This version:      2024 02 05
Newest version:  https://urisohn.com/44

# Stimulus Sampling Reimagined: Designing Experiments with Mix-and-Match, Analyzing Results with Stimulus Plots

Uri Simonsohn
ESADE Business School
urisohn@gmail.com

Andres Montealegre
Cornell University
am2849@cornell.edu

Ioannis Evangelidis
ESADE Business School
ioannis.evangelidis@esade.edu

**ABSTRACT.**

Psychology experimenters choose stimuli to indirectly manipulate latent variables that cannot be directly manipulated (e.g., trust, impatience, and arousal). Stimulus selection is typically unsystematic, undocumented, and irreproducible. This makes confounds likely to arise. Study results, in turn, are typically reported at the aggregate level, averaging across stimuli. This makes confounds unlikely to be detected. Here we propose changing both the design and analysis of psychology experiments. We introduce "Mix-and-Match", a procedure to systematically and reproducibly stratify-sample stimuli, and "Stimulus Plots", a visualization to report stimulus-level results, contrasting observed with expected variation. We apply both innovations to published studies demonstrating how things would be different with our reimagined approach to stimulus sampling.

It is tempting to assume that random assignment justifies making causal claims based on results from psychology experiments. This, however, is generally not the case, at least not for the causal claims of interest to psychologists. The reason is that, in contrast to the hard and some applied sciences, where experimenters can directly manipulate the independent variable of interest (e.g., physicists can directly manipulate an object's mass, economists can directly manipulate what the default option is on a tax form), many psychology experiments examine hypotheses about latent variables that are neither observable nor directly manipulable (e.g., trust, self-worth, impatience, risk-tolerance, arousal). Psychology researchers, therefore, typically indirectly manipulate variables of interest by randomly assigning participants to conditions with different stimuli. Given that stimuli are multidimensional, any two stimuli that participants are randomly assigned to will typically differ not only on the focal dimension the experimenter wishes to manipulate (e.g., the emotional reaction they induce), but also on other dimensions the experimenter does not wish to manipulate. In psychology we randomly assign stimuli to participants, but we seldom randomly assign attributes to stimuli.

For example, in his influential article on the analysis of experiments with multiple (word) stimuli, Clark (1973) discusses experiments by Rubenstein, Lewis, and Rubenstein (1971) which contrasted how long it took participants to recognize words as valid, when the words had homophones (e.g., 'maid' , 'made') vs when they did not (e.g., 'pest'). Clark noted that words have many attributes that impact how long it takes to recognize them as valid, such as length, meaning, spelling difficulty, etc. Comparisons between words with and without homophones are confounded.

Rubenstein et al. randomly assigned participants to words with vs without homophones, but obviously did not randomly assign words to have or not have a homophone, thus the correlation

between whether a word has a homophone and participants' time to recognize it is just that, a correlation; one which does not warrant causal interpretation, because words with and without homophones likely differ on other dimensions too.

Clark (1973) proposed, as have many methodologists in the decades since (e.g., Baribault et al., 2018; Judd, Westfall, & Kenny, 2012, 2017; Wells & Windschitl, 1999), that the way around this problem involves using many rather than few stimuli.[1] The idea is that selecting a large enough sample of stimuli will guard against the possibility that the results are due to the particular stimuli that were chosen. This recommendation follows from these authors having diagnosed the issue as a problem of external validity[2].

We propose here that external validity is the wrong diagnosis.

We believe the issue is not whether the stimuli that were chosen have the same effect as do the stimuli that were not chosen, but rather, whether the stimuli that were chosen have an effect *for the hypothesized reason*. The correct diagnosis, in our view, is that poorly selected stimuli, whether few or many, challenge internal rather than external validity.

Once we accept that diagnosis, that the challenge is to internal validity, the approach to choosing stimuli, to analyzing data from experiments with multiple stimuli, and to interpreting those results, changes. So, *everything*, changes.

Let's focus first on that consensual view we challenge here, the need to run *many* stimuli[3]. The *number* of stimuli used in an experiment *does not* actually matter very much for internal

---

[1] This literature, in turn, is related to an earlier debate in psychology on whether it is important for paradigms and stimuli to be ecologically valid by representing the context in which the studied phenomena occur. See for instance the article Brunswik (1955) and the rest of the special issue published in *Psychological Review* V62(3).

[2] Wells and Windschitl (1999) write that "failure to sample stimuli also can threaten ***construct validity***." (emphasis added; their abstract). But as we document in Supplement 5, all arguments in their article involve external rather than construct validity.

[3] Clark (1973) calls for many more than 20 words as stimuli, Judd et al. (2012) for 30 or 50 or more stimuli, Baribault et al. (2018) considers experiments with 100s of stimuli.

validity. There is no reason to expect that, in the population of all words, those with vs without homophones are matched on all confounds that impact how easy it is to recognize a word (e.g., that they have the same average length, the same average pronounceability, etc.). Therefore, there is no reason to expect that a sufficiently large sample of words with vs without a homophone differ, even on average, only in having a homophone. There is no reason for the first 10 words Rubenstein et al. chose to be more biased than the next 10 words, nor to expect the bias of the first 10 words to cancel out the bias of the next 10. A sample of 10 basketball players over-estimates human height. A sample of 1000 basketball players does also.

Even if Rubenstein et al. (1971) had included every word in the English Oxford Dictionary as stimuli in their study, the causal inference problem would remain unchanged. We still would not know if observed differences between all words with vs all words without a homophone occur *because* some words have homophones. To address bias we don't need bigger samples of stimuli, we need better samples of stimuli.

A key realization is that psychologists do not run studies to learn about the properties *of the stimuli* they use, they run studies to learn about *people*. Stimuli are the means, not the end. Rubenstein et al. cared about how language is encoded and retrieved by people, they did not care about the average time it takes to recognize a homophone as a valid word; probably nobody cares about that.

We now switch our working example from homophones to disgusting videos. Several experimenters have examined the causal impact of incidental disgust by having participants watch a toilet scene from the film "Trainspotting", sometimes using sadness as a control condition, e.g., watching a scene from the film "The Champ"*,* where a kid cries over his dead father's body.[4] If

---

[4] Landy and Goodwin (2015), identify four articles that have used the Trainspotting clip to induce disgust in the context of moral judgments. In addition, Lerner, Small, and Loewenstein (2004) use it in an endowment effect study.

these two scenes differed on anything other than the disgusting aspects of the Trainspotting scene, which they obviously do, the disgust manipulation would be confounded. Again, we randomly assign participants to watch a clip, we don't randomly assign disgust to a given movie scene. And, again, simply collecting a large sample of stimuli does not solve the problem, for there is no reason to expect that, on average, disgusting and non-disgusting scenes are matched on all (or any) other attribute that could impact moral judgments. Figure 1 depicts this situation, showing two of many possible confounds in each condition. And again, psychologists do not run studies with disgusting scenes to estimate the average effect of all possible disgusting scenes they could have chosen, instead, they run studies with disgusting scenes seeking to study how the mind reacts to experiencing disgust through an (assumed to be) clean manipulation of disgust.
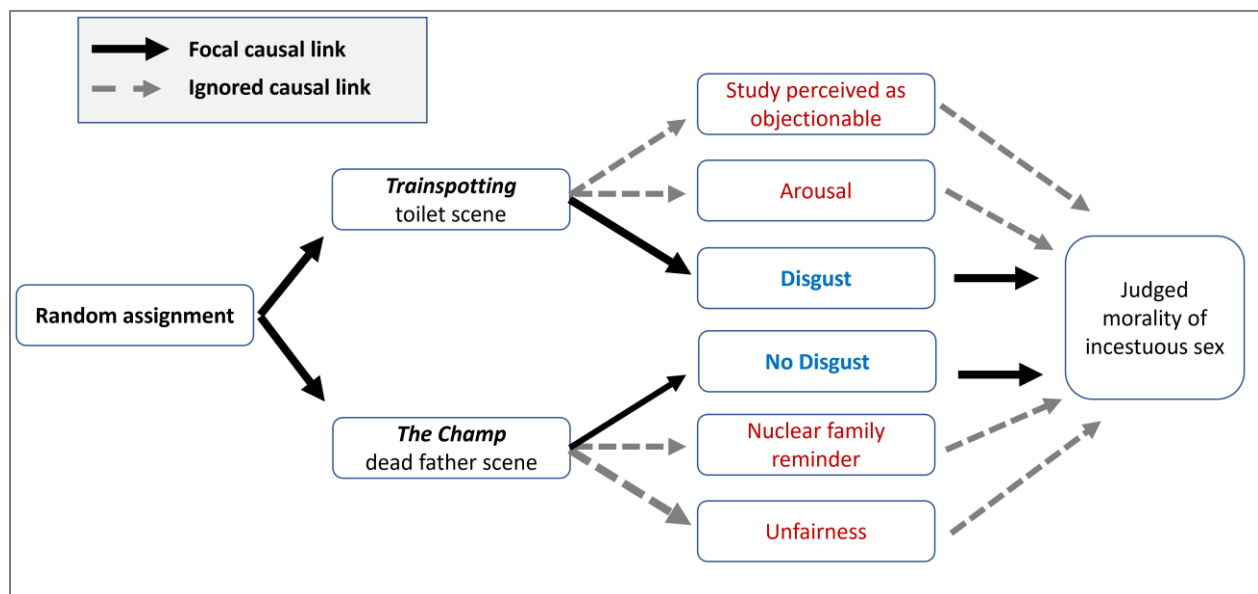


**Figure 1. Example of focal vs confounded causal links in psychology experiments**

In light of this fundamental and ubiquitous challenge to the validity of psychology experiments posed by the fact that stimuli are often confounded, we believe confound management should be at the center of experimental design and analysis.

