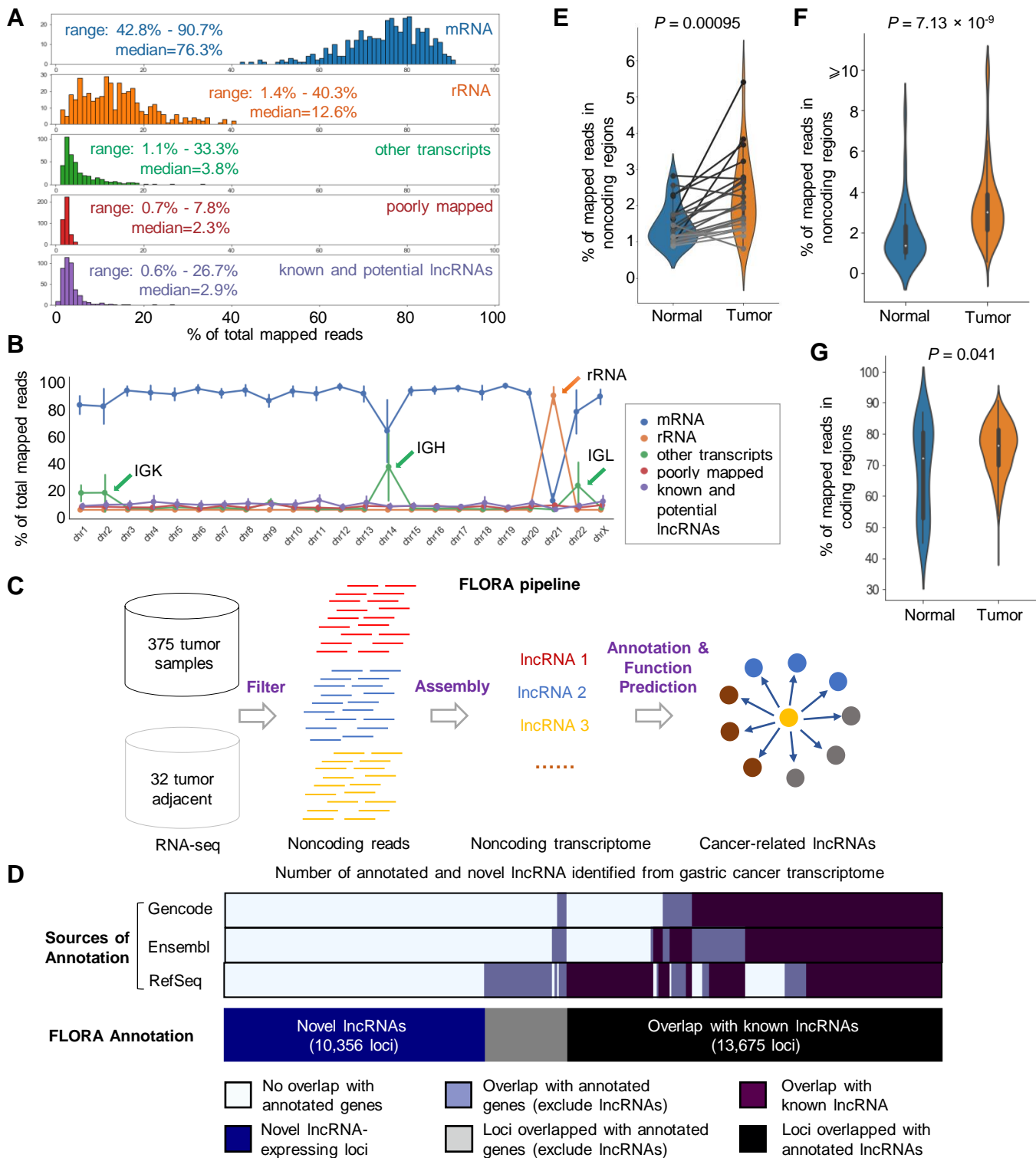


Fig. S1



Supplementary Figure S1. Identification of lncRNA by FLORA pipeline.

A The range and median fraction of reads from different types of transcripts in the whole transcriptome sequencing data of normal gastric samples and gastric cancer in the TCGA cohort.

B Fraction of reads mapped to different RNA species and reads with low mapping quality on each chromosome across all normal and tumor tissues. Genes that constitute large fractions of total mapped reads on several chromosomes are marked. The mean values are marked as dot, and standard derivations are marked as ticks.

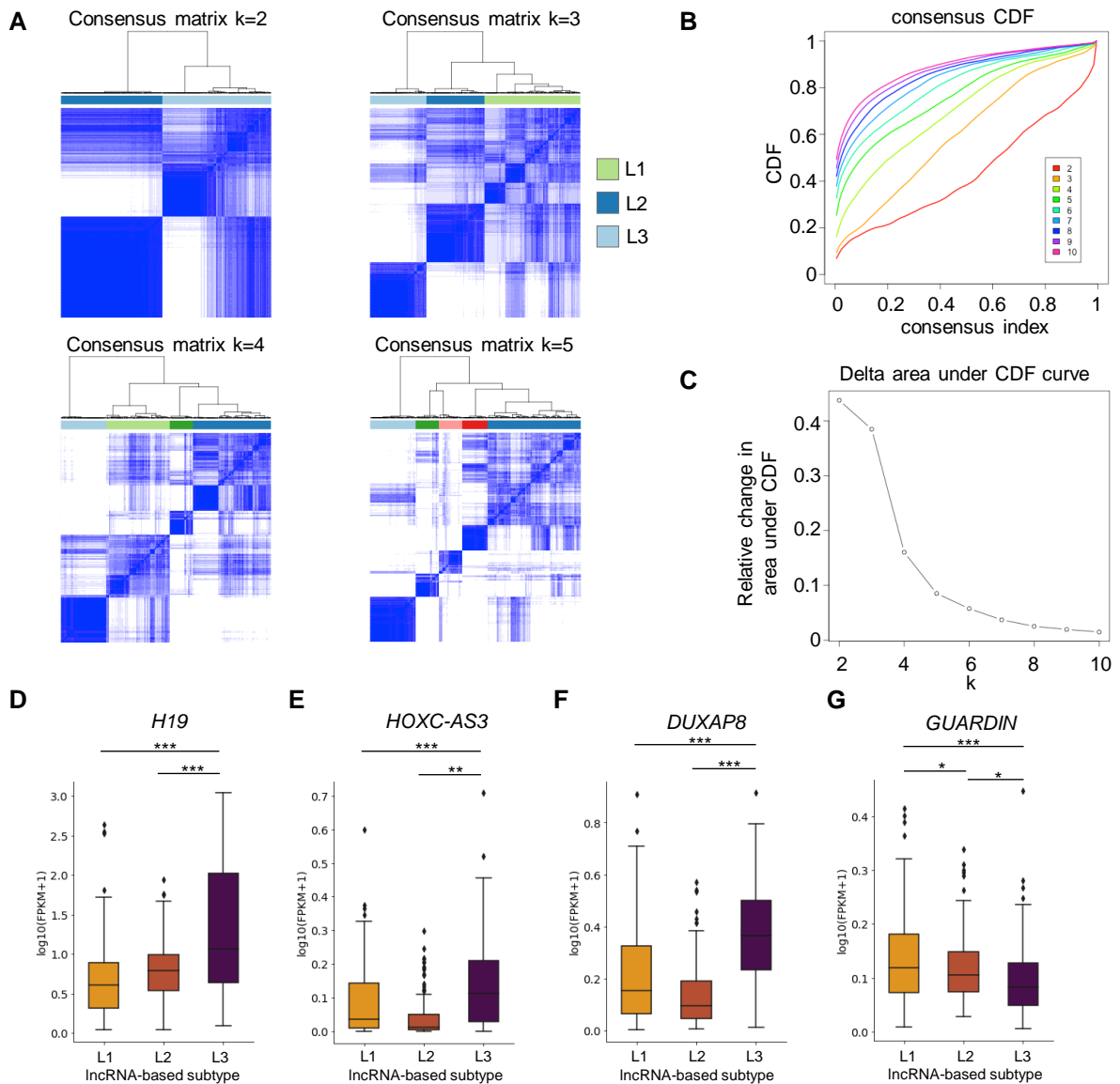
C Overview of the workflow of the FLORA pipeline.

