A CAUSAL INFERENCE APPROACH OF MONOSYNAPSES FROM SPIKE TRAINS

Zach Saccomano School of Neuroscience, Virginia Tech zachsaccomano@vt.edu Sam McKenzie Health Science Center, University of New Mexico samckenzie@salud.unm.edu

Horacio G. Rotstein

Federated Department of Biological Sciences, New Jersey Institute of Technology & Rutgers University

Asohan Amarasingham

Department of Mathematics, The City College of NY Depts. of Biology and Computer Science, The Graduate Center City University of New York aamarasingham@ccny.cuny.edu

ABSTRACT

Neuroscientists have worked on the problem of estimating synaptic properties, such as connectivity and strength, from simultaneously recorded spike trains since the 1960s. Recent years have seen renewed interest in the problem, coinciding with rapid advances in the technology of high-density neural recordings and optogenetics, which can be used to calibrate causal hypotheses about functional connectivity. Here, a rigorous causal inference framework for pairwise excitatory and inhibitory monosynaptic effects between spike trains is developed. Causal interactions are identified by separating spike interactions in pairwise spike trains by their timescales. Fast algorithms for computing accurate estimates of associated quantities are also developed. Through the lens of this framework, the link between biophysical parameters and statistical definitions of causality between spike trains is examined across a spectrum of dynamical systems simulations. In an idealized setting, we demonstrate a correspondence between the synaptic causal metric developed here and the probabilities of causation developed by Tian and Pearl [1]. Since the probabilities of causation are derived under distinct assumptions and include data from experimental randomization, this opens up the possibility of testing the synaptic inference framework's assumptions with juxtacellular or optogenetic stimulation. We simulate such an experiment with a biophysically detailed channelrhodopsin model and show that randomization is not achieved; strong confounding persists even with strong stimulations. A principal goal is to ask how carefully articulated causal assumptions might better inform the design of neural stimulation experiments and, in turn, support experimental tests of those assumptions.

Keywords functional connectivity · causal inference · spike trains · dynamical systems

1 Introduction

Various lines of experimental evidence suggest that, in some neuronal pairs, monosynaptic input can reliably produce a postsynaptic spike response *in vivo* [2, 3] with a delay and precision that acts on millisecond timescales. Moreover, it appears that the most plausible explanation for corresponding observations of appropriately-timed millisecond-timescale correlations is the presence of a monosynaptic connection between the two cells. What is more, there is evidence that the magnitude of such fine-timescale correlations co-vary with the synapse's strength [4]. This suggests that a

careful study of millisecond-timescale correlations in simultaneously-recorded spike trains might be a tool for studying synaptic dynamics during behavior [5].

In practice, the hypothesis that monosynaptic effects can act on millisecond timescales is often incorporated into their analysis by a statistical formulation of a separation of timescale hypothesis. This formulation can be appreciated by looking at anecdotal examples of cross-correlograms (CCGs) from studies [2] that offer support for a causal interpretation by juxtacellular and optogenetic stimulation of putative presynaptic neurons *in vivo* (see Figure 1). Nevertheless, causal claims in highly connected systems ought to be treated delicately. It has been suggested that isolating fine timescale effects might be a way to sidestep such concerns [6, 7]. While many methods have been proposed for monosynaptic inference, relatively few have modeled causal relationships explicitly (but see [8, 7]). Furthermore, many such methods operate on the CCG, and even under a timescale-separation assumption, the CCG is insufficient to identify synaptic properties even in quite simple models (Figure 1) [9].

The primary focus of this study is to contribute to the development of robust and rigorous approaches to monosynaptic inference in which the causal inference is explicit. We develop a causal inference framework for monosynaptic interactions that is based on separation of timescale hypotheses that are robust to strong forms of nonstationarity in the background dynamics (the concern we have most heavily emphasized in the context of spike train analysis more generally [10, 11, 12]), among other forms of model misspecification. Unbiased estimators and confidence intervals for causal quantities are derived under rigorously-articulated statistical assumptions, and we develop accurate, efficient algorithms for computing these quantities. The performance of causal inference is then examined over broad parameter ranges in simulations of increasing complexity, ranging from point process models to adaptive exponential integrate-and-fire (AdEx) neuron models.

We also use simulations to examine the correspondence between the causal metrics for monosynaptic interaction developed here and the probabilities of causation, as developed in Tian and Pearl [1] and elsewhere, which quantify the necessity and sufficiency of causation probabilistically. While the correspondence is studied in a setting that relies on strong idealizations, variations that are more finely tuned to experimental work that incorporates system-specific constraints might use an analogous correspondence to test the model's assumptions *in vivo* or calibrate its free parameters via stimulation. Toward this goal, we use a biophysically detailed opsin model to simulate such an experiment *in silico*. Using a theoretically motivated stimulation paradigm, we demonstrate that common input correlations might be difficult to disentangle from common causal influences with current experimental technologies, motivating future research on that point.

2 Preliminary considerations and general architecture

We begin with a general discussion of causal models to facilitate a uniform comparison between models and simulations instantiated at different levels of abstraction. As has been widely discussed, the key motivation for explicitly modeling causation is to distinguish association from causation. An intuitive model for doing so can be described by potential outcomes [13, 14, 15, 16]. We write $Y^{(X=x)}$ as the 'potential outcome' of the random variable Y if the variable X is 'forced' to take the value x. The random variables X, Y, as well as those of the form $Y^{(X=x)}$ are presumed defined on a common probability space. We distinguish observational from experimental trials in this way. In experimental trials, the behavior of an agent external to the system (i.e., an agent that intervenes on the system) is explicitly-modeled; in observational trials, there is no such intervention. For example, in a drug efficacy trial, let $\{X_k = 0\}$ represent the event that patient k takes the treatment and let $\{X_k = 1\}$ represent the event that patient k takes a placebo. Then, if Y_k represents the measured outcome (mortality, for example) for patient k in an observational trial, then $Y_k^{(X_k=1)}$ represents the measured outcome for patient k in an experimental trial, such as a randomized control trial (RCT), in which patient k has been assigned to take the treatment by a mechanism or agent external agent intervened on the system?' and are used to define causal relations. In this case, the causal effect of the drug on the measured outcome, for patient k, is $Y_k^{(X_k=1)} - Y_k^{(X_k=0)}$. It is commonly pointed out that the challenge of causal inference is that one of $Y_k^{(X_k=1)}$ and $Y_k^{(X_k=0)}$ is unobservable. The potential outcomes notation makes it straightforward to demonstrate that an RCT is designed to infer average causal effects across a population, e.g., $\mathbb{E} \left[Y_R^{(X_R=1)} - Y_R^{(X_R=0)} \right]$, where R is a patient chosen in a simple random sample. (The latter demonstration

A *constructive* way to model potential outcomes is by explicitly modeling interventions in terms of how a structural causal model [17] is simulated. In this approach, the structure of a system is modeled via the relationships among its variables, specified by a set of functions, as in a dynamical system. This set of functions can be put in correspondence

with a directed graph by associating each variable with a vertex: a source vertex and a target vertex has a directed edge if one of the functions has the source in its domain and the target in its range. It is required that the directed graph is acyclic, which is equivalent to requiring that there is a consistent (sequential) method of simulating the system, whose ordering respects the graph. The *background* variables (noise variables) – those variables whose corresponding vertices do not have incoming edges – are instantiated as independent random variables and represent the influence of the world external to the system, in the absence of interventions. We can then think of the simulation as a closed system. Random variables can be sampled by simulation. The background variables are sampled as noise terms. Probability propagates via the functional relationships to induce a joint probability distribution on the entire system of variables. This joint distribution specifies all probability distributions (conditional and marginal distributions) of interest. All such probability distributions can then, in principle, be estimated by simulation. This is a more or less standard probabilistic point of view. The language of association is the language of conditional distributions. The association between random variables X and Y describes the distribution P(X|Y); there is no association if P(X|Y) = P(X).

We can describe *interventions* explicitly in the sequential simulation just identified. In this description, intervening on some variables means explicitly resetting their values before the functions that call them are evaluated (in the sequential method of simulation). These resets are the interventions; interventions model the action of agents external to the system. The random variables sampled in simulating the intervened system are potential outcomes. The $Y^{(X=k)}$ encodes the outcome for variable Y but in the *intervened* system in which X is reset to the value k before its use in function calls. As before, the probabilities propagate via the functional relationships and the interventions to induce a new joint probability distribution on the entire system. This joint distribution specifies all probability distributions of interest in the intervened system. All such probability distributions can again, in principle, be estimated by simulation. Causal language can then be understood as a vocabulary for discussing how interventions modify probabilities of interest. Pearl [17] uses the term $do(\cdot)$ to specify probability distributions for intervened systems. P(Y|do(X = x), Z) represents the conditional distribution P(Y|Z) when the system is intervened upon by assigning random variable X to the value x, where assigning is taken in the sense of 'resetting' above. Thus P(Y|do(X = x), Z) is another way of writing $P(Y^{(X=x)}|Z^{(X=x)})$. In what follows, we use either notation freely, for convenience.

Where does this language – in which interventions by agents are explicitly modeled – improve upon more familiar statistical modeling? A textbook example is "Simpson's Paradox", a phenomenon where the observed statistical association between variables in a population is the opposite of that observed within each subgroup of a partition of the population (see [16, 18] for more explanation). In neurophysiology, perhaps the most familiar object used in the laboratory for synaptic inference is the cross-correlogram (CCG). To motivate the framework of this article, we construct examples where the CCG and presynaptic ACG [19, 20] are insufficient for causal inference. The CCG hides information about time-dependent background correlations and presynaptic bursts that produce associations in the CCG that are not causal. The situation is analogous to *Simpson's Paradox* where statistical associations must be isolated in strata of confounders that have specific properties. The CCG collapses over these strata of the confounders. Figure 1 illustrates these ideas in some point process examples, with a detailed mathematical explanation of the simulations given in Appendix 6.1.



Fig 1: Toy examples of confounding in monosynaptic interactions. Simulations for examples are explained in detail in Appendix 6.1. A: A CCG from a hippocampal pyramidal cell and interneuron *in vivo* from the study of English et al. [2]. The pyramidal neuron (hypothesized to be presynaptic) spontaneously spikes (black) or spikes in response to experimental juxtacellular stimulation (green). Data like these motivate timescale separation assumptions. B: CCG from Example 2 for Situation A where $\epsilon_A = 0.2$, $\lambda_{R,A} = 0.2$, $\lambda_{T,A} = 0.7$. Left panel is a simulation, right panel is analytic for all rows. C: Situation B in Example 2 where $\epsilon_B = 0.8$, $\lambda_{R,B} = 0.8$, $\lambda_{T,B} = 0.1$. Note the observed CCGs are the same in B and C but the causal stories generating the data are quite different. D: Situation A in Example 3 E: Situation B in Example 3. For D-E the parameters are $\sigma_A^2 = 10 \text{ ms}$, $\sigma_B^2 = 90 \text{ ms}$, $\sigma_s^2 = 2.5 \text{ ms}$, d = 4 ms, $\alpha = 0.005$, $\omega = 20 \text{ Hz}$. The observed CCGs in the left panels of D and E are quite different but in the right panels the number of causal spikes (the areas of the green regions) are identical and the presynaptic trains, and thus their ACGs, are also identical.

Symbol	Description
$(A_k)_{k\in\mathbb{Z}^*}$	An ordered sequence $(A_0, A_1,)$
\boldsymbol{A}	A set of numbers $\{A_1, A_2,\}$
$ec{A}$	A matrix
\vec{a}	A vector
A	The cardinality of a set A or A 's absolute values if A is a scalar
$\mathbb{1}\{A\}$	The indicator function of the event A
$oldsymbol{A}^{[n]}$	The set of all subsets of cardinality n from a set \boldsymbol{A}
$oldsymbol{X},oldsymbol{Y}$	Generic sets of points (e.g., spike times) often reused
$N_A(\boldsymbol{X})$	Number of spikes from X in the temporal regions A
$\gamma(t)$	The Δ coarse temporal interval containing t
$\Gamma(\mathbf{X})$	An abbreviation for $(N_{\gamma(\Delta k)}(X))_{k\in\mathbb{Z}^*}$
$S(\mathbf{X}, \delta, \tau)$	The union of all δ length intervals centered around the τ shifted elements of X
δ, τ, Δ	The Dime date function
O_d V (x)	The pirac delta function \mathbf{X} is a site to be \mathbf{X} if a solution \mathbf{X} is found to be
Y ^(w) DDIT	The reference heaterward interaction and target quarter as respectively.
$\mathbf{n}, \mathbf{D}, \mathbf{I}, \mathbf{I}$	The number of interactions caused by a superse
$a(\mathbf{R} x)$	The conditional probability some $x \in \mathbf{B}$ will be found in $S(\mathbf{B} \ \delta \ \tau)$
$q(\mathbf{R}, x)$ $q(\mathbf{R}, \mathbf{T}(\mathbf{R}))$	The collection $(q(\mathbf{R}, x))$
G	The set of reference spikes with no target spike found in their interaction regions
K	An index set for T
J	The indices of K labeling which elements of T equal B when $\theta_{even} > 0$
L	The indices of K labeling synchronous elements of T when $\theta_{sum} > 0$
$oldsymbol{U}_h,oldsymbol{V}_h$	Indices of <i>L</i> corresponding to two limiting cases for the hypothesis $\theta_{syn} = h \ge 0$
α	One minus the confidence level
$oldsymbol{J}_h^-,oldsymbol{J}_h^+$	Hypotheses for J given the hypothesis $\theta_{syn} = h \ge 0$
Z	Candidate points that might be near inhibitory events $T^{(\emptyset)} \setminus T^{(R)}$
$ ilde{K}$	An index set for Z
$ ilde{J}$	The indices in $ ilde{K}$ for points in $oldsymbol{Z}$ near $oldsymbol{T}^{(\emptyset)} \setminus oldsymbol{T}^{(oldsymbol{R})}$
$ ilde{m{J}}_h^-, ilde{m{J}}_h^+$	Hypotheses for \tilde{J} given the hypothesis $\theta_{sun} = h$
$\chi(\ell, oldsymbol{X}, oldsymbol{Y})$	The sample cross-correlation function between spike trains X and Y
$\mathcal{C}(oldsymbol{R},oldsymbol{T},lpha)$	$1 - \alpha$ confidence interval for $\theta_{syn} \ge 0$

Table 1: Description of symbols and notation in monosynaptic causal inference model

3 Monosynaptic causal inference model

3.1 Formulation for primary model

Let X be a finite set of spike times for an experiment of duration D. For any $A = \bigcup_i [a_i, b_i)$, let $N_A(X) := |X \cap A|$, termed the *increment* of the point process X in A [21, 22]. Assuming the time origin is randomized in the experimental sense, implicitly define a partition of \mathbb{R}^+ into Δ length intervals with a function that for any time point $t \in \mathbb{R}^+$ retrieves the unique coarse interval containing t,

$$\gamma(t) \coloneqq \left\{ I \in \{ [k\Delta, k\Delta + \Delta) : k \in \mathbb{Z}^* \} : t \in I \right\}.$$
⁽¹⁾

Thus the k-th interval will frequently be written by taking t at its left endpoint, $\gamma(k\Delta)$. However, in future sections, it will be equally convenient to access interval k by letting t equal any spike time in interval k. Abbreviate the sequence of spike counts for some spike times X in the coarse intervals with the special symbol

$$\boldsymbol{\Gamma}(\boldsymbol{X}) \coloneqq (N_{\gamma(k\Delta)}(\boldsymbol{X}))_{k \in \mathbb{Z}^*} = (N_{[0,\Delta)}(\boldsymbol{X}), N_{[\Delta,2\Delta)}(\boldsymbol{X}), \dots)$$
(2)

and for any X denote the union of all δ length intervals centered around the τ shifted elements of X as

$$S(\boldsymbol{X}, \delta, \tau) \coloneqq \bigcup_{x \in \boldsymbol{X}} \left\{ s : |x + \tau - s| \le \frac{\delta}{2} \right\}$$
(3)

where the second two arguments will often be suppressed from the notation when the context is clear, i.e., S(X). Let $Y^{(x)}$ be the potential outcome of a spike train Y if a spike train X is forced to spike at a set of times x with otherwise fixed background conditions. That is, introduce the counterfactual notion that $Y^{(x)}$ would have been the set of times Y spiked if X had been x [15]. As previewed above, we will work with a reference spike train, R, and a target train, T, and the scientific question is to ask if R acts on T with a monosynapse. Further suppose the target train is constructed from latent point processes B and I, termed background events and interactions, respectively. A causal model is then defined with the deterministic relation,

$$\forall \boldsymbol{r}, \boldsymbol{T}^{(\boldsymbol{r})} \coloneqq \left\{ \begin{array}{l} \boldsymbol{B} \cup \left(\boldsymbol{I}^{(\boldsymbol{r})} \setminus \bigcup_{r \in \boldsymbol{r}} \{ S(r) : N_{S(r)}(\boldsymbol{B}) > 0 \} \right), \text{ excitatory model} \\ \boldsymbol{B} \setminus \bigcup_{r \in \boldsymbol{r}} \{ S(r) : N_{S(r)}(\boldsymbol{I}^{(\boldsymbol{r})}) > 0 \}, \text{ inhibitory model} \end{array} \right.$$
(4)

where $T^{(r)}$ references the intervention $do(\mathbf{R} = \mathbf{r})$. Notice, by construction, the background events \mathbf{B} are invariant to any action $do(\mathbf{R} = \mathbf{r})$. While one might argue this is a strong simplification, it is appropriate to compare it to the assumption that smooth features in the CCG are non-causal [23] which, for Poisson-based models, is necessarily a subset of the simplification just made (Figure 1C-D). We will be concerned with estimation of the parameter $\theta_{syn} = N_{S(\mathbf{R})}(\mathbf{T}) - N_{S(\mathbf{R})}(\mathbf{T}^{(\emptyset)})$ where $\theta_{syn} > 0$ for excitatory interactions, $\theta_{syn} < 0$ for inhibitory interactions, and $\theta_{syn} = 0$ for non-interacting neurons. In the following, it will be useful to define,

$$q(\mathbf{R}, x) \coloneqq \frac{1}{\Delta} \int_{t \in \mathbb{R}^+} \mathbb{1}\{t \in S(\mathbf{R}) \cap \gamma(x)\} dt.$$
(5)

That is, $q(\mathbf{R}, x)$ is the proportion of times $t \in \gamma(x)$ that are within a distance $\delta/2$ of a point in $\{r + \tau : r \in \mathbf{R}\}$. Constructing confidence intervals for θ_{syn} will require some additional notation. First, let us set up objects that will be used for an exact excitatory confidence interval. Denoting $T_k = f_0(k)$, fix any bijective mapping $f_0 : \mathbf{K} \mapsto \mathbf{T}^{(\mathbf{R})}$ that satisfies

$$q(\boldsymbol{R}, T_1) \le q(\boldsymbol{R}, T_2) \le \dots \le q(\boldsymbol{R}, T_{|\boldsymbol{T}|}).$$
(6)

We will write

$$q(\mathbf{R}, \mathbf{T}^{(\mathbf{R})}) = (q(\mathbf{R}, T_1), q(\mathbf{R}, T_2), ..., q(\mathbf{R}, T_{|\mathbf{T}|})).$$
(7)

Let J denote the subset of K indexing the true background events, B, in the sense that J satisfies $\{T_j : j \in J \text{ and } J \subseteq K\} = B$.

With this notation fresh in mind, we define a background model using the principle of conditional uniformity, which has been motivated and developed as a canonical assumption in previous work [24, 25]. For our purposes here, the following technical definition is sufficient (see [22] for more on point processes and their characterization).

Definition 1. Conditionally uniform point process: Define $g(\mathbf{Y}, A) = |\mathbf{Y} \cap \gamma(\inf A)|$, where \mathbf{Y} is a point process and A is a subset of \mathbb{R} . A point process \mathbf{Y} is conditionally uniform, conditioned on $\Gamma(\mathbf{Y})$ and \mathbf{X} , if

$$\mathbb{P}\left(\bigcap_{k=1}^{m}\{|\boldsymbol{Y}\cap A_{k}|=n_{k}\}|\boldsymbol{\Gamma}(\boldsymbol{Y}),\boldsymbol{X}\right)=\prod_{k=1}^{m}\binom{g(\boldsymbol{Y},A_{k})}{n_{k}}\left(\frac{|A_{k}|}{\Delta}\right)^{n_{k}}\left(1-\frac{|A_{k}|}{\Delta}\right)^{g(\boldsymbol{Y},A_{k})-n_{k}},$$
(8)

if $n_k \leq g(\mathbf{Y}, A_k)$ for all $k \in \{1, 2, ..., m\}$, for any disjoint finite collection $A_1, A_2, ..., A_m$ of subsets of \mathbb{R} that satisfies: i) for each $j, A_j \subseteq [k_j \Delta, k_j \Delta + \Delta)$ for some integer k_j and ii) $\gamma(\inf A_{j_1}) \neq \gamma(\inf A_{j_2})$ whenever $j_1 \neq j_2$.

A common way of modeling point processes is with *conditional intensity functions* [21]. While the formulation just outlined does not make use of them, later we will simulate from conditional intensity function models to demonstrate this formulation is compatible. In continuous time, the conditional intensity function $\lambda_{\mathbf{X}}(t)$ for a point process \mathbf{X} is, $\lambda_{\mathbf{X}}(t) := \lim_{\Delta t \to 0} \mathbb{E}[N_{[t,t+\Delta t)}(\mathbf{X})|\mathcal{H}_t]/\Delta t$ where \mathcal{H}_t is the history of the system prior to time t.

Remark 1. In neuroscience, the term "rate" might refer to one of several ideas [12]. In the current work, we use the word rate with regard to samples of a conditional intensity function and highlight the normalization when used.

3.2 Assumptions

Assumption 1. Conditional uniformity: B is a conditionally uniform point process, conditioned on $\Gamma(B)$ and J.

Assumption 2. *Timescale separation:* For some τ , δ , and Δ , $I^{(r)} \subset S(r)$, for all r, where $\delta < \Delta$. (τ , δ , and Δ are model parameters.)

Assumption 3. *Positivity:* $0 < q(\mathbf{R}, k\Delta) < 1$, for all $k \in \mathbb{Z}^*$. Assumption 4. *Consistency:* $I = I^{(\mathbf{R})}$ and $T = T^{(\mathbf{R})}$.

Assumption *i* will be abbreviated as $\mathcal{A}.i$. There have been debates about whether $\mathcal{A}.4$ (consistency) is an assumption or axiom of causal inference [26, 27, 28]. Here, we take it as an assumption in the sense of highlighting where it is invoked or self-evident. Similarly, one can view $\mathcal{A}.3$ (positivity) as an identifiability condition, and in our case, its validity can be determined with observational data. For this reason, one could simply define θ_{syn} in terms of regions of an experiment that provide identifiable causal information. However, following the causal inference literature [29] and for full conceptual clarity, we make it an assumption which more easily accommodates an explanation of both perspectives, leaving scientists to make their own judgment within the context of specific questions. In particular, the assumption interrelates with various other issues, including choosing free parameters, which will be discussed at length. For this purpose, the following will be useful.

