



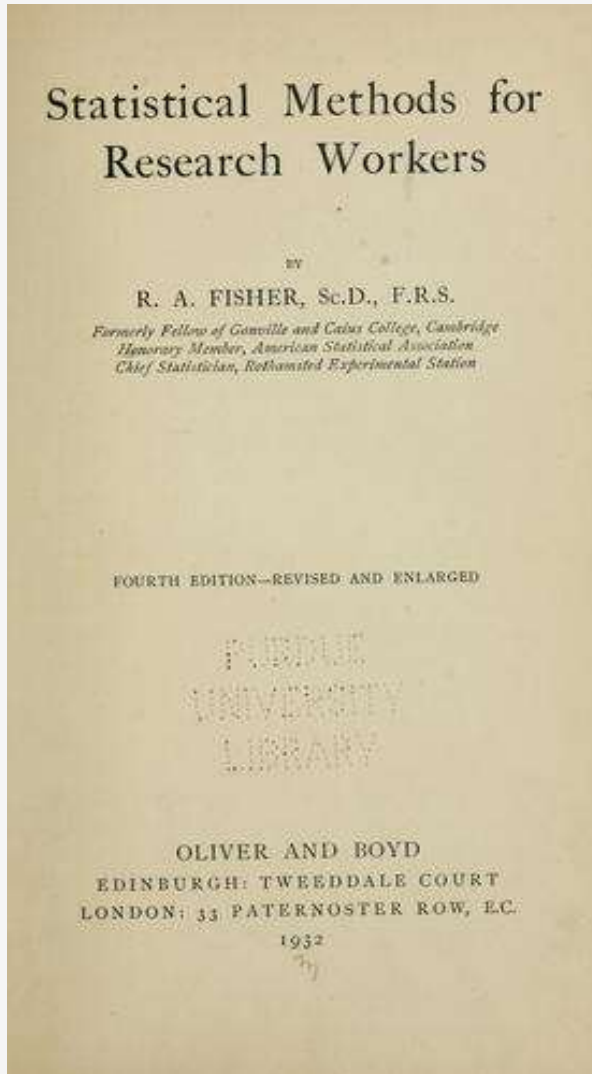
# Do Cochrane and non-Cochrane editors & authors prefer reporting statements based on statistically significant differences or do they prefer non-binary options?

**Ciapponi A, Glujovsky D, Bardach A.**  [aciapponi@iecs.org.ar](mailto:aciapponi@iecs.org.ar)

**No conflict of interest**

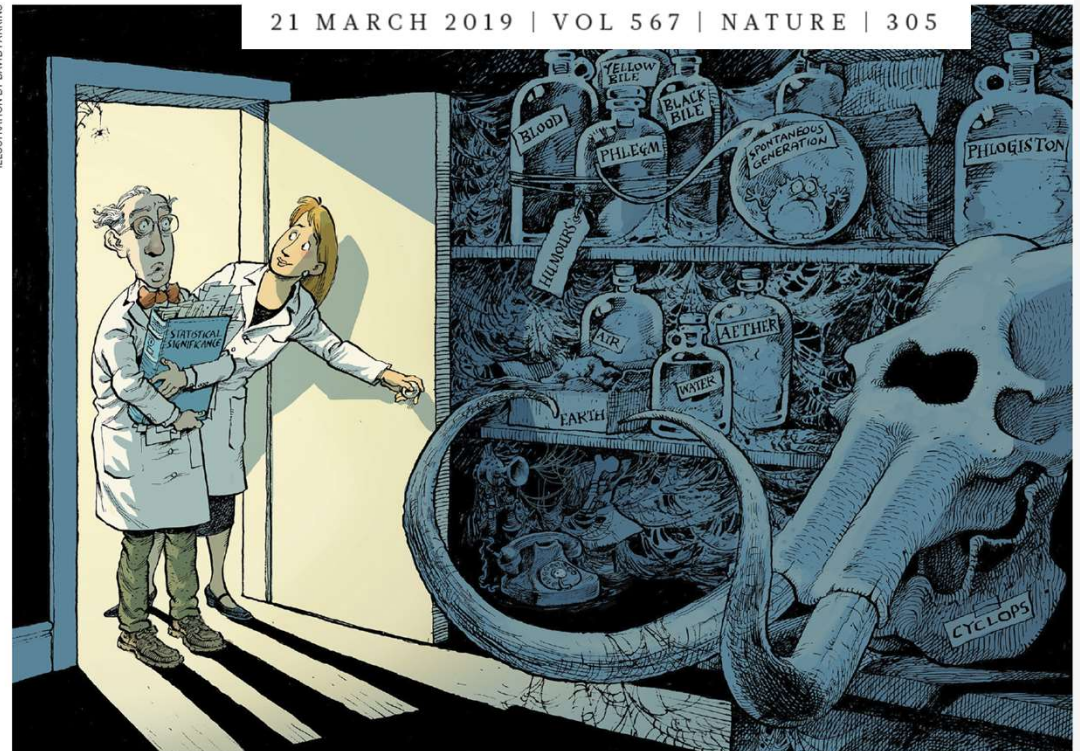
Trusted evidence.  
Informed decisions.  
**Better health.**





P

ILLUSTRATION BY DAVID PARKINS



21 MARCH 2019 | VOL 567 | NATURE | 305

# Retire statistical significance

Valentin Amrhein, Sander Greenland, Blake McShane and more than 800 signatories call for an end to hyped claims and the dismissal of possibly crucial effects.

