

Health Utility Adjusted Survival: a Composite Endpoint for Clinical Trial Designs

Yangqing Deng¹, PhD, John R. de Almeida^{2,3}, MD, MSc, FRCSC, Wei Xu^{1,4*}, PhD

¹Department of Biostatistics, University Health Network, Toronto, ON, Canada

²Department of Otolaryngology—H&N Surgery, University Health Network, Toronto, ON, Canada

³Institute of Health Policy, Management and Evaluation, University of Toronto, Toronto, ON, Canada

⁴Dalla Lana School of Public Health, University of Toronto, Toronto, ON, Canada

*Corresponding author: Wei Xu, Email: wei.xu@uhnresearch.ca

Address: 10-511, 610 University Ave, Toronto, M5G 2M9

Tel: (416)946-4501

ABSTRACT

Many randomized trials have used overall survival as the primary endpoint for establishing non-inferiority of one treatment compared to another. However, if a treatment is non-inferior to another treatment in terms of overall survival, clinicians may be interested in further exploring which treatment results in better health utility scores for patients. Examining health utility in a secondary analysis is feasible, however, since health utility is not the primary endpoint, it is usually not considered in the sample size calculation, hence the power to detect a difference of health utility is not guaranteed. Furthermore, often the premise of non-inferiority trials is to test the assumption that an intervention provides superior quality of life or toxicity profile without compromising the survival when compared to the existing standard. Based on this consideration, it may be beneficial to consider both survival and utility when designing a trial. There have been methods that can combine survival and quality of life into a single measure, but they either have

1 strong restrictions or lack theoretical frameworks. In this manuscript, we propose a method
2 called HUS (Health Utility adjusted Survival), which can combine survival outcome and
3 longitudinal utility measures for treatment comparison. We propose an innovative statistical
4 framework as well as procedures to conduct power analysis and sample size calculation. By
5 comprehensive simulation studies involving summary statistics from the PET-NECK trial,¹ we
6 demonstrate that our new approach can achieve superior power performance using relatively
7 small sample sizes, and our composite endpoint can be considered as an alternative to overall
8 survival in future clinical trial design and analysis where both survival and health utility are of
9 interest.

10

11 *Keywords:* Health utility; Overall survival; Time-to-event data; Hazard ratio; Proportional

12 hazards; Randomized controlled trials

1 1. INTRODUCTION

2 In many clinical studies, overall survival (OS) is used as the primary endpoint to assess efficacy
3 of treatments. Superiority trials are used to test whether a new treatment is better than a standard
4 or control treatment, while non-inferiority trials are used to test whether the new treatment is not
5 unacceptably worse than control. Non-inferiority trials are especially important in circumstances
6 where the new treatment may have other benefits (e.g., lower costs, fewer side effects, improved
7 quality of life, or is easier to implement) compared to control, and people are only interested in
8 showing the new treatment is not worse than control in terms of OS. When non-inferiority has
9 been established, clinicians may be interested in further examining whether the new treatment
10 can benefit patients more in terms of health utility.¹ Health utility is a construct, usually ranging
11 from 0 to 1 (although theoretically can also have negative values), that quantifies the preference
12 for a given health state experienced by a patient at a certain time point. A higher value means a
13 healthier state, while death usually corresponds to 0. Using health utility scores at different time
14 point during the treatment and post-treatment, statistical analysis may be performed to compare
15 different treatment groups' utility scores.²⁻⁴ However, given that the study design is usually
16 based on the primary endpoint of OS without considering health utility, whether there will be
17 enough power for health utility analysis is uncertain. Also, conducting tests for OS and health
18 utility separately may not be the most efficient, because it involves multiple testing adjustment
19 and can lose statistical power. Hence, it may be beneficial to consider using a composite
20 endpoint that combines survival and utility, which may lead to increased statistical power and
21 smaller required sample sizes.

22

1 Creation of a composite endpoint of survival and utility, can also aid in clinical interpretation of
2 non-inferiority trials where non-inferiority of survival is not the only acceptable outcome. For
3 example, a new therapeutic intervention may be purported as offering improvements in quality of
4 life or toxicity. However, clinicians may not be willing to sacrifice disease control to provide
5 these other benefits. In this case, testing this new intervention in phase 3 non-inferiority trial
6 where overall survival is the primary outcome and quality of life or toxicity is a secondary
7 outcome may establish the intervention as non-inferior from a survival perspective and then
8 falsely identify the new intervention as a standard of care without appropriate consideration of
9 quality of life and toxicity. On the other hand, one may consider a situation where a patients'
10 preference for improved quality of life (or utility) may outweigh their desire to have non-inferior
11 survival. In this instance, demonstration of superiority of utility may not be enough if it is
12 associated with a significant loss of survival and the two outcomes cannot be interpreted in
13 isolation. In this instance, a combination of both survival and utility endpoints may be needed to
14 declare a new intervention superior.

15
16 Some methods that can combine survival and utility have been proposed and used to analyze
17 clinical trial data, and the most commonly used method is called Q-TWiST (Quality-adjusted
18 Time Without Symptoms of disease or Toxicity).⁵⁻¹² Though Q-TWiST has not been commonly
19 seen as a primary endpoint for designing new studies, researchers have derived its statistical
20 properties as well as formulas for sample size calculations.⁸ That being said, one major issue
21 about Q-TWiST is that it divides each patient's status into three states (toxicity, time without
22 symptoms and toxicity, and relapse) and uses pre-selected weights for different states. In many
23 scenarios, with utility scores measured as continuous variables at different time points

1 throughout the trials, it may be much more desirable to analyze them in their original scales
2 rather than forcing to have three categories, which may likely result in loss of information and
3 decreased statistical power.

4
5 QALY (Quality-Adjusted Life Years), of which Q-TWiST can be considered as a special case, is
6 the most intuitive way to combine survival with utility when comparing different treatments.^{5, 13-}

7 ¹⁷ It has also been used in the field of cost utility analysis, where similar methods have been
8 proposed and compared.¹⁸⁻²¹ Quality-adjusted progression-free survival, a similar concept with
9 slightly different focus, has also been used to assess the benefits of different treatments in
10 randomized trials.²²⁻²⁴ However, such measures have rarely been considered as a primary
11 endpoint for designing new trials, and we are not aware of any detailly developed statistical
12 framework or comprehensive simulation studies that demonstrate the advantages and feasibility
13 of a quality-adjusted survival endpoint compared to the traditional survival endpoint.

14
15 With these limitations and considerations, we propose an innovative composite endpoint for
16 combining longitudinal health utility and survival, called HUS (Health Utility adjusted Survival),
17 with a detailed statistical testing framework as well as procedures to perform power analysis and
18 sample size calculations. By assigning weights to health utility and survival, HUS can be
19 modified to suit different scenarios with increased power.