Definition 2. Synchrony saturation: Synchrony saturation refers to an observation of the model where A.1 (conditional uniformity) and A.2 (timescale separation) are true but A.3 (positivity) is violated. That is, we say an observation is synchrony saturated if there is an interval identified by $k \in \mathbb{Z}^*$ such that $q(\mathbf{R}, k\Delta) = 1$.

The primary motivation for A.2 (timescale separation) is empirical [3, 5, 2, 4] although this assumption will be investigated in simulations of dynamical systems in later sections. Finally, A.1 (conditional uniformity) is motivated by the observation that *in vivo* spike trains are nonstationary [30, 31], and likely rapidly-varying. Hence, distinct points in time cannot be averaged to estimate conditional intensity functions or their variants, such as the cross-correlation function [12, 24], a matter made worse by confounding. As discussed in past work [25], the particular use of *uniformity* is motivated by the fact that the uniform distribution is the maximum entropy distribution on a finite interval.

3.3 Point estimation

Perhaps the key task of causal inference is to identify confounders and adjust for them. In the assumptions just put forth, it is conceived that the processes B, R, and $I^{(R)}$ may have non-trivial correlations that confound θ_{sun} on a Δ timescale. In causal inference, adjustment often ensues by stratifying the probability of the outcome variable conditioned on the treatment variable into different levels of the confounder. Similarly, here, the key to estimation will be to stratify time into Δ length neighborhoods and perform statistical adjustments locally (i.e., in time) via conditioning on the spike counts in those intervals. Note that a quite different approach would be to use the CCG for estimation where the spiking activity has already been averaged across levels of confounding. Figure 1, and its associated examples, essentially demonstrated that the decomposition of the CCG into causal parts is an ill-posed problem since information about confounding is hidden after averaging. Notice in the formulation section no assumptions were made about R and no assumptions were made about $I^{(R)}$ except for A.2 (timescale separation). In the following theorem, an unbiased estimator is provided for θ_{syn} . This precise deconfounding of the synaptic effect comes at the cost of not modeling the time-dependent shape of the synaptic gain onto the postsynaptic neuron. Instead, A.2 (timescale separation) simply requires interactions $I^{(R)}$ to be a subset of S(R). This does not require that no shape exists, it is simply not inside the model. Another potential source of confusion is that the idealization that there exist two processes B and $I^{(R)}$ does not mean there are two levels of synaptic efficacy. Since the events **B** are invariant under all the actions $do(\mathbf{R} = \mathbf{r})$, they indeed have zero effective synaptic weight. However, every event in $I^{(R)}$ may be generated from a different state-dependent effective synaptic weight. That is, A.2 (timescale separation) is a statement about timescale, and no assumptions about synaptic gain were made or its dependence on other factors in the model.

Theorem 1. Under A.1 (conditional uniformity), A.2 (timescale separation), A.3 (positivity) and A.4 (consistency) an unbiased point estimate of θ_{sun} in the excitatory and inhibitory models Eq. (4) is given by one expression,

$$\hat{\theta}_{syn} = \sum_{k \in \mathbb{Z}^*} \frac{N_{\gamma(k\Delta) \cap S(\mathbf{R})}(\mathbf{T}) - q(\mathbf{R}, k\Delta) N_{\gamma(k\Delta)}(\mathbf{T})}{1 - q(\mathbf{R}, k\Delta)}.$$
(9)

Proof Idea: The observed (confounded) synchrony in Δ -length temporal intervals can be expressed as a function of the hidden variables for each outcome. An appropriate conditional expectation yields calculations that isolate the causal effect as a function of observational data. Linearity of expectation across temporal intervals then recovers θ_{syn} (see Appendix 6.2).

Under violations of the identifiability assumption A.3 (positivity), one can easily salvage estimates from segments of the observation from which causal information is available.

Corollary 1. Suppose synchrony saturation occurs such that $q(\mathbf{R}, k\Delta) = 1$ where $k \in X, X \subseteq \mathbb{Z}^*$. Define the set, $\nu = \bigcup_{k \in X} \gamma(k\Delta)$ and the parameter $\theta'_{syn} = N_{S(\mathbf{R})\cap\nu}(\mathbf{T}) - N_{S(\mathbf{R})\cap\nu}(\mathbf{T}^{(\emptyset)})$. Then,

$$\hat{\theta}'_{syn} = \sum_{k \in \mathbb{Z}^* \setminus X} \frac{N_{\gamma(k\Delta) \cap S(\mathbf{R})}(\mathbf{T}) - q(\mathbf{R}, k\Delta) N_{\gamma(k\Delta)}(\mathbf{T})}{1 - q(\mathbf{R}, k\Delta)}$$
(10)

is unbiased under A.1 (conditional uniformity), A.2 (timescale separation), and A.4 (consistency).

Proof Idea: The proof (not shown) is exactly as before while highlighting the use of linearity of expectation mentioned in the previous proof idea.

3.4 Confidence intervals

The intuition behind the confidence intervals proposed here can be understood by explaining a naive algorithm for computing them. The algorithm's task is to explain the monosynaptic synchrony in a spike train pair in terms of the model. Consider the classical technique of obtaining a confidence interval by inverting a hypothesis test [32]. Intuitively, a confidence interval for θ_{syn} is the set of hypotheses for which we fail to reject the null hypothesis $H_0: \theta_{syn} = j$. A naive algorithm for calculating this interval would be to start with the hypothesis $H_0: \theta_{syn} = 0$. In this case, we assume all the observed spikes in the target train arise from the process B. Taking monosynaptic synchrony as our test statistic, the test might reject the null hypothesis that is placed under the supposition that all spikes are non-causal and thus conditionally uniform (A.1). In that case, for an excitatory interval, we proceed to conduct more hypothesis tests j = (1, 2, 3, ...). Computationally, one can imagine that before each new hypothesis test, we delete synchronous target spikes from the target train, subtract away their contribution to the test statistic, and calculate the null distribution of the test statistic under the supposition that the remaining data are conditionally uniform (A.1). We continue this process until we fail to reject the null (i.e. until the observed synchrony is explained).

The key question is, which spikes should be iteratively deleted from the target spike train in this naive algorithm? As will be shown analytically, we want to iteratively select synchronous target spikes that minimize or maximize the change in the tail probability of the test statistic. Under the model's assumption, this corresponds to removing synchronous target spikes that occur when the reference neuron's spike counts are either lowest or highest. Intuitively, suppose confounding events, B, tend to occur when the reference neuron's firing is highest; the confounding synchrony will be maximal. In that case, θ_{syn} needs to be minimal to explain the synchrony. At the other extreme, if confounding events, B, tend to occur when the reference neuron's firing is lowest, the confounding synchrony will be minimized, and thus θ_{syn} needs to be maximal to explain the synchrony. This is how θ_{syn} is bounded. We will now proceed to a more formal explanation of these intervals and prove they are exact. While the naive algorithm just described was for the purpose of intuition, a more sophisticated algorithm will also be developed for implementation later.

3.4.1 Formulation and derivation for exact excitatory confidence interval

Define $L := \{l \in K : T_l \in S(R)\}$ and let U_h be any set that satisfies $U_h \subset L$ such that $|U_h| = N_{S(R)}(T) - h$ and $\max_i \{q(R, T_i) : i \in U_h\} \leq \min_i \{q(R, T_i) : i \in L \setminus U_h\}$. Similarly, let V_h be any set that satisfies $V_h \subset L$ such that $|V_h| = N_{S(R)}(T) - h$ and $\min_i \{q(R, T_i) : i \in V_h\} \geq \max_i \{q(R, T_i) : i \in L \setminus V_h\}$. U_h and V_h identify $N_{S(R)}(T) - h$ indices (specified by f_0) of the synchronous target spikes with the smallest and largest $q(R, \cdot)$ values, respectively. Then define

$$J_h^- := U_h \cup (K \setminus L) \tag{11}$$

$$J_h^+ := V_h \cup (K \setminus L). \tag{12}$$

The conditional *pmf* for $N_{S(\mathbf{R})}(\mathbf{T}) - \theta_{syn} = N_{S(\mathbf{R})}(\mathbf{B})$, conditioned on $q(\mathbf{R}, \mathbf{T})$ and \mathbf{J} , is

$$\mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{B})=n \middle| \boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J}\right) = \sum_{Q \in \boldsymbol{J}^{[n]}} \prod_{i \in Q} q(\boldsymbol{R},T_i) \prod_{k \in \boldsymbol{J} \setminus Q} (1-q(\boldsymbol{R},T_k)).$$
(13)

Let $c^{-}(q(r, t), j)$ specify a lower (conditional) critical threshold for $N_{S(R)}(B)$ in the sense,

$$c^{-}(\boldsymbol{q}(\boldsymbol{r},\boldsymbol{t}),\boldsymbol{j}) \coloneqq \max_{k} \left\{ k : \mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{B}) \le k \middle| \boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}) = \boldsymbol{q}(\boldsymbol{r},\boldsymbol{t}), \boldsymbol{J} = \boldsymbol{j} \right) \le \alpha/2 \right\}$$
(14)

$$= \max_{k} \bigg\{ k : \sum_{n=0}^{k} \sum_{Q \in \boldsymbol{j}^{[n]}} \prod_{i \in Q} q(\boldsymbol{r}, t_i) \prod_{m \in (\boldsymbol{j} \setminus Q)} (1 - q(\boldsymbol{r}, t_m)) \le \alpha/2 \bigg\}.$$
(15)

In the same sense, let $c^+(q(r, t), j)$ specify an upper (conditional) critical threshold for $N_{S(R)}(B)$,

$$c^{+}(\boldsymbol{q}(\boldsymbol{r},\boldsymbol{t}),\boldsymbol{j}) \coloneqq \min_{k} \left\{ k : \mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{B}) \ge k \middle| \boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}) = \boldsymbol{q}(\boldsymbol{r},\boldsymbol{t}), \boldsymbol{J} = \boldsymbol{j} \right) \le \alpha/2 \right\}.$$
(16)

We now develop confidence intervals for θ_{syn} .

Lemma 1. Let

$$\mathcal{D}(\boldsymbol{R},\boldsymbol{T}) = \left\{ \boldsymbol{j} : |\boldsymbol{j}| = |\boldsymbol{J}|, (\boldsymbol{K} \setminus \boldsymbol{L}) \subseteq \boldsymbol{j} \right\}.$$
(17)

Abbreviate $\mathcal{D}(\mathbf{R}, \mathbf{T})$ as \mathcal{D} . Under $\mathcal{A}.1$ (conditional uniformity)

$$J_{\theta_{syn}}^{-} \in \operatorname*{arg\,min}_{\boldsymbol{j}\in\mathcal{D}} \left\{ c^{-}(\boldsymbol{q}(\boldsymbol{r},\boldsymbol{t}),\boldsymbol{j}) \right\}$$
(18)

and

$$\boldsymbol{J}_{\boldsymbol{\theta}_{syn}}^{+} \in \operatorname*{arg\,max}_{\boldsymbol{j} \in \mathcal{D}} \bigg\{ c^{+}(\boldsymbol{q}(\boldsymbol{r},\boldsymbol{t}),\boldsymbol{j}) \bigg\}.$$
(19)

Proof Idea: The *cdf* corresponding to the critical regions can be explicitly differentiated with respect to the labeling. Proving that $J_{\theta_{syn}}^{-}$ minimizes Eq. (18) follows from contradiction. Eq. (19) is shown in the same way (see Appendix 6.2).

An elementary idea embedded in Lemma 1 above is that if $X_1, X_2, ..., X_n$ are independent Bernoulli random variables with parameters $q_1, q_2, ..., q_n$, respectively, and those parameters satisfy $q_i \ge r_i$, for all $i \in \{1, 2, ..., n\}$, then $\mathbb{P}(\sum_{i=1}^N X_i \le c)$ is a lower bound for the same tail probability derived from a sum of independent Bernoulli random variables with parameters $r_1, r_2, ..., r_n$, respectively. All of these Bernoulli random variables are presumed to be mutually independent. In fact, there is a more general result [33] which implies that the confidence intervals developed are robust to certain violations of \mathcal{A} .1 (conditional uniformity).

Lemma 2. Denote $X_1, X_2, ..., X_{i-1}, X_{i+1}, ..., X_n$ as ${}_iX$. Then suppose $X_1, X_2, ..., X_n$ are Bernoulli random variables (not necessarily independent) satisfying

$$\mathbb{P}(X_i = 1|_i X, Z) \ge p_i(Z),\tag{20}$$

for all $i \in \{1, 2, ..., n\}$, and for some random variable Z. Then

$$\mathbb{P}\left(\sum_{i=1}^{n} X_{i} \le k \middle| Z\right) \le \mathbb{P}\left(\sum_{i=1}^{n} Y_{i} \le k\right), \forall k$$
(21)

if $Y_1, Y_2, ..., Y_n$ are conditionally independent Bernoulli random variables, conditioned on Z, with parameters $p_1(Z), p_2(Z), ..., p_n(Z)$ respectively, and $X_{1:n}$ and $Y_{1:n}$ are conditionally independent, conditioned on Z.

Proof Idea: The proof is by induction (see Appendix 6.2).

First we first demonstrate that an exact hypothesis test for $H_0: \theta_{syn} = h$ follows from the previous results.

Proposition 1. Under A.1 (conditional uniformity), A.2 (timescale separation), and A.4 (consistency) $\{N_{S(\mathbf{R})}(\mathbf{T}) - h \leq c^{-}(\mathbf{q}(\mathbf{R},\mathbf{T}),\mathbf{J}_{h}^{-})\}$ and $\{N_{S(\mathbf{R})}(\mathbf{T}) - h \geq c^{+}(\mathbf{q}(\mathbf{R},\mathbf{T}),\mathbf{J}_{h}^{+})\}$ are $\alpha/2$ -level critical region for all \mathbb{P} in H_0 : $\theta_{syn} = h$. That is, for all \mathbb{P} in H_0 : $\theta_{syn} = h$,

$$\mathbb{P}\left(N_{S(\mathbf{R})}(\mathbf{T}) - h \le c^{-}(\mathbf{q}(\mathbf{R}, \mathbf{T}), \mathbf{J}_{h}^{-})\right) \le \alpha/2$$
(22)

and

$$\mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{T}) - h \ge c^{+}(\boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J}_{h}^{+})\right) \le \alpha/2.$$
(23)

Proof: Appendix 6.2.

Finally, an exact confidence interval is constructed by inverting the hypothesis tests established in Proposition 1. **Theorem 2.** [Confidence interval for θ_{syn} .] Define

$$\mathcal{C}(\boldsymbol{R},\boldsymbol{T},\alpha) := \left\{ h : c^{-}(\boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J}_{h}^{-}) \leq N_{S(\boldsymbol{R})}(\boldsymbol{T}) - h \leq c^{+}(\boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J}_{h}^{+}) \right\}.$$
(24)

Then under A.1 (conditional uniformity), A.2 (timescale separation), and A.4 (consistency) $C(\mathbf{R}, \mathbf{T}, \alpha)$ is a $(1 - \alpha)$ -level confidence interval for $\theta_{syn} \ge 0$. That is,

$$\mathbb{P}(\theta_{syn} \in \mathcal{C}(\boldsymbol{R}, \boldsymbol{T}, \alpha)) \ge 1 - \alpha.$$
(25)

Proof. By construction,

$$\{\theta_{syn} \in \mathcal{C}(\boldsymbol{R}, \boldsymbol{T}, \alpha)\} = \left\{ c^{-}(\boldsymbol{q}(\boldsymbol{R}, \boldsymbol{T}), \boldsymbol{J}_{\theta_{syn}}^{-}) \leq N_{S(\boldsymbol{R})}(\boldsymbol{T}) - \theta_{syn} \leq c^{+}(\boldsymbol{q}(\boldsymbol{R}, \boldsymbol{T}), \boldsymbol{J}_{\theta_{syn}}^{+}) \right\}.$$
 (26)

Therefore,

$$\mathbb{P}(\theta_{syn} \in \mathcal{C}(\boldsymbol{R}, \boldsymbol{T}, \alpha)) = \mathbb{P}\left(c^{-}(\boldsymbol{q}(\boldsymbol{R}, \boldsymbol{T}), \boldsymbol{J}_{\theta_{syn}}^{-}) \le N_{S(\boldsymbol{R})}(\boldsymbol{T}) - \theta_{syn} \le c^{+}(\boldsymbol{q}(\boldsymbol{R}, \boldsymbol{T}), \boldsymbol{J}_{\theta_{syn}}^{+})\right) \ge 1 - \alpha.$$
(27)

The final inequality follows, under the model, from Proposition 1.

3.4.2 Sketch of inhibitory confidence interval and computational implementation

In Section 3.1, we modeled inhibition as a process that censors the elements of B where B is given the same properties as in the excitatory case. This approach is inspired in part by Spivak et al. [19] and, because of the superposition principle of Poisson processes [22], similar models are implied by CCG-based methods with heuristic reliance on Poisson assumptions that identify inhibition via short-latency troughs in the CCG [34, 35, 23]. However, we should regard this model with much greater skepticism, as inhibitory neurons may play a greater role in regulating downstream spike timing [36, 37] and few experimental studies exist that include causal manipulations of inhibitory neurons [38] in a manner relevant to functional connectivity. Nonetheless, we will sketch an algorithm for computing a bound for inhibition, particularly since it has been hypothesized that axo-axonic cells may function to precisely censor principal cell output in vivo [38]. Since the problem has a tenuous empirical foundation, we forgo any rigorous probabilistic interpretation of the algorithm's output. Hence, we prioritize assumptions here that permit the simplest, rather than the most precise, articulation of a concept. Perhaps experimentalists who wish to provide ground truth data for this problem may find this sketch useful while designing their experiments. Varying assumptions and comparing the resulting bound with measurements taken under experimental interventions might help to refine the assumptions needed for an inhibitory model (see Section 5.1). We emphasize that, for now, this sketch is wholly supported by intuition and simulation, guided by the philosophy that mathematical precision should yield to scientific constraints, which, as just explained, are scarcely available for inhibition.

As in Section 3.4.1, we must again consider what beliefs about B reasonably explain how the data were generated for inhibition. In particular, here we will suppose some elements of B are censored in empty synchrony regions and we must calculate null distributions corresponding to those suppositions and identify limiting cases to bound $\theta_{syn} < 0$. The problem is not directly analogous to before because, in the excitatory case, the candidate hypotheses for B include all possible subsets of T. However, for inhibition, the candidate hypotheses for inhibition include all possible subsets of S(G) where $G := \{r : r \in R, N_{S(r)}(T) = 0\}$; that is, all time points in empty synchrony regions in the observed train T. It is here that we will make a significant simplification and suppose, by assumption, that censored elements of B are simply a subset of G, and compute intervals in simulation thinking of this as an approximation. Inherent in this approximation is we ignore edge effects and assume that no more than one event of B is censored per synchrony region.

Define $Z := T \cup G$. Z now represents approximate candidate locations of background events for the excitatory and inhibitory models simultaneously (no generality is lost for the excitatory case in this notation). As earlier let an index set be $\tilde{K} = \{1, 2, ..., |Z|\}$ and define a bijective mapping $f_1 : \tilde{K} \mapsto Z$ as follows. Referring to $f_1(k)$ as Z_k , let f_1 be any such mapping that satisfies,

$$\underbrace{q(\boldsymbol{R}, Z_{1}) \geq q(\boldsymbol{R}, Z_{2}) \geq \dots \geq q(\boldsymbol{R}, Z_{|\boldsymbol{G}|})}_{Z_{k} \in \boldsymbol{G} \text{ for } 1 \leq k \leq |\boldsymbol{G}|} (28)}_{Z_{k} \in \boldsymbol{T} \cap S(\boldsymbol{R}) \text{ for } |\boldsymbol{G}|+1 \leq k \leq |\boldsymbol{G}|+N_{S(\boldsymbol{R})}(\boldsymbol{T})}}.$$

Notice there are no constraints on the mapping for the largest $|\mathbf{T}| - N_{S(\mathbf{R})}(\mathbf{T}) - 1$ elements of $\tilde{\mathbf{K}}$ which f_1 maps to the (non-synchronous) points in $\mathbf{T} \setminus S(\mathbf{R})$. (A mapping such as f_1 always exists since \mathbf{G} and $\mathbf{T} \cap S(\mathbf{R})$ are mutually exclusive.) With these simplifications, we can once again imagine two limiting cases of how the data might have been generated given a hypothesis that $\theta_{syn} = h$,

$$\tilde{J}_{h}^{-} := \begin{cases} (\begin{matrix} |G| + N_{S(R)}(T) - h & |G| + |T| \\ (\bigcup & i \end{pmatrix} \cup (\bigcup & |G| + |T| \\ i = |G| + 1 & i = |G| + N_{S(R)}(T) + 1 \\ \begin{matrix} |G| + |T| \\ \bigcup & i, \text{ if } h \le 0 \\ i = |G| + h + 1 \end{matrix}$$
(29)

$$\tilde{J}_{h}^{+} := \begin{cases} \bigcup_{\substack{i=|G|+h+1 \\ i=|G|+h+1 \\ (\bigcup_{i=1}^{-h} i) \cup (\bigcup_{i=|G|+1}^{|G|+|T|} i), \text{ if } h \ge 0. \end{cases}$$
(30)

For a hypothesis of the form $\theta_{syn} = h$ we will posit the existence of some censored background events, with associated probabilities that will then need to be convolved with a function representing a proposal about the distribution of $N_{S(\mathbf{R})}(\mathbf{B})$. As before, these hypotheses are made at limiting cases where the spike counts of \mathbf{R} on Δ timescales are either minimal or maximal (this is builtin to Eqs. (11)-(12)). Noting that we give no rigorous interpretation of these probability statements, let us use the notation $(*_{i \in \mathbb{N}} \vec{v}_i)(k) = \vec{v}_1 * \vec{v}_2 * \vec{v}_3...$ to denote the convolution of many vectors where k runs over the support of the resulting vector. In particular, consider the vectors $\vec{v}_i = (1 - q(Z_i), q(Z_i))$ for $i \in \tilde{K}$ and define,

$$\tilde{c}^{-}(\boldsymbol{q}(\boldsymbol{r},\boldsymbol{z}),\boldsymbol{j}) \coloneqq \max_{k} \left\{ k : \sum_{j=0}^{k} (*_{i \in \boldsymbol{j}} \vec{v}_{i})(j) \le \alpha/2 \right\}$$
(31)

$$\tilde{c}^{+}(\boldsymbol{q}(\boldsymbol{r},\boldsymbol{z}),\boldsymbol{j}) \coloneqq \min_{k} \bigg\{ k : \sum_{j=k+1}^{|\boldsymbol{j}|+1} (*_{i \in \boldsymbol{j}} \vec{v}_{i})(j) \le \alpha/2 \bigg\}.$$
(32)

Then let an approximate bound be,

$$\tilde{\mathcal{C}}(\boldsymbol{R},\boldsymbol{Z},\alpha) := \left\{ h : \tilde{c}^{-}(\boldsymbol{q}(\boldsymbol{R},\boldsymbol{Z}),\tilde{\boldsymbol{J}}_{h}^{-}) \leq N_{S(\boldsymbol{R})}(\boldsymbol{T}) - h \leq \tilde{c}^{+}(\boldsymbol{q}(\boldsymbol{R},\boldsymbol{Z}),\tilde{\boldsymbol{J}}_{h}^{+}) \right\}.$$
(33)

While presented under distinct notation to incorporate inhibition, this set is identical to the rigorous confidence interval derived previously for excitation when $\theta_{syn} \ge 0$. Using this notation, we present two algorithms for efficient computation of these confidence intervals. The algorithms use various tricks to minimize redundant computations and leverage state-of-the-art methods for fast and accurate tail probability computations for a sum of independent random variables. A detailed description of this approach, along with its rationale, is provided in Appendix 6.4. For the reader interested in direct application, we state Algorithm 1 and 2 immediately without explanation. Algorithm 2 is the main algorithm that computes, as an example, the lower bound for θ_{syn} in the excitatory case. For clarity, Algorithm 1 is separated but repeatedly called from Algorithm 2 and houses machinery for accurately computing tail areas for sums of independent random variables.