EDITORIAL

Open Access

There is life beyond the statistical significance

Ciapponi *et al. Reprod Health* (2021) 18:80  
<https://doi.org/10.1186/s12978-021-01131-w>



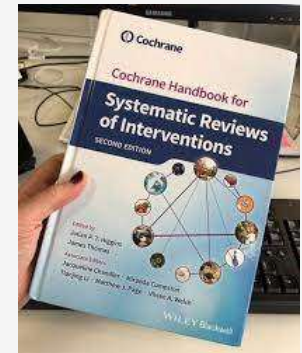
 BMC

- For decades the p value-based interpretation and reporting of results dominated the publications, but scientific community agrees that this binary approach is not enough and suggested a systemic reform to change this paradigm.
- The Cochrane Handbook, recommends reporting the point estimate, the CI + exact P-value, MIDs, some narrative statements, and against binary approaches:

Review authors should not describe results as ‘statistically significant’, ‘not statistically significant’ or ‘non-significant’ or unduly rely on thresholds for P values, but report the confidence interval together with the exact P value. Chapter 15



- **Which is the approach of Cochrane and non-Cochrane editors and authors for interpretation and reporting this case?**



NEW SECTION

# Methods

## Stakeholders surveyed

- **Cochrane editors** (N=65)
- **Cochrane authors**, that published reviews from 1/1/23 to 7/25/23 (N=321)

Source: Archie

- **Non-Cochrane editors** (N=20)
- **Non-Cochrane authors** (N=322)

Source: the 20 highest impact factor in the "General Medicine" and "Internal Medicine" categories in 2021 (edition of Clarivate Analytics Journal Citation Report) with available e-mail



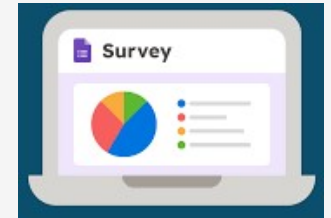
## Stakeholders had to choose the binary or non-binary option that better expresses the results for the following scenario:

“After exhaustive literature searches, a systematic review identified only two pivotal RCTs that evaluated the mortality of drug X versus placebo (P) in patients with a rare genetic disease. The risk of bias for all domains was low in both RCTs (assessed using the Cochrane RoB-2 tool), and there was no methodological, clinical, or statistical heterogeneity between studies. The meta-analysis showed the following results:”

	X drug	Placebo
<b>Mortality risk</b>	<b>26%</b> (10/39)	<b>45%</b> (18/40)
<b>Risk difference</b>	<b>With X 19% lower mortality</b> (95% CI <b>40% lower</b> to <b>1% higher</b> )	
<b>Risk Ratio</b>	<b>0.57</b> (95% CI <b>0.30</b> a <b>1.08</b> )	
<b>P value</b>	<b>0.0721</b>	

**Clinical important difference** with CI crossing the null effect

Please, select only the statement that better reflects the interpretation of the results, even if more than one is correct:













1. Mortality with **X is lower** than with **Placebo (P)**
2. Mortality with **X is probably lower** than **P**, **but no statistically significant** differences were found
3. Mortality with **X is probably lower** than with **P**, **but** the probability that the difference is **due to chance is 7%**
4. Mortality with **X is possibly lower** than with **P**, but the possibility that the difference is due to **chance is 7%** (this one is similar to statement 3, but replacing probably by possibly)
5. Mortality with **X is probably lower** than with **P**, **but** the **confidence interval (CI)** is **compatible with** both a **reduction** and an **increase** in mortality.
6. Mortality with **X is possibly lower** than with **P**, but the confidence interval is compatible with both a **reduction** and an **increase** in mortality (this one is similar to statement 5, but replacing probably by possibly)
7. **No differences were found** between **X** and **P**
8. Intervention **X did not show higher mortality** than **P**
9. **No statistically significant differences were found** between **X** and **P**

Could you, please, justify your answer?

▼ Should drug X vs. placebo be used for rare condition Y?

Bottom panel  Explanations 

Drug X compared to placebo for rare condition Y 

Certainty assessment						Summary of findings				Importance 	
No of studies 	Study design 	Risk of bias 	Inconsistency 	Indirectness 	Imprecision 	No of patients		Effect			Certainty 
						Drug X 	Placebo 	Relative (95% CI) 	Absolute (95% CI) 		

Mortality (follow-up: mean 12 months) 

2	randomised trials	not serious	not serious	not serious		none	10/39 (25.6%)	18/40 (45.0%)	RR 0.57 (0.30 to 1.08)	194 fewer per 1,000 (from 315 fewer to 36 more)	-	CRITICAL
---	-------------------	-------------	-------------	-------------	--	------	---------------	---------------	---------------------------	--	---	----------

Imprecision 

- not serious
- serious
- very serious
- extremely serious





## Updates on rating imprecision



ELSEVIER



Journal of Clinical Epidemiology 150 (2022) 225–242

Journal of  
Clinical  
Epidemiology

### GRADE GUIDANCE SERIES

#### GRADE guidance 35: update on rating imprecision for assessing contextualized certainty of evidence and making decisions

Holger J. Schünemann<sup>a,b,c,q,\*</sup>, Ignacio Neumann<sup>b,d</sup>, Monica Hultcrantz<sup>e</sup>, Romina Brignardello-Petersen<sup>b</sup>, Linan Zeng<sup>b,f</sup>, M Hassan Murad<sup>g</sup>, Ariel Izcovich<sup>h</sup>, Gian Paolo Morgano<sup>b</sup>, Tejan Baldeh<sup>b</sup>, Nancy Santesso<sup>a,b</sup>, Carlos Garcia Cuello<sup>a,b</sup>, Lawrence Mbuagbaw<sup>a,b</sup>, Gordon Guyatt<sup>b,c</sup>, Wojtek Wiercioch<sup>a,b</sup>, Thomas Piggott<sup>a,b</sup>, Hans De Beer<sup>i</sup>, Marco Vinceti<sup>j</sup>, Alexander G. Mathioudakis<sup>k</sup>, Martin G. Mayer<sup>l,m,n</sup>, Reem Mustafa<sup>o</sup>, Tommaso Filippini<sup>i</sup>, Alfonso Iorio<sup>b,c</sup>, Robby Nieuwlaat<sup>a,b</sup>, Maura Marcucci<sup>b,c</sup>, Pablo Alonso Coello<sup>p</sup>, Stefanos Bonovas<sup>q,r</sup>, Daniele Piovani<sup>q,r</sup>, George Tomlinson<sup>s,t</sup>, Elie A. Akl<sup>b,u</sup>, for the GRADE Working Group



