1
2
3              **POWER CALCULATIONS for LLNA Protocols**

4   **1.0      LLNA:BrdU-ELISA**

5   During their review of the LLNA: BrdU-ELISA test method, some members of the
6   ICCVAM LLNA expert peer review panel requested information on statistical power vs.
7   number of animals used for this assay. They wanted know how many animals would be
8   adequate for detecting the associated threshold stimulation index for a positive response
9   (e.g., SI > 3).

10  This required power calculations to determine the number of animals needed to
11  demonstrate statistical significance control and treatment groups. According, Dr.
12  Haseman was provided vehicle control data (spectrophotometer absorbance values) from
13  11 different experiments with the same vehicle in order to establish the variability among
14  these animals. Within each experiment, there were four animals and three replicates per
15  animal. For each animal, the three replicates were averaged, and then the four individual
16  animal means were averaged ("a mean of the means") to obtain overall control means and
17  standard deviations for that experiment. The data were also log transformed and the
18  transformed data were averaged. The summary statistics are given in **Table 1-1**.

19  **Table 1-1      Summary of the Control Absorbance Data for the LLNA: BrdU-**
20  **ELISA**

| Experiment | Original Scale | | Log Scale | |
|---|---|---|---|---|
| | Mean | SD | Mean | SD |
| 1 | 0.0676 | 0.0051 | -2.70 | 0.077 |
| 2 | 0.1197 | 0.024 | -2.14 | 0.221 |
| 3 | 0.1068 | 0.0425 | -2.29 | 0.367 |
| 4 | 0.0982 | 0.0216 | -2.34 | 0.212 |
| 5 | 0.0696 | 0.0275 | -2.73 | 0.410 |
| 6 | 0.0766 | 0.0329 | -2.64 | 0.457 |
| 7 | 0.0687 | 0.0062 | -2.68 | 0.092 |
| 8 | 0.4833 | 0.0681 | -0.74 | 0.151 |
| 9 | 0.4516 | 0.110 | -0.82 | 0.249 |
| 10 | 0.2479 | 0.1425 | -1.52 | 0.590 |
| 11 | 0.2252 | 0.1044 | -1.58 | 0.491 |

21

22  Several comments on the data:

23          • Note that there is considerable study-to-study variability. For example,
24            note that if Experiments 8 and 9 were actually a "treatment", then it would
25            be declared active relative to most if not all of the first 7 control groups
26            (treated/control ratio >3).

27      • There is much less within-study variability. Note also that on the original
28         scale, the SD tends to increase with increasing means. This suggests that a
29         log transformation will help to stabilize the variability, which in fact was
30         the case.

31      • Another important advantage of taking logs is that the apparent variable of
32         interest is the ratio of the treated to control response. Testing the null
33         hypothesis that this ratio is one is equivalent to testing the null hypothesis
34         that the difference in the logs is zero, which was the test chosen for the
35         power calculations.

36  The first step in the power calculation was to use the data from the 11 experiments to
37  derive a representative mean and SD for the control response. Although alternative
38  approaches are certainly possible, only the mean mean and mean SD were calculated for
39  simplicity (on the log scale). These were mean=-2.02 and SD=0.302. The corresponding
40  control mean on the original scale is 0.133.

41  Three hypothetical changes to the decision criteria when then evaluated:  a tripling of the
42  control response (on the original scale), a doubling of the control response, and a 1.3-fold
43  increase in the control response. Although more elegant tests may be possible, I chose to
44  base my power calculations on a simple one-sided Student's t test applied to the log-
45  transformed data. The calculations that are given below assume the same design that was
46  used in the 11 experiments (i.e., three replicates per animal).  I focused on an N of 4, but
47  also looked at other sample sizes as well. The results are summarized in **Table 1-2**
48  assuming a control response of -2.02 (log scale) and an SD of 0.302.

49  **Table 1-2       Treatment Group (Rx) Response Relative to Controls**

| Parameter | 3-fold Increase | 2-fold Increase | 1.3-fold Increase |
|---|---|---|---|
| Mean Rx response | 0.399 | 0.266 | 0.173 |
| Log (mean Rx response) | -0.92 | -1.32 | -1.76 |
| Difference from control (log scale) | 1.10 | 0.70 | 0.26 |
| Difference/SD | 3.64 | 2.32 | 0.88 |
| Power for N=4 | 99% | 80-90% | <50% |
| Other power | 95% (N=3) | 95% (N=5) | 50% (N=8) |
| Other power | | 50-80% (N=3) | 80% (N=16) |
| Other power | | | 90% (N=22) |

50

51  Therefore, four animals per group with three replicates per animal is sufficient to detect a
52  three-fold increase in the control response and would likely (with reasonable power)
53  detect a two-fold increase (an additional animal would give 95% power; N=3 would be
54  more problematic). However, it would not be realistic to expect to detect a 1.3 fold
55  increase in the control response without a significant addition of animals. Slight changes
56  in the underlying assumptions would not change the results of these power calculations in
57  any meaningful way.

58

58  **2.0    LLNA: BrdU-FC**

59  This set of power calculations is based on vehicle control data (flow cytometry BrdU
60  absorbance values) 64 experiments with four to five animals per experiment. Separate
61  power calculations were carried out for five different vehicles. There were four additional
62  experiments with other vehicles, (acetone, PEG 400 and 1% L92/dH20) but since these
63  vehicles involved only one or two studies, there was insufficient data to carry out a
64  meaningful power calculation. The data are summarized in **Tables 2-1** to **2-6**.

65  **Table 2-1      Summary of the Control Data for the LLNA: BrdU-FC (DAE Vehicle)**

| Experiment | Original scale | | Log scale | |
|---|---|---|---|---|
| | Mean | SD | Mean | SD |
| 1 | 11564.6 | 7776.85 | 8.9124 | 1.3722 |
| 2 | 7420.2 | 2387.47 | 8.8702 | 0.3228 |
| 3 | 4949.4 | 2273.08 | 8.4040 | 0.5330 |
| 4 | 8169.4 | 3838.27 | 8.8964 | 0.5612 |
| 5 | 18143.0 | 5594.13 | 9.7644 | 0.3316 |
| 6 | 7860.6 | 6780.59 | 8.6538 | 0.9457 |
| 7 | 11551.2 | 4883.84 | 9.2772 | 0.4474 |
| 8 | 7524.6 | 5591.07 | 8.7500 | 0.6241 |
| 9 | 17610.8 | 14954.73 | 9.5542 | 0.6937 |
| 10 | 22822.4 | 11361.37 | 9.9076 | 0.6001 |
| 11 | 3759.25 | 2862.25 | 7.9983 | 0.8003 |
| 12 | 14580.2 | 5268.96 | 9.5270 | 0.4045 |

