

# A physical map of the human genome – supplemental information and detail.

The International Human Genome Mapping Consortium\*

[\\*Complete author list.](#)

## Regional approach to large-scale map physical map construction.

Following the positional cloning model, most centers employed sequence tagged site(STS)<sup>1</sup> landmark-content maps to identify regional markers suitable for isolating clones. Once identified, regional clones were characterized by either STS-content alone or in conjunction with restriction enzyme digest fingerprinting methods and assembled into contigs. As an example, the availability of a high-quality, high-resolution yeast artificial chromosome (YAC)<sup>2</sup> STS-content map of chromosome 7<sup>3</sup> made it a logical first target for the Washington University Genome Sequencing Center (WUGSC) human genome sequencing pilot project. In early 1996 the WUGSC began experimenting with an STS-driven clone identification paradigm using collective markers from approximately one megabase intervals simultaneously for clone identification. Initially, clones from a Genome Systems BAC library (Human BAC Release I; <http://www.incyte.com/reagents/index.shtml>) and from the early California Technical Institute BAC library ([http://informa.bio.caltech.edu/Bac\\_info.html](http://informa.bio.caltech.edu/Bac_info.html)) were identified by the PCR using specific STSs in the laboratory of Eric Green, (NHGRI) and the well locations of STS-positive clones were communicated to the mapping group at the WUGSC. Simultaneously, methods for high-throughput hybridization-based BAC filter screening were developed at the WUGSC (J. McPherson). The general approach involved screening genomic BAC and PAC libraries by hybridization using overgo probes<sup>4</sup> to identify clones corresponding to specific STS markers. Overgo probes are made by filling in the single-stranded overhangs of two overlapping oligonucleotides using radiolabeled nucleotides and Klenow polymerase. Typically, two 24mers overlapping by 8 bp were used to generate a radiolabeled double-stranded 40mer. Overgo probes were arranged in three-dimensional arrays with 6 probes on each axis (6x6x6=216 probes each). A five directional pooling strategy allowed resolution of 80-90% of all markers with only 30 hybridizations. To date, greater than 25,000 human and mouse markers have been associated with BACs using this probe type at the WUGSC (J. McPherson). These efforts targeted the Genome Systems BAC library and the RPCI-4 and RPCI-5 PAC libraries

constructed by P. de Jong (<http://www.chori.org/bacpac/>). Evolution of the BAC filter screening method resulted in a procedure capable of simultaneous multiplex hybridization of 36 overgo<sup>4</sup> probes at a rate of 450 markers weekly. Once identified, BAC clones were retrieved from 384-well glycerol stocks and colony purified by re-streaking on agar plates. Individual colonies were cultured and DNA purified in 96-well format. DNA was digested with HindIII, the resulting restriction products resolved by agarose gel electrophoresis and analyzed using the FPC fingerprint assembly software package<sup>5,6</sup>. Subsequent efforts utilizing this strategy focused on chromosomes 2 and Y.

### **Whole Genome BAC map**

To generate fingerprint maps of the *Caenorhabditis briggsae* and *Arabidopsis thaliana* genomes, we had evolved the capacity to produce weekly approximately 2,000 BAC fingerprints with a team of six technicians. Each gel analyzed consisted of 40 BACs and 11 standard “marker” DNA lanes. We calculated that to achieve our aim of 300,000 human BAC fingerprints using this gel technology would require 7,500 gels. We devised a tandem 121-lane agarose gel format, allowing the simultaneous electrophoresis of 50 standard “marker” DNA lanes and of 192 BAC restriction digests, thereby reducing the number of gels required. This system, with the increased density of fingerprint information on the gels, would produce the target 300,000 fingerprints in only 1,600 gels. We verified that the new agarose gel format produced accurate fingerprints by comparing fragment sizes between clones prepped and digested in duplicate and by fingerprinting previously sequenced clones and comparing the *in silico* digests of the sequence to the agarose gel fingerprints as shown in [Figures 1a and 1b](#), respectively. With these and other improvements in the fingerprinting technology, and with addition of staff, throughput rose from 2,000 fingerprints weekly to a high of more than 20,000 weekly, the latter number representing approximately one human genome equivalent of fingerprinted DNA produced each week. The majority of the clones fingerprinted were from 3 libraries, RPCI-11, RPCI-13 and CT-C/D1. The fingerprint characteristics of the clones from these libraries are listed in [Table 1](#).

### **Manipulation of fingerprint data**

Due to considerable variability in the mobility of small fragments and that these small fragments were inconsistently identified during subsequent gel image analysis. Hence, fragments smaller than 600 base pairs (corresponding to a an Image mobility value greater than 3590) were removed prior to assembly. In addition, because of the difficulty in determining multiplets (cases where more than one fragment is located at nearly the same position on the gel), all fragments within mobilities of 8 were collapsed to only a single

band in the resulting fingerprint. This “sanitizing” process resulted in assemblies of increased reliability. Assembly parameters were determined empirically by changing various build parameters to obtain assemblies and then evaluating them for consistency using other mapping data associated with clones in the assembly (primarily radiation hybrid (RH) chromosomal localization data from the Stanford Human Genome Center (SHGC; D.R. Cox)). After assembly of the fingerprints with each set of parameters, the number of chimeric assemblies was detected by the presence of conflicting map information affixed to clones in the assembly. Optimal assembly parameters minimizing both the total number of assemblies and the number of chimeric assemblies are given in [Table 2](#). Incremental assemblies of the fingerprints accumulated throughout the project are shown in [Figure 2](#).

