

# Generalized Pairwise Comparisons: a statistical method for patient-centric medicine

*Marc Buyse, ScD  
IDDI and U Hasselt, Belgium*

*QuanTIM Webinar*

*17 February 2023*

# Agenda

- Theory
  - Generalized Pairwise Comparisons
  - Net Treatment Benefit
- Applications
  - Augmenting power *and* clinical relevance
  - Benefit / risk analyses
  - Multiple testing procedures
- Conclusions

# Theory

# Wilcoxon rank-sum test

TREATMENT  
GROUP (T)



CONTROL  
GROUP (C)



$X_1$

$X_2$

...

$X_n$

$Y_1$

$Y_2$

...

$Y_m$

# Wilcoxon rank-sum test

TREATMENT  
GROUP (T)



$X_1$   
 $X_2$

...

$X_n$

CONTROL  
GROUP (C)



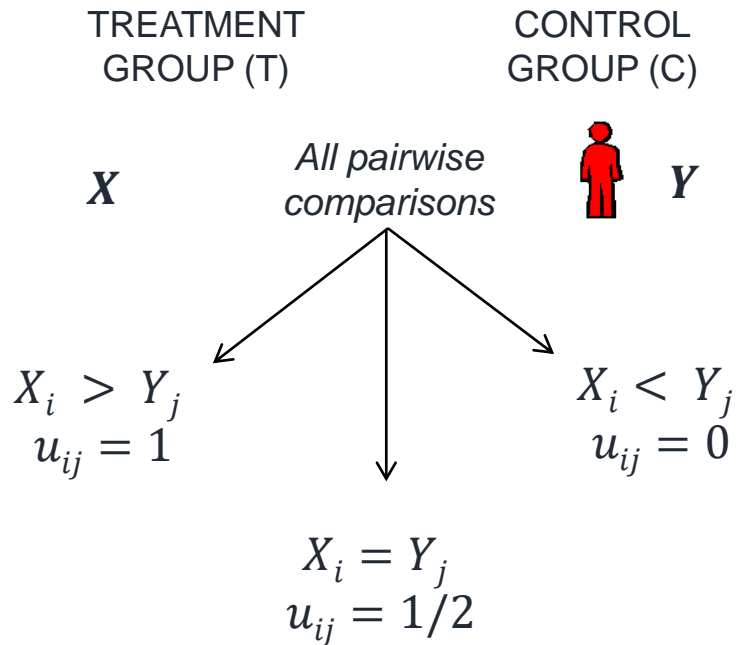
$Y_1$   
 $Y_2$

...

$Y_m$

1. Order the  $(n + m)$  elements of  $\mathbf{X} \cup \mathbf{Y}$
2. Let  $R_i$  be the rank order of the  $i^{\text{th}}$  element
3. For groups of tied values, assign a rank equal to the midpoint of the unadjusted ranks
4. Calculate  $U = \sum_{i=1}^n R_i$ , the sum of ranks of the elements of  $\mathbf{X}$
5. The statistic  $U$  has a known distribution under  $H_0$

# Mann-Whitney test



1. Perform pairwise comparisons between all elements of  $X$  and  $Y$

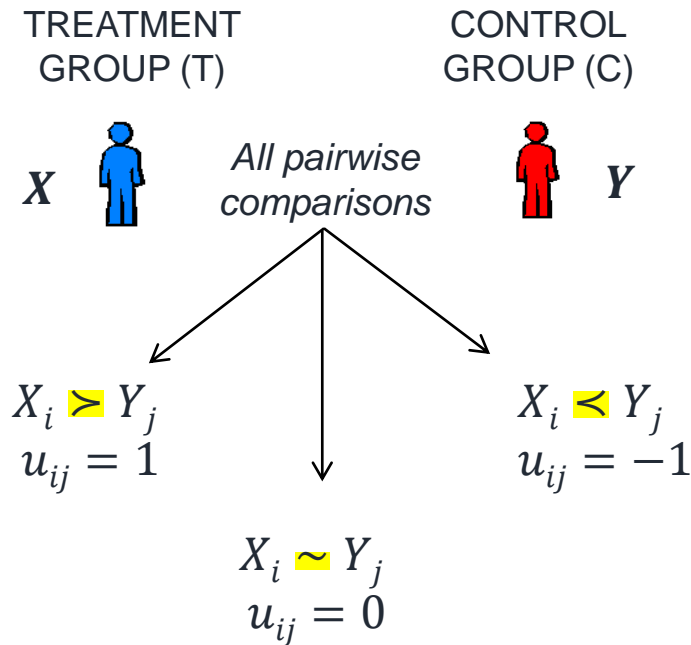
2. Calculate  $u_{ij} = \begin{cases} 1 & \text{if } X_i > Y_j \\ 0 & \text{if } X_i < Y_j \\ 1/2 & \text{if } X_i = Y_j \end{cases}$

3. The statistic

$$U = \frac{1}{m \cdot n} \sum_{i=1}^m \sum_{j=1}^n u_{ij}$$

has a known distribution under  $H_0$

# Generalized Pairwise Comparisons (GPC)



1. Perform pairwise comparisons between all elements of  $X$  and  $Y$

2. Calculate  $u_{ij} = \begin{cases} +1 & \text{if } X_i \succ Y_j \\ -1 & \text{if } X_i \preccurlyeq Y_j \\ 0 & \text{if } X_i \sim Y_j \end{cases}$

3. The statistic

$$U = \frac{1}{m \cdot n} \sum_{i=1}^m \sum_{j=1}^n u_{ij}$$

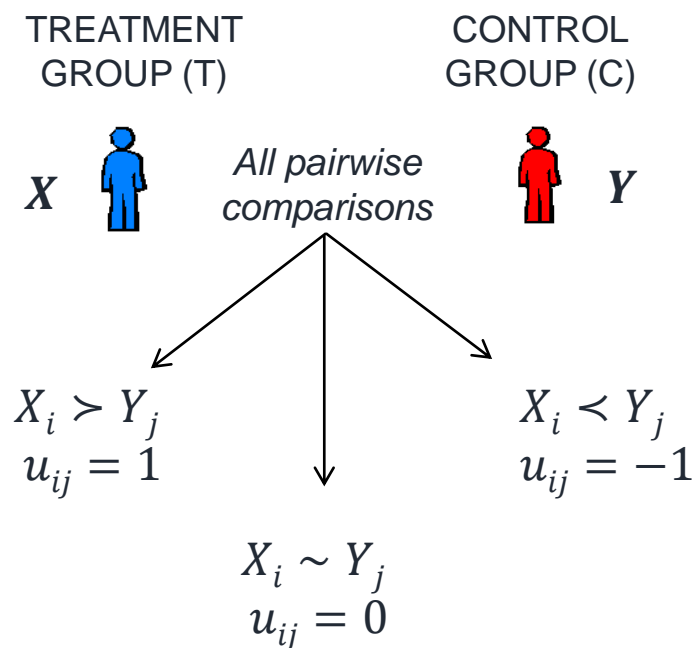
has a known distribution under  $H_0$

where  $\succ$  stands for “better” (win)

$\preccurlyeq$  stands for “worse” (loss)

$\sim$  stands for “similar” (tie) or “unclassified” (?)

# GPC – Outcome of any type

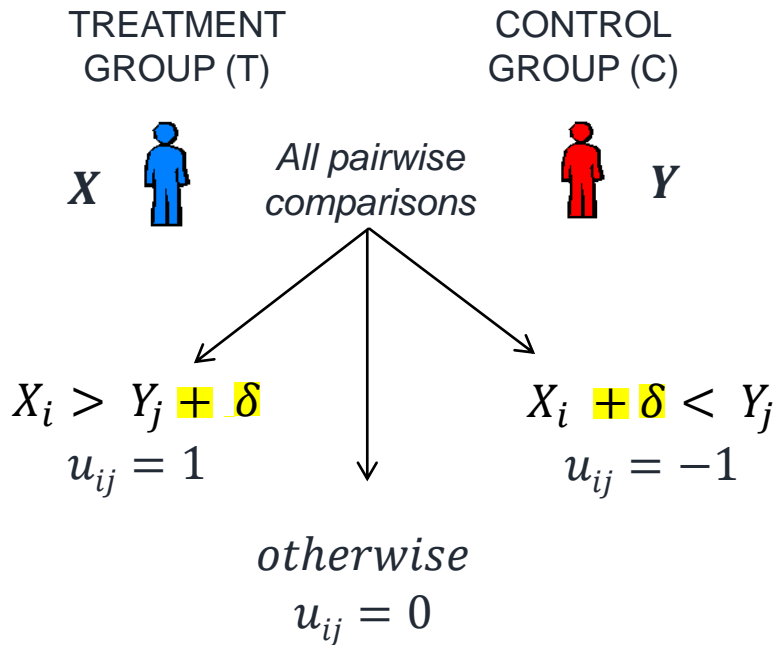


“ $X_i$  better than  $Y_j$ ” (wins):

- For ordered outcomes, with larger values preferable:  $X_i > Y_j$
- For binary outcomes, with 1 denoting success and 0 failure,  $X_i > Y_j$
- For time-to-event outcomes, with larger values preferable,  $X_i > Y_j$  unless  $Y_j$  censored
- For all outcome types, arbitrary definition



