

Additional File 1

Interpreting whole genome and exome sequencing data of individual gastric cancer samples

Daniela Esser^{1,4}, Niklas Holze², Jochen Haag², Stefan Schreiber^{1,3}, Sandra Krüger², Viktoria Warneke², Philip Rosenstiel^{1†} and Christoph Röcken^{2*†}

- ¹ Institute for Clinical Molecular Biology, Christian-Albrechts-University, 24105 Kiel, Germany
- ² Institute for Pathology, Christian-Albrechts-University, 24105 Kiel, Germany
- ³ Department of General Internal Medicine, University Hospital Schleswig-Holstein, 24105 Kiel, Germany
- ⁴ Institute for Experimental Medicine, Christian-Albrechts-University, 24105 Kiel, Germany

Figure S1

Coverage distribution from (A) WGS and (B) WES (target region). Number of called SNVs in each sample using different variant caller (C).

Figure S2

Estimation of sample purity. Peak shows the mutant allele fractions for heterozygous SNVs. This results in a tumor ratio, which is twice as high as the peak.

Figure S3

Comparison between SNVs called in WES data and SNVs called in WGS data. Note that in this case, “somatic” refers to SNVs, which were not detected in the corresponding non-tumor samples sequenced with the same technology.

Figure S4

Ratio of SNVs, which were called in one sequencing data set and could be validated in the second. (A) Total number of SNVs. (B) Exclusively SNVs, which were covered in both data sets. Note that in this case, “somatic” refers to SNVs, which were not detected in the corresponding non-tumor samples sequenced with the same technology.

Figure S5

Rainfall plots based on all somatic SNVs. The x-axis shows the chromosomal position. The distance between each mutation and the prior variant is plotted on the y-axis. The colors of the dots indicate the SNV-type.

Figure S6

Somatic SNV signatures estimated with NMF (non-negative matrix factorization).

Figure S7

(A) Description of SNV signatures identified in the investigated tumor samples. (B) Contribution of the identified signatures to the investigated samples.

Figure S8

Somatic SNV spectrum of TCGA WES samples and investigated GC tumor samples. Somatic signatures were estimated with NMS.

Figure S9

Number of somatic SNVs and somatic indels after different steps.

Figure S10

Protein network for MSI (A) and MSS (B) tumors based on all somatic SNVs, which were a stopgain mutation, predicted as damaging, at a conserved position or in a conserved gene. The connections were based on medium stringency according to the STRING database. Proteins without any connection are not displayed.

Figure S11

Enriched processes of the two protein subnetworks identified in the first patient. (A) Larger network. (B) Smaller network.

Figure S12

Enriched GO terms in MSS tumor. GO terms, which were more often affected in the tumor sample than in all individuals of the 1000 Genomes project, were included in the figure. To enable comparability between the samples, all called variants (germline + somatic) were used. (A) Analysis based on number of genes affected by at least one SNV with damaging prediction. (B) Analysis based on number of SNVs predicted as damaging.

Figure S13

Enriched GO terms in MSI tumor. GO terms, which were more often affected in the tumor sample than in all individuals of the 1000 Genomes project, were included in the figure. To enable comparability between the samples, all called variants (germline + somatic) were used. (A) Analysis based on number of genes affected by at least one SNV with damaging prediction. (B) Analysis based on number of SNVs predicted as damaging.

Figure S14

Algorithm to detect and filter large insertions

Figure S15

Overlap of SNVs called by different programs

Figure S16

The figure shows the number of verified SNVs for each SNV caller. Exclusively positions covered in both data sets were included in the analysis. In the exome data are from outside to the inside the results of the SNV callers GATK, Samtools, and DiBayes. In the genome data the results of GATK are displayed outside and the one of Samtools inside. The percentage of SNVs, which have less than 20% read support in the other data set are marked in light, the remaining ones in dark grey.

Figure S17

The figures show for each SNV-Caller the percentage of called SNVs out of all cross platform variants. (A) WGS and (B) WES. The SNVs were called confirmed, if in the data of the second technology (i) 10%, (ii) 20% or (iii) 50% of the reads support the variant. (C) Table summarizing the confirmed SNV counts.

Figure S18

Comparison between strict filtered (detected with all callers) SNVs in the WES data (called with Samtools and GATK and DiBayes) and those detected in the WGS data (called with Samtools and GATK).

Figure S19

False positive and false negative rate in WES and WGS sequencing data. All exonic SNVs called with Samtools were investigated in all sample pairs and compared between WES and WGS. (A) Percentage of SNVs, which were called in WGS and either uncovered or not supported in the WES data. (B) SNV type distribution for three classes of SNVs called in WGS: all SNVs called in WGS data, SNVs uncovered in WES, SNVs covered but not supported in WES. (C) Percentage of SNVs, which were called in WES and either uncovered or not supported in the WGS data. (D) SNV type distribution for three classes of SNVs called in WES: all SNVs called in WES data, SNVs uncovered in WGS, SNVs covered but not supported in WGS

Figure S20

Relation between the novel developed exonic gene conservation score and known cancer associations: (A) Scatter plot comparing the exonic gene conservation score with the cancer proliferation index, whereby tumor suppressors exhibit significant negative cPI values and oncogenes positive cPI values. (B) Enlargement of the small value region of figure part A. (C) Scatter plot comparing the exonic gene conservation score with the number of samples having a mutation with a FATHMM cancer prediction in the COSMIC database. (D) Enlargement of the small value region of figure part C.

Figure S1

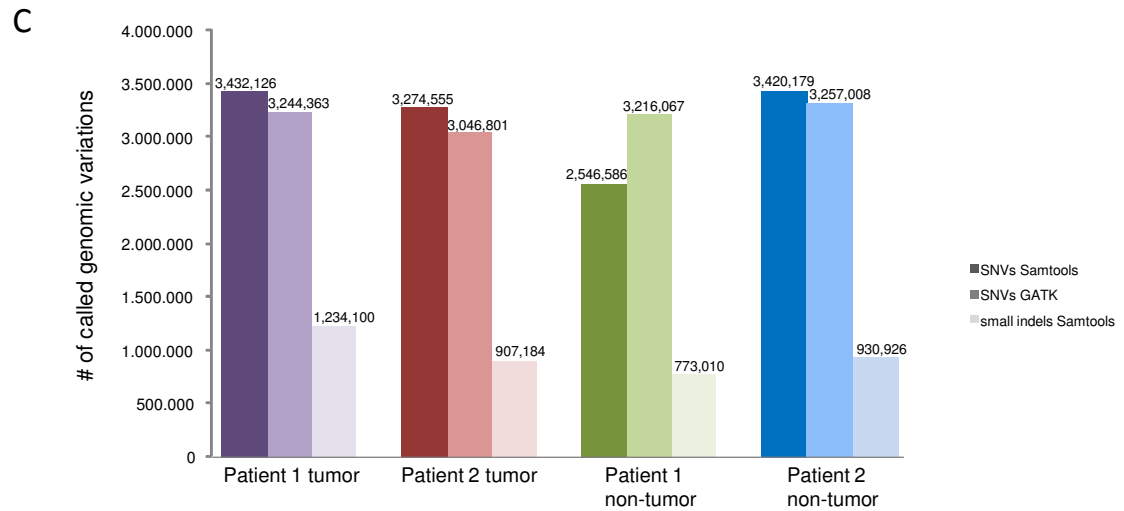
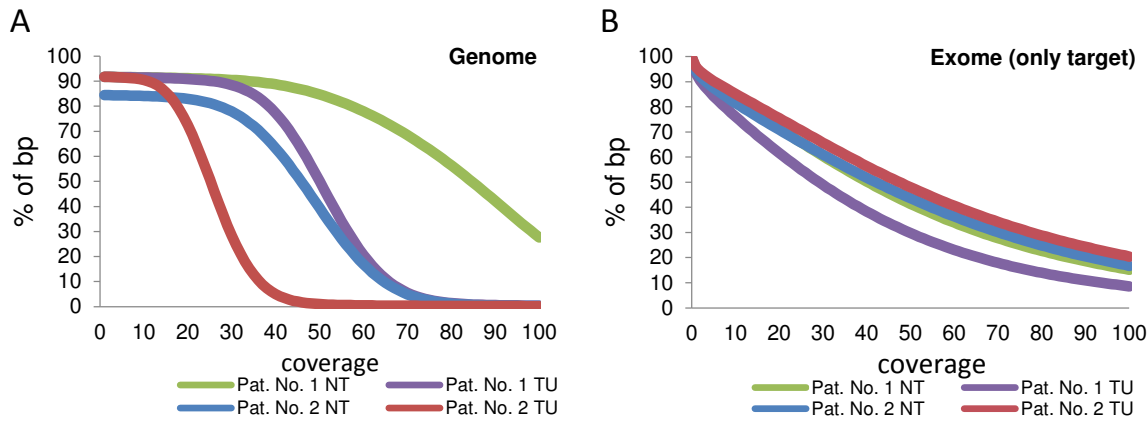


