

Plural causes
(c. 15,500 words)

Can Konuk¹, Tadek Quillien², and Salvador Mascarenhas¹

¹Ecole Normale Supérieure, Department of Cognitive Studies, Institut Jean-Nicod

²University of Edinburgh, School of Informatics

Author Note

All data and experimental materials for the studies in this manuscript are available on the Open Science Framework at

https://osf.io/43m5d/?view_only=3d26a80b8e394fa9ad7792a690de8fe6. The authors have no conflict of interest to disclose. This work was supported by *Agence Nationale de la Recherche* grant ANR-18-CE28-0008 LANG-REASON (PI: Mascarenhas). CRediT: Conceptualization CK, TQ, SM; Data curation CK, SM; Formal analysis CK TQ; Funding acquisition SM; Investigation CK; Methodology CK, TQ, SM; Visualization CK; Writing CK, TQ, SM.

Abstract

Causal selection is the process underlying our intuition that an outcome happened *because of* a given event, or that an event is *the cause* of an outcome. When a forest catches fire after a lightning strike, for example, people tend to say that the lightning bolt was the cause of the fire, not mentioning the presence of oxygen in the air, although they are well aware that the latter was no less indispensable for the fire to occur. We argue that the extant literature on causal selection has so far operated on the implicit premise that the only relevant variables for causal selection are *individual variables*, corresponding to distinct nodes in the relevant network of causes. Ours is the first systematic study of plural causes in the context of causal selection. First, we establish by means of two behavioral experiments the psychological reality and non-triviality of plural causes, ruling out potential deflationary explanations. Second, we show that state-of-the-art models of causal selection based on counterfactual dependence can be extended to make non-trivial predictions about plural causes consistent with our experimental findings. Third, we show that surprising logical properties of plurals *in natural language interpretation* can be found in causal reasoning with plural causes. We argue that the mental representations involved in causal-selection judgments show marks of the representations formal semanticists have proposed for plurals in natural language.

Significance Statement

Why is my child not doing well in school? *Why* are stars moving so fast in the outer perimeters of galaxies? Reasoning about causation is central to human intellectual life, because causes provide the answers to *why* questions, both in scientific and everyday discourse. Yet, causal structures are usually too complex to grasp and talk about, so a key task for human understanding is to find ways of *selecting* from a pool of causally active facts the *main cause*. In this work we demonstrate, first, that this task of

causal selection doesn't need to pick a single individual event as the main cause, but instead can pick a *plurality* of events. Second, we argue that these plural causes are represented in human minds with the same resources as pluralities in natural language, suggesting that we can take elements of our theories of natural language as theories of mental representations.

Keywords: causal selection, counterfactual theories of causation, plurals

Plural causes
(c. 15,500 words)

Introduction

Causal selection is the process underlying our intuition that an outcome happened *because of* a given event, or that an event is *the cause* of an outcome. Causal selection judgments go further than judgments of *actual causation* (Halpern, 2016; Halpern & Pearl, 2005; Hitchcock, 2001), whereby people merely identify which events can be counted as causes of an outcome. They cast a hierarchy over these events, singling out some as being more important than others in bringing about the outcome under consideration. When a forest catches fire after a lightning strike, for example, people tend to say that the lightning bolt was the cause of the fire, not mentioning the presence of oxygen in the air, although they are well aware that the latter was no less indispensable for the fire to occur. Causal selection in this sense is crucially distinct from *causal inference*, the problem of learning the relevant causal facts about the world. Causal selection concerns how we judge the relative importance of the many causes of an event, given that we already have a causal model of the situation.

A considerable literature has developed around what factors underly our preference for certain causal explanations of an outcome over others (Icard et al., 2017; Knobe & Fraser, 2008; Morris et al., 2019; Quillien & Barlev, 2022). Although theories diverge as to what the drivers of causal-selection judgments are, they all agree that the outcome of causal-selection judgments depends crucially on the initial pool of candidates under consideration to begin with.

Before the lightning bolt can be viewed as *the cause* of the fire, the events *lightning*, *oxygen*, *dry season*, and others must first be flagged by the mind as relevant candidates for causal selection, whose relative importance in bringing about the outcome will be compared. We argue here that the extant literature on causal selection has had a blind spot regarding that initial pool of candidates: it operates on the implicit

premise that the only relevant variables for causal selection are *individual variables*, corresponding to distinct nodes in the relevant network of causes.

Instead, we argue that causal selection judgments can recognize *plural causes*, featuring more than one variable, as when we say that “the dryness of the season and the strength of the wind” caused the uncontrollable spread of the fire. We argue that such plural causes are treated by the mind as candidate explanations on the same footing as the singular causes that compose them. The same factors that drive the attractiveness (or lack thereof) of singular cause explanations drive that of plural causes.

The idea that causal cognition admits causes featuring several variables is not in itself new. In causal inference, researchers have studied how people infer conjunctive causes, that is factors that act in concert to produce an effect (Novick & Cheng, 2004). The notion of a multivariate cause also plays a role in some theories of actual causation (e.g. Halpern, 2015), and, in a different way, in philosophers’ and economists’ concept of *collective responsibility* (e.g. Arendt, 1987; Miller, 2001). To our knowledge, however, the literature on causal selection judgments has yet to engage with the notion of plural causes.

We present the first systematic study of plural causes in the context of causal selection.¹ This study has three objectives. First, we want to empirically establish the psychological reality and non-triviality of plural causes. We show that people’s judgments about plural causes are sensitive to the prior probabilities of events, a key signature of causal-selection judgments. More importantly, we rule out a possible deflationary explanation for plural causes’ sensitivity to probabilities: that subjects might formulate a judgment about a plural cause like $A \wedge B$ simply by combining in some direct way their judgment about the importance of the individual events A and B that compose it. In so doing we provide evidence that people treat plurals as full-fledged candidates for causal

¹ Our Experiment 1 was presented at the Forty-fifth Annual Meeting of the Cognitive Science Society and published in the society’s non-archival proceedings (Konuk et al., 2023).

selection and engage with them in a holistic fashion.

Second, we want to show how considering plural causes can expand our understanding of the role of counterfactual reasoning in causal judgments. We show that models of causal selection based on the notion of counterfactual dependence can straightforwardly be extended to make non-trivial predictions about plural causes consistent with our findings. Counterfactual models consider that the causal impact of an event A on an outcome E is a function of the extent to which E depends on A across counterfactual worlds sampled in a certain way. We show that, similarly, people's intuitions as to the causal impact of an event $A \wedge B$ is to a large extent captured by the extent to which E counterfactually depends on it. At the same time, we highlight some ways in which other factors might contribute to the attractiveness of a plural explanation, above and beyond its mere counterfactual-dependence profile. These suggest ways in which current theories of causal selection could be improved to accommodate the new reality of plural causes.

Third, we develop a perspective on the nature of the *representations* involved in subjects' causal selection judgments. Specifically, we propose that humans represent multivariate causes in causal reasoning in a way entirely analogous to how they represent the meanings of pluralities in natural language. The study of natural language reveals that plural entities such as can be expressed by the English noun phrases "Ann and Bill" or "the boys" are not represented as simple conjunctions of the atomic entities that constitute them. Instead, they possess non-classical semantic properties, in particular with respect to how they interact with negation. We provide evidence in the present study that causal judgments for multivariate events are best explained when we suppose that such events are treated as plural entities, with the same mathematical properties as the plural representations that underly our faculty for language. We conclude that theories from formal natural-language semantics offer great promise as mathematically rigorous theories of mental representations in the general-purpose

language of thought.

Causation and causal selection

Humans are adept at representing the world through a web of causal relations between events. Representing causal relations allows people to make sense of what they observe, make predictions about what's to come, and influence the future in some cases (Chater & Oaksford, 2013; Gerstenberg & Tenenbaum, 2017; Pearl & Mackenzie, 2018; Sloman & Lagnado, 2015).

In the psychological literature, people's causal knowledge is usually modeled through formalisms such as Causal Bayes' Nets or Structural Causal Models. These systems represent aspects of the world with variables, causal relations between these variables, and probability distributions (Pearl, 2000). They appear as integral parts of accounts of psychological faculties and functions related to causation, such as causal inference and counterfactual reasoning.

One such causation-related function is causal selection: faced with a complex causal structure, humans will gladly *select* one cause (or, as we will show, more than one) as being more important than others. Moreover, they will assign different scores to different causal variables depending on how they perceive each of those variables as being *the* driver of the observed outcome.

Knowledge of the causal rule in the relevant system is of course one of the main factors determining the explanation humans will favor in causal-selection judgments. The other main driver of causal selection is the *normality* attached to events, a notion that combines the extent to which an event abides by moral or conventional rules, and the extent to which it was expected to happen, before it did happen (Icard et al., 2017; Morris et al., 2019; Quillien & Lucas, 2023).

The relationship between the causal rule that entangles events with the outcome, their normality, and causal-selection judgment can be complex. In a situation where several different variables are each individually *necessary* for an outcome, people tend

to think of the *least expected* variables (the lightning bolt) as *the cause*, and comparatively disregard the importance of the most expected variables (the presence of oxygen), a pattern of judgment known as *abnormal inflation*. The converse tendency is observed in situations where all of the variables considered are each individually *sufficient* for the outcome to occur. In this case, people tend instead to think of the most normal events as the most important causes of the outcome (Icard et al., 2017).

Defining the candidates for causal selection

Causal selection is determined by an amalgam of the system's underlying causal rule and the normality of events. Thus, a standard procedure for formulating theories about participants' causal-selection judgments starts by building a causal model that formalizes their causal knowledge of the system.

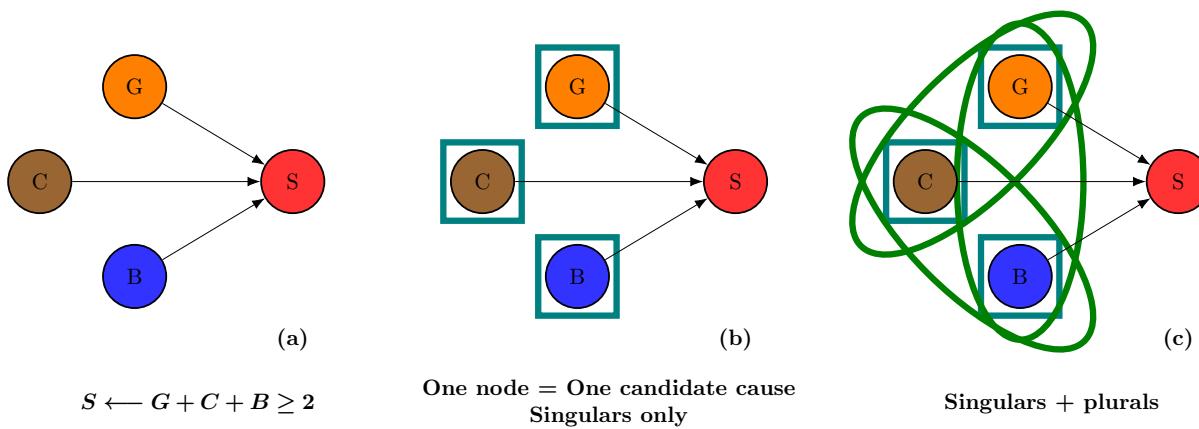
Suppose for example that I get a stomachache shortly after having eaten a piece of Gouda cheese and a plate of pudding containing chocolate cake and blueberry pie. A causal model of this situation would feature one variable for each of the causes of my stomach ache (i.e. one variable each for "eating the Gouda cheese," "eating the chocolate cake," and "eating the blueberry pie") as well as a variable for the effect ("having a stomachache"). The model also specifies a functional relationship between the variables, for example representing the fact that one develops stomach issues after eating too much, as schematized in Figure 1a.

As illustrated by this representational format, it is natural to think of the candidates for causal selection as particular realizations of the individual variables. If an individual equipped with the causal knowledge encapsulated by the model in Figure 1a wonders what *the cause* of their stomachache was, it may seem like they have to make a choice between the variables *G*, *C*, and *B*. This would directly identify the candidates for causal selection as the individual moving parts of the causal model, as represented in Figure 1b.

A striking feature of the psychological literature on causal selection is indeed that causal selection judgments are only ever queried at the level of singular variables

Figure 1

A causal model for the relations between various dishes and my stomachache. (a) I develop a stomachache if and only if I eat two kinds of pudding or more. (b) The standard implicit assumption in the literature is that only single variables are candidates for causal selection. (c) We propose instead that causal judgment can also target plurals, for example pairs of variables.



(Kominsky et al., 2015; Morris et al., 2019; Quillien & Barlev, 2022; Quillien & Lucas, 2023; Sytsma, 2020). Kinney and Lombrozo (2024) deserve an honorable mention in this connection however, since they compared participants' preferences for causal generics ("X causes Y") mentioning one vs several variables. But their work was on type-causation, while here we discuss token (actual) causation.

Concretely, when experimental participants are presented with a situation where an outcome depends on three different events A , B , or C , they are never asked to what extent a *plural* event like $A \wedge B$ can be considered the cause of the outcome. Intuitively though, causal explanations that mention combinations of variables can also be appealing. In our example above, saying that I got a stomachache "because I ate the entire dessert plate" might appear to be a better explanation than either "because I ate the chocolate cake" or "because I ate the blueberry pie" each on its own.

Note that allowing for many variables to feature in causal explanations does not

eliminate the need for causal selection: one might want to mention several causes of an event without mentioning *all* of them. For example, one might think that “because I ate the blueberry pie” is a better explanation for my stomachache than “because I ate the entire dessert plate” if for example I eat chocolate cake at every meal, but add a blueberry pie on top of it only exceptionally. Ultimately, the best candidates for causal selection are those causes that participants see as most *crucial* in bringing about the outcome, whether these be singular or plural, and in principle we can only know what the best causal explanations are after considering the entire set of possible candidates, including plural causes, as illustrated in Figure 1c.

Counterfactual theories

To properly argue the point above, we first need to outline a notion of what it means for a cause to be of a more or less crucial importance in bringing about an outcome. The notion we will rely on throughout this paper is rooted in counterfactual theories of causal selection (Icard et al., 2017; Quillien & Lucas, 2023).

Counterfactual theories of causal cognition in general build on the premise that humans represent causal relationships between variables in terms of counterfactual dependence (Gerstenberg & Tenenbaum, 2017; Halpern & Pearl, 2005; Krasich et al., 2024; Lewis, 1973; Woodward, 2003, 2006). The notion that “*C* caused *E*” is taken to be roughly equivalent to the notion that “had *C* not happened, *E* would not have happened either.”

In the case of causal-selection judgments, this is enriched with the important idea that the counterfactual dependence between *C* and *E* is evaluated not just in the actual world but in other possible worlds as well. Of particular relevance to this evaluation will be the possible worlds that are most *normal*, or *closest to the actual world* in which we are to select a cause (Lewis, 1973). Evaluating counterfactual dependence in these worlds is what allows a causal-selection judgement to provide explanations that are not just relevant to the situation under consideration, but also generalizable to other contexts

(Hitchcock, 2012; Lombrozo, 2010).

We will limit our discussion in this article to two counterfactual theories (and accompanying models) of causal-selection judgment that (1) have been stated in full mathematical rigor and (2) have been submitted to experimental scrutiny, the Necessity and Sufficiency Model (Icard et al., 2017, NSM) and the Counterfactual Effect Size Model (Quillien & Lucas, 2023, CESM). We chose to focus on these two theories because of their good track record in predicting participants causal selection judgments across a wide variety of tasks (Gerstenberg & Icard, 2020; Gill et al., 2022; Henne et al., 2019, 2021; Kirfel et al., 2021; Morris et al., 2019; O'Neill et al., 2022; Quillien & Barlev, 2022).

The two theories see causal selection as a two step process, the first of which identical, the second divergent. The procedure is as follows.

First, randomly sample a large number of counterfactual worlds. The sampling process operates at the level of the individual exogenous variables of the relevant causal model, that is the variables that have no parent in the causal graph. Across worlds, each of these variables is sampled with a frequency that is a function of two elements:

1. Its value in the actual world. Given a causal model with a set of exogenous variables E , and a valuation function $\llbracket \cdot \rrbracket^w$ that maps each variable in E to the value it has in a world of evaluation w , we can define a special world constant $w_@$ designating the actual world, that is the set of circumstances that in fact took place. The model includes a stability parameter s taking a value between 0 and 1, such that in every counterfactual world $w_i \in \{w_1, \dots, w_n\}$ that it samples, each variable in E will have in w_i the same value that it has in the actual world $w_@$, with probability s (see Lucas & Kemp, 2015; Quillien et al., 2023). This parameter is not present in the original version of the Necessity and Sufficiency Model, having been introduced in models of causal selection by Quillien and Lucas (2023). It can however be straightforwardly introduced into it, as we will do in this article.

2. Its prior probability of occurrence. When a variable's value is not directly

mapped to its actual world value, as will happen with probability $1 - s$, it is resampled from its prior probability distribution. This is where an event's sampling propensity (and from there, its causal score) gets to be sensitive to the normality attached to that event.

Consider the example of the causal system presented above relating variables G , C , and B to my stomachache. In the actual world, the variable G that encodes my eating Gouda cheese has value 1, meaning that event actually took place in the present world. As a result, when I sample counterfactuals to the present worlds (as in Figure 2), I will with a probability s automatically represent myself eating cheese also in each of these worlds, as depicted in the left side of the tree in Figure 2a. In the worlds where I don't do that (the right side of the decision tree), the variable will be resampled from the prior probability on that event. Suppose for example that my eating cheese is a rather exceptional occurrence, such that $P(G) = 0.1$. In this case, I am much more likely to travel along the rightmost sub-branch of the tree in Figure 2a, and simulate worlds in which I didn't eat Gouda cheese than if I were accustomed to the fact and ascribed a higher prior probability to the event.

