

Supplementary

Table S1 Details of kits used in the present study

Name	Company	Country
QIAamp DNA FFPE Tissue Kit	Qiagen	USA
DNeasy Blood and tissue Kit	Qiagen	USA
dsDNA HS Assay Kit	ThermoFisher Scientific	USA
KAPA Hyper Prep Kit	KAPA Biosystems	USA
xGen Exome Research Panel and Hybridization and Wash Reagents Kit	Integrated DNA Technology	USA
RNeasy Plus Universal Kit	Qiagen	USA
Qubit™ RNA HS Assay Kit	ThermoFisher Scientific	USA
Take 3	BioTek	USA
RNA Cartridge kit of the Qseq100 Bio-Fragment Analyzer	Biopic	China
VAHTS mRNA-seq V3 Library Prep Kit	Vazyme	China

Table S2 Details of software used in the present study

Name	Version
SOAPnuke	1.5.6
Burrows-Wheeler Alignment tool	0.7.12
SAMtools	1.3
SAMBLASTER	0.1.22
VarScan	2.4.1
Snpeff	4.3
CNVkit	0.8.1
ascaNgs	3.1.0
POLYSOLVER	1.0
Bwakit	0.7.11
trim galore	0.6.7
Kallisto	0.46.2
Gencode	38.0
MiXCR	2.1.10
VDJtools	1.2.1
PyClone	0.13.0

Table S3 Detailed information on hyperparameter combinations per model

Model	Hyperparameter combinations
Support vector machine	kernel: ['linear'], C: [0.1, 1, 10]
Random forest	n_estimators: randint(100, 1000), max_depth: randint(5, 20), min_samples_split: randint(2, 10), min_samples_leaf: randint(1, 10), max_features: ['auto', 'sqrt'], bootstrap: [True, False]
Gradient boosting classifier	n_estimators: [50, 100, 200], learning_rate: [0.1, 0.01, 0.001], max_depth: [4]
Decision tree classifier	max_depth: [None, 5, 10, 15], min_samples_split: [2, 5, 10], min_samples_leaf: [1, 2, 3]
Extra tree classifier	n_estimators: [100, 200, 300], max_depth: [None, 5, 10], min_samples_split: [2, 5, 10]
Gaussian process classifier	kernel = 1.0 *RBF(1.0, length_scale_bounds=(1e-3, 1e3)) n_restarts_optimizer=10, max_iter_predict=100
K-nearest neighbors	n_neighbors=9, weights='uniform'

For unspecified model parameters, default values are used.

