# Causal Emergence 2.0: Quantifying emergent complexity

Erik Hoel[*1]

[1]Allen Discovery Center, Tufts University, Medford, MA, USA

March 18, 2025

## Abstract

Complex systems can be described at myriad different scales, and their causal workings often have multiscale structure (e.g., a computer can be described at the microscale of its hardware circuitry, the mesoscale of its machine code, and the macroscale of its operating system). While scientists study and model systems across the full hierarchy of their scales, from microphysics to macroeconomics, there is debate about what the macroscales of systems can possibly add beyond mere compression. To resolve this longstanding issue, here a new theory of emergence is introduced wherein the different scales of a system are treated like slices of a higher-dimensional object. The theory can distinguish which of these scales possess unique causal contributions, and which are not causally relevant. Constructed from an axiomatic notion of causation, the theory's application is demonstrated in coarse-grains of Markov chains. It identifies all cases of macroscale causation: instances where reduction to a microscale is possible, yet lossy about causation. Furthermore, the theory posits a causal apportioning schema that calculates the causal contribution of each scale, showing what each uniquely adds. Finally, it reveals a novel measure of emergent complexity: how widely distributed a system's causal workings are across its hierarchy of scales.

# 1 Introduction

Complex systems operate across scales, and therefore evince a huge number of possible descriptions [1]. This embarrassment of multiplicity necessitates

---
[*]erik.hoel@tufts.edu

a formal mathematical theory of emergence. Such a theory should explain and quantify how macroscales (broadly, dimension reductions) contribute to a system's causal workings. A theory of emergence may even explain the large-scale spatiotemporal hierarchy of the sciences themselves, beyond their function of just useful compressions [2].

It might be protested that there is no room for emergence in science, as presumably the future of any given system can be predicted with full knowledge of its microscale, and presumably any given system can be reduced to its microscale. However, prediction is not the same thing as causation [3]. A toy example is that of a thermostat and room system [4, 5]. While in theory the microscale of all the individual particles in the room could be used to predict the thermostat's reading, in terms of causal understanding it represents a poor answer to the question of "What caused the thermostat to read $20°C$?" In fact, the exact microstate of all the particles is not causally necessary for the reading, since many other configurations could lead to it. Meanwhile, the macrostate (the temperature of the room) has a direct causal relationship to the thermostat's reading, in that it is necessary for any given value.

In another example, an incoming signal to a neuron's dendrites could be used to predict a downstream action potential will occur. Yet, in terms of causal analysis, the incoming signal would be insufficient to trigger, as an effect, some exact exchange of ions (as these would evolve unpredictably due to noise, such as from Brownian motion [6] or quantum effects [7]). Meanwhile, the incoming signal could still be deterministically sufficient to trigger, as an effect, the downstream neuronal macrostate of "firing."

Following these intuitions, in 2013 I and my co-authors introduced the theory of causal emergence [8]. The theory made use of discrete causal models (in the class of logic gate networks, DAGs, and Markov chains), and a measure of causation, the effective information (EI), defined as the mutual information following a uniform intervention distribution of the $do(x)$ operator [9, 10]. The theory provided a toolkit to search across all possible dimension reductions of such systems to find the one that maximized the EI. The total increase in EI quantified the degree of *causal emergence* in the system. Such causal emergence explains why macroscales of a system can have stronger causal relationships even while being reducible to their underlying macroscales: since macroscales are multiply-realizable, they can minimize the uncertainty in causal relationships, which a measure of causation like the EI is sensitive to. This is mathematically similar to how coding over an information channel can minimize the noise of communication [2, 9].

The theory of causal emergence has spawned a large amount of research,

such as measuring causal emergence in data spanning from cellular automata [11] to fMRI data [12] to gene regulatory networks[13, 14], as well as developing heuristics [15], like detecting causal emergence with trained artificial neural networks [16]. It's been related to phenomena like scale-freeness and robustness in network theory [17, 18], been adapted within Integrated Information Theory [19, 20], and there have been proposals for alternatives of what measure should be used for quantifying causal emergence, such using the reversibility of a system to approximate the EI [21]. For a full review, see Yuan et al. 2024 [22].

However, so far the understanding of causal emergence has been incomplete, due to two outstanding issues. The first is the reliance on the EI and its approximates to detect causal emergence. While the EI is a relatively well-constructed measure of causation [23], it makes strong background assumptions in its calculations (such as requiring a uniform distribution of interventions, which some have questioned) [24, 25, 26]. Additionally, as demonstrated here in Section 6.2, the use of EI actually underestimates causal emergence.

The second issue is that the theory only identified a single causally-relevant scale (the maximum of EI), ignoring all multiscale structure. Yet many systems seem to operate with different scales contributing; one prominent example would be proposals that the brain has different functional scales, ranging from the neuronal up to cortical minicolumns up to entire brain regions [27]. Another example is how a computer can be described at the microscale of its hardware circuitry, the mesoscale of its machine code software, or at the macroscale of its operating system and applications [28]; indeed, even what computations are occurring may change depending on the scale of description [29, 30].

To develop a universal and well-grounded theory of causal emergence that resolves these issues, here I introduce a novel formalization: Causal Emergence 2.0. The fundamental intuition of CE 2.0 is that a system is not bound to one particular scale of description, rather, it is best described by the set of scales that contribute to the system's causal workings. Any one scale (even the microscale) is much like taking a 2D slice of a 3D object, and therefore cannot fully capture the causation of the system. CE 2.0 posits a *causal apportioning schema* across scales that detects their causal contributions (if any) to the multiscale whole.

As an updated theory of causal emergence, CE 2.0 uses some of the same mathematical terms as CE 1.0, but it grounds the theory in an axiomatic notion of causation, which allows the theory to capture all cases of macroscale causation, and unfolds and quantifies the multiscale causal structure of sys-

tems in ways previously impossible. In turn, this provides a quantitative measure of emergent complexity itself: how widely distributed a system's causal workings are across its scales, wherein systems that possess many contributing scales are more complex.

In what follows, CE 2.0 is outlined, first by defining a measure of causation that is axiomatic and robust to background assumptions, then by using that measure to calculate the degree of macroscale causation in coarse-grains of model Markov chains and so to quantify causal emergence, then by detailing how causal contributions are assigned across scales, and furthermore exploring how this leads naturally to a novel measure of emergent complexity. Finally, CE 2.0 is directly compared to other related theories of emergence, demonstrating its advantages and outlining its conceptual implications, and how the theory's practical limitations can be addressed by future research is suggested.

## 2 The axioms of causation

### 2.1 Sufficiency and necessity

Scientists distill and extract causal knowledge about the systems they study [31]. Breakthroughs in the scientific understanding of causation include things like RA Fisher's formalization of randomized controlled trials [32], as well as Judea Pearl's more recent introduction of the $do(x)$ operator [10].

A number of researchers have gone beyond this, introducing specific probabilistic measures of causality to capture the degree of causation between a cause and an effect. The aim of such measures can be described in various ways, e.g., as capturing the power of a particular cause, the strength of a particular causal relationship, the amount by which one variable causally controls another, etc. Applying such measures of causation involve specifying a causal model, then using counterfactuals [33] or interventions [10] to separate causal knowledge from mere observation.

Recent work analyzing over a dozen proposed measures of probabilistic causation [26] showed that in the scientific literature there is *causal consilience*: measures of causation, independently introduced across fields from psychology to statistics to philosophy [34], all set in relation two basic terms. We therefore dubbed these terms "causal primitives" [26]—more commonly, they are known as the *sufficiency* and the *necessity*. Consilience held true for measures of causation ranging from those proposed by philosopher David Lewis [33], to to mathematician Judea Pearl [10], to myself and co-author's recent definition of actual causation [35]. Rediscovered many times inde-

pendently, the primitives form an axiomatic foundation for any measure of causation, and ensure that measures of causation have significant overlap in their mathematical behavior. Ultimately, each term represents an inverse of uncertainty: sufficiency is the certainty about an effect, given a cause, whereas necessity is the certainty about a cause, given an effect.

As will be shown, the causal primitives of sufficiency and necessity have a further information-theoretic generalization, the determinism and degeneracy (respectively). CE 2.0 is based explicitly on these primitives and their further generalization.

First, to define the causal primitives formally, some background terminology is required. For a discrete system like a Markov chain, DAG, or set of logic gates, assessing the causal primitives (hereafter CP) involves specifying some space associated with the system, $\Omega$, which defines its set of states or events or variables entered into the causal analysis. Then, for any occurrence, we can define the potential causes, $C \in \Omega$, and the potential effects, $E \in \Omega$.

