# Information Decomposition Diagrams Applied beyond Shannon Entropy: a Generalization of Hu's Theorem

Leon Lang[1], Pierre Baudot[2], Rick Quax[1], and Patrick Forré[1]

[1]Informatics Institute, University of Amsterdam

[2]Median Technologies, France

In information theory, one major goal is to find useful functions that summarize the amount of information contained in the interaction of several random variables. Specifically, one can ask how the classical Shannon entropy, mutual information, and higher interaction information relate to each other. This is answered by Hu's theorem, which is widely known in the form of information diagrams: it relates shapes in a Venn diagram to information functions, thus establishing a bridge from set theory to information theory. In this work, we view random variables together with the joint operation as a monoid that acts by conditioning on information functions, and entropy as a function satisfying the chain rule of information. This abstract viewpoint allows to prove a *generalization of Hu's theorem*. It applies to Shannon and Tsallis entropy, (Tsallis) Kullback-Leibler Divergence, cross-entropy, Kolmogorov complexity, submodular information functions, and the generalization error in machine learning. Our result implies for Chaitin's Kolmogorov complexity that the *interaction complexities* of all degrees are in expectation close to Shannon interaction information. For well-behaved probability distributions on increasing sequence lengths, this shows that the per-bit expected interaction complexity and information asymptotically coincide, thus showing a strong bridge between algorithmic and classical information theory.

## Contents

Leon Lang: l.lang@uva.nl, 0000-0002-1950-2831. Main contributing author.

Pierre Baudot: pierre.baudot@mediantechnologies.com, 0000-0002-5574-6809

Rick Quax: r.quax@uva.nl, 0000-0002-0299-0074

Patrick Forré: p.d.forre@uva.nl, 0000-0003-4663-3842

# 1  Introduction

Information diagrams, most often drawn for two or three random variables (see Figures 1 and 2), provide a concise way to visualize information functions. Not only do they show (conditional) Shannon entropy [48], mutual information, and interaction information — also called co-information [8] — of several random variables in one overview, they also provide an intuitive account of the *relations* between these functions.

This well-known fact goes beyond just three variables: diagrams with four (see Figure 3) and more variables exist as well. Hu's theorem [32, 63, 64] renders all this mathematically precise by connecting the set-theoretic operations of union, intersection, and set difference to joint information, interaction information, and conditioning of information functions, respectively. The map from sets to information functions is then a *measure* and turns disjoint unions into sums. Certain summation rules of information functions then follow visually from disjoint unions in the diagrams.

Our work is concerned with the question of whether Hu's theorem can be generalized to other information functions than entropy, such as Kullback-Leibler divergence and cross-entropy. Such functions are important in the context of statistical modeling of multivariate data, in which one aims to find a probabilistic model able to reproduce the information structure of the data. For instance, an information diagram for cross-entropy would then allow to visualize how the cross-entropy between a model probability distribution and the data distribution is decomposed into higher-order terms. [16] used these higher-order terms (which they called cluster (cross)-entropies) in their adaptive cluster expansion approach to statistical modeling of data with Ising models.

Kullback-Leibler divergence has been studied in the context of decompositions of joint entropy and information [1] and is often minimized in machine learning and deep learning [11, 12]. This becomes especially interesting for graphical methods, including diffusion models [50], which form the basis for widespread text-to-image generation methods like Dalle [41], Imagen [44], and stable diffusion [42]. Diffusion models involve a decomposition of a joint Kullback-Leibler divergence over a Markov chain. Once information diagrams are established in a generalized context, this might facilitate to study decompositions of loss functions for more general graphical models.

Our claim is that the language employed in the foundations of information cohomology [5] gives the perfect starting point for generalizing Hu's theorem. Namely, by replacing discrete random variables with partitions on a sample space, they give random variables the structure of a *monoid* that is commutative and idempotent. Furthermore, conditional information functions are formally described by a *monoid action*. And finally, the most basic information function that generates all others, Shannon entropy, is fully characterized as the unique function that satisfies the *chain rule of information*. We substantially generalize Hu's theorem by giving a proof only based on the properties just mentioned, leading to new applications to Kolmogorov complexity, Kullback-Leibler divergence, and beyond.

To clarify, the main contribution of this work is not to provide major previously unknown ideas — indeed, our proof is very similar to the original one given in [63] — but instead, to place and prove this result in its proper abstract context. This then reveals information diagrams for new information measures.

Section 2 summarizes classical definitions and results for Shannon information theory, generalized to *countable* discrete random variables to be later applied to Kolmogorov complexity. Section 3 — which can be read independently of the preceding section — contains our main result, the generalized Hu theorem. In Section 4, we prove a Hu theorem for Kolmogorov complexity. We also combine Hu's theorems for Shannon entropy and Kolmogorov complexity to generalize the well-known result that "expected Kolmogorov complexity is close to entropy" [28]: general *interaction complexity* is close to interaction information. For the case of well-behaved sequences of probability measures on binary strings with increasing length, this leads to an asymptotic result: in the limit of infinite sequence length, the *per-bit* interaction complexity and interaction information coincide. In Section 5, we consider further examples of Hu's theorem, including Kullback-Leibler divergence and the generalization error in machine learning. We conclude with a discussion in Section 6, followed by proofs in the appendices.

## Preliminaries and Notation

We mainly assume the reader to be familiar with the basics of measure theory and probability theory. They can be learned from any book on the topic, for example [45] or [54]. The main concepts we assume to be known are $\sigma$-algebras, the Borel $\sigma$-algebra on $\mathbb{R}^n$, measurable spaces, measures, measure spaces, probability measures, probability spaces, and random variables. We assume some very basic familiarity with abelian groups, (commutative, idempotent) monoids, and additive monoid actions. In contrast, we carefully define all basic notions from (algorithmic) information theory from scratch.

On notation: to aid familiarity, we will start writing the Shannon entropy with the symbol $H$, but then switch to the notation $I_1$ once we embed Shannon entropy in the concept of interaction information, Definition 2.7. Instead of the typical notation $H(Y \mid X)$ for the conditional entropy, we will use $X.H(Y) = X.I_1(Y)$. This is the general notation of monoid actions and is thus preferable in our abstract context. Furthermore, for two disjoint sets $A$ and $B$, we write their union as $A \dot\cup B$. The number of elements in $A$ is written as $|A|$. The power set of $A$, i.e., the set of its subsets, is denoted $2^A$. And finally, the natural and binary logarithms of $x$ are denoted $\ln(x)$ and $\log(x)$, respectively.

## 2    Preliminaries on Shannon Entropy of Countable Discrete Random Variables

In this technical introduction, we explain preliminaries on discrete random variables, entropy, mutual information, and interaction information. Our treatment will also emphasize abstract structures that lead us to the generalizations in Section 3. The goal is to arrive at Summary 2.17, which summarizes the properties of classical information functions in an abstract way suitable for our generalizations. We will omit many proofs of elementary and well-known results. When we say a set is *countable*, then we mean it is finite or countably infinite. Whenever we talk about *discrete measurable spaces*, we mean countable measurable spaces in which all subsets are measurable. Some technical considerations related to the measurability of certain functions in the infinite, discrete case are found in Appendix A.

### 2.1   Entropy, Mutual Information, and Interaction Information

We fix in this section a discrete sample space $\Omega$. We define

$$\Delta(\Omega) := \left\{ P : \Omega \to [0,1] \;\middle|\; \sum_{\omega \in \Omega} P(\omega) = 1 \right\} = \left\{ (p_\omega)_{\omega \in \Omega} \in [0,1]^\Omega \;\middle|\; \sum_{\omega \in \Omega} p_\omega = 1 \right\}.$$

If $\Omega$ is finite, we view it as a measurable space with the $\sigma$-algebra of Borel measurable sets. When $\Omega$ is infinite and discrete, we equip $\Delta(\Omega)$ with the smallest $\sigma$-algebra that makes all evaluation maps

$$\mathrm{ev}_A : \Delta(\Omega) \to \mathbb{R}, \quad P \mapsto \mathrm{ev}_A(P) := P(A)$$

for all subsets $A \subseteq \Omega$ measurable. In the finite case, this definition is equivalent to the one given before. We remark that we do not distinguish between probability measures and their mass functions in the notation or terminology: for a subset $A \subseteq \Omega$ and a probability measure $P : \Omega \to [0,1]$, we simply set $P(A) := \sum_{\omega \in A} P(\omega)$.

Our aim is the study of discrete random variables $X : \Omega \to E_X$. Here, being discrete means that $E_X$ — next to $\Omega$ — is discrete. For any probability measure $P$ on $\Omega$ and any random variable $X : \Omega \to E_X$, we define the *distributional law* $P_X : E_X \to [0,1]$ as the unique probability measure with

$$P_X(x) := P(X^{-1}(x)) = \sum_{\omega \in X^{-1}(x)} P(\omega)$$

for all $x \in X$. Clearly, $P_X \in \Delta(E_X)$.

For the following definition of Shannon entropy, introduced in [48, 49], we employ the convention $0 \cdot \infty = 0 \cdot (-\infty) = 0$ and $\ln(0) = -\infty$. Furthermore, set $\overline{\mathbb{R}} := \mathbb{R} \cup \{+\infty\}$.

**Definition 2.1** (Shannon Entropy). *Let $P \in \Delta(\Omega)$ be a probability measure. Then the* Shannon entropy *of $P$ is given by*

$$H(P) := -\sum_{\omega \in \Omega} P(\omega) \ln P(\omega) \in \overline{\mathbb{R}}.$$

*Here, $\ln : [0, \infty) \to \mathbb{R} \cup \{-\infty\}$ is the natural logarithm. Now, let $X : \Omega \to E_X$ be a discrete random variable. The* Shannon entropy *of $X$ with respect to $P \in \Delta(\Omega)$ is given by*

$$H(X; P) := H(P_X) = -\sum_{x \in E_X} P_X(x) \ln P_X(x) \in \overline{\mathbb{R}}.$$

For the identity $\mathrm{id}_\Omega : \Omega \to \Omega,\ \omega \mapsto \omega$ we then have $P_{\mathrm{id}_\Omega} = P$ and therefore $H(\mathrm{id}_\Omega; P) = H(P)$. For discrete probability distributions with infinite Shannon entropy, see [2].

Now, set

$$\Delta_f(\Omega) := \Delta(\Omega) \setminus \{P \in \Delta(\Omega) \mid H(P) = \infty\}.$$

$\Delta_f(\Omega)$ is the measurable space of probability measures with finite entropy. We restrict entropy functions to this space for algebraic reasons:

**Definition 2.2** (Entropy Function of a Random Variable). *Let $X : \Omega \to E_X$ be a discrete random variable. Then its* entropy function *or* Shannon entropy *is the measurable function*

$$H(X) : \Delta_f(\Omega) \to \mathbb{R}, \quad P \mapsto H(X; P)$$

*defined on probability measures with finite entropy. Its measurability is proven in Corollary A.3.*

Let $P : \Omega \to \mathbb{R}$ be a probability measure and $X : \Omega \to E_X$ a discrete random variable. Then we define the conditional probability measure $P|_{X=x} : \Omega \to \mathbb{R}$ by

$$P|_{X=x}(\omega) := \begin{cases} \frac{P\left(\{\omega\} \cap X^{-1}(x)\right)}{P_X(x)}, & P_X(x) \neq 0; \\ P(\omega), & P_X(x) = 0.^1 \end{cases} \tag{1}$$

For all $A \subseteq \Omega$, we then have

$$P|_{X=x}(A) = \begin{cases} \frac{P\left(A \cap X^{-1}(x)\right)}{P(X^{-1}(x))}, & P_X(x) \neq 0; \\ P(A), & P_X(x) = 0. \end{cases}$$

For the following definition, recall that a series of real numbers converges absolutely if the series of its absolute values converges. It converges unconditionally if every reordering of the original series still converges with the same limit. According to the Riemann series theorem [36], these two properties are equivalent.

**Definition 2.3** (Conditionable Functions, Averaged Conditioning). *Let $F : \Delta_f(\Omega) \to \mathbb{R}$ be a measurable function. $F$ is called* conditionable *if for all discrete random variables $X : \Omega \to E_X$ and all $P \in \Delta_f(\Omega)$, the sum*

$$(X.F)(P) := \sum_{x \in E_X} P_X(x) F(P|_{X=x}) \tag{2}$$

*converges unconditionally. Note that $P|_{X=x} \in \Delta_f(\Omega)$, which makes $F(P|_{X=x})$ in Equation (2) well-defined.*

*For all conditionable measurable functions $F : \Delta_f(\Omega) \to \mathbb{R}$ and all discrete random variables $X : \Omega \to E_X$, the function $X.F : \Delta_f(\Omega) \to \mathbb{R}$ is a measurable function by Corollary A.5, which we call the* averaged conditioning *of $F$ by $X$. The space of all conditionable measurable functions $F : \Delta_f(\Omega) \to \mathbb{R}$ is denoted by $\mathrm{Meas}_{\mathrm{con}}(\Delta_f(\Omega), \mathbb{R})$.*

---

[1] Note that the precise definition for the case $P_X(x) = 0$ does not matter since it almost surely does not appear. However, defining the conditional also in this case makes many formulas simpler since we do not need to restrict sums involving $P|_{X=x}$ to the case $P_X(x) \neq 0$.

If $X : \Omega \to E_X$ and $Y : \Omega \to E_Y$ are two (not necessarily discrete) random variables, then their (Cartesian) product, or joint variable, $XY : \Omega \to E_X \times E_Y$ is defined by

$$(XY)(\omega) \coloneqq \big(X(\omega), Y(\omega)\big) \in E_X \times E_Y.^2 \tag{3}$$

If we have two discrete random variables $X$ and $Y$ and a probability measure $P \in \Delta(\Omega)$, then this allows to consider $(P|_{X=x})_Y(y)$ for $(x, y) \in E_X \times E_Y$. In order to not overload notation, we will write this often as $P(y \mid x)$. Similarly, we will often write $P(x) \coloneqq P_X(x)$ and $P(\omega \mid x) \coloneqq P|_{X=x}(\omega)$. We obtain the following elementary lemma and corollary whose proofs are left to the reader:

**Lemma 2.4.** *Let $Y$ be a discrete random variable on $\Omega$. Then $H(Y)$ is conditionable. More precisely, for another discrete random variable $X$ on $\Omega$ and $P \in \Delta_f(\Omega)$, $H(X; P)$ and $H(XY; P)$ are finite and we have*

$$\big[X.H(Y)\big](P) = H(XY; P) - H(X; P),$$

*which results in $\big[X.H(Y)\big](P)$ converging unconditionally.*

**Corollary 2.5.** *The following chain rule*

$$H(XY) = H(X) + X.H(Y)$$

*holds for arbitrary discrete random variables $X : \Omega \to E_X$ and $Y : \Omega \to E_Y$.*

We will also write $Y.F(P) \coloneqq (Y.F)(P)$. For example, if $F = H(X)$ is the Shannon entropy of the discrete random variable $X$, we write

$$Y.H(X; P) = Y.H(X)(P) = [Y.H(X)](P) = \sum_{y \in E_Y} P_Y(y) H(X; P|_{Y=y}).$$

We emphasize explicitly that $Y$ can not act on $H(X; P)$ since this is only a number, and not a measurable function. Nevertheless, we find the notation $Y.H(X; P)$ for $[Y.H(X)](P)$ convenient. We obtain the following properties resembling those of an additive monoid action:

**Proposition 2.6.** *Let $X, Y$ be two discrete random variables on $\Omega$, $\mathbf{1} : \Omega \to * \coloneqq \{*\}$ a trivial random variable, and $F, G : \Delta_f(\Omega) \to \mathbb{R}$ two conditionable measurable functions. Then the following hold:*

1. *$\mathbf{1}.F = F$;*

2. *$Y.F$ is also conditionable, and we have $X.(Y.F) = (XY).F$;*

3. *$F + G$ is also conditionable, and we have $X.(F + G) = X.F + X.G$.*

*Proof.* Properties 1 and 3 are elementary and left to the reader to prove. 2 follows from $P(x, y) = P(x) \cdot P(x \mid y)$ and $(P|_{X=x})|_{Y=y} = P|_{XY=(x,y)}$. $\qquad\square$

Next, we define mutual information and, more generally, interaction information — also called co-information [8]. As we want to view interaction information as a "higher degree generalization" of entropy and treat both on an equal footing in Hu's theorem, we now change the notation: for any discrete random variables $X$, we set $I_1(X) \coloneqq H(X)$.

**Definition 2.7** (Mutual Information, Interaction Information)**.** *Let $q \in \mathbb{N}$ and assume that $I_{q-1}$ is already defined. Assume also that $Y_1, \ldots, Y_q$ are $q$ discrete random variables on $\Omega$. Then we define $I_q(Y_1; \ldots; Y_q) : \Delta_f(\Omega) \to \mathbb{R}$, the interaction information of degree $q$, as the function*

$$I_q(Y_1; \ldots; Y_q) \coloneqq I_{q-1}(Y_1; \ldots; Y_{q-1}) - Y_q.I_{q-1}(Y_1; \ldots; Y_{q-1}).$$

*$I_2$ is also called mutual information.*

---

[2]In the case that $E_X = E_Y = \mathbb{R}$, there is some ambiguity of notation, as the reader could understand $XY$ to be given by $(XY)(\omega) = X(\omega) \cdot Y(\omega)$. This definition plays a role in the *algebra of random variables* [51]. In our work, we instead *always* mean the Cartesian product.

**Remark 2.8.** *What we call* interaction information *is in the literature sometimes called* (higher / multivariate) mutual information. *In that case, the term* $J_q(Y_1; \ldots ; Y_q) := (-1)^{q+1} I_q(Y_1; \ldots ; Y_q)$ *is called interaction information, see for example [4].*

**Proposition 2.9.** *For all $q \geq 1$ and all discrete random variables $Y_1, \ldots, Y_q$, $I_q(Y_1; \ldots ; Y_q) : \Delta_f(\Omega) \to \mathbb{R}$ is a well-defined conditionable measurable function.*

*Proof.* $I_1(Y_1)$ is conditionable by Lemma 2.4. Assuming by induction that $I_{q-1}(Y_1; \ldots ; Y_{q-1})$ is well-defined and conditionable, we obtain the following: $Y_q.I_{q-1}(Y_1; \ldots ; Y_{q-1})$ is well-defined and conditionable by Proposition 2.6, part 2, and $I_q(Y_1; \ldots ; Y_q)$ is well-defined and conditionable by Proposition 2.6, part 3. □

## 2.2   Equivalence Classes of Random Variables

Assume all random variables are discrete. For two random variables $X$ and $Y$ on $\Omega$, we write $X \precsim Y$ if there is a function $f_{XY} : E_Y \to E_X$ such that $f_{XY} \circ Y = X$. The definition of $\precsim$ is equivalent to a preorder put forward in the context of conditional independence relations [18–20]. The latter work defines in their Section 6.2: $X \precsim Y$ if for all $\omega, \omega' \in \Omega$, the following implication holds true:

$$Y(\omega) = Y(\omega') \implies X(\omega) = X(\omega').$$

It is straightforward to show that this coincides with our own definition.

   Clearly, our relation is reflexive and transitive and thus a *preorder*. We define the equivalence relation $\sim$ by $X \sim Y$ iff $X \precsim Y$ and $Y \precsim X$. We denote by $[X]$ the equivalence class of $X$.

**Proposition 2.10** (See Proof 1)**.** *Let $Y \precsim X$ be two discrete random variables on $\Omega$. Then we have $I_1(Y) \leq I_1(X)$ as functions on $\Delta_f(\Omega)$, meaning that $I_1(Y; P) \leq I_1(X; P)$ for all $P \in \Delta_f(\Omega)$. In particular, if $X$ and $Y$ are equivalent (i.e., $X \precsim Y$ and $Y \precsim X$), then $I_1(X) = I_1(Y)$.*

**Proposition 2.11** (See Proof 2)**.** *Let $X \sim Y$ be two equivalent discrete random variables on $\Omega$. Then for all conditionable measurable functions $F : \Delta_f(\Omega) \to \mathbb{R}$ we have $X.F = Y.F$.*

**Proposition 2.12.** *Let $q \geq 1$ and $Y_1, \ldots, Y_q$ and $Z_1, \ldots, Z_q$ be two collections of discrete random variables on $\Omega$ such that $Y_k \sim Z_k$ for all $k = 1, \ldots, q$. Then $I_q(Y_1; \ldots ; Y_q) = I_q(Z_1; \ldots ; Z_q)$.*

*Proof.* For $q = 1$, this was shown in Proposition 2.10. The case $q > 1$ can be shown by induction using Definition 2.7 and Proposition 2.11. □

   This proposition shows that interaction information is naturally defined for collections of *equivalence classes of random variables*, instead of the random variables themselves.

## 2.3   Monoids of Random Variables

Again, assume all random variables to be discrete.

**Lemma 2.13.** *Let $X, Y, Z, X'$, and $Y'$ be random variables on $\Omega$. Let $\mathbf{1} : \Omega \to *$ be a trivial random variable, with $* = \{*\}$ a measurable space with one element. Then the following properties hold:*

   *0. If $X \sim X'$ and $Y \sim Y'$, then $XY \sim X'Y'$;*

   *1. $\mathbf{1}X \sim X \sim X\mathbf{1}$;*

   *2. $(XY)Z \sim X(YZ)$;*

   *3. $XY \sim YX$;*

   *4. $XX \sim X$.*

*Additionally, we have $X \precsim Y$ if and only if $XY \sim Y$.*

*Proof.* All of these statements are elementary and left to the reader to prove. □

Recall that a monoid is a tuple $(M, \cdot, \mathbf{1})$ with $M$ a set, $\cdot$ a multiplication, and $\mathbf{1} \in M$, such that $\mathbf{1}$ is neutral and the multiplication is associative. A monoid is commutative and idempotent if $m \cdot n = n \cdot m$ and $m \cdot m = m$ for all $m, n \in M$. Notice that rules 1 to 4 in the lemma resemble the properties of a commutative, idempotent monoid.

We remark that a commutative, idempotent monoid is algebraically the same as a join-semilattice (sometimes also called *bounded* join-semilattice), i.e., a partially ordered set which has a bottom element (corresponding to $\mathbf{1} \in M$) and binary joins (corresponding to the multiplication in a monoid). The partial order can be reconstructed from a commutative, idempotent monoid $M$ by writing $m \leq n$ if $m \cdot n = n$, which corresponds to the last statement in Lemma 2.13. The language of join-semilattices is, for example, used in the development of the theory of conditional independence [18].

**Proposition 2.14** (See Proof 3). *Let $\widehat{M} = \{X : \Omega \to E_X\}_X$ be a collection of random variables with the following two properties:*

*a) There is a random variable $\mathbf{1} : \Omega \to *$ in $\widehat{M}$ which has a one-point set $* = \{*\}$ as the target;*

*b) For every two $X, Y \in \widehat{M}$ there exists a $Z \in \widehat{M}$ such that $XY \sim Z$.*

*Let $[X]$ denote the equivalence class of $X$ under the relation $\sim$. Define $M := \widehat{M}/\sim$ as the collection of equivalence classes of elements in $\widehat{M}$. Define $[X] \cdot [Y] := [Z]$ for any $Z \in \widehat{M}$ with $XY \sim Z$. Then the triple $(M, \cdot, [\mathbf{1}])$ is a commutative, idempotent monoid.*

We note that the monoid of equivalence classes of discrete random variables is isomorphic to the monoid of partitions on $\Omega$, which is the formalization used in [5].

