

SOHC 2023 – Grundkurs Onkologie

Statistische Überlegungen in der Onkologie – eine Auswahl

Stefan Aebi

Luzerner Kantonsspital

Tumorzentrum, Medizinische Onkologie

6000 Luzern 16

stefan.aebi@onkologie.ch

The Large Print Giveth and the Small Print Taketh Away

Tom Waits

Potentielle Interessenkonflikte: Keine

Ich bin nicht Statistiker. Wenn Sie eine Studie planen, suchen Sie die Zusammenarbeit mit einem Biostatistiker!

„To consult the statistician after an experiment is finished is often merely to ask him to conduct a post mortem examination. He can perhaps say what the experiment died of.“

R.A. Fisher

Consulting Services

Support for your clinical trials and registries

The SAKK CC offers statistical and data management consulting services

Planning, analysis, and reporting of your studies or registries

Statistical support for publications

CRF and data base development in secuTrial and data base hosting

Data cleaning

Statistical education (e. g. Oncolunches, workshops), basic or tailored to your needs



Go to the SAKK website or contact Dr. Stefanie Hayoz (Stefanie.Hayoz@sakk.ch) for more information!



Themen

- ◆ Studiendesign
- ◆ Randomisation
- ◆ Endpunkte
- ◆ «Statistisch signifikant», α und β
- ◆ Stichprobenumfang
- ◆ Effektstärke, absolut und relativ
 - Odds ratio, relatives Risiko
 - Hazard ratio
- ◆ Überlebensanalyse
- ◆ Multivariable Analysen und Subgruppen-Analysen
- ◆ Meta-Analysen
- ◆ Vertrauensintervalle
- ◆ Überlegenheit und Nicht-Unterlegenheit
- ◆ «Komparatoren»
- ◆ Interne und externe Validität



THINGS GOT REALLY INTERESTING WHEN THE STATISTICIAN STARTED DOING WARD ROUNDS.

Endpunkte, Messgrößen

endpoints

Was wollen Sie messen?

- ◆ Tumormasse
- ◆ Überleben bis Ereignis
 - Rezidiv, Fernmetastase
 - Wiederauftreten einer malignen Neoplasie
 - Progression
 - Tod
- ◆ Befinden des Patienten

Endpunkte, Messgrößen

endpoints

Was wollen Sie messen?

- ◆ Tumormasse
- ◆ Überleben bis Ereignis
 - Rezidiv, Fernmetastase
 - Wiederauftreten einer malignen Neoplasie
 - Progression
 - Tod
- ◆ Befinden des Patienten

Messgröße in Studien

- ◆ Anteil Ansprechen (RECIST)
- ◆ «Time to event»-Daten
 - RFS, DDFS
 - DFS, iDFS
 - PFS
 - OS
- ◆ Lebensqualität, z.B. QLQ C-30, SF36, ...

Endpunkte, Messgrößen

endpoints

Was wollen Sie messen?

- ◆ Tumormasse
- ◆ Überleben bis Ereignis
 - Rezidiv, Fernmetastase
 - Wiederauftreten einer malignen Neoplasie
 - Progression
 - Tod
- ◆ Befinden des Patienten

Messgröße in Studien

◆ Anteil Ansprechen (RECIST)

◆ «Time to event»-Daten

• RFS, DDFS

• DFS, iDFS

• PFS

**Medikamenten-
Entwicklung,
SURROGATE**

• OS

«patientenrelevant»

◆ Lebensqualität,
z.B. QLQ C-30, SF36, ...

Endpunkte, Messgrößen

endpoints

Was wollen Sie messen?

- ◆ Tumormasse
- ◆ Überleben bis Ereignis
 - Rezidiv, Fernmetastase
 - Wiederauftreten einer malignen Neoplasie
 - Progression
 - Tod
- ◆ Befinden des Patienten

Messgröße in Studien

◆ Anteil Ansprechenden (RECIST)

◆ «Time to event»-Daten

• PFS, DDFS

• DFS, iDFS

• PFS

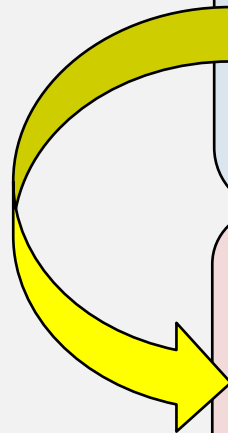
• OS

◆ Lebensqualität,
z.B. QLQ C-30, SF36, ...

**Positive
Korrelation?**

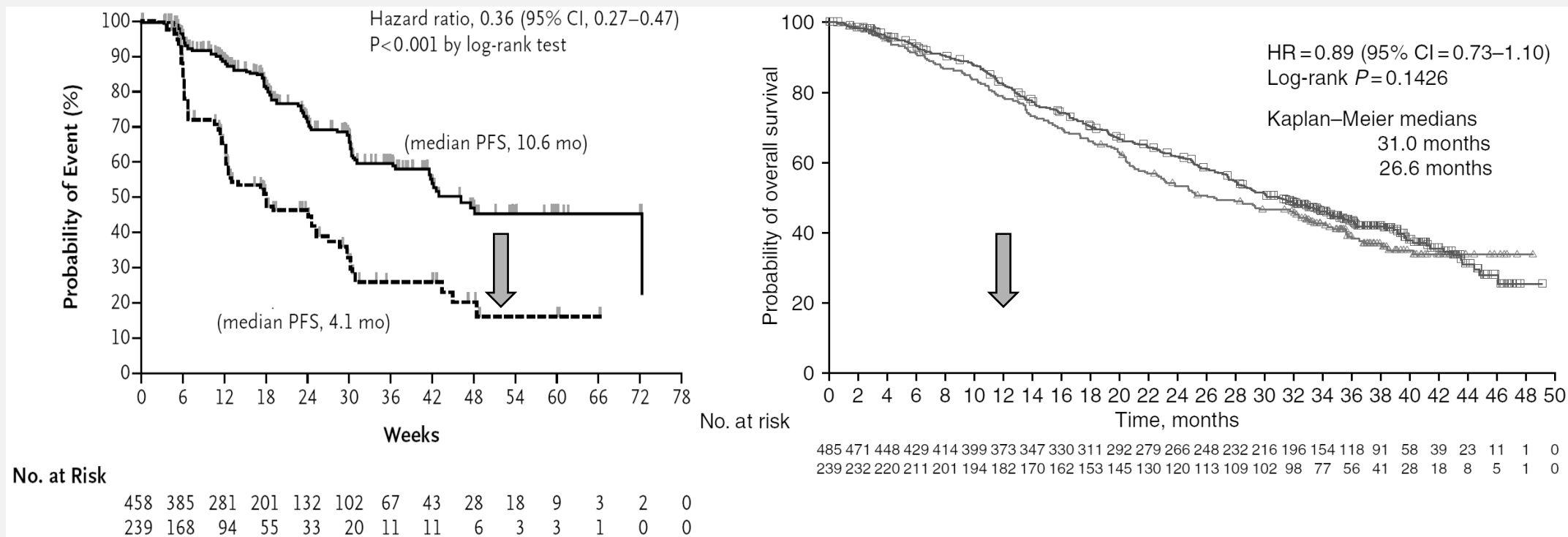
**Medikamenten-
Entwicklung,
SURROGATE**

«patientenrelevant»



Beispiel: PFS → OS?

BOLERO-2, Mammakarzinom Stadium IV. Exemestan±Everolimus



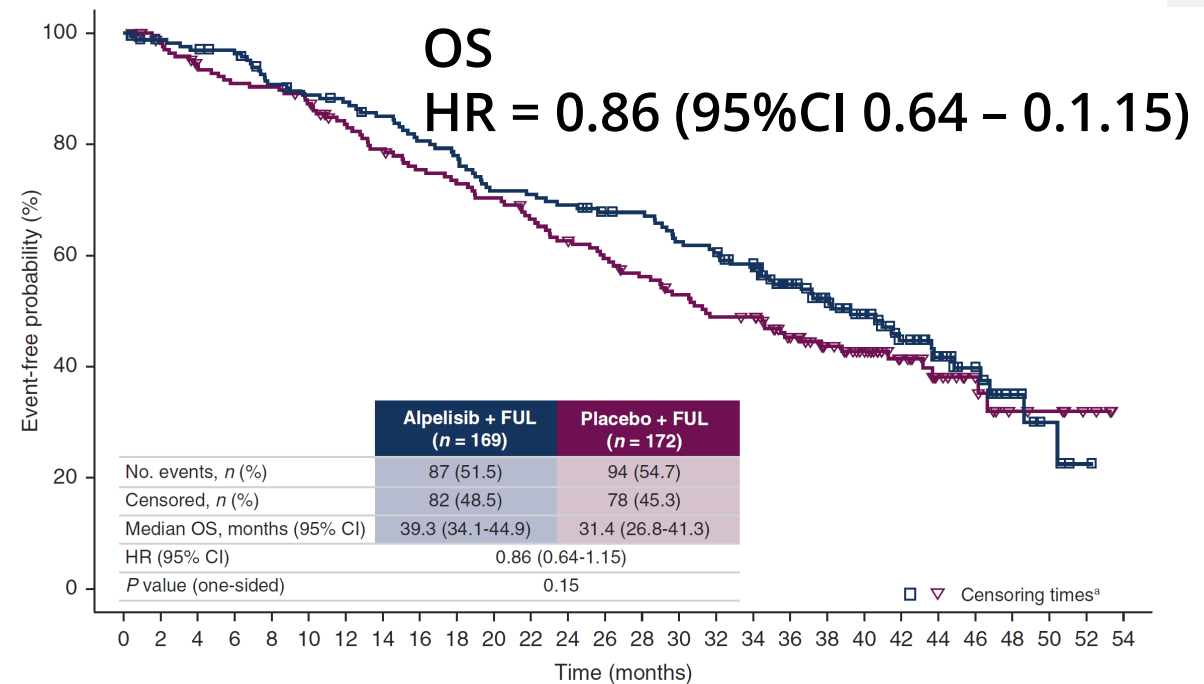
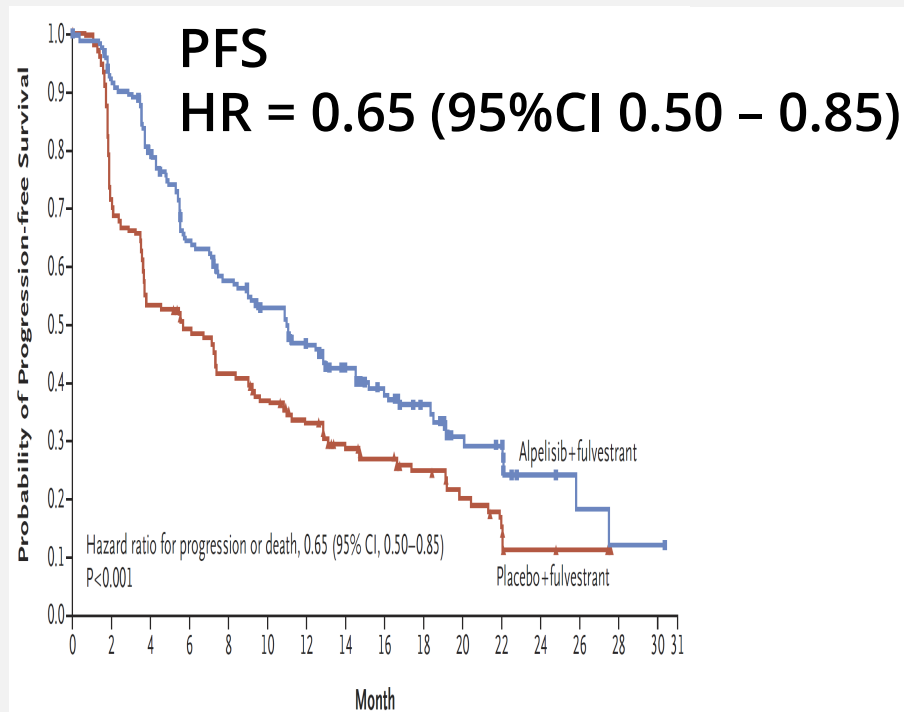
New England Journal of Medicine → Annals of Oncology

Verbesserung von PFS ist in diesem Beispiel kein Prädiktor von OS

PFS ist nicht generell ein Surrogat für OS

Beispiel: PFS → OS?

SOLAR-1, Mammakarzinom Stadium IV nach AI. Fulvestant ± Alpelisib

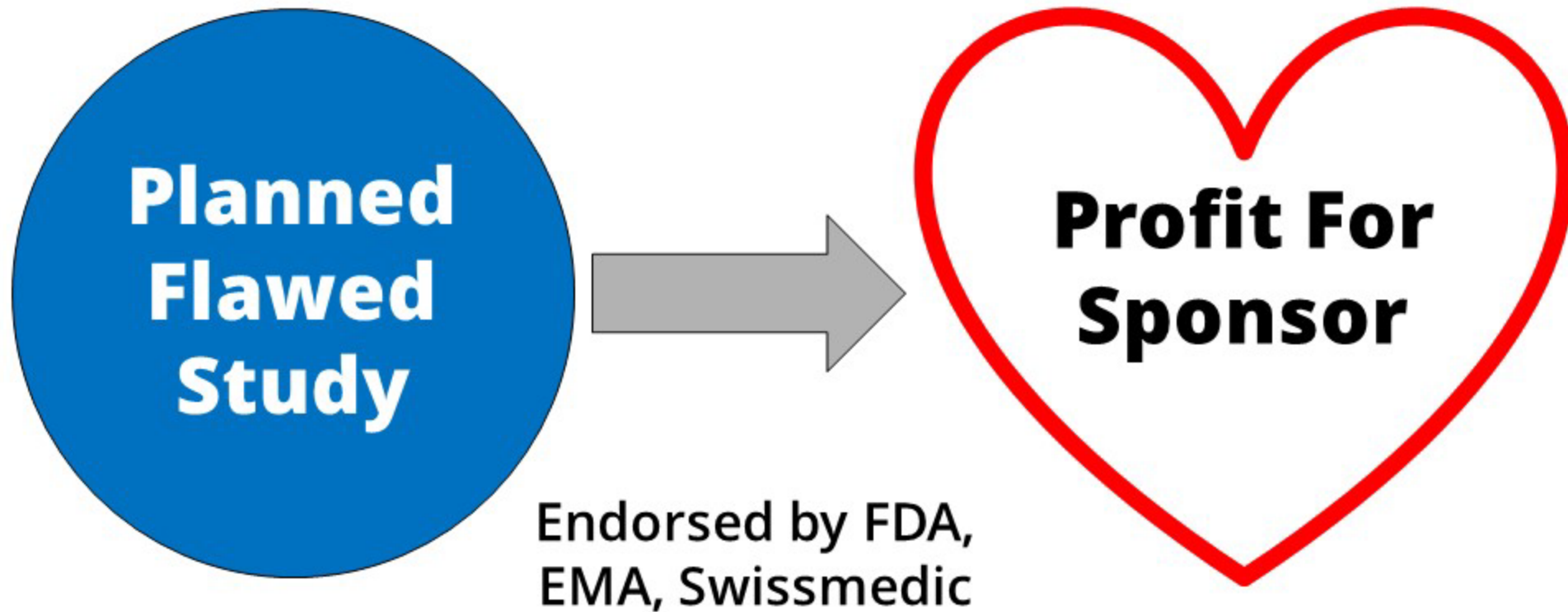


New England Journal of Medicine → Annals of Oncology

Verbesserung von PFS ist in diesem Beispiel kein Prädiktor von OS

PFS ist nicht generell ein Surrogat für OS

So what does PFS really mean?



