

# Deterministic Compression with Uncertain Priors\*

Elad Haramaty<sup>†</sup>

Madhu Sudan<sup>‡</sup>

September 12, 2015

## Abstract

Communication in “natural” settings, e.g., between humans, is distinctly different from that in classical designed settings, in that the former is invariably characterized by the sender and receiver not being in perfect agreement with each other. Solutions to classical communication problems thus have to overcome an extra layer of uncertainty introduced by this lack of prior agreement. One of the classical goals of communication is compression of information, and in this context lack of agreement implies that sender and receiver may not agree on the “prior” from which information is being generated. Most classical mechanisms for compressing turn out to be non-robust when sender and receiver do not agree on the prior. Juba et al. (Proc. ITCS 2011) showed that there do exist compression schemes with *shared randomness* between sender and receiver \*that do not share a prior\* that can compress information down roughly to its entropy.

In this work, we explore the assumption of *shared randomness* between the sender and receiver and highlight why this assumption is problematic when dealing with natural communication. We initiate the study of deterministic compression schemes amid uncertain priors, and expose some of the mathematical facets of this problem. We show some non-trivial deterministic compression schemes, and some lower bounds on natural classes of compression schemes. We show that a full understanding of deterministic communication turns into challenging (open) questions in graph theory and communication complexity.

**Keywords:** Source coding, communication complexity, graph coloring

## 1 Introduction

In this work we consider the task of compressing information deterministically, in settings where the sender and receiver are not in agreement with each other on the distribution from which the information is being generated. We start by first describing the general motivation for the study of this problem before formally describing the problem and our results.

---

\*An extended abstract of this paper appeared in the Proceedings of *Innovations in Theoretical Computer Science, ITCS'14*, Princeton, NJ, USA, January 12-14, 2014.

<sup>†</sup>Department of Computer Science, Technion, Haifa. [eladh@cs.technion.ac.il](mailto:eladh@cs.technion.ac.il). Work done in part when this author was visiting Microsoft Research New England.

<sup>‡</sup>Microsoft Research New England, One Memorial Drive, Cambridge, MA 02139, USA. [madhu@mit.edu](mailto:madhu@mit.edu).

## 1.1 Natural Communication: Context and Uncertainty

Natural communication, say between two humans, differs in significant ways from “designed communication”, say between a cell phone and its nearby tower. The latter is carefully engineered to optimize use of the channel of communication, while introducing careful redundancy to overcome any unreliability of the channel. The resulting problems and solution concepts, such as compression schemes and error-correcting encoding schemes are well understood by now.

Natural communication, as characterized by the vagaries of natural language, is much less understood mathematically. Natural communication, often characterized by dictionaries and grammars, does not follow the rules they prescribe. They tend to be ambiguous locally, and seemingly needlessly redundant in other cases; without offering the same reliability as error-correcting codes do. At the same time, natural communication has a remarkable capacity to overcome a lack of “perfect engineering” of the sender or receiver of information, and in particular does not seem to require perfect agreement between the two (on the design of the protocol). Understanding this resilience mathematically would be a fruitful pursuit and forms the motivation for this and prior works in this stream.

The main goal here is to understand what “language” looks like and what parameters it tries to optimize. Roughly, language takes an intended message in the sender’s mind and attempts to describe how to convert this message to a sequence of, say, bits. This part is similar to the encoding map of an error-correcting scheme. However to model language more precisely, one should really take into account its context-sensitivity. A more precise description of language would be as a *encoding* map from a pair  $(context, message)$  to a *word* which is a sequence of bits. For the time being, one may view context and message as elements of some abstract space (each being possibly a countable set). The language also gives rules on how a receiver should recover the message from the word: it does so by applying a *decoding* map that maps a pair  $(context, word)$  to a *message*. The goal hopefully is to make sure the receiver decodes to a message that is somehow *compatible* with the sender’s intent. Defining compatibility is complex, and we will skirt this issue in this paper. (Such issues are considered, e.g., in [7, 4].) We will simply require that the receiver decodes to the *same* message as the sender wished to send (so in particular the message spaces for the sender and the receiver are identical). Even with this simplification, there remains a major hurdle to communication, namely that sender’s context and receiver’s context may not be identical! Remarkably natural communication manages to function reasonably well even when these contexts are not the same, but close enough to each other, and it is this aspect that is the focus of this work.

In order to model contexts that are not identical, but are reasonably close, we need a more structured view of contexts (than merely as elements of an abstract set). Juba et al. [6] proposed a natural view: namely a context is simply a probability distribution on messages: The distribution that describes the message that the receiver is expecting to receive, or the distribution that the sender thinks the receiver is expecting to receive. In this framework, the encoding function becomes a map from a distribution  $P$  supported on some set  $U$  and message  $m \in U$  to a word  $w \in \{0, 1\}^*$ ; the decoding function becomes a map from a distribution  $Q$  also supported on  $U$  and a word  $w \in \{0, 1\}^*$  to a message  $\hat{m} \in U$ ; and the goal of the language is to ensure that  $\hat{m} = m$  even if  $P \neq Q$ , provided they are reasonably close. At the same time, the language would like to reduce the expected length of the communication, assuming say that the messages are generated from the distribution  $P$ . (More generally we could consider a setting where the messages are generated from some distribution  $R$ ; our assumption that  $R = P$  is mostly for simplicity.)

Classical communication deals with the setting where  $P = Q$ . In the classical setting, it

turns out that the entropy of the distribution  $P$  precisely describes the expected length of the transmission. Does entropy still give a good measure characterizing the expected length of the transmission in the natural setting, when  $P \neq Q$ ? (We note that classical communication also considers the setting where  $R \neq P = Q$  and this leads to notions such as KL-divergence, but the important aspect is that it do not consider settings where sender and receiver disagree on the prior.)

Juba et al. [6] consider the following notion of distance between distributions (a notion which also is used commonly in differential privacy):  $\delta(P, Q) = \max_{m \in U} \{|\log_2 P(m) - \log_2 Q(m)|\}$ . (If  $\delta(P, Q) \leq \Delta$  then for every  $m \in U$  we have  $2^{-\Delta}Q(m) \leq P(m) \leq 2^{\Delta}Q(m)$ .) With this notion of distance in place, Juba et al. roughly show that entropy is a good measure of the compression length: Specifically they show that there exist encoding and decoding schemes that use shared randomness between sender and receiver and manage to compress information to a length of roughly  $H(P) + 2\delta(P, Q)$  where  $H(\cdot)$  denotes the binary entropy function. (We note that works of Harsha et al. [5] and Braverman and Rao [1] also explore questions somewhat similar to the ones considered by Juba et al., though their motivations were quite different. Both these latter works focus on the setting when the sender and receiver have different priors and are trying to generate a random variable that is maximally correlated under their priors. In our case, the sender gets a concrete message from its prior and wishes to communicate it. The focus in both works is on randomized solutions that get the communication complexity down to the minimum possible amount, whereas our thrust is to use less (or no) randomness at the expense of slightly larger communication complexity.)

The main goal of this paper is to explore the need for shared randomness in their work. The assumption of shared randomness takes the solution further away from the motivation of natural communication. Roughly their solution suggests that if the “dictionary”, or any other codebook associated with language, is random, then information can be compressed. However in natural usage, the dictionary remains somewhat static whereas the contexts for communication vary vastly. Furthermore, it is even very plausible that the dictionary influences our communication and its context. We thus feel that a compression scheme based on shared randomness is not sufficient evidence to suggest that entropy is the natural complexity measure for capturing the complexity of communicating in natural settings. This leads us to study deterministic compression schemes in this paper, and as we stress below, this turns out to be surprisingly challenging to analyze. We illustrate the complexity by considering a rather simple problem in this space.

## 1.2 A toy problem

The following example illustrates the questions studied in this paper: Suppose Alice and Bob have a ranking of a set  $U$  of  $N$  elements, say, movies. Specifically Alice’s rank function is  $A: [N] \rightarrow U$  and Bob’s rank function is  $B: [N] \rightarrow U$  where  $[N] = \{1, \dots, N\}$  and  $A$  and  $B$  are bijections with  $A(i)$  naming the  $i$ th ranked movie in Alice’s ranking. Suppose further that Alice and Bob know that their rankings are “close”, specifically for every  $x \in U$ ,  $|A^{-1}(x) - B^{-1}(x)| \leq 2$ . How many bits does Alice have to send to Bob so that Bob knows her top-ranked movie, i.e.,  $A(1)$ ?

On the one hand Bob knows  $A(1)$  is one of the three element set  $S_1 = \{B(1), B(2), B(3)\}$  and so the information-content from his point of view is bounded by  $\log_2 3$  bits. Indeed this leads to a randomized communication scheme, with Alice and Bob sharing common randomness with  $O(1)$  bits of communication: Alice and Bob use their randomness to get a hash function hashing their universe to a constant number of bits. Alice sends the hash of  $A(1)$ , and Bob recovers the name of the movie provided the elements of  $S_1$  hash to distinct values. However the deterministic

communication complexity of the question is not as easily settled. Part of the reason is that Alice does not know  $S_1$  and so has to “guess” it to communicate  $A(1)$ . Still she is not clueless: She knows it is contained in  $T_2 = \{A(1), \dots, A(5)\}$  and perhaps this can help her communicate  $A(1)$  efficiently to Bob. The question of interest to us in this work is: Can Alice communicate  $A(1)$  to Bob with a number of bits that is independent of  $N$ ? (Unfortunately, we do not answer this question, though we do give a non-trivial upper bound. We will elaborate on this later.)

The question above is a prototypical example of “communication amid uncertainty”, where the communicating players have fairly good information about each other (in the example above Alice and Bob know each others ranking of each movie to within  $\pm 2$ ), but are not sure of each other’s information and do not have a common-ground to base communication on. As we elaborate on in this paper, solutions to this problem influence solutions to the general problems of communication amid uncertainty, while this problem is itself a special case (when Alice and Bob’s distributions are geometric).

We describe our problems and solutions shortly, but to give a gist of the findings: We show in this work that there is a solution to the above toy problem communicating roughly  $2^{O(\log^* N)}$  bits, which is a very slowly growing, but nevertheless a growing, function of  $U$ . We are unable to resolve if a growing function of  $N$  bits are necessary for this problem; however we do show that any solution that looks at only a constant number of Alice’s top movies (or a constant number of Bob’s top movies) must communicate  $\log^{(\Omega(1))} N$  bits (where  $\log^{(i)}$  denotes the  $i$ th iterated logarithm function).

### 1.3 Formal definitions and main results

We start by defining the notion of an “uncertain compression scheme”.

We let  $\{0, 1\}^*$  denote the set of all finite length binary strings. For  $x \in \{0, 1\}^*$ , let  $|x|$  denote its length. Throughout  $U$ , the set of all messages, will be a finite set of size  $N$ . Let  $\mathcal{P}(U)$  denote the space of all probability distributions over  $U$ . Again, we say that two distributions  $P, Q$  are  $\Delta$ -close, if for every  $m \in U$  we have  $2^{-\Delta}Q(m) \leq P(m) \leq 2^{\Delta}Q(m)$ .

