

# Exploring the benefits of adaptive sequential designs in time-to-event endpoint settings

Sarah C. Emerson,<sup>a,\*†</sup> Kyle D. Rudser<sup>b</sup> and Scott S. Emerson<sup>c</sup>

Sequential analysis is frequently employed to address ethical and financial issues in clinical trials. Sequential analysis may be performed using standard group sequential designs, or, more recently, with adaptive designs that use estimates of treatment effect to modify the maximal statistical information to be collected. In the general setting in which statistical information and clinical trial costs are functions of the number of subjects used, it has yet to be established whether there is any major efficiency advantage to adaptive designs over traditional group sequential designs. In survival analysis, however, statistical information (and hence efficiency) is most closely related to the observed number of events, while trial costs still depend on the number of patients accrued. As the number of subjects may dominate the cost of a trial, an adaptive design that specifies a reduced maximal possible sample size when an extreme treatment effect has been observed may allow early termination of accrual and therefore a more cost-efficient trial. We investigate and compare the tradeoffs between efficiency (as measured by average number of observed events required), power, and cost (a function of the number of subjects accrued and length of observation) for standard group sequential methods and an adaptive design that allows for early termination of accrual. We find that when certain trial design parameters are constrained, an adaptive approach to terminating subject accrual may improve upon the cost efficiency of a group sequential clinical trial investigating time-to-event endpoints. However, when the spectrum of group sequential designs considered is broadened, the advantage of the adaptive designs is less clear. Copyright © 2010 John Wiley & Sons, Ltd.

**Keywords:** adaptive design; clinical trial; group sequential test; group sequential trial; statistical efficiency; survival analysis

## 1. Introduction

In designing a clinical trial, investigators typically determine the statistical information required to discriminate between a null hypothesis of no treatment effect and some alternative hypothesis representing the minimal clinically important difference. When evaluating a proposed clinical trial design, the sponsor must balance the scientific goals with the cost of the trial and the ethical issues of human experimentation.

The number of subjects involved and the duration of a clinical trial are two key factors in determining the overall cost of the trial. Requiring fewer patients will decrease the trial costs associated with screening, treatment, and follow-up. Shorter trials allow earlier profit from treatments that are effective and waste fewer resources (investigator time, cost of following patients over time, cost of money, etc.) on ineffective treatments. Trial duration also plays a role in ethical considerations: Ending a trial sooner rather than later will expose fewer patients to an ineffective or harmful treatment and allow the broader population of patients earlier access to new effective treatments. This also frees up patients for the investigation of other treatments, which in turn can speed up the process of new treatment discovery in a disease for which clinical trial participants may be in short supply.

<sup>a</sup>Department of Statistics, Oregon State University, U.S.A.

<sup>b</sup>Division of Biostatistics, University of Minnesota, U.S.A.

<sup>c</sup>Department of Biostatistics, University of Washington, U.S.A.

\*Correspondence to: Sarah C. Emerson, Kidder Hall 44, Oregon State University, Corvallis, OR 97331, U.S.A.

†E-mail: scemerson@gmail.com, emersosa@stat.oregonstate.edu

When measurements of treatment effect are based on comparisons of means or proportions, statistical information is generally directly proportional to the number of subjects accrued. Hence, any ability to decrease calendar time is generally related to an ability to increase the number of centers recruiting patients. However, in a survival analysis/time-to-event setting with potentially censored observations, statistical information is more directly related to the number of observed events rather than the number of subjects. Thus, in such a setting there are more tradeoffs possible between sample size and calendar time, as decreasing the number of subjects accrued generally increases the calendar time required to observe the requisite number of events and vice versa.

Historically, group sequential tests have been the primary statistical method used to address the ethical and efficiency concerns in clinical trials. In a typical group sequential design, a rule is specified for determining the maximal statistical information  $N_J$ . Then, at periodic intervals during the conduct of the study, up to  $J$  interim analyses are performed to determine whether the trial should stop early. More recently there has been much interest in the statistical literature related to 'adaptive designs'. Such adaptation typically takes the form of decisions to extend a clinical trial beyond some previously planned maximal stopping time. The advantages of such an approach over the group sequential design have not been established in general. Tsiatis and Mehta [1] and Jennison and Turnbull [2] have not found any efficiency gains of the adaptive approach over the more standard group sequential designs. However, such explorations have focused on the setting in which statistical information was directly proportional to the number of subjects accrued to the study.

In a survival analysis/time-to-event setting, there may be a clearer advantage to adaptive designs due to the need to consider both the number of patients accrued and the calendar time of follow-up necessary to observe the desired number of events. The appeal of an adaptive design in this setting is that it offers the possibility that early trends in the estimated treatment effect may suggest a modification of the number of subjects that need to be accrued. For instance, suppose a clinical trial is designed based on a maximal statistical information of  $N_J$ . Suppose further that at the  $j$ th analysis, the estimated treatment effects were so extreme that it seemed likely that the ultimate decision for efficacy or futility could be precisely determined prior to observing all  $N_J$  events. Then, it might seem advantageous to consider a re-designed trial in which the revised maximal number of observed events  $N_{j*}^*$  was strictly less than originally planned, i.e.  $N_{j*}^* < N_J$ . Such a re-design might allow the number of patients accrued to the study to be similarly decreased, thereby possibly reducing the number of subjects involved. Alternatively, if at an early interim analysis of the data a less extreme treatment effect were observed than was initially anticipated, the sponsor might want to increase the maximal number of events  $N_{j*}^* > N_J$  in order to increase the conditional power of the study to attain statistical significance. In such a setting, it may be necessary to increase the number of subjects accrued to the clinical trial in order to observe the increased number of events in an acceptable interval of calendar time.

Here, our goal is to explore the potential for such adaptation to decrease the total number of subjects required and/or the calendar time necessary to complete a trial with a time-to-event endpoint. As the focus of this manuscript is the relative flexibility of standard group sequential designs and more recently described adaptive designs to meet the optimality criteria of the collaborators in a clinical trial, we restrict attention to the case study of a single hypothetical clinical trial setting. This exploration is not intended to be exhaustive, but rather illustrative of the possibilities and the issues that might arise in designing such an adaptive trial and comparing it to standard group sequential methods. We consider a range of stopping rules, both group sequential and adaptive, that address the types of operating characteristics most often addressed in the statistical design of a clinical trial. In particular, we consider a setting that is a slight modification of a design proposed for an industry sponsored clinical trial. In that setting, the sponsor adopted a group sequential clinical trial design to detect a specified design alternative. In order to protect against the possibility that the observed treatment effect might be less than that indicated by the specified design alternative, the sponsor also incorporated an adaptive modification of both the maximal number of subjects to be accrued and the maximal number of events to be observed. The conditions under which the sampling rule was modified were defined based on an interim estimate of the treatment effect. In this manuscript, we model such a modification of the clinical trial design through an adaptive switch between two group sequential stopping rules. We note that the adaptation is completely prespecified as required for a confirmatory trial, and the statistic used to define the adaptive design stopping rule is a minimal sufficient statistic. Furthermore, because we consider a prespecified stopping rule, there is no need to address the worst-case scenarios that must be considered when any adaptation is not completely prespecified. Statistical inference can be based on

the distribution of the minimal sufficient statistic under the sampling plans specified by the sequential stopping rules.

In Section 2, we describe the survival analysis setting that will be used to compare the group sequential approach to a more adaptive approach. We define the particular form of the adaptive designs considered in this manuscript, as well as the design parameters that are held constant between the group sequential and adaptive designs. The relative behavior of the group sequential and adaptive designs are then investigated in the absence of censoring in Section 3 and in the presence of censoring in Section 4. Section 5 introduces a simple cost model used to summarize the comparisons between pairs of adaptive and group sequential designs. Although this model is undoubtedly overly simplistic, it is adequate to demonstrate that considerations of patient-related costs and time costs of money may lead to the selection of different trial designs. We conclude in Section 6 with a discussion of the impact of the particular design parameters and operating characteristics that were constrained to be equal in our comparisons, and demonstrate our ability to find more efficient group sequential designs when those constraints are relaxed, thereby illustrating the need for careful evaluation of a broad spectrum of group sequential rules when attempting to adaptively improve design operating characteristics.

## 2. Background

In designing a clinical trial, there are several competing concerns. Emerson *et al.* [3,4] discuss the breadth of operating characteristics commonly considered when comparing candidate sampling plans for a clinical trial. Efficiency, statistical power to detect an effect of interest and ethical considerations are key factors in assessing the suitability of a proposed design. Efficiency generally refers to the number of subjects or events required, and is often measured as the expected sample size or *average sample number* (ASN). The maximal possible sample size is also frequently considered in efficiency comparisons. Power to detect an effect is the probability of deciding to reject the null hypothesis at a given effect size. Both power and ASN are functions of the true effect size. Ethical considerations are also addressed by minimizing the number of subjects and the time required to complete the trial.

Sequential analysis is a tool that is often employed to address these tradeoffs. The basic idea is that if results are convincing early on, there is no need to increase costs and lose efficiency by continuing with more subjects. At the  $j$ th of  $J$  potential interim analyses, a test statistic  $T_j$  is computed and compared to stopping boundaries. Following Kittelson and Emerson [5], it is generally sufficient to define at each analysis up to four stopping boundaries  $a_j \leq b_j \leq c_j \leq d_j$ , with early termination if  $T_j \leq a_j$ , if  $b_j < T_j < c_j$ , or if  $T_j \geq d_j$ . Designs with fewer than four early stopping boundaries can be obtained by setting  $a_j = -\infty$ ,  $b_j = c_j$ , or  $d_j = \infty$ , as appropriate for the setting. Ensuring that at the  $J$ th analysis  $a_J = b_J$  and  $c_J = d_J$  guarantees termination of the study. Group sequential clinical trial design typically involves choosing stopping boundaries, which will maintain a desired type I error, and choosing a maximal statistical information  $N_J$  such that the study will have adequate power to detect a specified design alternative under a schedule of analyses occurring when statistical information is  $N_1, N_2, \dots, N_J$ .