Figure S2

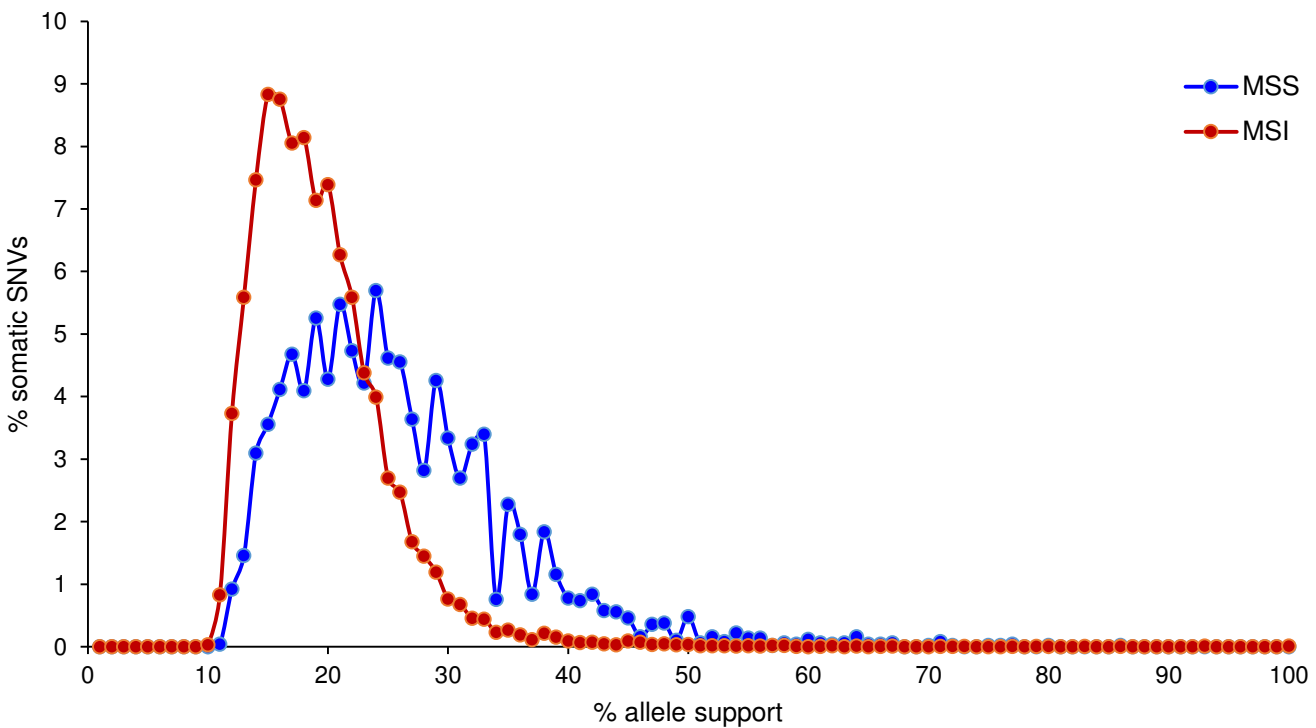


Figure S3

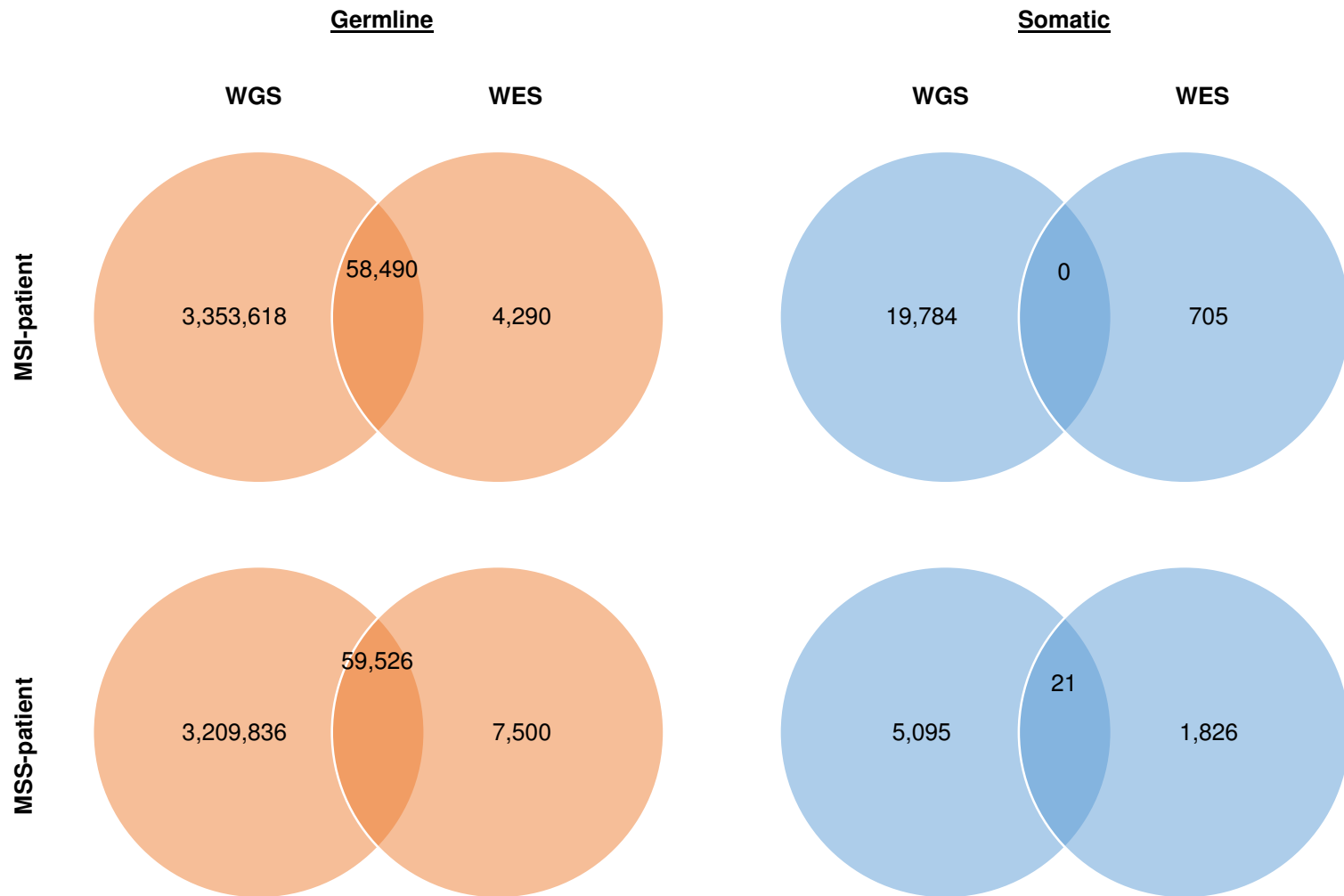


Figure S4

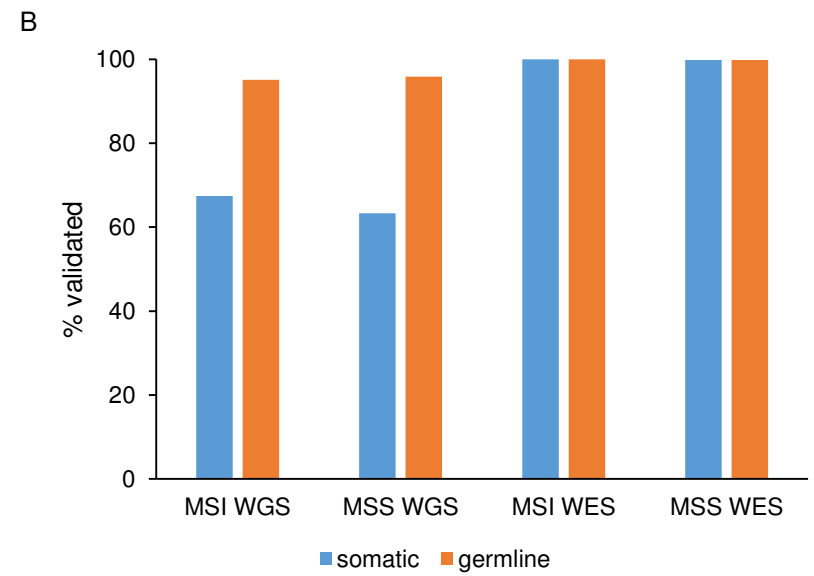