As the theory will be specified in simulated Markov chains, here $\Omega$ is just a system's statespace. Occurrences are then a state transition, in that for a Markov chain there is always some individual preceding state of the system, $c$, and its effect $e$, the next state. These transitions each have some probability, $P$.

Given an occurrence (here, a state transition), the sufficiency of a cause is then:

$$suff(e, c) = P(e \mid c)$$

which increases as $c$ is more probabilistically likely to bring about $e$, and reaches 1 when $c$ is fully sufficient for $e$.

The necessity of a cause is:

$$nec(e, c) = 1 - P(e \mid C, \neg c)$$

which specifies the probability of $e$ occurring without $c$. That is, given some set of causes $C$ within the system, and given that, within that set, $c$ itself did not occur, what is the probability of $e$? This probability is low when there are many common causes, and is 1 only when $c$ is absolutely necessary for $e$, in that no other members of $C$ could produce it (in a Markov chain, this would mean no other states lead to the state $e$).

Herein, these values will be calculated in Markov chains defined by some transition probability matrix (TPM). Since the goal is to analyze causation comparatively across scales, the terms "sufficiency" and "necessity" (and

their joint description as CP) will henceforth imply their system-wide average across all transitions.

In order to fully define these values, a final background condition must be specified, which is some distribution over the set of causes $C$. E.g., calculating the necessity (or the system-wide average sufficiency) requires defining:

$$P(e \mid C) = \sum_{c \in C} P(c)P(e \mid c)$$

This can be conceptualized as identifying a set of viable counterfactuals or, alternatively and identically, specifying some set of possible interventions over the states of the system [36].

Herein, the set of possible causes is the full set of states of the system. Conceptually, this just means that states of the system are treated equally when it comes to being viable interventions or counterfactuals to assess the causal relationships of other states. E.g., given a COPY gate with a self-loop of $p = 1$, this would imply we consider both 0 and 1 equally in $P(C)$, and therefore could correctly say that the state COPY $= 1$ (at some timestep) is sufficient and necessary for COPY $= 1$ in the next timestep (note that with merely an observed distribution such a judgment would be impossible).

## 2.2   Determinism and degeneracy

Two mathematical terms were introduced in CE 1.0: the *determinism* and the *degeneracy* of a system [8]. These are, it turns out, information-theoretic generalizations of the the sufficiency and necessity, but from the perspective of the probability distribution over sets of transitions.

Specifically, the determinism is the inverse of noise (or randomness) in the probability distributions of state transitions in a given system, and therefore is an information-theoretic generalization of the sufficiency. For a particular cause $c$ it can be defined as a coefficient, based on the entropy of $c$'s transition probability distribution over the set of effects $E$ in the system (normalized to provide a coefficient ranging from $[0, 1]$):

$$det(c) = 1 - \frac{H(E \mid c)}{\log_2 n}$$

wherein the central entropy term is:

$$H(E \mid c) = \sum_{e \in E} suff(e, c) \log_2 \frac{1}{suff(e, c)}$$

In what follows, the term "determinism" is reserved for the system-wide determinism coefficient, such that determinism is 1 only when the TPM is a permutation matrix, which is assessed via its average across all possible causes, given some $P(C)$:

$$det = \sum_{c \in C} P(c)\ det(c) = \sum_{e \in E,\ c \in C} P(e,c)\ det(e,c) = 1 - \frac{\sum_{c \in C} P(c)\ H(E \mid c)}{\log_2(n)}$$

"Degeneracy" is, just the same, based on the system-wide degeneracy coefficient, defined as the inverse of the entropy over the probability distribution of the full set of effects, $P(E)$, given the full set of causes, $P(C)$.

$$deg = 1 - \frac{H(E \mid C)}{\log_2(n)}$$

Degeneracy is high if causes have many similar effects. Degeneracy is zero when all causes have a unique effect. In Markov chains, a fully deterministic and non-degenerate system would be one in which every state transitions with $p = 1$ to some unique next state, with zero overlap. Degeneracy acts as an inclusive inverse of the necessity, in that

$$H(E \mid C) = \sum_{e \in E} P(e \mid C) \log_2 \frac{1}{P(e \mid C)}$$

wherein $P(e \mid C)$ is the inclusive form of the central term of the necessity calculation $P(e \mid C, \neg c)$. It calculates the necessity without the removal of $\neg c$ (essentially asking, for any given $e$, how necessary are the set of causes that lead to it). Degeneracy is 0 only when the TPM is a permutation matrix.

In order to avoid linguistic confusions around inverses (since a low degeneracy indicates a stronger causal relationship), here during calculations the degeneracy is often reversed into the *specificity*, for which increasing values indicate stronger causal relationships.

$$specificity = 1 - deg$$

To summarize: as can be seen by their construction from the same basic probability terms, determinism is essentially the normalized entropy of the sufficiences and degeneracy is essentially the normalized entropy of the necessities (as its inverse and calculated inclusively). In this way, determinism and degeneracy are information-theoretic generalizations of the sufficiency and necessity.

7

## 2.3 Causal primitives are sensitive to uncertainty

The causal primitives are sensitive to the uncertainty inherent to a causal model about causes and effects. Indeed, previous research on causal consilience has already shown that, due to their close mathematical relationship, the sufficiency and necessity over a set of two transitions behaved similarly to the determinism and degeneracy in conditions of increasing uncertainty [26].

Here it is further shown that, as implied by their shared terms, the system-wide determinism and degeneracy (represented by its inverse, specificity) behave similarly to the system-wide average sufficiency and necessity in larger systems with multiple possible causes under conditions of increasing uncertainty.

To show this, a system $S$ was modeled, defined as a set of 8 states with self-loops of $p = 1$ (see See Fig. 1, top left). In order to vary the uncertainty around causes and effects, noise (uncertainty about effects) and common causes (uncertainty about causes), were introduced along two separate axes. The first axis increased uncertainty about effects by shifting the system down to a condition of randomness in terms of its transitions (an all-to-all Markov chain wherein all transitions $= 1/n$, which means the system behaves as unpredictably as possible). The second axis, uncertainty about causes, moved the system to the condition wherein all transitions had an identical set of effects (thus increasing the number of common causes).

For every step along the axis that increased the noise, the probability from each self-loop was redistributed equally across the other states in the system, with the total amount of probability redistributed being $1/steps$ each step. For every step along the axis that increased the number of common causes, the full set of transitions for a state (a row in the TPM) were replaced one at a time with a duplication of the first row until all distributions were the same. Changing the model along just this latter axis began with all states having unique state-transitions and ended with all states transitioning to a single state (see Fig.1, bottom left). Finally, these two axes of changes to the system were combined such that at every step, both more noise in effects was introduced, and at the same step, more common causes were introduced (see Fig.1, middle diagonal).