It should be noted that  $N_J$  can be specified in units of some unknown variance of individual observations, in which case only the maximal sample size is prespecified. Alternatively, the maximal statistical information is prespecified and the actual maximal sample size is determined using estimates of the variance observed during the conduct of the clinical trial. Emerson [6] discusses further the scientific and statistical validity of both approaches to the prespecification of the rule for determining maximal statistical information in the setting of group sequential clinical trials. In the analysis of time-to-event data that are potentially right-censored, it is most common to design a clinical trial based on the maximal statistical information, which is proportional to the number of events. The number of subjects required is then determined based on accrual rate, accrual time, and the distribution of failure times.

By way of example, we will consider a group sequential design appropriate for testing for a lessened instantaneous risk of death (decreased hazard) when administering some new treatment or placebo to a population of severely ill patients. We will let  $\theta(x)$  denote the hazard ratio  $h_T(x)/h_C(x)$ , where  $h_T(x)$  is the hazard at time  $x$  for the treatment arm and  $h_C(x)$  is the hazard at time  $x$  for the control arm. We use a proportional hazards model, meaning that we assume that  $\theta(x) \equiv \theta$  is constant for all values of  $x$ .

As noted previously, we consider a setting that is a slight modification of a design proposed for an industry sponsored clinical trial. In keeping with the sponsor's initial exploration of designs, we consider in depth comparisons in which the accrual rate is the same for a group sequential design

and an adaptive design, and interim analyses are to be performed at identical intervals of accrual of statistical information. We start with an initial simple group sequential design, Design A, with two analysis points, at 100 and 200 events. The test statistic for this design is the estimated hazard ratio, and we desire a statistical sampling plan based on a one-sided level 0.025 type I error to discriminate between a null hypothesis of no effect (a hazard ratio of 1.0) and a design hypothesis of improved survival on the treatment arm relative to the control arm (hazard ratio less than 1.0).

## 2.1. Specification of the initial group sequential design: Design A

In selecting a stopping rule to be used as a guideline for early termination of the clinical trial, we take the common approach of selecting an efficacy (lower) boundary that is relatively conservative at the earliest interim analyses. The motivation for such early-conservatism is that the standards of evidence for adoption of a new, unproven treatment dictate that the drawbacks of having less available data to examine safety, longer term survival, and other secondary endpoints need to be counterbalanced by a marked benefit on survival over the shorter period of observation prior to the first interim analysis. That is, with lesser follow-up, we need to be confident of a highly effective treatment, and we should focus on stopping rules that would stop early only if, say, a 95 per cent confidence interval includes only hazard ratios that correspond to strong effect of the new treatment. We note that it is common for clinical trialists to focus instead on the criterion that early stopping should occur only if we are highly confident of an effective treatment, e.g. perhaps a 99.9 per cent confidence interval that excludes 1.0. However, such a criterion may not be measuring the most important scientific, clinical, and ethical issues that relate to ensuring that early-occurring effects are of sufficient benefit to outweigh uncertainty about safety and long-term effects.

We also consider the specification of a futility (upper) boundary that would correspond to a decision that the treatment effect is not sufficiently beneficial to warrant continued study. The advantages of early termination of a study for futility are that it avoids continued exposure of patients to an unproven therapy that is unlikely to be adopted and that it avoids continued delay of investigating other, more promising therapies. If there is no important secondary information that can be obtained from studying a therapy that we have with high confidence determined is not associated with a clinically important effect, there is also no need to build in the early-conservatism desirable for the efficacy boundary. Hence, a futility boundary might be chosen to afford more efficiency.

In the unified family described by Kittelson and Emerson [5], there are typically several different parameterizations that lead to nearly the same stopping boundary or, in the case of a stopping rule with a maximum of two analyses, the exact same boundary. We prefer parameterizations that maintain the same level of confidence when rejecting the null hypothesis (when setting the efficacy boundary) and rejecting the design alternative (when setting the futility boundary). Hence, in this case we choose to define the boundary to have type I error  $\alpha=0.025$  and to reject the design alternative with 97.5 per cent power. We further restrict attention to boundary shapes within the extended Wang and Tsatis [7] family. In the unified family, these boundary shape functions have parameters  $A=0$  and  $R=0$ , with  $P$  allowed to vary. In this setting, the  $P$  parameter measures the early-conservatism of the boundary, with  $P=\infty$  corresponding to no early stopping. Two important cases within this family are the Pocock [8] boundary shape function when  $P=0.5$  and the O'Brien-Fleming [9] boundary shape function when  $P=1$ . Each boundary may have its own parameterization, and when these parameters differ we will let  $P_a$  denote the  $P$  parameter for the efficacy ( $a$ ) boundary, and  $P_d$  denote the  $P$  parameter for the futility ( $d$ ) boundary. Emerson and Fleming [10] found that boundary shape functions similar to the Pocock [8] boundaries tended to nearly minimize the ASN under the hypothesis being accepted. That is, the ASN-efficient efficacy boundary shape for rejecting the null when the design alternative is true is close to a Pocock boundary, as would be the ASN-efficient futility boundary shape for rejecting the design alternative when the null is true. In accordance with the above criteria, we use the popular O'Brien-Fleming [9] boundary for the efficacy boundary ( $P_a=1$ ), but we choose a less conservative, more efficient Pocock [8] futility boundary ( $P_d=0.5$ ). In considering the efficiency of the stopping boundary, we follow the most common approach based on examination of the ASN.

To summarize, Design A is specified by the following constraints: two analyses are performed at sample sizes  $N_1=100$  events and  $N_2=200$  events; the type I error is set at  $\alpha=0.025$ ; the power at the design alternative is set at 97.5 per cent; the efficacy boundary is determined by the O'Brien-Fleming  $P_a=1$  boundary shape; and the futility boundary is determined by the Pocock  $P_d=0.5$  boundary shape. Using S+SeqTrial [11], we compute the design boundaries giving boundaries at the first analysis



of  $a_1=0.5792$ ,  $b_1=c_1$ ,  $d_1=0.8645$ . The study stops at the first analysis (100 observed events) if the estimated hazard ratio  $T_1 \leq a_1=0.5792$  (in which case the null hypothesis  $\theta \geq 1.0$  is rejected) or  $T_1 \geq d_1=0.8645$  (in which case the null hypothesis  $\theta \geq 1$  is not rejected). At the second (and maximal) analysis, a result corresponding to  $T_2 \leq a_2=d_2=0.7611$  would correspond to a rejection of the null hypothesis. With these boundary shape parameters, this design provides 97.5 per cent power to detect a hazard ratio of  $\theta_1=0.5596$ . Hence, the futility boundary can be interpreted based on a rejection of that design alternative, just as the efficacy boundary rejects the null hypothesis of a hazard ratio of 1.0. A more detailed examination of the power curve reveals that this design has 80 per cent power against a hypothesis  $\theta_1=0.6646$ .

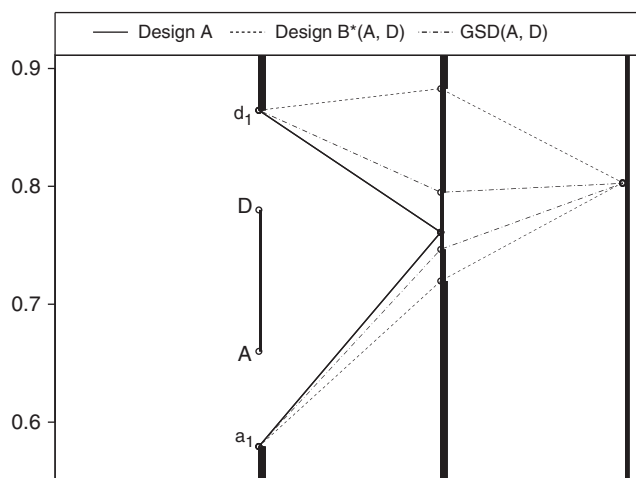
## 2.2. Specification of adaptive design

Next, we consider an adaptive design to increase the power of Design A at various effect sizes  $\theta$ . The basic idea is that at the first interim analysis we might (1) observe results so extreme that the stopping rule defined by Design A would suggest early termination of the study, (2) observe results that did not exceed a stopping boundary, but were close enough to the boundary as to suggest the eventual decision reached at the next analysis (no adaptation), or (3) observe results that were sufficiently far from our expectations that additional data might be desired to increase the power to obtain a statistically significant result (adapt by increasing the maximal sample size and adding an interim analysis). The form of adaptive design that we consider throughout is a prespecified adaptation based on an interim estimate of the effect size, and the statistic used in the specification of the stopping rule is a sufficient (and unweighted) statistic. This differs from the adaptive designs described in, for example, Proschan and Hunsberger [12] and Cui, Hung, and Wang [13], which are not necessarily prespecified and are based on a weighted combination of the independent intervals rather than on a minimal sufficient statistic.

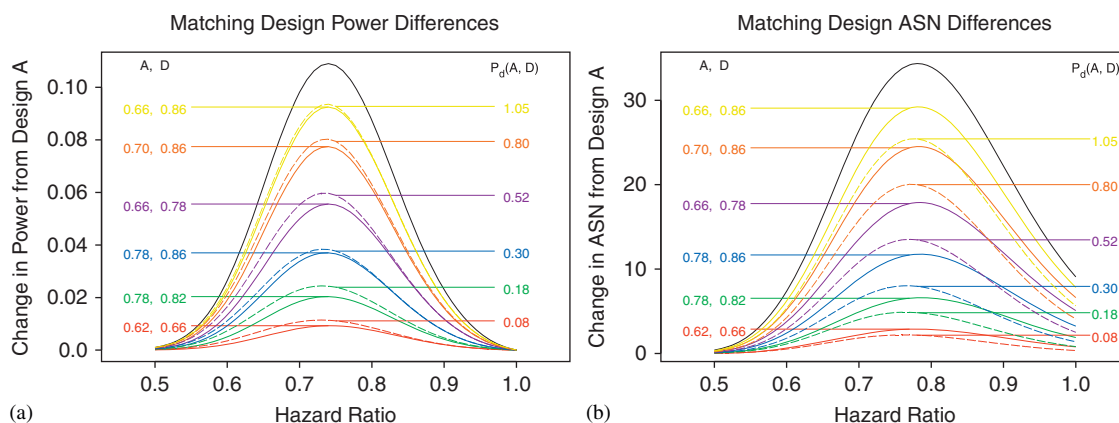
We presume that the stopping boundaries at the first interim analysis were chosen based on the results that would be judged clinically important and statistically credible. Let Design B be an extension of Design A, constrained to have the same boundary at the first analysis, but with three analyses at 100, 200, and 300 events instead of two analyses at 100 and 200 events. The same parameters used in construction of Design A are used to define the boundaries at the 2nd and 3rd analyses for Design B: the design has level  $\alpha=0.025$  with an O'Brien-Fleming efficacy boundary and a Pocock futility boundary. Such an extension is easily obtained using the constrained boundary approach described by Burington and Emerson [14] and implemented in S+SeqTrial, as illustrated in Appendix B.

