

1 **Title:** Robust inference and errors in studies of wildlife control

2
3 **Authors:** Adrian Treves^{1*}, Igor Khorozyan²

4 **Affiliations:**

5 ¹University of Wisconsin, Madison, USA.

6 ²Independent consultant, Göttingen, Germany.

7 *Corresponding author. Email: atreves@wisc.edu

8
9 **Abstract:** Randomized, controlled trials (RCT) are seen as the strongest basis for causal
10 inference, but their strengths of inference and error rates relative to other study have never been
11 quantified in wildlife control and rarely in other ecological fields. We simulate common study
12 designs from simple correlation to RCT with crossover design. We report rates of false positive,
13 false negative, and over-estimation of treatment effects for five common study designs under
14 various confounding interactions and effect sizes. We find non-randomized study designs mostly
15 unreliable and that randomized designs with suitable safeguards against biases have much lower
16 error rates. One implication is that virtually all studies of lethal predator control interventions
17 appear unreliable. Generally, applied fields can benefit from more robust designs against the
18 common confounding effects we simulated.
19

20 **Main Text:** Identifying the cause of a phenomenon often holds the key to developing an
21 effective intervention to interrupt the cause-and-effect connections or improve outcomes. The
22 stakes increase whenever an intervention risks counter-productive effects on the target or side-
23 effects for another valued entity. Therefore, scientific and public scrutiny of outcomes rather
24 than intentions is intensifying in many applied fields [1]. For example, as societies attach more
25 value to wild animals, scrutiny has intensified for interventions aimed at controls intended to
26 protect human interests from wild animals. Recognition of ineffective or counter-productive
27 effects of lethal wildlife control has exposed an alternative to the traditional hypothesis that
28 removing wild animals, e.g., killing gray wolves (*Canis lupus*), might prevent damage to assets
29 or resources [2]. The more recent hypothesis predicts that removing wild animals might
30 exacerbate the losses of property or threats to safety resources [2]. Hence, the field of wildlife
31 control has become increasingly introspective about robust study designs to evaluate the
32 effectiveness of interventions [2-5]. Resolving these uncertainties about wildlife control
33 interventions would advance the fields of human-animal interactions and ethics, including
34 subfields of biodiversity conservation, agricultural or other property protection, and animal
35 welfare. Other applied fields whose interventions may backfire might also benefit from such
36 introspection.

37 ***Quantifying the strengths of inference across study designs***

38 Most investigators advocate the so-called ‘gold-standard’ of randomized, controlled trials (RCT)
39 without biases [6-8]. Yet the urgency of problems may rule against using RCT, exposing tension
40 between swift action and well-informed action [9]. Moreover, RCT can also be infeasible or

41 opposed by interest groups [10, 11], let alone higher standard designs with crossover (within-
42 subject analysis including the reversal of treatment and control conditions for all subjects) and
43 other blinding steps to avoid research and publication biases [2]. Therefore, evaluations of the
44 effectiveness of interventions in many fields often rely on lower standards of evidence than RCT
45 [1, 11, 12]. Drawing inferences from studies with less robust designs than RCT is the norm in
46 studies of wildlife or ecosystems [3, 11, 13], including our field of wildlife control [2-5].
47 Approximately 75% of studies in one review of North American and European wildlife control
48 interventions [5], and an unquantified majority of studies in global reviews of wildlife control [3,
49 14, 15] were non-randomized. Lower standard study designs produce weaker inference because
50 they lack random assignment of treatments and controls or even strict observational controls.

51 Employing the convenient shorthand and ranking RCT as the gold-standard, we refer to the
52 platinum-standard for crossover designs defined as above, and we hypothesize that one could
53 improve the strength of inference in RCT by employing a within-subjects before-and-after
54 intervention[rBACI, for "before-after-control" impact or intervention, depending on how the
55 authors name it [2, 5, 16, 17]]. When non-randomized, we refer to nBACI or the 'silver
56 standard'.

57 The lowest standard in this study is the 'bronze standard' of simple correlation, which compares
58 different doses of intervention and outcomes. This so-called bronze-standard lacks within-
59 subjects comparisons so it introduces additional confounding variables of pre-existing
60 differences between subjects. Therefore, some authors [2, 5] predicted that the gold-standard and
61 higher would outperform the silver- and bronze-standards in strength of inference by a factor of
62 two or more. They further predicted that nBACI would outperform simple correlations and
63 rBACI would outperform RCT, but did not estimate by how much [2].

64 However, randomized designs are not free of concerns [6]. Murtaugh [17] simulated how
65 temporal autocorrelations confounded the interpretation of a treatment effect. Among the
66 concerns, false positive rates (FPR, inferring a treatment effect when none exists) figure
67 prominently, e.g., electric fences are routinely deemed effective in wildlife control when the
68 evidence is fairly weak [4]. FPR are usually under-estimated due to confusion with p-values
69 which do not tell us how often a test or intervention will fail [8, 18]. Also, "new discoveries" in
70 which the null hypothesis of no effect of an intervention is rejected, under the traditional $p=0.05$
71 threshold for statistical significance, have been producing high levels of spurious findings that
72 fail replication attempts, whether or not they use randomized study designs [1]. A short-term
73 remedy might be to lower the threshold for significance to $p=0.005$ for new discoveries. But
74 more importantly, Benjamin et al. [1] urge all applied fields to strengthen inference through more
75 robust study designs with safeguards against research and publication biases.

76 *Simulations to quantify error rates*

77 Here we quantify error rates to compare five study designs and their strengths of inference about
78 the effectiveness of lethal wildlife control interventions, following [11, 12]. The simulations in
79 [12] revealed that sample size and study design interact in a complex fashion to influence the
80 probability of detecting true effects on population density change. Here we extend that study by
81 holding sample size constant and investigating two sources of confounding effects. First, we
82 investigate the influence of background interactions arising from correlations between baseline
83 state and intervention (i.e., in our context, property loss and wildlife removal), which is
84 analogous to self-selection or treatment bias. This is a very common interaction in our subfield.

85 Second, we investigate the confounding effect of correlation between baseline property loss and
86 subsequent property loss in the absence of intervention (temporal autocorrelation). Third, we
87 extend [8, 11, 12] by measuring error rates in simulations of study designs that use Pearson
88 correlation coefficients when treatment effects vary in size and stochasticity. We use simple
89 simulations that expose the rates of Type I errors, Type II errors, and spurious correlations in
90 which the direction of the sign of correlation is reversed when compared to the true direction of
91 the cause and effect. We calculate FPRs and over-estimation bias.

92 Our approach applies generally to many or all fields that investigate systems characterized by the
93 baseline-intervention-outcome or state-stimulus-reaction causal relationships, including so-called
94 natural experiments. Our simulations model only three parameters and their interactions: (1) loss
95 of asset or resource prior to intervention, analogous to the baseline/state; (2) removal of wildlife,
96 shortly after time t , analogous to the intervention/stimulus; and (3) loss after intervention,
97 analogous to the outcome/reaction.

98 **Methods**

99 All variable names and definitions are presented in SM Table S1 along with definitions of study
100 designs and models.

101
102 To test the traditional wildlife control hypothesis (negative effect of treatment) and more recent
103 hypothesis (positive effect of treatment), we simulated losses of property such as the number of
104 domestic animals L_t lost at time t , followed by the intervention as people removed W wild
105 animals, and then we simulated losses in the next time step (L_{t+1}). To simulate crossover
106 designs, we added W at time $t + 1$ resulting in L_{t+2} . We modeled all W and L as independent,
107 normally distributed random, real numbers from zero to one inclusive, hereafter R . We varied
108 background interactions (B) to mimic potential conditions in the real world (see **Credibility of**
109 **models** below).