The system-wide sufficiency plus the necessity was calculated for each state, as well as the determinism plus the specificity, along the increasing uncertainty in effects axis (Fig. 2A), along the increasing common causes axis (Fig. 2B), as well as along both axes combined (Fig. 2C).
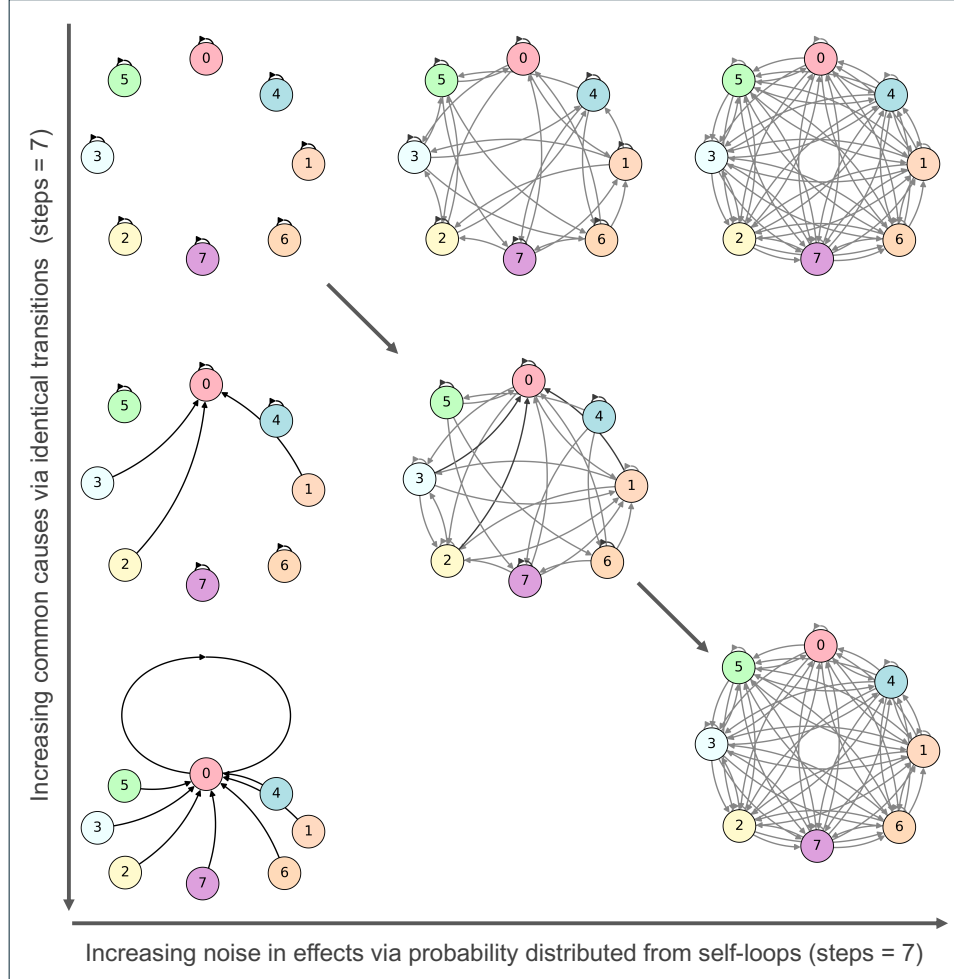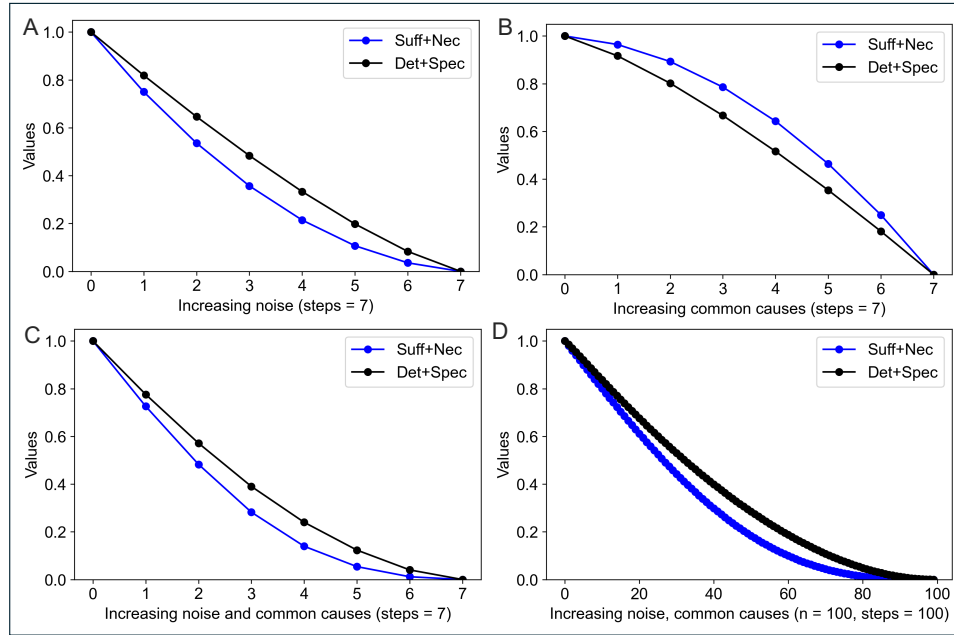
Figure 1: **Increasing uncertainty in the causal relationships of an 8-state system**. Starting in a state of self-loops with $p = 1$, states in the network were changed over a set number of steps equal to the size (number of states) of the system in three ways. Along the x-axis, self-loop probabilities were reduced by $1/steps$ and distributed equally to the other states (thus increasing the uncertainty of a particular effect, given a cause) until the system was an all-to-all network with random transitions. Along the y-axis, at each step a state was replaced with the transition distributions of another state (increasing the number of common causes and thus increasing the uncertainty of a cause, given an effect), until all states in the system shared the same transition. The system was also subjected to both changes at each step (the middle diagonal), ending again in an all-to-all state of random transitions.

9

At each step, the causal primitives (hereafter, CP) of the sufficiency and necessity pair and the determinism and degeneracy pair were calculated. This done in a way to ensure the same $[0,1]$ bounds (here and throughout). Specifically, both the sufficiency and the necessity were added together and 1 was subtracted. For the determinism and specificity (the inverse of degeneracy) pair, they were also added and 1 was subtracted; this latter case is equivalent to the *effectiveness* from CE 1.0, wherein the two values together were set in relation as $det - deg$.



Figure 2: **Causal primitives vary together with uncertainty.** (A) The sufficiency plus necessity value and the determinism plus specificity value are shown to behave similarly across increasing uncertainty over effects (noise) as the probabilities of self-loops in the system from Figure 1 are redistributed across the system in increasing steps. (B) Similar behavior in response to increases in common causes (overlap). (C) Similar behavior in response to increases in both noise and overlap. (D) The measures still behave similarly as the system becomes larger and more steps are added.

Notably, all CP values changed similarly along both axes of steps that increased uncertainty about causes and effects, even in larger systems (values for 100 states and 100 steps are also shown in Fig. 2D).

# 3   Quantifying macroscale causation in CE 2.0

CE 2.0 defines macroscale causation as when causal primitives (CP) are greater at a macroscale of a system, indicating that the macroscale has reduced the uncertainty about causes and effects. Note that not all dimension reductions always lead to gains in CP—some lead to zero gains or even decreases (cases of causal reduction). To identify positive gains, CE 2.0 makes use of an ordered micro→macro path across the hierarchy of scales of a system, revealing its multiscale structure, and then assesses each scale in terms of its unique CP gains (its causal contribution). The degree of causal emergence in a system is the total gain in CP along a micro→macro path, representing the sum of macroscale causation across all scales.

## 3.1   Formalizing macroscales

Macroscales are defined in CE 2.0 identically to previous research on causal emergence in systems of Markov chains (for details on how to derive a macroscale's TPM given some microscale TPM, see [8, 17]). Broadly, for a given system $S$, which represents the microscale, some new system $S_M$ is defined with macrostates replacing a set of microstates, wherein the transitions to and from a given macrostate are a summary statistic of the underlying microstates.

For simplicity, here I only consider dimension reductions that are coarse-grains of microstates. Coarse-grains are macroscales like $(0, 1, 2), (3)$, which would indicate that for some 4-state system the microstates $(0, 1, 2)$ have been coarse-grained into a macrostate, while $(3)$ remains as it was at the microscale. However, CE 2.0 can also be applied over other types of dimension reductions similar to coarse-graining, like black-boxing [9, 37] or higher-order macrostates [17].

Importantly, not all macroscales are sensible summaries of their underlying microscale; indeed, some macroscales may be dynamically inconsistent when defined [38]. Much as in Klein and Hoel [17], here I deem a macroscale valid if it is *consistent* with its underlying microscale, with consistency being defined as whether the path of random walkers on the Markov chain is the same at both the microscale and the macroscale (i.e., whether or not the macroscale acts as an accurate summary statistic for the microscale's dynamics).

Specifically, for a Markov chain the inconsistency can be defined as the Kullback-Leibler divergence [39], taken between an expected distribution of random walkers, across timesteps $t \to t_n$, in $S$ vs. $S_M$, given an identical

starting state on each scale. While previous research checked only the stationary distributions for such inconsistency [17], here I enforce a strict notion of consistency, wherein a random walker is dropped at every possible state, and inconsistency between the macroscale and the microscale is summed over all of their moves for the next 5 timesteps. Any non-zero values imply inconsistent macroscales, which are discarded. Therefore, all macroscales considered herein are fully consistent with the dynamics of their underlying microscales.

## 3.2 Defining a micro→macro path

A micro→macro path is the set of valid (i.e., consistent) scales that leads from the microscale up to a final macroscale, which acts as the endpoint of the path. Conceptually, a path is simply specifying what coarse-grains of the system are "on the way" to other coarse-grains across the hierarchy of scales that spans the system. Each step in the path is a single "slice" of the higher-dimensionality object, which the micro→macro path traverses from bottom (full dimensionality) to top (final dimension reduction).

