

# Compact representations as a search strategy: Compression EDAs

Marc Toussaint

*Institute for Adaptive and Neural Computation,  
University of Edinburgh, 5 Forrest Hill, Edinburgh EH1 2QL, Scotland, UK*

---

## Abstract

The choice of representation crucially determines the capability of search processes to find complex solutions in which many variables interact. The question is how good representations can be found and how they can be adapted online to account for what can be learned about the structure of the problem from previous samples. We address these questions in a scenario that we term indirect Estimation-of-Distribution: We consider a decorrelated search distribution (mutational variability) on a variable length genotype space. A one-to-one encoding onto the phenotype space then needs to induce an adapted phenotypic search distribution incorporating the dependencies between phenotypic variables that have been observed successful previously. Formalizing this in the framework of Estimation-of-Distribution Algorithms, an adapted phenotypic search distribution can be characterized as minimizing the Kullback-Leibler divergence to a population of previously selected samples (parents). The paper derives a relation between this Kullback-Leibler divergence and the description length of the encoding, stating that compact representations provide a way to minimize the divergence. A proposed class of Compression Evolutionary Algorithms and experiments with an grammar-based compression scheme illustrate the new concept.

*Key words:* Estimation-of-Distribution Algorithms, factorial representations, compression, minimal description length, Evolutionary Algorithms, genotype-phenotype mapping.

---

## 1 Introduction

The complexity of a problem largely depends on the interactions between the variables of a solution. A stochastic search process will perform well on a

---

*Email address:* mtoussai@inf.ed.ac.uk (Marc Toussaint).

complex problem only when the search distribution is adapted to these interactions, i.e., when the search distribution obeys these dependencies between variables.

One approach to design a search distribution is to choose a suitable representation. The questions are how to find representations that induce the desired dependencies on the search space, and how they can be adapted online to account for what can be learned about the structure of the problem from previous samples. There have been various approaches to characterize what a good representation is, considering, for example, all possible representations, Gray vs. binary codes, redundant representations, and recursive encodings (Liepins and Vose, 1990; Whitley et al., 1997; Barbulescu et al., 2000; Rothlauf and Goldberg, 2003; Hornby and Pollack, 2001). Theoretical approaches concerning the adaptation of the search distribution based on the specific current population include, for example, Estimation-Of-Distribution Algorithms (Pelikan et al., 1999), Walsh analysis (Heckendorn and Wright, 2004), and Maximal Entropy principles (Wright et al., 2004). Also using Minimum Description Length principles for building search distributions has been discussed before (Harik, 1999).<sup>1</sup>

Our approach is to consider a specific scenario that we term indirect Estimation-of-Distribution: We assume that the search distribution on a variable length coding space (genotype space) is decorrelated. A one-to-one encoding onto the search space (phenotype space) then needs to induce a properly structured phenotypic search distribution. The idea of this scenario is that the encoding receives all responsibilities to induce the structural properties of the phenotypic search distribution, leaving a simple problem of unstructured (decorrelated) adaptation on the genotype level.

We formalize the scenario within the framework of Estimation-of-Distribution Algorithms, where an adapted phenotypic search distribution can be characterized as minimizing the Kullback-Leibler divergence to a population of previously selected samples (parents). Our core result is a relation between the Kullback-Leibler divergence and the description length of the encoding in the specific scenario of indirect Estimation-of-Distribution based on a variable length genotype space. The result states that compact representations provide a way to minimize this divergence—and may thus be seen as transferring similar results on Minimum Description Length in the context of modeling, in particular model selection, (e.g., (Barron et al., 1998; Vitányi and Li, 2000)) to the specific domain of stochastic search based on adaptive representations.

---

<sup>1</sup> In the ECGA proposed by Harik (1999) MDL refers to finding a partitioning of the phenotypic variables that minimizes the sum of entropies in each partition, which is trivially equivalent to maximizing the sum of mutual informations in each partition (i.e., finding a high linkage partitioning). Finding or adapting a concrete mapping to another coding space is not considered.

An intuitive way to grasp these results might be the following: Consider a set of good (selected) samples in phenotype space. In this selected population there will generally exist dependencies between phenotypic variables, measurable as mutual information between them. This is the information that should be extracted and exploited for further search. Assume we can map these samples on variable size strings such that the average string length is minimized. Before the compression there was mutual information between the phenotypic variables. After the compression, there should be no mutual information between the (genetic) symbols that describe the samples because otherwise the mapping would not be a minimum description length compression. Thus, a compression is one way (among others) to map on a representation in which symbols are decoupled (i.e., to map on a factorial representation). A compression can also be considered as an implicit analysis of the dependencies that have been present in the parent population because it is able to dissolve them by introducing new symbols. Eventually, the key idea is that *inverting* the compression is a mechanism to induce exactly these dependencies. For instance, when there is noise (decorrelated mutational variability) on the genetic symbols, this should translate to a phenotypic variability that obeys these dependencies.

The following two sections will introduce the theoretical framework, including the basics of Estimation-of-Distribution Algorithms and the indirect induction of a search distribution. Section 4 derives the main results on the relation between compact representations and Estimation-of-Distribution. Section 5 aims to illustrate the implications by proposing a class of Compression Evolutionary Algorithms and presenting experiments with a Compression EDA and GA based on a simple grammar-based compression scheme (well comparable to the ideas of Nevill-Manning and Witten (1997)).

## 2 Estimation-Of-Distribution

Let  $P$  be the search space. A heuristic search scheme is a stochastic process in the space  $\Lambda$  of distributions over the search space in which a search distribution  $q \in \Lambda$  is propagated iteratively. In each iteration, samples from  $q$  are drawn, evaluated, and the outcome of evaluation is used to design, according to some heuristic, a new search distribution in the next step.

For instance, in ordinary Evolutionary Algorithms (EAs) the search distribution is given by a finite parent population and recombination and mutation operators (leading to a mixture of mutation kernels as search distribution). The search distribution is sampled, leading to the finite offspring population, which is then evaluated, leading to the selection probability distribution over these offspring. The heuristic to generate the new search distribution in sim-

ple EAs is to sample the selection probability distribution, leading to a new parent population which in turn induces a new search distribution. We do not need to specify these operators here explicitly. We develop the theory on the abstract level of search distributions.

What is a reasonable heuristic to design a search distribution given the results of evaluation of previous samples? We will follow here the idea of Estimation-Of-Distribution Algorithms (Pelikan et al., 1999) which can be described as follows.

We assume that the outcome of evaluation is given as a distribution  $p$  over  $P$ , which might represent the parent population.<sup>2</sup> Given that the space of feasible search distributions is limited to  $Q$ , a simple heuristic to choose the new search distribution is to pick the one that is most similar to  $p$ . Similarity can be measured by the Kullback-Leibler divergence (KLD)  $D(p : q) = \mathbb{E}_p \left\{ \log \frac{p(x)}{q(x)} \right\}$  between two distributions, which also captures the structural similarity between two distributions in the sense of the similarity of different order dependencies (see (Amari, 2001), please note the relations between the KLD, log-likelihood, free energy, and mean energy<sup>3</sup>). Thus, one heuristic to design the new search distribution  $q'$  reads

$$q' = \underset{q \in Q}{\operatorname{argmin}} D(p : q) . \quad (1)$$

We term this specific kind of an EDA *KL-search*. In general, the crucial parameter of KL-search is the choice of the set  $Q$  of feasible search distributions. On the one hand, the choice of  $Q$  determines the computational cost of the minimization (1) in every step. On the other hand, it determines the algorithm's capability to express and exploit the structure observable in  $p$ .

Some algorithms of the class of EDAs are exact instantiations of KL-search: MIMIC (de Bonet et al., 1997) chooses  $Q$  to be the set of Markov chains, PBIL (Baluja, 1994) chooses  $Q$  as the set of factorized distributions. Other EDAs differ from KL-search in the choice of the similarity measure (they use alternatives to the KLD, e.g., BOA (Pelikan et al., 2000) takes a Bayesian

---

<sup>2</sup> Generally, in this formalism  $p$  is meant to encode any information that we receive from evaluations. Typically,  $p$  is non-vanishing only on a finite set of samples (the offspring population) and the values of  $p$  might be (in a normalized way) the fitness values of these offspring. Alternatively,  $p$  might represent a resampling of such a selection distribution, which corresponds to the parent population.