Algorithm 1 Hybrid convolution power with shift plus sparse exceptions

Input: Scalar L_{div} for L_{div} -fold convolution power, $\vec{p} = (p_0, p_1, ..., p_N)$ probability vector to apply L_{div} -fold convolution power, $\vec{g} = (g_0, g_1, ..., g_M)$ residual probability vector, y_0 the observed test statistic 1: function SDPNT (s, \vec{p}, \vec{q}) $\kappa'_{\vec{n}}(s), \kappa'_{\vec{a}}(s) \leftarrow$ the derivatives of the CGFs of \vec{p} and \vec{q} evaluated at s 2: 3: ▷ Note: CGF stands for cumulant generating function return $L_{div}\kappa'_{\vec{n}}(s) + \kappa'_{\vec{a}}(s) - y_0$ 4: 5: end function 6: $\hat{s} \leftarrow \text{compute } s \text{ such that } Sdpnt(s, \vec{p}, \vec{q}) = 0$ 7: for all $x \in \{0, ..., N\}$ do \triangleright apply exponential tilts $\vec{p}_{\hat{s}}(x) \leftarrow \exp\left[\hat{s}x - L_{div}\kappa'_{\vec{n}}(\hat{s}) + \kappa'_{\vec{a}}(\hat{s})\right]\vec{p}(x)$ 8: 9: end for 10: for all $x \in \{0, ..., M\}$ do $\vec{g}_{\hat{s}}(x) \leftarrow \exp\left[\hat{s}x - L_{div}\kappa_{\vec{p}}'(\hat{s}) + \kappa_{\vec{q}}'(\hat{s})\right]\vec{g}(x)$ 11: 12: end for 13: $N_{\vec{b}} \leftarrow 2^{\lceil \log_2(K(N-1) + (M-1) + 1) \rceil}$ 14: $\vec{p}_{\hat{s}} \leftarrow Concat(\vec{p}_{\hat{s}}, \vec{0}_{L_{\vec{v}}})$ pad with zero vector of dimension $L_{\vec{p}} = N_{\vec{b}} - (N+1)$ 15: $\vec{g}_{\hat{s}} \leftarrow Concat(\vec{g}_{\hat{s}}, \vec{0}_{L_{\vec{q}}})$ where $L_{\vec{g}} = N_{\vec{b}} - (M+1)$ 16: $\vec{b}_{\hat{s}} \leftarrow D^{-1}((D\vec{g}_{\hat{s}}) \odot (D\vec{p}_{\hat{s}})^{\odot L_{div}})$ ▷ computed via FFT, IFFT $\triangleright D, D^{-1}, \& \odot$ are the DFT, IDFT, and pointwise product respectively 17: 18: for all $x \in \{0, ..., N_{\vec{h}}\}$ do \triangleright reverse tilt $\vec{b}(x) \leftarrow \exp\left[\kappa_{\vec{a}}'(\hat{s}) + \kappa_{\vec{a}}'(\hat{s}) - \hat{s}x\right]\vec{b}_{\hat{s}}(x)$ 19: 20: **end for** 20: end for 21: pval $\leftarrow \sum_{i=y_0}^{N_{\vec{b}}} b_i$ Output: pval

Algorithm 2 Coarse-to-fine lower confidence bound computation

Input: Spike trains R and T, α (1-desired confidence level), z_0 (the measured value of $N_{S(R)}(T)$), L_{div} (hyperparameter)

- 1: $(q(Z_1), q(Z_2), ...) \leftarrow$ compute from R & T for an arbitrary $G \in \mathcal{G}$
- 2: Implement BinarySearch for the case \tilde{J}_h^+ on $h \in (1, 2, ..., z_0)$ with search query: $h^* = \min_h \left\{ h : \left[f(h) + \mathbb{1}\{\sum_{i \in \tilde{J}_h^+} q(Z_i) > y\} \left(1 - 2f(h)\right) \right] > \alpha/2 \right\}$
- where $f(h) = \exp\left[-\sum_{i \in \tilde{J}_{h}^{+}} q(Z_{i}) + y + y \ln(\frac{1}{y} \sum_{i \in \tilde{J}_{h}^{+}} q(Z_{i}))\right] \& y = z_{0} h$
- 3: $CF_L \leftarrow h^*$ the result of BinarySearch
- 4: $CF_U \leftarrow h^{**}$ the result of analogous BinarySearch for the case \tilde{J}_h^-
- 5: $\vec{u} = (u_1, u_2, ...) \leftarrow \{q(Z_i) : i \in J^+_{CF_{U}}\}$
- 6: for all $i \in \{1, 2, ..., \dim(\vec{u})\}$ do 7: $n_i \leftarrow \sum_{j \in \tilde{J}^+_{CF_{IJ}}} \mathbb{1}\{q(Z_j) = u_i\}$
- 8: $m_i \leftarrow \lfloor n_i / L_{div} \rfloor$
- 9: $w_i \leftarrow n_i - m_i L_{div}$
- 10: end for
- 11: $\vec{p} \leftarrow \text{apply DC}$ to binomials with $n = m_i \& p = u_i$ for $i \in \{1, 2, ..., \dim(\vec{u})\}$
- 12: $\vec{a} \leftarrow \text{apply DC}$ to binomials with $n = w_i \& p = u_i$ for $i \in \{1, 2, ..., \dim(\vec{u})\}$
- 13: Implement BinarySearch on $h \in (CF_L, CF_L + 1, ..., CF_U)$ with search query: $h^{***} = \min_{h} \left\{ h : \mathbb{P}(N_{S(\mathbf{R})}(\mathbf{T}) - h \ge z_0 - h | (q(X_1), q(X_2), ...), \mathbf{J} = \mathbf{J}_h^+) > \alpha/2 \right\}$ where given each h tested the tail probability is computed via sub-steps: $(13.1) y_0 \leftarrow z_0 - h$ (13.2) $\vec{q} \leftarrow \text{apply DC}$ to \vec{a} and the distributions w/ success prob. $q(Z_i)$ for $i \in \tilde{J}_h^+ \setminus \tilde{J}_{CFu}^+$ (13.3) tail probability \leftarrow Pass \vec{p} , \vec{g} , y_0 , and L_{div} to Algorithm 1 14: Lower confidence bound $\leftarrow h^{***}$ the result of BinarySearch
- Output: Lower confidence bound

 \triangleright Unique probabilities for $i \in \tilde{J}^+_{CFu}$

4 Causal inferences in simple simulated systems

In the following sections, we will study the monosynaptic causal inference model in simulation. The non-parametric nature of the model possesses some features that may seem foreign to those accustomed to modeling point processes with objects such as conditional intensity functions and generalized linear models (GLMs). For example, to define a background timescale, we partitioned time into arbitrarily phased intervals, each of duration Δ . The model makes no use of conditional intensity functions and few assumptions were made about the interaction process $I^{(R)}$. At first, we will demonstrate the features just mentioned are appropriate in a conditional intensity model ensuring the process follows the monosynaptic causal inference model's assumptions only at the level of analogy. Later, we will test inferences in neural dynamical systems where it is much less clear if the assumptions are appropriate, and thus, various validations will be necessary.

4.1 Causal inferences in a conditional intensity function model

The conditional intensity model here will exhibit rapid nonstationarities with random phases, but the nonstationary fluctuations will have timescales with Δ as a lower bound. This will suggest the idealized construction of $\gamma(t)$ with fixed Δ is scientifically appropriate. One can also find theoretical arguments supporting this construction in past work [39]. The conditional intensity functions are made smooth by generating them from normalized Ornstein-Uhlenbeck processes with non-stationary means that define the background timescales. The synaptic coupling is generated by convolving the presynaptic spike trains with a truncated exponential kernel. The notion of a synapse with finite strength at infinite decay time will be abstracted out now and will naturally reenter the study of dynamical systems models later. In addition to confounding background excitability fluctuations, here, the synaptic kernel has a private background excitability fluctuations, confounding the causal relationship. These are all ways to confound the relationship between spike trains while maintaining a type of separation of timescale.

Motivated by the observation that populations of neurons have downstates and upstates [40, 41], let us introduce nonstationary fluctuations between neurons on coarse timescales with varying degrees of skew. Furthermore, consider that monosynaptically-interacting neurons might be at different phases of an oscillation in the local field potential (e.g., hippocampal gamma oscillations [42]). This example is concrete and empirical, but it is easy to imagine complex neural computations generate other types of confounding in pairwise interactions. Coarse timescale nonstationarities with complex dependence structure can then generate confounding in the CCG even when a separation of timescales assumption is true.

To generate some limiting cases in simulation, consider a sequence of multivariate skew random variables,

$$\vec{\boldsymbol{M}}_{n,k} = \begin{bmatrix} m_{0,k} \\ m_{1,k} \\ \vdots \\ m_{n-1,k} \end{bmatrix} \sim 2\phi_n(\vec{m}_{n,k}; \vec{\boldsymbol{\Omega}}) \Phi(\vec{\alpha}^\top \vec{m}_{n,k}), \vec{m}_{n,k} \in \mathbb{R}^n, k \in \mathbb{Z}^*$$
(34)

where $\phi_n(\vec{m}_{n,k}; \vec{\Omega})$ is a zero mean *n*-dimensional normal density with correlation matrix $\vec{\Omega}$, $\Phi(\cdot)$ is the standard normal distribution function, and $\vec{\alpha}$ is an *n*-dimensional vector that controls skewness [43]. Let $L_i \sim U(\Delta, a\Delta)$ ms, a > 1 and partition \mathbb{R}^+ into contiguous disjoint intervals $\{[x_{k-1}, x_k)\}_{k=0}^{\infty}$ such that $x_k = x_{k-1} + L_k$ and $x_0 = 0$. Associate with each L_k a sample from the multivariate skew distribution, $\vec{M}_{i,k}$, with dimension n = 3 and define a set of background excitability functions as,

$$b_i(t) = \sum_{k \in \mathbb{Z}^*} m_{i,k} \mathbb{1}\{t \in [x_{k-1}, x_k)\}, \text{ for } i = 0, 1, 2.$$
(35)

In simulation, we thus imagine that the excitability of the reference neuron, $b_0(t)$, target neuron, $b_1(t)$, and dendritic compartment of a synapse between them, $b_2(t)$, might have arbitrary confounding fluctuations on coarse timescales generated from a multivariate skew. That is, $b_0(t)$, $b_1(t)$, $b_2(t)$ might be skewed - i.e., rare up or down states [44] - and these states may have positive or negative correlations with each other.

Define τ_d as a conduction delay, τ_s as a phenomenological synaptic relaxation time, and $v(\tau) = \exp(-\tau/\tau_s)\mathbb{1}\{0 \le \tau < \tau_{mx}\}$ as a synaptic kernel zero everywhere but $\tau \in [0, \tau_{mx})$. The model is then defined by the conditional intensity

functions,

$$\lambda_R(t)|U_0(t) = \rho_0 U_0(t)$$
(36)

$$\lambda_T(t)|U_1(t), U_2(t), do(\boldsymbol{R} = \boldsymbol{r}) = \rho_1 U_1(t_j) + \epsilon U_2(t) \int_{-\infty}^{\infty} \upsilon(t - \tau - \tau_d) \sum_{r \in \boldsymbol{r}} \delta_d(t - r) d\tau$$
(37)

where ρ_0, ρ_1 are normalization factors, ϵ is a coupling constant,

$$\tau_I \frac{dI_i(t)}{dt} = -I_i(t) + b_i(t) + \sigma_I \sqrt{2\tau_I} \xi_i(t), \xi_i(t) \sim \mathcal{N}(0, 1)$$
(38)

is used to obtain $U_i(t) = a_i I_i(t) + b_i$, a_i and d_i are chosen to min-max normalize $U_i(t)$, and σ_I^2 and τ_I are the variance and timescale of the smoothing agent, respectively. For clarity, we restate the causal interpretation of the model in Eqs. (36)-(37) in relation to the general model stated in Section 2 and Example 1. The system is governed by a common probability space and under the subcausal model induced by $do(\mathbf{R} = \emptyset)$ Eq. (37) reduces to its first term.

The model is simulated in Figure 2. Figure 2A depicts the coupled conditional intensity model in a cartoon fashion. Each point in Figure 2B represents a distinct simulation as follows. For each simulation, the Vine Beta method, with its parameter fixed to 0.1, is used to generate a random covariance matrix with strong correlations [45] scaled to set $\vec{\Omega}$ and we sample $\vec{\alpha} \sim U([0, 100(2a-1))^n])$ where $a \sim Be(0.5)$. We then sample the sequence of independent and identically distributed (within but not across simulations) multivariate skew vectors $\vec{M}_{n,k}$ and construct the background excitability functions. As mentioned previously, the motivation for the particular form and parameters of the simulation is to induce confounding background fluctuations between the simulated neurons, with skewed up and down states, and confounded state-dependent synaptic efficacy as well. Normalization factors ρ_0 and ρ_1 are all sampled so that the average firing rates of all spike trains in the absence of coupling are uniformly distributed between 50 and 200 spikes/second. Finally, spike are simulated from $\lambda_R(t)|U_0(t), \lambda_T(t)|U_1(t), U_2(t), R$, and $\lambda_T(t)|U_1(t), U_2(t), do(\mathbf{R} = \emptyset)$. From these we obtain the spike trains \mathbf{R}, \mathbf{T} , and $\mathbf{T}^{(\emptyset)}$. We then compute the simulated ground truth θ_{syn} , along with $\hat{\theta}_{syn}$. Here we assume τ_{mx} is known, and so the statistical free parameter δ was made one time bin larger. We also assume knowledge of Δ and hence set the statistical parameter equal to the value used to simulate the process (which is a lower bound). For each simulation a point estimate is then computed with Eq. (9) and 95% confidence intervals are computed with Algorithms 1 & 2. For 101 simulations the empirical coverage probability of the confidence intervals is 0.98.

4.2 Mapping the statistical model onto a dynamical system

Thus far, it has been assumed that a postsynaptic spike train is derived from a latent mixture of background events, B, and interactions, $I^{(R)}$. It should be clear from the assumptions and from the demonstration in a conditional intensity model that the division into two classes does not at all mean a constant synaptic weight; the model houses conditional intensity models where events analogous to $I^{(R)}$ might arise from different state-dependent probabilities. Rather, the idealization lies in positing that there exist some events B that may be assumed to have zero effective causal weight, and both classes have timescale assumptions that make θ_{syn} identifiable. This might be described as a type of causal coarsening. However, despite its clear merits in terms of analytic tractability, the clean division of the postsynaptic train into two classes gave rise to three free parameters δ , τ , and Δ . By assumption, all interaction points are confined to be members of the set $S(\mathbf{R}, \delta, \tau)$ whereas Δ defines the background timescale. While constructed from conditional intensity functions, the simulated model of the previous section more or less ensured by construction that causal spikes would be confined to a set $S(\mathbf{R}, \delta, \tau)$ and non-causal spikes would possess no temporal structure for timescales smaller than Δ where $\delta < \Delta$.

It is now natural to challenge aspects of that idealization in some settings even more foreign to the one in which the model was derived. A sensible choice is dynamical neuron models, which well-capture features of cortical neurons [46] and where ground truth causal information is available by recycling the concept of frozen noise to be applied to stochastic input currents. We first must ask if a τ and δ can be chosen such that the simple division of a postsynaptic train into events B and $I^{(R)}$ might approximate the causal action of a presynaptic input through a dynamical system. The second question to consider is if, given knowledge of δ and τ , a Δ can be chosen to recover causal counterfactual quantities in the midst of confounding.

We do not provide a method to choose δ , τ , and Δ from first principles with observational data and are skeptical that the task is even possible, particularly in the case of Δ . Rather, we regard them as free parameters in the physicist's sense. So the task here is bent toward understanding the qualitative mapping of these free parameters onto some mechanistic features. As such, in contrast to our previous work [7], as a matter of interpretation and robustness, we study commonly



Fig 2: Point process demonstration. A: A cartoon depiction of the conditional intensity model described in Eqs. (36) - (37). $\lambda_R(t)$ and $\lambda_T(t)|do(\mathbf{R} = \emptyset)$ are denoted in black. The synaptic excitability function $b_2(t)$ is not depicted but modulates the height of the synaptic gain (the green gain on top of $\lambda_T(t)|do(\mathbf{R} = \emptyset)$). To generate confounding, the coarse timescale amplitudes of background rates and synaptic efficacy are generated from a multivariate skew distribution with strong correlations. The synaptic gain modulation is represented as the gray colormap. **B**: One hundred and one simulations from the model of Eqs. (36) - (37) with empirical coverage probability 0.98. Point estimates are always shown, whereas confidence intervals are not drawn if the null hypothesis $H_0 : \theta_{syn} = 0$ fails to reject.

used dynamical mechanisms (e.g., LIF, EIF, AdEx) throughout future sections. This will demonstrate that the statistical framework's validity is also a matter of qualitative considerations. For example, when one asserts that the synaptic process $I^{(R)}$ is fast, the word *fast* clearly has meaning only relative to the background timescale and hence what is of true theoretical interest is the ratio of the effective synaptic and background input timescales. The reader should make judgments about the quantitative plausibility from real data, for example, by consideration of the fine-timescale effects studied by English et al. [2], for which standard integrate-and-fire type models might be insufficient [7].

4.2.1 System of feedforward leaky integrate-and-fire (LIF) neurons

Consider first a system of standard leaky integrate-and-fire (LIF) neurons with instantaneous conduction delay driven by background input currents and a synaptic conductance from a presynaptic neuron (i = 0) into a postsynaptic neuron (i = 1),

$$C_m \frac{dV_i}{dt} = -g_l (V_i - E_l) - g_s g_0 (V_i - E_{syn}) \mathbb{1}\{i = 1\} + I_i, \text{ for } i = 0, 1$$
(39)

$$\tau_{syn}\frac{dg_s}{dt} = -g_s + \sum_{r \in \mathbf{R}} (1 - g_s)\delta_d(t - r) \tag{40}$$

where V_i is the voltage (mV) of neuron i, g_l and g_s are the leak and synaptic conductances (mS/cm²), E_l and E_{syn} are the leak and synaptic equilibrium potentials (mV), C_m is the specific capacitance (μ F/cm²), g_0 and τ_{syn} are the peak synaptic conductance (mS/cm²) and timescale (ms), and I_i is the background input (μ A/cm²) to neuron i.¹ If at time t_0 , $V(t_0) = V_T$, a spike is tabulated followed by the reset condition $\lim_{\{\epsilon \to 0:\epsilon > 0\}} V(t_0 + \epsilon) = E_l$ and clamped there for a refractory period of τ_r . We map the system Eqs. (39)-(40) onto the monosynaptic causal model by identifying $\mathbf{R} = \{t : V_0(t) = V_T\}$, $\mathbf{T} = \{t : V_1(t, \mathbf{R}) = V_T\}$, and $\mathbf{T}^{(\emptyset)} = \{t : V_1(t, do(\mathbf{R} = \emptyset)) = V_T\}$ where the last line refers to the trajectory $V_1(t)$ under the deterministic modification of Eq. (40) $do(\mathbf{R} = \emptyset) \implies g_s = 0, \forall t \in \mathbb{R}$ in the sense that the voltage trajectories may change but the I_i remain constant (i.e., frozen) for all realizations of the stochastic input current. In other words, unlike in Remark 1, here, the probability space is defined directly over the background input current, and the dynamical system determines the functional relationship between the reference and target train. To study the interpretation of the monosynaptic causal inference model in terms of dynamical mechanisms, here we simplify the form of the background model while maintaining confounding common input,

$$I_i = U_i + U_2, \text{ for } i \in \{0, 1\}$$
 (41)

$$\tau_{I,i} \frac{dU_i}{dt} = U_i + \mu_i + \sigma_i \sqrt{2\tau_{I,i}} \xi_i(t), \text{ for } i \in \{0, 1, 2\}$$
(42)

where $\xi_i(t) \sim \mathcal{N}(0, 1)$ and U_0, U_1, U_2 are all independent Ornstein-Uhlenbeck processes with means μ_i , variances σ_i^2 , and timescales $\tau_{I,i}$. It is not entirely clear from Eq. 40 where one ought to look for spiking events that may be well-attributed to a synaptic process like $I^{(R)}$. Certainly, we may first assume that they tend to occur after and not before the spikes R. Furthermore, we would imagine that their tendency to occur decreases as a function of distance from the spikes R given the exponential decay model of synaptic conductances. These are elementary and routinely applied assumptions but do not yet appeal to causal inference concepts. Let us make the simplification of instantaneous biophysical conduction delay, $\tau_d = 0$, and measure counterfactual spike counts in the spectrum of sets, $\{S(R, 2\delta, \delta) : \delta \in [0, \infty)\}$. For this purpose define the function,

$$g(\delta) = N_{S(\mathbf{R}, 2\delta, \delta)}(\mathbf{T}) - N_{S(\mathbf{R}, 2\delta, \delta)}(\mathbf{T}^{(\emptyset)}).$$
(43)

Assume for now $g(\delta)$ is monotone and that

$$\bar{\beta} = \lim_{\delta \to \infty} |g(\delta)| \tag{44}$$

exists. Then for any $0 \le \beta_0 < 1$ let,

$$\delta_{0} = \begin{cases} \min_{\delta} \left\{ \frac{\delta}{2} : \frac{|g(\delta)|}{\bar{\beta}} = \beta_{0} \right\} & \text{if } \exists \delta \in \mathbb{R}^{+} \text{ s.t. } \frac{|g(\delta)|}{\bar{\beta}} = \beta_{0} \\ 0 & \text{otherwise.} \end{cases}$$
(45)

Intuitively, δ_0 describes how large the time interval after the reference spikes must be to capture some proportion, β_0 , of the causal difference in spike counts between the counterfactuals relative to the causal difference at some long-term value $\bar{\beta}$.