In this paper, therefore, we reimagine stimulus sampling, the selection of stimuli for a given study (Wells & Windschitl, 1999), focusing on confound management. We propose a concrete procedure for choosing stimuli, along with a simple approach for analyzing stimulus-level results. We believe both are applicable to most psychology experiments.

In terms of generating stimuli: reading papers today, one seldom knows why the specific stimuli used were selected, how they were selected, and what other stimuli the authors would have considered valid substitutes. Papers often discuss confounds of chosen stimuli as afterthoughts that motivate the next study, or in the Limitations sections, or perhaps more often, not at all. Our proposal for generating stimuli, Mix-and-Match, changes all of this.

Mix-and-Match is a systematic and documentable process of stimuli generation which helps researchers be transparent about how and why they operationalize their latent constructs with the chosen stimuli, disclosing the confounds they considered, and how they attempted to address them. Confound management is moved to the earliest part of the discussion of experiments: the design section.

In terms of analysis, we propose "Stimulus Plots", which depict results at the individual stimulus level, identifying which stimuli show the effect and which contribute more or less than expected to the overall effect. Stimulus plots also contrast variability in effects obtained across different stimuli, with the variability that would be expected by chance alone. We demonstrate their use and value reanalyzing data from published papers.

We write this paper with four main goals: (1) that researchers who run studies with only one stimulus per condition, will consider running them with a few stimuli instead, (2) that researchers who run studies with multiple stimuli, will more purposefully, systematically and transparently choose their stimuli (using Mix-and-Match), (3) that authors and readers will no

longer act as if internal (or external) validity have been addressed by the mere fact that a significant overall result is obtained having used many stimuli, and (4) that authors and readers of studies with multiple stimuli will actively explore variation in the results across carefully chosen stimuli, through Stimulus Plots, to explicitly assess internal validity.

**Mix-and-Match: Systematically Generating Stimuli for Psychology Experiments**

We designed Mix-and-Match following three guiding principles. The first principle is that *stimuli* should be blind to hypothesis. It is widely accepted that *participants* should be blind to hypothesis, due in part to concerns of demand effects (see e.g., Rosenthal, 2009). But the notion that *stimuli* (selection) should be blind to hypothesis is seldom if ever considered. The concern we have is that when psychologists choose stimuli, they can often mentally simulate the experiment they are designing, and anticipate whether a particular stimulus is likely "to work". At the same time, it may be difficult to anticipate *why* it may work. This can lead researchers to (possibly unintentionally) be disproportionately likely to select stimuli that work for the wrong reasons (see e.g., Strickland & Suben, 2012). If, instead, experimenters chose stimuli by following a stated and reproducible rule, the stimuli become less *individually* selectable, and thus closer to, if not strictly, blind to hypothesis. Writing down a reproducible rule for selecting stimuli is thus part of Mix-and-Match.

The second principle is that stimuli should be diverse in ways that could help diagnose overlooked confounds. This involves varying stimuli on dimensions directly related to the operationalization of the latent variable of interest. For example, if visual stimuli are chosen to trigger disgust, variation should be along the ways in which disgust can be triggered visually (bodily fluids, pests, rot, etc.). This is the 'mixing' in Mix-and-Match.

The third principle is that there should be an explicit and defensible reason to expect stimuli across conditions to differ only, or at least primarily, on the attribute of interest. This is the 'matching' in Mix-and-Match.

From a confound management perspective, matching seeks to deal with confounds researchers anticipate, by controlling for them, and mixing seeks to deal with confounds they do not anticipate, by exploring variation across diverse stimuli.

We now discuss how to implement Mix-and-Match. First mixing, then matching.

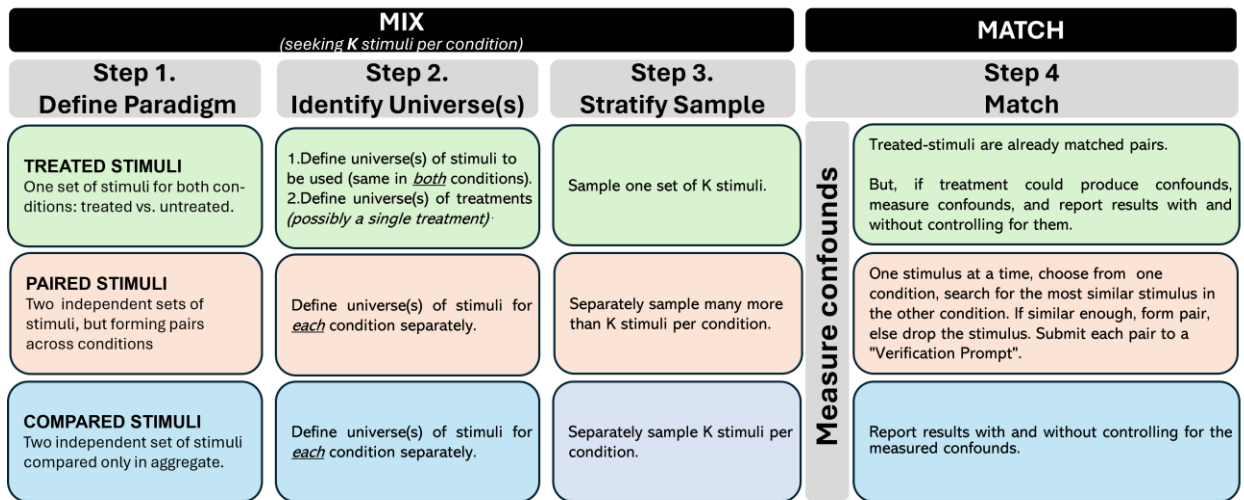Figure 2 provides a flowchart that overviews Mix-and-Match.



**Figure 2.** Overview of Mix-and-Match

*Mixing stimuli*

Mixing puts the sampling in "stimulus sampling". We propose the following three step procedure for sampling stimuli: (i) defining the "experimental paradigm" that will be used, (ii) identifying the universe(s) of stimuli that could be selected or generated for such paradigm, and (iii) stratify-sampling stimuli from those universes.

*(i) Identifying the experimental paradigm.* We define the term 'paradigm' as the description of an experimental procedure that constitutes a practical and valid test of a hypothesis of interest, where every specified design element in the paradigm is necessary for the experiment to be a valid and practical test. For example, an experimental procedure could be described simply as "disgust will be induced, and moral judgments will be elicited." But such level of (un)specificity allows for too diverse a set of stimuli, say, based on disgusting book passages, disgusting videos, and week-long internships in a slaughterhouse. It is *impractical* to include such broad range of stimuli in the same experiment, thus the experimental paradigm should, for practical considerations, entail more narrowly defining how disgust will be induced, e.g., that participants shall read texts of a certain length that describe disgusting scenes. Similarly, moral judgments can be elicited over too broad a range of targets (e.g., vignettes, videos, and in-person biblical reenactments), combining such diverse set of stimuli in a single study would be impractical, thus the paradigm would specify how the ambiguously immoral behavior is presented to participants.

The experimental paradigm, then, needs to be actionably specific. Something like: "participants will read paragraph-long texts, extracted from published books, that induce either disgust or sadness, and will then evaluate the morality of an ambiguous act described in a short vignette, providing their moral judgments in a 1-(very immoral) to 7-(completely moral) scale."

