

EXTREMUM-SEARCHING HIERARCHICAL PARALLEL PROBABILISTIC ALGORITHMS

IVAN KRAMOSIL

Considering a large finite set A together with a mapping f which takes A into a set with linear ordering, the extremum-search problem for $\langle A, f \rangle$ consists in searching for an element of A for which f attains the minimum (or maximum) value. Simple randomization or "non-idealized" parallelism are proved not to improve the situation substantially when compared with a systematic exhaustive inspection. Therefore, hierarchical parallel probabilistic algorithms for the problem in question are suggested and their time computational complexity is investigated and minimized.

1. INTRODUCTION

Many algorithms of artificial intelligence as well as of other domains of applied mathematics contain sub-algorithms for various extremum-search problems. In its full generality, the extremum-search problem can be described in the following very simple way.

Let $A = \{a_1, a_2, \dots, a_N\}$ be a non-empty finite set of abstract elements (i.e. their nature is not important for what follows), let B a non-empty set equipped by a linear ordering \leq . Hence, for all $x, y, z \in B$, $x \leq x$, $x \leq y$ & $y \leq x \Rightarrow x = y$, $x \leq y \Rightarrow (y \leq z \Rightarrow x \leq z)$ and $x \leq y \vee y \leq x$. Let us mention that throughout all this paper the interpretation of B as the real line $(-\infty, \infty)$ equipped by the usual ordering can be kept in mind as a useful illustration. Let f be a mapping defined on A and taking this set into B , let $A_0(f)$ be the set of elements in A on that f takes its minimum value, i.e.

$$(1) \quad A_0(f) = \{a: a \in A, f(a) \leq f(b) \text{ for all } b \in A\} = \\ = \bigcap_{b \in A} \{a: a \in A, f(a) \leq f(b)\}.$$

In the sequel, the analogous formulations and solutions for maximum values will be strictly dual and will not be presented explicitly. Extremum (hence, minimum) – search problem for $\langle A, f \rangle$ consists in finding at least one element from $A_0(f)$; as A is finite, $A_0(f)$ is trivially non-empty.

Clearly, our abilities to solve this problem as well as time, space or other demands and expenses connected with a solution, decisively depend on a priori knowledge and testing devices being at our disposal. Let us consider a situation when these tools are rather limited. Namely, take an oracle which contains an element $a \in A$ (this initialization is not substantial, as shall be easily seen) and works as follows: when $b \in A$ is put on the input, the device compares $f(a)$ with $f(b)$ and replaces a by b iff $f(b) \leq f(a)$; if this is not the case, a is left unchanged. The work finishes by putting at the output the instantaneous value of a , i.e. such an element among the tested ones, for which the value of f is minimum. The device works on the “black-box” principle, so that no insight into the way in which the device decides is possible.

Under the condition that the time necessary to test an element of A is the same for all elements and does not depend on N , and taking into consideration the usual “worst-case” analysis, the time complexity of an exhaustive systematic search for an element from $A_0(f)$ is in the $\mathcal{O}(N)$ -class. Clearly, this worst case occurs if $A_0(f)$ is a singleton, its element being tested as the last. As can be easily seen, the most elementary randomization does not improve the situation. Or, applying the Laplace principle and taking the uniform probability distribution over A (each element can be sampled with the same probability N^{-1}), we obtain this probabilistic algorithm: take M statistically independent random samples from the uniform distribution over A and find an element among them with minimal value of f . The probability of error, i.e. the probability that this element is not in $A_0(f)$, can be easily computed to be

$$(2) \quad PE = (1 - (\text{card } A_0)/(\text{card } A))^M.$$

Hence, for $\text{card } A_0(f) = 1$, M must be at least $(\ln(1/\varepsilon))N$ in order to obtain $PE < \varepsilon$. So, the time complexity is again in $\mathcal{O}(N)$, and this result holds in general, if $\text{card } A_0(f)$ is a constant independent of N . However, one advantage of this probabilistic algorithm is worth to mention for what follows: using the systematic exhaustive searching procedure we must assure, somehow, that each element of A is tested just once. When discussing the time complexity of this algorithm, we have not considered the time consumptions necessary to satisfy this demand. In the case of the probabilistic algorithm in question, on the other hand, statistical samples from A have been taken independently and with possible repetitions. All the “past” of the procedure is “hidden” in the instantaneous value $f(a)$ with which the function value for the element just tested is compared. In what follows, we shall try to improve the probabilistic algorithm outlined above by using the fact that random samples and tests of the sampled elements can be taken and done simultaneously by a number of identical copies of the testing oracle (processors). On the other hand, we shall accept rather realistic assumptions concerning the limited possibilities in inspection

and cummulation of outputs of parallel processors and the corresponding time demands. These assumptions are realistic particularly in comparison with the usual conventions in the theory of non-deterministic algorithms.

