# Accepted Manuscript

An evaluation of the fixed concentration procedure for assessment of acute inhalation toxicity

Fiona Sewell, Ian Ragan, Ian Indans, Tim Marczylo, Nigel Stallard, David Griffiths, Thomas Holmes, Paul Smith, Graham Horgan

Please cite this article as: Sewell, F., Ragan, I., Indans, I., Marczylo, T., Stallard, N., Griffiths, D., Holmes, T., Smith, P., Horgan, G., An evaluation of the fixed concentration procedure for assessment of acute inhalation toxicity, *Regulatory Toxicology and Pharmacology* (2018), doi: 10.1016/j.yrtph.2018.01.001.

1 **An evaluation of the Fixed Concentration Procedure for assessment of acute**

2 **inhalation toxicity**

3 Fiona Sewell[a] Ian Ragan[b], Ian Indans[c], Tim Marczylo[d], Nigel Stallard[e], David Griffiths[f],

4 Thomas Holmes[g], Paul Smith[h], Graham Horgan[i]

5 [a,*]National Centre for the Replacement, Refinement and Reduction of Animals in Research

6 (NC3Rs), UK

7 [b]Board member, NC3Rs

8 [c]Health and Safety Executive, UK

9 [d]Public Heath England, UK

10 [e]University of Warwick, UK

11 [f]Envigo, UK

12 [g]Exponent International Ltd, UK

13 [h]Charles River Laboratories Edinburgh Ltd., UK

14 [i] Biomathematics & Statistics Scotland (BioSS), UK

15 *Corresponding author. Email address: Fiona.Sewell@nc3rs.org.uk

16

17

18 **Abstract**

19 Acute inhalation studies are conducted in animals as part of chemical hazard identification

20 and for classification and labelling. Current methods employ death as an endpoint (OECD

21 TG403 and TG436) while the recently approved fixed concentration procedure (FCP[1])

22 (OECD TG433) uses fewer animals and replaces lethality as an endpoint with evident

23 toxicity. Evident toxicity is the presence of clinical signs that predict that exposure to the

24 next highest concentration will cause severe toxicity or death in most animals. Approval of

25 TG433 was the result of an international initiative, led by the National Centre for the

26 Replacement, Refinement and Reduction of Animals in Research (NC3Rs), which collected

27 data from six laboratories on clinical signs recorded for inhalation studies on 172

28 substances. This paper summarises previously published data and describes the additional

29 analyses of the dataset that were essential for approval of the TG.

30

31 **Highlights:**

32 The FCP for acute inhalation toxicity has been accepted by OECD as TG433.

33 TG433 uses evident toxicity while other approved methods use lethality.

34 A sighting study with 1 M and 1 F animal reliably identifies the more sensitive sex.

35 The three methods ($LC_{50}$, ATC, FCP) showed good agreement in a retrospective analysis.

36

37 **Keywords:**

---

38 Acute inhalation studies; 3Rs; Evident toxicity; Fixed concentration procedure (FCP);

39 Refinement; Regulatory toxicology; TG403; TG436; TG433.

## 1. Introduction

Acute inhalation studies are conducted in animals as part of chemical hazard identification and for classification and labelling purposes. There has been considerable work towards refining the existing methods so that 'evident toxicity' rather than death can be used as an endpoint, through the use of the fixed concentration procedure (FCP) (OECD, 2004). This has recently been accepted as OECD test guideline (TG) 433 as an alternative to the currently accepted $LC_{50}$[2] and the Acute Toxic Class (ATC) methods (OECD TGs 403 and 436 respectively) (OECD, 2009a; OECD, 2009b). The FCP also has the potential to use fewer animals, due to the use of a single sex, and fewer studies overall, as it will obviate the need to test at the next concentration up in some cases. The principles of the three methods are summarised in Table 1 and are described in more detail in Sewell *et al.* (2015). In brief, the $LC_{50}$ method involves testing at three or more concentrations to enable construction of a concentration-mortality curve and a point estimation of the $LC_{50}$ which allows classification into one of five toxic classes using the globally harmonised system (GHS) of classification and labelling of chemicals (OECD, 2001) (Table 2). The ATC method is a refinement of the $LC_{50}$. Rather than a point estimate of the $LC_{50}$, this method estimates which toxic class the $LC_{50}$ falls within, so that classification can be assigned. It uses an 'up-and-down' procedure to test up to four fixed concentrations from the boundaries of the categories (or toxic classes) in the GHS classification system. Depending on the number of deaths at each concentration further testing may be required, or a classification can be made. The FCP uses a similar up-and-down approach to the ATC, but instead identifies an exposure concentration that causes evident toxicity rather than death, so that the $LC_{50}$ can be inferred (based on the prediction of death at the next fixed higher concentration). Classification can then be assigned according to the GHS criteria using the predicted $LC_{50}$. Figures 1, 2 and 3 summarise the possible study outcomes and the resulting classifications for the $LC_{50}$, ATC and FCP methods

---

[2] the concentration that is expected to result in the death of 50% of the animals

65 respectively, using a starting concentration of 5mg/L for dusts and mists as an example

66 (Price et al., 2010).

67 The FCP was removed from the OECD work plan in 2007 because of three main concerns:

68 the ill-defined and subjective nature of evident toxicity; the lack of evidence for comparable

69 performance to the $LC_{50}$ and ATC methods; and suspected sex differences (the FCP

70 originally proposed the default use of females).  Concerns about the definition of 'evident

71 toxicity' were raised despite its long use in the Acute Oral Fixed Dose Procedure (OECD

72 TG420) without guidance on what constitutes evident toxicity, nor in the dermal toxicity

73 equivalent of this TG (OECD TG434) which was approved in 2017 without similar guidance.

74 However, all the concerns about the FCP have been resolved through the work of a global

75 initiative led by the UK National Centre for the Replacement, Refinement and Reduction of

76 Animals in Research (NC3Rs) resulting in its acceptance in April 2017.

77 Some of the work that led to this decision has already been published (Sewell *et al.*, 2015).

78 This previous paper described analyses of a large data set of acute inhalation studies using

79 the $LC_{50}$ or ATC methods in which signs predictive of death at the next highest concentration

80 (i.e. evident toxicity) were identified. Further analyses were needed to address fully the

81 points noted above and to satisfy concerns raised by the OECD national coordinators during

82 the consultation process, and were therefore vital for the final acceptance of the FCP

83 method by OECD. These included further support for the robustness of the signs previously

84 identified, new statistical calculations to support the value of the sighting study in choosing

85 the most sensitive sex, and retrospective classifications to compare outcomes obtained

86 using the three methods. This paper summarises the previously published data and presents

87 the new analyses that formed the basis for acceptance of the new test guideline.

88

89 **2. The robustness of evident toxicity as an endpoint**

90 **2.1 Definitions**

91  Evident toxicity is an accepted endpoint in the fixed dose procedure for acute oral toxicity

92  studies (OECD TG420) (OECD, 2002a).  Here evident toxicity is defined as "*a general term*

93  *describing clear signs of toxicity following the administration of test substance, such that at*

94  *the next highest fixed dose either severe pain and enduring signs of severe distress,*

95  *moribund status or probable mortality in most animals can be expected.*"  However, for this

96  accepted test guideline, no further guidance has been provided on what constitutes 'evident

97  toxicity', and it is not clear how often this test guideline is being used in practice.

98  Although evident toxicity was already accepted as an endpoint for this existing test guideline,

99  criticism of this endpoint was a major factor for the withdrawal of the FCP from the OECD

100  work plan in 2007, due to concerns around subjectivity.  With the aim of making evident

101  toxicity more objective and transferable between laboratories, the NC3Rs working group

102  collected data on the clinical signs observed in individual animals during acute inhalation

103  studies on 172 substances (Sewell *et al.*, 2015).  Because data was collected from a number

104  of laboratories, there was some variation in terminology, requiring retrospective

105  harmonisation by the working group leading to an agreed lexicon of signs (Sewell *et al.*,

106  2015).  These data were analysed to identify signs that could predict lethality would occur if

107  the animals were exposed to the next highest concentration, lethality here being defined as

108  the death, or severe toxicity requiring euthanasia, in two or more animals in a group of five.

