

ALL-IN meta-analysis: breathing life into living systematic reviews

Judith ter Schure¹ and Peter Grünwald²

¹Corresponding author: Judith.ter.Schure@cw.nl, CWI, Science Park 123, 1098XG Amsterdam, the Netherlands

²Peter.Grunwald@cw.nl, CWI, Science Park 123, 1098XG Amsterdam, The Netherlands/Leiden University, the Netherlands

Abstract Science is idolized as a cumulative process (“standing on the shoulders of giants”), yet scientific knowledge is typically built on a patchwork of research contributions without much coordination. This lack of efficiency has specifically been addressed in clinical research by recommendations for living systematic reviews and against research waste. We propose to further those recommendations with *ALL-IN meta-analysis: Anytime Live and Leading INterim meta-analysis*. ALL-IN provides statistical methodology for a meta-analysis that can be updated at *any time*—reanalyzing after each new observation while retaining type-I error guarantees, *live*—no need to prespecify the looks, and *leading*—in the decisions on whether individual studies should be initiated, stopped or expanded, the meta-analysis can be the leading source of information. We illustrate the method for time-to-event data, showing how synthesizing data at *interim* stages of studies can increase efficiency when studies are slow in themselves to provide the necessary number of events for completion. The meta-analysis can be performed on interim data, but does not have to. The analysis design requires no information about the number of patients in trials or the number of trials eventually included. So it can breathe life into living systematic reviews, through better and simpler statistics, efficiency, collaboration and communication.

Keywords

Anytime, Live, Leading, Interim, Meta-analysis, Efficiency, Collaboration, Communication, Research Waste

The scientific response to the Covid-19 pandemic constitutes a major gamble. In the US, for example, the funding program for vaccine development did not put money on a single vaccine, but on six different ones. They purposely took “multiple shots on goal” according to Larry Corey of the NIH Covid-19 Prevention Network in an interview with STAT (Branswell, 2021). Vaccine development is not a sure thing, and so their strategy needed to be robust enough to just “let the chips fall”. Also in the search for treatments, the scientific community had to hedge its bets. Clinical trials competed for resources and patients, and had to continuously change course when new information arrived. In contrast to vaccines, however, in most countries a strategy to find treatments was lacking. Many clinical trials suffered from “poor questions, poor study design, inefficiency of regulation and conduct, and non or poor reporting of results”: research waste (Glasziou et al., 2020). We believe that more strategic thinking can benefit a future pandemic response as well as non-pandemic evidence-based medicine, as uncertainty is often a given. Honest scientific bets can breathe life into the approach called *living systematic reviews* that aims to keep the evidence record up-to-date (Elliott et al., 2017) and the medical guidelines current (Akl et al., 2017). We propose to make those bets by using *ALL-IN meta-analysis* in clinical trial design, monitoring and reporting.

ALL-IN meta-analysis stands for *Anytime Live* and *Leading INterim meta-analysis*. The *Anytime* aspect provides analysis that controls type-I error in testing and coverage in interval estimation regardless of the decision making along the way, and so regardless of any stopping rules or accumulation bias processes (Ter Schure and Grünwald, 2019). The *Live* aspect prevents research waste caused by meta-analyses that are out-of-date, which is often the case in retrospective meta-analysis. The synthesis can be a bottom-up collaboration of trials, as well as a prospective top-down statistical analysis for decision making. The *Leading* aspect allows the systematically collected evidence included so far to drive the necessity and design of new trials. Finally, the *INterim* aspect is new in meta-analysis and makes for effortless combination of trials while they are still ongoing. What is more, *ALL-IN meta-analysis* is also literally *ALL-IN* since any number of new studies can be included; it has an unlimited horizon. We illustrate this in the setting of time-to-event data, where waiting for new events is an inherent challenge of clinical trials. Combining trials early can prevent delays if studies are slow in themselves to complete the necessary number of events. *ALL-IN* has advantages in four categories: statistics, efficiency, collaboration and communication. We introduce all four briefly (page 4-5) before we go into more detail, but first illustrate the language of betting for single trials studying a Covid-19 vaccine.

A single trial: the FDA game

On June 30th, 2020, the US Food and Drug Administration (FDA) published its guidance document on “Development and Licensure of Vaccines to Prevent Covid-19” (FDA, 2020). This set the goals for any Phase III clinical

trial betting on a protective effect of a vaccine against Covid-19. The guidance document advised on the definition of events of confirmed (symptomatic) SARS-CoV-2 infection for the trials to be counting. And in counting those, the document prescribed the two things to achieve: (1) at least a vaccine efficacy (VE) of 50% and (2) evidence against a null hypothesis of $\leq 30\%$ VE (FDA, 2020, p. 14). Most Covid-19 vaccine trials randomized large numbers of participants 50:50 vaccine:placebo such that we can assume that also throughout the trial the participants at risk stayed approximately balanced. According to the definition of SARS-Cov-2 infections, we start counting once a participant has a confirmed infection after being fully vaccinated for at least a number of days (e.g. 7 days in the Pfizer-BioNTech trial (Polack et al., 2020)). This is also when a (virtual) bet could start. In the following we reinterpret the design for the Covid-19 vaccine trials in the language of betting.

Each new event carries evidence that we express by a betting score. We make a (virtual) investment on one of the two outcomes: either the next event occurs in the vaccine group or it occurs in the placebo group. If there is no effect of the vaccine whatsoever, the 50:50 risk set ensures that the infected participant has 0.5 a chance to be vaccinated and 0.5 a chance to be a placebo. Yet, following the FDA, we do not only want to rule out an ineffective vaccine, but also reject the hypothesis that the vaccine has an effect that is too small—set as the null hypothesis of (at most) 30% VE. In that case each newly observed infection has slightly smaller chance to be a vaccinated participant. That probability to be in the vaccine group is 0.41, since each placebo group member has a 100% risk of Covid-19 and a vaccine group member has $100 - 30 = 70\%$ of the risk, which is a fraction 0.41 of the total risk ($70/(100 + 70)$). So if the VE is too small to be of interest we expect (at least) a fraction 0.41 of Covid-19 events to occur in the vaccine group and (at most) 0.59 in placebo.

How do we bet against that and win if the vaccine has a much larger protective effect? We are betting *against* the probability 0.41 of the next Covid-19 event to occur in the vaccine group. If this probability actually is that large (the vaccine is not very protective; the null hypothesis) we do not want the game to be favorable under any strategy, just like the casino does not want any gambler to earn a salary playing the roulette wheel. On the other hand, we are betting *in favor* of a much smaller probability for the vaccine group. If this probability is smaller (the vaccine is protective; the alternative hypothesis) we do want to win money, just like a professional poker player who makes a salary out of gambling well. We use the betting scores to decide whether the vaccine is a real deal-breaker (the scores behave like the salary of a professional poker player) or whether it is not effective enough (the scores behave like anyone playing the roulette wheel). To ensure that our betting scores can show either case, we first *design* the game such that it is fair—under the null hypothesis—and then *optimize playing* the game with a strategy that is profitable—under the alternative.

Designing a fair game under the null hypothesis Consider gambling at the roulette table where the vaccine trial analogy is like betting on red (vaccine) or black (placebo). Betting correctly doubles your investment, betting incorrectly loses everything you risked. Assuming no house edge and an initial €100 you do not expect to increase your investment, since you have 0.5 a chance of doubling ($2 \cdot €100$) and 0.5 a chance of losing all ($0 \cdot €100$). Whether you bet everything on black or red, in expectation the betting score after one round is $(0.5 \cdot 2 + 0.5 \cdot 0) \cdot €100$, which is the initial investment €100. To achieve the same thing betting against the 0.41:0.59 probabilities instead of 0.5:0.5, your investment needs to multiply by 2.4 ($1/0.41$) for vaccine and 1.7 ($1/0.59$) for placebo. If you bet everything on vaccine you have 0.41 chance of multiplying by 2.4 ($2.4 \cdot €100$) and 0.59 chance of losing all ($0 \cdot €100$) and if you bet everything on placebo you have 0.59 chance of multiplying by 1.7 ($1.7 \cdot €100$) and 0.41 chance of losing all ($0 \cdot €100$). The expected betting score after one round is again the initial investment for both: $(0.41 \cdot 1/0.41 + 0.59 \cdot 0) \cdot €100$ and $(0.59 \cdot 1/0.59 + 0.41 \cdot 0) \cdot €100$. Hence, at either the roulette table or in this FDA game, by design the game is fair and not favorable us. After all, if our observed infections land on the vaccine and control group with the probabilities 0.41:0.59—like a spin of the roulette wheel on black and red with 0.5:0.5—we do not expect to claim an effective vaccine.

Optimize playing the game under the alternative hypothesis How do we win as fast and as much as possible if our observed infections do not behave like a roulette wheel? It has been known since the work of Kelly (1956) and Breiman (1961) that the best way to increase your capital in the long run is to not bet all your (virtual) investment €100 on one of the two possible outcomes (red/vaccine or black/placebo) but to divide it based on the odds that make the game favorable to you. So our focus needs to be on the minimal VE of 50% from the FDA guidance. In the scenario of 50% VE, the probability that the next Covid-19 case is in the vaccine group is 1/3: if we set the risk of Covid-19 for a placebo group member to 100%, a vaccine group member has $100 - 50 = 50\%$ of that risk, which is 1/3 of the total risk ($50/(100 + 50)$). Kelly (1956) and Breiman (1961) urge us to invest one-third ($1/3 \cdot €100$) on observing the next infection in the vaccine group and two-thirds ($2/3 \cdot €100$) on placebo.

Likelihood ratios If we bet this way we can rewrite our betting scores in terms of a likelihood ratio. We first show this for the red-black roulette game where we double what we had put at risk on either black or red if the spin of the roulette wheel outputs the color we bet on. Just like in our strategy in the FDA game, we put $1/3 \cdot €100$ on red and $2/3 \cdot €100$ on black, so we win the following if the ball X lands on either red or black:

$$\begin{aligned}
 X = \text{red} & \quad 2 \cdot \frac{1}{3} \cdot €100 = \frac{\mathcal{L}(1/3 | X)}{\mathcal{L}(1/2 | X)} \cdot €100 \\
 X = \text{black} & \quad 2 \cdot \frac{2}{3} \cdot €100 = \frac{\mathcal{L}(1/3 | X)}{\mathcal{L}(1/2 | X)} \cdot €100
 \end{aligned}$$

The Bernoulli 1/3-likelihood $\mathcal{L}(1/3 | X)$ assigns likelihood 1/3 when is $X = \text{red}$ and 2/3 when is $X = \text{black}$. So if our strategy is to invest 1/3-2/3 in roulette, our payout is our initial investment €100 multiplied by the likelihood ratio, whether X is red or black.

$$\begin{aligned}
 X = \text{vaccinated} & \quad 2.4 \cdot \frac{1}{3} \cdot €100 = \frac{\mathcal{L}(50\% \text{ VE} | X)}{\mathcal{L}(30\% \text{ VE} | X)} \cdot €100 \\
 X = \text{placebo} & \quad 1.7 \cdot \frac{2}{3} \cdot €100 = \frac{\mathcal{L}(50\% \text{ VE} | X)}{\mathcal{L}(30\% \text{ VE} | X)} \cdot €100
 \end{aligned}$$

The likelihood for 50% VE ($\mathcal{L}(50\% \text{ VE} | X)$) assigns likelihood 1/3 when is $X = \text{vaccine}$ and 2/3 when is $X = \text{placebo}$. Similarly, the likelihood for 30% VE ($\mathcal{L}(30\% \text{ VE} | X)$) assigns likelihood 0.41 when is $X = \text{vaccine}$ and 0.59 when is $X = \text{placebo}$. Hence if our strategy is to invest 1/3-2/3 in the FDA game, our payout is also our initial investment €100 multiplied by the likelihood ratio, whether X is vaccine or placebo.

A winner The Pfizer/BioNTech trial observed 8 cases of Covid-19 among participants assigned to receive the vaccine and 162 cases among those assigned to placebo (Polack et al., 2020). We assume now that we start with an initial (virtual) investment of €1 instead of €100. At the first observation we bet €0.33 on vaccine and €0.66 on placebo. After we observe the event in the placebo group we lose our €0.33 bet on vaccine and multiply our €0.66 on placebo by 1.7 to €1.13. The likelihood ratio between our 30% VE alternative hypothesis and our 50% VE null hypothesis—so $\mathcal{L}(50\% \text{ VE} | X)/\mathcal{L}(30\% \text{ VE} | X)$ —is also about 1.13, so multiplying our initial investment of €1 into €1.13. On the other hand, if we observe the event in the vaccine group we lose our €0.66 bet on a placebo event and multiply our €0.33 on vaccine by 2.4 to €0.81. The likelihood ratio of a vaccine event multiplies our investment by 0.81. After each observed event we reinvest what we have left in the new bet, so multiply that with the next likelihood ratio. With 8 cases vaccine and 162 cases placebo the Pfizer/BioNTech publication could report a total betting score of $0.81^8 \cdot 1.13^{162} \cdot €1$, which is about €118 million (note that 1.13 is really 1.13333...). If someone wins that at the poker table, we have good reason to consider her a professional poker player with a favorable strategy, rather than a lucky beginner (Konnikova, 2020).

