

Manuscript version: Author's Accepted Manuscript

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

Persistent WRAP URL:

<http://wrap.warwick.ac.uk/163406>

How to cite:

Please refer to published version for the most recent bibliographic citation information. If a published version is known of, the repository item page linked to above, will contain details on accessing it.

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Publisher's statement:

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk.

Efficiency of a Randomized Confirmatory Basket Trial Design Constrained to Control the
Family Wise Error Rate by Indication

Linchen He^{1,2}, Yuru Ren¹, Han Chen¹, Daphne Guinn^{3,4}, Deepak Parashar^{5,6}, Cong Chen⁷,
Shuai Sammy Yuan^{7,8}, Valeriy Korostyshevskiy¹, Robert A. Beckman^{1,9}

¹Department of Biostatistics, Bioinformatics and Biomathematics, Lombardi
Comprehensive Cancer Center, Georgetown University Medical Center, Washington, DC
USA

²Division of Biostatistics, Department of Population Health, New York University School of
Medicine, New York, NY, USA

³Program for Regulatory Science and Medicine, Georgetown University, Washington,
District of Columbia, USA

⁴Department of Pharmacology and Physiology, Georgetown University, Washington,
District of Columbia, USA

⁵Statistics and Epidemiology Unit & Cancer Research Centre, Warwick Medical School,
University of Warwick, Coventry UK

⁶The Alan Turing Institute for Data Science and Artificial Intelligence, The British Library,
London UK

⁷Biostatistics and Research Decision Sciences, Merck & Co., Inc., Kenilworth, NJ USA

⁸Kite Pharma, a Gilead Company, 2400 Broadway. Santa Monica, CA 90404

⁹Department of Oncology, Lombardi Comprehensive Cancer Center and Innovation
Center for Biomedical Informatics, Georgetown University Medical Center, Washington,
DC USA

Corresponding author:

Linchen He, Department of Biostatistics, Bioinformatics and Biomathematics, Lombardi
Comprehensive Cancer Center, Georgetown University Medical Center, Washington, DC
USA

Email: lh1790@nyu.edu

Abstract

Basket trials pool histologic indications sharing molecular pathophysiology, improving development efficiency. Currently basket trials have been confirmatory only for exceptional therapies. Our previous randomized basket design may be generally suitable in the resource-intensive confirmatory phase, maintains high power even with modest effect sizes, and provides nearly k -fold increased efficiency for k indications, but controls false positives for the pooled result only. Since family-wise error rate by indications (FWER) may sometimes be required, we now simulate a variant of this basket design controlling FWER at $0.025k$, the total FWER of k separate randomized trials. We simulated this modified design under numerous scenarios varying design parameters. Only designs controlling FWER and minimizing estimation bias were allowable. Optimal performance results when $k = 3,4$. We report efficiency (expected # true positives/expected sample size) relative to k parallel studies, at 90% power (“uncorrected”) or at the power achieved in the basket trial (“corrected”, because conventional designs could also increase efficiency by sacrificing power). Efficiency and power (percentage active indications identified) improve with higher percentage of initial indications active. Up to 92% uncorrected and 38% corrected efficiency improvement is possible.

Even under FWER control, randomized confirmatory basket trials substantially improve development efficiency. Initial indication selection is critical.

Keywords: *confirmatory basket trial, adaptive design, family-wise error rate, power by indication, cost-effectiveness*

INTRODUCTION

Molecular oncology has led to increasingly numerous biomarker-defined niche indications.¹ For example, lung cancer now includes several small subgroups with distinct therapies. This may lead to clinical trial enrollment challenges, and to additional development expense and delay. Conversely, a targeted therapy may have potential application in numerous indications, as well as in multiple combination settings, creating a large number of potential clinical hypotheses for testing with finite resources.

We previously discussed the need for increased development efficiency in design of proof of concept studies and associated Go-No Go decisions, given the large number of potential hypotheses worthy of testing, which may strain available resources.² Efficiency is a measure of the utility of a study per unit of resource (number of trial participants and/or financial cost) expended in the study or as a *consequence of the study*. For example, in the instance where a Phase 2 proof of concept study gives a false positive result, Phase 3 resources will likely be expended as a consequence of the proof of concept study in a futile effort to confirm the false positive result, and this outcome can be used to penalize the

1
2
3
4
5
6
7
8
9 efficiency value of the proof of concept study design, probability weighted by its
10 Type I error rate).² The efficiency of a proof of concept study design has therefore
11 been expressed as the probability of finding a true positive proof of concept (utility)
12 divided by two cost terms: the number of trial participants utilized in the proof of
13 concept study itself and the probability weighted number of Phase 3 trial
14 participants for Phase 3 trials resulting from true and false positive results. To
15 calculate this quantity, one needs an estimated probability distribution that a
16 therapy entered into the proof of concept trial from a population of therapies will be
17 active in the state of nature.²
18
19
20
21
22
23
24
25
26
27
28
29

30 Efficiency can be defined differently in different contexts in a fit for purpose
31 fashion. It is a fundamental metric for judging the value of clinical trial designs. In
32 contrast to power, it incorporates cost (measured financially or in trial participants).
33 It is always possible to increase the study power, by simply increasing its sample
34 size, but the cost of doing so continually rises with increasing power, and the point
35 where one reaches diminishing returns in seeking power is subjective, although
36 governed by traditions. By incorporating a probabilistic estimate that the therapy is
37 truly active, as well as setting a defined Type I error threshold and a defined
38 target, efficiency metrics blend such concepts as positive and negative predictive
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 value into a summary of the return (in useful drug development knowledge) on
10 investment.
11
12

13
14 Importantly efficiency can be calculated for a single study or for an ensemble of
15 studies, or even across a clinical development portfolio. In the work on proof of
16 concept studies, the realistic case was considered in which budgetary constraints
17 were applied to a portfolio of proof of concept opportunities of equal merit. The
18 budget (in available dollars or trial participants) was insufficient to perform
19 traditional proof of concept studies testing all of the proof of concept hypotheses.
20 In this setting, a surprising result was found: efficiency was higher if the proof of
21 concept studies were reduced in size, cost, and power compared to the traditional
22 80% power, allowing a larger number of hypotheses to be tested within the fixed
23 budget. Otherwise, there was too high an opportunity cost of not testing credible
24 hypotheses under the budgetary constraint. There was the possibility of not
25 testing hypotheses that might have resulted in true positives, an event which we
26 termed a Type III error.¹
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

45 This paper is primarily concerned with the efficiency advantages of basket trial
46 designs in the resource intensive confirmatory phase of development. In all
47 phases of development, less efficient approaches may contribute to the very high
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 cost of therapy, to long development times delaying availability of therapy, as well
10 as to decisions not to develop drugs for niche indications. Given that the
11 confirmatory phase is the most resource-intensive, efficiency improvements in this
12 phase may have the greatest practical impact.
13
14
15
16
17
18

19 Master protocols^{3,4} can potentially increase the efficiency of drug development,
20 and facilitate development of niche indications. Master protocols include platform
21 trials, umbrella trials, and basket trials. In platform trials, different therapies are
22 perpetually cycled through an ongoing trial, resulting in notable operational
23 efficiencies due to the existence of a common infrastructure, the “platform”. In
24 umbrella trials, multiple therapies are matched to multiple biomarker subgroups
25 within a single traditional organ-system-based indication. Enhanced efficiency
26 comes primarily from the ability to share a common standard of care control over
27 multiple experimental arms. Adaptive randomization, in which randomization
28 probabilities are adjusted frequently based on interim results, may also improve
29 efficiency. These efficiencies are not primarily operational, but rather are inherent
30 in the design.
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45

46
47 In a basket trial, traditional indications are grouped together in a basket based on
48 a shared molecular or pathophysiologic characteristic thought to predict utility of a
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 therapy. These indications may borrow information from each other, or may be
10 frankly pooled, leading to large improvements in development efficiency. This
11 development efficiency is again inherent in the design, arising directly from the fact
12 that multiple indications may contribute to the total sample size if they may be fully
13 pooled. In principle, if k indications are fully pooled, they may all be tested in the
14 pool for the same sample size N that would have been used for the testing of only
15 one indication in a traditional design. Under ideal conditions, a k -fold increase in
16 development efficiency may thus be achieved, a far greater efficiency increment
17 than is available from platform trials or umbrella trials. The challenge in optimizing
18 basket trials to approach this substantial benefit comes from the risk of
19 heterogeneity between indications despite shared biomarkers sometimes
20 conferring similar benefits. Active indications (defined as indications in which the
21 therapy is active) may carry inactive indications along with them to create a
22 positive pooled result. Conversely, inactive indications may dilute the effectiveness
23 of active indications, leading to a negative pooled result. These effects of
24 heterogeneity must be accounted for, and thus the efficiency gain of basket trial
25 designs will be less than k -fold, compared to the ideal case.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 Several efficient basket trial approaches use response rate data in the exploratory
10 setting. Most of these are based on Bayesian hierarchical models⁵⁻⁷ with the
11 exception of one that considers the likelihood that the indications all come from
12 one statistical distribution, versus the likelihood they are best modeled
13 individually.⁸
14
15
16
17
18
19

20
21 The first oncology basket trial in a regulatory approval setting was for imatinib,
22 which had already demonstrated extraordinary value in chronic myelogenous
23 leukemia, and had been rationally designed based on considerable scientific
24 evidence.⁹ The study was non-randomized and based on response rate, with very
25 small sample sizes. Forty indications were evaluated in less than 200 patients.
26
27 One approval resulted from 1 response in 5 patients. Note this design did not
28 utilize pooling. In this case, operational efficiencies were realized by consolidating
29 what would have been 40 tiny studies. However, in some settings basket designs
30 have operational challenges. For example, a basket design may require
31 operational cooperation between GI oncology, lung oncology, and breast
32 oncology, three divisions which may lack experience working together or uniform
33 standard operating procedures for tissue collection⁴. Related designs resulted in
34 approvals for transformational drugs designed for alterations in b-raf and ntrk
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 oncogenes.^{10,11} Recently, the immune checkpoint inhibitor pembrolizumab was
10 approved in multiple solid tumors, based on a basket trial pooling patients from
11 these indications with a DNA repair defect resulting in microsatellite instability,
12 utilizing response rate as a primary endpoint.¹² In all of these cases, the drugs and
13 biomarkers were supported by unusually strong scientific evidence, had previously
14 achieved transformational results, and were being studied in underserved
15 indications, and thus were able to merit approval based on single-arm response
16 rate data in relatively small populations.
17
18
19
20
21
22
23
24
25
26

27 28 **Motivation for the Research** 29

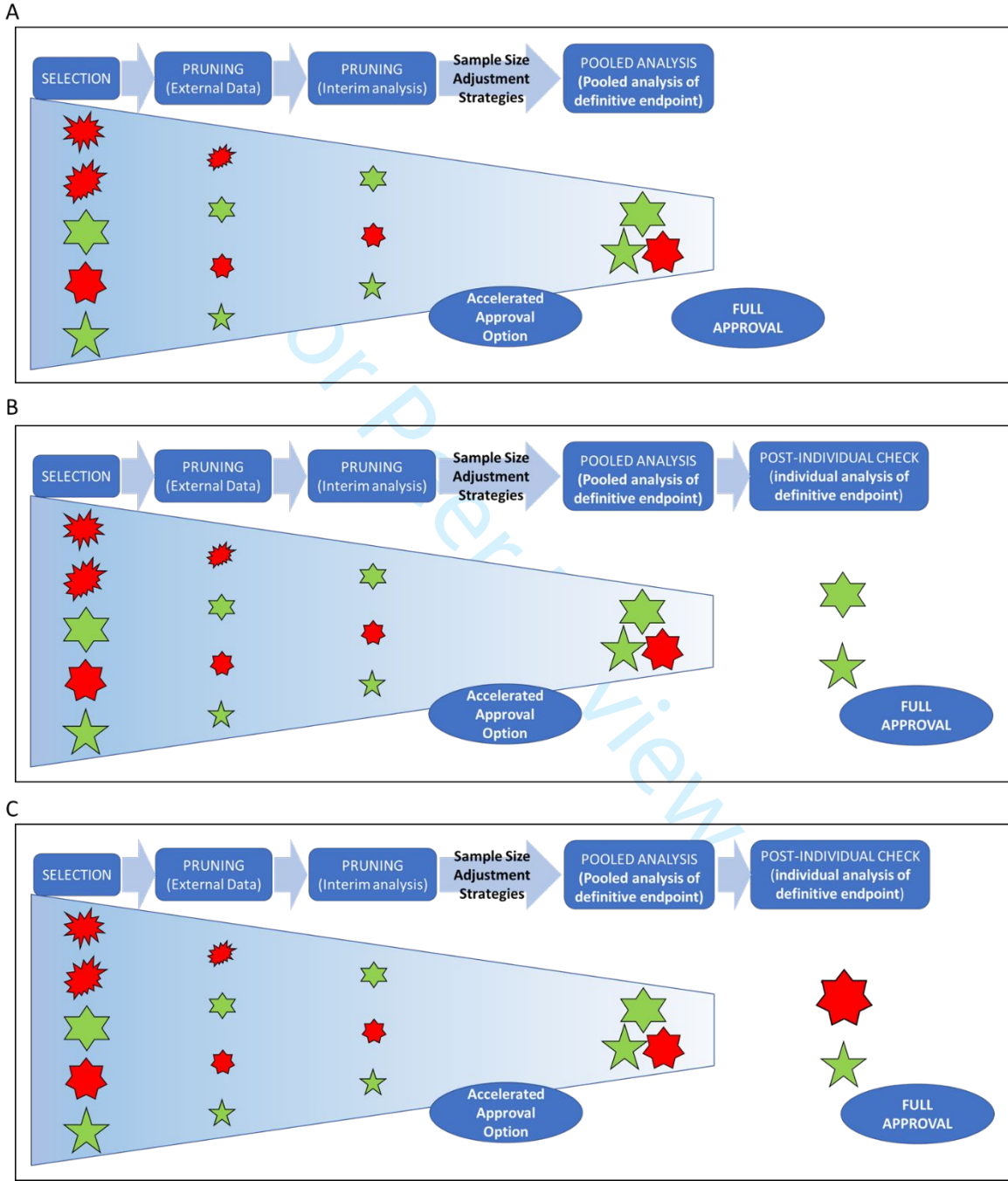
30
31 In contrast, many effective oncology drugs have required rigorous randomized
32 designs in the confirmatory setting. We developed a confirmatory basket trial
33 design^{13, 14} which, in addition to being applicable in a single-arm fashion or with
34 response rate endpoints, can also be utilized in a randomized controlled setting
35 using time to event (TTE) endpoints such as progression free survival (PFS) and
36 overall survival (OS). Randomization is generally important in approval of agents
37 whose clinical benefit must be measured based on TTE endpoints, unless the
38 effect size is transformative. We have been interested in developing a basket trial
39 design potentially suitable for the majority of confirmatory settings. We are
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 particularly interested in the confirmatory phase of development. Because the
10 confirmatory phase is the most resource and time-intensive phase, the marked
11 improvements in efficiency that are potentially available from basket trials can
12 have the greatest impact in efficiently bringing beneficial therapies to patients.
13
14
15
16
17
18

