# APPLICATION OF BIOMOLECULAR COMPUTING TO MEDICAL SCIENCE: A BIOMOLECULAR DATABASE SYSTEM FOR STORAGE, PROCESSING, AND RETRIEVAL OF GENETIC INFORMATION AND MATERIAL

John H. Reif, Michael Hauser, Michael Pirrung, and Thomas LaBean

*Departments of Computer Science, Ophthalmology, and Chemistry, Duke University, Durham, North Carolina*

A key problem in medical science and genomics is that of the efficient storage, processing, and retrieval of genetic information and material. This chapter presents an architecture for a Biomolecular Database system that would provide a unique capability in genomics. It completely bypasses the usual transformation from biological material (genomic DNA and transcribed RNA) to digital media, as done in conventional bioinformatics. Instead, biotechnology techniques provide the needed capability of a Biomolecular Database system without ever transferring the biological information into digital media. The inputs to the system are DNA obtained from tissues: either genomic DNA, or reverse-transcript cDNA. The input DNA is then tagged with artificially synthesized DNA strands. These "information tags" encode essential information (e.g., identification of the DNA donor, as well as the date of the sample, gender, and date of birth) about the individual or cell type that the DNA was obtained from. The resulting Biomolecular Database is capable of containing a vast store of genomic DNA obtained from many individuals (multiple army divisions, etc.). For example, the DNA of a million individuals requires about 6 pedabits ($6 \times 10^{15}$ bits), but due to the compactness of DNA a volume the size of a conventional test tube with a few milliliters of solution could contain that entire Biomolecular Database. Known procedures for amplification and reproduction of the

Address correspondence to: John H. Reif, Department of Computer Science, Room D223, LSRC Building, Duke University, Durham, NC 27708 (reif@cs.duke.edu).

resulting Biomolecular Database are discussed. The Biomolecular Database system has the capability of retrieval of subsets of stored genetic material, which are specified by associative queries on the tags and/or the attached genomic DNA strands, as well as logical selection queries on the tags of the database. We describe how these queries can be executed by applying recombinant DNA operations on the Biomolecular Database, which have the effect of selection of subsets of the database as specified by the queries. In particular, we describe how to execute these queries on this Biomolecular Database by the use of biomolecular computing (also known as DNA computing) techniques, including execution of parallel associative search queries on DNA databases, and the execution of logical operations using recombinant DNA operations. We also utilize recent biotechnology developments (recombinant DNA technology, DNA hybridization arrays, DNA tagging methods, etc.), which are quickly being enhanced in scale (e.g., output via DNA hybridization array technology). The chapter also discusses applications of such a Biomolecular Database system to various medical sciences and genomic processing capabilities, including: (a) rapid identification of subpopulations possessing a specific known genotype, (b) large-scale gene expression profiling using DNA databases, and (c) streamlining identification of susceptibility genes (high-throughput screening of candidate genes to optimize genetic association analysis for complex diseases). Such a Biomolecular Database system may provide a revolutionary change in the way that these genomic problems are solved.

## 1.  INTRODUCTION

### 1.1.  Motivation: The Need for a Compact Database System for Storing, Processing, and Retrieving Genetic Information/Material

Recent advances in biotechnology (recombinant DNA techniques such as rapid DNA sequencing, cDNA hybridization arrays [see also this volume, Part IV, chapter 1, by Meinhart and Wereley], cell sorters, etc.) have resulted in many benefits in health fields. However, these advances in biotechnology have also brought risks and considerable further challenges. The *risks* include the use of biotechnology for weaponry, for example, diseases (or environmental stresses) engineered to attack and disable military personnel. The *challenges* include the difficulties associated with the acquisition, storage, processing, and retrieval of individual genetic information. In particular, it is apparent that the sequencing of the human genome is not sufficient for many medical therapies, and one may instead require information about the specific DNA of the diseased individual, as well as information concerning the expression of genes in various tissue and cell types. In the scenario of biological warfare, such individual-specific information can be essential for therapies or risk mitigation (e.g., identification of individuals likely to be susceptible to a particular biological attack). To do this, there must be a capability to store this biological information, and also a capability to execute queries that identify individuals who contain certain selected subsequences in their DNA (or transcribed RNA). Hence, what is needed is essentially a database system capable of storing and retrieving biological material and information.

This biological information is quite data-intensive; the DNA of a single human contains about 6 gigabits of information, and the number of genes that potentially may be expressed may total approximately 30,000 (up to 15,000 genes may be expressed in each particular cell type, and there are thousands of cell types). The DNA of a single individual contains about $3 \times 10^9$ bases, which (with 4 bases) is $6 \times 10^9$ bits. The DNA of a million individuals (e.g., a large military force) therefore requires 6 pedabits (a pedabit is $10^{15}$ bits). The expression information for a few dozen cell types in each of a million individuals may also require multiple pedabits. Although the acquisition of such a vast DNA databank may be feasible via standard biotechnology, the rapid transfer of the DNA of such a large number of individuals into digital media seems infeasible, due to the tedious and time-consuming nature of DNA sequencing. Even if this large amount of information could be transferred into digital media, it certainly would not be compact: current storage technologies require considerable volume (at least a few dozen cubic meters) to store a pedabit. Furthermore, even simple database operations on such a large amount of data require vast computational processing power (if executed in a few minutes).

## 1.2. Overview of the Biomolecular Database System

This chapter presents the architecture of a Biomolecular Database system for the efficient storage, processing, and retrieval of genetic information and material. It completely bypasses the usual transformation from biological material (genomic DNA and transcribed RNA) to digital media, as done in conventional bioinformatics. Instead, biotechnology techniques provide the needed capability of a Biomolecular Database system, without ever transferring the biological information into a digital media. It may provide a potentially unique and revolutionary capability in genomics.

### 1.2.1. DNA: An Ultra-Compact Storage Media

The storage media of this database system is comprised by the strands of DNA, which are (in comparison to RNA) relatively stable and non-reactive: they can be stored for a number of years without significant degradation. In particular, the genetic information can be stored in the form of DNA strands containing fragments of genomic DNA as well as appended strands of synthesized DNA ("information tags") encoding information relevant to the genomic DNA. This Biomolecular Database is capable of containing a vast store of genomic DNA obtained from many individuals (e.g., multiple divisions of an army). We can provide the store with a redundancy (i.e., a number of copies of each DNA in the database) that ranges from a few hundred or thousand downwards to perhaps 10,

as the stringency of the methods increase. As mentioned above, the DNA of a million individuals contains 6 pedabits, but due to the compactness of DNA a volume the size of a conventional test tube can contain the entire Biomolecular Database. A pedabit of information can be stored (with tenfold redundancy) in less than a few milligrams of dehydrated DNA, or when hydrated may be stored within a test tube containing a few milliliters of solution.

### 1.2.2. Construction of the Biomolecular Database System

The inputs to the system are DNA obtained from tissues: either genomic DNA, or reverse-transcript cDNA obtained from mRNA expressed from the DNA of a particular cell type. The Biomolecular Database is constructed as follows:

a.  The input DNA strands are first fragmented, e.g., they may be partially digested into moderate-length sequences by the use of restriction enzymes. We describe a variety of methods for fragmentation protocols, and compare them by their distribution of strand lengths, and the predictability of the end sequences of the fragmented DNA.

b.  The DNA are then tagged with artificially synthesized DNA strands. These "information tags" encode essential information (e.g., identification of DNA donor, as well as the date of the sample, gender, date of birth) about the individual or cell type that the DNA was obtained from. These "information tags" are represented by a sequence of distinct DNA words, each encoding variables over a small domain. We describe and test tagging protocols based on primer extension and utilizing the predictability of the end sequences of the fragmented DNA.

### 1.2.3. Processing Queries in the Biomolecular Database System

The chapter then discusses how to execute queries on the resulting Biomolecular Database. The system makes use of biomolecular computing (also known as DNA computing) methods to execute these queries, including the execution of parallel associative search queries on DNA databases, and the execution of logical operations using recombinant DNA operations. We also describe the use of conventional biotechnology (recombinant DNA technology, DNA hybridization arrays, DNA tagging methods, etc.), for example, output is via DNA hybridization array technology.

These queries include retrieval of subsets of the stored genetic material, which are specified by associative queries on the tags and/or the attached genomic DNA strands, as well as logical selection queries on the tags of the database. These queries are executed by applying recombinant DNA operations on this Biomolecular Database, which have the effect of selection of subsets of the database as specified by the queries. We describe two distinct methods for processing logical queries: a surface-based primer-extension method, as well as a solution-based PCR method. Query processing is executed with vast molecular-level parallelism by a sequence of biochemical reactions requiring a length of time that remains nearly invariant with respect to the size of the database up to extremely large numbers (e.g., up to $10^{15}$). This is because the key limitation is the time for DNA hybridization, which is done in parallel on all the DNA. Output of the queries would be accomplished via DNA hybridization array technology.

### 1.2.4. Computer Simulations and Software

We describe computer simulations and software that can be used for the analysis and optimization of the experimental protocols. In particular, we describe the use of computer simulations for the design of hybridization targets for readout of information tags and SAGE tags by microarray analysis. We also discuss the scalability of these methods to do logical query processing within Biomolecular Databases of various sizes.

### 1.2.5. Applications

The chapter also discusses applications of a Biomolecular Database system to provide various genomic processing capabilities, including: (a) rapid identification of subpopulations possessing a specific known genotype, (b) large-scale gene expression profiling using Biomolecular Databases, and (c) streamlining identification of susceptibility genes: high-throughput screening of candidate genes to optimize genetic association analysis for complex diseases (see, e.g., this volume, Part III, section 6, on cancer). Such a Biomolecular Database system provides a revolutionary change in the way that these genomic problems can be solved, with the following advantages: (i) avoidance of sequencing for conversion from genomic DNA to digital media, (ii) extreme compactness and portability of storage media, (iii) use of vast molecular parallelism to execute operations, and (iv) scalability of the technology, requiring a volume that scales linearly with the size of the database, and a query time that is nearly invariant of that size. These unique advantages may potentially provide a number of opportunities for a variety of applications beyond medicine, since they also impact

defense and intelligence in the biological domain. Applications discussed include reasonable scenarios in (a) medical applications (e.g., in oncology, rapid screening, among a selected set of individuals, for expressed genes characteristic of specific cancers; see also preceding chapter 2 by Stolovitzky), (b) biological warfare (e.g., for biological threat analysis, rapid screening of a large selected set of personnel for possible susceptibility to natural or artificial diseases or environmental stresses, via their expressed genes), and (c) intelligence (e.g., identification of an individual, out of a large selected subpopulation, from small portions of highly fragmented DNA).

