

## SPECIAL ARTICLE

# Making Sense of Noninferiority: A Clinical and Statistical Perspective on Its Application to Cardiovascular Clinical Trials

Sanjay Kaul and George A. Diamond

Active control noninferiority trials are being used with increasing frequency in new drug or device development when standard placebo-controlled trials are considered unethical. Nevertheless, the design and analysis of these trials are founded on a number of assumptions and arbitrary criteria that are generally not well understood or justifiable. Trials designed to show noninferiority require an appropriate reference population, a proven active control and dose, an appropriate margin of noninferiority that is clinically relevant and statistically justifiable, a high level of adherence to treatment, and adequate statistical power to reliably conclude that a treatment is truly noninferior and therefore effective. Accordingly, if noninferiority trials are to be applied to clinical and regulatory decisions regarding the marketing and use of new treatments, the assumptions must be made explicit and their influence on the resultant conclusions must be assessed rigorously. When conservative criteria were applied to each of the key assumptions underlying 2 representative noninferiority trials, they materially undermined the conclusions regarding noninferiority failing to confirm reported conclusions regarding noninferiority despite enthusiastic dissemination and acceptance of the results. Because the clinical, regulatory, and economic impact of

active control noninferiority trials is substantial, robust criteria should be used routinely in their design, analysis, and interpretation to reach their intended objectives and to keep them from becoming wasted efforts.

© 2007 Elsevier Inc. All rights reserved.

Noninferiority—showing a treatment is good because it is not bad.<sup>1</sup>

WC Blackwelder

**I**n contrast to a conventional clinical trial, which is usually designed to show that a new treatment is superior to placebo, an active control noninferiority trial is designed to show that the new treatment is not inferior to standard treatment by some clinically acceptable difference. If noninferiority is thereby established, the utility of the new treatment can be based on ancillary advantages in safety, convenience, or cost. These noninferiority trials are used when placebo-controlled trials are considered unethical.

As discussed by several authors,<sup>2-11</sup> there are issues with these trials that make them less credible than superiority trials. The formal analysis of such trials is founded on a number of assumptions that cannot be justified or verified. These include reliance upon external information (historical placebo-controlled trials), arbitrary thresholds to characterize the degree of clinically important difference, the so-called assay sensitivity (the ability to discriminate effective from ineffective therapies) and constancy (the applicability of the historical data to the current trial). If such trials are to be relied upon for clinical and regulatory decision making, these assumptions must

---

From the Division of Cardiology, Cedars-Sinai Medical Center, Los Angeles, CA, and David Geffen School of Medicine, University of California, Los Angeles, CA.

Address reprint requests to Sanjay Kaul, MD, Division of Cardiology, Cedars-Sinai Medical Center, 8700 Beverly Blvd, Los Angeles, CA 90048.

E-mail: kaul@cshs.org

0033-0620/\$ - see front matter

© 2007 Elsevier Inc. All rights reserved.

doi:10.1016/j.pcad.2006.10.001

**Table 1. Primary ITT Results of REPLACE-2 Trial**

	Bivalirudin (n = 2975)	GPI + Heparin (n = 2991)	Absolute Difference (95% CI)	OR (95% CI)	<i>P</i>
Death, MI, urgent TVR or major bleeding	275 (9.2%)	299 (10.0%)	−0.75 (−2.25 to 0.74)	0.92 (0.77 to 1.09)	.32
Death, MI, or urgent TVR	227 (7.6%)	211 (7.1%)	0.57 (−0.75 to 1.90)	1.09 (0.90 to 1.32)	.40
Death	7 (0.2%)	12 (0.4%)	−0.17 (0.12 to −0.45)	0.59 (0.23 to 1.49)	.26
MI	207 (7.0%)	185 (6.2%)	0.77 (2.03 to −0.40)	1.13 (0.92 to 1.39)	.23
Q-wave MI	12 (0.4%)	13 (0.4%)	−0.03 (−0.36 to 0.30)	0.93 (0.42 to 2.04)	.43
Non-Q-wave MI	195 (6.6%)	172 (5.8%)	0.80 (−0.42 to 2.02)	1.15 (0.93 to 1.42)	.43
Urgent TVR	35 (1.2%)	42 (1.4%)	−0.19 (−0.77 to 0.38)	0.86 (0.55 to 1.35)	.44
Major bleeding	71 (2.4%)	123 (4.1%)	−1.72 (−2.61 to −0.82)	0.57 (0.42 to 0.77)	<.001
TIMI major bleeding	19 (0.6%)	26 (0.9%)	−0.23 (−0.67 to 0.21)	0.73 (0.40 to 1.33)	.30

ITT, intention to treat; OR, indicates odds ratio. *P* values for superiority are shown.

be made explicit, their basis must be sufficiently justified, and their influence on the resultant decisions must be assessed rigorously and expressed unambiguously in the published reports.

In this essay, we will describe the key issues underlying the design and analysis of these trials using 2 representative cardiovascular clinical trials (REPLACE-2 and ACUITY) as exemplars and explore the robustness of the investigators' conclusions with respect to these issues. Our explicit goals are to (i) review the key issues related to the design of noninferiority trials, (ii) review the statistical approaches and illustrate the degree to which the various assumptions in the statistical methodology influence the noninferiority conclusions, and (iii) suggest practical standards for reporting of these trials that improve the accuracy of their interpretation by clinicians and regulators alike.

### A Typical Noninferiority Trial

The REPLACE-2 trial is presented as representative of the typical noninferiority trial. This trial was a prospective randomized double-blind trial comparing bivalirudin, a direct thrombin inhibitor, plus provisional platelet glycoprotein (GP) IIb/IIIa inhibitor (new treatment) with unfractionated heparin plus planned GP inhibitor (standard treatment) during elective or urgent percutaneous coronary intervention (PCI) that was unrelated to acute myocardial infarction (MI) or acute coronary syndrome (ACS).<sup>12</sup> A placebo-controlled randomized trial of the new treatment could not be justified because withholding the established standard treatment would be deemed unethical.

In designing the trial, the investigators determined that a sample of 3000 patients was required for each treatment group to detect a 12.5% relative risk difference from a baseline quadruple end point incidence of 8% with a 2-sided  $\alpha$  level of 5.0% and power of 92% to establish noninferiority. The primary evaluation was based on a quadruple composite end point composed of 3 efficacy end points (death, MI, urgent revascularization) and 1 safety end point (in-hospital major bleeding). The secondary evaluation was based on the efficacy criterion alone (triple end point). Results are summarized in Table 1. Accordingly, the REPLACE-2 investigators concluded that “bivalirudin with provisional GP IIb-IIIa blockade was statistically not inferior to heparin plus planned GP IIb-IIIa blockade in terms of suppression of acute ischemic end points and bivalirudin was associated with less bleeding.” But is this conclusion sufficiently justified?

### Approaches to Noninferiority Analysis

There are 2 basic approaches to noninferiority analysis. The first approach seeks to determine whether the new treatment is inferior to the standard treatment by no more than some predefined margin (“fixed margin” analysis).<sup>4-11</sup> The second approach seeks to demonstrate indirectly whether the new treatment would be superior to placebo, had a placebo arm been used in the trial (“putative placebo” analysis). The putative placebo approach can also be used to determine whether the new treatment retains some predefined fraction, *f*, of the standard

treatment's effect ("fraction preservation" analysis).<sup>4-11</sup> Whereas the former addresses "relative efficacy" between the new and the standard treatment, the latter is required to establish "absolute efficacy" of the new treatment, that is, superiority over placebo.

We will explore the results of these approaches on the analysis of REPLACE-2 and ACUITY trials. Our review of this methodology is intended to be pedagogical and not exhaustive. Those seeking more detailed discussions are referred to a variety of technical sources.<sup>4-9,13-18</sup> Because the use of a combined efficacy and safety outcome as the primary evaluation criterion is unusual in noninferiority assessment from a regulatory perspective (separate assessments of efficacy and safety being the norm), we will focus our analysis on their secondary outcome (efficacy alone) and will compare the conclusions with those based on their primary outcome (efficacy + safety).

