

**Advice to statistical  
consultants contributing and  
evaluating evidence for a  
nutrition intervention**

**GGB Seminar**

**November 17, 2015**

**Daniel J. Tancredi, PhD**

*Associate Professor in Residence of Pediatrics,  
UC Davis*

# Objectives

Review statistical concepts, principles and methods relevant to evaluating evidentiary quality of nutrition studies

# **1. STATISTICAL PRELIMINARIES**

# **1A. P VALUES & CONFIDENCE INTERVALS**

# Which of the following things does a report of $P < 0.05$ allow you to know?

1. The probability that the null hypothesis is true.
2. The probability that the alternative hypothesis is true.
3. The probability that the observed effect is real.
4. The probability that a claim of a positive result is a false positive claim.
5. The probability that the result can be replicated.
6. The strength of the evidence in the data against the null hypothesis.

Sources: Lew MJ. (2012) “Bad statistical practice in pharmacology...: you probably don’t know P”. *BJP* 166: 1559-1567

Haller H, Krauss S (2002) “Misinterpretation of significance: a problem students share with their teachers.” *Methods Psych Res* 7: 1-20

# Definition of a P-value

*Suppose you have*

a null hypothesis and

a method for converting sample data into a test statistic that has the property that extreme values constitute evidence against the null hypothesis.

## Example

**$H_0$ : Mean of outcome is equal in two comparison groups**

**Test statistic is between-group difference in sample means**

# Definition of a P-value (*cont.*)

*Only then* can you define the p-value associated with the value for the test statistic observed in the given sample.

The p-value is the conditional probability, under (the data generating model associated with) the **null hypothesis**, of obtaining a value for the test statistic that is as least as extreme as **the observed value in the sample**.

# P-values and “significance testing”

- R.A. Fisher promoted the P-value as a measure of the strength of the evidence within the observed data against a null hypothesis and introduced the word “significant”
- Fisher’s rivals Jerzy Neyman and Egon Pearson introduce an alternative inferential approach that uses
  - long-term error rates,
  - appropriately powered experiments
  - binary decision making



# P-values and “significance testing”

Unfortunately, both approaches use term “*significant*”, leading to confusing hybrid approaches seen in practice (i.e. same paper using  $p < 0.05$ ,  $p < 0.01$ , etc.)

Sources: Fisher RA (1925): Statistical Methods for Research Workers. Oliver and Boyd: Edinburgh.

<http://psychclassics.yorku.ca/Fisher/Methods/>

Neyman J, Pearson ES (1933): On the problem of the most efficient test of statistical hypotheses. Philos Trans R Soc Long A 231: 289-337.

# Distribution of p values

Suppose you want to perform a two-group comparison of means using Student's t-test.

What's the shape of the (theoretical) distribution of the p-values under

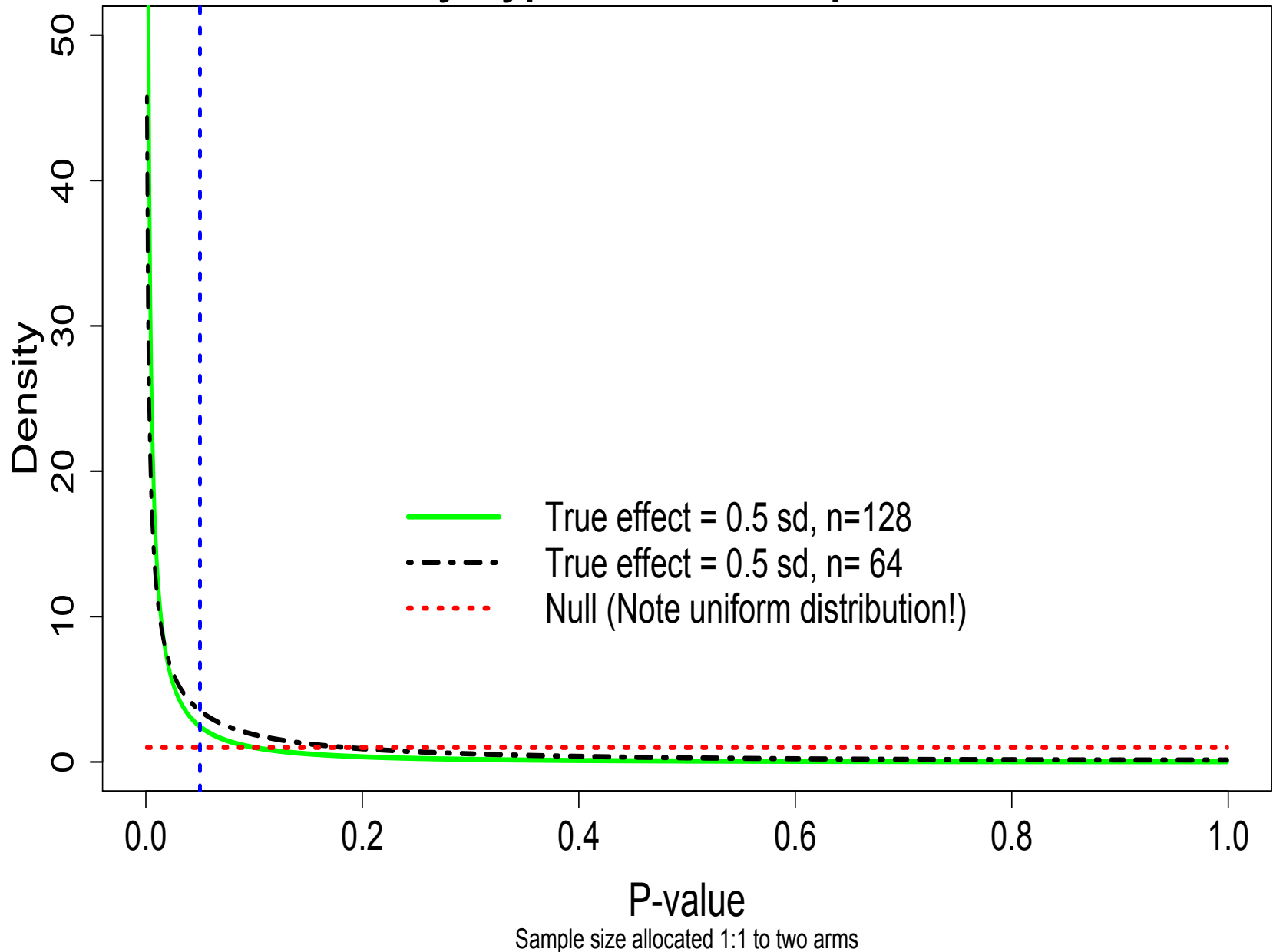
Null hypothesis?

When true effect size (difference in means) is 0.5 standard deviation and power is

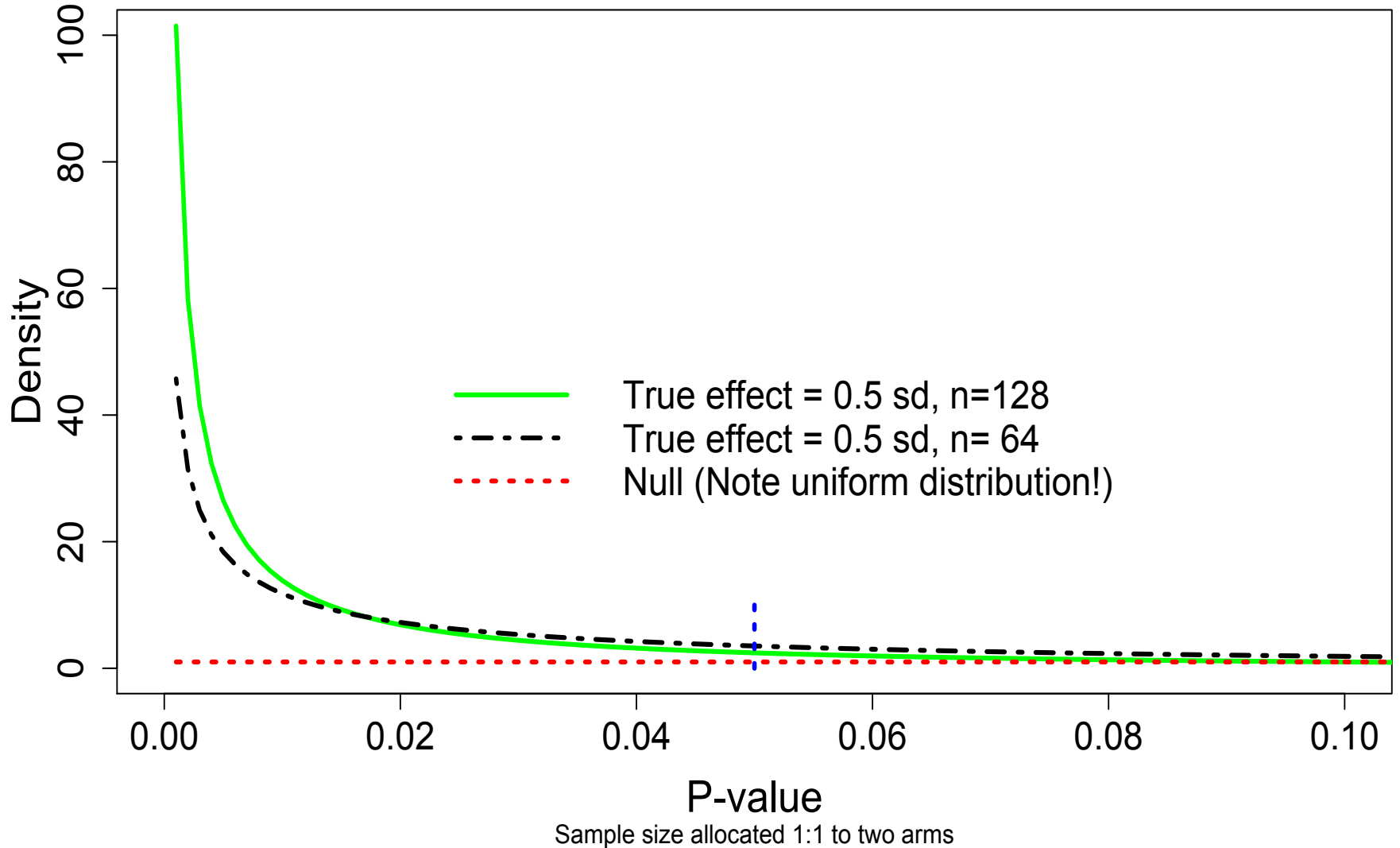
50%?

