

## Shannon theory beyond quantum: Information content of a source

Paolo Perinotti,<sup>\*</sup> Alessandro Tosini,<sup>†</sup> and Leonardo Vaglini<sup>‡</sup>

*QUIT Group, Physics Department, Pavia University, and INFN Sezione di Pavia, via Bassi 6, 27100 Pavia, Italy*



(Received 31 December 2021; accepted 6 May 2022; published 25 May 2022)

The information content of a source is defined in terms of the minimum number of bits needed to store the output of the source in a perfectly recoverable way. A similar definition can be given in the case of quantum sources, with qubits replacing bits. In the mentioned cases the information content can be quantified through Shannon's and von Neumann's entropy, respectively. Here we extend the definition of information content to operational probabilistic theories, and prove relevant properties such as the subadditivity and the relation between purity and information content of a state. We prove the consistency of the present notion of information content when applied to the classical and the quantum case. Finally, the relation with one of the notions of entropy that can be introduced in general probabilistic theories, the maximum accessible information, is given in terms of a lower bound.

DOI: [10.1103/PhysRevA.105.052222](https://doi.org/10.1103/PhysRevA.105.052222)

### I. INTRODUCTION

The birth of information theory, marked by Shannon's pioneering work [1], is represented by a thorough definition of information content of an information source, along with its quantification through a suitable quantity—the celebrated Shannon entropy. Since its early times, the information content was identified with the minimum number of elementary information carriers—bits—needed to encode messages from the source in a perfectly recoverable way.

The notion of information content was then lifted to the quantum scenario, where the elementary carrier is a qubit. Schumacher [2] proved that the quantum information content of a quantum source can be quantified through its von Neumann's entropy. The typical setting for the proof of the above mentioned results entails a regularization procedure, i.e., considering an arbitrarily large number of uses of the source, and encoding schemes that are not perfect, but arbitrarily accurate. This definition makes the notion of information content depend on the choice of the figure of merit used to evaluate accuracy [3–5]. While in the classical case one uses the (complement of) error probability, Schumacher's theorem uses the entanglement fidelity, which evaluates how well the encoding scheme preserves entanglement of the source with an external system. As we show here, this figure of merit is equivalent to the average input-output fidelity for an arbitrary decomposition of the state representing the source. One big lesson of quantum information theory is thus that preserving coherence is equivalent to preserving entanglement, in a motto: *quantum information is entanglement* [6].

In a broader sense, the above argument teaches us to judge the action of a transformation looking not only at the way

it transforms input states, but also correlations with remote systems. While in quantum theory the two perspectives are somehow interchangeable, there are other information theories where the second one becomes mandatory, such as real quantum theory [7] or fermionic theory [8,9]. The latter two theories are examples of operational probabilistic theories [10–13] (OPTs) or more generally of generalized probabilistic theories [14–16]. The definition of OPTs sets a framework of alternative theories of elementary information carriers (or possibly elementary physical systems) and their processes, that can be thought of as all the conceivable information theories. In the present work we show that it is actually possible to extend the notion of information content discussed above to the framework of OPTs, and compare it to some entropic functions introduced in this context by analogy with known entropies [17–19].

The requirement that we impose on compression schemes featuring in the definition of information content is that their effect on any preparation of ensembles that average to the considered state must be indistinguishable from leaving the preparation untouched. Thus, besides considering any refinement of the state under discussion, we consider the action of the compression scheme on decompositions of its *dilations*, i.e., joint states of the system and an arbitrary external system such that the state that one obtains after averaging and discarding the external system is precisely the one of interest.

The importance of considering the effect of transformations on external systems was recently discussed in other contexts, e.g., in assessing information vs disturbance [20].

It is worth mentioning one of the main assumptions made in the present work. In classical and quantum Shannon theories the amount of information is measured in bits and qubits respectively, as we have recalled just above. For a generic theory, with no further restrictions on its structure, in principle, one may not be able to identify an elementary information carrier. For this reason, we will consider theories that we name “digitizable.” Roughly speaking, we assume the existence of

<sup>\*</sup>paolo.perinotti@unipv.it

<sup>†</sup>alessandro.tosini@unipv.it

<sup>‡</sup>leonardo.vaglini01@universitadipavia.it

at least one elementary system, which we call orbit, such that an agent can always encode an arbitrary but finite number of copies of her or his system into an array made of an integer number of orbits. This feature is evidently satisfied by classical and quantum theory, and it does not rule out scenarios that are relevant from a foundational perspective, such as nonlocal boxes [21,22]. The latter is the prototypical example of a theory where the various notions of entropy exhibit odd features, such as violation of strong subadditivity, subadditivity, and concavity, and where they are also proven to be not equivalent [17–19]. Moreover, this assumption can also be applied to theories without local tomography, such as real quantum theory and quantum theory with superselection (e.g., fermionic theory).

After introducing the notion of information content in a general OPT, we prove some of its main properties and show that the optimized accessible information [17–19] generally provides a lower bound for it. As special cases, we then analyze classical and quantum theory, where our definition boils down to Shannon’s and von Neumann’s entropy, respectively. As a consequence of the present definition, finally, fermionic information content can be proved to coincide with von Neumann’s entropy of the fermionic state [23].

The paper is organized as follows. In Sec. II we give an account of the OPT framework. We set up the basic terminology and we provide the reader with the relevant definitions and assumptions that are necessary in the present work. In Sec. III we give the formal definition of information content and we show that such a definition is well posed. We also prove properties that do not require assumptions on the structure of the theory: Subadditivity and invariance under reversible transformations acting on the state at hand. In Sec. IV we explore the consequence of the steering assumption, namely, the possibility to steer any ensemble of the state by means of one of its dilations. We show that it is possible to assess the reliability of a compression protocol by taking into account the dilations of the state only, instead of considering also all the possible decompositions of them. We generalize the entanglement fidelity and we prove that it can be equivalently used as a figure of merit for defining the information content. In Sec. V we investigate the relation between state purity and the vanishing of the information content. The main result of this section is that all the states with vanishing information content must be necessarily pure. Moreover, we show that the converse is not generally true. Indeed, we first prove that sufficient conditions are the atomicity of parallel composition of states and essential uniqueness of purification. Then, we show that atomicity of parallel composition of states turns out to be a necessary condition for having a null information content on pure states. Thus, any theory violating this property must have pure states with a strictly positive information content. Finally, we show that the accessible information is a lower bound for the information content here introduced. In Sec. VI we show that the information content is simply given by Shannon and von Neumann entropies in classical and quantum information theories, respectively. The point here is to check the collapse of our figure of merit and the one usually used in those theories. We conclude summarizing the results of this work and discussing some open questions in Sec. VII.

## II. OPT FRAMEWORK

In this section we briefly review the framework of operational probabilistic theories.

### A. General description

The primitive notions of an operational theory are those of test, event and system. A test  $\{\mathcal{A}_i\}_{i \in X}$  is given by a collection of events, where  $i$  labels the elements of the outcome space  $X$ . The systems allow for the connection between different tests, and are denoted by capital Roman letters  $A, B, \dots$ . Therefore, a test is completely determined by its input and output systems, and the events associated with the outcome space  $X$ . In order to represent a test and its events  $\{\mathcal{A}_i\}_{i \in X}$  we use the usual diagrammatic notation

$$\begin{array}{c} A \\ \hline \boxed{\{\mathcal{A}_i\}_{i \in X}} \\ \hline B \end{array}, \quad \begin{array}{c} A \\ \hline \boxed{\mathcal{A}_i} \\ \hline B \end{array},$$

and we will call  $A, B$  the input and the output system of the test, respectively. If  $\{\mathcal{A}_i\}_{i \in X}$  and  $\{\mathcal{B}_j\}_{j \in Y}$  are two tests, one can define their sequential composition as the test  $\{\mathcal{C}_{i,j}\}_{(i,j) \in X \times Y}$ , with events  $\mathcal{C}_{i,j}$  that are diagrammatically represented by

$$\begin{array}{c} A \\ \hline \boxed{\mathcal{C}_{i,j}} \\ \hline C \end{array} := \begin{array}{c} A \\ \hline \boxed{\mathcal{A}_i} \\ \hline B \\ \hline \boxed{\mathcal{B}_j} \\ \hline C \end{array}.$$

Notice that this definition requires the output system of the events on the left to be necessarily the input system of the events on the right. A singleton test is a test whose outcome space set  $X$  is a singleton, and the unique event contained in it is called deterministic. For any system  $A$  there exists a unique identity test  $\{\mathcal{I}_A\}$  such that  $\mathcal{A} \mathcal{I}_A = \mathcal{A}$  ( $\mathcal{I}_A \mathcal{A} = \mathcal{A}$ ) for any event  $\mathcal{A}$ . Another operation that can be performed on tests for defining a new test is parallel composition. Given two systems  $A$  and  $B$  we call  $AB$  the composite system of  $A$  and  $B$ . Then, if  $\{\mathcal{A}_i\}_{i \in X}$  and  $\{\mathcal{B}_j\}_{j \in Y}$  are two tests, their parallel composition is the test  $\{\mathcal{A}_i \boxtimes \mathcal{B}_j\}_{(i,j) \in X \times Y}$ . Diagrammatically

$$\begin{array}{c} A \quad B \\ \hline \boxed{\mathcal{A}_i \boxtimes \mathcal{B}_j} \\ \hline C \quad D \end{array} := \begin{array}{c} A \quad B \\ \hline \boxed{\mathcal{A}_i} \\ \hline C \quad D \\ \hline \boxed{\mathcal{B}_j} \\ \hline \end{array}.$$

The parallel composition operation commutes with the sequential one, namely,  $(\mathcal{E}_h \boxtimes \mathcal{F}_k)(\mathcal{C}_i \boxtimes \mathcal{D}_j) = (\mathcal{E}_h \mathcal{C}_i) \boxtimes (\mathcal{F}_k \mathcal{D}_j)$ .

There is a special kind of system, the trivial system  $I$ , satisfying  $AI = IA = A$  for every system  $A$ . Tests with  $I$  as input system and  $A$  as the output one are called preparation tests of  $A$ , while tests with input system  $A$  and  $I$  as output are named observation tests of  $A$ . The events of a preparation test  $\{\rho_i\}_{i \in X}$  and of an observation test  $\{a_j\}_{j \in Y}$  are represented through the following diagrams:

$$\begin{array}{c} A \\ \hline \boxed{\rho_i} \end{array} := \begin{array}{c} I \\ \hline \boxed{\rho_i} \\ \hline A \end{array},$$

$$\begin{array}{c} A \\ \hline \boxed{a_j} \end{array} := \begin{array}{c} A \\ \hline \boxed{a_j} \\ \hline I \end{array}.$$

In the following we will always use Greek letters to denote preparation tests and Latin letters for the observation test.

Preparation and observation events will also be denoted by using round brackets, respectively,  $|\rho\rangle_A$  and  $\langle a|_A$ , and we will not make explicit the system whenever it is clear from the context.

A circuit is a diagram representing an arbitrary test that is obtained by sequential and parallel composition of other tests. We say that a circuit is *closed* when the input and output systems are both the trivial one, namely, when it starts with a preparation test and it ends with an observation test. An *operational probabilistic theory* is an operational theory where any closed circuit (equivalently, any test from the trivial system to itself) is given by a joint probability distribution conditioned by the tests building the circuit. Moreover, compound tests from the trivial system to itself are independent, namely, the joint probability distribution is simply given by the product of the probability distributions of the composing tests. A simple example is given by a preparation test  $\{\rho_i\}_{i \in X}$  sequentially followed by an observation test  $\{a_j\}_{j \in Y}$ :

$$\{p(i, j | \{\rho_i\}, \{a_j\})\} := \left\{ \left( \rho_i \right) \xrightarrow{A} \left( a_j \right) \right\},$$

with  $\sum_{i,j} p(i, j) = 1$ . Thus, one has a joint probability distribution, which is conditioned by the chosen tests  $\{\rho_i\}_{i \in X}$  and  $\{a_j\}_{j \in Y}$ . From now on we will simply omit this dependence. The probability associated with the closed circuit where a preparation  $\rho_i$  is followed by an observation  $a_j$  will also be denoted by a pairing,  $p(i, j) = \langle a_j | \rho_i \rangle$ .

Given any system A of an OPT, One can define an equivalence relation on the set of preparation events by declaring that  $\rho \sim \sigma$  iff  $\langle a | \rho \rangle \neq \langle a | \sigma \rangle$  for any observation event  $a$ . The set of equivalence classes with respect to this relation is called the set of *states* of system A, and it is denoted by  $\text{St}(A)$ . Similarly, one can define the set of *effects* as the set of equivalence classes of the observation events such that  $\langle a | \rho \rangle \neq \langle b | \rho \rangle$  for any preparation event  $\rho$ , and this is denoted by  $\text{Eff}(A)$ . The sets of deterministic states and effects will be denoted by  $\text{St}_1(A)$  and  $\text{Eff}_1(A)$ , respectively.