We can now study finite monoids of random variables as instances of the construction in Proposition 2.14. Let $n \geq 0$ be a natural number. Let $X_1, \ldots, X_n$ be fixed random variables on $\Omega$. Define $[n] := \{1, \ldots, n\}$. For arbitrary $I \subseteq [n]$, define $X_I := \prod_{i \in I} X_i$, the joint of the variables $X_i$ for $i \in I$. For $X_J$ and $X_I$, we have the equivalence $X_J X_I \sim X_{J \cup I}$. Note that $X_\emptyset : \Omega \to * = \{*\}$ is a trivial random variable.

**Definition 2.15** (Monoid of $X_1, \ldots, X_n$). *The monoid $\mathrm{M}(X_1, \ldots, X_n)$ of the variables $X_1, \ldots, X_n$ consists of the following data:*

1. *The elements are equivalence classes $[X_I]$ for $I \subseteq [n]$.*

2. *The multiplication is given by $[X_J] \cdot [X_I] = [X_{J \cup I}]$.*

3. *$\mathbf{1} := [X_\emptyset]$ is the neutral element with respect to multiplication.*

*This is a well-defined commutative, idempotent monoid by Proposition 2.14.*

Recall that an additive monoid action is a triple $(M, G, .)$, where $M$ is a monoid, $G$ is an abelian group, and $. : M \times G \to G$ is a function such that $\mathbf{1} \in M$ acts neutrally, with associativity (meaning $m.(n.g) = (m \cdot n).g$), and distributivity over addition in $G$.

**Proposition 2.16.** *Let $M$ be a monoid of (equivalence classes of) discrete random variables on $\Omega$ as in Proposition 2.14. Let $G = \mathrm{Meas}_{\mathrm{con}}\big(\Delta_f(\Omega), \mathbb{R}\big)$ be the group of conditionable measurable functions from $\Delta_f(\Omega)$ to $\mathbb{R}$. Then the averaged conditioning $. : M \times G \to G$ given by*

$$\big([X].F\big)(P) := \big(X.F\big)(P) = \sum_{x \in E_X} P_X(x) F(P|_{X=x})$$

*is a well-defined monoid action.*

*Proof.* The action is well-defined by Proposition 2.11 and Proposition 2.6, part 2. It is a monoid action by Proposition 2.6. $\square$

**Summary 2.17.** *We now summarize the abstract properties of interaction information $I_q$. Let $M$ be a commutative, idempotent monoid of discrete random variables as in Proposition 2.14. By abuse of notation, we do not distinguish between random variables and their equivalence classes,*

*i.e., we write $Y$ instead of $[Y]$. Denote by $G := \mathrm{Meas}_{\mathrm{con}}\left(\Delta_f(\Omega), \mathbb{R}\right)$ the group of conditionable measurable functions from $\Delta_f(\Omega)$ to $\mathbb{R}$. By Proposition 2.16, averaged conditioning $. : M \times G \to G$ is a well-defined monoid action.*

*By Proposition 2.12, we can view $I_q$ as a function $I_q : M^q \to G$ that is defined on tuples of* equivalence classes *of discrete random variables. By Proposition 2.5, entropy $I_1$ satisfies the equation*

$$I_1(XY) = I_1(X) + X.I_1(Y)$$

*for all $X, Y \in M$, where $X.I_1(Y)$ is the result of the action of $X \in M$ on $I_1(Y) \in G$ via averaged conditioning. Finally, by Definition 2.7, for all $q \geq 2$ and all $Y_1, \ldots, Y_q \in M$, one has*

$$I_q(Y_1; \ldots; Y_q) = I_{q-1}(Y_1; \ldots; Y_{q-1}) - Y_q.I_{q-1}(Y_1; \ldots; Y_{q-1}).$$

## 3  A Generalization of Hu's Theorem

In this section, we formulate and prove a generalization of Hu's theorem. Our treatment can be read mostly independently from the previous sections, but is motivated by Summary 2.17. First, in Section 3.1, we formulate the main result of this work, Theorem 3.2, together with its Corollary 3.3 that allows it to be applied to Kolmogorov complexity in Section 4 and the generalization error in Section 5. The formulation relies on a group-valued measure whose construction we motivate visually in Section 3.2. Afterwards, in Section 3.3, we deduce some general consequences on how (conditional) interaction terms of different degrees can be related to each other. The proofs can be found in Appendix C.

### 3.1  A Formulation of the Generalized Hu Theorem

Let $M$ be a commutative, idempotent monoid. We assume that $M$ is finitely generated, meaning there are elements $X_1, \ldots, X_n \in M$ such that all elements in $M$ can be written as arbitrary finite products of the elements $X_1, \ldots, X_n$. Since $M$ is commutative and idempotent, all elements in $M$ are of the form $X_I = \prod_{i \in I} X_i$ for some subset $I \subseteq [n] = \{1, \ldots, n\}$, and $X_I X_J := X_I \cdot X_J = X_{I \cup J}$. Additionally, fix an abelian group $G$ and an additive monoid action $. : M \times G \to G$.

For each $\emptyset \neq I \subseteq [n]$, we denote by $p_I$ an abstract atom. The only property we require of them is to be pairwise different, i.e., $p_I \neq p_J$ if $I \neq J$. Then, set $\widetilde{X}$ as the set of all these atoms:

$$\widetilde{X} := \left\{ p_I \mid \emptyset \neq I \subseteq [n] \right\}. \tag{4}$$

The atoms $p_I$ represent all smallest parts (the intersections of sets with indices in $I$ minus the sets with indices in $[n] \setminus I$) of a general Venn diagram for $n$ sets.

For $i \in [n]$, we denote by $\widetilde{X}_i := \left\{ p_I \in \widetilde{X} \mid i \in I \right\}$ a set which we can imagine to be depicted by a "disk" corresponding to the variable $X_i$, and we denote by $\widetilde{X}_I := \bigcup_{i \in I} \widetilde{X}_i$ the union of the "disks" corresponding to the joint variable $X_I$. Clearly, we have $\widetilde{X} = \widetilde{X}_{[n]}$. This is actually the simplest construction that leads to the $\widetilde{X}_i$ being in general position, as we have the following for all $\emptyset \neq I \subseteq [n]$:

$$\bigcap_{i \in I} \widetilde{X}_i \quad \setminus \bigcup_{j \in [n] \setminus I} \widetilde{X}_j = \{p_I\}. \tag{5}$$

We remark that $\widetilde{X}$ depends on $n$ and could therefore also be written as $\widetilde{X}(n)$. We will in most cases abstain from this to not overload the notation. In general, $\widetilde{X}$ has $2^n - 1$ elements. Therefore, for $n = 2$, $n = 3$ and $n = 4$, $\widetilde{X}$ has 3, 7, and 15 elements, respectively, see Figures 1, 2 and 3.

Remember that for a set $\Sigma$, $2^\Sigma$ is its powerset, i.e., the set of its subsets.

**Definition 3.1** ((*$G$-Valued) Measure*)**.** *Let $G$ be an abelian group and $\Sigma$ a set. A $G$-valued measure (on $\Sigma$) is a function $\mu : 2^\Sigma \to G$ with the property*

$$\mu(A_1 \cup A_2) = \mu(A_1) + \mu(A_2)$$

*for all disjoint $A_1, A_2 \subseteq \Sigma$.*

Figure 1: The generalized Hu theorem, visualized for a commutative, idempotent monoid $M$ generated by $X, Y$, and for $F_1$ and $F_2$. The measure $\mu$ turns sets into elements of the abelian group $G$ and disjoint unions into sums.

**Theorem 3.2** (Generalized Hu Theorem; See Section C.1 and Proof 4). *Let $M$ be a commutative, idempotent monoid generated by $X_1, \ldots, X_n$, $G$ an abelian group, $. : M \times G \to G$ an additive monoid action, and $\widetilde{X} = \widetilde{X}(n)$.*

1. *Assume $F_1 : M \to G$ is a function that satisfies the following chain rule: for all $X, Y \in M$, one has*

$$F_1(XY) = F_1(X) + X.F_1(Y). \tag{6}$$

*Construct $F_q : M^q \to G$ for $q \geq 2$ inductively by*

$$F_q(Y_1; \ldots; Y_q) := F_{q-1}(Y_1; \ldots; Y_{q-1}) - Y_q.F_{q-1}(Y_1; \ldots; Y_{q-1}) \tag{7}$$

*for all $Y_1, \ldots, Y_q \in M$.*

*Then there exists a $G$-valued measure $\mu : 2^{\widetilde{X}} \to G$ such that for all $q \geq 1$ and $J, L_1, \ldots, L_q \subseteq [n]$, the following identity holds:*

$$X_J.F_q(X_{L_1}; \ldots; X_{L_q}) = \mu\left( \bigcap_{k=1}^{q} \widetilde{X}_{L_k} \setminus \widetilde{X}_J \right). \tag{8}$$

*Concretely, one can define $\mu$ as the unique $G$-valued measure that is on individual atoms $p_I \in \widetilde{X}$ defined by*

$$\mu(p_I) := \sum_{\emptyset \neq K \supseteq I^c} (-1)^{|K|+|I|+1-n} \cdot F_1(X_K), \tag{9}$$

*where $I^c = [n] \setminus I$ is the complement of $I$ in $[n]$.[3]*

2. *Conversely, assume that $\mu : 2^{\widetilde{X}} \to G$ is a $G$-valued measure. Assume there is a sequence of functions $F_q : M^q \to G$ that satisfy Equation (8). Then $F_1$ satisfies Equation (6) and $F_q$ is related to $F_{q-1}$ as in Equation (7).*

*Sketch of Proof.* Part 1 can be shown as follows: When specializing Equation (8) to the case $X_J = \mathbf{1}$ and $q = 1$, one obtains

$$F_1(X_K) = \mu(\widetilde{X}_K) = \sum_{I : I \cap K \neq \emptyset} \mu(p_I),$$

---

[3]Alternatively, noting that $F_1(X_\emptyset) = 0$ and writing $K = K' \cup I^c$ for some unique $K' \subseteq I$, we can also write $\mu(p_I) = \sum_{K \subseteq I} (-1)^{|K|+1} \cdot F_1(X_K X_{I^c})$.

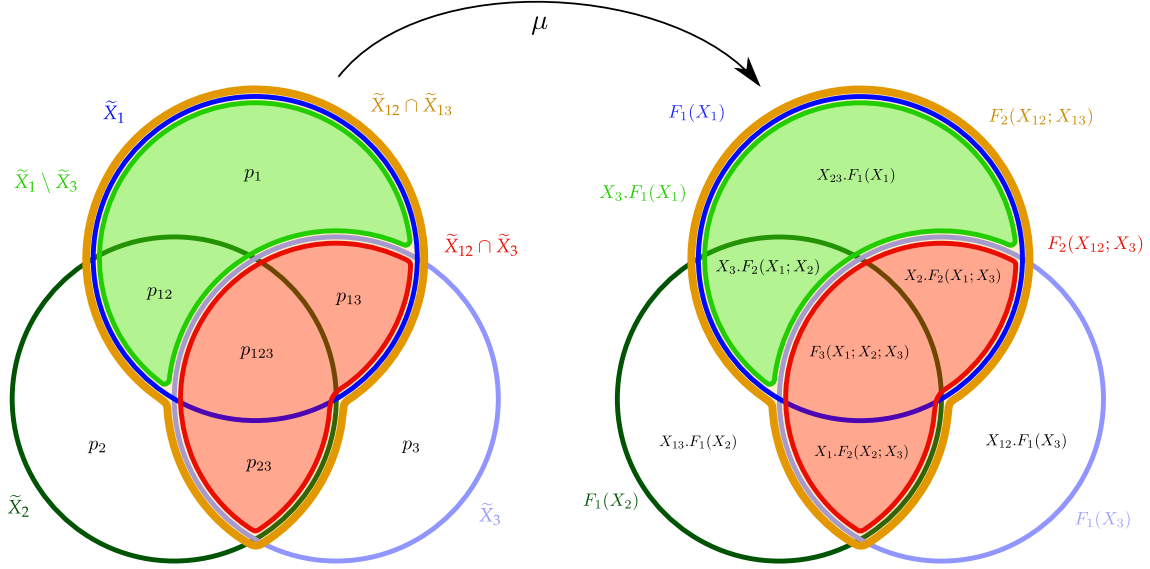Figure 2: A visualization of the generalized Hu theorem for a commutative, idempotent monoid generated by $X_1, X_2, X_3$. On the left-hand-side, three subsets of the abstract set $\widetilde{X}$ are emphasized, namely $\widetilde{X}_{12} \cap \widetilde{X}_{13}$, $\widetilde{X}_1 \setminus \widetilde{X}_3$, and $\widetilde{X}_{12} \cap \widetilde{X}_3$. On the right-hand-side, Equation (8) turns them into elements of the abelian group $G$, namely $F_2(X_{12}; X_{13})$, $X_3.F_1(X_1)$, and $F_2(X_{12}; X_3)$, respectively. Many decompositions of information functions into sums directly follow from the theorem by using that $\mu$ turns disjoint unions into sums, as exemplified by the equation $F_2(X_{12}; X_{13}) = X_3.F_1(X_1) + F_2(X_{12}; X_3)$.

which follows from the Möbius inversion formula on a poset [52, 3.7.1 Proposition] from Equation (9). The general formula for $q > 1$ then follows by induction using the properties of the monoid action. Part 2 follows by a direct computation.

More details can be found in Appendix C.1.                                                                    □

The following corollary will be applied to Kolmogorov complexity in Section 4 and the generalization error in machine learning in Section 5.

**Corollary 3.3** (Hu's Theorem for Two-Argument Functions; see Proof 5). *Let $M$ be a commutative, idempotent monoid generated by $X_1, \ldots, X_n$, $G$ an abelian group, and $\widetilde{X} = \widetilde{X}(n)$. Assume that $K_1 : M \times M \to G$ is a function satisfying the following chain rule:*

$$K_1(XY) = K_1(X) + K_1(Y \mid X), \tag{10}$$

*where we define $K_1(X) \coloneqq K_1(X \mid \mathbf{1})$ for all $X \in M$. Construct $K_q : M^q \times M \to G$ for $q \geq 2$ inductively by*

$$K_q\big(Y_1; \ldots; Y_q \mid Z\big) \coloneqq K_{q-1}\big(Y_1; \ldots; Y_{q-1} \mid Z\big) - K_{q-1}\big(Y_1; \ldots; Y_{q-1} \mid Y_q Z\big). \tag{11}$$

*Then there exists a $G$-valued measure $\mu : 2^{\widetilde{X}} \to G$ such that for all $L_1, \ldots, L_q, J \subseteq [n]$, the following identity holds:*

$$K_q(X_{L_1}; \ldots; X_{L_q} \mid X_J) = \mu\left(\bigcap_{k=1}^{q} \widetilde{X}_{L_k} \setminus \widetilde{X}_J\right). \tag{12}$$

*Concretely, one can define $\mu$ as the unique $G$-valued measure that is on individual atoms $p_I \in \widetilde{X}$ defined by*

$$\mu(p_I) \coloneqq \sum_{\emptyset \neq K \supseteq I^c} (-1)^{|K|+|I|+1-n} \cdot K_1(X_K), \tag{13}$$

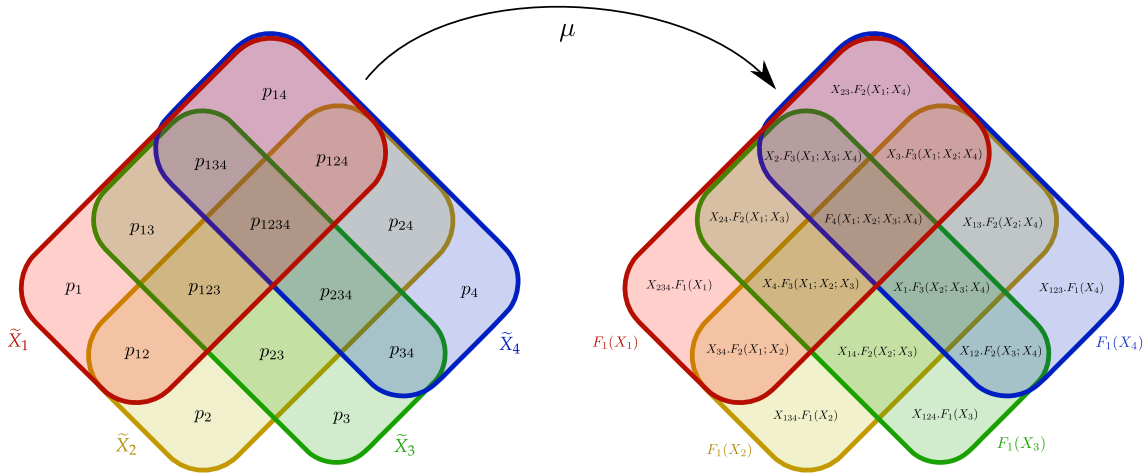*where $I^c = [n] \setminus I$ is the complement of $I$ in $[n]$.*

Figure 3: A visualization of the generalized Hu theorem for a commutative, idempotent monoid $M$ generated by $X_1, X_2, X_3, X_4$. To reduce clutter, we restrict to a visualization of the abstract sets $\widetilde{X}_i$ and the atoms $p_I$, as well as the corresponding information functions. On the right-hand-side, for computing $\mu(p_I)$ for the 15 atoms $p_I$, we use Lemma 3.4.

We conclude by discussing how Hu's theorem can be visualized, for which we will prove one further elementary lemma. For $I = \{i_1, \ldots, i_q\} \subseteq [n]$, set

$$\eta_I := X_{[n] \setminus I}.F_q(X_{i_1}; \ldots; X_{i_q}). \tag{14}$$

For the special case that $F_q = I_q$ is interaction information, these functions were discussed in [6] as generators of all information functions of the form $X_J.I_q(X_{L_1}; \ldots; X_{L_q})$. The following lemma gives an explanation for this: the functions $\eta_I$ generate the information measure (or, more generally: $G$-valued measure) $\mu$, which in turn generates all other information functions:

**Lemma 3.4.** *Let $\emptyset \neq I \subseteq [n]$ be arbitrary. Then $\eta_I = \mu(p_I)$.*

*Proof.* According to Equation (5), we have

$$\bigcap_{i \in I} \widetilde{X}_i \setminus \widetilde{X}_{[n] \setminus I} = \{p_I\}. \tag{15}$$

Thus, the lemma follows from Theorem 3.2.                                                     □

Thus, Theorem 3.2 can be visualized as follows: For each element $X_1, \ldots, X_n$, draw a disk $\widetilde{X}_i$ such that they intersect "in general position", meaning that all intersections of (part of) the disks are present. Assign the function $\eta_I$ to each atom $p_I$, as in the preceding lemma. Furthermore, assign subsets of $\widetilde{X}$ to information functions according to Equation (8). See Figures 1 and 3 for examples. In Figure 2, we exemplify how to use these diagrams to visually represent and prove identities of information functions. Note that in all figures, we write sets $I = \{i_1, \ldots, i_k\}$ for simplicity just as the sequence $i_1 i_2 \ldots i_k$.

## 3.2   Explicit Construction of the $G$-Valued Measure $\mu$

Assume all notation as in part 1 of Theorem 3.2. In this subsection, we explain how one could "guess" Equation (9) without knowledge of Möbius inversion theory. This section is meant as motivation, and other sections do not depend on it.

The high-level idea is the following: we have the sequence of functions $F_1, F_2, \ldots$, as our data to work with. We also know that $F_q$ is constructed from $F_{q-1}$ for all $q \geq 2$, which means that we should be able to express the measure $\mu$ in terms of $F_1$ alone. Additionally, we must have $F_1(X_K) = \mu(\widetilde{X}_K)$ in the end. Thus, our aim is to explain how, for arbitrary $\emptyset \neq I \subseteq [n]$, we can express $\mu(p_I)$ using only terms $\mu(\widetilde{X}_K)$ with $K \subseteq [n]$. This idea, while carried out differently, is also at the heart of the proof of the existence of information diagrams given in [63, 64].

We now look at some examples for $n$ and $I$ and derive $\mu(p_I)$ from the $\mu(\widetilde{X}_K)$. In the following visual computations, each Venn diagram always depicts the measure of the grey area. We frequently make use of the fact that $\mu$ is a $G$-valued measure. For $n = 1$ and $I = \{1\} = 1$,[4] we obtain:

$$\mu(p_1) = \boxed{p_1} = \mu(\widetilde{X}_1).$$

For $n = 2$ and $I = \{1\} = 1$, we have:

$$\mu(p_1) = \;\;\; = \;\;\; - \;\;\; = \mu(\widetilde{X}_{12}) - \mu(\widetilde{X}_2).$$

For $n = 2$ and $I = \{2\} = 2$, we get the same situation with 1 and 2 exchanged:

$$\mu(p_2) = \;\;\; = \;\;\; - \;\;\; = \mu(\widetilde{X}_{12}) - \mu(\widetilde{X}_1).$$

Next, we look at the case $n = 2$, $I = \{1, 2\} = 12$:

$$\mu(p_{12}) = \;\;\; = \;\;\; - \;\;\;$$

$$= \;\;\; + \;\;\; - \;\;\; = \mu(\widetilde{X}_1) + \mu(\widetilde{X}_2) - \mu(\widetilde{X}_{12}).$$

Finally, for $n = 3$ and $I = \{1, 2\} = 12$, we obtain:

$$\mu(p_{12}) = \;\;\; = \;\;\; - \;\;\;$$

$$= \;\;\; - \;\;\; - \;\;\; + \;\;\;$$

_____

[4]For simplicity, we write sets as a sequence of their elements.

$$= \quad \mu\big(\widetilde{X}_{23}\big) \quad - \quad \mu\big(\widetilde{X}_3\big) \quad - \quad \mu\big(\widetilde{X}_{123}\big) \quad + \quad \mu\big(\widetilde{X}_{13}\big).$$

In all cases, we managed to achieve our goal to only use terms of the form $\mu(\widetilde{X}_K)$. Additionally, a close look at the coefficients shows that these examples obey Equation (9), as desired.

## 3.3   General Consequences of the Explicit Construction of $\mu$

Assume the setting as in part 1 of Theorem 3.2, which we now consider proven. In this section, we consider general consequences of Hu's theorem that specifically use the explicit construction, Equation (9), of the $G$-valued measure $\mu : 2^{\widetilde{X}} \to G$. Corollary 3.5 explains how three different information functions can be expressed with respect to each other.

