

On the entropy of a noisy function

1

Alex Samorodnitsky

Abstract

Let $0 < \epsilon < 1/2$ be a noise parameter, and let T_ϵ be the noise operator acting on functions on the boolean cube $\{0, 1\}^n$. Let f be a nonnegative function on $\{0, 1\}^n$. We upper bound the entropy of $T_\epsilon f$ by the average entropy of conditional expectations of f , given sets of roughly $(1 - 2\epsilon)^2 \cdot n$ variables.

In information-theoretic terms, we prove the following strengthening of "Mrs. Gerber's lemma": Let X be a random binary vector of length n , and let Z be a noise vector, corresponding to a binary symmetric channel with crossover probability ϵ . Then, setting $v = (1 - 2\epsilon)^2 \cdot n$, we have (up to lower-order terms):

$$H(X \oplus Z) \geq n \cdot H_2 \left(\epsilon + (1 - 2\epsilon) \cdot H_2^{-1} \left(\frac{\mathbb{E}_{|B|=v} H(\{X_i\}_{i \in B})}{v} \right) \right)$$

Assuming $\epsilon \geq 1/2 - \delta$, for some absolute constant $\delta > 0$, this inequality, combined with a strong version of a theorem of Friedgut, Kalai, and Naor, due to Jendrej, Oleszkiewicz, and Wojtaszczyk, shows that if a boolean function f is close to a characteristic function g of a subcube of dimension $n - 1$, then the entropy of $T_\epsilon f$ is at most that of $T_\epsilon g$.

Taken together with a recent result of Ordentlich, Shayevitz, and Weinstein, this shows that the "Most informative boolean function" conjecture of Courtade and Kumar holds for high noise $\epsilon \geq 1/2 - \delta$.

Namely, if X is uniformly distributed in $\{0, 1\}^n$ and Y is obtained by flipping each coordinate of X independently with probability ϵ , then, provided $\epsilon \geq 1/2 - \delta$, for any boolean function f holds $I(f(X); Y) \leq 1 - H(\epsilon)$.

I. INTRODUCTION

This paper is motivated by the following conjecture of Courtade and Kumar [7].

Let (X, Y) be jointly distributed in $\{0, 1\}^n$ such that their marginals are uniform and Y is obtained by flipping each coordinate of X independently with probability ϵ . Let H_2 denote the binary entropy function $H_2(x) = -x \log_2 x - (1 - x) \log_2(1 - x)$. The conjecture of [7] is:

Conjecture 1.1: For all boolean functions $f : \{0, 1\}^n \rightarrow \{0, 1\}$,

$$I(f(X); Y) \leq 1 - H_2(\epsilon)$$

■

This inequality holds with equality if f is a characteristic function of a subcube of dimension $n - 1$. Hence, the conjecture is that such functions are the "most informative" boolean functions.

Following [9], we express $I(f(X); Y)$ in terms of the 'value of the entropy functional of the image of f under the noise operator' (all notions will be defined shortly). The question then becomes:

Which boolean functions are the "stabllest" under the action of the noise operator? That is, for which functions the entropy functional decreases the least under noise.

One can also consider a more general question of how the noise operator affects the entropy of a nonnegative function.

School of Engineering and Computer Science, The Hebrew University of Jerusalem, Jerusalem 91904, Israel. Research partially supported by ISF grants 1241/11 and 1724/15, and by BSF grant 2010451. Copyright (c) 2015 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

Our main result is that for a nonnegative function f on $\{0, 1\}^n$, the entropy of the image of f under the noise operator with noise parameter ϵ is upper bounded by the average entropy of conditional expectations of f , given sets of roughly $(1 - 2\epsilon)^2 \cdot n$ variables.

As an application, using the recent strengthening [6] of a theorem of [4], we show that for ϵ close to $1/2$ characteristic functions of $(n - 1)$ -dimensional subcubes are at least as stable under the noise operator as functions which are close to them.

This, in conjunction with [4] and a recent result of [14] which can be used to show that, for high noise levels $\epsilon \sim 1/2$, boolean functions, which are potentially as stable as the characteristic functions of $(n - 1)$ -dimensional subcubes, have to be close to these functions, implies the validity of Conjecture 1.1 for high noise levels.

A. Entropy of nonnegative functions and the noise operator

We introduce some relevant notions.

For a nonnegative function $f : \{0, 1\}^n \rightarrow \mathbb{R}$, we let the *entropy* of f to be defined as

$$\text{Ent}(f) = \mathbb{E}_x f(x) \log_2 f(x) - \mathbb{E}_x f(x) \cdot \log_2 \left(\mathbb{E} f(x) \right)$$

We note for future use that entropy is nonnegative, homogeneous $\text{Ent}(\lambda f) = \lambda \cdot \text{Ent}(f)$ and convex in f [8].

Given $0 \leq \epsilon \leq 1/2$, we define the *noise operator* acting on functions on the boolean cube as follows: for $f : \{0, 1\}^n \rightarrow \mathbb{R}$, we let $T_\epsilon f$ at a point x be the expected value of f at y , where y is ϵ -correlated with x . That is,

$$(T_\epsilon f)(x) = \sum_{y \in \{0, 1\}^n} \epsilon^{|y-x|} \cdot (1 - \epsilon)^{n-|y-x|} \cdot f(y) \quad (1)$$

Here $|\cdot|$ denotes the Hamming distance.

Note that $T_\epsilon f$ is a convex combination of shifted copies of f . Hence, convexity of entropy implies that the noise operator decreases entropy. Our goal is to quantify this statement.

1) *Connection between notions:* Let f be a nonnegative function on $\{0, 1\}^n$. Let X be a random variable on $\{0, 1\}^n$ distributed according to $f / \sum f$. Let Z be an independent noise random variable on $\{0, 1\}^n$. That is, $\Pr\{Z = z\} = \epsilon^{|z|} \cdot (1 - \epsilon)^{n-|z|}$, and X and Z are statistically independent. Then

$$\text{Ent}(f) = \mathbb{E} f \cdot (n - H(X))$$

$$\text{Ent}(T_\epsilon f) = \mathbb{E} f \cdot (n - H(X \oplus Z))$$

Let now $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a boolean function, let X be uniformly distributed in $\{0, 1\}^n$, let Z be an independent noise random variable, and let $Y = X \oplus Z$. Then

$$H(f(X)) = \text{Ent}(f) + \text{Ent}(1 - f)$$

We also have the following simple claim (proved in Section VI below)

Lemma 1.2.: For a boolean function $f : \{0, 1\}^n \rightarrow \{0, 1\}$,

$$I(f(X); Y) = \text{Ent}(T_\epsilon f) + \text{Ent}(T_\epsilon(1 - f))$$

Therefore, Conjecture 1.1 translates as follows:

Conjecture 1.3.: (An equivalent form of Conjecture 1.1)

For any boolean function $f : \{0, 1\}^n \rightarrow \{0, 1\}$ holds

$$\text{Ent}(T_\epsilon f) + \text{Ent}(T_\epsilon(1-f)) \leq 1 - H_2(\epsilon)$$

■

B. Mrs. Gerber's function and Mrs. Gerber's lemma

We describe a result from information theory, and a related function, which will be important for us ¹.

Let f_t be a function on the two-point space $\{0, 1\}$, which is t at zero and $2-t$ at one. We have

$$\text{Ent}(f_t) = 1 - H_2\left(\frac{t}{2}\right)$$

Let $\phi(x, \epsilon)$ be a function on $[0, 1] \times [0, 1/2]$ defined as follows:

$$\phi(x, \epsilon) = \text{Ent}(T_\epsilon f_t) \tag{2}$$

where t is chosen so that $\text{Ent}(f_t) = x$.

This function was introduced in [21]. We will now describe some of its properties.

Note that ϕ is increasing in x , starting from zero at $x = 0$.

In fact, it is easy to derive the following explicit expression for ϕ :

$$\phi(x, \epsilon) = 1 - H_2\left((1-2\epsilon) \cdot H_2^{-1}(1-x) + \epsilon\right)$$

A key property of ϕ is its concavity.

Theorem 1.4.: ([21]) The function $\phi(x, \epsilon)$ is concave in x for any $0 \leq \epsilon \leq 1/2$.

We mention a simple corollary.

Corollary 1.5.: For all $0 \leq \epsilon \leq 1/2$,

$$(1 - H_2(\epsilon)) \cdot x \leq \phi(x, \epsilon) \leq (1 - 2\epsilon)^2 \cdot x \tag{3}$$

Proof: It's easy to check $\phi(0, \epsilon) = 0$ and $\phi(1, \epsilon) = 1 - H_2(\epsilon)$. And, it's easy to check that $\frac{\partial \phi}{\partial x}$ at $x = 0$ is $(1 - 2\epsilon)^2$. ■

From now on, when the value of ϵ is clear from the context, we omit the second parameter in ϕ and write $\phi(x)$ instead of $\phi(x, \epsilon)$.

We now describe an inequality of [21], which is known as Mrs. Gerber's lemma. Following this usage, we will refer to the function ϕ as Mrs. Gerber's function.

This inequality upperbounds the entropy of the image of a nonnegative function under the action of the noise operator. We present it in terms of the entropy functional and the noise operator².

Theorem 1.6.: ([21]) Let f be a nonnegative function on $\{0, 1\}^n$. Then

$$\text{Ent}(T_\epsilon f) \leq n \mathbb{E} f \cdot \phi\left(\frac{\text{Ent}(f)}{n \mathbb{E} f}, \epsilon\right) \tag{4}$$

¹We are grateful to V. Chandar [3] for explaining the relevance of this result in connection to our previous work [18] on the subject.

²As pointed out to us by Chandar [3], this is equivalent to the standard information-theoretic formulation: Let X be a random binary vector of length n distributed according to $f/\sum f$, and let Z be a noise vector, corresponding to a binary symmetric channel with crossover probability ϵ . Then $H(X \oplus Z) \geq nH_2\left(\epsilon + (1-2\epsilon) \cdot H_2^{-1}\left(\frac{H(X)}{n}\right)\right)$.

C. Main results

For $A \subseteq [n]$ and for a nonnegative function $f : \{0, 1\}^n \rightarrow \mathbb{R}$, we denote

$$\mathbb{E}(f | A) = \mathbb{E}(f | \{x_i\}_{i \in A})$$

Here \mathbb{E} is the conditional expectation operator. That is, $\mathbb{E}(f | A)$ is the function of the variables $\{x_i\}_{i \in A}$, defined as the expectation of f given the values of $\{x_i\}$.³

We write

$$\text{Ent}(f | A) = \text{Ent}\left(\mathbb{E}(f | A)\right)$$

To connect notions, observe that if X is a random variable on $\{0, 1\}^n$ distributed according to $f / \sum f$, then the distribution of $\{X_i\}_{i \in A}$ on the $|A|$ -dimensional cube is given by $\frac{1}{2^{|A|} \mathbb{E}f} \cdot \mathbb{E}(f | A)$ and that

$$\text{Ent}(f | A) = \mathbb{E}f \cdot \left(|A| - H(\{X_i\}_{i \in A})\right) \quad (5)$$

Our main claim is that the entropy of a nonnegative function f under noise is upper bounded by the average entropy of conditional expectations of f , given certain random subsets of variables. We present several results which illustrate this fact.

Theorem 1.7: Let f be a nonnegative function on the cube with $\mathbb{E}f = 1$.

Let $0 < \epsilon < 1$ be a noise parameter. Let T be a random subset of $[n]$ generated by sampling each element $i \in [n]$ independently with probability $(1 - 2\epsilon)^2$. Then

$$\text{Ent}(T_\epsilon f) \leq \mathbb{E}_T \left(\text{Ent}(f | T) - \sum_{i \in T} \text{Ent}(f | \{i\}) \right) + \sum_{i=1}^n \phi \left(\text{Ent}(f | \{i\}) \right)$$

Remark 1.8: We are grateful to O. Ordentlich for suggesting this formulation for the claim of this theorem, as well as for Theorem 1.12 below (in earlier versions the average on the RHS was taken over sets of a fixed cardinality $\sim (1 - 2\epsilon)^2 \cdot n$, which led to more cumbersome calculations.)

Let us also mention that Polyanskiy and Wu [17] came up with a new and direct proof of the key claim, Proposition 4.1, which does not rely on linear programming, and this was used by Ordentlich [12] to give direct proofs for Theorems 1.7 and 1.12. ■

Applying the inequality $\phi(x, \epsilon) \leq (1 - 2\epsilon)^2 \cdot x$ (see (3)) to the claim of the theorem, gives the following, more streamlined claim. (However, the somewhat stronger claim of the theorem is needed for the applications.)

Corollary 1.9: In the notation of Theorem 1.7,

$$\text{Ent}(T_\epsilon f) \leq \mathbb{E}_T \text{Ent}(f | T)$$

Specializing to boolean functions, this implies the following claim.

Corollary 1.10: In the notation of Conjecture 1.1 and of Theorem 1.7, for a boolean function $f : \{0, 1\}^n \rightarrow \{0, 1\}$ holds

$$I(f(X); Y) \leq \mathbb{E}_T I(f(X); \{X_i\}_{i \in T})$$

Remark 1.11: Let B be a random subset of $[n]$ generated by sampling each element $i \in [n]$ independently with probability $1 - 2\epsilon$.

³We also may (and will) view $\mathbb{E}(f | A)$ as a function on $\{0, 1\}^n$, which depends only on variables with indices in A .

As pointed out by Or Ordentlich [12], it seems instructive to compare the bound in Corollary 1.10 to the weaker bound

$$I(f(X); Y) \leq \mathbb{E}_B I(f(X); \{X_i\}_{i \in B})$$

which can be obtained by the following information-theoretic argument.

An equivalent way to obtain Y from X is to replace each coordinate of X independently with a random bit, with probability 2ϵ .

Let S be the set of indices where the input bits were replaced with random bits, and let $B = S^c$.

Using the chain rule of mutual information we have

$$I(f(X); Y) = I(f(X); Y, S) - I(f(X); S | Y) = I(f(X); Y | S) - I(f(X); S | Y)$$

where the last equality follows since $I(f(X); S) = 0$.

In particular, by non-negativity of mutual information

$$I(f(X); Y) \leq I(f(X); Y | S) = \mathbb{E}_B I(f(X); \{X_i\}_{i \in B})$$

■

We also show a somewhat different strengthening of Corollary 1.9, which gives a stronger version of Mrs. Gerber's lemma (Theorem 1.6).

Theorem 1.12.: In the notation of Theorem 1.7, setting $t = (1 - 2\epsilon)^2 \cdot n$, the following is true:

$$\text{Ent}(T_\epsilon f) \leq n \cdot \phi \left(\frac{\mathbb{E}_T \text{Ent}(f | T)}{t}, \epsilon \right)$$

In the standard information-theoretic notation, this could be restated as follows. Let X be a random binary vector of length n , and let Z be an independent noise vector, corresponding to a binary symmetric channel with crossover probability ϵ . Then

$$H(X \oplus Z) \geq n \cdot H_2 \left(\epsilon + (1 - 2\epsilon) \cdot H_2^{-1} \left(\frac{\mathbb{E}_T H(\{X_i\}_{i \in T})}{t} \right) \right) \quad (6)$$

We refer to [13] for an application of (6).

Remark 1.13.:

Up to a negligible error term, the claim of the theorem is stronger than that of Theorem 1.6, since the sequence $a_t = \frac{\mathbb{E}_{|T|=t} H(\{X_i\}_{i \in T})}{t}$ is increasing, by Han's inequality [5].

■

We now return to Conjectures 1.1 and 1.3.

Let us first describe a family of functions for which these conjectures are known to hold with equality. Let $1 \leq k \leq n$ be an index, and let $g_k(x) = 1$ if and only if $x_k = 0$. (That is, g_k is a characteristic function of the $(n - 1)$ -dimensional subcube $\{x_k = 0\}$.)

It is easy to verify that $\text{Ent}(T_\epsilon g_k) = \frac{1}{2} \cdot (1 - H_2(\epsilon))$ and $\text{Ent}(T_\epsilon g_k) + \text{Ent}(T_\epsilon (1 - g_k)) = 1 - H_2(\epsilon)$.

