# Exact inference for adaptive group sequential designs

## Ping Gao,[a] Lingyun Liu[b] and Cyrus Mehta[b,c*†]

**Methods for controlling the type-1 error of an adaptive group sequential trial were developed in seminal papers by Cui, Hung, and Wang (Biometrics, 1999), Lehmacher and Wassmer (Biometrics, 1999), and Müller and Schäfer (Biometrics, 2001). However, corresponding solutions for the equally important and related problem of parameter estimation at the end of the adaptive trial have not been completely satisfactory. In this paper, a method is provided for computing a two-sided confidence interval having exact coverage, along with a point estimate that is median unbiased for the primary efficacy parameter in a two-arm adaptive group sequential design. The possible adaptations are not only confined to sample size alterations but also include data-dependent changes in the number and spacing of interim looks and changes in the error spending function. The procedure is based on mapping the final test statistic obtained in the modified trial into a corresponding backward image in the original trial. This is an advance on previously available methods, which either produced conservative coverage and no point estimates or provided exact coverage for one-sided intervals only. Copyright © 2013 John Wiley & Sons, Ltd.**

**Keywords:** estimation in adaptive design; exact adaptive confidence intervals; adaptive median unbiased estimates; group sequential estimation

## 1. Introduction

Group sequential designs are widely used in randomized clinical trials intended to demonstrate the efficacy and safety of new medical compounds. In a classical two-arm group sequential trial, key design parameters such as the number and spacing of the interim looks, the corresponding early stopping boundaries, and the maximum sample size are pre-specified. They may only be altered through a *blinded* analysis of the accumulating data, that is, by examining the data pooled over both treatment arms. Possible reasons for such design alterations might be slow patient accrual, unanticipated variability in the data, new results from external sources, or a combination of such factors, none of which require the data of the trial to be unblinded. In contrast, an adaptive group sequential trial permits data-dependent alterations of the key design parameters. It is thus permissible to alter the sample size, skip or add interim looks, alter the error spending function, or even alter the inclusion/exclusion criteria of the remainder of the trial after examining the interim data, *unblinded by treatment arm*. A recent survey was conducted by the Adaptive Design Scientific Working Group of the Drug Information Association [1] to document the perception and use of adaptive designs in industry and academia. Nine pharmaceutical/biotechnology companies, six contract research organizations, and one academic institution responded to the survey. Between them, they identified 51 confirmatory trials involving sample size re-estimation, 30 of them based on an unblinded analysis of accumulating data. Given that only 20% of the organizations contacted actually responded to the survey, it may be conjectured that unblinded sample size re-estimation is an important recent innovation influencing the practice of clinical trials. The primary motivation for unblinded sample size re-estimation and related adaptive modifications is the uncertainty regarding the efficacy of the new treatment relative to the control. Often, this efficacy parameter is chosen on the basis

[a]*The Medicines Company, Parsippany, New Jersey 07054, U.S.A.*
[b]*Cytel Corporation, 675 Massachusetts Avenue, Cambridge, Massachusetts 02139, U.S.A.*
[c]*Harvard School of Public Health, Huntington Avenue, Boston, Massachusetts, U.S.A.*
*\*Correspondence to: Cyrus Mehta, Cytel Corporation, 675 Massachusetts Avenue, Cambridge, Massachusetts 02139, U.S.A.*
†*E-mail: mehta@cytel.com*

of limited data from small pilot studies, making it desirable to consider a midcourse correction to the sample size at an interim analysis when a substantial amount of data are available for inspection from the trial itself. Mehta and Pocock [2], and Mehta [3] present several case studies of actual trials in which provision was made for such adaptive modifications.

Data-dependent modifications to an ongoing trial raise operational and statistical concerns. Operational issues, such a who may have access to the unblinded data, how such unblinded access could lead to operational biases, and the regulatory implications of such biases are discussed in The Guidance for Industry on Adaptive Design for Clinical Trials published by the Food and Drug Administration [4] and are outside the scope of this paper. The two major statistical problems for an adaptive group sequential trial are hypothesis testing and parameter estimation. Specifically, how can we prevent inflation of the type-1 error, and how can we obtain valid $p$-values, confidence intervals, and point estimates in an adaptive group sequential trial?

Cui, Hung, and Wang [5], and Lehmacher and Wassmer [6] showed that the type-1 error of an adaptive group sequential trial can be preserved by combining the independent data from the different stages of the trial with pre-specified weights. This approach is, however, only applicable for sample size alterations. A more general approach that permits, among other options, changes in the sample size, the number of interim looks, the spacing of interim looks, the error spending function, and subgroup selection, was proposed by Müller and Schäfer [7]. Their method is based on the principle of preserving the type-1 errors of the original and altered trials, conditional on the data obtained up to the time of the adaptation.

So far, no satisfactory method has been published for the related problem of parameter estimation. Cui, Hung, and Wang [5], and Müller and Schäfer [7] did not address this question. Lehmacher and Wassmer [6] proposed extending Jennison and Turnbull's [8] repeated confidence intervals method by applying it to the inverse-normal weighted statistic. Repeated confidence intervals do not exhaust the entire type-1 error and hence produce conservative coverage of the efficacy parameter (see [9] page 198, Table 9.3 and discussion). The simulation results in Section 5 demonstrate that the coverage of the Lehmacher and Wassmer [6] method is far in excess of what was requested. Mehta, Bauer, Brannath, and Posch [10] also proposed an approach on the basis of extending Jennison and Turnbull's [8] repeated confidence intervals. Their solution, on the basis of a generalization of the hypothesis testing procedure of Müller and Schäfer [7], is applicable to a broader class of adaptive changes than the method of Lehmacher and Wassmer [6]. However, their approach too produces conservative coverage and, moreover, has only been developed for one-sided confidence bounds. Furthermore, neither of the two proposed methods can provide a valid point estimate for the efficacy parameter. More recently, Brannath, Mehta, and Posch [11] proposed a one-sided lower confidence bound for the efficacy parameter, on the basis of extending the stagewise adjusted confidence intervals of Tsiatis, Rosner, and Mehta [12]. They were able to prove that their method provides exact coverage for the special case in which the adaptive alteration occurs at the penultimate look and is followed by the final analysis. For all other cases, a formal proof of exact coverage relied on a monotonicity assumption that they were unable to demonstrate mathematically. Nevertheless, they were able to claim near-exact coverage of the lower confidence bound and median unbiasedness of the point estimate through extensive simulation experiments. Brannath, Mehta, and Posch [11] did not provide a method for two-sided confidence intervals.

The present paper provides a method for obtaining median-unbiased point estimates and exact two-sided confidence intervals for adaptive group sequential designs. We are not aware of published inference methods that have these operating characteristics. Our method generalizes the stagewise adjusted confidence intervals developed by Tsiatis, Rosner, and Mehta [12] for classical group sequential designs, and the hypothesis tests developed by Müller and Schäfer [7] for adaptive group sequential designs, and combines these two ideas in a novel manner to produce what we refer to as *backward image* confidence intervals (BWCI). Section 2 is a brief review of classical group sequential inference. Section 3 describes the Müller and Schäfer [7] method for performing valid hypothesis tests in an adaptive setting. The main results of this paper are presented in Section 4 where the backward image method for computing $p$-values point estimates and confidence intervals is developed. Section 5 presents extensive simulation results that demonstrate median unbiasedness and exact coverage. Section 6 illustrates the method through a worked example of a clinical trial of deep brain stimulation for Parkinson's disease. This example was first provided by Müller and Schäfer [7]. We conclude with some final remarks in Section 7. Proofs of various technical propositions are given in Appendices A.1 to A.7.

## 2. Inference for the classical group sequential design

Consider a two-arm randomized clinical trial comparing a new treatment to an active control. The treatment effect is captured by a single parameter $\theta$ that might denote the difference of means for two normal distributions, the difference of proportions for two binomial distributions, the log hazard ratio for two survival distributions, or more generally, the coefficient of the treatment effect in a regression model. The accumulating data are captured by the efficient score statistic

$$W(t) = \hat{\theta} t$$

where $\hat{\theta}$ is the maximum likelihood estimate of $\theta$ and

$$t = [\text{se}(\hat{\theta})]^{-2}$$

is the Fisher information for $\theta$ obtained from the available data. Because $t$ depends on unknown parameters, it is replaced, in practice, by its large sample estimate. Furthermore, as is well known (e.g., Jennison and Turnbull, [8]), $W(t)$ converges in distribution to a Brownian motion with drift $\theta$. That is,

$$W(t) \xrightarrow{D} B(t) + \theta t \tag{1}$$

where $B(t) \sim N(0, t)$ and for any $t_2 > t_1$, $\text{cov}\{B(t_1), B(t_2)\} = t_1$.