# GPC – clinical threshold



1. Perform pairwise comparisons between all elements of ordered outcomes  $X$  and  $Y$

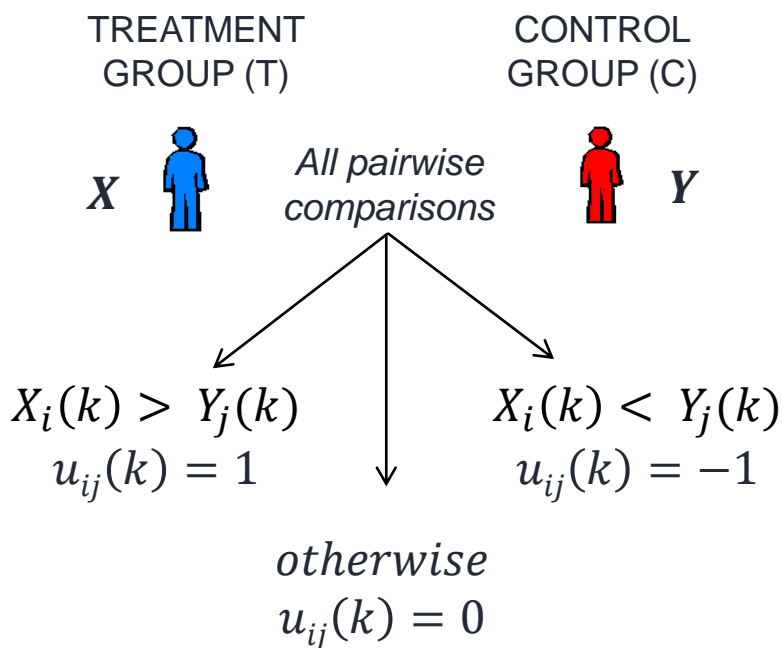
2. Calculate  $u_{ij} = \begin{cases} +1 & \text{if } X_i > Y_j + \delta \\ -1 & \text{if } X_i + \delta < Y_j \\ 0 & \text{otherwise} \end{cases}$

3. The statistic

$$U = \frac{1}{m \cdot n} \sum_{i=1}^m \sum_{j=1}^n u_{ij}$$

has a known distribution under  $H_0$

# GPC – multiple weighted outcomes



1. Perform pairwise comparisons between all elements of  $X$  and  $Y$

2. Calculate  $u_{ij}(k) = \begin{cases} +1 & \text{if } X_i(k) > Y_j(k) \\ -1 & \text{if } X_i(k) < Y_j(k) \\ 0 & \text{otherwise} \end{cases}$

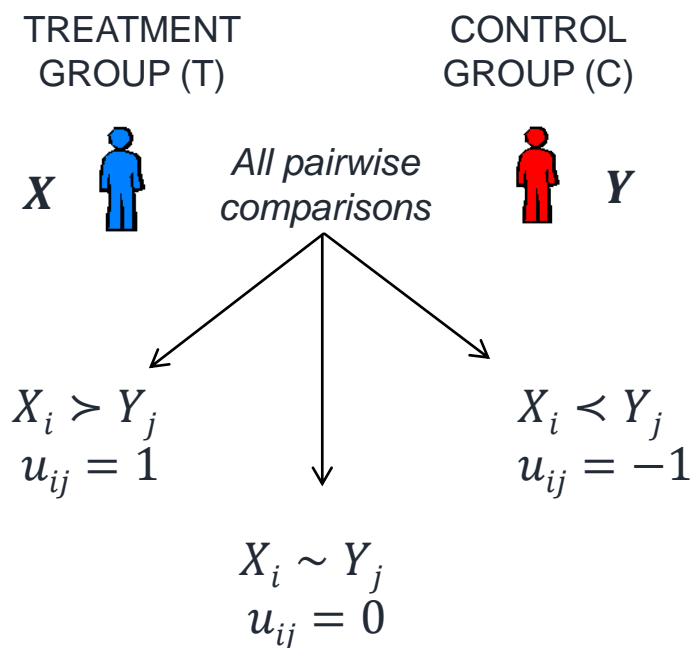
3. The statistic

$$U = \frac{1}{m \cdot n} \sum_{k=1}^K \sum_{i=1}^m \sum_{j=1}^n w(k) u_{ij}(k)$$

has a known distribution under  $H_0$

*Note: weights  $w(k)$  are arbitrary, usually chosen so that  $\sum_{k=1}^K w(k) = 1$*

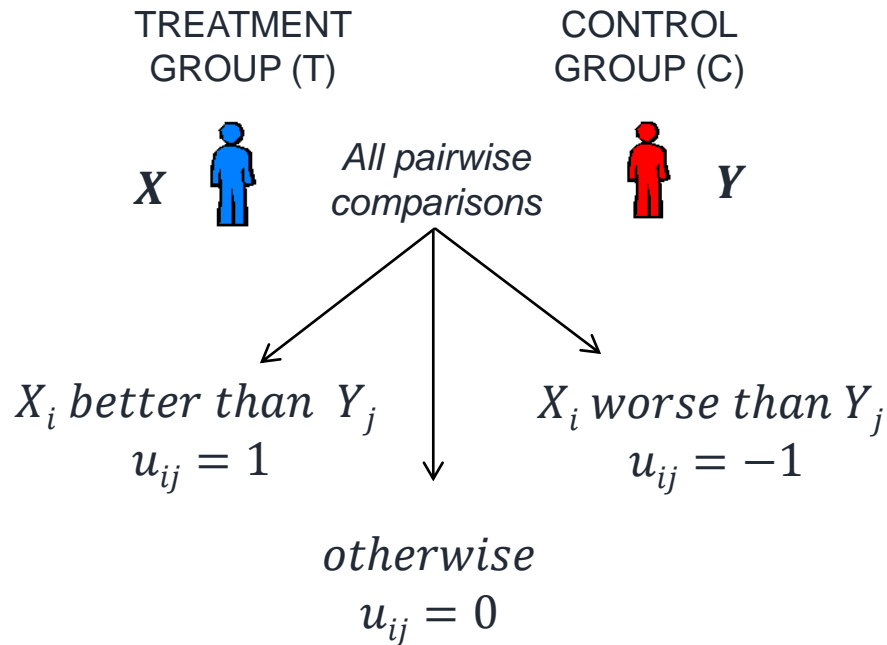
# GPC – multiple prioritized outcomes



Outcome of 1 <sup>st</sup> priority	Outcome of 2 <sup>nd</sup> priority	Overall
Win	-	Win
Loss	-	Loss
Tie or ?	Win	Win
	Loss	Loss
	Tie or ?	Tie or ?

*Note: priorities may be patient-centric*

# Net Treatment Benefit (*NTB*)



The Net Treatment Benefit (*NTB*) is a *U*-statistic

$$U = \frac{1}{m \cdot n} \sum_{i=1}^m \sum_{j=1}^n u_{ij}$$

$$= \frac{\#Wins - \#Losses}{\#Pairs}$$

# Measures of treatment effect

*Finkelstein-Schoenfeld statistic*<sup>1</sup> = #Wins – #Losses

$$NTB^2 = \frac{\#Wins - \#Losses}{\#Pairs}$$

$$Win Ratio^3 = \frac{\#Wins}{\#Losses}$$

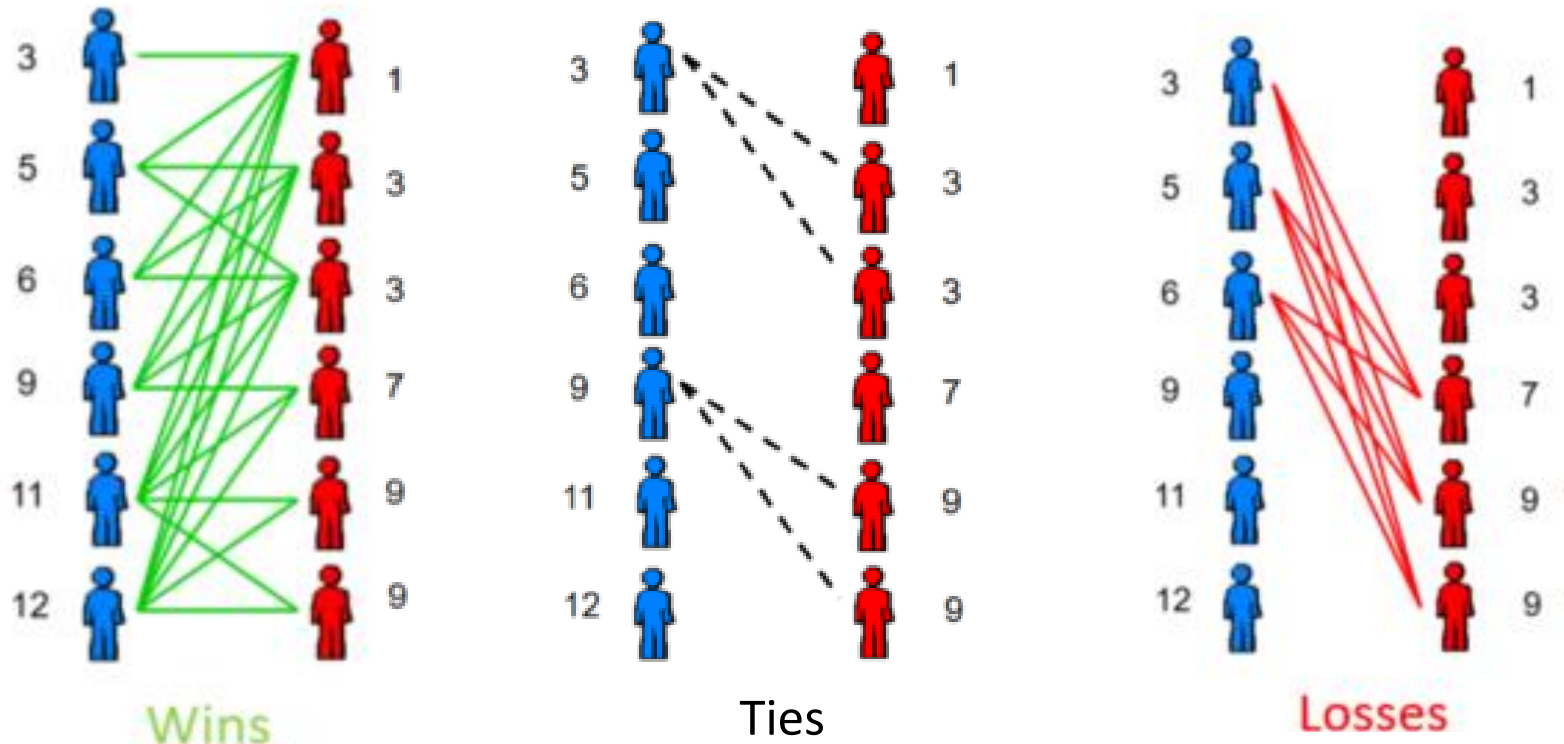
$$Win Odds^{4,5} = \frac{\#Wins + \frac{1}{2}\#(Ties or ?)}{\#Losses + \frac{1}{2}\#(Ties or ?)}$$

