# Introduction: Adaptive trials & sample size re-calculation

**Marc Vandemeulebroecke**
**Part of the BBS training:** *Advanced group-sequential and adaptive confirmatory clinical trial designs*
**Basel, 13 Sep 2022**

# Learning objectives

Participants should understand:

- ## What are adaptive clinical trials
  - Major subtypes, distinctions, definitions

- ## Essential statistical methodology of adaptive trials
  - p-value combinations
  - conditional error functions
  - CRP principle
  - Conditional power and sample size adjustment

# What are adaptive designs?

- There are many different types of adaptive trial designs
  - (Group sequential designs, early stopping)
  - Adaptive randomization
  - Adaptive dose escalation
  - Adaptive dose finding
  - Sample size re-estimation
  - Treatment arm selection
  - Enrichment designs
  - ...

- Various «schools» of adaptive designs have developed in parallel, depending on the application area

# What are adaptive designs?

- Key distinctions:        Our focus:
  - Exploratory or confirmatory?    → Confirmatory
  - Adaptations of which trial features?    → Any
  - Using unblinded data?    → Yes*
  - Predetermined adaptations or ad-hoc?    → Both
  - Based on interim data or external information?   → Both

- Excluded here:
  - Blinded design modifications (e.g. blinded sample size re-estimation; generally not controversial)
  - Bayesian designs (frequent in early development phases)
  - Response-adaptive randomization

- Our focus:
  - Frequentist confirmatory adaptive designs

* But possibly sponsor-blind                          Vandemeulebroecke (2008)

# Some definitions of adaptive designs

- ## Dragalin (PhRMA), 2006:

  - A multistage study design that uses accumulating data to decide how to modify aspects of the study without undermining the validity and integrity of the trial. [...] preplanning, as much as possible, based on intended adaptations.

- ## FDA draft guidance, 2010:

  - A study that includes a prospectively planned opportunity for modification of one or more specified aspects of the study design and hypotheses based on analysis of data (usually interim data) from subjects in the study.

- ## EMA reflection paper, 2007:

  - A study is called 'adaptive' if statistical methodology allows the modification of a design element [...] at an interim analysis with full control of the type I error.

- ## FDA guidance, 2019:

  - A clinical trial design that allows for prospectively planned modifications to one or more aspects of the design based on accumulating data from subjects in the trial.

Dragalin (2006); FDA (2010, 2018); EMA (2007)

# Why adaptive designs

- In the 1980's, group sequential designs were introduced and grew popular. They provided a rigorous theory for early stopping but no other adaptations.

- In practice, however, adaptations of running trials were sometimes needed and done. Their impact on the inference was unclear and often ignored.

- **Quiz**: What is the maximal Type I error for a two-stage group-sequ. test with nominal level 5%* if $n_2$ is chosen in light of the observed first stage effect?

  - 5%?          8.2%?          11.5%?

* Note: I use one-sided tests and p-values unless otherwise specified.     Proschan, Hunsberger (1995)
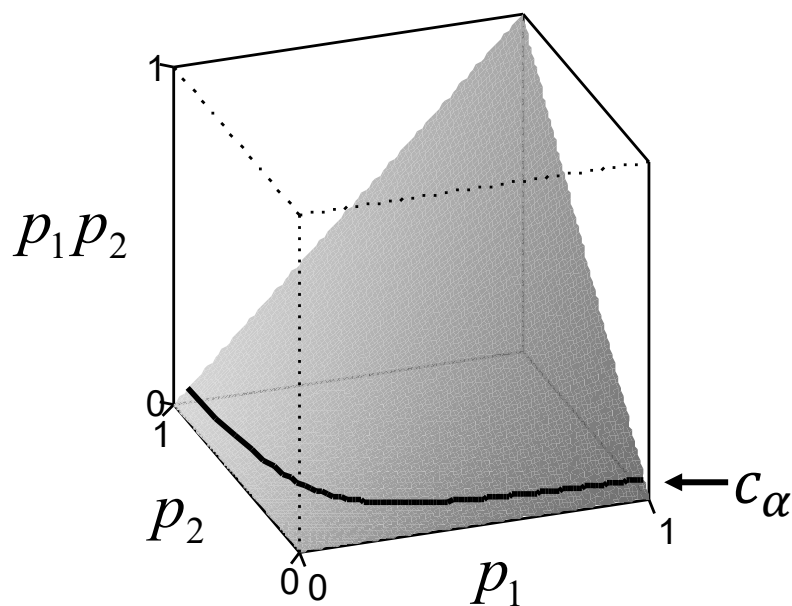
# Ignition: Bauer (1989)

- ## Idea borrowed from meta-analysis (MA):
  - MA combines the inference from separate *trials*
  - Now: <span style="color:red">combine the inference from separate stages of one trial</span>
  - This also allows adapting the second stage based on the first

- ## Method as well:
  - Take the product of the p-values from both trial stages
  - If $p_1 p_2$ is «too small» then reject $H_0$.
  - **Quiz**: What is «too small»?
    - Hint: How are $p_1$ and $p_2$ distributed under the null hypothesis?

Bauer (1989)

# Fisher's product test

- $p_1, p_2 \sim_{H_0} U[0,1]$ iid

- $-2\ln(p_1), -2\ln(p_2) \sim_{H_0} \chi^2_2$ iid

- $-2(\ln(p_1) + \ln(p_2)) \sim_{H_0} \chi^2_4$

- Rejecting $H_0$ when $-2(\ln(p_1) + \ln(p_2)) \geq \chi^2_{4,1-\alpha}$ is a level $\alpha$ test

- Equivalently, rejecting $H_0$ when
$p_1 p_2 \leq c_\alpha = \exp\left(-\frac{1}{2}\chi^2_{4,1-\alpha}\right)$
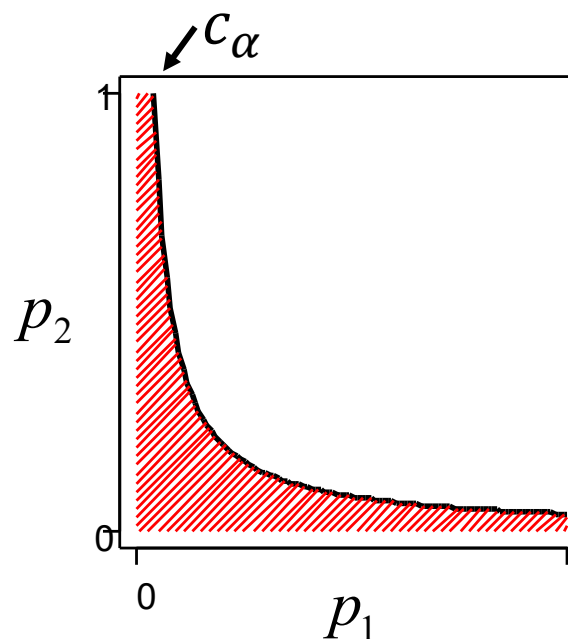
Fisher (1932)