We shall be interested in testing the null hypothesis $H_0 : \theta = 0$ versus the one-sided alternative $\theta > 0$ and will assume throughout that a positive value of $\theta$ indicates a better prognosis for the treatment arm relative to the control arm. The following group sequential trial will be employed to test $H_0$. Analyses are planned at information times $t_1^{(1)}, t_2^{(1)}, \ldots t_{K_1}^{(1)}$ with corresponding critical values $c_1^{(1)}, c_2^{(1)}, \ldots c_{K_1}^{(1)}$. (The superscript $^{(1)}$ is employed in anticipation of the next section where we will distinguish between interim analyses before and after a trial modification by superscripts $^{(1)}$ and $^{(2)}$, respectively.) The trial is terminated, and null hypothesis $H_0$ is rejected at the first information time, $t_j^{(1)}$ say, such that $W\left(t_j^{(1)}\right) \geqslant c_j^{(1)}$. If $W\left(t_j^{(1)}\right) < c_j^{(1)}$ for all $j = 1, 2, \ldots K_1$, then $H_0$ is retained. For a one-sided level-$\alpha$ test of $H_0$, the critical values, $c_1^{(1)}, c_2^{(1)}, \ldots c_{K_1}^{(1)}$, must satisfy the relationship

$$P_0\left(\bigcup_{i=1}^{K_1}\left[W\left(t_i^{(1)}\right) \geqslant c_i^{(1)}\right]\right) = \alpha, \tag{2}$$

where $P_\delta(.)$ represents probability under the assumption that $\theta = \delta$. The recursive integration algorithm of Armitage, McPherson, and Rowe [13] combined with the $\alpha$-spending methodology of Lan and DeMets [14] may be used to find the critical values, $c_1^{(1)}, c_2^{(1)}, \ldots c_{K_1}^{(1)}$, that satisfy (2). Group sequential clinical trials of normal, binomial, and time-to-event endpoints are important special cases of this general formulation.

Suppose that the trial is terminated at information time $t_I^{(1)}$ with $W\left(t_I^{(1)}\right) = x_I^{(1)}$. We have thus observed the event

$$A\left(t_I^{(1)}, x_I^{(1)}\right) = \bigcap_{i=1}^{I-1}\left[W\left(t_i^{(1)}\right) < c_i^{(1)}\right] \cap \left[W\left(t_I^{(1)}\right) = x_I^{(1)}\right].$$

To test the null hypothesis $H_\delta : \theta = \delta$ versus the one-sided alternative $\theta > \delta$, we must identify all events that are at least as extreme as $A\left(t_I^{(1)}, x_I^{(1)}\right)$ and sum their probabilities under $H_\delta$. On the basis of the *stagewise ordering of events* ([9], page 179), an event $A\left(t_J^{(1)}, x_J^{(1)}\right)$ is at least as extreme as an event $A\left(t_I^{(1)}, x_I^{(1)}\right)$ if either $J < I$ or $J = I$ and $x_J^{(1)} \geqslant x_I^{(1)}$. The one-sided $p$-value of the observed event $A\left(t_I^{(1)}, x_I^{(1)}\right)$ for the test of $H_\delta$ is thus

$$f_\delta\left(t_I^{(1)}, x_I^{(1)}\right) = P_\delta\left(\bigcup_{i=1}^{I-1}\left[W\left(t_i^{(1)}\right) \geqslant c_i^{(1)}\right] \cup \left[W\left(t_I^{(1)}\right) \geqslant x_I^{(1)}\right]\right), \tag{3}$$

and $H_\delta$ is rejected at level $\alpha$ if and only if $f_\delta\left(t_I^{(1)}, x_I^{(1)}\right) \leq \alpha$. This is a valid level-$\alpha$ test of $H_\delta$ because, as proven in Appendix A.1, $f_\delta\left(t_I^{(1)}, x_I^{(1)}\right)$ satisfies the defining property of a $p$-value,

$$P_\delta\left\{f_\delta\left(t_I^{(1)}, x_I^{(1)}\right) \leq p\right\} = p, \tag{4}$$

for any $\delta$ and any $p \in (0, 1)$. Note that in equation (4), we are treating $\left(t_I^{(1)}, x_I^{(1)}\right)$ as a random variable that assumes different values in hypothetical repetitions of the group sequential trial.

Equation (3) shows that, for a fixed outcome $\left(t_I^{(1)}, x_I^{(1)}\right)$, $f_\delta\left(t_I^{(1)}, x_I^{(1)}\right)$ is a monotone increasing function of $\delta$. Thus, for any $p \in (0, 1)$, there exists a unique $\delta_p$, such that $f_{\delta_p}\left(t_I^{(1)}, x_I^{(1)}\right) = p$. Therefore, in hypothetical repetitions of the group sequential trial where $\left(t_I^{(1)}, x_I^{(1)}\right)$ is treated as a random variable,

$$P_\theta(\theta \leq \delta_p) = P_\theta\left\{f_\theta\left(t_I^{(1)}, x_I^{(1)}\right) \leq f_{\delta_p}\left(t_I^{(1)}, x_I^{(1)}\right)\right\} = P_\theta\left\{f_\theta\left(t_I^{(1)}, x_I^{(1)}\right) \leq p\right\} = p .$$

The first equality in the expression arises from the monotonicity of $f_\delta\left(t_I^{(1)}, x_I^{(1)}\right)$ with respect to $\delta$ for any fixed $\left(t_I^{(1)}, x_I^{(1)}\right)$. It follows that the interval $(\delta_{\alpha/2}, \delta_{1-\alpha/2})$ is a $100\times(1-\alpha)\%$ confidence interval for $\theta$. A median-unbiased point estimate for $\theta$ is given by $\delta_{0.5}$. These results, presented initially by Tsiatis, Rosner, and Mehta [12], pertain only to classical group sequential trials. In this paper, we will extend them to the adaptive setting.

## 3. Adaptive alteration of statistical information

At any look $L < K_1$, with $W\left(t_L^{(1)}\right) = x_L^{(1)}$, it is possible to alter the number and spacing of the future looks on the basis of an examination of the data already obtained. Suppose it is decided to take $K_2$ future looks, at information times $t_1^{(2)}, t_2^{(2)}, \ldots t_{K_2}^{(2)}$. Let $c_1^{(2)}, c_2^{(2)}, \ldots c_{K_2}^{(2)}$ be corresponding critical values, so selected that

$$P_0\left\{\bigcup_{j=L+1}^{K_1} W\left(t_j^{(1)}\right) \geq c_j^{(1)} | W\left(t_L^{(1)}\right) = x_L^{(1)}\right\} = P_0\left\{\bigcup_{j=1}^{K_2} W\left(t_j^{(2)}\right) \geq c_j^{(2)} | W\left(t_L^{(1)}\right) = x_L^{(1)}\right\} . \tag{5}$$

We will continue to monitor the accumulating data and will reject $H_0$ at the first information time $t_I^{(2)} > t_L^{(1)}$ such that $W\left(t_I^{(2)}\right) \geq c_I^{(2)}$. If $W\left(t_i^{(2)}\right) < c_i^{(2)}$ for all $i = 1, 2, \ldots K_2$, then we will retain $H_0$ and set $t_I^{(2)} = t_{K_2}^{(2)}$. Müller and Schäfer [7] have shown that, despite this data-driven modification of the trial, the unconditional probability that such a procedure will reject $H_0$ remains $\alpha$. Equation (5) is referred to by Müller and Schäfer [7] as the principle of preserving the conditional rejection (CRP) probability (the CRP principle). It is based on the intuitive notion that if the future course of a trial is altered in such a way that the type-1 error conditional on the data observed so far remains the same for the original and altered trials, then the unconditional type-1 error of the original and altered trials is also preserved. Note that because $W(t)$ has independent increments, its stochastic behavior beyond look $L$ depends only on $x_L^{(1)}$ and not on earlier realizations of $W\left(t_i^{(1)}\right)$. Also, it is not necessary to pre-specify $K_2$ or the modified information times $t_1^{(2)}, t_2^{(2)}, \ldots t_{K_2}^{(2)}$. These modified design parameters can be chosen after examining the data that have accumulated up to and including information time $t_L^{(1)}$. The corresponding critical values $c_1^{(2)}, c_2^{(2)}, \ldots c_{K_2}^{(2)}$ in equation (5) are evaluated by recursive integration.

The setting in which the trial design is altered at the penultimate look, $L = K_1 - 1$, with a single future look at $K_2 = 1$, is an important special case. It covers, for example, two-stage designs ($K_1 = 2$), still the most common class of phase 3 designs with a sample size adaptation. It is possible to study the statistical properties of these designs in greater detail because, unlike the general case, closed-form formulae are available for the necessary computations. Suppose the sample size is modified at information time $t_{K_1-1}^{(1)}$, with $W\left(t_{K_1-1}^{(1)}\right) = x_{K_1-1}^{(1)}$, and a single future analysis at information time $t_1^{(2)}$ is proposed.

