

# The timing mega-study: comparing a range of experiment generators, both lab-based and online

David Bridges<sup>1</sup>, Alain Pitiot<sup>2</sup>, Michael R. MacAskill<sup>3,4</sup>, Jonathan Westley Peirce<sup>1</sup>

<sup>1</sup> School of Psychology, University of Nottingham, Nottingham, UK

<sup>2</sup> Laboratory of Image and Data Analysis, Ilixa Ltd., London (UK)

<sup>3</sup> Department of Medicine, University of Otago, Christchurch, Christchurch, New Zealand

<sup>4</sup> New Zealand Brain Research Institute, Christchurch, New Zealand

Corresponding Author:

Jonathan Peirce<sup>1</sup>

School of Psychology, University of Nottingham, Nottingham, NG7 2RD, UK

Email address: jonathan.peirce@nottingham.ac.uk

## Abstract

Many researchers in the behavioral sciences depend on research software that presents stimuli, and records response times, with sub-millisecond precision. There are a large number of software packages with which to conduct these behavioural experiments and measure response times and performance of participants. Very little information is available, however, on what timing performance they achieve in practice. Here we report a wide-ranging study looking at the precision and accuracy of visual and auditory stimulus timing and response times, measured with a Black Box Toolkit. We compared a range of popular packages: PsychoPy, E-Prime®, NBS Presentation®, Psychophysics Toolbox, OpenSesame, Expyriment, Gorilla, jsPsych, Lab.js and Testable. Where possible, the packages were tested on Windows, MacOS, and Ubuntu, and in a range of browsers for the online studies, to try to identify common patterns in performance.

Among the **lab-based experiments**, Psychtoolbox, PsychoPy, Presentation and E-Prime provided the best timing, all with mean precision under 1 millisecond across the visual, audio and response measures. OpenSesame had slightly less precision across the board, but most notably in audio stimuli and Expyriment had rather poor precision. Across **operating systems**, the pattern was that precision was generally very slightly better under Ubuntu than Windows, and that Mac OS was the worst, at least for visual stimuli, for all packages.

**Online studies** did not deliver the same level of precision as lab-based systems, with slightly more variability in all measurements. That said, PsychoPy and Gorilla, broadly the best performers, were achieving very close to millisecond precision on a number of browser configurations. For response times (using a high-performance button box), most of the packages achieved precision at least under 10 ms in all browsers, with PsychoPy achieving a precision under 3.5 ms in all.

42 There was considerable variability between operating systems and browsers,  
43 especially in audio-visual synchrony which is the least precise aspect of the  
44 browser-based experiments. Nonetheless, the data indicate that online methods  
45 can be suitable for a wide range of studies, with due thought about the sources of  
46 variability that result.

47 The results, from over 110,000 trials, highlight the wide range of timing qualities  
48 that can occur even in these dedicated software packages for the task. We stress  
49 the importance of scientists making their own timing validation measurements for  
50 their own stimuli and computer configuration.

## 51 Introduction

52 Many scientists need high-precision timing of stimuli and responses in their behavioral  
53 experiments and rely on software packages to provide that precise timing. Indeed, we often hear  
54 people state that they use particular software packages because they “need sub-millisecond  
55 timing”. Yet there is a lack of information in the literature about what is actually possible to achieve  
56 with different packages and operating systems, and very few labs report testing the timing of their  
57 studies themselves.

58 Before going further, we should establish the distinction we draw between *accuracy* and *precision*.  
59 In general, precision is the more important issue for a behavioral scientist. Precision refers to the  
60 trial-to-trial variability of the measures around the true value: the jitter of the timing measurement  
61 or its “variable error”. Accuracy refers to the “constant error” of a measurement which, in timing  
62 terms, is often referred to as the “lag”, “offset” or “bias” from the true value. Accuracy issues  
63 commonly arise from hardware characteristics and represent physical limitations of the setup, like  
64 a stimulus at the bottom of the screen typically appearing several milliseconds after a stimulus at  
65 the top of the screen, due to pixels being rendered sequentially from top to bottom. If its magnitude  
66 is known, a constant offset in time (poor accuracy) can be corrected for by simply subtracting it  
67 from each measured value. Alternatively, in many studies the ultimate outcome measure is a  
68 difference between two or more conditions, and hence any constant error is cancelled out by the  
69 taking that difference. A variable error (poor precision) cannot be corrected afterwards, as by its  
70 nature, its value is not known on any given trial.

71 Here we compare the timing performance, as directly as possible, of several commonly-used  
72 behavioral science software packages, on various operating systems, and in both laboratory-  
73 based “native” systems and on studies conducted remotely via web-browsers. The aims were to  
74 a) determine the range of timing performance that we encounter across platforms, packages and  
75 stimuli; b) identify commonalities in performance issues that need to be considered; c) assess  
76 whether online systems are technically capable of achieving sufficiently good timing. We also  
77 hope the data will encourage users to test the timing performance of their own experiments  
78 directly, using measurement validation hardware. A study like this can only show what  
79 performance it is *possible* to achieve in a given setting rather than what is likely to occur in a  
80 standard experiment. Note, for instance, that we use a high-performance button box for our tests  
81 in order to minimize the timing errors caused by external factors (the keyboard). In contrast, many  
82 laboratory-based studies, and nearly all web-based studies, are run with a standard USB  
83 keyboard, which can add further latencies of 20-40 ms, variable by keyboard (Neath et al., 2011).

84 There are also very few papers measuring **timing across packages**, allowing direct comparisons  
85 to be made based on similar hardware. Although comparisons across packages are no  
86 replacement for testing the timing on the system being used, the data presented in the current

87 study do highlight a number of consistent effects that should be informative to users and might  
88 also provide an incentive for software authors to improve their packages. The only study we are  
89 aware of that compared timing across multiple software packages is that of Garaizar *et al.* (2014)  
90 and the follow-up paper (Garaizar & Vadillo, 2014) in which they corrected an initial error. That  
91 study compared DMDX (Forster & Forster, 2003), E-Prime (Psychology Software Tools,  
92 Pittsburgh, PA) and PsychoPy but did so only on Windows 7, and only measured the precision of  
93 visual stimulus duration, without considering stimulus onset times, audio-visual asynchrony or  
94 response time measurements.

95 There also remains some confusion about the quality of **timing in online studies**, which are  
96 increasingly popular. As noted by other authors (see for example Reimers & Stewart, 2015), the  
97 rise in popularity is driven by the increasing ease with which participants can be collected by  
98 recruitment tools such as Amazon Mechanical Turk (or MTurk) or Prolific Academic, and partly  
99 by improvements in web technology, especially in timing (Reimers & Stewart, 2015). To date,  
100 studies that have explicitly tested performance, using dedicated hardware to measure stimulus  
101 onset and generate responses at precisely known times, have aimed to test generic software  
102 technologies, such as the use of the JavaScript versus Flash, rather than comparing software  
103 packages that have been specifically written for the purpose of behavioural testing (such as  
104 jsPsych, Gorilla, or PsychoJS). These have shown that when used for stimulus presentation in  
105 web browsers, HTML5 has a slightly higher tendency to drop frames (Garaizar, Vadillo & López-  
106 de-Ipiña, 2014; Reimers & Stewart, 2015) than studies run in desktop (non-browser) software.  
107 Web technology is also currently improving at a dramatic rate; there have been a number of  
108 improvements since 2015 that suggest the need for newer measurements.

109 Measured reaction time errors in these online studies have been found to consist of a lag beyond  
110 the native applications of roughly 25-45 ms, depending on the system, and an inter-trial variability  
111 (standard deviation) of 5-10 ms (Neath *et al.*, 2011, studies 5 and 6; Schubert *et al.*, 2013, study  
112 2; Reimers & Stewart, 2015, study 1). These studies did not compare any of the more recent  
113 online services such as Gorilla (Anwyl-Irvine *et al.*, 2019), jsPsych (de Leeuw, 2015; de Leeuw &  
114 Motz, 2016), PsychoPy/PsychoJS (Peirce *et al.*, 2019) or Lab.js (Henninger *et al.*, 2019).

115 A very recent paper (Pronk *et al.*, 2019) and another currently in pre-print (Anwyl-Irvine *et al.*,  
116 2020) have pointed to additional encouraging results in browser-based studies using  
117 touchscreens and keyboards. Pronk *et al.* used their own show response-time lags of 50-70 ms in  
118 most configurations (133 ms in one outlier) and inter-trial variability of 5-10 ms depending on the  
119 browser. Anwyl-Irvine *et al.* (2020) report longer lags and greater variability than some of the other  
120 papers, but they analyzed their data by including the variability between browsers and operating  
121 systems, with the aim of providing the sort of data that real experiments would contain. By  
122 contrast, the aim of the current study was to isolate the performance of the software packages  
123 themselves, and also to be comparable with measurements of lab-based experiment software  
124 packages, so we opted to use a high-performance button box in all measurements, even though  
125 in web studies this would typically not be expected.

126 While some authors have consistently pointed out the need for higher precision and more testing  
127 of timing (Plant, Hammond & Turner, 2004; Plant & Turner, 2009; Plant & Quinlan, 2013; Plant,  
128 2016), other authors have questioned whether sub-millisecond timing, of responses at least, is  
129 strictly necessary. Their point is that variability in response time measurements, once several  
130 trials have been averaged, should have relatively little impact on the statistical outcomes. For  
131 instance, Brand and Bradley (2012) modelled the effect of variability in a simulation study based  
132 on the known variability of participants and technical noise. They found, with a study of 158  
133 simulated participants (in keeping with online studies), that the addition of “technical” noise to the

134 simulated variability within participants made very little difference. Previous modelling work has  
135 also shown that timing errors, or at least lags, can be partially corrected with post-hoc calculations  
136 (Ulrich & Giray, 1989) although we aren't aware of this being common practice.

137 Partly this result is due to the inherently high trial-by-trial variability in individual participants'  
138 response times, which is on the order of several tens of milliseconds, depending on factors such  
139 as attention and motivation. Consider for instance the data from Reimers & Stewart (2007) where  
140 they compared response times (from real participants rather than a robot as in the current study),  
141 measured using code written in C versus Adobe/Macromedia Flash. Participants made a binary  
142 decision as fast as they could, which yielded a mean reaction time of between 375 ms and 400  
143 ms (depending on the software setup) but the response times of individual trials had a range of  
144 over 300 ms (interquartile range of roughly 100 ms, SD of over 80 ms). That level of variability is  
145 almost certainly driven primarily by variance in human response times rather than in the software  
146 or hardware, which are generally of considerably smaller magnitude.