2. TWO-LEVEL EXTREMUM-SEARCHING HIERARCHICAL PARALLEL PROBABILISTIC ALGORITHM

Consider sets A , B , a mapping $f: A \rightarrow B$, and $A_0(f)$ as above. The work of the testing oracle described informally in the previous section can be formalized by defining the mapping MIN which ascribes, to the given f and to each finite sequence $\langle a_i \rangle_{i=1}^n$ of elements of A , an element $MIN(f, \langle a_i \rangle_{i=1}^n)$ defined by induction as follows:

$$(3) \quad \begin{aligned} & \text{if } n = 1. \text{ then } MIN(f, \langle a_1 \rangle) = a_1, \\ & MIN(f, \langle a_1, \dots, a_n, a_{n+1} \rangle) = a_{n+1}, \text{ if } f(a_{n+1}) \leq f(MIN(f, \langle a_1, \dots, a_n \rangle)), \\ & MIN(f, \langle a_1, \dots, a_n, a_{n+1} \rangle) = MIN(f, \langle a_1, \dots, a_n \rangle) \text{ otherwise.} \end{aligned}$$

Definition 1. Let A , B , f , $A_0(f)$ be as above, let m , n , k be positive integers. Let $\{X(i, j)\}_{i=1}^m, \{Z_l\}_{l=1}^k$ be an $(m \times n)$ -tuple of statistically independent and identically distributed random variables defined on an abstract probability space $\langle \Omega, \mathcal{S}, P \rangle$, taking their values in A and such that, for each $i \leq m, j \leq n$ and $a \in A$,

$$(4) \quad P(\{\omega \in \Omega, X(i, j)(\omega) = a\}) = (\text{card } A)^{-1} (=N^{-1}).$$

Let $\{Z_l\}_{l=1}^k$ be a k -tuple of statistically independent and identically distributed random variables defined on $\langle \Omega, \mathcal{S}, P \rangle$, taking their values in the set $1, 2, \dots, m$ of integers and such that, for each $l \leq k, s \leq m$,

$$(5) \quad P(\{\omega \in \Omega, Z_l(\omega) = s\}) = m^{-1}.$$

Then the pair $\mathcal{X} = \langle \{X(i, j)\}_{i=1}^m, \{Z_l\}_{l=1}^k \rangle$ is called *two-level extremum searching hierarchical parallel probabilistic algorithm* (2-ESHPPA, abbreviately) for the extremum search problem $\langle A, B, f \rangle$. \mathcal{X} can be identified with a random variable taking $\langle \Omega, \mathcal{S}, P \rangle$ into A and defined in this way:

$$(6) \quad \mathcal{X}(\omega) = \mathcal{X}(A, B, f, \omega) = MIN(f, \{MIN(f, \{X(Z_l(\omega), j)(\omega)\}_{j=1}^n)\}_{l=1}^k).$$

In spite of its rather difficult form, the intuition behind (6) is quite straightforward. Each random variable $X(i, j)$ samples an element from A and the i th processor finds such one among the elements $X(i, 1)(\omega), X(i, 2)(\omega), \dots, X(i, n)(\omega)$ for which the value of f is minimum, putting this elements as its own output. Random variables Z_1, Z_2, \dots, Z_k take random samples from the set of all processors and the (unique) second level processor (supervisor) finds such an element among the output values of the sampled processors, which minimizes the value of f . This value, denoted in (6) by $\mathcal{X}(\omega)$, is called *solution to the extremum search problem* $\langle A, B, f \rangle$ given by the

2-ESHPPA \mathcal{X} . The solution is *correct*, if $\mathcal{X}(\omega) \in A_0(f)$; it is *wrong* otherwise. Probability of error $PE(\mathcal{X}(A, B, f))$ is defined as the probability of wrong solution, i.e.,

$$(7) \quad PE(\mathcal{X}(A, B, f)) = P(\{\omega: \omega \in \Omega, \mathcal{X}(A, B, f, \omega) \in A - A_0\}),$$

and the value of this probability serves as an important quantitative characteristic of statistical qualities of the 2-ESHPPA in question.

Theorem 1. Let $\mathcal{X} = \langle \{X(i, j)\}_{i=1}^m \{j=1}^n, \{Z_l\}_{l=1}^k \rangle$ be a 2-EHSPPA with the parameters m, n, k for the extremum search problem $\langle A, B, f \rangle$, let $\varepsilon > 0$ be given. If $m \cdot n \geq (\ln(2/\varepsilon)) (\text{card } A / \text{card } A_0(f))$, and $k \geq \ln(2/\varepsilon) m$, then $PE(\mathcal{X}(A, B, f)) < \varepsilon$.

Proof. Random variables $X(i, j)$ take $m \cdot n$ independent random samples from A . When computing the probability that at least one sample is in $A_0(f)$, we obtain

$$(8) \quad \begin{aligned} & P(\{\omega: \omega \in \Omega, \sum_{i=1}^m \sum_{j=1}^n \chi_{A_0(f)}(X(i, j)(\omega)) > 0\}) = \\ & = 1 - P(\bigcap_{i=1}^m \bigcap_{j=1}^n \{\omega: \omega \in \Omega, X(i, j)(\omega) \in A - A_0(f)\}) = \\ & = 1 - \prod_{i=1}^m \prod_{j=1}^n P(\{\omega: \omega \in \Omega, X(i, j)(\omega) \in A - A_0(f)\}) = \\ & = 1 - (P(\{\omega: \omega \in \Omega, X(i, j)(\omega) \in A - A_0(f)\}))^{m \cdot n} = \\ & = 1 - (1 - (\text{card } A_0(f) / \text{card } A))^{m \cdot n} \end{aligned}$$

