



King's Research Portal

DOI:

[10.1111/ejn.13847](https://doi.org/10.1111/ejn.13847)

Document Version

Early version, also known as pre-print

[Link to publication record in King's Research Portal](#)

Citation for published version (APA):

Jeffery, N. D., Bate, S. T., Safayi, S., Howard, M., Moon, L., & Jeffery, U. (2018). When neuroscience met clinical pathology: partitioning experimental variation to aid data interpretation in neuroscience. *European Journal of Neuroscience*. Advance online publication. <https://doi.org/10.1111/ejn.13847>

Citing this paper

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

General rights

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

Take down policy

If you believe that this document breaches copyright please contact librarypure@kcl.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.

1 ***TECHNICAL SPOTLIGHT***

2

3

4 **When neuroscience met clinical pathology: partitioning experimental variation to aid data**
5 **interpretation in neuroscience**

6

7

8 Nick D Jeffery^{1*}, Simon T Bate², Sina Safayi³, Matt Howard⁴, Lawrence Moon⁵, Unity Jeffery⁶

9

10

11 ¹Department of Small Animal Clinical Sciences, Texas A&M University, TX 77843, USA

12

13 ²Statistical Sciences, GlaxoSmithKline, Medicines Research Centre, Gunnels Wood Road,
14 Stevenage, Hertfordshire, SG1 2NY, UK

15

16 ³The University of Texas Graduate School of Biomedical Sciences at Houston, 6767 Bertner
17 Ave, Houston, TX 77030, USA

18

19 ⁴Department of Neurosurgery, University of Iowa Hospitals and Clinics, 200 Hawkins Drive,
20 Iowa City, IA 52242, USA

21

22 ⁵Neurorestoration Group, Wolfson Centre for Age-Related Diseases, King's College London,
23 University of London, London SE1 1UL, UK

24

25 ⁶Department of Veterinary Pathobiology, College of Veterinary Medicine, Texas A&M

26 University, College Station, TX 77840

27

28 * **Corresponding author**; address as above; email: njeffery@cvm.tamu.edu

29

30 **Running title:** Analysis by reference change value

31

32 **TOTALS:**

33

34 Pages: 38

35 Figures: 3

36 Tables: 3

37 Equations: 4

38 Words: a) whole manuscript: 4739 (excluding references, abstract
39 and figure legends)

40 b) Abstract: 228

41

42

43 **Keywords:** translation, reference change value, intra-individual, partition, variability

44 **Abstract**

45 Neuroscientists typically assess the effectiveness of interventions by comparing the average
46 response of groups of treated and untreated animals. These group comparisons provide useful
47 insights into the effectiveness of interventions, but focusing only on group effects risks
48 overemphasis of small, statistically significant but physiologically unimportant differences. Such
49 differences could be created by analytical variability or by physiological within-individual
50 variation, especially when the number of animals in each group is so small that one or two outlier
51 values have considerable impact on the summary measures for the group.

52

53 Physicians face a similar dilemma when comparing two results from the same patient. To
54 determine if the change between two values reflects disease progression or expected analytical
55 and physiologic variation, the magnitude of the difference between two results is compared to
56 the reference change value. These values are generated by quantifying analytical and within-
57 individual variation, and differences between two results from the same patient are considered
58 clinically meaningful only if they exceed the combined effect of these two sources of “noise”.

59

60 In this article we describe how this technique can be extended into neuroscience. This form of
61 analysis provides a measure of outcome at an individual level that complements traditional
62 group-level statistics and therefore we hope that introduction of this technique into neuroscience
63 will enrich interpretation of experimental data. In addition, it can also safeguard against some of
64 the possible misinterpretations that may occur during analysis of the small experimental groups
65 that are common in neuroscience and, by illuminating analytical error, may aid in design of more
66 efficient experimental methods.

67 **Introduction**

68 The repeated failure of neuroscience to generate novel clinical therapies has been a source of
69 unease for many years (Garner, 2014) and this disquiet has peaked recently as deficiencies in
70 experimental design and analysis have become more widely apparent (Begley, 2013; Button *et*
71 *al.*, 2013; Steward, 2016; Academy of Medical Sciences, 2015). Many remedies for
72 methodological shortcomings are available, focusing on both animal model relevance and
73 experimental design (Garner, 2014, Begley, 2013; Button *et al.*, 2013; Steward, 2016; Academy
74 of Medical Sciences, 2015; Landis, 2012; ARRIVE, 2017). Here we take another approach, by
75 considering how data analysis strategies that are routinely used in hospital laboratories might
76 complement traditional evaluation of neuroscience data.

77

78 In biomedical studies it is important to consider the magnitude of an effect because interventions
79 with greater effect will likely have greater impact in the clinic. Routinely, the effect is
80 determined by assessing an outcome measure before and after the intervention (within animal),
81 or by comparing treated *versus* untreated groups of animals. Such analyses imply comparison,
82 usually by statistical tests, of the population distribution and its central tendency (*i.e.* mean,
83 median *etc*) between groups. Unfortunately, while widely used, the derived P values provide
84 only an estimate of the consistency of the data with the null hypothesis; they do not directly
85 provide information on the magnitude of an intervention effect and are difficult to interpret
86 without knowledge of pre-study power calculations (Shaver, 1993; Halsey *et al.*, 2015).

87

88 In part, these problems arise from the considerable variability in data distribution, and group-
89 level analyses do not provide the complete picture. For instance, as highlighted previously in this

90 journal, specific individuals may show changes in outcome that are not reflected in the summary
91 mean or median values (Rousselet et al., 2016). In addition, a commonly overlooked problem is
92 that random sources of variation inherent in animal experiments, such as physiological
93 variability and intrinsic measurement inaccuracies, may combine to produce outcome effects that
94 might erroneously be attributed to the intervention itself, especially in the small sample
95 populations typically used in neuroscience (Button *et al.*, 2013). Currently established statistical
96 analyses in neuroscience do not always efficiently dissect the difference between inherent
97 random sources of variation and the effects induced by an experimental intervention. Individual-
98 level analysis, as described below, helps to guard against drawing inappropriate conclusions
99 from underpowered studies.

100

101 A further aim in translational biomedical research is to identify interventions that will have real-
102 life impact for each treated individual. Whilst it is of course useful and important to show, as
103 conventional group-level analysis can, that, on average, a treated individual will have a better
104 outcome than a non-treated individual it is also important to examine effects at an individual
105 level, because each individual patient wants their life to be meaningfully enhanced. The two
106 approaches are complementary: conventional group-level analysis can provide the information as
107 to whether, on average, one treatment is better than another and provide an estimate of how
108 much better. Individual-level analysis provides an estimate of: a) how many individuals attain a
109 benefit that is greater than that which can expected through variation attributable to experimental
110 inaccuracy and physiological variation; and also, b) the extent that the improvement for each
111 individual is beyond those limits. In this *Technical Spotlight* we explain how an individual-level

112 analysis can be derived and how it aids in dissecting intervention signal from experimental noise,
113 thereby providing a complement to standard group-level analyses.

