1       ***NICEATM Draft Technical Summary***

2

3       **Reduction of the Number of Animals Required for Ocular Irritancy Testing:**

4       **A Brief Review of the Literature**

5

6       **Introduction**

7       In the interest of reducing the number of animals required for regulatory safety testing,

8       regulatory authorities have revised testing procedures to reduce the minimum number of

9       animals required for ocular and dermal corrosivity and irritancy testing by 50-83%.  This

10      reduction was accomplished by changing the requirement for the routine use of six

11      animals to sequential testing of 1-3 animals.  The revised procedures now allow for

12      testing to stop and for a substance to be classified as an ocular/dermal corrosive or severe

13      ocular irritant when the corresponding injury occurs in any one animal during the

14      sequential testing.

15

16      With the objective of further reducing the number of animals used for ocular

17      corrosion/irritancy testing, the National Toxicology Program Interagency Center for the

18      Evaluation of Alternative Toxicological Methods (NICEATM) and the Interagency

19      Coordinating Committee on the Validation of Alternative Methods (ICCVAM) are

20      currently reviewing the validation status of *in vitro* methods proposed for identifying

21      ocular corrosives and severe ocular irritants.  In the course of this review, questions arose

22      regarding the reproducibility of the *in vivo* ocular test.  The decision by regulatory

23      authorities to reduce the minimum number of animals from six to three (or less) was

24      based on the fact that the animal test was considered sufficiently reproducible such that

25      using fewer animals would not significantly alter the accuracy of the test method for

26      hazard classification and labeling purposes.  This draft technical summary briefly reviews

27      the relevant available scientific literature on this topic.

28

29      **History and Background**

30      The original protocol of Draize et al. (1944) served as the basis for initial regulatory

31      requirements for eye irritation testing that mandated the use of at least six rabbits.  In

32   1981, the U.S. Environmental Protection Agency (EPA) published a report entitled " Eye

33   Irritation Testing: An Assessment of Methods and Guidelines for Testing Materials for

34   Eye Irritancy" (EPA, 1981).   This report concluded that initial testing for eye irritation

35   with three animals normally would be sufficient to identify substances that are non-

36   irritating or maximally irritating. However the report noted that testing with additional

37   animals might be necessary to reliably characterize substances of intermediate degrees of

38   irritancy.  In 1981, the Organization for Economic Co-operation and Development

39   (OECD) published Test Guideline (TG) 405, which proposed the use of as few as three

40   animals, but also included the provision that additional animals might be necessary in

41   order to clarify equivocal responses.

42

43   Several analyses were subsequently published that assessed the consequences of reducing

44   the number of rabbits per test from six to as few as two animals (Guillot et al. 1981;

45   DeSousa et al. 1984; Solti and Freeman 1988; Talsma et al. 1988; Springer et al. 1993;

46   Dalbey et al. 1993; Berdasco et al. 1996).  With the exception of Dalbey et al. 1993, each

47   study concluded that reducing the number of rabbits from six to three would not have an

48   unacceptable reduction in the accuracy of ocular irritancy classification/categorization.

49   These analyses were performed using maximum average Draize scores (MAS), internal

50   irritancy classification schemes, and/or regulatory classification schemes as endpoints for

51   comparison.  Several of these studies (DeSousa et al. 1984; Talsma et al. 1988; Dalbey et

52   al. 1993) confirmed that correlations between three-animal and six-animal classifications

53   were the highest among substances classified on the extreme ends of the irritancy range

54   (i.e., non-irritants and severe irritants), and that the majority of variability was seen

55   among substances classified in the middle range of irritation.  However, Dalbey et al.

56   1993 was the only study that concluded that this effect justified the continued routine use

57   of six animals.  EPA (EPA 1998), the European Union (EU 2004), and the OECD (in

58   revised TG 405) now recommend the use of a maximum of three animals, although

59   additional animals may be tested under certain circumstances (e.g., to confirm weak or

60   moderate responses).  The different evaluations on the numbers of rabbits appropriate per

61   study are summarized chronologically in the following sections.

62

63    Due to the lack of individual rabbit data in each of these reports, it was not possible to

64    calculate the impact of reduced animal numbers on hazard classification according to

65    EPA (1996), EU (2001), or the United Nations (UN) Globally Harmonized System

66    (GHS) of Classification and Labeling of Chemicals (UN 2003) classification systems.