Using the FPC parameter values listed in [Table 2](#), an automated assembly of the 283,287 clones resulted in 7,133 assemblies that contained 93% of all fingerprints in the database (December 1999). The remaining unincorporated clones (*i.e.* singletons) were excluded as they contained fewer than the 3 fragments specified as the minimum number required for assembly (min band = 3) or contained fewer than the total of 6 fragments that were empirically determined necessary for automated assembly under these conditions. Although the initial automated assemblies produced when using the FPC parameter values in [Table 2](#) were generally reliable with respect to the “binning” of related clones, determination of the correct relative order of the clones within the assemblies remained a manual task.

### **Achieving map continuity by manual editing of assemblies**

The initial goals of the manual assembly editing were to refine the relative ordering of the clones within contigs, to identify joins between contigs and to disassemble larger chimeric contigs. This process involved first editing the fingerprint assemblies, using the tools encapsulated in FPC, to ensure that every clone within a contig was properly situated with respect to its most highly related neighbors. A detailed discussion of the manual editing process can be found in Marra *et al.*, 1999.<sup>7</sup> Briefly, this was determined by minimizing the “Sulston” score<sup>8</sup> (a statistical measure of the of coincidental overlap) between adjacent clones and inspecting the fingerprint data extracted from the original gel image. Redundant clones contributing no unique fragment data to the assembly were also “buried” under their parent clones to simplify the contig view if this had not already been done in the automated process. During this manual inspection, additional clones were incorporated into the contig from the remaining pool of previously unassembled clones and newly fingerprinted clones that were not yet part of any assembly. Once the clone order had been established within each contig, clones at the extreme ends of each contig were used to query the FPC database at a reduced required fingerprint overlap stringency from that used to perform the initial assembly. In this way, potential joins between contigs were identified. Fingerprints of clones involved in potential joins were visually inspected to confirm that all restriction fragments were logically consistent and the joins were made if appropriate. The results of this editing activity are illustrated in [Table 3](#). The most notable effect of the intensive editing effort was the nearly 5-fold reduction of total assemblies. The large increase in the number of singletons was due to the continued addition of new clones to the database without a total assembly being performed. Clones from the singleton pool were only manually incorporated into assemblies as needed to extend the ends of the contigs. As a rule, clones from the singleton pool were not incorporated into the interior of contigs where additional redundancy was not needed.

### **Integration of map data with the BAC contig database**

The main purpose of generating a whole genome BAC contig map was to coordinate activities of the larger sequencing groups by making the FPC database available and to specifically select clones to feed the human sequencing capacities of the WUGSC, the Whitehead Institute for Biomedical Research (WIBR) and the Stanford Genome and Technology Center (SGTC). An essential component of the working draft strategy was that the selected clones for these centers be restricted as much as possible to chromosomes 2, 4, 7, 8, 11, 15,

17,18 and Y. Early on, mapped BACs were identified primarily from the hybridization of 13,695 overgo probes generated from sequences mapped to these chromosomes. These consisted predominantly of selected markers from the CEPH Généthon genetic map<sup>9</sup>, the GeneMap'99 genome-wide RH map (<http://www.ncbi.nlm.nih.gov/genemap/>).<sup>10-12</sup> and from plasmid library sequences prepared from flow-sorted chromosomes (Sanger Centre, unpublished data). The 13,695 markers identified 96,283 unique clones, providing many anchor points for the assembled contigs. As the BAC assemblies enlarged and the fingerprint resource gained popularity as a tool for coordinating large-scale sequencing efforts, markers across the genome were sought to anchor all contigs. The remainder of the genome was to be sequenced by other members of the International Sequencing Consortium;<sup>13</sup> however, it was realized that the WG BAC map would also be a useful resource for other centers. To enhance the utility of the map to all centers, as many markers as possible were integrated into the FPC database, resulting in the chromosomal localization of most BAC contigs. Because the RPCI-11 library was being used for many genome initiatives there was a wealth of marker information available. A significant source of such information came from the inclusion of chromosomal assignment data for 9,018 STS derived from BAC end sequences (D.R. Cox, unpublished data). Although greater than 15% of the available BAC end sequences of the RPCI-11 library are apparently mislabeled with respect to the correct BAC name<sup>14</sup>, the number of accurately mapped BAC ends facilitated the correct chromosomal assignment of contigs by considering the consistent trend of these data for a given contig. As the working draft sequence accumulated, known markers were readily identified utilizing electronic PCR (ePCR; a program that searches sequence for STS by identifying the associated primer sequences in the correct orientation and with correct spacing),<sup>15</sup> which identifies primers and these data provided for inclusion in the FPC database (G. Schuler, NCBI). The combined ePCR and hybridized data sets contained 69,507 markers, including 1,659 polymorphic markers from the Généthon genetic map. In addition, chromosomal assignment and integration of cytogenetic map positions were achieved by identification of 3,412 BACs mapped by fluorescence *in situ* hybridization (FISH) data as described in an accompanying manuscript.

The GeneMap'99 map was chosen for the purpose of anchoring and ordering contigs as it has a substantial marker set (> 50,000), is well-integrated with the Généthon genetic map and provides local ordering at <1 Mb resolution. It was also the most widely known genome map at that time. To ensure that markers assigned to BACs provided consistent localizations, single markers associated with clones in multiple contigs were not used for contig placement. For each contig, all chromosome assignment data associated with its clones were tallied and the contig was assigned to the chromosome with the majority of supportive data. Each chromosomal assignment datum was given equal weight with the exception of FISH data, which was considered the strongest evidence for placement. A clone could have more than one piece of supportive datum. For example, a clone with one chromosome 7 STS marker, chromosome 7 FISH data and BAC end sequence RH data for chromosome 2 would contribute two supportive data points for contig assignment to chromosome 7 and one for chromosome 2. In the case of a tie in the cumulative data for the contig, no initial chromosomal assignment was made.