## Design

### *Estimation of the margin*

The critical step in noninferiority design is the selection of the marginal difference ( $d$ ) upon which the noninferiority judgment is to be based. The margin quantifies the degradation in efficacy that is clinically acceptable considering the ancillary advantages of the new treatment. The International Conference on Harmonization guidelines emphasize that "the margin should be specified a priori, based on both clinical judgment and statistical reasoning and should be suitably conservative to reflect the uncertainty in evidence."<sup>13</sup>

### *Clinical judgment*

The margin is the maximum clinically acceptable difference that one is willing to give up in return for the secondary benefits of the new therapy. In general, noninferiority margin should be smaller than the minimal clinically important difference used in routine sample size estimations for superiority trials (typically 15%-25% proportional difference).<sup>13,14</sup> However, the choice of the margin varies on a case-by-case basis often being influenced by the seriousness of the clinical outcome (narrow margin for mortality or irreversible morbidity), the magnitude of standard

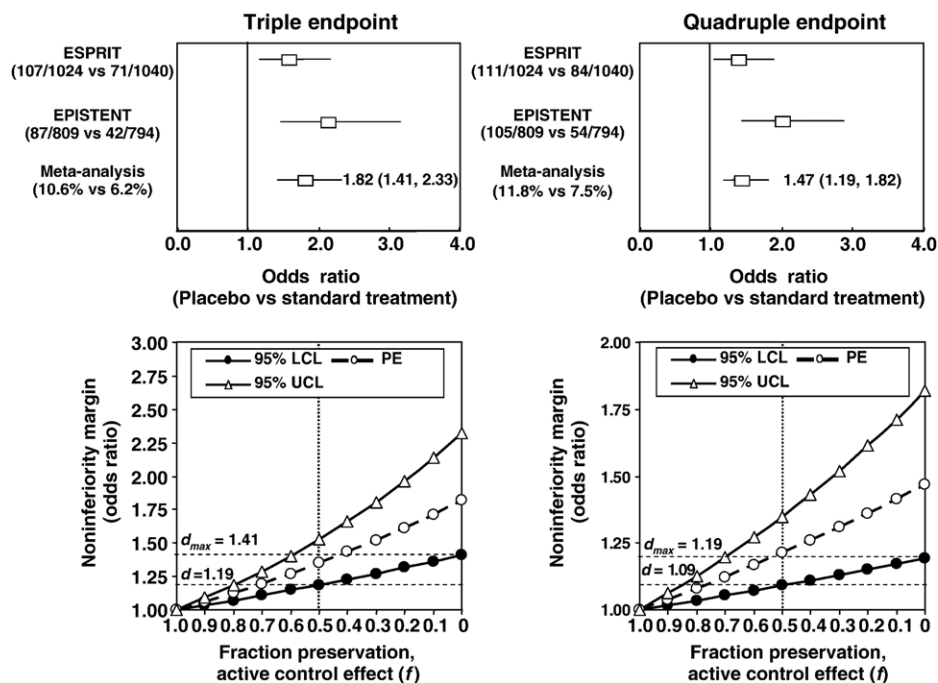
treatment effect (narrow margin for large treatment effect and vice versa), and the overall benefit-risk-cost profile. Given the subjective and somewhat arbitrary nature of these clinical judgments, it is difficult to obtain a consensus margin among the principal stakeholders (investigator, sponsor, or regulator).

### *Statistical reasoning*

The statistical determination of noninferiority margin is based on obtaining a precise estimate of the therapeutic effect of the active control (standard treatment). This is accomplished via a formal meta-analysis of previously available placebo-controlled studies using a fixed or random effect estimator to account for variability in the data (the latter being preferred as it accounts for variance between and within trials).<sup>5-11</sup> Although no formal guidelines exist for recommendation of an appropriate choice of margin, there is general agreement that it should be smaller than the active control effect, preferably a fraction ( $f$ ) of the active control effect, to gain the assurance that it is "suitably conservative."<sup>13,14</sup> Margins expressed on a relative difference scale (odds, risk, or hazard ratio) are preferred over absolute difference to "fix" the margin in case of unanticipated dissimilarities in observed and expected event rates.<sup>4,10,11</sup>

Thus, setting an appropriate margin that is not unduly wide (liberal) or restrictive (stringent) is critical in the design of noninferiority trials. Because the choice of margin has a critical impact on determination of sample size which varies inversely as the square root of the margin, some have argued that this decision should primarily be statistical. One reasonable compromise between the clinical relevance and statistical justification basis for the margin, the so-called minimum  $\Delta$  approach, is to estimate the margin ( $\Delta$ ) using both clinical judgment and statistical reasoning and then to choose the smaller of the 2 values.<sup>5</sup>

Fig 1 illustrates the estimation of the noninferiority margin in REPLACE-2 trial using the methods described above. The active control effect is derived from a random-effect meta-analysis of 2 previous placebo-controlled studies in the modern interventional era, one comparing abciximab (EPISTENT)<sup>19</sup> and the other eptifibatide (ESPRIT)<sup>20</sup> plus heparin



**Fig 1.** Estimation of noninferiority margin. The active control effect for both the triple (left upper panel) and quadruple end point (right upper panel) is derived from analysis of 2 historical trials (ESPRIT, EPISTENT) and expressed as odds ratio of placebo relative to control. Summary effects are estimated using the DerSimonian/Laird method for random-effect meta-analysis. The maximum noninferiority margin ( $d_{max}$ ) is derived as the 95% LCL of the meta-analytic estimate. The more restrictive margin ( $d$ ) is based on a 0.5 fraction preservation of active control. See text for details.

(standard treatment) vs heparin (placebo) alone. Summary effects of the standard treatment over placebo odds ratio ( $O_S/O_P$ ) were estimated as 0.55 (95% confidence interval [CI], 0.43-0.71) for the triple and 0.68 (95% CI, 0.55-0.84) for the quadruple end point. The corresponding placebo over standard treatment odds ratio,  $O_P/O_S$  (derived by inverting  $O_S/O_P$ ), is shown in the figure. The maximum noninferiority margin ( $d_{max}$ ) was derived as the 95% lower confidence limit (LCL) of  $O_P/O_S$ . A more restrictive margin ( $d$ ) was derived mathematically as  $d_{max}^{(1-f)}$ , where  $f$  represents the fraction of the standard treatment preserved by the new treatment (typically 0.50 but varies from 0 to 1.0 according to the benefit-risk-cost profile of the new treatment). Thus, for an  $f = 0.5$ ,  $d$  was estimated as an odds ratio of 1.19 ( $1.41^{0.5}$ ) for the triple end point and 1.09 ( $1.19^{0.5}$ ) for the quadruple end point. The corresponding risk ratio margins are 1.32 ( $d_{max}$ ) and 1.15 ( $d$ ) for the triple and 1.12 ( $d_{max}$ ) and 1.06 ( $d$ ) for the quadruple endpoint. Margins based on point estimate of  $O_P/O_S$  are generally

not recommended because they do not account for the variance in standard treatment effect, and those based on 95% upper CI are deemed to be too liberal (at best) and not valid (at worst). Margins based on point estimate may be justifiable (a) if the standard treatment effect has been reliably and repeatedly estimated in multiple historical placebo-controlled trials (minimal variance); and (b) if the value of  $f$  is chosen to be greater than 0.5 but less than 1.0 (typically 0.8). Hence, the greater the active control effect to be preserved, the smaller the margin, and the more robust the noninferiority inference.

