


# Reward Enhances Online Participants' Engagement With a Demanding Auditory Task

Trends in Hearing  
Volume 25: 1–9  
© The Author(s) 2021  
DOI: 10.1177/23312165211025941  
journals.sagepub.com/home/tia  


Roberta Bianco<sup>1</sup> , Gordon Mills<sup>1,2</sup>, Mathilde de Kerangal<sup>1,\*</sup>,  
Stuart Rosen<sup>3</sup>, and Maria Chait<sup>1</sup> 

## Abstract

Online recruitment platforms are increasingly used for experimental research. Crowdsourcing is associated with numerous benefits but also notable constraints, including lack of control over participants' environment and engagement. In the context of auditory experiments, these limitations may be particularly detrimental to threshold-based tasks that require effortful listening. Here, we ask whether incorporating a performance-based monetary bonus improves speech reception performance of online participants. In two experiments, participants performed an adaptive matrix-type speech-in-noise task (where listeners select two key words out of closed sets). In Experiment 1, our results revealed worse performance in online ( $N = 49$ ) compared with in-lab ( $N = 81$ ) groups. Specifically, relative to the in-lab cohort, significantly fewer participants in the online group achieved very low thresholds. In Experiment 2 ( $N = 200$ ), we show that a monetary reward improved listeners' thresholds to levels similar to those observed in the lab setting. Overall, the results suggest that providing a small performance-based bonus increases participants' task engagement, facilitating a more accurate estimation of auditory ability under challenging listening conditions.

## Keywords

speech-in-noise, coordinate response measure, CRM, remote testing, bonus

Received 10 February 2021; Revised 6 May 2021; accepted 28 May 2021

There is a growing interest in remote testing, both in the context of basic research (Anwyl-Irvine et al., 2020; Backx et al., 2020; Hartshorne et al., 2019; Shapiro et al., 2020) and clinical screening (Paglialonga et al., 2020; Sevier et al., 2019; Shafiro et al., 2020; Sheikh Rashid et al., 2017; Swanepoel & Clark, 2019; Swanepoel et al., 2019; Watson et al., 2012). The ability to conduct experiments online facilitates rapid data acquisition and provides access to a larger and more diverse subject pool than that available for lab-based investigations (Casey et al., 2017). However, in contrast to the lab setting, online experiments are associated with a lack of control over participants' equipment, environment, and engagement (Chandler & Paolacci, 2017; Clifford & Jerit, 2014). These limitations may be particularly detrimental to auditory assessments that often rely on highly controlled stimulus delivery and

necessitate focused engagement from the participant (e.g., Harrison & Müllensiefen, 2018).

Tasks that require effortful listening (e.g., when trying to estimate performance at threshold, or the just noticeable difference in a particular acoustic feature) may be

<sup>1</sup>UCL Ear Institute, University College London, London, United Kingdom

<sup>2</sup>National Institute for Health Research UCL Hospitals Biomedical Research Centre, Deafness and Hearing Problems Theme, London, United Kingdom

<sup>3</sup>UCL Speech, Hearing and Phonetic Sciences, University College London, London, United Kingdom

\*Mathilde de Kerangal is now at the Department of Bioengineering and Centre for Neurotechnology, Imperial College London, London, UK.

### Corresponding author:

Roberta Bianco, UCL Ear Institute, 332 Gray's Inn Road, London WC1X 8EE, UK.

Email: r.bianco@ucl.ac.uk



particularly susceptible to issues related to task engagement (including attention, motivation, and commitment). In laboratories or clinics, engagement is controlled by creating a “sterile environment” that isolates the participants from potential sources of distraction (e.g., their mobile phone, software notifications, doorbell, housemates, etc.). Compliance and motivation are promoted through face-to-face interaction with the experimenter (Guéguen & Pascual, 2000; Karakostas & Zizzo, 2016). To understand how these factors affect data obtained from online participants, in this series of experiments, we investigated how performance on one version of widely used auditory speech-in-noise perception tasks differs between in-lab and online settings and whether monetary reward may be used as a mean to encourage participant engagement.

We used an adaptive speech-in-noise task based on target materials similar to the Coordinate Response Measure (CRM) corpus of Bolia et al. (2000). The CRM measures the ability to identify two keywords (color and number words) in a spoken target sentence always cued by a so-called call sign. Participants are instructed to attend to the target sentence while ignoring a masker. The CRM is part of a family of adaptive speech reception in noise tests (see also digit-in-noise test commonly used in audiology practice; De Sousa et al., 2019). These paradigms have been shown to be powerful tests of listening in complex environments because of their sensitivity to small intelligibility changes in highly noisy backgrounds, their applicability to testing with different maskers, and their relative independence from semantic/syntactic cues (Brungart, 2001; De Sousa et al., 2020; Eddins & Liu, 2012; Humes et al., 2017). Accumulating work demonstrates that speech reception thresholds (SRTs) estimated with an adaptive CRM task correlate with audiometric thresholds and with age (de Kerangal et al., 2020; Schoof & Rosen, 2014; Venezia et al., 2020), rendering it a potentially efficient proxy of hearing ability (Semeraro et al., 2017). An additional advantage is that the task relies on manipulating the relative intensity of the target and the masker, and performance is largely independent of overall level over a reasonable range. Outcomes are therefore less affected by calibration of equipment compared with other tasks that rely on absolute sound level. These considerations make the CRM, as well as other similar speech-in-noise tasks (De Sousa et al., 2019, 2020), particularly attractive for estimating auditory abilities in online settings.