Meta-analysis: bottom-up collaboration

The Pfizer/BioNTech trial included more than 43 thousand participants (Polack et al., 2020), which is quite unique for a clinical trial. Usually trials are much smaller, and scientific consensus is built through systematic reviews and retrospectively combining trials in a meta-analysis. ALL-IN is a way to do so by collaborating bottom-up in a strategic way that can be live instead of retrospective. It has advantages in four categories that we will first briefly introduce and then further elaborate on in this paper: statistics, efficiency, collaboration and communication.

Statistics

Not all mRNA vaccines showed such favourable results as the Pfizer/BioNTech vaccine. In a press release [CureVac AG \(2021\)](#) announced that the final analysis of their clinical trial observed 83 events in the vaccinated group and 145 in placebo, so only a 43% VE¹ (our calculations assuming a 50:50 balanced risk set ($r = 1$ in [CureVac AG \(2020\)](#), p. 124)). Their protocol states the FDA goal in terms of a confidence interval that excludes a VE of 30%, adjusted for two interim analyses. That adjusted confidence interval at the final analysis is based on $Z_{\alpha/2}$ -statistic for the nominal level $\alpha/2 = 0.02281$ ([CureVac AG, 2020](#), Table 8). That interval is [25.3%, 57.1% VE] (our calculations; normal approximation interval) and, regrettably, does not exclude 30%. When the chips fell, this trial lost.

Statistical analyses like these are essentially *all-or-nothing*, just as any other $p < \alpha$ analysis. As soon as all the α is spent—either on a few interims and a final analysis or just on one fixed sample size—we cannot continue the trial and perform subsequent analyses without violating the type-1 error rate. This might be a reasonable price to pay in the urgency of a pandemic when multiple vaccines are competing, but it is a very inconvenient property for clinical trials in general. Usually, we do want to reanalyze a clinical trial in combination with other similar trials in a meta-analysis. Yet any $p < \alpha$ procedure is equivalent to setting a rejection region for the test statistic and checking whether the value for the statistic falls within that region. This rejection region is based on a sampling distribution that assumes the number of studies in the meta-analysis, and the number of patients within each study to be fixed in advance. Given such a fixed sample size (but also for any sequential stopping rule that sets a maximum sample size in advance, such as α -spending), there is only one region, and your test statistic is either in it or not. If it is not, you are not allowed to redo the analyses with an increased sample size. This problem is recognized in approaches to control type-I error for living systematic reviews ([Simmonds et al., 2017](#)). But also if the meta-analysis is not updated (often), the α is essentially already spent on the individual trial analyses. A different way to see this is by the actual sampling distribution of trials in a meta-analysis: any data-driven decision within the series—whether to accumulate more studies and when to perform the meta-analysis—changes the sampling distribution and invalidate the fixed-sample-size distribution assumed for $p < \alpha$. Hence hardly any meta-analysis has valid type-I error control, when the accumulation of trials is based on strategic decisions that introduce accumulation bias ([Ter Schure and Grünwald, 2019](#)).

ALL-IN meta-analysis is not *all-or-nothing* but can still combine all available studies. In fact, it allows any number of new studies or patients to be included without ever spending all α . In terms of gambling, we can keep betting our virtual investment because we never lose everything.

¹The [CureVac AG \(2021\)](#) press release reports a VE of 48%, so uses a different r (ratio of follow-up time in the two groups). In such large trials r can often be assumed to stay close to 1, so we set it to 1 to make all calculations simpler. All our calculations can be found via [Appendix B](#).

The [CureVac AG \(2021\)](#) results, for example, would have accumulated a betting score of $0.81^{83} \cdot 1.13^{145}$. $\text{€}1 = \text{€}1.84$. This single trial is not very profitable, but at least it still preserves some evidence to reinvest in the next trial, such that we can continue to observe evidence and express it by betting on additional observations in a new trial. An ALL-IN meta-analysis can always continue testing the null hypothesis—with type-I error control—and estimating the confidence interval—with coverage guarantees. Importantly, for these tests and intervals the procedures are exactly the same no matter what decisions—so-called stopping rules, or accumulation bias processes ([Ter Schure and Grünwald, 2019](#))—are at play.

Efficiency

Lack of efficiency has been addressed in clinical research in many ways. Not only in the proposal of living systematic reviews ([Elliott et al., 2017](#)), but also in encouragements to present new studies in the context of existing evidence ([Young and Horton, 2005](#)), in advice to design new trials based on systematic reviews and meta-analysis ([Chalmers and Lau, 1993](#); [Lau et al., 1995](#); [Sutton et al., 2007](#); [Goudie et al., 2010](#); [Lund et al., 2016](#)) and in pleading to prevent the “scandal” of wasteful research into clinical questions that are already answered or not of primary importance ([Altman, 1994](#); [Chalmers and Glasziou, 2009](#); [Glasziou and Chalmers, 2018](#); [Glasziou et al., 2020](#), “research waste”). These calls have not been completely ignored, since clinical research has seen an increase in efficiency—e.g. in platform trials or adaptive meta-analysis ([Tierney et al., 2021](#))—whenever collaboration is deemed possible prospectively. Nevertheless, most clinical trial data is synthesized retrospectively, and still deserves all of the above recommendations. ALL-IN meta-analysis enables these data-driven decisions that can make science more efficient. New studies can be easily informed by the synthesis of all data so far such that exactly the right number of patients are randomized to answer a research question, no more and no less. Moreover, an ALL-IN meta-analysis can give an account of the evidence at anytime and therefore facilitate prioritizing new studies, if more than one line of research needs additional data, but not all can be funded.

Collaboration

ALL-IN meta-analysis can be a *live* meta-analysis, since it does not matter how many studies will eventually be combined or which study will contribute most data. Whether it is based on summary statistics ([Tierney et al., 2021](#)) or on individual patient (IPD) data ([Polanin and Williams, 2016](#)), involvement in the same meta-analysis facilitates discussion between those running trials in the same line of research; especially if the line of research can be concluded early. Trial protocols and statistical analysis plans can be exchanged and scrutinized, to identify discrepancies between the design and the conduct of trials. In an ongoing meta-analysis, trials can be selected for inclusion before investigators are unblinded to the results, which helps to mitigate the problems of publication bias and

p-hacking. If IPD analysis is possible, intense collaboration might also prevent mistakes and fraudulent data that would otherwise depreciate the meta-analysis.

A meta-analysis benefits from homogeneity. With too much heterogeneity, it can be very disheartening to update a random-effects meta-analysis, since many trials are needed to precisely estimate the between trial variation and overcome it (Sutton et al., 2007; Kulinskaya and Wood, 2014; Jackson and Turner, 2017). Close collaboration might prevent unnecessary heterogeneity, if trial investigators are involved in the selection of trials in the meta-analysis; especially if they can advise on the design and conduct of new trials and align inclusion criteria and endpoint definitions. A fixed-effects meta-analysis can conclude the research effort early. Sufficient homogeneity may be possible in close collaboration.

Communication

The language of betting The interpretation of evidence in terms of a betting score might help to communicate the uncertainty in statistical results. As Shafer (2021) puts it: “When statistical tests and conclusions are framed as bets, everyone understands their limitations. Great success in betting against probabilities may be the best evidence we can have that the probabilities are wrong, but everyone understands that such success may be mere luck.” Thinking in terms of bets also helps to understand when statistical analyses can be *anytime-valid*. If they are of the *all-or-nothing* kind, but reanalyzed in a meta-analysis, they are gambling while broke. (This intuition can be made mathematically precise; see the description of Neyman-Pearson testing in terms of betting Shafer (2021) and Grünwald et al. (2019).) Yet if we add new studies to an ALL-IN meta-analysis, we are reinvesting the betting score that we saved from earlier studies, to evaluate whether the strategy in those earlier studies continues to succeed. Just like when reinvesting your profits in a casino from one slot machine into another, the notion of winning stays the same. Our evidence against the hypothesis of a *fair* casino does not change when we alternate slot machines. It does not change if we use the score so far to decide on alternating them or to decide when to cash out. If the slot machines are fair, any strategy of playing them is not expected to make money, and our notion of type-I error control holds under any dependency on past results (stopping rules or accumulation bias processes).

Other communication Those uncomfortable with the language of betting can also easily resort to any of three more familiar notions of statistical communication. Firstly, the likelihood ratios/betting scores and their generalizations, so-called *e-values* (Grünwald et al., 2019; Vovk and Wang, 2021), can be interpreted as conservative *p*-values by taking their inverse. If we denote any betting score or *e-value* by ϵ (e.g. $\epsilon = \text{€}1.84$ for the CureVac trial data), then $p < 1/\epsilon$ is a conservative *p*-value (e.g. $p = 1/1.84 = 0.54$ for the CureVac trial data). If we communicate the *p*-value $p = 1/\epsilon$ anyone can test by comparing $p < \alpha$ but with the addition that this conservative *p*-value

is anytime valid² and so $p < \alpha$ can never spend all α (it is never an *all-or-nothing* test). Secondly, the likelihood ratios have their own notion of evidence in the likelihood paradigm (Royall, 1997). Just as well as stating that the Pfizer/BioNTech trial (Polack et al., 2020) multiplied €1 to almost €118 million and the CureVac AG (2021) trial multiplied €1 to €1.84, we can state that their data was almost 118 million times and 1.84 times more likely if we assume the FDA’s goal of 50% VE in comparison to assuming only 30% VE. For Pfizer, that sounds very good, for CureVac, not so much, and so these numbers have an interpretation of their own without imposing any α level. Thirdly, likelihood ratios can be accepted by the Bayesian paradigm, as Bayes factors, and possibly combined with prior odds. Grünwald et al. (2019) and Grünwald (2021) show that betting scores/*e-values* and Bayes factors are closely related, although not all Bayes factors are betting scores/*e-values*. The bottom-line for communication purposes is that the reporting by ALL-IN meta-analysis can be interpreted in many ways—*p*-values, likelihood ratios, Bayes factors—but regardless of the interpretation provide fully frequentist type-I error control for tests and coverage for confidence sequences.

The remainder of this paper discusses the four categories of advantages in more detail: *Statistics* in Section 1, *Efficiency* in Section 2, *Collaboration* in Section 3 and *Communication* in Section 4. We use the Covid-19 vaccine trials as running examples, based on the FDA game described already, but also in terms of the *safe/e-value logrank test* (Ter Schure et al., 2020b). We also briefly discuss an actual ALL-IN meta-analysis in Section 3.1, that used this *safe/e-value logrank test* to study whether the BCG vaccine could protect against Covid-19 (Van Werkhoven et al., 2021, ALL-IN-META-BCG-CORONA). In the concluding section we will provide some broader context, with an overview of all the methods already developed—*e-values*, *safe tests* (Grünwald et al., 2019) and *anytime-valid* confidence sequences—methods already available in software—notably safestats R package (Ly et al., 2021)—and future work. R code for all calculations, simulations and plots is online available via Appendix B.

1 Statistics

The language of betting comes with the intuition that winning a large betting score has small probability if the null hypothesis is generating our observations (e.g. the roulette wheel is fair). We will make this intuition precise and show how to control the type-I error by bounding this probability by Markov’s inequality and Ville’s inequality. Crucial here is that the betting score underlying our test is an *e-value*. The language of betting also comes with the intuition that when playing a game that is favorable to us in principle, we can use strategies of different quality: even among all strategies under which we expect to get richer, some of them can be expected to earn us much

²Such conservative *p*-values cannot be pictured as the tails of a sampling distribution since such a picture needs a sample size. Appendix A gives more details.