**Definition 1.1** ((Basic) Uncertain Compression Scheme). *For positive real  $\Delta$ , an Uncertain Compression Scheme (UCS) for distance  $\Delta$  over the universe  $U$  is given by a pair of functions  $E: \mathcal{P}(U) \times U \rightarrow \{0, 1\}^*$  and  $D: \mathcal{P}(U) \times \{0, 1\}^* \rightarrow U$  that satisfy the following correctness condition: For every pair of distributions  $P, Q \in \mathcal{P}(U)$  that are  $\Delta$ -close and for every  $m \in U$ , we have  $D(Q, E(P, m)) = m$ . The performance of a UCS  $(E, D)$  is given by the function  $L: \mathcal{P}(U) \rightarrow \mathbb{R}^+$ , where  $L(P) = \mathbf{E}_{m \leftarrow P}[|E(P, m)|]$ , i.e., the expected length of the encoding under the distribution  $P$ . We refer to such a scheme as a  $(\Delta, L)$ -UCS.*

In English, the definition above explicitly provides the distribution as input to the encoding and decoding schemes, and expects the schemes to work correctly even if the distributions used by the encoder and decoder are not the same, as long as they are  $\Delta$ -close to each other. While in general we would like compression schemes which work for all possible distributions  $P, Q$  that are within  $\Delta$  of each other, and with no error (as expected in the definition above), some of our schemes are weaker and work with some error, or only for some class of distributions. We define such general UCS’s below.

**Definition 1.2** ((General) Uncertain Compression Scheme). *For positive real  $\Delta$  (for distance),  $\epsilon \in [0, 1]$  (for error), a class of distributions  $\mathcal{F} \subseteq \mathcal{P}(U)$ , and performance function  $L: \mathcal{F} \rightarrow \mathbb{R}^+$*

a  $(\Delta, \epsilon, \mathcal{F}, L)$ -Uncertain Compression Scheme (UCS) over the universe  $U$  is given by a pair of  $E: \mathcal{F} \times U \rightarrow \{0, 1\}^* \cup \{\perp\}$  and  $D: \mathcal{F} \times \{0, 1\}^* \cup \{\perp\} \rightarrow U \cup \{\perp\}$  that satisfy the following conditions:

1. For every pair of distributions  $P, Q \in \mathcal{F}$  that are  $\Delta$ -close and for every  $m \in U$ , it is the case that if  $E(P, m) \neq \perp$  then  $D(Q, E(P, m)) = m$ . Furthermore  $D(Q, \perp) = \perp$ .
2. For every  $P \in \mathcal{F}$ , we have  $\Pr_{m \leftarrow P} U[E(P, m) = \perp] \leq \epsilon$ .
3. For every  $P \in \mathcal{F}$ , we have  $\mathbf{E}_{m \leftarrow P} U[|E(P, m)|] \leq L(P)$ .

Note that we do not distinguish the two definitions above by name, but rather just by the number of parameters. So if the number of parameters is just two, then it is assumed that there is no error, and the performance holds for all distributions.

We note that the definitions above only cover deterministic compression schemes. A compression scheme with shared randomness can be defined analogously, but we do not do so here. We also stress that the choice of  $P$  and  $Q$  is “worst-case” within the family  $\mathcal{F}$  (as formalized by the universal quantifier in the correctness condition). There are no assumptions that  $\mathcal{F}$  is small (has only finitely many elements), which tends to be the setting for universal compression. Similarly, we do not consider a sequence of messages that need to be transmitted: Rather, we are considering one-shot communication with no assumptions on the distributions  $P$  and  $Q$ , other than that they are from  $\mathcal{F}$  and  $\Delta$ -close.

We recall that Juba et al. present a  $(\Delta, H(P) + 2\Delta + c)$ -UCS (with shared randomness) for some constant  $c \leq 3$ . We give two *deterministic* schemes in this paper, both having complexity depending on the universe size,  $N$ , but both using substantially less than  $\log N$  bits.

**Theorem 1.3.** *For every  $\Delta \geq 0$ , there exists a  $(\Delta, O(H(P) + \Delta + \log \log N))$ -UCS, i.e., a deterministic universal compression scheme that works for all pairs  $P, Q$  that are within distance  $\Delta$  of each other, and where the expected length of encoding is at most  $O(H(P) + \Delta + \log \log N)$ .*

The dependence on  $N$  of this scheme is non-trivial and thus may even be reasonable in “natural circumstances”. However it is not clear if such a dependency on  $N$  is necessary. Motivated by the quest to understand the dependence on  $N$  more closely, we explore schemes whose performance is not necessarily linear in  $H(P)$ . Simultaneously we relax our schemes to allow them to “drop” messages with  $\epsilon$  probability. We note that if we do not do the latter, then the former is not really a relaxation: Any error-free scheme with superlinear dependence on  $H(P)$  can be converted to one with linear dependence on  $H(P)$  by a simple reduction (see Lemma 3.13).

Our next theorem gives a scheme that is weaker than the one from Theorem 1.3 in its dependence on the entropy  $H(P)$  and in that it errs with non-zero probability. But it does achieve significantly better dependence on  $N$ .

**Theorem 1.4.** *For every  $\epsilon > 0$  and  $\Delta \geq 0$  there exists a  $(\Delta, \epsilon, \mathcal{P}(U), \exp(H(P)/\epsilon + \Delta \log^* N))$ -UCS, i.e., the scheme has error probability at most  $\epsilon$ , it works for all pairs of distributions  $P, Q$  within distance  $\Delta$  and the expected length of the encoding is at most  $\exp(H(P)/\epsilon + \Delta \log^* N)$ .*

In the above the notation  $\exp(x)$  denotes a function of the form  $c^x$  for some universal constant  $c$ , and  $\log^* N$  denotes the minimum integer  $i$  such  $\log^{(i)} N \leq 1$  and  $\log^{(i)}$  is the logarithm function iterated  $i$  times.

An alternate way to get around the barrier of Lemma 3.13, which insists that schemes must have linear dependence on  $H(P)$  or make some error, is to have schemes that do not work for

all possible pairs of distributions  $P$  and  $Q$ . As it turns out the scheme from Theorem 1.4 does have this behavior for many natural distributions. In Theorem 3.7 we show that our scheme from Theorem 1.4 works without error and with same performance as long as  $P$  (or  $Q$ ) are close to a “flat distribution” (uniform over a subset), or a geometric distribution, or a binomial distribution. We stress that the scheme is not particularly carefully tailored to the class of distributions (though of course the encodings and decodings do depend on the distributions), but naturally adapts to being error-free for the above classes.

## 1.4 Techniques: Graph Coloring

While the most natural framework for studying our problem is as a question of communication complexity of a relational problem (as in [8]), this turns out not to be the most useful for studying the deterministic communication complexity. Indeed, as pointed out earlier, the modern stress in communication complexity is often on designing and understanding the limits of protocols that are interactive and use shared randomness, while in our case the thrust is in the opposite direction.

It turns out that our questions are naturally captured as graph-coloring questions. Furthermore such questions (or related ones) have been studied in the literature on distributed computing in the attempt to color graphs in a local distributed manner. In particular, the work of Linial [9] shows that a “local” algorithm for 3-coloring a cycle, due to Cole and Vishkin [2], implies that a large “high-degree graph” is 3-colorable. The ideas of Cole and Vishkin [2] and Linial [9] turn out to be quite useful in our context. Our work abstracts some of these techniques, and extends them to get combinatorial results, which we then convert to efficient compression schemes.

**Uncertainty graphs and Chromatic number** We start by defining a class of structured combinatorial graphs whose chromatic number turns out to be central to our problems. Recall,  $[N] = \{1, \dots, N\}$ . Let  $S_N$  denote the set of all permutations on  $N$  elements, i.e., the set of all bijections from  $[N]$  to itself. For  $\pi, \sigma \in S_N$ , let  $\delta(\pi, \sigma) = \max_{i \in [N]} |\pi^{-1}(i) - \sigma^{-1}(i)|$ .

**Definition 1.5** (Uncertainty graphs). *For integer  $N, \ell$  the uncertainty graph  $\mathcal{U}_{N, \ell}$  has elements of  $S_N$  as its vertices, with edges  $\pi \leftrightarrow \sigma$  if (1)  $\pi(1) \neq \sigma(1)$  and (2)  $\delta(\pi, \sigma) \leq \ell$ .*

It turns out that the chromatic number of the uncertainty graphs have a close connection to uncertain communication schemes. Roughly these graphs emerge from a very restricted version of the communication problem, where the distributions  $P$  and  $Q$  are geometric distributions (giving probability proportional to  $\beta^{-\pi^{-1}(i)}$  and  $\beta^{-\sigma^{-1}(i)}$  to the element  $i \in [N]$ ). It follows that if  $\delta(\pi, \sigma)$  is small, then  $P$  and  $Q$  are close to each other. Furthermore, for simplicity these graphs only consider the case that the message is the element with maximal probability under  $P$ . To understand how the chromatic number plays a role, fix a receiver with distribution  $Q$  and consider two possible senders  $P$  and  $P'$  that could communicate with this receiver. Consider coloring  $P$  and  $P'$  by  $E(P, \operatorname{argmax}_m \{P(m)\})$  and  $E(P', \operatorname{argmax}_m \{P'(m)\})$  respectively. This would lead to distinct colors on pairs  $P$  and  $P'$  that are too close to each other, provided their messages, i.e.,  $\operatorname{argmax}_m \{P(m)\}$  and  $\operatorname{argmax}_m \{P'(m)\}$  are different. This exactly corresponds to adjacency in our graph: the underlying permutations  $\pi$  and  $\sigma$  are close, and the top ranked elements are different.

The results of Juba et al. imply that the “fractional chromatic number” of  $\mathcal{U}_{N, \ell}$  is bounded by  $O(\ell)$ .<sup>1</sup> The (integral) chromatic number on the other hand does not immediately seem to be

---

<sup>1</sup>The fractional chromatic number of a graph  $G$  is the smallest positive real  $w$  such that there exists a collection

bounded as a function of  $\ell$  alone. The implication of the low fractional chromatic number is that the chromatic number of  $\mathcal{U}_{N,\ell}$  is at most  $O(\ell N \log N)$ , but this is worse than the naive upper bound of  $N$ , which can be obtained by setting the color of  $\pi$  to be  $\pi^{-1}(1)$ . (By definition of adjacency this is a valid coloring.) Our main technical contribution is in obtaining some non-trivial upper bounds on the chromatic number of this graph.

To derive our upper bounds, we look at “coarsened” versions of the graph  $\mathcal{U}_{N,\ell}$ . For positive integer  $k \leq N$ , we say that  $\pi: [k] \rightarrow [N]$  is a  $k$ -subpermutation if  $\pi$  is injective. We let  $S_{N,k}$  denote the set of all  $k$ -subpermutations on  $[N]$ . For  $k' \geq k$ , we say subpermutation  $\pi: [k'] \rightarrow [N]$  extends the subpermutation  $\sigma: [k] \rightarrow [N]$  if  $\sigma(i) = \pi(i)$  for all  $i \in [k]$ . For  $k$ -subpermutations  $\pi$  and  $\sigma$ , we let  $\delta(\pi, \sigma) = \min_{\pi', \sigma' \in S_N \text{ extending } \pi, \sigma} \{\delta(\pi', \sigma')\}$ .

