ELSEVIER

# Can Amazon's Mechanical Turk be used to recruit participants for internet intervention trials? A pilot study involving a randomized controlled trial of a brief online intervention for hazardous alcohol use

CrossMark

John A. Cunningham[a,b,c,*], Alexandra Godinho[a,d], Vladyslav Kushnir[a,e]

[a] Centre for Addiction and Mental Health, Toronto M5S 2S1, Canada
[b] Department of Psychiatry, University of Toronto, Toronto M5T 1R8, Canada
[c] Research School of Population Health, The Australian National University, Canberra 2601, Australia
[d] Dalla Lana School of Public Health, University of Toronto, M5S 1A1, Canada
[e] Leslie Dan Faculty of Pharmacy, University of Toronto, Toronto M5S 3M2, Canada

## ARTICLE INFO

## ABSTRACT

*Objectives:* To determine whether Amazon's Mechanical Turk (MTurk) might be a viable means of recruiting participants for online intervention research. This was accomplished by conducting a randomized controlled trial of a previously validated intervention with participants recruited through MTurk.

*Methods:* Participants were recruited to complete an online survey about their alcohol use through the MTurk platform. Those who met eligibility criterion for age and problem drinking were invited to complete a 3-month follow-up. Those who agreed were randomized to receive access to an online brief intervention for drinking or were assigned to a no intervention control group (i.e., thanked and told that they would be re-contacted in 3 months).

*Results:* A total of 423 participants were recruited, of which 85% were followed-up at 3-months. All participants were recruited in 3.2 h. Only 1/3 of participants asked to access the online brief intervention did so. Of the 4 outcome variables (number of drinks in a typical week, highest number on one occasion, number of consequences, AUDIT consumption subscale), one displayed a significant difference between conditions. Participants in the intervention group reported a greater reduction between on the AUDIT consumption subscale between baseline and 3-month follow-up compared to those in the no intervention control group ($p = 0.004$).

*Conclusions:* Despite the current pilot showing only limited evidence of impact of the intervention among participants recruited through MTurk, there is potential for conducting trials employing this population (particularly if methods are employed to make sure that participants receive the intervention). This potential is important as it could allow for the rapid conduct of multiple trials during the development stages of online interventions.

## Trial registration

ClinicalTrials.gov # NCT02905123

## 1. Introduction

Amazon's Mechanical Turk (MTurk) is an online platform in which more than half a million people have registered as 'workers' (www.mturk.com). The worker then chooses tasks (often surveys) to complete through MTurk. Amazon provides the platform for this work and acts as the mediator for secure payment to workers.

MTurk has become a popular means of collecting survey data in some areas of psychology (Buhrmester et al., 2011; Chandler and Shapiro, 2016; Daly and Nataraajan, 2015; Litman et al., 2016; Shapiro et al., 2014; Wiens and Walker, 2015). Further, participants with problem drinking, gambling, or even illicit drug use have been recruited through MTurk (Kim and Hodgins, 2017; Kristan and Suffoletto, 2015). There is also the possibility that participants for online longitudinal studies could be identified through MTurk, including for brief intervention research. The potential to quickly and easily identify large numbers of participants for online trials is important for research evaluating online interventions during the period that these interventions are being developed and refined. This is because such a study participant pool could then be repeatedly tapped to test the impact of

different versions of an intervention (e.g., treatment dismantling studies to identify active ingredients of an intervention).

However, before proposing MTurk workers as a viable source for participants in such trials, it is important to evaluate the feasibility of using MTurk for such a purpose. This pilot study proposed to test this feasibility by systematically replicating a trial of an extensively evaluated brief online intervention for hazardous alcohol use (CheckYourDrinking.net; CYD) employing participants recruited through MTurk. The goals of the pilot were: 1) to establish whether it is possible to recruit participants quickly using MTurk and to then obtain a good follow-up rate; 2) to examine whether participants would access the intervention; and 3) to test whether a significant impact of the intervention could be observed.