# Let's look at it geometrically

- p-value combination
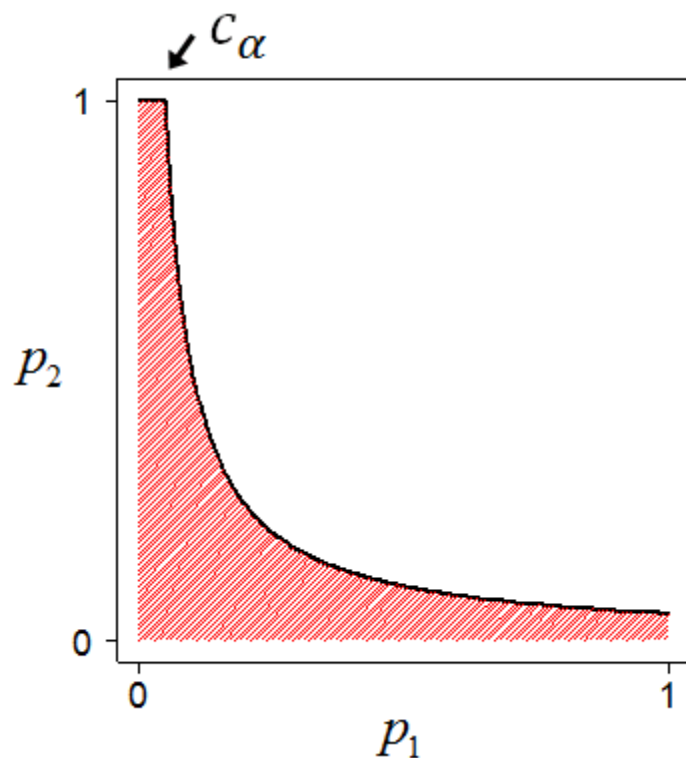


Reject if $p_1 p_2 \leq c_\alpha$

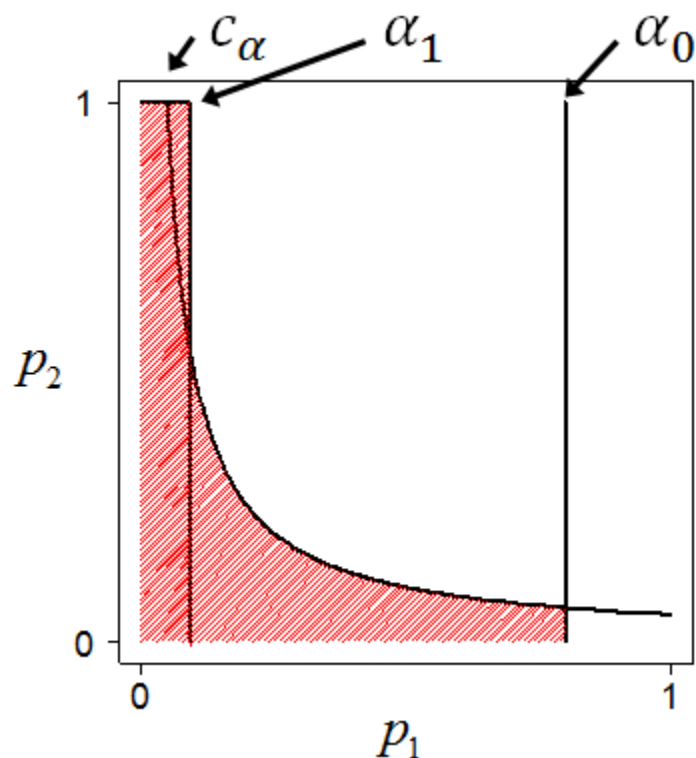- Projection onto the plane



Reject if $p_2 \leq c_\alpha / p_1$

- **Quiz**: How large is the red area?
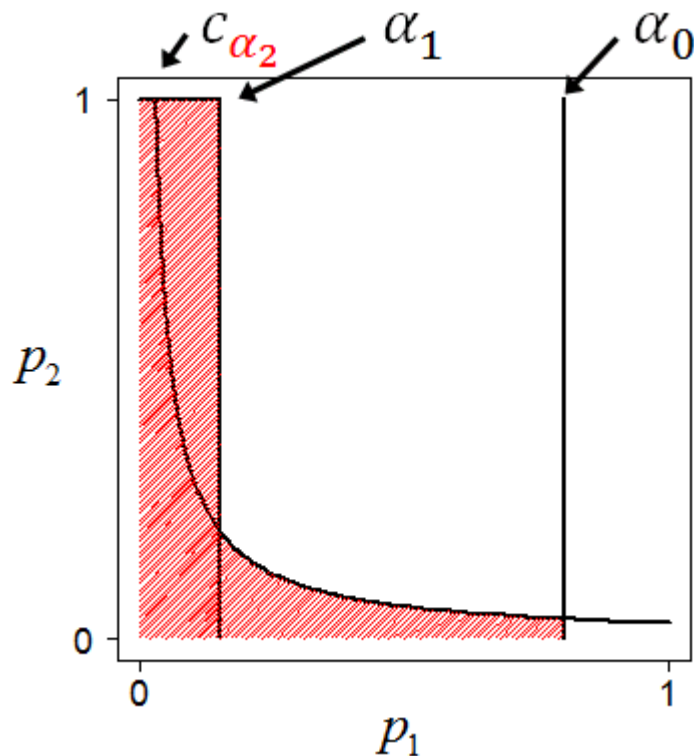
# The conditional error function



- Area of rejection region:
  $$\int_0^{c_\alpha} 1 \, dp_1 + \int_{c_\alpha}^1 c_\alpha/p_1 \, dp_1 = c_\alpha - c_\alpha \ln(c_\alpha)$$

- But we know this must be $\alpha$!
  - As $p_1, p_2 \sim_{H_0} U[0,1]$ iid, **areas correspond to probabilities**.
  - The rejection region has proba $\alpha$.

- This level curve **defines** a level $\alpha$ test of $H_0$. It is called a conditional error function (c.e.f.).

- Every p-value combination defines a family of c.e.f.'s that fills the unit square, and vice versa.

Proschan, Hunsberger (1995); Posch, Bauer (1999); Wassmer (1999); Vandemeulebroecke (2006).