### 1.3. Organization of the Chapter

In this section we have provided a brief medical science motivation for a Biomolecular Database system, and a brief overview of the system. In §2 we briefly discuss relevant conventional biotechnologies and briefly overview the biomolecular computing (also known as DNA computing) field. In §3 we describe in detail our Biomolecular Database system. In that section we make use of various relevant biomolecular computing methods, including the use of word designs for synthetic DNA tags, execution of parallel associative search queries on DNA databases, and the execution of logical operations using recombinant DNA operations. In §4 we discuss a number of genomic processing applications of Biomolecular Database systems. In §5 we conclude with a review of potential advantages of Biomolecular Database systems.

## 2.  REVIEW OF BIOTECHNOLOGIES FOR GENOMICS AND THE BIOMOLECULAR COMPUTING FIELD

### 2.1. Conventional Biotechnologies for Genomics

There have been considerable commercial biotechnological developments in the last few decades, and many further increases in scale can reasonably be expected over the next five years. For example, the DNA hybridization array technology developed by Affymetrix Inc. (the capability is currently up to 400,000 output spots, and within 5 years a projected 1,000,000 outputs) can be adapted for output of queries to conventional optical/electronic media. Other biotechnology firms (e.g., Genzyme Molecular Oncology Inc.) have developed competing biotechnologies.

#### 2.1.1.  Genomics

In the research field known as *genomics*, there are a number of main areas of focus, each with somewhat different goals. These include:

1. **DNA sequencing**. *Sequencing* is the determination of the specific base pair sequence making up the DNA. This tells us all the possible genes that a given organism may express—its genetic makeup. In conventional bioinformat-

ics, it is generally assumed that the genes discussed have been previously sequenced and placed in a computer database.

2. **Gene expression analysis**. *Expression analysis* attempts to determine which genes are being expressed in a given tissue or cell type at a specific moment in time. The objective, to identify all the genes that are being expressed, is challenging because of the great complexity of the mixture of mRNA being analyzed—each cell may express as many as tens of thousands of genes. SAGE Tagging and cDNA hybridization arrays, as discussed below, are techniques for determining comprehensive gene expression data for a given cell type or tissue. The technique of differential expression analysis compares the level of gene expression between two different samples. Variations in the level of expression of individual genes or groups of genes can provide valuable clues to the underlying mechanism of the disease process. There are a number of methods currently used to obtain comprehensive gene expression data.

a. *cDNA hybridization arrays*. A cDNA hybridization array is composed of distinct DNA strands arrayed at spatially distinct locations. A cDNA hybridization array operates by hybridizing the array with fluorescent-tagged probes made from mRNA, which anneal to its DNA strands. This generates a fluorescent image-defining expression, which provides a very rapid optical readout of expressed genes. However, cDNA hybridization arrays are generally manufactured for use with a given set of expressed genes, for example, those of a given cell type. The design and manufacture of cDNA hybridization arrays for a given expression library of a size over 10,000 can be quite costly and lengthy. Affymetrix has recently developed an oligonucleotide array, known as a UniversalChip, that is not specialized to any gene library; it consists of 2000 unique probe sequences that exhibit low cross-hybridization and broad sampling of sequence space It can be used with fluorescent-tagged probes made from DNA rather than mRNA. This technology can be used for output in a Biomolecular Database system.

b. *Serial analysis of gene expression* (SAGE) is a technique for profiling the genes present in a population of mRNA. By the use of various restriction enzymes, SAGE generates, for each mRNA, a 10-base tag that usually uniquely identifies a given gene. In the usual SAGE protocol, the resulting SAGE tags are blunt-end ligated together and the results are sequenced. The sequencing is faster than sequencing the entire set of expressed genes because the tags are much shorter than the actual mRNA they represent. Once sequencing is complete, the tag sequences can be looked up in a public database to find the corresponding gene. Using the sequence data and the current UniGene clusters, a computer processing stage determines the genes that have been expressed. SAGE can be used on any set of expressed genes and it is not specialized to a particular set. This technology can be adapted for use for additional information tags appended to the DNA in a Biomolecular Database. Genzyme Molecular Oncology Inc. is the developer of this SAGE technology.

c. *Differential expression analysis* is a technique for finding the difference in gene expression, for example, between two distinct gene types. Lynx Therapeutics Inc. has developed a randomized tagging technique for differential expression analysis. The randomized tagging techniques of Lynx Therapeutics may be adapted to determine the difference between two Biomolecular Database subsets.

## 2.2. Relevant Biomolecular Computing Techniques: Biomolecular Computing

In the field known as *Biomolecular Computing* (and also known as *DNA Computing*), computations are executed on data encoded in DNA strands, and computational operations are executed by use of recombinant DNA operations. Surveys of the entire field of DNA-based computation are given in (50,52).

The first experimental demonstration of Biomolecular Computing was offered by Adleman (1), who solved a small instance of a combinatorial search problem known as the Hamiltonian path problem. Considerable effort in the field of Biomolecular Computing methods has been made to solve Boolean satisfiability problems (SAT) problems, that is, the problem of finding Boolean variable assignments that satisfy a Boolean formula. Frutos and colleagues (24), Faulhammer et al. (23), and Liu et al. (37) applied surface-chemistry methods and Pirrung et al. (47) improved their fidelity. Adleman's group (11) recently solved an SAT problem with 20 Boolean variables using gel-separation methods. While the 20 Boolean variables size problem is impressive, Reif (52,53) has pointed out that the use of Biomolecular Computing to solve very large SAT problems is limited to at most approximately 80 variables, so is not greatly scalable in terms of number of variables.

The use of Biomolecular Computing to store and access large databases, in contrast, appears to be a much more scalable application. Baum (7) first discussed the use of DNA for information storage and associative search; Lipton (36) and Bancroft and coworkers (6) also discussed this application. Reif and LaBean (54) developed and Reif et al. (55) experimentally tested the synthesis of very large DNA-encoded databases with the capability of storing vast amounts of information in very compact volumes. Reif et al. (55) tested the use of DNA hybridization to do fast associative searches within these DNA databases. Reif (50) also developed theoretical DNA methods for executing more sophisticated database operations on DNA data, such as database join operations and various massively parallel operations on the DNA data. Gehani and Reif (27) investigated methods for executing DNA-based computation using microfluidics technologies. In addition, Gehani et al. (28) describe a number of methods for DNA-based cryptography and countermeasures for DNA-based steganography systems as well as discuss various modified DNA steganography systems that appear to have improved security. Kashiwamura and colleagues (34) describe the use of nested PCR to do hierarchical memory operations.

Suyama et al. (60) and Sakakibara and Suyama (59) have developed biomolecular computing methods for gene expression analysis. Garzon and colleagues (25) recently analyzed the efficiency and reliability of associative search in DNA databases, and Chen and colleagues (13) discussed DNA databases with natural DNA based on the prior work of Reif et al. (55) and the present work.

### 3. A BIOMOLECULAR DATABASE SYSTEM

### 3.1. Overview

The inputs to the system are natural DNA obtained from tissues: either genomic DNA or reverse-transcript cDNA obtained from mRNA expressed from the DNA of a particular cell type. A short piece of synthetic DNA is added to each natural DNA strand. This piece of synthetic DNA, called an information tag, is used to code information about the original piece of DNA. This information can include the age or gender of the person from whom the DNA came, or the clinical symptoms of individuals suffering from a disease. In a typical application, the Biomolecular Database consists of a mixture of DNA strands from many different people (or other organisms). This Biomolecular Database system is capable of storage, processing, and retrieval of genetic information and material. Individual molecules of DNA in the Biomolecular Database can be selected and removed from the mixture on the basis of the information encoded in their information tag. This chapter describes several innovative biological applications for Biomolecular Databases; in particular, we discuss the application of our Biomolecular Database system to a number of genomic information processing applications.

### 3.2. Biological Inputs

The inputs to the system are DNA obtained from tissues. This input DNA is typically either (i) genomic DNA, or (ii) reverse-transcript cDNA obtained from mRNA expressed from the DNA of a particular cell type. (To ensure stability and non-reactivity, we suggest that the database be composed of DNA rather than RNA.)

### 3.3. Preprocessing the DNA

Biochemical operations can be used to partially digest the DNA by restriction enzymes (ensuring the resulting DNA strands are of modest size), and then label the resulting genomic DNA fragments with synthetic DNA information tags.

