# The genomic and transcriptomic landscape of a HeLa cell line

**Author list:**

Jonathan Landry*[#], Paul Theodor Pyl*[#], Tobias Rausch*, Thomas Zichner*, Manu M. Tekkedil*, Adrian M. Stütz*, Anna Jauch§, Raeka S. Aiyar*, Gregoire Pau*[1], Nicolas Delhomme*[2], Julien Gagneur*[3], Jan O. Korbel*, Wolfgang Huber*, Lars M. Steinmetz*

**Affiliations:**

* European Molecular Biology Laboratory, Genome Biology Unit, Heidelberg, 69117 Germany

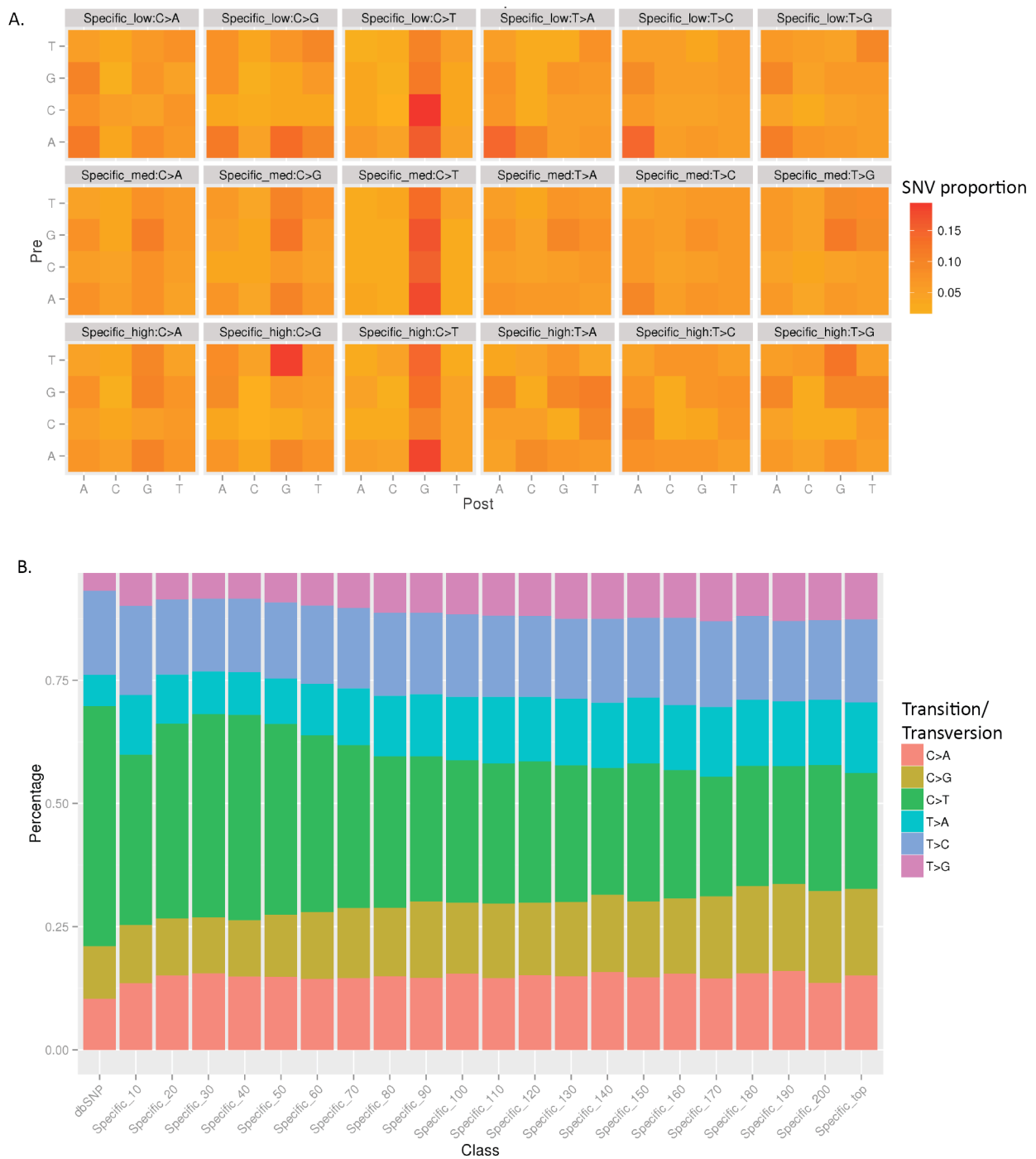§ University Hospital Heidelberg, Institute of Human Genetics, Heidelberg, 69120 Germany