To test $H_0$ at level $\alpha$, we must preserve the conditional type-1 error of the altered test. This is achieved by finding the value of $c_1^{(2)}$ that satisfies the CRP condition

$$P_0 \left\{ W\left(t_{K_1}^{(1)}\right) \geq c_{K_1}^{(1)} \mid W\left(t_{K_1-1}^{(1)}\right) = x_{K_1-1}^{(1)} \right\} = P_0 \left\{ W\left(t_1^{(2)}\right) \geq c_1^{(2)} \mid W\left(t_{K_1-1}^{(1)}\right) = x_{K_1-1}^{(1)} \right\} . \tag{6}$$

We can invoke the results in Gao, Ware, and Mehta [15] to obtain

$$c_1^{(2)} = \left[ \frac{\sqrt{t_1^{(2)} - t_{K_1-1}^{(1)}}}{\sqrt{t_{K_1}^{(1)} - t_{K_1-1}^{(1)}}} \left( c_{K_1-1}^{(1)} - x_{K_1-1}^{(1)} \right) + x_{K_1-1}^{(1)} \right] . \tag{7}$$

## 4. *P*-value, confidence interval, and point estimate for $\theta$

If the trial terminates at some information time $t_I^{(1)}$ without an adaptive alteration, the classical *p*-value, confidence interval, and point estimate are computed as described in Section 2. So let us suppose that at information time $t_L^{(1)}$, with $W\left(t_L^{(1)}\right) = x_L^{(1)} < c_L^{(1)}$, there is an adaptive alteration such that there are potentially $K_2$ future analyses at information times $t_1^{(2)}, t_2^{(2)}, \ldots t_{K_2}^{(2)}$ having corresponding critical values $c_1^{(2)}, c_2^{(2)}, \ldots c_{K_2}^{(2)}$ that satisfy the CRP condition (5). Suppose the trial terminates at information time $t_I^{(2)}$ with observed statistic $x_I^{(2)}$. We will then have observed the event

$$A\left(x_L^{(1)}, t_I^{(2)}, x_I^{(2)}\right) = \bigcap_{i=1}^{L-1} \left[ W\left(t_i^{(1)}\right) < c_i^{(1)} \right] \cap \left[ W\left(t_L^{(1)}\right) = x_L^{(1)} \right] \bigcap_{i=1}^{I-1} \left[ W\left(t_i^{(2)}\right) < c_i^{(2)} \right] \cap \left[ W\left(t_I^{(2)}\right) = x_I^{(2)} \right] .$$

To test the null hypothesis $H_\delta$, we must compute the *p*-value or probability of obtaining an event at least as extreme as $A\left(x_L^{(1)}, t_I^{(2)}, x_I^{(2)}\right)$ under $H_\delta$. We next describe how to identify events that are at least as extreme as $A\left(x_L^{(1)}, t_I^{(2)}, x_I^{(2)}\right)$. Consider, for instance, an alternative event

$$A\left(\tilde{x}_{\tilde{L}}^{(1)}, \tilde{t}_{\tilde{I}}^{(2)}, \tilde{x}_{\tilde{I}}^{(2)}\right) = \bigcap_{i=1}^{\tilde{L}-1} \left[ W\left(t_i^{(1)}\right) < c_i^{(1)} \right] \cap \left[ W\left(t_{\tilde{L}}^{(1)}\right) = \tilde{x}_{\tilde{L}}^{(1)} \right] \bigcap_{i=1}^{\tilde{I}-1} \left[ W\left(\tilde{t}_i^{(2)}\right) < \tilde{c}_i^{(2)} \right] \cap \left[ W\left(\tilde{t}_{\tilde{I}}^{(2)}\right) = \tilde{x}_{\tilde{I}}^{(2)} \right]$$

in which $\tilde{L} \neq L, \tilde{x}_{\tilde{L}}^{(1)} \neq x_L^{(1)}, \tilde{t}_i^{(2)} \neq t_i^{(2)}, \tilde{c}_i^{(2)} \neq c_i^{(2)}, \tilde{I} \neq I$ and $\tilde{x}_{\tilde{I}}^{(2)} \neq x_I^{(2)}$. It is not obvious whether $A\left(\tilde{x}_{\tilde{L}}^{(1)}, \tilde{t}_{\tilde{I}}^{(2)}, \tilde{x}_{\tilde{I}}^{(2)}\right)$ is less extreme, as extreme, or more extreme than the observed event $A\left(x_L^{(1)}, t_I^{(2)}, x_I^{(2)}\right)$ in terms of deviations from the null hypothesis $H_\delta$. Stagewise ordering is not directly applicable in this setting because the number, spacing, and critical values of the analysis time points after adaptation differ between the two events. For a meaningful comparison, we need to measure the extremeness of each event with a common yardstick. This is achieved by transforming the event that was actually obtained in the adaptive trial into an equivalent event that might have been obtained in the original trial had there been no adaptation. To this end, we compute $\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right)$, the backward image of the observed outcome $\left(t_I^{(2)}, x_I^{(2)}\right)$, such that

$$P_\delta \left\{ \bigcup_{i=1}^{I-1} \left[ W\left(t_i^{(2)}\right) \geq c_i^{(2)} \right] \cup \left[ W\left(t_I^{(2)}\right) \geq x_I^{(2)} \right] \mid x_L^{(1)} \right\} = P_\delta \left\{ \bigcup_{i=L+1}^{J_\delta-1} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] \cup \left[ W\left(t_{J_\delta}^{(1)}\right) \geq x_{J_\delta}^{(1)} \right] \mid x_L^{(1)} \right\} . \tag{8}$$

We show in Appendix A.2 that the backward image of any observed outcome in the adaptive trial is unique and can easily be computed. This computation is the key to comparing outcomes after a trial modification. It implies, as shown in Appendix A.3, that

$$P_\delta \left\{ \bigcup_{i=1}^{L} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] \cup \left( \bigcap_{i=1}^{L} \left[ W\left(t_i^{(1)}\right) < c_i^{(1)} \right] \cap \left\{ \left[ \bigcup_{i=1}^{I-1} W\left(t_i^{(2)}\right) \geq c_i^{(2)} \right] \cup \left[ W\left(t_I^{(2)}\right) \geq x_I^{(2)} \right] \right\} \right) \right\} \tag{9}$$

and

$$P_\delta \left\{ \bigcup_{i=1}^{J_\delta-1} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] \cup \left[ W\left(t_{J_\delta}^{(1)}\right) \geq x_{J_\delta}^{(1)} \right] \right\} \tag{10}$$

are equal. Now, (9) is the total probability under $H_\delta$ of all the events that are at least as extreme as the event

$$A\left(x_L^{(1)}, t_I^{(2)}, x_I^{(2)}\right) = \bigcap_{i=1}^{L-1}\left[W\left(t_i^{(1)}\right) < c_i^{(1)}\right] \cap \left[W\left(t_L^{(1)}\right) = x_L^{(1)}\right] \bigcap_{i=1}^{I-1}\left[W\left(t_i^{(2)}\right) < c_i^{(2)}\right] \cap \left[W\left(t_I^{(2)}\right) = x_I^{(2)}\right] \quad (11)$$

in the modified trial, whereas (10) is the total probability under $H_\delta$ of all the events that are at least as extreme as the corresponding event

$$A\left(x_L^{(1)}, t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right) = \bigcap_{i=1}^{L-1}\left[W\left(t_i^{(1)}\right) < c_i^{(1)}\right] \cap \left[W\left(t_L^{(1)}\right) = x_L^{(1)}\right] \bigcap_{i=L+1}^{J_\delta-1}\left[W\left(t_i^{(1)}\right) < c_i^{(1)}\right] \cap \left[W\left(t_{J_\delta}^{(1)}\right) = x_{J_\delta}^{(1)}\right] \quad (12)$$

in the original trial, *in terms of stagewise ordering*. Because of the equality of (9) and (10), we can say that the events (11) and (12) are equally extreme. Therefore, we can convert the problem of computing the probability of all events at least as extreme as the observed event in the modified trial into the equivalent problem of computing the probability of all events at least as extreme as its backward image in the unmodified trial. It follows that the one-sided $p$-value of the observed event $A\left(x_L^{(1)}, t_I^{(2)}, x_I^{(2)}\right)$ can be computed from its backward image $A\left(x_L^{(1)}, t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right)$ as

$$f_\delta\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right) = P_\delta\left\{\bigcup_{i=1}^{J_\delta-1}\left[W\left(t_i^{(1)}\right) \geq c_i^{(1)}\right] \cup \left[W\left(t_{J_\delta}^{(1)}\right) \geq x_{J_\delta}^{(1)}\right]\right\}. \quad (13)$$

To show that this definition results in a valid level-$\alpha$ test of $H_\delta$, we must prove that, for any $p \in (0, 1)$, $f_\delta\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right)$ satisfies

$$P_\delta\left\{f_\delta\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right) \leq p\right\} = p, \quad (14)$$

the defining property of a $p$-value. This is proven in Appendix A.4.