### 3.3.1. Fragmentation of the input DNA

The creation of a Biomolecular Database must involve some degree of fragmentation of genomic DNA. While this may at first seem a very simple step, it is in fact critical to later processing that this fragmentation step be done in a highly controllable way. We describe several methods to produce DNA strands of the desired length. The methods are required to produce a predictable distribution of lengths, and to ensure that at least one of the resulting ends has a defined sequence.

a. *Mechanical shearing*. This is a method that produces a certain size distribution; however, it is not so useful in our context since the resulting ends have undefined sequences.

b. *Reagent-less methods to create breaks*. Pirrung, Zhao, and Harris (48) developed a nucleoside analogue whose backbone can be cleaved by long-wavelength UV light, and specific photocleavable T analogues could be used (analogous to the dUTP method). However, again it is not so useful in our context since the resulting ends have undefined sequences.

c. *Controlled digestion of high-MW DNA by DNAse I*. This is another method that can be used to produce DNA of a specific size range. It relies on careful monitoring of reaction progress and does not produce specific sequences at the ends of the fragments to enable ligation or PCR processes.

d. *Digestion of DNA with restriction endonucleases*. This offers the advantage that known sticky ends are generated.

e. *"Rare Cutting" endonucleases*. These can be used to produce DNA fragments of larger size. The recognition sequence of such enzymes is as large as 8 bp, meaning that, on average, DNA is cut to $1/(0.25^8)$ or 65 kb. In many situations fragments larger than 65 kb may be desired; for example, complicated loci with many introns might comprise as much as 100 kb, and that is just for one gene.

f. *PCR methods for fragmentation*. One attractive alternative is to use PCR. Random-primed PCR has been used to amplify the whole genome of a single sperm (4,35,67). The challenge in using this strategy is to create long amplicons. In principle, amplicon size in random-primed PCR is a function only of the average distance between two inward-facing hybridized primers, which is then a function only of the primer concentration and temperature. Modest flexibility exists in the hybridization temperature in PCR, so a fruitful strategy to make long amplicons is to lower the primer concentration. In order to efficiently amplify with a low primer concentration, the primers should have a high melting temperature ($T_m$). Increasing length and G/C content increase primer $T_m$. Random-primed PCR can therefore be examined with novel conditions and primer designs to maximize the amplicon length. Lengthening the random primers by oligonucleotide synthesis is straightforward. Making the primers G/C-rich is challenging, as G/C-rich templates are known to be more difficult to amplify

owing to increased secondary structure. Substitution of *nonstandard bases* such as deazaG for G alleviates this difficulty.

g. *UV-sensitive nucleoside analogues in cell growth for fragmentation.* Another approach is to grow immortalized lymphoblast cell lines in the presence of UV-sensitive nucleoside analogues Pirrung et al. (48). These analogues can be incorporated into the DNA of the cells, which could subsequently be cleaved by exposure to UV light. The concentration of the analogues determines the frequency of their incorporation, and the size of the resulting fragments.

## 3.4. Creation of the Tagged Biomolecular Database

These DNA tags are composed of a concatenation of short subsequences, which encode scalar data values. For example, the information tags may contain the individual's unique ID and the cell type (in the case of reverse-transcript cDNA obtained from the RNA expressed by a particular cell type) of the genomic DNA and may also encode other useful information (e.g., sex and birthdate of the individual). The tagging can be done using known methods, for example, a primer-extension reaction, using the fact that one of the ends of the genomic cDNA can be predicted by the use of the appropriate initial fragmentation process, and further designing the tags with ends complementary to these sticky ends resulting from the fragmentation process. The resultant database elements have tags on each 5'- and 3'-end. This can be done so that each Biomolecular Database strand bears a universal amplification (primer) sequence at the extreme 5'- and 3'-ends.

## 3.5. DNA Word Design for the Information Tags

A key problem is the design of a lexicon of short DNA sequences (DNA words) for the information tags in our Biomolecular Database. (Our DNA "information tag" sequences are in general a subset of such a lexicon). Careful word design is crucial for optimizing error control in the queries executed within the Biomolecular Database. Good word design can be used to minimize unwanted secondary structure and to minimize mismatching by maximizing binding specificity. There are conflicting requirements on word design: as strand length decreases (which is desirable), the difference between distinct information words decreases (not desirable). Prior work in DNA word design includes a four-base mismatch word design used for surface-based DNA computing (29), and Frutos et al. (24) showed that surface morphology may be an important factor for discrimination of mismatched DNA sequences. A three-base design was used by Cukras and coworkers (17). Evolutionary search methods for word designs are described by Deaton et al. (18). Other DNA word designs are de-

scribed by Baum (8), Deaton and colleagues (19), Mir (39), and Garzon et al. (25). Laboratory experiments of word designs are described by Kaplan et al. (33), and ligation experiments are described by Jonoska and Karl (32). Wood (66) considers the use of error-correcting codes for word design and to decrease mismatch errors. One can utilize and improve on these methods for DNA word design, including evolutionary search methods and error-correcting codes. Hartemink and colleagues (30) describe an automated constraint-based procedure for nucleotide sequence selection in word design. In designing the DNA tags used in the database, one needs to determine how many residues should be used for each data block of the tag sequence (the tag sequence on the database strands binds to the probe sequence on the query strand), and then decide how many words are required at each block position (determined by the number of values available to the variable).

The range of possible sentences entailed by a word-block construction scheme is shown in Figure 1. For each block position in the sequence one word is chosen from the word set and synthesized on the growing DNA strand. Separate reaction vessels are used for each word in the block, so that all word choices are utilized but only one is present on any particular strand. For example, the arrows indicate the trace that results in the sentence: word1A–word2D–word3A–word4B. A particular bead is drawn through a particular path in the set of possible word choices, but all possible paths are populated with beads, so all possible DNA sentences are synthesized. Each bead contains multiple copies of a single DNA sequence that can be synthesized by the well known mix-and-split synthesis scheme. Figure 2 shows the scaling of library diversity with increasing sentence length (block count) and increasing number of available words within each block. Diversity is calculated by raising the word count to the exponential power given by block count (i.e., diversity = [word count]$^{\text{block count}}$). As a simple example, to achieve a total diversity greater than one million with sentences containing 6 blocks, for example, one would require a set of 10 word choices per block. Also, to achieve a total diversity of $12^{14}$ with sentences containing 14 blocks, for example, also requires a set of 12 word choices per block. In designing a DNA-encoded database one must consider several important factors including the following. (i) The overall length of the oligonucleotide sequences used for matching is critical because sequence length directly affects the fidelity and melting temperature of DNA annealing. (ii) The Hamming distance (the number of changes required to morph one sequence into another) is another critical consideration. One would like to maximize the Hamming distance between all possible pairs of encodings in the database in order to minimize near-neighbor false-positive matching. One strategy for maintaining sequence distance is to assign block structures to the sequences with sets of allowed words (subsequences) defined for each block. (iii) Another important consideration is the choice of the words themselves and the grouping of words into sets within the blocks. Sentence length, desired library diversity, and word-pair distance
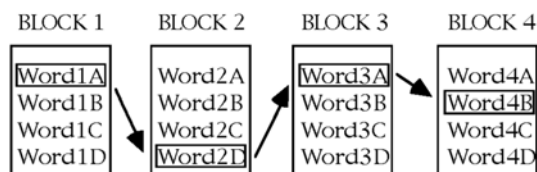
**Figure 1**. Range of possible sentences entailed by a word-block construction scheme.
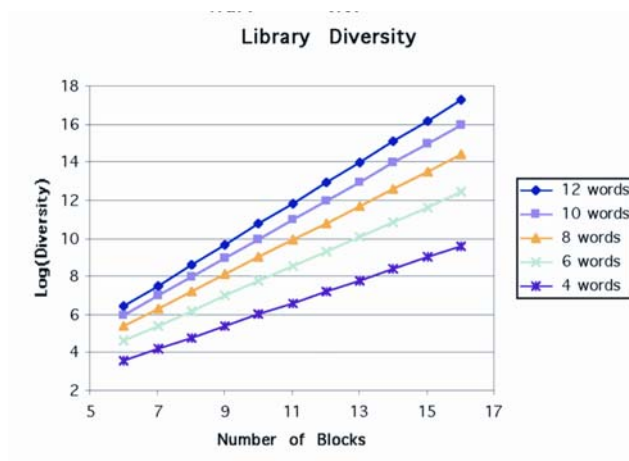


**Figure 2**. Scaling of library diversity with increasing sentence length (block count) and increasing number of available words within each block.

constraints all affect the choices of words in the lexicon. Word design can be maximized by careful design of words, the lexicon, and database elements, as well as experimental tuning of annealing conditions (e.g., temperature-ramp rate, pH, and buffer and salt concentrations). A useful tool in this task is the computer simulations of DNA hybridization, known as BIND, developed by Hartemink et al. (30).

## 3.6. Additional Tagging Methods for the DNA Strands

Various sophisticated tagging techniques have been developed by the biotechnology industry for expression analysis and differential expression analysis..

These include the SAGE tagging of Genzyme Molecular Oncology Inc. and the randomized tagging techniques of Lynx Therapeutics Inc.

a. *Serial analysis of gene expression (SAGE)* is a technique developed by *Genzyme Molecular Oncology Inc.*, for profiling the genes present in a population of mRNA. By the use of various restriction enzymes, SAGE generates, for each mRNA, a 10-base tag that usually uniquely identifies a given gene. In the standard SAGE protocol, the resulting SAGE tags are blunt-end ligated and the results are sequenced. Such sequencing is faster than sequencing the entire expressed genes because the tags are much shorter than the actual mRNA they represent. Once sequencing is complete, one may look up the tag sequences in a public database to find the corresponding gene. Using the sequence data and the current UniGene clusters, a computer processing stage determines the genes that have been expressed. SAGE can be used on any set of expressed genes; it is not specialized to any particular set. This technology can be adapted for use as additional information tags appended to the DNA in our database.

b. *Differential expression analysis* is a technique developed by Lynx Therapeutics Inc. for finding the difference in gene expression, for example, between two distinct cell types. The randomized tagging techniques of Lynx Therapeutics Inc. can be adapted to determine the difference between two DNA database subsets.

c. *Hybrid methods*. One can modify these methods and extend them to apply to the tagged DNA strands of our database. This requires considerable changes in the protocols, due to unwanted hybridization that may occur as a result of combination of synthetic tags with genomic DNA in our database strands. However, these modified methods can provide further powerful capabilities, for example, the capability for fingerprinting (creating short DNA tags that are nearly unique IDs for longer DNA strands of the database), identification of expressed genes of selected DNA strands, and also the capability for differential expression analysis of distinct selected subsets of the Biomolecular Database.

### 3.7. Amplification and Reproduction of Biomolecular Databases

Once a Biomolecular Database is created, it is important to be able to accurately replicate it, as it may be consumed during the course of interrogation. Prudence suggests maintaining each database in an archive, and querying only daughter databases prepared from the archival forms. Since each database member is designed to bear a universal amplification (primer) sequence at the extreme 5'- and 3'-ends, database replication can be performed using PCR. Because the length of the DNA strands in the database might be quite substantial, including both biological DNA information and many flanking tag sequences, the ability to produce full-length amplicons with long templates is

crucial to maintaining the fidelity of the databases. "Long accurate" PCR techniques (62,63), using novel thermostable proofreading polymerase enzymes such as *Pfu*, are currently capable of amplifying loci of up to ~40 kb. While powerful, the database design should not be limited to this length by the method of database replication, and it may be easier to enable PCR to produce amplicons of somewhat longer length. One simply needs to enhance by a moderate multiple the (amplicon) length that can be reliably amplified.