147 Furthermore, de Leeuw and Motz (2016) compared participants responding in a browser-based  
148 task with those in a lab-based version and found that, where there is a measurable difference in  
149 timing, this was predominantly in the form of an increased lag (decreased accuracy) but not  
150 increased variability. Given that many studies seek to measure effects based on a *difference in*  
151 *response* between conditions, it is only the variability that is usually a concern, as any constant  
152 lag is cancelled out by taking a difference. That might largely explain the findings of Miller et al.  
153 (2018) who made measurements of various standard psychology effects, both online and offline,  
154 and found essentially no discernible difference in data quality between the online and lab-based  
155 data.

156 There are some forms of study, however, where sub-millisecond precision really is essential. For  
157 example, electroencephalography event-related potentials can have components that are very  
158 brief and more consistently timed than behavioural responses. Analysis of these can be  
159 dramatically impacted by a variability of only 10 ms in the trigger pulses with respect to the  
160 stimulus, or in the measured timing of the response. Even in behavioural tasks, in the absence of  
161 large numbers of trials or participants over which to average (unlike the ample 158 simulated  
162 participants in the Brand and Bradley study), high precision may be required.

163 Here, we quantify the technological variability (precision) and lag (accuracy) using dedicated  
164 testing hardware (Black Box Toolkit; Plant, Hammond & Turner, 2004), aiming to understand what  
165 precision can be achieved in ideal circumstances for a range of software. We suspect that the  
166 vast majority of studies will not achieve this level of precision due to, for instance, using a  
167 keyboard instead of a button box, or by not accounting for the display introducing timing errors.

168 We tested the fidelity with a range of timing measures that scientists often require. We measured  
169 *stimulus duration* to test whether a 200 ms stimulus really lasts for 200 ms. We measured the  
170 *stimulus onset*, relative to a TTL pulse (although this was not possible for browser-based studies)  
171 as would be needed to tag the stimulus onset with a trigger. Third, we measured the absolute  
172 timing of an *audio onset* relative to the same TTL pulse (in lab-based studies) and the *audiovisual*  
173 *synchrony* (in both lab-based and online studies). Lastly, we measured the *reaction time* to a  
174 visual stimulus using a robotic responder (the Black Box Toolkit key actuator), eliminating the  
175 physiological variability of human responses.

## 176 **Materials & Methods**

177 We collected timing data for a range of common software packages in standard ‘native’,  
178 laboratory-based setups, for which we opted to use *PsychoPy* (v2020.1), *Psychophysics Toolbox*  
179 (v3.0.16 beta running on MATLAB R2018b), *OpenSesame* (v3.2.8 using *PsychoPy* backend  
180 (v1.85.3), *Expyriment* (v0.9.0), and NBS *Presentation* (v21 Build 006.06.19) and E-Prime  
181 (v3.0.3.8) run using E-Studio (v3.0.3.82).

182 We sought to compare the timing of those lab-based setups with several commonly-used online  
183 packages; PsychoPy/PsychoJS (v2020.1), Gorilla (Build 20190828), jsPsych (v6.0), lab.js  
184 (v2.4.4) and Testable.

185 For the lab-based applications, a trigger pulse was generated by the software at the time at which  
186 it was intended for a visual stimulus to be displayed, or for an audio stimulus to start playing. The  
187 actual time of stimulus onset or offset was then measured via a hardware detector. We could  
188 therefore test absolute visual onset timing (compared to a hardware trigger from the package),  
189 absolute auditory timing (compared to the same trigger), visual duration precision for a 200 ms  
190 stimulus, audiovisual synchrony (attempting to present the two stimuli with simultaneous onset),  
191 and the measurement error of response time to a visual stimulus. For browser-based packages,  
192 this was not possible because web scripts do not have access to parallel or USB ports and so  
193 cannot generate a trigger signal, but all other measures were collected.

194 We created the experiments in the manner that might be expected from a normal user of each  
195 package (as described by the package documentation), and therefore excluded advanced,  
196 undocumented code additions to optimize performance. For example, the PsychoPy scripts were  
197 automatically generated by the graphical Builder interface, and were not supplemented with any  
198 custom-written Code Components to optimize performance. Since we are the authors of that  
199 package, we are more knowledgeable about potential optimizations than most users and it would  
200 be inappropriate for this package to receive any advantage from that additional experience.

201 As a caveat to the general rule of creating studies exactly as a typical user would, at times we  
202 weren’t sure what “typical users” would be aware of. For instance, in Presentation and in  
203 Expyriment, achieving a stimulus duration of 200 ms is certainly possible but to do so requires  
204 setting the requested duration to just under 200 ms (say, 195 ms). The timing mechanisms of  
205 those applications appear not to take into account the time to render the final frame, such that  
206 when requesting exactly 200 ms, the stimulus will actually overshoot by 1 screen refresh (typically  
207 16.7 ms). While a naïve user *might not* take this into account, we considered it easy enough to  
208 apply that we should do so. Certainly, anyone validating the timing independently would notice  
209 the error, and be able to verify simple that the fix works in a reliable manner. Therefore on those  
210 packages, we specified the stimulus duration to be slightly less than the intended duration.

## 211 Equipment

212 **Linux and Windows** were tested on the same (dual boot) PC, with an AMD Ryzen 5 2600 6-core  
213 3850 MHz central processing unit on a B450 Tomahawk motherboard, with 16 GB of DDR4  
214 2133 MHz RAM, a Samsung 860 EVO 500 GB SATA III Solid State Disk, and a Gigabyte GeForce  
215 GTX 1050 Ti 4 GB graphics card. For **Windows** we used version 10 (10.0.18362), running the  
216 NVIDIA 417.01 graphics driver, and the REALTEK HD audio driver (6.0.1.8549). For **Linux**, we  
217 used the Ubuntu 18.04 operating system (Linux 5.0.0-31-generic), running the proprietary NVIDIA  
218 430.26 graphics driver, with Advanced Linux Sound Architecture (ALSA) audio driver (1.1.3).

219 The **Apple Macintosh** hardware was a 2019 Mac Mini 64-bit 3.2 GHz Intel Core i7 with 16 GB of  
220 DDR4 2667 MHz RAM, with an integrated Intel Ultra High Definition (UHD) 630 1536 MB graphics  
221 processing unit. Testing was done on Mac OS X 10.14.5. The built-in Core Audio drivers were  
222 used for audio output.

223 The same monitor was used for presenting stimuli throughout: an AOC 238LM00023 / I2490VXQ  
224 23.8" 60Hz LED Backlight LCD monitor with 1920 × 1080 pixel resolution and 4 ms response time  
225 (<https://eu.aoc.com/en/monitors/i2490vxq-bt/specs>). We confirmed that this model had no options  
226 to perform any 'optimizations' on the frames generated by the graphics card.

227 For online studies we used a range of browsers, as shown in Table 1.

Table 1: Browsers used for testing across the different operating systems. *Safari* and *Edge* are specific to *MacOS* and *Windows*, respectively. *Edge* was tested in 2 versions because Microsoft recently (with Edge version 78) changed the underlying engine to use *Chromium* (the open source engine behind the *Chrome* browser).

OS	64-bit Browsers				
	FireFox	Chrome	Safari	Edge	Edge Chromium
Mac	68.0.2	76.0.3809.1	12.1.1		
Win10	69.0.0	77.0.3865		44.18362.387.0	78.0.276.19
Linux	69.0.2	76.0.3809.1			

## 228 Measurement hardware

229 We used a Black Box Toolkit v2 (BBTK) to measure the onset and offset of trigger pulses, audio  
230 and visual stimuli. We also used it to trigger responses to the visual stimuli to test the response  
231 time measurements made by the software packages. Although that can be done by the BBTK all  
232 at once, in its Digital Stimulus Capture and Response (DSCAR) mode, that limited the number of  
233 trials we could include in a single run. We wanted to run 1000 trials continuously and therefore  
234 opted to run the trials once to collect the trigger, visual and auditory onset/offsets, using BBTK's  
235 Digital Stimulus Capture (DSC) mode, and then a second time using the response actuator to test  
236 the response timing of the software, in Digital Stimulus Response Echo (DSRE) mode.

237 Trigger (TTL) pulses were sent from all test systems using a LabHackers USB2TTL8 connected  
238 to the BBTK's TTL 25-way ASC/TTL breakout board. A BBTK opto (photodiode), positioned at  
239 the center top of the display, was used to provide information about the visual stimulus. Audio  
240 onsets were recorded from the 3.5 mm speaker jack on the back of the computer.

## 241 Response time to visual stimuli timing

242 Responses to visual stimuli were created using the BBTK's robotic response key actuator (RKA).  
243 The RKA was configured using the BBTK's TTL 25-way ASC/TTL breakout board and the BBTK  
244 software RKA calibration tools in order to determine the onset and duration times of the RKA  
245 device. To achieve the desired onset times for the RKA, accounting for its solenoid response  
246 times, a 16 ms offset was taken from the intended response times.

247 The response actuator was positioned over button 1 of a LabHackers MilliKey, a 1 kHz USB  
248 response box that was used to collect responses on all platforms. Note that this is likely to provide

249 a more precise measurement than in many lab scenarios, where standard consumer-grade  
250 computer keyboards are still commonly used. Only standard keyboards or touchscreens are used  
251 in nearly all online studies, but for comparison purposes, we considered it useful to measure a  
252 consistent high-precision response across all platforms.

## 253 Procedure

### 254 Response time latencies

255 For the response time measurements, we created an experiment in each package that simply  
256 presented a black screen for 300 ms, followed by a white square (positioned at the center top of  
257 the screen) for 200 ms. The experiment was programmed to measure the response time of the  
258 actuator, which was programmed through the BBTK to respond precisely 100 ms after the onset  
259 of the white square, with a keypress 50 ms in duration. This trial sequence was repeated 1000  
260 times in quick succession, following an initial pause of 5 s to give time for the BBTK to initialize in  
261 its DSRE mode.

### 262 Stimulus latencies

263 To measure the absolute and relative latencies of the visual and auditory stimuli, we programmed  
264 an almost identical task that would present a similar black screen for 300 ms, followed  
265 simultaneously by the onset of a TTL pulse sent via the LabHackers USB2TTL8 trigger box, a  
266 white square at the top of the screen, and a simple audio tone, all lasting 200 ms. This simple trial  
267 sequence was again repeated 1000 times for each package, following a 10 s initial blank screen  
268 while the BBTK initialized into DSC (Digital Stimulus Capture) mode. In some instances using  
269 online software, it was necessary to present 1 trial of the auditory stimuli before main trials, in  
270 order to initialise the audio drivers and eliminate start-up delay in audio presentation for the first  
271 trial. If required, it is reported below for the relevant software.

### 272 PsychoPy implementation

273 The aim was to mimic what relatively naïve users would normally do. To this end, the experiment  
274 in PsychoPy was created entirely in the Builder interface, except for Code Components used  
275 solely to automatically detect the LabHackers USB2TTL8 trigger box, set the status of the TTL  
276 object (e.g., started, stopped) and write triggers to the serial port in synchrony with the visual and  
277 audio stimuli presentation. The triggers were synchronised with the screen refresh using the  
278 PsychoPy Window method, *callOnFlip()*, which allows a call to be scheduled to run at the time of  
279 the next screen refresh, rather than immediately.

