

1
2
3 **BRIEFING**

4 **<1032> Design and Development of Biological Assays.** This General Chapter
5 *Design and Development of Bioassays* <1032> is one of an integrated group of new
6 General Chapters that provide guidance across several complementary bioassay topics.
7 The others include *Biological Assay Validation* <1033>; *Analysis of Biological Assay*
8 <1034>; and an as-yet unnumbered “roadmap” chapter that will include a glossary
9 applicable to General Chapters <1032>–<1034>. The group arose from one source,
10 General Chapter *Design and Analysis of Biological Assays* <111>, that is now official in
11 the *US Pharmacopeia (USP)*. The partitioning of chapters to different aspects of
12 biological assays allows both greater focus and clarity and the opportunity to expand on
13 issues. As the group of General Chapters evolves, General Chapter <111> will remain
14 in modified form. Thus the entire group when completed will consist of five *USP* General
15 Chapters.

16 The Bioassay Development ad hoc Advisory Panel of the Statistics Expert Committee
17 encourages input from all interested parties regarding the proposed <1032>. USP’s
18 intent is to reflect the best contemporary practices regarding the design and
19 development of bioassays. This will be achieved when members of the bioassay
20 community take advantage of this opportunity to engage in <1032>’s development by
21 responding to this In-Process Revision. Comments regarding <1032> should be sent to
22 Tina S. Morris, PhD (tsm@usp.org).

23
24 **<1032> DESIGN AND DEVELOPMENT OF BIOLOGICAL ASSAYS**

25
26 **1. INTRODUCTION**

27
28 **1.1 Purpose**

29
30 General Chapter *Design and Development of Biological Assays* <1032> presents
31 methodology for the development of bioassay procedures that have sound experimental
32 design, that provide data that can be analyzed using well-founded statistical principles,
33 and that are fit for their specific use.

34
35 General Chapter <1032> is one of a group of five planned General Chapters that focus
36 on relative potency assays, in which the activity of a Test material is quantified by
37 comparison to the activity of a Standard material. However, many of the principles can
38 be applied to other assay systems.

39
40 This General Chapter is intended to guide the design and development of a bioassay for
41 a drug product intended for commercial distribution. Although adoption of this chapter’s
42 recommended methods may be resource intensive during initial drug development,
43 early implementation may yield benefits.

44
45 **1.2 Audience**

47 This chapter is intended for both the practicing bioassay analyst and for the statistician
48 who is engaged in bioassay efforts. The former will find guidance for implementing
49 bioassay structure and methodology to achieve analytical goals while reliably
50 demonstrating the biological activity of interest, and the latter will gain insights regarding
51 the constraints of biology that can prove challenging to balance with the most
52 straightforward practice of statistics.

53
54 ***Focus on relative potency.*** Common applications of biological assays include
55 characterization of a biological or biotechnology product, stability testing, and lot
56 release. Potency assays help provide assurance of the quality and consistency of the
57 product. Because of the inherent variability in test systems (from animals, cells,
58 instruments, and reagents, and day-to-day and between-lab variation), an absolute
59 measure of potency (e.g., a dose of erythropoietin that increases hemoglobin content by
60 1 g) may not be available. This has led to the adoption of the relative potency
61 methodology. Assuming that the Standard and Test materials contain samples with
62 biologically similar (if not identical) activity, the Test sample can be expected to behave
63 like a concentration or dilution of the Standard, and *parallelism* (the quality of Test and
64 Standard concentration–response segments being parallel) should be present. For
65 samples that are biologically similar, the horizontal displacement between the curves is
66 interpreted as log relative potency. Relative potency is determined by comparison of
67 Test to Standard assay response, which means that the assay does not need to
68 achieve a specific observable response. The potency of the Standard is usually
69 assigned a value of 1 (or 100%). The Standard could be a material established as such
70 by a national (e.g., USP) or international (e.g., WHO) organization, or it could be an
71 internal Standard.

72 73 **2. BIOASSAY FITNESS FOR USE**

74
75 To evaluate whether an assay is fit for use, analysts must clearly specify the purposes
76 for performing the bioassay. Common uses for a bioassay include lot release of drug
77 substance (active pharmaceutical ingredient) and drug product; stability; qualification of
78 Standard sample and other critical reagents; characterization of process intermediates
79 and formulations, contaminants, and degradation products; and support of changes in
80 the product production process. The accuracy, specificity, precision, and robustness
81 requirements may be different for each of these potential uses. It is a good strategy to
82 develop and validate a bioassay to support multiple intended uses. Decisions about
83 fitness for use are based on scientific and statistical considerations, as well as practical
84 considerations such as cost, turnaround time, and throughput requirements for the
85 assay.

86
87 When assays are used for lot release, a linear-model bioassay may allow sufficient
88 assessment of similarity. For bioassays used to support stability or comparability, to
89 evaluate changes in production processes, or to qualify reference materials, critical
90 reagents, or changes in the assay process, it may be useful to assess similarity using
91 the asymptote of maximum response.

92

2.1 Process Development

Bioassays are often required in the development and optimization of product manufacturing, including formulation and scale-up processes. Bioassays can be used to evaluate purification strategies, optimize product yield, and measure product stability. Because samples taken throughout the process are often analyzed and compared, sample matrix effects should be carefully studied to determine an assay's fitness for use. For relative potency measures, the Standard material may require dilution into suitable matrices for quantitation. The bioassay's precision should be suitable for measuring process performance or for assessing and comparing the stability of candidate formulations.

2.2 Process Characterization

In order to assess an effect on drug potency specific to different stages of drug manufacture or following changes in the manufacturing process (i.e., to demonstrate product equivalence before and after process changes are made), analysts may perform bioassays distinct from routine product lot-release assays. These assays can provide information about product characteristics in the presence of process changes. Bioassays used in this type of application may be qualitative or quantitative. The assay's fitness for use depends on its sensitivity to unexpected changes in the product profile.

2.3 Product Release

Bioassays are used to evaluate the potency of the drug before product release and may predict clinical efficacy. To this end, the assay should reflect or mimic the product's known or intended mechanism of action. If the bioassay does not include the functional biology directly associated with the mechanism of action, it may be necessary to demonstrate a relationship between this bioassay's estimated potency results and those of some other assay that better or otherwise reflects biomimetic activity.

For product-release testing, a release specification is established to define the range of potency values that are acceptable for product stored under labeled conditions until the expiration date. The relative accuracy and the precision of the reportable value from the bioassay must support the number of significant digits listed in the specification (see *Biological Assay Validation* <1033>), as well as the specification range. These specifications should include the loss of potency (instability) that will be tolerated. In order to meet these specifications, manufacturing quality control will have sufficiently narrow release assay specifications in order to accommodate any loss of activity due to instability and uncertainty in the release assay. The expiration period may be calculated as the point at which the 95% confidence bound of the predicted value crosses the product specification limit. The more variable the assay and the steeper the slope of decay, the nearer in time the assigned expiration date will be.

2.4 Process Intermediates

139
140 Bioassay assessment of process intermediates can provide information regarding
141 interference and specificity. Downstream process decisions pertaining to formulation
142 and fill strategies may rely on bioassays in order to ensure that product will meet its
143 established specifications. For example, unformulated bulk materials may be held and
144 evaluated for potency. Bulks may be pooled with other bulk lots, diluted, or reworked
145 based on the potency results. For these types of applications, the bioassay must be
146 capable of measuring product activity in different matrices. In some cases, a separate
147 Standard material is made and is used to calculate relative potency for the process
148 intermediate. Thus fitness for use must demonstrate that decisions based on the
149 evaluation of process intermediates do in fact yield the desired outcomes for
150 downstream product, including final container.

151 152 **2.5 Stability**

153
154 Potency is a stability-indicating parameter of biotechnology and vaccine products.
155 Stability studies are performed during development to establish shelf life as well as to
156 identify and estimate the degradation products and rates of the product. Post licensure
157 stability studies are used to monitor and maintain product stability. Stability studies are
158 designed by selecting the number of observations and the testing intervals, as well as
159 the potency assay format that is predicted to achieve a suitable level of uncertainty in
160 the stability parameter of interest. Knowledge of both short-term and long-term
161 variability of the bioassay is important for a product's stability evaluation. In addition, a
162 trend in bias across levels in the bioassay can influence the measurement of potency
163 over time. A limit on the trend in bias can help ensure accurate assessment of product
164 stability.

165 166 **2.6 Qualification of Reagents**

167
168 The quantitative characterization of a new Standard material requires a highly precise
169 measurement of the new material's biological activity. This value is used either to
170 establish that the new lot is equivalent to the previous lot or to assign it a label potency
171 to which Test samples can be compared. Replication beyond routine testing may be
172 required to achieve a greater confidence in measuring the potency of the Standard
173 material. In some cases, such as maintenance of the cell seed train, the bioassay may
174 be used to qualify a reagent such as fetal bovine serum. The fitness for use in such
175 cases is tied to the ability of the assay to screen reagent lots and to ensure that lots that
176 may degrade the relative potency measurements are not accepted.

177 178 **2.7 Product Integrity**

179
180 Biotechnology, biological, and vaccine products may contain a population of
181 heterogeneous material, including the intended predominant product material.
182 Additionally, the product is subject to variation in manufacture, as well as degradation
183 because of stress and storage conditions. Some process impurities and degradation
184 products may be active, partially active, or inactive in the bioassay. Biological assays

185 can be used to differentiate intended product from product variants or derivatives in
186 which changes in structure or relative composition may be manifested in subtle yet
187 characteristic changes in the bioassay response (e.g., change in slope or asymptote).
188 Studies that identify characteristic changes associated with variants of the intended
189 product help ensure consistent product performance. Whenever possible, the bioassay
190 should be accompanied by orthogonal methods that are capable of separating and
191 detecting product variants, process impurities, and degradation products.

192 193 **3. BIOASSAY FUNDAMENTALS**

194 195 **3.1 In Vivo Bioassays**

196
197 In vivo potency assays are bioassays in which a set of dilutions of each of the Standard
198 and Test materials is administered to animals and the dose–response relationships are
199 used to estimate potency. For some animal assays, the end point is simple (e.g., rat
200 body weight gain assay for human growth hormone or rat ovarian weight assay for
201 follicle stimulating hormone), but others require further processing of samples collected
202 from treated animals (e.g., reticulocyte count for erythropoietin, steroidogenesis for
203 gonadotropins, neutrophil count for granulocyte colony stimulating factor, or antibody
204 titer after administration of vaccines). With the advent of cell lines specific for the
205 putative physiological mechanism of action (MOA), the use of animals for the
206 measurement of potency has greatly diminished. Cost, low throughput, ethics, and other
207 practical issues argue against the use of animal bioassays. Regulatory agencies have
208 encouraged the responsible limitation of animal use whenever possible (see The
209 Interagency Coordinating Committee on the Validation of Alternative Methods, Mission,
210 Vision, and Strategic Priorities; February 2004). When in vitro activity is not strongly
211 associated with in vivo activity (e.g., EPO), the combination of in vitro cell-based assay
212 and a suitable physicochemical method (e.g., IEF, glycan analysis) may substitute for in
213 vivo assays. However, a need for in vivo assays may remain when in vitro assays
214 cannot detect differences that may affect in vivo activity.

215
216 Animals' physiological responses to biological drugs (including vaccines) may predict
217 patients' responses. Selection of animal test subjects by species, strain, gender, and
218 maturity or weight range seeks to provide a representative model with which to assess
219 the activity of Test materials. In order to minimize assay variability, bioassays involve
220 comparison of a Test to a Standard material of defined potency, together with the
221 application of statistical tools, specific lab techniques, and rigorous adherence to
222 protocol.