19 We have previously published a randomized confirmatory basket trial design.^{13, 14}
20 Although efficiency was not formally evaluated, an example application was given
21 in which a confirmatory study evaluated 6 indications for a sample size which
22 would be typical for one or at most two typical Phase 3 studies, i.e. a 3-fold or 6-
23 fold increase in efficiency. This original adaptive design resembles a funnel (Figure
24 1A, see Methods for more details). Indications are carefully selected, and then
25 filtered (removed) in several “pruning steps”, first with data external to the study
26 (i.e. maturing Phase 2 data from the same drug, data from other agents in the
27 class), and then with data from an interim analysis, which may be based on a
28 surrogate endpoint considered predictive of the definitive endpoint (i.e. PFS
29 predicting OS) or on early analysis of the definitive endpoint. The interim analysis
30 is performed on each indication individually to facilitate the pruning of inactive
31 indications. After pruning, a sample size adjustment may be applied to the
32 remaining indications, which are then pooled. The study concludes with a pooled
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 analysis of the remaining indications based on the definitive endpoint, and is
10 positive if statistically significant benefit is shown. Descriptive statistics including
11 hazard ratio, confidence intervals, and safety data are presented by indication for
12 informal benefit/risk analysis by health authorities. Individual indications may be
13 removed from the pool if their results are inconsistent with the pooled result, at the
14 discretion of health authorities. Ideally, the extent of and defining criteria for this
15 individual indication analysis should be negotiated with health authorities in
16 advance.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



1
2
3
4
5
6
7
8
9 Figure 1. Confirmatory basket trial design (A) with pruning and pooling as in previous studies^{13,14},
10 and (B-C) with pruning, pooling, and post-individual check simulated in this study. The designs (1)
11 conduct a basket trial that consists of k tumor indications; (2) prune (remove) indications based on
12 external data; (3) conduct an interim analysis independently for each tumor indication. Indications
13 that meet these interim criteria may in some cases be eligible for accelerated approval, indications
14 that do not are pruned; (4) adjust sample size of remaining indications as needed; (5) conduct a
15 pooled analysis of the remaining indications; (6) in the current design (B-C) only, conduct a
16 prospectively defined post-individual check for each indication involved in the pooled analysis. In
17 previous studies (A), indications that passed the pooled analysis may potentially be eligible for full
18 approval, whereas in the new design utilized in this study (B-C), indications must also pass a
19 simultaneous post-individual check to be potentially eligible for full approval. (B) shows the final
20 checking step successfully removing an inactive indication from the pool as intended. (C) shows a
21 case where the final checking step fails to remove an inactive indication, and instead mistakenly
22 removes an active indication. Active indications are in green, inactive indications in red.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

This original design demonstrated a dramatic increase in development efficiency, which we define in this work as (expected number of true positives/expected sample size subject to control of the false positive rate), and maintained acceptable power over a variety of scenarios even with inactive indications in the basket.^{13, 14} However, it was designed to control the false positive rate only in the

1
2
3
4
5
6
7
8
9 pool as a whole (Figure 2A, column 3), i.e. a pool which contains one or more truly
10 active indications is considered a true positive. The design does not control the
11
12
13 **family wise error rate (FWER)** by traditional indication subgroups (Figure 2A,
14 column 4; Methods, Supplemental Methods). In this paper, we define active and
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
inactive indications as indications in which the test drug does or does not provide
clinical benefit, in the unknown state of nature. FWER by indication subgroup,
simply called FWER in this paper, is defined as the probability that one or more
inactive indications will be approved by the design. **To further illustrate these
principles, we show analogous definitions for negative and positive predictive
value (NPV and PPV) in Figure 3 (see also Table 1 for definitions). When the
indications are carefully selected, as recommended for these designs, the majority
should be true positives. Under these conditions, we found that due to the low
prevalence of true negatives, NPV was low, and PPV was higher than one minus
the FWER (not shown).**

As we elaborate in the discussion, formal control of the **family-wise error rate** by
subgroups is not normally required in a confirmatory study, although informal ad
hoc subgroup analyses may be performed in the approval setting. Nonetheless, for
basket trials, in which the subgroups correspond to traditional indications,

1
2
3
4
5
6
7
8
9 published opinions of authors with past or present affiliation with the European
10 Medicines Agency and other European health authorities¹⁵, as well as our own
11 informal interactions with health authorities, indicate that control of the FWER by
12 indication may at times be recommended for basket trials in confirmatory settings.
13
14
15
16
17

18 While this is a subject of ongoing discussion, at the present time there is a
19 practical need for a variation of our original design which provides control of the
20 indication-specific FWER, in case it is required for approval in a particular setting.
21
22
23
24

25 In the current environment, discussion with relevant health authorities is
26 recommended before deciding whether to utilize the original randomized
27 confirmatory basket design^{13, 14} or the new variant to be presented in this paper.
28
29
30
31

32 The aims of this present paper are to introduce a variant of the original
33 randomized confirmatory basket design,^{13, 14} demonstrate that it controls
34 indication-specific FWER, and characterize its performance in terms of efficiency
35 gains relative to traditional designs (“relative efficiency”), power, FWER, estimation
36 bias, and confidence interval estimation. We observe some power losses with the
37 new design. We note that the efficiency of a traditional design may be improved
38 simply by reducing the power, due to the increasing cost of incremental gains in
39 power as power **increases**. Hence, if our new design results in some power loss,
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 the appropriate efficiency comparator is a traditional design with the same reduced
10 power (“corrected relative efficiency”).
11
12
13

14 Cunanan et al documented that a well-known exploratory basket trial design did
15 not control FWER, and, without opining whether this was acceptable in the
16 exploratory setting, advocated for disclosure of performance properties for
17 complex designs.¹⁶ We agree, and characterized our original randomized
18 confirmatory basket trial in this respect, finding that it did not control FWER with
19 acceptable power levels (data not shown). It may still be suitable in those
20 confirmatory settings where control of FWER is not recommended.
21
22
23
24
25
26
27
28
29
30

31 In order to control FWER, and also maintain acceptable power levels, we
32 implemented a modification of the initial design in which, whenever the pooled
33 result is positive, each indication is re-checked at low to medium stringency for
34 statistical significance before final approval (Figure 1B and C, Methods). The
35 series of tests, each at lower stringency (higher alpha) is sufficient to control
36 FWER to a prespecified level more stringent than each of the individual checks. As
37 has been shown in biological systems where high fidelity with minimum energy
38 expenditure is required^{17,18} repeated testing at lower stringency is more efficient
39 than a single highly stringent test.
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 We describe the performance of this modified design in extensive simulations,
10 focusing on a scenario with the same TTE endpoint at interim and final analyses.
11
12 By varying design parameters (Table 1, Methods, Supplemental Methods), we
13
14 create numerous design variants. We select the variant that maximizes
15
16 performance, generally judged by relative efficiency. Acceptable design
17
18 parameters must control FWER to the same level as a system of individual
19
20 randomized studies for each indication in parallel, i.e. approximately 0.025
21
22 multiplied by the number of indications. Further, acceptable design parameters
23
24 must not introduce bias of greater than 10% in the effect size estimate, and the
25
26 estimated 95% confidence interval of the effect size must cover at least 90% of
27
28 simulation runs. For the input parameters to be considered acceptable, these
29
30 constraints must be met regardless of the number of inactive indications within the
31
32 basket.
33
34
35
36
37
38
39

40 For design parameters meeting these constraints we characterize development
41
42 efficiency and power as estimated by simulation. We define development
43
44 efficiency as the expected number of true positive indications identified divided by
45
46 the expected sample size. “Uncorrected” relative development efficiency is the
47
48 efficiency of the basket design divided by that of a group of parallel traditional
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 confirmatory studies, powered at 90%, investigating the same indications. As we
10 do not achieve 90% power by indication in the basket trial, we present a
11 “corrected” relative development efficiency adjusting for the power losses as
12 described above. As enhanced development efficiency can be achieved in
13 conventional designs by reducing power alone,² corrected relative development
14 efficiency compares the efficiency of the basket design to parallel conventional
15 studies at the same power.
16
17
18
19
20
21
22
23
24
25

26 Power is evaluated by indication (Figure 2C), a more stringent criterion than
27 traditional confirmatory studies, where subgroups are usually not formally
28 powered. Power is therefore the fraction of active indications expected to be
29 qualified for approval by the basket trial. For comparison, we also present the
30 traditional power of the pooled analysis.
31
32
33
34
35
36
37

38 We present in Results selected designs that optimize corrected relative efficiency,
39 either subject to a minimum power constraint (recommended development
40 scenarios) or irrespective of power. These designs are characterized as a function
41 of the number of total indications, the number of active indications, and the degree
42 of activity (hazard ratio 0.5 – 0.8).
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 We discuss overall utility of the design and key learnings for performance
10 optimization. We propose criteria for when control of FWER should and should not
11 be recommended by health authorities in confirmatory studies. Finally, we outline
12 future research topics aimed at addressing other potential concerns with
13 randomized confirmatory basket trial designs.
14
15
16
17
18
19
20

21 Methods

22 Study Design Overview and Design Parameters

23
24
25 Consider a randomized confirmatory basket trial of an experimental therapy that
26 consists of k tumor indications. For each indication, we perform 1:1 randomization
27 (experimental vs. control), with a TTE variable as the primary endpoint of interest
28 and n events per indication.
29
30
31
32
33
34
35
36
37

38 Figure 1B presents the current study design. We assume an interim analysis is
39 conducted on each tumor indication individually at a common information time t
40 $\in (0,1)$ based on nt events for all tumor indications, which we assume is also
41 the same actual time. The study designer should choose sample sizes to make
42 this approximately true, but may need to conduct several interim analyses and
43 sample size adjustments in practice. We assume a common interim bar α_t (one-
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

sided nominal Type I error rate) across tumor indications for simplicity. A tumor indication is “pruned” from the study at interim analysis if it does not meet the bar for pooling. Remaining indications proceed to the pooled analysis. After the interim analysis, we adjust the sample size for the remaining tumor indications to account for lost sample size due to pruning. We consider three sample size adjustment designs^{13, 14}:

Table 1. Glossary of Terms

Design parameters	Descriptions
The state of nature	
g	The number of active indications at the beginning
θ_i	The hazard ratio of experimental arm vs. control arm for the definitive endpoint for indication i
Study design input parameters	
k	The number of tumor indications at the beginning
D	Sample size adjustment strategies
α_t	A common bar for the interim analysis to prune inactive indications
α	False positive rate for the pooled analysis
β	False negative rate for the pooled analysis
α_{post}	A common bar for the post-individual tests
Outcome measurements	
m	The number of indications in the pooled analysis. $m \leq k$
α_{net}	The probability of the basket trial passing the pooled test and at least one false positive indication passing the post-individual test for a given value of g
Family-wise error rate (FWER)	The maximum probability of the basket trial passing the pooled test and at least one false positive indication passing the post-individual test for any g . The possible values of g indicate the “family” for which FWER is defined.

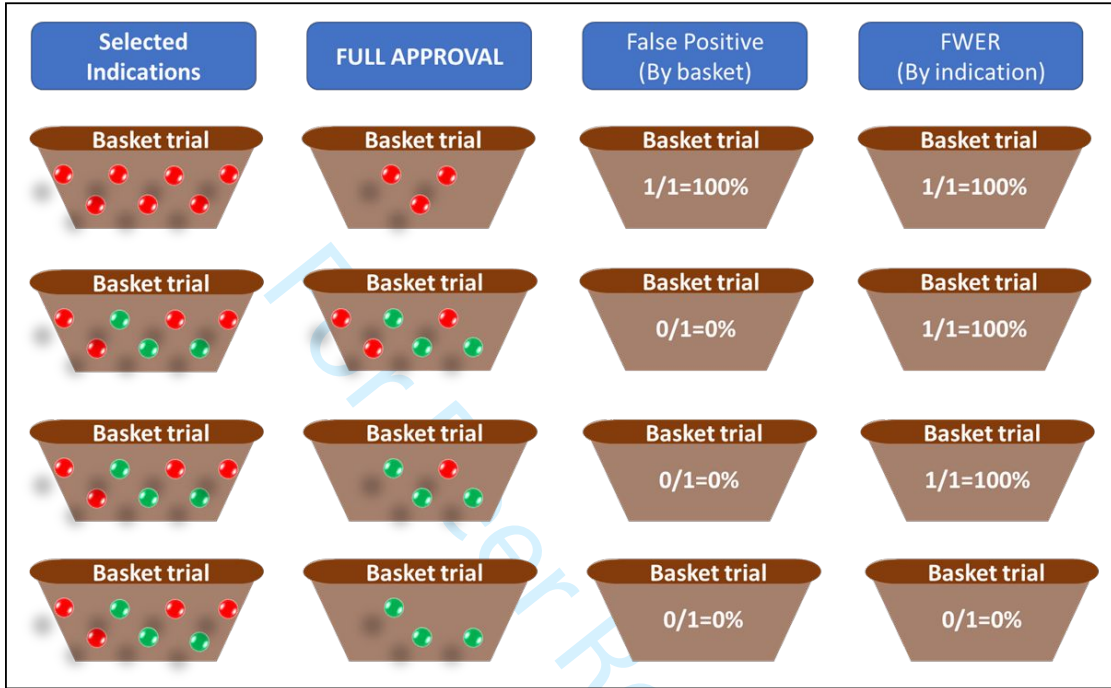
Power by indication	The probability of an active indication passing the interim test, the pooled test, and the post individual test
Power by basket	The probability of a true positive basket (with one or more active indications) passing the pooled analysis
Efficiency	The ratio of average number of active indications that pass the post-individual tests divided by the average sample size
Uncorrected relative efficiency	The efficiency of the basket design divided by that of a group of parallel traditional confirmatory studies, powered at 90%, investigating the same indications.
Corrected relative efficiency	The efficiency of the basket design divided by that of a group of parallel conventional studies at the same power by indication as observed for the basket design, investigating the same indications.
Negative predictive value (NPV) by indications	The proportion of true negative indications among all negative indications not passing the interim test, the pooled test, or the post individual test
Negative predictive value (NPV) by basket	The proportion of true negative baskets among all negative baskets without any indication passing the interim test, the pooled test, or the post individual test
Positive predictive value (PPV) by indications	The proportion of true positive indications among all positive indications passing the interim test, the pooled test, and the post individual test
Positive predictive value (PPV) by basket	The proportion of true positive baskets among all positive baskets with at least one indication passing the interim test, the pooled test, and the post individual test

1. Design one (D1): No sample size adjustment.
2. Design two (D2): Aggressive sample size adjustment to replace all originally planned events in the pruned indications.
3. Design three (D3): Moderate sample size adjustment to replace future originally planned events after the interim analysis in the pruned indications.

