

A new approach to designing phase I-II cancer trials for cytotoxic chemotherapies

Jay Bartroff,^{a,*†} Tze Leung Lai^b and Balasubramanian Narasimhan^b

Recently, there has been much work on early phase cancer designs that incorporate both toxicity and efficacy data, called phase I-II designs because they combine elements of both phases. However, they do not explicitly address the phase II hypothesis test of $H_0 : p \leq p_0$, where p is the probability of efficacy at the estimated maximum tolerated dose $\hat{\eta}$ from phase I and p_0 is the baseline efficacy rate. Standard practice for phase II remains to treat p as a fixed, unknown parameter and to use Simon's two-stage design with all patients dosed at $\hat{\eta}$. We propose a phase I-II design that addresses the uncertainty in the estimate $p = p(\hat{\eta})$ in H_0 by using sequential generalized likelihood theory. Combining this with a phase I design that incorporates efficacy data, the phase I-II design provides a common framework that can be used all the way from the first dose of phase I through the final accept/reject decision about H_0 at the end of phase II, utilizing both toxicity and efficacy data throughout. Efficient group sequential testing is used in phase II that allows for early stopping to show treatment effect or futility. The proposed phase I-II design thus removes the artificial barrier between phase I and phase II and fulfills the objectives of searching for the maximum tolerated dose and testing if the treatment has an acceptable response rate to enter into a phase III trial. Copyright © 2014 John Wiley & Sons, Ltd.

Keywords: cancer trials; generalized likelihood ratio; group sequential; isotonic regression; maximum tolerated dose; phase I; phase II

1. Introduction

In typical phase I studies in the development of relatively benign drugs, the drug is initiated at low doses and subsequently escalated to show safety at a level where some positive response occurs, and healthy volunteers are used as study subjects. This paradigm does not work for diseases like cancer, for which a nonnegligible probability of severe toxic reaction has to be accepted to give the patient some chance of a favorable response to the treatment. Therefore, patients (rather than healthy volunteers) are used as study subjects, and it is widely accepted that some degree of toxicity must be tolerated to experience any substantial therapeutic effects. Hence, an acceptable proportion q of patients experiencing *dose limiting toxicities* (DLTs) is generally agreed on before the trial, which depends on the type and severity of the DLT; the dose resulting in this proportion is thus referred to as the *maximum tolerated dose* (MTD). In addition to the explicitly stated objective of determining the MTD, a phase I cancer trial also has the implicit goal of safe treatment of the patients in the trial. However, the aims of treating patients in the trial and generating an efficient design to estimate the MTD for future patients often run counter to each other. Commonly used designs in phase I cancer trials implicitly place their focus on the safety of the patients in the trial, beginning from a conservatively low starting dose and escalating cautiously.

In [1, 2], Bartroff and Lai have given a review of model-based methods to design phase I cancer trials and proposed a general framework that incorporates both 'individual' and 'collective' ethics into the design of the trial. We have also developed a new design, which minimizes a risk function composed of two terms, with one representing the individual risk of the current dose and the other representing the

^aDepartment of Mathematics, University of Southern California, 3620 South Vermont Avenue, KAP 108, Los Angeles, CA 90089, U.S.A.

^bDepartment of Statistics, Sequoia Hall, Stanford University, Stanford, CA 94305, U.S.A.

*Correspondence to: Jay Bartroff, Department of Mathematics, University of Southern California, 3620 South Vermont Avenue, KAP 108, Los Angeles, CA 90089, U.S.A.

†E-mail: bartroff@usc.edu

collective risk, and have shown that it performs better than existing model-based designs in accuracy of the MTD estimate at the end of the trial, and toxicity and overdose rates of patients in the trial, and loss functions reflecting the individual and collective ethics.

The MTD determined from a phase I study is used in a subsequent phase II study, in which

a cohort of patients is treated, and the outcomes are related to the prespecified target or bar. If the results meet or exceed the target, the treatment is declared worthy of further study; otherwise, further development is stopped. This has been referred to as the “go/ no go” decision ([3], p. 927).

The most widely used designs for these single-arm phase II trials are Simon’s two-stage design [4], which allows early stopping of the trial if the treatment has not shown beneficial effect, that is measured by a Bernoulli proportion. Simon considered the design that stops for futility (i.e., accepts the null hypothesis H_0 in (1)) after n_1 patients if the number of patients exhibiting positive treatment effect is $r_1 (\leq n_1)$ or fewer and otherwise treats an additional n_2 patients and rejects the treatment (again, accepts H_0) if and only if the number of patients exhibiting positive treatment effect is $r (\leq n_1 + n_2)$ or fewer. Simon’s design requires that a null proportion p_0 , representing some ‘uninteresting’ level of positive treatment effect, and an alternative $p_1 > p_0$ be specified. The null hypothesis is

$$H_0 : p \leq p_0, \tag{1}$$

where p denotes the probability of positive treatment effect. The type I and II error probabilities $\alpha = P_{p_0}(\text{Reject } H_0)$, $\beta = P_{p_1}(\text{Accept } H_0)$, and the expected sample size $E_{p_0}N$ can be computed for any design of this form, which can be represented by the parameter vector (n_1, n_2, r_1, r) . By using computer search over these integer-valued parameters, Simon [4] tabulated the optimal designs in his Tables I and II for different values of (p_0, p_1) . Simon’s design has been generalized by Jung *et al.* [5, 6] who also give a graphical method of selecting from among the admissible designs, Simon’s original procedure being one of them, and by Lu *et al.* [7] to allow for partial responses. Whether the new treatment is declared promising in a phase II trial depends strongly on the prescribed p_0 and p_1 . The sample size m of a typical phase I trial and the maximum sample size $M = n_1 + n_2$ of a typical phase II trial are relatively small, 20–30 for phase I and no more than 60 for phase II. Vickers *et al.* [3] conclude that uncertainty in the choice of p_0 and p_1 can increase the likelihood that (i) a treatment with no viable positive treatment effect proceeds to phase III, or (ii) a treatment with positive treatment effect is abandoned at phase II.

1.1. An integrated approach to dose finding and testing for efficacy

In Sections 2 and 3, we address these issues concerning the design of early phase single-arm cancer trials by developing a novel seamless phase I-II trial design that uses efficient statistical methods for the design and analysis of the integrated trial, subject to ethical and sample size constraints. The data from the trial are toxicity and efficacy outcomes at various doses and consist of $(x_i, y_i, z_i), i = 1, \dots, N$, where N is the phase I-II total sample size, x_i denotes the dose given to the i th subject, $y_i = 1$ or 0 according to whether a DLT occurs or not, and $z_i = 1$ or 0 according to whether the subject responds to the treatment. For cytotoxic treatments, both the dose-toxicity curve $P(y_i = 1 | x_i = x)$ and the dose-response curve $P(z_i = 1 | x_i = x)$ increase with the dose x , and therefore, the MTD is the most efficacious dose subject to a prespecified probability q of severe toxic reaction. Whereas the objective of a traditional phase I cancer trial is to estimate the MTD, denoted by η , from $(x_i, y_i), i = 1, \dots, m$, and that of the ensuing phase II trial with maximum sample size M is to test if the response rate exceeds some prespecified level p_0 when all patients in the trial are assigned dose $\hat{\eta}$, which is the MTD estimate from the phase I trial, our integrated design continues sequential estimation of η throughout the trial with total maximum sample size $m + M$ and uses an efficient group sequential test of the null hypothesis that the response rate at η does not exceed p_0 . In Section 2, we consider commonly used logistic regression models for dose-toxicity and dose-response relationships to pinpoint the basic ideas. Section 3 removes the parametric assumptions and extends the methodology to dose-toxicity and dose-response relationships that are only assumed to be monotone. Simulation studies in Section 4 demonstrate the advantages of the integrated design, and Section 5 describes the underlying theory and implementation details.

1.2. Review of current methods using toxicity and efficacy/response data

Gooley *et al.* [8] suggested using efficacy and toxicity data together and performed simulations to compare the operating characteristics of three *ad hoc* designs. Thall and Russell [9] proposed a design

combining binary toxicity data y_i and trinomial response data $z_i = 0, 1$, or 2 for no, moderate, or severe response, respectively, into a single trinomial variable

$$w_i = \begin{cases} 0, & \text{if } z_i = 0 \text{ and } y_i = 0 \\ 1, & \text{if } z_i = 1 \text{ and } y_i = 0 \\ 2, & \text{if } z_i = 2 \text{ or } y_i = 1. \end{cases} \quad (2)$$

Using a proportional odds regression model for w_i on dose x_i with a prior distribution on its unknown parameters, a Bayesian posterior calculation along the lines of O’Quigley, Pepe, and Fisher’s [10] continual reassessment method (CRM) is performed to calculate the acceptability of the available discrete dose levels and escalate or de-escalate the current dose level. For a similar setting, O’Quigley *et al.* [11] proposed a phase I design for HIV trials in which binary efficacy z_i and toxicity y_i variables are combined into a single trinomial variable (2) in which we now set $w_i = 2$ if $y_i = 1$. A CRM-like calculation is used to treat the current patient at the posterior estimate of the dose maximizing the probability of simultaneous efficacy and nontoxicity.