## 2. Methods

### 2.1. Recruitment

Potential participants were recruited using a three stage process. The study was approved by the CAMH Research Ethics Board.

#### 2.1.1. Stage 1 of recruitment

Participants were recruited through Amazon's MTurk crowdsourcing platform. A brief description of the survey was posted on MTurk, "The Centre for Addiction and Mental Health is conducting a survey on people's drinking. Only people who currently drink alcohol are asked to participate," with a link that interested potential participants could click on to access the online consent process and complete the study survey. The advertisement of the survey on MTurk was restricted to workers from Canada or the US, who had MTurk reputations of 95% or higher, and those who had completed at least 100 hits to ensure data quality (i.e., completed 100 tasks on MTurk and did not have their work rejected and returned for at least 95 of these tasks) (Peer et al., 2014). Potential participants who clicked on the link were sent to a webpage providing a brief description of the survey. Those who clicked on the link at the bottom of the brief description completed a brief eligibility screener (eligibility questions comprised of being 18 years of age or older and having consumed alcohol weekly or more often in the last year). Those who were found eligibile were sent to an electronic consent form. Those not found eligible were thanked for their participation.

#### 2.1.2. Stage 2 of recruitment

Participants identified as eligible confirmed their willingness to participate by accepting that they had read and understood the research and their rights as described on the consent form. The Stage 2 consent form contained the information that some participants would be invited to take part in another study. Those consenting then completed the the baseline survey. This survey assessed demographics (age, sex, education, marital, family income, employment status and ethnic origin), the Alcohol Use Disorder Identification Test (AUDIT; with drinking items framed to ask about the last three months) (Saunders and Conigrave, 1990), number of drinks in a typical week and highest number of drinks on one ocassion during the last three months, and number of consequences associated with drinking in the last three months (10 items adapted from Wechsler et al., 1994 with one item added asking about driving under the influence of alcohol) (Wechsler et al., 1994). The survey included a picture that showed standard drink sizes for beer, wine, and liquor (based on a drink size of 13.6 g of alcohol). Any use of alcohol related treatment access was measured using the single item screener taken from the National Epidemiological Survey on Alcohol and Related Conditions (Grant et al., 2003). In addition, four attention check questions were asked, nested within the other survey items. Participants were paid US$1.50 for completing the 10 min Stage 1 survey, in the form of an MTurk payment (note: Amazon charges a 40% fee on top of the $1.50 paid to each participant). This honorarium amount is in line with what has been collectively deemed as a fair

reward rate/amount by MTurk participants. No personally identifying information was collected within the survey, as MTurk prohibits the collection of this information from workers (please see https://requester.mturk.com/help/faq#restrictions_use_mturk for full policies). Workers' identification numbers were collected and visible on MTurk for the purposes of compensating individuals, however this ID does not grant researchers access to any identifying information.

#### 2.1.3. Stage 3 of recruitment

Upon completing the Stage 2 baseline survey, all participants were thanked for completing the survey and paid. Participants who scored 8 or more on the AUDIT (indicating current hazardous alcohol use), who reported that they had provided accurate answers and that we should keep their data, and who endorsed all four attention check questions correctly, were then sent to a page inviting them to take part in another study. These participants were asked if they would be willing to complete another survey in three months' time that asked about their drinking experiences during that time period. Further, they were told that some people would also be provided access to some more information about drinking, but that we did not know if they would receive this information at this time. However, if they did receive access to this additional information, they would be asked their impressions of it as part of the three-month survey. Finally, participants were informed that they would be paid US$10 through the MTurk portal upon completion of the three month follow-up survey. The MTurk portal allowed for sending the three-month follow-up survey to the specific participants who had agreed to take part in the follow-up. Researchers had no access to any information that could lead to personal identification of participants.