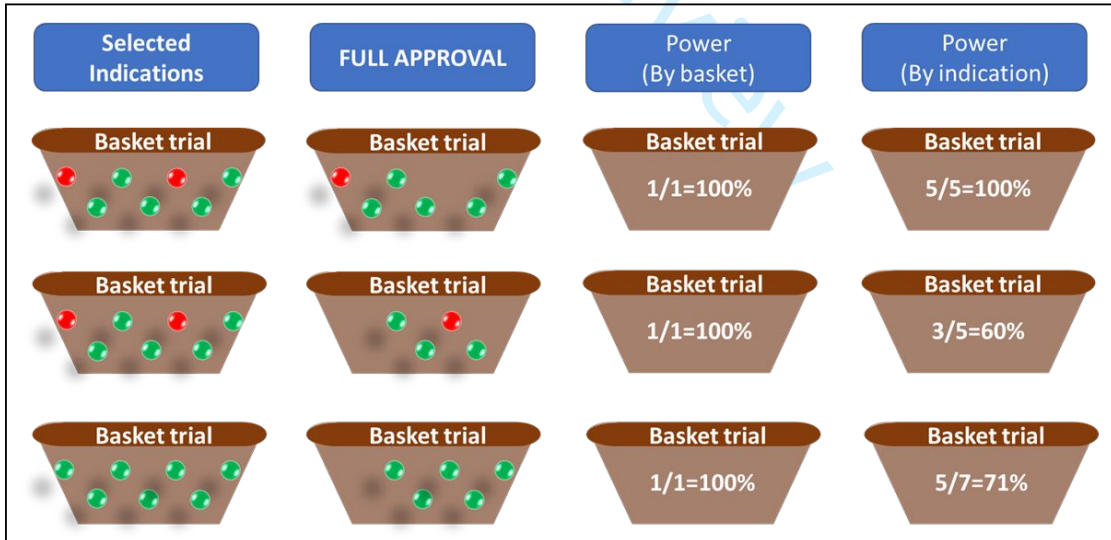
1
2
3
4
5
6
7
8
9 Although endpoints for pruning and pooling may be different, in this work we
10 consider the same endpoints only. Denote (α, β) as the false positive and negative
11 rates for one-sided hypothesis testing in the pooled population. The adjusted false
12 positive rate α^* to control the false positive rate for the pooled analysis at the
13 desired level will be calculated (see Supplemental Methods) for each strategy with
14 respect to the global null hypothesis that all indications are inactive (Figure 2A,
15 Supplemental Methods).^{13,14} Suppose m tumor indications are included in the
16 pooled analysis ($m \geq 1$). When a basket passes the pooled analysis, a
17 prospectively defined individual post-pool analysis, examining each of the m tumor
18 indications remaining in the pool after all pruning is complete, is conducted,
19 termed a “post-individual test”. We also assume a common prospective post-
20 individual bar α_{post} , varied independently of α_t . Indications that survive the post-
21 individual test may be eligible for full approval.
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

A



B

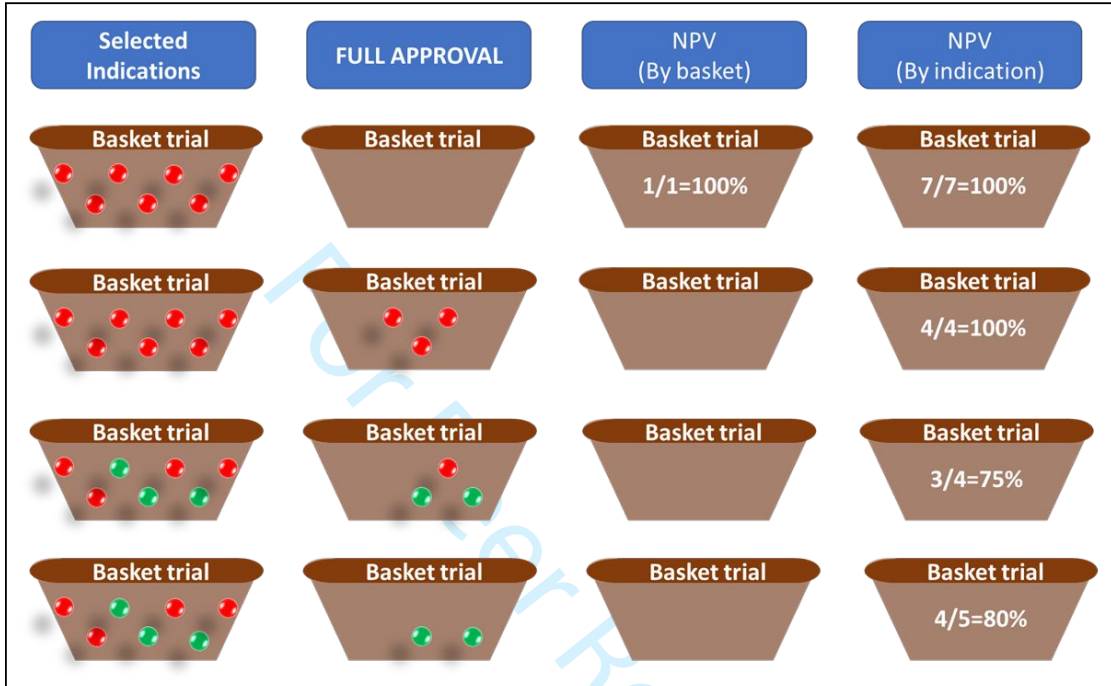


1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

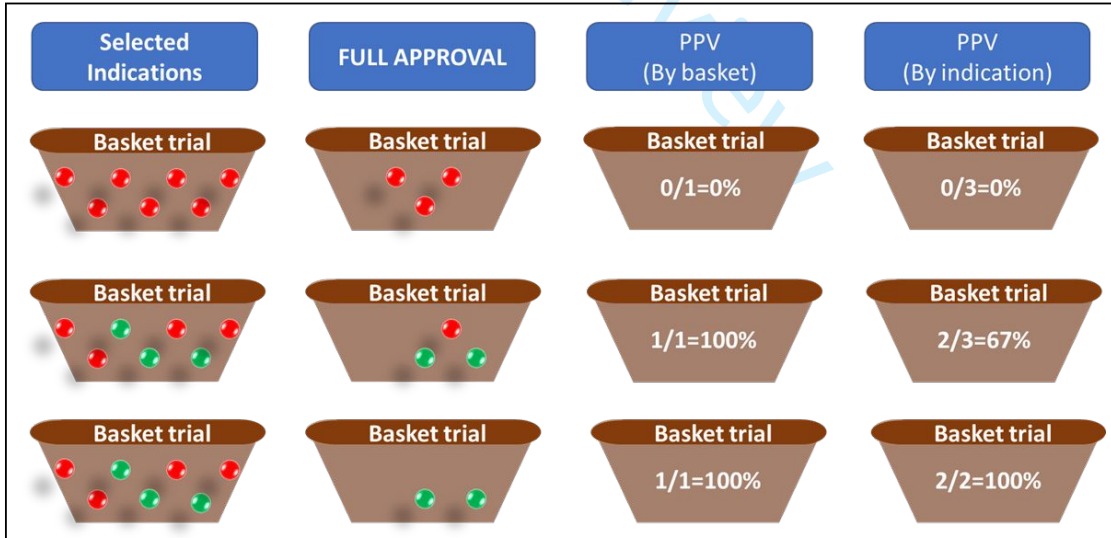
Figure 2. Measurements of false positive rate (type I error) (A); and power (B) showing the difference between criteria used in this study and less stringent criteria. Examples are shown for illustration. Each row represents an example and marbles represent the active (green) and inactive (red) indications. The left-hand column represents the initial selected indications. The second column represents the approved indications. The third column represents traditional criteria (by basket), and the fourth column represents the more stringent criteria (by indication) used in this study. (A) This study is designed to control the family-wise error rate (FWER) by indication subgroup (fourth column) rather than the false positive rate in the pool against a null hypothesis where all indications are inactive. In the third column, a false positive is scored only when a basket is approved containing only inactive indications (false positive rate by basket, as in the original design^{13,14}. Thus, the numerator for FWER is the number of approved baskets with zero active indications and the denominator is the number of approved baskets. In the fourth column, a false positive is declared if a single inactive indication is approved. Thus the numerator for FWER is the number of approved baskets with one or more inactive indications and the denominator is the number of approved baskets. Note the term “basket” means the basket which contains a collection of marbles (indication). (B) Power results are evaluated by indication in this study rather than by basket. Considering the power by basket (third column), the approved basket can be considered as a true positive if it has at least one active indication. Thus numerator is the number of “selected indications” baskets containing one or more active indications that result in an approval, and the denominator is the number of selected indications baskets containing one or more active indications. In the fourth column, the power by indication is defined as the percentage of active indications approved, and can be calculated as 100%, 60%, and 71% for the three examples, respectively.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

A



B



1
2
3
4
5
6
7
8
9 Figure 3. Measurements of negative predictive value (NPV) (A); and positive predictive value (PPV) (B)
10 showing the difference between criteria by indication and by basket. See definitions of NPV and PPV in Table 1.
11 Examples are shown for illustration. Each row represents an example and marbles represent the active (green)
12 and inactive (red) indications. Baskets denote the collection(s) of indications within the container(s) shown. The
13 left-hand column represents the initial selected indications. The second column represents the approved
14 indications. The third column represents the NPV/PPV (by basket), and the fourth column represents the
15 NPV/PPV (by indication) used in this study. (A) In the third column, a NPV is scored only when a basket is not
16 approved, and a basket is considered negative only if all indications are negative. Thus the numerator is the
17 number of baskets not approved that have all indications negative and the denominator is the number of
18 baskets not approved. In the fourth column, a NPV is calculated based on the indications that are not approved
19 and whether they are active or inactive. For baskets that are not approved, the numerator is the number of
20 inactive indications in the basket, and the denominator is the number of indications in the basket. (B) PPV
21 results are evaluated by indication in this study rather than by basket. In the third column, a PPV by basket is
22 scored only if a basket is approved. A basket is considered positive if it has at least one active indication. Thus
23 the numerator is the number of approved baskets that have at least one active indications and the denominator
24 is the number of approved baskets. A PPV by indication is defined if one or more indications are approved. In
25 the fourth column, the PPV by indication is defined as the percentage of active indications approved among all
26 approved indications, and can be calculated as 0%, 67%, and 100% for the three examples, respectively.
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Type I error evaluation

In this basket trial setting (Figure 1B) designed to examine multiple indications (k in number), there is no clear analog of conventional Type I error rates if we consider Type I error by indication, since the possibility of committing a Type I error may occur for tests of each indication. We consider the familywise error rate (FWER) by indication, which is defined as the probability of at least one false positive indication getting approved irrespective of the number of active and inactive indications, defined respectively as indications in which the drug provides or does not provide clinical benefit in the unknown state of nature. This FWER considers a family of null hypotheses in which one or more of the k indications are inactive, i.e. $2^k - 1$ null hypotheses. The FWER is progressively controlled by three successive pruning steps, none of which individually provide Type I error control as stringently as all three in sequence. Considering a basket trial that consists of k tumor indications, these three analyses must be passed for an indication to be approved:

1. Interim analysis: For each of k tumor indications independently, an interim analysis prunes (removes) inactive indications. We assume a common bar

1
2
3
4
5
6
7
8
9 α_t (in terms of the one-sided nominal Type I error rate) for all tumor
10 indications for simplicity.

- 11
12
13
14 2. Pooled analysis: For all remaining indications that pass the interim analysis,
15 a pooled analysis is performed relative to the null hypothesis that all
16 indications in the pool are inactive. The adjusted nominal level α^* is used
17 (further details in Supplemental Methods).^{13,14}
18
19
20
21
22
23 3. Post-individual check: For each indication that passes the pooled analysis,
24 a prospectively defined post-individual check determines whether an
25 indication may be eligible for full approval. We assume a common bar α_{post} ,
26 which is varied independently of α_t .
27
28
29
30
31
32

33 The three analyses are conducted at the nominal levels α_t , α^* , and α_{post} , however,
34 none of these nominal levels quantifies the Type I error of the entire trial. Rather
35 than control the false positive rate of any of three tests at the level of 0.025 one-
36 sided per indication, we evaluate the FWER of the entire trial according to the
37 cumulative effect of three sequential tests. To evaluate the FWER, we design
38 comprehensive simulation studies, where a range of nominal levels α_t and α_{post} is
39 selected (Supplemental Table S1) and the adjusted nominal level for testing the
40 pool α^* is calculated for each strategy. For a given set of design parameters and
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 hazard ratio of active indications, the overall Type I error resulting from the
10 cumulative effect of the 3 steps for a given value of g is $\alpha_{net}(g)$, where g is the
11 number of active indications ($0 \leq g \leq k$). We then vary g from 0 to k with all
12 other parameters constant to get the corresponding value of FWER:
13
14
15
16
17
18

$$19 \quad \text{FWER} = \max\{\alpha_{net}(g): g = 0, 1, \dots, k\}$$

20
21
22 FWER is considered controlled to the level of α_{target} if $\alpha_{net} \leq \alpha_{target}$ for a given set
23 of input parameters and all g from 0 to k , specifically, $\alpha_{target} = 0.025k$ in this study.
24
25
26

27 **Outcome Measurements**

28
29
30 We consider the power by indication, which evaluates the proportion of active
31 indications passing the post-individual tests. We also examine the efficiency,
32 defined as the ratio of average number of true positive indications passing the
33 post-individual tests divided by the average sample size, when subject to control of
34 the **FWER** by indication. Other measurements include the coverage of the
35 confidence interval (CI) of the hazard ratio (HR) and bias of the estimated HR. In
36 evaluating outcome measurements, we threshold the FWER at the level $0.025k$,
37 the coverage of the 95% CI of the HR as greater than 90% of the simulation runs,
38 and the bias of estimated HR less than 10%. Design parameter combinations
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 which cannot meet these criteria irrespective of the value of g in our simulation are
10 not allowed. Finally, we examine power by indication and efficiency for each
11 allowed design and utilize a ratio to compare efficiency relative to a reference
12 design.
13
14
15
16
17
18

19 The reference design for the uncorrected relative efficiency calculation assumes
20 parallel, independent Phase 3 designs planned for each indication with the false
21 positive and false negative rates $(\alpha_{ref}, \beta_{ref})$ set to $(0.025, 0.1)$.
22
23
24
25

26 For calculation of relative efficiency corrected for power losses, we first determine
27 the power by indication of a basket study with the same design parameters and
28 inputs, then set β_{ref} equal to 1 minus this power, while maintaining α_{ref} at 0.025,
29 and finally proceed as for the calculation of relative efficiency above. This
30 correction is necessary because the efficiency of the reference designs are higher
31 when run at lower power². Therefore, comparison to a reference design at the
32 same power is a more stringent, balanced relative efficiency comparison.
33
34
35
36
37
38
39
40
41
42

43 Other measurements include the coverage of the confidence interval for HR and
44 estimation bias for HR. Specifically, the CI coverage for HR is defined as the
45 probability that the estimated 95% CI for the HR covers the true HR given that the
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 individual test passed for that indication. Representing the relative difference in the
10 true HR and the estimated HR, the estimation bias of HR is defined as the ratio of
11 estimated average HR and true pooled HR for those indications that pass the
12 individual tests minus 1.
13
14
15
16
17

18 **Simulation study**

19
20 In our simulation study, we used parameter values that are summarized in
21 supplemental Table S1. We assume that for each indication the true hazard ratio
22 (HR) $\theta_i; i = 1, \dots, k$ is either at a null value $\theta_0 = 1$ or at an active value $\theta_a \in$
23 $\{0.5, 0.6, 0.7, 0.8\}$ and the true number of active indications is $g = 0, 1, \dots, k$. For
24 simplicity, we consider an exponential model for the distribution of event times and
25 do not consider censoring. We fix $\alpha = 0.025$ and vary design parameters
26 $k (3, 4, 5, 6)$ and $\beta (0.025, 0.05, 0.1, 0.2)$. Consequently, the total sample size in the
27 pooled population (kn) can be calculated as $kn = 4(Z_{1-\alpha} + Z_{1-\beta})^2 / [(\log \theta)^2]$. For
28 simplicity, we set a common information time for the interim analysis $t = 0.5$, so
29 that each indication should accrue nt patients for interim analysis. We note that the
30 in practice the indications should be chosen so that they are projected to reach nt
31 events at approximately the same time to avoid the operational inconvenience of
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

multiple interim analyses and sample size adjustments. Denote n_{i1} ($n_{i1} = nt$; $i = 1, \dots, k$) as the sample size for i -th indication at the interim analysis. The sample size for each indication at pooled analysis, denoted as n_{i2} ($i = 1, \dots, m$), should accrue as follows:

1. $n_{i2} = n$ under D1;
2. $n_{i2} = \frac{kn}{m}$ under D2, which is greater than the sample size under D1;
3. $n_{i2} = n\left(t + \frac{k(1-t)}{m}\right)$ under D3, which is greater than the sample size under D1 but smaller than the sample size under D2.

With these specifications, the total number of patients enrolled in the study can be calculated as $(k - m)n_{i1} + mn_{i2}$ if each indication has the same planned number of patients. We further explore the design for values of α_t and α_{post} of (0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4), setting these two parameters independently. For each setting, we generate 10000 simulated trials for the evaluation and comparison.

