

1 A primer on running human behavioural experiments online

2

3 Tijl Grootswagers*

4 *School of Psychology, University of Sydney, NSW, 2006, Australia, tijl.grootswagers@sydney.edu.au

5

6 Abstract

7 Moving from the lab to an online environment opens up enormous potential to collect behavioural data
8 from thousands of participants with the click of a button. However, getting the first online experiment
9 running requires familiarisation with a number of new tools and terminologies. There exist a number of
10 tutorials and hands-on guides that can facilitate this process, but these are often tailored to one specific
11 online platform. The aim of this paper is to give a broad introduction to the world of online testing. This
12 will provide a high-level understanding of the infrastructure before diving into specific details with more
13 in-depth tutorials. Becoming familiar with these tools allows moving from hypothesis to experimental
14 data within hours.

15

16 Introduction

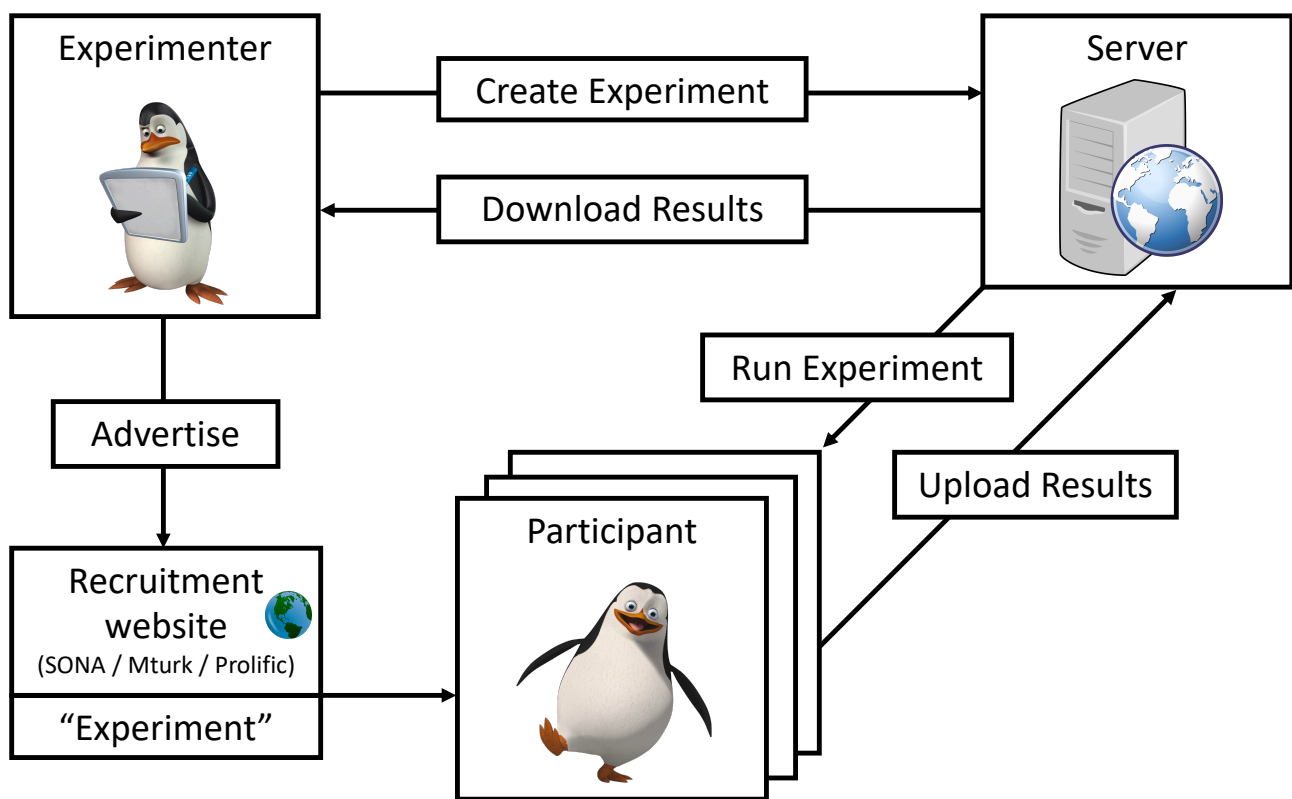
17 Lightning fast internet speeds and significant technological improvements have made it possible to
18 perform complex experiments within a modern web-browser. It is becoming increasingly popular to
19 combine browser-based experiments with recruiting participants on platforms such as Amazon's
20 Mechanical Turk (MTurk) or Prolific Academic (Palan & Schitter, 2018). There are several reasons why
21 researchers opt for online instead of lab-based testing. The first is efficiency. The recruitment platforms
22 (e.g., MTurk) have access to large numbers of participants, allowing to test many (thousands) of
23 participants simultaneously, which would not be possible in a lab-based setting. They are also not
24 restricted to office hours or teaching schedules, and do not require an on-campus presence for
25 participants or researchers. Secondly, participants from the online platforms are a better reflection of the
26 general population than the undergraduate students who typically participate in experiments on campus
27 (Berinsky et al., 2012). Finally, online experiments are more economical¹, because there is no need to
28 spend time recruiting, scheduling, and testing participants.