The adaptive design that we consider performs an analysis at 100 subjects according to Design A. We define parameters  $A$  and  $D$ , with  $a_1 \leq A \leq D \leq d_1$ , to be the values that define the adaptive behavior of the design. Based on the statistic  $T_1$  at the first analysis, the design either proceeds with Design A, or switches to a larger maximum sample size and an additional interim analysis. If  $T_1 \in (a_1, A)$  or  $T_1 \in (D, d_1)$  then Design A boundaries are used for the remainder of the study. If, however,  $T_1 \in (A, D)$ , then the boundaries of Design  $B^*(A, D)$ , a modified version of Design B, are used, and the maximal possible sample size is therefore increased from 200 to 300. In either case, we continue to use the same sufficient statistic, the estimated hazard ratio, as the test statistic at all analyses. Design  $B^*(A, D)$  is modified from Design B to maintain the overall experiment-wise type I error. Design  $B^*(A, D)$  has the same specifications and constraints as Design B, except that the type I error  $\alpha_B(A, D)$  is calculated, as shown in Appendix A, based on the values of  $A$  and  $D$  to guarantee the desired experiment-wise error  $\alpha=0.025$  for the adaptive procedure. This modified type I error and design can be found using standard sequential design software, as described in Appendix B. Figure 1 illustrates the boundaries corresponding to this adaptive design procedure.

The parameters  $A$  and  $D$  that define the adaptive behavior can be chosen based on several different criteria, such as conditional power or symmetry considerations. Note that if  $A=D$ , there is no adaptation and the adaptive design reduces to Design A, and if  $A=a_1$  and  $D=d_1$ , the adaptive design reduces to the unmodified Design B. Different choices of  $A$  and  $D$  move the power curves and ASN curves of the resulting adaptive designs smoothly between the respective curves for Design A and Design B. In our investigations, we consider a full range of possible  $A$  and  $D$  values to explore the whole space of adaptive designs defined in this way. As noted by the anonymous referees, some of the considered choices of  $A$  and  $D$  result in designs that appear sub-optimal or less than reasonable. However, it is worth noting that when the operating characteristic curves of all of these designs are considered, as shown in Figure 2, each one represents just a different tradeoff between power and ASN between Designs A and B, and none displays a clear loss on both criteria.



**Figure 1.** Stopping boundaries for an example of an adaptive design comprised of Design A and Design  $B^*(A, D)$ , as well as for the power-matched group sequential design  $GSD(A, D)$ .



**Figure 2.** Matched designs  $ASD(A, D)$  and  $GSD(A, D)$  operating characteristics comparison: (a) Change in power from Design A, matched designs shown in same shade with a solid line for  $ASD(A, D)$ , and a dashed line for  $GSD(A, D)$  and (b) Change in ASN from Design A, matched designs shown in same shade with a solid line for  $ASD(A, D)$ , and a dashed line for  $GSD(A, D)$ . Note that we were able to find group sequential designs that have higher power and lower ASN over the range of alternatives that would typically be considered in the design of a clinical trial.

### 2.3. Specification of a nonadaptive group sequential design for comparison

Let  $ASD(A, D)$  denote the adaptive sequential design resulting from a particular choice of  $A$  and  $D$ . For fixed parameters  $A$  and  $D$  we then seek to identify a comparable, nonadaptive group sequential design  $GSD(A, D)$ . The spectrum of group sequential designs is quite extensive and flexible, and, as noted by Tsiatis and Mehta [1] and Jennison and Turnbull [2], it is in general possible to find a group sequential design that matches the efficiency of an adaptive design. Our goal, however, is to determine the extent to which a traditional group sequential design may equally well satisfy the sponsor's constraints, which extend beyond the usual statistical power and sample size considerations. As such, these considerations, rather than error-spending functions, guide our choice of competing group sequential designs to consider. We also explored the 'optimal' matching group sequential designs of Tsiatis and Mehta [1], which match the error-spending function of the adaptive design. While these designs improve on the power across the range of  $\theta$  considered, they also resulted in an increase in ASN for the same values of  $\theta$ . Thus, these designs are difficult to fairly compare to the adaptive design, since they are better in one respect (power) but worse in another (ASN), and increasing ASN generally results in an increase in power. Instead, we attempted to find group sequential designs that had almost

identical power curves to the adaptive design under consideration while simultaneously satisfying the scientific and ethical constraints that guided the choice of the adaptive design parameters.

To that end we presume that the stopping boundaries at the first interim analysis (when  $N_1 = 100$  events) were chosen to guarantee the scientific and statistical credibility of results should the study be terminated early. We further presume that the timing of possible interim analyses was fixed by logistical constraints, and thus consider only group sequential tests having analyses at 100, 200, and 300 events. Because we want to demonstrate the comparability of the group sequential design to the adaptive design, we also restrict attention to those group sequential test designs that have power curves that closely match the power curve of the corresponding adaptive design. We define the power curves

$$\gamma_{\text{ASD}(A,D)}(\vartheta) = P(\text{Reject } H_0: \theta = 1 \text{ using ASD}(A, D) | \theta = \vartheta),$$

$$\gamma_{\text{GSD}(A,D)}(\vartheta) = P(\text{Reject } H_0: \theta = 1 \text{ using GSD}(A, D) | \theta = \vartheta)$$

and we want to minimize

$$\max_{\vartheta \in [0.5, 1]} |\gamma_{\text{ASD}(A,D)}(\vartheta) - \gamma_{\text{GSD}(A,D)}(\vartheta)| \quad (1)$$

subject to

$$\gamma_{\text{GSD}(A,D)}(\vartheta) \geq \gamma_{\text{ASD}(A,D)}(\vartheta) \quad \text{for all } \vartheta \in [0.5, 1]. \quad (2)$$

That is, we want the power curve of  $\text{GSD}(A, D)$  to be as close as possible to, without being less than, the power curve of  $\text{ASD}(A, D)$  for all values  $\theta = \vartheta$  in a range that includes the null and alternative hypotheses  $[0.5, 1]$ . Note that the adaptive design  $\text{ASD}(A, D)$  and the power-matched group sequential design  $\text{GSD}(A, D)$  have the same maximal possible sample size.

Even under the above constraints, there are likely many different group sequential stopping boundaries that could be considered. In searching for a group sequential design that matched the power curve of a given adaptive design  $\text{ASD}(A, D)$ , we found that a specially modified version of Design  $B^*(A, D)$  worked remarkably well in the sense that it tended to match the power of the adaptive designs with increases in power no more than 0.004 and tended to lead to a decrease in ASN. We defined  $\text{GSD}(A, D)$  to be a standard group sequential design with the same boundaries as Design  $B^*(A, D)$  at the first and third analyses ( $N_1 = 100$  events and  $N_3 = 300$  events). Then, we modified the boundary at the second analysis by changing the value of the  $P$  parameter for the design to be zero for the efficacy boundary ( $P_a = 0$ ), and to be some appropriately chosen positive number (in the range 0.05–1.25 for the examples we consider) for the futility boundary ( $P_d = P_d^*(A, D)$ ). The optimal value  $P_d^*(A, D)$  was chosen from a grid of possible values to optimize the criterion in (1) subject to the constraint (2). This modification effectively shrinks the stopping boundary at the second analysis from that of Design  $B^*(A, D)$  toward the boundary of Design A, and the value of  $P_d^*(A, D)$  controls the degree of shrinkage: smaller values of  $P_d^*(A, D)$  result in boundaries closer to Design A, while larger values of  $P_d^*(A, D)$  result in boundaries closer to Design B. An example with further details is provided in Appendix B. Figure 1 displays one example of the resulting group sequential design boundary,  $\text{GSD}(A, D)$ .

We note that these matching group sequential designs display some unusual and perhaps undesirable nonmonotonic boundary behavior for some choices of  $A$  and  $D$ . For instance, the futility boundary of  $\text{GSD}(A = 0.62, D = 0.66)$  with  $P_d^*(A = 0.62, D = 0.66) = 0.08$  is  $d_1 = 0.8645$ ,  $d_2 = 0.7665$ ,  $d_3 = 0.8025$  on the estimated hazard ratio scale. Using this design, the study might stop for futility with an estimated hazard ratio of 0.78 at the second analysis, but at the third analysis this same estimated value would lead to a decision for efficacy. This behavior is not restricted to the matching group sequential designs that we consider: the same phenomenon occurs in the ‘optimal’ group sequential design of Tsiatis and Mehta [1] for this same error-spending function. Such nonmonotonic behavior would seem to indicate that more efficient boundaries are likely possible, and that the adaptive design that is being matched is not as efficient as a better-designed group sequential design might be. However, as noted previously in Section 2.2, to the extent that both Design A and Design B are reasonable designs, all of the adaptive designs could be considered reasonable compromises between the two. Furthermore, when we examine the stopping probabilities for  $\text{GSD}(A = 0.62, D = 0.66)$  across the range of  $\theta \in [0.5, 1]$ , we find that the probability of proceeding to the third analysis is at most 0.0220, with this maximum occurring when  $\theta = 0.76$ . Therefore, this nonmonotonic boundary is very unlikely to have an effect on the outcome of the trial.