We apply Theorem 1.7 to show that, for $\epsilon \sim 1/2$, the conjectures also hold for functions which are close to characteristic functions of subcubes.

To make the notion of proximity more precise, recall (see [11]) that any function $f : \{0, 1\}^n \rightarrow \mathbb{R}$ can be expanded in terms of the Walsh-Fourier basis: $f(x) = \sum_{S \subseteq [n]} \widehat{f}(S) \cdot W_S(x)$. Here $W_S(x) = (-1)^{\sum_{i \in S} x_i}$.

The Walsh-Fourier expansion of g_k is especially simple: $\widehat{g}_k(0) = \mathbb{E} g_k = 1/2$, $\widehat{g}_k(\{k\}) = 1/2$, and $\widehat{g}_k(S) = 0$ for all other $S \subseteq [n]$.

It follows from [6] and [4] that a boolean function whose Walsh-Fourier expansion is close to that of g_k , in that it has a large (i.e., close to 1/2) Fourier coefficient at $\{k\}$, has to be very close, in the appropriate sense, to g_k .

The next claim shows the conjectures to hold for such functions.

Theorem 1.14.: There exists an absolute constant $\delta > 0$ such that for any noise $\epsilon \geq 0$ with $(1 - 2\epsilon)^2 \leq \delta$ and for any boolean function $f : \{0, 1\}^n \rightarrow \{0, 1\}$ such that

- $\frac{1}{2} - \delta \leq \mathbb{E} f \leq \frac{1}{2}$;
- There exists $1 \leq k \leq n$ such that $|\widehat{f}(\{k\})| \geq (1 - \delta) \cdot \mathbb{E} f$

Holds

1)

$$\text{Ent}(T_\epsilon f) \leq \frac{1}{2} \cdot (1 - H_2(\epsilon))$$

2)

$$\text{Ent}(T_\epsilon f) + \text{Ent}(T_\epsilon(1 - f)) \leq 1 - H_2(\epsilon)$$

This, in conjunction with [4] and [14], which can be used to show that, for noise parameter close to 1/2, boolean functions, which are potentially as stable as the characteristic functions of $(n - 1)$ -dimensional subcubes, have to satisfy the constraints of Theorem 1.14, implies the validity of Conjecture 1.1 for high noise levels.

Theorem 1.15.: There exists an absolute constant $\delta > 0$ such that for any noise $\epsilon \geq 0$ with $(1 - 2\epsilon)^2 \leq \delta$ and for any boolean function $f : \{0, 1\}^n \rightarrow \{0, 1\}$ holds

$$I(f(X); Y) \leq 1 - H_2(\epsilon)$$

D. More on Theorems 1.7 and 1.12

In this subsection we give a high-level description of the proofs of these theorems and argue that both their claims may be viewed as strengthenings of Mrs. Gerber's lemma.

Notation: For a direction $1 \leq i \leq n$ we define the noise operator in direction i as follows:

$$(T_{\epsilon_{\{i\}}} f)(x) = \epsilon \cdot f(x + e_i) + (1 - \epsilon) \cdot f(x)$$

where e_i is the i^{th} unit vector. The operators $\{T_{\epsilon_{\{i\}}}\}$ commute and, for $R \subseteq [n]$, we define T_{ϵ_R} to be the composition of $T_{\epsilon_{\{i\}}}$, $i \in R$. Note that the noise operator T_ϵ would be written in this notation as $T_{\epsilon_{[n]}}$.

We start with the proof of Mrs. Gerber's lemma (4). Since both sides of the inequality are homogeneous in f , we may assume $\mathbb{E} f = 1$.

By the chain rule for entropy, for any permutation σ in the symmetric group S_n holds

$$\text{Ent}(T_\epsilon f) = \sum_{i=1}^n \left(\text{Ent}(T_\epsilon f \mid \{\sigma(1), \dots, \sigma(i)\}) - \text{Ent}(T_\epsilon f \mid \{\sigma(1), \dots, \sigma(i-1)\}) \right) =$$

$$\sum_{i=1}^n \left(Ent\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i)\}\right) - Ent\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i-1)\}\right) \leq^7$$

$$\sum_{i=1}^n \phi \left(Ent\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i)\}\right) - Ent\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i-1)\}\right) \right) \quad (7)$$

Let us explain the last inequality. Let $y \in \{0, 1\}^{i-1}$. Let \tilde{f}_y be a function on $\{0, 1\}$ defined by the restriction of the function $\mathbb{E}\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i)\}\right)$, which we view as a function on the i -dimensional cube, to the points in which the coordinates $\sigma(k)$, $k = 1, \dots, i-1$ are set to be y_k . Then, it is easy to see that

$$Ent\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i)\}\right) - Ent\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i-1)\}\right) =$$

$$\mathbb{E}_y Ent\left(T_{\epsilon} \tilde{f}_y\right) = \mathbb{E}_y \left(\mathbb{E} \tilde{f}_y \cdot \phi \left(Ent \left(\frac{\tilde{f}_y}{\mathbb{E} \tilde{f}_y} \right) \right) \right) \leq \phi \left(\mathbb{E}_y Ent(\tilde{f}_y) \right) =$$

$$\phi \left(Ent\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i)\}\right) - Ent\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i-1)\}\right) \right)$$

The first equality in the second row follows from (2) and the linearity of entropy. The inequality follows from concavity of the function ϕ and the fact that $\mathbb{E}_y \mathbb{E} \tilde{f}_y = \mathbb{E}\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i)\}\right) = \mathbb{E} f = 1$.

We now continue from (7).

For $y \in \{0, 1\}^{i-1}$, let f_y be a function on $\{0, 1\}$ defined by the restriction of the function $\mathbb{E}\left(f \mid \{\sigma(1), \dots, \sigma(i)\}\right)$ to the points in which the coordinates $\sigma(k)$, $k = 1, \dots, i-1$ are set to be y_k .

Since the noise operator $T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}}$ is stochastic, the functions $\{\tilde{f}_y\}$ are a stochastic mixture of the functions $\{f_y\}$. Hence, since the Ent functional is convex, for any $0 \leq \epsilon \leq 1$ holds

$$Ent\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i)\}\right) - Ent\left(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}}\right) f \mid \{\sigma(1), \dots, \sigma(i-1)\}\right) =$$

$$\mathbb{E}_y Ent(\tilde{f}_y) \leq \mathbb{E}_y Ent(f_y) = \quad (8)$$

$$Ent(f \mid \{\sigma(1), \dots, \sigma(i)\}) - Ent(f \mid \{\sigma(1), \dots, \sigma(i-1)\})$$

And hence (7) is upper bounded by

$$\sum_{i=1}^n \phi \left(Ent(f \mid \{\sigma(1), \dots, \sigma(i)\}) - Ent(f \mid \{\sigma(1), \dots, \sigma(i-1)\}) \right) \leq n \cdot \phi \left(\frac{Ent(f)}{n} \right)$$

where in the last inequality the concavity of ϕ is used again.

1) *Our improvement:* We attempt to quantify the loss in inequality (8).

Let us introduce some notation. For a nonnegative function g on the cube, for a subset $A \subset [n]$, and for an element $m \notin A$, we define

$$I_g(A, m) = \text{Ent}(g \mid A \cup \{m\}) - \text{Ent}(g \mid A) - \text{Ent}(g \mid \{m\})$$

This quantity is always nonnegative. In fact, let X be distributed on $\{0, 1\}^n$ according to $g/\sum g$. Assume $\mathbb{E}g = 1$ and note that in this case, by Subsection I-A1 and by (5), we have $I_g(A, m) = H(\{X_i\}_{i \in A}) + H(X_m) - H(\{X_j\}_{j \in A \cup \{m\}}) = I(\{X_i\}_{i \in A}; X_m)$.

Coming back to (8), observe that $\text{Ent}(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}} f \mid \{\sigma(i)\}) = \text{Ent}(f \mid \{\sigma(i)\})$.

Hence, taking $A = \{\sigma(1), \dots, \sigma(i-1)\}$ and $m = \{\sigma(i)\}$, the decrease in (8) is from $I_f(A, m)$ to $I_{T_{\epsilon_A} f}(A, m)$. Therefore, our goal is to quantify the decrease in mutual information in the presence of noise.

In the next two sections we consider a somewhat more general question of upper bounding $I_{T_{\epsilon_A} f}(A, m)$, given f , A , and m . In Section II we upper bound $I_{T_{\epsilon_A} f}(A, m)$ by the value of a certain linear program. In Section III we introduce a symmetric version of this program and a symmetric solution for the symmetric program, and show its value to be at least as large as that of the original program.

We then find the value of the symmetric solution, as a function of f , A , and m . This value provides an upper bound on the noisy mutual information (see Proposition 4.1).

In order to prove Theorems 1.7 and 1.12 we apply the improved bound in (8), averaging the chain rule for the entropy of $T_{\epsilon} f$ over all permutations $\sigma \in S_n$.

This improvement in (8) is the reason we suggest to view both these claims as stronger versions of Mrs. Gerber's lemma.

On the other hand, strictly speaking, this line of argument does not necessarily provide a direct improvement of (4), since in the averaging step we have to replace $\phi(x, \epsilon)$ by a larger linear function $(1 - 2\epsilon)^2 \cdot x$, in order to be able to come up with manageable estimates.

In fact, the difference between the two claims stems from the different ways in which we apply this "linearization" of the function $\phi(x, \epsilon)$ during averaging. The bounds they give are incomparable, though Theorem 1.12 is a more evident improvement of (4).

We note that the two functions $\phi(x, \epsilon)$ and $(1 - 2\epsilon)^2 \cdot x$ almost coincide for small values of x , and, loosely speaking, if the entropy of f is not too large, as is the case, say, for boolean functions, all the arguments of ϕ should lie very close to zero, meaning not much lost in the linear approximation. In this case, the bounds in Theorems 1.7 and 1.12 are very close to that in Corollary 1.9.

2) *Related work:* Y. Polyanskiy [15] has pointed out to us that the related question of upper bounding $I_{T_{\epsilon_A} f}(A, m)$ given $I_f(A, m)$ belongs to the area of *strong data processing inequalities* (SDPI) in information theory (see [16], [17] for pertinent results, and, in particular, for a new proof of Proposition 4.1).

Organization of the paper: This paper is organized as follows. The proof of Theorem 1.7 is given in Sections II to IV. Theorem 1.14 is proved in Section V. The remaining proofs are presented in Section VI.

II. A LINEAR PROGRAMMING BOUND FOR NOISY MUTUAL INFORMATION

In this section we upper bound the noisy mutual information $I_{T_{\epsilon_A} f}(A, m)$ by the value of a certain linear program.

Let f be a nonnegative function on the cube. Let A be a subset of $[n]$ and let $m \notin A$.

Let $|A| = k$. We will assume, without loss of generality, that $A = [k]$ and that $m = k + 1$.

Notation: From now on, we write λ for $(1 - 2\epsilon)^2$.

Discussion: Before going into details, let us give a high-level description of what the linear program attempts to capture. For ease of discussion the notation we use here is slightly different from that in the definition of the program below (they are the same up to scaling).

Given a random variable X on $\{0, 1\}^n$ distributed according to $f/\sum f$, consider a function I on the k -dimensional boolean cube, defined for $S \subseteq [k]$ by the mutual information $I(S) = I(\{X_i\}_{i \in S}; X_{k+1})$.

For $S \subseteq [k]$ and for $i \in S$, let $y_{S,i} = I(S) - I(S \setminus \{i\})$ be the "discrete derivative" of I at S in direction i . Note that $y_{S,i} \geq 0$, since this is the mutual information between X_i and X_{k+1} , given $\{X_j\}_{j \in S \setminus \{i\}}$. We view y as a function on the edges of the cube. Note also that, for any S , the value of the summation of y on the edges of any path from \emptyset to S is $I(S)$.

For $R \subseteq [k]$, applying noise in directions in R to f leads to a new distribution $T_{\epsilon_R} f / (\sum T_{\epsilon_R} f)$ on $\{0, 1\}^n$. This defines a new random variable X^R , a mutual information function I^R and discrete derivative functions $x_{S,i}^R = I^R(S) - I^R(S \setminus \{i\})$. (Note that $x^\emptyset = y$).

Observe that noise decreases mutual information, and hence $I^R \leq I$. However, the discrete derivatives x^R do not necessarily decrease. With that, and this is a key fact, by the *strong data processing inequality* [2], noise in direction i decreases the discrete derivative in direction i (i.e., the conditional mutual information between X_i^R and X_{k+1}^R) by a factor of at least λ .

The variables in the linear program below are the values of the discrete derivatives x^R , while we consider the discrete derivatives $y = x^\emptyset$ related to the initial function f to be the boundary data of the program. We note that the noisy mutual information $I([k]) = I_{T_{\epsilon_{[k]}} f}([k], k+1)$ is a linear combination of the variables, and that the strong data processing inequality provides linear local constraints on the variables.

Finally, we would like to explain the intuition behind the symmetrization procedure in Section III. The fact that for any R and S the value of the summation of x^R on the edges of any path from \emptyset to S is $I^R(S)$ provides a family of "symmetric" linear constraints on the variables. This makes it natural to look for a symmetric feasible solution to the linear program (symmetrizing the boundary data accordingly), one in which $x^R(S, i)$ depends only on $|S|$ and on $|R \cap S|$.

We were led to expect that this symmetric solution would be an optimal one by the following informal speculation. It turns out that the strong data processing inequality $x_{S,i}^R \leq \lambda \cdot (x_{S,i}^{R \setminus i})$ may be replaced by a stronger inequality $x_{S,i}^R \leq \phi(x_{S,i}^{R \setminus i})$ (see (3)).⁴ This turns the program into a *strictly concave* optimization problem, for which optimality of a feasible symmetric solution might be anticipated. It might also be hoped for that replacing the concave constraint by a linear one would preserve this property, and this is indeed turns out to be true.

More to the point, it turns out that for the symmetric solution we define, all the inequalities $x_{S,i}^R \leq \lambda \cdot (x_{S,i}^{R \setminus i})$ hold with equality.

The resulting argument is straightforward, most of the work going into setting up notation, and verifying feasibility of the symmetric solution. The key step, relying on symmetric properties of the discrete cube, is made in Lemma 3.6.

Linear program: Boundary data: For $S \subseteq [k]$ and for $i \in S$, we write

$$y_{S,i} = \text{Ent}(f \mid S \cup \{k+1\}) - \text{Ent}(f \mid S \setminus \{i\} \cup \{k+1\}) - \text{Ent}(f \mid S) + \text{Ent}(f \mid S \setminus \{i\})$$

The numbers $\{y_{S,i}\}$ are the boundary data for this problem.

We note that $y_{S,i} \geq 0$ for all S and i . In fact, the value of $y_{S,i}$ is proportional to a certain conditional mutual information. To see this, let X be distributed on $\{0, 1\}^n$ according to $f/\sum f$. Assume $\mathbb{E} f = 1$ and note that,

⁴This was shown in [19] if f is monotone (which suffices for applications) and in [17] for general functions.

by Subsection I-A1 and by (5), $y_{S,i}$ is given by

$$H(\{X_i\}_{i \in S \setminus \{i\} \cup \{k+1\}}) + H(\{X_i\}_{i \in S}) - H(\{X_i\}_{i \in S \cup \{k+1\}}) - H(\{X_i\}_{i \in S \setminus \{i\}}) = I(X_i; X_{k+1} | \{X_j\}_{j \in S \setminus \{i\}}).$$

Variables: $x_{S,i}^R$ for $R, S \subseteq [k]$ and $i \in S$.

The optimization problem: Given the boundary data, we want to upper bound μ , where

$$\mu = \text{Max} \sum_{i=1}^k x_{\{1, \dots, i\}}^{[k]}; i \quad (9)$$

under the following constraints.

Constraints:

1)

$$x_{S,i}^\emptyset = y_{S,i}$$

2)

$$x_{S,i}^R = x_{S,i}^{R \cap S}$$

3) For all $\sigma, \tau \in S_k$ holds

$$\sum_{i=1}^k x_{\{\sigma(1), \dots, \sigma(i)\}, \sigma(i)}^R = \sum_{i=1}^k x_{\{\tau(1), \dots, \tau(i)\}, \tau(i)}^R$$

4) If $i \in R$ then

$$x_{S,i}^R \leq \lambda \cdot (x_{S,i}^{R \setminus i})$$

We then have the following claim.