29 Our lab has had an overwhelmingly positive experience with running online studies (Grootswagers et al.,
30 in press, 2017, 2018). While early days involved extensive JavaScript programming for relatively simple
31 online studies, recent advancements have made it much easier to get complex studies up and running
32 (Anwyl-Irvine, Massonnié, et al., 2020; Barnhoorn et al., 2014; De Leeuw, 2015; Henninger et al., 2019;
33 Peirce et al., 2019). These generally come with associated tutorials and hands-on guides, but these are
34 often specific to a single platform or method. Therefore, it can be a challenge to get familiar with the
35 infrastructure, tools, and terminology, especially when starting out from scratch. This document aims to
36 facilitate this process by introducing the basics to online testing. It is intended to serve as a high-level
37 overview, and guide the reader to relevant in-depth literature, reviews, and tutorials.

¹ There has been discussion about online studies being exploitative but the experimenter can pay participants a fair compensation in accordance with institutional ethics review boards (c.f. Crump et al., 2013; Mason & Suri, 2012; Shank, 2016)

38 The basics

39 The core infrastructure needed for online experiments consists of: (1) a browser-based experiment (2) a
 40 server to host the experiment, and (3) a participant recruitment tool. Figure 1 illustrates the general
 41 infrastructure and workflow for online experiments. Experiments are programmed to run in a browser
 42 and hosted on a server. Participants are recruited from online marketplaces and perform the task on their
 43 local machine. The data is uploaded to the hosting server where the experimenter can collect the results.



44

45 **Figure 1. Infrastructure model for online experiments.**

46 Creating the experiment

47 The experiment needs to be able to run in a web-browser (e.g., Safari, Google Chrome, Internet
 48 Explorer). It therefore needs to be programmed in a browser-compatible programming language (e.g.,
 49 JavaScript, PHP). The most popular language for online experiments is JavaScript, and there exist several
 50 JavaScript modules (e.g., JsPsych, PsychoJS, OSWeb, Lab.js) tailored to behavioural experiments. The

51 libraries provide a number of high-level functions to facilitate experiment-specifics, such as presenting
52 stimuli, control timing, randomisation, and collecting responses. Some (e.g., Lab.js) are accompanied by
53 web-based task builders that allow creating experiments without the need for any programming. Several
54 free and open-source graphical experiment builders can export experiments as browser-compatible
55 JavaScript code. For example, Psychopy (Peirce et al., 2019) can export to PsychoJS, and OpenSesame
56 (Mathôt et al., 2012) to OSWeb. There also exist commercial solutions that provide experiment builders
57 as part of a complete experiment hosting infrastructure, such as Testable, Inquisit, and Gorilla (Anwyl-
58 Irvine, Massonnié, et al., 2020).

59 Deciding on a suitable experiment creation method is often a matter of personal preference. Experiment
60 builders are easy to use but can lack flexibility. Some JavaScript modules are also easier to use than others
61 and can be guided by previous experience in experiment programming. For example, PsychoJS has a
62 similar code structure to its Psychopy counterpart and may therefore be well suited for those already
63 experienced with coding Psychopy experiments in python.