114

115 **Analysis of individual responses in clinical practice**

116 In clinical medicine the focus is on individual patients and there is frequently a need to determine
117 whether a patient's condition is deteriorating or whether a treatment is having a beneficial effect.

118 Such monitoring involves obtaining serial samples, which then raises the question as to how
119 large the difference between a pair of samples obtained from the same individual has to be before

120 it represents a meaningful change. In clinical pathology laboratories this problem has been

121 addressed by partitioning the various sources of variation and using this to derive boundaries that

122 encompass all sources of extraneous variation. This can then be used to assess the results for

123 each patient individually (Harris & Yasaka, 1983; Fraser & Harris, 1989; Walton, 2012). By

124 comparing the changes in an individual's laboratory test results to these boundaries 'real change'

125 can be inferred. Below we outline how the same concepts can be applied to experimental studies

126 in neuroscience. The combination of group- and individual-level analysis together provides a

127 more comprehensive examination of the data and may therefore aid in establishing which

128 laboratory interventions are most likely to have sufficiently robust effects to translate into

129 effective clinical therapy. This may also aid in design of appropriate pre-clinical functional tests.

130

131 **What is variation?**

132 Variation refers to random fluctuations in repeated results generated for a particular sample or

133 individual. This is distinct from systematic error, also termed bias, which describes consistent

134 under- or over-estimation of the true value (Theodorsson *et al.*, 2014). The standard deviation

135 (SD), or the *coefficient of variation* (CV) (= SD/mean) provide summary measures of variation
136 of the standard mean.

137

138 **Sources of variation**

139 In clinical pathology, when determining the value of a specific analyte the various possible
140 sources of variation are identified and partitioned. There are similar sources of variation in
141 neuroscience.

142

143 ***Investigator-derived pre-analytical variation*** encompasses all the investigator-dependent
144 sources of variation that influence the level of an analyte before measurements are made. For
145 example, variation in how long, or in what conditions, a sample has been stored before testing
146 may influence the final results of the experiment. In general, pre-analytical variation is not of
147 interest to the experimenter, instead forming a source of ‘noise’ in data interpretation, and so
148 should be eliminated as far as possible. In clinical pathology, investigator-dependent pre-
149 analytical variation is minimized by strict handling protocols for sample collection, handling and
150 storage.

151

152 ***Analytical variation*** arises because every analytical technique has intrinsic sources of laboratory
153 variability that can lead to variation in replicate results obtained from a single sample (this is also
154 sometimes referred to as measurement, technical or instrumental error). For instance, small
155 volume errors in loading a sample into an automated analyzer or simple variation in reaction
156 kinetics can produce different results. While analytical variation can be reduced, for example by
157 instrument calibration or cleaning, it cannot be completely eliminated. Nevertheless, by

158 performing multiple replicate analyses of each sample, or more generally by repeatedly
159 measuring the same physical material, it can be quantified and this estimate used when deducing
160 the sources of variation in the final results (see below).

161

162 **Biological variation** can be divided into intra-individual and inter-individual variation. Both
163 forms of biological variation include predictable sources of variation (*e.g.* variation in blood
164 hormone concentration over the reproductive cycle of an individual or variation in blood
165 hormone concentration between individual mice at corresponding points in the reproductive
166 cycle). Researchers typically already pay close attention to these known sources of variation and
167 standardize where possible, but unknown sources of biological variation also have important
168 implications for data interpretation.

169

170 *Intra-individual variation* arises because, as would be expected, there is some variation of
171 analyte levels within each individual between samples obtained at different times that is above
172 and beyond analytical variability. For instance, blood cholesterol concentration varies (even at
173 the same time each day) between one day and another, even in healthy individuals (Rotterdam *et*
174 *al.*, 1987). The magnitude of this variation varies between analytes: there is much less intra-
175 individual variation for some (*e.g.* chloride) than others (*e.g.* triglycerides: Nunes *et al.*, 1993).
176 Importantly, for some analytes the intra-individual variation may be very much less than the
177 inter-individual variation, implying that what might appear to be only a subtle change when
178 viewed in the context of the population as a whole may be critical to the health of a specific
179 individual. For such analytes it is critical to define limits of expected physiological change for
180 each individual and not rely on group-level assessment.

181
182 *Inter-individual variation* arises from the familiar biological observation that, even within a
183 defined population, individuals vary. Indeed, this is the rationale for carrying out experiments
184 using groups of animals, summarizing the outcomes across all members of the group(s) and
185 carrying out the statistical assessment at the population level. While this approach remains valid,
186 in view of the other sources of variation outlined above, apparent intervention effects may be
187 obscured - or magnified - by intra- or inter-individual variation, as well as by analytical and pre-
188 analytical variation. The risk of drawing false conclusions is greatest when the number of
189 experimental subjects within the groups is small, as is common in neuroscience. In addition,
190 because these random sources of variation may not be evenly distributed between groups, they
191 may generate unreliable P-values that may give a false impression of the statistical significance
192 of the overall difference between groups, or between repeated measures within groups (Button *et*
193 *al.*, 2013) (see Figure 1, Table 1). The addition of an analysis of individual animal outcomes, as
194 explained below, can help avoid this problem.

195

196 **How variation is dealt with in clinical pathology laboratories**

197 In clinical pathology laboratories it is common for an individual patient to undergo repeated
198 testing over time. The most useful approach for determining whether there has been a
199 meaningful change between measurements made at two different time points in a single
200 individual is to calculate a *reference change interval* (RCI), derived from calculation of the
201 *reference change value* (RCV) (Fraser, 2001):

202

$$203 \quad RCI = \text{Baseline} \pm (RCV * \text{Baseline}) \quad (1)$$

204 where $RCV = z_p * \sqrt{2} * \sqrt{(CV_I^2 + CV_A^2)}$ (2)

205

206 and CV_I is the intra-individual coefficient of variation and CV_A is the analytical coefficient of
207 variation. The z-score in (2), z_p , can be used to define how confident the experimenter wishes to
208 be that the new result is really different from the previous value. Conventionally, a 5% false
209 positive rate is selected, corresponding to a z-score of 1.96.

210

211 The reference change interval defines the boundaries (usually as a percentage change from a
212 previous result) within which measurements of a single analyte might be expected to vary within
213 a normal individual. By comparing pre- and post-intervention data, the reference change value
214 can also put the effect of an intervention into the context of normal biological variation for an
215 individual, thereby allowing more individualized assessment of the magnitude of the effect.