All in all, the sampling propensity (Icard, 2016) of any given exogenous variable V can be reconstructed as a function of the stability parameter s and that variable's prior probability $P(V)$, following the equation below, in the special case of interest here (binary variables, registering whether an event happens or not).

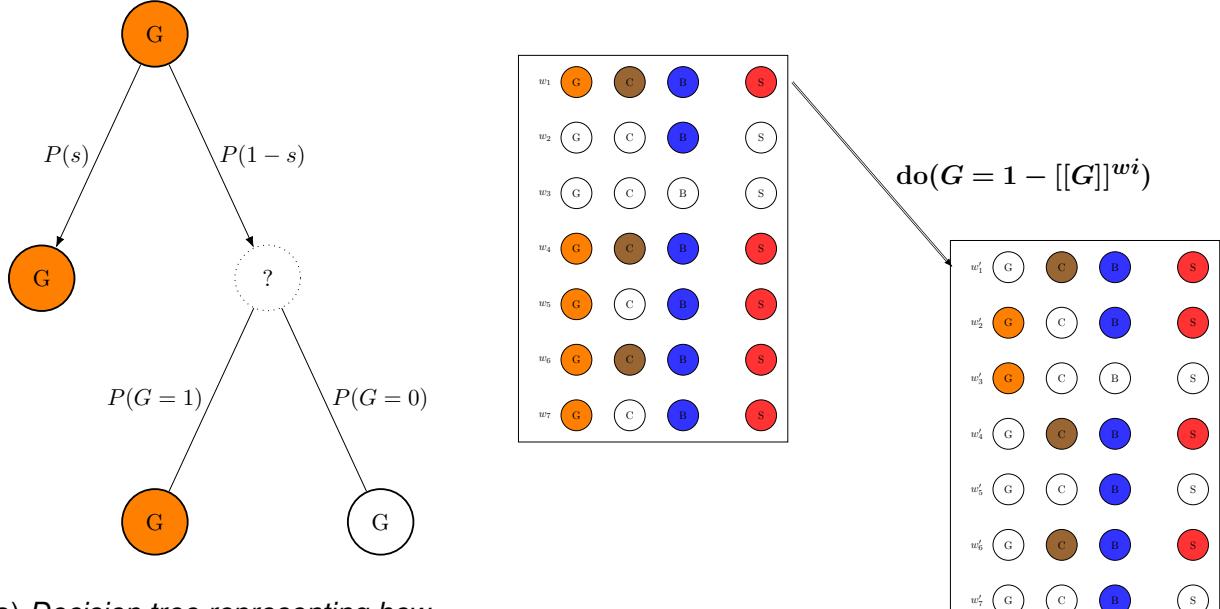
$$SP(V) = s \cdot \llbracket V \rrbracket^w + (1 - s) \cdot P(V)$$

Once the exogenous variables of the system have been sampled in this way, one can simulate the outcome, which follows from the variables via the causal rule underlying this particular system, as in Figure 2b.

Second, compute the causal impact of a given variable V across those counterfactual worlds. The precise way to measure this impact is different in the two theories under consideration. In the NSM, causal impact is scored as the weighted sum of the following two factors.

Figure 2

Sampling counterfactual worlds.



1. A Necessity score: In each world w in which $\llbracket V \rrbracket^w \neq \llbracket V \rrbracket^{w@}$, sample the outcome O from a probability distribution $P^\nu(O)$, where the value of each variable $V_j \in E$ other than V is switched to its actual world value $\llbracket V_j \rrbracket^{w@}$. Then, count one point for the necessity score if the value of the outcome in the resulting world is different from the one that it had in the actual world (i.e. $P^\nu(O) \neq \llbracket O \rrbracket^{w@}$), and zero points otherwise.

2. A Sufficiency score: For every world w in which $\llbracket V \rrbracket^w = \llbracket V \rrbracket^{w@}$, sample the outcome from the probability of Sufficiency $P_{V=\llbracket V \rrbracket^{w@}}^\sigma(O)$ of the event $V = \llbracket V \rrbracket^{w@}$ for the outcome O . There is more than one way to define $P_{V=\llbracket V \rrbracket^{w@}}^\sigma$, but they make extremely similar predictions. The definition that turned out to have the best fit with the data from

our experiments is the following.

$$P_{V=\llbracket V \rrbracket^{w@}}^\sigma(O) = SP(O \mid do(V = \llbracket V \rrbracket^{w@}), \neg\llbracket V \rrbracket^{w@}, \neg O)$$

From here, count one point for the sufficiency score if the value of the outcome thus sampled in w is the same as the value of the outcome in the actual world, and zero points otherwise.

Then, divide the number of points scored this way by the total count of worlds sampled. The dynamics of necessity and sufficiency scoring make it such that the necessity score is all the more important when the prior probability of an event is low, making it more likely to switch value across counterfactuals, whereas the sufficiency score is all the more important when this prior probability is high.

In the CESM, causal impact is computed using the same process in every world w_i , as follows.

1. Switch the value v of the variable V to a new, randomly sampled value v' . Then reevaluate the outcome in the new world w'_i where the value was switched. A representation of this resampling process is given in Figure 2b. The impact $K(V \rightarrow O_i, w_i)$ of V in the world w_i is then evaluated as

$$K(V \rightarrow O_i, w_i) = \frac{\Delta O}{\Delta V} = \frac{\llbracket O \rrbracket^{w_i} - \llbracket O \rrbracket^{w'_i}}{\llbracket V \rrbracket^{w_i} - \llbracket V \rrbracket^{w'_i}}.$$

This equation can be glossed as follows. Whenever the outcome is switched in the same direction as the target variable (both from 1 to 0, or both from 0 to 1), V scores a point; when the outcome is unaffected, it scores none. When it moves in the opposite direction, it scores a negative point.

2. The causal impact is then normalized by the ratio of the standard deviations $\frac{\sigma_O}{\sigma_V}$, and averaged across worlds to get the causal impact score $K(V \rightarrow O_i, 0, \llbracket \cdot \rrbracket^{w@})$ of the target variable for the target outcome in the actual world. In simple conditions like the ones we will deal with in this paper, the causal impact of V can be equated to the correlation coefficient between V and O across counterfactuals sampled at the first step.

Plurals in natural language

We hypothesize that humans represent and manipulate multivariate causes in causal selection in a way analogous to how they represent and manipulate the interpretations of plural noun phrases in natural language. Consider the following sentences.

- (1) Who lifted this piano?
- a. It was Ann.
 - b. It was Ann and Bill.

In response to a question as in (1), the answer in (1b) is just as natural as that in (1a): both sentences contain a perfectly coherent entity identified as the piano lifter, in (1a) a singular individual, in (1b) a plurality, or a collective. Accordingly, theories of plurality from linguistic semantics *generalize to the worst case*, taking the singular as the special case of the plural when cardinality is one, allowing for a unified account of singular and plural predication (Link, 1983). Notice also that, in a situation where (1b) is the complete answer, a speaker might still accept the (partial) truth of a sentence like (1a), by virtue of the fact that Ann participated in the lifting event, especially if Ann's role was particularly important.

Analogously, in a causal system we can answer questions like "What caused my stomachache?" by pointing to individual variables ("It was the chocolate cake") or to multiple variables in one go ("It was the chocolate cake and the blueberry pie"). Just as in the language example, the null hypothesis is to expect the exact same mechanism to handle the singular and the plural case, since both the singular and the plural are perfectly coherent entities of the same type: causal entities which can be the target of causal selection.

And again just like the language case, in a situation where the multivariate cause

“chocolate cake and blueberry pie” is the real culprit, one may still be inclined to (partially) accept the singular cause “chocolate cake,” on the grounds that the chocolate cake participated in engendering my stomachache, especially if the role of the chocolate cake in bringing about this outcome was particularly important. We submit that this is why, although causal selection is at its heart sensitive to plural causes, experimental methods that haven’t countenanced this possibly have still been very informative, just like judgments involving natural-language singular entities can be coherent and systematic and therefore informative, despite the fact that the fuller story involves plural entities.

Now, at first glance, in language in general as in causation, one would be tempted to propose that a plural entity should be understood as the mere *conjunction* of its individual constituents. That is, in the causal case, to claim that the plural variable “chocolate cake and blueberry pie” is to blame for my stomachache is actually mere shorthand for saying that the singular variable “chocolate cake” is to blame for my stomachache *and* the singular variable “blueberry pie” is to blame for my stomachache. The linguistic example in (1) already suggests that this might be missing something: if in actual fact it was Ann and Bill *together* that lifted the piano, then the plural in (1b) is the only precise way to describe the situation, a conjunction of the shape “Ann lifted this piano and Bill lifted this piano” would be quite inaccurate. Using language then as a source of conceptual possibilities with testable consequences, we might expect that plural causes should also display *togetherness* effects of this kind. The behavioral experiments we report on shortly will address this possibility head-on.

But plurals in natural language differ even more sharply from mere conjunction than by displaying this kind of holistic effect. Of special relevance to what’s to come is the

fact that the negation of a plural does *not* correspond to the negation of a conjunction.

(2) Ann and Bill don't speak German.

- a. Either Ann doesn't speak German, or Bill doesn't, or neither does.
- b. Neither Ann nor Bill speaks German.

The most natural reading of (2) is as in (2b), and not as in (2a) as the standard Boolean semantics for “and” and “not” might lead us to expect (Krifka, 1996; Lappin, 1989; Löbner, 2000; Szabolcsi & Haddican, 2004). Thus, if multivariate causes are *plurals* in any interesting sense, we must consider the possibility that they behave like linguistic plurals, suggesting perhaps a non-flat weighting of the three in-principle ways of falsifying a conjunction, strongly favoring the possibility where all variables constitutive of the plural are individually negated.

Naturally, this last possibility will be particularly relevant in situations where subjects are tasked with formulating judgments about *negative* outcomes. That is, for example, when they have to explain *losing* a round of a game, in a situation where they have only explicitly been instructed in the conditions for *winning*. For producing such judgments will require them to internally use some equivalent of negation on their mental representation of the winning conditions. Our Experiment 2 will address this question, among others.

Lastly, a word is warranted on the theories of plurality from linguistics which describe and partly explain these facts. Since Link's (1983) seminal work on the logic of plurality, the consensus view has been that plural entities require their own dedicated mathematical structures for representation, and cannot be subsumed as special cases of Boolean conjunction or Boolean disjunction. Specifically, plural entities are formed out of singular entities (“atoms” in the formal semantics jargon) by a *mereological sum* operation. Mathematically, this is a join operation, giving rise to a join semilattice. Mixing and matching elements of the algebraic and order-theoretic characterizations of join

semilattices for the sake of clarity and brevity, these are structures $\mathfrak{A} = \langle A, \oplus \rangle$, where A is the set of singular and plural entities and \oplus is the plural-forming (join) operator, such that for any $a, b \in A$, we will also find $a \oplus b \in A$, representing the smallest plural containing both a and b . \mathfrak{A} also comes with a partial order \leq over A which formalizes the notion of containment, so that, for example $a < a \oplus b$. From this point onward the details differ in various accounts, but all share this algebraic structure, crucially distinct from the algebra induced by the more familiar Boolean connectives. This structure provides the mathematical degrees of freedom required to now define both *cumulative* and *distributive* predication (Ann and Bill lifting a piano together vs. Ann lifting a piano and Bill lifting a piano), and to describe the observed *homogeneous* interaction with negation (Krifka, 1996; Križ & Spector, 2021; Löbner, 2000).

This is all the detail we can allow ourselves on this topic, short of reviewing an extensive and technically involved literature from formal semantics. The readers whose curiosity we've managed to whet will find a concise but mathematically rigorous introduction to these tools and their applications in an excellent handbook chapter by Champollion and Krifka (2016). Our goal with this brief illustration of the algebraic structures universally used in plural semantics is only to give a taste of what it might possibly mean to define a non-classical, non-Boolean, conjunction-like operation for plurality which is nevertheless entirely mathematically rigorous and intelligible. We will address why and how this kind of mathematical clarity about natural-language meaning matters in psychology in our general discussion.

Extending causal selection to plurals

The idea that people represent plural events produces testable predictions. For starters, an alternative view on such plural causal judgments could contend that people only ever have direct intuitions about the causal responsibility of the *individual* variables in their causal model. People might make judgments about a plural cause say by adding up or averaging the individual causal strengths of its constituent variables. For example,

to compute how much they agree that “eating the chocolate cake and the blueberry pie caused the stomachache,” people might first compute their agreement with “eating the chocolate cake caused the stomachache,” “eating the blueberry pie caused the stomachache,” and so on. Then they might somehow aggregate the causal strength of each individual variable. We will call this the *linear combination* hypothesis: a systematic simple combination of the scores for singular variables might underwrite participants’ judgments about plural causes. This hypothesis is deflationary with respect to the psychological reality of plural causes in that it holds that people can make plural-cause judgments when prompted to do so, but they cobble them together from more primitive representations of causal strength at the level of individual variables. This means that the cognitive process underlying causal selection is ultimately still only ever deployed at the level of singular causes. In keeping with our language analogy, this would mean that the holistic, *togetherness* effects we find with linguistic plurals do not exist, or are trivial, in the domain of plural causes.

In contrast, we consider the possibility that plural causal judgments are the output of a *holistic* computation. Under this possibility, the same cognitive process that allows people to formulate causal judgments about singular variables is deployed at the level of combinations of variables, yielding quantities that can on occasion diverge significantly and non-linearly from the causal judgments for constituent variables. We provide more details about how such a hypothesis could be implemented in models of causal selection in the next section below. But the general idea can be explained rather simply: to assess the causal importance of a plural event $A \wedge B$, which happens in every world where both A and B happen, is to look at the causal impact that this event has on the outcome of interest, using the same measures of causal impact that were detailed above for singular variables. This holistic computation will sometimes lead to different predictions than the hypothesis that consists in simply computing the impact of A , of B , and then combining them. For example, for a plural event $A \wedge B$ to be the most highly rated among possible

candidate causes, the singular events A and B need not necessarily be the most highly rated among singular causes.

In order to tease apart these two hypotheses in the case of causal selection, we need to identify contexts where they make different predictions about causal-selection judgments. This is what we do in the two experiments presented in this paper.

Experiment 1

Our first experiment has the following goals.² First, if plural causes are processed as genuine causes by the mind, factors that are known to affect causal-selection judgments should influence judgments about a plural cause. In particular, the probability of an event is known to affect judgments about whether that event caused an outcome (Morris et al., 2019). We expect analogous patterns for plurals: varying the probability of events should affect causal judgments about whether a conjunction of these events caused the outcome. Second, we aim to rule out a deflationary *linear combination* account of the impact of probability on participants' judgments. Evidence of non-linearity in causal judgments would constitute stronger evidence for the psychological reality of plural causes in human causal selection. We design a situation where both the CESM and the NSM predict that the causal strength of plural variables will not be a linear combination of the score of individual variables. We compare their predictions to those of a null model that tries to predict the score of plural causes as a linear combination of the scores of individual variables.