We first asked whether performance among young listeners recruited “blindly” online is consistent with that observed in the highly controlled laboratory setting. Results suggested poorer performance by online listeners. We hypothesized that reduced performance in the online compared with the in-lab sample may reflect a

lack of task engagement or motivation among the online cohort. Therefore, building on existing evidence that monetary reward can improve performance in tasks that involve executive or perceptual functions (Libera & Chelazzi, 2006; Plain et al., 2020; Shen & Chun, 2011), we asked whether incorporating a performance-based monetary bonus in a group of online participants could improve speech reception performance relative to an online group that does not receive a bonus. Our results revealed that a monetary bonus improved listeners’ threshold and that the resulting SRT distribution was similar to that observed in the lab setting. Overall, the results confirm that providing a small performance-based bonus increases participant task engagement (i.e., the readiness to exert effort and/or allocate sufficient attention to the task), facilitating a more accurate estimation of auditory ability.

## Experiment I

### Methods

**Participants.** Two participant groups ranging in age between 25 and 32 years were tested. An *in-lab group* (data pooled from de Kerangal et al., 2020 and an additional unpublished study) comprised 81 participants (59 females, mean age  $25 \pm 3$  years) who completed the task as part of a test battery. An age-matched *online group* of 49 participants (35 females, mean age  $26 \pm 3$  years) was recruited and compensated via the Prolific crowdsourcing platform. All listeners were young, native speakers of British English and reported no known hearing problems. The online sample was not formally tested for hearing problems. We assumed that this cohort of young listeners would exhibit a similar hearing profile to the aged-matched in-lab participants. Experimental procedures were approved by the research ethics committee of University College London, and informed consent was obtained from each participant.

**Stimuli and Procedure.** An SRT for each participant was obtained using target sentences introduced by Messaoud-Galusi et al. (2011)—the Children’s Coordinate Response Measure (CCRM), which is a modified version of the CRM corpus described by Bolia et al. (2000). The modifications were made to be able to embed the materials in the task as a straightforward command, and using call signs (here the animal name) that would be more appropriate for use with children, without precluding the use of the material in adults, nor changing the essential properties of the corpus. Note that the CCRM as used here is likely to be at least as difficult as the original CRM (both requiring the identification of a color and a number), but here there are six colors rather than four. On each trial,

participants heard a target sentence of the form “show the dog where the [color] [number] is.” The number was a digit from 1 to 9, excluding the number 7 (due to its bisyllabic phonetic structure, which would make it easier to identify). The colors were black, white, pink, blue, green, or red. Thus, there were a total of 48 combinations (6 colors  $\times$  8 numbers). Participants were instructed to press on the correct combination of color and number on a visual interface showing an image of a dog and a list of the digits in the different colors.

The target sentences were spoken by a single female native speaker of Standard Southern British English that was presented simultaneously with a two male-speaker babble that the participants were instructed to ignore. Each talker in the babble was recorded reading two five- to six-sentence passages that were concatenated together once passages were edited to delete pauses of more than 100 ms. The two talkers were then digitally mixed together at equal levels, with random sections of the appropriate duration from this 30-s long masker chosen for each trial.

The overall level of the mixture (target speaker + babble background) was kept fixed, with only the ratio between the target and masker changing on each trial. The signal-to-noise ratio (SNR) between the babble and the target speaker was initially set to 20 dB and was adjusted using a one-up one-down adaptive procedure, tracking the 50% correct threshold (Levitt, 1971). Initial steps were of 9 dB SNR, decreasing by 2 dB following the first two reversals and then fixed at a step size of 3 dB SNR for all subsequent trials. The procedure terminated after 7 reversals or after a total of 25 trials (the latter was never reached). The SRT for one run was calculated as the mean of the SNRs in the last four reversals. Each participant performed the test in four consecutive runs of approximately 2 min each. To allow a stable measure of a listener’s threshold, the SNR was averaged over the last four reversals within each run and then across the last three runs (Run 1 was used as practice). In all individual runs, a stable threshold was achieved within <20 trials. The in-lab data for this experiment are drawn from de Kerangal et al. (2020), and it was therefore important to use the parameters used in that study. de Kerangal et al. demonstrated that this parameter set produces reliable thresholds and yields the expected difference in SRT between young and old adults and a correlation between SRT and audiometric measures.

The in-lab test was conducted in a double-walled soundproof booth (IAC, Winchester). The task was implemented in MATLAB using a calibrated sound delivery system. Sounds were presented with a Roland Tri-capture 24-bit 96 kHz soundcard over headphones (Sennheiser HD 595) at a comfortable listening level of 70 dB sound pressure level (SPL).