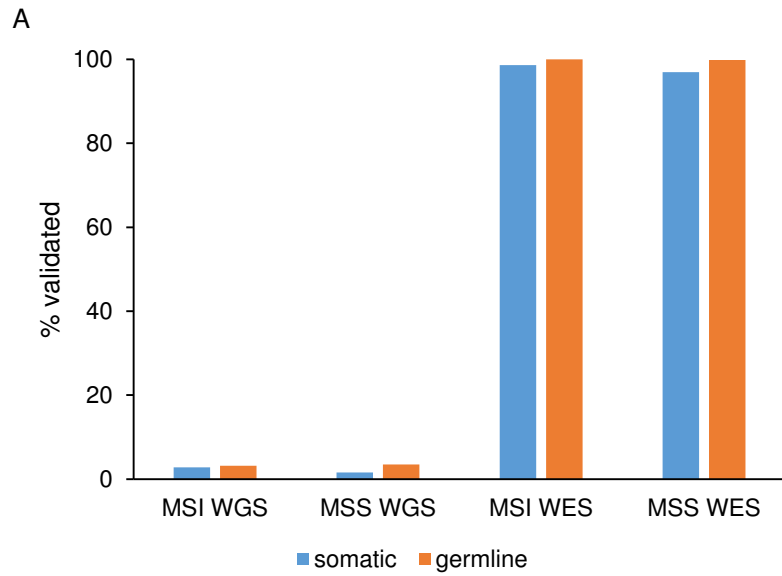
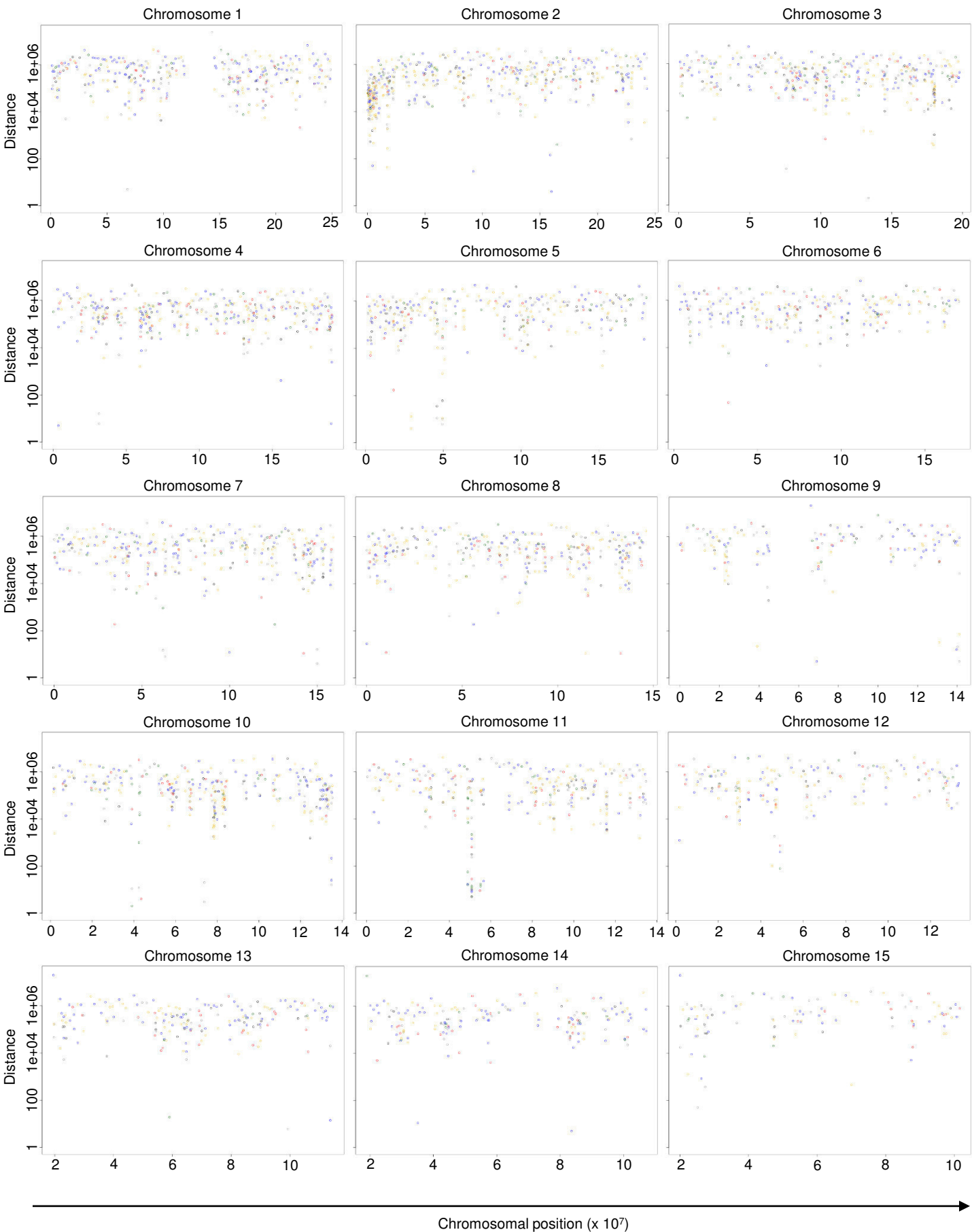
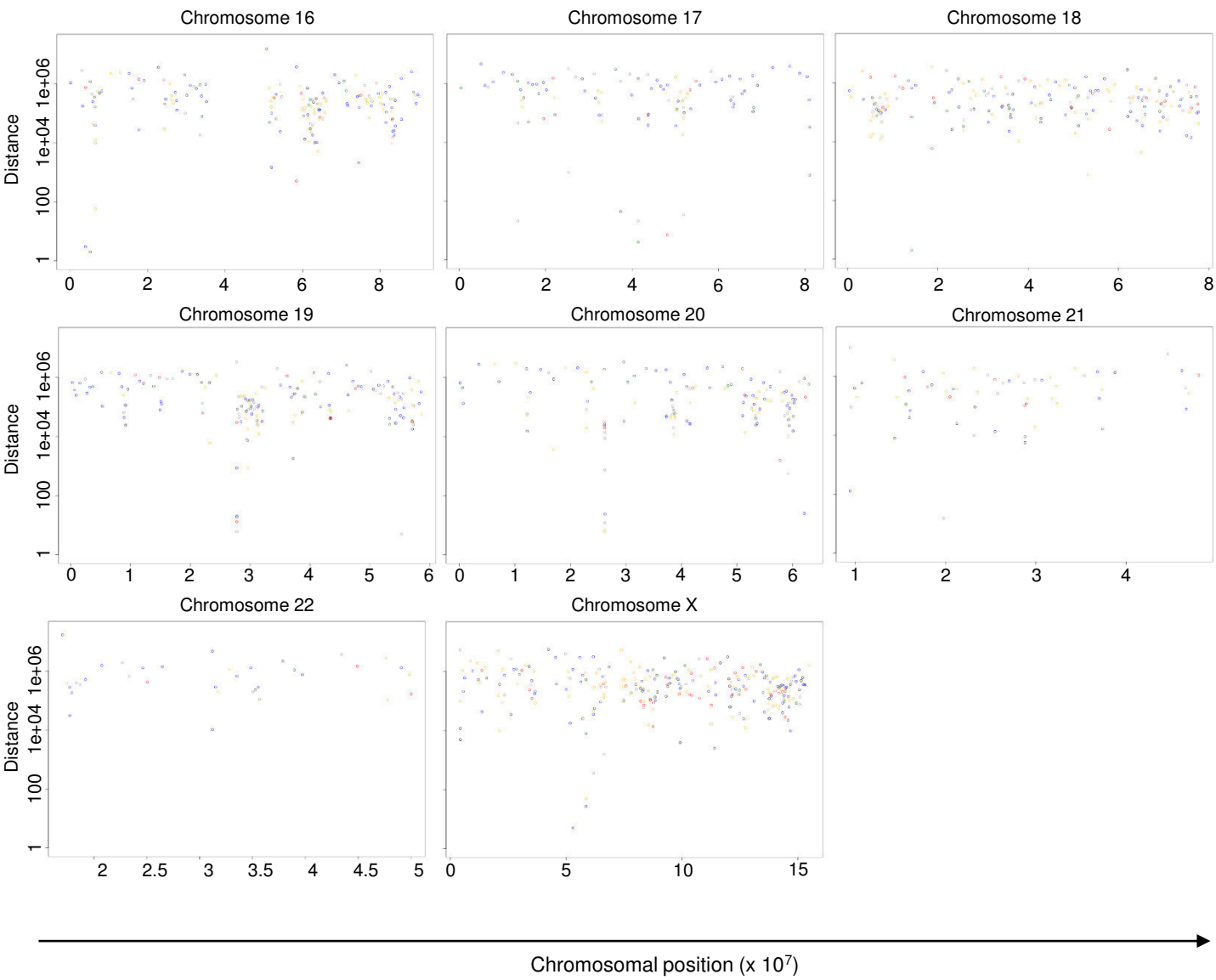