### 2.2. Randomization, experimental conditions and follow-up

Participants who agreed to complete the three-month follow-up survey were randomized (1:1 ratio with no stratification) to receive access to the Check Your Drinking screener (CYD condition) or to a no additional information condition (control condition). Those assigned to the CYD condition were told that they would be sent an email through the Mturk portal with a link to a website that would let them see how their drinking compared with others and that they would be asked their impressions of this website on the next survey. The email (sent the same day as the completion of the baseline survey) contained a link and password to a study portal that recorded which passwords had been used and provided each participant with a study specific version of the CYD. Those participants who did not use their password within one week were recontacted by email to request that they access the portal. Participants in the no intervention control condition were thanked for their participation and told that they would be contacted by email in three months' time to complete the follow-up survey. At the three-month follow-up, the MTurk portal was used to send invitation emails that contained a link to the survey. If the participant did not respond, this email was resent as a prompter 3 and 7 days later. The three-month follow-up survey asked the same drinking and drinking consequence items as the baseline survey, as well as any use of treatment services (all framed for the last three months).

### 2.3. The check your drinking (CYD) intervention

The CYD is a brief, personalized feedback intervention (Cunningham et al., 2009). Participants provide some brief demographic information about their age, sex, weight, typical cost of a drink, and country of residence, as well as 18 questions about their drinking (AUDIT, drinking in a typical week, highest number on one occasion, experience of consequences). The participant is provided with a final report that summarizes their drinking and compares it with others of the same age group, sex, and country of residence (at least for participants from Canada, the USA, and the U.K.). The efficacy of the CYD
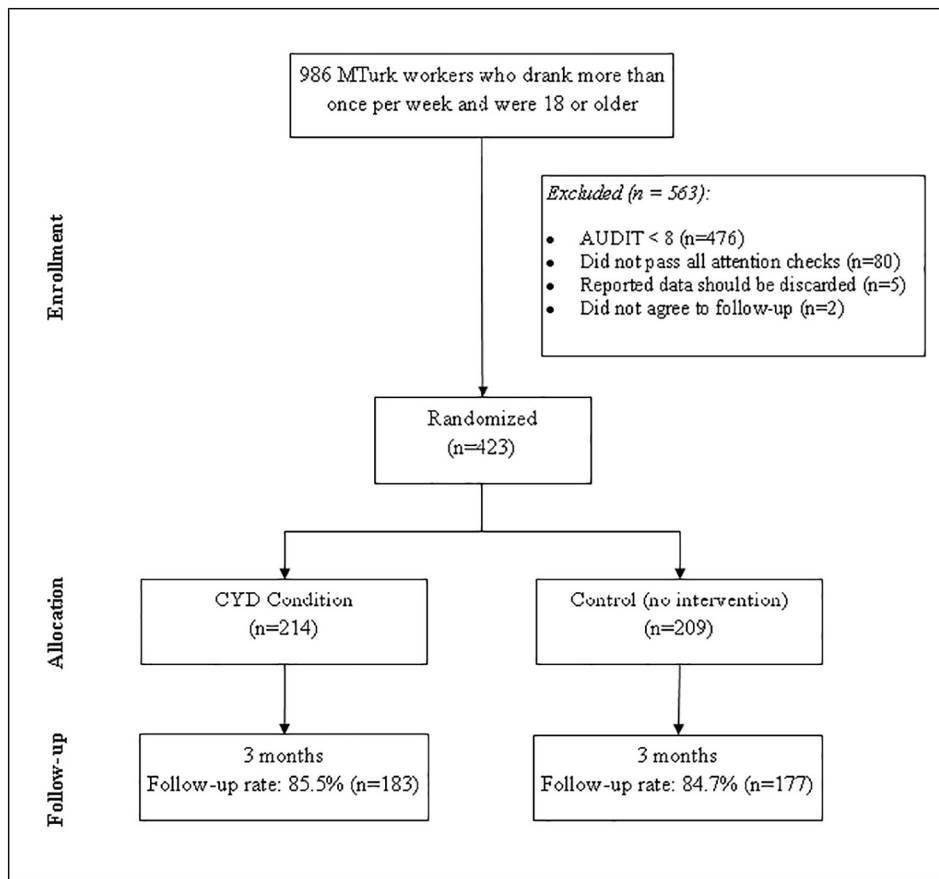
has been established through five randomized controlled trials conducted by two independent research groups (Cunningham et al., 2014; Cunningham et al., 2009; Doumas and Hannah, 2008; Doumas and Haustveit, 2008; Doumas et al., 2009) and displayed a consistent impact on reducing participants' drinking. It was chosen as the intervention of choice for the current pilot because of its brevity, and because it could reasonably be assumed to be an active intervention given the existing evidence base (i.e., if there was no observed impact of the intervention among participants from MTurk, then it could reasonably be assumed that it was something to do with the participant source, rather than the intervention itself, that lead to the observed lack of impact).

