

Supplemental material

Materials and methods

Preparation of BAC library CHORI-507

The CHORI-507 BAC library was prepared as described [Osoegawa et al. 1998]. Briefly, high-molecular-weight DNA from WAV 17 cells was prepared in agarose using DNA Plug Kits (BioRad). The DNA was partially digested with MboI (New England Biolabs) to generate fragments between 100 kb and 400 kb, which were size-fractionated by pulsed-field electrophoresis [Osoegawa et al. 1998]. DNA fractions in the 150 to 250 kb range were sliced from the gel and DNA was electro-eluted from the gel matrix. Vector DNA (pTARBAC1.3) was prepared by treatment with a combination of restriction enzymes (ApaI and BamHI) and alkaline phosphatase and then purified by gel-electrophoresis. *E. coli* DH10B T1 phage-resistant cells (Invitrogen) were transformed by electroporation with ligation products from vector and genomic DNA. Prior to full-scale library production, 379 random clones were picked for DNA size analysis. BAC DNA, prepared using an Autogen960 robot, was digested with NotI (New England Biolabs), and analyzed by pulsed-field electrophoresis in a 1% agarose gel. The library clone size distribution is given in Supplemental Fig. 2, and the average insert size was estimated to be 141 kb. There is a small fraction (1 out of 379) of non-recombinant clones. Following this quality assessment, 221 184 BAC clones (approximately 10-fold coverage) were arrayed into 576 384-well microtiter plates. Duplicate sets of twelve high-density colony filters were created for hybridization screening of the library with various human-specific repeat probes.

BAC end sequencing

Positive clones from CHORI-507 were end sequenced using the following protocol. BAC clones were inoculated into 1.5ml of 2xTY containing 12.5µg/ml chloramphenicol in 2ml 96 well growth boxes (Costar) and incubated for 22 hours at 37°C with shaking at 320rpm. Plates were centrifuged at 4000rpm for 3 min to obtain pellets, the supernatant was discarded and the cells were re-suspended in 100µl of GTE+RNaseA, 100µl of NaOH/SDS was added and mixed before adding 100µl of 3M KOAc and further mixing. After filtration through two Millipore plates (MADV6550/MANUBA50) on a vacuum manifold, the resultant material was dried and re-suspended in 35µl of 10mM Tris pH8.0. The resulting DNA was sequenced with BigDye Terminator Ready Mix v3.0 and BigDye Bufferx5 (Applied Biosystems) using T7 and SP6 primers. The sequenced products were cleaned up by washing with 5µl 3M NaOAc and 125µl 96% ethanol. The DNA was precipitated by centrifugation at 4000rpm for

15 mins, then washed with a further 100 μ l 70% ethanol and the centrifugation repeated. The products were re-suspended in 10 μ l of 0.1M EDTA, pH 7.4 and loaded on to ABI3700 capillary sequencers.

BAC clones were fingerprinted using HindIII restriction enzyme fingerprinting [Humphray et al. 2001; Schein et al. 2004]. BAC clones were directly cultured in 170 μ l 2xTY in 96 well plates. After overnight growth the BAC DNA was extracted by alkaline lysis, and digested with HindIII. Following electrophoresis on 121 lane 1% agarose gels the data are collected using a Typhoon 8600 fluorimager, raw images were entered into the fingerprint database using IMAGE software (<http://www.sanger.ac.uk/Software/Image>). The output of normalized band values, sizes and gel traces were analyzed in FPC [Soderlund et al. 2000], which bins and orders clones on the basis of shared bands. Sequence tilepaths were identified following inspection of assembled contigs. End sequences were aligned to the human and mouse genomes using BLAST and ssaha2 (www.ensembl.org).

BAC sequencing

For the shotgun phase [Bankier et al. 1987], pUC plasmids with inserts of mostly 1.4-2 kb were sequenced from both ends using the dideoxy chain termination method [Sanger et al. 1977] with different versions of big dye terminator chemistry [Rosenblum et al. 1997]. The resulting sequencing reactions were analyzed on ABI 3730 sequencing machines and the generated data were processed (<http://www.sanger.ac.uk/Software/sequencing/>) and assembly with PHRED [Ewing and Green 1998; Ewing et al. 1998] and PHRAP (<http://www.phrap.org/>) algorithms. For the finishing phase, we used GAP4 [Bonfield et al. 1995] to help assess, edit and select reactions to eliminate ambiguities and close sequence gaps. Sequence gaps were closed by a combination of primer walking, PCR, short/long insert sublibraries [McMurray et al. 1998], oligo screens of such sublibraries and transposon sublibraries.

Gene model copy number

We also performed absolute quantification [Hijri and Sanders 2005] using one assay for gene models 25.1, 28.1, 6.2 and two assays for gene model 15.1. For each gene model, a known amount of plasmid DNA containing a 50mer insert was used as a standard. Two-fold serial dilutions (ranging between 1863 and 3.63 pg) of the purified plasmid DNA were included in the experiment to generate the standard curve. Real-time PCR on the plasmid DNA was performed in six technical replicates (10 dilutions per replicate). The Ct values were then plotted against the log of the initial vector concentration containing the amplicon insert. All replicates offered very similar results producing a standard curve with regression coefficients

(R^2) >0.99. This regression was then used for predictions of all gene models copy number for a given Ct value. In the same experiments, two-fold serial dilutions of genomic DNA (ranging between 3200 and 6.25 pg) were included in each experiment to generate the Ct values for nuclear DNA copy number estimation. The number of copies of each gene model could then be calculated for each sample of target DNA, based on the regressions using the plasmid DNA as standards. From this we calculated the copy number of the each gene model per 3.3pg of gDNA (the amount of DNA per each copy of the entire human genome). As a control, we used herring sperm DNA as a carrier to yield 10 ng of total DNA per sample in a total volume of 40 μ l, irrespective of the dilution of the target DNA to show that the estimation of copy number is not biased by the amount of non-target DNA added to the reaction.