80%?

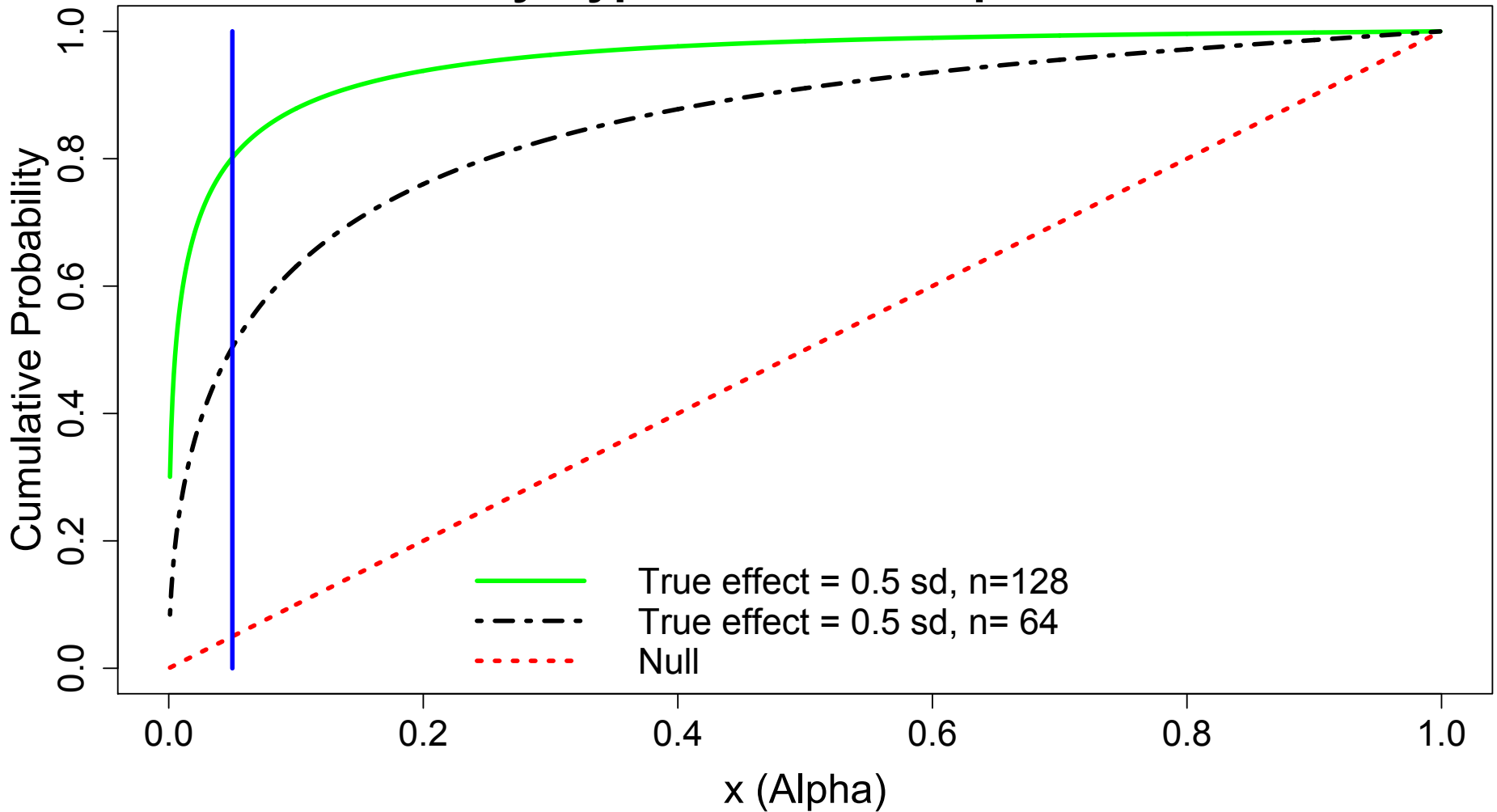
# Density of Student T-test p-values By Hypothesis & Sample size



# Density of Student T-test p-values (below 0.10) By Hypothesis & Sample size



## CDF of Student T-test p-values By Hypothesis & Sample size



Power is value of CDF at given alpha (5%, typically).

# So, which of the following things does a report of $P < 0.05$ allow you to know?

1. The probability that the null hypothesis was true.
2. The probability that the alternative hypothesis was true.
3. The probability that the observed effect was real.
4. The probability that a claim of a positive result is a false positive claim.
5. The probability that the result can be replicated.
6. The strength of the evidence in the data against the null hypothesis.

Sources: Lew MJ. (2012) “Bad statistical practice in pharmacology...: you probably don’t know P”. *BJP* 166: 1559-1567

Haller H, Krauss S (2002) “Misinterpretation of significance: a problem students share with their teachers.” *Methods Psych Res* 7: 1-20

# Problems with p-values and so-called Null Hypothesis Significance Testing

1. Failing to reject  $H_0$  is not proof that  $H_0$  is true (“absence of evidence is not evidence of absence”).
2. P value is very likely to be quite different if experiment is repeated, particularly for underpowered (most!) studies
3.  $H_0$  is almost never true (strictly), anyway. As  $n$  grows, so does probability of rejecting  $H_0$ .
- 4. P value does not give an estimate of the effect size.**
- 5. P value does not give information on precision.**

Sources: Cumming G (2008) “Replication and p Intervals: p Values predict the future Only Vaguely but Confidence Intervals do Much Better”. *Persp of Bio Sci* 3:286-300

Tressoldi PE et al (2013) “High Impact = High Statistical Standards? Not Necessarily So”. *PLOS ONE* 8:e56180

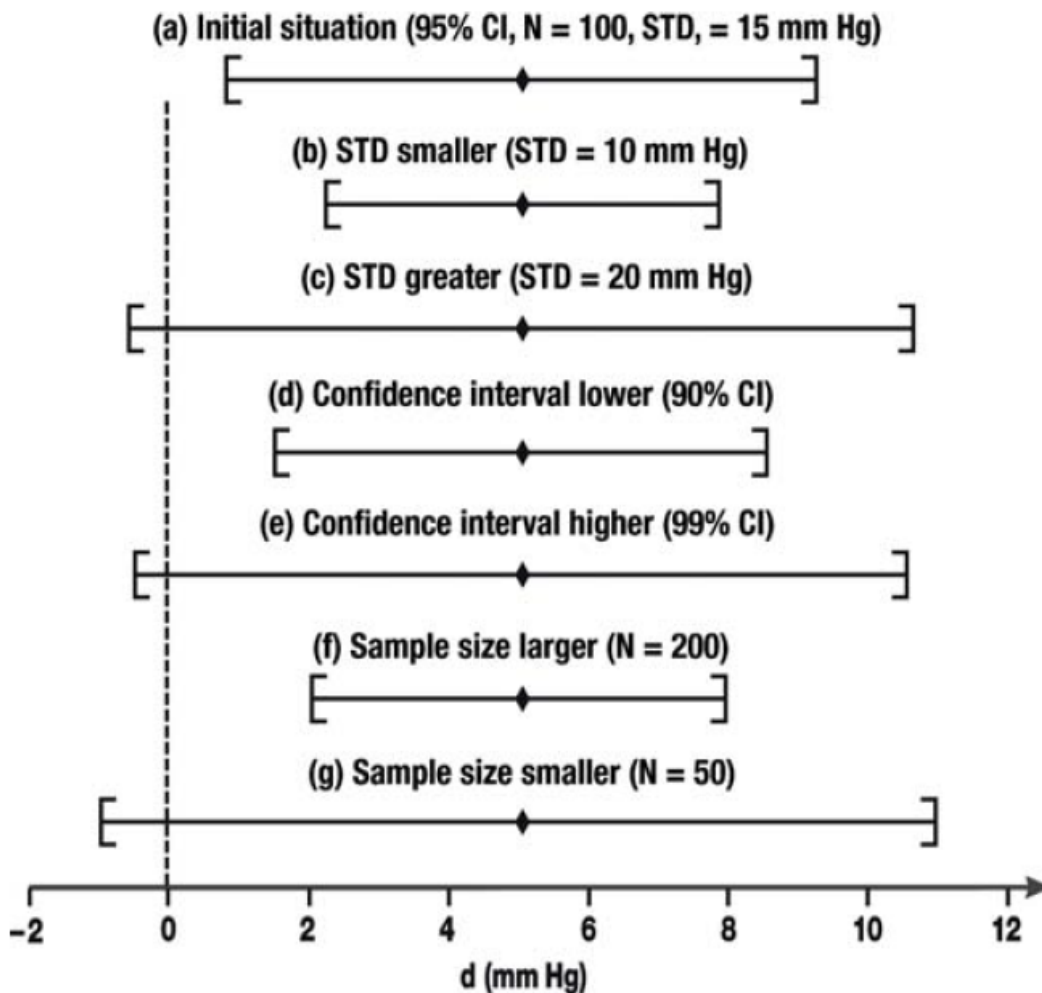
# Confidence intervals rather than P values: estimation rather than hypothesis testing

Well known article by Martin Gardner and Doug Altman [[Br Med J \(Clin Res Ed\)](#). 1986 Mar 15;292(6522):746-50] led to recent high-profile recommendations from CONSORT, APA, & ICMJE ([http://www.icmje.org/manuscript\\_1prepare.html](http://www.icmje.org/manuscript_1prepare.html)), like this:

“...When possible, quantify findings and present them with appropriate indicators of measurement error or uncertainty (such as confidence intervals). Avoid relying solely on statistical hypothesis testing, such as *P* values, which fail to convey important information about effect size...”