**Corollary 3.5** (See Proof 6). *Recall the functions $\eta_I$ from Equation (14). We obtain the following identities:*

1. *Let $1 \le q \le n$ and $\emptyset \ne I = \{i_1, \ldots, i_q\} \subseteq [n]$. Then*

$$\eta_I = \sum_{\emptyset \ne K \supseteq I^c} (-1)^{|K|+|I|+1-n} \cdot F_1(X_K).$$

2. *Let $K \subseteq [n]$ arbitrary. Then*

$$F_1(X_K) = \sum_{\substack{I \subseteq [n] \\ I \cap K \ne \emptyset}} \eta_I.$$

3. *Let $1 \le q \le n$ and $\emptyset \ne J = \{j_1, \ldots, j_q\} \subseteq [n]$ be arbitrary. Then*

$$F_q(X_{j_1}; \ldots; X_{j_q}) = \sum_{I \supseteq J} \eta_I.$$

4. *For $\emptyset \ne I \subseteq [n]$, we have*

$$\eta_I = \sum_{J \supseteq I} (-1)^{|J|-|I|} \cdot F_{|J|}(X_{j_1}; \ldots; X_{j_{|J|}}).$$

5. *Let $K \subseteq [n]$ arbitrary. Then one has*

$$F_1(X_K) = \sum_{\emptyset \ne J \subseteq K} (-1)^{|J|+1} \cdot F_{|J|}(X_{j_1}; \ldots; X_{j_{|J|}}).$$

6. *Let $1 \le q \le n$ and $\emptyset \ne J = \{j_1, \ldots, j_q\} \subseteq [n]$. Then one has*

$$F_q(X_{j_1}; \ldots; X_{j_q}) = \sum_{\emptyset \ne K \subseteq J} (-1)^{|K|+1} \cdot F_1(X_K).$$

## 4   Hu's Theorem for Kolmogorov Complexity

In this section, we establish the generalization of Hu's theorem for two-argument functions, Corollary 3.3, for different versions of Kolmogorov complexity. All of these versions satisfy a chain rule up to certain error terms. These can all be handled in our framework, but the most exact chain rule holds for *Chaitin's prefix-free Kolmogorov complexity*, on which we therefore focus our attention. Our main references are [15, 28, 34]. In this whole section, we work with the binary logarithm, which we denote by log, instead of the natural logarithm ln.

   This section is written with minimal prerequisites on the reader. We proceed as follows: in Section 4.1, we explain the preliminaries of prefix-free Kolmogorov complexity. Then in Section 4.2, we state the chain rule of Chaitin's prefix-free Kolmogorov complexity, which holds up to an additive constant. We reformulate this chain rule in Section 4.3 to satisfy the general assumptions

of Corollary 3.3 for two-argument functions. In Section 4.4, we then define interaction complexity analogously to interaction information, and make the resulting Hu theorem explicit.

Then in Section 4.5, we combine the two Hu theorems for interaction complexity and Shannon interaction information and show that expected interaction complexity is up to an error term equal to interaction information. This leads to the remarkable result that in all degrees, the "per-bit" expected interaction complexity equals interaction information for sequences of well-behaved probability measures on increasing sequence lengths.

Finally, the Sections 4.6 and 4.7 then summarize the resulting chain rules for standard prefix-free Kolmogorov complexity and plain Kolmogorov complexity, leaving more concrete interpretations of the resulting Hu theorems to future work.

Most proofs for this section can be found in Appendix D.

## 4.1   Preliminaries on Prefix-Free Kolmogorov Complexity

Let the *alphabet* be given by $\{0, 1\}$. The set of *binary strings* is given by

$$\{0, 1\}^* := \{\epsilon, 0, 1, 00, 01, 10, 11, 000, \dots\},$$

where $\epsilon$ is the empty string. The above lexicographical ordering defines a bijection $\mathbb{N} \to \{0, 1\}^*$ that we use to identify natural numbers with binary strings. Concretely, this identification maps

$$0 \mapsto \epsilon, \quad 1 \mapsto 0, \quad 2 \mapsto 1, \quad 3 \mapsto 00, \quad 4 \mapsto 01, \quad 5 \mapsto 10, \dots \tag{16}$$

We silently switch between viewing natural numbers as "just numbers" and viewing them as binary strings and vice versa.

If $x, y \in \{0, 1\}^*$ are two binary strings, then we can concatenate them to obtain a new binary string $xy \in \{0, 1\}^*$. A string $x \in \{0, 1\}^*$ is a proper prefix of the string $y \in \{0, 1\}^*$ if there is a string $z \in \{0, 1\}^*$ with $z \neq \epsilon$ such that $y = xz$. A set $\mathcal{A} \subseteq \{0, 1\}^*$ is called *prefix-free* if no element in $\mathcal{A}$ is a proper prefix of any other element in $\mathcal{A}$.

Let $\mathcal{X}$ and $\mathcal{Y}$ be sets. A *partial function* $f : \mathcal{X} \dashrightarrow \mathcal{Y}$ is a function $f : \mathcal{A} \to \mathcal{Y}$ defined on a subset $\mathcal{A} \subseteq \mathcal{X}$. A *decoder* for a set $\mathcal{X}$ is a partial function $D : \{0, 1\}^* \dashrightarrow \mathcal{X}$.[5] A decoder can be thought of as *decoding* the *code words* in $\{0, 1\}^*$ into *source words* in $\mathcal{X}$. A decoder $D : \{0, 1\}^* \dashrightarrow \mathcal{X}$ is called a *prefix-free decoder* if its domain $\mathcal{A} \subseteq \{0, 1\}^*$ is prefix-free.[6]

For a binary string $x$, $l(x)$ is defined to be its *length*, meaning the number of its symbols. Thus, for example, we have $l(\epsilon) = 0$ and $l(01) = 2$. Let $D : \{0, 1\}^* \dashrightarrow \mathcal{X}$ be a decoder. We define the length function $L_D : \mathcal{X} \to \mathbb{N} \cup \{\infty\}$ via

$$L_D(x) := \min \{l(y) \mid y \in \{0, 1\}^*, \ D(y) = x\},$$

which is $\infty$ if $D^{-1}(x) = \emptyset$.

In the following, we make use of the notion of a *Turing machine*. This can be imagined as a machine with very simple rules that implements an algorithm. We will not actually work with concrete definitions of Turing machines; instead, we let Church's Thesis 4.1 do the work, which we describe below — it will guarantee that any function that *intuitively* resembles an algorithm could equivalently be described by a Turing machine. Concrete definitions can be found in Chapter 1.7 of [34].

A *partial computable function* is any partial function $T : \{0, 1\}^* \dashrightarrow \{0, 1\}^*$ that can be computed by a Turing machine. The Turing machine *halts* on precisely the inputs on which $T$ is defined. We do not distinguish between Turing machines and the corresponding partial computable functions: If $T$ is a partial computable function, then we say that $T$ *is* a Turing machine. If $x \in \{0, 1\}^*$ is in the domain of the Turing machine $T$, we say that $T$ *halts* on $x$ and write $T(x) < \infty$. If $T$ does not halt on $x$, we sometimes write $T(x) = \infty$.

By the Church-Turing thesis, partial computable functions are precisely the partial functions for which there is an "algorithm in the intuitive sense" that computes the output for each input. We reproduce the formulation from [34]:

---

[5]Often, the word *code* is used instead of *decoder*. We find "decoder" less confusing.

[6]In the literature, this is often called a *prefix code*. We choose the name "prefix-free" as it avoids possible confusions.

**Thesis 4.1** (Church's Thesis). *The class of algorithmically computable partial functions (in the intuitive sense) coincides with the class of partial computable functions.*

We now define two prefix-free decoders for binary sequences. To do that, we first define the corresponding *encoders*: define the encoder $(\cdot)' : \{0,1\}^* \to \{0,1\}^*$ by

$$x' := 1^{l(l(x))}0l(x)x.^{[7]} \tag{17}$$

Note that the natural number $l(x)$ is viewed as a binary string using the identification in Equation (16).

The decoder corresponding to $(\cdot)'$ is a partial computable function $D' : \{0,1\}^* \dashrightarrow \{0,1\}^*$ that is only defined on inputs of the form $x'$. The underlying algorithm reads until the first 0 to know the length of the bitstring representing $l(x)$. Then it reads until the end of $l(x)$ to know the length of $x$. Subsequently, it can read until the end of $x$ to know $x$ itself, which it then outputs. This decoder is prefix-free: if $x'$ is a prefix of $y'$, then $l(x) = l(y)$ and $x$ is a prefix of $y$, from which $x = y$ and thus $x' = y'$ follows.

Let a pairing function $\{0,1\}^* \times \{0,1\}^* \to \{0,1\}^*$ be given by

$$(x,y) \mapsto x'y.$$

Note that we can algorithmically recover both $x$ and $y$ from $x'y$: reading the string $x'y$ from the left, the algorithm first recovers $l(x)$ and then $x$, after which the rest of the string automatically is $y$.

A Turing machine $T : \{0,1\}^* \dashrightarrow \{0,1\}^*$ is called a *prefix-free machine* if it is a prefix-free decoder. The input is then imagined to be a code word encoding the output string. There is a bijective, computable enumeration, called *standard enumeration*, $T_1, T_2, T_3, \ldots$, of all prefix-free machines ([34], Section 3.1). Computable here means the following: if we would encode the set of *rules* of any Turing machine as a binary sequence, then the map from natural numbers to binary sequences corresponding to the standard enumeration is itself computable.

A Turing machine $T : \{0,1\}^* \dashrightarrow \{0,1\}^*$ is called a *conditional Turing machine* if for all $x$ such that $T$ halts on $x$ we have $x = y'p$ for some elements $y, p \in \{0,1\}^*$; $p$ is then called the *program*, and $y$ the *input*. A *universal* conditional prefix-free machine is a conditional prefix-free machine $U : \{0,1\}^* \dashrightarrow \{0,1\}^*$ such that for all $i \in \mathbb{N}$ and $y, p \in \{0,1\}^*$, we have $U(y'i'p) = T_i(y'p)$, and $U$ does not halt on inputs of any other form. Here, again, $i$ is viewed as a binary string via Equation (16). One can show that such universal conditional prefix-free machines indeed do exist ([34], Theorem 3.1.1).

For the rest of this article, let $U$ be a fixed universal conditional prefix-free machine.

**Definition 4.2** (Prefix-Free Kolmogorov Complexity). *The* conditional prefix-free Kolmogorov complexity *is the function $K : \{0,1\}^* \times \{0,1\}^* \to \mathbb{N}$ given by*

$$K(x \mid y) := \min \left\{ l(p) \;\middle|\; p \in \{0,1\}^*, \; U(y'p) = x \right\}$$

$$= \min \left\{ l(i') + l(q) \;\middle|\; i \in \mathbb{N}, \; q \in \{0,1\}^*, \; U(y'i'q) = x \right\}$$

$$= \min \left\{ l(i') + l(q) \;\middle|\; i \in \mathbb{N}, \; q \in \{0,1\}^*, \; T_i(y'q) = x \right\}$$

$$< \infty.$$

*We define the* non-conditional *prefix-free Kolmogorov complexity by $K : \{0,1\}^* \to \mathbb{N}$, $K(x) := K(x \mid \epsilon)$. As $\epsilon' = 1^{l(l(\epsilon))}0l(\epsilon) = 0$,[8] we obtain*

$$K(x) = \min \left\{ l(p) \mid U(0p) = x \right\}.$$

*Here, the 0 can be thought of as simply signaling that there is no input, while each "actual" input starts with a 1 due to the definition of $y'$.*

---

[7]In the literature, this is viewed as a code for the *natural numbers* instead of $\{0,1\}^*$. But both viewpoints are equivalent due to the bijection $\mathbb{N} \cong \{0,1\}^*$.

[8]Here, we used $l(\epsilon) = 0$, which is a *natural number* corresponding to the *string $\epsilon$* that is plucked back into the formula.

**Definition 4.3** (Joint Conditional Prefix-Free Kolmogorov Complexity). *Define* $\mathrm{Concat} : (\{0,1\}^*)^n \to \{0,1\}^*$ *by* $\mathrm{Concat}(x_1, \ldots, x_n) \coloneqq x_1' \cdots x_{n-1}' x_n$. *For* $x_1, \ldots, x_n \in \{0,1\}^*$ *and* $y_1, \ldots, y_m \in \{0,1\}^*$, *we define the* (joint conditional) *prefix-free Kolmogorov complexity by*

$$K\big(x_1, \ldots, x_n \mid y_1, \ldots, y_m\big) \coloneqq K\big( \mathrm{Concat}(x_1, \ldots, x_n) \mid \mathrm{Concat}(Y_1, \ldots, y_m)\big).$$

*We then simply set* $K\big(x_1, \ldots, x_n\big) \coloneqq K\big(x_1, \ldots, x_n \mid \epsilon\big)$.

## 4.2　The Chain Rule for Chaitin's Prefix-Free Kolmogorov Complexity

Let $f, g : \mathcal{X} \to \mathbb{R}$ be two functions on a set $\mathcal{X}$. We adopt the following notation from [28]: $f \overset{+}{<} g$ means that there is a constant $c \geq 0$ such that $f(x) < g(x) + c$ for all $x \in \mathcal{X}$. We write $f \overset{+}{>} g$ if $g \overset{+}{<} f$. Finally, we write $f \overset{+}{=} g$ if $f \overset{+}{<} g$ and $f \overset{+}{>} g$, which means that there is a constant $c \geq 0$ such that $\big|f(x) - g(x)\big| < c$ for all $x \in \mathcal{X}$. If we want to emphasize the inputs, we may, for example, also write $f(x) \overset{+}{=} g(x)$.

Let $x \in \{0,1\}^*$ be arbitrary and $K(x)$ its prefix-free Kolmogorov complexity. Let $x^* \in \{0,1\}^*$ be chosen as follows: we look at all $y \in \{0,1\}^*$ of length $l(y) = K(x)$ such that $U(0y) = x$. Among those, we look at all $y$ such that $U$ computes $x$ on input $0y$ with the smallest number of computation steps. And finally, among those, we define $x^*$ to be the lexicographically first string. Based on this, Chaitin's prefix-free Kolmogorov complexity is given by

$$Kc : \{0,1\}^* \times \{0,1\}^* \to \mathbb{R}, \quad Kc(x \mid y) \coloneqq K(x \mid y^*)$$

and $Kc(x) \coloneqq Kc(x \mid \epsilon)$.

Clearly, there is a program that, on input $x'K(x)$, outputs $x^*$ — we simply run $U(0y)$ for all programs $y$ of length $K(x)$ in parallel, and the one that outputs $x$ the fastest and is lexicographically first among those is the output $x^*$. Vice versa, given $x^*$, one can compute $x'K(x)$ by simply computing $U(0x^*)'l(x^*)$. In this sense, $x^*$ and $x'K(x)$ can be said to "contain the same information". In the literature, Chaitin's prefix-free Kolmogorov complexity is, for this reason, also often defined by $Kc(x \mid y) \coloneqq K(x \mid y, K(y))$.

The following result might have for the first time been written down in [26], and was attributed therein to Leonid Levin.

**Theorem 4.4** (Chain Rule for Chaitin's Prefix-Free Kolmogorov Complexity). *The following identity holds:*

$$Kc(x, y) \overset{\pm}{=} Kc(x) + Kc(y \mid x). \tag{18}$$

*Here, both sides are viewed as functions* $(\{0,1\}^*)^2 \to \mathbb{R}$ *that map inputs of the form* $(x, y)$.

*Proof.* See [34], Theorem 3.8.1 for the proof of the inequality $Kc(x, y) \overset{+}{<} Kc(x) + Kc(y \mid x)$. The proof of the other direction, namely $Kc(y \mid x) \overset{+}{<} K(x, y) - K(x)$, in [34] seems incorrect to us, as it only seems to show that the constant is independent of $x$ and not of $y$. See the proof in [15] for that direction. □

## 4.3　A Reformulation of the Chain Rule in Terms of Our General Framework

Our goal is to express the result, Equation (18), in terms of the assumptions of Corollary 3.3. To do this, we need to find a framework under which the chain rule becomes *exact* instead of correct up to a constant, and in which the inputs come from a monoid. We will solve this by identifying functions whose difference is bounded by a constant.

For $n \geq 0$ any fixed natural number, we define $\mathrm{Maps}\big((\{0,1\}^*)^n, \mathbb{R}\big)$ as the abelian group of functions from $(\{0,1\}^*)^n$ to $\mathbb{R}$. We define the equivalence relation $\sim_{Kc}$ on $\mathrm{Maps}\big((\{0,1\}^*)^n, \mathbb{R}\big)$ by

$$F \sim_{Kc} H \quad :\Longleftrightarrow F \overset{\pm}{=} H.$$

The reason we put $Kc$ in the subscript of $\sim_{Kc}$ is that later, we will investigate different equivalence relations $\sim_K$ and $\sim_C$ for prefix-free and plain Kolmogorov complexity. Note that the functions $F$

with $F \sim_{Kc} 0$, i.e., $F \overset{\pm}{=} 0$, form a subgroup of $\text{Maps}\big(({\{0,1\}^*})^n, \mathbb{R}\big)$. Consequently, we obtain an abelian group $\text{Maps}\big(({\{0,1\}^*})^n, \mathbb{R}\big)/\sim_{Kc}$ with elements written as $[F]_{Kc}$.

Now, let the variables $X_1, \ldots, X_n$ be defined as the following projections:

$$X_i : ({\{0,1\}^*})^n \to {\{0,1\}^*}, \quad \boldsymbol{x} = (x_1, \ldots, x_n) \mapsto x_i.$$

Then, for any $i_1, \ldots, i_k \in [n]$, we can form the product variable $X_{i_1} \cdots X_{i_k}$:

$$X_{i_1} \cdots X_{i_k} : ({\{0,1\}^*})^n \to ({\{0,1\}^*})^k, \quad \boldsymbol{x} = (x_1, \ldots, x_n) \mapsto (x_{i_1}, \ldots, x_{i_k}).$$

These strings of projections form the elements of the monoid $\widetilde{M} = \{X_1, \ldots, X_n\}^*$, with multiplication simply given by concatenation. Then from $Kc : \{0,1\}^* \times \{0,1\}^* \to \mathbb{R}$, we can define the function

$$[Kc]_{Kc} : \ \widetilde{M} \times \widetilde{M} \to \text{Maps}\big(({\{0,1\}^*})^n, \mathbb{R}\big)/\sim_{Kc},$$
$$(Y, Z) \mapsto [Kc(Y \mid Z)]_{Kc},$$

with $Kc(Y \mid Z)$ simply being the function that inserts tuples from $({\{0,1\}^*})^n$ into the variables $Y$ and $Z$:

$$Kc(Y \mid Z) : ({\{0,1\}^*})^n \to \mathbb{R}, \quad \boldsymbol{x} \mapsto Kc(Y(\boldsymbol{x}) \mid Z(\boldsymbol{x})).$$

Similarly as before, one can then define $Kc(Y) : ({\{0,1\}^*})^n \to \mathbb{R}$ by $Kc(Y) \coloneqq Kc(Y \mid \epsilon)$ with $\epsilon \in \widetilde{M}$ being the empty string of variables. In the same way, $[Kc]_{Kc}(Y) \coloneqq [Kc]_{Kc}(Y \mid \epsilon) = [Kc(Y)]_{Kc}$. Since $\epsilon(\boldsymbol{x}) = \epsilon$ for all $\boldsymbol{x} \in ({\{0,1\}^*})^n$, these definitions are compatible with the earlier definition $Kc(x) \coloneqq Kc(x \mid \epsilon)$ for $x \in \{0,1\}^*$: we have $\big(Kc(Y)\big)(\boldsymbol{x}) = Kc\big(Y(\boldsymbol{x})\big)$.

**Proposition 4.5** (See Proof 7)**.** *For arbitrary $Y, Z \in \widetilde{M}$, we have the* exact *equality*

$$[Kc]_{Kc}(YZ) = [Kc]_{Kc}(Y) + [Kc]_{Kc}(Z \mid Y) \tag{19}$$

*of elements in* $\text{Maps}\big(({\{0,1\}^*})^n, \mathbb{R}\big)/\sim_{Kc}$.

To obtain a commutative, idempotent monoid, we show that we can permute and "reduce" the elements in $\widetilde{M}$ without affecting the resulting functions in $\text{Maps}\big(({\{0,1\}^*})^n, \mathbb{R}\big)/\sim_{Kc}$: for arbitrary $Y = X_{i_1} \cdots X_{i_k} \in \widetilde{M}$ we define the reduction $\overline{Y} \in \widetilde{M}$ by

$$\overline{Y} \coloneqq X_I \coloneqq \prod_{i \in I} X_i, \quad \text{with } I \coloneqq \Big\{ i \in [n] \ \Big| \ \exists s \in [k] : \ i_s = i \Big\}. \tag{20}$$

Here, the factors $X_i$ with $i \in I$ are assumed to appear in increasing order of the index $i$.

**Lemma 4.6** (See Proof 8)**.** *For all $Y, Z \in \widetilde{M}$, we have the equality*

$$[Kc]_{Kc}\big(Y \mid Z\big) = [Kc]_{Kc}\big(\overline{Y} \mid \overline{Z}\big)$$

*in* $\text{Maps}\big(({\{0,1\}^*})^n, \mathbb{R}\big)/\sim_{Kc}$.

Now, define the equivalence relation $\sim$ on $\widetilde{M}$ by $Y \sim Z$ if $\overline{Y} = \overline{Z}$, with $\overline{(\cdot)} : \widetilde{M} \to \widetilde{M}$ defined as in Equation (20). We define $M \coloneqq \widetilde{M}/\sim$. Each element $[Y] \in M$ is then represented by $\overline{Y}$ since $\overline{\overline{Y}} = \overline{Y}$; it is of the form $\overline{Y} = X_I$ for some $I \subseteq [n]$. Additionally, if $I \neq J$, then obviously we have $X_I \not\sim X_J$, and consequently, there is a one-to-one correspondence between representatives of the form $X_I$ and elements in $M$. Therefore, we can write elements in $M$ for convenience, and by abuse of notation, simply as $[Y] = X_I$. We then define the multiplication in $M$ by $[Y] \cdot [Z] \coloneqq [YZ]$, which in the new notation can be written as $X_I \cdot X_J = X_{I \cup J}$ and thus makes $M$ a well-defined commutative, idempotent monoid generated by $X_1, \ldots, X_n$. We define, by abuse of notation, $[Kc]_{Kc} : M \times M \to \text{Maps}\big(({\{0,1\}^*})^n, \mathbb{R}\big)/\sim_{Kc}$ in the obvious way on representatives, which is well-defined by Lemma 4.6. Overall, we obtain by Corollary 3.3 a Hu theorem for Chaitin's prefix-free Kolmogorov complexity, which we next explain in more detail.

## 4.4　Hu's Theorem for Chaitin's Prefix-Free Kolmogorov Complexity

We now deduce a Hu theorem for Chaitin's prefix-free Kolmogorov complexity. We formulate it without the abstraction of equivalence classes from the previous subsection (which is mainly important for the proof), with the goal to obtain an intrinsically more interesting version. For formulating the result, we first name the higher-degree terms analogously to the interaction information from Definition 2.7:

**Definition 4.7** (Interaction Complexity)**.** *Define $Kc_1 \coloneqq Kc : \{0,1\}^* \times \{0,1\}^* \to \mathbb{R}$ and $Kc_q :$ $(\{0,1\}^*)^q \times \{0,1\}^* \to \mathbb{R}$ inductively by*

$$Kc_q(y_1;\dots;y_q \mid z) \coloneqq Kc_{q-1}(y_1;\dots;y_{q-1}|z) - Kc_{q-1}(y_1;\dots;y_{q-1} \mid y_q,z).$$

*We call $Kc_q$ the* interaction complexity of degree $q$.

For example, $Kc_2(x;y) = Kc_1(x) - Kc_1(x \mid y)$ measures the reduction of the encoding length of $x$ when having access to $y$. E.g., if $x$ is thought of as "data" and $y$ thought of as a "theory", then $Kc_2(x;y)$ measures the extent to which $y$ helps in compressing $x$. See also the last paragraph in Section 6.3 for more interpretation of the potential meaning of these quantities. The interpretation of higher-order terms is future work.

