

S1 Appendix

1 Information decomposition in large multivariate systems

To formulate our theory of causal emergence for arbitrary order k , in Section *Partial information decomposition* we introduced definitions of k^{th} -order synergy and unique information. In this appendix we complement these with a matching definition of k^{th} -order redundancy, and show that these provide a full-fledged information decomposition for any $k = \{1, \dots, n-1\}$. For completeness, we present all definitions and examples here – including those that were necessary for the exposition of the main text and were previously presented in Section *Partial information decomposition*.

We begin by (re-)introducing the notion of k^{th} -order synergy between n variables, defined as

$$\text{Syn}^{(k)}(\mathbf{X}; Y) := \sum_{\alpha \in \mathcal{S}^{(k)}} I_{\partial}^{\alpha}(\mathbf{X}; Y) ,$$

with $\mathcal{S}^{(k)} = \{\{\alpha_1, \dots, \alpha_L\} \in \mathcal{A} : |\alpha_j| > k, \forall j = 1, \dots, L\}$. Intuitively, $\text{Syn}^{(k)}(\mathbf{X}; Y)$ corresponds to the information about the target that is provided by the whole \mathbf{X} but is not contained in any set of k or less parts when considered separately from the rest. Accordingly, $\mathcal{S}^{(k)}$ only contains collections of more than k sources. For example, for $n = 2$ we obtain the standard synergy $\mathcal{S}^{(1)} = \{\{12\}\}$, and for $n = 3$ we have $\mathcal{S}^{(1)} = \{\{12\}, \{13\}, \{23\}, \{12\}\{13\}, \{12\}\{23\}, \{13\}\{23\}, \{12\}\{13\}\{23\}, \{123\}\}$.

Similarly, the k^{th} -order unique information of \mathbf{X}^{β} with $\beta \subset [n]$ is calculated as

$$\text{Un}^{(k)}(\mathbf{X}^{\beta}; Y | \mathbf{X}^{-\beta}) := \sum_{\alpha \in \mathcal{U}^{(k)}(\beta)} I_{\partial}^{\alpha}(\mathbf{X}; Y) ,$$

with $\mathcal{U}^{(k)}(\beta) = \{\alpha \in \mathcal{A} : \beta \in \alpha, \forall \alpha \in \alpha \setminus \beta, \alpha \subseteq [n] \setminus \beta, |\alpha| > k\}$, and $\mathbf{X}^{-\beta}$ being all the variables in \mathbf{X} the indices of which are not in β . This corresponds to all the atoms where β is the only source of size k or less – which, importantly, is in general not just I_{∂}^{β} . Intuitively, this is the information that \mathbf{X}^{β} has access to and no other subset of parts has access to *on its own* (although bigger groups of other parts may). And again, for $n = 2$ we recover $\mathcal{U}^{(1)}(\{i\}) = \{\{i\}\}$, and for $n = 3$ we have e.g. $\mathcal{U}^{(1)}(\{1\}) = \{\{1\}, \{1\}\{23\}\}$.

Finally, the k^{th} -order redundancy is given by

$$\text{Red}^{(k)}(\mathbf{X}; Y) := \sum_{\alpha \in \mathcal{R}^{(k)}} I_{\partial}^{\alpha}(\mathbf{X}; Y) ,$$

with $\mathcal{R}^{(k)} = \{\alpha \in \mathcal{A} : \exists i \neq j, |\alpha_i|, |\alpha_j| \leq k\}$. Intuitively, $\text{Red}^{(k)}(\mathbf{X}; Y)$ is the information that is held by at least two different groups of size k or less. Again, in the $n = 2$ case we recover the standard redundancy $\mathcal{R}^{(1)} = \{\{1\}\{2\}\}$; and as an example for $n = 3$ we have $\mathcal{R}^{(1)} = \{\{1\}\{2\}, \{1\}\{3\}, \{2\}\{3\}, \{1\}\{2\}\{3\}\}$.

With the definitions above, we can build a coarse-grained PID which generalises the well-known construction for $n = 2$. This allows us to formulate decompositions with a small number of atoms that scale gracefully with system size, and, more interestingly, preserve the intuitive meaning that synergy, redundancy, and unique information have for $n = 2$.