67    This shortcoming emphasizes the importance of reporting individual animal scores in

68    publications.

69

70    **Guillot et al. (1981)** *as discussed in EPA (1981)*

71    "Guillot et al (1981) compared the mean ocular irritation ratings in two groups of three

72    rabbits each with those obtained using six rabbits. Classification differences due to test

73    group size resulted for 25 of the 56 substances tested. In only two cases, however, was a

74    test substance classified as non-irritating based on results in one group of three rabbits

75    while testing with three additional rabbits and with six rabbits resulted in a rating of

76    "slightly irritating".  Thus, the data showed that the use of three animals in a preliminary

77    test was adequate in differentiating a positive from a negative response for roughly 96%

78    (54/56) of  a wide variety of substances."

79

80    **DeSousa et al. (1984)**

81    These authors examined the statistical consequences of reducing the number of rabbits

82    per test to five, four, three, or two animals. Data included in the analysis was obtained

83    from three separate laboratories; one laboratory tested 55 chemicals, another tested 11

84    chemicals, and another tested one chemical, for a total of 67 substances.  No substances

85    were tested more than once.  The substances spanned a wide range of chemical classes

86    and the full range of irritancy potential (based on Draize scores).  In their analysis, the

87    authors used the maximum average Draize score (obtained at 1 hour or at 1, 2, 3, 7, 14,

88    and 21 days post treatment) from in-house six-rabbit tests. From the 67 six-rabbit test

89    results, scores for all possible subsample combinations of two to five rabbits were

90    calculated.  The original maximum average six-rabbit score was subtracted from each

91    corresponding subsample score, and the difference plotted versus the original six-rabbit

92    score to provide a measure of variability of the subsample scores.  From these plots,

93    subsample prediction intervals were calculated for the six-rabbit scores.

94   The authors found that prediction intervals for two, three, and four animals were

95   comparable, although decreasing the sample size did increase the prediction interval.  The

96   authors also noted that the greatest discordance from the six-animal test occurred in the

97   middle of the Draize scale (i.e., extremes at the high or low end of Draize scores

98   produced very little difference in prediction intervals).  An additional analysis was

99   performed based on an in-house classification system (Texaco Single-Digit Toxicity

100  Classification – SDTC) that classified the ocular irritancy of test substances according to

101  their maximum average Draize score into five categories (**Table 1**).

102

103  **Table 1.  Texaco Single-Digit Toxicity Classification for Eye Irritation**

| Single-Digit Toxicity Classification (SDTC) | Explanation | Draize rabbit scores |
|---|---|---|
| 0 | Minimally irritating | 0-15 |
| 1 | Slightly irritating | >15-25 |
| 2 | Moderately irritating | >25-50 |
| 3 | Severely irritating | >50-80 |
| 4 | Extremely irritating | >80-110 |

104  From DeSousa et al. (1984)

105

106  Overall, for the three-rabbit subsamples, there was 93% (1241/1340) agreement with the

107  classification obtained when the maximum average score for the six-rabbit test was used.

108  When the analysis was limited to the 10 substances classified as severely irritating (those

109  with maximum average Draize scores from >50 to 80) and extremely irritating (those

110  with maximum average Draize scores from >80 to 110), there was 84% (168/200)

111  agreement.  However, the authors recognized that the overall analysis could be biased by

112  the limited number of substances inducing mid-ranged Draize scores, where the greatest

113  variability would be expected to occur.  Based on the results of their analysis, the authors

114  concluded that a three-animal test system would be acceptable for most ocular irritant

115  classification systems used in the petrochemical industry, with additional three-animal

116  testing being performed in the case of tests with high variation in results.

117

118  **Solti and Freeman (1988)**

119  This poster presentation describes a subsample analysis of petrochemicals classified

120  based on six-animal tests. Based on a scatterplot of eye irritation scores, it appears that

121     the full range of irritancy responses were represented.  From each of the original test

122     groups (n = 6 animals) data for three animals was randomly selected for comparison to

123     the original classification.  A correlation coefficient of 0.975 was calculated from the

124     resulting analysis of 30 studies when comparing the irritation scores among three animals

125     versus six animals.  Both the U.S. Occupation Safety and Heath Administration (OSHA)

126     and the European Economic Community (EEC, now referred to as the European Union -

127     EU) labeling schemes were used to evaluate the studies.  Because the OSHA regulations

128     did not define criteria for a three-animal test, the six-animal rules were extrapolated down