Methods

Design and materials

We adapted a paradigm developed by Quillien and Lucas (2023). Participants made judgments about a game of chance, in which one randomly draws balls from a set

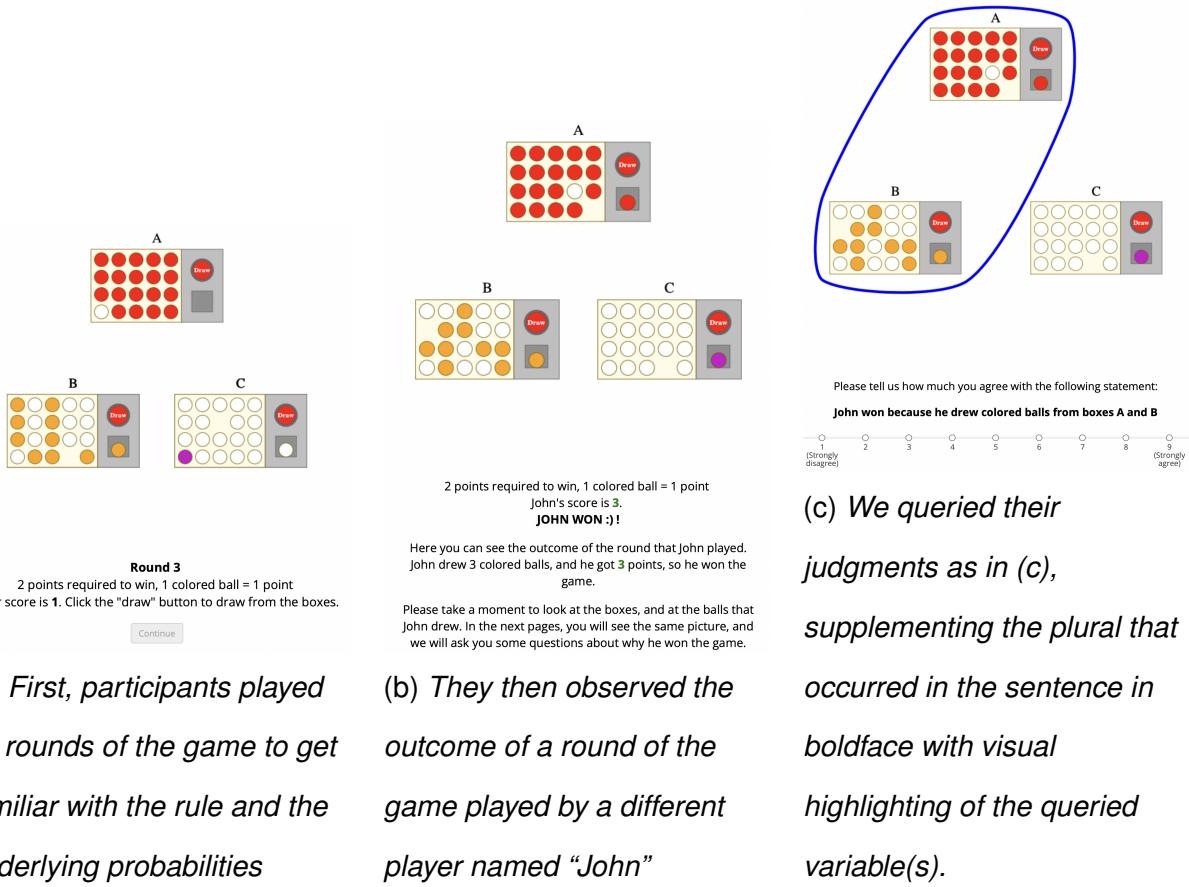
² This study was reported at the Forty-fifth Annual Meeting of the Cognitive Science Society and published in the society's non-archival conference proceedings (Konuk et al., 2023). Our writing in this section borrows directly from this preliminary report.

of urns, and wins by getting enough colored balls (see Figure 3 for illustrations). Participants observed a fictitious player draw a colored ball from each of three urns (labelled *A*, *B*, and *C*) and win the game as a result. Then they were asked to make a causal judgment about each singular cause (e.g. whether getting a colored ball from urn *A* caused the player to win the game), and about each pair of causes (e.g. whether getting a colored ball from urns *A* and *B* caused the player to win the game). For exploratory purposes, we also asked participants to make a causal judgment about the triplet (getting a colored ball from *A*, *B*, and *C*). We manipulated the prior probability of each outcome within participants by varying the proportion of colored balls in each urn, with probabilities of 0.05, 0.5, and 0.95 (Figure 3). We will refer to the three different urns as the LOW, INTERMEDIATE, and HIGH urns, respectively. The rule of the game, which was directly revealed to the participant at the outset, was that the player wins if they get two colored balls or more. This corresponds to the causal model below.

$$\text{WIN} := A + B + C \geq 2$$

Predictions

This paradigm provides a context where the linear and the holistic extensions of the models we outlined above make clearly different predictions. The CESM predicts that participants' singular causal-strength estimates should follow a particular ranking: INTERMEDIATE > LOW > HIGH, for any value of the *s* parameter. This is because, across possible counterfactual alternatives to what happened, there is a high correlation between getting a colored ball from the intermediate probability urn and winning the game. These predictions partially match participants' judgments collected in the previous iteration of this paradigm run by Quillien and Lucas (2023), where judgments were collected for singular variables only, and in which participants' responses followed the ranking: INTERMEDIATE probability urn > LOW = HIGH (there was no significant difference between the ratings for the low probability and the high probability urn).

Figure 3*The three phases of Experiment 1*

If one considers a linear extension of the CESM, where participants simply combine the causal strength of individual variables to make plural cause judgments, or simply a linear combination of participants own recorded judgment, they should consider that the pair **LOW & INTERMEDIATE** should have greater-than-or-equal causal strength to the pair **HIGH & INTERMEDIATE**, because the singular **LOW** has higher causal strength than **HIGH**.

In contrast, if participants judge the causal strength of plurals via a holistic computation, they should rate the pair **INTERMEDIATE & HIGH** as highest. For, across possible counterfactuals, there is a high correlation between getting a colored ball from these two urns and winning the game. Intuitively, since drawing a colored ball from the

low-probability urn is rare, and given that at least two balls are needed to win, most worlds where the player wins the game will be worlds in which they do so by getting a colored ball from the INTERMEDIATE and HIGH urns. This prediction is true for any value of the s parameter in the holistic version of the CESM. It is also true for the holistic version of the NSM, although in that case it does reverse the ranking that the NSM expects for singulars.

Procedure

Participants first completed ten rounds of the game, presented with urns as in Figure 3a. We pseudo-randomized the draws in such a way as to get participants to internalize the probabilities associated with each urn and how they connected to the outcome. Then participants saw the outcome of a round of the game played by another fictitious player, who drew a colored ball from all three urns, thereby winning with 3 points, as in Figure 3b. They were asked to rate the causal strength of each individual draw, as well as that of every combination of two or three draws for the winning outcome, on a Likert scale from 1 to 9 (strongly disagree to strongly agree), as in Figure 3c. For the singulars, participants were asked to rate their agreement with the statement “John won because he drew a colored ball from box [urn].” For the plurals, they rated their agreement with “John won because he drew colored balls from boxes [urn1] and [urn2].” Each question was displayed on a separate page, next to the urns that displayed the outcome of the fictitious player’s draw. The letters indexing the urns, as well as the colors of the balls, were randomized across participants to avoid confounding but were kept the same within a participant. Half of the participants were asked about the singulars first, and then about the pairs. The other half were asked about the pairs first, and then about the singulars. All participants were asked about the triplet at the very end. Within one class of questions (singulars vs. plurals) we randomized the order of presentation of questions. Finally, participants completed a brief demographic questionnaire and were redirected to Prolific for payment. We coded the experiment in the jsPsych library

(De Leeuw, 2015), with custom plugins for displaying urns developed in our lab.

Participants

We recruited 400 participants from all English-speaking countries from Prolific.

This sample size was inspired by the one used by Quillien et al. (2023), who used a comparable sample size (290 participants), for a study with similar design. We excluded from subsequent analysis 44 participants who failed to answer either of two elementary comprehension questions that checked their understanding of the rules of the game, leaving a total of 356 participants for analysis.

Transparency and openness

We report how we determined our sample size, all data exclusions (if any), all manipulations, and all measures in the study. All data, analysis code, and research materials are available at

https://osf.io/43m5d/?view_only=3d26a80b8e394fa9ad7792a690de8fe6. Data were analyzed using R (R Core Team, 2013), version 4.3.3 and the package ggplot2, version 3.5.0. All studies we report in this article received ethics approval by the *Comité d'évaluation de l'éthique de l'INSERM*, under research protocol *Le langage et les capacités cognitives connexes*. All studies were conducted entirely in English.

Results

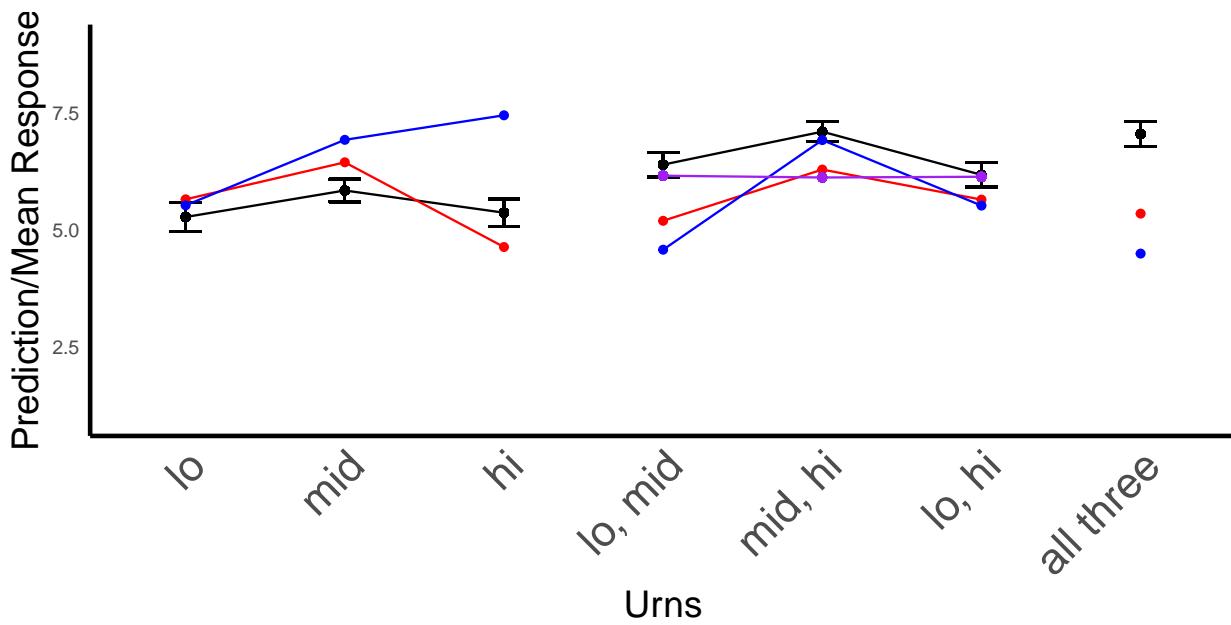
We first report analyses using standard statistical tests. Then we report the fit of computational models of causal judgment.

Basic results patterns

Prior probability affects both singular and plural causal judgments. Results are plotted in Figure 4. We ran two two-factor repeated-measure ANOVAs, one for each main type of cause queried (*singulaires* and *pairs*), using urn probabilities and order of presentation as predictor variables, and participants' responses as the dependent variable. Results are in Tables 1 and 2. There was a main effect of prior probability on

Figure 4

Mean ratings by question type, along with predictions for each theory under consideration. The error bars represent the standard error of the mean. The linear combination theory (purple) predicts that the score of the LOW & INTERMEDIATE and INTERMEDIATE & HIGH pairs should be equivalent, when in fact we see a significant difference between them, in accordance with the holistic versions of the two counterfactual models: the Counterfactual Effect Size Model (Quillien & Lucas, 2023) (in Red on the plot) and the Necessity and Sufficiency Model (Icard et al., 2017) (in Blue).



participants' causal judgments, for singular as well as for plural causes ($p < 0.02$ in both cases), consistent with our expectation that participants' judgments for plural causes should be sensitive to probabilities just like for other actual cause judgments.

The order of presentation/querying singular and plural selection judgments also had a significant effect ($p < 0.001$) on the ratings for singulars: singular causal judgments were lower when presented after the plurals. There was however no interaction effect between urn probability and order of presentation, suggesting that the impact of probability on causal estimates did not vary depending on the order in which

Table 1

ANOVA for singular causal-selection judgments, predicting urn ratings from urn probabilities and urn order of presentation.

Factors	Mean Sq	F-score	p-value
Probabilities of the urns	32.98	4.492	< 0.012
Order of presentation	138.81	18.904	< 0.0001
Probabilities:order	3.18	0.433	> 0.64

questions were asked. Therefore we drop this variable (order of presentation) from later analyses.

The causal strength of plural causes is not a linear combination of the causal strength of individual variables. The pattern of responses for singular variables replicated the patterns obtained by Quillien and Lucas (2023). Judgments for the INTERMEDIATE urn were higher than judgments for the LOW urn, $t(315.41) = -2.70$, $p = 0.007$, and the HIGH urn, $t(325.85) = -2.08$, $p = 0.038$. The difference between the LOW and HIGH urns was not significant, $t(350.59) = -0.63$, $p = 0.53$.

We can use these results to test the *linear combination* hypothesis, according to which participants derive their plural-cause strength estimates by adding up or averaging their estimates for the individual variables that compose a given plural. If this were correct, participants should give the same causal strength estimate for the two plural causes LOW & INTERMEDIATE and INTERMEDIATE & HIGH, since their estimates for the singular causes LOW and HIGH are not significantly different from each other. By contrast, both the holistic CESM and the holistic NSM predict a sharp difference between these two kinds of plurals, with the plural cause INTERMEDIATE & HIGH being rated higher (Figure 4).

Consistent with the holistic CESM, judgments about the INTERMEDIATE & HIGH pair were higher than for the LOW & INTERMEDIATE pair, $t(355) = -4.67$, $p < 0.001$, and

Table 2

ANOVA for pair causal-selection judgments.

Factors	Mean Sq	F-score	p-value
Probabilities of the urns	83.02	14.646	< 0.000001
Order of presentation	21.75	3.837	> 0.05
Probabilities:order	2.30	0.406	> 0.66

higher than for the LOW & HIGH pair, $t(355) = 6.858$, $p < 0.001$. In slight deviation from the CESM and NSM's predictions however, judgments for the *low, intermediate* pair were higher than for the *lo, high* pair, $t(355) = 2.3691$, $p = 0.02$ (Figure 4).

We conducted two more analyses to rule out the linear combination model. First, we ran a one-way repeated-measure ANOVA, predicting judgments for the pairs (LOW & INTERMEDIATE, INTERMEDIATE & HIGH, and LOW & HIGH) from judgments for the singulars (LOW, INTERMEDIATE, and HIGH), as well as their interactions, as within-participant factors. Each plural pair was regressed only on the values of the two singulars that comprised it.

The linear combination theory predicts that there should be no significant interaction: a participant's causal judgment for a given *singular* variable should have the same impact on every plural cause in which it features. One's estimate for the singular INTERMEDIATE, for example, should have an equal impact on one's estimate for INTERMEDIATE & HIGH and for LOW & INTERMEDIATE.

We find evidence against the linear combination theory (Table 3). There was a significant interaction between the INTERMEDIATE and HIGH urns, $p = 0.004$. In addition, the main effects of the singular judgments were not significant, for all but the LOW urn.

Second, we fitted linear multilevel regression models on participants' responses for pairs. Specifically, we compared the predictive performance of two different models on participants' plural-cause estimates. The first one used as predictor the average of

Table 3

Results of the ANOVA: estimate for pairs ~ est. for singular-1 × est. for singular-2.

Factors	Mean Sq	F-value	p-value
LOW	100.69	17.743	< 0.00001
INTERMEDIATE	20.80	3.664	0.05586
HIGH	15.93	2.808	0.094
LOW:INTERMEDIATE	6.64	1.169	0.2798
LOW:HIGH	0.41	0.073	0.7872
INTERMEDIATE:HIGH	46.64	8.219	0.00423

the two singular-cause estimates for the variables contained in a given plural (computed on a per-participant basis), plus a random intercept. The second model also included the question asked (that is, the specific plural being queried) as predictor. A likelihood ratio test shows that adding question as a predictor significantly improves the fit of the model, $\chi^2(5) = 27.53$, $p < 0.001$ (Table 4). Again, this is inconsistent with the linear combination account.

Computational modeling

We computed the predictions of two recent counterfactual models of causal selection, the Counterfactual Effect Size Model (Quillien & Lucas, 2023) and the Necessity and Sufficiency Model (Icard et al., 2017), presented in the introduction. Our implementation follows the one given in Quillien and Lucas (2023).

For each question we report on below, we generated causal judgments for the CESM using a process of counterfactual sampling. We generated predictions for the CESM by simulating 10^5 possible rounds of the game according to the rule, what was the case in the situation described to participants, and the sampling model described by the CESM. We computed CESM judgments for an event as the correlation between that

Table 4

Comparison between two models: the linear combination model of plurals (means of singulars + intercept), and the means of singulars + question model.

Models	LogLik	Df	χ^2	p-value	BIC
Means sing	-891.89	3			1804.709
+ Question	-878.13	5	27.53	< 0.00001	1791.126

event (for instance, whether the player draws a colored ball from urn *A*) and the outcome of the game (whether the player wins the game), across simulations. We computed NSM judgments analytically, as the sum of the variables' sufficiency and necessity scores across worlds.

We fit the value of the stability parameter s for both models by finding the value of s that results in the best fit between model judgments and average participant judgments across all seven questions. We quantified model fit by looking at the likelihood of mean answers per question under a normal distribution centered on the model's predictions, with a standard deviation fitted across questions.

We identified the best fit value via a grid search, exploring a wide range of values for the parameter s , crossed with different values for a scaling parameter γ (applied to a model's predictions as an exponent prediction^γ). The point of γ was to avoid situations where one model would systematically overshoot or undershoot actual participant answers, as the models are not meant to predict the exact value of participants' judgments, but only the relative difference between one variable and the next. Our technique here was analogous to that of Griffiths and Tenenbaum (2005).

For the CESM, the best fitting value was $s = 0.89$, with $\gamma = 0.26$. For the NSM, the best fitting value was $s = 0.71$, with $\gamma = 2.93$.

In our implementation, to assess the causal strength of plural causes a model assumes that people compute the causal strength of the conjunction of all variables

Table 5

Table of model comparison, Study 1, excluding the triple. The AIC and BIC values are computed for mixed effects models, including group and a random effect for participants.

Model	AIC	BIC	Cor.
CESM	9929.31	9957.65	0.12
NSM	9962.06	9990.4	0.06
Considering only the pairs			
CESM	4692.11	4716.978	0.12
NSM	4674.355	4699.222	0.16
Empirical average	4692.942	4717.81	-0.11

contained within that plural. For instance, the CESM computes the causal strength of LOW & HIGH by computing the correlation between the binary variable LOW \wedge HIGH (which has value 1 if both LOW and HIGH have value 1, and 0 otherwise) and the outcome.

The predictions of the models are plotted in Figure 4. Table 5 details the comparison. Overall, the CESM's prediction had the best fit to human judgments in this experiment, although the NSM had the best fit when models were compared on pairs of variables only. We also compared the models' performance on the judgments for pairs to a null model that used as predictor for each pair the average of mean human judgments for each singular variable contained within a given pair, as plotted in Figure 4. Both counterfactual models proved significantly better than this linear predictor (Table 5).

Discussion

We find evidence that, when people make a judgment about whether events *A* and *B* caused an outcome, their judgments track the correlation between the conjunction

of *A* and *B* and the outcome, across counterfactuals. Concretely, in our experiment, winning the game is in general strongly associated with getting a ball from both the intermediate- and high-probability urns, and people judged that combination of events to be highly causal. Importantly, this effect is inconsistent with a simpler account, according to which people's judgments about plurals are cobbled together from their causal intuitions about each individual variable in the plural.