66

67  **Table 2-2      Summary of the Control Data for the LLNA: BrdU-FC (AOO**
68  **Vehicle)**

| Experiment | Original scale | | Log scale | |
|---|---|---|---|---|
| | Mean | SD | Mean | SD |
| 1 | 2328.4 | 1566.27 | 7.5122 | 0.8425 |
| 2 | 17079.0 | 9402.10 | 9.6138 | 0.5903 |
| 3 | 11277.6 | 6872.04 | 9.1858 | 0.5980 |
| 4 | 17932.8 | 14014.27 | 9.3336 | 1.2341 |
| 5 | 8187.6 | 4714.16 | 8.8978 | 0.5121 |
| 6 | 34472.5 | 10504.11 | 10.4082 | 0.3370 |
| 7 | 14813.0 | 5897.59 | 9.5208 | 0.4876 |
| 8 | 14020.8 | 9854.00 | 9.2056 | 1.0883 |
| 9 | 19897.2 | 11461.51 | 9.7562 | 0.6043 |
| 10 | 17975.8 | 3813.69 | 9.7756 | 0.2400 |
| 11 | 6631.8 | 5725.49 | 8.4558 | 0.9473 |
| 12 | 15472.2 | 8093.26 | 9.5202 | 0.5829 |
| 13 | 8749.4 | 5702.84 | 8.8432 | 0.8431 |
| 14 | 11794.6 | 2858.56 | 9.3484 | 0.2688 |
| 15 | 20898.6 | 10979.71 | 9.7754 | 0.7342 |
| 16 | 10648.0 | 1927.73 | 9.2612 | 0.1749 |
| 17 | 16180.0 | 7711.57 | 9.5848 | 0.5393 |
| 18 | 6204.6 | 3877.74 | 8.5434 | 0.7277 |
| 19 | 9628.8 | 5075.28 | 9.0446 | 0.5858 |
| 20 | 7637.6 | 4022.84 | 8.8060 | 0.6072 |

69

70

70 **Table 2-3**     **Summary of the Control Data for the LLNA: BrdU-FC (DMSO**
71                      **Vehicle)**

| Experiment | Original scale | | Log scale | |
|---|---|---|---|---|
| | Mean | SD | Mean | SD |
| 1 | 11892.8 | 4239.52 | 9.3338 | 0.3499 |
| 2 | 17427.0 | 7999.14 | 9.6654 | 0.5283 |
| 3 | 8148.75 | 3707.66 | 8.9220 | 0.4842 |
| 4 | 8031.4 | 1939.59 | 8.9676 | 0.2428 |
| 5 | 40758.25 | 12831.56 | 10.5765 | 0.3238 |
| 6 | 28371.8 | 14171.47 | 10.1586 | 0.4781 |
| 7 | 46420.8 | 18065.75 | 10.6844 | 0.3918 |
| 8 | 24726.0 | 5326.84 | 10.0974 | 0.2146 |
| 9 | 14027.4 | 3476.44 | 9.5208 | 0.2729 |
| 10 | 15314.5 | 9320.34 | 9.5210 | 0.5276 |
| 11 | 13386.0 | 5516.88 | 9.4284 | 0.4399 |
| 12 | 24955.6 | 9786.46 | 10.0250 | 0.5643 |
| 13 | 19335.2 | 7644.20 | 9.8158 | 0.3544 |
| 14 | 41366.4 | 14242.19 | 10.5892 | 0.3088 |
| 15 | 26519.8 | 10408.41 | 10.1218 | 0.4048 |
| 16 | 52644.0 | 17384.31 | 10.8276 | 0.3306 |
| 17 | 21824.8 | 9779.87 | 9.9156 | 0.4243 |
| 18 | 21865.4 | 9182.90 | 9.8892 | 0.5617 |
| 19 | 29371.2 | 6978.60 | 10.2632 | 0.2539 |
| 20 | 22575.4 | 9225.93 | 9.9564 | 0.4170 |
| 21 | 11929.2 | 6187.36 | 9.2744 | 0.5411 |
| 22 | 22382.6 | 8667.60 | 9.9672 | 0.3325 |
| 23 | 22221.0 | 15029.10 | 10.1200 | 0.4161 |
| 24 | 17486.2 | 4157.51 | 9.7444 | 0.2531 |

72

73 **Table 2-4**     **Summary of the Control Data for the LLNA: BrdU-FC (DMF**
74                      **Vehicle)**

| Experiment | Original scale | | Log scale | |
|---|---|---|---|---|
| | Mean | SD | Mean | SD |
| 1 | 5728.8 | 3829.90 | 8.3252 | 1.1704 |
| 2 | 16018.4 | 4502.49 | 9.6438 | 0.1034 |
| 3 | 11607.4 | 9643.83 | 9.0312 | 0.8762 |
| 4 | 35928.2 | 25375.35 | 10.2938 | 0.4949 |

75

76 **Table 2-5**     **Summary of the Control Data for the LLNA: BrdU-FC (ETOH**
77                      **Vehicle)**

| Experiment | Original scale | | Log scale | |
|---|---|---|---|---|
| | Mean | SD | Mean | SD |
| 1 | 4096.2 | 2343.60 | 8.2070 | 0.5064 |
| 2 | 6636.5 | 4310.69 | 8.6040 | 0.7779 |
| 3 | 18806.4 | 5220.25 | 9.8122 | 0.2697 |
| 4 | 6920.4 | 3307.72 | 8.6970 | 0.6828 |

78
79

79 **Table 2-6    Average Means and Standard Deviations for Each Vehicle**

| Vehicle | N | Original Scale Averages | | Log-transformed Scale Averages | | Converted Control Mean[1] | Maximum Difference[2] |
|---|---|---|---|---|---|---|---|
| | | Mean | SD | Mean | SD | | |
| DAE 433 | 12 | 11329.6 | 6131.05 | 9.043 | .6364 | 8459 | 29-fold |
| AOO | 20 | 13591.5 | 6703.74 | 9.220 | .6273 | 10093 | 15-fold |
| DMSO | 24 | 23457.6 | 8969.54 | 9.891 | .3924 | 19753 | 4-fold |
| DMF | 4 | 17320.7 | 10837.89 | 9.324 | .6612 | 11198 | 14-fold |
| EtOH | 4 | 9114.9 | 3795.57 | 8.830 | .5592 | 6837 | 7-fold |

80 [1] Anti-log of the log transformed scale average (used in the power calculations).
81 [2] Maximum difference among animals within an experiment using this vehicle.
82
83 Note the large SD for every group except for the DMSO control. The power calculations
84 given in **Tables 2-7** to **2-12** are based on a one-sided $p < 0.05$ Student's t test applied to
85 the log-transformed data (just as in the previous power calculations; for completeness, the
86 power calculations are included for the acetone vehicle as well, although only two
87 experiments used this vehicle, as noted above). ***It should be noted that these calculations***
88 ***make the additional assumption that any "treatment effect" produced will have***
89 ***essentially the same SD (on a log-transformed scale) as the control data, i.e., that the***
90 ***treatment will change only the mean response and not the variability.***