In the following section we explore the behavior of this adaptive design and matching group sequential design in the (unrealistic) setting of no censoring, and find no advantage to this adaptive design. With no censoring, the statistical information is proportional to the number of subjects accrued. As such, the same results will hold for comparisons of means or differences of proportions. Then, in the following section we add in censoring and explore the tradeoffs of number of subjects versus calendar time, where we do find instances in which the adaptive design exhibits some advantages over a traditional group sequential design.

### 3. No censoring

In Figure 2 we show power and ASN comparisons for a selection of six adaptive designs and the corresponding matched group sequential designs. As these figures show, the group sequential designs  $GSD(A, D)$  are more efficient than the adaptive designs  $ASD(A, D)$  in terms of ASN, with equal or slightly superior power across the range of true effect size considered. Similar results were obtained as we explored the behavior of the adaptive designs and the group sequential designs over the complete range of possible values of  $A$  and  $D$ . Thus in settings where information is measured by number of subjects, these group sequential designs are observed to be uniformly superior to the corresponding adaptive designs over the range of alternatives that would typically be considered during the design of the study. In explorations not shown here, these results were found to generalize to clinical trial settings using means or binomial proportions. This is in keeping with the findings of Tsiatis and Mehta [1] and Jennison and Turnbull [2].

### 4. Censoring

Censoring occurs when we accrue subjects and only follow them for a certain amount of time, as opposed to following them indefinitely until an event occurs. Generally, the amount of follow-up time is determined by a certain date at which the study ends, at which time subjects who have not yet had an event are censored. Note that in this administrative censoring scenario the follow-up time for individual subjects may differ, depending on when they were accrued to the study.

In a setting with censoring, we now have to consider the costs associated with accrual of subjects and follow-up time. Follow-up time can be reduced by accruing more subjects, and conversely, the number of subjects required may be reduced by extending the follow-up time. Here, the adaptive design has the advantage of allowing accrual to stop earlier when it is determined that the maximum number of events needed is only 200 instead of 300. We explore the behavior of study duration versus the number of subjects required under a variety of accrual patterns and event rates. Following Schoenfeld and Richter [15], as implemented in S+SeqTrial [16], the accrual patterns we consider have a constant number of subjects accrued per time unit, and event rates are modeled as exponential with  $h_T = \theta h_C$ , where  $h_T$  is the hazard rate for the treatment arm,  $h_C$  is the hazard for the control arm, and  $\theta$  is the hazard ratio.

To reduce the dimension of the space of parameters that we consider, we fix the control group event rate to have a median of 1. Note that this reduction still allows us to explore the complete space of accrual patterns and event rates, as it is equivalent to changing the unit of time used. For instance, an accrual rate of 20 patients per month with a median control event time of 6 months is equivalent to an accrual rate of 120 patients per six-month period with a median control event time of 1 six-month period. Similarly, any combination of accrual rate  $r$  and median control event time  $m_C$  expressed in time units  $u$  may also be expressed as an accrual rate of  $r^* = rm_C$  and a median control event time of 1 in time units  $u^* = um_C$ . We choose to explore accrual rates in  $\{40, 60, 100, 150, 200, 250\}$  as a reasonably comprehensive representation of possible accrual scenarios.

To compare the behavior of the adaptive designs  $ASD(A, D)$  to the matching group sequential designs  $GSD(A, D)$  under a particular accrual rate  $r$  and effect size  $\theta$ , we considered the range of possible accrual times and the resulting estimated number of subjects and trial duration. We will assume that if accrual ends before the first analysis then the accrual time  $t$  must have been sufficient to obtain at least 300 subjects to ensure that analyses at 300 events will be possible (otherwise, if fewer than 300 subjects were accrued and if Design  $B^*(A, D)$  boundary were selected at the first analysis, it could prove necessary to restart accrual—a practice that is generally avoided). In such a case, the adaptive design clearly offers no benefit of curtailed accrual, and is therefore slightly less efficient than the



matching group sequential design. In the setting we consider here, accrual rates higher than 250 subjects per unit time were not explored, as they tend to result in accrual ending before the first analysis. For the adaptive designs, if accrual continues beyond the first analysis, we must consider two independent accrual times  $t_A$  and  $t_B$  depending on which boundary is adaptively chosen for the later analyses. These accrual times  $t_A$  and  $t_B$  are independent parameters that may be chosen by the investigators, subject to the constraints discussed below. Accrual time  $t_A$  is the total time that subjects will be accrued if Design A is adaptively selected, and similarly accrual time  $t_B$  is the total time that subjects will be accrued if Design B\*(A, D) is adaptively selected. Each combination of  $t_A$  and  $t_B$  produces an estimated number of subjects and trial duration, upon which the costs of a particular design ultimately depend. For a given adaptive design defined based on the number of events, different choices of  $t_A$  and  $t_B$  might be preferable depending on the relative contributions of per subject costs and calendar time costs to the total cost of the clinical trial.

The factors involved in choosing accrual times  $t_A$  and  $t_B$  for the adaptive design are:

- Each of  $t_A$  and  $t_B$  are constrained to be larger than the time of the first analysis. We consider the case of accrual finishing before the first analysis separately.
- $t_A$  and  $t_B$  are required to be large enough to obtain at least 200 and 300 subjects respectively, in order to ensure that we will be able to observe the necessary number of events.
- The maximum accrual time considered for each design was chosen to be the accrual time that resulted in zero follow-up time after the end of accrual for that design. Accrual times greater than this maximum would be pointless, as this would mean accruing subjects after the study was finished.
- Since accrual is fixed such that it never ends before the first analysis, we can freely decide which combination of accrual times we will use to achieve the adaptive design. Therefore, we consider all the possible combinations of accrual times  $t_A$  for Design A and accrual times  $t_B$  for Design B.

We explored a range of values for each of  $t_A$  and  $t_B$  subject to the above constraints. For each combination of  $t_A$  and  $t_B$ , we computed the estimated study duration and estimated number of subjects as follows. Define the following quantities:

$r$  = rate of accrual

$p_1$  = probability of stopping at the first analysis

$p_{2A}$  = probability of stopping at the second analysis, using Design A

$p_{2B}$  = probability of stopping at the second analysis, using Design B\*(A, D)

$p_3$  = probability of stopping at the third analysis

$\tau_1(t_A) = \tau_1$  = estimated time of the first analysis

$\tau_{2A}(t_A) = \tau_{2A}$  = estimated time of the second analysis for Design A

$\tau_{2B}(t_B) = \tau_{2B}$  = estimated time of the second analysis for Design B\*(A, D)

$\tau_3(t_B) = \tau_3$  = estimated time at the third analysis

$S$  = number of subjects accrued

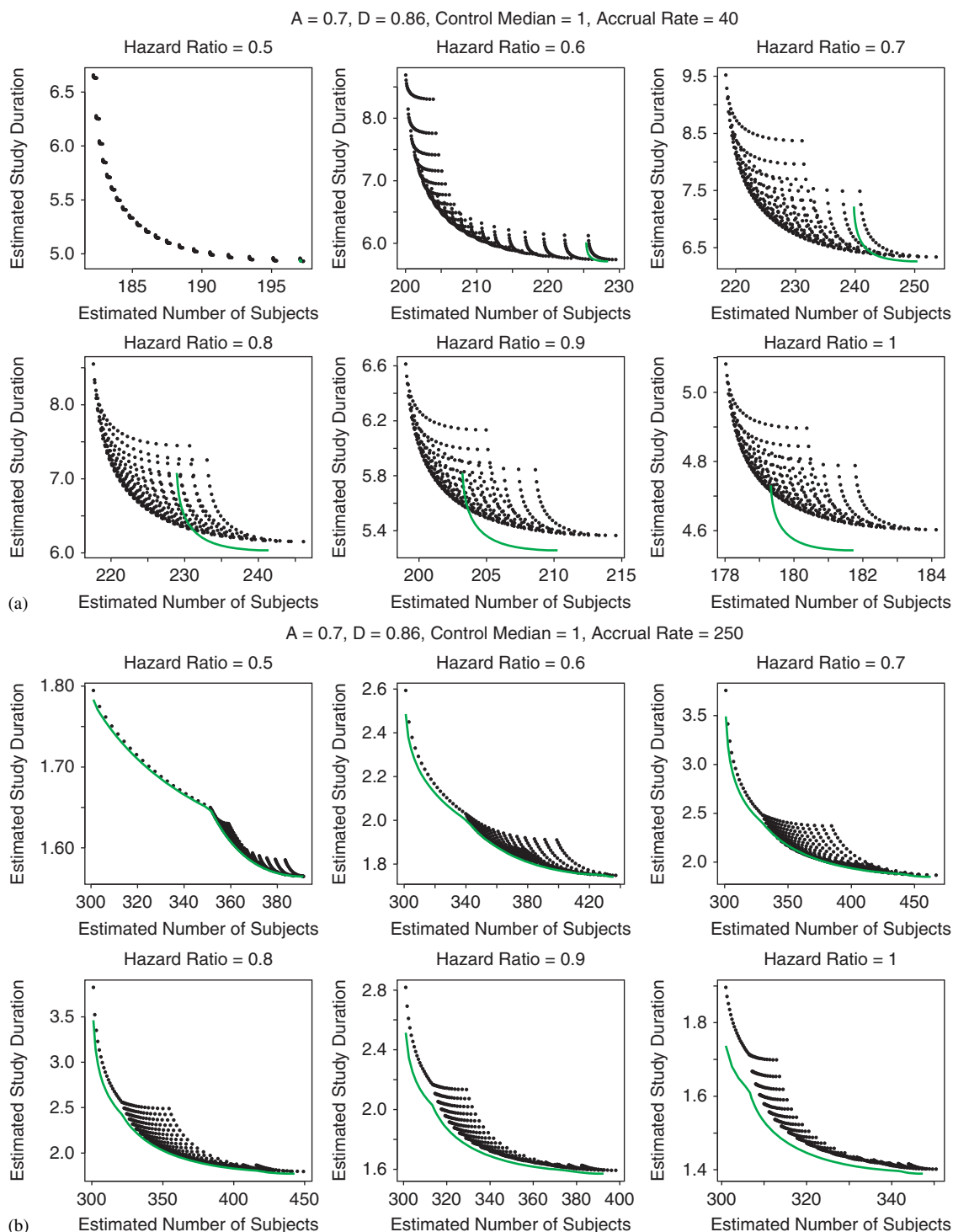
$T$  = total study duration.