Judgments about plural causes are affected by the prior probability of their constituent variables, but cannot be derived from the causal strength of these individual variables. As such, our results are in general consistent with the predictions of simple extensions of recent counterfactual models of causal selection (Icard et al., 2017; Quillien, 2020; Quillien & Lucas, 2023), augmented with the assumption that people judge plural causes in a holistic manner.

At the same time, these findings raise new questions about the psychology of causation. Presently we highlight two of these questions, which we investigate in Experiment 2.

First, participants in this study found the plurals overall more appealing than the singulars, a tendency which the counterfactual models we considered did not capture. Participants might have felt that plurals provided more exhaustive descriptions of the event: they give more complete information about what happened, in addition to why it happened. We also find that this effect is accentuated when singulars are presented after plurals. Making judgments about plurals first might highlight to participants the descriptive incompleteness of singulars. This finding suggests an interesting tension between two potential desiderata of causal judgment: highlighting the variables that were most causally important to the outcome, and providing an exhaustive list of the causal factors. If so, this calls for an investigation into the relative importance of these two pressures in participants' causal-selection judgments. When, for example, adding a variable weakens the counterfactual dependence profile of the resulting plural, such as

when the plural doesn't explain the outcome appreciably better than one of the singular variables within it, will participants still show a preference for plurals, on account of their greater completeness?

Second, a notable feature of this first experiment is that the causal structure used a simple additive rule (i.e. the player wins the game if their score is above a certain threshold). As such, there is a sense in which the variables each have an independent incremental causal effect on the outcome.

What will participants' plural-cause judgments look like in a causal structure where some conjunctions of events directly feature as such in the causal model that generates the outcome? Consider for example the causal rule $(A \wedge B) \vee C$. Here the urns A and B are specifically connected in the logical structure. Generalizing somewhat, our question here is: when an outcome specifically depends on the joint occurrence of A and B , should that make the plural cause $A \wedge B$ a more natural causal explanation than a potential alternative $A \wedge C$, even if C also makes an important contribution to the outcome?

Experiment 2

Experiment 1 established the psychological reality and relevance of plural causes for causal selection judgments. Building on its findings, Experiment 2 expands our exploration of plural causes in four directions.

First, we provide additional evidence against deflationary interpretations of plural causes. We give more examples of situations in which people's plural-cause judgments cannot be straightforwardly derived from a linear combination of their singular-cause judgments, to confirm the results obtained in the first experiment.

Second, we explore whether there is a robust bias toward preferring causes that contain more variables. In Experiment 1, participants gave overall stronger scores to plural causes than singular ones. Experiment 2 investigates whether this pattern always holds.

Third, we explore a richer causal structure. Here, two urns contain purple balls, and two urns contain orange balls. The player can win the game by getting either two purple balls or two orange balls, where “or” is meant inclusively. Formally, winning can be triggered by either of two distinct sufficient conditions $A \wedge B$ and $C \wedge D$, each a conjunction of two variables. This corresponds to the rule

$$\text{WIN} := (A \wedge B) \vee (C \wedge D) \quad (1)$$

Fourth, we explore participants’ judgments in situations where they have to explain a *negative* outcome. In the context of our experiment, this amounts to explaining a *loss*, in a situation where participants were instructed in the conditions for *winning*. When the task is to explain a loss, participants presumably must first form a representation of the losing conditions, before they can compute the effect of each event on losing. As we detail presently, this requirement opens interesting possibilities about the representations participants might deploy.

Finally, Experiment 2 was also designed to collect many more data points per participant, increasing our statistical power compared to Experiment 1. We ask each participant about the outcome of four possible rounds of the game (as opposed to just one outcome in Experiment 1), collecting a total of 36 causal judgments per participant.

Plural negation in loss contexts

From a logical standpoint, it seems natural to assume that the conditions for losing the game can be read off the *classical logical negation* of the conditions for winning the game, represented in equation 1 above, as follows.

$$\begin{aligned} \text{LOSS} &:= \neg((A \wedge B) \vee (C \wedge D)) \\ &\equiv \neg(A \wedge B) \wedge \neg(C \wedge D) \\ &\equiv (\neg A \vee \neg B) \wedge (\neg C \vee \neg D) \\ &\equiv (\neg A \wedge \neg C) \vee (\neg A \wedge \neg D) \\ &\quad \vee (\neg B \wedge \neg C) \vee (\neg B \wedge \neg D) \end{aligned} \quad (2)$$

While arguably the most appropriate strategy from a normative standpoint, this is not the only plausible option. Indeed, continuing with our theme of taking inspiration from the semantics of natural-language plurality for generating hypotheses, it is possible that participants consider a much stronger set of losing conditions than are presented in equation 2.

Plural entities in natural language have the logically surprising feature that negation applies homogeneously to each individual in the plurality. Consider the examples of plurals in (3) and (4), and the putative interpretations for the negated plural (4) in (4a) and (4b).

(3) The boys did their homework.

(4) The boys didn't do their homework.

- a. None of the boys did their homework.
- b. At least one of the boys didn't do his homework.

Sentence (3) means that *every* boy did his homework, with some tolerance for exceptions which needn't concern us here (Križ & Spector, 2021). Sentence (4) then ought to be simply the negation of (3), which would amount to the interpretation paraphrased in (4b). Yet, the negated plural in (4) has a much stronger interpretation, to the effect paraphrased in (4a). In general, negated plurals are interpreted in this unexpected way, from the standpoint of classical logic (Krifka, 1996; Lappin, 1989; Löbner, 2000). This observation applies to plurals as in (4), generated by a noun phrase with plural morphology “the boys,” but also to plurals formed by means of an explicit conjunction: a sentence like “John and Mary didn't countenance this hypothesis” means that neither John nor Mary considered this hypothesis, not merely that at least one of John or Mary failed to consider it (but see Szabolcsi & Haddican, 2004, for evidence of cross-linguistic variation on the available interpretations).

In light of these observations, we hypothesize that participants in our experiment might negate the rule in equation 1, which is a disjunction of plural terms, in this non-standard way. This would amount to the stronger loss conditions at the end of equation 3 below, which we preface with ‘ $\not\equiv$ ’ to indicate that it violates classical-logical equivalence.

$$\begin{aligned} \text{LOSS} &:= \neg((A \wedge B) \vee (C \wedge D)) \\ &\equiv \neg(A \wedge B) \wedge \neg(C \wedge D) \\ &\not\equiv \neg A \wedge \neg B \wedge \neg C \wedge \neg D \end{aligned} \tag{3}$$

We refrain for now from any discussion of the possible *reasons* and *mechanisms* whereby participants might engage in this language-like treatment of negated plural causes, at this point we mean only to point out that this is a plausible hypothesis worthy of testing. We will address the theory questions in the general discussion.

Methods

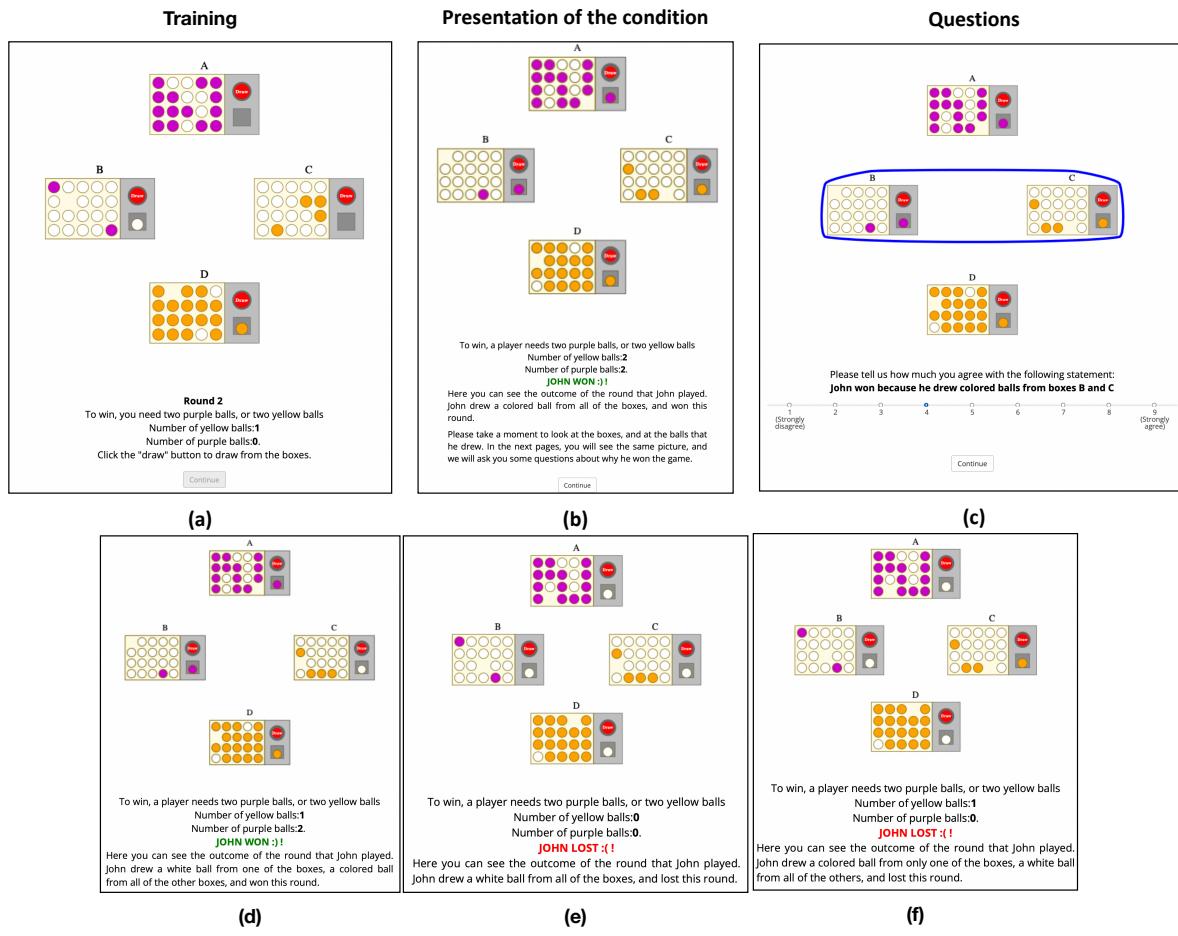
Design and materials

The methodology was similar to that of Experiment 1. We presented participants with a simple game of chance. This time, the game involved four urns, with two different colors, purple and yellow (Figure 5). We randomized the assignment of colors, but always in such a way that urns *A* and *B* were of one color, and urns *C* and *D* of the other color. To win a round of the game, one needed to draw “two purple balls or two yellow balls.”

While we randomized the specific urns’ indices and their spatial arrangement for each participant, for simplicity here we refer to a consistent arrangement as depicted in Figure 5, where urn *A* has 14 colored balls, urn *B* 2, urn *C* 4, and urn *D* 19 colored balls. These induce different prior probabilities of drawing a colored ball out of each urn, such that, respectively, $P(A) = 0.7$, $P(B) = 0.1$, $P(C) = 0.2$, and $P(D) = 0.9$. Throughout the experiment, the urn containing 14 colored balls and the urn containing 2 colored balls

Figure 5

Top: The three phases of Experiment 2, analogous to those of Experiment 1, by order of presentation to the participants. The condition presented here as example is the OVERDETERMINED NEGATIVE condition. Bottom: the three other conditions presented to participants: (d) TRIPLE-1 ; (e) OVERDETERMINED NEGATIVE ; (f) TRIPLE-0



were always of the same color, while the other two urns (19 and 4 colored balls) were of the other color, so that each color would contain one high-probability and one low-probability urn.

Procedure and participants

As in Experiment 1, participants first had the opportunity to familiarize themselves with the game and the rule determining a winning outcome, as well as with the

underlying probabilities, by playing the game for ten rounds, as in Figure 5a. Urn draws and outcomes at this stage were pseudo-randomized in such a way as to reflect the underlying probabilities.

After they played ten rounds of the game, they saw the outcomes of rounds played by another player named ‘John’ (as in Figure 5b) and were asked to rate on a Likert scale from 1 to 9 the causality of certain events, both singular and plural. Specifically, we queried their causal judgments by asking them the extent to which they agreed (on a 1–9 scale) with a sentence that followed the template: “John won (/lost) because he drew colored (/white) balls from box(es) [XYZ].” Figure 5c shows an example.

Note that although participants judgments were queried on a 1 to 9 scale, they were recorded in our data on a 0–8 scale; hence the figures below plot participants answers on that latter scale. All participants saw four different rounds of the game played by John, one at a time, and provided their judgments after each round. All the rounds were played with the same underlying rule and the same urns in the same display as the one with which participants had previously been familiarized with. Each trial differed only in the outcome of the draw made by John. We presented all participants with the four following rounds, in random order:

1. OVERDETERMINED POSITIVE: John drew a colored ball from each of the four urns — John won (Figure 5b)
2. TRIPLE-1: John drew a colored ball from urns *A*, *B*, and *D*, but not from urn *C* — John won (Figure 5d)
3. TRIPLE-0: John drew a white ball from urns *A*, *B*, and *D*, but not from urn *C* — John lost (Figure 5e)
4. OVERDETERMINED NEGATIVE: John drew a white ball from all four urns — John lost (Figure 5f).

Within each round, we asked participants about every singular event that featured a colored ball in the winning rounds, and every singular event that featured a white ball in

the losing rounds. We also asked about every plural combination of these singulars, with the exception of four-variable plurals (we considered those questionable candidates for causal selection judgments, since they provided an exhaustive description of all drawing events in a given round) and other plurals which we considered redundant with some that we already asked. The questions were presented in random order, with no separation between singulars and plurals.

We recruited a total of 368 participants (153 male, 215 female, mean age: 37.3) from all English-speaking countries on Prolific. We excluded from analysis 57 participants who failed to correctly answer either one of our two elementary comprehension questions, yielding a final sample of 311 participants whose data we analyzed. Each participant answered all of the questions of the four conditions in this experiment.

Computational modeling

We computed the predictions of the CESM and the NSM following the same procedure as in Experiment 1. We fitted the value of s and γ for both models by finding the parameter values that resulted in the best fit between model judgments and average participant judgments across all four conditions. As in Experiment 1, we used a grid search, exploring a wide range of values for the parameter s , crossed with different values for a scaling parameter γ . For the CESM, the best-fitting value was $s = 0.21$ (with $\gamma = 0.39$). For the NSM we find $s = 0.02$ (with $\gamma = 0.28$).

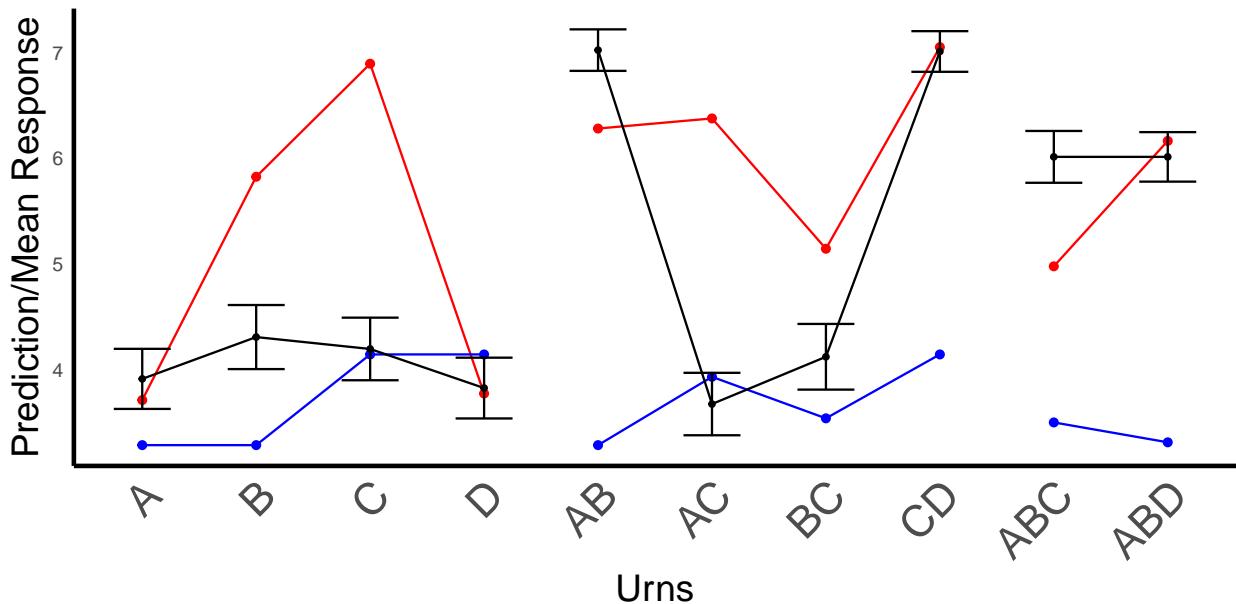
We also explored a variant of the computational models that allows for the possibility that participants interpret the condition for losing the game in a non-classical way, as described in sub-section X. We will describe this variant just before discussing the results for the loss events.

Results

We first go through the results for each round separately. We start each section by a brief exposition of the predictive performance of the CESM and NSM models for the

Figure 6

Participants' responses, along with model predictions, for the OVERDETERMINED POSITIVE round. The red line represents the CESM predictions, the blue line the NSM predictions, and the black line represents the mean of participants' responses.



round before delving into a qualitative analysis of the relevant patterns of judgments observed for that round. Note however that none of the patterns we identify or the interpretation we provide of them are dependent on the models considered, unless explicitly specified otherwise. We provide these predictions mainly for readers interested in how state-of-the art counterfactual models fare at predicting this new data.