### 2.4. Sample size estimate

This pilot study employed a sample size estimate generated from one of our previous trials which proposed that the inclusion of the intervention condition would result in a 3% increase in the explained variance of the outcome variable (Cunningham et al., 2009). Following the convention that studies should be designed to have a statistical power of at least 80%, and that hypotheses be tested at the 0.05 level of significance, a final sample size (after attrition) of $N = 170$ (85 participants per condition) was required. A 20% attrition rate was estimated, thus 204 participants were needed to be found eligible and agree to the Stage 2 recruitment. At the time of planning this study, it was unknown how many participants would be needed to complete the Stage 2 survey in order to obtain this number of participants. Thus, 1000 baseline surveys were released with the possibility of releasing another 1000 if needed.

### 2.5. Data analysis

Prior to data analysis of the primary outcome variables, outliers

found to be $> 3.29$ standard deviations above the sample mean were removed and variables with skewed distributions that were not bounded by a scale maximum were log-transformed. The primary hypothesis to be tested for the RCT component of this pilot was that participants receiving access to the CYD intervention would report a greater level of reduction in number of drinks in a typical week between the baseline survey and three-month follow-up, as compared to participants in the no information control condition. To test this, a mixed-effects model with random intercept was used to estimate the fixed effects of time, intervention group and their interaction, on changes in the number of drinks consumed during a typical week between baseline and 3 months. Analyses were conducted using an intent to treat approach, with all those assigned to access the CYD intervention retained in the intervention condition, whether they accessed the intervention or not. Secondary analyses similarly employed mixed-effects models with random intercepts to examine the fixed effects of time, intervention group and their interaction on the remaining outcome variables (i.e., AUDIT-C; the sum of the three alcohol consumption items from the AUDIT) (Dawson et al., 2005), highest number of drinks on one occasion, and number of consequences experienced). For outcome variables not normally distributed, general estimated equations with negative binomial with loglink or binary logistic models were conducted in addition to mixed-effects models. For ease of interpretation, mixed-effect models were reported when model outcomes did not differ and residuals of the model were normally distributed. The mixed-effects approach allowed us to use all available participant data in the models, by using restricted maximum likelihood to account for missing data. All analyses were two-tailed and carried out at an alpha level of 0.05 using IBM SPSS, version 24.0.

## 3. Results

A total of 1252 people accessed the Stage 1 eligibility screener, of which 986 participants were found eligible and completed the Stage 2 recruitment baseline survey in a period of 3.2 h. Of these, 423 were eligible and agreed to participate in the follow-up survey (Stage 3 recruitment). These participants ($n = 423$) were randomized to condition (214 in the CYD condition and 209 in the no intervention control condition). A total of 360 (85.1%) participants completed the 3 month follow-up. See Fig. 1 for a Consort chart of the trial.

