



Transposon-Mediated cDNA Sequencing at the British Columbia Cancer Agency, Genome Sciences Centre

D. Smailus, J. Asano, Y. Butterfield, N. Girn, R. Guin, M. Krzywinski, S. Lee, K. MacDonald, T. Olson, P. Pandoh, P. Saedi, U. Skalska, L. Spence, J. Stott, S. Taylor, K. Teague, G. Yang, J. Schein, S. Jones and M. Marra.

British Columbia Cancer Agency
Vancouver, British Columbia, Canada

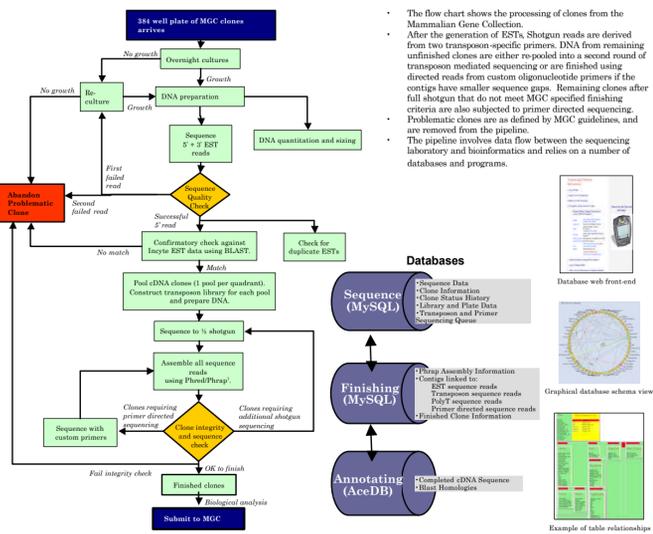
Genome Sciences Centre

www.bcgsc.ca
info@bcgsc.ca

Abstract

We have developed an efficient, high-throughput method for accurate DNA sequencing of entire cDNA clones through our participation in the NCI-sponsored Mammalian Gene Collection. Sequencing is accomplished through the insertion of Mu transposon into cDNAs, followed by sequencing reactions primed with Mu-specific sequencing primers. Transposon insertion reactions are not performed with individual cDNAs but rather on pools of up to 96 clones. Accurate clone insert size and DNA quantitation data are used to ensure proportional representation of each cDNA clone in the pool. This pooling strategy reduces the number of transposon insertion sequencing libraries that would otherwise be required, reducing the costs and enhancing the efficiency of the transposon library construction procedure. Sequences are assembled using Phred, Phrap, and Consed to yield the full-length cDNA sequence, with sequence editing and other sequence finishing activities performed as required to resolve sequence ambiguities. We are currently in our second year of the MGC project and have used the method to generate more than 7.5 Mb of finished sequence from 3,956 candidate full-length cDNAs. Analysis of 22,785 sequenced Mu transposon insertion events revealed a weak sequence preference for Mu insertion. However, the insertion pattern deviates only slightly from random and does not adversely affect the efficacy of our method. A detailed description of our transposon-mediated sequencing methodology and analysis of Mu transposon insertion events will be presented.

Pipeline Flow Chart



Using the database, the entire process of DNA preparation, sequencing, data analysis and storage are tracked and stored. All aspects of run conditions and analysis are kept in database tables to facilitate locating trends in performance, keeping quality control and compiling statistics.

DNA Isolation for EST Sequencing



DNA Quantitation, EST Analysis and Clone Sizing

DNA purifications are quantitated in a 96-well spectrophotometer and the data stored for automatic retrieval. Generate 3' and 5' ESTs / Screen Clones. 5' and 3' reads for each clone are generated by Big Dye Terminator or cycle sequence analysis using -21M13 Forward primer, T7 and Oligo-dT Plus primers respectively on ABI 3700 DNA analyzers. We currently perform 0.25X chemistry in 4ul reaction volumes containing approximately 40ng of DNA. Successful reads are checked against Incyte or Agilent data and MGC IDs. Clones with failed reads are re-cultured; DNA isolated, sequenced and checked. Clones failing MGC ID checks twice are abandoned.