Given the probabilistic structure, states can be seen as functionals on the set of effects and vice versa, and then one can consider linear combinations of them, thus defining two linear spaces,  $\text{St}(A)_{\mathbb{R}}$  and  $\text{Eff}(A)_{\mathbb{R}}$ , which are dual to each other assuming that they are finite-dimensional. The *size*  $D_A$  of a given system A is simply the dimension of the linear space  $\text{St}_{\mathbb{R}}(A)$ . A transformation event from system A to system B induces a linear map from  $\text{St}_{\mathbb{R}}(AC)$  to  $\text{St}_{\mathbb{R}}(BC)$  for any ancillary system C. Also the set of transformation events can be endowed with an equivalence relation. Indeed, given  $\mathcal{A}$  and  $\mathcal{B}$ , we say that they are *operationally equivalent* ( $\mathcal{A} \sim \mathcal{B}$ ) if the following identity holds:

$$\left( \Psi \right) \xrightarrow{A} \left( \begin{array}{c} A \\ \left( \begin{array}{c} C \\ A \end{array} \right) \\ B \end{array} \right) = \left( \Psi \right) \xrightarrow{A} \left( \begin{array}{c} A \\ \left( \begin{array}{c} C \\ B \end{array} \right) \\ B \end{array} \right),$$

for any  $\Psi \in \text{St}(AC)$ ,  $A \in \text{Eff}(BC)$  and any ancillary system C. In other words, two transformation events are operationally equivalent if they induce the same linear map for any ancillary system C. We then denote the set of all the equivalence classes with  $\text{Tr}(A \rightarrow B)$ , whose elements are simply called *transformations*. As for states and effects, the set of

deterministic transformations will be denoted by  $\text{Tr}_1(A \rightarrow B)$ . If  $\mathcal{U} \in \text{Tr}(A \rightarrow B)$  and there exists  $\mathcal{V} \in \text{Tr}(B \rightarrow A)$  such that  $\mathcal{V}\mathcal{U} = \mathcal{I}_A$  and  $\mathcal{U}\mathcal{V} = \mathcal{I}_B$ , we say that  $\mathcal{U}$  is *reversible*. Accordingly, two systems are called *operationally equivalent* if there exists a reversible transformation  $\mathcal{U} \in \text{Tr}_1(A \rightarrow B)$ . A notion that will be useful in the following is that of *asymptotically equivalent systems*.

*Definition II.1 (Asymptotical equivalence).* Two systems  $A_1$  and  $A_2$  are asymptotically equivalent if

- (1) There exists a pair of integers  $k_1, k_2 < \infty$ ,  $\mathcal{E} \in \text{Tr}_1(A_1^{\boxtimes k_1} \rightarrow A_2^{\boxtimes k_2})$  and  $\mathcal{D} \in \text{Tr}_1(A_2^{\boxtimes k_2} \rightarrow A_1^{\boxtimes k_1})$  such that  $\mathcal{D}\mathcal{E} = \mathcal{I}_{A_1^{\boxtimes k_1}}$ ;
- (2) There exists a pair of integers  $h_1, h_2 < \infty$ ,  $\mathcal{G} \in \text{Tr}_1(A_2^{\boxtimes h_2} \rightarrow A_1^{\boxtimes h_1})$  and  $\mathcal{F} \in \text{Tr}_1(A_1^{\boxtimes h_1} \rightarrow A_2^{\boxtimes h_2})$  such that  $\mathcal{F}\mathcal{G} = \mathcal{I}_{A_2^{\boxtimes h_2}}$ ;
- (3) Let  $M_2^{\min}(k_1)$  be the smallest  $k_2$  such that item 1 is satisfied for a given  $k_1$ , and similarly for  $M_1^{\min}(h_2)$  with reference to item 2. The following assumption is made:

$$\lim_{k_1 \rightarrow \infty} \frac{M_2^{\min}(k_1)}{k_1} = k, \quad \lim_{h_2 \rightarrow \infty} \frac{M_1^{\min}(h_2)}{h_2} = k^{-1}. \quad (1)$$

Now we set up some terminology and we introduce pure and mixed states, as well as the definition of state dilation.

*Definition II.2 (Refinement of an event).* Let  $\mathcal{C} \in \text{Tr}(A \rightarrow B)$ . A refinement of  $\mathcal{C}$  is given by a collection of events  $\{\mathcal{D}_i\}_{i \in Y} \subseteq \text{Tr}(A \rightarrow B)$  such that there exists a test  $\{\mathcal{D}_i\}_{i \in X}$  with  $Y \subseteq X$  and  $\mathcal{C} = \sum_{j \in Y} \mathcal{D}_j$ . We say that a refinement  $\{\mathcal{D}_i\}_{i \in Y}$  is trivial if  $\mathcal{D}_i = \lambda_i \mathcal{C}$ ,  $\lambda_i \in [0, 1]$  for every  $i \in Y$ . Conversely,  $\mathcal{C}$  is called the coarse graining of the events  $\{\mathcal{D}_i\}_{i \in Y}$ .

*Definition II.3.* Given two events  $\mathcal{C}, \mathcal{D} \in \text{Tr}(A \rightarrow B)$  we say that  $\mathcal{D}$  refines  $\mathcal{C}$ , and write  $\mathcal{D} < \mathcal{C}$ , if there exist a refinement  $\{\mathcal{D}_i\}_{i \in X}$  of  $\mathcal{C}$  such that  $\mathcal{D} \in \{\mathcal{D}_i\}_{i \in X}$ .

*Definition II.4 (Atomic and refinable events).* An event  $\mathcal{C}$  is called atomic if it admits only trivial refinements. An event is refinable if it is not atomic.

The notion of refinement and refinable events give rise to the definitions of *pure* and *mixed* states.

*Definition II.5 (Pure and mixed states).*  $\rho \in \text{St}(A)$  is called pure if it is atomic, or mixed otherwise. We will denote by  $\text{PurSt}(A)$  the set of all the pure states of system A.

*Definition II.6.* Let  $\rho \in \text{St}(A)$  and  $\Psi \in \text{St}(AB)$ . We say that  $\Psi$  is a dilation of  $\rho$  if there exists a deterministic effect  $e \in \text{Eff}(B)$  such that

$$\left( \rho \right) \xrightarrow{A} = \left( \Psi \right) \xrightarrow{A} \left( B \right) \xrightarrow{e}$$

We denote by  $D_\rho$  the set of all dilations of the state  $\rho$ . If  $\Psi$  is also pure, then we say that it is a purification of  $\rho$  and B is called the purifying system. Finally, we denote by  $P_\rho$  the set of all the purifications of  $\rho$ .

Trivially one has that  $P_\rho \subseteq D_\rho$ . Moreover, if  $\Omega \in D_\rho$ , then one has  $D_\Omega \subseteq D_\rho$ .

Two special instances of this framework are classical and quantum theory. In classical theory, the systems are associated with real vector spaces  $\mathbb{R}^{d_A}$ , and different systems are associated with different values of  $d_A$ . The set of states is made of substochastic vectors in these spaces, namely, by vectors

$\mathbf{x}$  satisfying  $\|\mathbf{x}\|_1 := \sum_{i=1}^{d_A} |x_i| \leq 1$ , and therefore it is a simplex. The pure states are then represented by the canonical basis vectors  $\mathbf{e}_i$  of  $\mathbb{R}^{d_A}$ . The convex set of effects is given by unit-dominated positive vectors, i.e., those  $\mathbf{x} \in \mathbb{R}^{d_A}$  such that  $0 \leq x_i \leq 1$  for any  $i = 0, \dots, d_A$ . Transformations from system A to system B are represented by  $d^B \times d^A$  substochastic matrices  $\mathbf{M}$  acting on the probability vectors by multiplication  $\mathbf{x} \rightarrow \mathbf{M}\mathbf{x}$ . Recall that a matrix is substochastic when each column is a substochastic vector. Sequential and parallel composition are trivially given by matrix multiplication and tensor product, respectively.

In quantum theory systems are associated with Hilbert spaces, and different systems are represented by spaces of different dimension  $d$  (we will assume  $d < \infty$  whenever we will refer to results relative to quantum theory). The convex set of states of a system A is given by subnormalized density matrices  $\rho$  on the associated space  $\mathcal{H}_A$ , i.e., matrices such that  $\rho > 0$  and  $\text{Tr}(\rho) \leq 1$ . The convex set of effects is given by functionals (acting on the set of states) of the form  $\text{Tr}(-E)$ , where  $E$  is a positive matrix dominated by the identity,  $0 \leq E \leq I$ . Transformations in  $\text{Tr}(A \rightarrow B)$  are mathematically represented by completely positive trace nonincreasing maps. These have a Kraus decomposition, namely, if  $\mathcal{E} \in \text{Tr}(A \rightarrow B)$  there exists a set of operators  $\{E_i\} \subseteq \mathcal{L}(\mathcal{H}_B, \mathcal{H}_A)$  such that  $\mathcal{E}(\cdot) = \sum_i E_i^\dagger \cdot E_i$ . Sequential composition is simply given by composition of maps, and the parallel one is represented by the tensor product, as in the classical case.

The linear space  $\text{St}_{\mathbb{R}}(A)$  can be endowed with a metric structure by means of the following norm, which has an operational meaning related to optimal discrimination schemes [11].

*Definition II.7 (Operational norm).* The norm of an element  $\rho \in \text{St}(A)_{\mathbb{R}}$  is defined as

$$\|\rho\|_{\text{op}} := \sup_{a \in \text{Eff}(A)} (2a - e_A|\rho),$$

where  $e_A$  is the deterministic effect obtained by the coarse graining of the observation test containing  $a$ .

This norm satisfies a monotonicity property, as stated in the following lemma.

*Lemma II.1 (Monotonicity of the operational norm).* For any  $\delta \in \text{St}_{\mathbb{R}}(A)$  and  $\mathcal{C} \in \text{Tr}(A, B)$  the following inequality holds:

$$\|\mathcal{C}\delta\|_{\text{op}} \leq \|\delta\|_{\text{op}}, \tag{2}$$

with the equality holding iff  $\mathcal{C}$  is reversible.

This notion of norm, which is valid for any OPT, reduces to the trace norm in quantum theory.

### B. Restricting the class of theories

Upon marginalization over the observation test, one can define the preparation probability conditioned by the test  $\{a_j\}_{j \in Y}$  as  $p(\rho_i | \{a_j\}) := \sum_j p(i, j)$ . Generally, the preparation probability is not one, unless the preparation test  $\{\rho_i\}_{i \in X}$  is the singleton, i.e., the state is deterministic. Moreover, as it is clear by its definition, it can also depend on the observation test we are marginalizing over. Usually, the causality condition is expressed as a *no signaling from the future* principle, namely, by saying that preparation probabilities are actually

independent of the chosen observation test, which is equivalent to state the uniqueness of the deterministic effect. In this paper we will adopt a stronger form of causality, that is the following one.

*Assumption 1 (Causal theories).* An OPT satisfies strong causality if for every test  $\{\mathcal{A}_i\}_{i \in X}$  and every collection of tests  $\{\mathcal{B}_j^i\}_{j \in Y}$  labeled by  $j \in Y$ , the collection of events  $\{\mathcal{C}_{i,j}\}_{(i,j) \in X \times Y}$  with

$$\text{---} \boxed{\mathcal{C}_{i,j}} \text{---} := \text{---} \boxed{\mathcal{A}_i} \text{---} \text{---} \boxed{\mathcal{B}_j^i} \text{---} \text{---},$$

is a test of the theory.

One can show that the above statement implies uniqueness of the deterministic effect [10,11].

Another assumption that we will use in some sections of this work stems from the steering property of quantum theory. This asserts that, given a state  $\rho \in \text{St}(A)$  and a purification  $\Phi \in \text{PurSt}(AB)$  of  $\rho$ , for any decomposition  $\sum_{i \in X} p_i \sigma_i$  of  $\rho$  there exists an observation test  $\{b_i\}_{i \in X} \subseteq \text{Eff}(B)$  such that

$$\boxed{\sigma_i} \text{---} \text{---} \text{---} \boxed{A} = \left( \begin{array}{c} \boxed{A} \\ \Phi \\ \boxed{B} \end{array} \right) \text{---} \boxed{b_i}. \quad \forall i \in X.$$

This feature cannot be assumed as it stands: The first reason is that a generic OPT may not encompass the existence of a purification for any state of the theory (and classical theory is a trivial example), therefore we are led to consider dilations instead of purifications. Second, there is no reason why one should be able to steer any decomposition by means of the same dilation. Thus, for a generic theory, one can state the steering as follows.