Lemma 1. *The k^{th} -order synergy, redundancy, and unique information defined above provide an exact decomposition of mutual information:*

$$I(\mathbf{X}; Y) = \text{Red}^{(k)}(\mathbf{X}; Y) + \text{Syn}^{(k)}(\mathbf{X}; Y) + \sum_{\substack{\beta \subset [n]: \\ |\beta| \leq k}} \text{Un}^{(k)}(\mathbf{X}^\beta; Y | \mathbf{X}^{-\beta}). \quad (1)$$

Proof. We will prove this by showing that the sets $\mathcal{R}^{(k)}$, $\mathcal{S}^{(k)}$ and $\mathcal{U}^{(k)}(\beta)$ are a partition of \mathcal{A} . We will do this in two steps: first, we show that their intersection is empty; and second, that their union is \mathcal{A} .

Let us show that the intersections between every pair of sets is empty:

- $\mathcal{R}^{(k)} \cap \mathcal{S}^{(k)} = \emptyset$, since if $\alpha \in \mathcal{R}^{(k)}$ it must contain at least one $\alpha \in \alpha : |\alpha| \leq k$, and therefore $\alpha \notin \mathcal{S}^{(k)}$.
- $\mathcal{U}^{(k)}(\gamma) \cap \mathcal{U}^{(k)}(\beta) = \emptyset$ if and only if $\gamma \neq \beta$, since every $\alpha \in \mathcal{U}^{(k)}(\gamma)$ has either no other elements apart from γ (in which case $\beta \notin \alpha$ and thus $\alpha \notin \mathcal{U}^{(k)}(\beta)$), or other elements of cardinality greater than k (in which case, again, $\beta \notin \alpha$ and thus $\alpha \notin \mathcal{U}^{(k)}(\beta)$).
- $\mathcal{S}^{(k)} \cap \mathcal{U}^{(k)}(\beta) = \emptyset$ for all $|\beta| \leq k$, since every $\alpha \in \mathcal{U}^{(k)}(\beta)$ contains at least one element with cardinality less than or equal to k (specifically, β), and therefore $\alpha \notin \mathcal{S}^{(k)}$.
- $\mathcal{R}^{(k)} \cap \mathcal{U}^{(k)}(\beta) = \emptyset$ for all $|\alpha| \leq k$, since every $\alpha \in \mathcal{R}^{(k)}$ contains at least two sets with cardinality less than or equal to k , while by the definition of $\mathcal{U}^{(k)}(\beta)$ every element other than β must have cardinality greater than k , and thus $\alpha \notin \mathcal{U}^{(k)}(\beta)$.

This concludes the first part of the proof. Next, we need to prove that the union of those sets is indeed \mathcal{A} . We will show this by proving that every $\alpha \in \mathcal{A}$ is in one of those sets:

- If $\nexists \alpha \in \alpha : |\alpha| \leq k$, then $\alpha \in \mathcal{S}^{(k)}$.
- If there is exactly one $\alpha \in \alpha : |\alpha| \leq k$, then $\alpha \in \mathcal{U}^{(k)}(\alpha)$.
- If there are more than one $\alpha \in \alpha : |\alpha| \leq k$, then $\alpha \in \mathcal{R}^{(k)}$.

□

The two possible decompositions for $n = 3$, together with the standard PID lattice, are shown in Figure A.

2 Properties of high-order PI atoms

As discussed in Section *A formal theory of causal emergence*, our theory does not depend on a specific functional form of PID. Instead, it applies to any PID that satisfies the following properties:

- *Deterministic equality:* if there exists a function $f(\cdot)$ such that $X^i = f(X^j)$ with $i \neq j$, then the information decomposition of \mathbf{X} is isomorphic to that of \mathbf{X}^{-j} (see below).
- *Non-negativity:* $\text{Syn}^{(k)}(\mathbf{X}; Y) \geq 0$, and $\min\{I(\mathbf{Z}; Y), I(\mathbf{Z}; Y | \mathbf{X})\} \geq \text{Un}^{(k)}(\mathbf{Z}; Y | \mathbf{X}) \geq 0$.