1
2
3
4
5
6
7
8
9 Simulations and analysis are performed using R (version 3.6). Source code and
10 details of simulation settings and outcome measurements are available in
11 Supplemental R codes.
12
13
14
15

16 RESULTS

17
18
19 We explore power and efficiency with various input parameters (Table 1,
20 Supplemental Table S1). We control FWER at $\text{FWER} \leq 0.025k$ for all values of
21 g , the number of active indications, from 0 to k inclusive. FWER is more difficult to
22 control with increasing number of tumor indications and increased therapeutic
23 effect size, reflecting greater heterogeneity (examples shown in Figure 4; more
24 details included in Supplemental Table S2). **FWER** α_{net} is generally maximal when
25 g is approximately half of k , again reflecting maximal heterogeneity.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

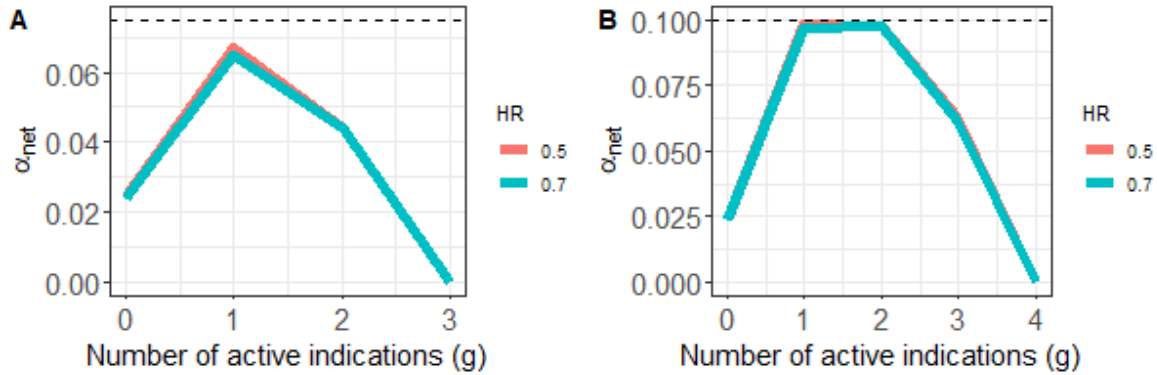


Figure 4. Controlled FWER at $\text{FWER} \leq 0.025k$ for all values of g ($g = 0, \dots, k$). (A) α_{net} is shown on the y axis, and the number of active indications on the x axis. $\alpha_{net} = \text{FWER}$ is controlled at $\leq 0.025k$ (dotted horizontal line) for $k=3$, $HR = 0.5$ (orange; Design 3, $kn = 128$, $\beta = 0.025$, $\alpha_t = 0.35$, $\alpha_{post} = 0.05$, $\alpha^* = 0.0089$) and $HR = 0.7$ (blue; Design 3, $kn = 409$, $\beta = 0.05$, $\alpha_t = 0.4$, $\alpha_{post} = 0.05$, $\alpha^* = 0.009$). (B) same as (A), for $k = 4$, $HR = 0.5$ (orange; Design 2, $kn = 128$, $\beta = 0.025$, $\alpha_t = 0.2$, $\alpha_{post} = 0.1$, $\alpha^* = 0.0094$) and $HR = 0.7$ (blue; Design 3, $kn = 483$, $\beta = 0.025$, $\alpha_t = 0.2$, $\alpha_{post} = 0.1$, $\alpha^* = 0.0075$).

We provide recommended optimal design parameters and associated performance results (Figure 5, Supplemental Figure S1). These optimal design parameters determined by simulation depend on the scenario studied and thus vary from panel to panel in the Figures. The legends list these parameters. These recommendations consider both corrected relative efficiency and power. Although corrected relative efficiency is arguably

1
2
3
4
5
6
7
8
9 the most fundamental metric, many sponsors will prefer to have power within an
10 “acceptable” range for a confirmatory study. We also provide design parameters that
11 optimize corrected relative efficiency irrespective of power (Figure 6, Supplemental Figure
12 S2) and show associated performance results. All simulation results are available in
13 Supplemental Table S2.
14
15
16
17
18
19

20 **Decreased power and/or efficiency if the majority of indications are inactive**

21
22
23 Power and/or efficiency increases with the proportion of indications that are active.
24 (Figures 5-6, Supplemental Figures S1-S2). Performance deteriorates when
25 multiple inactive indications are present in the basket.
26
27
28
29

30
31 As the hazard ratio of positive indications increases, indicating a decreased
32 therapeutic effect size, the corrected relative efficiency of the recommended
33 development scenario generally increases, but the power of the recommended
34 development scenario generally decreases. This pattern suggests a tradeoff
35 between power and efficiency.
36
37
38
39
40
41
42

43 Maximal corrected relative efficiency is seen at $k = 3,4$.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

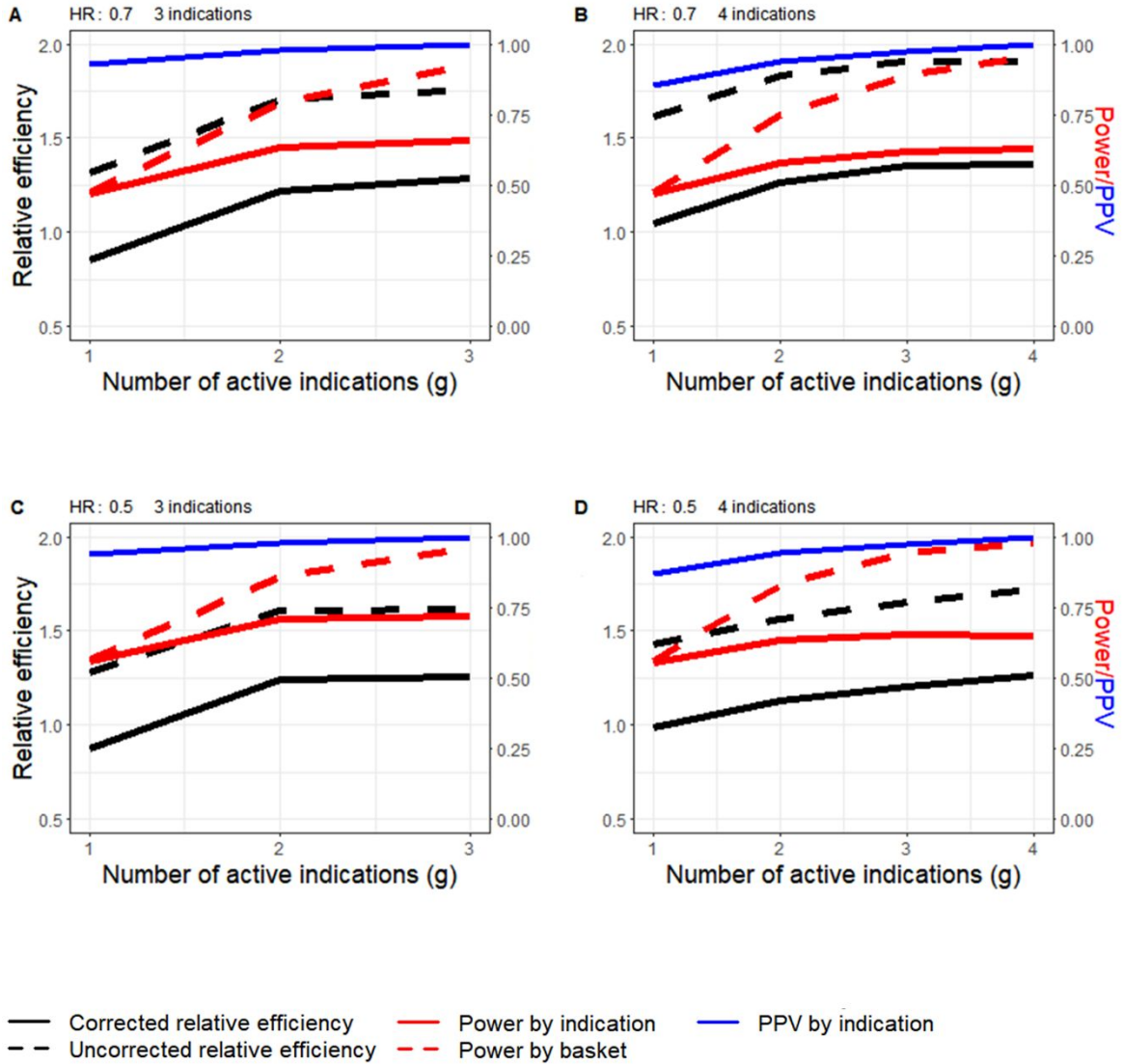


Figure 5. Recommended development approaches for (A) 3 indications with $HR = 0.7$ (Design 3, $kn = 409$, $\beta = 0.05$, $\alpha_t = 0.4$, $\alpha_{post} = 0.05$, $\alpha^* = 0.009$), (B) 4 indications with $HR = 0.7$ (Design 3,

1
2
3
4
5
6
7
8
9 $kn = 483, \beta = 0.025, \alpha_t = 0.2, \alpha_{post} = 0.1, \alpha^* = 0.0075$), (C) 3 indications with $HR = 0.5$ (Design 3,
10 $kn = 128, \beta = 0.025, \alpha_t = 0.35, \alpha_{post} = 0.05, \alpha^* = 0.0089$), and (D) 4 indications with $HR = 0.5$
11 (Design 2, $kn = 128, \beta = 0.025, \alpha_t = 0.2, \alpha_{post} = 0.1, \alpha^* = 0.0094$). The x-axis represents the
12 number of active indications (indications in which the drug provides clinical benefit), the primary y-
13 axis (left) represents the uncorrected/corrected relative efficiency, and the second y-axis (right)
14 represents the power (red) by indication and by basket, and the positive predictive value by
15 indication (blue).
16
17
18
19
20
21
22
23
24
25
26

27 **Recommended development scenarios**

28
29 We have previously emphasized that this design requires careful initial indication
30 selection.^{13,14} Users should strive to include only zero or one inactive indications in
31 the study (Discussion). To determine recommended development scenarios, we
32 considered these upside scenarios, specifically requiring power by indication
33 greater than 60% when all indications are active, assuming that lower power might
34 not be considered acceptable by some practitioners for a confirmatory study. We
35 assumed that for the very modest therapeutic effect sizes considered herein,
36 corresponding to effective but not transformational therapies, that a reduction in
37 power from the traditional 80 – 90% to 60% would be potentially applicable given
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 the marked increase in development efficiency offered, in settings where financial
10 resources or available trial participants were limiting. In practice, the practitioner
11 may use these simulation methods to set their own desired power cutoff, and
12 higher power cutoffs may be possible for larger effect sizes. We simulate a range
13 of values for k , β , D , α_t , and α_{post} . We note that α^* is not independently
14 adjustable, but is calculated based on the other input parameters and the
15 assumed hazard ratios using previous methods (Supplementary Material).¹³ After
16 implementing this filtering criterion for power in addition to the requirements for
17 control of FWER, estimation bias, and confidence interval coverage at all values of
18 g , only the scenarios with $k = 3, 4$ remain. Figure 5 summarizes the scenarios
19 meeting this minimal power criterion with optimal upside corrected relative
20 efficiency for $k = 3, 4$ and hazard ratio values of 0.5 and 0.7. Figure S1 presents
21 the same information for hazard ratios 0.6 and 0.8. Power by basket increases
22 from 50% to greater than 90% for any scenario as the number of active
23 indications increases from 1 to k , while the power by indications increases to 63%
24 –72% when all indications are active, and 62% –71% when there is one inactive
25 indication. NPV by indication decreases from 91% at $g = 1$ to 71% and 55% for
26 $g = k - 1$ for $k = 3$ and 4, respectively. PPV by indication increases from 86%

1
2
3
4
5
6
7
8
9 $(k = 4) - 93\%$ $(k = 3)$ to 98% as g increases from 1 to $g = k - 1$ for $k = 3, 4$, but
10
11 remains stable whether the hazard ratio of positive indications increases or
12
13 decreases. For the uncorrected relative efficiency, the relative efficiency is
14
15 compared to a group of parallel independent studies at 90% power, yielding 90%
16
17 average power overall. Scenarios with 4 indications exhibit about 43% - 92%
18
19 efficiency improvement depending on g , 23% - 78% efficiency improvement for
20
21 three indications. For the corrected relative efficiency, the basket scenarios are
22
23 compared to k independent trials with the same power as achieved in the basket
24
25 scenario, so that the average power is the same across the comparison.
26
27
28 Corresponding ranges of corrected relative efficiency improvement are - 17% - 29
29
30 %, and - 1% - 38%, for $k = 3, 4$, respectively. Note if most of the indications are
31
32 inactive, the design is inefficient. When all indications are active, or only one
33
34 indication is inactive, the uncorrected efficiency improvement ranges from 61%
35
36 - 78% and 65% - 92% for 3 and 4 indications, respectively, while the corrected
37
38 improvement in efficiency ranges from 21% - 29% and 20% - 38% for 3 and 4
39
40 indications, respectively.
41
42
43
44
45
46

47 **Optimal corrected relative efficiency scenarios**

48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 Without filtering for power, Figure 6 and S2 summarize scenarios with maximal
10 corrected relative efficiency for $k = 3, 4, 5$ and different values of the hazard ratio.
11
12 Scenarios with $k = 6$ have inferior performance (Supplemental Table S2) and are
13
14 not recommended. In the discussion, we consider explanations for the optimal
15
16 indication number range.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

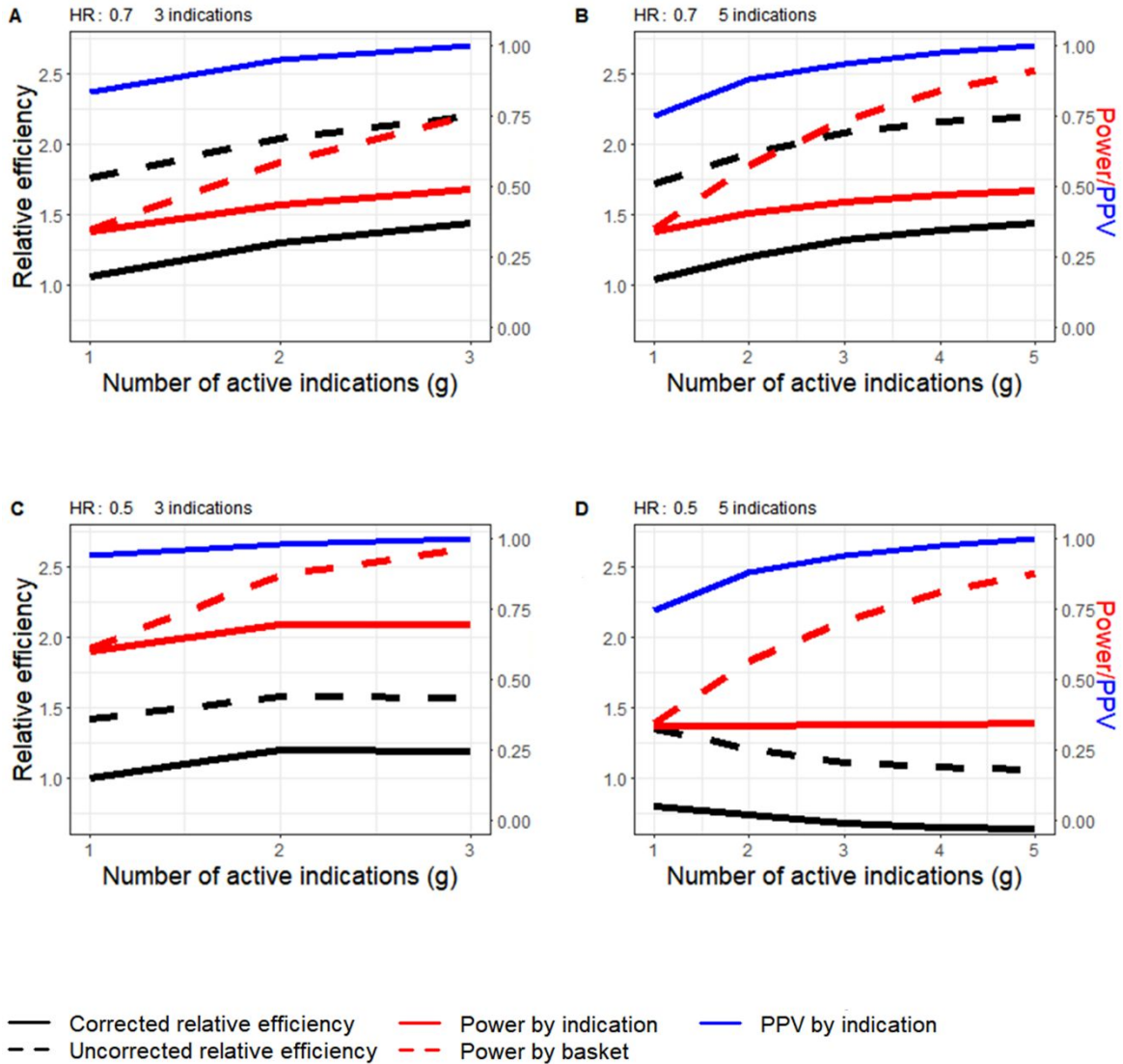


Figure 6. Cases with maximum corrected relative efficiency for (A) 3 indications with $HR = 0.7$ (Design 3, $kn = 247$, $\beta = 0.2$, $\alpha_t = 0.2$, $\alpha_{post} = 0.15$, $\alpha^* = 0.0101$), (B) 5 indications with $HR = 0.7$