Then, the expected number of subjects  $S$ , and the expected trial duration  $T$ , for the adaptive design with accrual times  $t_A$  and  $t_B$  are given by:

$$E[S] = [p_1 \times \tau_1 \times r] + [p_{2A} \times t_A \times r] + [p_{2B} \times \min(\tau_{2B}, t_B) \times r] + [p_3 \times t_B \times r],$$

$$E[T] = [p_1 \times \tau_1] + [p_{2A} \times \tau_{2A}] + [p_{2B} \times \tau_{2B}] + [p_3 \times \tau_3].$$

Each panel in Figure 3 illustrates examples of plots resulting from these calculations. Each dot in the figures represents the expected number of subjects and expected trial duration resulting from one combination of  $t_A$  and  $t_B$ . Points that correspond to a constant value of  $t_A$  tend to lie on smooth curves



**Figure 3.** Number of subjects versus study duration for the adaptive design with  $A=0.70$ ,  $D=0.86$  and the matching group sequential design, under two different accrual rates. The points represent adaptive design results, and the solid line represents the group sequential design results: (a) Accrual rate = 40 subjects per unit time and (b) Accrual rate = 250 subjects per unit time.

that are predominantly vertical, and points that correspond to a constant value of  $t_B$  tend to lie on smooth curves that are predominantly horizontal.

The individual panels in Figure 3 demonstrate the results for the adaptive design with  $A=0.70$ ,  $D=0.86$ , for accrual rates of 40 and 250 in panels of Figures 3(a) and (b), respectively. In each figure, the dots correspond to results for the adaptive design, and the solid line corresponds to the matching

group sequential design. The number of subjects versus study duration curves for the matching group sequential designs  $GSD(A, D)$  were similarly obtained by considering a range of accrual times  $t_G$ , subject to constraints analogous to those that governed the choices of  $t_A$  and  $t_B$ .

For the adaptive design and accrual rate of 40 in Figure 3(a), there is a potential for benefit using the adaptive design over the group sequential design, particularly for hazard rates less than  $\theta=0.8$ . With this slow accrual rate, there is a limited range for  $t_B$  (and also for  $t_G$ ), and thus the dots corresponding to the same value of  $t_B$  are very close together. The lines corresponding to the group sequential design are at the far right end of the plots, demonstrating the benefit of accrual modification in Design A of the adaptive design, for this scenario. The adaptive design allows the possibility of reducing the number of subjects required by 10–30 depending on the effect size  $\theta$ . Of course, there is a tradeoff of increasing the trial duration, but a moderate reduction in number of subjects does not produce a dramatic increase in the study time. For instance, at a hazard ratio of  $\theta=0.7$ , the number of subjects can be reduced from a minimum of 240 for the group sequential design to 230 for the adaptive design, while only increasing the expected study duration from 6.5 to 6.75 time units. As the hazard rate increases, the benefit of the adaptive design becomes less pronounced, and eventually disappears.

At an accrual rate of 250 subjects per time unit (median survival time on the control arm), the same adaptive design shows no potential for subject reduction. The curves for the group sequential design and adaptive design become quite close, with the group sequential design dominating particularly at hazard rates close to the null hypothesis  $\theta=1$ . For  $\theta$  in the range 0.8–1, the same group sequential design achieves a shorter average duration than the adaptive design for the same average number of subjects, across the entire range of possible number of subjects.

Clinical trialists would choose from among the spectrum of adaptive designs considered based on the efficiency and power curves desired. We present the example in Figure 3 to demonstrate the patterns of behavior resulting from the range of adaptive designs and accrual scenarios, with similar trends observed for other choices of  $A$ ,  $D$ , and accrual rates. The following general trends were observed: As accrual rate increases, the difference between an adaptive design and the corresponding group sequential design tends to disappear, with the group sequential design tending to be slightly more efficient. For lower accrual rates, the ability of either the adaptive design or the group sequential design to improve upon the other, in terms of expected trial duration at a given expected number of subjects, will depend on the choices of  $A$ ,  $D$ , and the effect size. More generally, there are tradeoffs between the adaptive designs and the matching group sequential designs that depend on the relative importance of minimizing the number of subjects versus minimizing trial duration. In the following section, we attempt to explore these tradeoffs.

## 5. Cost estimation

In order to explore the tradeoffs between increased sample size and decreased study duration, we consider the cost to the sponsor using a simple discrete time economic model. We presume the setting of the design of a Phase III clinical trial. At the start of the trial, the sponsor will have incurred costs related to treatment development and early-phase clinical trials, and there are costs to the sponsor associated with the money invested in that development program. For instance, the cost of prior development might total to \$10 million. Then, the Phase III trial might engender costs on the order of, say, \$10 000 per patient. In our simple model, we consider the cost of that prior investment, as well as the cost of the patients accrued. We then further consider the cost of study duration by allowing for interest to be paid by the sponsor on its investment. Letting  $n_t$  represent the number of subjects accrued between time  $t-1$  and time  $t$  and letting  $p$  be the per patient costs, we can then calculate the total cost  $C(t)$  up to calendar time  $t$  as  $C(t)=n_t \times p + (1+\omega) \times C(t-1)$ , where the interest rate  $\omega$  is the cost of money per unit time. Without loss of generality, it is sufficient for us to consider merely the ratio  $C(0)/p$ . The interest rate is used to represent the cost of time, i.e. study duration, and may serve as a surrogate for all time-related costs such as the expenses related to maintaining databases, personnel, and borrowing money. The ratio of the prior costs to the per-patient costs determines the direction of the tradeoff between trial duration and number of subjects. For relatively higher per-patient costs, the total trial cost is minimized when fewer patients are used and a longer study duration is permitted.

With this cost model, for each design and each accrual scenario, we can calculate the ratio of the optimal trial cost for  $GSD(A, D)$  to the optimal trial cost for  $ASD(A, D)$  for a range of  $\theta$  values. As an example, we consider a time cost of money based on  $\omega=0.005$ , which if the unit of time is

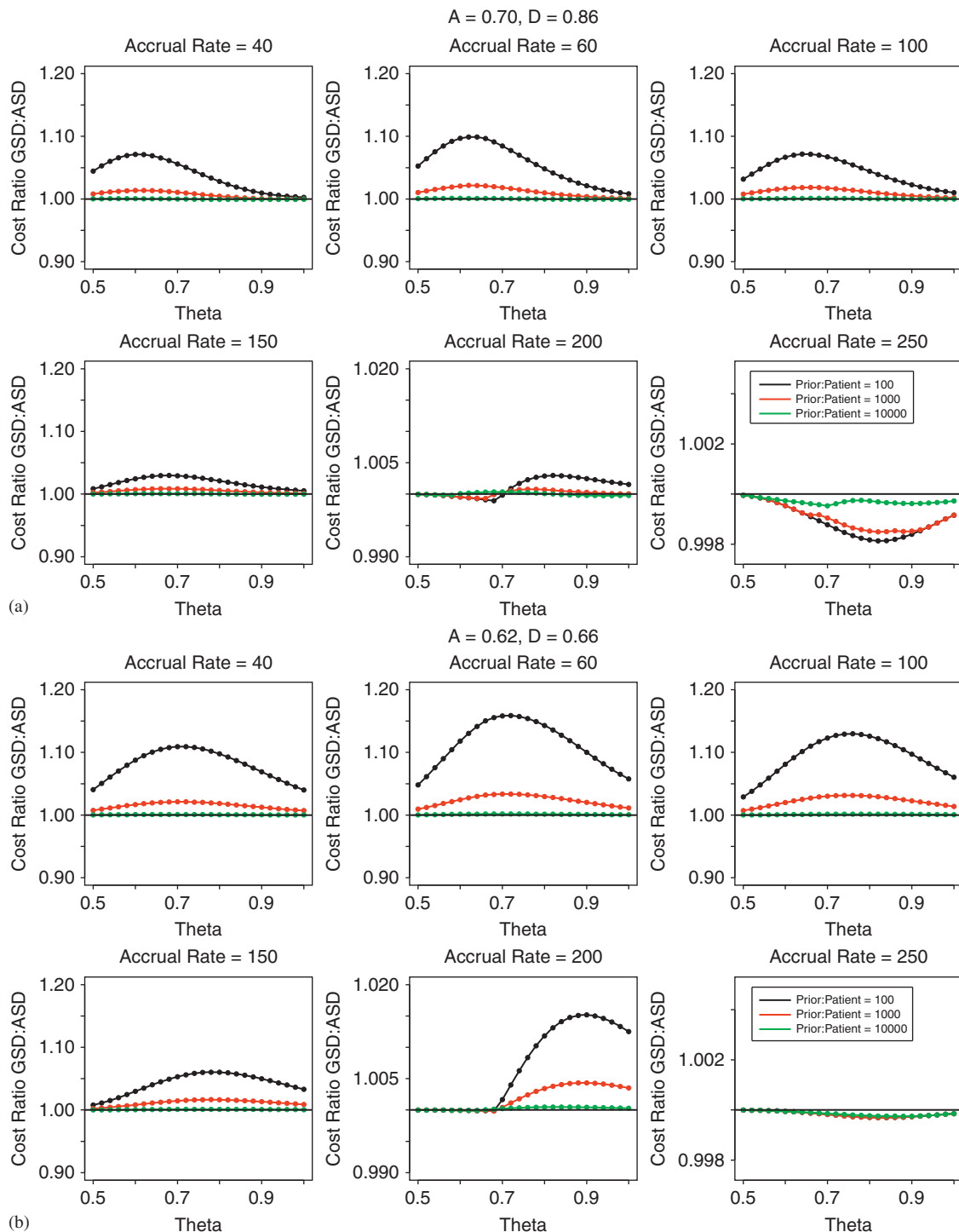
one month would correspond to a 6 per cent yearly interest rate. Higher interest rates tend to favor shorter trial duration, which would cause the group sequential designs to be more advantageous, as indicated in Figure 3. The resulting cost ratio plots at this interest rate for two of the possible adaptive designs are shown in Figure 4. As demonstrated by these plots, the adaptive design typically offers some cost improvements over the group sequential design for slow accrual rates, but as the accrual rate increases the cost improvement disappears. Figure 4(a) shows the cost ratio plots corresponding to an adaptive design with  $A=0.70$ ,  $D=0.86$ , the adaptive design of Figure 3. A range of accrual rates are considered, with prior to patient cost ratios of 100, 1000, and 10 000. We note that an anonymous referee commented that the prior to patient cost ratio may well exceed 10 000, which we acknowledge. The general trend presented here indicates that as the prior to patient cost ratios increase, the ratio of the adaptive to group sequential trial costs will tend toward 1. As the ratio of prior costs to per patient costs increases with the simple economic model considered here, the cost difference, which may be of greater importance than the cost ratio, tends to favor designs which minimize calendar time regardless of the sample size.