**Definition 1.6** (Restricted Uncertainty graphs). *For integers  $N, \ell$  and  $k$  the  $k$ -restricted uncertainty graph  $\mathcal{U}_{N,\ell,k}$  has elements of  $S_{N,k}$  as its vertices, with  $\pi \leftrightarrow \sigma$  if (1)  $\pi(1) \neq \sigma(1)$  and (2)  $\delta(\pi, \sigma) \leq \ell$ .*

Note that  $\mathcal{U}_{N,\ell,N} = \mathcal{U}_{N,\ell}$ . We derive our upper bounds on the chromatic number of  $\mathcal{U}_{N,\ell}$  by giving non-trivial upper bounds on the chromatic number of  $\mathcal{U}_{N,\ell,k}$ .

**Lemma 1.7.** 1. For every  $k \leq k'$ ,  $\chi(\mathcal{U}_{N,\ell,k'}) \leq \chi(\mathcal{U}_{N,\ell,k})$ .

2. For every  $N, \ell$ ,  $\chi(\mathcal{U}_{N,\ell,\ell+1}) \leq O(\ell^2 \log N)$ .

3. For every  $N, \ell$  and  $k \leq \ell \log^* N$  that is an integral multiple of  $\ell$ , we have  $\chi(\mathcal{U}_{N,\ell,k}) \leq 2^{O(k \log \ell)} \log^{(k/\ell)} N$ .

4. For every  $N, \ell$  and  $k \leq \frac{1}{2} \ell \log^* N$  that is an integral multiple of  $\ell$ , we have  $\chi(\mathcal{U}_{N,\ell,k}) \geq \log^{(2k/\ell)}(N/\ell)$ .

As an immediate application we get the following theorem.

**Theorem 1.8.** For every  $N$  and  $\ell$ , we have  $\chi(\mathcal{U}_{N,\ell}) \leq O(\min\{\ell^2 \log N, 2^{O(\ell \log \ell \log^* N)}\})$ .

Unfortunately, the lower bound from Part (4) of Lemma 1.7 becomes trivial as  $k \rightarrow N$  and so we don't get a growing function of  $N$  as a lower bound. However, it does rule out most natural strategies for coloring  $\mathcal{U}$ , and shows limitations of the *intuition* that suggests  $\mathcal{U}$  may be colorable with  $f(\ell)$  colors independent of  $N$ . This is so since the intuition, as well as most natural strategies, only use the top  $O(\ell)$  ranking elements of a permutation  $\pi$  to determine its color; and such the lower bound shows that such strategies are inherently limited. In particular, it shows that there is no hope to extend the methods of Juba et al. (which were based on this intuition) in a simple way to get a deterministic UCS.

## 1.5 Directions for further work

Given that most of the questions raised in this work have not found tight answers, there is an obvious number of natural questions to resolve here — the most fundamental one being whether one can compress information down to its entropy (to within constant multiplicative factors) deterministically (or even just with private randomness) in the uncertain setting.

---

of independent sets  $I_1, \dots, I_t$  in  $G$  with weights  $w_1, \dots, w_t$  such that  $\sum_{j=1}^t w_j = w$  and for every vertex  $u \in V(G)$  it is the case that  $\sum_{j: I_j \ni u} w_j \geq 1$ .

In addition to resolving open questions, a number of modelling challenges remain in trying to understand natural modes of communication. Language is an organically evolved concept with communicational, computational and societal pressures acting on it. The game that was played out with natural languages over the past millenia is now getting played out at a faster pace among computer networks: Protocols evolve, compete for survival, and develop a strange mix of tolerance for errors with intolerance for others. Most of the mechanics have not been studied mathematically and indeed little is known as to what the evolution process is trying to achieve, and what the steady state might look like, if at all one exists. Understanding aspects of language and its evolution definitely seem to be worthy causes.

One aspect in particular that we have not explored is the impact of “computational efficiency” of the encoding or decoding procedures. One of the reasons to set aside this concern for the time being is that ingredients like the dictionary suggest that natural language seems not to pay serious attention to the complexity of encoding/decoding relying instead on table lookup for much of its performance; and tables do not appear to be particularly compact. Nevertheless, efficiency perhaps does play a significant role in the evolution of languages since some changes are more easy for the humans to adapt to, as opposed to others. Understanding this aspect of efficiency is probably another challenge for the future.

**Organization of this paper.** We start with the analysis of the chromatic number in Section 2. We then use the methods to build uncertain compression schemes in Section 3.

## 2 Uncertainty Graphs

We start with some elementary material in Section 2.1 that already allows us to prove Parts (1) and (2) of Lemma 1.7. The lower bound mentioned in Part (4) of Lemma 1.7 follows also relatively easily from a result of Linial [9] and we show this in Section 2.2. Our main contribution, in Section 2.3, gives the upper bound from Part (3) of Lemma 1.7.

### 2.1 Preliminaries

We recall the concept of a homomorphism of graphs: For graph  $G = (V, E)$  and  $G' = (V', E')$ , we say that  $\phi: V \rightarrow V'$  is a homomorphism from  $G$  to  $G'$  if  $(u, v) \in E \Rightarrow (\phi(u), \phi(v)) \in E'$ . We say  $G$  is homomorphic to  $G'$  if there exists a homomorphism from  $G$  to  $G'$ .

**Proposition 2.1.** *For every  $N, \ell \geq 1$  and  $k' \leq k \leq N$ , the  $k$ -restricted uncertainty graph  $\mathcal{U}_{N,\ell,k}$  is homomorphic to the  $k'$ -restricted uncertainty graph  $\mathcal{U}_{N,\ell,k'}$ .*

*Proof.* We construct the homomorphism  $\phi$  from  $\mathcal{U}_{N,\ell,k}$  to  $\mathcal{U}_{N,\ell,k'}$  as follows: For  $\pi = \langle \pi(1), \dots, \pi(k) \rangle \in S_{N,k}$  let  $\phi(\pi) = \langle \pi(1), \dots, \pi(k') \rangle \in S_{N,k'}$ . From the definitions it follows that this is a homomorphism.  $\square$

**Proposition 2.2.** *For every  $G$  and  $G'$  such that  $G$  is homomorphic to  $G'$ , we have  $\chi(G) \leq \chi(G')$ .*

*Proof.* Follows from the composability of homomorphisms and the fact that  $G$  is  $k$ -colorable if and only if it is homomorphic to  $K_k$ , the complete graph on  $k$  vertices.  $\square$

Part (1) of Lemma 1.7 follows immediately from Propositions 2.1 and 2.2.

**Proposition 2.3.** *For every  $N, \ell$ , and  $k \geq \ell + 1$  the fractional chromatic number of the restricted uncertainty graph  $\mathcal{U}_{N,\ell,k}$  is at most  $4\ell$ .*

*Proof.* For every function  $f: [N] \rightarrow [2\ell]$  we associate the set  $I_f = \{\pi \in S_{N,k} \mid f(\pi(1)) = 1 \text{ and } f(\pi(j)) \neq 1 \forall j \in \{2, \dots, \ell + 1\}\}$ .

We claim that  $I_f$  is an independent set of  $\mathcal{U}_{N,\ell,k}$  for every  $f$ . To see this consider an edge  $(\pi, \sigma)$  and suppose  $\pi \in I_f$ . Then  $\sigma(1) \in \{\pi(2), \dots, \pi(\ell + 1)\}$  and so  $f(\sigma(1)) \neq 1$  and so  $\sigma \notin I_f$ .

Next we note that for every  $\pi$ , the probability that  $\pi \in I_f$  for  $f$  chosen uniformly at random is  $1/(2\ell) \cdot (1 - 1/(2\ell))^\ell \geq 1/(4\ell)$ .

Thus if we give each  $I_f$  a weight of  $4\ell/(2\ell)^N$ , then we have that the weight of independent sets containing any given vertex  $\pi$  is at least one, while the sum of all weights is  $4\ell$ , thus yielding the claimed bound on the fractional chromatic number.  $\square$

The following is a well-known connection between fractional chromatic number and chromatic number.

**Proposition 2.4.** *For every graph  $G$ ,  $\chi(G) \leq \chi_f(G) \cdot \ln |V(G)|$ .*

We are now ready to prove part (2) of Lemma 1.7.

**Lemma 2.5.**  $\chi(\mathcal{U}_{N,\ell}) \leq \chi(\mathcal{U}_{N,\ell,\ell+1}) \leq 4\ell(\ell + 1) \ln N$

*Proof.* The first inequality follows from Propositions 2.1 and 2.2. The second one follows from Proposition 2.4 and 2.3 and the fact that  $\mathcal{U}_{N,\ell,\ell+1}$  has at most  $N^{\ell+1}$  vertices.  $\square$

## 2.2 Lower Bound on Chromatic Number

We now prove Part (4) of Lemma 1.7 giving a lower bound on  $\chi(\mathcal{U}_{N,\ell,k})$ . We use a lower bound on a somewhat related family of graphs due to Linial [9].

**Definition 2.6** (Shift graphs). *For integers  $N$  and  $k < N$ , we say that  $\pi \in S_{N,k}$  is a left shift of  $\sigma \in S_{N,k}$  if  $\pi(i) = \sigma(i + 1)$  for  $i \in [k - 1]$  and  $\pi(k) \neq \sigma(1)$ . We say  $\pi$  is a right shift of  $\sigma$  if  $\sigma$  is a left shift of  $\pi$ , and we say  $\pi$  is a shift of  $\sigma$  if  $\pi$  is a left shift or a right shift of  $\sigma$ . For integers  $N$  and  $k$ , the shift graph  $\mathcal{S}_{N,k}$  is given by  $V(\mathcal{S}_{N,k}) = S_{N,k}$  with  $(\pi, \sigma) \in E(\mathcal{S}_{N,k})$  if  $\pi$  is a shift of  $\sigma$ .*

**Theorem 2.7** (Linial [9, Proof of Theorem 2.1]). *For every odd  $k$ ,  $\chi(\mathcal{S}_{N,k}) \geq \log^{(k-1)} N$ .*

(We note that the notation in [9] is somewhat different: The graph  $\mathcal{S}_{N,k}$  is denoted  $B_{N,t}$  for  $t = (k - 1)/2$  in [9].)

We show that the uncertainty graphs contain a subgraph isomorphic to the shift graph. This gives us our lower bound on the chromatic number of uncertainty graphs.