Given a final outcome $\left(t_I^{(2)}, x_I^{(2)}\right)$ in the adaptive trial, we compute $(\delta_{\alpha/2}, \delta_{1-\alpha/2})$, the $100 \times (1-\alpha)\%$ two-sided confidence interval for $\theta$ and $\delta_{0.5}$ the median-unbiased point estimate for $\theta$ by the following procedure:

Find $\delta_{\alpha/2}$ and corresponding backward image $\left(t_{J_{\delta_{\alpha/2}}}^{(1)}, x_{J_{\delta_{\alpha/2}}}^{(1)}\right)$ such that

$$f_{J_{\delta_{\alpha/2}}}\left(t_{J_{\delta_{\alpha/2}}}^{(1)}, x_{J_{\delta_{\alpha/2}}}^{(1)}\right) = \alpha/2. \quad (15)$$

Next, find $\delta_{1-\alpha/2}$ and corresponding backward image $\left(t_{J_{\delta_{1-\alpha/2}}}^{(1)}, x_{J_{\delta_{1-\alpha/2}}}^{(1)}\right)$ such that

$$f_{J_{\delta_{1-\alpha/2}}}\left(t_{J_{\delta_{1-\alpha/2}}}^{(1)}, x_{J_{\delta_{1-\alpha/2}}}^{(1)}\right) = 1 - \alpha/2. \quad (16)$$

Finally, find $\delta_{0.5}$ and corresponding backward image $\left(t_{J_{\delta_{0.5}}}^{(1)}, x_{J_{\delta_{0.5}}}^{(1)}\right)$ such that

$$f_{J_{\delta_{0.5}}}\left(t_{J_{\delta_{0.5}}}^{(1)}, x_{J_{\delta_{0.5}}}^{(1)}\right) = 0.5. \quad (17)$$

For this procedure to produce a confidence interval that has exact $100 \times (1 - \alpha)\%$ coverage of $\theta$ and a point estimate that is median unbiased, it is necessary to show that the $p$-value $f_\delta\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right)$ generated by the backward image of the observed outcome $\left(t_I^{(2)}, x_I^{(2)}\right)$ is a monotone increasing function of $\delta$ for any fixed value of $\left(t_I^{(2)}, x_I^{(2)}\right)$. This is proven in Section 4.1 for a special case. We were, however, unable to construct a mathematical proof for the general case because, unlike the classical case discussed in Section 2, where the argument of $f_\delta(.)$ does not change with $\delta$, here, the backward image $\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right)$ is a function of $\delta$. An operational proof of monotonicity is, however, possible. Once the two-sided interval $(\delta_{\alpha/2}, \delta_{1-\alpha/2})$ has been obtained, one may use well-established one-dimensional

search techniques (for example the book *Numerical Recipes* by Press *et al.* [16], provided a fast routine for initially bracketing a minimum) to ascertain if the function $f_\theta(.)$ increases monotonically inside this interval. This monotonicity check should be implemented not just for the one interval that was derived from the data actually obtained but also for additional intervals generated by simulating the design a large number of times over a range of values for $\theta$. If monotonicity is established for every one of these simulated intervals and if, moreover, these intervals can be shown to cover the underlying parameter $\theta$ at the desired confidence level, one may conclude that the procedure has worked accurately for the trial under consideration. In this sense, the proposed approach may be regarded as an operational proof of monotonicity for a specific trial.

In Section 5, we provide an operational proof of monotonicity by the aforementioned approach for three different adaptive group sequential designs. Each design is simulated 100,000 times, with each of five distinct values of $\theta$. The monotonicity check was successful in every one of these $100,000 \times 5 \times 3 = 1,500,000$ intervals, and furthermore, median unbiasedness and exact $100 \times (1 - \alpha)\%$ coverage up to Monte Carlo accuracy were obtained for each value of $\theta$ in each design. Although these do not constitute a mathematical proof, they provide a practical way to verify that the procedure produces a valid confidence interval and point estimate for any specific adaptive clinical trial under consideration. Under the monotonicity assumption, it follows that

$$P_\theta(\theta \leqslant \delta_{\alpha/2}) = P_\theta \left\{ f_\theta \left( t_{J_\theta}^{(1)}, x_{J_\theta}^{(1)} \right) \leqslant f_{J_{\delta_{\alpha/2}}} \left( t_{J_{\delta_{\alpha/2}}}^{(1)}, x_{J_{\delta_{\alpha/2}}}^{(1)} \right) \right\}$$

$$= P_\theta \left\{ f_\theta \left( t_{J_\theta}^{(1)}, x_{J_\theta}^{(1)} \right) \leqslant \alpha/2 \right\} = \alpha/2 \,,$$

$$P_\theta(\theta \leqslant \delta_{1-\alpha/2}) = P_\theta \left\{ f_\theta \left( t_{J_\theta}^{(1)}, x_{J_\theta}^{(1)} \right) \leqslant f_{J_{\delta_{1-\alpha/2}}} \left( t_{J_{\delta_{1-\alpha/2}}}^{(1)}, x_{J_{\delta_{1-\alpha/2}}}^{(1)} \right) \right\}$$

$$= P_\theta \left\{ f_\theta \left( t_{J_\theta}^{(1)}, x_{J_\theta}^{(1)} \right) \leqslant 1 - \alpha/2 \right\} = 1 - \alpha/2,$$

and therefore, $P_\theta(\delta_{\alpha/2} \leqslant \theta \leqslant \delta_{1-\alpha/2}) = 1 - \alpha$.

### 4.1. Adaptation at look $K_1 - 1$ with $K_2 = 1$

For the special case that the adaptation occurs at the penultimate look and is followed by a single further analysis, the confidence interval based on the backward image is available in closed form and guarantees exact coverage. The point estimate is likewise guaranteed to be median unbiased. To see this, suppose we observe $x_{K_1-1}^{(1)}$ at the penultimate look. After the adaptation at the penultimate look, suppose we observe $x_1^{(2)}$ at $t_1^{(2)}$. Then, the backward image of $\left( t_1^{(2)}, x_1^{(2)} \right)$ satisfies the following equation:

$$P_\delta \left\{ W \left( t_{K_1}^{(1)} \right) > x_{K_1}^{(1)} \mid W \left( t_{K_1-1}^{(1)} \right) = x_{K_1-1}^{(1)} \right\} = P_\delta \left\{ W \left( t_1^{(2)} \right) > x_1^{(2)} \mid W \left( t_{K_1-1}^{(1)} \right) = x_{K_1-1}^{(1)} \right\} . \quad (18)$$

By the property of independent increments, the equation can be rewritten as

$$P_\delta \left\{ W \left( t_{K_1}^{(1)} \right) - W \left( t_{K_1-1}^{(1)} \right) > x_{K_1}^{(1)} - x_{K_1-1}^{(1)} \right\} = P_\delta \{ W \left( t_1^{(2)} \right) - W \left( t_{K_1-1}^{(1)} \right) > x_1^{(2)} - x_{K_1-1}^{(1)} \} . \quad (19)$$

Note that $W \left( t_{K_1}^{(1)} \right) - W \left( t_{K_1-1}^{(1)} \right)$ is normally distributed with mean $\delta \left( t_{K_1}^{(1)} - t_{K_1-1}^{(1)} \right)$ and variance $t_{K_1}^{(1)} - t_{K_1-1}^{(1)}$. Therefore, the backward image satisfies the following equation:

$$x_{K_1}^{(1)} = \frac{\sqrt{t_{K_1}^{(1)} - t_{K_1-1}^{(1)}}}{\sqrt{t_1^{(2)} - t_{K_1-1}^{(1)}}} \left( x_1^{(2)} - x_{K_1-1}^{(1)} \right) + x_{K_1-1}^{(1)} + \delta \sqrt{t_{K_1}^{(1)} - t_{K_1-1}^{(1)}} \left( \sqrt{t_{K_1}^{(1)} - t_{K_1-1}^{(1)}} - \sqrt{t_1^{(2)} - t_{K_1-1}^{(1)}} \right) . \quad (20)$$

## 5. Simulation experiments

We evaluated the operating characteristics of the backward image method for estimating $\theta$ by repeatedly simulating a number of adaptive group sequential designs. In this section, we report the results of three such simulation experiments. (Several additional simulation experiments were performed with similar

conclusions.) Each experiment involved simulating an adaptive group sequential design with five different values $\theta$. We simulated the adaptive group sequential trial 100,000 with each value of $\theta$, thereby producing 100,000 confidence intervals whose coverage of $\theta$ we then assessed. All the simulations utilized normally distributed data with mean $\theta$ and $\sigma = 1$ (assumed known).

