# Supplemental information

# Diversity of ribosomes at the level of rRNA

# variation associated with human

# health and disease

Daphna Rothschild, Teodorus Theo Susanto, Xin Sui, Jeffrey P. Spence, Ramya Rangan, Naomi R. Genuth, Nasa Sinnott-Armstrong, Xiao Wang, Jonathan K. Pritchard, and Maria Barna
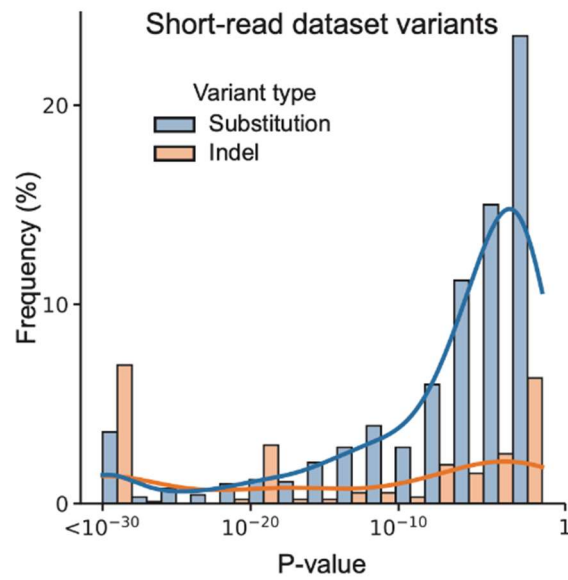
# Supplementary Figures



**Figure S1.  Indels are associated with lower P-values compared to SNVs, related to Figure 1.**

Distribution of Mutect2 false discovery rate (FDR)-corrected log 10 likelihood ratio scores of variant existence measured for substitutions and indels.
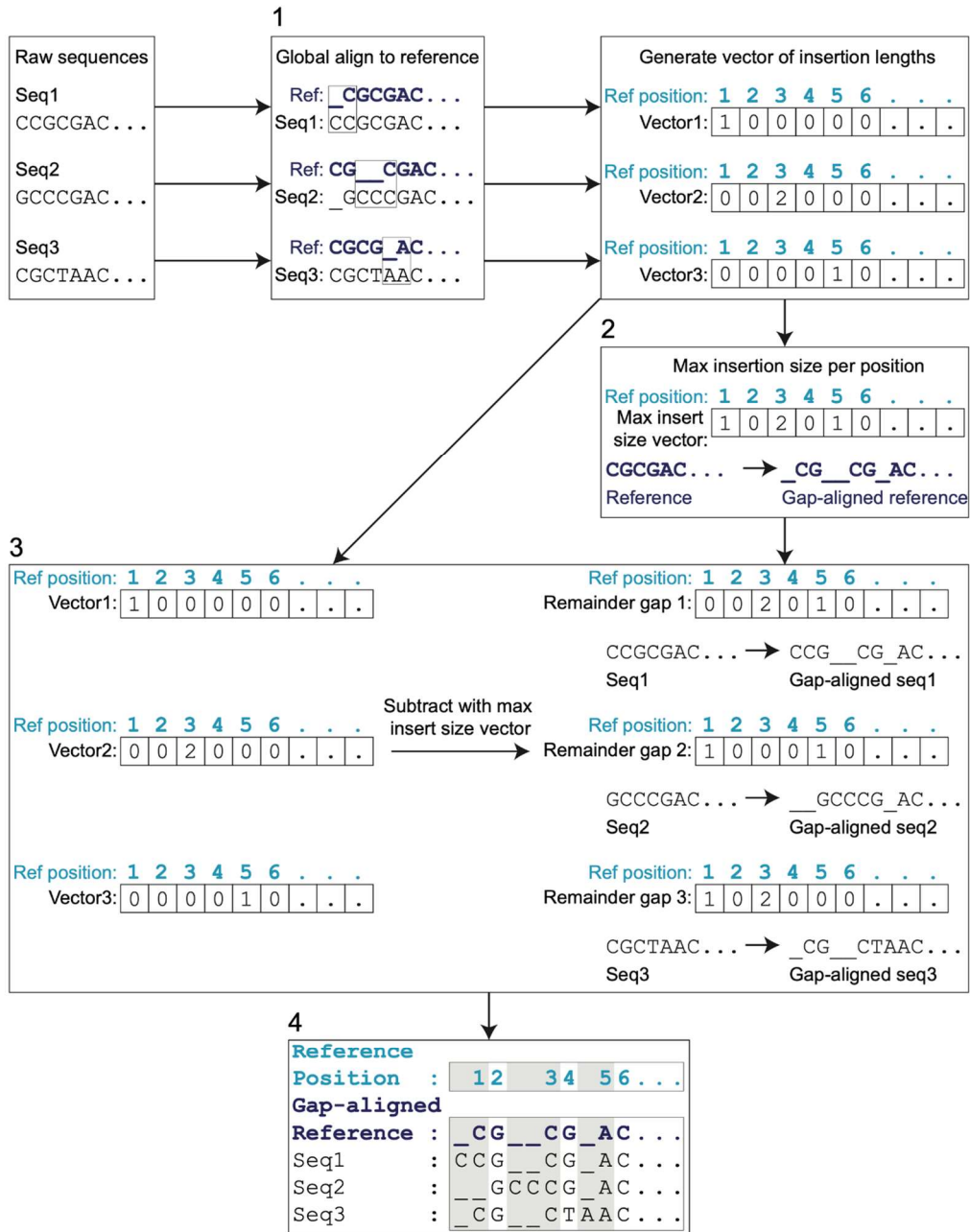
**Figure S2. Reference Gap Alignment (RGA) steps, related to Figure 1.**