64 Hosting the experiment

65 The experiment needs to be accessible to the world. This involves *hosting* the experiment code, stimuli,
66 and libraries on a server. This allows a participant to access the experiment code from their web browser.
67 The experiment then runs in the browser on the participant's computer. The participant completes the
68 experiment, and the script sends the participant's experimental data back to the server. This means that
69 the server should be able to receive and store the experiment data. Several paid hosting services exist that
70 are specifically aimed at collecting behavioural data online, such as Pavlovia, Gorilla, or Inquisit.
71 Alternatively, experiments can be directly hosted on a web server (or a cloud service such as Google or
72 Amazon). This requires knowledge of servers and security technology, but is flexible and allows for secure
73 and private data storage. JATOS (Lange et al., 2015) is an example of a free and open-source application
74 that facilitates the setting up and running of a web server for hosting online studies.

75 When choosing a hosting solution, factors to consider are the cost, flexibility, and ease of use.
76 Commercial services (e.g., Gorilla or Inquisit) are generally very user-friendly but also the most expensive
77 option and use their own experiment builders. Pavlovia is a non-commercial low-cost hosting service
78 that is still user-friendly and accommodates different types of JavaScript experiments. These hosting
79 services all charge a fee per participant or have limited term usage licenses. In contrast, JATOS is free
80 and open source software to host experiments that is flexible but requires more technical skills to set up
81 on a server.

82 Recruiting participants

83 The final step is to recruit participants. What is needed for this is a marketplace (on the web) where
84 participants can view and sign up for experiments. When they decide to participate, they get the link
85 (URL) to the experiment server and complete the task. Examples of such marketplaces are SONA
86 systems (often used for undergraduate testing at universities), MTurk, or Prolific (Palan & Schitter, 2018).
87 To be able to give participants compensation (e.g., course credits or payment) for their participation,
88 online experiments often display a unique code that participants can enter in the recruitment system so
89 the experimenter can verify their participation. It is useful to note the time zone of the participants, for
90 example, MTurk workers (based in the US) will be more likely to be online and see the experiment if it
91 is posted during their daytime. The recruitment systems will have the option to specify how many
92 participants are needed, and some provide additional screening criteria. When all participants have
93 completed the experiment, the researcher can simply download the data from the server and start
94 analysing.

95 Frequently asked questions

96 The basic infrastructure needed for online testing is not overly complex, as described in the previous
97 section. In addition, the available infrastructure has improved significantly in recent years with the
98 development of more sophisticated hosting solutions and programming libraries. Once familiar with

these powerful tools, it is extremely easy to go from hypothesis to experimental data within hours. The remainder of this paper will cover a number of frequently asked questions with regards to online testing.

How good are the data?

Several studies have compared data from online markets to data collected in the lab (Barnhoorn et al., 2014; Crump et al., 2013; de Leeuw & Motz, 2016; Simcox & Fiez, 2014; Zwaan & Pecher, 2012), with overall positive results. Tutorials and reviews have suggested that online experiment data is generally better when experiments are short, pay well, are fun, and have clear instructions. It is good to keep in mind that participants from online marketplaces (e.g., MTurk) are not as familiar with psychology experiments compared to undergraduate students. Therefore, it is essential to make very clear instructions and sometimes include a number of practice trials to ensure they understand the task.

How good is the timing?

Despite the progress in web-based technology, stimulus and response timing will be less reliable than the commercial equipment used in the lab. In general, latencies and variabilities are higher in web-based compared to lab-environments. Several studies have assessed the quality of timing in online studies, with encouraging results (Anwyl-Irvine, Dalmaijer, et al., 2020; Bridges et al., 2020; Pronk et al., 2019; Reimers & Stewart, 2015). An online evaluation of a masked priming experiment showed that very short stimulus durations (i.e., under 50ms) can be problematic (but see Barnhoorn et al., 2014), but other classic experimental psychology paradigms that rely on reaction times (e.g., Stroop, flanker, and Simon tasks) were successfully replicated (Crump et al., 2013).

What are the limitations?

Online experiments only work for some stimulus modalities. While the online approach is well suited for experiments consisting of visual stimuli and keyboard or mouse responses (but see previous question on timing), other paradigms are harder or impossible to move online. For example, studies requiring auditory

122 stimuli are possible (Cooke et al., 2011; Gibson et al., 2011; Schnoebelen & Kuperman, 2010; Slote &
123 Strand, 2016), but may necessitate a more extensive set-up procedure, such as procedures to make sure
124 the participant's set-up works. Presenting stimuli in other modalities, such as tactile or olfactory stimuli,
125 are impossible to achieve in an online environment.