216

217 Here we show how the same approach can be applied to measures of neurologic function in
218 laboratory animals.

219

220 **Translating clinical pathology analytical methods into neuroscience research**

221 *Pre-analytical variation relating to study design*

222 First, pre-analytical variation should be eliminated as far as possible because its reduction will
223 increase precision. For factors other than homeostatic variation this can largely be achieved by
224 following consistent protocols that should be fully reported and documented to aid transparency
225 (ARRIVE guidelines, Kilkenny *et al.*, 2012). This is most straightforward when applied to
226 familiar laboratory procedures such as immunohistochemistry, for which it is essential to use the

227 same tissue handling techniques and the same batches of antibodies and other reagents between
228 compared samples. The same concept can be applied to behavioral experiments. For instance, as
229 noted elsewhere (National Center for 3Rs, 2017), behavioral tests should be carried out at the
230 same time every day in the same environment with the same schedules and same experimenter.
231 Although attempts should be made to completely eliminate all possible sources of pre-analytical
232 variation (ARRIVE guidelines), it may be impossible (*e.g.* because of extreme weather events,
233 fire alarms *etc.*), in which case undefined sources of pre-analytical variation will usually then be
234 incorporated into other sources of variation (see below).

235

236 ***Derivation of analytical variation***

237 Analytical, or technical, variation is usually determined in clinical pathology by repeatedly
238 measuring analyte levels in one sample by splitting it into multiple aliquots. This approach can
239 easily be applied to routine laboratory techniques like ELISA or PCR, but is rarely appropriate in
240 behavioral neuroscience because repeated examination of one individual may change measured
241 outcomes through training. Although not ideal, it may be possible to obtain a surrogate measure
242 of inherent variation by viewing ‘aliquots’ of a single segment of video footage of various
243 individuals or repeatedly deriving values for kinematic variables from a single original source
244 dataset (and this process is illustrated in Example 1 below).

245

246 Quantification of analytical variation can be an extremely valuable exercise in itself, because it
247 will often provide insights into how the technique and precision of measurements can be
248 improved. However, not all tests in behavioral neuroscience are amenable to similar methods of
249 estimating analytical variation as described above, implying that this source of variation cannot

250 always be partitioned from that derived from intra-individual variation. Despite this apparent
251 obstacle the reference change value can still be calculated (see below) and, importantly, even
252 when the analytical component cannot be explicitly estimated, the reference change interval can
253 still be narrowed (and thus its sensitivity to intervention-induced change increased) by taking
254 steps to minimize analytical error. For example, eliminating any potential handling or
255 environmental conditions that trigger ‘freezing’ of movement in mice (Gouveia & Hurst, 2013)
256 reduces analytical error (and therefore the width of the reference change interval) for motor task
257 outcomes.

258

259 ***Derivation of intra-individual and inter-individual variation***

260 Many experiments involve repeatedly measuring outcomes in a group, or groups, of animals. The
261 inter- and intra-individual variation can be estimated using a nested random effects ANOVA
262 (Fraser and Harris, 1989) or repeated measures mixed models (Bate and Clark, 2014;
263 Marechenko, 2006), which provide flexibility when modelling the correlation between repeated
264 measurements taken on an individual. Typically, estimates of intra- and inter- animal variation
265 are derived from pre-intervention data because it is important that variability is estimated in
266 individuals at a steady plateau of disease or function. Fortunately, experience in clinical
267 pathology has indicated that, for most outcomes, patients with steady-state diseases or lesions
268 have similar variability to normal individuals (Carmen *et al.*, 2007).

269

270 However, in experimental studies an alternative approach might be feasible: instead of using
271 control animals to define the intra-animal variability for all animals in an experiment these
272 measures could instead be derived for each individual. Thus, before making a lesion or

273 implementing a therapeutic intervention it would be possible, through repeated testing, to
274 accumulate sufficient data to calculate the intra-animal variability for each individual animal,
275 thus providing an individualized measure of the intervention effect (Safayi *et al.*, 2015). For
276 instance, a z-score for each outcome can be derived for each animal by re-arranging equation (2)
277 so that:

$$278 \quad z\text{-score} = \textit{proportional change from baseline} / \sqrt{2} * \sqrt{(CV_I^2 + CV_A^2)} \quad (3)$$

279
280 This can then be used to define the magnitude of intervention effect in each individual which is
281 linked to its own biological variation. An intervention that has an effect considerably greater than
282 biological variation is more likely to produce a meaningful improvement in an individual
283 patient's quality of life than an intervention that has effects of similar magnitude to normal day-
284 to-day fluctuations in performance.

285
286 ***Defining the boundaries of outcome that might result from experimental and biological***
287 ***variation alone***

288 For each individual, if a new test result is to represent a 'real' change compared with a previous
289 test result it must lie outside the range that could be expected to arise through the combination of
290 pre-analytical, analytical and intra-individual variation alone. To determine if this is the case, the
291 previously derived values for CV_A and CV_I are combined to derive the reference change value,
292 which can then be applied, as given in equation (2), to define reference change interval
293 boundaries.

294

295 In many neuroscience experiments it may not be possible to derive (and therefore partition) the
296 component of variability that contributes to CV_A , in which case the analytical variation will
297 become a component of CV_I and equation (1) will reduce to:

$$298 \quad RCV = z_p * \sqrt{2} * CV_I \quad (4)$$

299 As before, the boundary values beyond which real change can be deduced to have occurred in an
300 individual (the RCI) is calculated using equation (2).

301

302 **Examples**

303 *Example 1:*

304 *Derivation and use of measures of analytical (CV_A), intra-individual (CV_I) and inter-individual*
305 *(CV_G) variability*

306 Here we use real and simulated data from an experiment to investigate the effect of a drug on the
307 gait of sheep as they walked on a treadmill at constant velocity; the outcome of interest was
308 hindlimb stride duration. Using a nested design analysis module (InVivoStat;
309 <http://invivostat.co.uk/>) the intra- and inter-animal coefficients of variation, denoted by CV_I and
310 CV_G respectively, were estimated from serial measurements made on trained normal animals
311 before drug administration (the raw data is available as Supplementary Material).

312

313 Sheep were trained to walk on a treadmill at a steady speed until they attained a plateau of
314 proficiency. Stride duration was then measured repeatedly on each of 4 days and on each of these
315 days a long recording segment (~ 4 minutes) was divided into four parts, with each part
316 considered as a separate ‘trial’. This is a nested design in which trials are nested within days and
317 days are nested within sheep; the sources of variability were apportioned into analytical

318 ('between-trial, within-day'), intra-individual ('between-day, within-sheep') and inter-individual
319 (between-sheep) to derive means and standard deviations which were used to derive the
320 respective coefficients of variation.

