisure, so rely more on calculated, rational thought"). Figure 3
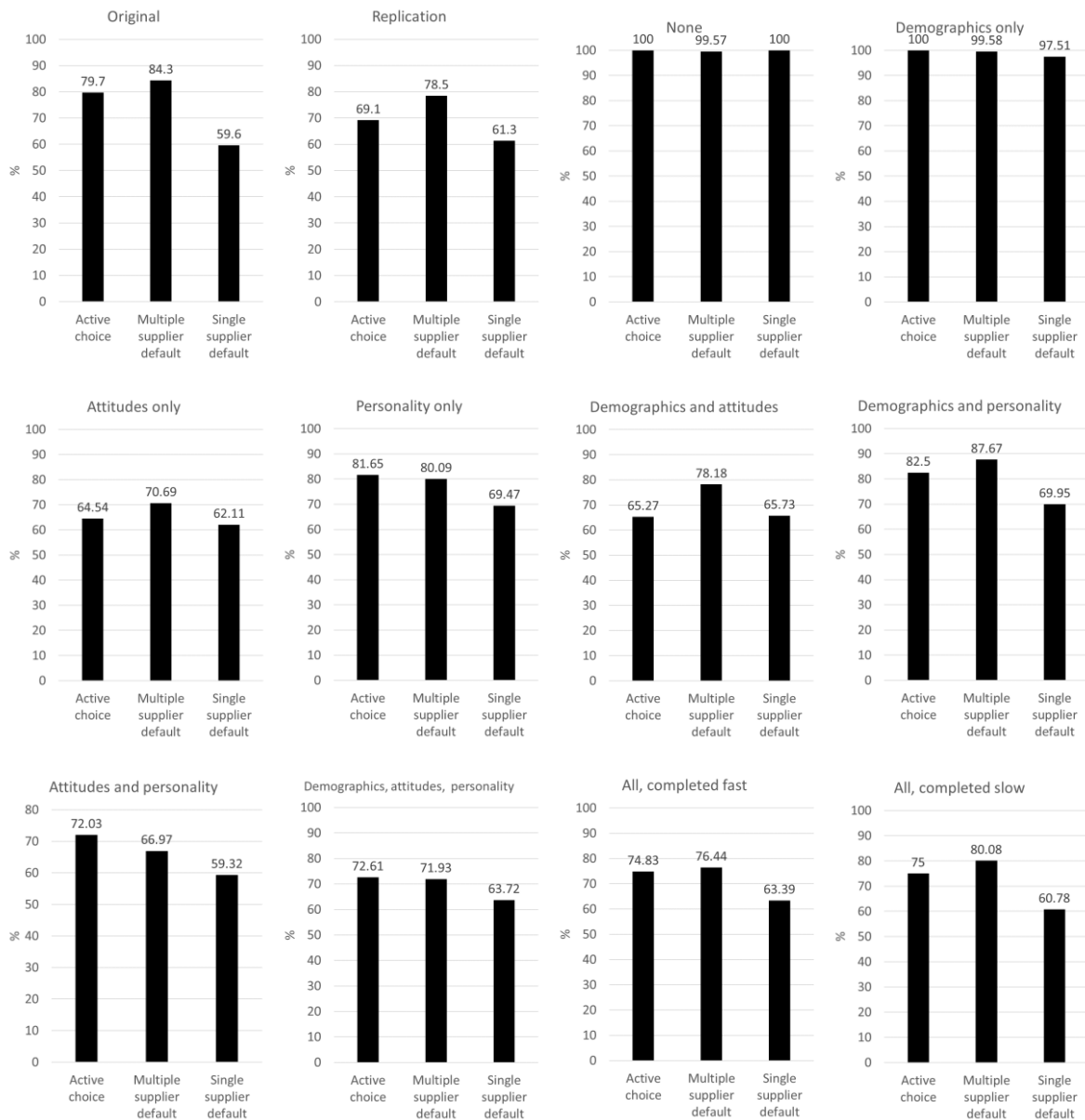201    shows the results of the sensitivity analysis.

202

203

204

*Figure 3: Sensitivity analysis of study 2. Bars indicate those respondents selecting the multiple supplier option in each case.*
*The first two graphs show the original results and my main replication, for comparative purposes.*

Providing the LLM with no varying characteristics, or varying only demographic factors, results in negligible outcome variation with practically all respondents selecting the multiple supplier option. All other combinations result in some level of variation. The combination with the closest pattern of resemblance to the original study includes all characteristics and slow completion, suggesting this could be a preferable option to the mix of fast (33%) and slow (67%) used in the main replications.

The original study yielded the unexpected finding that the preference for adding the local supplier was strong enough to override the default effect. It is striking that this finding is replicated here. For this study, to test the effectiveness of using a multiple-agent approach as compared to simply asking ChatGPT for a prediction, I submitted all three scenario descriptions directly to ChatGPT-4 in the form of a zero shot prompt, and requested it to predict the results of a nationally representative survey. Its prediction is shown in Figure 4 (see supplementary material for full prompt and response). It differs substantially from the actual findings, and includes the expected default effect.