due to the supposed statistical independence of random variables $X(i, j)$ and their uniform probability distribution over the set A . If $m \cdot n \geq (\ln(2/\varepsilon)) \cdot (\text{card } A / \text{card } A_0(f))$, then

$$(9) \quad \begin{aligned} & 1 - (1 - (\text{card } A_0(f) / \text{card } A))^{m \cdot n} \geq \\ & \geq 1 - (1 - (\text{card } A_0(f) / \text{card } A))^{(\ln(2/\varepsilon)) (\text{card } A / \text{card } A_0(f))} = \\ & = 1 - ((1 - (\text{card } A_0(f) / \text{card } A))^{\text{card } A})^{(\ln(2/\varepsilon)) / \text{card } A_0(f)} \geq \\ & \geq 1 - (e^{-\text{card } A_0(f) (\ln(2/\varepsilon)) / \text{card } A_0(f)}) = \\ & = 1 - e^{-\ln(2/\varepsilon)} = 1 - e^{\ln(\varepsilon/2)} = 1 - (\varepsilon/2). \end{aligned}$$

Now, if at least one sample $X(i, j)(\omega)$ is in $A_0(f)$, then at least one among the m processors contains a value from $A_0(f)$ as its output value. $\mathcal{X}(\omega) = 1$ iff at least one such output is sampled by some $Z_l, l \leq k$; due to the assumptions imposed on Z_l we obtain

$$(10) \quad \begin{aligned} & P(\{\omega: \omega \in \Omega, \mathcal{X}(A, B, f, \omega) \in A_0(f)\} | \{\omega: \omega \in \Omega, \\ & \sum_{i=1}^m \sum_{j=1}^n \chi_{A_0(f)}(X(i, j)(\omega)) > 0\}) \geq 1 - (1 - m^{-1})^k \geq 1 - (1 - m^{-1})^{(\ln(2/\varepsilon)) m} \geq \\ & \geq 1 - (e^{-1})^{\ln(2/\varepsilon)} = 1 - (\varepsilon/2) \end{aligned}$$

analogously as in (9). Combining (9) and (10),

$$\begin{aligned}
(11) \quad & P(\{\omega: \omega \in \Omega, \mathcal{X}(A, B, f, \omega) \in A_0(f)\}) = \\
& = P(\{\omega: \omega \in \Omega, \mathcal{X}(A, B, f, \omega) \in A_0(f)\} \mid \{\omega: \omega \in \Omega, \sum_{i=1}^m \sum_{j=1}^n \chi_{A_0(f)}(X(i, j)(\omega)) > 0\}) \cdot \\
& \quad P(\{\omega: \omega \in \Omega, \sum_{i=1}^m \sum_{j=1}^n \chi_{A_0(f)}(X(i, j)(\omega)) > 0\}) \cong \\
& \quad \geq (1 - \varepsilon/2)(1 - \varepsilon/2) > 1 - \varepsilon,
\end{aligned}$$

so that $PE(\mathcal{X}(A, B, f)) < \varepsilon$ and the theorem is proved. \square

3. OPTIMIZATION OF COMPUTATIONAL COMPLEXITIES 2-ESHPPA'S

The conditions of Theorem 1 bind the product mn , hence, the natural question arises, which ratio of m and n is the optimal one from the point of view of minimization of the corresponding time computational complexity. To be able to answer the question we are to precise it, i.e. to define quantitative criteria for time computational complexity, namely, the two following most simple criteria originating from [2] will be taken into consideration. Both these criteria will be defined as linear functions of the number of random samples which need to be successively executed, but the time complexity of particular samples will be treated in a different way. If \mathcal{X} is a 2-ESHPPA with parameters m, n and k , then its *unit time (computational) complexity* $TCU(\mathcal{X})$ is defined by $\alpha_1 n + \alpha_2 m + \alpha_3$ with appropriate $\alpha_1, \alpha_2 > 0$, $\alpha_3 \geq 0$ independent of A, B , and f . The *logarithmic time (computational) complexity* $TCL(\mathcal{X})$ is defined by $\beta_1 n \log_2 \text{card } A + \alpha_2 k \log m + \beta_3$ with appropriate $\beta_1, \beta_2 > 0$, $\beta_3 \geq 0$ independent of A, B , and f . Hence, in the first case the random sample from A is considered to be of constant time complexity no matter which the sampled element and the cardinality of the sample space may be (but for different set, say, A and $\{1, 2, \dots, m\}$, the constants may be different). On the other hand, in the logarithmical case the complexity of random samples increases linearly with the binary logarithm of the cardinality of the sample space, as if the samples were realized through a sequence of independent regular coin tosses (so that the length of this sequence must be $\lceil \log_2 \text{card } A \rceil$, with $\lceil \cdot \rceil$ denoting the upper integer part).