Note

$$NTB = \frac{Win Odds - 1}{Win Odds + 1}$$

<sup>1</sup> Finkelstein & Schoenfeld. *Stat Med* 1999;18:1341. <sup>2</sup> Buyse. *Stat Med* 2010;29:3245. <sup>3</sup> Pocock et al. *Eur Heart J* 2012;33:176. <sup>4</sup> Dong et al. *Stat Biopharm Res* 2020;12:99. <sup>5</sup> Brunner et al. *Stat Med* 2021;40:3367.

# Measures of treatment effect



$$NTB = \frac{23-9}{36} = 0.39$$

$$Win\ Ratio = \frac{23}{9} = 2.6$$

$$Win\ Odds = \frac{25}{11} = 2.3$$

# *NTB* – interpretation

*NTB* ranges from -1 to +1, with 0 indicating no overall treatment effect

$$NTB = P(X > Y) - P(Y > X)$$

*NTB* is the *net* probability of a better outcome in one treatment group than in the other

More precisely, *NTB* is the probability that a patient taken at random in the treatment group has a better outcome than a patient taken at random in the control group, minus the probability of the opposite situation.

# *NTB* – relationships

*NTB* is a linear transformation of the probabilistic index *PI*

$$NTB = 2 \cdot PI - 1$$

where

$$PI = P(X > Y) + \frac{1}{2}P(X = Y)$$

*PI* ranges from 0 to 1, with ½ indicating no overall treatment effect

*PI* is closely related to the proportion of similar responses <sup>1</sup>, the concordance index <sup>2</sup> the probability of overlap <sup>3</sup>, and the area under the ROC curve <sup>4</sup>.

<sup>1</sup> Rom & Wang. *Stat Med* 1996;15:1489. <sup>2</sup> Harrell. *Regression Model Strategies*, Springer 2001.

<sup>3</sup> Stine & Heyse. *Stat Med* 2001;20:215. <sup>4</sup> Brumback et al. *Stat Med* 2006;25:575.



# *NTB* – inference and estimation

For testing  $H_0: NTB = 0$ , estimation of *NTB* and confidence limits of *NTB* <sup>1</sup>:

- Exact permutation and bootstrap distribution of the *NTB* statistic<sup>2,3</sup>
- Re-randomization tests <sup>4</sup>
- Bootstrapping for confidence intervals <sup>5</sup>
- Asymptotic distribution of *U*-statistics <sup>6-8</sup>

<sup>1</sup> Verbeek et al. *J Biopharm Stat* 2020;30:765. <sup>2</sup> Finkelstein & Schoenfeld. *Stat Med* 1999;18:1341.

<sup>3</sup> Anderson & Verbeek. <https://arxiv.org/pdf/1901.10928.pdf>, 2019. <sup>4</sup> Buyse. *Stat Med* 2010;29:3245.

<sup>5</sup> Pocock et al. *Eur Heart J* 2012;33:176. <sup>6</sup> Dong et al. *Pharm Stat* 2016;15:430.

<sup>7</sup> Bebu & Lachin. *Biostatistics* 2016;17:178. <sup>8</sup> Ramchandani et al. *Biometrics* 2016;72:926

# *NTB* – adjustment for censoring

*NTB* (Gehan Wilcoxon test) is biased in the presence of censoring <sup>1</sup>.  
The bias can be removed through different approaches <sup>2</sup>

- Naïve, using the proportion of informative pairs <sup>3,4</sup>
- Imputations using the survival distribution <sup>1,5,6</sup>
- Inverse probability of censoring weighting <sup>7,8</sup>

<sup>1</sup> Efron. *Proc 5<sup>th</sup> Berkeley Symp* 1967;4:831. <sup>2</sup> Deltuvaite-Thomas et al. *Biometrical J* 2022.

<sup>3</sup> Harrell et al. *J Am Med Ass* 1982;247:2543. <sup>4</sup> Buyse. *Clin Trials* 2008;5:641.

<sup>5</sup> Latta. *Biometrika* 1977;63:633. <sup>6</sup> Péron et al. *Stat Meth Med Res* 2016;27:1230.

<sup>7</sup> Datta et al. *Scand J Stat* 2010;37:680. <sup>8</sup> Dong et al. *Stat Biopharm Res* 2020;30:882

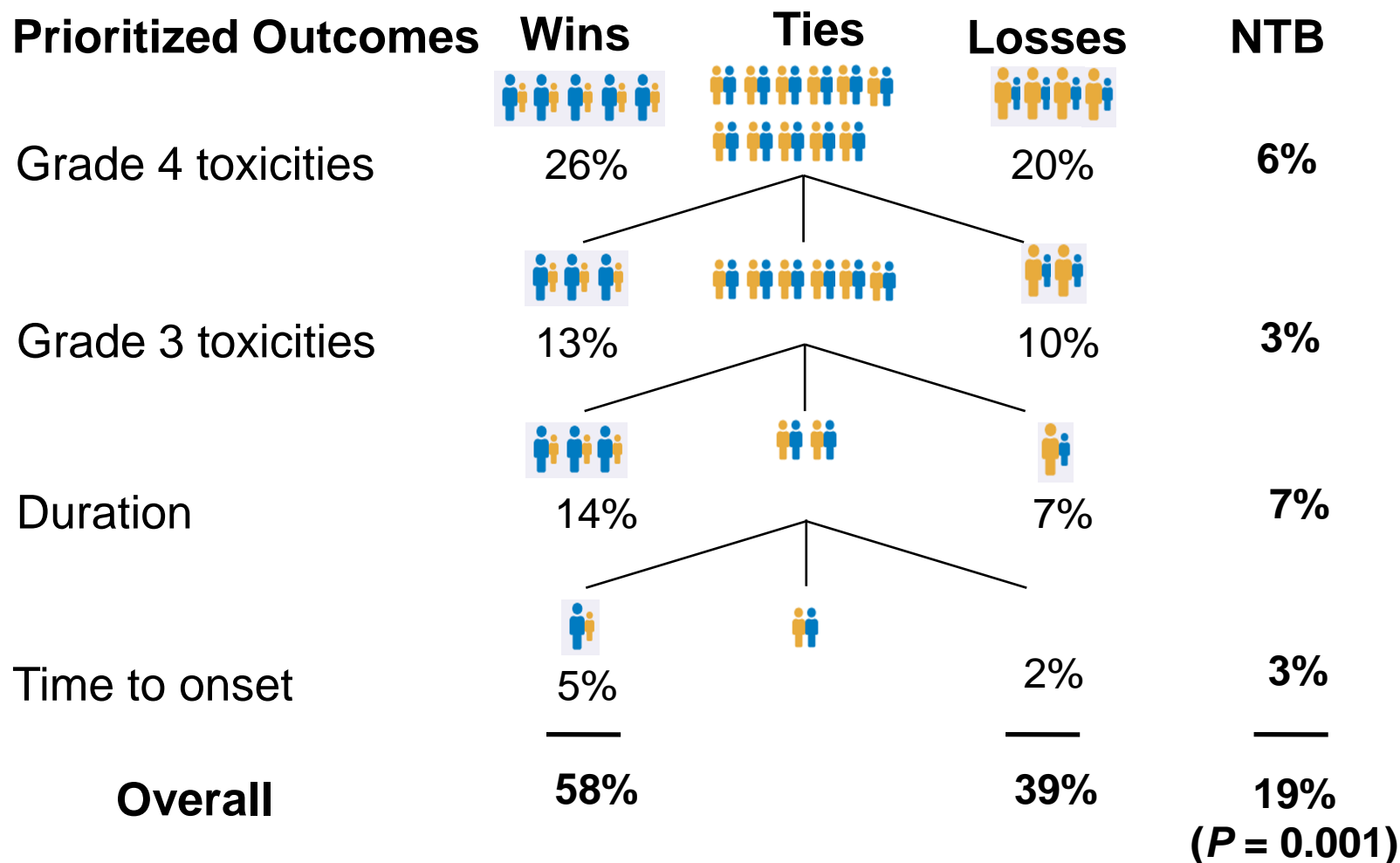
# Applications

# Augmenting power *and* clinical relevance

- Patients with cancer treated aggressively may experience severe toxicities
  - WHO grade 3: severe
  - WHO grade 4: life-threatening
  - WHO grade 5: lethal
- The traditional primary endpoint for comparing an experimental treatment with a control is incidence of WHO grade 3 or worse toxicity
- The analysis should take multiple prioritized outcomes into account:
  1. Severity (lower WHO grade better)
  2. Duration of severe toxicity (shorter better)
  3. Time to onset (later better)