20
21 This new endpoint may help better interpret the findings in clinical trials. Often non-inferiority
22 trials are plagued with uncertainty of the efficacy of a new intervention that is statistically
23 deemed non-inferior based mainly on survival estimates but that has not been clearly shown to

1 be more effective from a toxicity reduction or quality of life improvement perspective. In Table 1,
2 we provide several scenarios of how the new composite endpoint of HUS may improve
3 interpretation of clinical trial findings if this composite endpoint was used in place of standard
4 primary endpoints. For example, one may consider three scenarios in which a new treatment is
5 deemed non-inferior based on a primary outcome of survival in a typical non-inferiority design
6 where different utility scores may produce drastically different trial conclusions if a composite
7 HUS endpoint were used. If a new intervention had lower utility than the comparator, a non-
8 inferior trial would declare the new intervention non-inferior, when in fact, a HUS endpoint
9 would appropriately declare the new intervention inferior. In addition, as we will show in the
10 simulations, sufficient power may be achieved with smaller sample sizes to make statistical
11 inferences than non-inferiority trials based on a non-inferiority margin of survival. This feature
12 may improve the efficiency of trial conduct and arriving at meaningful conclusions with smaller
13 samples.

14
15 This manuscript is structured as follows. In section 2, we present the methodology of the HUS
16 endpoint, including its construction, sample size calculation and power analysis. In section 3, we
17 use comprehensive simulation studies with various settings, including scenarios incorporating
18 parameter estimates based on the PET-NECK trial¹ to demonstrate the power advantage of HUS
19 when analyzing study data and its potential to reduce required sample sizes when designing new
20 trials. At last, we provide a discussion on the advantages, limitations and future directions for
21 HUS in section 4.

22

1 2. METHODS

2 2.1. HEALTH UTILITY ADJUSTED SURVIVAL (HUS)

3 In this section, we describe the basic framework of Health Utility adjusted Survival (HUS). In
4 many clinical studies, overall survival is chosen as the primary endpoint, which determines the
5 sample size, while health utility scores are usually analyzed in the secondary analyses. To
6 construct a composite endpoint combining survival and health utility, we can take the product of
7 the survival curve and the utility curve, as illustrated by Figure 1.

8

9 Suppose the total length of the study follow up time is T , and we are interested in comparing
10 survival and health utility between treatment groups 1 and 2. We define a Q statistic to represent
11 the HUS of each treatment group as

$$Q_1 = \int_0^T S_1(t) \bar{U}_1(t) dt, \quad (1)$$

$$Q_2 = \int_0^T S_2(t) \bar{U}_2(t) dt, \quad (2)$$

12 where $S_1(t)$ and $\bar{U}_1(t)$ are the survival function (proportion of patients alive at t) and average
13 utility score of those alive at t for group 1. $S_2(t)$ and $\bar{U}_2(t)$ are the survival function and average
14 utility score of those alive at t for group 2. We propose to use the Kaplan Meier (KM) estimated
15 survival functions $\hat{S}_1(t)$, $\hat{S}_2(t)$ to substitute $S_1(t)$, $S_2(t)$.

16

17 We can also assign weights to the survival and utility separately by defining

$$Q_1 = \int_0^T [S_1(t)]^{\lambda_1} [\bar{U}_1(t)]^{\lambda_2} dt, \quad (3)$$

$$Q_2 = \int_0^T [S_2(t)]^{\lambda_1} [\bar{U}_2(t)]^{\lambda_2} dt. \quad (4)$$

1 If $\lambda_1 = 0$ and $\lambda_2 = 1$, then Q_1 and Q_2 only consider the utility functions without including
2 survival. If $\lambda_1 = 1$ and $\lambda_2 = 0$, then Q_1 and Q_2 simply calculate the areas under the survival
3 curves without adjusting for utility. For simplicity, we suggest fixing the weight λ_1 as 1, since
4 survival is usually considered as important. λ_2 can be chosen from different values (e.g., 0.5, 1,
5 2), and $\lambda_2 = 1$ leads to the standard definition of HUS. The higher λ_2 is, the more importance is
6 assigned to health utility. For the rest of this manuscript, we focus on $\lambda_1 = 1$ and $\lambda_2 = 1$ unless
7 otherwise specified. We will also show some results with various λ_2 in our simulation studies
8 and discuss its effect.

9

10 2.2. HYPOTHESIS TESTING

11 To examine the difference of HUS between two treatment groups, we can define the test statistic
12 as

$$\mathcal{T} = Q_1 - Q_2. \quad (5)$$

13 To perform a one-sided test on whether group 1 has better HUS than group 2, we can either use
14 the bootstrap method to obtain the confidence interval of \mathcal{T} , or use the permutation method to
15 obtain the distribution of \mathcal{T} under the null hypothesis.²⁵ We can reject or accept the null
16 hypothesis ($H_0: \mathcal{T} \leq 0$) based on bootstrap confidence intervals. Suppose groups 1 and 2 have n_1

1 and n_2 subjects respectively, and the chosen significance threshold is α . The bootstrap procedure
2 can be described as follows:

- 3 1. For iteration b ($b = 1, \dots, B$), take a bootstrap dataset from the original samples, meaning
4 that we randomly sample n_1 subjects with replacement from treatment group 1 to be group 1
5 in the new sample, n_2 subjects with replacement from treatment group 2 to be group 2 in the
6 new sample.
- 7 2. Calculate the \mathcal{T} test statistic for the new sample, denoted by $\mathcal{T}^{(b)}$.
- 8 3. After obtaining $\mathcal{T}^{(b)}$'s ($b = 1, \dots, B$), calculate the $(1 - \alpha)$ confidence interval based on
9 these B bootstrap samples. If the confidence interval does not contain 0, reject the null
10 hypothesis. Note that the confidence interval should be constructed based on the test of
11 interest (one-sided or two-sided).

12 The permutation procedure can be described as follows:

- 13 1. For iteration b ($b = 1, \dots, B$), permute on the original samples to get a new permutation
14 dataset, meaning that we randomly reassign all of the subjects into two groups with sample
15 sizes n_1 and n_2 .
- 16 2. Calculate the \mathcal{T} test statistic for the new sample, denoted by $\mathcal{T}^{(b)}$.
- 17 3. After obtaining $\mathcal{T}^{(b)}$'s ($b = 1, \dots, B$), calculate the $(1 - \alpha)$ confidence interval based on
18 these B permutation samples. If the observed test statistic \mathcal{T} is outside the confidence interval,
19 reject the null hypothesis.