Once a chromosomal assignment was made, the majority of contigs could be further localized to the median RH map position of all markers associated with the contig and assigned to that chromosome. To determine the median map position of a contig: 1. Each clone was allowed only one map position. This was the median position of all markers associated with that clone; 2. Markers that are likely outliers were

removed. A contig may have outliers if the map distance covered by the contig is greater than expected for a contig with the given number of clones. The outlier markers likely represent local mapping errors; 3. The median position was taken from the markers that remain.

The orientation of the contig with respect to the map was determined next. The marker content of accessioned sequences associated with the clones from the contig was determined by ePCR. Those markers found on one of five maps were examined (Généthon genetic map, Marshfield genetic map<sup>16</sup>, WIBR YAC STS-content ([http://carbon.wi.mit.edu:8000/cgi-bin/contig/phys\\_map](http://carbon.wi.mit.edu:8000/cgi-bin/contig/phys_map)), GeneMap'99, SHGC G3 RH map (<http://www-shgc.stanford.edu/Mapping/rh/index.html>), NCBI framework map (<http://www.ncbi.nlm.nih.gov/genome/guide/>)) and all mapped markers were then compared with their adjacent neighboring markers to determine if the order of markers within a contig for any given map were ascending or descending with respect to that map. Contigs were flipped as needed to orient them with respect to the majority consensus of all maps examined.

Lastly, finished map data provided by other groups for several chromosomes was examined. Although the contigs ordered by the methods above were largely consistent with these data, a few manual adjustments were made to more accurately reflect these well-characterized maps. The maps considered at this stage were for chromosomes 14 (Genoscope, J. Weissenbach)<sup>17</sup>, 19 (Lawrence Livermore National Laboratory, [http://www-bio.llnl.gov/bbrp/genome/html/chrom\\_map.html](http://www-bio.llnl.gov/bbrp/genome/html/chrom_map.html)), X (Sanger Centre, D. Bentley) and a 20 Mb segment of chromosome 15 (University of Washington, L. Rowen). Telomeric contigs were identified and positioned where possible. The data for the telomeric contigs is discussed in a separate manuscript in this issue of *Nature*.<sup>18</sup>

### **Automated clone selection for sequencing**

To maximize the amount of sequence produced while still ensuring the integrity of the clones selected, seed clones were chosen from contigs consisting of at least eight clones by finding the largest clone in each contig that had no more than 2 bands that were not confirmed by neighboring clones. Two fragments in each clone were assumed to be the insert-vector junction fragments as HindIII was used to generate fingerprints from this BAC library constructed using a partial EcoRI restriction digest. Once selected, each seed clone was compared to *in silico* digests of all of the human clones in GenBank to identify any

clones that may overlap previously generated sequence. If clones were found to excessively overlap an existing sequence by a Sulston score of  $<3 \times 10^{-6}$ , or if the overlap suggested localization to a chromosome other than those specifically targeted, then the seed clone was discarded. Based on the estimated sizes of the chromosomes and the number of seed clones selected from contigs localized to each chromosome, it was determined whether each chromosome had been seeded to a point of saturation at approximately one contig per 500 kb. When this point was reached, seed clones were no longer chosen from that chromosome and these nucleation points were extended to avoid excessive gaps when completing the sequence map. Since the fingerprint database was growing weekly and all data were not yet incorporated into existing contigs, a set of tools was developed that allowed extension from seed clones by identifying overlapping fingerprints in the rest of the database independent of contig assembly. Once manually ordered contigs became available, the automated clone selection tools attempted to walk in a minimal tiling path, except in cases where a clone chosen for a walk overlapped heavily with another sequenced clone. In this case, a bridging clone with significant overlap to both neighbors was chosen to complete the coverage. For choosing gap-spanning clones and for choosing minimal tiling paths, BAC end data, when available and consistent, were used in conjunction with fingerprint data. To complete the working draft, the end of each contig was extended where possible using fingerprint analysis, even if the degree of overlap was higher than our allowed degree of overlap for automated clone selection (25%). Clones were automatically selected from the ends of contigs when at least 30 kb of new sequence data would likely be obtained.

Throughout the clone selection process for the working draft, a central registry was maintained at the NCBI to track clones being sequenced by all groups in the consortium. Incorrect clone names and the depth of clone coverage for any given region limited the usefulness of this registry but it did provide a means for examining maps for conflicts. In addition, the first set of 96 sequence reads for any clone was used to check for overlap with known sequences. Initially, a clone was discarded at this stage if greater than 25% of sequences matched another clone. This overlap cutoff was increased as the working draft accumulated to allow for gap-spanning clones to enter the pipeline. Even with all these efforts, redundancy within the working draft sequence is higher than the  $<20\%$  that was targeted; however, this redundancy in fact may have benefits for analysis of sequence variation within the genome.