For $Y_1,\dots,Y_q, Z \in \widetilde{M} = \{X_1,\dots,X_n\}^*$, define $Kc_q(Y_1;\dots;Y_q \mid Z) \in \mathrm{Maps}\left((\{0,1\}^*)^n, \mathbb{R}\right)$ by

$$Kc_q(Y_1;\dots;Y_q \mid Z): \quad \boldsymbol{x} \mapsto Kc_q\big(Y_1(\boldsymbol{x});\dots;Y_q(\boldsymbol{x}) \mid Z(\boldsymbol{x})\big).$$

One can easily inductively show that

$$Kc_q(Y_1;\dots;Y_q \mid Z) \overset{\pm}{=} Kc_{q-1}(Y_1;\dots;Y_{q-1} \mid Z) - Kc_{q-1}(Y_1;\dots;Y_{q-1} \mid Y_q Z). \tag{21}$$

The full proof of the following theorem can be found in Appendix D, Proof 9. The main ingredient is the chain rule, Proposition 4.5, together with Corollary 3.3.

**Theorem 4.8** (See Proof 9)**.** *Let $\widetilde{X} = \widetilde{X}(n)$. There exists a measure $\mu : 2^{\widetilde{X}} \to \mathrm{Maps}\left((\{0,1\}^*)^n, \mathbb{R}\right)$ such that for all $L_1,\dots,L_q, J \subseteq [n]$, the relation*

$$Kc_q\big(X_{L_1};\dots;X_{L_q} \mid X_J\big) \overset{\pm}{=} \mu\left(\bigcap_{k=1}^q \widetilde{X}_{L_k} \setminus \widetilde{X}_J\right) \tag{22}$$

*of functions $(\{0,1\}^*)^n \to \mathbb{R}$ holds. Concretely, $\mu$ can be defined as the unique measure that is on individual atoms $p_I \in \widetilde{X}$ defined by*

$$\mu(p_I) \coloneqq \sum_{\emptyset \neq K \supseteq I^c} (-1)^{|K|+|I|+1-n} \cdot Kc_1(X_K), \tag{23}$$

*where $I^c = [n] \setminus I$ is the complement of $I$ in $[n]$.*

**Remark 4.9.** *In Theorem 4.8, the equality holds up to a constant independent of the input in $(\{0,1\}^*)^n$. However, there is a dependence on $q$, the degree, and $n$, the number of generating variables. We now briefly analyze this.*

*For analyzing the dependence on $q$, we note that the inductive step of the proof of the generalized Hu Theorem 3.2 uses the theorem for degree $q-1$ twice. That means that the number of comparisons doubles in each degree, leading to a dependence of $q$ of the form $O(2^q)$. Can one do better than this? One idea might be to not define $Kc_q$ inductively, but with an inclusion-exclusion–type formula motivated by Corollary 3.5, part 6. One sensible definition is the following:*

$$Kc_q(y_1;\dots;y_q \mid z) \coloneqq \sum_{K \subseteq [q]} (-1)^{|K|+1} \cdot Kc_1(\boldsymbol{y}_K z)$$

*with*

$$\boldsymbol{y}_K \coloneqq \prod_{k \in K} y'_k. \tag{24}$$

*However, this now leads to $2^q$ summands, which one would, for a proof of Hu's theorem, individually compare with the evaluation of $\mu$ on a "disk" in $\widetilde{X} = \widetilde{X}(n)$. As in the general definition Equation (24), the order of the factors in $\boldsymbol{y}_K$ does not follow the ordering of the generators $x_1, \ldots, x_n$, we expect there a reordering of the factors to be necessary for the comparison. This has each time a cost of $O(1)$, thus again leading to a dependence of the form $O(2^q)$. We currently do not see a way to improve this.*

*Now, for each of the $2^q$ comparisons, we would like to know the dependence on $n$. One possible algorithm for bringing $\boldsymbol{y}_K z$ "in order" works as follows: assuming that all of $y_k$, $k \in K$, and $z$ are given by a permutation (with omissions) of $x_1, \ldots, x_n$, then we have to specify $q + 1$ permutations, which each involves to specify the position of $n$ elements. The position is one of $1, \ldots, n$ plus "omission", which together has a cost of $\log(n + 1)$. Overall, this leads to a dependence on $n$ of $O\big((q + 1) \cdot n \cdot \log(n + 1)\big)$.*

*Overall, the dependence on $q$ and $n$ together is thus $O\big(2^q \cdot (q + 1) \cdot n \cdot \log(n + 1)\big)$.*



Figure 4: A visualization of Hu's theorem for Kolmogorov complexity for three variables $X, Y, Z$. On the left-hand-side, three subsets of the abstract set $\widetilde{XYZ}$ are emphasized, namely $\widetilde{XY} \cap \widetilde{XZ}$, $\widetilde{X} \setminus \widetilde{Z}$, and $\widetilde{XY} \cap \widetilde{Z}$. On the right-hand-side, Equation (22) turns them up to a constant error into the Kolmogorov complexity terms $Kc_2(XY; XZ)$, $Kc(X \mid Z)$, and $Kc_2(XY; Z)$, respectively. Many decompositions of complexity terms into sums directly follow from the theorem by using that $\mu$ turns disjoint unions into sums, as exemplified by the equation $Kc_2(XY; XZ) \stackrel{\pm}{=} Kc(X \mid Z) + Kc_2(XY; Z)$.

As an Example, we recreate Figure 2 for the case of Kolmogorov complexity in Figure 4. We can also translate back from the notation with variables to the more familiar notation in which elements of $\{0, 1\}^*$ are inserted in the formulas. If we do this, then the example equation from Figure 4 becomes

$$Kc_2(x, y; x, z) \stackrel{\pm}{=} Kc(x \mid z) + Kc_2(x, y; z),$$

where both sides are viewed as functions $(\{0, 1\}^*)^3 \to \mathbb{R}$.

## 4.5 Expected Interaction Complexity is Interaction Information

Recall Definition 2.7 of the interaction information of $q$ discrete random variables $Y_1, \ldots, Y_q$, denoted $I_q(Y_1; \ldots; Y_q)$. Additionally, recall that for another discrete random variable $Z$ defined on the same sample space, we can define the averaged conditioning $Z.I_q(Y_1; \ldots; Y_q)$, see Definition 2.3, which is again an information function. Its evaluation on a probability measure $P$ on the sample space is denoted $Z.I_q(Y_1; \ldots; Y_q; P)$.

In this section, we want to establish a relationship between interaction information of random variables defined on $(\{0, 1\}^*)^n$ with values in $(\{0, 1\}^*)^k$ for some $k$ on the one hand, and the

expectation of interaction complexity as defined in Definition 4.7 on the other hand. The deviation from an equality between interaction information and interaction complexity will be quantified by the Kolmogorov complexity of probability mass functions.

For this aim, we first need to interpret outputs of Turing machines as rational numbers: If $T$ is a Turing machine and $T(x) = m'n$ for some $m, n \in \{0,1\}^*$, then interpret $m, n$ as natural numbers via the identification map in Equation (16), and consequently $m'n$ as the rational number $m/n$, see also Li and Vitányi [34], Section 1.7.3. Interpret the output as 0 if it is not of the form $m'n$.

**Definition 4.10** (Kolmogorov Complexity of Probability Mass Functions). *Let* $P : (\{0,1\}^*)^n \to \mathbb{R}$ *be a probability mass function. Its* Kolmogorov complexity *is defined by*

$$K(P) := \min_{p \in \{0,1\}^*} \left\{ l(p) \ \bigg| \ \forall q \in \mathbb{N}, \ \ \forall \boldsymbol{x} \in (\{0,1\}^*)^n : \ \ \left| T_p(\boldsymbol{x}'q) - P(\boldsymbol{x}) \right| \leq 1/q \right\},$$

*where $T_p$ is the p'th prefix-free Turing machine.*

**Definition 4.11** (Computability of Probability Mass Functions). *A probability mass function* $P : (\{0,1\}^*)^n \to \mathbb{R}$ *is called* computable *if $K(P) < \infty$.*

In other words, a probability mass function $P$ is computable if there exists a prefix-free Turing machine $T_p$ that can, for all natural numbers $q$, approximate $P$ up to precision $1/q$.

We now unify the viewpoint of the variables $X_i$ as "placeholders" with the viewpoint that they are random variables: remember that the $X_i : (\{0,1\}^*)^n \to \{0,1\}^*$ are given by projections: $X_i(\boldsymbol{x}) = x_i$. They form the monoid $\widetilde{M} = \{X_1, \ldots, X_n\}^*$, with multiplication given by concatenation. Furthermore, we defined an equivalence relation $\sim$ with $Y \sim Z$ if $\overline{Y} = \overline{Z}$.

Now, interpret $(\{0,1\}^*)^n$ as a discrete sample space. Then the strings in $Y \in \widetilde{M}$ can be interpreted as random variables on $(\{0,1\}^*)^n$ with values in $(\{0,1\}^*)^k$ for some $k$. The concatenation of these strings is identical to the product of random variables defined in Equation (3). Now, remember that in Section 2.2 we also defined an equivalence relation for random variables, which we now call $\sim_r$ to distinguish it from $\sim$. For $Y : (\{0,1\}^*)^n \to (\{0,1\}^*)^{k_y}$ and $Z : (\{0,1\}^*)^n \to (\{0,1\}^*)^{k_z}$, we have $Y \sim_r Z$ if there exist functions $f_{ZY} : (\{0,1\}^*)^{k_y} \to (\{0,1\}^*)^{k_z}$ and $f_{YZ} : (\{0,1\}^*)^{k_z} \to (\{0,1\}^*)^{k_y}$ such that $f_{ZY} \circ Y = Z$ and $f_{YZ} \circ Z = Y$.

**Lemma 4.12** (See Proof 10). *For all $Y, Z \in \widetilde{M}$, we have*

$$Y \sim Z \quad \Longleftrightarrow \quad Y \sim_r Z.$$

*That is, the equivalence relations $\sim$ and $\sim_r$ are identical.*

This shows that the commutative, idempotent monoids $M = \{X_1, \ldots, X_n\}^*/\sim$ and $\mathrm{M}(X_1, \ldots, X_n)$ from Definition 2.15 are the same. The only difference is simply that the neutral element in $\{X_1, \ldots, X_n\}^*/\sim$ was denoted $\epsilon$, whereas the one of $\mathrm{M}(X_1, \ldots, X_n)$ was denoted $\mathbf{1}$. We denote both monoids simply by $M$ from now on. For the following theorems, recall that a probability measure $P \in \Delta(\Omega)$ has a Shannon entropy $I_1(P)$ which equals $I_1(\mathrm{id}_\Omega; P)$, see Definitions 2.1, 2.2. Our aim is to generalize the following theorem:

**Theorem 4.13** ([34], Theorem 8.1.1). *We have*

$$0 \leq \left( \sum_{\boldsymbol{x} \in (\{0,1\}^*)^n} P(\boldsymbol{x}) Kc(\boldsymbol{x}) - I_1(P) \right) \overset{+}{<} K(P),$$

*where both sides are viewed as functions in computable probability measures $P : (\{0,1\}^*)^n \to \mathbb{R}$ with finite entropy $I_1(P) < \infty$. That is, up to $K(P) + c$ for some constant $c$ independent of $P$, entropy equals expected Kolmogorov complexity.*

In the following theorem, if we write $f = g + O(h)$ for functions $f, g, h : \mathcal{X} \to \mathbb{R}$, we mean that there exists a $c \geq 0$ such that $|f(x) - g(x)| < c \cdot h(x)$ for all $x \in \mathcal{X}$. This is in contrast to our use of that notation in the following parts, Sections 4.6 and 4.7, where the inequality only needs to hold starting from some threshold value $x_0 \in \mathcal{X}$.

We prove the result in Appendix D, Proof 11, with the main ingredients being Hu's theorems for both Shannon entropy — which follows using Summary 2.17 from Theorem 3.2 — and Chaitin's prefix-free Kolmogorov complexity (Theorem 4.8). Both together allow a reduction to the well-known special case, Theorem 4.13.

**Theorem 4.14** (See Proof 11). *Let $X_1, \ldots, X_n : (\{0,1\}^*)^n \to \{0,1\}^*$ be the (random) variables given by $X_i(\boldsymbol{x}) = x_i$. Let $M = \{X_1, \ldots, X_n\}^*/\sim \; = \; \mathrm{M}(X_1, \ldots, X_n)$ be the idempotent, commutative monoid generated by $X_1, \ldots, X_n$, with elements written as $X_I$ for $I \subseteq [n]$. Then for all $q \geq 1$ and $Y_1, \ldots, Y_q, Z \in M$, the following relation holds:*

$$\sum_{\boldsymbol{x} \in (\{0,1\}^*)^n} P(\boldsymbol{x}) \cdot \big(Kc_q(Y_1; \ldots; Y_q \mid Z)\big)(\boldsymbol{x}) = Z.I_q(Y_1; \ldots; Y_q; P) + O\big(K(P)\big), \qquad (25)$$

*where both sides are viewed as functions in computable probability mass functions $P : (\{0,1\}^*)^n \to \mathbb{R}$ with finite entropy $I_1(P) < \infty$.*

**Remark 4.15.** *Similar to Remark 4.9, one can also for this theorem wonder about the dependence on $n$ and $q$. A similar analysis shows that our techniques lead to a dependence of the form*

$$O\Big(2^q\big((q+1)n \log(n+1) + K(P)\big)\Big).$$

**Corollary 4.16.** *Assume that $(P_m)_{m \in \mathbb{N}}$ is a sequence of computable probability mass functions $P_m : (\{0,1\}^*)^n \to \mathbb{R}$ with finite entropy. Additionally, we make the following two assumptions:*

- *$P_m$ has all its probability mass on elements $\boldsymbol{x} = (x_1, \ldots, x_n) \in (\{0,1\}^*)^n$ with sequence lengths $l(x_i) = m$ for all $i \in [n]$;*

- *$K(P_m)$ grows sublinearly with $m$, i.e.,*

$$\lim_{m \to \infty} \frac{K(P_m)}{m} = 0.$$

*Let $q \geq 1$ and $Y_1, \ldots, Y_q, Z \in M$ be arbitrary. Then the "per-bit" difference between expected interaction complexity and interaction information goes to zero for increasing sequence length:*

$$\lim_{m \to \infty} \frac{\sum_{\boldsymbol{x} \in (\{0,1\}^m)^n} P_m(\boldsymbol{x}) \cdot \big(Kc_q(Y_1; \ldots; Y_q \mid Z)\big)(\boldsymbol{x}) - Z.I_q(Y_1; \ldots; Y_q; P_m)}{m} = 0.$$

*Proof.* This follows immediately from Theorem 4.14.                                    □

**Example 4.17.** *As an example to Corollary 4.16, consider the case that we have $n$ parameters $p_1, \ldots, p_n \in (0,1)$ for Bernoulli distributions. Let $P_m$ be the probability mass function given on $\boldsymbol{x} \in (\{0,1\}^m)^n$ by*

$$P_m(\boldsymbol{x}) \coloneqq \prod_{i=1}^n P_m^{p_i}(x_i) \coloneqq \prod_{i=1}^n \prod_{k=1}^m p_i^{x_i^{(k)}} \cdot (1-p_i)^{1-x_i^{(k)}}.$$

*That is, $P_m$ consists of $n$ independent probability mass functions $P_m^{p_i}$ that correspond to $m$ independent Bernoulli distributions with parameter $p_i$. We have $K(P_m) = O(\log m)$ since $m$ is the only moving part in the preceding description for $P_m$, with $p_1, \ldots, p_n$ being independent of $m$. Consequently, Corollary 4.16 can be applied, meaning that the per-bit difference between an expected interaction complexity term and the corresponding interaction information goes to zero. This generalizes the observation after [28], Theorem 10, to $n > 1$ and more complicated interaction terms.*

## 4.6   Hu's Theorem for Prefix-Free Kolmogorov Complexity

We now argue that there is also a Hu theorem for prefix-free Kolmogorov complexity. It requires a logarithmic error term and is therefore less strong than the corresponding theorem for Chaitin's prefix-free Kolmogorov complexity. Additionally, we need to now use $O$-notation, since the equalities only hold for *almost all* inputs: for three functions $f, g, h : (\{0,1\}^*)^n \to \mathbb{R}$, different from

Section 4.5, we now write $f = g + O(h)$ if there is a constant $c \geq 0$ and a threshold $\boldsymbol{x}_0 \in (\{0,1\}^*)^n$ such that

$$\big| f(\boldsymbol{x}) - g(\boldsymbol{x}) \big| \leq c \cdot h(\boldsymbol{x})$$

for all $\boldsymbol{x} \geq \boldsymbol{x}_0$. The latter condition means that $\boldsymbol{x}$ is greater than or equal to $\boldsymbol{x}_0$ in at least one entry, where $\{0,1\}^*$ is ordered lexicographically.

[34], Exercise 3.9.6, shows the following relation:

$$K(y \mid x^*) = K(y \mid x) + O\big( \log K(x) + \log K(y) \big). \tag{26}$$

Overall, this results in the following chain rule for prefix-free Kolmogorov complexity:

**Theorem 4.18** (Chain Rule for Prefix-Free Kolmogorov Complexity)**.** *The following identity holds:*

$$K(x,y) = K(x) + K(y \mid x) + O\big( \log K(x) + \log K(y) \big). \tag{27}$$

*Here, both sides are viewed as functions $\{0,1\}^* \times \{0,1\}^* \to \mathbb{R}$ that map inputs of the form $(x,y)$.*

*Proof.* Combine Theorem 4.4 with Equation (26). $\qquad\square$

To get a precise chain rule, we can, similarly to the case of Chaitin's prefix-free Kolmogorov complexity and motivated by Equation (26), define a new equivalence relation $\sim_K$ on $\mathrm{Maps}\big((\{0,1\}^*)^n, \mathbb{R}\big)$ by

$$F \sim_K H \quad :\Longleftrightarrow \quad F(\boldsymbol{x}) = H(\boldsymbol{x}) + O\left( \sum_{i=1}^{n} \log K(x_i) \right), \quad \text{where } \boldsymbol{x} = (x_1, \ldots, x_n) \in (\{0,1\}^*)^n.$$

We denote the equivalence class of a function $F$ by $[F]_K \in \mathrm{Maps}\big((\{0,1\}^*)^n, \mathbb{R}\big) / \sim_K$. Then, we again use the monoid $M = \{X_1, \ldots, X_n\}^* / \sim$ and define

$$[K]_K : M \times M \to \mathrm{Maps}\big((\{0,1\}^*)^n, \mathbb{R}\big) / \sim_K,$$
$$(Y, Z) \mapsto [K(Y \mid Z)]_K$$

with

$$K(Y \mid Z) : \boldsymbol{x} \mapsto K\big( Y(\boldsymbol{x}) \mid Z(\boldsymbol{x}) \big).$$

Again, this is well-defined by the same arguments as in Lemma 4.6, only that this time, we don't need to use the chain rule in the proof. Furthermore, we can prove an analog of the chain rule given in Proposition 4.5.

**Proposition 4.19** (See Proof 12)**.** *For arbitrary $Y, Z \in M$, the following equality*

$$[K]_K(YZ) = [K]_K(Y) + [K]_K(Z \mid Y)$$

*of elements in $\mathrm{Maps}\big((\{0,1\}^*)^n, \mathbb{R}\big) / \sim_K$ holds.*

Thus, $[K]_K : M \times M \to \mathrm{Maps}\big((\{0,1\}^*)^n, \mathbb{R}\big) / \sim_K$ satisfies all conditions of Corollary 3.3 and we obtain a corresponding Hu theorem for prefix-free Kolmogogorov complexity. This could be worked out similarly to Theorem 4.8, which we leave to the interested reader.

## 4.7  Hu's Theorem for Plain Kolmogorov Complexity

Here, we briefly consider Hu's theorems for plain Kolmogorov complexity $C : \{0,1\}^* \times \{0,1\}^* \to \mathbb{R}$. Recall the $O$-notation from Section 4.6.

The plain Kolmogorov complexity $C : \{0,1\}^* \times \{0,1\}^* \to \mathbb{R}$ is defined in the same way as prefix-free Kolmogorov complexity, but it allows the set of halting programs to not form a prefix-free set, see [34], Chapter 2. This version satisfies the following chain rule:

**Theorem 4.20** (Chain Rule for Plain Kolmogorov Complexity)**.** *The following identity holds:*

$$C(x,y) = C(x) + C(y \mid x) + O\big( \log C(x,y) \big). \tag{28}$$

*Here, both sides are viewed as functions $\{0,1\}^* \times \{0,1\}^* \to \mathbb{R}$ that are defined on inputs of the form $(x,y)$.*

*Proof.* This is proved in [34], Theorem 2.8.                                        □

To get a precise chain rule, we can, similarly as for (Chaitin's) prefix-free Kolmogorov complexity, define a new equivalence relation $\sim_C$ on Maps$\big((\{0,1\}^*)^n, \mathbb{R}\big)$ by

$$F \sim_C H \quad :\Longleftrightarrow \quad F(\boldsymbol{x}) = H(\boldsymbol{x}) + O\big(\log C(\boldsymbol{x})\big), \quad \text{where } \boldsymbol{x} = (x_1, \ldots, x_n) \in (\{0,1\}^*)^n.$$

We denote the equivalence class of a function $F$ by $[F]_C \in \text{Maps}\big((\{0,1\}^*)^n, \mathbb{R}\big)/\sim_C$. Using again the monoid $M = \{X_1, \ldots, X_n\}^*/\sim$, one can define

$$[C]_C : M \times M \to \text{Maps}\big((\{0,1\}^*)^n, \mathbb{R}\big)/\sim_C$$
$$(Y, Z) \mapsto [C(Y \mid Z)]_C$$

with

$$C(Y \mid Z) : \boldsymbol{x} \mapsto C\big(Y(\boldsymbol{x}) \mid Z(\boldsymbol{x})\big).$$

Again, this is well-defined by the same arguments as in Lemma 4.6, and as for prefix-free Kolmogorov complexity, we do not need to use the chain rule in the proof. Furthermore, we can prove an analog of the chain rules given in Proposition 4.5 and Proposition 4.19:

**Proposition 4.21** (See Proof 13). *For arbitrary $Y, Z \in M$, the equality*

$$[C]_C(YZ) = [C]_C(Y) + [C]_C(Z \mid Y)$$

*of elements in* Maps$\big((\{0,1\}^*)^n, \mathbb{R}\big)/\sim_C$ *holds.*

Thus, $[C]_C : M \times M \to \text{Maps}\big((\{0,1\}^*)^n, \mathbb{R}\big)/\sim_C$ satisfies all conditions of Corollary 3.3, and we obtain a corresponding Hu theorem for plain Kolmogogorov complexity. This could again be worked out similarly to Theorem 4.8.

# 5    Further Examples of the Generalized Hu Theorem

In this section, we establish further examples of the premises of Theorem 3.2 and Corollary 3.3, which essentially boils down to finding a chain rule for a function with the correct type signature. For the case of Shannon entropy, the premises were summarized in Summary 2.17. We mostly leave investigations of the specific *meaning* of the resulting higher-order terms to future work, though we do briefly look at the second degree terms for both Kullback-Leibler divergence and the generalization error in machine learning. To keep things simple, we diverge from Sections 2 by only working with *finite* discrete random variables, in the cases where the monoid is based on random variables. As a result, we do not have to worry about questions of convergence and can replace $\Delta_f(\Omega)$ by $\Delta(\Omega)$ and Meas$_{\text{con}}$ by Meas everywhere.

Concretely, we investigate Tsallis $q$-entropy (Section 5.1), Kullback-Leibler divergence (Section 5.2), $q$–Kullback-Leibler divergence (Section 5.3), and cross-entropy (Section 5.4). We also study arbitrary functions on commutative, idempotent monoids (Section 5.5), the special case of submodular information functions (Section 5.6), and the generalization error from machine learning (Section 5.7). Some of the proofs for chain rules are found in Appendix E. The whole section is written in a self-contained way that requires minimal knowledge from the reader.