### 5.1. First simulation experiment

In this simulation experiment, the original trial is designed for up to four equally spaced looks with the Lan and DeMets [14] O'Brien–Fleming-type error spending function (LD(OF) error spending function). The total sample size of 480 subjects provides slightly over 90% power to detect $\delta = 0.3$ with a one-sided level-0.025 group sequential test. At look 1, with 120 subjects enrolled, the conditional power under the estimated value of $\theta$ is evaluated, and if it falls between 30% and 90%, the so-called 'promising zone' [2], the sample size is increased by the amount necessary to boost the conditional power up to 90%, subject to a cap of 1000 subjects. The trial then proceeds with the new sample size, up to three additional equally spaced looks, and new stopping boundaries derived from the LD(OF) error spending function. The $\alpha$ error of the new stopping boundaries for the adaptive extension is derived from equation (5) so as to preserve the unconditional type-1 error of the trial despite the data-dependent adaptation. This trial is simulated 100,000 times with a fixed value of $\theta$. At the end of each simulation, the point estimate of $\theta$, $\delta_{0.5}$, and the corresponding 95% two-sided confidence interval, $(\delta_{0.025}, \delta_{0.975})$, are computed. If the trial crosses the stopping boundary at look 1, there is no adaptation, and the classical stagewise adjusted point and interval estimates are obtained as described in Section 2. If, however, there is a sample size adaptation at look 1, the point and interval estimates for $\theta$ are computed by the backward image method by using equations (15), (16), and (17). Simulation results for $\theta = -0.15, 0, 0.15, 0.3$ and $0.45$ are presented in Table I. Column 1 contains the true value of $\theta$ that was used in the simulations. Column 2 contains the median of the 100,000 $\delta_{0.5}$ estimates and demonstrates that $\delta_{0.5}$ is indeed a median-unbiased point estimate for $\theta$. Column 3 contains the proportion of the 100,000 confidence intervals that contain the true value of $\theta$. These intervals demonstrate 95% coverage up to Monte Carlo accuracy. Columns 4 and 5 display the proportion of intervals that exclude the true value of $\theta$ from below and above, respectively.

### 5.2. Second simulation experiment

In this simulation experiment, the original trial is designed for up to three equally spaced looks with the LD(OF) error spending function. The total sample size of 390 subjects provides about 90% power to detect $\delta = 0.3$ with a one-sided level-0.05 group sequential test. If the trial does not cross an early stopping boundary at look 1 or look 2, then at look 2, with 240 subjects enrolled, the conditional power under the estimated value of $\theta$ is evaluated, and if it falls in the promising zone, here specified to be between 20% and 90%, the sample size is increased by the amount necessary to boost the conditional power up to 90%, subject to a cap of 780 subjects. The trial then proceeds with the new sample size for up to three additional equally spaced looks with new stopping boundaries derived from the Lan and DeMets [14] Pocock-type error spending function (the LD(PK) error spending function). This trial was simulated 100,000 times with different values of $\theta$. The median of the 100,000 point estimates for $\theta$ and the coverage proportion of the corresponding 90% confidence intervals for $\theta$ are reported in Table II. It is seen that the point estimates are median unbiased and the confidence intervals have exact 90% coverage up to Monte Carlo accuracy.

**Table I.** Results from 100,000 simulations of a 4-look LD(OF) group sequential design (GSD) with adaptation at look 1 to a 3-look LD(OF) GSD, demonstrating that the point estimate is median unbiased and the two-sided 95% confidence intervals provide exact coverage of the true value of $\theta$ up to Monte Carlo accuracy.

| True value of $\theta$ | Median of 100,000 point estimates | Proportion intervals containing $\theta$ | Proportion of intervals that exclude $\theta$ | |
|---|---|---|---|---|
| | | | From below | From above |
| -0.15 | -0.14971 | 0.94893 | 0.02568 | 0.02539 |
| 0.0 | 0.000363 | 0.94976 | 0.02486 | 0.02538 |
| 0.15 | 0.149574 | 0.94939 | 0.02484 | 0.02577 |
| 0.3 | 0.30028 | 0.95111 | 0.02442 | 0.02447 |
| 0.45 | 0.44996 | 0.95017 | 0.02489 | 0.02494 |

**Table II.** Results from 100,000 simulations of a 3-look LD(OF) group sequential design (GSD) with adaptation at look 2 to a 3-look LD(PK) GSD demonstrating that the point estimate is median unbiased and the two-sided 90% confidence intervals provide exact coverage of the true value of $\theta$ up to Monte Carlo accuracy.

| True value of $\theta$ | Median of 100,000 point estimates | Proportion intervals containing $\theta$ | Proportion of intervals that exclude $\theta$ | |
| --- | --- | --- | --- | --- |
| | | | From below | From above |
| -0.15 | -0.14972 | 0.90007 | 0.05022 | 0.04971 |
| 0.0 | 0.00027 | 0.90073 | 0.04920 | 0.05007 |
| 0.15 | 0.14986 | 0.89866 | 0.04955 | 0.05179 |
| 0.3 | 0.2999 | 0.90087 | 0.04940 | 0.04973 |
| 0.45 | 0.44963 | 0.89929 | 0.05083 | 0.04988 |

### 5.3. Third simulation experiment—comparison with Lehmacher and Wassmer

An alternative two-sided confidence interval was proposed by Lehmacher and Wassmer [6] on the basis of extending the repeated confidence intervals of Jennison and Turnbull [8]. It is well known that these repeated confidence intervals provide conservative coverage for classical group sequential designs because of the possibility that the trial might stop early and not exhaust all the available $\alpha$. It would therefore be instructive to assess the extent to which these repeated confidence intervals are conservative in the adaptive setting. Accordingly, we created a design with three equally spaced looks derived from the LD(OF) spending function and a planned adaptation at the end of look 1. The total sample size of 480 subjects has 90.44% power to detect $\theta = 0.3$ with a one-sided test operating at significance level $\alpha = 0.025$. If the trial does not cross the early stopping boundary at look 1 then, with 160 subjects enrolled, the conditional power under the estimated value of $\theta$ is evaluated, and if it falls in the promising zone, here specified to be between 30% and 90.44%, the sample size is increased by the amount necessary to boost the conditional power up to 90%, subject to a cap of 960 subjects. The trial then proceeds with the new sample size for up to two additional equally spaced looks with new stopping boundaries derived from the LD(OF) error spending function. This trial was simulated 100,000 times with different underlying values of $\theta$. Table III compares the actual coverage of $\theta$ by 100,000 95% confidence intervals obtained by the backward image method (BWCI) and the repeated confidence intervals method (RCI) due to Lehmacher and Wassmer [6]. The median of the 100,000 point estimates generated by the BWCI method and by the stagewise adjusted confidence interval method (SWCI) due to Brannath, Mehta, and Posch [11] methods is also reported. No corresponding method for obtaining a point estimate from the RCI method was developed by Lehmacher and Wassmer [6]; hence, none is reported.

As expected, the BWCI method produces median-unbiased point estimates and 95% confidence intervals with exact coverage up to Monte Carlo accuracy. The SWCI method also produces median-unbiased point estimates but does not provide two-sided confidence intervals. The RCI method does not provide valid point estimates and produces confidence intervals with increasingly conservative coverage as $\theta$ increases. The reason for the increase in conservatism is that as $\theta$ increases, the probability of stopping early and hence of not exhausting the entire $\alpha$ increases.

It is also informative to examine the extent of the one-sided coverage by the three methods. This is shown in Table IV. The BWCI interval excludes the true value for $\theta$ with 0.025 probability symmetri-

**Table III.** Comparison of the coverage of 100,000 simulated 95% backward image confidence intervals (BWCI), stagewise adjusted confidence intervals (SWCI), and repeated confidence intervals (RCI). The underlying design is a 3-look LD(OF) group sequential design (GSD) with adaptation at look 1 to a 2-look LD(OF) GSD.

| True value of $\theta$ | Median of 100,000 Point Estimates | | | Actual Coverage of 95% CIs | | |
| --- | --- | --- | --- | --- | --- | --- |
| | BWCI | SWCI | RCI | BWCI | SWCI | RCI |
| -0.15 | -0.15027 | -0.149794 | NA | 0.95062 | NA | 0.95771 |
| 0.0 | 0.000118 | -0.000421 | NA | 0.95014 | NA | 0.95213 |
| 0.15 | 0.150858 | 0.149064 | NA | 0.95016 | NA | 0.95017 |
| 0.3 | 0.300286 | 0.301016 | NA | 0.95062 | NA | 0.97597 |
| 0.45 | 0.449971 | 0.451704 | NA | 0.94936 | NA | 0.9875 |

**Table IV.** Comparing the backward image confidence intervals (BWCI), stagewise adjusted confidence intervals (SWCI), and repeated confidence intervals (RCI) in terms of the probability that the lower and upper bounds, respectively, of a 95% confidence interval will exclude $\theta$. The underlying design, a 3-look LD(OF) group sequential design (GSD) with adaptation at look 1 to a 2-look LD(OF) GSD, is simulated 100,000 times.

| True value of $\theta$ | Probability of low CL > $\theta$ | | | Probability of up CL < $\theta$ | | |
|---|---|---|---|---|---|---|
| | BWCI | SWCI | RCI | BWCI | SWCI | RCI |
| -0.15 | 0.02505 | 0.0256 | 0.01905 | 0.02529 | NA | 0.02324 |
| 0.0 | 0.02462 | 0.0251 | 0.02448 | 0.02524 | NA | 0.02339 |
| 0.15 | 0.02473 | 0.0256 | 0.02585 | 0.02511 | NA | 0.02238 |
| 0.3 | 0.02411 | 0.0253 | 0.00654 | 0.02527 | NA | 0.01749 |
| 0.45 | 0.02470 | 0.0259 | 0.00075 | 0.02594 | NA | 0.01050 |

CL, control limit.

cally from below and above, whereas the RCI method is both extremely asymmetric as well as extremely conservative. The SWCI method excludes the true value for $\theta$ with probability 0.025 from below but is not applicable for exclusion from above.