111 Estimating Type I and II error rates

112 Type I errors create false positives (we infer an effect of treatment when none exists) and Type II
113 errors lead to false negatives (we infer no treatment effect when one exists). We simulated
114 separately for each type of error. Separately with new iterations of simulations, we examined
115 extreme Type I error when the sign of correlation was reversed over the true sign of correlation.
116 In that simulation, we also examined extreme overestimation of treatment effects by $>2SD$ above
117 a positive mean treatment effect or $>2SD$ below a negative mean treatment effect.

118
119 In step one, we set $T = 0$ for no treatment effect ($W \times T$) and assigned $B = 0, -1.16, +1.16, -2.32,$
120 or $+2.32$. We combined different background interactions for Models 0-8 to estimate rates of
121 Type I errors (Table 1, Panels A–D). We set the coefficients empirically to yield an average
122 Pearson $r = 0.50$ ($n=1000$ replicates, 10 iterations) so there would be an equal space in either tail
123 for errors. We simulated 200 sets of 20 correlation coefficients with $n=50$ replicates each (400
124 iterations per scenario) for each of the 9 model permutations (3600 iterations per scenario-
125 model).

126
127 In step two, we repeated the same number of independent simulations as in step one. We
128 simulated cause-and-effect relationships between W and L_{t+1} (i.e., we set $T = \pm 0.58$, Table 1,
129 Panels E–H), to estimate rates of Type II errors (Table 1, Panels E–H).

131 For step three, we estimated false positive rates (FPR) following [8] as Type I error rate/[Type I
132 error rate + (1- Type II error rate)] using data from Table 1 to construct Table 2.

133
134 In step four, we produced five new independent simulations (400 iterations each) to investigate
135 variations of the Type I error in which the lack of a treatment effect changed from a constant $T =$
136 0 to a normally distributed random variable centered on zero but with more or less variability per
137 subject from -0.5 to $+0.5$, -1 to $+1$, -2 to $+2$, -4 to $+4$, and finally -8 to $+8$. Operationally, we
138 created that random T by subtracting two random numbers of equal magnitude from each other
139 for every replicate. This is analogous to a treatment effect that varies by subject (see Credibility
140 of models below). We estimated Type I error rates again as above. We modeled with a
141 generalized linear mixed model those error rates with four predictors (study design, variable
142 treatment effect for each replicate, background interactions from Models 3 and 4, and the
143 direction of the Type I error (i.e., whether a spurious significant result emerged for a positive or a
144 negative correlation).

145
146 In steps five and six, we explored the extreme Type II errors. We ran seven simulations
147 independent of those above (400 iterations each). For sign reversal, we counted the number of
148 correlation coefficients that had an opposite sign as the real correlation regardless of the
149 magnitude. In step 5, for extreme errors we repeated the procedure in steps 1-2 but counted the
150 number of treatment effect size estimates that exceeded the mean $+2SD$ for a positive treatment
151 effect or fell below the mean $-2SD$ for a negative treatment effect. For both steps 3 and 4,
152 temporal autocorrelation (B) varied from -2.32 to $+2.32$ independently of study design. We
153 estimated mean and standard deviations of error rates in both steps (Figs. 1 and 2).

154
155 In all steps, we chose deterministic and probabilistic scenarios in preference to empirical
156 domestic animal loss rates from the literature, because the latter would include unmeasured
157 background interactions and unreported treatment (e.g., poaching), which would undermine our
158 effort at measuring the odds of Type I and II errors.

159 Credibility of models

160 Background interactions simulate common situations in wildlife control. A positive correlation
161 between W and L_t (Models 1 and 2, Table S1) mimics a common background interaction in
162 which people kill more predators if losses were high in the past [19]. Probably uncommon, a
163 negative correlation between W and L_t mimics when people kill fewer predators after high
164 losses, e.g., when people and wildlife separate spatially after high losses [20, 21]. A positive
165 correlation between L_t and L_{t+1} (Models 3 and 4, Table S1) without intervention mimics a
166 common temporal autocorrelation, in which sites with high losses one year have high losses the
167 next year [22, 23]. Possibly less common, a negative temporal autocorrelation mimics cyclical
168 patterns of damage in non-sequential years. For example, when wild food availability influences
169 bear damage to crops and human foods, one may see a negative temporal autocorrelation of
170 losses from year to year [24, 25]. Or, if predators switch from domestic to wild prey selection
171 based on their relative scarcity or vulnerability varying over time, we can see prey switching
172 from season to season that might produce negative autocorrelations of losses in sequential time
173 steps [26-29].
174
175

176 These first four background interactions create univariate permutations. In the last four bivariate
177 permutations (Models 5–8, Table S1), we simulated both sets of interactions occurring
178 simultaneously in a two by two matrix of positive or negative interactions. For step four, when
179 we varied the treatment effect size in every replicate, we mimicked a situation in which the same
180 dose had variable effects on different replicates. For example, an individual predator may
181 respond differently than its neighbor or the composition of social groups may affect how the
182 survivors respond to removal of a group member, e.g., removing alpha individuals from a wolf
183 pack is expected to have different effects than removing subordinate adults or pups from a pack,
184 and even packs experiencing the same removal of dominant breeders might have different effects
185 depending on timing and availability of replacement breeders [30]. Hence, the same dose (W)
186 could have different treatment effect (T) depending on the idiosyncrasies of different replicates.
187 Similarly, some individual predators might be attracted or repelled by vacancies left by removals
188 of other predators [31]. Alternately, any of the individuals involved might respond differently to
189 lethal treatments. Theory provides five potential explanations for why the traditional hypothesis
190 may fail [31]. In brief, the wrong predators may be killed, e.g., [32]; the survivors may prey on
191 livestock that are more predictable than wild prey after the predators' social group has been
192 disrupted, e.g., pack hunting carnivores that rely on teamwork to hunt or reproduce successfully,
193 e.g., [33]; more immigrants may replace fewer residents that were killed, e.g., [34]; smaller-
194 bodied predator species at higher densities may refill the vacancies left by larger, scarcer
195 predator species that died, e.g., [35]; or humans and domestic animals may change their behavior
196 after lethal intervention. When we consider the entire set of actors, predators, humans, and
197 domestic animals, one can imagine interindividual differences in response to lethal interventions.
198 For example, some bold and tolerant individuals might explore wilder habitat after predator
199 removal while others might continue to avoid those areas [31]. In short, the same treatment of
200 different actors could result in diametrically opposed consequences even though the treatment
201 did have an effect on a subset of replicates. Despite different effects on different subjects, across
202 replicates, the general effect of treatment approximates zero so we estimated Type I error rates.

203 Analysis

204 We calculated Pearson's correlation coefficient r in JMP Pro V15.0.0 (SAS 2019). Pearson's r is
205 easily interpretable, dimensionless, and suitable for normally distributed, random variables [36].
206 With normally distributed response variables like L and change in L , Pearson's r is unbiased,
207 normal (Anderson-Darling test $A = 0.78$, $p = 0.05$ and $A = 0.37$, $p = 0.38$, respectively). We
208 calculated r in 20 batches of 50 replicates (analogous to independent sites or populations), a
209 larger sample size than most studies of wildlife control. We used the Pearson r standard critical
210 value of $|r| = 0.273$ (two-tailed test at $\alpha = 0.05$, $n = 50$ calculated from
211 [https://www.statisticssolutions.com/free-resources/directory-of-statistical-analyses/pearsons-
212 correlation-coefficient/table-of-critical-values-pearson-correlation/](https://www.statisticssolutions.com/free-resources/directory-of-statistical-analyses/pearsons-correlation-coefficient/table-of-critical-values-pearson-correlation/), accessed 5 July 2023) in 400
213 iterations of each combination of scenarios (Table S1) for a total of 108,000 independent
214 combinations. We calculated 400 correlations per simulation (108 scenarios in Table 1, 25
215 scenarios for the mixed model of Type I errors, and 35 scenarios for extreme Type II errors) for a
216 total of 67,200 Pearson r values including 50 independent replicates each. There were fewer
217 scenarios for randomized designs because the background interactions of L correlated with W
218 were eliminated by random assignment procedures (Table S1).