ELSEVIER



Journal of Clinical Epidemiology 150 (2022) 216–224

Journal of  
Clinical  
Epidemiology

### GRADE GUIDANCE SERIES

#### GRADE Guidance 34: update on rating imprecision using a minimally contextualized approach

Linan Zeng<sup>a,b,\*</sup>, Romina Brignardello-Petersen<sup>b</sup>, Monica Hultcrantz<sup>c</sup>, Reem A. Mustafa<sup>d</sup>, Mohammad H. Murad<sup>e</sup>, Alfonso Iorio<sup>b,f</sup>, Gregory Traversy<sup>g</sup>, Elie A. Akl<sup>h</sup>, Martin Mayer<sup>i,j,k</sup>, Holger J. Schünemann<sup>b,f</sup>, Gordon H. Guyatt<sup>b,f</sup>



ELSEVIER



Journal of Clinical Epidemiology 147 (2022) 69–75

Journal of  
Clinical  
Epidemiology

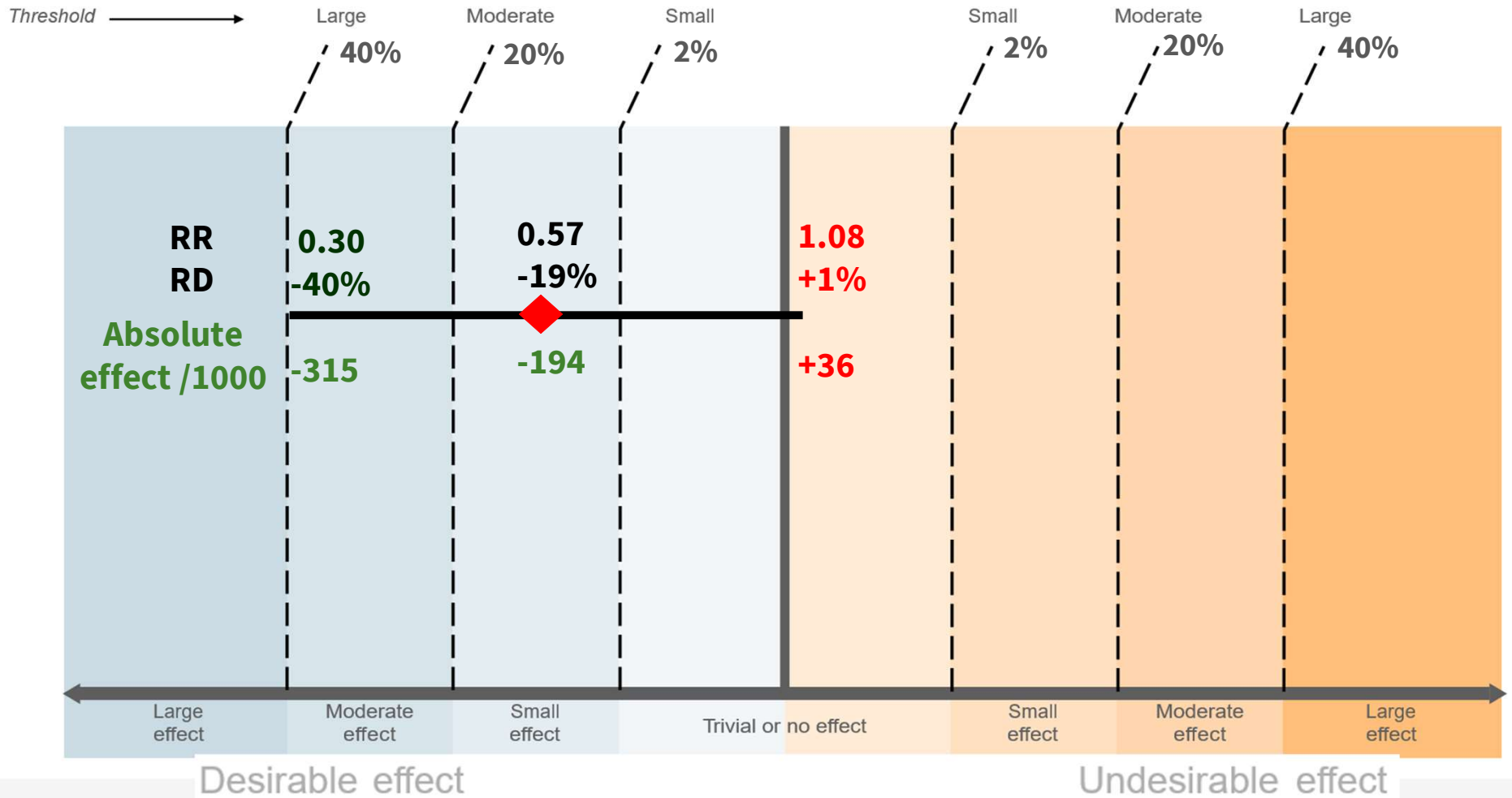
### Other GRADE Papers

#### Using Explicit Thresholds were valuable for judging Benefits and Harms in partially contextualized GRADE Guidelines

Ignacio Neumann<sup>a,b,c,\*</sup>, Eduardo Quiñelen<sup>b</sup>, Paula Nahuelhual<sup>b</sup>, Pamela Burdiles<sup>b</sup>, Natalia Celedón<sup>b</sup>, Katherine Cerda<sup>b</sup>, Paloma Herrera-Omegna<sup>b</sup>, Patricia Kraemer<sup>b</sup>, Karen Dominguez Cancino<sup>b,c,d</sup>, Juan Pablo Valenzuela<sup>b</sup>, Dino Sepúlveda<sup>b</sup>, Gian Paolo Morgano<sup>c</sup>, Elie A. Akl<sup>c,e</sup>, Holger J. Schünemann<sup>c</sup>

