



(19) **United States**

(12) **Patent Application Publication**  
**Walker**

(10) **Pub. No.: US 2012/0016819 A1**

(43) **Pub. Date: Jan. 19, 2012**

(54) **DISTRIBUTED MULTIMEDIA DOCUMENT INDEXING STRATEGIES**

(52) **U.S. Cl. .... 706/12**

(75) **Inventor: Reginald L. Walker, Los Angeles, CA (US)**

(57) **ABSTRACT**

(73) **Assignee: Tapicu, Inc., Los Angeles, CA (US)**

A method for a system that indexes/ranks/clusters multimedia documents using hybrids of information retrieval algorithms and the stochastic optimization techniques of evolutionary computation (EC) that optimizes parameter sets comprising of object parameters. The method creates a plurality of individual parameter sets, the parameter sets comprising information sharing system object parameters for describing a model, structure, shape, design, process, search query set, or dynamic search space to be optimized and setting the initial population as a current (static parent) population. These parameters are required to filter, organize, and index any large-scale data set—information stored on a single computer, a local area network (LAN), and a wide area network (WAN) that encompasses the whole Internet—that may consist of constantly fluctuating information content over relatively short periods of time

(21) **Appl. No.: 13/135,943**

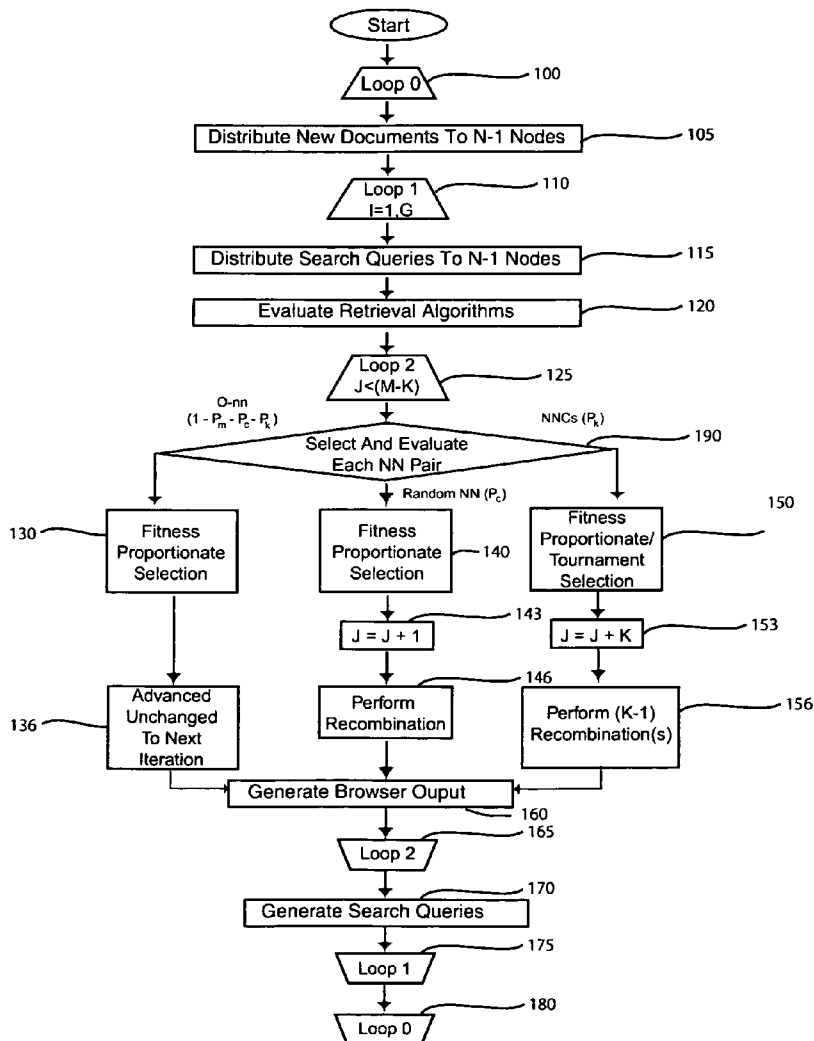
(22) **Filed: Jul. 19, 2011**

**Related U.S. Application Data**

(60) **Provisional application No. 61/399,961, filed on Jul. 19, 2010.**

**Publication Classification**

(51) **Int. Cl. G06F 15/18 (2006.01)**



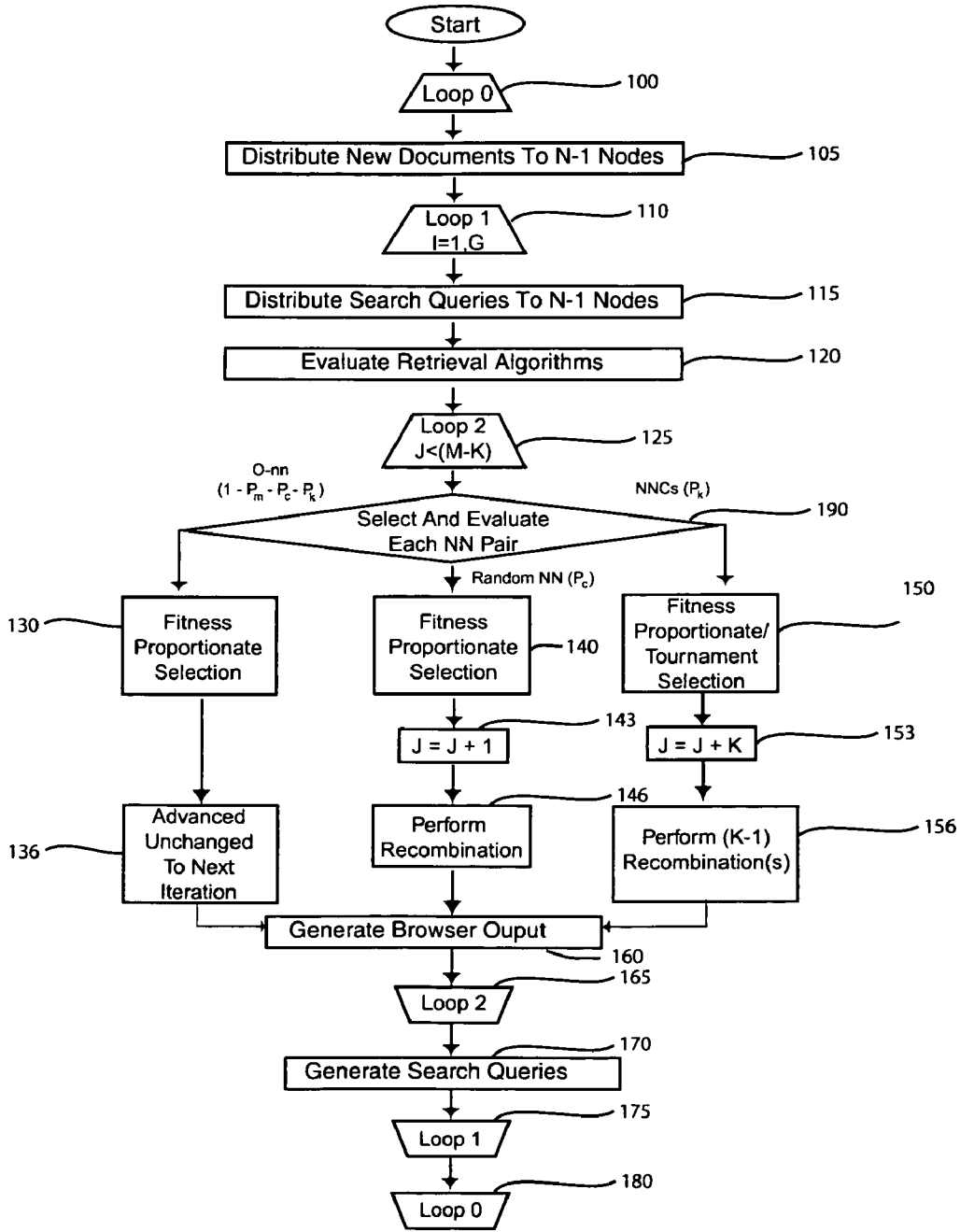


FIG. 1

**DISTRIBUTED MULTIMEDIA DOCUMENT INDEXING STRATEGIES**

**CROSS-REFERENCE TO RELATED APPLICATIONS**

**[0001]** This application claims benefit of provisional application Ser. No. 61/399,961, filed on Jul. 19, 2010 by the present inventor.