129     to three animals (i.e., positive response in $\geq 2/6$ animals = 1/3 animals).  Accordingly,

130     using a random subpopulation of three animals instead of six animals for labeling would

131     have resulted in different OSHA labeling in 10% (3/30) of the studies (one overclassified

132     [3% false positive rate], two underclassified [7% false negative rate]), while labeling

133     would have been different in 7% (2/30) of the studies (both overclassified [7% false

134     positive rate; 0% false negative rate]).   The authors concluded that using a reduced

135     number of animals in safety evaluations  of petrochemicals would not greatly impact on

136     hazard labeling decisions.

137

138     **Talsma et al. (1988)**

139     These authors performed a subsample analysis of 155 chemical and petroleum products

140     that covered the full range of irritancy responses, in which they also evaluated the ability

141     of irritation scores derived from two-, three-, four-, or five-animal tests to predict a six-

142     animal Draize score.  Similar to the approach used by DeSousa et al. (1984), the authors

143     used in their analysis the maximum average Draize score (obtained at 1, 2, 3, 7, 14, and

144     21 days post treatment) from in-house six-rabbit tests.  Also, similar to DeSousa et al.

145     (1984), the authors applied scores from each subsample to an in house (Amoco) four-

146     category classification system based on specific ranges of Draize scores (**Table 2**).  For

147     comparison, they also classified each substance according to the five-category SDTC

148     system of DeSousa et al. (1984) (see **Table 1**).

149

150    **Table 2.  Amoco Eye Irritation Classification System**

| Rating | Explanation | Draize rabbit scores |
|--------|-------------|----------------------|
| 0 | Minimally irritating | 0-15 |
| 1 | Mildly irritating | 16-30 |
| 2 | Irritating | 31-50 |
| 3 | Extremely irritating/eye damage | >50-80 |

151    Modified from Talsma et al. (1988)

152

153    Talsma et al. found that, for a three-rabbit test, the correlation of randomly selected

154    subset scores (i.e., maximum average score) with the six-rabbit Draize score was 0.99

155    and, overall, the number of correct classifications achieved was 94% (2912/3100)

156    accurate in predicting the six-rabbit classification using the Amoco classification system,

157    and 91% (2813/3100) using the SDTC classification system.  When the analysis was

158    limited to the 23 substances classified in the SDTC system as severely irritating (those

159    with maximum average Draize scores from >50 to 80) and extremely irritating (those

160    with maximum average Draize scores from >80 to 110), there was 89% (408/460)

161    agreement.

162

163    Similar to DeSousa et al. (1984), they also found that the width of the prediction interval

164    for a subsample was inversely proportional to the number of animals evaluated.  Talsma

165    et al. (1988) also pointed out that their database was weighted heavily toward minimally

166    irritating substances, which may have affected the outcome of the analysis, given that

167    they too noted that the greatest disagreement occurred among substances classified in the

168    middle range of irritation.   The authors concluded that their results indicate that a high

169    level of accuracy can be obtained with reduced numbers of rabbits per test.

170

171    **Springer et al. (1993)**

172    This report this is part of the published proceedings of the 1991 IRAG Workshop on

173    Updated Eye Irritation Methods.  This report detailed an analysis of eye irritation tests on

174    pesticides (n = 48, data submitted to the EPA), cosmetics/consumer products (n = 53;

175    data submitted to the U.S. Food and Drug Administration [FDA] or the U.S. Consumer

176    Products Safety Commission [CPSC]), cleaning products/ingredients (n = 30; data

177    submitted to the FDA or the CPSC), and unspecified chemicals (n = 12; data from

178    Marzulli and Ruggles 1973).  The substances spanned the full range of irritancy potential

179    (based on FHSA and/or EPA classification systems).  The only dataset from which

180    substances were tested multiple times was the Marzulli and Ruggles (1973) database.

181    These substances were all borderline compounds with respect to ocular irritancy (i.e.,

182    between nonirritating and irritating) which were tested in twelve different laboratories to

183    yield a total of 139 tests.