Bivariate comparisons found no significant ($p > 0.05$) differences in baseline demographic and drinking variables. Participant characteristics were as follows. Average age was 34.9 (9.4 standard deviation; SD), 57% were male, the majority were white (83%), half were married or living in a common law relationship, 71% had some post-secondary education, 73% were full-time employed (including full-time self-employed), and 29% reported a family income of less than US $20,000 per year. Participants had a mean AUDIT score of 14.3 (6.3 SD), drank 18.6 (12.3 SD) drinks in a typical week, reported 8.8 (4.4 SD) drinks as the most they drank on one occasion in the past 3 months, and experienced 2.7 (2.0 SD) consequences. A total of 13.5% of participants said that they had ever accessed treatment in relation to their alcohol use.

Of participants assigned to access the CYD intervention, only 38.3% ($n = 82$) actually used their password and accessed the intervention. For the primary outcome variable, number of drinks in a typical week, a mixed-effects model was conducted with three predictors: time, intervention group, and the time by intervention interaction (Table 1, model 1). The model revealed that the sample as a whole significantly reduced the number of reported drinks in a typical week from baseline to the 3 month follow-up (time estimate $= -0.13$, 95% confidence interval [CI] $= -0.18$ to $-0.09$ $p < 0.001$), however no differences in the level of reduction was observed across interventions (time-intervention interaction estimate $= 0.03$, 95% confidence interval [CI] $= -0.03$ to $0.10$, $p = 0.315$).

Secondary analyses examined changes over time and between intervention groups for the AUDIT-C scores, highest number of drinks reported on one occasion, and number of consequences experienced, by fitting mixed-effects models (Table 1, models 2,3,4 respectively). Overall, the models revealed that the entire sample experienced significant reductions in all three outcomes from baseline to the 3 month follow-up ($p < 0.05$). In addition, participants who received access to the CYD intervenion experienced a greater level of reduction in their AUDIT-C scores from baseline to 3 months, as compared to individuals in the control intervention ($p = 0.004$; Mean [Standard Error; SE]: CYD intervention: Baseline = 7.5 [0.16], Follow-up = 6.1 [0.17]; Control group: Baseline = 7.3 [0.16], Follow-up = 6.5 [0.17]).

## 4. Discussion

As a pilot test of recruiting participants through MTurk, the results of this trial are encouraging. It was possible to recruit a sample quickly, cost efficiently (about US$2100 to identify the 423 participants in the trial), and to obtain a good follow-up rate (85% at three-months using a US$10 payment for completion of the survey). However, while it is common to have some participants not access the intervention in online trials, the compliance rate in this sample was especially low, with only one-third accessing the CYD through the provided password portal. Perhaps because of this, and despite recruiting a sample about twice as large as that estimated in the power analysis, there was no evidence of a significant impact of the CYD intervention on the primary outcome variable – number of drinks in a typical week. Among the secondary variables, only the AUDIT-C displayed a significant impact ($p < 0.05$), with participants requested to access the CYD reporting significantly greater reductions between baseline and three-month follow-up compared to participants in the no intervention control condition. Future trials employing MTurk to recruit participants should consider means to increase this compliance rate, such as only randomizing participants after they access an intervention portal, or paying and requesting proof that the participant actually accessed the intervention website (e.g., by providing a screenshot).

It was also notable that, while we recruited participants with an AUDIT score of 8 or more, alcohol consumption among a substantial portion of the participants was fairly low (43% of participants stated that they typically drank < 15 drinks per week at baseline) compared to other trials using this same eligibility criteria but employing other forms of online advertisements (Cunningham et al., in press). Given that the outcome of the trial was to measure reductions in typical weekly drinking, the fact that so many participants were not drinking > 15 drinks per week may have limited the potential to see an impact of the intervention because of the potential for a floor effect (Cunningham, 2017). In retrospect, this relatively low drinking level is predictable, as while the MTurk sample cannot be taken as representative of the general population, previous publications have reported on the extent to which an MTurk sample mimics the distribution of demographic characteristics observed in the general population (and a general population sample with an AUDIT score of 8 or more contains many people who do not drink > 15 drinks per week) (Berinsky et al., 2012). It is suggested that future trials employ a heavy drinking inclusion criterion in addition to, or in place of, the AUDIT score of 8 or more criterion.