For efficacy and toxicity measurements, Ivanova [12] proposed an up-and-down design, which assigns doses in pairs on a discrete set of dose levels. Braun [13] proposed a bivariate version of CRM in which a bivariate joint distribution is chosen for (y_i, z_i) , and the target dose is defined to be the one minimizing the expected Euclidean distance to prespecified toxicity and efficacy rates, with respect to a chosen noninformative posterior distribution. In particular, the bivariate distribution of Arnold and Strauss [14], which gives Bernoulli conditional distributions of y_i given z_i , and vice versa, was recommended. Thall and Cook [15] proposed a different method for combining efficacy and toxicity responses. First, marginal efficacy and toxicity curves are assumed, which are then combined using a Gaussian or Gumbel copula; this approach differs from Braun’s method that specifies the conditional distributions rather than the marginals. Doses are then selected using ‘trade-off contours’ in the two-dimensional space of outcome probabilities on which the outcomes are equally desirable. Thall *et al.* [16] extended this method to allow for the inclusion of patient-specific covariates.

Even when the designs summarized earlier are called ‘phase I-II’ designs, it is because they incorporate efficacy (or tumor response) data. They do not address testing the efficacy hypothesis that is the purpose of typical phase II cancer studies, for which the standard practice is to use Simon’s two-stage design following the dose-finding portion. Moreover, this skirts the issue of uncertainty in the estimated MTD used in the null hypothesis in phase II, as well as ignores toxicity outcomes that are available during phase II, which could help improve this estimate, especially because the phase I sample size is usually small. The innovative phase I-II design proposed herein aims at rectifying these issues and hence provides a common framework that can be used all the way from phase I through the final accept/reject decision about the null hypothesis on efficacy in the phase II portion of the study, utilizing both toxicity and efficacy data for dose finding while performing efficient group sequential testing of the null hypothesis.

2. An integrated approach to designing early phase cancer clinical trials

A widely used model for the dose-toxicity curve is the logistic regression model

$$P(y_i = 1 | x_i = x) = F(x; \theta) := 1 / \left(1 + e^{-(\theta_1 + \theta_2 x)} \right), \quad (3)$$

where $\theta = (\theta_1, \theta_2)$, and it is assumed that $\theta_2 > 0$ (i.e., probability of toxicity increases with dose), for which the MTD is given by $\eta = [\log(q/(1 - q)) - \theta_1] / \theta_2$. Under (3), the estimate $\hat{\eta}$ based on $(x_i, y_i), i = 1, \dots, m$, can be obtained by maximum likelihood, which is equivalent to logistic regression. Similarly, we can model the dose-response curve by

$$P(z_i = 1 | x_i = x) = p(x; \psi) := 1 / \left(1 + e^{-(\psi_1 + \psi_2 x)} \right), \quad (4)$$

under which the probability p of the response in the null hypothesis $H_0 : p \leq p_0$ of the traditional phase II cancer trial is actually $p(\hat{\eta}; \psi)$. The difference between $\hat{\eta}$ and η is completely ignored in currently used designs, and the toxicity outcomes in the phase II trial are also ignored. Combining the toxicity outcomes in phase II with those in phase I can improve the estimate of η , especially because the phase I sample size is small. Changing the null hypothesis to

$$H_0 : p(\eta; \psi) \leq p_0 \tag{5}$$

not only takes into consideration the uncertainties in $\hat{\eta}$ as an estimate of η but also leads to continual updating of η with toxicity outcomes in the phase II trial if one uses a generalized likelihood ratio (GLR) test. Moreover, the GLR test also uses the phase I efficacy outcomes $z_i, i = 1, \dots, m$.

2.1. The first phase of the phase I-II trial

The first phase of the new phase I-II (or dose-finding) design involves only the dose-toxicity data but not the responses z_i . We can use traditional methods or recent advances in phase I cancer trial designs to perform dose escalation; see Section 5.2 and the references therein for details. At the end of the phase I trial, we compute the maximum likelihood or Bayes estimates $\tilde{\theta}, \tilde{\eta}$, and $\tilde{\psi}$ of θ, η , and ψ . Let \mathcal{F}_0 denote the phase I data $(x_1, y_1, z_1), \dots, (x_m, y_m, z_m)$.

2.2. The ensuing group sequential design to test efficacy and re-estimate η

After this initial group of m patients, the proposed design switches to a group sequential scheme, with specified group sizes m_1, \dots, m_K (e.g., $m_1 = \dots = m_K$ gives constant group size sampling). The group sequential scheme updates the MTD estimate via MLE at the k th interim analysis with an additional batch of size m_k of dose-toxicity data $(x_{\tau_{k-1}+1}, y_{\tau_{k-1}+1}), \dots, (x_{\tau_k}, y_{\tau_k})$, where

$$\tau_k = m + \sum_{i=1}^k m_i, \quad k = 1, \dots, K. \tag{6}$$

It also uses all the observed data $(x_i, y_i, z_i), 1 \leq i \leq \tau_k$, to perform a group sequential GLR test of $H_0 : p(\eta; \psi) \leq p_0$ at the k th interim analysis, where $p(x; \psi)$ is defined by (4). Lai and Shih [17] have developed a methodology of nearly optimal group sequential tests, which use versatile and asymptotically efficient GLR test statistics and stopping boundaries. In conjunction with GLR statistics, maximum likelihood (rather than Bayes) estimates of η are used for sequential updating of the estimated MTD.

To simplify the description, we begin by assuming that y_i and z_i are independent; this assumption will be removed in Section 2.3. Let $\ell_k(\psi)$ denote the log-likelihood function for ψ at the k th interim analysis, which because of the independence assumption only depends on the z_i and not the y_i :

$$\ell_k(\psi) = \log \left\{ \prod_{i=1}^{\tau_k} p(x_i; \psi)^{z_i} [1 - p(x_i; \psi)]^{1-z_i} \right\} = \sum_{i=1}^{\tau_k} \{z_i(\psi_1 + \psi_2 x_i) - \log(1 + e^{\psi_1 + \psi_2 x_i})\}.$$

Let $\hat{\psi}_k$ be an MLE maximizing this, $\hat{\theta}_k = (\hat{\theta}_{k,1}, \hat{\theta}_{k,2})$ be an MLE of θ based on the data up to and including the k th interim analysis, $\hat{\eta}_k = (\text{logit}(q) - \hat{\theta}_{k,1}) / \hat{\theta}_{k,2}$, and

$$S_k^j = \{\psi : p(\hat{\eta}_k; \psi) = p_j\} \quad \text{for } 0 \leq k \leq K, j = 0, 1, \tag{7}$$

where $\hat{\eta}_0 = \tilde{\eta}$, $p_1 > p_0$, and $H_1 : p(\eta; \psi) \geq p_1$ is the alternative hypothesis. The choice of p_1 will be discussed in Section 5.2.

As will be explained in Section 5.1, we can compute at the k th interim analysis the test statistics

$$\ell_{k,j} = \min_{\psi \in S_k^j} [\ell_k(\hat{\psi}_k) - \ell_k(\psi)], \quad j = 0, 1, \tag{8}$$

so that the group sequential test stops and rejects H_0 at interim analysis $k < K$ if

$$p(\hat{\eta}_k, \hat{\psi}_k) > p_0 \quad \text{and} \quad \ell_{k,0} \geq b, \tag{9}$$

and early stopping for futility (accepting H_0) at analysis $k < K$ can also occur if

$$p(\hat{\eta}_k, \hat{\psi}_k) < p_1 \quad \text{and} \quad \ell_{k,1} \geq \tilde{b}. \tag{10}$$

The test rejects H_0 at the K th analysis if

$$p(\widehat{\eta}_K, \widehat{\psi}_K) > p_0 \quad \text{and} \quad \ell_{K,0} \geq c. \tag{11}$$

The thresholds b, \widetilde{b} , and c are chosen so that

$$\max_{\psi \in \mathcal{S}_0^0} P_{\widetilde{\theta}, \psi}(H_0 \text{ rejected} | \mathcal{F}_0) = \alpha \tag{12}$$

and the power

$$\min_{\psi \in \mathcal{S}_0^1} P_{\widetilde{\theta}, \psi}(H_0 \text{ rejected} | \mathcal{F}_0) \tag{13}$$

is close to $1 - \beta$, as in [17, 18] and [19]. Details and software for implementation are given in Section 5.2.