D Number of known and novel lncRNAs in TCGA GC by FLORA pipeline. The 28,507 loci encoding potential lncRNAs are represented as columns. Loci that overlap with known lncRNAs or other RNA species in Gencode, Ensembl or RefSeq data are marked by purple and light blue, respectively. 13,675 loci overlap with known lncRNAs, and 10,356 loci that do not overlap with any annotated genes are defined as novel lncRNAs.

E Fraction of reads mapped to noncoding regions of the genome in the TCGA GC patients with paired normal and tumor samples (N=27). The normal and tumor samples from the same patient are connected with lines. The p-value is calculated by two-tailed paired t-test.

F-G Fraction of reads mapped to noncoding (F) and coding regions (G) of the genome in all tumor (N=375) and normal (N=32) TCGA samples. P-values are calculated by the Wilcoxon's rank-sum test.

Fig. S2



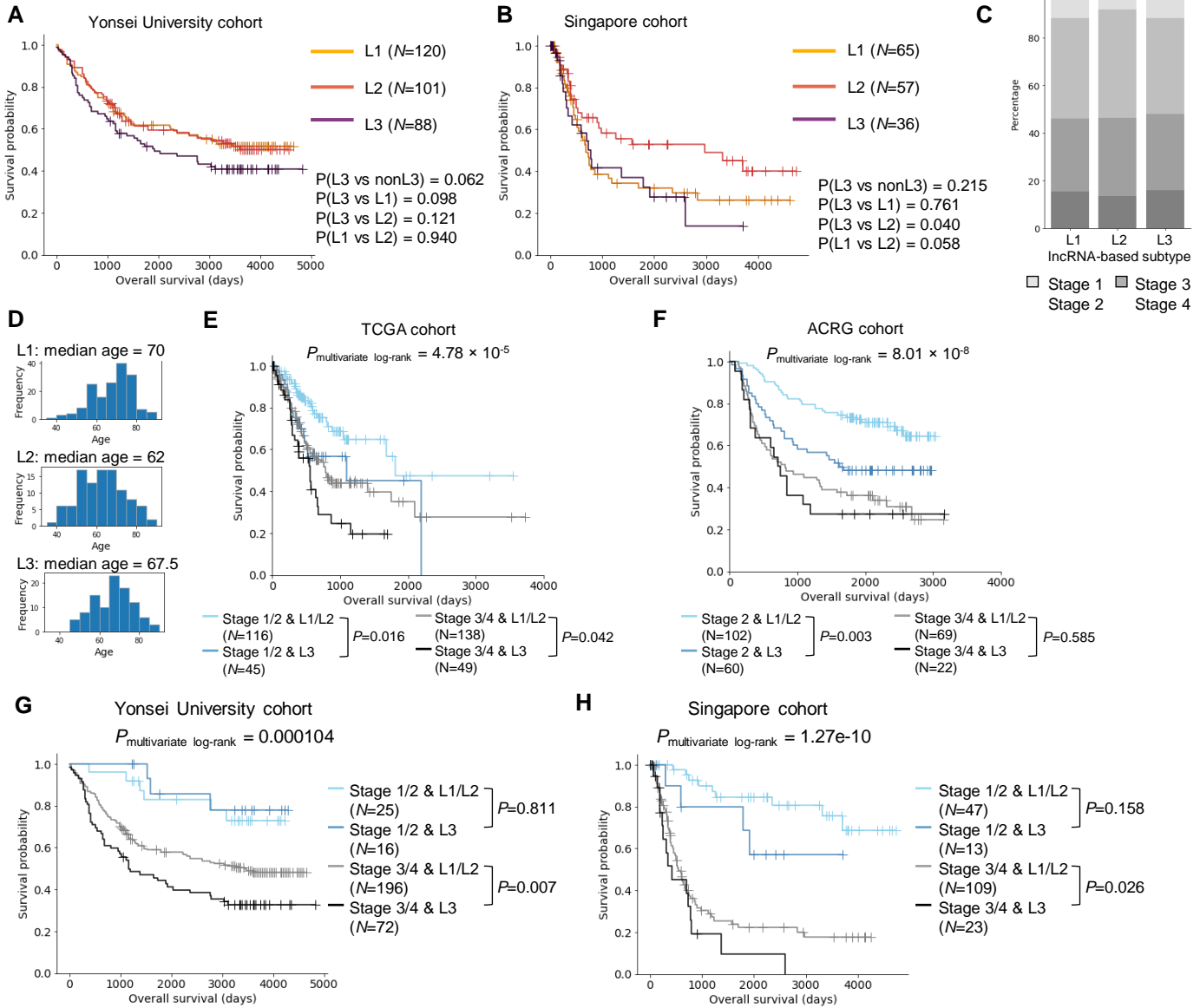
Supplementary Figure S2. Identification of lncRNA expressed-based GC subtypes.

A Consensus matrices of lncRNA expression-based clustering with the number of classes k ranging from 2 to 5. **B** Cumulative distribution function (CDF) curve of consensus index for different clustering results with the number of classes k ranging from 2 to 10.

C Relative change in area under CDF curve shown in panel B.

D-G Examples of differentially expressed lncRNAs in L3 subtype compared with L1 and L2 subtypes, in GC samples from the TCGA cohort. The midline, boxes, whiskers and dots show the median, quartiles, ranges and outliers of the distribution, respectively. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$ by unpaired two-tailed t-test.

Fig. S3



Supplementary Figure S3. Comparisons of the survival differences in different GC subtypes.

A-B Kaplan-Meier curve illustrating the overall survival for GC patients of the three subtypes from the Yonsei University (A) and Singapore (B) cohort. P-values are calculated by log-rank test.

C Distribution of tumor stages in the three lncRNA-based subtypes.

D Distribution of age at diagnosis in the three lncRNA-based subtypes.