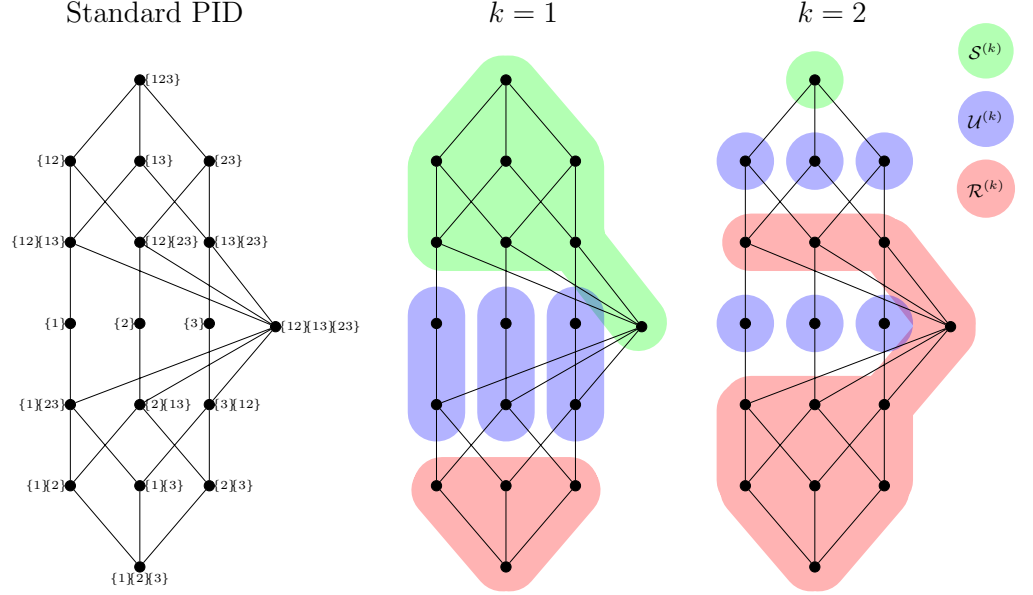


Fig A. Coarse-grained partial information decomposition of order k . (*left*) Standard PID lattice for $n = 3$, shown for reference. Node labels are omitted from the other lattices for clarity. (*middle*) Coarse-graining for $k = 1$, superimposed on the PID lattice. (*right*) Coarse-graining for $k = 2$. For both values of k , Lemma 1 guarantees that the k^{th} -order atoms provide an exact decomposition of mutual information.

- *Source data processing inequality:* $\text{Un}^{(k)}(\mathbf{W}; Y | \mathbf{X}) \leq \text{Un}^{(k)}(\mathbf{Z}; Y | \mathbf{X})$ for all $\mathbf{W} - \mathbf{Z} - (\mathbf{X}, Y)$ Markov chains.

To formulate the causal decoupling and downward causation indices in Section A *taxonomy of emergence*, we make use of ΦID , a recent extension of PID to multi-target settings [1]. As with PID, our theory does not require a particular functional form of ΦID , only the following property:

- $\sum_{|\alpha|=k} \text{Syn}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}^\alpha) \geq \mathcal{D}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) \geq \text{Syn}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}^\beta)$ for all $|\beta| = k$.

Finally, these properties are required to formulate the practical criteria for emergence in Section A *taxonomy of emergence*:

- *Whole-minus-sum:* $\text{Syn}^{(k)}(\mathbf{X}; Y) \geq I(\mathbf{X}; Y) - \sum_{|\alpha|=k} I(\mathbf{X}^\alpha; Y)$.
- *Target data processing inequality:* for all $\mathbf{X} - Y - U$ Markov chains, $\text{Syn}^{(k)}(\mathbf{X}; U) \leq \text{Syn}^{(k)}(\mathbf{X}; Y)$.

For completeness, we present a precise definition of the deterministic equality property, as previously introduced in the PID literature [2].

Definition 1. A PID satisfies deterministic equality if $I_\partial^\alpha(\mathbf{X}; Y) = I_\partial^{g_j(\alpha)}(\mathbf{X}^{-j}; Y)$ for all $\alpha \in \mathcal{A}$ whenever there exists a function $f(\cdot)$ such that $f(X^j) = X^i$ with $i \neq j$, with $g_j(\alpha)$ removes j from all the sets of indices in α .

It is direct to check that a number of well-known information decompositions, including the Minimum Mutual Information PID [3], and the corresponding ΦID [1], satisfy these requirements.

With these properties at hand we can prove the following results, used in Sections *Defining causal emergence* and *A taxonomy of emergence*:

Lemma 2. *If $X^{n+1} = \mathbf{X}$, then the following holds:*

$$\mathbf{Syn}^{(k)}(\mathbf{X}; Y) = \mathbf{Un}^{(k)}(X^{n+1}; Y | \mathbf{X}) . \quad (2)$$

Above, the second term corresponds to a PID over a system of $n + 1$ elements.

Proof. Begin by considering the PID of n sources X^1, \dots, X^n on the lattice \mathcal{A}^n , and define its set $\mathcal{S}^{(k)}$ as above. Now we add an additional $n + 1^{\text{st}}$ variable that is simply all of them concatenated, $X_{n+1} = \mathbf{X}^n$, and build a PID on the lattice \mathcal{A}^{n+1} . In the following, we consider the set $\mathcal{S}^{(k)}$ to belong to the lattice \mathcal{A}^n , and the set $\mathcal{U}^{(k)}(\beta)$ to belong to the lattice \mathcal{A}^{n+1} . To prove the lemma we need four ingredients, which we provide in the four paragraphs below.