Slide stolen from Prof. Ian Tannock, 2019

Hypothesentests

«Statistisch signifikant», α und β

- ◆ Vergleich von 2 Therapien
- ◆ Messgrößen der Studiengruppen A und B werden als Stichproben der Messgrößen aller potentieller Patienten betrachtet, die mit A oder B behandelt werden.
- ◆ Hypothese: Behandlungen A und B haben dieselbe Wirkung, die Messgrößen sind gleich verteilt. Die Differenz der Mittelwerte* der Messgrößen A und B ist 0.

* Differenzen und Mittelwerte als Beispiel, ebenfalls möglich sind z.B. Quotienten, Anzahl, etc.

Hypothesentests

«Statistisch signifikant», α und β

- ◆ Vergleich von 2 Therapien
- ◆ Messgrößen der Studiengruppen A und B werden als Stichproben der Messgrößen aller potentieller Patienten betrachtet, die mit A oder B behandelt werden.
- ◆ Hypothese: Behandlungen A und B haben dieselbe Wirkung, die Messgrößen sind gleich verteilt. Die Differenz der Mittelwerte* der Messgrößen A und B ist 0.



**H₀,
Nullhypothese**

* Differenzen und Mittelwerte als Beispiel, ebenfalls möglich sind z.B. Quotienten, Anzahl, etc.

«Statistisch signifikant», α und β

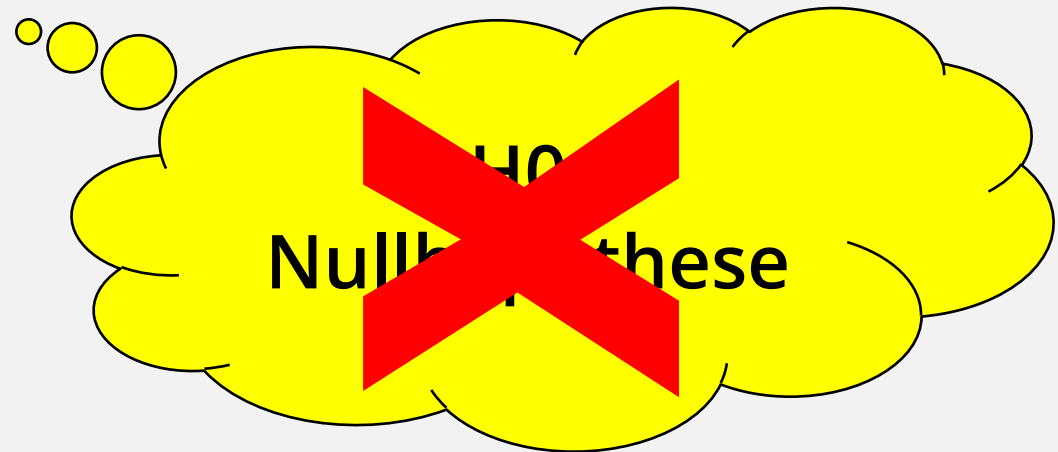
- ◆ Messgrößen der Studiengruppen A respektive B werden als Stichproben der Messgrößen aller potentieller Patienten betrachtet, die mit A oder B behandelt werden.
- ◆ Hypothese: Behandlungen A und B haben dieselbe Wirkung, die Messgrößen sind gleich verteilt. Die Differenz der Mittelwerte* der Messgrößen A und B ist 0.
- ◆ P = Wahrscheinlichkeit eine Differenz zwischen den Mittelwerten von A und B zu finden, die mindestens gleich gross ist wie die beobachtete *unter Annahme der Null-Hypothese*.
- ◆ Ist $P < \alpha$, dann ist die Wahrscheinlichkeit, dass sich die Mittelwerte von A und B in Wahrheit nicht unterscheiden, kleiner als α .

* Differenzen und Mittelwerte als Beispiel, ebenfalls möglich sind z.B. Quotienten, Anzahl, etc.

«Statistisch signifikant», α und β

- ◆ Messgrößen der Studiengruppen A respektive B werden als Stichproben der Messgrößen aller potentieller Patienten betrachtet, die mit A oder B behandelt werden.
- ◆ Hypothese: Behandlungen A und B haben dieselbe Wirkung, die Messgrößen sind gleich verteilt. Die Differenz der Mittelwerte* der Messgrößen A und B ist 0.
- ◆ P = Wahrscheinlichkeit eine Differenz zwischen den Mittelwerten von A und B zu finden, die mindestens gleich gross ist wie die beobachtete *unter Annahme der Null-Hypothese*.
- ◆ Ist $P < \alpha$, dann ist die Wahrscheinlichkeit, dass sich die Mittelwerte von A und B in Wahrheit nicht unterscheiden, kleiner als α .

* Differenzen und Mittelwerte als Beispiel, ebenfalls möglich sind z.B. Quotienten, Anzahl, etc.



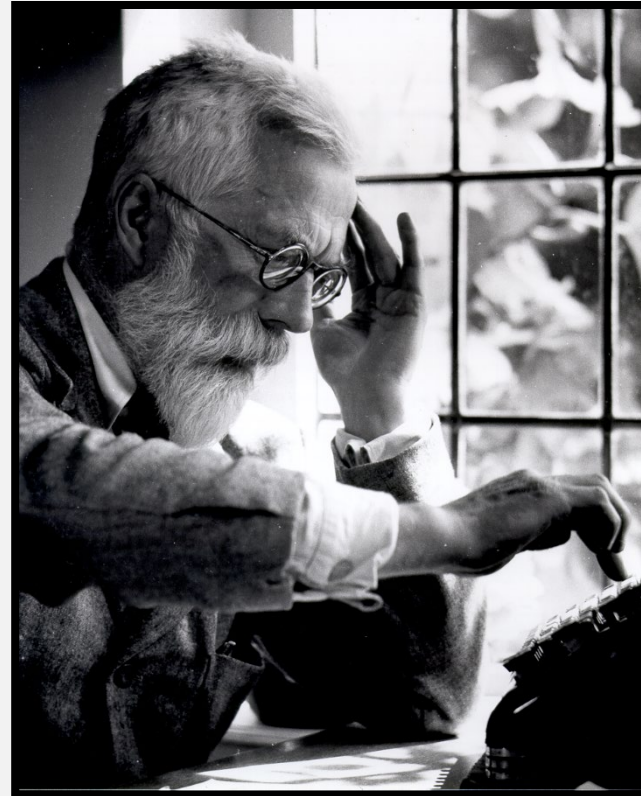
«Statistisch signifikant», α und β

- ◆ Der Wert der «Signifikanzschwelle» α ist arbiträr, konventionell ist 0.05 (5%)
- ◆ $P < \alpha \rightarrow$ wir postulieren, dass sich die Mittelwerte* von A und B unterscheiden – «statistisch signifikant» – und wissen dabei, dass wir uns selten irren.

*Differenzen und Mittelwerte als Beispiel, ebenfalls möglich sind z.B. Quotienten, Anzahl, etc.

«Statistisch signifikant», α und β

- ◆ Der Wert der «Signifikanzschwelle» α ist arbiträr, konventionell ist 0.05 (5%)
- ◆ $P < \alpha \rightarrow$ wir postulieren, dass sich die Mittelwerte von A und B unterscheiden – «statistisch signifikant» – und wissen dabei, dass wir uns selten irren.



Sir Ronald A. Fisher: «We are not often astray if we draw a conventional line at 0.05».

«Statistisch signifikant», α und β

- ◆ Der Wert der «Signifikanzschwelle» α ist arbiträr, konventionell ist 0.05 (5%)
- ◆ $P < \alpha \rightarrow$ wir postulieren, dass sich die Mittelwerte von A und B unterscheiden – «statistisch signifikant» – und wissen dabei, dass wir uns selten irren.



H1, Alternativhypothese

«Statistisch signifikant», α und β

- ◆ Der Wert der «Signifikanzschwelle» α ist arbiträr, konventionell ist 0.05 (5%)
- ◆ $P < \alpha \rightarrow$ wir postulieren, dass sich die Mittelwerte von A und B unterscheiden – «statistisch signifikant» – und wissen dabei, dass wir uns selten irren.
- ◆ Den Irrtum, «falsch positiver» Schluss, bezeichnen wir als **Fehler der ersten Art.**
- ◆ β = Wahrscheinlichkeit einen wahren Unterschied nicht zu finden: «falsch negativer» Schluss, **Fehler der zweiten Art.**
- ◆ $1 - \beta$ = Trennschärfe oder **Power**

Hypothesentests

«Statistisch signifikant», α und β

		«Wirklichkeit»	
		H0 ist wahr	H1 ist wahr
Entscheidung des Tests	für H0	richtige Entscheidung Wahrscheinlichkeit $1-\alpha$	Fehler 2. Art Wahrscheinlichkeit β
	für H1	Fehler 1. Art Wahrscheinlichkeit α	richtige Entscheidung Wahrscheinlichkeit $1-\beta$ (Power)

Effektgrösse relativ und absolut

effect size

Statistisch signifikant
ist nicht notwendig
klinisch signifikant

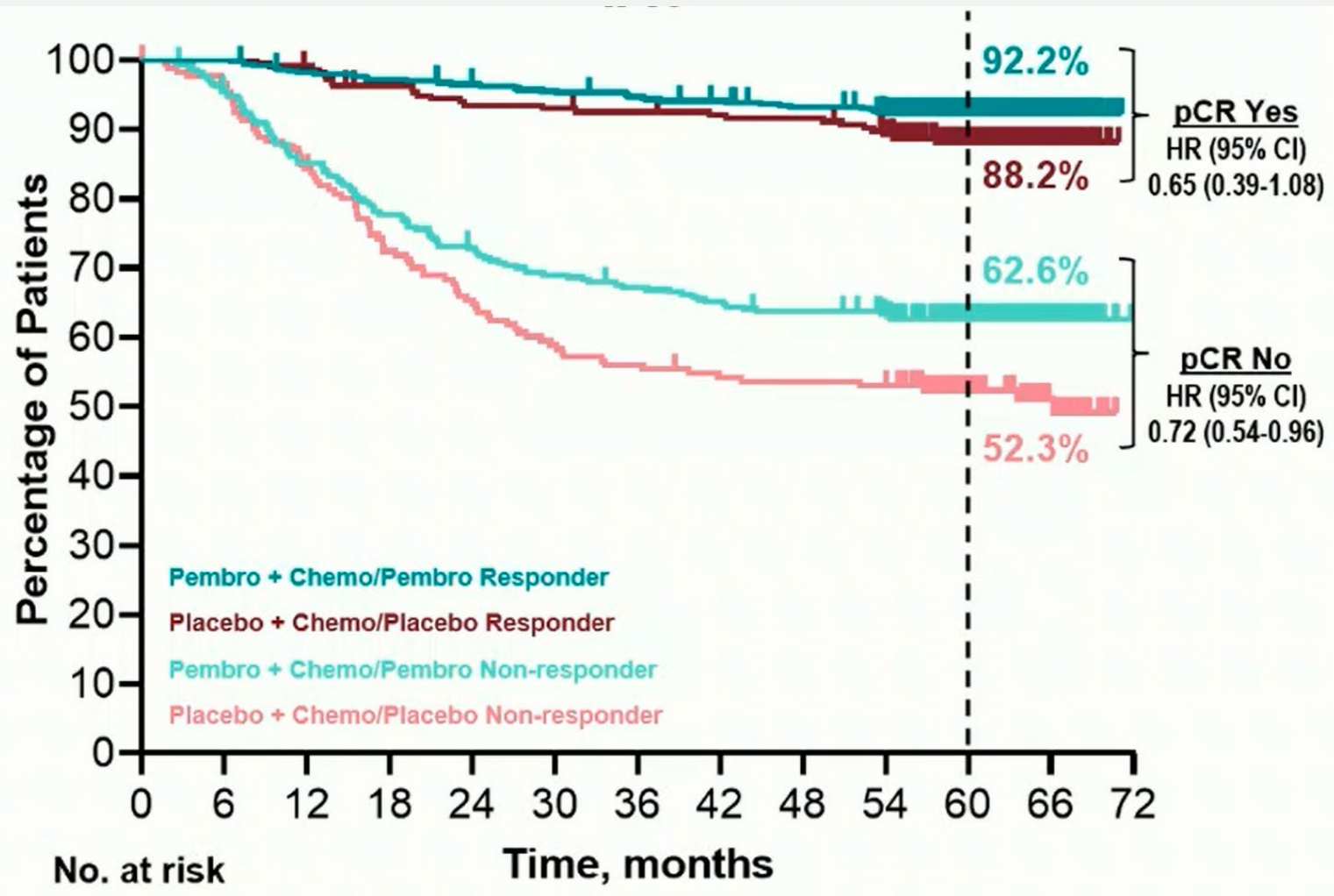
... aber ...