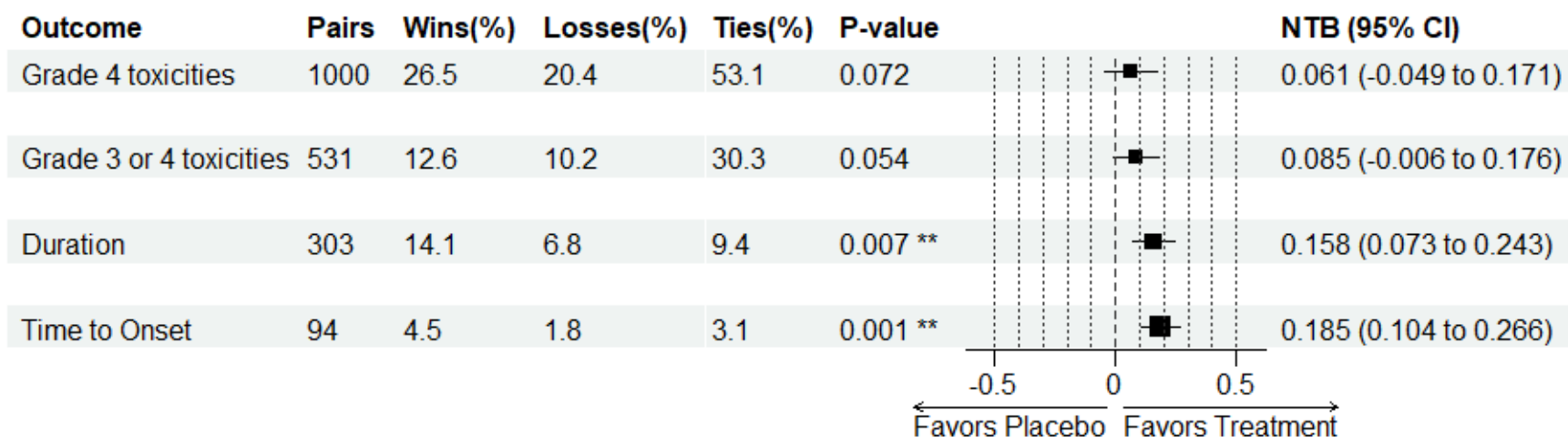
# Augmenting power *and* clinical relevance

Placebo controlled trial of experimental treatment protecting against a specific toxicity



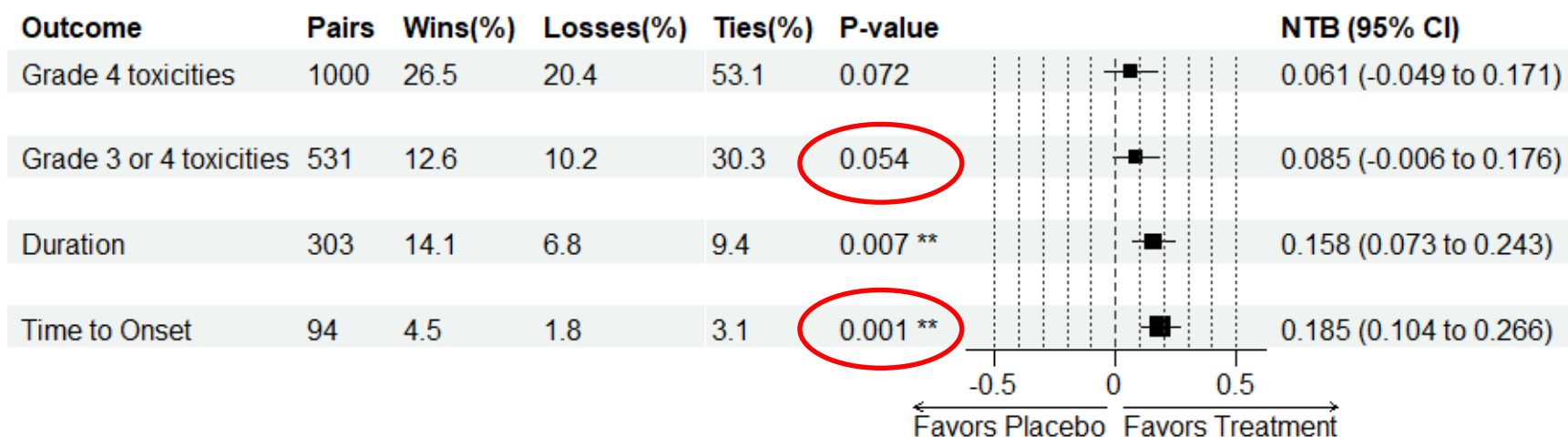
# Augmenting power *and* clinical relevance

Placebo controlled trial of experimental treatment protecting against a specific toxicity



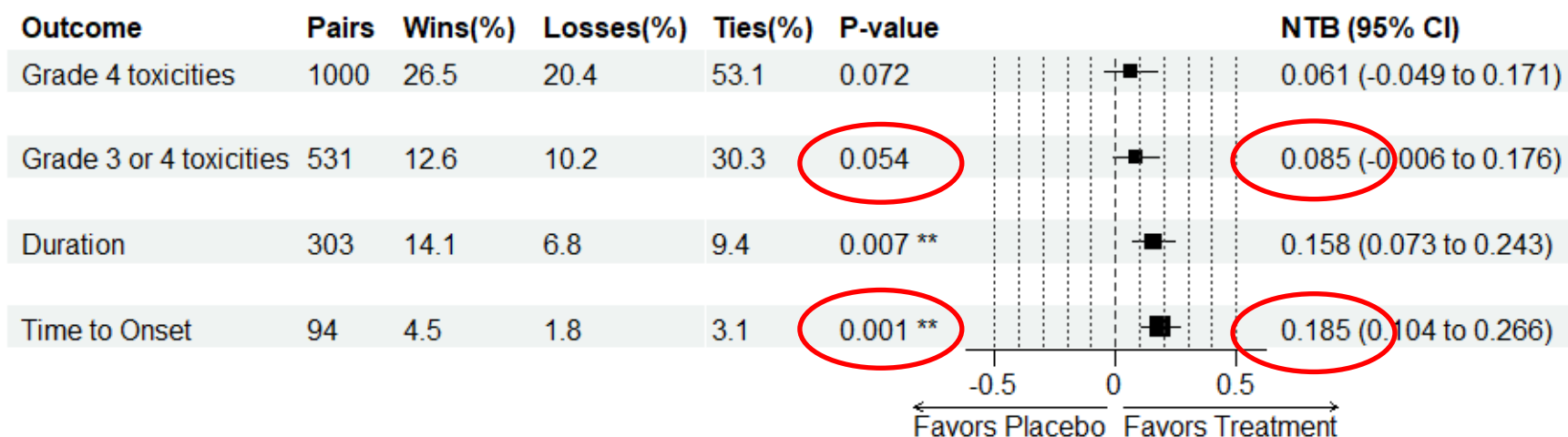
# Augmenting power *and* clinical relevance

Placebo controlled trial of experimental treatment protecting against a specific toxicity



# Augmenting power *and* clinical relevance

Placebo controlled trial of experimental treatment protecting against a specific toxicity