<sup>3</sup> With the definitions of the entropy  $H(p) = -\mathbb{E}_p \{ \log p(x) \}$  and the log-likelihood  $\mathcal{L}(q) = \mathbb{E}_p \{ \log q(x) \}$  we have  $D(p : q) = -H(p) - \mathcal{L}(q)$ . One could roughly say, "minimizing the KLD means maximizing the log-likelihood *and* the entropy". Further, when defining an energy functional  $E(x) = -\log q(x)$ , the mean energy  $E = \mathbb{E}_p \{ E(x) \} = -\mathcal{L}(q)$  is the negative log-likelihood while the free energy  $F = E - H(p)$  is the KLD.

Dirichlet Metric). But all of them can distinctly be characterized by their choice of  $Q$ , which may also be the set of dependency trees (COMIT, Baluja and Davies, 1997), Bayesian networks (BOA, Pelikan et al., 2000), or Bayesian networks with local structures (hBOA, Pelikan and Goldberg, 2003).

Despite its conceptual simplicity, the minimization required in each iteration of KL-search can be computationally very expensive, depending on the complexity of the distributions in  $Q$ . When only simple distributions, like factorized distributions (PBIL) or Markov chains (MIMIC) are allowed, the minimization can be calculated directly. For more complex distributions (like Bayesian networks, BOA), the minimization itself requires an iterative procedure.

Finally note that distributions in  $Q$  should typically be constrained to have a minimum entropy. In that way, a repeated cycle of entropy decrease (in the course of evaluation) and entropy increase (when picking a new search distribution) ensures exploration and prevents the algorithms from early convergence.

### 3 Indirect induction of search distributions

In this paper we are interested in indirect codings of search points. In the heuristic search framework, this means that a distribution  $\tilde{q} \in \tilde{\Lambda}$  over a coding space  $G$ , the *genotype space*, is maintained. The search distribution  $q$  over the actual search space  $P$  (*phenotype space*) is then given indirectly via a *coding*  $\phi: G \rightarrow P$ ,

$$q(x) = \sum_{g \in [x]_\phi} \tilde{q}(g) , \quad \text{where } x \in P, [x]_\phi = \{g \in G \mid \phi(g) = x\} . \quad (2)$$

We also use the short notation  $q = \tilde{q} \circ \phi^{-1}$  for this projection of  $\tilde{q}$  under  $\phi$ . The set  $[x]_\phi$  of all genotypes mapping to the same phenotype  $x$  is an equivalence class under  $\phi$ , also called *neutral set* of  $x$ .

Three additional constraints define the “indirect encoding case” considered in this paper: **First**, we impose that  $\phi: G \rightarrow P$  shall be bijective (one-to-one). We denote the space of all bijective codings  $G \rightarrow P$  by  $\Phi$ . The discussion in section 6 will establish a relation to non-bijective codings.

**Second**,  $G$  is the space of variable length strings over some finite alphabet  $\mathcal{A}$ ,

$$G = \bigcup_{l=1}^{\infty} \mathcal{A}^l . \quad (3)$$

It is technically unclear how to define a marginal over the  $i$ th symbol (or mu-

tual information between symbols) directly for a variable length distribution  $\tilde{q} \in \tilde{\Lambda}$ . Thus we will consider the decomposition

$$\tilde{q}(g) = \tilde{q}(g|l) \tilde{q}(l) , \quad \text{where } g \in G, l = \text{length}(g) \in \mathbb{N} . \quad (4)$$

Here,  $\tilde{q}(l)$  is a distribution over the genotype length  $l \in \mathbb{N}$ , and  $\tilde{q}(g|l)$  the conditional distribution over a fixed length alphabet  $\mathcal{A}^l$ . We use the short notation  $\tilde{q}^l \equiv \tilde{q}(\cdot|l)$  for this length-conditioned distribution. The marginal  $\tilde{q}_i^l$  over the  $i$ th symbol ( $i \leq l$ ) and the mutual information  $I(\tilde{q}^l) = \sum_i H(\tilde{q}_i^l) - H(\tilde{q}^l)$  can then be defined as usual in terms of the marginal entropies and total entropy.<sup>4</sup>

As the **third** constraint we limit the space of possible search distributions in a certain way: We impose that the length-conditioned distributions  $\tilde{q}^l$  on the genotype space have to factorize, i.e.,

$$\tilde{q}^l(g) = \prod_{i=1}^l \tilde{q}_i^l(g_i) , \quad (5)$$

or, equivalently, that the mutual information  $I(\tilde{q}^l)$  vanishes. We denote the set of feasible genotype distributions by

$$\tilde{Q} \subseteq \{\tilde{q} \in \tilde{\Lambda} \mid \forall l : I(\tilde{q}^l) = 0\} . \quad (6)$$

The space  $Q \subseteq \Lambda$  of feasible search distribution over  $P$  then is

$$Q = \{\tilde{q} \circ \phi^{-1} \mid \phi \in \Phi, \tilde{q} \in \tilde{Q}\} . \quad (7)$$

Table 1 recollects the basic notation we have introduced. In summary, indirect induction of the search distribution means that, in order to design a search distribution  $q \in \Lambda$  we have to pick a bijective coding  $\phi \in \Phi$  and a decorrelated distribution  $\tilde{q} \in \tilde{Q}$  on the genotype space. In other terms, every feasible search distribution  $q$  corresponds to a pair  $(\phi, \tilde{q})$ .

#### 4 Indirect Estimation-Of-Distribution via compression

KL-search proposes how to pick a new search distribution out of  $Q$  incorporating the knowledge on the evaluation of previous samples. In the previous section we specified a  $Q$  that defines the indirect coding case, where a choice

<sup>4</sup> For a distribution  $q$  over some product space  $\mathcal{A} \times \dots \times \mathcal{A}$ , we generally denote the marginal over the  $i$ th variable by  $q_i$ . The mutual information  $I(q) = \sum_i H(q_i) - H(q)$  measures all dependencies between variables of any order (not only pair-wise dependencies).

Table 1  
Basic notations.

in $P$	in $G$	description
$x$	$g = (g_1, \dots, g_l)$	a sample (a string in the case of $G$ )
$\Lambda$	$\tilde{\Lambda}$	space of distributions
$Q$	$\tilde{Q}$	space of feasible search distributions
$p$	$\tilde{p} = p \circ \phi$	distribution to be estimated (selected population)
$q$	$\tilde{q} = q \circ \phi$	search distribution
	$\tilde{q}(g) = \tilde{q}(l) \prod_i \tilde{q}_i^l(g_i)$	factorization of feasible search distribution on $G$

of  $q$  means to pick a bijective coding  $\phi$  and a factorized distribution  $\tilde{q}$  on  $G$ . Combining this, KL-search amounts to a heuristic to pick a coding  $\phi$  (and a  $\tilde{q}$ ) such that knowledge on previous evaluations is incorporated. In this section we will derive results on how this heuristic to pick a coding reads more explicitly. We first discuss the simpler fixed length case before addressing the general one:

*Fixed length case.* Assume that  $G$  contains only strings of fixed length  $l$ ,  $G = \mathcal{A}^l$ . Then the marginals  $\tilde{q}_i$  are straight-forward to define and  $I(\tilde{q}) = \sum_i H(\tilde{q}_i) - H(\tilde{q}) = 0$  constrains  $\tilde{q}$  to vanishing dependencies between genes.

Given definition (7) of  $Q$  for the case of indirect codings, we have

$$D(p : q) = \sum_x p(x) \ln \frac{p(x)}{\sum_{g' \in [x]_\phi} \tilde{q}(g')} = \sum_g \tilde{p}(g) \ln \frac{\tilde{p}(g)}{\tilde{q}(g)}. \quad (8)$$

Here we defined  $\tilde{p}$  as the back-projection of  $p$  onto the coding space  $G$ ,  $\tilde{p} = p \circ \phi$ . The last step uses that  $\phi$  is bijective such that there exists exactly one  $g \in [x]_\phi$ . Next we expand equation (8) into two terms and use that  $\tilde{q}$  has to factorize (5),

$$\begin{aligned} D(p : q) &= \sum_g \tilde{p}(g) \ln \frac{\tilde{p}(g)}{\prod_i \tilde{p}_i(g_i)} + \sum_g \tilde{p}(g) \ln \frac{\prod_i \tilde{p}_i(g_i)}{\prod_i \tilde{q}_i(g_i)} \\ &= I(\tilde{p}) + D(\tilde{p}_{(1)} : \tilde{q}). \end{aligned} \quad (9)$$

Here we defined  $\tilde{p}_{(1)}(g) = \prod_i \tilde{p}_i(g_i)$  as the “factorized reduction” of  $\tilde{p}$  (see also (Amari, 2001) on how to generally define the  $k$ th order reduction  $\tilde{p}_{(k)}$  of  $\tilde{p}$  containing all and only the dependencies of order  $\leq k$ ).