91 **Table 2-7    Treatment Group (Rx) Response Increase Relative to Controls for**
92 **DAE 433**

| | 3.0-fold increase | 2.5-fold increase | 2.0-fold increase | 1.5-fold increase | 1.3-fold increase |
|---|---|---|---|---|---|
| Mean Rx response | 25377 | 21147.5 | 16918 | 12688.5 | 10996.7 |
| Log (Mean Rx response) | 10.142 | 9.959 | 9.736 | 9.448 | 9.305 |
| Difference (log scale) | 1.099 | 0.916 | 0.693 | 0.405 | 0.262 |
| Difference/SD | 1.73 | 1.44 | 1.09 | 0.64 | 0.41 |
| Power for N=5 | nearly 80% | 50-80% | <50% | <50% | <50% |
| Power for N=4 | 50-80% | 50% | <50% | <50% | <50% |
| Power for N=3 | 50% | <50% | <50% | <50% | <50% |
| Other Power | 95% (N=9) | 95% (N=12) | 95%(N=19) | 95% (N=54) | 95% (N>100) |
| Other Power | 90% (N=7) | 90% (N=10) | 90% (N=15) | 90% (N=43) | 90% (N>100) |

93

94 **Table 2-8    Treatment Group (Rx) Response Increase Relative to Controls for**
95 **AOO**

| | 3.0-fold increase | 2.5-fold increase | 2.0-fold increase | 1.5-fold increase | 1.3-fold increase |
|---|---|---|---|---|---|
| Mean Rx response | 30279 | 25232.5 | 20186 | 15139.5 | 13120.9 |
| Log (Mean Rx response) | 10.318 | 10.136 | 9.913 | 9.625 | 9.482 |
| Difference (log scale) | 1.098 | 0.916 | 0.693 | 0.405 | 0.262 |
| Difference/SD | 1.75 | 1.46 | 1.10 | 0.65 | 0.42 |
| Power for N=5 | 80% | 50-80% | <50% | <50% | <50% |
| Power for N=4 | 50-80% | 50% | <50% | <50% | <50% |
| Power for N=3 | 50% | <50% | <50% | <50% | <50% |
| Other Power | 95% (N=9) | 95% (N=12) | 95%(N=19) | 95% (N=52) | 95% (N>100) |
| Other Power | 90% (N=7) | 90% (N=10) | 90% (N=15) | 90% (N=42) | 90% (N>100) |

96

97
98 **Table 2-9      Treatment Group (Rx) Response Increase Relative to Controls for DMF**

|                        | 3.0-fold increase | 2.5-fold increase | 2.0-fold increase | 1.5-fold increase | 1.3-fold increase |
|------------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| Mean Rx response       | 33594             | 27995             | 22396             | 16797             | 14557.4           |
| Log (Mean Rx response) | 10.422            | 10.240            | 10.017            | 9.729             | 9.586             |
| Difference (log scale) | 1.098             | 0.916             | 0.693             | 0.405             | 0.262             |
| Difference/SD          | 1.66              | 1.39              | 1.05              | 0.61              | 0.40              |
| Power for N=5          | 50-80%            | 50-80%            | <50%              | <50%              | <50%              |
| Power for N=4          | 50-80%            | 50%               | <50%              | <50%              | <50%              |
| Power for N=3          | 50%               | <50%              | <50%              | <50%              | <50%              |
| Other Power            | 95% (N=10)        | 95% (N=12)        | 95%(N=21)         | 95% (N=63)        | 95% (N>100)       |
| Other Power            | 90% (N=8)         | 90% (N=10)        | 90% (N=17)        | 90% (N=48)        | 90% (N>100)       |

99
100 **Table 2-10      Treatment Group (Rx) Response Increase Relative to Controls for ETOH**

|                        | 3.0-fold increase | 2.5-fold increase | 2.0-fold increase | 1.5-fold increase | 1.3-fold increase |
|------------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| Mean Rx response       | 20511             | 17092.5           | 13674             | 10255.5           | 8888.1            |
| Log (Mean Rx response) | 9.929             | 9.746             | 9.523             | 9.236             | 9.092             |
| Difference (log scale) | 1.099             | 0.916             | 0.693             | 0.406             | 0.262             |
| Difference/SD          | 1.97              | 1.64              | 1.24              | 0.73              | 0.47              |
| Power for N=5          | 80-90%            | 50-80%            | 50%               | <50%              | <50%              |
| Power for N=4          | 80%               | 50-80%            | <50%              | <50%              | <50%              |
| Power for N=3          | 50-80%            | 50%               | <50%              | <50%              | <50%              |
| Other Power            | 95% (N=7)         | 95% (N=10)        | 95%(N=15)         | 95% (N=42)        | 95% (N=100)       |
| Other Power            | 90% (N=6)         | 90% (N=8)         | 90% (N=12)        | 90% (N=33)        | 90% (N=80)        |

102
103 **Table 2-11      Treatment Group (Rx) Response Increase Relative to Controls for DMSO**

|                        | 3.0-fold increase | 2.5-fold increase | 2.0-fold increase | 1.5-fold increase | 1.3-fold increase |
|------------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| Mean Rx response       | 59259             | 49382.5           | 39506             | 29629.5           | 25678.9           |
| Log (Mean Rx response) | 10.990            | 10.807            | 10.584            | 10.297            | 10.153            |
| Difference (log scale) | 1.099             | 0.916             | 0.693             | 0.406             | 0.262             |
| Difference/SD          | 2.80              | 2.33              | 1.77              | 1.03              | 0.67              |
| Power for N=5          | 95-99%            | 95%               | 80%               | <50%              | <50%              |
| Power for N=4          | 90-95%            | 80-90%            | 50-80%            | <50%              | <50%              |
| Power for N=3          | 80-90%            | 50%               | 50-80%            | <50%              | <50%              |
| Other Power            |                   |                   | 95%(N=8)          | 95% (N=22)        | 95% (N=49)        |
| Other Power            |                   |                   | 90% (N=7)         | 90% (N=17)        | 90% (N=39)        |

105
106

106 **Table 2-12**     **Treatment Group (Rx) Response Increase Relative to Controls for**
107               **ACE**

|  | 3.0-fold increase | 2.5-fold increase | 2.0-fold increase | 1.5-fold increase | 1.3-fold increase |
|---|---|---|---|---|---|
| Mean Rx response | 25881 | 21567.5 | 17254 | 12940.5 | 11215.1 |
| Log (Mean Rx response) | 10.161 | 9.979 | 9.756 | 9.468 | 9.325 |
| Difference (log scale) | 1.098 | 0.916 | 0.693 | 0.405 | 0.262 |
| Difference/SD | 1.70 | 1.42 | 1.07 | 0.63 | 0.41 |
| Power for N=5 | 50-80% | 50-80% | <50% | <50% | <50% |
| Power for N=4 | 50-80% | 50% | <50% | <50% | <50% |
| Power for N=3 | 50% | <50% | <50% | <50% | <50% |
| Other Power | 95% (N=9) | 95% (N=12) | 95%(N=20) | 95% (N=56) | 95% (N>100) |
| Other Power | 90% (N=7) | 90% (N=10) | 90% (N=16) | 90% (N=45) | 90% (N>100) |

108
109 It is important to understand that the primary factor that influences power (in addition to
110 sample size) is the variability in response among control animals in a given study: the
111 greater the variability, the lower the power. Using this assay, for four of the five vehicles,
112 the variability among animals is so great, that it is unlikely that even a 3-fold increase in
113 response will be detected statistically, with 3-5 animals per group. For example, if the
114 controls show a range of variability similar to that seen in the first DAE 433 study,
115 ranging from 694 to 20171, a 29-fold difference, how realistic would it be to expect to
116 detect a much smaller (3-fold) increase in the treated group response relative to the
117 response seen in that control group?