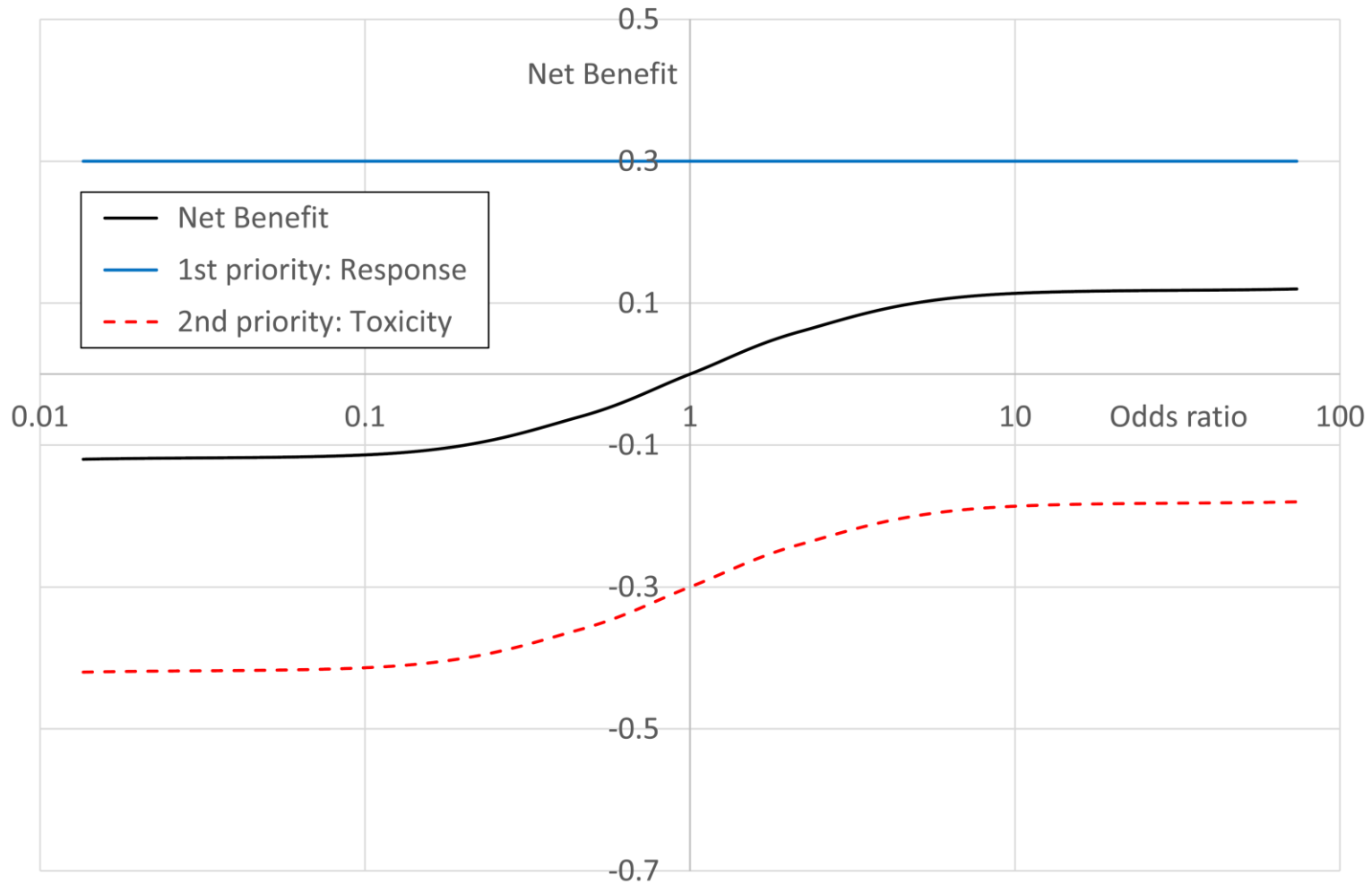
# Benefit / risk analysis

- Simple situation of binary efficacy outcome (1 = response, 0 = no response) and binary safety outcome (1 = no toxicity, 0 = toxicity)

<b>Outcomes</b>	<b>Treatment</b>	<b>Control</b>	<b>Difference</b>
Response rate (benefit)	0.5	0.2	0.3
Toxicity rate (risk)	0.6	0	0.6
Marginal benefit / risk difference			-0.3

- Naïve analysis suggests negative benefit / risk of -0.3
- What would GPC analysis show, assuming achievement of response is preferred to avoidance of toxicity?

# Benefit / risk analysis



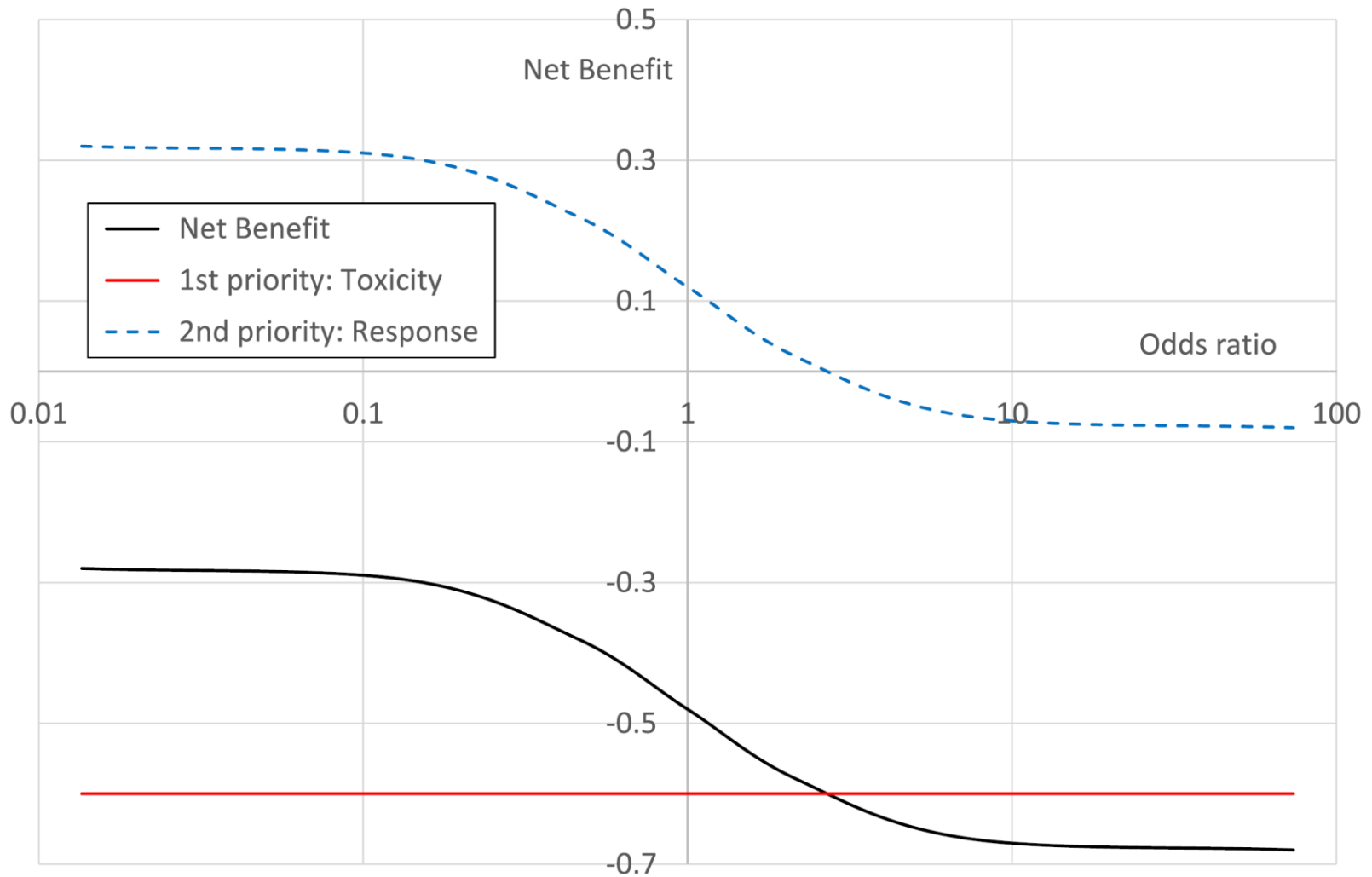
# Benefit / risk analysis

- $NTB$  depends on the association (odds ratio,  $OR$ ) between response and toxicity
  - If  $OR > 1$ ,  $NTB > 0$  : patients who respond also have toxicity (e.g., skin rash for inhibitors of the EGFR pathway)
  - If  $OR = 1$ ,  $NTB = 0$  : response is independent of toxicity (e.g., cardiac toxicities of anthracyclins)
  - If  $OR < 1$ ,  $NTB < 0$  : patients who do not respond have toxicity (e.g., toxicities to irinotecan in patients with enzyme deficiencies)

# Benefit / risk analysis

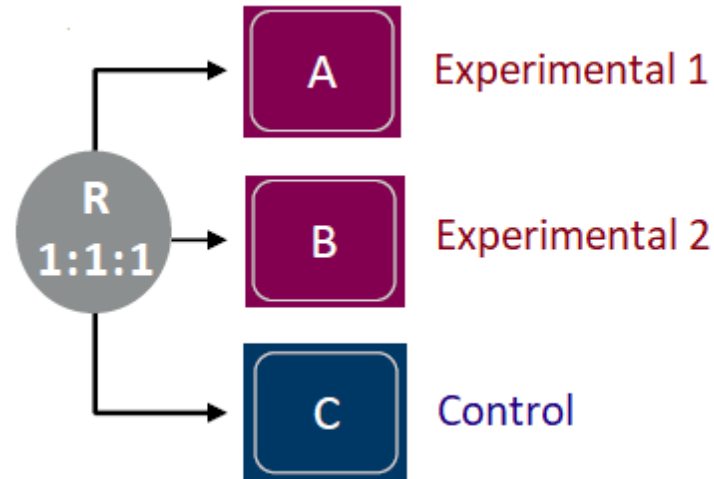
- *NTB* depends on the association (odds ratio, *OR*) between response and toxicity
  - If  $OR > 1$ ,  $NTB > 0$  : patients who respond also have toxicity (e.g., skin rash for inhibitors of the EGFR pathway)
  - If  $OR = 1$ ,  $NTB = 0$  : response is independent of toxicity (e.g., cardiac toxicities of anthracyclins)
  - If  $OR < 1$ ,  $NTB < 0$  : patients who do not respond have toxicity (e.g., toxicities to irinotecan in patients enzyme deficiencies)
- *NTB* would be quite different if avoidance of toxicity was preferred to achievement of response, allowing for patient-centric treatment choices

