

DOE
Human Genome
Contractor/Grantee
Workshop

November 3-4, 1989

Santa Fe Institute
Hilton Hotel

Santa Fe, NM

Speaker Abstracts

**CONSTRUCTION OF PARTIAL DIGEST LIBRARIES FROM FLOW
SORTED CHROMOSOMES.**

L.L. Deaven, C.E. Hildebrand, J.L. Longmire, and R.K. Moyzis, Los Alamos
National Laboratory, Los Alamos, NM 87545

The goal of the National Laboratory Gene Library Project is the production of chromosome-specific human DNA libraries and their distribution to the scientific community for studies of the molecular biology of genes and chromosomes, for the study and diagnosis of genetic disease, and for the physical mapping of chromosomes. This is a cooperative project between the Los Alamos and Lawrence Livermore National Laboratories. At Los Alamos, a set of complete digest libraries have been cloned into the EcoRI insertion site of Charon 21A. These libraries are available from the American Type Culture Collection, Rockville, MD. We are currently constructing sets of partial digest libraries in the cosmid vector, sCos1, and in the phage vector, Charon 40, for human chromosomes 4,5,6,8,10,13,14,15,16,17,20, and X. Individual human chromosomes are sorted from rodent-human cell lines until approximately 1 μ g of DNA has been accumulated. The sorted chromosomes are examined for purity by in situ hybridization, DNA is extracted, partially digested with Sau3A1, dephosphorylated, and cloned into sCos1 or Charon 40. Partial digest libraries have been constructed for chromosomes 4,5,8,16, and X. Purity estimates from sorted chromosomes, flow karyotype analysis and plaque or colony hybridization indicate that these libraries are 90-95% pure. A 10X portion of the chromosome 16 cosmid library has been arrayed in microtiter plates. Extensive analysis of this library indicates that it contains many sequences known to be localized on chromosome 16 and that it is very useful for physical map construction. Further library constructions and arrays of those libraries are in progress. Supported by the US Department of Energy.

CONSTRUCTION OF PARTIAL DIGEST HUMAN DNA LIBRARIES FROM FLOW-SORTED CHROMOSOMES.

Marvin A. Van Dilla, Anthony V. Carrano, Mari L. Christensen, Pieter J. de Jong, Joe Gray, Jennifer McNinch, Barbara Trask, Ger van den Engh and Kathy Yokobata, Biomedical Sciences Division, Lawrence Livermore National Laboratory, P.O. Box 5507, Livermore, CA 94550

The goal of the National Laboratory Gene Library Project at the Los Alamos and Lawrence Livermore National Laboratories is the production of chromosome-specific human gene libraries and their distribution to the scientific community for studies of the molecular biology of genes and chromosomes, for the study and diagnosis of genetic disease, and for the physical mapping of chromosomes. The specific aim of Phase 1 of the project was the production of complete digest libraries in the lambda vector Charon 21A for each of the human chromosomal types purified by flow sorting. This has been accomplished, and the libraries are deposited in a repository at American Type Culture Collection, Rockville, MD. Phase 2, the construction of partial digest libraries with large inserts in lambda replacement vectors (accept about 9-23 kb) and in cosmid vectors (accept about 33-46 kb), is underway. Livermore is cloning 12 chromosomal types (1,2,3,7,9,11,12,18,19,21,22, and Y) and Los Alamos the other 12. Thus, each chromosomal type will be cloned into both lambda and cosmid vectors. Livermore is using the lambda vectors Charon 40 (accepts 10-25 kb inserts), and GEM11, with about the same acceptance range but particularly suited to efficiently clone, map, sequence, and "walk" along a chromosome. At Livermore, the cosmid vector is Lawrist 5 (accepts 34-46 kb inserts), which has the same advantageous features for users as GEM11 plus double the insert size. We have constructed 2 large Charon 40 libraries for chromosome 19, large GEM11 libraries for chromosomes 11, 21, 22, and Y, and Lawrist 5 libraries for all these chromosomal types except 11. The cosmid libraries are amplified as individual clones picked into 96-well microtiter dishes. These libraries are being characterized before general release; for the #19 library in Charon 40, this process has been completed and the library is available from ATCC. Characterization results and plans for the future will be given.

Work performed under the auspices of the U.S. Department of Energy by the Lawrence Livermore National Laboratory under contract number W-7405-ENG-48.

Structural and functional map of human chromosome 11.

G.A. Evans¹, G. Hermanson¹, K. Lewis¹, L. Selleri¹, P. Lichter², D. Ward², C. Richard³, D. Cox³, B.E. Rothenberg¹, M. Saleh¹, S. Maurer¹, J. Eubanks¹, D. McElligott¹, D. Kaufman¹, K. Pomeroy¹, C. Landel¹, J. Zhao¹, S. Chen¹, G. Andreason¹, K. Pischel¹, C. Toraya¹, R. Hart¹, and M. Roman¹. Molecular Genetics Laboratory¹, The Salk Institute for Biological Studies, La Jolla, California; Department of Human Genetics², Yale University School of Medicine, New Haven, CT; and ³Department of Psychiatry, University of California, San Francisco, California.

We are constructing a physical map and ordered clone set for human chromosome 11 using a combination of traditional and novel technologies and including the analysis of sites of chromosomal pathology, the derivation of uniformly spaced DNA markers for clinical analysis, and the localization of sites of putative genes. The approach includes chromosome analysis by 1) the production of region-specific arrayed cosmid libraries, 2) determining all of the overlapping cosmids in the clone collection by a novel technique of cosmid multiplex analysis, 3) the construction of "contigs" by multiplex analysis and chromosome walking 4) ordering cosmid "contigs" by high resolution *in situ* hybridization to metaphase and interphase chromosomes using confocal laser scanning microscopy, 5) locating Not-1 "linking clones" and clones containing HTF islands by robotic processing of the arrayed cosmid clones, 6) linking "contigs" and "linking clones" in a long range physical map by pulsed field gel analysis, 7) locating clones which contain sequences coding for structural motifs found frequently expressed in families of mammalian proteins, 8) determining orientation and linkage of genes and anonymous probes by analysis of radiation-fragmentation hybrids of chromosome 11, 9) the isolation of yeast artificial chromosomes encoding portions of chromosome 11, and 10) the generation of mapped STS (sequence tagged sites) for the eventual integration of this map with the physical and linkage maps generated by other groups. The set of ordered landmarks which have been mapped to chromosome 11 thus far have proven useful in the definition of sites of chromosome pathology based on localization between DNA landmarks rather than traditional cytogenetics. This work forms the basis of a database and template set for the eventual complete structural characterization of chromosome 11.

X-LINKED DISEASES - A GENOME STRATEGY C. Thomas Caskey,

David L. Nelson, Andrea Ballabio, Thomas D. Webster, Laura Corbo, Richard A. Gibbs,
Jeffrey S. Chamberlain, Albert Edwards, Marcus Grompe and David H. Ledbetter.
Howard Hughes Medical Institute, Institute for Molecular Genetics, Baylor College of
Medicine, One Baylor Plaza, Houston, TX 77030

The human X chromosome is composed of 150 Mbp of DNA and is linked to nearly 100 defined genetic disorders. Characterization of the mutations responsible for these diseases requires efficient means for both the isolation of the genes involved and the development of mutation detection. We have utilized yeast artificial chromosome (YAC) vectors to isolate large regions of the X in order to identify candidate genes. Libraries have been constructed from two somatic cell hybrids, one retaining the entire X, the other the Xq24-qter region. Two hundred human YAC clones have been isolated from each library and contain inserts averaging 150 kb. Assignment of clones to specific regions of the X chromosome has been performed using a new technique allowing generation of probes from the YACs by PCR with *Alu* primers. 65 clones have been mapped to specific deletion intervals and YAC contigs have been established by cross-hybridization of the *Alu* PCR probes. The deletion mapping panel developed for this purpose allows the identification of clones close to disease genes involved in the contiguous gene syndromes in Xp22.3 and Xp21 and Fragile X in Xq27.3. A method for isolation of human candidate transcribed regions directly from somatic cell hybrids utilizing human-specific *Alu* primers and the mRNA has also been developed. A library of 100 such clones has been constructed from Xq24-qter, and contains sequences from the human HPRT gene, confirming the ability of the technique to identify human transcripts. To demonstrate the utility of large scale sequencing, we have determined the sequence of the human HPRT locus (60 kb) using automated methods. From the sequence we have identified a new class of minisatellite polymorphism and developed a PCR-based assay capable of distinguishing seven alleles. This sequence information has provided the ability to detect point mutations leading to Lesch-Nyhan syndrome by PCR and direct sequencing of exons. Precise definition of point mutations in the human OTC gene has been achieved using mismatch chemical cleavage. For more complex gene rearrangements, multiplex PCR primer sets have been developed for HPRT, DMD and STS. In the case of DMD, scanning nine regions of this 2 Mbp gene with this method provides carrier detection in 80% of cases with gene rearrangements. These new approaches promise to facilitate our isolation and characterization of genes responsible for numerous human genetic disorders, and to apply this knowledge to detection, prevention and therapy.

Physical Mapping Methods and Progress on Chromosome 19

Anthony V. Carrano, Elbert W. Branscomb, Pieter J. de Jong, Emilio Garcia, Harvey W. Mohrenweiser, Anne Olsen, Thomas Slezak, and Marvin A. van Dilla. Biomedical Sciences Division, Lawrence Livermore National Laboratory, Livermore, CA 94550

The initial goals of our effort are to create physical maps of human chromosomes, to correlate them with the genetic map, and to sequence selected regions of the chromosomes. The physical maps spanning the chromosomes will consist of overlapping cloned DNA fragments (contigs), contained in phage, cosmid and yeast vectors. In the past two years, we have made progress in several areas: vector and library construction, DNA fingerprinting chemistry, algorithm development, automation of the map construction tasks, and contig construction.

Vector and Library Construction. We constructed vectors to: 1) facilitate cloning small amounts of DNA in cosmids; 2) clone NotI linking probes in lambda and plasmids; and 3) clone large fragments of DNA as yeast artificial chromosomes (YACs). The Lawrist series of cosmid vectors were used to construct chromosome 19-specific libraries from flow-sorted chromosomes from a monochromosomal #19 hybrid. About 20,000 cosmids (~9-fold redundancy) in two different bacterial hosts have been arrayed in microtiter trays to form a reference library. We are currently expanding and characterizing a library of over 23,000 presumptive YAC clones derived from the same chromosome 19 hybrid cell line. New plasmid and lambda vectors were used to create a NotI linking library of chromosome 19 and about 30 clones have been isolated.

Fingerprinting Chemistry. In order to construct a set of cosmid contigs for chromosome 19, we developed an automated fluorescence-based strategy for fingerprinting each clone. For this procedure, fluorophors are attached to the ends of restriction fragments from each cosmid clone using a robotic system. Fragment lengths are determined using a commercially available laser scanning device to acquire electrophoretic mobility data in real time. Up to four different fluorophors (i.e. four DNA types) can be run in each gel lane. In the present configuration, this permits us to fingerprint up to 48 cosmids per gel run. We expect to double this throughput.

Algorithm Development. We developed software to process the acquired fluorophor signals and convert the signal data to restriction fragment lengths for each cosmid. The fragment length data is used to compute a statistical measure of overlap between cosmids. The overlap statistic is the basis for assembling the cosmids into contigs and determining the "minimal" spanning set of cosmids for the contig. The assembled contigs are then presented graphically to the user for interaction and database query.

Map Construction. Over 2500 cosmids have been processed to date, 163 from a YAC clone derived from chromosome 14 and the remainder from chromosome 19. For chromosome 14, the established contigs span at least 553 kbp of the 600 kbp YAC insert length. For chromosome 19, more than 320 contigs were formed from randomly selected cosmids. Several of the chromosome 19 contigs represent known gene loci. Contigs are validated by restriction enzyme site mapping and/or by *in situ* hybridization to metaphase chromosomes. Contig sizes range from 2 to 7 members with a few very large pseudo-contigs, including one whose members share a putative alphoid DNA repeat. Sequence tagged sites are being developed for selected contigs. Both YACs and pulse-field maps will be used to close the inevitable gaps in the contig map.

(This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract No. W-7405-ENG-48.)

PHYSICAL MAPPING OF HUMAN CHROMOSOME 16. C.E. Hildebrand, R.L. Stallings, D.C. Torney, J.H. Jett, N.A. Doggett, J.L. Longmire, L.L. Deaven, and R.K. Moyzis. Los Alamos National Laboratory, Los Alamos, NM 87545.

Construction of a complete physical map of an entire human chromosome will require multiple strategies to localize and access rapidly any region of DNA along the chromosome. Assembly of ordered, overlapping clone maps (or contigs) represents a key step in the global physical mapping process. We have developed a new approach for identifying overlapping cosmid clones by exploiting the high density of repetitive sequences in complex genomes. Primary recombinants (~26,000) from a flow-sorted, chromosome 16-specific cosmid library were arrayed in 96-well microtiter plates and subsets were replicated to membrane filters. Using these master filters, clones containing the highly conserved repetitive sequence (GT:AC)_n were identified by hybridization with ³²P labeled oligo (GT:AC). Individual clones were fingerprinted using a combination of restriction enzyme digestions followed by hybridization with selected classes of ³²P labeled repetitive sequences. Information on lengths of all restriction fragments (EcoRI, HindIII, and EcoRI/HindIII double digest) coupled with the occurrence of repeat sequences on restriction fragments was acquired by image capture and the data were entered into a relational database (Nelson, et al., these proceedings). Based on this fingerprinting data, pairs of overlapping cosmid clones were identified. Our results have validated this approach following the fingerprinting of 2,441 individual cosmid clones, most containing one or more GT:AC repeats. A statistical model (D.C. Torney, these proceedings) was used to identify pairs of overlapping clones, each pair overlapping with a probability >0.95. A subset of pairs and contigs has been confirmed via macrorestriction analysis of genomic DNA with multiple rare-cutters using clones from the ends of contigs as probes (N.A. Doggett, these proceedings). By "nucleating" at specific regions in the human genome, and exploiting the high density of interspersed repetitive sequences in the human DNA, this approach allows, 1) rapid progress in the early phases of contig mapping, 2) the detection of small regions of overlap between clones, and 3) the production of a contig map with useful polymorphic landmarks for rapid integration of the genetic and physical maps. The reference arrays of primary recombinants coupled with the relational database for tracking clone-specific fingerprints, and other mapping information (e.g., genetic linkage, defined sequence-tagged-sites, etc.) provide important resources to expedite the cost-effective acquisition of extended physical maps for human chromosomes, or any complex genome. (This work was accomplished through the dedicated expert technical assistances of D. A. Nelson, J. G. Tesmer, L. M. Clark, A. A. Ford, A. C. Munk, L. M. Meincke, N. C. Brown, E. P. Saunders, J. K. Graeber, and S. H. Cox.)

Progress towards construction of a long range physical map of human chromosome 21.

H. Fang, J. P. Abad, A. Saito, D. Wang, J. F. Cheng, Y. Wu, W. Michals, C. L. Smith and C. R. Cantor. Department of Genetics and Development, Microbiology and Psychiatry, College of Physicians and Surgeons, Columbia University, New York, NY10032; Department of Immunology and Virology, Saitama Cancer Center Research Institute, Ina-Machi, Kitaadachi-gun, 362 Saitama-ken, Japan.

A complete low resolution physical map of human chromosome 21 is being constructed. Chromosome 21 is the smallest human chromosome with an estimated size of 50 Mb DNA. Several strategies have been used to achieve this goal. Isolation of a human telomere YAC clone enabled the ends of the map to be defined. About 30 single copy DNA probes, with previously assigned genetic map locations along the length of the chromosome, are being employed as anchor points. These probes were used to identify corresponding large *Not*I and *Mlu*I DNA fragments by hybridization to pulsed field gel fractionated DNA restriction digests. The physical map is being constructed by assigning fragments using single copy probes with known regional locations and assigning neighboring bands using partial digests, chromosome specific *Not*I linking probes and polymorphism among different cell lines. These approaches have allowed us to identify about 40 Mb of DNA that derives from chromosome 21 and to link up *Not*I fragments in several regions along the q arm, including a continuous map from the q telomere which contains 7 *Not*I fragments and covers 8.5 Mb DNA. The accuracy of the map is 0.1 Mb for DNA fragments larger than 1 Mb and 0.02 Mb for those smaller than 1 Mb. Although the map is not yet complete, it reveals interesting features such as the uneven distribution of putative genes along the chromosome and a greater than expected gradient of enhanced recombination near the q telomere.

Special-Purpose VLSI-Based System for the Analysis of Genetic Sequences

T. Hunkapiller, M. Waterman, R. Jones, M. Eggert, E. Chow, J. Peterson, and L. Hood
Department of Biology, California Institute of Technology, Pasadena, CA 91125
(818) 356-6408

The current size of genetic sequence databases means that extensive similarity analyses based on robust mathematical models require large, expensive computer hardware not generally available to most investigators. We are currently involved in developing a hardware alternative that will give most laboratories access to these rigorous algorithms at the level of affordable workstations and/or PCs. We are designing a board-level biological information signal processor (BISP) coprocessor assembly. BISP is a systolic implementation of dynamic programming methods for the determination of local sequence similarities and alignments. Each BISP board will include an array of BISP processor chips responsible for performing the dynamic programming methods and a tightly coupled coprocessor (i860) that will provide complete alignments to the host CPU from the scoring information and locations provided by the systolic array. Difference tables and gapping penalties are completely user-definable. Also, there are no arbitrary limits on sequence lengths of either the query or database sequences, and searches can be made simultaneously for multiple query sequences (depending on the size of the sequences and the systolic array). The array will process at up to 20 megacharacters per second. At this maximum speed, a query sequence the size of the systolic array could be compared against all of the current GenBank (both orientations) in under 3 s. The size of the array is extensible and will range from about 500 processor cells to several thousand. The BISP chip is being implemented in 1 micron CMOS technology, based on custom standard cell methods. Each chip will have multiple identical processor cells, local control circuits, and a results cache. The nearly full-custom processor cell design is nearing completion while the control circuit logic and floor plan have been generally detailed. Our present schedule calls for tested design vectors to be completed by late fall and handed over to the fabrication service for prototype chip delivery by winter 1990.

MATHEMATICS, COMPUTATIONS, AND DATABASING IN THE LIVERMORE PHYSICAL MAPPING EFFORT

Elbert Branscomb, Tom Slezak, David O. Nelson, Mark Wagner, Robert Pae, and Anthony V. Carrano. Biomedical Sciences Division L-452, Lawrence Livermore National Laboratory, Livermore, CA 94550.

The effort underway at Livermore to construct physical maps of human chromosomes in the form of spanning ordered clone libraries presents computational problems in four general areas: (1) signal processing and feature extraction, (2) statistical analysis of overlap between cloned segments and automated contig assembly, (3) data visualization and computer assisted editing and analysis, and (4) data base support. A brief outline of these problems and of the approaches taken to address them will be presented (posters dealing with most of the issues involved are being presented separately). We will also summarize the results of a number of tests designed to assess the reliability and quality of the results produced by our mapping strategies. These tests are based primarily on the experience gained thus far in analyzing over 2500 cosmid-cloned fragments from chromosome 19, along with 163 cosmids from a 600 Kb YAC-cloned chromosome 14 fragment. The tests address: (a) the rate and causes of false positive overlap detection, (b) the role in overlap detection of experimental imprecision and error, (c) the efficiency of overlap detection, (d) the role of repetitive DNA elements, (e) the representativeness of the cosmid libraries, and (f) the deviations from statistical ideality of the fingerprinting methods being used to characterize the clones. (This work was performed under the auspices of the U.S. Department of Energy by the Lawrence Livermore National Laboratory under Contract No. W-7405-ENG-48.)

ABSTRACT

Computational Algorithms for the Construction of Physical Maps

David C. Torney, T-10
Los Alamos National Laboratory

A first stage in construction of a cosmid physical map is the identification of pairs of cosmids containing cloned DNA likely to overlap in the target genome. The following techniques were developed to predict the probability of overlap for the human DNA in a pair of cosmid clones, based on the LANL fingerprint data. A separate overlap statistic is computed for the comparison of fingerprint data of each digest. Since the LANL fingerprint currently employs three digests, three overlap statistics are computed for each pair of clones. To determine the probability of overlap, these results are used in conjunction with a statistical model in which the three overlap statistics are computed for simulated fingerprint data generated from clones known either to overlap or to not overlap. The overlap statistic is based on an analytic formula and it properly reflects the contributions from disparate fingerprint data types, namely, apparent fragment size and fragment hybridization status. The overlap statistic is large when there are many apparent pairs of nearly identical fragments in the digest of both clones. The overlap statistic is constructed from a comparison matrix, with the matrix elements containing the result of pairwise comparison of digest fragments in the first clone with digest fragments in the second clone. The overlap statistic results from the symmetric sum of products of the aforementioned matrix elements, taken in all possible ways 1, 2, 3, ..., n at a time -- under the restriction that no fragment be "paired" more than once in any product, that is only one element can be used from any row or column. These symmetric sums of products are evaluated approximately, with a complexity bounded by the cube of the maximum of the number of fragments in the two digests. For a randomly selected pair of clones' fingerprints, it takes approximately 1.0×10^4 cpu seconds to compute an overlap statistic on one processor of the IBM 3090 computer.

Many features of this approach to overlap detection can be generalized to other fingerprint data. Furthermore, we are currently extending this approach to contig construction. This analysis allows us to make predictions for the results of variations of the fingerprint data, such as changing the number of digests or number and type of hybridizations.

Acknowledgements:

This algorithm was developed on the IBM/Los Alamos joint study IBM 3090 computer, with the assistance of Doyce Nix and Steven White, of IBM. This work benefited from discussions with David J. Balding, University of London, and with Randall Dougherty, Ohio State.

LBL Human Genome Center Informatics Strategies

W. Johnston, M. Hutchinson, S. Lewis, V. Markowitz,
J. McCarthy, F. Olken, D. Robertson, S. Spengler, E. Theil, M. Zorn

Our current computing and information management strategy involves an image handling system, an integrated Chromosome 21 information system, an interactive mapping tool, data interchange standards, thesauri, Extended Entity Relationship (EER) database design and query tools, string matching and overlap detection algorithms, and a standard and extensible computing environment.