118 Thus, based on these data, the only way to assure decent power for this assay is to use
119 DMSO as the vehicle. If this vehicle is used, there is an excellent chance of detecting a
120 2.5-fold or a 3-fold increase in response if 4 or 5 animals per group are used. Another
121 advantage of using DMSO is that the variability within a study among control animals is
122 very reproducible, and thus predictable. In 24 studies using DMSO as the vehicle, the
123 within study variability among animals never exceeded a 4-fold difference. DMSO does
124 not show the wild fluctuations seen for the other vehicles in which one experiment can
125 show a 29-fold variation among control animals and the next experiment show only a 2-
126 fold variation.

127 It should be noted that the mean control response using the DMSO vehicle is much
128 greater than the mean control response using the other vehicles, so a 3-fold increase
129 relative to the DMSO control reflects a much larger actual dosed group response than a 3-
130 fold increase relative to a smaller control response. For example, the mean DMSO
131 control response is almost a 3-fold increase relative to the mean EtOH control response,
132 so a 3-fold increase relative to a DMSO control group would be almost a 9-fold increase
133 relative to the EtOH control group.

134 Finally, regardless of vehicle, it is unlikely that the assay can detect statistically a 2-fold
135 or less increase in response, with only 3-5 animals per group.
136

137

137  **3.0    LLNA: DA**

138  This analysis was based on vehicle control data (ATP levels) from 18 different
139  experiments. Within each experiment, there were three or four animals and two replicates
140  per animal. For each animal, the two replicates were averaged, and then the individual
141  animal means were averaged to obtain overall control means and SD's for that
142  experiment. The data were also log-transformed data and then averaged.  The summary
143  statistics are given in **Table 3-1**, and **Table 3-2** summarizes the average means and
144  standard deviations for each vehicle.

145

146  **Table 3-1       Summary of the Control Absorbance Data for the LLNA: DA**

147

| Experiment | Original scale | | Log scale | | Vehicle |
|---|---|---|---|---|---|
| | Mean | SD | Mean | SD | |
| 1 | 4410 | 752.6 | 8.381 | 0.1801 | AOO |
| 2 | 3871 | 343.8 | 8.258 | 0.0890 | AOO |
| 3a | 3014 | 435.9 | 8.003 | 0.1475 | AOO |
| 3b | 6674 | 1526.6 | 8.785 | 0.2384 | DMSO |
| 4a | 2580 | 517.7 | 7.838 | 0.2229 | AOO |
| 4b | 3465 | 888.9 | 8.124 | 0.2737 | DMF |
| 5 | 5168 | 4579.3 | 8.260 | 0.8812 | AOO |
| 6 | 3528 | 1880.8 | 8.040 | 0.6654 | AOO |
| 7 | 1509 | 455.0 | 7.275 | 0.3666 | AOO |
| 8 | 2668 | 1019.7 | 7.835 | 0.3804 | DMF |
| 9 | 2077 | 95.0 | 7.638 | 0.0452 | AOO |
| 10 | 3129 | 848.7 | 8.023 | 0.2537 | AOO |
| 11 | 2818 | 567.4 | 7.928 | 0.2010 | AOO |
| 12 | 2151 | 376.9 | 7.662 | 0.1740 | AOO |
| 13 | 1611 | 423.7 | 7.362 | 0.2419 | ACE |
| 14 | 3362 | 736.3 | 8.103 | 0.2083 | AOO |
| 15a | 10204 | 2765.9 | 9.203 | 0.2727 | DMSO |
| 15b | 4907 | 656.4 | 8.491 | 0.1422 | AOO |
| 16 | 2710 | 822.5 | 7.875 | 0.2716 | ACE |
| 17 | 64899 | 18696.8 | 11.047 | 0.3063 | DMSO |
| 18 | 2894 | 954.5 | 7.932 | 0.3165 | AOO |

148

149  **Table 3-2       Average Means and Standard Deviations for Each Vehicle**

| Vehicle | N | Original Scale Averages | | Log-transformed Averages | | Converted Control Mean[1] |
|---|---|---|---|---|---|---|
| | | Mean | SD | Mean | SD | |
| AOO | 14 | 3244 | 942.9 | 7.988 | 0.2781 | 2945 |
| DMSO | 3 | 27259 | 7663.1 | 9.678 | 0.2725 | 15968 |
| ACE | 2 | 2160 | 623.1 | 7.619 | 0.2568 | 2036 |
| DMF | 2 | 3066 | 954.3 | 7.980 | 0.3271 | 2920 |

150  [1]Used in power calculations and based on log-transformed scale average.

151

152  Clearly, the DMSO vehicle produces responses totally inconsistent with the other three
153  vehicles (which are reasonably similar among themselves). The power calculations given
154  in **Tables 3-3** to **3-6** are based on a one-sided $p<0.05$ Student's t test applied to the log
155  transformed data (just as in the previous power calculations). ***It should be noted that***
156  ***these calculations make the additional assumption that any "treatment effect"***

157    **produced will have essentially the same SD (on a log-transformed scale) as the control**
158    **data, i.e., that the treatment will change only the mean response and not the variability.**
159    The data in the table above are consistent with this assumption, since although the mean
160    response for the DMSO vehicle is a sizable increase over the mean response for the other
161    controls, the underlying variability (on a log scale) is very similar. The power
162    calculations are summarized below by vehicle.

163

164    **Table 3-3**      **Treatment Group (Rx) Response Increase Relative to Controls for**
165                   **AOO**

| Parameter | 3-fold Increase | 2.5-fold Increase | 2.0-fold Increase | 1.5-fold Increase | 1.3-fold Increase |
|---|---|---|---|---|---|
| Mean Rx response | 8835 | 7362.5 | 5890 | 4417.5 | 3828.5 |
| Log (mean Rx response) | 9.086 | 8.904 | 8.681 | 8.393 | 8.250 |
| Difference from control (log scale) | 1.098 | 0.916 | 0.693 | 0.405 | 0.262 |
| SD of the difference from control | 3.95 | 3.29 | 2.49 | 1.46 | 0.94 |
| Power for N=5 | 99% | 99% | 95% | 50-80% | <50% |
| Power for N=4 | 99% | 95-99% | 90% | 50% | <50% |
| Power for N=3 | 95% | 90-95% | 80% | <50% | <50% |
| Other power | | | | 95% (N=11) | 95% (N=25) |
| Other power | | | | 90% (N=9) | 90% (N=20) |