Theorem 2.1: The noisy mutual information $I_{T_{\epsilon_{[k]}} f}([k], k+1)$ is upperbounded by the value of the optimization problem (9).

Proof:

First, consider the boundary data. We claim that for any permutation $\sigma \in S_k$ holds

$$\sum_{i=1}^k y_{\{\sigma(1), \dots, \sigma(i)\}, \sigma(i)} = I_f([k], k+1) \quad (10)$$

In fact, it is easy to see that the LHS is a telescopic sum, summing to

$$\text{Ent}(f | [k+1]) - \text{Ent}(f | [k]) - \text{Ent}(f | \{k+1\}) = I_f([k], k+1)$$

Next we define a feasible solution for (9) whose value is $I_{T_{\epsilon_{[k]}} f}([k], k+1)$.

Fix $R \subseteq [k]$. Write f^R for $T_{\epsilon_R} f$. For $S \subseteq [k]$ and $i \in S$ set

$$x_{S,i}^R = \text{Ent}(f^R | S \cup \{k+1\}) - \text{Ent}(f^R | S \setminus \{i\} \cup \{k+1\}) - \text{Ent}(f^R | S) + \text{Ent}(f^R | S \setminus \{i\})$$

Clearly, $x_{S,i}^\emptyset = y_{S,i}$ and hence the first constraint of the program is satisfied.

As above, for any permutation $\sigma \in S_k$ holds

$$\sum_{i=1}^k x_{\{\sigma(1), \dots, \sigma(i)\}, \sigma(i)}^R = I_{T_{\epsilon_R} f}([k], k+1)$$

Hence, the third constraint is satisfied as well.

In particular,

$$\sum_{i=1}^k x_{\{1, \dots, i\}}^{[k], i} = I_{T_{\epsilon_{[k]}}} f([k], k+1)$$

so, the value given by this solution is indeed $I_{T_{\epsilon_{[k]}}} f([k], k+1)$.

We continue to prove its feasibility. We claim that for any $A \subseteq [k]$ holds $Ent(f^R | A) = Ent(f^{R \cap A} | A)$.

To see this, note that the noise operators commute with the conditional expectation operators, and hence

$$\mathbb{E}(T_{\epsilon_R} f | A) = T_{\epsilon_R} \mathbb{E}(f | A) = T_{\epsilon_{R \cap A}} T_{\epsilon_{R \setminus A}} \mathbb{E}(f | A) = T_{\epsilon_{R \cap A}} \mathbb{E}(f | A) = \mathbb{E}(T_{\epsilon_{R \cap A}} f | A)$$

Hence, by definition, $x_{S,i}^R = x_{S,i}^{R \cap S}$ for any $R, S \subseteq [k]$, and the second constraint holds.

To conclude the proof of the theorem, it remains to show that for any $R \subseteq S \subseteq [k]$ and $i \in R$ holds

$$x_{S,i}^R \leq \lambda \cdot (x_{S,i}^{R \setminus i}) \quad (11)$$

Recall that the strong data processing inequality [2] for a binary symmetric channel with crossover probability ϵ states that if V is a random variable with values in $\{0, 1\}$, and U is any random variable; and if $Y = V \oplus Z$, where Z is a Bernoulli random variable with parameter ϵ , statistically independent of U and V , then $I(U; Y) \leq \lambda \cdot I(U; V)$.

Let X be distributed on $\{0, 1\}^n$ according to $f^{R \setminus \{i\}} / \sum f^{R \setminus \{i\}}$. Assuming, as we may, $\mathbb{E} f = \mathbb{E} f^{R \setminus \{i\}} = 1$, we can rewrite (11) as

$$I(X_i \oplus Z; X_{k+1} | \{X_j\}_{j \in S \setminus \{i\}}) \leq \lambda \cdot I(X_i; X_{k+1} | \{X_j\}_{j \in S \setminus \{i\}})$$

which follows from applying the strong data processing inequality with $U = X_{k+1}$ and $V = X_i$, both conditioned on $\{X_j = x_j\}_{j \in S \setminus \{i\}}$, for all values of x_j .

III. THE OPTIMIZATION PROBLEM AND ITS SYMMETRIC VERSION

In this section we introduce a symmetric version of the optimization problem (9) and a specific symmetric feasible solution for the symmetric problem. We then argue that the value of this solution for the symmetric problem is at least as large as the optimal value for the original problem. Hence this value provides an upper bound on the noisy mutual information.

A. The symmetric problem and solution

Let $\{x_{S,i}^R\}$ be a feasible solution to the optimization problem (9) with boundary data $\{y_{S,i}\}$.

We define numbers y_1, \dots, y_k as follows. For $1 \leq s \leq k$ let

$$y_s = \mathbb{E}_{(S,i)} y_{S,i} \quad (12)$$

where the expectation is taken over all pairs (S, i) such that $|S| = s$ and $i \in S$.

For $0 \leq r < s \leq k$ we define x_s^r recursively in the following manner:

$$x_s^r = \begin{cases} y_s & \text{if } r = 0 \\ \lambda \cdot x_s^{r-1} + (1 - \lambda) \cdot x_{s-1}^{r-1} & \text{otherwise} \end{cases} \quad (13)$$

We now define the *symmetric version* of (9), by replacing the boundary data by a new, symmetric one. We set,¹² for all $i \in S \subseteq [k]$ with $|S| = s$:

$$\bar{y}_{S,i} = y_s$$

Next, we define the *symmetric solution* for the symmetric problem, in the following way. For $R \subseteq S$ with $|R| = r$, we set

$$\bar{x}_{S,i}^R = \begin{cases} \lambda \cdot x_s^{r-1} & \text{if } i \in R \\ x_s^r & \text{otherwise} \end{cases}$$

and for general R, S we set

$$\bar{x}_{S,i}^R = \bar{x}_{S,i}^{R \cap S}$$

Proposition 3.1.: The solution above is a feasible solution of the symmetric version of (9).

Moreover, for any $R \subseteq [k]$ of cardinality r and for any $\tau \in S_k$ holds

$$\sum_{i=1}^k \bar{x}_{\{\tau(1), \dots, \tau(i)\}, \tau(i)}^R = \sum_{j=1}^{k-r} y_j + \lambda \cdot \sum_{t=0}^{r-1} x_{k-r+t+1}^t \quad (14)$$

Proof:

The constraints 1 and 2 of (9) hold, by the definition of $\bar{x}_{S,i}^R$. We pass to constraint 4. Clearly, because of constraint 2, it suffices to prove it for $R \subseteq S$. In this case, taking $i \in R$, we have, by the definition of $\bar{x}_{S,i}^R$

$$\bar{x}_{S,i}^R = \lambda \cdot x_s^{r-1} = \lambda \cdot \bar{x}_{S,i}^{R \setminus \{i\}}$$

Next, we note that (14) will imply validity of constraint 3, since the RHS of (14) does not depend on τ .

It remains to prove (14). Let $i_1 < i_2 < \dots < i_r$ be such that $R = \{\tau(i_1), \tau(i_2), \dots, \tau(i_r)\}$. Then

$$\begin{aligned} \sum_{i=1}^k \bar{x}_{\{\tau(1), \dots, \tau(i)\}, \tau(i)}^R &= \sum_{j=1}^{i_1-1} + \sum_{j=i_1}^{i_2-1} + \dots + \sum_{j=i_r}^k = \\ &= \sum_{j=1}^{i_1-1} y_j + \left(\lambda \cdot y_{i_1} + \sum_{j=i_1+1}^{i_2-1} x_j^1 \right) + \left(\lambda \cdot x_{i_2}^1 + \sum_{j=i_2+1}^{i_3-1} x_j^2 \right) + \dots + \left(\lambda \cdot x_{i_r}^{r-1} + \sum_{j=i_r+1}^k x_j^r \right) \end{aligned}$$

Expanding $x_s^t = \lambda \cdot x_s^{t-1} + (1-\lambda) \cdot x_{s-1}^{t-1}$, we have the following exchange rule:

Two adjacent summands of the form $\lambda \cdot x_j^t + x_{j+1}^{t+1}$ can always be replaced by $x_j^t + \lambda \cdot x_{j+1}^t$. Applying this appropriate number of times in each bracket transforms the expression above into

$$\sum_{j=1}^{i_1-1} y_j + \left(\sum_{j=i_1}^{i_2-2} y_j + \lambda \cdot y_{i_2-1} \right) + \left(\sum_{j=i_2}^{i_3-2} x_j^1 + \lambda \cdot x_{i_3-1}^1 \right) + \dots + \left(\sum_{j=i_r}^{k-1} x_j^{r-1} + \lambda \cdot x_k^{r-1} \right)$$

Next we observe that the following rules apply in the original ordering of the summands: To the right of x_j^t is always either x_{j+1}^t or $\lambda \cdot x_{j+1}^t$. To the right of $\lambda \cdot x_s^r$ is always either x_{s+1}^{r+1} or $\lambda \cdot x_{s+1}^{r+1}$.

Moreover, this is easily verified to be preserved by the exchange rule above, by checking the four arising cases.

This means that applying the exchange rule as many times as needed, we can ensure all the summands multiplied by λ to be on the last r places on the right. Since the first summand is always either y_1 or $\lambda \cdot y_1$, these invariants guarantee that by doing so we obtain (14).

■

B. Optimality of the symmetric solution

Theorem 3.2.: Let $\{x_{S,i}^R\}$ be a feasible solution to the linear optimization problem (9). Let $\{\bar{x}_{S,i}^R\}$ be the symmetric solution for the symmetric version of this problem.

Then, for any $0 \leq r \leq k$ holds:

$$\mathbb{E}_{|R|=r} \sum_{i=1}^k x_{\{1,\dots,i\},i}^R \leq \mathbb{E}_{|R|=r} \sum_{i=1}^k \bar{x}_{\{1,\dots,i\},i}^R$$

Corollary 3.3.: The optimal value of (9) is upper bounded by the value of the symmetric solution to the symmetric version of the problem, which is given by

$$\lambda \cdot \sum_{t=0}^{k-1} x_{t+1}^t$$

Proof: Apply the theorem with $r = k$ and use (14). ■

Proof: (Of the theorem).

We proceed by double induction - on k and on $0 \leq r \leq k$. For $k = 1$ the claim is easily seen to be true.

Note also that the claim is true for any k and $r = 0$. This follows from constraints 1 and 3 of the linear program (9) and the definition of the symmetric boundary data. In fact, we have

$$\begin{aligned} \sum_{j=1}^k y_{\{1,\dots,j\},j} &= \mathbb{E}_{\sigma \in S_k} \sum_{j=1}^k y_{\{\sigma(1),\dots,\sigma(j)\},\sigma(j)} = \sum_{j=1}^k \mathbb{E}_{\sigma \in S_k} y_{\{\sigma(1),\dots,\sigma(j)\},\sigma(j)} = \\ & \sum_{j=1}^k \mathbb{E}_{|S|=j, i \in S} y_{S,i} = \sum_{j=1}^k y_j = \sum_{j=1}^k \bar{y}_{\{1,\dots,j\},j} \end{aligned}$$

Let now numbers r and k , with $0 < r \leq k$ be given. Assume the claim holds for $k - 1$, and also for k , for all $0 \leq t \leq r - 1$. We will argue it also holds for k and r .

We start with some simple properties of the linear program (9). We assume to be given the boundary data and a specific feasible solution to (9), and the symmetric solution to the symmetric version of (9), as in Theorem 3.2.

Lemma 3.4.: Let $M \subseteq [k]$. Let $\{y_{K,i}\}_{i \in K \subseteq M}$ be the restriction of the boundary data to subsets of M . For $R \subseteq M$, let $\{x_{K,i}^R\}_{i \in K \subseteq M}$ be the restriction of the feasible solution to subsets of M .

Then $\{x_{K,i}^R\}_{i \in K \subseteq M}$ is a feasible solution to the appropriate (smaller) optimization problem on M .

Proof:

Constraints 1, 2, and 4 are easy to check. As for constraint 3, let σ and τ be two permutations from M to itself. Extend them in the same way to permutations σ' and τ' on $[k]$. It is then easy to see that constraint 3 holds for σ and τ in the smaller problem, since it holds for σ' and τ' in the larger one.

■

Lemma 3.5.: Let $M \subseteq [k]$, with $|M| = m$ and let $R \subseteq [k]$. Let τ be a bijection from $[m]$ to M . Let

$$F(M, R, \tau) = \sum_{j=1}^m \bar{x}_{\{\tau(1),\dots,\tau(j)\},\tau(j)}^R$$

Then $F(M, R, \tau)$ depends only on m and $|R \cap M|$.

Proof:

Since the symmetric solution $\{\bar{x}_{S,i}^R\}$ satisfies constraint 2 of (9), we have

$$F(M, R, \tau) = \sum_{j=1}^m \bar{x}_{\{\tau(1), \dots, \tau(j)\}, \tau(j)}^R = \sum_{j=1}^m \bar{x}_{\{\tau(1), \dots, \tau(j)\}, \tau(j)}^{R \cap M} = F(M, R \cap M, \tau)$$

Let $r = |R \cap M|$.

Proceeding exactly as in the proof of Proposition 3.1, we get that

$$F(M, R \cap M, \tau) = \sum_{j=1}^{m-r} y_j + \lambda \cdot \sum_{t=0}^{r-1} x_{m-r+t+1}^t$$

That is, $F(M, R, \tau)$ depends only on m and $r = |R \cap M|$, as claimed. ■

Next, we introduce some notation.

1) *Notation:*

1) Let $M \subseteq [k]$. Let $\{y_{K,i}\}_{i \in K \subseteq M}$ be the restriction of the boundary data to the subsets of M .

We will denote by $\left\{ \mathcal{S}_M \left[x_{K,i}^R \right] \right\}$ the symmetric solution to the symmetric version of the smaller problem with this boundary data.

2) Let $L \subseteq [k]$, with $L = \{i_1, \dots, i_\ell\}$, so that $i_1 < i_2 < \dots < i_\ell$. Let $R \subseteq [k]$. Write

$$\mu^R(L) = \sum_{j=1}^{\ell} x_{\{i_1, \dots, i_j\}, i_j}^R$$

For $L \subseteq M \subseteq [k]$, and $R \subseteq M$, we denote

$$\mathcal{S}[\mu_M^R](L) = \sum_{j=1}^{\ell} \mathcal{S}_M \left[x_{\{i_1, \dots, i_j\}, i_j}^R \right]$$

Note that this quantity depends on M . With that, by Lemmas 3.4 and 3.5, given M , it depends only on the cardinalities $|L|$ and $|R \cap L|$.

3) Using the observation in the preceding paragraph, given $R \subseteq L \subseteq M \subseteq [k]$, with $|L| = \ell$, and $|R| = r$, we may also write $\mathcal{S}[\mu_M^r](\ell)$ for $\mathcal{S}[\mu_M^R](L)$.

In particular, note that the proof of Lemma 3.5 gives, in this notation

$$\mathcal{S}[\mu_{[k]}^r](m) = \sum_{j=1}^{m-r} y_j + \lambda \cdot \sum_{t=0}^{r-1} x_{m-r+t+1}^t \quad (15)$$

4) Finally, for $M \subseteq [k]$ and $0 \leq r \leq |M|$, we write

$$\mu_M^r = \mathbb{E}_{|R|=r, R \subseteq M} \mu^R(M) \quad \text{and} \quad \mathcal{S}[\mu_M^r] = \mathbb{E}_{|R|=r, R \subseteq M} \mathcal{S}[\mu_M^R](M)$$

We have completed introducing the new notation. In this notation the claim of the theorem amounts to:

$$\mu_{[k]}^r \leq \mathcal{S}[\mu_{[k]}^r] \quad (16)$$

We start with a lemma connecting the value of a solution of the optimization problem to these of smaller problems.