## 6. Deep brain stimulation for Parkinson's disease

We illustrate our estimation methods with a clinical trial of Parkinson's disease. This example was first introduced by Müller and Schäfer [7] to illustrate their method for adaptive sample size re-estimation and was subsequently used by Brannath, Mehta, and Posch [11] to obtain a one-sided lower confidence bound for the treatment effect. Patients were randomized to either the experimental arm (deep-brain stimulation) or the control arm (standard of care) in equal proportions. The primary endpoint was the quality of life as measured by the 39-item Parkinson's Disease Questionnaire (the PDQ-39). The investigators wished to design the trial to have 90% power to detect an improvement of $\theta = 6$ points in PDQ-39 with a one-sided level-0.05 test of significance. The standard deviation was assumed to be $\sigma = 17$. Because the actual conduct of this trial has not been reported, all the design and monitoring assumptions in the remainder of this section are hypothetical and are used mainly to illustrate the estimation procedure.

The trial is designed initially with a maximum sample size of 282 subjects and up to three equally spaced analyses by using stopping boundaries derived from the $\gamma(-4)$ error spending function proposed by Hwang, Shih, and DeCani [17]. Such a design would call for monitoring the data after enrolling $n_1^{(1)} = 94$, $n_2^{(1)} = 188$, and $n_3^{(1)} = 282$ subjects, respectively. The corresponding stopping boundaries for the Wald statistic, $Z\left(n_i^{(1)}\right) = \hat{\theta}_i \sqrt{n_i^{(1)}}/(2\hat{\sigma}_i)$, $i = 1, 2, 3$, are $b_1^{(1)} = 2.794$, $b_2^{(1)} = 2.289$, and $b_3^{(1)} = 1.680$. It is convenient to use the Wald statistic rather than the score statistic for this example because it has a more familiar interpretation as a standardized treatment effect. Also, most software packages monitor data on the Wald scale. The two statistics are linked by the relationship $W\left(t_i^{(j)}\right) = \sqrt{n_i^{(1)}} Z\left(n_i^{(j)}\right)/2\sigma$.

Suppose that at the first interim analysis, when 94 subjects have been evaluated, the estimate of $\theta$ is $\hat{\delta}^{(1)} = 4.5$ with estimated standard deviation $\hat{\sigma} = 20$ so that $Z_1^{(1)} = 1.091$. At this point, it is decided to increase the sample size because, if in truth, $\theta = 4.5$ and $\sigma = 20$, the conditional power is only about 60%, whereas we would prefer to proceed with at least 80% conditional power. It is permissible to use any decision rule to increase the sample size for the remainder of the trial. However, to protect the type-1 error in the face of a data-dependant sample size alteration, we must preserve the conditional type-1 error of the original and adapted trials as depicted by equation (5). The conditional type-1 error of the original design is

$$P_0\left\{\bigcup_{i=2}^{3}\left[Z\left(n_i^{(1)}\right) \geq b_i^{(1)}|Z_1^{(1)} = 1.091\right]\right\} = 0.1033$$

Therefore, 0.1033 is the amount of type-1 error permissible for the adaptive extension of the trial conditional on $Z_1^{(1)} = 1.0901$. Now, it is convenient for design and monitoring purposes to think of this

**Table V.** Comparison of estimates generated by different methods.

| lccc Method | Low CL | Up CL | Estimate |
|---|---|---|---|
| BWCI | 1.43237 | 9.5224 | 5.53591 |
| Mehta 2008 | 1.191284 | NA | 4.314697 |
| Brannath 2009 | 1.43224 | NA | 5.53607 |

CL, control limit; BWCI, backward image confidence intervals.

adaptive extension as a separate secondary trial with an unconditional type-1 error of 0.1033. This follows from the independent increments structure of the score statistic. One can then use standard group sequential software to design the secondary trial with a type-1 error of 0.1033. After a thorough examination of all available efficacy and safety data, it is decided to enroll 300 subjects to the secondary trial, thereby increasing the total sample size of the combined trial by 40%—from 282 subjects to 394 subjects. It is further decided to monitor the secondary trial up to three times at $n_1^{(2)} = 100$, $n_2^{(2)} = 200$, and $n_3^{(2)} = 300$. The corresponding stopping boundaries must satisfy the CRP requirement

$$P_0 \left\{ \bigcup_{i=1}^{3} \left[ Z\left(n_i^{(2)}\right) \geq b_i^{(2)} \right] \right\} = 0.1033 , \tag{21}$$

for the adaptive procedure to preserve the unconditional type-1 error at level 0.05. It is decided to generate stopping boundaries that satisfy (21) with the $\gamma(-2)$ error spending function. Thereby, we obtain $b_1^{(2)} = 2.162$, $b_2^{(2)} = 1.781$ and $b_3^{(2)} = 1.351$. Such a design has 84% power to reject $H_0$ if $\theta = 4.5$ and $\sigma = 17$.

Suppose that the secondary trial proceeds to the second look after the recruitment of $n_2^{(2)} = 200$ subjects and a treatment effect of $\hat{\delta}_2^{(2)} = 6.6$ and a standard deviation of $\hat{\sigma}_2^{(2)} = 19.5$ are obtained. This leads to $z_2^{(2)} = (6.6\sqrt{200})/(2 \times 19.5)) = 2.393$. Because $z_2^{(2)}$ exceeds the critical value $b_2^{(2)} = 1.781$, the trial is stopped with rejection of the null hypothesis $\theta = 0$. By applying the backward image estimation method discussed in Section 4 the two-sided 90% confidence interval for $\theta$ is (1.43237, 9.5224), and the median-unbiased estimate is 5.53591.

It is instructive to compare these estimates with those produced by the alternative approaches referenced in this article. Accordingly, Table V compares the BWCI results with those produced by the SWCI method of Brannath *et al.* [11] and the RCI method of Mehta *et al.* [10]. The RCI method of Lehmacher and Wassmer [11] cannot be used because it is only applicable for adaptations of the statistical information, whereas in this example, we have also altered the number of future looks and the error spending function at the time of the interim analysis. The BWCI and Brannath *et al.* [11] results are extremely similar, suggesting that the two methods might represent different ways of performing the same underlying computation. We could not prove that this is the case, but it is plausible. Intuitively, the Brannath *et al.* [11] method computes the conditional error probability of the unmodified design by looking ahead from the look $L$ at which the adaptive change occurs, whereas the BWCI method begins with the results of the modified design and searches backward for the conditional error probability of the unmodified design. The comparisons are limited to the lower confidence bound and the point estimate because only the BWCI method produces a two-sided confidence interval. The Mehta *et al.* [10] result differs from the other two. Because it has been derived by the RCI principle, we may conclude, on the basis of the simulation results discussed in Section 5 and displayed in Tables (III) and (IV), that it has produced a conservative lower confidence bound and a negatively biased point estimate.

## 7. Concluding remarks

We have presented a new method for computing confidence intervals and point estimates for an adaptive group sequential trial. The confidence intervals are shown to produce exact coverage, and the point estimates are median unbiased. These results close an important gap that previously existed for inference on adaptive group sequential designs. Hypothesis tests that control the type-1 error have been available for over a decade (Cui, Hung, and Wang [5]; Lehmacher and Wassmer [6]; Müller and Schäfer [7]). The development of procedures to produce valid confidence intervals and point estimates proved to be much more challenging. The first methods to guarantee two-sided coverage (Lehmacher and Wassmer

[6]; Mehta, Bauer, Posch, and Brannath [10]) were shown to be conservative and did not produce valid point estimates. Subsequently, Brannath, Mehta, and Posch [11] proposed a procedure that does produce exact coverage and valid point estimates. However, it only produces one-sided intervals. Like the procedure presented here, the method of Brannath, Mehta, and Posch [11] depends for its validity on a monotonicity property. This property was difficult to verify in a one-sided setting because one end of the interval extends to infinity. In contrast, the two-sided interval discussed here provides a bounded region within which it is possible to verify monotonicity with standard search procedures. This has enabled us to provide an operational proof that the intervals have exact coverage and the point estimates are median unbiased.