20 Note that the distribution generated by bootstrap is under the alternative hypothesis, whereas the
21 distribution generated by permutation is under the null hypothesis, which is why the former is
22 compared with 0, while the latter is compared with the observed test statistic. Based on our

1 experience, both bootstrap and permutation methods can control type I errors, but bootstrap tends
2 to have slightly higher power than permutation. Hence, we focus on the bootstrap method by
3 default. Some simulation results comparing bootstrap and permutation can be found in the
4 supplementary materials (Table S4, Figure S1). Besides, as hinted by Glasziou et al.,¹³ Jackknife
5 resampling can also be used to obtain the distribution of \mathcal{T} under the alternative hypothesis.^{26, 27}
6 However, our past experience shows that there is little difference in terms of type I error and
7 power when comparing the bootstrap method with Jackknife, while the distribution of \mathcal{T} based
8 on bootstrap samples tends to be closer to normal. As a result, we suggest using the bootstrap
9 method as default. In terms of the number of resamples, $B = 500$ is usually sufficient for
10 controlling type I errors and obtaining decent power. Examples showing the performance of
11 Jackknife and evaluating the choice of B are also provided in the supplementary materials
12 (Tables S4-S5, Figure S1).

13

14 2.3. THEORETICAL PROPERTIES

15 If we assume the survival time follows a piecewise exponential distribution, we can derive a
16 Monte Carlo approach to calculate the variance of the test statistic, which can be used for power
17 analysis and sample size calculation.²⁸ A similar idea was used by Royston and Parmar²⁹ to
18 calculate the variance for restricted mean survival time (RMST).³⁰⁻³²

19

20 We consider a simple case with three key time points: 0 (baseline), C (end of surgery) and T (end
21 of study). Focusing on one treatment group, suppose the survival time is piecewise exponential,
22 with piecewise constant hazards h_1, h_2 for time periods $0 \sim C, C \sim T$ respectively. The utility
23 function is piecewise linear, which starts from A_1 at time 0, changes to A_2 at time C , and then

1 goes to A_3 at time T . Let $X = \min(\xi, T)$, where ξ is the survival time with cumulative hazard
 2 function $H(t)$ and survival function $S(t)$. We can decompose X as $X_1 + X_2$, where

$$X_1 = \begin{cases} \xi & (0 \leq \xi \leq C) \\ C & (\xi > C) \end{cases}, \quad (6)$$

$$X_2 = \begin{cases} 0 & (0 \leq \xi \leq C) \\ \xi - C & (C < \xi \leq T) \\ T - C & (\xi > T) \end{cases}. \quad (7)$$

3 Denote $M = \int_{t=0}^T S(t)U_0(t)dt$ where $U_0(t)$ is the base utility function for the currently
 4 considered treatment group. Write its statistic of HUS as $Q = \int_{t=0}^T \hat{S}(t)\bar{U}(t)dt$. If we define

$$X^* = A_1X_1 + \frac{TA_2 - CA_3}{T - C}X_2 + \frac{A_2 - A_1}{2C}X_1^2 + \frac{A_3 - A_2}{2(T - C)}X_2^2 + \frac{A_3 - A_2}{T - C}X_1X_2, \quad (8)$$

5 we can derive that $M = E(X^*)$. Following Royston & Parmar (2013), we can assume that for a
 6 specific scenario, we have

$$SE(Q) = \phi \frac{SD(X^*)}{\sqrt{n}}, \quad (9)$$

7 where ϕ is a factor no less than 1 and n is the sample size for the group we are currently looking
 8 at. For convenience, we call ϕ the variance balance factor, which takes account of the extra
 9 variance introduced into HUS by missing utility, censored survival, KM estimation, etc. $SD(X^*)$
 10 can be calculated using the parameters, while ϕ can be estimated by Monte Carlo sampling.
 11 More details including the derivations are provided in the supplementary materials (Tables S2-
 12 S3). We will demonstrate in our simulations that ϕ is robust to different sample sizes.

13

1 Note that when two treatment groups are compared, they should have their own variance balance
2 factors, which we denote as ϕ_1 and ϕ_2 . Applying our assumed property to each of the groups,
3 we have

$$SE(Q_1) = \phi_1 \frac{SD(X^*_1)}{\sqrt{n_1}}, \quad (10)$$

$$SE(Q_2) = \phi_2 \frac{SD(X^*_2)}{\sqrt{n_2}}, \quad (11)$$

4 where Q_1 and Q_2 are the statistics of HUS for treatment groups 1 and 2, and n_1, n_2 are the
5 sample sizes of the two groups. X^*_1 and X^*_2 are constructed separately for the two groups using
6 their own parameter settings. Hence, the variance of $\mathcal{T} = Q_1 - Q_2$ is

$$\text{var}(\mathcal{T}) = [SE(Q_1)]^2 + [SE(Q_2)]^2. \quad (12)$$

7 For the one-sided test, we can reject the null hypothesis if $\mathcal{T} - z_{1-\alpha}\sqrt{\text{var}(\mathcal{T})} > 0$.

8

9 2.4. POWER ANALYSIS AND SAMPLE SIZE CALCULATION

10 In any scenario with prespecified parameters, given different sample size, we can calculate the
11 corresponding power of HUS using simulations. Then we can obtain a table showing different
12 power under different sample sizes, which can be used to determine the sample size needed to
13 achieve specific power (e.g., 80%) for a new trial. Detailed examples are provided in section 3.1.

14

15 If we assume that the special case described in section 2.3 is true, then we only need to run one
16 simulation given a fixed sample size (e.g., 200 subjects per treatment group), which can give us

1 estimates of ϕ_1 and ϕ_2 . For the one-sided test where we reject the null hypothesis if $\mathcal{T} -$

2 $z_{1-\alpha}\sqrt{\text{var}(\mathcal{T})} > 0$, the power is

$$\omega = P\left(\mathcal{T} - z_{1-\alpha}\sqrt{\text{var}(\mathcal{T})} > 0\right) = P\left(\mathcal{T} > z_{1-\alpha}\sqrt{\text{var}(\mathcal{T})}\right). \quad (13)$$

3 Assume \mathcal{T} follows $N(\mathcal{T}_{\text{true}}, \text{var}(\mathcal{T}))$ and denote the power by ω , we have

$$\omega = \Phi\left(\frac{\mathcal{T}_{\text{true}}}{\sqrt{\text{var}(\mathcal{T})}} - z_{1-\alpha}\right), \quad (14)$$

4 where Φ is the cumulative distribution function of the standard normal distribution. On the other

5 hand, to achieve power ω , the required sample sizes should satisfy

$$\phi_1^2 \frac{\text{SD}(X^*_1)^2}{n_1} + \phi_2^2 \frac{\text{SD}(X^*_2)^2}{n_2} = \left(\frac{\mathcal{T}_{\text{true}}}{\Phi^{-1}(\omega) + z_{1-\alpha}}\right)^2. \quad (15)$$

6 If we assume $n_1 = n_2$, then the required sample size per arm is

$$n_1 = \frac{(\Phi^{-1}(\omega) + z_{1-\alpha})^2 [\phi_1^2 \text{var}(X^*_1) + \phi_2^2 \text{var}(X^*_2)]}{\mathcal{T}_{\text{true}}^2}. \quad (16)$$

7 Note that it is difficult to calculate $\mathcal{T}_{\text{true}}$ based on the setting of parameters. However, we can

8 estimate it by using the average of the observed \mathcal{T} from our simulated samples. To summarize, in

9 the special situation with simplified settings described in section 2.3, we can use the following

10 procedure to calculate power yielded by a specific sample size:

11 1. Calculate $\text{var}(X^*_1)$, $\text{var}(X^*_2)$ based on parameter settings.