219
220
221 We involved neither animals nor human subjects in this research.

222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264

Results

False Positive Rates (FPR)

As predicted in [2], study designs differed noticeably in Type I and II error rates (Table 1) and therefore, in FPR (Table 2). As predicted by [8], FPRs exceeded Type I error rates based on p values in 93% (100/108) of our simulations (Table 2). None of the scenarios had FPR <1%. Therefore, we echo calls for lowering the statistical threshold for new discoveries [1].

The lowest FPR was 3.9% for rBACI when there were no background interactions (Table 2). In 8 scenarios, the FPR was 5.0% or less (4 scenarios with rBACI and 4 with crossover). Although rBACI had two of the lowest FPR (Table 2), it was outperformed by crossover when we introduced temporal autocorrelation in either direction, i.e., background interaction between B due to correlation between L_t and L_{t+1} . Indeed, crossover designs had a lower average FPR across 12 scenarios (6.1%, SD 1.4%) than RCT (6.4%, SD 1.0%) and rBACI (6.5%, SD 2.6%). Although these differences in FPR among randomized designs are small, the case for crossover design strengthened as we explain next.

We used a generalized linear mixed equation to model the interactions between confounding effects and study design on Type I error rates when treatment effects were centered on zero, but random in each replicate, i.e., no treatment effect in general (see Methods for examples of when this might arise). The mixed model revealed significant fixed effects only for study design (df=4, F=78, p<0.00001) and variable treatment effect for each replicate (df=1, F=31, p<0.0001); neither direction of error (df=1, F=0.2, p=0.62) nor the magnitude of temporal autocorrelation (df=6, F=1, p=0.44) were predictive of error. Also, study design and variable treatment effect for each replicate interacted significantly to predict the Type I error (df=4, F = 64, p<0.0001). Crossover performed best, because RCT and rBACI were somewhat vulnerable to randomly varying treatment effects (0.8% higher error rates), probably because the crossover design exposes each replicate to both control (treatment $T = 0$) and treatment (T varies randomly around zero) conditions. Because Type I error rates contribute to FPR directly, crossover design (platinum-standard) provided a stronger inference than the other study designs we tested [2].

Given FPR >1% seem risky to us, we recommend lowering the threshold for significance even when randomized designs are employed. Our results also corroborate prior cautions to measure and account for temporal autocorrelation [17]. We believe that temporal autocorrelation is a common condition in our field because of the widespread and frequent reports of 'hot spots' of damage by wild animals year after year [22, 37-40].

By comparison to the randomized study designs, we cannot recommend simple correlation or nBACI (bronze- and silver-standard, respectively) because their FPR ranged from 5.2-42% and 5.8-88%, respectively (Table 2). Negative temporal autocorrelation (Model 4) made these designs particularly vulnerable with FPR two to three times higher than for positive temporal autocorrelation. The highest FPR arose in Models 5-8 (Table 2). Although nBACI was somewhat resistant to Models 5 and 8 when the background interactions were strong (2.32), nBACI failed in most cases, including several with only one background interaction (Table 2). Although simple correlations yielded consistent FPR of 5-12.5% when we introduced only one background interaction, their FPR rose above 20% whenever we included two background interactions.

Although one might be tempted to look at a few low Type I error rates in Table 1 for simple correlation and nBACI, and declare these study designs viable in many circumstances, the FPR

265 in Table 2 warn against such confidence. Also, with FPR for simple correlation averaging 16%
266 (SD 12%) and nBACI averaging 29% (SD 25%), in the absence of good evidence about
267 background interactions, one should not credit these study designs. Indeed, in many experimental
268 situations, particularly under field conditions surrounding wildlife control, researchers will have
269 little or no evidence to dismiss background interactions. Even when such evidence for
270 background interactions is robust and well-accounted in the analyses, few researchers in our field
271 can build a sample size of 50 on which our simulations depend. Therefore, FPR values in Table 2
272 are likely under-estimates of what others will encounter with smaller samples, imperfect
273 randomization, variable treatment effect for each replicate, deviations from the assumptions of
274 Pearson correlations, and measurement error [8].

275 *Severe Type II errors: overestimation and sign reversal*

276 Some of the simulated Type II error rates were very high (Table 1), which by itself may not raise
277 concern because Type II error conservatively leads us to infer no effect when one exists in
278 reality. However, reporting the opposite sign of correlation than the real direction of correlation
279 when a treatment is effective would be an extreme form of Type II error that merits concern (Fig.
280 1). Also, when we overestimate the real effect substantially (e.g., >2SD above a positive mean or
281 below a negative mean), exaggerated claims about treatment effectiveness can mislead users,
282 payers, and distributors of that treatment (Fig. 2). As temporal autocorrelation increased, the rate
283 of sign reversal increased and simple correlation was more strongly affected than nBACI (Fig.
284 1). The converse was true for overestimation error, which declined among the non-randomized
285 study designs. Simple correlation was less prone to these errors than nBACI (Fig. 2).
286

287 Compared to randomized designs, the rates of sign reversal for simple correlation and nBACI
288 were higher (8% and 0.8% respectively; only simple correlation differed significantly from every
289 other design, each t-test pairwise comparison $p < 0.0001$) than randomized designs (RCT – 0.09%,
290 rBACI – 2%, crossover – 0.08%, which did not differ among randomized designs, Welch test
291 unequal variances, F ratio = 2, $p = 0.15$).

292 Similarly, non-randomized designs had higher rates of overestimating treatment effect sizes (8%
293 for simple correlation and 31% for nBACI), which differed significantly from randomized
294 designs ($p < 0.0001$ for each pairwise comparison with nBACI, $p < 0.009$ for pairwise comparisons
295 of simple correlation to each randomized design). Also, randomized study designs were
296 statistically different in rates of overestimation error (RCT – 0.2%, rBACI – 1%, crossover – 3%,
297 F ratio = 31, $p < 0.0001$).

298 In sum, our predictions of the relative strength of inference among study designs were only
299 partly supported [2]. The predicted difference between simple correlation (bronze-standard) and
300 nBACI (silver-standard) held for sign reversal (Fig. 1), but not for overestimation bias (Fig. 2) or
301 most FPR (Table 2). Similarly, the so-called gold+ of rBACI compared to gold-standard RCT
302 did not play out as we predicted [2]. Yet, our predictions about crossover design (platinum-
303 standard) producing stronger inference than RCT and rBACI (gold-standards) were supported.
304 Therefore, we revised our first hypotheses [2] by producing a schematic graph of relative
305 strengths of inference estimated for five study designs (Fig. 3).

306 **Discussion**

307 Some public authorities may not test treatments with randomized, controlled experiments
308 because they perceive intervening as infeasible or impractical, perhaps in part because they

309 believe the treatments will be popular and the placebo controls will be unpopular, e.g., [41].
310 Therefore, authorities may prefer to intervene in ways they consider less controversial, such as
311 treating all subjects or serving the loudest complainants [[5], see webpanel 1]. Such steps that
312 lead to non-randomized study designs risk backfiring or wasting time and resources.