Theorem 2. Consider an extremum search problem $\langle A, B, f \rangle$ and the class \mathcal{H} of all 2-ESHPPA's satisfying the conditions of Theorem 1 with respect to this problem. If $\mathcal{X} \in \mathcal{H}$ with parameters $\tilde{m}_N, \tilde{n}_N, \tilde{k}_N$ ($N = \text{card } A$) is such that $TCU(\mathcal{X}) \leq TCU(\mathcal{X}')$ for all $\mathcal{X}' \in \mathcal{H}$, then, for $v = \text{card } A_0(f)$,

$$\begin{aligned}
(12) \quad & \tilde{m}_N = \lceil \sqrt{(\alpha_1 \alpha_2^{-1} v^{-1})} \sqrt{N} \rceil, \\
& \tilde{n}_N = \lceil \ln(2/\varepsilon) \sqrt{(\alpha_2 \alpha_1^{-1} v^{-1})} \sqrt{N} \rceil, \\
& \tilde{k}_N = \lceil \ln(2/\varepsilon) \tilde{m}_N \rceil.
\end{aligned}$$

If $\mathcal{X} \in \mathcal{H}$ with parameters $\tilde{m}_N, \tilde{n}_N, \tilde{k}_N$ is such that $TCL(\mathcal{X}) \leq TCL(\mathcal{X}')$ for all $\mathcal{X}' \in \mathcal{H}$, then

$$(13) \quad \begin{aligned} \tilde{m}_N &= \lceil \sqrt{(\beta_1 \beta_2^{-1} (\ln 2) v^{-1})} \sqrt{N} + f_1(N), \\ \tilde{n}_N &= \lceil \ln(2/\varepsilon) \sqrt{(\beta_2 (\log_2 e) \beta_2^{-1} v^{-1})} \sqrt{N} + f_2(N), \\ \tilde{k}_N &= \lceil \ln(2/\varepsilon) \tilde{m}_N + f_3(N), \end{aligned}$$

where $f_i(N)$, $i = 1, 2, 3$, are appropriate $o(\sqrt{N})$ -functions.

Proof. For each $x', y', z' \in \langle 0, \infty \rangle$, satisfying $x'y' \geq v^{-1} \ln(2/\varepsilon) N$, $z' \geq \ln(2/\varepsilon) x'$, and minimizing the value of $\alpha_1 x' + \alpha_2 y' + \alpha_3$ ($= TCU(\mathcal{X})$) there are x, y, z such that $xy = \ln(2/\varepsilon) N$, $z = \ln(2/\varepsilon) x$, and $\alpha_1 x + \alpha_2 y + \alpha_3 = \alpha_1 x' + \alpha_2 y' + \alpha_3$. Putting $y = x^{-1} v^{-1} \ln(2/\varepsilon) N$ and $z = \ln(2/\varepsilon) x$ into $\alpha_1 x + \alpha_2 y + \alpha_3$, we obtain

$$(14) \quad f(x) = \alpha_1 x^{-1} v^{-1} \ln(2/\varepsilon) N + \alpha_2 \ln(2/\varepsilon) x + \alpha_3.$$

Hence,

$$(15) \quad (d/dx)f(x) = -\alpha_1 v^{-1} \ln(2/\varepsilon) N x^{-2} + \alpha_2 \ln(2/\varepsilon) = 0$$

yields

$$(16) \quad x^2 = \alpha_1 \alpha_2^{-1} v^{-1} N, \quad x = \sqrt{(\alpha_1 \alpha_2^{-1} v^{-1})} \sqrt{N}$$

so that, after an easy calculation,

$$(17) \quad y = (\ln(2/\varepsilon)) \sqrt{(\alpha_2 \alpha_1^{-1} v^{-1})} \sqrt{N}, \quad z = (\ln(2/\varepsilon)) x,$$

and taking the lowest integer values, we obtain (12).

In the case of the logarithmic criterion we have to minimize the expression

$$(18) \quad \beta_1 y \log_2 N + \beta_2 z \log_2 x + \beta_3$$

under the same constraints, so setting $y = v^{-1} \ln(2/\varepsilon) N x^{-1}$, $z = \ln(2/\varepsilon) x$, and differentiating with respect to x , we need to solve the equation

$$(19) \quad \begin{aligned} (d/dx) (\beta_1 x^{-1} v^{-1} \ln(2/\varepsilon) N + \beta_2 \ln(2/\varepsilon) x \log_2 x + \beta_3) &= \\ = -\beta_1 v^{-1} x^{-2} \ln(2/\varepsilon) N + \beta_2 \ln(2/\varepsilon) \log_2 e ((d/dx)x \ln x) &= \\ = -\beta_1 v^{-1} x^{-2} \ln(2/\varepsilon) N + \beta_2 \ln(1/\varepsilon) \log_2 e (\ln x + 1) = 0. \end{aligned}$$

Hence,

$$(20) \quad \ln x + x^2 = AN/B$$

with

$$(21) \quad A = -\beta_1 v^{-1} \ln(2/\varepsilon), \quad B = \beta_2 \ln(2/\varepsilon) \log_2 e.$$