Winning rounds

Overdetermined positive round. In this round, the player drew a colored ball from each of the four urns (as in Figure 5b) and therefore won the game.

Figure 6 summarizes the results. The CESM had a moderate but positive fit to participants' average judgments, $r(8) = .45$, while the NSM predictions were un-correlated with participants' judgments, $r(8) = -.18$.

Participants' judgments also reveal the following patterns.

Non-linearity. Participants judged that B and C were the most important singular causes. Therefore, a linear combination approach would predict that they should also view the pair $B \wedge C$ as the best plural cause. In fact, participants judged that the pairs $A \wedge B$ and $C \wedge D$ were significantly better causes than $B \wedge C$, in clear opposition to the predictions of the linear combination hypothesis.

Participants preferred non-crossing over crossing pairs. There was a clear preference for the pairs that did not cross the disjunction ($A \wedge B$, $C \wedge D$) over those that featured one variable on each side of the disjunction (e.g. $A \wedge C$, $B \wedge C$), (mean non-crossing: 7.02, mean crossing: 4.90; $t(df) = 23.971, p < 0.001$).

Weak abnormal inflation for the singular variables. We observed an abnormal-inflation effect at the level of singulars, meaning that participants deemed the urns B and C, which contained the lowest proportion of colored balls, more important for bringing about the outcome. Formally, judgments for B and C were higher than for A and D, $t(1239.3) = -2.56, p < 0.011$. This qualitative pattern aligned with the predictions of the CESM, but not with the predictions of the NSM, which predicted abnormal deflation in this context. No significant difference could be observed however between the two low-probability singulars, contrary to the CESM's expectations (means: 5.31, 5.19; $t(619.67) = 0.52, p > 0.6$).

The CESM overestimates the attractiveness of some plurals. The CESM mistakenly predicted that $A \wedge C$ should be rated higher than $A \wedge B$, and $A \wedge B \wedge D$ higher than $A \wedge B \wedge C$. In both cases, the predictions come from a tendency of the model to give a very similar rating to the singular X and the pair $X \wedge Y$ if Y is a high-probability variable. This is because if $Pr(Y)$ is high, the correlation between $X \wedge Y$ and the outcome is very similar to the correlation between X and the outcome. This property often results in erroneous predictions, not only in this particular round, but also in the TRIPLE 1 round below, where plurals containing the variable D are overestimated. We come back to this pattern in the Discussion section for this experiment.

Triple-1 round. In this round, the player drew a colored ball from urns *A*, *B* and *D* (as in Figure 5d) and therefore won the game. In such a draw, the win is not overdetermined like it was in the previous round, but clearly it is caused by the player’s getting a colored ball from both *purple* urns *A* and *B*. Urn *D*, on the other hand, is not an active cause of the win in the present world, because it has no effect on winning in the absence of *C*.

Notice that, in the particular context of this round, drawing a colored ball from urn *D* does not simply have a low impact on the win, but in a categorical sense it is not at all a cause of the outcome in the actual world. A standard view on how causal selection judgments work holds that only the events that can be counted as *actual causes* (Halpern, 2016) of the outcome qualify as candidates for causal selection in the first place (see for example Gerstenberg et al., 2021; Quillien & Lucas, 2023). Following this logic, the causal impact score of the event “drawing a colored ball from urn *D*” should simply be zero, and it is unclear if plural events that contain *D* (such as “drawing colored balls from urns *A* and *D*”) should count as actual causes or not. For simplicity, we gloss over this issue, allowing the model to give non-zero causal responsibility to *D* or plurals that feature *D*.

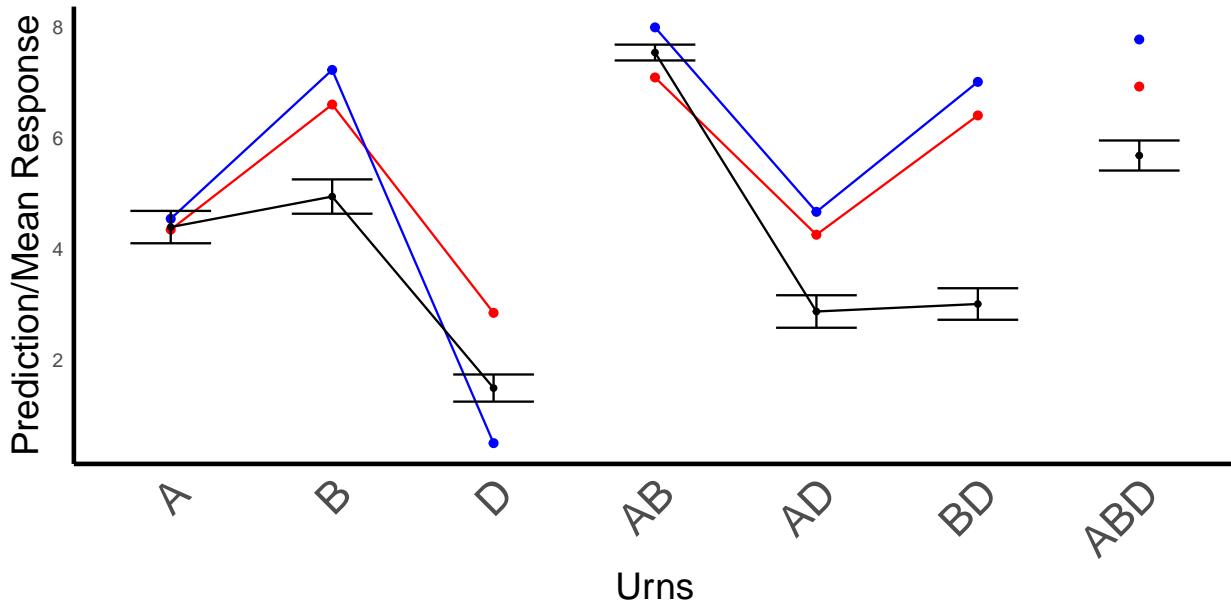
Figure 7 summarizes the results for the TRIPLE-1 rounds. Both counterfactual models give a good account of participants’ judgments: model predictions are correlated with average human judgments are $r(5) = .78$ (CESM) and $r(5) = .79$ (NSM). We now highlight the most significant patterns.

Abnormal inflation effect for singulars. We did observe an abnormal inflation effect, with the low-probability urn *B* being ranked significantly higher than high-probability urn *A* ($t(617.88) = -2.5363, p < 0.012$), in line with the predictions of both the CESM and the NSM.

Ceiling-high ratings for the pair $A \wedge B$. Participants were almost unanimous in giving ceiling-high ratings to $A \wedge B$. Only 51 participants (out of 311) in total gave it

Figure 7

Participants' responses, along with model predictions, for the TRIPLE-1 round. The red line represents the CESM predictions, the blue line the NSM predictions, and the black line represents the mean of participants' responses.



ratings different from the maximal value of the Likert scale.

Low ratings for D , and plurals containing D . Ratings for the idle variable D were very low. More than half of participants (171 out of 311) gave it maximally low ratings. Interestingly however, the ratings weren't as low as they were high for $A \wedge B$, suggesting that the fact that D does make a contribution to the win in other possible configurations still had some residual influence on participants' ratings.

Plurals containing D , such as the mixed pairs $A \wedge D$, $B \wedge D$, and the triple $A \wedge B \wedge D$, were systematically rated somewhere between the best cause that they contained and the low ratings of D . They were systematically rated lower than predicted by the models, which didn't penalize strongly enough the inclusion of the idle variable D . However, participants didn't seem to systematically disqualify a plural just for including the variable D (for example, by giving it ratings as low as those of D alone).

Losing rounds

The first two conditions just discussed collected judgments about the contribution of *colored ball* draws to a player’s *win* in a given round of the game. The two conditions we present next instead queried participants’ judgments on the contribution of *white ball* draws to a player’s *loss*.

We find that, unlike in the winning rounds, participants’ judgments about the losing rounds do not seem to be sensitive to the grouping suggested by the structure of the causal rule. That is, participants’ judgments do not appear to be sensitive to the fact that urns *A* and *B* are on one side of the disjunction, while *C* and *D* are on the other side.

Inspired by work on plural negation in natural language, we explore one possible explanation for these patterns: when participants make causal judgments about losing rounds, they might be representing the causal rule for losing the game as

$\text{LOSS} := \neg A \wedge \neg B \wedge \neg C \wedge \neg D$ (equation 3), that is “you lose if you don’t get any colored ball.”

In order to formalize this hypothesis in a counterfactual framework, we consider a variant of our computational models featuring a parameter w , which encodes participants’ propensity to represent the losing conditions in the non-classical, language-like way depicted in equation 3, as opposed to the classical, normative negation of the winning conditions. Concretely, when the outcome under consideration is a loss, the subject makes a random decision in each world, where

- with probability w , the loss is determined non-classically (equation 3 on page 35);
- otherwise, with probability $1 - w$, the loss is determined by the classical negation of the original rule, (equation 2 on page 33). This entails that the original models can be understood as a special case of the w models where $w = 0$;
- once it has been determined whether a given world is an instance of a win or a loss, the worlds that are not losses are recorded as wins. The impact of each variable on

the models is then computed exactly as before.

We fitted the models again in this new version (using data from all four conditions), via a three-dimensional grid search ($s, w; \gamma$). The best fitting values were respectively $s = 0.21$ and $w = 0.77$ (with $\gamma = 0.41$) for the CESM, and $s = .5$ and $w = 0.77$ (with $\gamma = 1.17$) for the NSM. For simplicity, all model predictions we report use the values of the s and γ parameter fitted conjointly with w , even for the base versions. Using the original fitted parameters for the base versions yields virtually identical results.

Adding the w parameter significantly improved the fit of both models, even accounting for differences in degrees of freedom (Table 6). For the negative conditions below, we report both versions of the models, to showcase the impact of the new parameter.

Overdetermined negative condition. The OVERDETERMINED NEGATIVE condition is the mirror image of the OVERDETERMINED POSITIVE condition. Here, the player drew a white ball from all four urns, and consequently lost, as pictured in Figure 5e.

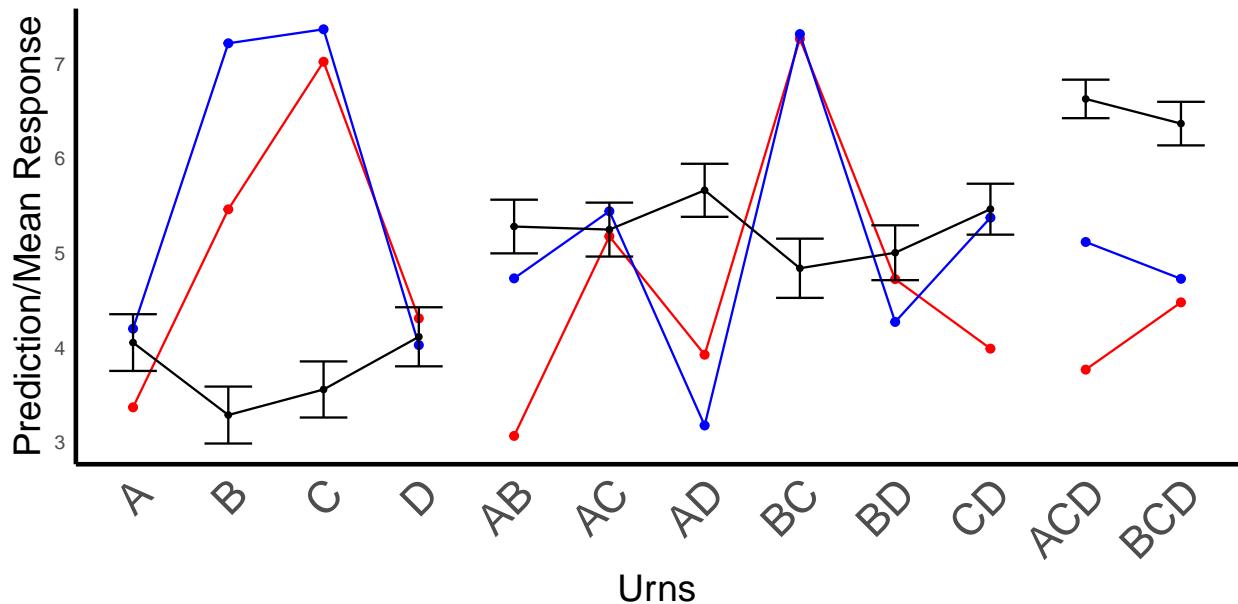
The results are summarized in Figures 8 and 9. The base versions of the CESM and NSM have a poor fit to participants' average judgments, $r(10) = -.38$ (CESM) and $r(10) = -.42$ (NSM). In contrast, the versions of the models featuring the w parameter provide a good account of the data, $r(10) = .82$ (CESM) and $r(10) = .87$ (NSM); see also Table 7. Presently we discuss our most telling findings.

Urns with the lowest number of white balls are given higher scores. This effect can be observed both for singulars and for plural causes, with combinations featuring urns *A* or *D* scoring higher than those featuring *B* or *C*. This pattern runs completely contrary to the predictions of counterfactual models under the classical representation of losing conditions from equation 2, but is captured by the version that assumes a non-classical representation of the rule.

Indeed, if participants are representing losing conditions as a disjunction of

Figure 8

Participants' responses, along with model predictions, for the OVERDETERMINED NEGATIVE round. The red line represents the CESM predictions, the blue line the NSM predictions, and the black line represents the mean of participants' responses.



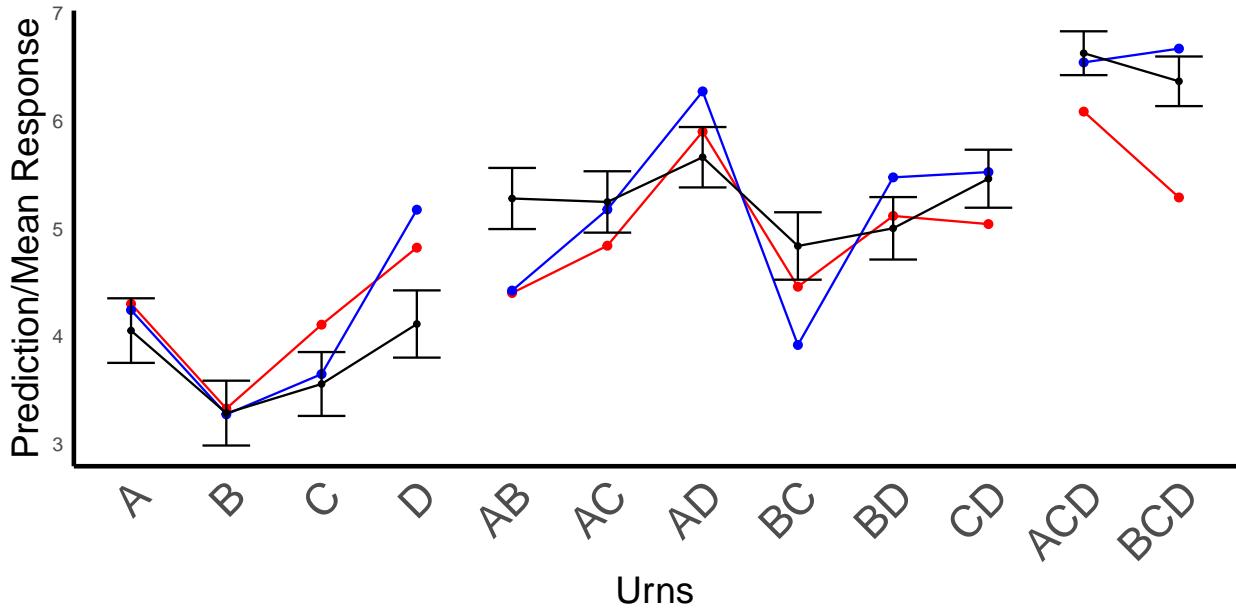
minimally sufficient conditions of that shape, we would expect their judgments to follow the logic of abnormal *deflation* and ascribe a greater causal impact to those urns out of which one is most likely to get a white ball, that is urns *B* and *C*. Instead, their judgments seem to follow a logic of abnormal *inflation*, with a preference for the urns that contain the lowest number of white balls, i.e. *A*, *D*, consistent with a representation of the losing conditions as a conjunction of necessary events as in equation 3.

No significant difference between pairs that cross the disjunction and those that do not. Participants' judgments for pairs that crossed the disjunction (e.g. *A* and *C*) were not significantly different than for pairs that did not cross the disjunction (mean crossing: 6.19; mean noncrossing: 6.37; $t(1311.9) = 1.47, p > 0.14$).

This finding is consistent with the idea that participants represent the rule for losing the game in the format $\text{LOSS} := \neg A \wedge \neg B \wedge \neg C \wedge \neg D$, i.e. in a format where there

Figure 9

Participants' responses, along with models predictions, for the OVERTERMINED NEGATIVE round, with w parameter, encoding participants' tendency to represent the losing conditions non-classically.



is no natural grouping of the variables. In contrast, a classical representation of the losing conditions would have predicted that any pair of events on the same side of the [Purple] vs. [Yellow] divide should be redundant, since a single white ball on either side is sufficient to cancel any contribution that this side could have made to a win. There is no such redundancy however if the representation is non-classical, where each white-ball-drawing event makes a crucial contribution to the outcome.

Triples are rated higher than pairs. (Mean pairs: 6.25; Mean triples: 7.37; $t(487.6) = 8.47, p < 0.001$). Here again, while triples would have been redundant under a classical representation, each element of the triple makes a non-zero contribution to the outcome if the representation is non-classical.