For online testing, the task was implemented in JavaScript, and the Gorilla Experiment Builder platform ([www.gorilla.sc](http://www.gorilla.sc)) was used to host the experiment (Anwyl-Irvine et al., 2020). Participants were recruited and pre-screened by the Prolific platform. Otherwise, the same stimuli and test heuristics were used as in the in-lab settings. As is common practice in online auditory experiments, participants were screened for headphone use. We used a strict version of the approach introduced and validated by Milne et al (2020) which yields a 7% false positive rate. In brief, this test uses a combination of Huggins pitch stimuli (Cramer & Huggins, 1958) which are only detectable when L and R channels are presented separately to each ear, and a pair of tones ( $f_1 = 1800\text{--}2500$  Hz;  $f_2 = f_1 + 30$  Hz) presented binaurally that sound smooth when listening dichotically or to each channel alone but contain a beat when the channels are mixed (Oster, 1973). Together, these probes allow us to identify those participants who are listening dichotically through separate L and R channels (i.e., using headphones) from those listening over a single channel or over speakers. The test was validated in a large group of normal-hearing listeners. For full details, information about validation, and the links to experience the task, see Milne et al. (2020).

The CCRM task took approximately 10 min to complete. It began with a volume calibration to make sure that stimuli were presented at an appropriate level. A target sentence without a masker was used for this purpose. Participants were instructed to play the sound and adjust the volume to as high a level as possible without it being uncomfortable.

At the end of the experiment, participants completed a short questionnaire about their listening environment and equipment. We encouraged honest reports by stressing that “your answers will not affect your payment but will help us to get the best quality data.” In particular, participants were asked about how much background noise they experienced during the experiment ( $0 = \text{not at all}$ ,  $10 = \text{a lot}$ ). This measure was used as a potential exclusion criterion to make sure that group differences in performance were not explained by mere differences in environmental noise. The experiment was piloted to take about 15 min. We thus set the base-pay rate to £2, corresponding to an hourly wage of £8.

We have made our implementation openly available and ready for use via Gorilla (<https://gorilla.sc/openmaterials/171870>).

**Statistical Analysis.** We used the two-sample Kolmogorov–Smirnov (KS) test (Conover, 1998) to ascertain the existence of a statistically significant difference between the (unknown) distributions of the two groups of interest. The KS test is a commonly used nonparametric test of the equality of continuous unidimensional probability distributions, based on the maximum distance between the

cumulative distributions of the two samples. Analyses were conducted in the R environment Version 0.99.320.

## Results

Figure 1 shows the probability density function (Panel A) and the cumulative distribution function (Panel B) of the SRT obtained from the in-lab (mean SRT =  $-16.2$  dB,  $SD = 2.08$ ) and online groups (mean SRT =  $-15.1$  dB,  $SD = 2.21$ ; mean difference in-lab—online =  $-1.1$  dB). A KS test indicated a significant difference between the two distributions ( $D = .347$ ,  $p = .001$ ). The maximal difference occurred at  $-16.9$  dB, which was reached by 47% of the in-lab group and only by the 12% of the online group. Despite the low level of background noise reported by the online sample ( $1.77 \pm 2.51$  from a range of 0 to 10), we repeated the analysis by excluding those participants who reported a high level of noise ( $\geq 5$ ; final sample  $N = 42$ ). The difference between groups was unaltered ( $D = .350$ ,  $p = .002$ ).

The overall pattern of results demonstrates that, relative to the in-lab cohort, fewer people in the online group achieved very low thresholds, suggesting that online testing may provide a less accurate measure of listeners' speech-in-noise detection performance. The differences between the online and in-lab groups may arise due to a poorer control of participants' listening environment and/or motivation.

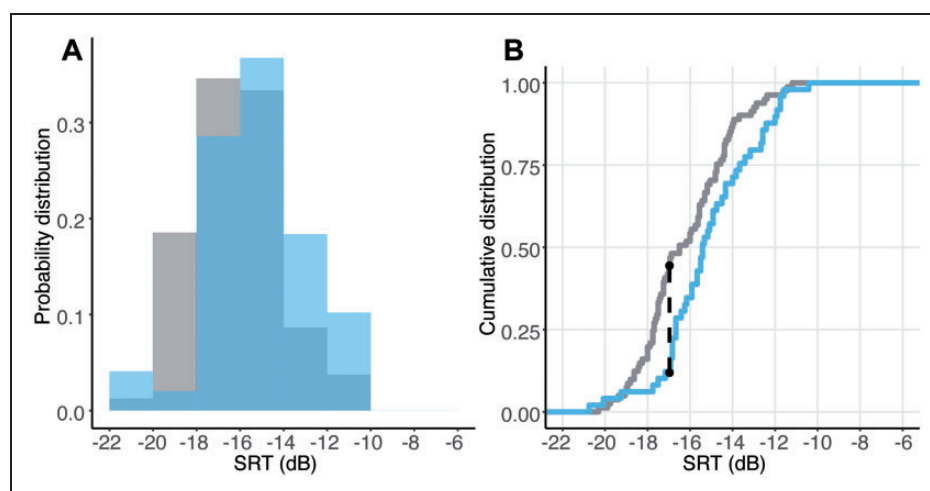
## Experiment 2

### Methods

**Participants.** Two hundred young, normal-hearing listeners ranging in age from 22 to 30 years (128 females, mean

age  $26 \pm 2.5$ ) were recruited online as described in Experiment 1. They were randomly assigned to one of two experimental groups. All participants received a fair base payment for the time spent on the experiment (Prolific recommends £8 per hour; £2 for 15 min). One group ( $N = 100$ , 62 females) additionally received a performance-based monetary bonus (up to £5) on top of the base pay (**BONUS+**). The other group ( $N = 100$ , 66 females) received no bonus (**BONUS-**).