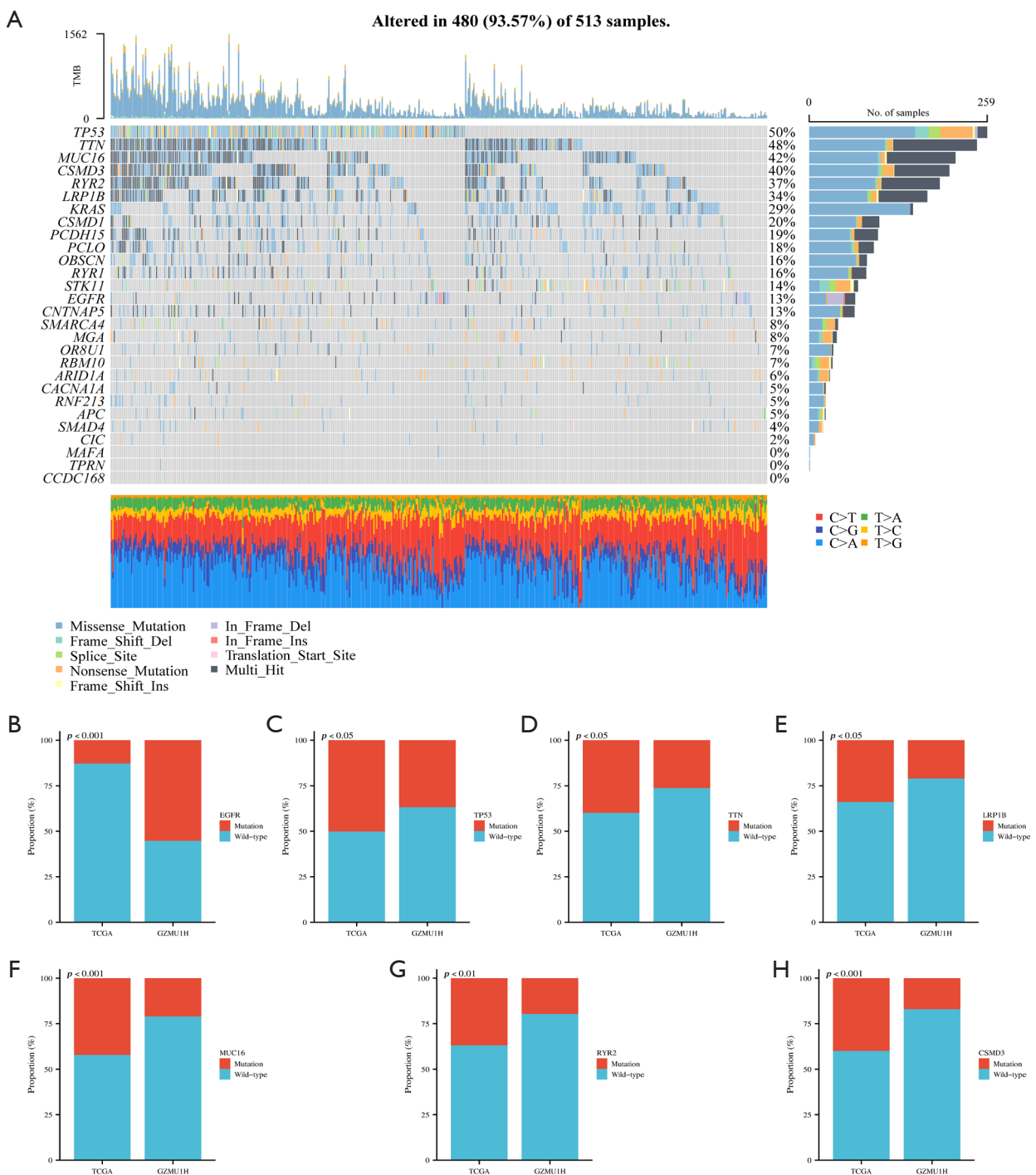


Figure S1 Mutation landscape of lung adenocarcinoma in the TCGA-LUAD cohort. Oncoplot demonstrating the highly mutant genes in the TCGA-LUAD dataset (A). Among the top ten mutated genes in the GZMU1H cohort, seven genes with significantly different mutant frequencies than the TCGA-LUAD cohort, including *EGFR* (B), *TP53* (C), *TTN* (D), *LRP1B* (E), *MUC16* (F), *RYR2* (G), and *CSMD3* (H), were found.