E-H Survival probability of patients segregated by the combination of tumor staging and lncRNA-based subtypes. The P-value across all subgroups is calculated by multivariate log-rank test, and the P-value of pairwise comparison by log-rank test.

Fig. S4

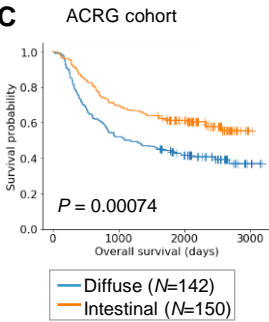
A

ACRG cohort			Univariate analysis		Multivariate analysis	
Factor	Variable	Comparator	HR (95% CI)	P-value	HR (95% CI)	P-value
AJCC stage	Stage 3/4	Stage 1/2	2.37 (1.72 – 3.25)	0.0000001	2.65 (1.85 – 3.80)	0.0000001
lncRNA-based subtype	L3	L1/L2	1.44 (1.26 – 3.18)	0.043204	1.72 (1.18 – 2.51)	0.004551
Age	Age ≥ 65	Age < 65	1.94 (1.22 – 3.09)	0.007213	1.69 (1.18 – 2.42)	0.003858
Histology	Diffuse	Intestinal	1.54 (0.92 – 2.57)	0.000865	1.47 (1.14 – 1.90)	0.002800
EBV	EBV	No EBV	1.06 (0.56 – 2.01)	0.867122		

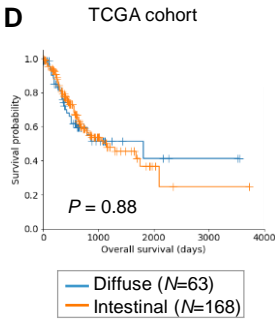
B

Yonsei University cohort			Univariate analysis		Multivariate analysis	
Factor	Variable	Comparator	HR (95% CI)	P-value	HR (95% CI)	P-value
AJCC stage	Stage 3/4	Stage 1/2	3.25 (1.65 – 6.37)	0.000615	3.41 (1.74 – 6.71)	0.000375
lncRNA-based subtype	L3	L1/L2	1.38 (0.98 – 1.93)	0.062949	1.48 (1.05 – 2.07)	0.023779
Age	Age ≥ 65	Age < 65	1.48 (1.07 – 2.04)	0.01782	1.42 (1.03 – 1.97)	0.032415

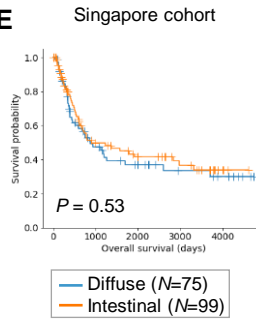
C



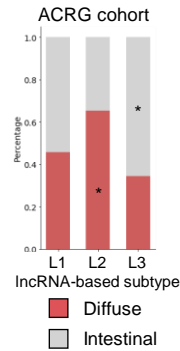
D



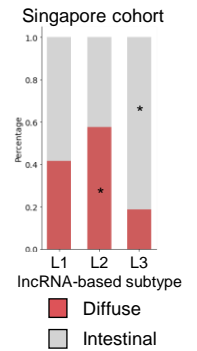
E



F



G



Supplementary Figure S4. The lncRNA-based subtype is an independent prognostic factor in GC.

A-B Value and 95% confidence intervals (95% CI) of the hazard ratios (HR) of different factors considered in the univariate and multivariate cox regression analysis of ACRG (A) and Yonsei University cohort (B).

C-E Survival probability of ACRG (C), TCGA (D) and Singapore cohort (E) patients with diffuse and intestinal subtype. *P*-value between subtypes are calculated by log-rank test.

F-G Distribution of diffuse and intestinal histology in the three lncRNA-based subtypes from the ACRG (F) and Singapore cohort (G). Asterisks indicate significantly enriched histological subtype in each of the three lncRNA-based subtypes by one-sided Fisher's exact test (*P* < 0.05).

Supplementary Figure S5. Genomic alterations, copy number variations, gene expression and survival in GC patients of lncRNA-based subtypes.

A Frequency of somatic mutations which are significantly enriched in lncRNA-based subtype ($P < 0.01$ by Fisher's exact test). Each dot represents one frequently mutated gene in GC, with size proportional to its mutation frequency.

B-C Ploidy level of GC samples of L1, L2 and L3 subtypes in the TCGA (B) and ACRG cohort (C). * $P < 0.05$, *** $P < 0.001$ by Wilcoxon's rank-sum test.

D Percentage of L1, L2 and L3 subtype samples from TCGA with whole-genome duplications (WGD1 and WGD2) and non-WGD, defined in Liu et al. * $P < 0.05$ by Fisher's exact test.

E The fraction of genome with copy number variations in different lncRNA-based subtypes in the ACRG cohort. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$ by Wilcoxon's rank-sum test.

F-G Recurrently amplified (F) and deleted (G) regions observed in lncRNA-based subtypes L1, L2 and L3 (in orange, red and purple, respectively) from the TCGA cohort. The x-axis of circles represents the chromosomal location and the y-axis represents the G-score of each location (ranging from 0-1.5 for amplifications and 0-1 for deletions). Genes located at significantly amplified/deleted peaks were marked, with genes that are more frequently amplified or deleted in L3 highlighted in red and blue, respectively.

H-I Recurrently amplified (H) and deleted (I) regions observed in lncRNA-based subtypes L1, L2 and L3 (in orange, red and purple, respectively) from the ACRG cohort. The x-axis of circles represents the chromosomal location and the y-axis represents the G-score of each location (ranging from 0-2.0 for amplifications and 0-2.0 for deletions).