WAV17

The somatic cell hybrid WAV 17 [Slate et al. 1978] was obtained from Coriell (ccr.coriell.org/nigms; GM08854) and cultured according to the recommended protocols. The monoallelic composition of HSA21 in WAV17 was assessed by genotyping with markers D21S11, D21S1270, D21S1435, D21S1411, D21S226 and *IFNAR*.

Screening BAC library CHORI-507

To identify human clones, CHORI-507 was screened with 32 P labeled α -satellite and Alu probes, and total human genomic DNA. This resulted in 2400 (1.09%) positive clones after screening using α -satellite and Alu probes, and 2752 (1.24%) positives with total human genomic DNA. This is consistent with the fact that the HSA21 content of WAV 17 is approximately 1-3% in an aneuploid mouse background. The non-redundant positive clones (3208) were re-arrayed into 384-well plates and gridded on nylon membrane for screening with 21p-specific probes. To identify specific 21p BACs, overgo probes were designed to STSs corresponding to the YAC ends 2E4R, 2E4L, 4E9L, 4E9R, 1B8L, 2C2R, 2C9R from Wang *et al.* [Wang et al. 1999]. A summary of the library screens is given in Supplemental Table 1 and Supplemental Table 2.

FISH

Cell lines were grown in RPMI 1640 with Glutamax I medium (Invitrogen) supplemented with 10% fetal calf serum and a 1% penicillin/streptomycin mix. Cultures were exposed to colcemid (0.1 μ g/mL; Invitrogen) for 1.5 h at 37°C and harvested according to routine cytogenetic protocols. G-banding was performed following standard cytogenetic methods.

The following probes were used in FISH to check the integrity of HAS21 in the WAV 17 cell line: HSA21 paint (Vysis), cosmid c55A10 [Chen et al. 1999], D21Z1 [Maratou et al. 1999], human and mouse Cot-1 DNA (Invitrogen). Interphase nuclei and metaphase spreads were counterstained with DAPI (Vysis) diluted in Vectashield antifade (Vector Labs). BAC and cosmid DNA was labeled with either a Biotin or Digoxigenin-Nick Translation kit (Roche) and detected with anti-dig fluorescein (green) or avidin (red). Cells were viewed under a Nikon fluorescence microscope equipped with appropriate filter combinations. Monochromatic images were captured and superimposed using the Applied Imaging automated imaging system.

Gene predictions

EST sequences from dbEST [Boguski et al. 1993] were aligned to the HSA21p sequence using the program EXONERATE [Slater and Birney 2005]. We considered only spliced alignments with sequence identity of 90% or higher and at least 88% coverage, and discarded those for which there existed a better spliced alignment elsewhere in the genome. Selected EST alignments were divided into four groups according to the properties of the alignment: best, the best alignment for this EST is in HSA21p sequence; pseudo, the best alignment for this EST is in euchromatin, but this is not spliced, indicating a possible processed pseudogene in the euchromatin; paralogue, the coverage of the EST alignment is 100% but the percentage identity is not as high as the best match in euchromatin, hence a potential paralogue; random, there is a better alignment for this EST in an unassembled euchromatin contig. From these alignments we calculated the set of unique exon-exon pairs (75 in total).

Using RepeatMasker (www.repeatmasker.org) without the low complexity filter (option -nolow) we masked the HSA21p sequences and performed in silico gene prediction using the two *ab initio* gene prediction programs, GeneID [Parra et al. 2000] and GENSCAN [Burge and Karlin 1997], and the comparative gene predictor SGP [Parra et al. 2003], in this case using mouse genome sequence (version mm5) as reference. We thus obtained 3 sets of predictions which we merged together into a single set of predicted exon-exon pairs removing duplicates. We further removed those exon-exon pairs that overlap with any of the EST exon-exon pairs calculated before. We obtained a set of 182 unique exon-exon pairs obtained by *ab initio* and comparative gene prediction and not occurring in the EST-predicted set. We classified the 182 exon-exon pairs into those being predicted by all three gene predictors, by two of them and those that are specific to one gene predictor. This classification, the coordinates and the sequences of these predictions are available through the on-line Supplemental material website at <http://genome.imim.es/datasets/hsa21p>.

Exon-exon pairs were assembled into 26 gene models (21pGMs) by aligning them with the genomic sequence and combining exon pairs conforming to the splice-site consensus sequence.

Sequence analysis

For general sequence analysis programs within the EMBOSS package were used [Rice et al. 2000]. BLAST searches were carried out at Ensembl (www.ensembl.org). Segmental duplications were detected using a BLAST-based detection scheme (WGAC) 45 to identify all pairwise similarities representing duplicated regions (≥ 1 kb and $\geq 90\%$ identity) within the finished BAC sequence of chromosome 21 p arm. The 21p BAC sequences were manually joined in their natural order and compared to all other chromosomes in the NCBI genome assembly (hg17 and 18). Satellite repeats were detected using RepeatMasker (version: 2002/05/15) on sensitive settings (www.repeatmasker.org) [Smit et al. 1996]. The program PARASIGHT (Bailey, unpublished) was used to generate images of pairwise alignments.

RT-PCR

Gene models were tested in 24 human cDNAs (brain, heart, kidney, spleen, liver, colon, small intestine, muscle, lung, stomach, testis, placenta, skin, PBLs, bone marrow, fetal brain, fetal liver, fetal kidney, fetal heart, fetal lung, thymus, pancreas, mammary glands, prostate). All amplicons spanned introns. Each cDNA was mixed with JumpStart REDTaq ReadyMix (Sigma) and 4 ng/ μ l primers (Sigma-Genosys) using a BioMek 2000 robot (Beckman) as described and modified [Reymond et al. 2002]. The first 10 cycles of PCR amplification were performed with a touchdown annealing temperature decreasing from 60 to 50°C; annealing temperature of the next 30 cycles was carried out at 50°C. Amplimers were separated on Ready-to-Run precast gels (Pharmacia) and sequenced.