8

219    While it is feasible that prompt chaining approaches could improve the likelihood of achieving a
220    more faithful replication (see [6]), this illustrates the value of the multiple LLM agent approach.
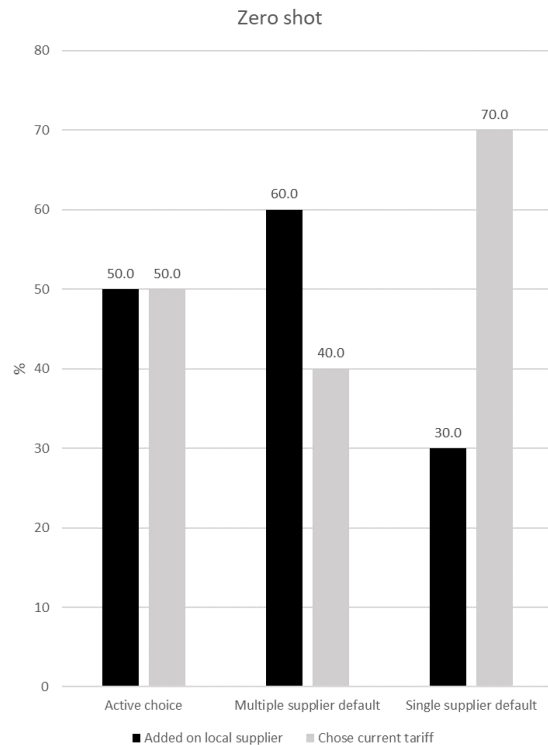


221

222    *Figure 4: ChatGPT (GPT-4) prediction for the results of Watson et al. (2020) when prompted directly with a zero-shot*
223    *prompt.*

## Study 3

224    Study 3
225    The third and final study I selected for replication was Kuhn et al. (2023) [21]. I wanted to use a
226    recent, rigorous study (which would not appear in the LLM training data), on a topic distinct from
227    the previous two, conducted outside of the UK, and selected this paper based on searches on a
228    convenience basis. It focuses on the role of defaults and trust in a housing provider to explore the
229    effect of these factors on stated uptake of a range of energy efficient appliances. The 956
230    participants in the US branch of the survey were experimentally shown either positive or negative
231    news headlines regarding real estate, and then presented with a choice task where either standard
232    or efficient appliances were the default condition, with the outcome being the number of efficient
233    appliances chosen.

234    Figure 5 shows the replication results for the replication of [21] based on 942 PolyPersona runs.
235    Again, a broadly similar pattern of responses can be observed between the original and replication.
236    A greater number of efficient appliances are chosen where efficient appliances are the default in
237    both case, while promoting mistrust of the housing provider results in lower efficient appliance
238    choice. The default effect is less pronounced in the replication, perhaps because it was not possible
239    to use the pre-checked box approach as was the case in the original study.
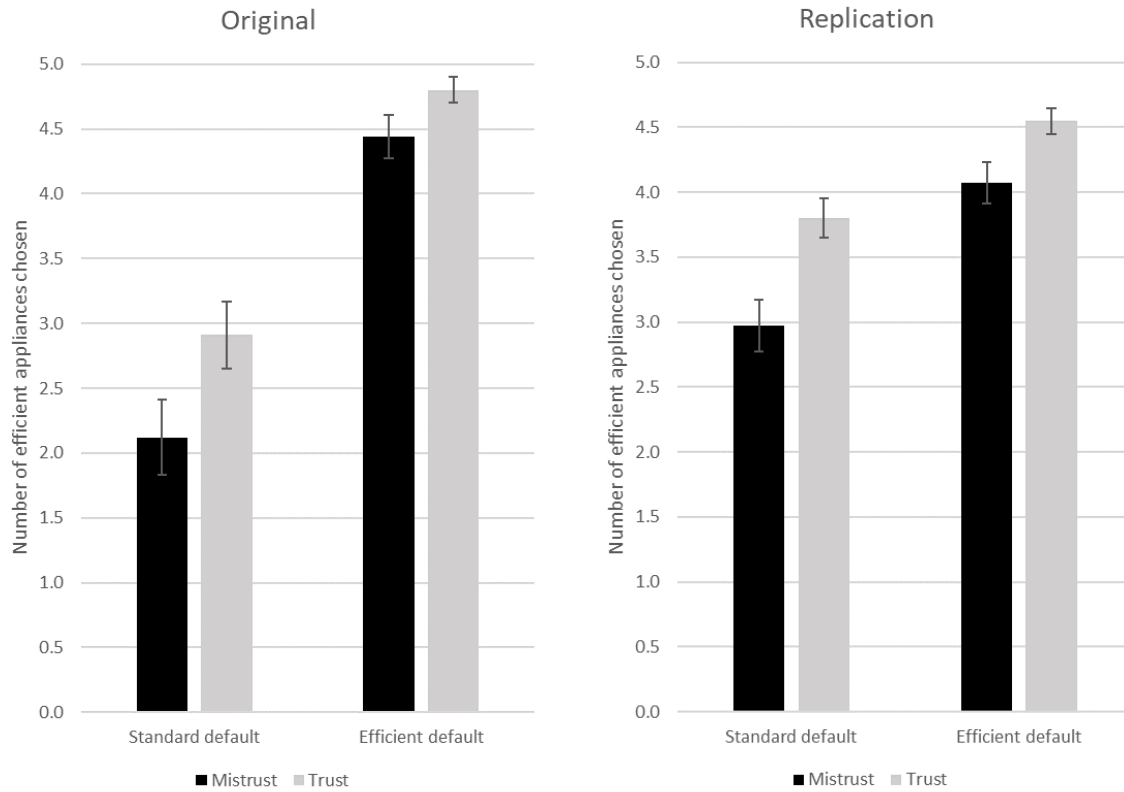
240  *Figure 5: Number of energy efficient appliances chosen by respondents in each experimental condition, for the original*
241  *study (left) and replication (right). Error bars show 95% confidence intervals. Original data based on Kuhn et al., 2023.*

242  Multiple linear regression analysis with the dependent variable of number of efficient appliances
243  chosen showed statistically significant positive main effects of default and trust, and a negative
244  interaction effect, as detailed in Table 3. These statistically significant effects were also detected in
245  the original study, although with a minor variation (the interaction effect was only significant when
246  also controlling for covariates).