Another issue to consider when interpreting the results of the trial is the nature of the sample. As the researcher only knows the participants' MTurk ID, the participants are functionally anonymous. While self-reports have generally been found to be reliable, there is no way to confirm them with this sample. Perhaps more important, is the fact that the MTurk sample consists of what are essentially professional survey takers, leading to the possibility that the results generated from the

**Table 1**
Mixed-effect models results of time, intervention, and time by intervention interaction on outcome variables.

| Effect | Model 1: Drinks in a typical week | | Model 2: AUDIT - C | | Model 3: Highest # drinks in a day | | Model 4: # of consequences | |
|---|---|---|---|---|---|---|---|---|
| | Estimate ± SE | p - value | Estimate ± SE | p - value | Estimate ± SE | p - value | Estimate ± SE | p - value |
| Intercept | 1.19 ± 0.02 | **< 0.001** | 7.50 ± 0.16 | **< 0.001** | 0.89 ± 0.02 | **< 0.001** | 0.37 ± 0.02 | **< 0.001** |
| Time (Ref: Baseline) | −0.13 ± 0.02 | **< 0.001** | −1.43 ± 0.16 | **< 0.001** | −0.06 ± 0.02 | **< 0.001** | −0.06 ± 0.02 | **0.005** |
| Intervention (Ref: CYD group) | −0.01 ± 0.03 | 0.706 | −0.18 ± 0.23 | 0.416 | −0.008 ± 0.02 | 0.739 | −0.03 ± 0.03 | 0.294 |
| Time x intervention (Ref: Baseline x CYD group) | 0.03 ± 0.03 | 0.315 | 0.65 ± 0.22 | **0.004** | 0.005 ± 0.02 | 0.848 | −0.0008 ± 0.03 | 0.978 |

**Note:** SE: Standard Error.
Models 1, 3, and 4 were conducted using the log transformation of the outcome variable.
AUDIT-C: Sum of the three consumption items on the AUDIT – frequency of drinking, drinks per drinking day, frequency of 5 + days.

present sample might not be generalizable to other groups who are not experienced with completing many questionnaires. This perhaps makes an MTurk sample more like a sample of university students than one from the general population. Finally, there is the issue that, while one of the goals of the study for the researchers was to assess the impact of the CYD on participants' drinking, the goal for the participants in taking part in the trial was to be paid (and perhaps to take part in something that interested or concerned them). Such a situational dynamic could perhaps be framed as a type of workplace health intervention trial, but realistically, it is challenging to know how to assess the generalizability of the sample. Nevertheless, given the speed and cost of conducting trials with MTurk participants, there may be a place for recruiting participants from this and similar websites in order to conduct quick evaluations of components of an intervention, or to pilot the intervention as a whole, before proceeding to a full-scale trial.

## Acknowledgements

## Conflict of interest

The authors have no conflicts of interest to declare.

## References

Berinsky, A.J., Huber, G.A., Lenz, G.S., 2012. Evaluating online labor markets for experimental research: Amazon.com's Mechanical Turk. Polit. Anal. 20, 351–368.

Buhrmester, M., Kwang, T., Gosling, S.D., 2011. Amazon's Mechanical Turk: A new source of inexpensive, yet high-quality, data? Perspect. Psychol. Sci. 6 (1), 3–5. http://dx.doi.org/10.1177/1745691610393980.

Chandler, J., Shapiro, D., 2016. Conducting clinical research using crowdsourced convenience samples. Annu. Rev. Clin. Psychol. 12, 53–81. http://dx.doi.org/10.1146/annurev-clinpsy-021815-093623.

Cunningham, J.A., 2017. Unintended impact of using different inclusion cut-offs for males and females in intervention trials for hazardous drinking. Addiction 112 (5), 910–911. http://dx.doi.org/10.1111/add.13760.

Cunningham, J.A., Wild, T.C., Cordingley, J., van Mierlo, T., Humphreys, K., 2009. A randomized controlled trial of an internet-based intervention for alcohol abusers.