**FEDERALLY SPONSORED RESEARCH**

**[0002]** Not Applicable

**SEQUENCE LISTING OR PROGRAM**

**[0003]** Not Applicable

**BACKGROUND OF THE INVENTION**

**[0004]** The invention relates generally to the optimizing of object parameters for describing a model, structure, shape, design, or process, for an information sharing indexer system. In particular, it relates to the stochastic optimization of evolutionary computation (EC) search strategy parameters for multimedia indexers for information sharing indexer systems such as search engines, data warehouses, and service oriented architectures (SOAs). The field of evolutionary computation encompasses stochastic optimization techniques, such as randomized search strategies, in the form of evolutionary strategies (ES), evolutionary programming (EP), genetic algorithms (GA), classifier systems, evolvable hardware (EHW), and genetic programming (GP).

**[0005]** There has always been a need to iteratively improve the clustering and ranking of multimedia documents. The stochastic optimization techniques of evolutionary computation (EC) contain mechanisms which enable the representation of certain unique aspects of individual behavior to improve document clustering. Principles of the stochastic optimization techniques of EC can be found for example in Reginald Louis Walker (2003) *“Tocorime Apicu: Design of an Experimental Search Engine Using an Information Sharing Model”*, University of California Dissertation, UMI Dissertation Publishing, Ann Arbor, Mich. 48106-1346 (www.proquest.com) or rwalker@cs.ucla.edu, which is incorporated by reference herein in its entirety.

**[0006]** The chief differences among the various types of EC stemming from: 1) the representation of solutions (known as individuals in EC), 2) the design of the variation operators (mutation and/or recombination—also known as crossover), and 3) selection mechanisms. A common strength of these optimization approaches lies in the use of hybrid algorithms derived by combining two or more of the evolutionary search methodologies. The underlying optimization methodologies of EC are used to implement unique stochastic aspects of search strategies that are combined with information retrieval methodologies. This mapping is extended by supplementing the search strategies with finding hidden knowledge in a collection of multimedia documents—related and/or unrelated—using search query sets. Canonical multimedia documents are generated to reduce the workload and storage requirements of the system, resulting in a set of condensed multimedia documents forming the data store. The system continuously repartitions the stored document space among a set of nodes whose goal is to form subclusters of nodes for redistributing the workload. The subclusters are formed by using the information retrieval (IR) algorithm metrics

coupled with two or more evolutionary search strategies as the basis of nearest neighbor clusters (NNC) among multimedia indexers. Fitness proportionate and tournament selection in this application forms the basis of nearest neighbor clustering, providing the mechanism for selecting nodes that will share information. Mutations and recombinations are implemented as random change (or multiple changes) of the description of the finite state machine (FSM) according to five different modifications: change of an output symbol, change of a state transition, addition of a state, deletion of a state, or change of the initial state.

**BACKGROUND OF THE INVENTION—OBJECTIVES**

**[0007]** Accordingly, the objectives and advantages of the invention are as follows:

**[0008]** It is an objective of the present invention to use hybrid algorithms derived by combining one or more of the information retrieval methodologies with one or more of the evolutionary computation search methodologies.

**[0009]** It is another objective of the present invention is to provide a stochastic selection process that iteratively improves a population of solutions—evolving sets of competing solutions over the space being searched. The components of an optimization application are:

1. Terminal set. Input variables or constants.
2. Function set. Domain-specific functions that construct potential solutions.
3. Fitness measure(s). Function(s) that assign numeric values to the individuals associated with a population (set of solutions that comprise the solution space).
4. Algorithm control parameters. Settings dependent on population size and workload redistribution (recombination and mutation) rates.
5. Termination criterion. Predicate that uses fitness measures to determine the appropriateness of a population based on tolerances or limits on the number of allowable generations/iterations.