247  *Table 3: Multiple regression results for replication of Kuhn et al., 2023. Adjusted $R^2$=.182.*

| Variable | B | Std Error | p |
|---|---|---|---|
| Efficient appliance default | 0.923 | 0.077 | < .001 |
| Trust-promoting headlines | 0.658 | 0.077 | < .001 |
| Interaction | -0.345 | 0.155 | 0.026 |

248

## Discussion and conclusions

250  In this study, I used LLM-based agents to replicate the findings of previous social surveys in energy.
251  The replication of [19] showed some similar trends to the original, such as a generally high
252  acceptability of P2P schemes and a tendency to higher uptake of schemes operating at a regional
253  level. Statistical test results were partially replicated. However, some key observations from the
254  original study were not apparent. The replications of [20] and [21] showed results that were much
255  more consistent with those of the original study, including capturing unanticipated results. Most

10

256  statistically significant findings were replicated. Overall, similar or the same conclusions could have
257  been drawn from the replications as the original studies. Viewed as a Turing Experiment [10], this
258  attempt to replicate the behaviour of human samples using AI could therefore be considered
259  reasonably successful.

260  Nevertheless, given evident variability and the fact that this work is still at an early developmental
261  stage, it is far from possible to conclude that LLM-based agents could reliably replace human
262  respondents at this time (nor would this be desirable in many cases – see discussion below).
263  However, the results are promising enough to suggest that further exploration of possible
264  applications of the approach in energy research is warranted (users can experiment with the
265  PolyPersona tool available open source at https://github.com/mikefsway/polypersona).

266  First, it is important to stress that LLMs do not model human thought processes, but rather replicate
267  a human-like response based on a vast body of training data. The explanation provided by the LLM,
268  and its subsequent decision, is probabilistically selected from all possible responses based on the
269  words produced so far and given the constraints presented [25]. It represents a highly sophisticated
270  black-box method of text sampling. None of my analyses were preregistered, although I report the
271  first and only replication attempt unless otherwise stated (i.e. replication of [19], and sensitivity
272  analysis).

273  It is also necessary to highlight important ethical considerations surrounding the use of LLMs in this
274  kind of research context. A key concern is around algorithmic bias, whereby biases present in
275  training data are replicated in LLM outputs [2]. This is a nuanced problem. In many cases, surveys are
276  explicitly attempting to understand the existence and implications of bias present in the population,
277  so in this respect the presence of bias (if appropriately controlled) can be helpful – what Argyle et al.
278  [9] refer to as algorithmic fidelity. It is striking that sensitivity analysis showed that varying only
279  demographic characteristics led to negligible variation in results across condition in study. This may
280  be due to developer efforts to mitigate bias based on demographic factors [26]. However, it still
281  likely that systematic underrepresentation of certain viewpoints in training data could lead to
282  misleading outputs. Similarly, certain biases of interest to researchers may not manifest in survey
283  outputs unless they are explicitly endowed with them [14], as may have been the case for risk
284  aversion as used here. A further ethical risk is the possibility that LLMs are used to produce fake
285  survey responses, either by malicious respondents or by researchers themselves. However, the
286  approach could also bring a range of benefits, including ethical ones.

287  The major advantage of "silicon sampling" is its cost- and time-effectiveness compared to human
288  surveys. The full replications included here cost a US dollar or less each using the GPT-3.5 API –
289  several orders of magnitude less than a comparable human survey – and depending on the number
290  of participants could be run in a timeframe of minutes to hours. Using PolyPersona, it is quick and
291  easy to set up and modify a variety of respondent characteristics, survey questions, and
292  experimental conditions. As such, it is reasonable to expect the approach to be used at the scoping
293  stage of research, for example to help identify conditions that have the biggest effects, and what the
294  sizes of these effects are. Such approaches, if they provide reliable, could help improve research
295  efficiency by focusing it in areas more likely to yield the biggest impacts. This has benefits to
296  researchers, research funders, and wider society. By analysing the explanations returned, it is also
297  possible to check if questions are being misinterpreted, and thereby make them clearer (if
298  necessary) for a human audience. This could help act as a supplement to human pilot testing. The

11

299  approach could also be helpful in producing pseudo data to help plan analyses before human data
300  are collected.

301  It is possible that LLM-augmented approaches could be useful where samples of interest may be
302  harder to access using standard survey techniques. For example, online surveys cannot represent
303  the input of non-Internet users. While LLM training data is drawn from the Internet, this does not
304  mean it is only representative of Internet users as a lot of content available online was previously
305  available in non-Internet formats. Especially if combined with fine-tuning approaches (such as
306  training on oral history databases), the approach could provide a powerful way to provide insight
307  into underrepresented viewpoints. This is a nuanced question which opens an important debate, as
308  there is also a risk that the approach could be used to substitute for genuine engagement with
309  underrepresented groups.

