

# Assessing the reliability and validity of pictorial-assisted 24-hour recall for measuring hand hygiene and child faeces disposal: a cross-sectional study in Malawi

Olivier Rizk<sup>1\*\*</sup>, Sarah Bick<sup>1\*\*</sup>, Blessings White<sup>2</sup>, Kondwani Chidziwisano<sup>2</sup>, Robert Dreibelbis<sup>1\*</sup>

1. Department of Disease Control, London School of Hygiene and Tropical Medicine, London, United Kingdom
2. WASHTED Centre, Malawi University of Business and Applied Sciences, Blantyre, Malawi

\*Corresponding author, [Robert.Dreibelbis@lshtm.ac.uk](mailto:Robert.Dreibelbis@lshtm.ac.uk)

\*\*These authors contributed equally to the work

Abbreviations<sup>1</sup>

---

<sup>1</sup>Abbreviations: P24hR – pictorial 24-hour recall; HWWs – handwashing with soap; TA – traditional authority; NGO – non-governmental organisation; LOA – limits of agreement; PPV – positive predictive value; NPV – negative predictive value

# 1 Abstract

2 Whilst improving hygiene and sanitation behaviours is key to cost-effective and sustainable  
3 WASH interventions, measuring behaviour change remains a challenge. This study  
4 assessed the validity and reliability of pictorial 24-hour recall (P24hR), a novel method using  
5 unprompted recall of past activities through pictures, compared to structured observation for  
6 measuring handwashing with soap (HWWS) and safe child faeces disposal in rural Malawi.  
7 Data were collected from 88 individuals across 74 households in Chiradzulu district using  
8 both methods over a two-day period, with the recall period of the P24hR corresponding to  
9 the period of structured observation completed the previous day. Results showed poor  
10 agreement between P24hR and observations in detection of hygiene opportunities and  
11 behaviours. P24hR under-reported handwashing opportunities when frequency was high  
12 and over-reported them when frequency was low. The 95% limits of agreement for  
13 handwashing opportunities estimated through Bland-Altman analysis (-7.62 to 4.89) were  
14 unacceptably wide given median 5 opportunities observed per participant. P24hR also over-  
15 reported HWWS and safe child faeces disposal, and kappa statistics indicated agreement no  
16 better than by chance. Structured observation remains the preferred method for measuring  
17 hygiene behaviours despite its known limitations, including potential reactivity bias.

# 18 Keywords

19 Measurement, validity, handwashing, pictorial assisted recall, hygiene

20

## 21 Introduction

22 Interventions to improve hand hygiene in domestic settings are associated with a 30%  
 23 reduction in diarrheal diseases among children under the age of five (Wolf et al., 2022) and a  
 24 17% reduction in acute respiratory infections (Ross et al., 2023). Our estimates of the  
 25 potential health benefits of hygiene interventions, however, are associated with exposure to  
 26 – rather than adoption of – hand hygiene interventions. This is because measuring hygiene  
 27 behaviour remains a challenge (Egreteau, 2017; Schmidt et al., 2019), with few validated  
 28 and reliable methods for measuring behaviours available. Contaminated hands are a critical  
 29 pathway for exposure to a range of environmentally transmitted pathogens (Wagner et al.,  
 30 1958) and quantifying and measuring hand hygiene behaviour is a key part of exposure risk  
 31 assessment (Kwong et al., 2020), intervention design and evaluation (Amon-Tanoh et al.,  
 32 2021), and understanding individual and population-level health risks (Wolf et al., 2019).

33 Among the methods used to measure behaviour, structured observation is often considered  
 34 the gold standard due to its ability to measure behaviours as they occur (Biran et al., 2008;  
 35 Schmidt et al., 2019). However, it is resource-intensive and can be seen as intrusive or  
 36 inappropriate for certain behaviours. Most importantly, direct observation can result in  
 37 reactivity from participants, in which case the validity of estimates is limited (Ram et al.,  
 38 2010).

39 Proxy measures – or indirect measures of behaviour – can also be used in hygiene and  
 40 sanitation research. They are often operationalized as the presence of necessary materials  
 41 or infrastructure to enable a specific behaviour (Schmidt et al., 2019). For example, the  
 42 presence of a handwashing facility is used as a proxy measure for hand hygiene behaviour  
 43 (Biran et al., 2008; Joint Monitoring Program, 2022) based on global estimates suggesting  
 44 individuals are almost two times as likely to wash hands with soap after faecal contact  
 45 events when both soap and water are available. Proxy measures are convenient as data can  
 46 be rapidly collected and at a low cost, but their accuracy may be limited (Biran et al., 2008;  
 47 Briceño et al., 2014).

48 Self-reporting tools are commonly used to measure behaviour. These tools are inexpensive,  
49 quick and require little expertise to put in place or use (Schmidt et al., 2019). However, self-  
50 reported hygiene and sanitation behaviours are often unreliable due to biases, including  
51 recall bias. Several studies have shown poor agreement between reported and observed  
52 hygiene behaviour (Chidziwisano et al., 2020; Curtis et al., 1993; Manun'Ebo et al., 1997;  
53 Stanton et al., 1987).

