Corresponding Author: Mr. Christopher D Cabrall, M. Sc.

Corresponding Author's Institution: Delft University of Technology

First Author: Christopher D Cabrall, M. Sc.

Order of Authors: Christopher D Cabrall, M. Sc.; Zhenji Lu; Miltos Kyriakidis; Laura Manca; Chris Dijksterhuis; Riender Happee; Joost de Winter

Manuscript Region of Origin: NETHERLANDS

Abstract: A common challenge with processing naturalistic driving data is that humans may need to categorize great volumes of recorded visual information. By means of the online platform CrowdFlower, we investigated the potential of crowdsourcing to categorize driving scene features (i.e., presence of other road users, straight road segments, etc.) at greater scale than a single person or a small team of researchers would be capable of. In total, 200 workers from 46 different countries participated in 1.5 days. Validity and reliability were examined, both with and without embedding researcher generated control questions via the CrowdFlower mechanism known as Gold Test Questions (GTQs).

By employing GTQs, we found significantly more valid (accurate) and reliable (consistent) identification of driving scene items from external workers. Specifically, at a small scale CrowdFlower Job of 48 three-second video segments, an accuracy (i.e., relative to the ratings of a confederate researcher) of 91% on items was found with GTQs compared to 78% without. A difference in bias was found, where without GTQs, external workers returned more false positives than with GTQs. At a larger scale CrowdFlower Job making exclusive use of GTQs, 12,862 three-second video segments were released for annotation. Infeasible (and self-defeating) to check the accuracy of each at this scale, a random subset of 1,012 categorizations was validated and returned similar levels of accuracy (95%).

In the small scale Job, where full video segments were repeated in triplicate, the percentage of unanimous agreement on the items was found significantly more consistent when using GTQs (90%) than without them (65%). Additionally, in the larger scale Job (where a single second of a video segment was overlapped by ratings of three sequentially neighboring segments), a mean unanimity of 94% was obtained with validated-as-correct ratings and 91% with non-validated ratings. Because the video segments

overlapped in full for the small scale Job, and in part for the larger scale Job, it should be noted that such reliability reported here may not be directly comparable. Nonetheless, such results are both indicative of high levels of obtained rating reliability.

Overall, our results provide compelling evidence for CrowdFlower, via use of GTQs, being able to yield more accurate and consistent crowdsourced categorizations of naturalistic driving scene contents than when used without such a control mechanism. Such annotations in such short periods of time present a potentially powerful resource in driving research and driving automation development.

1     **Abstract**

2     A common challenge with processing naturalistic driving data is that humans may need to categorize
3     great volumes of recorded visual information. By means of the online platform CrowdFlower, we
4     investigated the potential of crowdsourcing to categorize driving scene features (i.e., presence of other
5     road users, straight road segments, etc.) at greater scale than a single person or a small team of
6     researchers would be capable of. In total, 200 workers from 46 different countries participated in 1.5
7     days. Validity and reliability were examined, both with and without embedding researcher generated
8     control questions via the CrowdFlower mechanism known as Gold Test Questions (GTQs).

9

10    By employing GTQs, we found significantly more valid (accurate) and reliable (consistent)
11    identification of driving scene items from external workers. Specifically, at a small scale
12    CrowdFlower Job of 48 three-second video segments, an accuracy (i.e., relative to the ratings of a
13    confederate researcher) of 91% on items was found with GTQs compared to 78% without. A
14    difference in bias was found, where without GTQs, external workers returned more false positives
15    than with GTQs. At a larger scale CrowdFlower Job making exclusive use of GTQs, 12,862 three-
16    second video segments were released for annotation. Infeasible (and self-defeating) to check the
17    accuracy of each at this scale, a random subset of 1,012 categorizations was validated and returned
18    similar levels of accuracy (95%).

19

20    In the small scale Job, where full video segments were repeated in triplicate, the percentage of
21    unanimous agreement on the items was found significantly more consistent when using GTQs (90%)
22    than without them (65%). Additionally, in the larger scale Job (where a single second of a video
23    segment was overlapped by ratings of three sequentially neighboring segments), a mean unanimity of
24    94% was obtained with validated-as-correct ratings and 91% with non-validated ratings. Because the
25    video segments overlapped in full for the small scale Job, and in part for the larger scale Job, it should
26    be noted that such reliability reported here may not be directly comparable. Nonetheless, such results
27    are both indicative of high levels of obtained rating reliability.

28

29    Overall, our results provide compelling evidence for CrowdFlower, via use of GTQs, being able to
30    yield more accurate and consistent crowdsourced categorizations of naturalistic driving scene contents
31    than when used without such a control mechanism. Such annotations in such short periods of time
32    present a potentially powerful resource in driving research and driving automation development.

33

Highlights

- Naturalistic driving scene annotation was accomplished with crowdsourcing.
- 12,862 annotations were completed in 1.5 days by 200 external workers from 46 countries.
- Embedded control questions enhance annotation validity and reliability.

1 # VALIDITY AND RELIABILITY OF NATURALISTIC DRIVING SCENE
2 # CATEGORIZATION JUDGMENTS FROM CROWDSOURCING

3 Christopher D. D. Cabrall[1]
4 Mekelweg 2, 2628 CD, Delft, the Netherlands
5 Phone: + 31 (0)152785608 E-mail: c.d.d.cabrall@tudelft.nl

6 Co-authors(s): Zhenji Lu[1], Miltos Kyriakidis[1,2], Laura Manca[1], Chris Dijksterhuis[1,3], Riender Happee[1],
7 Joost de Winter[1]

8 [1]Delft University of Technology; [2]ETH Zurich; [3]Hanze University of Applied Sciences

9

10 **Abstract**

11 A common challenge with processing naturalistic driving data is that humans may need to categorize
12 great volumes of recorded visual information. By means of the online platform CrowdFlower, we
13 investigated the potential of crowdsourcing to categorize driving scene features (i.e., presence of other
14 road users, straight road segments, etc.) at greater scale than a single person or a small team of
15 researchers would be capable of. In total, 200 workers from 46 different countries participated in 1.5
16 days. Validity and reliability were examined, both with and without embedding researcher generated
17 control questions via the CrowdFlower mechanism known as Gold Test Questions (GTQs).
18
19 By employing GTQs, we found significantly more valid (accurate) and reliable (consistent)
20 identification of driving scene items from external workers. Specifically, at a small scale
21 CrowdFlower Job of 48 three-second video segments, an accuracy (i.e., relative to the ratings of a
22 confederate researcher) of 91% on items was found with GTQs compared to 78% without. A
23 difference in bias was found, where without GTQs, external workers returned more false positives
24 than with GTQs. At a larger scale CrowdFlower Job making exclusive use of GTQs, 12,862 three-
25 second video segments were released for annotation. Infeasible (and self-defeating) to check the
26 accuracy of each at this scale, a random subset of 1,012 categorizations was validated and returned
27 similar levels of accuracy (95%).
28
29 In the small scale Job, where full video segments were repeated in triplicate, the percentage of
30 unanimous agreement on the items was found significantly more consistent when using GTQs (90%)
31 than without them (65%). Additionally, in the larger scale Job (where a single second of a video
32 segment was overlapped by ratings of three sequentially neighboring segments), a mean unanimity of
33 94% was obtained with validated-as-correct ratings and 91% with non-validated ratings. Because the
34 video segments overlapped in full for the small scale Job, and in part for the larger scale Job, it should
35 be noted that such reliability reported here may not be directly comparable. Nonetheless, such results
36 are both indicative of high levels of obtained rating reliability .
37
38 Overall, our results provide compelling evidence for CrowdFlower, via use of GTQs, being able to
39 yield more accurate and consistent crowdsourced categorizations of naturalistic driving scene contents
40 than when used without such a control mechanism. Such annotations in such short periods of time
41 present a potentially powerful resource in driving research and driving automation development.

42

# 1. INTRODUCTION

Further knowledge specifically of (background) driving scene contexts could benefit transportation research and ultimately road safety. This study presents and evaluates a new method using crowdsourcing to provide content characterizations of natural driving video footage. Brief descriptions of both topics are provided in the following introductory sections.

## 1.1. Naturalistic driving and driving videos

Naturalistic driving studies (NDS) have been growing in popularity with much success over the last few decades. NDS offer advantages with respect to other traditional driving safety research methods such as eye witness recall (often being inaccurate or unavailable) within crash data evidence approaches and driving simulators (often causing artificial participant behavior) (Regan et al., 2012). However, a lack of experimental control (where extraneous variables except that of manipulative interest are held constant), has been a commonly recognized detriment to NDS. Thus, the accurate annotation of the situational aspects and conditional characteristics that freely vary in NDS becomes all the more important for the identification and understanding of potential causal factors. Augmented by accelerating developments in audio-visual technology, computing, and networking resources, blended research designs are emerging wherein stimuli can be naturally sourced from the real world, reproduced, and mixed with more controlled laboratory conditions.

Due to reductions both in size and costs of cameras, real life driving video is an increasingly accessible data resource that may allow recordings at a large scale and could help enrich other sources of data with otherwise missed contextualized information. However, so much video data might be recorded in naturalistic driving research and field operational tests that research resources are often overwhelmed to process such data libraries through pre-requisite rounds of organization and labeling (e.g., data reduction) towards fuller potentials of use. For example, challenges can arise regarding the availability of confederate researchers for laborious manual annotation or transcription tasks. Unfortunately for driving safety research, the use of real-life driving video footage has remained a relatively low-tapped exception (e.g., Crundall, Underwood, & Chapman, 1999; Chapman et al., 2007; Borowsky, Shinar, & Oron-Gilad, 2010) rather than a common resource, despite inherent strengths in face validity and generalizability of results.