310  Related to the question of sampling, the approach raises the intriguing possibility of being able to
311  explore the impact of respondent characteristics which are otherwise hard to measure or select for.
312  Personality traits, for example, while commonly measured in some disciplinary areas, are not
313  routinely captured in energy-related surveys and would be quite resource-intensive to do so using
314  established measurement scales. But it is straightforward to endow LLM-based agents with them,
315  and it is possible that these or other uncommonly measured characteristics could have unexpectedly
316  large impacts on outcomes. Preliminary LLM-based work could help researchers decide which might
317  be the most important characteristics to measure and select for, given limited resources.

318  There may be circumstances where the use of LLM-based respondents in place of a human survey
319  could be justified, although extreme caution should be exercised here. An example could be where
320  there is no data on a topic of interest, no resources to conduct a survey, and a "best guess" is
321  needed to facilitate further work. In such a case, an LLM-based approach may provide a starting
322  point for further sensitivity analysis. Use of such an approach should always be clearly highlighted
323  and justified, given the risk of reading too much validity into LLM-generated datasets that can
324  realistically resemble human-generated datasets.

325  The above discussion sets out some possible applications of an LLM-based survey approach in
326  energy social studies, but there are several specific areas where the work presented here could be
327  developed. Further work could explore what affects whether and why different characteristics are
328  prioritised in formulating explanations and decisions. Combinations of characteristics could be
329  identified that maximise general applicability with efficiency of token use. The impact of fine tuning
330  could be investigated, especially in the case of harder-to-access samples (as mentioned above) or
331  more niche topics. There is also a question as to whether simulating survey responses is the most
332  useful end point. Surveys have well-known limitations, especially where the research interest is in
333  behaviours where what is done (revealed preference) may differ significantly from what is said
334  (stated preference). It would be fruitful to explore whether LLM agents could be used in such a way
335  as to better reflect real behaviours in relation to energy.

336

## Acknowledgements

340 the purpose of open access, the author has applied a creative commons attribution (CC BY) licence
341 (where permitted by UKRI, 'open government licence' or 'creative commons attribution no-
342 derivatives (CC BY-ND) licence' may be stated instead) to any author accepted manuscript version
343 arising.

344

## Data availability

346 Datasets for the replications are available at https://osf.io/g7ua8/.

347

## References

349 1. Boyko J, Cohen J, Fox N, Veiga MH, Li JI-H, Liu J, et al. An Interdisciplinary Outlook on Large
350    Language Models for Scientific Research. arXiv; 2023. doi:10.48550/arXiv.2311.04929

351 2. Demszky D, Yang D, Yeager DS, Bryan CJ, Clapper M, Chandhok S, et al. Using large language
352    models in psychology. Nat Rev Psychol. 2023;2: 688–701. doi:10.1038/s44159-023-00241-5

353 3. Ziems C, Shaikh O, Zhang Z, Held W, Chen J, Yang D. Can Large Language Models Transform
354    Computational Social Science? Comput Linguist. 2023; 1–53. doi:10.1162/coli_a_00502

355 4. Bano M, Hoda R, Zowghi D, Treude C. Large language models for qualitative research in
356    software engineering: exploring opportunities and challenges. Autom Softw Eng. 2023;31: 8.
357    doi:10.1007/s10515-023-00407-8

358 5. Chopra F, Haaland I. Conducting Qualitative Interviews with AI. Rochester, NY; 2023.
359    doi:10.2139/ssrn.4583756

360 6. Rosenbusch H, Stevenson CE, van der Maas HLJ. How Accurate are GPT-3's Hypotheses About
361    Social Science Phenomena? Digit Soc. 2023;2: 26. doi:10.1007/s44206-023-00054-2

362 7. Jansen BJ, Jung S, Salminen J. Employing large language models in survey research. Nat Lang
363    Process J. 2023;4: 100020. doi:10.1016/j.nlp.2023.100020

364 8. Carter R. Google Bard vs ChatGPT: Battle of the AI Chatbots. In: UC Today [Internet]. 29 Dec
365    2023 [cited 4 Jan 2024]. Available: https://www.uctoday.com/unified-communications/google-
366    bard-vs-chatgpt-battle-of-the-ai-chatbots/

367 9. Argyle LP, Busby EC, Fulda N, Gubler JR, Rytting C, Wingate D. Out of One, Many: Using
368    Language Models to Simulate Human Samples. Polit Anal. 2023;31: 337–351.
369    doi:10.1017/pan.2023.2

370 10. Aher G, Arriaga RI, Kalai AT. Using Large Language Models to Simulate Multiple Humans and
371    Replicate Human Subject Studies. arXiv; 2023. Available: http://arxiv.org/abs/2208.10264

372 11. Jiang H, Zhang X, Cao X, Kabbara J. PersonaLLM: Investigating the Ability of Large Language
373    Models to Express Big Five Personality Traits. arXiv; 2023. doi:10.48550/arXiv.2305.02547

374 12. Park PS, Schoenegger P, Zhu C. Diminished Diversity-of-Thought in a Standard Large Language
375    Model. arXiv; 2023. doi:10.48550/arXiv.2302.07267

376　13.　Suguri Motoki FY, Monteiro J, Malagueño R, Rodrigues V. Charting New Territory: Using
377　　　　ChatGPT to Enhance Survey Instruments for Management Accounting Research. Rochester, NY;
378　　　　2023. doi:10.2139/ssrn.4595896

379　14.　Horton JJ. Large Language Models as Simulated Economic Agents: What Can We Learn from
380　　　　Homo Silicus? arXiv; 2023. doi:10.48550/arXiv.2301.07543