321

322 Inter-individual, intra-individual and analytical standard deviations were 0.039, 0.033 and 0.045
323 respectively and the overall mean stride duration was 0.862, resulting in the following
324 coefficients of variation:

325 $CV_G = 0.039 / 0.862 = 4.5\%$;

326 $CV_I = 0.033 / 0.862 = 3.8\%$;

327 $CV_A = 0.045 / 0.862 = 5.2\%$

328

329 The **reference change value** was 17.8%, calculated using equation (1) above, therefore, any
330 sheep in which the post-drug measurement was changed by more than 17.8% of the pre-drug
331 value (*i.e.* outside the upper or lower boundary of the reference change interval) can be
332 considered to have undergone 'real' change as a consequence of drug administration. It is of
333 course possible that the intervention may have had an effect that produces an outcome that does
334 not lie outside the reference change interval but such an effect can be inferred to be negligible in
335 terms of biomedical meaning (and, in clinical medicine, a drug having an effect of such small
336 magnitude would not be considered valuable).

337

338 To illustrate the different, and complementary, perspectives provided by population-level and
339 RCV within-individual analysis, see Figure 2, which depicts the difference in stride duration in
340 these same sheep before and after drug administration. A Wilcoxon matched pairs signed rank

341 test (post-drug population is not normally distributed) indicates a significant increase in stride
342 duration following this intervention ($p=0.023$).

343

344 When the complementary individual analysis is applied to the same dataset (Table 2), in only 3
345 of 8 sheep does the post-drug value exceed the upper boundary of the reference change interval
346 and indicate a real difference from pre-drug values. Furthermore, the largest drug-induced
347 change was a 27% change from baseline (Sheep 8). When applied in equation (3) this change
348 corresponds to a z-score of 2.97, indicating a value that is in the top 0.3 percentile of the
349 expected distribution of baseline values (see <https://measuringu.com/pcalcz/>).

350

351 Therefore, although there is some evidence of an overall drug effect at the population level, this
352 drug at this dose evoked a real (*i.e.* biomedically meaningful) change in less than 50% of
353 subjects and the largest effect in any single individual was moderate (*i.e.* within the top 0.3
354 percentile of the range of values that could be attributed to biological and analytical variation).

355 Despite the statistically significant result, there is evidence that the drug effect would require
356 considerable augmentation to achieve a worthwhile benefit that would be discernible in a large
357 proportion of subjects. Reporting both analyses alongside each other provides a more
358 comprehensive picture of drug effect. Although tempting, we are not recommending comparison
359 of the proportion of animals that respond to an intervention, because such analysis will inevitably
360 have very low power (Snapinn & Jiang, 2010), and does not take account of the magnitude of
361 effect beyond the reference change interval for each animal.

362

363

364 **Example 2:**

365 This example concerns bladder function in dogs in a clinical trial of a putative therapy for
366 chronic spinal cord injury. A parallel group study was performed with dogs randomized at time 0
367 to receive either percutaneous intraspinal chondroitinase injection or a sham injection in which
368 the skin was pierced with a hypodermic needle. A bladder compliance index (change in pressure
369 with increasing volume of urine) was measured by cystometry at baseline (before intervention)
370 and at 1, 3 and 6 months after intervention (or sham). In this example it is not possible to
371 segregate analytical variation from the inter- and intra-animal variation (because it is not possible
372 to repeatedly measure compliance on 1 day for each animal because this in itself would be
373 expected to cause changes) and so this becomes incorporated into other sources of variability.

374

375 a) First, the data from the **Control** animals over the entire observation period was analyzed using
376 a nested design analysis module (InVivoStat; <http://invivostat.co.uk/>) to segregate intra- and
377 inter-individual variability (the raw data is available as Supplementary Material). With the mean
378 compliance value this was used to derive the intra- and inter- animal coefficients of variation
379 (CV_I and CV_G respectively), where in this example

380

$$\begin{aligned} 381 \quad CV_I &= 0.865 / 1.508 \\ 382 &= 0.574 \end{aligned}$$

383

384 b) The **reference change value** was then calculated from the CV_I using equation (4)

385

$$386 \quad RCV = 1.96 * \sqrt{2} * CV_I$$

387 = 2.764 * 0.574

388 = 1.587

389

390 This implies that, for each dog, for a post-intervention value to be considered beyond the
391 threshold that could be attributed to variability caused by physiological and experimental
392 measurement variability, it would have to be more than 158.7% different from baseline.

393

394 c) We then applied this reference change value to the observed changes in dogs in the
395 *Intervention* group (Table 3). There are two main findings from this analysis:

396

397 i) Two dogs (of 21) showed a change in compliance index that was outside the reference change
398 interval, and all of these were increases. Furthermore, their z-scores were high, lying within the
399 highest 0.5 percentile of the expected distribution, suggesting that it was extremely probable that
400 real change in compliance had occurred during this time period in these dogs.

401

402 ii) No dog showed a reduction in compliance index. However, the analysis suggests that it is
403 probable that this outcome measure is insensitive to that direction of change. The reference
404 change value is large for this test (~159%) which indicates that real change can only be deduced
405 with a large change from baseline. Many of the dogs have low values at baseline meaning that
406 attaining values outside the lower boundary of the reference interval for those individuals is
407 biologically extremely implausible. Therefore this analysis also reveals a floor effect (limitation)
408 in this outcome measure.

409

410 d) This individualized analysis can be compared with conventional group-level analysis: a paired
411 Student's t-test (datasets for baseline and month 1 [above] are both normally distributed)
412 produces a *P* value of 0.139, conventionally regarded as not significantly different. The data
413 from individual dogs is shown in Figure 3.

414

415 The interpretation of these data is therefore that there is no overall evidence of effect of the
416 intervention at a group level but there is evidence that some animals exhibit changes in bladder
417 function that cannot reasonably be attributed to experimental noise (*i.e.* physiological variation
418 and measurement inaccuracies) alone. This might suggest that there are specific individuals that
419 respond to the intervention, which might prompt further investigation of the specific
420 characteristics of these individuals, for instance detailed investigation of the character of their
421 spinal cord injuries.

422

423 However, this individual-level analysis also reveals that there is considerable variability in this
424 outcome measure, which might be a consequence of specific physiological attributes (intra- and
425 inter- animal variability) or of difficulties in making repeatable measurements (analytical
426 variation), each of which will limit the ability to detect real change. Furthermore, there is
427 evidence of a floor effect, which limits the ability to detect reduction from baseline in all
428 animals. Both these features might suggest the need to stratify animals for entry into a study on
429 this outcome to increase the sensitivity of the experiment to physiologically important change.

430

431

432

433 **Discussion: How can RCV-based analysis aid in neuroscience research?**

434

435 *In translational research*

436 Confidence in neuroscience data has been eroded recently, mainly because of the problems that
437 have become apparent in translating apparently positive laboratory results into the clinic and
438 because of widespread difficulties in reproducibility (Garner, 2014). In order to make a
439 successful transition from laboratory to clinic an intervention needs to have a large enough effect
440 to change an individual patient's life in some meaningful way but must also be effective in a
441 suitably large proportion of treated individuals. The analysis strategy that we describe here can
442 aid in overcoming this confidence gap because it provides a more detailed description of the
443 magnitude of the intervention effect in each individual alongside defining the proportion of
444 animals within a treated cohort that show a response beyond a pre-defined level.