109

110  There are three important quantities derived from the analysis. The positive predictive value

111  (PPV) is defined as the percentage of times that the presence of a sign correctly predicts

112  lethality at the next highest concentration. A value less than 100% indicates some false

113  positives that would result in over-classification of the substance, undesirable from a

114  business perspective, but erring on the side of caution for human safety. Sensitivity is

115  defined as the proportion of lethality predicted by the presence of the sign at the lower

116  concentration.  There is no expectation that a single sign would predict 100% of toxicity at

117  the next higher concentration, but signs with very low levels of sensitivity are less useful

118  because of their rarity and their small contribution to overall evident toxicity. Less than 100%

119  sensitivity indicates some false negatives, that is, lethality occurs at the higher concentration

120  even though the sign was absent at the lower concentration. This does not result in incorrect

121  classification as testing would be carried out at the higher concentration anyway.  Specificity

122  is the measure of the percentage of non-lethality at the higher concentration associated with

123  the absence of the sign at the lower concentration. The individual signs focussed upon were

124  those with high PPV and specificity, with appreciable sensitivity.

125

126  In the absence of any deaths at the lower concentration, toxicity occurred at the higher

127  concentration in 77% of the studies (95% CI 72-82%), hence this value was used to set a

128  threshold for use of a sign as an indicator of toxicity.  Consequently, those signs with PPV's

129  not only in excess of this value, but whose lower value of the 95% confidence limits of the

130  PPV also exceeded 77% were selected.

131

132  **2.2 Death as a predictor of toxicity at the next highest concentration**

133  In the Sewell *et al.* (2015) dataset, death or euthanasia was found in the majority of studies

134  at one or more concentrations.  The PPV of a single death at the lower concentration was

135  93% (95% CI 84-98%) i.e. a single death is a strong predictor of lethality at the higher

136  concentration. Although evident toxicity is the intended endpoint for the FCP method, and

137  severe toxicity and death are to be avoided where possible, if death does occur this endpoint

138  can therefore also be used to make decisions concerning classifications (Figure 1). But

139  interestingly, since death is used as an objective endpoint for $LC_{50}$ and ATC methods, it

140  should also be noted that when two deaths occurred at the lower concentration this too was

141  only 97% (95% CI 91-99%) predictive of lethality at the next higher concentration.  That is to

142  say, for a small number of the studies conducted, fewer deaths occurred at the higher

143  concentration than at the lower. For the ATC method in particular, this could lead to an

144  inaccurate classification.

145

**2.3 Signs observed on day 0**

147

148 Signs seen on the day of the test cannot unambiguously be ascribed to the chemical and

149 may have resulted from handling, restraint or the inhalation procedure. Some signs such as

150 wet coat and writhing were only observed on day 0, but some of the common and severe

151 signs were seen both on day 0 and on subsequent days. For two such signs, irregular

152 respiration and hypoactivity, the effect of discounting the day 0 observations increased the

153 PPV and specificity (Sewell *et al.,* 2015) showing that signs that persist for more than 24h

154 after exposure are better predictors of toxicity. However, as pointed out in this paper and in

155 the new TG, severe signs seen on day 0 should be a signal to halt the study or possibly

156 euthanize the animals so affected.

157

**2.4 Signs of evident toxicity**

159

160 In the case of one death at the lower concentration, a number of signs observed in the

161 surviving animals increased the PPV of the single death (Sewell *et al.,* 2015). Some of these

162 also had high sensitivity. Most importantly, a subset of these were also seen to be highly

163 predictive in the absence of death at the lower level. The four signs in this subset were:

164 hypoactivity, tremors, bodyweight loss (>10%), and irregular respiration (Table 3). The data

165 showed that if any of these signs were observed in at least one animal from the day after

166 exposure, animals were highly likely to die if exposed to the next higher concentration.

167 Where any animals experienced tremors or hypoactivity this was 100% predictive of lethality

168 at the next higher concentration. If any animal experienced body weight loss in excess of

169 10% of their pre-dosing weight, this was predictive of death at the higher concentration in

170 94% of cases. Similarly, body weight loss has previously been shown to be a reliable and

171 frequent objective marker for the determination of the maximum tolerated dose (MTD) in

172 short term toxicity tests in animals (Chapman *et al.,* 2013). Irregular respiration was also

173 highly predictive, being indicative of lethality in 89% of cases.

174

175 These four signs were chosen to represent evident toxicity since they had lower 95%

176 confidence interval limits in excess of the 77% threshold detailed above. However, there

177 were other signs that were also highly predictive of lethality at the next higher concentration,

178 albeit with wider confidence intervals often due to their infrequent occurrence in the dataset.

179 For example, oral discharge occurred rarely (sensitivity 2.4%), but was 100% (95%

180 confidence interval (CI) 54.9 -100%) predictive of lethality at the next highest concentration.

181 Therefore the signs used to guide the decision of evident toxicity should not necessarily be

182 restricted to the four signs named in Table 3. Information on the pred

183 ictivity and sensitivity of each of the clinical signs observed in the dataset has been made

184 available in Supplementary Data File 1. Information on subclasses of the dataset for dusts

185 and mists, males and females is also available. This is intended to complement and add to

186 study director judgement and experience so that a decision can be made on the recognition

187 of evident toxicity in the absence of death or the four named signs.

188

189 The definition of 'evident toxicity' used for the purpose of the analysis was conservative

190 when considering the accepted definition of evident toxicity in TG420, since it was based

191 simply on the prediction of actual mortality or euthanasia at the higher concentration (in the

192 absence of death at the lower), and did not also include 'severe distress or moribund status'

193 at the higher concentration. However, this definition was chosen to reflect the different

194 outcomes used for decision making in the protocol, so that 'evident toxicity' could be used to

195 predict 'outcome A' (the death of 2 animals at the higher concentration), and therefore avoid

196 the need for testing at that level (Figure 1). By using evident toxicity, classification can be

197 made based on the *prediction* of death at the higher concentration. The method therefore

198 has the potential to minimise the number of studies (i.e. concentrations tested) that will be

199    required to make a classification and reduce the overall degree of suffering of animals in the

200    study.

201

202    **2.5 Severity and duration of signs**

203    Severity of signs was not recorded consistently in the dataset, only whether a sign was

204    present or not, and as the data had been generated in a number of different laboratories, the

205    grading of severity may have had a strong subjective element. Therefore in the previous

206    publication, only the severity of bodyweight loss was examined in more detail as it had been

207    recorded as either unspecified, mild (reduced weight gain), moderate (10-20% compared

208    with day 0) or substantial (>20% compared with day 0). In fact, PPV was largely unaffected

209    by dividing body weight loss into these subcategories, but sensitivity declined because of the

210    smaller numbers in each category.

211    Another way of looking at severity was to examine whether the sign was present in more

212    than one animal. In the previous paper (Sewell *et al.,* 2015), it was shown that for irregular

213    respiration (the sign for which there are the largest number of observations), the impact on

214    PPV and specificity of increasing numbers of animals showing the sign was very small.

215    However, because seeing the sign in a majority of animals was less common, the sensitivity

216    declined accordingly.

217    **2.6 Combinations and co-occurrence of signs (including signs in isolation)**

218    Sewell *et al.,* (2015) considered whether combinations of signs would increase sensitivity,

219    and thereby improve prediction of lethality at the higher concentration.  However, the gains

220    in sensitivity of all pairwise combinations were small because of the strong co-occurrence of

221    signs, and inclusion of third or fourth signs had progressively less impact.