[#] These authors contributed equally to this work.
[1] Department of Bioinformatics and Computational Biology, Genentech Inc., South San Francisco, California, 94080
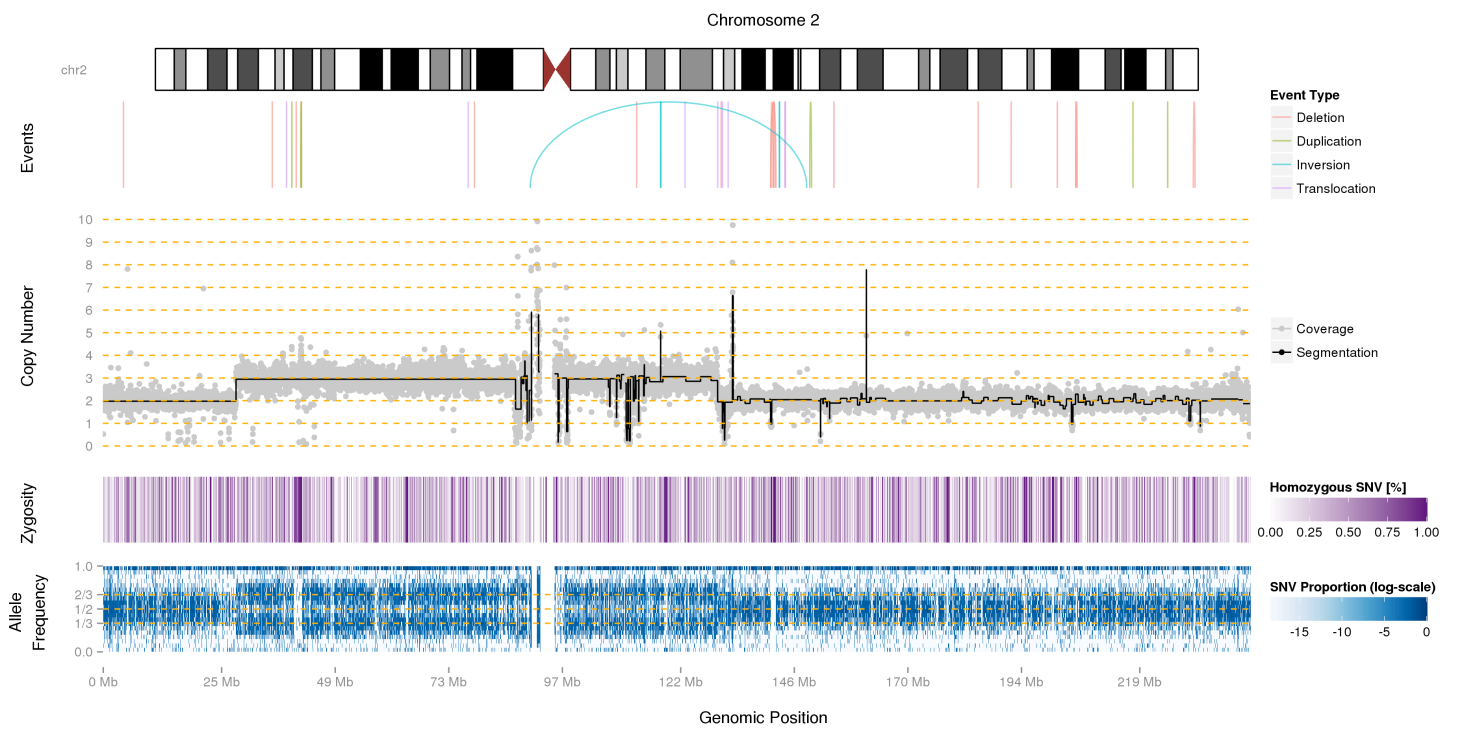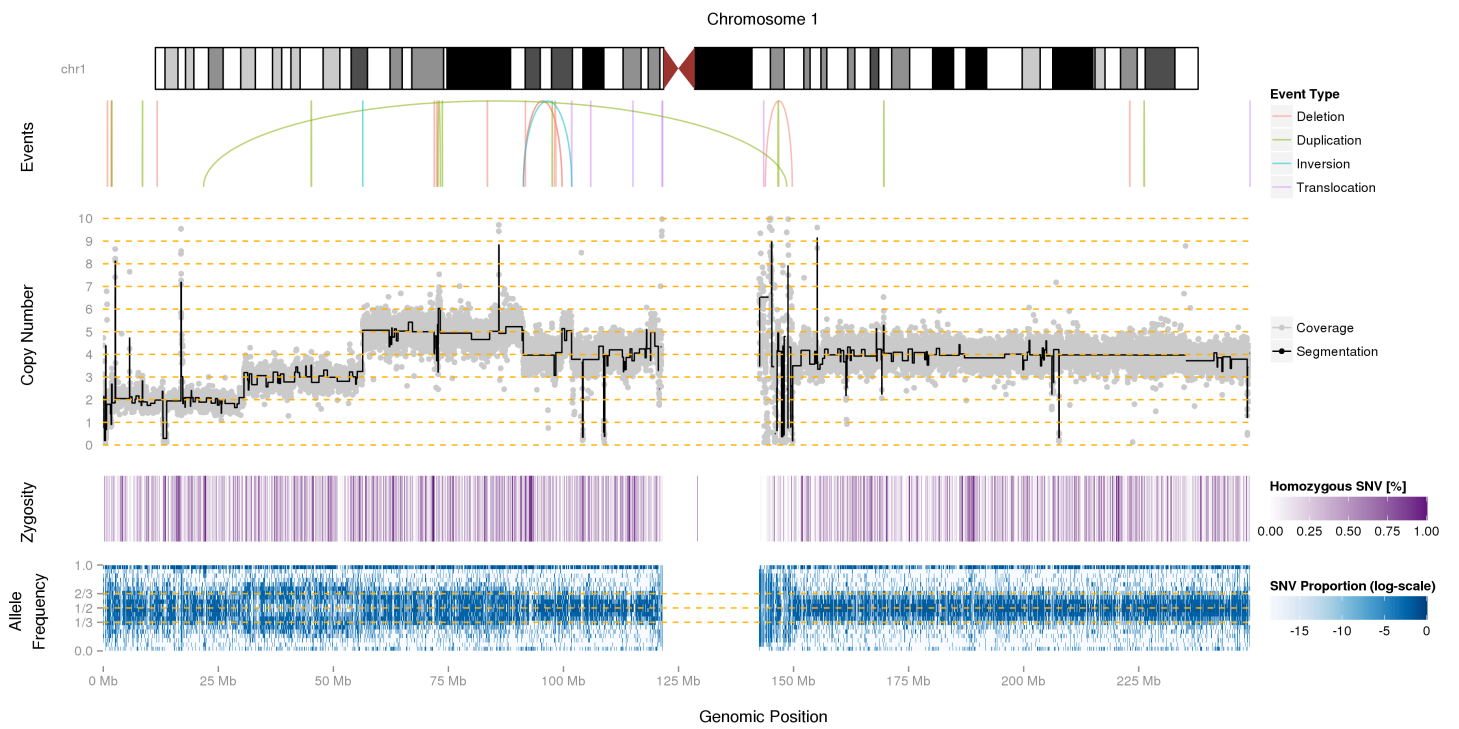[2] Department of Plant Physiology, Umeå Plant Science Center, Umeå, Sweden, S-901 87
[3] Department of Chemistry and