# Early stopping



- Impose bounds $\alpha_1$ and $\alpha_0$
  - Assume $c_\alpha \leq \alpha_1 < \alpha_0$
  - $p_1 \leq \alpha_1 \to$ stop for efficacy
  - $p_1 > \alpha_0 \to$ stop for futility
  - Otherwise, perform second stage and reject $H_0$ if $p_2 \leq c_\alpha / p_1$

- Red area must remain $\alpha$
  - $\alpha_1 + c_\alpha\big(\ln(\alpha_0) - \ln(\alpha_1)\big) = \alpha$
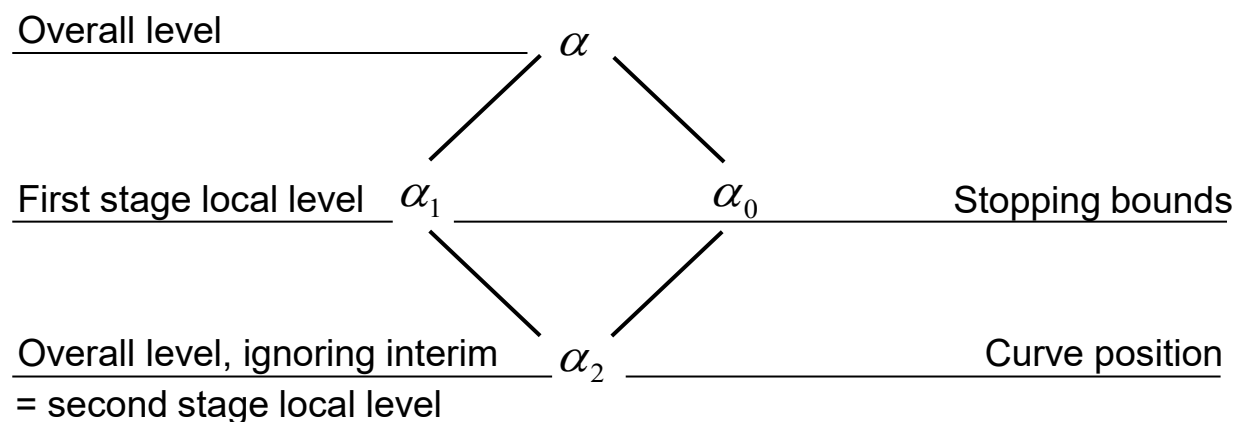
Bauer, Köhne (1994)

# Change height of curve



- Reject after second stage if $p_2 \leq c_{\alpha_2}/p_1$
  - This uses a different c.e.f. of the same family
  - The final test is performed at the local level $\alpha_2$

- Red area must remain $\alpha$
  $$\alpha_1 + c_{\alpha_2}\big(\ln(\alpha_0) - \ln(\alpha_1)\big) = \alpha$$

# The «alpha calculus»

- Four parameters are interdependent

Overall level — $\alpha$

First stage local level $\alpha_1$ — $\alpha_0$ — Stopping bounds

Overall level, ignoring interim $\alpha_2$ — Curve position
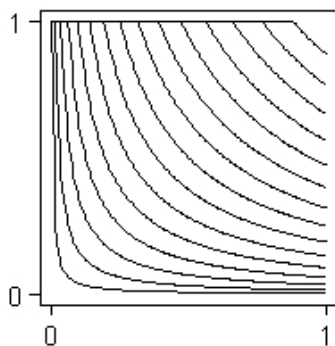= second stage local level

- Level condition: $\alpha_1 + c_{\alpha_2}\big(\ln(\alpha_0) - \ln(\alpha_1)\big) = \alpha$
- **Quiz**:
  - How would you specify a futility stop when control looks better?
  - How would you specify a «Pocock-type» test?
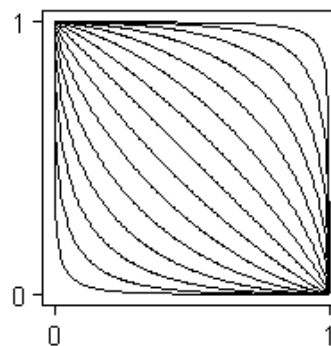
Vandemeulebroecke (2006)

# Inverse normal method & more

- Another natural way to combine p-values:

$$\frac{1}{\sqrt{2}}\left(\Phi^{-1}(1-p_1)+\Phi^{-1}(1-p_2)\right) \sim_{H_0} N(0,1)$$

- Same mechanism, with a different family of c.e.f.'s



Product test             Inverse normal method (INM)

- In principle, any such family defines an adaptive test by this mechanism

  - In practice, mainly these two are used. And out of these, mostly INM.

Lehmacher, Wassmer (1999); Vandemeulebroecke (2006)

# How do trial adaptations fit into this?



- This height is the Type I error probability given the first stage data

- We could now change the second stage into any design that respects this level

- The resulting overall procedure remains a level $\alpha$ test

# Why does this work?

- For continuously distributed test statistics based on separate stages, $p_1, p_2$ will generally be iid $U[0,1]$ under $H_0$ even if the second stage is modified based on the interim analysis

- More generally, it still works if $p_1, p_2$ are only «p-clud»
  - $P_{H_0}(p_1 \leq u) \leq u$ and $P_{H_0}(p_2 \leq u \mid p_1) \leq u$ for all $u \in [0,1]$

- For more details on probabilistic foundations, see Brannath et al. 2012.

Brannath et al. (2002); Brannath et al. (2012)

# Conditional Rejection Principle (CRP)

- Start with a (classical) level $\alpha$ test

- At an IA, review the data and possibly external information

- No reason to adapt → Continue as planned

- Reason to adapt

  → Compute cond. Type I error of the pre-defined design:

  $$P_{H_0}(\text{reject } H_0 | \text{ interim data})$$

  And choose (based on all info) a new design at **this** level to finish the trial

- This is a level $\alpha$ test, and the IA need not be preplanned

# Conditional Rejection Principle (CRP)

- ## How could that new second-stage design look like?

  - Increase the remaining sample size (e.g., to achieve a desired conditional power → *see later*)

    - Note: Health authorities view sample size **reductions** more critically

  - Replace the second stage by another two-stage design → multistage designs by «recursive combination»

  - ...and more

- ## Caveat

  - Adaptations must not jeopardize **interpretability** of results or **credibility** of the trial!

Müller, Schäfer (2004); Brannath et al. (2002)

# Relation: Group sequ. ↔ adaptive

- Group sequential designs follow a <span style="color:red">cumulative philosophy</span>: their test statstics are cumulative

- Adaptive designs follow a <span style="color:red">stagewise philosophy</span>: they use stagewise inferences (test statistics, p-values)
  - However, the decision rules of adaptive designs <span style="color:red">combine</span> the stagewise inferences – so overall **they do provide cumulative inference**
  - For example, Fisher's product test rejects $H_0$ if $p_2 \leq c_{\alpha_2}/p_1$

- The INM in particular reduces **exactly** to the group sequential test **if** no adaptations are done*. The test statistics, critical values and decision rules are identical.