223
224 Some methods lend themselves to the use of colony versus naïve animals. For
225 example, pyrogen and insulin testing benefit from using experienced colony rabbits that
226 provide a reliable response capacity. If animals recently introduced to the colony fail to
227 respond normally to the Test compound after several doses, they should be culled from
228 the colony so they do not cause repeated future invalid or indeterminate assay results.
229 In the case of assaying highly antigenic Test compounds for pyrogens, however, naïve
230 animals should be used to avoid generating inaccurate or confounded results from

231 given samples. Other colony advantages include common controlled environmental
232 conditions (macro/room, and micro/rack), consistent feeding schedule, provision of
233 water, and husbandry routine.

234
235 When a historical database of colony records and assay data is available, valuable
236 method optimization can be obtained. The influence of biasing factors can be reduced
237 by applying randomization principles such as distribution of weight ranges across dose
238 groups, group assignments from shipping containers to different cages, or use of
239 computer-generated or deck patterns for injection/dosing. A test animal must be healthy
240 to be suitable for use in a bioassay and must have time to stabilize in its environment.
241 Factors that combine to influence an animal's state of health include proper nutrition,
242 hydration, freedom from physical and psychological stressors, adequate sanitization of
243 housing components, controlled light cycle (diurnal/nocturnal), experienced handling,
244 skillful injections and bleedings, and absence of noise or vibration. Daily observation of
245 test animals is essential, and veterinary care must be available to minimize stress to the
246 animals and to evaluate potential issues that might challenge the validity of bioassay
247 results.

248 249 **3.2 Ex Vivo Bioassays**

250
251 Cells or tissues from human or animal donors can be cultured in the laboratory and
252 used to assess the activity of a Test article. In the case of cytokines, the majority of
253 assays use cells from the hematopoietic system and subsets of hematopoietic cells
254 from peripheral blood such as peripheral blood mononuclear cells or peripheral blood
255 lymphocytes. For proteins that act on solid tissues, such as growth factors and
256 hormones, specific tissue on which they act can be removed from animals, dissociated,
257 and cultured for a limited period either as adherent or semi-adherent cells. Although an
258 ex vivo assay system has the advantage of similarity to the natural milieu, it may also
259 suffer from substantial donor-to-donor variability, as well as challenging availability of
260 appropriate cells.

261
262 Bioassays that involve live tissues or cells from an animal (e.g., rat hepatocyte glucagon
263 method) require process management similar to that of in vivo assays to minimize
264 assay variability and bias. The level of effort to manage bias (e.g., randomization)
265 should be appropriate for the purpose of the assay. Additional factors that may affect
266 assay results include time of day, weight or maturity of animal, anesthetic used, buffer
267 components/reagents, incubation bath temperature and position, and cell viability.

268 269 **3.3 In Vitro (Cell-based) Bioassays**

270
271 Cell-based bioassays using clonal cell lines that respond to specific ligands or infectious
272 agents can be used for lot-release assays of therapeutic proteins and vaccines. These
273 cell lines can be derived from tumors, immortalized as factor-dependent cell lines, or
274 engineered cell lines transfected with appropriate receptors. Such cells can be banked
275 and used with an expectation of uniform response through some number of continuous
276 passages. The selection of appropriate cell lines with desirable magnitude of response

277 and stability may be challenging. Advances in recombinant DNA technology and the
278 understanding of cellular signaling mechanisms have allowed the generation of
279 engineered cell lines with improved response and longer stability. The cellular
280 responses to the protein of interest depend on the drug's MOA and the duration of
281 exposure. Such responses include cell proliferation, cell killing, antiviral activity,
282 differentiation, cytokine/mediator secretion, and enzyme activation. Assays involving
283 these responses may require incubation of the cells over several days, during which
284 time contamination, uneven evaporation, or edge effects may arise. Comparatively rapid
285 responses based on an intracellular signaling mechanism—such as second
286 messengers, protein kinase activation, or reporter gene expression—have proven
287 acceptable to regulatory authorities. Lastly, most cell lines used for bioassays express
288 receptors for multiple cytokines and growth factors. This lack of specificity may not be
289 detrimental if other assays demonstrate the Test material's specificity.

290
291 Cell-based bioassay design should reflect knowledge of the factors that influence
292 response of cells to the Test agent. Response variability is often reflected in parameters
293 such as slope, EC_{50} of the dose–response curve, or the ratio of maximum to minimum
294 response. Even though the adoption of relative potency methodology alleviates the
295 effect of these parameters on potency estimates, such response variability could result
296 in significant effects on system suitability, precision, and accuracy. Therefore sources of
297 variability relevant to a given bioassay method should be identified, and control
298 strategies to minimize their effects should be implemented.

299
300 The development of a cell-based bioassay begins with the selection or generation of a
301 cell substrate. The cell substrate is usually the most critical reagent in a cell-based
302 assay. To ensure an adequate and consistent supply of cells for product testing, the cell
303 line(s) used in the assay should be banked. Several considerations arise in developing
304 the assay cell banks.

305
306 The source of the cell line, whether generated by the product manufacturer or acquired
307 from a collaborator, academic institution, or culture collection, must have documentation
308 that details the cell line's history from origin to banking and supports its application for
309 commercial use. The origin, generation, and propagation of the cell line through to the
310 point at which the cell line is selected for use must be documented and described in
311 detail sufficient to permit recreation of a similar cell line if necessary. Information
312 pertaining to the cell line should be documented during assay development and before
313 banking. This information may include but is not limited to identity (e.g., isoenzyme,
314 phenotypic markers, genetic analysis); morphology (e.g., archived photographic
315 images); purity (e.g., mycoplasma, bacteria, fungus and virus testing); cryopreservation;
316 thaw and culture conditions (e.g., media components, thaw temperature and method,
317 methods of propagation, seeding densities, harvest conditions); thaw viability
318 (immediately after being frozen and after time in storage); growth characteristics (e.g.,
319 cell doubling times); and functional stability (e.g., ploidy).

320
321 Extensive characterization is necessary to ensure the quality and longevity of cell banks
322 for use in the QC environment. The general health and metabolic state of the cells at

323 the time of bioassay can greatly influence the test results. After a cell line has been
324 characterized and is ready for banking, analysts typically prepare a two-tiered bank
325 (Master and Working). A Master Cell Bank is created as the source layer for the
326 Working Cell Bank. The Working Cell Bank is derived by expansion of one or more vials
327 of the Master Cell Bank. The size of the banks depends on the growth characteristics of
328 the cells, the number of cells required for each assay, and how often the assay will be
329 performed. Some cells may be sensitive to cryopreservation, thawing, and culture
330 conditions, and the banks must be carefully prepared and characterized before being
331 used for validation studies and for regular use in the QC laboratory.

332
333 The following factors that may affect bioassay performance are common to many cell-
334 based bioassays. This list is not exhaustive, and analysts with comprehensive
335 understanding and experience with the cell line should be involved during assay
336 development. These experienced individuals should identify factors that might influence
337 assay outcomes and establish strategies for an appropriate level of control whenever
338 possible. Relevant factors include cell type (adherent or nonadherent); cell thawing;
339 plating density (at thaw and during seed train maintenance) and confluence (adherent
340 cells); culture vessels; growth, staging, and assay media; serum requirements (source,
341 heat inactivation, gamma irradiation); incubation conditions (temperature, CO₂,
342 humidity, culture times from thaw); cell harvesting reagents and techniques (for
343 adherent cells, method of dissociation); cell sorting; cell counting; determination of cell
344 health (growth rate, viability, yield); cell passage number and passaging schedule; cell
345 line stability (genetic, receptor, marker, gene expression level); and starvation or
346 stimulation steps.

347 348 **3.4 Standard** 349

350 The Standard is a critical reagent in bioassays because of its quality as a reliable
351 material to which a Test preparation can be quantitatively compared in a relative
352 potency assay. The Standard may be assigned a unitage or specific activity that
353 represents fully (100%) potent material. Where possible, a Standard should be prepared
354 using the same manufacturing process as the drug substance. If this is not possible
355 (e.g., lyophilization is the preferred Standard storage), verify the capacity of the assay
356 for providing accurate relative potency assessment of samples. Testing performed to
357 qualify a Standard may be more rigorous than the routine testing used for lot release.
358

359 A Standard must be stored under conditions that preserve its full potency during the
360 intended course of its use. To this end, the Standard may be stored under conditions
361 that are different from the normal storage of the drug substance or drug product. These
362 could include a different temperature (e.g., -70 °C or -20 °C instead of 2–8 °C), a
363 different container (e.g., plastic vials instead of syringes), a different formulation (e.g.,
364 lyophilizable formulation or addition of carrier proteins such as human serum albumin,
365 stabilizers, etc.). The Standard material should be tested for stability at appropriate
366 intervals. System suitability criteria of the bioassay such as maximum or background
367 response, EC₅₀, slope, or potency of assay control may be used to detect change in the

368 activity of the Standard. Accelerated stability studies can be performed to predict
369 degradation rates and establish meaningful markers of Standard instability.

370
371 At later stages in clinical development, the Standard may be prepared using the
372 manufacturing process employed in pivotal clinical trials. If the Standard formulation is
373 different than that used in the drug product process, demonstration is required of
374 comparable (similar) assay performance. An initial Standard may be referred to as the
375 *primary Standard*. Subsequent Standards can be prepared using current manufacturing
376 processes and can be designated *working Standards*. Separate SOPs are usually
377 required for establishing these standards for each product. Bias in potency
378 measurements sometimes can arise when the activity of the Standard gradually
379 changes. Nonequivalence in assay response may be observed if the Standard and the
380 Test sample are not biologically similar (e.g., if there are differences in glycosylation). In
381 such cases, careful attention should be paid to changes in critical reagents and/or
382 media components that differentially affect the activities of the Standard and Test
383 sample. It is prudent to archive aliquots of each Standard lot for assessment of
384 comparability with later Standards and for the investigation of assay drift. In this context,
385 trend charts may be useful in identifying the cause of assay drift.

386 **4. STATISTICAL ASPECTS OF BIOASSAY FUNDAMENTALS**

387 **4.1 Data**

388
389
390
391 The statistical elements of bioassay development include consideration of data type and
392 the bioassay model, along with statistical criteria for assessing and ensuring the quality
393 of bioassay results. These form the framework for the bioassay system that will be used
394 to estimate the potency of a Test article.

395
396 There are fundamentally two bioassay data types: quantitative and quantal
397 (categorical). Quantitative data can be either continuous (e.g., collected from an
398 instrument), counts (e.g., plaque-forming units), or discrete (e.g., endpoint dilution
399 titers). Quantal data are typically dichotomous, e.g., survival in an animal model that
400 uses challenge with a pathogen to measure the protection afforded by a Test article or
401 positivity in a plate-based infectivity assay that results in destruction of a cell monolayer
402 following administration of an infectious agent. Quantitative data can be transformed to
403 quantal data by selecting a threshold that statistically distinguishes a positive response
404 from a negative response. Such a threshold can be calculated from data acquired from
405 a negative control, e.g., adding (or subtracting) a measure of uncertainty, such as two
406 times the standard deviation of negative control responses, to the negative control
407 average. Analysts should be cautious about transforming quantitative data to quantal
408 data because this results in a loss of information that may affect bioassay
409 measurement.

410 **4.2 Assumptions**

411
412

413 The key assumption for the analysis of most bioassays is that the Standard and Test
414 samples contain the same effective analyte or population of analytes and thus may be
415 expected to behave similarly in the bioassay. This is termed *similarity*. As will be shown
416 in more detail in Chapter <1034> for specific statistical models, biological similarity
417 implies that statistical similarity is present (the Standard and Test curves are parallel, for
418 parallel-line and parallel-curve models; the Standard and Test curves have a common
419 intercept for slope-ratio models). The reverse is not true. Statistical similarity (parallel
420 lines, parallel curves, or common intercept, as appropriate) implies, but does not
421 ensure, biological similarity (same effective analyte). However, failure to satisfy
422 statistical similarity may be taken as evidence against biological similarity. The
423 existence of a Standard–Test sample pair that passes the assessment of statistical
424 similarity is thus a necessary but not sufficient test that the key assumption of biological
425 similarity is satisfied. Biological similarity thus remains, unavoidably, an assumption.
426 Departures from statistical similarity that are consistent in value across many assays
427 could be indicative of matrix effects or of real differences between Test and Standard
428 materials. This is true even if the departure from statistical similarity is sufficiently small
429 to support determination of a relative potency.