**Stimuli and Procedure.** The procedure was similar to that described for Experiment 1. The **BONUS+** and **BONUS-** groups received identical instructions and feedback. Encouraging language was used to maximize participant motivation. After each run, the achieved threshold was displayed, and participants were challenged to try to "beat their score" in the next run. The **BONUS+** group was additionally informed that each threshold was linked to a monetary bonus. They were told that at the end of the experiment, they would receive the bonus (up to £5) associated with the best threshold reached (e.g., if they reached thresholds  $-10$ ,  $-15$ ,  $-17$ ,  $-10$  over the runs, they were paid a bonus linked to threshold  $-17$ ). At the end of each run, participants were shown the current threshold and the bonus, but also the bonus they could receive if they improve their threshold in the following run. The bonus was preassigned to SNR values from  $-1$  to  $-28$  (in steps of 1) through an exponential function so that improvements at lower, more difficult thresholds were rewarded more than improvements at levels expected to be easily reached by young normal-hearing listeners. As in Experiment 1, following the main task, participants answered a set of questions about their listening environment. They were also asked to answer on a



**Figure 1.** A: Probability density distributions of the in-lab (gray) and online (blue) groups. B: Cumulative distribution of the in-lab and online groups. The black dashed line indicates the SRT at which the greatest distance between the two distributions was observed. Overall, the data pattern is consistent with a rightward shift (toward higher SRTs) of the online distribution. SRT = speech reception threshold.



scale from 0 to 10 (0 = *not at all*, 10 = *a lot*) how motivated they were in performing the task, and how engaging they found the task to be.

The base pay was set to £2 (for 15 min) for all participants. The average obtainable bonus for the BONUS+ group was £2 (range £0–5), therefore allowing them to double their pay. The BONUS+ group was only informed of the bonus at the instructions stage. To avoid bias in the selection process, participants were unaware of the possibility of being assigned to one or the other group when they signed up to the study.

## Results

Both the BONUS+ and BONUS– groups reported a similar level of environmental noise (BONUS+ =  $1.87 \pm 2.75$ ; BONUS– =  $1.44 \pm 2.27$ ; *t*-test:  $t(2,198) = 1.203$ ,  $p = .230$ ). However, to focus on the effect of bonus on performance, we excluded those participants who reported a level of noise  $\geq 5$  (on a scale from 0 to 10) resulting in the exclusion of  $\sim 15$  participants from each group (final numbers: BONUS+  $N = 84$ ; BONUS–  $N = 90$ ). Figure 2 shows the probability density function (Panel A) and the cumulative distribution function (Panel B) of the SRT obtained for the BONUS+ (mean SRT =  $-16.1$  dB,  $SD = 2.54$ ) and BONUS– (mean SRT =  $-15.1$  dB,  $SD = 2.33$ ) groups. Data from the in-lab group (see Experiment 1) are also provided as a benchmark.

KS tests indicated a significant difference between the BONUS+ versus BONUS– distributions ( $D = .276$ ,  $p = .003$ ), revealing better performance in the BONUS+ compared with the BONUS– group. The maximum

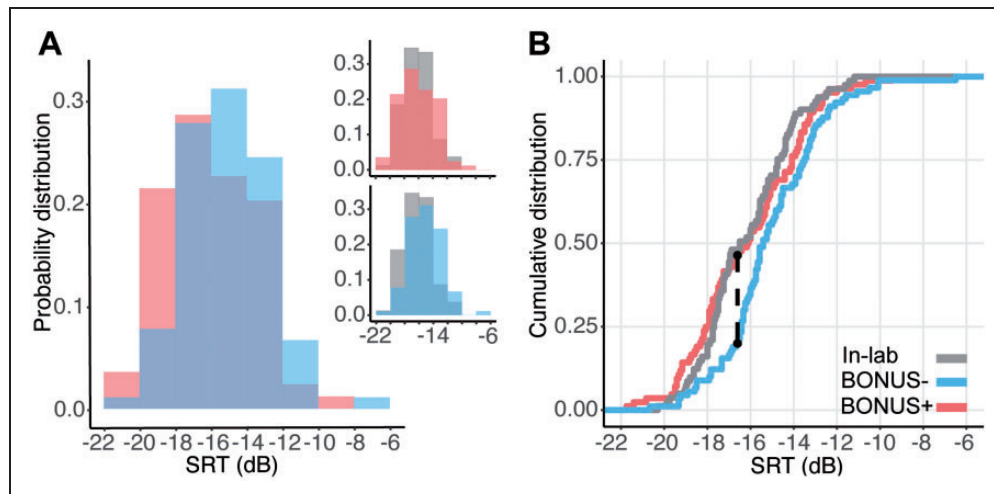
difference occurred at  $-16.6$  dB, which was reached by 47% of the BONUS+ and only by the 21% of the BONUS– group. The comparison of these two distributions with the in-lab distribution indeed showed that the BONUS+ performance was similar to the in-lab one ( $D = .127$ ,  $p = .519$ ), whilst the BONUS– was different ( $D = .304$ ,  $p = .001$ ). The results thus indicate that the provision of a bonus increased the proportion of high-performing participants in the online group to the levels exhibited by the in-lab cohort.

In an additional analysis, we compared the in-lab group with the online data pooled from the online group of Experiment 1 and the BONUS– group of Experiment 2 (for a total of  $N = 132$ , excluding participants who reported a level of background noise  $\geq 5$ ; note all results hold even without excluding participants based on noise reports). A KS test confirmed that online performance in the absence of an additional bonus is worse than that obtained in-lab ( $D = .318$ ,  $p < .001$ ), in line with what was observed in Experiment 1.