Biochemistry, Ludwig-Maximilians-Universität München, Munich, Germany, 81377
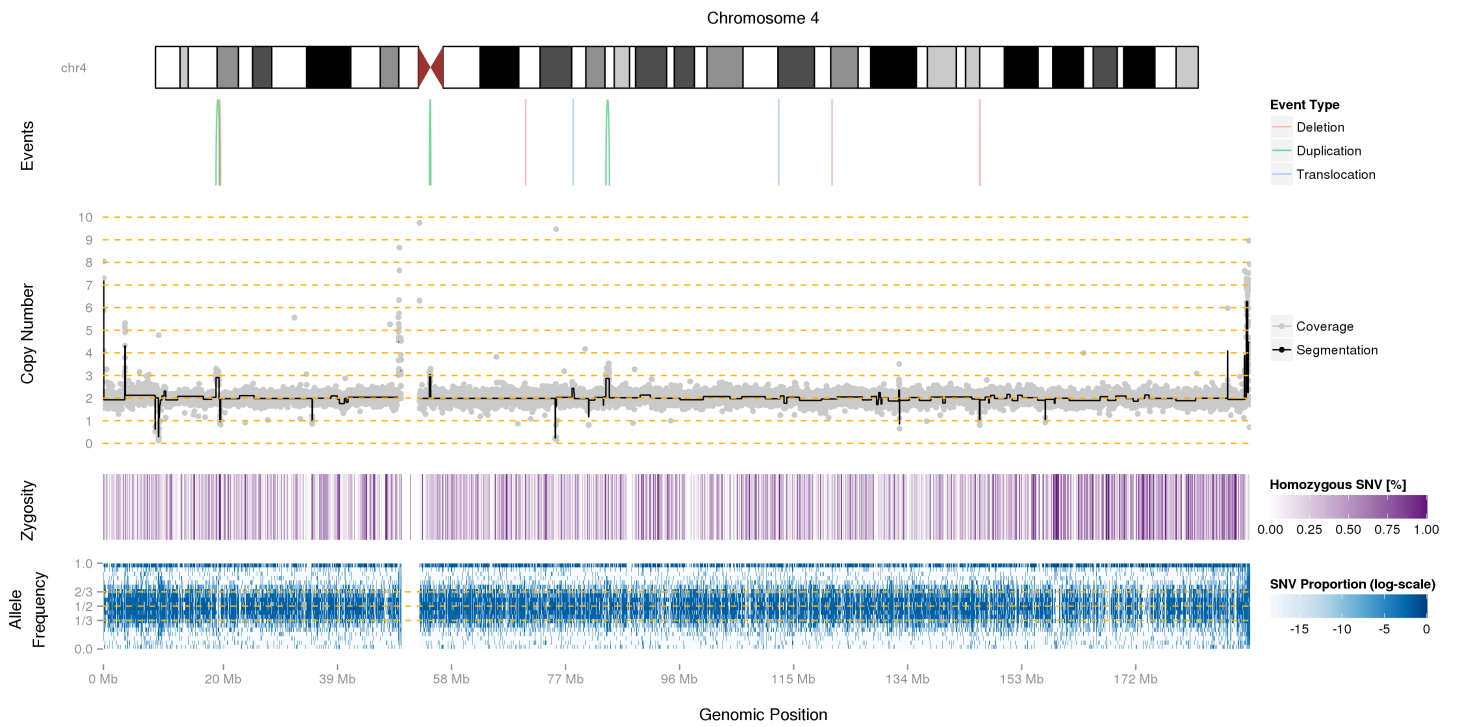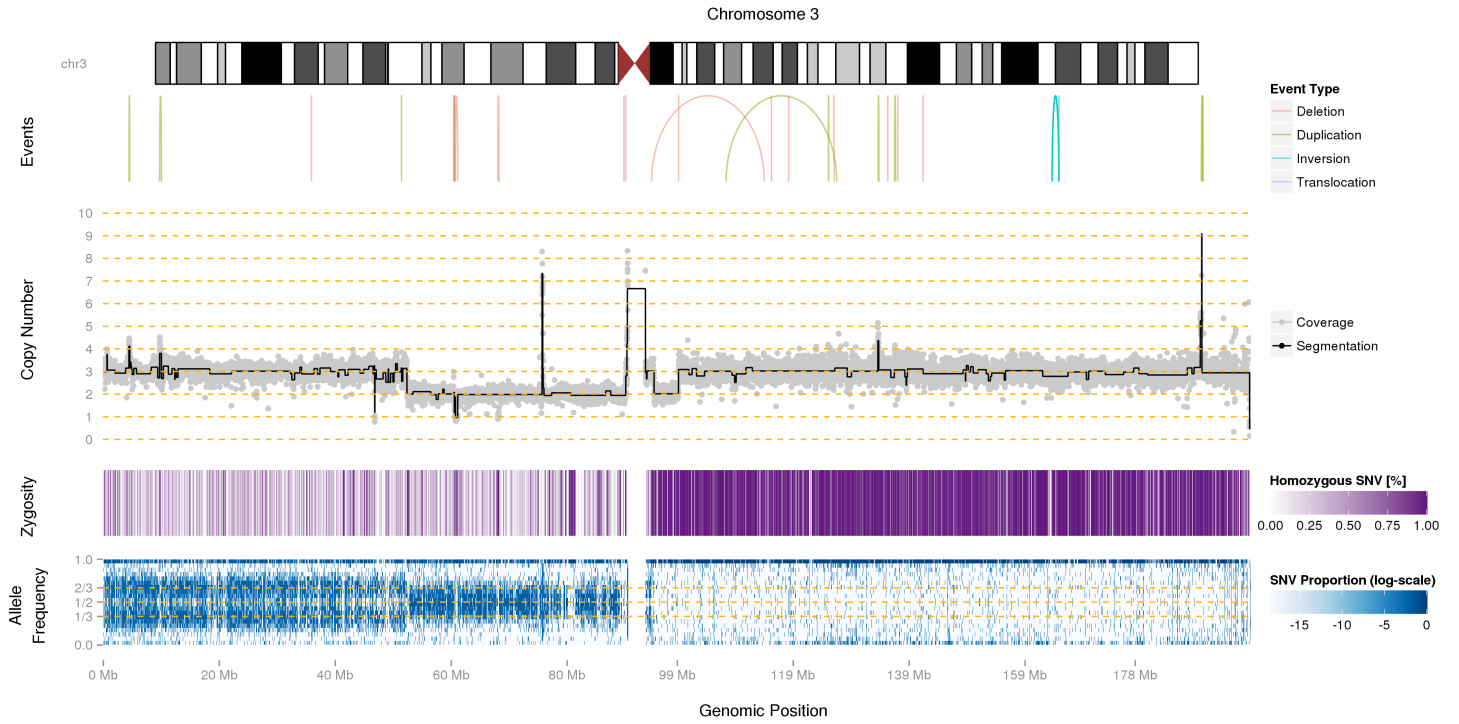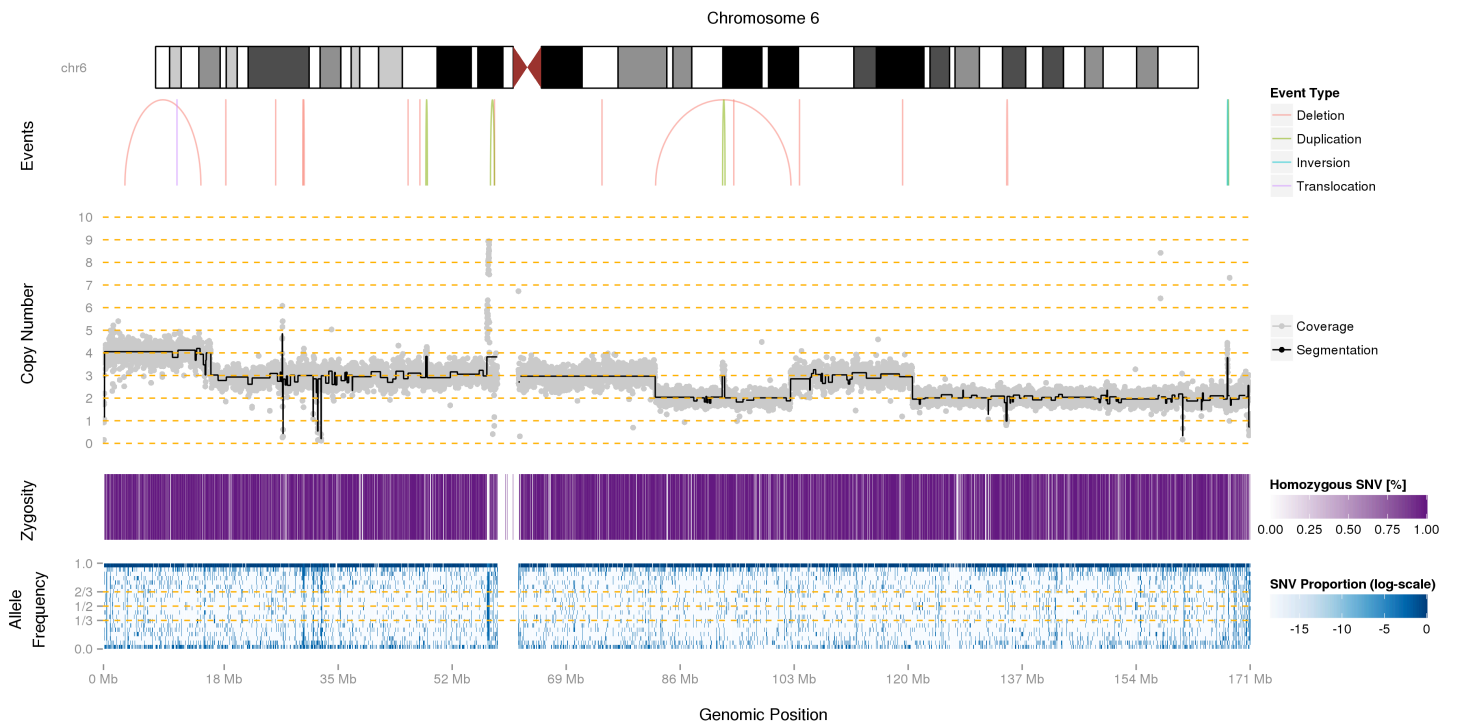