280 In the Experiment Settings, the audio library was set to be PTB (i.e. Psychtoolbox's  
281 PsychPortAudio engine, ported to Python) with audio latency mode set to Level 3 ("Aggressive  
282 low-latency"). The sound waveform was generated by PsychoPy (i.e. an 'A' tone was requested  
283 in the Sound Component settings, rather than a 'wav' file being loaded). On MacOS, a lower audio  
284 latency mode (level 1: "Share low latency driver") was required to achieve clean sound (i.e.,  
285 without crackling) on this Mac-mini, although that has not been the case on other Mac hardware  
286 that we have tested.

287 Components of PsychoPy experiments, in this case the visual and auditory stimuli and trigger  
288 pulses, can be run simultaneously simply by setting them to start and stop at the same time  
289 (whereas some of the packages only allow stimuli to be displayed sequentially). PsychoPy also

290 has a check box for non-visual components to determine whether they should be synchronized  
291 with the visual stimulus (i.e. starting and stopping at the same time as the screen refresh, rather  
292 than as soon as possible) and this was set to be *on* for the Keyboard and Sound Components.

## 293 Psychtoolbox implementation

294 For Psychtoolbox there are multiple techniques one *might* use to control the stimulus timing and  
295 sequencing and so this is probably the package with greatest scope for users to get different  
296 timing than that described here.

297 In our implementation we timed the visual stimulus by actively drawing a fixed number of frames  
298 (rather than, say, flipping a frame to the display just once and then waiting for a fixed period until  
299 the next scheduled stimulus change). Many studies use dynamic stimuli that need to be updated  
300 continuously (i.e. on every screen refresh interval). That requires this active drawing and timing  
301 mechanism rather, than a flip-once-and-wait method, and this was also a close match to the  
302 method used in the PsychoPy script. The trigger was sent by calling `fprintf(usb2ttl, 'WRITE`  
303 `255\n');` immediately after the flip of the first frame.

304 The sound stimulus was queued up before the first flip of the visual stimulus, by determining the  
305 time of the next screen flip using the `PredictVisualOnsetForTime()` function. This was then  
306 used as the *when* argument for `PsychPortAudio('Start',...)`. As with the PsychoPy  
307 implementation, the sound library was set to “aggressive low-latency” audio mode (level 3).  
308 Further, to play the sound synchronously with the visual stimulus rendering loop, we set the values  
309 of `waitForEndOfPlayback` and `blockUntilStopped` to be zero (off).

## 310 NBS Presentation implementation

311 Presentation is designed for sequential presentation of visual stimuli, but does allow for parallel  
312 presentation of audio and visual stimuli.

313 In the Port menu, the output port was given an “Init sequence” of “`WRITE 0\n`”, a “Code  
314 sequence” of “`WRITE 255 200000 0\n`” and an “End sequence” of “`WRITE 0\n`”. The “Init” and  
315 “End” sequences are called at the beginning and end of the task. The “Code” sequence is called  
316 every time a “code” parameter is specified in the task script. This “Code” string sent the instruction  
317 to the LabHackers device to set the TTL pulse ON for 200 ms, and set it OFF (zero) at the end of  
318 this sequence. Port device properties for the USB2TTL8 interface were set using Rate (155200),  
319 Parity (Even) Data bits (8), Stop Bits (1), clear-to-send, data-set-ready out/In set to ON, and data-  
320 terminal-ready and request-to-send were set to “enabled”. Also, the “FIFO Interrupt” checkbox  
321 was deselected.

322 In the Response menu, we added a keyboard device with the “1” button activated (the button  
323 used on the LabHackers Millkey response box). In the Video menu, the primary display driver  
324 was selected. In the Audio menu, we used the primary sound driver. The Presentation Mixer  
325 Settings were set to the low latency “exclusive” mode, according to NBS Presentation Audio Mixer  
326 Recommendations, with duplicate channels on load selected.

327 Both the response and stimuli timing tasks were coded in the Presentation script editor. A blank  
328 screen was generated using a blank text object, positioned in a picture object. To generate the  
329 visual stimulus we used a polygon graphic object, which defined a white 400 × 400 pixel rectangle,  
330 which was positioned at top and center of the screen, using a picture object. For the stimulus  
331 timing, an audio stimulus was created using a wavefile object, used to load a 200 ms long 440 Hz



332 wav file, with the preload parameter set to true. The wavfile object was added to a sound object,  
333 ready for presentation.

334 For both tasks, the trial timeline was generated using the trial object. In the trial object, we used  
335 the *trial\_duration* variable to set the trial duration to 500 ms. For the response task, setting  
336 *trial\_type* as “*first\_response*” ended the trial on the first recorded response. The  
337 stimulus\_event objects were used to present each of the following events. A blank screen starting  
338 from time zero for a duration of 300 ms, followed by a white stimulus, starting at 300 ms from  
339 zero, for a duration of 200 ms. In the stimulus timing task, we also presented an audio stimulus  
340 at 300 ms from zero, for a duration of 200 ms. For corrected onsets and durations, see the PCL  
341 code explanation below. The TTL trigger was added to the visual stimulus event only. To achieve  
342 parallel audio visual stimulus presentation, we followed the Parallel Stimulus Events guidelines  
343 on the NBS Presentation website. This only required that we set the parallel parameter in the  
344 audio stimulus event to true.

345 PCL code was used to define the trial presentation, where the stimulus events were presented  
346 for 1000 trials. Both tasks started and finished with a black screen for 1000 ms. Note, we used  
347 Presentations black “ready” screen (see Settings tab) to provide the Black Box Toolkit initialization  
348 time. The offset of the blank screen and visual stimulus, as well as the onset of the visual stimulus,  
349 were corrected so that blank offset and visual stim onset duration was shortened by half a screen  
350 refresh (i.e.,  $200 - \text{screen refresh} / 2$ ). Also, the visual stimulus onset began half a screen refresh  
351 before its desired onset of 300 ms, and thus started before the onset of the sound was scheduled  
352 (i.e.,  $300 - \text{screen refresh} / 2$ ).

## 353 E-Prime implementation

354 E-Prime is also inherently a sequential stimulus presenter but can achieve simultaneous audio-  
355 visual stimuli using the Slide object. E-Prime (version 3.0.3.80) is not compatible with version  
356 1903 of Windows 10, as used in this study, causing E-Prime to report the “Display is too busy”  
357 runtime error or freeze. To work around this error, the full-screen mode had to be disabled in  
358 favour of windowed mode, as recommended on the E-Prime website. This potentially would have  
359 reduced timing performance although the visual stimulus timing showed no signs of dropped  
360 frames or lags, so it appears that the computer was sufficiently powerful that this did not have any  
361 effect.

362 Both tasks (stimulus measurement and response measurement) used the default experiment  
363 settings, with the exception of the Devices settings, where we set the Display to a specific refresh  
364 rate of 60, giving minimum acceptable refresh rate of 59, and a maximum acceptable refresh rate  
365 of 61.

366 The overall layout of both tasks was controlled using a main Procedure object, which started and  
367 ended with an Inline script object for setting the TTL trigger to its OFF state, and a black  
368 TextDisplay screen that ended on a keypress – useful to await BBTK initialisation. No Inline code  
369 was required for the response timing task. The LabHackers USB2TTL8 was set up as a serial  
370 device in the experiment properties, where information can be sent to the serial port using the  
371 Serial object e.g., *Serial.WriteString* “WRITE 0\n” to set an OFF signal at the start of every  
372 trial, and *Serial.WriteString* “WRITE 255 200000 \n” to set an ON signal for 200 ms,  
373 simultaneously with the stimulus. For the main trials, we added a List to the main procedure. The  
374 List acted as the main loop, where we created 1000 samples (trials) by setting 100 cycles of 10  
375 samples per cycle, using a sequential selection. To the main loop, we added another Procedure  
376 object for setting the trial timeline.

377 For the stimulus timing task, the trial procedure began with an empty Slide object, set to a black  
378 background. We used default Duration/Input properties of the Slide object. The duration of the  
379 Slide was set to 300 ms and PreRelease was “same as duration”, where this PreRelease setting  
380 allows E-Prime to preload the following stimulus immediately, during the presentation of the  
381 current stimulus.

382 For the stimulus presentation, we used another Slide object. The Slide was given a white  
383 background, and a SoundOut component, playing a 440 Hz wav file. We used default  
384 Duration/Input properties of the Slide object, where duration was 200 ms and PreRelease was  
385 “same as duration”. The PreRelease setting allowed immediate processing of the TTL code,  
386 positioned after the stimulus Slide object in the trial procedure.

387 For the response timing task, we needed only two Slide objects. The trial was the same as the  
388 stimulus timing task, without the sound component in the second Slide object, and without the  
389 InLine script for setting the TTL trigger. In addition the stimulus Slide object had additional settings  
390 in the Device/Input properties. Specifically, a keyboard was added as a device, with any key  
391 allowable, a time limit the same as the duration of the stimulus, and an End Action of “Terminate”,  
392 to end the trial on a keypress.

### 393 Expyriment implementation

394 Expyriment is also structured around sequential presentation. Each stimulus was preloaded  
395 during the ISI to prepare the stimuli for fast presentation, then presented using the flip-then-wait  
396 method. We shortened the requested visual stimulus duration to 195 ms, which reliably achieved  
397 an actual duration of 200 ms. For both visual and audio stimuli, the *present()* method was called  
398 using the default values, which took approximately 1 ms to process, according to the returned  
399 time from the *present()* method. On each trial, the TTL pulse was fired immediately after the initial  
400 call to present the visual stimulus on screen.

### 401 OpenSesame implementation

402 OpenSesame is also structured around sequential presentation. In order to present the audio  
403 stimulus synchronously with the visual stimulus and TTL pulse we tested a number of ordering  
404 combinations. The best timing was achieved by a configuration in which the audio stimulus was  
405 presented first, followed by the visual stimulus, both with a notional duration of zero, such that the  
406 next object began immediately. These were followed by the TTL pulse, and then a call to *sleep*  
407 for the duration of the stimulus, so that the stimuli remained on screen for the correct duration.  
408 This is essentially a flip-then-wait method of stimulus presentation. We shortened the requested  
409 *sleep* duration to 191 ms, to achieve an actual duration of 200 ms.

410 The TTL pulse and calls to *sleep* were coded using inline script, an OpenSesame component for  
411 inserting custom code into the experiment. The LabHackers USB2TTL8 was connected using the  
412 Python *pySerial* module, whereas *sleep* refers to a method of the OpenSesame experiment  
413 class, e.g., *self.sleep(191)*.