#### Estimation of sample size

Compared to placebo-controlled trials, noninferiority trials typically have larger sample sizes because the margin is much smaller than the treatment difference ( $\Delta$ ) for which a placebo-controlled trial is powered. In addition, the sample size of a noninferiority trial is very sensitive to the expected effect of the new

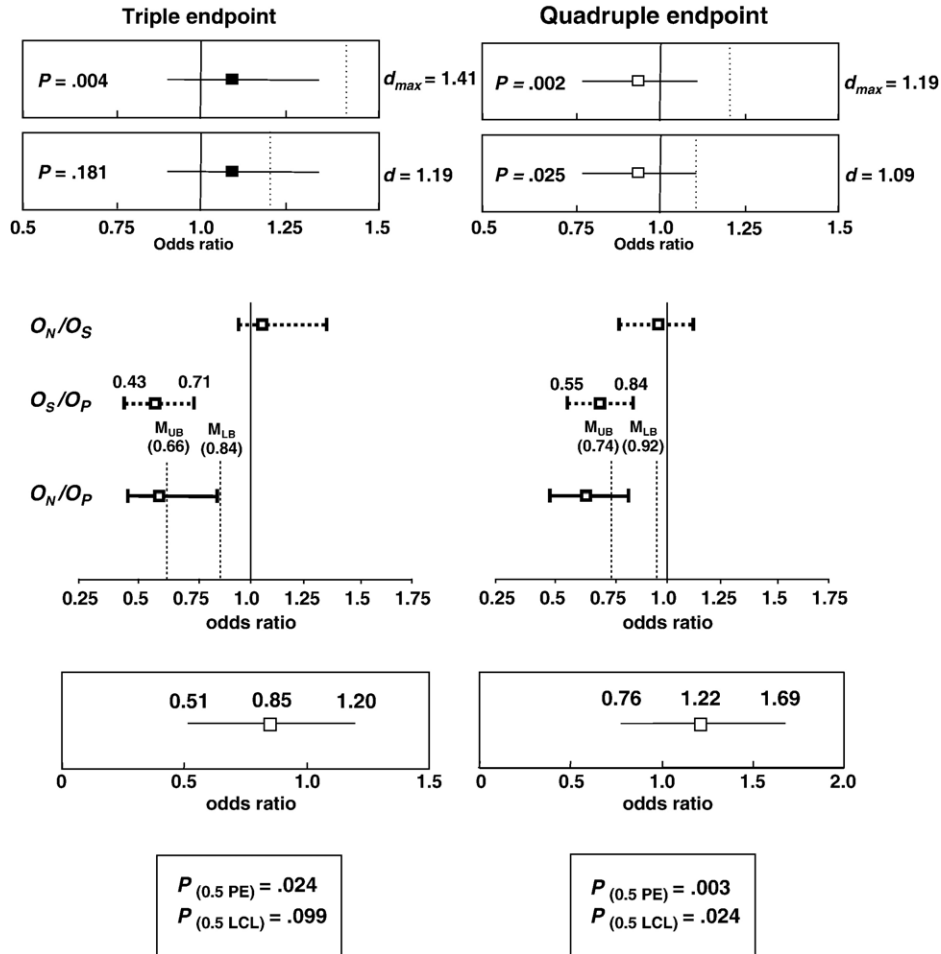


Fig 2. Top panel, Determination of therapeutic noninferiority. Indirect confidence interval comparison approach is illustrated in which noninferiority is established when the upper bound of the 1-sided 97.5% CI (corresponding to 2-sided 95% CI) lies within the noninferiority margin. The analysis is shown for a maximum ( $d_{max}$ ) and restrictive margin ( $d$ ). Results of hypothesis testing are shown as  $P$  values with 1-sided  $P \leq .025$  (corresponding to 1-sided 97.5% CI) as criterion for noninferiority. Middle panel, Determination of therapeutic efficacy over a putative placebo. The putative placebo approach is shown in the middle panel. The effect of the new treatment vs placebo ( $O_N/O_P$ ) is synthesized from the effect of the new vs the standard treatment ( $O_N/O_S$ ) observed in the current trial and the standard treatment vs placebo ( $O_S/O_P$ ) in the historical trials. Superiority over putative placebo is established if the derived  $O_N/O_P$  is less than 1.0. Noninferiority is established if the worst limit of  $O_N/O_P$  does not exceed the marginal threshold. Two estimates of marginal threshold are used: a liberal  $M_{LB}$  (50% of the lower bound of  $O_S/O_P$ ) and a stringent  $M_{UB}$  (50% of the upper bound of  $O_S/O_P$ ). Lower panel, Determination of fraction preservation of standard treatment effect. The fraction of standard treatment effect preserved by the new treatment is shown in the lower panels using the Hasselblad and Kong<sup>17</sup> method. Noninferiority is established if the 95% lower limit of the CI exceeds a target fractional threshold (typically 0.5). Noninferiority  $P$  values for 50% preservation are also shown based on the point estimate (PE) or LCL criteria.

treatment relative to the active control. For example, using a baseline event rate of 6.2% for the standard treatment group (triple end point) as observed in the historical trials, a noninferiority margin of 18% (equivalent to a risk ratio of 1.18 or an absolute risk difference of 1.08%, representing 50% preservation of the standard treatment's

worst effect), a 1-sided  $\alpha$  of 2.5%, and a  $\beta$  of 10% (90% power), a sample size of nearly 10 500 patients per treatment group would be required for noninferiority assessment in REPLACE-2 on the assumption that the 2 treatments were equivalent. However, based on the assumption that the new treatment was superior to the

standard by 0.8%, the sample size would be nearly 3500 per group. Thus, the sample size can be considerably larger if the 2 treatments are assumed to be equivalent than if the new treatment is assumed to be slightly more effective than the active control. By contrast, a placebo-controlled trial intended to demonstrate a 25% reduction (from 6.2% to 4.6%) with 90% power would require approximately 4200 patients per treatment group. It is important to note that both noninferiority and superiority can be assessed in the same clinical trial without any need for statistical adjustment.<sup>5-7</sup> Thus, superiority can be inferred in the setting of a noninferiority trial even without specifying a superiority hypothesis a priori. However, the converse, that is, adding a noninferiority hypothesis after a superiority trial is completed, would be subject to unacceptable bias owing to post hoc “gerrymandering” of the noninferiority margin.<sup>13,14</sup>

## Analysis and Interpretation

### Fixed margin analysis

A frequently used method of analyzing noninferiority is to estimate the effect of the new treatment relative to the standard and then to compare the results with that margin. If the CI (1-sided 97.5% CI or 2-sided 95% CI) of this estimate is entirely below this margin, then the new treatment is declared “noninferior.”<sup>4-11</sup> This is often referred to as the indirect CI comparison (ICIC) method. A hypothesis-testing approach can also be applied in which the null hypothesis of “inequality” (risk difference > margin) can be rejected in favor of the alternative hypothesis of “equality” (risk difference < margin) if the 1-sided  $P \leq .025$ .<sup>5-9,15</sup>

The results of fixed margin analysis of the REPLACE-2 trial using the ICIC and hypothesis testing methods are shown in Fig 2 (top panel). Noninferiority is established for the triple and quadruple end points using the maximum margin ( $d_{\max}$ ). When the more restrictive criterion ( $d$ ) is used, noninferiority is barely ( $P = .025$ ) inferred for the quadruple end point but not established for the triple end point.

### Putative placebo analysis

For an experimental therapy to be declared noninferior, it is necessary to demonstrate with statis-

tical certainty that it is also superior to placebo to support a claim of “therapeutic efficacy” required for regulatory approval.<sup>4-11,16-18</sup> Although active control trials lack a placebo arm by design, one can deduce or “impute” the effect of the new treatment relative to placebo by the so-called putative placebo approach wherein the odds ratio for the new treatment vs placebo ( $O_N/O_P$ ) is derived from the odds ratio for the new treatment vs the standard treatment observed in the current noninferiority trial ( $O_N/O_S$ ) multiplied by the odds ratio for the standard treatment vs placebo based on the historical trials ( $O_S/O_P$ )<sup>10,11,16-18</sup>:

$$O_N/O_P = O_N/O_S \times O_S/O_P$$

Therapeutic efficacy is established if the new treatment is adjudged superior to the putative placebo ( $O_N/O_P, <1.0$ ). By incorporating the historical placebo comparison, the putative placebo approach fulfills the requirement of *assay sensitivity* or internal validity. However, it relies on an arbitrary assumption that the active control effect size is the same in the current and the nonconcurrent historical studies, that is, the *constancy* assumption.

