

Online Data Supplement

Figure E1

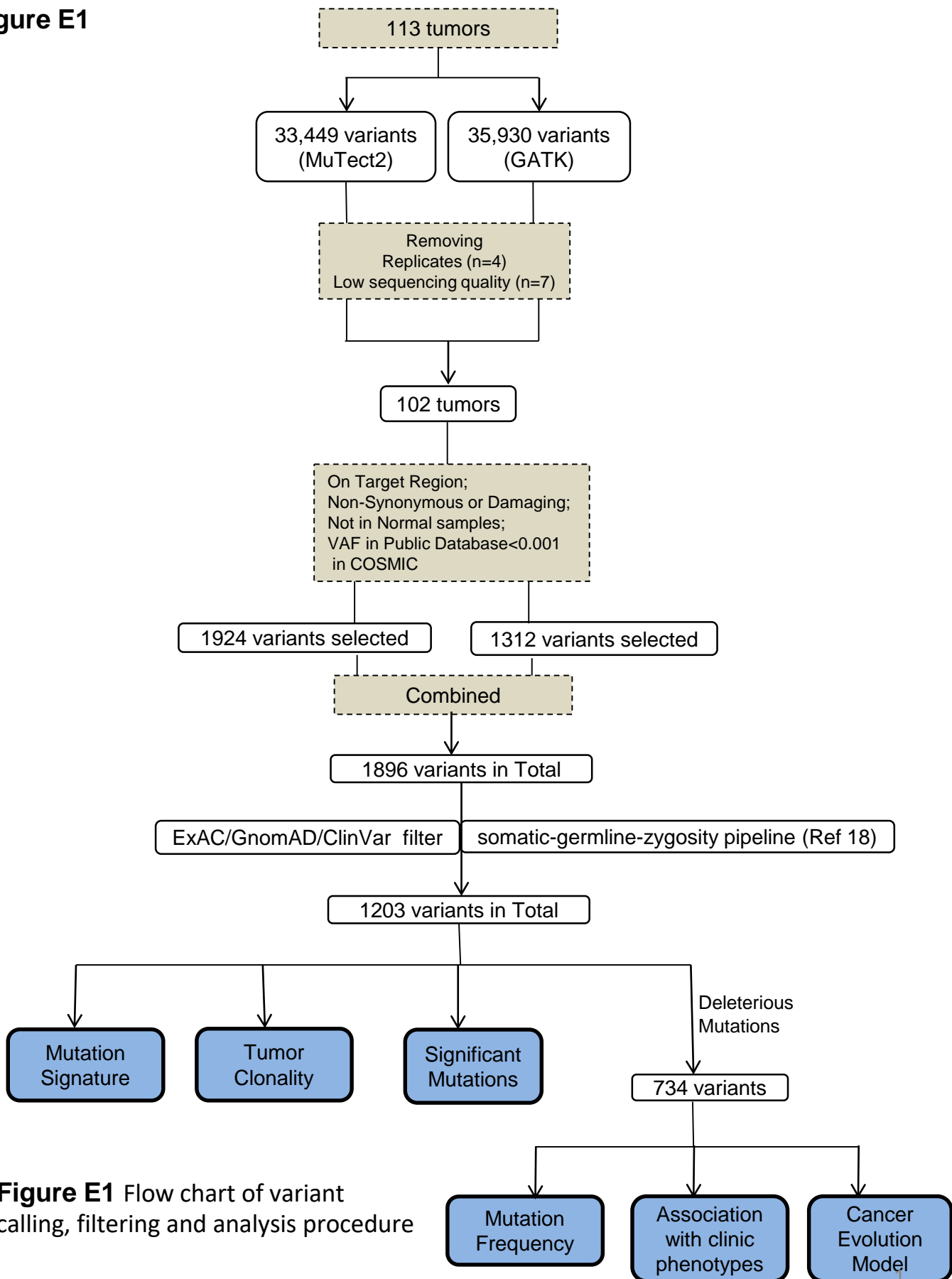


Figure E1 Flow chart of variant calling, filtering and analysis procedure

Figure E2. Somatic mutations and copy number in AIS, MIA and ADC. (A) The proportion of all mutation types in our cohort. (B) Nine genes harboring significant hot spot mutations were identified using the oncodriveCluster approach. (C) Frequencies of significantly amplified and deleted regions or genes identified using RUBIC tool. (D) Overview of top gene mutations and copy number alterations.

A

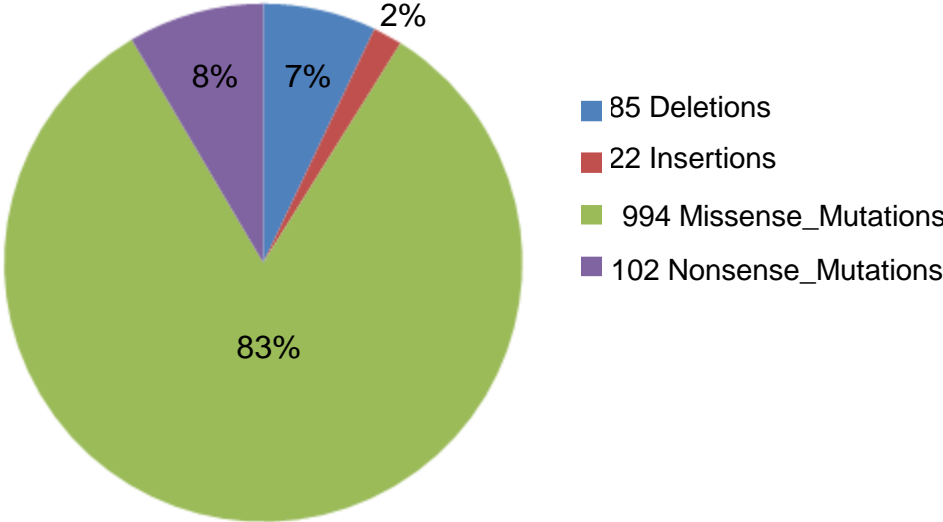
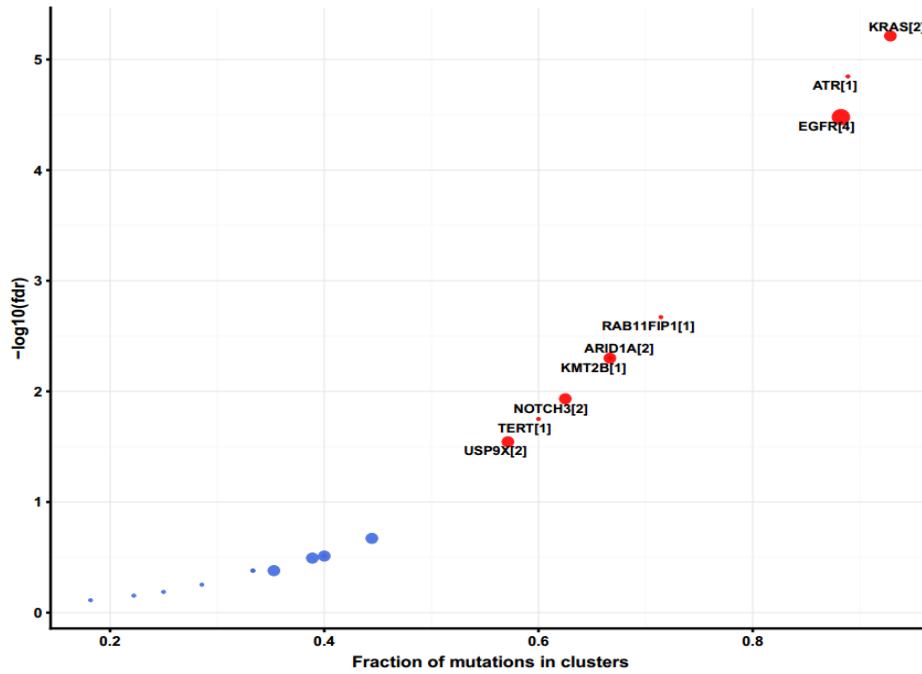


Figure E2

B



C

Top CNV in all

Altered in 83 (91.21%) of 91 samples.