12 2. Simulate samples to estimate ϕ_1 , ϕ_2 and $\mathcal{T}_{\text{true}}$.

13 3. For each new sample size combination n_1, n_2 , calculate $\text{SE}(Q_1)$, $\text{SE}(Q_2)$ using the estimated

14 ϕ_1, ϕ_2 .

1 4. Calculate $\text{var}(\mathcal{T})$ and power ω .

2

3 On the other side, we can use the following procedure to calculate the sample size required to
4 achieve specific power:

5 1. Calculate $\text{var}(X^*_1)$, $\text{var}(X^*_2)$ based on parameter settings.

6 2. Simulate samples to estimate ϕ_1 , ϕ_2 and $\mathcal{T}_{\text{true}}$.

7 3. Calculate n_1 using the sample size formula.

8

9 2.5. HANDLING MISSING UTILITY SCORES

10 In clinical studies, utility scores may not be available at each time point for all subjects, while the
11 current framework of HUS requires complete utility profiles to calculate the test statistic. The
12 most intuitive way is to impute the utility scores. We use linear functions to fill in the utility
13 scores using the available data. If a subject's utility score is only available at one time point, then
14 we use that score as the imputed utility at all other time points. This approach may seem simple,
15 but it can be quite effective. Another method we consider is to impute the group average at each
16 key time point (i.e., each time point at which at least one subject has their utility score recorded),
17 and then use linear functions to fill in the other missing scores. This approach can be regarded as
18 a combination of the cross-mean and linear interpolation methods³³. While imputing the group
19 average, we can also add some variation using a normal distribution with mean zero and its
20 standard deviation equal to the standard deviation of the recorded scores at that time point. In
21 this way, the imputed values may be closer to the true values, which may lead to increase of
22 statistical power. It is also worth noting that many other methods are available for imputing
23 longitudinal data, and a very recent study has compared the effects of different imputation

1 methods and shown that most of them are similar in various scenarios, whereas trajectory mean
2 single imputation has the best overall performance³³. Hence, we consider trajectory mean
3 imputation as a third method. A comparison of the three methods using simulation results is
4 provided in the supplementary materials (Table S1), which shows that method 1 has much worse
5 performance when the missing rate is higher, while methods 2 and 3 are not affected as much.
6 For convenience, we use method 2 by default.

7 3. RESULTS

8 3.1. SIMULATIONS WITH SIMPLIFIED SETTINGS

9 3.1.1 POWER COMPARISON

10 We conduct simulations in various scenarios to assess the performances of HUS. Suppose we are
11 designing a randomized clinical trial with two treatment arms. The total length of study is 36
12 months ($T = 36$), and each patient receives surgery at 3 months ($C = 3$). The two arms are
13 assigned to different treatment strategies to help them recover, and we are interested in
14 comparing the two treatments in terms of both survival and health utility. Denote the true
15 survival time, observed survival time and survival status for patient i from group g as T_{gi} , X_{gi}
16 and δ_{gi} respectively. Group 1 and group 2 have sample sizes n_1 and n_2 . The survival data is
17 simulated using

$$T_{gi} \sim \text{Exp}(h_g),$$

$$\xi_{gi} \sim \text{Unif}(0, \zeta),$$

$$X_{gi} = \min(T_{gi}, \xi_{gi}, T),$$

$$\delta_{gi} = \begin{cases} 1 & (T_{gi} < \xi_{gi} \text{ and } T_{gi} < T) \\ 0 & (\text{otherwise}) \end{cases}'$$

1 where ζ is chosen to control the censoring rate, denoted by $p_{\text{censoring}}$. The hazard ratio of
2 treatment 1 against treatment 2 is h_1/h_2 .

3

4 To simulate the health utility score, we first define base utility functions for the two groups. The
5 base utility at time t for group g can be written as

$$U_{g0}(t) = \begin{cases} A_{g1} + \frac{A_{g2} - A_{g1}}{C} t & (0 \leq t \leq C) \\ \frac{TA_{g2} - CA_{g3}}{T - C} + \frac{A_{g3} - A_{g2}}{T - C} t & (C < t \leq T) \end{cases}.$$

6 This definition means the average utility for group g starts from A_{g1} at baseline, changes to A_{g2}
7 at 3 months, and then changes to A_{g3} at the end of the study. The change is piecewise linear. Our
8 motivation for this setting is that usually a cancer patient's health utility reaches the lowest at the
9 end of treatment and gradually recovers after that. For patient i from group g , the health utility
10 score at time t , denoted by $U_{gi}(t)$, follows a normal distribution with mean $U_{g0}(t)$ and standard
11 deviation 0.1.

12

13 In practice, we do not expect health utility scores to be collected at each time point. Furthermore,
14 some of the scores scheduled to be collected may be missing. For our main simulation study, we
15 assume that the health utility scores are only collected at $t = 1, C$ and T . When $t = 1$, all
16 subjects have their utility scores collected. When $t = C$ or T , the subjects that are still being
17 followed have their utility scores collected, while there is a p_{missingU} chance that the score is
18 missing.

19

1 In this section, we focus on the situation where the two treatment groups do not have a difference
2 in OS, which is the situation that motivated our HUS framework. Other situations (e.g., the two
3 treatment groups differ in both OS and health utility) are explored in section 3.2 and the
4 supplementary materials (Tables S6-S9). Table 2 shows a summary of our major scenarios. In
5 each scenario, we compare the theoretical rejection rate using our results from section 2.4 and
6 the empirical rejection rates of HUS using bootstrap with $B = 500$. We consider three choices of
7 λ_2 : $\lambda_2 = 1$ corresponds to the standard HUS approach; $\lambda_2 = 0.5$ means giving utility less weight
8 than survival; $\lambda_2 = 2$ means giving utility more weight than survival. We also examine the
9 performance of OS-based tests. In the tables, “sup” represents the log-rank test that tests whether
10 group 1 is superior to group 2 in terms of OS using KM estimates. “5%” and “10%” correspond
11 to the inferiority test using the hazard ratio with margins 5% and 10% respectively. For instance,
12 a 5% margin means we establish non-inferiority (treatment 1 is non-inferior to treatment 2 in
13 terms of OS) if the upper bound of the 95% CI of the hazard ratio is smaller than 1.05.

14

15 In scenario 0, we examine the rejection rates of different methods when the two treatment groups
16 have the same OS and health utility. As shown in Table 3, all of the superiority tests are able to
17 control type I errors at 0.05. The rejection rates of the non-inferiority tests are power instead of
18 type I errors, since the alternative is true (treatment 1 is not inferior to treatment 2). This is why
19 they may be higher than 0.05.