**FIGURE 1**



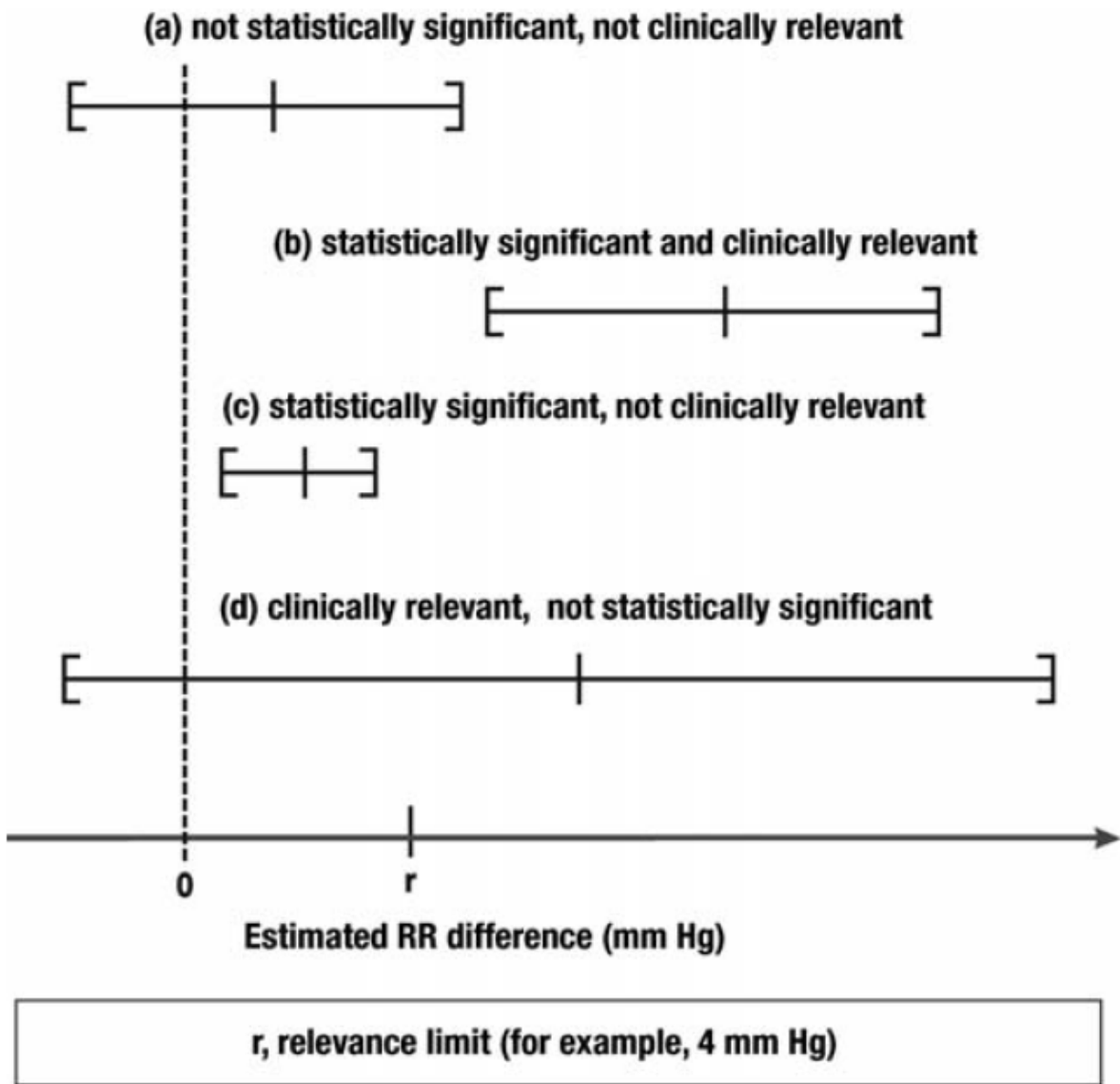
CI, confidence interval; N, sample size per group;  
STD, standard deviation per group as  
measure of variability (same in both groups);  
d, difference in the systolic pressure  
between the two groups

The confidence interval is a range of values with the property that it includes the true value of the parameter with a probability defined in advance. (The probability is a property of the procedure used to convert sample data into interval estimates.)

Describes hypothesized values of the true parameter that would be considered “plausible”.

*Source: (see next slide)*

**FIGURE 2**



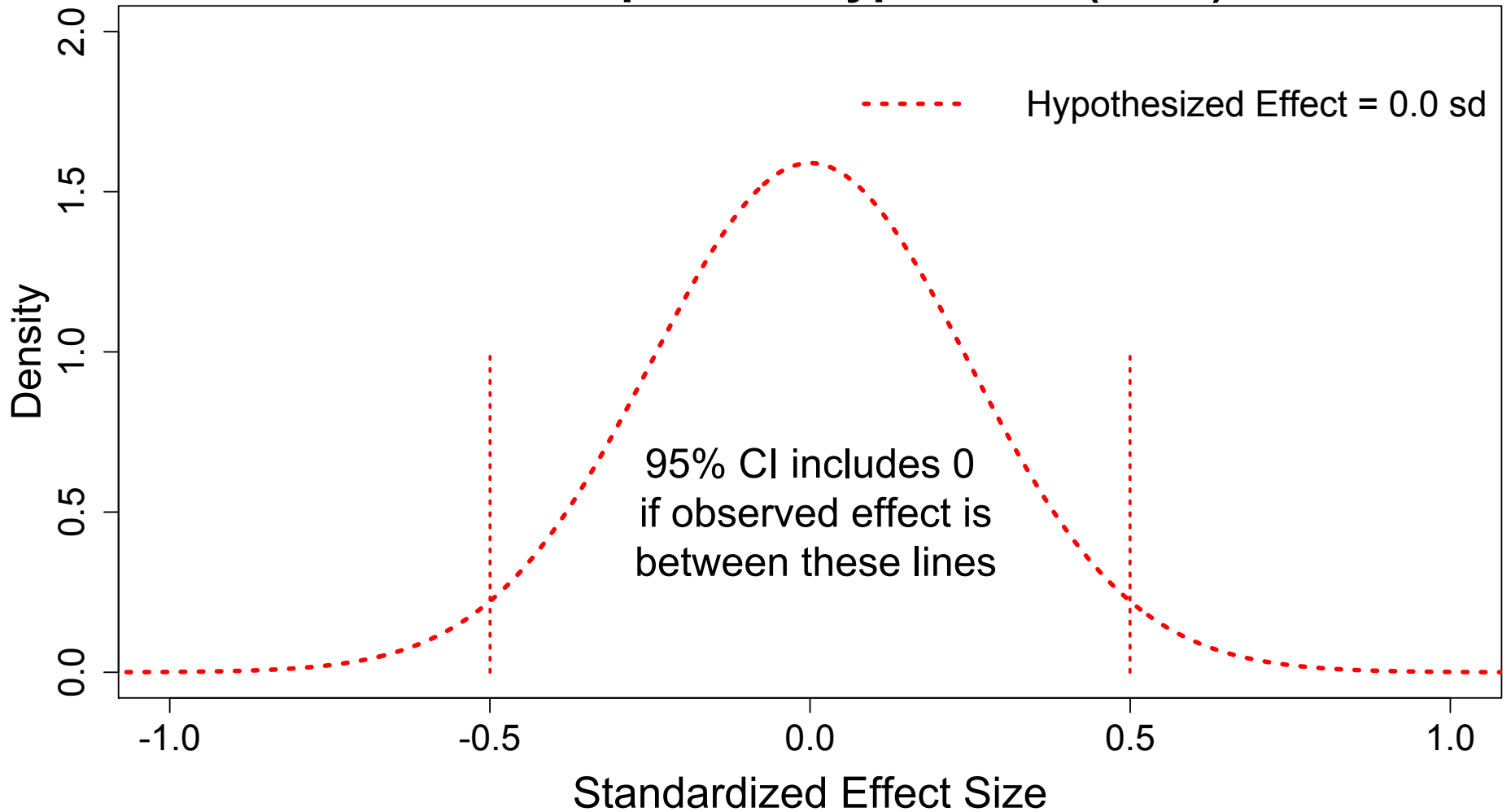
# Confidence Intervals are harder for lay readers to misinterpret disastrously!

Source for Figures 1-2: du Prel J-B, Hommel G, et al. (2009) "Confidence Interval of P-value?" *Dtsch Arztebl Int. May; 106(19): 335-339.*