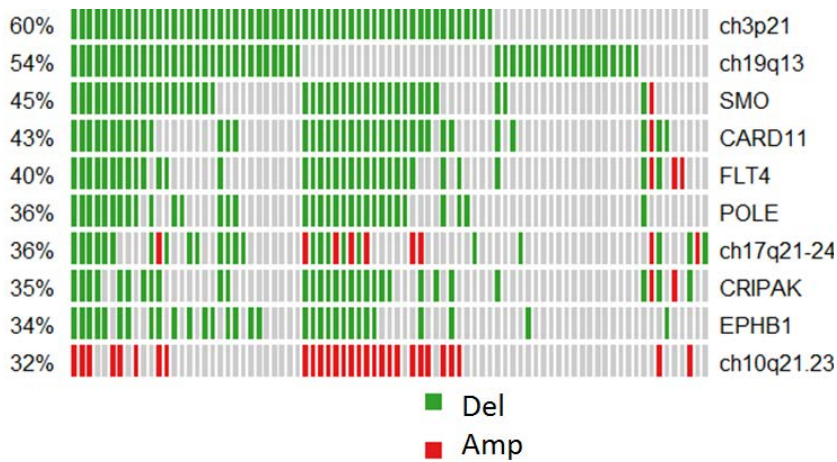


Figure E2

D

AIS

MIA

ADC

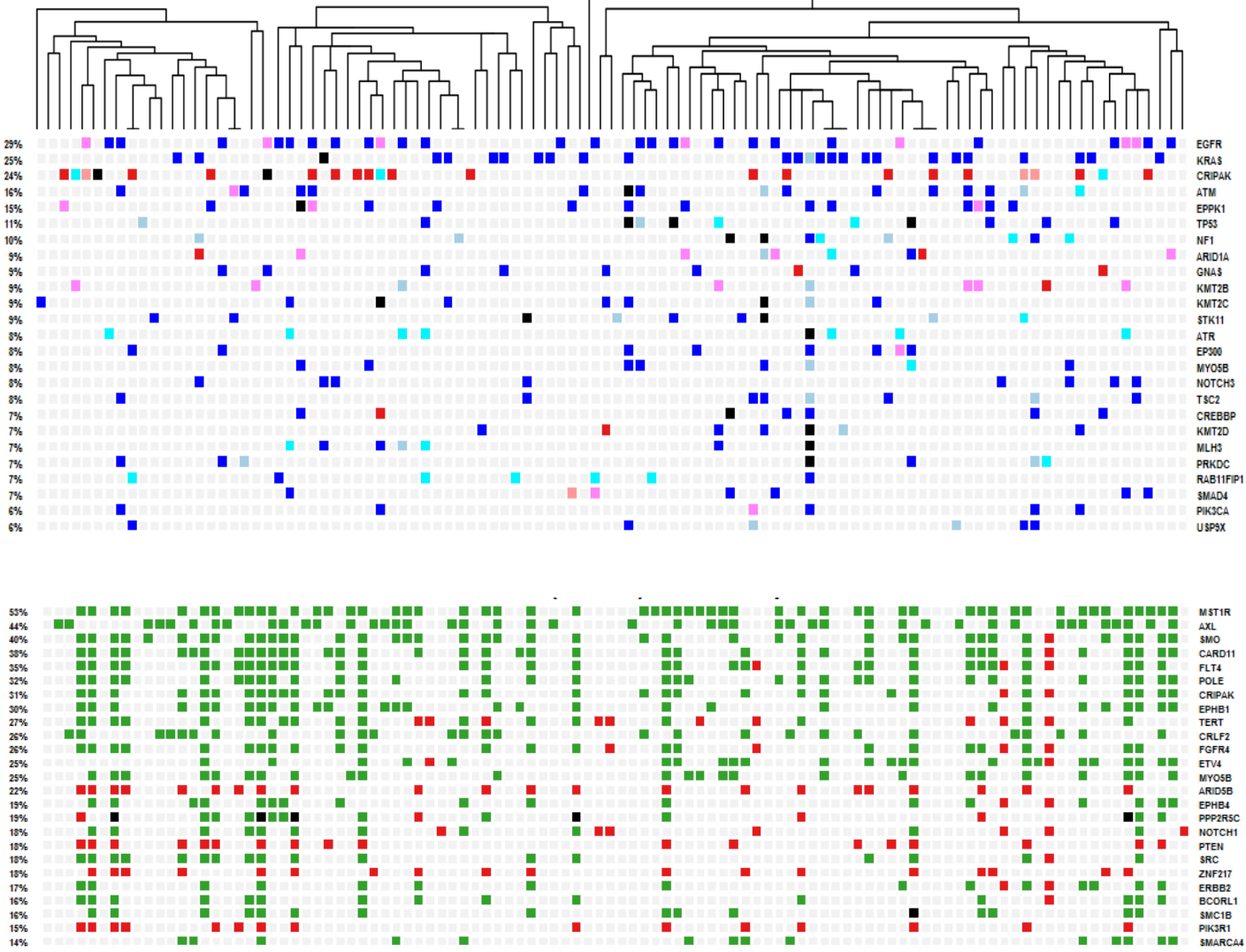


Figure E3. Mutational signature and tumor heterogeneity analysis. (A) Five mutation signatures identified by non-negative matrix factorization (NMF) approach on 96 trinucleotide mutation types across all tumors. Five mutation signatures include smoke exposure (signature 4), APOBEC activation (signature 2), aging (signature 1), mismatch repair (MMR) (signature 6) and unknown (signature 18). Each bar corresponds to the probability of observing a particular mutation in a trinucleotide context within each signature. **(B)** Relative contribution of each mutation signature in each tumor. **(C)** Mutation burden was significantly higher in tumors enriched in APOBEC and Sig18.unknown when compared to MMR signature enriched tumors. P values were calculated using negative binomial regression model. **(D)** Mutant-allele tumor heterogeneity (MATH) score tended to be positively associated with mutation burden ($r=0.19$, $P=0.07$) and was not associated with smoking status, age, stage, gender or tumor size. **(E)** Two representative cancer cell fraction (CCF) density plot using Sciclone showed one multi-clonal tumor (T45) and one mono-clonal tumor (T40). **(F)** The proportions of clonal/subclonal mutations were not significantly associated with the histology. **(G)** Tumors harboring any EGFR mutations or only subclonal EGFR mutations were not associated with overall survival. P values shown were calculated using Kaplan-Meier estimate. **(H)** VAF plot for EGFR mutations. The majority of subclonal mutations of EGFR were p.L858R.