20

21 In scenario 1, we compare the power of different methods when treatment group 1 has better
22 health utility than treatment group 2. For the theoretical power analysis, firstly, we run one
23 simulation with $n_1 = n_2 = 200$ and 4000 replications to obtain the estimates $\phi_1 = 1.07$, $\phi_2 =$

1 1.12, $\mathcal{T}_{\text{true}} = 3.11$. Then we can calculate the power of different sample sizes. For the other
2 methods such as bootstrap $\lambda_2 = 1$, $\lambda_2 = 0.5$, and $\lambda_2 = 2$, we need to simulate new datasets (200
3 replications) with different sample sizes to get the empirical power. Note that ϕ_1 , ϕ_2 and $\mathcal{T}_{\text{true}}$
4 are quite robust to different sample sizes. For example, if we use $n_1 = n_2 = 500$, the obtained
5 estimates are $\phi_1 = 1.06$, $\phi_2 = 1.11$, $\mathcal{T}_{\text{true}} = 3.11$, which is very close to the scenario of $n_1 =$
6 $n_2 = 200$. More results regarding the variance balance factors are provided in the supplementary
7 materials (Tables S2-S3).

8
9 As shown in Table 4, the bootstrap method with $\lambda_2 = 1$ performs close to the theoretical results,
10 which makes sense since the theoretical results are based on the standard HUS with $\lambda_2 = 1$.
11 Larger λ_2 tends to lead to higher power by giving more weight to utility than survival. This is
12 also expected because the two groups only differ in terms of utility. Meanwhile, the superiority
13 and non-inferiority tests based on OS have little power since there is no real difference in the two
14 group's OS. We also calculate the power corresponding to different sample sizes using our
15 theoretical results and plot the power curves in Figure 2. For scenario 1, to achieve 80% power,
16 using HUS as the endpoint only requires 85 subjects per arm.

17
18 In scenario 2, we increase the censoring rate to 60% and missing rates to 60%, and reduce the
19 difference between the two group's health utility scores. As shown in Table 4 and Figure 2,
20 results are very similar to those in scenario 1. Again, HUS is able to obtain decent power with
21 relatively small sample sizes while the superiority and non-inferiority tests struggle to find
22 enough evidence to show treatment 1's benefit compared to treatment 2. If we design a trial

1 based on HUS with the assumptions in scenario 2, we only need to have 151 patients in each
2 treatment group.

3
4 We would like to point out that even though choosing a larger λ_2 may seem to have higher
5 power in the above scenarios, it may not always be a good choice, especially when there is a
6 difference in OS. We recommend using $\lambda_2 = 1$ as default, though it can be modified depending
7 on the knowledge of the two treatments (e.g., whether treatment 1 is likely to have better OS than
8 treatment 2).

9
10 *3.1.2 SAMPLE SIZE CALCULATION*
11 In this subsection, we use our developed sample size calculation formulas to calculate sample
12 sizes needed for the composite endpoint, and the standard formulas to calculate sample sizes
13 needed for basic survival endpoint (implemented in PASS 2023, v23.0.2 with the one-sided log-
14 rank test), to further demonstrate the advantage of HUS. Following scenario 1 from 3.1.1, where
15 treatment 1 has better utility than treatment 2, we consider four different cases. In the first case,
16 there is no survival difference, which is consistent with the focus of this manuscript, and the
17 endpoint overall survival does not have power. In the second case, we assume that treatment 1
18 has better survival than treatment 2, while in the third case, we assume that treatment 2 has better
19 survival. In the last case, we assume treatment 1 has better survival, but there is no difference in
20 utility, and the utility function is the same as that in scenario 0. As shown in Table 5, with
21 scenario 1's utility functions, when h_1 is smaller than h_2 , meaning that treatment 1 has better OS
22 than treatment 2, the required sample size for HUS is decreased, which makes sense because the
23 difference in HUS is larger. When h_1 is larger than h_2 , meaning that treatment 1 has worse OS

1 than treatment 2, the required sample size for HUS is increased. Nevertheless, the numbers are
 2 still much smaller than those calculated for overall survival. If there is no utility difference, HUS
 3 will require more subjects than overall survival, which is expected, though the difference is not
 4 as big. These results show again that using the composite endpoint may help greatly reduce the
 5 required sample size to detect a significant difference when two treatments differ in utility.

6
 7 **Table 5.** Sample size calculation under a significance level of 0.05. When there is a utility
 8 difference, utility functions from scenario 1 are used. When there is no utility difference, utility
 9 functions from scenario 0 are used.

Utility: 1>2; survival: 1=2 ($h_1/h_2 = 1$)		
Targeted power	Sample size requirement for each arm	
	HUS	OS
70%	65	/
80%	85	/
90%	118	/
Utility: 1>2; survival: 1>2 ($h_1/h_2 = 0.9$)		
Targeted power	Sample size requirement for each arm	
	HUS	OS
70%	47	1710
80%	62	2247
90%	85	3112
Utility: 1>2; survival: 1<2 ($h_1/h_2 = 1.1$)		
Targeted power	Sample size requirement for each arm	
	HUS	OS
70%	78	2243
80%	103	2946
90%	143	4080

Utility: 1=2; survival: 1>2 ($h_1/h_2 = 0.7$)		
Targeted power	Sample size requirement for each arm	
	HUS	OS
70%	184	138
80%	242	180
90%	335	249

1

2

3 3.3. SIMULATIONS WITH REAL DATA ESTIMATES

4 To demonstrate the benefit of HUS in a more practical scenario, we conduct additional
5 simulations with average utility scores and the hazard ratio mimicking the summary data
6 provided in a real randomized trial PET-NECK.¹ PET-NECK is a randomized phase III non-
7 inferiority trial that compares Positron emission tomography-computerized tomography-guided
8 watch-and-wait policy (PET-CT) with planned neck dissection (planned ND) for head and neck
9 cancer patients. The two-year overall survival rates of the two treatment groups (PET-CT and
10 planned ND) with 282 subjects per arm, are 84.9% and 81.5% respectively, which leads to a
11 hazard ratio of 0.80. We conducted simulations utilizing the parameter setting to emulate the
12 survival times in PET-NECK. Figure 3 shows the average utility scores at different time points in
13 the study, with the maximum time being 24 months. Hence, in this scenario, we set $T = 24$ and
14 define the base utility functions following the observed average utility scores. We also record the
15 utilities at baseline and months 1, 3, 6, 12, 24 with 30% missing rate. Note that this scenario does
16 not fall into the framework of section 2.3, and thus we cannot apply our theoretical results
17 directly to calculate the power and sample sizes. However, obtaining the empirical results is
18 similar to what we describe in section 3.1.