¹We continue to work in these units throughout.



Fig 3: Mapping the statistical monosynaptic causal inference model onto a leaky integrate-and-fire (LIF) system. Functions $g(\delta)$ and $N_{S(\mathbf{R},2\delta,\delta)}(\mathbf{T})$ normalized by a factor $|\mathbf{R}| \cdot dt$ for $\delta \in [0,20]$ msec are plotted in green lines and black lines respectively. Vertical lines are δ_0 for $\beta_0 = 0.95$. Different line styles correspond to three variations of one postsynaptic parameter per plot. Each biophysical parameter setting corresponds to a long simulation of 27.77 simulated hours. The left column is for an excitatory synapse ($E_{syn} = 0 \text{ mV}$), and the right is for an inhibitory synapse ($E_{syn} = -70 \text{ mV}$). A: The synaptic decay time constant $\tau_{syn} \in \{1, 5, 9\}$ msec. B: The peak synaptic conductance $g_0 \in \{0.01, 0.02, 0.03\} \text{ mS/cm}^2$. C: The membrane time constant $\tau_m \in \{6, 10, 14\}$ msec. D: The postsynaptic noise amplitude $\sigma_{I,post} \in \{0.5, 1, 2\} \mu \text{A/cm}^2$ (referred to as $\sigma_{I,post}$ in the figure but defined as σ_1 in the text).

4.2.2 Simulations validating the mapping

Figure 3 shows simulations of Eq. (39)-(40) and plots $N_{S(\mathbf{R},2\delta,\delta)}(\mathbf{T})$ and $g(\delta)$ for $\delta \in [0, 20]$ ms. Both these functions are normalized by a factor $|\mathbf{R}| \cdot dt$. Each panel displays six lines: $N_{S(\mathbf{R},2\delta,\delta)}(\mathbf{T})$ and $g(\delta)$ for three different values of a dynamical parameter. Vertical lines mark δ_0 , which is obtained by setting $\beta_0 = 0.95$ and where $\bar{\beta}$ is approximated by taking $|g(\delta)|$ at $\delta = 1000$ ms. Plots are shown for both excitatory ($E_{syn} = 0$ mV) and inhibitory ($E_{syn} = -70$ mV) synapses. In the LIF neuron, $g(\delta)$ tends to rise roughly monotonically and saturate quite quickly, although different dynamical parameters have a significant impact, particularly the synaptic timescale τ_{syn} .

Figure 3A demonstrates simulations of the LIF model system with $\tau_{syn} \in \{1, 5, 9\}$ ms. For both the excitatory and inhibitory synapses, at $\tau_{syn} = 19$, $\delta_0 \approx 2$ ms whereas at $\tau_{syn} = 9$, $\delta_0 \approx 12$ ms. That long synaptic decay times increase the timescale of causal action on the postsynaptic neuron coincides with intuition. Figure 3B shows the effect of peak synaptic conductance, g_0 , on δ_0 ; the effect on δ_0 is less dramatic than τ_{syn} . For $g_0 \in \{0.01, 0.02, 0.03\}$ mS/cm², δ_0 clusters around $\delta_0 \approx 7$ ms for an excitatory synapse and $\delta_0 \approx 6$ ms for an inhibitory synapse. Figure 3C shows the analogous plots for membrane time constant, τ_m . The effect on δ_0 is slightly more pronounced than peak synaptic conductance but still less than synaptic decay time. In the excitatory case, for $\tau_m = 6$ ms, $\delta_0 \approx 4$ ms whereas for $\tau_m = 14$ we observe $\delta_0 \approx 6.6$ ms. In the inhibitory case δ_0 is clustered around $\delta_0 \approx 6$ ms with a slight positive trend with τ_m . Figure 3D again shows the analogous plots for the postsynaptic Gaussian noise amplitude which is known to influence postsynaptic response dynamics [47]. A negative trend is seen between the postsynaptic noise $\sigma_1 \in \{0.5, 1, 2\} \mu A/cm^2$ (refered to as $\sigma_{I,post}$ in the figure) and δ_0 . No trend is detected in the inhibitory case.

Figure 4 plots monosynaptic point estimates and confidence intervals for the LIF system varying the same parameters as the previous plot. To make sensible comparisons across plots, here parameters are normalized as $g(\delta_0)/(|\mathbf{R}| \cdot dt)$, termed *causal rate* in the figure. Estimates are $\hat{\theta}'_{syn}/(|\mathbf{R}| \cdot dt)$ For each plot, different levels of causal rate are produced by varying $g_0 \in [0, 0.1] \text{ mS/cm}^2$ with eleven equally spaced values for both excitatory ($E_s = 0 \text{ mV}$) and inhibitory $(E_s = -70)$ synapses. For estimation, we must also choose a value for the statistical parameter Δ . Recall that A.2(timescale separation) requires that some τ , δ , and Δ exist such that $I^{(r)} \subset S(r)$, for all r. In simulation δ_0 is chosen such that $\beta_0 = 0.95$ to approximate this assumption. On the other hand A.3 (positivity) requires $0 \le q(\mathbf{R}, k\Delta) < 1$, for all $k \in \mathbb{Z}^*$. Once δ is chosen to be δ_0 , an observer can choose some Δ_0 and at least verify $0 \le q(k\Delta_0) \le 1$, for all $k \in \mathbb{Z}^*$ because this condition only requires access to the observed data to verify. This is not as useful as it may seem, however, because, of course, Δ is unknown, and an appropriate selection of it determines the validity of A.1(conditional uniformity) which cannot be assessed from observational data. Furthermore, A.2 and A.3 are orthogonal assumptions: one can be true while the other is false. In these simulations, the common background input timescale was chosen to be rather large ($\tau_{I,2} = 50$ ms) so that a simple heuristic might automate the choice of Δ_0 given δ_0 which strongly biases this inquiry toward assessing c. In this figure, without too much thought, we choose $\Delta_0 = \delta_0 + 4$ ms, which typically well-approximates the $\mathcal{A}.3$ (positivity) in the regimes explored here, although some misestimation arises from violations. To study causal identifiability, in all plots to follow, we use $\hat{\theta}'_{sun}$ from Corollary 1 as an estimate and always define θ_{syn} as the ground truth value.

It is worth reminding the reader that Δ is unknown, and perhaps unknowable, in these simulations. The timescales of the membrane, background input, and synapse likely interact and might even produce a statistical background timescale that is smaller than the timescales of the physiological variables involved. Similarly, the assumption that the inequality $\delta < \Delta$ can be true while satisfying the other assumptions cannot be known in the simulation and, in fact, is one of the primary motivations for testing estimation in dynamical systems models while varying physiological parameters. That is, good estimation is regarded as evidence for the fulfillment of the assumptions.

The qualitative results of Figure 4 can, for the most part, be predicted by the results of the previous figure. That is, the estimation procedure provides highly accurate estimates to the degree that the model's assumptions are approximated in the sense of the mapping proposed earlier. Figure 4A shows point and interval estimates for $\tau_{syn} \in \{1, 5, 9\}$ ms. As Figure 3A predicts, as τ_{syn} increases δ_0 also increases which puts stress both on $\mathcal{A}.1$ (conditional uniformity) and $\mathcal{A}.3$ (positivity). One ought to heed the point just made about Δ being unknown in dynamical simulations.

Even at an unrealistically small synaptic timescale of $\tau_{syn} = 1$ ms, the magnitude of the inhibitory causal rate is slightly underestimated with empirical coverage probability 0.82 for the confidence intervals. As τ_{syn} increases, both the magnitude of excitatory and inhibitory causal rates are underestimated, however, the confidence intervals behave more conservatively in this regime and have empirical coverage probability 1 for $\tau_{syn} \in \{5,9\}$ across all simulations.

In 4B we observe that the qualitative behavior of estimation with respect to synaptic timescale τ_{syn} is recapitulated for membrane timescale τ_m although to a less pronounced degree. That is, Figure 3C indicates δ_0 will increase as membrane timescale τ_m increases and accordingly in 4B the magnitude of causal rate is slightly underestimated for excitatory and inhibitory interactions as τ_m increases. A notable feature of Figure 3 was that all else being equal, increasing τ_{syn} or τ_m increased the causal rate for all δ and for the most part trended positively with δ_0 . More intuitively, as the temporal scale of causal effect (i.e., δ_0) increased more spikes were causal, naturally. But this appears not to be the case for σ_1 (termed $\sigma_{I,post}$ in the figure). For excitatory interactions, in Figure 3D increasing σ_1 increased δ_0 , however, the magnitude of causal rate decreased as δ_0 increased unlike in the case of synaptic timescale τ_{syn} and membrane timescale τ_m . Yet, like the case with τ_{syn} and τ_m , Figure 3D and Figure 4C in combination show that estimation is accurate to the degree δ_0 is made small by a small σ_1 . This all suggests, not at all in conflict with intuition, that the model's validity has some mechanistic independence so long as the causal behavior at the level of spiking abides by the formal assumptions proposed earlier.

The idea that causal effects of inhibitory synapses can be estimated from spike trains has far less existing support from *in vivo* experiments. Furthermore, here, the estimation of inhibitory synapses is only highly accurate for unrealistic parameters for inhibitory synapses [48]. For these reasons, the latter figures primarily focus on the study of excitatory interactions, except in a few idealized settings where the math quickly implies an inhibitory solution. In that case, this cautionary statement still applies. However, it should be noted that physiological parameters interact, and they are typically measured in settings where the interactions may not be present, such as *in vitro* studies. Thus one cannot rule out the possibility that, at the level of spiking, the temporal scale of causal action for inhibitory synapses is still small *in vivo*, meaning the deficiency resides in the biophysical models. However, until more basic evidence of such mechanisms exists, we must remain skeptical that the inhibitory model proposed in this study has any relevance to neuroscience.

Parameter Name	Symbol	Unit	Value/Distribution		
Cellular Properties					
Membrane Capacitance	C_m	μ F/cm ²	1		
Leak Conductance	E_1 g_1	mv mS/cm ²	-65 0.1		
Spike Threshold	V_T	mV	-50		
Voltage Reset	V_R	mV	E_l		
	7r	IIIS	2		
Synapse					
Peak Synaptic Conductance	g_0	mS/cm^2	0.04		
Synaptic Reversal Potential	E_{syn}	mV	0		
Synaptic Time Constant	$ au_{syn}$	ms	3		
Conduction Delay	$ au_d$	ms	0		
Background Input					
Input Timescale	τ_{I_i} *	ms	50		
Input Mean	μ_i *	μ A/cm 2	0		
Input SD	σ_i *	μ A/cm 2	1		

Table 2: LIF base circuit parameters

*: for $i \in \{0, 1, 2\}$

4.3 Causality and spike history in feedforward adaptive exponential (AdEx) integrate-and-fire neurons

After developing a causal inference model in Section 3.1, we proceeded to test the model in a series of numerical experiments that challenged the model's assumptions. The point process experiments of Section 4.1 challenged aspects of how we constructed the background process to account for confounding; namely, the definition of $\gamma(t)$ and $\mathcal{A}.1$ (conditional uniformity). The LIF system experiments of Section 4.2, with intrinsic dynamics and conductance-based synapses, challenged aspects of how we constructed the interaction process to account for coupling effects; namely, the definition of $I^{(R)}$ and $\mathcal{A}.2$ (timescale separation). In this final subsection, we take this challenge further and try to identify a case where the causal effect of synaptic input is perhaps more complex than a transient increase in postsynaptic spiking probability followed by exponential decay (Section 4.2).



Fig 4: Causal inferences in the LIF model as a function of various dynamical parameters. Inference of postsynaptic parameters studied in the previous figure. δ_0 is assumed to be known. The true causal rate is defined as $g(\delta_0)/(|\mathbf{R}| \cdot dt)$ and the estimate $\hat{\theta}'_{syn}/(|\mathbf{R}| \cdot dt)$ for δ_0 . For each plot, different levels of causal rate are generated by varying $g_0 \in \{0, 0.01, 0.02, ..., 0.1\}$ mS/cm² both for excitation ($E_{syn} = 0$ mV) and inhibition ($E_{syn} = -70$ mV). Across all panels, the empirical coverage probability equals 0.9894. Parameters sweep left to right. A: $\tau_{syn} \in \{1, 5, 9\}$ ms. B: $\tau_m \in \{6, 10, 14\}$ ms. C: $\sigma_{I,post} \in \{0.5, 1, 2\} \mu A/cm^2$ (referred to as σ_1 in the text).

Consider a system of AdEx model neurons [49]. As before, let a presynaptic neuron (i = 0) drive a postsynaptic neuron (i = 1),

$$C_m \frac{dV_i}{dt} = -g_l(V_i - E_l) + g_l k_{a,i} \exp\left(\frac{V_i - V_{T,i}}{k_{a,i}}\right) - g_s g_0(V_i - E_{syn}) \mathbb{1}\{i = 1\} - I_{w,i} + I_i$$
(46)

$$\tau_w \frac{dI_{w,i}}{dt} = -I_{w,i} + a_i(V_i - E_l) + \mathbb{1}\{i = 1\} \sum_{y \in \mathbf{T}^{(\mathbf{R})}} b_i \delta_d(t - y) + \mathbb{1}\{i = 0\} \sum_{r \in \mathbf{R}} b_i \delta_d(t - r)$$
(47)

$$\tau_{syn}\frac{dg_s}{dt} = -g_s + \sum_{r \in \mathbf{R}} \delta_d(t-r) \tag{48}$$

where the LIF model has been embellished with a nonlinearity with activation slope k_a and with an adaptation current I_w with subthreshold adaptation coupling parameter a and spike-triggered adaptation parameter b. A spike is triggered when V(t) obtains the value $V_T + 5k_a$ at which time the voltage V(t) is as before reset to E_l for a refractory period of τ_r . The counterfactual interpretation of the system Eqs. (46)-(48) is exactly analogous to the LIF system of Eqs. (39)-(40) as already discussed noting that under the intervention $do(\mathbf{R} = \emptyset)$ the spike-triggered adaptation trigger times become

 $T^{(\emptyset)}$. As alluded to in the previous section, here we focus on excitatory synapses only. Parameters were chosen separately for each neuron so that the presynaptic and postsynaptic cells emulate a neocortical pyramidal neuron and fast-spiking interneuron, respectively [50]. Technically, under these parameters, the AdEx model for the postsynaptic neuron reduces to the exponential integrate-and-fire neuron (EIF), as it is a special case of the former.

As earlier with the LIF model, the synaptic conduction delay is set at $\tau_d = 0$ and Figure 5A plots normalized versions of $N_{S(\mathbf{R},2\delta,\delta)}(\mathbf{T})$ and $g(\delta)$ for $\delta \in [0, 20]$ ms. The AdEx system produces an apparent non-monotonic behavior in this regime in $g(\delta)$. This observation should be examined in the context of functional connectivity methods that use the CCG as the primary object of inference. For example, Spivak et al. [19] argues that presynaptic autocorrelation can produce secondary oscillations in the CCG and should be corrected for by a deconvolution procedure. We plot the CCG from the simulation of Figure 5A in Figure 5B. The secondary oscillations seen here are characteristic of those thought to arise from presynaptic autocorrelation. Under the assumption $g(\delta)$ is monotonic, secondary oscillations in the CCG would indeed be an artifact manifesting when finely-timed presynaptic bursts coincide with finely-timed postsynaptic spikes arising causally from one of the presynaptic spikes in the burst (see Example 3). The causal postsynaptic spike then contributes to the mass of the CCG in at least two places: the large primary short-latency CCG peak [2] as well as in one of the secondary oscillations. Whether the causal postsynaptic spike arises from the first or second spike in the presynaptic burst dictates whether it contributes to the duplicate mass in the secondary oscillation residing in the region of negative or positive lag.

Here, we have observed that $g(\delta)$ is not monotonic, indicating, by this fact alone, that part of the secondary oscillations is causal and not an artifact due to duplicate mass in the CCG. Yet, the model neocortical pyramidal neuron does have regular bursting as well. To tease apart the contribution of each factor, we append another simulation to the AdEx system simulation as follows. Let us reuse the AdEx simulated presynaptic train \mathbf{R} to keep presynaptic autocorrelation constant and reuse $\mathbf{T}^{(\emptyset)}$ to keep confounding partially constant. We take the counterfactual target spike train $\mathbf{T}^{(\emptyset)}$ and add $|\mathbf{T}| - |\mathbf{T}^{(\emptyset)}|$ spikes to it via a conditional intensity model of synaptic gain, taking the union of spikes induced by that model synapse with $\mathbf{T}^{(\emptyset)}$ (this is termed "artificial synapse" for short). More precisely, define the synaptic gain function

$$\lambda_A(t)|do(\boldsymbol{R}=\boldsymbol{r}) = \epsilon \int_{-\infty}^{\infty} \upsilon_0(t-\tau-\tau_d) \sum_{r \in \boldsymbol{r}} \delta_d(t-r) dt \qquad (\text{generates spikes } \boldsymbol{I}^{(\boldsymbol{R})})$$

where making a modification from before the kernel is not truncated in the direction of positive infinity: $v_0(\tau) = \exp(-\tau/\tau_s)\mathbb{1}\{\tau \ge 0\}$. Here, τ_s is chosen to maximize the correlation of the resulting CCG with the CCG obtained from the initial AdEx simulation. ϵ is also chosen to produce approximately $|\mathbf{T}| - |\mathbf{T}^{(\emptyset)}|$ spikes (from the initial simulation) and then some interactions, $\mathbf{I}^{(\mathbf{R})}$, simulated from this gain function are randomly omitted so the number of causal spikes in the second simulation exactly equal $|\mathbf{T}| - |\mathbf{T}^{(\emptyset)}|$ from the initial AdEx simulation. The resulting CCG from the artificial synapse is also displayed in Figure 5B. Secondary oscillations persist in the CCG due to presynaptic autocorrelation which is confirmed by the fact that for the artificial synapse $g(\delta)$ is now monotone in Figure 5C as expected. However, this does not capture the whole behavior of the initial AdEx system CCG with a biophysical synapse in Figure 5B. This indicates that these secondary oscillations are not pure epiphenomena but instead include some causal effect that is, in fact, confounded by presynaptic autocorrelation. This was already clear by the definition of $g(\delta)$ as well, which indicates that some fraction of the causal spikes contributing to the secondary oscillations is, in fact, comprise to presynaptic input (i.e., not "second spikes" in a rapid burst).

While the synapse is excitatory, the non-monotonic behavior of $g(\delta)$ in the AdEx system also implies some negative gain at some points on the curve. This is likely due to a combination of the refractory periods and bias selection that causes some spikes not to occur that would have happened if the synapse had not existed. This highlights several reasons why unbiased causal effects cannot be obtained from correlation functions, including deconvolution of the CCG with the presynaptic auto-correlogram (ACG) outside neatly controlled cases. Furthermore, it must be stressed that there might exist many other causes, including network oscillations, that give rise to secondary oscillations in CCG, and so spiking correlation functions, in general, fail to address the fundamental problem of causal inference: confounding.

Estimation of the AdEx system ensues exactly as before. In Figure 5D-F point and interval estimates are plotted for all simulations just explored using eleven equally-spaced values for $g_0 \in [0, 0.1] \text{ mS/cm}^2$ to generate different levels of causal rate. Figure 5A displays estimates for the full AdEx system defining θ_{syn} as $g(\delta_0)$. While all the confidence intervals still cover the true parameter, there is a clear underestimation. However, Figure 5E is obtained from the artificial synapse with the same presynaptic spike trains as Figure 5D and the bias vanishes. Thus, we may deduce that the bias observed in Figure 5A is not due to presynaptic autocorrelation. This is expected, as no assumptions were made about R in the theoretical development of the monosynaptic causal inference model. One possibility is that the monosynaptic causal inference model does not best approximate this dynamical system using δ_0 . Instead, we tried using $\delta_1 = \arg \max_{\delta} g(\delta)$. While the resulting estimates are not as precise as in the LIF, it appears in this model that



Fig 5: Causality and spike history effects in an AdEx system of a neocortical pyramidal cell driving a fast-spiking interneuron. The model neurons exhibit strong spike history effects. A: $g(\delta)$ and $N_{S(R,2\delta,\delta)}(T)$ are plotted for the AdEx system in a simulation lasting 27.77 simulated hours. Note the non-monotone fluctuations in $g(\delta)$. B: The green CCG is the observation $\chi(R, T)$ from the same simulation. The gray-filled CCG is $\chi(R, T^{(\emptyset)})$ from the corresponding frozen noise simulation with the synapse removed. The black line is a CCG constructed by adding synchronous spikes to $T^{(\emptyset)}$ (termed "artificial synapse") such that the total spike count equals |T|; the synchronous spike times are added with a time constant chosen to maximize correlation with $\chi(R, T)$. Note that secondary oscillations persist in the black CCG due to presynaptic autocorrelation, but the full behavior remains unexplained by assuming independence of the causal spikes in the postsynaptic train. C: $g(\delta)$ for R and T is plotted in green. Plotted in black are $g(\delta)$ for R and the modified target train constructed by adding the artificial synapse to $T^{(\emptyset)}$. Note the non-monotone fluctuations vanish. D: Estimates are biased for $g(\delta_0)$. E: The bias vanishes with the artificial synapse. F: Bias is also slightly reduced by estimating $g(\delta)$ at $\delta_1 = \arg \max_{\delta} g(\delta)$, however perhaps at the cost of less precision and a Type 1 error for $g(\delta_1) = 0$. Different levels of causal rate are generated by varying $g_0 \in \{0, .01, ..., 0.1\}$ mS/cm² as before.

 $g(\delta_1)$ is better identified than $g(\delta_0)$ for most coupling strengths as shown in Figure 5F. As a tangential point, this also shows that estimation, in general, might be reasonable across some range of δ .