Figure E3

A

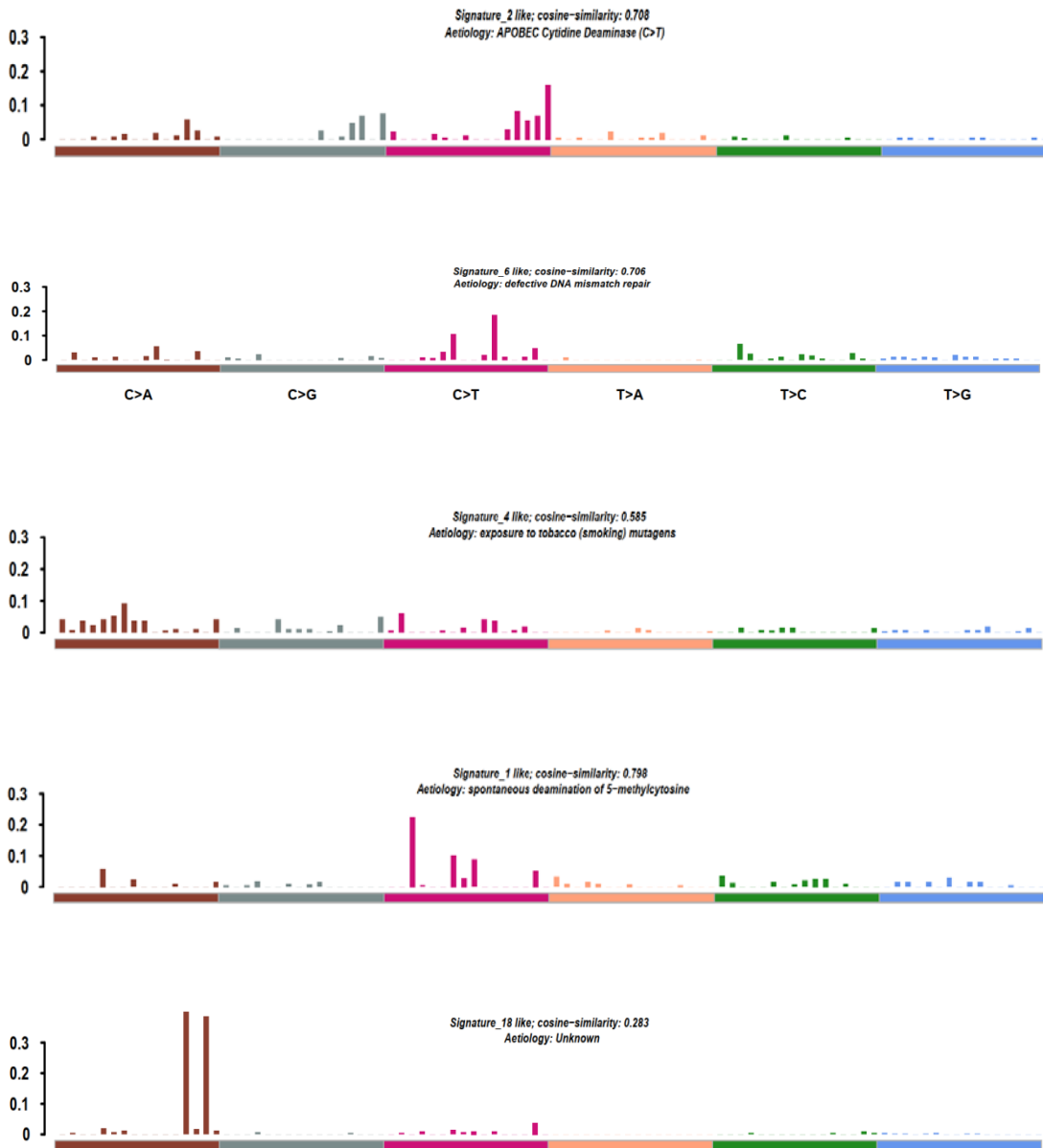
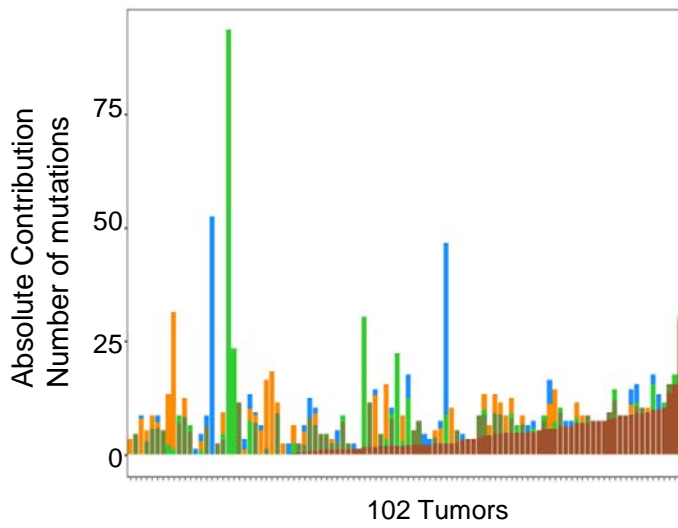
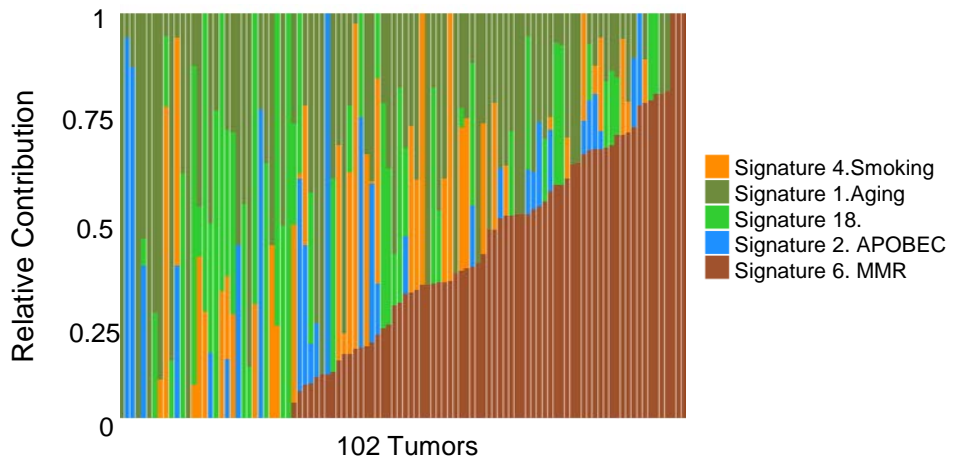