➔ *Next slide*

* ...and given some distributional assumptions.

# Relation: Group sequ. ↔ adaptive

- Test active vs. placebo with normally distr. endpoint

- Group sequential: $X_{ki} \sim N(\mu, \sigma^2)$ iid, $Y_{ki} \sim N(\nu, \sigma^2)$ iid
  - $k = 1,2$ (stage); $i = 1, \ldots, n_k$; $\sigma^2$ known
  - $n = n_1 + n_2$ total sample size per arm; $n_1 = n_2$ without loss of generality

- The Z-test:

  - Overall: $Z = \sqrt{\frac{n}{2}} \frac{\bar{X} - \bar{Y}}{\sigma} \sim_{H_0} N(0,1)$

  - Per stage: $Z_k = \sqrt{\frac{n}{4}} \frac{\bar{X}_k - \bar{Y}_k}{\sigma} \sim_{H_0} N(0,1); \; p_k = 1 - \Phi(Z_k)$

  - Group sequential: Using $Z_1$ and $Z$

  - Inverse normal method:

    Combining $p_1$ and $p_2$ to $\frac{1}{\sqrt{2}}\left(\Phi^{-1}(1 - p_1) + \Phi^{-1}(1 - p_2)\right) = \frac{1}{\sqrt{2}}(Z_1 + Z_2) = Z$

# Relation: Group sequ. ↔ adaptive

- The INM therefore generalizes the group sequ. test
  - Standard group sequential software can be used

- It is easily communicated with commonly used (Z-) statistics

- It is also the uniformly most powerful test if no adaptations are done

➢ All this is why the INM is often the <span style="color:red">preferred method</span>

# Weights

- More general version of the INM
  - Combine stagewise statistics using $w_1 Z_1 + w_2 Z_2$ instead of $\frac{1}{\sqrt{2}}(Z_1 + Z_2)$, with weights $w_k$
  - Weights can be freely chosen under the constraint $w_1^2 + w_2^2 = 1$
  - But they must be prespecified and remain fixed regardless of adaptations
    - Otherwise, the type I error may be inflated
  - Natural choice: $w_k = \sqrt{\dfrac{n_k}{n_1 + n_2}}$
  - Then all patients carry equal weight, and again we have $w_1 Z_1 + w_2 Z_2 = Z$
  - The case $n_1 = n_2$ above was a special case of this

# Efficiency vs. flexibility

- **Quiz:** What happens to the INM if we change the remaining sample size at the IA?
  - Not all patients carry equal weight → <span style="color:red">inefficient</span>

- A curious debate
  - Tsiatis, Mehta (2003): "On the **inefficiency** of the adaptive design [...]"
  - Brannath et al. (2006): "On the **efficiency** of adaptive designs [...]"

- What do **you** think?

- In my view, trialists should weigh **efficiency** (power) against **flexibility** (adaptation)

# Conditional power

- The conditional power is the power of the trial (at some alternative), given interim data

- Let's look at the inverse normal method

- Situation as before: $X_{ki} \sim N(\mu, \sigma^2)$ iid, $Y_{ki} \sim N(\nu, \sigma^2)$ iid
  - $k = 1,2$ (stage); $i = 1, \dots, n_k$
  - $n = n_1 + n_2$ total sample size per arm
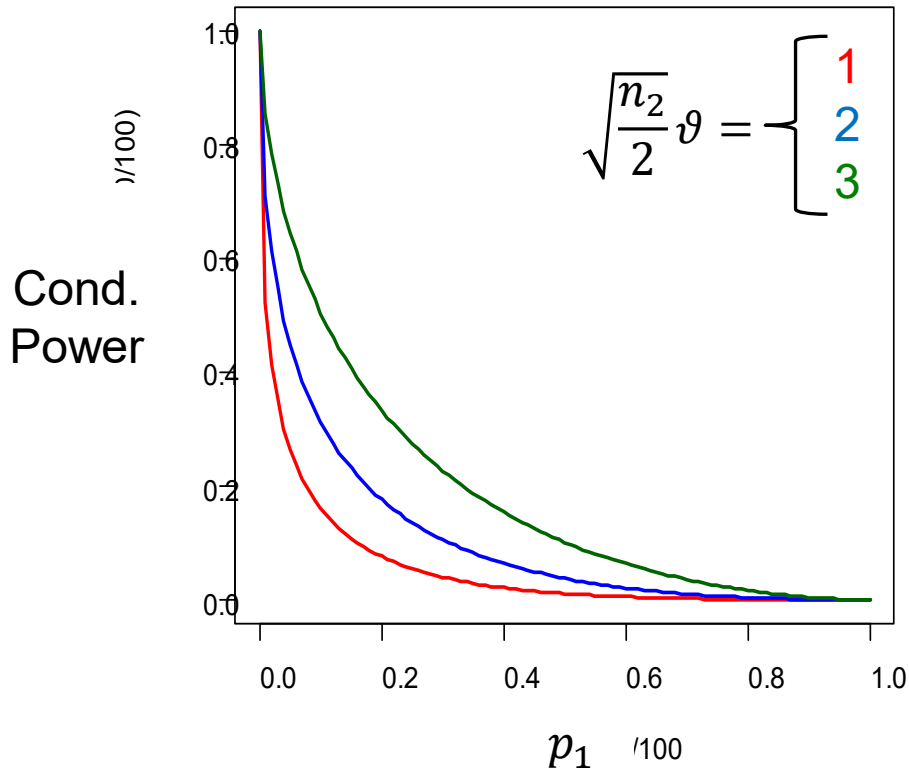  - Denote $\vartheta = \frac{\mu - \nu}{\sigma}$

→ *Next slide*

Proschan, Hunsberger (1995)

# Conditional power

- $CP_\vartheta = P_\vartheta \left( \frac{1}{\sqrt{2}} \left( \Phi^{-1}(1-p_1) + \Phi^{-1}(1-p_2) \right) \geq u_\alpha \,\middle|\, p_1 \right)$

$$= P_\vartheta \left( \frac{1}{\sqrt{2}} (Z_1 + Z_2) \geq u_\alpha \,\middle|\, Z_1 = z_1 \right)$$

$$= P_\vartheta \left( Z_2 \geq \sqrt{2} u_\alpha - z_1 \right)$$

$$= P_\vartheta \left( Z_2 - \sqrt{\frac{n_2}{2}} \vartheta \geq \sqrt{2} u_\alpha - z_1 - \sqrt{\frac{n_2}{2}} \vartheta \right)$$

$$= 1 - \Phi \left( \sqrt{2} u_\alpha - z_1 - \sqrt{\frac{n_2}{2}} \vartheta \right)$$

Here, $u_\alpha$ is the $(1-\alpha)$-quantile of $N(0,1)$.