*Assumption 2 (Steering).* Let  $\rho \in \text{St}(A)$  and  $\{\sigma_i\}_{i \in X} \subseteq \text{St}(A)$  be a refinement of  $\rho$ . Then there exist a system B, a state  $\Psi \in \text{St}(AB)$ , and an observation test  $\{b_i\}_{i \in X}$  such that

$$\boxed{\sigma_i} \text{---} \text{---} \text{---} \boxed{A} = \left( \begin{array}{c} \boxed{A} \\ \Psi \\ \boxed{B} \end{array} \right) \text{---} \boxed{b_i}, \quad \forall i \in X.$$

Notice that the state  $\Psi$  in the steering assumption must be a dilation of  $\rho$ , as one can easily verify upon summing over  $i \in X$ . The stronger steering feature satisfied by quantum theory can actually be proven to hold in any OPT satisfying atomicity of parallel composition of states, existence and uniqueness (up to reversible channels) of purification and perfect discriminability as axioms. For the present purposes we choose to state steering as a property that an OPT may or may not satisfy rather than discussing the conditions under which it holds.

Another property that we will assume throughout the paper is digitizability.

*Assumption 3 (Digitizability).* We say that an OPT is digitizable if there exists a system B (called *obit*) such that for any system X there exists  $k < \infty$  and a pair of maps  $\mathcal{C} \in \text{Tr}_1(X \rightarrow B^{\boxtimes k})$  and  $\mathcal{F} \in \text{Tr}_1(B^{\boxtimes k} \rightarrow X)$  such that  $\mathcal{F} \circ \mathcal{C} = \mathcal{I}_X$ . Moreover, if  $B_1$  and  $B_2$  are two such systems, then they are asymptotically equivalent.

The above assumption holds in quantum theory, since any qudit system can be taken as elementary. Moreover, let us consider two different qudits with dimension  $d_1$  and  $d_2$ , respectively. Generally, the equation  $d_1^N = d_2^M$  may have no integer solutions (for instance, when both  $d_1$  and  $d_2$  are prime). However, the smallest integer  $M$  such that we can isometrically embed  $\mathcal{H}_1^{\otimes N}$  into  $\mathcal{H}_2^{\otimes M}$  is given by  $\lceil N \log_{d_2} d_1 \rceil$  which is such that  $\lim_{N \rightarrow \infty} \frac{\lceil N \log_{d_2} d_1 \rceil}{N} = \log_{d_2} d_1$ . Similarly,  $\lceil M \log_{d_1} d_2 \rceil$  is needed for an isometric embedding of  $\mathcal{H}_2^{\otimes M}$  into  $\mathcal{H}_1^{\otimes N}$ , and  $\lim_{M \rightarrow \infty} \frac{\lceil M \log_{d_1} d_2 \rceil}{M} = \log_{d_1} d_2 (= 1/\log_{d_2} d_1)$ . In any theory satisfying the assumption of digitizability we can always encode the state of our system on the parallel composition of a sufficiently large number of elementary systems, which we can think of as a generalization of the qubit system for quantum theory. The request of digitizability comes from the need of a unit for the amount of information required for storing a given source. In classical information theory we use bits, in the quantum counterpart the qubits, and for a generic OPT satisfying digitizability we use obits, whose existence must then be postulated.

We want to stress the fact that the assumption of digitizability is extremely weak, to the extent that every theory in the literature abides by it, and it is very hard to imagine a theory that violates it. Indeed, a nondigitizable theory should contain infinitely many inequivalent system types, even asymptotically (see Definition II.1), and this immediately brings us into an unexplored territory of wild theories.

The nonlocal boxes [21] provide us with a non trivial example of an OPT satisfying the digitizability assumption, and with a strong departure from the quantum one. There exists a unique single system, whose state space is described by a square, and multipartite systems are obtained by using only this one, therefore it is trivially digitizable. In the literature—with a few remarkable exceptions [16,24]—nonlocal boxes are presented in terms of the geometry of their state space, and focusing on the correlations that measurements can produce, disregarding the behavior of transformations. However one can straightforwardly make them into an OPT by assuming that every collection of linear maps on the state space that map preparations to preparations is allowed. A similar construction was carried out, e.g., in Ref. [16]. Nonlocal boxes provide a scenario where the wealth of the known entropy notions is manifest. Moreover, the fact that we are only referring to the conversion of finitely many copies of the system at hand is not constraining from a conceptual point of view. For the present purposes, namely, taking a first step towards a Shannon theory for generic physical systems, this level of analysis is sufficient. However, the composition of a countable number of systems can be suitably defined (see [25]) opening the route to a generalization of this property in the infinite case.

### III. INFORMATION CONTENT IN OPT

Let  $\rho \in \text{St}_1(A)$  and consider  $N$  copies of the system on which we have prepared the same state  $\rho$  and let  $M$  be a positive integer. A compression scheme is then a pair of maps  $\mathcal{E} \in \text{Tr}_1(A^{\boxtimes N} \rightarrow B^{\boxtimes M})$ ,  $\mathcal{D} \in \text{Tr}_1(B^{\boxtimes M} \rightarrow A^{\boxtimes N})$ . Sometimes we will denote by  $\mathcal{C}$  the composition  $\mathcal{D}\mathcal{E}$ .

*Definition III.1.* An  $(\varepsilon, N)$ -reliable compression scheme  $(\mathcal{E}, \mathcal{D})$  is such that

$$\sup_{C, \{\Psi_i\}} \sum_{i \in X} \|[(\mathcal{D}\mathcal{E}) \boxtimes \mathcal{I}_C] \Psi_i - \Psi_i\| < \varepsilon,$$

where  $\{\Psi_i\}_{i \in X} \subseteq \text{St}_1(A^{\boxtimes N}C)$  denotes a refinement of any dilation of  $\rho^{\boxtimes N}$ . For fixed  $N, M$  we denote with  $E_{N,M,\varepsilon}(\rho)$  the set of  $\varepsilon$ -reliable compression schemes.

*Definition III.2 (Information Content).* Let  $\rho \in \text{St}_1(A)$ . We define the *smallest achievable compression ratio* for length  $N$  to tolerance  $\varepsilon$  as follows:

$$R_{N,\varepsilon}(\rho) := \frac{\min\{M : E_{N,M,\varepsilon}(\rho) \neq \emptyset\}}{N}. \tag{3}$$

The *information content* of the state  $\rho$  is defined as

$$I(\rho) := \lim_{\varepsilon \rightarrow 0} \limsup_{N \rightarrow \infty} R_{N,\varepsilon}(\rho). \tag{4}$$

*Proposition III.1.*  $I(\rho)$  is well defined for every  $\rho \in \text{St}(A)$  and every system  $A$ .

*Proof.* First, we show that for any choice of the elementary system,  $I(\rho)$  is a finite number for any state  $\rho$ . By the digitizability assumption we know that for any  $N$  there exists a positive integer  $K < \infty$  and a pair of maps  $\mathcal{E} \in \text{Tr}_1(A^{\boxtimes N} \rightarrow B^{\boxtimes K})$ ,  $\mathcal{D} \in \text{Tr}_1(B^{\boxtimes K} \rightarrow A^{\boxtimes N})$  such that  $\mathcal{D}\mathcal{E} = \mathcal{I}_{A^{\boxtimes N}}$ . Therefore, for any  $N, \varepsilon$  the set  $E_{N,K,\varepsilon}(\rho)$  is not empty, and the minimum in Eq. (3) is always finite. Moreover, it is immediate to realize that  $K$  does not need to grow more than linearly versus  $N$ , just considering  $N$  repetitions of the encoding for one copy  $A$ . Thus, the ratio in Eq. (3) is bounded, and one can take the  $\limsup_{N \rightarrow \infty}$  safely. The existence of  $\lim_{\varepsilon \rightarrow 0}$  follows by the fact that  $E_{N,M,\varepsilon}(\rho) \subseteq E_{N,M,\varepsilon'}(\rho) \neq \emptyset$  whenever  $\varepsilon \leq \varepsilon'$ , which, in turn, implies monotonicity of the function  $\limsup_{N \rightarrow \infty} R_{N,\varepsilon}(\rho)$  versus  $\varepsilon$ .

What is left to prove is that using two different obits we are not led to two incomparable notions of information content. First, fix  $N, \varepsilon$ , let  $\rho \in \text{St}_1(A)$ , and let

$$M_{1,N} := \min\{M : E_{N,M,\varepsilon}^1(\rho) \neq \emptyset\},$$

$$M_{2,N} := \min\{M : E_{N,M,\varepsilon}^2(\rho) \neq \emptyset\}$$

be the minimum number of obits  $B_1$  and  $B_2$  needed for an  $\varepsilon$ -optimal encoding, respectively. Rephrasing the first equation in (1), there exists a sequence  $\delta_1(M_1)$  such that  $M_2^{\min}(M_1) = kM_1 + \delta_1(M_1)$  with  $\lim_{M_1 \rightarrow \infty} \frac{\delta_1(M_1)}{M_1} = 0$ . Given the encoding  $\mathcal{E}$  of item 1 in Definition II.1 from  $M_{1,N}$  to  $M_2^{\min}(M_{1,N})$  (see item 3 in Definition II.1), we have an  $\varepsilon$ -optimal encoding of  $\rho^{\boxtimes N}$  onto  $M_2^{\min}(M_{1,N})$  obits  $B_2$ , therefore  $M_{2,N} \leq M_2^{\min}(M_{1,N})$  and this implies

$$\begin{aligned} \limsup_{N \rightarrow \infty} \frac{M_{2,N}}{N} &\leq \limsup_{N \rightarrow \infty} \frac{M_2^{\min}(M_{1,N})}{N} \\ &= \limsup_{N \rightarrow \infty} \left[ \frac{kM_{1,N}}{N} + \frac{\delta_1(M_{1,N})}{N} \right] \\ &\leq \limsup_{N \rightarrow \infty} \frac{kM_{1,N}}{N} + \limsup_{N \rightarrow \infty} \left| \frac{\delta_1(M_{1,N})}{M_{1,N}} \frac{M_{1,N}}{N} \right| \\ &= k \limsup_{N \rightarrow \infty} \frac{M_{1,N}}{N}. \end{aligned}$$

The last line follows by  $\lim_{M_1 \rightarrow \infty} \frac{\delta_1(M_1)}{M_1} = 0$  along with the fact that  $M_{1,N+1} \geq M_{1,N}$ . Taking  $\lim_{\varepsilon \rightarrow 0}$  we end up with  $I_2(\rho) \leq kI_1(\rho)$ . A similar argument can be used to show the reverse inequality, and we have that, for any  $\rho \in \text{St}(A)$ ,  $I_2(\rho) = kI_1(\rho)$ . ■

*Proposition III.2 (Subadditivity).* Let  $\Psi \in \text{St}_1(AB)$  and let  $\rho \in \text{St}_1(A)$ ,  $\sigma \in \text{St}_1(B)$  be its marginals. Then the following property holds:

$$I(\Psi) \leq I(\rho) + I(\sigma).$$

*Proof.* Let  $(\mathcal{E}^\rho, \mathcal{D}^\rho) \in E_{N, \bar{K}, \varepsilon}(\rho)$ ,  $(\mathcal{E}^\sigma, \mathcal{D}^\sigma) \in E_{N, \bar{L}, \varepsilon}(\sigma)$  with  $\bar{K} := \min\{M : E_{N, M, \varepsilon}(\rho) \neq \emptyset\}$  and similarly for  $\bar{L}$ . Now let  $\{\Gamma_i\}_{i \in X}$  be such that  $\sum_{i \in X} \Gamma_i \in D_{\Psi^{\boxtimes N}}$  and consider  $\{(\mathcal{C}^\rho \boxtimes \mathcal{J})(\Gamma_i)\}_{i \in X}$ , where  $\mathcal{C}^\rho := \mathcal{D}^\rho \mathcal{E}^\rho$ . Since  $\mathcal{C}^\rho$  is a channel and  $D_{\Psi^{\boxtimes N}} \subseteq D_{\rho^{\boxtimes N}}$ ,  $D_{\Psi^{\boxtimes N}} \subseteq D_{\sigma^{\boxtimes N}}$  we have that  $\sum_{i \in X} (\mathcal{C}^\rho \boxtimes \mathcal{J})(\Gamma_i)$  is a dilation of both  $\rho^{\boxtimes N}$  and  $\sigma^{\boxtimes N}$ . Similarly for  $\mathcal{C}^\sigma := \mathcal{D}^\sigma \mathcal{E}^\sigma$ . This implies the following bound:

$$\begin{aligned} & \sum_{i \in X} \|(\mathcal{C}^\rho \boxtimes \mathcal{C}^\sigma \boxtimes \mathcal{J})(\Gamma_i) - \Gamma_i\|_{\text{op}} \\ &= \sum_{i \in X} \|(\mathcal{C}^\rho \boxtimes \mathcal{C}^\sigma \boxtimes \mathcal{J})(\Gamma_i) - (\mathcal{C}^\rho \boxtimes \mathcal{J})(\Gamma_i) \\ & \quad + (\mathcal{C}^\rho \boxtimes \mathcal{J})(\Gamma_i) - \Gamma_i\|_{\text{op}} \\ &< 2\varepsilon, \end{aligned}$$

where we used the triangle inequality for the operational norm. Thus  $E_{N, \bar{K} + \bar{L}, 2\varepsilon}(\Psi) \neq \emptyset$  and this implies that

$$\frac{\min\{M : E_{N, M, 2\varepsilon}(\Psi) \neq \emptyset\}}{N} \leq \frac{\bar{K}}{N} + \frac{\bar{L}}{N}.$$