The four steps in the RGA method are illustrated with matching numbers to the steps in the main text.
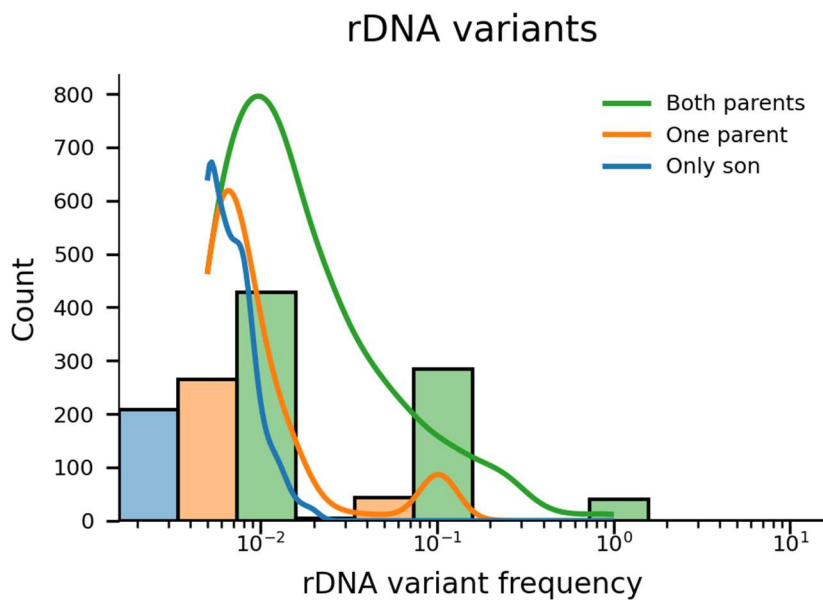
**Figure S3. rDNA variants frequecies from a progeny sample support variant heritability, related to Figure 1.**

rDNA variant frequencies of the GIAB Ashkenazy son are shown in a histogram and are color coded in blue if the are not found in neither parents, orange if found in only one of the parents and in green if found in both.
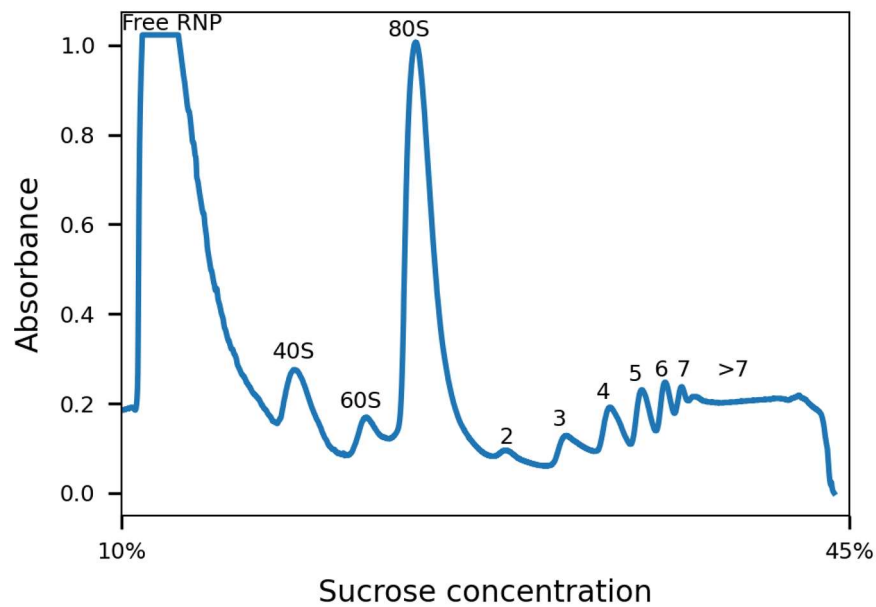
**Figure S4. polysome profile from H7-hESC 10-45% sucrose gradient fractionation, related to Figure 1.**

H7-hESC A260 trace showing the free ribonucleoproteins (RNP), free 40S and 60S subunits, 80S monosomes, and polysomes (marked with 2-7 and >7).
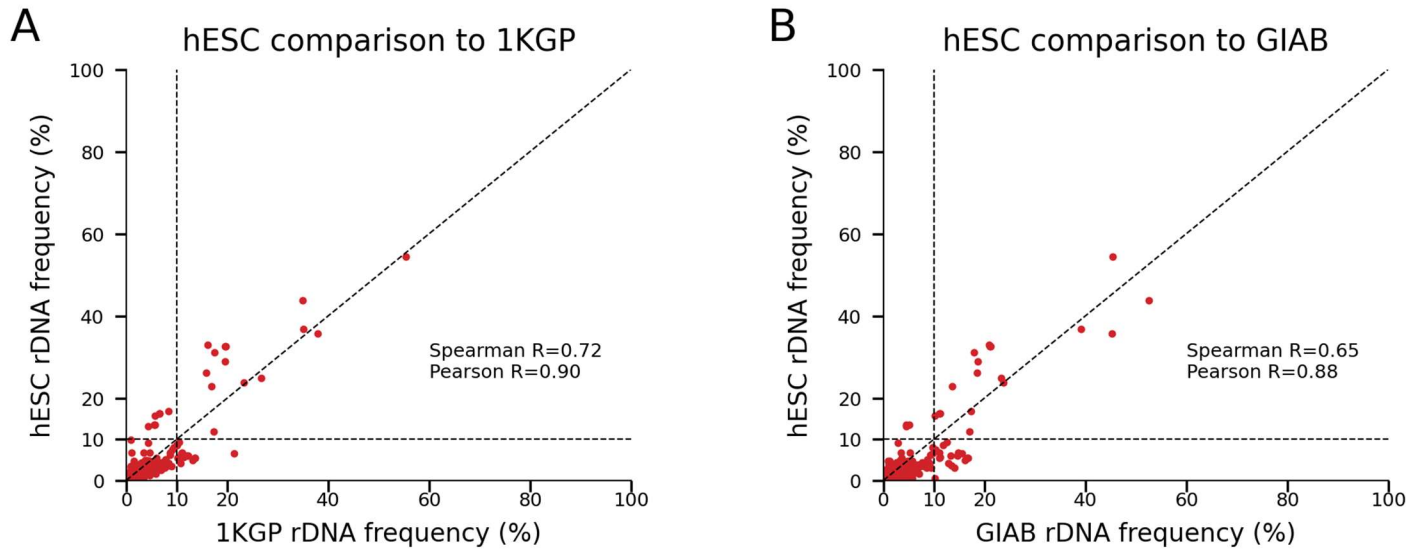


**Figure S5: 28S rDNA variant frequencies comparing 1KGP and GIAB to the H7-hESC, related to Figure 1.**

A. Scatter plot comparing 28S rDNA variant frequencies found in long-read of the 1KGP (x-axis) and in the H7-hESC (y-axis). Pearson and Spearman correlations are indicated.