Consistent with the interpretation that offering a bonus increased motivation, the BONUS+ group reported higher ratings of task engagement (BONUS+ =  $9.1 \pm 1.2$ ; BONUS– =  $8.4 \pm 1.5$ ; from range of 0–10;  $t(2,98) = 2.68$ ,  $p = .009$ ) and motivation (BONUS+ =  $9.2 \pm 1.12$ ; BONUS– =  $8.4 \pm 1.3$ ; from range of 0 to 10;  $t(2,98) = 3.212$ ,  $p = .002$ ) compared with the BONUS– group.

## Discussion

We report two main findings. First, we showed that the SRT of blindly recruited online participants was poorer



**Figure 2.** A: Probability density distributions (relative proportion) of the online BONUS+ (pink) versus the BONUS– (blue) groups. The insets show the probability density distributions of the BONUS+ (top) and the BONUS– (bottom) groups against the in-lab sample. B: Cumulative distribution of the BONUS+ and BONUS– groups. The data from the in-lab (gray) are plotted as benchmark. The black dashed line indicates the SRT at which the greatest distance between the BONUS+ and BONUS– distributions is observed. Overall, the data pattern is consistent with a leftward shift (toward better SRTs) of the BONUS+ relative to the BONUS– groups. SRT = speech reception threshold.

than that observed among an age-matched control group in the lab setting. Second, we demonstrated that the provision of a small performance-based monetary bonus improved online listeners' speech-in-noise performance to levels similar to those observed in the lab setting.

The results from Experiment 1 revealed that the distribution of the SRT in the online group differed from that obtained from the in-lab cohort: In the lab, 47% of listeners achieved an SRT below  $\sim -17$  dB. In contrast, within the online cohort, only 12% of participants reached that threshold. This discrepancy is relevant to consider when using remote testing to build normative data, or to accurately estimate hearing loss across the population.

In Experiment 2, we showed that a performance-based monetary bonus increased the proportion of highly performing participants up to levels similar to those observed in the lab. This suggests that the difference in performance between the online and in-lab groups observed in Experiment 1 is not mainly driven by constraints to the sound environment but rather associated with reduced task engagement among the online participants.

With the blooming of online experiments, it is important to understand how we can improve the quality of data obtained in remote auditory assessments (Leensen et al., 2011; Milne et al., 2020; Slote & Strand, 2016). Our finding that reward increased the proportion of participants who achieved low SRTs demonstrates that participant attention, motivation, and commitment are important factors to consider when auditory tests involving effortful listening are conducted online.

Higher task engagement in the in-lab than in the online population probably results from several factors that characterize the laboratory experience: the authority of the experimenter, the absence of temptation/distractions, the effort taken to come to the lab, and so forth. All these factors are likely to make in-lab participants already quite motivated. Similar considerations may apply to certain online testing situations. For example, participants in remote clinical assessments are likely to be highly intrinsically motivated to do their best, as revealed by studies reporting similar results between testing in the clinic and at home (de Graaff et al., 2018; Whitton et al., 2016). However, in many cases, online participants are unsupervised and anonymous, and often mainly motivated by financial incentives (Buhrmester et al., 2011; Litman et al., 2015). In a recent in-house survey conducted by Prolific.co, approximately 50% of the surveyed users stated that the amount of pay is the factor that most motivates them to take part in a study (<https://prolific2.typeform.com/report/PoUZHEmk/ttebnlTEllbRvdcg>). Therefore, a monetary bonus is an efficient method for increasing task engagement. This

consideration is also supported by the fact that the BONUS+ group in Experiment 2 reported higher ratings of task engagement and motivation compared with the BONUS- group.

Previous studies suggest that performance on many crowdsourcing tasks does not differ, and sometimes even exceeds that measured in the lab (Hauser & Schwarz, 2016, but see Harrison & Müllensiefen, 2018; Slote & Strand, 2016). In addition, a monetary incentive (amount of pay) does not always affect performance (van den Berg et al., 2019): For example, previous studies reported no modulatory effects of amount of monetary incentive on the quality of online performance in tasks such as speech transcription (Marge et al., 2010). Internal consistency in psychological surveys and attention in following instructions were also unaffected by different levels of payment (Buhrmester et al., 2011, but see Litman et al., 2015). However, the impact of incentive may depend on the kind of task under investigation. Financial incentives may have little effect on performance when the task is too easy or when return on effort is low, for example, when it is hard to improve performance (Camerer & Hogarth, 1999). Our finding that reward influences performance in the CCRM task is possibly linked to the fact that the return on effort is high: The task relies on attention to fine perceptual details, and increasing effort has the potential to lead to a notable improvement in performance.

The effect of incentives on performance may also be nuanced by how the reward is operationalized, in particular whether it is fixed or adaptive. For example, recent studies using demanding auditory tasks and where reward was fixed at a high or low value have reported no effect of reward on behavioral measures such as accuracy or response time (Koelewijn et al., 2018, 2021; Richter, 2016; see also Carolan et al., 2021). In contrast, Shen and Chun (2011), using a range of executive and perceptual tasks, demonstrated that reward can encourage participants to perform better when it is progressively increased from trial to trial, but not when the same high reward level is maintained. Furthermore, the effect of reward in Shen and Chun (2011) appeared to persist even when the ultimate outcome was success in a competition (e.g., a monetary reward assigned to the top 10% participants based on performance) rather than money itself (e.g., performance-based earning with no competition). Therefore, particularly in experiments that require many trials, a competitive setting may be a more effective incentive than a small (a few cents) reward per trial.