# GRADE

## Thresholds and ranges for trivial, small, moderate and large effects



Calculating the review informatin size (RIS) when effects are large		Prop	Prop (%)
User-entered data			
Threshold for small effect (MID)		0.02	2%
Threshold for moderate effect		0.2	20%
Threshold for large effect		0.4	40%
Baseline risk		0.45	45%
Type I error (alpha)		0.05	5%
Power		0.8	80%
Review size (total number of participants)		800	

### Results of RIS calculation

Threshold	No of participants (studies)	Sample size to rule out	GRADE a
Small (MID)			958 <b>Rate down</b>
Moderate effect	79 (2 RCTs)		958 <b>Rate down</b>
Large effect			94 <b>Do not rate down</b>

Should drug X vs. placebo be used for rare condition Y?

Bottom panel

Explanations

Drug X compared to placebo for rare condition Y

Certainty assessment							Summary of findings				Importance	
No of studies	Study design	Risk of bias	Inconsistency	Indirectness	Imprecision	Other considerations	No of patients		Effect			Certainty
							Drug X	Placebo	Relative (95% CI)	Absolute (95% CI)		

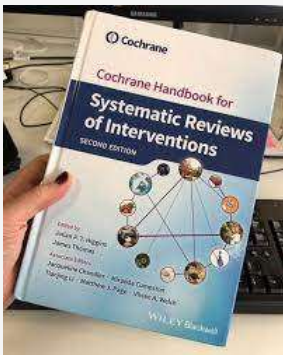
Mortality (follow-up: mean 12 months)

2	randomised trials	not serious	not serious	not serious	very serious <sup>a</sup>	none	10/39 (25.6%)	18/40 (45.0%)	RR 0.57 (0.30 to 1.08)	194 fewer per 1,000 (from 315 fewer to 36 more)	⊕⊕○○ Low	CRITICAL
---	-------------------	-------------	-------------	-------------	---------------------------	------	---------------	---------------	---------------------------	--	-------------	----------

<sup>a</sup> Downgraded to levels because the confidence interval crossed two effect thresholds

**Table 1** Suggested narrative statements for phrasing conclusions

Certainty of the evidence	Effect size	Suggested statements for conclusions (replace X with intervention, choose 'reduce' or 'increase' depending on the direction of the effect, replace 'outcome' with name of outcome, include 'when compared with Y' when needed)
⊕⊕○○ Low	Moderate	X <u>may</u> reduce/increase outcome The evidence <u>suggests</u> X reduces/increases outcome X may result in a reduction/increase in outcome The evidence suggests X results in a reduction/increase in outcome
<del>Probable</del>	<b>Possible</b>	



## Statement assessment (not exclude the reporting of numbers)

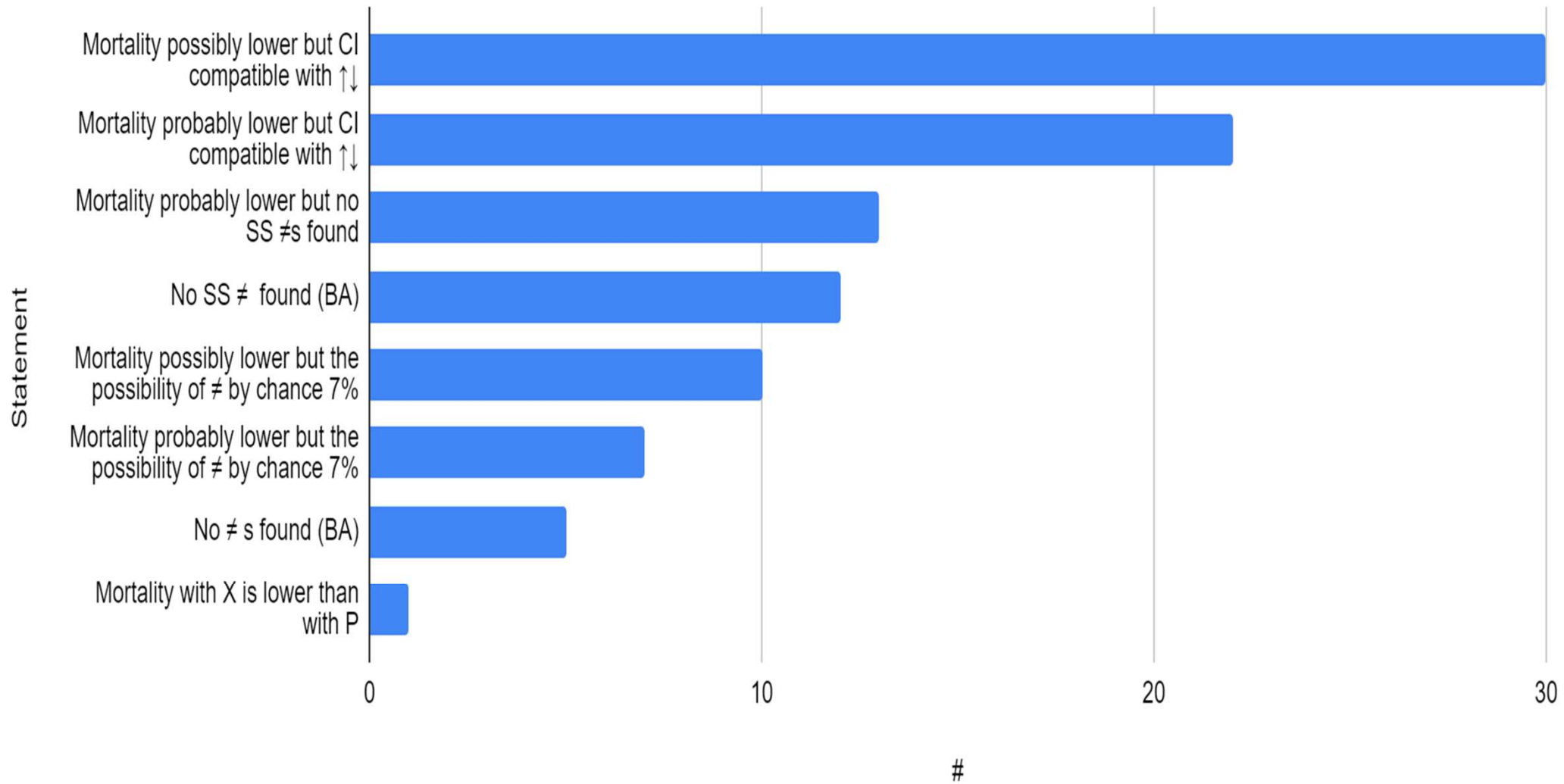
1. Mortality with X is **lower** than with P (**Binary approach**)
2. Mortality with X is **probably lower** than P, **but no statistically significant** differences were found
3. Mortality with X is **probably lower** than with P, **but** the probability that the difference is **due to chance is 7%**
4. Mortality with X is **possibly lower** than with P, but the possibility that the difference is due to **chance is 7%**
5. Mortality with X is **probably lower** than with P, **but** the **confidence interval (CI)** is **compatible with** both a **reduction** and an **increase** in mortality.
6. Mortality with X is **possibly lower** than with P, **but** the CI is compatible with both a **reduction** and an **increase** in mortality
7. **No differences were found** between X and P (**Binary approach** not high CoE: serious imprecision)
8. Intervention X **did not show higher mortality** than P (**Binary approach** not high CoE: serious imprecision)
9. **No statistically significant differences were found** between X and P (**Binary approach** not high CoE: serious imprecision)