In a hypothetical 4-state system, the coarse-grain of $(0, 1), (2), (3)$ would on a path to the lower-dimensionality coarse-grain of $(0, 1, 2), (3)$. That is, along such a path, the microstates $(0)$ and $(1)$ are first coarse-grained together into a single macrostate, $(0, 1)$, and then further coarse-grained into $(0, 1, 2)$. A micro→macro path for such a hypothetical 4-state system, starting at a microscale, might be $(0), (1), (2), (3) \rightarrow (0, 1), (2), (3) \rightarrow (0, 1, 2), (3)$ $\rightarrow (0, 1, 2, 3)$, ending with all states coarse-grained into one macrostate. Formally, this is just:

$$\pi^{(1)} \longrightarrow \pi^{(2)} \longrightarrow \cdots \longrightarrow \pi^{(k)},$$

wherein each $\pi^{(i)}$ is some valid partition (representing a coarse-grain) of the original $n$ microstates, and $\pi^{(i+1)}$ is a coarse-grain in turn of $\pi^{(i)}$, concluding at the endpoint, partition $\pi^{(k)}$.

As an example, a pre-chosen micro→macro path is plotted for the 8-state Markov chain visualized in Figure 3. Progress along the path is shown in Figure 3 via color contagion. Beginning at $(0), (1), (2), (3), (4), (5), (6), (7)$, the microscale, microstates coarse-grained together along the path are changed to be the same color at each coarse-graining step until the final end-point of $(0), (1, 2, 3, 4, 5, 6, 7)$, wherein all microstates have been coarse-grained together, except $(0)$.
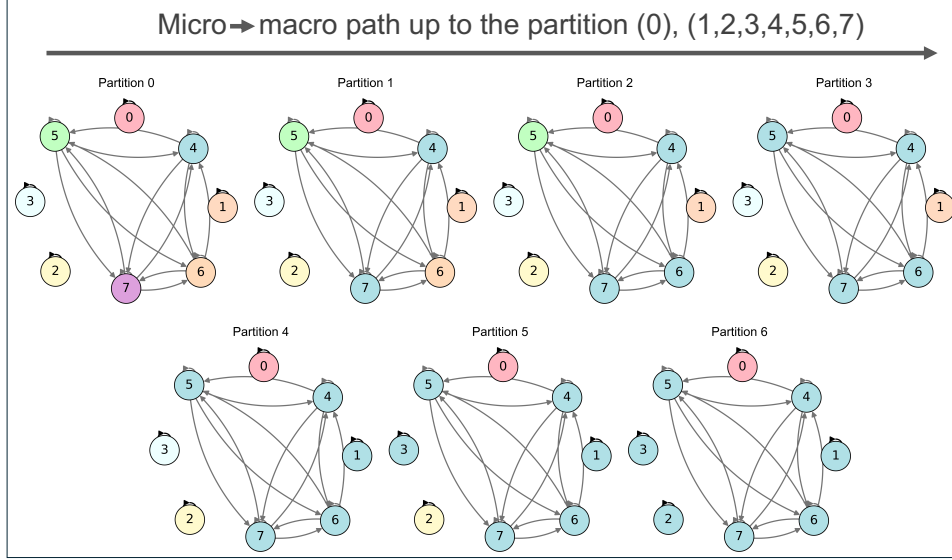
12

Figure 3: **A micro→macro path visualized**. An example 8-state Markov chain, with the probabilities of transitions represented in gray scale (for its TPM, see Fig.4A). Starting at the microscale, represented by the partition of $(0), (1), (2), (3), (4), (5), (6), (7)$, states are coarse-grained together, with each further partition being a step in the path (and thus a scale in the system), ending at $(0), (1, 2, 3, 4, 5, 6, 7)$. Changes to being the same color indicate when states get coarse-grained together along the path (color contagion).

## 3.3 Causal emergence as the total gain in CP along a path

For a system and a pre-chosen micro→macro path, CE 2.0 calculates the causal primitives along each step, tracking $\Delta$CP at each scale.

Calculating the $\Delta$CP along a path is demonstrated in an 8-state Markov chain, the starting microscale TPM of which is shown in Fig. 4A. It has an intuitive macrostate in the form of an equivalency class across microstates $(4, 5, 6, 7)$, which all share identical probability distributions of transitions. Its connectivity (state transitions) is the same as visualized in Figure 3, and the same pre-chosen micro→macro path is used.

Notably, in this system CP shows consistent gains until $(4, 5, 6, 7)$ are coarse-grained together into a single macrostate with a self-loop (the TPM of this macroscale is shown in Fig.4B, and is visualized in Fig.4C). But then the path immediately transitions to a domain of zero gains (plotted in Fig.4D).

Expressed in terms of determinism plus specificity, the microscale starts at CP $= 0.66$, and at the macroscale prior to the transition to zero gains this has increased by 0.33, attaining a maximum value of CP$= 1$, indicating that causal relationships become maximally deterministic and non-degenerate at that macroscale, and further coarse-graining entails no further gains. The degree of causal emergence in this system is therefore CE $= 0.33$, reflecting the total gains along the path.
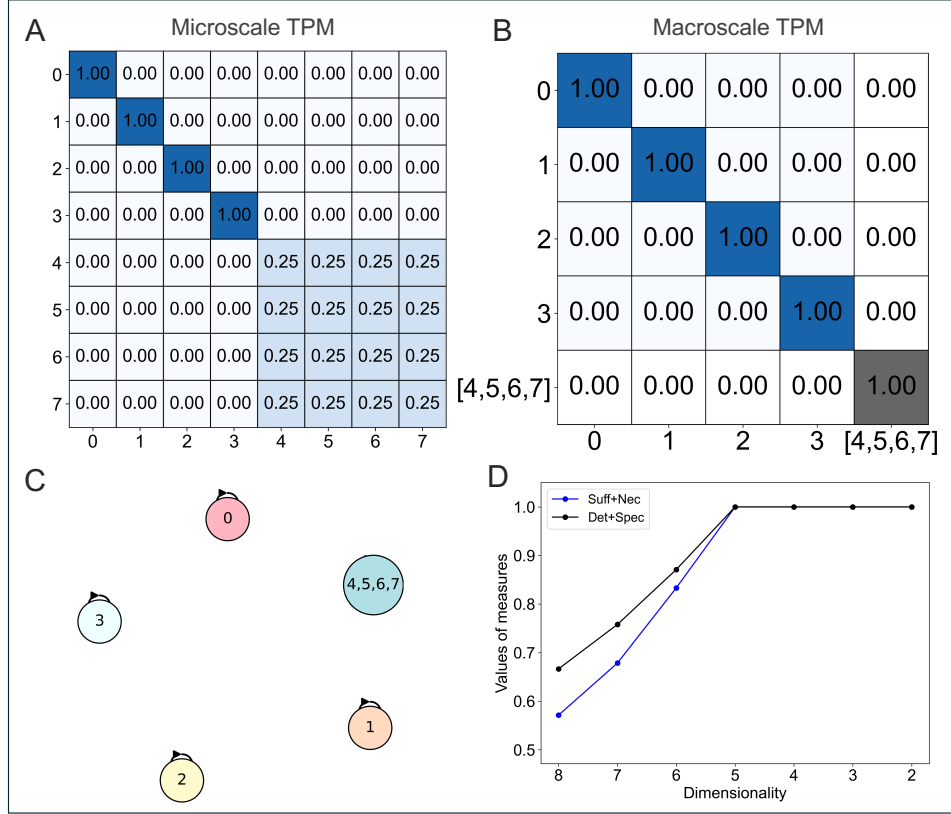


Figure 4: **Causal primitives along a micro→macro path**. (A) The TPM of the microscale, with cells colored based on their probability ($p = 1$ being a darker blue). (B) The TPM of the macroscale past which $\Delta$CP transitions abruptly to zero. (C) The same macroscale visualized as a network, with the macrostate labeled (its self-loop of $p = 1$ is not shown). (D) The change in causal primitives across the path of increasing dimension reduction, with the total gain in CP of 0.33 reflecting the degree of causal emergence.

14

## 3.4 Choosing micro→macro paths

When analyzing the causal emergence of a system, there may be prior reasons to pick a certain path; however, an appropriate micro→macro path can also be identified in a first-principles manner.

Specifically, the endpoint of the micro→macro path is the macroscale which entails the highest total gain in CP possible (the most causal emergence). In the case where there are multiple possible endpoints tied for the highest gain, the macroscale representing the lowest amount of dimensionality reduction is an optimal endpoint, as it indicates the point past which dimensionality reduction does not lead to gains in CP. Once the endpoint is identified, the maximally-informative micro→macro path to analyze $\Delta$CP along is the longest path across the set of consistent macroscales (with consistency defined as in Section 4.1) ranging from the microscale to the endpoint macroscale.

The macroscale which serves as the endpoint can be identified in a brute-force search by first generating all possible macroscales, discarding those which are inconsistent, calculating their CP, and then choosing the macroscale with the highest CP and highest dimensionality as the endpoint to the path.