# Benefit / risk analysis



# Multiple Testing Procedures

Assume several treatments are compared to a standard of care



Comparisons: A vs. C (Experimental 1, preferred)

B vs. C (Experimental 2)

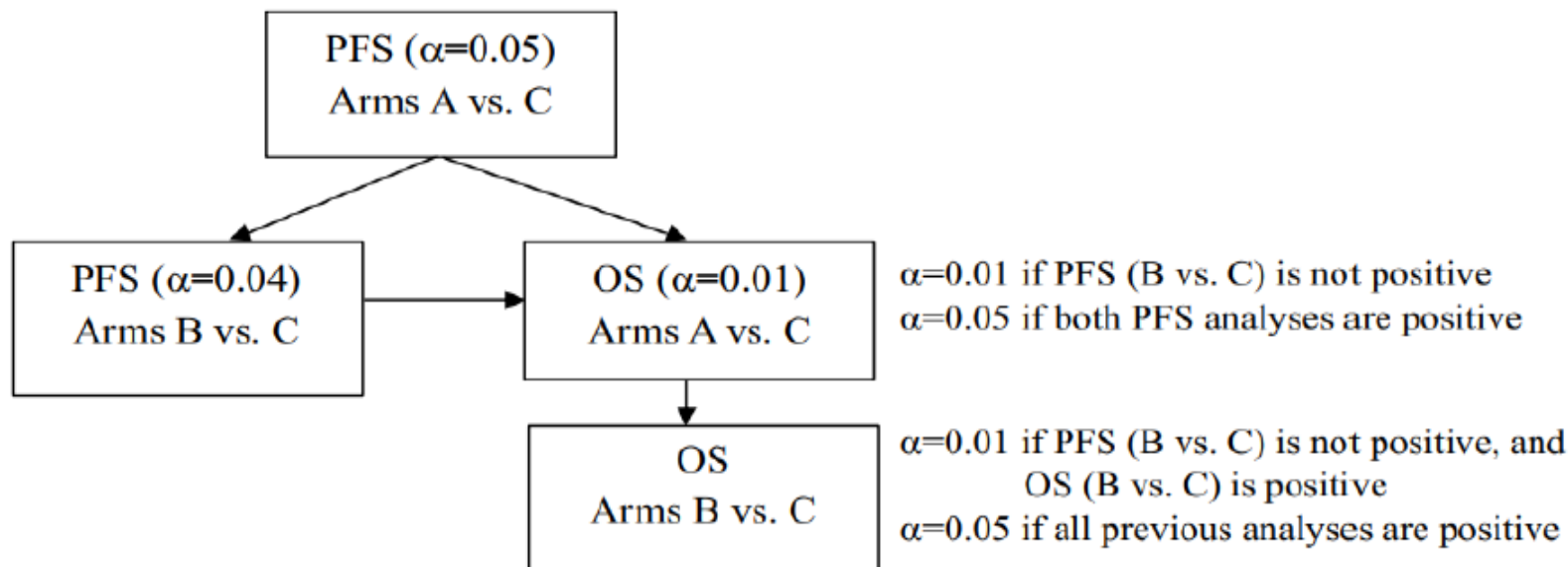
A vs. B (Not powered)

Outcomes: PFS (« Primary »)

OS (« Key secondary »)

# Multiple Testing Procedures

Testing procedure with strict control of type I error rate



OS of the preferred experimental arm is tested at full level of significance (0.05) *only if PFS of the other (non preferred) experimental arm reaches statistical significance !*