Much of our original data will consist of images from autoradiograms, confocal microscopy, and scanning tunneling electron microscopy. In order to deal with this data in a consistent and archivable way we are integrating image acquisition hardware, image processing, application specific image analysis, an image data management system, image compression, archiving, and a mass storage system.

We are assembling an integrated Chromosome 21 information system which will encompass all traditionally and electronically published data concerning Chromosome 21, e.g., physical, genetic, and cytogenetic maps, sequences, and bibliographic citations, as well as providing for incorporating local experimental results. A thesaurus is being constructed to facilitate translation among various nomenclatures. Sequence tagged sites will be used to integrate the various maps. This database will be accessible by other Chromosome 21 researchers.

Interactive map drawing and editing software, and map alignment and construction software (for our protocols) will be provided in a Unix workstation environment. The map editor will provide an interface to the C21 database that permits both graphical and parametric editing, and will update a map database.

For traditional DB design we have developed techniques based on Data Flow Diagrams and Extended Entity Relationship schemas. We use Data Flow Diagrams to model protocols, which we then translate to EER schemas. From here we have developed an EER schema to relational schema translator, which we are presently testing. This methodology is being used for the design of a Lab Information Management System modeled on LANL's lab notebook work.

We are working with the U.C. Berkeley and U.C. Davis Computer Science Departments to develop very efficient special purpose algorithms that search for patterns with certain characteristics (motifs). Examples include the problems of identifying inverted, separated or complemented repeat patterns in a string, where we do not put any a priori bound on the separation between the two strings. In the area of string matching for map alignment, the currently used Dynamic Programming techniques are likely to be impracticable for problems of the size that will be created by the massive

amounts of data generated by the Human Genome Initiative. These collaborations seek to develop algorithms which are more efficient than DP for many sequence matching problems. This work has already yielded efficient algorithms for certain approximate sequence matching problems.

In addition to the general informatics work we are constructing a modern, network distributed, long term, and expandable computing environment for LBL Human Genome Center. Part of this strategy includes using and developing software that is either portable and vendor system independent, or available via a standard server interface. To accomplish this, our software is comprised of appropriate academic and scientific prototype software, as well as commercial and locally developed software, all of which is Unix and X window system based, or in the case of the servers, TCP/IP or ISO protocol based. As a result of this, the environment extends over the large number of currently available Unix systems, from IBM PCs to Crays. For example, the imaging software is an extension of the New York University, Unix based HIPS image processing package. HIPS provides the foundation for the locally developed image enhancement and analysis software. The image database facility is provided by Image-Query (U.C. Berkeley), data archival and retrieval will be done by a locally developed archiving system, and an IEEE model mass storage system. For data management in this environment, we are investigating both relational (Sybase and Ingres) and object oriented (Iris and Objectivity) commercial data management systems. Due to their expense, these systems will likely be configured as remotely accessible servers. Our user interfaces are being built on NASA's user interface builder (TAE), Unidraw and Interviews (Stanford U.), and the HP window application library. MacIntoshes will be supported as X-window terminals. Initially, we may use map manipulation and display software like the Gene Construction Kit as a user interface on Macintoshes with the data being managed elsewhere. For sequence analysis we are currently using EuGene (Baylor), UCSF, and Lipman-Pearson (NIH) codes. Data interchange standards based on ASN.1 are being developed with NLM and NCBI. All of this software is, or soon will be, Unix and/or X-window system based.

Supporting a Curator:

Investigating Potential Applications of
Logic Programming and Parallel Computation

Ross Overbeek

Mathematics and Computer Science Division
Argonne National Laboratory

Researchers at Argonne National Laboratory have been investigating potential applications of logic programming and parallel computation to problems in sequence analysis. We are working closely with Carl Woese's group at the University of Illinois in Urbana, developing tools to support their research on the structure of the ribosome. We have written programs to create and maintain multiple sequence alignments and to predict secondary structure from such alignments. Using these tools, we have obtained alignments of as many as 400 sequences of 16S rRNA, each approximately 1600 nucleotides long; more importantly, the alignments closely resemble those produced manually. From these alignments, we have produced estimations of secondary structure that appear somewhat better than the computer-based estimations reported in the literature. We have also developed tools to search for covariance; these tools have afforded new insights into probable secondary/tertiary structure.

Our initial efforts have convinced us to help the group at Urbana produce a complete "biologists' workbench," initially supporting the functionality required for research on the ribosome. The software will be written in Strand (a dialect of logic programming) and C. The central benefits from using a bilingual approach based on Strand and C are as follows:

1. The use of a higher-level language like Strand allows rapid prototyping of algorithms. In our work on multiple sequence alignment and secondary structure prediction, we found that it was far easier to implement changes and specialized heuristics using logic programming.
2. The system remains portable over a wide class of hardware (from Next machines in the \$6500 range to large multiprocessors).
3. The capabilities of multiprocessors can be exploited with no source changes, and reasonable efficiency. For most of the computationally intensive operations currently requested, linear speedups are easily achievable.

Optimizing Procedures for Genomic Repositories

William C. Nierman and Donna R. Maglott, American Type Culture Collection, Rockville, MD 20852

An important activity of the American Type Culture Collection (ATCC) is to provide the scientific community with well-characterized, cloned DNA segments from a variety of taxa to support gene structure and function studies and genomic mapping efforts. Having established and operated for several years a repository of human and mouse clones and libraries, the ATCC has been evaluating methods to manage a repository of a complete genome. As a model system, a set of 852 minimally overlapping genomic clones belonging to 255 different contigs from the yeast Saccharomyces cerevisiae has been provided by Dr. Maynard Olson at Washington University (Olson et al., Proc Natl Acad Sci USA 83:7826-7830). Information describing each clone and its position was also provided in an electronic format. Procedures being evaluated in the laboratory include automation of sample preparation and handling, use of the DNA sequencer from Applied Biosystems, Inc. (Model 370A) for restriction fragment size analysis, and preservation of clones both for storage and for distribution to investigators. Computerized methods to manage information about repository materials are being developed for laboratory notebook data as well as for integrating genetic and physical maps.

Techniques for Determining the Physical Structure of Entire Human Chromosomes

Cassandra L. Smith
Human Genome Center, Lawrence Berkeley Laboratory

In order to construct a map, it is necessary to have a means of uniquely identifying each DNA fragment; our strategy is to use DNA sequences from the ends of the fragments. Knowing a small sequence (50-100) bp) from each end of a large DNA fragment will permit each fragment to be uniquely identified. Furthermore, matching these DNA sequences with similarly sized DNA sequences from linking clones (clones that contain the DNA from the ends of two adjacent large fragments) will facilitate map construction. The advantages of this protocol over existing ones is the speed of generating results, the precision of the ordering, the simplicity of data analysis, and the fact that the mapping process generates sequence data as well.

In addition to the traditional means of accomplishing the above tasks, we are also investigating ways of using amplification via the polymerase chain reaction (PCR) to obtain DNA sequences from the ends of large fragments and from linking clones. Ultimately, we plan to generate linking clones directly via PCR, thereby avoiding some of the pitfalls of traditional cloning methods that make it difficult to complete a map. The PCR-based sequencing strategies are also attractive because they can be readily automated and adapted to existing automated technologies, such as DNA sequencers.

THE IMPORTANCE OF NEW TECHNOLOGY FOR THE HUMAN GENOME PROGRAM.
Leroy Hood, Division of Biology, 147-75, California Institute of Technology, Pasadena,
CA 91125.

There are a variety of different visions of what the Human Genome Initiative entails. My own feeling is that for the first ten year period, there should be an enormous emphasis on the development of new technologies for mapping, sequencing and analyzing DNAs. Indeed, I think the mapping and sequencing projects that are carried out should be done on a regional basis and oriented toward developing new strategies, technologies and instrumentation for these objectives.

A variety of new approaches is being developed for DNA mapping procedures. Certainly, some of the most effective of these entail the use of fluorescent markers and multi-unknown analysis in a single channel. Likewise, a variety of new procedures has been proposed for DNA sequencing including those that envision sequencing a single DNA molecule after processive enzymatic degradation, tunneling electron microscopy and mass spectrometry. It is also possible to envision improvements in the current automated DNA sequencing technologies that will over a period of ten years will achieve a 100-fold increase in throughput, perhaps with increased accuracy and decreased cost. Finally, certainly some of the most challenging problems that face the Genome Initiative are those of acquiring, storing, analyzing and distributing the information that will come from this project. Our own belief is that a multiplicity of approaches will have to be directed at these problems including specialized co-processors, parallel processing hardware and software, the use of neural-net analysis as well as even more sophisticated techniques and of course the development of more effective algorithms for analyzing fundamental patterns in DNA and proteins. The question of how to handle the multiplicity of biological databases that will exist in ten years, making them transparent to scientists so that appropriate data can be acquired from each in the course of a particular project, is also a problem of fundamental importance.

These issues will be discussed in the context of where we stand regarding these technologies today.

ISOLATION OF REGION-SPECIFIC PROBES BY ALU-PCR AND COINCIDENCE CLONING.

Pieter J. de Jong, Chira Chen, Patricia Wilkie, Frans Lohman, Charalampos Aslanidis, Heinz-Ulrich Weier, Anne Fertitta and Anthony V. Carrano, Biomedical Sciences Division, Lawrence Livermore National Laboratory, Livermore, CA 94550

As part of our effort to prepare an ordered clone map of human chromosome 19, we are isolating dense sets of probes for specific chromosomal regions. Oligonucleotides hybridizing to the ALU repeat family are used as PCR primers. Such primers allow the *in vitro* enzymatic amplification of sequences flanked by two ALU repeats, provided that the distance between the repeats is in the range of a few hundred bp up to a few kbp. Amplification products have been characterized using pilot studies with a set of yeast artificial chromosomes containing human DNA fragments (120 – 600 kbp). The PCR primers (#33 and 34) correspond to positions (47–13) and (226–260) of the consensus ALU sequence, except for nucleotide degeneracy at sites in the consensus sequence known to be hyper-mutable. On average, one or two ALU-PCR products are generated per 100 kbp of human DNA (averaged over 7 YAC clones with a total of ~ 2Mbp), somewhat dependent on the particular choice of primers (only #33, only #34 or both). The PCR products appear to be useful as unique sequence probes. Yeast DNA and Chinese hamster DNA do not generate PCR products with primer 34 under conditions which generate a large set of PCR products (a smear in an ethidium bromide stained agarose gel) from human DNA or hybrid cell DNA (about 1% human in hamster background). Using shorter (non-degenerate) primers internal to primer 34, the total number of different PCR products is considerably reduced. As a result, hybrid cells (e.g. radiation hybrids) generate characteristic patterns of PCR products rather than a smear when analyzed by agarose gel electrophoresis. The combined PCR products derived from hybrid cell lines have been used successfully to preferentially stain the corresponding region of the genome by *in situ* hybridization. Regions used for the ALU-PCR probe isolation are limited by the availability of corresponding hybrid cell lines. To extend the usefulness of this technique, we are exploring a procedure to specifically clone (human)PCR products in common between two hybrid cells ("coincidence cloning"). This allows the boundaries for the regional probe isolation to be defined by combinations of hybrids rather than single hybrid cell lines.

Work performed by the Lawrence Livermore National Laboratory under the auspices of the U.S. Department of Energy under contract No.W-7405-ENG-48.

Thermal Stability Mapping of DNA by Random Fragmentation and Two-Dimensional Denaturing Gradient Electrophoresis

L. S. Lerman, Nashua Gabra, Eric Schmitt and Ezra Abrams
Department of Biology, Massachusetts Institute of Technology, Cambridge, MA 02139
(617) 253-6658

The thermal stability of the double-helix in a standard solvent is fully determined by the base sequence. Within a long DNA molecule, each local region (ranging in length from a few dozen bases up to several hundred base pairs) undergoes a transition from an ordered helix to a disordered, randomized configuration (melting) within a narrow temperature span, typically from 1 to 3 degrees for the change from 95% helical to 5% helical. It is convenient to characterize the transition in each region, or domain, by the T_m , the temperature at which there is a 50-50 equilibrium between the helical and melted forms. Within human genomic DNA there is substantial variation in the local T_m , as much as about 35 degrees, often with distinct and sharp boundaries between adjacent domains. While the pattern and characteristics of this sequence of domains in long DNA molecules is inferred principally by statistical-mechanical theory, their reality is more directly observable in short DNA molecules by means of absorption spectroscopy as a function of temperature, or by denaturing gradient electrophoresis. Since the T_m of each domain is changed only very slightly by the substitution, addition, or deletion of one or a very few bases, the sequence of domains provides a robust counterpart to the base sequence. It is less sensitive to trivial individual variation, also including methylation, than a map of restriction sites, and it reflects biological function more closely than the map of restriction sites.

In two-dimensional separation of randomly fragmented genomic DNA, each random fragment is identified by X, Y coordinates representing its length and the T_m of the domain with the lowest T_m in the fragment. All fragments in which that domain is the lowest will find a similar gradient level, regardless of their length. This distribution and the response to specific sequence probes provides a means, in principle, for determining the spacing, order, and T_m among those domains that have a lower T_m than the average and those with the highest T_m . Also, it provides measurements, in principle, of the nucleotide distances between each of these domains and any arbitrary set of sequence identifiers or probes.

Our current effort is concerned with refining and calibrating various aspects of the two dimensional denaturing gradient technique and related procedures. These include: 1) the preparation of a fully random distribution of fragments from lambda DNA and yeast artificial chromosomes containing long human genomic inserts using iron-peroxide nicking and S1 cleavage, 2) reducing the breadth of bands produced by very long DNA molecules in the denaturing gradient, 3) analysis of the band broadening observed when the domain with lowest T_m is surrounded by higher melting domains, and 4) development of optical, mathematical, and computing procedures for calibrating gel photographs or autoradiographs in terms of quantitative, point-to-point distributions of DNA.

A random and directed priming strategy for high-volume DNA sequencing using libraries of oligonucleotides

F. William Studier and John J. Dunn
Biology Department,
Brookhaven National Laboratory,
Upton, New York 11973

Direct sequencing of DNAs using libraries of oligonucleotide primers of length 8, 9, or 10 could greatly reduce the cost and effort of high-volume sequencing. Analysis of the statistics of priming indicates that libraries containing as few as 10,000 octamers, 14,200 nonamers, or 44,000 decamers would have the capacity to determine the sequence of almost any cosmid DNA with only about a 5% penalty of redundant sequencing dictated by lack of primers in the library (assuming 500 nt of sequence is determined from a single priming). A combination of random and directed priming with oligonucleotides of these lengths could determine the sequence of a cosmid DNA in 1.2-1.5 times the minimum number of sequencing reactions required, in contrast to random cloning techniques, which might require 5-10 times as many.

The sequence of each cosmid DNA could be determined from a single DNA preparation, eliminating the need for mapping or subcloning or preparation of multiple DNA samples. Primers would be instantly available, and since each preparation of oligonucleotide would be used repeatedly to prime sequencing reactions in different DNA molecules, the cost of primers would become well below \$0.001 per nucleotide of sequence information obtained.

The success of this strategy requires that a considerable fraction of octamers, nonamers, or decamers be able to prime selectively in double-stranded DNAs 45,000 base pairs long. Initial results indicate that this is likely to be the case.

Computer-Assisted Multiplex DNA Sequencing

G.M. Church, G. Gryan, S. Kieffer-Higgins, L. Mintz, M.J. Rubenfield, and M. Temple. Department of Genetics, Howard Hughes Medical Institute, Harvard Medical School, Boston, MA 02138-3800. (617) 732-7562.

Several laboratories are sequencing genomes (ranging from 1 to 15 Mbp) from each phylogenetic kingdom. The genome closest to completion is *E. coli* (20% of 4.7 Mbp). These sequences will define consensuses for classes of protein domains, evolutionary conservation, and change. While participating in this quest, we have developed a new multiplex DNA sequencing method [Church et al., *Science* 240, 185 (1988)]. In multiplex DNA sequencing, 480 sequencing reaction sets, each tagged with specific oligonucleotides, are run on a single gel in 12 pools of 40 and transferred to a membrane. We hybridize 75 such membranes simultaneously. The resulting sequence film images are digitized, and sequence interpretations are superimposed on the enhanced 2-D images for editing. The computer program (REPLICA) uses internal standards from multiplexing to establish lane alignment and lane-specific reaction rules by discriminant analysis. The automatic reading phase takes one hour per film (3kb) on a Vaxstation. Images with overlapping data can be viewed side by side to facilitate decision making. Hash-table-based routines for linking up shotgun sequences in the megabase range are compatible in speed with the rest of the software.

Multiplex DNA Sequencing. Robert Weiss and Raymond F. Gesteland, Dept. of Human Genetics, University of Utah Medical Center, Salt Lake City, UT 84132.

We have developed a rapid method for sequencing cosmid DNA that employs multiplex probing as originally described by Church and Gilbert. A large number (25-50) of DNA samples each cloned in an equal number of special vectors are prepared together and sequenced as a mixture using dideoxy sequencing primed from a common sequence. After conventional separation by gel electrophoresis, the mixed DNA pattern is transferred to a membrane. Each of the individual sequences is then revealed by repetitive rounds of probing, washing, reading, stripping and reprobing with labeled oligonucleotide each of which is unique for a sequence in one of the vectors. In a test project involving cosmids from the NF-1 locus on human chromosome 17, six membranes containing 1440 clones were processed in a large drum so that the probing and film exposure could be done without handling of the membranes. Consistent recoveries of 10-20,000 bases of sequence were obtained in each cycle with a cycle time of 1-2 days including the autoradiography from ^{32}P -labeled probe with very little labor. We, with help from Diane Dunn, are developing a method for efficiently creating in vivo random subclones that are ready for multiplex sequencing by using a set of Tn3-based transposons containing appropriate primer and identifier sequences. An optical lab has been set up by Dr. Jeff Ives and Dr. Achim Karger with the help of Dr. Joel Harris (Chemistry Dept.) to compare the feasibility of fluorescent and chemiluminescent tags for the multiple probes. Alternative membranes with low fluorescent backgrounds are being developed with Dr. Karin Caldwell (Bioengineering). Using radio frequency plasmid discharge in ammonia gas, a polypropylene membrane with higher DNA binding affinity and lower background fluorescence than common nylon membranes has been developed. The efficiency of multiplex sequencing has created a bottle neck at the step of reading autoradiograms and a research group including Ives, Mike Murdock and Drs. Tom Stockham (Electrical Engineering) and Neil Cotter (also Electrical Engineering) has been assembled that is trying to solve this problem and to deal with data that might come from CCD images of the fluorescent or chemiluminescent membranes. Harold Swerdlow is investigating the feasibility of using fluorescent detection of sequence ladders during gel electrophoresis in microbore capillaries (70 microns).

DNA SEQUENCE ANALYSIS WITH MODIFIED BACTERIOPHAGE T7 DNA POLYMERASE. Stanley Tabor, Hans E. Huber, John Rush, and Charles C. Richardson, Department of Biological Chemistry and Molecular Pharmacology, Harvard Medical School, Boston, MA 02115.

The 3' to 5' exonuclease activity of phage T7 DNA polymerase (gene 5 protein) can be inactivated selectively by reactive oxygen species. The chemically modified enzyme is highly processive in the presence of *E. coli* thioredoxin and discriminates against dideoxynucleoside triphosphates (ddNTPs) only four to six fold. Consequently, dideoxynucleotide-terminated fragments have highly uniform radioactive intensity throughout the range of a few to thousands of nucleotides in length. There is virtually no background due to terminations at pause sites or secondary-structure impediments in the template. Chemically modified gene 5 protein, by virtue of having low exonuclease activity, has enzymatic properties that distinguish it from native gene 5 protein. We have exploited these properties to show by a chemical screen that modification of a histidine residue reduces selectively the exonuclease activity. *In vitro* mutagenesis of histidine 123, and of the neighboring residues, results in varying reduction of the exonuclease activity. A deletion of 28 amino acids that encompasses His¹²³ eliminates all exonuclease activity ($< 10^{-6}$ %).

Incorporation of ddNTPs by T7 DNA polymerase and *E. coli* DNA polymerase I is more efficient when Mn^{2+} rather than Mg^{2+} is used for catalysis. Substituting Mn^{2+} for Mg^{2+} reduces the discrimination against ddNTPs approximately 100-fold for DNA polymerase I, and 4-fold for T7 DNA polymerase. With T7 DNA polymerase and Mn^{2+} , ddNMPs and dNMP are incorporated at virtually the same rate. Mn^{2+} also reduces the discrimination against other analogs with modifications in the furanose moiety, the base, and the phosphate linkage. The lack of discrimination against ddNTPs using the genetically modified T7 DNA polymerase and Mn^{2+} results in uniform terminations of DNA sequencing reactions, with the intensity of adjacent bands on polyacrylamide gels varying in most instances by less than 10%. A novel procedure that exploits the high uniformity of bands can be used for automated DNA sequencing. A single reaction with a single labeled primer is carried out using four different ratios of ddNTPs to dNTPs; after gel electrophoresis in a single lane, the sequence at each position is determined by the relative intensity of each band.

Abstract for Santa Fe Institute Workshop

New approaches for constructing expression maps of complex genomes.

Kandpal, R.P., Parimoo, S., Swaroop, A., Gruen, J., Arenstorf, H.P., Shukla, H., Ward, D.C. and Weissman, S.M.

An approach to mapping complex genomes is to obtain a set of cDNA clones corresponding to most, if not all, of the functional genes of the organism and to arrange these cDNAs in a linear order based on their position along the chromosomes. Correlation of such an ordered set with known genetic and physical markers should provide a valuable data base for localizing genetic disorders and for complementing other approaches of generating ordered and ultimately sequenced clones of genomic DNA. Towards this goal, we are developing a series of selection and cloning methods to be used in conjunction with the recent procedures of "in situ" hybridization and chromosome mapping of DNA probes. These methods are designed to first derive a set of rare-cutter linking fragments ordered along a chromosome. To expedite ordering these clones, we are exploring a "son-of-jumping" approach in which the principle of chromosome jumping is combined with sensitive polymerase chain reaction methods. These clones are then used to selectively enrich genomic DNA located in macro-restriction fragments adjacent to the chosen linking clone using the DNA "fishing" method. Fragments of cDNA corresponding to the "fished" genomic fragment are then isolated by a sensitive hybridization-selection procedure. We are also exploring methods to further enrich the specific cDNAs which exploit the second order kinetics of double-stranded DNA reannealing, and to use these fragments to directly obtain full-length cDNA clones from a "normalized" cDNA library.