**Lemma 2.8.** *For every  $N$ ,  $\ell$  and  $k$  that is an integral multiple of  $\ell$ , we have  $\chi(\mathcal{U}_{N,\ell,k}) \geq (\log^{(2k/\ell)}(N/\ell))$ .*

*Proof.* First without loss of generality we only consider the case of even  $\ell$ . Then we reduce to the case  $\ell = 2$ , by considering only those permutations  $\pi$  which fix  $\pi(i) = i$  if  $\ell/2$  does not divide  $i$ . This still leaves us with  $2N/\ell$  unfixed elements and subpermutations from  $S_{2N/\ell, 2k/\ell}$  that are within distance 2 of each other are within distance  $\ell$  when mapped back to  $S_{N,k}$ .

So we assume  $\ell = 2$  and show that  $\mathcal{U}_{N,2,k}$  contains a subgraph isomorphic to the shift graph  $\mathcal{S}_{N,k}$ . Consider the map  $\phi$  from  $V(\mathcal{S}_{N,k})$  to  $V(\mathcal{U}_{N,2,k})$  which send  $\pi = \langle \pi(1), \dots, \pi(k) \rangle$  to  $\phi(\pi) = \sigma = \langle \sigma(1), \dots, \sigma(k) \rangle$  as follows: Let  $t = \lfloor k/2 \rfloor$ . Then  $\sigma(2i) = \pi(t+i)$  and  $\sigma(2i+1) = \pi(t-i)$   $\sigma(t+i) = \pi(2i)$  and  $\sigma(t-i) = \pi(2i+1)$ . It is easy to verify that the map is a bijection and if  $\pi$  and  $\pi'$  are shifts of each other, then  $\phi(\pi)$  and  $\phi(\pi')$  are within distance 2 of each other. It follows that  $\mathcal{U}_{N,2,k}$  contains a copy of  $\mathcal{S}_{N,k}$  and so  $\chi(\mathcal{U}_{N,2,k}) \geq \chi(\mathcal{S}_{N,k}) \geq \log^{(k-1)} N$ .  $\square$

### 2.3 Upper Bound on Chromatic Number

In this section we give an upper bound on the chromatic number of the uncertainty graphs. We first describe our strategy. Fix  $N$  and  $\ell$ . Now for every  $k$ , we know that there is a homomorphism from  $\mathcal{U}_{N,\ell,k}$  to  $\mathcal{U}_{N,\ell,k-1}$ . However we note that if we jump from  $\mathcal{U}_{N,\ell,k}$  to  $\mathcal{U}_{N,\ell,k-\ell}$  then the homomorphism has an even nicer property. To describe this property, we introduce a new parameter associated with the homomorphism from  $\mathcal{U}_{N,\ell,k}$  to  $\mathcal{U}_{N,\ell,k-\ell}$ . Let us denote this homomorphism  $\phi_k$ . For  $\pi \in \mathcal{S}_{N,k}$  let  $d_k(\pi) = |\{\phi_k(\sigma) \mid (\pi, \sigma) \in E(\mathcal{U}_{N,\ell,k})\}|$ . Note that  $d_k(\pi)$  is independent of  $\pi$  and so we just denote it  $d_k$ . We note first that  $d_k$  is small.

Recall that  $\phi_k: \mathcal{S}_{N,k} \rightarrow \mathcal{S}_{N,k-\ell}$  and maps  $\pi: [k] \rightarrow [N]$  to  $\pi': [k-\ell] \rightarrow [N]$  by setting  $\pi'(i) = \pi(i)$ .

**Claim 2.9.** *For every  $k$ ,  $d_k \leq (2\ell + 1)^k$ .*

*Proof.* Let  $(\sigma, \pi) \in E(\mathcal{U}_{N,\ell,k})$  then  $\delta(\sigma, \pi) \leq \ell$ . In particular for every  $i \in [k-\ell]$ , we have that there exists  $j(i) \in \{-\ell, \dots, \ell\}$  such that  $\sigma(i) = \pi(i + j(i))$ . Thus the sequence  $j(1), \dots, j(k-\ell)$  completely specifies  $\phi_k(\sigma)$ . Since the number of such sequences is at most  $(2\ell + 1)^{k-\ell}$ , we get our claim.  $\square$

To show that a homomorphism with a small  $d$ -value yields especially good colorings we need the notion of an isolating hash family, that defined as follows. For positive integers  $\ell$ ,  $N$  and  $m \in [N]$  and  $S \subseteq [N] - \{m\}$ , we say that a function  $h: [N] \rightarrow \{0, 1\}^\ell$  isolates  $m$  from  $S$  if  $h(m) \notin \{h(m') \mid m' \in S\}$ . We say that a hash family  $\mathcal{H}_\ell = \{h_{1,\ell}, \dots, h_{M,\ell}\}$  is  $(N, \ell)$ -isolating if for every  $S \subseteq [N]$  with  $|S| \leq 2^{\ell-1}$ , and for every  $m \in [N] - S$ , there exists  $j = j(m, S)$  such that  $h_{j,\ell}(m) \notin h_{j,\ell}(S) \triangleq \{h_{j,\ell}(m') \mid m' \in S\}$ .

We note first that small isolating families exist and then give a coloring scheme based on small isolating families.

**Lemma 2.10.** *For every  $\ell$  and  $N$ , there exists an  $(N, \ell)$ -isolating family of size at most  $2^\ell \cdot \log N + 1$ .*

*Proof.* The proof is straightforward application of the probabilistic method. We pick  $\mathcal{H} = \{h_1, \dots, h_M\}$  by picking  $h_i$  uniformly and independently from the set of all functions from  $[N]$  to  $\{0, 1\}^\ell$ . Fix  $m \notin S \subseteq [N]$ . The probability that a randomly chosen  $h$  isolates  $m$  from  $S$  is at least  $1/2$ . Thus the probability that some  $h_i$  in  $\mathcal{H}$  does not isolate  $m$  from  $S$  is at most  $2^{-M}$ . Taking the union bound over all  $m, S$  we find that the probability that  $\mathcal{H}$  does not isolate some  $m$  from  $S$  is at most  $N^{2^\ell} / 2^M$ . We conclude that  $M \leq 2^\ell \cdot \log N + 1$  suffices for the existence of such a  $\mathcal{H}$ .  $\square$

We now ready to show that a homomorphism with a small  $d$ -value yields good colorings.

**Lemma 2.11.** *Let  $\phi$  be a homomorphism from  $G$  to  $H$  and let  $c = \chi(H)$  and  $d = \max_{v \in V(G)} |\{\phi(w) \mid (v, w) \in E(G)\}|$ . Then  $\chi(G) \leq 4d(4d \log c + 1) = O(d^2 \log c)$ .*

*Proof.* Let  $\mathcal{H}$  be a  $(c, \lceil \log d \rceil + 1)$ -isolating family of size  $M \leq 4d \log c + 1$ , which exists by Lemma 2.10.

We claim there is a coloring of  $G$  with  $4d \cdot M$  colors. To get such a coloring, let  $\chi'$  be a coloring of  $H$  with colors  $[c]$ . Now, consider  $v \in V(G)$  and let  $S_v = \{\chi'(\phi(w)) \mid (w, v) \in E(G)\}$ . By the definition of  $d$ , we have  $|S_v| \leq d$ . Also since  $\chi'$  is a coloring of  $H$  and  $\phi$  is a homomorphism, we have  $\chi'(\phi(v)) \notin S_v$ . Thus by the property of  $\mathcal{H}$ , we have that there exists a  $j = j(v)$  such that  $h_j(\chi'(\phi(v))) \notin \{h_j(i') \mid i' \in S_v\}$ . We let the coloring  $\chi$  of  $G$  be  $\chi(v) = (j(v), h_{j(v)}(\chi'(\phi(v))))$ . Syntactically it is clear that this is a  $2^{\lceil \log d \rceil + 1} \cdot M$  coloring of  $G$ , as required. To see it is valid, consider  $(v, w) \in E(G)$ . If  $j(v) \neq j(w)$  then we are done. Else, suppose  $j(v) = j(w) = j$ . Then by definition of  $S_v$  we have  $\chi'(\phi(w)) \in S_v$  and so  $h_j(\chi'(\phi(w))) \in \{h_j(i) \mid i \in S_v\}$ , and thus  $\chi(v) \neq \chi(w)$  as desired.  $\square$

We are now ready to prove Part (3) of Lemma 1.7, restated below.

**Lemma 2.12.** *There exists a constant  $c$  such that for every  $N, \ell, k$  such that  $k < \ell \log^* N$ , we have  $\chi(\mathcal{U}_{N, \ell, k}) \leq 2^{ck \log \ell} \log^{\lfloor (k-1)/\ell \rfloor - 1} N$ .*

*Proof.* We prove the lemma by induction on  $k$ . For notational simplicity assume  $k - 1$  is a multiple of  $\ell$ . For  $k \leq \ell$  the lemma is immediate from the fact that  $\chi(\mathcal{U}_{N, \ell, 1}) \leq N$ . Assume the lemma is true for  $k - \ell$ . Then, by Lemma 2.11 we have that for  $\chi(\mathcal{U}_{N, \ell, k}) \leq 2d_k(d_k + 1) \cdot \log(\chi(\mathcal{U}_{N, \ell, k-\ell})) \leq 4d_k^2 \log \chi(\mathcal{U}_{N, \ell, k-\ell})$ . By Claim 2.9,  $d_k \leq (2\ell + 1)^k \leq (4\ell)^k$  and so  $\chi(\mathcal{U}_{N, \ell, k}) \leq 4(4\ell)^{2k} \log(2^{c(k-\ell)} \log^{\ell} \log^{(k-\ell)/\ell} N) \leq 2^{ck \log \ell} \log^{(k-1)/\ell} N$  for a suitably large  $c$ .  $\square$

### 3 Uncertain Communication

We now convert some of the methods from Section 2.3 into schemes for uncertain compression. In Section 3.1 we derive a simple compression scheme based on the isolating hash families. We then use the “nested series of homomorphisms” to derive a second compression scheme in Section 3.2. The compression scheme of Section 3.2 can make errors with positive probability and has a non-linear dependence on entropy. In Section 3.3 we show that for some natural distributions, this scheme is error-free. In Section 3.4 we show how an error-free scheme working for all distributions would automatically have linear dependence on the entropy, suggesting some of the weaknesses in Section 3.2 are necessary.

#### 3.1 A simple, zero-error compression scheme

Our first construction uses the notion of an isolating hash family, as defined in Section 2.3. In this section we assume that for each  $\ell < N$ , both the encoder and decoder agree on a predetermined  $(N, \ell)$ -isolating family  $\{h_{j, \ell}\}_{j=1}^M$  of size  $M \leq 2^\ell \log N + 1$ , that its existence guaranteed by Lemma 2.10.

**Encoding:** Given  $m, P$  let  $S = \{m' \in [N] \setminus \{m\} \mid P(m') \geq P(m)/2^{2\Delta}\}$  and let  $\ell = \lceil \log 1/P(m) + 2\Delta \rceil$ . Let  $j \in [M]$  be such that  $h_{j, \ell}(m) \notin \{h_{j, \ell}(m') \mid m' \in S\}$ . The encoding  $E(P, m)$  is defined to be  $(j, h_{j, \ell}(m))$ .