«statistically not significant
but clinically meaningful»
ist Unsinn



Effektgrösse relativ und absolut

effect size



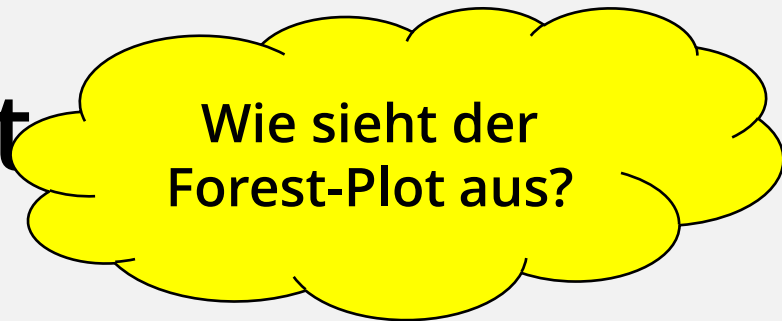
Relativ: Hazard ratio

pCR Yes HR = 0.65

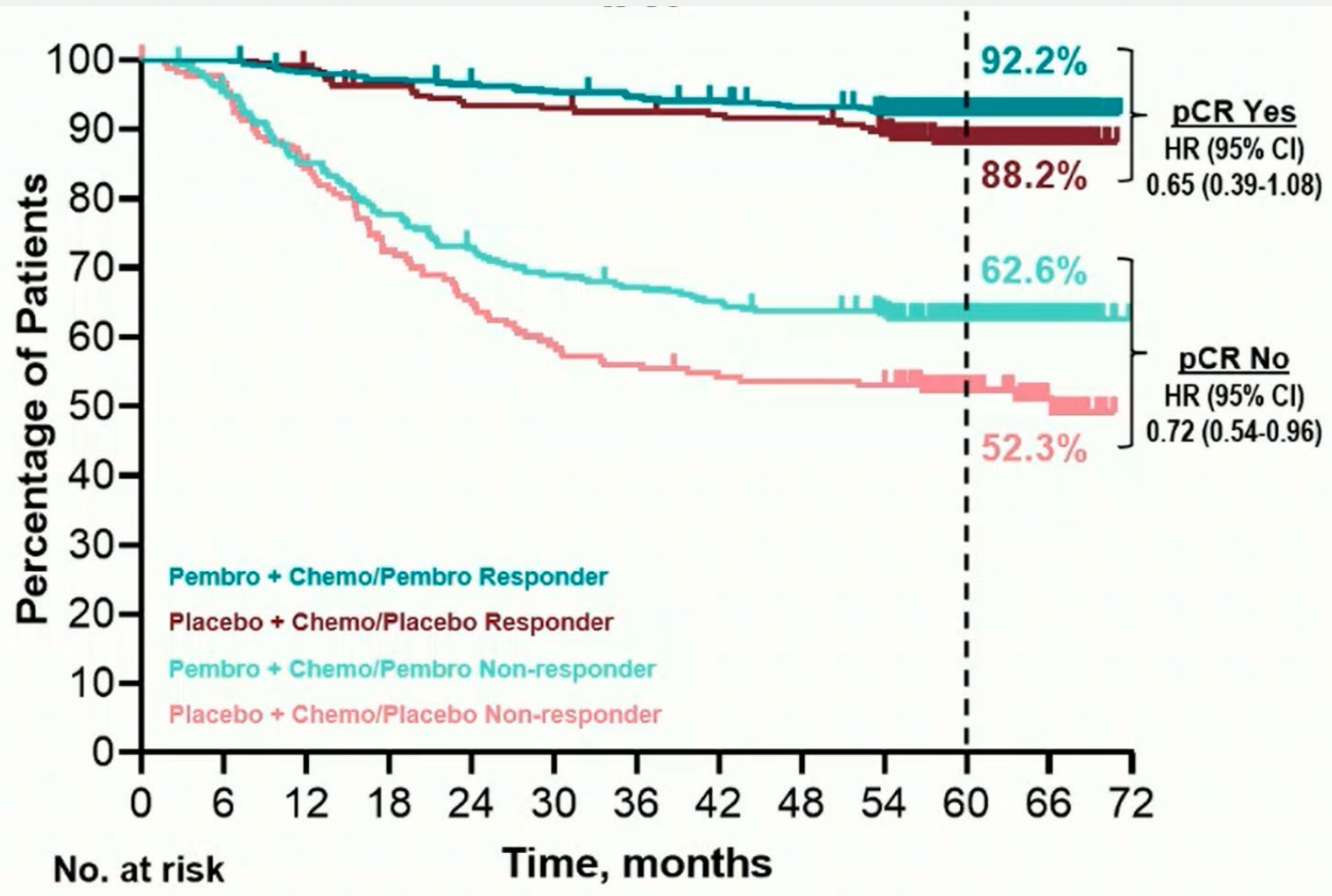
pCR No HR = 0.72

Effektgrösse relativ und absolut

effect size



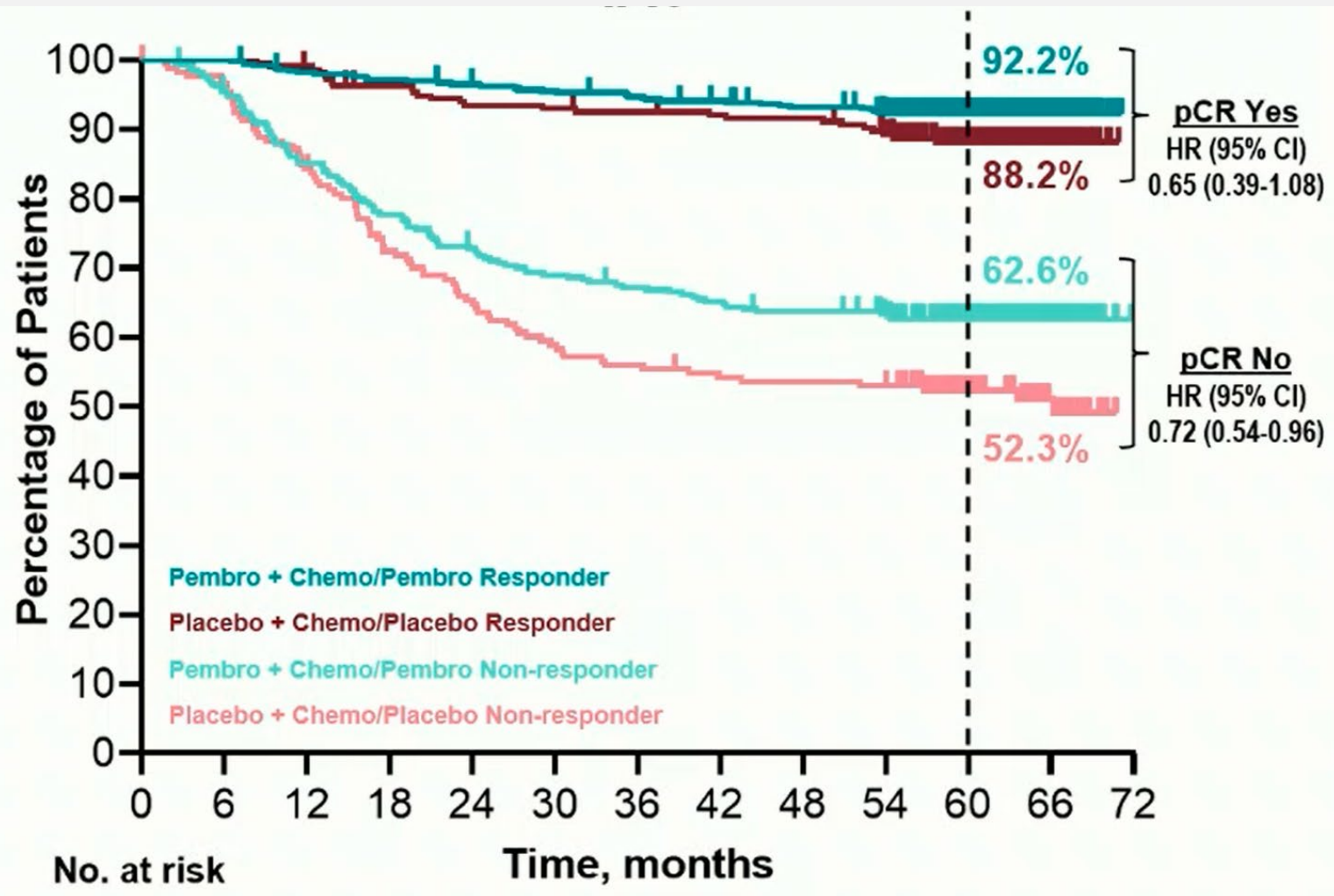
Wie sieht der Forest-Plot aus?