One could add further specificity to this description, e.g., indicating that the experimental paradigm involves reading a passage from the book 'Trainspotting', or evaluating the morality of president Trump kissing his daughter on national TV, but these additional specifications are not justifiable by theoretical concerns (e.g., other inductions of disgusts are equally justifiably ex-ante) nor by practical concerns (it is easy to implement an experiment where different participants read different book segments) thus this description is too specific to meet the definition of 'paradigm'.

We have discussed the importance of sampling stimuli for the independent variable. In terms of the dependent variable, combining results across dependent variables often imposes substantive practical challenges and thus, absent explicit interest in assessing the properties of a dependent variable, the paradigm could specify a single (rather than a set of alternative) dependent variables.

*(ii) Universe of stimuli.* The set(s) of stimuli that meet the description of the experimental paradigm constitutes what we refer to as the universe(s) of stimuli. In our working example, one universe of stimuli involves every passage of text, across all published books, that induces disgust on the reader. Another universe of stimuli is the infinite and uncountable set of vignettes that could be generated to describe a morally ambiguous act.

*(iii) Stratify sampling.* Given our emphasis on internal rather than external validity, we don't propose sampling the universe of stimuli in a representative fashion; in fact, it is often unfeasible and even meaningless to speak of representative samples from a universe with infinite, uncountable, and sometimes simply undefined units (e.g., one cannot draw a representative sample of all possible vignettes that could be written to depict a morally ambiguous acts)[5]. What we propose, instead of random sampling, is stratified sampling. We propose creating possibly arbitrary categories in the universe of stimuli, strata, that are meaningfully different from each other. Categories should differ on a dimension that corresponds to the instantiation of the latent construct of interest rather than a secondary attribute. For example, creating categories of disgust-inducing passages of text that differ *in the origin* of such disgust: sexual, rot, pest, etc., rather than in the

---

[5] A sample of vignettes may not be representative of a population neither in the general statistical sense of the word representative (i.e., used to define a random sample), nor in the sense proposed by Brunswik (1955), where the distribution of stimuli in psychology experiments represents the distribution of stimuli participants might encounter in everyday life.

length of the text or some other superficial feature. Mixing involves only alternative operationalizations of the latent independent variable.

As a default, we propose creating 5 strata but if authors have a reason to choose a different number they should. From each stratum, in turn, experimenters generate (sample) a number of stimuli, we propose 1 or 2 stimuli per stratum as a default, but if authors have a reason to choose another number they should.

It is *not* a problem if the strata aren't exhaustive (e.g., that the 5 strata that are defined do not capture all operationalizations of the construct), nor if different researchers would produce a different stratification. The goal, remember, is not to produce a representative sample of stimuli, the goal is producing meaningfully diverse stimuli selected (largely) blind to hypothesis.

We next propose concrete steps to implement mixing, for categorical stimuli (e.g., videos, faces, names, scenarios) and then for numerical stimuli (e.g., probabilities and monetary amounts).

*Categorical stimuli.* There are multiple approaches that could be relied upon to stratify categorical stimuli, such as relying on categories from a third party (e.g., consumer good categories at Amazon.com) or prior research (e.g., the disgust categorization by Haidt, McCauley, and Rozin (1994)). But, we propose relying on generative artificial intelligence (AI) tools like ChatGPT as a default, whenever another categorization approach is not explicitly preferred by experimenters. AI provides an easy, stimulus blind, and documentable approach for implementing the stratification, and sampling, of categorical stimuli: the approach requires specifying a 'prompt'. We propose the following template 'stimulus prompt': "*please generate 5 categories of <stimulus universe> that differ in <dimension used to create categories> and provide two specific examples of <stimuli> for each.*" Sometimes it helps to provide an example of one category and stimulus within it.

That second placeholder, *'dimension used to create categorie*s', is the more challenging one to specify; experimenters need to consider on what aspect they want the multiple stimuli to vary, striving to generate stimuli that are meaningfully different, hopefully entailing quite different instantiations of the latent variable of interest (rather than the same instantiation differing across stimuli on a superficial attribute). The examples below and in the supplement may be useful for appreciating its role.   Applying this 'stimulus prompt' to our working examples:

Prompt 1: *Please generate 5 categories of <u>homophones</u> that differ in their <u>etymological</u> origin, and provide two examples of specific <u>homophones</u> for each.*

Prompt 2: *Please generate 5 categories of <u>book scenes</u> that may induce disgust that differ <u>in the origin of the disgust</u> being induced, and provide two examples of specific <u>books of fiction</u> containing such scenes (e.g., the category of bodily fluids could contain a passage from the toilet scene in Trainspotting).*

The categories produced by prompt 1 involved homophones that: originate in a different language, have different roots in the same language, involve different parts of speech, have different derivational processes, and were impacted by different sound changes. The examples were: "flour/flower", "knight/night", "rays/raise", "maid/made ", and "son/sun".

Prompt 2 led to stratifying disgust by its source: bodily fluids, filth, putrefaction, gross-out horror, and moral repugnance. The examples provided involved segments from the books "The Road", "The Sisters Brothers", "The Shining", "Haunter", and "Lolita", respectively.

Answers provided by AI include a random component, and the algorithms and training data are often updated, thus the same prompt may lead to different results over time. Moreover, different researchers may choose different stratification strategies/prompts. This idiosyncratic variability is again fine because the goal is internal validity, not generalizability.

*Numerical stimuli.* For numerical stimuli, for example monetary outcomes or probabilities, we propose including in the paradigm definition the set of numbers that would be considered a practical and valid test of the hypothesis of interest (e.g., that to facilitate mental calculations the numerical stimuli need to be multiples of 100, but smaller than 10,000, and the probabilities should be multiples of 10% and smaller than 100%). For stratified sampling one could then choose a diverse set of numbers spanning the range of the consideration set. We exemplify this in Supplement 4, by providing a Mix-and-Match based design of the classic "Asian Disease" problem (Tversky & Kahneman, 1981).

*Matching stimuli*

The "Match" in Mix-and-Match involves striving to generate stimuli that across conditions differ only, or at least primarily, on the focal attribute of interest to the experimenter; striving to match stimuli on all identified potential confounds across conditions. Ideally stimuli are individually matched, so that every stimulus in one condition is paired with a matched stimulus in another condition, providing multiple mini-replications within a study. Matching can be achieved through three alternative study designs: (i) treated-stimuli design, (ii) paired-stimuli design, and (iii) compared-stimuli design. The first two obtain that ideal of individually matched stimuli.

In treated-stimuli designs, stimuli are selected for one condition, and those stimuli are either treated (modified) to be used in the other condition, or used in both conditions in the presence vs. absence of the treatment of interest. For example, experiments examining how a given item is valued when being bought vs sold, how a given fake story is treated when it has been previously encountered vs when it is encountered for the first time, or how the same computer-generated face is treated when depicted as having black vs white skin, are all experiments that *treat* stimuli. In treated-stimuli design, the stimuli are naturally paired: treated vs untreated versions of the same

stimulus. It's worth noting that *treatments* may be confounded; the question of whether pairs of treated/untreated stimuli differ only on the dimension of interest should be explicitly argued for by experimenters, and evaluated by readers. The "verification prompt" we propose later can be used for such purposes.

In paired-stimuli and compared-stimuli designs, stimuli are sampled separately across conditions. Examples include experiments examining how participants respond to photographs of White vs Black faces, male vs female names, experiential vs material purchases, disgusting vs sad videos, words with vs without homophones, and verbal vs math problems. These designs are naturally more challenging from an internal validity perspective than are treated-stimulus designs, because stimuli can differ on many, possibly infinite, non-focal attributes across conditions.

To match stimuli in such designs requires identifying confounding variables (ways in which the stimuli may differ in their impact on the dependent variable other than the focal mechanism), and then measuring those confounding variables for candidate stimuli. For example, for homophones, one identifies other word attributes that may influence how quickly people can recognize them as valid words, and measures those attributes: say, word frequency, language origin of the word, spelling difficulty, etc.

We propose relying on AI for identifying confounding variables as well, relying on this template 'confounds prompt': *"what variables might you expect to predict variation in <dependent variable> across <class of stimuli>?".* For instance, the 'confounds prompt' *"what variables might you expect to predict variation in reaction time to recognize a word as valid, across different words?"* led to identifying 10 confounders, including word length, frequency, phonological regularity, and semantic transparency. Researchers may need to apply (disclosed) judgment to

filter the suggestions (e.g., excluding the manipulation of interest, far-fetched suggestions, and unimplementable suggestions). See examples of this in the supplement.