381　15.　Arora A. 6 AI tools for user research. In: Bootcamp via Medium [Internet]. 2 Oct 2023 [cited 4
382　　　　Jan 2024]. Available: https://bootcamp.uxdesign.cc/5-ai-tools-for-user-research-3365891a78c3

383　16.　Buster G. Large Language Models (LLMs) for Energy Systems Research. National Renewable
384　　　　Energy Laboratory (NREL), Golden, CO (United States); 2023 Nov. Report No.: NREL/PR-6A20-
385　　　　87896. Available: https://www.osti.gov/biblio/2216925

386　17.　Zhou M, Li F, Zhang F, Zheng J, Ma Q. Meta In-Context Learning: Harnessing Large Language
387　　　　Models for Electrical Data Classification. Energies. 2023;16: 6679. doi:10.3390/en16186679

388　18.　Anik SMH, Gao X, Meng N. Towards automated occupant profile creation in smart buildings: A
389　　　　machine learning-enabled approach for user persona generation. Energy Build. 2023;297:
390　　　　113485. doi:10.1016/j.enbuild.2023.113485

391　19.　Fell MJ, Schneiders A, Shipworth D. Consumer Demand for Blockchain-Enabled Peer-to-Peer
392　　　　Electricity Trading in the United Kingdom: An Online Survey Experiment. Energies. 2019;12:
393　　　　3913. doi:10.3390/en12203913

394　20.　Watson NE, Huebner GM, Fell MJ, Shipworth D. Two energy suppliers are better than one:
395　　　　Survey experiments on consumer engagement with local energy in GB. Energy Policy.
396　　　　2020;147: 111891. doi:10.1016/j.enpol.2020.111891

397　21.　Kuhn S, Thøgersen J, Kutzner F. No trust in the choice architect? No problem! On the minor
398　　　　role of trust for the effectiveness of default interventions promoting the choice of energy-
399　　　　efficient appliances. J Environ Psychol. 2023;91: 102115. doi:10.1016/j.jenvp.2023.102115

400　22.　Mammadov S. Big five personality traits and academic performance: A meta-analysis. J Pers.
401　　　　2022;90: 222–255. doi:10.1111/jopy.12663

402　23.　Zell E, Lesick TL. Big five personality traits and performance: A quantitative synthesis of 50+
403　　　　meta-analyses. J Pers. 2022;90: 559–573. doi:10.1111/jopy.12683

404　24.　Statista. England: homeowners by age 2022. In: Statista [Internet]. 2023 [cited 2 Jan 2024].
405　　　　Available: https://www.statista.com/statistics/321065/uk-england-home-owners-age-groups/

406　25.　Reynolds L, McDonell K. Prompt Programming for Large Language Models: Beyond the Few-
407　　　　Shot Paradigm. arXiv; 2021. doi:10.48550/arXiv.2102.07350

408　26.　Ferrara E. Should ChatGPT be Biased? Challenges and Risks of Bias in Large Language Models.
409　　　　First Monday. 2023 [cited 13 Jan 2024]. doi:10.5210/fm.v28i11.13346

410

411

## Supplementary material for "Energy social surveys replicated with Large Language Model agents", Fell 2024

### Preliminary runs of the replication of Fell et al., 2019



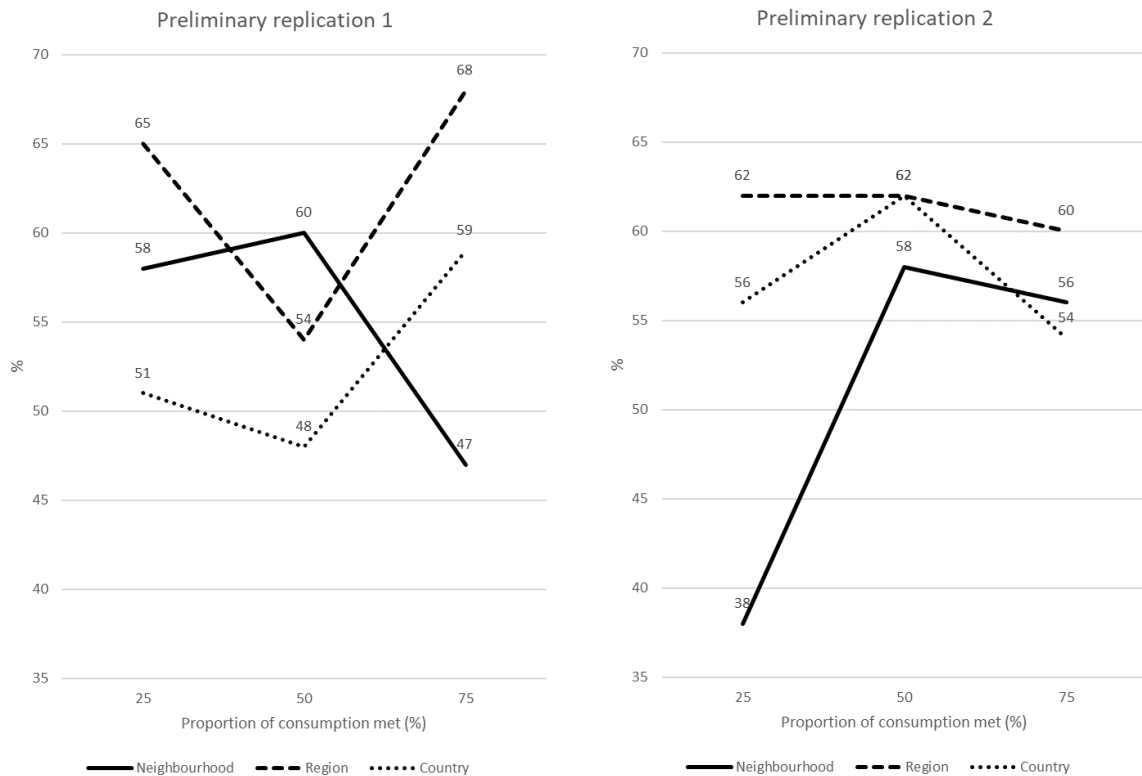| Preliminary replication 1 | Preliminary replication 2 |