Relativ: Hazard ratio
pCR Yes HR = 0.65
pCR No HR = 0.72

Effektgrösse relativ und absolut

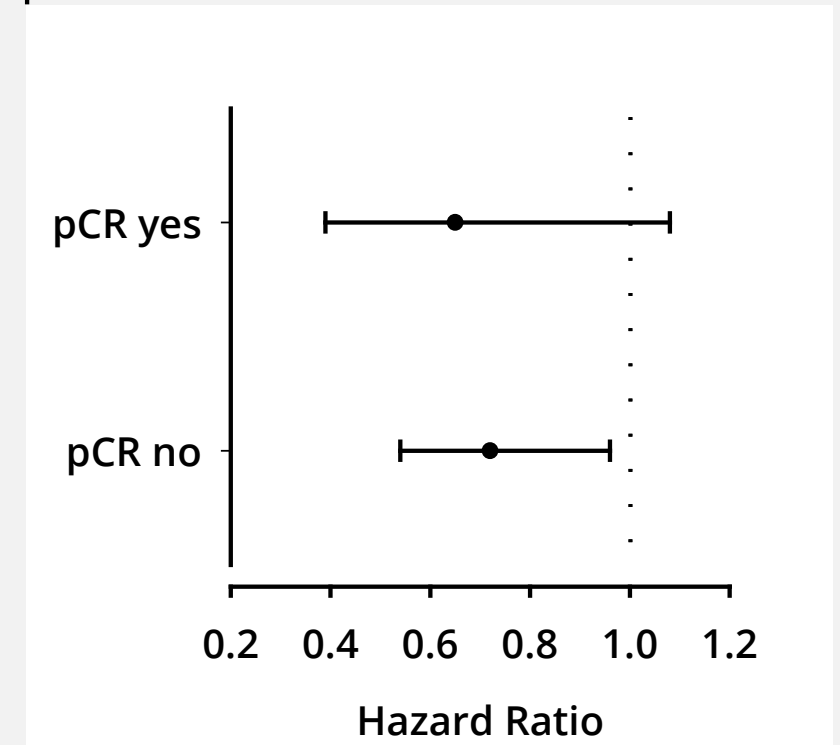
effect size



Relativ: Hazard ratio

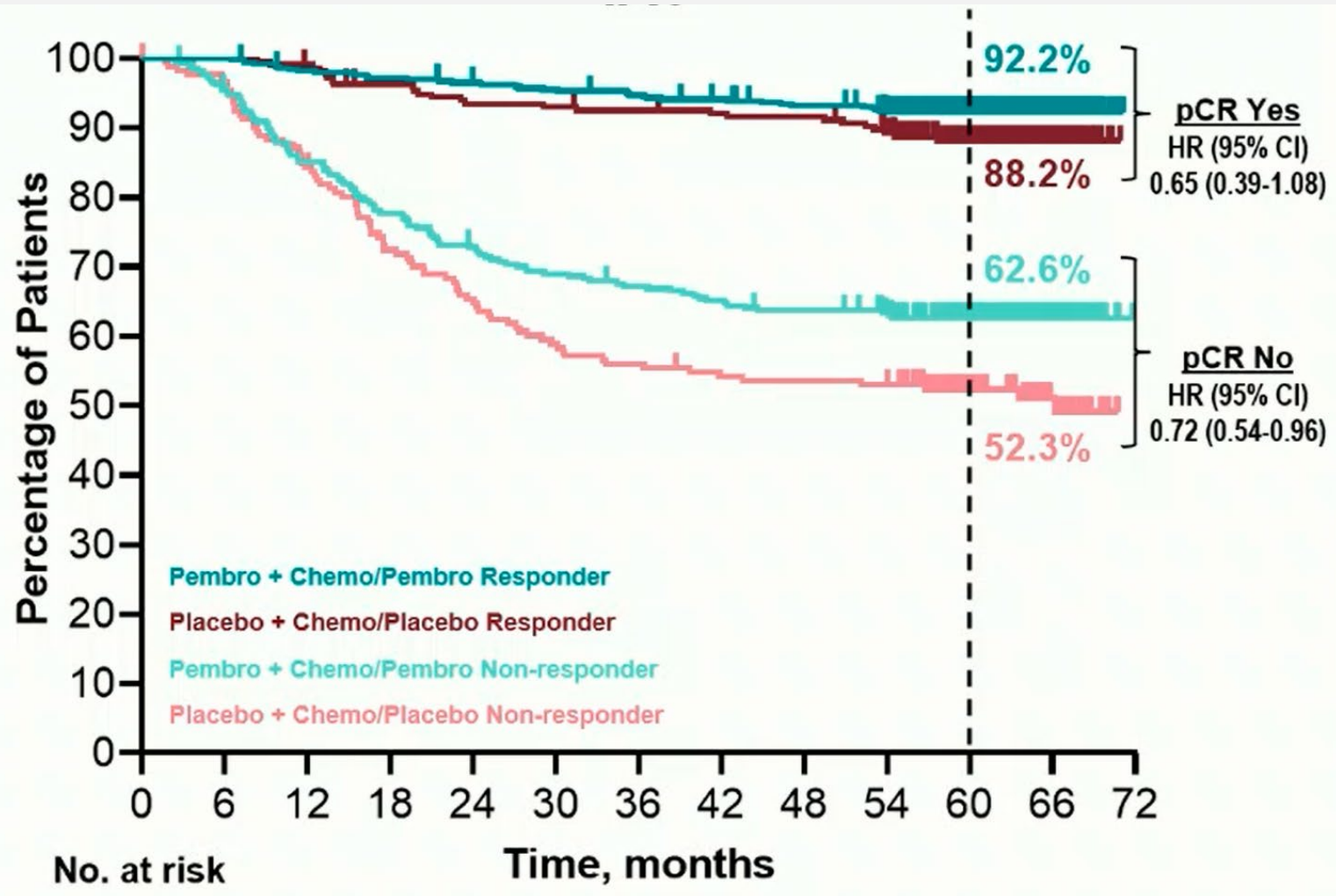
pCR Yes HR = 0.65

pCR No HR = 0.72



Effektgrösse relativ und absolut

effect size



Relativ: Hazard ratio

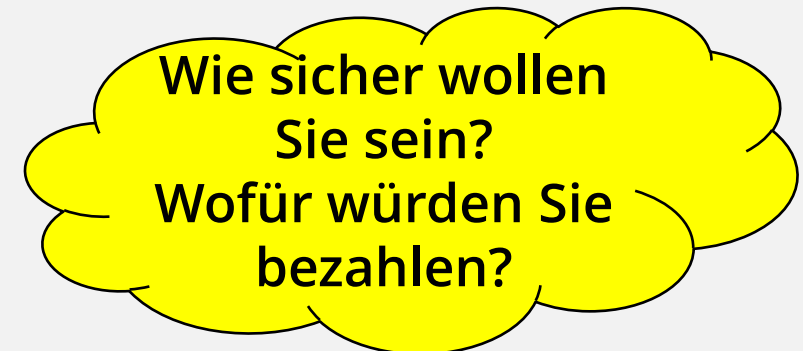
pCR Yes HR = 0.65

pCR No HR = 0.72

Absolut: Δ nach 3 Jahren

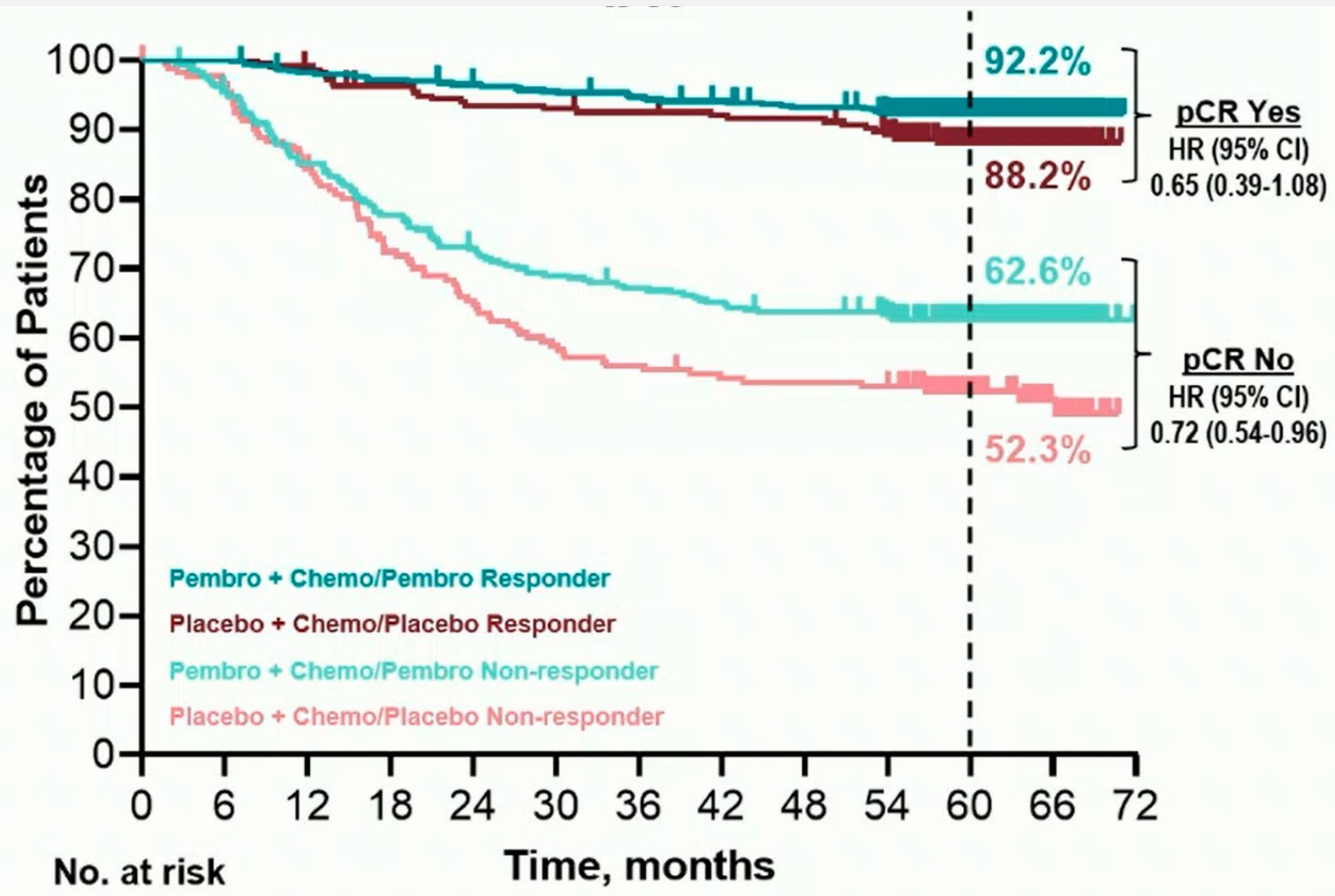
pCR Yes $\Delta_{\bullet} = 4\%$

pCR No $\Delta_{\bullet} = 10\%$



Effektgrösse relativ und absolut

effect size



Relativ: Hazard ratio

pCR Yes HR = 0.65

pCR No HR = 0.72

Absolut: Δ nach 3 Jahren

pCR Yes $\Delta = 4\%$

pCR No $\Delta = 10\%$

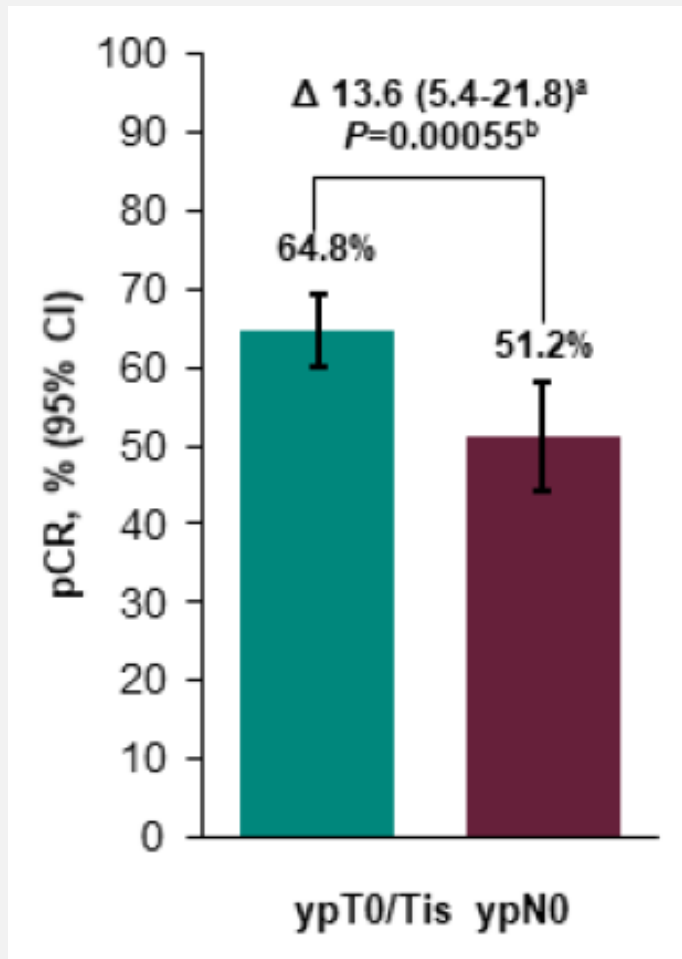
Number needed to treat

pCR Yes $1/0.04 = 25$

pCR No $1/0.10 = 10$

Effektgrösse relativ und absolut

effect size



Absolut

$\Delta = 13.6\%$

Relativ

Odds ratio = 1.75

Relatives Risiko = 1.27

Odds Pembro = $0.648/(1-0.648) = 1.84$

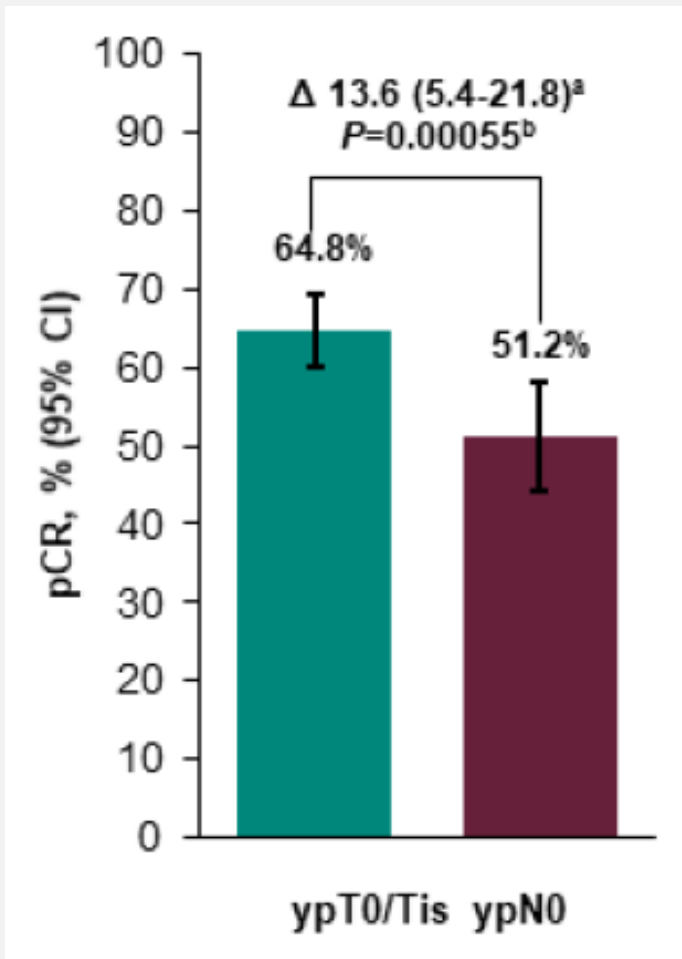
Odds Placebo = $0.512/(1-0.512) = 1.05$

Odds ratio = $1.84/1.05 = 1.75$

Relatives Risiko = $0.648/0.512 = 1.27$

Effektgrösse relativ und absolut

effect size



Absolut

$\Delta = 13.6\%$

Relativ

Odds ratio = 1.75

Relatives Risiko = 1.27

«Fast doppelt so gute Chance einer pCR!»

«27% bessere Chance einer pCR!»

Vertrauensintervalle

confidence intervals

- ◆ Messgrößen der Studiengruppe werden als Stichprobe der Messgrößen einer Population betrachtet.
- ◆ Die Messgröße der gesamten Population hat einen wahren Mittelwert
- ◆ Das 95%-Vertrauensintervall wird so berechnet, dass es für eine grosse Zahl Stichproben bei 95% den Mittelwert* der Population enthält

*Mittelwert als Beispiel, ebenso möglich sind andere Masszahlen wie odds ratio, hazard ratio, Überleben nach einer bestimmten Zeit

Vertrauensintervalle

confidence intervals

- ◆ Messgrößen der Studiengruppe werden als Stichprobe der Messgrößen einer Population betrachtet.
- ◆ Die Messgröße der gesamten Population hat einen wahren Mittelwert
- ◆ Das 95%-Vertrauensintervall wird so berechnet, dass es für eine grosse Zahl Stichproben bei 95% den Mittelwert* der Population enthält