1
2
3
4
5
6
7
8
9 (Design 3, $kn = 409$, $\beta = 0.05$, $\alpha_t = 0.15$, $\alpha_{post} = 0.15$, $\alpha^* = 0.0071$), (C) 3 indications with $HR = 0.5$
10 (Design 3, $kn = 128$, $\beta = 0.025$, $\alpha_t = 0.25$, $\alpha_{post} = 0.05$, $\alpha^* = 0.0093$), and (D) 5 indications with
11 $HR = 0.5$ (Design 2, $kn = 128$, $\beta = 0.025$, $\alpha_t = 0.05$, $\alpha_{post} = 0.4$, $\alpha^* = 0.0278$). The x-axis represents
12 the number of active indications (indications in which the drug provides clinical benefit), the primary
13 y-axis (left) represents the uncorrected/corrected relative efficiency, and the second y-axis (right)
14 represents the power (red) by indication and by basket, and the positive predictive value by
15 indication (blue).
16
17
18
19
20
21
22
23
24
25
26

27 Comparing to the recommended development scenarios, the optimal corrected
28 relative efficiency scenarios present higher efficiency, lower power, and lower
29 NPV/PPV by indication. Power by basket ranges from 34% –98%, 26% –98%,
30 and 21% –98%, depending on g , for $k = 3,4,5$, respectively. Corresponding power
31 ranges for power by indication are 34% –71%, 25% –65%, and 20% –59%. PPV
32 by indication ranges from 83% –100%, 71% –100%, and 55% –100%, for
33 $k = 3,4,5$, respectively, which are lower than those of recommended scenarios.
34
35
36
37
38
39
40
41
42

43 Uncorrected relative efficiency improvement ranges from 25% – 125%,
44 47% – 163%, and 6% – 172% depending on the value of g for $k = 3,4,5$
45 respectively. Corresponding ranges of corrected relative efficiency improvement
46 are - 9% –47%, - 6% –67%, and - 36% –66%. When all indications are active or
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 only one indication is inactive, the uncorrected efficiency improvement ranges from
10 46% –125%, 65% –163%, and 6% –172%, for $k = 3, 4, 5$, respectively, while the
11
12 corresponding corrected relative efficiency improvement ranges from 14% –47%,
13
14
15 20% –66%, and - 36% –66% . Removing the constraint on power increases
16
17 relative efficiency gains, again indicative of the tradeoff between power and
18
19 efficiency. In most cases, the basket trial design improves the efficiency, but this is
20
21 not always the case (Figure 6D). This figure considers a case in which the
22
23 potential level of heterogeneity is too great for the design to efficiently compensate
24
25 for, in that there are a large number of indications (5) and active indications are
26
27 postulated to be quite different from inactive ones (hazard ratio of 0.5 compared to
28
29 1.0). Controlling the Type I error by indication for all possible values of g in this
30
31 situation forces a very stringent interim check, resulting also in the elimination of
32
33 many true positives.
34
35
36
37
38
39
40
41
42
43
44

45 DISCUSSION

46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 Confirmatory basket trials potentially provide remarkable improvements in drug
10 development efficiency. With pooling, a fold-improvement in efficiency comparable
11 to the number of indications is possible while retaining high power for the pool, and
12 controlling alpha with respect to a global null hypothesis, in the randomized
13 confirmatory basket design.^{13, 14}

14
15
16
17
18
19
20
21 We studied a modification of the randomized confirmatory basket design^{13, 14} in a
22 more rigorous setting where control of FWER by indication subgroup is
23 recommended, as may be needed for health authority approval in some
24 instances.¹⁵ We have further studied the most challenging scenario, i.e. a slowly
25 maturing TTE endpoint without a highly predictive surrogate. Performance
26 characteristics for innovative designs should be publicly disclosed in detail,
27 including in challenging settings, but this is not always the case.¹⁶ Under the
28 conditions of the simulation, up to 92% improvement in relative efficiency, or 38%
29 improvement corrected for reduced power by indication, is still possible, while
30 controlling the **FWER** at a rate comparable to an equal number of parallel single-
31 indication studies. Power of the pooled analysis (“power by basket”) remains high,
32 and estimation bias and confidence interval coverage may also be well controlled.
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49 Power by indication can be as high as 60% – 73%, but declines if there are two or
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 more inactive indications in the basket. These results are possible only by using
10 simulation to extensively optimize various design parameters, as described in this
11 work. The results apply only to the conditions of the simulation. As there are an
12 infinite number of cases to simulate, a sponsor would need to reach agreement
13 with health authorities on the range of scenarios to be simulated, as has been
14 articulated previously for complex innovative designs.¹⁹

15
16
17
18
19
20
21
22
23 We previously published a related design which controls Type I error *by basket* in
24 the pooled data only^{13, 14} and while a formal efficiency analysis was not performed,
25 an application example suggested 200-500% improvement in efficiency. The
26 current variant, that considers FWER by indication subgroup (see Figure 2 for the
27 definition of control by indication subgroup as distinguished from control by
28 basket), is less efficient than the original, but still substantially improved compared
29 with independent parallel studies of each indication. Both the original and current
30 designs contain one or more checkpoints that operate on individual indications, but
31 also contain a checkpoint that operates on pooled data, and this checkpoint, the
32 most stringent one, is far more efficient due to the pooling, which allows multiple
33 indication subgroups to be considered for the sample size that would normally be
34 required for one indication. In the present variant, the Type I error is actually

1
2
3
4
5
6
7
8
9 controlled over the whole study by indication subgroup at 0.025k, reflecting
10 checkpoints some of which operate on pooled data, rather than individually at
11 0.025 per indication subgroup. Only in the example in which all indications are
12 treated identically would the indication subgroups all be controlled at 0.025
13 individually.
14
15
16
17
18
19
20

21 The optimal range for indication number is narrow. For recommended drug
22 development scenarios, which seek optimal corrected relative efficiency with a
23 minimum constraint on power by indication, either 3 or 4 indications are optimal. If
24 only corrected relative efficiency is optimized, one may also consider 5 indications.
25
26 Too few indications and one does not get the benefit of a basket trial. Too many,
27 and compensating for the large number of potential heterogeneity scenarios
28 involved in controlling FWER becomes challenging. Interestingly, earlier work
29 determining an optimal indication number in exploratory basket trials also
30 recommended 3–5 indications despite a very different scenario and approach to
31 optimization.²⁰
32
33
34
35
36
37
38
39
40
41
42
43
44

45 Optimal stringency of filtering varies, but generally greater stringency is applied in
46 post-individual tests than at interim. Higher stringency is perhaps better applied to
47 more mature datasets. Optimal sample size adjustment (SSA) strategies varied.
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 Sometimes moderate SSA (D3) was better than aggressive SSA (D2). Aggressive
10 SSA may optimize power at the expense of overall development efficiency,
11 another example of a case where focusing on power leads to diminishing returns
12 when considered from a portfolio perspective.²¹ The optimal nominal power for the
13 pooled analysis also varies according to the power/efficiency tradeoff. Higher
14 power requires disproportionate investment, as previously shown for proof of
15 concept studies.²

16
17
18
19
20
21
22
23
24
25
26 A study of ten oncology drugs found an average cost of \$648 million for their
27 clinical development, much of which spent in the resource-intensive confirmatory
28 phase.²² Expense and prolonged clinical development time delays availability of
29 therapies, and contributes to their high cost and potentially to unequal treatment
30 access, major public health issues. We believe the confirmatory basket design or
31 modifications thereof could contribute to the solution of these problems. Moreover,
32 in oncology niche indications, this design could make drug development
33 economically feasible even in the absence of a transformational therapy.

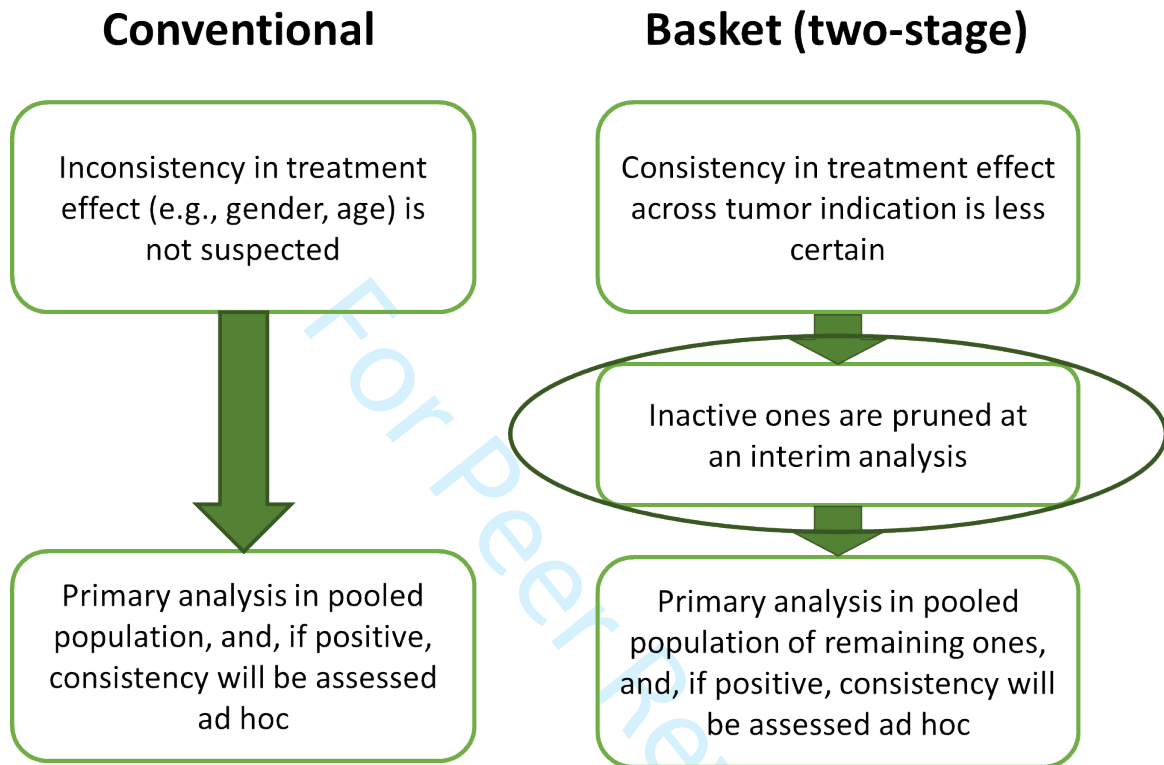
34
35
36
37
38
39
40
41
42
43
44 Transformational approaches such as immunotherapy must be further optimized,
45 creating a broad universe of potential combination studies across populations
46 sharing common characteristics such as high tumor mutation burdens. In this
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 example, immune checkpoint therapy could be combined in the experimental arm
10 with another drug designed to improve its performance in multiple tumors. These
11 therapy combinations may lead to important but incremental improvements that
12 require a randomized TTE approach for confirmation. Potential applications are
13 not limited to oncology. We are currently investigating the use of real world data of
14 off-label use to design and simulate a basket trial in multiple autoimmune diseases
15 in which rituximab is added to a standard of therapy arm consisting of steroid
16 therapy.^{23, 24}
17
18
19
20
21
22
23
24
25
26
27

28 We have previously suggested that the performance of randomized confirmatory
29 basket trials depends on careful indication selection.¹⁴ The importance of
30 indication selection is even greater when control of FWER is recommended, as is
31 documented in the results, in which performance improves with the proportion of
32 indications that are active. The more inactive indications are present, the more
33 stringent pruning and/or post check steps are required to reliably eliminate them
34 and control the **FWER** by indication. Highly stringent pruning and/or post check
35 steps run a greater risk of also eliminating some active indications, reducing power
36 and efficiency. The more inactive indications are present, the greater the risk that
37 they will dilute the signal in a pool containing positive indications.
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 These designs are therefore inappropriate for a collection of miscellaneous
10 uncharacterized indications. Any proposed indication should, when feasible, be
11 supported by preclinical studies, and Phase 2 clinical and biomarker data, ideally
12 from randomized biomarker-guided Phase 2 studies.¹⁴ Alternatively, it may be
13 helpful to filter these indications in Phase 2 with one of several exploratory basket
14 trial designs,⁵⁻⁸ especially if the indications have enrollment challenges. Real world
15 data/evidence, especially concerning off-label use, may be of value in screening
16 potential indications.^{23, 24} Ideally, one lead indication should have been confirmed
17 with the biomarker in a biomarker-guided traditional Phase 3 study, and the basket
18 trial would then be confirming supplemental indications using the same biomarker
19 assay adjusted as needed for different tissue types.¹⁴
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

35 For the vast majority of therapies, a positive randomized Phase 2 proof of concept
36 study is followed by a randomized Phase 3 confirmatory study. In analogous
37 fashion, this design and its original predecessor^{13, 14} are therefore expected to be
38 applied for the majority of effective therapies, ideally supported by randomized
39 Phase 2 data. In contrast, previous applications of confirmatory basket trials have
40 been limited to a small number of exceptional therapies.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 7. Parallels between a conventional Phase 3 study (left) and the randomized confirmatory basket trial design considered in previous work without strong control of the **FWER** by subgroup (right). Both designs formally test a hypothesis for a main group and do not formally test hypotheses involving subgroups. In the conventional design, organ site is the defining characteristic of the main group, and both known and unknown subgroups are present, the former perhaps undergoing more informal subgroup analyses. In the original randomized confirmatory basket trial design^{13,14} the traditional organ site classification is only a known subgroup, and a

1
2
3
4
5
6
7
8
9 biomarker or other pathophysiologic feature defines the main group. Organ site classification is
10 subjected to informal analysis only.
11
12
13
14
15

16
17 It is important to consider if and when strong control of FWER by indication
18 subgroup should be recommended for a confirmatory basket trial with pooling, an
19 area of debate which will likely evolve. In cases where control of Type I error by
20 basket is permitted by health authorities, our original design^{13, 14} should be used.
21
22 However, should health authorities require control of Type I error by indication, the
23 design in the current study would be suitable and still provide significantly
24 improved development efficiency relative to conventional approaches. In both
25 cases, the indication would be for a group of organ-specific indications sharing a
26 common biomarker or pathophysiologic mechanism.
27
28
29
30
31
32
33
34
35
36
37