**[0010]** It is another objective of the present invention to represent solutions as memes to reduce in the computational effort to achieve the periodic optimal document clusters. The fitness of a species (adaptive and iterative grouping of the solutions from selective indexers) can be improved by the non-genetic transmission of cultural information that uses a meme as the transmission mechanism rather than the genetically based gene. The difference between the two includes the fact that genetic transmissions (stochastic selection process) evolve over a period of generations, whereas cultural transmissions result from an educational process.

**[0011]** It is another objective of the present invention to use a function set that consists of a multimedia parser that works as a two-pass parser. The initial pass occurs as a component of the system that applies document layout analysis for its automated retrieval component. The second pass applies a full set of text-processing modules consisting of syntactic analysis, lexical analysis, layout analysis, and feature recognition. Layout analysis transforms a raw document into an application-specific document by saving the canonical format structural information as necessary. The syntactic analysis component verifies that the canonical structure adheres to a suitable format. The lexical analysis module is combined with the feature recognition module. These modules remove stop words, identify and record word boundaries, and index words for retrieval. Additionally, this component is respon-

sible for converting hyphenated and sequences of capitalized words into proximity constraints, and case conversions into compressed inverted files.

**[0012]** It is another objective of the present invention to continuously apply algorithm control parameters to improve the subclustering of documents in distributive applications leading to disjoint nodes for chosen sets of search queries.

**[0013]** It is another objective of the present invention to continuously adjust the operational parameters required to filter, organize, and index any large-scale data set—information stored on a single computer, a local area network (LAN), and a wide area network (WAN) that encompasses the whole Internet—that may consists of constantly fluctuating information content over relatively short periods of time.

#### SUMMARY OF THE INVENTION

**[0014]** The invention is a system and method for indexing/ranking and clustering multimedia documents using hybrid search strategies and the stochastic optimization techniques of evolutionary computation (EC). These stochastic optimization techniques form the basis of a regulatory mechanism for sharing information document clustering and ranking which leads to the migration of multimedia documents between multimedia indexers. The iterative application of these mechanisms improves the subclustering of multimedia documents in distributive applications leading to disjoint nodes for chosen sets of search queries.

**[0015]** It is to be understood that both foregoing general description and the following detailed description of the present invention are exemplary and explanatory and are extended to provide further explanation of the invention as claimed.

#### DETAILED DESCRIPTION OF THE DRAWINGS—FIGURES

**[0016]** FIG. 1 is a schematic flow diagram of the optimization method of the present invention.

#### DETAILED DESCRIPTION—PREFERRED EMBODIMENTS

**[0017]** A preferred embodiment of the present invention is now described with reference to the figures where like reference numbers indicate identical or functionally similar elements.

**[0018]** Some portions of the detailed descriptions that follow are presented in terms of algorithms and symbolic representations of operations on data bits within a computer memory. These algorithmic descriptions and representations are the means used by those skilled in data processing art to most effectively convey the substance of their work to others skilled in the art. Algorithms are here, and generally, conceived to be self-consistent sequence of steps (instructions) leading to desired results. The steps are those requiring physical manipulations of physical quantities.

**[0019]** Certain aspects of the present invention include process steps and instructions described herein in the form of an algorithm. It should be noted that the process steps and instructions of the present invention could be embodied in software, could be downloaded to reside on and be operated from different platforms used by a variety of operating systems.

**[0020]** The present invention also relates to an apparatus for performing the operations herein. This apparatus may be

specially constructed for the required purposes, or it may comprise a general-purpose computer selectively activated or reconfigured by a computer program stored in a computer. Furthermore, the computers referred to in the specifications may include a single processor or may be architectures employing multiple processors designed for increased computing capability.

