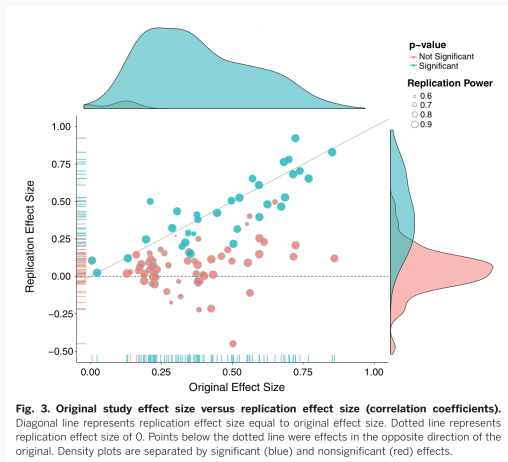


A Coherent Framework for Many Replicability Questions

Ester Alongi

Joint work with Gianmarco Altoè & Giovanni Parmigiani

Original vs. replication effect sizes (Nosek & Errington, Fig. 3)



Five success criteria for replication (Nosek & Errington, Fig. 3)

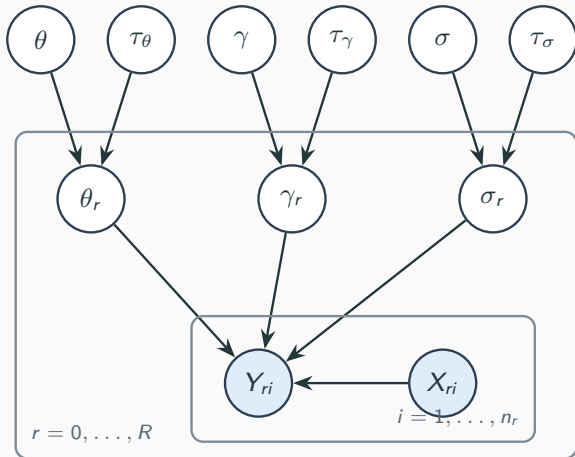
#	Indicator	Question it answers
1	Replication $p < .05$, same direction	Did the replication itself reject the null?
2	Original effect inside replication 95% CI	Is the original estimate compatible with the replication?
3	Replication effect size vs. original	How much smaller (or larger) is the new estimate?
4	Meta-analytic combination of the two	What is the pooled evidence?
5	Subjective "did it replicate?" rating	Holistic judgement by the replicating team

Some Challenges with Existing Methodologies

- **Fragmentation:** Over 50 metrics requiring different models currently exist (Hayes, 2013);
- **Dichotomization:** Replication assessment is collapsed into binary “success/failure” thresholds (Pawel et al., 2024);
- **Empirical Focus:** Studies are compared through horizontal empirical comparisons (Machery, 2020).

- Predictions about replicability are important prospectively, as we plan new studies; reflections on whether replication has occurred are best carried out conditionally on data, based on statements on effects or other parameters
- Whether replication experiments are extensions or not is also an empirical question.
- The same model can provide the backbone for a answering, in an integrated way, many questions of replicability.

Bayesian Hierarchical Model

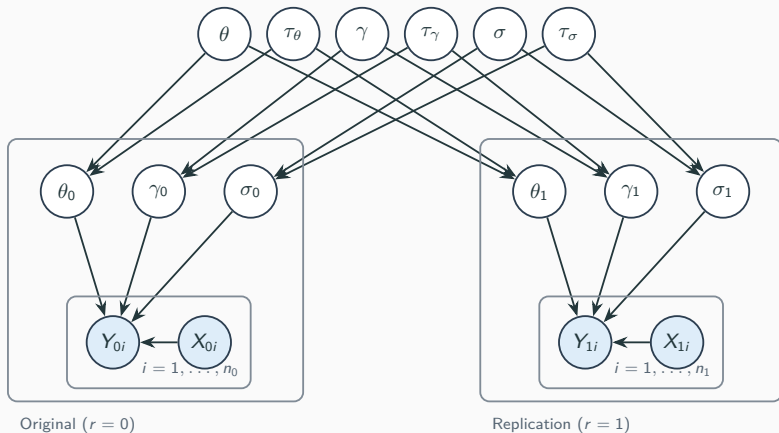


Generative model (generic). $Y_{ri} | \theta_r, \gamma_r, \sigma_r, X_{ri} \sim p_Y(\cdot | \gamma_r + \theta_r X_{ri}, \sigma_r^2)$, with population-level draws

$$\theta_r | \theta, \tau_\theta \sim p_\theta, \quad \gamma_r | \gamma, \tau_\gamma \sim p_\gamma, \quad \sigma_r | \sigma, \tau_\sigma \sim p_\sigma.$$

Hyperpriors on $\theta, \tau_\theta, \gamma, \tau_\gamma, \sigma, \tau_\sigma$ are set by empirical Bayes (next slide).