## 1.2. Crowdsourcing

Compared to less than 1% in 1995, about 48% of the world population has an Internet connection to date, placing the approximate number of Internet users in excess of 3.5 billion people (www.InternetLiveStats.com/internet-users/). Online crowdsourcing services make use of this extensive connectivity to create an on-call global workforce to complete large projects in small chunks (a.k.a., micro-task workers). Gosling and Mason (2015) review a broad and growing use of Internet resources in recent psychological research. They conclude that harnessing large, diverse, and real-world data sets presents new opportunities that can increase the societal impact of psychological research. In the automated driving domain, research has recently begun to emerge utilizing crowdsourcing resources through global survey initiatives to capture large scale international public opinion (Bazilinskyy & De Winter, 2015; Kyriakidis, Happee, & De Winter, 2015). In regards to crowdsourcing as a research method, investigation into the differences between laboratory participants versus crowdworkers has found faster responses but higher false alarms with crowdsourcing (Smucker & Jethani, 2011). Additional methodological research has revolved around the assurance of quality from the quick and inexpensive results typically returned by crowdsourcing and have recommended predetermined answer sets for use both in the screening of unethical workers as well as for the effective training of ethical workers (Le et al., 2010; Soleymani & Larson, 2010).

## 1.3. Present study

Real-world driving datasets come with large labor challenges in terms of data reduction like manual annotation and categorization. Pairing together expansive datasets of naturalistic driving video footage with crowdworkers may be a powerful method for progressing driving safety research. As a

94   prototypical example of the power of crowdsourcing, the online platform known as CrowdFlower can
95   accomplish routine categorization work at relatively low cost and at high speed by distributing the
96   work around the world, taking advantage of both differences in time zones and hourly wages.
97   However, such new methods require an investigation of validity and reliability to ensure trustworthy
98   results might still be retained when scaling up beyond a single researcher or small research team. The
99   present study investigated the use of CrowdFlower in the categorization of large amounts of videos
100  with diverse driving scene contents (i.e., presence of another vehicle, straight road segments, etc.)
101  through manipulation of one of its central quality control mechanisms to ascertain the quality and
102  capability of such a method.

103  **2.      METHODS**

104  2.1.     Quality Control Settings

105  Within its documentation, the CrowdFlower system promotes Gold Test Questions (GTQ) as its most
106  important quality control mechanism. By configuring this setting, we enforced that a set of
107  categorizations with known answers (i.e., given by the experimenters) were randomly intermixed with
108  the experimental categorizations of interest. Thresholds of performance on these GTQs were set in an
109  attempt to reduce the amount of indiscriminate responses that may occur within the results due to the
110  remotely distributed nature of work under unsupervised conditions.

111  2.2.     Participants/Workers

112  Participants in this research consisted of external micro-task workers from the online CrowdFlower
113  contributor community. From this network, workers were prescreened by a number of criteria
114  selectable within the CrowdFlower interface. Specifically, within CrowdFlower, performance levels
115  are automatically awarded based on CrowdFlower's criteria of accuracy across a variety of different
116  Job types. We selected a performance setting of Level 2 workers from a three-level scale, representing
117  the midpoint between anchors of "highest speed" (Level 1) and "highest quality" (Level 3). Moreover,
118  across all 51 of its current possible Channels for sourcing external workers (e.g. BitcoinGet,
119  ClixSense, CoinWorker.com, etc.), CrowdFlower was set to include workers only from those retaining
120  a ratio of Trusted to Untrusted Judgments greater or equal to 80% (39 Channels were left toggled on
121  and 12 set to off). All countries were permitted within the Geography setting, and no additional
122  Language Capability requirements were selected.

123  Table 1 lists the countries and source Channels of workers obtained across different sets of
124  categorizations performed within the present study along with distributions of unique worker IP
125  addresses and CrowdFlower worker IDs while Figure 1 depicts the country distribution of the workers.
126  For external crowdworkers, identification of country was determined by CrowdFlower based on IP
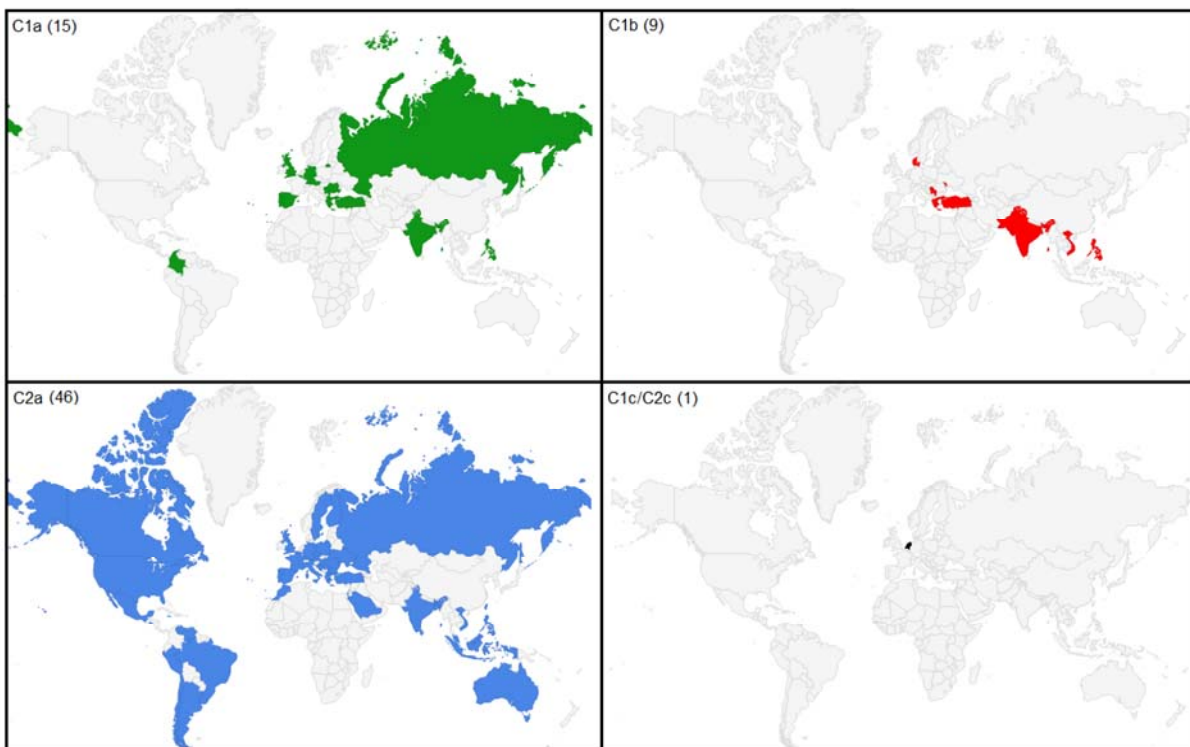127  address.
128
129  Table 1
130  *Overview of the five different sets of categorizations. These sets included differences in the amount of*
131  *video segments to be categorized (C1 = 48 segments, C2 = 12,862 segments), the use of Gold Test*
132  *Questions (C1b had none) and the relation of the annotators to the research (external =*
133  *CrowdFlower workers; internal = confederate research team).*

| Condition | Countries (ISO 3166-1 alpha-3) | Channels | Unique IP's | Unique ID's |
|-----------|-------------------------------|----------|-------------|-------------|
| C1a | 15 = AUT, BEL, COL, DEU, ESP, GBR, GRC, IND, MKD, PHL, PRT, ROU, RUS, SRB, TUR | 5 = clixsense, coinworker, elite, prodege, tremorgames | 18 | 18 |
| C1b | 9 = DNK, GRC, IND, MDA, PAK, PHL, SRB, TUR, VNM | 3 = clixsense, elite, tremorgames | 13 | 13 |

| | | | | |
|---|---|---|---|---|
| C1c | 1= NLD | 1 = n/a (internal) | 1 | 1 |
| C2a | 46 = ARG, AUS, AUT, BEL, BGD, BGR, BIH, BRA, CAN, CHL, CZE, DEU, ESP, FIN, FRA, GBR, GRC, HRV, HUN, IDN, IND, ISR, ITA, JAM, LKA, MAR, MDA, MEX, MKD, MYS, PER, PHL, POL, PRT, ROU, RUS, SAU, SRB, SWE, TUR, TWN, UKR, URY, USA, VEN, VNM | 16 = clixsense, coinworker, fusioncash, gifthulk, hiving, indivillagetest, instagc, personaly, pocketmoneygpt, points2shop, prodege, superrewards, surveymad, tremorgames, yute_jamaica, zoombucks | 247 | 200 |
| C2c | 1 = NLD | n/a (internal) | 12 | 7 |

*Note.* Country abbreviations are according to ISO 3166-1 alpha-3.



*Figure 1*. Annotator country locations by condition.

### 2.3. Apparatus and stimuli

To support projects oriented around the human factors of automated driving (i.e., exposing participants to various HMI/functional research concepts, measuring constructs of vigilance, situation awareness, mental models, reaction time, eye tracking behavior, etc.), a set of stimulus material was desired that had both qualities of high visual realism and controllable levels of uncertainty in repetition, freeze-ability, etc. Initial searches of YouTube with the keyword "dash cam" were conducted to compile a sample database of naturalistic driving video footage. Videos had to feature relatively high and consistent visual quality, a large and consistent field of view, and uninterrupted driving in order to be included. Candidate videos were selected from the search results in order to acquire nominal driving footage (i.e., excluding violations and crashes). We collected a set of 10 freely available YouTube videos ranging between 1 minute and 1 hour duration (but of bimodal typicality of about 3 or 13 minutes length) for a total of 6,934 seconds of driving footage. The countries in which the recordings

4

151 were filmed were not known, but driving was always on the right hand side. Audio was removed from
152 the videos.