A

40S with annotated 18S atlas variants

Decoding Center

Head

es6s

Shoulder

B

60S with annotated 28S atlas variants

Exit Tunnel

es27l

es39l

P-stalk

L7/L12
stalk
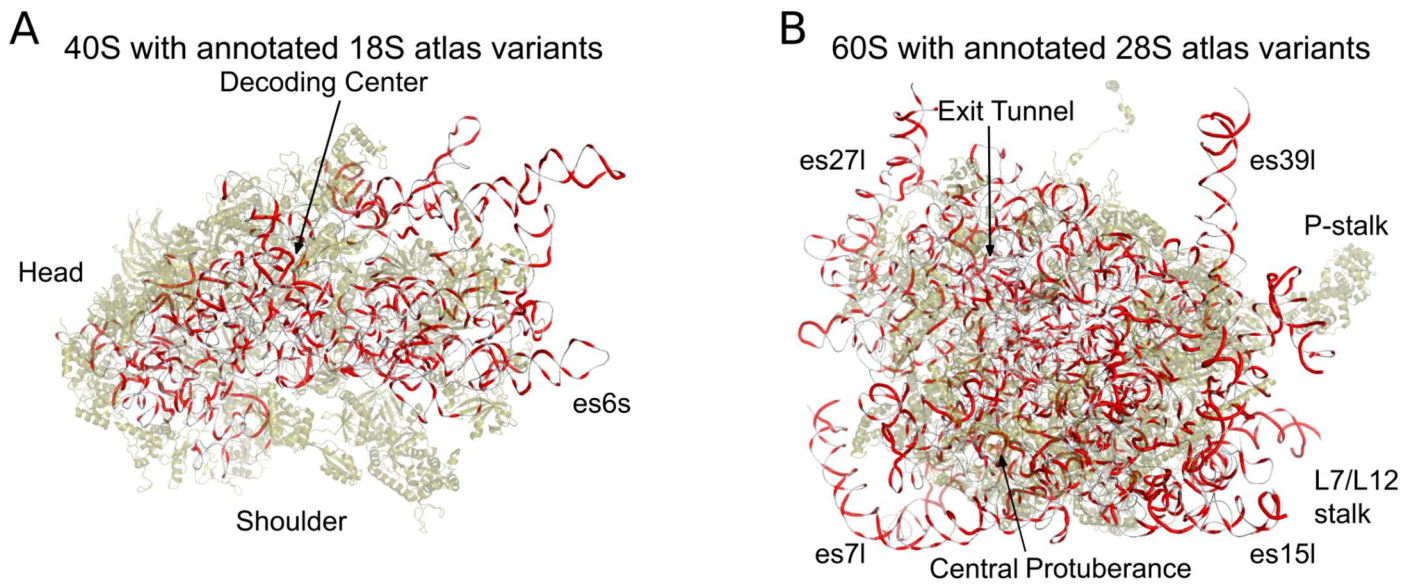
es7l

Central Protuberance

es15l

**Figure S6. 3D structure of the 40S and 60S with variants, related to Figure 1.**

A. 40S 3D structure with 18S atlas variants. Ribosomal proteins are presented in semitransparent green, 18S rRNA in gray. Positions with variants are highlighted in red.
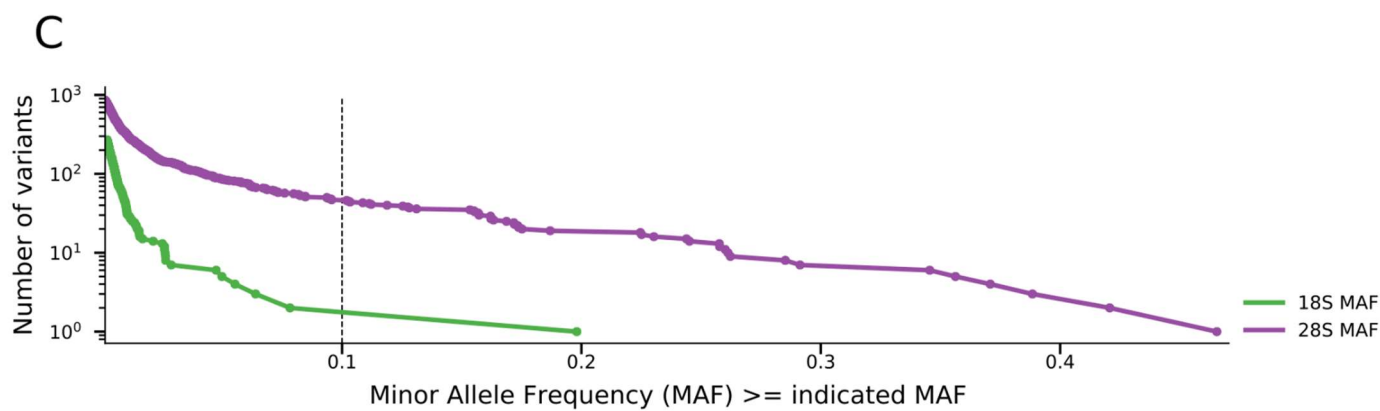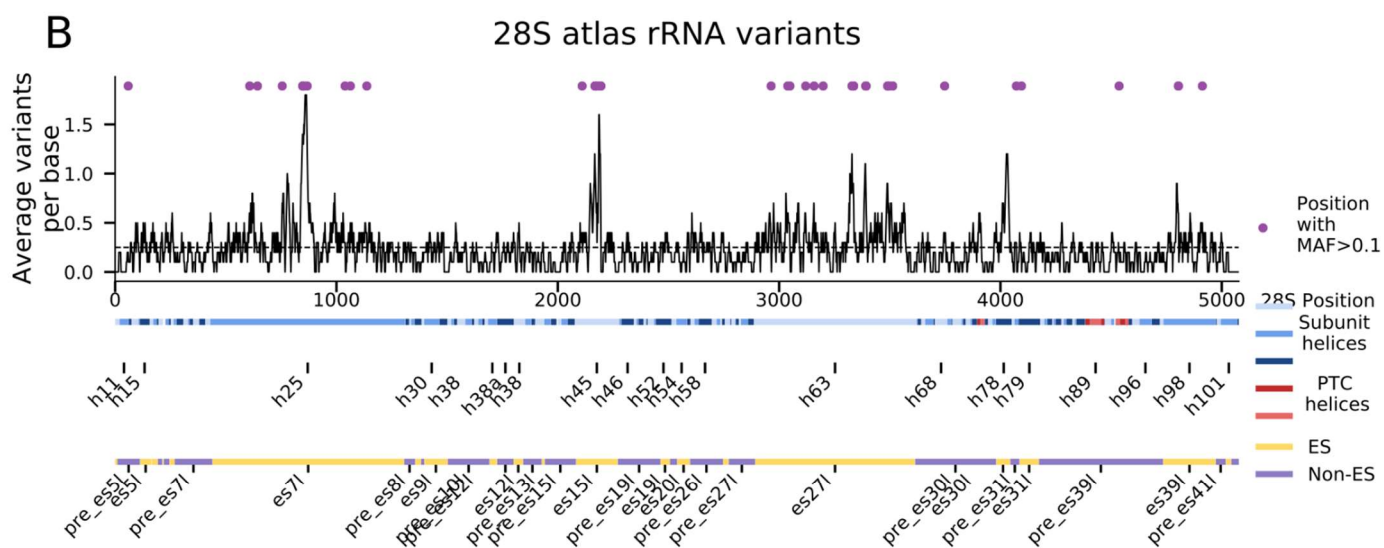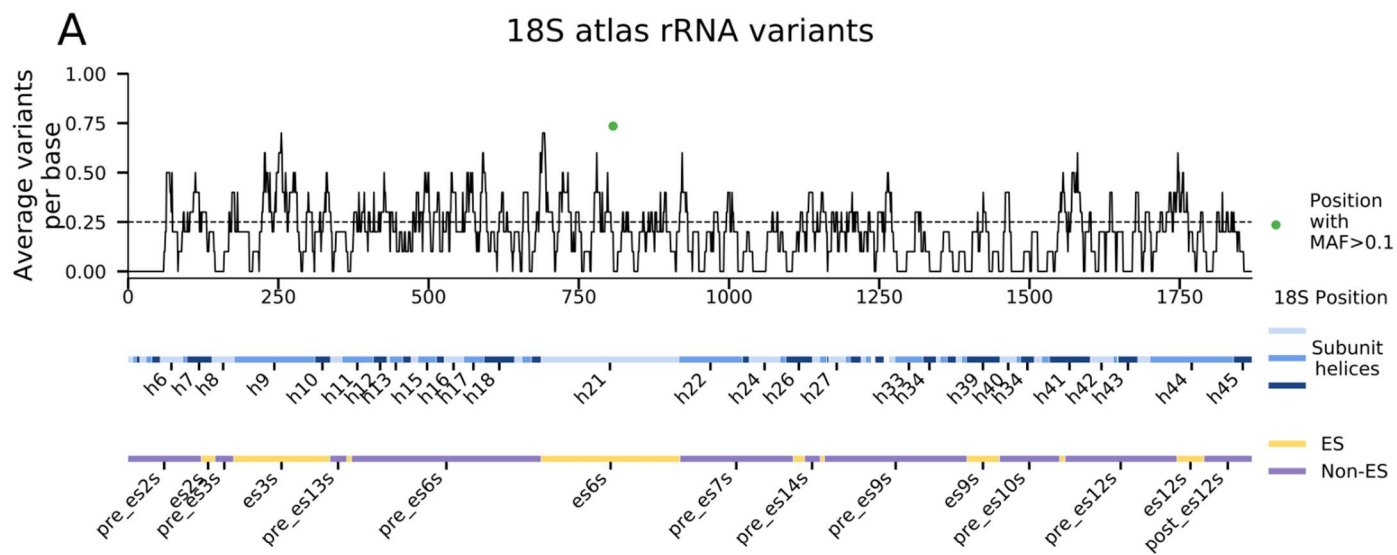
B. Same as (A) for the 60S and the 28S atlas variants.