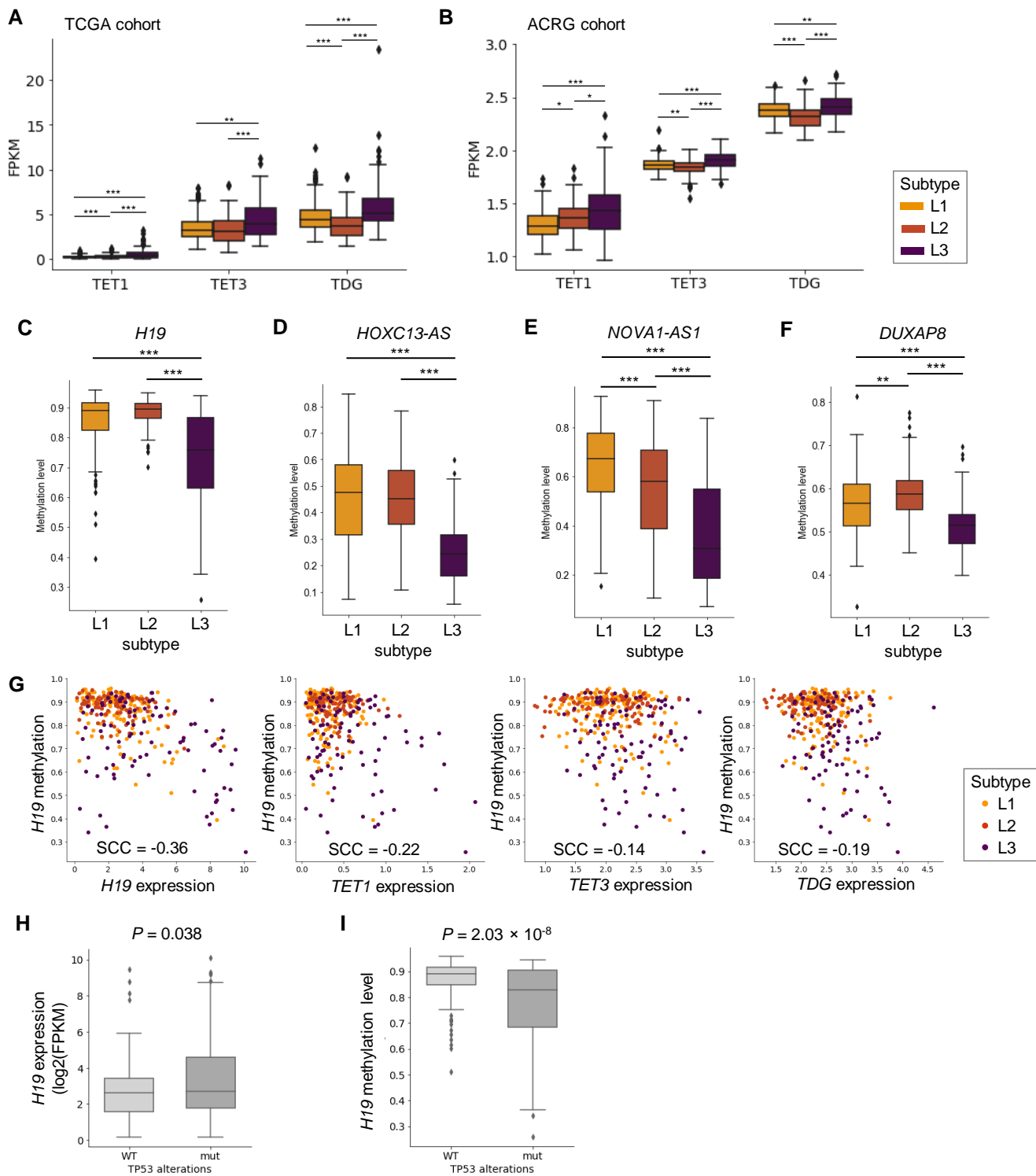
J Over-expression of *CCNE1* was observed in L3 subtype compared with L1 and L2. *** $P < 0.001$ by Wilcoxon's rank-sum test.

K-L Survival probability of TCGA GC patients in different molecular subtypes defined by Bass et al., (K) and Liu et al., (L). P -value is calculated by multivariate log-rank test.

M Distribution of TCGA molecular subtypes (defined by Bass et al., 2014) in the three lncRNA-based subtypes. * $P < 0.05$ by Fisher's exact test.

N Survival probability of L1-CIN and L3-CIN subgroups. P -value is calculated by log-rank test.

Fig. S6



Supplementary Figure S6. DNA hypomethylation and *TP53* mutations in L3 subtype are associated lncRNA overexpression.

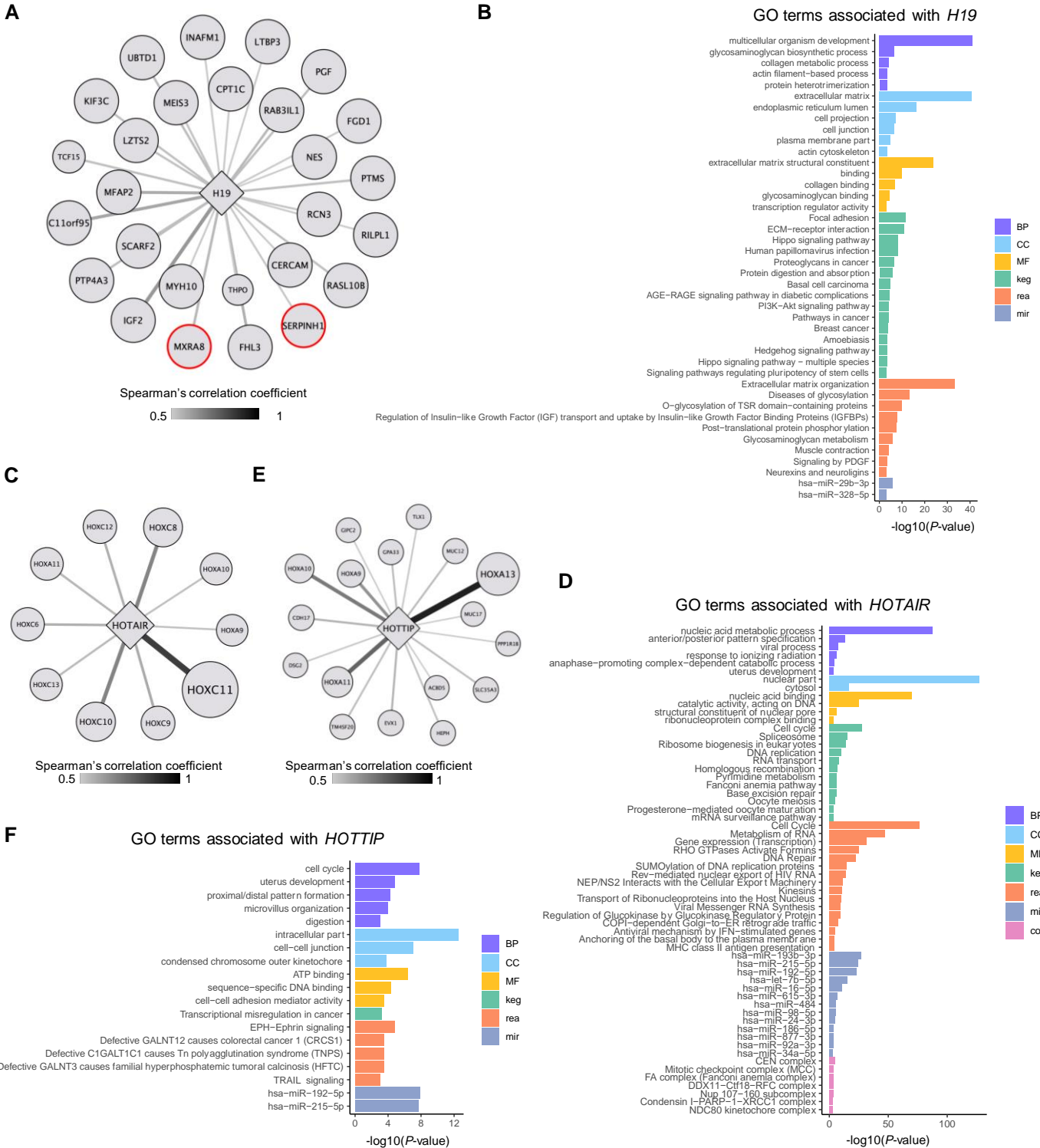
A-B Expression of *TET1*, *TET3* and *TDG* in lncRNA-based GC subtypes in the TCGA (A) and ACRG (B) cohorts.