# Conditional power

- Properties



$$\sqrt{\frac{n_2}{2}}\vartheta = \begin{cases} 1 \\ 2 \\ 3 \end{cases}$$

- Conditional power
  - Increases with $n_2$
  - Increases with $\vartheta$
  - Decreases for increasing $p_1$

Cond. Power

$p_1$

Here, $\alpha = 0.025$

# Conditional power

- Common applications
  - Stopping for futility if $CP_\vartheta$ is «too small» (e.g. below 20%)
  - Adjusting the second stage size to achieve a desired $CP_\vartheta$ (e.g. 90%)

    In the example, solve $0.9 = 1 - \Phi\left(\sqrt{2}u_\alpha - z_1 - \sqrt{\frac{n_2}{2}}\vartheta\right)$ for $n_2$

    Conduct the second stage and perform the final inference as planned through the adaptive design
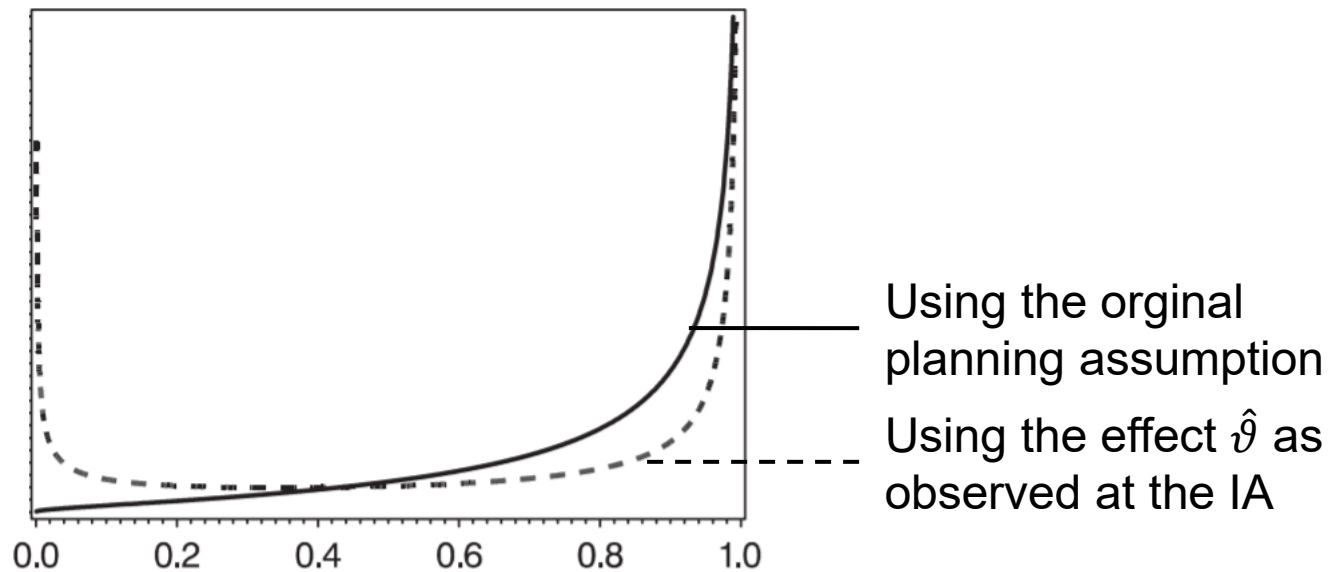
- **Quiz**: What $\vartheta$ would you use in $CP_\vartheta$?

# Conditional power

- Several options for $\vartheta$ in $CP_\vartheta$
  - The originally assumed effect size for sample size calculation (minimally clinically relevant effect – should not have changed!)
  - The effect size $\hat{\vartheta}$ as observed at the interim analysis (hoping that this comes closer to the «truth»)
    - **Caution**: Interim estimates such as $\hat{\vartheta}$ are notoriously volatile! → *Next slide*
  - Averaging across several choices
    - Weighted average of originally assumed and observed effect size
    - Integrating over some distribution for $\vartheta$ («predictive power»)

# Conditional power

- ## Using the interim effect estimate is risky
  - Because we rely **doubly** on little data: through $z_1$ and through $\hat{\vartheta}$
  - The density of $CP_\vartheta$ tends towards extremes if we use $\hat{\vartheta}$

Using the orginal planning assumption

Using the effect $\hat{\vartheta}$ as observed at the IA

Here, $n_1 = n_2$, $\alpha = 0.025$, with 80% desired power. Bauer, König (2006)

# The «Constrained Promising Zone» (CPZ) Approach

- A recent proposal for a more refined use of conditional power to re-calculate the sample size

  - Builds upon the previously proposed «Promising Zone» approach by Mehta and Pocock (2011) which had been shown to be (overly) conservative (Glimm 2012, Jennison and Turnbull 2015)