54 Pictorial 24h recall (P24hR) has been suggested as a novel methods to measure hygiene  
55 and sanitation behaviours (Schmidt et al., 2019). P24hR measures behaviours through  
56 facilitated recall of past activities with pictures and a diary sheet. P24hR is a validated  
57 method to measure dietary intake, with photos and pictures assumed to increase the  
58 accuracy of reporting compared to unfacilitated recall (Lazarte et al., 2012; National Cancer  
59 Institute, 2023). P24hR has been used to evaluate various handwashing interventions  
60 (Tidwell et al., 2019; Tidwell et al., 2020). In a study in India (Schmidt et al., 2019),  
61 researchers found that P24hR of handwashing behaviour was more closely aligned with  
62 direct observation data than reported hand hygiene. However, comparisons were made  
63 between two different study groups rather than compared among the same individuals; more  
64 information on the validity and reliability of P24hR is needed to further assess its utility in  
65 measuring hygiene behaviours.

66 The aim of this study was to assess the agreement of P24hR compared to direct observation  
67 for pre-selected hygiene and sanitation behaviours and determine the validity and reliability  
68 of P24hR. By measuring the same set of behaviours with different methods, we provided  
69 useful information regarding the measurement properties of P24hR compared to structured  
70 observation.

# 71 **Materials and methods**

72 This field-based cross-sectional study compared the agreement between measured  
73 prevalence of hand washing with soap (HWWS) and safe child faeces disposal practices in a  
74 sample of rural households in Chiradzulu district, Malawi. Target behaviours were measured  
75 using both structured observations and P24hR in the same participants over a two-day  
76 period, with the recall period of the P24hR corresponding to the period of direct observation  
77 completed the previous day.

## 78 **Study setting and sampling**

79 This study was conducted as part of a larger research and learning collaboration between  
80 the London School of Hygiene and Tropical Medicine and the Malawi University of Business  
81 and Applied Sciences. Chiradzulu is situated in the southern region of Malawi and is sub-  
82 divided into 8 Traditional Authorities (TA) (Figure 1). Villages included in this study were  
83 selected from a roster of villages present in TA-Likoswe and TA-Mpama. Both TAs were part  
84 of a community-based sanitation promotion programme implemented by the NGOs World  
85 Vision and Water For People the year before data collection.

86 We aimed to enrol approximately 75 households based on the range typically used for  
87 agreement studies (Han et al., 2022). Sampling was completed in 13 villages across the  
88 study area, six from TA-Likoswe and seven from TA-Mpama, with each village contributing  
89 six households to the final sample size. Within villages, we approached every sixth  
90 household from a pre-defined starting point. The household inclusion criterion was the  
91 presence of a child under five years of age at the time of the observations, to ensure the  
92 possibility of observing child faeces disposal practices.

93 Within each household, we recruited up to 3 individuals to participate in study activities.  
94 Participants were adult (over 18 years old) residents of the household. In instances where  
95 households contained more than three adult residents present, the adults contributing most  
96 to childcare and household activities were selected in priority.

## 97 **Tool development and implementation**

98 Detailed tool development is described in Appendix A. In brief, a list of daily activities was  
 99 developed and adapted to local context, resulting in a list of 41 discrete activities. Each  
 100 activity was translated into a pictorial image reflecting that specific activity (Appendix B).  
 101 Daily routines were organized around 5 temporal periods – early morning (waking until  
 102 breakfast), morning (breakfast through lunch), afternoon (after lunch until sunset), evening  
 103 (sunset until the evening meal) and night (evening meal until bedtime). Pilot testing found  
 104 that participants were able to map their reported activities to using the activity cards to the  
 105 organized daily diary.

## 106 **Data collection**

107 Data collection consisted of 8 staff who had prior experience with direct observation and  
 108 water, sanitation and hygiene research. Field teams were organized into two teams: six  
 109 observers and two enumerators conducting P24hR (henceforth ‘interviewers’). Observers  
 110 were different from interviewers to reduce the risk of bias. The six observers were female, as  
 111 it was easier for them to be allowed in homes where females were mostly present. After  
 112 obtaining approval by village chiefs and collecting appropriate consent from participating  
 113 household members, each household was visited twice over a two-day consecutive period to  
 114 conduct direct observations (day 1) and P24hR (day 2).

115 Observations lasted six hours and began in the morning (around 7:30 am), when most  
 116 household activities took place in the study population. Participants were all observed at the  
 117 same time. Observers would generally sit in the yard, where many activities take place. If a  
 118 participant left the household, observers would remain with the other participants.

119 Opportunities for handwashing and their associated behaviour were recorded for all  
 120 participants. Handwashing opportunities were pre-defined as: after going to the toilet, after  
 121 taking children for defecation, after cleaning children after defecation, after disposing of  
 122 children’s faeces, before washing food, before preparing food, before serving food, after  
 123 tending to animals (Appendix A). For each opportunity, observers could record one of the  
 124 following hand hygiene activities: no handwashing; handwashing with ash, mud or soil;

handwashing with water only or handwashing with soap. Observers also recorded any child defecation event and the faeces disposal method: in the latrines, buried, in the open or in the garbage. Each of the six observers visited one household per day for a total of six households observed per day.