C-F Hypomethylation of lncRNA species that are over-expressed in L3. The midline, boxes, whiskers and dots show the median, quartiles, ranges and outliers of the distribution, respectively.

G *H19* methylation level is anti-correlated with the expression level of *H19*, *TET1*, *TET3* and *TDG*. Probe cg11716026 was selected to represent *H19* gene methylation level. Each dot represents one GC sample in the TCGA dataset. GC samples of lncRNA expression-based subtype L1, L2 and L3 are shown in yellow, red and purple, respectively.

H-I Level of *H19* expression (H) and methylation (I) in TCGA GC samples with no *TP53* alterations and with *TP53* mutations and/or deletions. The midline, boxes, whiskers and dots show the median, quartiles, ranges and outliers of the distribution, respectively. * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$ by two-tailed unpaired t-test.

Fig. S7

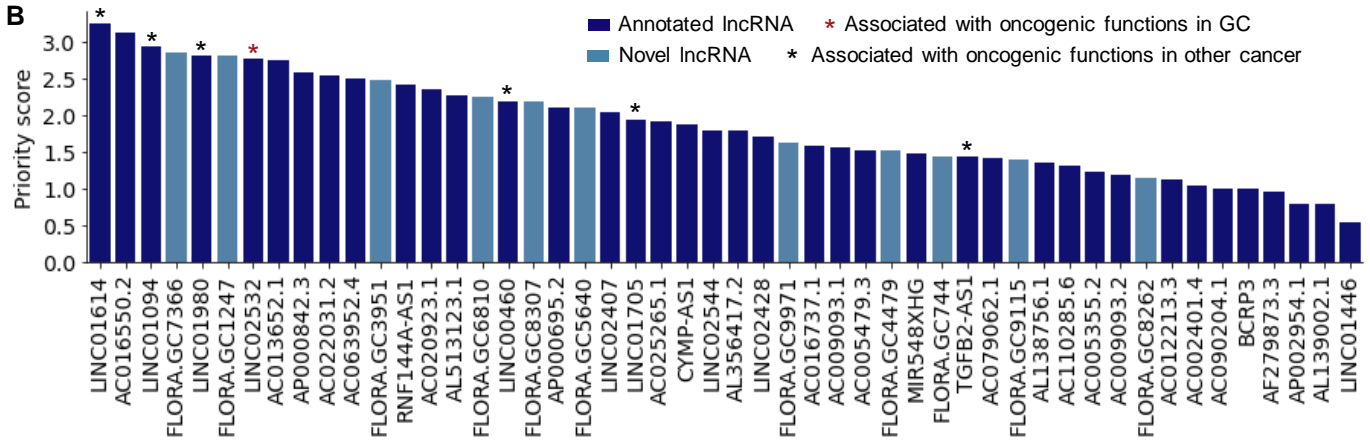
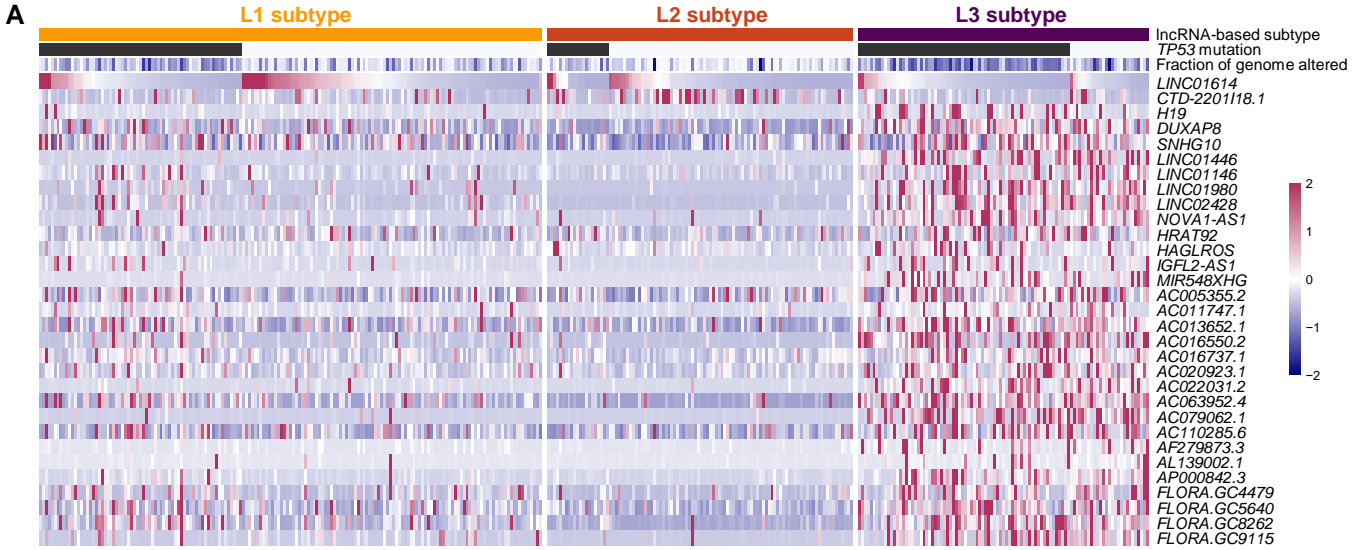


Supplementary Figure S7. Functional prediction and experimental validation of oncogenic lncRNAs.

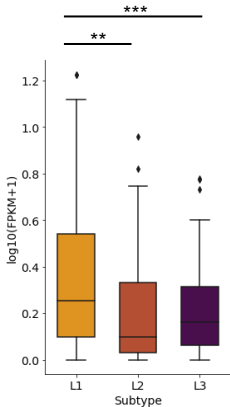
A,C,E Co-expression network of *H19* (A), *HOTAIR* (C) and *HOTTIP* (E) in GC. The genes with Spearman's correlation coefficient above 0.5 are showed in the network. Node size is proportional to $-\log_{10}(P\text{-value})$ of Spearman's correlation coefficient. The edge color represents the Spearman's correlation coefficient between the expression of lncRNA and its co-expressed genes.

B,D,F GO terms associated with *H19* (B), *HOTAIR* (D) and *HOTTIP* (F). The color represents three sub-ontologies of GO terms (BP: biological process; CC: cellular component; MF: molecular function; keg: KEGG pathways; rea: reactome; mir: miRNA targets).