313
314 When subjects are self-selected (self-selection bias), vulnerable subjects receive higher doses
315 (treatment bias), or baseline conditions affect outcomes and not just treatments (e.g., temporal
316 autocorrelation), we can expect high false positive rates (FPR, Table 2), especially for non-
317 random before-and-after comparisons of interventions (nBACI). When background interactions
318 are strong, FPR rise sharply (Table 2). When both sets of background interactions coincide, we
319 estimated that wrong conclusions would be drawn in 18–42% of simple correlation studies and
320 even more variably in 8–88% of nBACI (Table 2). Also, when temporal autocorrelation is
321 present, the results of non-randomized study designs will produce additional errors even if the
322 study is designed to minimize false positives. Non-randomized designs pose a considerable risk
323 of the reversal of the sign of correlation, which can substantially mislead researchers and
324 practitioners about the treatment effect (Fig. 1). If sign reversal does not occur, overestimation of
325 treatment effects is also possible (Fig. 2). These compounding errors associated with non-
326 randomized study designs can be visualized as a hierarchy of study designs (Fig. 3).

327 Overall, the compounding errors weigh heavily against non-randomized designs (Fig. 3). Unlike
328 randomized designs, non-randomized designs produce errors asymmetrical with regard to
329 positive or negative background interactions (Figs. 1, 2). Namely, positive temporal
330 autocorrelations produced more sign reversal errors and fewer overestimation errors in non-
331 randomized designs than did negative temporal autocorrelations. That asymmetry would tend to
332 confuse the direction of the treatment effect more often when outcomes correlate positively to
333 baseline conditions (Fig. 1); that situation is common in our subfield where hot spots of wildlife
334 damage recur annually (*SM*).

335 Regrettably, predator control has been dominated by unreliable, non-randomized studies. Hence,
336 predictably, there is no scientific consensus about the effects of predator control on subsequent
337 domestic animal losses, particularly in case of lethal treatments [3, 14, 15]. For example, non-
338 randomized study designs have produced equivocal results for lethal control including recurrent
339 findings of counter-productive increases in domestic animal losses following killing gray wolves
340 [42, 43], bears (*Ursus* spp.) [25, 44, 45] and cougars (*Puma concolor*) [46, 47]. Theory provides
341 five potential explanations for why the traditional hypothesis may fail, cf. [31] and described
342 with references in Methods. In brief, the wrong predators may be killed; the survivors' behaviors
343 may change if they relied on group-mates that were killed; immigrants of the same species or
344 smaller-bodied predatory species may refill in greater numbers the vacancies left after killing; or
345 survivors of any species may change behavior after predators are removed.

346 Even well-financed RCT across broad areas may be hard to interpret, e.g., U.K.-funded RCT of
347 badger (*Meles meles*) killing to prevent bovine tuberculosis documented variable effects of this
348 intervention that can be difficult to detect [48-53]. Even methods considered politically
349 unpalatable but highly effective, such as poisoning red foxes (*Vulpes vulpes*) in Australia to
350 protect sheep, when tested with RCT prove highly variable in effect [54]. The latter research
351 team concluded from an RCT that poisoning foxes wasted much effort and was ineffective
352 because it produced very slight decreases in lamb mortality. Despite these doubts, lethal methods
353 are rarely subjected to RCT. Most randomized studies of predator control have been conducted

354 on non-lethal methods to prevent predators from damaging property [41, 55, 56]. An analogy
355 would be to ignore experiments on handgun control [57] while subjecting pepper spray to robust
356 RCT. Moreover, in the absence of scientific consensus the historical intervention of killing
357 predators continues unabated despite years of criticism [5, 48].

358 The resilience of lethal treatments in policy circles may reflect a perceptual bias of “cherry
359 picking” arising from the adoption of a few effective cases and the dismissal of more numerous
360 ineffective cases [33, 42, 43, 58]. Our mixed models show that treatments that help some
361 replicates and harm others will raise FPR with worrying frequency in non-randomized studies. In
362 addition, animal killing may fall into another perceptual bias because either humans cannot
363 recognize individual animals, some of which are culprits and some of which are not [32, 33], or
364 some persons may claim a lethal treatment has succeeded because the death of a competitor
365 might have been their primary goal regardless of its culpability.

366 If a non-randomized design is analyzed in spite of our cautions above, researchers should
367 account for potential self-selection bias, treatment bias, and temporal autocorrelation. For
368 example, lethal wildlife control studies should measure (a) killing and property losses before that
369 killing occurred, and (b) property losses from year to year in the absence of intervention [17, 43].
370 The absence of intervention includes unplanned or unregulated interventions by the people
371 participating or using the same areas. This is a very difficult hurdle to overcome without strict
372 control of participant actions because predator killing can still be present as an illicit behavior
373 and hushed up [59-61]. Therefore, we suggest randomized designs in smaller, well-controlled
374 sites are likely to be more feasible than strict control over potentially confounding variables
375 across entire landscapes. Even for randomized designs, we counsel care because FPR does not
376 diminish to zero. To lower the risk of FPR, we recommend the platinum-standard crossover
377 design RCT (all subjects receive both treatment and placebo in random order), lowering the
378 significance threshold [1], and other safeguards against bias [2].

379 A common argument for drawing inference from non-randomized studies has been that experts
380 can infer accurately despite confounding variables [17]. For example, expert-based adaptive
381 managers claim they can intervene, learn, and revise without exacerbating the problems at hand
382 and without exposing hypotheses to experimental test [62, 63]. That argument depends on
383 learning correctly. The counter-argument is that biased designs and lower standards hinder
384 learning with false information and can produce inferences diametrically opposed to the actual
385 effect of interventions [6, 64]. Our results of sign reversal in treatment effects support that
386 concern. Therefore, prioritizing randomized designs for urgent and important policy decisions
387 may avoid the age-old problem that haste makes waste. The reasoning here provides a guide to
388 donors, regulators, and the public to anticipate situations in which RCT becomes a prerequisite
389 for reliable inference and sound policy.

390
391 **Acknowledgments:** We thank RJ Treves for statistical advice.

392 **Data and materials availability:** For scripts and a full spreadsheet with 1000 rows of data
393 for a single iteration of each simulation, see
394 [https://faculty.nelson.wisc.edu/treves/data_archives/Simulate_study_designs_scripts_data_ar](https://faculty.nelson.wisc.edu/treves/data_archives/Simulate_study_designs_scripts_data_archive.zip)
395 [chive.zip](https://faculty.nelson.wisc.edu/treves/data_archives/Simulate_study_designs_scripts_data_archive.zip), accessed 21 October 2023.