Figure S5A

MSS tumor





C>A



C>G



C>T



T>G



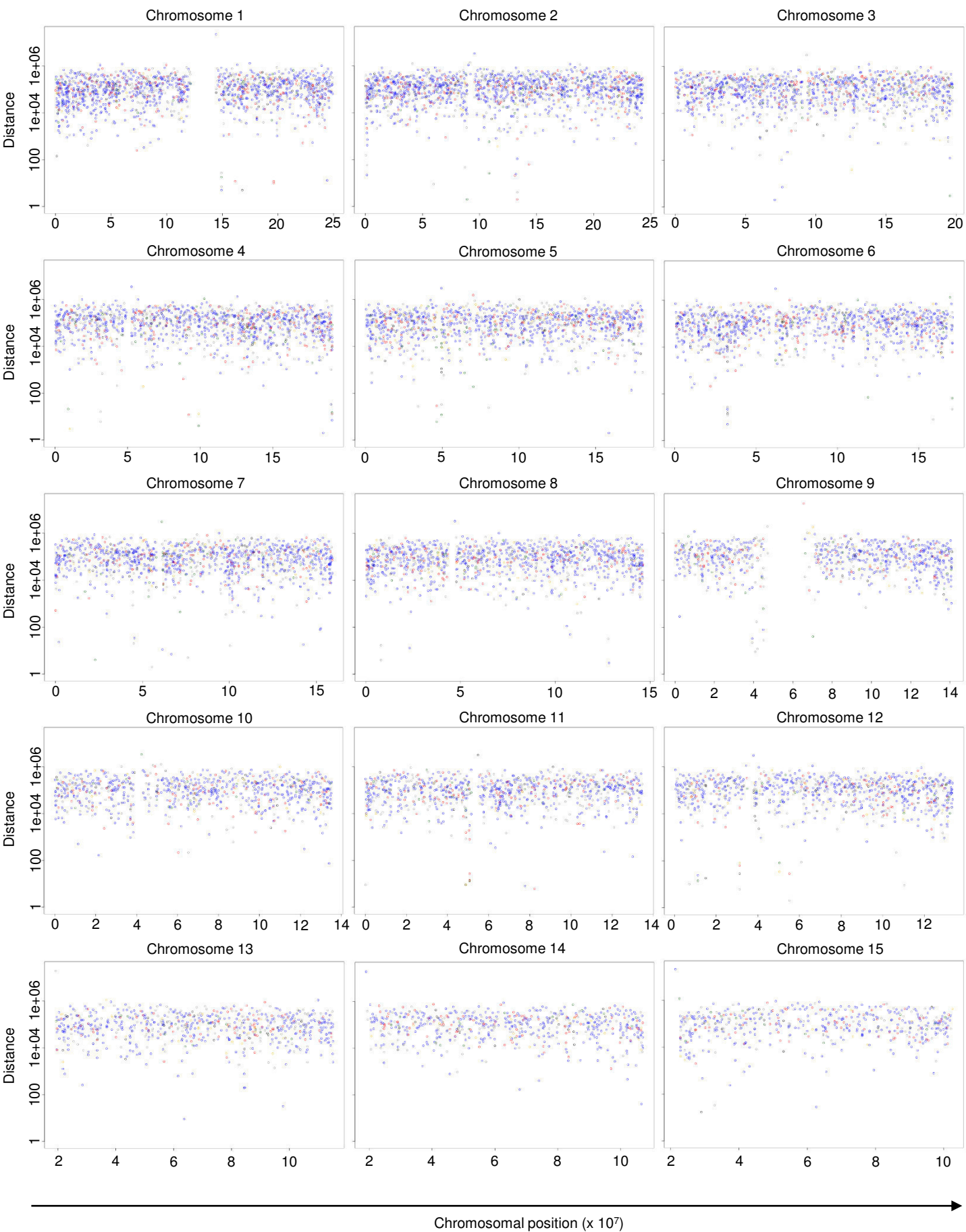
T>C



T>A

Figure S5B

MSI tumor



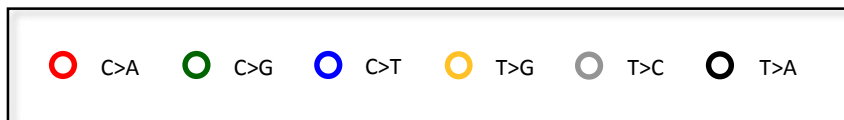
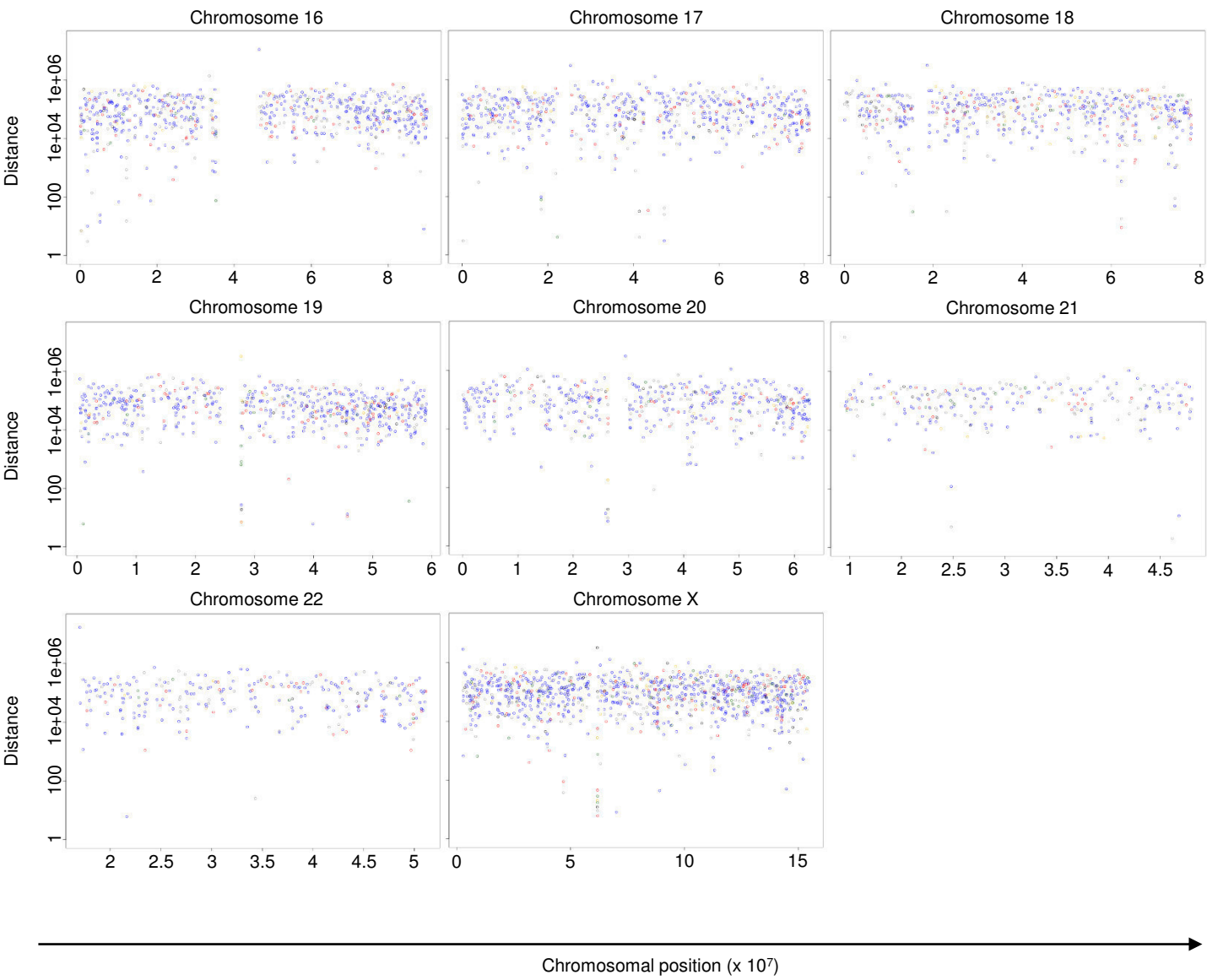
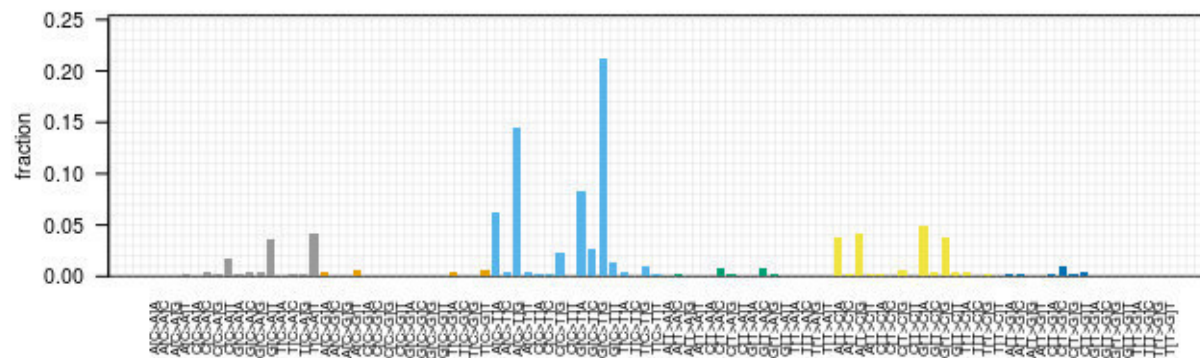
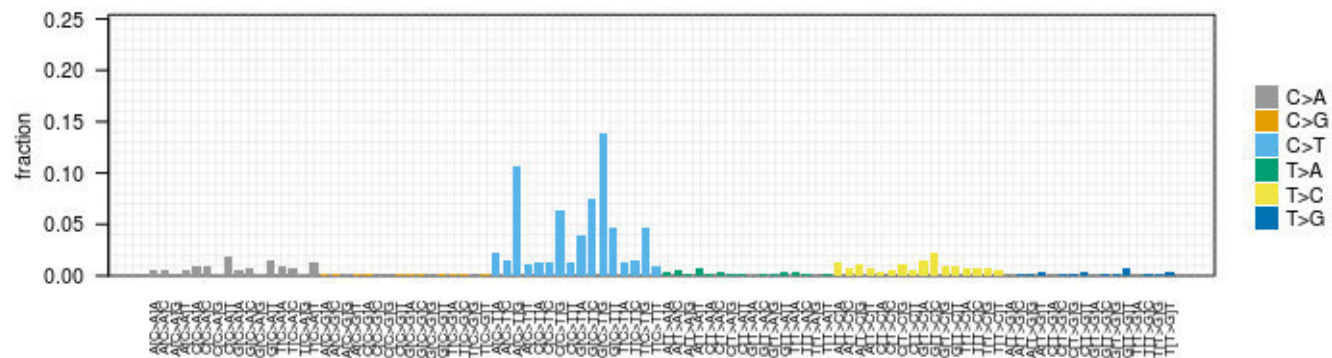