- Idea: *Boost the sample size within reasonable limits when the interim treatment effect appears «promising»*
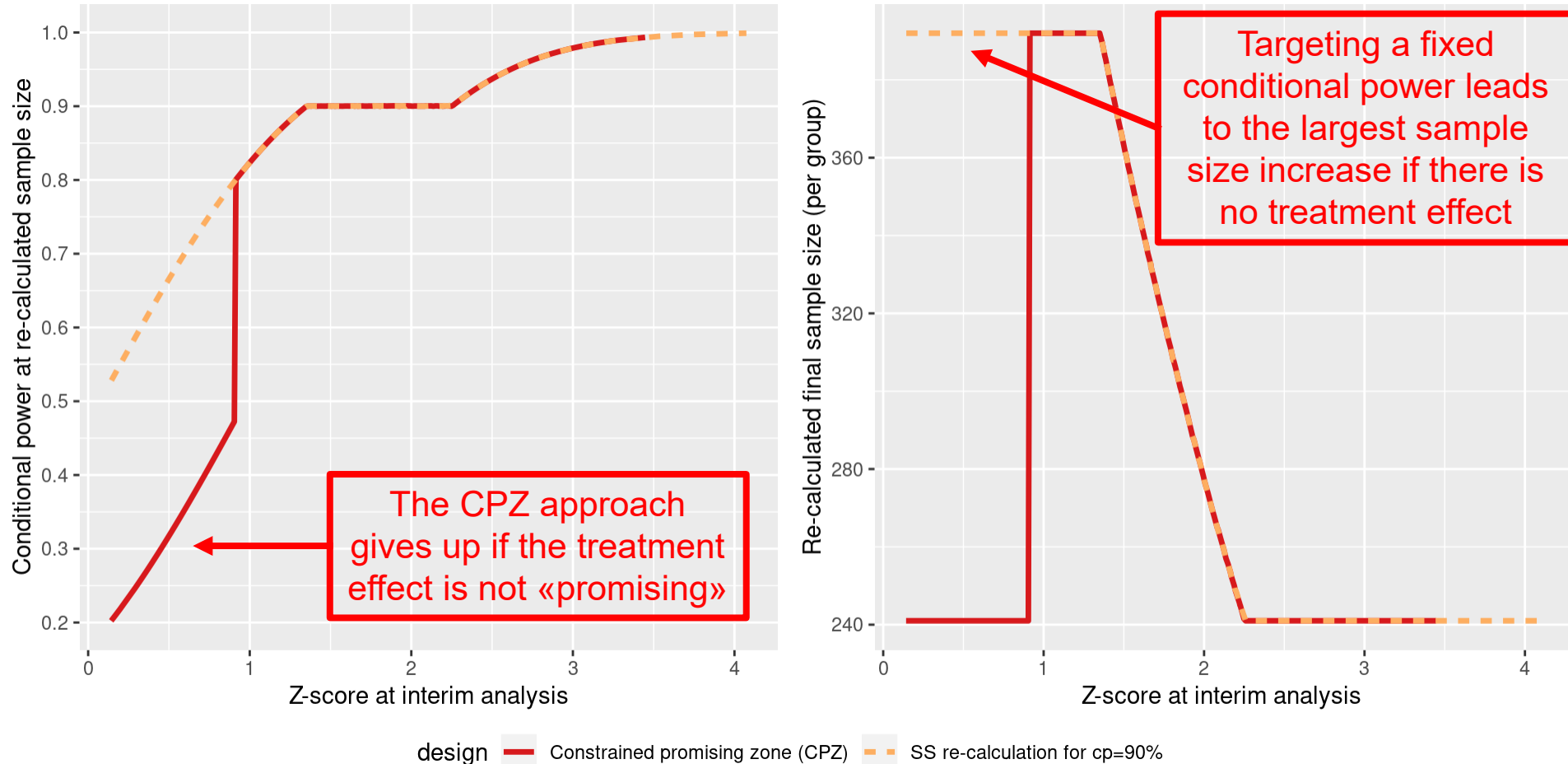
Hsiao et al., 2019. Credit to Marcel Wolbers and Kaspar Rufibach

# The «Constrained Promising Zone» (CPZ) Approach

- Concretely, pre-specify:
  - Impose limits to allowed total sample size per arm: $n_{min}, n_{max}$
  - Set smallest clinically meaningful effect size $\vartheta_{min}$, and smallest / largest desired conditional power at this point: $CP_{min}, CP_{max}$
  - Choose a combination test, e.g. INM with $w_1 = \sqrt{\frac{n_1}{n_{min}}}, w_2 = \sqrt{\frac{n_{min}-n_1}{n_{min}}}$

- Then re-calculate the sample size at the IA:
  - If $n^*$ exists between $n_{min}$ and $n_{max}$ such that $CP_{\vartheta_{min}}(z_1, n^*) = CP_{max}$, then set the total sample size (per arm) to $n^*$
  - Otherwise, if $CP_{\vartheta_{min}}(z_1, n_{max}) \geq CP_{min}$, then set it to $n_{max}$
  - Finally, otherwise, set it to $n_{min}$ because the IA is not «promising»

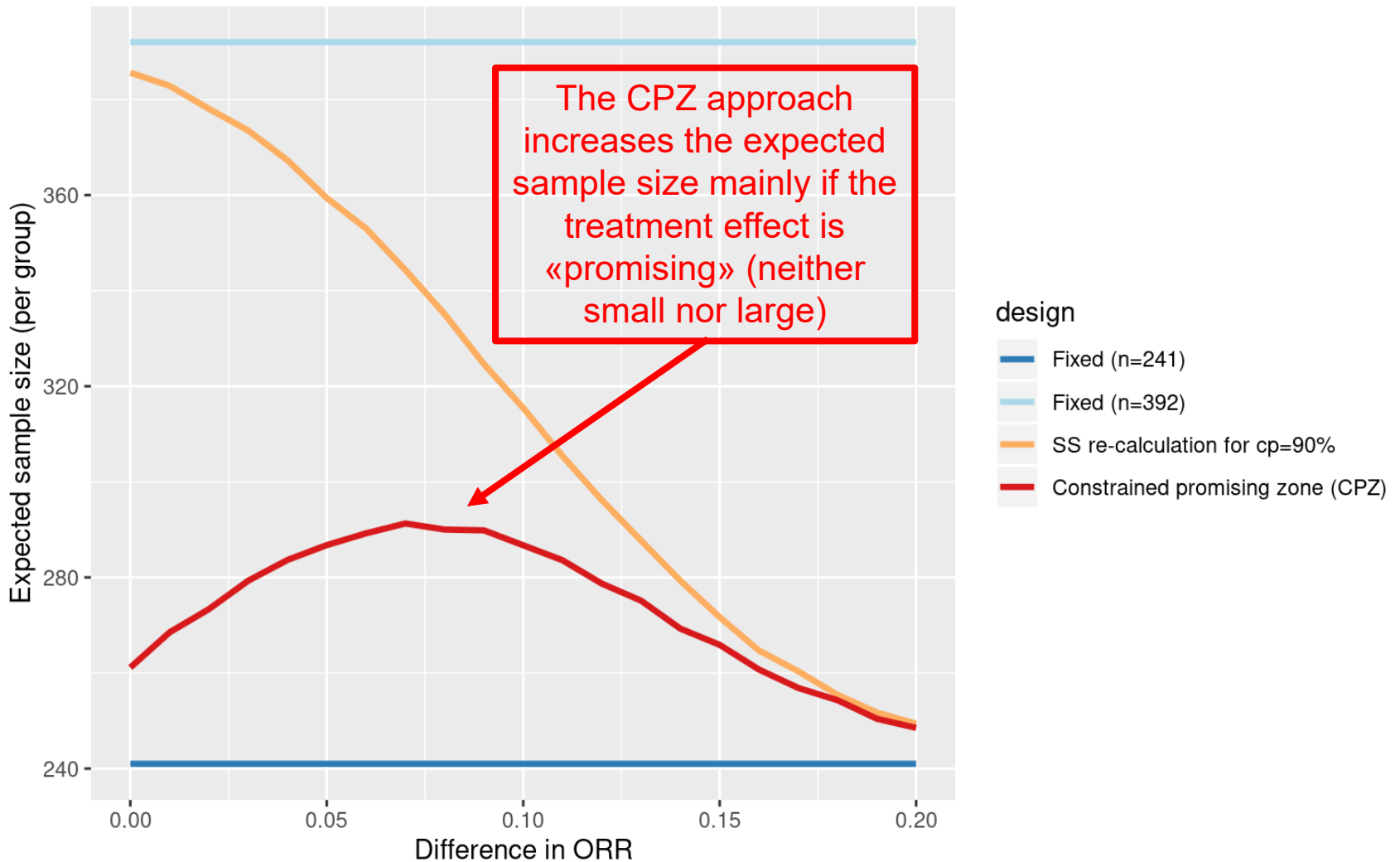Hsiao et al., 2019. Credit to Marcel Wolbers and Kaspar Rufibach

# The «Constrained Promising Zone» (CPZ) Approach - Example

- 1:1 randomization with Overall Response Rate (ORR) as primary endpoint
    - ORR=20% on Control; Drug increases this by 10-13%
    - 2.5% significance level (one-sided)
    - $n_1 = 120$
    - $n_{min} = 241$, $n_{max} = 392$ (90% power for Δ=13% and Δ=10%, resp.)