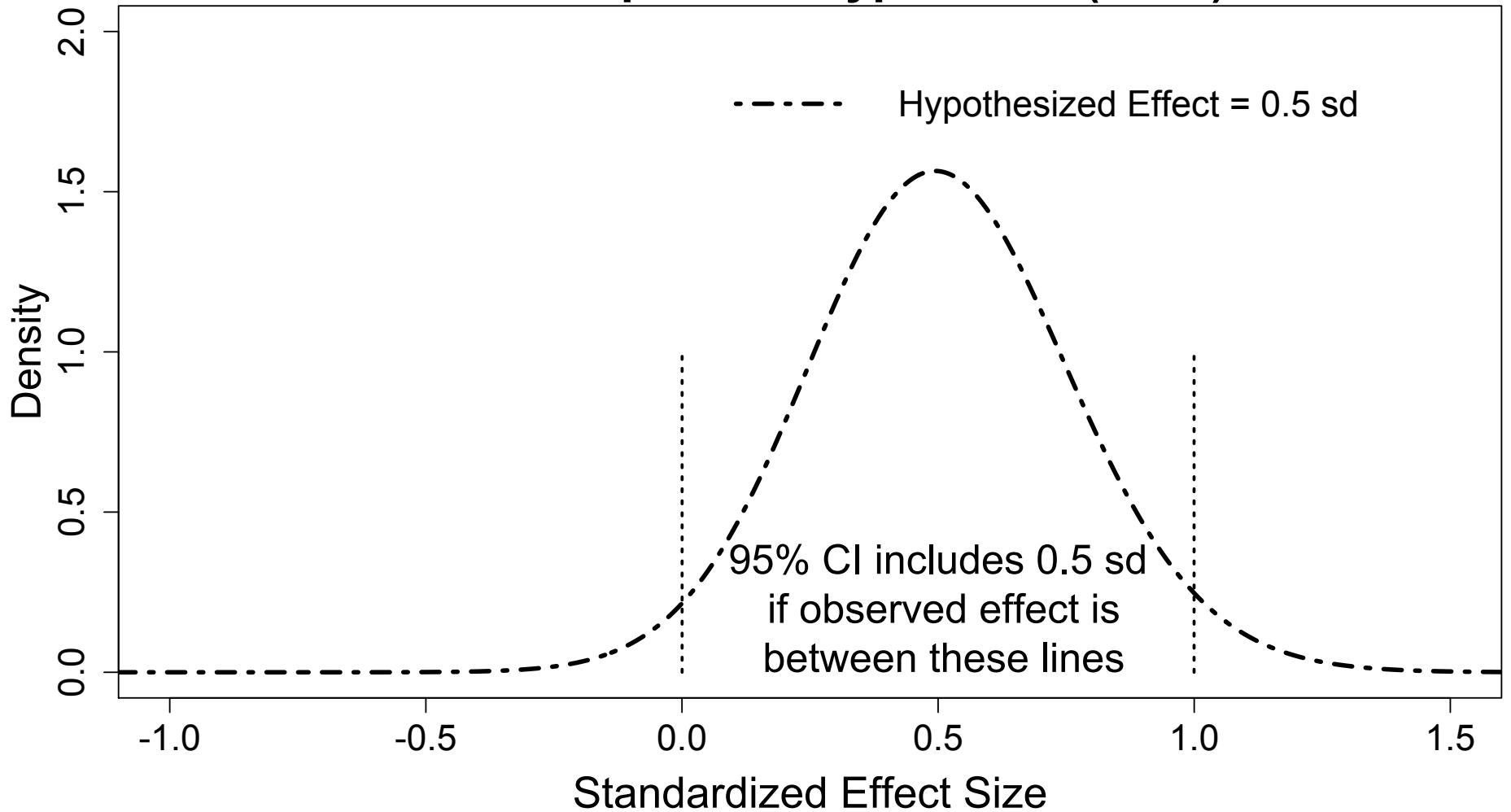
Also see: Hoenig JM and Heisey (2001). The abuse of Power: The Pervasive Fallacy of Power Calculations for Data Analysis. *Am Stat 55: 1-6.*

# **1B. DUALITY BETWEEN CONFIDENCE INTERVALS AND SIGNIFICANCE TESTING**

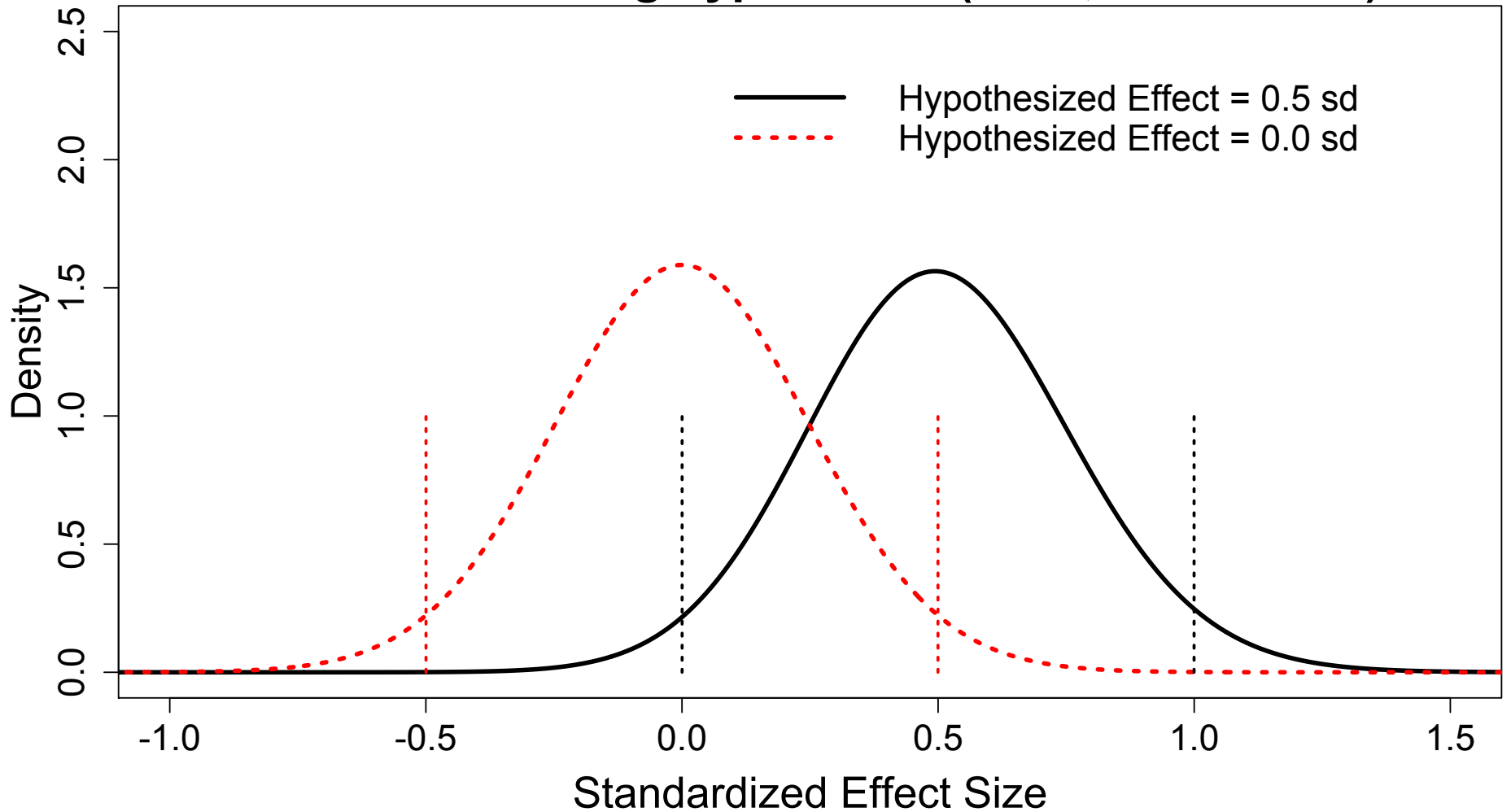
# Sampling density of estimated effect sizes under a specified hypothesis (n=64)