This result states that, in order to follow the KL-search scheme (1) in the indirect coding case, one should find a coding  $\phi$  and a search distribution  $\tilde{q}$  such that  $I(\tilde{p}) + D(\tilde{p}_{(1)} : \tilde{q})$  is minimized.

Here, I would like to distinguish two cases. In the first case we assume that  $\tilde{Q}$  comprises *all* factorized distributions without a bound on the entropy (equality in equation (6)). In this case, no matter which  $\phi$  is chosen, one can always minimize  $D(\tilde{p}_{(1)} : \tilde{q})$  to zero by picking  $\tilde{q} = \tilde{p}_{(1)}$ . Since  $I(\tilde{p})$  is independent of  $\tilde{q}$ , the minimization (1) can be realized by first optimizing  $\phi$  and then picking  $\tilde{q}$ :

$$\operatorname{argmin}_{q \in Q} D(p : q) = (\phi, \tilde{q}) , \quad \text{where } \phi = \operatorname{argmin}_{\phi} I(\tilde{p}) \text{ and } \tilde{q} = \tilde{p}_{(1)} . \quad (10)$$

We call this procedure (first optimizing  $\phi$ , then picking  $\tilde{q}$ ) the *two step procedure*. Note that  $\tilde{p}_{(1)}$  in the last equation depends on the  $\phi$  chosen before.

However, in a realistic algorithm,  $\tilde{Q}$  should not comprise all factorized distributions but obey a lower bound on the entropy of these distributions to ensure exploration.<sup>5</sup> Hence, in the second case, when  $\tilde{Q}$  is only a subset of all factorized distributions,  $D(\tilde{p}_{(1)} : \tilde{q})$  can generally not be minimized to zero and the minimization of (9) can *not exactly* be decomposed in the two steps of first minimizing  $I(\tilde{p})$  w.r.t.  $\phi$ , and then, for a fixed  $\phi$ , minimizing  $D(\tilde{p}_{(1)} : \tilde{q})$  w.r.t.  $\tilde{q}$ . The exact minimization of (9) remains a coupled problem of finding a pair  $(\phi, \tilde{q})$ .

For completeness, let us estimate a bound on the “error” made when still adopting the two step procedure of minimization. Let  $(\phi^*, \tilde{q}^*)$  be a coding and genotype distribution that indeed minimize (9); and let  $(\phi', \tilde{q}')$  be the result of the two step procedure, i.e.,  $\phi'$  minimizes  $I(\tilde{p}')$  and  $\tilde{q}'$  minimizes  $D(\tilde{p}'_{(1)} : \tilde{q}')$  for the given coding  $\phi'$ . Here,  $\tilde{p}' = p \circ \phi'$  and  $\tilde{p}^* = p \circ \phi^*$ . A rough bound for the error made can be estimated as follows, to be explained in detail below,

$$\begin{aligned} D(p : q') - D(p : q^*) &= D(\tilde{p}'_{(1)} : \tilde{q}') - D(\tilde{p}^*_{(1)} : \tilde{q}^*) + I(\tilde{p}') - I(\tilde{p}^*) \\ &\leq D(\tilde{p}'_{(1)} : \tilde{q}') - D(\tilde{p}^*_{(1)} : \tilde{q}^*) \leq D(\tilde{p}'_{(1)} : \tilde{q}') \\ &= \sum_{i=1}^l D(\tilde{p}'_i : \tilde{q}'_i) \\ &\leq - \sum_{i=1}^l \log \tilde{q}'_i(a_i) = -l \log(1-\alpha) \leq l \hbar . \end{aligned} \quad (11)$$

The first inequality stems from the fact that  $\phi'$  minimizes  $I(\tilde{p}')$  and thus  $I(\tilde{p}') \leq I(\tilde{p}^*)$ . Since both,  $\tilde{p}'_{(1)}$  and  $\tilde{q}'$ , are factorized distributions, their Kullback-Leibler divergence decomposes into a sum. For each marginal, when there is a lower bound  $\hbar$  on the entropy of  $\tilde{q}'_i$ , the divergence  $D(\tilde{p}'_i : \tilde{q}'_i)$  is particularly large when  $\tilde{p}'_i$  has very low entropy. In the worst case,  $\tilde{p}'_i$  has zero entropy, i.e., is

<sup>5</sup> Recall that for  $\tilde{q} = \tilde{p}_{(1)}$ , and since  $\phi$  is bijective,  $H(p) = H(\tilde{p}) = H(\tilde{p}_{(1)}) - I(\tilde{p}) = H(\tilde{q}) - I(\tilde{p}) = H(q) - I(\tilde{p})$ . Therefore, the entropy of search  $H(q)$  would only be greater than  $H(p)$  by the amount of  $I(\tilde{p})$ , which is minimized.

non-zero only for a single symbol  $a_i \in \mathcal{A}$ . In that case  $D(\tilde{p}'_i : \tilde{q}'_i) = -\log \tilde{q}'_i(a_i)$ . In order to minimize  $D(\tilde{p}'_i : \tilde{q}'_i)$ ,  $\tilde{q}'_i$  is chosen to have the form of the typical symbol mutation distribution with mutation rate  $\alpha$ , where  $\tilde{q}'_i(a_i) = 1 - \alpha$  and  $\tilde{q}'_i(a) = \frac{\alpha}{|\mathcal{A}|-1}$  for  $a \neq a_i$ . Then,  $H(\tilde{q}_i) = -(1 - \alpha) \log(1 - \alpha) - \alpha \log \frac{\alpha}{|\mathcal{A}|-1}$ . Given the lower bound  $\tilde{h}$  on the entropy, the minimal mutation rate  $\alpha$  can be chosen to ensure  $H(\tilde{q}_i) = \tilde{h}$  and  $D(\tilde{p}_i : \tilde{q}_i) = -\log(1 - \alpha) \leq \tilde{h}$ . Thus, in the worst case, the “error” made when using the two step procedure instead of the exact minimization of (9) is at most  $l \tilde{h}$ .

*Variable length case.* Let us repeat the above derivations in the general case when  $G = \bigcup_{l=1}^{\infty} \mathcal{A}^l$  comprises strings of any length over the alphabet  $\mathcal{A}$ . The constraint of vanishing mutual information in  $\tilde{q}$  now refers to the length-conditioned distributions  $\tilde{q}^l$ , i.e., we impose  $I(\tilde{q}^l) = 0$  while  $\tilde{q}(l)$  is unconstrained. (Recall  $\tilde{q}(g) = \tilde{q}^l(g) \tilde{q}(l)$  where  $l = \text{length}(g)$ .) Equation (8) now leads to

$$\begin{aligned}
D(p : q) &= \sum_l \sum_g \tilde{p}^l(g) \tilde{p}(l) \ln \frac{\tilde{p}^l(g) \tilde{p}(l)}{\tilde{q}^l(g) \tilde{q}(l)} \\
&= \sum_l \tilde{p}(l) \sum_g \tilde{p}^l(g) \ln \frac{\tilde{p}^l(g)}{\tilde{q}^l(g)} + \sum_l \tilde{p}(l) \ln \frac{\tilde{p}(l)}{\tilde{q}(l)} \\
&= \sum_l \tilde{p}(l) \sum_g \tilde{p}^l(g) \ln \frac{\tilde{p}^l(g)}{\tilde{p}_{(1)}^l(g)} + \sum_l \tilde{p}(l) \sum_g \tilde{p}^l(g) \ln \frac{\tilde{p}_{(1)}^l(g)}{\tilde{q}^l(g)} + D(\tilde{p}(l) : \tilde{q}(l)) \\
&= \text{E}_l \left\{ I(\tilde{p}^l) \right\} + \text{E}_l \left\{ D(\tilde{p}_{(1)}^l : \tilde{q}^l) \right\} + D(\tilde{p}(l) : \tilde{q}(l)) , \tag{12}
\end{aligned}$$

where we introduced  $\text{E}_l \{ \cdot \}$  as the expectation w.r.t.  $\tilde{p}(l)$  (which depends on  $\phi$ ). The entropy of  $p$  can be written as

$$\begin{aligned}
H(p) &= -\sum_g p(g) \ln p(g) = -\sum_l \tilde{p}(l) \sum_g \tilde{p}^l(g) \ln [\tilde{p}^l(g) \tilde{p}(l)] \\
&= -\sum_l \tilde{p}(l) \left[ \sum_g \tilde{p}^l(g) \ln \tilde{p}^l(g) + \sum_g \tilde{p}^l(g) \ln \tilde{p}(l) \right] \\
&= \sum_l \tilde{p}(l) H(\tilde{p}^l(g)) + H(\tilde{p}(l)) \\
&= \sum_l \tilde{p}(l) \left[ \sum_{i=1}^l H(\tilde{p}_i^l) - I(\tilde{p}^l) \right] + H(\tilde{p}(l)) \\
&= \text{E}_l \left\{ \sum_{i=1}^l H(\tilde{p}_i^l) \right\} - \text{E}_l \left\{ I(\tilde{p}^l) \right\} + H(\tilde{p}(l)) . \tag{13}
\end{aligned}$$