222    At the other extreme, we examined whether misclassification was likely if a sign was the only

223    one reported (i.e. seen in isolation), and occurred only once and in only one animal. Irregular

224 respiration and body staining were the most commonly observed signs in isolation (42% and

225 28% respectively of those animals that showed the sign) (Table 4). However, of the 268

226 pairs of studies[3] analysed, there were only 5 in which irregular respiration was recorded in

227 the absence of other signs, and only once in only one animal. In each case, at least two

228 animals died at the next higher concentration showing that the single sign was predictive

229 (Table 5). Admittedly this is a small data set, but the finding supports the general robustness

230 of the sign which is typically seen in more than one animal, and rarely occurs in isolation.

231 **2.7 Varying concentration ratios**

232 An odd feature of the GHS classification system is that the ratios of $LC_{50}$ concentrations

233 defined for each grade 1-5 are not of equal size but vary from 2 to 10. For example, for

234 dusts and mists the concentrations tested are 0.05, 0.5, 1.0 and 5.0 mg/l (Table 2). Sewell *et*

235 *al.* (2015) considered how this would affect classifications by the FCP method. It seemed

236 possible that lethality at the higher concentration would be more likely if the concentration

237 ratio was larger and that conversely, a smaller change in concentration might lead to a

238 greater number of false positives i.e. lethality not seen at the higher concentration despite

239 evident toxicity at the lower. This has now been looked at in two ways. Sewell *et al.* (2015)

240 found that, for a small number of signs, the average concentration ratio for false positives

241 was smaller than for true positives, in agreement with this idea. However, of the four signs

242 selected as markers of evident toxicity, two were never associated with false positives (PPVs

243 of 100%) and in the other two cases, the effect of concentration ratio did not reach statistical

244 significance.

245 A further analysis was undertaken to look at the effect of the ratio of the higher to lower

246 concentration on the PPV. In Table 6, PPVs are shown for a number of signs with >2 to <5,

247 >5 to <10 or >10-fold ratios between the lower and higher concentrations. As anticipated,

248 PPVs are higher for the larger concentration ratios, but since the majority of the studies used

---

[3] A pair of studies indicates a set of data from five animals, either all male or all female, exposed at two concentrations differing by at least a factor of two and in which no deaths occurred at the lower concentration.

249 the >2 to <5 fold ratio, the lower numbers in the remaining studies resulted in wider 95%

250 confidence limits of the PPV values. The conclusion is that the main signs of evident toxicity

251 were equally predictive regardless of the ratio of the higher to lower concentration.

252 **3. Default sex and sighting studies**

253 For the $LC_{50}$ procedure, since males and females are treated identically and classifications

254 are based on the sex that is most sensitive, sex differences generally do not have any

255 impact on classification. For the ATC procedure, since males and females are not treated

256 separately and the endpoints are based on the total number of deaths, irrespective of sex,

257 differences in sensitivity have more of an impact and make the test less stringent. For

258 example, where there is a 10-fold difference in sex sensitivity, simulations (Price *et al.,* 2011)

259 showed that substances where the $LC_{50}$ value of the most sensitive sex falls within GHS

260 class 3 (the narrowest GHS classification band), these are almost always incorrectly

261 classified as GHS class 4 (i.e. as less toxic). However, the guideline suggests that testing

262 should be conducted in the more sensitive sex alone if a sex difference is indicated, which

263 may mitigate this if sex differences are correctly identified in practice.

264 The original FCP method proposed the use of females as the default, as these were thought

265 to be the more sensitive sex, and males only used if they were known to be more sensitive.

266 In practice, significant differences in sensitivities between the sexes are fairly uncommon.

267 Price *et al.,* (2011) showed a significant statistical difference between the $LC_{50}$ values of

268 males and females for 16 out of 56 substances examined (29%), females being the more

269 sensitive in 11 of these. The dataset in Sewell *et al.* (2015) revealed little difference in

270 sensitivity between the sexes. There was no difference in the prevalence of death or

271 animals requiring euthanasia between the sexes, though some clinical signs were more

272 prevalent in one sex than the other (ano-genital staining was more prevalent in females than

273 males (p = 0.0002), whereas facial staining and gasping were marginally more common in

274 males (p = 0.028 and 0.044 respectively). However, the predictivity of these signs did not

275 differ between males and females, but the smaller numbers of studies in this analysis led to

276 wider confidence intervals.

277

278 The statistical simulations carried out by Price *et al.* (2011) showed that where there was an

279 unanticipated sex difference and testing was carried out in the less sensitive sex, this would

280 usually result in misclassification, regardless of the method used. Consequently, the new

281 test guideline proposes that a sighting study should be performed not only to determine a

282 suitable starting concentration for the main study but to also identify whether there is a more

283 sensitive sex. The sighting study is not recommended if there is existing information on

284 which to base these two decisions. Despite the earlier proposal that females should be the

285 default sex, the more recent data that failed to show any difference, and the general view of

286 the OECD coordinators, and their nominated inhalation experts, that males were potentially

287 more sensitive for inhaled substances, led to the proposal that males should be used in

288 preference.

289

290 The new sighting study uses a single male and a single female at one or more of the fixed

291 concentrations, depending on the outcome at each concentration as described by Stallard *et*

292 *al.* (2011) (Figure 4). If there is no difference in sensitivity between the sexes, then the

293 choice of sex for single sex studies for the FCP is irrelevant, and will not affect the

294 classification. Since males are now the default sex, if they are the more sensitive, correct

295 classification will still be made, since this is correctly based on the more sensitive sex. It is

296 only if females are the more sensitive sex and this is <u>not</u> correctly identified, that there is

297 potential for incorrect classification.

298

299 Though the risk of a sex difference is low, the new sighting study must be robust enough

300 despite using only one male and one female to identify the large differences in sensitivities

301 that might risk misclassification. To demonstrate this, we have carried out statistical

302  calculations of the probability of choosing the most sensitive sex, with varying ratios of male

303  and female sensitivity (i.e. $LC_{50}$ values) (Figure 5). The methods are similar to those

304  described by Stallard *et al.* (2011). Figure 5 shows the classification probabilities using the

305  new sighting study for dusts and mists with a concentration-response curve slope of 4 and *R*

306  (the ratio of the $LC_{50}$ and $TC_{50}$, the concentration expected to cause death or evident toxicity)

307  of 5 for both sexes, assuming a sighting study starting at 0.05mg/L. The heavy solid line

308  gives the probability of the correct classification given the $LC_{50}$. The heavy dashed line gives

309  the probability that the main study is conducted in females rather than males.

310

311  The first plot of Figure 5 corresponds to the case of no difference between the sexes (i.e.

312  males and females have identical $LC_{50}$ values).  In this case, the probability of the main

313  study being carried out in females varies around 0.25, and since there is no difference in

314  sensitivity this will not affect the classification. The other plots show what happens with

315  increasingly large sex differences, with the females becoming more susceptible.  In these

316  cases the $LC_{50}$ on the *x*-axis is that for the females, as this is the true value on which

317  classification should be based (since females are more sensitive), and the dashed line gives

318  the probability that the main study is conducted in the females. When the sex difference is

319  small, there is quite a high chance of erroneously testing in the males when the females are

320  marginally more sensitive.  For example, for a $LC_{50}$ ratio 1.5 the probability of incorrectly

321  testing in the males is more than 0.5 in many cases.  However, since the sex difference is

322  small this is unlikely to impact the classification.  As the sex difference increases, the chance

323  of seeing the sex difference in the sighting study and doing the main test in the females

324  correctly also increases.  For a ratio of $LC_{50}$ values of 10 or more the probability of choosing

325  females for the main test exceeds 0.9 except for the least toxic substances, when no effects

326  are seen in either sex even at the highest test concentration, or extremely toxic substances,

327  when deaths are seen in both sexes at the lowest test concentration.  The probability of

328  misclassification is higher therefore for GHS classes 3 and 4.