# Sampling density of estimated effect sizes under a specified hypothesis (n=64)

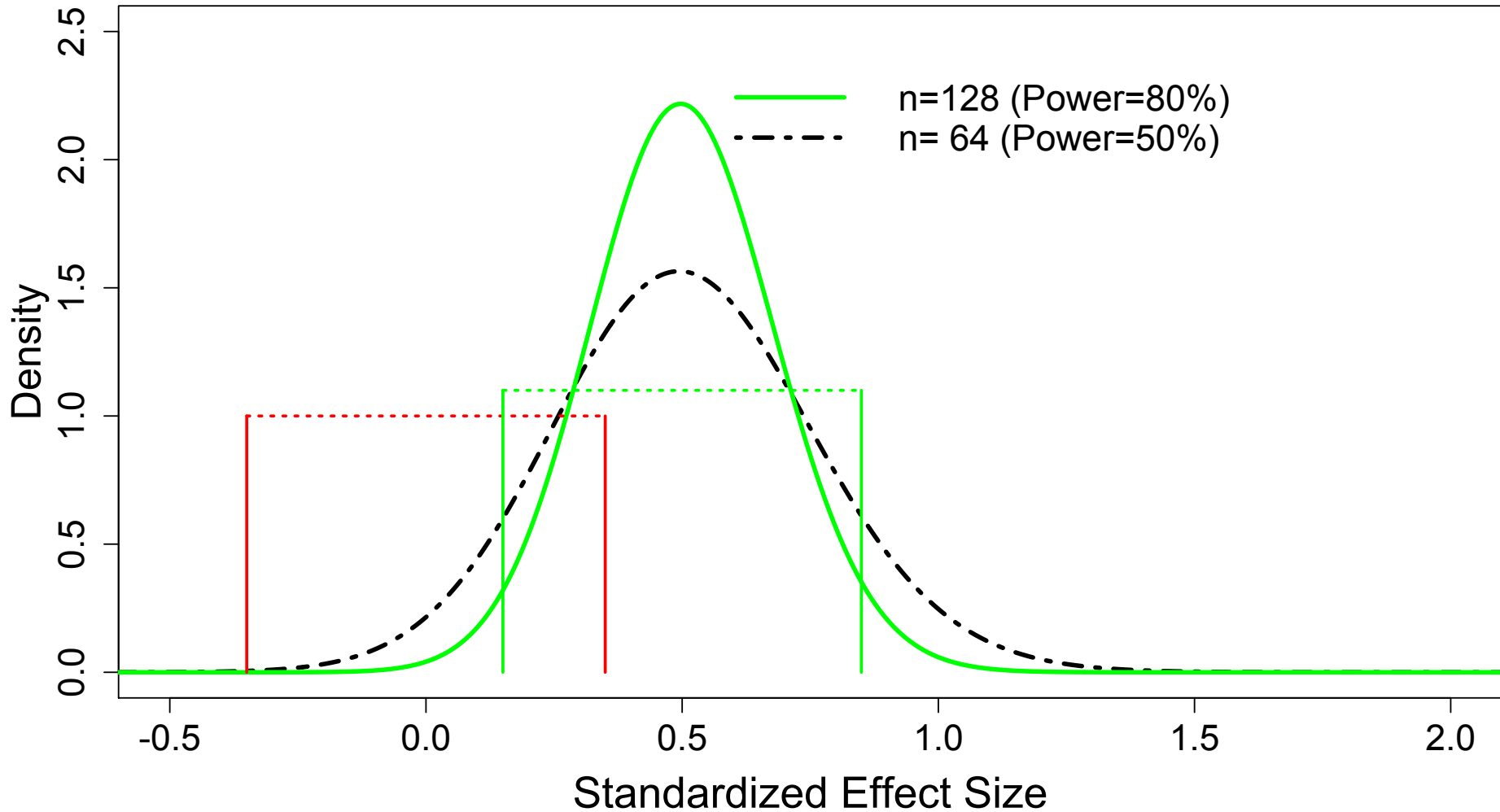


# Sampling density of estimated effect sizes under contrasting hypotheses (n=64, 50% Power)



# Sampling density of estimated effect sizes by sample size

## True Effect Size=0.5

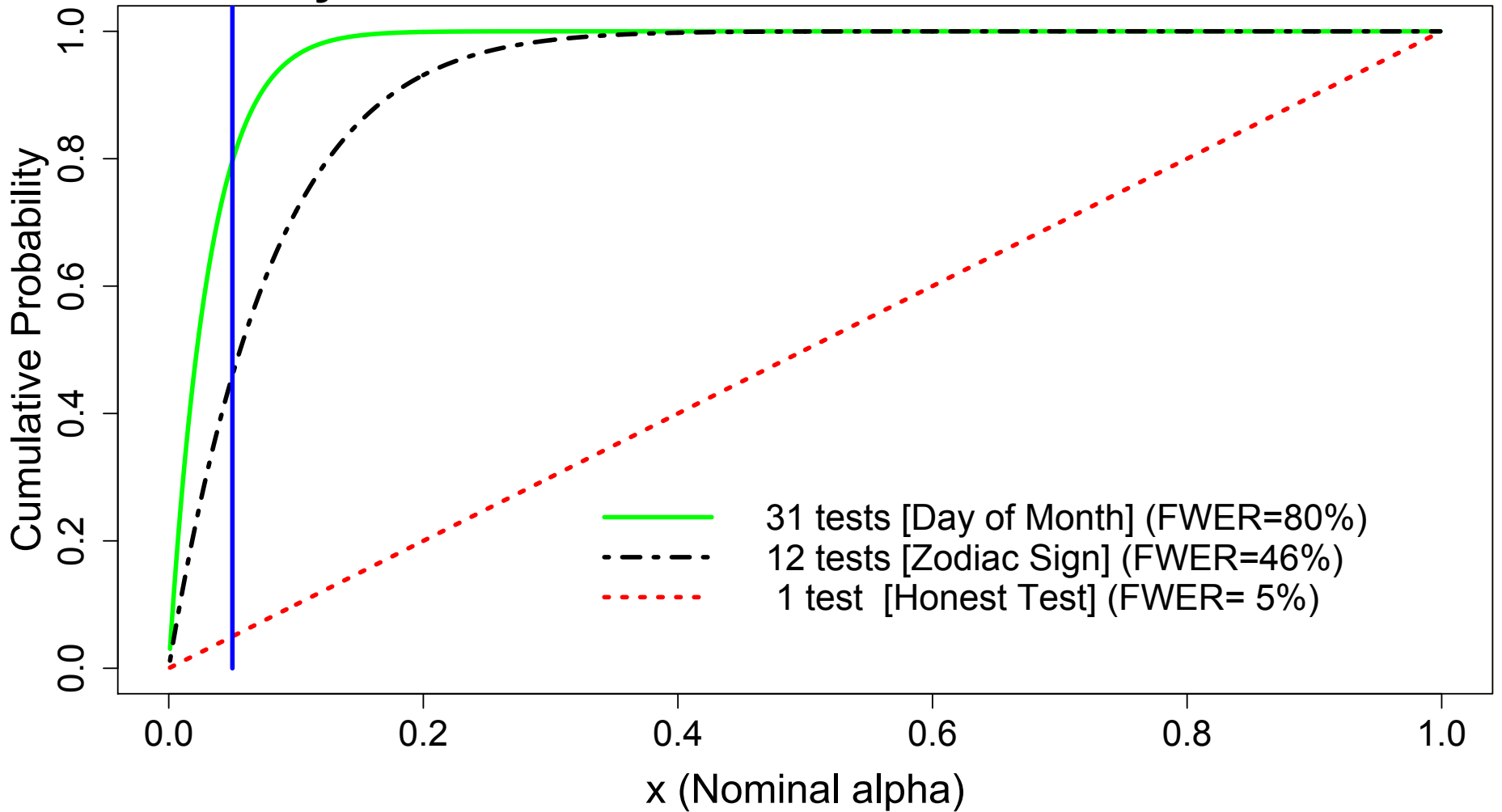


# **1C. MULTIPLE TESTING**



# CDF of minimum p-value

## By number of tests when all concern true nulls



# Bonferonni & Holm-Bonferonni Adjustments

Suppose 3 tests were performed and order p-values from smallest to largest, P<sub>1</sub>, P<sub>2</sub>, P<sub>3</sub>

	Bonferroni		Holm-Bonferroni	
Ordered	Adj. alpha	Adj. p	Adj. alpha	Adj. p
P <sub>1</sub>	0.05 / 3	3*P <sub>1</sub>	0.05 / 3	3*P <sub>1</sub>
P <sub>2</sub>	0.05 / 3	3*P <sub>2</sub>	0.05 / 2	2*P <sub>2</sub>
P <sub>3</sub>	0.05 / 3	3*P <sub>3</sub>	0.05 / 1	1*P <sub>3</sub>

## **1D. EFFECT SIZES**

# Important effect sizes

(Standardized) differences in means:

Differences in proportion (aka Risk Reduction)

Number Needed to Treat (aka NNT)

$$= 1 / \text{Risk Reduction}$$

Risk ratios (or variants involving Odds, Hazards)

[Difference in group mean log-transformed values is  
a log geometric mean ratio]

[Regression-based estimates of above]

## **2. QUALITY OF EVIDENCE**

**What do/should we mean  
when we talk about the  
quality of a study or a group  
of studies in how it addresses  
a research question?**

# Bias

A single study can provide an estimate of the true effect of an intervention (on average, in the sampled population):

*Study estimate = True Effect + Study error*

*Study error = Systematic error + sampling error*

*Bias = **Long run average**( Study error ),*  
over hypothetical repetitions of the study.

Hence, bias is essentially **systematic error** arising from such features as subject recruitment & retention, treatment assignment, measurement procedures and analysis

# Internal Validity

The extent to which the observed results of a clinical research study are not biased.

“Were the comparison groups similar in all important characteristics that may affect the measurements?”

“Were the data measured and compared using accurate methods?”

For causal claims, an internally valid study would:

Show association

Show temporal precedence

**Rule out plausible alternative explanations**

Source for definition of *Internal Validity*:

<http://www.effectivehealthcare.ahrq.gov/index.cfm/glossary-of-terms/>



# Internal validity & research designs

Quality of evidence depends crucially on level of internal validity associated with study

True experiments typically the preferred (primary) study design

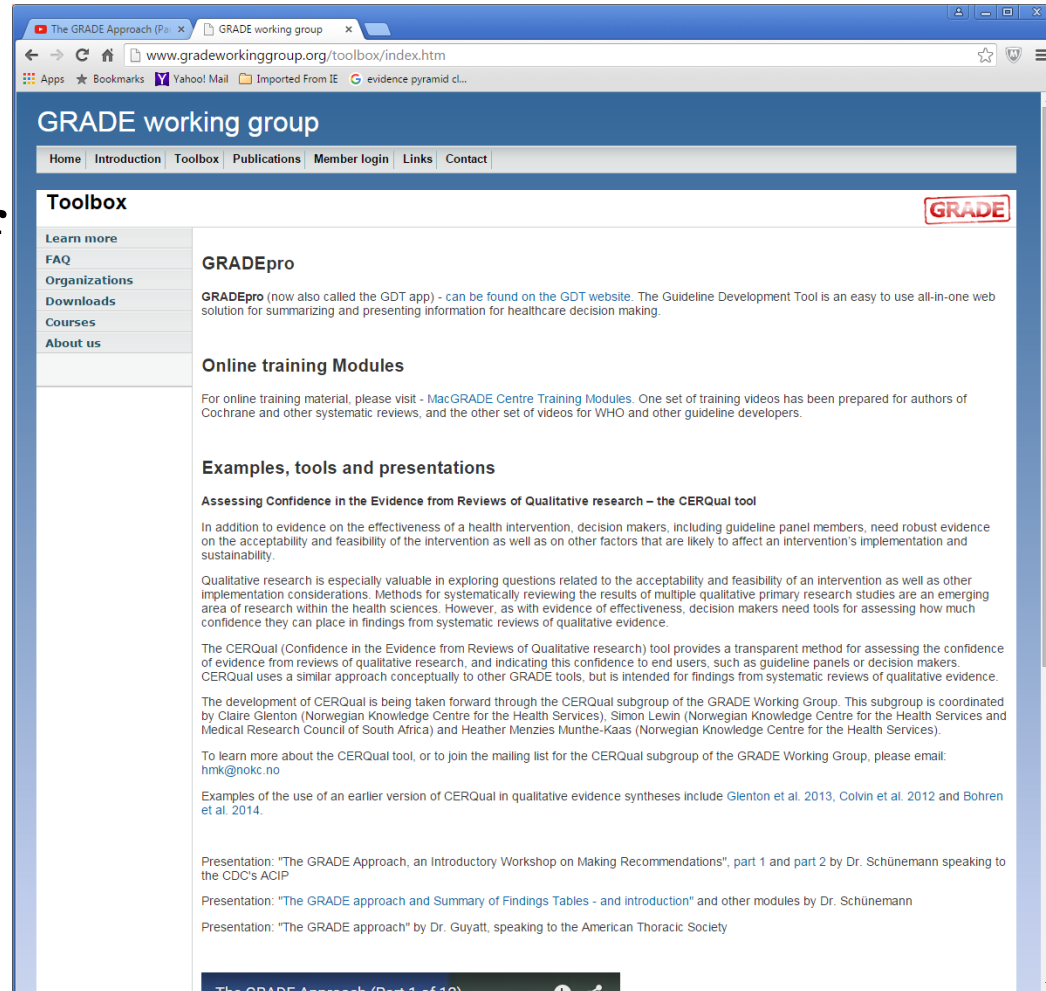
See, for example,

Puddy, R. W. & Wilkins, N. (2011). *Understanding Evidence Part 1: Best Available Research Evidence. A Guide to the Continuum of Evidence of Effectiveness*. Atlanta, GA: Centers for Disease Control and Prevention(  
[http://www.cdc.gov/violenceprevention/pdf/understanding\\_evidence-a.pdf](http://www.cdc.gov/violenceprevention/pdf/understanding_evidence-a.pdf) )

## **2. GRADE WORKING GROUP'S APPROACH TO QUALITY OF EVIDENCE**

# Key sources of material

gradeWorkingGroup.org,  
especially Dr. Guyatt's  
presentation to  
American Thoracic  
Society and the series of  
BMJ papers



The screenshot shows a web browser window displaying the GRADE working group website. The browser's address bar shows the URL [www.gradeworkinggroup.org/toolbox/index.htm](http://www.gradeworkinggroup.org/toolbox/index.htm). The website has a blue header with the text "GRADE working group" and a navigation menu with links for Home, Introduction, Toolbox, Publications, Member login, Links, and Contact. The main content area is titled "Toolbox" and features a sidebar with links for Learn more, FAQ, Organizations, Downloads, Courses, and About us. The main content includes sections for "GRADEpro", "Online training Modules", and "Examples, tools and presentations". The "Examples, tools and presentations" section is further divided into "Assessing Confidence in the Evidence from Reviews of Qualitative research – the CERQual tool", which includes text about the tool's purpose and a list of references. At the bottom, there are three presentation titles: "The GRADE Approach, an Introductory Workshop on Making Recommendations", "The GRADE approach and Summary of Findings Tables - and introduction", and "The GRADE approach".