### 414 PsychoPy online (PsychoJS) implementation

415 The same study as used in the lab-based implementation was used to generate the PsychoJS  
416 script, which was then pushed to Pavlovia.org for running (all of which is done automatically by  
417 PsychoPy).

418 Parts of the code that were not needed for, or compatible with, the online version of the study,  
419 such as the connection to the hardware triggers, are automatically skipped by PsychoPy Builder  
420 during JavaScript code generation. No further customizations to the experiment were required for  
421 the study to run in the browsers.

422 PsychoJS uses WebGL where possible (unlike most of the other JavaScript packages, as far as  
423 we know). In just one configuration that we tested - Firefox on Linux – WebGL was supported but  
424 needed to be explicitly enabled on the PC we used. This is because Mozilla blacklists certain  
425 GPUs based on driver numbers to ensure that WebGL does not crash the browser if runs on  
426 insufficient hardware. To turn off this blacklisting we opened the settings of Firefox and set  
427 *Layers.acceleration.force-enabled* to *true*. This ensured WebGL compatibility of FireFox  
428 on Linux. Until we had made this adjustment, a warning message was provided that prevented  
429 the experiment from starting.

### 430 Gorilla implementation

431 We created a Gorilla Project, with separate Tasks for stimulus timing tests and response timing  
432 tests. Each Task was added to its own Experiment, where Task nodes were positioned between  
433 the start and finish nodes. For both tasks, we created an initial and final blank black screen, each  
434 5000 ms duration. For stimulus timing, we added an audio tone to the start screen in order  
435 initialize the audio drivers, ready for the task, followed by the main trials presenting a blank screen  
436 for 300 ms, followed by a stimulus screen for 200 ms, containing separate Zones for synchronous  
437 image and audio content presentation. For response timing, the main trials consisted of a 300 ms  
438 black screen, followed by a stimulus screen for 200 ms, containing separate Zones for image  
439 content presentation and keyboard responses. The task was run using the “Preview Task” option,  
440 used for piloting the task.

### 441 jsPsych implementation

442 The jsPsych task was coded using pure JavaScript, which consists of creating a *timeline* (array)  
443 of elements using jsPsych plugins (i.e., JavaScript objects with jsPsych compatible parameters  
444 used for presenting stimuli, recording responses etc.) and passing the *timeline* array as a  
445 parameter to the *init* method of the jsPsych experiment object. The stimuli used in both tasks  
446 were requested to be preloaded via the *preload* parameter in the jsPsych *init* method. To code  
447 the experiments, we created the timeline array, and added the Pavlovia.org connection object to  
448 the array. Finally, the command was sent to Pavlovia to finish the task, and save any data.

449 For the response time task, we began with a welcome screen containing text, using an “html-  
450 keyboard-response” object. Pressing any key would begin the experiment. For each trial, we  
451 presented a black image to using “image-keyboard-response” for 300 ms, with no key options  
452 given, rendering the keyboard ineffective. Then, a stimulus screen presented a white stimulus on  
453 black background, using “image-keyboard-response” for 200 ms, with all keys allowed as a valid  
454 response. A response ended the trial.

455 The stimulus timing task was identical to the response time task, with the exception of the audio  
456 stimulus. On the stimulus presentation screen, the audio stimulus was presented simultaneously  
457 with the image stimulus using the “audio-keyboard-response” plugin, where the audio is passed  
458 to the stimulus parameter, and the visual stimulus is presented via the “prompt” parameter.

## 459 Testable implementation

460 We were informed that Safari cannot handle fast audio stimulus presentation via Testable  
461 (personal correspondence). Therefore we did not assess Testable for synchronous sound and  
462 visual stimuli presentation using Safari.

463 Testable experiments are created using a comma-separated values (CSV) file, or Excel  
464 spreadsheet, where each row contains stimulus presentation configurations for each trial. For  
465 both response and stimuli timing tasks, tasks started and finished with a 5 second black screen.  
466 For response times, each main trial consisted of an ISI of 300 ms, followed by the 200 ms  
467 presentation of visual stimuli, with a keyboard response defined for each trial. For stimulus timing,  
468 a start-up trial was presented, preceding the 5 second start screen, containing an ISI of 300 ms,  
469 followed by the 200 ms presentation of audio and visual stimuli, in order to initialise the audio  
470 drivers, ready for the task. Each main trial consisted of an ISI of 300 ms, followed by the 200 ms  
471 presentation of audio and visual stimuli.

## 472 Lab.js implementation

473 The Lab.js task was built using the lab.js Builder. Both stimulus and response timing tasks had  
474 the same structure, and only differed with the addition of a sound oscillator for stimulus timing.  
475 The task started and ended with a black Canvas-based display, presented for durations of 5  
476 seconds. We then added a Loop and set the sample to 1000, where the sample denotes the  
477 number of loop iterations that will occur. To the Loop, we added a Frame component, which acts  
478 as the container for the stimuli. Frames contain the area occupied by the stimuli and only update  
479 that contained area on each screen refresh, thus potentially enhancing performance. This was  
480 advantageous with our stimuli since they occupied only a small portion of the screen. To the  
481 Frame, we added a Sequence, containing the trial events. To the Sequence, we added a black  
482 Canvas-based display presented for 300 ms for the ISI, and another Canvas-based display shown  
483 for 200 ms to present the stimuli. The stimulus canvas contained a white rectangle at the top-  
484 center location of the screen. The Behaviour of the stimulus canvas was set to have a timeout, or  
485 duration, of 200 ms, but was also able to record a key-down event, if required. For the stimulus  
486 timing task, sound was created using an Oscillator, added to the Behaviour timeline, with an onset  
487 at time 0, relative to the onset of the stimulus canvas, and a duration of 200 ms.

## 488 Results

489 Note that we have deliberately not included significance testing on any of the measures  
490 presented below. Such tests would give a false impression of the results. The reason to provide  
491 tests would be to give a sense for the reader of whether this would generalize to their own  
492 hardware and environment but that is not something we can address. We have tested a large  
493 number of trials on a single machine and the variance we measured in that single machine is  
494 likely to be small compared with the variance between machines. We would therefore  
495 dramatically overestimate the significance of the differences with reference to the reader's own  
496 configuration.

497 To create the tables below we have calculated a mean precision score for each row (each  
498 combination of package, operating system and browser where appropriate) and sorted the table  
499 according to that mean precision.

**Table 2: Timing summaries of desktop software by package and platform.** The Var(iability) measures are the inter-trial standard deviations of the various latencies for that configuration. The table is sorted by the mean of those variabilities (Mean Var). The Lag/Bias measures are the mean latencies for that configuration. In the case of audiovisual synchrony, a negative bias indicates the audio lead the visual stimulus, a positive bias means the visual lead the audio. Each of the values with a hyperlink will lead to a plot of the distribution of values leading to that summary value. An interactive version of the table can be found at <https://psychopy.org/timing/2020/table2.html>

Package	Platform	Mean Var (ms)	Reaction times			Visual durations		Visual onset		Audio onset		Audiovisual sync	
			Var (ms)	Lag (ms)	Var (ms)	Bias (ms)	Var (ms)	Lag (ms)	Var (ms)	Lag (ms)	Var (ms)	Bias (ms)	
PsychToolBox	Ubuntu	0.18	<a href="#">0.31</a>	<a href="#">12.3</a>	<a href="#">0.15</a>	<a href="#">2.05</a>	0.18	4.53	<a href="#">0.17</a>	-0.74	<a href="#">0.11</a>	<a href="#">-5.27</a>	
Presentation	Win10	0.29	<a href="#">0.35</a>	<a href="#">11.48</a>	<a href="#">0.23</a>	<a href="#">-1.83</a>	0.34	7.07	<a href="#">0.31</a>	0.56	<a href="#">0.19</a>	<a href="#">-6.51</a>	
PsychToolBox	MacOS	0.39	<a href="#">0.44</a>	<a href="#">22.27</a>	<a href="#">0.12</a>	<a href="#">-2.15</a>	0.41	21.52	<a href="#">0.53</a>	0.09	<a href="#">0.43</a>	<a href="#">-21.43</a>	
PsychoPy	Ubuntu	0.46	<a href="#">0.31</a>	<a href="#">8.43</a>	<a href="#">1.19</a>	<a href="#">3.49</a>	0.34	4.71	<a href="#">0.31</a>	-0.71	<a href="#">0.16</a>	<a href="#">-5.43</a>	
E-Prime	Win10	0.57	<a href="#">0.53</a>	<a href="#">9.27</a>	<a href="#">0.18</a>	<a href="#">2.51</a>	0.18	4.41	<a href="#">0.98</a>	5.08	<a href="#">0.97</a>	<a href="#">0.67</a>	
PsychToolBox	Win10	0.67	<a href="#">0.42</a>	<a href="#">10.49</a>	<a href="#">0.75</a>	<a href="#">2.24</a>	0.19	4.56	<a href="#">0.99</a>	0.77	<a href="#">0.98</a>	<a href="#">-3.79</a>	
PsychoPy	Win10	1	<a href="#">0.35</a>	<a href="#">12.05</a>	<a href="#">2.42</a>	<a href="#">-1.97</a>	0.35	7.1	<a href="#">0.96</a>	0.85	<a href="#">0.93</a>	<a href="#">-6.25</a>	
PsychoPy	MacOS	2.75	<a href="#">0.4</a>	<a href="#">22.02</a>	<a href="#">11.56</a>	<a href="#">1</a>	0.55	18.24	<a href="#">0.7</a>	0.54	<a href="#">0.52</a>	<a href="#">-17.7</a>	
Open Sesame	MacOS	<b>3.14</b>	<a href="#">0.54</a>	<a href="#">21.21</a>	<a href="#">1.65</a>	<a href="#">18.94</a>	<b>0.79</b>	18.10	<a href="#">6.40</a>	9.46	<a href="#">6.30</a>	<a href="#">-8.64</a>	
Open Sesame	Ubuntu	<b>3.41</b>	<a href="#">0.45</a>	<a href="#">9.68</a>	<a href="#">9.16</a>	<a href="#">32.29</a>	<b>0.50</b>	2.35	<a href="#">3.45</a>	2.05	<a href="#">3.48</a>	<a href="#">-0.30</a>	
Open Sesame	Win10	<b>4.02</b>	<a href="#">1.22</a>	<a href="#">8.27</a>	<a href="#">1.12</a>	<a href="#">17.04</a>	<b>0.72</b>	3.85	<a href="#">8.56</a>	47.24	<a href="#">8.50</a>	<a href="#">43.39</a>	
Expyriment	Win10	<b>6.22</b>	<a href="#">2.90</a>	<a href="#">10.76</a>	<a href="#">0.55</a>	<a href="#">-0.08</a>	<b>0.19</b>	5.98	<a href="#">13.72</a>	106.83	<a href="#">13.72</a>	<a href="#">100.85</a>	
Expyriment	Ubuntu	<b>7.75</b>	<a href="#">2.73</a>	<a href="#">23.45</a>	<a href="#">8.31</a>	<a href="#">12.08</a>	<b>0.73</b>	16.75	<a href="#">13.49</a>	118.67	<a href="#">13.50</a>	<a href="#">101.92</a>	
Expyriment	MacOS	<b>9.05</b>	<a href="#">4.84</a>	<a href="#">33.83</a>	<a href="#">7.04</a>	<a href="#">-1.13</a>	<b>4.82</b>	29.02	<a href="#">13.84</a>	42.81	<a href="#">14.72</a>	<a href="#">13.79</a>	

## 500 Lab-based package results

501 Table 2 shows the timing performance of packages running lab-based studies (not via a web  
502 browser). Timing on the lab-based systems was generally impressive. Most of the packages  
503 tested were capable of sub-millisecond precision in the visual, audio and response timing tests  
504 used here. For MacOS, performance was less precise for visual presentations. This is due mostly  
505 to a known issue introduced in version 10.13 of MacOS whereby a delay of 1 frame is imposed  
506 when updating the display. For most of the packages, this caused a relatively constant delay, and  
507 did lead to reduced *precision*. The resulting lag did, however, have a knock-on effect for other  
508 measurements, such as the visual response time measurement. On Windows and Linux, the  
509 Windows Desktop Manager and the Linux Compositor respectively also have the potential to  
510 introduce presentation delays, but did not have an effect upon the data collected in this study.