153 Subsequently, new self-recorded dash cam driving recordings (6,026 seconds) were filmed in the
154 United States and saved as 39 different files (typically less than 3 minutes in length, but ranging up to
155 15 minutes). This complemented the videos collected from YouTube in order to exhibit a broader
156 range of real-life and experimentally interesting driving situations. These additional recordings
157 included driving at night, on mostly empty desert roads, in a visually complex metropolis, and via
158 multi-lane freeways, as well as at different driving speeds.

159 Driving videos from both sources were uploaded as 49 new private link-only access YouTube videos
160 ($M$ = 264 seconds duration) with an aggregate of 12,960 seconds of near driver point-of-view video
161 footage. Through a combination of MATLAB script and an online tool from www.tech-tipsforall.com
162 (ttfaloopandrepeat.appspot.com), auto-cueing URL links were generated to access each of the 12,862
163 possible 3-second segments from each of these 49 video. These URL links were embedded as text
164 only in our CrowdFlower surveys with one URL per Judgment. The video segments overlapped in a
165 manner such that a randomly selected worker categorized seconds one to three from video 1, another
166 randomly selected worker categorized seconds two to four from video 1, a third randomly selected
167 worker categorized seconds three to five from video 1, etc., for all videos 1 through 49. Example
168 screenshots from the driving video segments are shown in Figures 2a, 2b, and 2c.

169

170

171

172 *Figure 2.* Example screenshots from driving video segments a) recorded from within a publically
173 posted dash cam YouTube video, b) recorded by the experimenters within a visually complex
174 metropolis (i.e., Las Vegas strip), and c) recorded by the experimenters in a visually simple
175 environment (i.e, Nevada desert backroad). Video resolution/quality here is only approximately
176 representative as that initially made available to participants because differences in devices and
177 browsers, full-screen viewing, etc. were not controlled for in the online survey.
178

179 A coding scheme was created wherein each video segment categorization (i.e., Judgment) contained
180 two groups of questions. The first group consisted of 21 checkbox items pertaining to the non-
181 mutually exclusive presence of others, namely, (1) cars/trucks/vans/buses, (2)
182 motorcycles/scooters/mopeds, (3) bicycles, and (4) pedestrians. Each of these four categories
183 contained additional possible sub-specification of their position/direction of travel, namely, (5–8)
184 leading, (9–12) oncoming, (13–16) passing or being passed, and (17–20) crossing; all relative to the
185 present point-of-view vehicle. Additionally, there was a checkbox item which should be ticked for
186 (21) no one else was present.

187 The second group consisted of 10 checkbox items pertaining to presence of miscellaneous
188 infrastructural elements and aspects of vehicle behavior. These were: (1) straight road, (2) more than

6

189 one lane per direction of travel, (3) signs/signals facing the driver, (4) road surface markings other
190 than lane boundaries (e.g., crosswalks, arrows, writing, etc.), (5) lane change by this driver, (6) lane
191 change by another vehicle, (7) turning by this driver, (8) turning by another vehicle, (9) this driver
192 slowing to a stop, and (10) none of the above. In the second round of categorizations (C2, see Tables 1
193 & 2), the coding scheme was extended to include a position/direction item across all road user
194 categories (i.e., of being parked/stationary), plus a miscellaneous item for overt video edits/alterations.
195 Consequently, these extensions (for further data enrichment value) raised the total checkbox count per
196 video segment to 36. The full coding scheme of annotation items (as well as the specific full training
197 instructions given to annotators) is provided in Appendix A.

198 2.4.     GTQ video segments: multiple purposes and representative examples
199 GTQ videos were selected from the full pool of video segments under the criteria to serve as effective
200 screening and training devices. For the purpose of screening indiscriminate respondents, some of the
201 easiest and most unambiguous scenes were selected, as for example a video segment where only an
202 empty desert road is shown (see Example 1).

203 Example 1: https://www.youtube.com/embed/eS79DG08idY?start=12&end=15

204 For the purpose of explicating various annotation labels (e.g., surface paint markings, signage facing
205 the driver), video segments were selected that contained certain items of interest, such as a segment
206 where a railroad crossing sign appears on the side of the road as well as surface markings in the lane of
207 travel (see Example 2)

208 Example 2: https://www.youtube.com/embed/vA5AiKbzIww?start=82&end=85

209 2.5.     Conditions
210 Three different external CrowdFlower Jobs were conducted in two different rounds (C1 and C2), as
211 shown in Table 2. In the first round, C1, a set of 48 unique three-second long video segments
212 (randomly selected from the larger full dataset of collected video footage) were categorized by
213 external CrowdFlower workers with GTQs either turned on (C1a) or turned off (C1b). In C1a and
214 C1b, the default triplicate redundancy setting in CrowdFlower was kept on and so the Job ran until
215 three Judgments were collected for each video segment. Additionally, the same 48 segments were
216 categorized offline by an individual internal worker (i.e., a confederate researcher) in C1c.

217 In the second round, C2, Judgments were performed on CrowdFlower across all 12,862 possible 3-
218 second video segments of the full video dataset via external CrowdFlower workers (C2a) and over a
219 subset of these video segments by an internal worker team comprised of multiple confederate
220 researchers (C2c) using the same CrowdFlower structure as the external workers. Within the C2c
221 round of internal team ratings, one team member accomplished a high volume of Judgments ($n = 638$)
222 under two separate CrowdFlower accounts such that 38 different Judgments of the same driving scene
223 segment from the same person were available to establish intra-rater reliability.
224
225 The required set of Judgments ordered for each CrowdFlower Job was specified at Job launch and
226 included a redundancy option through a multiplier setting (x3 was used in C1, x1 was used in C2).
227
228 Table 2
229 *Categorization conditions*

| Condition | Workers | Video segments categorized | Redundancy | Gold Test Questions | Video segments per Page | Worker payment per Page | Total CrowdFlower Cost |
|---|---|---|---|---|---|---|---|
| C1a | external | 48 | 3 | 12 | 10 | $0.50 | $10.80 |
| C1b | external | 48 | 3 | 0 | 10 | $0.50 | $9.00 |
| C1c | internal | 48 | 1 | 12 | n/a | n/a | n/a |
| C2a | external | 12,862 | 1 | 53 | 11 | $0.25 | $349.32 |
| C2c | internal | 1,012 | 1 | 42 | 11 | n/a | n/a |

7

230     *Note.* The total worker payment differs from the total CrowdFlower costs because CrowdFlower
231     retained a margin of about 20%. Video segments per Page refers to the amount of videos the worker
232     was assigned at a time (i.e., stacked vertically, with a scrollbar); total Pages completed varied between
233     workers. A single Page consisted of 10 (C1) or 11 (C2) Judgments, that is, different driving video
234     segments to be annotated.

235     2.6.     Analyses
236     In the investigation of the utility of CrowdFlower for annotating driving video content, multiple
237     analyses from two different rounds of Jobs (Table 1) were undertaken to cover the separate but related
238     psychometric aspects of validity (i.e., accuracy) as well as reliability (i.e., consistency).

239     In terms of validity, we ascertained to what extent categorizations returned from external
240     CrowdFlower workers reflect what is actually visible in a given driving video segment. At an initial
241     reduced Job scale, the same set of video segments was repeated with and without GTQs (Table 1, C1a
242     vs. C1b) and compared to a reference set of categorizations of these same segments generated by a
243     confederate researcher (C1c). For subsequent accuracy analyses at the greater Job scale (where GTQs
244     were retained), ground truth was created by a team of internal confederates for a random subset due to
245     the infeasibility (and self-defeating purpose) of checking the accuracy of each annotation at this scale.

246     In terms of reliability, we assessed how consistent categorizations of the driving video segments were
247     when repeatedly administered. Supporting this aim, three analyses were conducted. First, from the
248     second round of confederate categorizations (C2c) one internal team member was given a subset to
249     categorize in duplicate to himself (i.e., randomly intermixed among his other categorizations, see 2.5
250     Conditions). Second, at the small scale Job (C1), each video segment was rated by three different
251     external CrowdFlower workers (both in C1a and in C1b). Third, the full dataset categorizations of C2a
252     provided an account of consistency due to the fact that the video segments overlapped such that any
253     second of driving video footage was categorized three times. That is, for any second "x" bounded by
254     start/end points [start, end] there existed a first segment: [x, x+2], a second segment: [x−1, x+1], and a
255     third segment: [x−2, x].

256     2.7.     Procedure
257     All workers were provided with a set of instructions and examples regarding the driving video
258     segment categorization coding scheme that remained available for consultation throughout their work
259     (Appendix A). A single Judgment consisted of a set of 31 (C1) or 36 (C2) checkboxes pertaining to
260     features visible within a randomly selected 3-second long driving video segment (Section 2.3). A
261     single Page consisted of 10 (C1) or 11 (C2) Judgments, that is, different driving video segments to be
262     annotated.