GRADE working group

Home Introduction Toolbox Publications Member login Links Contact

**Toolbox**

Learn more  
FAQ  
Organizations  
Downloads  
Courses  
About us

**GRADEpro**

GRADEpro (now also called the GDT app) - can be found on the GDT website. The Guideline Development Tool is an easy to use all-in-one web solution for summarizing and presenting information for healthcare decision making.

**Online training Modules**

For online training material, please visit - [MacGRADE Centre Training Modules](#). One set of training videos has been prepared for authors of Cochrane and other systematic reviews, and the other set of videos for WHO and other guideline developers.

**Examples, tools and presentations**

**Assessing Confidence in the Evidence from Reviews of Qualitative research – the CERQual tool**

In addition to evidence on the effectiveness of a health intervention, decision makers, including guideline panel members, need robust evidence on the acceptability and feasibility of the intervention as well as on other factors that are likely to affect an intervention's implementation and sustainability.

Qualitative research is especially valuable in exploring questions related to the acceptability and feasibility of an intervention as well as other implementation considerations. Methods for systematically reviewing the results of multiple qualitative primary research studies are an emerging area of research within the health sciences. However, as with evidence of effectiveness, decision makers need tools for assessing how much confidence they can place in findings from systematic reviews of qualitative evidence.

The CERQual (Confidence in the Evidence from Reviews of Qualitative research) tool provides a transparent method for assessing the confidence of evidence from reviews of qualitative research, and indicating this confidence to end users, such as guideline panels or decision makers. CERQual uses a similar approach conceptually to other GRADE tools, but is intended for findings from systematic reviews of qualitative evidence.

The development of CERQual is being taken forward through the CERQual subgroup of the GRADE Working Group. This subgroup is coordinated by Claire Glenton (Norwegian Knowledge Centre for the Health Services), Simon Lewin (Norwegian Knowledge Centre for the Health Services and Medical Research Council of South Africa) and Heather Menzies Munthe-Kaas (Norwegian Knowledge Centre for the Health Services).

To learn more about the CERQual tool, or to join the mailing list for the CERQual subgroup of the GRADE Working Group, please email: [hmk@nokc.no](mailto:hmk@nokc.no)

Examples of the use of an earlier version of CERQual in qualitative evidence syntheses include [Glenton et al. 2013](#), [Colvin et al. 2012](#) and [Bohren et al. 2014](#).

Presentation: "The GRADE Approach, an Introductory Workshop on Making Recommendations", [part 1](#) and [part 2](#) by Dr. Schünemann speaking to the CDC's ACIP

Presentation: "The GRADE approach and Summary of Findings Tables - and introduction" and other modules by Dr. Schünemann

Presentation: "The GRADE approach" by Dr. Guyatt, speaking to the American Thoracic Society

# Frame research question

Explicit specification of *PICO(TS)*

Population, including settings/locations

Intervention(s), including vehicles/matrices

Comparison

Outcome(s), including Timing & how measured

Study types (designs & methodological quality)

Source: Counsell, Carl. Formulating questions and locating primary studies for inclusion in systematic reviews. *Ann Intern Med.* 1997 Sep 1;127(5):380-7.

# **PICOTS specify key elements for reviewing efficacy claims**

**P**opulation: Condition(s), comorbidities, patient demographics, diet, physical activity levels, etc.

**I**ntervention: Dosage, frequency, and method of administration.

**C**omparator: Placebo, usual diet, or active control.

**O**utcome: Health outcomes: morbidity, mortality, quality of life. **T**iming: Duration of follow-up.

**S**etting: Lab, home; co-interventions.

# Explicitly address each outcome's importance

Score each outcome:

**7-9) Critical for decision making**

4-6) Important but not critical

1-3) Limited importance

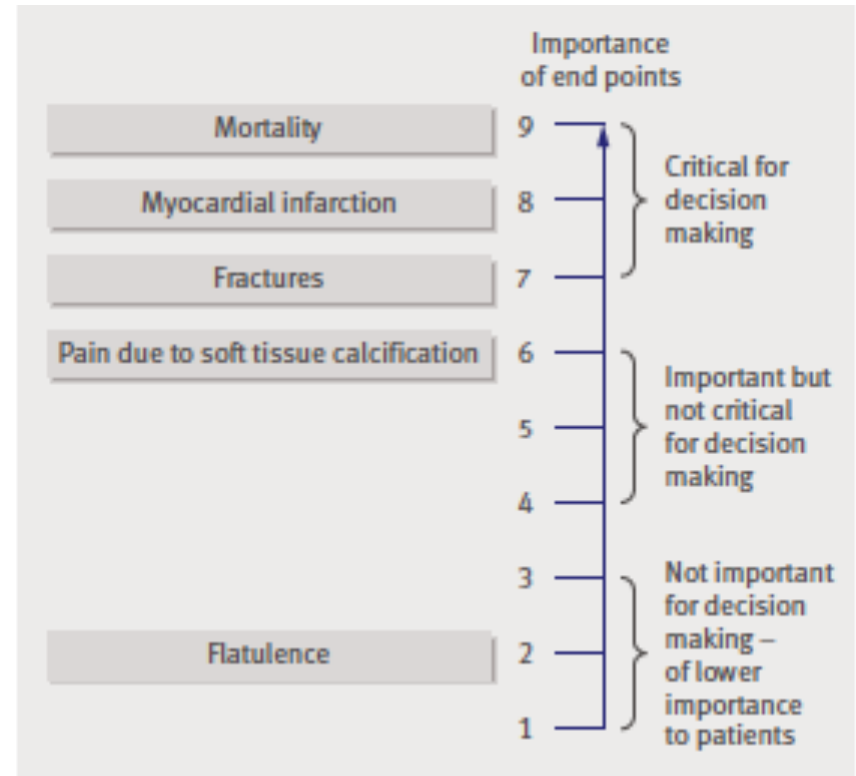


Fig 1 Hierarchy of outcomes according to importance to patients to assess effect of phosphate lowering drugs in patients with renal failure and hyperphosphataemia

# **Quality of evidence and context for recommendation**

Quality of evidence in a study is confidence that estimated effect size is close to true parameter

For decision making, quality is extent to which confidence in estimated effect is adequate to support decision

# GRADE Quality of evidence definitions

*High:* Further research (FR) unlikely to change confidence in estimated effect size (EES)

*Moderate:* FR can impact confidence in and may change EES

*Low:* FR very likely to impact confidence and likely to change EES

*Very low:* Any estimate of EES is uncertain

➤ Grading done for each important outcome!



# Classification of QoE

GRADE classifies QoE according to

- Study limitations
- Inconsistency of results
- Indirectness of evidence
- Imprecision
- Reporting bias

# Study design & limitations

RCTs presumed best, observational studies lower

For RCTs, assess (to lower quality rating)

- Random sequence generation/concealment
- Blinding
- Incomplete outcome data
- Selective reporting and other biases

For Observational studies, assess (to increase quality rating)

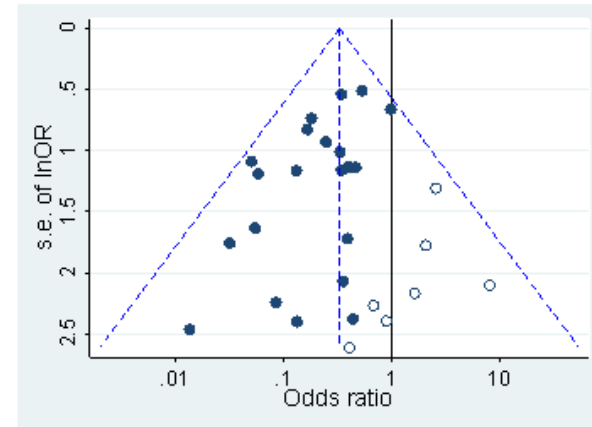
- Large magnitude of effect
- Size & direction of plausible confounding
- Dose-response gradient