445

446 *In laboratory science*

447 The partitioning of sources of variation can be helpful in redesigning pre-clinical outcome
448 measures and defining good laboratory practice to minimize unnecessary variation. For instance,
449 awareness of the concept of pre-analytical variation can aid in eliminating possible causes from
450 laboratory testing procedures, such as those associated with experimenter-centered effects, such
451 as the time of day at which tests are performed. A recent example of the importance of these
452 effects is the recognition that the gender of the researcher can alter mouse behavior (Sorge *et al.*,
453 2014). The large CV_I we describe above for bladder compliance in spinal cord-injured dogs may
454 partly be a result of an unidentified source of pre-analytical and, or, analytical variation, such as
455 urinary tract infection or poor pre-trial management of bladder emptying.

456

457 Analysis of the various sources of variation can also help in determining whether a functional
458 test is sufficiently sensitive to detect a change in the outcome that is being measured. As we
459 show, the CV_I for bladder compliance in spinal cord-injured dogs is large, meaning that, at a
460 group-level, this outcome measure is relatively insensitive to intervention-evoked change. In
461 clinical pathology the aim is for tests on patients to have an analytical component of variation
462 that is too small to have a major impact on clinical decision-making. Although the very small
463 analytical errors required in clinical pathology may well be unattainable for functional tests in
464 neuroscience, tests should be redesigned with the aim to improve precision where possible.
465 Another problem with functional testing that may be revealed is the possibility of floor or ceiling
466 boundaries on the responses that may limit the sensitivity of the assay for the outcome of
467 interest. This was also apparent in the bladder compliance data (example 2 above).

468

469 *Limitations*

470 Unfortunately, not all current outcome measures used in laboratory animals will be appropriate
471 for the analysis strategy outlined in this article. For example, behavioral analytical methods that
472 attribute numeric scores but are not truly numeric because the scores do not represent equally-
473 spaced intervals on a linear scale. This implies that calculation of standard deviation (and
474 therefore the CV) and, consequently, the RCV, can be problematic. Similarly, for some outcome
475 measures the data for computation of the RCV may not be normally distributed. Both of these
476 obstacles can be circumvented. First, log-transformation of positively-skewed data allows
477 application of the techniques we describe above, followed by anti-logging to convert the
478 mathematical answer into the appropriate clinical units. Second, measurement of test-retest

479 differences in a large population (>120) of individuals at a functional plateau can be used to
480 derive the 95% reference change interval of data of any distribution through non-parametric
481 analysis (Friedrichs et al., 2012). In practice, this method would be difficult to apply because of
482 the need for static functional status and large animal numbers but could be sufficiently valuable
483 to be worthwhile for some commonly applied testing methods.

484

485 **Conclusion**

486 Analysis of sources of variation and application of reference change intervals in individual
487 animals do not replace a conventional assessment of group-level outcomes, but can provide an
488 additional layer of analysis that complements and extends such findings. Therefore there would
489 appear to be considerable merit in reporting both types of analysis alongside each other.
490 Nevertheless, the clinical effectiveness of an intervention can be usefully estimated by
491 determining the proportion of individual outcomes that fall, or are expected to fall, outside
492 specific z-scores and, as we show above, the same analytical technique can aid in redesign of
493 functional tests to maximize their precision. Although we have focused on application of this
494 analytical technique to behavioral data, because these provide the most transparent examples, the
495 same methods can be applied to many other types of quantitative or semi-quantitative
496 neuroscience data including analysis of cell activity or production of specific molecules.

497

498 **Acknowledgements**

499 The International Spinal Research Trust (grant ref STR116) and University of Iowa have
500 provided funds for work with dog and sheep models of spinal cord injury, respectively, to NDJ.
501 LM receives funding from the European Research Council under the European Union's Seventh

502 Framework Programme (FP/2007-2013) / ERC Grant Agreement n. 309731. We thank Dr Hilary
503 Hu and Josh Bratsch-Prince for collecting raw data.

504

505 **Animal care and use**

506 All animal studies were reviewed by the Institutional Animal Care and Use Committee and by
507 the ethical review process at the institution where the work was performed.

508

509 **Conflict of interest**

510 The authors declare no conflict of interest.

511

512 **Author contributions**

513 NJ, UJ, SB designed article approach and analyzed data

514 SS, MH, LM collected and analyzed data

515 NJ, UJ, SB wrote the paper

516

517 **Data Accessibility**

518 The primary data is available as Supplementary Material.

519

520

521

522

523

524

525

526 **References**

527

528 ARRIVE Guidelines: [http://www.equator-network.org/reporting-guidelines/improving-](http://www.equator-network.org/reporting-guidelines/improving-bioscience-research-reporting-the-arrive-guidelines-for-reporting-animal-research/)
529 [bioscience-research-reporting-the-arrive-guidelines-for-reporting-animal-research/](http://www.equator-network.org/reporting-guidelines/improving-bioscience-research-reporting-the-arrive-guidelines-for-reporting-animal-research/)

530

531 Bate, S. T., & Clark, R. A. (2014). *The design and statistical analysis of animal experiments.*

532 Section 6.14: Nested design analysis module. Cambridge University Press. p285.

533

534 Begley, C.G. (2013) Six red flags for suspect work. *Nature* **497**, 433-4.

535

536 Biering-Sørensen, F., Craggs, M., Kennelly, M., Schick, E. & Wyndaele, J.J.

537 (2008) International urodynamic basic spinal cord injury data set. *Spinal Cord*

538 **46**, 513-516.

539

540 Button, K.S., Ioannidis, J.P., Mokrysz, C., Nosek, B.A., Flint, J., Robinson, E.S. & Munafò,

541 M.R. (2013) Power failure: why small sample size undermines the reliability of neuroscience.

542 *Nat. Rev. Neurosci.* **14**, 365-76.

543

544 Carmen, R., Iglesias, N. & Garcia-Lario, J. V., *et al.* (2007) Within-subject biological variation

545 in disease: collated data and clinical consequences. *Ann. Clin. Biochem.* **44**, 343-352.