184

185    Unlike DeSousa et al. (1984) and Talsma et al. (1988), this study did not perform a

186    subsample analysis.  Data from eye irritation tests from each group of substances were

187    examined to estimate the distribution of positive animal responses for substances

188    classified as irritant or nonirritant.  An animal was classified as positive for eye irritancy

189    if any score attained or exceeded the criterion for a positive response for corneal opacity

190    ≥1, conjunctival redness ≥2 or conjunctival chemosis ≥2.  There was no attempt to limit

191    the evaluation to substances classified as ocular corrosives or severe irritants.  Using data

192    from the six-rabbit tests, probability calculations were performed based on a three-rabbit

193    test (either a one-stage, three-rabbit test, or two-stage approach which sequentially tests

194    up to three rabbits) to determine the likelihood of correctly identifying a substance as

195    irritant or nonirritant.  This analysis showed that a high level of accuracy (≥ 94%; actual

196    numbers not available) could be obtained from a sample size of three rabbits in which

197    two positive responses were required to assign an irritant classification (false positive

198    rates ≤ 5% and false negative rates of 1%).

199

200    However, applying the EPA classification system (in effect at the time of this evaluation)

201    to a three-animal test, where only one animal is required to assign an irritant

202    classification, resulted in much higher false positive rates (20% to 50%).  Based on this

203    evaluation, the authors recommended revising the *in vivo* eye irritation protocol to

204    include testing of only three animals using either a one-stage or two-stage approach.

205

206    **Dalbey et al. (1993)**

207    Similar to DeSousa et al. (1984) and Talsma et al. (1988), Dalbey et al. (1993) evaluated

208    mean weighted Draize eye scores from subsets of two, three, four, or five rabbits and

209    their predictivity of the score produced by a six-rabbit test.  The database consisted of

210    data from 185 six-rabbit eye irritation studies conducted in-house with petroleum-based

211    products.  No indication of substances being tested more than once is provided.  The

212    substances spanned the full range of irritancy potential (based on Draize scores).  The

213    authors used in their analysis average Draize score for cornea, iris, and conjunctiva

214    calculated separately over three days after dosing, or mean weighted Draize scores

215    calculated for cornea, iris, and conjunctiva combined.  This study sought to confirm the

216    earlier conclusions of DeSousa et al. (1984) and Talsma et al. (1988), that a three-rabbit

217    test was suitable for classification of eye irritation and thus the classification resulting

218    from each subset was compared to that resulting from a six-animal test.  In this

219    comparison, the European Commission (EC, now referred to as the EU] and U.S. Federal

220    Hazardous Substances Act (FHSA) classification systems were considered.  In addition, a

221    "workshop classification" (based on recommendations from the 1991 Interagency

222    Regulatory Alternatives Group [IRAG] workshop) was considered.  Similar to previous

223    studies, the agreement between subsets and the original six-animal Draize score was

224    directly proportional to the number of animals.  Dalbey et al. found that, overall for a

225    three-rabbit test, there was approximately 90% agreement with the Draize scores

226    produced by a six-rabbit test.

227

228    With regard to EC classification system, there was 96% (3158/3280) agreement for

229    nonirritants, and 98% (412/420) agreement for irritants.  However, it is noteworthy that

230    only 11% (21/185) of the substances considered in this evaluation were classified as

231    irritants.  Upon classifying the same data according to the "workshop" classification

232    system, only 41% (76/185) were labeled as nonirritants, along with 42% (77/185) severe

233    irritants, and 17% (32/185) irritants.  For the irritant category, there was only 29%

234    (183/640) agreement between the three-animal and six-animal classification.  For

235    substances classified as severe irritants, there was 75% (1152/1540) agreement, while

236    there was 100% (1520/1520) agreement for the nonirritants.  Using the FHSA

237    classification scheme, there was 88% (1340/1520) agreement for nonirritants, 97%

238    (1488/1540) agreement for severe irritants, but only 55% (351/640) agreement for

239    irritants.

240

241    Based on these results, and unlike DeSousa et al. (1984) and Talsma et al. (1988), Dalbey

242    et al. (1993) concluded that the six-rabbit test should continue to be used (at least for the

243    purposes of classifying substances according to the FHSA system), although a three-

244    animal test could be used to screen for nonirritants or the most severe irritants, as these

245    types of substances produced the greatest agreement.  They emphasized the finding that

246    the greatest variability was noted among the middle range of irritation, and only the

247    extremes of the scoring scale were most accurate.  This observation is consistent with

248    previous evaluations that much of the variability lies within the mid-range

249

250    **Berdasco et al. (1996)**

251    These authors also performed a subset analysis of ocular irritation tests for 118

252    substances, by generating scores for five-, four, three, and two-rabbit subsets.  The