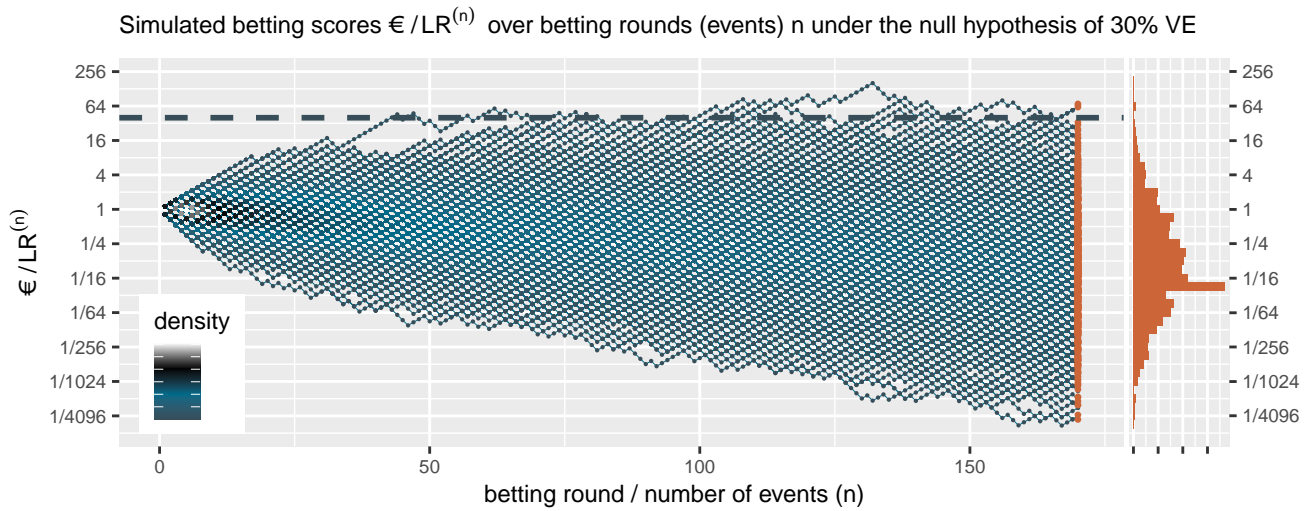


Figure 1. 1000 simulated betting scores in the FDA game over betting rounds n assuming a probability of 0.41 (70/170) for each event to occur in the vaccine group (the null hypothesis of 30% VE). The dashed line is the threshold $1/\alpha = 40$ one-sided. The histogram at the right shows the betting score/ $LR^{(170)}$ after 170 events. Note that the expectation of 1 of the scores is not the mode of its distribution nor its median and that the y-axis is on a log scale.

more than others. We will relate the more well known notion of power to such a different notion of *optimality*. In the following we discuss both *e-values* and *optimality* first for a single trial (in the FDA game and more generally) and then for ALL-IN meta-analysis. We conclude by a generalization of optimal *e-value* tests to confidence sequences.

1.1 Under the null: e-values in a single trial

To make the FDA game fair we imposed a multiplication by 2.4 (or 170/70) if we observe the event in the vaccine group and 1.7 (or 170/100) if we observe it in the placebo group. This multiplication has expectation 1 (or smaller) if we assume the null hypothesis of a vaccine with negligible VE of 30% (or smaller). In case of 30%, we have probability 0.41 (or 70/170) to observe a vaccine event and probability 0.59 (or 100/170) to observe placebo, so in expectation we multiply our investment by: $70/170 \cdot 170/70 + 100/170 \cdot 170/100 = 1$. No matter how we invest in the two outcomes, (e.g. putting 1/3 on vaccine and 2/3 on placebo, or something different) in expectation under the null we multiply the initial investment by 1. This means that our betting score is an *e-value*, since by definition an *e-value* is the outcome of a nonnegative random variable with expectation 1 under the null hypothesis (Grünwald, 2021).

Our betting score could also be rewritten as a likelihood ratio, so the expectation of the likelihood ratio ($\mathcal{L}(50\% \text{ VE} | X) / \mathcal{L}(30\% \text{ VE} | X)$) is 1 as well. We henceforth write the likelihood ratio after n rounds of betting (or after observing n events) as $LR^{(n)}$, with for the FDA game

$$LR^{(n)} = \prod_{i=1}^n \frac{\mathcal{L}(50\% \text{ VE} | X_i)}{\mathcal{L}(30\% \text{ VE} | X_i)}. \quad (1)$$

Using its expectation of 1, Markov’s inequality bounds the probability of observing a large multiplication of our in-

vestment ϵ (a large likelihood ratio) by α after $n = 170$ rounds as follows:

$$P_{30\% \text{ VE}} [LR^{(170)} \geq 1/\alpha] \leq \frac{E_{30\% \text{ VE}} [LR^{(170)}]}{1/\alpha} = \frac{1}{1/\alpha} = \alpha.$$

Figure 1 shows at the right side the histogram of betting scores in the FDA game after 170 events when we simulate events under the null hypothesis, with probability 0.41 to occur in the vaccine group, corresponding to 30% VE. A line is shown at 40, and indeed no more than $\alpha = 1/40 = 2.5\%$ of the scores seem to be larger than that threshold. In fact, in these 1000 runs of simulation only 0.3% of the runs have betting score larger than 40; Markov’s inequality is a loose bound. We also have a stronger result because we obtained our betting score over events by multiplying the score of the rounds (see (1), corresponding to reinvesting our winnings), called Ville’s inequality. We get the following from Ville (1939):

$$P_{30\% \text{ VE}} [LR^{(n)} \geq 1/\alpha \text{ for some } n] \leq \alpha.$$

Ville’s inequality is also illustrated in Figure 1: if we take the sequence of rounds into account, still only a few out of the 1000 simulations ever reach a betting score larger than 40. In fact, in these 1000 runs of simulation only 1.1% of the runs have a betting score that is larger at any round in the game, such that our type-I error is controlled at $\alpha = 2.5\%$ at any time. Moreover, this type-I error control is not tied to this maximum number of 170 events, but continues to hold with an unlimited horizon. If the game is fair, no strategy exists that would allow the player to win forever (Ville (1939, p. 89, “ne gagne pas indéfiniment”) as cited by (Shafer, 2020)). Making a large profit in such a fair game casts doubt on the null hypothesis and is captured by a likelihood ratio that grows away from 1: a large betting profit is obtained if the null likelihood is performing worse than alternative.

When trials can be summarized as bets Before they can be combined in a meta-analysis, individual trials are often characterized by the summary statistics from trial publications. Conventional meta-analysis combines these statistics (e.g. mean differences and standard deviations) in a Z -statistic (Borenstein et al., 2009). Unlike the vaccine/placebo outcomes that we have seen so far, such a Z -statistic has a continuous density and cannot be summarized by separately dealing with all possible outcomes. Fortunately, Shafer (2021) shows that any likelihood ratio of distributions can be viewed as a betting score in a game with initial investment €1. This is possible because likelihood ratios have expectation 1 in general if we assume the null hypothesis in the denominator of the ratio. For a Z -statistic we have two normal distributions with variance 1, one with mean μ_0 under the null hypothesis, and one with μ_1 under the alternative. If the data is generated by the null model, the expectation of the likelihood ratio is

$$E_{Z \sim \phi_{\mu_0}} \left[\frac{\phi_{\mu_1}(Z)}{\phi_{\mu_0}(Z)} \right] = \int_z \phi_{\mu_0}(z) \frac{\phi_{\mu_1}(z)}{\phi_{\mu_0}(z)} dz = \int_z \phi_{\mu_1}(z) dz, \tag{2}$$

which is 1 since $\phi_{\mu_1}(z)$ is a probability density that integrates to 1. This means that any such likelihood ratio for a Z -statistic is an e -value and can be used to construct tests by betting.

Not all summary statistics can be assumed to form a Z -statistic with a normal distribution. Fortunately for the logrank statistic this is reasonable (Ter Schure et al., 2020b) if studies are large and the effect size not too extreme (hazard ratios not too far away from 1). We will use the logrank Z -statistic as a running example for meta-analysis on summary statistics. For an IPD meta-analysis (on individual patient data), however, we recommend to use the exact safe/ e -value logrank test from Ter Schure et al. (2020b) that is valid regardless of the randomization (e.g. 1:1 balanced or 1:2 unbalanced), the number of participants at risk, the number of events or the size of the effect—so also for a hazard ratio 0.05 that corresponds to a VE of 95%.

1.2 Under the null: e -values in a (live) meta-analysis

Assume we want to perform a meta-analysis and we collect a Z -statistic Z_i from each trial i , e.g. a logrank statistic. Before observing Z_i we construct an honest bet $\mathbf{LR}_i = \phi_{\mu_1}(Z_i)/\phi_{\mu_0}(Z_i)$ for each trial that is an e -value and thus has type-I error control under the null hypothesis ϕ_{μ_0} —for a default logrank statistic this is always $\mu_0 = 0$ corresponding to hazard ratio of 1. If we think of the betting score from the first study and invest it in the second study, we are in fact multiplying likelihood ratios. We need to have a notion of time t , such that at each time we know the number of studies $k(t)$ so far and the number of observations $n_i(t)$ in each study i . If we assume that all studies are completed at time t with n_1, n_2, \dots, n_i events summarized by logrank Z -statistics $z_1^{(n_1)}, z_2^{(n_2)}, \dots, z_k^{(n_k)}$ we

can construct our ALL-IN bet as follows:

$$\mathbf{LR}_{\text{META}}^{(t)} = \prod_{i=1}^{k(t)} \mathbf{LR}_i^{(n_i)} = \prod_{i=1}^{k(t)} \frac{\phi_{\mu_1 \sqrt{n_i}}(z_i^{(n_i)})}{\phi_0(z_i^{(n_i)})}. \tag{3}$$

The global null hypothesis Each trial bet is testing the same null hypothesis $\mu_0 = 0$ in (3), such that the ALL-IN meta-analysis bet tests a *global null hypothesis* of no effect (0% VE) in all trials. Such a global null hypothesis can be rejected with a contribution from each trial, but also in case only one trial observes a large score betting against the hypothesis and no other trial observes a very small betting score that loses those winnings again. After all, the null in each trial is rejected as soon as the null is rejected in one of the trials.

Meta-analysis on interim data We can generalize this ALL-IN meta-analysis bet of completed trials to bets on interim data by assuming that we only have an interim logrank Z -statistic $z_1(t), z_2(t), \dots, z_k(t)$ for the $n_1(t), n_2(t), \dots, n_k(t)$ events observed so far at time t ; $k(t)$ still represents the number of studies so far at time t , but now these studies are not (all) completed. We construct our ALL-IN bet in a similar way:

$$\mathbf{LR}_{\text{META}}^{(t)} = \prod_{i=1}^{k(t)} \mathbf{LR}_i^{(n_i(t))} = \prod_{i=1}^{k(t)} \frac{\phi_{\mu_1 \sqrt{n_i(t)}}(z_i(t))}{\phi_0(z_i(t))}. \tag{4}$$

From the perspective of Ville’s inequality, the analysis on completed trials and the one on interim data are indistinguishable. The only thing that matters is that we include all the data we have so far at time t , such that we have type-I error control

$$P_0 [\mathbf{LR}_{\text{META}}^{(t)} \geq 1/\alpha \text{ for some } t] \leq \alpha, \tag{5}$$

for the global null hypothesis probability P_0 with an unlimited horizon over time t .

1.3 Under the alternative: optimality in a single trial

A power analysis sets a very specific goal for a trial, usually to detect an effect of minimal clinical relevance. This is the effect we would not like to miss if it were there, although we hope that the real effect is larger. We nevertheless use this smallest effect of interest to decide on the sample size of the trial, otherwise we risk a futile trial. The FDA was clear on what this minimal effect should be for the Covid-19 vaccine trials: a VE of 50% (FDA, 2020). This is the effect we used to bet in the FDA game.

Our strategy in the FDA game, however, was not trying to achieve optimal power. If we compare the *all-or-nothing* confidence interval for CureVac AG (2021) from the introduction—the final analysis on 83+145 events—we notice that this confidence interval [25.3%, 57.1% VE] is smaller than the final anytime valid interval we show in Figure 3 in Section 1.5, which is Note that we are comparing a $Z_{\alpha/2}$ -confidence interval for $\alpha/2 = 0.02281$ with $\alpha/2 = 0.05$, so the wider interval cannot be attributed to the level of α . The difference is that the former one is

optimized to have spent all α at the final analysis, while the latter one is optimized to continue data collection. Power is the probability of finding the desired result using the specified analysis at a sample size or stopping rule. So for an analysis that is intended to have unlimited horizon, power is not a well-defined concept. Instead Grünwald et al. (2019) introduced the concept of *growth-rate optimality in the worst case*, or *GROW*. Here, the goal is to optimize the expected rate at which the evidence grows (or the interval shrinks) for each new data point, not at a specific sample size. The worst case here is the 50% VE for a one-sided alternative hypothesis $H_1 = \{P_{VE} : 50\% \leq VE \leq 100\%\}$. We optimized the FDA bet in the introduction by putting this 50% VE in the alternative likelihood. This can be rewritten in terms of a likelihood ratio for the logrank statistic Z as follows:

$$\begin{aligned} \text{LR}^{(n)} &= \prod_{i=1}^n \frac{\mathcal{L}(50\% \text{ VE} | X_i)}{\mathcal{L}(30\% \text{ VE} | X_i)} = \frac{\mathcal{L}(50\% \text{ VE} | X_1, \dots, X_n)}{\mathcal{L}(30\% \text{ VE} | X_1, \dots, X_n)} \\ &\approx \frac{\phi_{\mu_{\min} \sqrt{n}}(Z^{(n)})}{\phi_{\mu_0 \sqrt{n}}(Z^{(n)})}, \end{aligned} \tag{6}$$

with $\mu_{\min} = \log(0.5)/4$ and $\mu_0 = \log(0.7)/4$ with 0.5 and 0.7 the hazard ratios corresponding to VE of 50% and 30% respectively (see Ter Schure et al. (2020b)). So our one-sided alternative hypothesis for the logrank Z -statistic is a Z -distribution with a mean representing an effect that is at least μ_{\min} :

$$H_1 = \{\phi_{\mu_1} : \mu_1 \leq \mu_{\min}\}$$

(since positive VE corresponds to a negative μ_1). Our choice of the parameter of the alternative likelihood μ_{\min} follows directly from the minimal effect set by the FDA. Kelly (1956) already showed that this way of betting optimizes the way our betting score grows if the true VE is 50% (our worst-case scenario). Breiman (1961) showed that this approach also minimizes the number of events we need to reach a given betting score set in advance (e.g. $\epsilon 1/\alpha$), for which some intuition is given in Figure 2. Grünwald et al. (2019), Shafer (2021) and Ter Schure et al. (2020b, Online Appendix) give various other reasons why this is the best way to bet, relating it to data compression, information theory, Neyman-Pearson testing, Gibb’s inequality, and Wald’s identity. Most crucial property for the purposes of ALL-IN meta-analysis is that the alternative likelihood puts some money on each possible outcome, such that no matter what outcome we observe, we keep some of the money we risk. This contrasts the approach with a classic $p < \alpha$ test that essentially puts all money on the rejection region, such that if the outcome is not in it, we lose all and cannot continue betting. A thorough interpretation of Neyman-Pearson testing and p -values in terms of betting is given by both Grünwald et al. (2019) and Shafer (2021).