The administration of P24hR was completed the next day and took on average 20-30 minutes per participant. Participants were introduced to the 41 pictures and the diary sheet, explained how to use them to describe their activities in the past 24 hours, and then given time to complete the diary sheet independently. Interviewers would help if participants had difficulty identifying pictures or time periods. After the diary sheet was completed, interviewers would go through the participants' day, one activity at a time, asking if they had forgotten anything. Finally, interviewers would manually record each activity and take a picture of the completed diary sheet. Each of the two interviewers visited three households per day for a total of six households per day.

The data from direct observations, pictorial 24h recall and household surveys were recorded on Android tablets with forms produced using the online platform Kobo Toolbox. Data were encrypted and uploaded daily to a secured server.

## **Statistical analysis**

All statistical analyses were conducted in Stata version 18 (StataCorp, College Station, TX, USA). The data collected during observations and P24hR were matched for each participant using a unique identifying code and the corresponding 6-hour observation period was isolated within the 24 hours of pictorial recall data for comparison (early morning and morning). Observed and reported opportunities and behaviours were extracted for each participant, for both HWWS and safe child faeces disposal.

The number of handwashing opportunities was a count variable totalling all opportunities for handwashing defined above. Handwashing opportunities that occurred in rapid succession in either the observation or P24hR data were treated as a single hand hygiene opportunity, for example 'Washing food' immediately followed by 'Preparing food'. The number of HWWS

events associated with an opportunity was originally constructed as a count variable, however, given the low rates of HWWS in both methods, we constructed a binary variable of any recorded HWWS associated with a handwashing opportunity during the period of interest.

Due to the low number of child defecation events, the count of events was converted into a binary variable representing any child defecation event during the period of interest. The binary variable of safe child faeces disposal was also defined as any safe disposal practices following child defecation as defined by the World Health Organization (buried or disposed in latrines) reported or observed during the period of interest (World Health Organization and United Nations Children's Fund, 2006).

Inter-method agreement for the count outcome (number of handwashing opportunities, modelled as a continuous variable) was evaluated using the Bland-Altman method (Bland and Altman, 1986). Bland-Altman analyses plot the differences in values obtained by two methods against the respective mean values. The mean difference between the two methods, referred to as the bias, indicates the extent to which the methods diverge. The standard deviation of the bias is used to estimate limits of agreement (LOA), which act as a reference interval between which 95% of the data should lie. An advantage of Bland-Altman plots is that they allow to simultaneously assess reliability and validity of the methods relatively to each other (Montenij et al., 2016) and standard approaches are recommended for when data violate distributional assumptions (Bland and Altman, 1999).

Inter-method agreement for binary variables was evaluated using kappa statistics (Cohen, 1960). This method measures agreement between two methods compared to expected agreement by chance alone. Kappa statistics below 0 indicate agreement worse than by chance; values equal to 0 indicate agreement no better than chance, and values between 0 and 1 reflective of increasing agreement (Landis and Koch, 1977). Additionally, results were analysed using McNemar's test to assess the symmetry in performances between the two



178 methods based on marginal totals, providing an estimate or over- or under-reporting (Curtis  
179 et al., 1993; Manun'Ebo et al., 1997; McNemar, 1947).

180 Using direct observation at the reference group, we also compared the sensitivity and  
181 specificity of P24hR methods. True positives were defined as target behaviours reported by  
182 both P24hR and direct observations; target behaviours reported by P24hR but not observed  
183 were considered false positives. True negatives were defined as the absence of target  
184 behaviour in both P24hR and observation data; target behaviours not reported by P24hR but  
185 capturing during observations were classified as false negative. Sensitivity, specificity, and  
186 positive and negative predictive values were calculated to compare the two methods (Guitart  
187 et al., 2021; Trevethan, 2017).

# 188 **Ethical considerations**

189 This study was approved by National Committee on Research in the Social sciences and  
190 Humanities in Malawi (Protocol No. P.01/23/718) as well as the Ethical Review Committee at  
191 the London School of Hygiene and Tropical Medicine (LSHTM MSc Ethics Ref: 28743).  
192 Informed consent was obtained in all households before beginning direct observations and  
193 confirmed either through their signature or a thumbprint if the participant was illiterate. In the  
194 case of an illiterate participant, the presence of a literate individual co-signing as an  
195 independent witness was also required.

## Results

### Sample Characteristics

In total, 88 individuals across 74 households participated in both structured observations and pictorial 24h recalls. Selected characteristics are presented in Table 1. In some of the smaller villages (<35 households), finding six households with a child under five was not always possible. Due to the time limitations, seven households without children under five were included to meet sample size requirements.

### Measurement of handwashing opportunities and behaviours

P24hR detected 412 (median 4 per participant) total handwashing opportunities compared to 531 in structured observations (median 5 per participant) (Table 2). Differences between the two methods in counts of total opportunities per participant were consistent with a normal distribution (Figure 2;  $p = 0.20$ ).