Figure S8B

MSI



Signature.1 : 0.491 & Signature.6 : 0.129 & Signature.15 : 0.381



error = 0.154

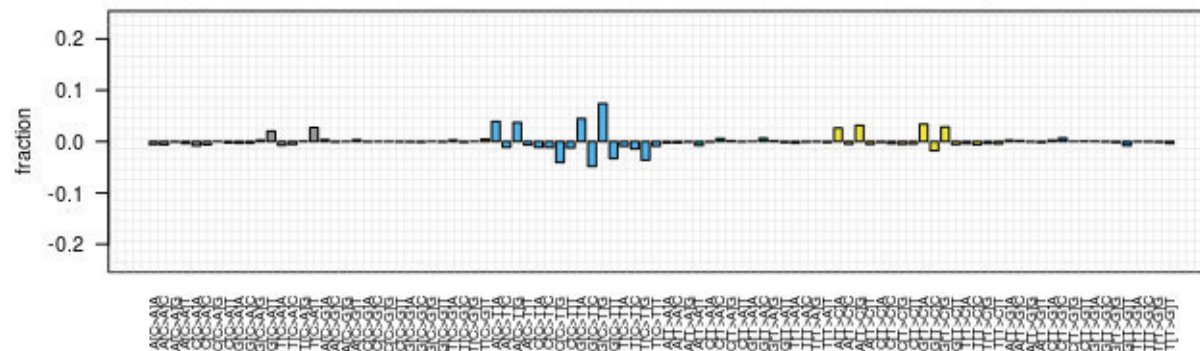


Figure S9

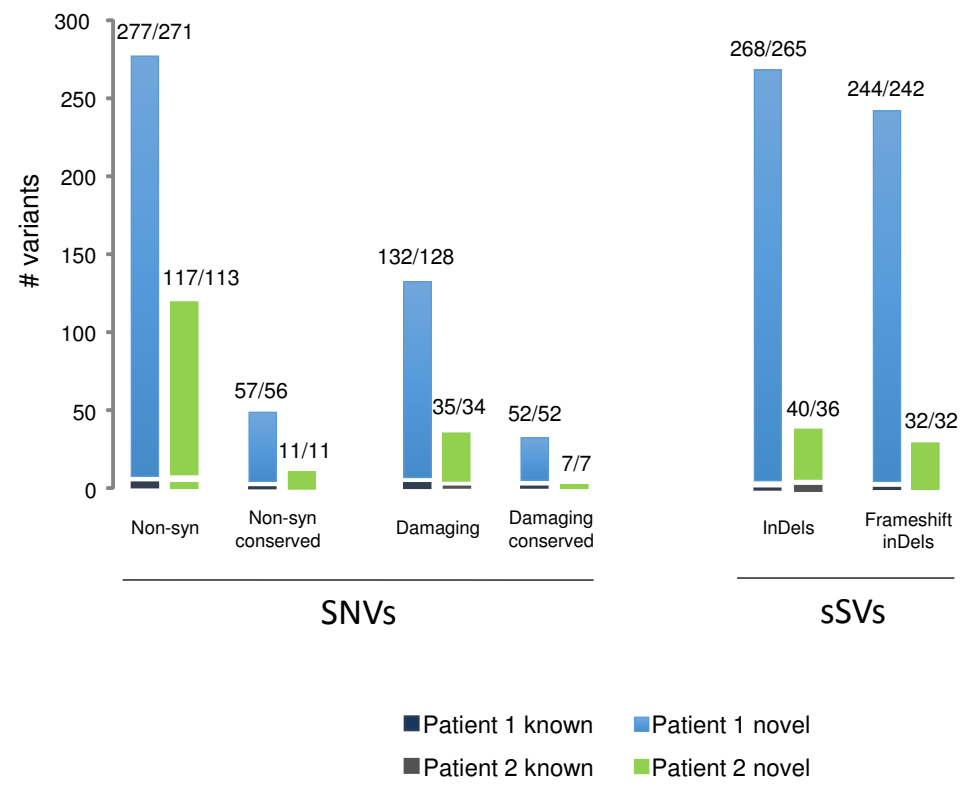


Figure S10 B

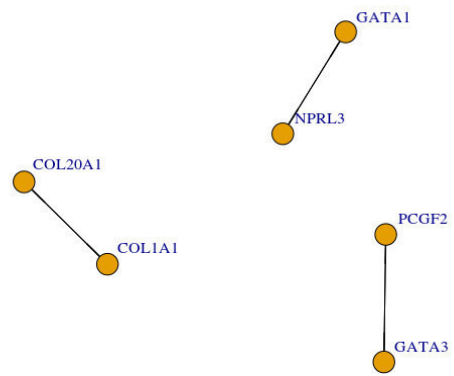
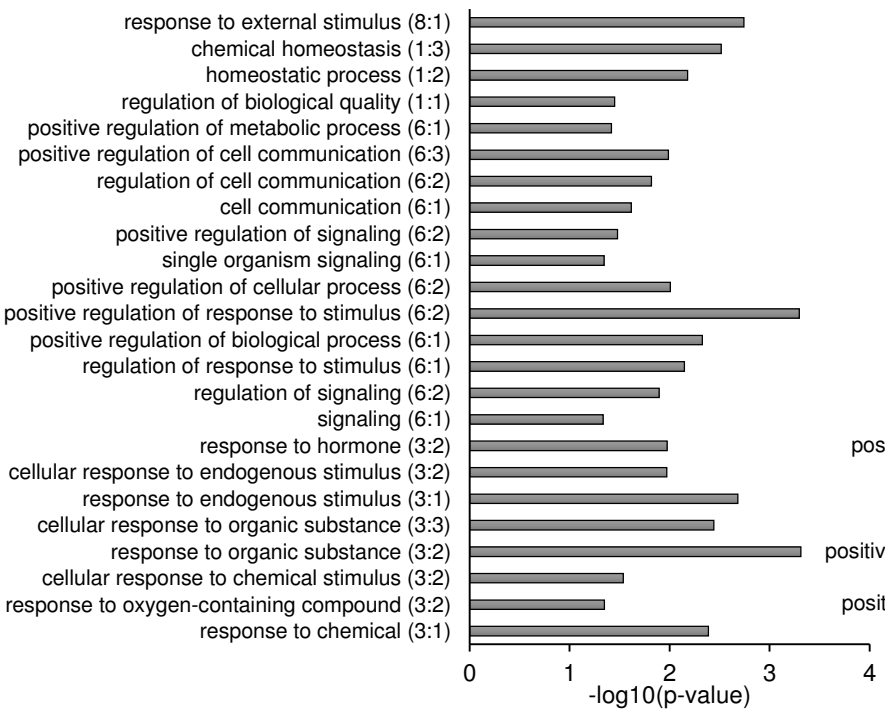


Figure S11

A



B

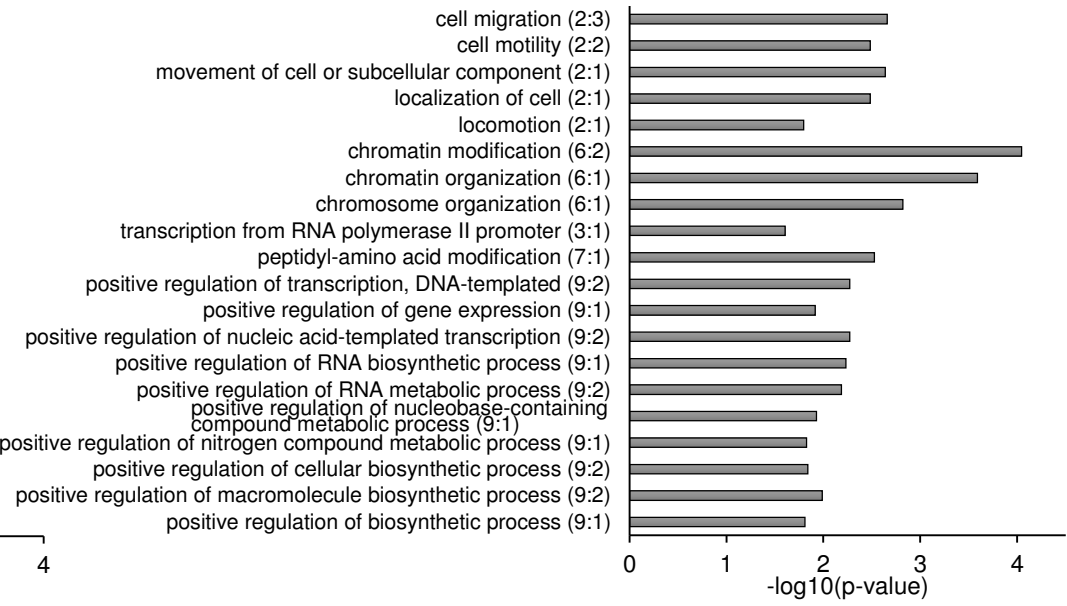
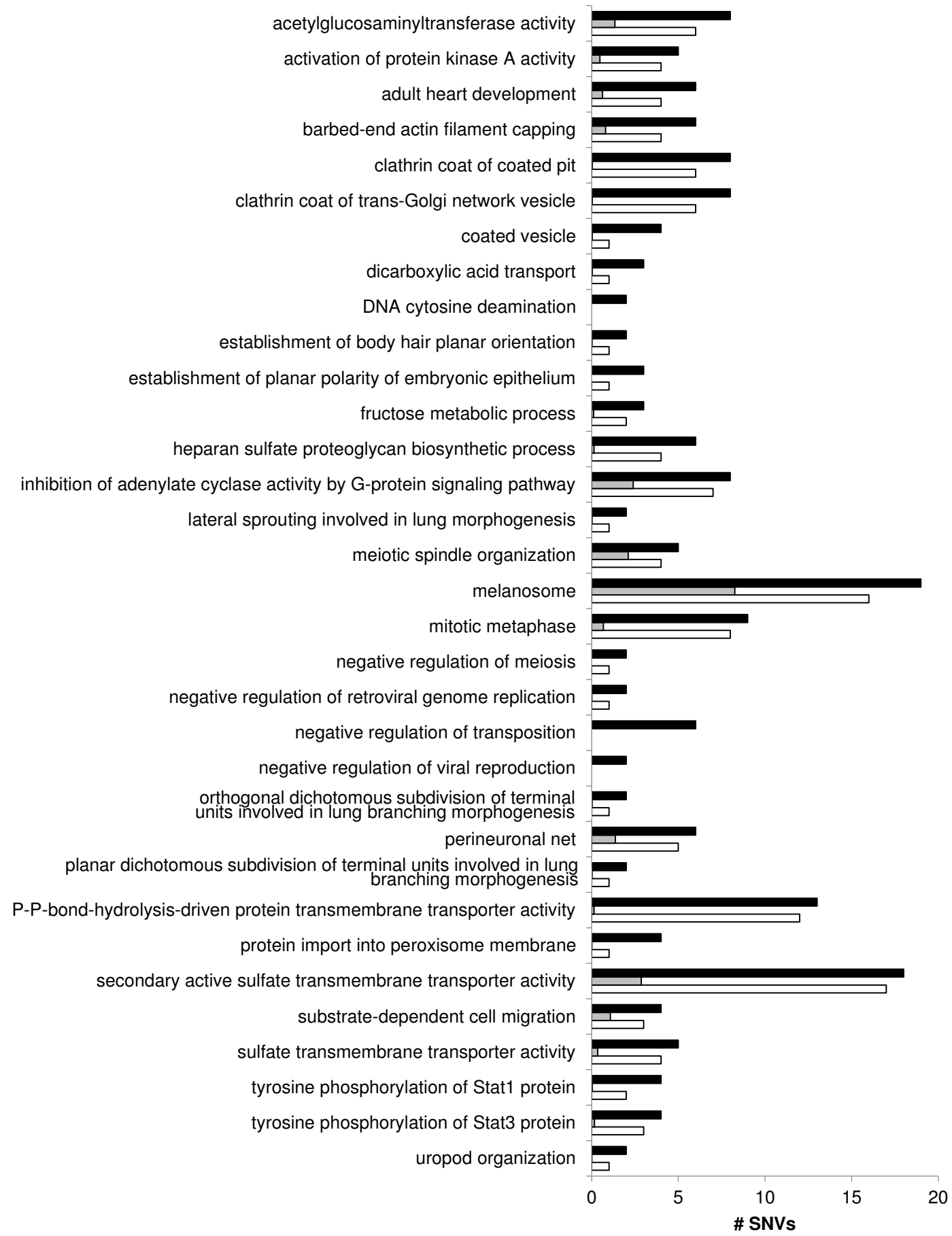


Figure S12

A



B

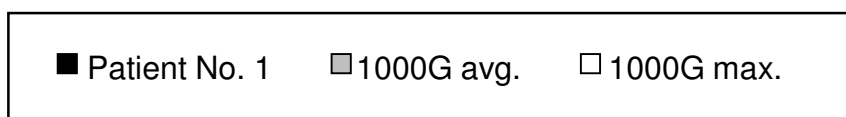
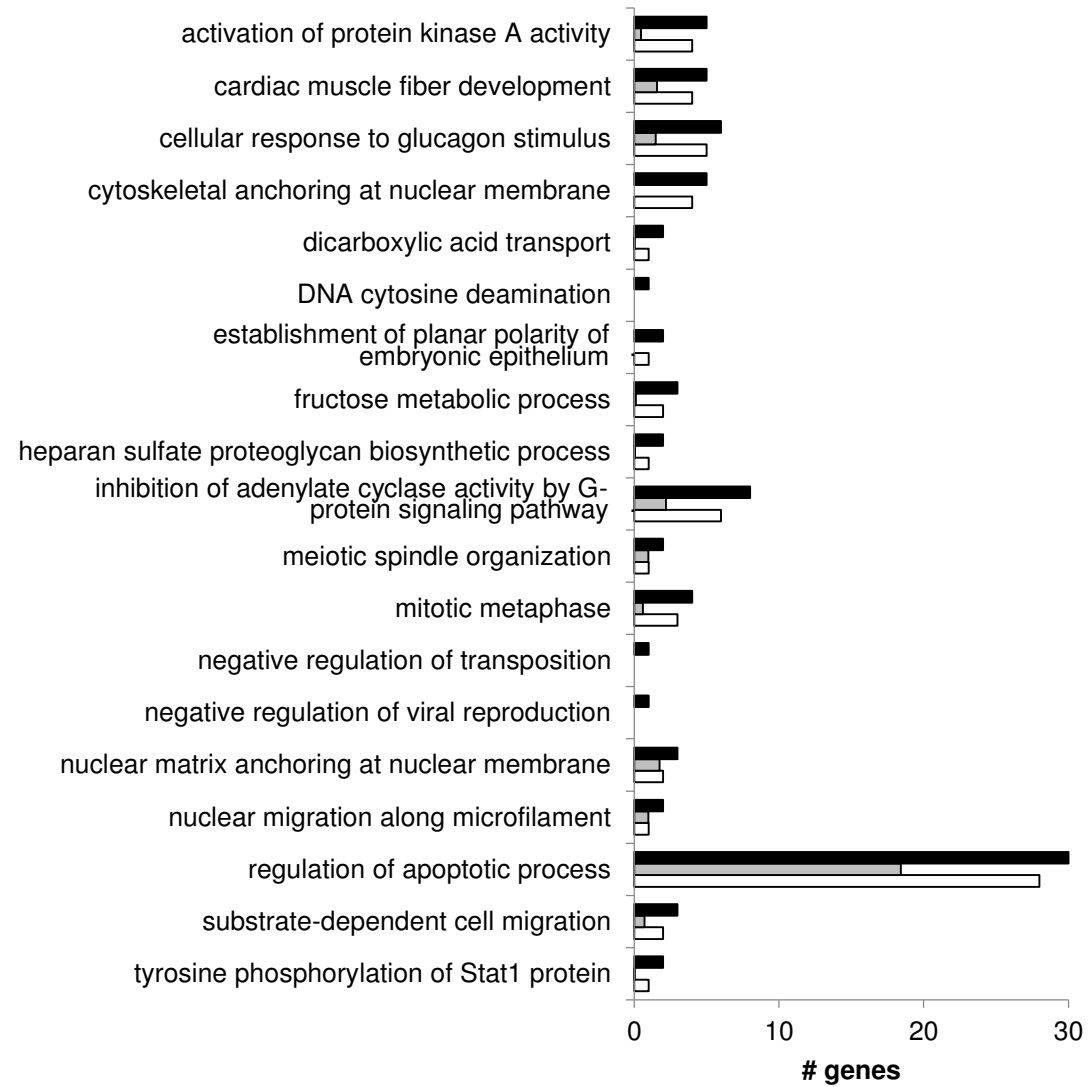