263     In the conditions where GTQs were active (C1a, C2a, C2c), task workers were first given a single
264     page of Quiz Mode GTQs Judgments to complete. Because of constraints of CrowdFlower, a GTQ
265     Judgment had to be answered perfectly in order to be scored as correct, with no partial credit given
266     (i.e., all 31 or 36 checkboxes had to be checked correctly against predetermined answers constructed
267     by the experimenters). If workers achieved a threshold correctness Trust Score on these GTQs of 70%
268     [i.e., 7 out of 10 Judgements] in C1, and 25% [i.e., 3 out of 11 Judgments] in C2, then workers were
269     automatically allowed by CrowdFlower to continue through as many more Pages of Work Mode as
270     they would like. Through trial and error, the set threshold was lowered from 70% in C1 to 25% in C2,
271     because it turned out to be often highly difficult to obtain a perfect answer on each of the checkboxes
272     of a Judgment. Additionally, in C2, participants were supported with further detailed feedback
273     explaining the correct answers. For an incorrect answer to any checkbox item of a GTQ during Quiz
274     Mode, workers were shown the correct answers of all checkboxes for that Judgment along with a brief
275     justification. Each Page of Work Mode had one new not-yet-seen GTQ randomly presented within the
276     other Judgments such that a worker was unable to identify which Judgments had a priori answers that
277     their own answers would be scored against. As long as workers maintained a running average Trust

278 Score above the set threshold (i.e., 70% in C1, 25% in C2), and there were still GTQs remaining that
279 they had not yet seen, they were allowed to continue.

280 In the CrowdFlower condition without GTQs (C1b), workers were allowed to enter Work Mode
281 straightaway without real-time screening criteria barring them from submitting Judgments. On a first-
282 come-first-serve (optionally screened) basis, Jobs in CrowdFlower are run until a pre-determined
283 amount of Judgments are completed by an indeterminate amount of workers.

284 In summary, the GTQ condition included further screening and training to enhance the responses of
285 task workers than the condition without GTQs.

286 **3. RESULTS**
287 The utility of the crowdsourcing platform CrowdFlower in the content categorization of naturalistic
288 driving video footage was investigated through multiple analyses concerning both validity and
289 reliability. Overall, the supposed utility of CrowdFlower in the present tasks was found to be
290 supported (see Table 3). Results were indicative of significantly increased utility both in terms of
291 validity and reliability in the presence of GTQs as compared to without GTQs. Results were obtained
292 both in the preliminary round of a reduced scale (C1: 48 video segments) and in the subsequent round
293 conducted at a larger scale (C2: 12,862 video segments).

294

295 Table 3
296 *Summary of analyses*

| Section | Analysis aim | Relative Job size | Analysis outcome |
|---|---|---|---|
| 3.1.1 | Validity | Small | The GTQ condition yielded more accurate Judgments than the No GTQs condition. Accuracy was assessed by using the Judgments of a single internal confederate rater as ground truth. |
| 3.1.2 | Validity | Large | The GTQ condition yielded accurate Judgments. Accuracy was assessed by using the Judgments of a small team of internal confederate raters as ground truth. |
| 3.2.1 | Reliability | Small | A single confederate rater was found to be consistent to himself. |
| 3.2.2 | Reliability | Small | The GTQ condition yielded more consistent Judgments than the No GTQ condition, for full Judgments and at the item level. |
| 3.2.3 | Reliability | Large | The GTQ condition yielded Judgments of high inter-rater consistency for overlapping video segments. Consistency was assessed for known-to-be-accurate Judgments. |
| 3.2.4 | Reliability | Large | The GTQ condition yielded high inter-rater consistency for overlapping video segments. Consistency was assessed for unknown-to-be-accurate Judgments. |

297

298 3.1. Validity
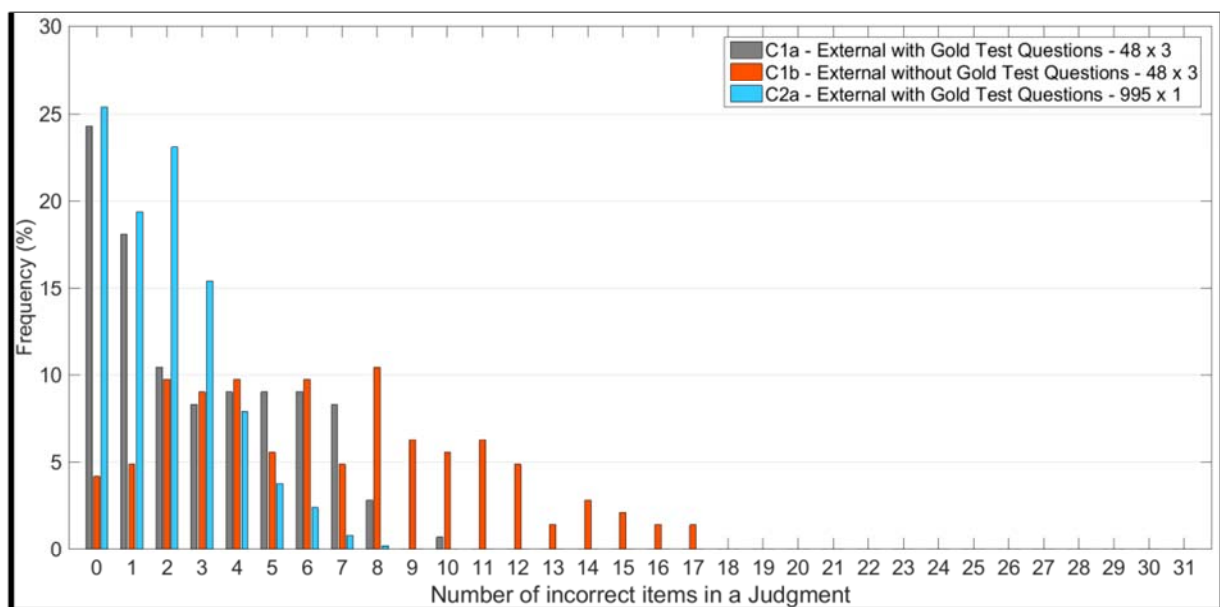
299 3.1.1. 48 Judgments, comparing GTQ with no GTQ
300 Results showed that there were 35 of 144 (24%) and 6 of 144 (4%) exact matches from C1a (with
301 GTQs) and C1b (without GTQs) respectively, relative to C1c (taken as a measure of ground truth).
302 Results thus indicated inaccuracies in the Judgments from both C1a and C1b (Fig. 3).

303 However, these inaccuracies occurred in different specificity/sensitivity biases. Phi correlation
304 coefficients were computed between each full Judgment (i.e., an array of 31 binary checkboxes) from

9

a condition (C1a or C1b) against the ground-truth Judgment returned by an internal confederate rater (C1c) matched for a specific video segment. The median across all 144 (48 x 3) correlation coefficients of the GTQ condition (C1a; $r = 0.78$) was significantly higher than for the No GTQ condition C1b ($r = 0.39$) (Mann-Whitney $U = 3756$, n1 = n2 = 144, $p < 0.001$ two tailed). Furthermore, greater total item accuracy across all 4,464 (31 x 48 x 3) categorized items was found in C1a (4,051 = 91%) than in C1b (3,504 = 78%).

Among the 4,464 categorized items in C1b (i.e., without GTQs), there were 396 false positives (i.e., items marked present but which were absent in the video segment according to the confederate researcher), yielding a false positive rate of 11% (396/3,519). Furthermore, there were 564 misses (i.e., items marked absent that were present in the video segment according to the confederate researcher), yielding a miss rate of 60% (564/945). In C1a (with GTQs), the false positive rate was 1.6% (57/3,519) and the miss rate was 38% (356/945). In other words, GTQs contributed to a reduction of both false positives and false negatives.



*Figure 3.* Distribution of the number of errors per Judgment at the smaller C1 Job scale of 144 Judgments (with and without GTQs) and for a subset of 995 Judgments from the larger C2 Job scale (with GTQs). Errors were determined against known answers (C1c or C2c). A score of 0 signifies a perfectly correct Judgment.

### 3.1.2.    1,012 Judgments, comparing external versus internal workers

The confederate research team (C2c) performed 995 Judgments of video segments (17 video segments were removed due to video playback errors) which were randomly selected from C2a. Results showed that there were 257 (26%) exact matches between the Judgments from C2a and C2c. Phi correlations with the ground truth for both the smaller scale Job (correlation between C1a and C1c: median $r = 0.78$, see also Section 3.2.1) and the larger scale Job (correlation between C2a and C2c: median $r = 0.80$) were not found to significantly differ (Mann-Whitney $U = 65298.5$, n1 = 144, n2 = 995, $p = 0.083$).

From the 35,820 C2a items re-rated within C2c (995 Judgments x 36 items per Judgment) the false positive rate was 2.1% (682/31,564) and the miss rate was 27.6% (1,176/4,256).

## 3.2. Reliability

### 3.2.1. 38 Judgments, comparing confederate to himself

In condition C2c, one confederate performed 638 Judgments about evenly split under two different CrowdFlower accounts, with an approximate 10% subset of his Judgments from each account coded in duplicate ($n = 38$). Intra-individual test-retest reliability results for this same rater using the same software settings but across different sessions were: 34 (89%) exact matches, an average phi correlation of 0.98 across the 38 Judgments, and an overall item accuracy of 99.5% (i.e., 1,361 out of 1,368).