396

References

- [1] D. Benjamin *et al.*, "Redefine statistical significance," *Nature Human Behaviour*, vol. 2, pp. 6–10, 2018.
- [2] A. Treves, M. Krofel, O. Ohrens, and L. M. Van Eeden, "Predator control needs a standard of unbiased randomized experiments with cross-over design," *Frontiers in Ecology and Evolution*, vol. 7 pp. 402-413, 2019, doi: 10.3389/fevo.2019.00462.
- [3] I. Khorozyan, "Defining practical and robust study designs for interventions targeted at terrestrial mammalian predators," *Conserv. Biol.*, vol. 36, p. e13805, 2022, doi: 10.1111/cobi.13805.
- [4] I. Khorozyan, "Dealing with false positive risk as an indicator of misperceived effectiveness of conservation interventions," *PLoS One*, vol. 16, no. 5, p. e0255784, 2021. [Online]. Available: <https://doi.org/10.1371/journal.pone.0255784>.
- [5] A. Treves, M. Krofel, and J. McManus, "Predator control should not be a shot in the dark," *Front. Ecol. Environ.*, vol. 14, pp. 380-388, 2016.
- [6] J. P. Ioannidis, "Why most published research findings are false," *PLOS Medicine*, vol. 2, no. 8, p. e124, 2005.
- [7] M. Baker and K. Brandon, "Is there a reproducibility crisis? A Nature survey lifts the lid on how researchers view the 'crisis' rocking science and what they think will help," *Nature*, vol. 533, pp. 452-454, 2016. [Online]. Available: <https://www.nature.com/articles/533452a>.
- [8] D. Colquhoun, "An investigation of the false discovery rate and the misinterpretation of p-values," *Royal Society Open Science*, vol. 1, p. 140216, 2014. [Online]. Available: <http://dx.doi.org/10.1098/rsos.140216>.
- [9] N. Oreskes, *Why Trust Science?* Princeton, NJ: Princeton University Press, 2019.
- [10] J. R. Platt, "Strong inference," *Science*, vol. 146, pp. 347–353, 1964.
- [11] A. P. Christie *et al.*, "Quantifying and addressing the prevalence and bias of study designs in the environmental and social sciences," *Nature Communications*, vol. 11, p. 6377, 2020. [Online]. Available: <https://doi.org/10.1038/s41467-020-20142-y> | www.nature.com/naturecommunications.
- [12] A. P. Christie, T. Amano, P. A. Martin, G. E. Shackelford, B. I. Simmons, and W. J. Sutherland, "Simple study designs in ecology produce inaccurate estimates of biodiversity responses," *J. Appl. Ecol.*, vol. 56, pp. 2742–2754, 2019, doi: 10.1111/1365-2664.13499.
- [13] T. J. Clark and M. Hebblewhite, "Predator control may not increase ungulate populations in the future: A formal meta-analysis," *J. Appl. Ecol.*, vol. 58, no. 4, pp. 812-824, 2021, doi: 10.1111/1365-2664.13810.
- [14] A. Eklund, J. V. López-Bao, M. Tourani, G. Chapron, and J. Frank, "Limited evidence on the effectiveness of interventions to reduce livestock predation by large carnivores," *Scientific Reports*, vol. 7, pp. 2097 | DOI:10.1038/s41598-017-02323-w, 2017.
- [15] L. M. van Eeden *et al.*, "Carnivore conservation needs evidence-based livestock protection," *PLoS Biol.*, vol. 16, no. 9, p. e2005577, 2018, doi: e2005577.
- [16] A. J. Underwood, "Beyond BACI: The detection of environmental impacts on populations in the real, but variable, world," *J. Exp. Mar. Biol. Ecol.*, vol. 161, pp. 145-178, 1992.

- 442 [17] P. A. Murtaugh, "On rejection rates of paired intervention analysis," *Ecology*, vol. 83, pp.
443 1752–1761, 2002.
- 444 [18] D. Colquhoun, "The reproducibility of research and the misinterpretation of p-values,"
445 *Royal Society Open Science*, vol. 4, p. 171085, 2017. [Online]. Available:
446 <http://dx.doi.org/10.1098/rsos.171085>.
- 447 [19] M. R. Conover, "Effect of hunting and trapping on wildlife damage " *Wildl. Soc. Bull.*,
448 vol. 29, no. 2, pp. 521-532, 2001.
- 449 [20] J. Knight, *Waiting for Wolves in Japan*. Oxford: Oxford University Press, 2003.
- 450 [21] L. Naughton-Treves, "Whose animals? A history of property rights to wildlife in Toro,
451 western Uganda," *Land Degradation and Development*, vol. 10, pp. 311-328, 1999.
- 452 [22] A. Treves, K. A. Martin, A. P. Wydeven, and J. E. Wiedenhoef, "Forecasting
453 environmental hazards and the application of risk maps to predator attacks on livestock,"
454 *Bioscience*, vol. 61, pp. 451-458, 2011.
- 455 [23] J. R. B. Miller, "Mapping attack hotspots to mitigate human–carnivore conflict:
456 approaches and applications of spatial predation risk modeling," *Biodivers. Conserv.*, vol.
457 24, no. 12, pp. 2887-2911, 2015.
- 458 [24] D. L. Garshelis, "Nuisance bear activity and management in Minnesota," in *Bear - People*
459 *Conflicts - Proceedings of a Symposium on Management Strategies*, M. Bromley Ed.
460 Yellowknife, Canada: Northwest Territories Department of Renewable Resources, 1989,
461 pp. 169-180.
- 462 [25] J. M. Northrup, E. J. Howe, J. Inglis, E. Newton, M. E. Obbard, B. Pond, and D. Potter,
463 "Experimental test of the efficacy of hunting for controlling human–wildlife conflict," *J.*
464 *Wildl. Manage.*, vol. 87, no. 3, p. e22363, 2022, doi: 10.1002/jwmg.22363.
- 465 [26] A. Janeiro-Otero, T. M. Newsome, L. M. Van Eeden, W. J. Ripple, and C. F. Dormann,
466 "Grey wolf (*Canis lupus*) predation on livestock in relation to prey availability," *Biol.*
467 *Conserv.*, vol. 243, 2020, doi: 10.1016/j.biocon.2020.108433.
- 468 [27] I. Khorozyan, A. Ghoddousi, M. Soofi, and M. Waltert, "Big cats kill more livestock
469 when wild prey reaches a minimum threshold," *Biol. Conserv.*, vol. 192, pp. 268–275,
470 2015.
- 471 [28] I. Laporte, T. B. Muhly, J. A. Pitt, M. Alexander, and M. Musiani, "Effects of wolves on
472 elk and cattle behaviors: implications for livestock production and wolf conservation,"
473 *PLoS One*, vol. 5, no. 8, p. e11954, 2010, doi: 10.1371/journal.pone.0011954.
- 474 [29] J. Odden, I. Herfindal, J. D. C. Linnell, and R. Andersen, "Vulnerability of domestic
475 sheep to lynx depredation in relation to roe deer density " *J. Wildl. Manage.*, vol. 72, no.
476 1, pp. 276-282, 2008.
- 477 [30] S. M. Brainerd *et al.*, "The effects of breeder loss on wolves," *J. Wildl. Manage.*, vol. 72,
478 no. 1, pp. 89-98, 2008.
- 479 [31] L. Elbroch and A. Treves, "Why might removing carnivores maintain or increase risks
480 for domestic animals? ," *Biol. Conserv.*, vol. 283, p. 110106, 2023, doi:
481 10.1016/j.biocon.2023.110106.
- 482 [32] L. Plumer, T. n. Talvi, P. Männil, and U. Saarma, "Assessing the roles of wolves and
483 dogs in livestock predation and suggestions for mitigating human-wildlife conflict and
484 conservation of wolves," *Conserv. Genet.*, vol. 19, pp. 665–672, 2018. [Online].
485 Available: <https://doi.org/10.1007/s10592-017-1045-4>