Triple 0 condition. In the TRIPLE 0 round, the player drew white balls from every urn except for urn C, as in Figure 5f. This makes it a mirror image of the TRIPLE 1 round,

where white balls are substituted for colored balls. In this round, the white ball from urn D is indispensable for the loss, whereas urns A and B are redundant with one another.

The same contrast between classical and non-classical representations of losing conditions applies in this round. Here, the w parameter that we enriched our models with encodes participants propensity to reinterpret the rule of the game as below.

$$\text{LOSS} := \neg A \wedge \neg B \wedge \neg D$$

We take it that the non-classical representation of the losing conditions in this round is slightly different from the OVERDETERMINED NEGATIVE round because, in the actual world, a colored ball was drawn from urn C . This makes the negation of the plural entity $C \wedge D$ in our rule harder to interpret as the strong plural negation $\neg C \wedge \neg D$, since the situation at hand is known to be one where the player in fact drew a colored ball from urn C . In other words, the player cannot possibly have lost *because* they drew a *white* ball from C , since they in fact drew a *colored* ball from C .

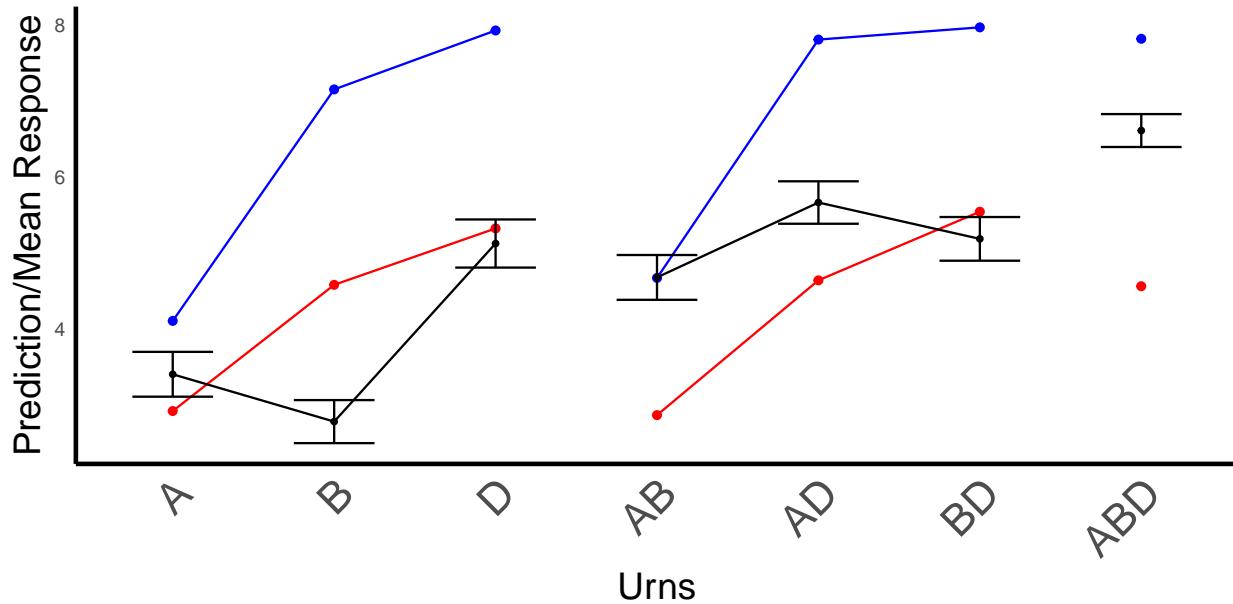
Results are summarized in Figures 10 and 11, and in Table 7. The base versions of the CESM and NSM have a moderate fit to participants' average judgments, $r(10) = .34$ (CESM) and $r(10) = .52$ (NSM). In contrast, the versions of the models featuring the w parameter provide a good account of the data, $r(10) = .94$ (CESM) and $r(10) = .99$ (NSM); see also Table 7. We highlight some of the most important qualitative patterns below.

Participants prefer urns with a lower number of white balls. Causal judgments for $\neg A$ were higher than $\neg B$ ($t(620) = 2.98, p < 0.01$), and causal judgments for $\neg A \wedge \neg D$ were higher than $\neg B \wedge \neg D$ ($t(620) = 2.34, p < 0.02$). preference for the urns featuring a lower number of white balls is similar as what we find in the OVERDETERMINED NEGATIVE round.

Again this pattern is most coherent with a non-classical representation of the losing conditions.

Figure 10

Participants' responses, along with models predictions, for the TRIPLE 0 round. The red line represents the CESM predictions, the blue line the NSM predictions, and the black line represents the mean of participants' responses.



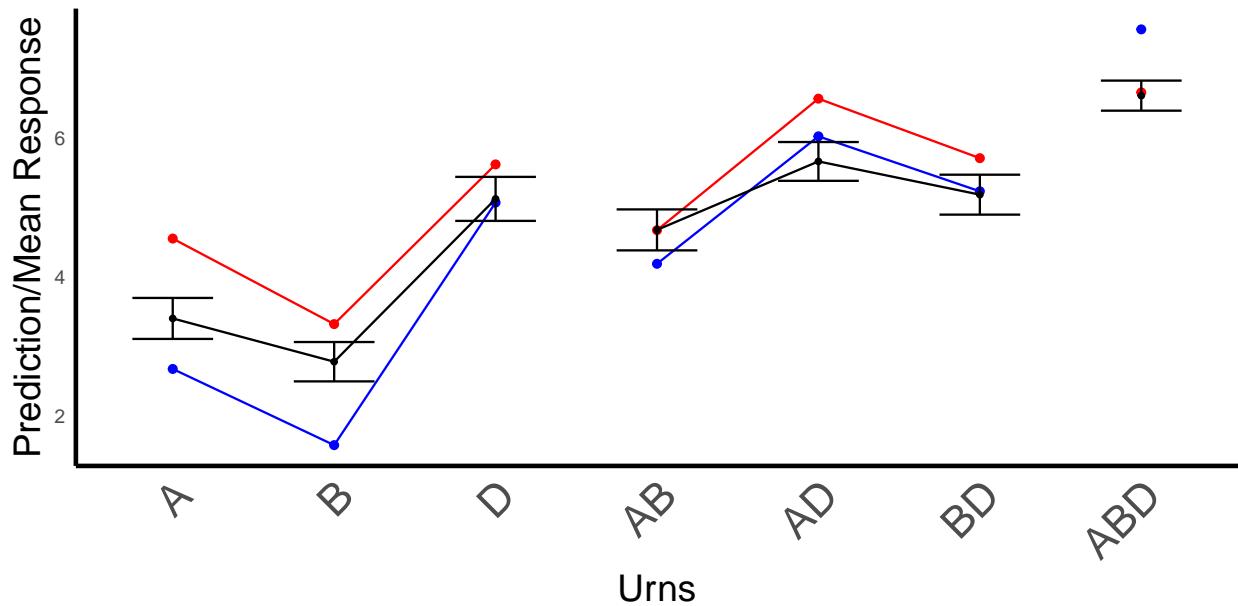
The pair $\neg A \wedge \neg B$ rates higher than either of its constitutive singulars, and the triple $\neg A \wedge \neg B \wedge \neg D$ rates higher than its constitutive pairs ($t(576) = 7.77$, $p < 0.0001$). Both patterns are examples of plurals whose effect in the outcome under the classical representation is redundant with that of one of the events (singular or plural) contained within it, which should lead them to be rated at most as high as the sufficient event in question. The fact that these are rated higher by participants is again suggestive of their representing the losing conditions non-classically.

Overall model comparison

Table 6 summarizes the comparison between the models at the global level (all conditions combined). The version of the CESM that includes the w parameter has the best fit overall (BIC = 53835.36; correlation with means: $r(35) = 0.67$, $p < 0.001$). This is better than the fit of the model without the w parameter (BIC = 55926.62; correlation with

Figure 11

Participants' responses, along with models predictions, for the TRIPLE 0 round, including the W parameter. The red line represents the CESM predictions, the blue line the NSM predictions, and the black line represents the mean of participants' responses.



means: $r(35) = 0.26, p < 0.001$, or than any of the versions of the NSM model (with w : BIC = 55295.31, cor.: $r(35) = 0.57, p < 0.001$; without w : BIC = 56967.07, cor.: $r(35) = 0.02$). In general, the versions of the models that include the w parameter are better than the versions without, by all metrics.

We also compared these models with a constant baseline model, which constantly made the same predictions about every question in every condition. The prediction was fitted to the data via the scaling parameter γ only. All counterfactual models had a better fit than the baseline model when assessed in terms of their correlations with mean human judgments, but only the version of the CESM that included the w parameter had a better BIC score than the baseline model.

Table 6

Model comparisons for Experiment 2, across all conditions.

Model	LogLik	χ^2	p-value	BIC	Cor.
Baseline	-27494			54997.32	0
CESM, with w	-26905	1175.91	< 0.0001 ***	53840.57	0.67
CESM, no w	-27963	2023.44	< 0.0001 ***	55889.97	0.2609
NSM, with w	-27634	642.55	< 0.0001 ***	55295.31	0.57
NSM, no w	-28390	11511.80	< 0.0001 ***	56967.07	0.02

Discussion

This second experiment provides more evidence in favor of the psychological reality of plural causes in the context of causal selection judgments.

Just like in the first experiment, participants' judgments for plural causes across all four rounds of the game were clearly sensitive to the probabilities attached to the corresponding events. Participants' judgments in the OVERDETERMINED POSITIVE round corroborate the non-linearity between participants' judgments for plurals and their judgments for the singular causes that constitute them. Given the pattern of abnormal inflation observed for singular variables, favoring B and C over A and D , a *linear* reconstruction of participants' judgments for plurals would have us expect the pair $B \wedge C$ to rank above all others pairs, when in fact it ranks much lower than the $A \wedge B$ and $C \wedge D$ pairs.

The winning rounds of the experiment also demonstrate that plural causes featuring more variables are *not necessarily* rated higher than proper subsets of the variables they contain. The TRIPLE 1 round shows a clear pattern in this regard: every time a plural features the variable D , its rating is systematically lower than that of the same cause (singular or plural), minus the variable D . This contradicts the hypothesis that adding more variables always makes an explanation more attractive, which the

Table 7

Table of model fits per model and condition. The Cor. column indicates the item-level correlation between model predictions and mean participant responses per question.

Condition	Model	BIC	AIC	Cor.
OVERDETERMINED	CESM	15009.14	14991.01	0.45
POSITIVE	NSM	15278.45	15260.32	-0.18
TRIPLE 1	CESM	10335.21	10318.16	0.78
	NSM	10505.71	10488.66	0.79
OVERDETERMINED	CESM, no <i>w</i>	19238.14	19225.69	-0.38
	CESM, <i>w</i>	17731.07	17718.62	0.82
	NSM, no <i>w</i>	19164.41	19151.96	-0.42
	NSM, <i>w</i>	17703.51	17684.84	0.87
TRIPLE 0	CESM, no <i>w</i>	10853.57	10842.19	0.34
	CESM, <i>w</i>	10335.21	10318.16	0.94
	NSM, no <i>w</i>	10655.7	10644.33	0.52
	NSM, <i>w</i>	10422.88	10405.82	0.99

results from Experiment 1 could not rule out. The phenomenon is not limited to the situation where an idle variable like *D* features in a plural: a similar observation can be made about the triplets $A \wedge B \wedge C$ and $A \wedge B \wedge D$ in the OVERDETERMINED POSITIVE condition, both of which are rated lower than the best pair that they contain, $A \wedge D$. Thus, although plurals featuring more variables might be descriptively more thorough, they can still be unappealing if their overall counterfactual dependence profile drops as a result of the variables added.

We also uncover properties of plural causal judgments that go beyond what is expected based purely on patterns of counterfactual dependence. First, in the winning

rounds of our second experiment, participants dislike causal explanations that “cross” the disjunction $(A \wedge B) \vee (C \wedge D)$, above and beyond what is predicted by counterfactual models. The fact that the causal rule features two clearly distinct sufficient conditions plausibly exerts an influence on participants’ explanatory preferences not fully captured by the counterfactual dependence profile of the variables in question.

Second, the following property of the CESM was not reflected in participants’ judgments. The model tends to give a very similar rating to a singular X and the pair $X \wedge Y$ if Y is a high-probability variable. This is because if $Pr(Y)$ is high, the correlation between $X \wedge Y$ and the outcome is very similar to the correlation between X and the outcome. This property often results in erroneous predictions, like in the OVERDETERMINED POSITIVE round where the model predicts (against participants’ judgments) that the pair $A \wedge C$ should rate higher than the pair $B \wedge C$, and the triplet $A \wedge B \wedge D$ higher than $A \wedge B \wedge C$.

There is likely more than one way to resolve this discrepancy. One avenue we think is worth exploring is to re-examine the assumptions we have made about the way people simulate alternatives to a plural event when they judge whether that plural event caused E . Here we made the conservative assumption that people sample a counterfactual alternative to the plural event by using the same procedure they use to sample the ‘background variables’ in the causal system (i.e. the variables that are not the current focus of causal judgment). Under this assumption when people judge whether a plural like $A \wedge C$ caused E , they sample alternatives to $A \wedge C$ in a way that is sensitive to the probabilities of both A and C . If $Pr(A)$ is high, then re-sampling $A \wedge C$ tends to have very similar effects as just re-sampling C , in that most of the $\llbracket A \wedge C \rrbracket = 0$ worlds will be $\llbracket C \rrbracket = 0$ worlds. Future research should explore the possibility that people in fact re-sample the candidate plural cause in a different way.

Finally, we found that counterfactual models could only account for participants’ judgments in the losing rounds if we assume that participants are re-interpreting the

causal rule for losing the game in a way that is not consistent with the classical-logical negation of the causal rule for winning the game, but that is the expected representation if something like the logic of natural-language plurality is what is operative here. This hypothesis is supported by the preference for urns with a lower number of white balls, for plurals containing a higher number of variables, and the absence of a preference for plurals that cross the $(A \wedge B) \vee (C \wedge D)$ disjunction in the negative rounds of the experiment.

General discussion

Humans make systematic judgments regarding which of several events influencing an outcome should be considered as *the cause*, or the most important cause of that outcome. These *causal selection* judgments are the object of a rich and actively expanding section of the psychological literature on actual causation. So far, however, that literature has been exclusively focused on *singular* events, identified with the distinct nodes of the relevant causal system. In this paper, we argue that its scope should be extended to include *plural* events, featuring multiple variables.

Our experiments present strong evidence that judgments about plural events cannot be captured in terms of linear combinations of the judgments for the events that constitute them. There appears to be no obvious way of combining participants' causal judgments regarding any two events *A* and *B* that would predict their judgment for the event "*A and B*." Our results hence establish the psychological reality of plural causes: plural causes are treated by the mind as causal entities in their own right, and their impact on the outcome is apprehended in a *holistic* fashion. We also uncovered patterns in participants' judgments that are difficult to explain under a naive view of what a plural is and how it interacts with negation, but are readily accounted for under the *sui generis* yet mathematically rigorous theories of plurality from natural-language formal semantics.

Summary of our findings and their immediate consequences

Plural cause judgments cannot be reconstructed as linear combinations of singular judgments.

It seems *prima facie* plausible that, when people make a causal judgment about whether “*A and B caused E*,” they might judge how much *A* caused *E*, judge how much *B* caused *E*, and then combine these two judgments into a single judgment for the plurality. Under this view, plural causal selection would be entirely predictable from facts about singular causal selection. One of our main goals was to rule out this null hypothesis.

In our two experiments we designed situations in which computational models predict that, if plurals are processed in a holistic manner, judgments about plural causes should not be simple combinations of judgments about their constituent singular variables. Participants’ judgments in these situations supported this prediction.

This finding is key to establishing the relevance of our object of study. Were participants’ evaluation of “*A and B*” systematically proportional to their evaluations of *A* and *B* taken separately, there would be no reason to study judgments specifically about plurals, or to build theories around such judgments.

Counterfactual models can account for plural causation judgments

A growing body of research provides strong evidence that causal judgment involves counterfactual thinking (e.g. Gerstenberg et al., 2017; Kahneman & Miller, 1986; Krasich et al., 2024; Quillien & Lucas, 2023). At the same time, there are debates about what phenomena counterfactual theories can explain (Hall, 2004; Henne, 2023; Lombrozo, 2010; Rose et al., 2021; Sytsma, 2020), and about the computations that counterfactuals might be an input to (Icard et al., 2017; Quillien, 2020).

Our experiments provide a rich opportunity to probe the out-of-distribution generalizability of counterfactual theories. None of the counterfactual theories that we are aware of were developed with the goal of explaining data about how people make plural causation judgments. Consequently, accounting for these judgments off-the-shelf

would constitute important evidence in favor of these theories.

We found that two recent counterfactual models of causal judgment (Icard et al., 2017 and in particular Quillien and Lucas, 2023) can quantitatively account for many features of participants' judgments. In particular, when participants' judgments about plurals diverge from a linear combination of their constituent singular causes, they typically do so in the way that is predicted by the counterfactual models. As such, our results appear to strengthen the case for counterfactual theories.

Yet, both counterfactual theories failed to predict the shape of people's judgments in the rounds of the game where the player *loses*. Counterfactual theories are only able to account for these data if we make an additional assumption about how participants might represent the causal structure when evaluating counterfactuals: they use representations with properties related to those found in the representations of natural-language plurals. The plausibility of a counterfactual account of these data thus depends on the plausibility of this additional assumption, which we address in detail later in this section.