**Decoding:** Given  $Q$  and  $y = (j, z) \in \mathbb{Z}^+ \times \{0, 1\}^*$ , let  $\ell = |z|$  and let  $\hat{m} = \operatorname{argmax}_{m \in [N]: h_{j, \ell}(m) = z} \{Q(m)\}$ . The decoding of the pair  $(Q, y)$  is given by  $D(Q, y) = \hat{m}$ .

Our next proposition verifies the correctness of the compression scheme.

**Proposition 3.1.** *For every pair of distributions  $P, Q$  such that  $\delta(P, Q) \leq \Delta$ , and for every message  $m \in [N]$ , it is the case that  $D(Q, E(P, m)) = m$ .*

*Proof.* Fix  $P, Q$  and  $m$  such that  $\delta(P, Q) \leq \Delta$ . Let  $E(m, P) = (j, z)$  with  $\ell = |z|$  and let  $D((j, z), Q) = \hat{m}$ . We will show that  $\hat{m} = m$ . By definition of  $E$ , we have  $h_{j, \ell}(m) = z$  and by definition of  $D$  we have  $h_{j, \ell}(\hat{m}) = z$ . Thus, by the condition that  $\hat{m}$  maximizes probability under  $Q$  of messages satisfying  $h_{j, \ell}(m') = z$ , we have  $Q(\hat{m}) \geq Q(m)$ . Since the distance of  $P$  and  $Q$  is at most  $\Delta$ , we have  $P(m) \leq Q(m)2^{2\Delta}$  and  $P(\hat{m}) \geq Q(\hat{m})/2^{2\Delta}$ . Combining the inequalities we get  $P(\hat{m}) \geq P(m)/2^{2\Delta}$ . Now let  $S = \{m' \in [N] - \{m\} \mid P(m') \geq P(m)/2^{2\Delta}\}$ . We have  $\hat{m} \in S \cup \{m\}$ . But by definition of  $j$ , we have  $h_{j, \ell}(m) \notin \{h_{j, \ell}(m') \mid m' \in S\}$  and since  $h_{j, \ell}(m) = h_{j, \ell}(\hat{m})$ , we must have  $m = \hat{m}$ .  $\square$

Finally we analyze the performance of our scheme.

**Lemma 3.2.** *The expected length of the encoding  $E$  is  $O(H(P) + \Delta + \log \log N)$ .*

*Proof.* Fix  $m \in S$ . Then we have  $\ell \leq 1 + \log 1/P(m) + 2\Delta$  and  $M \leq 2^\ell \log N + 1$ . Thus, the length of  $E(P, m)$  is at most  $2\ell + \log \log N = O(\log 1/P(m) + \Delta + \log \log N)$ . Taking expectation over  $m$  drawn from  $P$ , we have the expected length of the encoding is at most  $O(H(P) + \Delta + \log \log N)$ .  $\square$

Theorem 1.3 follows immediately from Proposition 3.1 and Lemma 3.2.

## 3.2 Compression with error in the low entropy setting

Our compression for the low entropy setting (with better dependence on  $N$ ) relies on an extension of our coloring scheme for the uncertainty graphs. We describe this extension in the next section and then use that to present our compression scheme afterwards.

### 3.2.1 Compression for chains

We start with some terminology. We say that a finite sequence of sets  $A_0, \dots, A_k$  with  $A_i \subseteq [N]$  is a *chain* in  $[N]$  if  $|A_0| = 1$  and  $A_i \subseteq A_{i+1}$  for every  $i$ . We say that  $w$  is the *leader* of the chain if  $A_0 = \{w\}$ . We use  $\text{Chain}(N)$  to denote the set of all chains in  $[N]$ .

In this section we will show how to compress the leader of a chain so that it is unambiguous relative to “nearby” chains. This is in the spirit of the coloring of uncertainty graphs. Indeed vertices of the uncertainty graph  $\mathcal{U}_{N, \ell, k}$  correspond to chains with the vertex  $\langle \pi(1), \dots, \pi(k) \rangle$  corresponding to the chain  $\mathcal{A}$  with  $A_0 = \{\pi(1)\}$  and  $A_i = \{\pi(1), \dots, \pi(\ell \cdot i)\}$  for  $i \geq 1$ . The compressing scheme will thus be similar to the coloring scheme, however there are two distinguishing factors: We will want to compress some chains more than others - a notion that would correspond to asking some vertices to use small colors while allowing others to use larger ones. Furthermore our chains will now grow arbitrarily fast (and not just in steps of 1 or more generally  $\ell$ ). We now describe the precise problem.

For a chain  $\mathcal{A} = \langle A_0, \dots, A_k \rangle$  we say the length of the chain, denoted  $\text{lgt}(\mathcal{A})$ , is the parameter  $k$ . We use  $\text{sz}(\mathcal{A})$  denote the size of the final set  $|A_k|$ . For a chain  $\mathcal{A}$  of length at least  $i$ , we let  $\mathcal{A}_i$  denote its prefix of length  $i$ , i.e.,  $\mathcal{A}_i = \langle A_0, \dots, A_i \rangle$ .

For chain  $\mathcal{A} = \langle A_0, \dots, A_k \rangle$  and chain  $\mathcal{B} = \langle B_0, \dots, B_{k-d} \rangle$ , we say  $\mathcal{B}$  is within distance  $d$  from  $\mathcal{A}$  if for all  $i \in \{0, \dots, k-d\}$ ,  $A_{i-d} \subseteq B_i \subseteq A_{i+d}$  (where we consider sets with negative index to be the empty set). We denote the set of all chains that are within  $d$  distance from  $\mathcal{A}$  by  $S^d(\mathcal{A})$ . Our

goal next is to compress the leader of chains so that the length of the compression is small as a function of  $\text{sz}(\mathcal{A})$ , while it remains unambiguous to chains that are nearby.

**Lemma 3.3.** *There exists a coloring scheme  $\text{Col}: \mathbb{Z}^+ \times \text{Chain}(N) \rightarrow \mathbb{Z}^+$  with the following properties:*

1. *If  $\text{lgt}(\mathcal{A}) \geq 2k$ , then for every  $s \geq \text{sz}(\mathcal{A}_{2k})$ ,  $\text{Col}(s, \mathcal{A}_{2k}) \leq 2^{6(s+1)} \log^{(k)} N$ .*
2. *Let  $\mathcal{A}$  and  $\mathcal{A}'$  be chains of the same length, with  $\text{lgt}(\mathcal{A}) \geq 2k$  and of size at most  $s$ . Then, if  $S^1(\mathcal{A}) \cap S^1(\mathcal{A}') \neq \emptyset$  and  $A_0 \neq A'_0$ , then  $\text{Col}(s, \mathcal{A}_{2k}) \neq \text{Col}(s, \mathcal{A}'_{2k})$ .*

*Proof.* Let  $c_{k,s} = 2^{6(s+1)} \log^{(k)} N$ . Fix  $s \geq \text{sz}(\mathcal{A})$ . We now describe a coloring scheme of a chain  $\mathcal{A}_{2k}$  with  $c_{k,s}$  colors, using induction on  $k$ .

For the base case  $k = 0$ , Let  $w$  be the leader of  $\mathcal{A}$ . Then  $\mathcal{A}_0$  gets the color  $\text{Col}(s, \mathcal{A}_0) = w$ , so clearly  $\text{Col}(s, \mathcal{A}_0) \leq N = \log^{(0)} N$ .

For  $k \geq 1$ , let  $\ell = 2.5s$  and let  $\mathcal{H}$  be an  $(\ell, c_{k-1,s})$ -isolating family (where isolating families were defined as in Section 3.1). By Lemma 2.10 such a family of size  $M = 2^\ell \log c_{k-1,s}$  exists, so let  $\mathcal{H} = \{h_i\}_{i=1}^M$  be one. Let  $T = \{\mathcal{B} \mid \text{lgt}(\mathcal{B}) = 2k - 2, \mathcal{B} \in S^2(\mathcal{A}_{2k}), \text{Col}(s, \mathcal{B}) \neq \text{Col}(s, \mathcal{A}_{2k-2})\}$ . Let  $j \in [2^\ell \log c_{k-1,s}]$  be such that  $h_j(\text{Col}(s, \mathcal{A}_{2k-2})) \neq h_j(\text{Col}(s, \mathcal{B}))$  for all  $\mathcal{B} \in T$ . With these definitions in place, we define  $\text{Col}(s, \mathcal{A}_{2k})$  to be  $(j, h_j(\text{Col}(s, \mathcal{A}_{2k-2})))$ . We verify below that this is a “small” coloring and a valid one.

Let us identify the set  $[2^\ell \log c_{k-1,s}] \times \{0, 1\}^\ell$  with  $[2^{2\ell} \log c_{k-1,s}]$ . The bound on  $c_{k,s}$  follows from the fact that

$$\begin{aligned} & 2^{2\ell} \log c_{k-1,s} \\ & \leq 2^{5s} \log \left( 2^{6(s+1)} \log^{(k-1)} N \right) \\ & \leq 2^{5s} \left( 6(s+1) + \log^{(k)} N \right) \\ & \leq 2^{6(s+1)} \log^{(k)} N, \end{aligned}$$

where the final inequality follows from the fact that  $2^s \cdot 2^6 \geq 6(s+1)$  which is true for every  $s \geq 0$ .

We now verify that the coloring satisfies the requirement in Part (2) of the lemma statement, i.e., that for chains  $\mathcal{A}$  and  $\mathcal{A}'$  of the same length and size at most  $s$ , if their prefixes have the same colors, then they have the same leader. Again we proceed by induction on  $k$ . Assume  $\text{Col}(s, \mathcal{A}_{2k}) = \text{Col}(s, \mathcal{A}'_{2k})$ .

For  $k = 0$ , by assumption we have  $\text{Col}(s, \mathcal{A}_0) = \text{Col}(s, \mathcal{A}'_0)$ . But by definition  $\text{Col}(s, \mathcal{A}_0) = w$  where  $w$  is the leader of  $\mathcal{A}_0$ . It follows thus that  $w$  is also the leader of  $\mathcal{A}'_0$  as claimed.

Now consider  $k \geq 1$ . Let  $\text{Col}(s, \mathcal{A}_{2k}) = (j, h_j(\text{Col}(s, \mathcal{A}_{2k-2})))$  and  $\text{Col}(s, \mathcal{A}'_{2k}) = (j', h_{j'}(\text{Col}(s, \mathcal{A}'_{2k-2})))$ . Since  $\text{Col}(s, \mathcal{A}_{2k}) = \text{Col}(s, \mathcal{A}'_{2k})$ , we have  $j = j'$ . Moreover,  $h_j(\text{Col}(s, \mathcal{A}_{2k-2})) = h_{j'}(\text{Col}(s, \mathcal{A}'_{2k-2}))$ .