Figure S13

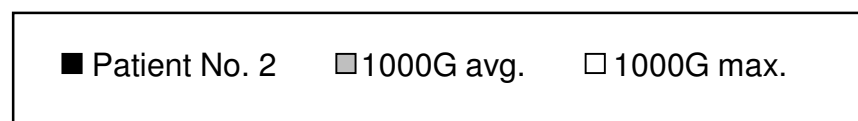
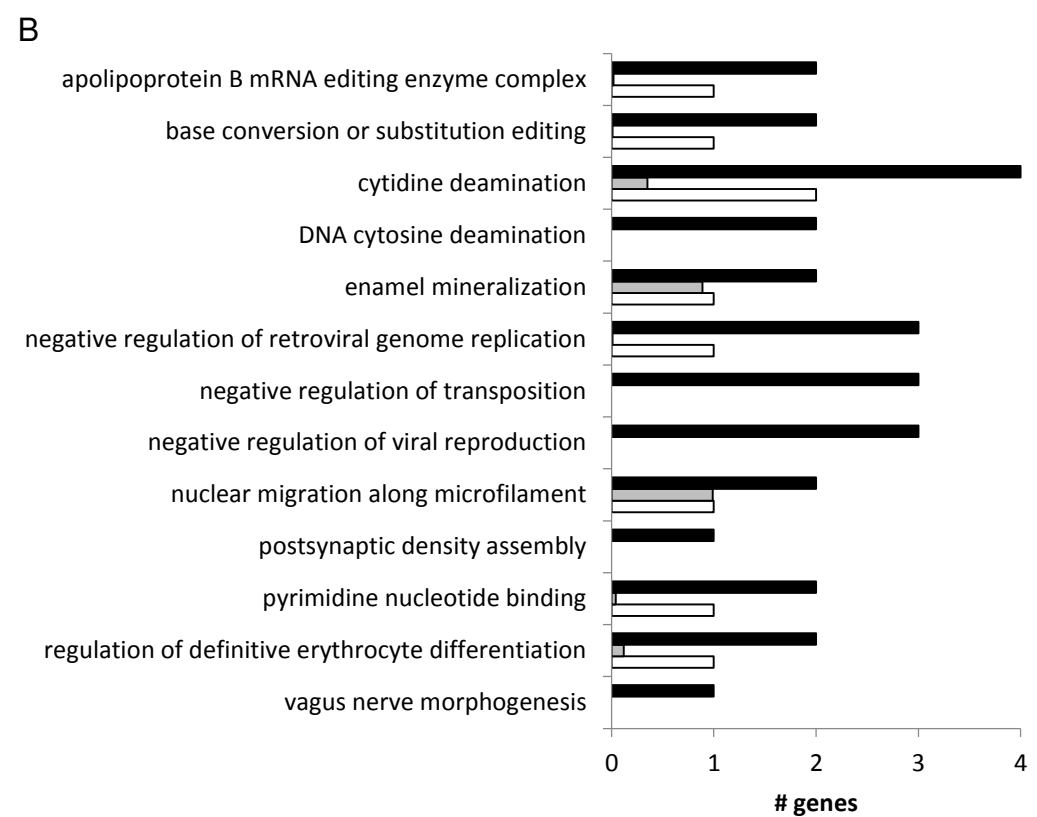
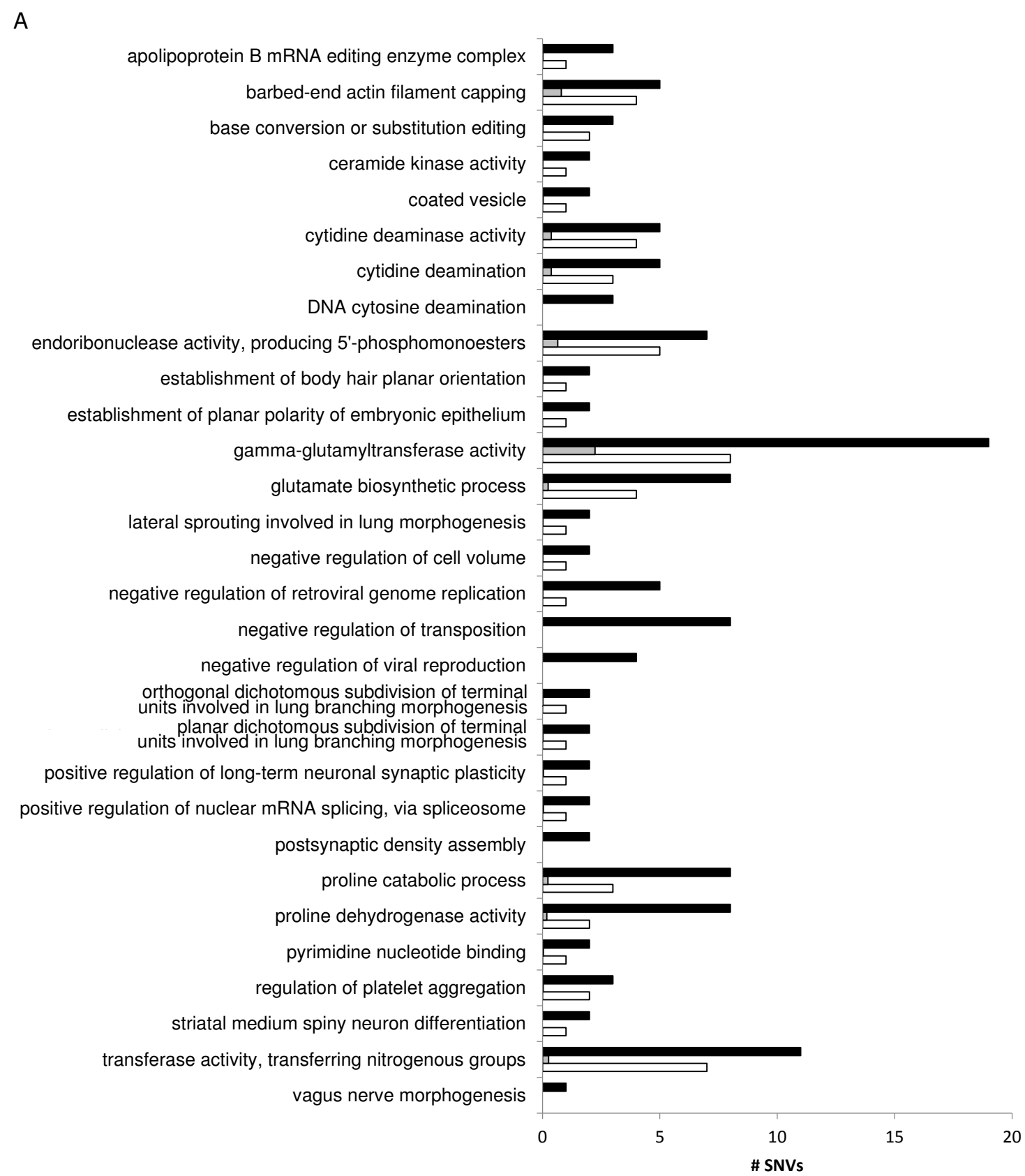