1. First, note that the nodes in \mathcal{A}^{n+1} that precede $\{n + 1\}$ are those in \mathcal{A}^n , but with the singleton $\{n + 1\}$ appended to them. More specifically, the mapping $f : \mathcal{A}^n \rightarrow \mathcal{A}^{n+1}$ of the form $f(\alpha) = \alpha \cup \{\{n + 1\}\}$ is such that for any $\alpha \in \mathcal{A}^n$ then $f(\alpha) \prec \{\{n + 1\}\}$. Additionally, due to the properties of the partial order, $\alpha \preceq \alpha'$ if and only if $f(\alpha) \preceq f(\alpha')$. This shows that \mathcal{A}^n is isomorphic to a sublattice of \mathcal{A}^{n+1} .
2. Next, by the *deterministic equality* property it is direct to check that $I_{\cap}^{f(\alpha)} = I_{\cap}^{\alpha}$. Since this equality holds for all $\alpha \in \mathcal{A}^n$, then applying a Möbius inversion we directly obtain that $I_{\partial}^{f(\alpha)} = I_{\partial}^{\alpha}$.
3. Additionally, by construction of $\mathcal{U}^{(k)}$ and $\mathcal{S}^{(k)}$, for all $\gamma \in \mathcal{S}^{(k)}$ one has $f(\gamma) \in \mathcal{U}^{(k)}(\{n + 1\})$. In other words, the set $\mathcal{U}^{(k)}(\{n + 1\})$ includes all atoms in $\mathcal{S}^{(k)}$, plus $\{n + 1\}$.
4. Finally, note that the node $\{12\dots n\}\{n + 1\}$ is the only direct predecessor of $\{n + 1\}$, since there exists no node $\beta \in \mathcal{A}^{n+1}$ such that $\beta \prec \{n + 1\}$ and not $\beta \preceq \{12\dots n\}\{n + 1\}$. By the *deterministic equality* property $I_{\cap}^{\{12\dots n\}\{n+1\}} = I_{\cap}^{\{12\dots n\}}$ and, therefore, $I_{\partial}^{\{n+1\}} = 0$.

With all of these, it is direct to see that

$$\begin{aligned} \mathbf{Un}^{(k)}(X^n; Y | \mathbf{X}) &= \sum_{\alpha \in \mathcal{U}^{(k)}(\{n+1\})} I_{\partial}^{\alpha}(\mathbf{X}; Y) \\ &= \sum_{\alpha \in \mathcal{S}^{(k)}} I_{\partial}^{\alpha}(\mathbf{X}; Y) = \mathbf{Syn}^{(k)}(\mathbf{X}; Y) . \end{aligned}$$

□

Corollary 1. *If $X^{n+1} = \mathbf{X}$ and $Y^{n+1} = \mathbf{Y}$, then the following holds:*

$$\mathcal{G}^{(k)}(\mathbf{X}; \mathbf{Y}) = \mathbf{Un}^{(k)}(X^{n+1}; Y^{n+1} | \mathbf{X}, \mathbf{Y}) \quad (3)$$

Above, the second term corresponds to a Φ ID over a system of $n + 1$ elements.

Proof. Follows from a direct Φ ID extension to the proof of Lemma 2, by formulating a decomposition

$$I(\mathbf{X}; \mathbf{Y}) = \sum_{\alpha, \beta \in \mathcal{A}^n} I_{\partial}^{\alpha \rightarrow \beta}(\mathbf{X}; \mathbf{Y}) ,$$

and applying the proof above to both α and β . Strictly speaking, this also requires a natural multi-target extension of the Deterministic Equality property, namely that $I_{\cap}^{\alpha \rightarrow \beta}(\mathbf{X}; \mathbf{Y}) = I_{\cap}^{\alpha \rightarrow \beta}(\mathbf{X}^{-j}; \mathbf{Y})$ if $X^i = f(X^j)$ with $j \neq i$, for any $\beta \in \mathcal{A}^n$; as well as the symmetric $I_{\cap}^{\alpha \rightarrow \beta}(\mathbf{X}; \mathbf{Y}) = I_{\cap}^{\alpha \rightarrow \beta}(\mathbf{X}; \mathbf{Y}^{-j})$ if $Y^i = f(Y^j)$ with $j \neq i$, for any $\alpha \in \mathcal{A}^n$. □

3 Mathematical properties of causal emergence

Proof of Lemma 1. The first property can be easily proven by noting that there can be no synergy in a univariate system: i.e. if $n = 1$, then $\mathcal{S}^{(k)} = \emptyset$ and therefore $\text{Syn}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) = 0$.