Having identified confounds, researchers then measure the candidate stimuli on those attributes (e.g., with a pilot study where participants rate the stimuli). For a paired-stimuli design, pairs of stimuli across conditions are formed by matching a target stimulus, say the word "bear", to a matched control stimulus, the words most similar to "bear" on all measured attributes (a 'nearest neighbor' approach). If a particular target stimulus lacks a sufficiently similar control based on the measured covariates, then it probably should not be used at all; otherwise, it introduces an unsolvable confound. If there is no close enough non-homophonic word neighbor to "bare", then it is not used in the study.

Sometimes such paired-stimuli designs may be unfeasible, e.g., stimuli are not selectable or modifiable at a sufficiently granular level to allow forming pairs that differ only in the focal attribute (e.g., it may be unfeasible to create pairs of videos that differ only on whether they are sad vs disgusting). In such cases, we would recommend that experimenters consider changing the paradigm (e.g., inducing emotion with vignettes instead of videos). If the paradigm must be used (e.g., because the manipulations are of intrinsic interest, such as assessing the impact of violent videos), then we have a 'compared stimuli' design, where a set of stimuli in one condition are compared to a set of stimuli in the other. Here experimenters may rely on a statistical model (e.g., linear regression) to control for the confounding variables. For example, this could involve running an emotion induction task using various disgust and sadness videos (say, 10 of each), and reporting the effect of disgust vs sadness controlling vs not-controlling for other attributes identified as potential confounds, measured for each video. Intuitively, one looks for *absence* of mediation for the confounds.

For paired-stimuli designs, we propose a final check to validate pairs. We once again rely on AI for this, and propose submitting each stimuli pair to a 'verification prompt': "*I am going to describe two <stimuli>, please identify 5 consequential differences between them that may impact <the dependent variable>*". If none of the 5 consequential differences are deemed plausible confounds by the experimenter, the stimuli-pair is ready for use. In some paradigms this final check may be redundant and thus unnecessary.
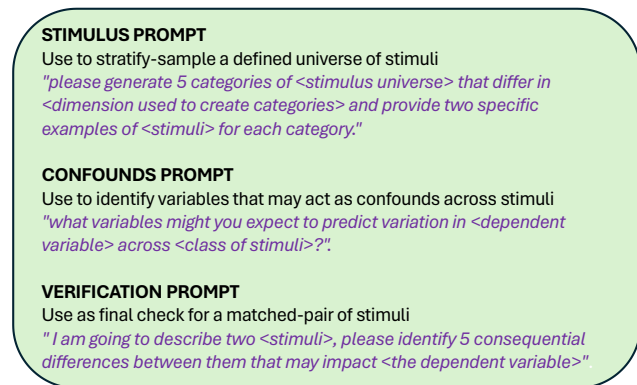
**STIMULUS PROMPT**
Use to stratify-sample a defined universe of stimuli
*"please generate 5 categories of <stimulus universe> that differ in <dimension used to create categories> and provide two specific examples of <stimuli> for each category."*

**CONFOUNDS PROMPT**
Use to identify variables that may act as confounds across stimuli
*"what variables might you expect to predict variation in <dependent variable> across <class of stimuli>?".*

**VERIFICATION PROMPT**
Use as final check for a matched-pair of stimuli
*" I am going to describe two <stimuli>, please identify 5 consequential differences between them that may impact <the dependent variable>"*

**Figure 3.** *Three AI Prompts to aid in Mix-and-Match*

We believe this 'verification prompt' may be useful not only for authors who are designing new experiments, but also for readers evaluating previously published work. We provide an example for a study by Salerno and Slepian (2022) in the next section of this paper.

In Supplements 1-4 we apply Mix-and-Match to four different study designs from published papers. It is applied for choosing vignettes for a moral psychology study (Salerno & Slepian, 2022), photographs for a power posing and race study (Karmali & Kawakami, 2023), social media posts for a misinformation study (Pretus et al., 2023), and the classic Asian Disease problem by Tversky and Kahneman (1981).

**Stimulus Plots**

Only by analyzing data at the individual stimulus level can the main goal of stimulus sampling be achieved: assessing internal validity.[6] Estimates are necessarily noisier when based on subsets of data, therefore, the expectation should not be that every stimulus is individually statistically (or practically) significant (or even with estimates of the same sign). Even if stimuli had the same true effect, because of sampling error, different stimuli will have different effect size estimates. Rather than conducting confirmatory analysis on each individual stimulus, the idea is to conduct exploratory analysis across them. To enable answering questions like: Is the effect evident only for a small subset of stimuli? Does a surprising share of stimuli show an effect in the opposite direction? Are there outlier stimuli with surprisingly big or small effects that may shed light on confounds, mediators, or moderators?

We propose analyzing individual stimuli relying on what we refer to as "Stimulus Plots", plotting stimuli-level results side-by-side, sorted by effect size. We suggest two figures: one that plots the means by condition, and another that plots the differences of means across conditions (note: proportions are also means). While these plots are exploratory, we propose also visually contrasting the observed heterogeneity of effect size across stimuli, with that which would be expected if all stimuli had the same effect size (under 'homogeneity'). This contrast helps calibrate the meaningfulness of differences in observed effect sizes, preventing researchers from over-interpreting random noise, and assessing if a pattern of interest is actually surprising. We will provide an R package, 'stimulus', with a function that makes Stimulus Plots in one line of code.

---

[6] We have come across some papers that report stimulus-level results (see e.g., Bar-Hillel, Maharshak, Moshinsky, & Nofech, 2012; Bartels, Li, & Bharti, 2023; Dias & Lelkes, 2022; Evangelidis, Levav, & Simonson, 2023; Novoa, Echelbarger, Gelman, & Gelman, 2023). But, it does not seem that this was done with the goal of assessing internal validity, and they did not contrast observed with expected variation. We believe our proposed stimulus plots would have added to the informativeness of even these papers that already reported stimuli-level results.

We next re-analyze data from three published papers relying on stimulus plots.


*Example 1. Some stimuli show no effect, some show huge effects*

In their Study 4, Salerno and Slepian (2022) examine whether people report that revealing another person's secret as punishment is more acceptable when the secret involves an intentional rather than unintentional transgression. The authors tested this prediction with a treated-stimulus design involving 20 vignettes. Each vignette had an intentional and an unintentional version. For example, in one vignette (referred to as *'drug'* in our Figure 4 below), the intentional condition read "*Ross brought illegal party drugs to a party, which he then took when he got there.*", while the unintentional one read "*Ross went to a party and, and although he had decided beforehand, he would not take any illegal party drugs, a friend offered him some, and in the heat of the moment, he said yes.*" (see their Appendix C; p.24).[7]

The authors report only an overall effect across all stimuli: higher average acceptability of revealing secrets of intentional acts, $M_1=2.55$ vs $M_2=3.20$, $p<.001$. We reanalyzed their posted data and created the stimulus plots reported in Figure 4 below.  The left panel shows that about 7 stimuli exhibit not difference across conditions, while several stimuli show very large differences. The right panel contrasts this variation with what would be expected if all stimuli were equally effective, and all variation were due to sampling error. The data exhibit much more heterogeneity than expected. For instance, we see that under homogeneity (see line in purple region), we expect the single smallest effect to be larger than what the smaller *seven* stimuli are in actuality.

---

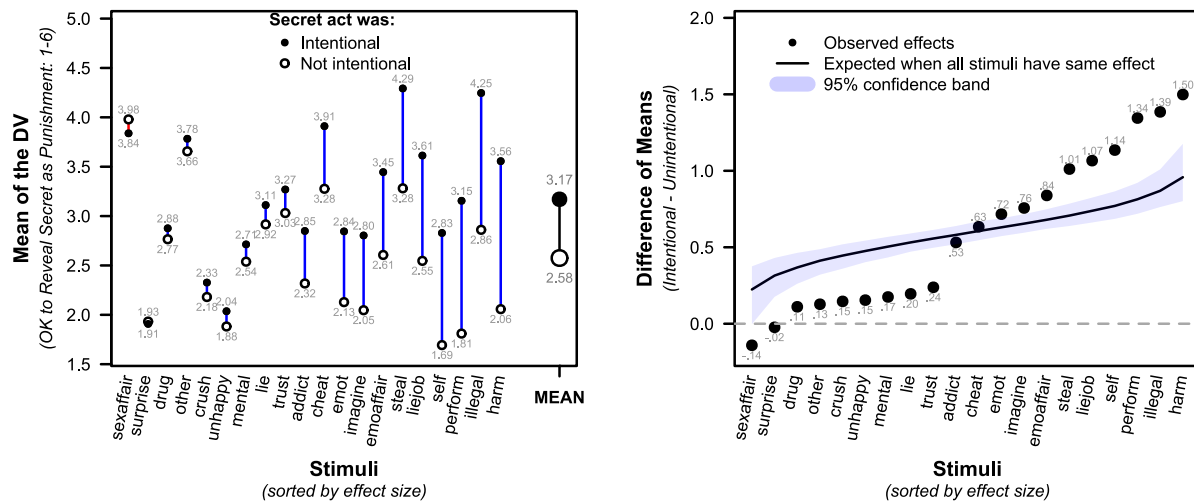[7] The word "and" appears twice in a row in the original text we quote from.