Lemma 3.6.:

$$\mu_{[k]}^r \leq \lambda \cdot \mu_{[k]}^{r-1} + (1 - \lambda) \cdot \mathbb{E}_{i \in [k]} \mu_{[k] \setminus \{i\}}^{r-1} \quad (17)$$

Proof:

Since the feasible solution $\{x_{S,i}^R\}$ satisfies constraints 2 and 3 of (9), for any $i \in R \subseteq [k]$ holds $\mu^R([k]) = \mu^{R \setminus \{i\}}([k] \setminus \{i\}) + x_{[k],i}^R$.

Similarly, $\mu^{R \setminus \{i\}}([k]) = \mu^{R \setminus \{i\}}([k] \setminus \{i\}) + x_{[k],i}^{R \setminus \{i\}}$.

Hence, by constraint 4,

$$x_{[k],i}^R \leq \lambda \cdot (x_{[k],i}^{R \setminus \{i\}}) = \lambda \cdot (\mu^{R \setminus \{i\}}([k]) - \mu^{R \setminus \{i\}}([k] \setminus \{i\}))$$

Averaging,

$$\begin{aligned} \mu_{[k]}^r &= \mathbb{E}_{R \subseteq [k], |R|=r} \mu_{[k]}^R = \mathbb{E}_{R, i \in R} (\mu^{R \setminus \{i\}}([k] \setminus \{i\}) + x_{[k],i}^R) \leq \\ &\mathbb{E}_{R, i \in R} \mu^{R \setminus \{i\}}([k] \setminus \{i\}) + \lambda \cdot \mathbb{E}_{R, i \in R} (\mu^{R \setminus \{i\}}([k]) - \mu^{R \setminus \{i\}}([k] \setminus \{i\})) = \\ &\lambda \cdot \mathbb{E}_{R, i \in R} \mu^{R \setminus \{i\}}([k]) + (1 - \lambda) \cdot \mathbb{E}_{R, i \in R} \mu^{R \setminus \{i\}}([k] \setminus \{i\}) \end{aligned}$$

It remains to note

$$\mathbb{E}_{R, i \in R} \mu^{R \setminus \{i\}}([k] \setminus \{i\}) = \mathbb{E}_{i \in [k]} \mathbb{E}_{|T|=r-1, T \subseteq [k] \setminus \{i\}} \mu^T([k] \setminus \{i\}) = \mathbb{E}_{i \in [k]} \mu_{[k] \setminus \{i\}}^{r-1}$$

and, similarly, $\mathbb{E}_{R, i \in R} \mu^{R \setminus \{i\}}([k]) = \mu_{[k]}^{r-1}$.

■

We now prove (16), starting from (17).

First, note that, by Lemma 3.4 and by the induction hypothesis for $k - 1$, we have $\mu_{[k] \setminus \{i\}}^{r-1} \leq \mathcal{S}[\mu]_{[k] \setminus \{i\}}^{r-1}$, for all $i \in [k]$.

Next, note that, by the induction hypothesis for k and $r - 1$, we have $\mu_{[k]}^{r-1} \leq \mathcal{S}[\mu]_{[k]}^{r-1}$.

This gives

$$\mu_{[k]}^r \leq \lambda \cdot \mathcal{S}[\mu]_{[k]}^{r-1} + (1 - \lambda) \cdot \mathbb{E}_{i \in [k]} \mathcal{S}[\mu]_{[k] \setminus \{i\}}^{r-1}$$

This implies that to prove (16) it suffices to show the following two identities:

1)

$$\mathbb{E}_{i \in [k]} \mathcal{S}[\mu]_{[k] \setminus \{i\}}^{r-1} = \mathcal{S}[\mu]_{[k]}^{r-1}(k - 1)$$

2)

$$\mathcal{S}[\mu]_{[k]}^r = \lambda \cdot \mathcal{S}[\mu]_{[k]}^{r-1} + (1 - \lambda) \cdot \mathcal{S}[\mu]_{[k]}^{r-1}(k - 1)$$

Lemma 3.7.:

$$\mathbb{E}_{i \in [k]} \mathcal{S}[\mu]_{[k] \setminus \{i\}}^{r-1} = \mathcal{S}[\mu]_{[k]}^{r-1}(k - 1)$$

Proof: We introduce the following notation. For $i = 1, \dots, k$ and for $0 \leq r < s \leq k - 1$, let

$$y_{s,i} = y_{s, [k] \setminus \{i\}} \quad \text{and} \quad x_{s,i}^r = x_{s, [k] \setminus \{i\}}^r$$

The values on the RHS of these identities are defined as in (12) and in (13) for the corresponding restricted problems.

We start with observing that $\mathbb{E}_{i \in [k]} y_{s,i} = y_s$. In fact, by definition,

$$\mathbb{E}_{i \in [k]} y_{s,i} = \mathbb{E}_{i \in [k]} \mathbb{E}_{|S|=s, S \subseteq [k] \setminus \{i\}, j \in S} y_{S,j} = \mathbb{E}_{|S|=s, j \in S} y_{S,j} = y_s$$

Next, we claim that for all $0 \leq r < s \leq k - 1$ holds $\mathbb{E}_{i \in [k]} x_{s,i}^r = x_s^r$.

This is easy to verify by induction on r . Note that we already know the claim holds for $r = 0$, and the induction step follows directly from the definitions and the induction hypothesis.

We now apply (14) to the restricted problems, to obtain that, for each $1 \leq i \leq k$ holds

$$\mathcal{S}[\mu]_{[k] \setminus \{i\}}^{r-1} = \sum_{j=1}^{k-r} y_{j,i} + \lambda \cdot \sum_{t=0}^{r-2} x_{k-r+t+1,i}^t$$

Hence, we have:

$$\mathbb{E}_{i \in [k]} \mathcal{S}[\mu]_{[k] \setminus \{i\}}^{r-1} = \sum_{j=1}^{k-r} \mathbb{E}_{i \in [k]} y_{j,i} + \lambda \cdot \sum_{t=0}^{r-2} \mathbb{E}_{i \in [k]} x_{k-r+t+1,i}^t = \sum_{j=1}^{k-r} y_j + \lambda \cdot \sum_{t=0}^{r-2} x_{k-r+t+1}^t$$

This, by (15), equals to $\mathcal{S}[\mu]_{[k]}^{r-1}(k-1)$, completing the proof of the lemma.

■

Lemma 3.8.:

$$\mathcal{S}[\mu]_{[k]}^r = \lambda \cdot \mathcal{S}[\mu]_{[k]}^{r-1} + (1 - \lambda) \cdot \mathcal{S}[\mu]_{[k]}^{r-1}(k-1)$$

Proof:

The proof of this lemma is similar to that of Lemma 3.6.

Since the symmetric solution $\mathcal{S}_{[k]}[x_{S,i}^R]$ (which is the same as $\{\bar{x}_{S,i}^R\}$) satisfies constraints 2 and 3 of (9), for any $i \in R \subseteq [k]$ holds

$$\mathcal{S}[\mu]_{[k]}^R([k]) = \mathcal{S}[\mu]_{[k]}^{R \setminus \{i\}}([k] \setminus \{i\}) + \mathcal{S}_{[k]}[x_{[k],i}^R]$$

Consider the notation we have introduced above. Using items 3 and 4 in the description of this notation, and recalling $\mathcal{S}_{[k]}[x_{[k],i}^R] = \lambda \cdot x_k^{r-1}$, we can rewrite this equality as

$$\mathcal{S}[\mu]_{[k]}^r = \mathcal{S}[\mu]_{[k]}^{r-1}(k-1) + \lambda \cdot x_k^{r-1}$$

On the other hand, we have, for $i \in R \subseteq [k]$:

$$\mathcal{S}[\mu]_{[k]}^{R \setminus \{i\}}([k]) = \mathcal{S}[\mu]_{[k]}^{R \setminus \{i\}}([k] \setminus \{i\}) + \mathcal{S}_{[k]}[x_{[k],i}^{R \setminus \{i\}}]$$

which is the same as

$$\mathcal{S}[\mu]_{[k]}^{r-1} = \mathcal{S}[\mu]_{[k]}^{r-1}(k-1) + x_k^{r-1}$$

Combining these two identities immediately implies the claim of the lemma.

■

This completes the proof of (16) and of the theorem.

■

C. The value of the symmetric optimization problem

Let $\{\bar{x}_{S,i}^R\}$ be the symmetric solution for the symmetric version of (9). By Corollary 3.3, its value depends linearly on the symmetric boundary data y_1, \dots, y_k , since $\{x_t^r\}$ are fixed linear functions of y_1, \dots, y_k . Let us denote this value by $V(y_1, \dots, y_k)$.

For $1 \leq s \leq k$, let e_s be the initial data vector with $y_s = 1$ and all the remaining y_t vanishing. Then $V(y_1, \dots, y_k) = \sum_{s=1}^k y_s \cdot V(e_s)$.

Next, we find the values of the parameters x_t^r for initial data given by a unit vector.

Lemma 3.9.: Let the initial data be given by the unit vector e_s , for some $1 \leq s \leq k$. Then the values of the parameters x_t^r , for $0 \leq r < t \leq k$, are as follows.

$$x_t^r = \begin{cases} \binom{r}{t-s} \cdot \lambda^{r-(t-s)} \cdot (1-\lambda)^{t-s} & \text{if } s \leq t \leq s+r \\ 0 & \text{otherwise} \end{cases}$$

(We use the convention $\binom{0}{0} = 1$.)

Proof: The claim of the lemma is easily verifiable by induction on r , or by directly verifying that (13) holds. ■

Corollary 3.10.:

$$V(e_s) = \lambda^s \cdot \sum_{m=0}^{k-s} \binom{s+m-1}{m} \cdot (1-\lambda)^m = 1 - \sum_{j=0}^{s-1} \binom{k}{j} \lambda^j (1-\lambda)^{k-j}$$

Proof: The first equality follows from Corollary 3.3. For the second equality, we proceed as follows

$$\begin{aligned} V(e_s) &= \frac{\lambda^s}{(s-1)!} \cdot \frac{\partial^{s-1}}{\partial x^{s-1}} \left[(1+x+\dots+x^{k-1}) \right]_{x=1-\lambda} = \\ &= \frac{\lambda^s}{(s-1)!} \cdot \left(\frac{\partial^{s-1}}{\partial x^{s-1}} \left[\frac{1}{1-x} \right]_{x=1-\lambda} - \frac{\partial^{s-1}}{\partial x^{s-1}} \left[\frac{x^k}{1-x} \right]_{x=1-\lambda} \right) = \\ &= 1 - \frac{\lambda^s}{(s-1)!} \cdot \frac{\partial^{s-1}}{\partial x^{s-1}} \left[\frac{x^k}{1-x} \right]_{x=1-\lambda} \end{aligned}$$

We have

$$\begin{aligned} \frac{\partial^t}{\partial x^t} \left[\frac{x^k}{1-x} \right] &= \sum_{i=0}^t \binom{t}{i} \frac{\partial^i}{\partial x^i} \left[\frac{1}{1-x} \right] \cdot \frac{\partial^{t-i}}{\partial x^{t-i}} [x^k] = \\ &= \sum_{i=0}^t \binom{t}{i} \cdot i! \cdot \frac{k!}{(k-t+i)!} \cdot x^{k-t+i} \cdot \frac{1}{(1-x)^{i+1}} \end{aligned}$$

Substituting $j = t - i$ and rearranging, this is

$$\frac{t!}{(1-x)^{t+1}} \cdot \sum_{j=0}^t \binom{k}{j} (1-x)^j \cdot x^{k-j}$$

Substituting $t = s - 1$, $x = 1 - \lambda$, and simplifying, we get

$$V(e_s) = 1 - \sum_{j=0}^{s-1} \binom{k}{j} \lambda^j (1-\lambda)^{k-j}$$

■

Corollary 3.11.:

$$V(y_1, \dots, y_k) = \sum_{s=1}^k \left(1 - \sum_{j=0}^{s-1} \binom{k}{j} \lambda^j (1-\lambda)^{k-j} \right) \cdot y_s$$

IV. PROOF OF THEOREM 1.7

We start with introducing some more notation.

1) *Notation:*

- For a subset S of $[n]$ of cardinality at most $n-2$, and for distinct $i, j \notin S$, we set

$$Z_{S;i,j} = \text{Ent}(f | S \cup \{i,j\}) - \text{Ent}(f | S \cup \{i\}) - \text{Ent}(f | S \cup \{j\}) + \text{Ent}(f | S)$$

- For $s = 1, \dots, n-1$, let $t_s = \mathbb{E}_{S;i,j} Z_{S;i,j}$.

Here the expectation is taken over all subsets S of $[n]$ of cardinality $s-1$, and, given S , over all distinct i, j not in S .

- Let A be a subset of $[n]$ of cardinality $k < n$ and let $m \notin A$. For $1 \leq s \leq k$, let

$$Y(A, m, s) = \mathbb{E}_{S,i} Z_{S;i,m}$$

where the expectation goes over subsets $S \subseteq A$ of cardinality $s-1$, and over $i \in A \setminus S$.

- For $1 \leq s \leq k \leq n$ let

$$\Lambda(k, s, \lambda) = 1 - \sum_{j=0}^{s-1} \binom{k}{j} \lambda^j (1-\lambda)^{k-j}$$

Proposition 4.1.: Let f be a nonnegative function on $\{0, 1\}^n$. Let A be a subset of $[n]$ of cardinality $k < n$ and let $m \notin A$.

Then

$$I_{T_{e_A} f}(A, m) \leq \sum_{s=1}^k \Lambda(k, s, \lambda) \cdot Y(A, m, s)$$

Proof:

By Theorem 2.1, the value of $I_{T_{e_A} f}(A, m)$ is bounded by the value of the linear optimization problem (9), with appropriate changes of indices.

By Theorem 3.2, this last value is upperbounded by the value of the symmetric version of the problem, which, according to Corollary 3.11, and tracing out the appropriate changes in indices and notation, is given by $\sum_{s=1}^k \Lambda(k, s, \lambda) \cdot Y(A, m, s)$.

■

Proof: (Of the theorem)

The proof relies on several lemmas. We start with a technical claim.

Lemma 4.2.: Let $1 \leq s \leq n-1$ be integer parameters. Let $0 < \lambda < 1$. Then

$$\sum_{k=s}^{n-1} \Lambda(k, s, \lambda) = \left(n - \frac{s}{\lambda} \right) + \frac{1}{\lambda} \cdot \sum_{j=0}^{s-1} \sum_{t=0}^j \binom{n}{t} \lambda^t (1-\lambda)^{n-t}$$

Proof:

$$\sum_{k=s}^{n-1} \Lambda(k, s, \lambda) = \sum_{k=s}^{n-1} \left(1 - \sum_{j=0}^{s-1} \binom{k}{j} \lambda^j (1-\lambda)^{k-j} \right) =$$

$$\binom{n-s}{k} - \sum_{k=s}^{n-1} \sum_{j=0}^{s-1} \binom{k}{j} \lambda^j (1-\lambda)^{k-j} = \binom{n-s}{k} - \sum_{j=0}^{s-1} \lambda^j \cdot \sum_{k=s}^{n-1} \binom{k}{j} (1-\lambda)^{k-j}$$

A simple calculation, similar to that in the proof of Corollary 3.10, gives

$$\lambda^j \cdot \sum_{k=s}^{n-1} \binom{k}{j} (1-\lambda)^{k-j} = \frac{1}{\lambda} \cdot \left(\sum_{t=0}^j \binom{s}{t} \lambda^t (1-\lambda)^{s-t} - \sum_{t=0}^j \binom{n}{t} \lambda^t (1-\lambda)^{n-t} \right)$$

The proof of the lemma is completed by summing the RHS over j , and observing

$$\sum_{j=0}^{s-1} \sum_{t=0}^j \binom{s}{t} \lambda^t (1-\lambda)^{s-t} = (1-\lambda) \cdot s$$

■

Lemma 4.3: Let f be a nonnegative function on $\{0, 1\}^n$ with expectation 1. Then

$$\text{Ent}(T_\epsilon f) \leq \sum_{i=1}^n \phi\left(\text{Ent}(f \mid \{i\})\right) + \sum_{s=1}^{n-1} w_s \cdot t_s$$

where

$$w_s = (\lambda n - s) + \sum_{j=0}^{s-1} \sum_{t=0}^j \binom{n}{t} \lambda^t (1-\lambda)^{n-t}$$

Lemma 4.4: Let f be a nonnegative function on $\{0, 1\}^n$. For any $0 \leq u \leq n-1$ holds

$$\mathbb{E}_{|B|=u+1} \text{Ent}(f \mid B) - (u+1) \cdot \mathbb{E}_{i \in [n]} \text{Ent}(f \mid \{i\}) = \sum_{s=1}^u (u-s+1) \cdot t_s$$

Next, we derive the theorem, assuming Lemmas 4.3 and 4.4 to hold.