**[0021]** The algorithms and displays presented herein are not inherently related to any particular computer of other apparatus. Various general-purpose systems may also be used with programs in accordance with the teaching herein, or it may prove convenient to construct more specialized apparatus to perform the required method steps. The required structure for a variety of these systems will appear from the description below. In addition, the present invention is not described with reference to any programming language. It will be appreciated that a variety of programming languages may be used to implement the teachings of the present invention as describe herein, and any references below to specific languages are provided for disclosure of enablement and best mode of the present invention.

**[0022]** In addition, the language used in the specification has been principally selected for readability and instructional purposes, and may not have been selected to delineate or circumscribed the invention subject matter. Accordingly, the disclosure of the present invention is intended to be illustrative, but not limiting, of the scope of the invention, which is set forth in the claims.

#### 1. Notational Conventions

**[0023]** a. Applying the EC Search Strategies to Stochastic Information Fluctuations

**[0024]** FIG. 1 shows an embodiment of the invention which performs the tasks associated with regulating the formulation of NNCs and adapting to information fluctuations. The tasks performed are:

**[0025]** 1. Periodic partitioning of the multimedia document dataset among indexer nodes **105**

**[0026]** 2. Generating workload assignments (resulting from fitness proportionate selection steps) for each node **105**

**[0027]** 3. Distributing dynamic workload assignments **105** and dynamic search query sets **115**

**[0028]** 4. Formulating NNCs using fitness proportionate selection **190**

**[0029]** 5. Selecting source of dynamic search query sets **170**

**[0030]** 6. Repeating step 1 through 5 **100-180**

**[0031]** The traditional EC approach for the recombination and mutation operators, as well as the normal (steady-state) approach, is restricted to one application per iteration for a single set of solutions. The load-balancing model of the multimedia document indexing system uses the EC recombination operator by restricting information sharing between members of disjoint node sets (species) which are chosen in a process that selects and evaluate each nearest neighbor (NN) pair **190**.

**[0032]** NNCs **190** can occur as one of three types based on the number of neighborhood seeds: 1) random seeds, 2) multiple seeds, or 3) overlapping seeds. The occurrence of multiple and overlapping seeds enhances the quality of the total cluster's solution space via the modification of the workload assignments of several nodes during one iteration (superstep).

**[0033]** The iterative formulation of NNCs **125,190** was implemented using the notion of an expandable search space which facilitates adaptive subclusters on an iteration-by-iteration basis. The selection process **190** can be applied multiple times **153**, where one node is the NN seed for one or more nodes—thus providing a stochastic hybrid of the recombination and mutation operators **130,140,150**.

b. Formulation of Nearest Neighbor Clusters (NNCs)

**[0034]** K-nearest neighbors (K-nn) **190,130,136** is implanted as the mutation operator when  $K=0$ .

**[0035]** Random NN **190,140,143,146** are implemented as follows: 1) the first node is randomly chosen, and 2) the second node is chosen by incrementing the node ID of the first node **190**, thereby mimicking the ring communication pattern based on the rank in order to determine adjacent nodes. Recombination is applied to the selected nodes **140,143,146** for each iteration **125-165**. The proportionate fitness method **140** assigns a random number to each neighborhood seed and selects individuals by repeatedly choosing various random numbers until one matches a node's random number.

**[0036]** Multiple neighborhoods (NNCs) **190,150,153,156** exists when there are at least one or more NNCs in which neighborhoods do not overlap. When a single node is a nearest neighbor of two disjoint NNCs, this node may be selected **150** as a NN one or more times based on the existence of one or more completing nodes in the disjoint neighborhoods. The selection of a node when two or more are present in a single neighborhood occurs via proportionate fitness selection **150**.

**[0037]** Overlapping neighborhoods **190,150,153,156** occur when two or more NNCs are formed from the seeds overlapping neighborhoods. The selection of one of the NNCs **150** from overlapping of neighborhoods occurs via two "popular" selection methodologies: 1) the proportionate fitness or roulette wheel selection, and 2) the tournament selection. The proportionate fitness method **150** assigns a random number to each node and selects individuals choosing various random numbers which may match an individual's random number.