329 These simulations show that the use of a single male and a single female in the sighting

330 study should be sufficient to identify broad differences in sensitivities. Since the effect of sex

331 differences is less when the concentration-response curve is steeper, these simulations

332 represent a worst-case scenario when based on a slope of 4, as it is estimated that only 1%

333 of substances have a concentration-response curve slope of less than this (Greiner, 2008).

334 Again, it is important to note that sex differences are relatively uncommon and only

335 unanticipated greater sensitivity in females is likely to influence classification. Furthermore,

336 for many substances prior knowledge may be also available (e.g. from the oral route) which

337 can be used to verify or indicate any suspected or apparent differences in sensitivity.

338 For the FCP method, the purpose of the sighting study is also to identify the starting

339 concentration for the main study where existing information is insufficient to make an

340 informed decision. A starting concentration should be chosen that is expected to cause

341 evident toxicity in some animals, and the use of two animals, one male and one female,

342 should be sufficient to determine whether this estimation is too high and allow a lower dose

343 to be used in the main study, particularly if existing data are available. The ATC method

344 does not include a sighting study and the choice of starting concentration is based on prior

345 knowledge or experience, or use of the suggested default starting concentrations of 10 mg/L,

346 1 mg/L or 2500 ppm for vapours, dusts/mists and gases, respectively. This is also an option

347 for the FCP method, since the sighting study is not compulsory. However, without the aid of

348 a sighting study, it is possible that an inappropriate starting concentration may be chosen,

349 which could result in testing at more concentrations and using more animals.

350 **4. Comparability to existing methods and retrospective analyses**

351 A number of publications have addressed the comparability of the three methods using

352 statistical calculations or simulations to compare the classifications made by each of the

353 three methods and the likelihood of misclassification (under or over) (Price *et al.* 2011;

354 Stallard *et al.* 2011; Stallard *et al.* 2003). The calculations described above were based on

355    hypothetical mortality concentration curves (with varying steepness) for a range of $LC_{50}$

356    values covering all five toxic classes to represent a wide range of hypothetical substances.

357    These include substances that clearly fall within a specific toxic class, (i.e. $LC_{50}$ within the

358    mid-range of the class bracket) as well as those on the class border (i.e. the most or least

359    toxic substances in each class) where there is greater potential for misclassification. The

360    simulations also took into account the potential for variation between the actual

361    concentration tested and the intended fixed concentration. For the calculations, a variation

362    of +/- 25% was used although this is greater than that permitted in the TG (+/- 20%) so these

363    represent worst-case examples.

364    The statistical calculations showed that the three methods were comparable, although each

365    of the methods did have the potential to misclassify even though the risk of this was low

366    overall (Price *et al.,* 2011). If anything, the FCP tended to over-classify and the other two

367    methods to under-classify. The impact of misclassification (over or under) and the choice of

368    inhalation test method may raise some diversity of opinion depending on safety, commercial

369    and 3Rs (Replacement, Refinement and Reduction) perspectives. The tendency of the $LC_{50}$

370    and ATC methods to *under*-classify is more of a concern to human health than the FCP

371    tendency for *over*-classification. However, it is worth highlighting that the statistical models

372    that these conclusions were based on used a conservative 'worst-case' scenario, with a low

373    concentration-response slope of four, and the potential to over-classify becomes less with a

374    steeper concentration–response curve. Moreover, the models used a greater than permitted

375    variation of the actual concentration from that intended.

376    The statistical calculations described above show that the three methods are comparable,

377    particularly in the absence of sex differences, or where these have been taken into account

378    with the use of the sighting study. However, all these methods rely on the assumption of

379    correct identification or prediction of the $LC_{50}$ value and the corresponding GHS class and

380    are not based on real data. We have therefore undertaken further analysis of the data set of

381    178 dusts and mists to make retrospective classifications by all three methods and to

382  compare their performance. For each method, the classifications were established using the

383  protocols and flow charts in their corresponding test guidelines, based on the order the

384  studies were carried out in practice (i.e. using the default or otherwise determined starting

385  concentration). Supplementary Data File 2 contains information on the 'classification rules'

386  for each method.  For the $LC_{50}$ method, rather than establish an $LC_{50}$ value from the data, a

387  flowchart method was used based on whether more or less than 50% animals died at each

388  concentration (as in Figure 1 in Price *et al*. 2011). Only 'valid' concentrations corresponding

389  to within ±20% of the four fixed concentrations for dusts and mists in the ATC and FCP

390  protocols (0.05, 0.5, 1 and 5 mg/L) were included, to comply with the guidelines.

391  Retrospective classifications could only be made for substances where all the necessary and

392  valid concentrations were available.  For example, in the FCP method, where testing started

393  at 1mg/L and there was no death or evident toxicity in any animal, further testing would be

394  required at 5mg/L.  If this concentration had not been tested or fell outside of the ±20%

395  criterion, then this substance could not be classified.

396  Retrospective classifications were made for 77 substances via the $LC_{50}$ method, 57

397  substances via ATC, and 124 substances for FCP.  For the FCP, classifications were

398  generally able to be made using one or two concentrations requiring five to ten animals

399  (Table 7). For the ATC and $LC_{50}$ methods, classifications were generally made after two

400  concentrations, requiring 12 animals and 20 animals respectively.

401  There were 42 substances for which a retrospective classification was made via all three

402  methods (including based on females and males separately), and for 35 of these (83.3%) all

403  classifications were in agreement (Table 8). If using the $LC_{50}$ as the 'reference' method

404  (though as described above there are limitations for this method and potential for

405  misclassification), the ATC method under-classified by one class on three occasions.  For

406  the FCP method, when conducted in males only, there was one occasion of over-

407  classification, and one of under-classification, both by one class.  When the FCP was

408  conducted in females only, there was also one occasion of over-classification, in the

409 adjacent more stringent class, but three occasions of under-classification, one of these by

410 two classes (class 4 *vs.* class 2). The reasons for these differences could be because the

411 retrospective classification method was not able to take sex differences into account, or

412 because the $LC_{50}$ value falls near a class border where there is greater potential for

413 misclassification. Table 9 shows that for 6 of these 7 substances there appears to be a

414 more sensitive sex. If for the FCP, the classification is made according to the most sensitive

415 sex, there are fewer disagreements with the classifications from the $LC_{50}$ method. For

416 example, instead there are now three occasions where classification made via FCP differs

417 from $LC_{50}$, and these are all over-classifications into the adjacent more stringent class.

418 Whereas the three occasions where the ATC method differed from the $LC_{50}$ method were

419 under-classifications into the less stringent adjacent class. This supports the conclusions

420 from the statistical calculations that show the FCP is comparable to the existing methods if

421 sex differences are taken into account.

422 Often it was not possible to make a retrospective classification using all three methods (e.g.

423 due to a missing concentration), and there are more examples of the classifications made by

424 two of the methods. Table 10 shows the agreement between any two of the methods. With

425 the exception of the male and female comparisons, which had an agreement of 76.5% and

426 87.0% for the FCP and $LC_{50}$ methods respectively, there was over 90% agreement with all

427 combinations of the other methods. Supplementary Tables S1 -S7 compare the

428 classifications made by each of these methods. The difference between the male and

429 female comparisons may reflect differences in sensitivities between sexes and the fact that

430 for the other comparisons the same animals will have been used to make the classification,

431 which could not be done for the male and female comparisons. It is vital for the acceptance

432 of the new TG that there is strong agreement between the classifications made by the FCP

433 and the two accepted methods, irrespective of the sex used by the FCP.

434 However, as previously pointed out, a major difference between the three methods is the

435 number of studies required to make a classification and consequently the numbers of

436 animals used (Table 7).