Using a classical BA approach, bias was -1.36 (95% confidence interval (95%CI): -2.04, -0.69), indicating under-reporting by P24hR, while the limits of agreement (LOA) extended from -7.62 to 4.89. Testing for the required assumptions for a classical analysis revealed that 7/88 observations (8.0%) laid beyond LOA and proportional bias was present as shown in Figure 3. Consequently, LOA were calculated using standard regression methods (Figure 3). For low average values, P24hR over-reported opportunities for handwashing, while for high average values, the method under-reported opportunities. P24hR was more precise when the average of opportunities was low compared to high averages as indicated by narrower LOA and data points closer to the line of equality.

Handwashing with soap was observed at 7 of the 531 opportunities (1.3%) while participants reported HWWS at 29/412 (7%) of opportunities identified in the P24hR (Table 2). Kappa statistic of presence of any HWWS was close to zero, indicating agreement no better than by chance (Table 3). Due to low rates of observed behaviour, a binary variable of any reported or observed HWWS was created for each participant. McNemar's test of this binary variable

222 gave strong evidence that the marginal prevalence of HWWS at any key moment differed  
223 between the two methods.

224 Using structured observation as the reference group, sensitivity of P24hR was low for  
225 HWWS (14%), while specificity was much higher (75%). This resulted in a very low PPV of  
226 P24hR compared to direct observation but high NPV (Table 3).

## 227 **Measurement of child faeces disposal opportunities and behaviours**

228 P24hR detected 16 total opportunities for child faeces disposal compared to 6 in structured  
 229 observations, and safe disposal was recorded at all but one of these opportunities for each  
 230 method (Table 2). Similarly to HWWS, the kappa statistics of the presence of any  
 231 opportunities for child faeces disposal and presence of any safe child faeces disposal were  
 232 both close to zero, indicating agreement no better than by chance, and a similar pattern of  
 233 sensitivity and specificity resulted in a very low PPV (6.7% and 7.1%) and high NPV (93%  
 234 and 95%) (Table 3).

## Discussion

This study evaluated the agreement between P24hR and structured observations and provides estimates of the reliability and validity of pictorial 24h recall as a novel method to measure hygiene and sanitation behaviours. Our findings suggest that P24hR has low agreement with direct observation, resulting in under-reporting of high frequency events, such as opportunities for handwashing, and over-reporting of “proper” or socially desirable behaviours, such as HWWS and safe child faeces disposal. Over-reporting by a self-reporting tool like P24hR is consistent with results from previous studies in Malawi (Chidziwisano et al., 2020) and other parts of the world (Curtis et al., 1993; Manun'Ebo et al., 1997). We found that P24hR tended to over-report handwashing opportunities when the average number of opportunities measured between methods was low, and under-report opportunities when the average number was high. This biphasic relationship illustrates that P24hR is a blunt instrument. The 95% LOA calculated (-7.62 to 4.89) are unacceptably wide considering the median 5 opportunities per participant observed. The high NPV suggests that P24hR is better suited for assessing the absence of specific behaviours rather than their presence, although the conceptual and practical utility of this may be limited.

Pictorial recall has been used in various other fields of research. In the field of nutrition, images representing different food groups and portion sizes have been widely used to facilitate dietary recall (Bulungu et al., 2021). Various tools have been validated to measure dietary diversity (Bulungu et al., 2021) and intake (Bulungu et al., 2021; Lazarte et al., 2012). Pictorial-assisted recall has been found to have high agreement with other methods of measuring time use in low resource settings (Masuda et al., 2014). In water, sanitation, and hygiene research, pictorial aids have been used as a way to facilitate recall data, for example to measure water use (Esrey et al., 1992; Wright et al., 2006) or in daily diaries to measure diarrhoea episodes (Rego et al., 2021; Wright et al., 2006), but their measurement properties have not been fully evaluated.

Our study aligns with previous research that demonstrated that alternative methods for collecting self-reported hand hygiene behaviours are also subject to over-reporting. After adjusting for confounders, Schmidt and colleagues found that pictorial assisted estimates of HWWS were 13 percentage points higher than measuring through direct observation in a similar study population and 24 percentage points higher for post-defecation HWWS (Schmidt et al., 2019). However, differences between pictorial assisted recall compared to observation was smaller than the difference between traditional self-report and observations. In Ethiopia, Cotzen et al. (Contzen et al., 2015) compared covert script-based methods, in which respondents describe the sequence of actions between two events, to direct observation and traditional self-report methods to direct observation. While covert-script based methods had a higher correlation with observed behaviours than traditional self-report, they still over-estimated behaviours by 16 – 22 percentage points. Our study was not intended to compare P24hR against traditional self-reported behaviour; however, P24hR's poor performance against structured observation by a variety of measures in this study makes any potential improvement against self-report of limited utility.