2.3. Modeling the dependence between y_i and z_i

We can model the dependence between y_i and z_i by replacing the marginal model (4) by the following model for the conditional distribution of z_i given y_i :

$$P(z_i = 1 | y_i = 0, x_i) = 1 / (1 + e^{-(\psi_1^0 + \psi_2^0 x_i)}) \tag{14}$$

$$P(z_i = 1 | y_i = 1, x_i) = 1 / (1 + e^{-(\psi_1^1 + \psi_2^1 x_i)}) \tag{15}$$

with parameters $\psi^0 = (\psi_1^0, \psi_2^0)$ and $\psi^1 = (\psi_1^1, \psi_2^1)$. Generalizing (5) to include (14)–(15), the null hypothesis is that the probability of efficacy at dose $x = \eta$ is less than or equal to p_0 , that is,

$$H_0 : \frac{1 - q}{1 + e^{-(\psi_1^0 + \psi_2^0 \eta)}} + \frac{q}{1 + e^{-(\psi_1^1 + \psi_2^1 \eta)}} \leq p_0, \tag{16}$$

noting that $F(\eta; \theta) = q$. This null hypothesis is an extension of that in Section 2.2 and can again be tested by sequential GLR theory.

2.4. Modifications for discrete dose levels

In practice, the dose levels in dose-finding studies of cancer drugs are usually chosen before the trial from a finite set

$$\Lambda = \{\lambda_1, \dots, \lambda_d\}, \quad \text{where} \quad \lambda_1 < \lambda_2 < \dots < \lambda_d, \tag{17}$$

unlike the continuous doses we have assumed so far. In this case, the MTD has to be redefined as

$$\eta = \begin{cases} \max\{\lambda \in \Lambda : F(\lambda; \theta) \leq q\}, & \text{if } F(\lambda_i; \theta) \leq q \text{ for some } i \\ \lambda_1, & \text{otherwise.} \end{cases} \tag{18}$$

Putting this modified definition of η in (5) or (16), we can still apply the group sequential GLR test of Section 2.2 or 2.3, in which we also modify the definition of $\widehat{\eta}_k$ accordingly to be Λ -restricted. That is, $\widehat{\eta}_k$ is the smallest $\lambda_j \in \Lambda$ maximizing the likelihood ℓ_k up through the k th interim analysis, and we set $x_i = \widehat{\eta}_k$ for $i = \tau_k + 1, \dots, \tau_{k+1}$. Note that the group sequential GLR test is based on all the observed data (x_i, y_i, z_i) up to the time of an interim analysis, irrespective of how the x_i are chosen, and therefore, no additional modifications are needed.

Because Λ is discrete, one can use more robust specification of the dose-toxicity and/or dose-response curve than the logistic regression models (3) and (4). Details are given in the next section. For samples of the size typically used in early phase cancer trials, however, one usually does not have enough data to detect departures from these ‘working models.’ In addition, the initial phase of a dose-finding study for cytotoxic chemotherapies is often very conservative, to avoid causing harm to patients before observing how the new treatment actually works in human subjects. This explains the popularity of the widely used, although inefficient, 3+3 designs. A more efficient alternative is to use a two-stage phase I design in which a more cautious design is used for the first stage before switching to a parametric model-based design in the second stage; see [1, Section 4.2]. Once we have zoomed in on a range around the MTD

that is narrow relative to the original dose range, the logistic model is actually quite robust because it can be viewed as a locally linear regression model around the MTD, adjusted with the logit link for Bernoulli outcomes. What this means is that one only needs to be concerned with the choice of the design levels x_i to ensure such robustness in the locally logit-linear model. Thus, the GLR test statistic can be restricted only to those x_i that are within a certain distance from $\hat{\eta}_k$ at the k th interim analysis.

3. Extension to monotone dose-toxicity and dose-response relationships

In many dose-finding trials, the number of discrete dose levels (17) is relatively small. For this situation, in this section we develop an approach similar to the Bayesian models of Yin *et al.* [20] and Yin and Yuan [21] where the probabilities of toxicity and efficacy are order restricted but in a frequentist setting. Assume for now that y_i and z_i are independent; the general case will be covered in Section 3.2. Because the number of dose levels is small, we also assume that all the levels have been used at least once during phase I; if this does not hold, then only the used dose levels are carried forward into phase II. Instead of the parameterization by the toxicity and efficacy parameters θ and ψ , we parameterize by the toxicity and efficacy probabilities

$$\phi_i = P(y = 1|x = \lambda_i), \quad \pi_i = P(z = 1|x = \lambda_i), \quad i = 1, \dots, d. \quad (19)$$

The MTD (18) can then be written

$$\eta = \lambda_{i^*} \quad \text{where} \quad i^* = \begin{cases} \max\{i : \phi_i \leq q\}, & \text{if } \phi_i \leq q \text{ for some } i \\ 1, & \text{otherwise} \end{cases}$$

so that the phase II null and alternative hypotheses can be expressed as

$$H_0 : \pi_{i^*} \leq p_0 \quad \text{vs.} \quad H_1 : \pi_{i^*} \geq p_1.$$

3.1. Order-restricted maximum-likelihood estimation and generalized likelihood ratio statistics

Letting $x_t = \lambda_{i_t}$ denote the t -th dose, $t = 1, \dots, \tau_k$ with τ_k given by (6), and $\boldsymbol{\pi} = (\pi_1, \dots, \pi_d)$, the log-likelihood at the k th interim analysis of phase II under the independence assumption is

$$\ell_k(\boldsymbol{\pi}) = \log \left\{ \prod_{t=1}^{\tau_k} \pi_{i_t}^{z_t} (1 - \pi_{i_t})^{1-z_t} \right\}. \quad (20)$$

The order-restricted MLE $\hat{\boldsymbol{\pi}}_k = (\hat{\pi}_{1,k}, \dots, \hat{\pi}_{d,k})$ maximizing (20) subject to $\pi_1 \leq \dots \leq \pi_d$ is given by the formula

$$\hat{\pi}_{i,k} = \min_{j' \geq i} \max_{j \leq i} \left(\frac{S_{j,k} + \dots + S_{j',k}}{v_{j,k} + \dots + v_{j',k}} \right), \quad i = 1, \dots, d, \quad (21)$$

where $S_{i,k} = \sum_{t=1}^{\tau_k} z_t 1\{i_t = i\}$ is the sum of the efficacy responses at level i and $v_{i,k} = \sum_{t=1}^{\tau_k} 1\{i_t = i\}$ is the number of patients that have been dosed at level i up through the k th analysis ([22], p. 52). An analogous formula holds for the order-restricted MLE of the toxicity probabilities $\hat{\boldsymbol{\phi}}_k = (\hat{\phi}_{1,k}, \dots, \hat{\phi}_{d,k})$. These order-restricted MLEs can be computed by solving the minimization-maximization problem in (21) or, equivalently, by using the well-known pool adjacent violators algorithm Pool Adjacent Violators Algorithm (PAVA); see [22, Section 2.4].

The order-restricted MLE of the MTD at the k th interim analysis can be defined as

$$\hat{\eta}_k = \lambda_{\hat{i}_k^*}, \quad \text{where} \quad \hat{i}_k^* = \begin{cases} \max\{i : \hat{\phi}_{i,k} \leq q\}, & \text{if } \hat{\phi}_{i,k} \leq q \text{ for some } i \\ 1, & \text{otherwise.} \end{cases} \quad (22)$$

Let $\tilde{\boldsymbol{\pi}}_k^j = (\tilde{\pi}_{1,k}^j, \dots, \tilde{\pi}_{d,k}^j)$, $j = 0, 1$, be the constrained order-restricted MLE, which maximizes (20) subject to the order-restriction $\pi_1 \leq \dots \leq \pi_d$ and the additional constraint that

$$\pi_{\hat{i}_k^*}^j \leq p_0 \quad \text{for } j = 0 \quad \text{and} \quad \pi_{\hat{i}_k^*}^j \geq p_1 \quad \text{for } j = 1, \quad (23)$$

which can be computed as follows. If $\widehat{\pi}_{i_k^*,k} \leq p_0$, then $\widetilde{\pi}_k^0 = \widehat{\pi}_k$. Otherwise, $\widehat{\pi}_{i_k^*,k} > p_0$, so suppose that $\widehat{\pi}_{i_k^*-r-1,k} \leq p_0 < \widehat{\pi}_{i_k^*-r,k}$, in which case we set $\widetilde{\pi}_{i_k^*-r,k}^0 = \dots = \widetilde{\pi}_{i_k^*,k}^0 = p_0$, and $\widetilde{\pi}_{i,k}^0$ coincides with $\widehat{\pi}_{i,k}$ for all other i . In other words, when $\widehat{\pi}_k$ falls outside H_0 , $\widetilde{\pi}_k^0$ is computed by setting the appropriate elements of $\widehat{\pi}_k$ to the boundary value p_0 , and $\widetilde{\pi}_k^1$ is computed similarly.