October 13, 1989

Uses of Stable Isotopes for DNA Sequencing
K. Bruce Jacobson, Biology Division, Oak Ridge National Laboratory,
Oak Ridge, TN 37831-8077

The advantages of stable isotopes over radioisotopes for DNA sequencing are several: 1) increased rate of sequencing is possible since many isotopes can be used simultaneously as labels for DNA fragments that are separated electrophoretically, 2) the labeled oligomers are stable, and 3) the hazards of exposure and disposal are eliminated. The two problems associated with stable isotopes are 1) how to label the DNA and 2) how to detect the label. To label oligomers with iron or tin methods have been devised that convert Fe_2O_3 to ferrocene carboxylic acid and SnO_2 to $(\text{C}_5\text{H}_5)_2\text{SnCH}_2\text{CH}_2\text{COOH}$. When the NHS-ester of each of these organometallic compounds reacted with oligomers that were synthesized to possess a hexylamine on the 5'-end, the oligo became labeled and was purified by HPLC. When the M13 primer (17-mer) contained either the Fe or Sn label it functioned normally as a primer when tested in the ^{32}P -based Sanger sequencing procedure. Detection of Fe- or Sn-labeled oligomers has been demonstrated using resonance ionization mass spectrometry (RIS). RIS may be used in three different modes and the relative merits of each procedure will be discussed. When ^{57}Fe -labeled 40-mer was placed on a gold-foil we detected 0.02 fmol with a 100 μm ionization beam but only 2 pmoles when it was placed on GeneScreen. The sensitivity was limited in both cases by surface contamination with iron. Contamination with tin is much lower than for iron; also tin has ten isotopes as compared to four for iron. By using all the readily available isotopes of elements that are easily detected by RIS over 50 labels are available to perform multiplex analysis of DNA sequences. As compared to standard methods of DNA sequencing the stable isotope method could be 10-100 times faster due to such multiplexing possibilities. (Research sponsored by OHER, U.S. DOE under contract DE-AC05-84021400 with the Martin Marietta Energy Systems, Inc.)

To be presented at the DOE Contractor-Grantee Workshop, November 3-4, 1989, at Santa Fe, New Mexico.

LASER-BASED DETECTION IN ELECTROPHORESIS. Edward S. Yeung, Ames Laboratory, Iowa State University, Ames, IA 50011.

Traditionally, visualization in electrophoresis and in blotting requires either a radioactive tag or a fluorescence tag. This complicates the mapping and sequencing process by introducing a derivatization step and increases cost by requiring expensive primers. Fluorescence tags can further affect the migration characteristics of the fragments during electrophoresis. We will show a detecting scheme based on "indirect fluorescence" to avoid derivatization. The normal electrophoresis buffer solution is replaced by one that contains a fluorescing ion. There is then a large and uniform fluorescence background in the gel. Where the DNA bands reside, due to charge displacement, there appear dark regions because fewer of the fluorescing ions are present. The separation process and the spatial information are identical to conventional electrophoresis; only the detection scheme is different.

The information on the gel is acquired in real time by a charged-coupled device array detector or by scanning a laser beam across the gel. Good sensitivity is obtained. The latter mode also provides for "smart scanning" to maximize the information content. The same concept has also been successfully applied to capillary zone electrophoresis, where extremely high separation powers are obtained. One can detect fragments in the 50 attomole range without tagging.

SINGLE MOLECULE DETECTION IN FLOWING SAMPLE STREAMS AS AN APPROACH TO DNA SEQUENCING

James H. Jett, Lloyd C. Davis, Jong Hoon Hahn, Richard A. Keller,
Letitia Krakowski, Babetta Marrone, John C. Martin, Robert Ratliff,
Newton K. Seitzinger, and E. Brooks Shera

Los Alamos National Laboratory
Los Alamos, NM 87545
(505) 667-3018

We are exploring a technique which has the potential to sequence large fragments of DNA at a rate of hundreds of bases per second. Our technique is based upon a projected ability to detect single chromophores by laser-induced fluorescence in flowing sample streams.¹ The technique involves: (1) labeling the nucleotides with base specific tags suitable for fluorescence detection, (2) selecting a desired fragment of DNA, (3) suspending the *single* DNA fragment in a flowing sample stream, (4) sequentially cleaving labeled bases from the free end of the DNA fragment using an exonuclease, and (5) detecting and identifying the cleaved, labeled bases as they flow through a focused laser beam.²

The rate that bases can be sequenced is determined by the kinetics of the exonuclease cleavage reaction and the time required for detection and identification of the labeled bases. Based upon our results for the detection of rhodamine-6G and studies of cleavage rates, we anticipate sequencing rates of several hundred bases per second on a single strand of DNA tens of kb in length.

References:

1. D. C. Nguyen, R. A. Keller, and M. Trkula, "Ultrasensitive Laser-induced Fluorescence Detection in Hydrodynamically Focused Flows", *J. Opt. Soc. Am. B.*, **4**, 138 (1987).
2. J. H. Jett, R. A. Keller, J. C. Martin, B. L. Marrone, R. K. Moyzis, R. L. Ratliff, N. K. Seitzinger, E. B. Shera and C. C. Stewart, "High-Speed DNA Sequencing: An Approach Based Upon Fluorescence Detection of Single Molecules", *J. Biomolecular Structure & Dynamics*, in press.

High-Sensitivity Single-Molecule Fluorescence Detection in Theory and Practice

Richard A. Mathies and Konan Peck, Chemistry Department,
University of California, Berkeley, CA 94720

The number of emitted photons that can be obtained from a fluorophore increases with the incident light intensity and the duration of illumination. However, saturation of the absorption transition and photodestruction place natural limits on the ultimate signal-to-noise ratio that can be obtained. Equations have been derived to describe the fluorescence-to-background-noise ratio in the presence of saturating light intensities and photodestruction.¹ The fluorescence lifetime and the photodestruction quantum yield are the key parameters that determine the optimum light intensity and exposure time. To test this theory we have performed single molecule detection of phycoerythrin (PE). The laser power was selected to give a mean time between absorptions approximately equal to the fluorescence decay time. The transit time was selected to be nearly equal to the photodestruction time of ~600 μ s. Under these conditions the photon count distribution function, the photon count autocorrelation function, and the concentration dependence clearly show that we are detecting bursts of fluorescence from individual fluorophores. A hard-wired version of this single-molecule detection system was used to measure the concentration of PE down to 10^{-15} M.² This single-molecule counter is three orders-of-magnitude more sensitive than conventional fluorescence detection systems. The approach presented here is now being applied to the optimization of fluorescence-detected DNA sequencing gels.

1. Mathies, R. A., Peck, K. & Stryer, L. (1989) *Biop. J.* in preparation.
2. Peck, K., Stryer, L., Glazer, A. & Mathies, R. A. (1989) *Proc. Natl. Acad. Sci. U.S.A.* **86**, 4087-4091.

CLEAVAGE OF DNA INTO BIG PIECES: Application of DNA methylases.

John Hanish, Mike Nelson, Mike Weil, Bo-Qin Qiang, Saibal Poddar, Yogesh Patel and Michael McClelland.

California Inst. of Biological Research, 11089 N. Torrey Pines Rd. La Jolla, CA 92037.

Three methods for altering the cleavage of restriction endonucleases are described. All these methods involve DNA methylation.

1) CONTROLLED PARTIAL DIGESTION BY METHYLASE/ENDONUCLEASE COMPETITION REACTIONS. A restriction endonuclease is used in excess of that required to produce complete cleavage of DNA embedded in an agarose plug and the endonuclease is competed with increasing amounts of a DNA methylase that completely protects DNA from cleavage by the endonuclease. Both enzymes diffuse together into the agarose, producing a defined partial digest that depends not on the amount of endonuclease, nor the amount of DNA, but merely on the ratio of methylase to endonuclease. Of particular interest are partial digests with endonucleases, such as *NotI*, *MluI* and *NruI*, that cleave infrequently in the human genome. We have applied this method for obtaining reliable partials to pulsed field gel electrophoresis (PFE) of model systems, genomic DNAs of bacteria, using *M-BspRI* (GG^mCC) versus *NotI* (GCGGCCGC) and *M-EnuDII* versus either *MluI* (ACGCGT) or *NruI* (TCGCGA). *M-BspRI* and *M-EnuDII* were highly purified from clones to permit their use with restriction endonucleases in the presence of Mg²⁺. In principle, the method can be applied to any endonuclease.

2) TWO-FOLD INCREASE IN THE SPECIFICITY OF *NotI*. The apparent specificity of the "rare-cutting" restriction endonuclease *NotI* (GCGGCCGC) has been increased by the use of the *M-EnuDII* methylase clone (^mCGCG). This methylase blocks *NotI* cleavage at overlapping GCGGC^mCGCG sequences. As a result of this "cross-protection" by *M-EnuDII* methylation, *NotI* only cleaves at the sequence (A/T/G)GCGGCCGC(A/T/G). The number of fragments produced by *NotI* is reduced about two-fold. The method has been applied to dramatically simplify the pattern produced by *NotI* cleavage of the genomes of *E. coli*, *B. subtilis* and human, as visualized by PFE. Cross-protection and enhancement of the apparent specificity of a restriction endonuclease is possible in most cases by selecting a known methylase with an overlapping specificity.

3) METHYLASE/*DpnI* CLEAVAGE AT A TWELVE-BASE-PAIR SEQUENCE. A highly specific enzymatic DNA cleavage strategy that cuts at TCTAGATCTAGA has been demonstrated. The method employs *M-XbaI* (TCTAG^mA) methylase (cloned by E. Van Cott and G.G. Wilson, New England Biolabs) and the methylation-dependent restriction endonuclease *DpnI* (G^mA/TC), (Cloned by Lacks et al., Brookhaven Natl. Lab.) This method should produce fragments averaging 4¹² (16,000,000) base pairs on a random DNA sequence. However, TCTAGATCTAGA may occur less than once every 1,000,000,000 base pairs in bacterial genomes because the tetranucleotide CTAG is very rare in these genomes. We are able to cleave exclusively at *M-XbaI/DpnI* sites that have been engineered into a transposon that has, in turn, been integrated into the genome. The twelve-base-pair cleavage specificity have been successfully applied to *E. coli* and *Salmonella typhimurium* by the use of *M-XbaI* methylation *in vivo* and *DpnI* cleavage *in vitro*. We can choose to cleave such a genome one or more times, depending on the number of transposition events. The resulting large DNA fragments have been visualized using PFE. Eight- and ten-base-pair methylase/*DpnI* cleavages have previously been demonstrated on PFE (Weil and McClelland).

Synthetic Endonucleases

Betsy M. Sutherland¹ and Gary Epling²

¹Biology Department, Brookhaven National Laboratory, Upton NY
and ²Chemistry Department, University of Connecticut,
Storrs CT

Recognition and mapping of functionally important DNA regions (e.g. regulatory and coding regions, initiation sequences) can be greatly facilitated by specific DNA cleavage at such sites. Synthetic endonucleases able to cleave at regions of functional importance will be created by coupling DNA site-specific binding proteins via linker arms to light-activatable cleaving moieties: specific binding function is provided by the DNA binding protein, and cleavage activity by the activatable cleaving molecules. A prototype system of Rose Bengal coupled via a hexanoic acid linker to a DNA lesion site-specific monoclonal antibody will be developed for other specific DNA binding proteins, including T7 RNA polymerase and mammalian transcription initiation factors. Additional cleaving groups, linkers and coupling procedures will be developed to optimize reaction with DNA binding proteins of differing surface groups and reactivities. The specificity and efficiency of binding and cleavage of each synthetic endonuclease will be using a new electronic imaging system. Their utility in mapping, cloning and sequencing human DNA will be evaluated.

Abstract: The Sequence-Selective Hydrolysis of Duplex DNA by an Oligonucleotide-Directed Nuclease. David R. Corey*, Dehua Pei, Peter G. Schultz

Sequence-selective hybrid nucleases have been synthesized by fusing staphylococcal nuclease to an oligonucleotide via a disulfide linkage. The resulting hybrid nucleases hydrolyze both single-stranded DNA and RNA adjacent to the site of hybridization. However, the selective cleavage of single stranded DNA only occurs within a narrow range of temperatures, reaction times, and substrate concentrations. Mutagenesis was done on staphylococcal nuclease to lower its V_{max}/K_M , which makes its activity more dependent on the attached oligonucleotide. A mutation, Y113A, greatly decreased the amount of nonselective cleavage, but did not reduce hydrolysis at the desired target sequence, presumably because hybridization keeps the local concentration of substrate high. Hybrid nucleases synthesized with this mutant nuclease hydrolyze substrates under a much wider range of conditions than those lacking the mutation. The optimized enzymes function catalytically at elevated temperatures, and selectively turnover substrate with rates as high as 30 m^{-1} . More recently this methodology has been applied to the cleavage of double stranded DNA, and hybrid nucleases have been shown to efficiently cleave supercoiled plasmids pUC 19 and M13mp19. This was accomplished by introducing the oligonucleotide-directed hybrid nuclease into double-stranded DNA via D-loop formation using 3.5 mM NaOH to partially denature substrate. The use of Topoisomerase I in conjunction with ethidium bromide introduces additional supercoiling into substrate plasmids and greatly increases the efficiency of the procedure. This method allows the stoichiometric cleavage of microgram amounts of substrate DNA, and leaves 3' and 5' termini which can be enzymatically manipulated in subsequent reactions. Further work on this project will involve developing conditions which permit the selective hydrolysis of chromosomal DNA. This requires either (a) the discovery of methods for the supercoiling of chromosomal DNA within an agarose matrix or (b) the discovery of methods which allow the efficient incorporation of the oligonucleotide-directed nuclease into relaxed DNA.



MOTION OF LARGE DNA FRAGMENTS IN CROSSED OSCILLATING ELECTRIC AND MAGNETIC FIELDS (COEMF).

Gunter A. Hofmann and Sukhendu B. Dev, BTX Inc., 3742 Jewell St., San Diego, CA 92109, USA.

Pulsed Field Gel Electrophoresis (PFGE) technique is a Coulomb force (acting on the electric charges of the DNA) based separation technique which appears to be limited to resolving DNA bands up to a maximum of 12 megabases and can be a very slow procedure for large DNA fragments. To overcome these limitations, BTX is developing a Lorentz force (acting on the movement of electric charges in a magnetic field) based separation technique. Its principle is based on the observation that DNA exhibits a large induced dipole moment at low frequencies with a relaxation time dependent on the length of the fragment. The polarizability of DNA fragments is utilized by subjecting them to an oscillating electric and perpendicularly superposed oscillating magnetic field in a liquid without gel matrix. The resulting unidirectional Lorentz force moves the DNA fragments perpendicular to both the electric and the magnetic field vectors with a drift velocity dependent on the DNA polarizability, which depends on DNA size and configuration. We demonstrated that DNA fragments can be polarized by an oscillating electric field and that a superposed magnetic field results in an unidirectional drift velocity of these fragments. Drift experiments have been run for DNA fragments varying from 48 Kb lambda, T2 and T5 phage DNA (166 Kb and 100 Kb), ultrapure genomic DNA (115 Kb) and *Candida albicans* (3 Mb). The drift velocity increases with the size of the DNA fragment, in contrast to conventional electrophoresis. For large DNA fragments, the drift velocity can be expected to be two to three orders of magnitude higher than the drift velocity observed in PFGE experiments. For a 1 Mb fragment we observed a COEMF mobility of $6.4 \times 10^{-4} \text{ cm}^2/\text{Vs}$ which is 30x higher than the corresponding mobility in PFGE. The COEMF mobility of DNA fragments, especially fragments with low molecular weights, can be greatly affected by the addition of detergent, poly-L-amino acids or a cationic polyelectrolyte. The fragment length of an unknown DNA, derived from the Human Carcinoma Tumor, has been correctly predicted from a 10 minute COEMF drift experiment (120-140 Kb) and confirmed by a 24 hour PFGE run.

*Research supported by DOE SBIR Contract No. DE-ACO3-88ER80655. Patent allowed.

DOE
Human Genome
Contractors/Grantee
Workshop

November 3-4, 1989

Santa Fe Institute

Santa Fe, NM

An Expanding Cosmid Contig Map for Chromosome 19.

L. Ashworth, E. Branscomb, L. Brown, C. Chen, P. de Jong, A. Fertitta, E. Garcia, J. Garnes, J. Lamerdin, F. Lohman, H. Mohrenweiser, D. Nelson, W. Nelson, A. Olsen, B. Perry, T. Slezak, K. Tynan, M. Wagner, P. Wilkie and A.V. Carrano. Biomedical Sciences Division, Lawrence Livermore National Laboratory, Livermore, CA.

High resolution gel electrophoresis provides a simple and versatile method for DNA fingerprinting and the creation of contigs or sets of overlapping genomic clones. Cosmid libraries are constructed from YAC clones or from flow-sorted chromosomes. Cosmid DNA is isolated by an alkaline-lysis procedure and each cosmid is digested with EcoRI to measure insert size and concentration. The isolated DNA is then cut with a combination of five restriction enzymes and the fragment ends labeled with one of four different fluorochromes. Our approach to contig construction: 1) uses a robotic system to label restriction fragments from cosmids with fluorochromes; 2) uses an automated DNA sequencer to capture fragment mobility data in a multiplex mode (i.e. three cosmids and a size standard in each lane); 3) processes the mobility data to determine fragment length and provide a statistical measure of overlap among cosmids; and 4) displays the contigs and underlying cosmids for operator interaction and access to a database. We have applied these methods to construct a cosmid contig map for a 600 kbp YAC clone from chromosome 14 and are currently analyzing cosmids to construct contigs for all of chromosome 19. Throughput rate is currently about 48 cosmids per day per machine but 96 cosmids per day is achievable. Resolution of fragment size is to within 1-1.5 bases over the range of 29-462 bases for which data are captured. The more than 2500 chromosome 19 cosmids analyzed to date assemble into over 320 contigs with an average contig length of 3.2 cosmids. Many of these contigs have been located to the chromosome by fluorescence *in situ* hybridization and also mapped to known genes. The "minimal" spanning sets of cosmids provide unique starting material for genome sequencing. This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract No. W-7405-ENG-48.

PROGRESS OF "TOP-DOWN" MAPPING APPROACHES TO CHROMOSOME X
David F. Barker, Arnold R. Oliphant, Pamela R. Fain, David E. Goldgar, Huntington F. Willard, Anne Vincent, Stephen Warren, Jennifer Puck, Robert L. Nussbaum and Christine Petit

The initial focus of work aimed at defining the genetic and molecular structure of the X chromosome has been the isolation and ordering of a set of clones defining RFLP markers. We have isolated more than 80 such markers for the X utilizing an X library, LAOXNL01, constructed by LANL under the auspices of the DOE human genome project. If combined with the over 100 other X RFLP markers that have been isolated in many laboratories, these markers would comprise a set with an average spacing better than 1 per megabase. The complete and reliable ordering of such a set would provide a useful "backbone" structure for both high resolution genetic-disease mapping and the ordering of sets of overlapping clones, "contigs", with respect to each other.

The ordering approaches we have used include genetic mapping in the CEPH linkage reference families and in genetic disease families. We have also tested a variety of human-rodent hybrid cell lines containing unique segments of the X chromosome to provide additional ordering information. The physical breakpoints have been generated by natural translocations or deletions or by strategies designed to select broken chromosomes in tissue culture lines. The latter include "pushmi-pullyu" hybrids generated by Brown et al. (HGM10) and "radiation hybrids" isolated as described by Cox et al. (AJHG 43: A141). The physical breakpoints characterized to date divide the chromosome into 25 regions, 8 on Xp and 17 on Xq. Nearly all of the 80 polymorphic markers which we have isolated have been mapped into one of these 25 regions and a summary map will be presented. The current status of the genetic linkage map will also be shown.

TOWARDS CONSTRUCTION OF A 2Mb PHYSICAL MAP OF HUMAN CHROMOSOME 16.

D.F. CALLEN, V.J. HYLAND, L.Z. CHEN, J.C. MULLEY, S. LANE, E.G. BAKER, G.R. SUTHERLAND

Cytogenetics Unit, Adelaide Children's Hospital, North Adelaide, South Australia 5006.

The use of mouse/human hybrids containing portions of human chromosome 16 provides a method for rapid physical mapping. We aim to extend our hybrid panel of this chromosome until the chromosome can be subdivided into approximately 50 intervals. At this point the average interval spanned by breakpoints will be 2Mb.

The human parent used in the construction of these hybrids are chromosome 16 translocations and deletions which have been identified as the result of cytogenetic investigations. These have been obtained from our laboratory, from the NIGMS cell repository and from the kind co-operation of many other cytogenetic laboratories. Each newly generated hybrid containing a derived human chromosome 16 is evaluated using a battery of probes which have been already physically mapped to chromosome 16. This allows the ordering of the breakpoints of the new hybrid in relation to the existing physical map.

In conjunction with the generation and evaluation of new hybrids further probes on chromosome 16 will be physically mapped. These include gene probes, probes that we have generated from a lambda library derived from the hybrid CY3, single-copy probes on chromosome 16 obtained from Dr. P. Harris, chromosome 16 cosmid clones from Dr. M. Breuning, and probes which have been genetically mapped in the CEPH pedigrees from Dr. C. Julier. It is planned to map other polymorphic probes which have been genetically mapped from Dr. P. O'Connell and to locate cosmid contigs in collaboration with Dr. E. Hildebrand. This will lead to a detailed correlation of the genetical and physical maps of chromosome 16.

At the present time we have constructed twenty hybrids of chromosome 16 which, with the use of the three rare fragile sites on this chromosome, can subdivide the chromosome to potentially 24 intervals. At the present time 83 protein markers, gene probes and anonymous DNA fragments have been physically mapped to all, or a subset of, these hybrids.

This hybrid panel can provide a useful resource for work involving other chromosomes. Breakpoints on chromosomes 1, 3, 4, 9, 10, 11, 12, 13 and 22 are included in the panel. All hybrids are available from Dr. David F. Callen on request.

This work is supported by the Department of Energy Grant DE-FG02-89 ER60863. This support does not constitute an endorsement by DOE of the views expressed in this abstract.

Generation of a Physical Map of the Long Arm of Human Chromosome 11.

G.G. Hermanson*, P. Lichter#, D.C. Ward#, and G.A. Evans*.