**Figure S7. 18S and 28S variant distributions, related to Figure 1.**

A. Number of variants at 18S positions. X axis is the nucleotide position along the 18S. Below the X axis are annotations for 18S helix regions and ES regions. Y axis is the average number of atlas variants at a given position for a window size of 20 bases. Green dot at x = 807 (at helix h21 or es6s) is a position with a rRNA variant with minor allele frequency (MAF) >0.1. This dot is also plotted in panel (C).
B. Same as (A) but for the 28S. Purple dots annotate all positions with a variant with MAF>0.1. All purple dots are presented in panel (C).
C. rRNA allele frequency spectrum plot. Values in the X-axis indicate a MAF, and the Y axis are the number of variants with at least the X-axis matched MAF. For example, there are two variants with MAF >0.4 and 6 variants with MAF>0.3. The green plot matches 18S variants and purple plot matches the 28S variants. Individual variants are marked with a dot on the line plot. Dashed line indicates MAF equal to 0.1. The green and purple dots which are in panels (A) and (B) are on the right side of the dashed black line marking MAF=0.1.
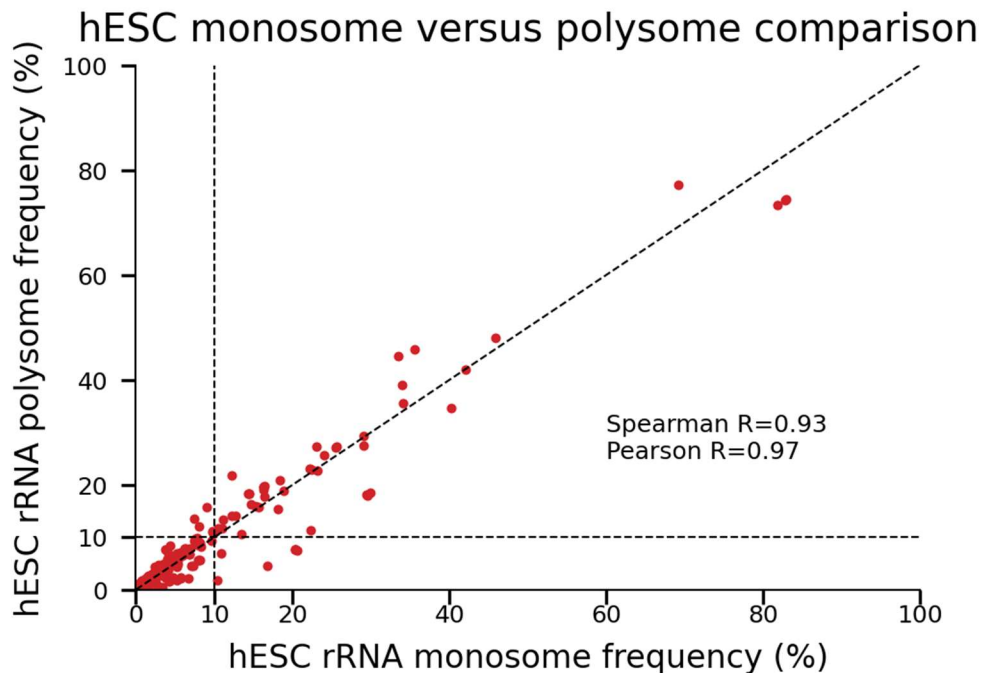