# Probiotics for the prevention of Clostridium difficile-associated diarrhea in adults and children

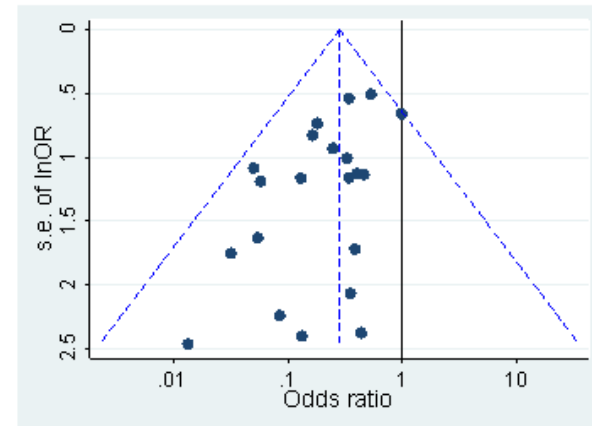
Study	Random sequence generation (selection bias)	Allocation concealment (selection bias)	Blinding of participants and personnel (performance bias): CDAD	Blinding of participants and personnel (performance bias): AE	Blinding of participants and personnel (performance bias): C. difficile incidence	Blinding of participants and personnel (performance bias): AAD	Blinding of outcome assessment (detection bias): CDAD	Blinding of outcome assessment (detection bias): AE	Blinding of outcome assessment (detection bias): C. difficile incidence	Blinding of outcome assessment (detection bias): AAD	Incomplete outcome data (attrition bias): CDAD	Incomplete outcome data (attrition bias): AE	Incomplete outcome data (attrition bias): C. difficile incidence	Incomplete outcome data (attrition bias): AAD	Selective reporting (reporting bias)	Other bias
Arvola 1999	+	?	+	+		+	+	+		+	-	-		-	+	+
Beausoleil 2007	?	?	+	+		+	+	+		+	-	+		+	+	+
Bravo 2008	?	?	+	+		+	+	+		+	-	+		+	+	+
Can 2006	?	?	+			+	+			+	+			+	?	+
Cindoruk 2007	+	?	+	+		+	+	+		+	-	?		?	+	?
Duman 2005	?	?	?	-		?	?	-		?	-	+		+	-	?
Gao 2010	+	+	+	+		+	+	+		+	+	+		+	+	?
Hickson 2007	+	?	+	?		+	+	?		+	+	-		+	-	?
Imase 2008	?	?			?	?			?	?			+	+	+	-
Klarin 2008	?	?		+	+		+	+	+		+	+	+		+	-
Koning 2008	?	?		+	+	+	+	+	+		+	+	+	+	+	?
Kotowska 2005	+	+	+	+		+	+	+		+	+	+		+	+	?
Lewis 1998	?	?			+	+			+	+			+	+	+	+
Lommermark 2010	+	+	+	+	+	+	+	+	+	+	?	?	?	?	+	-
McFarland 1995	?	?	+	+	+	+	+	+	+	+	-	+	+	-	+	-
Miller 2008a	+	?	+	+			+	+			+	+			?	-
Miller 2008b	+	?	+	+		+	+	+		+	?	?			?	-
Nord 1997	?	?		+	+			+	+			+			?	?
Plummer 2004	?	?	+		+	+	+		+		+		+		+	-
Pozzoni 2012	+	+	+	+		+	+	+		+	-	+		?	+	+
Psaradellis 2010	?	?	+	+		+	+	+	+	+	-	+		+	+	?
Rafiq 2007	?	?	?				+				?				?	?
Ruszczyński 2008	+	+	+	+			+	+			+	+			+	+

# Funnel Plots Example

Symmetrical plot in the absence of reporting bias



Asymmetrical plot in the presence of reporting bias



# **TRIAL Registration**

Registration of a clinical trial in a recognized trial registry is a crucial protection against reporting biases!

[http://www.ICMJE.org/recommendations/  
browse/publishing-and-editorial-issues/  
clinical-trial-registration.html](http://www.ICMJE.org/recommendations/browse/publishing-and-editorial-issues/clinical-trial-registration.html)

# Meta-Analysis

Meta-analysis may reduce imprecision, but it can't reduce biases

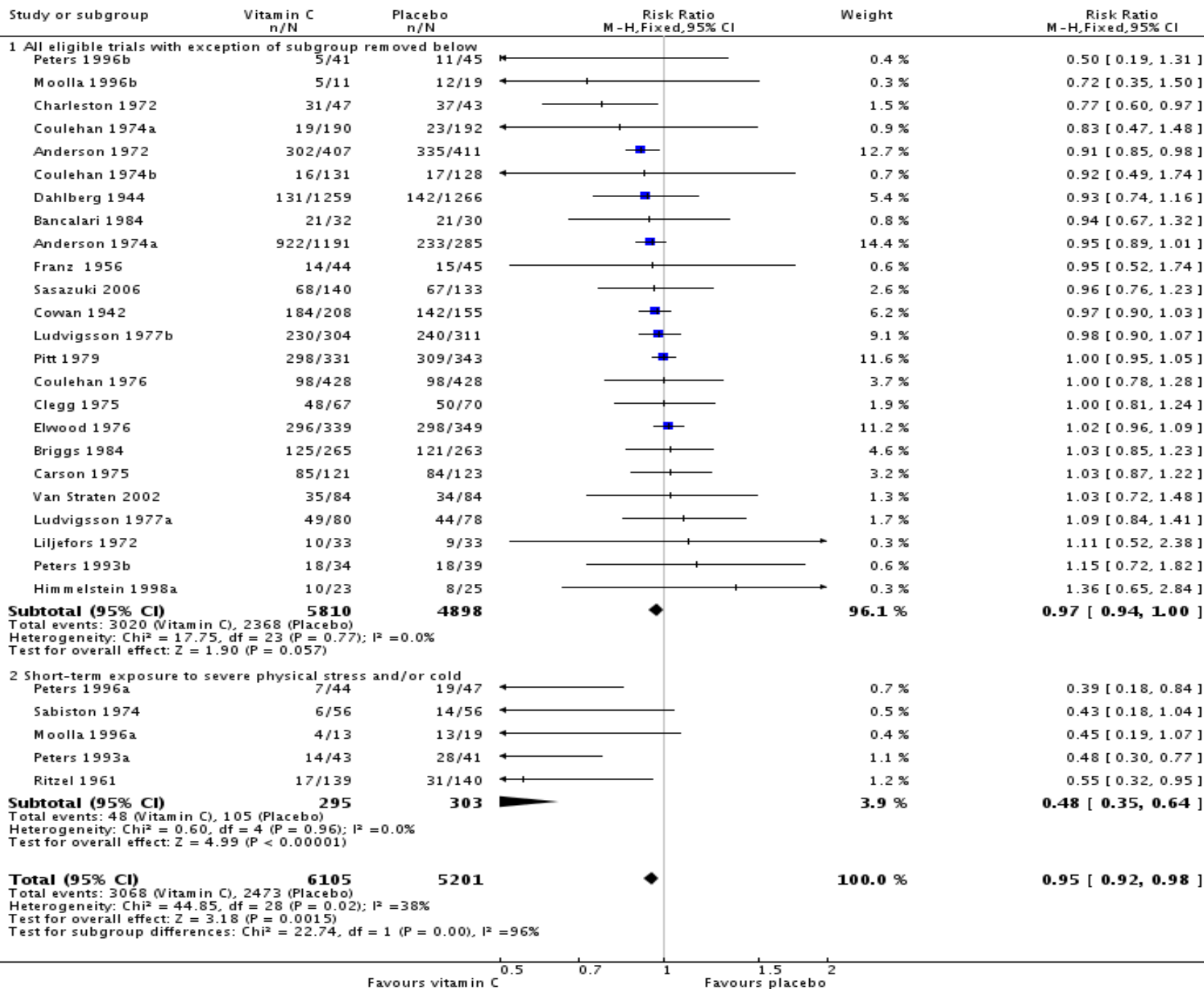
Heuristically, a pooled effect size is estimated as a weighted averages of sample effect sizes, resulting in a more precise estimate (with a smaller uncertainty interval)

Assess explainable/unexplainable heterogeneity in effect sizes, including subgroups

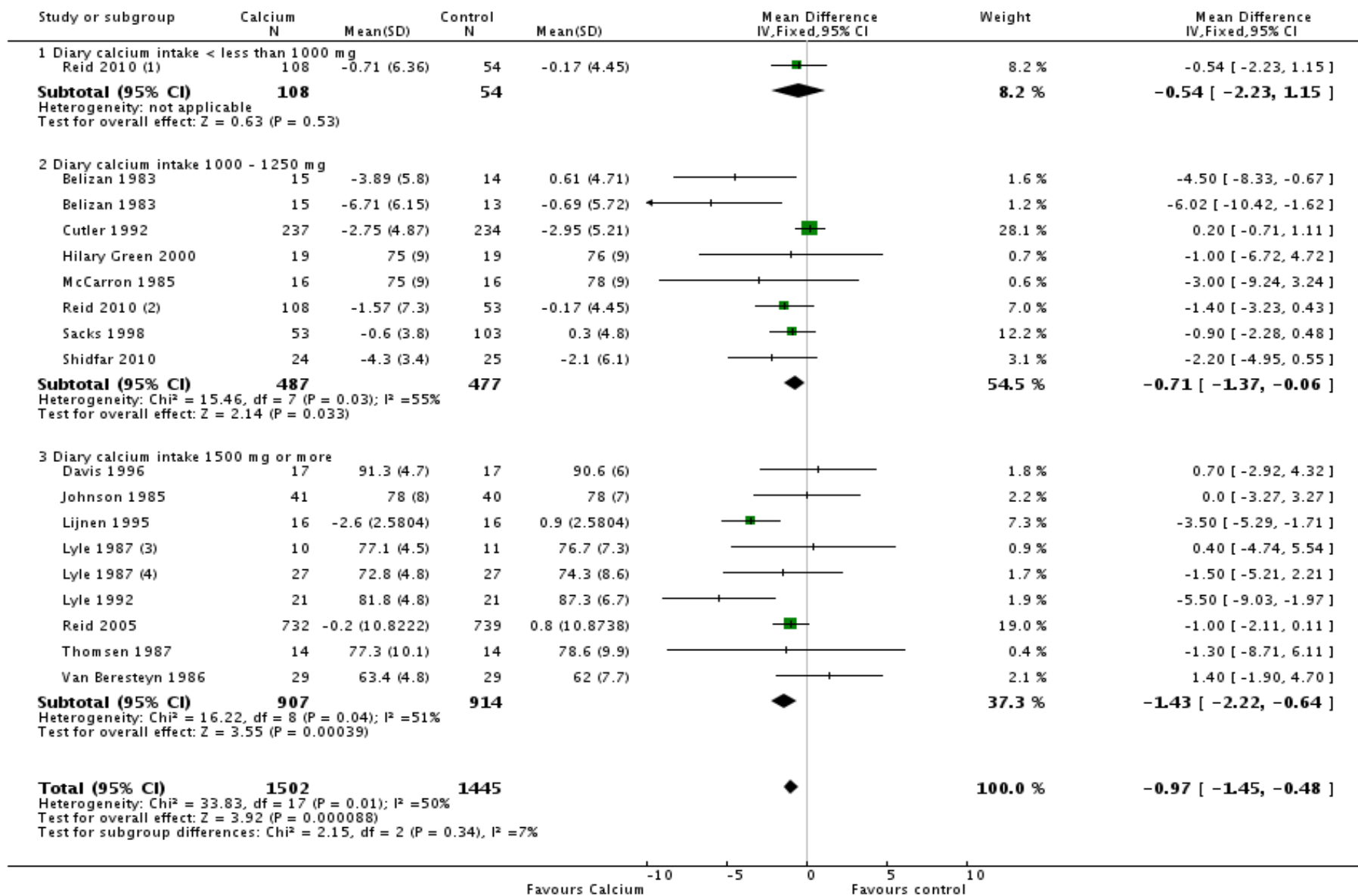
Sensitivity analyses to assess robustness