## 5.1    Tsallis $q$-Entropy

We now investigate the Tsallis $q$-entropy, which was introduced in [55]. We follow the investigations in [57] and translate them into our framework.

That is, assume a finite, discrete sample space $\Omega$, $n$ finite, discrete random variables $X_1, \ldots, X_n$ on $\Omega$, and the monoid $\text{M}(X_1, \ldots, X_n)$ generated by equivalence classes of these random variables, see Definition 2.15. Now, fix an arbitrary number $q \in \mathbb{R} \setminus \{1\}$. Then we define the monoid action

$$\cdot_q : \text{M}(X_1, \ldots, X_n) \times \text{Meas}\big(\Delta(\Omega), \mathbb{R}\big) \to \text{Meas}\big(\Delta(\Omega), \mathbb{R}\big),$$

which we define for $X \in \mathrm{M}(X_1, \ldots, X_n)$, $F \in \mathrm{Meas}\big(\Delta(\Omega), \mathbb{R}\big)$, and $P \in \Delta(\Omega)$ by

$$(X._q F)(P) := \sum_{x \in E_X} P_X(x)^q \cdot F(P|_{X=x}).$$

This is well-defined — meaning that equivalent random variables act in the same way — by the same arguments as in Proposition 2.11. That it is a monoid action can be proved as in Proposition 2.6. Now, define for arbitrary $q \in \mathbb{R} \setminus \{1\}$ the $q$-logarithm by

$$\ln_q : (0, \infty) \to \mathbb{R}, \quad \ln_q(p) := \frac{p^{q-1} - 1}{q - 1}.$$

We have $\lim_{q \to 1} \ln_q(p) = \ln(p)$, as can be seen using l'Hospital's rule. Finally, we can define the Tsallis $q$-entropy $I_1^q : \mathrm{M}(X_1, \ldots, X_n) \to \mathrm{Meas}\big(\Delta(\Omega), \mathbb{R}\big)$ by

$$\big[I_1^q(X)\big](P) := -\sum_{x \in E_X} P_X(x) \ln_q P_X(x) = \frac{\sum_{x \in E_X} P_X(x)^q - 1}{1 - q}.$$

This can be shown to be well-defined similarly as in Proposition 2.10. Since $\lim_{q \to 1} \ln_q p = \ln p$, we consequently also have $\lim_{q \to 1} I_1^q(X; P) = I_1(X; P)$. That is, the $q$-entropy generalizes the Shannon entropy.

The following chain rule guarantees the existence of a corresponding Hu theorem.

**Proposition 5.1** (See Proof 14). *$I_1^q : \mathrm{M}(X_1, \ldots, X_n) \to \mathrm{Meas}\big(\Delta(\Omega), \mathbb{R}\big)$ satisfies the chain rule*

$$I_1^q(XY) = I_1^q(X) + X._q I_1^q(Y)$$

*for all $X, Y \in \mathrm{M}(X_1, \ldots, X_n)$.*

## 5.2 Kullback-Leibler Divergence

In this section, we study the chain rule of Kullback-Leibler divergence. It resembles the one described in [57], chapter 3.7, in the language of information cohomology. A more elementary formulation of the chain rule can also be found in [17], Theorem 2.5.3, which is applied in their Section 4.4 to prove a version of the second law of thermodynamics. In the end, we will also briefly study and interpret *KL divergence of degree 2*, in analogy to mutual information $I_2$, in Example 5.3.

Let again the monoid $\mathrm{M}(X_1, \ldots, X_n)$ of $n$ discrete random variables on $\Omega$ be given. For $P, Q \in \Delta(\Omega)$, we write $P \ll Q$ if for all $\omega \in \Omega$, the following implication is true: $Q(\omega) = 0 \implies P(\omega) = 0$. In the literature, $P$ is then called *absolutely continuous* with respect to the measure $Q$. We set

$$\widetilde{\Delta(\Omega)^2} := \Big\{(P, Q) \in \Delta(\Omega)^2 \ \Big| \ P \ll Q\Big\}.$$

We will silently make use of the fact that $P \ll Q$ implies $P_X \ll Q_X$ and $P|_{X=x} \ll Q|_{X=x}$ for all discrete random variables $X : \Omega \to E_X$ and $x \in E_X$.

We now define $G := \mathrm{Meas}\Big(\widetilde{\Delta(\Omega)^2}, \mathbb{R}\Big)$. We write elements $F \in G$ applied to inputs $(P, Q)$ as $F(P\|Q)$. Define for $X \in \mathrm{M}(X_1, \ldots, X_n)$ and $F \in \mathrm{Meas}\Big(\widetilde{\Delta(\Omega)^2}, \mathbb{R}\Big)$, and $P \ll Q \in \Delta(\Omega)$ the monoid action by:

$$(X.F)(P\|Q) := \sum_{x \in E_X} P_X(x) F\big(P|_{X=x}\|Q|_{X=x}\big).$$

Similarly as before, this is a well-defined, additive monoid action. In the following, we use the convention that $0 \cdot x = 0$ for $x \in \mathbb{R} \cup \{\pm\infty\}$ and $\ln(0) = -\infty$. Finally, we define the function $D_1 : \mathrm{M}(X_1, \ldots, X_n) \to \mathrm{Meas}\Big(\widetilde{\Delta(\Omega)^2}, \mathbb{R}\Big)$ as the Kullback-Leibler divergence, given for all $X \in \mathrm{M}(X_1, \ldots, X_n)$ and $P \ll Q \in \Delta(\Omega)$ by

$$\big[D_1(X)\big](P\|Q) := D_1\big(X; P\|Q\big) := -\sum_{x \in E_X} P_X(x) \ln \frac{Q_X(x)}{P_X(x)}.$$
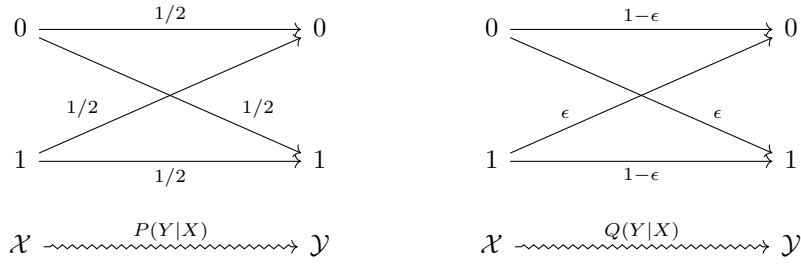
This is well-defined, and we obtain:

Figure 5: Binary symmetric channels for the joint distributions $P$ and $Q$ in Example 5.3. For a uniform prior $P(X) = Q(X)$, $P$ and $Q$ have the same marginals $P(Y) = Q(Y)$, but differ in their conditionals $P(Y \mid X)$ and $Q(Y \mid X)$. This leads for small $\epsilon > 0$ to an arbitrarily large negative mutual Kullback-Leibler divergence $\big[D_2(X;Y)\big](P\|Q)$.

**Proposition 5.2** (See Proof 15). $D_1 : \mathrm{M}(X_1, \ldots, X_n) \to \mathrm{Meas}\left(\widetilde{\Delta(\Omega)^2}, \mathbb{R}\right)$ satisfies the chain rule for all $X, Y \in \mathrm{M}(X_1, \ldots, X_n)$:

$$D_1(XY) = D_1(X) + X.D_1(Y).$$

**Example 5.3.** In [25], the following situation is discussed: $\mathcal{X}$ and $\mathcal{Y}$ are finite sets, and $\Omega = \mathcal{X} \times \mathcal{Y}$. One can consider the two marginal variables

$$X : \mathcal{X} \times \mathcal{Y} \to \mathcal{X}, \quad (x, y) \mapsto x,$$
$$Y : \mathcal{X} \times \mathcal{Y} \to \mathcal{Y}, \quad (x, y) \mapsto y.$$

A channel from $X$ to $Y$ is a conditional distribution $P(Y \mid X)$. Together with a prior distribution $P(X)$, it forms a joint $P(X, Y)$ over $\mathcal{X} \times \mathcal{Y}$. Now, take two distributions $P \ll Q \in \Delta(\mathcal{X} \times \mathcal{Y})$. Then, as noted in [25], the chain rule Proposition 5.2 shows the following:

$$D_1\big(P\|Q\big) = D_1\big(P(X)\|Q(X)\big) + \sum_{x \in \mathcal{X}} P(x) \cdot D_1\big(P(Y \mid x)\|Q(Y \mid x)\big).$$

Note that for ease of notation, we write $P(X)$ for $P_X$, $D_1\big(P(X)\|Q(X)\big)$ for $\big[D_1(X)\big](P\|Q)$, $P(x)$ for $P_X(x)$, $P(Y \mid x)$ for $(P|_{X=x})_Y$, etc.

In our context, the "mutual Kullback-Leibler divergence" $D_2(X;Y)$ is of interest. With respect to $P$ and $Q$, it is given according to Equation (7) and using symmetry of $D_2$ (which follows from Theorem 3.2 due to set operations being symmetric) as follows:

$$\big[D_2(X;Y)\big](P\|Q) = D_1\big(P(Y)\|Q(Y)\big) - \sum_{x \in \mathcal{X}} P(x) \cdot D_1\big(P(Y \mid x)\|Q(Y \mid x)\big).$$

It is well-known that a simple use of Jensen's inequality proves the non-negativity of the Kullback-Leibler divergence $D_1$. We also know that mutual information $I_2$ is non-negative. Can the same be said about the mutual Kullback-Leibler divergence $D_2$?

The answer is no. Consider the case $\mathcal{X} = \mathcal{Y} = \{0, 1\}$, and let the prior distributions $P(X) = Q(X)$ both be uniform. Furthermore, let $P(Y \mid X)$ and $Q(Y \mid X)$ be binary symmetric channels ( [17], Section 7.1.4), given as in Figure 5. Note that the marginal distributions $P(Y)$ and $Q(Y)$ are identical, and so

$$D_1\big(P(Y)\|Q(Y)\big) = 0.$$

We now work with binary logarithms $\log$. For the second term, we then obtain

$$\sum_{x \in \{0,1\}} P(x) \cdot D_1\big(P(Y \mid x)\|Q(Y \mid x)\big) = \sum_{x \in \{0,1\}} P(x) \sum_{y \in \{0,1\}} P(y \mid x) \log \frac{P(y \mid x)}{Q(y \mid x)}$$

$$= \frac{1}{4} \cdot \left[ \log \frac{P(0 \mid 0)}{Q(0 \mid 0)} + \log \frac{P(1 \mid 0)}{Q(1 \mid 0)} + \log \frac{P(0 \mid 1)}{Q(0 \mid 1)} + \log \frac{P(1 \mid 1)}{Q(1 \mid 1)} \right]$$

$$= \frac{1}{4} \cdot \Big[ -4 - 2\log(1-\epsilon) - 2\log(\epsilon) \Big]$$

$$= -1 - \frac{1}{2} \cdot \Big[ \log(1-\epsilon) + \log(\epsilon) \Big]$$

*Note that for very small $\epsilon$, $\log(1-\epsilon)$ becomes negligible and $\log(\epsilon)$ approaches $-\infty$, and so the term above approaches $+\infty$. Overall, this means that*

$$\big[ D_2(X;Y) \big](P\|Q) = - \sum_{x \in \{0,1\}} P(x) \cdot D_1\big( P(Y \mid x) \| Q(Y \mid x) \big) < 0$$

*is negative, and even unbounded, reaching $-\infty$ as $Q$ becomes deterministic. We can compare this conceptually to mutual information as follows: $I_2(X;Y)$ is the average reduction of uncertainty in $Y$ when learning about $X$. Similarly, we can interpret $D_2(X;Y)$ as the average reduction of Kullback-Leibler divergence between two marginal distributions in $Y$ when learning about $X$. However, in this case, the divergence only becomes visible when the evaluation of $X$ is known, since there is no difference in the marginals $P(Y)$ and $Q(Y)$. Thus, the "reduction" is actually negative.*

## 5.3   $q$–Kullback-Leibler Divergence

Similarly to the Tsallis $q$-entropy from Section 5.1, one can also define a $q$–Kullback-Leibler divergence, as is done in [57], Chapter 3.7.[9] The monoid action $._q : \mathrm{M}(X_1, \ldots, X_n) \times \mathrm{Meas}\left( \widetilde{\Delta(\Omega)^2}, \mathbb{R} \right) \to \mathrm{Meas}\left( \widetilde{\Delta(\Omega)^2}, \mathbb{R} \right)$ is now given by

$$(X._q F)(P\|Q) := \sum_{x \in E_X} P_X(x)^q Q_X(x)^{1-q} \cdot F\big( P|_{X=x} \| Q|_{X=x} \big).$$

Now, we define the $q$–Kullback-Leibler divergence $D_1^q : \mathrm{M}(X_1, \ldots, X_n) \to \mathrm{Meas}\left( \widetilde{\Delta(\Omega)^2}, \mathbb{R} \right)$ for all $X \in \mathrm{M}(X_1, \ldots, X_n)$ and $P \ll Q \in \Delta(\Omega)$ as the following generalization of standard Kullback-Leibler divergence:

$$\big[ D_1^q(X) \big](P\|Q) := \sum_{x \in E_X} P_X(x) \ln_q \frac{P_X(x)}{Q_X(x)} = \frac{\sum_{x \in E_X} P_X(x)^q Q_X(x)^{1-q} - 1}{q - 1}.$$

**Proposition 5.4** (See Proof 16). *$D_1^q : \mathrm{M}(X_1, \ldots, X_n) \to \mathrm{Meas}\left( \widetilde{\Delta(\Omega)^2}, \mathbb{R} \right)$ satisfies the chain rule*

$$D_1^q(XY) = D_1^q(X) + X._q D_1^q(Y)$$

*for all $X, Y \in \mathrm{M}(X_1, \ldots, X_n)$.*

## 5.4   Cross-Entropy

We choose the same monoid action $. : \mathrm{M}(X_1, \ldots, X_n) \times \mathrm{Meas}\left( \widetilde{\Delta(\Omega)^2}, \mathbb{R} \right) \to \mathrm{Meas}\left( \widetilde{\Delta(\Omega)^2}, \mathbb{R} \right)$ as for the Kullback-Leibler divergence. The cross-entropy $C_1 : \mathrm{M}(X_1, \ldots, X_n) \to \mathrm{Meas}\left( \widetilde{\Delta(\Omega)^2}, \mathbb{R} \right)$ is given by:

$$\big[ C_1(X) \big](P\|Q) := C_1\big( X; P\|Q \big) := - \sum_{x \in E_X} P_X(x) \ln Q_X(x).$$

**Proposition 5.5.** *$C_1$ satisfies the chain rule for all $X, Y \in \mathrm{M}(X_1, \ldots, X_n)$:*

$$C_1(XY) = C_1(X) + X.C_1(Y).$$

---

[9]Our definition differs from the one given in [57] by using a slightly different definition of the $q$-logarithm. We did this to be consistent with the definition of the Tsallis $q$-entropy above.

*Proof.* This follows with the same arguments as Proposition 5.2.                    □

**Remark 5.6.** *One can easily show the following well-known relation between cross-entropy $C_1$, Shannon entropy $I_1$, and Kullback-Leibler divergence $D_1$:*

$$\big[C_1(X)\big](P\|Q) = \big[I_1(X)\big](P) + \big[D_1(X)\big](P\|Q).$$

*This means that the study of $C_q$ is entirely subsumed by that of $I_q$ and $D_q$. Since we already looked at $D_2$ in Example 5.3, we omit looking at $C_2$ here.*

## 5.5   Arbitrary Functions on Commutative, Idempotent Monoids

Let $M$ be any commutative monoid, and $R : M \to G$ be *any* function into an abelian group $G$. Define the two-argument function $R_1 : M \times M \to G$ by $R_1(A \mid B) \coloneqq R(AB) - R(B)$. Set $R_1(A) \coloneqq R_1(A \mid \mathbf{1}) = R(A) - R(\mathbf{1})$, where $\mathbf{1} \in M$ is the neutral element. These definitions mean that the chain rule is satisfied *by definition*, making Hu's theorem a purely combinatorial fact. The reader can verify the following proposition:

**Proposition 5.7.** $R_1 : M \times M \to G$ *satisfies the chain rule*

$$R_1(AB) = R_1(A) + R_1(B \mid A)$$

*for all $A, B \in M$.*

Therefore, if $M$ is also idempotent and finitely generated, then $R_1 : M \times M \to G$ satisfies all conditions of Corollary 3.3, and one obtains a corresponding Hu theorem.

## 5.6   Submodular Information Functions

Using the framework of Section 5.5, we can study the submodular information functions from [22, 38, 53], which they use to formulate generalizations of conditional independence and the causal Markov condition. Alternatively, we could also analyze general submodular set functions [46], but decided to restrict to submodular information functions since they are closer to our interests.

Recall that a lattice is a tuple $L = (L, \vee, \wedge)$ consisting of a set $L$ together with commutative, associative, and idempotent operations $\vee, \wedge : L \times L \to L$ that satisfy the absorption rules $a \vee (a \wedge b) = a$ and $a \wedge (a \vee b) = a$. Given a lattice $L$, one can define a corresponding partial order on $L$ by $a \leq b$ if $a = a \wedge b$.

From now on, let $(L, \vee, \wedge)$ be a finite lattice, meaning that $L$ is a finite set. One can define $\mathbf{0} \coloneqq \bigwedge_{a \in L} a$, the meet of the finitely many elements in $L$. By the axioms above, this is neutral with respect to the join operation, that is $b \vee \mathbf{0} = b$. Note that $\mathbf{0} \wedge b = \mathbf{0}$ for all $b \in L$ due to the second absorption rule above. Consequently, $\mathbf{0} \leq b$ for all $b \in L$. [22, 38, 53] then study the concept of a submodular (information) function. We follow the version outlined in [53]:

**Definition 5.8** (Submodular Information Function). *Let $L$ be a finite lattice. Then a function $R : L \to \mathbb{R}$ is called a* submodular information function *if all of the following conditions hold for all $a, b \in L$:*

1. *normalization: $R(\mathbf{0}) = 0$;*

2. *monotonicity: $a \leq b$ implies $R(a) \leq R(b)$;*

3. *submodularity: $R(a) + R(b) \geq R(a \vee b) + R(a \wedge b)$.*

*In particular, the second property implies $R(b) \geq R(\mathbf{0}) = 0$, meaning $R$ is non-negative.*

They then define the conditional $R_1 : L \times L \to \mathbb{R}$ by $R_1(a \mid b) \coloneqq R(a \vee b) - R(a)$. Furthermore, to define conditional independence and obtain a generalized causal Markov condition, they define the conditional mutual information $I : L^2 \times L \to \mathbb{R}$ by

$$I(a; b \mid c) \coloneqq R(a \vee c) + R(b \vee c) - R(a \vee b \vee c) - R(c).$$

Now, note that $(L, \vee, \mathbf{0})$ is a finitely generated, commutative, idempotent monoid. Thus, Proposition 5.7 shows that $R_1$ gives rise to Hu's theorem for higher-order functions $R_2, R_3, \ldots$, as defined in Corollary 3.3. We can easily see that $R_2$ agrees with the definition of $I$ from above:

$$R_2(a; b \mid c) := R_1(a \mid c) - R_1(a \mid b \vee c)$$
$$= R(a \vee c) - R(c) - R(a \vee b \vee c) + R(b \vee c)$$
$$= I(a; b \mid c).$$

As special cases of submodular information functions, [53] consider Shannon entropy on sets of random variables, Chaitin's prefix-free Kolmogorov complexity, other compression-based information functions, period lengths of time series, and the size of a vocabulary in a text.

## 5.7 Generalization Error

Before coming to the generalization error, we briefly consider the situation dual to that of Section 5.5. Let $M$ be a commutative, idempotent monoid. Let $G$ be an abelian group and $\mathcal{E} : M \to G$ be any function. Define $\mathrm{Ad} : M \times M \to G$ by $\mathrm{Ad}(A \mid B) := \mathcal{E}(B) - \mathcal{E}(AB)$ Here, Ad stands intuitively for "advantage", a terminology that becomes clear in the machine learning example below. Similarly as in the case of Kolmogorov complexity, define $\mathrm{Ad}(A) := \mathrm{Ad}(A \mid \mathbf{1}) = \mathcal{E}(\mathbf{1}) - \mathcal{E}(A)$. The reader can easily verify the chain rule:

**Proposition 5.9.** $\mathrm{Ad} : M \times M \to G$ *satisfies the chain rule: one has*

$$\mathrm{Ad}(AB) = \mathrm{Ad}(A) + \mathrm{Ad}(B \mid A)$$

*for all* $A, B \in M$.

Consequently, $\mathrm{Ad} : M \times M \to G$ satisfies the assumptions of Corollary 3.3. One then obtains a corresponding Hu theorem.

We now specialize this investigation to the *generalization error* from machine learning [37, 47]. In this case, let $J = [n]$ be a finite set and the monoid be given by $2^J = (2^J, \cup, \emptyset)$, i.e., $\cup$ is the operation and $\emptyset$ the neutral element. This monoid is idempotent, commutative, and finitely generated by $\{1\}, \ldots, \{n\}$.

For all $j \in J$, let $\mathcal{X}_j$ be a measurable space. Let $(X_j)_{j \in J}$ be the random variable of feature tuples with values in $\prod_{j \in J} \mathcal{X}_j$. Similarly, let $\mathcal{Y}$ be another measurable space and $Y$ the random variable of labels in $\mathcal{Y}$. A typical assumption is that there exists a joint distribution $P := P\big((X_j)_{j \in J}, Y\big)$ from which "the world samples the data". Additionally, let $\Delta(\mathcal{Y})$ be the space of probability measures on $\mathcal{Y}$, and $L : \Delta(\mathcal{Y}) \times \mathcal{Y} \to \overline{\mathbb{R}} := \mathbb{R} \cup \{+\infty\}$ a loss function that compares a model distribution over labels to the true label.

For all $A \subseteq J$, assume that $\mathcal{F}(A) \subseteq \mathrm{Maps}\big(\prod_{a \in A} \mathcal{X}_a, \Delta(\mathcal{Y})\big)$ is a class of functions[10] that, given a feature tuple with indices in $A$, predicts a distribution over $\mathcal{Y}$. We call this the set of *hypotheses* for predicting the labels given features in $A$. For a hypothesis $q \in \mathcal{F}(A)$ and $x_A \in \prod_{a \in A} \mathcal{X}_a$, we denote the output by $q(Y \mid x_A) := q(x_A) \in \Delta(\mathcal{Y})$. A learning algorithm with access to features in $A$ is supposed to find a hypothesis $q \in \mathcal{F}(A)$ that minimizes the *generalization error*:

$$\mathcal{E}(A) := \inf_{q \in \mathcal{F}(A)} \mathrm{E}_{(\hat{x}, \hat{y}) \sim P} \Big[ L\big(q(Y \mid \hat{x}_A) \ \| \ \hat{y}\big) \Big].$$

Then, as above, define $\mathrm{Ad}_Y : 2^J \times 2^J \to \mathbb{R}$ by

$$\mathrm{Ad}_Y \big( X_A \mid X_B \big) := \mathcal{E}(B) - \mathcal{E}(A \cup B).\text{[11]}$$

From Proposition 5.9, we obtain the following chain rule:

$$\mathrm{Ad}_Y(X_{A \cup B}) = \mathrm{Ad}_Y(X_A) + \mathrm{Ad}_Y(X_B \mid X_A). \tag{29}$$

---

[10]Further below, we will make the assumption that $A \subseteq B$ implies $\mathcal{F}(A) \subseteq \mathcal{F}(B)$, in a suitable sense. We make no other assumptions on the collection of $\mathcal{F}(A)$ for $A \subseteq J$.