Finally, by taking the  $\limsup_{N \rightarrow \infty}$  and then  $\lim_{\varepsilon \rightarrow 0}$  on both sides we get the thesis. ■

We notice that, in order to compute the information content, one can test the compression schemes on pure decompositions  $\{p_i, \Phi_i\}$  only. More precisely, let  $E_{N, M, \varepsilon}^{\text{pur}}(\rho)$  be the set of schemes which are  $(\varepsilon, N)$ -reliable according to the following criterion:

$$\sup \left( \sum_i p_i \|(\mathcal{D} \mathcal{E} \boxtimes \mathcal{J})(\Phi_i) - \Phi_i\|_{\text{op}} \right) < \varepsilon, \quad (5)$$

where the supremum is taken on all the pure decompositions  $\{p_i, \Phi_i\}$  of any  $\Omega \in D_{\rho^{\boxtimes N}}$ . Let  $I^{\text{pur}}(\rho)$  be the information content computed restricting to such maps:

$$\begin{aligned} R_{N, \varepsilon}^{\text{pur}}(\rho) &:= \frac{\min\{M : E_{N, M, \varepsilon}^{\text{pur}}(\rho) \neq \emptyset\}}{N} \\ I^{\text{pur}}(\rho) &:= \lim_{\varepsilon \rightarrow 0} \limsup_{N \rightarrow \infty} R_{N, \varepsilon}^{\text{pur}}(\rho). \end{aligned}$$

Then one has  $I(\rho) = I^{\text{pur}}(\rho)$ .

*Lemma III.1.* Let  $\rho \in \text{St}_1(A)$ , then  $I(\rho) = I^{\text{pur}}(\rho)$

*Proof.* On the one hand, we trivially have  $I^{\text{pur}}(\rho) \leq I(\rho)$ . On the other hand, let  $\{\Psi_i\}_{i \in X}$  be a refinement of  $\Omega \in D_{\rho^{\boxtimes N}}$ . For any  $i$  we can further decompose  $\Psi_i$  in terms of pure states  $\{q_j^i, \Phi_{i,j}\}_{j \in Y}$ , with  $\sum_{i,j} q_j^i = 1$ . Therefore  $\{q_j^i, \Phi_{i,j}\}_{(i,j) \in X \times Y}$  is a pure state decomposition of  $\Omega$ , and by the triangle

inequality one has

$$\begin{aligned} & \sum_{i \in X} \|[(\mathcal{D} \circ \mathcal{C}) \boxtimes \mathcal{J}]\Psi_i - \Psi_i\|_{\text{op}} \\ & \leq \sum_{(i,j) \in X \times Y} q_j^i \|[(\mathcal{D} \circ \mathcal{C}) \boxtimes \mathcal{J}]\Phi_{i,j} - \Phi_{i,j}\|_{\text{op}}. \end{aligned}$$

This implies that  $E_{N, M, \varepsilon}^{\text{pur}}(\rho) \subseteq E_{N, M, \varepsilon}(\rho)$ , and in turns that  $I(\rho) \leq I^{\text{pur}}(\rho)$ . Therefore  $I(\rho) = I^{\text{pur}}(\rho)$ . ■

*Proposition III.3.* Let  $\rho \in \text{St}_1(A)$  and  $\mathcal{U} \in \text{Tr}_1(A)$  be a reversible channel, then  $I(\rho) = I(\mathcal{U}(\rho))$ .

*Proof.* We show that  $E_{N, M, \varepsilon}(\mathcal{U}(\rho)) \neq \emptyset \Rightarrow E_{N, M, \varepsilon}(\rho) \neq \emptyset$ . Let  $(\mathcal{E}, \mathcal{D}) \in E_{N, M, \varepsilon}(\mathcal{U}(\rho))$  and let  $\{\Psi_i\}_{i \in X}$  be such that  $\sum_{i \in X} \Psi_i \in D_{\rho^{\boxtimes N}}$ . It is clear that  $\sum_{i \in X} (\mathcal{U}^{\boxtimes N} \boxtimes \mathcal{J})(\Psi_i) \in D_{\mathcal{U}(\rho)^{\boxtimes N}}$  and therefore

$$\sum_{i \in X} \|[(\mathcal{D} \mathcal{E} - \mathcal{J}) \boxtimes \mathcal{J}](\mathcal{U}^{\boxtimes N} \boxtimes \mathcal{J})(\Psi_i)\|_{\text{op}} < \varepsilon.$$

Upon defining  $\tilde{\mathcal{E}} := \mathcal{E} \mathcal{U}^{\boxtimes N}$  and  $\tilde{\mathcal{D}} := (\mathcal{U}^{-1})^{\boxtimes N} \mathcal{D}$ , recalling that  $\mathcal{U}$  is reversible and that the operational norm is invariant under reversible transformations, the above inequality can be rewritten as follows:

$$\begin{aligned} \varepsilon &> \sum_{i \in X} \|(\mathcal{U}^{\boxtimes N} \boxtimes \mathcal{J})[(\tilde{\mathcal{D}} \tilde{\mathcal{E}} \boxtimes \mathcal{J})(\Psi_i) - (\Psi_i)]\|_{\text{op}} \\ &= \sum_{i \in X} \|[(\tilde{\mathcal{D}} \tilde{\mathcal{E}} \boxtimes \mathcal{J})(\Psi_i) - (\Psi_i)]\|_{\text{op}}, \end{aligned}$$

namely, since  $\{\Psi_i\}_{i \in X}$  is arbitrary,  $(\tilde{\mathcal{E}}, \tilde{\mathcal{D}}) \in E_{N, M, \varepsilon}(\rho) \neq \emptyset$ . This implies that  $R_{N, \varepsilon}(\rho) \leq R_{N, \varepsilon}(\mathcal{U}(\rho))$ , and then  $I(\rho) \leq I[\mathcal{U}(\rho)]$ . The reverse inequality is now trivial

$$I(\rho) = I[\mathcal{U}^{-1} \mathcal{U}(\rho)] \geq I[\mathcal{U}(\rho)],$$

where we have used the previous result along with the fact that  $\mathcal{U}^{-1}$  is also reversible. ■

#### IV. STEERING: INFORMATION CONTENT FROM DILATIONS

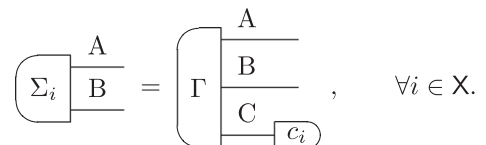
In this section we show that in an OPT satisfying the steering property (assumption 2) the information content of a state can be computed by considering the action of the compression schemes on the set  $D_{\rho^{\boxtimes N}}$  only.

*Lemma IV.1.* Let  $\rho \in \text{St}_1(A)$ ,  $\mathcal{C} \in \text{Tr}_1(A)$  and  $\varepsilon > 0$ . If  $\|\mathcal{C} \boxtimes \mathcal{J}(\Psi) - \Psi\| < \varepsilon$  for any  $\Psi \in D_\rho$  and Assumption 2 holds, then one has

$$\sum_{i \in X} \|\mathcal{C} \boxtimes \mathcal{J}(\Sigma_i) - \Sigma_i\|_{\text{op}} < \varepsilon,$$

for any refinement  $\{\Sigma_i\}_{i \in X}$  of an element of  $D_\rho$ .

*Proof.* Let  $\{\Sigma_i\}_{i \in X}$  be the refinement of an element  $\Omega$  of  $D_\rho$ . By Assumption 2 there exists  $\Gamma \in D_\Omega \subseteq D_\rho$  and an observation test  $\{c_i\}_{i \in X}$  such that



For any  $i \in X$ , let  $A_i \in \text{Eff}(AB)$  be the effect achieving the norm, namely, such that

$$\|\mathcal{C} \boxtimes \mathcal{I}(\Sigma_i) - \Sigma_i\|_{\text{op}} = (A_i | [(\mathcal{C} - \mathcal{I}) \boxtimes \mathcal{I}] | \Sigma_i).$$

Since  $\sum_{i \in X} A_i \boxtimes c_i$  is an effect (see Appendix A), we have that

$$\begin{aligned} & \sum_{i \in X} \|\mathcal{C} \boxtimes \mathcal{I}(\Sigma_i) - \Sigma_i\|_{\text{op}} = \\ & \sum_{i \in X} \left( \Gamma \begin{array}{c} \text{A} \\ \text{B} \\ \text{C} \end{array} \begin{array}{c} \mathcal{C} \\ \mathcal{I} \\ c_i \end{array} \begin{array}{c} \text{A} \\ \text{A}_i \\ \text{C}_i \end{array} \right) - \left( \Gamma \begin{array}{c} \text{A} \\ \text{B} \\ \text{C} \end{array} \begin{array}{c} \text{A}_i \\ \text{C}_i \end{array} \right) \\ & \leq \|\mathcal{C} \boxtimes \mathcal{I}(\Gamma) - \Gamma\|_{\text{op}} < \varepsilon, \end{aligned}$$

which straightforwardly leads to the thesis.  $\blacksquare$

*Proposition IV.1.* Let  $\rho \in \text{St}_1(A)$  and consider

$$I^{\text{dil}}(\rho) := \lim_{\varepsilon \rightarrow 0} \limsup_{N \rightarrow \infty} \frac{\min \{M : E_{N,M,\varepsilon}^{\text{dil}}(\rho) \neq \emptyset\}}{N}, \quad (6)$$

where  $E_{N,M,\varepsilon}^{\text{dil}}(\rho)$  is the set of the compression schemes  $(\mathcal{E}, \mathcal{D})$  such that  $\sup_{\Psi \in D_{\rho, \otimes^N}} \|\mathcal{C} \boxtimes \mathcal{I}(\Psi) - \Psi\| < \varepsilon$ . Then  $I(\rho) = I^{\text{dil}}(\rho)$ .

*Proof.* This is a straightforward consequence of the lemma. Indeed, all the dilations are a refinement of themselves, so that  $E_{N,M,\varepsilon}(\rho) \subseteq E_{N,M,\varepsilon}^{\text{dil}}(\rho)$  and therefore  $I^{\text{dil}}(\rho) \leq I(\rho)$ . The inclusion  $E_{N,M,\varepsilon}(\rho) \supseteq E_{N,M,\varepsilon}^{\text{dil}}(\rho)$  follows by Lemma IV.1, and this implies the thesis.  $\blacksquare$

We now prove some bounds that involve a quantity generalizing the classical and quantum fidelity. For this purpose, we consider a definition of fidelity [19] that can be adopted in the OPT framework, which reduces to the classical or quantum one in the in the respective theories.  $\blacksquare$

*Definition IV.1.* Let  $\rho, \sigma \in \text{St}_1(A)$ . For any observation test  $\{a_i\}_{i \in X} \subseteq \text{Eff}(A)$  denote by  $\mathbf{p} := p_i$  and  $\mathbf{q} := q_i$  the probability distributions defined by

$$\begin{aligned} p_i &:= \left( \rho \begin{array}{c} \text{A} \\ \text{A} \end{array} \begin{array}{c} a_i \end{array} \right), \\ q_i &:= \left( \sigma \begin{array}{c} \text{A} \\ \text{A} \end{array} \begin{array}{c} a_i \end{array} \right). \end{aligned}$$

Then one can define the fidelity between  $\rho$  and  $\sigma$  as

$$F(\rho, \sigma) := \inf_{\{a_i\}_{i \in X}} F_c(\mathbf{p}, \mathbf{q}), \quad (7)$$

where  $F_c(\mathbf{p}, \mathbf{q}) = \sum_i \sqrt{p_i q_i}$ .

Since the classical fidelity is bounded by 1, and is equal to 1 only for  $\mathbf{p} = \mathbf{q}$ , one clearly has  $F(\rho, \sigma) \leq 1$ , with equality if and only if  $\rho = \sigma$ . Fidelity satisfies the following property, that generalizes the Fuchs–van de Graaf inequality [26] and is relevant for the present work.