**Figure S8. H7-hESC 28S rRNA variant frequencies in monosomes and polysomes, related to Figure 1.**

Scatter plot of 28S rRNA variant frequencies from the H7-hESC comparing polysome fractions. In the x-axis rRNA variant frequencies are calculated from monosomes and in the y-axis from polysomes. Pearson and Spearman correlations are indicated.
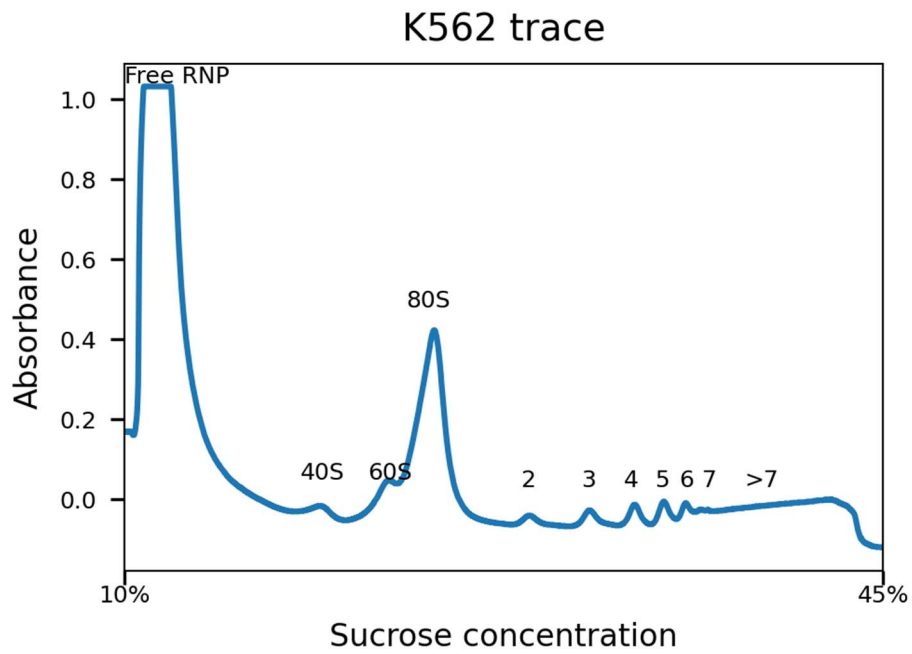
**Figure S9. polysome profile from K562 cell-line in 10-45% sucrose gradient fractionation, related to Figure 1.**

K562 cancer cell-line A260 trace showing the free ribonucleoproteins (RNP), free 40S and 60S subunits, 80S monosomes, and polysomes (marked with 2-7 and >7).
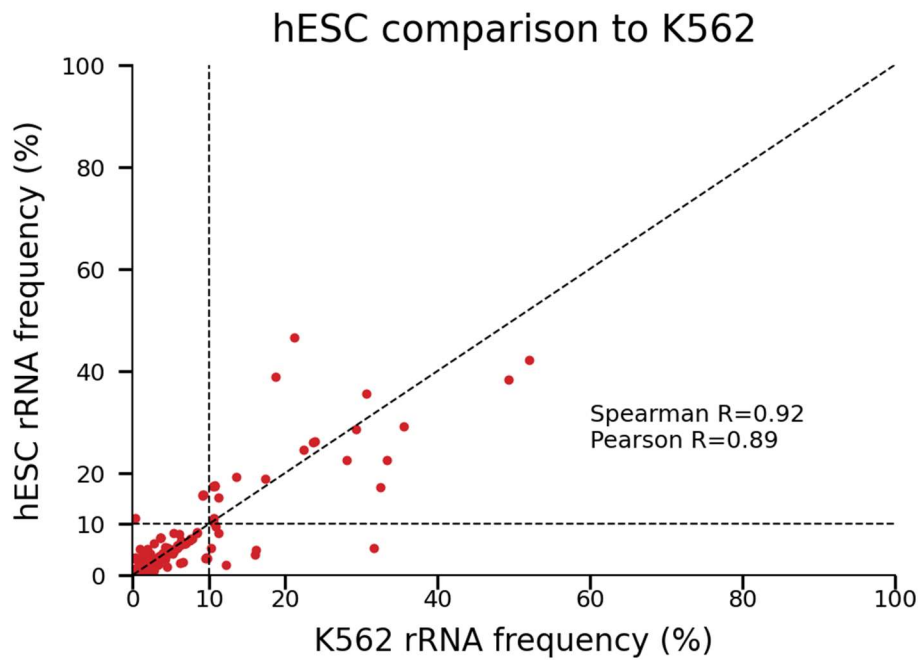
## hESC comparison to K562

Spearman R=0.92
Pearson R=0.89

**Figure S10. 28S rRNA variant frequencies comparing K562 to the H7-hESC, related to Figure 1.**

Scatter plot comparing 28S rRNA variant frequencies found in the K562 cancer cell-line (x-axis) and in the H7-hESC (y-axis). Pearson and Spearman correlations are indicated.
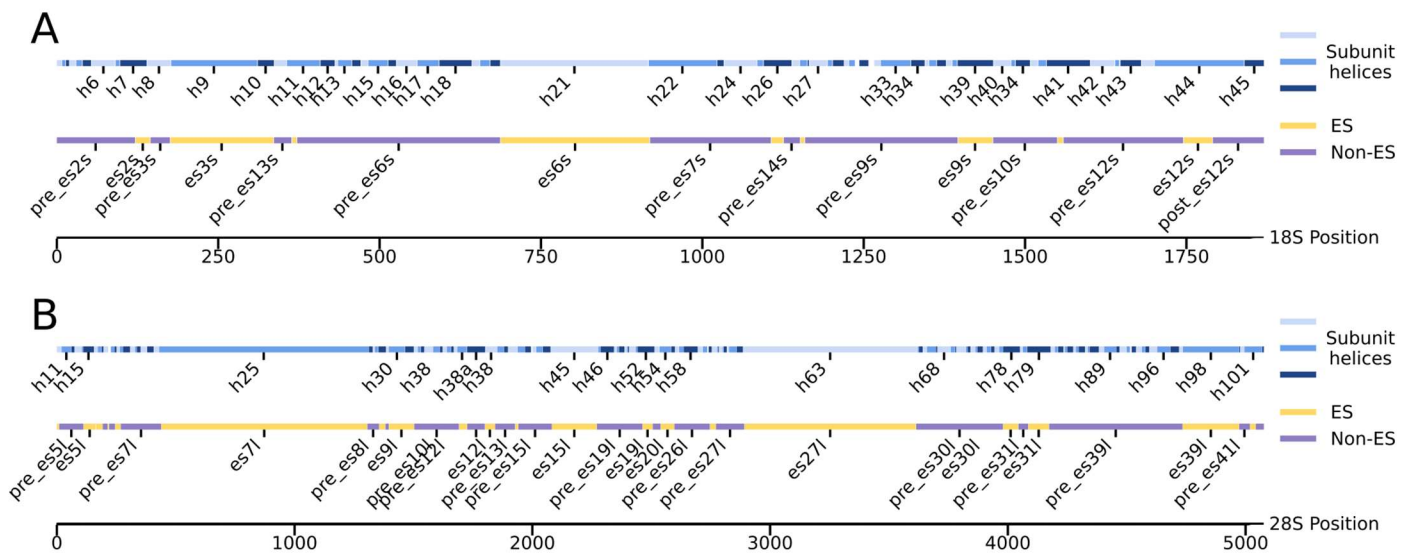