**Figure 4. Stimulus Plots for Study 4 in Salerno & Slepian (2022)**
The study involves a 2-cell treated-stimuli design, comparing participants' willingness to reveal another person's secret based on whether the transgression was intentional or not intentional. The expected line, and its 95% confidence interval, in the right panel, are obtained via resampling.
R Code to reproduce the figure: https://researchbox.org/2257.5.5 (code CXUWHS)

That a substantial share of stimuli "do not work" in this study does *not* necessarily invalidate the main conclusion, but does warrant a deeper exploration of the design and results than is provided in the article. For example, are there moderators or confounds that may explain why the effect is absent for several stimuli but quite large for other stimuli? Figure 4 drew our attention to the vignette leading to the largest effect, "harm", which involves a vignette about John cutting himself intentionally ("to deal with his emotional pain"), vs unintentionally ("while chopping vegetables"). See Appendix C in Salerno & Slepian 2022, p.24. We wonder whether the large difference in willingness to reveal that John cut himself across condition may arise because respondents wished to *help* John with his self-cutting problems rather than to *punish* him. Relying on our proposed verification prompt (see Figure 3), ChatGPT also noted that divulging the secret in the intentional harm condition could be "morally justified if it leads to him getting the help he needs". It is speculative of course, whether that's why that stimulus shows such a large effect. But

speculation is the goal of stimulus plots. Generating hypotheses about surprising variation in effect size that can be explored with more data either before or after the work gets published.

*Example 2. Almost half the stimuli show the opposite effect*

Karmali and Kawakami (2023) examine differences in how Black vs White people are perceived when assuming expansive vs constrictive poses. Their paper reports 4 studies, all relying on the same photographs of 20 Black and 20 White men assuming two different expansive and two different constrictive poses.[8] Study 3 is the one we focus on, because its stimulus plot revealed the most information left unexplored in the original paper.

In Study 3, n=105 undergraduates were asked to choose potential partners for an upcoming task. They saw 10 sets of 4 photographs of different people, and they chose one out of the four in each set as a potential partner. The study's key finding is that White partners were chosen more often when in an expansive than constrictive pose (Z=4.96, p<.001), but that this effect of pose was not observe for Black partners (Z=1.26, p=.208); a race *x* pose attenuated interaction (Z=2.47, *p*=.013).

The authors posted the raw data which we re-analyzed to create Figures 5 and 6 below. We observe that while *on average* Black potential partners are not more or less likely to be chosen in expansive rather than constrictive poses, posing has a highly heterogeneous effect. There are eight Black potential partners which exhibit a *negative* effect of expansive posing, seven of which show an effect *bigger* in magnitude than the average (positive) effect for White potential partners. There are, however, also several Black potential partners showing strong effects in the oppositive direction, cancelling out on average.

---

[8] The design involves 5 expansive and 5 contractive poses. Any given potential target was shown in 2 out of 5 poses of each kind.
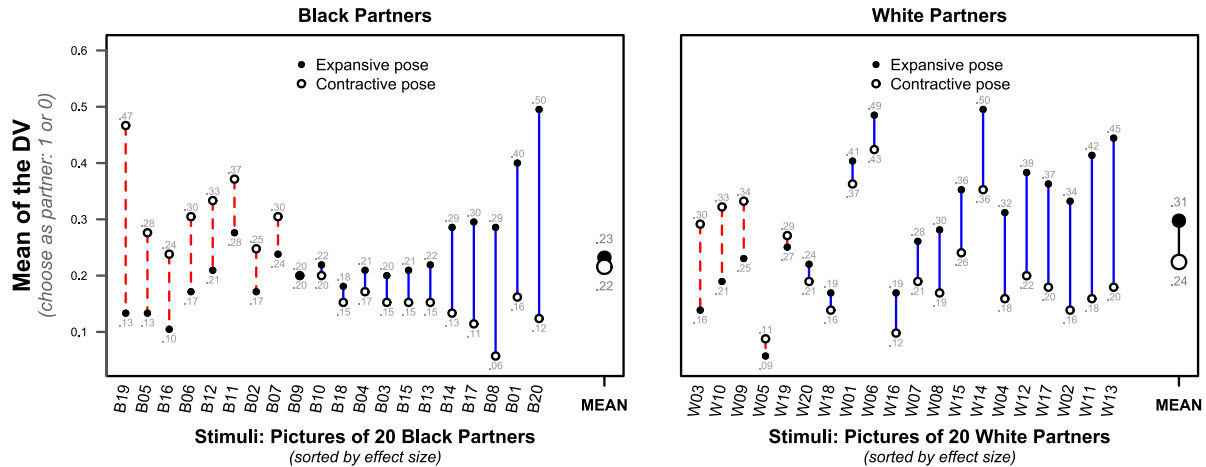
**Figure 5. Stimulus Plot for Means in Karmali & Kawakami – Study 3**
The study involves a 2 (race [compared]) x 2 (power posing [treated]) stimuli design. Participants chose 1 of 4 potential partners based on photographs where they were either in an expansive or a contractive pose. The figure depicts the percentage of times each stimulus (potential partner) was chosen.
R Code to reproduce the figure: https://researchbox.org/2257.8  (use code CXUWHS)

Figure 6 reports mean differences, observed and expected under homogeneity across stimuli. The eight Black targets with a negative effect are outside the 95% confidence band. We believe the stimulus plots from Figures 5 and 6 make clear that there is important heterogeneity to explore before interpreting the results from this study in the way they have been interpreted.
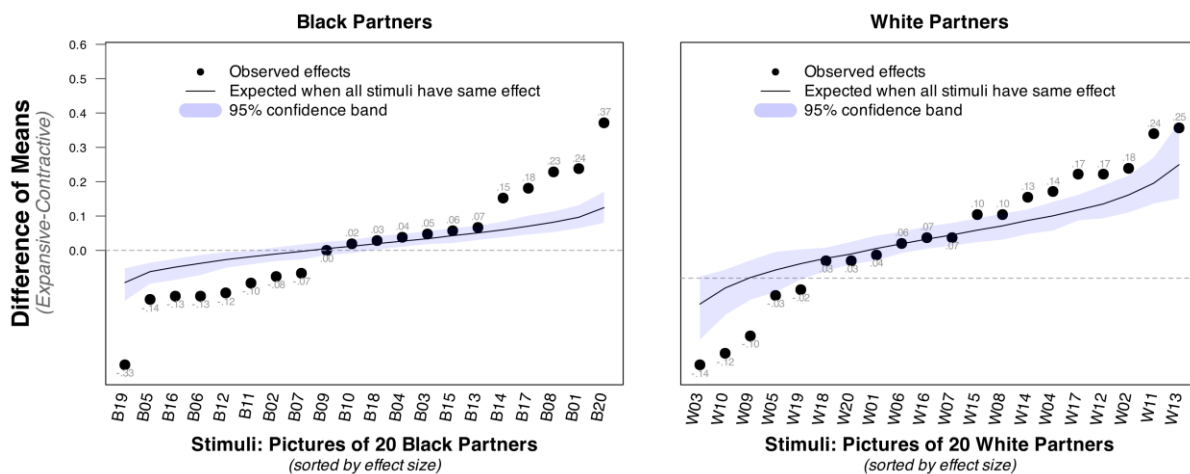


**Figure 6. Stimulus Plot for Effects in Karmali & Kawakami – Study 3**
Differences computed off means from Figure 5. The expected line, and its 95% confidence interval, are obtained via resampling.
R Code to reproduce the figure: https://researchbox.org/2257.8  (use code CXUWHS)

*Example 3. All stimuli are consistent*

Pretus et al. (2023) examine the psychological processes that underlie misinformation sharing. In Experiment 2 they asked N=797 participants how likely they would be to share a tweet, (which contained misinformation) on a 1-6 likert scale. The authors relied on 16 different tweets, and the manipulation of interest to us is whether the tweet was accompanied by a Twitter fact-check message (their design is more complex and includes additional manipulated and measured differences). The study involved a treated-stimulus design in that the same tweet was presented with or without a fact-check message. The paper reports an overall average effect of the fact-check of M=0.16, *p*=.006 (p.3124).