511 **PsychoPy** performed well on all the timing tests under Windows 10 and Linux. Since version 3.2,  
512 PsychoPy has used the same engine as Psychtoolbox (ported to Python by Mario Kleiner),  
513 enabling excellent audio and response timing. It should be noted that earlier versions of the  
514 software did not attain this level of performance, so upgrading to PsychoPy 3.2+ is strongly  
515 recommended. As with the other packages, performance on macOS was poorer. In PsychoPy's  
516 case, however, on this platform there appeared to be a reduced *precision* of visual stimulus  
517 presentation durations as well as the greater lag, which was not observed for the other packages.

518 **Psychophysics toolbox (PTB)** performance was excellent, at least on Linux and Windows.  
519 Achieving this precision does require more knowledge than when using PsychoPy's  
520 automatically-generated scripts. That is, there are many ways to get poorer performance  
521 unwittingly but, when programmed well, PTB can deliver excellent timing.

522 **E-Prime** performed very well out-of-the box, with no tweaking or effort. The audio stimulus had a  
523 slight (5 ms) lag compared to some of the other packages, but that is something that could  
524 presumably have been measured and corrected for, as was done for Presentation. Critically, the  
525 inter-trial variability (standard deviation) of the timings was sub-millisecond on every measure.

526 **NBS Presentation** timing was ultimately excellent, but this was not the case in the first instance.  
527 Initially we found duration measurements that overshot the desired 200 ms (which were corrected  
528 by requesting a duration of  $\frac{1}{2}$  frame less than the desired duration). We also initially found audio  
529 latencies to be both delayed and variable, having simply set the audio stimulus to play immediately  
530 after the visual stimulus. Detailed Presentation technical documentation on "Parallel Stimulus  
531 Events" describe a work-around that did allow the sound and visual stimulus to be prescheduled,  
532 if the user knows the latency that needs to be compensated for. Applying this compensation  
533 enabled the excellent timing shown in Table 2. To achieve this was rather more difficult than on  
534 other platforms, however, requiring familiarity with advanced documentation that many users will  
535 not have read.

536 **OpenSesame** timing performed well in the visual stimulus domain and response timing was also  
537 fairly good (an inter-trial variability of 1.16 ms on Windows was worse than the packages  
538 described above, but still adequate for most behavioral measurements). To get this response  
539 timing the package must constantly check the keyboard (it cannot do so asynchronously) which  
540 means that other stimulus updates can't be made at the same time (whereas PsychoPy and PTB  
541 allow checking the keyboard while presenting dynamic stimuli) but, again, this would be sufficient  
542 for many simple tasks. Audio timing was less good, with a lag of over 40 ms and an inter-trial  
543 variability of 3-9 ms, depending on operating system. This poorer performance is because, at the

544 time of writing, OpenSesame was using an older version of PsychoPy as its backend, which did  
545 not support the new PsychPortAudio library.

546 **Expyriment** had the worst performance in nearly all domains. Indeed in many instances it was  
547 out-performed by the packages that were running experiments online. Expyriment's stimulus  
548 presentation and response monitoring is built upon the Pygame Python library, which has not  
549 been optimized for low-latency, high-precision timing. We would not recommend the use of this  
550 package where precise stimulus/response timing is required.

## 551 Web-based package results

552 It is important to remember that the “online” response timing was measured as if the user had  
553 access to a low-latency button box. That is, we were using lab-based hardware from within a  
554 browser environment, which will not reflect the heterogeneous and generally low-spec commodity  
555 keyboards that will be used “in the wild” for online studies.

556 Table 3 shows the performance of packages in browser-based studies. Although the timing of the  
557 packages in online experiments did not match that of the lab-based packages, it was perhaps  
558 surprisingly good. The data were more mixed in terms of which packages performed the  
559 strongest, with some packages performing well on some browsers and poorly on others. Similarly,  
560 there was no clear winner in terms of operating system – Linux often performed poorly in these  
561 tests whereas it had generally been superior in the lab-based studies.

562 **PsychoPy/PsychoJS** version 2020.1 achieved an inter-trial variability under 5 ms in nearly all  
563 browsers for nearly all measures and often exceeded sub-millisecond precision. It should be  
564 noted that substantial timing improvements were made to this package in the 2020.1 release so  
565 users needing precise timing in their web experiments are strongly encouraged to upgrade and  
566 re-compile their JavaScript outputs from Builder.

567 PsychoPy had the lowest inter-trial variability in reaction times (under 4 ms on every browser/OS  
568 combination) with a mere 0.2 ms inter-trial variability for Ubuntu Chrome. Interestingly, the  
569 response time measure showed more *lag* under PsychoPy than some of the other packages, but  
570 better *precision*. We suspect that is due to PsychoPy using WebGL where available. That could  
571 well be introducing a 1-frame lag as the window is rendered, but then increases the certainty of  
572 when the rendering occurs. As discussed in this article, we consider constant lags of lesser  
573 importance than variability, so this may be an acceptable compromise, but we will certainly be  
574 trying to find ways using JavaScript of getting low lags at the same time as low variability.

575 **Gorilla** performed relatively well in the visual tests, with consistently low variability across the  
576 browsers and operating systems. Similarly, it performed well with visual reaction times with under  
577 6 ms inter-trial variability in all browsers and sub-millisecond in Chrome on Ubuntu. Where Gorilla  
578 struggled was with audio stimuli, with inter-trial variability over 10 ms in five of the browsers tested,  
579 and lags exceeding 100 ms in three cases.

580 **Lab.js** reaction time measures showed an inter-trial variability under 9 ms, with Firefox on macOS  
581 showing sub-millisecond precision. The notable thing about lab.js was that it showed surprisingly  
582 low *lag* values for measures like reaction time but not an improved precision. Indeed, on some  
583 trials, the lag was negative: it was reported as having a shorter response time should have been  
584 possible, not something we saw in any other package or configuration.

**Table 3: Timing summaries of web-based software by package, platform, and browser.** The Var(iability) measures are the inter-trial standard deviations of the various latencies for that configuration. The table is sorted by the mean of those variabilities (Mean Var). The Lag/Bias measures are the mean latency values, for that configuration. In the case of audiovisual sync, a negative bias indicates the audio lead the visual stimulus, a positive bias means the visual lead the audio. Each of the values with a hyperlink will lead to a plot of the distribution of values leading to that summary value. An interactive version of the table can be found at <https://psychopy.org/timing/2020/table3.html>

Package	Platform	Browser	Mean Var (ms)	Reaction times		Visual durations		Audiovisual sync	
				Var (ms)	Lag (ms)	Var (ms)	Bias (ms)	Var (ms)	Bias (ms)
PsychoPy	Win10	Chrome	<b>1.36</b>	<a href="#">0.39</a>	<a href="#">43.95</a>	<a href="#">0.67</a>	<a href="#">-2.08</a>	<a href="#">3.01</a>	<a href="#">65.32</a>
Gorilla	Win10	Firefox	<b>1.84</b>	<a href="#">1.11</a>	<a href="#">24.83</a>	<a href="#">2.67</a>	<a href="#">1.35</a>	<a href="#">1.73</a>	<a href="#">88.27</a>
Gorilla	MacOS	Firefox	<b>2.18</b>	<a href="#">4.47</a>	<a href="#">30.34</a>	<a href="#">0.94</a>	<a href="#">1.16</a>	<a href="#">1.12</a>	<a href="#">38.43</a>
PsychoPy	Win10	Edge (Standard)	<b>2.22</b>	<a href="#">2.03</a>	<a href="#">42.00</a>	<a href="#">0.93</a>	<a href="#">-2.28</a>	<a href="#">3.69</a>	<a href="#">56.19</a>
PsychoPy	MacOS	Firefox	<b>2.65</b>	<a href="#">1.17</a>	<a href="#">67.01</a>	<a href="#">3.38</a>	<a href="#">0.24</a>	<a href="#">3.40</a>	<a href="#">-10.21</a>
PsychoPy	MacOS	Safari	<b>2.66</b>	<a href="#">1.05</a>	<a href="#">33.50</a>	<a href="#">4.26</a>	<a href="#">0.49</a>	nan	nan
Gorilla	Ubuntu	Firefox	<b>2.76</b>	<a href="#">4.71</a>	<a href="#">24.71</a>	<a href="#">2.35</a>	<a href="#">2.05</a>	<a href="#">1.23</a>	<a href="#">-30.61</a>
PsychoPy	Win10	Firefox	<b>2.76</b>	<a href="#">1.96</a>	<a href="#">40.97</a>	<a href="#">2.42</a>	<a href="#">-2.61</a>	<a href="#">3.90</a>	<a href="#">58.93</a>
jsPsych	MacOS	Safari	<b>3.39</b>	<a href="#">0.66</a>	<a href="#">31.31</a>	<a href="#">4.39</a>	<a href="#">3.09</a>	<a href="#">5.11</a>	<a href="#">-23.48</a>
jsPsych	Win10	Edge (Chromium)	<b>3.85</b>	<a href="#">1.74</a>	<a href="#">15.19</a>	<a href="#">4.21</a>	<a href="#">2.97</a>	<a href="#">5.60</a>	<a href="#">44.51</a>
Testable.org	Win10	Firefox	<b>3.92</b>	<a href="#">3.87</a>	<a href="#">31.36</a>	<a href="#">2.94</a>	<a href="#">1.91</a>	<a href="#">4.95</a>	<a href="#">76.32</a>
PsychoPy	Ubuntu	Firefox	<b>3.97</b>	<a href="#">1.57</a>	<a href="#">42.50</a>	<a href="#">4.97</a>	<a href="#">1.02</a>	<a href="#">5.36</a>	<a href="#">190.45</a>
Testable.org	Ubuntu	Firefox	<b>4.05</b>	<a href="#">3.97</a>	<a href="#">31.57</a>	<a href="#">3.25</a>	<a href="#">1.84</a>	<a href="#">4.92</a>	<a href="#">-44.36</a>
PsychoPy	Ubuntu	Chrome	<b>4.14</b>	<a href="#">0.2</a>	<a href="#">66.98</a>	<a href="#">1.77</a>	<a href="#">1.77</a>	<a href="#">10.45</a>	<a href="#">187.19</a>
Lab.js	MacOS	Firefox	<b>4.20</b>	<a href="#">0.97</a>	<a href="#">16.38</a>	<a href="#">8.61</a>	<a href="#">19.51</a>	<a href="#">3.01</a>	<a href="#">4.26</a>
PsychoPy	Win10	Edge (Chromium)	<b>4.24</b>	<a href="#">1.04</a>	<a href="#">46.01</a>	<a href="#">3.03</a>	<a href="#">-3.36</a>	<a href="#">8.66</a>	<a href="#">63.30</a>