Figure E3

B



C

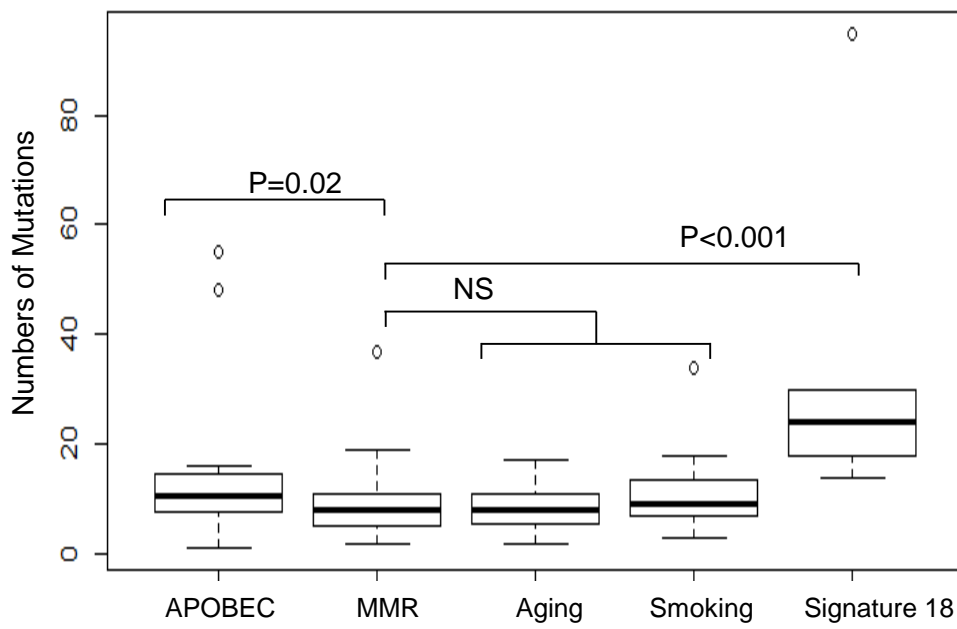


Figure E3

D

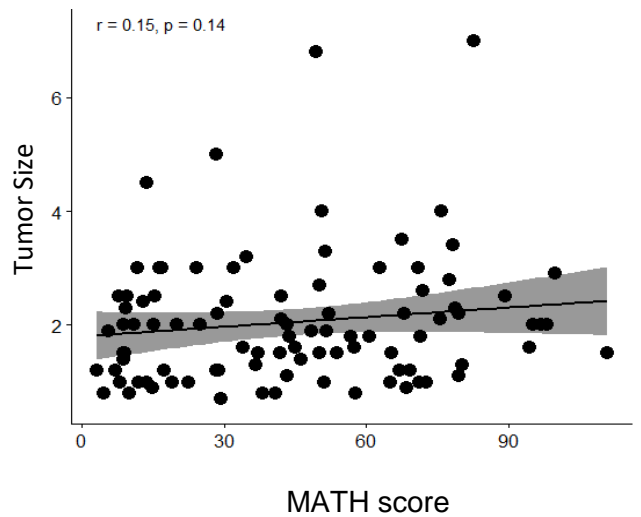
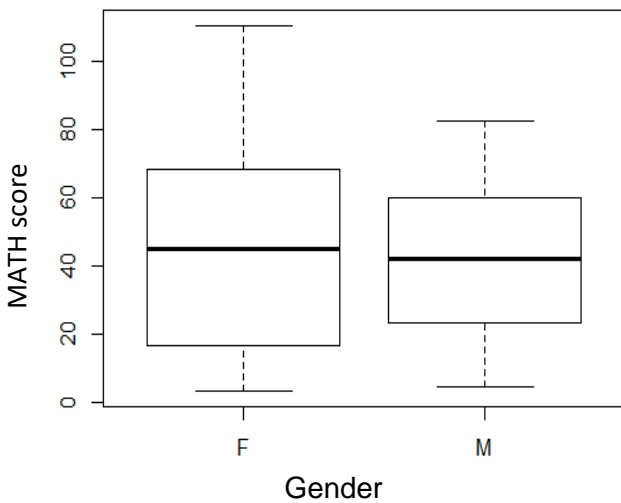
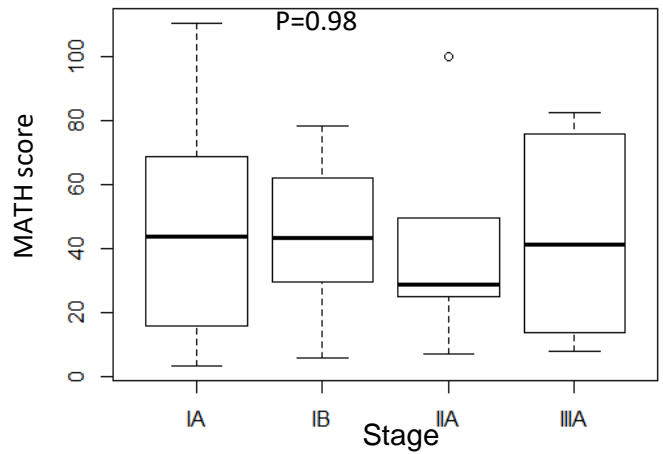
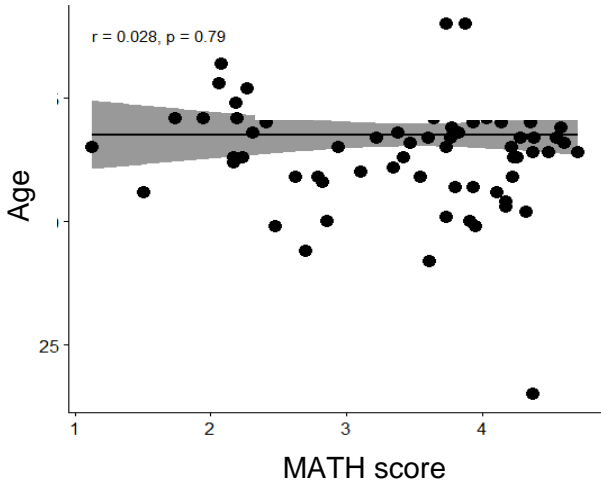
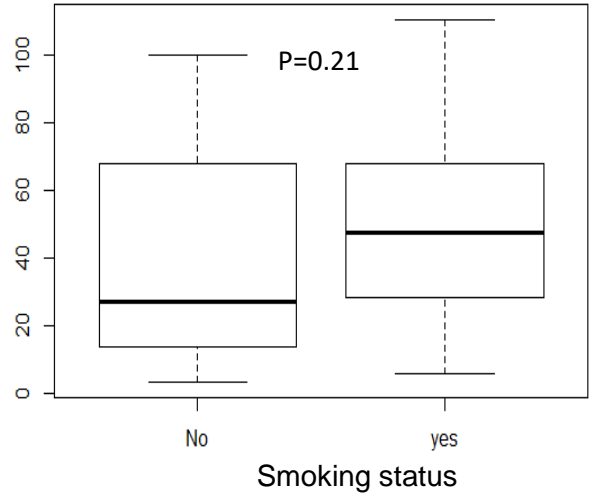
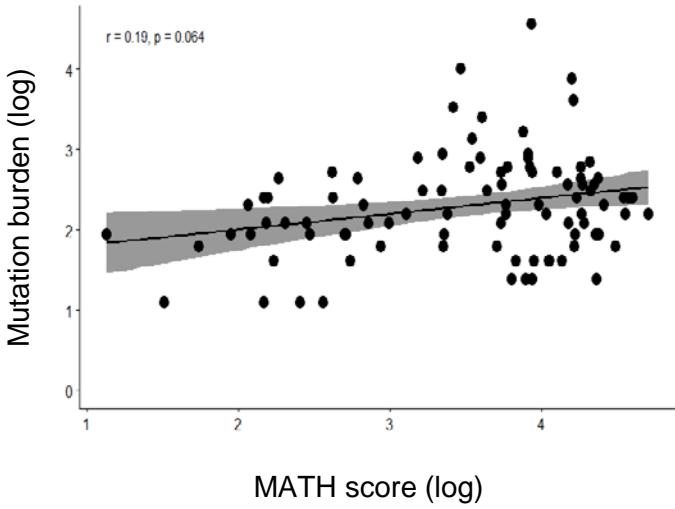
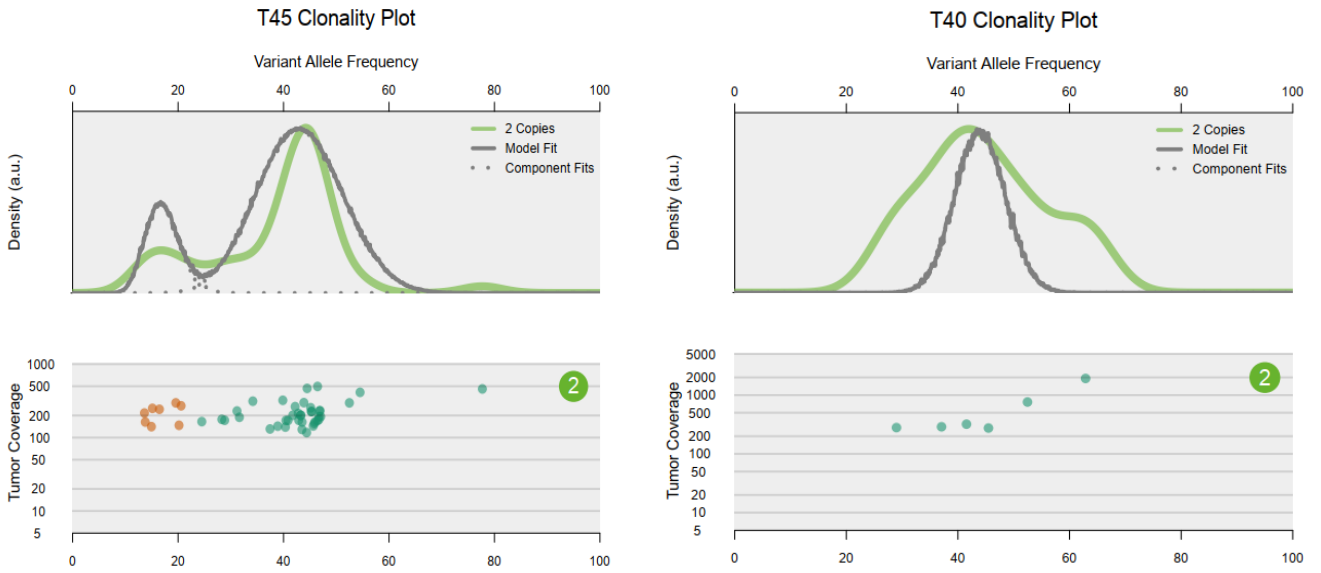