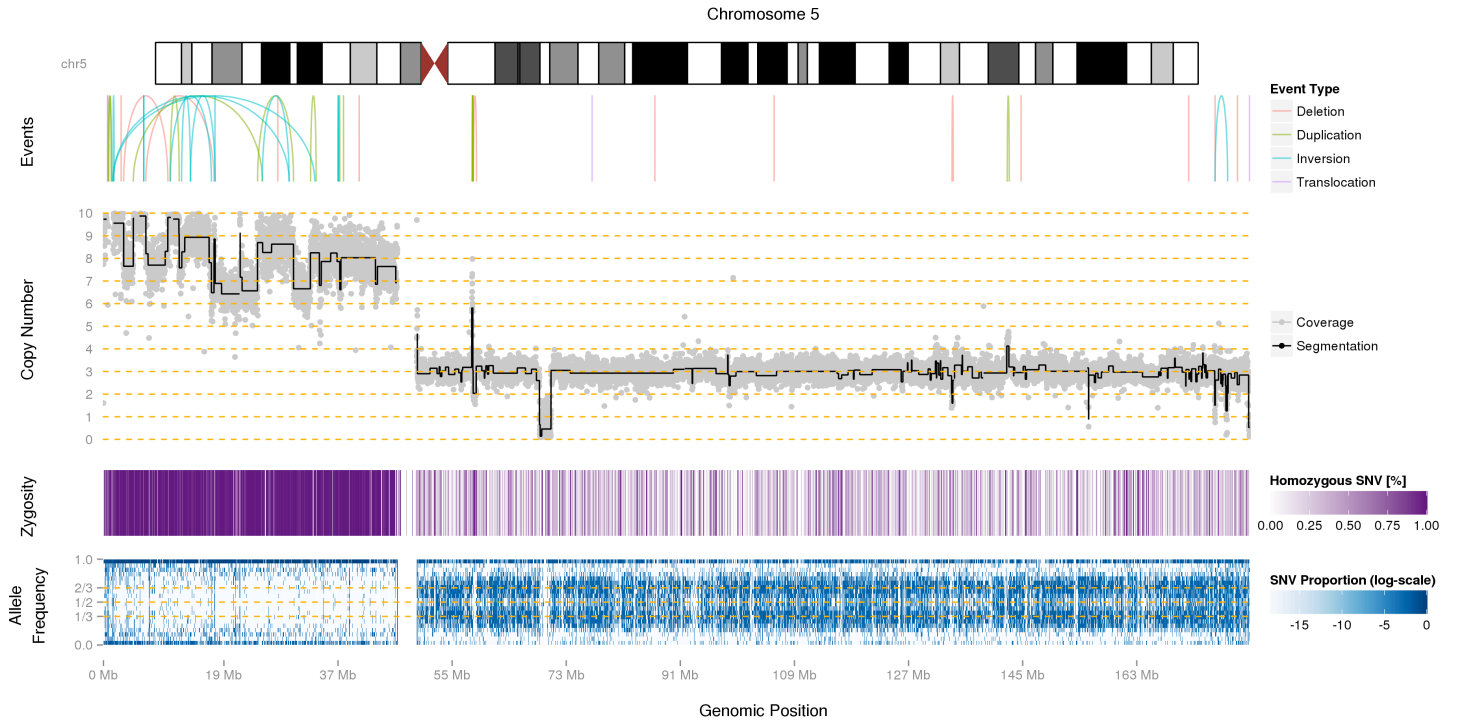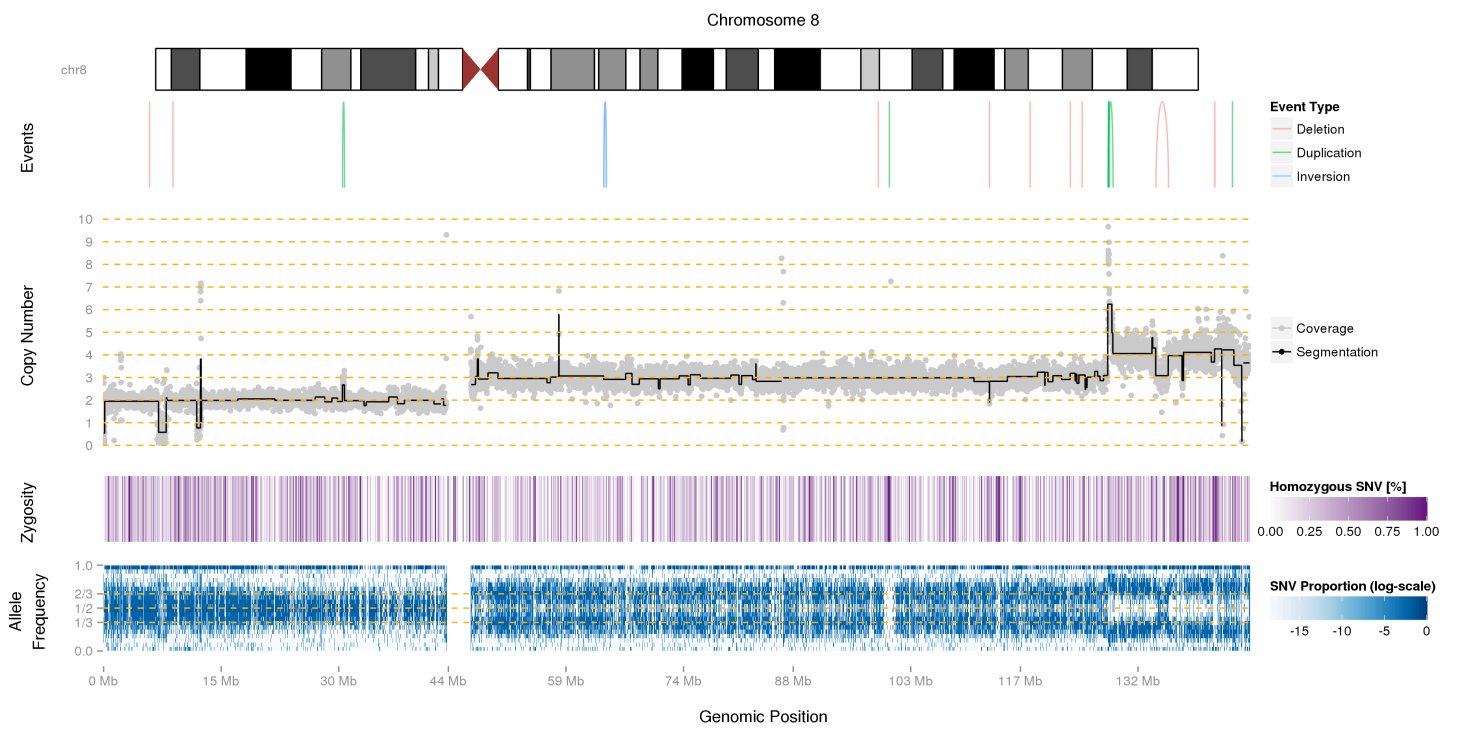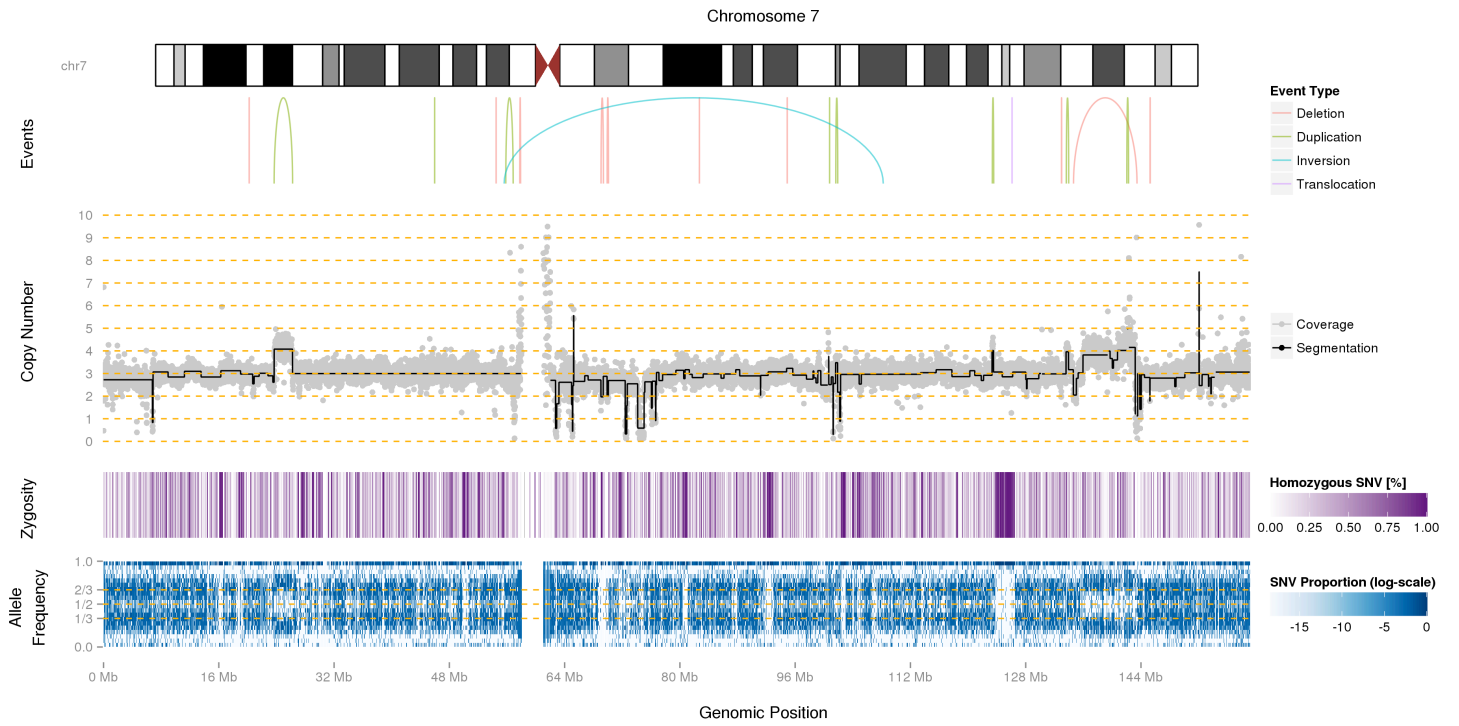
**Figure S1  Mutational spectra in the HeLa Kyoto genome.**