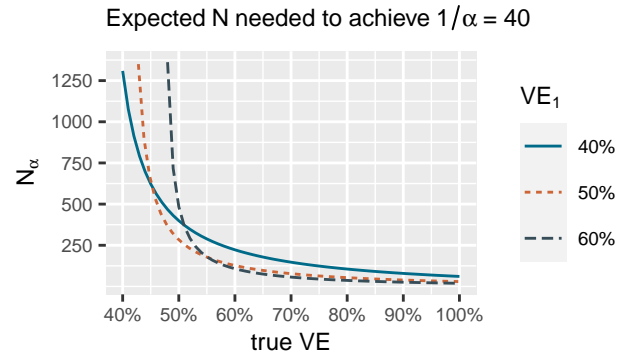


Figure 2. N_α is the expected number of events needed to reach a betting score of $1/\alpha = 40$ for $\alpha = 0.025$ if we bet according to VE_1 indicated by the three different lines, with bets each of the form $\prod_{i=1}^N \frac{\mathcal{L}(VE_1 | X_i)}{\mathcal{L}(30\% \text{ VE} | X_i)}$. The number of events we need decreases if the true VE underlying the data increases (the true difference in risk between vaccine and control is larger). The smallest number of events for a true VE of 40% is reached by betting VE_1 of 40% (blue solid line), the smallest number of events for a true VE of 50% by betting VE_1 of 50% (orange dotted line) and the smallest number for true VE of 60% by betting VE_1 of 60% (grey dashed line). Note that for the alternative in the FDA game $H_1 = \{P_{VE} : 50\% \leq VE \leq 100\%\}$ we are only interested in playing the game well if the true VE is 50% or larger. Since for larger true VE, taking $VE_1 = 50\%$ performs quite well, our strategy is to optimize for the worst case of 50% VE itself and use the bet with $VE_1 = 50\%$ in the FDA game.

1.4 Under the alternative: optimality in a meta-analysis

ALL-IN meta-analysis allows for a retrospective meta-analysis that is bottom-up. The betting score that we accumulate by reinvesting from one trial into the other (which is multiplying betting scores) has an interpretation without enforcing a common design or stopping rule on all included trials. This is especially important if trials have their own stopping rules, or if meta-accumulation processes are at play that influence the existence of trials based on earlier (trial) results in the same meta-analysis. While a meta-analysis can be bottom-up and each have its own design and effect of minimal interest, it can be advisable to agree on a μ_{\min} for the meta-analysis. However, the meta-analysis betting score can also allow each trial i to have its own alternative likelihood with parameter $\mu_{\min(i)}$. Then the following multiplication of those betting scores is still a valid meta score with type-I guarantees:

$$\text{LR}_{\text{META}}^{(t)} = \prod_{i=1}^{k(t)} \frac{\phi_{\mu_{\min(i)} \sqrt{n_i}}(z_i^{(n_i)})}{\phi_0(z_i^{(n_i)})}. \tag{7}$$

As long as $\phi_{\mu_{\min(i)} \sqrt{n_i}}$ is a probability density that integrates to 1, we have that each likelihood ratio integrates to 1 under the global null hypothesis, such that (5) holds. This means that trials can also learn their parameter $\mu_{\min(i)}$ from already completed trials. This is sometimes the case

if trials are not powered to detect an effect of minimal interest, but an effect that is plausibly true based on earlier research. [Kulinskaya et al. \(2016\)](#) shows that such use of existing studies to power new trials can actually bias conventional meta-analysis since it introduces yet another dependency between sample size and results that is unaccounted for in any analysis that assumes a fixed sample size. For ALL-IN meta-analysis this is no problem at all, and trials can learn from each other as long as the parameter $\mu_{\min(i)}$ is fixed before seeing new data that is evaluated using that parameter in (7). In [Ter Schure et al. \(2020b\)](#) we discuss the advantages of even learning the parameter within one trial using prequential plugins or Bayesian posteriors. In a game like the FDA game with a clear goal, this is inferior to the GROW approach, but in other situations it could be preferred.

1.5 Confidence sequences

The [CureVac AG \(2021\)](#) trial reached their final interim analysis but was not able to reject the null hypothesis of 30% VE. The trial had been too optimistic and powered for 60% instead of 50% VE ([CureVac AG, 2020](#)). If a trial is underpowered but still has a large number of participants in follow-up, there is good reason to continue the trial, or combine the trial with results from a new trial in a meta-analysis. However, with a total of 227 events this trial was not underpowered to detect an effect in the same ballpark as the Pfizer/BioNTech trial that reported 95% VE. In such a case it is very interesting to zoom in on the estimate for the effect, instead of its test.

A standard confidence interval can be seen as an inversion of a hypothesis test: if the null falls outside a two-sided 90%-confidence interval it can be rejected with a one-sided type-I error level of $\alpha/2 = 0.05$. In general, the interval excludes all the values for the parameter that can be rejected. Similarly, in our context, an anytime-valid confidence interval excludes all values of the parameter that can be rejected by an e -value test. So the interval is essentially tracking a whole range of bets, each against a different null hypothesis. [Figure 3](#) gives a sequence of anytime-valid confidence intervals for a random ordering of the [CureVac AG \(2021\)](#) data, one for each new observed event or betting round. It shows that the more events we observe, the more parameters (hazard ratios, or their corresponding VEs) we can exclude from the interval. Because these intervals are valid at any time, once we can exclude a value, we never have to include it again. So we also show a sequence of intervals that is the running intersection of all the previous intervals. This of course crucially depends on the ordering, so the one shown for the [CureVac AG \(2021\)](#) data is just an example, since the ordering is not real. Since these intervals are anytime valid, it is possible to further shrink the intervals by continuing follow-up and observing more events. The coverage of an anytime-valid confidence sequence—like an e -value test—has an unlimited horizon.

An ALL-IN meta-analysis confidence interval that is based on a running intersection is of course only possible in an IPD meta-analysis, and cannot be based on summary

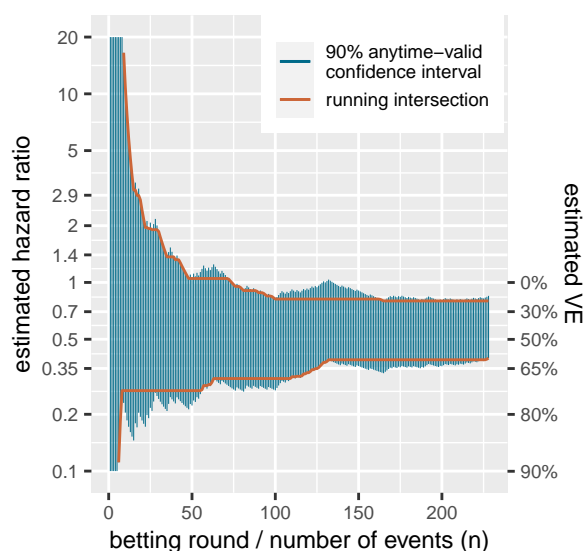


Figure 3. 90%-confidence sequence for a random ordering of the 83 events in the vaccine group and 145 events in placebo from the [CureVac AG \(2021\)](#) trial. Note that the y-axis is on the log scale.

statistics. The confidence interval shown in [Figure 3](#) is based on the logrank Z-statistic (by repeatedly calculating it after each event), which can also be a summary statistic to achieve a single interval that is anytime-valid. The interval follows from the likelihood ratio of normal densities from (6) and follows a general recipe for constructing confidence sequences from [Howard et al. \(2021\)](#) where the hazard ratio is obtained by means of the Peto estimator ([Peto, 1987](#)). The same approach can be used to obtain an ALL-IN meta-analysis confidence interval. A fixed-effects meta-analysis Z-statistic corresponds to a logrank statistic stratified by trial, and an estimate can be obtained from such a logrank statistic that [Peto \(1987\)](#) calls a *typical* hazard ratio. We discuss this approach a bit further in the final section.

2 Efficiency

Trials often suffer from recruitment difficulties ([McDonald et al., 2006](#)) and find themselves underpowered according to their own protocol: when they decide the stop the recruitment and obtain the final sample size for analysis, they have a high probability for their test statistic to fall outside the rejection region they set in advance. This is exactly the scenario where meta-analysis could rescue the line of research by combining multiple underpowered trials. However, the literature on *research waste* ([Chalmers and Glasziou, 2009](#)) and *Evidence-Based Research* ([Lund et al., 2016](#)) shows that we are not using the existing evidence base well to design the new trials needed for conclusion or to interpret new research. ALL-IN meta-analysis makes this very easy to do. It comes with a simple notion of the evidence already collected and what is still needed, and a notion of a new trial's ability to provide that: the implied target. The combination of the two has the capacity to make study design more hon-

Distribution of simulated final betting scores $\text{€} / \text{LR}^{(160)}$ after 160 events under the alternative hypothesis of 60% VE

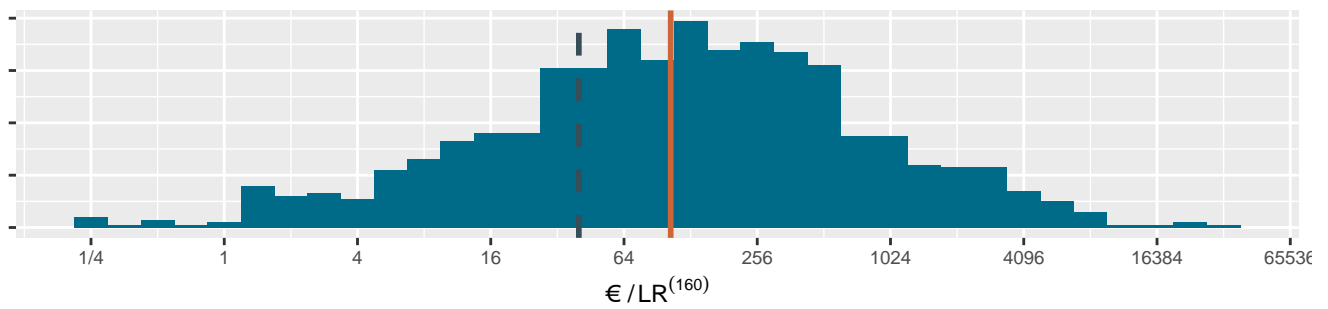


Figure 4. (and **Figure 5**) 1000 simulated sequences of betting scores by round in the FDA game after 160 events assuming a probability of 0.29 (40/140) for each event to occur in the vaccine group. This is the alternative hypothesis of 60% VE used to power the **CureVac AG (2020)** trial at a number of events of 160. The dashed line is the threshold $1/\alpha = 40$ one-sided and the solid line is the implied target of 112. Note that the x-axis is on a log scale.