Figure E3

E



F

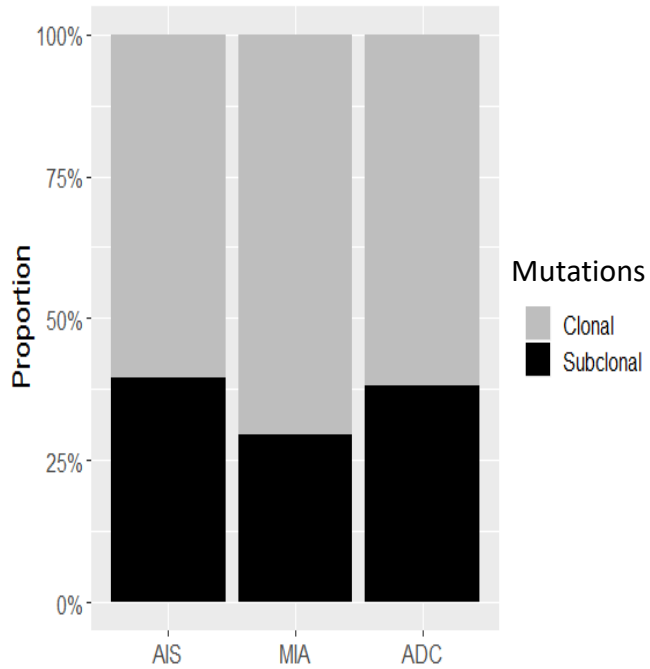


Figure E3

G

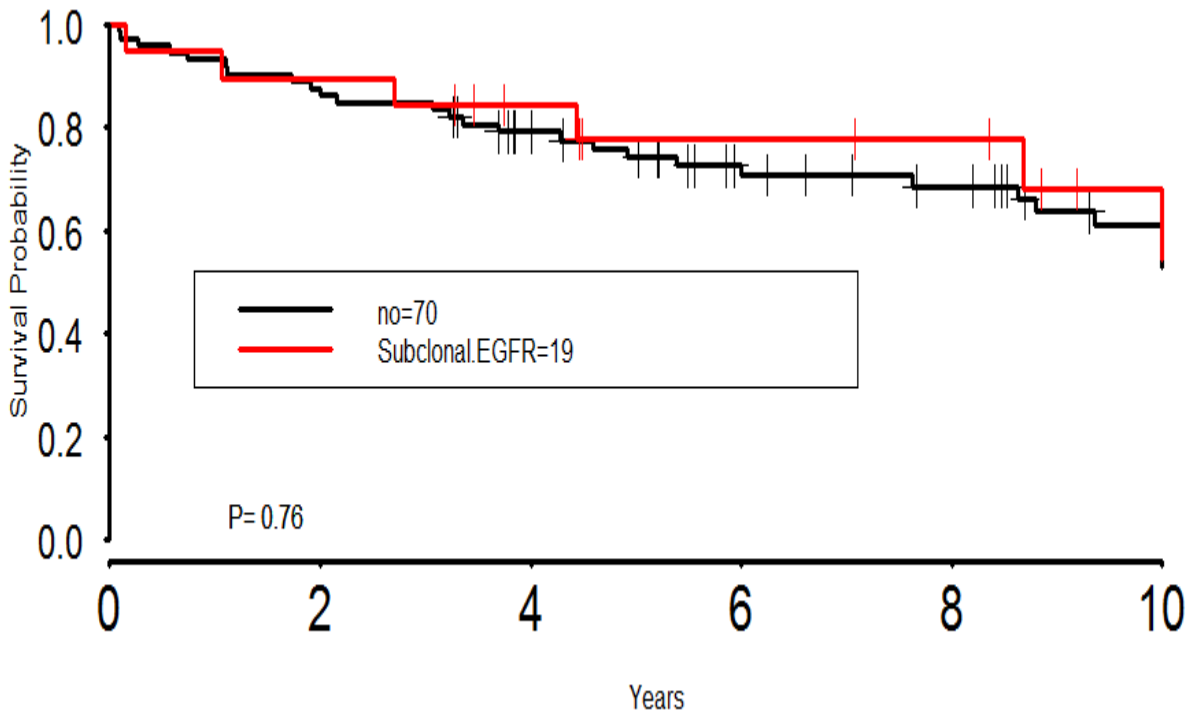
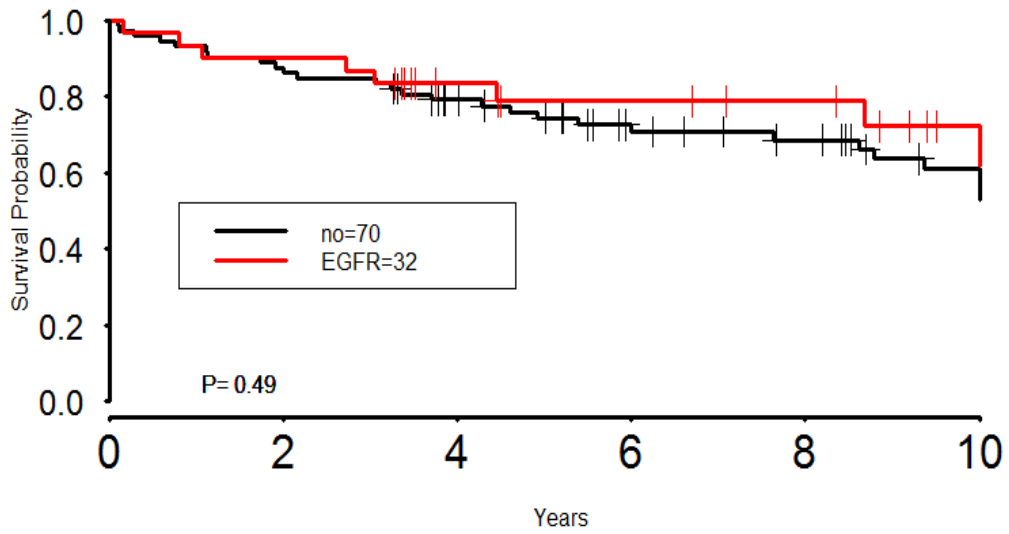


Figure E3

H

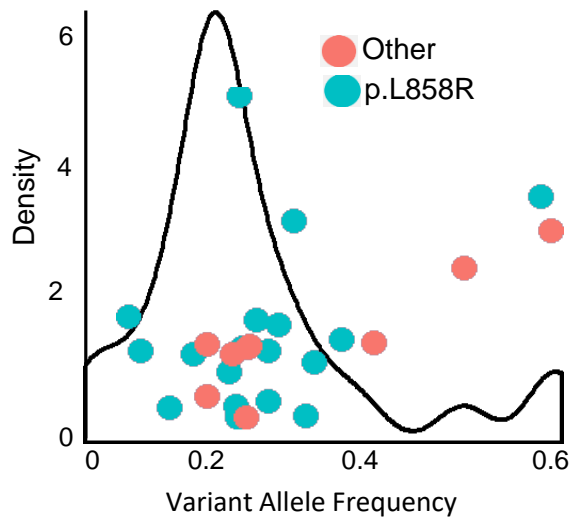
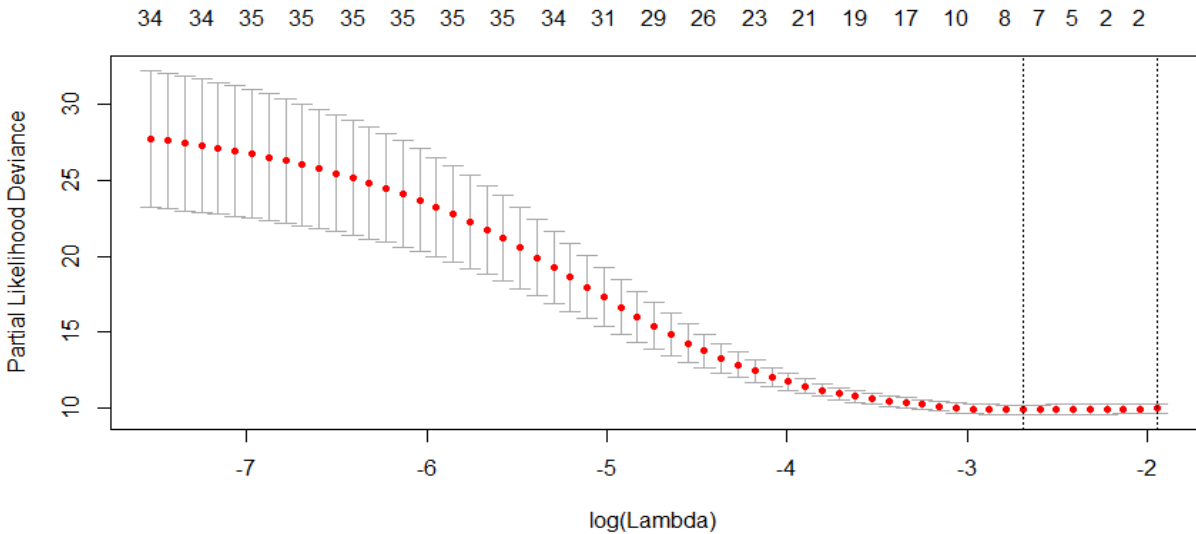
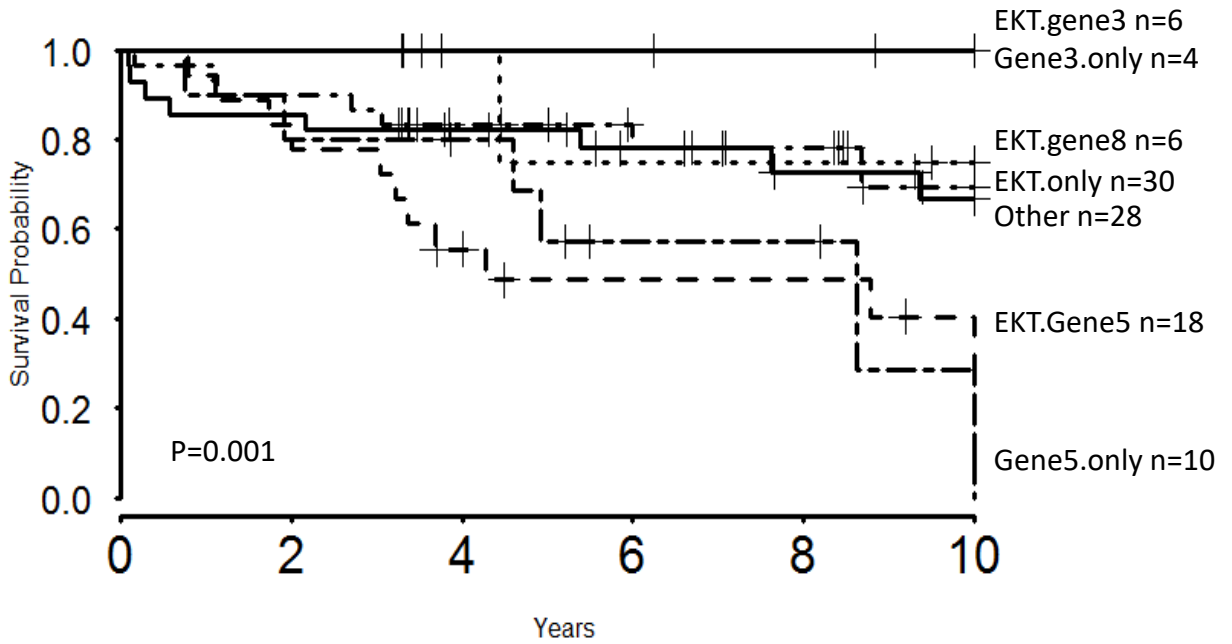