The log-likelihood ratio statistics at the k th interim analysis for testing $H_0 : \pi_{i^*} \leq p_0$ versus $H_1 : \pi_{i^*} \geq p_1$ are given by

$$\ell_{k,j} = \ell_k(\widehat{\pi}_k) - \ell_k(\widetilde{\pi}_k^j), \quad j = 0, 1, \quad (24)$$

with $\ell_k(\pi)$ defined by (20), and the group sequential test stops and rejects H_0 at interim analysis $k < K$ if

$$\widehat{\pi}_{i_k^*,k} > p_0 \quad \text{and} \quad \ell_{k,0} \geq b, \quad (25)$$

stops for futility if

$$\widehat{\pi}_{i_k^*,k} < p_1 \quad \text{and} \quad \ell_{k,1} \geq \widetilde{b}, \quad (26)$$

and otherwise rejects H_0 at the K th analysis if

$$\widehat{\pi}_{i_K^*,K} > p_0 \quad \text{and} \quad \ell_{K,0} \geq c. \quad (27)$$

As in Section 2.2, the thresholds b, \widetilde{b} , and c are chosen so that (12) holds and the power is close to $1 - \beta$. Details are given in Section 5.2.

3.2. Modeling the dependence between y_i and z_i

A flexible method for modeling the general case where the toxicity and efficacy observations may not be independent is to introduce d additional parameters in the form of the global cross ratios

$$\rho_i = \frac{\Pi_i(0,0)\Pi_i(1,1)}{\Pi_i(1,0)\Pi_i(0,1)}, \quad i = 1, \dots, d, \quad \text{where} \quad \widetilde{\Pi}_i(y, z) = P(y_t = y, z_t = z | x_t = \lambda_i).$$

Dale [23] proposed using the global cross ratio as a useful measurement of dependence in discrete ordered bivariate responses, and they have been recently used by Yin *et al.* [20] in a Bayesian phase I-II design. If the toxicity and efficacy responses are independent, then $\rho_i = 1$ for all $i = 1, \dots, d$. The complete joint distribution $\Pi_i(y, z)$ of the toxicity and efficacy responses can be recovered from the parameters $\pi_i, \phi_i, \rho_i, i = 1, \dots, d$ through the following formulas:

$$\Pi_i(1,1) = \begin{cases} (a_i - \sqrt{a_i^2 + b_i})/[2(\rho_i - 1)], & \text{if } \rho_i \neq 1 \\ \pi_i \phi_i, & \text{if } \rho_i = 1 \end{cases}$$

$$\Pi_i(1,0) = \pi_i - \Pi_i(1,1), \quad \Pi_i(0,1) = \phi_i - \Pi_i(1,1), \quad \Pi_i(0,0) = 1 - \pi_i - \phi_i + \Pi_i(1,1),$$

where $a_i = 1 + (\pi_i + \phi_i)(\rho_i - 1)$ and $b_i = -4\rho_i(\rho_i - 1)\pi_i\phi_i$. The log-likelihood at the k th interim analysis of phase II for this general case is

$$\ell_k(\pi, \phi, \rho) = \log \left\{ \prod_{t=1}^{\tau_k} \Pi_{i_t}(y_t, z_t) \right\} \quad \text{where} \quad x_t = \lambda_{i_t}, \quad (28)$$

and the log-likelihood ratio statistics at the k th interim analysis for testing $H_0 : \pi_{i^*} \leq p_0$ versus $H_1 : \pi_{i^*} \geq p_1$ are given by

$$\ell_{k,j} = \ell_k(\widehat{\pi}_k, \widehat{\phi}_k, \widehat{\rho}_k) - \ell_k(\widetilde{\pi}_k^j, \widetilde{\phi}_k^j, \widetilde{\rho}_k^j), \quad j = 0, 1, \quad (29)$$

with stopping rules as in (25)–(27), where $\widehat{\pi}_k, \widehat{\phi}_k, \widehat{\rho}_k$ are MLEs maximizing (28) subject to the order restrictions $\pi_1 \leq \dots \leq \pi_d$ and $\phi_1 \leq \dots \leq \phi_d$, and $\widetilde{\pi}_k^j, \widetilde{\phi}_k^j, \widetilde{\rho}_k^j$ maximize (28) subject to these order restrictions plus the constraints (23).

4. Simulation studies

4.1. Operating characteristics of the traditional and proposed phase I-II designs on a continuous dose space

To investigate the effect of uncertainty in the estimate $\hat{\eta}$ on the operating characteristics of the phase II hypothesis test that is used in current practice, we first simulated a phase I design, which we take to be escalation with overdose control (EWOC) introduced by Babb *et al.* [24], followed by Simon's optimal two-stage design. EWOC is a popular dose-finding method originally proposed for continuous dose spaces, which we consider here. Motivated by a real trial for 5-fluorouracil to treat solid colon tumors described in Babb *et al.* [24], we let $[x_{\min}, x_{\max}] = [140, 425]$ denote the known range of acceptable dose values and assume $m = 24$ patients are treated in phase I. We parametrize the toxicity responses' distribution $F(x; \cdot)$ by η and $\rho = F(x_{\min}; \theta)$ rather than $\theta = (\theta_1, \theta_2)$ and assume that (ρ, η) has the uniform distribution on $[0, q] \times [x_{\min}, x_{\max}]$ as its prior distribution; see [1, Section 2] for more details. Fixing $\eta = 250$, $q = 1/3$, and $\rho = .1$, Table I gives some operating characteristics of this phase I-II design using Simon's design for testing (1) with $p_0 = .1$ and $p_1 = .25$, with $\beta = .2$ and various values of α . These were evaluated from 100,000 simulations using the aforementioned values of η , x_{\min} , and x_{\max} , and under the efficacy parameter $\psi = (-3.895, .00679)$ chosen so that $p(\eta; \psi) = p_0 = .1$ and $p(x_{\max}; \psi) = .9$.

For several values of the parameters (n_1, n_2, r_1, r) of Simon's two-stage design [4, Table 2] of the phase II trial, Table I compares the prescribed type I error probability α of Simon's test with the actual probability of rejecting $H_0 : p(\eta; \psi) \leq p_0$, denoted by $P(-H_0|\cdot)$, for three choices of the MTD estimate $\hat{\eta}$ that is used as the dose for the phase II trial. The three types of estimation are the MLE, the final posterior mean of the phase I trial (which is what the original version of the Bayesian CRM [10] would use), and the dose recommended by EWOC that is used in the phase I design of this simulation study. Table I shows that the actual probability $P(-H_0|\cdot)$ of falsely rejecting H_0 is largely inflated over the prescribed value α of the type I error probability used for the phase II trial, especially when the posterior mean is used for the MTD estimate $\hat{\eta}$. The reason for this is the frequent overestimation of η by $\hat{\eta}$, as shown by the five-number summary (maximum, first quartile Q_1 , median, third quartile Q_3 , and maximum) of the

Method for $\hat{\eta}$	MLE	CRM	EWOC	
$\min(\hat{\eta})$	140.0	141.2	141.0	
$Q_1(\hat{\eta})$	226.3	246.9	229.1	
$\text{med}(\hat{\eta})$	244.7	264.7	246.9	
$Q_3(\hat{\eta})$	264.1	318.1	246.9	
$\max(\hat{\eta})$	425.0	391.6	362.7	
$E(\hat{\eta})$	252.6	276.7	239.8	
$\text{RMSE}(\hat{\eta})$	52.2	44.2	29.0	
α	$n_1/n_2/r_1/r$	$P(-H_0 \text{MLE})$	$P(-H_0 \text{CRM})$	$P(-H_0 \text{EWOC})$
.05	18/25/2/7	.180 (.001)	.479 (.002)	.100 (.0009)
.04	18/30/2/8	.176 (.001)	.476 (.002)	.094 (.0009)
.03	18/35/2/9	.170 (.001)	.470 (.002)	.088 (.0009)
.02	22/44/3/11	.167 (.001)	.464 (.002)	.083 (.0009)
.01	22/58/3/14	.156 (.001)	.458 (.002)	.074 (.0008)

The true MTD is $\eta = 250$; $\hat{\eta}$ denotes the final maximum tolerated dose estimate by either maximum-likelihood estimation (MLE), posterior mean (continual reassessment method [CRM]), or escalation with overdose control; and α denotes the prescribed type I error probability of Simon's phase II test of $p(\hat{\eta}; \psi) \leq p_0$ with design parameters n_1, n_2, r_1 , and r . The actual probability of rejecting $H_0 : p(\eta; \psi) \leq p_0$ is denoted by $P(-H_0|\hat{\eta})$, with standard errors in parentheses, where $\hat{\eta} = \text{MLE, CRM, or escalation with overdose control (EWOC)}$ is the final maximum tolerated dose estimate in phase I.

100,000 simulated values of $\hat{\eta}$ given in the table. Although underestimation of η by $\hat{\eta}$ also occurs, it is more often overestimated, which causes rejection of H_0 at rates higher than prescribed by the design parameters of Simon's test. Also given in the table are the mean $E(\hat{\eta})$ and the root mean square error (RMSE)($\hat{\eta}) = \{E(\hat{\eta} - \eta)^2\}^{1/2}$ of the estimated MTD. We comment that here we have only considered the most basic versions of CRM [10] and EWOC [24], and many variants have been proposed since then (e.g., [25, 26]). It seems likely that the properties of $\hat{\eta}$ could be improved using one of these variants of CRM or EWOC, but because our focus here is more on the interaction between phase I and phase II, we do not explore that option here.