Simulated betting scores $\text{€} / \text{LR}^{(n)}$ over betting rounds (events) n under the alternative hypothesis of 60% VE

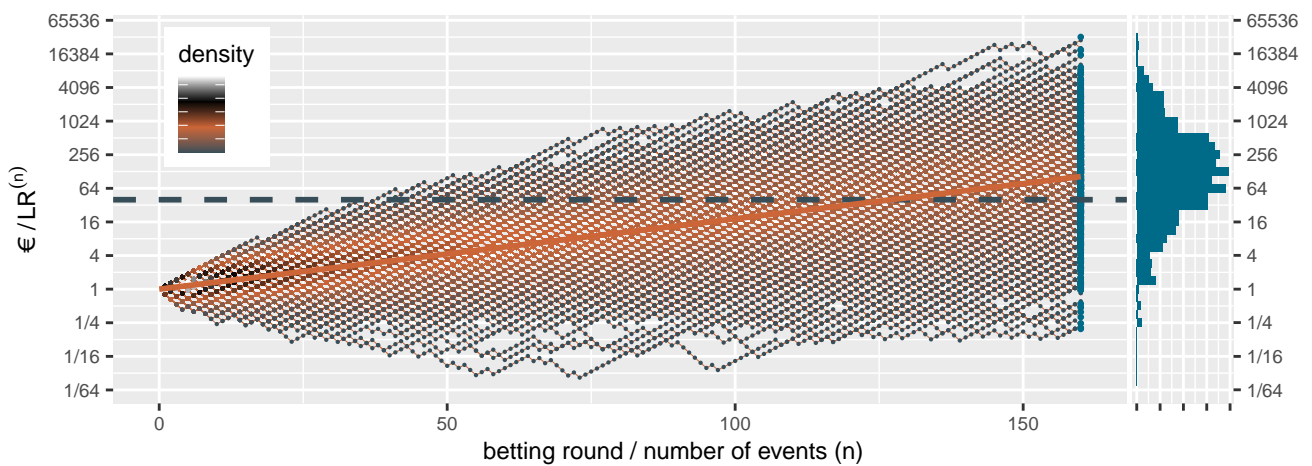


Figure 5. (See above at **Figure 4.**) The histogram for the final betting scores at the right shows the larger scores above and the smaller ones at the bottom, which means that if we turn it, it is the mirror image of the histogram in **Figure 4.** The dashed line is the threshold $1/\alpha = 40$ one-sided and the solid line shows each events'/betting round's contribution to the implied target at $n = 160$ of 112. In this figure, the design has an approximate 79% power to observe a betting score/ e -value larger than $1/\alpha = 40$ before 160 events and 72% power at exactly 160 events (better visible in **Figure 4.**) Note that the y-axis is on a log scale.

est, showing what a trial can add to the existing evidence base instead of just evaluating a misguided goal to single-handedly answer a research question.

2.1 The evidence so far and what is still needed

An ALL-IN meta-analysis can set a prospective goal for conclusion, e.g. $\alpha = 0.0025 = 0.05^2$ corresponding to the level of α required by authorities like the FDA that ask for two trials at the $\alpha = 0.05$ level. Following Ville's inequality (5) we need a betting score of $1/\alpha = \text{€}400$ if we start with $\text{€}1$ to reach a conclusion. Because an ALL-IN meta-analysis combines trials by reinvesting or multiplying betting scores, a very simple calculation gives the betting score we still need at any given point. If an initial trial is able to reach a score of $\text{€}8$, any new trial can be designed to multiply that by 50. So on its own, starting with $\text{€}1$ instead of $\text{€}8$, it would need a betting score of $\text{€}50$ to help the meta-analysis reach $\text{€}400$. We could evaluate

the sample size of the new trial on its ability to reach 50, which for a fixed sample size gives the conditional power of the meta-analysis once the new trial is added. However, if this second trial also foresees recruitment issues, it is more difficult to evaluate its planned contribution since it will probably not be the final trial in the meta-analysis. For this, Shafer (2021) proposed a new notion for the ability of a study, not to reach a specific target betting score, but as a continuous notion of how profitable it can be: the implied target.

2.2 The ability of a new trial: the implied target

The likelihood ratio summarizes the data not in just two categories—statistical significant or not statistical significant—but captures the evidence so far on its way to a certain threshold. Similarly we propose to not evaluate experimental design as all-or-nothing, but summarize its ability to build on what is already there and facilitate

future research. To capture a study's expected contribution to a series of studies, we formulate the *implied target* from Shafer (2021) as the multiplicative amount with which the combined evidence is expected to grow if the study—designed with a certain μ_{\min} and sample size n —is added. In general, the implied target E^* is defined as:

$$E^* = \exp\left(\mathbf{E}_{Z^{(n)} \sim \phi_{\mu_{\min}\sqrt{n}}}\left[\log(\mathbf{LR}^{(n)}(Z^{(n)}))\right]\right). \quad (8)$$

The logarithm appears in equation (8) because the distribution of a likelihood ratio based on n events is very non-symmetric and heavy tailed, with extremely large likelihood ratios occurring with only small probability (see Figure 4). So the expectation of the likelihood ratio is drawn very far from its typical values by these large likelihood ratios and is not a good expression of what to expect. The logarithm makes the distribution more symmetric (asymptotically (for large n) and for normal likelihood ratios even normally distributed), such that the expectation is a more meaningful summary of the evidence promised by the study. By exponentiation ($\exp()$) we bring this expectation back to the scale of the likelihood ratio, such that it can be interpreted as a betting score or e -value.

In the FDA game the expected growth rate per new event in the CureVac trial, assuming their effect of minimal interest of 60% VE is the following:

$$\begin{aligned} & \exp\left(\mathbf{E}_{60\% \text{ VE}}\left[\log\left(\frac{\mathcal{L}(50\% \text{ VE} | X)}{\mathcal{L}(30\% \text{ VE} | X)}\right)\right]\right) \\ &= \exp\left(\frac{40}{140} \cdot \log\left(\frac{50/150}{70/170}\right) + \frac{100}{140} \cdot \log\left(\frac{100/150}{100/170}\right)\right) \\ &= 1.029454. \end{aligned}$$

The cumulative contribution of each new event is shown as the linear line on the log scale in Figure 5. The CureVac AG (2020, Table 8) design planned a final analysis at $n = 160$ events, so their implied target was $1.029454^{160} \approx 104$. In comparison to the target score of €104 at 160 events, the actual betting score €1.84 after $83 + 145 = 228$ events in the press release is quite disappointing. Shafer (2021) gives more examples of how betting scores and implied target help to interpret study results in the context of study design.

2.3 Honest study design

An implied target does require an honest proposal of the effect of minimal clinical interest μ_{\min} , to evaluate the merits of the study. In regular power analysis, this parameter might be tweaked—e.g. setting an unrealistically large effect—to still argue for the study's advancement with only small sample size. Or the smallest effect size of interest analysis is set after data is observed (Wang et al., 2018). This behavior is incentivized by the *all-or-nothing* character of Neyman-Pearson tests that also make the power analysis all-or-nothing. If your desired sample size does not meet the power hoped-for, you need to either increase it or abandon the study. This aspect of traditional analyses fully ignores the ideal of cumulative science in which one study is not expected to single-handedly answer a research question and small increments in knowl-

edge are valuable, as long as they build towards a common goal. If they use e -values and the ALL-IN framework, researchers do not have to view their analysis as the final one, which helps them to evaluate their study more honestly (Lakens, 2021).

3 Collaboration

The *Evidence-Based Research Network* (Lund et al., 2016) aims to always inform new research by past results and to reduce research waste by separating research ideas that are necessary from those that are wasteful. This is not easy to do, however. Different communities might have different notions of necessity or even of what is ethical (a state of so-called *clinical equipoise* (Shamy et al., 2020)). It might therefore be very beneficial to have all those running new clinical trials in a field collaborate together in an ALL-IN meta-analysis.

3.1 ALL-IN-META-BCG-CORONA

We ran two ALL-IN meta-analyses during the Covid-19 pandemic with the involvement of seven trials in one and four in the other. All were designed to study whether the BCG vaccine, originally developed to protect against tuberculosis, could protect against Covid-19 (based on a theory of non-specific immune effects and innate immunity (Netea et al., 2020)). The two meta-analyses study different populations (healthcare workers and elderly) and two questions each: the effect of the BCG vaccine on Covid-19 infection (not necessarily symptomatic) and the effect on severe Covid-19 (indicated by hospitalizations). In the following description we will focus on the analysis of Covid-19 infections in the healthcare workers population.

ALL-IN-META-BCG-CORONA followed many of the steps also outlined by Tierney et al. (2021), that we will briefly discuss here: (1) Meta-analysis design, (2) Systematic search for trials, (3) Systematic review for trial inclusion (4) Data upload, and (5) Disseminating results.

(1) Meta-analysis design Early in the project we decided to aim for an IPD meta-analysis on interim data and wrote our protocols and statistical analysis plans. This timestamped two important decisions on the meta-analysis design: the hazard ratio of minimal interest of 0.8 (20% VE) for events of Covid-19 and the level of α set at 0.0025 so the threshold for the e -value was at $1/\alpha = 400$. For these decisions we set up a meta-analysis Steering Committee that was still fully blinded to any results at the time. The design was preregistered (Van Werkhoven et al., 2021) and all documentation and a webinar explaining the methodology were made available on a project website (Ter Schure et al., 2020a).

(2) Systematic search for trials We continuously searched for trials to include in the meta-analysis. Some were already known to our Steering committee before we started. They initiated a BCG trial of their own very early in the pandemic and shared their protocol with many of their contacts in the BCG research community. Other trials were found by a repeated systematic search of trial

ALL-IN-META-BCG-CORONA

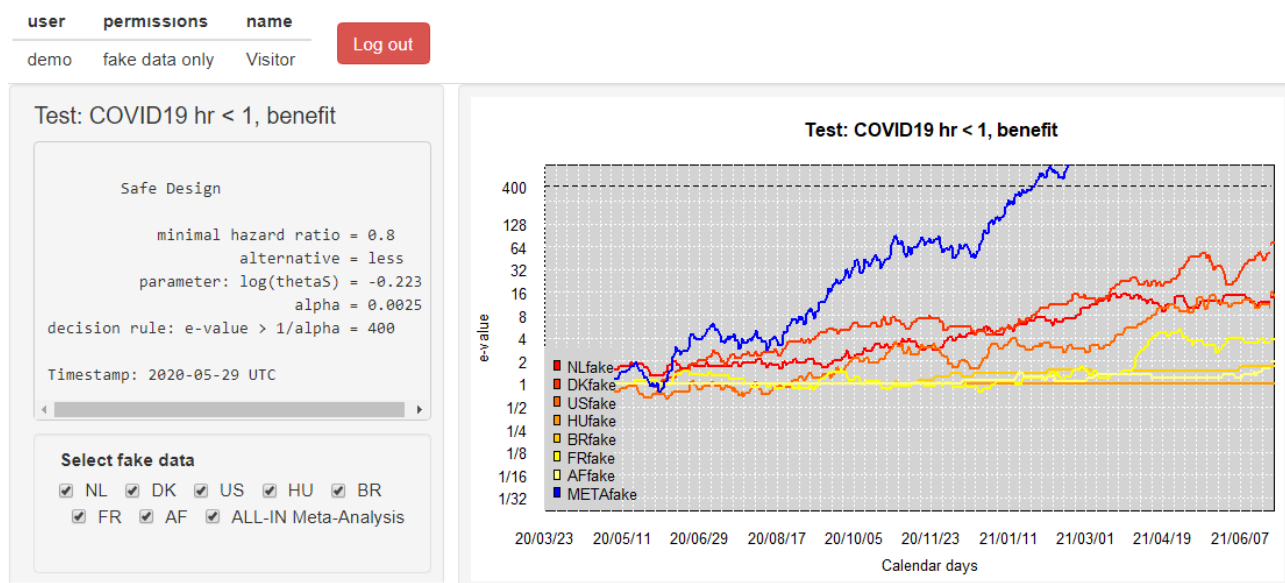


Figure 6. Dashboard used to communicate interim results in *ALL-IN-META-BCG-CORONA* to all data uploaders with a login. The involved trials were performed in the Netherlands (NL), Denmark (DK), the United States (US), Hungary (HU), Brazil (BR), France (FR) and Guinea-Bissau/Mozambique (AF). The dashboard is in demo mode and shows fake data. The option to (de)select trials is for plotting purposes of individual trial e -values; all trials in the dashboard stay included in the meta e -value, following the decision from the Steering committee on trial inclusion. Note that the y-axis is on the log scale.

registries. The trials that agreed to join the meta-analysis were each represented by a member in the Advisory Committee. Meetings of the Advisory committee were scheduled regularly and the trials involved could point us to any new developments. A major advantage of ALL-IN meta-analysis here is that the number of trials does not need to be specified in advance.

(3) Systematic review for trial inclusion We received external advice from Cochrane Netherlands on trial inclusion based on a thorough risk-of-bias assessment. For this assessment, each trial shared their protocols, and subsequently all Cochrane’s evaluations were shared and discussed with the Steering committee and Advisory committee (where trials were usually represented by their PI’s who were blinded to any trial results). Trials had multiple opportunities to answer questions—from Cochrane Netherlands as well as other trials involved—explain their trial and express other concerns about differences between the trials included. The Steering Committee made the final decision on including a trial, before any of that trial’s results were known to anyone part of the discussion. The decision of the Steering committee explicitly incorporated both trial quality and meta-analysis homogeneity.