### 3.2.2. 48 Judgments, comparing GTQ versus no GTQ

During C1a and C1b, each video segment collected three external worker Judgments and so allowed for a consistency measure of how many categorization ratings (both for full Judgments and/or across items within Judgments) were returned identically between external CrowdFlower task workers. Unanimous agreement on all 31 items of a Judgement was found in 7 of 48 Judgments in C1a (with GTQs) and in 1 of 48 Judgments in C1b (without GTQs). Per item, the unanimous agreement percentage across the 48 Judgments was computed, and was found to be significantly higher for C1a ($M = 90\%$, $SD = 13$) than for C1b ($M = 65\%$, $SD = 19$, $n1 = n2 = 31$, $t(60) = 5.85$, $p < 0.001$).

### 3.2.3. 257 Judgments, comparing ratings by unanimous voting

For the correct 257 Judgments in C2 (see Section 3.1.2), a reliability analysis was conducted by comparing overlapping categorizations across sequential seconds of video footage. For example, the correct true/false answer provided for an item in a video segment that began at time $x$, was compared with the answer received for that same item by another external worker whose video segment began at time $x-1$ and additionally by another external worker whose video segment began at time $x-2$. It should be noted that some variation between overlapping video segments would be expected to exist (e.g., a car seen only in the last second of a segment that starts at $x = 0$ might not be visible in the previous videos $x-1$ and $x-2$). Due to such uncertainty, somewhat less than perfect reliability may be expected even from perfectly reliable raters. This necessitates consideration of proportional consistency analysis across the entire array of 36 items contained within a Judgment. In other words, it is assumed that while one or a few aspects might vary between overlapping videos, the majority of aspects should remain the same.

Results showed that 74 of 257 correct Judgments (29%) received the same true/false rating across all 36 items by three different external workers who rated overlapping video segments. Figure 4 shows a distribution of the 257 Judgments according to the number of items yielding unanimous agreement. Judgments always had more than two-thirds (i.e., at least 25 out of 36 items) unanimous agreement, and the mean number of items yielding unanimous agreement was 33.9 out of a possible 36.

*Figure 4*. Frequency of validated (i.e., 257 fully correct) and all returned Judgments (originally 12,862) from C2a according to number of items yielding unanimous agreement from three independent raters.

### 3.2.4.   12,862 Judgments, comparing ratings by unanimous voting

For all 12,862 Judgments, a reliability analysis of unanimous answers was conducted with overlapping sequential seconds again as in Section 3.2.3, but now for the full dataset. The first and last two Judgments of each video required removal due to a logical lack of full overlap, resulting in a total of 12,670 Judgments ($12,862 - 4 \times 48$).

Regarding unanimity of full Judgments, 1,129 of 12,670 answers (9%) received the same true/false value across all 36 items by the three different external workers. The mean number of items with unanimous agreement per Judgment was 32.6 out of 36 possible.

The distributions of Judgments in Figure 4 shows that disagreement existed in the categorizations of overlapping sequential seconds of video footage; this occurred most frequently for two items.

## 4.      DISCUSSION, CONCLUSIONS AND RECOMMENDATIONS

The CrowdFlower crowdsourcing platform may present great potential for driving research by bringing task workers from across the world to categorize a rapidly growing resource of naturalistic driving video data. Due to its inherently distributed structure, CrowdFlower and online tools of similar kind may be more susceptible to fraudulent or non-discriminating responses as compared to locally administered and more tightly controlled traditional methods. Specifically, the utility of CrowdFlower with (and without) its self-purported most important quality control mechanism of GTQs was investigated in the objective categorization of driving video contents via binary presence/absence flagging of pre-specified driving items of interest both at a preliminary reduced and a subsequently increased Job scale.

Exhibiting credible signs of validity and reliability (Table 3), the potential for the method of crowdsourcing the categorization of driving video contents can be considered in a meaningful and valuable way. For example, as a result of our settings in the present study, 12,862 CrowdFlower annotation categorizations were completed in about one and a half days by 200 external workers from 46 different countries working at an hourly rate of 1.09 USD each (total cost of about 349.32 USD inclusive of a 20% transaction fee) with an average of 75 seconds per Judgment. Through volunteer confederate collaboration, 1,002 annotation categorizations were completed in about two weeks by six internal confederate workers from the Netherlands working between/around their other work duties at

12

404 a conservative estimated hourly rate around $20.25 USD each (total cost estimate of about $394.54
405 with an average of 70 seconds per Judgment). Thus, for the same approximate costs, the external
406 workers returned categorizations about ten times faster.
407
408 Several limitations exist within the present study and are worth mentioning. The first and foremost, is
409 that the GTQ mechanism is explicitly designed to work with objective tasks where there are clear and
410 definable right and wrong answers and so it may not be suitable for many otherwise desirable
411 subjective judgments from a distributed task worker network. A GTQ is constructed in CrowdFlower
412 to require pre-defined correct answers with as minimal ambiguity as possible as well as detailed and
413 documentable justification/motivation of that answer (similar to how both annotator screening and
414 training is used in more controlled laboratory experiments). It should be noted that the design of the
415 present study does not lend itself towards some other research questions that might be addressed from
416 pairing crowdsourcing to naturalistic driving data for example for purposes of investigating the
417 general human ability in perception/annotation of various aspects of driving scenes (inter-item
418 research questions) and/or the bearing of universal/local driving cultures on driving scene
419 interpretation (inter-cultural research questions). Instead, the present study aimed to eliminate
420 ambiguities on an equal par between conditions to test the principle manipulation of interest: the use or
421 not of GTQs.

422 Nonetheless, some of our requested annotation items appear to have contributed to some confusion
423 between some raters. The worst three annotation items, both in terms of accuracy and reliability,
424 pertained to identification of fully straight roads, signage/signals facing the driver, and number of
425 lanes per direction of travel. Overall, performance with these items averaged around 63% (reliability)
426 and 79% (accuracy) compared to averages taken across all the remaining items of 93% (reliability) and
427 96% (accuracy). Without proper hypotheses/controls in place, we cannot propose these as particularly
428 systematic nor meaningful results in human perception or suitability to crowdsourcing beyond our
429 own inabilities to more thoroughly formulate such desired details for our driving video data library
430 into more fully objective definitions/terms (see Appendix A). For example, while relative decreases in
431 miss rates were obtained through use of GTQs, the absolute levels of miss rates (38% and 28%, in C1a
432 and C2a respectively) might be indicative of annotation items requiring further scrutiny and/or ease in
433 task criteria definition. Our annotation task contained a combination of both demanding visual search
434 and items with low ground truth base rates. Thus, it would be logical or even possibly more natural for
435 a rater to adopt a conservative strategy when faced with annotation uncertainty (i.e., not checking a
436 box unless they have explicitly seen something). Relatedly, the high miss rates may reflect a bias due
437 to the fact that all items were by default unchecked (absent) requiring checking as needed, rather than
438 being checked (present) requiring unchecking as needed. Indeed, complexities in universal
439 instructions, clear coding rule descriptions, and controlled balancing of default absence/presence
440 question valences could be a relevant concern in crowdsourcing annotations from large, diverse, and
441 remote participant populations without local remediation of a real-time physically present
442 experimenter. However, it should be noted that we did not use any CrowdFlower geography/language
443 settings and thus kept this aspect equally random across our external worker conditions so as not to
444 confound our relative evaluations regarding potential benefits of GTQs.

445 Secondly, the specific items of the coding scheme created and used in the present study may be
446 challenged further than issues of clarity towards aspects of organization and inter-item independence.
447 The item checkboxes within a Judgment were pre-tested and arranged by probable frequencies of
448 occurrence such that categorization speeds might benefit from predictable and likely emergent patterns
449 of responses. Thus, the repetitive and non-random ordering of items may be a source of bias towards
450 consistency (although, again it should be noted that the same structure was presented to both GTQ and
451 non-GTQ condition groups).

13

452   Lastly, several dependency relations existed between items which may degrade the power of some of
453   the analyses of the present study. For example, several items pertained to the identification of object
454   classes (cars, motorcycles, bicycles, and pedestrians, respectively) that upon selection, each expanded
455   with sub-item location information (i.e., leading, oncoming, passing, crossing, parking). For cases
456   where only one object from the class was present, the sub-item location information thus became
457   mutually exclusive rather than independent. As another example, items pertaining to actions of other
458   vehicles such as "Lane change by another vehicle" and "Turning on/off between this and any other
459   road by another vehicle" logically depend on presence of another vehicle and thus retain relations to
460   ratings of item vehicle class identification.

461   More traditional and established methods for interrater reliability (e.g. Cohen's/Fleiss' kappa) were
462   not pursued. The reason for that is the difficulty of determining a chance agreement for our Judgments
463   that contained a composite of yes/no decisions with inter-item dependencies as described above.
464   Instead, simpler measures of consistency, such as the phi coefficient and the proportion of unanimous
465   Judgments, were used. Further studies with CrowdFlower more specific to questions of validity and
466   reliability might limit such complexities in advance, sacrificing some annotation meaning in favor of
467   stricter control, standard analyses, and afforded reflection regarding the broader annotation literature.
468   Additionally, further assessments of the ground truth reliability of our internal rating team (beyond the
469   single rater repetitions of the analysis in 3.2.1) would be desirable in future work. For now, the
470   reliability agreements observed in our approach (Fig. 4/5) appear qualitatively consistent with levels
471   from previous image annotation work (Nowak & Ruger, 2010; containing 53 annotations per image
472   across a set of 99 without presuming the existence of two persons that annotated the whole set of
473   images). Specifically, in comparison to the average identical accuracy they obtained of 0.906,
474   following their Equation 2, we computed our own average unanimous annotation accuracies
475   respectively as 0.941 (section 3.2.3, Fig. 4) and 0.906 (section 3.2.4, Fig. 5).