In practice though, this is not feasible for larger systems, and heuristics are required. One is to coarse-grain along a path until diminishing returns are reached. To put this heuristic formally, there is a sequence of incremental gains $\Delta$CP$_i$ at each step $i$ in the path. The system enters "diminishing returns" at step $i^*$ if $\Delta$CP$_{i^*} < \varepsilon$ (for some small threshold $\varepsilon > 0$), or if the ratio $\Delta$CP$_{i+1}$/$\Delta$CP$_i$ consistently decreases over a long path length.

In other words, once $\Delta$CP becomes negligible (below $\varepsilon$) or keeps shrinking step-by-step for some long portion of the path, that indicates a transition point to diminishing returns and an approximate endpoint that does not leave out substantial gains in CP. However, care must be taken to not simply stop at local maxima (see Section 5) by choosing a small enough $\varepsilon$ or a long enough path length to assess diminishing returns.

For a bound of the amount of causal emergence in a system prior to defining a path or even any macroscales, one can simply measure CP at the microscale. Its distance from 1 provides an upper bound for CE without the need to search across scales, allowing for quick estimations.

Indeed, if the CP of the microscale is significantly less than 1, it is likely that there is some dimension reduction (like a coarse-grain) that increases it to at or near maximum; therefore, the difference between the CP of the microscale and 1 may approximate the CE value for many systems (but

does not specify which macroscale the gain comes from). It is even possible that for large enough systems with enough structure there is always some macroscale where CP = 1 when taking into account the full set of dimension reductions, like Higher Order Macrostates [17], or when relaxing consistency assumptions.

# 4   Emergent complexity

Traditionally, the field of complexity science has been motivated by the qualitative notion that complexity "emerges" from system dynamics, either in that system dynamics become more complex over time, or in that there is intuitive macroscale or mesoscale structure to the system beyond the microscale. Specific quantitative methods in complexity science for detecting this latter case (that of macroscale or mesoscale structure) have focused mainly on compressibility or efficiency [40], or more recently, by assessing the info-theoretic surprise [41]. However, this means that the mesoscales detected may just be convenient compressions and have no causal relevancy.

In comparison, by analyzing the distribution of causal contributions, CE 2.0 can be used to quantify the emergent complexity that's actually causally-relevant for the workings of the system. CE 2.0 therefore provides a taxonomy for how complex the causal working's of a system are: if they are mostly confined to a single scale (like either the microscale, or in that it is dominated by a "top-heavy" macroscale) then the system is simple, whereas if the system has intermediate mesoscales that also have substantial causal contributions, it is complex. This distinguishes cleanly between when a system is emergent (here, a high CE value) and when that emergence evinces significant complexity.

Specifically, the $\Delta$CP for each step along a path represents that scale's causal contribution to the total CP (which fully determines the system's causal workings)—therefore, its distribution along a path can be assessed.

To demonstrate this novel aspect of CE 2.0, in Figure 5 two causally emergent systems are analyzed: one that has no mesoscale structure and one that does, but is otherwise as similar as possible.

Fig.5A shows the TPM of a system composed of an intuitive macroscale with no distinguishable mesoscale structure, wherein $(0, 1, 2, 3)$ and $(4, 5, 6, 7)$ have been coarse-grained into two respective macrostates. Indeed, the microstates coarse-grained into the two macrostates each make up an equivalency class (see Fig.5B for a visualization).

At each scale along the micro→macro path, $\Delta$CP is tracked. Notably,

along the micro→macro path gains in CP are clustered at the endpoint of the path (see Fig.5C). This indicates a "top-heavy" structure mostly composed of the causal contributions of the microscale (contributing 0.14 of the total CP) and the endpoint of a large macroscale over the two equivalency classes (contributing 0.18). Together the microscale and macroscale account for the majority of the full CP value (.41) for the system, quantifying that its causal workings are dominated by those two scales.

To formally capture the taxonomy of "top-heavy" or "bottom-heavy" systems on one side, and systems with substantial mesoscale structure on the other side, here I introduce a notion of *emergent complexity* (EC). It is based on the entropy of the causal contributions along a path of length $L$. Given a set of gains $\Delta\mathrm{CP}_i$ at each step $i = 1, 2, \ldots, L$, then

$$p_i = \frac{\Delta\mathrm{CP}_i}{\sum_{j=1}^{L} \Delta\mathrm{CP}_j} \quad \text{for } i = 1, \ldots, L$$

which ensures that $\{p_1, \ldots, p_L\}$ is a probability distribution over the $L$ steps. To measure how "spread out" (multiscale) the gains are, the entropy of these relative gains is calculated

$$\mathrm{EC} = \log_2(L) - \sum_{i=1}^{L} p_i \log_2(p_i)$$

which equals $\log_2(L)$ if all $\Delta\mathrm{CP}_i$ are equal (i.e., $p$ is uniform), and decreases when a small number of steps in the path dominate the total gains.

To demonstrate how this detects mesoscale structure, a similar system is modeled in Fig.5D-E, but with the change that (0) and (4) are distinguishable (in terms of their state transitions) from the other of members they are coarse-grained together with at their endpoint macrostate, $(0, 1, 2, 3)$ and $(4, 5, 6, 7)$, respectively. That is, for this system the ultimate largest dimension reduction is not over a pure equivalency class.

In this mesoscale system (with a total CE of 0.13 at the endpoint where $(0, 1, 2, 3)$ and $(4, 5, 6, 7)$ are macrostates), there is a visible earlier maxima when the causal contributions along the path are plotted (Fig.5F). The mesoscale is revealed by how $\Delta\mathrm{CP}$ peaks at a higher dimensionality than the endpoint.

Because of this earlier maximum in terms of $\Delta\mathrm{CP}$, the emergent complexity of the mesoscale system is 2.6 bits, whereas the emergent complexity of the "top-heavy" system is only 1.9 bits (to compare systems of different sizes, this value itself can be normalized).
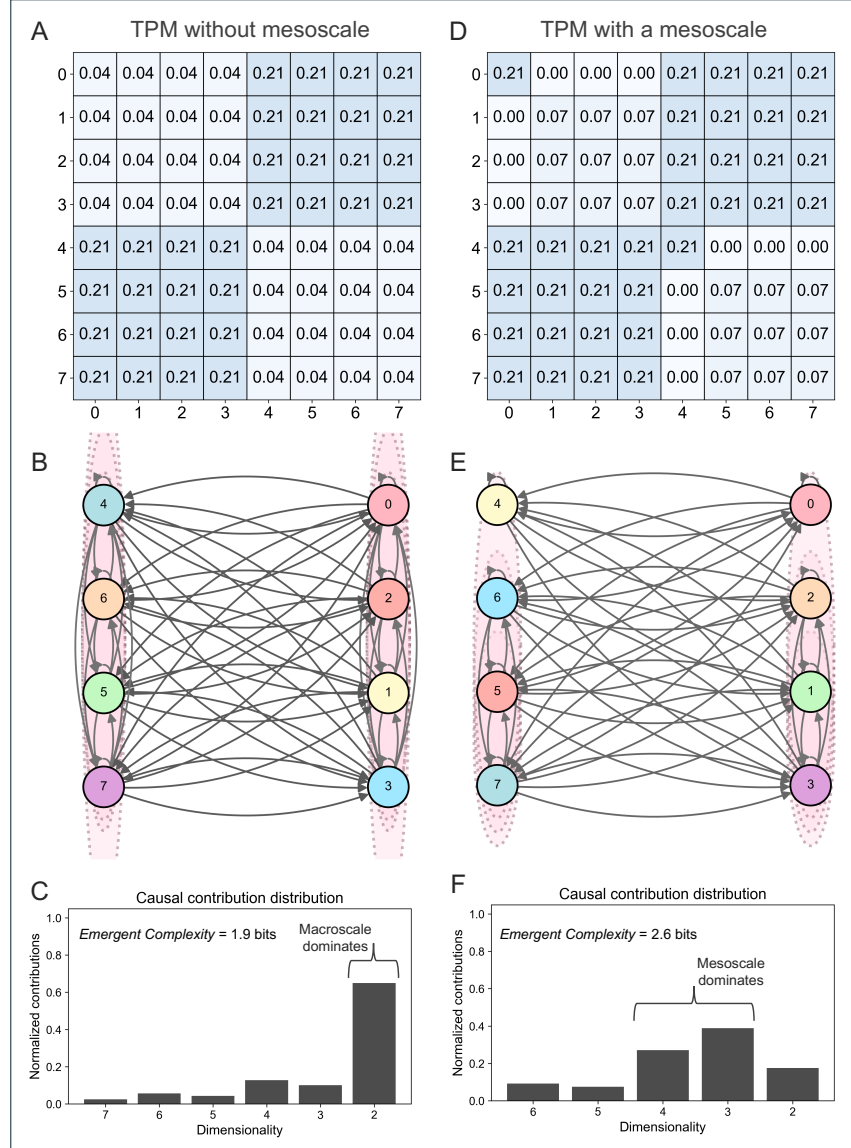
Figure 5: **Causal contributions across scales**. (A) Microscale TPM of a system with no mesoscale structure. (B) The same system visualized. (C) CE 2.0 identifies the causal contributions as "top-heavy," in that the last dimension reduction contributes the most. (D) Microscale TPM of an otherwise similar system with mesoscale structure. (E) The mesoscale system visualized. (F) Causal contributions are shifted toward lesser dimension reductions, indicating a predominately multiscale causal structure; it therefore possesses more emergent complexity.