### **Mapping of accessioned sequences.**

For the placement of each accessioned sequence, *in silico* digests of each sequence were used to query the FPC human BAC database. If the FPC database fingerprint with the highest restriction fragment similarity to the *in silico* digest corresponded to the clone referenced in the associated GenBank record, then the accessioned sequence was mapped to this position; however, if it corresponded to a different clone than that referenced in the GenBank record, additional information was used to position this sequence as accurately as possible. In many cases, the highest matching fingerprint was associated with a clone overlapping with, or completely redundant to, the referenced clone. This could occur when working draft sequence was used as all restriction fragments may not be adequately represented in the assembled sequence and thus would be absent from the *in silico* digest. In these cases, the accessioned sequence was positioned with respect to the GenBank referenced clone. In a third scenario, the highest matching fingerprint was associated with a clone mapping independently of that referenced in the GenBank record. If the *in silico* digest match to the fingerprint was sufficiently stringent (Sulston score  $<1 \times 10^{-9}$ ) then the accessioned sequence was mapped to the matching position. If the match was less stringent, then two pieces of additional information were required to confirm this placement. Additional information considered was the alignment of BAC end sequences to the accessioned sequence. If the majority of these BAC end alignments represented clones overlapping with the matching clone, then the accessioned sequence was placed in that position on the map. At the time of last analysis, there were 302,351 BAC end sequences that could be aligned to 32,235 HTGS human sequences. Of these alignments, 83,219 of them were from clones contained in the WG BAC map and confirmed or contributed to accessioned sequence placement. Additional evidence for placement of the sequence came from examining believed sequence overlaps with neighboring sequenced clones. The believed overlaps were generated by an all-against-all sequence comparison of all accessioned sequences (G. Schuler and J. Kent). The believed overlaps were ranked as either strong (includes additional supporting data such as paired BAC end linkages) or weak (sequence overlap alone) and were weighted accordingly in this consideration.

1. Olson, M., Hood, L., Cantor, C. & Botstein, D. A common language for physical mapping of the human genome. *Science* **245**, 1434-5 (1989).

2. Burke, D. T., Carle, G. F. & Olson, M. V. Cloning of large segments of exogenous DNA into yeast by means of artificial chromosome vectors. *Science* **236**, 806-12 (1987).

3. Bouffard, G. G. et al. A physical map of human chromosome 7: an integrated YAC contig map with average STS spacing of 79 kb. *Genome Res* **7**, 673-92 (1997).
4. Ross, M. T., LaBrie, S., McPherson, J. & Stanton, V. P. in *Current Protocols in Human Genetics* (eds. Dracopoli, N. C. et al.) 5.6.1-5.6.5 (John Wiley and Sons, New York, 1999).
5. Soderlund, C., Longden, I. & Mott, R. FPC: a system for building contigs from restriction fingerprinted clones. *Comput Appl Biosci* **13**, 523-35 (1997).
6. Soderlund, C., Humphray, S., Dunham, A. & French, L. Contigs Built with Fingerprints, Markers, and FPC V4.7. *Genome Res* **10**, 1772-1787 (2000).
7. Marra, M. et al. A map for sequence analysis of the Arabidopsis thaliana genome. *Nat Genet* **22**, 265-70 (1999).
8. Sulston, J. et al. Software for genome mapping by fingerprinting techniques. *Comput Appl Biosci* **4**, 125-32 (1988).
9. Dib, C. et al. A comprehensive genetic map of the human genome based on 5,264 microsatellites. *Nature* **380**, 152-4 (1996).
10. Schuler, G. D. et al. A gene map of the human genome. *Science* **274**, 540-6 (1996).
11. Deloukas, P. et al. A physical map of 30,000 human genes. *Science* **282**, 744-6 (1998).
12. The International Radiation Hybrid Mapping Consortium, A new gene map of the human genome: GeneMap'99. (cited October 2000) <http://www.ncbi.nlm.nih.gov/genemap/> (1999).
13. The International Human Genome Sequencing Consortium. Initial sequencing and analysis of the human genome. *Nature*, (this issue).
14. Zhao, S. et al. Human BAC ends quality assessment and sequence analyses. *Genomics* **63**, 321-32 (2000).



15. Schuler, G. D. Electronic PCR: bridging the gap between genome mapping and genome sequencing. *Trends Biotechnol* **16**, 456-9 (1998).
16. Broman, K. W., Murray, J. C., Sheffield, V. C., White, R. L. & Weber, J. L. Comprehensive human genetic maps: individual and sex-specific variation in recombination. *Am J Hum Genet* **63**, 861-9 (1998).
17. Bruls, T. et al. A physical map of human chromosome 14. *Nature* (this issue).
18. Reithman, H. C. et al. Integration of telomeric DNA sequences with the human working draft sequence. *Nature* (this issue).