126 A second limitation is the lack of experimental control. For example, while a participant's screen size is
127 reported by the browser, there is no way to know the participant's distance to screen. It is therefore
128 impossible to control the exact visual angle of stimuli, which can be a limiting factor for some
129 experiments. It is also hard to test whether participants are paying attention to the experiment. A common
130 approach is to exclude participants based on their performance on catch-trials (Mason & Suri, 2012). Still,
131 there can be a large amount of variability in attention amongst online participants and they could be
132 distracted by other sources while performing experiments, such as listening to radio, looking at their
133 phone, or watching their children.

134 Conclusion

135 Online experiments offer large-scale participant testing in a short time and are cheaper to run than their
136 lab-based counterparts. They can be a suitable option for many research questions but have some
137 limitations in the amount of experimental control. This manuscript has provided a high-level overview
138 of the infrastructure. For more in-depth reading, the reader is referred to the more specialised tutorials
139 and reviews cited above. The JavaScript experiment libraries (e.g., JsPsych, PsychoJS, Lab.js) also have
140 associated hands-on tutorials and contain many examples of classic cognitive science experiments, which
141 are a good place to start with programming the online experiment.

142 Open Practices Statement

143 Any relevant data and materials are available at <https://osf.io/xkdy4>

144 References

- 145 Anwyl-Irvine, A. L., Dalmaijer, E. S., Hodges, N., & Evershed, J. K. (2020). *Online Timing Accuracy and*
 146 *Precision: A comparison of platforms, browsers, and participant's devices* [Preprint]. PsyArXiv.
 147 <https://doi.org/10.31234/osf.io/jfec>
- 148 Anwyl-Irvine, A. L., Massonnié, J., Flitton, A., Kirkham, N., & Evershed, J. K. (2020). Gorilla in our
 149 midst: An online behavioral experiment builder. *Behavior Research Methods*, 52(1), 388–407.
 150 <https://doi.org/10.3758/s13428-019-01237-x>
- 151 Barnhoorn, J. S., Haasnoot, E., Bocanegra, B. R., & Steenbergen, H. (2014). QRTEngine: An easy
 152 solution for running online reaction time experiments using Qualtrics. *Behavior Research Methods*,
 153 47(4), 918–929. <https://doi.org/10.3758/s13428-014-0530-7>
- 154 Berinsky, A. J., Huber, G. A., & Lenz, G. S. (2012). Evaluating Online Labor Markets for Experimental
 155 Research: Amazon.com's Mechanical Turk. *Political Analysis*, 20(3), 351–368.
 156 <https://doi.org/10.1093/pan/mpr057>
- 157 Bridges, D., Pitiot, A., MacAskill, M. R., & Peirce, J. W. (2020). *The timing mega-study: Comparing a range of*
 158 *experiment generators, both lab-based and online* [Preprint]. PsyArXiv.
 159 <https://doi.org/10.31234/osf.io/d6nu5>
- 160 Cooke, M., Barker, J., Lecumberri, M. L. G., & Wasilewski, K. (2011). Crowdsourcing for word
 161 recognition in noise. *Twelfth Annual Conference of the International Speech Communication Association*.
- 162 Crump, M. J. C., McDonnell, J. V., & Gureckis, T. M. (2013). Evaluating Amazon's Mechanical Turk as
 163 a Tool for Experimental Behavioral Research. *PLoS ONE*, 8(3), e57410.
 164 <https://doi.org/10.1371/journal.pone.0057410>
- 165 De Leeuw, J. R. (2015). jsPsych: A JavaScript library for creating behavioral experiments in a Web
 166 browser. *Behavior Research Methods*, 47(1), 1–12.
- 167 de Leeuw, J. R., & Motz, B. A. (2016). Psychophysics in a Web browser? Comparing response times
 168 collected with JavaScript and Psychophysics Toolbox in a visual search task. *Behavior Research*
 169 *Methods*, 48(1), 1–12. <https://doi.org/10.3758/s13428-015-0567-2>
- 170 Gibson, E., Piantadosi, S., & Fedorenko, K. (2011). Using Mechanical Turk to Obtain and Analyze
 171 English Acceptability Judgments. *Language and Linguistics Compass*, 5(8), 509–524.
 172 <https://doi.org/10.1111/j.1749-818X.2011.00295.x>