Optimized choice of amplicons may be achieved by exploiting two principles: experimental design (9,10,21) and combinatorial chemistry (44,45). Continuous variables that affect PCR reactions include the temperature of the initiation (hot start), annealing, extension, and dissociation steps, and the concentration of buffer components, additives, nucleotides, primers, and template. These variables compose a multidimensional space. A pervasive challenge in science and technology is identifying specific values for each parameter affecting multivariable processes that result in globally optimum performance and avoid local maxima. Commercial software enables the design of experiments that much more reliably and quickly lead to the global optimum. Noncontinuous variables that affect PCR reactions include the identity of the template, primers, and polymerase. An optimum combination of these molecules can be found only by systematic screening for each. For tractable numbers of combinations, all can be examined explicitly. When the diversity space expands beyond that domain, "indexing" techniques are available that permit optimum performers to be identified even when in a mixture with lower performers (46). A selection of variable-length primers can be examined, including those incorporating modified bases (deazapurines, 2'-OMe RNA) that suppress primer consumption by dimerization. A selection of commercial polymerase enzyme systems can be examined, including MasterAmp™ Taq, ThermalAce,™ Advantage-Tth,™ AdvanTaq™, and KlenTaq™/Pfu. A selection of templates should be examined, including whole viral genomes, bacterial artificial chromosomes (BACs), yeast artificial chromosomes (YACs), and the smallest yeast chromosome (225 kb). Analysis of the products of these reactions is challenging due to shearing of large DNA molecules by conventional sieving matrices. Pulsed-field gel electrophoresis (12,41) can therefore be used with amplicons of this size.

## 3.8. Associative Search in Biomolecular Databases

### 3.8.1. DNA-Based Associative Search

Eric Baum (7) first proposed the idea of using DNA annealing to do parallel associative search in large databases encoded as DNA strands. The idea is very appealing since it represents a natural way to execute a computational task in massively parallel fashion. Moreover, the required volume scales only linearly

with the database size. However, there were further technical issues to be resolved. For example, the query may not be an exact match with any data in the database, but DNA annealing affinity methods work best for exact matches. Reif and LaBean (54) described improved biotechnology methods to do associative search in DNA databases. These methods adapted some information processing techniques (error-correction and VQ coding) to optimize input and output (I/O) to and from conventional media, and to refine the associative search from partial to exact matches.

Reif and colleagues (55) developed and then experimentally tested a method for executing associative searches in DNA databases of encoded images, and this method was tested using an artificially synthesized DNA database. Prior to that project, the idea of using DNA annealing to do parallel associative search in synthetic DNA databases had never been experimentally implemented. They detailed a study involving the design, construction, and testing of large databases for storage and retrieval of information within the nucleotide base sequences of artificial DNA molecules. The databases consisted of a large collection of single-stranded DNA molecules that was immobilized on polymer beads. Each database strand carried a particular DNA sequence, consisting of a number of sequence words drawn from a predetermined lexicon. They made a number of experimental databases of artificially synthesized DNA sequences designed for encoding digital data, scaled in increasing sizes. Each DNA strand of the database is single stranded, and encodes a number that provides the index to the database element. They used an extensive computer search for the design of the DNA word libraries, to ensure a significant Hamming distance between distinct words and allow for annealing discrimination. They constructed their largest synthetic databases in two phases. In the first phase they constructed an initial DNA database by combinatorial, mix-and-split methods on plastic microbeads. This constituted by far the largest artificially constructed synthetic databases of this sort. The next phase was the development of a construction method for much larger synthetic databases by combining pairs of the synthetic database strands so as to square the size of the database to approximately $10^{15}$ distinct data elements (each represented redundantly by over 10 identical strands of DNA). Even with this greater than tenfold redundancy, the DNA database using this construction method is extremely compact, and requires only 10 milligrams of DNA.

### 3.8.2. Associative Search via PCR

PCR methods can be used for associative search queries in Biomolecular Databases (in particular, on the words of the tagged portions of the Biomolecular Database strands), using known and modified PCR techniques previously developed by Reif and coworkers (55). They describe experiments for executing

associative search queries within the above described synthetic DNA databases. Associative search queries were executed by hybridization of a database DNA strand with a complementary query strand. Discrimination in annealing experiments was enhanced by the library design, which guaranteed a minimum Hamming distance between distinct sequences. In their initial annealing experiments for processing associative search queries, they employed fluorescently labeled query strands and then performed separation of fluorescent versus nonfluorescent beads by Fluorescence-Activated Cell Sorting (FACS). They also experimentally tested variants of conventional PCR techniques for executing associative search queries, and, in addition, developed a PCR technique for associative search in the pairwise constructed DNA database.

### 3.8.3. Analysis of Associative Search

Similar error analysis and experimental testing methods can be employed in our proposed generalizations of this prior work (55) to tagged genomic DNA. It would be informative to measure rates of various search errors including: false positives from near-neighbor mismatches, partial matches, and nonspecific binding, as well as false negatives from limit-of-detection problems. It is desirable to directly measure the limits of detection, and to measure the ability to retrieve rare sequences within databases of high strand diversity.

### 3.9. Logical Query Processing in Biomolecular Databases

Biochemical operations can be used to execute query operations on this Biomolecular Database, so as to retrieve subsets of the Biomolecular Database. Each of the information strands of the database encodes a sequence of data values $v_1$, $v_2$, ..., $v_k$, where the $i$th value $v_i$ ranges over a small finite domain $D_i$ (e.g., $D_i$ typically would range over 10 or less possible values, each encoded by a distinct fixed-length DNA sequence). Retrieval can be specified by logical queries on the tags of the database as well as associative queries on the attached genomic DNA strands. The associative searches can be executed by recombinant DNA operations, for example, variants of PCR combined with surface-chemistry methods and/or solution-based methods. The logical queries include the following: (i) SELECTION—selects DNA strands of a given ID or cell type; and (ii) Logical SELECTION—executes logical queries that select those genomic DNA strands whose information strands satisfy a specified logical query formula, whose logical conjunctives include AND as well as OR. These logical conjunctives are applied to selective predicates of the form "Tag($i$) = $v$", where Tag($i$) is the $i$th portion of the information tag of a DNA strand of the database, and $v$ is a fixed value over the domain $D_i$. (The Boolean NOT of a selective

predicate of the form "Tag$(i) = v$" is not applied directly (since PCR and similar methods do not allow this) by instead applying the OR of selective predicates of the form "Tag$(i) = u$" for all possible $u$ in $D_i$ that are not equal to $v$). These selection operations can be executed by the use of recombinant DNA operations, applying and improving on logical processing methods developed in the field of DNA computing. Furthermore, one can provide the additional operation of selective amplification of the DNA populations. If these amplification operations are also executed, the logical selection and amplification operations result in a test tube whose selected DNA is vastly amplified. After the amplification process is completed, the output strands should vastly predominate all other strands of the Biomolecular Database. Other database operations that can be implemented by biochemical operations include database unions and limited joins (50).

### 3.9.1.  Scalability of Our Query Processing

These operations can be executed in a scalable way. The required volume never grows significantly; the volume is a fixed linear function of the number of elements of the database. (The constant multiple here is the degree of redundancy with which DNA strands are used to store database elements; we expect that one can allow between a few hundred and possibly as few as 10 DNA strands to encode a given database element.) The number of required DNA hybridization steps grows only linearly with the size of the query formula. So the time for executing a query grows just linearly with the length of the query formula, which in practice is of very modest size (as compared to the size of the database, which can be enormous)—say 20 or so variables. Hence the key time limitation is the time for DNA hybridization. But DNA hybridization time is nearly invariant of the size of the database even if the hybridization is execution on an enormous amount of DNA (up to extremely large database sizes, say $10^{15}$). However, there are considerable technical challenges in the design of the protocols—for example, biological data strands may be originally dsDNA while search protocols would function best with ssDNA (hence the protocols need to either form ssDNA or be modified appropriately). A key additional technical challenge in scaling the technology is the scale and number of resulting molecular biology reactions, requiring many tedious laboratory steps, particularly in the case of extremely large database sizes. This can be addressed by subsequent automation. We discuss two distinct methods for logical query processing: the first uses primer-extension techniques on solid support, previously developed (47) to solve SAT problems, and the other uses solution-based PCR amplification techniques. The second has greater potential for scalability due to the fact that it is solution based (so the chemistry operates in 3D, rather than being con-

strained to a surface). In both cases, one can apply DNA hybridization array technology for output of query results.

### 3.9.2. Executing Queries into the Biomolecular Database via Primer-Extension Techniques on Solid Support

A number of DNA computing researchers have developed microarray methods for DNA-based computing, exploiting the high fidelity of the primer-extension reaction to detect complementarity between primer libraries of all solutions to SAT problems and logical queries as templates (see, e.g., the work of Faulhammer et al. (23) and Liu et al. (37), and also that of Pirrung et al. (47), which improved fidelity). Primer extension is a two-step process, involving first annealing of a template molecule to the primer, the efficiency of which is directly related to sequence complementarity throughout the primer/template complex (see Figure 3). Second, a polymerase enzyme binds to the primer/template complex and adds a nucleotide or nucleotides complementary to the base
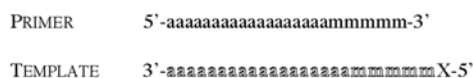
PRIMER     5'-aaaaaaaaaaaaaaaaaammmmm-3'

TEMPLATE   3'-aaaaaaaaaaaaaaaaaammmmmX-5'

**Figure 3**. A primer and template used in the primer-extension method for logical queries.

t         h         e              b         a         s         e
 X, the first unpaired base at the 5'-end of the template, only when there is a perfect match in the last portion of base pairs of the primer/template complex. It is important to emphasize that while primer extension was in this case performed on a DNA microarray, the elementary step of a polymerase chain reaction (PCR) is also a primer-extension process and thus is subject to the same stringent sequence requirements. The variables (primers) in the SAT computation of (47) were composed of two portions, which can be considered the message (the last few bases "m" at the 3'-end of the oligonucleotide) and the address (a sequence of bases "a" at the 5'-end). With the base-4 encoding of DNA, a message sequence is capable of encoding 10 Boolean variables. For the experiment of (47) all addresses were the same, as only one SAT problem was being addressed. However, this need not be the case. Using similar designs, it is possible to design up to approximately 20 blocks of distinct address sequences, which con-

catenated form the tag. Each of these blocks of distinct address sequences should exhibit no cross-hybridization under stringent conditions (a Hamming distance of at least 5), thereby enabling independent encoding and primer extension and therefore interrogation of up to 20 distinct attributes, each with up to 10 scalar values.