## **The International Human Genome Mapping Consortium**

Full author list

### **[Washington University School of Medicine, Genome Sequencing Center:](#)**

**<http://genome.wustl.edu>**

John D. McPherson<sup>1</sup>, Marco Marra<sup>1\*</sup>, LaDeana Hillier<sup>1</sup>, Robert H. Waterston<sup>1</sup>, Asif Chinwalla<sup>1</sup>, John Wallis<sup>1</sup>, Mandeep Sekhon<sup>1</sup>, Kristine Wylie<sup>1</sup>, Elaine R. Mardis<sup>1</sup>, Richard K. Wilson<sup>1</sup>, Robert Fulton<sup>1</sup>, Tamara A. Kucaba<sup>1</sup>, Caryn Wagner-McPherson<sup>1</sup>, William B. Barbazuk<sup>1</sup>, Amanda Aabbott<sup>1</sup>, Johar Ali<sup>1</sup>, Stephanie Andrews<sup>1</sup>, Larisa Belaygorod<sup>1</sup>, Gary Bemis<sup>1</sup>, Amy Berghoff<sup>1</sup>, Merry Brook<sup>1</sup>, Marco Cardenas<sup>1</sup>, Jason Carter<sup>1</sup>, Jim Cloud<sup>1</sup>, Marc Cotton<sup>1</sup>, Jye'mon Crockett<sup>1</sup>, Kristy Drone<sup>1</sup>, Feiyu Du<sup>1</sup>, Jen Edwards<sup>1</sup>, Jackie Fedele<sup>1</sup>, Dan Fischer<sup>1</sup>, Nat Florence<sup>1</sup>, Catherine Franklin<sup>1</sup>, Tony Gaige<sup>1</sup>, Diane Gaige<sup>1</sup>, Marilyn Gibbons<sup>1</sup>, Neenu Grewal<sup>1</sup>, Heather Grover<sup>1</sup>, Chris Gund<sup>1</sup>, Gwen Harmon<sup>1</sup>, Njata Harvey<sup>1</sup>, Shunfang Hou<sup>1</sup>, Sara Jaeger<sup>1</sup>, Corrie Joshu<sup>1</sup>, Sara Kohlberg<sup>1</sup>, Colin Kremitzki<sup>1</sup>, Dan Layman<sup>1</sup>, Shawn Leonard<sup>1</sup>, Jason Maas<sup>1</sup>, Ken MacDonald<sup>1</sup>, Catherine Marquis-Homeyer<sup>1</sup>, Rachel Maupin<sup>1</sup>, Ryan Mcadow<sup>1</sup>, Cindy McCabe<sup>1</sup>, Rebecca McGrane<sup>1</sup>, Kelly Mead<sup>1</sup>, Richard Morales<sup>1</sup>, Nancy Mudd<sup>1</sup>, Christine Nguyen<sup>1</sup>, Phil Ozersky<sup>1</sup>, Carrie Ragan<sup>1</sup>, Amy Reiley<sup>1</sup>, Kerry Robinson<sup>1</sup>, Ellen Ryan<sup>1</sup>, Samuel Sasso<sup>1</sup>, Debra Scheer<sup>1</sup>, Kelsi Scott<sup>1</sup>, Kelsi Scott<sup>1</sup>, Ryan Seim<sup>1</sup>, Karina Shapiro<sup>1</sup>, Proteon Shelby<sup>1</sup>, Aimee Smith<sup>1</sup>, Tamberlyn Stoneking<sup>1</sup>, Hui Sun<sup>1</sup>, Carrie Sutterer<sup>1</sup>, Elizabeth Sweet<sup>1</sup>, Brenda Theising<sup>1</sup>, Jane Threideh<sup>1</sup>, Rebecca Walker<sup>1</sup>, J. Patrick Woolley<sup>1</sup>, & Martin Yoakum<sup>1</sup>

### **[Wellcome Trust Genome Campus:](#)**

**<http://www.sanger.ac.uk/>**

Simon G. Gregory<sup>2</sup>, Sean J. Humphray<sup>2</sup>, Lisa French<sup>2</sup>, Richard S. Evans<sup>2</sup>, Graeme Bethel<sup>2</sup>, Adam Whittaker<sup>2</sup>, Jane L. Holden<sup>2</sup>, Owen T. McCann<sup>2</sup>, Andrew Dunham<sup>2</sup>, Carol Soderlund<sup>2\*</sup>, Carol E. Scott<sup>2</sup> & David R. Bentley<sup>2</sup>

### **[National Center for Biotechnology Information:](#)**

**<http://www.ncbi.nlm.nih.gov>**

Gregory Schuler<sup>3</sup>, Hsiu-Chuan Chen<sup>3</sup> & Wonhee Jang<sup>3</sup>

### **[National Human Genome Research Insititute:](#)**

**<http://genome.nhgri.nih.gov/chr7/>**

Eric D. Green<sup>4</sup>, Jacquelyn R. Idol<sup>4</sup> & Valerie V. Braden Maduro<sup>4</sup>

### **[Albert Einstein College of Medicine:](#)**

**<http://sequence.aecom.yu.edu/chr12/>**

Kate T. Montgomery<sup>5</sup>, Eunice Lee<sup>5</sup>, Ashley Miller<sup>5</sup>, Suzanne Emerling<sup>5</sup> & Raju Kucherlapati<sup>5</sup>

### **[Baylor College of Medicine, Human Genome Sequencing Center:](#)**

**<http://www.hgsc.bcm.tmc.edu/>**

Richard Gibbs<sup>6</sup>, Steve Scherer<sup>6</sup>, J. Harley Gorrell<sup>6</sup>, Erica Sodergren<sup>6</sup>, Kerstin Clerc-Blankenburg<sup>6</sup>, Paul Tabor<sup>6</sup>, Susan Naylor<sup>7</sup> & Dawn Garcia<sup>7</sup>

### **[Roswell Park Cancer Institute:](#)**

Pieter J. de Jong<sup>8\*</sup>, Joseph J. Catanese<sup>8\*</sup>, Norma Nowak<sup>8</sup> & Kazutoyo Osoegawa<sup>8\*</sup>

**Multimegabase Sequencing Center:**

Shizhen Qin<sup>9</sup>, Lee Rowen<sup>9</sup>, Anuradha Madan<sup>9</sup>, Monica Dors<sup>9</sup> & Leroy Hood<sup>9</sup>