- 486 [33] F. F. Knowlton, E. M. Gese, and M. M. Jaeger, "Coyote depredation control: An interface
487 between biology and management," *Journal of Range Management*, vol. 52, pp. 398-
488 412., 1999.
- 489 [34] H. S. Cooley, R. B. Wielgus, H. S. Robinson, G. M. Koehler, and B. T. Maletzke, "Does
490 hunting regulate cougar populations? A test of the compensatory mortality hypothesis,"
491 *Ecology*, vol. 90, pp. 2913-2921, 2009.
- 492 [35] K. R. Crooks and M. E. Soulé, "Mesopredator release and avifaunal extinctions in a
493 fragmented system," *Nature*, vol. 400, pp. 563-566, 1999.
- 494 [36] Open Science Collaboration, "Estimating the Reproducibility of Psychological Science,"
495 ed. <https://osf.io/ezcuj/>, 2015, p. <https://osf.io/ezcuj/>.
- 496 [37] L. Naughton-Treves, "Predicting patterns of crop damage by wildlife around Kibale
497 National Park, Uganda," *Conserv. Biol.*, vol. 12, pp. 156-168, 1998.
- 498 [38] K. K. Karanth, A. M. Gopalaswamy, P. K. Prasad, and S. Dasgupta, "Patterns of human-
499 wildlife conflicts and compensation: insights from Western Ghats protected areas," *Biol.*
500 *Conserv.*, vol. 166, pp. 175-185, 2013.
- 501 [39] J. R. B. Miller, Y. V. Jhala, J. Jena, and O. J. Schmitz, "Landscape-scale accessibility of
502 livestock to tigers: implications of spatial grain for modeling predation risk to mitigate
503 human–carnivore conflict," *Ecology and Evolution*, vol. 5, no. 6, pp. 1354-1367, 2015.
- 504 [40] A. Treves and M. F. Rabenhorst, "Risk map for wolf threats to livestock still predictive 5
505 years after construction," *PLoS One*, vol. 12, no. 6, p. e0180043, 2017, doi:
506 <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0180043>.
- 507 [41] O. Ohrens, C. Bonacic, and A. Treves, "Non-lethal defense of livestock against predators:
508 Flashing lights deter puma attacks in Chile.," *Front. Ecol. Environ.*, vol. 17, no. 1, pp. 32-
509 38, 2019, doi: 10.1002/fee.1952.
- 510 [42] O. Grente, "Présentation des objectifs et de la méthodologie de la thèse sur l'efficacité des
511 tirs de loup et la gestion adaptative du loup, menée conjointement par l'ONCFS et le
512 CEFE," OFB, Direction de la recherche et de l'appui scientifique, Unité Prédateurs,
513 animaux déprédateurs et exotiques, Équipe Loup et Lynx, Gières, France, 2021.
- 514 [43] F. J. Santiago-Avila, A. M. Cornman, and A. Treves, "Killing wolves to prevent
515 predation on livestock may protect one farm but harm neighbors," *PLoS One*, vol. 13, no.
516 1, p. e0189729 2018, doi: /10.1371/journal.pone.0189729.
- 517 [44] I. Khorozyan and M. Waltert, "Variation and conservation implications of the
518 effectiveness of anti-bear interventions," *Scientific Reports*, vol. 10, p. 15341, 2020, doi:
519 10.1098/rsos.190826.
- 520 [45] D. L. Garshelis, K. V. Noyce, and V. St-Louis, "Population reduction by hunting helps
521 control human-wildlife conflicts for a species that is a conservation success story," *PLoS*
522 *One*, vol. 15, no. 8, p. e0237274, 2020, doi: 10.1371/journal.pone.0237274.
- 523 [46] J. W. Laundré and C. Papouchis, "The elephant in the room: What can we learn from
524 California regarding the use of sport hunting of pumas (*Puma concolor*) as a management
525 tool?," *PLoS One*, vol. 15, no. 2, p. e0224638, 2020, doi:
526 <https://doi.org/10.1371/journal.pone.0224638>.
- 527 [47] K. Peebles, R. B. Wielgus, B. T. Maletzke, and M. E. Swanson, "effects of remedial sport
528 hunting on cougar complaints and livestock depredations.," *PLoS One*, vol. 8, no. 11, p.
529 e79713, 2013.

- 530 [48] J. Bielby, F. Vial, R. Woodroffe, and C. A. Donnelly, "Localised badger culling increases
531 risk of herd breakdown on nearby, not focal, land," *PLoS One*, vol. 11, no. 10, p.
532 e0164618, 2016, doi: 10.1371/journal.pone.0164618.
- 533 [49] C. Donnelly and R. Woodroffe, "Reduce uncertainty in UK badger culling," *Nature*, vol.
534 485, p. 582, 2012.
- 535 [50] C. Donnelly *et al.*, "Positive and negative effects of widespread badger culling on
536 tuberculosis in cattle," *Nature*, vol. 439, pp. 843-846, 2006.
- 537 [51] C. Donnelly *et al.*, "Impacts of widespread badger culling on cattle tuberculosis:
538 concluding analyses from a large-scale field trial.," *International Journal of Infectious*
539 *Diseases*, vol. 11, pp. 300-308, 2007.
- 540 [52] H. C. J. Godfray *et al.*, "A restatement of the natural science evidence base relevant to the
541 control of bovine tuberculosis in Great Britain," *Proceedings of the Royal Society B*, vol.
542 280, no. 1768, p. 20131634, 2013.
- 543 [53] F. Vial and C. Donnelly, "Localized reactive badger culling increases risk of bovine
544 tuberculosis in nearby cattle herds," *Biol. Lett.*, vol. 8, pp. 50-53, 2012.
- 545 [54] C. Greentree, G. Saunders, L. McLeod, and J. Hone, "Lamb predation and fox control in
546 south-eastern Australia," *J. Appl. Ecol.*, vol. 37, pp. 935-943, 2000.
- 547 [55] C. G. Radford, J. W. McNutt, T. Rogers, B. Maslen, and N. R. Jordan, "Artificial
548 eyespots on cattle reduce predation by large carnivores," *Communications Biology*
549 *Nature*, vol. 3, p. 430, 2020. [Online]. Available: [https://doi.org/10.1038/s42003-020-](https://doi.org/10.1038/s42003-020-01156-0)
550 [01156-0](https://doi.org/10.1038/s42003-020-01156-0) | www.nature.com/com.
- 551 [56] P. Fersterer, D. Nolte, G. Ziegltrum, and H. Gossow, "Effect of feeding stations on the
552 home ranges of American black bears in western Washington," *Ursus*, vol. 12, pp. 51-53,
553 2001.
- 554 [57] A. I. Leshner and V. J. Dzau, "Good gun policy needs research," *Science*, vol. 359, no.
555 6381, p. 1195, 2018.
- 556 [58] M. Krofel, R. Černe, and K. Jerina, "Effectiveness of wolf (*Canis lupus*) culling as a
557 measure to reduce livestock depredations," *Acta Silvae et Ligni*, vol. 95, pp. 11-22, 2011.
- 558 [59] O. Liberg, G. Chapron, P. Wabakken, H. C. Pedersen, N. T. Hobbs, and H. Sand, "Shoot,
559 shovel and shut up: cryptic poaching slows restoration of a large carnivore in Europe,"
560 *Proceedings of the Royal Society of London Series B*, vol. 270, pp. 91-98, 2012.
- 561 [60] G. Chapron and A. Treves, "Correction to 'Blood does not buy goodwill: allowing
562 culling increases poaching of a large carnivore'," *Proceedings of the Royal Society B*,
563 vol. 283, no. 1845, p. 20162577, 2016.
- 564 [61] F. J. Santiago-Ávila and A. Treves, "Poaching of protected wolves fluctuated seasonally
565 and with non-wolf hunting," *Scientific Reports*, vol. 12, p. e1738 2022, doi:
566 10.1038/s41598-022-05679-w.
- 567 [62] J. Hone, V. A. Drake, and C. J. Krebs, "the effort–outcomes relationship in applied
568 ecology: evaluation and implications," *Bioscience*, vol. 67, pp. 845–852, 2017.
- 569 [63] N. Salafsky, R. Margoluis, K. Redford, and J. G. Robinson, "Improving the practice of
570 conservation: a conceptual framework and research agenda for conservation science,"
571 *Conserv. Biol.*, vol. 16, no. 6, pp. 1469-1479, 2002.
- 572 [64] J. González-González, E. Dorsey-Treviño, N. Alvarez-Villalobos, F. Barrera-Flores, A.
573 Díaz González-Colmenero, C. Quintanilla-Sánchez, and e. al., "Trustworthiness of
574 randomized trials in endocrinology—A systematic survey," *PLoS One*, vol. 14, no. 2, p.
575 e0212360, 2019, doi: <https://doi.org/10.1371/journal.pone.0212360>.