**[0038]** The selection processes **190,150** for overlapping neighborhoods uses the radius of two or more nodes resulting in possibly K-nn per cluster by performing the following:

**[0039]** 1. Randomly selects one of the overlapping nodes as the seed of one of the NNCs using the tournament selection method **150**

**[0040]** 2. Using roulette wheel selection **150**

**[0041]** a. Randomly selects a node for recombination

**[0042]** b. Randomly selects a range for recombination

**[0043]** c. Performs recombination **156** on the two nodes only if they are NN using proportionate fitness method **150**

**[0044]** 3. If necessary, repeats step 2 **125-165**

The number of iterations **156** a selected node is used for recombination is random—this potentially providing the node with an emulator of the mutation operator **130** (occurring if the selected node was previously selected during an application of the recombination operator). However, the same node may be chosen for two or more iterations with the possibility of swapping previously exchanged recombinations. The system does not advance until k possible recombinations **156** have been completed. The occurrence of overlapping NNCs regulates the recombination rate and the selection rate. The recombination rate and the selection rate use the information retrieval algorithms to generate stochastic metrics for determining nearest neighbor (NN) resulting in the

emergence of subclustering within each cluster/subcluster since static parents (node existence and hierarchy) are maintained throughout this application.

**[0045]** Another component of the recombination rate and the selection rate stems from overlapping nearest neighbor clusters (NNCs) and is equivalent to sharing information between diverse set of computer processors and/or systems. This phenomenon adds random noise to the whole process by creating, at most K-nn in one component of a superstep based on overlapping NNCs—an event which is beneficial to the prevention of premature convergence and to the incorporation of various optimization techniques such as supersteps and dissassortive mating when selecting nodes from initial subclusters such subspecies A and B. Supersteps resulted from two or more applications of the recombination operator during one iteration (generation) via overlapping NNCs or multiple disjoint NNCs. Dissassortive selection is a results of selecting NN for the recombination operator from a disjoint list of disjoint subcluster members, as in the case of random NN using the even nodes as one cluster of individuals and the odd nodes as a subcluster.

c. Input Parameters

**[0046]** The methodology used in retrieval calculations **120**—computing the stochastic measurements—was based on: 1) generating the canonical representation of the raw multimedia documents—an application-specific document of structural information, and 2) applying the stochastic optimization retrieval algorithms to determine NNCs **190**—computing the raw fitness, standardized fitness, and adjusted fitness.

d. Synchronization Points

**[0047]** FIG. 1 provides periodic synchronization points **165,175,180** used for consistency restoration. Using a self-scheduling policy, the load-balancing model distributes the multimedia documents **105** that comprise the document dataset for each iteration. This random approach to the distribution of documents enables the system to adapt to each machine's characteristics at various stages of this iterative process **100-180**. By requiring that each node start each iteration **100,110,125** on the basis of a consistent state, the synchronization points are used to restore a consistent global state. FIG. 1 allows for continuous updates and redistribution of multimedia documents **105,115,160,170** which incorporate the local and system-wide computational parameter adjustments.

**[0048]** The need for synchronization points **165,175,180** can be traced to scientific applications that are known to exhibit a diverse set of I/O access patterns. These are known as:

**[0049]** 1. Compulsory

**[0050]** 2. Checkpoint/restart

**[0051]** 3. Regular snapshots of the computation's progress

**[0052]** 4. Out-of-core read/writes

**[0053]** 5. Continuous output of data for visualization and other post-processing

The variability in the canonical document size accounts for the seemingly high random file accesses. Combining the file access patterns of all the indexers in the system reflects their compulsory nature. The synchronization points **165,175,180** provide the I/O checkpoints. The regular snapshots of the computation's progress are reflected in the intermediate solutions **160,170** that are created at the end of each iteration **165,175,180**.

[0054] While particular embodiments and applications of the present invention have been illustrated and described herein, it is understood that the invention is not limited to the precise construction and components disclosed herein and that various modifications, changes, and variations may be made in the arrangement, operation, and details of the methods and apparatuses of the present invention without departing from the spirit and scope of the invention as it is defined in the appended claims.