476   Multiple ethical and privacy concerns can be raised in consideration of methods that employ
477   crowdworkers with human annotation of naturalistic driving video data. Some of these may not be
478   new and include attempting to anonymize video data in the sense that specific combinations of
479   sensitive information are not presented in combination to result in personably identifiable information
480   from both aspects of the drive (time, date, location, etc.) along with aspects of driver identity (name,
481   face, home/work address, etc.). A major difference between the present method and the classical way
482   of annotating naturalistic driving data is that in the present method the task is outsourced to
483   crowdworkers who are themselves anonymous and residing in different countries, while in the
484   classical way the annotation is done by trained team members who are typically local and
485   known/approved by the principal investigator(s). Aside from the annotation integrity
486   (accuracy/consistency) concerns specifically addressed in the experimental design and results of the
487   present study, other new challenges are worth discussing such as legal requirements of the handling of
488   data. In the present study, the video data were obtained from public sources, which is uncommon
489   within traditional NDS approaches. Thus, any terms and conditions regarding data sharing, ownership,
490   and viewership restrictions put in place a priori by the responsible parties would need to be considered
491   and respected so as not to be violated. Additionally, the regulations and policies pertaining to the
492   online reproduction/distribution of (video) data specific to each country or online hosting community
493   should be adhered to, and this includes the presentation of potentially disturbing images such as might
494   be the case with automobile crashes/accidents or illegal driving behavior.

495   A few positive privacy points regarding the present method are interesting to consider as well.
496   Because the annotating work is distributed across many crowdworkers in distal locations, a relatively
497   small amount of the total data is restrictively released to single/isolated persons at a time. For example,
498   in the present study, only random 3-second clips from randomly different drives and randomly
499   different drivers were distributed. Accordingly, it becomes much less likely that a crowdworker can
500   come to recognize a driver's travel patterns or other aspects that may pose risks to privacy. This

501  compares favorably in contrast to a classical annotation perspective where a single or smaller group of
502  annotators may more likely become familiar with the travel patterns contained within the data.
503  Additionally, the present study does not propose to share all data (e.g., geospecific, CANBUS, etc.) as
504  may be accessible to classical annotators in naturalistic research but to selectively distribute only
505  pieces of the full dataset (i.e., herein only video annotation was outsourced and only that of forward
506  facing cameras from public roads where filming is allowed). Lastly, crowdworkers themselves are
507  employed under certain terms of service to which they must accept and abide (e.g.,
508  https://www.crowdflower.com/legal/). If crowdworkers were to violate such terms (e.g., share
509  proprietary data) they would be subject to consequences not limited to but including the likes of losing
510  their worker privileges such as payment, membership, etc.

511  An increasing amount of real-life driving videos are being recorded both within naturalistic driving
512  studies as well as from public channels of user generated content. For example, at the start of
513  conducting the current research, there were approximately 795,000 returns for the term "dashcam" on
514  YouTube (November 19, 2015). Upon presenting this work at the international conference for Road
515  Safety on Five Continents (May 19, 2016), there were 1.13 million returns for the same search (i.e.,
516  +42% increase in about half a year), and by the time of manuscript revisions (August 8, 2017), a total
517  of 4.26 million were available (i.e., +436% increase in less than 2 years' time). Categorizing such
518  expansive data sets can be a costly and time-consuming manual process. One solution is to train
519  automated algorithms to conduct coding tasks such as in machine learning and classification.
520  However, such algorithms themselves often require some diligently pre-labeled examples for their
521  own accuracy and only through diverse training sets may overcome common challenges of overfitting.
522  Under the correct circumstances (e.g., open-access data) and quality control settings (i.e., the
523  construction and use of GTQs), Crowdsourcing tools like CrowdFlower appear to have the potential
524  for delivering equivalent accuracy and reliability utility as locally trained humans. It is therefore
525  recommended that future driving research and ultimately driving safety itself might benefit from
526  exploiting increasingly large scale and publically available data sets through embracing and
527  channeling a growing global pool of human resources.

531

532  REFERENCES

533  Bazilinskyy, P., & De Winter, J. C. F. (2015). Auditory interfaces in automated driving: An
534  international survey. *PeerJ Computer Science.* 1:e13.

535  Borowsky, A., Shinar, D., & Oron-Gilad, T. (2010). Age, skill, and hazard perception in driving.
536  *Accident Analysis and Prevention, 42*, 1240-1249.

537  Chapman, P., Van Loon, E., Trawley, S., & Crundall, D. (2007). A comparison of drivers' eye
538  movements in filmed and simulated dangerous driving situations. *Behavioral research in road safety,*
539  *Seventeenth seminar.* London: Department for Transport.

540  Crundall, D., Underwood, G., & Chapman, P. (1999). Driving experience and the functional field of
541  view. *Perception, 28*, 1075-1087.

542  Gosling, S., & Mason, W. (2015). Internet research in psychology. *Annual Review of Psychology, 66,*
543  877-902.

544  Kyriakidis, M., Happee, R., & De Winter, J. (2015). Public opinion on automated driving: Results of
545  an international questionnaire among 5,000 respondents. *Transportation Research Part F: Traffic*
546  *Psychology and Behaviour, 32*, 127-140.

547  Le, J., Edmonds, A., Hester, V., & Biewald, L. (2010). Ensuring quality in crowdsourced search
548  relevance evaluation: The effects of training question distribution. *Proceedings of the SIGIR 2010*
549  *Workshop on Crowdsourcing for Search Evaluation (CSE 2010).*

550  Nowak, S., & Ruger, S. (2010). How reliable are annotations via crowdsourcing: A study about inter-
551  annotator agreement for multi-label image annotation. *Proceedings of the International Conference on*
552  *Multimedia Information Retrieval*, ACM, 557-566.

553  Regan, M., Williamson, A., Grzebieta, R., & Tao, L. (2012). Naturalistic driving studies: literature
554  review and planning for the Australian naturalistic driving study. *Australasian College of Road Safety*
555  *Conference*, Sydney, New South Wales, Australia.

556  Smucker, M., & Jethani, C. (2011). The Crowd vs. the Lab: A comparison of crowd-sourced and
557  university laboratory participant behavior. *Proceedings of the SIGIR 2011 Workshop on*
558  *Crowdsourcing for Information Retrieval (CIR 2011).*

559  Soleymani, M., & Larson, M. (2010). Crowdsourcing for affective annotation of video: Development
560  of a viewer-reported boredom corpus. In V. Carvalho, M. Lease, & E. Yilmaz (Eds.). *Proceedings of*
561  *the SIGIR 2010 Workshop on Crowdsourcing for Search Evaluation (CSE 2010).*

562  University of Ottawa (2014). Validity and Reliability of Measurements: 1. Core knowledge. Retrieved
563  from http://www.med.uottawa.ca/sim/data/Measurement_validity.htm

564  **APPENDIX A (DOUBLE-CLICK TO OPEN)**

Adobe Acrobat
Document

565

# Watching Very Short Videos And Marking Yes/no Checkboxes Confeds

**Instructions ▲**

## Overview

Watch short driving video segments (3 seconds each) and select the elements that can be seen in the video segment. This should be a simple and easy objective job task and not much if any subjectivity, matter of opinion, or "thinking" involved at all.

We are trying to train a computer to recognize things in these videos and match like driving situations to each other but first we need to label what is actually in them so we are trying to use crowdsourced human eyes to objectively say what is there or not.

Each video categorization has been previously measured to take people about 1.5 minutes or less on average.

## We provide

URL links to specific driving video segments ("Play Video" button)

Lists of true/false items (checkboxes)

## Process/Procedure

Watch and review the video segment provided (3 seconds each) FULL SCREEN IS RECOMMENDED. Replay and pause video as often as you need. This is NOT a memory test. Check all that apply.

**For the first half of the questions (1 through 7)**, we would like categorizations of what other vehicles/pedestrians are present within that driving video segment based on 5 different possible locations (directions of travel) relative to the driver whose vehicle the video was filmed from. **Click on the check boxes to denote various vehicle/pedestrian presence by its location/position/direction-of-travel**.
a) in front of and in the same lane
b) traveling in an oncoming/opposite direction
c) traveling in the same direction whether alongside or ahead
d) traveling at any different angles
e) parked, parking, un-parking

**For the second half of the questions (8a through 8k)**, we would like categorizations of other miscellaneous aspects within that same video segment. **Click on the checkboxes to indicate which elements are "Obviously Visible" within that video segment.** This may include things the vehicles of the videos do not actually reach or complete, but which you don't have to squint for under a magnifying glass.

Please note:

 - "8b ... Just straight road (no angles/bends/curves in the entirety of the visible road of travel)" applies not only to the portion of the road driven but also that which is "Obviously Visible" anywhere within that video segment including the road ahead.

- Lane changes or turns don't have to be complete to count.

- Categorizations should span the full 3 seconds of the segment (watch out for things that are there in the first moments even if they shortly disappear due to motion of the video)

***Disclaimer***: Unfortunately some people cheat by clicking randomly or using computer programs to complete jobs. Please ***ALWAYS*** check the check box to confirm you are a diligent human contributor (it is located at the end of the "other vehicles/pedestrians" section). Also, please ***NEVER*** check any boxes that ask you to leave them empty, off, or unchecked (these occur at the beginning of each section and in the middle of the first section underneath group 3). Don't worry, these will be very very obvious! However, if too many are missed you will be excluded and not paid, so please no random clicking!