384 EcoRI restriction digested clones are sized by loading onto four 121-lane combs embedded into a single agarose gel (96 samples and 25 marker lanes per comb). The gel contains a DNA marker in every fifth lane. After electrophoresis, the gel is stained with SYBR Green and an image of the gel is collected on a Molecular Dynamics Fluorimager. Lanetracking with Image. Automated restriction fragment identification with BandReader.

EST Analysis

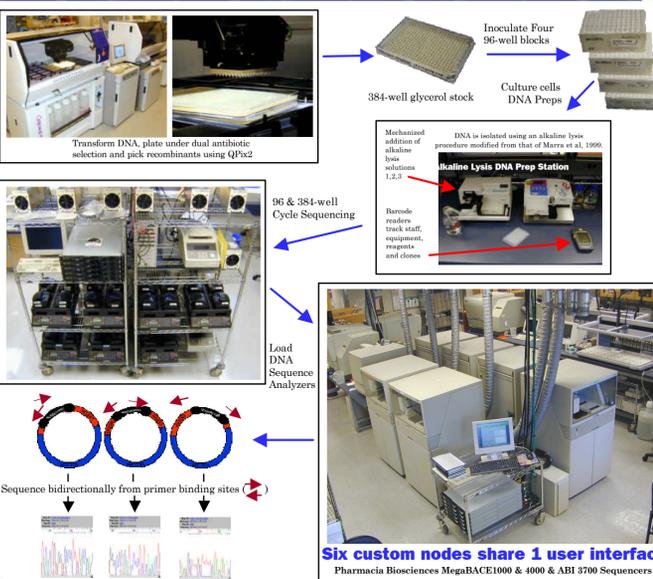
Analysis of data that are generated at various stages in the pipeline allows for both manual and automatic assessments of each clone. The first step is to confirm the identity and quality of clones that we have received by matching the EST reads against existing EST sequences generated from the pipeline and reported as problematic. Clones that do not result in a positive match are automatically removed from the pipeline and reported as problematic. Failure to generate a confirmatory EST read can result from no growth or cross-contaminated wells. The analyses are summarized in varying levels of detail on our cDNA web pages. Also, ESTs are checked for matches to other ESTs from the same quadrant in order to identify duplicate clones. Clone sizes are determined by restriction enzyme analysis and band called using Image (www.sanger.ac.uk/Software/Image) and BandReader (Fuhrman, Dan et al. 2001, Automated Image Analysis for DNA Fingerprinting unpublished). The accuracy of the sizing is exemplified by the graph (shown at right) that compares the completed sequence length against size as determined by the gel.

Table with columns: MGC ID, Location, Status, Contig, Get (date), Self Match, Incomplete, Score, Percent. Lists various MGC clones and their associated data.

Pool cDNA Clones and Add Entransposons™

Molar Ratio Calculator. Two oligonucleotides that prime from either end of the randomly incorporated Entransposons are used to shotgun sequence pooled clones. Total number of transformed clones to sequence for each shotgun library is calculated automatically at 1:1 read/bk of target DNA ("half shotgun"). Incorporate Entransposons™ into pooled cDNA clones. Transform, plate under dual antibiotic selection and pick recombinants using QPix2.

Transposon-Mediated cDNA Sequencing Pipeline



Web Tools

Sequence Integrity Checks. EST matches Incyte EST? Both 5' and 3' EST from clone? Restriction and linker sites on each end? Poly A tail found? Assembly size and gel size agree to within 10%? Error rate of 1/50000? Contig assembly information from the Phrap output file for each build is parsed into an SQL Finishing database. This database allows for the quick assessment of contigs and identification of clones. Web tools shown at right and below facilitate the automatic visualization of the required checks and finishing of clones without having to look at each individual contig. Information such as clone and assembly sizes, quality, the number of reads, and various sequence integrity checks are available. The interface also allows for manually changing the status of clones when needed. Other status changes are automatically made by various scripts.

Statistics

Summary statistics of the sequencing pipeline, including total clones, reads, and completion rates. Includes a graph showing the number of clones completed over time.