**NEW SECTION**

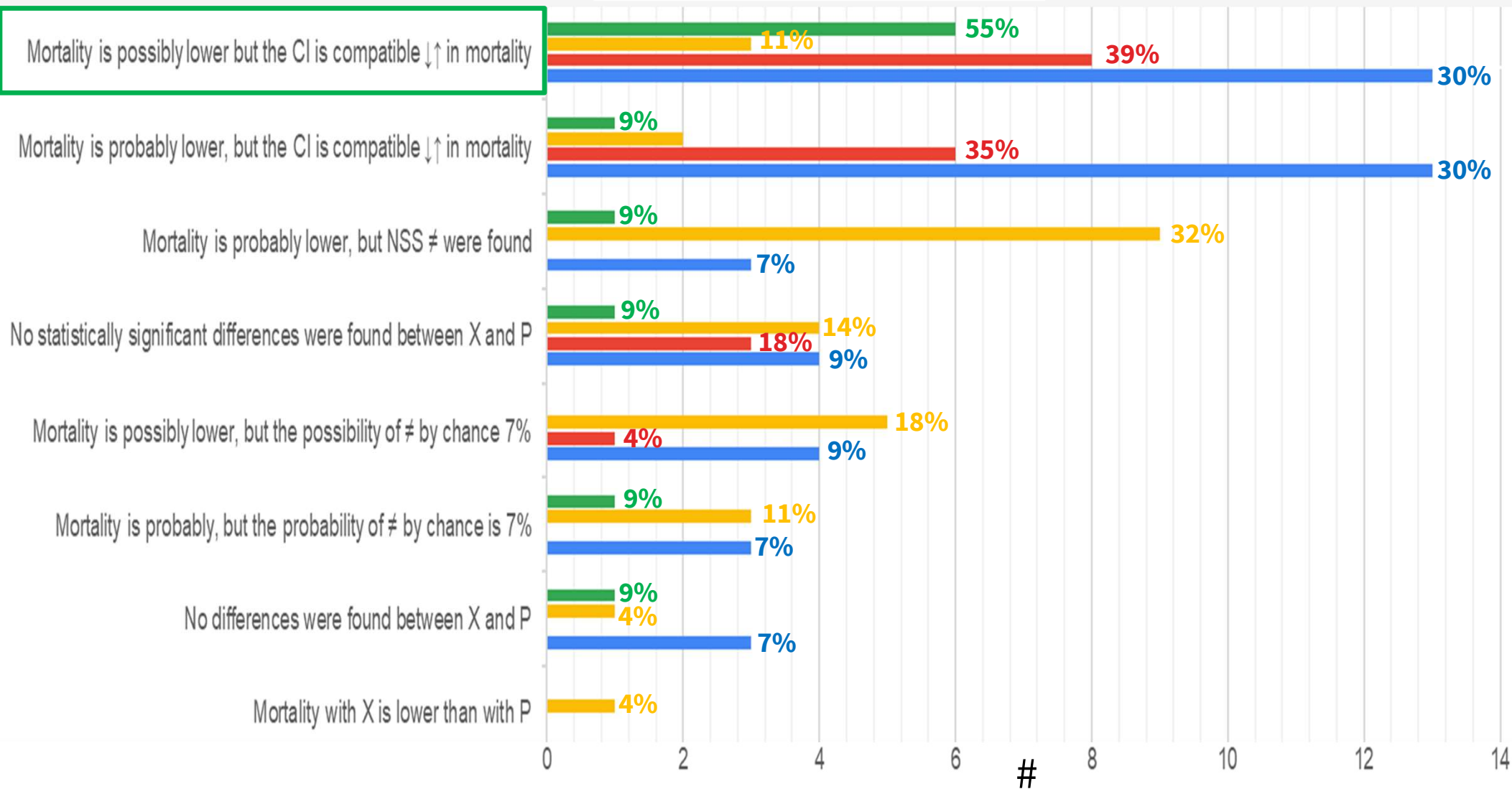
# Results



# Statement (N=101)



**Statement by stakeholder, % by:** ■ NCE ■ NCA ■ CE ■ CA





# Cochrane Authors (n=43, 13%)

the CI is wide, so benefit on mortality cannot be completely excluded

Low certainty, NSS, Wide CI

The # of participants and events so low that I'd rather say possibly than probably