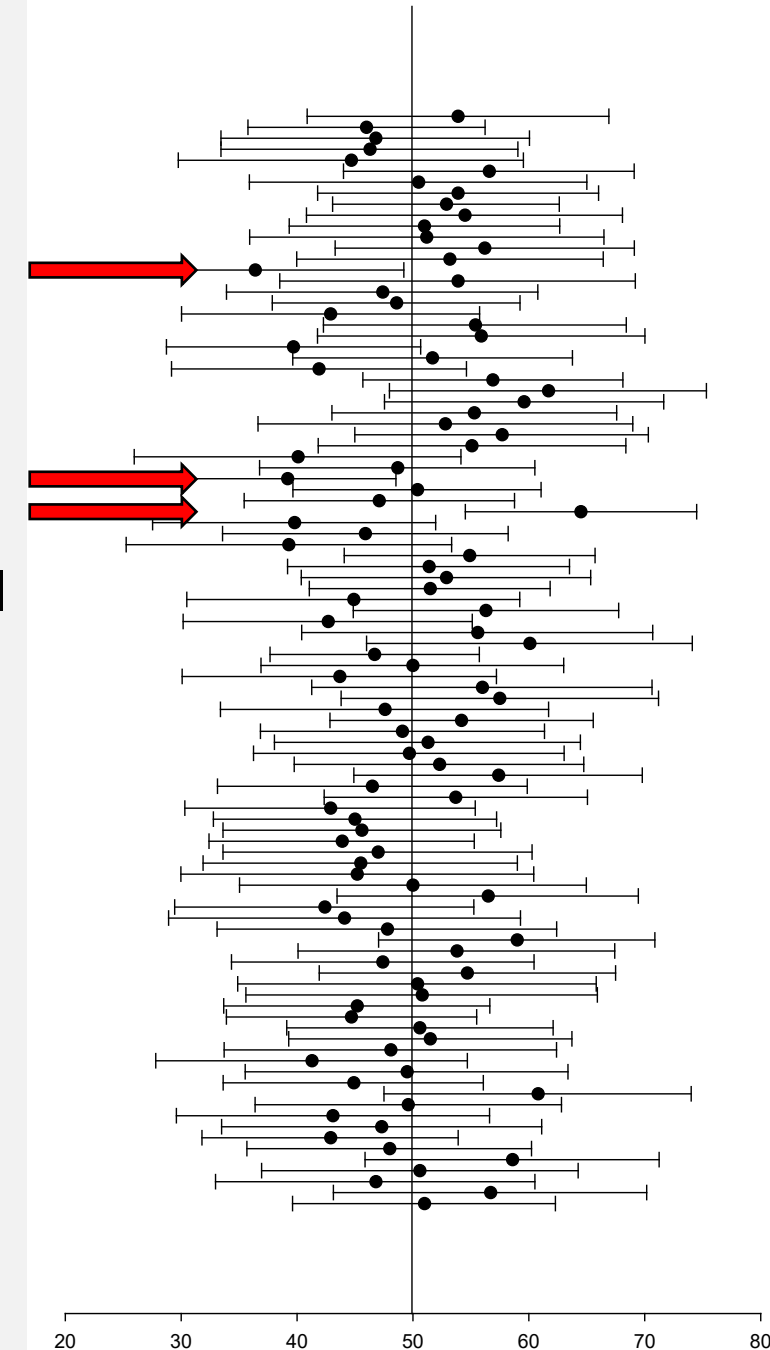
*Mittelwert als Beispiel, ebenso möglich sind andere Masszahlen wie odds ratio, hazard ratio, Überleben nach einer bestimmten Zeit

2000 Zufallszahlen
zwischen 0 und 100

Erwartetes
arithmetisches Mittel
= 50

100 Stichproben zu
20 Zufallszahlen

Arithmetisches
Mittel der
Stichproben und
95%
Vertrauensintervall

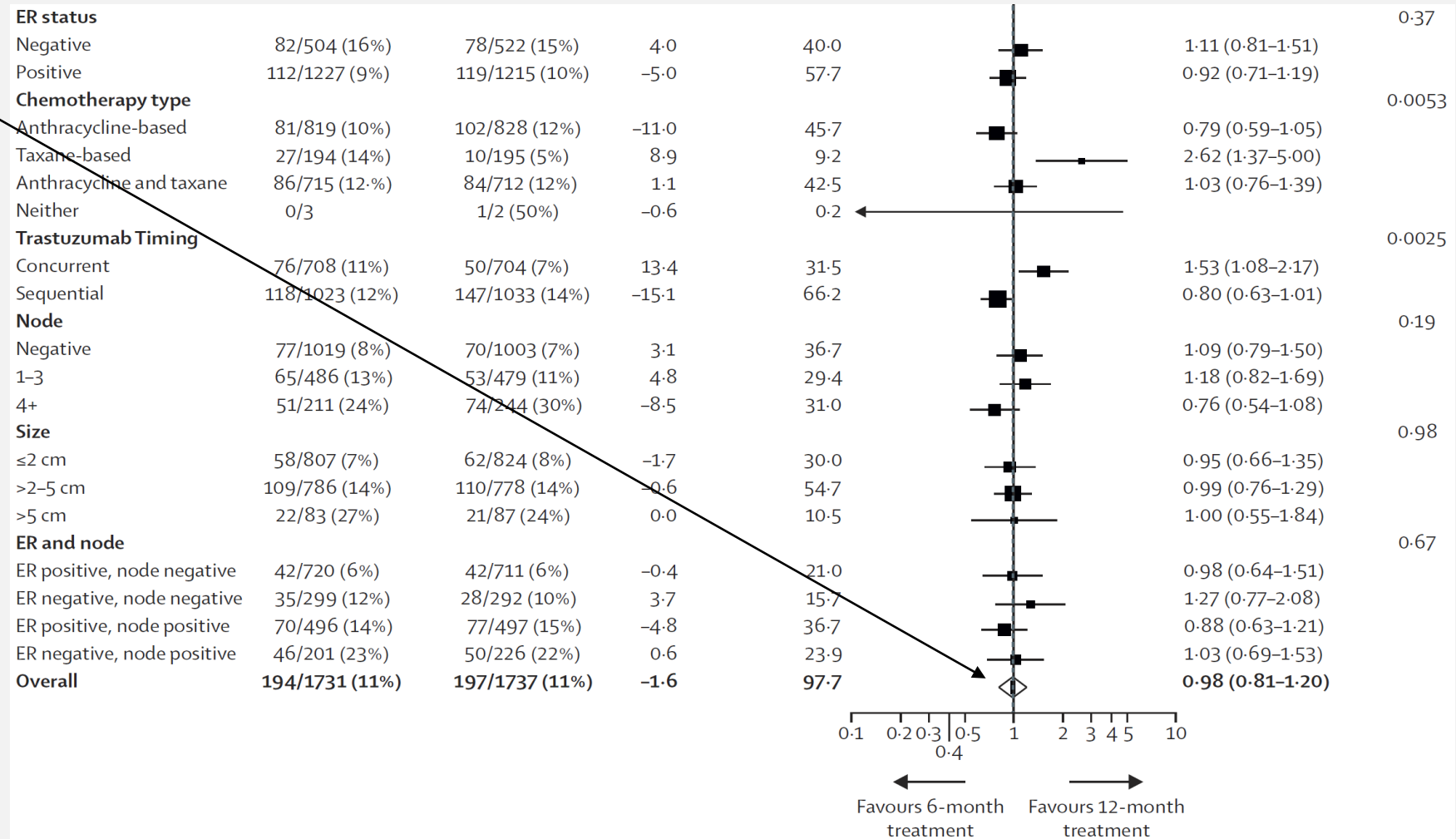


Vertrauensintervalle

confidence intervals

Wenn das 95% Vertrauensintervall der Gesamtwirkung die Null-Effekt-Linie nicht überlappt, ist die Gesamtwirkung unterschiedlich von null Effekt

*als Beispiel; rechts Hazard Ratio, M=1

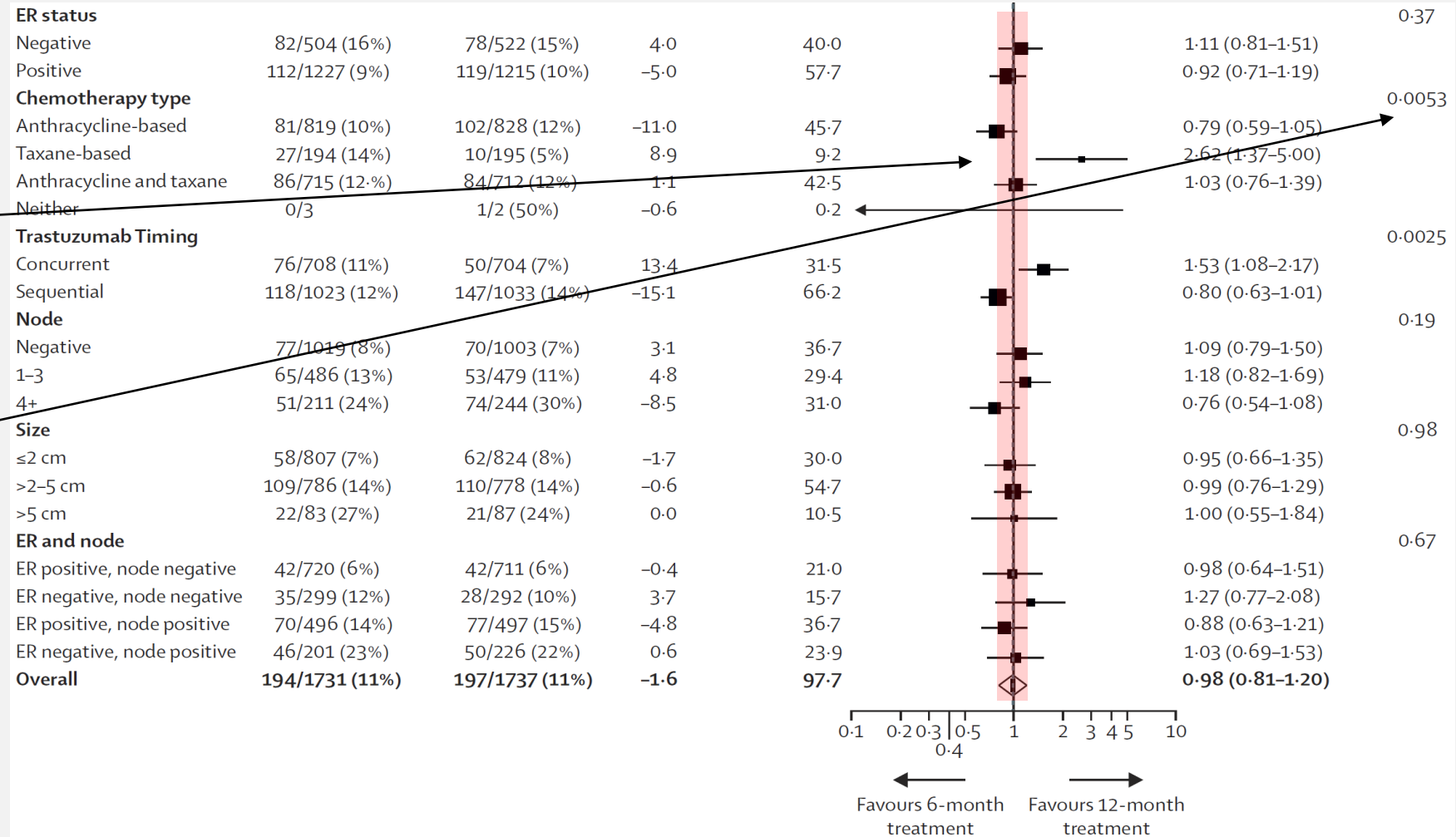


Vertrauensintervalle

confidence intervals

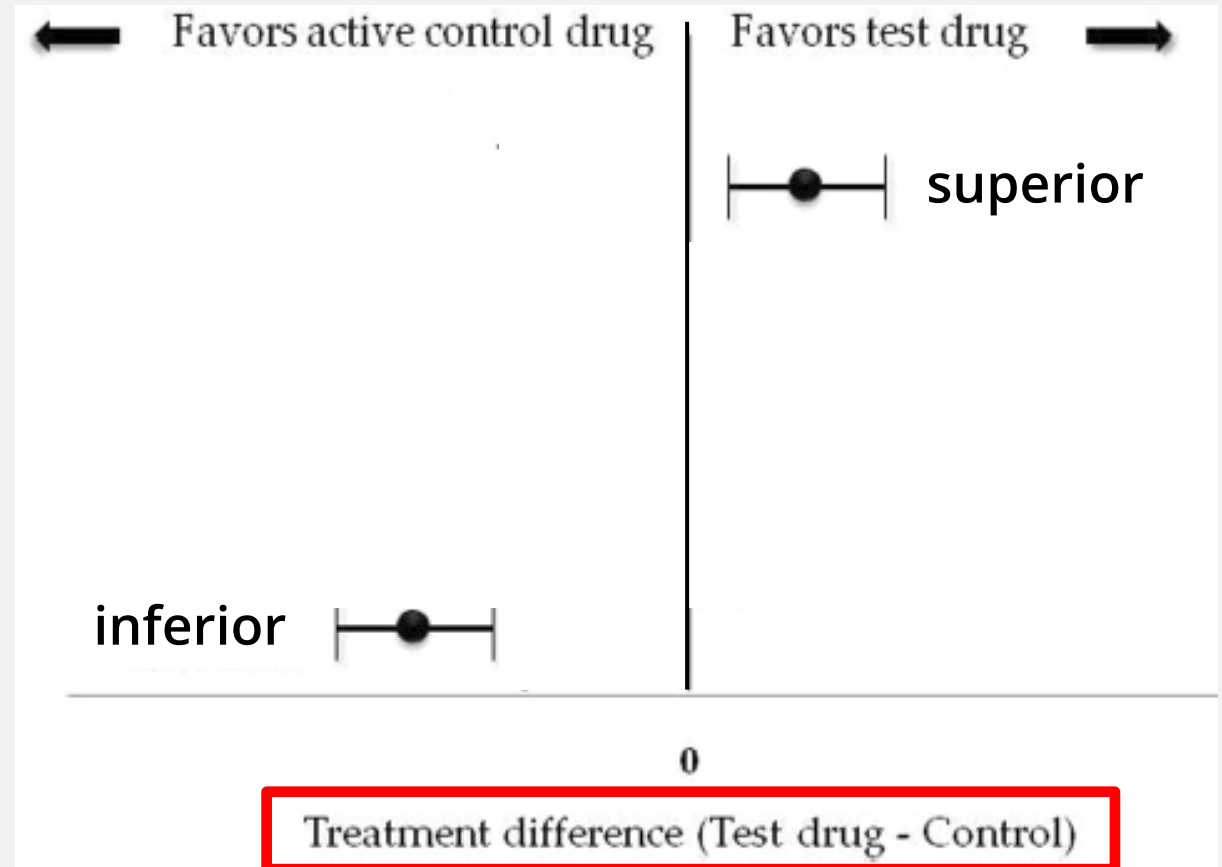
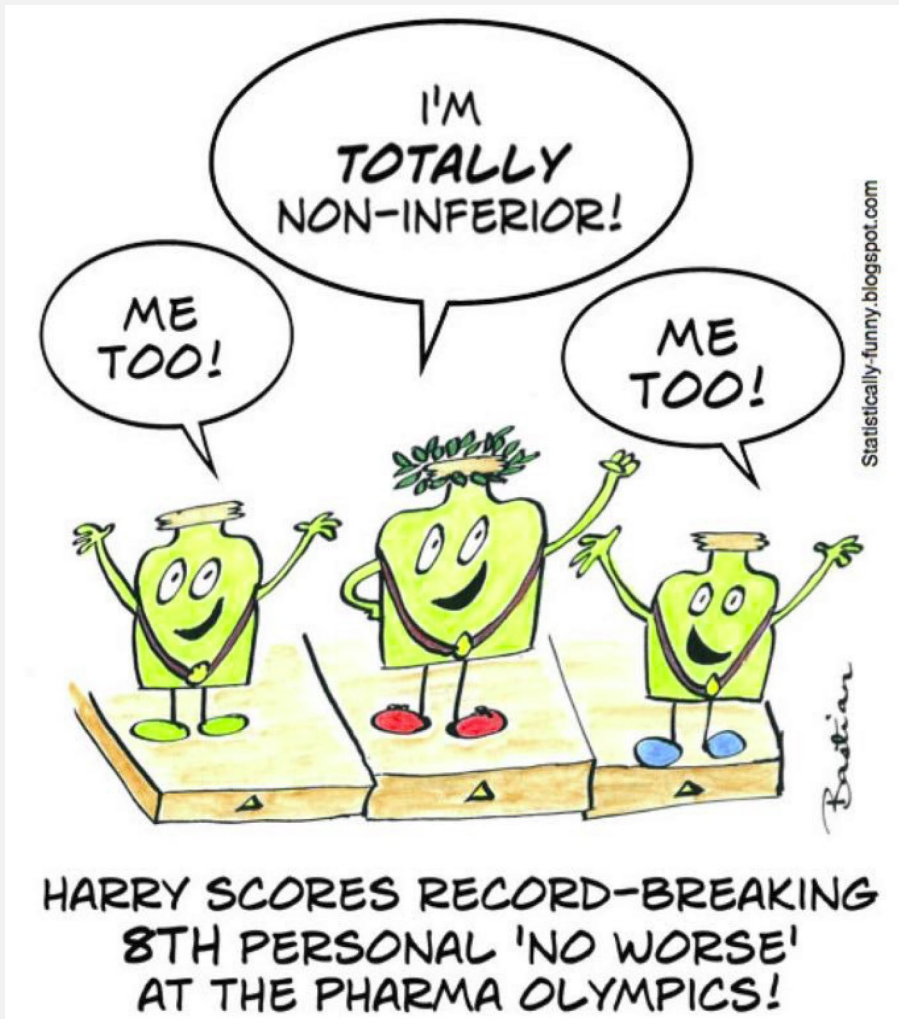
ABER: Achten Sie in Forest Plots auf