The strength of this study is that observations and P24hR were conducted on the same individuals only 24 hours apart, enabling direct comparison of two methods for measuring behaviour over the same approximate time period. A limitation of this study was the difficulty in accurately identifying the 6-hour observation period in 24h recall data. Despite collecting additional information to facilitate matching, some cut-off points had to be decided subjectively which may have resulted in misclassification of reported behaviours occurring before or after the time periods covered in the structured observations. The use of independent raters could be beneficial when isolating observation periods in recall data as well as measuring outcomes. Second, the schedule of data collection required P24hR to take place the day after observations, which lead to twelve participants being lost to follow-up. Given the poor performance of P24hR compared to structured observations in our analysis, it is unlikely that these 12 observations would have significantly improved the performance of P24HR. Third, child faeces disposal was rarely observed, which meant

assessments were done using very few data points. Additionally, the high prevalence of null values for child faeces disposal and HWWS prevented the analysis of results through the Bland-Altman method. While the transformation of outcomes into categorical variables still permitted a relevant analysis of the data (Green, 2021), future tool evaluations could use negative binomial regression instead, as used by Schmidt et al., to prevent this issue (Schmidt et al., 2019). Finally, this study used structured observations as a reference. Observations have certain limitations, especially reactivity which could lead participants to wash their hands more than usual in the presence of an observer. However, observers are capable of precisely recording series of events and the timeframe in which they occur unlike P24hR. This means that opportunities and behaviours can be measured with less uncertainty.

## Conclusions

This study assessed the potential of pictorial 24h recall as a novel method to measure hygiene and sanitation behaviour for future evaluations of WASH interventions. Overall, agreement with structured observation was poor: P24hR tended to under-report hygiene opportunities and over-report socially desirable, “correct” behavioural outcomes. The negative predictive value of P24hR was high, although the conceptual and practical utility of this may be limited. While structured observations remain both time and resource intensive and may still result in biases, they remain the best method for measuring hygiene behaviours.

## **Declarations**

### **Acknowledgements**

The authors would like to thank the participants of our study for giving generously of their time. We also thank the field data teams who participated in both tool development and final data collection.

### **Author contributions**

Conceptualization: KC, RD; Design: OR, SB, BW, KC, RD; Acquisition of data: OR, BW, Analysis and Interpretation: OR, SB, RD; First draft of the manuscript: OR; Review and editing: SB, BW, KC, RD

### **Funding sources**

OR was the recipient of a travel grant from LSHTM to support travel and field data collection costs. No other specific funding supported this work.

### **Declaration of interests**

The authors declare no competing interests.



## References

- 309 Amon-Tanoh, M.A., McCambridge, J., Blon, P.K., Kouame, H.A., Nguipdop-Djomo, P., Biran,  
310 A., Cousens, S., 2021. Effects of a social norm-based handwashing intervention including  
311 handwashing stations, and a handwashing station-only intervention on handwashing with  
312 soap in urban Cote d'Ivoire: a cluster randomised controlled trial. *Lancet Glob Health* 9,  
313 e1707-e1718, [https://doi.org/10.1016/S2214-109X\(21\)00387-9](https://doi.org/10.1016/S2214-109X(21)00387-9)
- 314 Biran, A., Rabie, T., Schmidt, W., Juvekar, S., Hirve, S., Curtis, V., 2008. Comparing the  
315 performance of indicators of hand-washing practices in rural Indian households. *Trop Med*  
316 *Int Health* 13, 278-285, <https://doi.org/10.1111/j.1365-3156.2007.02001.x>
- 317 Bland, J.M., Altman, D.G., 1986. Statistical methods for assessing agreement between two  
318 methods of clinical measurement. *Lancet* 1, 307-310,
- 319 Bland, J.M., Altman, D.G., 1999. Measuring agreement in method comparison studies. *Stat*  
320 *Methods Med Res* 8, 135-160, <https://doi.org/10.1177/096228029900800204>
- 321 Briceño, B., Colford, J.M., Gertler, P.J., Arnold, B.F., Chase, C., Sahli, M.W., Vidal, A.O.,  
322 Ram, P.K., 2014. Validity of rapid measures of hand washing behavior : an analysis of data  
323 from multiple impact evaluations in the global scaling up hand washing project.
- 324 Bulungu, A.L.S., Palla, L., Priebe, J., Forsythe, L., Katic, P., Varley, G., Galinda, B.D., Sarah,  
325 N., Namboozee, J., Wellard, K., Ferguson, E.L., 2021. Validation of a life-logging wearable  
326 camera method and the 24-h diet recall method for assessing maternal and child dietary  
327 diversity. *British Journal of Nutrition* 125, 1299-1309,  
328 <https://doi.org/10.1017/S0007114520003530>
- 329 Chidziwisano, K., Tilley, E., Morse, T., 2020. Self-Reported Versus Observed Measures:  
330 Validation of Child Caregiver Food Hygiene Practices in Rural Malawi. *International Journal*  
331 *of Environmental Research and Public Health* 17, 4498,  
332 <https://doi.org/10.3390/ijerph17124498>
- 333 Cohen, J., 1960. A Coefficient of Agreement for Nominal Scales. *Educational and*  
334 *Psychological Measurement* 20, 37-46, <https://doi.org/10.1177/001316446002000104>
- 335 Contzen, N., De Pasquale, S., Mosler, H.J., 2015. Over-Reporting in Handwashing Self-  
336 Reports: Potential Explanatory Factors and Alternative Measurements. *PLoS One* 10,  
337 e0136445, <https://doi.org/10.1371/journal.pone.0136445>
- 338 Curtis, V., Cousens, S., Mertens, T., Traore, E., Kanki, B., Diallo, I., 1993. Structured  
339 observations of hygiene behaviours in Burkina Faso: validity, variability, and utility. *Bull*  
340 *World Health Organ* 71, 23-32,
- 341 D-Maps.com, 2023. Malawi.
- 342 Egreteau, D., 2017. Behaviour Change determinants, the key to successful WASH  
343 strategies. *Humanitarian Aid on the Move* 14, 28-34,
- 344 Esrey, S.A., Habicht, J.P., Casella, G., 1992. The complementary effect of latrines and  
345 increased water usage on the growth of infants in rural Lesotho. *Am J Epidemiol* 135, 659-  
346 666, <https://doi.org/10.1093/oxfordjournals.aje.a116345>
- 347 Green, J.A., 2021. Too many zeros and/or highly skewed? A tutorial on modelling health  
348 behaviour as count data with Poisson and negative binomial regression. *Health Psychology*  
349 *and Behavioral Medicine* 9, 436-455, <https://doi.org/10.1080/21642850.2021.1920416>
- 350 Guitart, C., Robert, Y.-A., Lotfinejad, N., Fourquier, S., Martin, Y., Pires, D., Sauser, J.,  
351 Beuchat, R., Pittet, D., 2021. Assessing the accuracy of a new hand hygiene monitoring  
352 device (SmartRub®): from the laboratory to clinical practice. *Antimicrobial Resistance &*  
353 *Infection Control* 10, 158, <https://doi.org/10.1186/s13756-021-01026-2>