The backward image method can be generalized to handle multiple adaptations. Suppose the original trial is modified $N - 1$ times, resulting in interim analyses at time points $t_1^{(i)} < t_2^{(i)} < \cdots < t_{K_i}^{(i)}$, $i = 1, 2, \ldots N-1$, with modifications occurring at the observations $W\left(t_{L^{(m)}}^{(m)}\right) = x_{L^{(m)}}^{(m)}$, $m = 1, 2, \ldots, N-1$, and with final termination at $W\left(t_{I^{(N)}}^{(N)}\right) = x_{I^{(N)}}^{(N)}$. Then for any specific value of $\theta = \delta$, the successive backward images $\left(t_{J^{(N-1)}}^{(N-1)}, x_{J^{(N-1)}}^{(N-1)}\right), \left(t_{J^{(N-2)}}^{(N-2)}, x_{J^{(N-2)}}^{(N-2)}\right), \ldots \left(t_{J^{(1)}}^{(1)}, x_{J^{(1)}}^{(1)}\right)$ can be obtained, leading to the stagewise adjusted $p$-value

$$f_\delta\left(t_{J^{(1)}}, x_{J^{(1)}}\right) = P_\delta \left\{ \bigcup_{i=1}^{J^{(1)}-1} \left[W\left(t_i^{(1)}\right) \geq c_i^{(1)}\right] \cup \left[W\left(t_{J^{(1)}}\right) \geq x_{J^{(1)}}\right] \right\}$$

where, for notational convenience, we have suppressed the dependence of $t_J^{(i)}, x_J^{(i)}$ on $\delta$. The confidence interval and median-unbiased estimate can now be obtained in the usual way. The details of this generalization will be worked out and presented in a future paper.

Our method was developed for parameter estimation at the conclusion of a one-sided group-sequential design with early efficacy stopping. The method extends easily to one-sided group sequential designs with both an efficacy boundary and a non-binding futility boundary. We have not investigated how to extend the method to two-sided designs. In practice, however, hypothesis tests in most clinical trials are one sided. It is rarely of clinical interest to test the null hypothesis $H_0 : \theta = 0$ against the two-sided alternative hypothesis $H_1 : |\theta| > 0$. If a positive value for $\theta$ indicates good prognosis, then one is interested in powering the trial against the one-sided alternative hypothesis $H_1^+ : \theta > 0$. If a negative value for $\theta$ indicates good prognosis, then one is interested in powering the trial against the one-sided alternative hypothesis $H_1^- : \theta < 0$. Thus, even if a study protocol specifies that a two-sided level-$\alpha$ test will be performed, the actual power calculations are based on a one-sided test at level-$\alpha/2$. On the other hand, it is of clinical interest to bound the value of $\theta$ within a two-sided confidence $\theta$ after the trial concludes. Finally, the entire development in this paper was expressed in terms of score statistics and so is applicable to all types of efficacy endpoints including normal, binomial, and survival endpoints and model-based endpoints derived from contrasts of regression parameters and estimated by maximum likelihood methods.

## Appendix

### A.1. Distribution of $f_\delta\left(t_I^{(1)}, x_I^{(1)}\right)$ for a nonadaptive trial

Suppose the treatment parameter $\theta$ has the value $\delta$. Suppose that the trial terminates at look $I$ and $\left(t_I^{(1)}, x_I^{(1)}\right)$ is the test statistic at the end of the trial. In hypothetical repetitions of the trial, $\left(t_I^{(1)}, x_I^{(1)}\right)$ is a random variable. We wish to show that $P_\delta\left\{f_\delta\left(t_I^{(1)}, x_I^{(1)}\right) \leq p\right\} = p$ for any $p \in (0, 1)$.

For any $\delta$ and any $j = 1, 2, \ldots K_1$, define

$$a_j(\delta) = P_\delta \left\{ \bigcup_{i=1}^{j} \left[W\left(t_i^{(1)}\right) \geq c_i^{(1)}\right] \right\}$$

Let $a_0(\delta) = 0$ and $a_{K_1+1}(\delta) = 1$. Then, because the events in the aforementioned probability expression are nested such that $\left\{\bigcup_{i=1}^{j}\left[W\left(t_i^{(1)} \geq c_i^{(1)}\right)\right]\right\} \subset \left\{\bigcup_{i=1}^{j+1}\left[W\left(t_i^{(1)} \geq c_i^{(1)}\right)\right]\right\}$, it follows that $a_0(\delta) <$

$a_1(\delta) \cdots < a_{K_1}(\delta) < a_{K_1+1}(\delta)$. Thus, for any $p \in (0, 1)$, there exists a unique $L_\delta \in \{1, 2, \ldots K_1\}$ such that $a_{L_\delta-1}(\delta) < p < a_{L_\delta+1}(\delta)$, and there exists a unique $x_{L_\delta}^{(1)}$ such that

$$P_\delta \left\{ \bigcup_{i=1}^{L_\delta-1} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] \cup \left[ W\left(t_{L_\delta}^{(1)}\right) \geq x_{L_\delta}^{(1)} \right] \right\} = p \,.$$

Because of the stagewise ordering, the event $f_\delta\left(t_I^{(1)}, x_I^{(1)}\right) \leq p$ occurs if and only if the trial terminates at look $I < L_\delta$ or at look $I = L_\delta$ with $x_I^{(1)} \geq x_{L_\delta}^{(1)}$, that is, if and only if the event

$$\bigcup_{i=1}^{L_\delta-1} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] \cup \left[ W\left(t_{L_\delta}^{(1)}\right) \geq x_{L_\delta}^{(1)} \right]$$

occurs. Thus,

$$P_\delta \left\{ f_\delta\left(t_I^{(1)}, x_I^{(1)}\right) \leq p \right\} = P_\delta \left\{ \bigcup_{i=1}^{L_\delta-1} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] \cup \left[ W\left(t_{L_\delta}^{(1)}\right) \geq x_{L_\delta}^{(1)} \right] \right\} = p \,.$$

### A.2. Uniqueness of the backward image

The backward image $\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right)$ of the observed outcome $\left(t_I^{(2)}, x_I^{(2)}\right)$ satisfies the following equation.

$$P_\delta \left\{ \bigcup_{i=1}^{I-1} \left[ W\left(t_i^{(2)}\right) \geq c_i^{(2)} \right] \cup \left[ W\left(t_I^{(2)}\right) \geq x_I^{(2)} \right] | x_L^{(1)} \right\} = P_\delta \left\{ \bigcup_{i=L+1}^{J_\delta-1} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] \cup \left[ W\left(t_{J_\delta}^{(1)}\right) \geq x_{J_\delta}^{(1)} \right] | x_L^{(1)} \right\} \,. \tag{A1}$$

Let us denote the left-hand side of (A1) by $\alpha^*(\delta)$, that is,

$$\alpha^*(\delta) = P_\delta \left\{ \bigcup_{i=1}^{I-1} \left[ W\left(t_i^{(2)}\right) \geq c_i^{(2)} \right] \cup \left[ W\left(t_I^{(2)}\right) \geq x_I^{(2)} \right] | x_L^{(1)} \right\} \,.$$

Also, let us define the conditional rejection probabilities $\alpha_J(\delta)$ as

$$\alpha_J(\delta) = P_\delta \left\{ \bigcup_{i=L+1}^{J} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] | x_L^{(1)} \right\} \,, \tag{A2}$$

for $J = L + 1, L + 2, \ldots K_1$. Let $\alpha_L(\delta) = 0$ and $\alpha_{K_1+1} = 1$. Note that $0 = \alpha_L(\delta) < \alpha_{L+1}(\delta) < \ldots < \alpha_{K_1}(\delta) < \alpha_{K_1+1}(\delta)$, which implies that there must exist a unique $J_\delta$ with $L + 1 \leq J_\delta \leq K_1 + 1$ such that $\alpha_{J_\delta-1}(\delta) < \alpha^*(\delta) < \alpha_{J_\delta}$. Then, the backward image must satisfy the following equation.

$$P_\delta \left\{ \bigcup_{i=L+1}^{J_\delta-1} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] \cup \left[ W\left(t_{J_\delta}^{(1)}\right) \geq x_{J_\delta}^{(1)} \right] | x_L^{(1)} \right\} = \alpha^*(\delta) \,. \tag{A3}$$