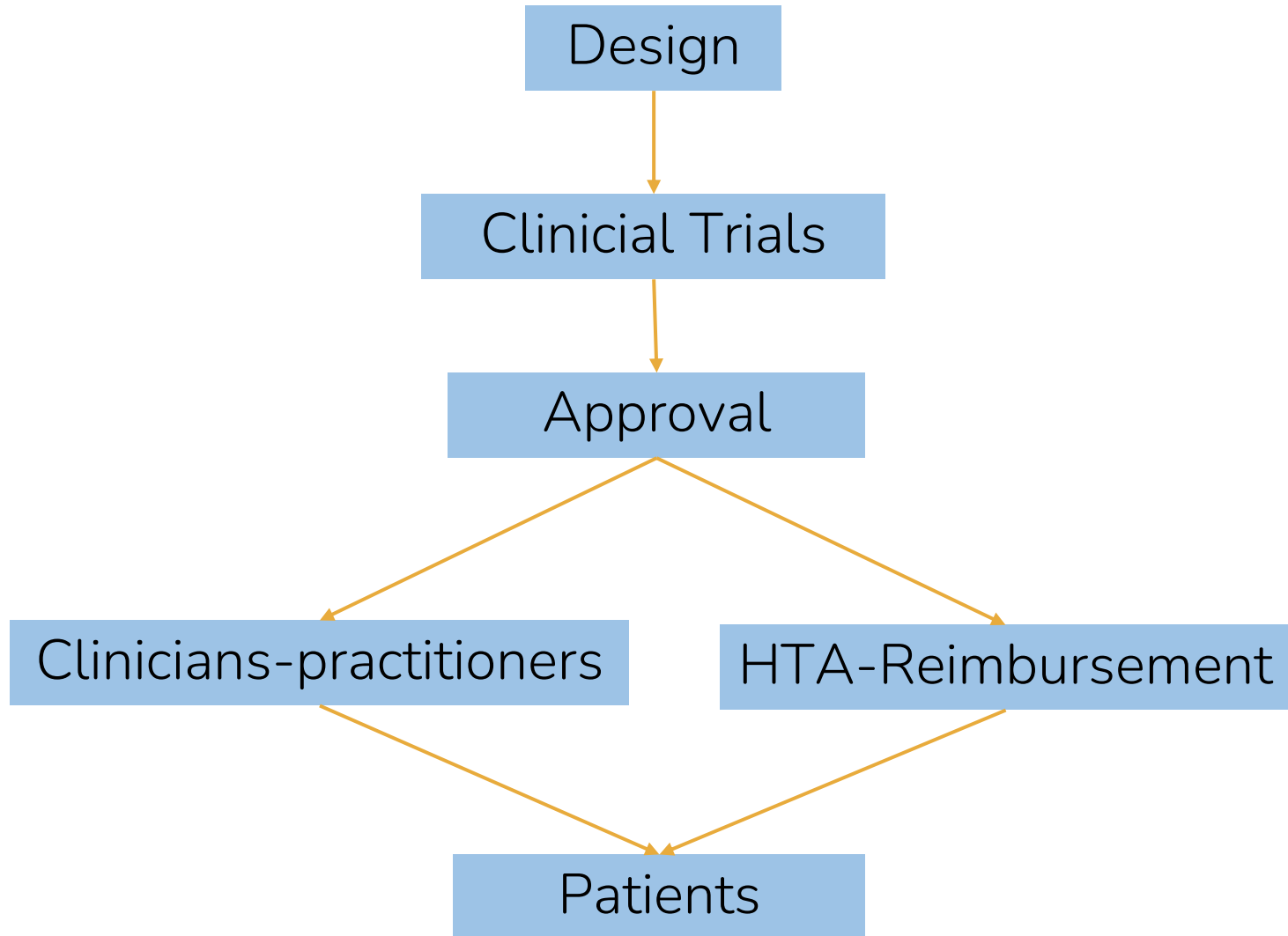
# Conclusions



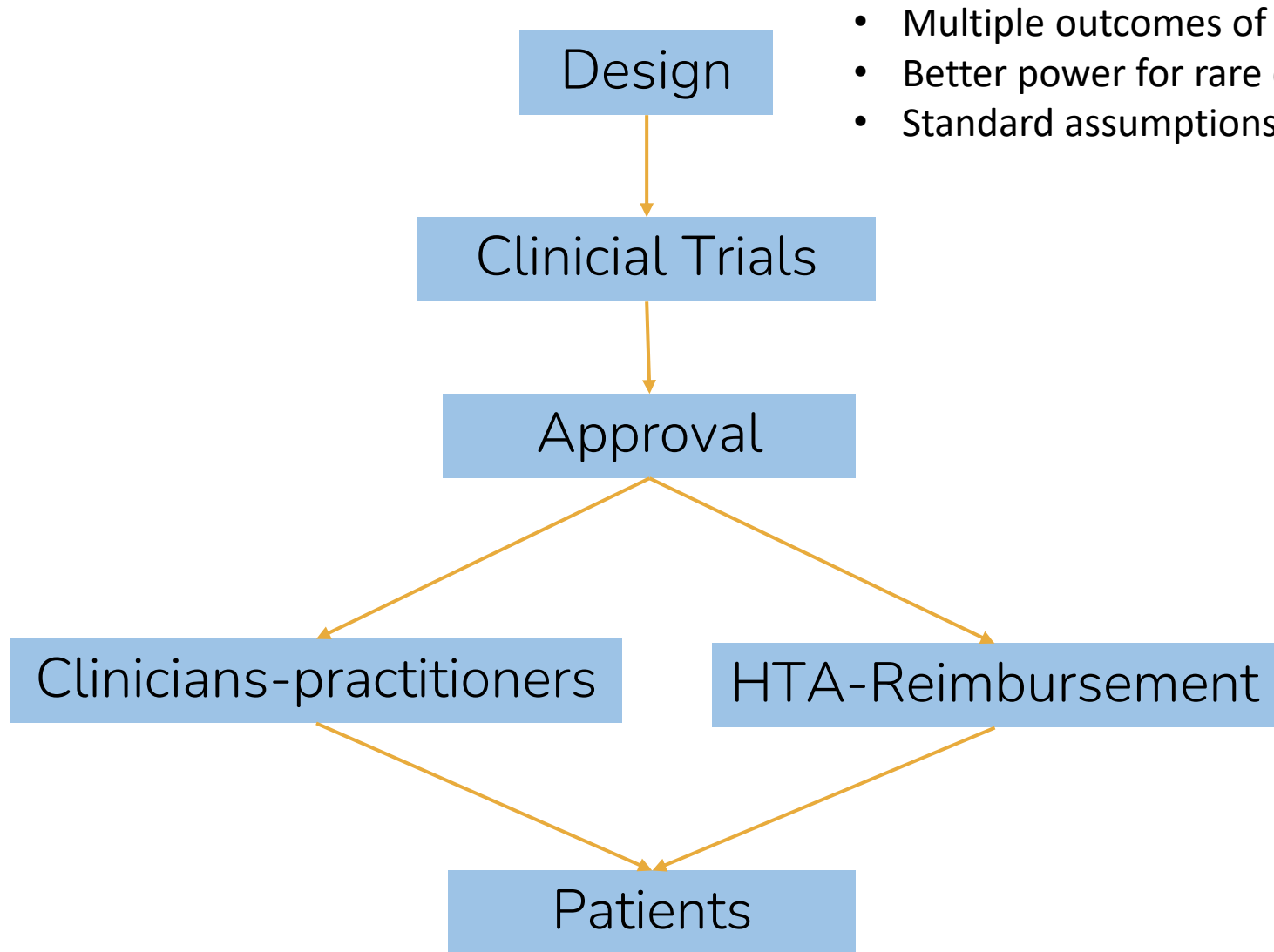
# GPC benefits

- Increases flexibility of analyses
- Incorporates multiple outcomes
- Incorporates thresholds of clinical relevance
- May increase power as compared with single outcome
- Can be adapted to individual patient preferences
- Provides unique measure of treatment effect that is meaningful to patients and caregivers

# GPC usage

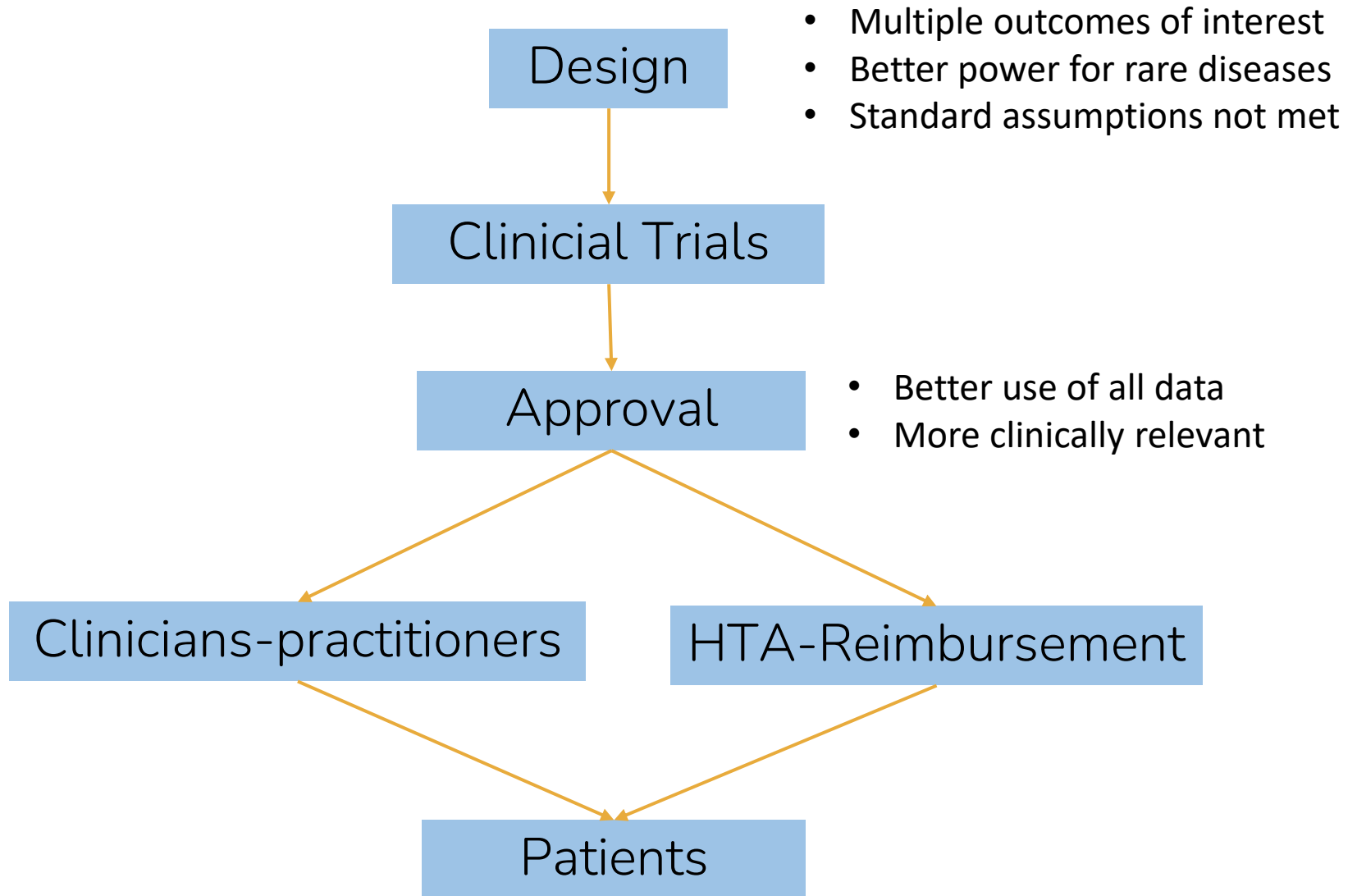


# GPC usage

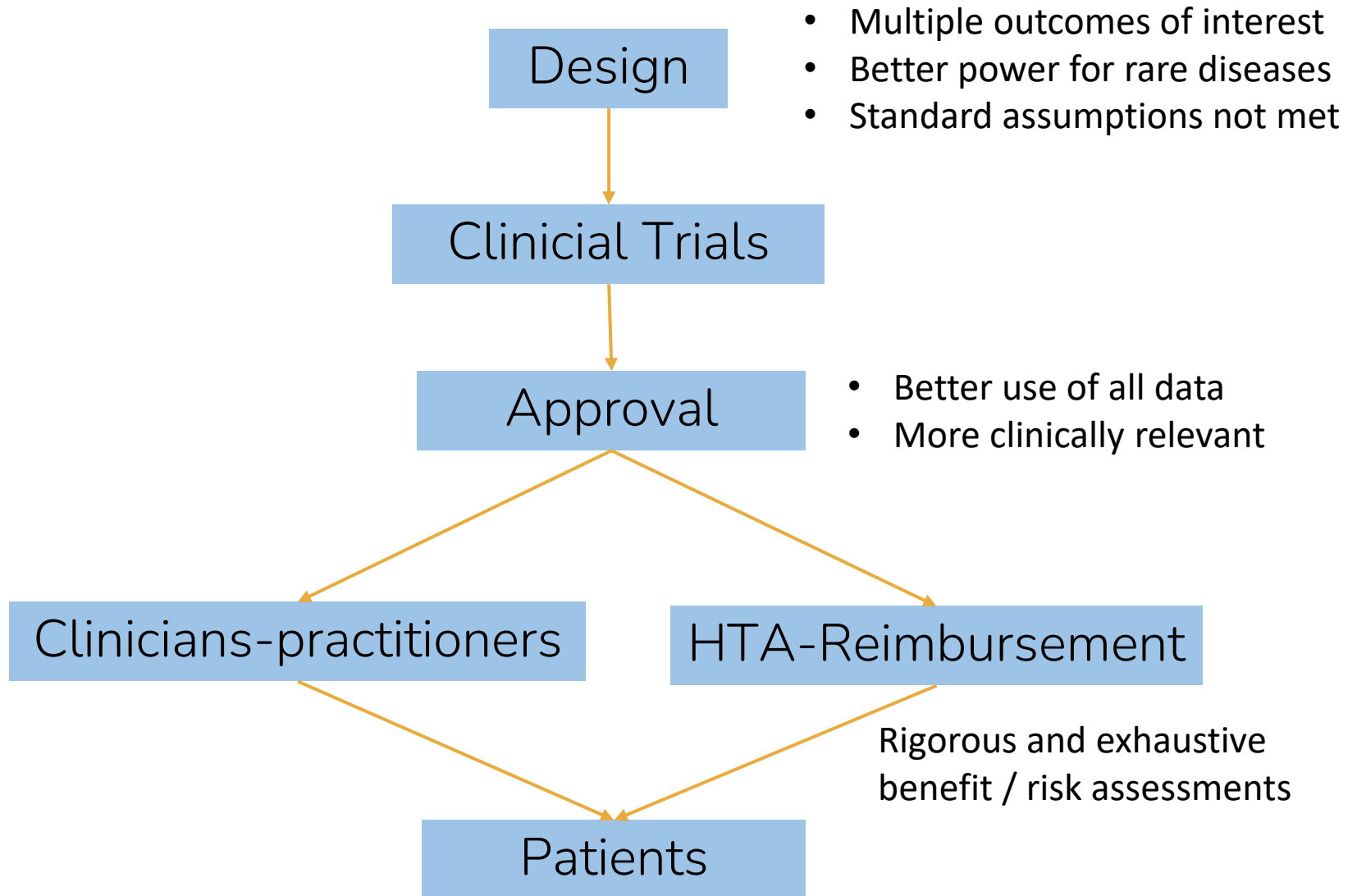


- Multiple outcomes of interest
- Better power for rare diseases
- Standard assumptions not met