546

547 Center for 3Rs: <http://3rs.ccac.ca/en/research/reduction/experimental-design.html>

548

549 Flatland, B., Freeman, K.P., Friedrichs, K.R., Vap, L.M., Getzy, K.M., Evans,
550 E.W. & Harr, K.E. (2010) ASVCP quality assurance guidelines: control of
551 general analytical factors in veterinary laboratories. *Vet. Clin. Path.* **39**, 264–277.
552

553 Fraser, C.G. & Harris, E.K. (1989) Generation and application of data on
554 biological variation in clinical chemistry. *Crit. Rev. Clin. Lab. Sci.* **27**, 409-437.
555

556 Fraser, C.G. (2001) Chapter 3. Changes in serial results. In Fraser, C.G. (ed) *Biological*
557 *variation: from principles to practice*. AACC Press, Washington DC, pp. 67-90.
558

559 Friedrichs, K.R., Harr, K.E., Freeman, K.P., Szladovits, B., Walton, R.M., Barnhart, K.F. &
560 Blanco-Chavez, J. American Society for Veterinary Clinical Pathology. (2012) ASVCP
561 reference interval guidelines: determination of de novo reference intervals in veterinary species
562 and other related topics. *Vet. Clin. Pathol.* **41**, 441-53.
563

564 Clinical and Laboratory Standards Institute. *Defining, Establishing, and Verifying Reference*
565 *Intervals in the Clinical Laboratory; Approved Guidelines*. 3rd ed. Wayne, PA: CLSI; 2008.
566

567 Garner, J.P. (2014) The significance of meaning: why do over 90% of behavioral neuroscience
568 results fail to translate to humans, and what can we do to fix it? *I.L.A.R. J.* **55**, 438-56
569

570 Gouveia, K. & Hurst, J.L. (2013) Reducing mouse anxiety during handling:
571 effect of experience with handling tunnels. *PLoS One* 20, e66401.
572

573 Halsey, L.G., Curran-Everett, D., Vowler, S.L. & Drummond, G.B. (2015) The fickle P value
574 generates irreproducible results. *Nat. Methods* 12, 179-85.
575

576 Harris, E.K. & Yasaka, T. (1983) On the calculation of a "reference change" for
577 comparing two consecutive measurements. *Clin. Chem.* 29, 25–30.
578

579 Kilkenny C, Browne WJ, Cuthi I, Emerson M, Altman DG. (2012) Improving bioscience
580 research reporting: the ARRIVE guidelines for reporting animal research. *Vet. Clin. Pathol.* 41:
581 27-31.
582

583 Landis, S.C., Amara, S.G., Asadullah, K., Austin, C.P., Blumenstein, R., Bradley, E.W., Crystal,
584 R.G., Darnell, R.B., Ferrante, R.J., Fillit, H., Finkelstein, R., Fisher, M., Gendelman,
585 H.E., Golub, R.M., Goudreau, J.L., Gross, R.A., Gubitza, A.K., Hesterlee, S.E., Howells,
586 D.W., Huguenard, J., Kelner, K., Koroshetz, W., Krainc, D., Lazic, S.E., Levine, M.S., Macleod,
587 M.R., McCall, J.M., Moxley, R.T. 3rd., Narasimhan, K., Noble, L.J., Perrin, S., Porter,
588 J.D., Steward, O., Unger, E., Utz, U. & Silberberg, S.D. (2012) A call for transparent reporting to
589 optimize the predictive value of preclinical research. *Nature* 490, 187-91.
590

591 Marchenko, Y. (2006) Estimating variance components in Stata. *The Stata*
592 *Journal* 6, 1-21.
593

594 Nunes, L.A., Brenzikofer, R. & de Macedo, D.V. (2010) Reference change values of blood
595 analytes from physically active subjects. *Eur. J. Appl. Physiol.* **110**, 191-8.
596

597 Rotterdam, E.P., Katan, M.B., & Knuiman, J.T. (1987) Importance of time interval
598 between repeated measurements of total or high-density lipoprotein cholesterol
599 when estimating an individual's baseline concentrations. *Clin. Chem.* **33**, 1913-
600 1915.
601

602 Rousselet, G.A., Foxe, J.J. & Bolam, J.P. (2016) A few simple steps to improve description of
603 group results in neuroscience. *Eur.J. Neurosci.* **44**, 2647-2651.
604

605 Safayi, S., Jeffery, N.D., Shivapour, S.K., Zamanighomi, M., Zylstra, T.J., Bratsch-Prince, J.,
606 Wilson, S., Reddy, C.G., Fredericks, D.C., Gillies, G.T., Howard, M.A. 3rd. (2015) Kinematic
607 analysis of the gait of adult sheep during treadmill locomotion: Parameter values, allowable total
608 error, and potential for use in evaluating spinal cord injury. *J. Neurol. Sci.* **358**, 107-12.
609

610 Shaver, J.P. (2013) What Statistical Significance Testing Is, and What It Is Not. *J. Exp. Educ.* **61**,
611 293-316.
612

613 Snapinn, S.M. & Jiang, Q. (2007) Responder analyses and the assessment of a clinically relevant
614 treatment effect. *Trials* **8**, 31.

615

616 Sorge, R.E., Martin, L.J., Isbester, K.A., Sotocinal, S.G., Rosen, S., Tuttle, A.H., Wieskopf,
617 J.S., Acland, E.L., Dokova, A., Kadoura, B., Leger, P., Mapplebeck, J.C., McPhail, M., Delaney,
618 A., Wigerblad, G., Schumann, A.P., Quinn, T., Frasnelli, J., Svensson, C.I., Sternberg, W.F.
619 & Mogil, J.S. (2014) Olfactory exposure to males, including men, causes stress and related
620 analgesia in rodents. *Nat. Methods* **11**, 629-632.

621

622 Steward, O. (2016) A Rhumba of "R's": Replication, Reproducibility, Rigor, Robustness: What
623 Does a Failure to Replicate Mean? *eNeuro* **7**, 3.

624

625 The Academy of Medical Sciences Reproducibility and reliability of biomedical research:
626 improving research practice: Symposium report, October 2015.
627 <https://acmedsci.ac.uk/viewFile/56314e40aac61.pdf>

628

629 Theodorsson, E., Magnusson, B. & Leito, I. (2014) Bias in clinical chemistry
630 *Bioanalysis* **6**, 2855-2875.

631

632 Walton, R.M. (2012) Subject-based reference values: biological variation,
633 individuality, and reference change values. *Vet. Clin. Pathol.* **41**, 175-181.

634 **Figure legends**

635

636 **Figure 1:** Illustration of the relationship between ‘experimental noise’ and group-level change.
637 In this theoretical example, we assume that previous analysis has revealed that the combination
638 of experimental inaccuracy and physiological variability can be responsible for changes of up to
639 20% between sequential test results (methods to derive these values are explained in the text).
640 Collectively these sources of variability are termed the reference change interval, which defines
641 limits of experimental noise that can be expected to arise through chance; methods to derive this
642 interval are shown in the text.

643 **a:** The pre-intervention test result is illustrated for each of eight experimental subjects, along
644 with bars depicting the range of values encompassed by the reference change interval of 20%
645 from baseline. **b:** Comparison of pre- and post-intervention results and group-level statistical
646 comparison by paired Student’s t-test, indicating a statistically significant difference. **c:** Post-
647 intervention results for each subject. **d:** Post-intervention test results for each subject, shown
648 alongside the pre-test result and the expected analytical variability for each test result (‘RCI’,
649 bars). Note that the post-intervention result for each subject is within the reference change
650 interval for that subject.