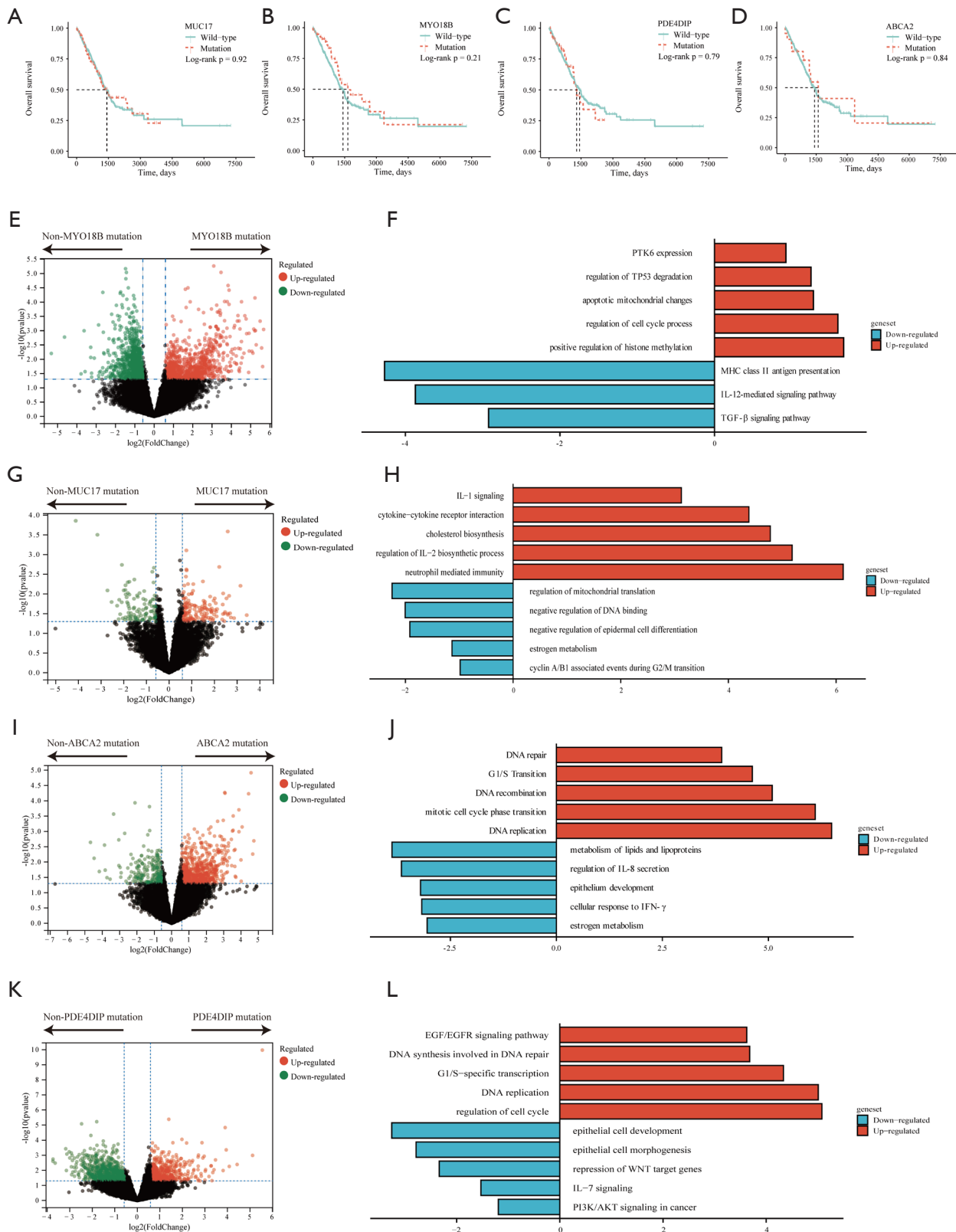


Figure S2 Transcriptomic spectrums and prognostic effects of specific gene mutations of early-stage non-squamous non-small cell lung cancer. Prognostic effects of monogenic mutation, including *MUC17* (A), *MYO18B* (B), *PDE4DIP* (C), and *ABCA2* (D), in the TCGA-LUAD cohort. Differentially expressed genes and corresponding enriched pathways of mutant versus wild-type *MYO18B* (E,F), *MUC17* (G,H), *ABCA2* (I,J), *PDE4DIP* (K,L) in the *GZM1H* cohort. Red and green dots refer to significantly up-regulated and down-regulated genes, respectively. Black dots represent genes with insignificant changes in expression levels.