*Figure S1: Test runs of the first replication. Preliminary replication 1 had 50 responses per experimental condition, while preliminary replication 2 had 100 responses.*

### Example user prompts

The following are example user prompts created for use in the replications. Italics are used to differentiate the response instructions from the survey question. For all experimental conditions, see https://osf.io/g7ua8/.

For the replication of Fell et al., 2019:

A new energy offer is available which allows you to buy energy directly from homes and businesses with their own solar panels. Energy bought this way is slightly cheaper than what you get from your usual supplier. If you participate, you could meet around a quarter of your household's electricity needs through the offer. You would buy electricity directly from homes and businesses located on the streets in your neighbourhood. You would continue to buy the rest of your energy from your current supplier. Would you sign up to participate in this offer if it was available to you today? *Respond in JSON format. The first JSON object should be a short (<50 word) explanation of your reasoning, called 'explanation', drawing on your demographic, attitudinal, and personality*

15

435　*characteristics. Then, output your decision on whether or not you would participate, in a JSON object*
436　*called 'decision', with response options yes=1 or no=0 (return an integer). Your decision must be*
437　*consistent with your explanation. Example output as follows: 'explanation': 'explanation text here',*
438　*'decision': integer*

439

440　For the replication of Watson et al., 2020:

441　Imagine that you receive the following letter from your current electricity supplier. Dear customer,
442　We are working in partnership with LocalEnergy, a new local energy company. You are currently
443　enrolled on our single-provider service. This means that all your electricity will continue to be
444　provided by us. Alternatively, you can switch to the new combined-provider service, which is offered
445　in partnership with LocalEnergy. Customers of the combined-provider service will get most of their
446　electricity from locally produced, renewable sources, such as solar and wind farms. If there is not
447　enough local energy to cover your needs, we will act as a back-up, so that you will always have the
448　power you need. The price of each service will be the same as your current tariff. If you would like to
449　stay with the single-provider service, you do not need to do anything. If you would rather change to
450　the combined-provider service, you need to go to our website, log in to your account, and check a
451　box indicating that you would like to switch. It will take a few days to switch. You can switch back at
452　any time. What would you choose to do in this scenario? Stay with the single-provider service, or
453　switch to the combined-provider service. *Respond in JSON format. The first JSON object should be a*
454　*short (<50 word) explanation of your reasoning, called 'explanation', drawing on your demographic,*
455　*attitudinal, and personality characteristics. Then, output your decision on whether you would choose*
456　*the single or combined provider option, in a JSON object called 'decision', with response options*
457　*combined=1 or single=0 (return an integer). Your decision must be consistent with your explanation.*
458　*Example output as follows: 'explanation': 'explanation text here', 'decision': integer*

459

460　For the replication of Kuhn et al., 2023:

461　Imagine you consider renting a newly constructed apartment rented out by the Better Living Real
462　Estate Group. Right now, this is your only available option. You have therefore decided to take it and
463　signed the contract. Better Living has a programme that supports your move to the new apartment
464　and is asking you to fill out an online form in which you choose amenities for your new home. Before
465　you fill out the form you come across some recent headlines from news websites that could be
466　relevant in this situation. Please read these headlines carefully: 'Award Winners: Affordable high-
467　quality housing and community project by Better Living wins the 2021 American Residential Real
468　Estate'; 'More transparency in real estate: Memorandum passed to combat corruption in U.S.
469　housing sector'; 'Buyers, Notice: Report shows real estate sector can be trusted'; 'The rise of the
470　housing market: What everyone gets wrong about housing in the U.S. Today, most public housing
471　developments do, in fact, deliver on their promise of providing adequate, affordable housing to
472　people in need, but they remain stigmatised because of the negative associations with distressed
473　public housing'; 'Americans deserve not just better, but the best: Administration is seeking largest
474　ever investment in building resilient, quality houses'. Consider how interesting you perceive each
475　article to be for your situation. Now read this letter from Better Living Real Estate Group: 'Dear
476　tenant, We are glad that you chose to rent with us. We want to make moving in to your new
477　apartment as easy and satisfying as possible for you. We have therefore preinstalled the following
478　amenities in your new home. They are standard amenities, but if you want we can replace them with