Figure E4. Association between gene mutations and 10-year overall survival. **(A)** Eight mutated genes selected by using an elastic-net penalized regression modeling (Glmnet R package). **(B)** Kaplan-Meier (KM) plot showed three patient groups had distinct survival outcomes. Group 1 includes patients having 5 gene (Gene5) mutation set and patients having concurrent EKT mutations (EGFR, KRAS and TP53); group2 includes patients having any of EKT mutations, patients having concurrent EKT with Gene8 mutations and patients having none of EKT or Gene8 mutations; group3 includes those patients having any of Gene3 mutations only or concurrent EKT with Gene3 mutations. **(C)** Group1 tumors harboring Gene5 only mutations or harboring concurrent EKT with Gene5 mutations were significantly associated with worse overall survival, after adjusting for age, gender, LVI and histology. **(D)** Tumors harboring any of Gene5 mutations, but not having any of EKT or Gene3 mutations were significantly associated with worse overall survival, after adjusting for age, gender, LVI or histology.

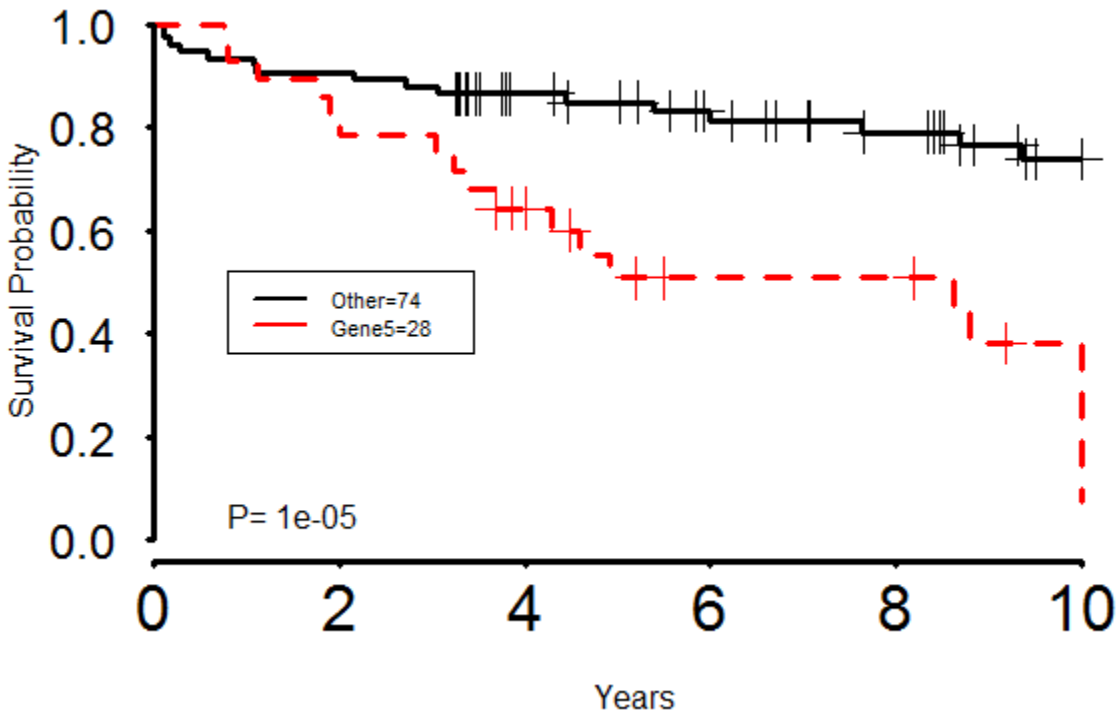
A.