We now show that  $\mathcal{A}'_{2k-2} \in S^2(\mathcal{A}_{2k})$ . Let  $\mathcal{B} \in S^1(\mathcal{A}) \cap S^1(\mathcal{A}')$  and consider its prefix  $\langle B_0, \dots, B_{2k-1} \rangle$ . So, for every  $i \in \{0, \dots, 2k-1\}$

$$A_{i-1} \subseteq B_i \subseteq A_{i+1} \text{ and } A'_{i-1} \subseteq B_i \subseteq A'_{i+1}.$$

In other words, for all  $i \in \{0, \dots, 2k-2\}$

$$A_{i-2} \subseteq B_{i-1} \subseteq A'_i \subseteq B_{i+1} \subseteq A_{i+2},$$

Hence  $\mathcal{A}'_{2k-2} \in S^2(\mathcal{A}_{2k})$ .

From our choice of  $j$ ,  $h_j(\text{Col}(s, \mathcal{A}_{2k-2})) = h_j(\text{Col}(s, \mathcal{A}'_{2k-2}))$  for  $\mathcal{A}'_{2k-2} \in S^2(\mathcal{A}_{2k})$  only if  $\text{Col}(s, \mathcal{A}'_{2k-2}) = \text{Col}(s, \mathcal{A}_{2k-2})$ . For conclusion,  $\mathcal{A}_{2k-2}$  and  $\mathcal{A}'_{2k-2}$  are both chains of size at most  $s$  of the same length, and have the same color. From the induction hypothesis they have the same leader.  $\square$

### 3.2.2 The Compression Scheme

We are now ready to define our final compression scheme.

**Encoding:** Given  $m, P$  define  $r = \lfloor -\log P(m) \rfloor$  and  $f = 2 \lceil \log^* N \rceil - 1$ . Further define the chain  $\mathcal{A}$  of length  $f$  as follows.  $A_0 = \{m\}$  and  $A_k = \{m' \in [N] \mid |\log 1/P(m') - r| \leq \Delta(k+1) + 1\}$  (so that  $A_k$  is the set of messages of probability roughly  $P(m)$  with the difference in logarithms being at most  $(k+1)\Delta + 1$ ). Let  $s = \text{sz}(\mathcal{A})$ . The encoding  $E_{\text{low}}(P, m) = E(P, m)$  is

$$E(P, m) = \begin{cases} (s, r, \text{Col}(s, \mathcal{A})) & \text{if } s \leq 2^{\frac{H(P)}{\epsilon} + 2\Delta \log^* N + 1} \\ \perp & \text{otherwise.} \end{cases}$$

(We assume that  $s$  and  $r$  above are encoded in some prefix-free encoding, so that the receiver can separate the three parts.)

**Decoding:** The decoding function  $D_{\text{low}}(Q, y) = D(Q, y)$  works as follows: If  $y = \perp$  then the decoder outputs  $\perp$ . Else let  $y = (s, r, c)$  and let  $f = 2 \lceil \log^* N \rceil - 1$ . Let  $\mathcal{B} = \langle B_0, \dots, B_{f-1} \rangle$  be as follows:  $B_0 = \{w\}$  for some  $w$  such that  $|\log 1/Q(w) - r| \leq \Delta + 1$ . For  $k \geq 1$ ,  $B_k = \{m' \mid |\log 1/Q(m') - r| \leq (k+1)\Delta + 1\}$ . Find a chain  $\mathcal{A}'$  with the following properties:  $\mathcal{B} \in S^1(\mathcal{A}')$ ,  $\text{lgt}(\mathcal{A}') = f$ ,  $\text{sz}(\mathcal{A}') \leq s$  and  $\text{Col}(s, \mathcal{A}') = c$ . Let  $\hat{m}$  be the leader of  $\mathcal{A}'$ . The decoding  $D(Q, y)$  is set to be  $\hat{m}$ .

We first analyze the correctness of the decoder.

**Lemma 3.4.** *For every pair of distributions  $P, Q$  such that  $\delta(P, Q) \leq \Delta$  and for every message  $m \in [N]$  such that  $E_{\text{low}}(P, m) \neq \perp$ , it holds that  $D_{\text{low}}(Q, E_{\text{low}}(P, m)) = m$ .*

*Proof.* Fix  $P \in \mathcal{P}([N])$  and a message  $m \in [N]$  such that  $E_{\text{low}}(P, m) \neq \perp$ . The following claims will show that the decoding process is well defined (and then correctness will be essentially be immediate).

**Claim 3.5.** *There exists  $w \in [N]$  such that  $|\log 1/Q(w) - r| \leq \Delta + 1$ .*

*Proof.* By our choice of  $r$ , we have  $|\log 1/P(m) - r| \leq 1$ . Now using  $\delta(P, Q) \leq \Delta$ , we have  $|\log 1/P(m) - \log 1/Q(m)| \leq \Delta$ , and so  $|\log 1/Q(m) - r| \leq \Delta + 1$ . So  $w = m$  gives an element in  $[N]$  with the desired property.  $\square$

Thus the chain  $\mathcal{B}$  is now well-defined. It remains to show that there exists a chain  $\mathcal{A}'$  satisfying the required properties. The next claim shows that  $\mathcal{B} \in S^1(\mathcal{A}')$ , therefore  $\mathcal{A}'$  is a candidate for the role of  $\mathcal{A}'$ .

**Claim 3.6.**  $\mathcal{B} \in S^1(\mathcal{A}')$ .

*Proof.* The proof follows easily from our choice of  $\mathcal{A}, \mathcal{B}$  and the fact that  $P$  and  $Q$  are  $\Delta$ -close. Let  $k \in \{0, \dots, f-1\}$ . We need to show that  $B_k$  is sandwiched between  $A_{k-1}$  and  $A_{k+1}$ .

First, We will show that  $B_k \subseteq A_{k+1}$ . When  $k = 0$ , we need to show that  $w \in A_1$ . Indeed,

$$\begin{aligned} & |\log 1/Q(w) - r| \leq \Delta + 1 \\ \Rightarrow & |\log 1/P(w) - r| \leq 2\Delta + 1 \\ \Rightarrow & w \in A_1 . \end{aligned}$$

Now consider  $1 \leq k \leq f - 1$ . We have,

$$\begin{aligned} B_k &= \{m' \in [N] \mid |\log 1/Q(m') - r| \leq (k + 1)\Delta + 1\} \\ &\subseteq \{m' \in [N] \mid |\log 1/P(m') - r| \leq (k + 2)\Delta + 1\} \\ &= A_{k+1} . \end{aligned}$$

This shows that  $B_k \subseteq A_{k+1}$ . Next we show that  $A_{k-1} \subseteq B_k$ , for  $2 \leq k \leq f - 1$ . We have

$$\begin{aligned} A_{k-1} &= \{m' \in [N] \mid |\log 1/P(m') - r| \leq k\Delta + 1\} \\ &\subseteq \{m' \in [N] \mid |\log 1/Q(m') - r| \leq (k + 1)\Delta + 1\} \\ &= B_k . \end{aligned}$$

The case where  $k = 1$  and  $m \in B_1$  was proved in Claim 3.5. So we are done.  $\square$

To conclude, the decoder can find a chain  $\mathcal{A}'$  such that  $\text{sz}(\mathcal{A}') \leq s$ ,  $\text{lgt}(\mathcal{A}') = \text{lgt}(\mathcal{A})$ ,  $\text{Col}(s, \mathcal{A}') = \text{Col}(s, \mathcal{A})$  and there exists a chain  $\mathcal{B} \in S^1(\mathcal{A}') \cap S^1(\mathcal{A})$ . From Lemma 3.3 the leader of  $\mathcal{A}'$  is  $m$  as required.  $\square$

We are now ready to prove Theorem 1.4.

*Proof.* We now estimate the probability that the encoder will fail. Fix some probability distribution  $P$  and a message  $m$  such that  $E(P, m) = \perp$ . We will first show that  $P(m) \leq 2^{-\frac{H(P)}{\epsilon}}$ . Later, we will bound the probability that “ $m$  has such small probability” by  $\epsilon$ .

Consider the chain  $\mathcal{A} = \langle A_0, \dots, A_f \rangle$  as defined by the encoder. In this case, the size of the largest set,  $|A_f|$ , is more than the threshold  $T = 2^{\frac{H(P)}{\epsilon} + 2\Delta \log^* N + 1}$ . So, there is some element  $m' \in A_f$  such that  $P(m') \leq \frac{1}{T}$ . By our choice of  $A_f$ ,  $P(m') \geq 2^{-\lfloor -\log P(m) \rfloor - (f+1)\Delta - 1} \geq P(m)2^{-2\Delta \log^* N - 1}$ . Calculating,

$$\frac{1}{T} \geq P(m)2^{-2\Delta \log^* N - 1} \Rightarrow P(m) \leq \frac{2^{2\Delta \log^* N + 1}}{T} = 2^{-\frac{H(P)}{\epsilon}} .$$

Therefore, we can bound the failure probability by the probability that  $P(m) \leq 2^{-\frac{H(P)}{\epsilon}}$ . Using the fact that  $\mathbf{E}_{m \leftarrow P[N]} \left[ \log \frac{1}{P(m)} \right] = H(P)$ , we deduce the following by Markov's inequality,

$$\Pr_{m \leftarrow P[N]} \left[ P(m) \leq 2^{-\frac{H(P)}{\epsilon}} \right] = \Pr_{m \leftarrow P[N]} \left[ \log \frac{1}{P(m)} \geq \frac{H(P)}{\epsilon} \right] \leq \epsilon$$

We will finish the proof by bounding the performance of the scheme. To this end consider a distribution  $P$  and a message  $m \in [N]$  such that  $E(P, m) \neq \perp$  (i.e.  $\text{sz}(\mathcal{A}) \leq T$ ). The encoder sends  $r = \lfloor -\log P(m) \rfloor$ ,  $s = \text{sz}(\mathcal{A})$  and  $\text{Col}(s, \mathcal{A})$ . We first analyze the contribution of sending  $r$  to

the performance. Because  $\log |r| = O\left(\log \frac{1}{P(m)}\right)$ , the accepted length of sending  $r$  in a prefix-free encoding is at most  $O\left(\mathbf{E}_{m \leftarrow P} \log \frac{1}{P(m)}\right) = O(H(P))$ .

Now we analyze the length of  $(s, \text{Col}(s, \mathcal{A}))$ . By Lemma 3.3:

$$\text{Col}(s, \mathcal{A}) \leq 2^{6(s+1)} \log^{(f)} N = 2^{O(s)}$$

Hence, the length of  $(s, \text{Col}(s, \mathcal{A}))$  is at most

$$O(\log s) + \log \text{Col}(s, \mathcal{A}) = O(s) = 2^{\frac{H(P)}{\epsilon} + 2\Delta \log^* N + O(1)}.$$

Thus, from the linearity of expectations, it follows that the total performance is at most  $2^{\frac{H(P)}{\epsilon} + 2\Delta \log^* N + O(1)}$ .  $\square$

### 3.3 Error-free Compression for Natural Distributions

In this section we will show that for a large class of natural distributions, the above scheme is error free. We start by describing the natural distributions we can capture.