Parameter Name	Symbol	Unit	Value		
Cellular Properties					
Leak Conductance	g_1	mS/cm ²	1/15		
Refractory Period	$ au_r$	ms	5		
Membrane Capacitance	C_m	μ F/cm 2	1		
Leak Reversal Potential	E_1	mV	-65		
Voltage Reset	V_{R}	mV	E_l		
Adaptation Current Timescale	$ au_w$	ms	500		
Spike Threshold	V_T	mV	-50		
Synapse					
Peak Synaptic Conductance	g_0	mS/cm ²	0.05		
Synaptic Reversal Potential	$E_{\rm syn}$	mV	0		
Synaptic Time Constant	$ au_{ m syn}$	ms	3		
Conduction Delay	$ au_d$	ms	0		
Pyramidal Neuron					
Activation Slope	$k_{a,0}$	mV	2		
Adaptation Conductance	a_0	mS/cm ²	2.04		
Adaptation Increment	b_0	μ A/cm 2	0.02		
Reset Condition	$V_{T,0}$	mV	$V_T + 5k_{a,0}$		
Interneuron					
Activation Slope	$k_{a,1}$	mV	0.5		
Adaptation Conductance	a_1	mS/cm ²	0		
Adaptation Increment	b_1	μ A/cm 2	0		
Reset Condition	$V_{T,1}$	mV	$V_T + 5k_{a,1}$		
Background Input Currents					
Input timescales	$\tau_{I,i}$ *	ms	50		
Input Mean	μ_i *	μ A/cm ²	0		
Input SD	$\sigma_i *$	$\mu A/cm^2$	1		

 Table 3: AdEx Circuit Parameters

*: for $i \in \{0, 1, 2\}$

5 Neural perturbations for testing assumptions and fitting free parameters

The frequently invoked separation of timescales hypothesis in monosynaptic inference [3, 10, 2, 51, 19] to some degree suggests we may learn something useful by studying a toy model of instantaneously coupled Bernoulli processes in discrete time. Importantly, this setting possesses the feature that presynaptic and postsynaptic spikes can be thought of as sequences of binary treatment and outcome variables. When the synaptic effect is very fine-timescale, as is often observed *in vivo* [7], and when firing rates are sparse, this might be a reasonable approximation. Of course, the analogy breaks in obvious ways including long synaptic decay times, temporal summation of PSPs, spike history effects, etc. But the toy model can clarify issues about causality and, fortunately for neuroscience, well-developed causal inference concepts for binary treatment and outcomes variables can then be applied to pairwise spike trains in a fairly straightforward way. In this section, we make this simplification to discuss how perturbation experiments (e.g., optogenetics) could test the monosynaptic model's assumptions or fit free parameters.

5.1 Monosynaptic model calibration in an ideal neural perturbation experiment

For simulations in this setting, we will also retreat back to point process simulations that are even simpler than the one of Section 4.1. As building blocks, piecewise constant excitability functions will be used for various purposes,

$$b_i(t_j) = \sum_{k \in \mathbb{Z}^*} m_{i,k} \mathbb{1}\{(k-1)\Delta \le t_j < k\Delta\}, \text{ for } i = 0, 1$$
(49)

where the $m_{i,k}$ are repurposed from a multivariate skew of dimension n = 2 (see Eq. 34), with discrete time points $\{t_j\}_{j \in \{0,1,\dots,D\}}$ for an experiment of duration D, and where Δ is the bandwidth of the amplitudes chosen as a constant equal to the statistical free parameter of the same name defined in previous sections. This will be used to construct conditional intensity functions in an idealized monosynapse model. Working in discrete time, sets of spike times in this section will be defined as sets of integers.

Here the relationship between θ_{syn} and the *probabilities of causation* of Tian and Pearl [1] is demonstrated in simulation. This provides an alternative set of assumptions to identify causal effects that utilize observational and experimental data. As before we work in the toy case of instantaneously coupled Bernoulli processes in simulation. For $0 \le j \le D$, define the conditional intensity functions,

$$\lambda_R(t_j)|b_0(t) = \rho_0 b_0(t_j) \tag{50}$$

$$\lambda_T(t_j)|b_1(t), do(\boldsymbol{R} = \boldsymbol{r}) = \rho_1 b_1(t_j) + \epsilon \mathbb{1}\{t_j \in \boldsymbol{r}\}$$
(51)

where ρ_0 and ρ_1 are normalization factors and ϵ is a fixed instantaneous coupling constant chosen such that $\lambda_T(t_j)|\mathbf{R}$ remains a proper intensity function. We map this toy model onto the monosynaptic causal model by identifying \mathbf{R} and $\mathbf{T}^{(\mathbf{R})}$ as the sets of spike times generated from $\lambda_R |b_0(t), \lambda_T(t_j)| b_1(t), \mathbf{R}$ respectively.

We implement the model structurally by extending the analogy of Example 1, with \mathbf{R} and \mathbf{T} conditionally independent given knowledge of the (causal) conditional intensity functions (51). Expanding on that, define an idealized neural intervention of the presynaptic neuron as one that causally induces a new reference train \mathbf{r}_0 ; $do(\mathbf{R} = \mathbf{r}_0)$ is implemented by independently sampling the reference train from a constant intensity function $\lambda_{opto}(t_j) = \lambda_0$ inducing an experimental version of the postsynaptic intensity, $\lambda_T(t_j)|b_1(t), do(\mathbf{R} = \mathbf{r}_0) = \rho_1 b_1(t_j) + \epsilon \mathbb{1}\{t_j \in \mathbf{r}_0\}$ with outcomes $\mathbf{T}^{(\mathbf{r}_0)}$. Here (in discrete time) this is effectively the common notion of experimental randomization whereby every time bin is assigned to spike by mechanisms that act independently and homogeneously across time.

Returning to *probabilities of causation*, in the general case of Bernoulli random variables X and Y, respectively, Tian and Pearl [1] define these probabilities as follows,

Ì

$$PN = \mathbb{P}\left(Y^{(X=0)} = 0 | X = 1, Y = 1\right)$$
 (probability of necessity) (52)

$$PS = \mathbb{P}\left(Y^{(X=1)} = 1 | X = 0, Y = 0\right)$$
 (probability of sufficiency) (53)

$$PNS = \mathbb{P}\left(Y^{(X=1)} = 1, Y^{(X=0)} = 0\right)$$
 (probability of necessity & sufficiency). (54)

For example, probability of necessity (PN) is the probability that $\{X = 1\}$ is a necessary cause of the effect $\{Y = 1\}$. It is the probability that, given the event that $\{X = 1\}$ and $\{Y = 1\}$ both occur, Y is 0 when X is forced (via intervention) to be 0. More loosely, X would be 0, were it not that Y is 1; that is, X is the *necessary* cause of Y. PS and PNS have similar interpretations. We refer the interested reader to Ch. 9 of Pearl [17] for a fuller review.

To map these probabilities into our experiments (e.g., spikes simulated from the structural causal model in the specification above, including Eqs. (50)-(51)), let V be a random time: $V \sim Uniform\{1, 2, ..., D\}$. Then $X := \mathbb{1}\{V \in \mathbf{R}\}$, and $Y := \mathbb{1}\{V \in \mathbf{T}\}$ and apply Eqs. (52-54). Hence, in simulation, we will identify the ground truth of PN with its intervention-inferred numerical estimate $PN = \theta_{syn}/|\mathbf{R} \cap \mathbf{T}|$. This is the true proportion of causal synchrony to observed synchrony. (Note that the noise processes are not iid, so there is an additional, implicit assumption that the noise processes are mixing quickly enough to make the error in this identification negligible. We do not analyze this error, or incorporate a variability assessment.)

Pearl identifies several ways to identify PN from observational and experimental data [17]. For our purposes, an acceptable assumption is *monotonicity*, which here simply requires a synapse to be strictly excitatory ($\epsilon > 0$) or strictly inhibitory ($\epsilon < 0$). Let us explain the excitatory case. The inhibitory case follows precisely the same logic but redefines the outcome variable as silence rather than a spike. We follow Pearl [17] and for finite data assert by hypothesis an alternative estimate for $PN = \theta_{syn}/|\mathbf{R} \cap \mathbf{T}|$ as,

$$\hat{PN}_{exp} := \left(\frac{|\boldsymbol{T} \cap \boldsymbol{R}|}{D}\right)^{-1} \left(\frac{|\boldsymbol{T}|}{D} - \frac{|\boldsymbol{T}^{(\boldsymbol{r}_0)} \setminus \boldsymbol{r}_0|}{D - |\boldsymbol{r}_0|}\right), \text{ if } \epsilon \ge 0$$
(55)

where as defined earlier D is the duration of the experiment. The estimator uses spontaneous and perturbation data as just outlined. Under the monosynaptic causal inference model, the analogous estimator is denoted $\hat{PN}_{obs} = \hat{\theta}_{syn}/|\mathbf{R} \cap \mathbf{T}|$ requiring only observational data under its assumptions. Likewise, we suggest $PNS = \epsilon$ in the toy model with the alternative estimator,

$$P\hat{N}S_{exp} := \frac{|\mathbf{T}^{(\mathbf{r}_{0})} \cap \mathbf{r}_{0}|}{|\mathbf{r}_{0}|} - \frac{|\mathbf{T}^{(\mathbf{r}_{0})} \setminus \mathbf{r}_{0}|}{D - |\mathbf{r}_{0}|}, \text{ if } \epsilon \ge 0.$$
(56)

The monosynaptic causal inference model's corresponding estimate will be $P\hat{N}S_{obs} = \hat{\theta}_{syn}/|\mathbf{R}|$. Notice this gives a more principled account of what neurophysiologists often call *efficacy* [52] or *spike transmission gain* [53], which are, loosely, the excess probability of a postsynaptic spike given that a presynaptic spike occurred. We use the word *loosely* because the word *excess* has no universal interpretation (see excellent review in Stevenson [20]), and to our knowledge, none have formally interpreted *excess* in terms of counterfactuals and potential outcome random variables. Finally, for the ground truth numerical *probability of sufficiency* let $\iota := (\mathbb{Z}^* \cap [0, D)) \setminus (\mathbf{R} \cup \mathbf{T})$. Then, $PS = |\mathbf{T}^{(\iota)} \cap \iota|/\iota$ with alternative estimate,

$$\hat{PS}_{exp} := \left(\frac{D - |\mathbf{T} \cup \mathbf{R}|}{D}\right)^{-1} \left(\frac{|\mathbf{T}^{(\mathbf{r}_0)} \cap \mathbf{r}_0|}{|\mathbf{r}_0|} - \frac{|\mathbf{T}|}{D}\right), \text{ if } \epsilon \ge 0$$
(57)

and estimated from observational data only by the monosynaptic causal inference model as,

$$\hat{PS}_{obs} = \left(\frac{D - |\mathbf{T} \cup \mathbf{R}|}{D}\right)^{-1} \left(\hat{PNS}_{obs} - \left(\frac{|\mathbf{T} \cap \mathbf{R}|}{D}\hat{PN}_{obs}\right)\right).$$
(58)

As mentioned before, these quantities can be obtained for inhibition in the exact same way where the queried postsynaptic outcome variable is silence. For visualization purposes, in the inhibitory case, we define the probabilities of causation through multiplication by -1 so that an estimate of inhibition can be plotted simultaneously with excitation and compared with θ_{syn} on its negative support. For example, in the inhibitory case, we will have,

$$\hat{PN}_{syn} = -1 \left(\frac{|\boldsymbol{R} \setminus \boldsymbol{T}|}{D} \right)^{-1} \left(\frac{D - |\boldsymbol{T}|}{D} - \frac{D - |\boldsymbol{T}^{(\boldsymbol{r}_0)} \cup \boldsymbol{r}_0|}{D - |\boldsymbol{r}_0|} \right), \epsilon < 0.$$
(59)

The significant observation is that the probabilities of causation are obtained from experimental and observational data without appeal to some of the assumptions that make θ_{syn} identifiable from observational data alone. Namely, $\mathcal{A}.1$ (conditional uniformity) and $\mathcal{A}.2$ (separation of timescales) are not required to identify the probabilities of causation. For this reason, if these idealized concepts could be extended to fit more realistic aspects spike trains recorded *in vivo*, we have here provided an experimental test of $\mathcal{A}.1$ and $\mathcal{A}.2$ that could be conducted in the laboratory. Essential in this endeavor would be confidence limits, say for PN, with finite data, which is research currently being pursued [54]. However, the final section of this study will argue that such an experiment is not easily achieved by current experimental technologies (e.g., optogenetic stimulation). These alternative estimators also might provide a route to estimate the free parameter δ , τ , and Δ .

We conclude this section by simulating spike trains from the toy model in Eqs. (50)-(51). Simulation details are exactly analogous to those in Figure 2. Figure 6 shows the results of forty-two simulations (twenty-one for excitatory and inhibitory estimates) as just described and each plot shows the corresponding point estimates for PNS (Figure 6A), PS (Figure 6B), and PN (Figure 6C). In each case, a tight correspondence is shown between the θ_{syn} -derived estimates, which come from observational data, and the alternative estimates, which use a combination of experimental and observational data.



Fig 6: An idealized experimental test of the monosynaptic model's assumptions. In sparse firing conditions where presynaptic and postsynaptic spikes can be approximated as binary treatment and outcome variables, the monosynaptic model can be related to the *probabilities of causation* of Tian and Pearl [1] in a toy spiking model. $P\hat{N}S_{exp}$, $P\hat{S}_{exp}$, and $P\hat{N}_{exp}$ are derived directly from Tian and Pearl [1] and provide alternative estimates for monosynaptic causal effects by combining neural intervention data with spontaneous data and thus require fewer assumptions on the background processes, providing an experimental test of the monosynaptic causal inference model's assumptions. $P\hat{N}S_{obs}$, $P\hat{S}_{obs}$, and $P\hat{N}_{obs}$ are different normalizations of $\hat{\theta}_{syn}$ obtained from observational data as described in the text. A: Probability of necessity and sufficiency (*PNS*) corresponds to ϵ in the toy model. B: Probability of sufficiency (*PS*). C: Probability of necessity (*PN*).

5.2 Even strong perturbations might quite strongly fail as randomized experiments.

In the previous section, we explored conceptually the notion of an ideal neural intervention in a toy spiking model. The purpose of this was to highlight connections between θ_{syn} and more well-established causal inference concepts and to speculate about avenues for future research that might make the interventions suitable for more realistic dynamical models. Another concern persists, which is the degree to which current experimental technologies actually achieve the theoretical notion of an *intervention* in causal inference. Recently, Lepperød et al. [8] fruitfully analyzed the confounding that arises from optogenetic stimulation activating many neurons that may be unobserved. However, while it is well-understood that stimulation often increases the empirical rate of the presynaptic neuron on a coarse timescale [2], it is not clear in a dynamical system how much deconfounding occurs at the level of voltage and hence what the proper interpretation of juxtacellular or optogenetic stimulation is. In this section, we caution against simple interpretations (and thus show the difficulty of obtaining the ideal intervention we proposed in the previous section) by injecting stochastic input currents into correlated but unconnected LIF neurons, as well as stimulating them with a biophysically detailed channelrhodopsin model [55].

A simple example to consider is two unconnected LIF neurons with common input that produces structure in the CCG. An ideal intervention, where each point in time is randomly assigned a presynaptic spike or not (see Section 5.1) should destroy all structure in the cross-correlogram during stimulation, yielding a flat histogram. Consider two LIF neurons,

$$C_m \frac{dV_i}{dt} = -g_l(V_i - E_l) + I_c(t) + I_i(t) + I_p(t)\{i = 1\}, \text{ for } i = 0, 1$$
(60)

where as before if at time t_0 , $V(t_0) = V_T$, the voltage is reset to E_l . With the same form as Eq. 42 but with a slight change in notation, $I_c(t)$ is common OU noise to both neurons, $I_i(t)$ for i = 0, 1 is independent OU noise for each neuron. $I_p(t)$ is either an injected current identical to the stimulus to be described momentarily or the same stimulus filtered by the channelrhodopsin (ChR2) model of Williams et al. [55].

Let S(t) be the stochastic stimulus, then

$$I_p(t) = S(t)$$
 (for current stimulation) (61)

$$I_p(t) = g_{ChB2}G(V)(O_1 + \gamma O_2)(V_0 - E_{ChB2})$$
 (for optogenetic stimulation) (62)

where in the notation of Williams et al. [55] g_{ChR2} is the max conductance of the photocurrent, E_{ChR2} is the reversal potential for channelrhodopsin, G(V) is a voltage-dependent rectification function, O_1, O_2 are open state probabilities,

and γ is a normalization factor. Eq. (62) is identical to Eq. 1 in Williams et al. [55], and we replace what in their notation is termed $S_0(\theta)$ in their Eq. 11 with our stimulus S(t). We refer the reader to the rest of that study since the channelrhodopsin model is rather complicated and has various state variables and parameters. We used identical parameters from the original study for channelrhodopsin.

In simulation, we take S(t) to be a special discrete construction of a Gauss-Markov process with Hurst or Hölder parameter $H \in [0, 1]$. The motivation for this is to generate repeatable spike patterns in the presynaptic neuron regardless of the level of other sources of noise [56], constituting the notion of experimental intervention. The parameter H plays the same role as in fractional Brownian motion, intuitively describing how rough (small H) versus smooth (high H) the trajectory is, however the process used here is colored (i.e., its power spectrum is not flat). The process was developed as an injected current in previous work to suggest that more reliable spiking patterns can be induced into a LIF neuron to the degree that H is small regardless of the neuron's level of independent noise [57]. The construction of the process is described in Appendix 6.3. In this setting, it is simply being employed as technology to produce reliable spiking responses to stimulation [58, 59], although the tenability of this very statement in this setting is what is being tested in the simulation. Consider that if a spiking pattern were perfectly reliable to a repeated stimulus, then an experimentalist would know that they are deconfounding in the sense of the $do(\cdot)$ operator of causal inference.

Figure 7 simulates the system in Eq. (60) for different parameters of the stochastic input current or light stimulus; the timescale τ_H and Hölder parameter H. In each simulation, the timescales of the intrinsic processes $I_0(t)$, $I_1(t)$, and $I_c(t)$ were set to 10 ms, and their amplitudes to unit variance. The amplitudes of the stimulations were then adjusted so that the reference neuron's empirical rate during stimulation was approximately 470% greater than the spontaneous rate as in the experiment of English et al. [2]. Equalizing firing rate across experimental conditions in this way, surprisingly quite strong common input correlations persist for current injection and optogenetic stimulation. Furthermore, varying the input parameters H and τ_H leads to hardly detectable differences in the deconfounding as measured through the CCG. If current or optogenetic stimulation fulfilled the notion of $do(\cdot)$ as applied to spike trains in Section 5.1, the CCG during stimulation should be flat.



Fig 7: Strong juxtacellular or optogenetic stimulation might fail to be randomized experiments. Two LIF neurons are driven by common inputs, and their CCG is plotted (black plots). With the same frozen noise input, the model system is subjected to either (1) current injection in the pattern of a Gaussian process (gray plots) or (2) photocurrent stimulation in the same pattern on a biophysically detailed opsin model affixed to one of the LIF neurons (green plots). A special Gaussian process is utilized from a theory that predicts reliable spike patterns to be produced to the degree that a parameter $H \in [0, 1]$ is small. The timescale of the stimulus, τ_H , is also split into two conditions, 5 ms and 50 ms simulations. The variance of the stimulations was adjusted for electric current injection or photostimulation such that the stimulated firing rate was approximately 470% greater than the spontaneous rate [2]. Even for this strong perturbation, confounding common synaptic input correlations persist and do not significantly differ given the character of the input when the variance of the stimulation is adjusted to produce equal firing rates across conditions.

6 Appendix

Abbreviation	Definition
ACG	Auto-correlogram
AdEx	Adaptive exponential integrate-and-fire neuron
CCG	Cross-correlogram
cdf	Cumulative distribution function
CGF	Cumulant generating function
ChR2	Channelrhodopsin-2
DC	Direct convolution
DFT/IDFT	Discrete-Fourier transform and its inverse
FFT/IFFT	Fast-Fourier transform and its inverse
iid	Independent and identically distributed
INT	Interneuron
LIF	Leaky integrate-and-fire neuron
OU	Ornstein–Uhlenbeck process
pmf	Probability mass function
PSP	Postsynaptic potential
PYR	Pyramidal neuron
SD	Standard deviation

Table 4: List of abbreviations

6.1 Examples of confounding and non-identifiability in the CCG

To prepare for examples that demonstrate this issue, let R and T be a finite set of spike times (a point process) for an experiment of fixed duration. We will, in general, consider a reference spike train R hypothesized to be presynaptic, and a target spike train T, hypothesized to be postsynaptic. A goal is to quantify the evidence for that hypothesis and for a number of its characteristics. We are thus interested in the potential outcome random variable, $T^{(R=r)}$, abbreviated $T^{(r)}$, which is the target train, in a causal model, induced by do(R = r). That is, we are interested in the causal influence of R on T. For any spike trains X and Y define the unnormalized sample cross-correlation function (sample CCF) as,

$$\hat{\chi}(\boldsymbol{X},\boldsymbol{Y},\tau) \coloneqq \int_{-\infty}^{\infty} \sum_{x \in \boldsymbol{X}} \delta_d(t-x) \sum_{y \in \boldsymbol{Y}} \delta_d(t-y+\tau) dt$$
(63)

where δ_d is the Dirac delta function. We will also write $\chi(\mathbf{X}, \mathbf{Y}, \tau) = \mathbb{E}[\hat{\chi}(\mathbf{X}, \mathbf{Y}, \tau)]$ and will occasionally assume the spike trains are discrete, reinterpreting the notation accordingly when specified. The term unnormalized cross-correlogram (CCG) likewise refers to a binned version of $\hat{\chi}(\mathbf{X}, \mathbf{Y}, \tau)$.

The following examples motivate the approach of this article. Figure 1 illustrates their simulation and the causal decompositions described in the examples. Example 1 presents an example of a causal model in terms of point process models, and subsequent simulations will utilize this definition of causality.

We start with the simplest model one might imagine.

Example 1 (Instantaneously-coupled Bernoulli processes with fixed coupling constant ϵ). Define a probability space which contains $\boldsymbol{\omega} = (\omega_1, ..., \omega_{2N})$, a vector of 2N independent uniform [0,1] random variables. Then consider the following potential outcomes model: $\mathbf{R}(\boldsymbol{\omega}) = \{t_j : \omega_j \leq \lambda_R\}$ and $\mathbf{T}^{(\mathbf{R}=\mathbf{r})}(\boldsymbol{\omega}) = \{t_j : \omega_{j+N} \leq \lambda_T + \epsilon \mathbb{1}\{t_j \in \mathbf{r}\}\}$. By independence, there is no confounding (of \mathbf{R} and \mathbf{T}). The average causal effect of the coupling at time t, $E[\mathbb{1}\{t \in \mathbf{T}\} - \mathbb{1}\{t \in \mathbf{T}^{(\mathbf{R}=\emptyset)}\}]$, is ϵ . Consider, for example, the intervention $do(\mathbf{R} = \emptyset)$: $\mathbf{T}^{(\mathbf{R}=\emptyset)}(\boldsymbol{\omega}) = \{t_j : \omega_{j+N} \leq \lambda_T + \epsilon \mathbb{1}\{t_j \in \emptyset\}\} = \{t_j : \omega_{j+N} \leq \lambda_T\}$. $T^{(\mathbf{R}=\emptyset)}(\boldsymbol{\omega})$ is defined as a function on the same probability space as the functions $\mathbf{R}(\boldsymbol{\omega})$ and $\mathbf{T}(\boldsymbol{\omega}) = \mathbf{T}^{(\mathbf{R})}(\boldsymbol{\omega})$. $(\omega_1, \omega_2, ..., \omega_{2N})$ are so-called 'background' variables. We think of a particular realization of $\boldsymbol{\omega}$ as encoding the state(s) of the 'external' world. Interventions modify the relations between \mathbf{R} and \mathbf{T} to define potential outcomes for \mathbf{T} , given that the state(s) of the world (i.e., the background variables $\boldsymbol{\omega}$) are fixed (i.e., 'frozen') over potential outcomes.