*Molecular Genetics Laboratory, The Salk Institute, La Jolla, CA 92138

#Department of Human Genetics, Yale University School of Medicine,
New Haven, CT 06510

Many loci implicated in human diseases have been mapped to the long arm of human chromosome 11 including: ataxia telangiectasia, tuberous sclerosis, multiple endocrine neoplasia type 1, and translocations found in Ewing's sarcoma, and acute leukemias. This region also contains genes which are members of the immunoglobulin superfamily: NCAM, CD3 γ , δ , and ϵ , and Thy-1, as well as the oncogene *c-ets-1* and the leukocyte marker CD5. Given that only a small proportion of human genes have been identified and mapped, it is clear that many more genes which might be important in human development or disease may be contained in this region. To generate a physical map and localize these potentially important genes, we have isolated linking clones containing multiple rare-cutting restriction enzyme sites. These clones can be used as probes for pulsed-field gel and fluorescent *in situ* hybridization techniques to physically map this region, and identify potential HTF islands that are associated with many gene coding regions. In order to identify these linking/HTF island clones, a cosmid library was constructed from a somatic cell hybrid line containing only the long arm of human chromosome 11 from 11q12-11qter. Cosmid clones containing human insert DNA were selected by hybridization to total human genomic DNA and picked into a total of ten 96-well plates for further analysis. A Beckman robotic workstation was used to prepare miniprep DNA from this cosmid collection, as well as to assay the 960 clones for the presence of the rare cutting restriction enzyme sites Not I, Sac II, BssH II, Pvu I, Mlu I, and Sfi I. A total of 175 cosmid clones contained at least one Not I site in their insert. Since Not I sites are rare in the genome, cosmids were initially selected for further analysis only if they contained at least one of these sites. Thirty-two unique cosmids were finally selected that contain at least one Not I site and sites for the majority of the other rare cutting enzymes. These 32 linking/HTF island cosmids, as well as cosmids containing previously identified genes are being ordered into a molecular and cytogenetic map of chromosome 11 by single copy fluorescent *in situ* hybridization and pulsed-field gel electrophoresis techniques.

CONSTRUCTION OF HUMAN CHROMOSOME 21 SPECIFIC YEAST CHROMOSOMES

Mary Kay McCormick^{1,2}, James H. Shero⁴, Mei Chi Chung³, Yuet Wai Kan³, Philip Hieter^{2,4}, Stylianos E. Antonarakis^{1,2}

¹Genetics Unit, Department of Pediatrics, ²Predoctoral Training Program in Human Genetics, Johns Hopkins University School of Medicine, Baltimore, MD 21205

³Howard Hughes Medical Institute, University of California, San Francisco, CA 94143

⁴Department of Molecular Biology and Genetics, Johns Hopkins University School of Medicine, Baltimore, MD 21205

Chromosome 21 specific yeast artificial chromosomes (YACs) have been constructed by a method that performs all steps in agarose, allowing size selection by pulsed field gel electrophoresis and the use of nanogram to microgram quantities of DNA. The DNA sources were hybrid cell line WAV-17, containing chromosome 21 as the only human chromosome, and flow sorted chromosome 21. The transformation efficiency of ligation products was similar to that obtained in aqueous transformations and yielded YACs with sizes ranging from 100 kilobases to over 1 megabase when polyamines were included in the transformation procedure. Twenty five YACs containing human DNA have been obtained from the mouse-human hybrid, ranging in size from 200 kb to > 1000 kb with an average size of 410 kb. Ten of these YACs were localized to sub regions of chromosome 21 by hybridization of riboprobes (corresponding to the YAC ends recovered in *E. Coli*) to a panel of somatic cell hybrid DNAs. Twenty one human YACs, ranging in size from 100 kb to 500 kb with an average size of 150 kb, were obtained from ~50ng of flow sorted chromosome 21 DNA. Three were localized to subregions of chromosome 21. Yeast artificial chromosomes will aid the construction of a physical map of human chromosome 21 and the study of disorders associated with chromosome 21 such as Alzheimer's disease and Down syndrome.

Assembly of Cosmid Contigs in the Region of the Nucleotide Excision Repair Genes on Human Chromosome 19. Mohrenweiser, H. W., de Jong, P. J., Perry, B. A., Tynan, K. T., Lohman, F. P. and Carrano, A. V. Biomedical Sciences Division, L-452, Lawrence Livermore National Laboratory, Livermore, CA 94550

The genes CKMM, ERCC1 and ERCC2 have been assigned to the region q13.2 - q13.3 of human chromosome 19 and have been localized to within ~250 kb of each other by PFG electrophoresis. A human chromosome 19 specific cosmid library was screened with a pool of probes for these genes. The DNA from each of the selected cosmids was analyzed using a high density restriction enzyme site mapping ("fingerprinting") strategy (Carrano et al., Genomics 4:129, 1989). Overlapping cosmids were detected, and contigs assembled, based upon commonality of restriction enzyme digest fragment sizes. Two contigs could be generated from the fingerprinting data obtained from analysis of the group of cosmids initially selected by probing. Upon reprobing of this group of selected cosmids with the individual probes, the 2 cosmids in one contig hybridized with the ERCC1 probe while the second contig of 3 cosmids contained the CKMM gene. Four additional cosmids, analyzed during the fingerprinting of ~2200 random cosmids (~1.4X library), have been linked to the most centromeric cosmid of the original CKMM contig yielding a contig with a tiling path of 6 cosmids. No cosmids containing the ERCC2 gene were isolated and no randomly selected cosmids have been linked to the most telomeric CKMM cosmid, thus a gap of ~10kb exists between the CKMM contig and a previously isolated set of cosmids containing the ERCC2 gene. The ERCC2 gene is ~150kb from the ERCC1 contig. An expressed sequence overlapping the ERCC1 gene is within this region, thus at least 4 expressed genes, comprising >60kb, are within this 250kb region. Closure of the gaps, as necessary to form a single contig containing the ERCC1, ERCC2 and CKMM genes should be attained with very limited walking, although it will apparently be necessary to screen additional libraries (YAC, lambda, cosmid) to complete the contig.

Work performed under auspices of the US DOE by the Lawrence Livermore National Laboratory; contract No.W-7405-ENG-48

Contig Assembly and Characterization of a Chromosome 19q Specific Minisatellite Element. Tynan, K. T., Mohrenweiser, H. W., Branscomb, E. W., deJong, P. J. and Carrano, A. V. Biomedical Sciences Division, L-452, Lawrence Livermore National Laboratory, Livermore, CA 94550

A minisatellite consisting of a repeat of six 37 bp elements was identified by Das et al. (J Biol Chem 262: 4787, 1987) in intron six of the ApoCII gene, a gene mapping to 19q13.2. Five elements of this minisatellite are present in the ERCC2 locus, another locus on chromosome 19q13.2 (Weber et al., EMBO J in press). It has been estimated that this minisatellite exists at ~60 "loci." Additional evidence suggests this minisatellite is localized to the q13 region of chromosome 19. Approximately 150 cosmids were identified as containing this minisatellite during the screening of a human chromosome 19-specific cosmid library (~7000 cosmids) with a probe containing the repeat element. Sixty seven of these cosmids have been analyzed by fingerprinting (Carrano et al., Genomics 4:129, 1989) during the analysis of ~2200 random cosmids. Twenty six of the 67 cosmids have been assigned as members of 16 contigs. Three of these contigs are comprised of only 2 cosmids, both members of each pair being minisatellite positive. Of the remaining 14 contigs, the minisatellite is present as 3 adjacent cosmids in 3 contigs, a pair of overlapping cosmids in one additional contig and only once in the 9 other contigs. As expected, this appears to be an efficient strategy for establishing "seed" contigs in this region of chromosome 19 and also assists in validating the contig building strategy. One of the minisatellite containing cosmids is within the PVS locus contig at 19q13.2, thus at least 3 of the minisatellite repeat "loci" are located within functional genes. A number of additional minisatellite containing cosmids have been mapped by *in situ* hybridization and all map to the 19q13.2-q13.4 region. Hybridization analysis of DNA from a human lymphoblastoid cell line and hamster cells containing all or parts of human chromosome 19 following either normal or PFG electrophoresis are consistent with the previous estimate of ~60 loci and localization of these elements to human chromosome 19q.

Work performed under auspices of the US DOE by the Lawrence Livermore National Laboratory; contract No.W-7405-ENG-48

Abstract for DOE Contractor-Grantee Workshop November 3-4, 1989

CHARACTERIZATION OF HUMAN CHROMOSOME-SPECIFIC PARTIAL DIGEST LIBRARIES IN LAMBDA AND COSMID VECTORS.

Kathy Yokobata, Jennifer McNinch, Lee Pederson, Marvin A. Van Dilla and Pieter J. de Jong, Biomedical Sciences Division, Lawrence Livermore National Laboratory, P. O. Box 5507, Livermore, CA 94550

As part of the National Gene Library Project we have constructed partial digest chromosome-specific human genomic libraries for studies of genetic disease, physical mapping of chromosomes, and other studies of the molecular biology of genes and chromosomes. New procedures were developed for isolating DNA of high molecular weight from flow sorted human chromosomes and for preparing large insert libraries in lambda replacement and cosmid vectors from small quantities (about 1 μ g) of DNA. These procedures have been used successfully for the preparation of lambda and cosmid libraries specific for chromosomes 19, 21, 22 and Y. A large insert lambda library has been prepared from sorted chromosome 11 DNA as well. The lambda vectors used were Charon 40 and GEM 11; the cosmid vector was Lawrist 5, a lambda origin cosmid vector.

For all libraries, we have estimated purity by flow karyotype analysis. We have also examined the purity of the lambda libraries by plaque hybridization and have examined the origin of clones which show no signal in these plaque hybridizations. The average insert size was determined by excising insert DNA from randomly selected clones for each lambda library. In addition, we have screened the chromosome-19 lambda library with probes known to map to the chromosome to establish representation of these sequences in the library.

We have begun characterization of the cosmid libraries by mapping randomly selected cosmid clones by fluorescent *in situ* hybridization. We have examined the stability of insert sequences in cosmid libraries in various bacterial hosts and describe studies with an unstable cosmid clone from the CKM locus on chromosome 19 and with insert sequences in the chromosome-Y library. Much of the characterization of the chromosome-19 cosmid library is being done under the auspices of the chromosome 19 ordering project.

The libraries for chromosomes 11 and 19 were made from chromosomes sorted from monochromosomal hybrid lines, whereas the libraries for 21, 22 and Y used human cell lines for the starting material. Because of the purity advantages of obtaining sorted chromosomes from monochromosomal hybrids and the new availability of sortable hybrids for these chromosomes, new 21, 22 and Y libraries are being constructed. We are currently isolating chromosomal DNA from chromosomes 3 and 12 for lambda and cosmid libraries.

Work performed under the auspices of the U.S. Department of Energy by the Lawrence Livermore National Laboratory under contract number W-7405-ENG-48.



October 12, 1989

Dr. Sylvia Spengler
LBL - Human Genome Center
459 Donner Laboratory
Berkeley, CA 94720

Dear Dr. Spengler:

The following is in response to your letter dated September 26, 1989. It is an abstract for the poster session at the DOE Human Genome Project Contractor/Grantee workshop.

Title: A Natural Language Query System for Genbank

Authors: Michael Cinkosky (LANL), Rowland R. Johnson (LLNL)
and Rob Pecherer (LANL)

Abstract:

The Genbank database contains a wide variety of information that is of interest to researchers involved in the Human Genome Project. The logical structure of the database is complex enough that procedures to query the database are required. Such procedures have been implemented or will be implemented for the most typical classes of queries. Atypical queries require a procedure to be constructed by the end user. It is unlikely that a Human Genome Project researcher will be able to construct a procedure that will satisfy an atypical query. Thus, the usefulness of Genbank to the Human Genome Project will be restricted.

A possible solution to this problem is a system that will accept a query stated in English and translate it to a procedure appropriate for the database system. Towards this end, a prototype of such a system was developed in July 1989 for the Human Genome Project at LANL. This system is based on a commercially available product that translates English into a sequence of relational database query commands.

The work undertaken in the development of the prototype was to customize the product to work in the Genbank domain. This customization is a continuing, but diminishing, effort. New users of the system will bring out unrealized English constructions that must be incorporated. As more and more people use the system, there will be fewer unrealized English constructions to incorporate.

Currently, the prototype contains a subset of the data that will be in Genbank. The system will be available at the conference so that participants may try queries on the system.

P09

Title: GnomeView: A Graphics Interface to the Human Genome

Authors: Richard J. Douthart, David A. Thurman and Victor B. Lortz

Abstract:

Pacific Northwest Laboratory is developing GnomeView, a software system that provides a graphical interface to the large quantities of data and information generated by the Human Genome Initiative. GnomeView allows the user to visually browse and manipulate color graphic representations of genetic maps, physical maps, and sequences. These representations provide the user with a sense of topology and reveal patterns in the data that are otherwise difficult to detect. This hierarchy of mappings can be traversed using the pan-and-zoom capabilities of GnomeView and all objects in maps will be queryable.

GnomeView uses the X Window System to render and display maps and sequences on a UNIX workstation. It is being developed on a Sun workstation in portable C and should be portable to any UNIX workstation. It includes db_VISTA, a network-model database, that permits natural and efficient representation of the relationships in genomic information. Landmarks, features, and blocks of sequence are stored in linked lists of objects in the database. Objects common to more than one mapping need only be stored once in the database, but may be referenced in many mappings. Specific attention has been paid to designing database algorithms and user interface techniques that scale well to high data volumes.

GnomeView accepts queries from the user and responds in a number of different ways. Depending on the query, these responses can take the form of maps, textual lists, or color-coded histograms. Further information can be obtained by querying these responses. As an example, the hierarchy of the superoxide dismutase loci on Chromosome 21, from band location, to restriction map, to base sequence, will be presented.

Identification of genes in anonymous DNA sequences

C. A. Fields and C. A. Soderlund

Computing Research Laboratory, New Mexico State University, Las Cruces, NM 88003-0001, USA. 505-646-5466.

The objective of this project is the development of practical software to automate the identification of genes in anonymous DNA sequences from the human, and other higher eukaryotic genomes. A prototype automated sequence analysis system, **gm**, has been implemented in C to run on Unix workstations. This system accepts as input: i) a DNA sequence, ii) consensus matrices for locating splice sites, translational start sites, and polyadenylation sites, iii) match-quality cutoffs for consensus searches, and iv) base frequency and codon usage standards for coding regions and introns. It produces as output both schematic models of possible genes contained in the sequence that show the locations of the coding sequences, introns, and control signals, and predicted amino-acid sequences for each of these possible genes. The models include the numerical results of evaluating each of the component exons and introns.

gm has been extensively tested on *C. elegans* sequences in the 10 kb size range containing known genes of up to 10 exons, and is capable of generating complete, correct analyses showing all possible alternative splicing patterns. Such analyses typically require a few minutes running time on a Sun 4/60 workstation, depending on the stringency of the search parameters used. Current effort is focussed on improving the pattern recognition and statistical analysis modules used by **gm**, and on implementing greedy algorithms for performing fast first-pass analyses using low stringency parameters.

Sequence Matching and Motif Identification

Abstract

Daniel Gusfield

Computer Science Division, U.C. Davis *and*
ICS Division, Lawrence Berkeley Laboratory

Eugene L. Lawler William I. Chang

Computer Science Division, U.C. Berkeley *and*
ICS Division, Lawrence Berkeley Laboratory

Frank Olken

ICS Division, Lawrence Berkeley Laboratory

Sequence Matching and Alignment

Dynamic Programming has been the technique of choice for sequence matching problems arising in biology. However, DP algorithms are likely to be impracticable for problems of the size that will be created by the massive amounts of data generated by the Human Genome Initiative. On the other hand, there are other well developed techniques based on the use of suffix trees or finite state machines, which are more efficient than DP for many sequence matching problems. Our work is focussed on extending the use of suffix trees to additional problems in computational biology, with the objective of developing computational methods that are faster and more practicable than those which currently exist.

The principal difficulty in using suffix trees in biological applications is that they are best adapted for problems in exact, rather than approximate, sequence matching. However, we have found ways to use suffix trees to improve existing efficient algorithms for certain approximate sequence matching problems. For example, we have been able to use suffix trees to replace hashing in the Lipman-Pearson algorithms. By first building a suffix tree of the shorter sequence we can, in a single left-to-right scan of the longer sequence, compute for each position the longest exact match with any portion of the shorter sequence. This yields information which encompasses the Lipman-Pearson hashing step for all choices of tuple size simultaneously. This subroutine also enables us to simplify and make practicable a theoretical result of Landau and Vishkin, namely that matching up to k indel/substitution errors can be accomplished in time a factor k worse than linear. Our results further indicate that when the error threshold is below 15 percent for nucleotides (in particular errors arising from sequencing) almost all mismatches can be eliminated from consideration very quickly so that in linear time we can expect to find every match. Applications include the problem of computing overlaps in sequence assembly when sequencing errors are not negligible.

Another area of interest to us is that of sensitivity analysis. Existing algorithms for approximate matching require the user to specify mismatch (indel/substitution) penalties. However, these penalties are not known exactly. We have already had some success in developing efficient parametric algorithms which reveal the optimal sequence matching solution as a function of the penalty costs.

Searching for Motif Patterns

It is often possible to develop very efficient special purpose algorithms that search for repeat patterns in a sequence. Examples include the identification of *inverted* and/or *complemented* repeats with gaps (there are no a priori upper or lower bounds placed on gap lengths). For concreteness, a pair $..A....A'..$ is a *maximal inverted pair* if A' is an inverted copy of A and the sequence does not look like $..AB..B'A'..$ or $..BA....A'B'..$. Note that the same subsequence can be part of more than one maximal inverted pair: in the sequence $..abc.xba.cba..$, both $..abc....cba..$ and $..ab..ba.....$ may be maximal. The problem of identifying all such maximal pairs (including positions) can be solved in quadratic time and space by dynamic programming methods. However, we have recently developed a very simple method which finds all maximal inverted pairs in linear space and in time proportional to the length of the sequence plus the number of such pairs. Importantly, our method can be modified to find only pairs above a certain length.

DNA Bend Detection

Marjorie S. Hutchinson *
Information and Computing Science Division
Lawrence Berkeley Laboratory
Berkeley, Ca, 94720

We will describe and demonstrate a program on a Sun workstation that examines DNA sequences for bends or kinks. The program has the capability of scanning either a single sequence or a specified portion of Genbank to identify likely candidates for a bend. We have developed a user-friendly interface that allows the user to examine a candidate or list of candidates and view the parameters that indicate whether a bend might be present. Some of the more important output parameters are plotted against residue number and these plots are analyzed for the likelihood of a bend. Threshold values for the various tests for a bend can be set by the user. We display the calculated 3D structure of the sequence, allowing the user to rotate the structure, and zoom in on interesting portions.

The structural properties of the DNA sequence are derived using software provided by Wilma Olson of Rutgers University. ¹ Her program utilizes the potential energies of interaction of free base pairs (as calculated by Srinivasan et al.) ² to calculate a static structure. This structure is determined by choosing the local conformational geometries which optimize the computed average orientations of adjacent residues.

The Olson program also calculates several measures of chain stiffness including the persistence length, which is a measure of the distance over which the initial direction of the DNA is preserved. For each residue in the chain, it also calculates: the average angular orientations of the sequence, the average vectorial displacement, and the mean-square end-to-end distance. Variations in these values over the chain length may indicate with bends. We display this data and utilize Fourier transform coefficients of the resulting curves as further evidence for or against bending.

Execution of the program in its current stage of development will be demonstrated and planned improvements will be described.

*email: margeh@csam.lbl.gov

¹W.K. Olson and A.R. Srinivasan, (1988), The translation of DNA primary base sequence into three-dimensional structure. *Cabios*, 4, 133-142.

²A.R. Srinivasan, R. Torres, W. Clark and W.K. Olson. Base sequence effects in double helical DNA. I. Potential energy estimates of local base morphology, (1987) *J. Biomol. Struct. Dynam.*, 5, 459-496.

ALGORITHM FOR SEQUENCE GENERATION FROM K-TUPLE WORDS
CONTENT: Labat, I., Drmanac, R., Crkvenjakov, R. Genetic
Engineering Center, PO Box 794, 11000 Belgrade, Yugoslavia

Any text as well as nucleotide sequence, can be represented as a set of the overlapping k-tuple words, similar to the methods applied in the most efficient sequence comparison algorithms used today. Our algorithm uses intrinsic nucleotide sequence informatics to regenerate the original sequence without k-tuple position and frequency information inherent to the former algorithms. K-tuples are ordered by maximal overlapping up to the moment when none, or two or more k-tuples overlap with the last one attached. Further ordering is ambiguous. A primary subfragment (SF') is thus defined. Number of the heuristic methods developed repairs and unambiguously connects SF' into the real subfragments (SF). Number and length dispersion of the SFs as sequence informatics entities depends on length and simplicity of the sequence as well as length of k-tuples and extent of mistakes in the set. The other part of the algorithm enables regeneration of the sequence of the length of the human genome fragmented in the suggested manner (Drmanac et al., GENOMICS 4, '89, 114). It consists of several k-tuples sets and SFs manipulations on different, overlapping sequence fragments. Our software is applied on the IBM PC/AT compatible. In simulation experiment on the 50kb sequence, complete k-tuples sets (k=8 to 12, depending on GC content) of the consecutive nucleotide sequence fragments up to 900 bp were handled. In over 91% of analyzed fragments the complete sequences were regenerated. In remaining cases, sequences were regenerated to the level of several (below 15) SFs. Also, 10% of false negative k-tuples in the set makes no problem in most of analysed sequences. Further improvements of algorithm are needed (and suggested) for complete regeneration of some specific sequences, and for using sets with more false positive and false negative k-tuples.

Robotic Control of the Laboratory Procedure for Clone Candidate Selection

S. Lewis, J. Gingrich, and J.C. Bartley
Engineering and Life Sciences Divisions
Lawrence Berkeley Lab
Berkeley, California

We have automated a laboratory procedure to detect and select potentially transformed yeast cultures, by adapting the standard laboratory technique to a standard 96 well microtiter plate format and by programming a general purpose robot to carry out the procedure. The robot runs unattended and can screen up to 960 candidates in less than an hour.