# 5 Comparison to other theories of emergence

CE 2.0 has many advantages as a theory: (a) it is axiomatically grounded in causal primitives, (b) it captures all possible cases of macroscale causation, which even CE 1.0 does not, and (c) it elucidates the multiscale causal structure of systems in a novel manner, resolving longstanding conflicts around over-determination and causal exclusion. Here, these advantages are considered and demonstrated.

## 5.1 An axiomatic grounding

As discussed, the EI's use in CE 1.0 has been previously criticized on the basis of being drawn from a maximum entropy (uniform) distribution [24, 25]. In some cases, like in Integrated Information Theory, the maximum entropy distribution can be justified as a function of taking the "intrinsic perspective" on a system [20]. However, this is reliant on accepting the postulates of IIT, including its analysis of consciousness.

While there are potential replies to this criticism, such as defending the EI as a well-tuned measure of causation [23], a theory of causal emergence should not be reliant on a single particular measure of causation and its assumptions: this would impact both its practical uses and also its theoretical foundations.

Comparatively, for CE 2.0, it is already known that gains in CP at the macroscale are robust to choice of intervention distribution or counterfactual choice (in the form of specifying a particular $P(C)$), to the degree of gains even occurring when using the observed distribution [26].

Therefore, CE 2.0 is more theoretically robust than CE 1.0 was, by virtue of being grounded in causal primitives that historically have shown to be fundamental to the nature of causation and more robust to assumptions in application.

## 5.2 CE 2.0 captures all macroscale causation

CE 2.0 can detect cases of macroscale causation that the CE 1.0 framework does not.

Examples of this detection are shown in an 8-state system is constructed of two equivalency classes, making it a "block model" at the microscale (see Fig. 6A, left). A single macroscale is specified, wherein the two equivalency classes are each coarse-grained into a respective macrostate with a self-loop, represented by the coarse-grain $(0, 1, 2, 3), (4, 5, 6, 7)$. Causal emergence as

calculated via CE 1.0 (based on the gain in EI at the macroscale) and the causal emergence as calculated via CE 2.0 (based on the gain in CP at the macroscale) is shown across a manipulation of that system.

Starting with the initial TPM shown in Fig.6A (left), the probabilities *within* each equivalency class for each state $s_i$ are manipulated in steps such that, over 50 steps, the probabilities that were previously transitioning to the rest of the equivalency class are added incrementally to the self-loop probability of $s_i$, eventually reaching $p = 1$. The midpoint of this manipulation is shown in Fig.6A (middle), and the final ending system is the TPM in Fig.6A (right). This progresses the system via discrete steps of probability redistribution from an initial "block model" configuration to a permutation matrix in the form of a set of 8 microstates with self-loops of $p = 1$.

During the manipulation of the system described above, comparative causal emergence values are calculated off a chosen macroscale. Specifically, the macroscale $(0, 1, 2, 3), (4, 5, 6, 7)$, is used as the coarse-grain for all comparisons across all system manipulations. At each step calculations from CE 1.0 are shown, compared to when causal emergence is calculated by the gain in CP as in CE 2.0 (with an endpoint of the fixed chosen macroscale).

Counterintuitively, the EI from CE 1.0 detects no causal emergence, even when the system is initially split into two equivalency classes. Meanwhile, when viewed from the new lens of CE 2.0, the system in Figure 6 sensibly starts with a significant degree of macroscale causation. This degree then decreases in accordance with the increasing self-loop probabilities and the increasing distinguishability of the microscale, becoming weaker and weaker as the macroscale contributes marginally less, until it vanishes altogether when the microscale becomes perfectly deterministic and non-degenerate.

However, in some cases, CE 1.0 and CE 2.0 will overlap (e.g., the scale right before $\Delta$CP transitions to being zero in Figure 4 is the same as would be picked out by searching for the maximum of EI). This is because CE 1.0 and CE 2.0 share a close mathematical connection in their terms, since the EI has a decomposition wherein:

$$EI = eff * log_2(n)$$

That is, the EI can be decomposed into the determinism minus the degeneracy (the *effectiveness*), which is then multiplied by a *size* term, $log_2(n)$, which in turn is just the dimensionality of the given scale [8]. However, in CE 2.0, the *size* term is rendered unnecessary by appropriate multiscale structure analysis and causal apportioning.
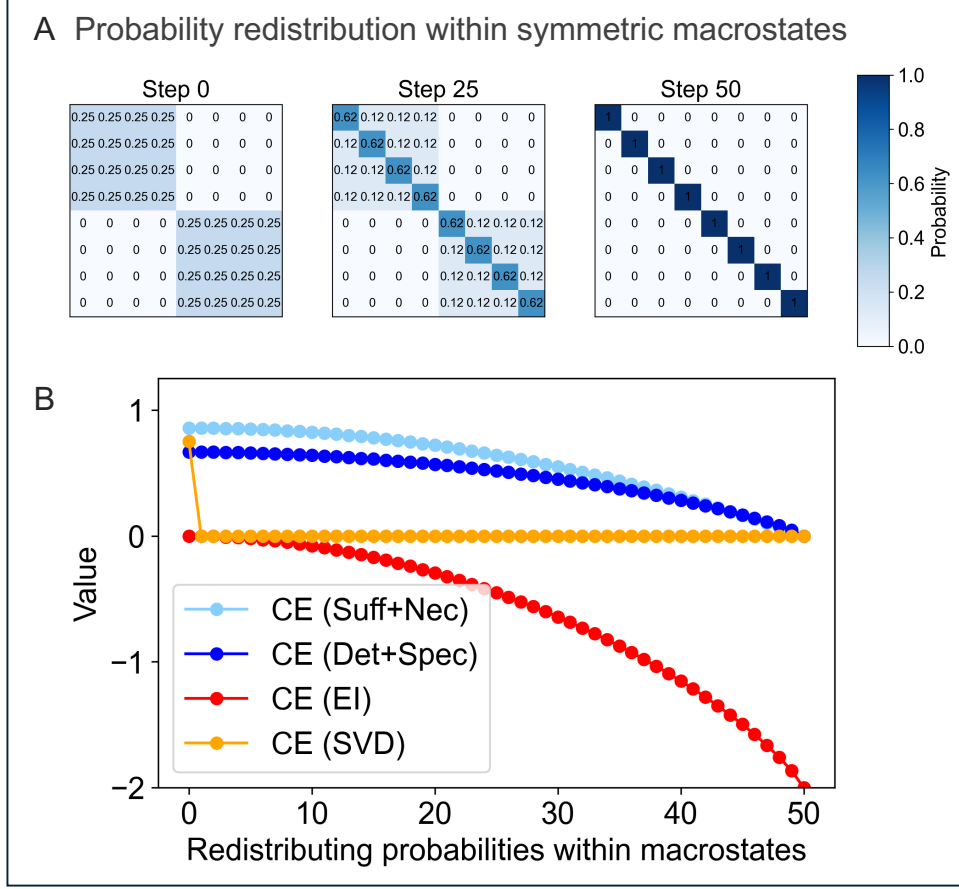
Figure 6: **CE 1.0 is sensitive to symmetries and mathematically unstable.** (A) TPMs (probabilities shown in bluescale) of a "block model" system with two macrostates over its equivalency classes at the beginning, midpoint, and end of increasing the self-loop probabilities of each state. Redistribution is performed by drawing away probability from its full set of transitions via increments of $1/steps$ until the microscale is entirely composed of states with self-loops of $p = 1$. (B) CE 2.0 detects the macroscale causation and decreases sensibly as the microscale becomes more causally distinguishable during the probability redistribution. However, causal emergence as calculated in the CE 1.0 framework (based on the difference in EI) is insensitive to the macroscale causation. Additionally, when applying the CE 1.0 framework to the singular value decomposition (SVD) method for measuring causal emergence [21] (discussed in Section 6), it is unstable when detecting macroscale causation, immediately collapsing to zero after any amount of probability redistribution of this type.