166

167    **Table 3-4**      **Treatment Group (Rx) Response Increase Relative to Controls for**
168                   **ACE**

| Parameter | 3-fold Increase | 2.5-fold Increase | 2.0-fold Increase | 1.5-fold Increase | 1.3-fold Increase |
|---|---|---|---|---|---|
| Mean Rx response | 6108 | 5090 | 4072 | 3054 | 2646.8 |
| Log (mean Rx response) | 8.717 | 8.535 | 8.312 | 8.024 | 7.881 |
| Difference from control (log scale) | 1.098 | 0.916 | 0.693 | 0.405 | 0.262 |
| SD of the difference from control | 4.28 | 3.57 | 2.70 | 1.58 | 1.02 |
| Power for N=5 | 99% | 99% | 95-99% | 50-80% | <50% |
| Power for N=4 | 99% | 99% | 90-95% | 50% | <50% |
| Power for N=3 | 99% | 95% | 80-90% | <50% | <50% |
| Other power | | | | 95% (N=10) | 95% (N=23) |
| Other power | | | | 90% (N=8) | 90% (N=18) |

169

170

170  **Table 3-5      Treatment Group (Rx) Response Increase Relative to Controls for**
171  **DMF**

| Parameter | 3-fold Increase | 2.5-fold Increase | 2.0-fold Increase | 1.5-fold Increase | 1.3-fold Increase |
|---|---|---|---|---|---|
| Mean Rx response | 8760 | 7300 | 5840 | 4380 | 3796 |
| Log (mean Rx response) | 9.078 | 8.896 | 8.672 | 8.385 | 8.242 |
| Difference from control (log scale) | 1.098 | 0.916 | 0.692 | 0.405 | 0.262 |
| SD of the difference from control | 3.36 | 2.80 | 2.12 | 1.24 | 0.80 |
| Power for N=5 | 99% | 95-99% | 90% | 50% | <50% |
| Power for N=4 | 95-99% | 90-95% | 80% | <50% | <50% |
| Power for N=3 | 90-95% | 80-90% | 50% | <50% | <50% |
| Other power | | | | 95% (N=15) | 95% (N=35) |
| Other power | | | | 90% (N=12) | 90% (N=28) |

172

173  **Table 3-6      Treatment Group (Rx) Response Increase Relative to Controls for**
174  **DMSO**

| Parameter | 3-fold Increase | 2.5-fold Increase | 2.0-fold Increase | 1.5-fold Increase | 1.3-fold Increase |
|---|---|---|---|---|---|
| Mean Rx response | 47904 | 39920 | 31936 | 23952 | 20758.4 |
| Log (mean Rx response) | 10.777 | 10.595 | 10.371 | 10.084 | 9.941 |
| Difference from control (log scale) | 1.099 | 0.917 | 0.693 | 0.406 | 0.263 |
| SD of the difference from control | 4.03 | 3.37 | 2.54 | 1.49 | 0.97 |
| Power for N=5 | 99% | 99% | 95% | 50-80% | <50% |
| Power for N=4 | 99% | 95-99% | 90% | 50% | <50% |
| Power for N=3 | 95% | 90-95% | 80% | <50% | <50% |
| Other power | | | | 95% (N=11) | 95% (N=24) |
| Other power | | | | 90% (N=9) | 90% (N=19) |

175
176  Therefore, using three to five animals per group (and two replicates per animal), there is a
177  very high probability that a 2.5-fold and a 3-fold increase will be detected and a good
178  chance that a 2-fold increase will be detected, regardless of vehicle. However, detecting a
179  1.3 to 1.5-fold increase may be too much to expect with only three to five animals per
180  group.

181

182  Note that all four vehicles produce similar power profiles. This is because the
183  transformed SDs in the table above are all very similar; if they were identical, so would
184  be the power profiles. However, the actual magnitudes of the treated group responses for
185  a given power will differ from vehicle to vehicle because the control responses
186  themselves differ (especially for DMSO). For example, a 3-fold increase in the control
187  response for the ACE vehicle would be an increase from 2036 to 6108, and would be
188  detected with approximately a 99% probability. However, a 6108 treatment response
189  relative to the AOO vehicle would only be approximately a 2-fold increase and would be

190    detected with only a 95% probability. A treatment response of 6108 for the DMSO
191    vehicle would actually be far below the DMSO control response.

192

193    Finally, Experiments 5-8 produced notably more variability (among and within animals)
194    than the other experiments. I cannot help but wonder if these four studies were done at a
195    different lab than the others. If so, then the power specific to that lab would be notably
196    lower than that currently reported, while the power associated with the other experiments
197    would be increased slightly if the four experiments were excluded.

198

199    **4.0    Traditional LLNA**

200    These control data come from three different labs, but the same vehicle was used. The
201    raw data are decays per minute (dpm) from a scintillation counter. Within each
202    experiment, there were five animals and one replicate per animal. For each animal, the
203    five animals were averaged to obtain overall control means and SD's for that experiment.
204    The log-transformed data were also averaged.  The summary statistics are given in **Table**
205    **4-1**.

206    **Table 4-1    Summary of the Control DPM Data for the LLNA**

| Experiment | Original scale | | Log scale | | Range of responses | Lab |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | **Mean** | **SD** | **Mean** | **SD** | | |
| 1A | 443.4 | 233.86 | 5.976 | 0.5531 | | 1 |
| 1B | 410.2 | 100.30 | 5.994 | 0.2421 | | 1 |
| 1C | 462.2 | 172.26 | 6.078 | 0.3874 | | 1 |
| 1D | 397.8 | 92.64 | 5.968 | 0.2092 | | 1 |
| 1E | 466.8 | 154.26 | 6.104 | 0.3262 | | 1 |
| 1F | 352.6 | 118.53 | 5.826 | 0.3211 | | 1 |
| 1G | 333.0 | 167.74 | 5.702 | 0.5336 | | 1 |
| 2A | 487.8 | 164.01 | 6.142 | 0.3649 | | 2 |
| 2B | 729.2 | 314.07 | 6.496 | 0.5214 | | 2 |
| 2C | 586.6 | 279.96 | 6.296 | 0.4252 | | 2 |
| 2D | 618.4 | 103.27 | 6.416 | 0.1644 | | 2 |
| 2E | 487.4 | 80.26 | 6.178 | 0.1585 | | 2 |
| 2F | 304.1 | 208.62 | 5.402 | 0.9937 | | 2 |
| 2G | 309.4 | 110.19 | 5.686 | 0.3512 | | 2 |
| 3A | 330.5 | 145.26 | 5.706 | 0.5184 | 137.67 to 515.98 | 3 |
| 3B | 288.5 | 229.15 | 5.338 | 1.0113 | 42.13 to 654.45 | 3 |
| 3C | 152.5 | 31.78 | 5.008 | 0.2275 | 103.56 to 189.17 | 3 |
| 3D | 296.2 | 126.07 | 5.604 | 0.4820 | 131.13 to 447.97 | 3 |
| 3E | 215.3 | 149.44 | 5.104 | 0.9148 | 38.62 to 437.46 | 3 |

207

208    Power calculations were carried out for each lab separately and for all labs combined.
209    The summary statistics are given below.