The present study relied on the data from Experiment 1 to adjust the bonus growth rate. It is important to acknowledge, however, that paying a bonus based on performance may disadvantage certain participants (e.g., hearing impaired individuals in the present case)

in that the maximum bonus amount will not be equally achievable by all participants despite comparable effort to perform the task. Online settings, where the researcher has no contact with the participants, make it particularly difficult to determine whether poor performance (associated with a low bonus) is due to inability to perform well (e.g., due to hearing impairment), poor understanding of instructions, or lack of engagement with the task. To mitigate this ethical concern, a fair base pay for the time spent on participation in the experiment is therefore critical.

## Conclusions

How reward might motivate performance is an empirical question and a long-standing object of debate. Accumulating evidence suggests that reward does seem to matter particularly in tasks where performance depends on effortful engagement (Camerer & Hogarth, 1999). The CCRM task used here is analogous to many threshold-based tasks commonly used in auditory research. The observed effect of bonus on performance should thus generalize to other auditory tasks, helping to motivate participants to exert the extra bit of effort that is needed when the task becomes just doable.

## Declaration of Conflicting Interests

The authors declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.



## Funding

The authors disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: This research was funded by a BBSRC grant (BB/P003745/1) to M. C. and the NIHR UCLH BRC Deafness and Hearing Problems Theme.

## Open Practice Statements

Stimuli and code implementing the CCRM test can be found in Gorilla Open Materials: <https://gorilla.sc/openmaterials/171870>.

## ORCID iDs

Roberta Bianco  <https://orcid.org/0000-0001-9613-8933>  
 Maria Chait  <https://orcid.org/0000-0002-7808-3593>

## References

- Anwyl-Irvine, A. L., Dalmaijer, E., Hodges, N., & Evershed, J. (2020). *Online timing accuracy and precision: A comparison of platforms, browsers, and participant's devices* (pp. 1–22). <https://doi.org/10.31234/osf.io/jfeca>
- Backx, R., Skirrow, C., Dente, P., Barnett, J. H., & Cormack, F. K. (2020). Comparing web-based and lab-based cognitive assessment using the Cambridge neuropsychological test automated battery: A within-subjects counterbalanced study. *Journal of Medical Internet Research, 22*(8), e16792. <https://doi.org/10.2196/16792>
- Bolia, R. S., Nelson, W. T., Ericson, M. A., & Simpson, B. D. (2000). A speech corpus for multitalker communications research. *The Journal of the Acoustical Society of America, 107*(2), 1065–1066. <https://doi.org/10.1121/1.428288>
- Brungart, D. S. (2001). Evaluation of speech intelligibility with the coordinate response measure. *The Journal of the Acoustical Society of America, 109*(5), 2276–2279. <https://doi.org/10.1121/1.1357812>
- Buhrmester, M., Kwang, T., & Gosling, S. D. (2011). Amazon's mechanical Turk: A new source of inexpensive, yet high-quality, data? *Perspectives on Psychological Science, 6*(1), 3–5. <https://doi.org/10.1177/1745691610393980>
- Camerer, C. F., & Hogarth, R. M. (1999). The effects of financial incentives in experiments: A review and capital-labor-production framework. *Journal of Risk and Uncertainty, 19*(1–3), 7–42. <https://doi.org/10.1023/a:1007850605129>
- Carolan, P. J., Heinrich, A., Munro, K. J., & Millman, R. E. (2021). Financial reward has differential effects on behavioural and self-report measures of listening effort. *International Journal of Audiology, 0*(0), 1–11. <https://doi.org/10.1080/14992027.2021.1884907>
- Casey, L. S., Chandler, J., Levine, A. S., Proctor, A., & Strolovitch, D. Z. (2017). Intertemporal differences among MTurk workers: Time-based sample variations and implications for online data collection. *SAGE Open, 7*(2), 1–15. <https://doi.org/10.1177/2158244017712774>
- Chandler, J. J., & Paolacci, G. (2017). Lie for a dime: When most prescreening responses are honest but most study participants are impostors. *Social Psychological and Personality Science, 8*(5), 500–508. <https://doi.org/10.1177/1948550617698203>
- Clifford, S., & Jerit, J. (2014). Is there a cost to convenience? An experimental comparison of data quality in laboratory and online studies. *Journal of Experimental Political Science, 1*(2), 120–131. <https://doi.org/10.1017/xps.2014.5>
- Conover, W. J. (1998). *Practical nonparametric statistics*. John Wiley & Sons.
- Cramer, E. M., & Huggins, W. H. (1958). Creation of pitch through binaural interaction. *Journal of the Acoustical Society of America, 30*(5), 413–417. <https://doi.org/10.1121/1.1909628>
- de Graaff, F., Huysmans, E., Merkus, P., Theo Goverts, S., & Smits, C. (2018). Assessment of speech recognition abilities in quiet and in noise: A comparison between self-administered home testing and testing in the clinic for adult cochlear implant users. *International Journal of Audiology, 57*(11), 872–880. <https://doi.org/10.1080/14992027.2018.1506168>
- de Kerangal, M., Vickers, D., & Chait, M. (2020). The effect of healthy aging on change detection and sensitivity to predictable structure in crowded acoustic scenes. *Hearing Research, 399*, 108074. <https://doi.org/10.1016/j.heares.2020.108074>
- De Sousa, K. C., Smits, C., Moore, D. R., Myburgh, H. C., & Swanepoel, D. W. (2020). Pure-tone audiometry without bone-conduction thresholds: Using the digits-in-noise test to detect conductive hearing loss. *International Journal of*