- Compare two approaches
    - Sample size increase for a conditional power of 90% (if true Δ=10%)
    - CPZ design with $CP_{min} = 80\%$, $CP_{max} = 90\%$

- Corresponding R-code is in this vignette
    - [Simulation of a Trial with a Binary Endpoint and Unblinded Sample Size Re-Calculation with rpact](#)
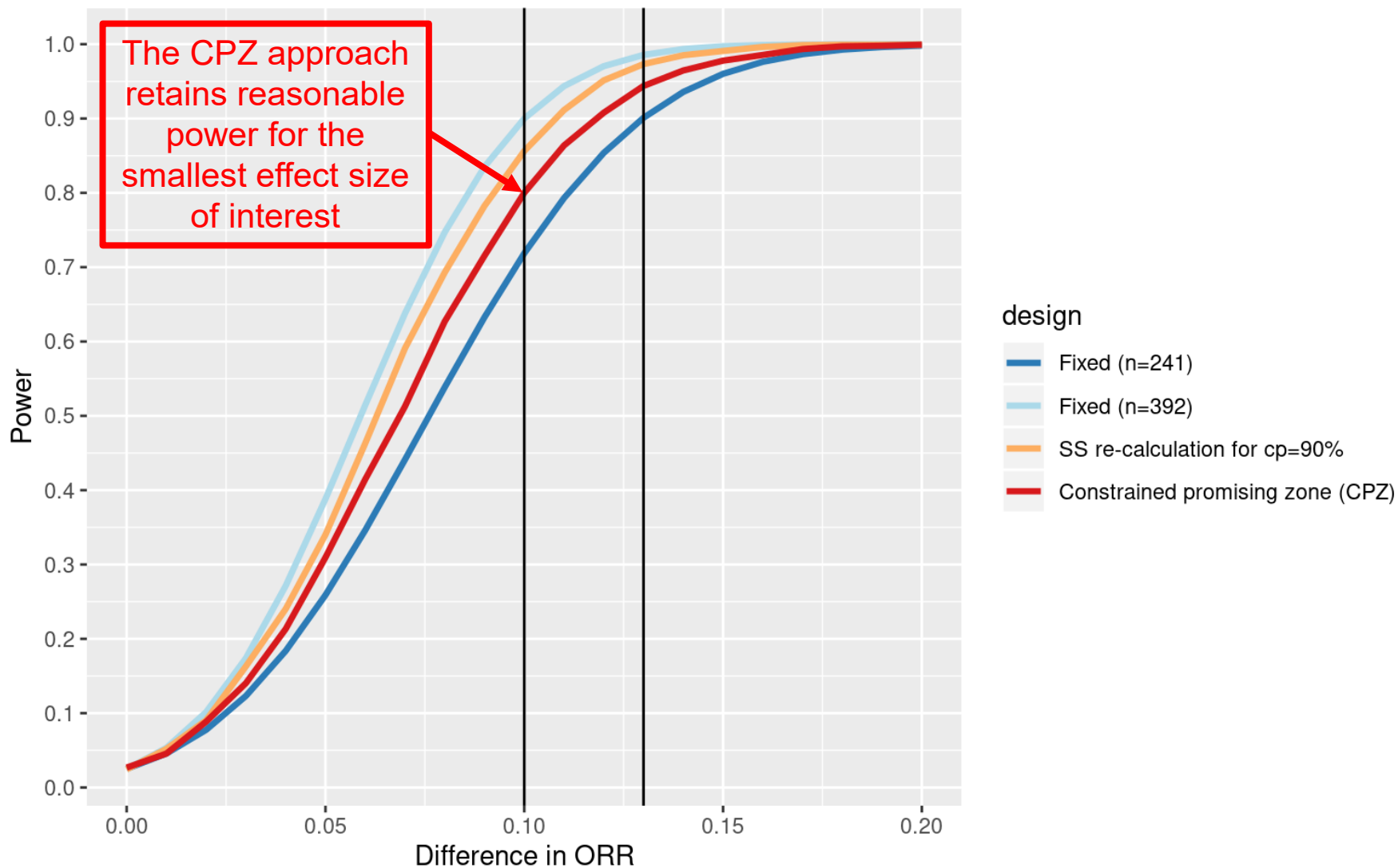
Hsiao et al., 2019. Credit to Marcel Wolbers and Kaspar Rufibach

# Cond. power and total sample size depending on the interim Z-score

Hsiao et al., 2019. Credit to Marcel Wolbers and Kaspar Rufibach

# Expected sample size depending on effect size

Hsiao et al., 2019. Credit to Marcel Wolbers and Kaspar Rufibach

# Power depending on effect size



The CPZ approach retains reasonable power for the smallest effect size of interest

design

Fixed (n=241)

Fixed (n=392)

SS re-calculation for cp=90%

Constrained promising zone (CPZ)

Hsiao et al., 2019. Credit to Marcel Wolbers and Kaspar Rufibach

# Regulatory guidance on unblinded sample size adaptation

**EMA guidance on adaptive designs 2007**

- The option to reassess sample size in an ongoing trial should not be seen as a substitute for careful planning. The relevance of a particular size of treatment effect should be discussed at the planning stage of the trial and not deferred to the point where interim results are already available.

- Whenever possible, methods for blinded sample size reassessment that properly control the type I error should be used […]. In cases where sample size needs to be reassessed based on unblinded data, sufficient justification should be made.

**FDA adaptive designs guidance 2019**

- [such designs] might be used when there is considerable uncertainty about the true treatment effect size.

- […] to appropriately control the Type I error […and] prospective planning […of] the statistical hypothesis testing method […and] the rule governing the sample size modification.