Relying on data provided by the authors upon request (they had posted the data, but not with individual stimuli identifiers), we created Figure 7 below. The left panel suggests substantial variation in effect size, but the right panel suggests differences in results across stimuli are entirely consistent with sampling error. If differences across stimuli of this magnitude were considered important, then a larger sample would be needed to explore them.
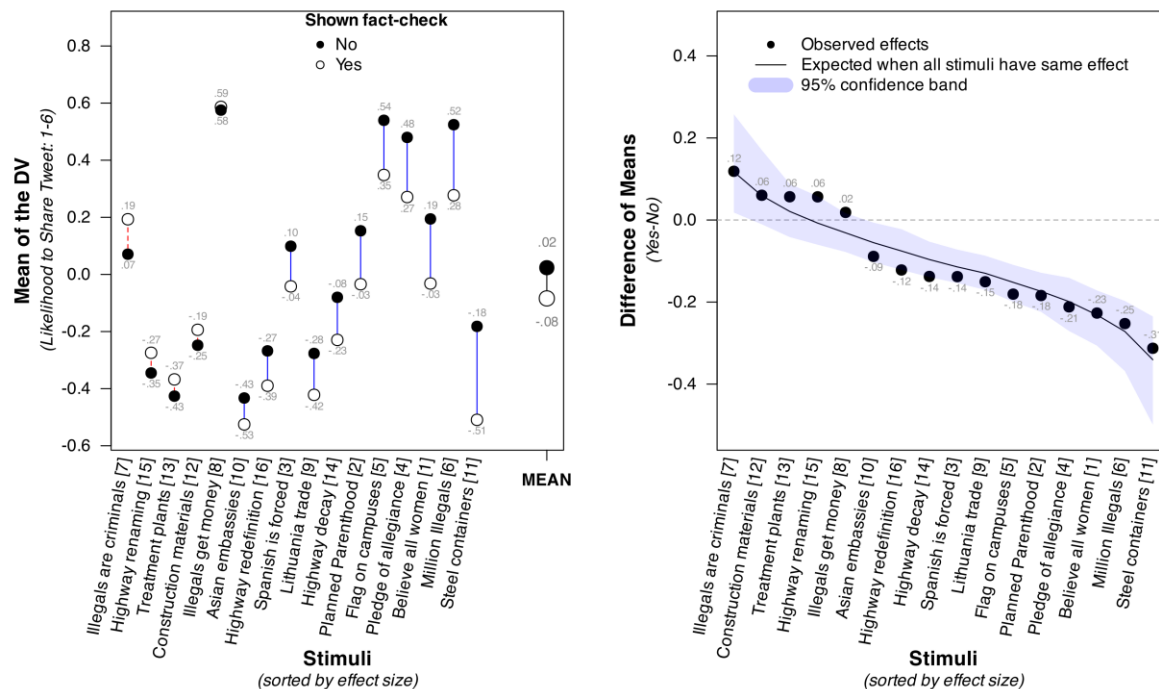
**Figure 7. Stimulus Plots for Pretus et al. (2023) – Study 2**

The study involves a two-cell treated-stimuli design, comparing participants' reported willingness to share a tweet containing information having been presented, or not, with a Twitter fact-check. The expected line, and its 95% confidence interval, are obtained via resampling.

R Code to reproduce the figure: https://researchbox.org/2257.9 (use code CXUWHS)

**Stimulus Plots Contribution to the Statistical Analyses of Clustered Data**

To appreciate how stimulus plots upend the conclusions we may draw from traditional statistical analyses of experimental results, in this section we carry out such analyses for our three example and contrast the conclusions with vs without considering stimulus plots.

Data from experiments with multiple stimuli are 'clustered': a given condition has data from multiple stimuli, and multiple observations are often collected from individual participants. Nested data are commonly analyzed in one of two ways: with (1) regressions that include fixed-effects, clustered standard errors, or both, or with (2) mixed-models models with random intercepts and possibly random slopes. McNeish (2023) provides a detailed overview of both approaches and

how they relate to each other and to the research question motivating data collection. We conducted our analyses with both approaches.

*Estimated models*

We estimated three models for each dataset. First, we estimated a regression model with stimuli and participant (ID) fixed effects, i.e., 1/0 dummies for each ID, to control for variability across them (this increases power), and clustering errors by participant, to account for any lack of independence that may remain across observations by the same ID beyond their fixed effect.

In R Code: `miceadds::lm.cluster(dependent variable~condition+factor(stimulus)+factor(ID), cluster=ID)`

Second, we estimated a mixed-effects model that included random intercepts for participants and for stimuli; these increase power by controlling for stimuli heterogeneity and (partially) address dependence across multiple observations by the same participant.[9]

In R Code: `lme4::lmer(dependent variable ~ condition + (1|stimulus) + (1|ID))`

Third, we estimated a mixed-effects model that includes also 'random slopes' for stimuli. These 'random slopes' correspond to the effect of the manipulation for each stimulus (e.g., what the effect of intentionality was on the tendency to divulge a secret for each of the vignettes).

In R Code: `lme4::lmer(dependent variable ~ condition + (1+condition|stimulus) + (1|ID))`

A mixed-model that includes random slopes typically enlarges the confidence interval for the main effect of the manipulation based on the variability of the estimated effect across stimuli, the standard deviation across slopes (e.g., the confidence interval around the overall estimate on

---

[9] We say "partially" accounting for dependence because a random intercept only accounts for dependence that arises from different participants having different mean evaluations. It's equivalent to the fixed effect but the mixed model lacks the catch-all ability to account for *other* dependency that the regression has with clustered errors (Abadie, Athey, Imbens, & Wooldridge, 2023; McNeish, 2023).

divulging secrets is increased by the variability in the observed effects across the 20 vignettes). Including random slopes, then, usually lowers statistical power.

It has been proposed that (these power lowering) random slopes for stimuli should be used whenever the data allow, in order to generalize findings obtained with some stimuli to results that could be obtained with other stimuli, and that without random slopes the overall test of the manipulation has an elevated false-positive rate (see e.g., Barr, Levy, Scheepers, & Tily, 2013; Brauer & Curtin, 2018; Judd et al., 2012; Oberauer, 2022; Wickens & Keppel, 1983). We believe these claims rest on rather unrealistic assumptions (Simonsohn, Montealegre, & Evangelidis, 2024) and are not persuaded random slopes are necessary for valid inferences. But we report results with and without random slopes, readers can individually decide which to focus on.

*Results*

Figure 8 shows the results of these three alternative models estimated on each of the datasets discussed above. First, note that for any given dataset, the point estimates of the average effect of the manipulation is essentially identical across the three procedure. Second, the confidence intervals are essentially identical for the regression and the mixed model with random intercepts.[10] Third, in Examples 1, 2a and 2b, the confidence interval gets meaningfully wider upon including random slopes, this occurs because there is meaningful heterogeneity in the effect size across stimuli for them (see again their stimulus plots). As mentioned earlier, random-slopes models add such heterogeneity to the confidence interval of the overall effect, hence lowering power. Because in Example 3 there is no more heterogeneity than expected by chance (see Figure 7), the random slopes leave the confidence interval unchanged.