Table S5 Gene expression levels with prognostic significance in the univariate Cox regression analysis

Gene	P value	HR	95% CI_L	95% CI_U
<i>CCR9</i>	0.017	3.28E-05	6.65E-09	0.162
<i>KLRC4</i>	0.023	1.66E-04	8.95E-08	0.309
<i>KIR2DL3</i>	0.050	0.004	1.69E-05	0.996
<i>KLRC1</i>	0.050	0.007	5.36E-05	0.997
<i>KLRD1</i>	0.004	0.009	3.46E-04	0.218
<i>NCR1</i>	0.016	0.009	2.06E-04	0.413
<i>IL27</i>	0.033	0.017	4.26E-04	0.718
<i>TXK</i>	0.023	0.130	0.022	0.758
<i>PTGDR</i>	0.040	0.137	0.021	0.915
<i>CD244</i>	0.025	0.160	0.032	0.793
<i>IL18RAP</i>	0.044	0.196	0.040	0.960
<i>DUOX2</i>	0.013	0.241	0.078	0.745
<i>IL18R1</i>	0.029	0.251	0.073	0.866
<i>LCN10</i>	0.049	0.266	0.071	0.995
<i>CXCR6</i>	0.007	0.266	0.102	0.698
<i>KLRK1</i>	0.018	0.290	0.105	0.805
<i>GPR17</i>	0.044	0.310	0.099	0.970
<i>PPP3CC</i>	0.045	0.318	0.104	0.974
<i>HDGFL3</i>	0.047	0.347	0.123	0.984
<i>PLCG2</i>	0.009	0.366	0.172	0.778
<i>TRAV3</i>	0.045	0.379	0.147	0.979
<i>IL34</i>	0.049	0.396	0.158	0.995
<i>SEMA6A</i>	0.042	0.415	0.178	0.969
<i>GIPR</i>	0.032	0.421	0.191	0.928
<i>ICAM2</i>	0.024	0.438	0.214	0.896
<i>PTK2B</i>	0.003	0.442	0.257	0.758
<i>JAK2</i>	0.035	0.447	0.211	0.944
<i>BMP6</i>	0.012	0.461	0.253	0.841
<i>TRIM22</i>	0.007	0.475	0.276	0.819
<i>S1PR1</i>	0.029	0.492	0.260	0.930
<i>EPOR</i>	0.047	0.494	0.246	0.991
<i>VIPR1</i>	0.042	0.526	0.283	0.977
<i>ADRB2</i>	0.030	0.624	0.408	0.955
<i>PI3</i>	0.030	0.637	0.424	0.957
<i>INHBB</i>	0.026	0.679	0.482	0.956
<i>DUOX1</i>	0.027	0.708	0.522	0.961
<i>CD79A</i>	0.028	0.724	0.542	0.966
<i>IGHG3</i>	0.041	0.805	0.655	0.991
<i>PLAU</i>	0.031	1.310	1.020	1.670
<i>SPP1</i>	0.005	1.347	1.090	1.660
<i>TUBB3</i>	0.050	1.381	1.000	1.910
<i>HTR3A</i>	0.046	1.453	1.010	2.100
<i>CRABP2</i>	0.001	1.532	1.190	1.970
<i>PLXNB3</i>	0.040	1.542	1.020	2.330
<i>F2RL1</i>	0.027	1.585	1.050	2.380
<i>PPIA</i>	0.035	1.609	1.030	2.500
<i>IL13RA2</i>	0.015	1.621	1.100	2.390
<i>EGFR</i>	0.028	1.643	1.060	2.560
<i>PLAUR</i>	0.017	1.663	1.090	2.530
<i>PGLYRP4</i>	0.043	1.715	1.020	2.890
<i>ULBP2</i>	0.008	1.747	1.160	2.640
<i>LGR4</i>	0.006	1.880	1.190	2.960
<i>IL31RA</i>	0.003	2.062	1.290	3.310
<i>IL11</i>	0.002	2.163	1.330	3.510
<i>PSMD4</i>	0.041	2.207	1.030	4.720
<i>GDF5</i>	0.016	2.486	1.180	5.230
<i>NGF</i>	0.038	3.060	1.060	8.800
<i>CCR8</i>	0.034	3.826	1.100	13.300
<i>HFE</i>	0.044	3.925	1.040	14.800
<i>EPGN</i>	0.038	4.016	1.080	14.900
<i>IL1A</i>	0.008	4.366	1.470	13.000
<i>CRHR1</i>	0.022	100.150	1.940	5.17E+03
<i>HTR3B</i>	0.044	1.03E+03	1.190	8.97E+05
<i>PPBPP2</i>	0.002	2.32E+04	39.500	1.36E+07
<i>MBL2</i>	0.010	1.42E+05	17.500	1.16E+09