354 Han, O., Tan, H.W., Julious, S., Sutton, L., Jacques, R., Lee, E., Lewis, J., Walters, S., 2022.  
355 A descriptive study of samples sizes used in agreement studies published in the PubMed  
356 repository. BMC Medical Research Methodology 22, 242, [https://doi.org/10.1186/s12874-](https://doi.org/10.1186/s12874-022-01723-5)  
357 [022-01723-5](https://doi.org/10.1186/s12874-022-01723-5)

358 Joint Monitoring Program, 2022. Progress on household drinking water, sanitation and  
359 hygiene 2000-2022: a special focus on gender. WHO and UNICEF, WHO and UNICEF Joint  
360 Monitoring program.

361 Kazembe, D., 2018. Barriers to Accessing WASH for Disabled People in Chiradzulu, Malawi.

362 Kwong, L.H., Ercumen, A., Pickering, A.J., Unicomb, L., Davis, J., Luby, S.P., 2020. Age-  
363 related changes to environmental exposure: variation in the frequency that young children  
364 place hands and objects in their mouths. J Expo Sci Environ Epidemiol 30, 205-216,  
365 <https://doi.org/10.1038/s41370-019-0115-8>

366 Landis, J.R., Koch, G.G., 1977. The Measurement of Observer Agreement for Categorical  
367 Data. Biometrics 33, 159-174, <https://doi.org/10.2307/2529310>

368 Lazarte, C.E., Encinas, M.E., Alegre, C., Granfeldt, Y., 2012. Validation of digital  
369 photographs, as a tool in 24-h recall, for the improvement of dietary assessment among rural  
370 populations in developing countries. Nutrition Journal 11, 61, [https://doi.org/10.1186/1475-](https://doi.org/10.1186/1475-2891-11-61)  
371 [2891-11-61](https://doi.org/10.1186/1475-2891-11-61)

372 Manun'Ebo, M., Cousens, S., Haggerty, P., Kalengaie, M., Ashworth, A., Kirkwood, B., 1997.  
373 Measuring hygiene practices: a comparison of questionnaires with direct observations in  
374 rural Zaire. Trop Med Int Health 2, 1015-1021, [https://doi.org/10.1046/j.1365-3156.1997.d01-](https://doi.org/10.1046/j.1365-3156.1997.d01-180.x)  
375 [180.x](https://doi.org/10.1046/j.1365-3156.1997.d01-180.x)

376 Masuda, Y.J., Fortmann, L., Gugerty, M.K., Smith-Nilson, M., Cook, J., 2014. Pictorial  
377 approaches for measuring time use in rural Ethiopia. Soc Indic Res 115, 467-482,  
378 <https://doi.org/10.1007/s11205-012-9995-x>

379 McNemar, Q., 1947. Note on the sampling error of the difference between correlated  
380 proportions or percentages. Psychometrika 12, 153-157,  
381 <https://doi.org/10.1007/BF02295996>

382 Montenij, L.J., Buhre, W.F., Jansen, J.R., Kruitwagen, C.L., de Waal, E.E., 2016.  
383 Methodology of method comparison studies evaluating the validity of cardiac output  
384 monitors: a stepwise approach and checklist. Br J Anaesth 116, 750-758,  
385 <https://doi.org/10.1093/bja/aew094>