Gene copy number

Relative DNA copy number was determined as described [Lyle et al. 2006] with some modifications. Genomic DNA was isolated from lymphoblastoid cell lines using DNeasy Tissue kits (Qiagen) according to the manufacturer's protocol. DNAs from a total of 13 individuals were used: 8 unrelated individuals, and mother and son from CEPH families (Coriell), and 3 unrelated caucasian individuals from the University of Geneva DNA bank. Oligonucleotides were designed using PrimerExpress (Applied Biosystems) with default parameters. For each gene prediction a minimum of two assays were designed. Seven assays were used for normalization of input DNA. Control assays were designed in regions that are present in one (*AKAP1*) and two copies (*HBG2*) per haploid genome. Assay sequences are

available in Supplemental Table 10.

Quantitative PCR reactions were performed using PowerSYBR Green PCR master Mix (Applied Biosystems) and each DNA sample were amplified in three replicates. Fluorescence data were collected using the ABI Prism 7900 Sequence Detection System (Applied Biosystems) and data analysis was performed with Qbase (Hellemans et al., unpublished).

Polymorphisms

RNA was obtained from human brain, testis, heart, spleen, thymus and lung (Ambion). RNA from somatic cell hybrids carrying human chromosome 13, 14, 15, 21, and 22 (GM10898, GM10479, GM11418, GM08854, GM10888; Coriell) was extracted with RNeasy mini kit (Qiagen). cDNA was synthesized using SuperScript II (Invitrogen) and each predicted gene model was tested with the appropriate primer set (Supplemental Table 10) designed to span the gene model. Platinum *Pfx* DNA Polymerase (Invitrogen) was used to perform the PCR. After purification with QIAquick PCR Purification Kit (Qiagen), PCR products were cloned into pCR4Blunt-TOPO vector (Invitrogen). Purification of clones was performed using Plasmid Miniprep⁹⁶ kits (Millipore) and sequencing reaction was carried out with ABI Prism BigDye Terminator Cycle Sequencing Ready Reaction Kit (Applied Biosystems). Lasergene v6.0 (DNASar) was used to quality check and align sequences. Alignments were edited and clustered using BioEdit (<http://www.mbio.ncsu.edu/BioEdit/bioedit.html>) and GeneDoc (<http://www.psc.edu/biomed/genedoc>).

References

- Bankier, A.T., Weston, K.M., Barrell, B.G. 1987. Random cloning and sequencing by the M13/dideoxynucleotide chain termination method. *Methods Enzymol* **155**: 51-93.
- Boguski, M.S., Lowe, T.M.J., Tolstoshev, C.M. 1993. dbEST - database for "expressed sequence tags". *Nature Genetics* **4**: 332-333.
- Bonfield, J.K., Smith, K., Staden, R. 1995. A new DNA sequence assembly program. *Nucleic Acids Res* **23**: 4992-4999.
- Burge, C., Karlin, S. 1997. Prediction of complete gene structures in human genomic DNA. *J Mol Biol* **268**: 78-94.
- Chen, H., Rossier, C., Morris, M.A., Scott, H.S., Gos, A., Bairoch, A., Antonarakis, S.E. 1999. A testis-specific gene, TPTE, encodes a putative transmembrane tyrosine phosphatase and maps to the pericentromeric region of human chromosomes 21 and 13, and to chromosomes 15, 22, and Y. *Hum Genet* **105**: 399-409.
- Ewing, B., Green, P. 1998. Base-calling of automated sequencer traces using phred. II. Error probabilities. *Genome Res* **8**: 186-194.
- Ewing, B., Hillier, L., Wendl, M.C., Green, P. 1998. Base-calling of automated sequencer traces using phred. I. Accuracy assessment. *Genome Res* **8**: 175-185.
- Hijri, M., Sanders, I.R. 2005. Low gene copy number shows that arbuscular mycorrhizal fungi inherit genetically different nuclei. *Nature* **433**: 160-163.
- Humphray, S.J., Knaggs, S.J., Ragoussis, I. 2001. Contiguation of bacterial clones. *Methods*

- Mol Biol* **175**: 69-108.
- Lyle, R., Radhakrishna, U., Blouin, J.L., *et al.* 2006. Split-hand/split-foot malformation 3 (SHFM3) at 10q24, development of rapid diagnostic methods and gene expression from the region. *Am J Med Genet A* **140**: 1384-1395.
- Maratou, K., Siddique, Y., Kessling, A.M., Davies, G.E. 1999. Novel methodology for the detection of chromosome 21-specific alpha-satellite DNA sequences. *Genomics* **57**: 429-432.
- McMurray, A.A., Sulston, J.E., Quail, M.A. 1998. Short-insert libraries as a method of problem solving in genome sequencing. *Genome Res* **8**: 562-566.
- Osoegawa, K., Woon, P.Y., Zhao, B., Frengen, E., Tateno, M., Catanese, J.J., de Jong, P.J. 1998. An improved approach for construction of bacterial artificial chromosome libraries. *Genomics* **52**: 1-8.
- Parra, G., Agarwal, P., Abril, J.F., Wiehe, T., Fickett, J.W., Guigo, R. 2003. Comparative gene prediction in human and mouse. *Genome Res* **13**: 108-117.
- Parra, G., Blanco, E., Guigo, R. 2000. GeneID in Drosophila. *Genome Res* **10**: 511-515.
- Reymond, A., Marigo, V., Yayiaoglu, M.B., *et al.* 2002. Human chromosome 21 gene expression atlas in the mouse. *Nature* **420**: 582-586.
- Rice, P., Longden, I., Bleasby, A. 2000. EMBOS: the European Molecular Biology Open Software Suite. *Trends Genet* **16**: 276-277.
- Rosenblum, B.B., Lee, L.G., Spurgeon, S.L., Khan, S.H., Menchen, S.M., Heiner, C.R., Chen, S.M. 1997. New dye-labeled terminators for improved DNA sequencing patterns. *Nucleic Acids Res* **25**: 4500-4504.
- Sanger, F., Nicklen, S., Coulson, A.R. 1977. DNA sequencing with chain-terminating inhibitors. *Proceedings of the National Academy of Science USA* **74**: 5463-5467.
- Schein, J., Kucaba, T., Sekhon, M., Smailus, D., Waterston, R., Marra, M. 2004. High-throughput BAC fingerprinting. *Methods Mol Biol* **255**: 143-156.
- Slate, D.L., Shulman, L., Lawrence, J.B., Revel, M., Ruddle, F.H. 1978. Presence of human chromosome 21 alone is sufficient for hybrid cell sensitivity to human interferon. *Journal of Virology* **25**: 319-325.
- Slater, G.S., Birney, E. 2005. Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics* **6**: 31.
- Smit, A.F.A., Hubley, R., Green, P. 1996. RepeatMasker Open-3.0.
- Soderlund, C., Humphray, S., Dunham, A., French, L. 2000. Contigs Built with Fingerprints, Markers, and FPC V4.7. *Genome Research* **10**: 1772-11787.
- Wang, S.Y., Cruts, M., Del Favero, J., *et al.* 1999. A high-resolution physical map of human chromosome 21p using yeast artificial chromosomes. *Genome Research* **9**: 1059-1073.