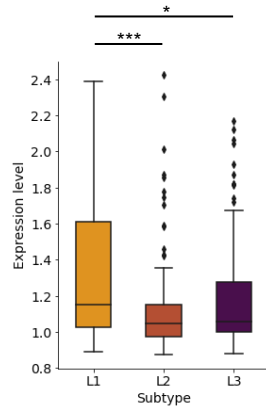
Fig. S8



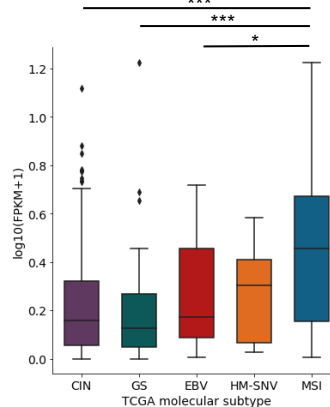
C *LINC01614* expression in TCGA cohort



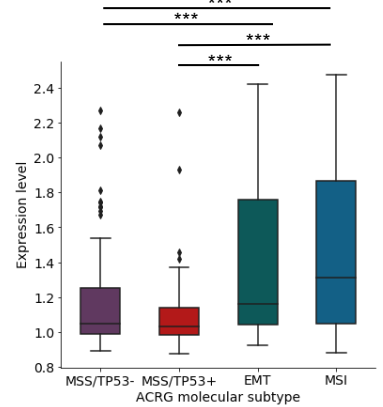
D *LINC01614* expression in ACRG cohort



E *LINC01614* expression in TCGA cohort



F *LINC01614* expression in ACRG cohort



Supplementary Figure S8. Expression of *LINC01614* in different subtypes of GC.

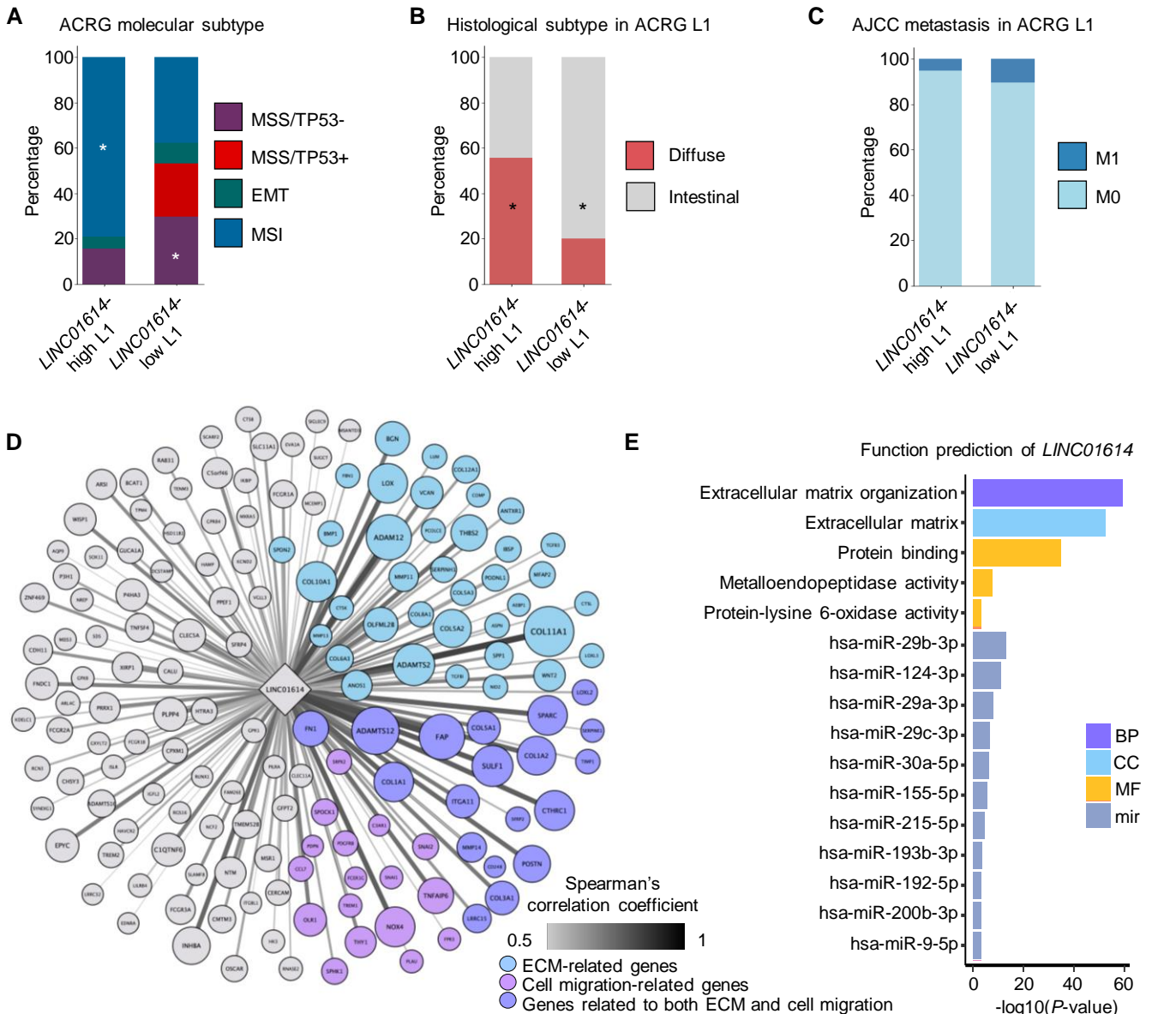
A Heatmap showing the molecular characteristics of the lncRNA expression-based subtypes, including prevalence of *TP53* mutations, fraction of genome with copy number alterations and expression of several subtype-specific lncRNAs.

B Prioritization of 50 GC-specific survival-related lncRNAs in GC. The annotated lncRNAs and novel lncRNAs are colored navy and blue, respectively. Asterisks indicate lncRNAs with experimental evidence of oncogenic functions.

C-D Expression of *LINC01614* in lncRNA expression-based subtypes in the TCGA (B) and ACRG cohort (C).

E-F Expression of *LINC01614* in molecular subtypes defined by the TCGA study (E) and by the ACRG cohort (F). * $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$ by two-tailed unpaired t-test.