(A) Heatmap representation of mutational spectra observed in HeLa Kyoto specific calls stratified by local coverage (low: < 10; med: >=10 and < 60; high: >= 60). Each row corresponds to one group of SNVs and each column to a type of mutation (e.g.: C>A). Each heatmap shows the observed proportion of SNVs for all possible combinations of preceding ('Pre') and following ('Post') bases.

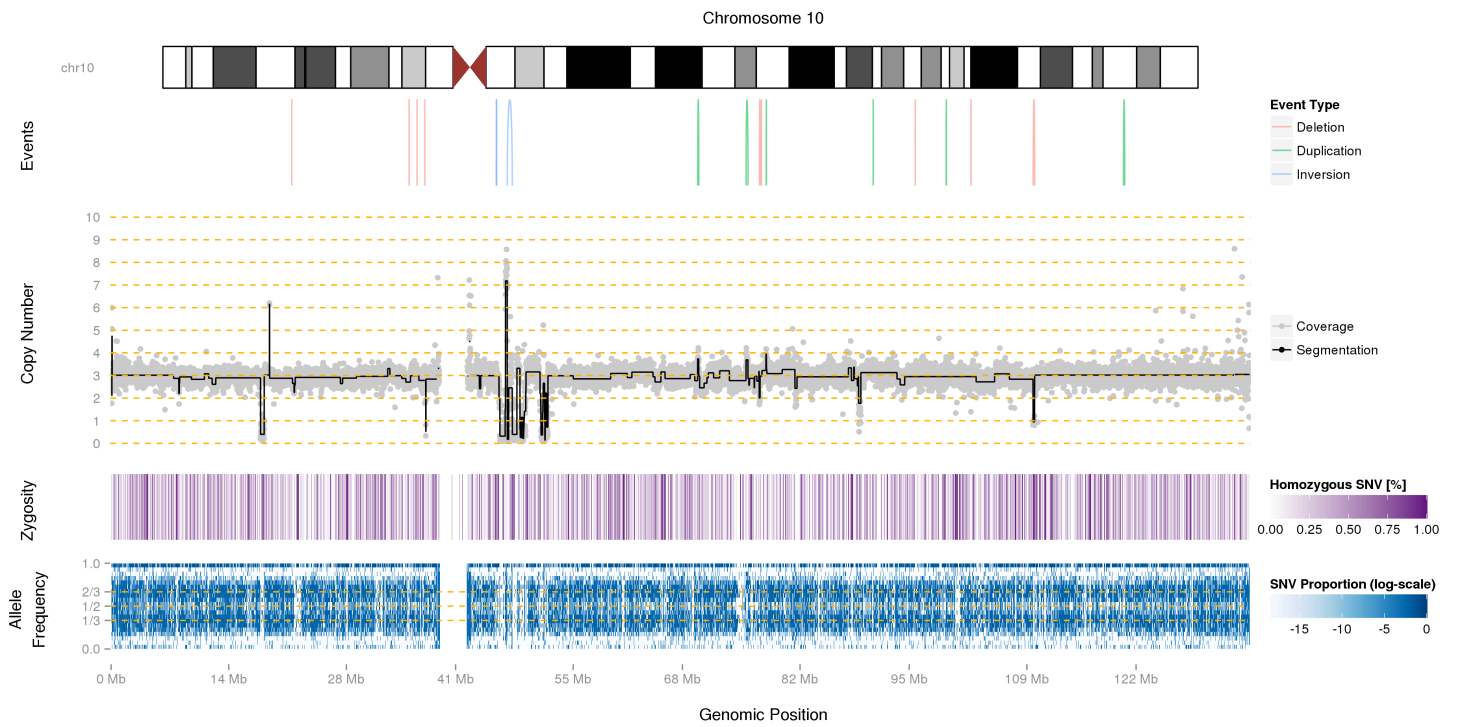(B) Stacked barplots of the distribution of mutation types (e.g.: C>T) in observed HeLa specific SNVs, stratified by local coverage: Column "Specific_x" represents calls with coverage between x-10 and x. Column "dbSNP" shows called SNVs that have dbSNP ids, column "Specific_top" shows calls with local coverage greater than 200.