Example 2 now examines the behavior of the CCF for Example 1 demonstrating that ϵ is not identifiable.

Example 2 (Identical CCGs with different coupling strength). One can verify from the independence relations that a normalized CCF for the model in Example 1 is $\chi(\mathbf{R}, \mathbf{T}^{(\mathbf{R})}, \tau)/N = \mathbb{E}[\sum_{t=1}^{N} \mathbb{I}\{t - \tau \in \mathbf{R}\}\mathbb{I}\{t \in \mathbf{T}^{(\mathbf{R})}\}]/N = \lambda_R(\lambda_T + \epsilon)\mathbb{I}\{\tau = 0\} + (\lambda_R\lambda_T + \lambda_R^2\epsilon)\mathbb{I}\{\tau \neq 0\}$ dismissing edge effects. Since in this model the CCF is flat everywhere but the coupling lag at $\tau = 0$, the CCF peak, $\rho = \chi(\mathbf{R}, \mathbf{T}^{(\mathbf{R})}, \tau = 0) - \chi(\mathbf{R}, \mathbf{T}^{(\mathbf{R})}, \tau = z)$ for any $z \neq 0$, can be related to the average causal effect as $\epsilon = \rho/(\lambda_R - \lambda_R^2)$ where we have normalized the peak by λ_R which is standard in functional connectivity studies. ϵ and ρ are not equal. Consider two situations. In Situation A the model has parameters $\lambda_{R,A}, \epsilon_A, \lambda_{T,A}$. In Situation B, an identical λ_R -normalized CCF is obtained by setting $\lambda_{R,B} = 1 - \epsilon_A, \epsilon_B = 1 - \lambda_{R,A}, \lambda_{T,B} = \lambda_{T,A} + \lambda_{R,A}\epsilon_A - \lambda_{R,B}\epsilon_B$. For example, we can have the average causal effects be $\epsilon_A = 0.2$, in Situation A, and $\epsilon_B = 0.8$, in Situation B, and yet their CCFs are identical.

In the previous example, ϵ is identifible if supplemented by one additional unknown, λ_R . Here that is trivial if the process is stationary, however neural data is known to be highly nonstationary [31, 30] leading to extreme difficulty in estimating analogous time-varying quantities [12]. When λ_R is unknown, one can verify through level set analysis that the model parameters in the example can vary widely for fixed ϵ , suggesting large relative bias if λ_R is slightly misestimated, even when λ_R is small.

Related to the observation that λ_R confounds estimation in Example 2, it has long been understood that presynaptic autocorrelation might in some way influence the CCG between two neurons [60] leading to suggestions that deconvolution of the CCG with the presynaptic ACG might help in deconfounding, particularly under stationarity assumptions [19]. In the next example, we examine the nature of this confounding in a nonstationary setting. Two neurons are given confounding oscillatory backgrounds. In addition, the presynaptic cell emits a burst of three spikes approximately every second. As a thought experiment, imagine these bursting events alternate such that the spike times in the bursts are either generated by a Gaussian with small variance or large variance. Causal conclusions from the CCG vary widely in their dependence on whether the causal interactions tend to occur among the bursts with small or large variances, highlighting the non-identifiability of causal inference from correlation functions altogether. Later, we use this intuition to construct confidence intervals by supposing causal events occur at these limiting cases of presynaptic firing, thus bounding the estimate over this uncertainty (Section 3.4.1).

Example 3 (Different CCGs with identical coupling strength). Consider the following generative model for an experiment of duration D. Let background intensity functions be $\lambda_R(t) = \lambda_T(t) = \alpha \cos(\omega t) + \alpha$ for $t \in [0, D), \alpha \in [0, 1/2)$ both generating sets of real-valued points \mathbf{R}_0 and \mathbf{T}_0 , respectively. Define a sequence of latent events $0 \leq \ell_1, \ell_2, ..., \ell_K \leq D$ such that $\ell_k - \ell_{k-1} \sim Uniform(0.8, 1.2)$ seconds. Let these latent events be the center of a burst of three spikes, $X_{k,i} \sim \mathcal{N}(\ell_k, \sigma_A \mathbbm{1}\{k \text{ is oven}\} + \sigma_B \mathbbm{1}\{k \text{ is odd}\}$ for $i \in \{1, 2, 3\}$. Collect these events in a set $\mathbf{R}_1 = \bigcup_{k \in \mathbb{N}} \bigcup_{i=1}^3 X_{k,i}$ and define the presynaptic spike train as $\mathbf{R} = \mathbf{R}_0 \cup \mathbf{R}_1$. Also, let $Y_{k,i} \sim \mathcal{N}(X_{k,i} + d, \sigma_s \mathbbm{1}\{k \text{ is even}\} + \sqrt{2\sigma_A^2} + \sigma_s^2 \mathbbm{1}\{k \text{ is odd}\}$ for $i \in \{1, 2, 3\}$ and $\mathbf{X} = \bigcup_{k \text{ even}} \bigcup_{i=1}^3 Y_{k,i}$ and $\mathbf{Y} = \bigcup_{\{k \text{ odd}\}} \bigcup_{i=1}^3 Y_{k,i}$. Now consider two situations. In Situation A, $\mathbf{T}^{(\emptyset)} = \mathbf{T}_0$ and $\mathbf{T}^{(\mathbf{R})} = \mathbf{T}_0 \cup \mathbf{X}$. Since cross-correlation is a linear operator and the constituent processes are in superposition, the unnormalized CCF can be expressed in approximate closed-form as $\chi_A(\mathbf{R}, \mathbf{T}, \tau) = \chi(\mathbf{R}_0, \mathbf{T}_0, \tau) + \chi(\mathbf{R}_1, \mathbf{T}_0, \tau) + \chi(\mathbf{R}_0, \mathbf{X}, \tau) + \chi(\mathbf{R}_1, \mathbf{X}, \tau) \approx D\alpha^2/2\cos(\omega t) + \alpha^2 + D^{-1}(\mathbb{E}[|\mathbf{R}_1|] \mathbb{E}[|\mathbf{R}_0|]) + \frac{1}{2} \mathbb{E}[|\mathbf{R}_1|] \mathcal{N}(\tau | d, \sqrt{2\sigma_A^2} + \sigma_s^2)$ where we have dismissed edge effects. In Situation B, $\mathbf{T}^{(\emptyset)} = \mathbf{T}_0$ and $\mathbf{T}^{(\mathbf{R})} = \mathbf{T}_0 \cup \mathbf{Y}$ and by the same logic $\chi_B(\mathbf{R}, \mathbf{T}, \tau) \approx D\alpha^2/2\cos(\omega t) + \alpha^2 + D^{-1}(\mathbb{E}[|\mathbf{R}_1|] \mathbb{E}[|\mathbf{T}_0|] + \mathbb{E}[|\mathbf{Y}|] \mathbb{E}[|\mathbf{R}_0|]) + \frac{1}{2} \mathbb{E}[|\mathbf{R}_1|] \mathcal{N}(\tau | d, \sqrt{2\sigma_A^2} + \sigma_s^2) + \frac{3-1}{2} \mathbb{E}[|\mathbf{R}_1|] \mathcal{N}(\tau | d, \sqrt{2\sigma_A^2} + \sigma_s^2)$. In both cases, we write $\frac{3-1}{2} \mathbb{E}[|\mathbf{R}_1|] \mathcal{N}(\tau | d, \sqrt{2\sigma_A^2} + \sigma_s^2) + \frac{3-1}{2} \mathbb{E}[|\mathbf{R}_1|] \mathcal{N}(\tau | d, \sqrt{2\sigma_A^2} + \sigma_s^2)$. In both cases, we write $\frac{3-1}{2} \mathbb{E}[|\mathbf{R}_1|]$ to highlight that there are three causal

These examples are intentionally dramatic to be instructive and often exceed plausible neurophysiological behavior. Namely, the firing rates are often quite high, and we used Gaussian functions for analytic tractability, although they allow for some causality to occur in reverse time. However, if the examples are understood in mathematical detail, it's straightforward to see that the degeneracy is quite general, especially if we are concerned with the relative error of estimates. Moreover, the regime of nonstationary, high presynaptic bursting coinciding with information transfer is perhaps the most biologically relevant [61, 47, 62]. A common way presynaptic bursting manifests in the CCG is as secondary oscillations visually distinct from the primary monosynaptic peak because of the refractory periods [19]. We explore this in Section 4.3 and Figure 5. Note, however, that in Example 3 the influence of bursting on the CCG depends on how the temporal resolution of bursting interacts with the temporal resolution of the causal interactions

they associate with; a concept that should generalize beyond the idealizations made here. So it is not guaranteed that refractory periods will dissolve the issue. In fact, statistical dependence between the temporal resolution of bursting and the causal interactions (e.g., from short-term plasticity) might likely make the interpretation even more subtle, leading to multiple sources of estimation error at different lags of the CCG (that is, a combination of Example 2 and Example 3). Complex dependencies between background input currents and other forms of nonstationarity suggest that the toy examples here may paint a forgiving picture of the confounding present in real neural data [6].

6.2 Monosynaptic causal inference model proofs

Theorem 1. Under A.1 (conditional uniformity), A.2 (timescale separation), A.3 (positivity) and A.4 (consistency) an unbiased point estimate of θ_{syn} in the excitatory and inhibitory models Eq. (4) is given by one expression,

$$\hat{\theta}_{syn} = \sum_{k \in \mathbb{Z}^*} \frac{N_{\gamma(k\Delta) \cap S(\mathbf{R})}(\mathbf{T}) - q(\mathbf{R}, k\Delta) N_{\gamma(k\Delta)}(\mathbf{T})}{1 - q(\mathbf{R}, k\Delta)}.$$
(9)

Proof.

Case 1: Excitation, $\theta_{syn} \ge 0$.

For an arbitrary $k \in \mathbb{Z}^*$ we have from the definition of the model Eq. (4) we have, for all r,

$$N_{\gamma(k\Delta)\cap S(\boldsymbol{r})}(\boldsymbol{T}^{(\boldsymbol{r})}) = N_{\gamma(k\Delta)}(\boldsymbol{I}^{(\boldsymbol{r})} \setminus \bigcup_{r \in \boldsymbol{r}} \{S(r) : N_{S(r)}(\boldsymbol{B}) > 0\}) + N_{\gamma(k\Delta)\cap S(\boldsymbol{r})}(\boldsymbol{B})$$
(64)

where $S(\mathbf{r})$ has been excluded from the subscript of the increment in the first term of the RHS by $\mathcal{A}.2$. Taking conditional expectations with respect to $(\Gamma(\mathbf{T}), \mathbf{R})$, and applying $\mathcal{A}.4$ (consistency) and linearity, we have

$$\mathbb{E}\left[N_{\gamma(k\Delta)\cap S(\mathbf{R})}(\mathbf{T})\middle|\mathbf{\Gamma}(\mathbf{T}),\mathbf{R}\right]$$
(65)

$$= \mathbb{E}\left[N_{\gamma(k\Delta)}\left(\boldsymbol{I}^{(\boldsymbol{r})} \setminus \bigcup_{\boldsymbol{r} \in \boldsymbol{R}} \{S(\boldsymbol{r}) : N_{S(\boldsymbol{r})}(\boldsymbol{B}) > 0\}\right) \middle| \boldsymbol{\Gamma}(\boldsymbol{T}), \boldsymbol{R}\right]$$

$$+ \mathbb{E}\left[N_{\gamma(k\Delta) \cap S(\boldsymbol{R})}(\boldsymbol{B}) \middle| \boldsymbol{\Gamma}(\boldsymbol{T}), \boldsymbol{R}\right].$$
(66)

By A.1 (conditional uniformity),

$$\mathbb{E}[N_{\gamma(k\Delta)\cap S(\mathbf{R})}(\mathbf{B})|\mathbf{\Gamma}(\mathbf{T}),\mathbf{R}] = q(\mathbf{R},k\Delta) \mathbb{E}[N_{\gamma(k\Delta)}(\mathbf{B})|\mathbf{\Gamma}(\mathbf{T}),\mathbf{R}].$$
(67)

Furthermore, Eq. (4) under $\mathcal{A}.2$ gives $N_{\gamma(k\Delta)}(\mathbf{I}^{(r)} \setminus \bigcup_{r \in \mathbf{r}} \{S(r) : N_{S(r)}(\mathbf{B}) > 0\}) = N_{\gamma(k\Delta)}(\mathbf{T}^{(r)}) - N_{\gamma(k\Delta)}(\mathbf{T}^{(\emptyset)})$, for all \mathbf{r} . Eq. (65) then becomes,

$$\mathbb{E}\left[N_{\gamma(k\Delta)\cap S(\boldsymbol{R})}(\boldsymbol{T})\middle|\boldsymbol{\Gamma}(\boldsymbol{T}),\boldsymbol{R}\right]$$
(68)

$$= \mathbb{E}\left[N_{\gamma(k\Delta)}(\boldsymbol{T}) - N_{\gamma(k\Delta)}(\boldsymbol{T}^{(\emptyset)}) \middle| \boldsymbol{\Gamma}(\boldsymbol{T}), \boldsymbol{R}\right] + q(\boldsymbol{R}, k\Delta) \mathbb{E}\left[N_{\gamma(k\Delta)}(\boldsymbol{B}) \middle| \boldsymbol{\Gamma}(\boldsymbol{T}), \boldsymbol{R}\right]$$
(69)

$$= \mathbb{E}\left[N_{\gamma(k\Delta)}(\boldsymbol{T}) - N_{\gamma(k\Delta)}(\boldsymbol{T}^{(\emptyset)}) \middle| \boldsymbol{\Gamma}(\boldsymbol{T}), \boldsymbol{R} \right]$$
(70)

$$+q(\boldsymbol{R},k\Delta) \mathbb{E}\left[N_{\gamma(k\Delta)}(\boldsymbol{T}) - \left(N_{\gamma(k\Delta)}(\boldsymbol{T}) - N_{\gamma(k\Delta)}(\boldsymbol{T}^{(\emptyset)})\right) \middle| \boldsymbol{\Gamma}(\boldsymbol{T}),\boldsymbol{R}\right]$$
(71)

where the substitution of $N_{\gamma(k\Delta)}(B)$ inside the expectation of the last term again results from Eq. (4) under A.2. Rearranging, we obtain

$$\mathbb{E}\left[N_{\gamma(k\Delta)}(\boldsymbol{T}) - N_{\gamma(k\Delta)}(\boldsymbol{T}^{(\emptyset)}) \middle| \boldsymbol{\Gamma}(\boldsymbol{T}), \boldsymbol{R}\right]$$
(72)

$$= \mathbb{E}\left[\frac{N_{\gamma(k\Delta)\cap S(\boldsymbol{r})}(\boldsymbol{T}) - q(\boldsymbol{R}, k\Delta)N_{\gamma(k\Delta)}(\boldsymbol{T})}{1 - q(\boldsymbol{R}, k\Delta)}\Big|\boldsymbol{\Gamma}(\boldsymbol{T}), \boldsymbol{R}\right],$$
(73)

using the general fact that $\mathbb{E}[f(Y) \mathbb{E}[X|Y]] = \mathbb{E}[\mathbb{E}[f(Y)X|Y]] = \mathbb{E}[Xf(Y)]$ for (measurable) random variables X, Y, and functions f. Summing over k and taking expectations on both sides of the above gives

$$\mathbb{E}\left[\hat{\theta}_{syn}\right] = \mathbb{E}\left[\sum_{k\in\mathbb{Z}^*} \frac{N_{\gamma(k\Delta)\cap S(\boldsymbol{r})}(\boldsymbol{T}) - q(\boldsymbol{R}, k\Delta)N_{\gamma(k\Delta)}(\boldsymbol{T})}{1 - q(\boldsymbol{R}, k\Delta)}\right]$$
(74)

$$= \mathbb{E}\left[\sum_{k \in \mathbb{Z}^*} N_{\gamma(k\Delta)}(\boldsymbol{T}) - N_{\gamma(k\Delta)}(\boldsymbol{T}^{(\emptyset)})\right]$$
(75)

$$= \mathbb{E}\left[N_{S(\boldsymbol{R})}(\boldsymbol{T}) - N_{S(\boldsymbol{R})}(\boldsymbol{T}^{(\boldsymbol{\emptyset})})\right] = \theta_{syn}.$$
(76)

Case 2: Inhibition, $\theta_{syn} \leq 0$.

By Eq. (4) and A.2 the following invariant holds for all realizations of the inhibitory model and every r,

$$N_{\gamma(k\Delta)\backslash S(\boldsymbol{r})}(\boldsymbol{T}^{(\boldsymbol{r})}) = N_{\gamma(k\Delta)\backslash S(\boldsymbol{r})}(\boldsymbol{T}^{(\emptyset)}).$$
(77)

In the inhibitory model, Eq. (4) under $\mathcal{A}.2$ gives $N_{\gamma(k\Delta)\setminus S(r)}(T^{(\emptyset)}) = N_{\gamma(k\Delta)\setminus S(r)}(B), \forall r$. Using this substitution and applying similar steps as in Case 1, including $\mathcal{A}.4$ and $\mathcal{A}.1$, we have

$$\mathbb{E}\left[N_{\gamma(k\Delta)\setminus S(\boldsymbol{R})}(\boldsymbol{T})\middle|\boldsymbol{R},\boldsymbol{\Gamma}(\boldsymbol{T})\right] = \left(1 - q(\boldsymbol{R},k\Delta)\right)\mathbb{E}\left[N_{\gamma(k\Delta)}(\boldsymbol{B})\middle|\boldsymbol{R},\boldsymbol{\Gamma}(\boldsymbol{T})\right].$$
(78)

In the inhibitory model, we again have from Eq. (4) under $\mathcal{A}.2 N_{\gamma(k\Delta)}(\mathbf{T}^{(r)}) = N_{\gamma(k\Delta)}(\mathbf{B}) - N_{\gamma(k\Delta)}(\mathbf{B} \cap \bigcup_{r \in \mathbf{r}} \{S(r) : N_{S(r)}(\mathbf{I}^{(r)}) > 0\})$. Using this relation to substitute $N_{\gamma(k\Delta)}(\mathbf{B})$ in Eq. (78) and again using $\mathcal{A}.4$ for terms inside the expectation,

$$\mathbb{E}\left[N_{\gamma(k\Delta)\setminus S(\boldsymbol{R})}(\boldsymbol{T})\middle|\boldsymbol{R},\boldsymbol{\Gamma}(\boldsymbol{T})\right]$$

$$= \left(1 - q(\boldsymbol{R},k\Delta)\right) \mathbb{E}\left[N_{\gamma(k\Delta)}(\boldsymbol{T}) + N_{\gamma(k\Delta)}(\boldsymbol{R}) - \left|\int S(r) \cdot N_{\sigma(\lambda)}(\boldsymbol{I}) > 0\right|\right] \boldsymbol{R} \boldsymbol{\Gamma}(\boldsymbol{T})\right]$$
(79)

$$= \left(1 - q(\boldsymbol{R}, k\Delta)\right) \mathbb{E}\left[N_{\gamma(k\Delta)}(\boldsymbol{T}) + N_{\gamma(k\Delta)}(\boldsymbol{B} \cap \bigcup_{r \in \boldsymbol{R}} \{S(r) : N_{S(r)}(\boldsymbol{I}) > 0\}) \middle| \boldsymbol{R}, \boldsymbol{\Gamma}(\boldsymbol{T})\right]$$
(80)

$$= \left(1 - q(\boldsymbol{R}, k\Delta)\right) \mathbb{E}\left[N_{\gamma(k\Delta)}(\boldsymbol{T}) - \left(N_{\gamma(k\Delta)}(\boldsymbol{T}^{(\boldsymbol{R})}) - N_{\gamma(k\Delta)}(\boldsymbol{T}^{(\boldsymbol{\theta})})\right) \middle| \boldsymbol{R}, \boldsymbol{\Gamma}(\boldsymbol{T})\right]$$
(81)

where the change of sign inside the expectation of the last line results from the fact that under $\mathcal{A}.2$ $N_{\gamma(k\Delta)}(\mathbf{T}^{(\mathbf{R})}) - N_{\gamma(k\Delta)}(\mathbf{T}^{(\emptyset)}) \leq 0$ in the inhibitory model. Rearranging and following analogous steps as used in Case 1,

$$\theta_{syn} = -\sum_{k} \mathbb{E}\left[\frac{N_{\gamma(k\Delta)\backslash S(\boldsymbol{R})}(\boldsymbol{T}) - (1 - q(\boldsymbol{R}, k\Delta))N_{\gamma(k\Delta)}(\boldsymbol{T})}{1 - q(\boldsymbol{R}, k\Delta)}\right].$$
(82)

Setting this expression equal to Eq. (74) all terms cancel.