We say that a distribution  $P \in \mathcal{P}([N])$  is *flat* if there exists a set  $S \subseteq [N]$  such that  $P$  is uniform on  $S$ . The distribution is called *geometric* if there exist parameter  $\alpha \in (0, 1)$  and a permutation  $\pi$  on  $[N]$  such that for all  $k \in [N-1]$  it holds that  $P(\pi(k+1)) = \alpha P(\pi(k))$ . We call  $P$  *binomial* if there exist a parameter  $p \in (0, 1)$  and a permutation  $\pi$  on  $[N]$  such that  $\forall k \in [N]$ ,  $P(\pi(k)) = \binom{N}{k} p^k (1-p)^{n-k}$ . The sets of all flat, geometric and binomial distributions over  $[N]$  are denoted by  $\text{Flat}_N$ ,  $\text{Geo}_N$  and  $\text{Bin}_N$  respectively.

The following theorem shows that the scheme  $(E_{\text{low}}, D_{\text{low}})$  performs well *without error* on all of the above natural distributions. Moreover, this theorem is stable in the sense that the guarantee on the performance holds even if a distribution is only close to one of the above-mentioned natural distributions.

**Theorem 3.7.** *Let  $\mathcal{F} \triangleq \text{Flat}_N \cup \text{Geo}_N \cup \text{Bin}_N$  and  $L(P) \triangleq 2^{H(P)} \lceil \Delta \log^* N \rceil$ . Then the scheme  $(E_{\text{low}}, D_{\text{low}})$  (with  $\epsilon$  set to 0) is a  $(\Delta, 0, \mathcal{F}, O(L(P)))$ -UCS. Moreover, if  $P \in \mathcal{P}([N])$  is  $\Delta \log^* N$ -close to a distribution  $\tilde{P} \in \mathcal{F}$  then the performance of the scheme on  $P$  is  $\mathbf{E}_{m \leftarrow P} [|E(P, m)|] = O(L(\tilde{P}))$ .*

We prove the theorem above by identifying a broad condition on distributions, which we call the *capacity*, and showing that the performance of our scheme is good if the capacity is small. We define this notion next, show that it is small for the distributions under consideration in Lemma 3.8 next, and finally bound the performance as a function of the capacity in Lemma 3.9 afterwards, thus leading to a proof of Theorem 3.7.

Let  $P \in \mathcal{P}([N])$  be a distribution and let  $S \subseteq [N]$  be its support. We say that  $U \subseteq S$  is a *unit set* of  $P$  if for any two elements  $m_1, m_2 \in U$  the distance  $|\log P(m_1) - \log P(m_2)| \leq 1$ . We define the *capacity* of  $P$ , denoted by  $\mathcal{C}(P)$ , to be the minimal  $c \in \mathbb{R}$  such that the size of every unit set of  $P$  is bounded by  $2^c$ .

Later, we will prove the next lemma, showing that for the previously discussed distributions, the capacity is roughly the entropy.

**Lemma 3.8.** *Let  $P \in \text{Flat}_N \cup \text{Geo}_N \cup \text{Bin}_N$ . Then  $\mathcal{C}(P) \leq H(P) + O(1)$ .*

Theorem 3.7 follows immediately from Lemma 3.8 combined with the following lemma.

**Lemma 3.9.** *For every  $P$  ( $E_{\text{low}}, D_{\text{low}}$ ) (with respect to  $\epsilon = 0$ ) is a  $(\Delta, O(\log(H(P)) + 2^{\mathcal{C}(P)} \lceil \Delta \log^* N \rceil))$ -scheme. Moreover, if  $P$  is  $\Delta \log^* N$ -close to a distribution  $\tilde{P}$ , then the performance of the scheme on  $P$  is  $O(\log(H(P)) + 2^{\mathcal{C}(\tilde{P})} \lceil \Delta \log^* N \rceil)$ .*

*Proof.* When setting  $\epsilon = 0$ , the encoder never outputs  $\perp$ . Lemma 3.4 already implies the correctness of the scheme. The only remaining task is to analyze the performance of the scheme.

Recall, the output of the encoder has three components:  $r$ ,  $s$  and  $(f, \mathcal{A})$ . From linearity of expectation it suffices to analyze the expected length of each component separately.

For a given word  $m \in [N]$ , the first component is  $r = \left\lfloor \log \frac{1}{P(m)} \right\rfloor$ . Its length is  $|r| = O(\log \log \frac{1}{P(m)})$ . Using the concavity of the function  $\log$  we can bound the expectation of  $|r|$  as follows:

$$\mathbf{E} \left[ \log \log \frac{1}{P(m)} \right] \leq \log \left( \mathbf{E} \left[ \log \frac{1}{P(m)} \right] \right) = \log(H(P)) .$$

Now consider the chain  $\mathcal{A}$  with size  $s$  and length  $f = \log^*(N) - O(1)$  as defined by the encoder. The second component is the size  $s$ . Clearly,  $|s| = O(\log s)$ .

The third component is  $\text{Col}(s, \mathcal{A})$ . By Lemma 3.3,  $\text{Col}(s, \mathcal{A}) = \exp(s)$  so  $|\text{Col}(s, \mathcal{A})| = O(s)$ .

Hence the expected length of the last two components is bounded by  $O(s)$ . Let  $\tilde{P} \in \mathcal{P}([N])$  be a distribution that is  $\Delta \log^* N$ -close to  $P$ . To achieve the results it is enough to show that the size  $s$  of the chain  $\mathcal{A}$  associated with  $P$  and  $m$  is bounded by  $O(2^{\mathcal{C}(\tilde{P})} \lceil \Delta \log^* N \rceil)$ .

The size  $s = \text{sz}(\mathcal{A})$  is the size of the following set,

$$A = \{m' \in [N] \mid |\log 1/P(m') - r| \leq 2 \lceil \Delta \log^* N \rceil + 1\} .$$

We will show that this set can be covered by  $O(\lceil \Delta \log^* N \rceil)$  unit sets of  $\tilde{P}$ . This will yield an upper bound on  $s$  of  $O(2^{\mathcal{C}(\tilde{P})} \lceil \Delta \log^* N \rceil)$  as required.

Let  $k = 3 \lceil \Delta \log^* N \rceil + 1$ . Define  $U_{-k}, \dots, U_{k-1}$  as

$$U_i = \{m' \mid i \leq r + \log \tilde{P}(m') \leq i + 1\} .$$

Clearly the  $U_i$ s are unit sets of  $\tilde{P}$ . Moreover, their union is the set

$$\bigcup_{i=-k}^{k-1} U_i = \{m' \mid |\log 1/\tilde{P}(m') - r| \leq 3 \lceil \Delta \log^* N \rceil + 1\} .$$

Let  $m' \in A$ . It remains to verify that  $m' \in \bigcup_{i=-k}^{k-1} U_i$ . Indeed,

$$\begin{aligned} \left| \log 1/\tilde{P}(m') - r \right| &\leq \left| 1/\log P(m') - r \right| + \Delta \log^* N \\ &\leq 3 \lceil \Delta \log^* N \rceil + 1 . \end{aligned}$$

Therefore,  $|A| \leq \sum |U_i| = O(2^{\mathcal{C}(\tilde{P})} \lceil \Delta \log^* N \rceil)$  as required.  $\square$

To complete the proof of Theorem 3.7, we will prove Lemma 3.8. The proof follows immediately from the next three claims.

**Claim 3.10.** *Let  $P \in \text{Flat}_N$ . Then  $\mathcal{C}(P) \leq H(P)$ .*

*Proof.* Let  $S \subseteq [N]$  be the support of  $P$ . Clearly,  $H(P) = \log |S|$ . For every  $U \subseteq S$  that is a unit set of  $P$ ,

$$|U| \leq |S| = 2^{H(P)}.$$

Thus,  $\mathcal{C}(P) \leq H(P)$ . □

**Claim 3.11.** *Let  $P \in \text{Geo}_N$ . Then  $\mathcal{C}(P) \leq H(P) + O(1)$ .*

*Proof.* The definitions of capacity and entropy depend only on the set  $\{P(k)\}_{k=1}^N$  (and not on its order) so we can assume, without loss of generality, that  $P(1) \geq P(2) \geq \dots \geq P(N)$ . Let  $\alpha \in (0, 1)$  be such that for all  $k \in [N-1]$ ,  $P(k+1) = \alpha P(k)$ . We will further assume that  $\alpha^N < \frac{1}{2}$ . Otherwise,

$$H(P) \geq \log(N) - 1 \geq \mathcal{C}(P) - 1,$$

and we are done.

Let  $U$  be the maximal unit set of  $P$ , i.e.  $|U| = u = 2^{\mathcal{C}(P)}$ . Let  $k \in U$  be the element with the highest probability in  $U$ . From maximality of  $U$  we can assume that  $U = \{k, k+1, \dots, k+u-1\}$ . Calculating,

$$1 \geq |\log P(k) - \log P(k+u-1)| = (u-1) \log \frac{1}{\alpha}$$

Therefore,  $u = \frac{1}{\log \frac{1}{\alpha}} + 1 = O(\frac{1}{1-\alpha})$ . To achieve the result it is enough to show that  $\frac{1}{1-\alpha} \leq 2^{H(P)+O(1)}$ , i.e.  $H(P) \geq \log \frac{1}{1-\alpha} - O(1)$ . Calculating the entropy, indeed,

$$H(P) = \log \left( \frac{1 - \alpha^N}{1 - \alpha} \right) + \left( \frac{1 - N\alpha^{N-1}}{1 - \alpha^N} \right) \alpha \log \frac{1}{\alpha} + \left( \frac{1 - \alpha^{N-1}}{1 - \alpha^N} \right) \alpha^2 \cdot \frac{\log \frac{1}{\alpha}}{1 - \alpha} \stackrel{(\alpha^N < \frac{1}{2})}{\geq} \log \left( \frac{1}{1 - \alpha} \right) - O(1),$$

as required. □

**Claim 3.12.** *Let  $P \in \text{Bin}_N$ . Then  $\mathcal{C}(P) \leq H(P) + O(1)$ .*

*Proof.* Assume without loss of generality that there is  $p \in (0, 1)$  such that  $P(k) = \binom{N}{k} p^k (1-p)^{n-k}$  (i.e.  $\pi$  from the definition of binomial distribution is the identity permutation). Let  $U$  be a unit set of  $P$  with size  $2^{\mathcal{C}(P)}$ . We will partition the codewords in  $[N]$  into three regions and bound the number of codewords from each region in  $U$ . The regions are:

1.  $\{k \in [N] \mid k > pN + \sqrt{pN}\}$ ,
2.  $\{k \in [N] \mid k < pN - \sqrt{pN} - 1\}$
3. and  $\{k \in [N] \mid pN - \sqrt{pN} - 1 \leq k \leq pN + \sqrt{pN}\}$ .

We will show that in any region, the number of elements from the region in  $U$  is bounded by  $O(\sqrt{pN})$ . This will yield a total bound of  $|U| = O(\sqrt{pN})$ . The entropy of  $P$  is  $H(P) = \frac{1}{2} \log(2\pi e N p(1-p))$ . Therefore  $\sqrt{pN} = 2^{H(P)+O(1)}$  and the result follows.