Absolute effects are as important as relative effects

CI compatible with ↓ & ↑ in mortality

Add p value

Moderate certainty

CI including null effect => No ≠

**1. Mortality possibly lower but CI compatible with ↑↓**

30.2%

**1. Mortality probably lower but CI compatible with ↑↓**

30.2%

**2. Mortality possibly lower but the possibility of ≠ by chance 7%**

9.3%

↓ 2 levels GRADE because the OIS is not reached (<100 events, and CI overlapping the line of no effect). The best "mortality with X may be lower than with placebo but..."

there is probably/may be little or no ≠ between groups

P

**3. No differences found (BA)**

7%

P

CIs include the null effect probably due to low power

**2. No statistically significant differences found (BA)**

9.3%

**3. Mortality probably lower but no statistically significant differences found**

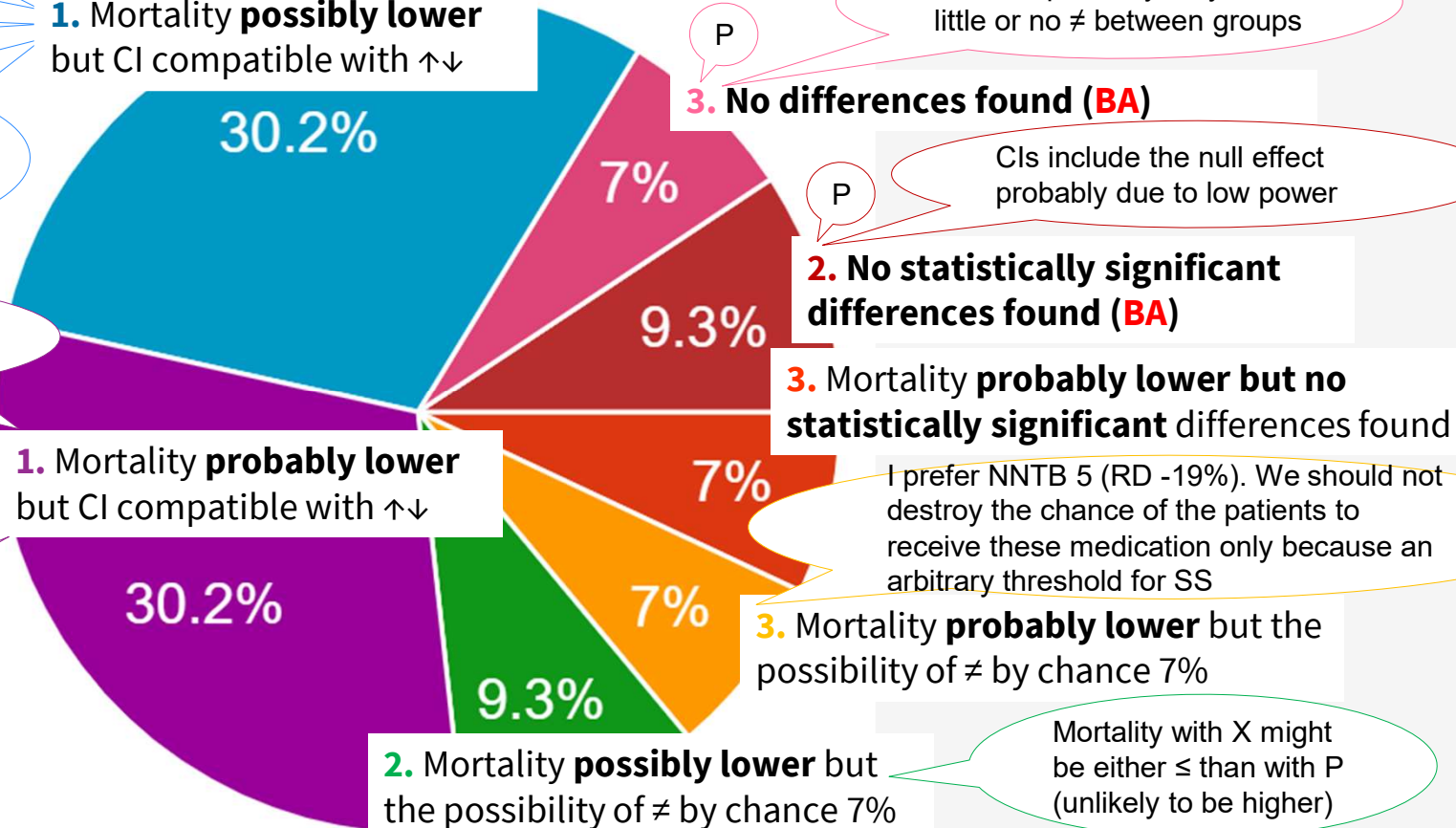
7%

I prefer NNTB 5 (RD -19%). We should not destroy the chance of the patients to receive these medication only because an arbitrary threshold for SS