386 National Cancer Institute, 2023. 24-hour Dietary Recall (24HR) At a Glance.

387 Ram, P.K., Halder, A.K., Granger, S.P., Jones, T., Hall, P., Hitchcock, D., Wright, R.,  
388 Nygren, B., Islam, M.S., Molyneaux, J.W., Luby, S.P., 2010. Is structured observation a valid  
389 technique to measure handwashing behavior? Use of acceleration sensors embedded in  
390 soap to assess reactivity to structured observation. Am J Trop Med Hyg 83, 1070-1076,  
391 <https://doi.org/10.4269/ajtmh.2010.09-0763>

392 Rego, R., Watson, S., Alam, M.A.U., Abdullah, S.A., Yunus, M., Alam, I.T., Chowdhury, A.,  
393 Haider, S.M.A., Faruque, A., Khan, A.I., Hofer, T., Gill, P., Islam, M.S., Lilford, R., 2021. A  
394 comparison of traditional diarrhoea measurement methods with microbiological and  
395 biochemical indicators: A cross-sectional observational study in the Cox's Bazar displaced  
396 persons camp. EClinicalMedicine 42, 101205, <https://doi.org/10.1016/j.eclim.2021.101205>

397 Ross, I., Bick, S., Ayieko, P., Dreibelbis, R., Wolf, J., Freeman, M.C., Allen, E., Brauer, M.,  
398 Cumming, O., 2023. Effectiveness of handwashing with soap for preventing acute  
399 respiratory infections in low-income and middle-income countries: a systematic review and  
400 meta-analysis. Lancet 401, 1681-1690, [https://doi.org/10.1016/S0140-6736\(23\)00021-1](https://doi.org/10.1016/S0140-6736(23)00021-1)

401 Schmidt, W.P., Lewis, H.E., Greenland, K., Curtis, V., 2019. Comparison of structured  
402 observation and pictorial 24 h recall of household activities to measure the prevalence of

403 handwashing with soap in the community. International Journal of Environmental Health  
404 Research 29, 71-81, <https://doi.org/10.1080/09603123.2018.1511772>

405 Stanton, B.F., Clemens, J.D., Aziz, K.M., Rahman, M., 1987. Twenty-four-hour recall,  
406 knowledge-attitude-practice questionnaires, and direct observations of sanitary practices: a  
407 comparative study. Bull World Health Organ 65, 217-222,

408 Tidwell, J.B., Gopalakrishnan, A., Lovelady, S., Sheth, E., Unni, A., Wright, R., Ghosh, S.,  
409 Sidibe, M., 2019. Effect of Two Complementary Mass-Scale Media Interventions on  
410 Handwashing with Soap among Mothers. Journal of Health Communication 24, 203-215,  
411 <https://doi.org/10.1080/10810730.2019.1593554>

412 Tidwell, J.B., Gopalakrishnan, A., Unni, A., Sheth, E., Daryanani, A., Singh, S., Sidibe, M.,  
413 2020. Impact of a teacher-led school handwashing program on children's handwashing with  
414 soap at school and home in Bihar, India. PLOS ONE 15,  
415 <https://doi.org/https://doi.org/10.1371/journal.pone.0229655>

416 Trevethan, R., 2017. Sensitivity, Specificity, and Predictive Values: Foundations, Plabilities,  
417 and Pitfalls in Research and Practice. Front Public Health 5, 307,  
418 <https://doi.org/10.3389/fpubh.2017.00307>

419 Wagner, E.G., Lanoix, J.N., World Health Organization, 1958. Excreta disposal for rural  
420 areas and small communities. World Health Organization.

421 Wolf, J., Hubbard, S., Brauer, M., Ambelu, A., Arnold, B.F., Bain, R., Bauza, V., Brown, J.,  
422 Caruso, B.A., Clasen, T., Colford, J.M., Jr., Freeman, M.C., Gordon, B., Johnston, R.B.,  
423 Mertens, A., Pruss-Ustun, A., Ross, I., Stanaway, J., Zhao, J.T., Cumming, O., Boisson, S.,  
424 2022. Effectiveness of interventions to improve drinking water, sanitation, and handwashing  
425 with soap on risk of diarrhoeal disease in children in low-income and middle-income settings:  
426 a systematic review and meta-analysis. Lancet 400, 48-59, [https://doi.org/10.1016/S0140-6736\(22\)00937-0](https://doi.org/10.1016/S0140-6736(22)00937-0)

428 Wolf, J., Johnston, R., Freeman, M.C., Ram, P.K., Slaymaker, T., Laurenz, E., Pruss-Ustun,  
429 A., 2019. Handwashing with soap after potential faecal contact: global, regional and country  
430 estimates. Int J Epidemiol 48, 1204-1218, <https://doi.org/10.1093/ije/dyy253>

431 World Health Organization, United Nations Children's Fund, 2006. Core questions on  
432 drinking water and sanitation for household surveys. World Health Organization, Geneva.