38 One might argue that Type I error control by basket should be adequate in many
39 cases. A conventional Phase 3 study is really quite heterogeneous with respect to
40 both known and unknown subgroups (Figure 7), and we do not control FWER in
41 assessing the vast majority of these subgroups, employing other approaches.^{25, 26}
42
43
44
45
46

47 In a basket trial, we invert the usual classification: molecular subgroups are now
48 indication-defining, and organ sites, formerly indication-defining, are now
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 subgroups. This alone may present a perceptual barrier to full acceptance of the
10 concept. Collignon et al.¹⁵ (an author group representing individuals with past or
11 current associations with European health authorities) argue that homogeneity of
12 outcomes may be difficult to interpret clinically since the populations are “different”,
13 illustrating the unproven perception that differences in organ sites are more
14 fundamental than the many other known and unknown differences between
15 subpopulations. Nonetheless they also provide a definition of subgroup
16 homogeneity in striking agreement with our thinking: “homogeneous if they share
17 important clinical characteristics such that, in light of the available scientific
18 evidence, the interpretation of treatment effect and the assessment of benefit/risk
19 are meaningful for the overarching target population...” This comes down in the
20 end to science and medicine, not statistics, and indeed the scientific justification
21 for pooling must be robust if we are to forego FWER control by subgroup. We
22 know that even in the classic case of an antagonist of a driver mutation in the b-raf
23 gene, drug effectiveness still depends on organ site.¹⁰ Increased understanding of
24 how driver gene mutations interact with tissue-specific gene expression programs
25 may be important. If we do not have a robust justification for pooling, control of
26 FWER by indication will be necessary. Collignon et al.^{15(SI)} evaluate our original
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9 confirmatory basket design,^{13,14} and our assertion that evidence supporting a
10 consistent benefit/risk assessment across traditional indications should be
11 provided at a level prospectively agreed with health authorities.¹⁴ They state that
12 availability of post-approval data would be important. This should become more
13 practical as electronic health records systems improve. Informal interactions with
14 health authorities further indicate that control of the **FWER** by indication is a
15 controversial and evolving issue and may be required in some instances, and
16 therefore there is a practical need for characterizing the performance of a
17 randomized confirmatory basket design constrained to control the **FWER** by
18 indication. The choice between our original design^{13, 14} and the current design can
19 thus be determined only based on discussion about required Type I error control
20 with health authorities.
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35

36
37 In ongoing research, we are considering important questions regarding the
38 suitability of these designs for the confirmatory phase, in particular how to deal
39 with differences between indications in endpoints and in safety issues. In future
40 work, we will consider the effects of different control therapies, enrollment rates,
41 and endpoint maturation rates. We are investigating real world data/evidence in
42 indication screening and parameter estimation in simulations. Finally, current
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

performance may be improved further by application of Bayesian techniques
previously devised for exploratory trials.^{5-8, 27}

For Peer Review

REFERENCES

1. Beckman RA, Clark J and Chen C. Integrating predictive biomarkers and classifiers into oncology clinical development programmes. *Nature Reviews Drug Discovery* 2011; 10: 735-749.
2. Chen C and Beckman RA. Maximizing return on socioeconomic investment in Phase II Proof-of-Concept trials. *Clinical Cancer Research* 2014; 20: 1730-1734.
3. Woodcock J and LaVange LM. Master protocols to study multiple therapies, multiple diseases, or both. *New England Journal of Medicine* 2017; 377: 62-70.
4. Antonijevic Z and Beckman RA. *Platform trial designs in drug development: Umbrella trials and basket trials*. Boca Raton: Chapman & Hall/CRC Biostatistics Series, Taylor & Francis Group, 2018.
5. Cunanan KM, Iasonos A, Shen R, et al. An efficient basket trial design. *Stat Med* 2017; 36: 1568-1579.
6. Berry SM, Broglio KR, Groshen S, et al. Bayesian hierarchical modeling of patient subpopulations: Efficient designs of phase II oncology trials. *Clinical Trials* 2013; 10: 720-734.
7. Chu Y and Yuan Y. A Bayesian basket trial design using a calibrated Bayesian hierarchical model. *Clinical Trials* 2018; 15: 149-158.
8. Simon R, Geyer S, Subramanian J, et al. The Bayesian basket design for genomic-variant driven phase II trials. *Semin Oncol* 2016; 43: 13-18.

- 1
2
3
4
5
6
7
8
9 9. Heinrich MC, Joensuu H, Demetri GD, et al. Phase II open-label study evaluating the
10 activity of imatinib in treating life-threatening malignancies known to be associated with
11 imatinib-sensitive tyrosine kinases. *Clin Cancer Research* 2008; 14: 2717-2725.
12
13
14
- 15 10. Hyman DM, Puzanov I, Subbiah V, et al. Vemurafenib in multiple nonmelanoma
16 cancers with BRAF V600 mutations. *New England Journal of Medicine* 2015; 373:
17 726-736.
18
19
20
- 21 11. Drlon A, Laetsch TW, Kummar S, et al. Efficacy of larotrectinib in TRK fusion-positive
22 cancers in adults and children. *New England Journal of Medicine* 2018; 378: 731-739.
23
24
- 25 12. Le DT, Uram JN, Wang H, et al. PD-1 blockade in tumors with mismatch-repair
26 deficiency. *New England Journal of Medicine* 2015; 372: 2509-2520.
27
28
- 29 13. Chen C, Li N, Yuan S, et al. Statistical design and considerations of a phase 3 basket
30 trial for simultaneous investigation of multiple tumor types in one study. *Statistics in*
31 *Biopharmaceutical Research* 2016; 8: 248-257.
32
33
34
- 35 14. Beckman RA, Antonijevic Z, Kalamegham R, et al. Adaptive design for a confirmatory
36 basket trial in multiple tumor types based on a putative predictive biomarker. *Clinical*
37 *Pharmacology and Therapeutics* 2016; 100: 617-625.
38
39
40
- 41 15. Collignon O, Gartner C, Haidich AB, et al. Current statistical considerations and
42 regulatory perspectives on the planning of confirmatory basket, umbrella, and platform
43 trial. *Clinical Pharmacology and Therapeutics* 2020; 107: 1059-1067.
44
45
46
- 47 16. Cunanan KM, Iasonos A, Shen R, et al. Specifying the true- and false-positive rates in
48 basket trials. *JCO Precision Oncology* 2017; 1:1-5.
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3
4
5
6
7
8
9 17. Freter RR and Savageau MA. Proofreading systems of multiple stages for improved
10 accuracy of biological discrimination. *J Theor Biol* 1980; 85: 95-123.
11
12 18. Beckman RA and Loeb LA. Multistage proofreading in DNA replication. *Quarterly*
13 *Reviews of Biophysics* 1993; 26: 225-331.
14
15 19. U.S. Food and Drug Administration. Interacting with the FDA on Complex Innovative
16 Trial Designs for drugs and biological products: Draft guidance for industry,
17 [https://www.fda.gov/regulatory-information/search-fda-guidance-](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/interacting-fda-complex-innovative-trial-designs-drugs-and-biological-products)
18 [documents/interacting-fda-complex-innovative-trial-designs-drugs-and-biological-](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/interacting-fda-complex-innovative-trial-designs-drugs-and-biological-products)
19 [products](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/interacting-fda-complex-innovative-trial-designs-drugs-and-biological-products) (2019, accessed 30 June 2020).
20
21
22
23
24
25
26 20. Chen C, Deng Q, He L, et al. How many tumor indications should be initially screened
27 in clinical development of next generation immunotherapies? *Contemporary Clinical*
28 *Trials* 2017; 59: 113-117.
29
30
31
32 21. Antonijevic Z. The impact of adaptive design on portfolio optimization. *Therapeutic*
33 *Innovation and Regulatory Science* 2016; 5: 615-619.
34
35
36 22. Prasad V and Mailankody S. Research and Development Spending to Bring a Single
37 Cancer Drug to Market and Revenues After Approval. *JAMA Intern Med* 2017;
38 177(11):1569-1575.
39
40
41
42
43 23. Guinn D, Ren Y, Wilhelm EE, et al. Utilizing real-world data to inform a confirmatory
44 basket trial design- studying use of rituximab in autoimmune diseases. *Value in Health*
45 2018; 21: S229-30.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3
4
5
6
7
8
9 24. Guinn D, Madhavan S and Beckman RA. *Harnessing Real-World Data to Inform*
10 *Platform Trial Design*. Platform Trial Designs in Drug Development: Umbrella Trials
11 and Basket Trials, 2018, p.55.
12
13
14
15 25. European Medicines Agency. Guideline on the investigation of subgroups in
16 confirmatory clinical trials, [https://www.ema.europa.eu/en/documents/scientific-](https://www.ema.europa.eu/en/documents/scientific-guideline/guideline-investigation-subgroups-confirmatory-clinical-trials_en.pdf)
17 [guideline/guideline-investigation-subgroups-confirmatory-clinical-trials_en.pdf](https://www.ema.europa.eu/en/documents/scientific-guideline/guideline-investigation-subgroups-confirmatory-clinical-trials_en.pdf) (2019,
18 accessed 30 June 2020).
19
20
21
22
23 26. ICH Expert Working Group. ICH Harmonised Tripartite Guideline: Statistical Principles
24 for Clinical Trials E9. *International Conference on Harmonisation (ICH)*, 1998.
25
26
27 27. Zhou H, Chen C, Sun L, et al. Bayesian optimal phase II clinical trial design with time
28 to event endpoint. *Pharmaceutical Statistics*. Epub ahead of print 10 June 2020.
29 <https://doi.org/10.1002/pst.2030>.
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

ACKNOWLEDGEMENTS

We thank Drs. Zoran Antonijevic, Lisa LaVange, Peter Mesenbrink, Martin Posch, William Rosenberger, and Alex Sverdlov for careful and insightful review of the manuscript.

DISCLAIMER

Dr. Daphne Guinn contributed to this article in her personal capacity. The views expressed are her own and do not necessarily represent the views of the Food & Drug Administration or the United States Government.

SUPPLEMENTAL METHODS

The adjusted nominal level α^* in the pooled analysis^{13, 14}

The pooled analysis aims to examine that there is treatment effect in at least one tumor indication. This analysis combined with the previous pruning and sample size readjustment steps is controlled at type I error of 0.025 for the global null hypothesis that all indications are inactive using methods from reference 13. Note that the definition of Type I error of the entire trial in reference 13 is different from the FWER in this study. Specifically, in reference 13, let Y_{i1} be the standardized test statistics based on the endpoint used for pruning at the interim analysis, and Y_{i2} be the standardized test statistics based on the endpoint for pooling for the i -th tumor indication at the final analysis ($i = 1, \dots, k$). Suppose that m tumor indications are included in the pooled analysis ($m \geq 1$). Let V_m be the corresponding standardized test statistics pooled from Y_{i2} , which can be written as $(\sum_{i=1}^m Y_{i2})/\sqrt{m}$.

Under the null hypothesis H_0 that there is no treatment effect in any of the tumor indications, the probability of incorrectly declaring activity in a basket, denoted as α , is

$$\alpha = \sum_{m=1}^k c(k,m) Q_0(\alpha^* | \alpha_t, m)$$

where $c(k,m) = k! / ((k-m)!m!)$ is the number of choices for selection of m tumor indications from k , and $Q_0(\alpha^* | \alpha_t, m)$ is the probability of V_m being statistically significant at the adjusted α^* level given m out of k indications in the pool, formulated as

$$Q_0(\alpha^* | \alpha_t, m) = P_{H_0}(\cap \{Y_{i1} > Z_{1-\alpha_t}; i = 1, \dots, m\}, \cap \{Y_{i1} < Z_{1-\alpha_t}; i = m+1, \dots, k\}, V_m > Z_{1-\alpha^*})$$

Assuming that $\{Y_{i1}; i = 1, \dots, k\}$ are *i.i.d.*, under the global null hypothesis we have

$$Q_0(\alpha^* | \alpha_t, m) = P_{H_0}(\cap \{Y_{i1} > Z_{1-\alpha_t}; i = 1, \dots, m\}, V_m > Z_{1-\alpha^*}) (1 - \alpha_t)^{(k-m)}$$

Consider $\{Y_{11}, \dots, Y_{m1}, V_m\}$ following a multivariate normal distribution $(\mathbf{0}, \Sigma)$, where

$$\Sigma = \begin{pmatrix} 1 & & 0 & \text{corr}(Y_{11}, V_m) \\ & \ddots & & \vdots \\ 0 & & 1 & \text{corr}(Y_{m1}, V_m) \\ \text{corr}(Y_{11}, V_m) & \dots & \text{corr}(Y_{m1}, V_m) & 1 \end{pmatrix}. \text{ Setting } \alpha = 0.025, \text{ the}$$

adjusted level α^* , at which the pooled analysis is nominally set, can be solved based on the correlation between Y_{i1} and V_m , $\text{corr}(Y_{i1}, V_m)^{13}$.

1
2
3
4
5
6
7
8
9 In this study, we consider the following three sample size adjustment strategies¹³
10 and the corresponding $\text{corr}(Y_{i1}, V_m)$ can be calculated as:

- 14 1. Design one (D1): The sample size for each tumor indication is fixed
15 upfront at n as planned. After pruning, the sample size will be less than or
16 equal to the originally planned kn . Under D1, $\text{corr}(Y_{i1}, V_m) = \sqrt{t/m}$. This
17 strategy corresponds to no sample size adjustment.
18
- 19 2. Design two (D2): The sample size for each tumor indication will increase
20 after the interim analysis so that the total sample size in the pooled analysis
21 remains kn , which is greater than the sample size under D1 after pruning.
22 Under D2, $\text{corr}(Y_{i1}, V_m) = \sqrt{t/k}$. This is the most aggressive sample size
23 adjustment strategy.
24
- 25 3. Design three (D3): The sample size for each tumor indication will
26 increase after the interim analysis so that the total sample size in the overall
27 study remains kn . Thus, the total sample size in the pooled analysis is
28 greater than the sample size under D1 but smaller than the sample size
29 under D2. Under D3, $\text{corr}(Y_{i1}, V_m) = \sqrt{t/(mt + k(1 - t))}$.
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

SUPPLEMENTAL TABLES

Table S1. Glossary of terms for simulation study

Design parameters	Value(s) in simulation study	Descriptions
The state of nature		
g	$0, \dots, k$	The number of active indications at the beginning
θ_i	At null value $\theta_0 = 1$, At active value $\theta_1 = 0.5, 0.6, 0.7, 0.8$.	The hazard ratio of experimental arm vs. control arm for the definitive endpoint of each indication
Study design input parameters		
k	3,4,5,6	The number of tumor indications at the beginning
D	D1, D2, D3	Sample size adjustment strategies
t	0.5	A common information time for the interim analysis
α_t	0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4	A common bar for the interim analysis to prune inactive indications
α	0.025	The false positive rate for pooled analysis together with the pruning step with respect to the global null hypothesis, after inflation from the nominal value α^*
β	0.025, 0.05, 0.1, 0.2	False negative rate for the pooled analysis

α_{post}	0.05,0.1,0.15,0.2,0.25,0.3,0.35,0.4	A common bar for the post-individual test, which is varied independently of α_t
α_{ref}	0.025	False positive rate for the reference design
$\beta_{uncorrected}$	0.1	False negative rate for the uncorrected reference design
S	10000	The number of simulated replications
Calculated parameters		
n	$n = \frac{4(Z_{1-\alpha} + Z_{1-\beta})^2}{k(\log\theta)^2}$	Planned sample size for each tumor indication
α^*	Calculated by numerically solving the equation $\sum_{i=1}^k c(k,m)Q_0(\alpha^*, \alpha_t, m) = \alpha$	The adjusted nominal level α^* , at which the pooled analysis is nominally set to control the false positive rate α for the pooled analysis combined with the pruning step.
$\beta_{corrected}$	1 – power by indication corresponding to the given input parameters as determined by simulation	False negative rate for the corrected reference design
Simulated parameters	Descriptions	
m	The number of tumor indications included in the pooled analysis	
n_{i1}	Sample size for each tumor indication at the interim analysis (nt)	
n_{i2}	Sample size for each tumor indication at the pooled analysis	
Y_{i1}	The standardized test statistics for the i -th tumor indication at the interim analysis	
Y_{i2}	The standardized test statistics for the i -th tumor indication at the final analysis	
V_m	The standardized test statistics pooled from Y_{i2}	
d	The number of false positive indications passing the final individual tests	
j	The number of true positive indications passing the final individual tests	