- Audiology*, 59(10), 801–808. <https://doi.org/10.1080/14992027.2020.1783585>
- De Sousa, K. C., Swanepoel, D. W., Moore, D. R., Myburgh, H. C., & Smits, C. (2019). Improving sensitivity of the digits-in-noise test using antiphasic stimuli. *Ear and Hearing*, 41(2), 442–450. <https://doi.org/10.1101/677609>
- Eddins, D. A., & Liu, C. (2012). Psychometric properties of the coordinate response measure corpus with various types of background interference. *The Journal of the Acoustical Society of America*, 131(2), EL177–EL183. <https://doi.org/10.1121/1.3678680>
- Guéguen, N., & Pascual, A. (2000). Evocation of freedom and compliance: The “but you are free of...” technique. *Current Research in Social Psychology*, 5(18), 264–270.
- Harrison, P. M. C., & Müllensiefen, D. (2018). Development and validation of the computerised adaptive beat alignment test (CA-BAT). *Scientific Reports*, 8(1), 1–19. <https://doi.org/10.1038/s41598-018-30318-8>
- Hartshorne, J. K., de Leeuw, J. R., Goodman, N. D., Jennings, M., & O'Donnell, T. J. (2019). A thousand studies for the price of one: Accelerating psychological science with Pushkin. *Behavior Research Methods*, 51(4), 1782–1803. <https://doi.org/10.3758/s13428-018-1155-z>
- Hauser, D. J., & Schwarz, N. (2016). Attentive Turkers: MTurk participants perform better on online attention checks than do subject pool participants. *Behavior Research Methods*, 48(1), 400–407. <https://doi.org/10.3758/s13428-015-0578-z>
- Humes, L. E., Kidd, G. R., & Fogerty, D. (2017). Exploring use of the coordinate response measure in a multitalker babble paradigm. *Journal of Speech, Language, and Hearing Research*, 60(3), 741–754. [https://doi.org/10.1044/2016\\_JSLHR-H-16-0042](https://doi.org/10.1044/2016_JSLHR-H-16-0042)
- Karakostas, A., & Zizzo, D. J. (2016). Compliance and the power of authority. *Journal of Economic Behavior and Organization*, 124, 67–80. <https://doi.org/10.1016/j.jebo.2015.09.016>
- Koelwijn, T., Zekveld, A. A., Lunner, T., & Kramer, S. E. (2018). The effect of reward on listening effort as reflected by the pupil dilation response. *Hearing Research*, 367, 106–112. <https://doi.org/10.1016/j.heares.2018.07.011>
- Koelwijn, T., Zekveld, A. A., Lunner, T., & Kramer, S. E. (2021). The effect of monetary reward on listening effort and sentence recognition. *Hearing Research*, 406, 108255. <https://doi.org/10.1016/j.heares.2021.108255>
- Leensen, M. C. J., De Laat, J. A. P. M., Snik, A. F. M., & Dreschler, W. A. (2011). Speech-in-noise screening tests by internet, part 2: Improving test sensitivity for noise-induced hearing loss. *International Journal of Audiology*, 50(11), 835–848. <https://doi.org/10.3109/14992027.2011.595017>
- Levitt, H. (1971). Transformed up-down methods in psychoacoustics. *The Journal of the Acoustical Society of America*, 49(2B), 467–477. <https://doi.org/10.1121/1.1912375>
- Libera, C.D., & Chelazzi, L. (2006). Visual selective attention and the effects of monetary rewards. *Psychological Science*, 17(3), 222–227. <https://doi.org/10.1111/j.1467-9280.2006.01689.x>
- Litman, L., Robinson, J., & Rosenzweig, C. (2015). The relationship between motivation, monetary compensation, and data quality among US- and India-based workers on Mechanical Turk. *Behavior Research Methods*, 47(2), 519–528. <https://doi.org/10.3758/s13428-014-0483-x>
- Marge, M., Banerjee, S., & Rudnicky, A. I. (2010, March). Using the Amazon Mechanical Turk for transcription of spoken language. In *2010 IEEE International Conference on Acoustics, Speech and Signal Processing* (pp. 5270–5273). IEEE, doi: 10.1109/ICASSP.2010.5494979.
- Messaoud-Galus, S., Hazan, V., & Rosen, S. (2011). Investigating speech perception in children with dyslexia: Is there evidence of a consistent deficit in individuals? *Journal of Speech, Language, and Hearing Research*, 54(6), 1682–1701. [https://doi.org/10.1044/1092-4388\(2011\)09-0261](https://doi.org/10.1044/1092-4388(2011)09-0261)
- Milne, A. E., Bianco, R., Poole, K. C., Zhao, S., Oxenham, A. J., Billig, A. J., & Chait, M. (2020). An online headphone screening test based on dichotic pitch. *Behavior Research Methods*, 1–12. <https://doi.org/10.1101/2020.07.21.214395>
- Oster, G. (1973). Auditory beats in the brain. *Scientific American*, 229(4), 94–102.
- Pagialonga, A., Polo, E. M., Zanet, M., Rocco, G., van Waterschoot, T., & Barbieri, R. (2020). An automated speech-in-noise test for remote testing: Development and preliminary evaluation. *American Journal of Audiology*, 29(3 Special Issue), 564–576. [https://doi.org/10.1044/2020\\_AJA-19-00071](https://doi.org/10.1044/2020_AJA-19-00071)
- Plain, B., Richter, M., Zekveld, A. A., Lunner, T., Bhuiyan, T., & Kramer, S. E. (2021). Investigating the Influences of Task Demand and Reward on Cardiac Pre-Ejection Period Reactivity During a Speech-in-Noise Task. *Ear and Hearing*, 42(3), 718. <https://doi.org/10.1097/aud.0000000000000971>
- Richter, M. (2016). The moderating effect of success importance on the relationship between listening demand and listening effort. *Ear and Hearing*, 37, 111S–117S. <https://doi.org/10.1097/AUD.0000000000000295>
- Schoof, T., & Rosen, S. (2014). The role of auditory and cognitive factors in understanding speech in noise by normal-hearing older listeners. *Frontiers in Aging Neuroscience*, 6(Oct), 1–14. <https://doi.org/10.3389/fnagi.2014.00307>
- Semeraro, H. D., Rowan, D., van Besouw, R. M., & Allsopp, A. A. (2017). Development and evaluation of the British English coordinate response measure speech-in-noise test as an occupational hearing assessment tool. *International Journal of Audiology*, 56(10), 749–758. <https://doi.org/10.1080/14992027.2017.1317370>
- Sevier, J. D., Choi, S., & Hughes, M. L. (2019). Use of direct-connect for remote speech-perception testing in cochlear implants. *Ear and Hearing*, 40(5), 1162–1173. <https://doi.org/10.1097/AUD.0000000000000693>
- Shafiro, V., Hebb, M., Walker, C., Oh, J., Hsiao, Y., Brown, K., Sheft, S., Li, Y., Vasil, K., & Moberly, A. C. (2020). Development of the basic auditory skills evaluation battery for online testing of cochlear implant listeners. *American Journal of Audiology*, 29(3S), 577–590. [https://doi.org/10.1044/2020\\_AJA-19-00083](https://doi.org/10.1044/2020_AJA-19-00083)
- Shapiro, M. L., Norris, J. A., Wilbur, J. C., Brungart, D. S., & Clavier, O. H. (2020). *TabSINT*: Open-source mobile software for distributed studies of hearing. *International Journal*