## Steps

### a) check all that apply ...



http://ttfaloopandrepeat.appspot.com/showVideo.html?st=90&et=93&vld=5HiykkcjruA&l=yes&lnf=1000000&ap=no

**Play Video** — **Step 1** Click button to open URL in new tab/window

**Step 2** Watch/repeat 3 second long driving video segment

☐ ***NEVER check this checkbox: leave it unchecked/empty/off***

☐ 1) I can see one or more Cars/Trucks/Vans/Buses in this video segment ...

☐ 2) I can see one or more Motorcycles/ScootersMopeds in this video segment ...

☐ 3) ***NEVER check any checkboxes in group 3: leave this one and its sub parts all unchecked/empty/off***

☐ 4) I can see one or more Bicycles in this video segment ...

☐ 5) I can see one or more Pedestrians in this video segment ...

☑ 6) No one else present, this driver is alone

☑ 7) ***ALWAYS check this box to confirm you are a diligent human contributor***

**Step 3**
Check all that apply regarding the obviously visible aspects of this driving video segment
- Don't check boxes that say "NEVER check ..."
- Be sure to check the confirm box that says "ALWAYS check ..."

8) Which elements are contained in THIS driving video segment?
☐ 8a ... ***NEVER check this checkbox: leave it unchecked/empty/off***
☐ 8b ... Just straight road (no bends/curves visible)
☐ 8c ... More than one lane per direction of travel
☐ 8d ... Any signs/signals facing driver (road signs, billboards, traffic lights, building names, ads, etc.)
☐ 8e ... Painted communication on any visible road surface (includes crosswalks, arrows, etc. but NOT lane boundary/edge info)
☐ 8f ... Lane change by this driver
☐ 8g ... Lane change by another vehicle
☐ 8h ... Turning on/off between this and any other road by THIS driver
☐ 8h ... Turning on/off between this and any other road by another vehicle
☐ 8i ... This driver slowing to a stop, being stopped, or pulling away from a stop
☐ 8j ... Editing alterations in the video file (discontinuity, added text, pauses, slow motion, sped up sections, etc.)
☑ 8k ... None of these misc. elements are present in this video segment

### b) check all that apply, including sub statements as they appear...

http://ttfaloopandrepeat.appspot.com/showVideo.html?st=120&et=123&vId=BAZirVeEEN0&l=yes&Inf=1000000&ap=no

Play Video

**Step 2**
Watch/repeat 3 second long driving video segment

☐ ***NEVER check this checkbox: leave it unchecked/empty/off***

☑ 1) I can see one or more Cars/Trucks/Vans/Buses in this video segment ...

☐ 2) I can see one or more Motorcycles/ScootersMopeds in this video segment ...

☐ 3) ***NEVER check any checkboxes in group 3: leave this one and its sub parts all unchecked/empty/off***

☐ 4) I can see one or more Bicycles in this video segment ...

☑ 5) I can see one or more Pedestrians in this video segment ...

☐ 6) No one else present, this driver is alone

☑ 7) ***ALWAYS check this box to confirm you are a diligent human contributor***

**8) Which elements are contained in THIS driving video segment?**
☐ 8a ... ***NEVER check this checkbox: leave it unchecked/empty/off***
☑ 8b ... Just straight road (no bends/curves visible)
☐ 8c ... More than one lane per direction of travel
☑ 8d ... Any signs/signals facing driver (road signs, billboards, traffic lights, building names, ads, etc.)
☐ 8e ... Painted communication on any visible road surface (includes crosswalks, arrows, etc. but NOT lane boundary/edge info)
☐ 8f ... Lane change by this driver
☐ 8g ... Lane change by another vehicle
☐ 8h ... Turning on/off between this and any other road by THIS driver
☐ 8h ... Turning on/off between this and any other road by another vehicle
☐ 8i ... This driver slowing to a stop, being stopped, or pulling away from a stop
☐ 8j ... Editing alterations in the video file (discontinuity, added text, pauses, slow motion, sped up sections, etc.)
☐ 8k ... None of these misc. elements are present in this video segment

**Step 3**
Check all that apply regarding the obviously visible aspects of this driving video segment
- Don't check boxes that say "NEVER check ..."
- Be sure to check the confirm box that says "ALWAYS check ..."
- Answer all sub statements as they appear

**1)**
☑ 1a ... they are leading (same lane ahead)
☑ 1b ... they are oncoming (traveling opposite direction)
☐ 1c ... they are passing or being passed (traveling same direction)
☐ 1d ... they are crossing (traveling on a road that intersects)
☑ 1e ... they are parked (parking, or un-parking/pulling out)

**5)**
☐ 5a ... they are leading (same lane ahead)
☐ 5b ... they are oncoming (traveling opposite direction)
☑ 5c ... they are passing or being passed (traveling same direction)
☐ 5d ... they are crossing (traveling on a road that intersects)
☐ 5e ... they are parked (parking, or un-parking/pulling out)

## Tips

The URL link may appear broken across multiple lines, so be sure to copy/paste the entire URL link all the way from "http://ttfa…" and ending with "…&" if needed

Use "full screen" to see the video in a larger view.

**If you miss any in QUIZ mode**, please read the provided answers and reasoning carefully as it should help clarify what we mean/expect by our various categorization items. Also refer to Extended Sample Set below as needed for visual examples of various items.

## Thank You!

Your help on this task is greatly appreciated!

## Extended Sample Set (reference as needed)
Note: more than one item may apply within the same video segment (**always check all that apply**). In these examples, we have used highlights only to show some of the possible clues or cues that should help you know to mark the current specific check box item as true.

**1a)** Car/Van/Truck/Bus ... they are traveling in the same direction in the same lane ahead (leading)

**1b)** Car/Van/Truck/Bus ... they are traveling in the opposite direction (oncoming)

**1c)** Car/Van/Truck/Bus ... they are traveling in the same direction (passing, being passed, pass-able)
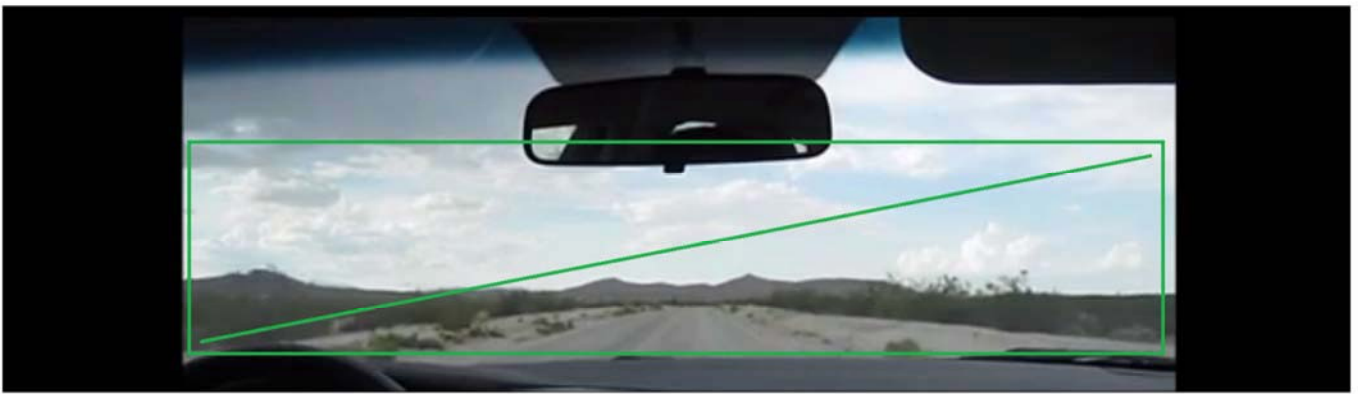




**1d)** Car/Van/Truck/Bus ... they are traveling on a road that intersects (crossing)
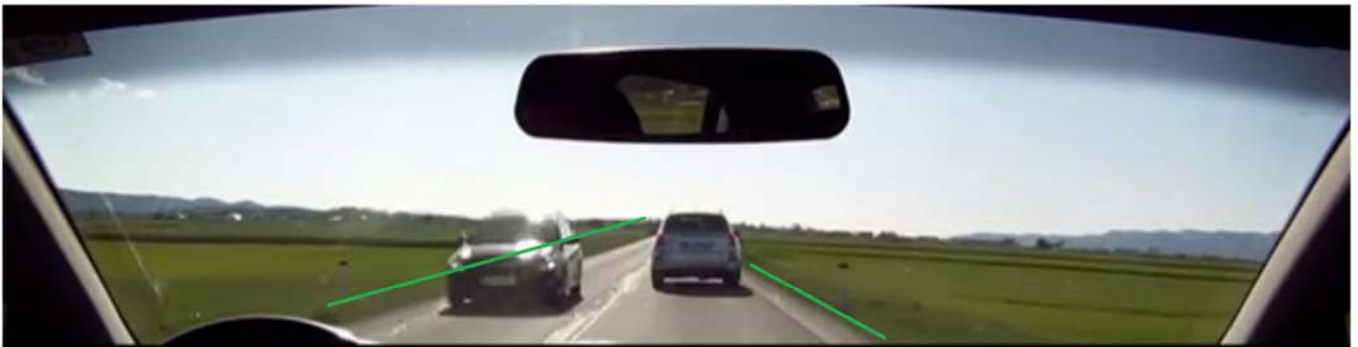
**1e)** Car/Van/Truck/Bus ... they are parked, parking, pulling out/unparking