**5. Summary and conclusions**

438 The new work described here strengthens and clarifies the conclusions of earlier

439 publications on the FCP method. In particular we have shown that evident toxicity can

440 reliably predict death or moribund status at the next highest fixed concentration irrespective

441 of the fold-change in concentration or the number of animals showing the sign of evident

442 toxicity, so demonstrating the robustness of the method.

443 As part of the OECD approval process, the simplicity of the definition of evident toxicity was

444 questioned (i.e. that evident toxicity is said to have been reached if only one of the four signs

445 is observed at least once in at least one animal). However, the dataset had been

446 extensively interrogated to look at multiple scenarios, including the effect of combinations of

447 signs, the duration of signs, and/or the number of animals displaying the sign(s) (see

448 sections 2.5 and 2.6 and Sewell et al., 2015). Whilst predictivity did increase to some extent

449 for some of these, these were associated with wider confidence intervals, since the pool of

450 data also decreased. Clearly, if other data sets become available, it might be possible to

451 confirm these trends more precisely. Therefore, increases in severity and/or the number of

452 animals displaying the sign may increase confidence in the decision, but the statistical

453 analysis of the dataset supports the simple definition regardless of any of such additional

454 information.

455 The change of the default sex from female to male was an unexpected outcome from the

456 consultation with the OECD national coordinators, but there was no evidence from the

457 analysis of Sewell *et al.* (2015) for a consistent bias one way or the other. The decision

458 therefore to adopt males as the default sex was based on the experience of the national

459 coordinators and their nominated inhalation experts. However, since use of the less sensitive

460 sex could result in misclassification, it was important to establish that the proposed sighting

461 study with one male and one female would have the power to identify the more sensitive

462 sex, at least under those circumstances where the difference in sensitivity was large enough

463 that it might have led to wrong classification and in the absence of existing information on

464 sex differences. The results of the statistical analysis confirms that a sighting study with one

465 male and one female has the power to identify the more sensitive sex.

466 The retrospective analysis of the dataset to classify the chemicals by all three methods

467 ($LC_{50}$, ATC and FCP) was especially important in gaining acceptance of TG 433 by OECD.

468 Agreement between the three methods is very good as only 7 out of 42 substances showed

469 any disagreement between the three methods and then by only one class if the most

470 sensitive sex was selected for the FCP method. All three methods have the potential to

471 misclassify so it is important that the advantages and limitations of each test method are

472 understood so that users can select the most appropriate test method for their needs.

473 However in the absence of any other considerations, the FCP method is to be preferred

474 since it offers animal welfare benefits through the avoidance of death as an endpoint, and

475 other 3Rs benefits through the use of fewer animals and fewer studies when compared to

476 the ATC and $LC_{50}$ methods.  We hope that these factors will encourage wide uptake and use

477 of the method in the future.

478 We attribute the reluctance to use the equivalent method for oral toxicity studies (TG 420) to

479 lack of guidance on evident toxicity and the absence of the detailed analyses described

480 here, that were needed to convince the OECD national coordinators that TG 433 was fit for

481 purpose. A similar exercise is therefore planned in collaboration with the European

482 Partnership for Alternatives to Animal Testing to examine clinical signs observed during

483 acute oral toxicity studies and to provide guidance that will encourage the use of TG 420.

484 The experience of gaining acceptance of the FCP method for acute inhalation has been both

485 positive and negative. The positive is the agreement to accept extensive retrospective

486 analysis as sufficient justification for a new test guideline without the need for prospective

487 validation studies which would have required further use of animals. This approach could no

488 doubt be used on other occasions. The negative is the inordinately long time it has taken to

489 get this method accepted even though the principle of evident toxicity had already been

490 accepted by OECD, and the cumbersome process of consultation and submission which

491 was required. Even now, the experience with the oral toxicity guideline TG 420 suggests that

492 there will still be work needed to ensure that TG 433 becomes the preferred method for

493 assessment of inhalation toxicity, and it is to be hoped that this will not take a further 13

494 years.

495

## 496 Acknowledgements

497 We would like to thank everyone who was involved in the Test Guideline Development

498 Process, including the OECD secretariat, the OECD national co-ordinators and their

499 nominated experts.

500

501

## 502 References

503 Chapman, K., F. Sewell, L. Allais, J.L. Delongeas, E. Donald, M. Festag, S. Kervyn, D.

504 Ockert, V. Nogues, H. Palmer, M. Popovic, W. Roosen, A. Schoenmakers, K.

505 Somers, C. Stark, P. Stei, and S. Robinson. 2013. A global pharmaceutical company

506 initiative: an evidence-based approach to define the upper limit of body weight loss in

507 short term toxicity studies. *Regulatory toxicology and pharmacology : RTP*. 67:27-38.

508 Greiner. 2008. *'Report on biostatistical performance assessment of draft TG436 acute toxic

509 class method for acute inhalation toxicity,'* 2008.

510 OECD. 1987. Deleted Test Guidline 401: Acute Oral Toxicity. LD50 method.

511 http://ntp.niehs.nih.gov/iccvam/docs/acutetox_docs/udpproc/udpfin01/append/appi.pd

512 f.

513 OECD. 2001. Harmonized integrated hazard classification system for human health and

514 environmental effects of chemical substances.

515 http://www.oecd.org/officialdocuments/publicdisplaydocumentpdf/?doclanguage=en&

516 cote=env/jm/mono%282001%296.

517 OECD. 2002a. Test Guideline 420: Acute Oral Toxicity - Fixed Dose Procedure.

518 http://www.oecd-ilibrary.org/environment/test-no-420-acute-oral-toxicity-fixed-dose-

519 procedure_9789264070943-en.

520 OECD. 2002b. Test guideline 423: Acute Oral toxicity - Acute Toxic Class Method.

521 http://www.oecd-ilibrary.org/environment/test-no-423-acute-oral-toxicity-acute-toxic-

522 class-method_9789264071001-en.

523 OECD. 2004. Draft proposal for a new guideline:433.

524 http://www.oecd.org/chemicalsafety/testing/32035886.pdf.

525 OECD. 2008. Test Guideline 425: Acute Oral Toxicity: Up-and-Down Procedure.

526 http://www.oecd-ilibrary.org/environment/test-no-425-acute-oral-toxicity-up-and-

527 down-procedure_9789264071049-en.

528 OECD. 2009a. OECD Test Guideline 403: Acute Inhalation Toxicity. http://www.oecd-

529 ilibrary.org/docserver/download/9740301e.pdf?expires=1433247636&id=id&accname

530 =guest&checksum=7012504CE687B2E5614DB637989CE606.

531 OECD. 2009b. Test Guideline 436: Acute Inhalation Toxicity - Acute Toxic Class Method.

532 http://www.oecd-

533 ilibrary.org/docserver/download/9743601e.pdf?expires=1433247696&id=id&accname

534 =guest&checksum=DF8991A8B9D88D13E75E914D1A2D5022.

535 Price, C., N. Stallard, S. Creton, I. Indans, R. Guest, D. Griffiths, and P. Edwards. 2011. A

536 statistical evaluation of the effects of sex differences in assessment of acute

537 inhalation toxicity. *Human & experimental toxicology*. 30:217-238.

538 Sewell, F., I. Ragan, T. Marczylo, B. Anderson, A. Braun, W. Casey, N. Dennison, D.

539      Griffiths, R. Guest, T. Holmes, T. van Huygevoort, I. Indans, T. Kenny, H. Kojima, K.

540      Lee, P. Prieto, P. Smith, J. Smedley, W.S. Stokes, G. Wnorowski, and G. Horgan.

541      2015. A global initiative to refine acute inhalation studies through the use of 'evident

542      toxicity' as an endpoint: Towards adoption of the fixed concentration procedure.