*Proposition IV.2.* Let  $\rho, \sigma \in \text{St}_1(A)$ . The following inequalities hold:

$$1 - F(\rho, \sigma) \leq \frac{1}{2} \|\rho - \sigma\|_{\text{op}} \leq \sqrt{1 - F(\rho, \sigma)^2}. \quad (8)$$

*Proposition IV.3 (Monotonicity).* Let  $\rho, \sigma \in \text{St}_1(A)$  and  $\mathcal{C} \in \text{Tr}_1(A, B)$ . The following inequality holds:

$$F(\mathcal{C}(\rho), \mathcal{C}(\sigma)) \geq F(\rho, \sigma). \quad (9)$$

In quantum theory we have a notion, the so called *entanglement fidelity*, which measures how well correlations with an

environment are preserved by a given channel acting on our local system. If  $\rho \in \text{St}_1(A)$  is the state of our local system,  $\Phi \in \text{PurSt}(AB)$  is a purification of  $\rho$  and  $\mathcal{C} \in \text{Tr}_1(A \rightarrow C)$  the channel locally applied to  $A$ , then the entanglement fidelity is defined as the square of the Uhlmann one between input and output,  $F[\Phi, \mathcal{C} \boxtimes \mathcal{I}(\Phi)]^2$ . This is a well defined quantity since it is independent of the chosen purification. By means of the generalized notion of fidelity in (7) we can define an analogous of the entanglement fidelity in the OPT framework. Again, we must be aware of the fact that in a generic OPT, there may be states that cannot be purified (mixed states in classical theories are a rather trivial example). In order to encompass the most general situation, we refer to dilations rather than focusing on purifications. Moreover, we want to define a quantity that is independent of the particular dilation, and we are thus led to the following definition.

*Definition IV.2.* Let  $\rho \in \text{St}(A)$  and  $\mathcal{C} \in \text{Tr}_1(A \rightarrow C)$ . We define the correlation fidelity as follows:

$$F(\rho, \mathcal{C}) = \inf_{\Psi \in D_\rho} F[\Psi, \mathcal{C} \boxtimes \mathcal{I}(\Psi)]^2. \quad (10)$$

By means of the generalized Fuchs–van de Graaf inequality (8) we can see that the correlation fidelity can be used as an equivalent figure of merit on  $D_{\rho, \otimes^N}$ . More precisely, the following proposition holds.

*Proposition IV.4.* Let  $\rho \in \text{St}_1(A)$  and define

$$I^F(\rho) := \lim_{\varepsilon \rightarrow 0} \limsup_{N \rightarrow \infty} \frac{\min \{M : E_{N,M,\varepsilon}^F(\rho) \neq \emptyset\}}{N},$$

where

$$E_{N,M,\varepsilon}^F(\rho) := \{(\mathcal{E}, \mathcal{D}) | F(\rho^{\boxtimes N}, \mathcal{D}\mathcal{E}) > 1 - \varepsilon\},$$

then  $I^{\text{dil}}(\rho) = I^F(\rho)$ .

*Proof.* This is a straightforward consequence of Proposition IV.2. Let  $(\mathcal{E}, \mathcal{D}) \in E_{N,M,\varepsilon}^{\text{dil}}(\rho)$ . By the first inequality in (8) we have that

$$F[\Psi, (\mathcal{C} \boxtimes \mathcal{I})(\Psi)]^2 \geq 1 - \varepsilon + \frac{\varepsilon^2}{4},$$

for any  $\Psi \in D_{\rho, \otimes^N}$ , and this implies

$$F(\rho^{\boxtimes N}, \mathcal{C}) > 1 - \varepsilon.$$

Therefore  $(\mathcal{E}, \mathcal{D}) \in E_{N,M,\varepsilon}^F(\rho)$ , namely,  $E_{N,M,\varepsilon}^{\text{dil}}(\rho) \subseteq E_{N,M,\varepsilon}^F(\rho)$ , whence  $I^F(\rho) \leq I^{\text{dil}}(\rho)$ .

Now let  $(\mathcal{E}, \mathcal{D}) \in E_{N,M,\varepsilon}^F$ , then by definition

$$F(\rho^{\boxtimes N}, \mathcal{C}) > 1 - \varepsilon,$$

By the second inequality in Proposition IV.2 we have that

$$\begin{aligned} & \|(\mathcal{C} \boxtimes \mathcal{I})(\Psi) - \Psi\|_{\text{op}} \\ & \leq 2\sqrt{1 - F[\Psi, (\mathcal{C} \boxtimes \mathcal{I})(\Psi)]^2}, \end{aligned}$$

for any  $\Psi \in D_{\rho, \otimes^N}$ , which means that  $(\mathcal{E}, \mathcal{D}) \in E_{N,M,2\sqrt{\varepsilon}}^{\text{dil}}(\rho)$ , and the reverse inequality  $I^F(\rho) \geq I^{\text{dil}}(\rho)$  follows.  $\blacksquare$

## V. INFORMATION CONTENT AND STATE PURITY

In the following we will assume that

$$\mathcal{E} : A^{\otimes N} \rightarrow B^{\otimes M}, \quad \mathcal{D} : B^{\otimes M} \rightarrow A^{\otimes N}.$$

Moreover, we will denote by  $\{\Psi_i\}$  any preparation test of  $A^{\otimes N}C$  such that

$$\sum_i e_C \circ \Psi_i = \rho^{\otimes N}.$$

Finally, for every observation test  $\{a_j\}$  of  $A^{\otimes N}C$ , and for any pair  $(\mathcal{E}, \mathcal{D})$ , let us define the two probability distributions

$$p_{i,j} := \left( \Psi_i \begin{array}{c} A^{\otimes N} \\ \hline C \\ \hline a_j \end{array} \right), \quad (11)$$

and

$$q_{i,j} := \left( \Psi_i \begin{array}{c} A^{\otimes N} \\ \hline \mathcal{E}^N \quad B^{\otimes M} \\ \hline C \quad \mathcal{D}^N \\ \hline A^{\otimes N} \\ \hline a_j \end{array} \right). \quad (12)$$

We can then introduce the following functions that represent the Shannon mutual information between classical random variables  $X$  and  $Y$ , distributed according to  $P(X = x_i, Y = y_j) = p_{i,j}$  or  $X$  and  $\tilde{Y}$ , distributed according to  $P(X = x_i, \tilde{Y} = y_j) = q_{i,j}$ :

$$I(X : Y) := \sum_{i,j} p_{i,j} \log_2 \frac{p_{i,j}}{p_i^X p_j^Y},$$

$$I(X : \tilde{Y}) := \sum_{i,j} q_{i,j} \log_2 \frac{q_{i,j}}{q_i^X q_j^{\tilde{Y}}},$$

where  $p_j^Y$ ,  $q_j^{\tilde{Y}}$ , and  $q_i^X = p_i^X$  denote the elements of the marginal distributions.

**Definition V.1.** Let  $\{\Psi_i\}$  denote a preparation test such that  $\sum_i \Psi_i \in D_{\rho^{\otimes N}}$ . We denote by  $E_{N,M,\delta}^C(\rho)$  the set of those compression schemes such that

$$\sup_{C, \{\Psi_i\}, \{a_j\}} L^{-1} |I(X : Y) - I(X, \tilde{Y})| < \delta,$$

with  $L = \log_2(mn - 1)$  where  $n$  is the cardinality of the preparation test  $\{\Psi_i\}$  and  $m$  that of the test  $\{a_j\}$ . We then define the following quantities:

$$R_{\delta,N}^C(\rho) := \frac{\min \{M \mid E_{N,M,\delta}^C(\rho) \neq \emptyset\}}{N}, \quad (13)$$

$$R_{\delta}^C(\rho) := \limsup_{N \rightarrow \infty} R_{\delta,N}^C(\rho), \quad (14)$$

$$I^C(\rho) := \lim_{\delta \rightarrow 0} R_{\delta}^C(\rho). \quad (15)$$

The last quantity above satisfies the following lemmas.

**Lemma V.1.** Let  $\rho \in \text{St}_1(A)$ . Then  $I(\rho) \geq I^C(\rho)$ .

The proof can be found in Appendix B.

In proving the following lemma and the subsequent proposition we assume that when we compose systems, the size does not increase more than exponentially. More precisely, we formulate the following assumption that will hold in the remainder.

**Assumption 4 (Regular scaling).** For every type of system  $A$ , there exists a constant  $k_A > 0$  such that the size  $D(N) := D_{A^{\otimes N}}$  of the compound system  $A^{\otimes N}$  satisfies  $D(N) \leq k_A D(1)^N$ .

**Lemma V.2.** Let  $\rho \in \text{St}_1(A)$  be a mixed state, then  $I^C(\rho) > 0$ .

The proof can be found in Appendix C.

**Proposition V.1.** Let  $\rho \in \text{St}_1(A)$ . If  $I(\rho) = 0$  then  $\rho$  is a pure state.

*Proof.* Let  $\rho$  be mixed. By Lemmas V.1 and V.2

$$I(\rho) \geq I^C(\rho) > 0,$$

whence the thesis.  $\blacksquare$

We now use the above results to prove some general facts about theories with *essentially unique purification* and *atomicity of parallel composition* for states. For this purpose, let us start reminding in the first place the definition of the latter requirements.

**Definition V.2 (Existence of purification).** We say that an OPT satisfies purification if for any  $\rho \in \text{St}(A)$  one has  $P_{\rho} \neq \emptyset$ .

**Definition V.3 (Essential uniqueness of purification).** We say that an OPT satisfies essential uniqueness of purification if, for any  $\rho \in \text{St}(A)$  such that  $P_{\rho} \neq \emptyset$ ,  $\forall \Phi, \Psi \in P_{\rho}$  with  $\Psi, \Phi \in \text{St}(AB)$ , there exists a reversible transformation  $\mathcal{U}$  such that

$$\left( \Psi \begin{array}{c} A \\ \hline B \end{array} \right) = \left( \Phi \begin{array}{c} A \\ \hline B \\ \hline \mathcal{U} \\ \hline B \end{array} \right). \quad (16)$$

**Definition V.4 (Atomicity of parallel composition of states).** We say that an OPT satisfies atomicity of parallel composition of states if for any pair  $\phi \in \text{PurSt}(A)$  and  $\psi \in \text{PurSt}(B)$  we also have  $\phi \boxtimes \psi \in \text{PurSt}(AB)$ .

We now prove that, in every strongly causal theory that satisfies regular scaling and uniqueness of purification, null information content of pure states is equivalent to atomicity of parallel composition of states. We stress that the requirement of *existence* of purification is not needed, but only its uniqueness. In other words, if a state has a purification, then the latter is unique, even though it needs not have one. The following result is particularly interesting because it provides an alternative way of understanding the operational content of atomicity of parallel composition.

**Proposition V.2.** Let us consider strongly causal OPT satisfying regular scaling. Then the requirements of essential uniqueness of purification and atomicity of parallel composition of states imply that  $I(\phi) = 0$  for any  $\phi \in \text{PurSt}(A)$ . Conversely, if  $I(\phi) = 0$  for any  $\phi \in \text{PurSt}(A)$ , atomicity of parallel composition of states holds.

*Proof.* By Lemma III.1 we have  $I(\phi) = I^{\text{pur}}(\phi)$ . Now, let us fix  $N$  and consider a dilation  $\Omega$  of  $\phi^{\otimes N}$ . Let  $\{\Psi_i\}_{i \in X}$  be a pure decomposition of  $\Omega$ , then by purity of  $\phi^{\otimes N}$  we must have

$$p_i \left( \phi^{\otimes N} \begin{array}{c} A^N \\ \hline \end{array} \right) = \left( \Psi_i \begin{array}{c} A^N \\ \hline B \\ \hline e \end{array} \right), \quad \forall i \in X. \quad (17)$$

Now, let  $\eta \in \text{PurSt}(B)$  and consider  $\phi^{\otimes N} \boxtimes \eta$ . This is still a pure state, hence a purification of  $\phi^{\otimes N}$ . Therefore, by essential



uniqueness of purifications

$$\begin{array}{c} \text{A}^{\otimes N} \\ \text{B} \\ \Psi_i \end{array} = p_i \begin{array}{c} \text{A}^{\otimes N} \\ \text{B} \\ \eta \end{array} \begin{array}{c} \text{A}^{\otimes N} \\ \text{B} \\ \mathcal{U}_i \end{array}, \quad \forall i \in \mathcal{X}. \quad (18)$$

where  $\mathcal{U}_i$  are reversible channels on B.