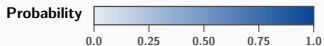
Hierarchical model with one replication ($R = 1$)



Each study draws its own $(\theta_r, \gamma_r, \sigma_r)$ from the shared population-level distributions; precise replication corresponds to $\tau_\theta, \tau_\gamma, \tau_\sigma \rightarrow 0$, so $\theta_0 = \theta_1$ a.s.

- Dependent vs Independent Analyses of Studies
- Replicability of Study Conclusions versus Agreement of Studies with generative mode (inference on which capture the field's knowledge about the effect).
- Prospective planning of new studies for replication or extension

Multiple Facets of Replicability



Replication probabilities for θ_r

$$+\infty \leftarrow \mu_{r\beta}, \mu_{r\alpha}, \mu_{r\sigma}$$

Independence

Dependence

Agreement between studies

$P_{overall}$
$P_{non-null}$
P_{pos}
P_{null}
P_{neg}

Conditional on one study

$P_{cond S_1}$
$P_{cond S_2}$
$P_{cond S_3}$

Agreement between studies

$P_{overall}$
$P_{non-null}$
P_{pos}
P_{null}
P_{neg}

Conditional on one study

$P_{cond S_1}$
$P_{cond S_2}$
$P_{cond S_3}$

Consistency with generative effect

$P_{gen,overall}$
$P_{gen,non-null}$
$P_{gen,pos}$
$P_{gen,null}$
$P_{gen,neg}$

Conditional on generative effect

P_{cpos}
P_{cnull}
P_{cneg}

How consistent is effect size knowledge across studies?

How consistent is effect size knowledge with the generative model?

What do we know about the generative effect

Inference at the meta-analytic level

P_{β}

Consistency between existing evidence and a new study

$P_{new, overall}$

$P_{new, non-null}$

$P_{new, pos}$

$P_{new, null}$

$P_{new, neg}$

$P_{new, cond}$

What is the probability that a new study, arising from the same superpopulation, is consistent with existing evidence?

Definition

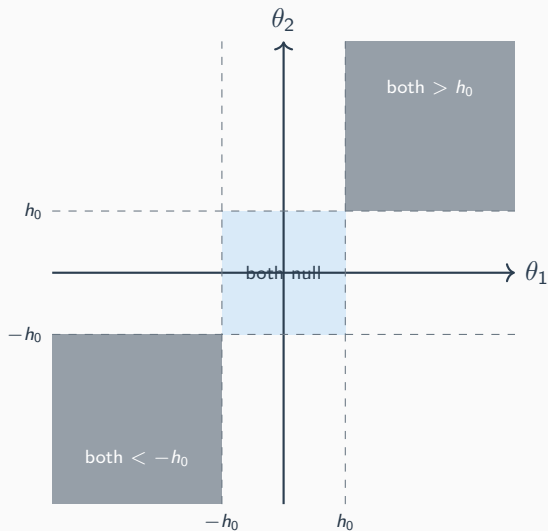
Given a threshold $h_0 > 0$, the **region of practical equivalence (ROPE)** is the interval $(-h_0, h_0)$. The null hypothesis

$$H_0 : |\theta| < h_0$$

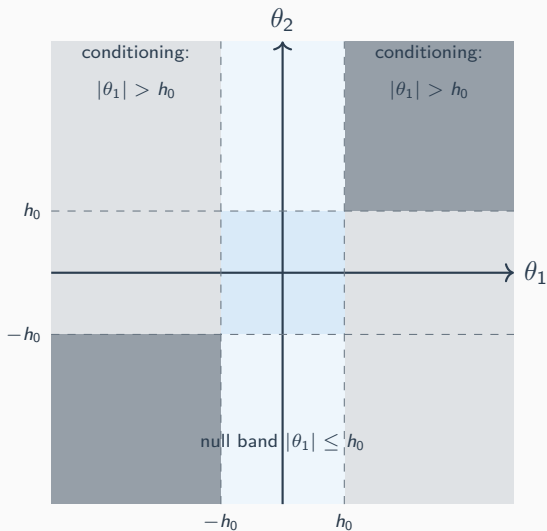
asserts that the effect is too small to matter scientifically; the alternative $H_1 : |\theta| \geq h_0$ asserts that it lies outside the ROPE.