Let T be a random subset of $[n]$ generated by sampling each element $i \in [n]$ independently with probability λ . We will show

$$\mathbb{E}_T \left(\text{Ent}(f \mid T) - \sum_{i \in T} \text{Ent}(f \mid \{i\}) \right) = \sum_{s=1}^{n-1} w_s \cdot t_s$$

Combining this with the claim of Lemma 4.3 will complete the proof.

For $0 \leq k \leq n$, let $p_k = \binom{n}{k} \lambda^k (1-\lambda)^{n-k}$. And, for $0 \leq u \leq n-1$, let

$$\mu_u = \mathbb{E}_{|B|=u+1} \left(\text{Ent}(f \mid B) - \sum_{i \in B} \text{Ent}(f \mid \{i\}) \right)$$

Then, using Lemma 4.4 and observing that $\mu_0 = 0$,

$$\begin{aligned} \mathbb{E}_T \left(\text{Ent}(f \mid T) - \sum_{i \in T} \text{Ent}(f \mid \{i\}) \right) &= \sum_{k=2}^n p_k \mu_{k-1} = \\ &= \sum_{k=2}^n p_k \sum_{s=1}^{k-1} (k-s) \cdot t_s = \sum_{s=1}^{n-1} \left(\sum_{k=s+1}^n (k-s) p_k \right) \cdot t_s \end{aligned}$$

We conclude by verifying the identity $w_s = \sum_{k=s+1}^n (k-s) p_k$, for $s = 1, \dots, n-1$.

In fact,

$$w_s = (\lambda n - s) + \sum_{j=0}^{s-1} \sum_{t=0}^j p_t = \sum_{k=0}^n (k-s) p_k + \sum_{t=0}^{s-1} (s-t) p_t = \sum_{k=s+1}^n (k-s) p_k$$

■

It remains to prove the lemmas.

Proof: (Of Lemma 4.3)

Recall that, by the chain rule for noisy entropy (7), for any permutation $\sigma \in S_n$ holds that $Ent(T_\epsilon f)$ is bounded from above by

$$\sum_{i=1}^n \phi \left(Ent(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}} f \mid \{\sigma(1), \dots, \sigma(i)\}) - Ent(T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}} f \mid \{\sigma(1), \dots, \sigma(i-1)\}) \right)$$

Using the notation introduced in Subsection I-D1, we can write this as

$$\sum_{i=1}^n \phi \left(Ent(f \mid \{\sigma(i)\}) + I_{T_{\epsilon_{\{\sigma(1), \dots, \sigma(i-1)\}}} f} \left(\{\sigma(1), \dots, \sigma(i-1)\}, \sigma(i) \right) \right)$$

Observe that the function ϕ is concave, and $\phi(0) = 0$. Hence $\phi(x+y) \leq \phi(x) + \phi(y)$ for any $0 \leq x, y \leq 1$. By this subadditivity of ϕ , the last expression is at most

$$\sum_{i=1}^n \phi \left(Ent(f \mid \{i\}) \right) + \sum_{k=2}^n \phi \left(I_{T_{\epsilon_{\{\sigma(1), \dots, \sigma(k-1)\}}} f} \left(\{\sigma(1), \dots, \sigma(k-1)\}, \sigma(k) \right) \right)$$

Averaging this expression over all $\sigma \in S_n$, we obtain

$$Ent(T_\epsilon f) \leq \sum_{i=1}^n \phi \left(Ent(f \mid \{i\}) \right) + \mu,$$

where

$$\mu = \mathbb{E}_\sigma \sum_{k=2}^n \phi \left(I_{T_{\epsilon_{\{\sigma(1), \dots, \sigma(k-1)\}}} f} \left(\{\sigma(1), \dots, \sigma(k-1)\}, \sigma(k) \right) \right)$$

Next, we upper bound μ . By transitivity of action of the symmetric group and by concavity of ϕ we have

$$\mu \leq \sum_{k=1}^{n-1} \phi(b_k) \quad \text{where} \quad b_k = \mathbb{E}_{A,m} T_{\epsilon_A f} (A, m)$$

where the expectation is over all $A \subseteq [n]$ of cardinality k and $m \notin A$.

Applying Proposition 4.1, we get

$$b_k \leq \mathbb{E}_{A,m} \sum_{s=1}^k \Lambda(k, s, \lambda) \cdot Y(A, m, s) = \sum_{s=1}^k \Lambda(k, s, \lambda) \cdot \mathbb{E}_{A,m} Y(A, m, s)$$

By the definition of $Y(A, m, s)$,

$$\mathbb{E}_{A,m} Y(A, m, s) = \mathbb{E}_{A,m} \mathbb{E}_{S,i} Z_{S;i,m} = \mathbb{E}_{S,i,m} Z_{S;i,m} \cdot \mathbb{E}_A 1 = \mathbb{E}_{S,i,m} Z_{S;i,m}$$

where in the second expression the first expectation is over k -subsets A of $[n]$ and $m \notin A$, and the second expectation is over $(s-1)$ -subsets S of A and over $i \in A \setminus S$. Rearranging, we get the third expression in

which the first expectation is over all subsets S of $[n]$ of cardinality $s - 1$ and over all distinct $i, m \notin S$, and the second expectation is over all supersets A of S of cardinality k with $i \in A$ and $m \notin A$.

Recalling the definition of t_s above, we deduce $b_k = \sum_{s=1}^k \Lambda(k, s, \lambda) \cdot t_s$.

Using the inequality $\phi(x) \leq \lambda x$, and Lemma 4.2, we have

$$\begin{aligned} \mu &\leq \lambda \cdot \sum_{k=1}^{n-1} b_k = \lambda \cdot \sum_{k=1}^{n-1} \sum_{s=1}^k \Lambda(k, s, \lambda) \cdot t_s = \\ &\sum_{s=1}^{n-1} t_s \cdot \left(\lambda \cdot \sum_{k=s}^{n-1} \Lambda(k, s, \lambda) \right) = \sum_{s=1}^{n-1} w_s \cdot t_s \end{aligned}$$

■

Proof: (of Lemma 4.4)

By (10), for any subset A of $[n]$ of cardinality $1 \leq k \leq n - 1$, for any $m \notin A$, and for any bijection $\tau : [k] \rightarrow A$ holds, in the notation of this section,

$$\sum_{s=1}^k Z_{\{\tau(1), \dots, \tau(s-1)\}; \tau(s), m} = I_f(A, m)$$

We now average over all the variables, setting

$$c_k = \mathbb{E}_{A, m, \tau} \sum_{s=1}^k Z_{\{\tau(1), \dots, \tau(s-1)\}; \tau(s), m}$$

On one hand, we have

$$\begin{aligned} c_k &= \mathbb{E}_{A, m} I_f(A, m) = \mathbb{E}_{A, m} \left(Ent(f | A \cup \{m\}) - Ent(f | A) - Ent(f | \{m\}) \right) = \\ &\mathbb{E}_{|B|=k+1} Ent(f | B) - \mathbb{E}_{|A|=k} Ent(f | A) - \mathbb{E}_{i \in [n]} Ent(f | \{i\}) \end{aligned}$$

On the other hand, similarly to the computation in the preceding lemma, we have

$$c_k = \sum_{s=1}^k \mathbb{E}_{A, m, \tau} Z_{\{\tau(1), \dots, \tau(s-1)\}; \tau(s), m} = \sum_{s=1}^k \mathbb{E}_{A, S, i, m} Z_{S; i, m} = \sum_{s=1}^k t_s$$

where the expectation in the third expression is over k -subsets A of $[n]$, over $(s - 1)$ -subsets S of A , over $m \notin A$ and $i \in A \setminus S$.

Hence, for any $1 \leq u \leq n - 1$ holds

$$\mathbb{E}_{|B|=u+1} Ent(f | B) - (u + 1) \cdot \mathbb{E}_{i \in [n]} Ent(f | \{i\}) = \sum_{k=1}^u c_k = \sum_{s=1}^u (u - s + 1) \cdot t_s$$

completing the proof of the lemma and of the theorem.

■

Let δ be the constant in the theorem. We will assume in the following argument that δ is sufficiently small.

Let $0 < \epsilon < 1/2$ be a noise parameter, such that $(1 - 2\epsilon)^2 \leq \delta$. Let $\lambda = (1 - 2\epsilon)^2$.

Let $f : \{0, 1\}^n \rightarrow \{0, 1\}$ be a boolean function, satisfying the constraints of the theorem. Let $1 \leq k \leq n$ be the coordinate such that $|\widehat{f}(k)|$ is large. W.l.o.g. assume that $k = 1$ and that $\widehat{f}(1)$ is positive.

We introduce some additional notation.

Notation:

- Let $0 \leq \alpha \leq \delta$ be such that $\widehat{f}(1) = (1 - \alpha) \cdot \mathbb{E} f$.
- Let $0 \leq \beta \leq \delta$ be such that $\mathbb{E} f = 1/2 - \beta$. Let $\gamma = \alpha + \beta$.
- If $\alpha \leq \lambda$, we define $\tau = \left(\frac{1-\lambda}{1-\alpha}\right)^2$, and define auxiliary noise ϵ_τ , such that $(1 - 2\epsilon_\tau)^2 = \tau$. If $\alpha > \lambda$, we set $\tau = 1$ and $\epsilon_\tau = 0$.
- Let ϵ_1 be such that $T_\epsilon = T_{\epsilon_1} T_{\epsilon_\tau}$. Let $\lambda_1 = (1 - 2\epsilon_1)^2$. Note that $\lambda = \tau \cdot \lambda_1$.
- Let $h = T_{\epsilon_\tau} f$. Note that $T_\epsilon f = T_{\epsilon_1} h$, and hence $\text{Ent}(T_\epsilon f) = \text{Ent}(T_{\epsilon_1} h)$.

A. Proof of the first claim of the theorem

We start with applying Theorem 1.7 to the function h with noise ϵ_1 . The theorem is stated for functions with expectation 1. We modify it, using the linearity of entropy, to obtain

$$\text{Ent}(T_{\epsilon_1} h) \leq \mathbb{E}_T \left(\text{Ent}(h | T) - \sum_{i \in T} \text{Ent}(h | \{i\}) \right) + \mathbb{E} h \cdot \sum_{i=1}^n \phi \left(\text{Ent} \left(\frac{h}{\mathbb{E} h} \mid \{i\} \right), \epsilon_1 \right)$$

Here T is a random subset of $[n]$ generated by sampling each element $i \in [n]$ independently with probability $(1 - 2\epsilon_1)^2$.

Since there are more than one noise parameters involved, we now write the function ϕ with the noise parameter stated explicitly.

Next, note that by (3), for any $1 \leq i \leq n$ holds

$$\mathbb{E} h \cdot \phi \left(\text{Ent} \left(\frac{h}{\mathbb{E} h} \mid \{i\} \right), \epsilon_1 \right) \leq \lambda_1 \cdot \text{Ent}(h | \{i\})$$

Hence the previous inequality implies

$$\begin{aligned} \text{Ent}(T_{\epsilon_1} h) &\leq \lambda_1 \cdot \mathbb{E}_{T, 1 \in T} \left(\text{Ent}(h | T) - \text{Ent}(h | \{1\}) \right) + \\ &(1 - \lambda_1) \cdot \mathbb{E}_{T, 1 \notin T} \text{Ent}(h | T) + \mathbb{E} h \cdot \phi \left(\text{Ent} \left(\frac{h}{\mathbb{E} h} \mid \{1\} \right), \epsilon_1 \right) \end{aligned} \quad (18)$$

The proof will be based on three lemmas, which upperbound each of the three summands on the RHS of (18).

Lemma 5.1.:

$$\mathbb{E}_{T, 1 \in T} \left(\text{Ent}(h | T) - \text{Ent}(h | \{1\}) \right) \leq O \left(\lambda_1 \cdot \gamma + \gamma^2 \ln \left(\frac{1}{\gamma} \right) \right)$$

Lemma 5.2.:

$$\mathbb{E}_{T, 1 \notin T} \text{Ent}(h | T) \leq O \left(\lambda_1^2 \cdot \gamma + \lambda_1 \cdot \gamma^2 \ln \left(\frac{1}{\gamma} \right) \right)$$

Lemma 5.3.:

$$\mathbb{E} h \cdot \phi \left(Ent \left(\frac{h}{\mathbb{E} h} \mid \{1\} \right), \epsilon_1 \right) \leq \frac{1}{2} \cdot (1 - H(\epsilon)) - \Omega(\lambda \cdot \gamma)$$

The asymptotic notation in each of the lemmas hides absolute constants.

Given the lemmas, the first claim of the theorem is easy to verify. Indeed, recall that λ_1 is a constant multiple of λ . Hence, the lemmas and (18) imply that

$$Ent(T_{\epsilon} f) = Ent(T_{\epsilon_1} h) \leq \frac{1}{2} \cdot (1 - H(\epsilon)) - \Omega(\lambda \cdot \gamma) + o_{\lambda, \gamma \rightarrow 0}(\lambda \cdot \gamma)$$

Therefore, for a sufficiently small $\delta > 0$, bearing in mind that $0 \leq \alpha, \beta, \lambda \leq \delta$, the claim holds.

It remains to prove the lemmas. For that purpose we will need the following version of the logarithmic Sobolev inequality for the boolean cube.

Lemma 5.4.: Let g be a nonnegative function on $\{0, 1\}^n$. Let $\mathcal{E}(g, g)$ be the *Dirichlet form*, given by $\mathcal{E}(g, g) = \mathbb{E}_{x \in \{0, 1\}^n} \mathbb{E}_{y \sim x} (g(y) - g(x))^2$. Then

$$\mathcal{E}(g, g) \geq 2 \ln 2 \cdot \mathbb{E} g \cdot Ent(g)$$

Proof:

We start with a simple auxiliary claim.

Let $x_1 \geq x_2 \geq \dots \geq x_N$ be nonnegative numbers summing to 1. Then the numbers $y_k = \frac{x_k^2}{\sum_{i=1}^N x_i^2}$, for $k = 1, \dots, N$, majorize $\{x_k\}$. That is,

$$y_1 \geq x_1, \quad y_1 + y_2 \geq x_1 + x_2, \quad \dots, \quad y_1 + \dots + y_N = 1 = x_1 + \dots + x_N$$

To see this, fix some $1 \leq t \leq N$. We have to show $\sum_{k=1}^t x_k^2 \geq \left(\sum_{k=1}^t x_k \right) \cdot \left(\sum_{k=1}^N x_k^2 \right)$.

We may and will assume that all of the x_k are strictly positive. After some rearrangement, the claim reduces to showing

$$\frac{\sum_{k=1}^t x_k^2}{\sum_{k=1}^t x_k} \geq \frac{\sum_{m=t+1}^N x_m^2}{\sum_{k=t+1}^N x_m}$$

This holds because the LHS is lowerbounded by x_t , and the RHS is upperbounded by x_{t+1} .

A simple corollary of this claim is that for any nonnegative not identically zero function g on a finite domain endowed with uniform measure, holds that $g^2 / \mathbb{E} g^2$ majorizes $g / \mathbb{E} g$.