Let \tilde{x}_N solve the equation (20), let x_N solve the equation $x^2 = AN/B$, then

$$(22) \quad \frac{\ln \tilde{x}_N + \tilde{x}_N^2}{x_N^2} = 1.$$

Evidently, $\tilde{x}_N < x_N$ for almost all N . Let there be $0 < K < 1$ such that, for infini-

tely many N 's, $\tilde{x}_N \leq Kx_N$, then

$$(23) \quad \frac{\ln \tilde{x}_N + \tilde{x}_N^2}{x_N^2} \leq \frac{\ln Kx_N + K^2x_N^2}{x_N^2} = K^2 + \frac{\ln K + \ln x_N}{x_N^2} =$$

$$= K^2 + \frac{\ln K + \ln \sqrt{(AB^{-1})} + (1/2) \ln N}{AB^{-1}N} - K^2$$

as $N \rightarrow \infty$, hence, for each $K', K < K' < 1$, and for almost all N 's, $(\ln \tilde{x}_N + \tilde{x}_N^2)/x_N^2 \leq K'$ and this inequality contradicts (22). So $\lim_{N \rightarrow \infty} (\tilde{x}_N/x_N) = 1$ and computing the values A and B we obtain (13). The assertion is proved. \square

4. MANY-LEVEL HIERARCHIES

Summing the results of Chapters 2 and 3, we can see that two-level hierarchies improve significantly the time complexity for the extremum-searching problem reducing this complexity from the $\mathcal{O}(N)$ -class into the $\mathcal{O}(\sqrt{N})$ -class (in the case of unit criterion), or from the $\mathcal{O}(N \log N)$ class into the $\mathcal{O}(\sqrt{N} \log N)$ -class (for logarithmic criterion). As shown in [3], this quadratic speed-up is, in a sense, the optimum which is reachable by two-level hierarchies with independent random choice. On the other hand, the straightforward idea arises to improve the results of foregoing chapters by an appropriate many-level application of the hierarchical principle. Let us begin with a simple example.

Again, consider an extremum-search problem $\langle A, B, f \rangle$ and suppose that $N = \text{card } A = 2^K$ for a positive integer K (the modifications necessary if it is not the case are not important for what follows). Given $\varepsilon > 0$, set $\delta = \varepsilon/K$, and $N_i = N/2^i$, $i = 1, 2, \dots, K$, so that $N_K = 1$. Consider N_1 processors of the first level, each of them taking $\lceil 2 \ln(1/\delta) \rceil$ independent random samples from the uniform probability distribution over A and putting on its output the latest among the sampled elements which minimizes the function f (according to the definition of the operation *MIN*). Consider N_2 second-level processors, each of them taking, again $\lceil 2 \ln(1/\delta) \rceil$ independent random samples from the uniform probability distribution over the $N_1 = 2N_2$ -element set of first-level processors, and putting on its output the last sampled output value, among the sampled first-order processor, which minimizes the function f (again, the operation *MIN* is applied). Consider N_3 third-level processors and proceed by induction. Finally, there is just one K th level processor, which takes $\lceil 2 \ln(1/\delta) \rceil$ random samples from the set of the two $(K-1)$ -st level processors with one-half probability of sampling for each of them in each step, and which outputs this one from the (one or two) sampled output values which minimizes the function f . Such a statistical decision procedure can be called *many-level extremum-searching hierarchical parallel probabilistic algorithm* (MLESHPA, abbreviately) for the problem $\langle A, B, f \rangle$; if this MLESHPA is

denoted by \mathcal{X} , its output (result) may be denoted by $\mathcal{X}(\omega)$ and we may ask whether this result is *correct* (i.e. whether $\mathcal{X}(\omega) \in A_0(f)$), or whether it is *wrong* (false, i.e. $\mathcal{X}(\omega) \in A - A_0(f)$).

Let us consider the probability $P(\{\omega: \omega \in \Omega, X(\omega) \in A_0(f)\})$. Multiplying the number of first-level processors by the number of samples made by each of them we can easily see that at least $N \ln(1/\delta)$ statistically independent random samples from the uniform probability distribution over A have been taken. As easy calculation analogous to that in the proof of Theorem 1 yields that with probability at least $1 - \delta$ at least one processor samples an element from $A_0(f)$ at least once. But, if a processor samples a value from $A_0(f)$, its output value must be also in $A_0(f)$ due to the definition of the *MIN*-operation. Hence, with a probability at least $1 - \delta$ there is at least one element from $A_0(f)$ among the output values of the N_1 first-level processors. This calculation can be repeated by induction and we obtain: if there is at least one element from $A_0(f)$ among the output values of the i th level processors, then with probability at least $1 - \delta$ there is a value from $A_0(f)$ also among the output values of the $(i + 1)$ -st level processors. Combining these conditional probabilities we obtain that the output of the K th level processor (i.e. $\mathcal{X}(\omega)$) is in $A_0(f)$ with probability at least

$$(24) \quad (1 - \delta)^K > 1 - K\delta = 1 - \varepsilon,$$

so that the probability of error is majorized by ε . At each of the K levels we have performed $2 \ln(1/\delta)$ samples, so the number of samples to be taken sequentially is

$$(25) \quad 2K \ln(1/\delta) = 2 \log N \ln(K/\varepsilon) = 2 \log N \ln((\log_2 N)/\varepsilon).$$