430
431 In many assays multiple compounds will yield similar concentration–response curves. It
432 may be reasonable to use a biological assay system to describe or even compare
433 response curves from different compounds. But it is not appropriate to report relative
434 potency unless the Standard and Test samples contain only the same active analyte or
435 population of analytes. Biological products typically exhibit lot-to-lot variation in the
436 distribution of analytes (i.e., most biological products contain an intended product and,
437 at suitably low levels, some process contaminants that may be active in the bioassay).
438 Assessment of similarity is then, at least partially, an assessment of whether the
439 distribution of analytes in the Test sample is close enough to that of the distribution in
440 the Standard sample for relative potency to be meaningful; i.e., the assay is a
441 comparison of like to like. When there is evidence (from methods unrelated to the
442 bioassay) that the Standard and Test samples do not contain the same compound(s),
443 the biological assumption of similarity is not satisfied, and it is not appropriate to report
444 relative potency.

445
446 Other common statistical assumptions in bioassay practice are constant variance of the
447 responses around the fitted model (see Section 4.3 for further discussion), normally
448 distributed residuals (see Section 4.4 for further discussion), and independence of the
449 observations. Constant variance, normality, and independence are interrelated in the
450 practice of bioassay. For bioassays with a quantitative response, a well-chosen
451 transformation may be used to achieve near constant variance and a nearly normal
452 distribution of residuals. For such assays, once transformation of data has been
453 imposed, the remaining assumption, independence, then remains to be addressed with
454 proper modeling.

455
456 Independence of measurements in bioassay is important for effective design and
457 appropriate modeling of data. Two events are independent if they have no common
458 factor that could cause them to vary together (co-vary). Independence can be assessed

459 with statistical tools, including variance component estimates and estimates of
460 correlation among observations. A proper understanding requires knowledge of assay
461 design and the analytical model so that experimental units are accurately identified. The
462 statistical model should capture the design structure imposed by the presence of blocks
463 or grouping of responses (e.g., due to serial dilution or the use of a multichannel pipette)
464

465 **4.3 Variance Heterogeneity, Weighting, and Transformation**

466
467 Analysis of bioassay data requires that the data be approximately normally distributed
468 with constant variance across the range of the data. For linear and nonlinear regression
469 models, the variance referred to here is the residual variance from the fit of the model.
470 Constant variance is often not observed, e.g., an increase in variability with increase in
471 response (variance heterogeneity). If the variances are not equal but the data are
472 analyzed as if they were, then the estimate of relative potency may be reasonable.
473 However, failure to address nonconstant variance around the fitted dose–response
474 model results in an unreliable estimate of within-assay variance. Further, the
475 assessment of statistical similarity may not be accurate, and standard errors and
476 confidence intervals for all parameters (including a Fieller’s Theorem–based interval for
477 the relative potency) should not be used. Confidence intervals for relative potency that
478 combine multiple assays will still be correct.

479
480 Constant variance can be assessed via residual plots, Box-Cox (or power law) analysis,
481 or Bartlett’s or Levene’s test. With either of the tests, improvement of the value of the
482 statistic obtained is useful as a basis for judging whether homogeneity is improved or
483 worsened; do not rely on the *P* value. Assess variance on a large body of assay data.
484 Using only the replication variances from the current assay is not appropriate because
485 there are too few data to properly determine truly representative variances specific to
486 each concentration. This approach is necessary during development, and it is prudent
487 to re-assess variance during validation and to monitor it periodically during ongoing use
488 of the assay.

489
490 Two methods used to mitigate variance heterogeneity are weighting and transformation.
491 Lack of constant variance can be addressed with a suitable transformation chosen for
492 the assay system (not on an assay-by-assay basis). If no suitable transformation is
493 identified, weighting can be explored. Weighting is a statistical method of emphasizing
494 precise data instead of imprecise data. Transformation is used to address nonconstant
495 variation in the response and normality of residuals and to improve the fit of a statistical
496 model to the data.

497
498 Variance may be proportional to concentration or may be some function of
499 concentration. This possibility can be examined by plotting the sample variance at each
500 concentration (preferably pooled across multiple assays) against concentration and then
501 against a function of concentration (e.g., concentration squared). Variance will be
502 proportional to the function of concentration where the plot follows approximately a
503 straight line. If the data fall along a horizontal line, no weighting is needed. If such a
504 function can be found, then the weights are taken as the reciprocal of that function.

505 There may be no such function, or at least not one that is easy to describe, particularly if
506 the variation is higher at both extremes of the concentration range studied. An
507 alternative is to develop an estimate of the variability at each concentration using
508 historical data for the assay as currently conducted. Using only the replication variances
509 from the current assay is not appropriate because there are too few data to properly
510 determine truly representative variances specific to each concentration. Whether a
511 model or historical data are used, the goal is to capture the relative variability at each
512 concentration. Even with use of historical data, one does not necessarily have to
513 assume that the absolute level of variability of the current assay is identical to that of the
514 historical data but only that the ratios of variances among concentrations are consistent.

515
516 Alternatively, an approach using a power of the mean (POM) response at each
517 concentration has been advocated for many bioassay systems.

518
519 Additionally, transformation of the bioassay data may be used to help satisfy the
520 requirements of normality and variance homogeneity. Bioassay data are commonly
521 displayed with concentration on a log scale (in parallel-line or parallel-curve analyses).
522 Slope-ratio assays are displayed with concentration on the original scale.

523 Transformation may be performed to the response data as well as to the concentration
524 data. Common choices for a transformation of the response include log, square root (for
525 counts), reciprocal, and, for count data with known asymptotes, logit of the percent of
526 maximum response. Log transformations are commonly used, partly because of the
527 ease of transforming back to the original scale for interpretation.

528
529 A log–log fit may be performed on data exhibiting nonlinear behavior. The log
530 transformation is used to linearize a portion of the concentration–response relationship,
531 but such a transformation often also normalizes the data, generating responses that are
532 more uniform in their variance. Other transformations may be performed in accord with
533 the POM principle cited for weighting; i.e., data may be transformed by the inverse of
534 the POM function. A POM coefficient $k = 2$ corresponds to a log transformation of the
535 data. Note that transformation of the data requires re-evaluation of the model used to fit
536 the data. If nontransformed data are expected to follow a four-parameter logistic model,
537 another nonlinear model might be considered to fit the transformed measurements.

538
539 In many bioassays a suitable transformation (often the logarithm) yields data that have
540 nearly constant variance, near normal distribution of the residuals, and a good fit to a
541 four-parameter logistic curve. From a statistical perspective there is nothing special
542 about the original scale of measurement. Any transformation that improves accordance
543 with assumptions is acceptable. Analysts should recognize, however, that
544 transformations, choice of statistical model, and choice of weighting scheme are
545 interrelated. If a transformation is used, that may affect the choice of model. That is,
546 transforming the response, e.g., by a log or square root, may change the shape of the
547 response curve, as may the range of concentrations for which the responses are nearly
548 straight and nearly parallel.

549

550 Appropriate training and experience in statistical methods are essential in determining
551 an appropriate variance-modeling strategy. Sources of variability may be misidentified if
552 the wrong variance model is used. For example, data can have constant variation
553 throughout a four-parameter logistic concentration–response curve but can also have
554 appreciable variation in the EC_{50} parameter from block to block within the assay, or from
555 assay to assay. This assay can appear to have large variation in the response for
556 concentrations near the long-term average value of ED_{50} , if the between-block or
557 between-assay variability is not recognized. A weighted model with low weights for
558 concentrations near the EC_{50} would misrepresent a major feature of such an assay
559 system.

560 **4.4 Normality**

561
562
563 Many statistical methods for the analysis of quantitative responses assume normality of
564 the residuals. If the normality assumption is not met, the estimate of relative potency
565 and its standard error may be reasonable, but the commonly calculated confidence
566 interval for the relative potency estimate will not be valid. Most methods used in this
567 chapter are reasonably robust to departures from normality, so the goal is to detect
568 substantial nonnormality. During assay development, in order to discover substantial
569 departure from normality, one can use graphical tools such as a normal probability plot
570 or a histogram (or something similar like stem-and-leaf or box plots) of the residuals
571 from the model fit. The histogram should appear unimodal and symmetric. The normal
572 probability plot should approximately follow a straight line. A pattern about a straight line
573 is one indication of nonnormality. Nonnormal behavior may be due to measurements
574 that are log normal and show greater variability at higher levels of response. This is
575 characterized by a concave pattern in the residuals in a normal plot.

576
577 Statistical tests of normality are not generally useful. If such tests are used, as for
578 homogeneity of variance, improvement of the value of the statistic obtained is useful for
579 judging whether normality is improved or worsened. Again, reliance should not be
580 placed on the P -value. Evaluate normality on as large a body of assay data as possible
581 during development, re-assess during validation, and monitor periodically during
582 ongoing use of the assay. Important departures from normality can often be mitigated
583 with a suitable transformation. Failure to assess and mitigate important departure from
584 normality carries the risks of disabling appropriate outlier detection and losing capacity
585 to obtain reliable estimates of variation. This assessment typically is performed during
586 assay development and is not routinely conducted with each assay.

587 **4.5 Linearity of Concentration–Response Data**

588
589
590 Some bioassay analyses assume that the shape of the concentration–response curve is
591 a straight line or approximates a straight line over a limited range of concentrations. In
592 those cases, a linear-response model may be assessed to determine if it is justified for
593 the data in hand. Difference testing methods for assessing linearity face the same
594 problems as do standard methods for difference testing of parallelism. More data and
595 better precision make it more likely to detect nonlinearity. Because situations when

596 nonlinearity does not affect potency generally are rare and occur only by chance,
597 analysts should routinely assess departure from linearity if they wish to use a linear-
598 response model to estimate potency.

599
600 If an examination of a data plot reveals gross departure from linearity, this is sufficient to
601 support a conclusion that linearity is not present. High data variability, however, may
602 mask departures from linearity. Thus a general approach for linearity can conform to
603 that for similarity, developed in detail in Section 4.7, *Implementing equivalence testing*
604 *for similarity (parallelism)*:

- 605
- 606 (1) Specify a measure of departure from nonlinearity. Possibilities are the
607 nonlinearity sum of squares, assuming similar slopes, or the quadratic
608 coefficients from a fitted quadratic model or from contrast estimates from a
609 means model.
 - 610 (2) Use one of the four approaches of Section (2) of *Implementing*
611 *equivalence testing for similarity (parallelism)* to determine a range of acceptable
612 values (acceptance interval) for the measure of nonlinearity. A common
613 alternative is to determine the acceptable range by visual examination. Because
614 departures from nonlinearity are typically evident visually, the acceptable range
615 consists of those values of the measure of nonlinearity for which the nonlinearity
616 is not evident.
 - 617 (3) Determine a 90% two-sided confidence interval, following the Two One-
618 Sided Test (TOST) procedure, and compare the result to the acceptance interval
619 determined in Step 2.

620

621 **Selecting a linear range.** Often a subset of the concentrations measured in the assay
622 will be selected in order to derive a linear concentration–response curve. The subset
623 can often be identified graphically. If, in the final assay, the intent is to use only
624 concentrations in the linear range, choose a range that will result in parallel straight
625 lines for the range of relative potencies expected during routine use of the assay
626 Otherwise, the assay will fail parallelism tests when the problem is actually that the
627 potency produces assay response values outside the linear range, and repeat assays
628 will be required. The repeat assays together with the valid assays may generate a
629 biased estimate of potency because of the selective process of repeating assays when
630 the response is in the extremes of the concentration–response curve.