479    energy-efficient alternatives. Compared to the standard amenities, the energy-efficient amenities
480    cost more, but due to the savings on the electricity bill, the operating costs are lower. We estimate
481    that the energy-efficient amenities save the equivalent of the initial higher costs in about three
482    years. If you want to keep the standard amenities, you do not need to do anything. If you want us to
483    replace one or more with energy-efficient alternatives, you can let us know.' The appliances are as
484    follows: freezer; refrigerator; washing machine; electric oven and stove; air conditioner. *Respond in*
485    *JSON format. The first JSON object should be a short (<25 word) summary of the newspaper*
486    *headlines, called 'summary'. Then, in a JSON object called 'better living', consider the headlines and*
487    *state to what extent you trust that Better Living has the customers' best interests in mind, providing*
488    *your answer as an integer on a scale from 0 to 10, where 0 means 'do not trust at all' and a 10 means*
489    *'completely trust'. The next object should be a short explanation (<50 word) explanation of your*
490    *reasoning for your response to the Better Living Real Estate Group letter, called 'explanation',*
491    *drawing on your demographic, attitudinal, and personality characteristics, and bearing in mind the*
492    *headlines. Then, considering each amenity separately, output your decision on whether you would*
493    *like the standard or efficient version, in JSON objects called 'freezer'; 'refrigerator'; 'washing*
494    *machine'; 'electric oven and stove'; and 'air conditioner',  with response options efficient=1 or*
495    *standard=0 (return an integer). Example output as follows: 'summary': 'headline summary text here',*
496    *'better living': integer, 'explanation': 'explanation text here', 'freezer': integer, 'refrigerator': integer,*
497    *'washing machine': integer, 'electric oven and stove': integer, 'air conditioner': integer.*

498

## Regression table for replication of Fell et al., 2019

500

501    *Table S1: Results of binary logistic regression for replication of Fell et al., 2019.*

| Model | Deviance | AIC | BIC | df | $X^2$ | p |
|---|---|---|---|---|---|---|
| H0 | 2487.213 | 2489.213 | 2494.708 | 1797 | | |
| H1 | 2452.985 | 2462.985 | 2490.457 | 1793 | 34.228 | < .001 |

502

| | Estimate | Standard Error | Odds Ratio | z | Wald Statistic | df | p |
|---|---|---|---|---|---|---|---|
| (Intercept) | 0.382 | 0.106 | 1.465 | 3.588 | 12.872 | 1 | <0.001 |
| 25% consumption | -0.27 | 0.117 | 0.763 | -2.309 | 5.331 | 1 | 0.021 |
| 75% consumption | 0 | 0.118 | 1 | <0.001 | <0.001 | 1 | 1 |
| neighbourhood | -0.01 | 0.117 | 0.99 | -0.084 | 0.007 | 1 | 0.933 |
| country | -0.533 | 0.116 | 0.587 | -4.595 | 21.113 | 1 | <0.001 |

503

504

## Direct prompting of ChatGPT as part of replication of Watson et al., 2020

506    Here are the GPT-4 prompt and response used to compare the multiple agent approach used in
507    replications with a direct prompt.

508 **Prompt**: I am running a survey. It presents each participant with one of the following three scenarios
509 at random: 1. "Imagine that you receive the following letter from your current electricity supplier.
510 Dear customer, We are working in partnership with LocalEnergy, a new local energy company. You
511 are currently enrolled on our single-provider service. This means that all your electricity will continue
512 to be provided by us. Alternatively, you can switch to the new combined-provider service, which is
513 offered in partnership with LocalEnergy. Customers of the combined-provider service will get most
514 of their electricity from locally produced, renewable sources, such as solar and wind farms. If there is
515 not enough local energy to cover your needs, we will act as a back-up, so that you will always have
516 the power you need. The price of each service will be the same as your current tariff. If you would
517 like to stay with the single-provider service, you do not need to do anything. If you would rather
518 change to the combined-provider service, you need to go to our website, log in to your account, and
519 check a box indicating that you would like to switch. It will take a few days to switch. You can switch
520 back at any time. What would you choose to do in this scenario? Stay with the single-provider
521 service, or switch to the combined-provider service." 2. "Imagine that you receive the following
522 letter from your current electricity supplier. Dear customer, We are working in partnership with
523 LocalEnergy, a new local electricity company. You have automatically been switched to our new
524 combined-provider service, which is offered in partnership with LocalEnergy. Customers of the
525 combined-provider service will get most of their electricity from locally produced, renewable
526 sources, such as solar and wind farms. If there is not enough local energy to cover your needs, we
527 will act as a back-up, so that you will always have the power you need. Alternatively, you can switch
528 back to our single-provider service. This means that all your energy will continue to be provided by
529 us. The price of each service will be the same as your current tariff. If you would like to stay with the
530 combined-provider service, you do not need to do anything. If you would rather change back to the
531 single-provider service, you need to go to our website, log in to your account, and check a box
532 indicating that you would like to switch. It will take a few days to switch. You can switch back at any
533 time. What would you choose to do in this scenario? Stay with the combined-provider service, or
534 switch back to the single-provider service.". 3. "Imagine that you receive the following letter from
535 your current electricity supplier. Dear customer, We are working in partnership with LocalEnergy, a
536 new local energy company. You now have the choice of two services: 1) Single-provider: This means
537 that all your electricity will be provided by us. 2) Combined-provider: This service is offered in
538 partnership with LocalEnergy. Customers of the combined-provider service will get most of their
539 electricity from locally produced, renewable sources, such as solar and wind farms. If there is not
540 enough local energy to cover your needs, we will act as a back-up, so that you will always have the
541 power you need. The price of each service will be the same as your current tariff. You need to go to
542 our website, log in to your account, and check a box indicating which service you would like to
543 choose. You can switch at any time. Which service would you choose in this scenario? The single-
544 provider service or the combined-provider service." Imagine I administer this survey to a
545 representative sample of 1200 UK citizens. Please consider what you think will affect the results,
546 then predict what proportion of people will choose each outcome option in each scenario, and
547 create a bar chart.