Causal judgments favor sets of variables belonging to the same disjunct

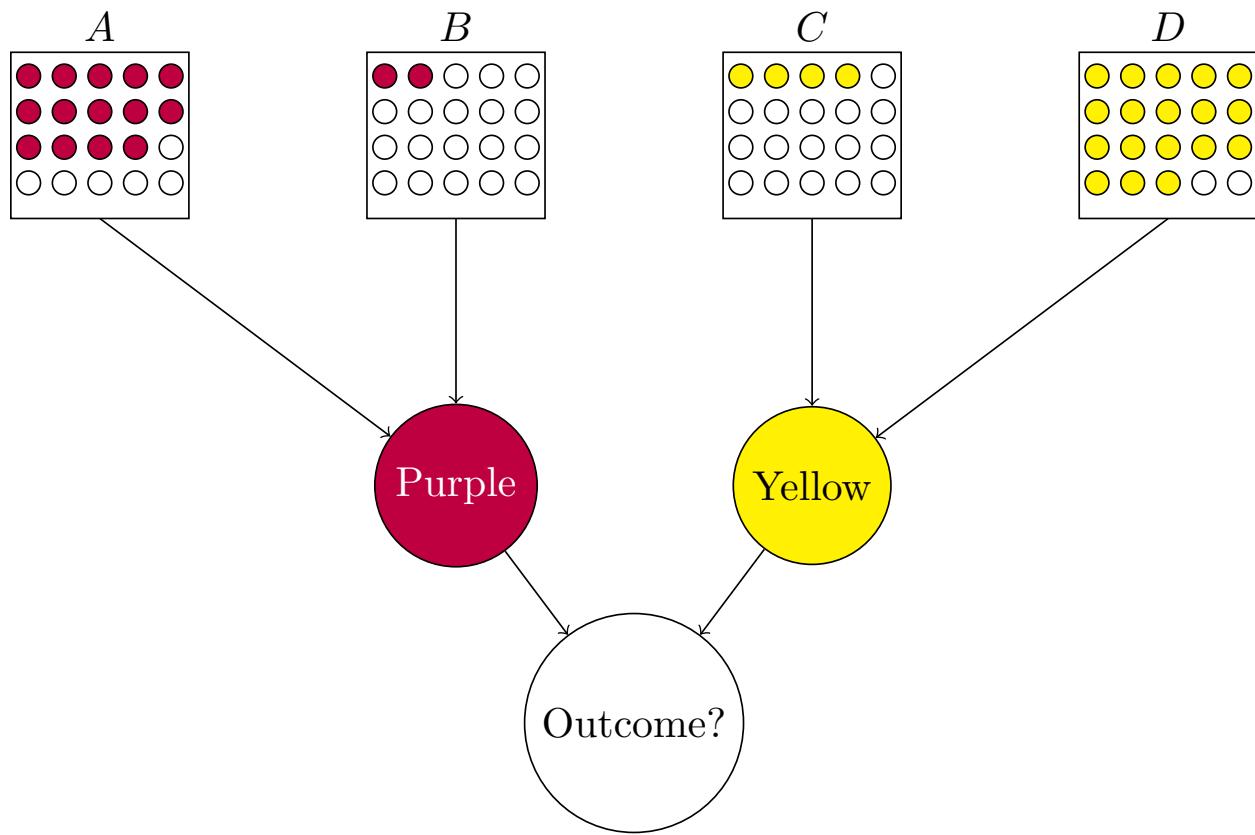
In the game that participants played in Experiment 2, the player needed to draw either two purple or two yellow balls in order to win the game. From a logical point of view, this rule is a *disjunction*: the player wins if either one of the conditions for victory is met; each condition is a *disjunct*.

Participants favored plural causes that do not cross the boundary between the two disjuncts. Suppose for example that the player drew purple balls from urns *A* and *B* and yellow balls from urns *C* and *D*. In this situation, participants would be reluctant to say that the player "won because he drew a colored ball from urn *B* and urn *C*." The counterfactual models also disfavored these boundary-crossing plurals, but participants did so to an even greater extent than predicted.

There are different possible explanations for this pattern. At a superficial level, for

Figure 12

Causal graph representing participants putative model of the situation.



example, participants might have preferred causal explanations that mentioned balls of the same color because of low-level perceptual biases.

At a deeper level, participants might have built an internal representation of the game in terms of a causal model with a particular structure. This causal model would contain intermediate variables (in the technical, causal-model sense of ‘variable’, Pearl, 2000) representing whether each condition for victory (‘getting two purple balls’) and (‘getting two yellow balls’) is met, see Figure 12. Such a model would be distinct from one without the intermediate variables, in that it would support new kinds of interventions and therefore inferences unavailable to the more straightforward model. It remains to be seen whether this difference in intervention potential makes the right predictions.

Another possibility is raised by research on the mental representation of

disjunction (Chung et al., 2022; Koralus & Mascarenhas, 2013; Walsh & Johnson-Laird, 2004). Studies of deductive reasoning have investigated how people reason about logical statements of the shape $(A \wedge B) \vee C$. In experiments replicated and varied multiple times, participants overwhelmingly conclude B from the two premises $(A \wedge B) \vee C$, and A (Koralus & Mascarenhas, 2018; Picat & Mascarenhas, 2020; Sablé-Meyer & Mascarenhas, 2021; Walsh & Johnson-Laird, 2004). But this is a fallacy: it is compatible with the premises but not the conclusion that A and C should be true while B is false. Koralus and Mascarenhas (2013) explain this fact in terms of question-answer dynamics: the disjunction in the first premise is naturally interpreted as demanding the participant *choose* between one of the two disjuncts. This in turn induces dependencies between propositions occurring *within disjuncts*: in the context of $(A \wedge B) \vee C$, the second premise A is seen as an answer in the $A \wedge B$ direction, introducing dependence between A and B . In general, this approach predicts that the conjuncts within each disjunct will be taken, as it were, to *hang together* in a cohesive way, so that learning about one will constitute evidence in favor of all of the others.

There is even evidence of such effects *absent* the language of disjunction, in experiments where the same information was conveyed by means of visual stimuli in the form of animations (Chung et al., 2022), indicating that this “packaged” way of representing a disjunction is not simply a fact about the interpretation of the word “or” and its equivalent locutions. Rather, these rich, structured disjunctive representations which induce dependencies not predicted by standard Boolean interpretations of logical connectives are available to human minds far more generally. In particular, they might have been available to participants in our Experiment 2, and may have played a part in shaping their causal judgments, by pushing them to associate A and B on the one hand and C and D on the other more tightly than is predicted by classical accounts of disjunction (whether deductive or probabilistic).

Plural causes featuring more events are not necessarily better

In Experiment 1, we found that participants preferred causal explanations that mentioned the most causes, and that this preference was stronger than predicted by counterfactual models. We probed the extent of this trend in Experiment 2, where we found that mentioning more causes does not always make a causal explanation better. For example, a causal explanation mentioning only two events *A* and *B* might be judged better than an explanation mentioning *A*, *B* and *C*.

These results suggest that causal judgments are subject to a trade-off between two different considerations. On the one hand, people might favor explanations that give detailed information about what events happened. Since every explanation of the shape “*X happened because Y*” comes with the implication that “*Y happened*” (Halpern, 2016), causal explanations that feature many causes offer more complete descriptions of what happened. With respect to this criterion, plural cause explanations are always more helpful than singular ones, since they highlight more true facts about the situation.

On the other hand, causal explanations convey information about patterns of counterfactual dependence (Quillien, 2020). Under this criterion, large plural causes can sometimes be worse. For example, an explanation mentioning three events *A*, *B*, and *C* might misleadingly suggest that the outcome strongly covaries with the conjunction of these three events, across counterfactuals.

In losing rounds, participants appear to simulate counterfactuals using a different representation of the rule

In the Experiment 2 trials where the player loses the game, we found it is difficult to account for the data if we assume that participants internally represent the conditions for losing the game as the classical complement of the conditions for winning. Instead, judgments can be captured quite adequately if we suppose that they simulate counterfactuals using a different representation of the rule of the game, in the particular case of losing rounds. Specifically, participants seem to represent the losing conditions

as “the player loses the game if they draw a white ball from all urns,” an interpretation not supported by standard Boolean logic.

The representation our participants seem to have for the losing conditions is however quite consistent with fundamental facts about the semantics of plurals. In natural language, the negation of a plural event is naturally represented as the negation of each of its singular constituents, rather than the negation of their conjunction. That is, for a sentence like “Mary and John didn’t come to the party,” we infer that Mary did not come and John did not come, rather than that at least one of them did not come (Krifka, 1996; Lappin, 1989; Löbner, 2000; Szabolcsi & Haddican, 2004). In our experimental setting, negating the conditions for winning the game yields the following losing conditions expressed in natural language: “you lose if you don’t get the purple balls and you don’t get the yellow balls.” Natural-language plural negation then predicts that these losing conditions correspond to not getting *any* purple balls and not getting *any* yellow balls. Consequently, our results provide suggestive evidence that a signature effect of the semantics of plurals, often dubbed *homogeneity* in the linguistics literature, affects people’s representations of events.

It is important to be more precise about the exact level at which we take this mental representation of the rule to occur. We are not claiming that our participants explicitly believed that one needs to draw a white ball from *every single urn* on the screen to lose. We did not directly probe participants’ judgments on their understanding of the rule, but they played the game of chance for 10 rounds before responding to any causal-judgment questions, and these draws included losing cases in which colored balls were drawn. Moreover, the effects we find are unlikely to stem from linguistic experimenter demands when interpreting the causal statements verbally. We asked participants about the causal impact of “drawing a white ball” on a “loss,” never about the impact of “not drawing a colored ball” on “not winning.”

Instead, we propose that a non-classical representation of the losing conditions is

deployed when participants implicitly simulate counterfactual possibilities. In other words, the effects we find stem from features of a (likely unconscious) process of counterfactual simulation, rather than as explicit and conscious misconceptions about the losing conditions of our game of chance.

Broader theoretical implications: natural-language semantics and causal cognition

Two of our proposals were inspired by work on natural language, specifically work on the formal semantics of plurals. This work inspired our proposals that i) plural causes are processed holistically, and ii) people deploy the equivalent of natural-language plural negation when simulating counterfactuals for losing events.

Why should work on natural language semantics be relevant to causal judgment? Here we present several considerations which we mean neither exclusively nor exhaustively.

First, participants in our experiments have to use natural language to read the description of the causal structure (i.e. the rules of the game) and the causal statements they have to evaluate. This stage of linguistic processing might “package” information into a particular representational format. For example, when participants read “the player won the game because he got a colored ball from urns *A* and *B*,” they might process the event “he got a colored ball from urns *A* and *B*” as one holistic entity, as is typically the case according to linguistic theories of plurality (Link, 1983). The causal judgment process then treats the plural event as a holistic entity because it received the input in that format. Consequently, it’s not that participants’ causal reasoning is manipulating language-like representations, it’s simply that the language suggested some degree of togetherness between the variables mentioned, by virtue of plural semantics, but the causal-reasoning system now handles this togetherness in its own proprietary way.

In favor of this hypothesis is the fact that our basic experimental findings regarding holistic evaluations of plural causes are largely predicted by the two causal-selection

theories we considered. Both causal-selection theories, as instantiated in our models, had a good fit with participants' judgments, especially in Experiment 1. On the other hand, our Experiment 2 provides evidence that subjects' judgments traffic in plural representations even when those are not directly prompted by the instructions of the experiment, in ways that go beyond the predictions of the causal-selection theories at hand. In our losing rounds, participants had to consider the conditions for *not winning*, in order to assess the extent to which "the player lost because" of this or that cause. Now, the conditions for *winning* were given linguistically and included two plurals connected by a disjunction ("colored balls from urns *A* and *B* or colored balls from urns *C* and *D*"), but crucially our instructions never presented the *negation* of this disjunction of plurals. Accordingly, and unlike the fundamental facts about holistic entities, the causal-selection theories in and of themselves accounted for the losing rounds very poorly indeed.

This more surprising fact can be interpreted in at least two different ways. One possibility is that some participants might be running an internal monologue when completing the task. That is, subjects might be *talking to themselves* in the course of their attempts to put together a representation of the relevant causal structure. For example, in the conditions where the player loses the game, they might be reconstructing the conditions for losing by saying to themselves "you lose if you don't get the purple balls and you don't get the yellow balls, this means that you lose if you don't get any colored balls."

Such a proposal is in principle testable. For example, the same sort of cognitive-load manipulations that are known to reduce people's tendency for pragmatic inferences (De Neys & Schaeken, 2007; Marty & Chemla, 2013; Picat & Mascarenhas, 2019, 2020; van Tiel et al., 2019) might also be used to interfere with the sort of self-talk that is hypothesized here. Similarly for verbal shadowing tasks, which can interfere with reasoning processes that plausibly rely on natural language (Carruthers, 2002). If such self-talk is the cause of the signature effects of plurals that we identified here, such as

subjects' non-classical reconstructions of the losing conditions, we would expect these manipulations to bring subjects' behavior back in line with the classical representation of the negation of the winning conditions. While we think such a study would be worth doing, we are skeptical of the prediction, which is that participants would have been perfectly capable of handling this complex rule had they not been engaging in deliberate reasoning.

Another hypothesis, which we favor, is that our judgments about causes display the same sort of effects found in language because the underlying cognitive processes operate on the same set of representations as our language faculty itself. That is, both natural language and causal judgment might depend on a shared *language and logic of thought* which supports many of our higher cognitive faculties.

These considerations resonate with the recent renaissance which Jerry Fodor's (1975, 2008) *language of thought* hypothesis has been enjoying. As Quilty-Dunn et al. (2023) observe in a recent target article in *Behavioral and Brain Sciences*, much current research embraces the idea that human cognition relies on symbolic representations of a language-like nature, of the kind that Fodor proposed were at the core of human thought.

While the classical illustrations of language of thought came chiefly from the domain of natural language and general purpose, integrative thought, current research on this program has been paying particular attention to areas of cognition that are minimally connected to language, if at all, and plausibly do not require integrative thought. Part of the reason for this is sound methodology: as noted by Quilty-Dunn et al. (2023), those investigations provide a new class of arguments in favor of language of thought as a general hypothesis, showing how language-like representations might pervade cognition across the board, and in plausibly domain-specific ways.

The work presented here has a mixed status in this regard. On the one hand, our participants seem to make use of language-like representations in a cognitive process that does not inherently depend on natural language, namely counterfactual sampling

and causal judgment. Indeed, no extant theory of causal judgment even suggests that the phenomenon might hinge on language in any appreciable way. This might suggest that we are seeing here yet another proprietary language of thought: a language of thought for causation.

On the other hand, causal reasoning applies to every walk of human life, from ecologically natural contexts like tool making or animal husbandry, to highly abstract modern contexts such as science or public policy. Given that we are arguing for representations entirely parallel to linguistic representations, the unbounded general-purpose representational system *par excellence* (Chomsky, 1965; Hockett, 1960), we are inclined to think that causal reasoning, whether computed by the general-purpose reasoning system or by a causation-specific system, taps into the language of thought in the broadest sense: general-purpose, integrative thought.

Our perspective also recovers Fodor's view that, significant differences between natural language and thought notwithstanding, natural language itself is an important tool for investigating thought (Fodor, 2001), in the sense we are concerned with here. Fodor's work came at the tail end of an earlier renaissance of Enlightenment philosophy originating in linguistics and most eloquently articulated in Noam Chomsky's (1966) essay on what he dubbed "Cartesian linguistics"; namely, the view that studying the properties of natural language provides a privileged window into the properties of human thought itself. This perspective aligns with Wellwood and Hunter's (2023) commentary to the BBS target article, where they submit that natural language remains "a particularly fruitful domain for formulating *specific* hypotheses about candidate [language of thought] representations."

We embrace this approach and would add that specifically *formal natural-language semantics* constitutes a particularly fruitful source of hypotheses about the most general-purpose language of thought. The present study can be taken (among other things) as a case study illustrating the power of this methodology. The formal

semantics of plurals allowed us to consider novel hypotheses about causal cognition. The result is a contribution not just to our understanding of causal judgment, but potentially also the representational structure of higher cognition more generally: alongside the nigh-universally recognized standard Boolean operations, human thought might involve the more general and more expressive lattice-theoretic operations found in linguists' theories of plurality. Crucially, this contribution gets to co-opt the mathematical rigor that characterizes the linguistic work which inspired it, providing not just a *language* of thought, but also a *logic* and a *model theory* of thought.

This last point deserves dwelling on. It is not uncommon in psychology to hear that mathematical logic and model theory are constitutionally inadequate to model thought and reasoning, due to their purported rigidity. If one looks closely however, the arguments do not support this strong indictment. Take the two dominating paradigms in the psychology of reasoning: mental-model theory and the probabilistic “New Paradigm.” Johnson-Laird et al. (2024) argue that the right basis of human reasoning is “models of possibilities instead of logic,” but the fine print makes it transparently clear that their argument is against *classical* logics (Bringsjord & Govindarajulu, 2020). And indeed it can be shown that a mental-model theory of reasoning can be fully formalized as a non-classical, non-standard logic (Koralus & Mascarenhas, 2013). This allows for metalogical results of real psychological import (Holliday & Icard, 2018), in particular *proving* what reasoning strategies in the theory will guarantee sound reasoning. The rational-analysis probabilistic theory of Oaksford and Chater (2007) in turn proposes that, rather than binary decisions of validity as classical logic would have it, humans make intrinsically probabilistic judgments, even in what ostensibly look like deductive-inference tasks. This may well be right, yet the Bayesian probabilistic calculus ubiquitously used in this work in no way does without a logic. In fact, it *presupposes* one! For probability measures apply to complex propositions built from (at the very least) propositional connectives. The probabilistic paradigm in the psychology of reasoning is

an *extension* of classical logic, not a repudiation, and accordingly other logics can be extended as well into non-classical probabilistic calculi (Narens, 2015).