631

632 The problem is more complex in assays where there is even modest variation in the
633 shape or location of the concentration–response curve from run to run or from block to
634 block within an assay. In such assays, which are quite common, it is appropriate to
635 choose subsets for each sample in each assay or even in each block within an assay.
636 Note that a fixed-effects model will mask any need for different subsets in different
637 blocks, but a mixed-effects model may reveal and accommodate different subsets in
638 different blocks (see Section 4.9).

639

640 Additional specific guidance about selection of data subset(s) for linear model
641 estimation of relative potency includes the following: Use at least three and preferably

642 four adjacent concentrations; require that the slope of the linear segment is sufficiently
643 steep; and require that the lines fit to Standard and Test samples are straight and that
644 the lines are parallel. One way to derive a steepness criterion is to compute a *t*-statistic
645 on the slope. If the slope is not significant the bioassay will have poor performance.
646 Another aspect that supports requiring adequate steepness of slope is the use of subset
647 selection algorithms. Without a suitable slope steepness criterion, a subset selection
648 algorithm that seeks to identify subsets of three contiguous data points that are straight
649 and parallel might select doses on an asymptote. These are obviously poor subsets to
650 use to estimate potency. Requiring a slope that has a small *P*-value is a simple and
651 appropriate way to steer the subset selection algorithm to more appropriate doses. How
652 steep or how significant the steepness of the slope should be depends on the assay.
653 This criterion should be set during assay development and possibly refined during
654 assay validation.

655 656 **4.6 Common Bioassay Models**

657
658 Most bioassays consist of a series of concentrations or dilutions of both a Test article
659 and a Standard material. A mathematical model is fit to the dose–response data, and a
660 relative potency is then calculated from the parameters of the model. Choice of model
661 may depend on whether quantitative or qualitative data are being analyzed.

662
663 For quantitative data, parallel-profile models are preferred because of their superior
664 statistical properties. If a parallel-profile model is used, concentrations or dilutions are
665 usually scaled geometrically, usually in 2-fold, log, or half-log increments. If a slope-ratio
666 model is used, concentrations or dilutions should be scaled arithmetically. Several
667 functions may be used for fitting a parallel profile model to quantitative data, including a
668 linear function, a higher-order polynomial function, a four-parameter logistic (symmetric
669 sigmoid) function, or a five-parameter function for asymmetric sigmoids. Such functions
670 require a sufficient number of concentrations or dilutions to fit the model. It is preferable
671 to have at least four more concentrations or dilutions than the number of parameters
672 that will be estimated in the model. Thus, a linear model with two parameters might
673 require six concentrations, a quadratic model with three parameters might require seven
674 concentrations, a four-parameter model may utilize eight concentrations, and a five-
675 parameter model may use nine concentrations. The four- and five-parameter sigmoid
676 models also require concentrations that help to establish the asymptotes of the
677 response profile. Two concentrations are recommended to support each asymptote.

678
679 A linear model is sometimes selected because of efficiency and ease of processing.
680 Because bioassay response profiles are usually nonlinear, the laboratory might perform
681 an experiment with a wide range of concentrations in order to identify the approximately
682 linear region of the concentration–response profile. For data that follow a four-
683 parameter logistic model, these are the concentrations near the center of the response
684 region, usually between 20% to 80% response when the data are rescaled to the
685 asymptotes. Caution is appropriate in using a linear model because for a variety of
686 reasons the apparently linear region may shift. A stable linear region may be identified
687 after sufficient experience with the assay and with the variety of samples tested in the

688 assay. Data following the four-parameter logistic function may also be linearized by
689 transformation. The lower region of the function is approximately linear when the data
690 are log transformed (log–log fit).

691
692 Quantal data are typically fit using more complex mathematical models. A probit or logit
693 model may be used to estimate a percentile of the response curve (usually the 50th
694 percentile) or, more directly, the relative potency of the Test to the Standard. Spearman-
695 Kärber analysis is a nonmodeling method for determining the 50th percentile of a
696 quantal dose–response curve. Note that the Spearman-Kärber analysis requires a full
697 range (0% through 100% response) in response to yield an accurate assessment of
698 potency.

700 4.7 Suitability Testing

701
702 System suitability and sample suitability assessment should be performed to ensure the
703 quality of bioassay results. System suitability in bioassay, as in other analytical
704 methods, consists of prespecified criteria by which the validity of an assay (or perhaps a
705 run containing several assays) is assessed. Analysts can assess system suitability by
706 determining that some of the parameters of the Standard response are in their usual
707 ranges and that some properties (e.g., residual variation) of all the data are in their
708 usual range. To achieve high assay acceptance rates it is advisable to accept large
709 fractions of these usual ranges (99% or more) and to assess system suitability with a
710 few relatively uncorrelated measures. System suitability parameters and their ranges
711 may also be selected on the basis of empirical or simulation studies that measure the
712 impact of changes in a parameter on potency estimation.

713
714 Sample suitability in bioassay is evaluated using prespecified criteria for the validity of
715 the potency estimate of an individual Test article. System and sample validity criteria
716 should be established after the bioassay has been developed and before bioassay
717 validation. Where there is limited experience with the bioassay these criteria may be
718 considered tentative.

719
720 **System suitability.** System suitability parameters may be selected based on the design
721 and the model of the procedure. Regardless of the design and model, however, system
722 suitability parameters should be directly related to the quality of the bioassay. In
723 parallel-line assays low values of the Standard slope typically yield estimates of potency
724 with low precision. Rather than reject assays with low slope, analysts may find it more
725 effective to temporarily use additional replicate assays until the assay system can be
726 improved to consistently yield higher-precision estimates of potency. It may be
727 particularly relevant to monitor the range of response levels and location of asymptotes
728 associated with controls or Standard sample to ensure appropriate levels of response. A
729 drift or a trend in some of the criteria may indicate a systematic degradation of a critical
730 reagent or Standard material. Statistical process control (SPC) methods should be
731 implemented to detect systematic trends in system suitability parameters.

732

733 Two common measures of system suitability are assessment of the adequacy of the
734 model (goodness of fit) and of precision. With replicates, a pure error term may be
735 separated from assessment of lack of fit. Care should be taken in deriving a criterion for
736 lack of fit. The use of the wrong error term may result in an artificial assessment. A
737 laboratory may choose alternative measures of system suitability or add others if
738 deemed necessary or useful to maintain control of the assay.
739

740 For adequacy of the model, a sound measure is the lack of fit sum of squares from the
741 model for the Standard alone. This value should be compared to the distribution of such
742 values obtained historically. If the data are greater than some threshold percentile (e.g.,
743 99th) of the historical distribution, then the data are not suitable. Note that the Test data
744 are not used here. Adequacy of the model for the Test is implicitly tested with the
745 parallelism assessment.
746

747 For assessment of precision, two alternatives can be considered: One follows an
748 approach similar to the assessment of model adequacy and uses the mean squared
749 error (residual variance) from the model fit to the Standard. This requires appropriate
750 replication. However, this may be associated with a small number of degrees of
751 freedom for the variance estimate. A second alternative is to use the mean squared
752 error from the separate model fits to Standard and Test. Once the measure is selected,
753 use historical data to determine a threshold for acceptance.
754

755 **Sample suitability.** Sample suitability or sample acceptance in bioassay consists of
756 two major components: assessment of similarity and assessment of range. Bioassays
757 can report relative potency only from samples that are similar to Standard and are within
758 the range of the assay system. Samples for which a relative potency estimate is within
759 range and the samples are not similar to the Standard should be reported as nonsimilar
760 to Standard without a quantitative potency estimate.
761

762 *Similarity.* For the assessment of similarity (parallelism), an equivalence testing
763 approach offers advantages. The equivalence test criterion may be formulated as a
764 hypothesis test where the null hypothesis of nonsimilarity is that a similarity measure
765 has an absolute value less than a critical value required to demonstrate similarity. When
766 the resulting P value is below a suitable threshold, such as 0.05, there is sufficient
767 evidence to reject the null hypothesis of nonsimilarity. Also, when a (given $P = 0.05$)
768 90% confidence interval on the similarity measure is entirely between the equivalence
769 bounds, there is sufficient evidence to reject the null hypothesis of nonsimilarity.
770

771 Equivalence testing is contrasted with classical hypothesis testing (also known as
772 difference testing) in which the null hypothesis is that there is no difference between
773 Standard and Test values of a similarity parameter, because an equivalence testing
774 approach states that sufficiently similar (at most an unimportant departure from
775 similarity) is the alternative statistical hypothesis, and nonsimilarity is the null
776 hypothesis. In this manner, equivalence testing controls (via alpha, the type I error rate)
777 the risk of falsely declaring samples to be equivalent. With difference testing the risk of
778 falsely declaring samples to be equivalent is a type II error. Without careful attention to

779 and monitoring of the variation of the similarity measure, the sample size, and the
780 critical value of the similarity measure, analysts run the risk of falsely declaring samples
781 similar. With equivalence testing, less information and monitoring are needed to control
782 the critical risk of falsely declaring a sample to be similar to Standard. However,
783 appropriately justified critical values of the similarity measure are needed. “Sufficiently
784 similar” must exclude important differences.
785

786 More simply expressed: traditionally, classical hypothesis testing (difference testing)—
787 not equivalence testing—has been used to establish parallelism between a Test sample
788 and the Standard sample. From such an approach the laboratory cannot conclude that
789 the slopes are equal. The data may be too variable, or the assay design may be too
790 weak to establish a difference. The laboratory can conclude that the slopes are
791 sufficiently similar using the equivalence testing approach.
792

793 Equivalence testing has practical advantages compared to difference testing, including
794 that improved assay replication or precision will increase the chances that samples will
795 pass the similarity criteria, that decreased assay replication or precision will decrease
796 the chances that samples will pass the similarity criteria, and that sound approaches to
797 combining data from multiple assays of the same sample to better understand whether
798 a sample is truly similar to Standard or not are obtained.
799

800 The challenge in implementing equivalence testing is in setting appropriate equivalence
801 bounds for the similarity measures. Ideally, information is available to link departures
802 from similarity to relevant outcomes in clinical or animal data. Information may be
803 available from evaluation of similar assays. Because different assays have responses
804 and doses measured on different scales, it is challenging to compare changes in slope,
805 intercepts (in slope-ratio assays), or asymptotes across assays. For perspective it can
806 be useful to consider slope differences as a percentage of the Standard slope. Similarly,
807 it can be useful to consider asymptote (or intercept) differences as a percentage of the
808 range of the Standard asymptotes. Though these may not be universal experiences,
809 some have observed that changes in slopes smaller than 30% and changes in
810 asymptotes smaller than 5% are rarely meaningful. Setting critical values for changes in
811 slopes > 50% or changes in asymptotes > 15% should be done with special care to
812 ensure that there is a solid scientific case that these limits are reasonable. In routine
813 use of the assay, the critical values of similarity parameters should not be on
814 percentage scale because these scales have more uncertainty (wider confidence limits).
815

816 Because of the advantages associated with the use of equivalence testing in the
817 assessment of similarity, analysts can transition existing assays to equivalence testing
818 or can implement equivalence testing methods when changes are made to existing
819 assays. It is informative to examine the risk that the assay will fail good samples. This
820 risk depends on the precision of the assay system, the amount of replication in the
821 assay system, and the critical values of the similarity parameters (a process capability
822 analysis is one approach to this risk analysis). One approach to transitioning an
823 established assay from difference testing to equivalence testing (for similarity) is to use
824 the process capability of the assay to set critical values for similarity parameters. This

825 approach is reasonable for an established assay because the risks (of falsely declaring
826 samples similar and falsely declaring samples nonsimilar) are implicitly acceptable
827 because of the assay's history of successful use. Equivalence testing improves the
828 relationships between these risks with changes in assay precision or assay sample size.
829 In contrast, a new assay's capacity for keeping the risks at acceptable levels is
830 unknown, and using a process capability analysis to set critical values of similarity
831 measures would yield an assay system with unknown and unmanaged risks.