Generating straightforwardly the definition of *TCU* and supposing that $\alpha_i = 1$ for all $i \leq K$ and $\alpha_{K+1} = 0$, the expression (25) yields immediately $TCU(\mathcal{X})$. In the case of logarithmic criterion with $\beta_i = 1$, $i \leq K$, and $\beta_{K+1} = 0$, we obtain that

$$(26) \quad \begin{aligned} TCU(\mathcal{X}) &= \sum_{i=0}^K \lceil 2 \ln(1/\delta) \rceil \log_2 N_i = \\ &= \sum_{i=0}^{\log_2 N} \lceil 2 \ln((\log_2 N)/\varepsilon) \rceil \log_2(N/2^i) = \sum_{i=0}^{\log_2 N} 2i \lceil \ln((\log_2 N)/\varepsilon) \rceil = \\ &= 2 \lceil \ln((\log_2 N)/\varepsilon) \rceil (1/2) \lceil \log_2 N \rceil (\lceil \log_2 N \rceil - 1) = \\ &= \lceil \ln((\log_2 N)/\varepsilon) \rceil ((\log_2 N)^2 - \log_2 N). \end{aligned}$$

Hence, $TCU(\mathcal{X})$ is in $\mathcal{O}(\log N \log \log N)$ and $TCL(\mathcal{X})$ in $\mathcal{O}((\log N)^2 \log \log N)$, so that the result for two-level hierarchies are substantially improved. A result showing the limitations of this improvement will be given later.

Generalizing this example, the formalized definition can read as follows.

Definition 2. *Many-level extremum-searching hierarchical parallel probabilistic algorithm \mathcal{X} for extremum-search problem $\langle A, B, f \rangle$ is defined by a sequence*

$$(27) \quad \mathcal{X} = \langle K, \{N_i\}_{i=1}^K, \{k_i\}_{i=1}^K, \{X_{ij}\}_{i=1}^{N_1} \}_{j=1}^{k_1}, \{Z(i, j, l)\}_{i=1}^{K-1} \}_{j=1}^{N_{i+1}} \}_{l=1}^{k_{i+1}} \rangle$$

where

$K \geq 2, N_1 > N_2 > \dots > N_K = 1, k_2, k_2, \dots, k_K$ are positive integers,

$\{X_{ij}\}_{i=1}^{N_1} \{j=1}^{k_1}$ is an $(N_1 \times k_1)$ -tuple of statistically independent and identically distributed random variables, each of them taking the probability space $\langle \Omega, \mathcal{S}, P \rangle$ into A and such that, for each $a \in A$,

$$(28) \quad P(\{\omega: \omega \in \Omega, X_{ij}(\omega) = a\}) = (\text{card } A)^{-1} (=N^{-1}).$$

For each $i \leq K-1, k \leq N_{i+1}, l \leq k_{i+1}, Z(i, k, l)$ are statistically independent random variables defined on $\langle \Omega, \mathcal{S}, P \rangle$. Each $Z(i, k, l)$ takes values in $\{1, 2, \dots, N_i\}$ and for each $j \leq N_i, l \leq k_{i+1}$,

$$(29) \quad P(\{\omega: \omega \in \Omega, Z(i, k, l)(\omega) = j\}) = N_i^{-1}.$$

\mathcal{X} can be taken as a random variable taking $\langle \Omega, \mathcal{S}, P \rangle$ into A and defined, by induction, as follows.

$$(30) \quad \begin{aligned} Y_k^1(\omega) &= \text{MIN}(f, \langle X_{kj}(\omega) \rangle_{j=1}^{k_1}), \quad k = 1, 2, \dots, N_1, \\ Y_k^{i+1}(\omega) &= \text{MIN}(f, \langle Y_{Z(i,k,j)}^i(\omega) \rangle_{j=1}^{k_{i+1}}), \quad k = 1, 2, \dots, N_{i+1}, \\ \mathcal{X}(\omega) &= \mathcal{X}(A, B, f, \omega) = Y_1^K(\omega). \end{aligned}$$

As in the case of 2-ESHPPA's the solution $\mathcal{X}(\omega)$ to the problem $\langle A, B, f \rangle$ is *correct*, if $\mathcal{X}(\omega) \in A_0(f)$, and it is *wrong (false)*, if $\mathcal{X}(\omega) \in A - A_0(f)$; the probability of error $PE(\mathcal{X}(A, B, f))$ is defined as the probability of wrong solution.

5. OPTIMIZATION OF COMPUTATIONAL COMPLEXITIES FOR HOMOGENEOUS MLESHPPA's

We shall not solve the optimization problem for MLESHPPA's in its full generality and we limit ourselves to certain subclasses of such algorithms.