(4) Data upload Parallel to the discussions on trial inclusion, data transfer agreements were signed and data was shared through a secure upload. Each trial had a data uploader that was in close contact with the ALL-IN meta-trial statistician (the first author of this paper) about data quality. The ALL-IN statistician did not attend the discussion meetings and kept the Steering committee and Ad-

visory committee blinded to any results before each trial inclusion decision.

(5) Disseminating results Each data-uploader received a dashboard account with permissions to inspect the meta-analysis e -value and their own trial contribution. Their access of interim meta-analysis results in the dashboard served as a motivator to keep their own trial data upload up-to-date and to check the sequence of e -values for errors. Figure 6 shows this dashboard based on a demo login with fake data (this demo was available for everyone involved to get an impression). After an initial period where the data-uploaders could only inspect their own trial results, they granted each other permission to inspect all the individual trial contributions. When the first trials were completed and the meta-analysis was approaching its conclusion, the results were also presented to the Advisory and Steering committees. Any interim decisions were planned based on the ALL-IN meta-analysis e -values, but results were also presented in the context of power, implied target per trial and confidence sequences.

3.2 Collaborative and bottom-up meta-analysis

In many aspects, our approach agrees with the *FAME* proposal from Tierney et al. (2021). *FAME* argues for prospective meta-analyses in close collaboration with ongoing trials to achieve the same advantages outlined in this paper, such as aligning trial characteristics (“minimize heterogeneity”) and reducing publication bias and “bias [in] both review and meta-analysis methods” introduced by “prior knowledge of trial results” (e.g. accumu-

lation bias (Ter Schure and Grünwald, 2019)). In alignment with the FAME recommendations, ALL-IN-META-BCG-CORONA was prospectively designed by preregistering an overall effect size of minimal interest and an α -level. However, ALL-IN meta-analysis in general also allows for a more bottom-up approach when each trial's e -value is based on that trial's own design (effect size of minimal interest, see Section 1.4) and trial evidence is synthesized more loosely without a strict decision rule. In comparison to FAME, ALL-IN meta-analysis is much more adaptive. FAME proposes to use conventional meta-analysis (with a fixed sample size) and optimize the timing of the meta-analysis “to anticipate the earliest opportunity for a potentially definitive meta-analysis”. In that sense, FAME can only adapt to the speed of recruitment, while ALL-IN allows to adapt to any information so far including the evidence in the trials and the synthesis of the meta-analysis itself. There is also a statistical inconsistency in the FAME approach that is concerned with “striking a balance between maximising the absolute and relative information size and producing a sufficiently timely review” but does not explicitly state that the last step of the meta-analysis is to “assess the value of updating the systematic review and meta-analysis”. A fixed sample-size statistical analysis should not be reanalyzed using the same statistical methodology. Any proposal to use conventional meta-analysis for efficiency purposes risks accumulation bias (Ter Schure and Grünwald, 2019) because the timing of the meta-analysis might be driven by some of the results part of that same analysis. Hence the best approach is to combine the recommendations from FAME (Tierney et al., 2021) with statistical approach from ALL-IN meta-analysis and the spirit of living systematic reviews.

Collaboration using a dashboard A dashboard for ALL-IN meta-analysis allows us to spot trends in the accumulating evidence, or allow other stakeholders to monitor. A dashboard like Figure 6 can give access to the accumulating e -values to those that need to prepare for crossing a threshold in the near future, e.g. for independent data monitoring committees of ongoing trials or for those considering new trials or preparing to update medical guidelines. On a log-scale, the increase in e -values is linear (in expectation) and the observed trends can be estimated, e.g. in Figure 6 as an increase in evidence per additional calendar day.

For ALL-IN-META-BCG-CORONA, the time unit t in the definition of $LR^{(t)}$ from (4) was set to calendar days and the e -values were updated at each calendar day with an event. The dashboard plots in Figure 6 horizontal lines at 1 for trials that do not observe any events yet: they have not started betting and are still at their initial investment of €1 contributing a neutral amount to the multiplication meta- e -value. ALL-IN meta-analysis monitors e -values as events come in, also when they do so from multiple trials simultaneously. In the language of betting, even the analysis of simultaneous events is considered a sequential bet. If the bet on the events from one trial pays out €4, it multiplies our initial capital by 4, and if the events

from another trial pay out €5, it does so by a factor 5. Yet if we actually consider those trials to be consecutive bets, we reinvest the €4 from the first into the second, and obtain $€1 \cdot 4 \cdot 5 = €20$, as follows from the definition of the meta-analysis e -value on interim data in (4).

Figure 6 illustrates what going ALL-IN means: the evidence in all studies can be monitored and compared to the required threshold at any time. The hypothesis test is carried out by comparing the meta e -value in blue to the threshold $1/\alpha$ of 400, plotted as a dotted line. Because the meta-analysis is *anytime* and *live* a conclusion is reached whenever the e -value sequence passes that threshold. The synthesis of studies can efficiently lead the decision to stop recruiting, treat the placebo group or discourage new trials to start, while encouraging inspection of each individual trial's contribution to the meta-analysis. Since each trial's contribution is a simple multiplication, their components can often be conveniently spotted in the agreement of the shape of the meta-analysis and individual trial lines in a dashboard like Figure 6 (as long as not too many trials are contributing simultaneously).

Collaboration in a competitive field or a pandemic

ALL-IN meta-analysis also prevents losing type-I error control when many trials compete for answers on the same research question, e.g. in an uncoordinated scientific response to a pandemic. If trials are only evaluated in isolation and a response follows the first positive result of a single trial, serious multiple testing issues arise that inflate the type-I error and result in unreliable inference and, subsequently, poor decisions. This happens especially if all trials perform interim analyses on their own, and a type-I error occurs at an interim analyses before any other trial results are published to refute it. The example dashboard also clearly demonstrates decreased type-II errors: synthesizing the evidence in a meta-analysis at interim stages of the trials, and not after trials are completed, improves the ability to find an effect early. Collaboration is indeed much more efficient.

3.3 Fixed-effects and random-effects meta-analysis

Sutton et al. (2007, p. 2491) note that “in a meta-analysis with considerable heterogeneity, the impact of a new (large) study will be (much) less in a random compared to fixed effect model”. This is due the incorporation of a parameter in the model that represents the between-study variation. Also Kulinskaya and Wood (2014) find that the goal of sequentially updating a random-effect meta-analysis might involve planning a large number of small trials to estimate the between-study variance well. Even if that is considered advisable, a random-effects model result might still be very difficult to interpret (Riley et al., 2011). Hence there are various reasons to prefer the fixed-effects model to monitor evidence efficiently and to ensure that the trials are sufficiently homogeneous.

Alongside ALL-IN-META-BCG-CORONA we initiated a second ALL-IN meta-analysis. While the first included trials on healthcare workers, the second included trials in the elderly. Early in the process, before seeing any data,

our Steering committee noticed that the two groups could be very different. Based on a theory of innate and trained immunity, they expected a different effect of the BCG vaccine on the younger immune system of healthcare workers than on the older immune system in the elderly. It could even be that the BCG vaccine effect was beneficial in the ability to fight off Covid-19 in one population but harmful in the other. In general, the differences between trials can be in three categories: heterogeneous effects, conflicting effect and multiple testing.

Heterogeneous effects Our Steering committee decided that to declare success, all included trials in healthcare workers should observe an effect of 20% VE or larger. If they indeed do, heterogeneity in their effect sizes (e.g. one 20%, one 50%, one 25%) does not matter for their joint ability to reject the *global null hypothesis* of no effect in all trials. So for testing the global null, trials are allowed to be heterogeneous in where they are in the space of the alternative hypothesis $H_1 = \{VE : 20\% \leq VE \leq 100\%\}$. For estimation, however, it is not clear what the ALL-IN confidence interval is estimating if we assume that the effects in the trials are very different. Still, as a first summary, a *typical effect size* (Peto, 1987) might be useful if we are unable to estimate a random effects model. The development of confidence sequences for random-effects meta-analysis is a major goal for future work. We do not, however, believe that the evidence in a line of research should be monitored based on whether this interval excludes the null hypothesis, or whether the *e-value* corresponding to the random-effects null model does: for testing, the global null is much more natural. Waiting for a random-effect model to reach a certain threshold is counter-intuitive, since it might require many small trials to estimate the between-trial variability instead of focusing on testing the treatment effect. Moreover, the goal of rejecting the null hypothesis corresponding to this model can be quite strange. When testing a zero-effect null hypothesis, it assumes that there are true effects of harm and true effects of benefit among the trials and that their mean is exactly zero.

Conflicting effects If one of the trials has an effect smaller than 20% or even a harmful effect, we should anticipate betting scores or *e-values* that are smaller than 1. So a meta-analysis multiplication of those *e-values* would reduce the evidence available from other trials. If we can identify groups for which we expect that the trials in each group have an effect in the same direction and of at least the minimal size, we can perform separate meta-analyses. This was the rationale behind grouping healthcare workers and the elderly each in their own ALL-IN-META-BCG-CORONA analysis.

Multiple testing When our analysis is exploratory, and we really have no idea how to group the various trials, we are faced with a multiple testing problem. Note that in this situation also no conventional meta-analysis method would be used to test a common null-hypothesis. We wonder whether any of the trials has the ability to reject the null hypothesis. In that case, we can divide our initial investment over the trials, and see if the totality of their

bet achieves a high betting score. Research into this use of *e-values* has shown that indeed averaging *e-values* is the optimal way to have type-I error control in a standard multiple testing setting (Vovk and Wang, 2021). We return to the notion of hedging bets and averaging *e-values* in Section 4.

Problems with heterogeneity in meta-analysis are not tied to the ALL-IN approach and familiar to anyone working with meta-analysis methods. ALL-IN-META-BCG-CORONA had the advantage that many of the trials that started later had drawn inspiration from the protocol of the first trial. The same sort of alignment of inclusion criteria and outcome definitions might be achieved in other lines of research as well. Hence close collaboration can be very important and the promise of an early conclusion of the research effort might keep a research field motivated to keep the goals aligned.

4 Communication

We have illustrated that the language of betting can be useful in interpreting results from an ALL-IN meta-analysis. Here we argue this further by giving extensions of our method that are very easily explained in terms of betting.

4.1 The language of betting for two-sided tests

Our examples so far covered one-sided tests, but those can be easily extended to two-sided tests, e.g. by taking

$$\text{LR}_{\text{two-sided}}^{(n)} = \frac{1}{2} \cdot (\text{LR}_{\text{left}}^{(n)} + \text{LR}_{\text{right}}^{(n)}),$$

with

$$\text{LR}_{\text{left}}^{(n)} = \frac{\phi_{\mu_{\min(\text{left})\sqrt{n}}}(z^{(n)})}{\phi_{\mu_0}(z^{(n)})} \quad \text{and} \quad \text{LR}_{\text{right}}^{(n)} = \frac{\phi_{\mu_{\min(\text{right})\sqrt{n}}}(z^{(n)})}{\phi_{\mu_0}(z^{(n)})},$$

to represent a two-sided alternative hypothesis

$$H_1 = \{\phi_{\mu_1} : \mu_1 \leq \mu_{\min(\text{left})} \text{ or } \mu_1 \geq \mu_{\min(\text{right})}\}.$$

Such a two-sided test is easy to interpret in the language of betting. We essentially split our initial investment (e.g. €1) between the two sides of the alternative hypothesis (e.g. by betting €0.50 on one side and €0.50 on the other). Any other weighting of the two sides is also possible and corresponds to a different division of the initial investment. The crucial thing is that each side tests the same null hypothesis $H_0 = \{\phi_{\mu_0}\}$ and has expectation 1 under the null hypothesis, such that any weighted average also has expectation 1 and is an *e-value*. Note that for a meta-analysis at time t with $k(t)$ studies this becomes:

$$\text{LR}_{\text{two-sided}}^{(t)} := \frac{1}{2} \left(\prod_{i=1}^{k(t)} \text{LR}_{i,\text{left}}^{(n_i(t))} + \prod_{i=1}^{k(t)} \text{LR}_{i,\text{right}}^{(n_i(t))} \right). \quad (9)$$

Usually one side of the bet is losing and the other is winning such that we do not want to reinvest (multiply) across sides but keep them separate for all trials. In our ALL-IN-META-BCG-CORONA dashboard we also visualized these two sides of the meta-analysis test separately; in Figure 6 we show only the left-sided test (for benefit) of the two.