**Figure S11. 18S and 28S region annotation, related to Figure 1.**

A.  18S helix and ES region annotation. Only helices and ES regions with at least 20 bases are labeled

B. 28S helix and ES region annotation. Only helices and ES regions with at least 40 bases are
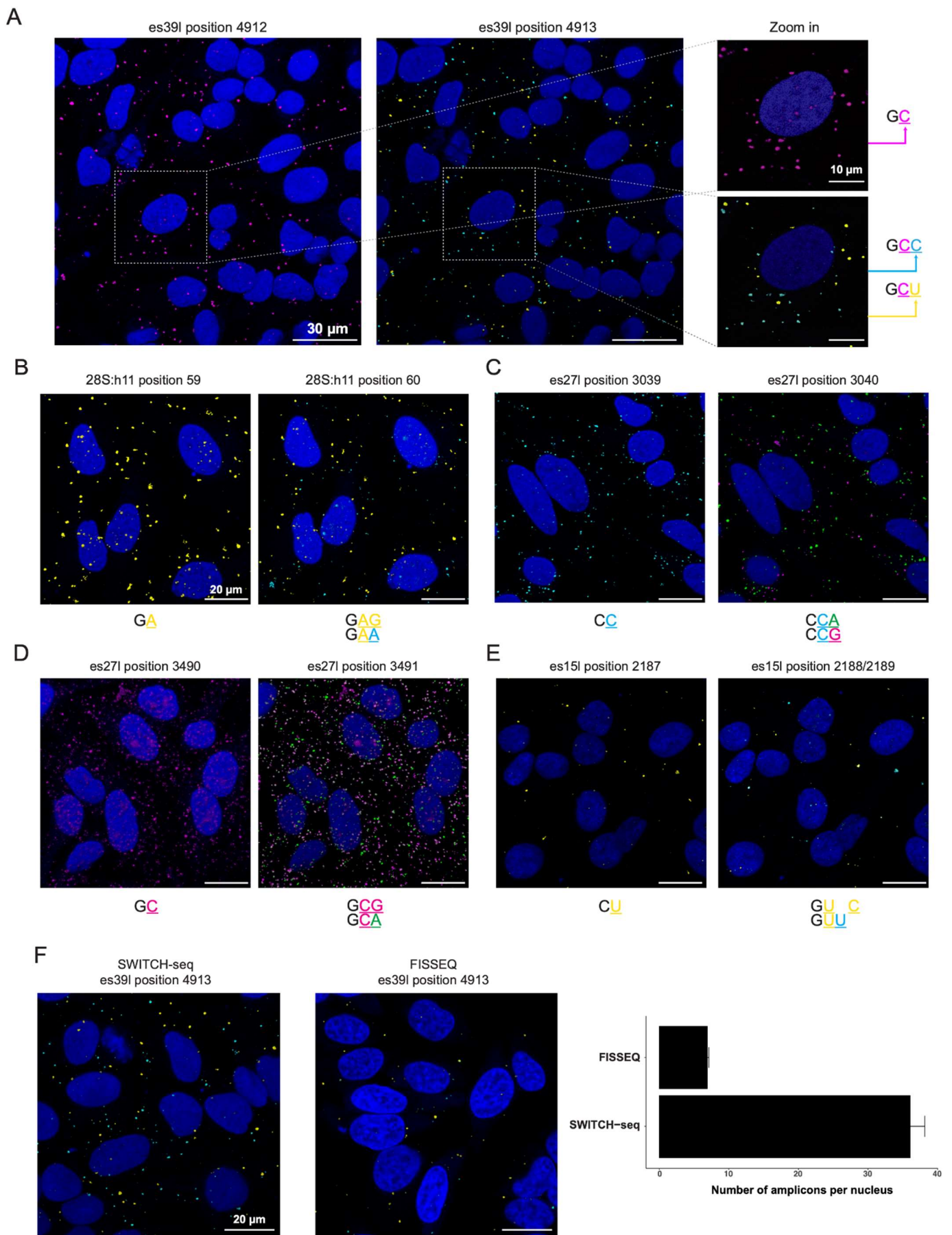   labeled

A

es39I position 4912          es39I position 4913          Zoom in

30 μm          10 μm          G**C**

GC**C**
GC**U**

B

28S:h11 position 59          28S:h11 position 60

20 μm

G**A**          G**AG**
          G**AA**

C

es27I position 3039          es27I position 3040

C**C**          CC**A**
          CC**G**

D

es27I position 3490          es27I position 3491

G**C**          GC**G**
          GC**A**

E

es15I position 2187          es15I position 2188/2189

C**U**          G**U**  **C**
          G**UU**

F

SWITCH-seq          FISSEQ
es39I position 4913          es39I position 4913

20 μm

FISSEQ

SWITCH−seq

0          10          20          30          40
Number of amplicons per nucleus

**Figure S12.** *In-situ* **sequencing of rRNA variants, related to Figure 2.**

A. Two rounds of representative fluorescent *in situ* sequencing images of HeLa cells (DAPI staining in blue) are presented for the es39l-probed region. We identified a non-variable base C (magenta) at position 4912. At position 4913, two alternative sequences were revealed: the known reference sequence C (cyan) and the alternative variant U (yellow).