Use of the primer-extension method for logical queries into the Biomolecular Database is most efficient if performed in solution rather than on a microarray. This creates a challenge in terms of product detection and identification. The following method enables both to be accomplished. An example is presented for the fate of one molecule, though it is appreciated that all molecules in the library are subjected to the same process in parallel. The database member is a DNA molecule that has been created by the methods described earlier, with a biological DNA sequence flanked by one or more created tag sequences, which are to be the templates in a primer-extension reaction. The "bottom" strand is interrogated in this example. Primers with the following structure (shown in expanded form below the database element in Figure 4) are created to interrogate each tag. Complements to the address sequence in the tag/template are the same in each primer. Also common to each primer is a "barbed tail" in the form of a 5'-psoralen group. The irradiation of psoralens with long-wavelength UV is widely used to cross-crosslink duplex DNA (31,43,64). The message sequence
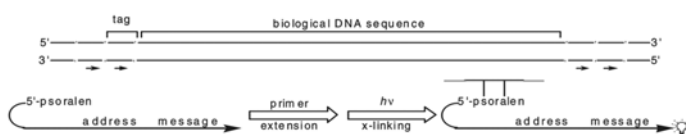


**Figure 4**. Obtaining a bottom strand in single-strand form, hybridized to a cDNA microarray, via photochemical cross-linkage of psoralen.

must be unique to each variable value, meaning that up to 10 primers are prepared per variable. The primer is also designed to address a unique X base in the tag/template to be interrogated. The primer-extension reaction is performed using a dideoxynucleotide terminator complementary to X and bearing a fluorescent dye with a unique and readily imaged emission spectrum. The dye color is specific to the variable, with the same dye/terminator being used for all interrogations of that variable. Multiple tags can be interrogated simultaneously because their dyes are different. The challenge at this stage is to read out the tags (based on the color(s) of the incorporated fluorescent dye(s)) in the context of the biological DNA. While the primer is still bound to the template, the psoralen

is photochemically cross-linked to the bottom strand of the library member, preserving the color of the dye. The bottom strand can then be obtained in single-strand form, which is hybridized to a cDNA microarray (see Figure 4). The color(s) of the array element complementary to the biological DNA identify the outcome of the queries of the tag sequences connected with it.

This concept could be applied in a similar fashion to a sequential (nested) PCR process by omitting the terminators and psoralen and providing one primer for the top strand and one for the bottom strand in each PCR. The eventual production of a full-length amplicon is dependent on the complementarity of *each* of the primers (logical AND) with its cognate tag sequence. This approach lends well to the use of DNA hybridization array technology for output of query results, providing distinct special locations for distinct outputs.

Another approach for executing Boolean queries on a Biomolecular Database is to use the gel separation-based method for SAT from Braich et al. (11), who succeeded in solving a 20-variable Boolean satisfiability problem. Although the queries would be executed on the tag portions of the DNA strands of the database, it is not clear how the efficiency of these separation methods would be affected by the genomic portion of the DNA strands in the databases.

### 3.9.3. Executing Queries in the Biomolecular Database via PCR Amplification Techniques

Another approach for Logical Query Processing is to use a variant of PCR amplification. The goal of this query processing is to selectively amplify only those DNA sequences (the *output strands*) whose information tags satisfy a given logical query. After the amplification process is completed, the output strands would vastly predominate all other strands of the Biomolecular Database.

### 3.9.4. Initialization Before Logical Query Processing

First, operations are executed that generate, from each DNA strand in the database, a new strand containing a concatenation of multiple copies of the Watson-Crick complement of the original strand. This can be done by a known sequence of routine recombinant DNA operations known as rolling circle replication (38). This begins by a circularization of each strand on the database, and then a primer-extension reaction on the circularized strand that repeatedly replicates the complement of the DNA strand to form a repeated sequence, followed by denature and separation of the result. The length of the resulting DNA strands is predictable (via the time duration and various parameters, including temperature) only to a degree, but it is predictable enough to allow us to con-

struct strands that we *expect* to have at least a given repeat length (as required by the below protocol). Recall that we have assumed that each database DNA strand is also redundantly represented by a number (ranging from up to a few thousand down to as few as 10) of identical DNA strands. This redundancy aids us here, since this initialization procedure results in a Biomolecular Database, where most of the redundant strands are lengthened by at least the given multiplying factor.

Figure 5 illustrates a scheme for processing genomic DNA into this database format which might include the following steps: (i) Cleave dsDNA into manageable pieces. (ii) Append prefixes to both ends of both strands. Heat denature dsDNA. Anneal to circularizing oligo. (iii) Ligate ssDNA circles. (iv) DNA polymerase reaction with circular templates to produce linear ssDNA containing multiple concatenated database entries. Note that the process of converting the DNA into database format may have unintended effects on the representation of entries in the database due to uneven amplification. Artificial bias may take the form of variations in the number of copies present on the average strand (distribution of strand lengths) or differences in the number of strands present for a given database entry. These protocols need to be optimized to take into account these possible affects.
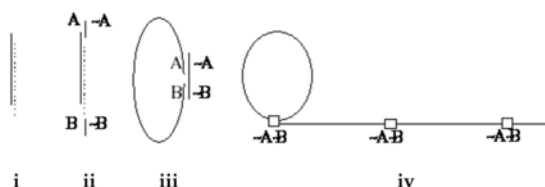


**Figure 5**. Scheme for processing genomic DNA into database format.

Multiple copies of the database entry are required on a single ssDNA strand, so that when Boolean variables recorded in the prefixes (A and B in Figure 5) are queried by primer binding and PCR, information recorded farther out toward the ends of the strand is not lost by failure to be copied (PCR only amplifies sequence physically between primer binding sites). The goal is to keep at least a few copies of the prefix information internally within the database strands so that information is not lost to subsequent rounds of query.

*3.9.5. Logical Query Processing Using Repeated PCR Operations*

We assume that the logical query is presented as the logical AND of a list of *K* logical clauses (each clause needs to be satisfied), where each clause consists

of the logical OR of a list of literals (the literals can be Boolean variables or their negation), one of which needs to be satisfied. Each clause $C$ in the formula is processed in turn, selectively amplifying only those DNA strands whose Boolean variables satisfy at least one literal of that given clause. To do so, one adds PCR primers encoding the literals of that clause $C$ and their Watson-Crick complements. Then a series of primer-extension reactions are executed that replicate only those DNA strands (or their Watson-Crick complements) that have subsequences that encode one of the literals of clause $C$. This process, applied as a series of PCR cycles, thus amplifies only those DNA strands whose Boolean variables satisfy at least one literal of that given clause, so that they vastly predominate all other strands of the Biomolecular Database. (**A technical note**: on each cycle, the amplified strands undergo loss of the material prefixing the primer's location, but the initial step of concatenating to each DNA strand in the database multiple copies of the strand ensure that is not a problem.) After the process is completed for each of the clauses in turn, the output strands that satisfy all the clauses would vastly predominate all other strands of the Biomolecular Database. This method for processing a logical query in the database is exquisitely sensitive: to get a result, one requires that the initial database have no more than 10 identical strands of DNA that satisfy the query. Again, DNA hybridization array technology can be used for output of query results, providing distinct special locations for distinct outputs.

### *3.9.6. Scalability*

As discussed above, our query processing is executed with vast molecular-level parallelism by a sequence of biochemical reactions requiring a time that remains nearly invariant of the size of the database up to extremely large database sizes (e.g, up to $10^{15}$). This is because the key limitation is the time for DNA hybridization, which is done in parallel for all the DNA.

### 3.10. Management of Errors

The logical and associative searches used to select specific molecules and sets of molecules from Biomolecular Databases are not 100% specific or effective. There may be several different kinds of errors: false negatives (appropriate DNA strands are present but not selected, either because of a lack of sensitivity or depletion of the relevant sequences from the database), false positives (inappropriate DNA strands are selected along with desired strands), errors based on degradation of the Biomolecular Database contents, and errors resulting from poorly designed queries, based on incomplete understanding of complex biological parameters. These kinds of errors can affect the results of the applica-

tions described here. For example, false negative errors can prevent finding existing individuals with the desired genetic variant. This is less serious than a false positive result, which could lead to sending a nonresistant individual into a contaminated area, under the false belief that he is genetically protected from a biological agent. A similar type of error could arise from database degradation (as, for example, from repeated error-prone duplication of the database). This type of error can be easily eliminated by follow-up confirmatory screening of that single individual's DNA. In general, it would be best to use the Biomolecular Database for very powerful and rapid selections based on genetic information, but then to confirm all results on individual DNA samples. This would require maintenance of individual stocks of DNA for each individual. This is a relatively large task, but well within current technology. An LIMS (laboratory information maintenance software) system and robotic liquid handling capacity are musts for this type of storage. There could also be errors of magnitude. These errors result from preferential amplification of one DNA strand over another. This kind of error is particularly troublesome, for it would skew allele frequencies. It may be necessary for us to monitor the frequency and extent of such errors and develop Boolean search strategies that minimize them. The final type of error, based on an incomplete understanding of the human genome, can only be rectified by continued research in other fields. This type of error could result from incomplete knowledge of the way in which genetic variants are distributed among different racial and ethnic groups. For example, the well-described *ccr5* variant that prevents HIV infection has been detected to date only among white males. If one were to select for nonexistent African American females expressing this variant, one might well obtain a small number of false positives. This type of error could also arise from mistaken assumptions. A given genetic variant might protect Hispanic females from infection by a given biologic agent, but oriental males carrying the same variant might be fully susceptible to infection because of another independent genetic variant.