Mu Transposon Target Site Analysis

Mu Transposon Target Site Base Composition. The generation of such a large amount of sequence data has also allowed for a comprehensive survey of the insertion profile of the Mu transposon in the insert. An analysis of the insert region that spans this target site describes a consensus sequence preference for Mu transposon insertion. The insertion site displays a symmetry that includes the target site consisting of pyrimidines followed by purines as shown in the table to the right and graphed above. Statistical analysis and our overall observations show that at least for the cDNA sequences used in this study, the effective insertion profile does not differ greatly from a random model and has allowed efficient sequencing of these templates.

Comparison of the frequency of 5-mers occurring cDNAs with the frequency of 5-mers utilized in 22,785 transposon insertion events. Mu transposon insertion deviates only slightly from random. The insertions of Mu into 1,242 cDNA clones were analyzed using the binomial test and assigned to bins (Materials and Methods). The resulting p-values reflect the likelihood that the observed insertion events were not random. Plotted are the numbers of bins grouped into p-value ranges of 0.01. P-values of greater than 0.05 correspond to bins for which the observed insertion events are likely to be random. P-values of 0.05 or less (indicated by shaded bars) correspond to bins for which the observed insertion events cannot be confidently described as random occurrences.

Analysis of transposon insertions in the pOTB7 vector. A set of sequences (5,552) were analyzed and the relative positions of transposon insertions within the vector were mapped. Of the sequences in this set, 22% (1,233) were observed to initiate within the vector. If transposon insertion into the vector was random and all insertion events were recovered, we would expect to observe 46% (2,547) of reads initiating in the vector. However, the observed results match more closely to that expected (18% insertion into vector) assuming zero insertions into the vector (chimeric acetyl acetyl) transposon of replication (ori). This is reflected in no notably fewer observed transposon-generated sequence reads initiating from these regions. Transposon insertions into the vector origin of replication are effectively lethal. When chloramphenicol resistance is included during selection of transposon-containing clones, insertions into the cat gene are also effectively lethal. This lethality results in a decreased number of recovered vector insertion events. See Y. Butterfield et al. 2002, NAR Vol30, #11.

Table comparing expected vs observed insertion events along the entire vector and assuming zero insertions in cat and ori.

Biological Checks and Annotating

Biological Analysis. Full length (poly A tail and start codon)? Potentially sequenced? Intragenomic sequences? No obvious ORF? Mitochondrial genomic DNA? Duplication? "Classic" clone? Insertion/deletion? The coding element and positive function of each finished clone is analyzed within ACEDB. Sequence similarity searching is done with each clone against protein and EST databases and alignment of cDNA to genomic sequence. Homologues to ESTs. Homologues to proteins and observed cDNAs. MGC 22887 Q16760 CARBONIC ANHYDRASE IX. TISSUE SPECIFICITY EXPRESSED PRIMARILY IN CARCINOMA CELLS LINES. PROTEOLYSIS IS RESTRICTED TO VERY FEW NORMAL TISSUES AND IS NOT ABUNDANTLY EXPRESSED IN THE EPITHELIAL CELLS OF GASTRIC MUCOSA.

Summary

Summary table showing pipeline progress by plate number, including received date, days in pipeline, status, and clone breakdown (analyzed, problematic, remaining).

16 plates n/a In Progress 90.8% completed 4521 1057 566. *Status - covers the progress of the plate in terms of completed clones, not the degree of sequencing done on the plate. *Analyzed - includes all completed clones which have been analyzed by a finisher and have also been identified as requiring no further sequencing effort to finish. These clones are either submitted or are in the queue to be finished and checked biologically. This category does not include clones that have been marked as problematic but may include clones that essentially may be identified as problematic from the finished sequence. *Problematic - includes all completed clones that have been marked as problematic and hence require no further sequencing or inspection. These clones are completed and have been or will be placed into a problem clone report. *Remaining - includes all clones of neither of the above two categories. Either the clones have not yet been through all stages of the pipeline yet or further sequencing/finishing is required to finish the clones or identify it as problematic.



Dr. Michael Smith 1932-2000 Founding Director