Supplemental figures and tables

Supplemental Fig. 1. Confirmation of HSA21 content of WAV 17 cell line. **(a)** G-banding of a single nucleus identifies that this cell contains four chromosomes 21 (red arrowheads). **(b)** FISH paint with HSA21 probe confirms the identity of the human chromosome in the mouse background. Counting of 100 nuclei showed that each cell contains 3-5 chromosomes 21. **(c)** Hybridization with the cosmid c55A10 confirms the presence of the short arm on HSA21 in the cell line. D21Z1 is a chromosome 21 centromere probe [Maratou et al. 1999]. **(d)** Hybridization with human Cot-1 DNA shows that HSA21 is the only human component of WAV 17 and that no mouse sequences could be detected on this chromosome.

Supplemental Fig. 2. Insert size distribution of CHORI-507 clones.

Supplemental Fig. 3. FISH of a 21p BAC clone to a metaphase spread from normal human diploid lymphocytes. Clone 68N6 hybridizes to 21p and 13p, as well as to the pericentromeric regions of 2 and 9. All clones hybridize to 21p, and to various extents to the other acrocentrics. +, positive hybridization signal; ?, chromosome not identified but in this group.

Supplemental Fig. 4. Dotplots of HSA21p BAC sequences.

Supplemental Fig. 5. Major alternative spliced isoforms of 21pGM6. Spliced exons are represented for spleen, thymus and HSA13 somatic cell hybrid. For each RNA source, the number of different alternatively spliced transcripts detected is given, and the group representing the highest percentage (determined as a fraction of the total number of clones sequenced) is shown. For HSA13 all variants show 2 nt substitutions (compared to the original gene model) represented by two white asterisks.

Supplemental Table 1. Results of CHORI-507 screen for HSA clones.

Supplemental Table 2. Results of CHORI-507 screen for 21p clones. Probes are STS probes derived from 21p YAC end sequences [Wang et al. 1999]. BACs were binned according to probe. BACs selected for sequencing are marked in bold.

Supplemental Table 3. Sequenced 21p BAC clones. Non-redundant length refers to contigs constructed from overlap of clones 71C21/54M18 and 201H5/216K13. *This clone is a human/mouse chimera, sequence length is for the human portion only.

Supplemental Table 4. BAC repeat content. Data for 21p BACs was obtained using RepeatMasker; data for the euchromatic proportion of the genome is from Lander et al.

(2001).

Supplemental Table 5. 21p bases involved in segmental duplication and pairwise alignment. The table summarizes the sequence properties of 21p segmental duplications (>90% identity >1kb) compared to the human genome (hg18), including the fraction of duplicated sequence (non-redundant duplication), the number of pairwise alignments, the average length of the pairwise alignments (kb) and their sequence identity. *21p duplication with random chromosomes.

Supplemental Table 6. Expression of 21p gene models. Method: e, based on 100% EST matches; p, *ab initio* prediction only. ^a, data from two alternatively spliced transcripts. ^b, number of PCR fragments sequenced: gene models with a sequenced fragment, but not marked as positive, indicates that RT-PCR expression was positive in at least one tissue but the fragment did not match 21p exactly.