576 accessed 19 August 2023.

577

578 **Author contributions:**

579 Conceptualization: AT

580 Methodology: AT, IK

581 Investigation: AT, IK

582 Visualization: AT

583 Funding acquisition: AT

584 Project administration: AT

585 Supervision: AT

586 Writing – original draft: AT

587 Writing – review & editing: AT, IK

588

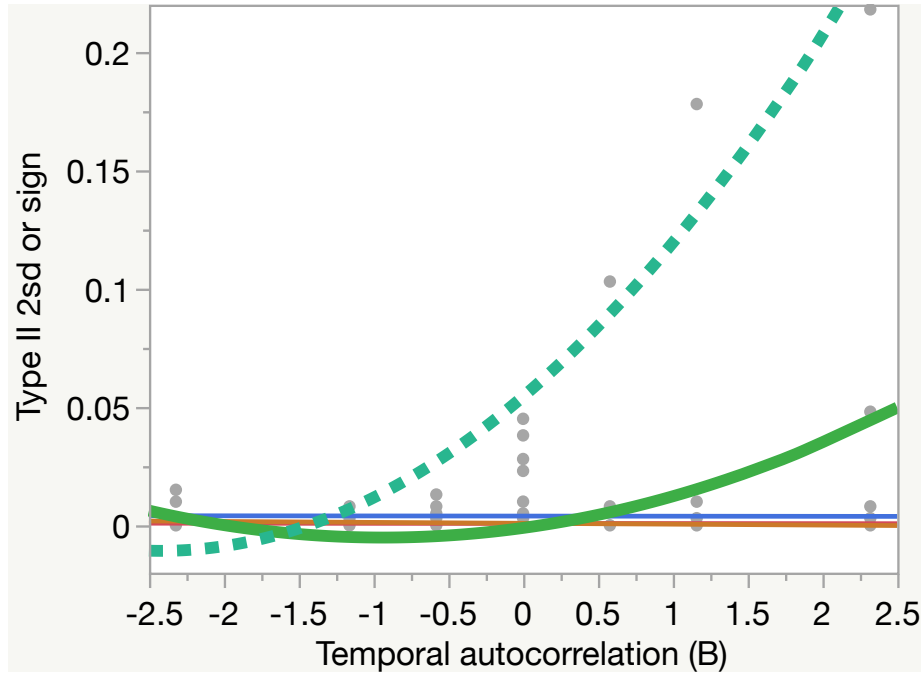
589 **Additional Information**

590 **Funding:** AT acknowledges the receipt of a fellowship from the OECD Co-operative
591 Research Programme: Sustainable Agricultural and Food Systems in 2022.

592 **Competing interests:** The authors declare no competing interest but readers can judge for
593 themselves by accessing a full statement of AT’s potentially competing interests at
594 http://faculty.nelson.wisc.edu/treves/archive_BAS/funding.pdf, accessed 13 August 2023,
595 with a complete CV at
596 http://faculty.nelson.wisc.edu/treves/archive_BAS/Treves_vita_latest.pdf, accessed 13
597 August 2023.