- h_0 is fixed a priori from substantive knowledge (clinical, biological, economic).
- Bayesian view: posterior mass inside $(-h_0, h_0)$ summarises support for the null.
- Frequentist counterpart: equivalence testing against the bounds $\pm h_0$.

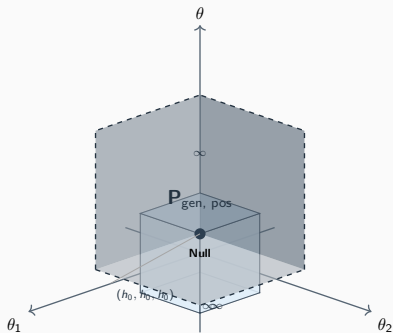
Overall Agreement







Conditional probabilities



$(P_{\text{gen, pos}})$



References

-  Consortium, G. (2020). **The gtex consortium atlas of genetic regulatory effects across human tissues.** *Science*, 369(6509), 1318–1330.
-  Hayes, A. F. (2013). ***Introduction to mediation, moderation, and conditional process analysis: A regression-based approach.*** Guilford Press.
-  Machery, E. (2020). **What is a replication?** *Philosophy of Science*, 87(4), 545–567.
-  National Academies of Sciences, E., & Medicine. (2019). ***Reproducibility and replicability in science.*** National Academies Press.



Pawel, S., Heyard, R., Micheloud, C., & Held, L. (2024). **Replication of null results: Absence of evidence or evidence of absence?** *Elife*, 12, RP92311.

Backup slides

Metrics explanation

Let $t = 1, \dots, T$ index the posterior samples generated by MCMC. For each study s , we obtain posterior draws $\theta_r^{(t)}$, $t = 1, \dots, T$, where $h_0 > 0$ is a prespecified theoretical threshold.

The counts of studies showing positive, negative, or null effects are

$$N_{pos}^{(t)} = \sum_{r=0}^R \mathbf{1}(\theta_r^{(t)} > h_0), \quad N_{neg}^{(t)} = \sum_{r=0}^R \mathbf{1}(\theta_r^{(t)} < -h_0), \quad N_{null}^{(t)} = \sum_{r=0}^R \mathbf{1}(|\theta_r^{(t)}| \leq h_0)$$

Agreement between studies

Overall agreement (P_{overall})

The posterior probability that a consensus majority k (e.g., $k \geq 2$ out of $S = 3$) agree on the direction, or agree on the absence of a relevant effect:

$$P_{\text{overall}, k} = \frac{1}{T} \sum_{t=1}^T \mathbf{1} \left(N_{\text{null}}^{(t)} \geq k \vee N_{\text{pos}}^{(t)} \geq k \vee N_{\text{neg}}^{(t)} \geq k \right)$$

Interpretation: How often do the populations qualitatively agree?

Signed agreement (P_{pos} , P_{neg} , P_{null})

For a specific direction (e.g., positive)

$$P_{\text{pos}, k} = \frac{1}{T} \sum_{t=1}^T \mathbf{1} \left(N_{\text{pos}}^{(t)} \geq k \right)$$

Conditional Replication Probability on one study

Conditional replication probability on one study ($P_{\text{cond}|O}$)

The conditional posterior probability that the effect discovered in the original study O replicates in at least $k - 1$ other studies:

$$P_{\text{cond}|O} = \frac{\sum_{t=1}^T \mathbf{1}\left(\left(\theta_O^{(t)} > h_0 \wedge N_{\text{pos}}^{(t)} \geq k\right) \vee \left(\theta_O^{(t)} < -h_0 \wedge N_{\text{neg}}^{(t)} \geq k\right)\right)}{\sum_{t=1}^T \mathbf{1}\left(|\theta_O^{(t)}| > h_0\right)}$$

Interpretation: Assuming the effect is practically significant in the discovery cohort, what is the probability it exists in at least $k - 1$ other populations?