The putative placebo approach can also be used for assessment of noninferiority if the worst limit (lower bound) of CI for  $O_N/O_P$  does not exceed a predefined threshold. This threshold is derived as a fraction (typically  $f = 0.5$ ) of the standard treatment ( $O_S/O_P$ ) that a new treatment is required to preserve. The best limit (upper bound) of this interval (equivalent to the worst limit of the  $O_P/O_S$  described in the fixed margin analysis) is generally recommended as the preferred margin for noninferiority assessment.<sup>16</sup> The REPLACE-2 investigators used an identical approach as their principal assessment for noninferiority but with one key distinction—instead of using the best limit, they used the worst limit (lower bound). As before, whereas the former can arguably result in a stringent marginal threshold, the latter is often viewed as too liberal with some even questioning its validity.<sup>8</sup>

The results of bivalirudin vs a putative placebo (heparin) in REPLACE-2 are shown for both the triple and quadruple end points in Fig 2 (middle panel). Superiority over heparin (derived  $O_N/O_P, <1.0$ ) is imputed for bivalirudin for both end points. Noninferiority is

**Table 2. Primary ITT Results of ACUITY Trial**

	Bivalirudin + GPI (n = 4604)	Bivalirudin (n = 4612)	GPI + Heparin/ Enoxaparin (n = 4603)	Absolute Difference (95% CI)	Risk Ratio (95% CI)	P-Value
Death, MI, urgent TVR, or major bleeding (%)	541 (11.8)	466 (10.1)	538 (11.7)			
Bivalirudin + GPI vs heparin/enoxaparin + GPI				0.06 (-1.25 to 1.38)	1.01 (0.90 to 1.12)	.93
Bivalirudin vs heparin/enoxaparin + GPI				-1.58 (-0.31 to -2.86)	0.86 (0.77 to 0.97)	.015
Death, MI, or urgent TVR (%)	356 (7.7)	360 (7.8)	334 (7.3)			
Bivalirudin + GPI vs heparin/enoxaparin + GPI				0.48 (-0.60 to 1.55)	1.07 (0.92 to 1.23)	.39
Bivalirudin vs heparin/enoxaparin + GPI				0.55 (-0.53 to 1.63)	1.08 (0.93 to 1.24)	.32
Major bleeding (%)	243 (5.3)	139 (3.0)	262 (5.7)			
Bivalirudin + GPI vs heparin/enoxaparin + GPI				-0.41 (0.52 to -1.34)	0.93 (0.78 to 1.10)	.38
Bivalirudin vs heparin/enoxaparin + GPI				-2.68 (-1.85 to -3.51)	0.53 (0.43 to 0.65)	<.001
Death (%)	70 (1.5)	74 (1.6)	62 (1.3)			
Bivalirudin + GPI vs heparin/enoxaparin + GPI				0.17 (-0.31 to 0.66)	1.13 (0.80 to 1.58)	.48
Bivalirudin vs heparin/enoxaparin + GPI				0.26 (-0.23 to 0.75)	1.19 (0.85 to 1.67)	.31
MI (%)	229 (5.0)	248 (5.4)	227 (4.9)			
Bivalirudin + GPI vs heparin/enoxaparin + GPI				0.04 (-0.84 to 0.93)	1.01 (0.84 to 1.21)	.93
Bivalirudin vs heparin/enoxaparin + GPI				0.45 (-0.46 to 1.35)	1.09 (0.92 to 1.30)	.33
Urgent TVR (%)	123 (2.7)	110 (2.4)	105 (2.3)			
Bivalirudin + GPI vs heparin/enoxaparin + GPI				0.39 (-0.24 to 1.03)	1.17 (0.91 to 1.51)	.23
Bivalirudin vs heparin/enoxaparin + GPI				0.10 (-0.51 to 0.72)	1.05 (0.80 to 1.36)	.74
TIMI major bleeding (%)	76 (1.7)	43 (0.9)	86 (1.9)			
Bivalirudin + GPI vs heparin/enoxaparin + GPI				-0.20 (-0.76 to 0.32)	0.88 (0.65 to 1.20)	.43
Bivalirudin vs heparin/enoxaparin + GPI				0.94 (-0.46 to -1.42)	0.50 (0.35 to 0.72)	<.001

GPI indicates glycoprotein IIb/IIIa inhibitor; heparin, unfractionated heparin; TIMI, Thrombolysis in Myocardial Infarction. Absolute difference = standard therapy (heparin/enoxaparin + GPI) minus new therapy (bivalirudin +/- GPI), negative values favor new therapy. *P* values for superiority are shown.

established for quadruple and triple end point using the liberal lower bound criterion ( $M_{LB}$ ) used in REPLACE-2. Had the investigators used the more stringent upper bound ( $M_{UB}$ ) criterion, however, noninferiority would not be established for either end point.

### Fraction preservation analysis

The putative placebo approach can also be used to estimate the fraction of the standard treatment effect preserved by the new treatment. This is determined as a ratio of the new vs standard treatment effect relative to the placebo vs standard treatment effect.<sup>9-11,16-18</sup> Two different estimates of the standard treatment effect are generally used,

one based on the point estimate (liberal criterion) and the other on the 95% LCL (stringent criterion). The point estimate is used in the method proposed by Hasselblad and Kong,<sup>17</sup> where

$$f = 1 + \{\log(O_P/O_N)/\log(O_S/O_P)\}$$

Noninferiority is inferred if the lower limit of the 95% 2-sided CI of this fraction exceeds a prespecified minimum threshold (arbitrarily 0.5, but could be higher for additional conservatism). Alternatively, noninferiority *P* value for testing 50% or higher preservation can be estimated and noninferiority concluded at a 1-sided  $P \leq .025$ .

The results of fractional preservation analysis for REPLACE-2 based on the Hasselblad and

**Table 3. Noninferiority Results of ACUITY Trial**

Trial	Observed RR, New vs Standard (95% CI)	Historical RR, Placebo vs Standard (95% CI)	Derived RR, New vs Placebo (95% CI)	Noninferiority Margin (RR)			Noninferiority Conclusion			50% Fraction Preservation ( <i>P</i> Value)	
				$d_{max}$	<i>d</i>	ACUITY Margin	$d_{max}$	<i>d</i>	ACUITY Margin	Point Estimate	Lower Confidence Limit
ACUITY, triple (bivalirudin vs heparin/ enoxaparin + GPI)	1.08 (0.93, 1.24)	1.72 (1.32, 2.27)	0.63 (0.46, 0.85)	1.32	1.15	1.25	Yes	No	Yes	.024	.082
ACUITY, triple (bivalirudin + GPI vs heparin/ enoxaparin + GPI)	1.07 (0.92, 1.23)	1.72 (1.32, 2.27)	0.62 (0.46, 0.85)	1.32	1.15	1.25	Yes	No	Yes	.020	.072
ACUITY, quadruple (bivalirudin vs heparin/ enoxaparin + GPI)	0.86 (0.77, 0.97)*	1.59 (1.12, 2.22)	0.54 (0.38, 0.78)	1.12	1.06	1.25	Yes	Yes	Yes	<.001	.005
ACUITY, quadruple (bivalirudin + GPI vs heparin/ enoxaparin + GPI)	1.01 (0.90, 1.12)	1.59 (1.12, 2.22)	0.64 (0.44, 0.91)	1.12	1.06	1.25	Yes	No	Yes	.015	.058

GPI indicates glycoprotein IIb/IIIa inhibitor; heparin unfractionated heparin. RR = risk ratio. Historical RR is based on placebo versus standard treatment effect derived from a random-effect meta-analysis of EPISTENT and ESPRIT trials. Noninferiority margins:  $d_{max}$  - 95% lower limit of historical RR; *d* - 50% of  $d_{max}$ ; ACUITY margin - 1.25 RR.

Criteria for noninferiority: 95% upper limit of new vs standard "observed" RR  $\leq$  noninferiority margin; *P* value for 50% fraction preservation  $\leq$  .025.

Criteria for efficacy: new vs placebo "derived" RR < 1.0.

Kong<sup>17</sup> method are shown in Fig 2 (bottom panel). The fraction of the standard treatment effect retained by the new treatment lies within a wide range from 0.51 to 1.20 (triple end point) and from 0.76 to 1.69 (quadruple end point). These data suggest that the new treatment is at least 51% and 76% as effective as the standard treatment effect. The corresponding noninferiority *P* values for 50% preservation of the standard treatment effect are significant for both the end points based on the point estimate criteria, and for the quadruple end point only using the 95% LCL criterion.