832
833 The similarity measures have, in many cases, interpretable, practical meaning in the
834 assay (these measures are based on the parameters of the dose–response curve and
835 include, e.g., slope for a straight parallel-line assay, intercept for a slope-ratio assay,
836 slope and asymptotes for a four-parameter logistic parallel-line assay, and slope,
837 asymptotes, and nonsymmetry parameter in a five-parameter sigmoid model). In some
838 cases these changes in curve shape predict certain classes of interaction between the
839 product and the bioassay. When possible, discussion of these interactions and their
840 likely impacts is a valuable part of setting appropriate equivalence boundaries.

841
842 *Implementing equivalence testing for similarity (parallelism).* Many statistical procedures
843 for assessing similarity are based on a statistical null hypothesis consistent with
844 similarity and an alternative hypothesis of some degree of nonsimilarity. Failure to find
845 statistically significant nonsimilarity is then taken as a conclusion of similarity. However,
846 failure to detect nonsimilarity does not prove similarity. Equivalence testing provides a
847 method for the analyst to proceed to a conclusion (if warranted by the data) of
848 sufficiently similar while controlling the risk of doing so inappropriately. Following is a
849 sequence for this process:

850
851 (1) Choose a measure of nonsimilarity. For the parallel-line case, this could
852 be the difference or ratio of slopes. (The ratio of slopes can be less sensitive to
853 the value of the slope. Also, framing the slope difference as a proportional
854 change from Standard rather than in absolute slope units has an advantage
855 because it is invariant to the units on the concentration and response axes.) For
856 a slope-ratio assay, the measure of nonsimilarity involves the difference in y-
857 intercepts between Test and Standard samples. Again, it can be advantageous
858 to frame this difference as a proportion of the (possibly transformed) response
859 range of Standard to make the measure invariant to the units of the response.
860 For the four-parameter logistic model, similarity between Standard and Test
861 samples must be assessed on the basis of three parameters: the upper
862 asymptote, the slope, and the lower asymptote. If sigmoid curves with additional
863 parameters are used to fit bioassay data, it is also important to consider
864 addressing similarity between Standard and Test preparations of the additional
865 curve parameters (e.g., asymmetry parameter of the five-parameter model). The
866 comparison could be based on the individual parameters, one at a time, or some
867 single composite measure of nonparallelism. One such composite measure is the
868 parallelism sum of squares. This is found as the difference in residual sum of
869 squares (RSSE) between the value obtained from fitting the Standard and Test
870 curves separately and the value obtained from imposing parallelism:

871 Parallelism sum of squares = $RSSE_P - RSSE_S - RSSE_T$
872 where the subscripts P, S, and T denote Parallel model, Standard model, and
873 Test model, respectively. With any composite measure, the analyst must
874 consider the implicit relative weighting of the importance of the three (or more)
875 parameters and whether the weighting is appropriate for the problem at hand.
876 For the parallelism sum of squares, for example, with nonlinear models, the
877 weighting given to the comparison of the asymptotes depends on the amount of
878 data in the current assay on and near the asymptotes.

879
880 (2) Specify a range of acceptable values, typically termed an equivalence
881 interval or “indifference zone,” for the measure of nonsimilarity. When a ratio of
882 slopes is used as a measure of nonsimilarity (nonparallelism), that measure is
883 free to vary above or below a ratio of 1.0 (a ratio of 1.0 indicates perfect
884 parallelism). The acceptable values for nonparallelism then constitute an interval.
885 Following are four approaches that can be used to determine this interval. If
886 pharmacopeial limits have been specified for a defined measure of nonsimilarity,
887 then the assay should satisfy these requirements.

888
889 a. The first approach is to compile historical data that compare the
890 Standard to itself and to determine, from the historical data, the
891 equivalence interval as a tolerance interval for the measure of
892 nonparallelism. The advantage of using historical data is that they give the
893 laboratory control of the false-failure rate (the rate of failing an assay that
894 is in fact acceptable). The disadvantage is that there is no control of the
895 false pass rate (the rate of passing an assay that may have an
896 unacceptable difference in performance relative to the Standard). The
897 equivalence interval specification is driven solely by assay capability.
898 Laboratories that use this approach should take care that an assay in
899 need of improvement is not driving overly wide equivalence intervals. Also
900 note that any change in assay capability means changing the equivalence
901 interval—the interval must be relevant to the assay as it is currently
902 conducted. Because of its limitations, this approach is best suited only for
903 early stages in assay development where historical data may be the only
904 information available.

905
906 b. The approach of (a) is simple to implement in routine use and can
907 be used with assay designs that do not provide reliable estimates of
908 within-assay variation and hence confidence intervals. However, there is a
909 risk that assays with larger than usual amounts of within-assay variation
910 can pass inappropriately. The preferable alternative to (a) is therefore to
911 determine a tolerance interval for the confidence interval for the measure
912 of nonparallelism. The following is particularly appropriate to transition an
913 existing assay with a substantial body of historical data on both Standard
914 and Test samples from a difference testing approach to an equivalence
915 approach:

- 916 i. For each value of the measure of nonparallelism from the historical
917 data, determine a 95% confidence interval, (a, b).
918 ii. For each confidence interval, determine its maximum departure
919 from perfect parallelism. This is $\max(|a|, |b|)$ for differences,
920 $\max(1/a, b)$ for ratios, and simply b for quantities that must be
921 positive, such as a sum of squares.
922 iii. Determine a tolerance interval for the maximum departures
923 obtained in (ii). This will be a two-sided tolerance interval for
924 differences and ratios and a one-sided tolerance interval for
925 necessarily positive quantities. A nonparametric tolerance interval
926 approach is preferred.
927 iv. "Sufficiently parallel" is concluded for new data if the confidence
928 interval for the measure of nonparallelism falls completely within the
929 interval determined in (iii).
930

931 As with (a), this approach is best suited for early stages in assay
932 development.
933

934 c. The third approach starts with historical data comparing the
935 Standard to itself and adds data comparing the Standard to known
936 failures, e.g., to degraded samples. Compare values of the measure of
937 nonparallelism for data for which a conclusion of parallelism is appropriate
938 (Standard against itself) and data for which a conclusion of parallelism is
939 not appropriate, e.g., degraded samples. Based on this comparison,
940 determine a value of the measure of nonparallelism that discriminates
941 between the two cases. If using this approach, take care to include more
942 than extreme examples. The work should include a range of samples for
943 which a conclusion of parallelism is not appropriate, and be sure to cover
944 samples with the minimal important change. For nonlinear models, this
945 comparison also can be used to determine which parameters should be
946 assessed. Some may not be sensitive to the failures that can occur with
947 the specific assay.
948

949 d. The fourth approach is based on what is known of the product and
950 the assay, much as (0.80, 1.25) is used for bioequivalence of generic
951 products, or based on general experience. For example, when the ratio of
952 slopes is not more than 1.3 and not less than 0.77, departure from
953 parallelism is difficult to discern. In contrast, ratios greater than 1.5 or less
954 than 0.67 are typically substantial. A related consideration could be the
955 therapeutic index of the drug
956

957 For approaches (a) and (b) of Step 2, compare the obtained value of the
958 measure of nonparallelism (a) or its confidence interval (b) to the interval
959 obtained in Step 2. The value must be within the limits if one uses (a), or the
960 confidence interval must be completely within the limits if one uses (b).
961

962 For approach (c) of Step 2, an alternative approach can be used. This approach
963 essentially treats the parallelism as a discrimination problem. The choice of the
964 cut point in (c) should take into account the rates of false positive and false
965 negative decisions (and the acceptable risks to the laboratory) and should reflect
966 the between-assay variability in precision. Thus it is reasonable to compare the
967 point estimate of the measure of nonparallelism to the cut point and to not use
968 confidence intervals. This approach is simpler to implement in routine use and
969 can be used with assay designs that cannot provide reliable estimates of within-
970 assay variation.

971
972 For approach (d) of Step 2, demonstrate that the measure of nonparallelism is
973 significantly greater than the lower endpoint of the acceptance interval and
974 significantly less than the upper endpoint. (If the acceptance interval is one-
975 sided, then apply only the single applicable test.) This is use of the TOST
976 procedure. For most situations, TOST can be most simply implemented by
977 calculating a 90% two-sided confidence interval, which corresponds to a 5%
978 equivalence test. If this confidence interval lies entirely within the equivalence
979 interval specified in Step 2, then similarity is sufficiently demonstrated. For
980 parallel-line models, one can use the confidence interval based on value $\pm k$
981 times the standard error of the value for the difference of slopes or Fieller's
982 Theorem for the ratio of slopes. (This also applies to the difference of intercepts
983 in slope-ratio models.) For nonlinear models, there is evidence that these simple
984 confidence interval methods do not attain the stated level of confidence, and
985 methods based on likelihood profile or resampling are more appropriate.

986
987 An alternative to the approach described above is to use an average (historical) value
988 for the variance of the ratio or difference in a similarity parameter—obtained from some
989 number of individual assays—to compute an acceptance interval for a point estimate of
990 the similarity parameter. This approach is simpler to implement in routine use and can
991 be used with assay designs that are unable to provide reliable estimates of within-assay
992 variation. However, there is a price. The equivalence testing approach that relies on
993 assay-specific (within-assay) measure(s) of variation (i.e., the confidence intervals) is
994 conservative in the sense that it will fail to pass similarity for samples from assays that
995 have larger than usual amounts of within-assay variation. Using an acceptance region
996 for a similarity parameter—rather than an acceptance region for confidence intervals for
997 the similarity parameter—loses this conservative property and hence is not preferred
998 where alternatives exist.

999
1000 **Range.** The range for a relative potency bioassay is the interval between the upper and
1001 lower relative potencies for which the bioassay is demonstrated to have a suitable level
1002 of similarity, precision, accuracy, and assay linearity. It is straightforward to determine
1003 whether or not a sample that is similar to Standard has a relative potency within the
1004 (validated) range of the assay system. For samples that are apparently not similar
1005 according to standard assay methodology, it is more challenging to determine whether a
1006 relative potency estimate for the sample might be obtained. In a nonlinear parallel-line
1007 assay a sample that does not have data on one asymptote can usually be assumed to

1008 be out of the potency range of the assay. In a parallel straight-line assay a sample that
1009 does not have three or more points on the steep portion of the response curve can also
1010 be assumed to be out of the potency range of the assay. In some assay systems fit with
1011 nonlinear models, it may be reasonable to declare a sample to be out of range based on
1012 an estimate of relative potency that is not predicated on a demonstration of similarity. In
1013 some assay systems fit with parallel straight-line models, it may be reasonable to
1014 declare a sample out of range using a crude potency estimate based on a subset of the
1015 Test sample that contains only two doses.

1016 1017 **4.8 Outliers** 1018

1019 Bioassay data should be screened for outliers before relative potency analysis. Outliers
1020 may be a simple random event or a signal of a systematic problem in the bioassay.
1021 Systematic error that generates outliers may be due to a dilution error at one or more
1022 concentrations of a Test article or the Standard or due to a mechanical error (e.g.,
1023 system malfunction). Several approaches for outlier detection can be considered. Visual
1024 inspection is frequently utilized but should be augmented with a more objective
1025 approach to avoid potential bias.

1026
1027 An outlier is a datum that appears not to belong among the other data present. An
1028 outlier may have a distinct, identifiable cause, such as a mistake in the bench work,
1029 equipment malfunction, or a data recording error, or it could just be an unusual value
1030 relative to the variability typically seen and may appear without an identifiable cause.
1031 The essential question pertaining to an outlier becomes: Is the apparent outlier sampled
1032 from the same population as the other, less discordant, data, or is it from another
1033 population? If it comes from the same population and the datum is, therefore, an
1034 unusual (yet still legitimate) value obtained by chance, then the datum should stand. If it
1035 comes from another population and the datum's excursive value is due to human error
1036 or instrument malfunction, then the datum should be omitted from calculations. In
1037 practice, the answer to this essential question is often unknown, and investigations into
1038 causes are often inconclusive. Outlier management relies on procedures and practices
1039 to yield the best answer possible to that essential question and to guide response
1040 accordingly.