548

549 **Response**: To predict the outcomes of your survey with these three scenarios, we need to consider
550 several factors that might influence the decisions of the participants:

551    1.   Status Quo Bias: People have a tendency to prefer the current state of affairs. In scenarios
552         where a default option is presented (Scenario 1 and 2), many might choose to stick with the
553         default due to inertia or the perceived effort of changing.

18

554      2. Environmental Concerns: Given the emphasis on renewable energy in the scenarios,
555          individuals who are more environmentally conscious may be more inclined to choose the
556          combined-provider service.
557      3. Perceived Effort and Convenience: The effort required to switch services (or stay with the
558          current one) might influence decisions. If switching seems cumbersome, people might avoid
559          it.
560      4. Trust in New Services: People's trust in the new combined-provider service, especially
561          regarding its reliability and the partnership with LocalEnergy, will play a role.
562      5. Price Sensitivity: Since the price remains the same for both services, this factor might not
563          heavily influence the decision.
564  Based on these considerations, here are the predictions:

565

566      1. Scenario 1 (Default is single-provider): A significant number of people might stay with the
567          single-provider service due to status quo bias. However, some environmentally conscious
568          individuals might switch to the combined-provider service.
569      2. Scenario 2 (Default is combined-provider): A large portion might stay with the combined-
570          provider service, again due to status quo bias. The effort to switch back to the single-
571          provider might deter changes.
572      3. Scenario 3 (No default, active choice required): This scenario is likely to show a more
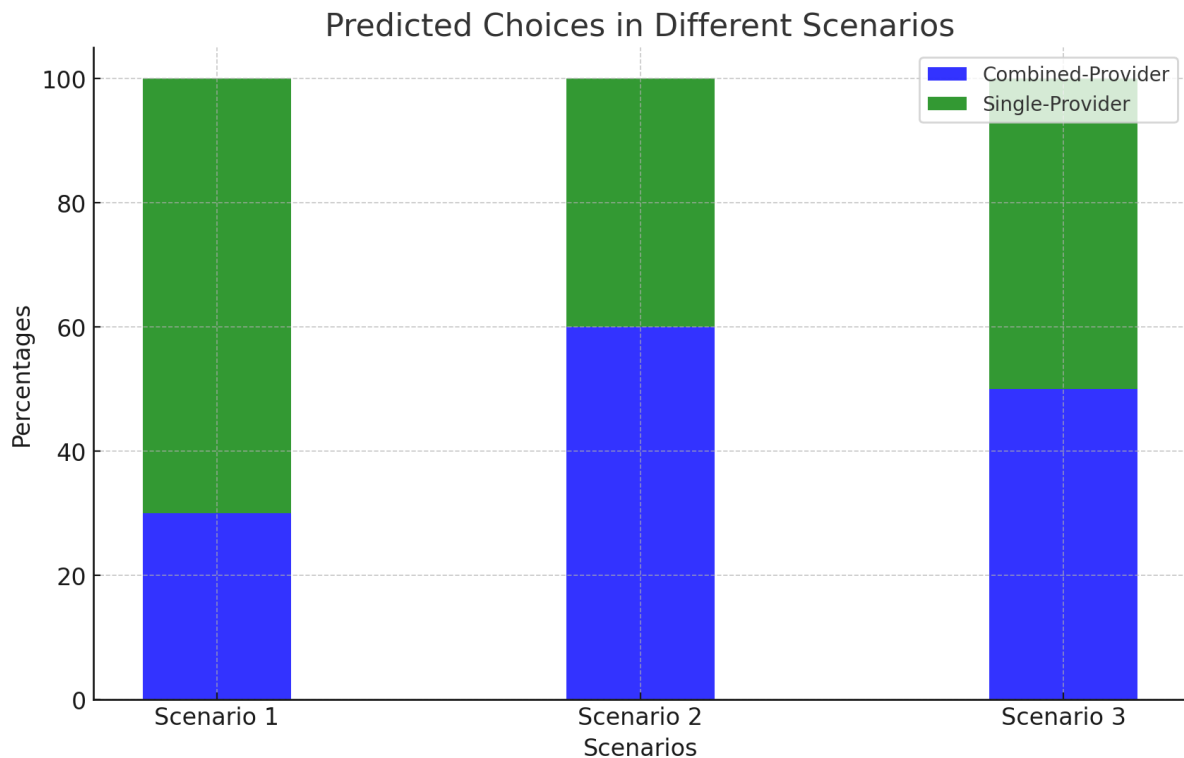573          balanced distribution, as people are required to make a deliberate choice without a default
574          option.
575  Now, I'll create a bar chart based on these predictions. Let's assume the following proportions for
576  each scenario:

577  Scenario 1: 30% choose combined-provider, 70% stay with single-provider.

578  Scenario 2: 60% stay with combined-provider, 40% switch back to single-provider.

579  Scenario 3: 50% choose each service.

580  Let's generate the bar chart with these assumptions.

## Predicted Choices in Different Scenarios



581

582    Here is the bar chart based on the predicted outcomes for each scenario in your survey:

583    Scenario 1: It shows that 30% of participants are predicted to choose the combined-provider service,
584    while 70% would stay with the single-provider service.

585    Scenario 2: Here, 60% are expected to stay with the combined-provider service, and 40% might
586    switch back to the single-provider service.

587    Scenario 3: This scenario predicts an even split, with 50% choosing each service.

588    These predictions are based on the factors discussed earlier, like status quo bias, environmental
589    concerns, perceived effort, and trust in new services. Keep in mind that the actual results might vary
590    based on specific attitudes and preferences of the surveyed UK citizens.

591

592

593

594

595

596