- den **Haupteffekt**
- Abweichungen vom Haupteffekt in Untergruppen
- Wechselwirkungen = Interaktionen zwischen Subgruppen und Behandlung



Überlegenheit und Nicht-Unterlegenheit

superiority, non-inferiority

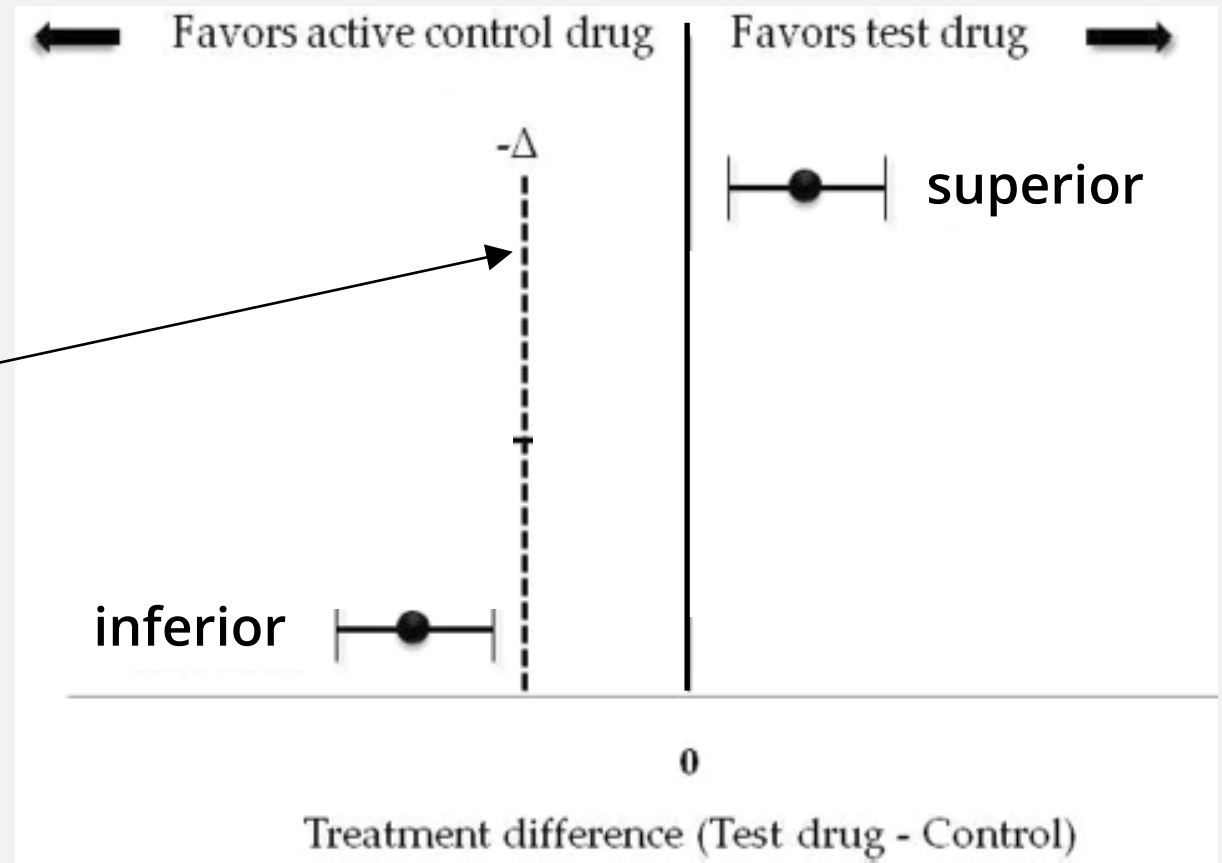


Überlegenheit und Nicht-Unterlegenheit

superiority, non-inferiority

Voraussetzungen einer Nicht-Unterlegenheits-Studie

1. **Aktive Kontrolle**
2. Relevante Messgrösse
3. **Nicht-Unterlegenheits-Grenze**
4. Absolute oder relative Formulierung der Grenze
5. Messgenauigkeit

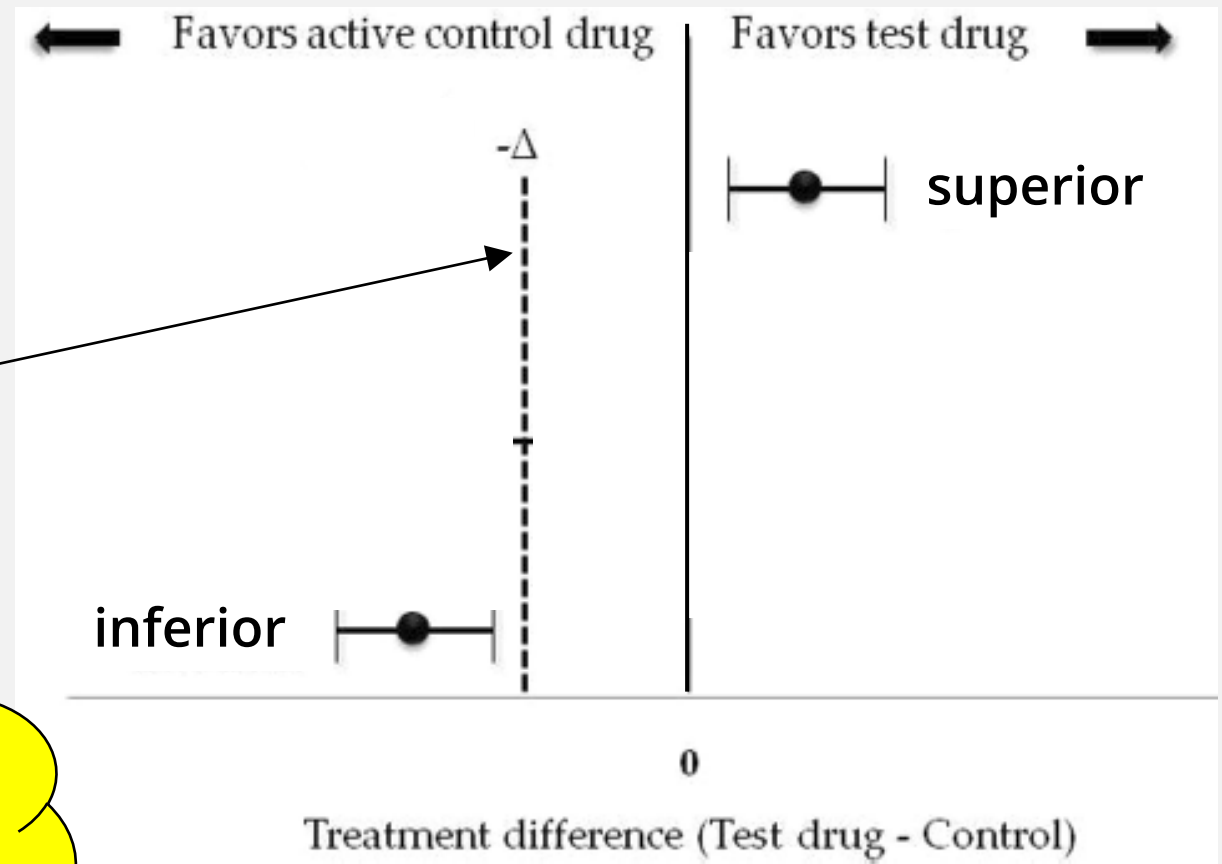


Überlegenheit und Nicht-Unterlegenheit

superiority, non-inferiority

Voraussetzungen einer Nicht-Unterlegenheits-Studie

1. Aktive Kontrolle
2. Relevante Messgrösse
3. **Nicht-Unterlegenheits-Grenze**
4. Absolute oder relative Formulierung der Grenze
5. Messgenauigkeit



Wie viel schlechter gilt noch als «nicht-unterlegen»?