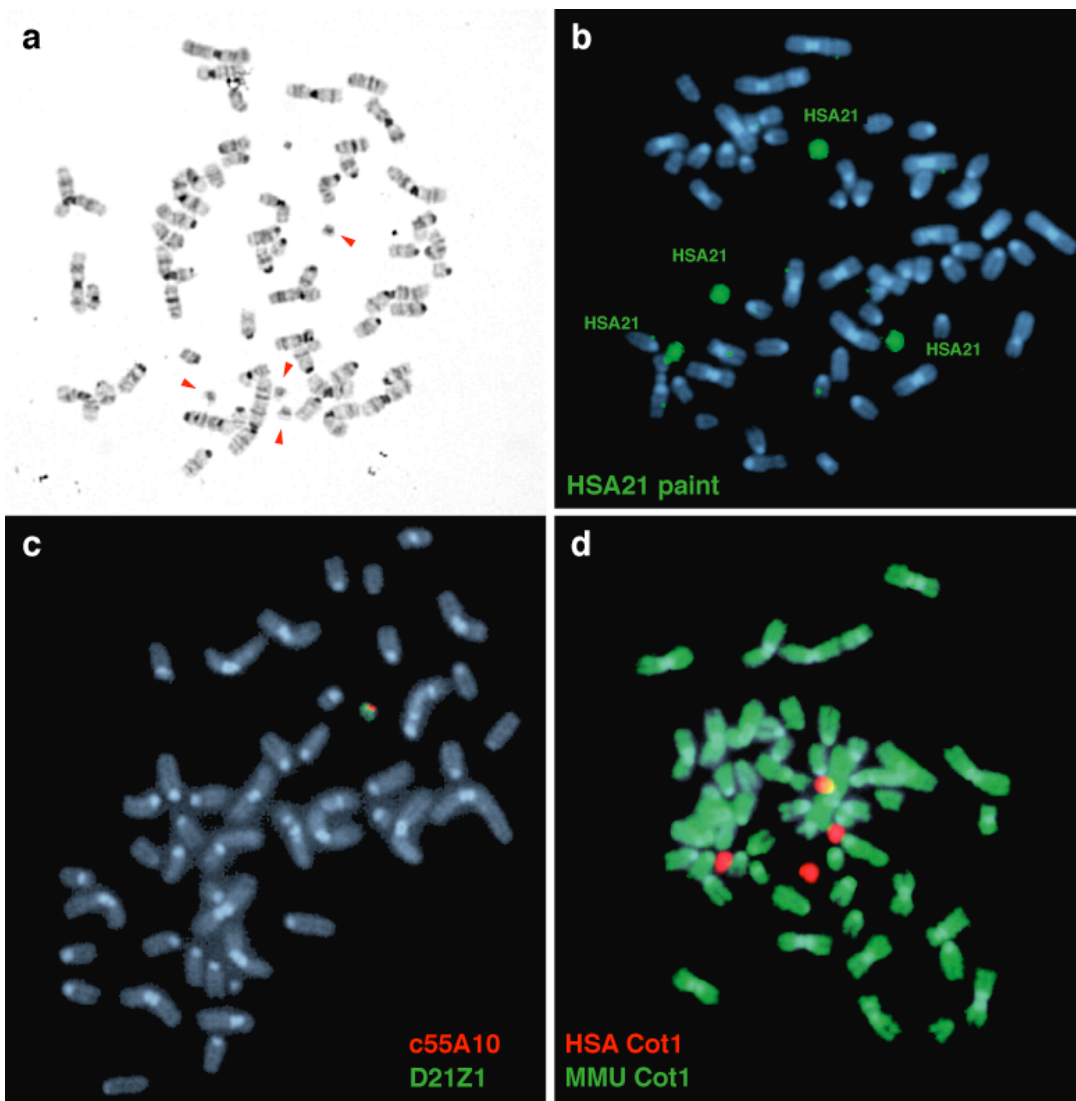
Supplemental Table 7. Homology and conservation of 21p gene models. ORFs of gene models were translated and searched against the databases CCDS Peptides and Peptides using ENSEMBL BLAST. In the case of gene models with no obvious single ORF (i.e. multiple short ORFs), all three forward frames were searched. For human homology, the best match gene is given. For conservation, results are presented as: name of gene with best match, length of match, % ID, e-value. PTR, Pan troglodytes; MMU, Mus musculus; GGA, Gallus gallus; DRE, Danio rerio; DME, Drosophila melanogaster; CEL, Caenorhabditis elegans. §, bp.

Supplemental Table 8. Copy Number results obtained by qRT-PCR. The numbers shown are the copy number of a particular assay for a gene model in each of 13 individuals. The IDs for the 13 individuals used in the analysis are given. ¹Different assays for the same gene model (see text).

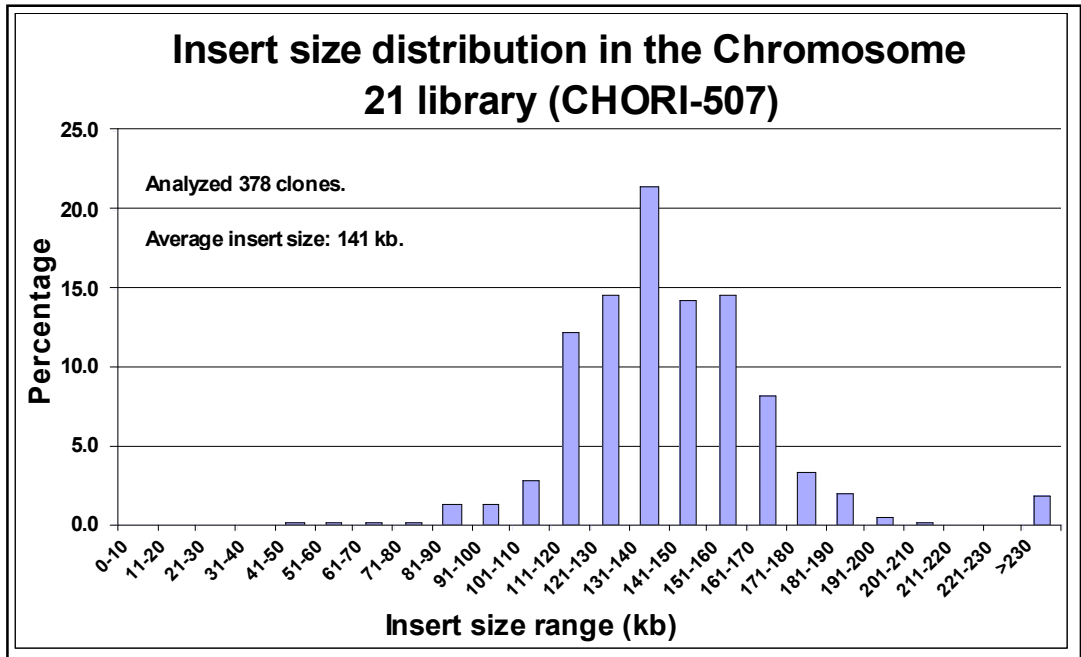
Supplemental Table 9. Polymorphism results. (a) The total number of sequence performed per gene model. (b) Number of sequence relative to different cDNAs derived from different from somatic cell hybrid (HSA) carrying acrocentric chromosomes and tissues. (c) Percentage of sequences that are identical to the original gene model predicted on the HAS21p. (d) The total number of groups obtained after the alignment and clustering of identical sequences. (e) Number of groups containing sequences with nucleotide and splicing variants. In parenthesis the number of sequences related to this group. (f) Number of groups containing splicing isoforms. (g) Percentage of groups containing protein truncation codon (PTC).

Supplemental Table 10. Oligonucleotides for expression, copy number and variation analyses.