Definition 3. Let $\varepsilon > 0, \lambda > 1$ be given, an MLESHPPA \mathcal{X} with K levels for an extremum search problem $\langle A, B, f \rangle$ is called (ε, λ) -homogeneous (is in the class $\mathcal{H}(\varepsilon, \lambda)$), if the following conditions hold

$$(31) \quad K = \lceil \log_\lambda \text{card } A \rceil = \lceil \log_\lambda N \rceil$$

$$(32) \quad N_i/N_{i+1} = \lambda \quad \text{for all } i = 0, 1, 2, \dots, K-1, \quad N_0 = N,$$

$$(33) \quad P(\{\omega: \omega \in \Omega, \bigcup_{j=1}^{N_{i+1}} \{Y_j^{i+1}(\omega)\} \cap A_0(f) \neq \emptyset\} \mid \{\omega: \omega \in \Omega, \bigcup_{j=1}^{N_i} \{Y_j^i(\omega)\} \cap A_0(f) \neq \emptyset\}) \geq 1 - (\varepsilon/K) \quad \text{for each } i = 1, 2, \dots, K-1. \quad \square$$

In fact, condition (31) is not necessary, being implied by (32), and is introduced here just because of a more lucidity. Also the condition (33) is rather simple, in spite of its difficult formalization, and it reads: Supposing that there is an element from

$A_0(f)$ on the output of some of the i th level processors, then with the probability at least $1 - (\epsilon/K)$ there is an element from $A_0(f)$ on the output of some of the $(i + 1)$ -st level processors. Combining these conditional probabilities we can easily deduce, using (24), that $PE(\mathcal{X}(A, B, f)) < \epsilon$ for all $\mathcal{X} \in \mathcal{H}(\epsilon, \lambda)$.

The optimization problem within the class $\bigcup_{\lambda > 1} \mathcal{H}(\epsilon, \lambda)$ is solved by the following assertion.

Theorem 3. Given $\epsilon > 0$, consider the class $\mathcal{H}_\epsilon = \bigcup_{\lambda > 1} \mathcal{H}(\epsilon, \lambda)$ of MLESHPPA's for an extremum search problem $\langle A, B, f \rangle$ such that $\text{card } A_0(f) = 1$. Let the homogeneous unit time (computational) complexity $TCUH(\mathcal{X})$ of \mathcal{X} be defined by $\sum_{i=1}^K k_i$ (generalization of TCU with $\alpha_1 = \alpha_2 = 1$, $\alpha_3 = 0$, to the many-level case). Then λ_N satisfying the equality

$$(34) \quad (\ln \ln N + \ln \epsilon^{-1})(\ln \lambda - 1) = (\ln \lambda - 1) \ln \ln \lambda - 1$$

has the property that for each $\mathcal{X} \in \mathcal{H}(\epsilon, \lambda_N)$, $\mathcal{X}^0 \in \mathcal{H}_\epsilon$, $TCUH(\mathcal{X}) \leq TCUH(\mathcal{X}^0)$, hence, the homogeneous unit time complexity is minimized when the ratio between N_i and N_{i+1} is λ_N .

Proof. If $A_0(f)$ is a singleton and if $\{Y_j^i(\omega)\} = A_0(f)$ for some $i \leq K - 1, j \leq N_i$, we need $[N_i \ln(1/(\epsilon/K))]$ independent samples made in total by the $(i + 1)$ -st level processors in order to assure the validity of (33). If, moreover, $N_i/N_{i+1} = \lambda$, each processor of the $(i + 1)$ -st level must take $[\lambda \ln(1/(\epsilon/K))]$ samples, hence,

$$(35) \quad k_i = \lceil \lambda \ln(1/(\epsilon/K)) \rceil = \lceil \lambda \ln(K/\epsilon) \rceil = \lceil \lambda \ln(\log_2 N/\epsilon) \rceil,$$

and, for $\mathcal{X} \in \mathcal{H}(\epsilon, \lambda)$,

$$(36) \quad \begin{aligned} TCUH(\mathcal{X}) &= \sum_{i=1}^K k_i = K \lceil \lambda \ln((\log_2 N)/\epsilon) \rceil = \\ &= \lambda (\log_2 N) \ln((\log_2 N)/\epsilon) = \lambda (\ln N / \ln 2) \ln((\ln N)/((\ln 2) \epsilon)). \end{aligned}$$

The only we need to find is the minimum value of the function

$$(37) \quad \begin{aligned} f(x) &= (x/\ln x)(\ln \ln N - \ln \ln x - \ln \epsilon) = \\ &= (x/\ln x)(\ln \ln N - \ln \epsilon) - ((x \ln \ln x)/\ln x). \end{aligned}$$

But,

$$(38) \quad (d/dx)(x/\ln x) = (\ln x - 1)(\ln x)^{-2},$$

$$(39) \quad \begin{aligned} (d/dx)(x \ln \ln x / \ln x) &= \\ &= (\ln x)^{-2} [(d/dx)(x \ln \ln x)] \ln x - x \ln \ln x [(d/dx) \ln x] = \\ &= (\ln x)^{-2} [\ln x (\ln \ln x + (\ln x)^{-1}) - \ln \ln x], \end{aligned}$$

so that

$$(40) \quad \begin{aligned} (df(x)/dx) &= (\ln x)^{-2} (\ln \ln N + \ln e^{-1}) (\ln x - 1) - \\ &\quad - (\ln x)^{-2} ((\ln x - 1) \ln \ln x + 1), \end{aligned}$$