[11]There is a one-to-one correspondence between all $A \in 2^J$ and all variables $X_A$ with $A \in 2^J$. We simply denote the monoid of all $X_A$ again by $2^J$, with the multiplication rule becoming $X_A X_B = X_{A \cup B}$.

To interpret this chain rule sensibly, we make one further assumption: namely that, when having access to *more features*, the learning algorithm can still use all hypotheses that simply *ignore these additional features*. More precisely, for $B \subseteq C \subseteq J$, let us interpret each map $q_B \in \mathcal{F}(B)$ as a function $\widetilde{q_B} : \prod_{c \in C} \mathcal{X}_c \to \Delta(\mathcal{Y})$ by

$$\widetilde{q_B}\big((x_c)_{c \in C}\big) \coloneqq q_B\big((x_b)_{b \in B}\big).$$

The assumption is that $\widetilde{q_B} \in \mathcal{F}(C)$, for all $B \subseteq C \subseteq J$ and $q_B \in \mathcal{F}(B)$. Overall, we can interpret this as $\mathcal{F}(B) \subseteq \mathcal{F}(C)$. It follows that $\mathcal{E}(B) \geq \mathcal{E}(C)$. Consequently, for all $A, B \subseteq J$ (without any inclusion imposed), it follows

$$\mathrm{Ad}_Y(X_A \mid X_B) = \mathcal{E}(B) - \mathcal{E}(A \cup B) \geq 0. \tag{30}$$

The meaning of this is straightforward: $\mathrm{Ad}_Y(X_A \mid X_B)$ measures what a perfect learning algorithm can gain from knowing all the features in $A$ if it already has access to all the features in $B$ — the *advantage* motivating the notation $\mathrm{Ad}_Y(X_A \mid X_B)$. The chain rule, Equation (29), thus says the following: for a perfect learning algorithm, the advantage from getting access to features in $A \cup B$ equals the advantage it receives from the features in $A$, plus the advantage it receives from $B$ when it already has access to $A$.

We can then ask: is then the "mutual advantage", as defined from Equation (11) by

$$\mathrm{Ad}_Y^2(X_A; X_B) \coloneqq \mathrm{Ad}_Y(X_A) - \mathrm{Ad}_Y(X_A \mid X_B),$$

necessarily also positive, as we expect from the case of entropy and mutual information? The answer is *no*, as the following simple example shows:

**Example 5.10.** *Let $J = \{1, 2\}$, $\mathcal{X}_1 = \mathcal{X}_2 = \mathcal{Y} = \{0, 1\}$, $X_1, X_2$ two independent Bernoulli distributed random variables, and $Y$ be the result of applying a XOR gate to $X_1$ and $X_2$. In other words, the joint distribution $P(X_1, X_2, Y) \in \Delta\big(\{0, 1\}^3\big)$ is the unique distribution with*

$$P\big(X_1 = 0, X_2 = 0, Y = 0\big) = 1/4,$$
$$P\big(X_1 = 0, X_2 = 1, Y = 1\big) = 1/4,$$
$$P\big(X_1 = 1, X_2 = 0, Y = 1\big) = 1/4,$$
$$P\big(X_1 = 1, X_2 = 1, Y = 0\big) = 1/4.$$

*We define the loss function $L : \Delta\big(\{0, 1\}\big) \times \{0, 1\} \to \overline{\mathbb{R}}$ as the cross-entropy loss: $L\big(q(Y) \parallel y\big) \coloneqq -\log q(y)$, where $\log$ is the binary logarithm. Furthermore, we define $\mathcal{F}(A) \coloneqq \big\{q : \mathcal{X}_A \to \Delta(\{0, 1\})\big\}$ as the space of* all *possible prediction functions with access to features in $A \subseteq J = \{1, 2\}$. Now, note that if one does not have access to both features, i.e. $A \neq \{1, 2\}$, then it is impossible to do better than random, since $X_1 \perp\!\!\!\perp Y$ and $X_2 \perp\!\!\!\perp Y$. Thus, in that case, the best prediction is $q(\hat{y} \mid \hat{x}_A) = 1/2$, irrespective of $\hat{x}$ and $\hat{y}$. If, however, one has access to both features, then perfect prediction is possible, since $Y$ is a deterministic function of $(X_1, X_2)$. Using $-\log(1/2) = 1$ and $-\log(1) = 0$, this leads to the following generalization errors:*

$$\mathcal{E}\big(\emptyset\big) = 1, \quad \mathcal{E}\big(\{1\}\big) = 1, \quad \mathcal{E}\big(\{2\}\big) = 1, \quad \mathcal{E}\big(\{1, 2\}\big) = 0.$$

*Consequently, the mutual advantage of $X_1$ with $X_2$ is given by*

$$\begin{aligned} \mathrm{Ad}_Y^2(X_1; X_2) &= \mathrm{Ad}_Y(X_1) - \mathrm{Ad}_Y(X_1 \mid X_2) \\ &= \mathcal{E}\big(\emptyset\big) - \mathcal{E}\big(\{1\}\big) - \mathcal{E}\big(\{2\}\big) + \mathcal{E}\big(\{1, 2\}\big) \\ &= -1 \\ &< 0. \end{aligned}$$

*Thus, in this example, the mutual advantage is negative. Rearranging the inequality, we can read this as $\mathrm{Ad}_Y(X_1) < \mathrm{Ad}_Y(X_1 \mid X_2)$. In general, beyond the specifics of this example, the inequality*

$$\mathrm{Ad}_Y(X_A) < \mathrm{Ad}_Y(X_A \mid X_B)$$

*means that features in $A \subseteq J$ are more predictive of $Y$ if we already have access to features in $B$. This indicates a case of* feature interaction *or* synergy: *the contribution of a set of features in predicting $Y$ is greater than the individual contribution of each single feature. Intuitively, we expect such situations in many machine learning applications, and think it might be worthwhile to investigate the meaning of the higher degree interaction terms* $\mathrm{Ad}_Y^q$ *appearing in Hu's theorem as in Corollary 3.3.*

# 6   Discussion

## 6.1   Major Findings: a Generalization of Hu's Theorem and its Applications

In this work, we have systematically abstracted away from the details of Shannon's information theory [48, 49] to generalize Hu's theorem [32] to new situations. To obtain information diagrams, one simply needs a finitely generated commutative, idempotent monoid $M$ — also known under the name of a join-semilattice — acting additively on an abelian group $G$, and a function $F_1 : M \to G$ satisfying the chain rule of information: $F_1(XY) = F_1(X) + X.F_1(Y)$. Alternatively, with $M$ and $G$ being as above, the additive monoid action and $F_1$ together can be replaced by a two-argument function $K_1 : M \times M \to G$ satisfying the chain rule: $K_1(XY) = K_1(X) + K_1(Y \mid X)$. The proof of the main result — Theorem 3.2 together with Corollary 3.3 — is very similar to the one given in [63] for the case of Shannon entropy; the main insight is that it is possible to express the basic *atoms* of an information diagram with an inclusion-exclusion type expression over "unions of disks":

$$\mu(p_I) = \sum_{\emptyset \neq K \supseteq I^c} (-1)^{|K|+|I|+1-n} \cdot F_1(X_K) = \sum_{K \subseteq I} (-1)^{|K|+1} \cdot F_1(X_K X_{I^c}).$$

This formula is visually motivated in Section 3.2. Relations to different interaction terms are explored in Section 3.3.

With the monoid given by equivalence classes of (countably infinite) discrete random variables, the abelian group by measurable functions on probability measures, and the additive monoid action by the conditioning of information functions, we recover information diagrams for Shannon entropy, see Summary 2.17. Beyond this classical case, we obtained Hu's theorems for several versions of Kolmogorov complexity [34] (Section 4), Tsallis $q$-entropy [55], Kullback-Leibler divergence, $q$–Kullback-Leibler divergence, cross-entropy [57], general functions on commutative, idempotent monoids, submodular information functions [53], and the generalization error from machine learning [37, 47] (all in Section 5). For Kolmogorov complexity, we generalized the well-known theme that "expected Kolmogorov complexity is close to Shannon entropy":

$$\text{"expected interaction complexity"} \quad \approx \quad \text{"interaction information"}.$$

For well-behaved probability distributions, this results in the limit of infinite sequence length in an actual *equality* of the per-bit quantities for the two concepts (Section 4.5).

## 6.2   The Cohomological Context of this Work

The main context in which our ideas developed is information cohomology [5, 9, 57, 58]. The setup of that work mainly differs by using partition lattices instead of equivalence classes of random variables and generalizing this further to so-called *information structures*. The functions satisfying the chain rule are reformulated as so-called "cocycles" in that cohomology theory, which are "cochains" whose "coboundary" vanishes:

$$(\delta F_1)(X;Y) := X.F_1(Y) - F_1(XY) + F_1(X) = 0.$$

That gives these functions a context in the realm of many cohomology theories that were successfully developed in mathematics. The one defined by Gerhard Hochschild for associative algebras is maybe most closely related [30]. For the special case of probabilistic information cohomology, [5, 57] were able to show that Shannon entropy is not only *a* cocycle, but is in some precise sense the *unique* cocycle generating all others of degree 1. Thus, Shannon entropy finds a fully cohomological

interpretation. Arguably, without the abstract nature of that work and the consistent emphasis on abstract structures like monoids and monoid actions, our work would not have been possible.

There is one way in which information cohomology tries to go beyond Shannon information theory: it tries to find higher degree cocycles that *differ* from the interaction terms $F_q$. This largely unsolved task has preliminary investigations in [57], Section 3.6, and [21]. In that sense, information cohomology can be viewed as a generalization of Hu's theorem. Since some limitations in the expressiveness of interaction information are well-known [33], we welcome any effort to make progress on that task.

## 6.3   Unanswered Questions and Future Directions

**Further generalizations**   On the theoretical front, it should be possible to generalize Hu's theorem further from commutative, idempotent monoids to what [57] calls *conditional meet semi-lattices*. As these *locally* are commutative, idempotent monoids, the generalization can probably directly use our result.

**A transport of ideas**   More practically, we hope that the generalization of Hu's theorem leads to a transport of ideas from the theory of Shannon entropy to other functions satisfying the chain rule. There are many works that study information-theoretic concepts based on the interaction information functions and thus ultimately Shannon entropy, for example O-information [27, 43], total correlation [60], dual total correlation [29], and information paths [3, 6]. All of these can trivially be defined for functions satisfying the chain rule that go beyond Shannon entropy, and can thus be generalized to all the example applications in Sections 4 and 5. Most of the basic algebraic properties should carry over since they often follow from Hu's theorem itself. It is our hope that studying such quantities in greater generality may lead to new insights into the newly established application areas of Hu's theorem.

Additionally, it should not be forgotten that even Shannon interaction information *itself* deserves to be better understood. Understanding these interaction terms in a more general context could help for resolving some of the persisting confusions about the topic. One of them surrounds the possible negativity of interaction information $I_3(X; Y; Z)$ of three (and more) random variables [4, 8], which is sometimes understood as meaning that there is more synergy than redundancy present [61, 62]. Similarly, we saw in Example 5.10 that the mutual feature advantage $I_Y^2(X_A; X_B)$ can be negative as well, which has a clear interpretation in terms of synergy. Example 5.3 shows that the mutual Kullback-Leibler divergence $D_2(X; Y)$ of two distributions $P \ll Q$ can be negative if knowing $X$ "reveals" the divergence of $P$ and $Q$ in $Y$. We would welcome more analysis in this direction, ideally in a way that transcends any particular applications and could thus shed new light on the meaning of classical interaction information.

**Further chain rules**   It goes without saying that we were likely not successful in finding *all* functions satisfying a chain rule. One interesting candidate seems to be differential entropy $h$ ([17], Theorem 8.6.2):

$$h(X, Y) = h(X) + h(Y \mid X).$$

However, it seems to us that differential entropy is not well-behaved. For example, if $X$ is a random variable with values in $\mathbb{R}$, then even if $h(X)$ exists, the differential entropy of the joint variable $(X, X)$ with values in $\mathbb{R}^2$ is negative infinity:

$$h(X, X) = -\infty.$$

In particular, we have $h(X) \neq h(X, X)$, and so Hu's theorem cannot hold.

As clarified, for example, in [59], differential entropy is measured *relative to a given base measure*. Given that $(X, X)$ takes values only in the diagonal of $\mathbb{R}^2$, which has measure 0, explains why the differential entropy degenerates. To remedy this, one would need to change the base measure to also live on the diagonal; it is unclear to us how to interpret this, or if a resulting Hu theorem could indeed be deduced.

Another possible candidate is quantum entropy, also called von Neumann entropy, which also allows for a conditional version that satisfies a chain rule ([14], Theorem 1). Interestingly, conditional quantum entropy, also called partial quantum information, can be negative [13, 31], which contrasts it from classical Shannon entropy.

In analogy to the Kullback-Leibler divergence (Section 5.2), also quantum entropy admits a relative version, which has many applications in quantum information theory [56]. In [23], a chain rule for quantum relative entropy was proven, which, however, is an *inequality*. In [39], Proposition 1 and Example 1, one can find a chain rule–type statement for quantum relative entropy that generalizes the one for non-relative quantum conditional entropy. We leave the precise meaning or interpretation of these results in the context of our work to future investigations.

**Kolmogorov complexity and information decompositions**    In the context of Kolmogorov complexity, we would welcome a more thorough analysis of the size of the constants involved in Theorems 4.8 and 4.14, potentially similar to [65]. More precisely, it would be worthwhile to improve on the dependence on $q$ or $n$ that we explain in Remarks 4.9 and 4.15.

More broadly, one could try to understand complex interactions that go beyond interaction information in the context of Kolmogorov complexity.[12] For example, partial information decomposition (PID) [61, 62][13] aims to complement the usual information functions with unique information, shared information, and complementary information. It argues that the mutual information of a random variable $Z$ with a joint variable $(X, Y)$ can be decomposed as follows:

$$I_2\big((X, Y); Z\big) = \underbrace{UI(X \setminus Y; Z)}_{\text{unique}} + \underbrace{UI(Y \setminus X; Z)}_{\text{unique}} + \underbrace{SI(X, Y; Z)}_{\text{shared}} + \underbrace{CI(X, Y; Z)}_{\text{complementary}}.$$

Here, $UI(X \setminus Y; Z)$ is the information that $X$ provides about $Z$ that is not also contained in $Y$; $SI(X, Y; Z)$ is the information that $X$ and $Y$ both share about $Z$; and finally, $CI(X, Y; Z)$ is the information that $X$ and $Y$ can *only together* provide about $Z$, but neither on its own. $SI$ is also called "redundant information", and $CI$ "synergistic information". This then leads to an interpretation of interaction information as a difference of shared and complementary information:

$$I_3(X, Y, Z) = \underbrace{SI(X, Y; Z)}_{\text{shared}} - \underbrace{CI(X, Y; Z)}_{\text{complementary}}.$$

While it is known that such functions exist, no proposals have yet satisfied all axioms that are considered desirable. In this sense, the search for shared, redundant, and synergistic information in the framework of PID is still ongoing [35]. See also [10, 24, 40] for related work.

We could imagine that attempting a similar decomposition for Kolmogorov complexity could provide new insights. To argue that this might be possible, we can look, for example, at the thought experiment of $x$ and $y$ being binary strings encoding physical theories, and $z$ being a binary string containing data about a physical phenomenon. Then a hypothesized "algorithmic complementary information" $CI(x, y; z)$ would intuitively be high if the theories $x$ and $y$ *only together* allow explaining (parts of) the data $z$; a high shared information $SI(x, y; z)$ would mean that $x$ and $y$ are theories that are *equally* able to explain (parts of) the data in $z$. One hope is that averaging such quantities leads to a partial information decomposition in the usual information-theoretic sense, thus providing a new bridge that helps with the transport of ideas between fields:

$$\text{"expected algorithmic PID"} \quad \overset{?}{\approx} \quad \text{"PID"}.$$

## 6.4   Conclusion

To restate our main finding, we can say: whenever you find a chain rule

$$F_1(XY) = F_1(X) + X.F_1(Y),$$

you will under mild conditions obtain information diagrams. Most of their implications are yet to be understood.

---

[12]Or in the context of any other of the application areas in Section 5 of our generalized Hu theorem.

[13]The only privately communicated version, [62], of [61], has a stronger emphasis on the axiomatic framework and is more up to date.

# Appendix

## A   Measure Theory for Countable Discrete Spaces

In this section, we investigate some technical details related to the measurability of certain functions. For more background on measure theory, any book on the topic suffices, for example [54] and [45]. As the results are elementary, we leave most of them to the reader to prove.

Recall that for a measurable space $\mathcal{Z}$, the space of probability measures $\Delta(\mathcal{Z})$ on $\mathcal{Z}$ carries the smallest $\sigma$-algebra that makes all evaluation maps

$$\mathrm{ev}_A : \Delta(\mathcal{Z}) \to [0,1], \quad P \mapsto P(A)$$

for measurable $A \subseteq \mathcal{Z}$ measurable. Also recall that discrete random variables are functions $X : \Omega \to E_X$ such that both $\Omega$ and $E_X$ are discrete, meaning they are countable and all of their subsets are measurable. Finally, recall that for a discrete sample space $\Omega$, $\Delta_f(\Omega)$ is the measurable subspace of probability measures $P \in \Delta(\Omega)$ with finite Shannon entropy $H(P)$.

**Lemma A.1.** *Let $X : \Omega \to E_X$ be a random variable. Then the function*

$$X_* : \Delta(\Omega) \to \Delta(E_X), \quad P \mapsto \Big( P_X : A \mapsto P\big(X^{-1}(A)\big) \Big)$$

*is measurable.*

*Proof.* This is elementary and left to the reader to prove.                                   □

To investigate the measurability of the Shannon entropy function and "conditioned" information functions, we need the result that pointwise limits of measurable functions are again measurable:

**Lemma A.2.** *Let $(f_n)_{n \in \mathbb{N}}$ be a sequence of measurable functions $f_n : \mathcal{X} \to \mathbb{R}$ from a measurable space $\mathcal{X}$ to the real numbers $\mathbb{R}$. Assume that the pointwise limit function*

$$f : \mathcal{X} \to \mathbb{R}, \quad x \mapsto \lim_{n \to \infty} f_n(x)$$

*exists. Then $f$ is also measurable.*

*Proof.* See [45], Corollary 8.10.                                                             □

**Corollary A.3.** *Let $X : \Omega \to E_X$ be a discrete random variable. Then the corresponding Shannon entropy function*

$$H(X) : \Delta_f(\Omega) \to \mathbb{R}, \quad P \mapsto H(X; P) := - \sum_{x \in E_X} P_X(x) \ln P_X(x)$$

*is measurable.*

*Proof.* We already know from Lemma A.1 that the function $P \mapsto P_X$ is measurable. Therefore, we can reduce to the case $X = \mathrm{id}_\Omega$, i.e.: we need to show that the function

$$H : \Delta_f(\Omega) \to \mathbb{R}, \quad P \mapsto - \sum_{\omega \in \Omega} P(\omega) \ln P(\omega)$$

is measurable. Note that $P(\omega) = \mathrm{ev}_\omega(P)$. $\mathrm{ev}_\omega$ is measurable by definition of the $\sigma$-algebra on $\Delta_f(\Omega)$. Also, $\ln : \mathbb{R}_{>0} \to \mathbb{R}$ is known to be measurable. Since also limits of measurable functions are measurable by Lemma A.2, the result follows.                                           □

**Lemma A.4.** *Let $X : \Omega \to E_X$ be a discrete random variable and $x \in E_X$ any element. Then the function*

$$(\cdot)|_{X=x} : \Delta(\Omega) \to \Delta(\Omega), \quad P \mapsto P|_{X=x},$$

*with $P|_{X=x}$ defined as in Equation (1), is measurable.*

*Proof.* This is elementary and left to the reader to prove.                                                                   □

**Corollary A.5.** *Let $\Omega$ be a discrete measurable space and $F : \Delta_f(\Omega) \to \mathbb{R}$ a conditionable measurable function, meaning that for all discrete random variables $X : \Omega \to E_X$ and all $P \in \Delta_f(\Omega)$, the series*

$$(X.F)(P) = \sum_{x \in E_X} P_X(x) \cdot F\big(P|_{X=x}\big)$$

*converges unconditionally. Then the function $X.F : \Delta_f(\Omega) \to \mathbb{R}$ is also measurable.*

*Proof.* We have

$$(X.F)(P) = \sum_{x \in E_X} (\mathrm{ev}_x \circ X_*)(P) \cdot \big(F \circ (\cdot)|_{X=x}\big)(P).$$

The result follows from the measurability of $\mathrm{ev}_x : \Delta(E_X) \to \mathbb{R}$, $X_*$ as stated in Corollary A.1, $F$, $(\cdot)_{X=x} : \Delta(\Omega) \to \Delta(\Omega)$ as proven in Lemma A.4, and finally the fact that limits of measurable functions are measurable, see Lemma A.2.                                                                   □

# B   Proofs for Section 2

**Proof 1 for Proposition 2.10.** Let $P : \Omega \to [0,1]$ be any probability measure with finite entropy. Since $Y \precsim X$, there is a function $f_{YX} : E_X \to E_Y$ such that $f_{YX} \circ X = Y$. We obtain

$$
\begin{aligned}
I_1(Y; P) &= - \sum_{y \in E_Y} P\big(Y^{-1}(y)\big) \ln P\big(Y^{-1}(y)\big) \\
&= - \sum_{y \in E_Y} P_X\big(f_{YX}^{-1}(y)\big) \ln P_X\big(f_{YX}^{-1}(y)\big) \\
&= - \sum_{y \in E_Y} \sum_{x \in f_{YX}^{-1}(y)} P_X(x) \ln \sum_{x' \in f_{YX}^{-1}(y)} P_X(x') \\
&\overset{(1)}{\leq} - \sum_{y \in E_Y} \sum_{x \in f_{YX}^{-1}(y)} P_X(x) \ln P_X(x) \\
&= I_1(X; P).
\end{aligned}
$$

In step (1) we use that $-\ln$ is a monotonically decreasing function and $\sum_{x' \in f_{YX}^{-1}(y)} P_X(x') \geq P_X(x)$ for each $x \in f_{YX}^{-1}(y)$.                                                                   □

**Proof 2 for Proposition 2.11.** We have $f_{XY}$ and $f_{YX}$ with $f_{XY} \circ Y = X$ and $f_{YX} \circ X = Y$. For every conditionable measurable function $F : \Delta_f(\Omega) \to \mathbb{R}$ and probability measure $P : \Omega \to \mathbb{R}$, we obtain

$$
\begin{aligned}
(X.F)(P) &= \sum_{x \in \mathrm{im}\, X} P_X(x) F(P|_{X=x}) \\
&\overset{(1)}{=} \sum_{y \in \mathrm{im}\, Y} P_X\big(f_{XY}(y)\big) F\big(P|_{X=f_{XY}(y)}\big) \\
&\overset{(2)}{=} \sum_{y \in \mathrm{im}\, Y} P_Y(y) F(P|_{Y=y}) \\
&= (Y.F)(P).
\end{aligned}
$$

In step (1), we use that $f_{XY} : \mathrm{im}\, Y \to \mathrm{im}\, X$ is a bijection. Step (2) can easily be verified.        □

**Proof 3 for Proposition 2.14.** All required properties follow from Lemma 2.13: first of all, the multiplication $\cdot : M \times M \to M$ is well-defined, i.e., does not depend on the representatives of the

factors $[X]$, $[Y]$ by property 0. We get $[\mathbf{1}] \cdot [X] = [X] = [X] \cdot [\mathbf{1}]$ from property 1. $[X] \cdot [Y] = [Y] \cdot [X]$ follows from property 3. We have $[X] \cdot [X] = [X]$ due to property 4.