B. Similar to (A) for 28S:h11 where G and A are detected at position 60

C. Similar to (A) for es27l where A and G are detected at position 3040

D. Similar to (A) for es27l where G and A are detected at position 3491

E. Similar to (A) for es15 where a U insertion is detected at position 2188.

F. Representative fluorescence images comparing SWITCH-seq and FISSEQ in the detection of the rRNA variant at es39l position 4913.
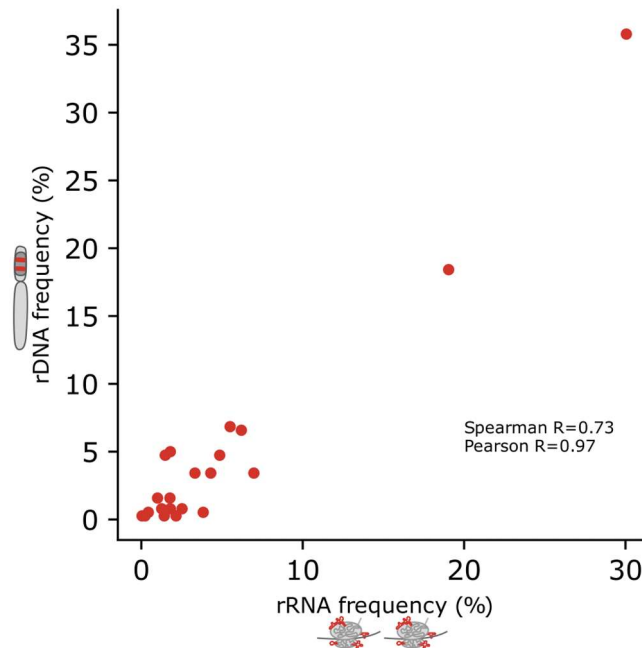


**Figure S13. 28S rRNA subtype frequencies in h7-hESC, related to Figure 3.**

Scatter plot showing the frequencies of 28S rRNA subtypes in rRNA (x-axis) and rDNA (y-axis) in H7-hESC
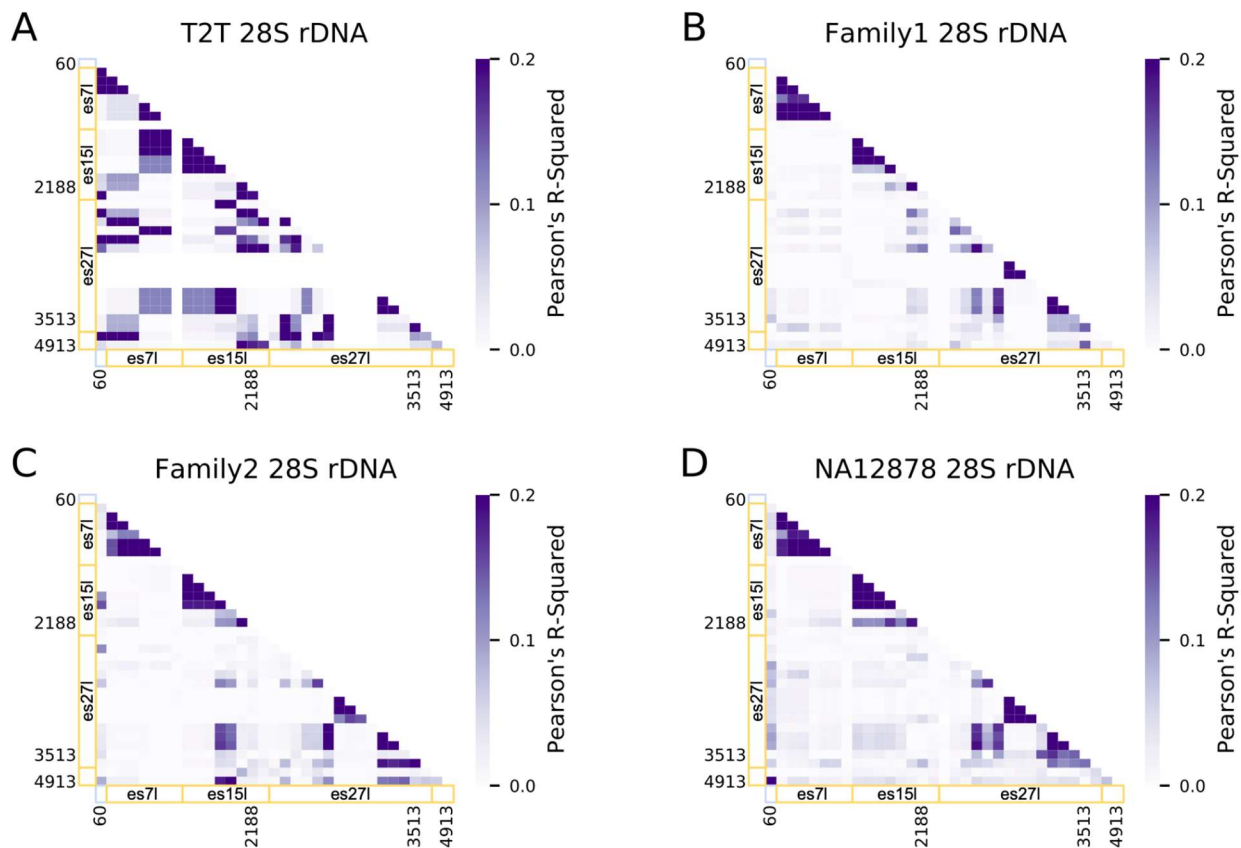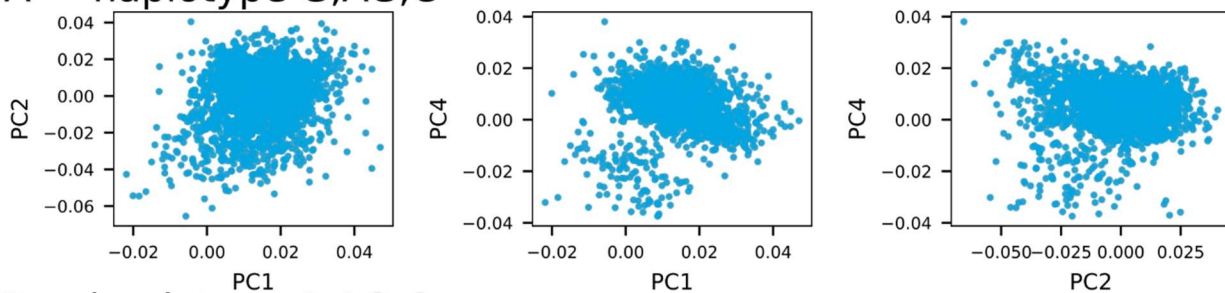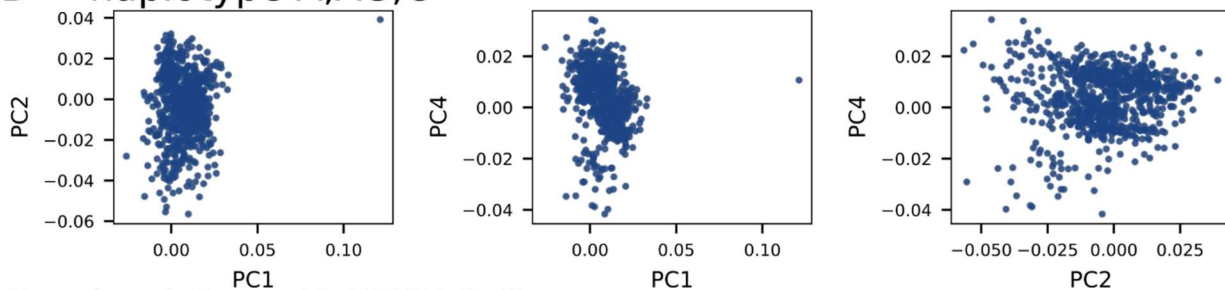
**Figure S14. 28S subtypes in T2T and GIAB, related to Figure 3.**