- of *Audiology*, 59(sup1), S12–S19. <https://doi.org/10.1080/14992027.2019.1698776>
- Sheikh Rashid, M., Dreschler, W. A., & de Laat, J. A. P. M. (2017). Evaluation of an internet-based speech-in-noise screening test for school-age children. *International Journal of Audiology*, 56(12), 967–975. <https://doi.org/10.1080/14992027.2017.1378932>
- Shen, Y. J., & Chun, M. M. (2011). Increases in rewards promote flexible behavior. *Attention, Perception, and Psychophysics*, 73(3), 938–952. <https://doi.org/10.3758/s13414-010-0065-7>
- Slote, J., & Strand, J. F. (2016). Conducting spoken word recognition research online: Validation and a new timing method. *Behavior Research Methods*, 48(2), 553–566. <https://doi.org/10.3758/s13428-015-0599-7>
- Swanepoel, D. W., & Clark, J. L. (2019). Hearing healthcare in remote or resource-constrained environments. *Journal of Laryngology and Otology*, 133(1), 11–17. <https://doi.org/10.1017/S0022215118001159>
- Swanepoel, D. W., De Sousa, K. C., Smits, C., & Moore, D. R. (2019). Mobile applications to detect hearing impairment: Opportunities and challenges. *Bulletin of the World Health Organization*, 97(10), 717–718. <https://doi.org/10.2471/BLT.18.227728>
- van den Berg, R., Zou, Q., & Ma, W. J. (2019). No effect of monetary reward in a visual working memory task. *BioRxiv*, 1–13. <https://doi.org/10.1101/767343>
- Venezia, J. H., Leek, M. R., & Lindeman, M. P. (2020). Suprathreshold differences in competing speech perception in older listeners with normal and impaired hearing. *Journal of Speech, Language, and Hearing Research*, 63(7), 2141–2161. [https://doi.org/10.1044/2020\\_JSLHR-19-00324](https://doi.org/10.1044/2020_JSLHR-19-00324)
- Watson, C. S., Kidd, G. R., Miller, J. D., Smits, C., & Humes, L. E. (2012). Telephone screening tests for functionally impaired hearing: Current use in seven countries and development of a US version. *Journal of the American Academy of Audiology*, 23(10), 757–767. <https://doi.org/10.3766/jaaa.23.10.2>
- Whitton, J. P., Hancock, K. E., Shannon, J. M., & Polley, D. B. (2016). Validation of a self-administered audiometry application: An equivalence study. *The Laryngoscope*, 126(10), 2382–2388. <https://doi.org/10.1002/lary.25988>