**Fred Hutchinson Cancer Research Institute:**

Barbara Trask<sup>10</sup>, Cynthia Friedman<sup>10</sup> & Hillary Massa<sup>10</sup> The Children's Hospital of Philadelphia: Vivian G. Cheung<sup>11</sup>, Ilan R. Kirsch<sup>12</sup>, Thomas Reid<sup>12</sup> & Raluca Yonescu<sup>12</sup>

**Genoscope:**

Jean Weissenbach<sup>13</sup>, Thomas Bruls<sup>13</sup> & Roland Heilig<sup>13</sup>

**US DOE Joint Genome Institute:**

[http://www.jgi.doe.gov/JGI\\_home.html](http://www.jgi.doe.gov/JGI_home.html)

Elbert Branscomb<sup>14</sup>, Anne Olsen<sup>14</sup>, Norman Doggett<sup>14</sup>, Jan-Fang Cheng<sup>14</sup> & Trevor Hawkins<sup>14</sup>

**Stanford Human Genome Center and Department of Genetics:**

<http://www-shgc.stanford.edu/>

Richard M. Myers<sup>15</sup>, Jin Shang<sup>15</sup>, Lucia Ramirez<sup>15</sup>, Jeremy Schmutz<sup>15</sup>, Olivia Velasquez<sup>15</sup>, Kami Dixon<sup>15</sup>, Nancy E. Stone<sup>15</sup> & David R. Cox<sup>15</sup>

**University of California, Santa Cruz:**

<http://genome.ucsc.edu>

David Haussler<sup>16,17</sup>, W. James Kent<sup>18</sup>, Terrence Furey<sup>17</sup>, Sanja Rogic<sup>17</sup> & Scot Kennedy<sup>19</sup>

**Department of Genome Analysis, Institute of Molecular Biotechnology:**

<http://genome.imb-jena.de/printHum.html>

André Rosenthal<sup>21</sup>, Gaiping Wen<sup>21</sup>, Markus Schilhabel<sup>21</sup>, Gernot Gloeckner<sup>21</sup>, Gerald Nyakatura<sup>21\*</sup>, Reiner Siebert<sup>22</sup> & Brigitte Schlegelberger<sup>22</sup>

**Departments of Human Genetics and Pediatrics, University of California:**

Julie Korenberg<sup>23</sup> & Xiao-Ning Chen<sup>23</sup>

**RIKEN Genomic Sciences Center:**

<http://hgp.gsc.riken.go.jp/>

Asao Fujiyama<sup>24</sup>, Masahira Hattori<sup>24</sup>, Atsushi Toyoda<sup>24</sup>, Tetsushi Yad<sup>24</sup>,

Hong-Seok Park<sup>24</sup> & Yoshiyuki Sakaki<sup>24</sup>

**Department of Molecular Biology, Keio University School of Medicine:**

<http://www.dmb.med.keio.ac.jp>

Nobuyoshi Shimizu<sup>25</sup>, Shuichi Asakawa<sup>25</sup>, Kazuhiko Kawasaki<sup>25</sup>, Takashi Sasaki<sup>25</sup>, Ai Shintani<sup>25</sup>, Atsushi Shimizu<sup>25</sup>, Kazunori Shibuya<sup>25</sup>, Jun Kudoh<sup>25</sup> & Shinsei Minoshima<sup>25</sup>

**Max-Planck-Institute for Molecular Genetics:**

<http://seq.molgen.mpg.de/hgs/>

Juliane Ramser<sup>26</sup>, Peter Seranski<sup>26,27</sup>, Celine Hoff<sup>26,27</sup>, Uwe Radelo<sup>26</sup>, Ralf Sudbrak<sup>26</sup>, Annemarie Poustka<sup>26,27</sup>, Richard Reinhardt<sup>26</sup> & Hans Lehrach<sup>26</sup>

1. Washington University School of Medicine, Genome Sequencing Center, Department of Genetics, 4444 Forest Park Boulevard, St. Louis, Missouri 63108, USA

2. Wellcome Trust Genome Campus, Hinxton, Cambridge, CB10 1SA, UK

3. National Center for Biotechnology Information, National Institutes of Health, Bethesda, Maryland 20894, USA

4. National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland 20892, USA

5. Department of Molecular Genetics, Albert Einstein College of Medicine, 1300 Morris Park Avenue, Bronx, New York 10461, USA

6. Baylor College of Medicine, Human Genome Sequencing Center, Houston, Texas, USA  
<http://www.hgsc.bcm.tmc.edu/>

7. University of Texas, San Antonio, Texas, USA

8. Roswell Park Cancer Institute, Buffalo, New York 14263, USA

9. Multimegabase Sequencing Center, Institute for Systems Biology, Seattle, Washington 98105, USA

10. Division of Human Biology, Fred Hutchinson Cancer Research Institute, Seattle, Washington 98109, USA

11. Department of Pediatrics, The Children's Hospital of Philadelphia, University of Pennsylvania, Philadelphia, Pennsylvania 19104, USA

12. Genetics Department, Medicine Branch, National Cancer Institute, Washington DC, USA

13. Genoscope, Centre National de Séquencage, 2 Rue Gaston Crémieux, CP 5706, 91057 Evry, France

14. US DOE Joint Genome Institute, Walnut Creek, California, USA

15. Stanford Human Genome Center and Department of Genetics, Stanford University School of Medicine, Stanford, California 94305, USA