543      *Regulatory toxicology and pharmacology.* 73:770-779.

544 Stallard, N., C. Price, S. Creton, I. Indans, R. Guest, D. Griffiths, and P. Edwards. 2011. A

545      new sighting study for the fixed concentration procedure to allow for sex differences.

546      *Human & experimental toxicology.* 30:239-249.

547 Stallard, N., A. Whitehead, and I. Indans. 2003. Statistical evaluation of the fixed

548      concentration procedure for acute inhalation toxicity assessment. *Human &*

549      *experimental toxicology.* 22:575-585.

550 **Figure Legends**

551 **Figure 1:** $LC_{50}$ test (OECD test guideline 403) for dusts and mists, using example

552 concentrations, starting at 5 mg/L (Price et al., 2010). Please note the $LC_{50}$ test method does

553 not require fixed concentrations, but specifies that 10 animals (5 males and 5 females)

554 should be exposed at three different concentration levels. The concentration levels should

555 be sufficiently spaced to enable construction of a mortality curve so that an estimation of the

556 $LC_{50}$ can be obtained.

557

558 **Figure 2:** Acute toxic class (ATC) method for dusts and mists for an example starting

559 concentration of 5 mg/L (Price et al., 2010). Please note, the ATC method specifies that 6

560 animals (3 males and 3 females) are tested at fixed concentrations that form the upper limit

561 of the GHS categories. The starting concentration is either the highest concentration, or that

562 which is expected to lead to mortality in some of the exposed animals, based on prior

563 information.

564

565 **Figure 3:** Fixed concentration procedure (FCP) method for dusts and mists for an example

566 starting concentration of 5 mg/L (Price et al., 2010). Please note, the draft test guideline

567 specifies that substances are tested at fixed concentrations that form the upper limit of the

568 GHS categories. The starting concentration is chosen to be the fixed concentration level that

569 is most likely to lead to evident toxicity but not death.

570

571 **Figure 4:** FCP sighting study for dusts and mists.

572

573 **Figure 5:** Classification probabilities for the fixed concentration procedure (FCP) with the

574 new sighting study for dusts and mists with concentration-response curve slope of 4 and R

575 (LC50/TC50) of 5 assuming sighting study starting at 0.05 mg/L. The different plots show

576 varying sex differences, to assess the impact of increased female sensitivity compared to

577 male (i.e. female $LC_{50}$ increasingly lower than male $LC_{50}$). The vertical dotted line in each

578 plot indicates the classification boundary concentrations and the light solid line indicates the

579 cumulative probabilities of classification (on left-hand axis scale) into each toxic class for

580 $LC_{50}$ values shown. The heavy solid line gives the probability of the correct classification

581 given the $LC_{50}$. The heavy dashed line gives the probability that the main study is conducted

582 in females rather than males. For more information on these plots please refer to Stallard *et*

583 *al*. (2011).

584

585 **Supplementary data**

586 **Supplementary Data File 1**: Information on the predictivity and sensitivity of each of the

587 clinical signs observed in the dataset.

588 **Supplementary Data File 2**: 'Classification rules' for each method.

**Table 1:** Comparison of LC$_{50}$, ATC and FCP methods.

| Parameter | LC$_{50}$ (concentration causing 50% lethality) | ATC (acute toxic class) | FCP (fixed concentration procedure) |
|---|---|---|---|
| **OECD Test Guideline** | 403 | 436 | 433 |
| **Endpoint** | Death | Death | Evident toxicity |
| **Sighting study** | No sighting study required. | No sighting study required. | A sighting study may be carried out to help inform the starting concentration and choice of sex, if deemed necessary. This is not compulsory.<br><br>1M+1F at one to four concentrations (usually only one or two concentrations required).<br><br>The starting concentration should be that which is most expected to produce evident toxicity in some animals. If no prior information is available this should be 10 mg/L, 1 mg/L or 2500 ppm for vapours, dusts/mists and gases, respectively. |
| **Number of animals** | 5M+5F per study.<br><br>Usually three studies required.<br><br>Min 10 – max 40 animals. | 3M+3F per study.<br><br>Usually at least two studies required (12 animals), though classification can sometimes be made based on one study, if testing at the lowest or highest concentrations (depending on the outcome).<br><br>Numbers of animals range from 6 to max 24 (depending on the number of studies). An inappropriate starting concentration (causing too much or too little toxicity) may require testing at additional concentrations and may therefore result in higher numbers of animals being used.<br><br>Where a marked sex difference is observed additional | Single (most sensitive) sex, or males only as default. 5 animals per study.<br><br>Classification can often be made after a single study (5 animals).<br><br>Numbers of animals range from 5 to max 20 (depending on the number of studies). Plus 2-8 in the sighting study (though the use of 8 animals in the sighting study would be very unusual, and only if the highest or lowest concentrations were chosen inappropriately as the starting concentration).<br><br>An inappropriate starting concentration (causing too much or too little toxicity) may require testing at additional concentrations and may therefore result in |

| | | animals may be required. | higher numbers of animals being used. However, a sighting study should avoid this. |
|---|---|---|---|
| **Number of concentrations** | At least three concentrations (to enable production of a concentration-mortality curve and estimation of $LC_{50}$). | An 'up and down method' is used, requiring 1 to 4 fixed concentrations (based on the upper limit of the GHS classification system) depending on the outcome at each concentration.<br><br>Generally at least two concentrations are required to make a classification. Sometimes a classification can be made based on only one study if starting at the highest or lowest fixed concentration, and depending on the outcome. | An 'up and down method' is used, requiring 1 to 4 fixed concentrations (based on the upper limit of the GHS classification system) depending on the outcome at each concentration.<br><br>A classification can often be made based on one study only. |
| **Starting concentration** | n/a<br><br>This is not a sequential method. At least three concentrations are required to enable production of a concentration-mortality curve and estimation of $LC_{50.}$ | Starting concentration level should be that which is most likely to produce toxicity in some animals.<br><br>If no prior information is available the starting concentration will be 10 mg/L, 1 mg/L or 2500 ppm for vapours, dusts/mists and gases, respectively.<br><br>An inappropriate starting concentration (causing too much or too little toxicity) may require testing at more concentrations than if a more appropriate concentration had been chosen. | Starting concentration level should be that which is most expected to produce evident toxicity in some animals. The sighting study may inform this choice, or prior information if available.<br><br>If a sighting study has not been conducted or is inconclusive, or if no prior information is available the starting concentration will be 10 mg/L, 1 mg/L or 2500 ppm for vapours, dusts/mists and gases, respectively.<br><br>An inappropriate starting concentration (causing too much or too little toxicity) may require testing at more concentrations than if a more appropriate concentration had been chosen. The use of a sighting study should avoid this. |
| **Classification Method** | Based on a point estimate of $LC_{50}$ which allows classification according to the GHS classification system. | Based on an interval estimate of $LC_{50}$, so that classification is based on the toxic class that the estimated $LC_{50}$ falls within, using the GHS classification system. | $LC_{50}$ is inferred through the use of evident toxicity to *predict* death at a higher dose, and classification made according to the inferred $LC_{50}$ using the GHS classification system. |

**Table 2:** GHS classification system for inhalation. For the $LC_{50}$ method, a point estimate of the $LC_{50}$ allows classification into the relevant GHS class according to the table. The ATC method estimates which class the $LC_{50}$ falls within and makes classification on that basis, whereas classifications made by FCP are based on the *inferred* $LC_{50}$.

| GHS category | Vapours (mg/L) | Dusts and mist (mg/L) | Gases (ppm) |
|---|---|---|---|
| 1 (most toxic) | ≤0.5 | ≤0.05 | ≤100 |
| 2 | >0.5 and ≤2 | >0.05 and ≤0.5 | >100 and ≤500 |
| 3 | >2 and ≤10 | >0.5 and ≤1 | >500 and ≤2,500 |
| 4 | >10 and ≤20 | >1 and ≤5 | >2,500 and ≤20,000 |
| 5 | >20 | >5 | >20,000 |

GHS, Globally Harmonised System; LC50, median concentration; ppm, parts per million.