Thus, noninferiority is consistently demonstrated with all 3 analytical approaches using rather liberal (wide margin) and unconventional (efficacy plus safety end point) interpretive criteria in REPLACE-2. However, when noninferiority is limited to the conventional efficacy end point, it fails to be established using reasonably

conservative marginal criteria. Our reanalysis of REPLACE-2 challenges the investigators' interpretation of the trial data and supports the Food and Drug Administration's conclusion that "statistical noninferiority was not demonstrated for the triple [ischemic] end point."<sup>21</sup>

#### Noninferiority Assessment in ACUITY Trial

We applied the methods described above for the assessment of noninferiority in ACUITY trial, a recently published open-label trial in which 13819 patients with moderate- to high-risk ACS undergoing early invasive strategy were randomized to 1 of 3 arms: unfractionated heparin or enoxaparin (a low-molecular-weight heparin) plus planned GP inhibitor; bivalirudin plus planned GP inhibitor; or bivalirudin alone with provisional GP inhibitor.<sup>22</sup> The primary data for ACUITY are summarized in Table 2. Treatment

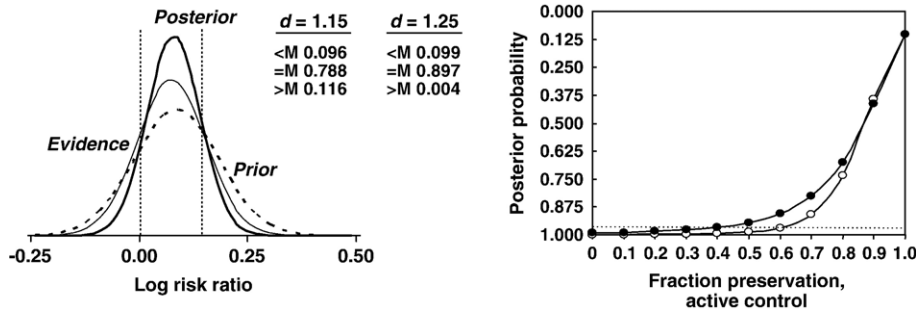
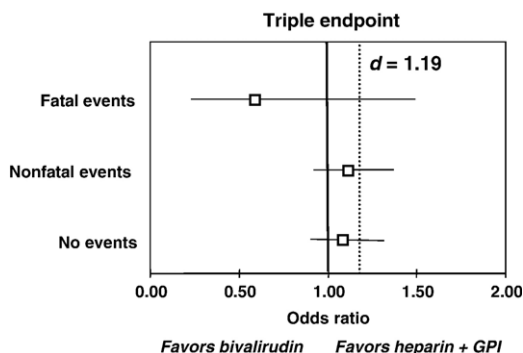


Fig 3. Left panel, Bayesian marginal analysis of ACUITY trial (left panel). Tri-plots showing posterior (thick line) distributions derived from integrating evidence or likelihood (thin line) from ACUITY trial and prior information (dashed line) from the REPLACE-2 trial, according to Bayes' theorem. The margin of noninferiority ( $M$ ) is indicated by the 2 vertical dotted lines and is equivalent to a log risk ratio of 0.14 (equivalent to a risk ratio of 1.18). Posterior probability of any effect size can be calculated by computing the area under the curve. The probabilities of falling below ( $<M$ ), within ( $=M$ ), or above  $M$  ( $>M$ ) are shown on the top right-hand corner. Probability of noninferiority is computed as the sum of probability of  $<M + =M$  (eg, the posterior probability of noninferiority is 0.096 [ $<M$ ] + 0.789 [ $=M$ ] = 0.884). Noninferiority is inferred at a posterior probability of 0.975 or higher (corresponding to a 1-sided  $P \leq .025$ ). Probability of noninferiority using the margin reported by the ACUITY investigators ( $d = 1.25$  risk ratio) is 0.996. Right panel, Bayesian fraction preservation analysis of ACUITY trial. The point estimate (open circle) and LCL (closed circle) of the historical standard treatment effect ( $O_S/O_P = 0.55$  [95% CI, 0.43-0.71]) are used to estimate the fraction preservation of standard treatment. The posterior observed odds ratio ( $O_N/O_S = 1.08$  [95% CI, 0.96-1.22]) is derived from integrating evidence or likelihood from ACUITY trial and prior information from REPLACE-2. Noninferiority is established at an  $f = 0.5$ . Superiority over putative placebo and standard treatment is established at  $f = 0$  and 1, respectively. Noninferiority is supported if posterior probability is greater than 0.975 (horizontal dashed line).

with bivalirudin alone was associated with a slight but nonsignificant 7% relative increase in ischemic events which was more than offset by a significant 47% relative decrease in major bleeding resulting in an overall significant 14% relative risk reduction in the quadruple end point compared to treatment with heparin or enoxaparin plus planned GP inhibitor. The bleeding advantage for bivalirudin alone treatment disappeared when combined with GP inhibitor.

The principal noninferiority results for the ACUITY trial using the ICIC, putative placebo, and fraction preservation analyses are shown in Table 3. The ACUITY investigators selected a noninferiority margin of 35% (1.25 risk ratio). Several features regarding the choice of margin deserve comment. First, the margin appears to be chosen entirely on the basis of expert consensus among the trial committee members without formal statistical justification.<sup>23</sup> Second, it is larger than the margins used in contemporary noninferiority acute coronary syndrome trials (10% in SYNERGY and 11% in A-to-Z trial).<sup>24,25</sup> Third, strictly from a statistical perspective, the choice of active control may not be justifiable because neither heparin plus GP inhibitor nor

enoxaparin plus GP inhibitor has been proven in placebo-controlled trials to reduce periprocedural ischemic events in the trial setting relative to heparin or enoxaparin alone, respectively. Demonstration of the active control to be a well-established effective standard therapy with predictable, quantifiable, and consistent treatment effects is a prerequisite for noninferiority assessment.<sup>2-9</sup> A meta-analysis of 6 ACS trials comparing heparin plus GP inhibitor with heparin reported negligible treatment effect on the composite of death, MI, or urgent reintervention at 30 days: odds ratio of 0.98 (95% CI, 0.93-1.02).<sup>26</sup> However, less than 15% of patients underwent early (<5 days) PCI in these studies, making these data ineligible (at worst) and questionable (at best) for estimation of noninferiority margin for ACUITY. Thus, these limitations call for exercising caution when evaluating noninferiority in ACUITY. In such cases where a noninferiority margin cannot be chosen properly because it has not been demonstrated under similar trial conditions that the active control would reliably show a treatment effect (rendering assay sensitivity untenable), a 3-arm trial design with a placebo arm or an active-control superiority trial should



**Fig 4. Information preserving composite end point analysis for the triple end point in REPLACE-2 trial.** The composite end point is classified as death as the worst outcome, no event as the best, and nonfatal events as intermediate. A margin of 1.19 odds ratio is used for assessment of noninferiority for each outcome.

be considered.<sup>2,3,11</sup> If noninferiority testing is still pursued, then the noninferiority margin must be very conservative so that the probability of making a false conclusion is extremely small.

Because nearly one third of patients enrolled in EPIDENT<sup>19</sup> and one fifth of patients enrolled in ESPRIT trials were patients with ACS undergoing PCI within 2 days,<sup>20</sup> we chose the REPLACE-2 margins derived from these trials for our analysis of ACUITY. We justified the choice of margin based on the assumption that the benefits of heparin plus GP inhibitor observed in the setting of nonurgent PCI would, at the very least, be preserved (if not enhanced) during urgent PCI for patients with moderate- to high-risk ACS. Furthermore, the enoxaparin and heparin therapy were considered interchangeable based on 2 recent large randomized trials reporting similar ischemic outcomes between enoxaparin (plus GP inhibitor) and heparin (plus GP inhibitor).<sup>24,25</sup>

The results of ICIC in Table 3 show that bivalirudin alone or bivalirudin plus GP inhibitor is noninferior to heparin/enoxaparin plus GP inhibitor for both end points according to the investigator predefined margin. For the analysis using REPLACE-2 margins, noninferiority is established for bivalirudin alone and bivalirudin plus GP inhibitor arm (both quadruple and triple end points) using the liberal margin  $d_{max}$ . However, using the more stringent margin  $d$ , noninferiority is established for 1 comparison only—bivalirudin alone with respect to the

quadruple end point. Superiority over putative heparin placebo (efficacy) is established for all 4 comparisons. Fraction preservation analysis is consistent with the marginal analysis and shows that noninferiority, using 50% preservation criterion, is established for all 4 comparisons using the liberal point estimate criterion and with only 1 comparison using the stringent LCL criterion (bivalirudin alone for quadruple end point).