HR, hazard ratio; CI_L, lower bounds of the 95% confidence interval; CI_U, upper bounds of the 95% confidence interval.

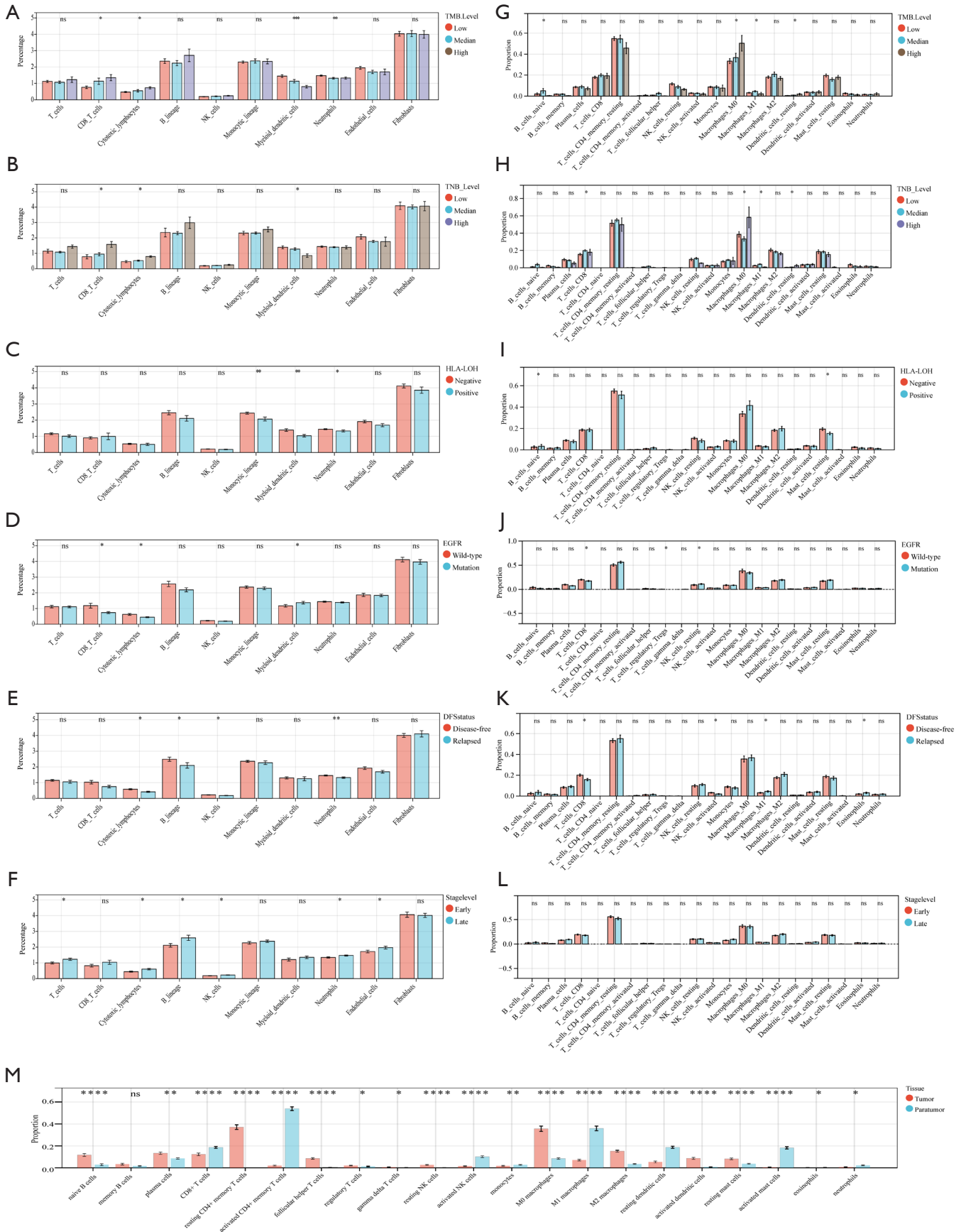


Figure S3 Immune infiltration differences between tumor nest and adjacent tissues. Intratumoral immune infiltration differences between different tumor mutational burden (TMB) levels (A), tumor neoantigen burden (TNB) levels (B), human leukocyte antigen loss of heterozygosity (HLA-LOH) status (C), EGFR mutation status (D), disease-free survival (DFS) status (E), and cTNM stage level (F), as evaluated by the MCP-counter algorithm. Intrastromal immune infiltration differences between different TMB levels (G), TNB levels (H), HLA-LOH status (I), EGFR mutation status (J), DFS status (K), and cTNM stage level (L), as evaluated by the CIBERSORT algorithm. Comparison of immune infiltration difference between tumor nest and paratumor tissue (M). Comparison of continuous data by Kruskal-Wallis test. *, $P < 0.05$; **, $P < 0.01$; ****, $P < 0.0001$; ns, non-significant.