Figure S14

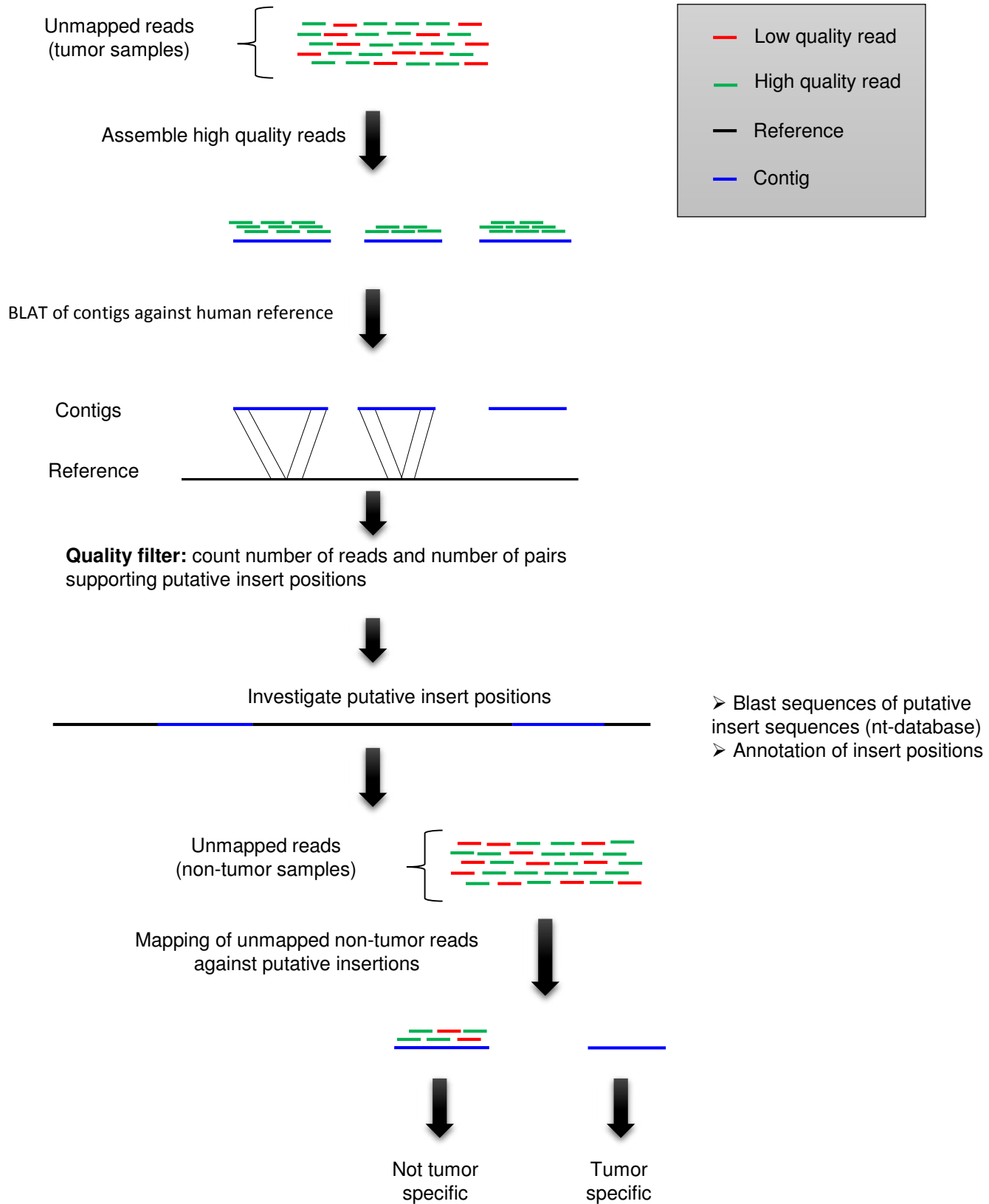


Figure S15

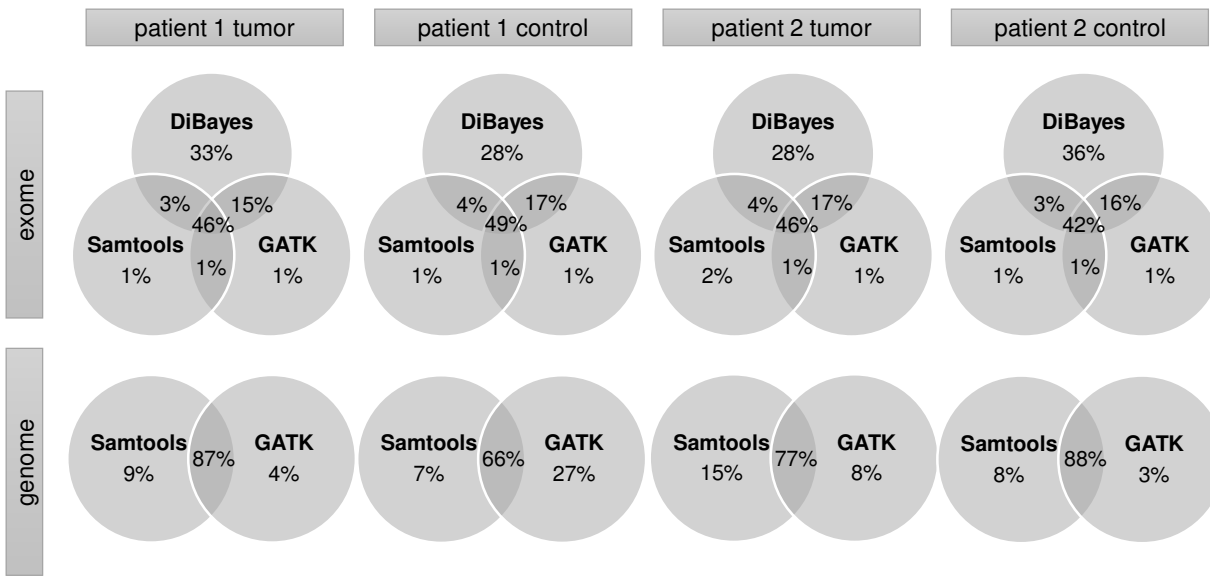


Figure S16

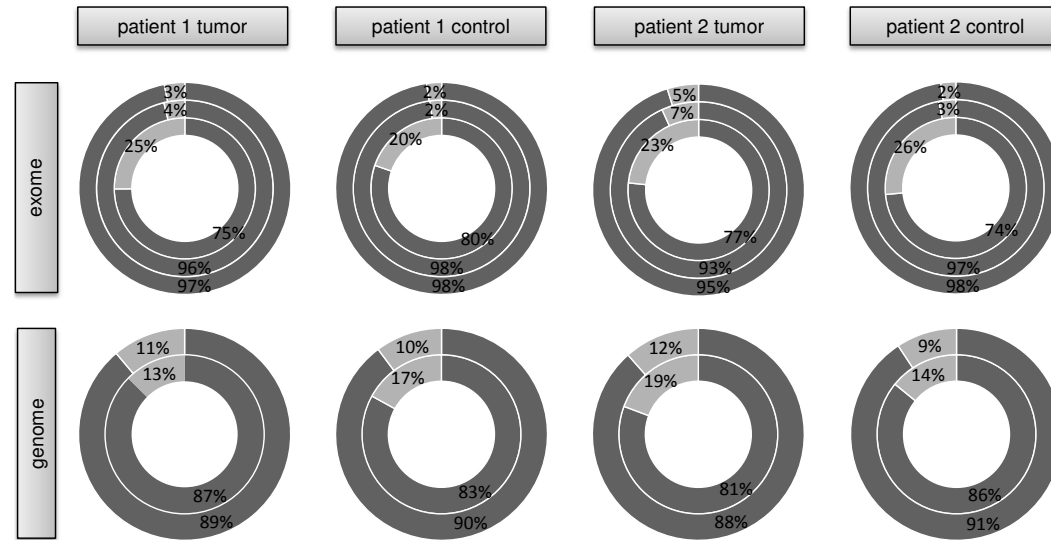
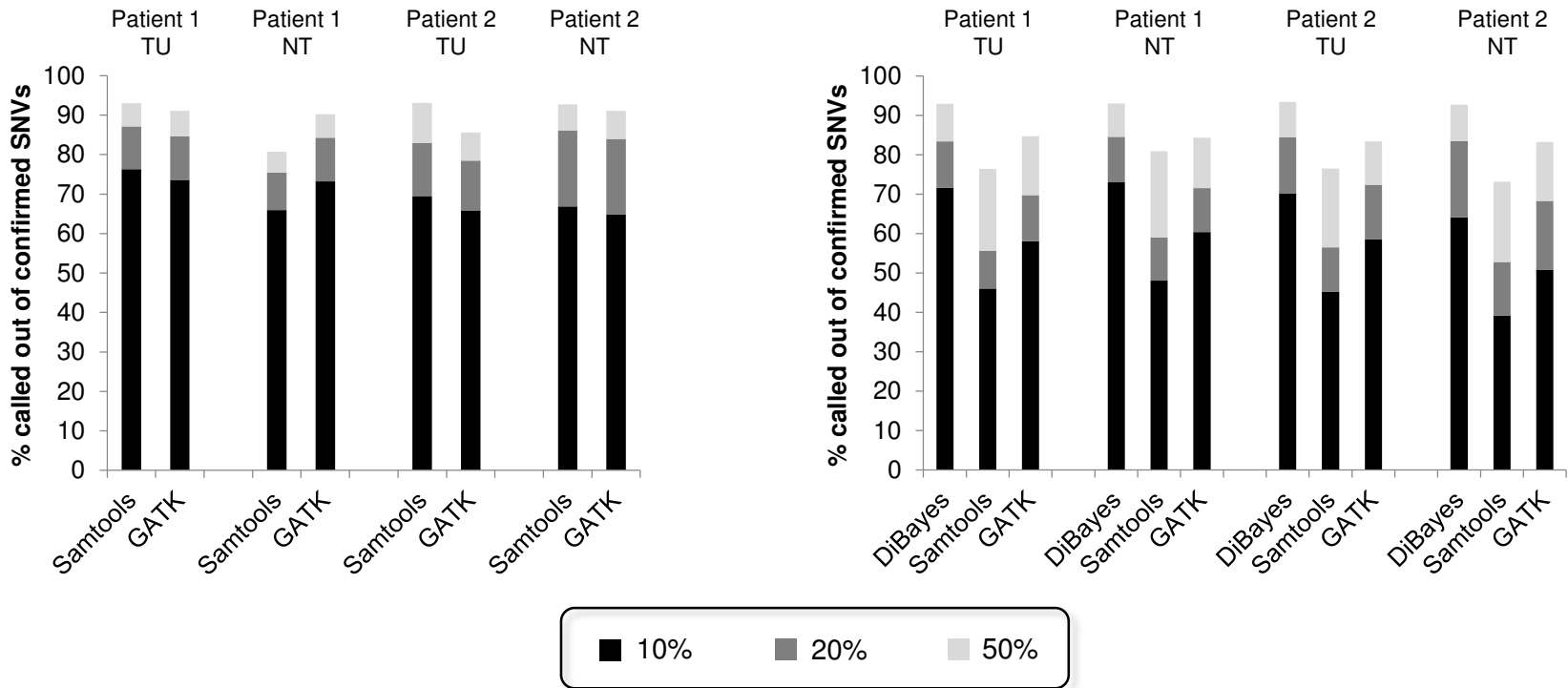


Figure S17



Number of confirmed SNVs (10%, 20%, 50%)

	Patient 1 TU	Patient 1 NT	Patient 2 TU	Patient 2 NT
WGS -Samtools	67426 / 63229 / 30967	65751 / 62894 / 29815	67969 / 64310 / 34095	78203 / 74402 / 34643
WGS - GATK	65010 / 61420 / 30324	72990 / 70246 / 33336	64420 / 60855 / 31352	75691 / 72530 / 34048
WES - DiBayes	63249 / 60504 / 30944	72779 / 70446 / 34944	68645 / 65459 / 34216	74895 / 72090 / 34653
WES - Samtools	40667 / 40363 / 25447	47894 / 49151 / 36259	44308 / 43826 / 28385	45759 / 45572 / 27353
WES - GATK	51276 / 50577 / 28209	60156 / 59624 / 31131	57220 / 56097 / 31468	59430 / 58950 / 31120