Thus, like in REPLACE-2, noninferiority of bivalirudin (with or without planned GP inhibitor) is only established for the quadruple end point using a liberal margin. In contrast, noninferiority is not supported for the conventional efficacy end point using conservative criterion.

### Bayesian Analysis of Noninferiority

A Bayesian approach can be adapted to both fixed margin and putative placebo noninferiority analyses.<sup>10,27-29</sup> Briefly, normal posterior distributions are derived using the log mean risk or odds ratio ( $\mu$ ), and its standard deviation ( $\sigma$ )? according to Bayes' theorem which states that:

$$\text{posterior} \propto \text{likelihood} * \text{prior}$$

where “posterior” is the probability for the hypothesis (noninferiority) given the evidence; “likelihood” is the probability for the evidence given the hypothesis; and “prior” is the probability for the hypothesis independent of the evidence. The advantages of this approach, and its applications to noninferiority trials, are reviewed elsewhere.<sup>27-29</sup>

Fig 3 (left panel) describes a Bayesian fixed margin analysis where the posterior probability of noninferiority for bivalirudin alone is derived by integrating the prior information from the REPLACE-2 trial (informative prior) with the evidence from the ACUITY trial. The probability of noninferiority (for a margin  $d = 1.15$  RR) increased from 75% in REPLACE-2 and 82% in ACUITY to a “posterior” probability of 88% (still less than the threshold probability of 97.5% required for noninferiority inference) with regard to the triple end point. With regard to the quadruple end point (data not shown), the noninferiority probability remained high (96% in REPLACE-2, 100% in ACUITY and 100% posterior).

Fig 3 (right panel) describes the results of a putative placebo analysis using the Bayesian

**Table 4. Composite Score for Grading the Quality of Noninferiority Trials**

Trial	Therapeutic Noninferiority	Therapeutic Efficacy	Nonefficacy Benefit	Composite Score
REPLACE-2 (bivalirudin alone)	0	1	1	2
ACUITY (bivalirudin alone)	0	1	1	2
ACUITY (bivalirudin + GP inhibitor)	0	1	0	1
REPLACE-2 + ACUITY (bivalirudin alone)	0	1	1	2

Each attribute is graded on a 0 (unestablished) to 1 (established) scale.  
Criteria for:

1. Therapeutic noninferiority: 95% upper confidence limit of risk difference between new and standard treatment < margin (based on 50% of the 95% LCL of the risk difference between placebo and standard treatment).
2. Therapeutic efficacy: risk difference (ratio) between new treatment and imputed placebo <1.0.
3. Nonefficacy benefit: preferably superior or at least acceptable safety or tolerability with cost and/or convenience advantage of the new over the standard treatment.

approach described previously (Simon). The posterior probability of noninferiority (derived from evidence from ACUITY and prior information from REPLACE-2) is plotted as a function of the fraction preservation of active control. Based on this analysis, noninferiority (posterior probability >0.975 assessed at  $f = 0.5$ ) is supported for triple end point using point estimate but not the 95% LCL criterion. Although the probability of superiority of bivalirudin over imputed heparin (assessed at  $f = 0$ ) is 100%, the probability of superiority over heparin + GP inhibitor (assessed at  $f = 1$ ) is only 10%. In contrast, not only is bivalirudin alone demonstrated to be noninferior, but also superior to heparin + GP inhibitor using the quadruple end point (data not shown).

Because noninferiority depends on one's choice of the threshold fraction ( $f$ )—the lower the fraction to be preserved, the easier it is to establish noninferiority—a sensitivity analysis is best performed to define the robustness of this choice. Based on this analysis (Fig 3, right panel), noninferiority ( $P < .025$ ) would have been supported for an  $f = 0.43$  (LCL) and 0.63 (point estimate).

These results are consistent with the conventional frequentist analyses and fail to establish noninferiority of the new treatment relative to the standard, using rather stringent interpretive criteria.

### Interpretation of Composite End Points

The analysis of composite end points poses a particular challenge to the interpretation of clinical trials whether they are designed as

superiority or noninferiority trials. The construction of the composite end point is generally based on the premise that each component end point is interchangeable.<sup>30</sup> However, for this assumption to be valid, each component should be of equal or comparable clinical importance, occur with similar frequency, and be equally responsive to treatment intervention.<sup>31</sup> This is seldom fulfilled. Examination of the data in Table 1 shows this assumption is not supported in REPLACE-2. For example, the contribution of each component to the triple composite end point is less than 3% from the most robust component (death), 16% from the least robust end point (urgent target vessel revascularization [TVR]), and 81% from an intermediate end point (MI, mostly non-Q-MI). In addition, the treatment differences are virtually entirely attributable to differences in non-Q-MI with little or no impact on death or urgent TVR. This fact renders the REPLACE-2 analysis less than ideal at best, and at worst misleading, because important information regarding the individual component end points may be obscured by combining them into a composite.

Utility ranking or weighted schemes may offer a viable approach to combining end points of different value and provide a potential solution against over-interpretation of composite end point results. One can create a suitable ordered categorical end point to retain more information in the analysis.<sup>32</sup> This information-preserving composite end point (IPCE) approach would have death as the worst outcome, no event as the best (akin to event-free survival), and nonfatal events as intermediate. The results of such an "IPCE" analysis applied post hoc to the RE-

PLACE-2 trial are shown in Fig 4. The triple end point data show numerical trends in opposite directions of benefit with death favoring bivalirudin and nonfatal events and no events favoring the standard treatment. When the components show numerical trends in opposite directions, it makes the composite end point even more difficult to interpret. Of note, noninferiority is not established with respect to any of the 3 outcomes (upper limits of 95% CI crossing the margin).

The unconventional use of a composite efficacy and safety outcome biased the assessment of noninferiority in favor of bivalirudin in REPLACE-2. Typically, the noninferiority claim is confined to efficacy alone. Although combining efficacy and safety into one composite outcome might be desirable to inflate the event rate and enhance trial feasibility, it can often be misleading because drugs that are relatively ineffective but safer can be made to appear as good as or even better than effective drugs.<sup>33</sup> This is illustrated in REPLACE-2 and ACUITY where the difference in major bleeding events (43% and 47% risk reduction in favor of bivalirudin) exceeded the difference in MIs (13% and 9% risk increase), thereby resulting in the quadruple composite end point favoring bivalirudin.

### A Composite Score for Assessment of Noninferiority

Analysis of noninferiority should ideally be founded on 3 prerequisite judgments<sup>11</sup>—that the new treatment (i) exhibits “therapeutic noninferiority” (relative efficacy) to the standard treatment, (ii) would exhibit “therapeutic efficacy” in a placebo-controlled trial, and (iii) offers ancillary “nonefficacy benefits” in safety, tolerability, convenience, or cost. The new treatment should offer some or all of the nonefficacy benefits to justify its use in lieu of the standard treatment. A formal assessment of superiority with respect to such matters, although not currently a regulatory requirement, may nevertheless be desirable and preferably specified prospectively within the active-control trial design as a secondary objective.<sup>5-7</sup> At the very least, the new treatment should have acceptable safety and tolerability, and the evidence in support of these ancillary benefits should be specified explicitly in the published report. Bivalirudin offers

significant nonefficacy advantages over heparin plus GP inhibitor in terms of superior safety—significantly reduced major bleeding complications in both REPLACE-2 and ACUITY (Tables 1 and 2), cost advantage (\$395 vs \$615 for eptifibatide or \$1400 for abciximab), and convenience (1-hour vs 12- to 18-hour infusion).