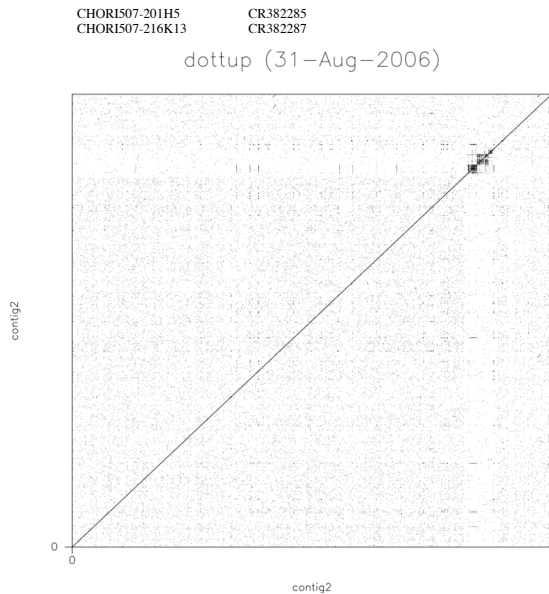
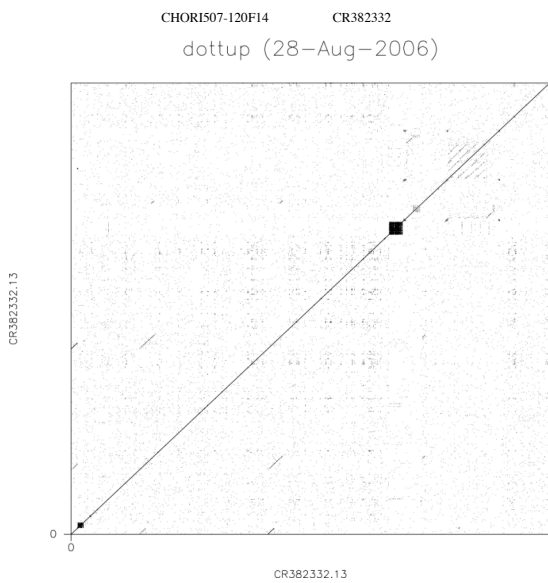
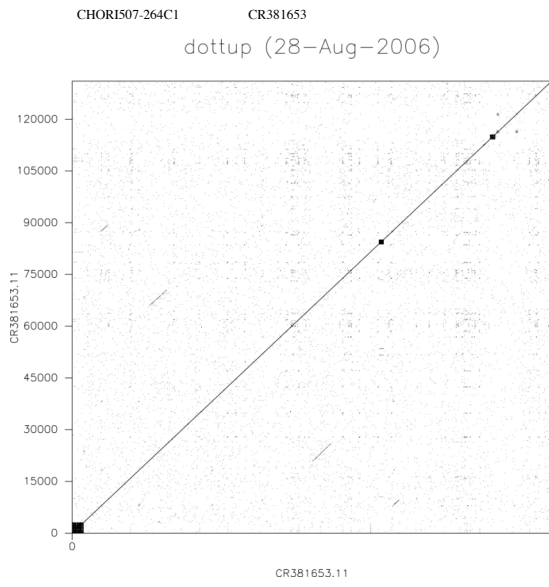
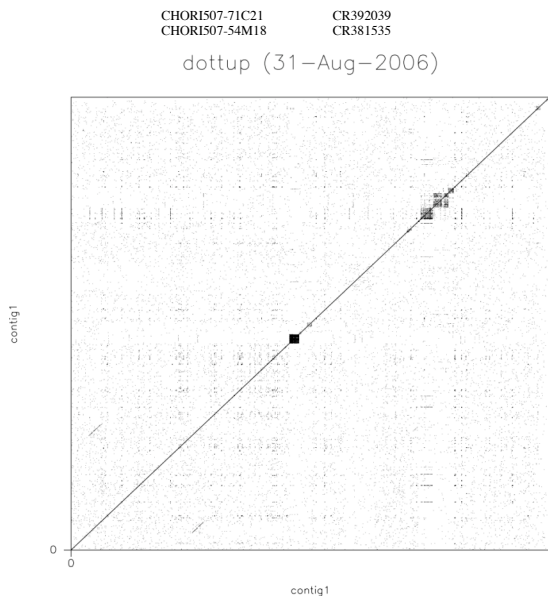
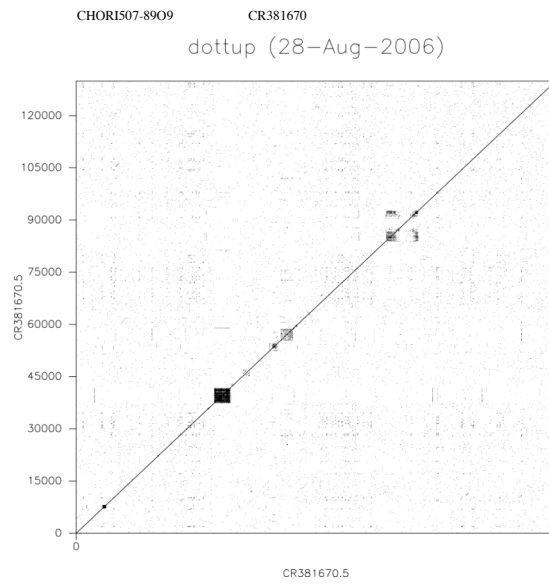
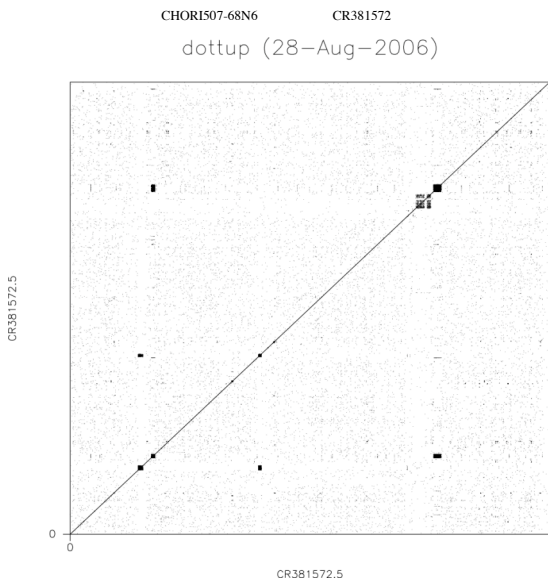
Supplemental Figure 1