This is well-known to imply (see [10]) that $g / \mathbb{E} g$ is a convex combination of permuted versions of $g^2 / \mathbb{E} g^2$. Since the entropy functional is linear and convex, this implies

$$Ent(g^2) \geq \frac{\mathbb{E} g^2}{\mathbb{E} g} \cdot Ent(g) \geq \mathbb{E} g \cdot Ent(g)$$

The claim of the lemma follows from this inequality combined with the logarithmic Sobolev inequality [8]:

$$\mathcal{E}(g, g) \geq 2 \ln 2 \cdot Ent(g^2)$$

■

In the following argument we are going to use the Walsh-Fourier expansion for functions on the boolean cube, writing a function g as $\sum_{S \subseteq [n]} \hat{g}(S) \cdot W_S$, where $\{W_S\}_{S \subseteq [n]}$ is the Walsh-Fourier basis.

In particular, for the Dirichlet form, we have $\mathcal{E}(g, g) = 4 \cdot \sum_{S \subseteq [n]} |S| \widehat{g}^2(S)$. Hence the preceding lemma²⁴ implies

$$Ent(g) \leq \frac{2}{\ln 2} \cdot \frac{1}{\mathbb{E}g} \cdot \sum_{S \subseteq [n]} |S| \widehat{g}^2(S) \quad (19)$$

We will also need the following precise version of an inequality of [4], due to [6]:

Theorem 5.5: There exists a universal constant $L > 0$ with the following property. For $g : \{0, 1\}^n \rightarrow \{-1, 1\}$, let $\rho = \left(\sum_{A \subseteq [n]: |A| \geq 2} \widehat{g}^2(A) \right)^{1/2}$. Then there exists some $B \subseteq [n]$ with $|B| \leq 1$ such that

$$\sum_{A \subseteq [n]: |A| \leq 1, A \neq B} \widehat{g}^2(A) \leq L \cdot \rho^4 \ln \left(\frac{2}{\rho} \right)$$

and $|\widehat{g}(B)|^2 \geq 1 - \rho^2 - L \cdot \rho^4 \ln \left(\frac{2}{\rho} \right)$.

Consider the function f and recall that it satisfies the assumptions of Theorem 1.14.

Let $g = 2f - 1$. Then $g : \{0, 1\}^n \rightarrow \{-1, 1\}$. Note that $\widehat{g}(0) = 2\widehat{f}(0) - 1$, and that $\widehat{g}(S) = 2\widehat{f}(S)$, for $|S| > 0$. In particular, $\widehat{g}(0) = 2\mathbb{E}f - 1 = -2\beta$, and $\widehat{g}(\{1\}) = 2(1 - \alpha)\mathbb{E}f = (1 - \alpha)(1 - 2\beta)$.

Recall that $0 \leq \alpha, \beta \leq \delta$, and that $\gamma = \alpha + \beta$. Hence, assuming δ is sufficiently small, we have

$$\sum_{|A| \geq 2} \widehat{f}^2(A) \leq \sum_{|A| \geq 2} \widehat{g}^2(A) \leq 1 - \widehat{g}^2(\{1\}) \leq L \cdot \gamma, \quad (20)$$

for some absolute constant L .

Applying Theorem 5.5 to the function g , we get, for a sufficiently large constant L_1 ,

$$\sum_{k=2}^n \widehat{f}^2(\{k\}) \leq \sum_{k=2}^n \widehat{g}^2(\{k\}) \leq L_1 \cdot \gamma^2 \ln \left(\frac{1}{\gamma} \right) \quad (21)$$

Proof of Lemma 5.2: Fix $T \subseteq [n]$. Let $g_T = \mathbb{E}(h \mid T)$.

Note that $g_T = \sum_{S \subseteq T} \widehat{h}(S) \cdot W_S$, and hence, by (19), we have

$$Ent(g_T) \leq \frac{2}{\ln 2} \cdot \frac{1}{\mathbb{E}g_T} \cdot \sum_{S \subseteq T} |S| \widehat{h}^2(S) = \frac{2}{\ln 2} \cdot \frac{1}{\mathbb{E}h} \cdot \sum_{S \subseteq T} |S| \widehat{h}^2(S)$$

Hence,

$$\mathbb{E}_{T, 1 \notin T} Ent(g_T) \leq \frac{2}{\ln 2} \cdot \frac{1}{\mathbb{E}h} \cdot \mathbb{E}_{T, 1 \notin T} \sum_{S \subseteq T} |S| \widehat{h}^2(S) = \frac{2}{\ln 2} \cdot \frac{1}{\mathbb{E}h} \cdot \sum_{S, 1 \notin S} |S| \lambda_1^{|S|} \widehat{h}^2(S)$$

Recall that $h = T_{\epsilon_r} f$. This means (see, e.g., [11]) that for any $S \subseteq [n]$, holds $\widehat{h}(S) = \tau^{|S|/2} \cdot \widehat{f}(S)$.

In particular, $|\widehat{h}(S)| \leq |\widehat{f}(S)|$. Applying (20) and (21), we have that, for a sufficiently large absolute constant L , the last expression is bounded by

$$L \cdot \left(\lambda_1^2 \cdot \gamma + \lambda_1 \cdot \gamma^2 \ln \left(\frac{1}{\gamma} \right) \right)$$

This concludes the proof of the lemma. ■

Proof of Lemma 5.3: Let $g = \mathbb{E}\left(\frac{f}{\mathbb{E}f} \mid \{1\}\right)$. Then g is a function on a 2-point space $\{0, 1\}$, with $g(0) = 2^{-\alpha}$ and $g(1) = \alpha$.

Observe that the noise operator commutes with the projection operator. Hence, since $h = T_{\epsilon_\tau} f$, we have $g_1 := \mathbb{E}\left(\frac{h}{\mathbb{E}h} \mid \{1\}\right) = T_{\epsilon_\tau} g$.

Observe also that, by the definition of Mrs. Gerber's function ϕ , we have

$$\phi\left(\text{Ent}\left(\frac{h}{\mathbb{E}h} \mid \{1\}\right), \epsilon_1\right) = \text{Ent}(T_{\epsilon_1} g_1) = \text{Ent}(T_{\epsilon_1} T_{\epsilon_\tau} g) = \text{Ent}(T_\epsilon g)$$

The last equality follows from the definition of ϵ_1 and ϵ_τ .

It is easy to verify that $T_\epsilon g(0) = 1 + (1 - \alpha) \cdot \lambda^{1/2}$ and that $T_\epsilon g(1) = 1 - (1 - \alpha) \cdot \lambda^{1/2}$.

Hence, $\text{Ent}(T_\epsilon g) = 1 - H_2\left(\frac{1 - (1 - \alpha) \cdot \lambda^{1/2}}{2}\right)$.

Recall that

$$H_2\left(\frac{1-x}{2}\right) = 1 - \frac{1}{\ln 2} \cdot \sum_{k=1}^{\infty} \frac{1}{2k(2k-1)} \cdot x^{2k}$$

with the series converging absolutely for $-1 \leq x \leq 1$.

Let $F(x) = 1 - H_2\left(\frac{1-\sqrt{x}}{2}\right)$, for $0 \leq x \leq 1$. Then $F(x) = \frac{1}{\ln 2} \cdot \sum_{k=1}^{\infty} \frac{1}{2k(2k-1)} \cdot x^k$.

This is a convex function on $[0, 1]$, and hence for any $0 \leq x < y \leq 1$ holds $F(y) - F(x) \geq (y - x) \cdot F'(x)$. The derivative F' is given by $F'(x) = \frac{1}{2 \ln 2} \cdot \sum_{k=1}^{\infty} \frac{1}{2k-1} \cdot x^{k-1}$, with the series converging for $0 \leq x < 1$.

Hence $F' \geq \frac{1}{2 \ln 2}$ on $(0, 1)$, and $F(y) - F(x) \geq \frac{1}{2 \ln 2} \cdot (y - x)$. Applying this with $y = \lambda$ and $x = (1 - \alpha)^2 \cdot \lambda$, we get

$$\left(1 - H_2(\epsilon)\right) - \text{Ent}(T_\epsilon g) = F(\lambda) - F\left((1 - \alpha)^2 \cdot \lambda\right) \geq c_1 \cdot \lambda \cdot \alpha$$

where $c_1 > 0$ is an absolute constant.

In other words,

$$\phi\left(\text{Ent}\left(\frac{h}{\mathbb{E}h} \mid \{1\}\right), \epsilon_1\right) = \text{Ent}(T_\epsilon g) \leq \left(1 - H_2(\epsilon)\right) - c_1 \cdot \lambda \cdot \alpha$$

To conclude the proof of the lemma, note that, for a sufficiently small λ , we have $\text{Ent}(T_\epsilon g) \geq c_2 \cdot \lambda$, for an absolute constant c_2 , and hence

$$\begin{aligned} \mathbb{E}h \cdot \phi\left(\text{Ent}\left(\frac{h}{\mathbb{E}h} \mid \{1\}\right), \epsilon_1\right) &= \left(\frac{1}{2} - \beta\right) \cdot \text{Ent}(T_\epsilon g) \leq \\ \frac{1}{2} \cdot \left(1 - H_2(\epsilon)\right) - c \cdot \lambda \cdot (\alpha + \beta) &= \frac{1}{2} \cdot \left(1 - H_2(\epsilon)\right) - c \cdot \lambda \cdot \gamma \end{aligned}$$

for an absolute constant c . For the inequality, note that $1 - H_2(\epsilon) = F(\lambda) \geq \frac{1}{2 \ln 2} \cdot \lambda$.

This completes the proof of the lemma. ■

The proof of Lemma 5.1 is somewhat harder. We present it in the next subsection.

1) *Proof of Lemma 5.1:* We proceed similarly to the proof of Lemma 5.2, and use the notation introduced²⁶ in that proof.

Given a function g on the boolean cube, we write $\mathbb{E}(g \mid x_1 = 0, x_2, \dots, x_k)$ for the restriction of $\mathbb{E}(g \mid x_1, x_2, \dots, x_k)$ on the subcube $x_1 = 0$, and similarly for $\mathbb{E}(g \mid x_1 = 1, x_2, \dots, x_k)$.

We note that for $g = \sum_{S \subseteq [n]} \widehat{g}(S) \cdot W_S$, we have

$$\mathbb{E}(g \mid x_1 = 0, x_2, \dots, x_n) = \sum_{R \subseteq [n], 1 \notin R} (\widehat{g}(R) + \widehat{g}(R \cup \{1\})) \cdot W_R$$

and

$$\mathbb{E}(g \mid x_1 = 1, x_2, \dots, x_n) = \sum_{R \subseteq [n], 1 \notin R} (\widehat{g}(R) - \widehat{g}(R \cup \{1\})) \cdot W_R$$

We will also use the following easily verifiable identity, holding for nonnegative functions g :

$$\text{Ent}(g) - \text{Ent}(g \mid \{1\}) = \frac{1}{2} \cdot \text{Ent}(g \mid x_1 = 0, x_2, \dots, x_n) + \frac{1}{2} \cdot \text{Ent}(g \mid x_1 = 1, x_2, \dots, x_n)$$

As before, let $g_T = \mathbb{E}(h \mid T)$, for a subset $T \subseteq [n]$. Note that if $1 \in T$, then $\mathbb{E}(g_T \mid \{1\}) = \mathbb{E}(h \mid \{1\})$.

Hence

$$\begin{aligned} \mathbb{E}_{T, 1 \in T} \left(\text{Ent}(h \mid T) - \text{Ent}(h \mid \{1\}) \right) &= \mathbb{E}_{T, 1 \in T} \left(\text{Ent}(g_T) - \text{Ent}(g_T \mid \{1\}) \right) = \\ &= \frac{1}{2} \cdot \mathbb{E}_{T, 1 \in T} \text{Ent}(g_T \mid x_1 = 0, x_2, \dots, x_n) + \frac{1}{2} \cdot \mathbb{E}_{T, 1 \in T} \text{Ent}(g_T \mid x_1 = 1, x_2, \dots, x_n) \end{aligned}$$

We will prove the lemma by showing that, for a sufficiently large absolute constant L , hold both

$$\mathbb{E}_{T, 1 \in T} \text{Ent}(g_T \mid x_1 = 0, x_2, \dots, x_n) \leq L \cdot \lambda_1 \cdot \gamma \quad (22)$$

and

$$\mathbb{E}_{T, 1 \in T} \text{Ent}(g_T \mid x_1 = 1, x_2, \dots, x_n) \leq L \cdot \left(\lambda_1 \cdot \gamma + \gamma^2 \ln \left(\frac{1}{\gamma} \right) \right) \quad (23)$$

Proof of (22): Fix a subset $T \subseteq [n]$, with $1 \in T$. Recall that $g_T = \sum_{S \subseteq T} \widehat{h}(S) \cdot W_S$, and hence

$$\mathbb{E}(g_T \mid x_1 = 0, x_2, \dots, x_n) = \sum_{R \subseteq T \setminus \{1\}} (\widehat{h}(R) + \widehat{h}(R \cup \{1\})) \cdot W_R$$

In particular,

$$\mathbb{E}(g_T \mid x_1 = 0) = \widehat{h}(0) + \widehat{h}(\{1\}) = \widehat{f}(0) + \tau^{1/2} \cdot \widehat{f}(\{1\}) \geq \mathbb{E} f$$

Applying (19), we have, for a sufficiently large constant L_1 ,

$$\begin{aligned} \text{Ent}(g_T \mid x_1 = 0, x_2, \dots, x_n) &\leq \frac{2}{\ln 2} \cdot \frac{1}{\mathbb{E} f} \cdot \sum_{R \subseteq T \setminus \{1\}} |R| \cdot (\widehat{h}(R) + \widehat{h}(R \cup \{1\}))^2 \leq \\ &L_1 \cdot \sum_{R \subseteq T \setminus \{1\}} |R| \cdot (\widehat{h}^2(R) + \widehat{h}^2(R \cup \{1\})) \end{aligned}$$

Averaging over T , we have

$$\mathbb{E}_{T,1 \in T} \text{Ent}(g_T \mid x_1 = 0, x_2, \dots, x_n) \leq L_1 \cdot \left(\sum_{R,1 \notin R} |R| \lambda_1^{|R|} \widehat{h}^2(R) + \sum_{R,1 \notin R} |R| \lambda_1^{|R|} \widehat{h}^2(R \cup \{1\}) \right)$$

Using the fact that $|\widehat{h}(S)| \leq |\widehat{f}(S)|$ for all $S \subseteq [n]$, and applying (20) and (21), we have, for a sufficiently large constant L_2 ,

$$\sum_{R,1 \notin R} |R| \lambda_1^{|R|} \widehat{h}^2(R) \leq L_2 \cdot \left(\lambda_1 \cdot \gamma^2 \ln \left(\frac{1}{\gamma} \right) + \lambda_1^2 \cdot \gamma \right)$$

and

$$\sum_{R,1 \notin R} |R| \lambda_1^{|R|} \widehat{h}^2(R \cup \{1\}) \leq L_2 \cdot \lambda_1 \cdot \gamma$$

Summing up, this gives (22).

Proof of (23): Similarly to the above,

$$\mathbb{E} \left(g_T \mid x_1 = 1, x_2, \dots, x_n \right) = \sum_{R \subseteq T \setminus \{1\}} \left(\widehat{h}(R) - \widehat{h}(R \cup \{1\}) \right) \cdot W_R$$

Which means that

$$\mathbb{E} \left(g_T \mid x_1 = 1 \right) = \widehat{h}(0) - \widehat{h}(\{1\}) = \widehat{f}(0) - \tau^{1/2} \cdot \widehat{f}(\{1\}) = \mathbb{E} f \cdot \left(1 - \tau^{1/2} \cdot (1 - \alpha) \right)$$

Recall that $\tau^{1/2} = 1$ if $\alpha \geq \lambda$ and $\tau^{1/2} = \frac{1-\lambda}{1-\alpha}$ otherwise. In both cases, note that we have $\mathbb{E} \left(g_T \mid x_1 = 1 \right) \geq \lambda \cdot \mathbb{E} f$.