Consistency with the Generative Effect

Overall consistency ($P_{\text{gen, overall}}$)

How often do the majority of studies and the population-level effect simultaneously agree in sign (or null)?

$$P_{\text{gen, overall, } k} = \frac{1}{T} \sum_{t=1}^T \mathbf{1} \left(\begin{array}{l} (N_{\text{pos}}^{(t)} \geq k \wedge \theta^{(t)} > h_0) \vee \\ (N_{\text{neg}}^{(t)} \geq k \wedge \theta^{(t)} < -h_0) \vee \\ (N_{\text{null}}^{(t)} \geq k \wedge |\theta^{(t)}| \leq h_0) \end{array} \right)$$

Signed consistency ($P_{\text{gen, pos}}$, $P_{\text{gen, neg}}$, $P_{\text{gen, null}}$)

For a specific direction (e.g., positive)

$$P_{\text{gen, pos, } k} = \frac{1}{T} \sum_{t=1}^T \mathbf{1} \left(N_{\text{pos}}^{(t)} \geq k \wedge \theta^{(t)} > h_0 \right)$$

Interpretation: Does the global generative effect θ represent a robust biological consensus across most populations?

Conditional Replication Probability on the generative effect

Conditional replication probability on the generative effect (P_{cpos} , P_{cneg} , P_{cnull})

Conditional on the generative effect θ being practically positive, the probability that this is robustly supported by a consensus of at least k studies is

$$P_{cpos, k} = \frac{\sum_{t=1}^T \mathbf{1}(N_{pos}^{(t)} \geq k \wedge \theta^{(t)} > h_0)}{\sum_{t=1}^T \mathbf{1}(\theta^{(t)} > h_0)}$$

Interpretation: When the true global effect is positive, how often do the specific populations correctly identify and replicate a positive effect?

Simulation Strategy: Isolating Sources of Heterogeneity

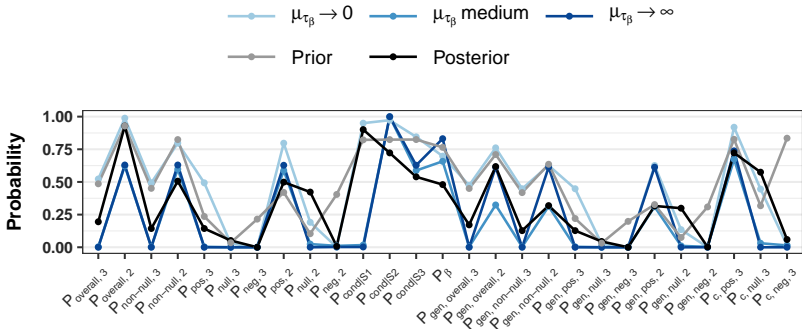
Simulation Scenario	Perturbed Heterogeneity	Fixed Parameters (Real Posterior Means)		
		Gen. Effect (θ)	Baseline Exp. (γ)	Noise (σ)
1. Effect	Effect Heterogeneity (μ_{τ_θ})	–	$\hat{\gamma}_{r,real}$	$\hat{\sigma}_{r,real}$
2. Baseline	Intercept Heterogeneity (μ_{τ_γ})	$\hat{\theta}_{r,real}$	–	$\hat{\sigma}_{r,real}$
3. Residual Variance	Noise Heterogeneity (μ_{τ_σ})	$\hat{\theta}_{r,real}$	$\hat{\gamma}_{r,real}$	–

Note: For the perturbed parameter, the 3 simulated studies are forced at the 10%, 50%, and 90% quantiles of the τ distribution.

Scenario 1: Genetic Effect Heterogeneity

Research Question

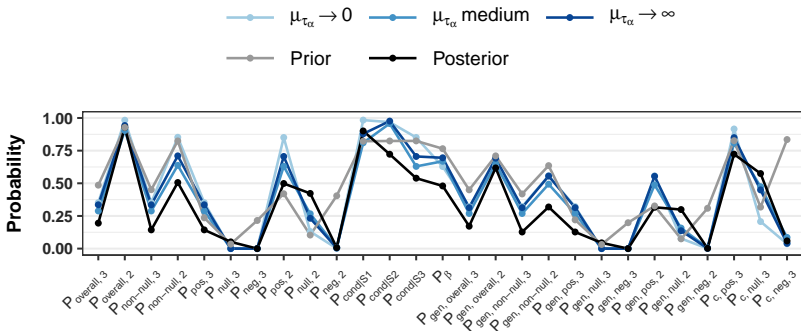
How would replication metrics change if the genetic effect was perfectly consistent across cohorts ($\mu_{\tau_{\theta}} \rightarrow 0$), assuming baselines (γ) and noise (σ) remain at their real-world values?