B



C



D

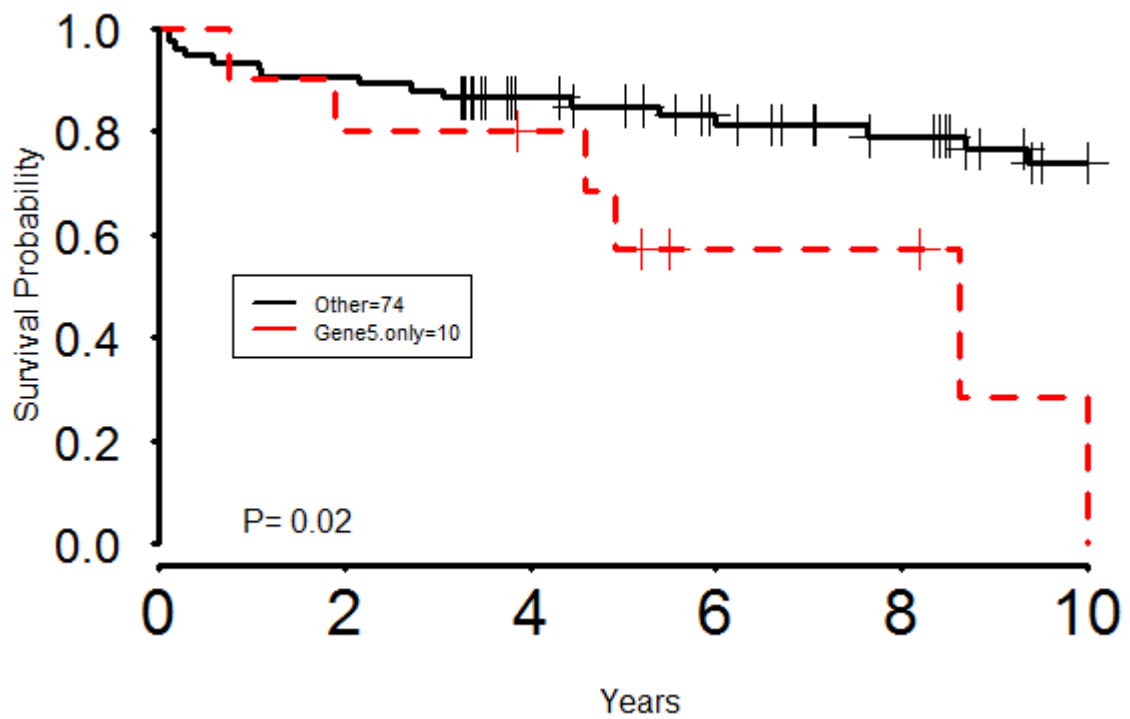


Figure E5. Cancer evolution model in early ADC detected by CAPRI algorithm. **(A)** Figure E5A has all tumors harboring a total of 59 significantly mutated genes plus four significantly altered CNVs with confidence shown as edge labels. The first, the second and third are p-values for temporal priority, probability raising and hypergeometric test. The last represents the relation confidence estimated with 100 non-parametric bootstrap (npb score) iterations. Red p-values are above the minimum significance threshold of 0.05. See Methods in the online supplement and Figure 4A in the Main Text for an interpretation of this model. **(B)** and **(C)** showed cancer evolution model in AIS/MIA combined and ADC using CAPRI algorithm. Only nonparametric bootstrap scores were shown. **(D)** Significantly mutually exclusive mutations between EGFR and KRAS ($P=0.0005$), EGFR and NF1 ($P=0.04$), and co-occurring mutations between EGFR and ATR ($P=0.0007$), EGFR and SMAD4 ($P=0.02$), STK11 and KEAP1 ($P=0.0001$). P values were calculated using Fisher exact test.