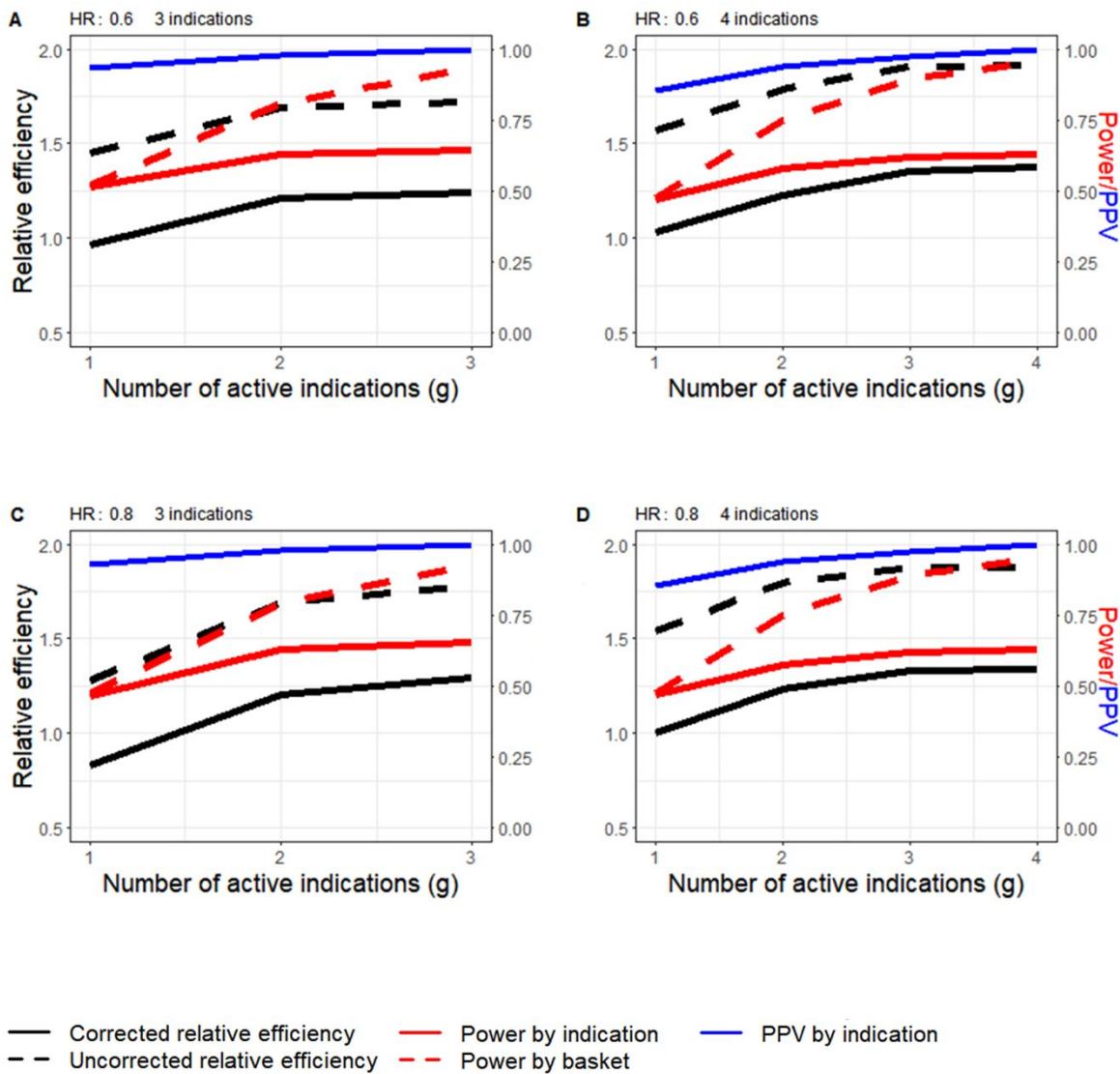
Outcome measurements	Estimates in simulation study	Descriptions
$\alpha_{net}(g)$	$\frac{1}{S} \sum_{s=1}^S I(V_m^{(s)} > Z_{1-\alpha^{(s)}}) I(d^{(s)} > 0)$	The probability of the basket trial passing the pooled test and at least one false positive indication passing the post-individual test for a given value of g .
Family-wise error rate (FWER)	$\max\{\alpha_{net}(g): g = 0, \dots, k\}$	For any g ($g = 0, \dots, k$), the maximum probability of the basket trial passing the pooled test and at least one false positive indication passing the post-individual test
Power by indication	$\frac{1}{S} \sum_{s=1}^S \frac{j_{new}^{(s)}}{g} I(V_m^{(s)} > Z_{1-\alpha^{(s)}})$	The proportion of true positive indications that pass the post-individual test (requires that the pooled test passed)
Power by basket	$\frac{1}{S} \sum_{s=1}^S [I(g > 0) I(V_m^{(s)} > Z_{1-\alpha^{(s)}})]$	The probability of an active basket (one that contains at least one active indication) passing the pooled analysis
Sample size	$(k - m)n_{i1} + mn_{i2}$	The sample size of a basket trial
Efficiency	$\frac{\text{power} \times g}{\frac{1}{S} \sum_{s=1}^S \left(\sum_{i=1}^m n_{i2}^{(s)} + \sum_{i=m^{(s)}+1}^k n_{i1}^{(s)} \right)} = \frac{\sum_{s=1}^S j_{new}^{(s)} I(V_m^{(s)} > Z_{1-\alpha^{(s)}})}{\sum_{s=1}^S \left(\sum_{i=1}^m n_{i2}^{(s)} + \sum_{i=m^{(s)}+1}^k n_{i1}^{(s)} \right)}$	The ratio of average number of active indications that passed the post-individual tests divided by the average sample size
Uncorrected reference efficiency	$\frac{g(1 - \beta_{uncorrected})(\log \theta_1)^2}{4k(Z_{1-\alpha_{ref}} - Z_{1-\beta_{uncorrected}})^2}$	The ratio of estimated number of active indications divided by the pre-defined total sample

		size in the reference study powered at 90%.
Uncorrected relative efficiency	Efficiency/Uncorrected reference efficiency	The ratio of efficiency and uncorrected reference efficiency
Corrected reference efficiency	$\frac{g(1 - \beta_{corrected})(\log \theta_1)^2}{4k(Z_{1-\alpha_{ref}} - Z_{1-\beta_{corrected}})^2}$	The ratio of estimated number of true positive indications divided by the pre-defined total sample size in the reference study at the same power by indication observed in the simulation in the basket trial with corresponding parameters, investigating the same indications.
Negative predictive value (NPV) by indications	$\frac{1}{S} \sum_{s=1}^S \frac{k - g - d^{(s)}}{k - j^{(s)} - d^{(s)}}$	The proportion of true negative indications among all negative indications not passing the interim test, the pooled test, or the post individual test
Negative predictive value (NPV) by basket	$\frac{1}{S} \sum_{s=1}^S I(g = 0)I(j^{(s)} + d^{(s)} = 0)$	The proportion of true negative baskets among all negative baskets without any indication passing the interim test, the pooled test, or the post individual test
Positive predictive value (PPV) by indications	$\frac{1}{S} \sum_{s=1}^S \frac{j^{(s)}}{j^{(s)} + d^{(s)}}$	The proportion of true positive indications among all positive indications passing the interim test, the pooled test, and the post individual test
Positive predictive	$\frac{1}{S} \sum_{s=1}^S I(g > 0)I(j^{(s)} + d^{(s)} > 0)$	The proportion of true positive baskets among all positive baskets with at

value (PPV) by basket		least one indication passing the interim test, the pooled test, and the post individual test
Corrected relative efficiency	Efficiency/Corrected reference efficiency	The ratio of efficiency and corrected reference efficiency
Coverage for hazard ratio (HR)	$\frac{1}{S} \sum_{s=1}^S [\sum_{i=1}^{m^{(s)}} I(Y_{i2}^{(s)} > Z_{1-\alpha_{post}}) I(Y_{i2}^{(s)} + \log(HR) > Z_{1-\alpha_{post}})]$	The probability that the estimated 95% CI for HR covers the true HR given the individual test passed
Bias of estimated HR	$\frac{1}{S} \sum_{s=1}^S \left\{ \frac{1}{m^{(s)}} \sum_{i=1}^{m^{(s)}} \exp\left(-Y_{i2}^{(s)} \sqrt{\frac{4}{n_{i2}^{(s)}}}\right) I(Y_{i2}^{(s)} > Z_{1-\alpha_{post}}) - HR \right\}$	The relative difference in the true HR and the estimated HR, defined as the ratio of estimated average HR and true pooled HR for those indications that pass the individual tests minus 1.

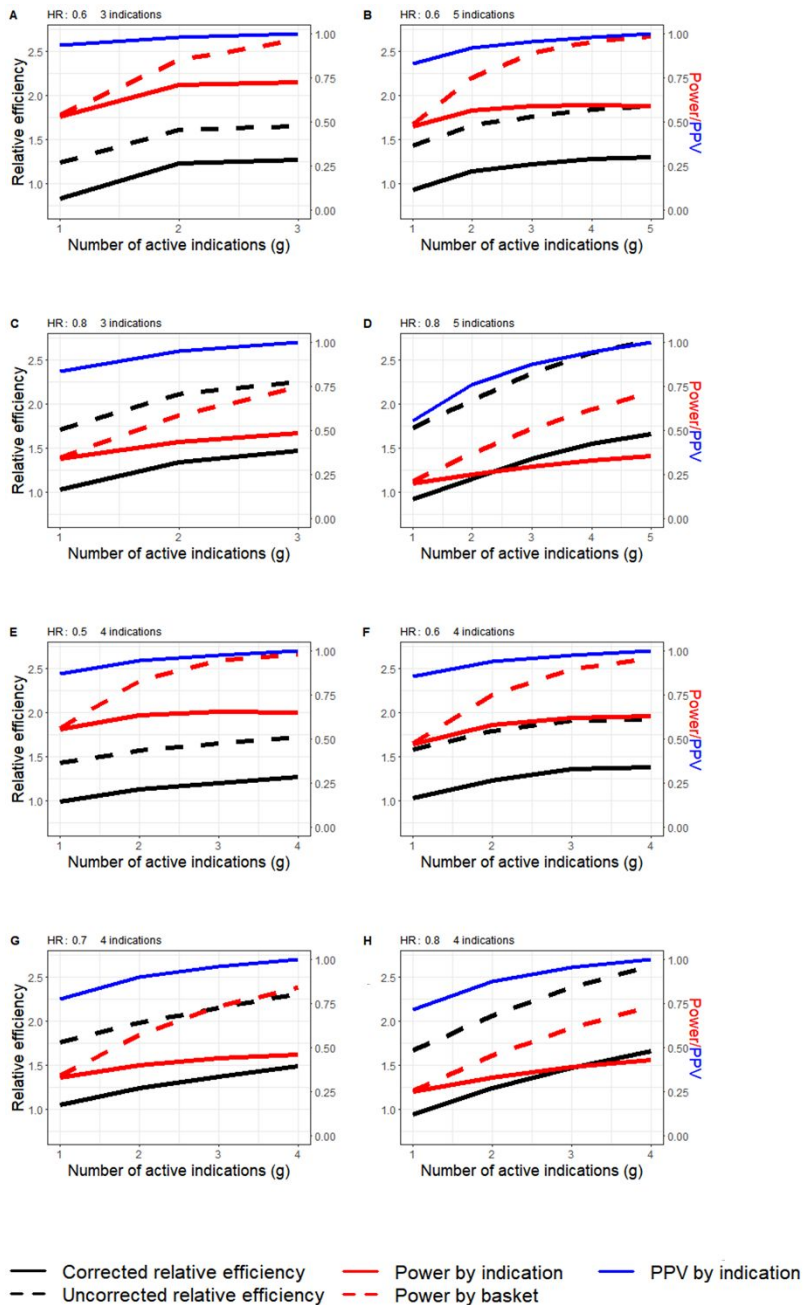
Table S2. Simulation results. Each row summarizes the results of 10000 simulations for a given scenario with input parameters: hazard ratio (HR), number of indications (k), β , α_t (α_t), α_{post} (α_{post}), number of active indications (g), and sample size adjustment strategies. The outcome measurements include: α_{net} (α_{net}), power by indication, power by basket, mean sample size over simulations, efficiency, 95% CI coverage for HR, bias, uncorrected reference efficiency, corrected reference efficiency, uncorrected relative efficiency, and corrected relative efficiency.

SUPPLEMENTAL FIGURES



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure S1. Recommended development approaches for (A) 3 indications with $HR = 0.6$ (Design 3, $kn = 199$, $\beta = 0.05$, $\alpha_t = 0.3$, $\alpha_{post} = 0.05$, $\alpha^* = 0.009$), (B) 4 indications with $HR = 0.6$ (Design 3, $kn = 236$, $\beta = 0.025$, $\alpha_t = 0.2$, $\alpha_{post} = 0.1$, $\alpha^* = 0.0075$), (C) 3 indications with $HR = 0.8$ (Design 3, $kn = 1044$, $\beta = 0.05$, $\alpha_t = 0.4$, $\alpha_{post} = 0.05$, $\alpha^* = 0.009$), and (D) 4 indications with $HR = 0.8$ (Design 3, $kn = 1234$, $\beta = 0.025$, $\alpha_t = 0.2$, $\alpha_{post} = 0.01$, $\alpha^* = 0.0075$). The x-axis represents the number of active indications (indications in which the drug provides clinical benefit), the primary y-axis (left) represents the uncorrected/corrected relative efficiency, and the second y-axis (right) represents the power (red) by indication and by basket, and the positive predictive value by indication (blue).



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure S2. Cases with maximum corrected relative efficiency for (A) 3 indications with $HR = 0.6$ (Design 3, $kn = 236$, $\beta = 0.025$, $\alpha_t = 0.4$, $\alpha_{post} = 0.05$, $\alpha^* = 0.009$), (B) 5 indications with $HR = 0.6$ (Design 2, $kn = 236$, $\beta = 0.025$, $\alpha_t = 0.2$, $\alpha_{post} = 0.1$, $\alpha^* = 0.0078$), (C) 3 indications with $HR = 0.8$ (Design 3, $kn = 631$, $\beta = 0.2$, $\alpha_t = 0.2$, $\alpha_{post} = 0.15$, $\alpha^* = 0.0101$), (D) 5 indications with $HR = 0.8$ (Design 3, $kn = 631$, $\beta = 0.2$, $\alpha_t = 0.15$, $\alpha_{post} = 0.4$, $\alpha^* = 0.0071$), (E) 4 indications with $HR = 0.5$ (Design 2, $kn = 128$, $\beta = 0.025$, $\alpha_t = 0.2$, $\alpha_{post} = 0.1$, $\alpha^* = 0.0094$), (F) 4 indications with $HR = 0.6$ (Design 3, $kn = 236$, $\beta = 0.025$, $\alpha_t = 0.2$, $\alpha_{post} = 0.1$, $\alpha^* = 0.0075$), (G) 4 indications with $HR = 0.7$ (Design 2, $kn = 247$, $\beta = 0.2$, $\alpha_t = 0.2$, $\alpha_{post} = 0.15$, $\alpha^* = 0.0094$), and (H) 4 indications with $HR = 0.8$ (Design 3, $kn = 631$, $\beta = 0.2$, $\alpha_t = 0.2$, $\alpha_{post} = 0.2$, $\alpha^* = 0.0075$). The x-axis represents the number of active indications (indications in which the drug provides clinical benefit), the primary y-axis (left) represents the uncorrected/corrected relative efficiency, and the second y-axis (right) represents the power (red) by indication and by basket, and the positive predictive value by indication (blue).