Focusing on the traditional two-stage design with $\alpha = .05$ in Table I (denoted here by Trad) and concentrating on MLE estimation for simplicity, we compare its operating characteristics with those of the new phase I-II design described in Section 2 (denoted by New). In order to match the Trad design's probability $P(\text{rej. } H_0) = .18$ of falsely rejecting $H_0 : p(\eta; \psi) \leq p_0$, where $p_0 = .1$, at the parameter values determined by $p(\eta; \psi) = .1$, we choose critical values $b = 3$, $\tilde{b} = 3.5$, and $c = .7$ in (9)–(11). Although a type I error probability of .18 is usually deemed too high, we can keep the probability of falsely rejecting H_0 close to .05 if we use $p_0 = .05$ instead, as shown in Table II, which compares the operating characteristics of the Trad and New designs based on 10,000 simulations. The two designs both have phase I sample size of $m = 24$ and maximum phase II sample size of 43, and the New design achieves this through phase II group sizes 10, 10, 10, 10, and 3. As in Table I, η is fixed at 250 and $\rho = .1$, whereas ψ is specified by fixing $p(x_{\max}, \psi) = .9$ and varying $p(\eta; \psi)$ over the values .05, .1, .2, .3, .4, and .5. For each scenario, Table II gives $P(\text{rej. } H_0)$, in which $p_0 = .1$, and the total expected sample size EN over the two phases. It shows that the new design has smaller $P(\text{rej. } H_0)$ than Trad for $p(\eta; \psi) = .05$ and larger $P(\text{rej. } H_0)$ for all values $p(\eta; \psi) > .1$, and uniformly smaller expected sample size, substantially so for parameter values $p(\eta; \psi) > .3$. In addition, Table II also gives the probability $p(\hat{\eta}_{\text{rec}}; \psi)$ of efficacious response at the recommended dose $\hat{\eta}_{\text{rec}}$, which for Trad is the MTD estimate at the end of phase I and for New is the final MLE at the end of phase II, the overall response rate (denoted Eff) for subjects in the study, the overall overdose rate (denoted OD) of subjects treated at doses above the true MTD, and the RMSE($\hat{\eta}_{\text{rec}}$) of the recommended dose. The RMSE of the recommended dose for New is substantially smaller than Trad throughout, which we attribute to its continued estimation of η during phase II. The values $p(\hat{\eta}_{\text{rec}}; \psi)$ and Eff are comparable with $p(\eta; \psi)$ throughout for New, whereas the corresponding values for Trad are larger, and Trad has larger OD values than New.

Table II shows a dramatic improvement of the New design relative to the Trad design in terms of both power and average sample size. In order to discern how much of this improvement is due to the group sequential sampling used (relative to Simon's two-stage design) versus how much is due to the continued estimation of the MTD during phase II that the proposed design allows, more simulation studies were performed whose results are in Tables III and IV. In addition, both of these simulation studies were performed under different parameter values than in Table II in order to see the proposed design's performance over a broad range of scenarios.

In Table III, the traditional phase I-II design (denoted Trad) was implemented but, instead of using Simon's two-stage design for phase II, the same group sequential sampling scheme that New used in Table II with group sizes 10, 10, 10, 10, and 3 was used. The proposed design (denoted New) was also implemented using these groups sizes and compared with Trad so that the only difference between the two designs is that Trad does not update the estimate $\hat{\eta}$ of the MTD during phase II. To see the performance of the proposed design in a different scenario than Table II, using the same dose range $[x_{\min}, x_{\max}] = [140, 425]$ and prior structure as there, the true MTD η was taken to be 350 and the probability of toxicity ρ at dose x_{\min} was taken to be .2. This scenario represents a much 'flatter' dose-toxicity curve than in Table II. In this set up, the phase II null hypothesis $H_0 : p(\eta; \psi) \leq p_0$ was tested with $p_0 = .5$, and Table III contains the operating characteristics of these designs at six different values of the response parameter ψ determined by $p(x_{\max}, \psi) = .95$ and $p(\eta, \psi) = .4, .5, .6, .7, .8$ and .9. Unlike the Trad design in Table II, which does not achieve the overall type I error probability $P(\text{rej. } H_0)$ at $p(\eta; \psi) = p_0$ equal to the prescribed value $\alpha = .05$ because of the variance of the MTD estimate used in phase II, here, the Trad design uses the stopping rule (9)–(11) with the critical values b, \tilde{b} , and c chosen so that this quantity is equal to α for $p_0 = .5$; they are $b = 24.1, \tilde{b} = 64.4$, and $c = 21.8$. The New design uses the values $b = 18.8, \tilde{b} = 54.4$, and $c = 11.7$, also chosen so that its type I error probability is α , and are slightly different than Trad's critical values because New continues to update $\hat{\eta}$ during phase II. Table III contains the operating characteristics of these designs based on 10,000 Monte Carlo replications at each parameter value. As might be expected from designs using the same sampling scheme, Trad and New have very similar expected sample size, and sample sizes are in general

Table II. Operating characteristics of the traditional (denoted Trad) and new (denoted New) designs described on Section 4.1. The toxicity parameter is fixed at $(\eta, \rho) = (250, .1)$, and the six values of the efficacy parameter ψ are determined by $p(x_{\max}; \psi) = .9$ and $p(\eta; \psi) = .05, .1, .2, .3, .4, .5$.

$p(\eta; \psi)$	5%		10%		20%		30%		40%		50%	
	Trad	New										
$p(\hat{\eta}_{rec}; \psi)$.101	.054	.150	.102	.233	.202	.319	.296	.409	.392	.499	.486
Eff	.096	.061	.140	.104	.219	.200	.310	.293	.405	.381	.498	.474
OD	.303	.291	.314	.312	.326	.289	.327	.256	.336	.252	.331	.249
RMSE($\hat{\eta}_{rec}$)	51.0	28.4	52.2	29.0	52.4	29.3	52.3	28.6	51.7	29.0	52.1	29.8
$P(\text{rej. } H_0)$.090	.051	.180	.180	.479	.645	.776	.923	.939	.989	.987	.999
E/N	45.9	40.2	49.8	47.3	57.7	51.0	63.2	43.7	66.0	37.0	66.7	34.6

All designs have phase I sample size $m = 24$ and maximum phase II sample size 43, for a maximum phase I-II sample size of 67. Eff is the overall response rate for subjects in the study, OD is the overall overdose rate of subjects treated at doses above the true maximum tolerated dose, and RMSE($\hat{\eta}_{rec}$) is the root mean square error of the recommended dose.

Table III. Operating characteristics of the traditional (denoted Trad) and new (denoted New) designs. The toxicity parameter is fixed at $(\eta, \rho) = (350, .2)$, and the six values of the efficacy parameter ψ are determined by $p(x_{\max}; \psi) = .95$ and $p(\eta; \psi) = .4, .5, .6, .7, .8, .9$ described on Section 4.1.

$p(\eta; \psi)$	40%		50%		60%		70%		80%		90%	
	Trad	New										
$p(\hat{\eta}_{rec}; \psi)$.183	.179	.207	.219	.250	.270	.324	.361	.474	.511	.798	.812
Eff	.129	.123	.153	.151	.195	.195	.270	.280	.432	.439	.785	.788
OD	.111	.108	.109	.108	.111	.107	.109	.111	.109	.106	.102	.093
RMSE($\hat{\eta}_{rec}$)	73.3	65.4	72.6	65.8	73.2	65.8	72.8	65.9	73.3	65.8	72.5	65.8
$P(\text{rej. } H_0)$.040	.043	.050	.050	.057	.062	.071	.088	.085	.134	.234	.590
E/N	61.2	60.7	62.8	62.2	64.6	64.0	66.2	65.7	66.8	66.6	65.9	63.7

All designs have phase I sample size $m = 24$ and maximum phase II sample size 43, for a maximum phase I-II sample size of 67. Eff is the overall response rate for subjects in the study, OD is the overall overdose rate of subjects treated at doses above the true maximum tolerated dose, and RMSE($\hat{\eta}_{rec}$) is the root mean square error of the recommended dose.

Table IV. Operating characteristics of the traditional (denoted Trad) and new (denoted New) designs. The toxicity parameter is fixed at $(\eta, \rho) = (200, .15)$, and the six values of the efficacy parameter ψ are determined by $p(x_{\max}; \psi) = .99$ and $p(\eta; \psi) = .025, .05, .25, .45, .65, .85$.