and this value is zero just when (34) holds, so the assertion is proved. \square

The only value λ_N which minimizes $TCUH(\mathcal{X})$, and which is in $\mathcal{O}(\ln N)$ -class, can be approximated by $\lambda_N \approx e = 2.718 \dots$ in the following sense. Take the function

$$(41) \quad h(\lambda, N) = (\lambda/\ln \lambda) \ln N \ln \ln N,$$

if $\mathcal{X} \in \mathcal{A}(\varepsilon, \lambda)$, then

$$(42) \quad \begin{aligned} \lim_{N \rightarrow \infty} \frac{TCUH(\mathcal{X})}{h(\lambda, N)} &= \lim_{N \rightarrow \infty} \frac{(\lambda/\ln \lambda) \ln N (\ln \ln N - \ln \ln \lambda_N \varepsilon)}{(\lambda/\ln \lambda) \ln N \ln \ln N} = \\ &= 1 - \lim_{N \rightarrow \infty} \frac{\ln \ln \lambda_N \varepsilon}{\ln \ln N} \geq 1 - \lim_{N \rightarrow \infty} \frac{\ln \ln \ln N}{\ln \ln N} = 1. \end{aligned}$$

Hence, $TCUH(\mathcal{X}) \approx h(\lambda, N)$ in the sense of (42) and $h(\lambda, N)$, taken as a function of λ , is minimized by $\lambda = e$ (as $(d/d\lambda)(\lambda/\ln \lambda) = (\ln \lambda)^{-2} (\ln \lambda - 1)$). The fact that only $\lambda_N \in \mathcal{O}(\ln N)$ are considered is not restrictive, as for $\lambda_N \in \mathcal{O}(\ln N)$ we obtain

$$(43) \quad TCUH(\mathcal{X}) \in \mathcal{O}((\ln N)^2)$$

and this result is qualitatively worse than that with $\lambda = const$ (when we were in $\mathcal{O}(\ln N \ln \ln N)$). For the same reasons also λ_N increasing more rapidly than $\ln N$ can be avoided from our considerations.

Let us close this paper by some remarks concerning the further possibilities of investigation and development in the domain in question. Some of them are inspired by [4], when we have investigated similar hierarchical parallel algorithmic structures for the searching problem; this problem consists in answering the question whether a subset of the basic set A is empty or not or not supposing that the only device being at hand is an oracle which tests whether at random sampled elements of V are in the subset in question or not. So the following problems may be considered:

- optimization of \mathcal{X} for broader classes than those investigated here.
- instead of probability of error in the sense defined above (i.e. $\mathcal{X}(\omega) \in A - A_0(f)$) some other criterion should be optimized, e.g., the expected distance $\mathbb{E}q(\mathcal{X}(\omega), A_0(f))$ between $\mathcal{X}(\omega)$ and $A_0(f)$ should be minimized supposing that A is equipped by a metric q and $q(\mathcal{X}(\omega), A_0(f)) = \min_{x \in A_0(f)} q(\mathcal{X}(\omega), x)$.
- instead of "classical" worst case analysis the Bayesian one can be applied, so that the extremum-search problem $\langle A, B, f \rangle$ is supposed to be sampled at random (by the so-called a priori distribution) and probabilities of error or other qualitative characteristics for the worst cases are replaced by their expected values with

respect to the apriori distribution. This involves a number of problems connected with the Bayesian approach in general, and we do not go here into details.

- the possibility of failure of the oracle which compares $f(a)$ and $f(b)$, for $a, b \in A$ can be incorporated into the model, i.e. we may suppose that if $f(a) \leq f(b)$, then the device gives this correct answer with probability $1 - \varepsilon < 1$, and gives the wrong answer $f(a) > f(b)$ with positive probability ε . In [4] we have investigated this possibility of failure for the searching problem, and the way in which the probability of failure combines with the other uncertainties within our model proved to lead to qualitatively different results than those in the case of failure-proof oracles.

In the list below, Feller's monography may serve for all references concerning the elementary notions, models and results of probability theory used in this paper.

(Received April 30, 1987.)

REFERENCES

- [1] W. Feller: *An Introduction to Probability Theory and its Applications*, Vol. I. Second edition. John Wiley and Sons — Chapman and Hall, New York—London 1957 (Russian translation: Mir, Moscow 1964).
- [2] U. Manber and M. Tompa: The complexity of problems on probabilistic, non-deterministic, and alternating decision trees. *J. Assoc. Comput. Mach.* 32 (1985), 3, 720—732.
- [3] J. Reif: On synchronous parallel computations with independent probabilistic choice. *SIAM J. Comput.* 13 (1984), 1, 46—55.
- [4] I. Kramosil: Hierarchical connection of probabilistic approach and parallelism in searching tasks of artificial intelligence (in Czech). In: *Aplikace umělé inteligence AP' 87*, 23—31.

RNDr. Ivan Kramosil, DrSc., Ústav teorie informace a automatizace ČSAV (Institute of Information Theory and Automation — Czechoslovak Academy of Sciences), Pod vodňárenskou věží 4, 182 08 Praha 8, Czechoslovakia.