jsPsych	Ubuntu	Chrome	<b>4.63</b>	<u>3.23</u>	<u>48.29</u>	<u>4.33</u>	<u>4.05</u>	<u>6.34</u>	<u>27.73</u>
PsychoPy	MacOS	Chrome	<b>4.84</b>	<u>3.22</u>	<u>35.29</u>	<u>6.47</u>	<u>-0.3</u>	<u>4.82</u>	<u>-6.86</u>
jsPsych	Ubuntu	Firefox	<b>5.12</b>	<u>4.11</u>	<u>31.38</u>	<u>4.84</u>	<u>4.05</u>	<u>6.43</u>	<u>10.55</u>
Lab.js	Ubuntu	Chrome	<b>5.12</b>	<u>8.27</u>	<u>31.79</u>	<u>1.19</u>	<u>2.34</u>	<u>5.91</u>	<u>198.05</u>
jsPsych	MacOS	Firefox	<b>5.16</b>	<u>6.85</u>	<u>53.70</u>	<u>2.57</u>	<u>3.27</u>	<u>6.05</u>	<u>-15.29</u>
Testable.org	Ubuntu	Chrome	<b>5.46</b>	<u>4.23</u>	<u>47.09</u>	<u>7.06</u>	<u>-11.92</u>	<u>5.09</u>	<u>94.18</u>
Testable.org	MacOS	Chrome	<b>5.52</b>	<u>3.99</u>	<u>41.94</u>	<u>7.43</u>	<u>6.34</u>	<u>5.13</u>	<u>67.2</u>
jsPsych	MacOS	Chrome	<b>5.62</b>	<u>5.68</u>	<u>41.42</u>	<u>5.74</u>	<u>3.25</u>	<u>5.45</u>	<u>-9.71</u>
Lab.js	Win10	Firefox	<b>5.78</b>	<u>7.88</u>	<u>8.22</u>	<u>4.21</u>	<u>14.25</u>	<u>5.25</u>	<u>70.93</u>
Lab.js	MacOS	Chrome	<b>5.79</b>	<u>3.26</u>	<u>20.3</u>	<u>8.78</u>	<u>6.51</u>	<u>5.31</u>	<u>-0.21</u>
Gorilla	Win10	Edge (Standard)	<b>5.96</b>	<u>4.95</u>	<u>40.23</u>	<u>7.15</u>	<u>4.95</u>	<u>5.79</u>	<u>70.93</u>
Testable.org	MacOS	Firefox	<b>6.01</b>	<u>5.20</u>	<u>57.66</u>	<u>7.63</u>	<u>22.83</u>	<u>5.48</u>	<u>40.73</u>
Lab.js	Ubuntu	Firefox	<b>6.19</b>	<u>2.91</u>	<u>25.54</u>	<u>10.22</u>	<u>13.66</u>	<u>5.44</u>	<u>185.35</u>
jsPsych	Win10	Chrome	<b>6.23</b>	<u>7.85</u>	<u>23.27</u>	<u>5.1</u>	<u>3.60</u>	<u>5.73</u>	<u>43.57</u>
Testable.org	Win10	Edge (Chromium)	<b>6.80</b>	<u>4.11</u>	<u>15.99</u>	<u>8.34</u>	<u>-5.39</u>	<u>7.94</u>	<u>73.79</u>
jsPsych	Win10	Firefox	<b>7.38</b>	<u>8.37</u>	<u>25.70</u>	<u>7.04</u>	<u>15.32</u>	<u>6.74</u>	<u>32.32</u>
Gorilla	Win10	Chrome	<b>7.89</b>	<u>3.58</u>	<u>25.60</u>	<u>5.03</u>	<u>4.24</u>	<u>15.06</u>	<u>98.84</u>
Testable.org	Win10	Chrome	<b>8.08</b>	<u>7.88</u>	<u>23.96</u>	<u>8.38</u>	<u>-5.90</u>	<u>7.98</u>	<u>72.57</u>
Lab.js	Win10	Edge (Chromium)	<b>8.57</b>	<u>4.22</u>	<u>17.14</u>	<u>8.03</u>	<u>-4.90</u>	<u>13.45</u>	<u>82.45</u>
Lab.js	Win10	Chrome	<b>9.48</b>	<u>8.44</u>	<u>19.20</u>	<u>7.76</u>	<u>-1.73</u>	<u>12.25</u>	<u>86.03</u>
Gorilla	MacOS	Chrome	<b>9.76</b>	<u>5.30</u>	<u>35.31</u>	<u>5.36</u>	<u>3.40</u>	<u>18.61</u>	<u>32.03</u>
Gorilla	Win10	Edge (Chromium)	<b>11.34</b>	<u>4.89</u>	<u>23.4</u>	<u>3.01</u>	<u>1.56</u>	<u>26.13</u>	<u>121.57</u>
Gorilla	Ubuntu	Chrome	<b>14.17</b>	<u>0.43</u>	<u>40.85</u>	<u>1.66</u>	<u>3.36</u>	<u>40.42</u>	<u>200.55</u>
Gorilla	MacOS	Safari	<b>19.16</b>	<u>1.53</u>	<u>29.65</u>	<u>30.11</u>	<u>22.25</u>	<u>25.83</u>	<u>285.81</u>

585 **jsPsych** and **Testable** both showed inter-trial variability in the range 3.2-8.4 ms in all  
586 configurations, slightly less precise than Gorilla and PsychoPy. Nonetheless that variability is  
587 still less than the typical physiological variability of human participants themselves, or the  
588 keyboards and touchscreens that are typically used for responses.

## 589 **Discussion**

590 From the data, it is clear that modern computers are capable of very precise timing in presenting  
591 audio and visual stimuli and in receiving responses, at least when used with a button box rather  
592 than a keyboard. There are a number of absolute lags that are common to all the software  
593 packages, which cannot be avoided, and there are differences between software packages but,  
594 in the best-performing packages on Windows and Linux, the inter-trial variability in the timing (the  
595 standard deviation of the measurements for a single configuration) was typically under a  
596 millisecond for all measures.

597 For lab-based studies, PsychoPy and Psychtoolbox were the most precise, outperforming even  
598 the proprietary NBS Presentation and E-Prime packages. OpenSesame and Expyriment followed  
599 in precision, in that order. For online studies the timing was less precise than in the native  
600 applications and was quite variable between browsers. Stimulus duration remained relatively  
601 precise on most of the packages. For response times, there were larger lags, and these varied  
602 between browsers, but the precision within a software/browser combination (as would be  
603 experienced by an individual participant) was relatively good, with an inter-trial variability in the  
604 range 5-10 ms in most cases and even less for PsychoJS. The findings lend support to the notion  
605 that online studies might well be sufficient for many studies, except where the utmost precision is  
606 required, or where *absolute* response times must be compared between individuals, which would  
607 be impacted by responding on different systems. Further details on the particular packages and  
608 operating systems are considered below.

609 It is very important to note that the timings measured here represent something approaching a  
610 best-case scenario. Although we deliberately tested using mid-spec, rather than high-end,  
611 computers there are various reasons that our timing may have been better than in many standard  
612 experimental set ups. We used very simple stimuli: a single white square and a sound. We  
613 ensured that our monitors were correctly configured and the graphics settings were appropriate,  
614 whereas some labs, and most online participants, will leave the monitor and graphics card in  
615 whatever their default settings are. We also used a button box to measure the response times,  
616 rather than the standard commodity-grade keyboard, in use by many labs, and nearly all online  
617 participants. Probably most importantly, however, by validating the timing independently with  
618 dedicated hardware, we could detect when timing was not as good as expected (whether because  
619 of hardware or software settings). For example, as mentioned above, we found that some of the  
620 software packages presented the visual stimulus for 1 extra frame (216.7 ms rather than the  
621 intended 200 ms) unless we reduced the requested duration to 195 ms. We suspect that a large  
622 number of studies are being conducted with timing that is considerably worse than reported here,  
623 by virtue of stimuli being incorrectly programmed, hardware being incorrectly configured, or by  
624 computers that aren't sufficient for the task. We discuss below a range of specific ways in which  
625 timing performance can be dramatically impaired. This should highlight the importance of testing  
626 timing for in every study, on its particular combination of operating system and hardware. Papers  
627 such as this one, which report "best case" timing performance, should not be used in lieu of study-  
628 specific validation and testing.

## 629 Comparing the lab-based packages

630 PsychoPy, Psychtoolbox, E-Prime and NBS Presentation all had precision that was below 1 ms  
631 on average across the measures. OpenSesame was slightly less precise across the board but  
632 most notably in audio. We were using the PsychoPy backend, but this in OpenSesame is currently  
633 using an older version of PsychoPy not supporting the new low-latency audio options. Expyriment  
634 had the worst timing and would not be recommended, particularly for studies needing precisely-  
635 timed audio stimuli. The exception in timing quality for all the packages was MacOS, where no  
636 package achieved good visual timing, as discussed below.

637 For visual stimulus durations, most of the packages showed similar timing, although with  
638 Presentation and Expyriment, the correct duration was only achieved by setting the requested  
639 duration to a shorter time (setting to 200 ms resulted in a 1-frame overshoot, whereas setting it to  
640 191 resulted in good timing). This is the sort of issue where, unless validating the timing with a  
641 hardware device, unsuspecting users would find it very easy to produce an incorrect stimulus  
642 duration. Surprisingly in OpenSesame we set the duration to be 191 ms but still observed an  
643 overshoot.