- […] additional challenges in maintaining trial integrity […]

# Our recommendations for unblinded sample size adaptation

- Approach is accepted by health authorities, but more justification is needed than for blinded sample size adaptation

- Main application: Considerable uncertainty about the size of the treatment effect and reluctance to fund a group-sequential trial powered to the smallest clinically relevant effect size

  - *«Start small and invest more resources if results look promising»*

- Extensive literature of such designs versus «more efficient» group-sequential designs

  - E.g., Liu et al, 2018: «*...under reasonable decision rules for increasing sample size […] there is little or no loss of efficiency for the adaptive designs in terms of unconditional power. The two approaches, however, have very different conditional power profiles.*»

- Extensive clinical trials simulations and comparisons to group-sequential designs are highly recommended

  - Can also help to explore potential bias in estimation
  - `rpact` can produce median unbiased estimators and other inference adjusted for the adaptive design

# Final thoughts on adaptive designs

- Allowance for adaptations of the trial design without inflating type I error
  - Adaptations should be pre-planned in most circumstances
  - ...but can be occasionally be used to react to unforeseen circumstances

- Can be extended to multi-arm and enrichment designs (covered later)

- Adaptive designs are more complicated than fixed or group-sequential designs in terms of trial planning, logistics, and regulatory requirements to ensure trial integrity and avoid operational bias

- Two attitudes:
  - The social event trial: «Let's come together, let's see and then adapt until significance» (Koch 2006)
  - Much better: «A multistage study design that uses accumulating data to decide how to modify aspects of the study without undermining the validity and integrity of the trial.» (Dragalin 2006)

- **Adaptive designs are not a remedy for sloppy planning!**

# References (1)

**Reference book**

- Wassmer, Brannath (2016): Group Sequential and Confirmatory Adaptive Designs. Springer

**Reviews**

- Todd (2007): A 25-year review of sequential methodology in clinical studies. Statistics in Medicine 26, 237–252
- Vandemeulebroecke (2008): Group Sequential and Adaptive Designs – A Review of Basic Concepts and Points of Discussion. Biometrical Journal 50, 541–557
- Bauer et al. (2016): Twenty-five years of confirmatory adaptive designs: opportunities and pitfalls. Statistics in Medicine 35, 325–347

**Regulatory guidance**

- EMA (2007). Reflection paper on methodological issues in confirmatory clinical trials planned with an adaptive design. European Medicines Agency, CHMP/EWP/2459/02
- FDA (2010): Draft guidance for industry. Adaptive design clinical trials for drugs and biologics. Rockville, MD
- FDA (2019): Guidance for industry. Adaptive designs for clinical trials of drugs and biologics. Silver Spring, MD
- Brannath et al. (2010): Comments on the Draft Guidance on "Adaptive Design Clinical Trials for Drugs and Biologics" of the U.S. Food and Drug Administration. Journal of Biopharmaceutical Statistics 20, 1125–1131

**Discussion & interpretation**

- Dragalin (2006): Adaptive Designs: terminology and classification. Drug Information Journal 40, 425–435
- Koch (2006): Confirmatory clinical trials with an adaptive design. Biometrical Journal 48, 574–585. Rejoinder pp. 616–622
- Burman, Sonesson (2006): Are flexible designs sound? (with discussion). Biometrics 62, 664–683
- Vandemeulebroecke (2008): *see above*

# References (2)

**Methodological contributions**

- Fisher (1932): Statistical methods for research workers. Oliver & Boyd, London
- Marcus, Peritz, Gabriel (1976): On closed testing procedures with special reference to ordered analysis of variance. Biometrika 63, 655–660.
- Bauer (1989): Multistage testing with adaptive designs (with discussion). Biometrie und Informatik in Medizin und Biologie 20, 130–148
- Bauer, Köhne (1994): Evaluation of experiments with adaptive interim analyses. Biometrics 50: 1029-1041. Correction in Biometrics 52 (1996): 380
- Bauer, Röhmel (1995): An adaptive method for establishing a dose-response relationship. Statistics in Medicine 14, 1595–1607
- Proschan, Hunsberger (1995): Designed extension of studies based on conditional power. Biometrics 51, 1315–1324.
- Posch, Bauer (1999): Adaptive two stage designs and the conditional error function. Biometrical Journal 41: 689–696.
- Wassmer (1999): Statistische Testverfahren für gruppensequentielle und adaptive Pläne in klinischen Studien. Verlag Alexander Mönch, Köln
- Lehmacher, Wassmer (1999): *Adaptive sample size calculations in group sequential trials.* Biometrics 55: 1286–1290
- Brannath et al. (2002): *Recursive combination tests.* JASA 97 (457): 236–244
- Müller, Schäfer (2004): A general statistical principle for changing a design any time during the course of a trial. Statistics in Medicine 23, 2497–2508
- Bauer, König (2006): The reassessment of trial perspectives from interim data – a critical view. Statistics in Medicine 25, 23–36
- Vandemeulebroecke (2006): An investigation of two-stage tests. Statistica Sinica 16, 933–951
- Brannath, Gutjahr, Bauer (2012): Probabilistic foundation of confirmatory adaptive designs. JASA 107: 824-832
- Xsiao, Liu, Mehta (2019): Optimal promising zone designs. Biometrical Journal 61: 1175-1186

# References (3)

**Clinical trial examples**

- Vandemeulebroecke, Bornkamp, Bretz, Pinheiro (2010): Adaptive dose-ranging studies. Chapter 11 in: Handbook of Adaptive Designs for Pharmaceutical and Clinical Development. Chapman & Hall
- Barnes et al. (2010): Integrating indacaterol dose selection in a clinical study in COPD using an adaptive seamless design. Pulmonary Pharmacology & Therapeutics 23: 165–171
- Schmoll et al. (2012): Cediranib with mFOLFOX6 versus bevacizumab with mFOLFOX6 as first-line treatment for patients with advanced colorectal cancer: A double-blind, randomized phase III study (HORIZON III). Journal of Clinical Oncology 30, 3588–3595
- Cuffe, Lawrence, Stone, Vandemeulebroecke (2014): When is a seamless study desirable? Case studies from different pharmaceutical sponsors. Pharmaceutical Statistics 13, 229–237

**Software reviews**

- Wassmer, Vandemeulebroecke (2006): A brief review on software developments for group sequential and adaptive designs. Biometrical Journal 48: 732–737
- Tymofyeyev (2014): A Review of Available Software and Capabilities for Adaptive Designs. Chapter in: Practical Considerations for Adaptive Trial Design and Implementation. Springer

**Commercial software**

- ADDPLAN: http://www.iconplc.com/innovation/addplan/
- EastAdapt and EastSurv: http://www.cytel.com/software/east

**R packages**

- adaptTest: https://cran.r-project.org/web/packages/adaptTest/index.html
- AGSDest: https://cran.r-project.org/web/packages/AGSDest/index.html
- asd: https://cran.r-project.org/web/packages/asd/index.html
- rpact: https://www.rpact.com/