---

[10] This implies there is no consequential dependence within participants after accounting for their mean response (their fixed/random effect).
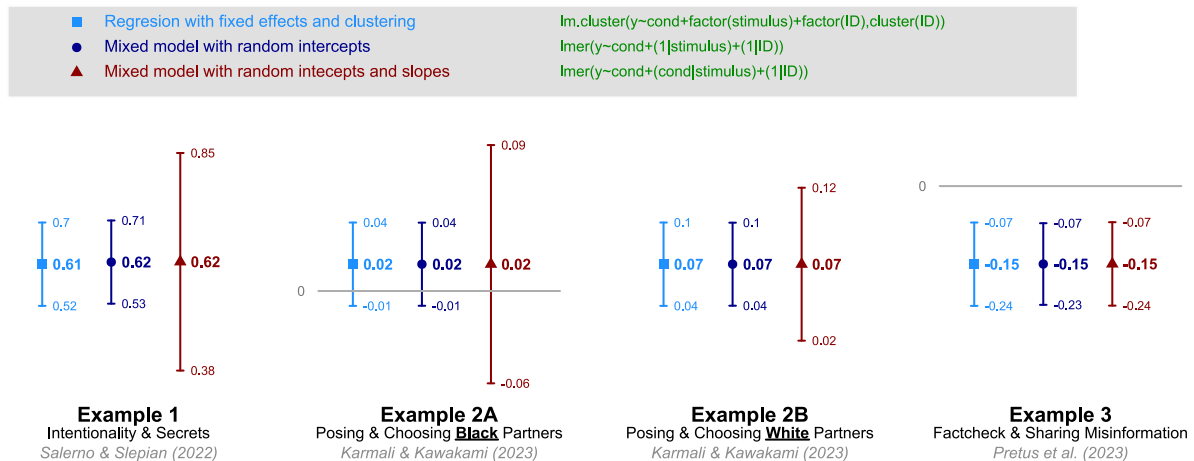
**Figure 8. Estimates of overall effect and its confidence interval for the three examples.**
R Code to reproduce the figure: https://researchbox.org/2257.10 (use code CXUWHS)

It's interesting to contrast the insights one arrives at when looking at these summary overall results alone, to those one arrives at when also inspecting the stimulus plots. In Example 1, the regression and mixed-models point to a robust, substantial, and significant effect of intentionality on divulging secrets. Only looking at the stimulus plots do we realize a surprising share of stimuli did not work, and only looking at the stimulus plots is our attention driven to outlier results that may have been impacted by confounding intentionality with other differences across the two scenarios (e.g., desire to help a person who cuts himself on purpose). In Example 2A, the regression and mixed results both lead to a relatively tight confidence interval for no effect. We would conclude, as the original authors did, that "pose did not impact choice of Black partners" (p.59). Only looking at the stimulus plot do we realize that, actually, pose impacted the *majority* of Black partners but that effects of opposite sign cancel out on average.

Finally looking at Example 3, on the one hand, the stimulus plot does not actively add information, for it shows lack of (surprising) variation across stimuli, and thus the overall summary from the regression/mixed models is interpretable as is. But on the other hand, only with the stimulus plot do we realize the overall average is an apt summary of the results.

## Choosing the statistical analysis

Because regressions with fixed effects or clustered errors are not currently commonly used in psychology, in Figure 9 below we provide a brief overview mapping design decisions to the elements that must be included or are beneficial to include in the regression. Specifically, clustered errors prevent inflated false-positive rates due to lack of dependence, while participant fixed effects can increase power, and stimulus fixed effect can also increase power if participants see a subset of all stimuli in a condition. We created this overview to help researchers determine the 'right analysis' for their designs. Which design to choose, however, should be based on substantive considerations related to the particular research question at hand. We note that for binary dependent variables, the mixed model is particularly sensitive to violations of the strong assumptions it is based on (Grilli & Rampichini, 2015; Heagerty & Kurland, 2001), providing a stronger justification to rely on regressions with clustered errors.

| DESIGN | | | | ANALYSIS | | | |
|---|---|---|---|---|---|---|---|
| Case | Design | Datapoints per participant | Participants assigned to | Cluster by Participant | Stimulus Fixed Effects | Participant Fixed Effects | Syntax in R |
| 1 | Paired stimuli | 1 | 1 condition | --- | Yes | --- | lm(data=df, y~condition+factor(stimulus)) |
| 2 | Paired stimuli | Many | 1 condition | Yes | Yes | --- | lm.cluster(data=df, y~condition+factor(stimulus),cluster=ID) |
| 3 | Paired stimuli | Many | ≥2 conditions | Yes | Yes | Yes | lm.cluster(data=df, y~condition+factor(stimulus) + factor(ID),cluster=ID) |
| 4 | Compared | 1 | 1 condition | -- | -- | --- | lm(data=df, y~condition) |
| 5 | Compared | Many | 1 condition | Yes | -- | --- | lm.cluster(data=df, y~condition,cluster=ID) |
| 6 | Compared | Many | ≥2 conditions | Yes | -- | Yes | lm.cluster(data=df, y~condition+factor(ID),cluster=ID) |

**Figure 9. Regression analysis for multi-stimuli studies**

For stimuli fixed effect, the same stimulus id must be used for a pair across conditions (e.g., stimulus[1]="face_23", or stimulus[1]="mug"). If authors feel the need to run a logistic regression with clustered errors, they can use: glm.cluster(…, cluster="ID", *family="binomial"*). The function lm.cluster() is included in the 'miceadds' package, it requires indicating the data.frame which is done by including its name as the first argument in the call (see "data" in each row). Other R packages that produce clustered errors include 'lfe', 'jtools', 'estimatr' and 'fixest'.

**General discussion**

In this paper we have proposed Mix-and-Match, a principled, documentable, and standardized process for selecting stimuli for psychology experiments, and have introduced Stimulus Plots, a graphical approach for presenting stimuli-level results. We believe these proposals could help increase, and help evaluate, the internal validity of psychology experiments.

In supplements 1-4 we redesign published studies relying on Mix-and-Match, providing concrete illustrations of how different psychology experiments look, when designed in this manner. In turn, the stimuli plots in Figures 4-7 concrete illustrations of how the inferences we draw from data from psychology change when we consider variation in results across stimuli.

We wrap up this paper touching on a series of issues we expect readers may be thinking about as they reach this paragraph.

*Isn't external validity also important?* Prior papers on the selection and analysis of experiments with multiple stimuli have focused on external validity. We have already argued in detail why the emphasis should instead be on internal validity. But to be clear, we do believe external validity is valuable. If something only happens in contrived lab environments, it is not clear psychologists should care about it, and in any case they should be aware that it does only happen in contrived lab environments. However, we don't think that external validity involves testing different stimuli (which may or may not be internally valid) within the same paradigm. Rather, external validity for an experimental paradigm can only be assessed by collecting data *outside* that paradigm; and to know if a finding is consequential in the real world, a perhaps more common definition of *external* validity, the hypothesis needs to be tested…   …in the real world.

*What about statistical power?* One concern we believe people may have with our call for routinely using multiple stimuli in experiments, is that doing so may lower power to detect an

overall effect. This concern is largely unfounded. Adding stimuli could indeed decrease power if authors knew which stimulus shows the largest effect and were to choose, in the absence of stimulus sampling, that single stimulus for their experiment. But, a more likely scenario is that experimenters don't know for sure which stimuli will show larger effects, and in that case adding stimuli will tend to increase power.[11] Additionally, if one is able to present more than one stimulus per participant, adding stimuli also increases power. Thus, power concerns, if anything generally provide additional justification for multiple stimuli. But in any case, power considerations are not a defensible justification for low internal validity.

*What about costs?* Another concern people may have is that there are experimental paradigms where stimulus sampling could be prohibitively expensive. This is indeed true, especially for field experiments where each stimulus has a large implementation cost. For instance, it may not be possible to attempt 20 or even 5 alternative implementations of a "nudge" in a field experiment. When practical circumstances prevent using multiple stimuli in a study, one could rely on other (presumably lab) studies, in the same paper, or elsewhere, that rely on stimulus sampling to validate the stimulus used in the field study.

*Why within a study?* An interesting question we have received is 'what is the benefit of running one study with many stimuli instead of many studies with one stimulus each?' First, running multiple stimuli with a given paradigm in one study, allows changing the paradigm across studies, which is valuable for internal and perhaps external validity. Second, running the multiple

---

[11] To get an intuition for this: imagine two stimuli, one has a very big effect, the other no effect. Using only one of them, blindly, expected power is 52.5%, if you choose the right one you find p<.05 for sure, if you choose the wrong one you only have a 5% chance. But if you use both, with a big enough sample, power is 100%. More generally, if only some stimuli show detectable effects, and ex-ante is hard to tell which (which we intuit is a quite common scenario), using multiple stimuli can dramatically increase power. This intuition we have that anticipating effect size is difficult may be at odd with our concern that experimenters choose confounded stimuli by simulating the experiment in their heads. The premise is that the effect size variation across stimuli with vs without blatant confounds is typically larger than the variation in effect size among ex-ante valid stimuli.

stimuli in the same study allows differences in results across stimuli to be causally interpretable (as they arise under random assignment and/or from the same participants). Third, transparent reporting of all stimuli attempted is verifiable if done in one study (that's pre-registered), but not across studies (which may be file-drawered).