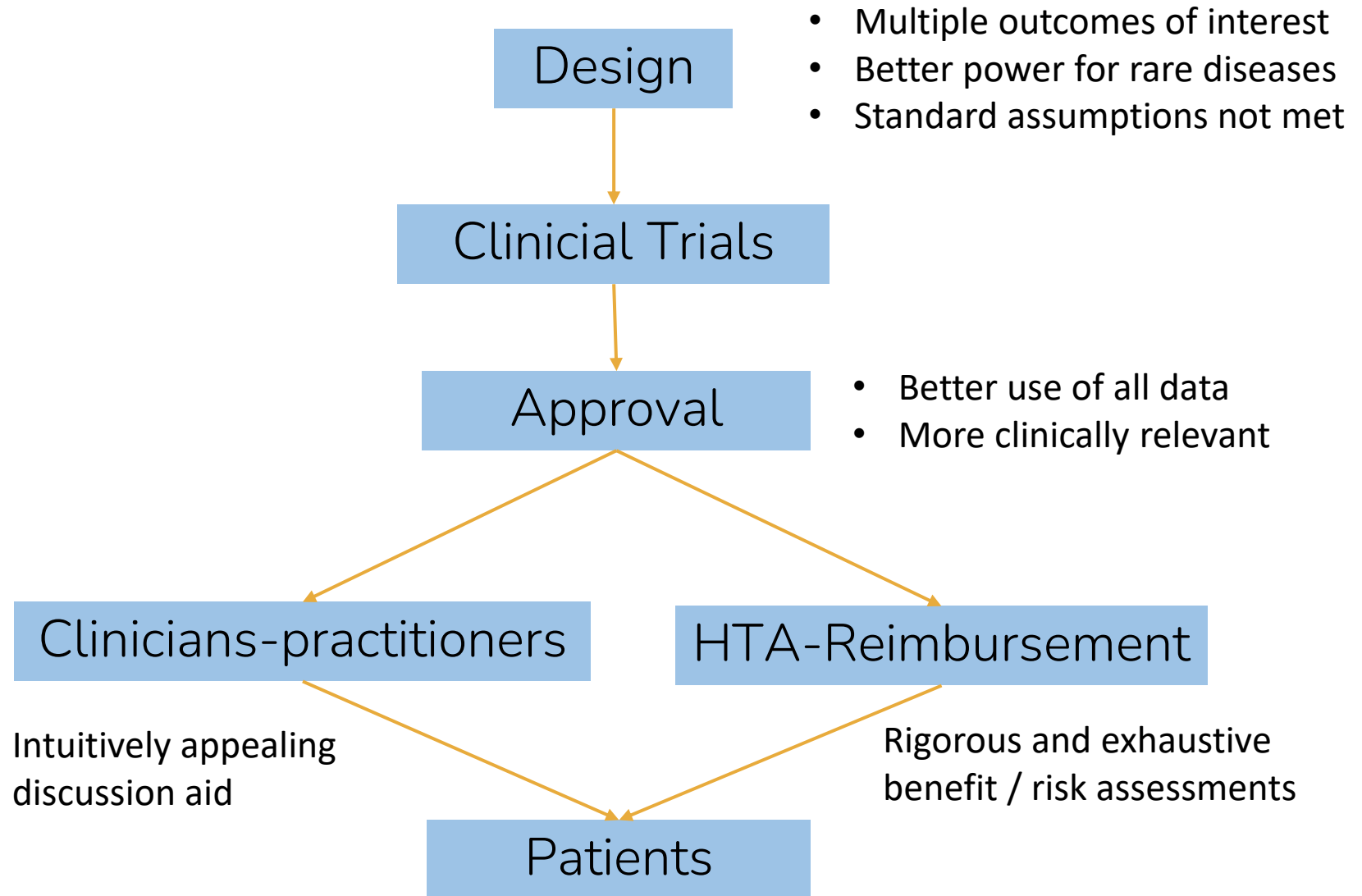
# GPC usage



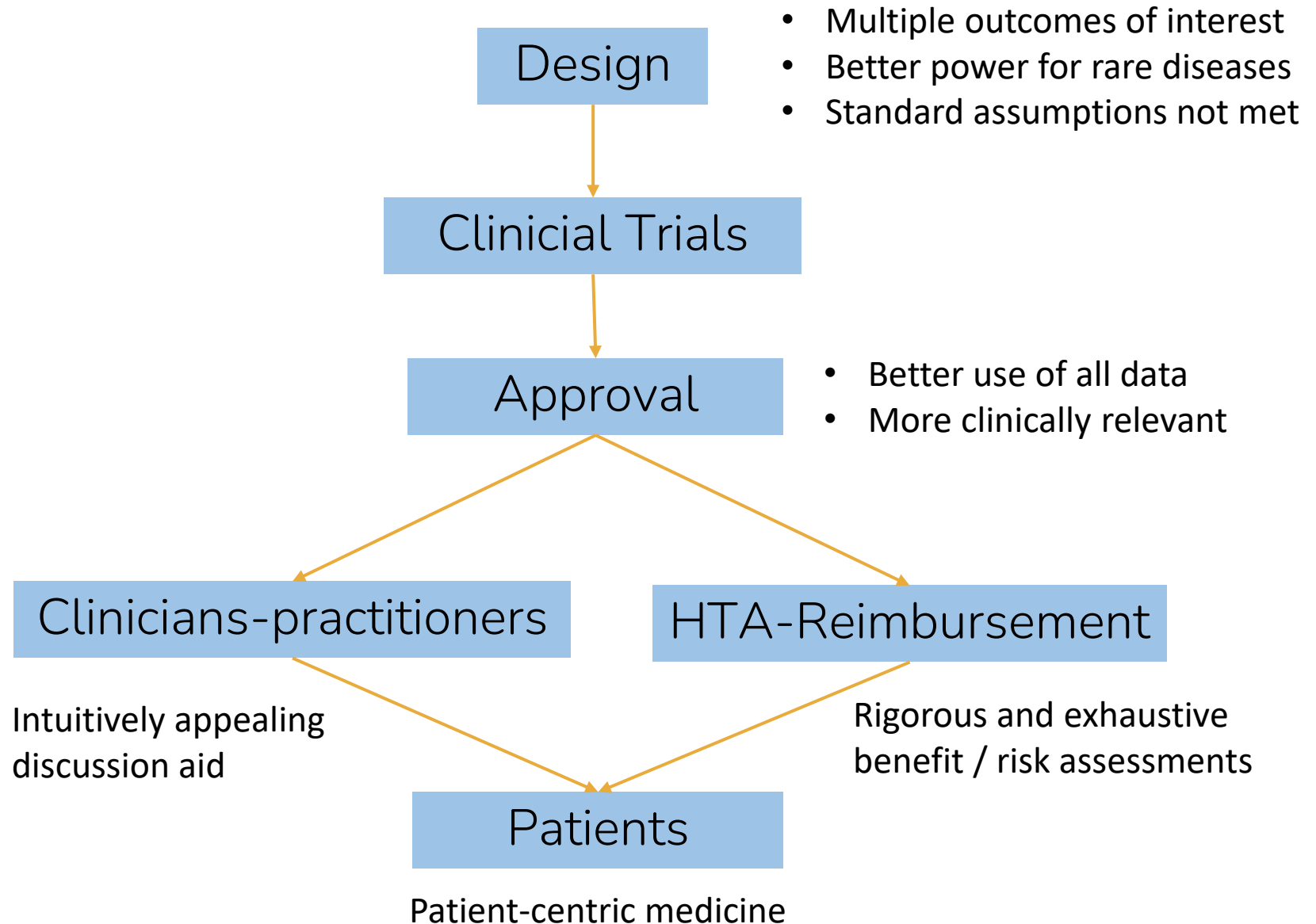
# GPC usage



# GPC usage



# GPC usage



# Questions / References

[marc.buyse@iddi.com](mailto:marc.buyse@iddi.com)