The most significant benefit of the adaptive design is seen for mid-range values of the true effect size ( $\theta \in (0.65, 0.8)$ ), when the patient costs are high (at a prior cost to per-patient cost ratio of 100), and for accrual rates near the low end of the spectrum (40–100 patients per unit time). In this range, the group sequential design may be 5–10 per cent more expensive than the matched adaptive design. However, when accrual rates are high, the group sequential design is actually very slightly more cost effective than the adaptive design, saving a small fraction of a percent over the adaptive design. Note that the y-axis scale in these plots changes for the two highest accrual rates in Figure 4, and also that the rougher behavior observed in these two high accrual rate scenarios is due to both the change-points observed in the subject versus duration plot of Figure 3 and the increased resolution. Figure 4(b) shows the cost ratio plots summarizing an adaptive design with  $A=0.62$ ,  $D=0.66$ . The results for this design are similar to those of Figure 4(a), though the cost benefit of the adaptive design is somewhat increased. Again, for the highest accrual rate of 250 subjects per unit time, the group sequential design offers a very slight improvement over the adaptive design when looking at the ratio of total costs. Depending upon the total prior costs that small advantage in relative costs may translate into a substantial additive benefit.

## 6. Discussion

In our comparisons considered here, we compared adaptive designs of a form similar to those initially proposed for an industry-sponsored clinical trial to traditional group sequential designs that might have had the same operating characteristics. There are many parameters that can be considered in a group sequential stopping rule including the type I error, the power under some design alternative, the number of interim analyses, the relative timing of the interim analyses, and boundary shape functions for each of the decisions that might be reached. The boundary shape functions can in turn be defined for any one of several different statistics: partial sum of (potentially transformed) observations, the maximum likelihood estimate, the standardized  $Z$  statistic, the fixed sample  $P$  value, the error-spending function, the conditional power under some hypothesized treatment effect, the Bayesian predictive power under some prior distribution for the true treatment effect, or the Bayesian posterior probability of some hypothesis. For each of these statistics, the boundary shape function relates the early conservatism with which the boundary would allow termination of the study at the earliest analyses. In the unified family of Kittelson and Emerson [5], a user may choose from a broad spectrum of boundary shape functions through the choice of three parameters that can be chosen separately for each stopping boundary.

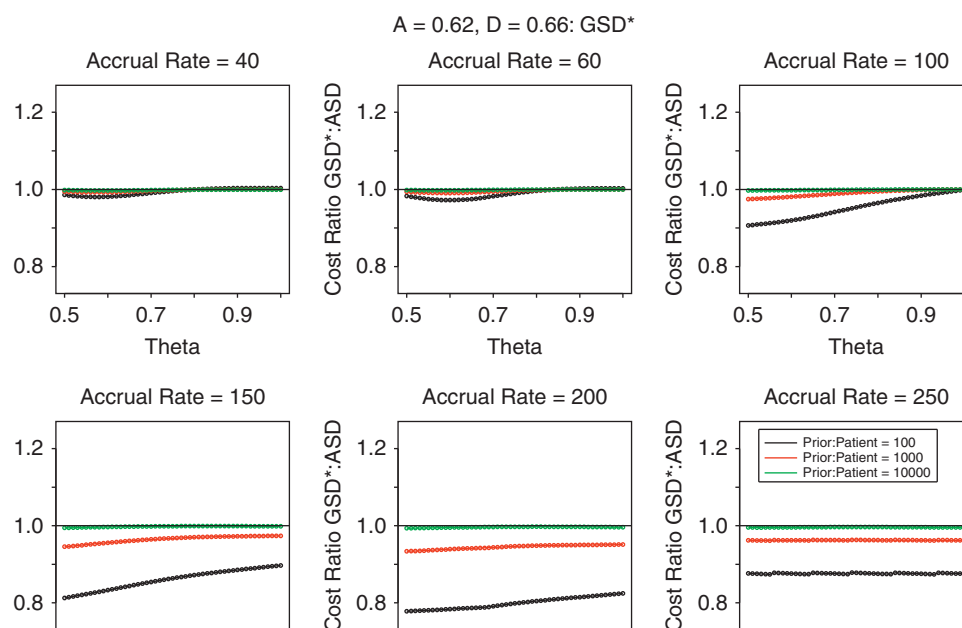
This wide flexibility of group sequential stopping rules means that there are likely many different group sequential designs with the same power curves and ASN curves, for instance. Hence, when evaluating the ability of adaptive designs to improve on standard group sequential methods, we must ensure that we understand the design constraints that are to be held constant, and those that are allowed to vary. In the investigations presented in this manuscript, we presumed that we needed to maintain the criteria for stopping at the earliest analyses, as well as the schedule and timing of interim analyses. As part of our interest was to see how easily we could match the adaptive designs, we considered only a single method of modifying the group sequential boundaries at the second analysis in order to achieve a comparable group sequential design.



**Figure 4.** Cost ratio plots (a) for the adaptive design considered in Figure 3, and (b) for another adaptive design with different values of  $A$  and  $D$ , for a range of accrual rates. Cost ratio is calculated as the ratio of the minimal cost for the group sequential design to the minimal cost for the adaptive design, where costs are estimated as described in Section 5. Three different prior to patient cost scenarios are considered, representing high patient costs (black line), mid-range patient costs, and low patient costs. Note that the y-axis scale changes for the two highest accrual rates.

While it is clear that the adaptive designs offer no advantages in uncensored settings such as when evaluating the difference of means or proportions, there are some distinct advantages to the adaptive design in certain survival analysis settings, because we may gain efficiency in the number of subjects





**Figure 5.** Cost ratio plots comparing a less-constrained, optimized group sequential design GSD\*(0.62, 0.66) to the adaptive design ASD(0.62, 0.66), for a range of accrual rates. Cost ratio is calculated as the ratio of the minimal cost for the group sequential design to the minimal cost for the adaptive design, where costs are estimated as described in Section 5. Three different prior to patient cost scenarios are considered, representing high patient costs (black line), mid-range patient costs, and low patient costs.

needed and/or the calendar time required for the study to complete. Such considerations will also be relevant in trials with a delayed response or in longitudinal studies, as observed by an anonymous referee. We found that the degree of benefit depends on the distributions of event times and accrual rate as well as on the particular adaptive design under consideration. It is therefore worth considering the cost-effectiveness of using such an adaptive design in time-to-event endpoints. Even the simple cost model considered here, when used with trial-specific values of prior development costs, per patient costs, accrual rates, and the current interest rates reflecting the time cost of money, could provide useful insight into the tradeoffs between potentially higher patient accrual or longer calendar time.

It is worth noting, however, that our decision to hold the number and schedule of interim analyses constant may represent an unreasonable restriction. To briefly explore the effects of relaxing these constraints, we consider the example with  $A=0.62$ ,  $D=0.66$ , where the adaptive design appeared to offer the most potential for benefit. We relaxed the constraint on the timing of the second and third analyses to explore a broader class of group sequential designs, but we continued to enforce the stopping boundary for the first analysis at 100 events. Having searched across a range of possible maximal sample sizes and  $P=(P_a, P_d)$  parameters to find an improved group sequential design within these relaxed constraints, we found that a design with maximum sample size of 210 events (analyses at 100, 155, and 210 events) and  $P=(1.3, 1.3)$  matched the power curve of ASD(0.62, 0.66) while dramatically improving ASN. The cost ratio plots resulting from comparing this design to ASD(0.62, 0.66) are shown in Figure 5, from which we can see that significant reductions in total trial cost are possible for certain accrual and cost scenarios.

We do acknowledge that the above exercise is not totally fair. We presumed that the adaptive person tipped his/her hand first. Thus, we only had to show we could improve over their choice. In the context of this example, which was based loosely on the type of design proposed for an industry sponsored study, we found that we could easily find a group sequential design that met the same general operating characteristics. Given the small number of operating characteristics that were actually considered relative to the large number of group sequential test parameters at our disposal, this is not surprising. It would be similarly unsurprising to find that a proponent of adaptive designs could match the specified constraints of any particular group sequential design, and perhaps improve on some others. However, we believe the advantages of the well-developed group sequential trial theory makes it advantageous to use the group sequential design whenever the two approaches are roughly comparable.

As noted by Emerson [6], there remain problems with inference following the use of such an adaptive design, so in cases where there is questionable or insignificant gain from the adaptive design it may be wiser to continue to use a standard group sequential design where inferential methods are readily available in commercially available statistical software. Thus, we would argue that the time of clinical trialists is probably better spent exploring the wide range of group sequential trials already described and implemented, rather than trying to find *ad hoc* adaptive designs. Our belief is that the careful evaluation of candidate group sequential designs can largely address the issues that have motivated the development of adaptive designs.