While CE 2.0 focuses on $\Delta$CP, if a single causally-relevant scale with high dimensionality is desired for causal modeling or explanation, its methods can be used to pick one out: the individual macroscale that maximizes CP via the smallest amount of dimensionality reduction. This can be done either by analyzing for diminish returns in $\Delta$CP, as described in Section 3.4, or even via explicit re-introduction of the *size* term against which to weigh gains in CP, flexibly re-capturing the analysis of CE 1.0.

## 5.3    Comparison to other related theories of emergence

Alternative proposals have been suggested for different ways to detect causal emergence, or emergence more generally, such as via integrated information decomposition [42], or via examining dynamical dependency [43]. However, both these theories make use of the mutual information. This is problematic for causal emergence, as it is definitional of causation that it is not dependent on simply the data distribution of the process being measured [36]. For example, in a cycle of logic gates performing the COPY function, the mutual information scales entirely off of how varied the initial state is, rather than capturing the fact that every gate is sufficient and necessary for the next step in the chain in the way the causal primitives do [21, 26].

Other recent work on emergence has focused on identifying cases wherein macroscales are consistent with their underlying microscales, but still independently describable in their dynamics, like the software of a computer [28]. Such efforts have examined whether or not a macroscale is "causally closed" in that it can can be thought of as being its own cause. This is closely related to the condition of macroscale consistency based on random walkers laid out here and (to a less strict degree) previously as well [17]. However, merely checking for consistency, lumpability, causal closure, etc, does not directly measure causal emergence, as it does not capture what a macroscale contributes to a system's causal workings above and beyond the microscale, which requires some specified measure of causation. Rather, it merely identifies which macroscales are valid descriptions of their microscale that preserve its dynamics and therefore are appropriate compressions. E.g. for the system analyzed in Figure 4, the largest dimension reduction at the end the pre-chosen path is a valid macroscale, completely consistent with its microscale, and yet has trivially zero causal contribution.

## 5.4    Conceptual implications of CE 2.0

All of science, outside of microphysics, implicitly operates as if there is causal emergence, in that it takes for granted the macroscale entities in its models and explanations and experiments are efficacious regarding a system's causal workings. This is contradicted by a nominal commitment to universal reductionism, which would seem to imply that all causal powers "drain away" to the bottom microscale of any system [44, 45]. This is due to the causal exclusion argument [46]—for any given supervening macroscale, its effect could also be described as a cause of its underlying microscale, which then renders the macroscale description unnecessary.

CE 1.0 flipped the exclusion argument on its head by noting that, according to the EI, the macroscale had greater causal power. A similar sort of thinking underlies the exclusion postulate in Integrated Information Theory, which is perhaps the most controversial of the postulates [47]. But it meant that in CE 1.0 there was the counterintuitive result that, even when macrostates were not over exact equivalency classes, the underlying microscale could itself be excluded; a surprising result with unclear epistemological and ontological implications.

Comparatively, in CE 2.0 causal exclusion is handled more gracefully. When viewed from the current analysis of a single path, macroscales do not over-ride the causation of the microscale (although the microscale's causal contribution can still be vanishingly small). Instead, they simply contribute additional causal power via the causal apportioning schema, leading to a more comprehensive mereology wherein individual scales are a lossy slice of a higher-dimensional object that contains all the relevant information about the system's causal workings. The source of this extra causal power beyond the microscale is non-mysterious, being based in uncertainty reduction [9, 48], which in turn stems from the multiple-realizability of macrostates [2].

Since in CE 2.0 emergence occurs via the minimization of uncertainty at macroscales, there is the broader question of whether uncertainty (in the form of noise or common causes) exists merely epistemically when it comes to the causal models of science and the systems they represent. Yet answering this involves speculation about unknowns like the scientific end-state of physics [2]. It is also worth noting even small sources of true uncertainty (like indeterminism) can be amplified via chaotic systems, and there may even be provably undecidable systems [49]. Even if all uncertainty inherent to the causal relationships of scientific causal models were only in principle epistemic, this would hold only for closed causal models that span the entire universe, and in such a universe-size causal model all notion of causation dis-

appears entirely anyway, as there are no definable interventions from outside the model [10]. Overall, while causal emergence vanishes in the condition of a microscale possessing zero uncertainty about the effect of causes, and also possessing no common causes, such conditions entail a far departure from most causal models in science.

# 6    Limitations and future research directions

A practical limitation for CE 2.0 is that searching the set of macroscales of a system entails a combinatorial explosion. While heuristics exist for CE 1.0 that could be adapted for CE 2.0 [15, 16] (including in continuous systems [50]), a further limitation of CE 2.0 is that, in its current formulation, it assumes a particular chosen micro→macro path to define a hierarchy of scales. This initial formulation makes practical and conceptual sense; especially since a single path might often be desired for analysis. However, work remains in developing a causal apportioning schema that can integrate the full set of non-commensurate micro→macro paths, which ideally should also address the issue of combinatorial explosion.

Recent research points a way forward. Work by Zhang et al. [21] introduced an alternative way to measure causal emergence in the CE 1.0 framework via the singular value decomposition (SVD) of a Markov chain. They define "vague" causal emergence as when the average of the resultant singular values ($\sigma$s), which are chosen based on being higher than some unspecified $\epsilon$, surpasses the total average, which they call $\gamma$ (and prove can approximate the determinism plus specificity used herein; for details, see [21]). Using a low threshold $\epsilon$ to include most non-zero $\sigma$ values then approximates the maximal increase available in EI at some possible macroscale, and so the SVD method is able to recapture the analysis of CE 1.0 without combinatorial explosions.

However, this method runs into the same limitations when applied through the lens the CE 1.0 framework, i.e., when used to approximate the macroscale with maximal EI (via a low $\epsilon$ to capture most non-zero $\sigma$s). This can be seen via the same system manipulation shown in Fig. 6A. During probability redistribution of this type within a macroscale, the measure is mathematically unstable, reducing to zero at the slightest amount of redistribution (see Fig. 6B for a plotted value, labeled "CE(SVD)"). This indicates that using dynamical reversibility based on the CE 1.0 framework also underestimates macroscale causation.

Yet this suggests a clear direction for future research wherein the SVD

method is adapted for the CE 2.0 framework (in which case, "vague" causal emergence would become a specific value). In such an adaptation, the set of singular values could be used to represent essentially *directionalities* of coarse-graining, which could capture the local and global maxima of $\Delta$CP beyond an individual path. The causal contributions of each directionality could then be calculated, not in a vague manner, but in a specific one via an adaptation of the causal apportioning schemes detailed in Section 4, assessing the differences of each $\sigma$ to $\gamma$ to locate positive contributions.

Future research will explore whether this possible connection allows CE 2.0's multiscale analysis to be applied to large systems without any combinatorial explosions.

# 7 Conclusion

CE 2.0 provides a conceptually and mathematically novel theory of emergence that treats systems as a hierarchy of scales. Individual scales, even in most cases the microscale, are simply slices of a higher-dimensional object— but only a slim minority of these scales are causally relevant, and the theory can identify them, revealing the hierarchy that matters to a system's causal workings. Specifically, macroscale causation is measurable via the gain in information-theoretic generalizations of the causal primitives (sufficiency and necessity) along a specified path that traverses the possible dimension reductions of a system. Causal emergence is the sum of this gain, and causal contributions can be apportioned out along the path. A taxonomy is revealed where some systems are "top-heavy" in terms of their causation, wherein a macroscale dominates, whereas others have more mesoscale structure (assessed by identifying the distribution of gains in the causal primitives along a traversing path). As a result, the theory allows for a novel definition of emergent complexity based on how spread out across the hierarchy of scales causal contributions are.

While the theory's initial formulation, based on a path that traverses dimension reductions, faces a limitation in the form of combinatorial explosions when applied to larger causal models, there are clear future research directions to address this.

CE 2.0 has critical applications in fields like physics, biology, neuroscience, and economics. Here, it's also worth highlighting one new particular use case: Given the axiomatic importance of causal primitives for understanding causation in complex systems, and given previous research showing the EI is responsive to changes following learning in artificial neu-

ral networks [51], CE 2.0 may also be primed to contribute to the growing field of AI interpretability [52] and AI safety [53], such as analyzing the multiscale structure of deep neural networks.

# 8    Acknowledgments

# References

[1] Geoffrey West. *Scale: The universal laws of life, growth, and death in organisms, cities, and companies.* Penguin, 2018.

[2] Erik Hoel. *The world behind the world: Consciousness, Free Will, and the Limits of Science.* Simon and Schuster, 2024.