### 3.11. Computer Simulations

Reif et al. (52,53) have made computer simulations of their methods for DNA-based associative search. They constructed computer software (viewable on the web) that provide a simulation of the entire experimental process, including conversion of this attribute database into a DNA database using DNA chips, the PCR method for associative search in the DNA database (using a software simulation of the kinetics of DNA hybridization), and, finally, conversion of the result of this query (using extensions of techniques described in (55)) into conventional media by use of a DNA expression array. Our computer simulation software to the above-described query processing provides a basis for future

software to simulate and optimize experimental protocols for query processing in Biomolecular Database systems.

## 4.  APPLYING OUR BIOMOLECULAR DATABASE SYSTEM TO EXECUTE GENOMIC PROCESSING

There is tremendous potential to apply Biomolecular Databases to the solution of a number of biological problems. The huge amount of data provided by the sequencing of the human genome has outstripped many conventional methods for DNA analysis.

### 4.1. Genomic Processing Applications

We now discuss applications of such a Biomolecular Database system to provide key genomic processing capabilities. Three basic kinds of applications are discussed, which demonstrate different ways in which the massive parallelism of Biomolecular Databases can be used: (1) *Rapid identification of individuals either susceptible to or resistant to chemical or biological agents*. We describe the selection of a group of DNA molecules based on a common property, and then use the information tags to identify the individuals selected. (2) *Large-scale gene expression profiling using Biomolecular Databases*. Expressed genes from multiple tissues are represented in a Biomolecular Database, from which they can be selected individually or in groups for subsequent expression analysis. (3) *High-throughput screening of candidate genes to optimize genetic association analysis for complex diseases such as heart disease or Parkinson's disease*. Pools of individuals are selected through use of the information tags appended to each DNA molecule in the database. The pools so selected are then subjected to genetic analysis. We describe in detail these three applications that concern genomic information processing, and constitute important genomic processing applications of Biomolecular Database systems for medical science.

### 4.1.1.  Rapid Identification of All Individuals Possessing a Specific Known Genotype

A single known genotype can confer properties making the individual either susceptible to or resistant to a particular chemical or biological agent found in the environment. It is certainly possible with existing biotechnology (e.g., hybridization experiments) to screen individuals for a given genotype. This is done one individual at a time, and is thus a relatively slow process. In addition, the cost of traditional genotyping of an individual ranges from $300 to over $1,000.

(At least one genotype databank has executed genotyping of approximately 1,000 individuals at considerable expense and time; however, for experimental purposes, that databank provides examples of previously executed individual genotyping at a cost of $0.50–1.00 per sample.) Clearly, an effort to screen a large number of individuals (say a million) would be slow and very expensive. In contrast, the methodology described herein is a *selection* for individuals with a certain genotype rather than a screen. It is correspondingly faster and less expensive. There is currently no available methodology for selection of specific genotypes. Many drugs that are quite effective in treating disease are very toxic to a small portion of the population. Currently, such drugs are removed from the market to avoid these rare but fatal adverse reactions. Such an approach is very costly from the standpoint of untreated disease. The removal of drugs from the marketplace because of rare fatal reactions is very costly in terms of untreated individuals, as well as the money spent on bringing those drugs to the market in the first place. Improved methods for identification of individuals at risk for adverse reactions would eliminate this cost. The capability of screening large numbers of individuals for a given genotype could also avoid a tremendous potential loss of life in the event of battlefield release of biological weapons or chemical agents.

As another example of a clinical application, one could construct a Biomolecular Database made from the blood samples of people suffering from Alzheimer's disease and their families, with the goal of finding genes that may increase people's risk of contracting the condition. The information tags could be used to select specific groups of molecules from this database. These molecules, which come from people with similar clinical symptoms, can then be used to test a large number of possible Alzheimer's disease genes. Genes that yield promising results could then be tested on the large number of individual samples from which the Biomolecular Database was made. The advantage of this approach is that it allows very efficient use of the limited DNA samples, and it is a good way to look at lots of different combinations of clinical features.

### 4.1.2. Large-Scale Gene Expression Profiling Using Biomolecular Databases

One may wish to determine the entire set of genes expressed by a particular cell type for a population of individuals who suffer debilitating effects due to a (perhaps unknown) chemical or biological agent. This may allow us to determine if there is a single or small number of genotypes that characterize susceptibility to that agent, over that population. Gene expression profiling is a labor-intensive and slow process. The conventional methods used are as follows: (a) *cDNA Hybridization Arrays*, which are 2D arrays of DNA spotted onto a solid support in an addressable way such that the spatial location of a spot identifies to the sequence of the DNA bound there. The input cDNA is labeled with a fluo-

rescent dye with a unique and readily imaged emission spectrum. After annealing on this array, the fluorescent cDNA provide a visual readout of the expression. (b) *SAGE libraries*, prepared by extraction from cDNA of very short tag sequences that characterize the expressed gene, followed by concatenation of a number of these tag sequences together for sequencing. Then computer software (using prior information on the relation between these tags and the original expressed cDNA) is used to determine which genes are being expressed. Gene expression profiling can require the development of new cDNA hybridization arrays, or the construction and sequencing of SAGE libraries. The methods for parallel analysis of large numbers of samples described here would streamline this process. In addition, the readout of SAGE data by microarray hybridization would result in significant savings of time and money as compared to the standard method of sequencing SAGE libraries. It would enhance our understanding of both complex disease processes and acute responses to biologic agents.

As another example of a clinical application, one could construct a Biomolecular Database made from a large group of healthy people, with the goal of finding people who are naturally resistant to certain germs, or who respond in certain ways to prescription drugs. One could study the selection of DNA strands from this mixture that have a specific sequence change in a specific gene that is known to change a person's resistance to germs or their response to drugs. Once these strands are isolated, the information tags would be examined to identify the people who have that change in their genes. This would be an extremely useful way to identify people who could have a bad reaction to a drug commonly used to treat disease. It could also be very useful in discovering people who are resistant to naturally occurring diseases or those caused by agents released during germ warfare.

### 4.1.3. Streamlining Identification of Susceptibility Genes

In terms of high-throughput screening of candidate genes to optimize genetic association analysis for complex diseases, consider the problem of genomic characterization of those individuals who first were infected by a biological agent, and then died. The death may often have been due to complications involving additional "complex" diseases, such as heart disease. Hence, mortality resulting from a chemical or biological agent attack may often have been due to complications involving a preexisting disease such as heart disease. Mortality can thus often only be predicted by considering both the individual's susceptibility to that agent, as well as that to various preexisting "complex" diseases. For many "complex" diseases, susceptibility often depends on a number of single-nucleotide polymorphisms (SNPs) in the human genome. Research into the genetic causes of complex disease is currently very expensive, and progress is slow, and complex diseases are quite common, affecting large propor-

tions of the population. Delays in understanding the genetic basis of these diseases slow the development of improved treatments at a significant financial and human cost. In the last 2 to 3 years, there have been several large-scale efforts to identify single-nucleotide polymorphisms in the human genome. The SNP consortium, a non-profit foundation composed of the Wellcome Trust and eleven pharmaceutical and technological companies, has agreed to deposit all SNPs they discover to public databases such as dbSNP, the SNP database maintained by the National Center for Biotechnology Information (NCBI). The number of entries in this database has increased from several thousand to over two million within the last 3 years. This sudden increase in the number of polymorphic markers has completely overwhelmed current methods for SNP genotyping and high-throughput screening. It has also become apparent that the incidence of single-nucleotide polymorphisms varies widely from one region of the genome to another, and large numbers of SNPs must be screened to analyze each candidate gene. Even with unlimited funds and the capacity for genotyping, serious challenges to the family-based association screening would remain, because the individual screening of a large number of SNPs would quickly exhaust the amount of DNA that can be easily obtained from a single individual. This problem is compounded by the sample cost of preparing pools of DNA from multiple individuals by simple mixing: once samples are mixed they cannot be separated again, and leftover, pooled DNA is wasted. Indexing of a Biomolecular Database can be of significant assistance in this regard. Large numbers of different groups of individuals can be selected from the Biomolecular Database by logical queries on the information tags. These pools can be used for allelic frequency determinations, and any remaining DNA can be added back to the remaining database.

As another example of a clinical application, one could use Biomolecular Databases to help discover what genes are turned on in a specific tissue of the body. Genes that are needed in the brain may not be expressed in the muscles, and genes needed in the muscles may not be needed in the liver. For this reason, measuring what genes are turned on in a specific tissue can help us understand what the possible functions of those genes might be. Biomolecular Databases would provide increased efficiency for these approaches.

## 4.2. Further Applications

The applications described above could be of critical value to the United States in the event of a terrorist release of a biological or chemical agent, as in the following brief scenario. A biological agent is released by a terrorist group in an American city or another populated area. The city is evacuated, but it becomes necessary to traverse a potentially contaminated area, or to revisit a known contaminated area. Clearly, any personnel sent into this area,

even with protective gear, are at risk for infection. A Biomolecular Database query would be initiated to identify personnel who possess a known genetic variation that prevents or mitigates infection. The personnel actually sent into the contaminated area could then be selected from the list of genetically resistant individuals.

As an alternative anti-terrorist application, suppose a large population (e.g., that of a city) has been exposed to a given biological or chemical agent. It then becomes apparent that a subgroup of individuals require significantly more aggressive medical therapy to survive, but for logistical reasons such aggressive therapy cannot be provided to ALL exposed individuals. Stored DNA from resistant and susceptible individuals can be used to determine the status of specific groups of genetic markers as described in application $C$ (markers are chosen based on biological and medical inferences). In this way, a series of markers diagnostic for increased susceptibility can be identified. This type of analysis is called class discovery, and has been applied to the treatment of breast cancer and leukemia, among other disorders. However, the use of Biomolecular Databases can greatly streamline this work. Once diagnostic markers have been identified, the techniques worked out in application 4.1.3 can identify individuals in need of more aggressive care.