Combining equations (12) and (13) we find

**Lemma 1** *In the indirect encoding case, for any  $p \in \Lambda$ , any bijective encoding*

$\phi : G \rightarrow P$ , and any factorized distribution  $\tilde{q} \in \tilde{Q}$ , we have

$$D(p : q) = \mathbb{E}_l \{ I(\tilde{p}^l) \} + \mathbb{E}_l \{ D(\tilde{p}_{(1)}^l : \tilde{q}^l) \} + D(\tilde{p}(l) : \tilde{q}(l)) \quad (14)$$

$$= \mathbb{E}_l \left\{ \sum_{i=1}^l H(\tilde{p}_i^l) \right\} + \mathbb{E}_l \{ D(\tilde{p}_{(1)}^l : \tilde{q}^l) \} + D(\tilde{p}(l) : \tilde{q}(l)) + H(\tilde{p}(l)) - H(p) . \quad (15)$$

The second RHS term in equation (15) is a comparison of only the marginals of  $\tilde{p}^l$  and  $\tilde{q}^l$ , the third term is a comparison of the length distributions, the fourth term is the entropy of the genotype length, which depends on  $\phi$ , and the last term depends only on  $p$ , not on  $\phi$  or  $\tilde{q}$ . The first term in equation (15) is of particular interest here. The following bounds show how it relates to the description length. Note that  $\log |\mathcal{A}|$  is the maximal entropy of a marginal:

$$\mathbb{E}_l \left\{ \sum_{i=1}^l H(\tilde{p}_i^l) \right\} \leq \log |\mathcal{A}| \mathbb{E}_l \left\{ \sum_{i=1}^l 1 \right\} = \log |\mathcal{A}| \mathbb{E}_l \{ l \} = L_p \log |\mathcal{A}| , \quad (16)$$

where we introduce  $L_p = \mathbb{E}_l \{ l \}$  as the expected description length of samples of  $p$  in the encoding  $\phi$ . On the other hand, from (13), we get

$$\mathbb{E}_l \left\{ \sum_{i=1}^l H(\tilde{p}_i^l) \right\} \geq H(p) - H(\tilde{p}(l)) . \quad (17)$$

Note that for an optimally compact coding  $L_p \log |\mathcal{A}| = H(p) - H(\tilde{p}(l))$ ,<sup>6</sup> and both bounds are tight.

We collect these findings in

**Lemma 2** *In the indirect encoding case, the expected description length  $L_p$  of samples from  $p$  (parents) gives an upper bound on  $D(p : q)$ ,*

$$\begin{aligned} 0 &\leq D(p : q) - \mathbb{E}_l \{ D(\tilde{p}_{(1)}^l : \tilde{q}^l) \} - D(\tilde{p}(l) : \tilde{q}(l)) \\ &\leq L_p \log |\mathcal{A}| - H(p) + H(\tilde{p}(l)) . \end{aligned} \quad (18)$$

*For an optimally compact encoding  $\phi$ , these bounds are tight, i.e., we have*

$$D(p : q) = \mathbb{E}_l \{ D(\tilde{p}_{(1)}^l : \tilde{q}^l) \} + D(\tilde{p}(l) : \tilde{q}(l)) , \quad (19)$$

*and  $D(p : q)$  can easily be minimized by adapting the genotype marginals  $\tilde{q}_i^l$  to estimate  $\tilde{p}_i^l$  and setting the length distribution  $\tilde{q}(l)$  equal to  $\tilde{p}(l)$ .*

<sup>6</sup> Here is a slight difference to the usually considered case of channel capacity: In our case, the length of a genome itself can carry information (even for  $\mathcal{A} = \{0\}$ ) of the amount  $H(\tilde{p}(l))$  such that the symbols only need to encode  $H(p) - H(\tilde{p}(l))$  entropy. In the channel capacity case, where a continuous stream of symbols is transmitted, the tight bound is  $L_p \log |\mathcal{A}| = H(p)$ , as for the Shannon-Fano code.

Let us briefly summarize and discuss these results by emphasizing certain aspects:

1. *Compression is doing “more than we need”*: Reconsider the exact expressions (14) and (15) of lemma 1. We related the term  $E_l \left\{ \sum_i H(\tilde{p}_i^l) \right\}$  to the description  $L_p$  via the bound (16) and showed that this bound is exact for an optimal compression. One should note though that compression is not the only way to minimize the term  $E_l \left\{ \sum_i H(\tilde{p}_i^l) \right\}$ ; it can equally be minimized by reducing the marginal entropies  $H(\tilde{p}_i^l)$ , which means not to exploit the expressional power of the alphabet. This can better be understood going back to expression (14) involving the mutual information: Ultimately, what matters is to reduce  $E_l \left\{ I(\tilde{p}^l) \right\}$ , i.e. finding a factorial code, which can also be done perfectly with very low marginal entropies, not exploiting the alphabet. By relating it to the description length we showed that a compression is reducing  $E_l \left\{ I(\tilde{p}^l) \right\}$  while *additionally* trying to exploit the alphabet optimally. Thus, compression is doing “more than we need” from the strict point of view of minimizing  $D(p : q)$ . Clearly, this also means that an optimally compact coding is not the only solution to minimize  $D(p : q)$  via indirect induction—but it is one.

2. *If there are no further constraints on  $\tilde{q}$*  (e.g., no entropy bound) then  $\tilde{q}^l$  and  $\tilde{q}(l)$  can always be set equal to  $\tilde{p}_{(1)}^l$  and  $\tilde{p}(l)$ , thus perfectly minimizing  $E_l \left\{ D(\tilde{p}_{(1)}^l : \tilde{q}^l) \right\} + D(\tilde{p}(l) : \tilde{q}(l))$ . In this case, we have

$$D(p : q) = E_l \left\{ I(\tilde{p}^l) \right\} = E_l \left\{ \sum_{i=1}^l H(\tilde{p}_i^l) \right\} + H(\tilde{p}(l)) - H(p) , \quad (20)$$

and the problem reduces to finding an encoding that extinguishes the mutual information  $E_l \left\{ I(\tilde{p}^l) \right\}$  or, as discussed, an optimal compression.

3. *The two step procedure*: If  $\tilde{Q}$  is additionally constrained by a bound on the entropy, minimizing  $D(p : q)$  remains a coupled problem (eq. (14)) of reducing  $E_l \left\{ I(\tilde{p}^l) \right\}$  by a proper choice of  $\phi$  and reducing  $E_l \left\{ D(\tilde{p}_{(1)}^l : \tilde{q}^l) \right\} + D(\tilde{p}(l) : \tilde{q}(l))$  by a proper choice of  $\tilde{q}$ , which though depends on  $\phi$ . We discussed this coupled problem already in the fixed length case. The two step procedure of first finding a compact coding  $\phi$  of  $p$  and then adapting the marginals of  $\tilde{q}$  to those of  $p \circ \phi$  is an approximate method for this minimization. In the worst case one may miss a reduction of  $D(p : q)$  by an amount  $\leq L_p \hbar$ , when  $\hbar$  is the lower bound on the entropy in each marginal  $\tilde{q}_i^l$ . Note that a compact coding minimizes this worst case error.

## 5 Compression EAs

### 5.1 General Approach

A general approach to design a new Evolutionary Algorithm, exploiting the results on compact representations, is to combine any compression technique with any standard EA. Such a *Compression EA* reads

- (1) Initialize a finite population  $q = \{x_1, \dots, x_\lambda\} \subset P$ .
- (2) Evaluate  $q$  and select a subpopulation  $p = \{x_1, \dots, x_\mu\} \subset q$ .
- (3) Find a compression  $\phi$  that (approximately) minimizes  $L_p = \frac{1}{\mu} \sum_{i=1}^{\mu} \text{length}(\phi^{-1}(x_i))$ .
- (4) On the compressed representation, apply standard heuristic operators (e.g., mixing or EDA-operators) to generate  $\lambda$  offspring from  $\mu$  parents.
- (5) Map the  $\lambda$  offspring from  $G$  back to  $P$  using  $\phi$ , yielding the new search samples  $q$ .
- (6) Repeat from step 2.