Lemma 1. Let

$$\mathcal{D}(\boldsymbol{R},\boldsymbol{T}) = \left\{ \boldsymbol{j} : |\boldsymbol{j}| = |\boldsymbol{J}|, (\boldsymbol{K} \setminus \boldsymbol{L}) \subseteq \boldsymbol{j} \right\}.$$
(17)

Abbreviate $\mathcal{D}(\mathbf{R}, \mathbf{T})$ as \mathcal{D} . Under $\mathcal{A}.1$ (conditional uniformity)

$$J_{\theta_{syn}}^{-} \in \operatorname*{arg\,min}_{j \in \mathcal{D}} \left\{ c^{-}(\boldsymbol{q}(\boldsymbol{r}, \boldsymbol{t}), \boldsymbol{j}) \right\}$$
(18)

and

$$\boldsymbol{J}_{\theta_{syn}}^{+} \in \operatorname*{arg\,max}_{\boldsymbol{j}\in\mathcal{D}} \left\{ c^{+}(\boldsymbol{q}(\boldsymbol{r},\boldsymbol{t}),\boldsymbol{j}) \right\}.$$
(19)

Proof. We will establish (18). (19) can be established in the same way.

$$\mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{B}) \le c \left| \boldsymbol{q}(\boldsymbol{R}, \boldsymbol{T}), \boldsymbol{J} \right.\right)$$
(83)

$$=\sum_{n=0}^{\infty}\sum_{Q\in \boldsymbol{J}^{[n]}}\prod_{i\in Q}q(\boldsymbol{R},T_i)\prod_{k\in(\boldsymbol{J}\backslash Q)}(1-q(\boldsymbol{R},T_k))$$
(84)

$$=\sum_{n=0}^{c-1}\sum_{Q\in\{\boldsymbol{J}\setminus\boldsymbol{x}\}^{[n]}}\prod_{i\in\{Q\cup\boldsymbol{x}\}}q(\boldsymbol{R},T_i)\prod_{k\in(\boldsymbol{J}\setminus(Q\cup\boldsymbol{x}))}(1-q(\boldsymbol{R},T_k))$$
$$+\sum_{n=0}^{c}\sum_{Q\in\{\boldsymbol{J}\setminus\boldsymbol{x}\}^{[n]}}\prod_{i\in Q}q(\boldsymbol{R},T_i)\prod_{k\in\{\boldsymbol{J}\setminus\boldsymbol{Q}\}}(1-q(\boldsymbol{R},T_k)),$$
(85)

using $\{Q \in J^{[n]} : x \in Q\} = \{Q \cup x : Q \in \{J \setminus x\}^{[n-1]}\}$. For an arbitrary $x \in J$, this implies

$$\frac{\partial}{\partial q(\boldsymbol{R}, T_x)} \mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{B}) \le c \left| \boldsymbol{q}(\boldsymbol{R}, \boldsymbol{T}), \boldsymbol{J} \right.\right) =$$

$$= \sum_{n=0}^{c-1} \sum_{Q \in \{\boldsymbol{J} \setminus x\}^{[n]}} \prod_{i \in Q} q(\boldsymbol{R}, T_i) \prod_{k \in \boldsymbol{J} \setminus (Q \cup x)} (1 - q(\boldsymbol{R}, T_k))$$

$$- \sum_{n=0}^{c} \sum_{Q \in \{\boldsymbol{j} \setminus x\}^{[n]}} \prod_{i \in Q} q(\boldsymbol{R}, T_i) \prod_{k \in \boldsymbol{J} \setminus (Q \cup x)} (1 - q(\boldsymbol{R}, T_k))$$

$$(87)$$

$$\sum_{\boldsymbol{Q} \in \{\boldsymbol{j} \setminus x\}^{[n]}} \prod_{i \in Q} q(\boldsymbol{R}, T_i) \prod_{k \in \boldsymbol{J} \setminus (Q \cup x)} (1 - q(\boldsymbol{R}, T_k))$$

$$(87)$$

$$= -\sum_{Q \in \{\boldsymbol{J} \setminus x\}^{[c]}} \prod_{i \in Q} q(\boldsymbol{R}, T_i) \prod_{k \in \boldsymbol{J} \setminus (Q \cup x)} (1 - q(\boldsymbol{R}, T_k)) \le 0,$$
(88)

because $q(\mathbf{R}, T_i) \in [0, 1]$ for all $i \in \mathbf{K}$. Note that $\partial/\partial q(\mathbf{R}, T_x)[\mathbb{P}(N_{S(\mathbf{R})}(\mathbf{B}) \leq c | \mathbf{q}(\mathbf{R}, \mathbf{T}), \mathbf{J})]$ is independent of $q(\mathbf{R}, T_x)$.

(Proof by contradiction.) Suppose $J_{\theta_{syn}}^- \not\in \arg\min_{j \in \mathcal{D}} \{c^-(q(R,T),j)\}$. Fix any $J^* \in \arg\min_{j \in \mathcal{D}} c^-(q(R,T),j)$. Accordingly, assume there exists an $x^* \in J_{\theta_{syn}}^-$ such that $x^* \notin J^*$ so that $q(R, T_{x^*}) < q(R, T_m)$ for some $m \in J^*$ by the definition of $J_{\theta_{syn}}^-$. (If such an x^* does not exist, then $c^-(q(R,T), J^*) = c^-(q(R,T), J_{\theta_{syn}}^-)$ by construction.) By Eq. 88, this implies $c^-(q(R,T), J^*) > c^-(q(R,T), J^* \cup \{x^*\} \setminus \{m\})$. This is a contradiction as $J^* \cup \{x^*\} \setminus \{m\} \in \mathcal{D}$.

Lemma 2. Denote $X_1, X_2, ..., X_{i-1}, X_{i+1}, ..., X_n$ as ${}_iX$. Then suppose $X_1, X_2, ..., X_n$ are Bernoulli random variables (not necessarily independent) satisfying

$$\mathbb{P}(X_i = 1|_i X, Z) \ge p_i(Z),\tag{20}$$

for all $i \in \{1, 2, ..., n\}$, and for some random variable Z. Then

$$\mathbb{P}\left(\sum_{i=1}^{n} X_{i} \leq k \middle| Z\right) \leq \mathbb{P}\left(\sum_{i=1}^{n} Y_{i} \leq k\right), \forall k$$
(21)

if $Y_1, Y_2, ..., Y_n$ are conditionally independent Bernoulli random variables, conditioned on Z, with parameters $p_1(Z), p_2(Z), ..., p_n(Z)$ respectively, and $X_{1:n}$ and $Y_{1:n}$ are conditionally independent, conditioned on Z.

Proof. The proof is by induction. The case n = 1 is self-evident. Observe that

$$\mathbb{P}\left(\sum_{i=1}^{n} X_i \le k | Z\right) = \mathbb{P}\left(\sum_{i=1}^{n-1} X_i \le k \left| Z, X_n = 0\right\right) \left(1 - \mathbb{P}(X_n = 1 | Z)\right)$$

$$(89)$$

$$+ \mathbb{P}\left(\sum_{i=1}^{n-1} X_i \le k - 1 \middle| Z, X_n = 1\right) \mathbb{P}(X_n = 1)$$

$$(90)$$

and

$$\mathbb{P}\left(\sum_{i=1}^{n} Y_i \le k\right) = \mathbb{P}\left(\sum_{i=1}^{n-1} Y_i \le k\right) + \mathbb{P}\left(\sum_{i=1}^{n-1} Y_i \le k-1\right) p_n.$$
(91)

Conditioning on $X_n = 0$, the induction hypothesis is satisfied for $X_1, X_2, ..., X_{n-1}$. Therefore,

$$\mathbb{P}\left(\sum_{i=1}^{n-1} X_i \le k | Z, X_n = 0\right) \le \mathbb{P}\left(\sum_{i=1}^n Y_i \le k\right)$$
(92)

and analogously,

$$\mathbb{P}\left(\sum_{i=1}^{n-1} X_i \le k \middle| Z, X_n = 1\right) \le \mathbb{P}\left(\sum_{i=1}^{n-1} Y_i \ge k - 1\right).$$
(93)

Also note,

$$\left\{\sum_{i=1}^{n} Y_i \le k\right\} \subseteq \left\{\sum_{i=1}^{n} Y_i \le k - 1\right\} \implies \mathbb{P}\left(\sum_{i=1}^{n-1} Y_i \le k\right) \le \mathbb{P}\left(\sum_{i=1}^{n-1} Y_i \le k - 1\right)$$
(94)

and,

$$\mathbb{P}(X_n = 1) = \sum_{iX} \mathbb{P}(X_n = 1 | Z_{,i} X) \mathbb{P}(_i X) \le \sum_{iX} p_n \mathbb{P}(_i X) = p_n.$$
(95)

From this reasoning we obtain,

$$\mathbb{P}\left(\sum_{i=1}^{n} Y_i \le k\right) = \mathbb{P}\left(\sum_{i=1}^{n-1} Y_i \le k\right) (1-p_n) + \mathbb{P}\left(\sum_{i=1}^{n-1} Y_i \le k-1\right) p_n \tag{96}$$

$$\geq \mathbb{P}\left(\sum_{i=1}^{n-1} Y_i \le k\right) \left(1 - \mathbb{P}(X_n = 1)\right) + \mathbb{P}\left(\sum_{i=1}^{n-1} Y_i \le k - 1\right) \mathbb{P}(X_n = 1)$$

$$\tag{97}$$

$$\geq \mathbb{P}\left(\sum_{i=1}^{n-1} X_i \leq k \Big| Z, X_n = 0\right) (1 - \mathbb{P}(X_n = 1)) + \mathbb{P}\left(\sum_{i=1}^{n-1} X_i \leq k - 1 \Big| Z, X_n = 1\right) \mathbb{P}(X_n = 1)$$
(98)

$$= \mathbb{P}\left(\sum_{i=1}^{n} X_i \ge k \middle| Z\right).$$
⁽⁹⁹⁾

Proposition 1. Under A.1 (conditional uniformity), A.2 (timescale separation), and A.4 (consistency) $\{N_{S(\mathbf{R})}(\mathbf{T}) - h \leq c^{-}(\mathbf{q}(\mathbf{R},\mathbf{T}),\mathbf{J}_{h}^{-})\}$ and $\{N_{S(\mathbf{R})}(\mathbf{T}) - h \geq c^{+}(\mathbf{q}(\mathbf{R},\mathbf{T}),\mathbf{J}_{h}^{+})\}$ are $\alpha/2$ -level critical region for all \mathbb{P} in H_0 : $\theta_{syn} = h$. That is, for all \mathbb{P} in H_0 : $\theta_{syn} = h$,

$$\mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{T}) - h \le c^{-}(\boldsymbol{q}(\boldsymbol{R}, \boldsymbol{T}), \boldsymbol{J}_{h}^{-})\right) \le \alpha/2$$
(22)

$$\mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{T}) - h \ge c^{+}(\boldsymbol{q}(\boldsymbol{R}, \boldsymbol{T}), \boldsymbol{J}_{h}^{+})\right) \le \alpha/2.$$
(23)

and

Proof. We will prove the first inequality Eq. 22. Eq. 23 can be proved in the same way. Note that

$$\mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{T}) - h \le c^{-}(\boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J}) \middle| q(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J}\right) \le \alpha/2,$$
(100)

by the definition of $c^{-}(\cdot)$, and symmetry. Thus, by the fact that $J \in \mathcal{D}(\mathbf{R}, \mathbf{T})$, under H_0 and by Lemma 1,

$$c^{-}(\boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J}_{h}^{-}) \leq c^{-}(\boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J})$$
(101)

for all realizations of the model. We have,

$$\mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{T}) - h \le c^{-}(\boldsymbol{q}(\boldsymbol{R}, \boldsymbol{T}), \boldsymbol{J}_{h}^{-}) \middle| q(\boldsymbol{R}, \boldsymbol{T}), \boldsymbol{J}\right)$$
(102)

$$\leq \mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{T}) - h \leq c^{-}(\boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J}) \middle| \boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J}\right) \leq \alpha/2,$$
(103)

(almost surely). Therefore

$$\mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{T}) - h \le c^{-}(\boldsymbol{q}(\boldsymbol{R}, \boldsymbol{T}), \boldsymbol{J}_{h}^{-})\right)$$
(104)

$$= \mathbb{E}\left[\mathbb{P}\left(N_{S(\boldsymbol{R})}(\boldsymbol{T}) - h \leq c^{-}(\boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J}_{h}^{-}) \left| \boldsymbol{q}(\boldsymbol{R},\boldsymbol{T}),\boldsymbol{J}\right)\right]$$
(105)

$$\leq \alpha/2.$$
 (106)

Remark 2. We found the following intuition useful regarding the conditional inference in Proposition 1. It is not necessary that the critical region's boundary term, $c^-(q(\mathbf{R}, \mathbf{T}), \mathbf{J}_h^-)$, be $(q(\mathbf{R}, \mathbf{T}), \mathbf{J})$ -measurable, because of the inequality Eq. 101. Yet, a key feature is that it is (\mathbf{R}, \mathbf{T}) -measurable, so that, speaking informally, its evaluation does not require "knowledge" of \mathbf{J} .

6.3 Construction of multivariate Ornstein-Uhlenbeck process

To try to create reliable spike patterns, we utilize the Ornstein-Uhlenbeck construction of Taillefumier and Magnasco [57] as an injected current or light stimulus. We describe a multivariate version of the process because it is far more efficient to sample a multivariate version and then concatenate the dimensions into a longer trial. The process used for the former task is easily recovered in the one-dimensional case. These authors construct the Ornstein-Uhlenbeck process U(t) from discrete Haar-like basis functions. We focus on the discrete representation since the primary goal is to simulate the process. The whole process is divided into dyadic segments, with the following basis functions tiling the dyadic segments at various resolutions,

$$\lambda_{n,k}(t) \coloneqq \begin{cases} \frac{\sigma_H \cdot \sinh\left(|\alpha_H|(t-2k \cdot 2^{-n})\right)}{\sqrt{\alpha_H \sinh\left(\alpha_H 2^{1-n}\right)}} & \text{if } (2k)2^{-n} \le t < (2k+1)2^{-n} \\ \frac{\sigma_H \cdot \sinh\left(|\alpha_H|(2(k+1)2^{-n}-t)\right)}{\sqrt{\alpha_H \sinh\left(\alpha_H 2^{1-n}\right)}} & \text{if } (2k+1)2^{-n} \le t < 2(k+1)2^{-n} \\ 0 & \text{otherwise,} \end{cases}$$
(107)

for $0 \le 2k < 2^n$ with timescale τ_H , scaling parameters σ_H and α_H , and where

$$\lambda_{0,0}(t) = \frac{\sigma_H \cdot \exp(-\alpha_H/2)\sinh(|\alpha_H|t)}{\sqrt{\alpha_H \cdot \sinh(\alpha_H)}}.$$
(108)

A parameter $H \in [0, 1]$ describes how the amplitude of these basis functions scale with the resolution of the support, Δ_H . For a multi-dimensional version, for $1 \le i \le 2n - 1$ and $1 \le j \le n$, define the matrices,

$$\dot{M}_{n}(i,j) = \mathbb{1}\{i = 2j-1\} + \mathbb{1}\{i/2 = j\} + \mathbb{1}\{i/2 = j-1\}$$
(109)

$$\vec{U}_n(i,j) = \mathbb{1}\{i \text{ is even}\}\tag{110}$$

$$\vec{V}_n(i,j) = \mathbb{1}\{i \text{ is odd}\}.$$
(111)

For a multivariate process of dimension M and 2^N time points, let the value of the process at each time point be $\vec{x}_t \in \mathbb{R}^m$. For a resolution n, arrange the values of the process at the dyadic points $(0 \cdot 2^N/2^n, 0 \cdot 2^N/2^n, ..., 2^n \cdot 2^N/2^n)$ into distinct matrices indexed by n,

$$\vec{X}_{n} = \begin{bmatrix} \vec{x}_{0.2N/2^{n}} \\ \vec{x}_{1.2N/2^{n}} \\ \vdots \\ \vec{x}_{2^{n}.2^{N}/2^{n}} \end{bmatrix}.$$
(112)

For numerical reasons, the endpoints are clamped such that $\vec{X}_0 = [0]_{2,m}$ and the full process \vec{X}_N can then be computed for $0 \le n \le N-1$ efficiently by the recursive matrix operations,

$$\underbrace{\overbrace{\substack{\vec{x}_{n+1}\\\vec{x}_{1\cdot2^{N}/2^{n+1}}\\\vec{x}_{2\cdot2^{N}/2^{n+1}}\\\vdots\\\vec{x}_{2^{n+1}\cdot2^{N}/2^{n+1}}\end{bmatrix}}_{\vec{x}_{2\cdot2^{N}/2^{n+1}}} = \left(\frac{1}{2}\cosh^{-1}\left(\Delta_{H}/2\tau_{H}\right)\vec{U}_{n+1} + \vec{V}_{n+1}\right)\odot\vec{M}_{n+1} \underbrace{\overbrace{\substack{\vec{x}_{0\cdot2^{N}/2^{n}}\\\vec{x}_{1\cdot2^{N}/2^{n}}\\\vdots\\\vdots\\\vec{x}_{2^{n}\cdot2^{N}/2^{n}}\end{bmatrix}}_{\vec{x}_{2\cdot2^{N}/2^{n}}} \right) (113)$$

$$+ \alpha_{H} [\tau_{H} \tanh (\Delta_{H}/2\tau_{H})]^{H} \sqrt{\frac{1}{\tau_{H}}} \underbrace{\begin{bmatrix} \vec{0}_{m} \\ \vec{\xi}_{1\cdot2^{N}/2^{n+1}} \\ \vec{0}_{m} \\ \vec{\xi}_{3\cdot2^{N}/2^{n+1}} \\ \vdots \\ \vec{\xi}_{(2^{n+1}-1)\cdot2^{N}/2^{n+1}} \\ \vec{0}_{m} \\ \vec{\Xi}_{n+1} \end{bmatrix}}_{\vec{\Xi}_{n+1}}$$
(114)

where for fixed $n \vec{\Xi}_n$ contains in its odd rows *iid* normal random vectors, $\vec{\xi}_t$, obeying $\mathcal{N}(\vec{0}_m, \vec{\Sigma}_{m,m}(n))$ and zeros otherwise. That is, each resolution n has a characteristic dependence structure parameterized by the covariance matrices $\vec{\Sigma}_{m,m}(n)$ for $1 \leq n \leq N$. $\vec{\Sigma}_{m,m}(0)$ is not defined since the endpoints are clamped. For background inputs into the HH-type model system, we used a constant covariance matrix for all n. As was previously done in the text, the Vine Beta method, with its parameter fixed to 0.1, was used to generate a random covariance matrix with strong positive and negative associations in three dimensions. Sampling is made much more efficient by the following. First, for a process of dimension M sample Eq.(114) instead for a process of $M \cdot N_{trial}$ with the covariance matrix for the M-dimensional process replaced with the block matrix $\vec{\Sigma}_{N_{trial}m,N_{trial}m}(n) = \vec{I} \otimes \Sigma_{m,m}(n)$ where \otimes is the Kronecker product and \vec{I} the identity matrix. Implementing Eq.(114) recursively is efficient up to matrices of moderate size, and finally, then every j-th column for $j \in \{1, 2, ...m\}$ of X_{2N} can be concatenated together N_{trial} times.

When we used the above process as a stimulus in the simulated neural perturbation experiments, both H and τ_H were varied as indicated in Section 5.2.

6.4 Rationale for monosynaptic confidence interval algorithm

Synaptic inference is commonly performed on low-dimensional and easy-to-visualize objects, the CCG in particular, and the most immediate objection to the theory outlined in previous sections is that it may be computationally prohibitive. Point estimation via Eq. (9) is quite simple, but for confidence intervals, this objection may be reasonable because inference is performed on the sequence of spike counts and synchrony in small temporal intervals, which is of much higher dimension and grows with the duration of the spike trains. Earlier, this objection was neutralized via a principled algorithm. Here, we describe the algorithms in prose and highlight their rationale.

Given an observation R and T all the probabilities in Eq. (28) can be obtained as well as \tilde{J}_h^- and \tilde{J}_h^+ for any h. As mentioned in the main text, a naive but easy-to-grasp strategy to compute confidence intervals would be to begin with

the two-tailed hypothesis $H_0: \theta_{syn} = 0$ enacted through convolution of the distributions with success probabilities $q(\mathbf{R}, Z_i)$ for $i \in \tilde{J}_0^-$ and then evaluation of tail areas given the observed value of $N_{S(\mathbf{R})}(\mathbf{T})$. Then, if we reject the null hypothesis at the upper tail proceed to test positive values for $H_0: \theta_{syn} = h$ in the sequence $h \in (1, 2, 3, ...)$ recomputing the tail probability each time for distributions arising from \tilde{J}_h^- until we fail to reject. If, instead, we fail to reject the null hypothesis at the left tail, proceed to test negative values $h \in (-1, -2, -3, ...)$ until we fail to reject. The process then needs to be repeated for \tilde{J}_h^+ except for the case h = 0. Each value h tested is expensive because we must compute a sum of independent random indicators and a central question is how to reuse computations most effectively across these tests. For most data of reasonable size a much faster strategy is to apply standard binary search [63] to the sequence $(0, 1, 2, ..., N_{S(\mathbf{R})}(\mathbf{T}))$ (i.e., the values of h to test) with the search query being the location of adjacent values in the sequence for which one and not the other fails to reject the null hypothesis $H_0: \theta_{syn} = h$.

The next question is how exactly to convolve the distributions that emerge from every step of binary search so that the tail area can be evaluated; that is, we must compute the *cdf* of a sum of independent but not necessarily identically distributed indicators. While in practice this is often computed with the *Fast-Fourier transform* (FFT), FFT can have very large relative errors for small tail probabilities [64]. In contrast, *direct convolution* (DC) is the most accurate method and uses only the convolution definition of the distribution function of a sum of random variables. While its runtime complexity is $O(N^2)$, typically motivating the use of FFT, Biscarri et al. [65] observe that DC is still most efficient when convolving a small number of vectors or many vectors of small dimensions (e.g., Bernoulli vectors) suggesting DC and FFT can work in concert in a divide-and-conquer scheme. Through both theoretical considerations and experimentation, we adopt a similar mixed approach, also taking inspiration from Peres et al. [66]. These studies address the general problem of convolving independent indicators, while the problem here is to construct confidence intervals, and there are various other domain-specific constraints that we can exploit. For example, a special feature of our problem is we may assume many of the random indicators to be convolved will have equal success probabilities [67].

The overview is as follows. First, as a first pass, we compute a conservative confidence interval using Chernoff bounds to estimate tail probabilities at every iteration of binary search. Then, in a second pass of binary search, we refine the Chernoff confidence interval on a much-narrowed search space by computing tail probabilities with a mix of DC and FFT accompanied by a method for recovering the relative accuracy of FFT via exponential tilting [68]. The strategy also exploits redundant structure in various ways to minimize computations. We digress to quickly introduce Chernoff's bound to those unfamiliar.

Remark 3. The following is a well-known result. Let $X_i \sim Be(p_i)$ for $i \in \{1, 2, ...n\}$. Consider the condition the sum is above some bound t,

$$\sum_{i} X_i \ge t. \tag{115}$$

By multiplying through by a constant λ , exponentiating, and applying Markov's inequality the generic Chernoff bound is obtained,

$$P\left(\sum_{i} X_{i} \ge t\right) \le exp(-\lambda t) \mathbb{E}[exp(\lambda \sum_{i} X_{i})].$$
(116)

In the case of sums of independent indicators a bit more work can yield the result,

$$P\left(\sum_{i} X_{i} \ge t\right) \le (\mu e/t)^{t} e^{-\mu}$$
(117)

such that $ln(t/\mu) = \lambda$ and $\mu = \sum_i p_i$ [69].

Denote \tilde{C}_{CF} as confidence intervals for θ_{syn} obtained by using Chernoff's bound as a tail probability estimate. If we substitute computation of the exact *cdf* with Chernoff's bounds, it is guaranteed that $\tilde{C}_{CF} \supseteq \tilde{C}$. Furthermore, unlike the exact *cdf*, iterative computation of Chernoff's bound is very cheap for binary search on the sequence $(1, 2, ..., N_{S(R)}(T))$ relative to its imprecision.