Scenario 2: Baseline Expression Shifts

Research Question

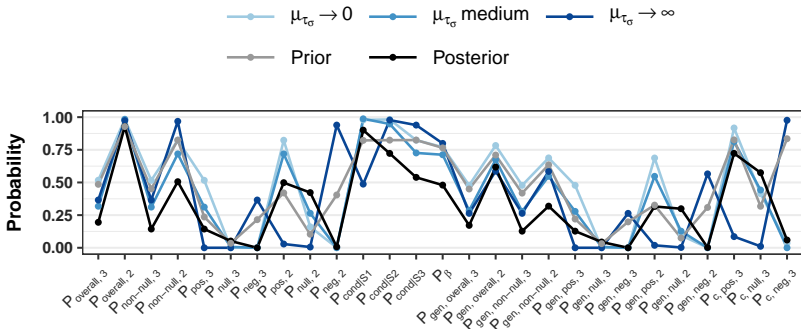
How sensitive is the replication of a genetic effect to systematic differences in baseline gene expression (γ) across populations, assuming constant θ and σ ?



Scenario 3: Unmeasured Noise Heterogeneity

Research Question

If cohorts have vastly different levels of unmeasured environmental/technical noise (σ), what impact does this have on the replicability of the effect of interest?



Empirical Bayes: Hyperpriors distributions

To match empirical variances to a Student-t distribution ($\nu = 3$), we define the scale conversion factor $c = \sqrt{(\nu - 2)/\nu}$.

Path Coefficients and Intercepts

For each triplet k , using weights $w_{rk} = 1/SE_{rk}^2$:

$$\theta_k = \frac{\sum_r w_{rk} \hat{\theta}_{rk}}{\sum_r w_{rk}} \quad \tau_{\theta,k} = c \cdot \sqrt{\frac{\sum_r w_{rk} (\hat{\theta}_{rk} - \theta_k)^2}{\sum_r w_{rk}}}$$

Residual Standard Deviations

For each triplet k , using weights $w_{rk} = df_{rk}$:

$$\sigma_k = \frac{\sum_r w_{rk} \hat{\sigma}_{rk}}{\sum_r w_{rk}} \quad \tau_{\sigma,k} = c \cdot \sqrt{\frac{\sum_r w_{rk} (\hat{\sigma}_{rk} - \sigma_k)^2}{\sum_r w_{rk}}}$$

Global Empirical Bayes Hyperparameters

$$\mu_{\theta} = \text{Mean}(\theta_k)$$

$$\eta_{\theta} = c \cdot \text{SD}(\theta_k)$$

$$\mu_{\tau_{\theta}} = \text{Mean}(\tau_{\theta,k})$$

$$\eta_{\tau_{\theta}} = c \cdot \text{SD}(\tau_{\theta,k})$$

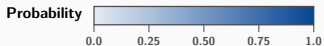
$$\mu_{\sigma} = \text{Mean}(\sigma_k)$$

$$\eta_{\sigma} = c \cdot \text{SD}(\sigma_k)$$

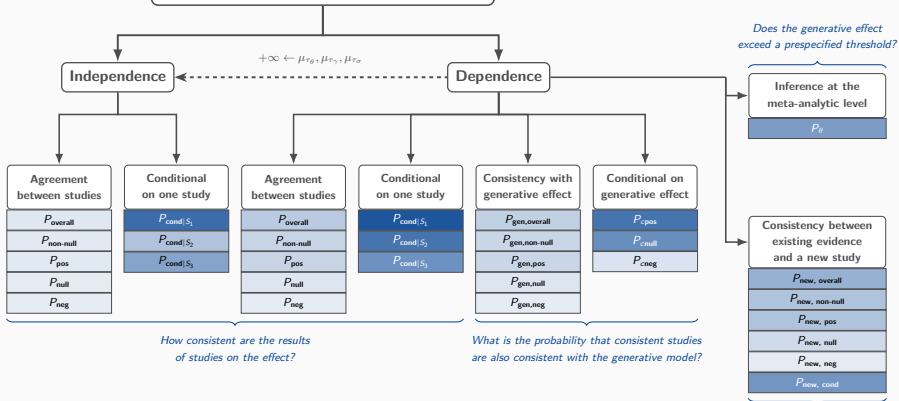
$$\mu_{\tau_{\sigma}} = \text{Mean}(\tau_{\sigma,k})$$

$$\eta_{\tau_{\sigma}} = c \cdot \text{SD}(\tau_{\sigma,k}).$$

Multiple Facets of Replicability



Replication probabilities for θ_r ($k = 3$)



How consistent are the results of studies on the effect?

What is the probability that consistent studies are also consistent with the generative model?

Does the generative effect exceed a prespecified threshold?

Consistency between existing evidence and a new study

What is the probability that a new study, arising from the same superpopulation, is consistent with existing evidence?