*Isn't the implementation of Mix-and-Match subjective and arbitrary?* In short. Yes. But… It is *less* subjective and *less* arbitrary than the status quo where researchers follow undisclosed and presumably unsystematic procedures of stimulus selection. Mix-and-Match does not eliminate idiosyncrasies in how psychologists operationalize hypotheses, but it reduces those idiosyncrasies, it highlights them, and it provides a framework for discussing them.

*Doesn't mediation analysis solve the internal validity problem?* The goal of mediation is indeed to ascertain whether a randomly assigned manipulation produces an observed effect through a hypothesized channel. But, it has long been recognized that mediation analysis does not deliver on its stated goal (Bullock & Green, 2021; Bullock, Green, & Ha, 2010; Judd & Kenny, 1981, pp. 607, last paragraph; Rohrer, Hünermund, Arslan, & Elson, 2022). Most notably, mediation analysis is biased towards finding mediation which does not exist under two likely scenarios. First, if the mediator is correlated with the dependent variable outside of the experiment (for the intuition, see Simonsohn, 2022), and second, if the stimuli across conditions differ in more than in the attribute of interest and those alternative mediators are not included in the analysis.

*Possible misuses.* In this paper we have proposed new tools, and all tools from pencils to rearview mirrors, can be misused. We believe misuse may involve mixing and/or matching over superficial dimensions, so that studies do not actually include truly diverse stimuli nor match stimuli on potential confounds. We don't believe it is possible to fully prevent this, but we hope

that the concrete recommendations from Mix-and-Match will reduce unintentional instances of this problem.

In terms of stimulus plots, a possible  misuse involves unreasonably expecting all stimuli to conform to predictions, be it with authors file-drawering results because some stimuli don't behave as expected, or reviewers encouraging authors to "explain" something they can't really explain. We hope the confidence band we include in stimulus plot, and the perspective we have given throughout the article will be effective protection against such misuse.

# References

Abadie, A., Athey, S., Imbens, G. W., & Wooldridge, J. M. (2023). When should you adjust standard errors for clustering? *The Quarterly Journal of Economics, 138*(1), 1-35.

Bar-Hillel, M., Maharshak, A., Moshinsky, A., & Nofech, R. (2012). A rose by any other name: A social-cognitive perspective on poets and poetry. *Judgment and Decision making, 7*(2), 149-164.

Baribault, B., Donkin, C., Little, D. R., Trueblood, J. S., Oravecz, Z., Van Ravenzwaaij, D., . . . Vandekerckhove, J. (2018). Metastudies for robust tests of theory. *Proceedings of the National Academy of Sciences, 115*(11), 2607-2612.

Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of memory and language, 68*(3), 255-278.

Bartels, D. M., Li, Y., & Bharti, S. (2023). How well do laboratory-derived estimates of time preference predict real-world behaviors? Comparisons to four benchmarks. *Journal of Experimental Psychology: General*.

Brauer, M., & Curtin, J. J. (2018). Linear mixed-effects models and the analysis of nonindependent data: A unified framework to analyze categorical and continuous independent variables that vary within-subjects and/or within-items. *Psychological methods, 23*(3), 389.

Brunswik, E. (1955). Representative design and probabilistic theory in a functional psychology. *Psychological review, 62*(3), 193.

Bullock, J. G., & Green, D. P. (2021). The failings of conventional mediation analysis and a design-based alternative. *Advances in Methods and Practices in Psychological Science, 4*(4), 25152459211047227.

Bullock, J. G., Green, D. P., & Ha, S. (2010). Yes, But What's the Mechanism?(Don't Expect an Easy Answer). *Journal of personality and social psychology, 98*(4), 550-558.

Clark, H. H. (1973). The language-as-fixed-effect fallacy: A critique of language statistics in psychological research. *Journal of verbal learning and verbal behavior, 12*(4), 335-359.

Dias, N., & Lelkes, Y. (2022). The nature of affective polarization: Disentangling policy disagreement from partisan identity. *American Journal of Political Science, 66*(3), 775-790.

Evangelidis, I., Levav, J., & Simonson, I. (2023). The upscaling effect: how the decision context influences tradeoffs between desirability and feasibility. *Journal of Consumer Research, 50*(3), 492-509.

Grilli, L., & Rampichini, C. (2015). Specification of random effects in multilevel models: a review. *Quality & Quantity, 49*, 967-976.

Haidt, J., McCauley, C., & Rozin, P. (1994). Individual differences in sensitivity to disgust: A scale sampling seven domains of disgust elicitors. *Personality and Individual differences, 16*(5), 701-713.

Heagerty, P. J., & Kurland, B. F. (2001). Misspecified maximum likelihood estimates and generalised linear mixed models. *Biometrika, 88*(4), 973-985.

Judd, C. M., & Kenny, D. A. (1981). Process analysis: Estimating mediation in treatment evaluations. *Evaluation review, 5*(5), 602-619.

Judd, C. M., Westfall, J., & Kenny, D. A. (2012). Treating stimuli as a random factor in social psychology: a new and comprehensive solution to a pervasive but largely ignored problem. *Journal of personality and social psychology, 103*(1), 54.

Judd, C. M., Westfall, J., & Kenny, D. A. (2017). Experiments with more than one random factor: Designs, analytic models, and statistical power. *Annual review of psychology, 68*(1), 601-625.

Karmali, F., & Kawakami, K. (2023). Posing while black: The impact of race and expansive poses on trait attributions, professional evaluations, and interpersonal relations. *Journal of personality and social psychology, 124*(1), 49-68.

Landy, J. F., & Goodwin, G. P. (2015). Does incidental disgust amplify moral judgment? A meta-analytic review of experimental evidence. *Perspectives on Psychological Science, 10*(4), 518-536.

Lerner, J., Small, D. A., & Loewenstein, G. F. (2004). Heart strings and purse strings - Carryover effects of emotions on economic decisions. *Psychological science, 15*(5), 337-341.

McNeish, D. (2023). A practical guide to selecting and blending approaches for clustered data: Clustered errors, multilevel models, and fixed-effect models. *Psychological methods*.

Novoa, G., Echelbarger, M., Gelman, A., & Gelman, S. A. (2023). Generically partisan: Polarization in political communication. *Proceedings of the National Academy of Sciences, 120*(47), e2309361120.

Oberauer, K. (2022). The Importance of Random Slopes in Mixed Models for Bayesian Hypothesis Testing. *Psychological science, 33*(4), 648-665.

Pretus, C., Servin-Barthet, C., Harris, E. A., Brady, W. J., Vilarroya, O., & Van Bavel, J. J. (2023). The role of political devotion in sharing partisan misinformation and resistance to fact-checking. *Journal of Experimental Psychology: General*.

Rohrer, J. M., Hünermund, P., Arslan, R. C., & Elson, M. (2022). That's a lot to PROCESS! Pitfalls of popular path models. *Advances in Methods and Practices in Psychological Science, 5*(2), 25152459221095827.

Rosenthal, R. (2009). Blind and Minimized Contact. In R. Rosenthal & R. L. Rosnow (Eds.), *Artifacts in Behavioral Research* (pp. 592-602): Oxford University Press.

Rubenstein, H., Lewis, S. S., & Rubenstein, M. A. (1971). Evidence for phonemic recoding in visual word recognition. *Journal of verbal learning and verbal behavior, 10*(6), 645-657.

Salerno, J. M., & Slepian, M. L. (2022). Morality, punishment, and revealing other people's secrets. *Journal of personality and social psychology, 122*(4), 606.

Simonsohn, U. (2022). [103] Mediation Analysis is Counterintuitively Invalid. Retrieved from https://datacolada.org/103

Simonsohn, U., Montealegre, A., & Evangelidis, I. (2024). *Do Random Slopes in Mixed Models Help with Generalizability? No (manuscript in preparation)*.

Strickland, B., & Suben, A. (2012). Experimenter philosophy: The problem of experimenter bias in experimental philosophy. *Review of Philosophy and Psychology, 3*, 457-467.

Tversky, A., & Kahneman, D. (1981). The framing of decisions and the psychology of choice. *Science, 211*(4481), 453-458.

Wells, G., & Windschitl, P. (1999). Stimulus sampling and social psychological experimentation. *Personality and Social Psychology Bulletin, 25*(9), 1115.

Wickens, T. D., & Keppel, G. (1983). On the choice of design and of test statistic in the analysis of experiments with sampled materials. *Journal of verbal learning and verbal behavior, 22*(3), 296-309.