After \tilde{C}_{CF} is calculated, the search space is significantly reduced and we can now choose a more accurate method to implement a second pass of binary search on the sequence $(CF_L, CF_L + 1, ..., CF_U - 1, CF_U)$ where CF_L and CF_U are the lower and upper confidence bounds for θ_{syn} obtained via Chernoff's bound. The first pass with Chernoff's bound also guarantees that we need the lower confidence interval to be at least CF_L , meaning we are assured that at least CF_L random variables will not need to be convolved in the second pass. On the other hand, the first pass assures that at least $N_{S(\mathbf{R})}(\mathbf{T}) - CF_U + |\mathbf{T} \setminus S(\mathbf{R})|$ random variables will need to be convolved and hence we can design an algorithm that only computes those convolutions once then reuses the result.

We will describe the process for computing exact confidence intervals that ensues for the case of h > 0 and for the tail corresponding to \tilde{J}_h^+ . The other cases are analogous, and we will conclude by highlighting alterations. For the remainder of this section, we will describe the strategy in words to supplement the description with rationale. Where ambiguity might be present here, the reader can refer to the precise description as summarized in Algorithm 1 and Algorithm 2.

Given CF_L and CF_U , the $N_{S(\mathbf{R})}(\mathbf{T}) - CF_U + |\mathbf{T} \setminus S(\mathbf{R})|$ random indicators that Chernoff tells us must be summed are those with success probabilities $q(Z_i)$ for $i \in \tilde{J}_{CF_U}^+$. The second pass of binary search will test values for h in $(CF_L, CF_L + 1, ..., CF_U - 1, CF_U)$ and, for each new value h, the distributions with success probabilities $q(Z_i)$ for $i \in \tilde{J}_h^+ \setminus \tilde{J}_{CF_U}^+$ will need to be appended to previous computations to obtain the tail probability.

Let us now ask how we can efficiently compute the initial convolution that Chernoff tells us must contribute to the final result and is quite likely to comprise the bulk of the final result, with the foresight that subsequent iterations will need to reincorporate this bulk into various distinct new tail probability computations until the second pass of binary search halts. We may ask if we can choose a number L_{div} such that the success probabilities $q(Z_i)$ for $i \in \tilde{J}_{CF_U}^+$ can be divided into L_{div} groups. In particular, we wish to divide them into L_{div} groups such that each group contains precisely the same mixture of constituent success probabilities. To reiterate once more, we require that each group is an identical mixture of perhaps unequal Bernoulli vectors. If we then calculate the sum of the random indicators within each of L_{div} groups, the result will be L_{div} *iid* distributions as an intermediary step. This strategy has three advantages that easily compensate for the additional overhead of determining the L_{div} groups. First, we only need to compute the within-group sum of indicators once, not L_{div} times, since each group is *iid*. Second, this one convolution needed will contain a small number of random variables of low dimensions to the extent L_{div} is large, and thus, we can exploit DC in the regime it is efficient [65] with great payoff. Third, since the L_{div} groups are *iid*, to obtain the distribution for the sum over groups, we can use highly efficient L_{div} -fold convolution power, meaning we need only compute one FFT rather than L_{div} FFTs to obtain the final result (L_{div} also cannot be too large to avoid numerical errors in FFT but the algorithm has robust performance for a significant range, say $L_{div} \in [2^2, 2^6]$).

Finding L_{div} *iid* groups is easy and useful to the extent that the initial mixture of success probabilities contains a small number of unique values relative to the total number of indicators. This condition is highly applicable to the scientific context. Clearly, if every initial success probability were unique, it would be impossible to divide them into L_{div} groups giving rise to L_{div} *iid* intermediary sums. Even if there are many repeated success probabilities in the initial mixture, it is unlikely to find a L_{div} in a helpful range to achieve the objective since that requires the frequencies of each unique occurring Bernoulli vector to have L_{div} as a common divisor. In contrast, it is easy to split the initial mixture into L_{div} *iid* groups if we tolerate one group of "residual" success probabilities to be triaged and dealt with as a special case. This is the approach we take.

From the probabilities $q(Z_i)$ for $i \in \tilde{J}_{CF_U}^+$, compute the unique success probabilities and the frequency of occurrence for each unique vector. Then divide each of the frequencies by L_{div} and round down. The resulting numbers will be how many random indicators of each unique success probability will contribute to one of L_{div} intermediary sums. Rather than using DC to convolve all the Bernoulli vectors in this one group, instead firstly directly compute binomial distributions for the underlying unique groups of *iid* Bernoulli vectors with the number of trials for each binomial as the frequency of the unique occurring Bernoulli vectors divided by L_{div} rounded down. Then, use DC on the resulting binomial vectors; at this stage, DC will still be efficient if L_{div} is large enough. Denote \vec{p} as this resulting probability vector that will be the input to L_{div} -fold convolution downstream. In the steps just described, the operation of dividing the frequency of each unique vector in the initial mixture by L_{div} and rounding down produces "leftovers" - to become the triaged group - because of rounding. Denote the probability vector of the triaged convolution sum as \vec{a} . This concludes the preparation for the second pass; the distributions with success probabilities $q(Z_i)$ for $i \in \tilde{J}_{CF_U}^+$ are now encoded in \vec{a} , \vec{p} , and L_{div} .

During the second pass of binary search, \vec{a} is itself immutable but convolved per iteration with distributions that arise from new hypotheses $H_0: \theta_{syn} = h$ producing as a result a temporary vector \vec{g} . Specifically, for each new h, \vec{g} is obtained by convolving \vec{a} with distributions with success probabilities $q(Z_i)$ for $i \in \tilde{J}_h^+ \setminus \tilde{J}_{CF_U}^+$. These new distributions are convolved with \vec{a} using DC since the dimensions are likely to still be small to the extent Chernoff's bound is close to the exact answer.

Finally, for each h tested on the second pass, \vec{p} and \vec{g} are sent to an FFT-based step with an exponential tilt to correct for the errors typically induced by FFT. Denote the exponentially tilted versions of \vec{p} and \vec{g} as $\vec{p_s}$ and $\vec{g_s}$, respectively. Specifically, the convolution for each h occurs via L_{div} -fold convolution power for $\vec{p_s}$ and point-wise multiplication of the result with $\vec{g_s}$ in the frequency domain (each embedded in the dimension of the final result). It is important to mention that one might argue we only need to compute the FFT of $\vec{p_s}$ once after the first pass, and hence $\vec{p_s}$ can then sit idly in the frequency domain during the second pass waiting to be point-wise multiplied with new versions of $\vec{g_s}$ for each new test *h* encountered by the second pass of binary search. However, this does not work because every new convolution requires a new exponential tilting parameter. That is, $\vec{p_s}$ is a temporary vector that changes each iteration because a new tilting parameter is required. While this prohibits such reuse [68], it offers extraordinarily large increases in accuracy in return. We prioritize accuracy over this decrease in speed at this juncture. Overall, the algorithm is still competitively fast.

The iterative FFT convolution with exponential tilting is detailed and given its own description in Algorithm 1. The overall procedure is then summarised in Algorithm 2 for an excitatory lower confidence bound. There are three more cases: the excitatory upper bound as well as inhibitory lower and upper bounds. The other cases are very similar but we quickly remark here on the differences. First, a lower bound of θ_{syn} corresponds to the right tail probability of the test statistic whereas an upper bound for θ_{syn} corresponds to a left tail probability of the test statistic. Exponential tilting is suitable for computing right tail probabilities, but for any random variable X the left tail probability can be obtained by computing the right tail probability of -X. Last, the initial implementation of binary search using Chernoff's bound for excitatory intervals searches h on the sequence $(1, 2, ..., N_{S(R)}(T))$ but for inhibition the sequence searched is (1, 2, ..., |G|).

Algorithm 2 begins with the assumption that the null hypothesis $H_0: \theta_{syn} = 0$ has already been rejected and code posted online tests $H_0: \theta_{syn} = 0$ and breaks if it fails to be rejected. The philosophy behind this decision is that in large-scale neural recordings, most neurons will be unconnected, and computing confidence intervals for all pairwise interactions will be expensive and not insightful. A very rapid approximation, such as the normal approximation, can be used for the test $H_0: \theta_{syn} = 0$, and a p-value can be calibrated for the precision a specific question requires. The confidence intervals can then be applied to significant pairs and still may include 0 when the exact method is employed. Code posted online also gives the option to return confidence intervals computed with p-values obtained from the normal approximation which is much faster for large spike trains. Preliminary numerical experiments suggest spike trains of moderate size will yield normal approximations very close to the exact solution. However, the proximity to the exact solution depends upon many factors and the tolerance for error depends on the scientific question, necessitating specific use case evaluations.

6.5 Simulation software and hardware

To ensure careful handling of frozen noise inputs, all numerical simulations of dynamical systems were programmed from scratch and integrated with modified Euler's method in Python. Long simulations were optimized often through parallelization and concatenation. This was done in Google Colab using the NVIDIA T4 GPU.

Acknowledgements

We thank Jonathan Platkiewicz, Gyorgy Buzsaki, Daniel English, Uri Keich, and Thibaud Taillefumier for helpful conversations that guided this research.

References

- [1] Jin Tian and Judea Pearl. Probabilities of causation: Bounds and identification. *Annals of Mathematics and Artificial Intelligence*, 28(1-4):287–313, 2000.
- [2] Daniel Fine English, Sam McKenzie, Talfan Evans, Kanghwan Kim, Euisik Yoon, and György Buzsáki. Pyramidal cell-interneuron circuit architecture and dynamics in hippocampal networks. *Neuron*, 96(2):505–520, 2017.
- [3] Jozsef Csicsvari, Hajime Hirase, Andras Czurko, and György Buzsáki. Reliability and state dependence of pyramidal cell-interneuron synapses in the hippocampus: an ensemble approach in the behaving rat. *Neuron*, 21 (1):179–189, 1998.
- [4] Jean-Sébastien Jouhanneau, Jens Kremkow, and James FA Poulet. Single synaptic inputs drive high-precision action potentials in parvalbumin expressing gaba-ergic cortical neurons in vivo. *Nature Communications*, 9(1): 1–11, 2018.
- [5] Shigeyoshi Fujisawa, Asohan Amarasingham, Matthew T Harrison, and György Buzsáki. Behavior-dependent short-term assembly dynamics in the medial prefrontal cortex. *Nature neuroscience*, 11(7):823–833, 2008.

- [6] David Marc Anton Mehler and Konrad Paul Kording. The lure of causal statements: Rampant mis-inference of causality in estimated connectivity. *arXiv e-prints*, pages arXiv–1812, 2018.
- [7] Jonathan Platkiewicz, Zachary Saccomano, Sam McKenzie, Daniel English, and Asohan Amarasingham. Monosynaptic inference via finely-timed spikes. *Journal of Computational Neuroscience*, 49(2):131–157, 2021.
- [8] Mikkel Elle Lepperød, Tristan Stöber, Torkel Hafting, Marianne Fyhn, and Konrad Paul Kording. Inferring causal connectivity from pairwise recordings and optogenetics. *PLOS Computational Biology*, 19(11):e1011574, 2023.
- [9] Ian H Stevenson, James M Rebesco, Lee E Miller, and Konrad P Körding. Inferring functional connections between neurons. *Current opinion in neurobiology*, 18(6):582–588, 2008.
- [10] Asohan Amarasingham, Ting-Li Chen, Stuart Geman, Matthew T Harrison, and David L Sheinberg. Spike count reliability and the poisson hypothesis. *Journal of Neuroscience*, 26(3):801–809, 2006.
- [11] Asohan Amarasingham, Matthew T Harrison, Nicholas G Hatsopoulos, and Stuart Geman. Conditional modeling and the jitter method of spike resampling. *Journal of Neurophysiology*, 107(2):517–531, 2012.
- [12] Asohan Amarasingham, Stuart Geman, and Matthew T Harrison. Ambiguity and nonidentifiability in the statistical analysis of neural codes. *Proceedings of the National Academy of Sciences*, 112(20):6455–6460, 2015.
- [13] Jerzy S Neyman. On the application of probability theory to agricultural experiments. essay on principles. section 9.(tlanslated and edited by dm dabrowska and tp speed, statistical science (1990), 5, 465-480). Annals of Agricultural Sciences, 10:1–51, 1923.
- [14] Donald B Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology*, 66(5):688, 1974.
- [15] Guido W Imbens and Donald B Rubin. Causal inference in statistics, social, and biomedical sciences. Cambridge University Press, 2015.
- [16] Larry Wasserman. All of statistics: a concise course in statistical inference. Springer Science & Business Media, 2013.
- [17] Judea Pearl. Causality. Cambridge university press, 2009.
- [18] Judea Pearl, Madelyn Glymour, and Nicholas P Jewell. *Causal inference in statistics: A primer*. John Wiley & Sons, 2016.
- [19] Lidor Spivak, Amir Levi, Hadas E Sloin, Shirly Someck, and Eran Stark. Deconvolution improves the detection and quantification of spike transmission gain from spike trains. *Communications Biology*, 5(1):520, 2022.
- [20] Ian H Stevenson. Circumstantial evidence and explanatory models for synapses in large-scale spike recordings. *arXiv preprint arXiv:2304.09699*, 2023.
- [21] Robert E Kass, Valérie Ventura, and Emery N Brown. Statistical issues in the analysis of neuronal data. *Journal* of neurophysiology, 94(1):8–25, 2005.
- [22] Daryl J Daley and David Vere-Jones. An Introduction to the Theory of Point Processes. Volume II: General Theory and Structure. Springer, 2008.
- [23] Ryota Kobayashi, Shuhei Kurita, Anno Kurth, Katsunori Kitano, Kenji Mizuseki, Markus Diesmann, Barry J Richmond, and Shigeru Shinomoto. Reconstructing neuronal circuitry from parallel spike trains. *Nature communications*, 10(1):1–13, 2019.
- [24] Asohan Amarasingham, Matthew T Harrison, Nicholas G Hatsopoulos, and Stuart Geman. Conditional modeling and the jitter method of spike re-sampling: Supplement. *arXiv preprint arXiv:1111.4296*, 2011.

- [25] Matthew T Harrison, Asohan Amarasingham, and Wilson Truccolo. Spatiotemporal conditional inference and hypothesis tests for neural ensemble spiking precision. *Neural computation*, 27(1):104–150, 2015.
- [26] Judea Pearl. Brief report: On the consistency rule in causal inference:" axiom, definition, assumption, or theorem?". *Epidemiology*, pages 872–875, 2010.
- [27] Tyler J VanderWeele. Concerning the consistency assumption in causal inference. *Epidemiology*, 20(6):880–883, 2009.
- [28] Stephen R Cole and Constantine E Frangakis. The consistency statement in causal inference: a definition or an assumption? *Epidemiology*, 20(1):3–5, 2009.
- [29] Miguel A Hernán and James M Robins. Causal inference, 2010.
- [30] Shigeru Shinomoto, Yutaka Sakai, and Shintaro Funahashi. The ornstein-uhlenbeck process does not reproduce spiking statistics of neurons in prefrontal cortex. *Neural Computation*, 11(4):935–951, 1999.
- [31] W. R. Softky and C. Koch. The highly irregular firing of cortical cells is inconsistent with temporal integration of random epsps. *J Neurosci*, 13(1):334–350, Jan 1993.
- [32] George Casella and Roger L Berger. Statistical inference, volume 2. Duxbury Pacific Grove, CA, 2002.
- [33] Asohan Amarasingham. *Statistical methods for assessment of temporal structure in the activity of the nervous system.* PhD thesis, Brown University, Providence, RI, 2004.
- [34] Peter Barthó, Hajime Hirase, Lenaïc Monconduit, Michael Zugaro, Kenneth D Harris, and György Buzsáki. Characterization of neocortical principal cells and interneurons by network interactions and extracellular features. *Journal of neurophysiology*, 92(1):600–608, 2004.
- [35] Hiroshi Tamura, Hidekazu Kaneko, Keisuke Kawasaki, and Ichiro Fujita. Presumed inhibitory neurons in the macaque inferior temporal cortex: visual response properties and functional interactions with adjacent neurons. *Journal of neurophysiology*, 91(6):2782–2796, 2004.
- [36] Michael Wehr and Anthony M Zador. Balanced inhibition underlies tuning and sharpens spike timing in auditory cortex. *Nature*, 426(6965):442–446, 2003.
- [37] Frédéric Pouille and Massimo Scanziani. Enforcement of temporal fidelity in pyramidal cells by somatic feed-forward inhibition. *Science*, 293(5532):1159–1163, 2001.
- [38] Barna Dudok, Miklos Szoboszlay, Anirban Paul, Peter M Klein, Zhenrui Liao, Ernie Hwaun, Gergely G Szabo, Tristan Geiller, Bert Vancura, Bor-Shuen Wang, et al. Recruitment and inhibitory action of hippocampal axoaxonic cells during behavior. *Neuron*, 2021.
- [39] Jonathan Platkiewicz, Eran Stark, and Asohan Amarasingham. Spike-centered jitter can mistake temporal structure. *Neural computation*, 29(3):783–803, 2017.
- [40] Artur Luczak, Peter Barthó, Stephan L Marguet, György Buzsáki, and Kenneth D Harris. Sequential structure of neocortical spontaneous activity in vivo. *Proceedings of the National Academy of Sciences*, 104(1):347–352, 2007.
- [41] Tomáš Hromádka, Anthony M Zador, and Michael R DeWeese. Up states are rare in awake auditory cortex. *Journal of Neurophysiology*, 109(8):1989–1995, 2013.
- [42] Kenneth D Harris, Jozsef Csicsvari, Hajime Hirase, George Dragoi, and György Buzsáki. Organization of cell assemblies in the hippocampus. *Nature*, 424(6948):552–556, 2003.
- [43] Adelchi Azzalini and Antonella Capitanio. Statistical applications of the multivariate skew normal distribution. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 61(3):579–602, 1999.

- [44] Michael R DeWeese and Anthony M Zador. Non-gaussian membrane potential dynamics imply sparse, synchronous activity in auditory cortex. *J Neurosci*, 26(47):12206–12218, 2006.
- [45] Daniel Lewandowski, Dorota Kurowicka, and Harry Joe. Generating random correlation matrices based on vines and extended onion method. *Journal of multivariate analysis*, 100(9):1989–2001, 2009.
- [46] Wulfram Gerstner and Richard Naud. How good are neuron models? Science, 326(5951):379–380, 2009.
- [47] Alix Herrmann and Wulfram Gerstner. Noise and the psth response to current transients: Ii. integrate-and-fire model with slow recovery and application to motoneuron data. *Journal of Computational Neuroscience*, 12(2): 83–95, 2002.
- [48] TS Otis and I Mody. Modulation of decay kinetics and frequency of gabaa receptor-mediated spontaneous inhibitory postsynaptic currents in hippocampal neurons. *Neuroscience*, 49(1):13–32, 1992.
- [49] Romain Brette and Wulfram Gerstner. Adaptive exponential integrate-and-fire model as an effective description of neuronal activity. *Journal of neurophysiology*, 94(5):3637–3642, 2005.
- [50] Yann Zerlaut, Sandrine Chemla, Frederic Chavane, and Alain Destexhe. Modeling mesoscopic cortical dynamics using a mean-field model of conductance-based networks of adaptive exponential integrate-and-fire neurons. *Journal of computational neuroscience*, 44(1):45–61, 2018.
- [51] Naixin Ren, Shinya Ito, Hadi Hafizi, John M Beggs, and Ian H Stevenson. Model-based detection of putative synaptic connections from spike recordings with latency and type constraints. *Journal of Neurophysiology*, 124 (6):1588–1604, 2020.
- [52] WR Levick, BG Cleland, MW Dubin, et al. Lateral geniculate neurons of cat: retinal inputs and physiology. *Invest Ophthalmol*, 11(5):302–311, 1972.
- [53] Moshe Abeles. Corticonics: Neural circuits of the cerebral cortex. Cambridge University Press, 1991.
- [54] Ang Li, Song Jiang, Yizhou Sun, and Judea Pearl. Learning probabilities of causation from finite population data. arXiv preprint arXiv:2210.08453, 2022.
- [55] John C Williams, Jianjin Xu, Zhongju Lu, Aleksandra Klimas, Xuxin Chen, Christina M Ambrosi, Ira S Cohen, and Emilia Entcheva. Computational optogenetics: empirically-derived voltage-and light-sensitive channelrhodopsin-2 model. *PLoS computational biology*, 9(9):e1003220, 2013.
- [56] Thibaud Taillefumier and Marcelo Magnasco. A transition to sharp timing in stochastic leaky integrate-and-fire neurons driven by frozen noisy input. *Neural computation*, 26(5):819–859, 2014.
- [57] Thibaud Taillefumier and Marcelo O Magnasco. A haar-like construction for the ornstein uhlenbeck process. *Journal of Statistical Physics*, 132:397–415, 2008.
- [58] Zachary F Mainen and Terrence J Sejnowski. Reliability of spike timing in neocortical neurons. Science, 268 (5216):1503–1506, 1995.
- [59] Amir Levi, Lidor Spivak, Hadas E Sloin, Shirly Someck, and Eran Stark. Error correction and improved precision of spike timing in converging cortical networks. *Cell Reports*, 40(12), 2022.
- [60] George P Moore, Jose P Segundo, Donald H Perkel, and Herbert Levitan. Statistical signs of synaptic interaction in neurons. *Biophysical journal*, 10(9):876–900, 1970.
- [61] Srdjan Ostojic, Germán Szapiro, Eric Schwartz, Boris Barbour, Nicolas Brunel, and Vincent Hakim. Neuronal morphology generates high-frequency firing resonance. *The Journal of Neuroscience*, 35(18):7056–7068, 2015.
- [62] Joanna Pressley and Todd W Troyer. The dynamics of integrate-and-fire: Mean versus variance modulations and dependence on baseline parameters. *Neural computation*, 23(5):1234–1247, 2011.

- [63] Derrick H Lehmer. Teaching combinatorial tricks to a computer. Combinatorial Analysis, pages 179–193, 1960.
- [64] Uri Keich. sfft: a faster accurate computation of the p-value of the entropy score. *Journal of Computational Biology*, 12(4):416–430, 2005.
- [65] William Biscarri, Sihai Dave Zhao, and Robert J Brunner. A simple and fast method for computing the poisson binomial distribution function. *Computational Statistics & Data Analysis*, 122:92–100, 2018.
- [66] Noah Peres, Andrew Ray Lee, and Uri Keich. Exactly computing the tail of the poisson-binomial distribution. *ACM Transactions on Mathematical Software (TOMS)*, 47(4):1–19, 2021.
- [67] Daniel Jeck and Ernst Niebur. Closed form jitter methods for neuronal spike train analysis. In 2015 49th Annual Conference on Information Sciences and Systems (CISS), pages 1–3. IEEE, 2015.
- [68] Huon Wilson and Uri Keich. Accurate pairwise convolutions of non-negative vectors via fft. *Computational Statistics & Data Analysis*, 101:300–315, 2016.
- [69] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.