651

652 **Figure 2:** Hindlimb stride duration in a group of eight trained sheep walking on a treadmill
653 before and after administration of drug X. Wilcoxon matched pairs signed rank test suggests a
654 significant effect of drug X on stride duration ($P=0.023$).

655

656 **Figure 3:** Baseline and 1-month post-intervention measures of bladder compliance index for
657 dogs that had received intraspinal injection of chondroitinase ABC at time 0. Group-level
658 analysis by paired Student's t-test reveals a non-significant difference between time points.
659

1 **Supplementary Material**

2

3 The data used in Examples 1 and 2 are supplied as Excel files ('sheep treadmill data' and 'dog
4 compliance data').

5 The nested design analysis module in the free-to-download statistical software package

6 InVivoStat (<http://invivostat.co.uk/>) can be used to derive the analytical, intra- and inter-animal
7 standard deviations for both sets of data.

8 Guidance for use of this module is available at <http://invivostat.co.uk/wp-content/user->

9 [guides/Nested-Design-Analysis.pdf](#) and also in Section 6.14 of the book '*The design and*

10 *statistical analysis of animal experiments*' by Simon T. Bate and Robin A. Clark (Cambridge
11 University Press).

12

13 **Source of data**

14 ***Example 1 – sheep treadmill data***

15 Sheep were trained to walk on a treadmill moving at constant speed until they were able to
16 maintain a constant walking gait without stopping or trotting (treadmill speed was ~4km/h). To
17 reach a plateau of steady performance took daily training sessions for a period of about 10 days
18 (Safayi *et al.*, 2015).

19 Retroflective markers were attached to the hooves and the skin overlying boney prominences on
20 the hindlimbs and used with a motion capture system (Vicon) to capture at least 4 minutes of
21 continuous walking gait on four individual days, each separated by a week. Each segment of
22 motion capture data was analyzed using a custom-written Matlab (MathWorks Inc) code to
23 extract the duration of each hindlimb step (determined by the timing to minimum hoof height).

24 Each 4 minute segment of motion capture data was divided into four individual segments of 1
25 minute each and these we termed ‘trials’. Analysis of each segment was used to correspond to
26 analysis of aliquots of a fluid sample in a clinical pathology laboratory.

27 The ‘post-drug’ data was simulated (to show specific aspects of the relationship between group-
28 level and individual-level analyses). It was derived by multiplying the ‘pre-drug’ data by a
29 restricted range of random numbers (90-130%) that would, on average, increase the ‘post-drug’
30 values by a small proportion but allow for individual variation.

31 Reference change value calculations were made on the pre-drug dataset and partitioned into
32 analytical, intra-animal and inter-animal variability as described in the main text. The calculated
33 reference change interval was then compared with ‘post-drug’ values. The Supplementary
34 Dataset ‘Sheep treadmill data’ contains the pre-drug values. The ‘post-drug’ values are available
35 in Table 2 in the main article.

36

37 ***Example 2 – dog bladder compliance data***

38 Dogs in this study were domestic pets that had incurred severe acute spinal cord injury at least 12
39 weeks previously and had not regained voluntary motor function or continence. They were
40 entered into a blinded randomized controlled trial on the efficacy of intraspinal injection of
41 chondroitinase ABC in reversing functional deficits associated with chronic spinal cord injury
42 (Hu *et al.*, in press). The bladder compliance was measured in all dogs at baseline (entry into the
43 study) and then again at 1, 3 and 6 months after intervention (or sham). Dogs were randomly
44 allocated to active intervention (chondroitinase) or sham at time 0.

45 Briefly, bladder compliance was measured by passing a dual lumen catheter into the bladder,
46 draining the urine and then infusing sterile saline into the bladder at a constant rate. Infusion was

47 terminated when saline leaked from the urethra, or when an intravesicular pressure of greater
48 than 40cmH₂O was attained. Compliance is measured as the change in volume (from baseline)
49 divided by the change in pressure (from baseline). The data distribution was not normally
50 distributed and so it was log-transformed for analysis in this article; therefore it is referred to as a
51 ‘compliance index’.

52 To calculate the reference change interval we used the dogs that were subsequently identified as
53 controls (in which the allocation was blinded for observers by making needle sticks in the dorsal
54 skin). We had repeated data from each time point that allowed calculation of intra-animal and
55 inter-animal variation in this outcome. Analytical variability could not be specifically partitioned
56 because repeat testing of compliance would likely lead to changes in the outcome measure.

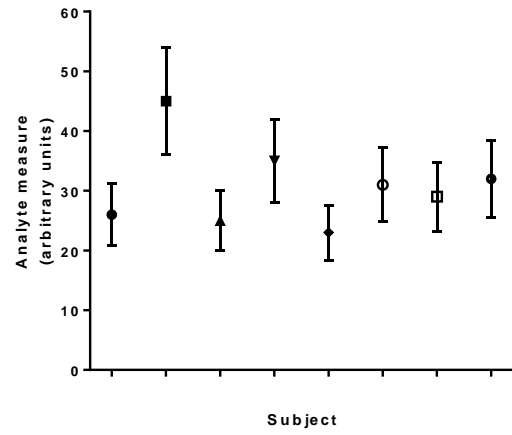
57 As before, the reference change interval calculated in the control dogs was then applied to the 1-
58 month data from the active intervention dogs to determine how much ‘real change’ had occurred.
59 The control dog data used to calculate the reference change value is contained in the
60 Supplementary Dataset ‘Dog compliance data’ and the data from dogs in the active intervention
61 group are shown in Table 3 in the main article.