The operators in step 4 may be any standard operators used in Evolutionary Algorithms. They have to be memory-less though, since the encoding will change in each iteration step and thus integrating knowledge from previous time steps may become futile.

One should emphasize that the crucial ingredient in a Compression EDA is the choice of the compression technique in step (3). This choice decisively determines the limitedness of the algorithm, i.e., kinds of problems it is able to solve. This is in strong analogy to the choice of the class of probabilistic models used in a traditional EDA. For instance, a simple model class like factorized distributions (PBIL) cannot express dependencies between solution variables and thus fails on problems with complex linkage, whereas more general classes of probabilistic models (e.g., Bayesian Networks, BOA) have a much wider range of problems they can be applied on. The same is true for different choices of compression techniques, although future research needs to clarify more specifically which kinds of compression techniques can express which kinds of dependency structures in problems.

In the remainder of this section, we present two Compression EAs. Both of these algorithms will use the same grammar-based compression technique, which clearly has limits w.r.t. the dependency structures it can capture. Thus, the presented algorithms should be understood as examples of Compression EAs. More specifically, the compression technique is limited in that it only detects dependencies between contiguous symbols and contiguous substrings of a string. This hierarchical notion of contiguity has its origin in the grammar-

based compression and allows the induced model to capture dependencies that go beyond pairwise dependencies on only primitive symbols (i.e., beyond a simple Markovian model). Still, the technique is not versatile to capture arbitrary dependency structures (see also (Toussaint, 2005)).

The first of the presented Compression EAs is an EDA that is a direct implementation of the theoretical framework introduced in the previous sections and is strongly related to a variable length version of PBIL (Baluja, 1994) on a compact representation. The second uses ordinary GA operators on the compact code. Before describing the grammar-based compression in section 5.3 we address some heuristics to cope with finite sampling effect. We will test both algorithms on the Hierarchical XOR problem which aims to demonstrate their scalability up to binary strings of length 16 384.

## 5.2 Coping with stochasticity

When following the above sketched algorithm to design a Compression EDA, there are two specific points at which finite sampling and other stochastic effects lead to problems: (1) The optimal compression of a finite population  $p = \{x_1, \dots, x_\mu\}$  is typically “degenerate”, by which we mean that it assigns a single symbol in  $\mathcal{A} = \{1, \dots, \mu\}$  to each sample. (2) Since the compression and thereby the model building is stochastic (steps 2 and 4 in table 3, see below) some generations might produce rather inferior offspring; a heuristic solution is that the statistics for model building “average” over a number of generations. We address these problems as follows.

### 5.2.1 Stochastic compression levels to avoid degenerate compression.

Consider a finite number of samples  $p = \{x_1, \dots, x_\mu\}$ . A mapping  $\phi$  that maps these samples on an alphabet  $\mathcal{A} = \{1, \dots, \mu\}$  in the manner  $x_i \mapsto i$  clearly is an optimal compression—with  $L_p = 1$ . We call this a *degenerate compression*. On the compressed representation the mutual information obviously vanishes. However, this representation is inapt as a basis for search. If the samples are disjoint, it is impossible to design a search distribution with  $\tilde{q}(l) = \tilde{p}(l)$  that has more entropy than  $\tilde{p}$  (which is  $\log \mu$ ). The degenerate compression violates possible constraints on the lower bound on entropy of  $\tilde{q}$ .

A heuristic to solve this problem is to choose the level of compression stochastically. Since the L-System compression is an iterative procedure, this can be realized by escaping the compression loop at each iteration with a fixed probability  $\beta$ , see step 2 in table 3.

### 5.2.2 Averaging statistics to build probabilistic models.

To build probabilistic models reliably in a finite sampling algorithm one needs sufficient statistics. There are two simple alternatives to collect them. First, by choosing the sample size  $\lambda$  large enough so that the parent population in a single generation is large enough to build a model. Second, by averaging statistics in some way over a number of generations. For standard EDAs both approaches are in result approximately equivalent (Harik et al., 1997) and one should tend to favour the simpler first approach.

However, since our compression technique is in itself stochastic (steps 2 and 4 in table 3) we prefer the second approach for the following reason. The offspring generated in a single generation can, due to the stochasticity of the model building, be inferior to previous generations and thus an unfortunate basis for model building, even when the sample size is chosen large enough. Averaging over generations helps to overcome this problem.

Standard averaging approaches in EDAs involve learning rates. For instance, in PBIL (Baluja, 1994), when at generation  $t$  a new population of selected samples  $p^{(t)} = \{x_1, \dots, x_\mu\}$  is available and its marginals are  $p_1^{(t)}, \dots, p_l^{(t)}$ , then the search distribution parameters (its marginals) are updated by  $q_i^{(t)} = (1 - \alpha) q_i^{(t-1)} + \alpha p_i^{(t)}$ . Similarly, the pair-statistics  $A_{ij}$  in the COMIT algorithm (Baluja and Davies, 1997), which are used to build the dependency tree model, are equally averaged by  $A_{ij}^{(t)} = (1 - \alpha) A_{ij}^{(t-1)} + \alpha B_{ij}^{(t)}$ , where  $B_{ij}^{(t)}$  is the pair-statistics of the current selected samples  $p^{(t)}$ .

We implement a similar heuristic in our Compression EDA—but a simple “decay” averaging of the distribution parameters is prohibitive because the representation changes in each generation and operators have to be memoryless. A simple solution is to store a larger pool  $K = \{x_1, \dots, x_\kappa\}$  of the last  $\kappa$  evaluated samples and use this as a basis to build the probabilistic model. The computational costs in space and time increase only by a constant factor. This is implemented in steps 3 and 4 in table 2, providing a larger selected population  $p$  as a basis to build the probabilistic model.

The algorithm given in table 2 describes the full L-System Compression EDA in detail.

### 5.3 Grammar-based compression

A simple technique of compression is to recursively analyze the parents for frequent pairs of neighboring symbols and replace such pairs by new symbols. We take this approach to compute a L-System compression following

Table 2

**Algorithm 1 (L-System Compression EDA (LC-EDA))**

*parameters:* sample, pool, and selection sizes  $\lambda$  [100],  $\kappa$  [1000],  $\mu$  [100], entropy rate  $\alpha$  [1]

- (1) Initialize a population  $q = \{x_1, \dots, x_\lambda\} \subset P$  with random length 2 bit strings, and a *pool*  $K = \emptyset \subset P$ .
- (2) Evaluate the samples in  $q$ .
- (3)  $K \leftarrow K \cup q$  (eliminate multiples), if  $|K| > \kappa$  delete the oldest samples from  $K$ .
- (4) Select a subset  $p \subset K$  of the  $\mu$  best samples in  $K$ .
- (5) Calculate a compression  $\phi$  for  $p$  (see table 3) yielding  $\tilde{p} = p \circ \phi$ .
- (6) Calculate the length distribution  $\tilde{p}(l)$ , max length  $l^*$ , and marginals  $\tilde{p}_i^l$  for  $\tilde{p}$ .
- (7) Set  $\tilde{q}(l) = (1 - \frac{\alpha}{l^*}) \tilde{p}(l) + \frac{\alpha}{l^*} \mathcal{U}(\{1, \dots, l^*+1\})$  and  $\tilde{q}_i^l = (1 - \frac{\alpha}{l^*}) \tilde{p}_i^l + \frac{\alpha}{l^*} \mathcal{U}(\mathcal{A})$ , where  $\mathcal{U}(\cdot)$  is the uniform distribution over a set.
- (8) Take  $\lambda$  samples  $\{g_1, \dots, g_\lambda\}$  from the distribution  $\tilde{q}$ .
- (9) Decode the samples via  $x_i = \phi(g_i)$ , giving a new set of search samples  $q \leftarrow \{x_1, \dots, x_\lambda\}$ .
- (10) Repeat from step 2.

(Nevill-Manning and Witten, 1997) (see also (Lehman and Shelat, 2002) for a brief review of grammar-based compression). An L-System is a sequence  $\Pi = \langle \pi_1, \dots, \pi_k \rangle$  of  $k$  productions  $\pi_i$ . Given the alphabet  $\mathcal{A}$ , each production  $\pi = (l \rightarrow r_1..r_m)$  consists of a LHS symbol  $l \in \mathcal{A}$  and a sequence  $r_1..r_m$  of RHS symbols. The L-System  $\Pi$  defines a mapping  $\phi$  from one sequence  $x$  to another by applying all productions  $\pi_i$ , in the given order, on  $x$ . Applying a production  $l \rightarrow r_1..r_m$  on  $x$  means to replace every  $l$  that occurs in  $x$  by the sequence  $r_1..r_m$ . The mapping  $\phi$  can be inverted by applying all productions, in reverse order, inversely on a sequence. Inverse application of a production  $l \rightarrow r_1..r_m$  on  $x$  means to replace every subsequence  $r_1..r_m$  that occurs in  $x$  by  $l$ .