Zooming out, by the terms “logic” and “model theory” we mean what logicians and mathematicians mean: mathematically rigorous frameworks in which we can give mathematically rigorous theories of inference (“logic,” “proof theory”) and meaning (“model theory”). This includes not just the familiar classical logical systems from Philosophy 101, but a host of non-classical systems which have none of the properties psychologists rightly criticize in classical logic. From the statement that one is interested in offering a proposal for logical or model-theoretic properties of thought, in no way does it follow that there will be only two truth values, that statements of the form $p \vee \neg p$ are trivial, that contradictions license any arbitrary conclusion, that conjunction is commutative, and so forth. All that follows is that one is interested in giving a mathematically rigorous formalization of inference and meaning, surely self-evidently something a theory of thought must aspire to.

In conclusion, we think that the time is ripe to formulate strong and mathematically rigorous hypotheses about the representational arsenal of human thought, which entails giving precise proposals for logics and model-theories of thought. Formal natural-language semantics offers a largely untapped fountain of such hypotheses, with few though notable exceptions (see in particular Phillips & Kratzer, 2024, for one of our favorites). We are sympathetic to the view that there may be many domain-specific languages of thought (Mandelbaum et al., 2022; Sablé-Meyer et al., 2022), but we find that currently there is far more that can be done and in fact has already been done on the domain-general language of thought. Unlike its special-purpose relatives, the general language of thought has been the exclusive focus of investigation for the past fifty years on the part of a small but dedicated community of linguists, who have been fastidiously building an impressively broad and deep body of mathematically precise scholarship on the model theory and the logic of general

language of thought in this sense. Many if not most would be loath to characterize their work this way, and would likely insist that their theories apply, at best, to *language*, thought being the purview of the philosophers and the psychologists. We disagree, and we believe that the present work illustrates how, if the language of thought is at least one of the best games in town to cash out a computational and representational theory of mind, then the *formal* study of natural language *meaning* offers the most productive and most woefully underexplored path toward building a mathematically rigorous theory of domain-general mental representations.

References

- Arendt, H. (1987). Collective responsibility. In J. W. S. Bernauer (Ed.), *Amor mundi: Explorations in the faith and thought of Hannah Arendt* (pp. 43–50). Springer, Dordrecht. https://doi.org/10.1007/978-94-009-3565-5_3
- Bringsjord, S., & Govindarajulu, N. S. (2020). Rectifying the mischaracterization of logic by mental model theorists. *Cognitive Science*, 44(12).
<https://doi.org/10.1111/cogs.12898>
- Carruthers, P. (2002). The cognitive functions of language. *Behavioral and Brain Sciences*, 25(6), 657–674. <https://doi.org/10.1017/S0140525X02000122>
- Champollion, L., & Krifka, M. (2016). Mereology. In P. Dekker & M. Aloni (Eds.), *The Cambridge handbook of formal semantics* (pp. 369–388). Cambridge University Press. <https://doi.org/10.1017/CBO9781139236157.014>
- Chater, N., & Oaksford, M. (2013). Programs as causal models: Speculations on mental programs and mental representation. *Cognitive Science*, 37(6), 1171–1191.
<https://doi.org/10.1111/cogs.12062>
- Chomsky, N. (1965). *Aspects of the theory of syntax*. Cambridge, Mass: MIT Press.
- Chomsky, N. (1966). *Cartesian linguistics: A chapter in the history of rationalist thought*. Harper; Row, New York.
- Chung, W., Bade, N., Blanc-Cuenca, S., & Mascarenhas, S. (2022). Question-answer dynamics in deductive fallacies without language. In J. Culbertson, A. Perfors, H. Rabagliati, & V. Ramenzoni (Eds.), *Proceedings of the 44th annual meeting of the cognitive science society*. <https://escholarship.org/uc/item/9711612q>
- De Leeuw, J. R. (2015). Jsppsych: A javascript library for creating behavioral experiments in a web browser. *Behavior research methods*, 47(1), 1–12.
<https://doi.org/10.3758/s13428-014-0458-y>

- De Neys, W., & Schaeken, W. (2007). When people are more logical under cognitive load: Dual task impact on scalar implicature. *Experimental Psychology*, 54(2), 128–133. <https://doi.org/10.1027/1618-3169.54.2.128>
- Fodor, J. (1975). *The language of thought*. Harvard University Press.
- Fodor, J. (2001). Language, thought and compositionality. *Mind & Language*, 16(1), 1–15. <https://doi.org/10.1111/1468-0017.00153>
- Fodor, J. (2008). *LOT 2: The language of thought revisited*. Oxford University Press.
- Gerstenberg, T., Goodman, N., Lagnado, D., & Tenenbaum, J. (2021). A counterfactual simulation model of causal judgments for physical events. *Psychological Review*, 128. <https://doi.org/10.1037/rev0000281>
- Gerstenberg, T., & Icard, T. (2020). Expectations affect physical causation judgments. *Journal of Experimental Psychology: General*, 149(3), 599.
- Gerstenberg, T., Peterson, M. F., Goodman, N. D., Lagnado, D. A., & Tenenbaum, J. B. (2017). Eye-tracking causality. *Psychological science*, 28(12), 1731–1744.
- Gerstenberg, T., & Tenenbaum, J. B. (2017). Intuitive theories. In *Oxford handbook of causal reasoning* (pp. 515–548). Oxford University Press.
- Gill, M., Kominsky, J. F., Icard, T. F., & Knobe, J. (2022). An interaction effect of norm violations on causal judgment. *Cognition*, 228, 105183.
- Griffiths, T., & Tenenbaum, J. (2005). Structure and strength in causal induction. *Cognitive psychology*, 51, 334–84.
<https://doi.org/10.1016/j.cogpsych.2005.05.004>
- Hall, N. (2004). Two concepts of causation.
- Halpern, J. Y. (2015). A modification of the Halpern-Pearl definition of causality. *arXiv preprint*. <https://doi.org/10.48550/arXiv.1505.00162>
- Halpern, J. Y. (2016). *Actual causality*. MIT Press.

- Halpern, J. Y., & Pearl, J. (2005). Causes and explanations: A structural-model approach. part i: Causes. *British Journal for the Philosophy of Science*, 56(4), 843–887. <https://doi.org/10.1093/bjps/axi147>
- Henne, P. (2023). Experimental metaphysics: Causation. *The compact compendium of experimental philosophy*. De Gruyter.
- Henne, P., Kulesza, A., Perez, K., & Houcek, A. (2021). Counterfactual thinking and recency effects in causal judgment. *Cognition*, 212, 104708.
- Henne, P., Niemi, L., Pinillos, Á., De Brigard, F., & Knobe, J. (2019). A counterfactual explanation for the action effect in causal judgment. *Cognition*, 190, 157–164.
- Hitchcock, C. (2001). The intransitivity of causation revealed in equations and graphs. *The Journal of Philosophy*, 98, 273–299. <https://doi.org/10.2307/2678432>
- Hitchcock, C. (2012). Portable causal dependence: A tale of consilience. *Philosophy of Science*, 79(5), 942–951. Retrieved April 15, 2024, from <http://www.jstor.org/stable/10.1086/667899>
- Hockett, C. F. (1960). The origin of speech. *Scientific American*, 203(3), 88–97. <https://doi.org/10.1038/scientificamerican0960-88>
- Holliday, W. H., & Icard, T. F. (2018). Axiomatization in the meaning sciences. In D. Ball & B. Rabem (Eds.), *The science of meaning: Essays on the metatheory of natural language semantics*. Oxford University Press. <https://doi.org/10.1093/oso/9780198739548.003.0002>
- Icard, T. F. (2016). Subjective probability as sampling propensity. *Review of Philosophy and Psychology*, 7(4), 863–903. <https://doi.org/10.1007/s13164-015-0283-y>
- Icard, T. F., Kominsky, J. F., & Knobe, J. (2017). Normality and actual causal strength. *Cognition*, 161, 80–93. <https://doi.org/10.1016/j.cognition.2017.01.010>
- Johnson-Laird, P. N., Byrne, R. M. J., & Khemlani, S. (2024). Models of possibilities instead of logic as the basis of human reasoning. *Minds and Machines*, 34(19). <https://doi.org/10.1007/s11023-024-09662-4>

- Kahneman, D., & Miller, D. T. (1986). Norm theory: Comparing reality to its alternatives. *Psychological review*, 93(2), 136.
- Kinney, D. B., & Lombrozo, T. (2024). Building compressed causal models of the world. *Cognitive Psychology*. <https://osf.io/preprints/psyarxiv/2f7x6>
- Kirfel, L., Icard, T., & Gerstenberg, T. (2021). Inference from explanation. *Journal of Experimental Psychology: General*, 151. <https://doi.org/10.1037/xge0001151>
- Knobe, J., & Fraser, B. (2008). Causal judgment and moral judgment: Two experiments. In W. Sinnott-Armstrong (Ed.), *Moral psychology*. MIT Press.
- Kominsky, J., Phillips, J., Gerstenberg, T., Lagnado, D., & Knobe, J. (2015). Causal superseding. *Cognition*, 137, 196–209.
<https://doi.org/10.1016/j.cognition.2015.01.013>
- Konuk, C., Goodale, M., Quillien, T., & Mascarenhas, S. (2023). Plural causes in causal judgment. In M. Goldwater, F. K. Anggoro, B. K. Hayes, & D. C. Ong (Eds.), *Proceedings of the 45th annual meeting of the cognitive science society* (pp. 3180–3186). <https://escholarship.org/uc/item/0014w3r1>
- Koralus, P., & Mascarenhas, S. (2013). The erotetic theory of reasoning: Bridges between formal semantics and the psychology of deductive inference. *Philosophical Perspectives*, 27, 312–365. <https://doi.org/10.1111/phpe.12029>
- Koralus, P., & Mascarenhas, S. (2018). Illusory inferences in a question-based theory of reasoning. In K. Turner & L. Horn (Eds.), *Pragmatics, truth, and underspecification: Towards an atlas of meaning* (pp. 300–322, Vol. 34). Leiden: Brill. https://doi.org/10.1163/9789004365445_011
- Krasich, K., O'Neill, K., & De Brigard, F. (2024). Looking at mental images: Eye-tracking mental simulation during retrospective causal judgment. *Cognitive Science*, 48(3), e13426.

- Krifka, M. (1996). Pragmatic strengthening in plural predication and donkey sentences. In T. Galloway & J. Spence (Eds.), *Proceedings of SALT 6* (pp. 136–153).
<https://doi.org/10.3765/salt.v6i0.2769>
- Križ, M., & Spector, B. (2021). Interpreting plural predication: Homogeneity and non-maximality. *Linguistics and Philosophy*, 44, 1131–1178.
<https://doi.org/10.1007/s10988-020-09311-w>
- Lappin, S. (1989). Donkey pronouns unbound. *Theoretical Linguistics*, 15(3), 263–289.
<https://doi.org/10.1515/thli.1988.15.3.263>
- Lewis, D. (1973). *Counterfactuals*. Oxford: Basil Blackwell.
- Link, G. (1983). The logical analysis of plural and mass terms: A lattice-theoretical approach. In R. Bäuerle, C. Schwarze, & A. von Stechow (Eds.), *Meaning, use, and interpretation of language* (pp. 302–323). Berlin: Walter de Gruyter.
- Löbner, S. (2000). Polarity in natural language: Predication, quantification and negation in particular and characterizing sentences. *Linguistics and Philosophy*, 23, 213–308. <https://doi.org/10.1023/A:1005571202592>
- Lombrozo, T. (2010). Causal-explanatory pluralism: How intensions, functions, and mechanisms influence causal ascriptions. *Cognitive Psychology*, 61(4), 302–332.
<https://doi.org/10.1016/j.cogpsych.2010.05.002>
- Lucas, C. G., & Kemp, C. (2015). An improved probabilistic account of counterfactual reasoning. *Psychological review*, 122(4), 700.
- Mandelbaum, E., Dunham, Y., Feiman, R., Firestone, C., Green, E., Harris, D., Kibbe, M. M., Kurdi, B., Mylopoulos, M., Shepherd, J., Wellwood, A., Porot, N., & Quilty-Dunn, J. (2022). Problems and mysteries of the many languages of thought. *Cognitive Science*, 46(12). <https://doi.org/10.1111/cogs.13225>
- Marty, P. P., & Chemla, E. (2013). Scalar implicatures: Working memory and a comparison with *only*. *Frontiers in Psychology*, 4(403).
<https://doi.org/10.3389/fpsyg.2013.00403>

- Miller, S. (2001). Collective responsibility. *Public Affairs Quarterly*, 15(1), 65–82.
<https://www.jstor.org/stable/40441276>
- Morris, A., Phillips, J., Gerstenberg, T., & Cushman, F. (2019). Quantitative causal selection patterns in token causation. *PLoS ONE*, 14(8).
<https://doi.org/10.1371/journal.pone.0219704>
- Narens, L. (2015). *Probabilistic lattices with applications to psychology*. World Scientific.
- Novick, L. R., & Cheng, P. W. (2004). Assessing interactive causal influence.
Psychological Review, 111(2), 455–485.
<https://doi.org/10.1037/0033-295X.111.2.455>
- Oaksford, M., & Chater, N. (2007). *Bayesian rationality: The probabilistic approach to human reasoning*. Oxford University Press.
- O'Neill, K., Quillien, T., & Henne, P. (2022). A counterfactual model of causal judgment in double prevention. *Conference in computational cognitive neuroscience*.
- Pearl, J. (2000). *Causality: Models, reasoning, and inference*. Cambridge University Press.
- Pearl, J., & Mackenzie, D. (2018). *The book of why: The new science of cause and effect*. Basic Books, Inc. <https://doi.org/10.5555/3238230>
- Phillips, J., & Kratzer, A. (2024). Decomposing modal thought. *Psychological Review*, 131(4), 966–992. <https://doi.org/10.1037/rev0000481>
- Picat, L., & Mascarenhas, S. (2019). Reasoning with disjunctions under cognitive load [Talk given at *Brain, Language and Learning* 2019, University of Siena].
- Picat, L., & Mascarenhas, S. (2020). *On the interplay between interpretation and reasoning in compelling fallacies* [Manuscript under review available on PsyArXiv].
<https://osf.io/preprints/psyarxiv/8kywa>
- Quillien, T. (2020). When do we think that X caused Y? *Cognition*, 205.
<https://doi.org/10.1016/j.cognition.2020.104410>

Quillien, T., & Barlev, M. (2022). Causal judgment in the wild: Evidence from the 2020 U.S. presidential election. *Cognitive Science*, 56(2).

<https://doi.org/10.1111/cogs.13101>

Quillien, T., & Lucas, C. G. (2023). Counterfactuals and the logic of causal selection. *Psychological Review*. <https://doi.org/10.1037/rev0000428>

Quillien, T., Szollosi, A., Bramley, N. R., & Lucas, C. (2023). Causal inference shapes counterfactual plausibility. *Proceedings of the Annual Meeting of the Cognitive Science Society*, 45(45).

Quilty-Dunn, J., Porot, N., & Mandelbaum, E. (2023). The best game in town: The re-emergence of the Language of Thought hypothesis across the cognitive sciences. *Behavioral and Brain Sciences*, 46(e292).

<https://doi.org/10.1017/S0140525X22002849>

R Core Team. (2013). *R: A language and environment for statistical computing* [ISBN 3-900051-07-0]. R Foundation for Statistical Computing. Vienna, Austria.

Rose, D., Sievers, E., & Nichols, S. (2021). Cause and burn. *Cognition*.

Sablé-Meyer, M., Ellis, K., Tenenbaum, J. B., & Dehaene, S. (2022). A language of thought for the mental representation of geometric shapes. *Cognitive Psychology*, 139(101527). <https://doi.org/10.1016/j.cogpsych.2022.101527>

Sablé-Meyer, M., & Mascarenhas, S. (2021). Indirect illusory inferences from disjunction: A new bridge between deductive inference and representativeness. *Review of Philosophy and Psychology*, 12(2). <https://doi.org/10.1007/s13164-021-00543-8>

Sloman, S. A., & Lagnado, D. A. (2015). Causality in thought. *Annual Review of Psychology*, 66, 223–247.

Sytsma, J. (2020). Causation, responsibility, and typicality. *Review of Philosophy and Psychology*, 12(4), 699–719. <https://doi.org/10.1007/s13164-020-00498-2>

- Szabolcsi, A., & Haddican, B. (2004). Conjunction meets negation: A study in cross-linguistic variation. *Journal of Semantics*, 21(3), 219–249.
<https://doi.org/10.1093/jos/21.3.219>
- van Tiel, B., Pankratz, E., Marty, P., & Sun, C. (2019). Scalar inferences and cognitive load. *Proceedings of Sinn und Bedeutung* 23, 2, 427–441.
- Walsh, C., & Johnson-Laird, P. N. (2004). Coreference and reasoning. *Memory and Cognition*, 32(1), 96–106. <https://doi.org/10.3758/BF03195823>
- Wellwood, A., & Hunter, T. (2023). Linguistic meanings in mind. *Behavioral and Brain Sciences*, 46(e289). <https://doi.org/10.1017/S0140525X23001887>
- Woodward, J. F. (2003). *Making things happen: A theory of causal explanation*. Oxford University Press.
- Woodward, J. F. (2006). Sensitive and insensitive causation. *The Philosophical Review*, 115. <https://doi.org/10.1215/00318108-115-1-1>