## 5. DISCUSSION AND CONCLUSIONS

We have described Biomolecular Databases constructed from DNA for rapid genetic analysis of large populations of individuals and complex diseases involving multiple genetic loci. They may improve on conventional methods in size of database and speed of search with the Biomolecular Databases system.

### 5.1. Comparison with Biomolecular Computing Methods for SAT Problems

As described above, these selection operations can be executed by the use of recombinant DNA operations, using logical processing methods developed in the field of DNA computing. The methods used in DNA computing to solve combinatorial search problems such as the Boolean satisfiability (SAT) problem have the disadvantage that they require a volume that scales exponentially with the size of the problem (i.e., the number of Boolean variables). This is because the search space of possible Boolean variable assignments scales exponentially. In contrast, logical queries are executed only on the information tags of the existing database, so the volume only scales linearly with the number of strands of the Biomolecular Database.

## 5.2. The Key Advantages of Biomolecular Databases

The key advantages of Biomolecular Databases appear to be:

a. *Bypassing of conventional impasses*: In particular, the avoidance of sequencing for conversion from DNA (genomic DNA and transcribed RNA) to digital media.

b. *Ultra-compact storage media*: The extreme compactness and portability of the storage media—a pedabit of information can be stored (with tenfold redundancy) in less than a few milligrams of dehydrated DNA, or when hydrated may be stored in a few milliliters of solution. A Biomolecular Database is capable of containing the DNA of a million individuals (6 pedabits of information) in a volume the size of a conventional test tube.

c. *Massive molecular parallelism*: Although one query may require a number of minutes, it is operated on vast numbers of data items (DNA strands), implying a processing power of vast molecular parallelism with at least a few hundred teraflops. The operations can operate in parallel on an entire population of DNA.

d. *Scalability*: The technology requires volume that scales linearly with the size of the database, and a query time that remains nearly constant up to extremely large database sizes.

e. *Limitations*: The Biomolecular Database technology is limited to applications of a biological nature (where the data are DNA or easily convertible to DNA), and the operations are limited to logical queries in the Biomolecular Database, associative searches, and some essential database operations. It is not intended that the technology compete in any direct way with conventional high-performance computers. Instead, the objective is to bypass conventional bioinformatics methodology by processing biological material (genomic DNA and transcribed RNA) in "wet" media, rather than digital media.

## 5.3. Scalability of Biomolecular Databases Systems

The key parameters of Biomolecular Database are: (a) $N$ = the number of distinct elements of the Biomolecular Database, (b) $v$ = the number of variables (each ranging over 10 possible values) used in queries, and (c) $k$ = the number of individuals in application studies.

For our practical genomic applications of Biomolecular Databases to be fully realized in practice: (i) the database size $N$ should grow to extremely large values (with a long-term goal of approximately $10^{15}$), but (ii) for these applications the number of variables $v$ needs only to grow to moderately small constant values (with a long-term goal of approximately $v = 14$), since for the genomic applications considered only a limited number of values need to be recorded in the information tag per database element. The relative difficulty of obtaining

human genomic material limits the number of individuals $k$ in possible studies to approximately 1000, which is the size of the largest genomic database we are aware of for which one can legally obtain samples of genomic DNA. However, this figure of $k = 1000$ is by no means a limit on the capability of Biomolecular Database technology. In particular, these genomic databases are quickly growing in size, and it may be projected to grow by a number of multiples in a few years. Furthermore, military sources of human genomic DNA may be obtainable, providing alternate routes to obtain the samples of genomic DNA required in large-scale studies.

## 6. ACKNOWLEDGEMENTS

## 7. REFERENCES

1. Adleman L. 1994. Molecular computation of solution to combinatorial problems. *Science* **266**:1021–27.
2. Alizadeh AA, Eisen MB, Davis RE, Ma C, Lossos IS, Rosenwald A, Boldrick JC, Sabet H, Tran T, Yu X, Powell JI, Yang L, Marti GE, Moore T, Hudson J Jr, Lu L, Lewis DB, Tibshirani R, Sherlock G, Chan WC, Greiner TC, Weisenburger DD, Armitage JO, Warnke R, Levy R, Wilson W, Grever MR, Byrd JC, Botstein D, Brown PO, Staudt LM. 2000. Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling. *Nature* **403**(6769):503–511.
3. Alon U, Barkai N, Notterman DA, Gish K, Ybarra S, Mack D, Levine AJ. 1999. Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proc Natl Acad Sci USA* **96**(12):6745–6750.
4. Arnheim, N, Li HH, Cui XF. 1990. PCR analysis of DNA sequences in single cells: single sperm gene mapping and genetic disease diagnosis. *Genomics* **8**:415–419.
5. Bach E, Condon A, Glaser E, Tanguay C. 1996. Improved models and algorithms for DNA computation. In *Proc. 11th annual IEEE conference on computational complexity*, *J Comp Syst Sci* ?
6. Bancroft C, Bowler T, Bloom B, Clelland CT. 2001. Long-term storage of information in DNA. *Science* **293**(5536):1763–1765.
7. Baum EB. 1995. How to build an associative memory vastly larger than the brain. *Science* **268**:583–585.
8. Baum EB. 1996. DNA sequences useful for computation. In *DNA sequences useful for computation*, *Proc. 2nd DIMACS workshop on DNA-based computing*, Princeton. AMS DIMACS Series, **44**:235–241. Ed. LF Landweber, E Baum. See (http://www.neci.nj.nec.com/homepages/eric/seq.ps.)

9.    Box, GEP. 1978. *Statistics for experimenters: an introduction to design, data analysis, and model building*. Wiley, New York.

10.   Box, GEP. 1987. *Empirical model-building and response surfaces*. Wiley, New York.

11.   Braich RS, Chelyapov N, Johnson C, Rothemund PWK, Adleman L. 2002. Solution of a 20-variable 3-SAT problem on a DNA computer. *Science* **296**(5567):499–502.

12.   Cantor CR, Smith CL, Mathew MK. 1988. Pulsed-field gel electrophoresis of very large DNA molecules. *Annu Rev Biophys Biophys Chem* **17**:287–304.

13.   Chen CJ, Deaton R, Wang Y.. 2003. A DNA-based memory with in vitro learning and associative recall, *Proc. 9th annual meeting on DNA-based computers*, pp. 127–136.

14.   Clayton SJ, Scott FM, Walker J, Callaghan K, Haque K, Liloglou T, Xinarianos G, Shawcross S, Ceuppens P, Field JK, Fox JC. 2000. K-ras point mutation detection in lung cancer: comparison of two approaches to somatic mutation detection using ARMS allele-specific amplification. *Clin Chem* **46**:1929–1938.

15.   Corder EH, Saunders AM, Strittmatter WJ, Schmechel DE, Gaskell PC, Small GW, Roses AD, Haines JL, Pericak-Vance MA. 1993. Gene dose of apolipoprotein E type 4 allele and the risk of Alzheimer's disease in late onset families. *Science* **261**(5123):921–923.

16.   Corder EH, Saunders AM, Risch NJ, Strittmatter WJ, Schmechel DE, Gaskell Jr PC, Rimmler JB, Locke PA, Conneally PM, Schmader KE, Small GW, Roses AD, Haines JL, Pericak-Vance MA. 1994. Protective effect of apolipoprotein e type 2 allele for late onset Alzheimer disease. *Nature Genet* **7**:180–184.

17.   Cukras AR, Faulhammer D, Lipton, RJ, Landweber LF. 2000. Molecular computation: RNA solutions to chess problems, *Proc Natl Acad Sci USA* **97**:1385–1389.

18.   Deaton R. Murphy RE, Rose JA, Garzon M, Franceschetti DR, Stevens Jr SE. 1997. A DNA-based implementation of an evolutionary search for good encodings for DNA computation. In *Proc. IEEE Conference on Evolutionary Computation*, ICEC-97, pp. 267–271.

19.   Deaton R, Garzon M, Rose JA, Franceschetti DR, Murphy RC, Stevens Jr SE. 1998. Reliability and efficiency of a DNA-based computation. *Phys Rev Lett* **80**:417–420.

20.   Deaton R, Murphy RC, Garzon M, Franceschetti DR, Stevens Jr SE. 1999. Good encodings for DNA-based solutions to combinatorial problems. In *Proc. DNA-based computers*, II: DIMACS Workshop 10–12 June. Ed LF Landweber and EB Baum. DIMACS Series in Discrete Mathematics and Theoretical Computer Science **44**:247–258.

21.   Deming SN. 1987. *Experimental design: a chemometric approach*. Elsevier, New York.

22.   DeRisi J, van den Hazel B, Marc P, Balzi E, Brown P, Jacq C, Goffeau A. 2000. Genome microarray analysis of transcriptional activation in multidrug resistance yeast mutants. *FEBS Lett* **470**(2):156–160.

23.   Faulhammer D, Cukras AR, Lipton RJ, Landweber. 2000. Molecular computation: RNA solutions to chess problems. *Proc Natl Acad Sci USA* 97:1385–1389.

24.   Frutos AG, Thiel AJ, Condon AE, Smith LM, Corn RM. 1997. DNA computing at surfaces: 4 base mismatch word design. In *Proc. 3rd DIMACS meeting on DNA-based computers*, University of Pennsylvania, Philadelphia, June.

25.   Garzon M, Deaton R, Neathery P, Murphy RC, Franceschetti DR, Stevens Jr SE. 1997. On the Encoding Problem for DNA Computing. In *Proc. 3rd DIMACS meeting on DNA-based computers*, University of Pennsylvania, Philadelphia, June.

26.   Garzon M, Neel A, Bobba K. 2004. Efficiency and reliability of semantic retrieval in DNA-based memories. In *DNA computing, 9th international workshop on DNA-based computers*. Ed. J Chen, JH Reif. *Lect Notes Comput Sci* **2943**:157–169.