**Table 5. Essentials of Noninferiority Assessment**

1. Ethical imperative:
  - (a) Placebo control cannot be used because effective standard treatment is available.
  - (b) New treatment should offer substantial benefits in safety, cost, or convenience over the standard treatment.
2. Choice of active control: best available comparator with large, reliable, and consistent treatment effect in placebo-controlled trials.
3. Noninferiority margin:
  - (a) Defined a priori based on clinical judgment and statistical reasoning.
  - (b) Relative risk difference scale (risk, odds, or hazard ratio) preferred over absolute risk difference.
4. Adequate power and sample size to minimize type II error (false negative).
5. Proper trial design and high quality of conduct:
  - (a) Identical patient population and protocol as in historical placebo-controlled trials
  - (b) Maximize protocol adherence.
6. Critical assumptions:
  - (a) Assay sensitivity (internal validity), assured if optimal choice for active control used in the current trial
  - (b) Constancy—active control effect is similar in current trial as in historical trials, assured by proper trial design and high quality of conduct.
7. Statistical analysis:
  - i. Fixed margin analysis
    - (a) Indirect CI comparison: upper limit of 2-sided 95% CI of treatment difference < margin
    - (b) Hypothesis testing:  $P \leq 0.025$  to reject the null hypothesis of inequality (risk difference  $\geq$  margin)
    - (c) Bayesian analysis: posterior probability of noninferiority  $\geq 0.975$
  - ii. Putative placebo analysis
    - (a) Superiority over imputed placebo: OR of new vs standard treatment < 1.0
    - (b) Fraction preservation of active control: at least 50% for noninferiority claim.
    - (c) Bayesian analysis: posterior probability of superiority over imputed placebo > 1.0 and 50% fraction preservation  $\geq 0.975$ .
8. Robust interpretive criteria for noninferiority
  - (a) Stringent marginal and fractional threshold and CI (2-sided 95% over 1-sided 95%)
  - (b) Stability of noninferiority inference for relative vs absolute outcomes, and for ITT vs per-protocol analysis
  - (c) Noninferiority claim for efficacy and superiority claim for safety/tolerability established in the same trial.

We hereby propose a composite score by which each of the 3 attributes of therapeutic noninferiority, therapeutic efficacy, and non-efficacy benefit is graded on a 0 (unestablished) to 1 (established) scale.<sup>34</sup> A score of 3 out of 3 thereby supports a judgment of so-called virtual superiority to justify consideration of the new over the standard treatment. As summarized in Table 4, virtual superiority is not established for either REPLACE-2 or ACUITY trial or for their combined analysis. It is also clear that the addition of GP inhibitor to bivalirudin eliminates the nonefficacy benefits associated with the latter.

## Discussion

Active control noninferiority trials are being used with increasing frequency in the cardiovascular arena. The interpretation of these trials poses a particular challenge to most clinicians. In this paper, we suggest practical standards for the analysis and reporting of these trials to improve the accuracy of their interpretation.

There are several key aspects of the non-inferiority inference (summarized in Table 5) that are critical for scientific credibility and regulatory acceptability. *First*, the ethical imperative requires a noninferiority assessment to be conducted only if a placebo-controlled trial (withholding an effective standard therapy) would be unjustifiable. An additional requirement that the new treatment offers some non-efficacy advantage over the standard treatment would be desirable to avoid the development of “me too” types of treatments.

*Second*, an optimal choice of active control is critical for the inference of efficacy of the new treatment. The active control should represent the best available comparator demonstrating large, reliable, and consistent treatment effect across a series of historical studies. This is crucial in establishing *assay sensitivity*.<sup>2-11</sup> Otherwise, the noninferiority assessment cannot be conducted with credibility, and therefore superiority of the new over the standard treatment should be the principal goal of the trial. Alternatively, a 3-armed trial that includes a placebo arm should be used wherever possible to establish direct evidence of efficacy of both new and standard treatment relative to placebo.

*Third*, for the historical comparison to have operational validity in the current trial, the critical assumption of *constancy* must be met.<sup>2-11</sup> A proper trial design and high quality of conduct (eg, maximal compliance, minimization of protocol deviations and outcome misclassifications, adherence to identical experimental protocol, etc) are crucial for this as well as assay sensitivity assumptions to hold. However, significant differences with respect to key design features might be unavoidable such as patient characteristics, concomitant disease-modifying medications, intensity of treatment, evolution of practice, etc, thereby invalidating these assumptions.<sup>4-8</sup> Because of this uncertainty, stringent interpretive criteria are required to raise the standard of evidence. In this regard, choosing the worst estimate of the active control or ensuring a fraction preservation of its effect may be considered to be forms of buffer or “discounting” to raise the strength of evidence.<sup>8-11</sup> The “50% rule” endorsed by the Food and Drug Administration (for thrombolytic trials) represents a form of “double discounting” in which preservation of 50% fraction is applied to the worst estimate of the active control effect to make it suitably conservative for both efficacy as well as noninferiority claim.<sup>8-11</sup> The noninferiority assessment is highly sensitive to the constancy assumption and is optimal when it holds (similar active control effect in the active control and historical trial populations). When the control effect is less in the active control relative to historical trials, setting a higher fraction preservation threshold (80% instead of 50%) might provide protection against inflation of type I (false-positive) error rate. For example, in REPLACE-2 trial, the observed active control event rate of 7.1% for triple end point was higher than the historical rate of 6.2%, thereby biasing the results toward noninferiority. Thus, a more conservative discounting (80% instead of 50% preservation) would have led the REPLACE-2 investigators away from an erroneous (false-positive) conclusion of noninferiority.

*Fourth*, the choice of the analytic strategy, that is, intention-to-treat (ITT) vs on-treatment or per-protocol, may have a substantial impact on the noninferiority inference.<sup>5-9</sup> The ITT is widely recognized as the most valid analytic approach for superiority trials because it is generally conservative and it preserves the advantages of

randomization. However, including data after drug discontinuation (as in ITT analysis) tends to favor a null difference biasing the results toward noninferiority. The per-protocol analysis excludes data from patients with major protocol violations and can potentially introduce substantial bias in either direction (mostly against noninferiority). In contrast to the superiority trial, which is conducted meticulously to increase the chances of detecting treatment differences, noninferiority trial by virtue of its nature tends to be sloppy in conduct to maximize the goal of detecting “no difference.” Therefore, a dual strategy incorporating both ITT and per-protocol approaches is recommended, and the inference strengthened only if both approaches support noninferiority. In REPLACE-2, the trial conduct was high with minimal protocol deviations or noncompliance, and noninferiority was therefore established for both ITT and per-protocol analyses (using rather liberal criteria). In ACUITY, neither the quality of trial conduct, nor the magnitude of protocol deviations and noncompliance was reported explicitly and noninferiority was established only for TT analysis.<sup>22</sup>

*Fifth*, proper sample size estimation is critical to detect a difference that might actually exist. We determined that a sample of nearly 10500 patients per treatment group would be required for noninferiority assessment (under the assumption that the 2 treatments are equivalent), three and a half times the number actually used in REPLACE-2 ( $n = 3000$  per group). Although the sample size estimated for ACUITY was appropriate for a margin of 1.25 risk ratio and an active control rate of 6.5% ( $n = 4800$  per group for an  $\alpha = .025$ ,  $\beta = .1$ ), using the REPLACE-2 margin (1.18 risk ratio) would require a sample of nearly 9300 per group. Thus, the relatively smaller sample size in REPLACE-2 and ACUITY underpowered the assessment of noninferiority with respect to stringent margins, thereby potentially inflating the type II (“false-negative”) error, that is, erroneous rejection of a truly noninferior treatment.

Although the choice of conservative margins can be justified by statistical reasoning that the “worst-case scenario” increases the robustness of the analysis, such narrow margins result in large sample sizes that often render the trials imprac-

tical as shown above.<sup>8-11,35</sup> Reconciling these offsetting considerations of trial feasibility and stringency poses a substantial challenge. To avoid these ambiguities, the statistician and the clinician together in consultation with the regulatory authorities should determine early on during the planning phase of investigation a consensus noninferiority limit that is clinically relevant and statistically feasible and that balances the contrasting perspectives of the principal stakeholders.