**2e)** Motorcycles/Scooters/Mopeds ... they are parked, parking, pulling out/unparking

**4c)** Bicycles ... they are traveling in the same direction (passing, being passed, pass-able)





**5b)** Pedestrians ... they are traveling in the opposite direction (oncoming)

**5c)** Pedestrians ... they are traveling in the same direction (passing, being passed, pass-able)

**5d)** Pedestrians ... they are traveling on a road that intersects (crossing)



**5e)** Pedestrians ... they are standing still or otherwise stationary





**6)** None of the above. No external vehicle/pedestrian present. This driver is alone.

**8B) ...** Just straight road (no angles/bends/curves in the entirety of the visible road of travel)





**8C) ...** More than one lane per either direction of travel
*Please note that some lane division markings may vary (white/yellow) between different roads. However, please **mark 8c only if there is little to no ambiguity in the multi lane situation.** For example, immediately above in the images of "8b)" at different times with other vehicles present or not, or even with and without a visible dividing lane line at all, it might be hard to tell if it should be considered either a single lane of travel in both directions (8c would not apply = false) or two lanes in the direction of travel of the driver (8c would apply = true). Please assume that the later second to be a very **RARE CASE** and you might only expect it when there is another parallel roadway for the **other** direction of*

*travel (i.e. a "divided highway" situation). Furthermore if the road has no lane division markings then it has no "lanes" and so 8c = false. Please, reserve the marking of 8c = true for clearly obvious and unambiguous multiple lanes per either direction of travel (see below with lanes counted out in green numbers ascending from the edge to the center of the road separately per direction of travel).*





**8D)** ... Any signs/signals facing driver (road signs, billboards, traffic lights, building names, ads, etc.)

*Note: if you can make out colors, text, symbols, pictures, etc. on the sign/signal then for our purposes here it is "facing the driver" even if it is not 100% legible; the only ones you don't count are those that are just the backs of signs (e.g. facing the opposite direction, oncoming traffic).*

**8E)** ... Painted communication on any visible road surface (includes crosswalks, arrows, etc. but NOT lane boundary/edge info

**8F)** ... Lane change by this driver



**8G)** ... Lane change by another vehicle

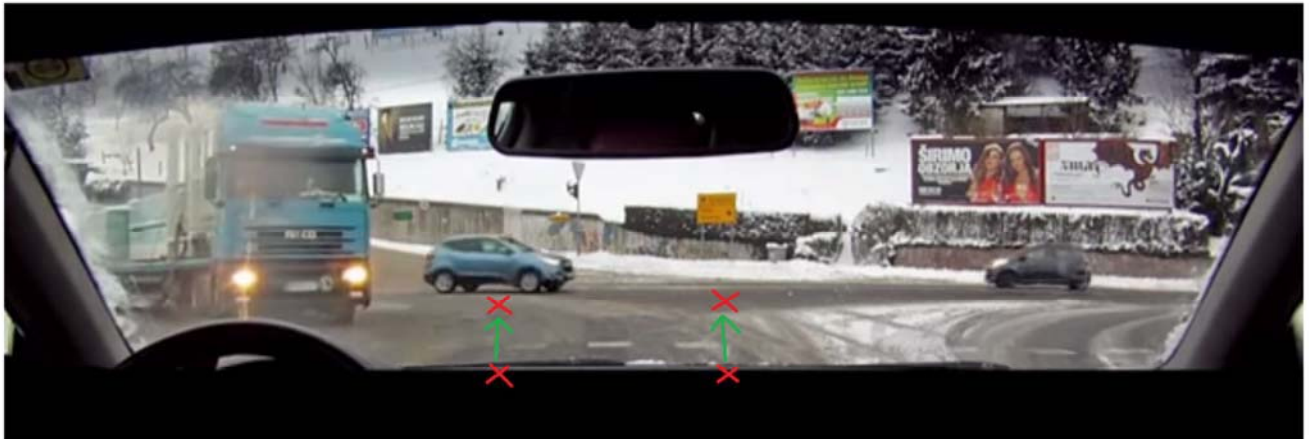**8H)** ... Turning on/off between this and any other road by THIS driver

**8I) ... Turning on/off between this and any other road by another vehicle**





**8J) ... This driver is slowing to a stop, is stopped, or pulling away from a stop.**
*Note: The red crosses here are meant as possible places you might notice deceleration or other motion patterns indicative of this item. The green circles are other possible/probable*
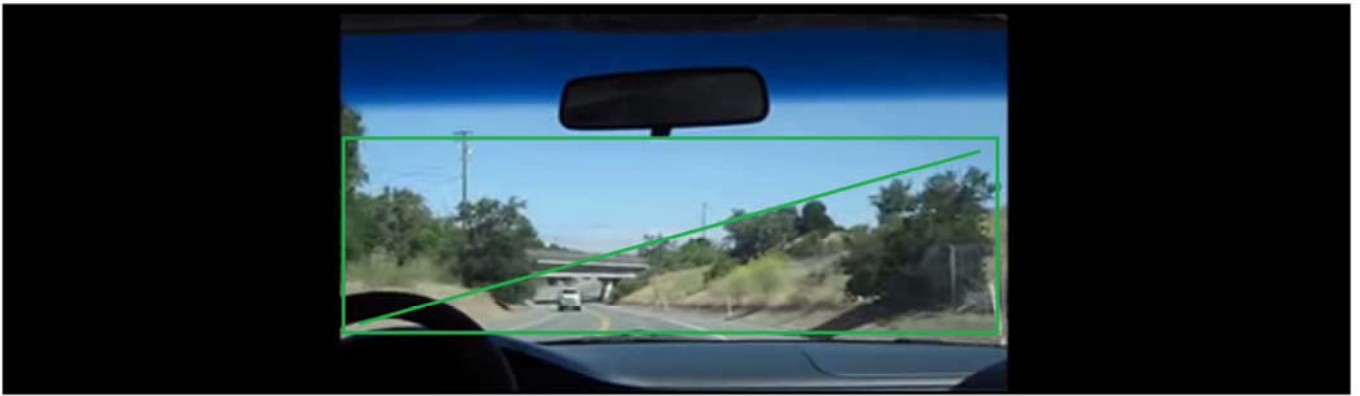
*clues of stopping contexts in common driving situations.*





**8K)** ... Editing alterations in the video file (discontinuity, added text, pauses, slow motion, sped up sections, etc.)



**8L)** ... None of these miscellaneous elements are present in this video segment

**Thank you.**

http://ttfaloopandrepeat.appspot.com/showVideo.html?
st=90&et=93&vId=5HiykkcjruA&l=yes&lnf=1000000&ap=no

Play Video (http://ttfaloopandrepeat.appspot.com/showVideo.html?
st=90&et=93&vId=5HiykkcjruA&l=yes&lnf=1000000&ap=no)

☐ ***NEVER check this checkbox: leave it unchecked/empty/off***

☑ 1) I can see one or more Cars/Trucks/Vans/Buses in this video segment ...

**1)**

☐ 1a ... they are traveling in the same direction in the same lane ahead (leading)
☐ 1b ... they are traveling in the opposite direction (oncoming)
☐ 1c ... they are traveling in the same direction (passing, being passed, pass-able)
☐ 1d ... they are traveling on a road that intersects (crossing)
☐ 1e ... they are parked (parking, or un-parking/pulling out)

☑ 2) I can see one or more Motorcycles/ScootersMopeds in this video segment ...

**2)**

☐ 2a ... they are traveling in the same direction in the same lane ahead (leading)
☐ 2b ... they are traveling in the opposite direction (oncoming)
☐ 2c ... they are traveling in the same direction (passing, being passed, pass-able)
☐ 2d ... they are traveling on a road that intersects (crossing)
☐ 2e ... they are parked (parking, or un-parking/pulling out)

☐ 3) ***NEVER check any checkboxes in group 3: leave this one and its sub parts all

unchecked/empty/off***

☐ 4) I can see one or more Bicycles in this video segment ...

☐ 5) I can see one or more Pedestrians in this video segment ...

☐ 6) None of the above. No external vehicle/pedestrian present. This driver is alone

☑ 7) ***ALWAYS check this box to confirm you are a diligent human contributor***

## 8) Which elements are contained in THIS driving video segment?

☐ 8a ... ***NEVER check this checkbox: leave it unchecked/empty/off***
☐ 8b ... Just straight road (no bends/curves in the entirety of the visible road of travel)
☐ 8c ... More than one lane per either direction of travel
☐ 8d ... Any signs/signals facing driver (road signs, billboards, traffic lights, building names, ads, etc.)
☐ 8e ... Painted communication on any visible road surface (includes crosswalks, arrows, etc. but NOT lane boundary/edge info)
☐ 8f ... Lane change by this driver
☐ 8g ... Lane change by another vehicle
☐ 8h ... Turning on/off between this and any other road by THIS driver
☐ 8i ... Turning on/off between this and any other road by another vehicle
☐ 8j ... This driver slowing to a stop, is stopped, or pulling away from a stop
☐ 8k ... Editing alterations in the video file (discontinuity, added text, pauses, slow motion, sped up sections, etc.)
☐ 8l ... None of these miscellaneous elements are present in this video segment

## Comments?

[ ]

http://ttfaloopandrepeat.appspot.com/showVideo.html?
st=99&et=102&vId=Ge0a27WR6WI&l=yes&lnf=1000000&ap=no

Play Video (http://ttfaloopandrepeat.appspot.com/showVideo.html?
st=99&et=102&vId=Ge0a27WR6WI&l=yes&lnf=1000000&ap=no)

☐ ***NEVER check this checkbox: leave it unchecked/empty/off***