253    substances included in this analysis included pesticides, antimicrobials, consumer

254    products and industrial chemicals. The substances spanned the full range of irritancy

255    potential.  Each substance was assigned an ocular irritancy category based on the EPA

256    (1989) classification system using the six-animal test results, and then according to each

257    subset result.  The accuracy of the *in vivo* ocular irritation test using three rabbits instead

258    of six was 96% (113/118) for Category I substances (EPA 1989), with a false negative

259    rate of 10% (5/48) and a false positive rate of 0% (0/48).  Based on these results, and

260    similar to Dalbey et al. (1993), the authors concluded that as few as three animals could

261    be used in an initial eye irritation test, with the provision that up to six rabbits might be

262    necessary to clarify equivocal (or disparate) results.

263

264    **References**

265

266    Berdasco N, Gilbert K, Lacher J, and Mattsson J.  1996.  Low rate of severe injury from

267    dermal and ocular irritation tests and the validity of using fewer animals.  Journal of the

268    American College of Toxicology.  15:177-193.

269

270   Dalbey W, Rodriguez S, Wilkins K, Cope C.  1993. Reducing the number of rabbits in

271   eye and skin irritancy tests.  Journal of the American College of Toxicology 12:347-357.

272

273   DeSousa D, Rouse A, Smolon W.  1984.  Statistical consequences of reducing the

274   number of rabbits utilized in eye irritation testing: Data on 67 petrochemicals.

275   Toxicology and Applied Pharmacology 76:234-242.

276

277   EPA.  1981.  Eye Irritation Testing: An Assessment of Methods and Guidelines for

278   Testing Materials for Eye Irritancy. EPA 560/11-82-001. Office of Pesticides and Toxic

279   Substances, EPA, Washington, D.C.

280

281   EPA.  1989.  Pesticide Assessment Guidelines, Subdivision F, Hazard Evaluation:

282   Human and Domestic Animals, Series 83-4, Primary Eye Irritation. U.S. Government

283   Printing Office, Washington, D.C.

284

285   EPA. 1998.  Health Effects Test Guideline, OPPTS 870.2400 Acute Eye Irritation. EPA

286   712-C-98-195. Washington, DC: U.S. Environmental Protection Agency.

287

288   EU.  2004.  Commission Directive 2004/73/EC of 29 April 2004 adapting to technical

289   progress for the 29th time Council Directive 67/548/EEC on the approximation of laws,

290   regulations and administrative provisions relating to the classification, packaging and

291   labelling of dangerous substances. Official Journal of the European Union L152, 1-316.

292

293   Guillot JP, Caillard L, Gonnet JF, Clement C.  1981.  Chemicals-Ocular and Cutaneous

294   Local Tolerance-"Cosmetic," A.F.N.O.R. and OECD Protocols. Institut Francais de

295   Rechserches et Essais Biologiques.  Centre de Lyon, Les Oncins. b.p. 109, 69210

296   L'Arbresle.

297

298   Marzulli F, Ruggles D.  1973.  Rabbit eye irritation test: collaborative study.  Journal of

299   the Association of Analytical Chemists 56:905-914.

300   OECD. 2002.  Guideline for testing of chemicals revised guideline 405: Acute Eye

301   Irritation/Corrosion. Available: http://www.oecd.org. [accessed 26 August 2004].

302

303   Solti J, Freeman JJ.  1988.  Effect of reducing the number of animals in acute

304   toxicity/irritation tests on U.S. and European labeling requirements.  The Toxicologist

305   8:263.

306

307   Springer J, Chambers W, Green S, Gupta K, Hill R, Hurley P, Lambert L, Lee C, Lee J,

308   Liu P, Lowther D, Roberts C, Seabaugh V, Wilcox N.  1993.  Number of animals for

309   sequential testing.  Food and Chemical Toxicology 31:105-109.

310

311   Talsma D, Leach C, Hatoum N, Gibbons R, Roger J-C, Garvin P.  1988.  Reducing the

312   number of rabbits in the Draize eye irritancy test: A statistical analysis of 155 studies

313   conducted over 6 years.  Fundamental and Applied Toxicology 10:146-153.

314

315   UN. 2003. Globally Harmonised System of Classification and Labeling of Chemicals

316   (GHS). New York & Geneva: United Nations Publications.

317