To prove the second property, let us assume that there exists a function $g(\cdot)$ such that $g(X_t^j) = V_t$ for some j . Then, one can show that

$$\text{Un}^{(1)}(V_t; \mathbf{X}_{t'} | \mathbf{X}_t) \leq I(V_t; \mathbf{X}_{t'} | \mathbf{X}_t) \quad (4)$$

$$\leq I(V_t; \mathbf{X}_{t'} | X_t^j) \quad (5)$$

$$\leq H(V_t | X_t^j) \quad (6)$$

$$= 0, \quad (7)$$

where (4) is due to the non-negativity property of $\text{Un}^{(k)}$ introduced in Section *Partial information decomposition*. This shows that V_t cannot exhibit emergent behaviour. \square

Proof of Theorem 1. If $\text{Syn}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) > 0$, then by Lemma 2 it is clear that the feature $V_t = \mathbf{X}_t$ exhibits causal emergence.

To prove the converse, note that all supervenient features follow the Markov chain structure $V_t - \mathbf{X}_t - \mathbf{X}_{t'}$. Therefore, for any supervenient feature $V_t = f(\mathbf{X}_t)$ the following holds:

$$0 \leq \text{Un}^{(k)}(V_t; \mathbf{X}_{t'} | \mathbf{X}_t) \leq \text{Un}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'} | \mathbf{X}_t) \quad (8)$$

$$= \text{Syn}(\mathbf{X}_t; \mathbf{X}_{t'}) , \quad (9)$$

where (8) is due to the data processing inequality of the unique information (c.f. Section *Partial information decomposition*), and (9) is due to Lemma 2. Using this result, is clear that if $\text{Syn}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) = 0$ then $\text{Un}^{(k)}(V_t; \mathbf{X}_{t'} | \mathbf{X}_t) = 0$ for any supervenient feature V_t . \square

The above proof implies that the system has emergent behaviour if and only if the system as a whole seen as a feature (i.e. $V_t = \mathbf{X}_t$) is causally emergent. Note that the fact that this trivial feature is helpful for detecting the presence of emergence doesn't imply that it is an appropriate way of representing it, as in most cases also carries non-interesting information, and in practice one may be interested in features that exhibit emergence but are shorter to describe than the microstate of the system, i.e. $H(V_t) < H(\mathbf{X}_t)$.

Proof of Theorem 2. Let us first assume that there exists a supervenient feature V_t such that $\text{Un}^{(k)}(V_t; \mathbf{X}_{t'}^\alpha | \mathbf{X}) > 0$ for some $\alpha : |\alpha| = k$. Then, one can find that

$$0 < \text{Un}^{(k)}(V_t; \mathbf{X}_{t'}^\alpha | \mathbf{X}_t) \leq \text{Un}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}^\alpha | \mathbf{X}_t) \quad (10)$$

$$= \text{Syn}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}^\alpha) \quad (11)$$

$$\leq \mathcal{D}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) . \quad (12)$$

Above, (10) is a consequence of the data processing inequality of the unique information, (11) comes from Lemma 2, and (12) is from the properties of $\mathcal{D}^{(k)}$ stated in Section *A taxonomy of emergence*.

To prove the converse, let us assume that all supervenient features V_t satisfy $\text{Un}^{(k)}(V_t; \mathbf{X}_{t'}^\alpha | \mathbf{X}_t) = 0$ for all k and $|\alpha| = k$. In particular, this is true for the feature $V_t = \mathbf{X}_t$. Then, another application of Lemma 2 shows that

$$\begin{aligned} \mathcal{D}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) &\leq \sum_{|\alpha|=k} \text{Syn}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}^\alpha) \\ &= \sum_{|\alpha|=k} \text{Un}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}^\alpha | \mathbf{X}_t) \\ &= 0 . \end{aligned} \tag{13}$$

Above, (13) is a consequence of the properties of $\mathcal{D}^{(k)}$. □

Proof of Theorem 3. We begin by proving that a system has a causally decoupled feature iff $\mathcal{G}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) > 0$. This proof follows a similar structure to that of Theorem 2 for downward causation.