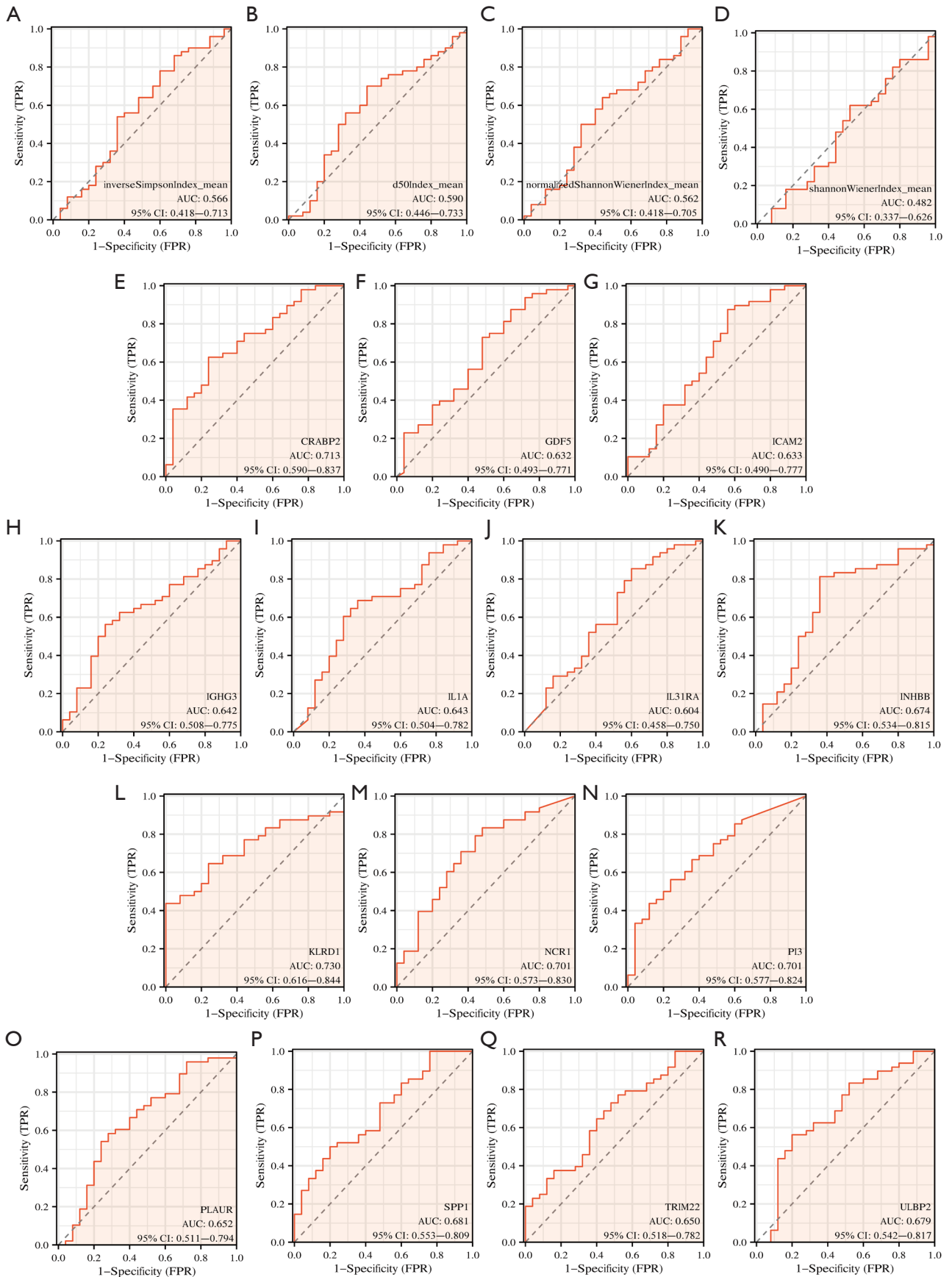


Figure S4 The predictive accuracy of single omics/biomarker in disease-free survival. T cell receptor repertoire diversity features (A-D), and transcriptomic characteristics (E-R) showed limited performance in predicting prognosis.

Table S6 Performance of the machine learning algorithms in the training cohort

	Accuracy	Precision	Recall	F1-score	AUC
Clinical + RNA					
SVM	0.833	0.813	0.650	0.722	0.818
DTC	0.933	0.900	0.900	0.900	0.978
ETC	1.000	1.000	1.000	1.000	1.000
GBC	0.950	1.000	0.850	0.919	0.998
GPC	0.975	0.975	0.975	0.975	0.999
KNN	0.767	0.800	0.400	0.533	0.815
RF	0.983	1.000	0.950	0.974	0.975
Clinical +DNA+RNA+TCR					
SVM	0.733	0.625	0.500	0.556	0.735
DTC	0.917	0.857	0.900	0.878	0.973
ETC	1.000	1.000	1.000	1.000	1.000
GBC	0.667	0	0	0	0.986
GPC	0.813	0.931	0.675	0.783	0.904
KNN	0.800	0.833	0.500	0.625	0.869
RF	0.933	1.000	0.800	0.889	1.000
Clinical + DNA					
SVM	0.767	0.750	0.450	0.563	0.829