Sackett DL, Glasziou P, Chalmers I. *Meta-analysis may reduce imprecision, but it can't reduce bias*. Unpublished commentary commissioned by the New England Journal of Medicine, 1997. (see SR in Health Care, p. xiv)

Review: Vitamin C for preventing and treating the common cold  
 Comparison: 1 Incidence of colds while taking  $\geq 0.2$  g/day vitamin C regularly  
 Outcome: 1 Proportion of participants developing  $\geq 1$  cold episodes during the trial



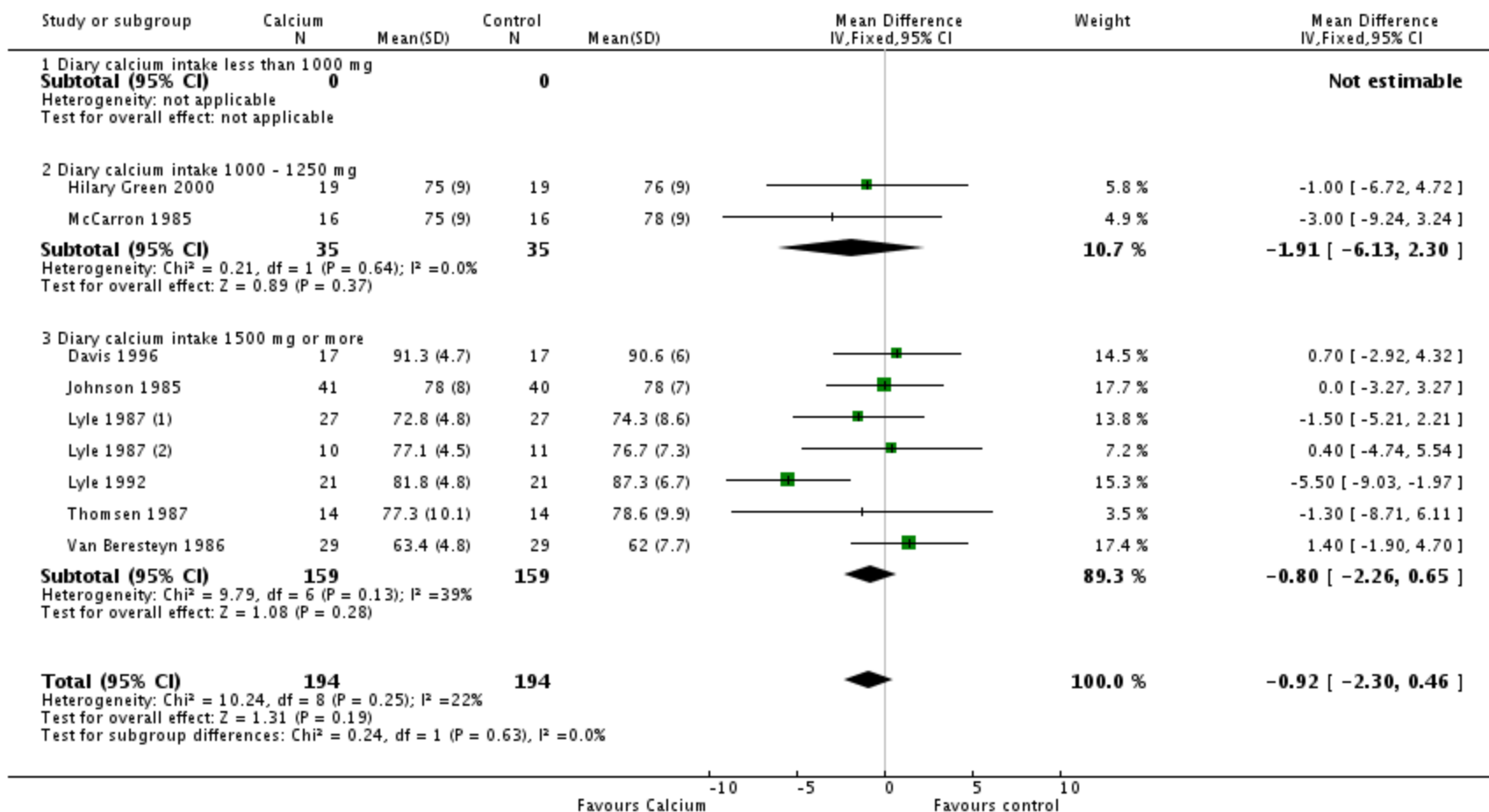
Review: Calcium supplementation for prevention of primary hypertension  
 Comparison: 1 Calcium supplementation/fortification vs control  
 Outcome: 20 Effect mean difference of diastolic blood pressure by dose



- (1) Intervention: elemental calcium 600 mg daily
- (2) Intervention: elemental calcium 1200 mg daily
- (3) Black men
- (4) White men



Review: Calcium supplementation for prevention of primary hypertension  
 Comparison: 1 Calcium supplementation/fortification vs control  
 Outcome: 24 Final value in diastolic blood pressure by dose



(1) White men  
 (2) Black men

# Factors that affect strength of a recommendation

GRADE considers

- Quality of evidence
- Uncertainty about the balance of desirable & undesirable effects
- Uncertainty or variability in values & preferences
- Uncertainty whether intervention represents a wise use of resources

Choices for recommendation: *weak* or *strong*

**Table 3. Quality assessment (GRADE evidence profile).**

No. studies	Design	Risk of bias	Inconsistency	Indirectness	Imprecision	Publication bias	Intervention group (n)	Control group (n)	Effect	Quality	Importance
<b>Emergency department visits for asthma</b>											
1	RCT	No serious limitations	n/a	No serious indirectness	No serious imprecision	n/a	50	50	P = 0.015 <sup>g</sup>	Moderate	Critical
<b>Asthma exacerbations (outcomes are number of participants experiencing asthma attacks and wheezing requiring beta2-agonists, use of beta2-agonists (puffs/day), and undefined asthma exacerbation attack)</b>											
6	RCT	Serious limitations <sup>b</sup>	No serious inconsistency	Serious indirectness <sup>c</sup>	No serious imprecision	None	257	250	0.41 (0.27 to 0.63) <sup>f</sup> ; P<0.05 <sup>g</sup> ; no effect <sup>d</sup>	Low	Critical
<b>Asthma symptoms (outcomes are scores (in points) based on ACT, ATAQ for children, daily diary card, ACQ and undefined asthma symptom scores)</b>											
6	RCT	Serious limitations <sup>a</sup>	No serious inconsistency	Serious indirectness <sup>c</sup>	No serious imprecision	None	117	114	No effect <sup>d</sup> ; P = 0.01(6mo follow-up) <sup>g</sup>	Low	Critical
<b>Lung function (outcomes are FEV1 (L in 1 sec or % of predicted value) and PEF (mL/min)</b>											
7	RCT	Serious limitations <sup>a</sup>	No serious inconsistency	No serious indirectness	No serious imprecision	None	167	164	0.00 (-3.17 to 3.18) <sup>g</sup> ; P<0.001 <sup>g</sup> ; no effect <sup>d</sup>	Low	Critical
<b>Serum 25(OH)D (nmol/L)</b>											
6	RCT	Serious limitations <sup>a</sup>	Serious inconsistency <sup>h</sup>	No serious indirectness	No serious imprecision	None	117	114	19.66 (5.96 to 33.37) <sup>f</sup> ; no effect <sup>d</sup>	Low	Important

Abbreviations: ACT, Asthma Control Test; ATAQ, Asthma Therapy Assessment Questionnaire; ACQ, Asthma Control Questionnaire; FEV1, forced expiratory volume in 1 second; PEF, peak expiratory flow rate.

<sup>a</sup>Unclear allocation concealment, blinding of participants and outcome assessors, accounting of patients and outcome events, and other risk of bias (carryover effects in crossover trial).

<sup>b</sup>Unclear allocation concealment, blinding of participants and outcome assessors, accounting of patients and outcome events.

<sup>c</sup>Differences in interventions and outcomes measured across studies.

<sup>d</sup>Non-significant effect across studies not included in the meta-analysis.

<sup>e</sup>Weighted difference in mean (WMD) change between intervention and control group.

<sup>f</sup>Risk ratio (RR): risk of experiencing asthma exacerbation in the intervention group as compared to the control group.

<sup>g</sup>Not included in the meta-analysis; favours intervention group.

<sup>h</sup>Significant statistical heterogeneity observed based on random effects meta-analysis.

<sup>i</sup>Weighted mean difference (WMD) at end of intervention between intervention and control group.

# Recommendations

- ✓ Frame research questions meaningfully
- ✓ Consider whether and how study will contribute to evidence synthesis
- ✓ Emphasize key determinants of study quality (adequate sample size, randomization, allocation concealment, objective measurement, complete follow-up and honest reporting) and of quality of evidence synthesis (study limitations, inconsistency, indirectness, imprecision, reporting biases)
- ✓ Report effect sizes and 95% CI
- ✓ Have analysis and interpretation strategies to account for multiple outcomes



# Questions and Comments

[djtancredi@ucdavis.edu](mailto:djtancredi@ucdavis.edu)