Applying (19), and averaging over T , we have, for a sufficiently large constant L_1 ,

$$\mathbb{E}_{T,1 \in T} \text{Ent}(g_T \mid x_1 = 1, x_2, \dots, x_n) \leq L_1 \cdot \frac{1}{\lambda} \cdot \sum_{R,1 \notin R} |R| \lambda_1^{|R|} \cdot \left(\widehat{h}(R) - \widehat{h}(R \cup \{1\}) \right)^2$$

Let $g = \mathbb{E} \left(h \mid x_1 = 1, x_2, \dots, x_n \right)$. Then $g = \sum_{R \subseteq [n], 1 \notin R} \left(\widehat{h}(R) - \widehat{h}(R \cup \{1\}) \right) \cdot W_R$. Hence

$$\mathbb{E}_{T,1 \in T} \text{Ent}(g_T \mid x_1 = 1, x_2, \dots, x_n) \leq L_1 \cdot \frac{1}{\lambda} \cdot \sum_{R,1 \notin R} |R| \lambda_1^{|R|} \cdot \widehat{g}^2(R) \quad (24)$$

Consider the function g . Since $h = T_{\epsilon_\tau} f$, we have

$$g = \epsilon_\tau \cdot T_{\epsilon_\tau} \left(\mathbb{E} \left(f \mid x_1 = 0, x_2, \dots, x_n \right) \right) + \left(1 - \epsilon_\tau \right) \cdot T_{\epsilon_\tau} \left(\mathbb{E} \left(f \mid x_1 = 1, x_2, \dots, x_n \right) \right)$$

For $i = 0, 1$, let $f_i = \mathbb{E} \left(f \mid x_1 = i, x_2, \dots, x_n \right)$, and let $t_i = T_{\epsilon_\tau} f_i$. Note that for $i = 0, 1$ and for any R , $1 \notin R$, holds $|\widehat{t}_i(R)| \leq |\widehat{f}_i(R)|$.

Therefore, since $g = \epsilon_\tau \cdot t_0 + (1 - \epsilon_\tau) \cdot t_1$, we have, for any R , $1 \notin R$ that

$$\widehat{g}^2(R) \leq \epsilon_\tau \cdot \widehat{t}_0^2(R) + (1 - \epsilon_\tau) \cdot \widehat{t}_1^2(R) \leq \epsilon_\tau \cdot \widehat{f}_0^2(R) + (1 - \epsilon_\tau) \cdot \widehat{f}_1^2(R)$$

Hence,

$$\sum_{R,1 \notin R} |R| \lambda_1^{|R|} \cdot \widehat{g}^2(R) \leq \epsilon_\tau \cdot \sum_{R,1 \notin R} |R| \lambda_1^{|R|} \widehat{f}_0^2(R) + (1 - \epsilon_\tau) \cdot \sum_{R,1 \notin R} |R| \lambda_1^{|R|} \widehat{f}_1^2(R) \quad (25)$$

Exactly as above, we have the following upper bound for the first summand: For a sufficiently large constant L_2 holds

$$\sum_{R, 1 \notin R} |R| \lambda_1^{|R|} \widehat{f}_0^2(R) = \sum_{R, 1 \notin R} |R| \lambda_1^{|R|} \left(\widehat{f}(R) + \widehat{f}(R \cup \{1\}) \right)^2 \leq L_2 \cdot \lambda_1 \cdot \gamma$$

Consider the second summand. The function f_1 is a boolean function, whose expectation equals $\widehat{f}(0) - \widehat{f}(\{1\}) = \alpha \cdot \mathbb{E} f \leq \alpha$. Similarly, $\mathbb{E} f_1^2 = \mathbb{E} f_1 \leq \alpha$.

We now apply the inequality of [20], which states that

For a boolean function $g : \{0, 1\}^m \rightarrow \{0, 1\}$ with expectation $\mu \leq 1/2$ holds $\sum_{k=1}^m \widehat{g}^2(\{k\}) \leq L_3 \cdot \mu^2 \cdot \ln(1/\mu)$, for a sufficiently large absolute constant L_3 .

In our case, this implies $\sum_{k=2}^n \widehat{f}_1^2(\{k\}) \leq L_3 \cdot \alpha^2 \cdot \ln\left(\frac{1}{\alpha}\right)$, for a sufficiently large constant L_3 .

This means that, for a sufficiently large constant L_4 , we can upperbound the second summand in (25) by

$$\sum_{R, 1 \notin R} |R| \lambda_1^{|R|} \widehat{f}_1^2(R) \leq L_4 \cdot \left(\lambda_1 \cdot \alpha^2 \ln\left(\frac{1}{\alpha}\right) + \lambda_1^2 \cdot \alpha \right)$$

Recall that for $\alpha < \lambda$, we have $\epsilon_\tau = \frac{1-\tau^{1/2}}{2} = \frac{1-(1-\lambda)/(1-\alpha)}{2} \leq L_5 \cdot \lambda$, for an absolute constant L_5 ; and that for $\alpha \geq \lambda$, we have $\epsilon_\tau = 0$. Plugging these estimates into (25), we have

$$\sum_{R, 1 \notin R} |R| \lambda_1^{|R|} \cdot \widehat{g}^2(R) \leq L_2 \cdot L_5 \cdot \lambda \cdot \lambda_1 \cdot \gamma + L_4 \cdot \left(\lambda_1 \cdot \alpha^2 \ln\left(\frac{1}{\alpha}\right) + \lambda_1^2 \cdot \alpha \right)$$

And hence, coming back to (24), and recalling that $\lambda = \tau \cdot \lambda_1$, we have, for sufficiently large absolute constants L, L' , that

$$\begin{aligned} \mathbb{E}_{T, 1 \in T} \text{Ent}\left(g_T \mid x_1 = 1, x_2, \dots, x_n\right) &\leq L' \cdot \left(\lambda_1 \cdot \gamma + \alpha^2 \ln\left(\frac{1}{\alpha}\right) + \lambda_1 \cdot \alpha \right) \leq \\ &L \cdot \left(\lambda_1 \cdot \gamma + \gamma^2 \ln\left(\frac{1}{\gamma}\right) \right) \end{aligned}$$

This completes the proof of (23), of Lemma 5.1, and of the first claim of the theorem.

■

B. Proof of the second claim of the theorem

First, note that if f is balanced, that is $\mathbb{E} f = \frac{1}{2}$, then so is $1 - f$, and the second claim of the theorem follows immediately from the first claim.

If $\mathbb{E} f \neq \frac{1}{2}$, some additional work is required. We only sketch the argument below, since it is very similar to the proof of the first claim.

Applying Theorem 1.7 to the function $1 - f$ gives (cf. (18))

$$\begin{aligned} \text{Ent}\left(T_\epsilon(1 - f)\right) &= \text{Ent}\left(T_{\epsilon_1}(1 - h)\right) \leq \lambda_1 \cdot \mathbb{E}_{T, 1 \in T} \left(\text{Ent}\left((1 - h) \mid T\right) - \text{Ent}\left((1 - h) \mid \{1\}\right) \right) + \\ &(1 - \lambda_1) \cdot \mathbb{E}_{T, 1 \notin T} \text{Ent}\left((1 - h) \mid T\right) + \mathbb{E}(1 - h) \cdot \phi\left(\text{Ent}\left(\frac{1 - h}{\mathbb{E}(1 - h)} \mid \{1\}\right), \epsilon_1\right) \end{aligned} \quad (26)$$

As in the proof of the first claim, we upperbound each of the three summands on the RHS of (26) separately.

Repeating the argument, with the necessary (minor) differences, leads to the same first two bounds:

•

$$\mathbb{E}_{T,1 \in T} \left(Ent\left((1-h) \mid T\right) - Ent\left((1-h) \mid \{1\}\right) \right) \leq O\left(\lambda_1 \cdot \gamma + \gamma^2 \ln\left(\frac{1}{\gamma}\right)\right)$$

•

$$\mathbb{E}_{T,1 \notin T} Ent\left((1-h) \mid T\right) \leq O\left(\lambda_1^2 \cdot \gamma + \lambda_1 \cdot \gamma^2 \ln\left(\frac{1}{\gamma}\right)\right)$$

Indeed, this should not be surprising since, roughly speaking, these two bounds for h are obtained by analysing the behavior of (the squares of) its non-trivial Fourier coefficients, and this is the same for h and for $1-h$.

As to the third summand, we will follow the argument in the proof of Lemma 5.3.

Let $g = \mathbb{E}\left(\frac{1-f}{\mathbb{E}(1-f)} \mid \{1\}\right)$. This is a function on a 2-point space $\{0,1\}$, with $g(0) = \rho$ and $g(1) = 2 - \rho$, where $\rho = 1 + \frac{(1-\alpha)(1-2\beta)}{1+2\beta}$.

Note that $\rho \leq 2 - c \cdot \gamma$, for some absolute constant $c > 0$. Hence, proceeding as in the proof of Lemma 5.3, gives

$$\phi\left(Ent\left(\frac{1-h}{\mathbb{E}(1-h)} \mid \{1\}\right), \epsilon_1\right) = Ent(T_\epsilon g) \leq (1 - H_2(\epsilon)) - c' \cdot \lambda \cdot \gamma,$$

for an absolute constant $c' > 0$.

Next, recall that in the proof of Lemma 5.3 we show $\phi\left(Ent\left(\frac{h}{\mathbb{E}h} \mid \{1\}\right), \epsilon_1\right) \leq (1 - H_2(\epsilon)) - c_1 \cdot \lambda \cdot \alpha$ for an absolute constant $c_1 > 0$.

Hence

$$\begin{aligned} \mathbb{E} h \cdot \phi\left(Ent\left(\frac{h}{\mathbb{E}h} \mid \{1\}\right), \epsilon_1\right) + \mathbb{E}(1-h) \cdot \phi\left(Ent\left(\frac{1-h}{\mathbb{E}(1-h)} \mid \{1\}\right), \epsilon_1\right) &\leq \\ &(1 - H_2(\epsilon)) - c_2 \cdot \lambda \cdot \gamma, \end{aligned}$$

for an absolute constant $c_2 > 0$.

We can now complete the proof of the second claim of the theorem.

Combining all bounds on the right hand sides of (18) and of (26) above gives

$$Ent(T_\epsilon f) + Ent(T_\epsilon(1-f)) \leq (1 - H_2(\epsilon)) - \Omega(\lambda \cdot \gamma) + o_{\lambda, \gamma \rightarrow 0}(\lambda \cdot \gamma)$$

Since $\lambda, \gamma \leq O(\delta)$, this implies that for a sufficiently small $\delta > 0$ holds

$$Ent(T_\epsilon f) + Ent(T_\epsilon(1-f)) \leq 1 - H_2(\epsilon)$$

■

A. Proof of Lemma 1.2

We have, for a boolean function f :

$$\begin{aligned} I(f(X); Y) &= H(f(X)) - H(f(X)|Y) = H(f(X)) - \mathbb{E}_y H(f(X)|Y=y) = \\ &H_2(\mathbb{E}f) - \mathbb{E}_y H_2((T_\epsilon f)(y)) \end{aligned}$$

We have $H_2(\mathbb{E}f) = \mathbb{E}f \log \frac{1}{\mathbb{E}f} + (1 - \mathbb{E}f) \log \frac{1}{1 - \mathbb{E}f}$.

We also have (all the logarithms are binary)

$$\begin{aligned} \mathbb{E}_y H_2((T_\epsilon f)(y)) &= \mathbb{E}_y \left((T_\epsilon f)(y) \log \frac{1}{(T_\epsilon f)(y)} + (1 - (T_\epsilon f)(y)) \log \frac{1}{1 - (T_\epsilon f)(y)} \right) = \\ &-\left(Ent(T_\epsilon f) + \mathbb{E} T_\epsilon f \log \mathbb{E} T_\epsilon f \right) - \left(Ent(T_\epsilon(1-f)) + \mathbb{E} T_\epsilon(1-f) \log \mathbb{E} T_\epsilon(1-f) \right) = \\ &-\left(Ent(T_\epsilon f) + Ent(T_\epsilon(1-f)) \right) + \mathbb{E}f \log \frac{1}{\mathbb{E}f} + (1 - \mathbb{E}f) \log \frac{1}{1 - \mathbb{E}f} \end{aligned}$$

In the last step we have used the fact $\mathbb{E} T_\epsilon g = \mathbb{E} g$ for any function g . The claim of the lemma follows. ■

B. Proof of Corollary 1.10

Applying Corollary 1.9 to the functions f and $1 - f$, we obtain, by Lemma 1.2:

$$I(f(X); Y) = Ent(T_\epsilon f) + Ent(T_\epsilon(1-f)) \leq \mathbb{E}_T \left(Ent(f | T) + Ent((1-f) | T) \right)$$

To conclude the proof of the corollary, it suffices to show that for any $T \subseteq [n]$ holds

$$Ent(f | T) + Ent((1-f) | T) = I(f(X); \{X_i\}_{i \in T})$$

To see this, we proceed exactly as in the proof of Lemma 1.2, observing that, by definition,

$$Pr\{f(X) = 1 \mid \{X_i\}_{i \in T}\} = \mathbb{E}(f | T)$$

Here we interpret both sides as functions of $\{x_i\}$, $i \in T$.

■

C. Proof of Theorem 1.12

The proof of this theorem is very similar to that of Theorem 1.7 and uses the notation and some of the results from that proof.

As in the proof of Lemma 4.3, our starting point is the chain rule for noisy entropy (7), which states that for any permutation $\sigma \in S_n$ the noisy entropy $Ent(T_\epsilon f)$ is bounded from above by

$$\sum_{i=1}^n \phi \left(Ent(f | \{\sigma(i)\}) + I_{T_{\{\sigma(1), \dots, \sigma(i-1)\}}} f \left(\{\sigma(1), \dots, \sigma(i-1)\}, \sigma(i) \right) \right)$$

Averaging over $\sigma \in S_n$ and using transitivity of action of the symmetric group and concavity of ϕ , this is at most ³¹

$$\sum_{k=0}^{n-1} \phi \left(\mathbb{E}_{i \in [n]} \text{Ent}(f | i) + b_k \right) \quad \text{where} \quad b_k = \mathbb{E}_{A, m} T_{\epsilon A} f(A, m)$$

and the expectation is over all $A \subseteq [n]$ of cardinality k and $m \notin A$. (In particular, we set $b_0 = 0$). Using the concavity of ϕ again, this is at most

$$n \cdot \phi \left(\mathbb{E}_{i \in [n]} \text{Ent}(f | i) + \frac{1}{n} \cdot \sum_{k=0}^{n-1} b_k \right)$$

The analysis in the proof of Theorem 1.7 shows that if T is a random subset of $[n]$ generated by sampling each element $i \in [n]$ independently with probability λ then

$$\begin{aligned} \sum_{k=0}^{n-1} b_k &= \frac{1}{\lambda} \cdot \sum_{s=1}^{n-1} w_s \cdot t_s = \frac{1}{\lambda} \cdot \mathbb{E}_T \left(\text{Ent}(f | T) - \sum_{i \in T} \text{Ent}(f | i) \right) = \\ &= \frac{1}{\lambda} \cdot \mathbb{E}_T \text{Ent}(f | T) - n \cdot \mathbb{E}_{i \in [n]} \text{Ent}(f | i) \end{aligned}$$

Substituting and simplifying, we get, setting $t = \lambda n$,

$$\text{Ent}(T_{\epsilon} f) \leq n \cdot \phi \left(\frac{\mathbb{E}_T \text{Ent}(f | T)}{t} \right)$$

which is the claim of the theorem. ■

1) *Proof of (6):* Let f be the distribution of X multiplied by 2^n . Then $\mathbb{E} f = 1$, and Theorem 1.12 can be applied.

By Section I-A1 and (5), we have $\text{Ent}(T_{\epsilon} f) = n - H(X \oplus Z)$ and $\text{Ent}(f | T) = |T| - H(\{X_i\}_{i \in T})$.

We also recall $\phi(x, \epsilon) = 1 - H_2(\epsilon + (1 - 2\epsilon) \cdot H_2^{-1}(1 - x))$.

Substituting in the claim of the theorem, and simplifying, gives

$$H(X \oplus Z) \geq n \cdot H_2 \left(\epsilon + (1 - 2\epsilon) \cdot H_2^{-1} \left(\frac{\mathbb{E}_T H(\{X_i\}_{i \in T})}{t} \right) \right)$$

which is the claim of (6).

■

D. Proof of Theorem 1.15

Let δ be the constant in the theorem. We will assume that δ is sufficiently small.