1
2 As shown in Table 6, with the two groups differing in both OS and health utility, the superiority
3 test based on HUS still has much higher power than the superiority and non-inferiority tests
4 based on OS. We would only need 200 subjects per arm to achieve 80% power of showing PET-
5 CT has better HUS than planned ND, which is fewer than the subjects in the original study which
6 were based on OS comparison. In terms of weighting, $\lambda_2 = 2$ again leads to higher power, while
7 $\lambda_2 = 0.5$ has lower power compared to the standard HUS. Nevertheless, in certain scenarios,
8 especially if the difference in health utility is small, using a larger λ_2 may not be as beneficial.
9 More results, including scenarios where there is no difference in health utility, are available in
10 the supplementary materials (Tables S6-S9, Figures S2-S4).

11

12 4. DISCUSSION

13 We have presented a methodological framework to compare two treatment groups using HUS as
14 a composite endpoint combining survival and health utility. As demonstrated by our
15 comprehensive simulation studies, when there is a difference in health utility, HUS has a
16 significant power advantage over the statistical tests based on OS endpoint, meaning that using
17 HUS as an endpoint for new trials may require much smaller sample sizes to achieve decent
18 power. We have also demonstrated two different procedures (theoretical and empirical
19 approaches) to conduct power analysis and sample size calculation with specified parameters.
20 When the model assumptions are met, the two procedures yield similar results.
21 There are several different options when applying HUS. We recommend using bootstrap given
22 its popularity as well as its convenience of constructing confidence intervals for the test statistic,

1 though permutation may be theoretically more appropriate for testing the null hypothesis, since it
2 can obtain the null distribution of the test statistic. In terms of weighting on survival and utility,
3 we recommend choosing weights $\lambda_1 = \lambda_2 = 1$ as default. Using a larger λ_2 may increase the
4 power in certain scenarios, especially when the two treatment groups differ in health utility but
5 not in OS. However, it may not be beneficial when there is no difference in health utility. A
6 possible way to combine different weighting options without having to choose one is to apply an
7 idea similar to the aSPU test.^{34, 35}

8
9 Note that the theoretical properties we have presented are based on assumptions by analogy with
10 the assumptions used by Royston and Parmar,²⁹ though we have shown the validity of our
11 theoretical results in our simulations. In the future, we may explore the asymptotic properties
12 with relaxed assumptions (e.g., the survival times do not have to be piece-wise exponential),
13 which would be helpful when designing or analyzing trials where our previously used
14 assumptions are likely to be violated. We would also like to point out that the linear imputation
15 method we use to fill in the utility scores may be problematic in some cases, especially if the
16 scores are only recorded at a few time points and the missing rate is high. We may consider other
17 imputation methods or modifying the definition of HUS so that it does not require complete
18 utility score profiles as input.^{36, 37} Besides that, given the possible drawbacks brought by KM
19 estimates, sometimes it may be beneficial to apply other models, including the flexible
20 parametric model for survival analysis.³⁸

21
22 Another possible direction worth exploring is to take different functions of the utility score into
23 consideration. One special case is that in many clinical studies, multiple measures of health

1 status are recorded. There are various ways to combine different measures into a single utility
2 score.³⁹⁻⁴¹ Extending HUS to be able to handle any function of utility may potentially increase
3 power and help us gain insight on how the utility is different in different treatment groups.
4 Furthermore, considering utility may have different importance at different time points, we may
5 assign different weights across time. For example, having a better utility score at the later stage
6 of the study, which means the patients have recovered better, may be more important than having
7 a better utility score at the end of surgery. In such case, we can consider giving higher weights to
8 later time points, and the resulted HUS may provide a clearer picture of which treatment is more
9 beneficial for recovery.

10

11 5. DATA AVAILABILITY

12 R code for our simulation studies and summary data of health utility are available at
13 <https://github.com/yangq001/HUS>.

14

15 CONFLICTS OF INTEREST

16 The authors declare that there are no conflicts of interest.

17

18 ACKNOWLEDGMENTS

19 The authors would like to acknowledge the contributions of Dr. Hisham Mehanna (Institute of
20 Head and Neck Studies and Education, University of Birmingham) and Dr. Sue Yom
21 (Department of Radiation Oncology, University of California) for clinical insights and discussion.

1

2 FUNDING

3 This work was supported by the Alan Brown Chair in Molecular Genomics, the Lusi Wong Family Fund,
4 and the Posluns Family Fund, all through the Princess Margaret Cancer Foundation.

5

6

REFERENCES

1. Mehanna H, McConkey CC, Rahman JK, et al: PET-NECK: a multicentre randomised Phase III non-inferiority trial comparing a positron emission tomography–computerised tomography-guided watch-and-wait policy with planned neck dissection in the management of locally advanced (N2/N3) nodal metastases in patients with squamous cell head and neck cancer. *Health Technol Assess* 21:1–122, 2017
2. Mathias SD, Bates MM, Pasta DJ, et al: Use of the Health Utilities Index With Stroke Patients and Their Caregivers. *Stroke* 28:1888–1894, 1997
3. Horsman J, Furlong W, Feeny D, et al: The Health Utilities Index (HUI®): concepts, measurement properties and applications. *Health Qual Life Outcomes* 1:54, 2003
4. Jewell EL, Smrtka M, Broadwater G, et al: Utility Scores and Treatment Preferences for Clinical Early-Stage Cervical Cancer. *Value in Health* 14:582–586, 2011
5. Glasziou PP, Simes RJ, Gelber RD: Quality adjusted survival analysis. *Statist Med* 9:1259–1276, 1990
6. Gelber RD: Quality-of-Life-Adjusted Evaluation of Adjuvant Therapies for Operable Breast Cancer. *Ann Intern Med* 114:621, 1991
7. Gelber RD, Goldhirsch A, Cole BF, et al: A Quality-Adjusted Time Without Symptoms or Toxicity (Q-TWiST) Analysis of Adjuvant Radiation Therapy and Chemotherapy for Resectable Rectal Cancer. *JNCI Journal of the National Cancer Institute* 88:1039–1045, 1996
8. Murray S, Cole B: Variance and Sample Size Calculations in Quality-of-Life-Adjusted Survival Analysis (Q-TWiST). *Biometrics* 56:173–182, 2000
9. Konski AA, Winter K, Cole BF, et al: Quality-adjusted survival analysis of Radiation Therapy Oncology Group (RTOG) 90-03: Phase III randomized study comparing altered fractionation to standard

fractionation radiotherapy for locally advanced head and neck squamous cell carcinoma. *Head Neck* 31:207–212, 2009

10. Zbrozek AS, Hudes G, Levy D, et al: Q-TWiST Analysis of Patients Receiving Temsirolimus or Interferon Alpha for Treatment of Advanced Renal Cell Carcinoma. *Pharmacoeconomics* 28:577–584, 2010