Now let us consider a compression scheme defined by a measure and prepare one, as follows:

$$\mathcal{E} := \begin{array}{c} \text{A}^{\otimes N} \\ e \end{array}, \quad \mathcal{D} := \begin{array}{c} \text{A}^{\otimes N} \\ \phi^{\otimes N} \end{array}. \quad (19)$$

It is clear that for any dilation  $\Omega$  of  $\phi^{\otimes N}$ , the above scheme is such that  $(\mathcal{D}\mathcal{E} \boxtimes \mathcal{I})(\Psi_i) - \Psi_i = 0$  and this implies that for any  $N$  and  $\varepsilon$  we have  $E_{N,0,\varepsilon}^{\text{pur}}(\rho) \neq \emptyset$  and then  $I(\rho) = I^{\text{pur}}(\rho) = 0$ .

Now, let us assume that for any A and for any  $\phi \in \text{PurSt}(A)$  we have  $I(\phi) = 0$ . Let  $\rho \in \text{PurSt}(A)$  and  $\sigma \in \text{PurSt}(B)$ . By Proposition III.2 we have that

$$I(\rho \boxtimes \sigma) \leq I(\rho) + I(\sigma) = 0.$$

Thus  $I(\rho \boxtimes \sigma) = 0$ , and by proposition V.1  $\rho \boxtimes \sigma$  is pure, namely, Assumption V.4 holds. ■

It is interesting to observe that, due to the above proposition, one can exhibit operational probabilistic theories having pure states with nonvanishing information content. In Ref. [27] the authors construct a theory, *bilocal classical theory*, where all systems are classical (the set of states is a simplex), but with a parallel composition rule that differs from the one of classical information theory, thereby violating Assumption V.4. Accordingly, bilocal classical theory must have pure states with nonnull information content.

## VI. INFORMATION CONTENT IN QUANTUM AND CLASSICAL INFORMATION THEORY

Before restricting to the quantum case, we prove the following lemma concerning the correlation fidelity defined in Definition IV.2.

*Lemma VI.1.* Let  $\rho \in \text{St}_1(A)$  and  $\mathcal{C} \in \text{Tr}_1(A)$ . If every state has a purification (Definition V.2), then one has

$$F(\rho, \mathcal{C}) = \inf_{\Phi \in P_\rho} F[\Phi, \mathcal{C} \boxtimes \mathcal{I}(\Phi)]^2. \quad (20)$$

Moreover, in an OPT with essential uniqueness of purification (Definition V.4) and atomicity of parallel composition of states (Definition V.3), for any  $\Phi \in P_\rho$  one has

$$F(\rho, \mathcal{C}) = F[\Phi, \mathcal{C} \boxtimes \mathcal{I}(\Phi)]^2. \quad (21)$$

*Proof.* If  $P_\tau \neq \emptyset$  for every state  $\tau$ , then for any  $\Psi \in D_\rho$  one has that there exists  $\Gamma \in P_\Psi \subseteq P_\rho$ . Therefore, by monotonicity of the fidelity (Proposition IV.3) we have

$$\begin{aligned} F[\Psi, \mathcal{C} \boxtimes \mathcal{I}(\Psi)]^2 &\geq F[\Gamma, \mathcal{C} \boxtimes \mathcal{I}(\Gamma)]^2 \\ &\geq \inf_{\Phi \in P_\rho} F[\Phi, \mathcal{C} \boxtimes \mathcal{I}(\Phi)]^2. \end{aligned}$$

Since this holds for any  $\Psi \in D_\rho$ , it implies  $F(\rho, \mathcal{C}) \geq \inf_{\Phi \in P_\rho} F[\Phi, \mathcal{C} \boxtimes \mathcal{I}(\Phi)]^2$ . The reverse inequality is trivial, since  $P_\rho \subseteq D_\rho$ .

If all the purifications of  $\rho$  with the same purifying system are connected through a reversible transformation  $\mathcal{U}$  and atomicity of parallel composition of pure states also hold (see Definitions V.4 and V.3), then, for any fixed purification  $\Phi$  in  $P_\rho$ , and any other  $\Gamma$  in  $P_\rho$ , there exists a channel  $\mathcal{A} \in \text{Tr}_1(B, C)$  such that

$$\begin{array}{c} \text{A} \\ \text{C} \\ \Gamma \end{array} = \begin{array}{c} \text{A} \\ \text{B} \\ \Phi \end{array} \begin{array}{c} \text{A} \\ \text{C} \\ \mathcal{A} \end{array}. \quad (22)$$

By monotonicity one has  $F(\rho, \mathcal{C}) \geq F[\Phi, \mathcal{C} \boxtimes \mathcal{I}(\Phi)]^2$  and the reverse inequality is trivial, as  $\Phi \in P_\rho$ . ■

We recall the statement of the Schumacher theorem. Let  $\rho \in \text{St}_1(\rho)$  and  $(\mathcal{E}, \mathcal{D})$  be a compression scheme

*Theorem VI.1 (Schumacher).* Let  $\rho \in \text{St}_1(A)$  with  $\mathcal{H}_A$  the Hilbert space corresponding to the quantum system A, let  $(\mathcal{E}, \mathcal{D})$  be a compression scheme and define its ratio  $R$  as

$$R := \frac{\log[\dim[\text{Supp}(\mathcal{E}(\rho^{\otimes N}))]]}{N}.$$

For every  $\varepsilon > 0$  and  $R > S(\rho)$  there exists  $N_0$  such that  $\forall N \geq N_0$  there exists a compression scheme with ratio  $R$  such that  $F(\rho^{\otimes N}, \mathcal{D}\mathcal{E}) > 1 - \varepsilon$ . Conversely, for every  $R < S(\rho)$  there is  $\varepsilon > 0$  such that for every compression scheme  $(\mathcal{E}, \mathcal{D})$  with ratio  $R$  one has  $F(\rho^{\otimes N}, \mathcal{D}\mathcal{E}) \leq \varepsilon$ .

*Proposition VI.1.* Let  $\rho \in \text{St}_1(A)$  be a quantum state and denote with  $S(\rho)$  its von Neumann entropy. Then  $I(\rho) = S(\rho)$ .

*Proof.* We start by showing that  $I^F(\rho) \leq S(\rho)$ . Let  $\delta > 0$ ,  $R \in (S(\rho), S(\rho) + \delta]$  and  $\varepsilon > 0$ . By the direct part of the Schumacher theorem there exists a  $N_0$  such that for any  $N \geq N_0$  there is a  $(N, \varepsilon)$ -reliable compression scheme with rate  $R$ . Thus, upon embedding  $\text{Supp}[\mathcal{E}(\rho^{\otimes N})]$  in  $[NR]$  qubits, by using a suitable isometry, we have  $E_{N,[NR],\varepsilon}^F(\rho) \neq \emptyset$  for any  $N \geq N_0$ . This implies

$$\begin{aligned} \limsup_{N \rightarrow \infty} \frac{\min \{M : E_{N,M,\varepsilon}^F(\rho) \neq \emptyset\}}{N} \\ \leq \lim_{N \rightarrow \infty} \frac{[NR]}{N} = R \leq S(\rho) + \delta. \end{aligned}$$

Since the argument holds for any  $\varepsilon > 0$ , we get  $I^F(\rho) \leq S(\rho) + \delta$ , and being delta arbitrary, we find  $I^F(\rho) \leq S(\rho)$ .

Now let  $\delta > 0$  and consider  $M/\bar{N}$  such that  $S(\rho) - \delta \leq M/\bar{N} < S(\rho)$ . By the converse part of Schumacher theorem there exists  $\bar{\varepsilon} > 0$  such that for any compression scheme with ratio  $M/\bar{N}$  one has  $F(\rho^{\otimes N}, \mathcal{D}\mathcal{E}) \leq 1 - \bar{\varepsilon}$ . In particular, since any compression scheme from  $k\bar{N}$  copies of the system to  $kM$  qubits has ratio  $M/\bar{N}$ , one has

$$E_{k\bar{N},kM,\bar{\varepsilon}}^F(\rho) = \emptyset, \quad \forall k \in \mathbb{N}.$$

Therefore, for any  $0 < \varepsilon < \bar{\varepsilon}$  and  $k$

$$\begin{aligned} S(\rho) - \delta &\leq \frac{M}{\bar{N}} < \frac{\min \{L : E_{k\bar{N},L,\bar{\varepsilon}}^F(\rho) \neq \emptyset\}}{k\bar{N}} \leq \\ &\leq \frac{\min \{L : E_{k\bar{N},L,\varepsilon}^F(\rho) \neq \emptyset\}}{k\bar{N}}. \end{aligned}$$

Thus, by taking the  $\limsup_{k \rightarrow \infty}$  we find that

$$\begin{aligned} S(\rho) - \delta &\leq \limsup_{k \rightarrow \infty} \frac{\min \{L : E_{kN, L, \varepsilon}^F(\rho) \neq \emptyset\}}{kN} \\ &\leq \limsup_{N \rightarrow \infty} \frac{\min \{L : E_{N, L, \varepsilon}^F(\rho) \neq \emptyset\}}{N}, \end{aligned}$$

for any  $0 < \varepsilon < \bar{\varepsilon}$ . By taking the  $\lim_{\varepsilon \rightarrow 0}$  and the arbitrariness of  $\delta$  we finally get  $S(\rho) \leq I^F(\rho)$ . The statement then follows by the fact that in quantum theory one has  $I^F(\rho) = I(\rho)$  (Propositions IV.1 and IV.4). ■

Let us now turn our focus to the classical case. In this setting, the input and the output of the compression scheme are given by strings of  $N$  letters drawn from an alphabet  $\mathcal{X}$ . Each letter  $x_i$  appears with a given probability  $p_i$  and the probability that the overall string  $x_{i_1} \dots x_{i_N}$  is emitted is given by the joint probability  $p_{i_1, \dots, i_N}$ . If we assume that each symbol is independently and identically distributed, then  $p_{i_1, \dots, i_N} = p_{i_1} \dots p_{i_N}$ . The probability  $p(e)$  of emitting an output string which is different from the input one is often considered in the literature as a figure of merit. More formally, let  $C$  be the Markov matrix representing the composition of the compression and decompression maps, and  $\mathbf{i} := i_1 \dots i_N$  define the input string, then the error probability is defined as

$$p^C(e) := \sum_{\mathbf{i}} \sum_{\mathbf{j} \neq \mathbf{i}} p(\mathbf{i} \neq \mathbf{j} | \mathbf{i}) = 1 - \sum_{\mathbf{i}} C_{\mathbf{i}, \mathbf{i}} p_{\mathbf{i}}.$$

The set of states of classical theory is given by a simplex, and any probability vector representing a state can be uniquely decomposed in terms of pure states  $\mathbf{e}_i$ , corresponding to vectors with all zero components except the one in the  $i$ th position:  $(e_i)_j = \delta_{i,j}$ . Since  $C$  is a stochastic matrix we have the following chain of equalities:

$$\begin{aligned} &\sum_{\mathbf{i}} p_{\mathbf{i}} \|C\mathbf{e}_{\mathbf{i}} - \mathbf{e}_{\mathbf{i}}\| \\ &= \sum_{\mathbf{i}} p_{\mathbf{i}} \sum_{\mathbf{j}} |C_{\mathbf{j}, \mathbf{i}} - \delta_{\mathbf{j}, \mathbf{i}}| \\ &= \sum_{\mathbf{i}} p_{\mathbf{i}} \left( \sum_{\mathbf{j} \neq \mathbf{i}} C_{\mathbf{j}, \mathbf{i}} + 1 - C_{\mathbf{i}, \mathbf{i}} \right) \\ &= \sum_{\mathbf{i}} 2p_{\mathbf{i}}(1 - C_{\mathbf{i}, \mathbf{i}}) \\ &= 2 \sum_{\mathbf{i}} p_{\mathbf{i}}(1 - C_{\mathbf{i}, \mathbf{i}}) = 2 \left( 1 - \sum_{\mathbf{i}} p_{\mathbf{i}} C_{\mathbf{i}, \mathbf{i}} \right) = 2p^C(e), \end{aligned}$$

namely,

$$p^C(e) = \sum_{\mathbf{i}} p_{\mathbf{i}} \frac{1}{2} \|C\mathbf{e}_{\mathbf{i}} - \mathbf{e}_{\mathbf{i}}\|_1.$$

Now consider the unique pure decomposition of some dilation  $\mathbf{\Pi} \in D_{\mathbf{p}^{\otimes N}}$ . This is given by  $\mathbf{\Pi} = \sum_{\mathbf{i}, \mathbf{j}} \Pi_{\mathbf{i}, \mathbf{j}} \mathbf{e}_{\mathbf{i}} \otimes \mathbf{e}_{\mathbf{j}}$  with  $\sum_{\mathbf{j}} \Pi_{\mathbf{i}, \mathbf{j}} = p_{\mathbf{i}}$  (the pure states of the composite system are the tensor product vectors of the pure ones of the composing systems). Then we find

$$\sum_{\mathbf{i}, \mathbf{j}} \Pi_{\mathbf{i}, \mathbf{j}} \| (C \otimes I) \mathbf{e}_{\mathbf{i}} \otimes \mathbf{e}_{\mathbf{j}} - \mathbf{e}_{\mathbf{i}} \otimes \mathbf{e}_{\mathbf{j}} \|_1$$

$$\begin{aligned} &= \sum_{\mathbf{i}, \mathbf{j}} \Pi_{\mathbf{i}, \mathbf{j}} \| (C - I) \mathbf{e}_{\mathbf{i}} \otimes \mathbf{e}_{\mathbf{j}} \|_1 \\ &= \sum_{\mathbf{i}, \mathbf{j}} \Pi_{\mathbf{i}, \mathbf{j}} \| (C - I) \mathbf{e}_{\mathbf{i}} \|_1 \\ &= \sum_{\mathbf{i}} p_{\mathbf{i}} \| (C - I) \mathbf{e}_{\mathbf{i}} \|_1 = 2p^C(e), \end{aligned}$$

having used the fact that for any  $\mathbf{j}$

$$\| (C - I) \mathbf{e}_{\mathbf{i}} \otimes \mathbf{e}_{\mathbf{j}} \|_1 = \| (C - I) \mathbf{e}_{\mathbf{i}} \|_1.$$

Summarizing, we have proved that for any dilation  $\mathbf{\Pi} \in D_{\mathbf{p}^{\otimes N}}$

$$p^C(e) = \frac{1}{2} \sum_{\mathbf{i}, \mathbf{j}} \Pi_{\mathbf{i}, \mathbf{j}} \| (C \otimes I) \mathbf{e}_{\mathbf{i}} \otimes \mathbf{e}_{\mathbf{j}} - \mathbf{e}_{\mathbf{i}} \otimes \mathbf{e}_{\mathbf{j}} \|_1.$$

Namely, in the classical case, the error probability is exactly our figure of merit in (5).