The procedure steps are as follows:

- Initial transformed spheroplasts which have been grown up suspended in solid media are picked by hand into selective liquid media in the wells of microtiter plates.
- These microtiter plates are placed into the robot's incubator station and the robot application is started.
- After incubation the robot transfers each plate in turn to a plate reader. Growth is determined according to the turbidity in each well.
- For each well in which there is growth the robot pipets an aliquot into a new well of selective media, which the robot has previously filled. Only those wells in which growth has occurred are transferred, resulting in fewer plates.
- These plates are then placed in the incubator by the robot for subsequent incubation.

Data describing the identities and characteristics of the selected cultures are recorded for subsequent inclusion in a laboratory notebook database. We also discuss specific problems that arose automating this procedure such as: selection of well bottom shape (flat, u, or v), or suspension and aeration of the yeast growing in liquid media. A video tape of the procedure will illustrate the robot's capabilities and each individual outlined above.

Automated Extraction of Band Data from Digitized Autoradiogram Images

S. Lewis, K. Gong, J. Jaklevic, W. Johnston, E. Theil
Engineering Division
Lawrence Berkeley Lab
Berkeley, California

We are developing computer methods to analyze electrophoresis gel and autoradiogram images generated in a production mapping laboratory. The objective is to distill from the digitized images essential band size and intensity information. As electrophoresis gels display a lot of operational differences, the analysis software must handle problems such as: diffuse or sharp bands, overlapping bands, crooked lanes, over and under exposure times and varying DNA migration rates. Though the program provides fully automatic analysis of the bands if desired, there is the opportunity at every step for the operator to substitute his own judgement. For instance, the correct position of crooked lanes is indicated by the operator pointing and clicking the length of the lane. By isolating each component step in the analysis the flexibility of substituting either user results or improved algorithms is easily achieved.

The filtering techniques which we are using to detect bands are described. Currently this is a simple single pass non-adaptive digital filter which removes the background and slopes in the original one-dimensional data. The filtered 1-D data is then analyzed for peaks indicating bands. Once the bands have been located the sizes are calibrated according to their relative position within the lane.

The sets of band data thus derived is subsequently used for algorithmic comparisons between lanes and gels, as well as in mapping algorithms. The band data and pointers to the original archive images are stored in a laboratory data base.

ImageQuery: An Interface to a Biological Images and Videos Database

S. Lewis and W. Johnston —
Engineering and Computer Science Divisions
Lawrence Berkeley Lab
B. Morgan and S. Jacobson
Advanced Technology Planning
University of California at Berkeley
Berkeley, California

ImageQuery is a software tool which combines textual and visual methods for organizing and querying an image database. It was developed by the Advanced Technology Planning group at UC Berkeley to provide online access to collections of primarily visual materials. One of the initial applications at UCB provided information on archeological artifacts. In our case it has been adapted to catalog autoradiograms.

The ImageQuery interface enables researchers to search through images based upon either selectable fields of descriptive information using boolean logic, or to browse through iconic representations of the images themselves. ImageQuery icons represent the type of image and a unique identifier.

Once a set of images has been selected from the database via any combination of selection methods, ImageQuery can invoke specific analysis programs directly. For example, the program AnGel is used to locate and determine the sizes of bands within autoradiograms, or the program HIPsTools which provides users with zoom, pan, distance measurements, and other standard image processing capabilities.

ImageQuery runs on Sun work stations and can be interfaced either to a flat file of descriptive parameters, or to Ingres, a relational database management system. A demonstration of the query capabilities will be given using selected electrophoresis gel image data.

Supercomputer Simulations and Experimental DNA Electrophoresis*

Lim, H.A., Burnette, D.E., and McMullen, D.F.

Supercomputer Computations Research Institute

Florida State University

Tallahassee, FL 32306-4052

The main objective of this project is the development and optimization of general purpose supercomputer algorithms so that the mobility of large DNA molecules in electrophoresis can be simulated. By working in parallel with laboratory experimentalists, this project should integrate new supercomputer-based DNA electrophoresis simulating programs with innovations in electrophoretic techniques, and biochemical manipulations of chromosomal DNA. Though electrophoresis has been successfully used to separate chromosomes from lower eucaryotes (*e.g. Saccharomyces cerevisiae*) and other simple organisms (*e.g. Drosophila melanogaster*), the amount of DNA in the average human chromosome (about 1.4×10^8 bp) is at least 10-fold to 100-fold greater than current record size separable by electrophoresis. The project will focus on three major aspects of the problems associated with the current state of the art of electrophoresis: (1) Size—the current size record separable is still at least a factor of 10 smaller than the average size of human chromosome; (2) Speed—the current rate of data collection and analysis can still be improved; and (3) Resolution—errors in DNA sequence determination occur most commonly by insufficient resolution of DNA fragments in electrophoresis due to band inversion or compression. Unless this can be overcome, the human genome sequence determined by electrophoresis will be of little use.

* Funded by DOE Energy Research through SCRI

Human Genome Management Information System

Betty K. Mansfield, Judy M. Wyrick, John S. Wassom, Po-Yung Lu, Mary A. Gillespie, and Sandy E. McNeill

Health and Safety Research Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831-6050
(615) 576-6669, FTS 626-6669

The Human Genome Management Information System (HGMIS), sponsored by the U.S. Department of Energy (DOE) at the Oak Ridge National Laboratory, has the following roles in the Human Genome Program: (1) to assist the DOE Office of Health and Environmental Research (OHER) in communicating issues relevant to human genome research to DOE contractors and grantees and to the public and (2) to provide a forum for exchange of information among individuals involved in genome research or the development of instrumentation and methodologies to implement genome research. To fulfill these communications goals, HGMIS is producing technical reports, DOE Human Genome Program reports, a quarterly newsletter, and an electronic bulletin board. These documents/facilities are available to all interested persons upon request. The first technical report will assess instrumentation and methodology development relevant to DNA mapping and sequencing. The DOE Human Genome Program reports will contain information on human genome research and development activities supported by the DOE program as well as background information. The *Human Genome Quarterly* newsletter features technical articles, meeting reports, a calendar of genome events, announcements, and other information relevant to genome research. Accessible via modem through direct dial or via the user's host mainframe computer network, the electronic bulletin board contains information organized by categories (e.g., menu, general information, news and comments from OHER; summaries and highlights of research projects; meeting announcements/calendar; literature highlights; DOE Human Genome Program contacts; and international activities). HGMIS welcomes comments, suggestions, and contributions from the genome research community.

SDT : A DATABASE SCHEMA DESIGN TOOL *

Victor M. Markowitz and Frank Olken
Computer Science Research and Development Department
Lawrence Berkeley Laboratory
1 Cyclotron Road, Berkeley, CA 94720

We present a database schema design tool (*SDT*) developed at Lawrence Berkeley Laboratory. The purpose of *SDT* is to provide a powerful and easy to use design interface for biologists, and to increase the productivity of the database design process. This entails insulating the schema designer from the underlying database management system (DBMS).

For the schema design interface we have chosen a version of the *Data-Flow* (DF) model for describing processes and process correlations, and a version of the *Extended Entity-Relationship* (EER) model for the specification of the static structure of information systems. The EER model we use includes, in addition to the basic construct of object (entity and relationship), both generalization and full aggregation abstraction capabilities. We have developed an integrated DF/EER schema design methodology. Following this methodology, DF specifications are represented by EER constructs, so that design of an information system results in an EER schema that captures both the structural and functional characteristics of the modeled system. Once an EER schema is specified, *SDT* is employed in order to generate the corresponding DBMS schema.

SDT consists of two main modules, *SDT_R* and *SDT_{DM}*. The first module, *SDT_R*, takes EER schemas as input and generates abstract relational schemas. *SDT_R* consists of three parts: the canonical mapping of EER schemas into normalized relational schemas; the assignment of names to relational attributes; and merging relations. The canonical mapping generates relational schemas, including key and referential integrity constraints. The high normal form (BCNF) of this schema ensures efficient update performance by the DBMS. Name assignment can be customized in order to meet the needs of the user (e.g. short names, minimum number of attributes, etc.). Finally,

merging of relations reduces the number of relations, thus improving query performance.

The second module, *SDT_{DM}*, takes abstract relational schemas as input and generates relational DBMS (e.g. SYBASE, DB2, INGRES) schemas. For a DBMS that supports the specification of *triggers*, such as SYBASE, the main part of *SDT_{DM}* consists of generating the appropriate *insert*, *delete*, and *update* triggers corresponding to the referential integrities associated with the abstract relational schema.

The research related to *SDT* is presented in [1] and [2]. The DF/EER design methodology is described in [3] and an example application is presented in [4]. The logical algorithms of *SDT* and their implementation are described in [5]; *SDT* was implemented using C, LEX, and YACC, on Sun 3 under Sun Unix OS 4.0.3.

References

- [1] V.M. Markowitz and A. Shoshani, "On the Correctness of Representing Extended Entity-Relationship Structures in the Relational Model", Proc. of 1989 SIGMOD Conference, June 1989.
- [2] V.M. Markowitz and A. Shoshani, "Name Assignment Techniques for Relational Schemas Representing Extended Entity-Relationship Structures", Proc. of 8th International Conference on Entity-Relationship Approach, Toronto, 1989.
- [3] V.M. Markowitz, "Representing Processes in the Extended Entity-Relationship Model", to appear in the Proc. of 6th International Conference on Data Engineering, February 1990.
- [4] V.M. Markowitz and F. Olken, "An Extended Entity-Relationship Schema for a Molecular Biology Laboratory Information Management System", Technical Report LBL-27042, May 1989.
- [5] V.M. Markowitz and W. Fang, "*SDT* Programmer's Manual", Technical Report LBL-27843, November 1989.

* This work was supported by the Office of Health and Environmental Research Program of the Office of Energy Research, U.S. Department of Energy under Contract DE-AC03-76SF00098.

Data Thesaurus for Physical Mapping

John L. McCarthy*

*Information and Computing Sciences Division
Lawrence Berkeley Laboratory*

Physical mapping of human chromosomes must integrate various types of data from different sources. One barrier to integration is that many pertinent biological entities (e.g., cell lines, DNA segments, genes, probes, restriction enzymes, etc.) are known by different names, abbreviations, and local laboratory codes. Multiplicity of names for the same entity (not to mention use of the same name for different entities) is a significant problem for both people and computer programs. Although a small, well-disciplined laboratory can address such problems by requiring that everyone use a single, standard local name for each entity, such an approach is unrealistic for a large laboratory that must interact with many external sources of information.

As an alternative to enforcing a single local naming standard, LBL's Human Genome Center plans to adapt a software mechanism called the data thesaurus to support systematic maintenance and automatic translation of synonyms and related information for major types of biological entities used in its laboratories. LBL's initial prototype data thesaurus will focus on entities pertaining to Chromosome 21, including genes, probes, restriction enzymes, and cell lines.

As demonstrated by its success in conjunction with another LBL scientific data project, the data thesaurus can serve a variety of purposes for both data administrators and end users. For data administrators, the data thesaurus will be a tool for maintaining, documenting and updating various controlled vocabularies. Names of genes, probes, cell lines, and so

on will be registered as either primary terms or synonyms in the thesaurus before they can be used in other LBL databases. The thesaurus can then be used to provide lists of allowable values for such controlled entities and/or to validate relevant fields for either batch or interactive data entry.

In addition to automatic translation of synonyms, the thesaurus paradigm also can support automatic "explosion" of hierarchical groups or classes of entities. For example, if restriction enzymes are classed by both cutting frequency and vendor, people and programs could automatically access sets of enzyme names belonging to larger classes such as "infrequent cutters" or "Merck."

The data thesaurus requires standard database capabilities for access control, data integrity constraints, and so on, with special emphasis on retrieval of nested and repeating text data structures. Since other computer programs, as well as human users, will depend on it for name-based information, the thesaurus must be easily accessible to both via standard interfaces over local and wide-area networks. This combination of requirements may be difficult to achieve with a commercial relational data management system. We are currently trying to identify an appropriate data management system that we can use for the data thesaurus in conjunction with other Human Genome Project software.

*Bldg 50B - 3238, LBL, Berkeley 94720; Internet
Email: JLMcCarthy@lbl.gov

Neural Net Applications to DNA Sequence Analysis *

McMullen, D. F., Lim, H.A., and Burnette, D.E.

Supercomputer Computations Research Institute

Florida State University

Tallahassee, FL 32306-4052

The main objective of this project is to develop neural network ("Connectionist") algorithms for performing several common operations in the analysis of DNA sequence data. Primarily these operations involve the comparison of a sequence with the contents of a data base or the alignment of a number of related DNA fragments. Current, nonconnectionist methods employ heuristic rules or a dynamic programming algorithm to define the degree of similarity desired. For long sequences or data with noise (point mutations or errors, or longer insertions or deletions) similarity searches and alignment can be performed more efficiently using a content-addressable memory scheme than with conventional methods. In order to test the efficacy of connectionist algorithms to the analysis of DNA sequence data, a simulated "chromosome" was constructed by overlaying a random sequence of bases with known sequence data associated with hc14 and regularly distributed features such as inverted palindromes and regions of known base concentration. The simulated chromosome ("simugen") is then "cleaved" by searching for restriction enzyme sites using a three layer back propagation network, and the fragments used in subsequent tests of alignment and data base search algorithms.

* Funded by DOE Energy Research through SCRI

Mapping Algorithms for the Probed Partial Digestion Problem

Abstract

Dalit Naor

Computer Science Division, U.C. Davis

and

ICS Division, Lawrence Berkeley Laboratory

The Probed Partial Digestion method partially digests the DNA with a restriction enzyme. A probe, known to be located between two RE cutting sites, is then hybridized to the partially digested DNA, and the sizes of fragments which the probe hybridizes to are measured. The objective is, then, to reconstruct the linear order of the RE cutting sites from the set of measured lengths.

A Backtracking algorithm, which runs in $O(2^{n-1})$ worst case, where n is the number of cutting sites, is given. The algorithm finds all possible orderings that are consistent with the data. It can be modified to handle inaccurate data. If the data set contains only the lengths but not their multiplicities (that is, band intensities are not taken into account) then the algorithm runs in $O(k^n 2^n)$, where k is the degree of multiplicity, which is believed to be small. A more general version of the problem that uses multiple probes is considered. Two special cases for which simplified solutions exist are pointed out.

The Backtracking algorithm has been implemented for the case where the data set is complete, and preliminary tests have shown that the stated worst case running time is over pessimistic.

Finally, we look at various deconvolution algebraic methods that were suggested in the literature to solve the Partial Digestion problem (when no probes are used). These methods are known to run in polynomial time and are based on algebraic algorithms for factoring polynomials; however, they are *only* applicable when the data is complete and accurate. We point out the difficulties in applying these methods to Probed Partial Digestion.

ROBUST METHODS FOR SIGNAL EXTRACTION AND CALIBRATION IN RESTRICTION FINGERPRINTS

David O. Nelson, Tom Slezak, Elbert W. Branscomb, and Anthony V. Carrano. Biomedical Sciences Division L-452, Lawrence Livermore National Laboratory, Livermore, CA 94550.

Analyzing restriction fingerprints for Chromosome 19 presents several difficulties not normally encountered in data generated, for example, for DNA sequencing. For instance, the data acquired is not in the typical "picket-fence" shape, but rather consists of a superposition of a random number of possibly overlapping peaks of varying sizes. In addition, the signal is corrupted by at least three distinct sources of noise: random noise that is uncorrelated from channel to channel, cross-talk from the other samples loaded in the same lane, but tagged with different colored dyes, and correlated noise corresponding to situations such as imperfections in the gel. We are developing robust, reliable methods for signal extraction and analysis in this complex environment.

The three major signal processing tasks consist of noise suppression, peak detection, and fragment size determination. We suppress the random, uncorrelated noise using a standard robust smoother ("4253H,twice"). We model the color bleed-through as a matrix equation $Ax = b$, where b is the observed data, A is an empirically-determined bleed-through transfer function, and x is the unknown signal we want. For each b , we color-correct by using Hanson's constrained least squares routine "SBOLS"¹ to minimize $\|b - Ax\|$, subject to $x \geq 0$. Enforcing non-negativity in the solution serves to ensure a physically meaningful result as well as to further suppress the noise in the signal. We have just begun to examine ways to suppress correlated "noise" due to gel imperfections and the like.

We are exploring two different, complementary ways to detect peaks. One way, based on smoothing splines, can quickly find peaks in parts of the signal where the structure is not too complex. The other approach, based on treating the peak detection problem as one of *deconvolution*, is about two orders of magnitude more computationally intensive, but can find peaks that simpler, more non-parametric methods cannot. In this model, we assume the given signal is a noisy convolution of a set of "impulse functions" (representing the peaks) by a blurring kernel (representing the diffusion process which occurs simultaneously with migration through the gel). If so, we can recover approximations to the impulse functions by deconvolving the signal with a representation of the blurring kernel. We are using Zhuang's basic approach to Maximum Entropy deconvolution.² However, instead of using his somewhat heuristic algorithm, we are solving his system of Differential-Algebraic Equations directly, using a state-of-the-art DAE solver called DASSL, developed at LLNL by L. Petzold.

In addition to the above noise suppression and data extraction tasks, we also must calibrate the peak locations to a standard. We determine fragment sizes from lane positions in two steps. First, we use a dynamic-programming algorithm which matches lengths from a known standard with peak positions corresponding to that standard (one standard is run in each lane). Finally, we interpolate intervening distances using Fritsch's monotone spline package.^{3 4}

¹RJ Hanson (1982). Linear least squares with bounds and linear constraints. SNLA Report SAND82-1517.

²Zhuang et al (1987). IEEE Trans. Acoustics, Speech, and Signal Processing. ASSP-35:2, 208-218.

³FN Fritsch, RE Carlson (1980). Siam J.Numer.Anal. 17:2, 238-246.

⁴This work performed under the auspices of the U.S. Department of Energy by the Lawrence Livermore National Laboratory under Contract No. W-7405-ENG-48.

The Laboratory Notebook: A Relational Database for the Management of Physical Mapping Data

Debra Nelson, Carmella M. Rodriguez, and Thomas G. Marr
Los Alamos National Laboratory

We have designed and implemented a relational database to manage the data produced by the physical mapping effort for human chromosome 16 at Los Alamos National Laboratory. The database design is sufficiently general to be useful, with little modification, for most mapping strategies. The "Laboratory Notebook" database is meant to be an electronic version of the ubiquitous lab notebooks found in all molecular biology experimental facilities and was designed not only for management of data resulting from experiments, but also to manage information associated with materials, methods, and procedures. The database resides in Sybase, a relational database management system, on a Sun 4 computer operating under UNIX.

We have developed tools to support the flow of data from the laboratory into the database, and once in the database into various forms suitable for analysis and presentation. Electrophoresis gel image information (whole band report) including DNA fragment sizes are transferred directly to the database from the BioImage Visage 110, an image-processing workstation, where the photographs of stained gels are digitized and analyzed. Fingerprint annotation is added to the DNA fragments through a form-based user interface developed for data entry, editing, and retrieval. Reports from these data are created for direct input to contig construction programs (e.g., programs which estimate the probability of pairwise overlap).

Because the users of the "Laboratory Notebook" database are in most cases molecular biologists and not computer scientists, the user interface was designed to present a conceptual view of the data without burdening the users with the need to understand the structure of the database and details of data storage. The user interface was developed using the Sybase application APT-forms and requires very little training to use.

We are in the process of extending the implementation of the database and enhancing the user interface to allow on-line access to other data entities. We will continue to develop tools to facilitate automatic transfer and integration of our local data as well as those data being generated by our collaborators.

SIZE AND DNA BASE COMPOSITION ANALYSIS OF DNA FRAGMENTS.
D. Peters* and J. Gray. Lawrence Livermore National Laboratory,
Livermore, CA

A dual laser, capillary electrophoresis apparatus has been assembled and used to classify DNA restriction fragments according to size and DNA base composition. In this system, DNA fragments stained with Hoechst 33258 (binds preferentially to AT rich DNA) and chromomycin A3 (binds preferentially to GC-rich DNA) are loaded electrophoretically into a 30 cm quartz capillary (50 μ m id) filled with 1% agarose and separated electrophoretically by applying a 1KV potential across the capillary. The fragments pass through laser beams that are focussed through the capillary near its end. The lasers are adjusted to the UV and 442 nm to excite Ho and CA3, respectively. The fragments are classified according to their time of arrival at the laser beams (i.e., according to molecular weight) and according to the ratio of their HO and CA3 fluorescence intensities (i.e., according to DNA base composition).

Work performed under the auspices of the U.S. Department of Energy by the Lawrence Livermore National Laboratory under contract number W-7405-ENG-48.

SIMULATION OF PHYSICAL MAPPING

Karl Sirotkin and Eric Fairfield

Los Alamos National Laboratory

We have created a set of modular programs to evaluate strategies for physical mapping of the human genome. A person using these programs can: create a test genomic segment, generate clones from this segment, extract fingerprint data from each clone, test various strategies for determining pairwise overlap between clones, reassemble the genome from these overlaps, display the resulting contigs, and evaluate the success of the reassembly.

The programs have been structured to evaluate many mapping strategies and to assemble contigs from real data as it becomes available. In anticipation of other users, both program and data files are readable by the user and the same symbolic parameter names are used throughout the data and program files. In addition, we have designed the program structure so that it is easy to tailor the installation for individual users.

These programs have been implemented in two phases. The first phase only used exact fingerprint data, while in the second phase there were controlled levels of error in the fingerprint data.

By using exact data, four different strategies reassembled similar percentages of a genome segment into contigs; although, the exact contigs were different. By combining information from all of these methods, the coverage of the genome by contigs increased. With errors in the fragment lengths, it is not yet clear whether particular strategies for reassembly of fragments into contigs make better use of the data. Small changes in the experimental errors seem to have large effects on contig generation.

The current set of programs have been delivered to a beta test site and will be available for testing at the conference.

HUMPTY: An Algorithm for Fully-Automated Contig Assembly

Tom Slezak, Elbert W. Branscomb, and Anthony V. Carrano. Biomedical Sciences Division, Lawrence Livermore National Laboratory, Livermore, CA 94550

Developing a physical map of a human chromosome in the form of an ordered library of cosmid clones involves assembling many thousands of cosmids. Manual contig assembly on this scale is neither practical nor necessary. We have developed an algorithm that automatically assembles contigs and determines a near-minimal spanning path with highly-confident determination of contig ends. Input to this algorithm is a sorted list of the LOD scores (odds ratio in favor of overlap) of all pair-wise combinations of DNA fragments. An optimized data structure is coupled with a depth-first greedy search strategy, yielding correct reconstruction of 8,000 simulated chromosome 19 cosmids in under 15 minutes on a Sun-4/260. The algorithm is coded in C for running under Unix. EcoR1 digests of selected chromosome 14 contigs have verified both contig membership and spanning path determination. Output from the HUMPTY algorithm can be viewed and manipulated graphically using the contig browser described in a separate poster by Mark Wagner of LLNL.