Let  $\mathcal{A}$  be the non-negative integer numbers,  $\mathcal{A} = \mathbb{N}_0$ . Starting with a population  $p = \{x_1, \dots, x_\lambda\}$  of integer sequences, there is a straight-forward way to construct an L-System that compresses the population by recursively extracting and encapsulating pairs of symbols that occur frequently in the population. Table 3 describes this scheme; the computational complexity of constructing such an L-System is only linear in the string length when using certain efficient data structures (Nevill-Manning and Witten, 1997).

To give an impression on the codings developed by the L-System Compression scheme, table 4 displays an L-System found to compress a population of identical solutions to the 1024-bit HXOR problem (see below). In this case, we did not exit the compression loop at some stochastic level but calculated the full, degenerate compression such that eventually every sample is represented by a

Table 3

**Algorithm 2 (L-System Compression)**

*input:* a set  $p = \{x_1, \dots, x_\lambda\}$  of integer sequences

*output:* the L-System  $\Pi$  (i.e.,  $\phi$ ) and the compressed set  $p$

*parameters:* the stochastic compression level parameter  $\beta$  [.1]

- (1) Initialize the L-System  $\Pi = \langle \rangle$  and the *new-symbol*  $c = 1 + \max\{\text{integers in } p\}$ .
- (2) With a probability  $\beta$ , exit the algorithm.
- (3) Calculate the frequency of every symbol pair that occurs in  $p$ .
- (4) For every pair  $r_1r_2$  of symbols that occurs more often than once, create a new production  $c \rightarrow r_1r_2$ , append it at the *beginning* of  $\Pi$ , and increment the *new-symbol*  $c \leftarrow c + 1$ . This is to be done in order, starting with the pair of highest frequency, and random order between pairs of equal frequency. If there are no such pairs, exit the algorithm.
- (5) Recode the population  $p \leftarrow \{\phi^{-1}(x_1), \dots, \phi^{-1}(x_\mu)\}$  (effectively, only the new productions will result in replacements).
- (6) Repeat from step 2.

single symbol  $c$ . We find that some productions represent the typical modules of HXOR solutions (like E:0110 or J:10010110) while others represent parts of these modules. Recall that the order in which productions are added to the L-System is stochastic (cf. step 4 in table 3). Consequently, the coding is not strictly hierarchical as a human might have designed it (or DEVREP, see below) and the L-System comprises more pair-productions than minimally necessary to encode the sequence of length 1024 (the minimum would be 19 productions).

#### 5.4 The Hierarchical XOR problem.

The fitness function we consider is the Hierarchical XOR (HXOR) function (Watson and Pollack, 1999; de Jong, 2003). For a string  $x \in \{0, 1\}^n$  we first define a boolean function  $h(x) \in \{0, 1\}$ , determining whether  $x$  is “valid” or not: Let

$$n = \text{length}(x) , \quad l = \lfloor \log_2(n-1) \rfloor \in \mathbb{N}_0 , \quad L = x_{1:2^l} , \quad R = x_{2^{l+1}:n} . \quad (21)$$

Note that  $l$  is an integer cutting the string into a left part  $L$  and a right part  $R$  at position  $2^l$ . We define

$$h(x) = \begin{cases} 1 & \text{if } n = 1 \\ 1 & \text{if } n = 2^l \wedge h(L) = 1 \wedge h(R) = 1 \wedge L = \overline{R} . \\ 0 & \text{otherwise} \end{cases} \quad (22)$$

Table 4

An L-System found to compress a population of identical solutions to the 1024-bit HXOR. The full L-System is composed of 27 productions and reads  $\langle C:10, D:1C, E:0D, F:CD, G:0F, H:EG, I:CH, J:CG, K:EI, L:EJ, M:LI, N:JK, O:NM, P:JM, Q:KO, R:PO, S:KR, T:QS, U:QP, V:UT, W:SV, X:US, Y:PT, Z:PW, a:XY, b:VZ, c:ba \rangle$ . In the table below, we expanded the RHS of these productions; for brevity only the first 15 productions are displayed. The single symbol  $c$  expands to the length 1024 solution.

C:10	K:0110 1001 1001 0110
D:110	L:0110 1001 0110
E:0110	M:0110 1001 0110 1001 1001 0110
F:1 0110	N:1001 0110 0110 1001 1001 0110
G:01 0110	O:1001 0110 0110 1001 1001 0110 0110 1001 0110 1001 1001 0110
H:01 1001 0110	P:1001 0110 0110 1001 0110 1001 1001 0110
I:1001 1001 0110	<i>etc...</i>
J:1001 0110	

Here  $\bar{R}$  is the bit-wise negation of the right part. The last condition  $L = \bar{R}$  means that the bit-wise XOR between left and right part must be true for each bit. For instance, the strings for which  $h(x) = 1$  are, up to length 16,  $\langle 0 \rangle$ ,  $\langle 01 \rangle$ ,  $\langle 0110 \rangle$ ,  $\langle 0110 1001 \rangle$ ,  $\langle 0110 1001 1001 0110 \rangle$ , and their bit-wise negations. Based on  $h$ , we define a fitness function  $H(x) \in \mathbb{N}$ , for  $n \geq 2$ ,

$$H(x) = H(L) + H(R) + \begin{cases} n & \text{if } h(x) = 1 \\ 0 & \text{else} \end{cases}, \quad (23)$$

and  $H(x) = 1$  if  $n = 1$ . To normalize and put a limit on the string length, we define the  $l$ th HXOR function  $H_l(x) \in [0, 1]$ : If  $x$  is longer than  $2^l$ , let  $x' = x_{1:2^l}$  and otherwise  $x' = x$ . Then,

$$H_l(x) = \frac{1}{2^l(l+1)} H(s'). \quad (24)$$

The normalization is given by the highest possible value  $2^l(1+l)$  of  $H(x)$  for a length  $2^l$  string. There exist two global optima of  $H_l$ , namely the two “valid” strings of length  $l$  for which  $h(x) = 1$ .

### 5.5 Results of the L-System Compression EDA

We tested the L-System Compression EDA by running 20 trials on each of the HXOR problems of lengths  $\{2^1, 2^2, \dots, 2^{14} = 16384\}$ . For all problem sizes we chose the same parameter settings of algorithm 1 and 2 as given in the brackets in tables 2 and 3. Clearly, the optimal choice of parameters depends