Chromosome 1

Chromosome 2

Chromosome 3

Chromosome 4

Chromosome 5



Chromosome 6

Chromosome 7

Chromosome 8

Chromosome 9



Chromosome 10

**Chromosome 11**

**Chromosome 12**

Chromosome 13

Chromosome 14

Chromosome 15



Chromosome 16

Chromosome 17

Chromosome 18

Chromosome 19

Chromosome 20

Chromosome 21

Chromosome 22

**Figure S2   Structural variants, copy number and loss of heterozygosity for chromosomes 1 to 22 and X.**
Arcs in the top panels labeled 'Events' represent the predicted connections between fragments derived from SV calls based on read pair orientation and spacing. Different read pair signatures indicate the following event types: deletions, tandem dup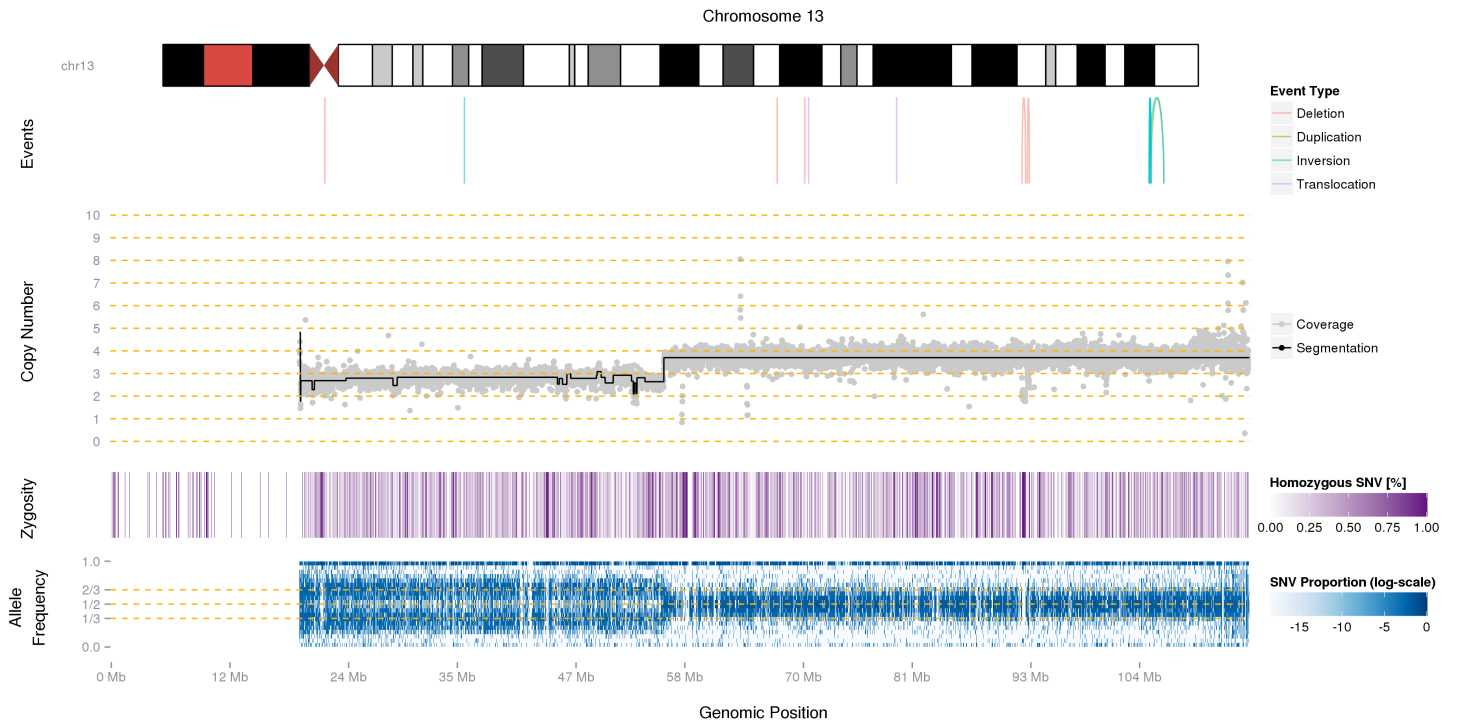lications, inversions, and interchromosomal translocations. The center panel (Copy Number) represents the copy number estimates in 10 kb bins (grey) overlaid with their segmentation (black). The associated CN is shown on the y-axis. The zygosity track shows the proportion of homozygous SNV calls in 10 kb bins, darker purple regions contain more homozygous calls (up to 100%) and indicate potential LoH. The bottom panel shows the allele frequency (AFS) distribution as a heatmap in 10 kb bins on the chromosome axis and 5% bins on the allele frequency axis; darker blue indicates more SNVs with the given AF in the corresponding 10 kb region. The color scale is according to the log of proportion of SNVs falling into the AFS bin (e.g. 10-15%, i.e. the row) in the 10 kb region (i.e. the column).