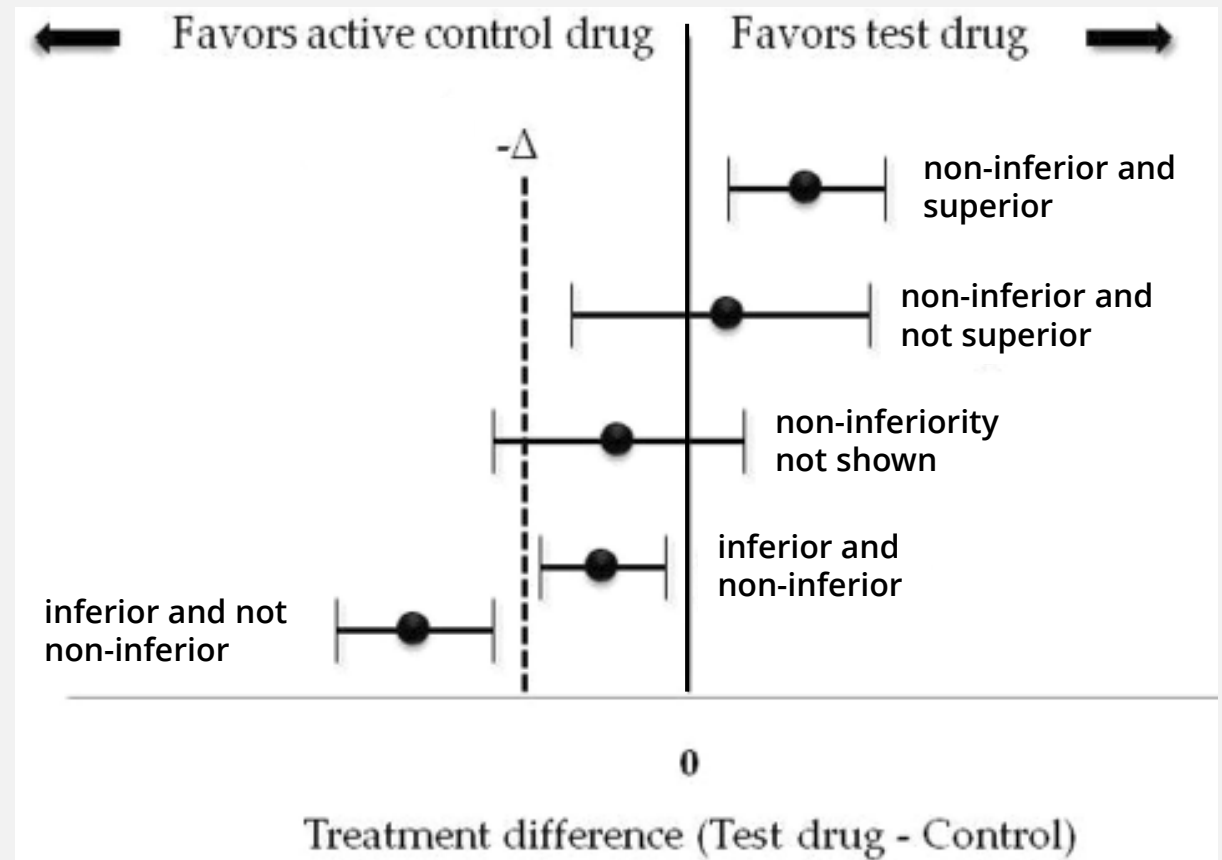
Überlegenheit und Nicht-Unterlegenheit

superiority, non-inferiority

Voraussetzungen einer Nicht-Unterlegenheits-Studie

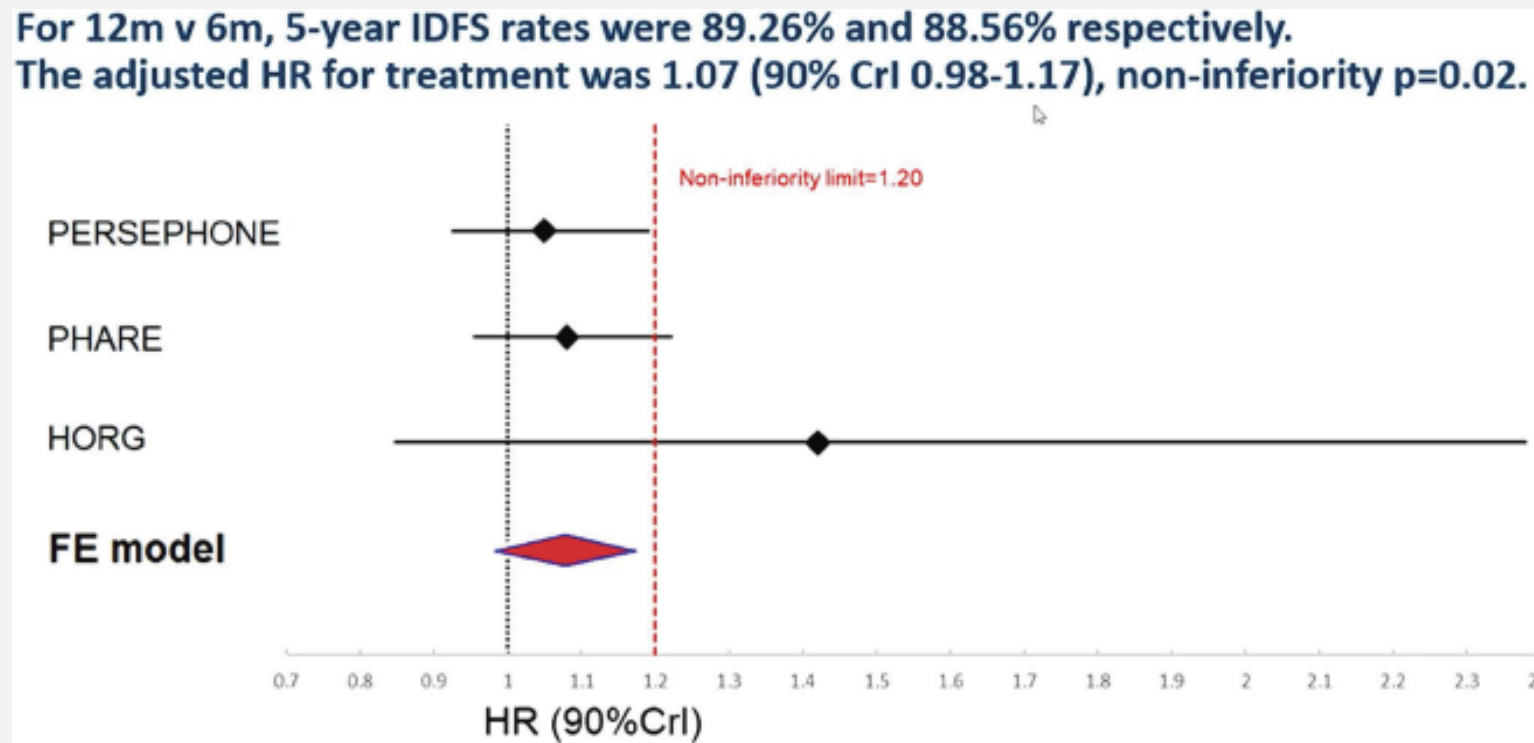
1. Aktive Kontrolle
2. Relevante Messgrösse
3. **Nicht-Unterlegenheits-Grenze**
4. Absolute oder relative Formulierung der Grenze
5. Messgenauigkeit

Nicht-Unterlegenheit: das 95% Vertrauensintervall der Messgrösse schneidet die Nicht-Unterlegenheits-Grenze nicht.



Nicht-Unterlegenheit

- ◆ Adjuvante Therapie mit Trastuzumab, 12 vs. 6 (oder 3) Monate
- ◆ Nicht-Unterlegenheits-Grenze für 5-Jahres-IDFS: **2% (absolut)**
«making sure that a little bit worse is not too much worse»



Nicht-Unterlegenheit – was kriegen wir dafür?

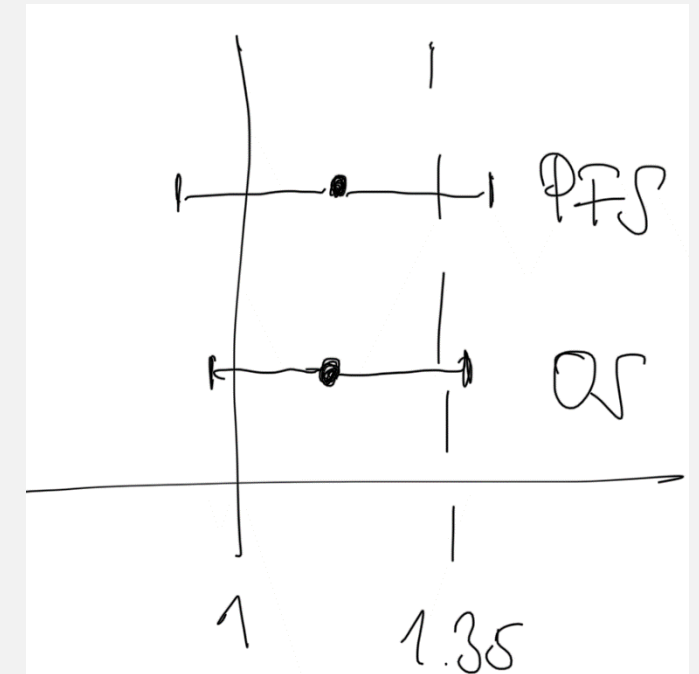
Nutzen in anderen Domänen

Vergleich	Endpunkt	Unerwünschte Wirkungen	Zeitliche Belastung	Kosten	Fortschritt oder «Me too»?
HCC Lenvatinib vs. Sorafenib	OS	--	--	--	Me too
Mamma Chemoendokrin vs. endokrin TAILORx	IDFS	↓	↓	↓	Fortschritt

Nicht-Unterlegenheit – aktive Kontrolle

- ◆ Docetaxel + Capecitabin ist überlegen gegenüber Docetaxel
(Capecitabin $2 \times 1250 \text{ mg/m}^2 \times 14/21$ Tage)
O'Shaughnessy J et al. J Clin Oncol. 2002;20(12):2812-23
- ◆ Docetaxel + Capecitabine ($2 \times 825 \text{ mg/m}^2$) ist nicht-unterlegen gegenüber
Docetaxel + Capecitabine ($2 \times 1250 \text{ mg/m}^2$)
Buzdar AU et al. Ann Oncol. 2012;23(3):589-97
- ◆ Nicht-Unterlegenheits-Grenze für $HR=1.35$,
95% VI schneiden 1 und 1.35
Weder unterlegen noch nicht-unterlegen!

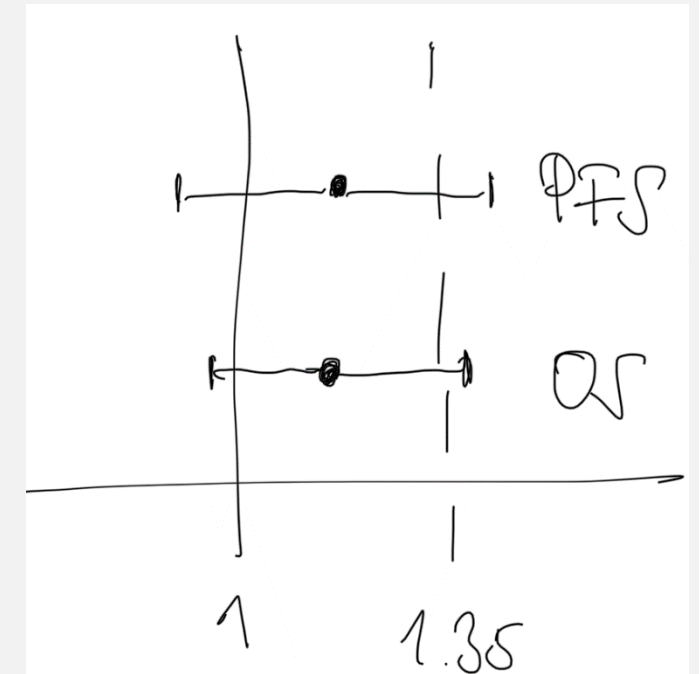
Problem ?



Nicht-Unterlegenheit – aktive Kontrolle

- ◆ Docetaxel + Capecitabin ist überlegen gegenüber Docetaxel
(Capecitabin $2 \times 1250 \text{ mg/m}^2 \times 14/21$ Tage)
O'Shaughnessy J et al. J Clin Oncol. 2002;20(12):2812-23
- ◆ Docetaxel + Capecitabine ($2 \times 825 \text{ mg/m}^2$) ist nicht-unterlegen gegenüber
Docetaxel + Capecitabine ($2 \times 1250 \text{ mg/m}^2$)
Buzdar AU et al. Ann Oncol. 2012;23(3):589-97
- ◆ Nicht-Unterlegenheits-Grenze für $HR=1.35$,
95% VI schneiden 1 und 1.35
Weder unterlegen noch nicht-unterlegen!

Problem: Ist die niedrige Dosis aktiv, wirksamer als nichts?



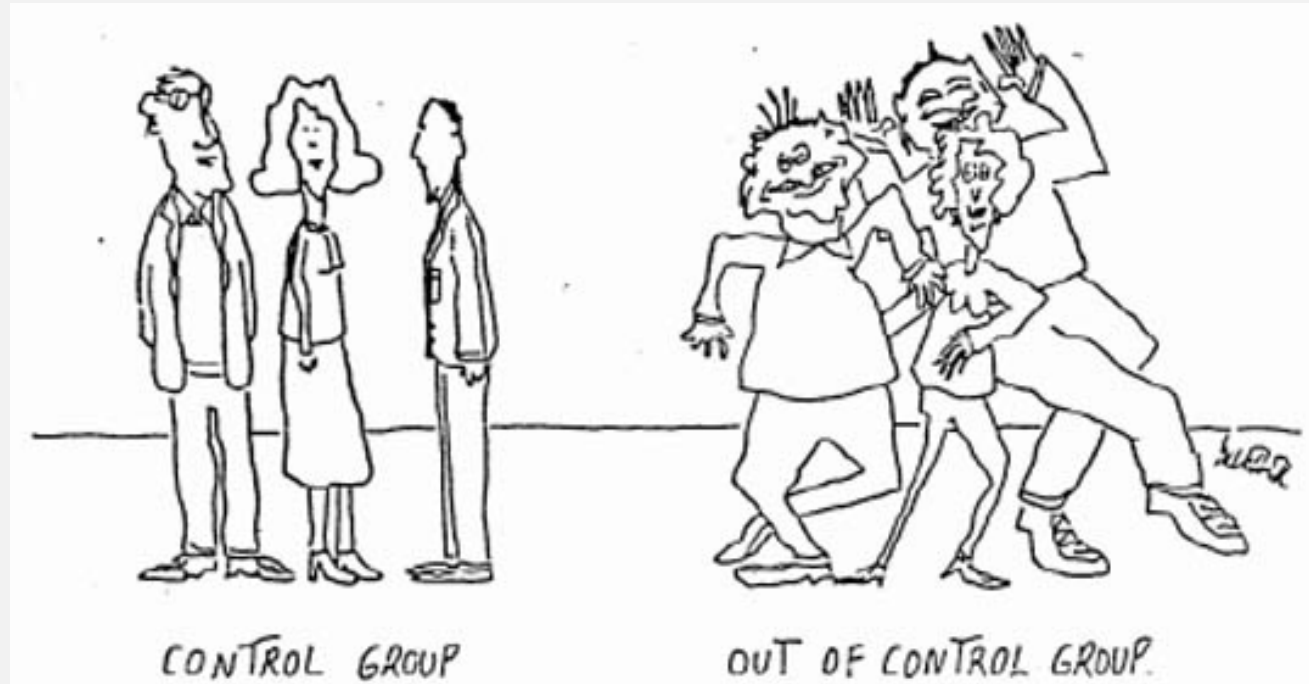
Vergleichstherapie

comparator therapy

Die experimentelle, 'neue', Therapie muss mit der besten etablierten Verglichen werden.

Endpunkt Gesamtüberleben

Wird eine Behandlung, die in einer späteren Therapie-'Linie' wirksam war, in eine frühere verschoben, muss für Teilnehmer der Vergleichsgruppe die Anwendung in der späteren Linie möglich sein.