**Table 3:** Clinical signs indicating evident toxicity (PPV, sensitivity and specificity)

| Clinical signs | PPV (95% CI) | | Sensitivity (95% CI) | | Specificity (95% CI) | |
|---|---|---|---|---|---|---|
| Hypoactivity | 100.0 | (92.4 - 100.0) | 18.4 | (13.6 - 24.1) | 100.0 | (95.2 - 100.0) |
| Tremors | 100.0 | (68.8 - 100.0) | 3.90 | (1.90 - 7.20) | 100.0 | (95.2 - 100.0) |
| Bodyweight loss | 94.0 | (84.6 - 98.4) | 22.7 | (17.4 - 28.8) | 95.1 | (87.2 - 98.7) |
| Irregular respiration | 89.0 | (80.9 - 94.5) | 35.3 | (29.0 - 42.0) | 85.2 | (74.7 - 92.5) |

CI, Confidence Interval; PPV, positive predictive value.

**Table 4:** Number of animals displaying a clinical sign in isolation, and the total number of animals displaying the sign.

| Clinical sign | No. animals displaying sign ONLY (%) | | Total no. animals displaying the sign |
|---|---|---|---|
| Irregular respiration | 137 | (42%) | 325 |
| Body staining | 27 | (27%) | 99 |
| Hypoactivity | 12 | (16%) | 77 |
| Laboured respiration | 12 | (16%) | 77 |
| Faeces reduced | 13 | (12%) | 107 |
| Hunched posture | 18 | (8%) | 227 |
| Ano-genital staining | 4 | (8%) | 51 |
| Naso-ocular discharge | 6 | (7%) | 89 |
| Congested respiration | 4 | (5%) | 87 |
| Facial staining | 3 | (5%) | 65 |
| >10% bodyweight loss | 2 | (2%) | 93 |
| Noisy respiration | 1 | (0.4%) | 267 |

**Table 5:** Studies where irregular respiration was observed only once in one animal at the lower concentration in females, with no other signs.

| Study | Concentration tested | Female observations | | Male observations | |
|---|---|---|---|---|---|
| | | Number of Deaths | Number with evident toxicity | Number of Deaths | Number with evident toxicity |
| 1 | 0.05 mg/L | 0 | 1 | 0 | 4 |
| | 0.5 mg/L | 5 | - | 3 | 2 |
| | 2 mg/L | 5 | - | 5 | 0 |
| 2 | 0.06 mg/L | 0 | 1 | 0 | 5 |
| | 0.5 mg/L | 2 | 3 | 3 | 2 |
| | 2 mg/L | 4 | 1 | 5 | - |
| 3 | 0.5 mg/L | 0 | 1 | 0 | 4 |
| | 2 mg/L | 2 | 3 | 2 | 3 |
| 4 | 0.05 mg/L | 0 | 1 | 0 | 2 |
| | 0.2 mg/L | 5 | - | 5 | - |
| | 2 mg/L | 5 | - | 5 | - |
| | 5 mg/L | 5 | - | 5 | - |
| 5 | 0.06 mg/L | 0 | 1 | n/a | n/a |
| | 0.5 mg/L | 2 | 3 | 0 | 5 |
| | 2 mg/L | 5 | - | 5 | 0 |

**Table 6:** PPV (95% confidence interval) for highly predictive signs with 2, 5 or 10-fold concentration change between the lower and higher concentration.

| Clinical sign | ≥2-fold (95% CI) | | ≥5-fold (95% CI) | | ≥10-fold (95% CI) | |
|---|---|---|---|---|---|---|
| Tremors | 100.0 | (68.8 - 100.0) | 100.0 | (5.0 - 100.0) | 100.0 | (5.0 - 100.0) |
| Hypoactivity | 100.0 | (92.0 - 100.0) | 100.0 | (47.3 - 100.0) | 100.0 | (47.3 - 100.0) |
| >10% bodyweight loss | 91.7 | (79.0 - 97.8) | 85.7 | (47.0 - 99.3) | 100.0 | (36.8 - 100.0) |
| Irregular respiration | 89.0 | (80.9 - 94.5) | 95.8 | (81.2 - 99.8) | 100.0 | (86.1 - 100.0) |
| Body staining | 88.5 | (71.8 - 97.0) | 100.0 | (60.7 - 100.0) | 100.0 | (22.4 - 100.0) |
| Ano-genital staining | 86.4 | (67.3 - 96.4) | 0.0 | (0.0 - 95.0) | 100.0 | (5.0 - 100.0) |
| Faeces reduced | 85.3 | (70.4 - 94.4) | 100.0 | (47.3 - 100.0) | 100.0 | (47.3 - 100.0) |
| Naso-ocular discharge | 84.2 | (70.1 - 93.3) | 100.0 | (74.1 - 100.0) | 100.0 | (65.2 - 100.0) |
| Noisy respiration | 80.5 | (70.9 - 88.0) | 94.1 | (74.3 - 99.7) | 100.0 | (68.8 - 100.0) |
| Hunched posture | 78.0 | (65.0 – 87.8) | 87.5 | (64.5 - 97.8) | 100.0 | (54.9 - 100.0) |
| Gasping | 76.5 | (52.5 – 92.0) | 100.0 | (22.4 - 100.0) | 100.0 | (22.4 - 100.0) |

**Table 7:** Number of studies required to make a classification, and the associated number of animals.

| No. studies to make a classification | FCP | | | ATC | | LC$_{50}$ | |
|---|---|---|---|---|---|---|---|
| | No. animals involved | No. studies | | No. animals involved | No. studies | No. animals involved | No. Studies |
| | | FCP-F | FCP-M | | | | |
| 1 study | 5 | 54 | 64 | 6 | 18 | 10 | 32 |
| 2 studies | 10 | 46 | 41 | 12 | 37 | 20 | 41 |
| 3 studies | 15 | 1 | 3 | 18 | 2 | 30 | 3 |
| 4 studies | 20 | 0 | 1 | 24 | 0 | 40 | 1 |

**Table 8:** Classifications made by all three methods, showing the number of substances classified into each class and the number of substances where there was a disagreement between the three methods (which is expanded on in Table 9).

| Classification | No. substances |
|---|---|
| Class 1 | 1 |
| Class 2 | 11 |
| Class 3 | 3 |
| Class 4 | 14 |
| Class 5 | 6 |
| Disagreements | 7 |

**Table 9:** Substances where there were differences in retrospective classifications made via the $LC_{50}$, ATC and FCP methods. FCP retrospective classifications were made for both females (F) and males (M) only. For each substance the concentrations tested, the number of deaths and/or animals with evident toxicity are indicated.

| Substance | Concentrations tested | No. deaths | | No. evident toxicity | | Classification | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | F | M | F | M | $LC_{50}$ | ATC | FCP(F) | FCP(M) | |
| 1 | 0.5mg/L | 0 | 0 | 0 | 0 | 3 | 4 | 3 | 4 | |
| | 1 mg/L | 4 | 1 | 1 | 0 | | | | | |
| 2 | 1 mg/L | 0 | 0 | 4 | 4 | 5 | 5 | 4 | 5 | |
| | 5 mg/L | 2 | 1 | 3 | 4 | | | | | |
| 3 | 1 mg/L – males | - | 0 | - | 0 | 5 | 5 | 5 | 4 | |
| | 5 mg/L | 0 | 2 | 5 | 3 | | | | | |
| 4 | 1 mg/L – males | - | 0 | - | 5 | 4 | 5 | 5 | 4 | |
| | 5 mg/L | 0 | 3 | 5 | 2 | | | | | |
| 5 | 0.05 mg/L | 0 | 0 | 0 | 0 | 2 | 2 | 4 | 2 | |
| | 0.5 mg/L | 3 | 4 | 0 | 0 | | | | | |
| | 1 mg/L | 1 | 4 | 2 | 0 | | | | | |
| 6 | 1 mg/L | 0 | 0 | 0 | 0 | 5 | 5 | 4 | 5 | |
| | 5 mg/L | 2 | 0 | 3 | 5 | | | | | |
| 7 | 0.5 mg/L | 0 | 0 | 0 | 0 | 3 | 4 | 4 | 3 | |
| | 1 mg/L | 1 | 3 | 4 | 2 | | | | | |
| | 5 mg/L | 5 | 5 | - | - | | | | | |