4.2 The language of betting for co-primary endpoints

Another way to hedge our bets is by considering multiple primary outcomes. In ALL-IN-META-BCG-CORONA, for example, not only the Covid-19 events were counted, but Covid-19 hospitalizations as well, as an indicator for severe disease. We started with $\alpha = 0.05$ and put 10% on Covid-19 ($\alpha = 0.0025$ on each of the two sides of a two-sided test) and 90% on hospitalisations ($\alpha = 0.0225$ on each of the two sides of a two-sided test). So the thresholds to achieve with the e -value for Covid-19 was set at $1/\alpha = 400$ and the one for hospitalization at $1/\alpha = 44.44$. A different way to formulate this is that each had to achieve $1/\alpha = 20$, but that the sequence of e -values for Covid-19 started with an initial investment of €0.05 for each side of the two-sided test (and had to multiply by 400 to reach €20) and that the e -value for hospitalization started with an initial investment of €0.45 for each side (and had to multiply by 44.44 to reach €20).

There are two ways to consider such a bet on two co-primary outcomes: separately and combined. If we evaluate the e -values for each primary outcome separately and reach the threshold with either of the two, we are rejecting the null for that outcome. We are doing two separate tests. If we evaluate the e -values combined, we average them weighted by their α , just as for the two sides of the two-sided test. In that case we have similar type-I error control, but reject the null hypothesis that both are a null effects in favor of the alternative hypothesis that one of them is not. Yet we cannot conclude which one is non-null with the same type-I error since our α level applies to the combined bet and the individual components to the averaged bet are essentially lost.

Concluding remarks

The novelty of this paper lies in a new method for meta-analysis. We do not claim any novelty for the underlying mathematics, though. The basic methods we describe can be viewed as relatively minor variations of the anytime-valid tests that are designed to preserve type-I error under optional stopping, as designed by H. Robbins and his students (Darling and Robbins, 1968; Robbins, 1970). Unfortunately and surprisingly, these tests have not caught on in statistics until a few years ago—right now they are thriving in work on so-called *safe tests*, *anytime-valid confidence sequences* and *e-values* e.g. Shafer et al. (2011); Johari et al. (2021); Pace and Salvan (2019); Howard and Ramdas (2019); Howard et al. (2021); Ramdas et al. (2020); Vovk and Wang (2021); Shafer (2021); Grünwald et al. (2019); Turner et al. (2021); Henzi and Ziegel (2021). As far as we know, it has never before been suggested to use such methods in a meta-analysis context. (Group sequential methods, which have originally also been inspired by the anytime-valid tests, have in turn spurred developments in meta-analysis, but these are substantially different from ALL-IN.) Also, the fact that the logrank test can give a likelihood ratio of the type needed for an anytime-valid test/an ALL-IN meta-analysis is a new finding described by Ter Schure et al. (2020b).

Likelihood ratios, E -variables and e -values

In this paper we presented betting scores/ e -values that are equivalent to likelihood ratios. In general though, betting scores and e -values are really generalizations of likelihood ratios that preserve the properties of likelihood ratios that give them a prominent role in statistic. Entire books have been written to advocate for summarizing evidence in observed data by a likelihood ratio (Edwards, 1974; Royall, 1997) and to separate the goal of measuring evidence from expressing posterior beliefs and making decisions. Likelihood ratios have the property that they can “favor a true hypothesis over a false one more and more strongly” and while a likelihood ratio can be misleading, “strong evidence cannot be misleading very often” (Royall, 1997, p. 14). This latter type-I error control is also referred to as a *universal bound* by Royall (1997) and, by recognizing Ville’s inequality, can be generalized to other betting scores and e -values.

A betting score ϵ is a random outcome of a bet and its random variable is an E -variable if it is nonnegative and for all $P \in H_0$, $E_P[\epsilon] \leq 1$. For a given outcome of the bet, the value of such a random variable is the betting score or e -value. Ville’s inequality relies on the multiplication of E -variables—forming a test martingale—which also has expectation smaller than 1 and thus is itself an E -variable. For the example e -values in this paper, the requirement on the expectation $E_0[\text{LR}] \leq 1$ holds for a simple null hypothesis, e.g. $H_0 = \{\phi_0\}$.

Apart from likelihood ratios of two simple hypothesis, e -values can also be defined for more complicated tests—e.g. a t-test with a nuisance parameter for the variance—in which case the unit expectation needs to hold not for a single mean-0-normal distribution with known variance, but for all mean-0-distributions with any variance. Grünwald et al. (2019) shows that it often is possible to construct E -variables for such composite testing problems, which is why we consider the e -value the right generalization of the likelihood ratio.

Anytime-valid confidence sequences

In this paper we briefly presented a confidence sequence (in Figure 3) for the hazard ratio or VE that was based on the Gaussian approximation to the logrank statistic and the Peto (1987) estimator. This estimator can be derived from summary statistics and is therefore still quite common in meta-analysis as a so-called *two-stage method*, although it is advised against for extreme hazard ratios (Simmonds et al., 2011). Research into other confidence sequences for the hazard ratio is still ongoing. For other estimation problems, confidence sequences already have been thoroughly studied, for example for medians and other quantiles (Howard and Ramdas, 2019), and odds ratios (Turner et al., 2021). These have not, however, been extended to meta-analysis, and especially for the random-effects meta-analysis model, research into confidence sequences is a major goal of future work.

Availability in software

The safestats R package (Ly et al., 2021) provides software to do an *e*-value analysis for the *t*-test, Z-test, logrank test and 2x2-tables. Also functions are available to calculate the power and implied target for these study designs. Confidence sequences can be calculated for the odds ratio in 2x2-tables and the hazard ratio in time-to-event data.

Competing interests

Both authors are proud to have received two shots of the Pfizer-BioNTech Covid-19 vaccine. No further competing interests were disclosed.

Grant information

This work is part of the NWO TOP-I research programme *Safe Bayesian Inference* assigned to Peter Grünwald, with project number 617.001.651, which is financed by the Netherlands Organisation for Scientific Research (NWO).

Acknowledgements

We acknowledge Henri van Werkhoven for his confidence and nerve. Amid the chaos of the initial Covid-19 pandemic months, he quickly grasped the subtleties and benefits of our ideas and committed to intensify a network of Covid-19 research to implement ALL-IN-META-BCG-CORONA. In this partnership, he structured our thinking for this paper. We further acknowledge Marc Bonten and Mihai Netea for the atmosphere of collaboration they put in place in BCG vaccine research and their public stance against “each-small-study-on-its-own” research culture. We also thank Daniël Lakens, Muriel Pérez and Alexander Ly for feedback and extensive discussions. We are especially grateful to Alexander for being the lead developer of the safe/*e*-value logrank test software and co-meta-statistician on ALL-IN-META-BCG-CORONA.

References

Elie A Akl, Joerg J Meerpohl, Julian Elliott, Lara A Kahale, Holger J Schünemann, Thomas Agoritsas, John Hilton, Caroline Perron, Elie Akl, Rebecca Hodder, et al. Living systematic reviews: 4. living guideline recommendations. *Journal of clinical epidemiology*, 91:47–53, 2017.

D. G. Altman. The scandal of poor medical research. *BMJ*, 308(6924):283–284, January 1994. ISSN 0959-8138, 1468-5833. doi: 10.1136/bmj.308.6924.283. URL <https://www.bmj.com/content/308/6924/283>. Publisher: British Medical Journal Publishing Group Section: Editorial.

Michael Borenstein, Larry V. Hedges, Julian P. T. Higgins, and Hannah R. Rothstein. *Introduction to Meta-Analysis*. John Wiley & Sons, Ltd, 2009. ISBN 978-0-470-74338-6. DOI: 10.1002/9780470743386.refs.

Helen Branswell. 12 lessons covid-19 taught us about developing vaccines during a pandemic. <https://www.statnews.com/2021/06/30/12-lessons-covid-19-developing-vaccines/>, 2021. Accessed: 12 July 2021.

Leo Breiman. Optimal gambling systems for favorable games. *Fourth Berkeley Symposium*, 1961.

Iain Chalmers and Paul Glasziou. Avoidable waste in the production and reporting of research evidence. *The Lancet*, 374(9683):86–89, July 2009. ISSN 0140-6736, 1474-547X. doi: 10.1016/S0140-6736(09)60329-9. URL [https://www.thelancet.com/journals/lancet/article/PIIS0140-6736\(09\)60329-9/abstract](https://www.thelancet.com/journals/lancet/article/PIIS0140-6736(09)60329-9/abstract). Publisher: Elsevier.

Thomas C Chalmers and Joseph Lau. Meta-analytic stimulus for changes in clinical trials. *Statistical Methods in Medical Research*, 2(2):161–172, 1993.

CureVac AG. Clinical trial protocol a phase 2b/3, randomized, observer-blinded, placebo-controlled, multicenter clinical study evaluating the efficacy and safety of investigational sars-cov-2 mrna vaccine cvn-cov in adults 18 years of age and older. https://www.curevac.com/wp-content/uploads/2021/06/HERALD_CV-NCOV-004-Protocol.pdf, 2020. Accessed: 16 July 2021.

CureVac AG. Curevac final data from phase 2b/3 trial of first-generation covid-19 vaccine candidate, cvn-cov, demonstrates protection in age group of 18 to 60. <https://www.curevac.com/en/2021/06/30/curevac-final-data-from-phase-2b-3-trial-of-first-generation-covid-19-vaccine-candidate>, 2021. Accessed: 16 July 2021.

DA Darling and Herbert Robbins. Some nonparametric sequential tests with power one. *Proceedings of the National Academy of Sciences of the United States of America*, 61(3):804, 1968.

A.W. Fairbank Edwards. *Likelihood: An account of the statistical concept of likelihood and its application to scientific inference*. Cambridge University Press, New York, 1974.

Julian H Elliott, Anneliese Synnot, Tari Turner, Mark Simmonds, Elie A Akl, Steve McDonald, Georgia Salanti, Joerg Meerpohl, Harriet MacLehose, John Hilton, et al. Living systematic review: 1. introduction—the why, what, when, and how. *Journal of clinical epidemiology*, 91:23–30, 2017.

FDA. Development and Licensure of Vaccines to Prevent COVID-19, 2020. URL <https://www.fda.gov/media/139638/download>. Accessed: 12 July 2021.

Paul Glasziou and Iain Chalmers. Research waste is still a scandal—an essay by Paul Glasziou and Iain Chalmers. *BMJ*, 363, November 2018. ISSN 0959-8138, 1756-1833. doi: 10.1136/bmj.k4645. URL <https://www.bmj.com/content/363/bmj.k4645>. Publisher: British Medical Journal Publishing Group Section: Feature.

Paul P Glasziou, Sharon Sanders, and Tammy Hoffmann. Waste in covid-19 research. *BMJ*, 369, May 2020. ISSN 1756-1833. doi: 10.1136/bmj.m1847. URL <https://www.bmj.com/content/369/bmj.m1847>. Publisher: British Medical Journal Publishing Group Section: Editorial.

Alison C. Goudie, Alexander J. Sutton, David R. Jones, and Alison Donald. Empirical assessment suggests that existing evidence could be used more fully in designing randomized controlled trials. *Journal of Clinical Epidemiology*, 63(9):983–991, September 2010. ISSN 0895-4356. doi: 10.1016/j.jclinepi.2010.01.022. URL <http://www.sciencedirect.com/science/article/pii/S089543561000140X>.