11. Seymour JF, Gaitonde P, Emeribe U, et al: A Quality-Adjusted Survival (Q-TWiST) Analysis to Assess Benefit-Risk of Acalabrutinib Versus Idelalisib/Bendamustine Plus Rituximab or Ibrutinib Among Relapsed/Refractory (R/R) Chronic Lymphocytic Leukemia (CLL) Patients. *Blood* 138:3722–3722, 2021

12. Jerusalem G, Delea TE, Martin M, et al: Quality-Adjusted Survival with Ribociclib Plus Fulvestrant Versus Placebo Plus Fulvestrant in Postmenopausal Women with HR+HER2– Advanced Breast Cancer in the MONALEESA-3 Trial. *Clinical Breast Cancer* 22:326–335, 2022

13. Glasziou PP, Cole BF, Gelber RD, et al: Quality adjusted survival analysis with repeated quality of life measures. *Stat Med* 17:1215–1229, 1998

14. Prieto L, Sacristán JA: Problems and solutions in calculating quality-adjusted life years (QALYs). *Health Qual Life Outcomes* 1:80, 2003

15. Whitehead SJ, Ali S: Health outcomes in economic evaluation: the QALY and utilities. *British Medical Bulletin* 96:5–21, 2010

16. Touray MML: Estimation of Quality-adjusted Life Years alongside clinical trials: the impact of ‘time-effects’ on trial results. *J Pharm Health Serv Res* 9:109–114, 2018

17. Chung C-H, Hu T-H, Wang J-D, et al: Estimation of Quality-Adjusted Life Expectancy of Patients With Oral Cancer: Integration of Lifetime Survival With Repeated Quality-of-Life Measurements. *Value in Health Regional Issues* 21:59–65, 2020

18. Laska EM, Meisner M, Siegel C: Power and Sample Size in Cost- Effectiveness Analysis. *Med Decis Making* 19:339–343, 1999
19. Willan AR, Lin DY: Incremental net benefit in randomized clinical trials. *Statist Med* 20:1563–1574, 2001
20. Hollingworth W, McKell-Redwood D, Hampson L, et al: Cost–utility analysis conducted alongside randomized controlled trials: Are economic end points considered in sample size calculations and does it matter? *Clinical Trials* 10:43–53, 2013
21. Bader C, Cossin S, Maillard A, et al: A new approach for sample size calculation in cost-effectiveness studies based on value of information. *BMC Med Res Methodol* 18:113, 2018
22. Billingham LJ, Abrams KR, Jones DR: Methods for the analysis of quality-of-life and survival data in health technology assessment. *Health Technol Assess* 3:1–152, 1999
23. Diaby V, Adunlin G, Ali AA, et al: Using quality-adjusted progression-free survival as an outcome measure to assess the benefits of cancer drugs in randomized-controlled trials: case of the BOLERO-2 trial. *Breast Cancer Res Treat* 146:669–673, 2014
24. Oza AM, Lorusso D, Aghajanian C, et al: Patient-Centered Outcomes in ARIEL3, a Phase III, Randomized, Placebo-Controlled Trial of Rucaparib Maintenance Treatment in Patients With Recurrent Ovarian Carcinoma. *JCO* 38:3494–3505, 2020
25. Good PI: *Permutation, parametric and bootstrap tests of hypotheses* 3rd ed. New York, Springer, 2005
26. Wu CFJ: Jackknife, Bootstrap and Other Resampling Methods in Regression Analysis [Internet]. *Ann Statist* 14, 1986[cited 2022 Dec 14] Available from: <https://projecteuclid.org/journals/annals-of-statistics/volume-14/issue-4/Jackknife-Bootstrap-and-Other-Resampling-Methods-in-Regression-Analysis/10.1214/aos/1176350142.full>

27. Shao J, Tu D: The Jackknife and Bootstrap [Internet]. New York, NY, Springer New York, 1995[cited 2022 Dec 14] Available from: <http://link.springer.com/10.1007/978-1-4612-0795-5>
28. Myers ND, Ahn S, Jin Y: Sample Size and Power Estimates for a Confirmatory Factor Analytic Model in Exercise and Sport: A Monte Carlo Approach. *Research Quarterly for Exercise and Sport* 82:412–423, 2011
29. Royston P, Parmar MKB: Restricted mean survival time: an alternative to the hazard ratio for the design and analysis of randomized trials with a time-to-event outcome. *BMC Med Res Methodol* 13:152, 2013
30. Irwin JO: The standard error of an estimate of expectation of life, with special reference to expectation of tumourless life in experiments with mice. *J Hyg* 47:188–189, 1949
31. Royston P, Parmar MKB: The use of restricted mean survival time to estimate the treatment effect in randomized clinical trials when the proportional hazards assumption is in doubt. *Statist Med* 30:2409–2421, 2011
32. Zhao L, Claggett B, Tian L, et al: On the restricted mean survival time curve in survival analysis: On the Restricted Mean Survival Time Curve in Survival Analysis. *Biom* 72:215–221, 2016
33. Jahangiri M, Kazemnejad A, Goldfeld KS, et al: A wide range of missing imputation approaches in longitudinal data: a simulation study and real data analysis. *BMC Med Res Methodol* 23:161, 2023
34. Pan W, Kim J, Zhang Y, et al: A Powerful and Adaptive Association Test for Rare Variants. *Genetics* 197:1081–1095, 2014
35. Kim J, Bai Y, Pan W: An Adaptive Association Test for Multiple Phenotypes with GWAS Summary Statistics. *Genet Epidemiol* 39:651–663, 2015

- 36.** Naeim A, Keeler EB, Mangione CM: Options for Handling Missing Data in the Health Utilities Index Mark 3. *Med Decis Making* 25:186–198, 2005
- 37.** Graham JW: Missing Data [Internet]. New York, NY, Springer New York, 2012[cited 2022 Dec 14]
Available from: <http://link.springer.com/10.1007/978-1-4614-4018-5>
- 38.** Lambert PC, Royston P: Further Development of Flexible Parametric Models for Survival Analysis. *The Stata Journal* 9:265–290, 2009
- 39.** Hawthorne G, Richardson J, Day NA: A comparison of the Assessment of Quality of Life (AQoL) with four other generic utility instruments. *Annals of Medicine* 33:358–370, 2001
- 40.** Fisk JD: A comparison of health utility measures for the evaluation of multiple sclerosis treatments. *Journal of Neurology, Neurosurgery & Psychiatry* 76:58–63, 2005
- 41.** Pickard AS, Ray S, Ganguli A, et al: Comparison of FACT- and EQ-5D–Based Utility Scores in Cancer. *Value in Health* 15:305–311, 2012

Table 1. Interpretations for different scenarios of survival and utility.