210    **Table 4-2        Average Means and Standard Deviations for Each Vehicle**

| Lab | N | Original Scale Averages | | Log-Transformed Scale Averages | | Converted Control Mean[1] |
|---|---|---|---|---|---|---|
| | | Mean | SD | Mean | SD | |
| 1 | 7 | 409.4 | 148.51 | 5.950 | 0.3675 | 383.8 |
| 2 | 7 | 503.3 | 180.05 | 6.088 | 0 .4256 | 440.5 |
| 3 | 5 | 256.6 | 136.34 | 5.352 | 0.6308 | 211.0 |
| All 3 | 19 | 403.8 | 156.93 | 5.843 | 0.4582 | 344.8 |

211    [1]Used in power calculations and based on log-transformed scale average.
212

213    The power calculations given in **Tables 4-3** to **4-6** are based on a one-sided $p<0.05$
214    Student's t test applied to the log-transformed data (just as in the previous power
215    calculations). *It should be noted that these calculations make the additional assumption*
216    *that any "treatment effect" produced would have essentially the same SD (on a log-*
217    *transformed scale) as the control data (i.e. that the treatment will change only the*
218    *mean response and not the variability)*.

219

220    **Table 4-3        Treated Group (Rx) Response Increase Relative to Controls: Lab 1**

| | 3.0-fold increase | 2.5-fold increase | 2.0-fold increase | 1.5-fold increase | 1.3-fold increase |
|---|---|---|---|---|---|
| Mean Rx response | 1151.4 | 959.5 | 767.6 | 575.7 | 498.94 |
| Log (Mean Rx response) | 7.049 | 6.866 | 6.643 | 6.356 | 6.212 |
| Difference (log scale) | 1.099 | 0.916 | 0.693 | 0.406 | 0.262 |
| Difference/SD | 2.99 | 2.49 | 1.89 | 1.10 | 0.71 |
| Power for N=5 | 99% | 95% | 80% | <50% | <50% |
| Power for N=4 | 95% | 90% | 50-80% | <50% | <50% |
| Power for N=3 | 90% | 80% | 50-80% | <50% | <50% |
| Other Power | | | | 95% (N=19) | 95% (N=45) |
| Other Power | | | | 90% (N=15) | 90% (N=36) |

221

222

222      **Table 4-4          Treated Group (Rx) Response Increase Relative to Controls: Lab 2**

|  | **3.0-fold increase** | **2.5-fold increase** | **2.0-fold increase** | **1.5-fold increase** | **1.3-fold increase** |
|---|---|---|---|---|---|
| Mean Rx response | 1321.5 | 1101.25 | 881.0 | 660.75 | 572.65 |
| Log (Mean Rx response) | 7.187 | 7.004 | 6.781 | 6.493 | 6.350 |
| Difference (log scale) | 1.099 | 0.916 | 0.693 | 0.405 | 0.262 |
| Difference/SD | 2.58 | 2.15 | 1.63 | 0.95 | 0.62 |
| Power for N=5 | 95% | 90% | 50-80% | <50% | <50% |
| Power for N=4 | 90% | 80% | 50% | <50% | <50% |
| Power for N=3 | 80% | 50-80% | <50% | <50% | <50% |
| Other Power |  |  |  | 95% (N=25) | 95% (N=57) |
| Other Power |  |  |  | 90% (N=20) | 90% (N=46) |

223

224      **Table 4-5          Treated Group (Rx) Response Increase Relative to Controls: Lab 3**

|  | **3.0-fold increase** | **2.5-fold increase** | **2.0-fold increase** | **1.5-fold increase** | **1.3-fold increase** |
|---|---|---|---|---|---|
| Mean Rx response | 633.0 | 527.5 | 422.0 | 316.5 | 274.3 |
| Log (Mean Rx response) | 6.450 | 6.268 | 6.045 | 5.757 | 5.614 |
| Difference (log scale) | 1.098 | 0.916 | 0.693 | 0.405 | 0.262 |
| Difference/SD | 1.74 | 1.45 | 1.10 | 0.64 | 0.42 |
| Power for N=5 | 80% | 50-80% | <50% | <50% | <50% |
| Power for N=4 | 50-80% | 50% | <50% | <50% | <50% |
| Power for N=3 | 50% | <50% | <50% | <50% | <50% |
| Other Power |  |  | 95% (N=19) | 95% (N=53) | 95% (N>100) |
| Other Power |  |  | 90% (N=15) | 90% (N=43) | 90% (N=100) |

225

226

226  **Table 4-6      Treated Group (Rx) Response Increase Relative to Controls:**
227  **                Combined Labs 1, 2, and 3**

|  | 3.0-fold increase | 2.5-fold increase | 2.0-fold increase | 1.5-fold increase | 1.3-fold increase |
|---|---|---|---|---|---|
| Mean Rx response | 1034.4 | 862.0 | 689.6 | 517.2 | 448.24 |
| Log (Mean Rx response) | 6.942 | 6.759 | 6.536 | 6.248 | 6.105 |
| Difference (log scale) | 1.099 | 0.916 | 0.693 | 0.405 | 0.262 |
| Difference/SD | 2.40 | 2.00 | 1.51 | 0.88 | 0.57 |
| Power for N=5 | 95% | 80-90% | 50-80% | <50% | <50% |
| Power for N=4 | 90% | 80% | 50% | <50% | <50% |
| Power for N=3 | 50-80% | 50-80% | <50% | <50% | <50% |
| Other Power |  |  | 95% (N=11) | 95% (N=29) | 95% (N=68) |
| Other Power |  |  | 90% (N=9) | 90% (N=23) | 90% (N=54) |

228

229  These data show considerable variability, and thus, the power is relatively low.  Labs 1
230  and 2 have a reasonably good (but not great) chance of detecting 3 and 2.5-fold increases
231  if N=4 or N=5 are used. Lesser increases will likely be missed. N=3 also appears to be
232  inadequate.

233

234  Lab 3 will likely be unable to detect any increase of 3-fold or less, even with N=5. The
235  best case is a power of approximately 80% for detecting a 3-fold increase with N=5. The
236  reason for the low power is the high within-study variability.  For example, 2 of the 5
237  experiments at this lab had 11-fold and 15-fold differences among the control responses.
238  If the control responses can differ by a factor of 15, how reasonable is it to expect to
239  detect a 3-fold increase in a treatment group with only 3-5 animals?  Because of Lab 3's
240  poor performance, the "all labs combined" performance suffers as well (**Table 4-6**).

241

242  The power calculations presented above assume that the data will be subjected to some
243  formal statistical test at a pre-specified level of significance (e.g., $p<0.05$).  However, it is
244  also possible for an interpretative strategy to adopt a strict decision rule, such as the
245  following, which I will refer to in this report as the "Ratio Rule":

246
247  "Declare the result positive if the ratio of mean treated response to mean control response
248  is greater than 3; otherwise, declare the response negative".

249
250  One advantage of the Ratio Rule is that it is easy to understand and to apply and requires
251  no statistical test, simply a calculation of means and a ratio.  One disadvantage of the
252  Ratio Rule is that the false positive rate (i.e., the "p value" associated with this decision
253  making strategy) is unknown and will vary from assay to assay, depending upon the
254  underlying variability among animals.  The associated power is also unknown.