- Peter Grünwald. Peter D. Grünwald's contribution to the Discussion of 'Testing by betting: A strategy for statistical and scientific communication' by Glenn Shafer. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 184(2):440–441, 2021.
- Peter Grünwald, Rianne de Heide, and Wouter Koolen. Safe Testing. *arXiv:1906.07801 [cs, math, stat]*, June 2019. URL <http://arxiv.org/abs/1906.07801>.
- Alexander Henzi and Johanna F. Ziegel. Valid sequential inference on probability forecast performance. *arXiv preprint arXiv:2103.08402*, 2021.
- Steven R Howard and Aaditya Ramdas. Sequential estimation of quantiles with applications to a/b-testing and best-arm identification. *arXiv preprint arXiv:1906.09712*, 2019.
- Steven R Howard, Aaditya Ramdas, Jon McAuliffe, and Jasjeet Sekhon. Time-uniform, nonparametric, nonasymptotic confidence sequences. *The Annals of Statistics*, 49(2):1055–1080, 2021.
- Dan Jackson and Rebecca Turner. Power analysis for random-effects meta-analysis. *Research synthesis methods*, 8(3):290–302, 2017.
- Ramesh Johari, Pete Koomen, Leonid Pekelis, and David Walsh. Peeking at a/b tests: Why it matters, and what to do about it. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1517–1525, 2017.
- Ramesh Johari, Pete Koomen, Leonid Pekelis, and David Walsh. Always valid inference: Continuous monitoring of a/b tests. *Operations Research*, 2021.
- J.L. Kelly. A new interpretation of information rate. *Bell System Technical Journal*, pages 917–926, 1956.
- Maria Konnikova. *The Biggest Bluff: How I Learned to Pay Attention, Master Myself, and Win*. Penguin, 2020.
- Elena Kulinskaya and John Wood. Trial sequential methods for meta-analysis. *Research synthesis methods*, 5(3):212–220, 2014.
- Elena Kulinskaya, Richard Huggins, and Samson Henry Dogo. Sequential biases in accumulating evidence. *Research synthesis methods*, 7(3):294–305, 2016.
- Daniel Lakens. Sample size justification. <https://psyarxiv.com/9d3yf/download?format=pdf>, 2021.
- Joseph Lau, Christopher H Schmid, and Thomas C Chalmers. Cumulative meta-analysis of clinical trials builds evidence for exemplary medical care. *Journal of clinical epidemiology*, 48(1):45–57, 1995.
- Hans Lund, Klara Brunnhuber, Carsten Juhl, Karen Robinson, Marlies Leenaars, Bertil F Dorch, Gro Jamtvedt, Monica W Nortvedt, Robin Christensen, and Iain Chalmers. Towards evidence based research. *Bmj*, 355:i5440, 2016.
- Alexander Ly, Rosanne Turner, Muriel F Pérez-Ortiz, Judith ter Schure, and Peter Grünwald. R-package *safestats*, 2021. Maintainer: Alexander Ly <a.ly@jasp-stats.org>, install in R by `devtools::install_github("AlexanderLyNL/safestats", ref = "logrank", build_vignettes = TRUE)`.
- Alison M McDonald, Rosemary C Knight, Marion K Campbell, Vikki A Entwistle, Adrian M Grant, Jonathan A Cook, Diana R Elbourne, David Francis, Jo Garcia, Ian Roberts, et al. What influences recruitment to randomised controlled trials? a review of trials funded by two UK funding agencies. *Trials*, 7(1):1–8, 2006.
- Mihai G Netea, Evangelos J Giamarellos-Bourboulis, Jorge Domínguez-Andrés, Nigel Curtis, Reinout van Crevel, Frank L van de Veerdonk, and Marc Bonten. Trained immunity: a tool for reducing susceptibility to and the severity of sars-cov-2 infection. *Cell*, 181(5):969–977, 2020.
- Luigi Pace and Alessandra Salvan. Likelihood, replicability and Robbins' confidence sequences. *International Statistical Review*, 2019.
- Richard Peto. Why do we need systematic overviews of randomized trials?(transcript of an oral presentation, modified by the editors). *Statistics in medicine*, 6(3):233–240, 1987.
- Fernando P Polack, Stephen J Thomas, Nicholas Kitchin, Judith Absalon, Alejandra Gurtman, Stephen Lockhart, John L Perez, Gonzalo Pérez Marc, Edson D Moreira, Cristiano Zerbini, et al. Safety and efficacy of the bnt162b2 mrna covid-19 vaccine. *New England Journal of Medicine*, 2020.
- Joshua R Polanin and Ryan T Williams. Overcoming obstacles in obtaining individual participant data for meta-analysis. *Research synthesis methods*, 7(3):333–341, 2016.
- Aaditya Ramdas, Johannes Ruf, Martin Larsson, and Wouter Koolen. Admissible anytime-valid sequential inference must rely on nonnegative martingales. *arXiv preprint arXiv:2009.03167*, 2020.
- Richard D Riley, Julian PT Higgins, and Jonathan J Deeks. Interpretation of random effects meta-analyses. *Bmj*, 342, 2011.
- Herbert Robbins. Statistical methods related to the law of the iterated logarithm. *Annals of Mathematical Statistics*, 41:1397–1409, 1970.
- Richard Royall. *Statistical evidence: a likelihood paradigm*, volume 71. CRC press, 1997.
- Judith ter Schure. Code for paper ALL-IN meta-analysis: breathing life into living systematic reviews, Sep 2021. URL <https://osf.io/d9jny/>.
- Judith ter Schure and Peter Grünwald. Accumulation Bias in meta-analysis: the need to consider time in error control [version 1; peer review: 2 approved]. *F1000Research*, 8:962, June 2019. ISSN 2046-1402. doi: 10.12688/f1000research.19375.1. URL <https://f1000research.com/articles/8-962/v1>.
- Judith ter Schure, Alexander Ly, and Peter Grünwald. Safestats and ALL-IN meta-analysis project page. <https://projects.cwi.nl/safestats/>, 2020a.
- Judith ter Schure, Muriel F Pérez-Ortiz, Alexander Ly, and Peter Grünwald. The safe logrank test: Error control under continuous monitoring with unlimited horizon. *arXiv preprint arXiv:2011.06931*, 2020b. URL <https://arxiv.org/abs/2011.06931>.
- Glenn Shafer. The language of betting as a strategy for statistical and scientific communication. <http://probabilityandfinance.com/articles/54.pdf>, 2020. Online; accessed 24 July 2020.

Glenn Shafer. Testing by betting: A strategy for statistical and scientific communication. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 184(2):407–431, 2021.

Glenn Shafer, Alexander Shen, Nikolai Vereshchagin, and Vladimir Vovk. Test martingales, Bayes factors and p-values. *Statistical Science*, 26(1):84–101, 2011.

Michel Shamy, Brian Dewar, and Mark Fedyk. Different meanings of equipoise and the four quadrants of uncertainty. *Journal of Clinical Epidemiology*, 127:248–249, 2020.

Mark Simmonds, Georgia Salanti, Joanne McKenzie, Julian Elliott, Thomas Agoritsas, John Hilton, Caroline Perron, Elie Akl, Rebecca Hodder, Charlotte Petridge, et al. Living systematic reviews: 3. statistical methods for updating meta-analyses. *Journal of clinical epidemiology*, 91:38–46, 2017.

Mark C Simmonds, Jayne Tierney, Jack Bowden, and Julian PT Higgins. Meta-analysis of time-to-event data: a comparison of two-stage methods. *Research synthesis methods*, 2(3):139–149, 2011.

Alexander J. Sutton, Nicola J. Cooper, David R. Jones, Paul C. Lambert, John R. Thompson, and Keith R. Abrams. Evidence-based sample size calculations based upon updated meta-analysis. *Statistics in Medicine*, 26(12):2479–2500, 2007. ISSN 1097-0258. doi: 10.1002/sim.2704. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/sim.2704>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/sim.2704>.

Jayne F Tierney, David J Fisher, Claire L Vale, Sarah Burdett, Larysa H Rydzewska, Ewelina Rogozińska, Peter J Godolphin, Ian R White, and Mahesh KB Parmar. A framework for prospective, adaptive meta-analysis (FAME) of aggregate data from randomised trials. *PLoS medicine*, 18(5):e1003629, 2021.

Rosanne Turner, Alexander Ly, and Peter Grünwald. Safe tests and always-valid confidence intervals for contingency tables and beyond. *arXiv preprint arXiv:2106.02693*, 2021.

Jean Ville. Etude critique de la notion de collectif. *Bull. Amer. Math. Soc*, 45(11):824, 1939.

Vladimir Vovk and Ruodu Wang. E-values: Calibration, combination, and applications. *Annals of Statistics*, 2021.

Min Qi Wang, Alice F Yan, and Ralph V Katz. Researcher requests for inappropriate analysis and reporting: A US survey of consulting biostatisticians. *Annals of internal medicine*, 169(8):554–558, 2018.

Cornelis H. van Werkhoven, Judith ter Schure, Marc Bonten, Mihai Netea, Peter Grünwald, and Alexander Ly. Anytime Live and Leading Interim meta-analysis of the impact of Bacillus Calmette-Guérin vaccination in health care workers and elderly during the sars-cov-2 pandemic (ALL-IN-META-BCG-CORONA). https://www.crd.york.ac.uk/prospero/display_record.php?RecordID=213069&VersionID=1473878, 2021.

Charles Young and Richard Horton. Putting clinical trials into context. *The Lancet*, 366(9480):107–108, July 2005. ISSN 0140-6736, 1474-547X. doi: 10.1016/S0140-6736(05)66846-8. URL <https://www.thelancet.com/journals/lancet/article/>

PIIS0140-6736(05)66846-8/abstract. Publisher: Elsevier.

A The inverse-conservative p-value

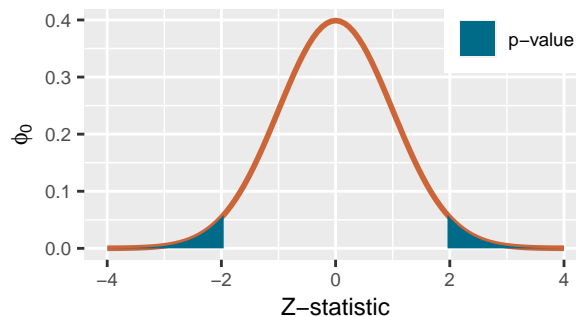


Figure 7. Two-sided Z-test against the null hypothesis $H_0 = \{\phi_0\}$ observing $z = 1.96$ with a p-value of 0.05.

The standard way to introduce p-values in introductory texts is to use a graph like Figure 7 that shades the tail area of a sampling distribution under the null hypothesis for an observed test statistic. In Figure 7 the observed test statistic is a Z-statistic of 1.96, such that the two-sided p-value is 0.05. The precise definition of this p-value for the random variable Z and any observation z (e.g. $z = 1.96$) is:

$$p\text{-value}(z) = \mathbf{P}_{Z \sim \phi_0} [|Z| \geq z] \tag{10}$$

$$\text{such that for all } 0 \leq \alpha \leq 1: \mathbf{P}_0[p\text{-value} \leq \alpha] = \alpha, \tag{11}$$

with \mathbf{P}_0 the probability under the null hypothesis. While explaining-p-by-picture is helpful at first, it also has strong limitations. For example, if the sample size n is not fixed in advance but determined by a fixed, a priori known stopping rule (e.g. stop as soon as the first z comes in that has value > 2.5), then n will itself be random and the threshold for the test statistic will depend on n; there is no evident visualization.³ For this reason, it is common in theoretical statistics to define p-values directly as random variables that have a uniform distribution under the null. From (10) and (11) we see that this is compatible with the introductory approach.

Conservative p-values Defined by property (11), the standard (strict) p-value can be generalized to a (conservative) p-value that does not have a uniform distribution, but whose distribution is stochastically dominated by the uniform distribution under the null hypothesis. This means that the equality in equation (11) is replaced by an inequality:

$$\mathbf{P}_0[p\text{-value} \leq \alpha] \leq \alpha. \tag{12}$$

³An exception is the case in which there are just two possible sample sizes, n_1 and n_2 , with the rule for choosing between n_1 and n_2 based on the first n_1 data points $z^{(n_1)}$ given (e.g. stop if $z^{(n_1)} > 2.5$, continue otherwise). We can then specify a bivariate distribution and evaluate the tail area for the combined observations (e.g. $(z^{(n_1)}, z^{(n_2)})$). This is sometimes illustrated in introductory texts on group-sequential methods, but cannot be extended beyond this (3D) bivariate case.

For such conservative p -values, it still holds that, for any fixed significance level α , the test that rejects if $p \leq \alpha$ has type-I error at most α . The inequality is thus in the right direction to preserve type-I error guarantees.

Our first observation is very simple: for any fixed n , the likelihood ratio $\mathbf{LR}^{(n)}$ is an inverse-conservative p -value for samples of that size n , i.e. (12) holds with ‘ p -value’ set to $1/\mathbf{LR}^{(n)}$. However, simply viewing LRs as inverse conservative p -values is not nearly doing them sufficient justice. A better (yet still incomplete) comparison is to ‘anytime-valid p -values’, a terminology that stems from (Johari et al., 2017, 2021). Whereas standard p -values only make sense for an a priori given, fixed sample size n or a priori given, fixed stopping rule, anytime-valid p -values are really sequences p'_1, p'_2, \dots , one for each sample size n . Even if the stopping rule for determining n is unknown in advance, or unclear, or greedy (‘stop at first n such that $p_n \leq \alpha$ ’), the type-I error guarantee is preserved. It turns out that the inverse of the LRs as defined for ALL-IN meta-analysis can be viewed in this way: the sequence $(1/\mathbf{LR}^{(1)}, 1/\mathbf{LR}^{(2)}, \dots)$ as defined in Section 1 is a sequence of anytime-valid p -values.

B R Code for calculations, simulations and plots

R code for the calculations, simulations and plots in this paper can be found on the Open Science Framework (Ter Schure, 2021, <https://osf.io/d9jny/>), including:

- Settings of the FDA game
- Pfizer/BioNtech results in the FDA game
- CureVac results in terms of confidence interval and the FDA game
- Plotting the FDA game betting scores under the null hypothesis (Figure 1)
- Plotting the expected sample size for various strategies in the FDA game (Figure 2)
- Plotting the anytime-valid confidence sequence for a random ordering of the CureVac data (Figure 3)
- Plotting the implied target of the CureVac design in the FDA game (Figure 4 & 5)
- Plotting a p -value of 0.05 (Figure 6; Appendix)