Recent systematic reviews of noninferiority trials revealed a variety of methodological flaws causing inflation of the type I or type II errors.<sup>36-39</sup> As a result of such errors, potentially inferior treatments might well become the active controls for future noninferiority trials with ever increasing frequency until the active control winds up being no better than a placebo, a phenomenon referred to as “biocreep” or “drift.”<sup>7</sup> Taken to the extreme, this could lead to the introduction of suboptimal treatments into routine clinical practice. The potential cost in dollars and lives is incalculable. In this paper, we have shown how attention to proper methods and the choice of operative thresholds, from liberal to clinically relevant to judiciously conservative (reflecting the core philosophies of the sponsor, practitioner, and regulator, respectively), help improve the interpretation of the data and help identify and reduce the potential for errors.

In summary, a number of inherent problems challenge the design, conduct, analysis, interpretation, and implementation of active control noninferiority trials. The design of these trials relies on various conventions (arbitrary marginal and fractional thresholds and historical controls) that would likely not be widely accepted as reasonable if it were not for the commercial implications associated with judgments of therapeutic noninferiority that derive from the trials.

Accordingly, if noninferiority trials are to be applied to clinical and regulatory decisions regarding the marketing and use of new treatments, the assumptions must be made explicit and their influence on the resultant conclusions must be assessed rigorously. Because the clinical, regulatory, and economic impact of such variation is substantial, conservative criteria should be used routinely in the design, analysis, and interpretation of noninferiority trials.

## References

1. Blackwelder WC: Showing a treatment is good because it is not bad: when DOES “noninferiority” imply effectiveness? *Control Clin Trials* 23:52-54, 2002
2. Temple R, Ellenberg SS: Placebo-controlled trials and active-control trials in the evaluation of new treatments: Part 1. Ethical and scientific issues. *Ann Intern Med* 133:455-463, 2000
3. Ellenberg SS, Temple R: Placebo-controlled trials and active-control trials in the evaluation of new treatments: Part 2. Practical issues and specific case. *Ann Intern Med* 133:464-470, 2000
4. Siegel JP: Equivalence and noninferiority trials. *Am Heart J* 139:S166-S170, 2000
5. Hung HMJ, Wang S-J, Tsong Y, et al: Some fundamental issues with noninferiority testing in active controlled trials. *Stat Med* 22:213-225, 2003
6. Hung HMJ, Wang S-J, O'Neill R: A regulatory perspective on choice of margin and statistical inference issue in noninferiority trials. *Biom J* 47:28-36, 2005
7. D'Agostino Sr RB, Massaro JM, Sullivan LM: Non-inferiority trials: design concepts and issues—the encounters of academic consultants in statistics. *Stat Med* 22:169-186, 2003
8. Snapinn SM: Alternatives for discounting in the analysis of noninferiority trials. *J Biopharm Stat* 14:263-273, 2004
9. Rothmann M, Li N, Chen G, et al: Design and analysis of non-inferiority mortality trials in oncology. *Stat Med* 22:239-264, 2003
10. Kaul S, Diamond GA, Weintraub WS: Trials and tribulations of noninferiority: the ximelagatran experience. *J Am Coll Cardiol* 46:1986-1995, 2005
11. Kaul S, Diamond GA: Good enough. A primer on the analysis and interpretation of noninferiority trials. *Ann Intern Med* 145:62-69, 2006
12. Lincoff AM, Bittl JA, Harrington RA, et al: Bivalirudin and provisional glycoprotein IIb/IIIa blockade compared with heparin and planned glycoprotein IIb/IIIa blockade during percutaneous coronary intervention: REPLACE-2 randomized trial. *JAMA* 289:853-863, 2003
13. International Conference on Harmonisation. Statistical principles for clinical trials (*ICH E 9*) (1998); International Conference on Harmonisation. Guidance on choice of control group and related design and conduct issues in clinical trials (*ICH E 10*) (2000). Food and Drug Administration, Department of Health and Human Services.
14. Committee for Proprietary Medicinal Products: Points to consider on the choice of noninferiority margin. The European Agency for the Choice of Medicinal Products. 2004 CPMP/EWP/2158/99 draft
15. Blackwelder W: Proving the null hypothesis in clinical trials. *Control Clin Trials* 3:345-353, 1982
16. Durrleman S, Chaikin P: The use of putative placebo in active control trials: two applications in a regulatory setting. *Stat Med* 22:941-952, 2003
17. Hasselblad V, Kong DF: Statistical methods for comparison to placebo in active-controlled trials. *Drug Inform J* 35:435-449, 2001
18. Holmgren EB: Establishing equivalence by showing that a specified percentage of the effect of the active control over placebo is maintained. *J Biopharm Stat* 9:651-659, 1999
19. The EPISTENT Investigators: Randomised placebo-controlled and balloon-angioplasty-controlled trial to assess safety of coronary stenting with use of platelet glycoprotein-IIb/IIIa blockade. *Lancet* 352:87-92, 1998
20. The ESPRIT Investigators: Novel dosing regimen of eptifibatide in planned coronary stent implantation (ESPRIT): a randomised, placebo-controlled trial. *Lancet* 356:2037-2044, 2000
21. Hughes S: FDA approves REPLACE-2 label for bivalirudin. *Heartwire news*, 2005 (<http://www.theheart.org>)
22. Stone GW, McLaurin BT, Cox DA, et al, for the ACUITY investigators: Bivalirudin for patients with acute coronary syndrome. *New Engl J Med* 355:2203-2216, 2006
23. Stone GW, Bertrand M, Colombo A, et al: Acute Catheterization and Urgent Intervention Triage strategy (ACUITY) trial: study design and rationale. *Am Heart J* 148:764-775, 2004
24. SYNERGY Trial Investigators: Enoxaparin vs unfractionated heparin in high-risk patients with non-ST-segment elevation acute coronary syndromes managed with an intended early invasive strategy: primary results of the synergy randomized trial. *JAMA* 292:45-54, 2004
25. Blazing MA, de Lemos JA, White HD, et al, for the A to Z Investigators: Safety and efficacy of enoxaparin vs unfractionated heparin in patients with non-ST-segment elevation acute coronary syndromes who receive tirofiban and aspirin: a randomized controlled trial. *JAMA* 292:55-64, 2004
26. Boersma E, Harrington R, Moliterno D, et al: Platelet glycoprotein IIb/IIIa inhibitors in acute coronary syndromes: a meta-analysis of all major randomised clinical trials. *Lancet* 359:189-198, 2002
27. Simon R: Bayesian design and analysis of active control clinical trials. *Biometrics* 55:484-487, 1999
28. Spiegelhalter DJ, Freedman LS, Parmar MKB: Bayesian approaches to randomized trials. *J Royal Stat Soc A Series* 157:357-416, 1994
29. Diamond GA, Kaul S: Prior convictions: Bayesian approaches to the analysis and interpretation of clinical megatrials. *J Am Coll Cardiol* 43:1929-1939, 2004
30. Lubsen J, Kirwan BA: Combined endpoints: can we use them? *Stat Med* 21:2959-2970, 2002
31. Montori VM, Gaieta P-M, Ignacio F-G, et al: Validity of composite endpoints in clinical trials. *Br Med J* 330:594-596, 2005
32. Bucher HC, Guyatt GH, Griffith LE, et al: The results of direct and indirect treatment comparisons in meta-analysis of randomized controlled trials. *J Clin Epidemiol* 50:683-691, 1997
33. Antman EM: Should bivalirudin replace heparin during percutaneous coronary interventions? *JAMA* 289:903-905, 2003

34. Kaul S, Diamond GA, Weintraub WS: Trials and tribulations of non-inferiority: the ximelagatran experience. Reply to the letter to the editor. *J Am Coll Cardiol* 48:1059, 2006
35. Ware JH, Antman EM: Equivalence trials. *N Engl J Med* 337:1159-1161, 1997
36. Greene WI, Concato J, Feinstein AR: Claims of equivalence in medical research: are they supported by the evidence? *Ann Intern Med* 132: 715-722, 2000
37. Le Henanff A, Giraudeau B, Baron G, et al: Quality of reporting of noninferiority and equivalence randomized trials. *JAMA* 295:1147-1151, 2006
38. Piaggio G, Elbourne DR, Altman DG, et al: Reporting of noninferiority and equivalence randomized trials: an extension of the CONSORT statement. *JAMA* 295:1152-1160, 2006
39. Gotsche PC: Lessons from and cautions about noninferiority and equivalence randomized trials. *JAMA* 295:1172-1174, 2006