We now prove the rule $([X] \cdot [Y]) \cdot [Z] = [X] \cdot ([Y] \cdot [Z])$. For any two random variables $U, V \in \widehat{M}$, we write $Z_{UV} \in \widehat{M}$ for a chosen random variable with $UV \sim Z_{UV}$. Then, we obtain:

$$\big( [X] \cdot [Y] \big) \cdot [Z] = [Z_{Z_{XY}Z}] \overset{(\star)}{=} [Z_{XZ_{YZ}}] = [X] \cdot \big( [Y] \cdot [Z] \big).$$

For step $(\star)$, one uses the equivalence $Z_{Z_{XY}Z} \sim Z_{XZ_{YZ}}$ that follows from Lemma 2.13. $\qquad\square$

# C  Proofs for Section 3

## C.1  Proof of the Generalized Hu Theorem 3.2 and Corollary 3.3

All notation and assumptions are as in Theorem 3.2.

**Lemma C.1.** *Let $\mu$ be the measure given on atoms by Equation (9). For all $I \subseteq [n]$, we have $F_1(X_I) = \mu(\widetilde{X}_I)$.*

*Proof.* This is an application of a version of the inclusion-exclusion principle [7], a special case of the Möbius inversion formula on a poset [52, 3.7.1 Proposition]. It states the following: For any two functions $f, g : 2^{[n]} \to G$, the following implication holds:

$$g(I) = \sum_{K \supseteq I} f(K) \quad \Longrightarrow \quad f(I) = \sum_{K \supseteq I} (-1)^{|K|-|I|} g(K).$$

Set $\mu(p_\emptyset) := -F_1(X_{[n]})$. We apply the principle to the functions $g(I) := (-1)^{|I|} \mu(p_{I^c})$ and $f(K) := (-1)^{|K|+1} F_1(X_K)$. Then Equation (9) implies the premise of the inclusion-exclusion principle. From the conclusion, we obtain:

$$(-1)^{|I|+1} F_1(X_I) = \sum_{K \supseteq I} (-1)^{|K|-|I|} \cdot (-1)^{|K|} \mu(p_{K^c}),$$

which implies

$$F_1(X_I) = - \sum_{K \supseteq I} \mu(p_{K^c}) = - \sum_{K:\, K \cap I = \emptyset} \mu(p_K) = F_1(X_{[n]}) - \sum_{\emptyset \neq K:\, K \cap I = \emptyset} \mu(p_K).$$

In the last step we used $\mu(p_\emptyset) = -F_1(X_{[n]})$. Thus, showing that $F_1(X_I) = \mu(\widetilde{X}_I) = \sum_{\emptyset \neq K:\, K \cap I \neq \emptyset} \mu(p_K)$ reduces to the following special case for $I = [n]$:

$$F_1(X_{[n]}) = \mu\big(\widetilde{X}_{[n]}\big).$$

To show this, note that Equation (9) implies

$$\mu\big(\widetilde{X}_{[n]}\big) = \sum_K (-1)^{|K|+1-n} \left[ \sum_{\emptyset \neq I:\, I \supseteq K^c} (-1)^{|I|} \right] F_1(X_K). \tag{31}$$

Ignoring that $\emptyset \neq I$, the inner coefficient is given by

$$\sum_{I \supseteq K^c} (-1)^{|I|} = (-1)^{n-|K|} \sum_{i=0}^{|K|} (-1)^i \binom{|K|}{i} = \begin{cases} 0, & K \neq \emptyset \\ (-1)^n, & \text{else.} \end{cases}$$

Note that $F_1(X_\emptyset) = 0$, so the last case is irrelevant. Also, note that the condition $I \neq \emptyset$ only restricts the inner sum in Equation (31) when $K = [n]$. Thus, in that case, we need to subtract 1 from the result just computed and obtain:

$$\mu\big(\widetilde{X}_{[n]}\big) = (-1)^{n+1-n} \cdot (-1) \cdot F_1(X_{[n]}) = F_1(X_{[n]}),$$

proving the claim. $\qquad\square$

**Proposition C.2.** *For all $n \in \mathbb{N}_{\geq 0}$, for $\mu$ being the $G$-valued measure constructed from $F_1$ as in Equation (9), for all $L_1, J \subseteq [n]$, the following identity holds:*

$$X_J.F_1(X_{L_1}) = \mu(\widetilde{X}_{L_1} \setminus \widetilde{X}_J)$$

*Proof.* This follows immediately from Lemma C.1 and the chain rule, Equation (6), together with the fact that $\mu$ is a measure.                                                                  □

We have now done all the hard work for finishing the proof of Theorem 3.2:

**Proof 4 for Theorem 3.2. Part 1.**          This is a simple inductive argument, using Proposition C.2 for $q = 1$, and Equation (7) for showing the step from $q - 1$ to $q$.

**Part 2.**      For part 2, using Equation (8), we observe

$$X_J.F_1(X_I) - F_1(X_{J \cup I}) + F_1(X_J) = \mu(\widetilde{X}_I \setminus \widetilde{X}_J) - \mu(\widetilde{X}_J \cup \widetilde{X}_I) + \mu(\widetilde{X}_J) = 0.$$

Thus, $F_1$ satisfies Equation (6). For $q \geq 2$, using Equation (8) again, we similarly observe

$$F_{q-1}(X_{L_1}; \ldots; X_{L_{q-1}}) - X_{L_q}.F_{q-1}(X_{L_1}; \ldots; X_{L_{q-1}}) = F_q(X_{L_1}; \ldots; X_{L_q}).$$

That finishes the proof.                                                                                        □

**Proof 5 for Corollary 3.3.** Define $\widetilde{G} := \mathrm{Maps}(M, G)$ as the group of functions from $M$ to $G$. Define, using *currying*, the function $\widetilde{K}_1 : M \to \widetilde{G}$ by

$$\big[\widetilde{K}_1(X)\big](Y) := K_1(X \mid Y).$$

Define the additive monoid action $. : M \times \widetilde{G} \to \widetilde{G}$ by

$$(X.F)(Y) := F(XY)$$

for all $X, Y \in M$. Note that we need $M$ to be commutative to show that this is indeed a monoid action. Then clearly, $\widetilde{K}_1$ satisfies the chain rule.

Define $\widetilde{K}_q : M^q \to \widetilde{G}$ as in Theorem 3.2 inductively by

$$\widetilde{K}_q(Y_1; \ldots; Y_q) := \widetilde{K}_{q-1}(Y_1; \ldots; Y_{q-1}) - Y_q.\widetilde{K}_{q-1}(Y_1; \ldots; Y_{q-1}).$$

By induction, one can show that $K_q(Y_1; \ldots; Y_q \mid Z) = \big[\widetilde{K}_q(Y_1; \ldots; Y_q)\big](Z)$ for all $Y_1, \ldots, Y_q, Z \in M$.

By the conclusion of Theorem 3.2, we obtain a $\widetilde{G}$-valued measure $\widetilde{\mu} : 2^{\widetilde{X}} \to \widetilde{G}$ with

$$\widetilde{\mu}\left(\bigcap_{k=1}^q \widetilde{X}_{L_k} \setminus \widetilde{X}_J\right) = X_J.\widetilde{K}_q(X_{L_1}; \ldots; X_{L_q}).$$

Now, define $\mu : 2^{\widetilde{X}} \to G$ by $\mu(A) := \big[\widetilde{\mu}(A)\big](1)$ for all $A \subseteq \widetilde{X}$. Clearly, since $\widetilde{\mu}$ is a $\widetilde{G}$-valued measure, $\mu$ is a $G$-valued measure. The results immediately follow from these definitions and Hu's Theorem.                                                                                        □

## C.2   Further Proofs for Section 3

**Proof 6 for Corollary 3.5.** We proceed as follows:

1. This follows from Lemma 3.4 and Equation (9).

2. By Lemma 3.4 and Theorem 3.2, we have

$$\sum_{\substack{I \subseteq [n] \\ I \cap K \neq \emptyset}} \eta_I = \mu\Big(\Big\{p_I \ \Big| \ I \subseteq [n], \exists k \in K : k \in I\Big\}\Big) = \mu(\widetilde{X}_K) = F_1(X_K).$$

3. Using Lemma 3.4 and Theorem 3.2 again, we obtain

$$F_q(X_{j_1}; \dots; X_{j_q}) = \mu\left(\bigcap_{j \in J} \widetilde{X}_j\right) = \sum_{I, \forall j \in J : j \in I} \eta_I = \sum_{I \supseteq J} \eta_I.$$

4. This is formally a consequence of 3 and the inclusion-exclusion principle [7].

5. This follows by combining results 2 and 4.

6. This follows by combining results 1 and 3, or by the inclusion-exclusion principle [7] applied to result 5.

<div style="text-align: right">□</div>

# D   Proofs for Section 4

**Proof 7 for Proposition 4.5.** Let $Y, Z \in \widetilde{M}$ be arbitrary. In the following, for functions $f : (\{0,1\}^*)^n \to \mathbb{R}$, we write $f = f(\boldsymbol{x})$ for simplicity, and mean by it the function mapping $\boldsymbol{x}$ to $f(\boldsymbol{x})$. We obtain:

$$\begin{aligned}
Kc(YZ) &= Kc\big((YZ)(\boldsymbol{x})\big) \\
&\stackrel{\pm}{=} Kc\big(Y(\boldsymbol{x})'Z(\boldsymbol{x})\big) \\
&\stackrel{\pm}{=} Kc\big(Y(\boldsymbol{x})\big) + Kc\big(Z(\boldsymbol{x}) \mid Y(\boldsymbol{x})\big) \\
&\stackrel{\pm}{=} Kc(Y) + Kc(Z \mid Y).
\end{aligned}$$

In the computation, the associativity rule in the second step holds as we can write a program of constant size that translates between the different nestings of the strings.[14] In the third step we use Theorem 4.4. The result follows. □

**Proof 8 for Lemma 4.6.** We have

$$\begin{aligned}
[Kc]_{Kc}(Y \mid Z) &\stackrel{(1)}{=} [Kc]_{Kc}(YZ) - [Kc]_{Kc}(Z) \\
&\stackrel{(2)}{=} [Kc]_{Kc}(\overline{Y}\ \overline{Z}) - [Kc]_{Kc}(\overline{Z}) \\
&\stackrel{(3)}{=} [Kc]_{Kc}(\overline{Y} \mid \overline{Z}).
\end{aligned}$$

In the computation, steps (1) and (3) follow from Proposition 4.5. For step (2) one can show that $Kc(YZ) \stackrel{\pm}{=} Kc(\overline{Y}\ \overline{Z})$ and $Kc(Z) \stackrel{\pm}{=} Kc(\overline{Z})$ in the same way as the associativity rule in Proposition 4.5 was shown. □

**Proof 9 for Theorem 4.8.** Remember $M = \widetilde{M}/\sim$ and the function $[Kc]_{Kc} : M \times M \to$ Maps$\big((\{0,1\}^*)^n, \mathbb{R}\big)/\sim_{Kc}$, which we now denote by $[Kc] = [Kc]_1 := [Kc]_{Kc}$. From this, we can inductively define $[Kc]_q : M^q \times M \to$ Maps$\big((\{0,1\}^*)^n, \mathbb{R}\big)/\sim_{Kc}$ as in Corollary 3.3 by

$$[Kc]_q(Y_1; \dots; Y_q \mid Z) := [Kc]_{q-1}(Y_1; \dots; Y_{q-1} \mid Z) - [Kc]_{q-1}(Y_1; \dots; Y_{q-1} \mid Y_q Z).$$

---

[14]For this, we use that we can algorithmically extract all $x_i$ for indices appearing in $Y$ and $Z$ from the strings $(YZ)(\boldsymbol{x})$ and also $Y(\boldsymbol{x})'Z(\boldsymbol{x})$. This argument uses that the encoding $x \mapsto x'$ is prefix-free.

From Equation (21), one can inductively show that

$$[Kc]_q(Y_1; \ldots; Y_q \mid Z) = [Kc_q(Y_1; \ldots; Y_q \mid Z)] \tag{32}$$

for all $Y_1, \ldots, Y_q, Z \in M$. Note that $Kc_q$ was defined on $\widetilde{M}$ and not $M$, which means that we plug in representatives of equivalence classes at the right-hand-side. Using Lemma 4.6 and induction, one can show that this is well-defined. Then, construct $\mu : 2^{\widetilde{X}} \to \mathrm{Maps}\left(((\{0,1\}^*)^n, \mathbb{R}\right)$ explicitly as in Equation (23). Define, now, the measure $[\mu] : 2^{\widetilde{X}} \to \mathrm{Maps}\left(((\{0,1\}^*)^n, \mathbb{R}\right)/ \sim_{Kc}$ by

$$[\mu](A) := [\mu(A)] \quad \forall A \subseteq \widetilde{X}. \tag{33}$$

Then Equation (32) shows that

$$[\mu](p_I) = \sum_{\emptyset \neq K \supseteq I^c} (-1)^{|K|+|I|+1-n} [Kc]_1(X_K). \tag{34}$$

Consequently, $[\mu]$ is the measure that results in Corollary 3.3, see Equation (13). We obtain for all $L_1, \ldots, L_q, J \subseteq [n]$:

$$\begin{aligned}
\left[\mu\left(\bigcap_{k=1}^q \widetilde{X}_{L_k} \setminus \widetilde{X}_J\right)\right] &= [\mu]\left(\bigcap_{k=1}^q \widetilde{X}_{L_k} \setminus \widetilde{X}_J\right) && \text{(Equation (33))} \\
&= [Kc]_q(X_{L_1}; \ldots; X_{L_q} \mid X_J) && \text{(Proposition 4.5, Corollary 3.3)} \\
&= \left[Kc_q(X_{L_1}; \ldots; X_{L_q} \mid X_J)\right] && \text{(Equation (32))}.
\end{aligned}$$

As two representatives of the same equivalence class in $\mathrm{Maps}\left(((\{0,1\}^*)^n, \mathbb{R}\right)$ differ by a constant, the result follows. $\qquad \square$

**Lemma D.1.** *Let $P : (\{0,1\}^*)^n \to \mathbb{R}$ be a computable probability mass function. Let $K \subseteq [n]$ a subset and $P_K$ the corresponding maginal distribution. Then $P_K$ is also computable, and the relation*

$$K(P_K) \overset{+}{<} K(P).$$

*between their Kolmogorov complexities holds.*

*Proof.* We know that $P$ is computable, and so there exists a prefix-free Turing machine $T_p$ of length $l(p) = K(P)$ such that

$$\left|T_p(\boldsymbol{x}'q) - P(\boldsymbol{x})\right| \leq 1/q$$

for all $q \in \mathbb{N}$ and $\boldsymbol{x} \in (\{0,1\}^*)^n$. Now, fix $q \in \mathbb{N}$. Let $(\boldsymbol{x}^i)_{i \in \mathbb{N}}$ be a computable enumeration of $(\{0,1\}^*)^n$. Define the approximation $P_q : (\{0,1\}^*)^n \dashrightarrow \mathbb{R}$ of $P$ by

$$P_q(\boldsymbol{x}^i) := T_p\left((\boldsymbol{x}^i)'(4q \cdot 2^i)\right).$$

Then for all subsets $I \subseteq \mathbb{N}$, we have

$$\begin{aligned}
\left|\sum_{i \in I} P_q(\boldsymbol{x}^i) - \sum_{i \in I} P(\boldsymbol{x^i})\right| &\leq \sum_{i \in I} \left|T_p\left((\boldsymbol{x}^i)'(4q \cdot 2^i)\right) - P(\boldsymbol{x}^i)\right| \\
&\leq \sum_{i=1}^{\infty} \frac{1}{4q \cdot 2^i} \\
&= \frac{1}{4q}.
\end{aligned} \tag{35}$$

Consequently, by setting $I = \mathbb{N}$ and using $\sum_{i \in \mathbb{N}} P(\boldsymbol{x}^i) = 1$, one can determine $i_q$ such that for the first time we have

$$\left|\sum_{i=1}^{i_q} P_q(\boldsymbol{x}^i) - 1\right| \leq \frac{1}{2q}. \tag{36}$$

Note that $i_q$ can be *algorithmically* determined by computing one $P_q(\boldsymbol{x}^i)$ at a time and checking when the condition holds. Now, for arbitrary $\boldsymbol{x}_K \in (\{0,1\}^*)^{|K|}$ and $q \in \mathbb{N}$, we define

$$T(\boldsymbol{x}'_K q) := \sum_{\substack{i=1 \\ (\boldsymbol{x}^i)_K = \boldsymbol{x}_K}}^{i_q} P_q(\boldsymbol{x}^i).$$

We now show that $T(\boldsymbol{x}'_K q)$ approximates $P_K(\boldsymbol{x}_K)$ up to an error of $1/q$:

$$\begin{aligned}
\left| T(\boldsymbol{x}'_K q) - P_K(\boldsymbol{x}_K) \right| &= \left| \sum_{\substack{i=1 \\ (\boldsymbol{x}^i)_K = \boldsymbol{x}_K}}^{i_q} P_q(\boldsymbol{x}^i) - P_K(\boldsymbol{x}_K) \right| \\
&\leq \left| \sum_{\substack{i=1 \\ (\boldsymbol{x}^i)_K = \boldsymbol{x}_K}}^{i_q} P_q(\boldsymbol{x}^i) - \sum_{\substack{i=1 \\ (\boldsymbol{x}^i)_K = \boldsymbol{x}_K}}^{i_q} P(\boldsymbol{x}^i) \right| + \left| \sum_{\substack{i=1 \\ (\boldsymbol{x}^i)_K = \boldsymbol{x}_K}}^{i_q} P(\boldsymbol{x}^i) - P_K(\boldsymbol{x}_K) \right| \\
&\overset{(35)}{\leq} \frac{1}{4q} + P_K(\boldsymbol{x}_K) - \sum_{\substack{i=1 \\ (\boldsymbol{x}^i)_K = \boldsymbol{x}_K}}^{i_q} P(\boldsymbol{x}^i) \\
&= \frac{1}{4q} + \left| 1 - \sum_{i=1}^{i_q} P(\boldsymbol{x}_i) \right| \\
&\leq \frac{1}{4q} + \left| 1 - \sum_{i=1}^{i_q} P_q(\boldsymbol{x}^i) \right| + \left| \sum_{i=1}^{i_q} P_q(\boldsymbol{x}^i) - \sum_{i=1}^{i_q} P(\boldsymbol{x}^i) \right| \\
&\overset{(36),(35)}{\leq} \frac{1}{4q} + \frac{1}{2q} + \frac{1}{4q} = 1/q.
\end{aligned}$$

Now, note that $T$ is computable, since it uses in its definition only the computable enumeration $(\boldsymbol{x}^i)_{i \in \mathbb{N}}$, the number $i_q$ for which we described an algorithm, and the Turing machine $T_p$ inside the definition of $P_q$. Thus, $T$ is a prefix machine $T_{p_K}$ for a bitstring $p_K$ of length $l(p_K) \leq l(p) + c = K(P) + c$, where $c \geq 0$ is some constant. It follows $K(P_K) \leq l(p_K) \leq K(P) + c$, and we are done. $\qquad\square$

**Proof 10 for Lemma 4.12.** Assume that $Y \sim Z$. Then Lemma 2.13 parts 3 and 4[15] show that $Y \sim_r \overline{Y} = \overline{Z} \sim_r Z$, and so $Y \sim_r Z$ by transitivity.

On the other hand, if $Y \sim_r Z$, then also $X_I = \overline{Y} \sim_r \overline{Z} = X_J$ for some $I, J \subseteq [n]$, again by Lemma 2.13 parts 3 and 4. Let $I = \{i_1 < \cdots < i_{|I|}\}$ and $J = \{j_1 < \cdots < j_{|J|}\}$. Then there exist functions $f_{JI}$ and $f_{IJ}$ such that $f_{JI} \circ X_I = X_J$ and $f_{IJ} \circ X_J = X_I$. That is, for all $\boldsymbol{x} \in (\{0,1\}^*)^n$ we have

$$\begin{aligned}
f_{JI}(x_{i_1}, \ldots, x_{i_{|I|}}) &= (x_{j_1}, \ldots, x_{j_{|J|}}), \\
f_{IJ}(x_{j_1}, \ldots, x_{j_{|J|}}) &= (x_{i_1}, \ldots, x_{i_{|I|}}).
\end{aligned}$$

The first equation shows $J \subseteq I$, as otherwise, changes in $\boldsymbol{x}_{J \setminus I}$ lead to changes in the right-hand-side, but not the left-hand-side. In the same way, the second equation shows $I \subseteq J$, and overall we obtain $I = J$. That shows $Y \sim \overline{Y} = X_I = X_J = \overline{Z} \sim Z$; due to transitivity, it follows $Y \sim Z$. $\qquad\square$

**Proof 11 for Theorem 4.14.** We generalize the proof strategy that [34] use for their Lemma 8.1.1, which is a special case of our theorem for $n = 2$, $q = 2$, $Y_1 = X_1, Y_2 = X_2$, and $Z = \epsilon = \boldsymbol{1}$.

---

[15]What we denoted by $\sim$ in that lemma is denoted $\sim_r$ here.