Addiction 104 (12), 2023–2032. http://dx.doi.org/10.1111/j.1360-0443.2009.02726.x.

Cunningham, J.A., Murphy, M., Hendershot, C.S., 2014. Treatment dismantling pilot study to identify the active ingredients in personalized feedback interventions for hazardous alcohol use: randomized controlled trial. Addict. Sci. Clin. Pract. e9 (22).

Cunningham, J.A., Shorter, G.W., Murphy, M., Kushnir, V., Rehm, J., Hendershot, C.S., 2017. Randomized controlled trial of a brief versus extended internet intervention for problem drinkers. Int. J. Behav. Med. http://dx.doi.org/10.1007/s12529-016-9604-5. (in press).

Daly, T.M., Nataraajan, R., 2015. Swapping bricks for clicks: crowdsourcing longitudinal data on Amazon Turk. J. Bus. Res. 68, 2603–2609.

Dawson, D.A., Grant, B.F., Stinson, F.S., Zhou, Y., 2005. Effectiveness of the derived Alcohol Use Disorders Identification Test (AUDIT-C) in screening for alcohol use disorders and risk drinking in the US general population. Alcohol. Clin. Exp. Res. 29 (5), 844–854.

Doumas, D.M., Hannah, E., 2008. Preventing high-risk drinking in youth in the workplace: a web-based normative feedback program. J. Subst. Abus. Treat. 34 (3), 263–271.

Doumas, D.M., Haustveit, T., 2008. Reducing heavy drinking in intercollegiate athletes: evaluation of a web-based personalized feedback program. Sport Psychol. 22, 213–229.

Doumas, D.M., McKinley, L.L., Book, P., 2009. Evaluation of two web-based alcohol interventions for mandated college students. J. Subst. Abus. Treat. 36 (1), 65–74.

Grant, B.F., Moore, T.C., Shepard, J., Kaplan, K., 2003. Source and Accuracy Statement. Wave 1. National Epidemiologic Survey on Alcohol and Related Conditions (NESARC) (Retrieved from Bethesda, MD).

Kim, H.S., Hodgins, D.C., 2017. Reliability and validity of data obtained from alcohol, cannabis, and gambling populations on Amazon's Mechanical Turk. Psychol. Addict. Behav. 31 (1), 85–94. http://dx.doi.org/10.1037/adb0000219.

Kristan, J., Suffoletto, B., 2015. Using online crowdsourcing to understand young adult attitudes toward expert-authored messages aimed at reducing hazardous alcohol consumption and to collect peer-authored messages. Transl. Behav. Med. 5 (1), 45–52. http://dx.doi.org/10.1007/s13142-014-0298-4.

Litman, L., Robinson, J., Abberbock, T., 2016. TurkPrime.com: a versatile crowdsourcing data acquisition platform for the behavioral sciences. Behav. Res. Methods 49 (2), 433–442. http://dx.doi.org/10.3758/s13428-016-0727-z.

Peer, E., Vosgerau, J., Acquisti, A., 2014. Reputation as a sufficient condition for data quality on Amazon Mechanical Turk. Behav. Res. Methods 46 (4), 1023–1031. http://dx.doi.org/10.3758/s13428-013-0434-y.

Saunders, J.B., Conigrave, K.M., 1990. Early identification of alcohol problems. Can. Med. Assoc. J. 143, 1060–1069.

Shapiro, D., Chandler, J., Mueller, P.A., 2014. Using Mechanical Turk to study clinical populations. Clin. Psychol. Sci. Pract. 1, 213–220.

Wechsler, H., Davenport, A., Dowdall, G., Moeykens, B., Castillo, S., 1994. Health and behavioral consequences of binge drinking in college: a national survey of students at 140 campuses. J. Am. Med. Assoc. 272, 1672–1677.

Wiens, T.K., Walker, L.J., 2015. The chronic disease concept of addiction: helpful or harmful? Addict. Res. Theory 23, 309–321.