Fig. S9

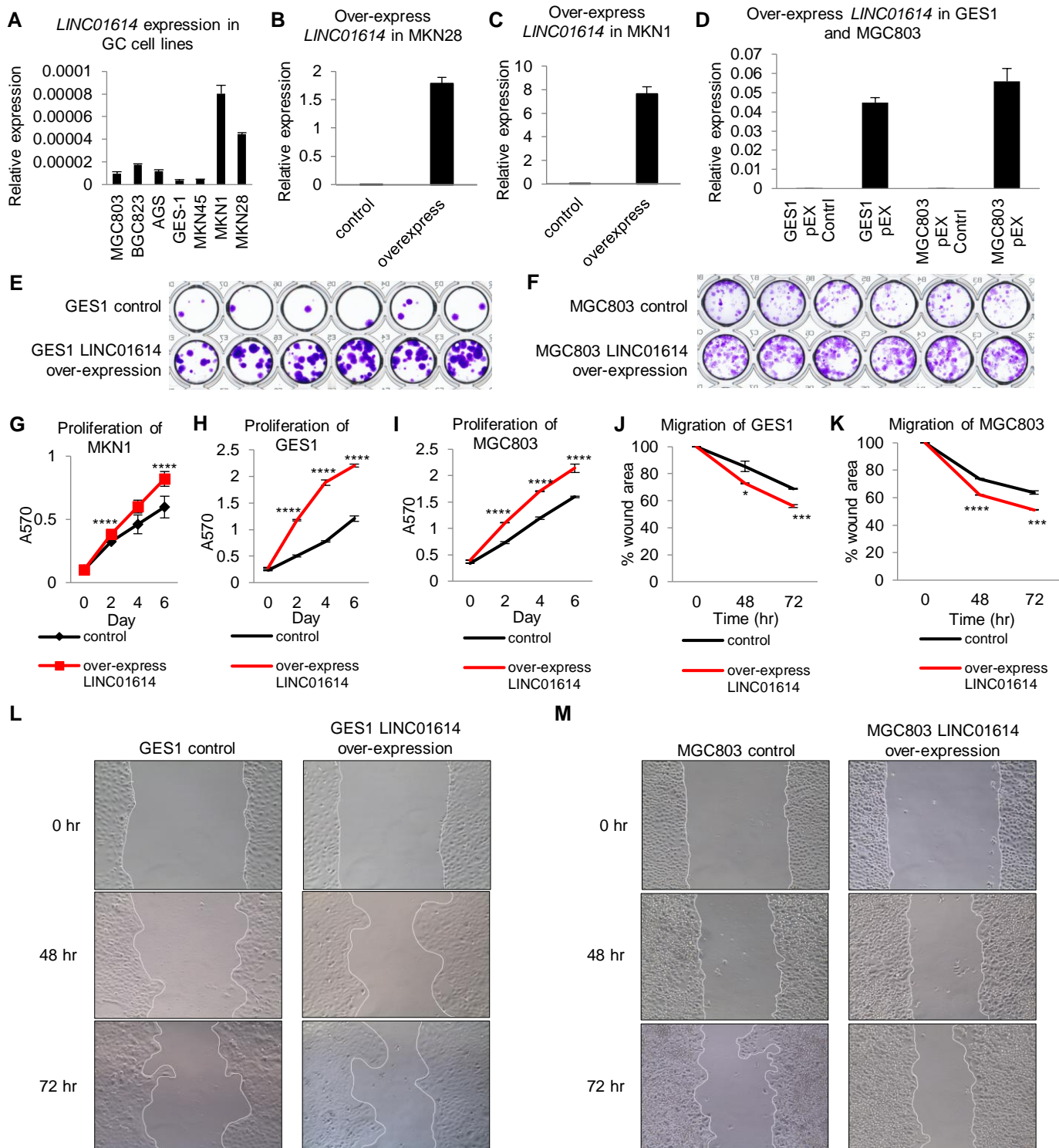


Supplementary Figure S9. Characteristic of *LINC01614*-high L1 tumor and functional prediction of *LINC01614* based on co-expression analysis in TCGA dataset.

A-C Distribution of ACRG *LINC01614*-high L1 and *LINC01614*-low L1 subgroups in molecular subtypes defined by the ACRG cohort (A), histological subtypes (B) and AJCC metastasis status (C). * $P < 0.05$ by Fisher's exact test.

D Co-expression network of *LINC01614*. Genes with Spearman's correlation coefficient > 0.5 were selected for the network visualization. Node size represents mutual information between the expression of *LINC01614* and the gene in the circle, while edge color represents the spearman correlation coefficient. Genes related to extracellular matrix (ECM), cell migration or both are colored blue, purple and orchid, respectively. **E** GO terms enriched in the co-expression network of *LINC01614* (BP: biological process; CC: cellular component; MF: molecular function; mir: miRNA targets)..

Fig. S10



Supplementary Figure S10. Over-expression of *LINC01614* in GC cell lines promotes proliferation and migration.

A Validation of *LINC01614* expression in normal (GES1) and gastric cancer cell lines by semi-quantitative PCR.

B-D Validation of *LINC01614* over-expression in MKN28 (B), MKN1 (C), GES1 and MGC803 (D) cell lines by qPCR.

E-F Colony formation in GES1 (E) and MGC803 (F) cell lines with and without *LINC01614* over-expression.

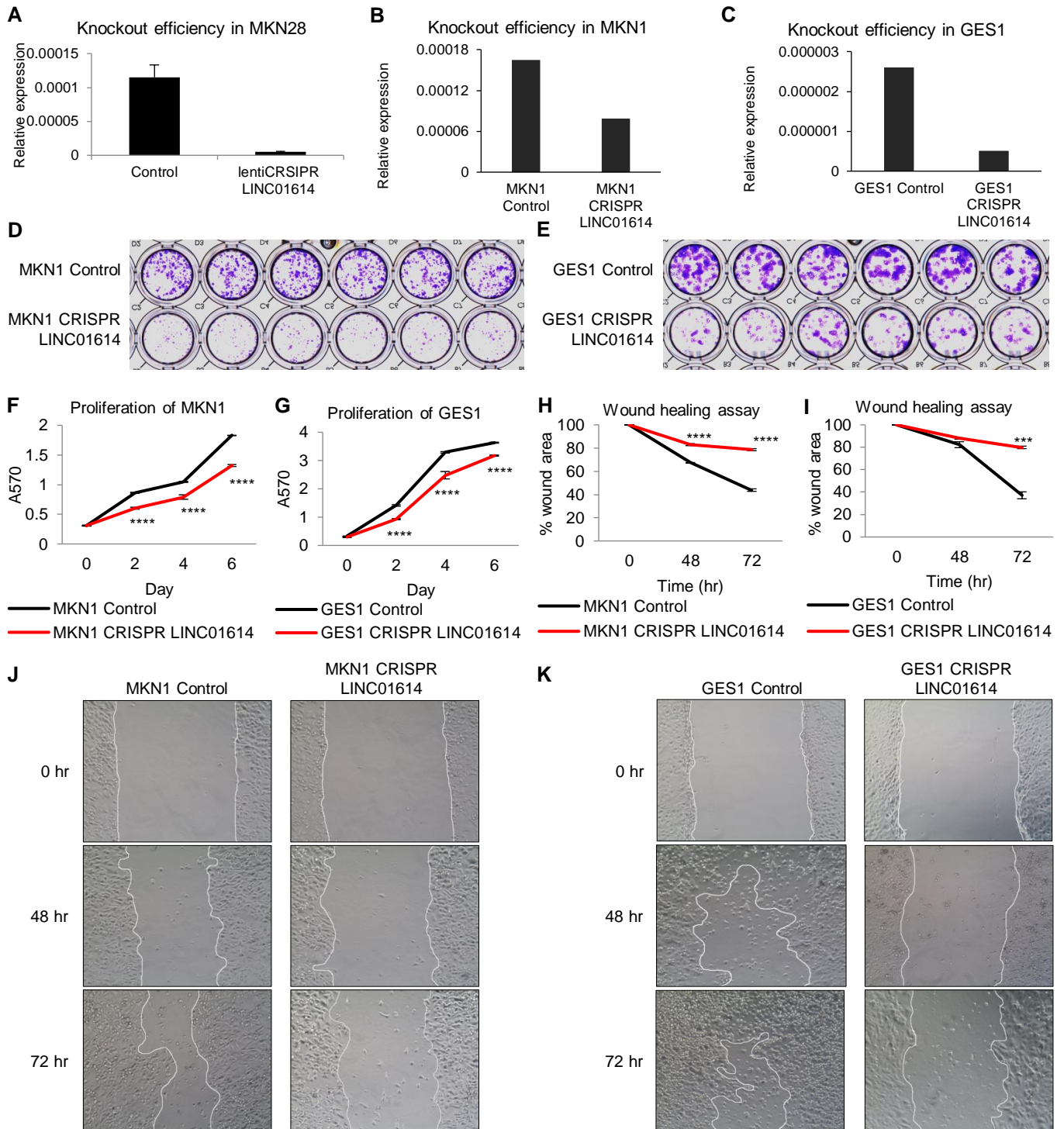
G-I Proliferation rate under *LINC01614* over-expression (red) and control (black) in MKN1 (G), GES1 (H) and MGC803 (I) cell lines.