SUPPLEMENTAL R CODES

```
#Calculate alpha* given alpha_t, information time, and the number of indications.  
#Reference: Cong Chen, Xiaoyun (Nicole) Li, Shuai Yuan, Zoran Antonijevic,  
#Rasika Kalamegham & Robert A. Beckman(2016) Paper: Statistical Design and  
#Considerations of a Phase 3 Basket Trial for Simultaneous Investigation of  
#Multiple Tumor Types in One Study, Statistics in Biopharmaceutical Research,  
#8:3, 248-257 DOI: 10.1080/19466315.2016.1193044. Function "mf2" is the  
original  
#function from Chen's online code, it is used to express formula (3) in  
#Chen-Beckman paper.
```

```
mf2 <- function(alphastar, alphas, t, m, k, d, Rho_for_endpoints=1) {  
  # test1, test2, and test3 denotes correlation matrices for D1, D2, and D3.  
  test <- matrix(0, ncol <- m + 1, nrow <- m + 1) # D1  
  low <- rep(qnorm(1 - alphas), m + 1)  
  low[m + 1] <- qnorm(1 - alphastar)  
  up <- rep(Inf, m + 1)  
  diag(test) <- 1  
  
  for (i in 1:m){  
    test[i, m + 1] <- test[m + 1, i] <- switch (d,Rho_for_endpoints * sqrt(t / m),
```

```

1
2
3
4
5
6
7
8
9           Rho_for_endpoints * sqrt(t / k),
10          Rho_for_endpoints * sqrt(t / (k * (1 - t) + m * t)))
11
12      }
13
14
15
16
17      # joint_probability1, joint_probability2, and joint_probability3 denotes the
18      # joint probability in equation (3) in Chen-Beckman paper for D1, D2, and D3.
19      joint_probability <- pmvnorm(lower <- low, upper <- up, mean <- rep(0, m + 1), corr
20      <- test)[1]
21
22
23
24
25      return(joint_probability)
26
27  }
28
29
30
31  # Function "type2" is the original function from Chen's online code, it is used
32  # to calculate alpha* by equation (3) in Chen-Beckman paper.
33
34  type2 <- function(alphastar, alphas, t, k, d, Rho_for_endpoints = .5) {
35
36      # joint_probability denotes the joint probability in equation(3) in Chen-Beckman
37      # paper.
38
39      joint_probability <- 0
40
41
42
43
44      for (i in 1:k) joint_probability <- joint_probability + factorial(k) / (factorial(i) *
45      factorial(k - i)) *
46
47      (1-alphas)^(k-i)*mf2(alphastar=alphastar, alphas=alphas, t=t, m=i, k=k,
48      d=d,Rho_for_endpoints=Rho_for_endpoints)
49
50
51
52
53
54
55
56
57
58
59
60

```

```
1
2
3
4
5
6
7
8
9
10
11     return(joint_probability - 0.025)
12
13 }
14
15
16
17 # "Simulation" function is denoted to calculate Type I error and power
18 # D denotes D1, D2, and D3, alpha_t denotes Type I error in interim stage,
19 # alpha_tt denotes Type I error after final stage for the post-trial test.
20
21 Simulation=function(alpha_t, alpha_tt, g, k, t = 0.5, design=c(D1=1,D2=2,D3=3),
22                    Rho_for_endpoints,hr, delta=-log(hr), n, simulation_times=10000) {
23
24
25
26
27
28     #set seed before simulation is to guarantee the results are reproducible.
29
30     set.seed(123)
31
32     # print(dummy_indication)
33
34     #t denotes information time
35
36     dummy_indication= sample(c(rep(1, g), rep(0, k-g)))
37
38     delta=delta * dummy_indication; hr = exp(-delta)
39
40     # print(delta);print(hr)
41
42     #alphastar denotes alpha* under D1, D2 and D3 by Chen-Beckman's formula (2),
43     the
44
45     #function of "uniroot" is from Chen's online code.
46
47     alphastar <- lapply(design,function(dd)
48
49         uniroot(type2, c(0, 1), alphas=alpha_t, t=t, k=k,
50                 d=dd,Rho_for_endpoints=Rho_for_endpoints)$root)
51
52
53
54
55
56
57
58
59
60
```

```

1
2
3
4
5
6
7
8
9
10
11 mean_n <- delta * sqrt(n * t / 4)
12
13 ##### Simulate result in before pruning
14
15 x1 <- t(sapply(1:k,function(ln) if(dummy_indication[ln]==1)
16 {rnorm(simulation_times,mean=mean_n[ln], sd = 1)
17
18 }else {rnorm(simulation_times, mean = 0, sd = 1)}))
19
20 passed_pruning <- (x1 > qnorm(1 - alpha_t))
21
22
23
24 #m denotes the number of remained indications after pruning
25
26 m=colSums(passed_pruning)
27
28
29
30 # Calculate sample size after pruning based on adjustment strategies.
31
32 sample_size <- lapply(design,function(dd) apply(passed_pruning,2,
33 function(pptmp)
34
35   switch (dd,n * pptmp, ceiling(k * n / ifelse(sum(pptmp)>0,sum(pptmp),Inf)) *
36   pptmp,
37
38     ceiling((n * t + k * n * (1 - t) / ifelse(sum(pptmp)>0,sum(pptmp),Inf))) *
39   pptmp)))
40
41
42
43 ## calculate total sample size in the trial
44
45 total_sample_size <- lapply(design,function(dd) apply(passed_pruning, 2,
46 function(l)
47
48   switch (dd,
49
50     sum(n*l)+sum(n*t*(1-l)),
51
52
53
54
55
56
57
58
59
60

```

```
1
2
3
4
5
6
7
8
9         ifelse(sum(l)>0,sum(n)+sum(n*t*(1-l)),sum(n*t*(1-l))),
10        ifelse(sum(l)>0,sum(n),sum(n*t*(1-l)))
11
12    )))
13
14
15
16
17
18
19    # Calculate correlation between Yi1 and Yi2 if not all indications are pruned after
20    pruning step.
21
22    rho_between_standardized_test_statistics <- lapply(design,function(dd)
23    sapply(m, function(mi)
24        ifelse(mi>0, switch(dd,
25            Rho_for_endpoints * sqrt(t),
26            Rho_for_endpoints * sqrt(t * mi / k),
27            Rho_for_endpoints * sqrt(t / (t + k * (1 - t) / mi))
28        ), 0)))
29
30
31
32
33
34
35
36    # Generate Yi2 based on Yi1, corr(Yi1, Yi2), and adjusted sample size.
37
38    #mean in the interim stage, mu1=sqrt(n*t/4)*delta*dummy_indication
39
40    #representing means in the interim stage
41
42    mu1 <- sqrt(n * t / 4) * delta * dummy_indication
43
44
45
46    #mean in the final stage, mu2=sqrt(sample_size*1/4)*delta*dummy_indication
47
48    #representing means in the final stage
49
50
51
52
53
54
55
56
57
58
59
60
```

```

1
2
3
4
5
6
7
8
9     mu2 <- lapply(design, function(dd) (sqrt(sample_size[[dd]] / 4) * delta *
10 dummy_indication))
11
12
13
14     #To generate Yi2, variance of Yi2 is sqrt(((1-corr(Yi1,Yi2)^2))*s2^2,s2=1 in
15     #our case(standardized normal distribution), given one of D1, D2, or D3.
16
17     sd2 <- lapply(design, function(dd) sqrt((1 -
18 rho_between_standardized_test_statistics[[dd]] ^ 2)))
19
20
21
22
23     #To generate Yi2, mean of Yi2(denoted as
24     #mean_x2)=mu2+(s2/s1)*corr(Yi1,Yi2)*(Yi1-mu1),s1=s2=1 in our
25     #case(standardized normal distribution) given one of D1, D2, or D3.
26
27     mean_x2 <- lapply(design, function(dd) mu2[[dd]] + (rep(1,k) %o%
28 rho_between_standardized_test_statistics[[dd]]) *
29     (x1 - mu1))
30
31
32
33
34
35
36     #Generate Yi2 based on mean and variance for 3 indications, given one of D1,
37     #D2, or D3.
38
39     x2 <- lapply(design, function(dd) sapply(1:simulation_times, function(nsim)
40     rnorm(mean_x2[[dd]][,nsim], mean_x2[[dd]][,nsim], sd = sd2[[dd]][nsim])))
41
42
43
44
45     #Calculate V_m statistics, denotes as the sum of Yi2 for those indications
46     #passed pruning at interim stage, divided by square root of number of
47     #indications remained after pruning
48
49
50
51
52
53
54
55
56
57
58
59
60

```

```
1
2
3
4
5
6
7
8
9 vm <- lapply(design, function(dd) sapply(1:simulation_times, function(nsim)
10   ifelse(m[nsim]>0,sum(x2[[dd]][passed_pruning[,nsim],nsim]) / sqrt(m[nsim]),0)))
11
12
13
14
15 # p_value_for_final_stage_testing denotes p-values for each simulation
16
17 #time to compare with alphastar later
18
19   p_value_for_final_stage_testing <- lapply(vm, pnorm)
20
21
22
23 #If pooled indications are positive, do a post-check on individual
24 #indications at alpha=post_alpha_t, and discard indications that do not
25 #achieve statistical significance.
26
27 #passed_pruning_post_trial denotes indications that passed pruning and pass
28 #post-trial test
29
30 passed_pruning_post_trial <-lapply(design, function(dd)
31   (x1 > qnorm(1 - alpha_t) & x2[[dd]] > qnorm(1 - alpha_tt)))
32
33
34
35
36
37
38 # calculate the coverage rate of 95% CI for individual indications and
39 # pooled indications.
40
41 # coverage of individual approved indications
42
43 coveragebias=lapply(design, function(dd)
44 sapply(1:simulation_times,function(nsim) {
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
```

```

1
2
3
4
5
6
7
8
9
10 pppttmp=passed_pruning_post_trial[[dd]][,nsim];wttmp=n[pppttmp]/sum(n[pppttmp]
11 )
12
13   pppttmp.tp=(passed_pruning_post_trial[[dd]][,nsim]) & (dummy_indication==1)
14   wttmp.tp=n[pppttmp]/sum(n[pppttmp])
15
16
17
18
19   hrtmp=exp(- x2[[dd]][,nsim] * sqrt(4 / sample_size[[dd]][,nsim]))[pppttmp]
20   hrtmp.tp=exp(- x2[[dd]][,nsim] * sqrt(4 / sample_size[[dd]][,nsim]))[pppttmp.tp]
21
22
23
24
25   c(individual=ifelse(sum(pppttmp)>0,
26       sum(abs(mu2[[dd]][pppttmp,nsim]-x2[[dd]][pppttmp,nsim])<
27         qnorm(p = 0.025, lower.tail = F))/sum(pppttmp),NA),
28
29
30
31
32   pooled=ifelse(sum(pppttmp)>0,
33       abs(sum(x2[[dd]][pppttmp,nsim]) / sqrt(sum(pppttmp))+
34
35
36
37   log(sum(hr[pppttmp]*wttmp))*sqrt(sum(sample_size[[dd]][pppttmp,nsim]) / 4))<
38       qnorm(p = 0.025, lower.tail = F),NA),
39
40
41
42
43   bias1=ifelse(sum(pppttmp)>0,mean(hrtmp)/ mean(hr[pppttmp])-1,NA),
44   bias2=ifelse(sum(pppttmp)>0,sum(hrtmp * wttmp) / sum(hr[pppttmp] * wttmp)-
45     1,NA),
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

```



```

1
2
3
4
5
6
7
8
9     indiv.tp=ifelse(sum(pppttmp.tp)>0,
10                   sum(abs(mu2[[dd]][[pppttmp.tp,nsim]]-x2[[dd]][[pppttmp.tp,nsim]])<
11                       qnorm(p = 0.025, lower.tail = F))/sum(pppttmp.tp),NA),
12                   qnorm(p = 0.025, lower.tail = F))/sum(pppttmp.tp),NA),
13
14     pooled.tp=ifelse(sum(pppttmp.tp)>0,
15                      abs(sum(x2[[dd]][[pppttmp.tp,nsim]]) / sqrt(sum(pppttmp.tp))+
16                          log(sum(hr[pppttmp.tp]*wttmp.tp))*sqrt(sum(sample_size[[dd]][[pppttmp.tp,nsim]]) /
17                              4))<
18                          qnorm(p = 0.025, lower.tail = F),NA),
19
20                      bias1.tp=ifelse(sum(pppttmp.tp)>0, mean(hrtmp.tp)/ mean(hr[pppttmp.tp])-
21                                  1,NA),
22
23                      bias2.tp=ifelse(sum(pppttmp.tp)>0,sum(hrtmp.tp * wttmp.tp) /
24                                  sum(hr[pppttmp.tp] * wttmp.tp)-1,NA)
25
26                      )))
27
28
29
30
31
32
33
34
35     #Record tp, fp in each simulation time after interim stage
36
37     #tp denotes the number of active remained indication after pruning, fp denotes
38     the number of remained
39
40     #indication after pruning.
41
42
43
44     fp= lapply(design,function(dd) colSums(dummy_indication==0 &
45     passed_pruning_post_trial[[dd]]==1))
46
47
48
49     #j, denotes the number of active indications remained after pruning and pass
50
51
52
53
54
55
56
57
58
59
60

```

```
1
2
3
4
5
6
7
8
9 #post-trial test
10
11 #(if we don't do post trial test, we change passed_pruning_post_trial to
12 #passed_pruning in this line.)
13
14 tp=lapply(design, function(dd) colSums(dummy_indication==1 &
15 passed_pruning_post_trial[[dd]]==1))
16
17
18
19
20 # Use the formula of Type I error and powers to get simulation results.
21
22
23
24 final_pooled_test=lapply(design, function(dd)
25   m > 0 & p_value_for_final_stage_testing[[dd]] > (1 - alphastar[[dd]]))
26
27
28
29 ftp <-lapply(design, function(dd)
30   c(type_I_error=sum(final_pooled_test[[dd]] & fp[[dd]] > 0) / simulation_times,
31     power1=ifelse(g>0,sum(tp[[dd]][final_pooled_test[[dd]]) / (g *
32 simulation_times),0),
33     power2=ifelse(g>0,sum(final_pooled_test[[dd])/simulation_times,0)))
34
35
36
37
38
39
40 ## Average total sample size
41
42 average_total_sample_size <- lapply(total_sample_size, mean)
43
44
45
46 ## Efficiency
47
48 efficiency <- lapply(design, function(dd) c(efficiency=g * ftp[[dd]]['power1'] /
49 average_total_sample_size[[dd]]))
50
51
52
53
54
55
56
57
58
59
60
```

```
1
2
3
4
5
6
7
8
9   inverse_efficiency <- lapply(design, function(dd) 1/efficiency[[dd]])
10
11
12
13   ## Coverage and bias
14
15   coveragebiasmean <-lapply(coveragebias,rowMeans,na.rm=TRUE)
16
17   cbmse <-lapply(coveragebias,function(cb) apply(cb, 1, function(cbtmp)
18     mean((cbtmp-mean(cbtmp,na.rm = TRUE))^2,na.rm = TRUE)))
19
20
21
22   output <- list(test=fptp, mean_samplesize=average_total_sample_size,
23     efficiency=efficiency,cbmse=cbmse,
24     mean_coveragebias=coveragebiasmean,coveragebias=coveragebias)
25
26
27
28
29   return(output)
30
31 }
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
```