Figure S18

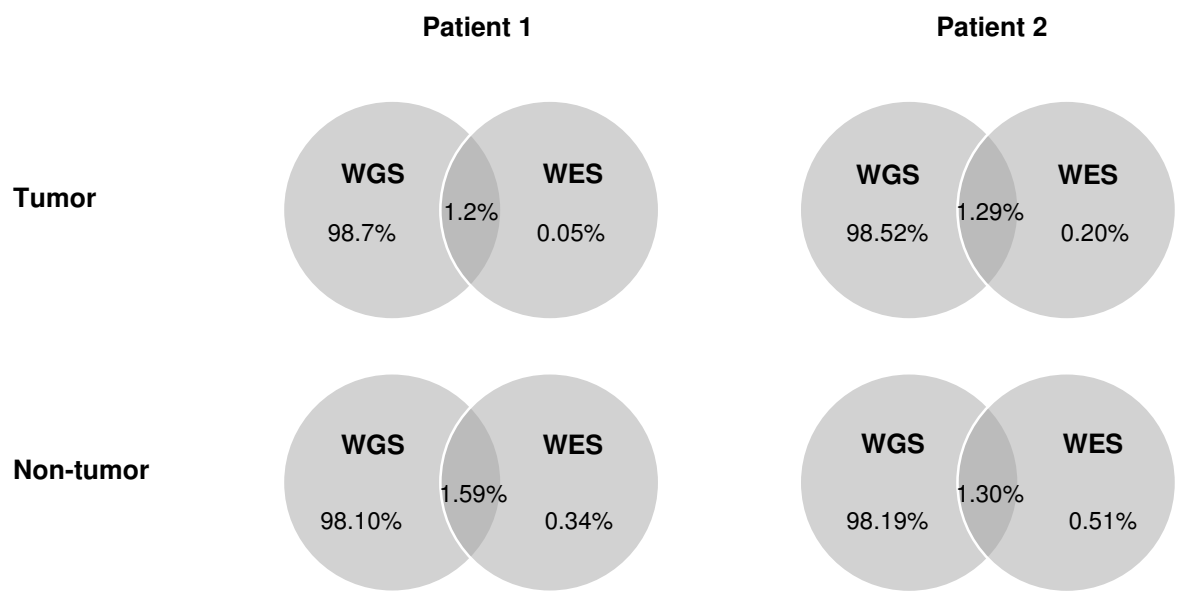


Figure S19

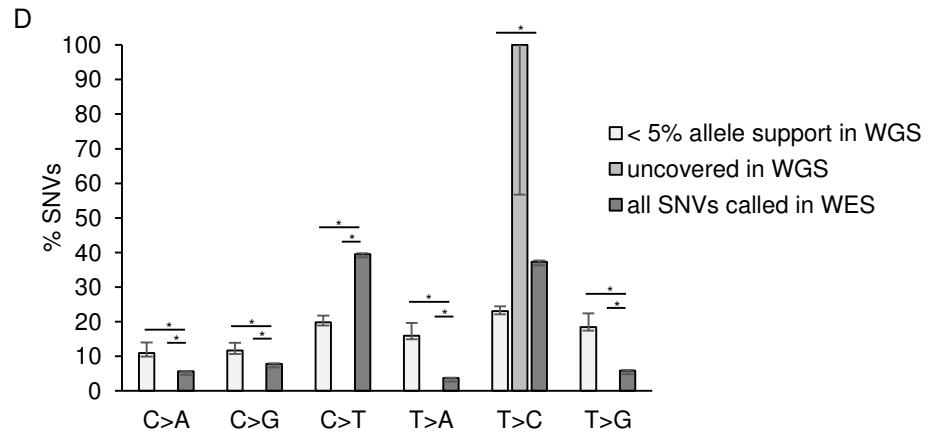
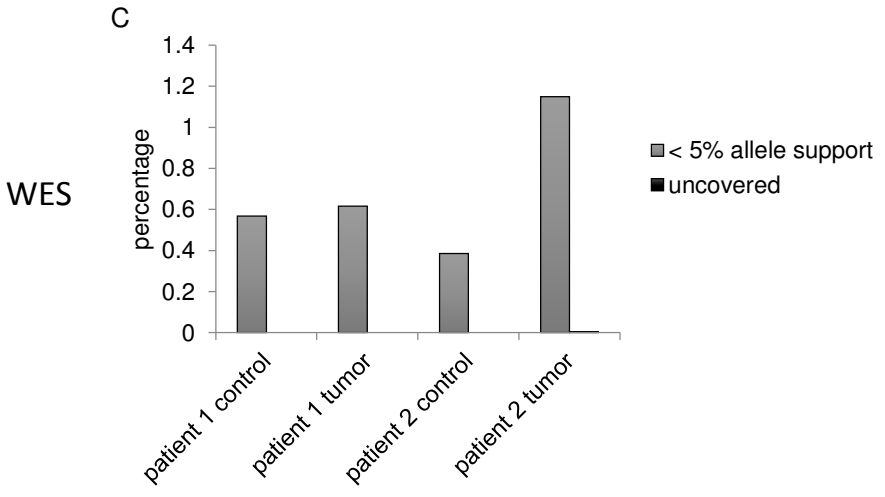
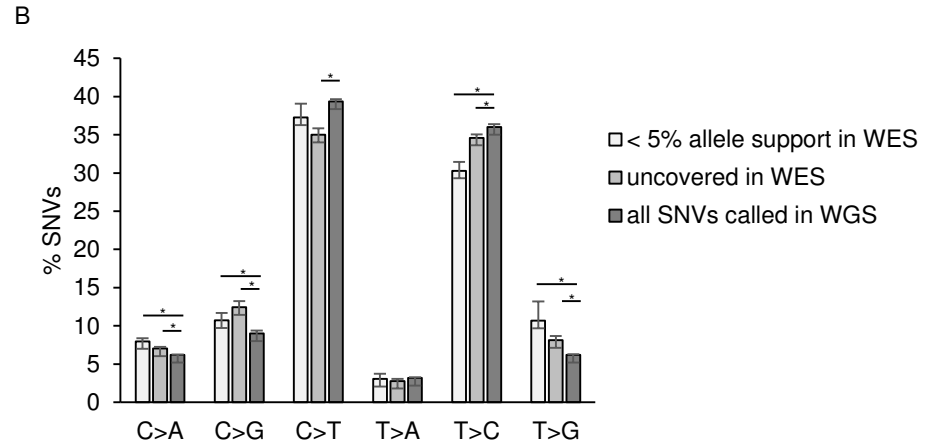
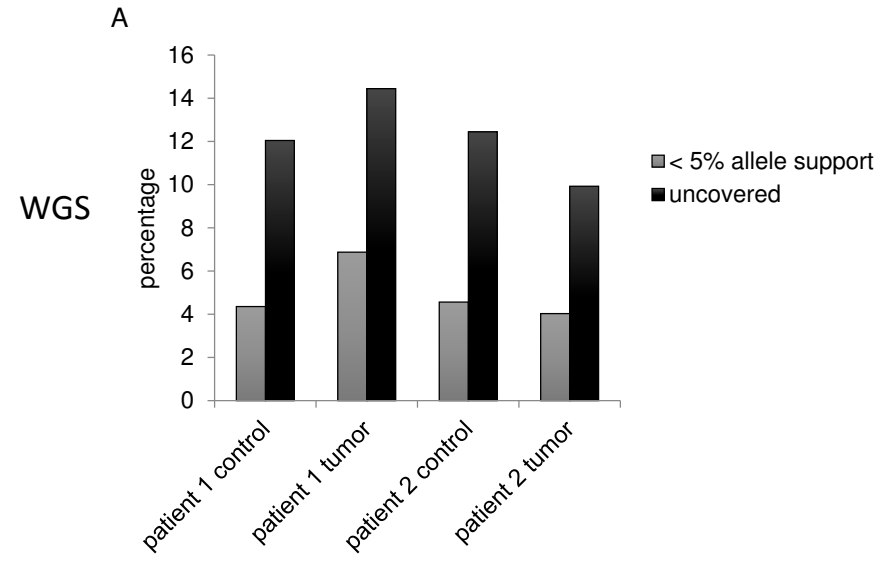
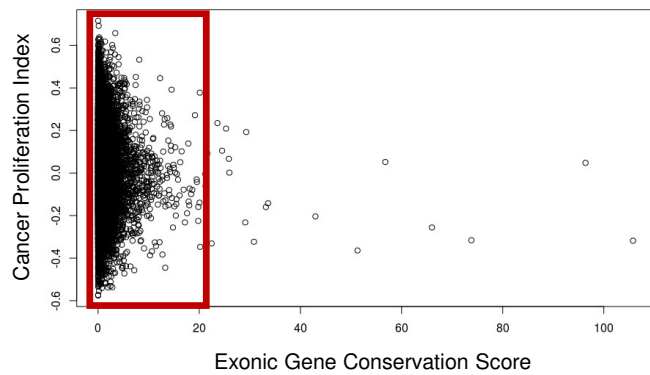
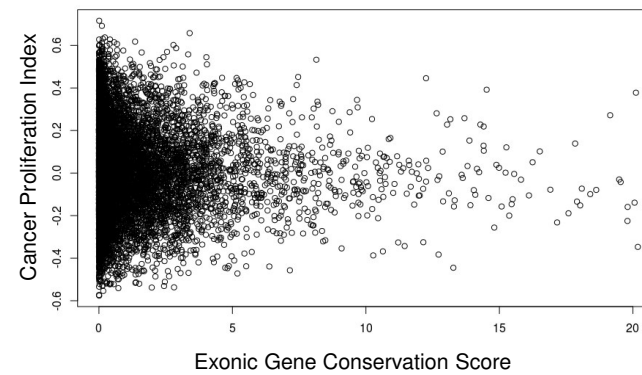


Figure S20

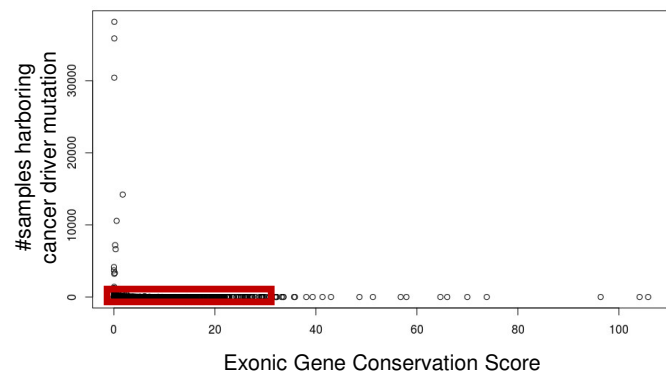
A



B



C



D