DTC	0.983	1.000	0.950	0.974	0.999
ETC	0.900	1.000	0.700	0.824	0.988
GBC	1.000	1.000	1.000	1.000	1.000
GPC	0.975	0.975	0.975	0.975	0.997
KNN	0.783	0.733	0.550	0.629	0.803
RF	0.833	0.917	0.550	0.687	0.936
Clinical + TCR					
SVM	0.717	0.600	0.450	0.514	0.685
DTC	0.917	0.941	0.800	0.865	0.968
ETC	1.000	1.000	1.000	1.000	1.000
GBC	1.000	1.000	1.000	1.000	1.000
GPC	0.875	0.917	0.825	0.868	0.890
KNN	0.767	0.800	0.400	0.533	0.804
RF	0.817	0.846	0.550	0.667	0.926
DNA+RNA					
SVM	0.783	0.733	0.550	0.629	0.860
DTC	0.933	0.833	1.000	0.909	0.972
ETC	0.917	1.000	0.750	0.857	0.992
GBC	0.667	0	0	0	0.973
GPC	1.000	1.000	1.000	1.000	1.000
KNN	0.733	0.667	0.400	0.500	0.819
RF	0.983	1.000	0.950	0.974	1.000
DNA+RNA+TCR					
SVM	0.650	0.455	0.250	0.323	0.695
DTC	0.967	0.950	0.950	0.950	0.997
ETC	0.983	1.000	0.950	0.974	1.000
GBC	1.000	1.000	1.000	1.000	1.000
GPC	0.788	0.829	0.725	0.773	0.869
KNN	0.800	0.900	0.450	0.600	0.850
RF	0.967	1.000	0.900	0.947	1.000