Let ϵ be a noise parameter, such that $(1 - 2\epsilon)^2 \leq \delta$. Denote $\lambda = (1 - 2\epsilon)^2$.

It is known (see [7]) that for any boolean function f holds $I(f(X); Y) \leq \lambda \cdot H_2(\mathbb{E} f)$. This immediately implies the validity of Conjecture 1.1 for boolean functions with expectation lying in $[0, c] \cup [1 - c, 1]$, for some absolute constant $0 < c < 1/2$.

In addition, we may assume, by symmetry, that $\mathbb{E} f \leq 1/2$. Combining these two observations, it remains to consider the case

$$c \leq \mathbb{E} f \leq 1/2 \tag{27}$$

Let f be a boolean function satisfying (27) with $I(f(X); Y) \geq 1 - H_2(\epsilon)$. This is the same as $Ent(T_\epsilon f) + Ent(T_\epsilon(1-f)) \geq 1 - H_2(\epsilon)$.

At this point, we need a technical lemma.

Lemma 6.1.: For any nonnegative non-zero function f holds⁵

$$Ent(T_\epsilon f) \leq \frac{1}{\mathbb{E} f} \cdot \left(\frac{1}{2 \ln 2} \cdot \sum_{k=1}^n \widehat{f}^2(\{k\}) \right) \cdot \lambda + O_{\lambda \rightarrow 0} \left(\frac{\mathbb{E} f^2}{\mathbb{E} f} \cdot \lambda^{4/3} \right) + O_{\lambda \rightarrow 0} \left(\frac{\mathbb{E}^2 f^2}{\mathbb{E}^3 f} \cdot \lambda^2 \right)$$

We will now proceed with the proof of the theorem, and prove the lemma below.

Applying Lemma 6.1 to functions f and $1-f$ and taking into account (27) gives

$$Ent(T_\epsilon f) + Ent(T_\epsilon(1-f)) \leq \frac{1}{\mathbb{E} f(1-\mathbb{E} f)} \cdot \left(\frac{1}{2 \ln 2} \cdot \sum_{k=1}^n \widehat{f}^2(\{k\}) \right) \cdot \lambda + O_{\lambda \rightarrow 0}(\lambda^{4/3})$$

Combining the two inequalities for $Ent(T_\epsilon f) + Ent(T_\epsilon(1-f))$, recalling $1 - H_2(\epsilon) \geq \frac{\lambda}{2 \ln 2}$, and using (27), we get

$$\sum_{k=1}^n \widehat{f}^2(\{k\}) \geq \mathbb{E} f \cdot (1 - \mathbb{E} f) - O(\lambda^{1/3})$$

For a boolean function $f : \{0, 1\}^n \rightarrow \{0, 1\}$ holds $\mathbb{E} f^2 = \mathbb{E} f$, and consequently $\sum_{S \neq \emptyset} \widehat{f}^2(S) = \mathbb{E} f(1 - \mathbb{E} f)$. Hence we have

$$\sum_{|S| \geq 2} \widehat{f}^2(S) \leq O(\lambda^{1/3})$$

Let $g = 2f - 1$. Then $g : \{0, 1\}^n \rightarrow \{-1, 1\}$, and $\sum_{|S| \geq 2} \widehat{g}^2(S) = 4 \cdot \sum_{|S| \geq 2} \widehat{f}^2(S) \leq O(\lambda^{1/3})$.

Note also that (27) implies $2c - 1 \leq \widehat{g}(0) = 2\mathbb{E} g - 1 \leq 0$.

Hence, assuming λ is sufficiently small, Theorem 5.5 implies that $|\mathbb{E} g| \leq O(\lambda^{1/3} \cdot \sqrt{\ln \frac{1}{\lambda}})$, and there exists an index $1 \leq k \leq n$ such that $\widehat{g}^2(\{k\}) \geq 1 - O(\lambda^{1/3})$.

This means that

$$\frac{1}{2} - O\left(\lambda^{1/3} \cdot \sqrt{\ln \frac{1}{\lambda}}\right) \leq \mathbb{E} f \leq \frac{1}{2} \quad \text{and} \quad |\widehat{f}(\{k\})| \geq \frac{1}{2} - O(\lambda^{1/3})$$

If λ is sufficiently small, f satisfies the conditions of Theorem 1.14. By the second claim of this theorem,

$$Ent(T_\epsilon f) + Ent(T_\epsilon(1-f)) \leq 1 - H_2(\epsilon),$$

completing the proof of Theorem 1.15.

■

⁵Asymptotic notation hides absolute constants independent of the remaining parameters.

1) *Proof of Lemma 6.1:* The argument below is a slight extension of an argument in [14].

In the following, we may and will assume, by homogeneity, that $\mathbb{E} f = 1$.

Let us introduce some notation. For $x \in \{0, 1\}^n$, let x^c be the complement of x , that is the element of $\{0, 1\}^n$ with $x_i^c = 1 - x_i$ for all $1 \leq i \leq n$.

For a nonnegative function g on $\{0, 1\}^n$, let g_0 be the 'even' part of g defined by $g_0(x) = (g(x) + g(x^c))/2$, and let $g_1 = g - g_0$ be the 'odd' part of g . By definition, $g_0(x) = g_0(x^c)$ and $g_1(x) = -g_1(x^c)$. Note also that $|g_1| \leq g_0$.

We will need the following well-known (and easy to verify) fact:

$$g_0 = \sum_{|S| \text{ even}} \widehat{g}(S) \cdot W_S \quad \text{and} \quad g_1 = \sum_{|S| \text{ odd}} \widehat{g}(S) \cdot W_S$$

We start with an auxiliary claim.

Lemma 6.2.: For a function g with $\mathbb{E} g = 1$ holds

$$Ent(g) = Ent(g_0) + \mathbb{E}_x g_0(x) \cdot \left(1 - H_2 \left(\frac{1 - |g_1(x)|/g_0(x)}{2} \right) \right)$$

Here for x such that $g_0(x) = g_1(x) = 0$, the expression $g_0(x) \cdot \left(1 - H_2 \left(\frac{1 - |g_1(x)|/g_0(x)}{2} \right) \right)$ is interpreted as 0.

Proof: We have

$$\begin{aligned} Ent(g) &= \mathbb{E}_x g(x) \log g(x) = \frac{1}{2} \cdot \mathbb{E}_x \left(g(x) \log g(x) + g(x^c) \log g(x^c) \right) = \\ &= \frac{1}{2} \cdot \mathbb{E}_x \left((g_0(x) + g_1(x)) \cdot \log (g_0(x) + g_1(x)) + (g_0(x) - g_1(x)) \cdot \log (g_0(x) - g_1(x)) \right) \end{aligned}$$

It is easy to verify that for any $0 \leq b \leq a$ holds

$1/2 \cdot ((a+b) \log(a+b) + (a-b) \log(a-b)) = a \log a + a \cdot \left(1 - H_2 \left(\frac{1-b/a}{2} \right) \right)$, where the last expression should be interpreted as 0 for $a = b = 0$.

Using this identity with $a = g_0(x)$ and $b = g_1(x)$ gives the claim of the lemma. ■

Next, as in [14], we upper bound $1 - H_2 \left(\frac{1-x}{2} \right)$ by $\frac{1}{2 \ln 2} \cdot x^2 + \left(1 - \frac{1}{2 \ln 2} \right) \cdot x^4$.

Substituting this bound in the claim of Lemma 6.2 gives

$$Ent(g) \leq Ent(g_0) + \frac{1}{2 \ln 2} \cdot \mathbb{E}_x \frac{g_1^2(x)}{g_0(x)} + \left(1 - \frac{1}{2 \ln 2} \right) \cdot \mathbb{E}_x \frac{g_1^4(x)}{g_0^3(x)}$$

Let $g = T_\epsilon f$. It is easy to verify $(T_\epsilon f)_i = T_\epsilon(f_i)$ for any function f and $i = 0, 1$. Consequently, we get the bound

$$Ent(T_\epsilon f) \leq Ent(T_\epsilon f_0) + \frac{1}{2 \ln 2} \cdot \mathbb{E}_x \frac{(T_\epsilon f_1(x))^2}{T_\epsilon f_0(x)} + \left(1 - \frac{1}{2 \ln 2} \right) \cdot \mathbb{E}_x \frac{(T_\epsilon f_1(x))^4}{(T_\epsilon f_0(x))^3}$$

We upperbound each of the summands on the RHS separately.

1) The first summand. Note that $\mathbb{E} T_\epsilon f_0 = \mathbb{E} f_0 = \mathbb{E} f = 1$. Recall also (see e.g., [14]) that for any function g on $\{0, 1\}^n$ holds $T_\epsilon g = \sum_S \lambda^{|S|/2} \widehat{g}(S) \cdot W_S$.

Hence, by Lemma 5.4,

$$\begin{aligned} Ent(T_\epsilon f_0) &\leq O\left(\sum_S |S| \cdot \widehat{T_\epsilon f_0}^2(S)\right) = O\left(\sum_S |S| \lambda^{|S|} \widehat{f_0}^2(S)\right) = \\ &O\left(\sum_{|S| \text{ even}} |S| \lambda^{|S|} \widehat{f}^2(S)\right) = O(\mathbb{E} f^2 \cdot \lambda^2) \end{aligned}$$

2) The second summand.

First, we argue that $T_\epsilon f_0$ is bounded away from 0 with high probability. Recall that $\mathbb{E} T_\epsilon f_0 = 1$, and note that $Var(T_\epsilon f_0) = \sum_{S \neq \emptyset} \widehat{T_\epsilon f_0}^2(S) = O(\lambda^2 \cdot \mathbb{E} f^2)$. Hence, by Chebyshev's inequality, for any $0 \leq \alpha < 1$ holds

$$Pr\{T_\epsilon f_0 \leq \alpha\} \leq O\left(\frac{\lambda^2 \cdot \mathbb{E} f^2}{(1 - \alpha)^2}\right) \quad (28)$$

Second, recall that for any x holds $|T_\epsilon f_1(x)| \leq T_\epsilon f_0(x)$, and hence $\frac{(T_\epsilon f_1(x))^2}{T_\epsilon f_0(x)} \leq T_\epsilon f_0(x)$.

Therefore, taking $\alpha = 1 - \lambda^{1/3}$ in (28), we have $\mathbb{E}_x \frac{(T_\epsilon f_1(x))^2}{T_\epsilon f_0(x)}$ bounded from above by

$$Pr\{f_0 \leq 1 - \lambda^{1/3}\} \cdot \alpha + \left(1 + O(\lambda^{1/3})\right) \cdot \mathbb{E}_x (T_\epsilon f_1(x))^2$$

Recalling that $T_\epsilon f_1 = \sum_{|S| \text{ odd}} \lambda^{|S|/2} \widehat{f}(S) \cdot W_S$, this is at most

$$\begin{aligned} O(\mathbb{E} f^2 \cdot \lambda^{4/3}) + \left(1 + O(\lambda^{1/3})\right) \cdot \left(\left(\sum_{k=1}^n \widehat{f}^2(\{k\})\right) \cdot \lambda + O(\mathbb{E} f^2 \cdot \lambda^3)\right) = \\ \left(\sum_{k=1}^n \widehat{f}^2(\{k\})\right) \cdot \lambda + O(\mathbb{E} f^2 \cdot \lambda^{4/3}) \end{aligned}$$

3) The third summand. Note that $\mathbb{E} f_1 = 0$. Hence, as in Lemma 1 in [14] (where the requirement on f to be boolean does not seem to be necessary) we have, for a sufficiently small λ , that

$$\mathbb{E}_x (T_\epsilon f_1(x))^4 \leq O\left(\left(\mathbb{E}_x f_1^2(x)\right)^2 \cdot \lambda^2\right) = O(\mathbb{E}^2 f^2 \cdot \lambda^2)$$

We can now upperbound the third summand using the Chebyshev inequality, as above. Taking $\alpha = 1/2$ in (28) upperbounds $\mathbb{E}_x \frac{(T_\epsilon f_1(x))^4}{(T_\epsilon f_0(x))^3}$ by

$$O(\mathbb{E} f^2 \cdot \lambda^2) \cdot \alpha + O(\mathbb{E}^2 f^2 \cdot \lambda^2) = O(\mathbb{E} f^2 \cdot \lambda^2) + O(\mathbb{E}^2 f^2 \cdot \lambda^2)$$

Combining these estimates leads to the claim of Lemma 6.1. ■

■

ACKNOWLEDGMENTS

We are grateful to Yuval Kochman, Or Ordentlich, and Yuri Polyanskiy for many very helpful conversations and valuable remarks. We also thank Venkat Chandar for valuable remarks.

We would also like to thank the anonymous referees for their numerous suggestions, which led to a significant improvement in the presentation of this paper.

- [1] N. Alon and J. Spencer, **The Probabilistic Method**, 3rd ed. Hoboken, NJ: Wiley, 2008.
- [2] R. Ahlswede and P. Gacs, *Spreading of sets in product spaces and hypercontraction of the Markov operator*, Ann. Probab., vol. 4, no. 6, pp. 925-939, 1976.
- [3] V. Chandar, personal communication, 2014.
- [4] E. Friedgut, G. Kalai, and A. Naor, *Boolean functions whose Fourier transform is concentrated on the first two levels*, Advances in Applied Mathematics, vol. 29, 3, pp. 427-437, 2002.
- [5] T. S. Han, *Nonnegative entropy measures of multivariate symmetric correlations*, Inform. Contr. 36, pp. 133-156, 1978.
- [6] J. Jendrej, K. Oleszkiewicz, and J. O. Wojtaszczyk, *On some extensions of the FKN theorem*, Theory of Computing, vol. 11 (18), pp. 445-469, 2015.
- [7] T. Courtade and G. Kumar, *Which boolean functions maximize mutual information on noisy inputs?* IEEE Transactions on Information Theory, vol. 60, no. 8, pp. 4515-4525, 2014.
- [8] M. Ledoux, **Concentration of measure and logarithmic Sobolev inequalities**, 1997.
- [9] E. Mossel and S. Sachdeva, personal communication, 2014.
- [10] A. W. Marshall, and I. Olkin, **Inequalities: Theory of Majorization and Its Applications**, Academic, New York, 1979.
- [11] R. O'Donnell, **Analysis of Boolean functions**, Cambridge U.P., 2014.
- [12] O. Ordentlich, personal communication, 2015.
- [13] O. Ordentlich, *On the entropy rate of binary hidden Markov processes*, 2016, arXiv preprint. [Online]. Available: <http://arxiv.org/abs/1601.06453>.
- [14] O. Ordentlich, O. Shayevitz, and O. Weinstein, *An Improved Upper Bound for the Most Informative Boolean Function Conjecture*, 2015, arXiv preprint. [Online]. Available: <http://arxiv.org/abs/1505.05794>.
- [15] Y. Polyanskiy, personal communication, 2015.
- [16] Y. Polyanskiy and Y. Wu, *Dissipation of information in channels with input constraints*, 2014, arXiv preprint. [Online]. Available: <http://arxiv.org/abs/1405.3629>.
- [17] Y. Polyanskiy and Y. Wu, *Strong data-processing inequalities for channels and Bayesian networks*, 2016, arXiv preprint. [Online]. Available: <http://arxiv.org/abs/1508.06025>.
- [18] S. Sachdeva, A. Samorodnitsky, and I. Shahaf, *On conjectures of Kumar and Courtade*, manuscript, 2014.
- [19] A. Samorodnitsky, manuscript, 2015.
- [20] M. Talagrand, *How much are increasing sets positively correlated?*, Combinatorica 16.2, 243-258, 1996.
- [21] A. D. Wyner and J. Ziv, *A theorem on the entropy of certain binary sequences and applications: Part I*, IEEE Trans. Inform. Theory, vol. 19, no. 6, pp. 769-772, 1973.

Alex Samorodnitsky received the B.A. degree in computer science and mathematics (1987), the M.Sc. degree in mathematics (1990), and the Ph.D. degree in mathematics (1998), all from the Hebrew University of Jerusalem.

He is now a Professor of Computer Science at the School of Engineering and Computer Science in the Hebrew University, which he joined in 2001 as an Assistant Professor. His research interests are in theoretical computer science, coding theory, and combinatorics.