**Figure S3  Principal component analysis of SNVs in HeLa Kyoto and 640 HapMap individuals from 8 different populations.**
The populations are:
CEU - Northern and Western European ancestry from Utah, USA;
CHB - Han Chinese from Beijing, China;
JPT - Japanese from Tokyo, Japan;
YRI - Yoruba from Ibadan, Nigeria;
ASW - African ancestry from Southwest USA;
LWK - Luhya from Webuye, Kenya;
MXL - Mexican ancestry from Los Angeles, USA;
TSI - Toscani from Italy.

**Figure S4  Effect of GC adjustment on sequencing coverage.**
(A) Mean coverage value per 10 kb bin (y-axis) as a function of the mean GC % per 10 kb bin, before any adjustment was applied. The black line represents a robust local regression fitted on the data.
(B) Adjusted mean coverage value per 10 kb bin (y-axis) as a function of the mean GC % per 10 kb bin.
(C) Distribution of adjusted mean coverage value per 10 kb bins along the genome. The green line corresponds to the adjusted coverage median value for chromosome 4, which was used for normalization as representative of regions with CN 2.

**Figure S5   Primer design for detection of genomic rearrangements.**
For deletions, the amplicons generated are smaller than expected size. The orientation of the primer targeting tandem duplication and inversion allow detection of an amplicon only if the event is present.

**Figure S6   Comparison of HeLa transcriptome profile to Illumina Body Map tissues and ENCODE cell lines.**
(A) Venn diagram of genes with nondetectable expression across all ENCODE cell lines (excluding HeLa-S3), all Body Map tissues, HeLa Kyoto and HeLa S3. The numbers indicate genes (listed in the human reference annotation file from ENSEMBL, see Methods).
(B) Dendrogram comparing the transcription profiles of all ENCODE cell lines and the HeLa Kyoto cell line. The y-axis is the Euclidean distance between cell lines.

**Table S1  Potential viral insertions**

| Chr | Start | End | Strand | Read support | Viral species |
|---|---|---|---|---|---|
| 1 | 10002 | 10118 | + | 12 | Human herpesvirus 6B<br>Equid herpesvirus 2<br>Human herpesvirus 6A<br>Human herpesvirus 7<br>Gallid herpesvirus 2<br>Gallid herpesvirus 3<br>Meleagrid herpesvirus 1<br>Ovine herpesvirus 2<br>Cyprinid herpesvirus 3<br>Saimiriine herpesvirus 1 |
| 8 | 128189764 | 128190032 | - | 552 | Human papillomavirus 18 and 32<br>(Alphapapillomavirus 7 and 1) |
| 8 | 128192272 | 128192499 | + | 23 | |
| 8 | 128193167 | 128193435 | + | 174 | Human papillomavirus 18<br>(Alphapapillomavirus 7) |
| 8 | 128200419 | 128200690 | + | 163 | |
| 12 | 49659073 | 49659136 | - | 6 | Human herpesvirus 5 |
| 12 | 95467 | 95567 | + | 9 | Human herpesvirus 6B<br>Human herpesvirus 6A<br>Human herpesvirus 7<br>Gallid herpesvirus 2<br>Gallid herpesvirus 3<br>Meleagrid herpesvirus 1<br>Cyprinid herpesvirus 3 |
| 13 | 19648667 | 19648780 | + | 7 | |
| 13 | 19649043 | 19649209 | - | 8 | Taterapox virus |
| X | 155185203 | 155185362 | - | 8 | Human herpesvirus 6B<br>Equid herpesvirus 2<br>Human herpesvirus 6A<br>Human herpesvirus 7<br>Gallid herpesvirus 2<br>Gallid herpesvirus 3<br>Meleagrid herpesvirus 1<br>Ovine herpesvirus 2 |