**Table 10:** Differences in classifications between the three methods, showing the numbers of substances for which pairwise comparisons were made, and the number for which there was agreement between the two methods.

| Comparison | | | No. classified | No. substances in agreement | % agreement |
|---|---|---|---|---|---|
| FCP-M | *vs.* | FCP-F | 85 | 65 | 76.5% |
| $LC_{50}$-M | *vs.* | $LC_{50}$-F | 46 | 40 | 87.0% |
| ATC | *vs.* | FCP-F | 46 | 42 | 91.3% |
| $LC_{50}$ | *vs.* | FCP-F | 43 | 40 | 93.0% |
| $LC_{50}$ | *vs.* | FCP-M | 44 | 41 | 93.2% |
| ATC | *vs.* | FCP-M | 51 | 48 | 94.1% |
| $LC_{50}$ | *vs.* | ATC | 46 | 44 | 95.7% |

Start

| 0.05 mg/L | 0.5 mg/L | 1 mg/L | 5 mg/L |
| 5 males | 5 males | 5 males | 5 males |
| 5 females | 5 females | 5 females | 5 females |

Ⓐ Ⓑ Ⓒ    Ⓐ Ⓑ Ⓒ *    Ⓐ Ⓑ Ⓒ *    Ⓐ Ⓑ Ⓒ *

GHS category    1  1  2        3        4        5

| 0.05 mg/L | 0.05 mg/L | 1 mg/L |
| 5 animals male or female | 5 animals male or female | 5 animals male or female |

Ⓓ Ⓔ    Ⓓ Ⓔ    Ⓓ Ⓔ

GHS category    1  2        3        4

Outcome

Ⓐ >50% deaths in both sexes

Ⓑ >50% deaths in one sex

Ⓒ <50% deaths in both sexes

Ⓓ >50% deaths

Ⓔ <50% deaths

*testing continues for the most sensitive sex

Starting at 5mg/L

Start

| 0.05 mg/L | 0.5 mg/L | 1 mg/L | 5 mg/L |
|---|---|---|---|
| 3 males | 3 males | 3 males | 3 males |
| 3 females | 3 females | 3 females | 3 females |

(A) (B)    (A) (B)    (A) (B)    (A) (B)

4-6 (at 1)  3 (at 1)

| GHS class | 1 | 2 | 3 | 4 | 4 | 5 |

(A) ≥50% deaths (3-6 animals)   (B) <50% deaths (0-2 animals)

Starting at
5 mg/L

Start

| 5 Animals 0.05 mg/L | 5 Animals 0.5 mg/L | 5 Animals 1 mg/L | 5 Animals 5 mg/L |

(A) (B) (C)　(A) (B) (C)　(A) (B) (C)　(A) (B) (C)

GHS
class

1　2　2　　　3　3　　　4　4　　　5　5

Outcome

(A) 2 or more deaths

(B) 1 or more with evident toxicity and/or 1 death

(C) Neither death nor evident toxicity

Sighting study starting at 0.05 mg/L

Start

| Male and female 0.05 mg/L | Male and female 0.5 mg/L | Male and female 1 mg/L | Male and female 5 mg/L |

A B C   A B C   A B C   A B C

Classify GHS class 1*

Main study starting concentration (mg/L): 0.05    0.05 0.5    0.5 1    1 5 5

Sighting study starting at 0.5 mg/L

Start

| Male and female 0.05 mg/L | Male and female 0.5 mg/L | Male and female 1 mg/L | Male and female 5 mg/L |

A B C   A2 A1 B C   A B C   A B C

Classify GHS class 1*

Main study starting concentration (mg/L): 0.05 0.05    0.5    0.5 1    1 5 5

Most sensitive sex 0.05 mg/L

A B C

Classify GHS class 1*

Main study starting concentration (mg/L): 0.05 0.05

Sighting study starting at 1 mg/L

Start

| Male and female 0.05 mg/L | Male and female 0.5 mg/L | Male and female 1 mg/L | Male and female 5 mg/L |

A B C   A2 A1 B C   A2 A1 B C   A B C

Classify GHS class 1*

Main study starting concentration (mg/L): 0.05 0.05    0.5 0.5    1    1 5 5

Most sensitive sex 0.05 mg/L     Most sensitive sex 0.5 mg/L

A B C   A B C

Classify GHS class 1*

Main study starting concentration (mg/L): 0.05 0.05    0.5    0.5

Sighting study starting at 5 mg/L

Start

| Male and female 0.05 mg/L | Male and female 0.5 mg/L | Male and female 1 mg/L | Male and female 5 mg/L |

A B C   A2 A1 B C   A2 A1 B C   A2 A1 B C

Classify GHS class 1*

Main study starting concentration (mg/L): 0.05 0.05    0.5 0.5    1    1 5 5

Most sensitive sex 0.05 mg/L     Most sensitive sex 0.5 mg/L     Most sensitive sex 1 mg/L

A B C   A B C   A B C

Classify GHS class 1*

Main study starting concentration (mg/L): 0.05 0.05    0.5    0.5    1    1

Outcome from 2 animals

A Any death
B No death, some evident toxicity
C Neither death nor evident toxicity
A1 1 death
A2 2 deaths

Outcome from 1 animal

A Death
B Evident toxicity
C Neither death nor evident toxicity

**LC50 (m)/LC50 (f) = 1**

Probability

0.8

0.4

0.0

0.05    0.50    5.00    50.00

$LC_{50}$ (females)

**LC50 (m)/LC50 (f) = 1.5**

Probability

0.8

0.4

0.05    0.50    5.00    50.00

$LC_{50}$ (females)

**LC50 (m)/LC50 (f) = 2**

Probability

0.8

0.4

0.05    0.50    5.00    50.00

$LC_{50}$ (females)

**LC50 (m)/LC50 (f) = 2.5**

Probability

0.8

0.4

0.0

0.05    0.50    5.00    50.00

$LC_{50}$ (females)

**LC50 (m)/LC50 (f) = 3**

Probability

0.8

0.4

0.0

0.05    0.50    5.00    50.00

$LC_{50}$ (females)

**LC50 (m)/LC50 (f) = 4**

Probability

0.8

0.4

0.0

0.05    0.50    5.00    50.00

$LC_{50}$ (females)

**LC50 (m)/LC50 (f) = 5**

Probability

0.8

0.4

0.0

0.05    0.50    5.00    50.00

$LC_{50}$ (females)

**LC50 (m)/LC50 (f) = 10**

Probability

0.8

0.4

0.0

0.05    0.50    5.00    50.00

$LC_{50}$ (females)

**LC50 (m)/LC50 (f) = 20**

Probability

0.8

0.4

0.0

0.05    0.50    5.00    50.00

$LC_{50}$ (females)

**Highlights:**

The FCP for acute inhalation toxicity has been accepted by OECD as TG433.

TG433 uses evident toxicity while other approved methods use lethality.

A sighting study with 1 M and 1 F animal reliably identifies the more sensitive sex.

The three methods ($LC_{50}$, ATC, FCP) showed good agreement in a retrospective analysis.