(This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract No. W-7405-ENG-48.)

Special Requirements: Poster Presentation, no equipment needed

A Graphical Contig Browser Tool for DNA Mapping

Mark C. Wagner, Thomas R. Slezak, Elbert W. Branscomb, and Anthony V. Carrano.
Biomedical Sciences Division, Lawrence Livermore National Laboratory, Livermore, CA 94550

We have developed a highly-automated system for the visualization of mapping information from cosmid clones. Input to this Browser tool is derived from HUMPTY, a program which takes the overlap probability for all pairwise combinations of cosmids based on fingerprint data and forms spanned sets of overlapping clones (contigs). This information is too complex and voluminous to be readily understood on paper. We developed a graphical contig browser to assist in comprehending and manipulating this data. Providing a color-coded visual representation of the data permits a better understanding of the relationships between the individual cosmids that comprise each contig. Our implementation is based on LLNL-developed graphics libraries that run on top of the industry-standard X11 graphics system (Sun and Stellar), and on Silicon Graphics Iris workstations. All coding is in C for use under Unix. Our prototype version is in daily use helping us to analyze the over 2,500 human chromosome cosmids fingerprinted and mapped to date. Work in progress will tie the contig browser to a Sybase relational database system, add the ability to view the raw cosmid fingerprint data from several cosmids simultaneously, and allow the ability to view and manipulate 2 contigs at once to allow gap closure from non-fingerprint data. This tool is our model for our planned graphical chromosome database browser that would allow ready access to all data generated on the Human Genome project.

(This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract No. W-7405-ENG-48.)

Special Requirements: Poster Presentation, software can be demonstrated live if a color Sun-4 system is available.

Abstract for poster presentation to DOE Contractor's Workshop, Nov 3 - 4, 1989 in Santa Fe, NM, by Mr. John West, Principal Investigator, BioAutomation, Inc.

CUSTOM CHIP TECHNOLOGY FOR IMAGE SCANNING

Numerous applications exist in modern biology for image scanning. The push for automation in the Human Genome Initiative will expand these applications. We report on experimental work done in our laboratory in 1989, investigating the use of application specific integrated circuit (ASIC) technology to implement circuitry for image scanning. In our work we have used a 1.2 micron CMOS gate array. With isolated flip flop toggle rate capability of 100 MHz, we find that the balance of routing and logic resources on the chip, combined with the capacitance induced routing delays, allows implementations to date with up to 6 MHz operating frequencies. This is a frequency which is easy to interface with other components of a system. The parallelism possible in such a hardware implementation provides for high performance at this frequency. With moderate volume, the costs of such an implementation can also be quite attractive.

The Human Genome Information Resource

Scott L. Williams,¹ Rena Whiteson, David C. Torney, Robert D. Sutherland, Karen R. Schenk, Alice Sandstrom-Bertini,¹ Carmella M. Rodriguez, Robert M. Pecherer, Debra Nelson, Frances A. Martinez, Thomas G. Marr, James H. Jett,² C. Edgar Hildebrand,² Michael J. Cinkosky, and Christian Burks.³ Theoretical Biology and Biophysics Group; T-10, MS K710; Los Alamos National Laboratory; Los Alamos, NM 87545; U.S.A.

¹Computer Research Laboratory; New Mexico State University; Las Cruces, NM.

²Life Sciences Division; Los Alamos National Laboratory; Los Alamos, NM 87545.

³Corresponding contact: telephone, 505-667-6683; e-mail, cb%intron@lanl.gov.

The Human Genome Information Resource is focused on the need for better information management and analysis tools for physical mapping data (e.g., the data currently being generated in the context of the effort to generate a complete physical map of human chromosome 16 [Hildebrand et al. (1989) these proceedings]), and reflects a long-term interest in extending research to other, related data sets such as nucleotide sequences and genetic maps.

Because of the desirability of having "real" experimental data to examine and manipulate, and because of the close ties with the experimental group in LS-Division at LANL, the project has focussed initially on developing strategies and tools for supporting the flow of data from DNA gels into computers and, once in the computer, into various forms for analysis and presentation. This flow of data begins with the digitization and processing of electrophoretic gel images. Data corresponding to the clones analyzed on the electrophoretic gels are then passed into a computerized laboratory notebook [Nelson et al. (1989) these proceedings] based on a relational database management system. These data are made available for subsequent analysis of the clone fingerprints, leading to the development of contig maps [Torney (1989) these proceedings] and -- eventually -- comparison to other, related data sets that will allow for higher-order assemblies and ordering of contigs.

As one gets to the end of this flow path, the need for more sophisticated data management, analysis, and interface tools becomes evident. The thrust of the HGIR project will shift to the development of these tools, and their use to facilitate the cross-linking among multiple levels of physical mapping data, as well as between physical maps and other, related data sets (e.g., sequences and genetic maps). The design work on data structures for physical map and sequence data we are doing is being undertaken with this emphasis (and the future extension to yet unrecognized data structures) in mind. Finally, we plan to design an on-line system and set of interfaces allowing for the provision of these data and tools to a much more broadly-defined user community than the current in-house activity.

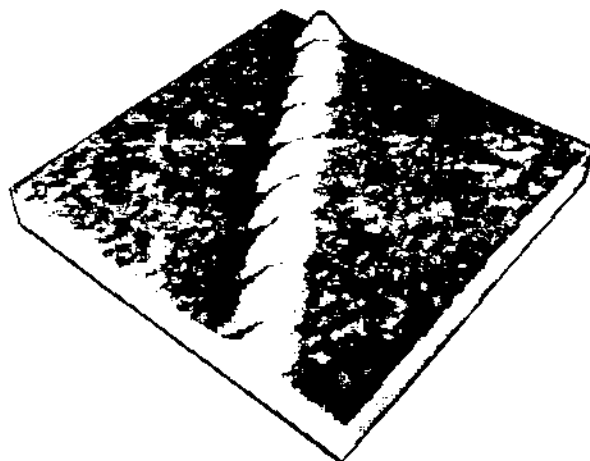
Scanning Tunneling Microscopy of Macromolecules

D. P. Allison, J. R. Thompson, K. B. Jacobson, R. J. Warmack,
and T. L. Ferrell

Oak Ridge National Laboratory
Oak Ridge, Tennessee

Beginning in 1982 when the first images were reported, the development of the scanning tunneling microscope (STM) has revolutionized microscopy in the 1980's. Although the primary applications are for research on metal and semiconductor surfaces, success in imaging a number of biological samples, mounted on a variety of conductive surfaces, have established STM as a valuable emerging technology in biological research. In 1987 we began our biological STM studies using tobacco mosaic virus (TMV), an easily identifiable rod shaped virus, to test the feasibility of using STM on naked biological samples. The STM operates by scanning a sharp tip a few atomic diameters away from a conductive surface and measuring changes in current or voltage as a function of changes in surface topography. We deposited TMV by spraying an aqueous solution of the virus on to evaporated or sputter-coated palladium-gold (Pd/Au) films supported on flat mica surfaces. Although TMV could be clearly identified its observed width was typically 70-100 nm instead of the known diameter of the virus, 18 nm. On evaporated Pd/Au substrates, tip traces revealed structures elevated above the substrate surface suggesting that the virus became coated with Pd/Au allowing normal conduction to occur. On sputter-coated substrates tip traces revealed depressed substructures. This is presumed due to the poor electrical conductivity of TMV causing the true differential tip motion to be recorded erroneously. Although not routinely obtainable, we have observed exceptional images of the virus revealing protein subunit structures separated by only 2.4 nm. We propose these images are the product of an exceptional tip, combined with an unusual conductivity of the virus.

In May of this year we obtained our first images of DNA. A circular plasmid of p BR 322 containing two genes for antibiotic resistance (tet^R and amp^R) was mounted on a graphite surface and imaged in air. A portion of one of these images is shown, clearly demonstrating the right-handed helix with an axial repeat spacing of 4.7 nm somewhat different than the 3.4 nm repeat expected of the B-form of DNA. With the routinely available resolution demonstrated in these images, the STM is currently capable of detecting conformational changes and binding of substances, such as repressor proteins, to DNA molecules. In the future, with improvements in methods for routinely attaching both single and double-stranded DNA to conducting substrate surfaces, we propose to use STM to identify nucleotide sequences in intact DNA molecules. This will be accomplished by hybridizing oligonucleotides of known sequences, suitably labeled for recognition by STM, to intact DNA molecules.



* Sponsored by the Office of Health and Environmental Research, U.S. Department of Energy, under contract DE-AC05-84OR21400 with Martin Marietta Energy Systems, Inc.

NOVEL DNA POLYMORPHIC SYSTEM: VARIABLE POLY A TRACT 3' TO ALU I REPETITIVE ELEMENTS.

S.E. Antonarakis, E.P. Economou, and A.W. Bergen.

Genetics Unit, Dept. of Pediatrics, The Johns Hopkins University School of Medicine, Baltimore, MD 21205.

DNA polymorphisms are extremely useful in the mapping of the human genome and the search for disease gene loci. Several categories of DNA polymorphisms include single nucleotide substitutions, variable number of tandem repeats, presence or absence of L1, Alu I elements or pseudogenes and variable number of dinucleotide repeats (CA)_n or (CT)_n.

We describe here a novel class of DNA polymorphisms: Variable number of adenylic acid residues (As) at the end of Alu I repetitive elements. To test the hypothesis that the poly A tract of Alu I sequences is polymorphic, three areas that contain the 3' end of an Alu I element have been selected for polymerase chain reaction (PCR) amplification. The first about 600 nt 5' to the β globin gene (β Alu) the second about 1100 nt 5' to the transcription initiation site of adenosine deaminase gene (ADA-Alu) and the third in IVS 1 of the factor VIII gene (F8 Alu). The PCR oligonucleotides were chosen from the "right arm" of the Alu I repetitive element and from a single copy sequence following the poly A tract and were used in a 1:10 concentration ratio respectively. The single copy oligonucleotides were end labelled with ³²P and the PCR product was electrophoresed in a 6% acrylamide sequencing gel. Due to a different length of the poly A tract in different alleles several polymorphic patterns were observed. Mendelian inheritance was demonstrated in CEPH families and nucleotide sequencing confirmed that the multiple allelism was due to the different number of As.

The β Alu showed 32% heterozygosity for different (multiple) alleles in the grandparents and parents of the CEPH families. The ADA Alu showed 6 different alleles with frequencies of 29%, 10%, 28%, 27%, 7% and 1% in 267 independent chromosomes examined in the CEPH families. The observed heterozygosity for this polymorphism was 75%. The F8 Alu showed no polymorphisms in 20 unrelated females from the CEPH families.

Since more than 10⁵ Alu I sequences exist in the human genome their variable poly A tract (Alu-VpA) may prove to be one of the most abundant and useful polymorphic systems.

Imaging of DNA Molecules Deposited on Graphite. R. Balhorn*, M. Allen*, B. Tensch**, J.A. Mazrimas*, M. Balooch†, and W. Siekhaust†. *Biomedical Sciences Division, **Department of Applied Science, and †Chemistry and Materials Science, Lawrence Livermore National Laboratory, Livermore, CA 94550.

The conditions required for imaging DNA near atomic resolution with the scanning tunneling microscope are being examined as the first step in our effort to devise an alternate, electronic method for sequencing DNA. Biotinylated lambda phage DNA and a defined length synthetic duplex DNA have been deposited on highly oriented pyrolytic graphite (HOPG) and imaged in air by scanning tunneling microscopy. The lambda phage DNA was tagged with streptavidin coated 20nm gold spheres and the spheres located in the initial 4000Å X 4000Å scans. High resolution images of three attached strands of DNA show the variability in level of detail that can be observed. Helical coiling is detected in several regions, but it is obscured in others. The end of one molecule is unwound and the last 100-120 base-pairs have separated and are visible as single strands. Reproducible images of a 47 base-pair (bp) DNA sequence have also been obtained. Each molecule exhibits a similar periodic structure and length compatible with expected values. Interesting structural details are revealed, including the left handedness of the DNA helix in the G-C rich region, the presence of short, single-stranded "sticky" ends, and the major and minor grooves. This work was funded by the U.S. D.O.E. by the Lawrence Livermore National Laboratory under Contract W-7405-ENG-48.

TRANSPOSON Tn5 FACILITATED DNA SEQUENCING

D. E. Berg, S. H. Phadnis, T. Tomcsanyi, H.V. Huang and C. M. Berg*

Depts. Molec. Micro., Washington Univ. Med. School, St. Louis, MO. 63110, and

*Cell & Molec. Biol., University of Conn., Storrs, CT. 06269

Bacterial transposons are being developed to place unique sites for DNA sequencing primers and multiplex probes throughout cloned DNAs, thereby minimizing the need for random DNA subcloning or repeated syntheses of new oligonucleotide primers. These experiments involve derivatives of Tn5, a transposon that inserts efficiently and quasi-randomly in diverse target DNAs.

Tn5supF. We constructed this 264 bp mini-transposon for insertion mutagenesis and sequencing of DNAs cloned in phage λ . Tn5supF is marked with the suppressor tRNA gene, *supF*. Insertions into amber mutant λ are selected by plaque formation on wild type *E. coli*. Insertions in non-amber λ are selected similarly using the *dnaB*-amber bacterial strain DK21 from David Kurmit, a selection that exploits the need for DnaB protein during λ DNA replication.

Saturation mutational analyses of entire genomes will grow in importance as genome sequencing projects near completion. We have begun recombining Tn5supF inserts made in the λ /*E. coli* hybrid phage of Kohara et al. (Cell 50:495-508) into the *E. coli* chromosome. Because these phage are defective in λ repressor synthesis (*cI*) phage λ b221 *cI857* was used as a co-infecting helper to supply repressor and thus permit survival of infected cells. Haploid bacterial recombinants were obtained readily with a *lacZ*::Tn5supF mutation. In contrast, only partial diploids (containing both mutant and wild type alleles) were obtained with *rpmD*::Tn5supF and *rpoA*::Tn5supF insertion mutations because *rpmD* and *rpoA* encode essential proteins. We conclude that Tn5supF-based reverse genetics, entailing first physical mapping and then phenotypic testing of new mutations, is well suited for analysing genes and sites found during DNA sequencing.

Deletion factory. Our recent experiments have shown that intramolecular Tn5 transposition generates deletions that place different regions of the target DNA close to a transposon end, analogous to exonuclease-based in vitro and IS1-based in vivo nested deletion strategies. For these tests we placed a synthetic Tn5 element next to a *sacB* (sucrose sensitivity) gene, so that deletions could be selected by sucrose-resistance. In contrast to results with IS1, the endpoints of deletions made by Tn5 transposition were widely distributed in target DNAs.

Supported by grants GM37138 and DE-FG02-89ER60862.

STRUCTURAL AND TRANSCRIPTIONAL ANALYSIS OF A CLONED HUMAN TELOMERE Jan-Fang Cheng+, Cassandra L. Smith* and Charles R. Cantor+, Departments of +Genetics and Development, *Microbiology, and *Psychiatry, College of Physicians and Surgeons, Columbia University, New York, NY 10032

Isolation of a human telomeric YAC clone, yHT1 (Nucleic Acid Res 17, 6109-6127), allows the characterization of a common DNA structure next to the telomeric TTAGGG repeats. Blot hybridizations using various portions of the yHT1 clone to probe against a panel of somatic hybrids indicate that this common subtelomeric structure spans at least 4 kb in length, and appears in multiple copies on some chromosomes but does not appear on the X chromosome. The complete DNA sequence of yHT1 has been determined. This clone contains an AT-rich region and a CpG island, separated by a human Alu repeat. Cross-hybridizations are weak in rodents when using various portions of the yHT1 clone as probes. However, the CpG island gives strong and distinct cross-hybridizing fragments. The significance of this evolutionarily conserved DNA sequence remains to be determined. Transcription patterns in this subtelomeric region have been examined by Northern blot analysis. Both polyA+ and polyA- transcripts were detected when probing with DNA isolated from the AT-rich region. Only polyA-RNA was detected when probing with DNA isolated from the CpG island.

MINIATURIZATION OF SEQUENCING BY HYBRIDIZATION (SBH):
A NOVEL METHOD FOR GENOME SEQUENCING

Crkvenjakov, R., Drmanac, R., Strezoska Z., Labat I.,
Genetic Engineering Center, PO Box 794, 11000 Belgrade,
Yugoslavia

Human genome sequencing based either on gel electrophoresis, or recently proposed hybridization (Drmanac et al. GENOMICS (1989) 4:114) methods requires automated equipment on macro scale and can not be imagined as a routine procedure. Macro scale is mandated due to the requirements of robotic positioning of samples on predetermined coordinates and polymer separation in gels. However, determination of oligonucleotide contents of DNA which underlies SBH theoretically allows the micro scale processes with micro separated samples altogether comprising a macro scale reaction. It is possible to use the determination which clone/probe is on which random micro position instead of placing clone/probe on predefined macro position or volume. We propose the use of micro discrete particles (DPs) as vehicles for samples/probes. The recognition of specific association of a DP and a clone/probe is achievable by premarking of DPs and/or determining characteristics of clone/probe in situ. The most obvious ways of marking DPs are shape, size or color, or attaching to it a specific combination of known oligonucleotides. We offer two possibilities for human genome sequencing (Drmanac et al., manuscript in preparation). For direct SBH 1×10^7 clones coming from 10 separate genome parts are bound to 1×10^6 different DPs in as many macro reactions (or eventually in a single macro reaction). 1×10^3 monolayers containing more than 1×10^7 individual previously mixed DPs are each after DP identification hybridized with groups of 100 differently labeled octamers. To this end we have developed conditions for reliable short oligonucleotide hybridizations. For inverse SBH $1 \times 10^{7-9}$ different DPs are prepared each carrying a unique 12-15mer and unique combination of 20 out of 40 marking oligos. No more than 5000 separate macro reactions are needed for their preparation. After 40 hybridizations with marker oligos to find association between 12-15mers and DPs in a monolayer, in 1-100 hybridizations with fragmented, end labeled human DNA data for sequencing are generated. The monolayer area covers at most 100 microscope slides. The data collection for both approaches needs automated image analysis giving speed of data bits acquisition of 1×10^6 /s. Finally a substantial computing has a major role to keep track of information and generate sequence (see accompanying abstract). The described miniaturization concept and ensuing savings make human genome sequencing immediately feasible in a laboratory pending technological development.

DNA fragment fingerprinting and/or sequence determination by Fourier analysis of coherent X-ray scattering

Joe W. Gray¹, James Trebes², James Brase², Daniel Pinkel¹, Thomas Yorkey², and Heinz-Ulrich Weier¹

¹Biomedical Sciences Division, ²Laser Program
Lawrence Livermore National Laboratory, Livermore, CA 94550

This report describes the theoretical basis for rapid characterization of the distribution of labels (e.g., high Z scatterers such as Au, I or Br) along DNA molecules. In this approach, the DNA fragment to be characterized is amplified to $\sim 10^7$ copies (e.g., by in vitro DNA amplification), labeled (e.g., by hybridization to labeled oligonucleotides for fingerprinting or by incorporation of labeled bases during in vitro DNA amplification for DNA sequence analysis) and arranged as an array of straight (but not necessarily parallel) DNA molecules. The distribution of labels in the ensemble of linear DNA fragments is determined by Fourier analysis of the scattering pattern formed during irradiation with coherent X-rays. The Fourier analysis is simplified by labeling one end of the test DNA fragments with a distinct scatterer (e.g., by hybridization with a gold microsphere labeled oligonucleotide). Preliminary analyses suggest that sufficient scattering for DNA sequence analysis can be obtained from $\sim 10^7$ I-labeled DNA fragments with an exposure of a few minutes to existing 8 KeV x-ray sources synchrotrons, laser plasmas, electron beam sources. The key to success in this process is the creation of arrays of *straight* DNA molecules. Constrained electrophoresis in low ionic strength buffer is being investigated as a way of accomplishing this. Initial experiments are directed toward analysis of the locations of iodine labeled thymines in the DNA sequence:

CCC CCC CCC CCC CCC TAA AAA AAA AAT AAA AAT.

This work performed under the auspices of the U.S. Department of Energy by the Lawrence Livermore National Laboratory under Contract number W-7405-ENG-48.

High Resolution DNA Mapping by STEM

James F. Hainfeld
Biology Department
Brookhaven National Laboratory
Upton, NY 11973

A new method has been developed to map DNA and RNA such that specific sequences might be visualized in the electron microscope to within 3 to 5 base pairs.

Preliminary results have been obtained using the following test system. A 622 base pair (bp) sequence from pBR322 was excised with restriction enzymes and purified. Next, a 128 bp T7 piece was inserted at position 276 (giving 720 bp total). Equal quantities of the 622 bp and 720 bp fragments were denatured and renatured. This resulted in 50% formation of heteroduplexes (one 622 strand paired with a 720 strand) leaving the extra bases as a single stranded loop. A 26-mer oligonucleotide was synthesized that was complementary to a region of the single stranded insert. A chemical modification added a sulfhydryl at the 5' end of the oligo and the undecagold cluster (with a 0.8 nm diameter gold core) was covalently attached to it. Next, the oligo and heteroduplexes were mixed under renaturing conditions and examined in the Scanning Transmission Electron Microscope (STEM). Gold clusters were observed at the expected positions.

The gold cluster is about 10 \AA from the base it labels (3 base pairs) and the accuracy of positioning a base from the end of DNA segments in the STEM is 2 bp, giving a total potential positional accuracy of 3-5 bp. This should prove useful in the physical mapping of genomes.

APPLICATION OF THE SCANNING TUNNELING MICROSCOPE TO STUDIES FOR DNA STRUCTURES

M. Salmeron, M. Bednarski, D.F. Ogletree, T. Wilson

Center for Advanced Materials
Materials and Chemical Sciences Division
Lawrence Berkeley Laboratory
Berkeley, California 94720

Scanning Tunneling Microscopy has great potential as a tool for the study of biological macromolecules, including DNA. The STM can be operated in air or in liquids, and image contrast does not depend on metal shadowing or replication methods. STM images of unshadowed DNA obtained in our laboratory and by other groups have demonstrated sub-nanometer spatial resolution.