$p(\eta; \psi)$	2.5%		5%		25%		45%		65%		85%	
	Trad	New										
$p(\hat{\eta}_{rec}; \psi)$.011	.032	.017	.061	.082	.263	.197	.454	.383	.648	.706	.849
Eff	.008	.038	.011	.067	.061	.268	.178	.458	.377	.656	.706	.850
OD	.612	.611	.613	.615	.367	.564	.595	.569	.575	.578	.570	.570
RMSE($\hat{\eta}_{rec}$)	31.3	25.8	31.7	27.6	30.1	34.7	30.8	35.4	30.9	33.7	30.7	35.3
$P(\text{rej. } H_0)$.007	.002	.010	.010	.209	.185	.636	.478	.940	.763	.999	.819
EN	33.9	34.1	34.8	36.8	42.7	52.8	49.4	54.0	53.2	53.9	53.9	53.9

All designs have phase I sample size $m = 24$ and maximum phase II sample size 30, for a maximum phase I-II sample size of 54. Eff is the overall response rate for subjects in the study, OD is the overall overdose rate of subjects treated at doses above the true maximum tolerated dose, and RMSE($\hat{\eta}_{rec}$) is the root mean square error of the recommended dose.

larger in this scenario than the one in Table II, which is also to be expected because of the flatness of the dose-toxicity curve, which makes η difficult to estimate accurately, reflected in the power of both designs being low until $p(\eta, \psi)$ reaches 90%, where the power of New is 59%, but Trad is still severely underpowered. Note also that even though the flatness of the dose-toxicity makes the MTD difficult to estimate accurately, the chance of overdose is relatively low. Overall, New is slightly but consistently more efficient with smaller RMSE despite having slightly smaller average sample size, and New has higher power and response probabilities $p(\hat{\eta}_{rec}; \psi)$ over the range of parameter values in the alternative. These results are consistent with the two designs using the same phase II sampling scheme but New using continued estimation of the MTD throughout phase II.

Table IV considers another scenario, with a smaller phase II sample size, in which both Trad and New use for phase II a two-stage design with early stopping only for futility. For Trad, this is Simon's two-stage design, and for New, this is the stopping rule (9)–(11) with $K = 2$ groups and b fixed at ∞ so that only early stopping for futility can occur. In this scenario, both Trad and New have maximum phase II sample size 30 (compared with 43 in Tables II and III) and phase I sample size $m = 24$. To achieve this, Trad uses Simon's [4, Table 1] design with $r_1 = 0, n_1 = 9, r = 3,$ and $n_2 = 21$ for $\alpha = .05$ and $\beta = .1$ at $p_0 = .05$ and $p_1 = .25$. As in Tables I and II, the Trad design using these parameters does not achieve the type I error probability at the prescribed value $\alpha = .05$ because of variance of the MTD estimate used in phase II. Indeed, Table IV shows its actual type I error probability to be .01 at $p(\eta; \psi) = p_0 = .05$. Unlike Table II that shows inflation of type I error probability, here, the type I error probability is substantially smaller than the prescribed value $\alpha = .05$. In order to make a meaningful comparison between designs, we choose the parameters of the New design to match this smaller value of the type I error probability, for which we use $b = \infty$ (to allow early stopping only for futility), $\tilde{b} = 2.1$ and $c = 19.3$ in (9)–(11) and phase II group sizes 9 and 21, the same as the Simon design. The operating characteristics of these designs are given in Table IV, based on 10,000 Monte Carlo replications each, in yet another scenario with $\eta = 200, \rho = .15,$ and six values of ψ determined by $p(x_{max}; \psi) = .99$ and $p(\eta; \psi) = .025, .05, .25, .45, .65,$ and $.85$. The dose range and prior structure are the same as in Table III. The response probabilities $p(\hat{\eta}_{rec}; \psi)$ at New's recommended dose stay much closer to the true values than at Trad's recommended dose, likely because of New's update of the MTD estimate during phase II. The overall response rate of subjects in the study is also substantially higher in New than in Trad. The two designs have similar average sample sizes, reflective of their similar sampling schemes, and Trad has higher power in the alternative. The RMSEs of the two designs are small and relatively close, with Trad's being slightly smaller. Note, however, that the squared error RMSE($\hat{\eta}_{rec}$) ignores the sign of $\eta - \hat{\eta}_{rec}$ and that the results on $p(\hat{\eta}_{rec}; \psi)$ show that $\hat{\eta}_{rec}$ tends to underestimate η .

4.2. Performance of the traditional and new phase I-II designs on discrete dose space under monotonicity constraints

To evaluate the performance of the phase II method proposed in Section 3 for monotonic efficacy and toxicity models on a discrete dose space, we performed a similar study to the one in Section 4.1, assuming independence of the toxicity and efficacy responses for simplicity; we have performed additional simulations under dependent responses using the model described in Section 3.2, and the performance of the new method is similar. Again focusing on the Trad design in Table I and using isotonic MLE estimation (22) for both the Trad and new (denoted by New) designs, the estimated operating characteristics are compared in Table V based on 10,000 simulated trials, wherein the phase I doses of the $m = 24$ patients are uniformly sampled from the dose set $\Lambda = \{140, 200, 250, 300, 350, 425\}$. In this setting, the Trad design with nominal level $\alpha = .05$ for testing $\pi_{i^*} \leq p_0$ actually has type I error probability $P(\text{rej. } H_0) = .211$ of falsely rejecting $H_0 : \pi_{i^*} \leq p_0 = .1$, and so in order to compare New and Trad in this setting, we choose critical values $b = .13, \tilde{b} = 3.3,$ and $c = .03$ in (9)–(11) in order to approximately match this, giving $P(\text{rej. } H_0) = .201$ for New at $\pi_{i^*} = .1$. In order to have the same maximum phase II sample size $M = 43$ as Trad, again New uses group sequential sampling with group sizes 10, 10, 10, 10, and 3. In this discrete nonparametric setting, the unknown parameters are the true toxicity and efficacy probabilities ϕ and π given by (19), and in order to compare Trad and New in a setting similar to the one in Section 4.1, we consider values of ϕ and π given by the corresponding parametric models $F(x; \theta)$ and $p(x; \psi)$ and parameter values given there: $\eta = \lambda_{i^*}$ is fixed at 250, $\rho = \phi_1 = .1,$ $p(x_{max}, \psi) = \pi_d = .9,$ and $p(\eta; \psi) = \pi_{i^*} = .05, .1, .2, .3, .4,$ and $.5$. The relative performance of Trad and New is very similar to that in the previous section: The new design has smaller $P(\text{rej. } H_0)$ than Trad for parameter values $\pi_{i^*} \leq .1$ in the null hypothesis, larger $P(\text{rej. } H_0)$ for all values $\pi_{i^*} > .1$ in

Table V. Operating characteristics of the traditional (denoted Trad) and new (denoted New) designs described in Section 4.2. The true toxicity and efficacy probabilities are determined by the same parameters as in Table II: $\lambda_{i^*} = 250$, $\phi_1 = .1$, $\pi_d = .9$, and the six cases $\pi_{i^*} = .05, .1, .2, .3, .4$, and $.5$ as described in Section 4.2.

π_{i^*}	5%		10%		20%		30%		40%		50%	
	Trad	New										
π_{i^*}	.072	.030	.116	.061	.194	.131	.274	.206	.357	.295	.449	.395
$\pi_{i^*}^*$.185	.196	.225	.231	.286	.296	.350	.364	.416	.441	.492	.524
Eff	.390	.376	.388	.363	.366	.361	.347	.370	.328	.390	.320	.406
OD	56.5	60.1	57.4	60.9	56.7	59.1	56.9	59.2	56.5	58.3	57.2	57.9
RMSE(λ_{i^*})	.117	.076	.211	.201	.410	.486	.615	.729	.805	.895	.931	.981
$P(\text{rej. } H_0)$	56.5	38.7	49.4	40.7	54.8	41.7	59.8	40.2	63.6	37.7	65.9	35.6

Eff is the overall response rate for subjects in the study, OD is the overall overdose rate of subjects treated at doses above the true maximum tolerated dose, and RMSE($\hat{\eta}_{rec}$) is the root mean square error of the recommended dose.

the alternative, and uniformly smaller expected sample size, substantially so when π_{i^*} is large or small relative to $p_0 = .1$. The other operating characteristics given in the table are the same as in Table II: The response rate $\pi_{\hat{\gamma}_i^*}$ at the final recommended dose, overall response rate (Eff) and overdose rate (OD) of patients in the study, and the RMSE of the final recommended dose $\lambda_{\hat{\gamma}_i^*}$. The Eff rate of New is larger than Trad at all parameter values considered, which we attribute to the proposed design's ability to vary the dose throughout phase II, and hence 'correct' for a poorly chosen MTD estimate at the end of phase I, to some measure. The OD rates of the two designs are close, with New being sometimes smaller and sometimes larger. The RMSE of New is slightly larger, but comparable with Trad, which we attribute to its markedly smaller average sample size.