27.   Gehani A, and Reif JH. 1999. Microflow bio-molecular computation. In *Proc. 4th DIMACS workshop on DNA-based computers*, University of Pennsylvania, June 1998. Series in Discrete Mathematics and Theoretical Computer Science. Ed. H Rubin. American Mathematical Society, Providence, RI. Also appeared in special issue of *Biosystems: J Biol Inform Processing Sci* 52:(1–3):197–216.

28.   Gehani A, LaBean TH, Reif JH. 2000. DNA-based cryptography. In *5th DIMACS workshop on DNA-based computers*, MIT, June 1999. Series in Discrete Mathematics and Theoretical Computer Science. Ed. E Winfree. American Mathematical Society, Providence, RI.

29.   Gray JM, Frutos TG, Berman AM, Condon AE, Lagally MG, Smith LM, Corn RM. 1996. *Reducing errors in DNA computing by appropriate word design*. Draft paper, University of Wisconsin, Department of Chemistry, October 9.

30.   Hartemink A, Gifford D, Khodor J. 1998. Automated constraint-based nucleotide sequence selection for DNA computation, In *Proc. 4th DIMACS workshop on DNA-based computers*, University of Pennsylvania, June 1998.

31.   Helene C, Thuong NT. 1991. Design of bifunctional oligonucleotide intercalator conjugates as inhibitors of gene expression. *Nucleic Acids Symp Ser* **24**:133–137.

32.   Jonoska N, Karl SA. 1997. Ligation experiments in computing with DNA. In *Proc. IEEE Conference on Evolutionary Computation*, ICEC-97, pp. 261–265.

33.   Kaplan P, Cecchi G, Libchaber A. 1996. DNA-based molecular computation: template–template interactions in PCR. In *Proc. 2nd DIMACS workshop on DNA-based computing*. Ed. LF Landweber and EB Baum. DIMACS Series in Discrete Mathematics and Theoretical Computer Science, **44**:94–102.

34.   Kashiwamura S, Yamamoto M, Kameda A, Shiba T, Ohuchi A. 2003. Hierarchical DNA memory based on nested PCR. In *Proc. 8th DIMACS workshop on DNA-based computing*, Sapporo, Japan, June 10–13. Ed. M Hagiya, A Ohuchi. *Lect Notes Comput Sci* **2568**:112–123.

35.   Li HH, Cui XF, Arnheim N. 1990. Analysis of DNA sequences in individual gametes: application to human genetic mapping. *Prog Clin Biol Res* **340C**:207–211.

36.   Lipton RJ. 1996. DNA computations can have global memory. In *Proc. 2nd DIMACS workshop on DNA-based computing*. Ed. LF Landweber and EB Baum. DIMACS Series in Discrete Mathematics and Theoretical Computer Science, **44**:259–266.

37.   Liu Q, Liman W. Frutos AG, Condon AE, Corn RM, Smith LM. 2000. DNA Computing on surfaces. *Nature* **403**:175–179.

38.   Lizardi P, Huang X, Zhu Z, Bray-Ward P, Thomas DC, Ward DC. 1998. Mutant detection and single molecule counting using isothermal rolling circle replication. *Nature Genet* **19**:225–232.

39.   Mir KU. 1996. A restricted genetic alphabet for DNA computing. In *Proc. 2nd DIMACS workshop on DNA-based computing*. Ed. LF Landweber and EB Baum. DIMACS Series in Discrete Mathematics and Theoretical Computer Science, **44**:???page nos???.

40.   Niculescu AB, Segal DS, Kuczenski R, Barrett T, Hauger RL, Kelsoe JR. 2000. Identifying a series of candidate genes for mania and psychosis: a convergent functional genomics approach. *Physiol Genomics* **4**(1):83–91.

41.   Olson MV. 1989. Separation of large DNA molecules by pulsed-field gel electrophoresis: a review of the basic phenomenology. *J Chromatogr* **470**:377–383.

42.   Perou CM, Sorlie T, Eisen MB, van de Rijn M, Jeffrey SS, Rees CA, Pollack JR, Ross DT, Johnsen H, Akslen LA, Fluge O, Pergamenschikov A, Williams C, Zhu SX, Lonning PE, Borresen-Dale AL, Brown PO, Botstein D. 2000. Molecular portraits of human breast tumours. *Nature* **406**(6797):747–752.

43.   Pieles U, Englisch U. 1989. Psoralen covalently linked to oligodeoxyribonucleotides: synthesis, sequence specific recognition of DNA and photo-cross-linking to pyrimidine residues of DNA. *Nucleic Acids Res* **17**:285–99.

44.   Pirrung MC. 1995. Combinatorial libraries: chemistry meets Darwin. *Chemtracts Org Chem* **8**:5.

45.   Pirrung MC. 1997. Spatially-addressable combinatorial libraries. *Chem Rev* **97**:473.

46.   Pirrung MC, Chau JH-L, Chen J. 1996. Indexed combinatorial libraries: non-oligomeric chemical diversity for the discovery of novel enzyme inhibitors. In *Combinatorial chemistry: a high-tech search for new drug candidates*, pp. 191–206. Ed. SR Wilson, R Murphy. John Wiley & Sons, New York.

47. Pirrung MC, Connors RV, Montague-Smith MP, Odenbaugh AL, Walcott NG, Tollett JJ. 2000. The arrayed primer-extension method for DNA microchip analysis: molecular computation of satisfaction problems. *J Am Chem Soc* **122**:1873.

48. Pirrung MC, Zhao X, Harris SV. 2001. A universal, photocleavable, DNA base: nitropiperonyl 2'-deoxyriboside (dP*). *J Org Chem* **66**:2067.

49. Quillent C, Oberlin E, Braun J, Rousset D, Gonzalez-Canali G, Metais P, Montagnier L, Virelizier JL, Arenzana-Seisdedos F, Beretta A. 1998. HIV-1-resistance phenotype conferred by combination of two separate inherited mutations of CCR5 gene. *Lancet* **351**(9095):14–18.

50. Reif, J.H. 1998. Paradigms for biomolecular computation. Paper presented at 1st international conference on unconventional models of computation, Auckland, New Zealand, January. In *Unconventional models of computation*, pp. 72–93. Ed. CS Calude, J Casti, MJ Dinneen. Springer, New York.

51. Reif JH. 1999. Parallel Molecular Computation: Models and Simulations. In *Proc. 7th annual ACM symposium on parallel algorithms and architectures* (SPAA'95), Santa Barbara, CA, July 1995, pp. 213–223. Published in *Algorithmica*, special issue on *Comput Biol* **25**(2):142–176.

52. Reif JH. 2002. The emergence of the discipline of biomolecular computation in the US. Invited paper presented in a special issue on *Biomolecular Computing, New Generation Computing*, ed. M Hagiya, M Yamamura, T Head, **20**(3):217–236.

53. Reif, JH. 2002. Perspectives: successes and challenges. *Science* **296**:478–479.

54. Reif JH. LaBean TH. 2001. Computationally inspired biotechnologies: improved dna synthesis and associative search using error-correcting codes and vector-quantization, In *Proc. 6th DIMACS workshop on DNA-based computers*, Leiden, The Netherlands, June 13–17, 2000. *Lect Notes Comput Sci* **2054**:145–172.

55. Reif JH, LaBean TH, Pirrung M, Rana VS, Guo B, Kingsford C, Wickham GS. 2002. Experimental construction of very large-scale DNA databases with associative search capability. In *Proc. 7th DIMACS workshop on DNA-based computers*, Tampa, FL, June 10–13, 2001. *Lect Notes Comput Sci* **2340**:231–247.

56. Risch N, Merikangas K. 1996. The future of genetic studies of complex human disorders. *Science* **273**(5281):1516–1517.

57. Robinson BH, Seeman NC. 1987. The design of a biochip: a self-assembling molecular-scale memory device. *Prot Eng* 1:295–300.

58. Roweis S, Winfree E, Burgoyne R, Chelyapov NV, Goodman MF, Rothemund PWK, Adleman LM. 1998. A sticker-based model for DNA computation, *J Comput Biol* 5:615–629.

59. Sakakibara Y, Suyama A. 2000. Intelligent DNA chips: logical operation of gene expression profiles on DNA computers. *Genome Informatics* **11**:33–42.

60. Suyama A, Nishida N, Kurata K, Omagari K. 2000. Gene expression analysis by DNA computing. *Curr Comput Mol Biol* **30**:12–13.

61. Szatmari I, Aradi J. 2001. Telomeric repeat amplification, without shortening or lengthening of the telomerase products: a method to analyze the processivity of telomerase enzyme. *Nucleic Acids Res* **29**:E3.

62. Taylor GR, Logan WP. 1995. The polymerase chain reaction: new variations on an old theme. *Curr Opin Biotechnol* **6**:24–29.

63. Taylor GR, Robinson P. 1998. The polymerase chain reaction: from functional genomics to high-school practical classes. *Curr Opin Biotechnol* **9**:35–42.

64. Wellinger RE, Lucchini R, Dammann R, Sogo JM. 1999. In vivo mapping of nucleosomes using psoralen-DNA crosslinking and primer extension. *Methods Mol Biol* **119**:161–173.

65. Winfree E. 1998. Whiplash PCR for O(1) computing. In *Proc. 4th DIMACS workshop on DNA-based computers*, University of Pennsylvania, June 1998.

66. Wood DH. 1998. Applying error-correcting codes to DNA computing. In *Proc. 4th DIMACS workshop on DNA-based computers*, University of Pennsylvania, June 1998,

67. Zhang L, Cui X, Schmitt K, Hubert R, Navidi W, Arnheim N. 1992. Whole genome amplification from a single cell: implications for genetic analysis. *Proc Natl Acad Sci USA* **89**:5847–5851.
68. Zhao R, Gish K, Murphy M, Yin Y, Notterman D, Hoffman WH, Tom E, Mack DH, Levine AJ. 2000. Analysis of p53-regulated gene expression patterns using oligonucleotide arrays. *Genes Dev* **14**(8):981–993.