1041
1042 General Chapter *Analytical Data—Interpretation and Treatment* <1010>, addresses
1043 outlier labeling, identification, and rejection, including statistical methods, and provides
1044 material that the bioassay practitioner will find useful. General Chapter <1010> also lists
1045 additional sources of information that can provide a comprehensive review of the
1046 relevant statistical methodology. General Chapter <1010> makes no explicit remarks
1047 regarding outlier analysis in linear or nonlinear regression. Outlier analysis techniques
1048 appropriate for data obtained from regression of response on concentration can be
1049 used. Some remarks about outliers are provided here in the context of bioassays to
1050 emphasize or complement the information in <1010>.

1051
1052 Of the procedures employed for analysis of drug compounds and biological drugs, the
1053 bioassay can be expected to be the most prone to outlying data. The management of

1054 outliers is appropriate with bioassay data on at least two levels: where an individual
1055 datum or a group of data (e.g., data at a concentration) can be checked against
1056 expected responses for the sample and concentration; and, separately, when estimates
1057 of relative potency from an assay can be checked for consistency with other
1058 independent estimates of the potency of the same material.

1059

1060 Three important aspects of outlier management are prevention, labeling, and
1061 identification.

1062

1063 Outlier prevention is obviously preferred and is facilitated by procedures that are less
1064 subject to error and by checks that are sensitive to the sorts of errors that, given the
1065 experience gained in assay development, may be expected to occur. In effect, the error
1066 never becomes an outlier because it is prevented from occurring.

1067

1068 Good practice calls for the examination of data for outliers and labeling (“flagging”) of
1069 the apparently outlying observation(s) for investigation. If investigation finds a cause,
1070 then the outlying datum may be excluded from analysis. Because of the ordinary
1071 occurrence of substantial variability in bioassay response, a laboratory’s investigation
1072 into the outlying observation is likely to yield no determinable cause. However, the lack
1073 of evidence regarding an outlier’s cause is not a clear indication that statistical outlier
1074 testing is warranted. Knowledge of the typical range of assay response variability should
1075 be the justification for the use of statistical outlier tests.

1076

1077 Outlier identification is the use of rules to confirm that the values are inconsistent with
1078 the known or assumed statistical model. For outliers with no determined cause, it is
1079 tempting to use statistical outlier identification procedures to discard unusual values.
1080 Discarding data solely because of statistical considerations should be a rare event.
1081 Falsely discarding data leads to overly optimistic estimates of variability and can bias
1082 potency estimates. The laboratory should monitor the failure rate for its outlier
1083 procedure and should take action when this is significantly higher than expected.

1084

1085 Statistical procedures for outlier identification depend on assumptions about the
1086 distribution of the data without outliers. Identification of data as outliers may mean only
1087 that the assumption about distribution is not correct. If dropping outliers because of
1088 statistical considerations is common, particularly if outliers tend to occur more often at
1089 high values or at high responses, then this may be an indication that the data require
1090 some adjustment, such as log transformation, as part of the assay procedure. Following
1091 are presentations of two approaches to statistical assessment of outlying data:
1092 replication based and model based.

1093

1094 **Replication-based approaches.** When replicates are performed at concentrations of a
1095 Test article and the Standard, an “extra variability” (EV) criterion may be employed to
1096 detect outliers. Historical data can be analyzed to determine the range in variability
1097 commonly observed among replicates, and this distribution of ranges can be used to
1098 establish an extreme in the range that might signal an outlier. Metrics that can be
1099 utilized are the simple range (maximum replicate minus minimum replicate), the

1100 standard deviation, or the CV or RSD among replicates. However, if the bioassay
1101 exhibits heterogeneity of variability, assumptions about uniform scatter of data are
1102 unsupported. Analysts can use a variable criterion across levels in the bioassay, or they
1103 can perform a transformation of the data to a scale that yields homogeneity of
1104 variability. Transformation can be performed with a POM approach as discussed
1105 previously. Where heterogeneity exists nonnormality is likely present, and the range
1106 rather than standard deviation or RSD should be used.

1107
1108 The actions taken upon detection of a potential outlier depend in part on the number of
1109 replicates. If EV is detected within a pair ($n = 2$) at a concentration of a Test article or
1110 the Standard, it will not always be clear which of the replicates is aberrant, and the
1111 laboratory should eliminate the concentration from further processing. If more than two
1112 replicates are performed at each dilution the laboratory may choose to adopt a strategy
1113 that identifies which of the extremes may be the outlier. Alternatively the laboratory may
1114 choose to eliminate the dilution from further processing.

1115
1116 Outliers may also be detected among replicate potency measurements on a Test article
1117 from a series of runs. Methods for outlier detection are described in <1010>, or analysts
1118 may use methods described above for replicates within a dilution of a Test article or the
1119 Standard.

1120
1121 **Model-based approaches.** Model-based approaches may be used to detect outliers
1122 within bioassay data. These approaches use the residuals from the fit of the appropriate
1123 bioassay model (the difference between the observed response and the response
1124 predicted by the model). Statistical methods are available for assessing outliers among
1125 residuals in this framework. Statistical support should be obtained in order to identify a
1126 method that is suitable for a particular application.

1127
1128 Lastly, an alternative to discarding outlying data is to use robust methods that are less
1129 sensitive to influence by outlying observations. Use of the median rather than the mean
1130 to describe the data's center exemplifies a robust perspective. Also, regression using
1131 the method of least squares, which underlies many of the methods in this chapter, is not
1132 robust in the presence of outliers. The use of contemporary methods of robust
1133 regression may be appropriate. This is left for investigation beyond this chapter.

1134 1135 **4.9 Fixed and Random Effects in Models of Bioassay Response**

1136
1137 The choice of treating design factors as fixed or random is important both to the design
1138 and the statistical analysis of the assay. Fixed effects are factors for which all levels, or
1139 all levels of interest, are discretely present, e.g., concentration. Fixed effects, such as
1140 temperature and duration of thaw, are expected to cause a consistent shift in
1141 responses. Analysts study fixed effects by controlling them in the design and examining
1142 changes in means across levels of the factor. In a bioassay, sample and concentration
1143 are fixed effects. Studies of other fixed effects are typically performed early in assay
1144 development as part of the process of optimizing assay performance.

1145

1146 Random effects are factors whose levels in a particular run of an assay are considered
1147 representative of levels that could have been present. That is, there is no expectation of
1148 a particular value for a random effect. Rather, that value may vary subject to some
1149 expected distribution of values and thus may be a source of variability. Examples of
1150 random effects include reagent lot, operator, or date of assay if there is no interest in
1151 specific reagent lots, operators, or dates as sources of variability. Analysts study
1152 random effects by measuring the variance components corresponding to each random
1153 effect. As with fixed effects, these are profitably studied during assay development.
1154

1155 The choice of treating a factor as fixed or random is important to the design of the assay
1156 and to proper reporting of its precision. Treating all factors as fixed, for example, leads
1157 to an understatement of assay variability because it ignores all sources of variability
1158 other than replication. The goal is to identify specific sources of variability that can be
1159 controlled, to properly include those factors in the design, and then to include other
1160 factors as random.

1161
1162 Additionally, if the factor may switch from random to fixed effect or vice versa, the factor
1163 should be modeled as a random effect. For example, reagent lots cannot be controlled,
1164 so different lots are typically considered to cause variability, and reagent lot would be
1165 considered a random effect. However, if a large shift in response values has been
1166 traced to a particular lot, a comparison among a set of lots could be performed using the
1167 predicted levels of each lot's fixed effect. Similarly, within-assay location (e.g., block,
1168 plate, plate row, plate column, or well) or sequence may be considered a source of
1169 random variation or a source of a consistent (fixed) effect.
1170

1171 Assay designs that consist of multiple factors are efficient, but they require
1172 corresponding statistical techniques that incorporate the factors as fixed or random
1173 effects in the analysis. If all factors are fixed, the statistical model is termed a fixed-
1174 effects model. If all are random, it is termed a random-effects model. If some factors are
1175 fixed and some random, the model is a mixed-effects model. Note that the concepts of
1176 fixed and random effects apply to models for qualitative and integer responses as well
1177 as for quantitative responses.
1178

1179 For assay designs that include multiple experimental units (e.g., samples assigned to
1180 sets of tubes and concentrations assigned to preplate tubes) a mixed-effects model in
1181 which the experimental units are treated as random effects is particularly effective.
1182

1183 **5. Stages in the Bioassay Process**

1184 **5.1 Design: assay layout, replication strategy, blocking, and randomization**

1185
1186 **Assay layout.** The assay layout depends on the type of bioassay being performed. In
1187 general all bioassays include operational units such as wells of a plate or individual
1188 animals to which one of a series of concentrations of a Test article or a Standard is
1189 administered. The numbers and levels of concentrations, as well as numbers and origin
1190 of the replicates of each concentration constitute a portion of the assay layout. Layout
1191

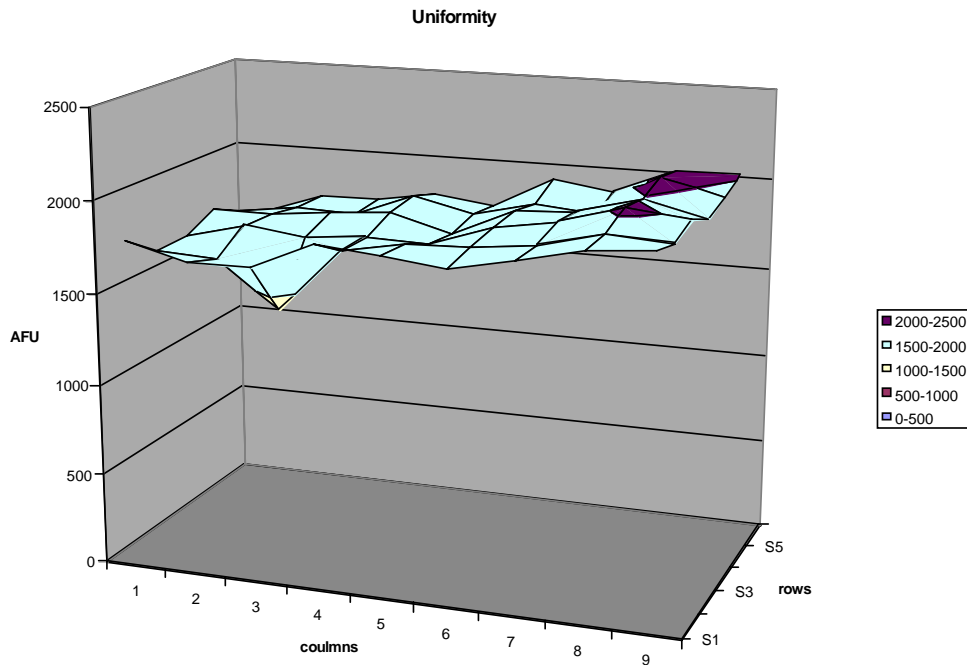
1192 considerations should be weighed against the potential for bias and their impact on
1193 bioassay variability. To the extent possible, operations that potentially can lead to bias
1194 should be managed so that the bias is translated into random variability. This is
1195 because although random variability can be effectively managed by replication, bias
1196 cannot.