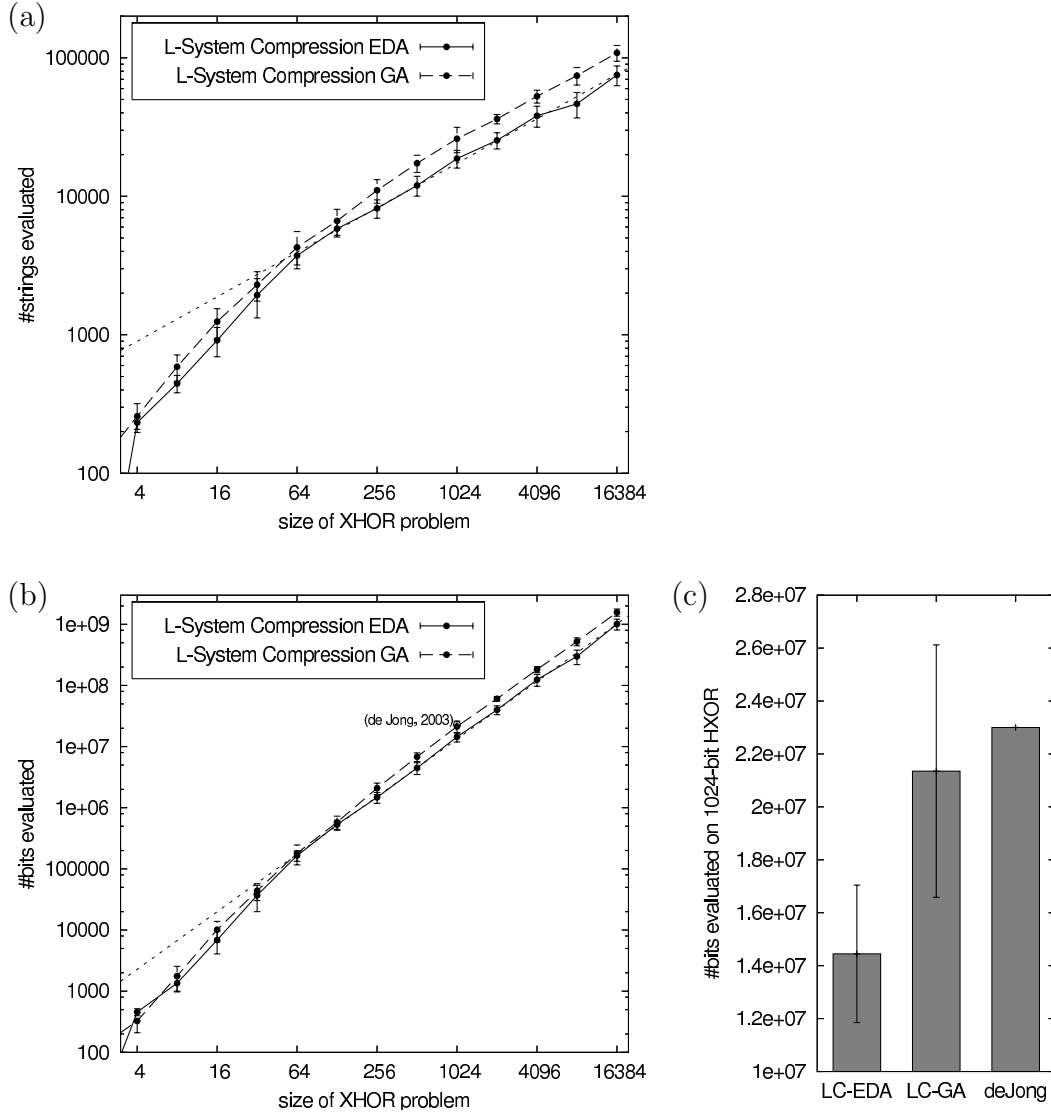


Fig. 1. The number of (a) string evaluations and (b) bit evaluations needed to find an optimum for HXOR problems of lengths  $\{2^1, 2^2, \dots, 2^{14}\}$ . For each problem size we collected 20 independent runs and display here the average and standard deviation. The dotted line are polynomials (a)  $y \propto x^{0.534 \pm 0.010}$  and (b)  $y \propto x^{1.564 \pm 0.012}$  fitted to the performance of the LC-EDA for problem sizes  $\geq 64$ . Figure (c) more clearly than the log-plots displays the differences in performance for 1024-bit HXOR problem.

on the problem size; we chose the same parameter settings for all problem sizes only for simplicity. As a cost measure we investigated the number of string evaluations and the number of bit evaluations (assuming the cost of each string evaluation were equal to its length) needed to find a solution, the latter to allow comparison with (de Jong, 2003).

All runs consistently found a solution. Figure 1 displays the average costs needed together with the standard deviations. Generally, the performance of

the algorithm—in terms of the generations needed—is of the same order as the DEVREP algorithm presented by de Jong (2003), which is the only previous algorithm we are aware of capable of solving large HXOR problems. (de Jong (2003) reported only one result for the 64-bit and 1024-bit HXOR problem, where about  $2.3 \cdot 10^7$  bit evaluations were needed in the single 1024-bit run presented.) The DEVREP algorithm is tailored to hierarchical problems, where different hierarchy levels are explicitly distinguished and mutational variations allowed only within a specific hierarchy level. It should be clear that plain variable length GAs fail on (reasonable length) HXOR problem because of the hierarchy of local minima when search is performed on a direct representation (see (de Jong, 2003) for experiments).

Figure 1 also displays polynomials that are fitted to the performance of the LC-EDA in figures based on the means and variations of the performance considering only problems of size 64 or larger. These suggest that the number of evaluations scales with  $l^{0.534 \pm 0.010}$  when  $l$  is the size of the problem (in bits) while the number of bit evaluations scales with  $l^{1.564 \pm 0.012}$ .

### 5.6 A L-System Compression GA and results

Genetic Algorithms (GAs) are heuristic search schemes, and as such they can be understood by investigating what kind of search distribution they induce depending on previously evaluated samples (the parent population). A traditional GA has independent symbol alterations as a mutation operator, which induces a factorized variability distribution  $\sigma$  for each individual. Considering the whole parent population, a mutation only GA has a mixture  $q = \frac{1}{\mu} \sum_{i=1}^{\mu} \sigma_i$  of factorized distributions as search distribution. Crossover is an operator that makes this search distribution “more factorized”. Actually, one can understand crossover as a move in distribution space from the parent population to a more factorized version of the parent population without changing the gene marginals. One can show that both, mutation and crossover, can only reduce the total mutual information that was present in the parent population (Toussaint, 2003c). In any case, although GAs do not induce a perfectly factorized search distribution, its inherent search mechanisms are of a factorial nature. This is also the reason why they have been modeled as factorial search, e.g., by PBIL or gene pool GAs.

We also implemented a L-System Compression GA, which uses a standard crossover and mutation operator to generate offspring on the compact representation in place of building a probabilistic model. The exact algorithm is given in table 5, the results also in figure 1. The mutation operator applies  $\xi$  mutations on a sample, where  $\xi$  is drawn from a Poisson distribution with mean 1. Possible mutations are (1) replacing, (2) inserting, or (3) deleting a

Table 5

**Algorithm 3 (L-System Compression GA (LC-GA))**

*parameters:* sample and selection sizes  $\lambda$  [100],  $\mu$  [30]

- (1) Initialize a population  $q = \{x_1, \dots, x_\lambda\} \subset P$  with random length 2 bit strings.
- (2) Evaluate the samples in  $q$ .
- (3) Select the subset  $p \subset q$  of the  $\mu$  best samples in  $q$ .
- (4) Calculate a compression  $\phi$  for  $p$  (see table 3) yielding  $\tilde{p} = p \circ \phi$ .
- (5) Apply crossover and mutations (as detailed in the text) on  $\tilde{p}$  to generate  $\lambda$  new offspring  $\tilde{q}$ .
- (6) Decode the offspring  $\tilde{q}$  to  $q$ .
- (7) Add an elitist (the best of  $p$ ) to  $q$ .
- (8) Repeat from step 2.

random symbol from/in the sample; each of these applies with equal probability. The crossover operator splits the parent population randomly into pairs and for each pair applies uniform crossover to generate two offspring. If one parent is longer than the other, one offspring inherits this extra sequence, the other doesn't.

## 6 Discussion

Beyond the introduction of Compression EAs, the main goal of this paper was to provide a theoretical grounding for the use of compact representations for stochastic search. The main results, subsumed in Lemmas 1 and 2, establish the relationship between the Kullback-Leibler divergence and the description length of a variable-length encoding. Namely, minimizing the Kullback-Leibler divergence between a selected distribution  $p$  and a search distribution  $q$  can be achieved by finding a representation that minimizes the description length of  $p$ -samples. Then a simple factorized search distribution  $\tilde{q}$  on the compact representation allows us to model  $\tilde{p}$ .

This result is one piece of a boarder perspective on the role of factorial representations for search. Generally, finding a factorial representation for data to a large extent means to understand structural characteristics of the data. Independent Component Analysis (Hyvärinen and Oja, 2000) and other latent variable models that try to construct factorial representations demonstrate the depth of this concept in the general Machine Learning context. In (Toussaint, 2005), various possibilities to construct factorial representations (in a concrete variable length string context) are discussed. The representations proposed there typically achieve independence between variables by *expanding* the representation. Finding factorial representations by compression seems a

comparably unique approach.

However, the choice of a specific compression technique is crucial for a compression-based search algorithm and decides on whether structural patterns in the data can be detected and exploited for further search. In fact, the problem of compression is strongly related to that of probabilistic modeling; any compression technique first needs to analyze dependencies in the data before it can introduce new symbols for dependent features. Hence it is a matter of the choice of compression technique whether a factorial representation can be found for a specific type of problem and the choice of compression technique is analogous to the choice of class of probabilistic models in traditional EDAs. Standard string compression algorithms, such as Lempel-Ziv compression as well as the L-System compression we investigated here, are limited in that they basically search for contiguous patterns in the original or partly compressed string. This was successful for the HXOR problem since the dependencies can (on various hierarchy levels) be described in terms of contiguous patterns. For many other hard optimization problems (e.g., MAXSAT) the dependencies will typically be between arbitrary variables and hardly detectable for standard string compression techniques. In its generality, we cannot answer the question of what kind of compression techniques are suited for a specific class of problems. More general analysis of dependencies between variables (as sketched in (Toussaint, 2005)) might lead a way to more general compression approaches for search.