Initially it was believed that the STM would have limited application to biology since most interesting materials are non-conductors. Experiment has shown that the STM can image many molecules on conductive substrates that are insulators as bulk materials. Two major problems remain to be solved before the STM can be used as a routine tool for molecular biology.

The first problem is contrast - in the STM contrast depends primarily on shape, so conductive and chemically inert substrates are required that are smooth on the sub-nanometer scale over areas of several microns.

The second major problem is fixation of molecules to the substrate. Tip-surface forces in STM are often sufficient to deform or displace molecules. Methods must be developed both to reduce tip-surface forces and to bond molecules to suitable substrates. We will show that in principle, this can be solved by constructing organic monolayers that possess reactive functional groups in the surface. Examples of such layers on boron doped Silicon substrates will be presented.

HUMAN REPETITIVE DNA SEQUENCES FOR USE AS MARKERS IN MAPPING THE
HUMAN GENOME

C.W. Schmid, E.P. Leeflang, G. Wang

A 480 clone library of repetitive human DNA sequences is being analyzed to generate potential probes for use in mapping the human genome.

The library was screened for known repeated sequences and of the remaining 264 clones, 23 clones have thus far been selected for further study including lambda clone base sequence analysis, copy number, genomic arrangement and homology to rodent DNA.

A Probe-Based Mapping Strategy for DNA Sequencing with Mobile Primers

Linda D. Strausbaugh, Michael T. Bourke, Martin T. Sommer and Claire M. Berg
Department of Molecular and Cell Biology, The University of Connecticut, Storrs, CT 06269.

Currently popular large-scale methods for DNA sequence acquisition require sets of short, often random, DNA fragments adjacent to primer binding sites. An alternative sequencing strategy utilizes mobile transposons whose ends are used as primer binding sites, thus permitting large clones to be sequenced without fragmentation. We demonstrate a novel and efficient probe-based method for the localization and orientation of such transposon-borne primer sites, which requires no prior restriction enzyme mapping or knowledge of the cloned sequence. This approach, which eliminates the inefficiency inherent in totally random sequencing methods, is applicable to mapping insertions of any transposon in plasmids and will be particularly valuable for larger recombinant molecules in vectors such as cosmids and P1.

The transposons gamma delta (Tn 1000) and Tn5 show considerable promise for large scale sequencing. Although not as intensively developed as some other elements, gamma delta (Tn1000) inserts quite randomly, and plasmids containing a single gamma delta insertion can be obtained readily. A 6.7 kb *EcoRI* fragment of *Drosophila melanogaster* DNA cloned in pBR325 was chosen as the target because partial sequence analysis had shown that this fragment contains regions of atypical base composition. In this model system, we have used existing features of wild type gamma delta and this particular recombinant DNA: *EcoRI* cuts gamma delta in its control region and cuts the plasmid at the two plasmid-vector junctions. DNAs from plasmids containing gamma delta insertions were digested with *EcoRI*, and the resulting fragments electrophoresed, transferred, hybridized to radioactive probes, and visualized by autoradiography. Fifty insertions were easily mapped and oriented using one probe specific for an end of gamma delta and a second probe specific for an end of the cloned fragment.

Primers specific for unique subterminal segments at each end of gamma delta were used to prime dideoxy double stranded sequencing. Each transposon yielded at least 200 bp of sequence information from each primer. These results confirm the random nature of gamma delta insertion and demonstrate the effectiveness of probe mapping. Since transposition and resolution functions can be provided in trans, mini-gamma delta derivatives designed especially for probe mapping will be easy to construct.

Transposon-based probe mapping and sequencing bridges the gap between large cloned segments and unordered subclones. The "duplex" sequencing strategy described can be adapted to multiplex sequencing by inserting heterologous probe/primer sites in gamma delta derivatives.

Characterization and use of linking libraries in Chromosome 21
restriction map construction. Denan Wang, Akihiko.Saito,
Jose P.Abad, William M.Michels, Cassandra L.Smith and
Charles R.Cantor. Lawrence Berkeley Laboratory, University of
California, Berkeley, CA 94720.

A top down approach to mapping and ultimately, sequencing
the human genome starts by the construction of a low resolution
restriction map of each chromosome. Effective procedures have
been developed for constructing *Not* I linking libraries starting
from chromosome-specific genomic libraries. Seventeen unique
single copy *Not* I linking clones from human chromosome 21
were identified in two libraries. Their chromosomal origin was
confirmed, and regional location established by using hybrid cell
panels. Hybridization experiments with these probes revealed
pairs of genomic *Not* I fragments and neighboring *Not* I sites.
Additionally, partial digestion as well as cell line polymorphism
strategies were use to see ncighboring fragments. These strategies
construct maps in regions for which linking clones are not identified.

Cloning of Yeast Artificial Chromosomes by electroporation. M. Bell and R. K. Mortimer, Lawrence Berkeley Laboratory and Department of Molecular and Cellular Biology, Division Of Biophysics, University of California, Berkeley, CA 94720. A strong bias against the cloning of larger Yeast Artificial Chromosomes (YACs) by spheroplast-PEG transformation has been observed. Although treatment with polyamines facilitates the cloning of larger YACs, the smaller Yeast Artificial Chromosomes in any given ligation mixture are cloned preferentially. In an attempt to eliminate this problem, we are exploring electroporation of both spheroplasted and intact yeast cells. Since electroporation of *S. cerevisiae* is a relatively new field, we are currently establishing basic electroporation transformation procedures with control plasmids.

References

- Burgers, P. M. J. and K. J. Percival (1987) *Analytical Biochemistry* 163: 391-397.
Delorme, E. (1989) *Applied and Environmental Microbiology* 55: 2242-2246.
Mc Cormick M. K., Shero, J. H., Antonarakis, S. E and P. H. Hieter (1989) *Technique*, in press.

Poster Presentation

DETECTION OF DNA SEQUENCES WITH CHEMILUMINESCENCE

Irena Bronstein, Tropix, Inc., 47 Wiggins Ave., Bedford, MA
01730

Non-isotopic detection of DNA sequences has most commonly been achieved with fluorophores and, in some cases, with alkaline phosphatase as the label which can be detected with a colorimetric substrate BCIP/NBT. We have developed a new substrate for alkaline phosphatase 3-(2'-spiro-adamantane)-4-methoxy-4-(3"-phosphoryloxy) phenyl-1,2-dioxetane (AMPPD™), which chemiluminesces upon enzymatic dephosphorylation. This substrate when coupled with suitably engineered oligonucleotide probes provides ultrasensitive detection of DNA in Southern blots, and rapid detection of DNA sequences in the genomic sequencing protocols. DNA probes which were labeled with biotin and incubated with streptavidin-alkaline phosphatase and AMPPD allowed the detection of subpicogram quantities of target DNA using Southern analysis. Chemiluminescent detection of DNA sequences using the genomic sequencing procedure of Church and Gilbert revealed images of sequence ladders on x-ray film with exposure times of less than 30 minutes, as compared to 40 hours for a similar exposure with a ³²P labeled oligomer. The demonstrated shorter exposure times would permit more frequent serial reprobings of DNA sequences.

The Separation of Non-Denatured DNA Fragments
with Electrophoresis Gels that are Easily Formed
by Crosslinking an Acrylamide-Rich Copolymer

A novel way of forming electrophoresis gels and separations that can be achieved with these gels are described. The gel-forming procedure is straightforward, begins with a stock solution of copolymer and does not involve toxic chemicals, oxygen exclusion, free radical polymerization or heating. Gels of polymer content greater than 2% are readily obtained and these provide media in which excellent separation and resolution of non-denatured DNA fragments up to 5,000 bp can be achieved.

Authors: Kenneth G. Christy, Jr., David B. LaTart,
Hans W. Osterhoudt and Ignazio S.
Ponticello. Life Sciences Research
Laboratories, Eastman Kodak Company,
Rochester, New York 14650-2122.

HWO:jrs/#081C

10/13/89

Instrumentation Development for Molecular Biology

J.B. Davidson
Instrumentation and Controls Division
Oak Ridge National Laboratory

Ultra low light level detection and imaging are techniques with potentially wide application in several areas of molecular biology. Among these are analysis of 2-D protein gels, sequencing gels and mapping blots. In addition, they can provide the basis for 2-D detectors in the important area of 3-D molecular structure determination by neutron and x-ray crystallography and small angle scattering. In these applications, film can be eliminated and images of the fluorescent and radioactive bands and diffraction spots can be accumulated in a digital memory for analysis. The principles can be applied at the microscopic level for neuronography and in situ hybridization studies.

Progress will be reported on:

- Electronic autofluorography developments
- A "lensless" radiation microscope development
- A P.C. based 2-D detection system for neutron and x-ray diffraction
- A novel viewing aid to reading sequencing films, the "Unsmiler" (Demonstration)

Advanced Concepts for Base Sequencing in DNA

J. H. Jett, R. A. Keller, J. C. Martin, E. B. Shera

Life Sciences Division, Los Alamos National Laboratory, Los Alamos, NM 87545

(505) 667-3843, FTS 843-3843

We are addressing the problem of rapidly sequencing the bases in large fragments of DNA. The ideas presented represent the combined effort of a multidisciplinary task force composed of physicists, physical chemists, cellular and molecular biologists, and organic chemists. To reduce mapping requirements, the emphasis is on sequencing methods that are rapid, require little DNA, and are capable of sequencing large fragments. After evaluation of several physical approaches to sequencing, the decision was made to proceed with a modified-flow cytometer approach that employs laser-induced fluorescence to detect individual fluorescent molecules. A large fragment of DNA, approximately 40 kb in length, will be labeled with base-identifying tags and suspended in the flow stream of a flow cytometer capable of single molecule detection. The tagged bases will be sequentially cleaved from the single fragment and identified as the liberated tag passes through the laser beam. We are projecting a sequencing rate of 100 to 1000 bases per s on DNA strands approximately 40 kb in length.

Experimental Apparatus for Pulsed Field Electrophoresis Research

W.F.Kolbe, J. E. Katz, S. E. Lewis and J. M. Jaklevic

Engineering Division and Human Genome Center

LBL etc.etc.

A test apparatus for research and development involving pulsed-field gel-electrophoresis has been constructed. The system includes a 24-node computer-controlled power supply capable of independent programming and sequencing of individual attached electrodes. An experimental gel-box associated with the system is designed to provide flexibility in terms of individual electrode design and array geometries. A closed loop cooling system maintains precise control of buffer temperature with minimum perturbation to the streamline flow across the gel. Computer calculations are used to define the electrode voltage distributions required to generate specific electric field profiles within the gel-region. A repetitive cycle of the calculated voltage distributions can be imposed on the electrode array under software control in order to implement a variety of existing and experimental pulsed-field electrophoresis protocols. A three dimensional precision manipulator allows computer controlled scanning of the gel with interchangeable probes in order to map the spatial distribution of voltage, buffer temperature, pH, etc. Details of the system will be described and preliminary results obtained using a linear variation in the dwell times of homogeneous, switched electric fields will be presented.

CLONING LARGE HUMAN INSERTS

Andrew A. Kumamoto, Ronald Law*, Robert Deans[†], and Philip Youderian. California Institute of Biological Research, La Jolla, CA 92037; *MBI, UCLA, Los Angeles, CA 90024; [†]Department of Microbiology, USC School of Medicine, Los Angeles, CA 90033.

We have devised a protocol for preparing large fragments of human DNA (250-450 kbp) based on a procedure for preparing phage P22 DNA. This procedure is rapid (about three hours), requires a minimum number of manipulations, and, most importantly, does not use phenol or lengthy dialysis steps. 44 kbp phage P22 DNA prepared in this manner has been shown to have 4 to 5-fold greater biological activity than phenol-extracted phage DNA.

We are developing a novel vector that will permit the cloning of these large (250-450 kbp) segments of human DNA. This vector will carry large human inserts as segments of *E. coli* F plasmids. Several human (and, for that matter, *E. coli*) sequences cannot be maintained on high copy number vectors, including traditional cosmid cloning vectors, and only can be recovered intact as single copy clones. Therefore, we are testing three different single-copy origins derived from F (oriV, oriVII, and oriS) as possible vector origins. To facilitate the preparation of large amounts of insert DNA, these vectors will be modelled after F'::Mud-P22 elements that may be amplified to very high copy number (e.g., 500 copies/cell of a 250 kbp plasmid) after a terminal induction event.

MANIPULATION OF SINGLE DNA MOLECULES IN A MICROSCOPE STAGE. M. F. Maestri‡, M. Miller‡, S. Goolsby‡, W. Johnston†, C. Bustamante*, Lawrence Berkeley Laboratory, Berkeley CA. and *Chemistry Dept. University of New Mexico. Supported by NIH Grant #AI08427 and *GM32543 and by the † Director, Office of Energy Research, Scientific Computing Staff, U.S. Dept. of Energy, contract DE-AC03-76SF00098, and by the ‡ † LBL Human Genome Center.

The aim of the project is to develop techniques that will permit the visualization and manipulation of a selected single DNA molecule for the purposes of mechanical or chemical alteration of the molecule. To this goal we have designed a microelectrode chamber with a spatially distributed electrode network of microscopic dimensions to be used in l with a fluorescent imaging. The electrode network consist of 24 electrodes of dimensions of 10 micrometers in thickness separated from each other by 10 micrometers. With the development of the technique, we hope to do controlled single molecule chemistry. By this we mean the ability to bring a DNA molecule to a particular enzyme which is immobilized in a specific region of the microchamber and after reaction to move the resultant products to further manipulation or reaction in other regions of the chamber.

The project can be divided into approximately four parts: a) design and operation of the microelectrode array, b) determination of the field strengths in the spaces in the microelectrode array c) manipulation strategies for the orientation, stretching, immobilizing and selection of single DNA molecule in the microscopic stage by the use of the forces induced by the microelectrode array. d) Measurement of the dynamics of electro-optic relaxation, field free relaxation and the viscoelastic properties of single DNA molecules. This will permit the measurement of the optical properties of single oriented DNA at the microscopic level. It allows the observation of alterations of local structure (at the resolution of the microscope objective) monitored through the changes in the polarization parameters or intensity of the fluorescent signals.

The measurement consist of labeling the DNA molecule with a fluorescent intercalating dye, e.g. acridine orange (AO), or ethidium bromide etc. The labeled DNA molecule is then placed on a slide that has micro-electrodes deposited on its surface. The position, motion and shape of the molecule is visualized by the technique of epi-fluorescence microscopy. The fluorescence emitted by the labeled DNA is visualized by and intensified video camera. The voltage of each of the 24 electrodes is controlled by a computer, allowing the manipulation of the DNA on a microscopic scale. The images are continuously recorded for processing by image analysis programs.

As an example of the possible measurements available with this technique, we have studied the field-free relaxation behavior of a single DNA molecule that has been extended by an electric field and allowed to relax at zero field strength under Brownian motion. The length of the stretched DNA is measured as it shortens as a function of time. The relaxation time is a measure of the hydrodynamic forces affecting the DNA molecule, the elasticity of the molecule, and the total length of the DNA (or molecular weight).

When the DNA molecule is totally stretched under the applied electric fields, waves can be seen propagating down the length of the molecule. These waves are analogous to the vibrations seen in a stretched string and are a function of the tension of the string and the modulus of the material. Consequently, the measurement of the velocity of propagation of the waves will give information on the rigidity of the DNA molecule (i.e. the Youngs modulus), a quantity heretofore inaccessible to direct measurement.

Title: Image Acquisition and Processing System for the Analysis of Fluorescence from Stained DNA Gels, Ronald A. McKean

Current methods for analyzing DNA in gels using fluorescence techniques are inadequate since results cannot be easily accessed by a computer and are difficult to reproduce. The lack of instrumentation for converting a fluorescing image into a digitized record for computer entry and the lack of techniques for standardizing analysis performed under varying conditions, severely limit usefulness of DNA separation in gels. Overcoming these problems is essential as needs increase for efficient DNA analysis.

KMS Fusion, Inc. is in the second phase of a DOE sponsored SBIR to develop a system for direct analysis of fluorescence from stained DNA gels. The system development consists of a versatile optical scanner, control and analysis software, and acquisition and control interfaces to an IBM-AT compatible personal computer.

The scanner employs a unique optical system capable of high spatial and photometric resolution, as well as low light level imaging. Scanning techniques are used to image stained DNA directly from agarose or polyacrylamide gels. The scanner incorporates many features, including optical/sensor calibration, modular filter and excitation designs, and selectable image apertures and magnifications. It is operational through front panel controls or an RS232 port. Image data is made immediately available to the computer for further processing.

The software controls the scan procedure, processes image data, creates compact data files, and presents results in a graphical manner. Lane data can be readily standardized to provide results in units of concentration and molecular weight. Statistical features are also provided that allow direct comparison of results from lanes contained in one or more gels.

The system utilizes the popular, and inexpensive, IBM-AT compatible personal computer. Data acquisition electronics and RS232 interface are placed within the computer.

This system offers a practical, inexpensive solution to problems long associated with gel analysis of DNA. Direct quantitation of fluorescence in stained DNA gels allows immediate access to the data by the computer. Computer analysis and processing allow data to be presented in units of concentration and molecular weight. Results from gels electrophoresced under varying conditions can be compared directly. The analyses are presented graphically as plots or reconstructed gel images. Unique data processing techniques allow gel data to be archived using minimal memory.

ESTIMATION OF THE DNA CONTENT OF HETEROMORPHIC AND ABERRANT CHROMOSOMES BY BIVARIATE FLOW KARYOTYPING. Barb Trask*, Ger van den Engh, Joe Gray; Lawrence Livermore National Laboratory, Livermore, CA

For flow karyotyping, chromosomes are analyzed for DNA content and relative base composition on a dual beam flow cytometer. Improvements in sample preparation, instrument accuracy and analysis software allow discrimination of all human chromosomes except 9-12. We have determined the relationship between peak location in a flow karyotype and chromosomal DNA content determined by quantitative microscopy (CYDAC). Maternal and paternal-derived homologs of many chromosomes can be distinguished on the basis of small differences in DNA content (3-5%). Heteromorphism in a population of normal donors was studied. The chromosomes showing the most variation are Y,21,22,13,14,15,16, and 9. The least heteromorphic chromosomes are X,2,4,7,8, and 17. Some variants could be correlated with variation in the size of regions identified by chromosome-specific repetitive sequence probes. DNA contents determined from flow measurements of heteromorphic chromosomes are correlated closely to earlier CYDAC measurements on the same individuals. Family studies show that heteromorphisms are faithfully inherited. Deletion and insertion detection using flow karyotyping will be discussed in light of normal heteromorphism. For example, the DNA content of chromosome 21 can differ by as much as 50% among normal individuals. A series of lines with X chromosome abnormalities with DNA contents ranging from 0.49 to 1.85 times that of a normal X was flow karyotyped. Measured DNA content change was linearly related to that predicted by cytogenetics. Small deletions in chromosome X of ≈ 2 Mbp, below the limit of banding resolution, were detected and quantified using flow karyotyping. Flow karyotyping is also a means to rapidly monitor somatic cell hybrids for the presence of intact human chromosomes.

Work performed under the auspices of the U.S. DOE (contract W-7405-ENG-48) with support from USPHS grant HD-17665.

ABSTRACT TITLE:

Chemiluminescent Imaging of DNA in Electrophoretic Agarose Gels

AUTHORS:

Doris Willis, B.S. Medical Technology,
Paul A. Gray, M.S. Biology,
Rosemarie F. Werba, M.S. Environmental Education,
Robert W. Coughlin, Ph.D. Chemical Engineering, P.E.,
Edward M. Davis, Ph.D. Biochemistry, Symbiotech, Inc.,
8 Fairfield Boulevard, Wallingford, CT 06492,
(203) 284-7465.

ABSTRACT:

We have developed a new chemiluminescent (C.L.) labeling procedure for visualizing DNA in electrophoretic (E.P.) gels that is safe and sensitive. C.L. imaging eliminates the danger of mutagenic dyes such as ethidium bromide (EtBr) and hazardous u.v. light. In addition, C.L. imaging of DNA in electrophoretic gels is more sensitive than EtBr staining. When Lambda DNA in agarose E.P. gels is stained in EtBr (5 ug/ml) for 30 minutes and then destained in 50 mM Tris buffer, pH 6.5, for 30 minutes only 3 ng of DNA is detectable. In contrast, 0.8 ng of DNA is detectable by C.L. imaging. Our procedure employs streptavidin-horseradish peroxidase (SA-HRP), luminol, and peroxide. DNA is biotinylated with photo-active biotin (Photoprobe, Vector Laboratories), electrophoresed in 1% agarose in TBE buffer at 10 v/cm for approximately one hour, affinity labeled with SA-HRP, and soaked in luminol-peroxide solution for 5 minutes. The luminous image of DNA in the E.P. gels is recorded with Polaroid 612 film using contact prints. Both Hind III Lambda DNA and KB DNA ladders have been visualized with this technique.

Single-stranded DNA Imaged using Scanning Tunneling Microscopy

by David Dunlap and Carlos Bustamante

Departments of Chemistry and Pathology, University of New Mexico,
Albuquerque, N.M., 87131

The scanning tunneling microscope has the potential to image biological macromolecules with atomic detail, but recent images of DNA have not realized such resolution. A poor understanding of the imaging mechanisms is partly to blame, but a mechanically stable sample preparation has also been difficult to achieve. We hypothesized that single-stranded DNA, with its exposed, uncharged bases, would adsorb more stably onto highly oriented pyrolytic graphite, making the adsorbed molecules less susceptible to perturbation by the tip. The coincident lateral spreading of the molecules also might make the bases accessible for imaging. After depositing polydeoxyadenylate, which is not self-complementary, we observed molecules aligned in parallel with their bases lying flat and the charged phosphodiester backbone raised upward, in a manner consistent with the hydrophobicity of graphite. A molecular model corresponds well with the images, and also indicates a hydrogen bond that could stabilize the parallel arrangement of the polymer molecules. These micrographs demonstrate the utility of the scanning tunneling microscope for structural studies of nucleic acids and provide evidence that it could be used to sequence DNA.

Imaging of Kinked Configurations of DNA Molecules Undergoing OFAGE Using Fluorescence Microscopy

by Carlos Bustamante and Sergio Gurrieri
Department of Chemistry University of New Mexico.