I claim:

1. A method for a system that indexes/ranks/clusters multimedia documents using hybrids of information retrieval algorithms and the stochastic optimization techniques of evolutionary computation (EC) that optimizes parameter sets comprising of object parameters, the steps of:

- a. Creating an initial population of a plurality of individual parameter sets, the parameter sets comprising information sharing system object parameters for describing a model, structure, shape, design, process, search query set, or dynamic search space to be optimized and setting the initial population as a current (static parent) population;
- b. For each individual parameter set in a current (static meme) population, mutating the parameters and optionally applying recombination via crossover and/or supercedure mating to improve the feasibility of the current (static parent) population of individual parameter sets, wherein the strength of an individual object parameter mutation is enlarged by decreasing a noise contribution to enhance the robustness of the optimization;
- c. Evaluating the quality of each in the offspring (static meme) population;
- d. Selecting individuals of the offspring (static population) for recombination/crossover/supercedure mating by tournament selection to be the current (static meme) population in the next generation; and
- e. Repeating steps (b) through (d) to decrease the noise contribution to enhance the robustness of the optimization—the termination criterion.

2. The method of claim 1 wherein said parameter sets include information retrieval indexing, evolutionary computation, and stochastic optimization search strategy parameters.

3. The method of claim 1 wherein the strength reduction of the noise contribution is adapted such that the estimated population variance is reduced substantially below or equal to a prescribed variance governed by the robustness criterion.

4. The method of claim 1 wherein the noise contribution varies for the different object parameters of a parameter set.

5. The method of claim 1 wherein the selection in step (d) is nondeterministic selection of stochastic optimization techniques of evolutionary computations.

6. The method of claim 1 wherein the center of mass recombination known in stochastic optimization techniques of evolutionary computations as nearest neighbor clusters (NNCs) is used in step (b).

7. The method of claim 1 wherein the random sources for mutation/recombination is randomly distributed.

8. The method of claim 1 wherein the estimation of the population variance is subject to hybrids of information retrieval algorithms and the stochastic optimization techniques of evolutionary computation (EC).

9. The method of claim 1 for optimizing the shape of the dynamic search space comprised of any large-scale data set—information stored on a single computer, a local area network (LAN), and a wide area network (WAN) that encompasses the whole Internet—that may consists of constantly fluctuating information content over relatively short periods of time.

10. A method for a system that indexes/ranks/clusters multimedia documents using hybrids of information retrieval algorithms and the stochastic optimization techniques of evolutionary computation (EC) that optimizes parameter sets comprising of object parameters, the steps of:

- a. Creating an initial population of a plurality of individual parameter sets, the parameter sets comprising information sharing system object parameters for describing a model, structure, shape, design, process, search query set, or dynamic search space to be optimized and setting the initial population as a current (static parent) population;
- b. For each individual parameter set in the current (static parent) population, mutating the parameters and optionally applying recombination via crossover and/or supercedure mating to create of the offspring (static parent) population of individual parameter sets;
- c. Evaluating the quality of each individual in the current (static) population;
- d. Selecting individuals of the current (static population) for mutation/recombination/crossover/supercedure mating by tournament selection to be the current (static parent) population in the next generation; and
- e. Repeating steps (b) through (d) to decrease the noise contribution to enhance the robustness of the optimization—the termination criterion.

11. The method of claim 10 wherein said parameter sets include information retrieval, evolutionary computation, and stochastic optimization search strategy parameters.

12. The method of claim 10 wherein the average change in quality is subject to strength reduction of the noise contribution before advancing to the next (static parents) generation.

13. The method of claim 10 wherein the selection in step (d) is nondeterministic selection of stochastic optimization techniques of evolutionary computations.

14. The method of claim 10 wherein the center of mass recombination known in stochastic optimization techniques of evolutionary computations as nearest neighbor clusters (NNCs) is used in step (b).

15. The method of claim 10 for optimizing the shape of the dynamic search space comprised of any large-scale data set—information stored on a single computer, a local area network (LAN), and a wide area network (WAN) that encompasses the whole Internet—that may consists of constantly fluctuating information content over relatively short periods of time.

\* \* \* \* \*