First we consider elements  $k$  from the first region. Let  $u_1$  be the number of words in  $U$  from this region. In this case

$$\begin{aligned} \frac{P(k+1)}{P(k)} &= \frac{\binom{N}{k+1} p^{k+1} (1-p)^{N-(k+1)}}{\binom{N}{k} p^k (1-p)^{N-k}} = \frac{(N-k)}{(k+1)} \cdot \frac{p}{1-p} \leq \frac{(N-k)}{k} \cdot \frac{p}{1-p} = \left(\frac{N}{k} - 1\right) \cdot \frac{p}{1-p} \\ &\leq \left(\frac{N}{pN + \sqrt{pN}} - 1\right) \cdot \frac{p}{1-p} \leq 1 - \frac{1}{\sqrt{pN} + 1}. \end{aligned}$$

In a similar way to the proof of Claim 3.11, we can conclude that  $u_1$  is bounded by  $O(\sqrt{pN} + 1) = O(\sqrt{pN})$

Now consider element  $k$  in the second region, similarly:

$$\begin{aligned} \frac{P(k+1)}{P(k)} &= \frac{(N-k)}{(k+1)} \cdot \frac{p}{1-p} \geq \frac{N-(k+1)}{k+1} \cdot \frac{p}{1-p} = \left(\frac{N}{k+1} - 1\right) \cdot \frac{p}{1-p} \\ &\geq \left(\frac{N}{pN - \sqrt{pN}} - 1\right) \cdot \frac{p}{1-p} \geq 1 + \frac{1}{\sqrt{pN}}. \end{aligned}$$

Therefore  $u_2$ , the number of elements from the second region in  $U$ , is bounded by  $O(\sqrt{pN})$

Clearly,  $u_3$ , the number of elements from  $U$  in the last region, is bounded by the size of the region. So  $u_3 = O(\sqrt{pn})$ .

Combining the above, we get

$$2^{\mathcal{E}(P)} = |U| = \sum_{i=1}^3 u_i = O(\sqrt{pn}) = 2^{H(P)+O(1)}$$

as required. □

### 3.4 Dependence of communication on entropy

In the previous sections we gave a scheme with performance that is exponential in the entropy. This scheme is error-free for some natural distributions and had positive error for general distributions. The next lemma shows that if we cannot find a scheme with performance that is linear in the entropy, then any scheme that we will find must have positive error for some distributions.

**Lemma 3.13.** *For every non-decreasing function  $L: \mathfrak{R}^+ \rightarrow \mathfrak{R}^+$  there exists a constant  $c = c_L$  such that the following holds: If there exists  $(\Delta, L(H(P)))$ -UCS for some  $\Delta > 0$ , then there exists a  $(\Delta, c \cdot (1 + H(P)))$ -UCS.*

*Proof.* We will prove the lemma for  $c = L(3) + 2$ . Let  $(E, D)$  be the  $(\Delta, L(H(P)))$ -UCS. We will construct a UCS  $(E', D')$  that has the required performance.

For every distribution  $P \in \mathcal{P}([N])$  and real number  $M > 1$ , we introduce a notion of an  $M$ -concentrated version of  $P$ , denoted  $P_M$ , to be:  $P_M(1) = 1 - 1/M + (1/M) \cdot P(1)$  and  $P_M(i) = (1/M) \cdot P(i)$  for  $i > 1$ . So  $P_M$  is mostly focused on a single point and so has small entropy, but it provides enough variability to capture the variation of  $P$ . In what follows, we will apply  $(E, D)$  to the distributions  $P_M$  and  $Q_M$  for an appropriate choice of  $M$ , chosen to reduce the entropy of  $P_M$  to be a constant and this will give the schemes  $E'$  and  $D'$ .

**The new scheme**  $(E', D')$ : On input  $P \in \mathcal{P}([N])$  and  $m \in [N]$ ,  $E'(P, m)$  is computed as follows: Let  $M = \lceil H(P) \rceil$ . Then  $E'(P, m) = (M, E(P_M, m))$ .

On input  $Q$  and received string  $y' = (M, y)$  the decoding is  $D'(Q, y') = D(Q_M, y)$ .

In what follows we argue that this is a valid zero-error UCS for uncertainty parameter  $\Delta$ , with performance  $c \cdot H(P)$ . We start by proving its validity.

**Claim 3.14.** *For every pair  $P, Q \in \mathcal{P}([N])$  such that  $\delta(P, Q) \leq \Delta$ , and for every  $m \in [N]$  we have  $D'(Q, E'(P, m)) = m$ .*

*Proof.* Fix  $M = \lceil H(P) \rceil$ . Since  $E'(P, m) = E(P_M, m)$  and  $D'(Q, (M, y)) = D(Q_M, y)$ , it suffices to prove that  $P_M$  and  $Q_M$  are  $\Delta$ -close, since then we can use the correctness of  $(E, D)$  on  $P_M$  and  $Q_M$  to conclude  $D(Q_M, E(P_M, m)) = m$ . Below we verify that  $P_M$  and  $Q_M$  are  $\Delta$ -close.

First we consider  $m \in [N] \setminus \{1\}$ . For such  $m$  we have  $P_M(m) = \frac{1}{M}P(m)$  and  $Q_M(m) = \frac{1}{M}Q(m)$  and so  $P_M(m)/Q_M(m) = P(m)/Q(m)$ . So  $|\log P_M(m)/Q_M(m)| = |\log P(m)/Q(m)| \leq \Delta$ .

Now, consider  $m = 1$ . In this case  $P_M(m) = \left(\frac{M-1}{M}\right) + \left(\frac{1}{M}\right) \cdot P(1)$  and  $Q_M(m) = \left(\frac{M-1}{M}\right) + \left(\frac{1}{M}\right) \cdot Q(1)$ . Assume  $P(1) \geq Q(1)$  (the other case is similar) and so  $0 \leq \log P(1)/Q(1) \leq \Delta$ . On the one hand we have  $P_M(1) \geq Q_M(1)$  and on the other hand we have  $P_M(1)/Q_M(1) \leq P(1)/Q(1)$  (which holds for every  $M > 0$ ). It follows that  $0 \leq \log P_M(1)/Q_M(1) \leq \log P(1)/Q(1) \leq \Delta$ .

It follows that  $\delta(P_M, Q_M) \leq \Delta$  and the claim follows.  $\square$

It remains to analyze the performance of the scheme.

**Claim 3.15.** *For every distribution  $P \in \mathcal{P}([N])$ , we have  $\mathbf{E} \left[ |E'_{m \sim P}[N](P, m)| \right] \leq c \cdot H(P)$ .*

*Proof.* Recall that the encoding of  $m \in [N]$  is the pair  $(M, E(P_M, m))$  where  $M = \lceil H(P) \rceil$ . It follows that the first part the encoding is always of length at most  $2 \cdot (1 + \log H(P))$  (allowing for prefix free encodings and rounding up of  $H(P)$  to its ceiling). We crudely bound the above by  $2(1 + H(P))$ .

We turn to the length of the second part, i.e.,  $E(P_M, m)$ . We first show that  $\mathbf{E}_{m \sim P[N]} [|E(P_M, m)|] \leq M \cdot \mathbf{E}_{m \sim P_M[N]} [|E(P_M, m)|]$ . We then bound  $\mathbf{E}_{m \sim P_M[N]} [|E(P_M, m)|]$  by  $L(3)$  thus giving us that total expected length of the encoding  $\mathbf{E}_{m \sim P[N]} [|E'(P, m)|] \leq (L(3) + 2) \cdot (1 + H(P)) = c(1 + H(P))$ .

We start by showing the first step. We have

$$\begin{aligned} & \mathbf{E}_{m \sim P_M[N]} [|E(P_M, m)|] \\ &= \frac{1}{M} \mathbf{E}_{m \sim P[N]} [|E(P_M, m)|] + \left(1 - \frac{1}{M}\right) \mathbf{E}_{m \sim P[N]} [|E(P_M, 1)|] \\ &\geq \frac{1}{M} \mathbf{E}_{m \sim P[N]} [|E(P_M, m)|]. \end{aligned}$$

It follows that  $\mathbf{E}_{m \sim P[N]} [|E(P_M, m)|] \leq M \mathbf{E}_{m \sim P_M[N]} [|E(P_M, m)|]$  as asserted.

By the performance of  $E$  on  $P_M$ , we have  $\mathbf{E}_{m \sim P_M[N]} [|E(P_M, m)|] \leq L(H(P_M))$ . So it suffices to show  $H(P_M) \leq 3$ . This is straightforward from the definition of  $P_M$  and the choice of  $M$ . We

have

$$\begin{aligned} H(P_M) &= \sum_{m \in [N]} P_M(m) \log \frac{1}{P_M(m)} \\ &\leq \left(1 - \frac{1}{M}\right) \log \frac{1}{P_M(1)} + \frac{1}{M} \cdot \sum_{m \in [N]} P(m) \log \frac{M}{P(m)} \\ &\leq 1 + \frac{1}{M} \cdot (H(P) + \log M) \\ &\quad \text{(Using } P_M(1) \geq 1/2 \text{ if } M \geq 2 \text{ and } 1 - \frac{1}{M} = 0 \text{ otherwise.)} \\ &\leq 1 + 1 + \frac{\log M}{M} \\ &\leq 3 \end{aligned}$$

as required.

The claim follows and so does the lemma. □

□

## References

- [1] Mark Braverman and Anup Rao. Information equals amortized communication. In Rafail Ostrovsky, editor, *FOCS*, pages 748–757. IEEE, 2011.
- [2] Richard Cole and Uzi Vishkin. Deterministic coin tossing with applications to optimal parallel list ranking. *Information and Control*, 70(1):32–53, 1986.
- [3] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley, 1991.
- [4] Oded Goldreich, Brendan Juba, and Madhu Sudan. A theory of goal-oriented communication. *J. ACM*, 59(2):8, 2012.
- [5] Prahladh Harsha, Rahul Jain, David A. McAllester, and Jaikumar Radhakrishnan. The communication complexity of correlation. *IEEE Transactions on Information Theory*, 56(1):438–449, 2010. Preliminary version in IEEE CCC 2007.
- [6] Brendan Juba, Adam Tauman Kalai, Sanjeev Khanna, and Madhu Sudan. Compression without a common prior: an information-theoretic justification for ambiguity in language. In Bernard Chazelle, editor, *ICS*, pages 79–86. Tsinghua University Press, 2011.
- [7] Brendan Juba and Madhu Sudan. Universal semantic communication I. In *Proceedings of the 2008 ACM International Symposium on Theory of Computing, Victoria, British Columbia, Canada, May 17-20, 2008*, pages 123–132. ACM, 2008.
- [8] Eyal Kushilevitz and Noam Nisan. *Communication Complexity*. Cambridge University Press, 1997.
- [9] Nathan Linial. Locality in distributed graph algorithms. *SIAM J. Comput.*, 21(1):193–201, 1992.

- [10] Claude E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, 27:379–423, 623–656, 1948.
- [11] Jacob Ziv and Abraham Lempel. A universal algorithm for sequential data compression. *IEEE Transactions on Information Theory*, 23(3):337–342, 1977.