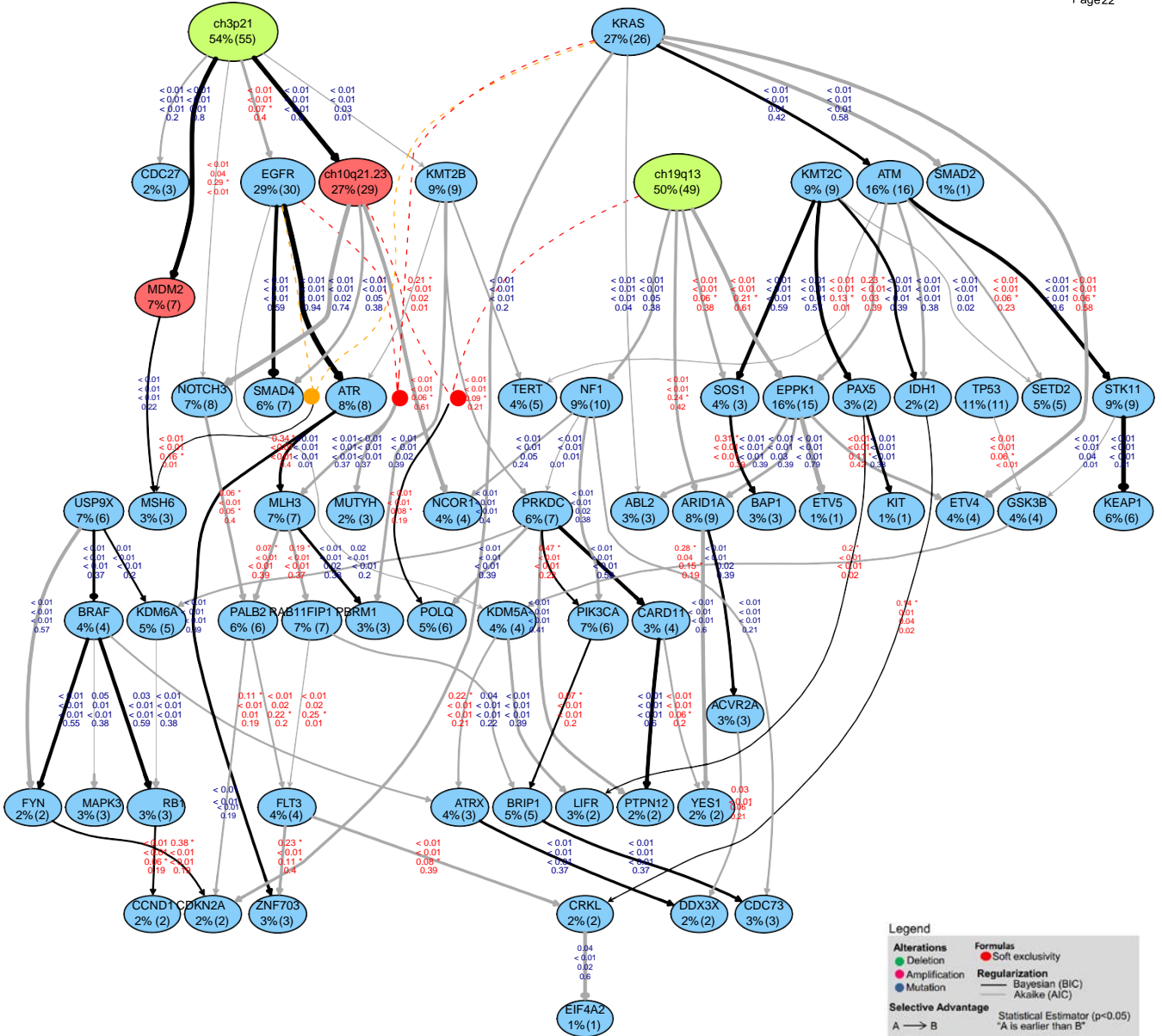


Figure E5B

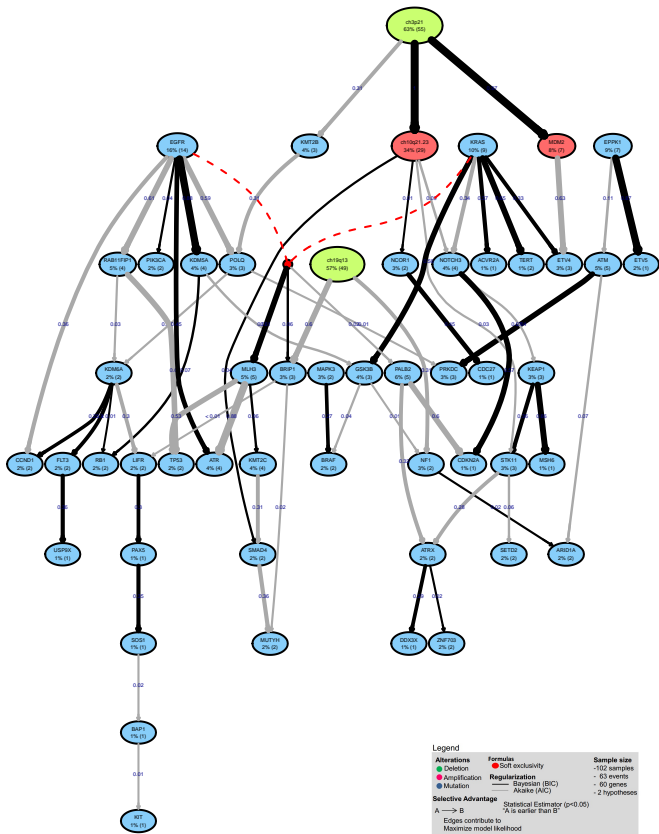


Figure E5C

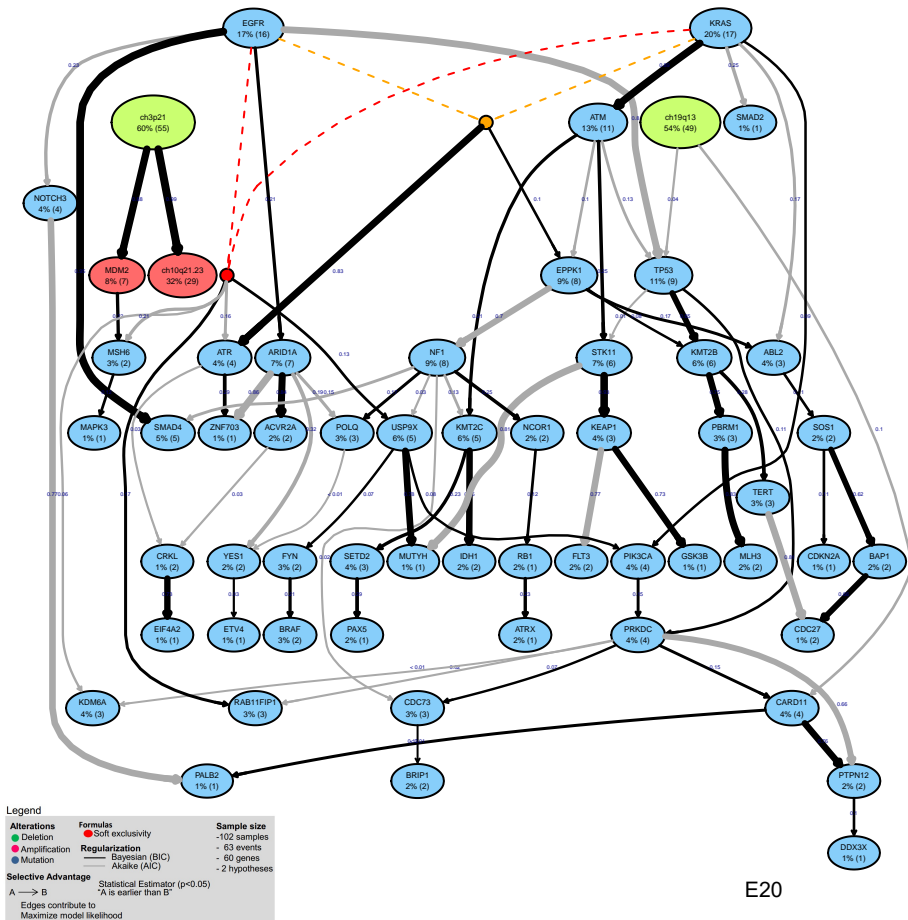


Figure E5

D

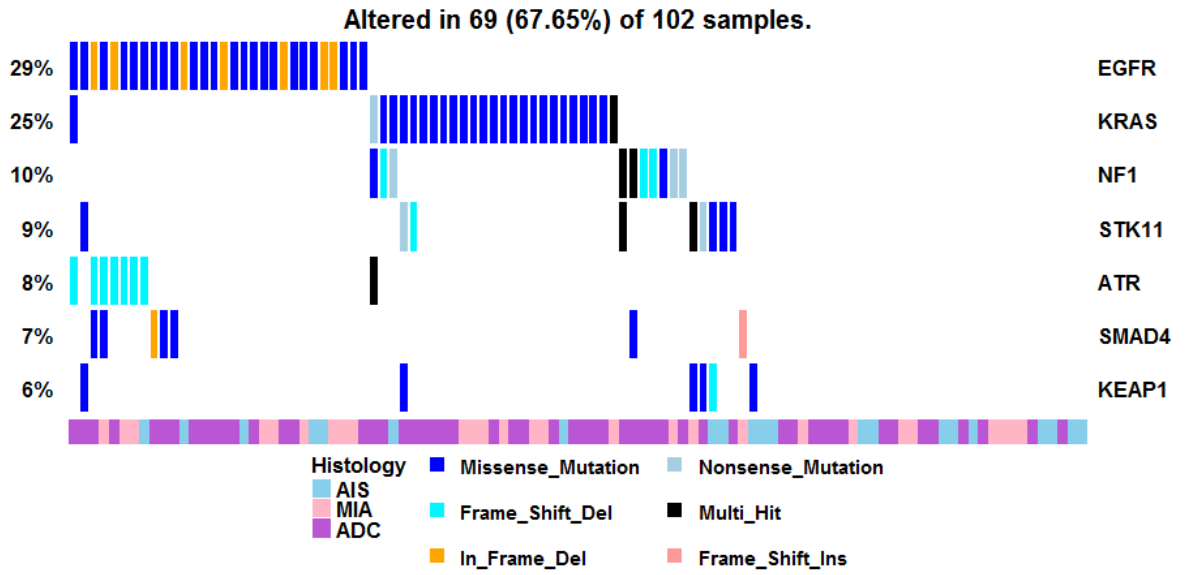


Figure E6

Figure E6. Non-supervised cluster analysis using genes that were mutated in greater than 5% tumors

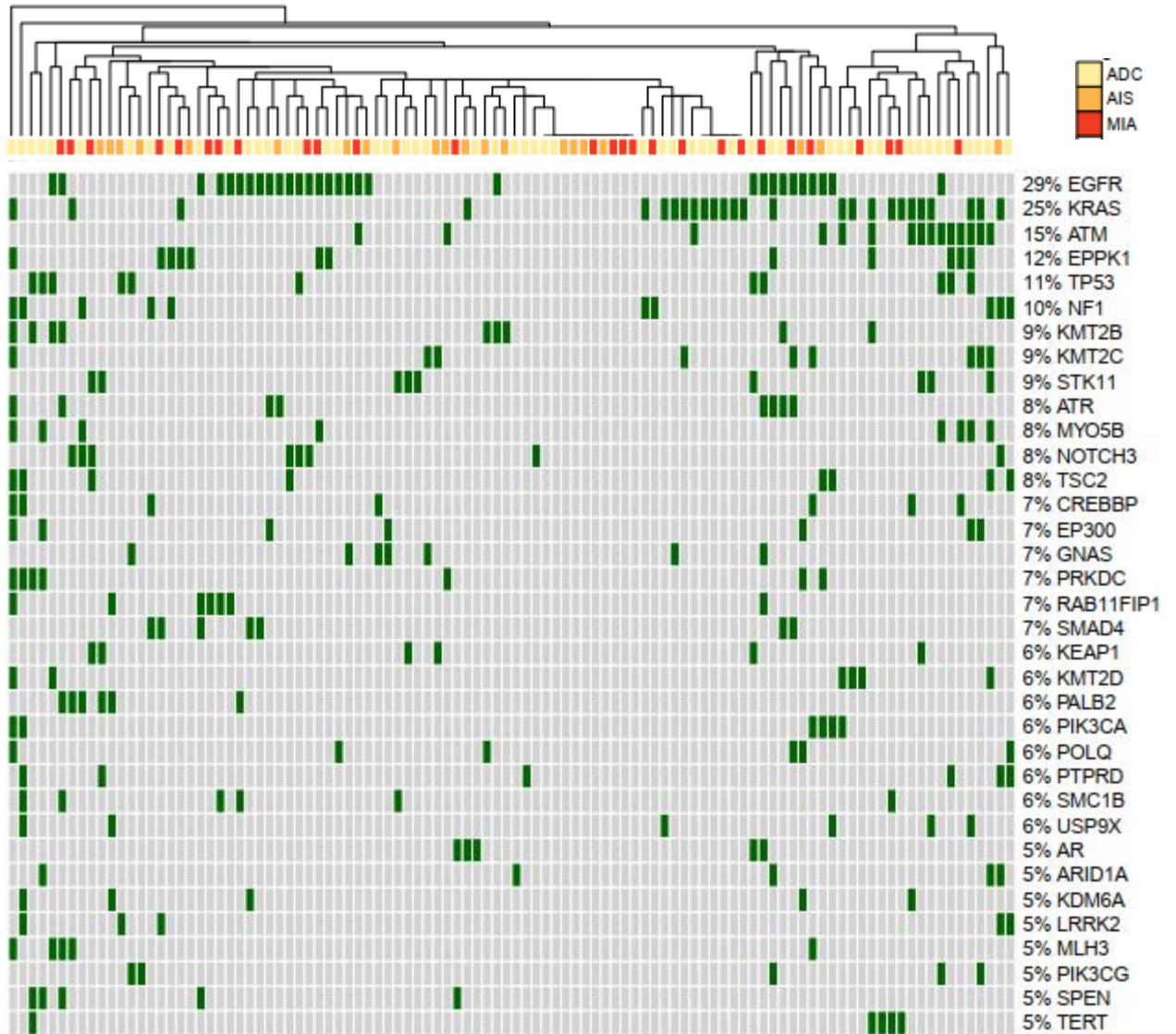


Figure E7

Figure E7. Clonality diversity is not associated with mutation signature