Now, one can use the first Shannon theorem in order to prove that the information content of a classical state is exactly its Shannon entropy

*Theorem VI.2 (Shannon).* Let  $\mathbf{p} \in \text{St}_1(\mathcal{A})$  be a classical state, let  $(\mathcal{E}, \mathcal{D})$  be a compression scheme and define its ratio  $R$  as

$$R := \frac{\log_2 |\mathcal{E}(\text{Rng}(X^N))|}{N}.$$

For every  $\varepsilon > 0$  and  $R > H(X)$  there exists  $N_0$  such that  $\forall N \geq N_0$  there exists a compression scheme with ratio  $R$  such that  $p^C(e) < \varepsilon$ . Conversely, for every  $R < H(X)$  there is  $\varepsilon > 0$  such that for every compression scheme  $(\mathcal{E}, \mathcal{D})$  with ratio  $R$  one has  $p^C(e) \geq \varepsilon$ .

A proposition analogous to Proposition VI.1 holds for the classical case, whose proof is essentially the same. The main issue in both cases is to correctly identify the figure of merit that must be adopted in order to define the information content.

## VII. CONCLUSION

We have defined the information content for a source of information of an arbitrary operational probabilistic theory. The only assumption needed is that of digitizability: A theory is digitizable if any system of the theory can be asymptotically perfectly mapped into finitely many copies of a reference system, called “obit,” playing the role that “bit” and “qubit” play in classical and quantum theory, respectively. The information content of a source is defined as the minimum number of obits needed to store the output of the source in a such a way that it can be recovered with arbitrary accuracy. The figure of merit for establishing accuracy, independently of the features of the theory, is robust against any distortion effect that a compression scheme could induce on the state of the source, on its admissible preparations and on the correlations with external systems. Accordingly, the figure of merit meets the following two criteria: (1) any preparation of ensembles that average to the considered state must be indistinguishable from leaving the preparation untouched and (2) the compression scheme must preserve decompositions of dilations of the state of interest, namely, joint states of the system and arbitrary external systems such that the state that one obtains after

averaging and discarding the external system is the one at hand.

We first proved that the information content is always a well defined quantity. Moreover, in the hypothesis of steering of ensembles, we show that the information content can be computed using simple figures of merit, e.g., a generalization of entanglement fidelity here denoted by correlation fidelity. Then we show that the present notion of information content coincides with the Shannon and von Neumann entropies in the classical and quantum case, respectively. In quantum theory both the entanglement fidelity and the average input-output fidelity for an arbitrary decomposition of the state representing the source identify the von Neumann entropy as the minimal compression rate. This opens a relevant question: Which are the minimal assumptions behind the collapse of “global” and “local” figures of merit, namely, quantifiers of the ability to recover the correlations of the source, and the local preparations of the source, respectively?

Like Shannon’s and von Neumann’s entropy, we proved that the information content is subadditive, and can be used to measure the purity of a state. Indeed both Shannon’s and von Neumann’s entropy vanish if and only if the state is pure, and here we show that the information content has this feature as well. While it is always true that a source with null information content corresponds to a pure state, the opposite implication is satisfied in the presence of atomicity of parallel composition (the parallel composition of any two pure states is pure) and unique purification (if a state has a purification, then the latter is unique up to reversible channels on the remote system).

In the light of the above results we propose the information content as a candidate entropic quantity generalizing Shannon entropy of classical systems and von Neumann entropy of quantum systems. A basic message across the literature on general probabilistic theories [17–19] is that a theory is usually not *monoentropic*, namely, multiple entropic quantities can be defined, each one reducing to Shannon’s and von Neumann’s entropy in classical and quantum theory, respectively. The notions of entropy usually considered are defined in terms of classical information quantities as follows: (1) The *measurement entropy* of a system, namely, the infimum Shannon entropy of any possible measurement on the system, quantifies the minimum measurement uncertainty, provided that the system is prepared in the state of interest. (2) The *decomposition (or mixing) entropy*, namely, the infimum of the Shannon entropies over all possible ways of preparing the system’s state as a mixture of pure states, quantifies the minimum uncertainty for a preparation of a state with respect to pure states. (3) The supremum of the Shannon mutual information between two random variables related, respectively, to measurements on the system and decompositions of the state of interest, quantifies the maximum *accessible information*.

For a general theory the above quantities can be very different and violate some of the typical features of Shannon and von Neumann entropies, such as *concavity* and *strong subadditivity*. It is known [17–19] that measurement entropy is both subadditive and concave but in general it is not strongly subadditive and does not provide a measure of purity, while decomposition entropy is generally neither subadditive nor concave. For example, in Ref. [17] it is shown that, for non-

local boxes, the decomposition entropy is not concave. Less is known about the third entropic quantity given in terms of the Shannon mutual information, which still could satisfy all features of Shannon and von Neumann entropies.

One of the main outcome of this manuscript are a series of results that can be used to clarify which of the possible entropies of a general probabilistic theory has operational meaning in terms of optimal compression ratio. We started here the analysis of the relation between information content and other entropies of a general probabilistic theory focusing on the accessible information. On one hand, both quantities provide a measure of purity of a state, and on the other hand we proved that the accessible information is a lower bound for information content. An important open question is under what conditions the information content, which by definition is the optimal compression ratio, coincides with the accessible information in general.

Finally, we leave the question open as whether nondigitizable theories exist, or one can figure out a counterexample. While we conjecture that the second choice is the case, it is very hard to exhibit a nondigitizable theory, precisely because it is hard to conceive a pair of systems that are not asymptotically equivalent, though it is intuitive that, e.g., a system with a state space that is a polyhedron cannot be asymptotically equivalent to a system whose state space is an ellipsoid. A nondigitizable theory should contain infinitely many system types, all pairwise asymptotically inequivalent.

## ACKNOWLEDGMENTS

A.T. acknowledges financial support from the Elvia and Federico Faggin Foundation through Silicon Valley Community Foundation, Grant No. 2020-214365.

## APPENDIX A: A SIMPLE LEMMA

*Lemma A.1.* Let  $\{c\}_{i \in X} \subseteq \text{Eff}(B)$  be an observation test and  $\{A\}_{i \in X} \subseteq \text{Eff}(A)$  a collection of effects. If causality holds, then  $\sum_{i \in X} A_i \boxtimes c_i \in \text{Eff}(AB)$ .

*Proof.* This is a straightforward consequence of causality. For any  $i \in X$ , there exists an observation test  $\{\tilde{A}_j^{(i)}\}_{j \in Y_i}$  such that  $A_i \in \{\tilde{A}_j^{(i)}\}_{j \in Y_i}$ . Thus, if we consider the test  $\{\mathcal{S}_B \boxtimes c_i\}_{i \in X}$  and the collection of effects  $\{\tilde{A}_j^{(i)} \boxtimes c_i\}_{(i,j) \in X \times Y}$  we have

$$\begin{array}{c} \text{A} \\ \hline \tilde{A}_j^{(i)} \\ \hline \text{B} \\ \hline c_i \end{array} = \begin{array}{c} \text{A} \\ \hline \mathcal{S}_i \\ \hline \text{B} \end{array} \begin{array}{c} \text{A} \\ \hline \mathcal{B}_j^i \\ \hline \text{I} \end{array}, \quad (\text{A1})$$

with  $\mathcal{S}_i := \mathcal{S}_A \boxtimes c_i$  and  $\mathcal{B}_j^i := \tilde{A}_j^{(i)}$ . Therefore, causality implies that  $\{\tilde{A}_j^{(i)} \boxtimes c_i\}_{(i,j) \in X \times Y}$  is an observation test, and  $\sum_{i \in X} A_i \boxtimes c_i \in \text{Eff}(AB)$ , being a coarse graining of effects from the same test. ■

## APPENDIX B: PROOF OF LEMMA VI

We start by defining the following number:

$$\zeta(N, \delta) := \sup \{ \varepsilon \mid E_{N,M,\varepsilon}(\rho) \subseteq E_{N,M,\delta}^C(\rho) \}. \quad (\text{B1})$$

First, we can observe that in the above definition we can safely take the maximum, since the following inclusion holds:

$$E_{N,M,\zeta(N,\delta)}(\rho) \subseteq E_{N,M,\delta}^C(\rho).$$

Indeed, let  $(\mathcal{E}, \mathcal{D}) \in E_{N,M,\zeta(N,\delta)}$ . By definition of  $E_{N,M,\zeta(N,\delta)}$  we of have

$$\sup_{\mathcal{C}, \{\Psi_i\}} \sum_i \|[(\mathcal{D}\mathcal{E} - \mathcal{I}) \boxtimes \mathcal{I}_C] \Psi_i\|_{\text{op}} < \zeta(N, \delta),$$

then there exists  $\varepsilon' < \zeta(N, \delta)$  such that  $\sup_{\mathcal{C}, \{\Psi_i\}} \sum_i \|[(\mathcal{D}\mathcal{E} - \mathcal{I}) \boxtimes \mathcal{I}_C] \Psi_i\|_{\text{op}} < \varepsilon'$ . Thus, by definition of  $\zeta(N, \delta)$ , one has  $(\mathcal{E}, \mathcal{D}) \in E_{N,M,\varepsilon''}(\rho)$  with  $\varepsilon' < \varepsilon'' < \zeta(N, \delta)$  and  $E_{N,M,\varepsilon''}(\rho) \subseteq E_{N,M,\delta}^C(\rho)$ . Finally, since  $E_{N,M,\varepsilon'}(\rho) \subseteq E_{N,M,\varepsilon''}(\rho)$ , we have  $(\mathcal{E}, \mathcal{D}) \in E_{N,M,\delta}^C(\rho)$ , and consequently  $E_{N,M,\zeta(N,\delta)}(\rho) \subseteq E_{N,M,\delta}^C(\rho)$ .

This inclusion has another consequence, which is our starting point for proving the lemma. Indeed, by definition one has

$$\limsup_{N \rightarrow \infty} R_{\zeta(N,\delta),N}(\rho) \geq R_{\delta}^C(\rho).$$

We now have the following two possibilities:

- (1)  $\exists \delta_0 > 0$  such that,  $\forall 0 < \delta < \delta_0$ ,  $\liminf_{N \rightarrow \infty} \zeta(N, \delta) = 0$
- (2)  $\forall \delta > 0$  one has  $\liminf_{N \rightarrow \infty} \zeta(N, \delta) =: \bar{\zeta}(\delta) > 0$ .