One such area of evaluation that we have not explored here, but one that should receive careful attention in a time-to-event setting, is that of the ability to assess time-varying treatment effects: In the setting of treatment effects that might be of greater magnitude either soon after randomization or after some delay, the tradeoffs between sample size and calendar time take on great importance. A study that terminates early with most events corresponding to short periods of treatment may not detect a clinically important difference in treatment behavior with additional follow-up. Hence, further scientific and statistical evaluation of the impact of patient accrual and calendar time issues are also important to address when considering adaptively changing accrual patterns.

## Appendix A: Theory

We consider the adaptive switching from a prespecified group sequential Design A to a prespecified group sequential Design  $B^*(A, D)$  according to whether the test statistic  $T_1$  computed at the first analysis is between the values of  $A$  and  $D$ .

Notationally, we define group sequential Design A as a level  $\alpha$  one-sided test of a lesser alternative having continuation sets  $\mathcal{C}_1 = (a_1, d_1)$  and  $\mathcal{C}_2 = \emptyset$  for  $T_1$  and  $T_2$ , respectively, computed at analyses performed when the accrued sample sizes are  $N_1$  and  $N_2$ , respectively, where  $N_1 < N_2$ . The threshold  $a_2$  for statistical significance at the second analysis is defined to guarantee an experiment-wise error of  $\alpha$ . Hence

$$\alpha = P(T_1 \leq a_1 | \theta = 1) + P(a_1 \leq T_1 \leq d_1, T_2 \leq a_2 | \theta = 1).$$

Now, suppose that if we do not terminate the clinical trial at the first analysis, we want to switch to an alternative stopping rule whenever  $T_1$  is observed between prespecified values of  $A$  and  $D$  satisfying  $a_1 \leq A \leq D \leq d_1$ . If  $a_1 < T_1 < A$  or  $D < T_1 < d_1$ , we will continue to use the sampling plan that specified a maximal sample size of  $N_2$ , with a threshold for statistical significance of  $a_2$  at that last analysis. Let  $\alpha_A(A, D)$  be the probability of stopping for efficacy using the Design A portion of the adaptive design, under the null hypothesis:

$$\alpha_A(A, D) = P(T_1 \leq a_1 | \theta = 1) + P(D \leq T_1 < d_1, T_2 \leq a_2 | \theta = 1) + P(a_1 < T_1 \leq A, T_2 \leq a_2 | \theta = 1). \quad (A1)$$

Based on the prespecified values of  $A$  and  $D$ , we further prospectively identify a group sequential Design  $B^*(A, D)$  having continuation sets  $\mathcal{C}_2^* = (a_2^*, d_2^*)$  and  $\mathcal{C}_3^* = \emptyset$  for test statistics  $T_2^*$  and  $T_3^*$ , respectively, computed at analyses performed when the accrued sample sizes are  $N_2^*$  and  $N_3^*$ , respectively, with  $N_1 < N_2^* < N_3^*$ . Values of  $T_2^* \leq a_2^*$  or  $T_3^* \leq a_3^*$  will be judged cause to reject the null hypothesis. Hence, we need to pre-specify values of  $a_2^*$ ,  $d_2^*$ , and  $a_3^* = d_3^*$  that, when used in conjunction with the group sequential Design A and the adaptation prespecified through the choice of  $A$  and  $D$ , will preserve the experiment-wise error of  $\alpha$ :

$$\begin{aligned} \alpha &= P(\text{Reject } H_0 | \theta = 1) \\ &= P(T_1 \leq a_1 | \theta = 1) + P(D \leq T_1 < d_1, T_2 \leq a_2 | \theta = 1) + P(a_1 < T_1 \leq A, T_2 \leq a_2 | \theta = 1) \\ &\quad + P(A < T_1 < D, T_2^* \leq a_2^* | \theta = 1) + P(A < T_1 < D, a_2^* < T_2^* < d_2^*, T_3^* \leq a_3^* | \theta = 1) \end{aligned} \quad (A2)$$

Substituting the result of equation (A1) into equation (A2), we therefore have

$$\alpha = \alpha_A(A, D) + P(A < T_1 < D, T_2^* \leq a_2^* | \theta = 1) + P(A < T_1 < D, a_2^* < T_2^* < d_2^*, T_3^* \leq a_3^* | \theta = 1). \quad (A3)$$

Thus, rearranging equation (A3), we only need to find  $a_2^*$ ,  $d_2^*$ , and  $a_3^*$  to satisfy

$$\alpha - \alpha_A(A, D) = P(A < T_1 < D, T_2^* \leq a_2^* | \theta = 1) + P(A < T_1 < D, a_2^* < T_2^* < d_2^*, T_3^* \leq a_3^* | \theta = 1). \quad (A4)$$

In particular, we can define a group sequential Design  $B_{comp}(A, D)$  using the constrained boundary approach of Burington and Emerson (2003) in which analyses are performed at sample sizes  $N_1$ ,  $N_2^*$ , and  $N_3^*$ , the continuation set  $(a_1^*, d_1^*)$  at the first analysis is constrained to be  $a_1^* = A$  and  $d_1^* = D$ , and  $a_2^*$ ,  $d_2^*$ , and  $a_3^* = d_3^*$  can be any values such that the resulting group sequential design has type I error of

$$\begin{aligned}\alpha_B(A, D) &= P(\text{Reject } H_0 \text{ using Design } B_{comp}(A, D) | \theta = 1) \\ &= P(T_1 \leq A | \theta = 1) + P(A < T_1 < D, T_2^* \leq a_2^* | \theta = 1) + P(A < T_1 < D, a_2^* < T_2^* < d_2^*, T_3^* \leq a_3^* | \theta = 1) \\ &= P(T_1 \leq A | \theta = 1) + \alpha - \alpha_A(A, D),\end{aligned}\quad (A5)$$

where the last equality follows from equation (A4). The stopping boundaries  $a_2^*$ ,  $d_2^*$  and  $a_3^* = d_3^*$  of any such Design  $B_{comp}(A, D)$  will thus preserve an experiment-wise error of  $\alpha$  for the adaptive procedure, and we will refer to the group sequential design with these boundaries at the second and third analyses and Design A boundary at the first analysis as Design  $B^*(A, D)$ .

It should be noted that the above approach will be valid regardless of the group sequential design family that is used to parameterize the specification of group sequential Designs A and  $B^*(A, D)$ . Hence, subject to the specification representing a valid design, the group sequential designs could be specified in the unified family of Kittelson and Emerson (1999), a family of error-spending functions, a specification of Bayesian posterior probabilities, or conditional or predictive power families. In this manuscript, we have chosen to specify Design  $B^*(A, D)$  using the same parameterization as was used for Design A; namely, to maintain the O'Brien-Fleming [9] efficacy boundary and the Pocock [8] futility boundary, subject to the constraints on the first analysis boundaries.

## Appendix B: Implementation

The operating characteristics of the adaptive design considered here can be calculated using standard group sequential software, so long as the software allows the specification of arbitrary boundaries and the calculation of stopping probabilities at each analysis. The following approach will work in the program S+SeqTrial.

Design A can be computed according to the specified design parameters, which might include specifying the desired type I error, the number of analyses, the boundary shape parameters for both efficacy and futility, and any two of the design alternative, the desired statistical power, and the maximal number of events. For example, using the unified family of [5], the following code specifies Design A to be a one-sided level 0.025 test of a lesser hazard ratio having a maximum of two analyses after 100 and 200 events have been observed and using an O'Brien-Fleming efficacy (lower) boundary (so boundary shape parameters of  $P_a = 1$ ,  $R_a = 0$ ,  $A_a = 0$ ) and a Pocock futility (upper) boundary (with boundary shape parameters of  $P_d = 0.5$ ,  $R_d = 0$ ,  $A_d = 0$ ), and also specifies Design B with three analyses at 100, 200, and 300 events using the same parameters, constrained to match Design A at the first analysis:

```
> designA <- seqDesign(prob.model="hazard", test.type="less",
  size=0.025, sample.size=c(100, 200), power=0.975, nbr.analyses=2,
  P=c(1, 0.5))
> bou <- seqBoundary(designA)
> bou <- matrix(NA, 3, 4)
> bou[1,] <- seqBoundary(designA)[1,]
> designB <- update(designA, sample.size=c(100, 200, 300), nbr.
  analyses=3, exact.constraint=bou)
```

The actual stopping boundaries are printed with the command:

```
> designA
```

PROBABILITY MODEL and HYPOTHESES:

Two arm study of censored time to event response variable

Theta is hazard ratio (Treatment : Comparison)

One-sided hypothesis test of a lesser alternative:

Null hypothesis : Theta >= 1 (size = 0.025)

Alternative hypothesis : Theta <= 0.5596 (power = 0.975)

```

STOPPING BOUNDARIES: Sample Mean scale
              a          d
Time 1 (N= 100) 0.5792 0.8645
Time 2 (N= 200) 0.7611 0.7611

```

Then, a modification of Design A, Design  $A_{comp}(A, D)$ , is specified in order to assist in computations reflecting the adaptive switching from Design A to a modification of Design B. For specified  $A$  and  $D$ , we modify the boundary of Design A to allow stopping and switching to the modified Design B if the observed hazard ratio is between  $A$  and  $D$  at the first analysis. We make use of the facility for constrained boundaries (Burlington and Emerson, 2003). For instance, if we choose  $A=0.62$  and  $D=0.66$ ,

```

> bouA <- seqBoundary(designA)
> bouA[1,2] <- 0.62
> bouA[1,3] <- 0.66
> designAcomp.AD <- update(designA, test.type="two.sided", exact.
  constraint=bouA)

```

The boundaries of the Design  $A_{comp}(A, D)$  are obtained as:

```

> seqBoundary(designAcomp.AD)

STOPPING BOUNDARIES: Sample Mean scale
              a      b      c      d
Time 1 (N= 100) 0.5792 0.62 0.66 0.8645
Time 2 (N= 200) 0.7611 NA  NA  0.7611

```

Stopping probabilities computed under Design  $A_{comp}(A, D)$  then reflect the probabilities of decisions made using Design A. For instance, under the null hypothesis of a hazard ratio of 1.0:

```

> seqOC(designAcomp.AD, theta=1)

```

```

Operating characteristics at theta= 1
ASN= 121.9633
Lower Power= 0.0218
Upper Power= 0.9677

```

```

Stopping Probabilities:
              Lower Null Upper Total
Analysis time 1 0.0032 0.0105 0.7668 0.7804
Analysis time 2 0.0187 0.0000 0.2010 0.2196

```

From the above, we see that for these values of  $A$  and  $D$  under the null hypothesis there is a probability of 0.0032 of stopping at the first analysis (when 100 events have been observed) with a decision for efficacy, a probability of 0.7668 of stopping at the first analysis with a decision for futility, a probability of 0.0187 of staying with Design A and then deciding for efficacy at the second analysis (when 200 events have been observed), and a probability of 0.2010 of staying with Design A and then deciding for futility at the second analysis. The remaining probability of 0.0105 corresponds to deciding to switch to Design  $B^*(A, D)$  based on the observation of an estimated hazard ratio between 0.62 and 0.66.