J-K Wound healing rate under *LINC01614* over-expression (red) and control (black) in GES1 (J) and MGC803 (K) cell lines.

L-M Representative images of wound healing assay in GES1 (L) and MGC803 (M) cell lines with *LINC01614* over-expression and control.

In all experiments, three biological replicates were performed for each group. Data are represented as mean \pm SD. * $P < 0.05$; ** $P < 0.01$; *** $P < 0.001$; **** $P < 0.0001$ by two-tailed unpaired t-test.

Fig. S11



Supplementary Figure S11. Experimental validation of *LINC01614* functions in driving GC cell proliferation and migration.

A-C Validation of effective CRISPR-Cas9 knockout of *LINC01614* and reduced expression in MKN28 (A), MKN1 (B) and GES1 (C) cell lines by qPCR.

D-E Colony formation in MKN1 (E) and GES1 (F) cell lines with and without CRISPR-Cas9 knockout of *LINC01614*.

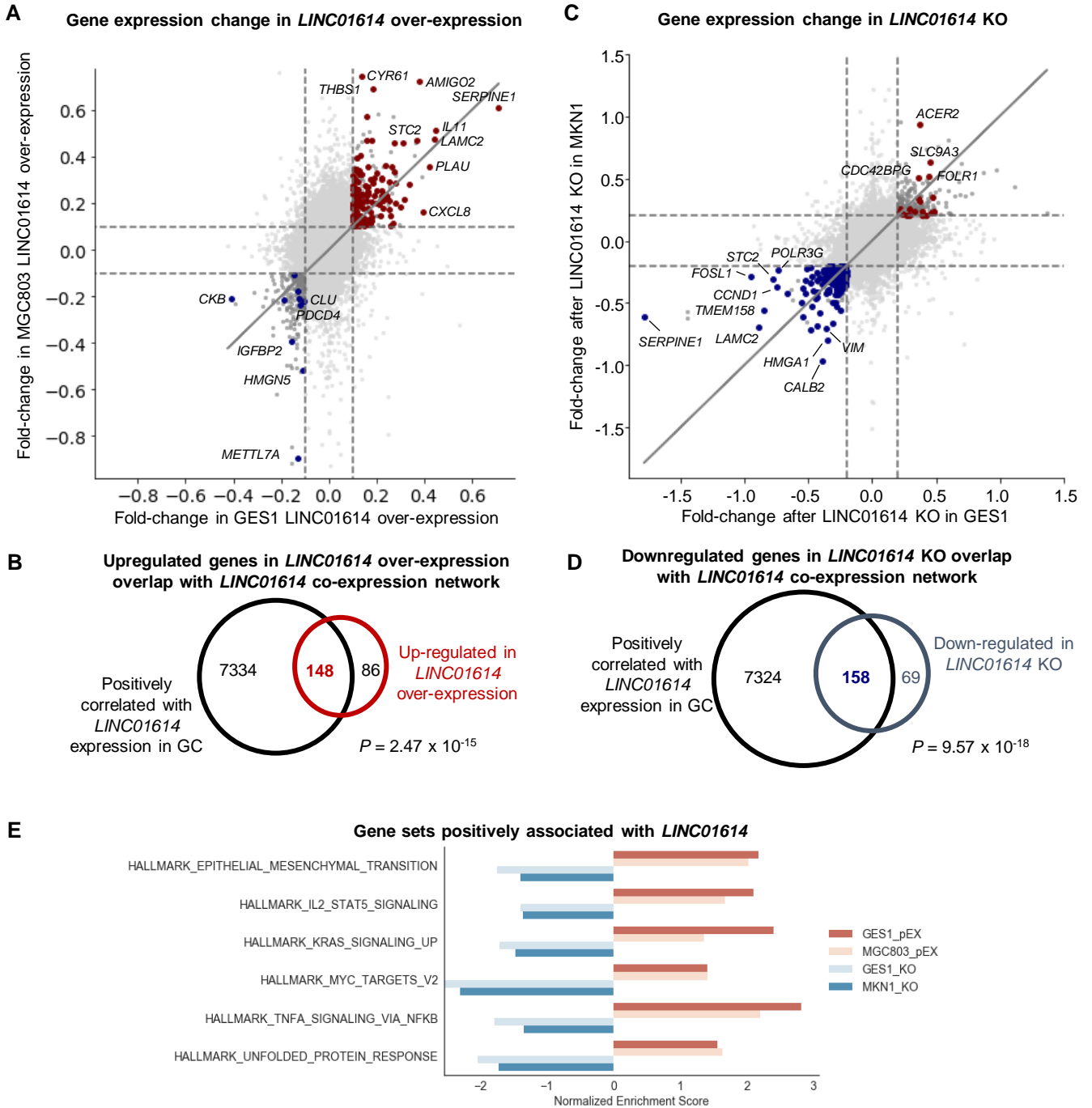
F-G Proliferation rate under CRISPR-Cas9 knockout of *LINC01614* (red) and control (black) in MKN1 (F) and GES1 (G) cell lines.

H-I Wound healing rate under CRISPR-Cas9 knockout of *LINC01614* (red) and control (black) in MKN1 (H) and GES1 (I) cell lines.

J-K Representative images of wound healing assay in MKN1 (L) and GES1 (M) cell lines with CRISPR-Cas9 knockout of *LINC01614* and control.

In all experiments, three biological replicates were performed for each group. Data are represented as mean \pm SD. * $P < 0.05$; ** $P < 0.01$; *** $P < 0.001$; **** $P < 0.0001$ by two-tailed unpaired t-test.

Fig. S12



Supplementary Figure S12. Analysis of transcriptome-wide changes in GC cell lines after *LINC01614* over-expression or CRISPR-Cas9 knockout.

A-B Fold change of gene expression in GC cell lines after *LINC01614* over-expression (A) or CRISPR-Cas9 knockout (KO) (B) compared to the controls.

C-D Overlap between the number of genes that are positively associated with *LINC01614* by coexpression analysis and the genes that are strongly affected by *LINC01614* over-expression (C) or *LINC01614* KO (D).

E Cancer hallmark gene sets that are significantly altered in all *LINC01614* over-expression and KO experiments (NOM p-value < 0.05 and FDR q-value < 0.25). The x-axis shows the normalized enrichment score by Gene Set Enrichment Analysis of each experiment.