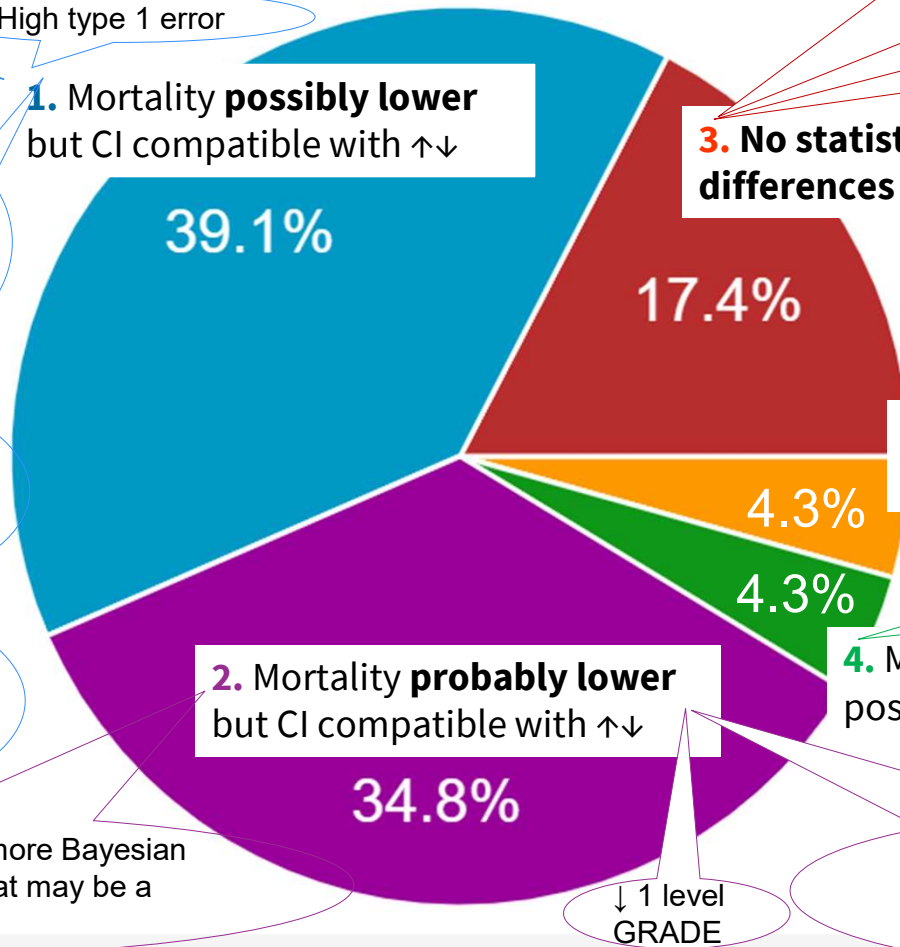
**3. Mortality probably lower but the possibility of ≠ by chance 7%**

7%

Mortality with X might be either ≤ than with P (unlikely to be higher)



# Cochrane Editors (n=23, RR 35%)



Only 79 participants, wide CI,  $\downarrow$  # events

High type 1 error

If the studies were assessed according to the trustworthiness screening tool or the Research Integrity Assessment tool or similar

$\downarrow$  2 levels GRADE guide 34 update on rating imprecision using a minimally contextualized approach

+ explicitly statement that there were NSS  $\neq$  & taking account other outcomes to assist with the interpretation

The probably is driven by a more Bayesian view...but is also possible that may be a small increase in mortality

I would prefer "uncertain" in here or "little or no difference"

Mortality is lower with X, although the evidence is not strong enough. It is more neutral.

3. No statistically significant differences found (**BA**)

4. Mortality **probably** lower but the possibility of  $\neq$  by chance 7%

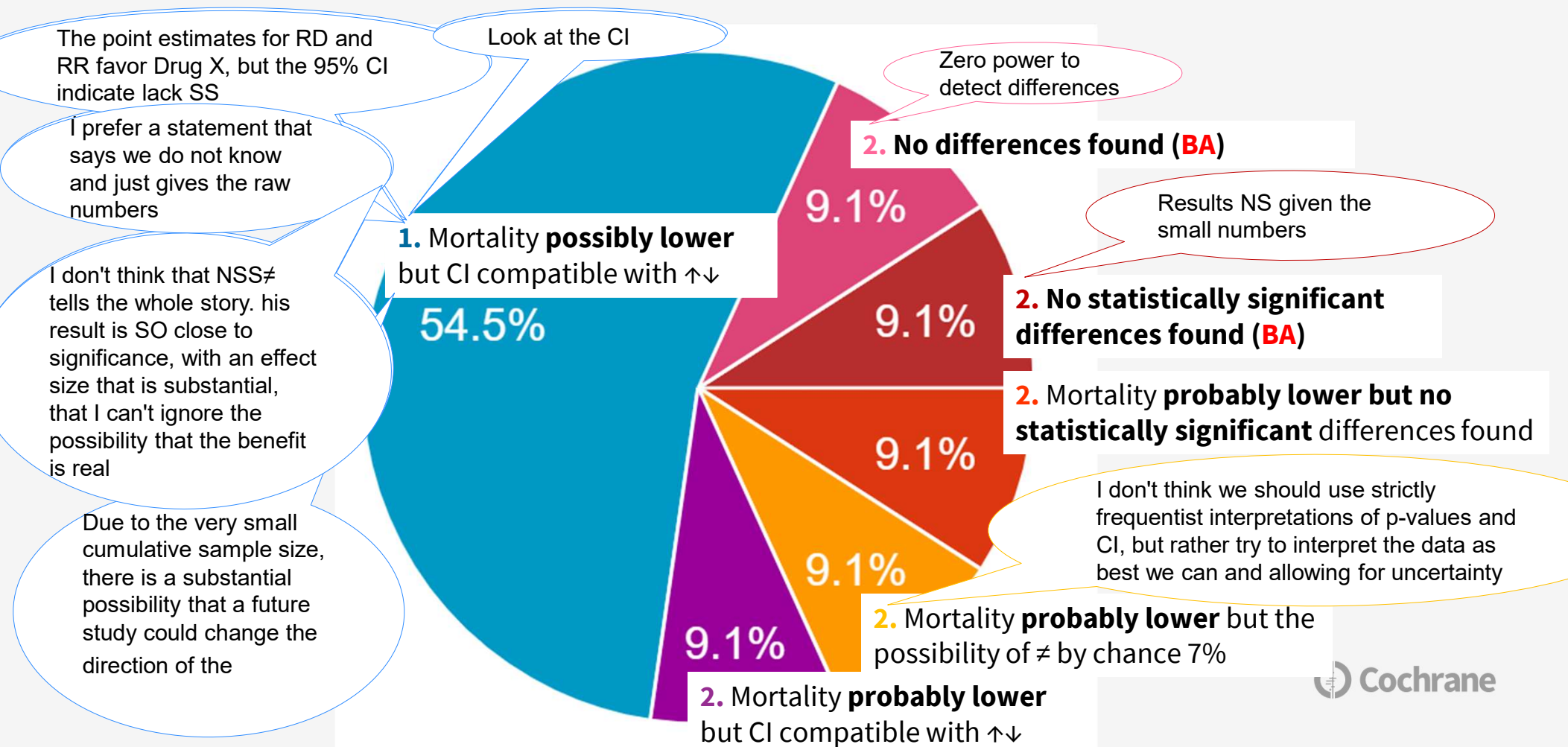
I would prefer insufficient evidence of a  $\neq$