Let us first assume that there exists a supervenient feature V_t with $\text{Un}^{(k)}(V_t; V_{t'} | \mathbf{X}_t, \mathbf{X}_{t'}) > 0$. Then,

$$\begin{aligned} 0 &< \text{Un}^{(k)}(V_t; V_{t'} | \mathbf{X}_t, \mathbf{X}_{t'}) \\ &\leq \text{Un}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'} | \mathbf{X}_t, \mathbf{X}_{t'}) \end{aligned} \tag{14}$$

$$= \mathcal{G}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) . \tag{15}$$

Above, (14) can be obtained through a combination of the data processing inequalities of the unique information and the synergy, as well as Lemma 2; and (15) is an application of Corollary 1.

To prove the converse, assume that for all supervenient features $\text{Un}^{(k)}(V_t; V_{t'} | \mathbf{X}_t, \mathbf{X}_{t'}) = 0$. As before, this also includes $V_t = \mathbf{X}_t$. Therefore, applying Corollary 1 we arrive at $\mathcal{G}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) = \text{Un}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'} | \mathbf{X}_t, \mathbf{X}_{t'}) = 0$, which concludes the proof.

We now move on to prove the results on perfectly causally decoupled systems, where $\mathcal{G}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) > 0$ and $\mathcal{D}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) = 0$. As $\text{Syn}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) = \mathcal{G}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) + \mathcal{D}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) > 0$, thanks to Theorem 1 this is equivalent to the existence of at least one emergent feature V_t . Additionally, due to Theorem 2, $\mathcal{D}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) = 0$ is equivalent to $\text{Un}^{(k)}(V_t; \mathbf{X}_{t'}^\alpha | \mathbf{X}_t) = 0$ for all emergent features and $\alpha : |\alpha| = k$ □

4 Mathematical properties of emergence criteria

In this appendix we provide the necessary proofs linking the practical criteria for emergence in Eqs. (10a-10c) from the main text with the definitions in Section *A formal theory of causal emergence*. For completeness, we provide formulae of Ψ , Δ , and Γ for arbitrary emergence order k :

$$\Psi_{t,t'}^{(k)}(V) := I(V_t; V_{t'}) - \sum_{|\alpha|=k} I(\mathbf{X}_t^\alpha; V_{t'}) , \tag{16a}$$

$$\Delta_{t,t'}^{(k)}(V) := \max_{|\alpha|=k} \left(I(V_t; \mathbf{X}_{t'}^\alpha) - \sum_{|\beta|=k} I(\mathbf{X}_t^\beta; \mathbf{X}_{t'}^\alpha) \right) , \tag{16b}$$

$$\Gamma_{t,t'}^{(k)}(V) := \max_{|\alpha|=k} I(V_t; \mathbf{X}_{t'}^\alpha) . \tag{16c}$$

Proof of Proposition 1. To start, note that a direct calculation shows that

$$\Psi_{t,t'}^{(k)}(V) \leq I(\mathbf{X}_t; V_{t'}) - \sum_{|\alpha|=k} I(\mathbf{X}_t^\alpha; V_{t'}) \quad (17)$$

$$\leq \text{Syn}^{(k)}(\mathbf{X}_t; V_{t'}) \quad (18)$$

$$\leq \text{Syn}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) . \quad (19)$$

Above, (17) is due to the data processing inequality applied over the Markov chain $V_t - \mathbf{X}_t - \mathbf{X}_{t'} - V_{t'}$; (18) is due to the whole-minus-sum property of the synergy; and (19) is due to the data processing inequality of the synergy. Therefore, it is clear that if $\Psi_{t,t'}^{(k)}(V) > 0$ for some feature V_t then $\text{Syn}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) > 0$. This, combined with Theorem 1, guarantees that the system exhibits causal emergence.

To check the condition for downward causation, a direct calculation shows that, for some $|\alpha| = k$,

$$\Delta_{t,t'}^{(k)}(V) \leq I(\mathbf{X}_t; \mathbf{X}_{t'}^\alpha) - \sum_{|\beta|=k} I(\mathbf{X}_t^\beta; \mathbf{X}_{t'}^\alpha) \quad (20)$$

$$\leq \text{Syn}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}^\alpha) \quad (21)$$

$$\leq \mathcal{D}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) . \quad (22)$$

Above, (20) is due to the data processing inequality applied over the Markov chain $V_t - \mathbf{X}_t - \mathbf{X}_{t'}^\alpha$; (21) to the whole-minus-sum property of the synergy; and (22) to the properties of $\mathcal{D}^{(k)}$. From here, is clear that if $\Delta_{t,t'}^{(k)}(V) > 0$ for some feature V_t then $\mathcal{D}^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) > 0$. This, combined with Theorem 1, guarantees that the system exhibits downward causation.