16. Howard Hughes Medical Institute, University of California, Santa Cruz, Santa Cruz, California 95064, USA

17. Department of Computer Science, University of California, Santa Cruz, Santa Cruz, California 95064, USA

18. Department of Biology, University of California, Santa Cruz, Santa Cruz, California 95064, USA

19. Department of Mathematics, University of California, Santa Cruz, Santa Cruz, California 95064, USA

20. British Columbia Cancer Research Centre, 600 West 10th Avenue, Room 3427, Vancouver, British Columbia V5Z 4E6, Canada

21. Dept. of Genome Analysis, Institute of Molecular Biotechnology, Beutenbergstrasse 11, D-07745 JENA, Germany

22. Institute of Human Genetics, University of Kiel, Germany

23. Departments of Human Genetics and Pediatrics, University of California, Los Angeles, California, USA

24. RIKEN Genomic Sciences Center, 1-7-22 Suehiro-cho, Tsurumi-ku, Yokohama 230-0045, Japan

25. Department of Molecular Biology, Keio University School of Medicine, 35 Shinanomachi Shinjuku-ku, Tokyo 160-8582, Japan

26. Max-Planck-Institute for Molecular Genetics, Ihnestrasse 73, D-14195, Berlin, Germany

27. Abteilung Molekulare Genomanalyse, Deutsches Krebsforschungszentrum, Im Neuenheimer Feld 280, 69120 Heidelberg, Germany

\*Present addresses: British Columbia Cancer Research Centre, 600 West 10th Avenue, Room 3427, Vancouver, British Columbia V5Z 4E6, Canada (M.M.); Clemson University Genome Institute, 100 Jordan Hall, Clemson University, Clemson, South Carolina 29634-5727, USA (C.S.); Children's Hospital Oakland Research Institute, BACPAC Resources, Oakland, California 94609, USA and Pfizer Global Research & Development, Alameda Laboratories, Alameda, California 94502 USA (P.J.d.J., J.J.C., K.O.); MWG-Biotech AG, Ebersberg, Germany (G.N.).

Figure 1: Assays of fingerprint reproducibility and fingerprint accuracy.

Figure 1a: Verification of fingerprint reproducibility.

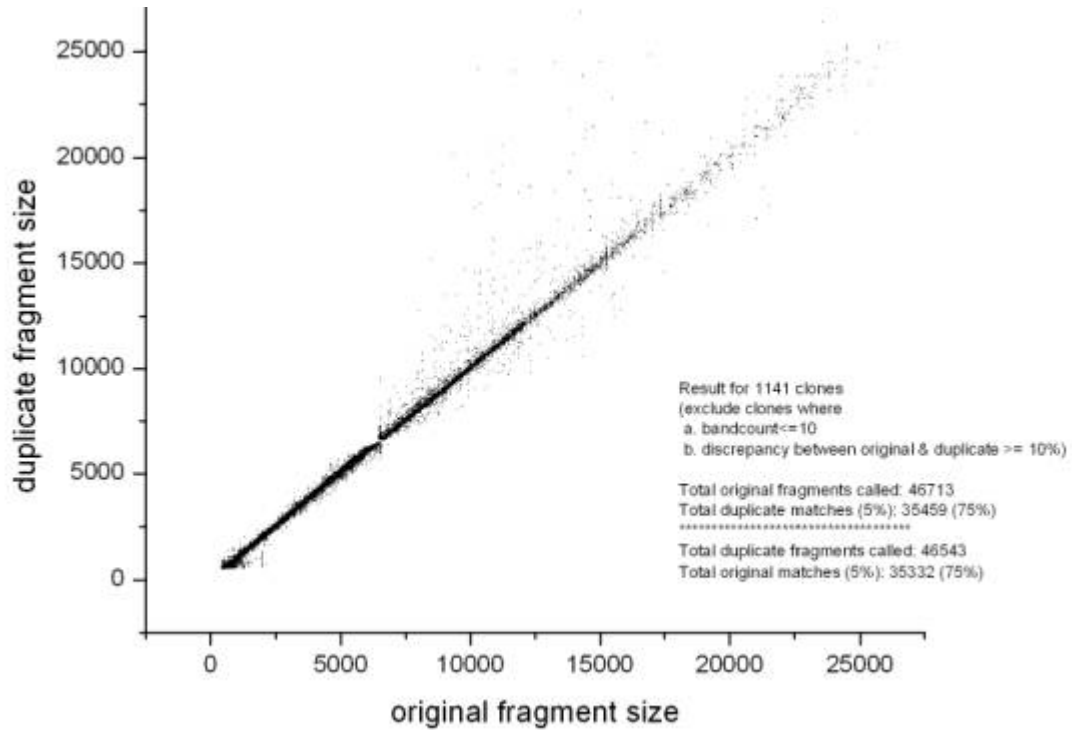
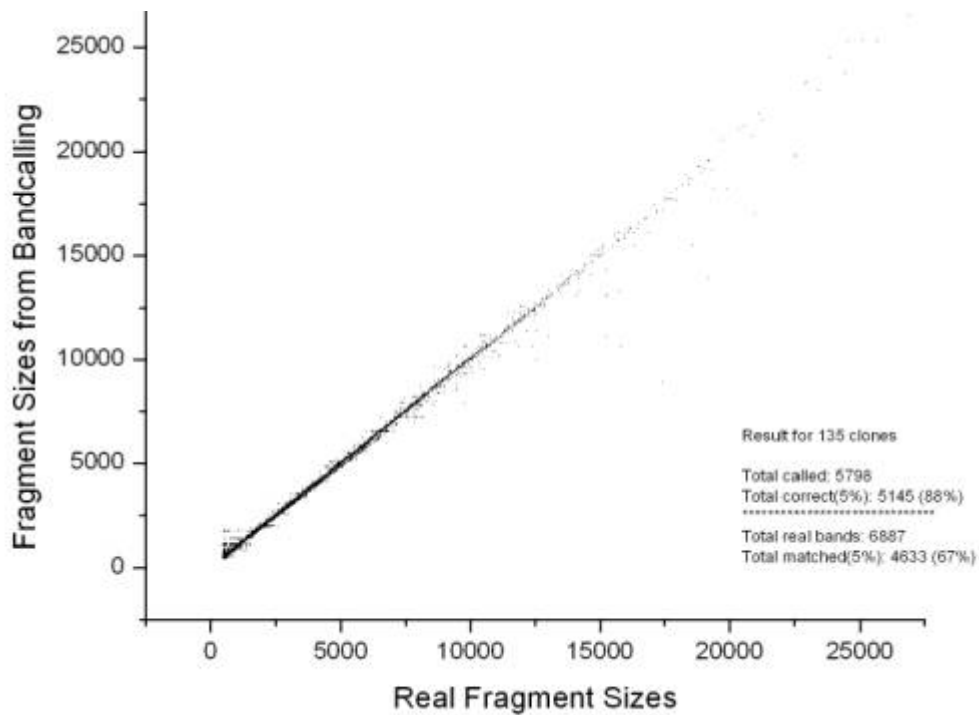