Design  $B^*(A, D)$  is a modification of Design B found in such a way as to ensure the experiment-wise type I error of 0.025. There are an infinite number of ways to proceed. For the purposes of this manuscript, we considered maintaining the parameterization of the boundary shapes within the unified family, but constraining the stopping boundaries at the first analysis to agree with the boundaries of Design A, as discussed previously. In order to perform the necessary computations, we must specify a design, Design  $B_{comp}(A, D)$  that will have the same boundary as Design  $B^*(A, D)$ , except that the continuation region at the first analysis will be defined by  $A$  and  $D$  rather than matching the Design A boundary. The operating characteristics of Design  $B_{comp}(A, D)$  will then correctly reflect the probabilities resulting from adaptively switching to Design  $B^*(A, D)$ . Design  $B_{comp}(A, D)$  was specified to guarantee the experiment-wise error, which is the sum of the probability of 0.0218 of deciding for efficacy using Design A( $A, D$ ) plus the probability of deciding for efficacy either at the second or third analysis using Design  $B_{comp}(A, D)$ . Therefore, the specified size for Design  $B_{comp}(A, D)$  should equal the desired experiment-wise type I error of 0.025 minus the probability of 0.0218 for declaring efficacy

when using Design  $A_{\text{comp}}(A, D)$  plus the probability that the statistic at the first analysis is less than  $A=0.62$  (found to be 0.0084), as described by equation (A5) in Appendix A. In the example presented in this Appendix, a type I error of  $0.0250-0.0218+0.0084=0.0116$  was found to be the necessary type I error for Design  $B_{\text{comp}}(A, D)$ . Hence, we used code:

```
> bouB <- matrix(NA, 3, 4)
> bouB[1,1] <- 0.62
> bouB[1,4] <- 0.66
> bouB <- seqBoundary(bouB)
> designBcomp.AD <- update(designA, size=0.0116, nbr.analyses=3,
  sample.size=c(100, 200, 300), exact.constraint=bouB)
```

The stopping boundaries when using Design  $B_{\text{comp}}(A, D)$  are found to be:

```
> seqBoundary(designBcomp.AD)

STOPPING BOUNDARIES: Sample Mean scale
              a          d
Time 1 (N= 100) 0.6200 0.6600
Time 2 (N= 200) 0.7283 0.9386
Time 3 (N= 300) 0.8095 0.8095
```

We can verify the experiment-wise error for the design resulting from adaptively switching from Design A to Design  $B^*(A, D)$  by using the operating characteristics of Design  $A_{\text{comp}}(A, D)$  and Design  $B_{\text{comp}}(A, D)$ :

```
> seqOC(designBcomp.AD, theta=1)

Operating characteristics at theta= 1
ASN= 101.8519
Lower Power= 0.0116
```

```
Stopping Probabilities:
              Lower Null    Upper    Total
Analysis time 1 0.0084      0  0.9811  0.9895
Analysis time 2 0.0018      0  0.0006  0.0024
Analysis time 3 0.0014      0  0.0067  0.0081
```

In the adaptive design based on Design A and Design  $B^*(A, D)$  with parameters  $A=0.62$ ,  $D=0.66$ , we thus find an experiment-wise error of 0.025: A probability of deciding for efficacy of 0.0032 at the first analysis and 0.0187 at the second analysis when using Design A and a probability of 0.0018 at the second analysis and 0.0014 at the third analysis when using Design  $B^*(A, D)$ . Using similar computations of stopping probabilities under the design alternative of a hazard ratio of 0.5596, we find an experiment-wise power of 0.9757: A probability of deciding for efficacy of 0.5684 at the first analysis and 0.3079 at the second analysis when using Design A and a probability of 0.0970 at the second analysis and 0.0024 at the third analysis when using Design  $B^*(A, D)$ . The average sample size (ASN) for the adaptive design can be found by multiplying the number of events at a given analysis time by probability of stopping at that analysis time, for each of Design  $A_{\text{comp}}(A, D)$  and Design  $B_{\text{comp}}(A, D)$ . The ASN at the null hypothesis  $\theta=1$  is given by

$$\begin{aligned} \text{ASN} &= 100 \times P(\text{Stop at Analysis Time 1 using Design } A_{\text{comp}}(A, D) | \theta=1) \\ &\quad + 200 \times P(\text{Stop at Analysis Time 2 using Design } A_{\text{comp}}(A, D) | \theta=1) \\ &\quad + 200 \times P(\text{Stop at Analysis Time 2 using Design } B_{\text{comp}}(A, D) | \theta=1) \\ &\quad + 300 \times P(\text{Stop at Analysis Time 3 using Design } B_{\text{comp}}(A, D) | \theta=1) \\ &= 100 \times (0.0032 + 0.7668) + 200 \times (0.0187 + 0.2010) + 200 \times (0.0018 + 0.0006) + 300 \\ &\quad \times (0.0014 + 0.0067) \\ &= 123.85. \end{aligned}$$



In order to find a comparable prespecified group sequential design, we can again use the constrained boundary approach in order to match the decision boundaries at the first analysis. There are then an infinite number of ways that the design parameters can be modified in future analyses in order to closely match the unconditional power curve or the ASN curve to that of the adaptive approach. For instance, the following code could be used if the boundary shape parameters (the  $P$  parameters of the unified group sequential design family described by Kittelson and Emerson [5]) were to be modified:

```
> bouGS <- seqBoundary(designB)
> bouGS[2, ] <- NA
> designGS <- update(designB, exact.constraint=bouGS, P=c(0, 0.08))
> designGS
```

#### PROBABILITY MODEL and HYPOTHESES:

Two arm study of censored time to event response variable

Theta is hazard ratio (Treatment : Comparison)

One-sided hypothesis test of a lesser alternative:

Null hypothesis :  $\Theta \geq 1$  (size = 0.025)

Alternative hypothesis :  $\Theta \leq 0.5617$  (power = 0.975)

#### STOPPING BOUNDARIES: Sample Mean scale

	a	d
Time 1 (N= 100)	0.5792	0.8645
Time 2 (N= 200)	0.7589	0.7665
Time 3 (N= 300)	0.8025	0.8025

This design has the same boundaries as Design B\*( $A, D$ ) at the first and third analyses, and the boundary shape at the second analysis is determined by the  $P$  parameter, which is chosen to be 0 for the efficacy boundary, and 0.08 for the futility boundary. These values of  $P$  were chosen as discussed in Section 2.3 to produce a group sequential design that has power as close as possible to, but always greater than or equal to, the power of the adaptive design for  $\theta \in (0.5, 1)$ . For this example, the group sequential design has marginally higher power and slightly lower ASN than the adaptive design given above for all alternatives corresponding to hazard ratios between 0.3 and 1.2.

## Acknowledgements

This work is supported by NIH T32NS048005.

## References

1. Tsiatis AA, Mehta CR. On the inefficiency of the adaptive design for monitoring clinical trials. *Biometrika* 2003; **90**:367–378.
2. Jennison C, Turnbull BW. Adaptive and nonadaptive group sequential tests. *Biometrika* 2006; **93**(1):1–21.
3. Emerson SS, Kittelson JM, Gillen DL. Frequentist evaluation of group sequential designs. *Statistics in Medicine* 2007; **26**(28):5047–5080.
4. Emerson SS, Kittelson JM, Gillen DL. Bayesian evaluation of group sequential designs. *Statistics in Medicine* 2007; **26**(7):1431–1449.
5. Kittelson JM, Emerson SS. A unifying family of group sequential test designs. *Biometrics* 1999; **55**:874–882.
6. Emerson SS. Issues in the use of adaptive clinical trial designs. *Statistics in Medicine* 2006; **25**(19):3270–3296.
7. Wang SK, Tsiatis AA. Approximately optimal one-parameter boundaries for group sequential trials. *Biometrics* 1987; **43**:193–199.
8. Pocock SJ. Group sequential methods in the design and analysis of clinical trials. *Biometrika* 1977; **64**:191–200.
9. O'Brien PC, Fleming TR. A multiple testing procedure for clinical trials. *Biometrics* 1979; **35**:549–556.
10. Emerson SS, Fleming TR. Symmetric group sequential test designs. *Biometrics* 1989; **45**:905–923.
11. Insightful Corporation, *S+SeqTrial*. Seattle, Washington, 2002.
12. Proschan MA, Hunsberger SA. Designed extension of studies based on conditional power. *Biometrics* 1995; **51**:1315–1324.
13. Cui L, Hung HMJ, Wang SJ. Modification of sample size in group sequential clinical trials. *Biometrics* 1999; **55**:853–857.
14. Burington BE, Emerson SS. Flexible implementations of group sequential stopping rules using constrained boundaries. *Biometrics* 2003; **59**:770–777.
15. Schoenfeld DA, Richter JR. Nomograms for calculating the number of patients needed for a clinical trial with survival as an endpoint. *Biometrics* 1982; **38**:163–170.
16. Emerson SS. *S+SeqTrial Technical Overview*. Insightful Corporation, 2003.