The results presented here also have implications for the understanding of natural evolution. In section 5.6 we briefly explained why standard crossover and mutation operators in GAs are comparable to a factorized search distribution. At least approximately, also in nature variability is generated by independent mutational incidents. So, assuming that nature must rely on roughly factorial search distributions, evolution would certainly profit from finding a factorial representation of solutions (organisms). This discussion might seem rather artificial—how could evolution construct itself factorial representations of the distribution of previously successful solutions? However, accounting for the possibility of neutral evolution (Huynen, 1996), this becomes very plausible.

Neutrality means that the genotype-phenotype mapping  $\phi : G \rightarrow P$  is non-injective (i.e., many-to-one). Many different genetic representations  $g \in [x]_\phi$  exist to encode the same phenotype  $x$  and neutral mutations can transform one representation  $g_1$  to another representation  $g_2$  of the same phenotype. It was shown that selection also induces an implicit selectional pressure on the choice of representation of a phenotype, which leads to a directed evolution also within the sets  $[x]_\phi$  (neutral evolution or drift). This selectional pressure turns out to be proportional to the (negative) Kullback-Leibler divergence (which is related to the effective fitness (Nordin and Banzhaf, 1995; Stephens and Vargas, 2000); see also the discussions in (Toussaint, 2003b) and (Burjorjee

and Pollack, 2005)).

So, the central term that was discussed throughout this paper also determines the evolution of genetic representations in the case of neutrality. When here we proved that minimization of the KLD is related to finding a compact representation, this implies that the implicit selection pressure tends to select compact representations within a neutral set. In fact, in (Toussaint, 2003a) genotypes were considered to be L-Systems and neutral mutations could transform one L-System into an equivalent one (which expresses to the same phenotype). The underlying selectional pressure on the genetic representation indeed lead to compact L-Systems to describe phenotypes. There, it was argued that the origin of compactness is the advantage in mutational robustness when every genetic symbol underlies a constant mutation rate. The results derived here additionally imply that compact representations are *structurally* more favorable—the phenotypic variability they induce tends to minimize the Kullback-Leibler divergence and follow the Estimation-of-Distribution search heuristic.

### *Acknowledgments*

I would like to thank Chris Williams for the inspiring discussions on compression and latent variable probabilistic models, and the German Research Foundation (DFG) for their funding of the Emmy Noether fellowship TO 409/1-1, allowing me to pursue this research.

### **References**

- Amari, S., 2001. Information geometry on hierarchy of probability distributions. *IEEE Transactions on Information Theory* 47, 1701–1711.
- Baluja, S., 1994. Population-based incremental learning: A method for integrating genetic search based function optimization and competitive learning. Tech. Rep. CMU-CS-94-163, Comp. Sci. Dep., Carnegie Mellon U.
- Baluja, S., Davies, S., 1997. Using optimal dependency-trees for combinatorial optimization: Learning the structure of the search space. In: *Proc. of Fourteenth Int. Conf. on Machine Learning (ICML 1997)*. pp. 30–38.
- Barbulescu, L., Watson, J.-P., Whitley, D., 2000. Dynamic representations and escaping local optima: Improving genetic algorithms and local search. In: *Seventeenth National Conference on Artificial Intelligence (AAAI)*. pp. 879–884.
- Barron, A., Rissanen, J., Yu, B., 1998. The minimum description length principle in coding and modeling. *IEEE Transactions on Information Theory* 44, 2743–2760.

- Burjorjee, K., Pollack, J., 2005. Theme preservation and the evolution of representation. In: Genetic and Evolutionary Computation Conference (GECCO 2005); Theory of Representations workshop.
- de Bonet, J. S., Isbell, Jr., C. L., Viola, P., 1997. MIMIC: Finding optima by estimating probability densities. In: Mozer, M. C., Jordan, M. I., Petsche, T. (Eds.), *Advances in Neural Information Processing Systems*. Vol. 9. The MIT Press, p. 424.
- de Jong, E. D., 2003. Representation development from Pareto-Coevolution. In: Genetic and Evolutionary Computation Conference (GECCO 2003). pp. 265–276.
- Harik, G., 1999. Linkage learning via probabilistic modeling in the ECGA. Tech. Rep. IlliGAL-97006, Illinois Genetic Algorithms Laborator.
- Harik, G., Lobo, F., Goldberg, D., 1997. The compact genetic algorithm. Tech. Rep. IlliGAL-97006, Illinois Genetic Algorithms Laborator.
- Heckendorn, R. B., Wright, A. H., 2004. Efficient linkage discovery by limited probing. *Evolutionary Computation* 12, 517–545.
- Hornby, G. S., Pollack, J. B., 2001. The advantages of generative grammatical encodings for physical design. In: Proc. of 2001 Congress on Evolutionary Computation (CEC 2001). IEEE Press, pp. 600–607.
- Huynen, M. A., 1996. Exploring phenotype space through neutral evolution. *Journal of Molecular Evolution* 43, 165–169.
- Hyvärinen, A., Oja, E., 2000. Independent component analysis: algorithms and applications. *Neural Networks* 13, 411–430.
- Lehman, E., Shelat, A., 2002. Approximation algorithms for grammar-based compression. In: Proceedings of the Thirteenth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA 2002). ACM/SIAM, pp. 205–212.
- Liepins, G. E., Vose, M. D., 1990. Representation issues in Genetic Algorithms. *Journal of Experimental and Theoretical Artificial Intelligence* 2.
- Nevill-Manning, C. G., Witten, I. H., 1997. Identifying hierarchical structure in sequences: A linear-time algorithm. *Journal of Artificial Intelligence Research* 7, 67–82.
- Nordin, P., Banzhaf, W., 15-19 1995. Complexity compression and evolution. In: Eshelman, L. (Ed.), *Genetic Algorithms: Proc. of Sixth International Conf. (ICGA 1995)*. Morgan Kaufmann, Pittsburgh, pp. 310–317.
- Pelikan, M., Goldberg, D. E., 2003. Hierarchical BOA solves Ising spin glasses and MAXSAT. In: Genetic and Evolutionary Computation Conference 2003 (GECCO 2003). Springer, pp. pp. 1271–1282.
- Pelikan, M., Goldberg, D. E., Cantú-Paz, E., 2000. Linkage problem, distribution estimation, and Bayesian networks. *Evolutionary Computation* 9, 311–340.
- Pelikan, M., Goldberg, D. E., Lobo, F., 1999. A survey of optimization by building and using probabilistic models. Tech. Rep. IlliGAL-99018, Illinois Genetic Algorithms Laboratory.
- Rothlauf, F., Goldberg, D. E., 2003. Redundant representations in Evolutionary Computation. *Evolutionary Computation* 11, 381–415.

- Stephens, C. R., Vargas, J. M., 2000. Effective fitness as an alternative paradigm for evolutionary computation I: General formalism. *Genetic Programming and Evolvable Machines* 1, 363–378.
- Toussaint, M., 2003a. Demonstrating the evolution of complex genetic representations: An evolution of artificial plants. In: *Genetic and Evolutionary Computation Conference (GECCO 2003)*. pp. 86–97.
- Toussaint, M., 2003b. On the evolution of phenotypic exploration distributions. In: Cotta, C., De Jong, K., Poli, R., Rowe, J. (Eds.), *Foundations of Genetic Algorithms 7 (FOGA VII)*. Morgan Kaufmann, pp. 169–182.
- Toussaint, M., 2003c. The structure of evolutionary exploration: On crossover, building blocks, and Estimation-Of-Distribution algorithms. In: *Genetic and Evolutionary Computation Conference (GECCO 2003)*. pp. 1444–1456.
- Toussaint, M., 25-29 June 2005. Factorial representations to generate arbitrary search distributions. In: *Genetic and Evolutionary Computation Conference (GECCO 2005); Theory of Representations workshop*. ACM Press, Washington, D.C., USA, pp. 339–345.
- Vitányi, P. M. B., Li, M., 2000. Minimum Description Length induction, Bayesianism, and Kolmogorov complexity. *IEEE Trans. Inform. Theory* IT-46, 446–464.
- Watson, R. A., Pollack, J. B., 1999. Hierarchically consistent test problems for genetic algorithms: Summary and additional results. In: *Late breaking papers at the Genetic and Evolutionary Computation Conference*. pp. 292–297.
- Whitley, D., Rana, S., Heckendorn, R., 1997. Representation issues in neighborhood search and evolutionary algorithms. In: *Genetic Algorithms and Evolution Strategy in Engineering and Computer Science*. John Wiley & Sons Ltd., pp. 39–58.
- Wright, A. H., Poli, R., Stephens, C. R., Langdon, W. B., Pulavarty, S., 2004. An Estimation of Distribution Algorithm based on maximum entropy. In: *Genetic and Evolutionary Computation Conference (GECCO 2004)*. Springer, Berlin, pp. 343–354.