644 For response times, PsychoPy, Psychtoolbox, E-Prime and Presentation all provided reaction  
645 times with an inter-trial variability of 0.5 ms or less on all 3 operating systems. OpenSesame  
646 reaction times were slightly poorer on Windows, with an inter-trial variability of 1.2 ms, but similar  
647 high-precision on Linux and macOS. Expyriment reaction times had an inter-trial variability of 2-  
648 5ms depending on the operating system. Due to the aforementioned visual lag on the Mac, the  
649 response times on that platform all appear to be 1 frame slower but, again, as that is roughly  
650 constant, the achieved precision is generally good.

651 The relatively poorer performance of Expyriment is likely to stem from its use of the Pygame  
652 library (a Python wrapper of SDL), which provides convenient features for programming, but sub-  
653 optimal performance.

## 654 Comparing Operating Systems

655 Mario Kleiner, on the Psychtoolbox forum, has long advocated the use of Linux for optimal timing  
656 and that is somewhat born out here. Timing was indeed nearly always better on Ubuntu than the  
657 other systems, but the difference for these particular tests was relatively small (compare for  
658 example an audio variability of roughly 0.2 ms on Linux with 0.5 ms on MacOS and 1.0 ms on  
659 Windows for both PsychoPy and Psychtoolbox). The difference may well be accentuated with  
660 tougher testing environments, such as testing whether the package still performs well under high  
661 computing loads.

662 The most notable poor performance overall was the lag of visual stimulus onset on Apple's  
663 MacOS. Attempting to sync a trigger pulse with a visual stimulus on a Mac revealed a 1-frame  
664 delay for most of the software packages, and on Expyriment the lag was longer and variable. This  
665 lag on the Mac is something that depends on the operating system version. Up to and including  
666 OS X 10.12, we could see the same high-precision visual timing as on the other operating systems  
667 but this changed in the system update from 10.13 onwards (persisting at least until MacOS  
668 10.14.5, as tested here). It appears that the system has added some additional buffering step  
669 ("triple buffering") into its rendering pipeline. Therefore when the experimental software regards  
670 the framebuffer as having 'flipped', it has actually just progressed to the next buffering stage and  
671 is not yet visible on the screen.

672 The same behavior occurs in Windows 10 if triple-buffering is enabled in the driver settings, or by  
673 turning on screen scaling (which appears to implicitly use triple-buffering). On Windows these can  
674 always be turned off if the user knows to search for them. On MacOS, however, there is currently  
675 no user-accessible way to disable this lag. In most other aspects, the Mac had good timing.

## 676 Comparing online packages

677 For the online packages, we could not measure the absolute lag as there is no means to send a  
678 hardware trigger synchronization pulse from within the browser environment. The only measures  
679 we could provide in that setting were the precision for a stimulus duration, the audio-visual onset,  
680 and the visual stimulus response. In many cases those are, in any case, the things that the  
681 scientist might need but they do little to inform us of the cause of any discrepancies. For example,  
682 when we find that the response time has a lag of 35 ms it isn't clear whether this is caused by a  
683 delay in the visual stimulus appearing, or a delay in detecting the keypress.

684 Those caveats aside, behavioral scientists might be reassured to find that most of the online  
685 packages were certainly capable of presenting visual stimuli for a relatively precise number of  
686 milliseconds. For the vast majority of packages and configurations we found an inter-trial  
687 variability of less than 5 ms for stimulus duration and rarely any consistent under/overshoots.  
688 There was more variability than in the native packages, but the effects on most experiments are  
689 likely to be very small.

690 For response timing, similarly, we were generally impressed by the performance, finding the inter-  
691 trial variability to be under 10 ms in all cases. PsychoPy/PsychoJS topped this table recording a  
692 precision of under 4 ms in every browser/OS combination, and with sub-millisecond precision  
693 using Chrome for both Windows and Linux. As noted by other authors, the absolute lags in the  
694 response times were longer than in desktop studies, but these are typically the less important  
695 measure. In a study where one takes a measure by comparing the response times of participants  
696 across multiple conditions (consider a Stroop task, for instance, or an Implicit Association Test)  
697 then this is of little consequence. As we take the difference in response times, the absolute lag is  
698 subtracted and the only value of relevance is the variability.

699 Note that the measured timings here are much higher precision than in the most recent other  
700 study to our knowledge (Anwyl-Irvine et al., 2020). For instance, they report standard deviations  
701 of over 10 ms for most packages (Table 2), but that value is the standard deviation across devices  
702 and browsers for a package, not the inter-trial variability within any single configuration as  
703 reported here. A second key difference is that they used a keyboard for all measurements, but  
704 this would not be sufficient to explain some of the extreme reaction time measures that they found:  
705 for all packages they report maximum errors of over 150 ms, which we have not encountered in  
706 any of our measurements and which should not have resulted simply from the use of a keyboard.  
707 A last key difference with their data is that, for PsychoPy, those authors used an older version  
708 (v3.1.5) than in the current study (v2020.1) and it is certainly the case that PsychoPy's timing has  
709 improved a great deal between those versions.

710 For audio-visual synchrony, the data are less encouraging; all the tested packages are currently  
711 struggling in that domain. There are some browsers where the synchrony is good, but others  
712 where it is extremely poor using the same software package. In some cases the sound failed to  
713 play reliably at all. The results indicate that JavaScript probably isn't currently a technology ready  
714 for precisely-timed audio and this is an area for all the software authors to work on.

## 715 The importance of making your own measurements

716 The comparisons made here generally represent best-case scenarios. There are several reasons  
717 that the timing would be poorer on a typical lab system than measured in the current study. The  
718 chief of these reasons is simply that we independently tested the timing with photodiodes and  
719 microphones and, in so doing, we found timing problems that could be addressed first. We  
720 consider below just some of the many factors that can cause timing quality to be reduced.

721 **Visual lags from monitor and operating systems.** The experimenter might experience lags of  
722 one or more refresh period for the visual stimulus (and that delay is effectively also added to your  
723 response times and to the audio-visual synchrony) in a manner that your stimulus presentation  
724 system cannot detect without dedicated hardware. This could be caused by your monitor itself.  
725 Monitors often now have multiple viewing “modes” such as movie mode or game mode, where  
726 the manufacturer seeks to optimize the picture for that particular activity. In optimizing the picture,  
727 they perform further processing on the pixel values that you have sent to the screen from your  
728 graphics card. As well as frustrating the careful scientist by subtly altering the images, the delay  
729 incurred can have a dramatic impact on timing, because the processing can take 10-20 ms,  
730 introducing a lag. Windows 10 can add a 1-frame lag if you turn on seemingly innocuous features  
731 like screen scaling. Unless you measure the physical timing against a hardware trigger signal,  
732 then you wouldn’t know that this was occurring. If you measure the system and establish that the  
733 timing is good but then stop measuring further, you are still susceptible to changes introduced as  
734 the operating system or the audio/graphics card drivers are updated (for example, Apple silently  
735 introducing a 1-frame lag at the level of the operating system in MacOS X 10.13). Note also that  
736 whereas some timing errors (such as a dropped frame) can be detected by the software packages  
737 themselves and will show up in log files and/or alert messages, others, such as image processing  
738 in the monitor, cannot be detected by any of the packages. Hardware tests are required to detect  
739 that if the monitor is adding delays to stimulus display.

740 **Audio timing.** Unlike the visual stimulus, for which we can at least detect when the graphics card  
741 has flipped a frame buffer (and then hope that the screen updated reasonably quickly after that),  
742 audio libraries do not report back on the progress they have made. When we request that a sound  
743 plays, we have no real information about when the speakers begin physically vibrating, except by  
744 recording it externally with a microphone. A software package might be able to use information  
745 about the size of the audio buffer and latency settings that are reported in the card to estimate  
746 when the sound is physically played, but it cannot detect any feedback signal that a sound has  
747 played. Again, the only way to know about the timing quality of your audio stimuli is to test with  
748 hardware devices.

749 **Stimulus differences.** The next reason to think that your timing performance might not match  
750 the performance described above is that, in nearly all cases, experimental stimuli will be more  
751 complex than those we presented here, and tests need to be carried out for *that stimulus*  
752 configuration to confirm that the quality of timing is still being maintained. Complexity can affect  
753 performance, and it might not be clear to many users what constitutes “complex” (for instance,  
754 rendering text stimuli is more computationally challenging for most packages than rendering a  
755 photographic image or a shape). Similarly, although we used computers with moderately high  
756 specifications rather than heavy-duty powerhouses, in reality very many experiments are run on  
757 relatively weak hardware. The characteristics of the computer and, especially its graphics card,  
758 can have a drastic effect on experimental timing. Especially in the days of high-resolution displays,  
759 many users are probably unaware of the demands these place on a graphics card and that they  
760 should only be used in conjunction with high-end dedicated graphics processors. A “4K” display,  
761 increasingly common as a computer monitor, typically has a resolution of 3840 × 2160, which is

762 roughly 8.3 million pixels, each consisting of red, green, blue and alpha values. At 60 frames per  
763 second the graphics card needs to update, and output, a staggering *40 billion values per second*.  
764 That is something that modern cards are capable of doing, but cheaper or older computers are  
765 not. When they encounter a high-resolution display they will typically just send fewer screen  
766 updates per second, so that a nominally 60 frames per second display runs at, say, 30 Hz or  
767 becomes irregular.

768 **Erroneous or inefficient code.** While graphical experiment builders do a great deal to help  
769 reduce errors in code, it is still easy for a user to make mistakes. In particular, many users have  
770 rather little knowledge of the underlying limitations of their computer, or at least haven't thought  
771 through the implications of those limitations. For example, it takes a relatively long time for an  
772 image to be loaded from disk and uploaded to the graphics card, whereas rescaling it and  
773 repositioning the same image has minimal overhead for a hardware-accelerated graphics  
774 application. At other times of course the scientist might have excellent knowledge but simply made  
775 a typographical error while creating the experiment. The best way to ensure that you haven't  
776 made a mistake with your coding is to test the physical output of that coding (i.e. by testing the  
777 stimulus appearance with hardware).

778 **Miscellaneous.** There are many additional recommendations, of course. While running your  
779 study you should turn off unnecessary network services, such as Dropbox and email applications,  
780 as well as any applications that might be using substantial amounts of memory. You should keep  
781 image files to roughly the number of pixels that they will need when rendered on-screen (rather  
782 than loading a 12-megapixel image to be displayed on a 2-megapixel screen). You should try to  
783 perform any time-consuming activities in your code during inter-trial intervals or similar. If any of  
784 these miscellaneous issues are causing timing errors, however, most of them should be  
785 detectable in your software log files, which you should also check as a part of your testing process.