1197
1198 Most cell-based assays are performed using a cell culture plate (6-, 12-, 96-, or 384-well
1199 microtiter plate) that not only must support cell culture in a uniform fashion from well to
1200 well but may also undergo a series of wash steps and incubations that may induce
1201 different cell responses in certain regions of the plate. One goal during development is
1202 to find operating conditions for the assay so that the cell responses to the analyte are
1203 free of systematic effects across the plate, across multiple plates, and across assays.
1204 Consistent performance is an important component of assay capability. Even assay
1205 conditions intended to minimize the potential for systematic bias (e.g., good analyst
1206 technique, careful calibration of pipettes, control of incubation time or temperature) often
1207 yield systematic gradients on the plate. These gradients may occur across rows, across
1208 columns, or from the edge to the center of the plate and are termed *plate effects*. Even
1209 modest-sized or inconsistent gradients (across plates, days, or analysts) should be
1210 addressed during assay development (plate layout, randomization, and replication).

1211
1212 Plate effects can be evaluated in a *uniformity trial* in which a single assay dose (often
1213 chosen from the portion of the assay curve that represents the greatest sensitivity to
1214 change in dose) is used across the entire plate. In the 3D plot seen in Figure 1, a trend
1215 of decreasing signal is evident from right to left. In this instance, the wash program for
1216 the platewasher had to be adjusted because plates were being washed harder on the
1217 left side of the plate. Another common plate effect is a differential cell-growth pattern in
1218 which the outer wells of the plate grow cells in such a way that the assay signal is
1219 attenuated. When this occurs, it is best to not utilize the outer wells of the assay plate.

1220
1221

1222 Figure 1. 3D plot of change in sensitivity across a plate (note decreasing sensitivity from
1223 right to left).



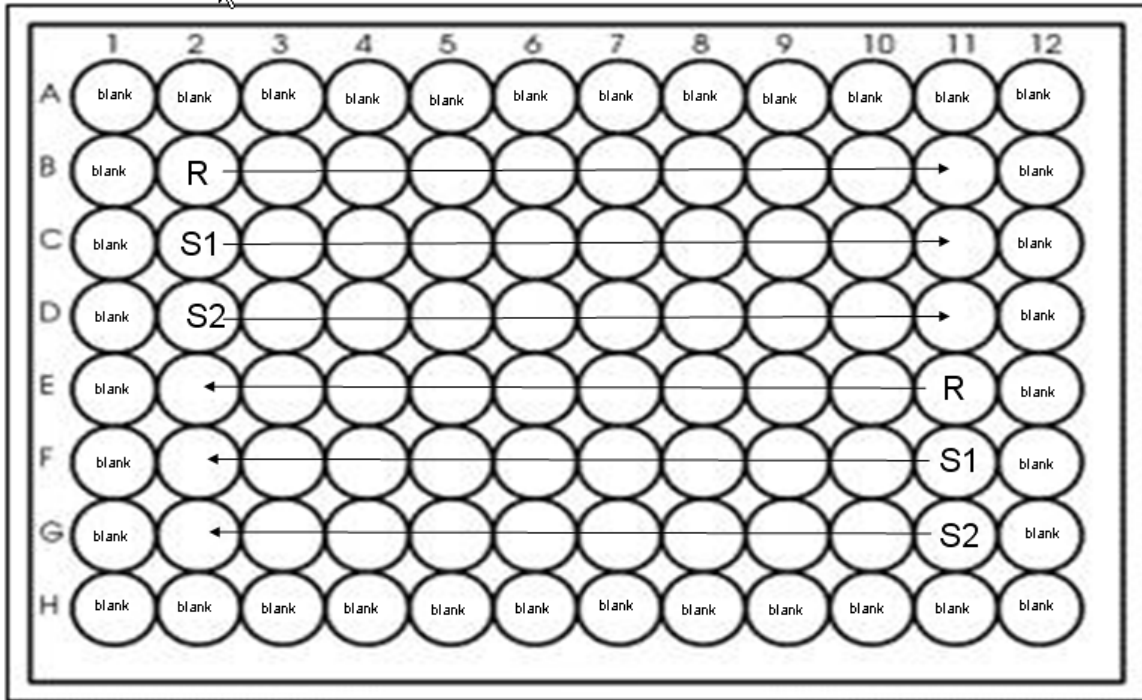
1224
1225
1226 Complete randomization of sample treatments and dilutions has been advocated as the
1227 best approach to minimizing assay bias or, more accurately, to protecting the assay
1228 results from known and unknown sources of bias by converting the latter into variance.
1229 Randomization constitutes a good strategy because it addresses systematic bias.
1230 Increased replication alone is likely to improve precision, but it cannot address
1231 systematic bias without some degree of randomization. Even when complete
1232 randomization is not practical, a plate layout can be designed to minimize plate effects
1233 by alternating sample positions across plates and the pattern of dilutions within and
1234 across plates. When no plate effects are evident, the plate layout design should, at
1235 minimum, alternate sample positions across plates within an assay run to accommodate
1236 possible bias introduced by the analyst or equipment on a given day. It is prudent to use
1237 a balanced rotation of layouts on plates so that the collection of replicates (each of
1238 which uses a different layout) provides some protection against likely sources of bias.
1239

1240 Figure 2 illustrates the design of a single assay plate wherein the samples are titrated in
1241 two different directions. This can help protect against plate column effects. Figure 3
1242 illustrates a simple alternation of Test (Test sample 1 = "1"; Test sample 2 = "2") and
1243 Standard positions across plates (across a single assay run). If the samples are
1244 alternated to different positions across an assay run, the design can help protect against
1245 plate row effects. Combining the two methods can effectively help convert plate bias
1246 into assay variance, which can be addressed by increased assay replication (increased
1247 number of plates in an assay).
1248

1249 For additional discussion of randomization, see Blocking and randomization, below.

1250
1251

Figure 2. A single assay plate in which samples are titrated in two different directions.



1252
1253
1254
1255

Figure 3. Schematic of a single assay plate with simple alternation of Test samples' and Standard positions.

| Plate Row | Plate 1 | Plate 2 | Plate 3 |
|-----------|---------|---------|---------|
| B | R | 2 | 1 |
| C | 1 | R | 2 |
| D | 2 | 1 | R |
| E | R | 2 | 1 |
| F | 1 | R | 2 |
| G | 2 | 1 | R |

1256
1257
1258
1259
1260
1261
1262
1263
1264
1265

Replication strategy. Concentrations of a Test article and the Standard can be derived in multiple ways. Laboratories that use manual multichannel pipettes often perform serial dilutions wherein each dilution is prepared from one or more previous dilutions of the material. Alternatively the laboratory may prepare independent dilutions from the Test article and Standard to obtain the concentration series. These two strategies result in the same nominal concentrations, but they have very different properties related to bias and noise. Serial dilutions are subject to propagation of error across the dilution series, and a dilution error made at an early dilution will result in a shift in the dilution

1266 profile throughout the series. Independent dilutions help mitigate the bias resulting from
1267 such a dilution error. Although serial dilutions are easier to perform in manual
1268 operations, independent dilutions can be readily obtained with robotics.

1269
1270 Replicates at each concentration similarly may be achieved in several ways. When
1271 multiple units of volume from a single point in a dilution series are dispensed into
1272 multiple wells, pseudoreplicates are present. Other strategies that provide more useful
1273 information are using independent dilution series of the Test article and Standard or
1274 making independent dilution replicates at each dilution point from each material. The
1275 reward for independent replicates is less potential bias and greater precision.

1276
1277 The number and levels of concentrations depend on desired throughput and the model
1278 that will be used to process bioassay data. The number of concentrations should
1279 increase as the number of model parameters increases. Linear modeling requires fewer
1280 concentrations than do nonlinear logistic regression models (see Section 3).

1281 Additionally, the levels of concentrations should be selected to support the model. For a
1282 linear model at least four concentrations are desirable to assess goodness of fit and
1283 obtain reliable estimates of relative potency. For four- or five-parameter logistic models
1284 it is useful to have at least four concentrations in the “linear” region of the curve and at
1285 least two concentrations in the minimum and maximum regions (asymptotes).

1286

1287 **Blocking and randomization.** Blocking and randomization are used to transform bias
1288 into random variability. Bias is introduced into bioassay measurement by operational
1289 factors related to instrumentation and the bioassay environment. Operational factors
1290 that may induce bias are edge effects in plate-based bioassays, reaction across a plate
1291 in a time-sensitive bioassay, and cage location in animal-based bioassays. The
1292 influence of these effects is mitigated by blocking and randomization, either separately
1293 or in combination. Blocking is the grouping of related experimental units in experimental
1294 designs, such as a dilution series of a Test article or the Standard. Randomization is the
1295 process of assigning treatment to experimental units based on chance so that all equal-
1296 sized groups of units have an equal chance of receiving a given treatment, e.g., the
1297 placement of a bioassay unit or a block in random configuration in a bioassay.

1298

1299 Different blocking and randomization strategies can be utilized by a bioassay laboratory.
1300 Figure 4 illustrates a poor plate-based assay design. Dilutions and replicates of the Test
1301 preparations (A and B) and the Standard (R) are placed together on the plate. Bias due
1302 to a plate or incubator effect can influence part or all the concentrations of one of the
1303 samples.

1304

1305 Figure 4. Schematic of a poorly designed plate (note that dilutions and replicates of the
 1306 Test preparations and the Standard are placed together on the plate).

| | | | | | | | | | | |
|--|----|----|----|----|----|----|----|----|----|-----|
| | | | | | | | | | | |
| | A1 | A2 | A3 | A4 | A5 | A6 | A7 | A8 | A9 | A10 |
| | A1 | A2 | A3 | A4 | A5 | A6 | A7 | A8 | A9 | A10 |
| | B1 | B2 | B3 | B4 | B5 | B6 | B7 | B8 | B9 | B10 |
| | B1 | B2 | B3 | B4 | B5 | B6 | B7 | B8 | B9 | B10 |
| | R1 | R2 | R3 | R4 | R5 | R6 | R7 | R8 | R9 | R10 |
| | R1 | R2 | R3 | R4 | R5 | R6 | R7 | R8 | R9 | R10 |
| | | | | | | | | | | |

1307

1308

1309 A split-plot design, an alternative that randomizes samples across plate rows and
 1310 dilutions (concentrations) within each row, is seen in Figure 5. Such a strategy is
 1311 tedious and requires robotics for routine reliability.

1312

1313 Figure 5. Schematic of a split-plot design that randomizes samples across plate rows
 1314 and dilutions (concentrations).

| | | | | | | | | | | |
|--|-----|----|-----|----|----|-----|----|-----|----|-----|
| | | | | | | | | | | |
| | A3 | A6 | A10 | A2 | A5 | A1 | A4 | A9 | A8 | A7 |
| | R6 | R9 | R1 | R2 | R5 | R10 | R3 | R8 | R7 | R4 |
| | B1 | B5 | B2 | B8 | B3 | B10 | B7 | B9 | B4 | B6 |
| | R5 | R1 | R9 | R4 | R2 | R8 | R3 | R10 | R6 | R3 |
| | A4 | A2 | A8 | A7 | A1 | A9 | A3 | A10 | A5 | A6 |
| | B10 | B3 | B7 | B2 | B4 | B5 | B9 | B6 | B1 | B10 |
| | | | | | | | | | | |

1315

1316

1317 A layout that provides some protection from plate effects and can be performed
 1318 manually is a strip-plot design, shown in Figure 6.

1319

1320 Figure 6. Schematic of a layout that provides some protection from plate effects and can
 1321 be performed manually.

| | | | | | | | | | | |
|--|-----|----|----|----|----|----|----|----|----|-----|
| | | | | | | | | | | |
| | A1 | A2 | A3 | A4 | A5 | A6 | A7 | A8 | A9 | A10 |
| | R1 | R2 | R3 | R4 | R5 | R6 | R7 | R8 | R9 | R10 |
| | B1 | B2 | B3 | B4 | B5 | B6 | B7 | B8 | B9 | B10 |
| | R10 | R9 | R8 | R7 | R6 | R5 | R4 | R3 | R2 | R1 |
| | A10 | A9 | A8 | A7 | A6 | A5 | A4 | A3 | A2 | A1 |
| | B10 | B9 | B8 | B7 | B6 | B5 | B4 | B3 | B2 | B1 |
| | | | | | | | | | | |

1322

1323

1324 Here samples are again randomized to rows of a plate, but dilution series are performed
 1325 in different directions in different locations on the plate to mitigate bias across columns
 1326 of the plate. An added advantage of the strip-plot design is the ability to detect a plate
 1327 effect by the interaction of sample and dilution direction (left-to-right or right-to-left).