### A.3. Equivalence of equations (9) and (10)

Equation (9) can be written as

$$P_\delta \left\{ \bigcup_{i=1}^{L} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] \right\} + P_\delta \left\{ \left\{ \bigcap_{i=1}^{L} \left[ W\left(t_i^{(1)}\right) < c_i^{(1)} \right] \right\} \cap \left\{ \bigcup_{i=1}^{I-1} \left[ W\left(t_i^{(2)}\right) \geq c_i^{(2)} \right] \cup \left[ W\left(t_I^{(2)}\right) \geq x_I^{(2)} \right] \right\} \right\} \tag{A4}$$

and equation (10) can be written as

$$P_\delta \left\{ \bigcup_{i=1}^{L} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] \right\} + P_\delta \left\{ \left\{ \bigcap_{i=1}^{L} \left[ W\left(t_i^{(1)}\right) < c_i^{(1)} \right] \right\} \cap \left\{ \bigcup_{i=1}^{J_\delta-1} \left[ W\left(t_i^{(1)}\right) \geq c_i^{(1)} \right] \cup \left[ W\left(t_{J_\delta}^{(1)}\right) \geq x_{J_\delta}^{(1)} \right] \right\} \right\} \,. \tag{A5}$$

The first term of these two equations is the same. The second term of (A4) can be factored as

$$\int_{-\infty}^{c_1^{(1)}} p\left(0; x_1^{(1)}; 0; t_1^{(1)}\right) dx_1^{(1)} \int_{-\infty}^{c_2^{(1)}} p\left(x_1^{(1)}; x_2^{(1)}; t_1^{(1)}; t_2^{(1)}\right) dx_2^{(1)} \cdots \int_{-\infty}^{c_L^{(1)}} p\left(x_{L-1}^{(1)}; x_L^{(1)}; t_{L-1}^{(1)}; t_L^{(1)}\right)$$

$$P_\delta\left\{\bigcup_{i=1}^{I-1}\left[W\left(t_i^{(2)}\right) \geq c_i^{(2)}\right] \cup \left[W\left(t_I^{(2)}\right) \geq x_I^{(2)}\right] | x_L^{(1)}\right\} dx_L^{(1)},$$

(A6)

where $p\left(x_{i-1}^{(1)}, x_i^{(1)}, t_{i-1}^{(1)}, t_i^{(1)}\right)$ is the probability of a transition from the score $W\left(t_{i-1}^{(1)}\right) = x_{i-1}^{(1)}$ to the score $W\left(t_i^{(1)}\right) = x_i^{(1)}$. Similarly, the second term of (A5) can be factored as

$$\int_{-\infty}^{c_1^{(1)}} p\left(0; x_1^{(1)}; 0; t_1^{(1)}\right) dx_1^{(1)} \int_{-\infty}^{c_2^{(1)}} p\left(x_1^{(1)}; x_2^{(1)}; t_1^{(1)}; t_2^{(1)}\right) dx_2^{(1)} \cdots \int_{-\infty}^{c_L^{(1)}} p\left(x_{L-1}^{(1)}; x_L^{(1)}; t_{L-1}^{(1)}; t_L^{(1)}\right)$$

$$P_\delta\left\{\bigcup_{i=L+1}^{J_\delta-1}\left[W\left(t_i^{(1)}\right) \geq c_i^{(1)}\right] \cup \left[W\left(t_{J_\delta}^{(1)}\right) \geq x_{J_\delta}^{(1)}\right] | x_L^{(1)}\right\} dx_L^{(1)}.$$

(A7)

Therefore, by (8), equations (A6) and (A7) yield the same probability.

*A.4. Distribution of $f_\delta\left(t_{J_\delta, x_{J_\delta}^{(1)}}^{(1)}\right)$ for an adaptive trial*

Suppose the treatment parameter $\theta$ has the value $\delta$, and let $\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right)$ denote the test statistic at the end of the trial. If the trial has terminated after an adaptation, $\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right)$ is the backward image of the final test statistic that was obtained in the modified trial. Otherwise, it is the actual test statistic observed at termination. In hypothetical repetitions of the trial, $\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right)$ is a random variable. We wish to show that $P_\delta\{f_\delta\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right) \leq p\} = p$ for any $p \in (0, 1)$.

*Case 1: An adaptation is planned at a fixed look L.* Here, $L$ may be any look between 1 and $K_1$, where the choice $L = K_1$ corresponds to having planned not to modify the trial. We have shown in Appendix A.1 that for a given $\delta$ and $p$, there exists a unique $\left(t_{L_\delta}^{(1)}, x_{L_\delta}^{(1)}\right)$ such that

$$P_\delta\left\{\bigcup_{i=1}^{L_\delta-1}\left[W\left(t_i^{(1)}\right) \geq c_i^{(1)}\right] \cup \left[W\left(t_{L_\delta}^{(1)}\right) \geq x_{L_\delta}^{(1)}\right]\right\} = p.$$

Suppose that either $J_\delta < L_\delta$ or $J_\delta = L_\delta$ and $x_{J_\delta}^{(1)} \geq L_\delta$. Then, by (13), $f_\delta\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right) \leq p$. Conversely, suppose $f_\delta\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right) \leq p$. Then it must be the case that $J_\delta < L_\delta$ or $J_\delta = L_\delta$ and $x_{J_\delta}^{(1)} \geq x_L^{(1)}$. Thus, the event $f_\delta\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right) \leq p$ occurs if and only if either $J_\delta < L_\delta$ or $J_\delta = L_\delta$ and $x_{J_\delta}^{(1)} \geq x_L^{(1)}$. It follows that

$$P_\delta\left\{f_\delta\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right) \leq p\right\} = P_\delta\left\{\bigcup_{i=1}^{L_\delta-1}\left[W\left(t_i^{(1)}\right) \geq c_i^{(1)}\right] \cup \left[W\left(t_{L_\delta}^{(1)}\right) \geq x_{L_\delta}^{(1)}\right]\right\} = p.$$

*Case 2: An adaptation is planned at a random look L.* In this case,

$$P_\delta\left\{f_\delta\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right) \leq p\right\} = \sum_{L=1}^{K_1-1} P_\delta\left\{f_\delta\left(t_{J_\delta}^{(1)}, x_{J_\delta}^{(1)}\right) \leq p | L\right\} P(L) = p \sum_{L=1}^{K_1} P(L) = p.$$

## Acknowledgements

## References

1. DIA-ADSWG Survey Subteam. Perception and use of adaptive designs in the industry and academia: persistent barriers and recommendations to overcome challenges. *Unpublished Manuscript presented at the DIA Adaptive Designs in Clinical Trials Meeting on Nov 30, 2012*, Washington DC, 2012.
2. Mehta CR, Pocock SJ. Adaptive increase in sample size when interim results are promising: a practical guide with examples. *Statistics in Medicine* 2011; **30**(28):3267–3284.
3. Mehta CR. Sample size re-estimation for confirmatory clinical trials. In *Chapter 4 of Designs for Clinical Trials*, Harrington D (ed.). Springer: New York, 2012; 81–108.
4. Food and Drug Administration. Guidance for industry-adaptive design clinical trials for drugs and biologics, 2010.
5. Cui L, Hung MJ, Wang S-J. Modification of sample size in group sequential clinical trial. *Biometrics* 1999; **55**:853–857.
6. Lehmacher W, Wassmer G. Adaptive sample size calculations in group sequential trials. *Biometrics* 1999; **55**:1286–1290.
7. Müller HH, Schäfer H. Adaptive group sequential designs for clinical trials: combining the advantage of adaptive and of classical group sequential approaches. *Biometrics* 2001; **57**:886–891.
8. Jennison C, Turnbull BW. Interim analyses: the repeated confidence interval approach (with discussion). *Journal of the Royal Statistical Society B* 1989; **51**(3):305–61.
9. Jennison C, Turnbull BW. *Group Sequential Methods with Applications to Clinical Trials*. Chapman and Hall/CRC: London, 2000.
10. Mehta CR, Bauer P, Posch M, Brannath W. Repeated confidence intervals for adaptive group sequential trials. *Statistics in Medicine* 2007; **26**(30):5422–5433.
11. Brannath W, Mehta C, Posch M. Exact confidence intervals following adaptive sequential tests. *Biometrics* 2009; **64**:1–22.
12. Tsiatis AA, Rosner GL, Mehta C. Exact confidence intervals following a group sequential test. *Biometrics* 1984; **40**:797–803.
13. Armitage P, McPherson CK, Rowe BC. Repeated significance tests on accumulating data. *Journal of the Royal Statistical Society A* 1969; **132**:232–44.
14. Lan KKG, DeMets DL. Discrete sequential boundaries for clinical trials. *Biometrika* 1983; **70**:659–663.
15. Gao P, Ware JH, Mehta C. Sample size re-estimation for adaptive sequential designs. *Journal of Biopharmaceutical Statistics* 2008; **18**:1184–1196.
16. Press WH, Flannery BP, Teukolsky SA, Vetterling WT. *Numerical Recipes*. Cambridge University Press: New York, 1986.
17. Hwang IK, Shih WJ, DeCani JS. Group sequential designs using a family of type I error probability spending functions. *Statistics in Medicine* 1990; **9**(12):1439–1445.