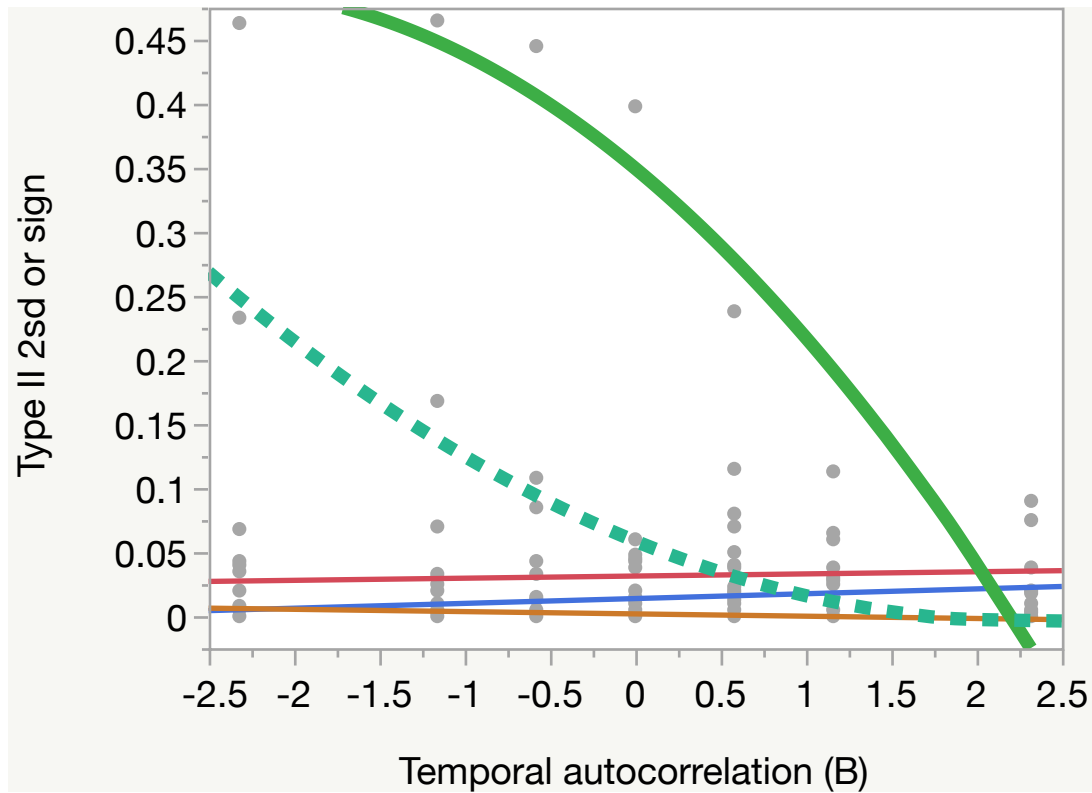
598

599



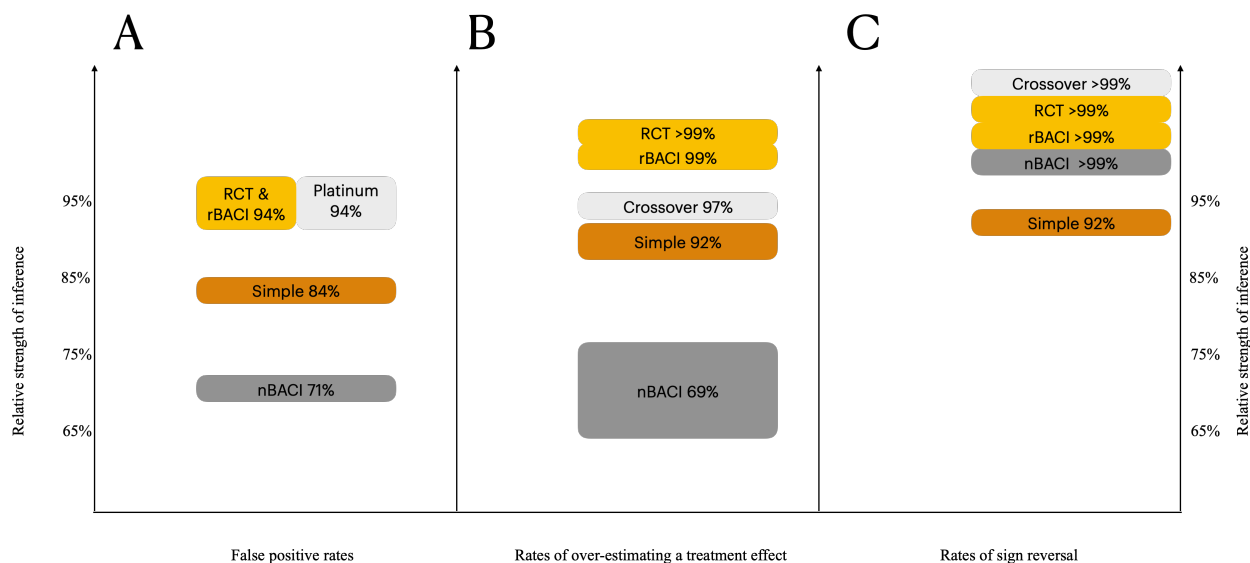
600
601
602
603
604
605
606
607
608
609

Figure 1. Severe Type II error resulting in reversal of the sign of correlation, in relation to temporal autocorrelation between L_t and L_{t+1} (B). We present a curve fit by second-order ordinary least squares regression for visualization purposes only for each study design (dashed green = simple correlation, solid, thick green = nBACI, gold = RCT, purple = rBACI, red = crossover). The x-axis presents varying levels of temporal autocorrelation from Models 3 and 4 (Table S1). The y-axis presents the frequency of reversal of the true sign of correlation to the opposite sign estimated from 400 iterations of each combination of study design and value of B.



610
 611
 612
 613
 614
 615
 616
 617
 618
 619
 620

Figure 2. Overestimation of treatment effect in relation to temporal autocorrelation between L_t and L_{t+1} (B). We present a curve fit by second-order ordinary least squares regression for visualization purposes only for each study design (dashed green = simple correlation, solid, thick green = nBACI, gold = RCT, purple = rBACI, red = crossover). The x-axis presents varying levels of temporal autocorrelation from Models 3 and 4 (Table S1). The y-axis presents the frequency of overestimation of treatment effect >2 SD below and above the mean estimated from 400 iterations per data point. Simulations are the same as in Fig. 1.



621
622
623
624
625
626
627
628
629

Figure 3. Relative strength of inference (100% - mean error rate) for crossover (platinum), RCT (gold), rBACI (gold), nBACI (silver), and simple correlation (bronze). The height of polygons is scaled to the 95% CI within each panel: (A) False positive rates, (B) Rates of overestimating as treatment effect, and (C) Rate of sign reversal. Side-by-side bars (e.g., panel A platinum and gold standards indicate identical mean and 95% CI but stacked bars indicate means were not identical (e.g., Panel C).

630
631
632

Table 1. Error rates estimated with and without background interactions: (A) B=1.16, (B) B=2.32, (C, D) T = 0 for Type I error; (E–H) are set to T = 0.58 x W.

	Simple correlation	nBACI	RCT	rBACI	Crossover design	Simple correlation	nBACI	RCT	rBACI	Crossover design
		I	†	†	†			†	I †	†
Models	A. Background interactions 1.16					B. Background interactions 2.32				
	C. Type I errors					D. Type I errors				
0	0.053	0.053	0.055	0.040	0.068	0.053	0.053	0.055	0.040	0.068
1	0.055	0.515				0.068	0.745			
2	0.068	0.548				0.060	0.718			
3	0.050	0.050	0.075	0.060	0.043	0.038	0.045	0.053	0.048	0.035
4	0.045	0.075	0.063	0.083	0.050	0.058	0.070	0.053	0.055	0.060
5	0.225	0.145				0.405	0.105			

6	0.223	0.595				0.435	0.743			
7	0.240	0.615				0.448	0.760			
8	0.218	0.158				0.455	0.088			
Models	E. Type II errors, positive treatment					F. Type II errors, positive treatment				
0	0.025	0.185	0.00 0	0.023	0.193	0.025	0.185	0.00 0	0.02 3	0.193
1	0.005	0.385				0.000	0.010			
2	0.000	0.000				0.000	0.000			
3	0.245	0.195	0.02 0	0.020	0.190	0.515	0.475	0.35 0	0.23 8	0.203
4	0.190	0.410	0.03 0	0.238	0.200	0.595	0.710	0.34 0	0.50 5	0.165
5	0.005	0.135				0.000	0.000			
6	0.210	0.915				0.365	0.890			
7	0.000	0.000				0.000	0.000			
8	0.190	0.000				0.350	0.015			
Models	G. Type II errors, negative treatment					H. Type II errors, negative treatment				
0	0.030	0.195	0.00 0	0.015	0.188	0.030	0.195	0.00 0	0.01 5	0.188
1	0.000	0.000				0.000	0.000			
2	0.000	0.440				0.000	0.005			
3	0.255	0.220	0.03 3	0.030	0.185	0.640	0.500	0.27 5	0.17 3	0.205
4	0.205	0.435	0.01 8	0.215	0.208	0.575	0.715	0.33 8	0.50 5	0.245
5	0.180	0.005				0.385	0.025			
6	0.005	0.005				0.005	0.005			
7	0.180	0.890				0.370	0.895			
8	0.000	0.075				0.000	0.000			

† Blank cells reflect that random assignment eliminates a correlation between W and L t.

633
634
635

636
637
638

Table 2. False positive rates (FPR) estimated from Type I and II error rates in Table 1 with background interactions: (A) B = 1.16 (B) B = 2.32, (C) positive treatment effect, (D) negative treatment effect.

Models	False positive rates (FPR) %									
	Simple correlation	nBACI	RCT †	rBACI †	Crossover design †	Simple correlation	nBACI	RCT †	rBACI †	Crossover design †
	A. Background interactions 1.16					B. Background interactions 2.32				
	C. Positive treatment ††									
0	5.2	6.1	5.2	3.9	7.8	5.2	6.1	5.2	3.9	7.8
1	5.5	45.6				6.4	42.9			
2	6.4	35.4				5.7	41.8			
3	6.2	5.8	7.1	5.8	5.0	7.3	7.9	7.5	5.9	4.2
4	5.3	11.3	6.1	9.8	5.9	12.5	19.4	7.4	10.0	6.7
5	18.4	14.4				28.8	9.5			
6	22.0	87.5				40.7	87.1			
7	19.4	38.1				30.9	43.2			
8	21.2	13.6				41.2	8.2			
	D. Negative treatment effect ††									
0	5.2	6.2	5.2	3.9	7.7	5.2	6.2	5.2	3.9	7.7
1	5.2	34.0				6.6	42.7			
2	6.4	49.5				5.7	41.9			
3	6.3	6.0	7.2	5.8	5.0	9.5	8.3	6.8	5.5	4.2
4	5.4	11.7	6.0	9.6	5.9	12.0	19.7	7.4	10.0	7.4
5	21.5	12.7				39.7	9.7			
6	18.3	37.4				30.4	42.8			
7	22.6	84.8				41.6	87.9			
8	17.9	14.6				31.3	8.1			
Minimum	5.2	5.8	5.2	3.9	5.0	4.35.2	6.1	5.2	3.9	4.2
95% CI of mean	9–15	17–41	5–7	4–8	5–7	13–27	18–42	6–8	5–9	5–7

639
640
641

† Blank cells reflect that random assignment eliminates a correlation between W and L t.
†† Simulated positive treatments may produce different FPR than negative treatments.