[3] Matteo Grasso, Larissa Albantakis, Jonathan P Lang, and Giulio Tononi. Causal reductionism and causal structures. *Nature neuroscience*, 24(10):1348–1355, 2021.

[4] Erik P Hoel. Agent above, atom below: how agents causally emerge from their underlying microphysics. *Wandering towards a goal: how can mindless mathematical laws give rise to aims and intention?*, pages 63–76, 2018.

[5] Jessica C Flack. Coarse-graining as a downward causation mechanism. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 375(2109):20160338, 2017.

[6] Hans Albert Braun. Stochasticity versus determinacy in neurobiology: From ion channels to the question of the "free will". *Frontiers in Systems Neuroscience*, 15:629436, 2021.

[7] Peter Jedlicka. Revisiting the quantum brain hypothesis: toward quantum (neuro) biology? *Frontiers in molecular neuroscience*, 10:366, 2017.

[8] Erik P Hoel, Larissa Albantakis, and Giulio Tononi. Quantifying causal emergence shows that macro can beat micro. *Proceedings of the National Academy of Sciences*, 110(49):19790–19795, 2013.

[9] Erik P Hoel. When the map is better than the territory. *Entropy*, 19(5):188, 2017.

[10] Judea Pearl. *Causality*. Cambridge university press, 2009.

[11] Thomas F Varley. Causal emergence in discrete and continuous dynamical systems. *arXiv preprint arXiv:2003.13075*, 2020.

[12] Mingzhe Yang, Zhipeng Wang, Kaiwei Liu, Yingqi Rong, Bing Yuan, and Jiang Zhang. Finding emergence in data by maximizing effective information. *National Science Review*, 12(1):nwae279, 2025.

[13] Federico Pigozzi, Adam Goldstein, and Michael Levin. Associative conditioning in gene regulatory network models increases integrative causal emergence.

[14] Erik Hoel and Michael Levin. Emergence of informative higher scales in biological systems: a computational toolkit for optimal prediction and control. *Communicative & Integrative Biology*, 13(1):108–118, 2020.

[15] Ross Griebenow, Brennan Klein, and Erik Hoel. Finding the right scale of a network: efficient identification of causal emergence through spectral clustering. *arXiv preprint arXiv:1908.07565*, 2019.

[16] Jiang Zhang and Kaiwei Liu. Neural information squeezer for causal emergence. *Entropy*, 25(1):26, 2022.

[17] Brennan Klein and Erik Hoel. The emergence of informative higher scales in complex networks. *Complexity*, 2020(1):8932526, 2020.

[18] Brennan Klein, Erik Hoel, Anshuman Swain, Ross Griebenow, and Michael Levin. Evolution and emergence: higher order information structure in protein interactomes across the tree of life. *Integrative Biology*, 13(12):283–294, 2021.

[19] Erik P Hoel, Larissa Albantakis, William Marshall, and Giulio Tononi. Can the macro beat the micro? integrated information across spatiotemporal scales. *Neuroscience of Consciousness*, 2016(1):niw012, 2016.

[20] William Marshall, Graham Findlay, Larissa Albantakis, and Giulio Tononi. From micro to macro units: a mathematical framework for identifying the causal grain of a system from its intrinsic perspective. *bioRxiv*, pages 2024–04, 2024.

[21] Jiang Zhang, Ruyi Tao, Keng Hou Leong, Mingzhe Yang, and Bing Yuan. Dynamical reversibility and a new theory of causal emergence based on svd. *npj Complexity*, 2(1):3, 2025.

[22] Bing Yuan, Jiang Zhang, Aobo Lyu, Jiayun Wu, Zhipeng Wang, Mingzhe Yang, Kaiwei Liu, Muyun Mou, and Peng Cui. Emergence and causality in complex systems: A survey of causal emergence and related quantitative studies. *Entropy*, 26(2):108, 2024.

[23] David Balduzzi. Information, learning and falsification. *arXiv preprint arXiv:1110.3592*, 2011.

[24] Frederick Eberhardt and Lin Lin Lee. Causal emergence: When distortions in a map obscure the territory. *Philosophies*, 7(2):30, 2022.

[25] Scott Aaronson. Higher-level causation exists (but I wish it didn't), jun 2017.

[26] Renzo Comolatti and Erik Hoel. Causal emergence is widespread across measures of causation. *arXiv preprint arXiv:2202.01854*, 2022.

[27] Rafael Yuste. From the neuron doctrine to neural networks. *Nature reviews neuroscience*, 16(8):487–497, 2015.

[28] Fernando E Rosas, Bernhard C Geiger, Andrea I Luppi, Anil K Seth, Daniel Polani, Michael Gastpar, and Pedro AM Mediano. Software in the natural world: A computational approach to emergence in complex multi-level systems. *arXiv preprint arXiv:2402.09090*, 2024.

[29] Andrés Gómez-Emilsson and Chris Percy. The "slicing problem" for computational theories of consciousness. *Open Philosophy*, 5(1):718–736, 2022.

[30] Joshua Bongard and Michael Levin. There's plenty of room right here: Biological systems as evolved, overloaded, multi-scale machines. *Biomimetics*, 8(1):110, 2023.

[31] Helen Beebee, Christopher Hitchcock, Peter Charles Menzies, and Peter Menzies. *The Oxford handbook of causation*. Oxford Handbooks Online, 2009.

[32] Lloyd D Fisher. Advances in clinical trials in the twentieth century. *Annual Review of Public Health*, 20(1):109–124, 1999.

[33] David Lewis. Causation. *The journal of philosophy*, 70(17):556–567, 1973.

[34] Branden Fitelson and Christopher Hitchcock. *Probabilistic measures of causal strength*.

[35] Larissa Albantakis, William Marshall, Erik Hoel, and Giulio Tononi. What caused what? a quantitative account of actual causation using dynamical causal networks. *Entropy*, 21(5):459, 2019.

[36] Judea Pearl and Dana Mackenzie. *The book of why: the new science of cause and effect*. Basic books, 2018.

[37] William Marshall, Larissa Albantakis, and Giulio Tononi. Black-boxing and cause-effect power. *PLoS computational biology*, 14(4):e1006114, 2018.

[38] Paul K Rubenstein, Sebastian Weichwald, Stephan Bongers, Joris M Mooij, Dominik Janzing, Moritz Grosse-Wentrup, and Bernhard Schölkopf. Causal consistency of structural equation models. *arXiv preprint arXiv:1707.00819*, 2017.

[39] Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.

[40] Carlos Gershenson and Nelson Fernández. Complexity and information: Measuring emergence, self-organization, and homeostasis at multiple scales. *Complexity*, 18(2):29–44, 2012.

[41] Emiliano Marchese, Guido Caldarelli, and Tiziano Squartini. Detecting mesoscale structures by surprise. *Communications Physics*, 5(1):132, 2022.

[42] Fernando E Rosas, Pedro AM Mediano, Henrik J Jensen, Anil K Seth, Adam B Barrett, Robin L Carhart-Harris, and Daniel Bor. Reconciling emergences: An information-theoretic approach to identify causal emergence in multivariate data. *PLoS computational biology*, 16(12):e1008289, 2020.

[43] Lionel Barnett and Anil K Seth. Dynamical independence: discovering emergent macroscopic processes in complex dynamical systems. *Physical Review E*, 108(1):014304, 2023.

[44] Thomas D Bontly. The supervenience argument generalizes. *Philosophical Studies*, 109:75–96, 2002.

[45] Ned Block. Do causal powers drain away? *Philosophy and Phenomenological Research*, 67(1):133–150, 2003.

[46] Jaegwon Kim. *Mind in a physical world: An essay on the mind-body problem and mental causation*. MIT press, 2000.

[47] Tim Bayne. On the axiomatic foundations of the integrated information theory of consciousness. *Neuroscience of consciousness*, 2018(1):niy007, 2018.

[48] Kaiwei Liu, Bing Yuan, and Jiang Zhang. An exact theory of causal emergence for linear stochastic iteration systems. *arXiv preprint arXiv:2405.09207*, 2024.

[49] Robert Cardona, Eva Miranda, Daniel Peralta-Salas, and Francisco Presas. Constructing turing complete euler flows in dimension 3. *Proceedings of the National Academy of Sciences*, 118(19):e2026818118, 2021.

[50] Pavel Chvykov and Erik Hoel. Causal geometry. *Entropy*, 23(1):24, 2020.

[51] Scythia Marrow, Eric J Michaud, and Erik Hoel. Examining the causal structures of deep neural networks using information theory. *Entropy*, 22(12):1429, 2020.

[52] Adly Templeton. *Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet*. Anthropic, 2024.

[53] Seth Lazar and Alondra Nelson. Ai safety on whose terms?, 2023.