255
256  To investigate this matter further, I looked at the ELISA data again, searching for an

257    example showing approximately 95% power based on a Student's t test, so that I could
258    investigate whether this power could be increased or decreased by application of the
259    Ratio Rule.  For the ELISA data, the N=3 case had approximately a 95% power
260    associated with a one-sided Student's t test for detecting a 3-fold increase in response
261    (see **Table 1-2**).  To compare this power with the "Ratio Rule", I made the following
262    assumptions/calculations.
263
264    I assumed that the mean logged response for ELISA was -2.02 and the mean SD response
265    was 0.302 (as before).  The standard error (SE) associated with N=3 is simply the
266    standard deviation divided by the square root of 3 or 0.1744.  This SE is the SD we would
267    expect to see among (logged) mean responses based on N=3.
268
269    I then enumerated (using the cumulative normal probability distribution at probability
270    intervals of 0.02) the approximate distribution of mean log responses consistent with an
271    SD of 0.1744 and a 3-fold increase in the ratio.  That is, I approximated the continuous
272    distribution of both the treated and control responses by a discrete distribution of 50 mean
273    responses, spaced so that each outcome has approximately a 2% probability of
274    occurrence.  These two distributions are given below.  If you calculate the summary
275    statistics, you will find that the mean of the logged control response is -2.02 and the mean
276    of the logged treated group response is -0.92 (a 3-fold increase on the original scale);
277    both have a SD of 0.1744.  Importantly, these are expected mean responses for a group of
278    3 animals, not individual animal responses, so the range of responses is relatively narrow.
279    These distributions formed the basis of the new power calculations.
280

| 281 | Control Mean | | Treated Mean | | Contribution to power |
|---|---|---|---|---|---|
| 282 | Logged | Response | Logged | Response | (Control) for detecting |
| 283 | Response | | Response | | a 3-fold increase |
| 284 | -1.617 | .198 | -0.517 | .596 | .02 |
| 285 | -1.69 | .185 | -0.59 | .554 | .02 |
| 286 | -1.73 | .177 | -0.63 | .533 | .06 |
| 287 | -1.76 | .172 | -0.66 | .517 | .08 |
| 288 | -1.79 | .167 | -0.69 | .502 | .10 |
| 289 | -1.81 | .164 | -0.71 | .492 | .10 |
| 290 | -1.82 | .162 | -0.72 | .487 | .14 |
| 291 | -1.84 | .159 | -0.74 | .477 | .14 |
| 292 | -1.85 | .157 | -0.75 | .472 | .18 |
| 293 | -1.87 | .154 | -0.77 | .463 | .20 |
| 294 | -1.88 | .153 | -0.78 | .458 | .20 |
| 295 | -1.89 | .151 | -0.79 | .454 | .24 |
| 296 | -1.90 | .150 | -0.80 | .449 | .24 |
| 297 | -1.91 | .148 | -0.81 | .445 | .28 |
| 298 | -1.92 | .147 | -0.82 | .440 | .28 |
| 299 | -1.93 | .145 | -0.83 | .436 | .32 |
| 300 | -1.94 | .144 | -0.84 | .432 | .32 |
| 301 | -1.95 | .142 | -0.85 | .427 | .36 |
| 302 | -1.96 | .141 | -0.86 | .423 | .36 |

| | | | | | |
|---|---|---|---|---|---|
| 303 | -1.97 | .139 | -0.87 | .419 | .40 |
| 304 | -1.98 | .138 | -0.88 | .415 | .42 |
| 305 | -1.99 | .137 | -0.89 | .411 | .42 |
| 306 | -2.00 | .135 | -0.90 | .407 | .46 |
| 307 | -2.01 | .134 | -0.91 | .403 | .48 |
| 308 | -2.016 | .133 | -0.916 | .400 | .50 |
| 309 | -2.024 | .132 | -0.924 | .397 | .52 |
| 310 | -2.03 | .131 | -0.93 | .395 | .54 |
| 311 | -2.04 | .130 | -0.94 | .391 | .56 |
| 312 | -2.05 | .129 | -0.95 | .387 | .56 |
| 313 | -2.06 | .127 | -0.96 | .383 | .60 |
| 314 | -2.07 | .126 | -0.97 | .379 | .62 |
| 315 | -2.08 | .125 | -0.98 | .375 | .62 |
| 316 | -2.09 | .124 | -0.99 | .372 | .64 |
| 317 | -2.10 | .122 | -1.00 | .368 | .68 |
| 318 | -2.11 | .121 | -1.01 | .364 | .70 |
| 319 | -2.12 | .120 | -1.02 | .361 | .72 |
| 320 | -2.13 | .119 | -1.03 | .357 | .72 |
| 321 | -2.14 | .118 | -1.04 | .353 | .74 |
| 322 | -2.15 | .116 | -1.05 | .350 | .78 |
| 323 | -2.16 | .115 | -1.06 | .346 | .80 |
| 324 | -2.17 | .114 | -1.07 | .343 | .82 |
| 325 | -2.19 | .112 | -1.09 | .336 | .82 |
| 326 | -2.20 | .111 | -1.10 | .333 | .84 |
| 327 | -2.22 | .109 | -1.12 | .326 | .86 |
| 328 | -2.23 | .108 | -1.13 | .323 | .88 |
| 329 | -2.25 | .105 | -1.15 | .317 | .92 |
| 330 | -2.28 | .102 | -1.18 | .307 | .94 |
| 331 | -2.31 | .099 | -1.21 | .298 | .96 |
| 332 | -2.35 | .095 | -1.25 | .287 | .98 |
| 333 | -2.423 | .089 | -1.323 | .266 | .98 |

334
335                                    Total=25.12
336                                    Power =  0.02 x Total or 50.24%
337

338   Thus, the power is reduced from 95% to 50% by using the Ratio Rule rather than a one-
339   sided p<0.05 Student's t test, although the "gain" is that the false positive rate is reduced
340   from 5% to essentially zero (note from the distributions given above that the overall
341   range of mean control responses is less than 3-fold, so the false positive rate is essentially
342   zero). This "tradeoff" is typical, even for a formal statistical test.  What is needed is a
343   reasonable balance between false positive and false negative rates, and the Ratio Rule
344   seems designed to sacrifice power for the sake of maintaining a low false positive rate.
345
346   One way to modify the Ratio Rule to increase its power would be change the critical
347   value of the ratio from 3 to some smaller number such as 2 or 2.5.  This would increase
348   power while still keeping the false positive rate low.