433 Wright, J.A., Gundry, S.W., Conroy, R., Wood, D., Preez, M.D., Ferro-Luzzi, A., Genthe, B.,  
434 Kirimi, M., Moyo, S., Mutisi, C., 2006. Defining episodes of diarrhoea: results from a three-  
435 country study in Sub-Saharan Africa. Journal of Health, Population and Nutrition, 8-16,

## Artwork

*Figure 1. Maps of Malawi and Chiradzulu district.*

Left: Map of Malawi with Chiradzulu district highlighted in red (Adapted from (D-Maps.com, 2023)). Right: Map of Chiradzulu district and its traditional authorities taken from (Kazembe, 2018).

*Figure 2. Histogram: difference in total number of handwashing opportunities measured by pictorial 24h recall and structured observation with normal density function overlaid.*

*Figure 3. Bland-Altman plot of the number of handwashing opportunities with regression-based bias and limits of agreement.*

Bland-Altman plot of difference in number of opportunities for handwashing measured by P24hR and observations against the mean number of opportunities recorded by the two methods. Bias represented by a solid green line. Limits of agreement (mean difference  $\pm$  2 SD) are shown by the shaded grey section. Bias is estimated by  $y = 2.70 - 0.757 * ((\text{observations} + \text{P24hR})/2)$ . Lower LOA is estimated by  $y = -0.719 - 1.35 * ((\text{observations} + \text{P24hR})/2)$ . Upper LOA is estimated by  $y = 4.67 - 0.164 * ((\text{observations} + \text{P24hR})/2)$ . Overlapping points separated by jitter effect.

## Tables with captions

*Table 1. Characteristics of participants and households*

| <b>Characteristics of participants</b>                     | <b>n (%) or mean (SD)</b> |
|--|---------------------------|
| n  | 88                        |
| Female   | 80 (91%)                  |
| Age  | 35.0 (14.5)               |
| Education  |                           |
| Primary  | 56 (63.6%)                |
| Secondary  | 29 (33.3%)                |
| <b>Characteristics of households</b>                       |                           |
| n  | 74                        |
| Household residents, median (IQR)                          | 5 (4, 6)                  |
| Presence of child < 5 years                                | 67 (90.5%)                |
| Age of child < 5 years, mean (SD)                          | 2.7 (1.2)                 |
| Electricity  | 11 (14.9%)                |
| Mobile phone   | 55 (74.3%)                |
| Monthly income (MWK)                                       |                           |
| <K10,000.00  | 9 (12.2%)                 |
| K10,000.00 to K20,000.00                                   | 19 (25.7%)                |
| K20,000.00 to K30,000.00                                   | 15 (20.3%)                |
| K30,000.00 to K40,000.00                                   | 15 (20.3%)                |
| K40,000.00 to K50,000.00                                   | 5 (6.7%)                  |
| >K50,000.00  | 11 (14.8%)                |
| Water source   |                           |
| Unprotected well   | 1 (1.3%)                  |
| Borehole or tubewell                                       | 71 (96.0%)                |
| Piped into compound, yard, plot                            | 2 (2.7%)                  |
| Sanitation facility  |                           |
| No toilet or neighbour's toilet (not shown)                | 8 (10.8%)                 |
| Flush / pour flush   | 1 (1.3%)                  |
| Pit latrine with slab                                      | 32 (43.2%)                |
| Pit latrine without slab                                   | 33 (44.6%)                |
| Handwashing facility                                       |                           |
| Mobile object reported (bucket / jug / kettle / tippy tap) | 11 (14.9%)                |
| No handwashing place in dwelling/yard/plot                 | 63 (85.1%)                |
| Presence of soap in the household                          | 46 (62.2%)                |

Table 2. Observed and reported hygiene and sanitation opportunities and practices.

| Measurement method      | Opportunities for HWWS | Opportunities per participant, median (IQR) | HWWS practiced (%) | Opportunities for child faeces disposal | Safe child faeces disposal (%) |
|-------------------------|------------------------|---|--------------------|---|--------------------------------|
| Structured observations | 531                    | 5 (3.5, 8.5)                                | 7 (1.3)            | 6                                       | 5 (83.3)                       |
| P24hR                   | 412                    | 4 (3, 6)                                    | 29 (7.0)           | 16                                      | 15 (93.8)                      |

Table 3. Evaluation of validity of pictorial 24h recall compared to structured observation for hygiene and sanitation behaviours (n pairs = 88)

| Behaviour                                   | Observed agreement | Kappa -score | McNemar's test <i>p</i> -value | Reported only (n) | Observed only (n) | Reported and observed (n) | Sensitivity | Specificity | PPV  | NPV |
|---|--------------------|--------------|--------------------------------|-------------------|-------------------|---------------------------|-------------|-------------|------|-----|
| Any handwashing with soap practiced         | 70.5%              | -0.054       | 0.009                          | 20                | 6                 | 1                         | 14%         | 75%         | 4.8% | 91% |
| Any opportunities for child faeces disposal | 78.4%              | -0.002       | 0.069                          | 14                | 5                 | 1                         | 17%         | 83%         | 6.7% | 93% |
| Any safe child faeces disposal practiced    | 80.7%              | 0.024        | 0.049                          | 13                | 4                 | 1                         | 20%         | 84%         | 7.1% | 95% |