ABSTRACT

The dynamics of individual DNA molecules undergoing OFAGE (Orthogonal Field Alternating Gel Electrophoresis) has been studied using T2 DNA molecules labelled with a dye and visualized with a fluorescence microscope. The mechanism of reorientation used by a molecule to align itself in the direction of the new orthogonal field, depends on the degree of extension of the chain immediately before the application of this field. The formation of kinks is promoted when time is allowed between the application of the two orthogonal fields so that the molecule attains a partially relaxed configuration. In this case, the chain appears bunched up in domains moving along the contour of the molecule. These regions are found to be the locations where the kinks are formed upon application of the second field perpendicular to the chain. The formation of kinks provide a significant retardation of the reorientation of the molecules, relative to molecules that do not form kinks and appear to play an important role in the fractionation attained with OFAGE. A classification of various reorientation mechanisms observed in molecules that form kinks is presented.

USE OF PULSED-FIELD GEL ANALYSIS TO CONFIRM COSMID OVERLAPS IN CONSTRUCTING AN ORDERED COSMID LIBRARY.

Norman A. Doggett, Lynn M. Clark, Carl E. Hildebrand, Raymond L. Stallings, and Robert K. Moyzis. Genetics Group, LS-3, Los Alamos National Laboratory, Los Alamos, NM 87545.

A chromosome 16 specific cosmid library, constructed from flow sorted material, is being ordered at Los Alamos National Laboratory by a restriction fragment-repetitive DNA hybridization fingerprinting strategy. To complement this "bottom-up" mapping approach we are using pulsed-field gel electrophoresis to confirm overlap between cosmids pairs with minimal overlap or within contigs containing multiple cosmids. DNA from CY18 (a mouse-human somatic cell hybrid from which the the sorted library was constructed) is digested in agarose plugs with several different infrequent cutting enzymes. The digested samples are run in duplicate on the left and right halves of a pulsed-field gel. After transfer, the nylon membrane is cut in half to separate the duplicate digests. These half blots are hybridized separately with either potentially overlapping cosmid pairs or with cosmids from each end of a contig in the presence of human placental competitor DNA. The same banding pattern for each blot confirms overlap. The use of different enzymes effectively eliminates the possibility of a false confirmation due to the coincidental migration of bands. In addition, the digestion with different enzymes ensures that at least a few lanes will have bands that are accurately resolved within the separation range of a given pulsed-field gel. This work was supported by the U.S. Department of Energy under contract W-7405-ENG-36.

U. S. Department of Energy
Human Genome Program
Contractor/Grantee Workshop
Santa Fe, New Mexico
November 3-4, 1989

CONTRACTORS

David Allison
Post Office Box 2009
Oak Ridge, TN 37831-8077

George Bell
T10 MS K710
Los Alamos National
Laboratory
Los Alamos, NM 87545

Martin Burschka
T10 MS K710
Los Alamos National
Laboratory
Los Alamos, NM 87545

Stylianos E. Antonarakis
The John Hopkins
University
CMCS 10-110
School of Medicine
600 N. Wolfe Street
Baltimore, MD 21205

Claire M. Berg
Molecular & Cell Biology
U-131
University of Connecticut
354 Mansfield Road
Storrs, CT 06269-2131

Carlos Bustamante
Chemistry Department
University of New Mexico
Albuquerque, NM 87131

Rod Balhorn
Biomedical Sciences
Division, L-452
Lawrence Livermore
National Laboratory
PO Box 5507
Livermore, CA 94550

Douglas E. Berg
Microbiology Box 8093
Washington University
Medical School
724 S. Euclid
St. Louis, MO 63110

David Callen
Cytogenetics Unit
Adelaide Children's Hospital
King William Road
North Adelaide, S.A. 5006
AUSTRALIA

Benjamin J. Barnhart
Program Manager,
Human Genome
ER-72, GTN
OHER, DOE
Washington, DC 20545

Tony Beugelsdijk
MAA3 MS J580
Los Alamos National
Laboratory
Los Alamos, NM 87545

Charles R. Cantor
Dept. of Genetics &
Development
Columbia University
701 168th St.
Room 1602
New York, NY 10032

Jack Bartley
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

Elbert Branscomb
Biomedical Sciences
Division, L-452
Lawrence Livermore
National Laboratory
PO Box 5507
Livermore, CA 94550

Anthony V. Carrano
Biomedical Sciences
Division L-452
Lawrence Livermore
National Lab
Livermore, CA 94550

Mark Bednarski
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

Christian Burks
T10 MS K710
Los Alamos National
Laboratory
Los Alamos, NM 87545

Thomas C. Caskey
Institute for Molecular
Genetics
Baylor College of Medicine
One Baylor Plaza
Houston, TX 77030

U. S. Department of Energy
Human Genome Program
Contractor/Grantee Workshop
Santa Fe, New Mexico
November 3-4, 1989

Jan-Fang Cheng
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

Edward M. Davis
Symbiotech, Inc.
8 Fairfield Blvd.
Wallingford, CT 06492

Radoje Drmanac
Center for Genetic
Engineering
Genome Structure Unit
283 Vojvode Stepe
P.O. Box 794
11000 Belgrade YUGOSLAVIA

George Church
Department of Genetics
Harvard Medical School
25 Shattuck Street
Boston, MA 02115

Larry L. Deaven
Life Sciences Division
MS M881
Los Alamos National
Laboratory
Los Alamos, NM 87545

David Dunlap
Chemistry Department
University of New Mexico
Albuquerque, NM 87131

Michael Cinkosky
T10 MS K710
Los Alamos National
Laboratory
Los Alamos, NM 87545

Pieter de Jong
Biomedical Sciences
Division, L-452
Lawrence Livermore
National Laboratory
PO Box 5507
Livermore, CA 94550

John J. Dunn
Biology Department
Brookhaven National
Laboratory
Upton, NY 11973

David Corey
Dept. of Chemistry
University of California,
Berkeley
Berkeley, CA 94720

Jeanne Dietz-Band
Life Sciences Division
MS M886
Los Alamos National
Laboratory
Los Alamos, NM 87545

Gary Epling
Dept. of Chemistry, U-60
University of Connecticut
Storrs, CT 06268

Radomir Crkvenjakov
Center for Genetic
Engineering
Genome Structure Unit
283 Vojvode Stepe
P.O. Box 794
11000 Belgrade YUGOSLAVIA

Norman Doggett
Life Sciences Division
MS M886
Los Alamos National
Laboratory
Los Alamos, NM 87545

Glen A. Evans
Salk Institute for Biological
Studies
PO Box 85800
San Diego, CA 92138

J.B. Davidson
Martin Marietta Energy
Systems, Inc.
PO Box 2008
Oak Ridge, TN 37830

Richard J. Douthart
Batelle Pacific Northwest
Laboratory
Battelle Blvd.
Richland, WA 99352

Eric Fairfield
T-10 MS K710
Los Alamos National
Laboratory
Los Alamos, NM 87545

U. S. Department of Energy
Human Genome Program
Contractor/Grantee Workshop
Santa Fe, New Mexico
November 3-4, 1989

Hong Fang
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

Deborah Grady
Life Sciences Division
MS M886
Los Alamos National
Laboratory
Los Alamos, NM 87545

Ed Hildebrand
Life Sciences Division
MS M886
Los Alamos National
Laboratory
Los Alamos, NM 87545

Christopher A. Fields
New Mexico State University
Computing Research
Laboratory
Box 30001
Las Cruces, NM 88003-0001

Joe Gray
Biomedical Sciences
Division, L-452
Lawrence Livermore
National Laboratory
PO Box 5507
Livermore, CA 94550

Diane Hinton
Genome Project
Howard Hughes Medical
Institute
6701 Rockledge Drive
10th Floor
Bethesda, MD 20817

Emilio Garcia
Biomedical Sciences
Division, L-452
Lawrence Livermore
National Laboratory
PO Box 5507
Livermore, CA 94550

Mark Guyer
National Center for Human
Genome Research
NIH
Building 1, Rm. 332
9000 Rockville Pike
Bethesda, MD 20892

Bob Hollen
MS J580
Los Alamos National
Laboratory
Los Alamos, NM 87545

Raymond F. Gesteland
Howard Hughes Medical
Inst.
Department of Genetics
743 Wintrobe Bldg.
University of Utah
Salt Lake City, UT 84132

James F. Hainfeld
Biology Department
Brookhaven National
Laboratory
Building 463
Upton, NY 11973

Leroy Hood
Division of Biology 147-75
California Institute of
Technology
Pasadena, CA 91125

Alex Glazer
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

Richard P. Haugland
Molecular Probes, Inc.
4849 Pitchford Avenue
Eugene, OR 97402

Tim Hunkapiller
Division of Biology 147-75
California Institute of
Technology
Pasadena, CA 91125

Gerald Goldstein
Physical & Technological
Research Division
ER-74 GTN
OHER, DOE
Washington, DC 20545

Gary Hermanson
Salk Institute for Biological
Studies
PO Box 85800
San Diego, CA 92138-9216

Marge Hutchison
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

U. S. Department of Energy
Human Genome Program
Contractor/Grantee Workshop
Santa Fe, New Mexico
November 3-4, 1989

K. Bruce Jacobson
Biology Division
Oak Ridge National
Laboratory
PO Box 2009
Oak Ridge, TN 37831-8077

Ivan Labat
Center for Genetic
Engineering
Genome Structure Unit
283 Vojvode Stepe
P.O. Box 794
11000 Belgrade YUGOSLAVIA

Po-Yung Lu
Biomedical and
Environmental Information
Analysis
Oak Ridge National Lab.
PO Box 2008, Bldg. 2001
Oak Ridge, TN 37831-6050

Joe Jaklevic
Instrumentation Division
70A-2205
Lawrence Berkeley
Laboratory
Berkeley, CA 94720

Esther Leeftang
Dept. of Chemistry
University of California,
Davis
Davis, CA 95616

Marcos Maestre
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

J. H. Jett
Life Sciences Division
Los Alamos National
Laboratory
PO Box 1663
Los Alamos, NM 87545

Leonard S. Lerman
Massachusetts Institute of
Technology 56-743
Department of Biology
77 Massachusetts Avenue
Cambridge, MA 02139

Donna Maglott
American Type
Culture Collection
12301 Parklawn Drive
Rockville, MD 20852-1776

Bill Johnston
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

Suzanna Lewis
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

Betty Mansfield
Human Genome Mgmt.
Information System
Oak Ridge National Lab.
PO Box 2008, Bldg. 2001
Oak Ridge, TN 37831-6050

Roland Johnson
L-157
Lawrence Livermore
National Lab.
PO Box 5507
Livermore, CA 94550

Hwa Lim
Supercomputer
Computations Research
Institute
467 SCL
Florida State University
Tallahassee, FL 32306-4052

John Martin
LS4 MS M888
Los Alamos National
Laboratory
Los Alamos, NM 87545

Richard Keller
MS G738
Los Alamos National
Laboratory
Los Alamos, NM 87545

Jonathan Longmire
Life Sciences Division
MS M881
Los Alamos National
Laboratory
Los Alamos, NM 87545

Richard A. Mathies
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

U. S. Department of Energy
Human Genome Program
Contractor/Grantee Workshop
Santa Fe, New Mexico
November 3-4, 1989

Mary Kay McCormick
Life Sciences Division
MS M886
Los Alamos National
Laboratory
Los Alamos, NM 87545

Robert K. Mortimer
Dept. of Molecular and
Cellular Biology Div.
UC Berkeley
Berkeley, CA 94720

Frank Ogletree
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

Donald McMullen
455 Supercomputer
Computations Research
Institute
Florida State University
Tallahassee, FL 52306-4052

R. K. Moyzis
Life Sciences Division
MS M886
Los Alamos National
Laboratory
Los Alamos, NM 87545

Arnold Oliphant
Medical Informatics
University of Utah
420 Chipeta Way, Rm. 180
Salt Lake City, UT 84108

Dr. Mortimer L. Mendelsohn
Biomedical & Sciences Div.
Lawrence Livermore
National Lab.
PO Box 5507
Livermore, CA 94550

David Nelson
Computations, L-305
Lawrence Livermore
National Laboratory
Livermore, CA 94550

Frank Olken
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road
Berkeley, CA 94720

Julianne Meyne
LS-3 MS M886
Los Alamos National
Laboratory
Los Alamos, NM 87545

David L. Nelson
Institute for Molecular
Genetics
Baylor College of Medicine
One Baylor Plaza
Houston, TX 77030

Anne Olsen
Biomedical Sciences
Division, L-452
Lawrence Livermore
National Laboratory
PO Box 5507
Livermore, CA 94550

Harvey Mohrenweiser
Biomedical Sciences
Division, L-452
Lawrence Livermore
National Laboratory
PO Box 5507
Livermore, CA 94550

Debra Nelson
T10 MS K710
Los Alamos National
Laboratory
Los Alamos, NM 87545

Ross Overbeek
MCS 221 C-216
Argonne National
Laboratory
Argonne, IL 60439

Frederick Morse
Associate Director for
Research
MS A-114
Los Alamos National
Laboratory
Los Alamos, NM 87545

William C. Nierman
American Type Culture
Collection
12301 Parklawn Drive
Rockville, MD 20852-17

Rob Pecherer
T10 MS K710
Los Alamos National
Laboratory
Los Alamos, NM 87545

U. S. Department of Energy
Human Genome Program
Contractor/Grantee Workshop
Santa Fe, New Mexico
November 3-4, 1989

Donald Peters
Biomedical Sciences
Division, L-452
Lawrence Livermore
National Laboratory
PO Box 5507
Livermore, CA 94550

Karl Sirotkin
T10 MS K710
Los Alamos National
Laboratory
Los Alamos, NM 87545

Raymond Stallings
Life Sciences Division
MS M886
Los Alamos National
Laboratory
Los Alamos, NM 87545

Robert Ratliff
Life Sciences Division
MS M886
Los Alamos National
Laboratory
Los Alamos, NM 87545

Tom Slezak
Biomedical Sciences
Division, L-452
Lawrence Livermore
National Laboratory
PO Box 5507
Livermore, CA 94550

M. Stodolsky
ER-72, GTN
OHER, DOE
Washington, DC 20545

Mark Salmeron
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road, 50B-3238
Berkeley, CA 94720

Cassandra Smith
1602 HHSC
College of P&S
701 W 168
New York, NY 10032

F. William Studier
Biology Department
Brookhaven National Lab
Upton, NY 11973

Karen Schenk
T10, MS K710
Los Alamos National
Laboratory
Los Alamos, NM 87545

David Smith
ER-72, GTN
OHER, DOE
Washington, DC 20545

Betsy Sutherland
Biology Department 463
Brookhaven National
Laboratory
Upton, NY 11973

Jeffrey Schmaltz
ER-72 GTN
OHER, DOE
Washington, DC 20545

Carol Soderlund
New Mexico State University
Computing Research
Laboratory
Box 30001
Las Cruces, NM 88003-001

Stanley Tabor
Dept. of Biological
Chemistry & Molecular
Pharmacology
240 Longwood Ave.
Harvard Medical School
Boston, MA 02115

Paul Silverman
Associate Laboratory
Director
Lawrence Berkeley Lab.
1 Cyclotron Road
50A, 5104
Berkeley, CA 94720

Sylvia Spengler
Human Genome Center
Lawrence Berkeley
Laboratory
1 Cyclotron Road, 1-213
Berkeley, CA 94720

Thomas Tenforde
Batelle Pacific Northwest
Laboratory
Batelle Blvd.
Richland, WA 99352

U. S. Department of Energy
Human Genome Program
Contractor/Grantee Workshop
Santa Fe, New Mexico
November 3-4, 1989

David Thurman
Batelle Pacific Northwest
Laboratory
Batelle Blvd.
Richland, WA 99352

Mark Wagner
Biomedical Sciences
Division L-228
Lawrence Livermore
National Laboratory
Livermore, CA 94550

Robert Weiss
Howard Hughes Medical
Inst.
Dept. of Human Genetics
743 Wintrobe Bldg.
University of Utah
Salt Lake City, UT 84132

David Torney
T10 MS K710
Los Alamos National
Laboratory
Los Alamos, NM 87545

Robert P. Wagner
LS3 MS M886
Los Alamos National
Laboratory
Los Alamos, NM 87545

Sherman Weissman
Yale University
Dept. of Human Genetics
PO Box 3333
333 Cedar Street
New Haven, CT 06510

Barbara Trask
Biomedical Sciences
Division L-452
Lawrence Livermore
National
Laboratory
Livermore, CA 94550

Ronald A. Walters
Program Director
Biological & Environment
MS A114
Los Alamos National
Laboratory
Los Alamos, NM 87545

Patricia Wilkie
Biomedical Sciences
Division L-452
Lawrence Livermore
National Laboratory
Livermore, CA 94550

James Treves
L-473
Lawrence Livermore
National Laboratory
Livermore, CA 94550

Denan Wang
Human Genome Center
Lawrence Berkeley Lab
1 Cyclotron Road
Berkeley, CA 94720

Gayle Woloschak
Argonne National
Laboratory
BIN-202
9700 S. Cass Avenue
Argonne, IL 60439

Kathryn Tynan
Biomedical Sciences
Division L-452
Lawrence Livermore
National Laboratory
Livermore, CA 94550

R. J. Warmack
MS 6123
PO Box 2008
Oak Ridge National Lab.
Oak Ridge, TN 37831-6123

John Wooley
Director
Info. & Resources Division
National Science Foundation
1800 G Street NW
Washington, D.C. 20550

Marvin Van Dilla
Biomedical Sciences
Division
Lawrence Livermore
National Laboratory
Livermore, CA 94550

John Wassom
Director
Human Genome &
Toxicology Program
Oak Ridge National Lab.
PO Box 2008, Bldg. 2001
Oak Ridge, TN 37831-6050

Judy Wyrick
PO Box 2008
Oak Ridge National
Laboratory
Oak Ridge, TN 37831-6050

U. S. Department of Energy
Human Genome Program
Contractor/Grantee Workshop
Santa Fe, New Mexico
November 3-4, 1989

E.S. Yeung
Dept. of Chemistry
85 E. Gilman Hall
Iowa State University
Ames, Iowa 50011

Edward David Hyman
Sybtrel
3330 N. Causeway Blvd.
Metairie, LA 70002

INDUSTRY

George Bers
Bio-Rad Labs
1414 Harbour Way S.
Richmond, CA 94804

Kathy Yokobata
Biomedical Sciences
Division L-452
Lawrence Livermore
National Laboratory
Livermore, CA 94550

Andrew Kumamoto
California Institute of
Biological Research
11099 North Torrey Pines
Rd.
La Jolla, CA 92037

David Botten
EG&G Biomolecular
4 Tech Circle
Natick, MA 01760

SBIR GRANTEES

Irena Bronstein
Tropix, Inc.
47 Wiggins Avenue
Bedford, MA 01730

Michael McClelland
California Institute of
Biological Research
11099 North Torrey Pines
Rd.
La Jolla, CA 92037

Thomas Brennan
Genomyx
460 Pt. San Bruno Blvd.
So. San Francisco, CA 94080

Jim Brulé
Coherent Research, Inc.
100 East Washington Street
Syracuse, NY 13202

Patricia McGrath
Tropix, Inc.
47 Wiggins Avenue
Bedford, MA 01730

Robert Brown
Oracle Corp.
120 Wood Avenue S.
Iselin, NJ 08830

Sukhendu Dev
BTX, Inc.
3472 Jewell Street
San Diego, CA 92109

Ronald McKean
KMS Fusion
700 KMS Place
PO Box 1567
Ann Arbor, MI 48106

John W. Chase
Scientific Director
U.S. Biochemical
Corporation
26111 Miles Road
Cleveland, OH 44128

Gunter Hofmann
BTX, Inc.
3472 Jewell Street
San Diego, CA 92109

Jeffrey M. Stiegman
Biophotonics Corporation
P.O. Box 3756
Ann Arbor, MI 48106

Stephen Dias
Natural Language Inc.
879 W. 190th Street
Los Angeles, CA 90248-4214

U. S. Department of Energy
Human Genome Program
Contractor/Grantee Workshop
Santa Fe, New Mexico
November 3-4, 1989

Robert Eades
Cray Research, Inc.
1333 Northland Drive
Mendota Heights, MN 55120

Robert Jones
Applied Biosystems
850 Lincoln Center Drive
Foster City, CA 94404

Bonner Nishida
IBM
Neighborhood Road
MS 228
Kingston, NY 12401

Steve Ferris
Bio-Rad Laboratories
1414 Harbour Way South
Richmond, CA 94804

Eva Juhos
Beckman Instruments Inc.
1050 Page Mill Road
Palo Alto, CA 94304

Hans Osterhoudt
Strategic Planning
Eastman Kodak Co.
Life Science Division
343 State Street
Rochester, NY 14650

David Fram
Oracle Corp.
901 Mariners Island Blvd.
San Mateo, CA 94404

Leonard Klevan
Life Technologies Inc.
8717 Grovemont Circle
PO Box 6009
Gaithersburg, MD 20877

Dieter Rabussay
Life Technologies Inc.
8717 Grovemont Circle
PO Box 6009
Gaithersburg, MD 20877

Okan Gurel
IBM
101 Main Street
Cambridge, MA 02142

Terry Lerner
Integrated Genetics, Inc.
Framingham, MA 01701

Ira Schildkraut
New England Biolabs Inc.
32 Tozer Road
Beverly, MA 01915

John Harding
Life Technologies Inc.
8717 Grovemont Circle
PO Box 6009
Gaithersburg, MD 20877

Michael Liebman
Amoco Technology
Biotechnology Division
PO Box 3021
Naperville, IL 60056

Larry Stewart
CRAY Research, Inc.
555 Oppenheimer
Suite 100
Los Alamos, NM 87544

W. Charles Johnson
Beckman Instruments
2500 Harbor Blvd.
Fullerton, CA 92634

Mike Nemzek
EG&G Biomolecular
4 Tech Circle
Natick, MA 01760

Mark Sullivan
Eastman Kodak
Life Science Division
343 State Street
Rochester, NY 14650

U. S. Department of Energy
Human Genome Program
Contractor/Grantee Workshop
Santa Fe, New Mexico
November 3-4, 1989

David Turek
IBM
Neighborhood Road
MS 228
Kingston, NY 12401

John West
BioAutomation Inc.
A-204
Front & Ford Streets
Bridgeport, PA 19405

Michael Zoccoli
Cetus Corporation
1400 Fifty-Third St.
Emeryville, CA 94608