5. Group sequential likelihood theory and implementation details

5.1. Theory of group sequential generalized likelihood ratio tests

We first assume independence between y_i and z_i given x_i as in Section 2.2. In this case, the likelihood function, based on a sample of size τ_k , is of the form $L_{1,k}(\theta)L_{2,k}(\psi)$, where

$$L_{1,k}(\theta) = \prod_{i=1}^{\tau_k} [F(x_i; \theta)]^{y_i} [1 - F(x_i; \theta)]^{1-y_i}, \quad L_{2,k}(\psi) = \prod_{i=1}^{\tau_k} [p(x_i; \psi)]^{z_i} [1 - p(x_i; \psi)]^{1-z_i}.$$

The GLR statistic for testing $p(\eta; \psi) = p_j$, which is the boundary of H_j , is

$$\log \left[\frac{\sup_{\theta} L_{1,k}(\theta) \times \sup_{\psi} L_{2,k}(\psi)}{\sup_{(\theta, \psi): p(\eta; \psi) = p_j} L_{1,k}(\theta)L_{2,k}(\psi)} \right], \quad (30)$$

and the signed root likelihood ratio statistic is approximately normal under $p(\eta; \psi) = p_j$; see [17, p. 513]. Note that $p(\eta; \psi) = p_j$ can be expressed as an equality constraint $\psi_1 + \eta\psi_2 = \text{logit}(p_j)$ on the linear function $\psi_1 + \eta\psi_2$ of ψ , and we can reparameterize ψ as $(\psi_1, \psi_1 + \eta\psi_2)$ and θ as (η, ρ) . Therefore, standard asymptotic analysis of GLR statistics shows that under $p(\eta; \psi) = p_j$, (30) has the same limiting distribution as

$$\log \left[\frac{\sup_{\psi} L_{2,k}(\psi)}{\sup_{\psi: p(\hat{\eta}_k; \psi) = p_j} L_{2,k}(\psi)} \right], \quad (31)$$

jointly over $1 \leq k \leq K$; see [27, Section 9.3(iii)]. Because the x_i are sequentially determined random variables (based on group sequential estimates of the MTD), we use the martingale central limit theorem ([28], p. 411) here instead of the traditional central limit theorem as in [17]. Note that (31) is the same as $\ell_{k,j}$ defined in (8). For the dependent case in Section 2.3, the likelihood function $L_{2,k}(\psi)$ involves both z_i and y_i in view of (14) and (15) but does not depend on η . A similar argument can be used to show that the GLR statistic at the k th interim analysis is still asymptotically equivalent to (31).

The group sequential GLR test of H_0 is much more flexible and efficient than Simon's two-stage likelihood ratio test [4] for phase II cancer trials. As noted in the last paragraph of Section 1.1, Simon's procedure actually tests $p(\hat{\eta}; \psi) \leq p_0$ with all doses set at the MTD estimate $\hat{\eta}$ from the phase I toxicity data, whereas the proposed test considers the more natural $H_0 : p(\eta; \psi) \leq p_0$ and uses all the observed (x_i, y_i, z_i) up to the time of interim analysis to test H_0 . Moreover, unlike Simon's two-stage design, which is actually a group sequential test with two groups and only allows futility stopping in the first stage, we use a more flexible group sequential design that allows early stopping for both efficacy and futility. In addition, the estimate of η of the phase I-II trial uses data up to the end of the trial. The group sequential GLR test uses the alternative p_1 implied by the maximum size τ_K (Section 5.2) to derive the futility stopping criterion, namely stopping when there is enough evidence against $H_1 : p(\eta; \psi) \geq p_1$. Similarly, it stops early for efficacy if the GLR statistics show enough evidence against $H_0 : p(\eta; \psi) \leq p_0$.

The group sequential GLR test in Section 3 that considers discrete dose levels also involves a finite number of parameters satisfying certain monotonicity constraints. Therefore, the theory of group sequential tests that we have applied to the logistic regression models in Section 2 can also be applied to

Section 3 that imposes certain structure on the parameter space. Lai and Shih [17, Section 3] have established the asymptotic efficiency of these group sequential GLR tests in terms of the expected sample size and power function. Here, we extend this theory in two ways. The first extension is from the independent and identically distributed model to the regression model, with sequentially determined regressors x_i . The second extension is to replace the GLR statistics by more easily computable and interpretable approximations that have the same asymptotic distributions. As noted earlier, martingale theory used in conjunction with likelihood theory provides the key tools for such extensions.

5.2. Implementation details

The MLEs of θ and ψ involved in the design proposed in Section 2 should be computed under the assumption of positive slope, that is, $\theta_2 > 0$ and $\psi_2 > 0$. In practice, this can be imposed by choosing a small value $\delta > 0$ and computing the MLEs under the constraint $\theta_2 \geq \delta$ and $\psi_2 \geq \delta$. A related issue is that the MLEs of θ and ψ may not exist in the first few stages of phase I (see p. 195 of [29]). In this case, their Bayes estimates from a Bayesian model-based design can be used instead.

The alternative $p_1 > p_0$ is implied by the maximum sample size τ_K and the desired type I and II error probabilities α and β , respectively. That is, for the GLR test that has fixed sample size τ_K and rejects H_0 if and only if

$$p(\hat{\eta}_K; \hat{\psi}_K) > p_0 \quad \text{and} \quad \min_{\psi \in \mathcal{S}_K^0} [\ell_K(\hat{\psi}_K) - \ell_K(\psi)] \geq C_\alpha, \quad (32)$$

let $p_1 > p_0$ be the alternative satisfying

$$\min_{\psi \in \mathcal{S}_0^1} P_{\theta, \psi}^{\sim} [(32) \text{ occurs} | \mathcal{F}_0] = 1 - \beta. \quad (33)$$

In (32), C_α is such that

$$\max_{\psi \in \mathcal{S}_0^0} P_{\theta, \psi}^{\sim} [(32) \text{ occurs} | \mathcal{F}_0] = \alpha \quad (34)$$

and the doses x_1, \dots, x_{τ_K} are chosen by some design. The computation of the left-hand sides of (33) and (34) will be described in the succeeding text.

The thresholds b, \tilde{b} , and c in (9)–(11) can be determined as follows. Let $0 < \varepsilon < 1/2$ and first choose \tilde{b} so that

$$\max_{\psi \in \mathcal{S}_0^1} P_{\theta, \psi}^{\sim} [(10) \text{ occurs for some } 1 \leq k < K | \mathcal{F}_0] = \varepsilon\beta. \quad (35)$$

Then choose b so that

$$\max_{\psi \in \mathcal{S}_0^0} P_{\theta, \psi}^{\sim} [(9) \text{ occurs for some } 1 \leq k < K, p(\hat{\eta}_{\tau_{k'}}, \hat{\psi}_{\tau_{k'}}) \geq p_1 \text{ and } \ell_{k',1} < \tilde{b} \text{ for all } k' < k | \mathcal{F}_0] = \varepsilon\alpha, \quad (36)$$

and finally choose c so that

$$\max_{\psi \in \mathcal{S}_0^0} P_{\theta, \psi}^{\sim} [(11) \text{ occurs and (9), (10) do not occur for any } 1 \leq k < K | \mathcal{F}_0] = (1 - \varepsilon)\alpha. \quad (37)$$

The determination of b, \tilde{b} , and c in (35)–(37) follows that in [17] and aims at controlling the type I error probability (12) and keeping the power (13) close to $1 - \beta$.

As in Section 3.4 of [17], we can use the joint asymptotic normality of the signed root likelihood ratio statistics to approximate the probabilities in (33)–(37). Because the GLR statistics are asymptotic pivots, the convergence in distribution holds uniformly over \mathcal{S}_0^1 or \mathcal{S}_0^0 , and therefore, the minimum (or maximum) over \mathcal{S}_0^1 or \mathcal{S}_0^0 in the left-hand sides of (33)–(37) poses no additional difficulty when we use the normal approximation. An alternative to normal approximation is to use Monte Carlo similar to that used in the bootstrap tests. Bootstrap theory suggests that we can simulate from the estimated distribution under an assumed composite hypothesis because the GLR statistic is an approximate pivot under that hypothesis. Thus, the bootstrap test chooses the $\psi \in \mathcal{S}_0^j$ in (33)–(37) to be the maximum-likelihood estimate (MLE) based on the phase I data \mathcal{F}_0 , of ψ under the constraint $p(\tilde{\eta}; \psi) = p_j$. In the simulation studies in Section 4, we use 10,000 bootstrap simulations to estimate the probabilities in (33)–(37). The implementation of the group sequential order-restricted GLR

test of $H_0 : \pi_{i^*} \leq p_0$ in Section 3 is similar, as we have explicit formulas (20) and (21). A software package to design the proposed phase I-II trial has been developed using R and is available at the website <http://med.stanford.edu/biostatistics/ClinicalTrialMethodology.html>.

6. Discussion

The simulation studies in Section 4, which are motivated by the trial in Babb *et al.* [24], show that the estimate $\hat{\eta}$ at the end of the phase I trial can substantially overestimate or underestimate η and therefore have a significantly higher or lower response rate than $p(\eta; \psi)$. Another situation in which the latter can occur is when using the 3+3 dose escalation scheme in phase I, which tends to produce a subtherapeutic dose $\hat{\eta}$ at the end of phase I. Continuing dose finding in phase II can add substantial information for estimating η , as Section 4 has shown.