- 173 Grootswagers, T., Cichy, R. M., & Carlson, T. A. (2018). Finding decodable information that can be read
174 out in behaviour. *NeuroImage*, 179, 252–262. <https://doi.org/10.1016/j.neuroimage.2018.06.022>
- 175 Grootswagers, T., Kennedy, B. L., Most, S. B., & Carlson, T. A. (in press). Neural signatures of dynamic
176 emotion constructs in the human brain. *Neuropsychologia*.
177 <https://doi.org/10.1016/j.neuropsychologia.2017.10.016>
- 178 Grootswagers, T., Ritchie, J. B., Wardle, S. G., Heathcote, A., & Carlson, T. A. (2017). Asymmetric
179 Compression of Representational Space for Object Animacy Categorization under Degraded
180 Viewing Conditions. *Journal of Cognitive Neuroscience*, 29(12), 1995–2010.
181 https://doi.org/10.1162/jocn_a_01177
- 182 Henninger, F., Shevchenko, Y., Mertens, U., Kieslich, P. J., & Hilbig, B. E. (2019). *lab.js: A free, open, online*
183 *experiment builder*. Zenodo. <https://doi.org/10.5281/zenodo.2775942>
- 184 Lange, K., Kühn, S., & Filevich, E. (2015). "Just Another Tool for Online Studies" (JATOS): An Easy
185 Solution for Setup and Management of Web Servers Supporting Online Studies. *PLOS ONE*,
186 10(6), e0130834. <https://doi.org/10.1371/journal.pone.0130834>
- 187 Mason, W., & Suri, S. (2012). Conducting behavioral research on Amazon's Mechanical Turk. *Behavior*
188 *Research Methods*, 44(1), 1–23. <https://doi.org/10.3758/s13428-011-0124-6>
- 189 Mathôt, S., Schreij, D., & Theeuwes, J. (2012). OpenSesame: An open-source, graphical experiment
190 builder for the social sciences. *Behavior Research Methods*, 44(2), 314–324.
- 191 Palan, S., & Schitter, C. (2018). Prolific.ac—A subject pool for online experiments. *Journal of Behavioral*
192 *and Experimental Finance*, 17, 22–27. <https://doi.org/10.1016/j.jbef.2017.12.004>
- 193 Peirce, J., Gray, J. R., Simpson, S., MacAskill, M., Höchenberger, R., Sogo, H., Kastman, E., & Lindeløv,
194 J. K. (2019). PsychoPy2: Experiments in behavior made easy. *Behavior Research Methods*, 51(1), 195–
195 203. <https://doi.org/10.3758/s13428-018-01193-y>
- 196 Pronk, T., Wiers, R. W., Molenkamp, B., & Murre, J. (2019). Mental chronometry in the pocket? Timing
197 accuracy of web applications on touchscreen and keyboard devices. *Behavior Research Methods*.
198 <https://doi.org/10.3758/s13428-019-01321-2>
- 199 Reimers, S., & Stewart, N. (2015). Presentation and response timing accuracy in Adobe Flash and
200 HTML5/JavaScript Web experiments. *Behavior Research Methods*, 47(2), 309–327.
201 <https://doi.org/10.3758/s13428-014-0471-1>

- 202 Schnoebelen, T., & Kuperman, V. (2010). Using Amazon Mechanical Turk for linguistic research1.
203 *PSIHOLOGIJA*, 43(4), 441–464.
- 204 Shank, D. B. (2016). Using Crowdsourcing Websites for Sociological Research: The Case of Amazon
205 Mechanical Turk. *The American Sociologist*, 47(1), 47–55. [https://doi.org/10.1007/s12108-015-](https://doi.org/10.1007/s12108-015-9266-9)
206 9266-9
- 207 Simcox, T., & Fiez, J. A. (2014). Collecting response times using Amazon Mechanical Turk and Adobe
208 Flash. *Behavior Research Methods*, 46(1), 95–111. <https://doi.org/10.3758/s13428-013-0345-y>
- 209 Slote, J., & Strand, J. F. (2016). Conducting spoken word recognition research online: Validation and a
210 new timing method. *Behavior Research Methods*, 48(2), 553–566. [https://doi.org/10.3758/s13428-](https://doi.org/10.3758/s13428-015-0599-7)
211 015-0599-7
- 212 Zwaan, R. A., & Pecher, D. (2012). Revisiting Mental Simulation in Language Comprehension: Six
213 Replication Attempts. *PLOS ONE*, 7(12), e51382.
214 <https://doi.org/10.1371/journal.pone.0051382>

215