Scenario	Non-Inferiority Interpretation	Health Utility Interpretation	Clinical Interpretation and Caveats
Survival non-inferior Improved Utility	New Treatment Non-Inferior	New Treatment Superior	With composite endpoint, patients and clinicians can be confident that weighted health utility adjusted survival is superior.
Survival non-inferior Worse Utility	New Treatment Non-Inferior	New Treatment Not Superior	With non-inferiority design, the new treatment may be falsely accepted as a treatment option despite worse utility
Survival non-inferior Similar Utility	New Treatment Non-Inferior	New Treatment Not Superior	As above
Survival Inferior Improved Utility	New Treatment Inferior	New Treatment may be Superior, Similar or Worse Depending on magnitude of effect	With non-inferiority design new option is rejected as non-inferior. However, if there is a large therapeutic benefit with the new intervention a composite endpoint

			may demonstrated this new treatment to be superior.
Survival Inferior Worse Utility	New Treatment Inferior	New Treatment Inferior	Non-Inferior design may appropriately declare new treatment as inferior
Survival Inferior Similar Utility	New Treatment Inferior	New Treatment Inferior	As above

Table 2. Simulation settings with different scenario.

Scenario	$p_{\text{censoring}}$	p_{missingU}	Average utility			
			Group	Baseline	3 months	36 months
0	30%	30%	1	0.8	0.4	0.7
			2	0.8	0.4	0.7
1	30%	30%	1	0.8	0.5	0.8
			2	0.8	0.35	0.7
2	60%	60%	1	0.8	0.5	0.8
			2	0.8	0.4	0.7

Table 3. Rejection rates of different methods in scenario 0 based on 1000 replications.

n_1, n_2	HUS				OS		
	Theoretical	Bootstrap ($\lambda_2 = 1$)	$\lambda_2 = 0.5$	$\lambda_2 = 2$	Superiority test	Non- inferiority test with margin 5%	Non- inferiority test with margin 10%
50	0.052	0.052	0.050	0.053	0.056	0.040	0.050
100	0.052	0.054	0.053	0.051	0.068	0.051	0.073
150	0.053	0.047	0.052	0.048	0.054	0.042	0.079
200	0.053	0.052	0.046	0.048	0.053	0.050	0.092
500	0.054	0.049	0.051	0.050	0.045	0.083	0.184

Table 4. Power comparison of different methods in scenarios 1 and 2 based on 200 replications.

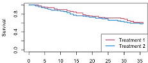
Scenario 1							
n_1, n_2	HUS				OS		
	Theoretical	Bootstrap ($\lambda_2 = 1$)	$\lambda_2 = 0.5$	$\lambda_2 = 2$	Superiority test	Non- inferiority test with margin 5%	Non- inferiority test with margin 10%

50	0.61	0.56	0.28	0.9	0.05	0.04	0.05
100	0.86	0.85	0.44	1	0.05	0.05	0.07
150	0.95	0.95	0.59	1	0.05	0.06	0.1
200	0.99	1	0.71	1	0.06	0.04	0.06
Scenario 2							
n_1, n_2	HUS				OS		
	Theoretical	Bootstrap ($\lambda_2 = 1$)	$\lambda_2 = 0.5$	$\lambda_2 = 2$	Superiority test	Non- inferiority test with margin 5%	Non- inferiority test with margin 10%
50	0.42	0.42	0.2	0.76	0.05	0.04	0.06
100	0.65	0.67	0.31	0.94	0.06	0.06	0.08
150	0.80	0.82	0.42	0.97	0.06	0.05	0.1
200	0.89	0.92	0.48	1	0.05	0.05	0.06

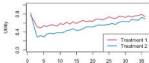
Table 6. Power comparison for simulations using real data estimates.

n_1, n_2	HUS	OS
------------	-----	----

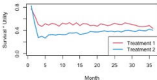
	Bootstrap ($\lambda_2 = 1$)	$\lambda_2 = 0.5$	$\lambda_2 = 2$	Superiority test	Non- inferiority test with margin 5%	Non- inferiority test with margin 10%
50	0.34	0.26	0.47	0.09	0.1	0.12
100	0.54	0.4	0.69	0.14	0.17	0.22
150	0.74	0.54	0.84	0.19	0.24	0.34
200	0.82	0.61	0.97	0.2	0.25	0.33
282	0.94	0.76	0.99	0.27	0.36	0.44

Survival

Month

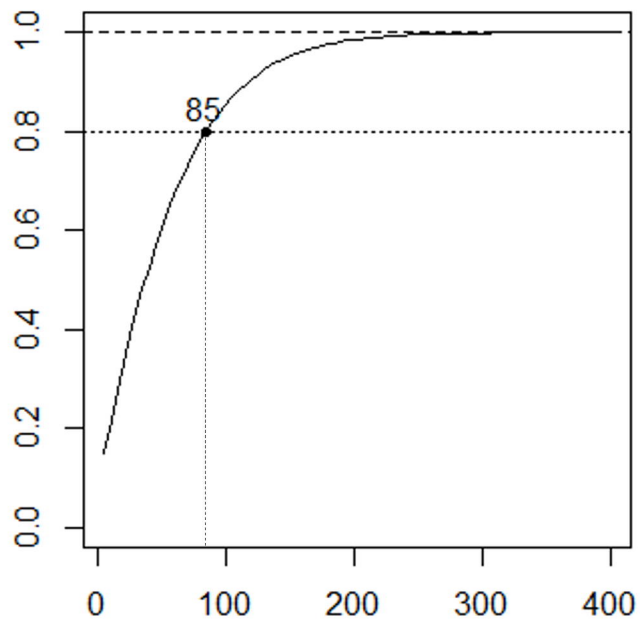
Utility

Month

Survival * Utility

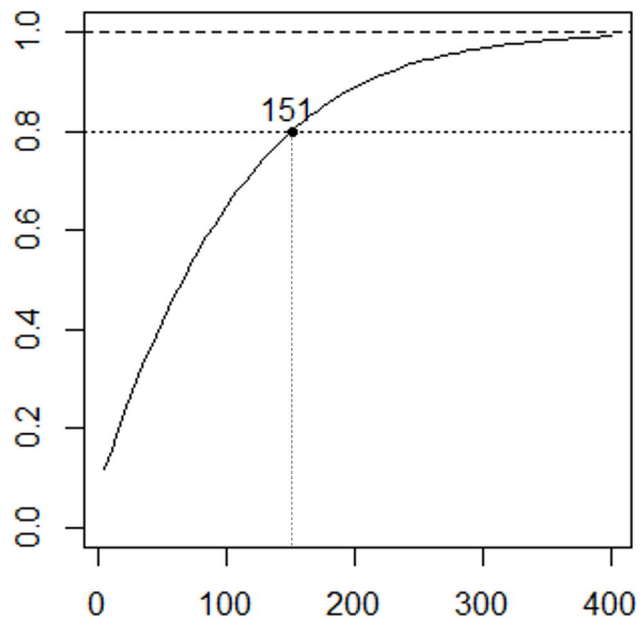
Month

Power



Sample size per arm

Power



Sample size per arm