Supplemental Figure 2



Supplemental Figure 4



Supplemental Figure 5

Source	No. groups with alternative splicing	% of the total splicing variants	Transcripts
21pGM6	-	-	
Spleen	4	37	
		39	
Thymus	5	63	
HSA13 somatic cell hybrid	5*	80*	

Supplemental Table 1

Filter	Positives		Percentage	
	HSA genomic	Alu/alpha	HSA genomic	Alu/alpha
1	316	194	1.71%	1.05%
2	215	232	1.17%	1.26%
3	241	190	1.31%	1.03%
4	254	128	1.38%	0.69%
5	179	155	0.97%	0.84%
6	278	215	1.51%	1.17%
7	212	219	1.15%	1.19%
8	247	201	1.34%	1.09%
9	189	207	1.03%	1.12%
10	219	262	1.19%	1.42%
11	210	200	1.14%	1.09%
12	192	197	1.04%	1.07%
Total	2752	2400		
Combined unique		3208		
Average	229	200	1.24%	1.09%

Supplemental Table 3

BAC	Acc number	Length	Status	Non-redundant length	% g/c
CHORI507-68N6	CR381572	184355	finished	184355	37.3
CHORI507-89O9	CR381670	129889	finished	129889	41.3
CHORI507-71C21	CR392039	158069	finished	209483	43.3
CHORI507-54M18	CR381535	152296	finished		
CHORI507-264C1	CR381653	131056	finished	131056	39.3
CHORI507-120F14	CR382332	166871	2 unordered pieces	166871	40.8
CHORI507-20IH5	CR382285	178865	finished	281920	36.1
CHORI507-216K13*	CR382287	150002	finished		
Total		1251403		1103574	39.7

Supplemental Table 4

	21p BACs			Euchromatic genome			χ^2 p value
	Number of elements	length occupied	% sequence	Number of elements	length occupied	% sequence	
SINEs	371	93798	8.5	1558000	359600000	13.1	0.00000
LINEs	309	252694	22.9	868000	558800000	20.4	0.10040
LTR elements	253	156953	14.2	443000	227000000	8.3	0.00000
DNA elements	80	29906	2.7	294000	77600000	2.8	0.83780
total length		1103312			2736681887		
bases masked		608555	55.2		1532256900	56.0	

Supplemental Table 5

	Non-redundant duplication (Mb)*			Number of pairwise alignment			Average length of pairwise alignment (kb)			Total aligned base pairs (Mb)			Average percent identity of pairwise alignment (%)		
	Inter	Intra	Both	Inter	Intra	Both	Inter	Intra	Both	Inter	Intra	Both	Inter	Intra	Both
21p	1.042	0.703	1.05	249	29	278	26.8	56.77	29.93	6.673	1.646	8.32	96.41	98.05	96.73
Random*	0.581	0.703	0.892	24	25	49	41.01	64.26	52.87	0.984	1.607	2.591	97.38	98.16	97.86

Supplemental Table 9

	Gene model				
	21pGM6	21pGM9	21pGM15	21pGM25	21pGM28
number of sequenced clones^a	121	48	185	69	31
somatic cell hybrid/tissues (n of sequences)^b	HSA13 (25) spleen (57) thymus (39)	heart (35) testis (13)	HSA13 (32) HSA15 (24) HSA21 (23) HSA22 (28) heart (33) brain (45)	HSA21 (16) HSA22 (18) lung (35)	brain (31)
sequences identical to HSA21p^c	0	0	38%	16%	84%
number of groups^d	25	26	34	35	5
nucleotide variation groups^e	19 (36%)	23 (85%)	29 (60%)	34 (84%)	4 (16%)
splicing isoform groups^f	6 (64%)	3 (15%)	4 (2%)	0	0
groups with PTC^g	25	3	15	1	-