Let us start analyzing case 2. In this case, by definition of limit inferior, one has

$$\forall \delta, \gamma > 0 \begin{cases} \exists N_0, \quad \forall N \geq N_0, & \zeta(N, \delta) > \bar{\zeta}(\delta) - \gamma, \\ \forall N_0, \quad \exists N \geq N_0, & \zeta(N, \delta) < \bar{\zeta}(\delta) + \gamma. \end{cases}$$

This implies that for every  $\delta > 0$  and every positive  $\gamma$ , for suitably large  $N$  it is  $R_{\bar{\zeta}(\delta)-\gamma,N}(\rho) \geq R_{\zeta(N,\delta),N}(\rho)$ , and consequently, for suitably large  $N$  it is  $R_{\bar{\zeta}(\delta)/2,N}(\rho) \geq R_{\zeta(N,\delta),N}(\rho)$ . In turn, this implies

$$R_{\bar{\zeta}(\delta)/2}(\rho) \geq \limsup_{N \rightarrow \infty} R_{\zeta(N,\delta),N}(\rho) \geq R_{\delta}^C(\rho),$$

and finally, being  $\bar{\zeta}(\delta)$  increasing as a function of  $\delta$ , taking the limit for  $\delta \rightarrow 0$  one has some value  $\varepsilon \geq 0$  such that

$$I(\rho) \geq \lim_{\zeta \rightarrow \varepsilon} R_{\zeta}(\rho) = \lim_{\delta \rightarrow 0} R_{\bar{\zeta}(\delta)}(\rho) \geq I^C(\rho).$$

We now turn to case 1 and show that this is not possible. The hypotheses imply indeed that there exists  $\delta_0 > 0$  such that  $\liminf_{N \rightarrow \infty} \zeta(N, \delta_0) = 0$ , and the same is then true of every  $0 < \delta \leq \delta_0$ . This means that for every  $\gamma > 0$  and every  $N_0$  there exists  $N \geq N_0$  such that  $\zeta(N, \delta) < \gamma$  for all  $0 < \delta \leq \delta_0$ . By definition, this means that for every  $\gamma$  there exists a scheme  $(\mathcal{E}, \mathcal{D}) \in E_{N,M,\gamma}(\rho)$  such that  $(\mathcal{E}, \mathcal{D}) \notin E_{N,M,\delta}^C(\rho)$ . More explicitly

$$\sup_{\mathcal{C}, \{\Psi_i\}} \sum_i \|[(\mathcal{D}\mathcal{E} - \mathcal{I}) \boxtimes \mathcal{I}_C] \Psi_i\|_{\text{op}} < \gamma, \\ \sup_{\mathcal{C}, \{\Psi_i\}, \{a_j\}} L^{-1} |I(X : Y) - I(X : \tilde{Y})| > \delta,$$

where  $L$  has been introduced in Definition V.1. First, we remark that if  $m = 1$  or  $n = 1$ , then  $H(X) = 0$  or  $H(Y) = H(\tilde{Y}) = 0$ , respectively, and thus  $I(X : Y) = I(X : \tilde{Y}) = 0$ , since  $I(A : B) \leq \min\{H(A), H(B)\}$ . The minimum relevant

value of  $L$  is thus  $\log_2 3$ . Now according to Theorem 2 in [28], for  $\|\mathbf{p} - \mathbf{q}\| < \gamma < 1 - 1/mn$  one has

$$L^{-1} |I(X : Y) - I(X' : Y')| \\ \leq 3\gamma + 3L^{-1} H_2(\gamma) \\ \leq 3\gamma + \frac{3}{\log_2 3} H_2(\gamma),$$

where  $X, Y$  and  $X', Y'$  are distributed according to  $p_{i,j}$  and  $q_{i,j}$ , respectively. We can then conclude that for every  $\gamma > 0$  one has

$$\delta < 3\gamma + \frac{3}{\log_2 3} H(\gamma).$$

However, our hypotheses imply that the latter condition must hold for some  $\delta > 0$ , which is absurd.

### APPENDIX C: PROOF OF LEMMA V.2

Let us take  $\delta > 0$ , and consider  $(\mathcal{E}, \mathcal{D}) \in E_{N,M,\delta}^C(\rho)$ . Let us consider first a single use of the source associated with  $\rho$  corresponding to the decomposition  $\{\Psi_i\}$ , and let  $\{a_j\}$  be the observation test such that  $I(X_0 : Y_0)$  is maximum, where  $X_0$  is the classical variable corresponding to the outcome  $i$  of the preparation test, and  $Y_0$  that of the observation test. Notice that by Krein-Millman's theorem and Caratheodory's theorem one can always find the supremum of mutual information considering atomic decompositions and observation-tests with a bounded number of elements, and thus the optimization problem has a compact domain. Let now  $\{\Psi_i\}$  be the decomposition of  $\rho^{\boxtimes N}$  defined by

$$\Psi_i := \Psi_{i_1} \boxtimes \Psi_{i_2} \boxtimes \dots \boxtimes \Psi_{i_N},$$

and  $\Psi_i$  be the decomposition that maximizes  $I(X_0 : Y_0)$ , with  $m_0$  outcomes. Let now  $\{b_j\}$  be the observation test on  $N$  copies of the system that maximizes  $I(X : Y)$  where  $X$  is the i.i.d. classical variable given by the preparation event  $\mathbf{i}$  and  $Y$  by the outcome  $j$ . Since  $\{b_j\}$  maximizes the mutual information it is clear that the test  $\{(b_j|\mathcal{D}\mathcal{E})\}$  will provide a mutual information  $I(X : \tilde{Y})$  no larger than  $I(X : Y)$ . Thus we can write

$$\delta > \frac{I(X : Y) - I(X : \tilde{Y})}{\log_2 m_0^N D(N) - 1} \\ \geq \frac{I(X : Y) - I(X : \tilde{Y})}{N \log_2 m_0 D_0 + \log_2 k},$$

where in the first bound we used the fact that the number of outcomes for the observation test maximizing the mutual information does not exceed the dimension of the space of effects  $D(N)$ , while in the second bound we used the hypothesis that there exist  $k, D_0$  such that  $D(N) \leq kD_0^N$ . Now, by definition of  $I(X : Y)$  we have  $I(X : Y) \geq NI(X_0 : Y_0)$ , while by the result of Theorem 2 in [29] we have

$$I(X : \tilde{Y}) \leq \log_2 D(M) \leq \log_2 k' + M \log_2 D_1,$$

where we think of the scheme given by the decomposition  $\{\mathcal{E}|\Psi_i\}$  and the observation test given by  $\{(b_j|\mathcal{D})\}$ , involving  $M$  obits. We can then write the following inequality:

$$\delta > \frac{NI(X_0 : Y_0) - M \log_2 D_1 - \log_2 k'}{N \log_2 m_0 D_0 + \log_2 k},$$

and consequently

$$\frac{M}{N} \frac{\log_2 D_1 + \log_2 k'/M}{\log_2 m_0 D_0 + \log_2 k/N} + \delta > \frac{I(X_0 : Y_0)}{\log_2 m_0 D_0 + \log_2 k/N}.$$

In particular, if the scheme  $(\mathcal{E}, \mathcal{D})$  has the minimum  $M$  for fixed  $N, \delta$  we can then conclude that

$$\begin{aligned} R_{\delta, N}^C \frac{\log_2 D_1}{\log_2 m_0 D_0 + \frac{\log_2 k}{N}} + \delta + \frac{1}{N} \frac{\log_2 k'}{\log_2 m_0 D_0 + \frac{\log_2 k}{N}} \\ > \frac{I(X_0 : Y_0)}{\log_2 m_0 D_0 + \frac{\log_2 k}{N}}. \end{aligned}$$

Taking the limit superior for  $N \rightarrow \infty$  on both sides we have

$$R_{\delta}^C \frac{\log_2 D_1}{\log_2 m_0 D_0} + \delta \geq \frac{I(X_0 : Y_0)}{\log_2 m_0 D_0},$$

and finally, in the limit  $\delta \rightarrow 0$  we obtain

$$I^C(\rho) \frac{\log_2 D_1}{\log_2 m_0 D_0} \geq \frac{I(X_0 : Y_0)}{\log_2 m_0 D_0},$$

namely,

$$I^C(\rho) \geq \frac{I(X_0 : Y_0)}{\log_2 D_1}.$$

For a mixed state,  $I(X_0 : Y_0) > 0$ , and this implies the thesis.

- 
- [1] C. E. Shannon, Communication in the presence of noise, *Proc. IRE* **37**, 10 (1949).
- [2] B. Schumacher, Quantum coding, *Phys. Rev. A* **51**, 2738 (1995).
- [3] M. Koashi and N. Imoto, Compressibility of Quantum Mixed-State Signals, *Phys. Rev. Lett.* **87**, 017902 (2001).
- [4] H. Barnum, C. M. Caves, C. A. Fuchs, R. Jozsa, and B. Schumacher, On quantum coding for ensembles of mixed states, *J. Phys. A: Math. Gen.* **34**, 6767 (2001).
- [5] Z. B. Khanian and A. Winter, Entanglement-assisted quantum data compression, in *2019 IEEE International Symposium on Information Theory (ISIT)* (IEEE, Piscataway, NJ, 2019), pp. 1147–1151.
- [6] B. Schumacher and M. Westmoreland, *Quantum Processes Systems, and Information* (Cambridge University Press, Cambridge, 2010).
- [7] L. Hardy and W. K. Wootters, Limited holism and real-vector-space quantum theory, *Found. Phys.* **42**, 454 (2012).
- [8] G. M. D'Ariano, F. Manessi, P. Perinotti, and A. Tosini, Fermionic computation is non-local tomographic and violates monogamy of entanglement, *EPL (Europhys. Lett.)* **107**, 20009 (2014).
- [9] G. M. D'Ariano, F. Manessi, P. Perinotti, and A. Tosini, The Feynman problem and fermionic entanglement: Fermionic theory versus qubit theory, *Int. J. Mod. Phys. A* **29**, 1430025 (2014).
- [10] G. M. D'Ariano, G. Chiribella, and P. Perinotti, *Quantum Theory from First Principles: An Informational Approach* (Cambridge University Press, Cambridge, 2017).
- [11] G. Chiribella, G. M. D'Ariano, and P. Perinotti, Probabilistic theories with purification, *Phys. Rev. A* **81**, 062348 (2010).
- [12] G. Chiribella, G. M. D'Ariano, and P. Perinotti, Informational derivation of quantum theory, *Phys. Rev. A* **84**, 012311 (2011).
- [13] G. M. D'Ariano, M. Erba, and P. Perinotti, Classical theories with entanglement, *Phys. Rev. A* **101**, 042118 (2020).
- [14] L. Hardy, Disentangling nonlocality and teleportation, [arXiv:quant-ph/9906123](https://arxiv.org/abs/quant-ph/9906123) (1999).
- [15] R. W. Spekkens, Evidence for the epistemic view of quantum states: A toy theory, *Phys. Rev. A* **75**, 032110 (2007).
- [16] J. Barrett, Information processing in generalized probabilistic theories, *Phys. Rev. A* **75**, 032304 (2007).
- [17] H. Barnum, J. Barrett, L. Orloff Clark, M. Leifer, R. Spekkens, N. Stepanik, A. Wilce, and R. Wilke, Entropy and information causality in general probabilistic theories, *New J. Phys.* **12**, 033024 (2010).
- [18] A. J. Short and S. Wehner, Entropy in general physical theories, *New J. Phys.* **12**, 3023 (2010).
- [19] G. Kimura, K. Nuida, and H. Imai, Distinguishability measures and entropies for general probabilistic theories, *Rep. Math. Phys.* **66**, 175 (2010).
- [20] G. M. D'Ariano, P. Perinotti, and A. Tosini, Information and disturbance in operational probabilistic theories, *Quantum* **4**, 363 (2020).
- [21] S. Popescu and D. Rohrlich, Quantum nonlocality as an axiom, *Found. Phys.* **24**, 379 (1994).
- [22] G. M. D'Ariano and A. Tosini, Testing axioms for quantum theory on probabilistic toy-theories, *Quant. Info. Proc.* **9**, 95 (2010).
- [23] P. Perinotti, A. Tosini, and L. Vaglini, Shannon theory for quantum systems and beyond: Information compression for fermions, [arXiv:2106.04964](https://arxiv.org/abs/2106.04964) (2021).
- [24] D. Gross, M. Müller, R. Colbeck, and O. C. O. Dahlsten, All Reversible Dynamics in Maximally Nonlocal Theories are Trivial, *Phys. Rev. Lett.* **104**, 080402 (2010).
- [25] P. Perinotti, Cellular automata in operational probabilistic theories, *Quantum* **4**, 294 (2020).
- [26] C. Fuchs and J. van de Graaf, Cryptographic distinguishability measures for quantum-mechanical states, *IEEE Trans. Inf. Theory* **45**, 1216 (1999).
- [27] G. M. D'Ariano, M. Erba, and P. Perinotti, Classicality without local discriminability: Decoupling entanglement and complementarity, *Phys. Rev. A* **102**, 052216 (2020).
- [28] Z. Zhang, Estimating mutual information via Kolmogorov distance, *IEEE Trans. Inf. Theory* **53**, 3280 (2007).
- [29] S. Fiorini, S. Massar, M. K. Patra, and H. R. Tiwary, Generalized probabilistic theories and conic extensions of polytopes, *J. Phys. A: Math. Theor.* **48**, 025302 (2015).