SVM, support vector machine; DTC, decision tree classifier; ETC, extra tree classifier; GBC, gradient boosting classifier; GPC, Gaussian process classifier; KNN, K-nearest neighbors; RF, random forest.

Table S7 Performance of the machine learning algorithms in the testing cohort

	Accuracy	Precision	Recall	F1-score	AUC
Clinical + RNA					
DTC	0.688	0.600	0.500	0.545	0.692
ETC	0.875	1.000	0.667	0.800	0.933
GBC	0.688	1.000	0.167	0.286	0.800
GPC	0.688	1.000	0.167	0.286	0.733
KNN	0.750	1.000	0.333	0.500	0.867
RF	0.813	1.000	0.500	0.667	0.833
SVM	0.813	1.000	0.500	0.667	0.900
Clinical +DNA+RNA+TCR					
DTC	0.625	0.500	0.333	0.400	0.667
ETC	0.688	0.667	0.333	0.444	0.833
GBC	0.625	0	0	0	0.717
GPC	0.750	0.625	0.833	0.714	0.817
KNN	0.625	0	0	0	0.933
RF	0.625	0	0	0	0.817
SVM	0.563	0.400	0.333	0.364	0.617
Clinical + DNA					
DTC	0.625	0.500	0.333	0.400	0.567
ETC	0.625	0.500	0.167	0.250	0.717
GBC	0.625	0.500	0.333	0.400	0.433
GPC	0.563	0	0	0	0.517
KNN	0.875	0.833	0.833	0.833	0.767
RF	0.625	0	0	0	0.733
SVM	0.688	1.000	0.167	0.286	0.750
Clinical + TCR					
DTC	0.563	0.400	0.333	0.364	0.508
ETC	0.688	1.000	0.167	0.286	0.650
GBC	0.625	0	0	0	0.650
GPC	0.813	0.714	0.833	0.769	0.850
KNN	0.750	1.000	0.333	0.500	0.833
RF	0.688	1.000	0.167	0.286	0.600
SVM	0.625	0.500	0.167	0.250	0.450
DNA+RNA					
DTC	0.750	0.750	0.500	0.600	0.683
ETC	0.625	0.500	0.167	0.250	0.700
GBC	0.625	0	0	0	0.833
GPC	0.625	0.500	0.333	0.400	0.817
KNN	0.688	0.667	0.333	0.444	0.817
RF	0.688	1	0.167	0.286	0.767
SVM	0.625	0.500	0.167	0.250	0.667
DNA+RNA+TCR					
DTC	0.625	0.500	0.167	0.250	0.608
ETC	0.625	0.500	0.167	0.250	0.683
GBC	0.688	1.000	0.167	0.286	0.633
GPC	0.875	0.833	0.833	0.833	0.900
KNN	0.625	0	0	0	0.758
RF	0.625	0	0	0	0.683
SVM	0.563	0	0	0	0.483

DTC, decision tree classifier; ETC, extra tree classifier; GBC, gradient boosting classifier; GPC, Gaussian process classifier; KNN, K-nearest neighbors; RF, random forest; SVM, support vector machine.