Supplemental Table 10

	Gene/gene model	Assay	F	R
Expression	21pGM1	18481p4.1.a	TCCAAATGGCATCTCTACC	CCTCCAGCTTTCAAAGTGTCT
	21pGM2	18477p3.1.a	CACATGGGGACACAGTCAAG	GCCTCCTCGGGTTTATGCG
	21pGM3	22703p7.5.a	GAAAGTGAGTAATGGCCACAGA	CATTTTATGCCAAGCTTCTTAGTTGC
	21pGM4	18463e6.16.a	TGGCTAGACACCCCTTCTG	AGGCCCTTCAGGGATAACAC
	21pGM5	18455e5.3.a	GGGTGTCAATGACCTTT	GCGGCTTTTGGCCA
	21pGM6	30204e3.1.d	GGAAAGACAAGAAAGGTGTGA	TTCTTGGGAAATGCTGAGG
	21pGM7	22705p7.14.c	GATCATGGCACAAGTTGCAG	GTTAATCGCAAACTGTACGAGAA
	21pGM8	18485p4.3.a	AAACCGGGTTCACATCAA	CTCTCCCTCCAGAAACT
	21pGM9	18409e2.1.a	CGGCTGTAGTTGCTCTCAC	GACTGTGGTCATGGCAGTG
	21pGM9	18419e2.6.a	ACCTGGGACTCCAGTGTGAC	GCTCATGTGGATGTCCTTGA
	21pGM10	18459e6.14.a	TCAGGCCACAAGTGAACATC	CACTGGGGGCTTGATTTTT
	21pGM11	30200e6.1.b	TGCTTTGGATTCACCCAT	CAAATGCCAAACACCAACC
	21pGM12	18473p1.4.a	TAAATGCAGCTGGGCACTTT	GCTGATCCAGTGGCCAGAA
	21pGM13	30202e6.12.a	TCAGACTTCAAATACTCCAGA	CAAITAAAATGTGTAATAGCCCACT
	21pGM14	22715p7.9.c	CAAGAACAGCCATGGAACA	TTTTCTTTCCAGTCCCATGC
	21pGM15	18393e1.1.a	CTTGGGGCTTGGTTCTATGT	TTCAITTTTGTGCTCTTGG
	21pGM16	18483p4.2.a	TGATCTTATGGCGGAGAAGG	CTGTGCAAAAGGCTTCTTA
	21pGM17	18457e6.11.b	ATCTGGAGATTTTCTAAGCTTCT	TTTGTAAATGACATAATTTTCAGCA
	21pGM18	18395e1.2.a	GCTTAAATGGGTTGTGTCT	CAATGCCTTCTAATGACCTT
	21pGM19	18447e3.2.a	TTGCTGGTCCATCTCTCT	TGAATTTCAACCAGCAAGAA
	21pGM20	18403e1.6.a	GAGGCAATACACTGGCATCA	GCTTCCCTTTTTCGTCTCT
	21pGM21	22717p7.17.c	CCAAGGTGCGAAAAGAAAG	CCAITGAGCAGCTGTTTGGT
	21pGM24	30917C21p-e6.4a	CAAGGTGAGGTGAGAGTCA	TGGCTTACCCTGGGATTTA
	21pGM25	18453e5.2.a	GCCATGGACTGTGATGCT	GCACCAGCCCTTTGATACTG
	21pGM26	18439e2.19.a	GGCTTGTCTTATGTCCCTGT	TCGGCAATCAGAAAGATCAGA
	21pGM27	22725p6.15.a	TGCCCCACAGAAGATTAGG	GCCTTTCCTCCCTGATTT
	21pGM28	18397e1.3.a	GCTCGCTCCTCTCACTTGT	GGCGACTACCATCGAAAGTT
	Copy number	21pGM6	Hs21pGM6.2-a	CCTTTACCTGGACCTTCTCATGA
		Hs21pGM6.2-b	TGCTGTCTGACAGATTAACCAAG	CTGAGGATACGTTTCCCAATCT
21pGM9		Hs21pGM9.1-i	CCTCCACCCAGAGACCC	GCGTGGCAGAGATTGTCT
		Hs21pGM9.1-e	AGGAACCTCTGCACATAGCCC	GCCGGTGTGAGAGCAACTACA
		Hs21pGM9.1-f	CTGCGCAACCTCAAGGCAT	GGCAGTGACGCTTCTTCCA
21pGM15		Hs21pGM15.1-b	GGCTTGGTCTATGTCCCTGC	TTCTCTTCGCCCTTGA
		Hs21pGM15.1-c	CGGACCAAGTGTCCCTAGTTG	CAATGTGCCACCGGTC
		Hs21pGM15.1-f	TTCCCTAGTTGTGGGAGCAGA	AACATAGCTGCACCTCAGGCC
21pGM25		Hs21pGM25.1-e	TCACCACTTCCATCCAATTC	TGCTGCATCTGTGAGTCCATG
		Hs21pGM25.1-h	CCCAGGAATTGATGTTGAATG	GAACCGTCCACGATCCGAGT
		Hs21pGM25.1-d	CACCCACAGTGCACACTCA	CCCTCAATCTTCTCTCCA
		Hs21pGM25.1-c	GAAGGACCTGGACTAGCCAT	TGTTGCTCCTGAGACAGCAC
21pGM28		Hs21pGM28.1-a	GGCTGACCCCTTCG	TGGTTTTGATCTGATAATGCAG
		Hs21pGM28.1-f	GGGCGCTGACCCCT	TGGTTTTGATCTGATAATGCAG
		Hs21pGM28.1-h	GCTGACCCCTTCGCG	CGGGTTGGTTTTGATCTGATAAT
		Hs21pGM28.1-n	TAGATAAACCCTGGGCGATCG	AAAGTTGATAGGGCAGAGCTCGAA
TBL2		15_TBL2	CTGTCTGTCTATGACCTCTGCA	CCAGGCCCTTCGGCA
MLXIP1		14_WBSCR14	CACGCATTTCTGATCCCAT	GCAGGCCGCTGCGC
RFC2	6_RFC2	TGACAAACCCCTTGCCAC	TGCCATATTTGGTAAATGCGC	
CYLN2	5_CYLN2	CTGCTCTTCTCTCTGTGGCT	AGCTGTTTGGCCAGGCTCT	
WBSCR16	4_WBSCR16	CATGCTTTTGGGAACTGCTT	AGAGGACAGCAGTGTGAATCCA	
AKAP1	HsAkap1_1	AAAGCACTTAGTTCGGTCCG	AAACTCACATAGCCGCCCTG	
AKAP1	HsAkap1_2	GATACCTGTGACTACGCCGA	CCCGGAGCAGCTCACTTTC	
HBG2	HsHBG2_1	CGGCTGGCTAGGGATGAG	TGTGGAACCTGTAAGGGTG	
SAMSN1	hsacgh-21-1A	TCTCGATCTTATAGGTCACCT	GCTCAGCTCAGGATCGAGAA	
NCAM2	hsacgh-21-3A	CCACCTCAAGTATACCGAAGTCTTA	GGAAGGCTGCATGTGAA	
-	hsacgh-21-4A	CAATTCAGGTGAGTGTGATACTAGTA	GCCAGGTTTGAATGTTTGTCTAAGTC	
-	hsacgh-21-6A	CGTCACATCACATCTCACTATTG	ACCACGTGAAAGGAGGTTTCC	
-	HsHUG1	GCCAGGAGCCCACTGTA	AGATTTCTGCAGCCACCTCAC	
ELN	9_ELN	GGCATGAGACGCTCCACAT	TTCTGAGCCAGGGCAACAC	
POR	0_1_400KbWBcentro	CACCAAAGACTGGCGTTTCC	GCGGACCCAGATTTCTTGAC	
Variation	21pGM6	21pGM6.2cDNA	CAGAGCCCTTACGTGGACCT	TGCGTATTTCTGCTGGAGT
	21pGM9	21pGM9.1cDNA	GGCGGTCTGTAGCTGAG	CTGTTTATTTGGGGGATTTGG
	21pGM15	21pGM15.1cDNA	GGCTTGGGGCTTGGTTCT	CCAGAGATTTATTCATTTCTTGTG
	21pGM25	21pGM25.1cDNA	ACAGGOCACCTCAGCTCACT	TTGAAGTGGAAACCTCCAC
	21pGM28	21pGM28.1cDNA	TGTGGTAATTTAGAGCTAATACATGC	GTGGCGACTACCATCGAAAG