62

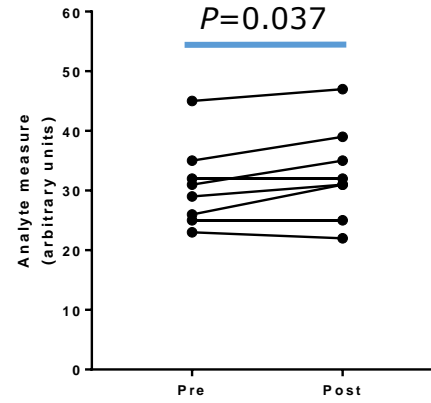
63

Figure 1

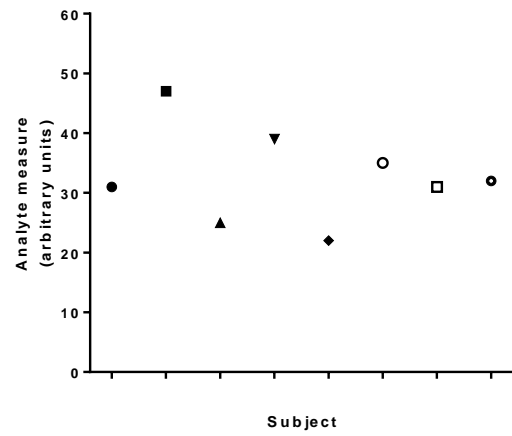
a



b



c



d

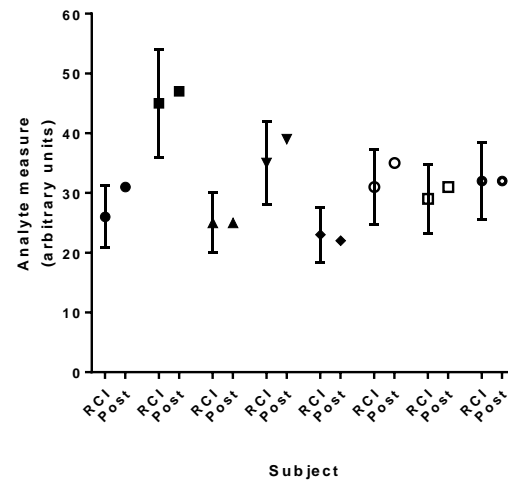


Figure 2

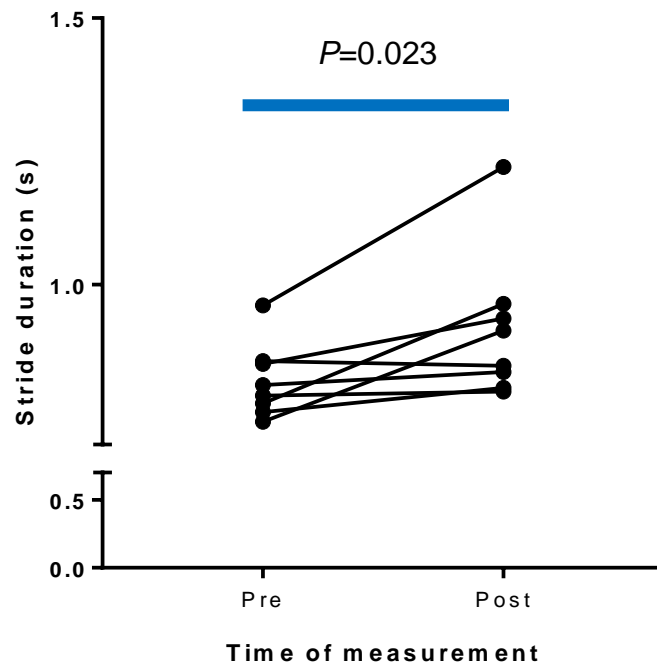


Figure 3

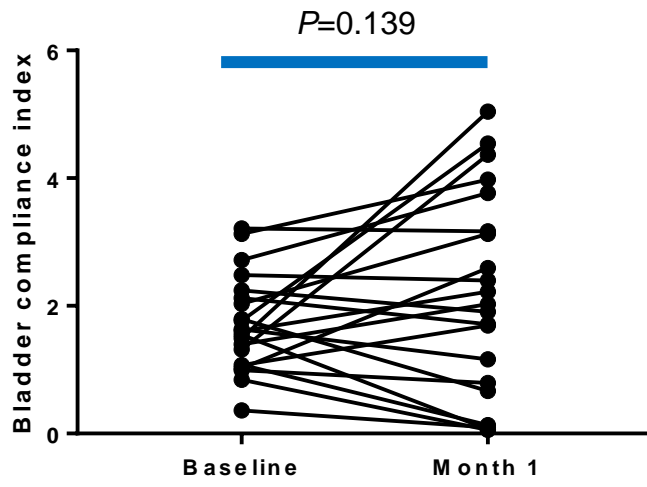


Table 1: Baseline and repeat testing results

Baseline	Repeat	Repeat / Baseline %
26	31	119
45	47	104
25	25	100
35	39	111
23	22	96
31	35	113
29	31	107
32	32	100

In this example we know that previous experiments have shown that the combination of various sources of 'experimental noise' (see text) can make retest results vary by as much as 20% from baseline values in normal, untreated animals.

In this example we use paired Student's t-test to compare baseline values with repeat test results after an experimental intervention, producing $t = 2.204$, $P = 0.037$.

Therefore population-level analysis alone suggests significant increase between baseline and repeat testing. However, inspection of the data above reveals that for none of the subjects does the retest value exceed what might be attributable to experimental noise alone.

This type of paradoxical outcome is most likely to occur in small sample size experiments as are frequently used in neuroscience.

Table 1: Calculation of the upper limit of the reference change interval for hindlimb stride duration and comparison with measured effect of drug administration

Condition	Sheep							
	1	2	3	4	5	6	7	8
Pre-drug	0.743	0.812	0.792	0.778	0.761	0.852	0.856	0.961
Upper RCI boundary	0.875	0.957	0.933	0.916	0.896	1.004	1.008	1.132
Post-drug	0.914	0.836	0.800	0.964	0.807	0.937	0.848	1.221
Post>upper RCI boundary?	Yes	No	No	Yes	No	No	No	Yes

For each sheep (columns) we show the duration of hindlimb stride before ('pre-drug') and after ('post-drug') administration of drug X. After calculating the reference change value .we have used this, in conjunction with the Pre-drug measurement to define the upper boundary of the reference change interval, see equation (2). In the bottom row we ask whether the Post-drug measurement exceeds the upper boundary of the reference change interval.

Table 3: Analysis of change in bladder compliance index at 1 month after injection of intraspinal chondroitinase ABC

Baseline	Month 1	Baseline + RCV	Baseline - RCV	Compliance increased?	Compliance decreased?	Z-score
1.609	2.223	4.162	-0.944	No	No	0.472
0.993	0.792	2.569	-0.583	No	No	-0.250
1.396	2.028	3.611	-0.819	No	No	0.560
2.485	2.398	6.429	-1.459	No	No	-0.043
3.214	3.169	8.313	-1.886	No	No	-0.017
2.121	1.715	5.486	-1.245	No	No	-0.236
1.492	5.044	3.859	-0.876	Yes	No	2.944
1.768	4.544	4.574	-1.038	No	No	1.940
2.716	3.767	7.026	-1.594	No	No	0.478
1.792	0.670	4.636	-1.052	No	No	-0.774
0.843	0.058	2.181	-0.495	No	No	-1.151
1.022	2.590	2.644	-0.600	No	No	1.896
1.070	1.692	2.768	-0.628	No	No	0.719
3.129	3.977	8.095	-1.837	No	No	0.335
0.365	0.095	0.944	-0.214	No	No	-0.914
2.031	3.129	5.254	-1.192	No	No	0.668
1.546	0.067	4.000	-0.908	No	No	-1.183
1.079	0.134	2.791	-0.633	No	No	-1.083
2.244	1.907	5.805	-1.317	No	No	-0.186
1.633	1.161	4.225	-0.959	No	No	-0.357
1.316	4.369	3.404	-0.772	Yes	No	2.868