## 786 How to test timing

787 There are a variety of options for this depending on your needs and budget. With an oscilloscope,  
788 a microphone, and a photodiode you can test audio-visual synchrony. If you don't have a  
789 traditional oscilloscope there are now cheap PC-based options available to help, such as the  
790 [BitScope](#), although those are more likely to require some programming (we haven't tested one).  
791 The downside of traditional oscilloscopes is that while you can easily visualize the offset between  
792 various signals such as triggers and visual stimuli, this can be difficult or impossible to automate  
793 unless your scope also supports computer communication.

794 If you want to test the absolute lag of a visual or auditory stimulus then you will also need  
795 something with which to generate a trigger pulse, such as a parallel port. Physical parallel ports  
796 are increasingly hard to source, hard to configure and may not be compatible with your computer  
797 (e.g. on a Mac or a laptop PC). There are other cheap cross-platform USB solutions though, such  
798 as the [LabHackers USB2TTL8](#) we used in this study and is extremely simple to set up and use,  
799 or a [LabJack](#), which has more channels and provides analogue as well as digital in/outputs, but  
800 requires a little more programming effort.

801 The [Labhackers Millikey](#) is a high-performance button box with the additional feature that it can  
802 be commanded via the USB port (with RS232 protocol) to send a virtual keypress and this can be  
803 used to test how quickly the software then detected the keyboard event. We have tested this and  
804 found a round-trip time of under 1 ms to send the command to the box and detect the resulting  
805 keypress. An upgrade called the [DeLux](#) allows you to fire the virtual keypress on the basis of a

806 visual stimulus onset, which you could also use to test how long it took to be detected by your  
807 software (after stimulus onset). Such a setup allows you to test the overall timing error from  
808 stimulus presentation to response collection (for instance, revealing the visual stimulus delay on  
809 recent MacOS) and it can be scripted.

810 For the ultimate range of measurements, a key device is the [Black Box Toolkit](#) (BBTK) (Plant,  
811 Hammond & Turner, 2004), as used in all the timing tests reported here. BBTK v2 has many  
812 inputs and outputs allowing you to test the timing over 4 photodiodes and 2 microphones and to  
813 send “responses” via sounders and keyboard actuators. You can use it to test the timing of stimuli,  
814 or to respond to your stimuli in numerous ways, all synchronized with trigger pulses over  
815 numerous TTL input/outputs. This is a more expensive device, costing between £1500 and £3000  
816 depending on the options, but is highly recommended as a resource for complete testing of your  
817 experimental setup.

## 818 **Conclusions**

819 We find that PsychoPy, Psychtoolbox, E-Prime® and NBS Presentation® are all capable of similar  
820 sub-millisecond precision in both stimulus and response timing. OpenSesame and Expyriment  
821 were slightly less precise, especially in terms of audio stimulus presentation. In comparing  
822 operating systems, the ideal system depends on the type of study: for visual stimuli Apple’s  
823 macOS suffers from a visual stimulus lag (since OS X version 10.13); Linux can be optimized to  
824 be extremely precise but its web browsers are seemingly poorly optimized, and Windows 10  
825 appears to have reasonable performance in all domains.

826 For online studies we report the fastest, least variable data yet measured in browser-based  
827 latency tests. All the packages we tested were also reasonably precise in visual stimulus  
828 presentation, with PsychoPy achieving particularly impressive reaction time precision of under 4  
829 ms on all browsers. That said, these packages remain not quite as precise as the lab-based  
830 equivalents. In particular, no online system can yet provide audio stimuli with precisely timed  
831 onsets. We also note, in agreement with previous studies (Reimers & Stewart, 2015), that quality  
832 of *absolute* timing (the *accuracy* rather than the *precision*) is poor in browsers. Comparisons  
833 between participants of their absolute response times rather than is unwise for web-based data.  
834 Studies should aim always to make comparisons with a control condition on the same  
835 browser/computer (such as in within-subject comparisons) so that these absolute lags are  
836 naturally removed.

## 837 **Acknowledgements**

838 We are grateful to the authors of the other software packages covered here. As well as working  
839 tirelessly, often unpaid, to create this array of packages we have had enjoyable and informative  
840 discussions with many of them. We are particularly grateful, of course, to Mario Kleiner for porting  
841 the PsychPortAudio and PsychHID libraries over to Python, work which was critical in Python  
842 users having access to the high-performance audio and keyboard timing reported here. The fact  
843 that Psychtoolbox and PsychoPy have such strong timing is largely due to his knowledge and  
844 work in the field.

## 845 **Competing interests**

846 Although we are authors of one of the packages being compared here (PsychoPy/PsychoJS),  
847 that software is provided free and open-source to all users. We are also shareholders of Open  
848 Science Tools Ltd, which receives income from studies being conducted on Pavlovia.org.

849 Pavlovia.org is agnostic to the software package creating the studies – it supports studies from  
850 PsychoPy, jsPsych and Lab.js – and all resulting income goes into further development of the  
851 open-source software packages on which the studies are created.

852



853 **References**

- 854 Anwyl-Irvine AL, Dalmaijer ES, Hodges N, Evershed J. 2020. Online Timing Accuracy and  
855 Precision: A comparison of platforms, browsers, and participant's devices. *PsyArXiv*. DOI:  
856 10.31234/osf.io/jfecfa.
- 857 Anwyl-Irvine AL, Massonnié J, Flitton A, Kirkham N, Evershed JK. 2019. Gorilla in our midst: An  
858 online behavioral experiment builder. *Behavior Research Methods*. DOI: 10.3758/s13428-  
859 019-01237-x.
- 860 Brand A, Bradley MT. 2012. Assessing the Effects of Technical Variance on the Statistical  
861 Outcomes of Web Experiments Measuring Response Times. *Social Science Computer  
862 Review* 30:350–357. DOI: 10.1177/0894439311415604.
- 863 Forster KI, Forster JC. 2003. DMDX: A Windows display program with millisecond accuracy.  
864 *Behavior Research Methods, Instruments, & Computers* 35:116–124. DOI:  
865 10.3758/BF03195503.
- 866 Garaizar P, Vadillo MA. 2014. Accuracy and Precision of Visual Stimulus Timing in PsychoPy: No  
867 Timing Errors in Standard Usage. *PLOS ONE* 9:e112033. DOI:  
868 10.1371/journal.pone.0112033.
- 869 Garaizar P, Vadillo MA, López-de-Ipiña D. 2014. Presentation Accuracy of the Web Revisited:  
870 Animation Methods in the HTML5 Era. *PLOS ONE* 9:e109812. DOI:  
871 10.1371/journal.pone.0109812.
- 872 Garaizar P, Vadillo MA, López-de-Ipiña D, Matute H. 2014. Measuring Software Timing Errors in  
873 the Presentation of Visual Stimuli in Cognitive Neuroscience Experiments. *PLOS ONE*  
874 9:e85108. DOI: 10.1371/journal.pone.0085108.
- 875 Henninger F, Shevchenko Y, Mertens UK, Kieslich PJ, Hilbig BE. 2019. lab.js: A free, open, online  
876 study builder. *PsyArXiv*. DOI: 10.31234/osf.io/fqr49.
- 877 de Leeuw JR. 2015. jsPsych: A JavaScript library for creating behavioral experiments in a Web  
878 browser. *Behavior Research Methods* 47:1–12. DOI: 10.3758/s13428-014-0458-y.
- 879 de Leeuw JR, Motz BA. 2016. Psychophysics in a Web browser? Comparing response times  
880 collected with JavaScript and Psychophysics Toolbox in a visual search task. *Behavior  
881 Research Methods* 48:1–12. DOI: 10.3758/s13428-015-0567-2.
- 882 Miller R, Schmidt K, Kirschbaum C, Enge S. 2018. Comparability, stability, and reliability of  
883 internet-based mental chronometry in domestic and laboratory settings. *Behavior  
884 Research Methods* 50:1345–1358. DOI: 10.3758/s13428-018-1036-5.
- 885 Neath I, Earle A, Hallett D, Surprenant AM. 2011. Response time accuracy in Apple Macintosh  
886 computers. *Behavior Research Methods* 43:353. DOI: 10.3758/s13428-011-0069-9.
- 887 Peirce JW, Gray JR, Simpson S, MacAskill M, Höchenberger R, Sogo H, Kastman E, Lindeløv  
888 JK. 2019. PsychoPy2: Experiments in behavior made easy. *Behavior Research Methods*  
889 51:195–203. DOI: 10.3758/s13428-018-01193-y.

- 890 Plant RR. 2016. A reminder on millisecond timing accuracy and potential replication failure in  
891 computer-based psychology experiments: An open letter. *Behavior Research Methods*  
892 48:408–411. DOI: 10.3758/s13428-015-0577-0.
- 893 Plant RR, Hammond N, Turner G. 2004. Self-validating presentation and response timing in  
894 cognitive paradigms: How and why? *Behavior Research Methods, Instruments, &*  
895 *Computers* 36:291–303. DOI: 10.3758/BF03195575.
- 896 Plant RR, Quinlan PT. 2013. Could millisecond timing errors in commonly used equipment be a  
897 cause of replication failure in some neuroscience studies? *Cognitive, Affective, &*  
898 *Behavioral Neuroscience* 13:598–614. DOI: 10.3758/s13415-013-0166-6.
- 899 Plant RR, Turner G. 2009. Millisecond precision psychological research in a world of commodity  
900 computers: New hardware, new problems? *Behavior Research Methods* 41:598–614.  
901 DOI: 10.3758/BRM.41.3.598.
- 902 Pronk T, Wiers RW, Molenkamp B, Murre J. 2019. Mental chronometry in the pocket? Timing  
903 accuracy of web applications on touchscreen and keyboard devices. *Behavior Research*  
904 *Methods*. DOI: 10.3758/s13428-019-01321-2.
- 905 Reimers S, Stewart N. 2007. Adobe Flash as a medium for online experimentation: A test of  
906 reaction time measurement capabilities. *Behavior Research Methods* 39:365–370. DOI:  
907 10.3758/BF03193004.
- 908 Reimers S, Stewart N. 2015. Presentation and response timing accuracy in Adobe Flash and  
909 HTML5/JavaScript Web experiments. *Behavior Research Methods* 47:309–327. DOI:  
910 10.3758/s13428-014-0471-1.
- 911 Schubert TW, Murteira C, Collins EC, Lopes D. 2013. ScriptingRT: A Software Library for  
912 Collecting Response Latencies in Online Studies of Cognition. *PLOS ONE* 8:e67769. DOI:  
913 10.1371/journal.pone.0067769.
- 914 Ulrich R, Giray M. 1989. Time resolution of clocks: Effects on reaction time measurement—Good  
915 news for bad clocks. *British Journal of Mathematical and Statistical Psychology* 42:1–12.  
916 DOI: 10.1111/j.2044-8317.1989.tb01111.x.
- 917