Finally, for the condition for causal decoupling it is sufficient to note that

$$\Gamma_{t,t'}^{(k)}(V) = \max_{|\alpha|=k} I(V_t; \mathbf{X}_{t'}^\alpha) \quad (23)$$

$$\geq \text{Un}^{(k)}(V_t; \mathbf{X}_{t'}^\alpha | \mathbf{X}_t) \geq 0 , \quad (24)$$

where the inequality is due to the bounds of the unique information. \square

Proof of Proposition 2. Let us consider V_t a k -synergistic observable. Due to stationarity, the fact that $V_t \in \mathcal{C}_k(\mathbf{X}_t)$ implies that $V_{t'} \in \mathcal{C}_k(\mathbf{X}_{t'})$. Using this fact, and noting that $V_t - \mathbf{X}_t - \mathbf{X}_{t'} - V_{t'}$ is a Markov chain, it is direct to check that

$$\mathcal{G}_\star^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) \geq I(V_t; V_{t'}) > 0 , \quad (25)$$

Additionally, since $\mathcal{D}_\star^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) \geq 0$, Eq. (25) implies that $\text{Syn}_\star^{(k)}(\mathbf{X}_t; \mathbf{X}_{t'}) > 0$, proving the desired result. \square

5 Simulation details

Let us focus first on the Game of Life (GoL) simulations illustrated by Figure 4 in the main text. For the initial state, two particles of three fixed types (nothing, a glider, or a lightweight spaceship) were selected at random, and placed at random positions in a 15x15 square cell array. The GoL evolution rule was run for 1000 steps, which, in most cases (judged by visual inspection) was enough for the system to settle on a stable configuration – which was typically either nothing, a small number of static structures, or a small number of particles in non-colliding tracks. To compute the emergent feature

V_t , particles were detected by simply pattern-matching the resulting system state against the known shapes of each particle. We considered five categories: still lifes, oscillators, gliders, lightweight spaceships, or nothing.

For static structures, we used a single symbol to represent all still lifes, and a single symbol for all oscillators. This particle detector was found to be very effective, with only 2% of runs resulting in unrecognised particles. A total of 5×10^4 independent runs were simulated and, due to the high number of possible states of V_t , to reduce bias we used the quasi-Bayesian estimator by Archer *et al.* [4].

For the boids simulation, $N = 10$ boids were simulated on a torus of side length $L = 200$. Boids are initialised with random positions, speeds, and head angles, and at each timestep each boid $i = 1, \dots, N$ is updated according to the equations:

$$\begin{aligned}x_{t'}^i &= x_t^i + s^i \cos(\alpha_t^i) \\y_{t'}^i &= y_t^i + s^i \sin(\alpha_t^i) \\ \alpha_{t'}^i &= \alpha_t^i + a_1\theta_1 + a_2(\pi + \theta_2) + a_3\theta_3 ,\end{aligned}$$

Where θ_1 is the bearing to the flock’s center of mass, θ_2 the bearing to the nearest boid, and θ_3 the mean alignment of all boids within a 20 unit radius. The scalars a_1, a_2, a_3 are the aggregation, avoidance, and alignment parameters, respectively.

The results in Figure 5 in the main text were obtained averaging Ψ across 25 independent runs of 5000 timesteps each, keeping $a_1 = 0.15$, $a_3 = 0.25$ fixed across all simulations. To compute Ψ , we pre-processed the trajectories using the same procedure as Seth [5] (i.e. each boid was described by its distance to the center of the environment and all time series were first-order differenced), and information-theoretic quantities were computed using the non-parametric estimator implemented in the JIDT toolbox [6, 7] using a dynamic correlation exclusion window of 10 samples [8].

To compute the uncertainties and error bars reported in the text and figures we used standard surrogate data methodology: first, system trajectories are time-shuffled to generate one set of surrogate time series, then the quantities of interest (e.g. $\Psi_{t,t'}^{(1)}$) are estimated on the surrogate data, and standard deviations over multiple realisations of the surrogates are reported.