Figure 1b: Verification of fingerprint accuracy.



Shown in Figure 1a is a comparison of clone fingerprints, generated in duplicate, in which each point represents a restriction fragment compared between duplicates. Shown in Figure 1b are fingerprint-derived restriction fragment sizes compared to computer-derived “*in silico*” virtual digests predicted from DNA sequence data. Restriction fragments were considered identical if their sizes (base pairs) were within 2% of each other.

**Table 1: Fingerprint characterization of the major clone classes used.\***

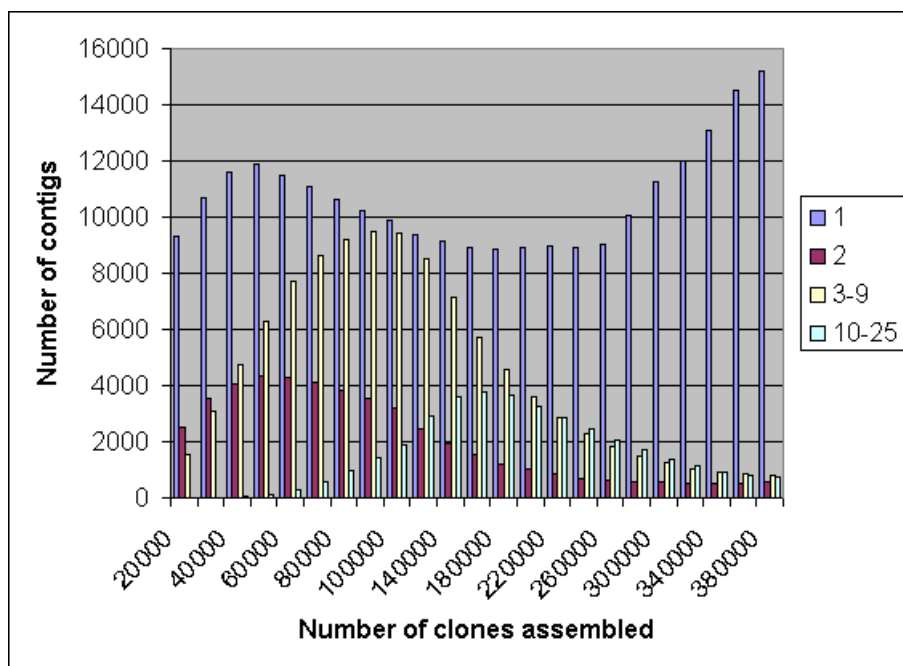
Library	Clones	Ave. insert (kb)	Ave. number of bands
		(std. dev.; median)	(std. dev.; median)
RPCI-11	265,619	167.7 (21.1; 167.4)	41.8 (7.9; 42)
RPCI-13	50,202	151.5 (27.9; 156.6)	38.6 (9.5; 40)
CTC-C/D1	49,349	123.4(29.1; 120.5)	30.9 (8.6; 30)

\*Only clones with an estimated insert size greater than 50 kb and less than 300 kb were included to minimize the effect of fingerprinting artifacts.

**Table 2: Optimal FPC parameters used for fingerprint assembly.**

FPC Parameter	Value
tolerance	7
cutoff	3x10 <sup>-12</sup>
automated bury	0.10
use CpM	no
min bands	3
best of	10

**Figure 2: Incremental assemblies of the BAC fingerprints.**



At the various points indicated, the database of BAC fingerprints was assembled into contigs using the FPC software suite. Assembly parameters are as in Table 2. The figure key indicates the range of the number of clones per contigs in each group considered (contigs containing 10-15 clones, 3-9 clones, two clones and remaining clones with no match to another clone or “singletons”).

**Table 3: Status of FPC database after automated assembly and manual editing.**

	Automated assembly	Manually edited database
Date	December, 1999	September, 2000
BAC clones in FPC	283,287	372,264
Number of assemblies	7,133	1,447
Clones in assemblies	264,555	295,828
Number of singletons*	11,022	76,436
Assemblies containing:		
>25 clones	2,652	912
10-25 clones	1,720	260
3 – 9 clones	1,886	204
2 clones	875	71

\* clones not incorporated into any assembly; see text.