We prove this in several steps by first handling convenient subcases. In the special case $q = 1$, $Z = \epsilon = \mathbf{1}$, and $Y_1 = X_K$ for some $K \subseteq [n]$, we can look at the marginal $P_K$ of $P$ and obtain

$$\sum_{\boldsymbol{x} \in (\{0,1\}^*)^n} P(\boldsymbol{x})\big(Kc(X_K)\big)(\boldsymbol{x}) = \sum_{\boldsymbol{x}_K \in (\{0,1\}^*)^{|K|}} P_K(\boldsymbol{x}_K)\big(Kc(X_K)\big)(\boldsymbol{x}_K)$$

$$= I_1(P_K) + O\big(K(P_K)\big) \qquad\qquad \text{(Theorem 4.13)}$$

$$= I_1(X_K; P) + O\big(K(P)\big), \qquad\qquad \text{(Lemma D.1)},$$

which is the wished result. Now, let

$$\mu : 2^{\widetilde{X}} \to \text{Maps}\,\big((\{0,1\}^*)^n, \mathbb{R}\big), \qquad\qquad \text{(Equation (23))}$$

$$\mu^r : 2^{\widetilde{X}} \to \text{Meas}_{\text{con}}\,\Big(\Delta_f\big((\{0,1\}^*)^n\big), \mathbb{R}\Big) \qquad\qquad \text{(Equation 9)}$$

be the measures corresponding to Chaitin's prefix-free Kolmogorov complexity $Kc : M \times M \to$ Maps $\big((\{0,1\}^*)^n, \mathbb{R}\big)$ and Shannon entropy $I_1 : M \to \text{Meas}_{\text{con}}\,\Big(\Delta_f\big((\{0,1\}^*)^n\big), \mathbb{R}\Big)$, remembering that $\Delta_f\big((\{0,1\}^*)^n\big)$ is the space of finite-entropy probability measures (or mass functions) on our *countable*[16] sample space $(\{0,1\}^*)^n$.[17] Let $I \subseteq [n]$ be any subset. Then we obtain

$$\sum_{\boldsymbol{x} \in (\{0,1\}^*)^n} P(\boldsymbol{x})\big(\mu(p_I)\big)(\boldsymbol{x}) \overset{(23)}{=} \sum_{\boldsymbol{x} \in (\{0,1\}^*)^n} P(\boldsymbol{x}) \sum_{\emptyset \neq K \supseteq I^c} (-1)^{|K|+|I|+1-n}\big(Kc(X_K)\big)(\boldsymbol{x})$$

$$= \sum_{\emptyset \neq K \supseteq I^c} (-1)^{|K|+|I|+1-n} \sum_{\boldsymbol{x} \in (\{0,1\}^*)^n} P(\boldsymbol{x})\big(Kc(X_K)\big)(\boldsymbol{x})$$

$$\overset{(\star)}{=} \sum_{\emptyset \neq K \supseteq I^c} (-1)^{|K|+|I|+1-n}\Big(I_1(X_K; P) + O\big(K(P)\big)\Big)$$

$$= \Bigg(\sum_{\emptyset \neq K \supseteq I^c} (-1)^{|K|+|I|+1-n} I_1(X_K)\Bigg)(P) + O\big(K(P)\big)$$

$$= \big(\mu^r(p_I)\big)(P) + O\big(K(P)\big),$$

using our earlier result in step $(\star)$ and the definition of $\mu^r$. Now, using that $\mu$ and $\mu^r$ are additive over disjoint unions, we can deduce for all $A \subseteq \widetilde{X}$ the equality

$$\sum_{\boldsymbol{x} \in (\{0,1\}^*)^n} P(\boldsymbol{x})\big(\mu(A)\big)(\boldsymbol{x}) = \big(\mu^r(A)\big)(P) + O\big(K(P)\big).$$

Now, let $Y_1 = X_{L_1}, \ldots, Y_q = X_{L_q}, Z = X_J$ for some $L_1, \ldots, L_q, J \subseteq [n]$. Then, using Hu's theorems for interaction information (Theorem 3.2 together with Summary 2.17) and Kolmogorov complexity 4.8, the result follows by setting $A := \bigcap_{k=1}^q \widetilde{X}_{L_k} \setminus \widetilde{X}_J$. $\qquad\square$

**Proof 12 for Proposition 4.19.** We have

$$K(YZ) = K\big((YZ)(\boldsymbol{x})\big)$$

$$\overset{(1)}{=} K\big(Y(\boldsymbol{x})'Z(\boldsymbol{x})\big) + O(1)$$

$$\overset{(2)}{=} K\big(Y(\boldsymbol{x})\big) + K\big(Z(\boldsymbol{x}) \mid Y(\boldsymbol{x})\big) + O\Big(\log K\big(Y(\boldsymbol{x})\big) + \log K\big(Z(\boldsymbol{x})\big)\Big)$$

$$\overset{(3)}{=} K(Y) + K(Z \mid Y) + O\Bigg(\sum_{i=1}^n \log K(x_i)\Bigg).$$

---

[16]The fact that $(\{0,1\}^*)^n$ is not finite but countably infinite is the main reason why we considered countable sample spaces in Summary 2.17.

[17]The superscript in $\mu^r$ is used to notationally distinguish it from $\mu$. $r$ can be thought of as meaning "random".

where step (1) follows as in Proposition 4.5, step (2) uses Theorem 4.18, and step (3) follows from the subadditivity of $K$[18] and the logarithm, which holds for large enough inputs. $\qquad\square$

**Proof 13 for Proposition 4.21.** Let $Y, Z \in M$ be arbitrary. Then, following the same arguments as in Proposition 4.5 and Proposition 4.19, we are only left with showing the following:

$$\log C\big(Y(\boldsymbol{x}), Z(\boldsymbol{x})\big) = O\big(\log C(\boldsymbol{x})\big),$$

where the left-hand-side is viewed as a function $(\{0,1\}^*)^n \to \mathbb{R}$. In fact, we even have

$$\log C\big(Y(\boldsymbol{x}), Z(\boldsymbol{x})\big) \le \log C(\boldsymbol{x}) + c$$

for some constant $c$ starting from some threshold $\boldsymbol{x}_0$: we can find a program in constant length that takes $\boldsymbol{x}$, extracts $x_1, \ldots, x_n$ from it, and rearranges and concatenates them in such an order to obtain $Y(\boldsymbol{x})'Z(\boldsymbol{x})$, and the logarithm, being subadditive for large enough inputs, preserves the inequality. $\qquad\square$

# E   Proofs for Section 5

**Proof 14 for Proposition 5.1.** For notational ease, we write $P(x) = P_X(x)$, $(P|_{X=x})_Y(y) = P(y \mid x)$ and $P(x,y) = P_{XY}(x,y)$ in this proof. We have

$$\big[I_1^q(X) + X._q I_1^q(Y)\big](P) = \big[I_1^q(X)\big](P) + \sum_{x \in E_X} P(x)^q \big[I_1^q(Y)\big](P|_{X=x})$$

$$= \frac{\sum_{x \in E_X} P(x)^q - 1}{1 - q} + \sum_{x \in E_X} P(x)^q \frac{\sum_{y \in E_Y} P(y \mid x)^q - 1}{1 - q}$$

$$= \frac{\sum_{x \in E_X} P(x)^q - 1 + \sum_{(x,y) \in E_X \times E_Y} \big(P(x)P(y \mid x)\big)^q - \sum_{x \in E_X} P(x)^q}{1 - q}$$

$$= \frac{\sum_{(x,y) \in E_X \times E_Y} P(x,y)^q - 1}{1 - q}$$

$$= \big[I_1^q(XY)\big](P).$$

$\qquad\square$

**Proof 15 for Proposition 5.2.** Let $X, Y \in M(X_1, \ldots, X_n)$ and $P \ll Q \in \Delta(\Omega)$. The following proof of the chain rule is similar to the one of Lemma 2.4 for Shannon entropy. For simplicity, we write $Q(x) = Q_X(x)$, $P(y \mid x) = (P|_{X=x})_Y(y)$ and $P(x,y) = P_{XY}(x,y)$ in this proof:

$$\big[X.D_1(Y) + D_1(X)\big](P\|Q)$$

$$= X.D_1(Y; P\|Q) + D_1(X; P\|Q)$$

$$= \sum_{x \in E_X} P(x) D_1(Y; P|_{X=x} \| Q|_{X=x}) - \sum_{x \in E_X} P(x) \ln \frac{Q(x)}{P(x)}$$

$$\overset{(1)}{=} -\sum_{x \in E_X} P(x) \sum_{y \in E_Y} P(y \mid x) \ln \frac{Q(y \mid x)}{P(y \mid x)} - \sum_{x \in E_X} P(x) \left(\sum_{y \in E_Y} P(y \mid x)\right) \ln \frac{Q(x)}{P(x)}$$

---

[18]The subadditivity property for $K$ says that $K(x,y) \le K(x) + K(y) + O(1)$: one can construct a prefix-free Turing machine that extracts $x^*$ and $y^*$ from $x^*y^*$, which is of length $K(x) + K(y)$, and outputs $x'y$. Note that since the set of halting programs of the universal Turing machine $U$ is *prefix-free*, one does not need to indicate the place of separation between $x^*$ and $y^*$.

$$= - \sum_{(x,y)\in E_X \times E_Y} P(x) \cdot P(y \mid x) \cdot \left[ \ln \frac{Q(y \mid x)}{P(y \mid x)} + \ln \frac{Q(x)}{P(x)} \right]$$

$$= - \sum_{(x,y)\in E_X \times E_Y} P(x,y) \ln \frac{Q(x,y)}{P(x,y)}$$

$$= \big[ D_1(XY) \big](P\|Q).$$

In step (1), we used for the second sum that $P(y \mid x)$ is a probability measure in $y$ and thus sums to 1. $\qquad \square$

**Proof 16 for Proposition 5.4.** Let $X, Y \in \mathrm{M}(X_1, \ldots, X_n)$ and $P \ll Q \in \Delta(\Omega)$ be arbitrary. The following proof of the chain rule is similar to the one for the $q$-entropy, Proposition 5.1. For simplicity, we write $Q(x) = Q_X(x)$, $P(y \mid x) = (P|_{X=x})_Y(y)$ and $P(x,y) = P_{XY}(x,y)$ in this proof:

$$\big[ D_1^q(X) + X._q D_1^q(Y) \big](P\|Q) = \big[ D_1^q(X) \big](P\|Q) + \big[ X._q D_1^q(Y) \big](P\|Q)$$

$$= \big[ D_1^q(X) \big](P\|Q) + \sum_{x\in E_X} P(x)^q Q(x)^{1-q} \big[ D_1^q(Y) \big]\big( P|_{X=x} \| Q|_{X=x} \big)$$

$$= \frac{\sum_{x\in E_X} P(x)^q Q(x)^{1-q} - 1}{q-1} + \sum_{x\in E_X} P(x)^q Q(x)^{1-q} \frac{\sum_{y\in E_Y} P(y \mid x)^q Q(y \mid x)^{1-q} - 1}{q-1}$$

$$= \frac{-1 + \sum_{(x,y)\in E_X \times E_Y} \big( P(x)P(y \mid x) \big)^q \big( Q(x)Q(y \mid x) \big)^{1-q}}{q-1}$$

$$= \frac{\sum_{(x,y)\in E_X \times E_Y} P(x,y)^q Q(x,y)^{1-q} - 1}{q-1}$$

$$= \big[ D_1^q(XY) \big](P\|Q).$$

$\qquad \square$

# References

[1] S.-I. Amari. Information Geometry on Hierarchy of Probability Distributions. *IEEE Transactions on Information Theory*, 47(5):1701–1711, 2001. https://doi.org/10.1109/18.930911.

[2] Valentina Baccetti and Matt Visser. Infinite Shannon entropy. *Journal of Statistical Mechanics: Theory and Experiment*, 2013(4), 2013. ISSN 17425468. https://doi.org/10.1088/1742-5468/2013/04/P04010.

[3] Pierre Baudot. The Poincare-Shannon Machine: Statistical Physics and Machine Learning Aspects of Information Cohomology. *Entropy*, 21(9), sep 2019. ISSN 10994300. https://doi.org/10.3390/E21090881.

[4] Pierre Baudot. On Information Links. *arXiv e-prints*, page arXiv:2103.02002, 2021. https://doi.org/10.48550/arXiv.2103.02002.

[5] Pierre Baudot and Daniel Bennequin. The Homological Nature of Entropy. *Entropy*, 17(5): 3253–3318, 2015. ISSN 10994300. https://doi.org/10.3390/e17053253.

[6] Pierre Baudot, Monica Tapia, Daniel Bennequin, and Jean Marc Goaillard. Topological Information Data Analysis. *Entropy*, 21(9):1–38, 2019. ISSN 10994300. https://doi.org/10.3390/e21090869.

[7] Robert A. Beeler. *How to Count: An Introduction to Combinatorics and Its Applications.* Springer International Publishing, 2015. ISBN 9783319138435. https://doi.org/10.1007/978-3-319-13844-2.

[8] Anthony J Bell. THE CO-INFORMATION LATTICE. *in Proc. 4th Int. Symp. Independent Component Analysis and Blind Source Separation*, pages 921–926, 2003.

[9] Daniel Bennequin, Olivier Peltre, Grégoire Sergeant-Perthuis, and Juan Pablo Vigneaux. Extra-Fine Sheaves and Interaction Decompositions. *arXiv e-prints*, page arXiv:2009.12646, sep 2020. https://doi.org/10.48550/arXiv.2009.12646.

[10] Nils Bertschinger, Johannes Rauh, Eckehard Olbrich, Jürgen Jost, and Nihat Ay. Quantifying Unique Information. *Entropy*, 16(4):2161–2183, apr 2014. https://doi.org/10.3390/e16042161.

[11] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, 1 edition, 2007. ISBN 0387310738. https://doi.org/10.1117/1.2819119.

[12] Christopher Michael Bishop and Hugh Bishop. *Deep Learning - Foundations and Concepts.* 1 edition, 2023. ISBN 978-3-031-45468-4. https://doi.org/10.1007/978-3-031-45468-4.

[13] N. J. Cerf and C. Adami. Entropic Bell Inequalities. *Physical Review A - Atomic, Molecular, and Optical Physics*, 55(5):3371–3374, 1997. ISSN 10941622. https://doi.org/10.1103/PhysRevA.55.3371.

[14] N. J. Cerf and C. Adami. Quantum extension of conditional probability. *Physical Review A - Atomic, Molecular, and Optical Physics*, 60(2):893–897, 1999. ISSN 10941622. https://doi.org/10.1103/PhysRevA.60.893.

[15] Gregory. J. Chaitin. *Algorithmic Information Theory.* Cambridge University Press, oct 1987. ISBN 9780521343060. https://doi.org/10.1017/CBO9780511608858.

[16] S. Cocco and R. Monasson. Adaptive Cluster Expansion for the Inverse Ising Problem: Convergence, Algorithm and Tests. *Journal of Statistical Physics*, 147(2):252–314, April 2012. https://doi.org/10.1007/s10955-012-0463-4.

[17] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory.* Wiley-Interscience, 2006. ISBN 9780471241959. https://doi.org/10.1002/047174882X.

[18] A P Dawid. Separoids: A Mathematical Framework for Conditional Independence and Irrelevance. *Annals of Mathematics and Artificial Intelligence*, 32(1):335–372, 2001. ISSN 1573-7470. https://doi.org/10.1023/A:1016734104787.

[19] A. Philip Dawid. Conditional Independence in Statistical Theory. *Journal of the Royal Statistical Society. Series B (Methodological)*, 41(1):1–31, 1979. ISSN 00359246. https://doi.org/10.1111/j.2517-6161.1979.tb01052.x.

[20] A. Philip Dawid. Conditional Independence for Statistical Operations. *The Annals of Statistics*, 8, 1980. https://doi.org/10.1214/aos/1176345011.

[21] Hubert Dubé. *On the Structure of Information Cohomology.* Phd dissertation, University of Toronto, 2023. URL https://hdl.handle.net/1807/130552.

[22] Jack Edmonds. *Submodular functions, matroids, and certain polyhedra*, page 11–26. Springer-Verlag, Berlin, Heidelberg, 2003. ISBN 3540005803. https://doi.org/10.1007/3-540-36478-1_2.

[23] Kun Fang, Omar Fawzi, Renato Renner, and David Sutter. Chain Rule for the Quantum Relative Entropy. *Physical Review Letters*, 124(10):100501, 2020. ISSN 10797114. https://doi.org/10.1103/PhysRevLett.124.100501.

[24] Conor Finn and Joseph Lizier. Pointwise Partial Information Decomposition Using the Specificity and Ambiguity Lattices. *Entropy*, 20(4):297, apr 2018. https://doi.org/10.3390/e20040297.

[25] James Fullwood. On a 2-Relative Entropy. *Entropy*, 24(1):74, dec 2021. https://doi.org/10.3390/e24010074.

[26] Peter Gacs. On the Symmetry of Algorithmic Information. *Soviet Math. Dokl.*, 15(January 1974):1477–1480, 1974.

[27] Marilyn Gatica, Rodrigo Cofré, Pedro A.M. Mediano, Fernando E. Rosas, Patricio Orio, Ibai Diez, Stephan P. Swinnen, and Jesus M. Cortes. High-Order Interdependencies in the Aging Brain. *Brain Connectivity*, 11(9):734–744, 2021. ISSN 21580022. https://doi.org/10.1089/brain.2020.0982.

[28] Peter D. Grünwald and Paul M B Vitányi. Algorithmic Information Theory. *arXiv e-prints*, abs/0809.2:arXiv:0809.2754, sep 2008. ISSN 21971765. https://doi.org/10.1007/978-981-15-0739-7_2.

[29] Te Sun Han. *Information and Control*, 36(2):133–156, 1978. ISSN 0019-9958. https://doi.org/10.1016/S0019-9958(78)90275-9.

[30] Gerhard Hochschild. On the Cohomology Groups of an Associative Algebra. *Annals of Mathematics*, 46(1):58–67, 1945. ISSN 0003486X. https://doi.org/10.2307/1969145.

[31] Michal Michał Horodecki, Jonathan Oppenheim, and Andreas Winter. Partial quantum information. *Nature*, 436(7051):673–676, 2005. ISSN 1476-4687. https://doi.org/10.1038/nature03909.

[32] Kuo Ting Hu. On the Amount of Information. *Theory of Probability & Its Applications*, 7(4): 439–447, 1962. https://doi.org/10.1137/1107041.

[33] Ryan G. James and James P. Crutchfield. Multivariate dependence beyond Shannon information. *Entropy*, 19(10), 2017. ISSN 10994300. https://doi.org/10.3390/e19100531.

[34] Ming Li and Paul Vitányi. *An Introduction to Kolmogorov Complexity and Its Applications*. Springer, 1997. ISBN 9783030112974. https://doi.org/10.1007/978-1-4757-3860-5.

[35] Joseph T. Lizier, Nils Bertschinger, Juergen Jürgen Jost, and Michael Wibral. Information Decomposition of Target Effects from Multi-Source Interactions: Perspectives on Previous, Current and Future Work. *Entropy*, 20(4), 2018. ISSN 10994300. https://doi.org/10.3390/e20040307.

[36] Thomas Murray MacRobert and Thomas John I'Anson Bromwich. *An Introduction to the Theory of Infinite Series*. Macmillan, London, 1926. https://doi.org/10.2307/3604161.

[37] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, 2 edition, 2018. ISBN 978-0-262-03940-6. https://doi.org/10.1007/s00362-019-01124-9.

[38] G. L. Nemhauser, L. A. Wolsey, and M. L. Fisher. An analysis of approximations for maximizing submodular set functions–I. *Math. Program.*, 14(1):265–294, dec 1978. ISSN 0025-5610. https://doi.org/10.1007/BF01588971.

[39] Arthur J. Parzygnat. Towards a functorial description of quantum relative entropy. *arXiv e-prints*, page arXiv:2105.04059, 2021. https://doi.org/10.1007/978-3-030-80209-7_60.

[40] Rick Quax, Omri Har-Shemesh, and Peter Sloot. Quantifying Synergistic Information Using Intermediate Stochastic Variables †. *Entropy*, 19(2):85, February 2017. https://doi.org/10.3390/e19020085.

[41] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-Shot Text-to-Image Generation. *arXiv e-prints*, art. arXiv:2102.12092, February 2021. https://doi.org/10.48550/arXiv.2102.12092.

[42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-Resolution Image Synthesis with Latent Diffusion Models. 2021. https://doi.org/10.48550/arXiv.2112.10752.

[43] Fernando E. Rosas, Pedro A.M. Mediano, Michael Gastpar, and Henrik J. Jensen. Quantifying High-Order Interdependencies via Multivariate Extensions of the Mutual Information. *Physical Review E*, 100(3):1–17, 2019. ISSN 24700053. https://doi.org/10.1103/PhysRevE.100.032305.

[44] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding. *arXiv e-prints*, art. arXiv:2205.11487, May 2022. https://doi.org/10.48550/arXiv.2205.11487.

[45] René L. Schilling. *Measures, Integrals and Martingales*. Measures, Integrals and Martingales. Cambridge University Press, 2017. ISBN 9781316620243. https://doi.org/10.1017/CBO9780511810886.

[46] A Schrijver. *Combinatorial Optimization - Polyhedra and Efficiency*. Springer, 2003.

[47] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning — From Theory to Algorithms*. Cambridge University Press, 2014. ISBN 978-1-10-705713-5. https://doi.org/10.1017/CBO9781107298019.

[48] C. E. Shannon. A Mathematical Theory of Communication. *The Bell System Technical Journal*, 27(3):379–423, 1948. ISSN 15387305. https://doi.org/10.1002/j.1538-7305.1948.tb01338.x.

[49] Claude E. Shannon and Warren Weaver. *The Mathematical Theory of Communication*. The University of Illinois Press, 1964. ISBN 9780252725487. https://doi.org/10.1063/1.3067010.

[50] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep Unsupervised Learning using Nonequilibrium Thermodynamics. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of

*Proceedings of Machine Learning Research*, pages 2256–2265, Lille, France, 07–09 Jul 2015. PMLR. https://doi.org/10.48550/arXiv.1503.03585.

[51] Melvin Dale Springer. *The Algebra of Random Variables*, volume 23 of *Wiley series in probability and mathematical statistics*. John Wiley & Sons, 1979. https://doi.org/10.2307/1268039.

[52] Richard P. Stanley. *Enumerative Combinatorics*. Cambridge Studies in Advanced Mathematics. Cambridge University Press, 2 edition, 2011. https://doi.org/10.1017/CBO9781139058520.

[53] Bastian Steudel, Dominik Janzing, and Bernhard Schölkopf. Causal Markov condition for submodular information measures. *COLT 2010 - The 23rd Conference on Learning Theory*, pages 464–476, 2010. https://doi.org/10.48550/arXiv.1002.4020.

[54] Terrence Tao. *An Introduction to Measure Theory*. American Mathematical Society, 2013. ISBN 9781470409227. https://doi.org/10.1090/gsm/126.

[55] Constantino Tsallis. Possible Generalization of Boltzmann-Gibbs statistics. *Journal of Statistical Physics*, 52(1):479–487, 1988. ISSN 1572-9613. https://doi.org/10.1007/BF01016429.

[56] V. Vedral. The Role of Relative Entropy in Quantum Information Theory. *Reviews of Modern Physics*, 74(1):197–234, 2002. ISSN 00346861. https://doi.org/10.1103/RevModPhys.74.197.

[57] Juan Pablo Vigneaux. *Topology of Statistical Systems — A Cohomological Approach to Information Theory*. Phd dissertation, Université de Paris, 2019.

[58] Juan Pablo Vigneaux. Information Structures and Their Cohomology. *Theory and Applications of Categories*, 35(38):1476–1529, 2020. ISSN 1201561X. https://doi.org/10.48550/arXiv.1709.07807.

[59] Juan Pablo Vigneaux. Entropy under Disintegrations. *arXiv e-prints*, 1:arXiv:2102.09584, 2021. https://doi.org/10.48550/arXiv.2102.09584.

[60] Satosi Watanabe. Information Theoretical Analysis of Multivariate Correlation. *IBM Journal of Research and Development*, 4(1):66–82, 1960. https://doi.org/10.1147/rd.41.0066.

[61] Paul L. Williams and Randall D. Beer. Nonnegative Decomposition of Multivariate Information. *arXiv e-prints*, page arXiv:1004.2515, 2010. https://doi.org/10.48550/arXiv.1004.2515.

[62] Paul L. Williams and Randall D. Beer. Decomposing Multivariate Information. *Privately Communicated*, 2011.

[63] Raumond W. Yeung. A New Outlook on Shannon's Information Measures. *IEEE Transactions on Information Theory*, 37(3):466–474, 1991. ISSN 21915776. https://doi.org/10.1109/18.79902.

[64] Raymond W Yeung. *A First Course in Information Theory*. Springer-Verlag, Berlin, Heidelberg, 2002. ISBN 9781461346456. https://doi.org/10.1007/978-1-4419-8608-5.

[65] A. K. Zvonkin and L. A. Levin. The Complexity of Finite Objects and the Development of the Concepts of Information and Randomness By Means of the Theory of Algorithms. *Russian Mathematical Surveys*, 25(6):83–124, 1970. ISSN 0036-0279. https://doi.org/10.1070/rm1970v025n06abeh001269.