6 ECoG preprocessing and decoding

ECoG signals were preprocessed following the steps presented in the original publication: specifically, the data was notch-filtered to remove line noise and band-pass filtered from 0.1 Hz to 600 Hz, and then downsampled to the same sampling frequency as the MoCap data (120 Hz) [9]. Then, data were divided into a 60/40 train/test split. To build the predictor, the training set was first standardised to zero mean and unit variance, and then it was dimensionality-reduced with a 20-component Partial Least Squares (PLS) regression using the three dimensions of the MoCap as dependent variables. We then trained a non-linear Support Vector Machine (SVM) [10] using a squared exponential kernel predicting the MoCap data from the 20-dimensional PLS scores. SVM and kernel hyperparameters were optimised through 5-fold cross-validation, and a model with the best hyperparameter configuration was trained on the full training set. This cross-validation and optimisation procedure was repeated for each dimension of the MoCap data, resulting in three separate SVMs.¹

¹Note that the accuracy of this model is lower than the model presented in Ref. [9], for a precise reason: to predict position at time t , Chao *et al.* use a wavelet transform to obtain features that have represent ECoG from $t-1.1$ s to t . Unlike Chao *et al.*, we work under the constraint of using supervenient features: for the axioms of the theory (in its current form) to hold, V_t needs to be a function of \mathbf{X}_t only. This limits the possibilities for feature extraction, and naturally may result in lower-accuracy decoders.

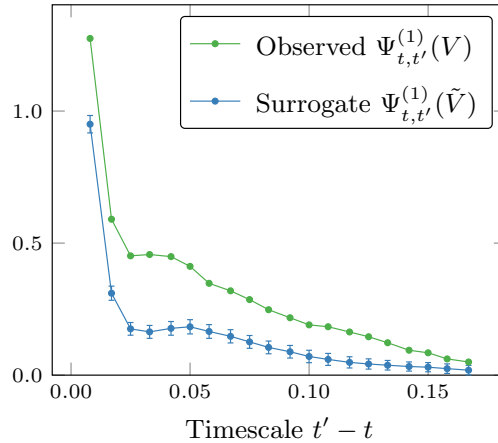


Fig B. Surrogate data test for emergence in monkey ECoG. As a control, we computed the emergence criterion Ψ on a surrogate feature \tilde{V} , and found that it yields significantly lower Ψ than the original data (error bars are standard error of the mean over 10 realisations of \tilde{V} ; see text for details).

The computation of Ψ was performed on the held-out test data set. ECoG signals were standardised and projected onto the PLS latent space (using the means, standard deviations, and projection matrix obtained from the training data), and mutual information terms involved in the computation of Ψ were calculated using the Gaussian estimator implemented in the open-source JIDT toolbox [7].

Additionally, we performed a control using surrogate data to confirm the results were not driven by the autocorrelation or filtering properties of the ECoG, or the regularisation of the SVM. For this, we time-shuffled the wrist position time series, trained an SVM, and evaluated it on a separate (un-shuffled) ECoG test set, resulting in a surrogate feature \tilde{V}_t that preserves the autocorrelation properties of the ECoG but does not meaningfully extract motor information. Then we computed $\Psi_{t,t'}^{(1)}(\tilde{V})$, repeated over 10 realisations (i.e. shuffling and training a separate PLS-SVM 10 times), and found that the observed $\Psi_{t,t'}^{(1)}(V)$ lies significantly above its surrogate, confirming the observed emergence criterion scores over and above what would be expected by the autocorrelation structure of the ECoG (Fig. B).

References

1. Mediano PA, Rosas F, Carhart-Harris RL, Seth AK, Barrett AB. Beyond integrated information: A taxonomy of information dynamics phenomena. arXiv preprint arXiv:190902297. 2019;
2. Kolchinsky A. A novel approach to multivariate redundancy and synergy. arXiv preprint arXiv:190808642. 2019;
3. Barrett AB. Exploration of synergistic and redundant information sharing in static and dynamical Gaussian systems. *Physical Review E*. 2015;91:052802. doi:10.1103/PhysRevE.91.052802.
4. Archer E, Park I, Pillow J. Bayesian and quasi-Bayesian estimators for mutual information from discrete data. *Entropy*. 2013;15(5):1738–1755.
5. Seth AK. Measuring autonomy and emergence via Granger causality. *Artificial Life*. 2010;16(2):179–196.

6. Kraskov A, Stögbauer H, Grassberger P. Estimating mutual information. *Physical Review E*. 2004;69(6):066138. doi:10.1103/PhysRevE.69.066138.
7. Lizier JT. JIDT: An Information-Theoretic Toolkit for Studying the Dynamics of Complex Systems. *Frontiers in Robotics and AI*. 2014;1:37. doi:10.3389/frobt.2014.00011.
8. Kantz H, Schreiber T. *Nonlinear Time Series Analysis*. Cambridge University Press; 2004.
9. Chao Z, Nagasaka Y, Fujii N. Long-term asynchronous decoding of arm motion using electrocorticographic signals in monkey. *Frontiers in Neuroengineering*. 2010;3:3. doi:10.3389/fneng.2010.00003.
10. Bishop CM. *Pattern Recognition and Machine Learning*. Springer; 2006.