349
350     For example, by my calculations, if the Ratio Rule applied to the distribution data above
351     was changed from "Ratio > 3" to "Ratio > 2", then the power would be approximately
352     95%, but the false positive rate would still be low (approximately 0.002).
353
354     The 50% power found by enumerating the entire distribution for the example above
355     simply confirms what should be intuitive for the Ratio Rule, namely, that if you have two
356     distributions for which the underlying means differ by a factor of 3, then approximately
357     half the time the ratio of means from sampled data will exceed 3 and approximately half
358     the time it will be less than 3.  So it is unnecessary to perform additional power
359     calculations for the Ratio Rule, at least for detecting an underlying 3-fold increase in
360     response.  Regardless of the underlying SD (and for that matter, regardless of the number
361     of animals used), the power of the Ratio Rule for detecting a 3-fold increase in response
362     will always be approximately 50%.  Of course, if the underlying ratio is greater than 3,
363     then the sample size and underlying variability do become important in the Power
364     Calculations for the Ratio Rule.
365
366     The power of Student's t test depends upon the sample size and the underlying
367     variability, but for the various cases considered (see tables above), the power was always
368     well above 50%.
369
370     My conclusion is that the "Ratio Rule" has a much lower false positive rate than a formal
371     statistical test, but it also has a much higher false negative rate (i.e., lower power).  This
372     reduced power can be considerable, and the Ratio Rule will always show approximately
373     50% power for detecting an underlying treatment effect that on average shows a 3-fold
374     increase relative to controls.  Moreover, the power of the Ratio Rule is less than 50% for
375     detecting increases in the ratio of 2.5, 2, 1.5, or 1.3, but use of the Ratio Rule implies that
376     such increases are likely not biologically important anyway, as discussed in more detail
377     below.
378
379     The ultimate objective of a decision strategy is to maximize the ability of an assay's
380     outcome to predict correctly the human response (positive or negative), and to achieve
381     this objective, a formal statistical test may or may not be necessary. It is my
382     understanding that the "Ratio Rule" was not established arbitrarily, but rather was
383     derived empirically, on the basis that 3 was the "cut-off ratio value" that provided the
384     optimal performance of the assay when differentiating "true" human positives from
385     "true" human negatives for one of the assays.  It is also my understanding that this Ratio
386     Rule has not been "validated" empirically for all of the various assays to which it is now
387     being routinely applied (ELISA, traditional, DA, FC, etc.).
388
389     If the Ratio Rule seems to "work" very well in practice in predicting the human response,
390     that is the ultimate goal, so there may be no need of a formal statistical test, as long as
391     everyone fully understands what the use of such a rule implies.  Since, based on the
392     control data provided to me, a false positive outcome is nearly impossible (or at least has
393     a very low probability) using the Ratio Rule, use of this rule implicitly assumes that are
394     some, perhaps even many, compounds that are "true positives" in the assay, but the

395    response that they produce (e.g., a 2-fold or 2.5-fold increase in the treated/control ratio),
396    while detectable statistically, should be considered a negative response, since it is of
397    insufficient magnitude for the compound in question to be positive in humans.  Even a 3-
398    fold increase in the ratio of treated to control mean response is considered relatively
399    unimportant, since it will be detected only approximately 50% of the time by the Ratio
400    Rule.  Is such a performance acceptable to the scientific community?  Are chemicals that
401    are truly active in the assay, but produce a ratio of <3, generally negative in humans and
402    thus can be discounted when a response of this magnitude is observed in an assay?  Use
403    of the Ratio Rule assumes that the answer to this question is Yes.
404
405    To summarize, use of the Ratio Rule assumes that there are compounds that are
406    statistically positive in the assay (and are actually producing an effect in the assay), but
407    the magnitude of the effect is insufficient for the compound to be positive in humans.
408    Unless this is known with certainty, I personally prefer using a formal statistical test
409    rather than a strict rule, a rule whose performance characteristics (power, false positive
410    rate) are unknown and vary from assay to assay.
411
412    **5.0 Final Comments and Summary**
413
414    (1)  One result of the data analyses presented above is that it reinforces the need for
415    concurrent control data.  **Table 5-1** below shows the variability observed in the mean
416    control responses across experiments.  Only for the traditional assay are the results
417    reasonably reproducible.  For the other assays, concurrent controls are clearly essential,
418    since the data are so variable across experiments, and I would recommend that concurrent
419    controls be routinely included in the study design of all assays.  Note that in many of the
420    assays, a control response in one experiment would clearly be considered "active"
421    relative to the control response in another experiment, since the ratio is far greater than 3.
422
423    **Table 5-1:  Variability in the mean control response across experiments**
424
425                         Maximum Difference Among
426                             Control Means
427
428    ELISA                    7-fold difference
429    FC:  DAE Vehicle         6-fold difference
430    FC:  AOO Vehicle         15-fold difference
431    FC:  DMSO Vehicle        7-fold difference
432    FC:  DMF Vehicle         6-fold difference
433    FC:  ETOH Vehicle        5-fold difference
434    DA:  AOO Vehicle         3.4-fold difference
435    DA:  DMSO Vehicle        10-fold difference
436    Traditional:  Lab 1      1.4-fold difference
437    Traditional:  Lab 2      2.4-fold difference
438    Traditional:  Lab 3      2.2-fold difference
439
440    (2)  A related point is that it is important to have individual animal control data, so that

441   the within-study among-animal variability can be assessed and factored into the data
442   evaluation.  Individual animal data are essential if the data are to be evaluated
443   statistically.  Even if the "Ratio Rule" is used, individual animal data are highly desirable.
444
445   (3)  Since a formal statistical test has so much more power than the Ratio Rule, it is
446   definitely of interest to examine those specific compounds for which there are
447   contradictory results, i.e., a statistically verified treatment effect in the assay, but the
448   Ratio Rule criterion is not met.  It is important to determine whether these chemicals are
449   positive or negative in the human setting, to understand if these "statistical positives" in
450   the assay are "false positives" or "true positives" in the human setting.
451
452   (4)  These analyses also reinforce the importance of the choice of vehicle in certain of the
453   assays.  For example, DMSO shows a response that is much higher than that seen with
454   the other vehicles.  The variability in response among animals may also be dependent
455   upon the vehicle in some experimental settings.  Similarly, the data suggest that some
456   labs are better than others in reproducing the control responses among animals within a
457   given experiment (a further argument for the routine reporting of individual animal
458   control data).
459
460   (5)  The decision to use 4 or 5 animals depends upon whether or not the gain in power
461   achieved by the extra animal is deemed sufficiently important to justify the extra time,
462   effort and cost.  In some cases (e.g., the ELISA assay for detecting a 3-fold increase), the
463   extra animal would add little, since the power for N=4 is already 99%.  For other assays
464   that show more variability (see tables above), the extra animal may be more important.  It
465   is a judgment call.
466
467   Importantly, this comparison of sample size is linked to use of a formal statistical test.  If
468   the Ratio Rule is used instead, and power is calculated for a true underlying response
469   ratio of 3, then sample size is irrelevant, since the power will always be approximately
470   50%, regardless of sample size.  Although I have not made sample size comparisons for
471   the power of the Ratio Rule applied when the true underlying ratio exceeds 3, the low
472   power of the Ratio Rule in general suggests the use of as many animals as are feasible, so
473   an N of 5 rather than 4 may be important if the Ratio Rule is used.
474
475   (6)  Finally, the decision whether to use a formal statistical test or the Ratio Rule is
476   beyond the scope of this evaluation.  Since the Ratio Rule has notably less power than a
477   formal statistical test, then the "default" approach in my opinion should be to use a
478   formal statistical analysis, unless it can be demonstrated that the "statistical positives"
479   that are identified in the assay but "missed" by the Ratio Rule are in fact negative in
480   humans.  If such compounds are in fact negative in humans, it would indicate that the
481   assay is "overly sensitive" and detects effects that are not relevant to humans, and this
482   needs to be understood by the scientific community.
483
484   Joe Haseman
485   2-14-08