1328

1329 For related discussion pertaining to plate effects in cell-based assays, see *Assay layout*,
1330 above.

1331

1332 5.2 Development

1333

1334 Assay development ends with the completion of the standard operating procedure
1335 (SOP) for the bioassay. The SOP is a documented protocol for bioassay
1336 implementation. The SOP should include enough detail so that a qualified laboratory
1337 with a trained analyst can perform the procedure in a routine manner. A strategic part of
1338 development is a look forward toward performance maintenance. Standard operating
1339 procedures for reagent and technician qualification, as well as for calibration of the
1340 working Standard, help complete the bioassay development package.

1341

1342 **One factor at a time vs design of experiments.** Bioassay development proceeds
1343 through a series of experiments in which conditions and levels of assay factors are
1344 varied to identify those that support a reliable and robust bioassay qualified for routine
1345 use. Those experiments may be conducted one factor at a time (OFAT), studying each
1346 parameter separately to identify ideal conditions, or through the use of multi-factor
1347 design of experiments (DOE). DOE is an efficient and effective strategy for developing a
1348 bioassay and optimizing bioassay performance, thus helping to obtain a measurement
1349 system that meets its requirements. In comparison to OFAT, DOE generally requires
1350 fewer experiments and also provides insight into interactions of factors that affect
1351 bioassay performance. Assay development and optimization using DOE proceeds
1352 through a series of steps: process mapping and risk analysis; screening; response
1353 optimization; and confirmation.

1354

1355 **Process mapping and risk analysis.** Bioassay optimization begins with a systematic
1356 examination of bioassay factors and a risk assessment to identify those factors that may
1357 influence performance. It is useful to visualize bioassay factors using a bioassay
1358 process map such as a cause-and-effect or fishbone diagram. Using the process map
1359 as a guide, the laboratory can examine assay factors that might affect assay
1360 performance, such as buffer pH, incubation temperature, and incubation time. Historical
1361 experience with one or several of the bioassay steps, along with sound scientific
1362 judgment, can identify key factors that require further evaluation. One tool that may be
1363 used to prioritize factors is a failure mode and effects analysis. Factors are typically
1364 scored by the combination of their potential severity and the likelihood that they will
1365 occur. The laboratory must be careful to recognize potential interactions between assay
1366 parameters.

1367

1368 **Screening.** Once potential key factors have been identified from process mapping and
1369 risk analysis, the laboratory should conduct an initial screening experiment to probe for
1370 effects that may require control. Screening designs such as factorial and fractional
1371 factorial designs are commonly used for this purpose. Software is available to assist the
1372 practitioner in the selection of the design and in subsequent analysis. Analysts should
1373 take care, however, to understand their assumptions about design selection and
1374 analysis to ensure accurate identification of experimental factors. Care should also be

1375 taken in defining design variables, which are measurements made from design runs and
1376 should be directly or indirectly related to bioassay performance.
1377

1378 **Response optimization.** A screening design will usually detect a few important factors
1379 from among those studied. Such factors can be further studied in a response-
1380 optimization design. Response-optimization designs such as central composite designs
1381 are performed to determine optimal settings for combinations of bioassay factors for
1382 achieving desired response. The information obtained from response optimization is
1383 typically depicted as a response surface and can be used to establish ranges that yield
1384 acceptable assay performance. With the advent of Quality by Design, that region is
1385 called the design space for the bioassay. This is also the set of ranges that will be
1386 incorporated into the bioassay SOP.

1387
1388 **Confirmation.** The mathematical model depicting assay performance as a function of
1389 changes in key assay factors is an approximation. Thus it is customary to confirm
1390 performance at the ideal settings of the bioassay. Confirmation can take the form of a
1391 qualification trial in which the assay is performed, preferably multiple independent times
1392 using optimal values for factors. Or the laboratory may determine that the bioassay has
1393 been adequately developed and may move to validation. Qualification is a good
1394 practice, not a regulatory requirement. The decision to perform confirmatory qualifying
1395 runs or to engage in validation is a business decision and depends upon the strength of
1396 the accumulated information obtained throughout development.

1397

1398 **5.3 Data Analysis during Assay Development**

1399

1400 Analysis of bioassay data during assay development enables analysts to make
1401 decisions regarding the statistical model that will be used for routine analysis, including
1402 transformation and/or weighting of data. The analysis also provides information
1403 regarding which elements of design structure should be used during outlier detection
1404 and the fitting of a full model. This may also include a plan for choosing subsets of data,
1405 such as a linear portion, for analysis or, for nonlinear bioassays, a model reduction
1406 strategy for samples similar to Standard. Once these decisions are made and checked
1407 during validation, they are not revisited for each bioassay analysis. These goals may be
1408 pursued in a stepwise manner, as follows:

1409

1410 (1) Choose an appropriate statistical model (see Section 4.6). Many
1411 considerations are involved in making this determination. First, the model should
1412 be appropriate for the type of assay endpoint—continuous, a count, or
1413 present/absent. Second, the model should incorporate the structure of the assay
1414 design. For any design other than completely randomized, there will be terms in
1415 the model for the structural elements. These could be, for example, within-plate
1416 blocking, location of cage in the animal facility, day, etc. A third consideration is
1417 applicable to continuous endpoints and involves whether to use a regression
1418 model or a means model (an analysis of variance model that fits a separate
1419 mean at each dilution level of each sample tested), with appropriate error terms.

1420

1421 (2) Fit the chosen statistical model to the data without the assumption of
1422 parallelism, and then assess the distribution of the residuals (differences
1423 between the data and the model), specifically examining them for departures
1424 from normality and constant variance. Transform the data as necessary or, if
1425 needed, choose a weighting scheme (Section 4.3). Use as large a body of assay
1426 data as possible for this step—at least several, and preferably dozens, of
1427 independent assays. The primary goal is to address any departure from constant
1428 variance of responses across the range of concentrations in the assay. Step 2
1429 will likely alternate between imposing a transformation and assessing the
1430 distribution of the residuals.

1431
1432 (3) Screen for outliers. This step normally follows the initial choice of a
1433 suitable transformation and/or weighting method. Outlier analysis should be
1434 conducted in a predefined and consistent manner and is best done on the
1435 residuals from a model that is fit to all the data within an assay. Ideally the model
1436 used for outlier detection contains the important elements of the assay design
1437 structure, allows nonsimilar curves, and makes fewer assumptions about the
1438 functional shape of the concentration–response curve than did the model used to
1439 assess similarity. Outlier analysis on a small number of replicates or
1440 pseudoreplicates is inadvisable. See Section 4.8 and General Chapter <1010>.
1441 In some cases outliers may be so severe that a reasonable model cannot be fit,
1442 and thus residuals will not be available. In such cases, it is necessary to screen
1443 the raw data for outliers before attempting to fit the model.

1444
1445 (4) Remove outliers as appropriate. See Section 4.8 for a discussion of
1446 methods. Routinely, before an observation is declared an outlier, an investigation
1447 should be conducted to determine whether a cause can be identified.

1448
1449 During development, prepare instructions for the investigation and treatment of
1450 an outlier observation, including any limits on how many outliers are acceptable.
1451 Include these instructions in the assay procedure. Good practice includes
1452 recording the result of this investigation, outlier test(s) applied, and results from
1453 the latter. Note that outlier procedures must be considered apart from the
1454 investigation and treatment of an out-of-specification (OOS) result (reportable
1455 value). Decisions to remove an outlier from data analysis should not be made on
1456 the basis of how the reportable value will be affected in terms of a potential OOS
1457 result.

1458
1459 Removing data as outliers should be rare. If many values from a run are removed
1460 as outliers, that run should be considered suspect.

1461
1462 (5) After outlier removal, refit the model with the transformation and/or
1463 weighting previously imposed (Step 2) and re-assess the appropriateness of the
1464 model.

1465

1466 (6) If necessary or desired, choose a scheme for identifying subsets of data to
1467 use for potency estimation, whether the model is linear or nonlinear (see 4.5).
1468

1469 **5.4 Bioassay Validation**

1470
1471 The bioassay validation is a protocol-driven study that demonstrates that the procedure
1472 is fit for use. The validation should not be carried out until there are no expectations of
1473 further changes in the procedure and an assay SOP has been developed. Preliminary
1474 system and sample suitability controls should be established and clearly written into the
1475 assay SOP. These may be finalized based on additional experience gained from the
1476 validation exercise. USP Chapter *Biological Assay Validation* <1033> provides details
1477 regarding the design and analysis of the validation.

1478 **5.5 Bioassay Maintenance**

1479
1480
1481 The development and validation of a bioassay, though discrete operations, lead to
1482 ongoing activities. Assay improvements may be implemented as technologies change
1483 and as the laboratory becomes more skilled with the procedure, and changes to
1484 bioassay methodology require re-evaluation of bioassay performance. Some of these
1485 changes may be responses to unexpected performance during routine processing.
1486 Corrective action should be monitored using routine control procedures. Substantial
1487 changes may require a study verifying that the bioassay remains fit for use. An
1488 equivalence testing approach can be used to show that the change has resulted in
1489 acceptable performance. A statistically oriented study can be performed to demonstrate
1490 that the change does not compromise the previously acceptable performance
1491 characteristics of the assay.

1492
1493 **Assay transfer.** Assay transfer assumes both a known intended use of the bioassay in
1494 the recipient lab and the associated required capability for the assay system. These
1495 implicitly, though perhaps not precisely, inform limits on the amount of bias and
1496 precision degradation allowed between labs. Consider, for example, an assay with an
1497 established potency range from 0.5 to 2.0 with validation data demonstrating bias less
1498 than 5% and intermediate precision from a single assay of 9% geometric relative
1499 standard deviation. If this assay is used to support a product with a potency
1500 specification of 0.71 to 1.41, replicate assays will be required on each lot to achieve a
1501 process capability of 1.3. Using two laboratories interchangeably to support one product
1502 will require considering variation between labs in addition to intermediate precision for
1503 sample size requirements to determine process capability. With only two laboratories it
1504 is not sensible to estimate a variance, and so an estimate of reproducibility is not
1505 available. It is reasonable to determine a 90% confidence interval on the bias between
1506 labs and to add the most extreme endpoint of this bias interval to the assay bias in the
1507 process capability calculation. Working through these calculations for various
1508 magnitudes of assumed bias will demonstrate the impact of various assay transfer bias
1509 acceptance limits on the future replication requirements of the assay system. For a
1510 discussion and example pertaining to the interrelationship of bias, process capability,

1511 and validation, see A Bioassay Validation Example in *Biological Assay Validation*
1512 <1033>.

1513

1514 **Improving or updating a bioassay system.** If a new version of a bioassay improves
1515 the bias, precision, range, robustness, or specificity, the new version may lower the
1516 operating costs or offer other compelling advantages. When improving or updating a
1517 bioassay system analysts can use a bridging study to compare the performance of the
1518 new vs established versions. A wide variety of samples (e.g., lot release, stability,
1519 stressed, critical isoforms) should be used for demonstrating equivalence of estimated
1520 potencies. Compare the performance of the new and existing assays using the method
1521 described above under assay transfer. If the assay systems are very different (e.g., an
1522 animal bioassay vs a cell-based bioassay) but use the same Standard and mechanism
1523 of action, one can reasonably expect comparable potencies. If the new assay uses a
1524 different Standard, the minimum requirement for a comparison is a unit slope of the log
1525 linear relationship between the estimated potencies. An important implication of this
1526 recommendation is that poor precision or biased assays used early can have lasting
1527 impact on the replication requirements, even if the assay is later replaced by an
1528 improved assay.