Recognizing that the dose chosen at the end of the phase I trial may not ensure safety, Bryant and Day [30] have extended Simon's two-stage design for the phase II trial to incorporate toxicity outcomes in the phase II trial by stopping the trial after the first stage if either the observed response rate is inadequate or the number of observed toxicities is excessive and by recommending the treatment at the end of the phase II trial only if there are both a sufficient number of responses and an acceptably small number of toxicities. Note that the Bryant–Day design still uses $\hat{\eta}$ determined from the phase I data to be the dose throughout the phase II trial. We have developed herein a novel methodology, which continues dose finding to estimate the MTD in phase II and which uses the toxicity outcomes throughout the trial in a natural way, while focusing on testing the efficacy hypothesis during the phase II component of the phase I-II design. The methodology enables the user to carry out the novel group sequential extensions, allowing early stopping not only for futility but also for efficacy of Simon's two-stage design that is widely used in phase II cancer trials. These group sequential tests use efficient GLR statistics, which we have extended herein from the traditional logistic regression models in Section 2 to robust isotonic regression models in Section 3.

Bayesian designs have been proposed for phase II trials, allowing early stopping for efficacy or futility, and rejecting (or accepting) the hypothesis $p \leq p_0$ if the posterior probability of $p > p_0$ exceeds some threshold (or falls below another threshold), thereby extending the Bayesian approach from phase I to phase II trials; see Chapter 4 of [31]. Yin *et al.* [20] and Yuan and Yin [21] have developed Bayesian phase I-II designs to incorporate the bivariate outcomes of toxicity and efficacy to determine the dose sequentially for the next cohort of patients in the trial. Their underlying philosophy is that 'with a very limited sample size in the (traditional) phase I trial, the MTD might not be obtained in a reliable way', and therefore, they aim instead at finding 'the optimal dosage of a drug which has the highest effectiveness as well as tolerable toxicity' [20, p. 777]. Two motivating trials that attempt to 'speed up the drug discovery and reduce the total cost' are given in [32, p. 925 and Section 3] and [20].

The trials that motivate the phase I-II design proposed herein are traditional phase I and phase II trials at cancer centers of most medical schools, such as the Norris Comprehensive Cancer Center at the University of Southern California and the Cancer Institute at Stanford University. The protocols usually have small sample sizes for phase I, followed by Simon's two-stage design for phase II that uses the MTD estimated from the phase I data. Simon's design has been popular because it allows interim analysis for a go/no go decision while preserving the type I error probability and power at the effect size used to justify the sample size specified in the protocol. The reason why investigators with whom we have worked adhere to this design although they recognize difficulties with the relatively small sample sizes for both phases is that they can publish the trial results in medical journals that prefer frequentist testing. The phase I-II design proposed herein is an attempt to enable the investigators to perform valid group sequential tests of efficacy while continuing estimation of the MTD during the entire course of the phase I-II trial. Even though pharmaceutical companies do not need to publish the results of phase II trials and can focus on dose finding that incorporates both toxicity and efficacy as in the Bayesian designs of [20] and [21], many industry-sponsored phase II trials are still conducted at academic centers where this innovative phase I-II design can allow investigators to carry out group sequential frequentist testing of efficacy at the MTD and update the MTD estimate during the entire course of the trial. While the present paper has established the basic methodology, much of the work for its adoption still lies ahead. This includes generating some experience in actual trials and their protocols, holding monthly forums and regular consulting sessions for clinical investigators at the University of Southern California Norris Cancer Center and the Stanford Cancer Institute, and developing user-friendly software based on this experience, which will facilitate its use by other academic centers.

Acknowledgements

Bartroff's work was supported by NSF grants DMS-0907241 and DMS-1310127 and NIH grant GMS-068968. Lai's work was supported by NSF grant DMS-1106535 and NIH grant 5P30CA124435. Narasimhan's work was supported by NCI Cancer Center Support grant 5P30CA124435.

References

1. Bartroff J, Lai TL. Approximate dynamic programming and its applications to the design of phase I cancer trials. *Statistical Science* 2010; **25**:245–257.
2. Bartroff J, Lai TL. Incorporating individual and collective ethics into phase I cancer trial designs. *Biometrics* 2011; **67**:596–603.
3. Vickers AJ, Ballen V, Scher HI. Setting the bar in phase III trials: the use of historical data for determining 'go/ no go' decision for definitive phase II trials. *Clinical Cancer Research* 2007; **13**:972–976.
4. Simon R. Optimal two-stage designs for phase II clinical trials. *Controlled Clinical Trials* 1989; **10**:1–10.
5. Jung S, Carey M, Kim K. Graphical search for two-stage designs for phase II clinical trials. *Controlled Clinical Trials* 2001; **22**:367–372.
6. Jung S, Lee T, Kim K, George S. Admissible two-stage designs for phase II cancer clinical trials. *Statistics in Medicine* 2004; **23**(4):561–569.
7. Lu Y, Jin H, Lamborn KR. A design of phase II cancer trials using total and complete response endpoints. *Statistics in Medicine* 2005; **24**(20):3155–3170.
8. Gooley TA, Martin PJ, Fisher LD, Pettinger M. Simulation as a design tool for phase I/II clinical trials: an example from bone marrow transplantation. *Controlled Clinical Trials* 1994; **15**(6):450–462.
9. Thall PF, Russell KE. A strategy for dose-finding and safety monitoring based on efficacy and adverse outcomes in phase I/II clinical trials. *Biometrics* 1998; **54**:251–264.
10. O'Quigley J, Pepe M, Fisher L. Continual reassessment method: a practical design for phase I clinical trials in cancer. *Biometrics* 1990; **46**:33–48.
11. O'Quigley J, Hughes MD, Fenton T. Dose-finding designs for HIV studies. *Biometrics* 2001; **57**(4):1018–1029.
12. Ivanova A. A new dose-finding design for bivariate outcomes. *Biometrics* 2003; **59**(4):1001–1007.
13. Braun T. The bivariate continual reassessment method: extending the CRM to phase I trials of two competing outcomes. *Controlled Clinical Trials* 2002; **23**:240–256.
14. Arnold BC, Strauss DJ. Bivariate distributions with conditionals in prescribed exponential families (Corr: V53 p700). *Journal of the Royal Statistical Society, Series B: Methodological* 1991; **53**:365–375.
15. Thall PF, Cook JD. Dose-finding based on efficacy-toxicity trade-offs. *Biometrics* 2004; **60**(3):684–693.
16. Thall PF, Nguyen HQ, Estey EH. Patient-specific dose finding based on bivariate outcomes and covariates. *Biometrics* 2008; **64**(4):1126–1136.
17. Lai TL, Shih MC. Power, sample size and adaptation considerations in the design of group sequential clinical trials. *Biometrika* 2004; **91**:507–528.
18. Bartroff J, Lai TL. Generalized likelihood ratio statistics and uncertainty adjustments in adaptive design of clinical trials. *Sequential Analysis* 2008; **27**:254–276.
19. Bartroff J, Lai TL. Efficient adaptive designs with mid-course sample size adjustment in clinical trials. *Statistics in Medicine* 2008; **27**:1593–1611.
20. Yin G, Li Y, Ji Y. Bayesian dose-finding in phase I/II clinical trials using toxicity and efficacy odds ratios. *Biometrics* 2006; **62**(3):777–787.
21. Yin G, Yuan Y. Bayesian model averaging continual reassessment method in phase I clinical trials. *Journal of the American Statistical Association* 2009; **104**(487):954–968.
22. Silvapulle MJ, Sen PK. *Constrained Statistical Inference: Inequality, Order, and Shape Restrictions*. Wiley-Interscience: Hoboken, New Jersey, 2004.
23. Dale J. Global cross-ratio models for bivariate, discrete, ordered responses. *Biometrics* 1986; **42**:909–917.
24. Babb J, Rogatko A, Zacks S. Cancer phase I clinical trials: efficient dose escalation with overdose control. *Statistics in Medicine* 1998; **17**:1103–1120.
25. Goodman SN, Zahurak ML, Piantadosi S. Some practical improvements in the continual reassessment method for phase I studies. *Statistics in Medicine* 1995; **14**:1149–1161.
26. Tighiouart M, Rogatko A. Dose finding with escalation with overdose control (EWOC) in cancer clinical trials. *Statistical Science* 2010; **25**(2):217–226.
27. Cox DR, Hinkley DV. *Theoretical Statistics*. Chapman and Hall: London, 1974.
28. Durrett R. *Probability: Theory and Examples*, 3rd edn. Thomson: Belmont, 2005.
29. Agresti A. *Categorical Data Analysis*. John Wiley & Sons: Hoboken, New Jersey, 2002.
30. Bryant J, Day R. Incorporating toxicity considerations into the design of two-stage phase II clinical trials. *Biometrics* 1995; **51**:1372–1383.
31. Berry SM, Carlin BP, Lee JJ, Muller P. *Bayesian Adaptive Methods for Clinical Trials*. CRC Press: Boca Raton, FL, 2010.
32. Yuan Y, Yin G. Bayesian phase I/II adaptively randomized oncology trials with combined drugs. *Annals of Applied Statistics* 2011; **5**(2A):924–942.