4. Mortality **possibly** lower but the possibility of  $\neq$  by chance 7%

$\downarrow$  1 level GRADE

Probably instead of possibly because of the very large possible benefit vs the very small possible harm (1% $\uparrow$ )

# Non-Cochrane Editors (n=11, RR 55%)



# Non-Cochrane authors (n=28, RR 9%)

↓↓ NNTB but low power

5. Mortality **probably lower** but CI compatible with ↑↓

4. Mortality **possibly lower** but CI compatible with ↑↓

GRADE guide 34

2. Mortality **possibly lower** but the possibility of ≠ by chance 7%

10.7%

6. No differences found (BA)

The 2 studies were described as being of sound methodology. It is not justified to imply that there are differences when the available high-quality data does not support this conclusion.

GRADE guide 34

The biggest concern here is the very small sample size

17.9%

14.3%

3. No statistically significant differences found (BA)

P

7. Mortality **is lower** with X

Some would say this shows a trend, but one would need more data

3.6%

4. Mortality **probably lower** but the possibility of ≠ by chance 7%

10.7%

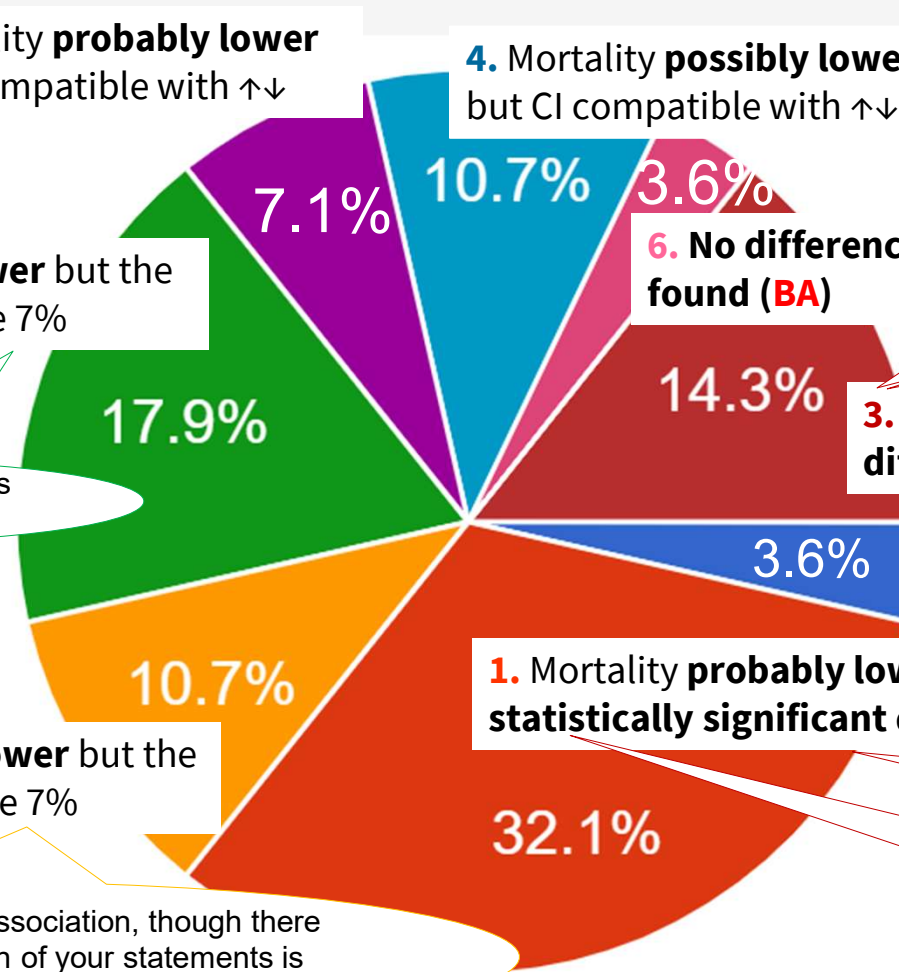
1. Mortality **probably lower** but no statistically significant differences found

CIs include the null effect probably due to low power

This likely represents a true association, though there is not power to conclude which of your statements is most accurate

32.1%

CI crosses 1; > 0.05. The outcome is tremendously important justifying clinical action



**NEW SECTION**

# Conclusions



- There is high heterogeneity of selected statements
- 1/5 Cochrane and non-Cochrane editors & authors still select binary approaches
- The GRADE approach is not always considered to define the certainty of evidence.
- A very low proportion (3%) explicitly considered the GRADE update for rating imprecision (not at all among non-Cochrane editors)
- Probable the best option: *“Mortality may be lower with X than with P, but the CI is compatible with both a reduction and an increase in mortality”*
- Including the probability of chance in the statement is better than only referring to the statistical significance, but it could be also informed by adding the p-value to other effect measures in numbers.

- The case of clinical important difference with CIs crossing the null effect is still an reporting and interpretation challenge.
- The moderate response rate does not warrant representativeness, but suggests that non responders could have a worse performance.
- There are several correct reporting statements and it would be desirable a higher consistency.
- The GRADE update should be strongly diffused.
- Further research should assess the interpretability of the reporting statements.

**Thank you!**