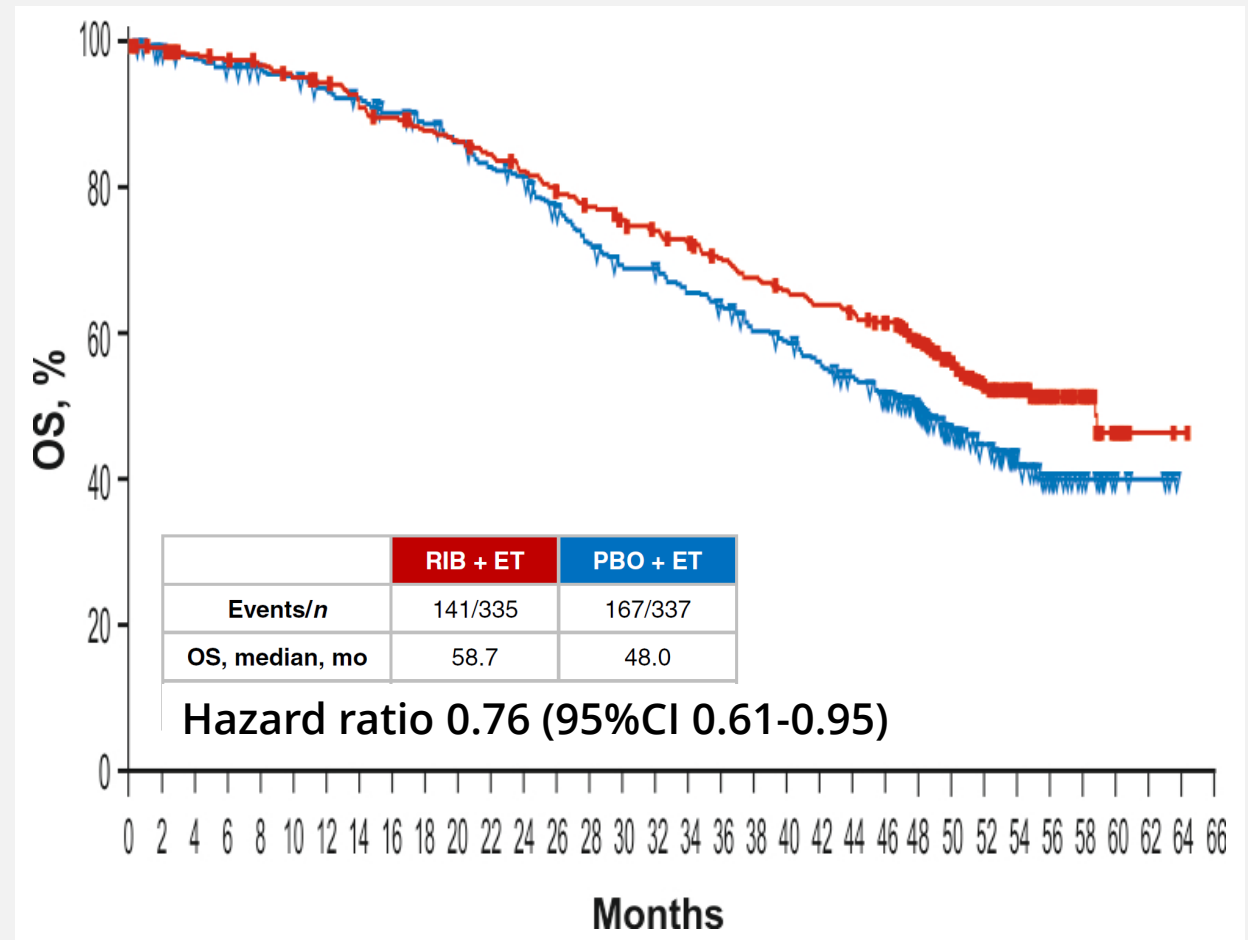


Vergleichstherapie

comparator therapy

MONALEESA-7 – Mammakarzinom

- ◆ Prämenopausal, ER+, HER2– M1 (41% primär)
- ◆ Erste Linie: GnRH-Analog + Aromatase-Inhibitor (oder Tamoxifen) + Ribociclib oder Placebo

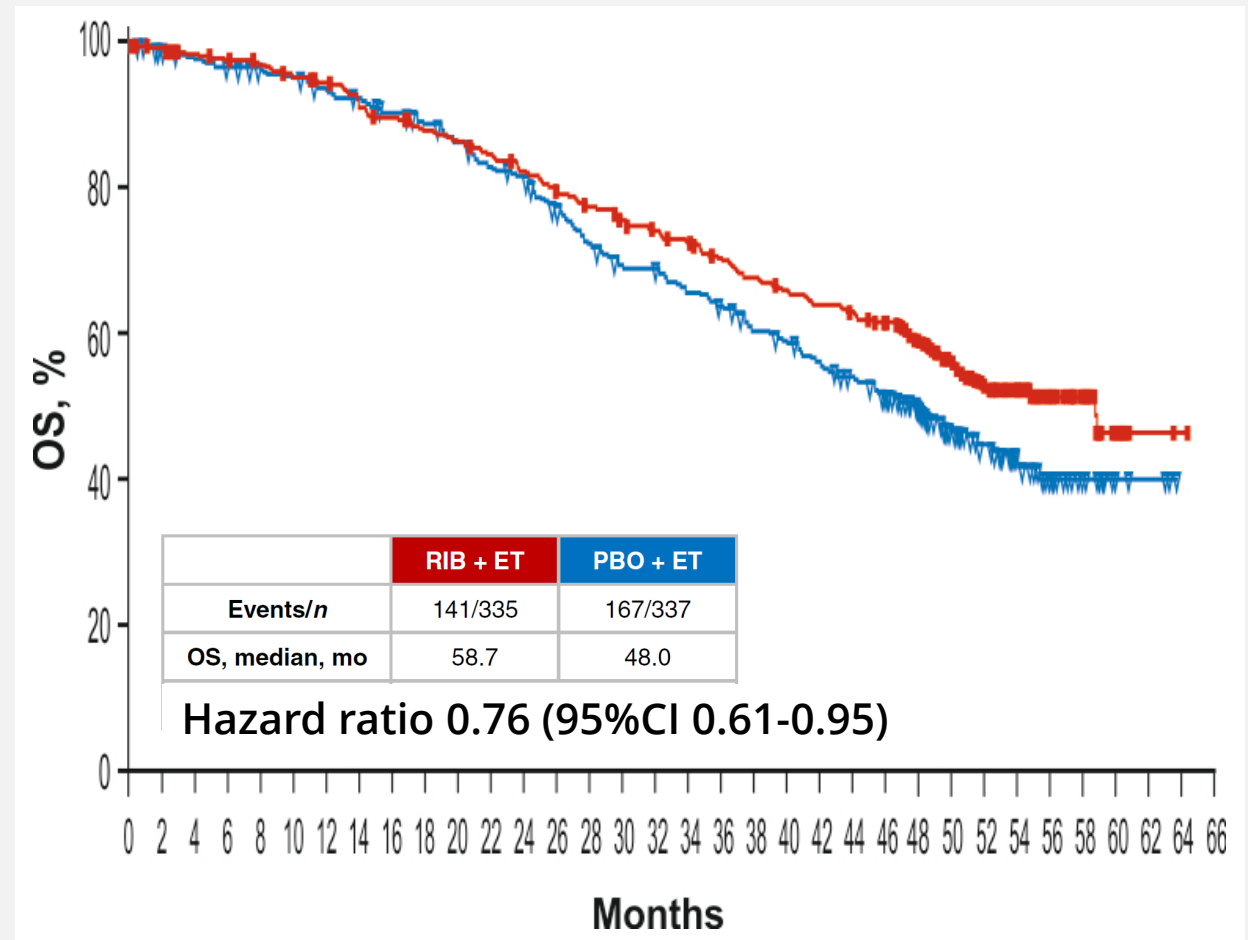


Vergleichstherapie

comparator therapy

MONALEESA-7 – Mammakarzinom

- ◆ Prämenopausal, ER+, HER2– M1 (41% primär)
- ◆ Erste Linie: GnRH-Analog + Aromatase-Inhibitor (oder Tamoxifen) + Ribociclib oder Placebo
- ◆ Therapie nach Progression in der Placebo-Gruppe
 - CDK4/6-Inhibitor 19%
 - Irgendeine 73%

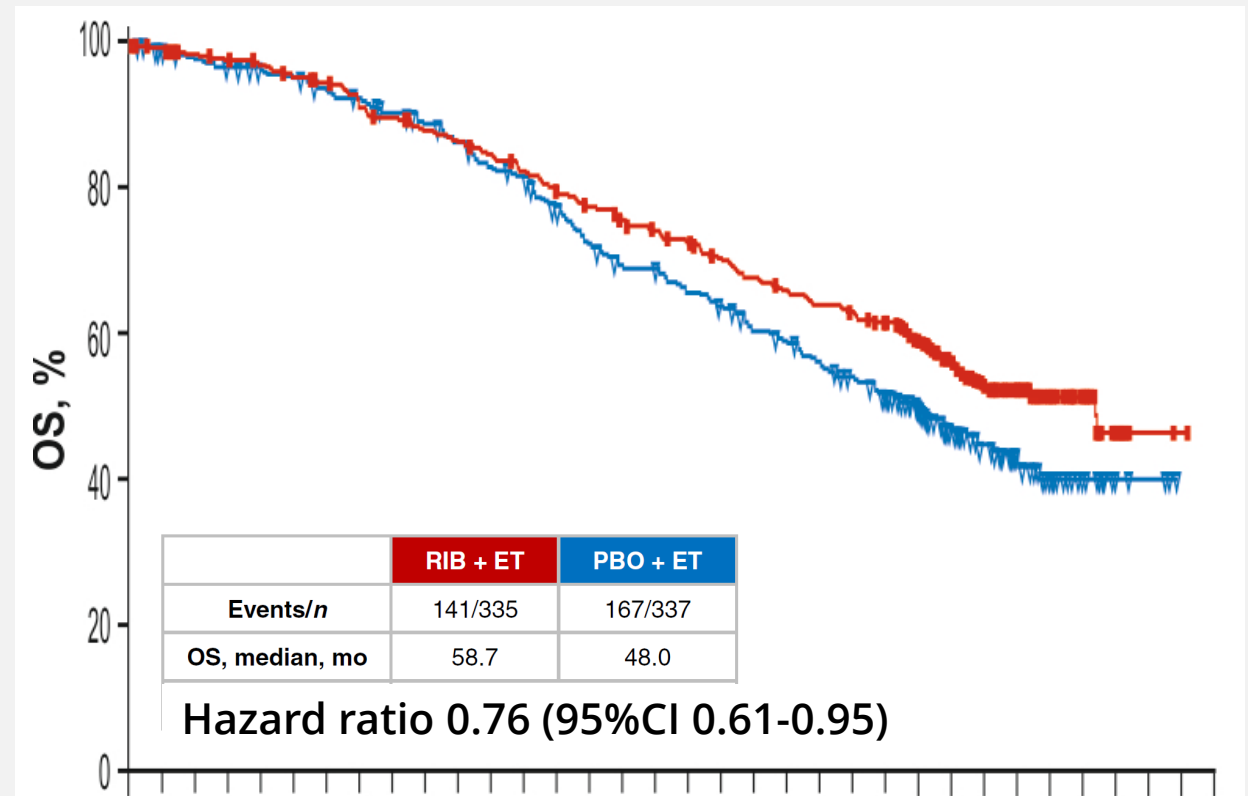


Vergleichstherapie

comparator therapy

MONALEESA-7 – Mammakarzinom

- ◆ Prämenopausal, ER+, HER2- M1 (41% primär)
- ◆ Erste Linie: GnRH-Analog + Aromatase-Inhibitor (oder Tamoxifen) + Ribociclib oder Placebo
- ◆ Therapie nach Progression in der Placebo-Gruppe
 - CDK4/6-Inhibitor 19%
 - Irgendeine 73%



Medianes Alter: 45 (29-58), ECOG PS 0+1: 76%+23%

Erhielten die Patientinnen der Vergleichsgruppe die beste Standardtherapie?

Validität – intern und extern

internal and external validity

Interne Validität: Die Schlüsse der Studie trifft zu auf die Population der Studie; Aussage über die «Strenge» der Studien-durchführung

- + Studienprotokoll
- + Repräsentative Population
- + Randomisation, Verblindung
- Nicht-zufälliges Ausscheiden (informative censoring)
- Folgetherapien

Wie gut wurde die Studie geplant und realisiert?

Externe Validität: Die Schlüsse der Studie können generalisiert werden über die Studienpopulation hinaus; Aussage über die Anwendbarkeit der Resultate

- + Pragmatische Auswahl der Teilnehmerinnen
- + Kulturell, ethnisch, sozial repräsentative Teilnehmerinnen
- + Vergleichstherapie nach aktuellem Stand des Wissens

Wie anwendbar sind die Resultate im Alltag («real world»)?

Validität – intern und extern

internal and external validity

Beispiel: Gesamtüberleben in MONALEESA-7

Intern ✓ Standard randomisierte Phase-3-Studie	Extern ∅ 40% Stadium IV bei Diagnose junge Patientinnen in gutem AZ, ¼ ohne Therapie der 2. Linie Nur 18% Cross-over
Ribociclib ist sicher aktiv	Widerspiegelt nicht die Praxis in der Schweiz

Arrog (Great World)

CALLING BULLSHIT

Carl Bergstrom & Jevin West

Calling Bullshit
*The Art of Skepticism
in a Data-Driven
World*



A PELICAN BOOK

The Art of Statistics
Learning from Data
David Spiegelhalter

'This marvellous book will transform your relationship with the numbers that swirl all around us'

TIM HARFORD

INTERNATIONAL EDITION

Not authorised for sale in United States, Canada, Australia, New Zealand, Puerto Rico or the U.S. Virgin Islands

Clinical Epidemiology The Essentials

Sixth Edition



Grant S. Fletcher

 Wolters Kluwer

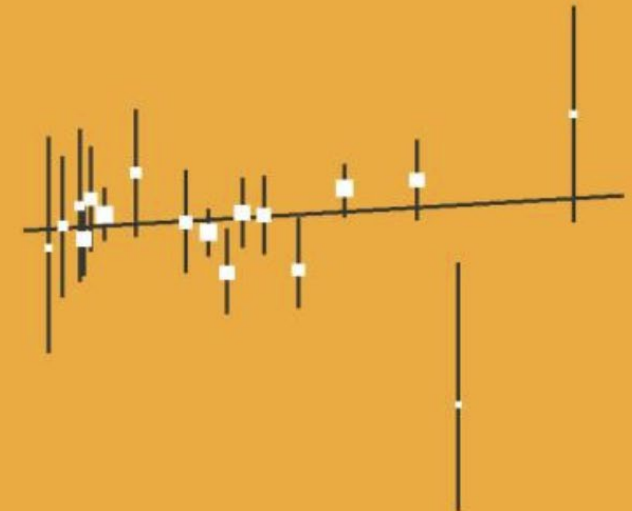
 Wolters Kluwer

Copyrighted Material

OXFORD

4TH EDITION

AN INTRODUCTION TO MEDICAL STATISTICS



MARTIN BLAND

Copyrighted Material

<https://statistically-funny.blogspot.com/>

$$\lim_{x \rightarrow \infty} \left(x + \frac{1}{x^2-3}\right) \frac{3\pi \cos^{-1}x}{\theta'} \neq \frac{4 \ln}{\theta'}$$

**Vielen Dank für Ihre Verbesserungs- und
Ergänzungsvorschläge!**

stefan.aebi@onkologie.ch