A. Correlation coefficient (Pearson's $r^2$) heatmap of T2T 28S rDNA for the same positions analyzed in the h7-hESC (**Figure 1F, 3A**). X-axis and Y-axis are annotated by regions. Helix regions are annotated by light blue and ES regions are annotated by yellow. H7-hESC positions with higher $r^2$ are annotated.

B. Same as (A) for the Chinese Han family trio from the GIAB dataset.

C. Same as (A) for the Ashkenazi family trio from the GIAB dataset.

D. Same as (A) for the HA12878 cell line from the GIAB dataset.
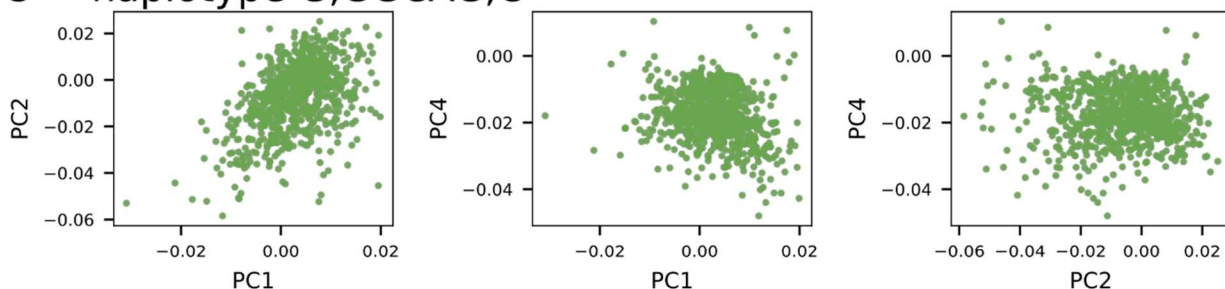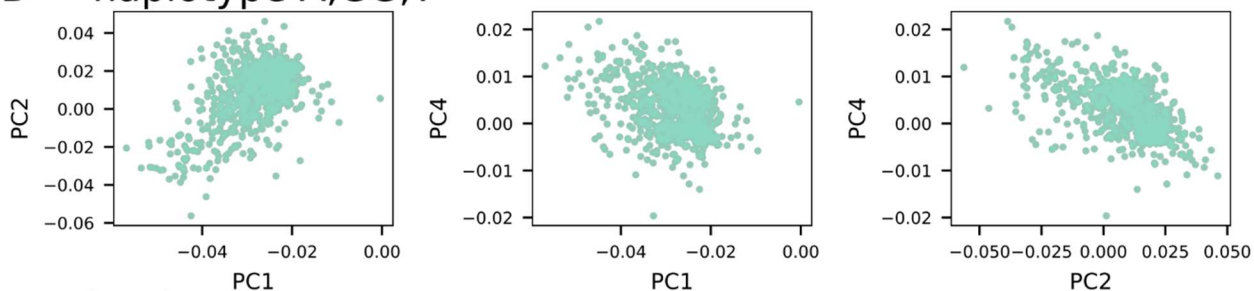
A  haplotype G,AG,C

B  haplotype A,AG,C

C  haplotype G,GGCAG,C
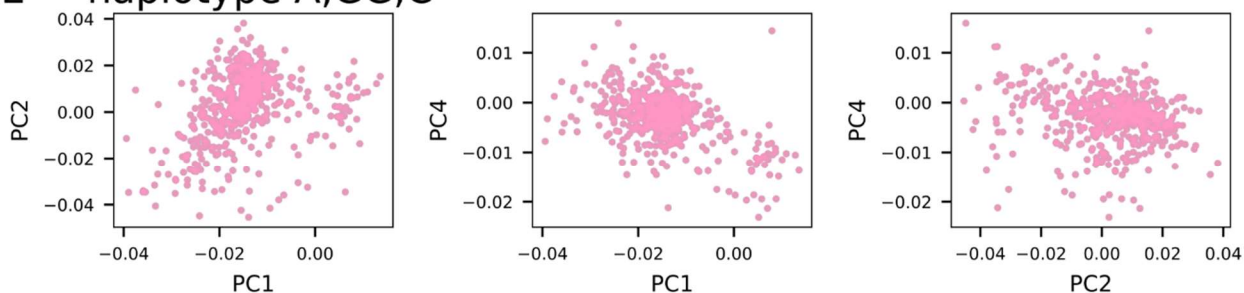
D  haplotype A,GG,T

E  haplotype A,GG,C

**Figure S15. 28S haplotype groups as shown by Bray-Curtis Principal Coordinate Analysis (PCoA), related to Figure 3.**

PCoA of 386 28S rDNA sequences from each GIAB sample. The first, second and fourth PCs are presented for 28S sequences that belong to different haplotypes where each haplotype is presented in a separte panel A-E. Each dot is a complete 28S rDNA sequence with similarity between sequences measured on 6-mers. The colors correspond to coloring an rDNA sequence by its 3 position haplotype shown in the main Figure 3D.
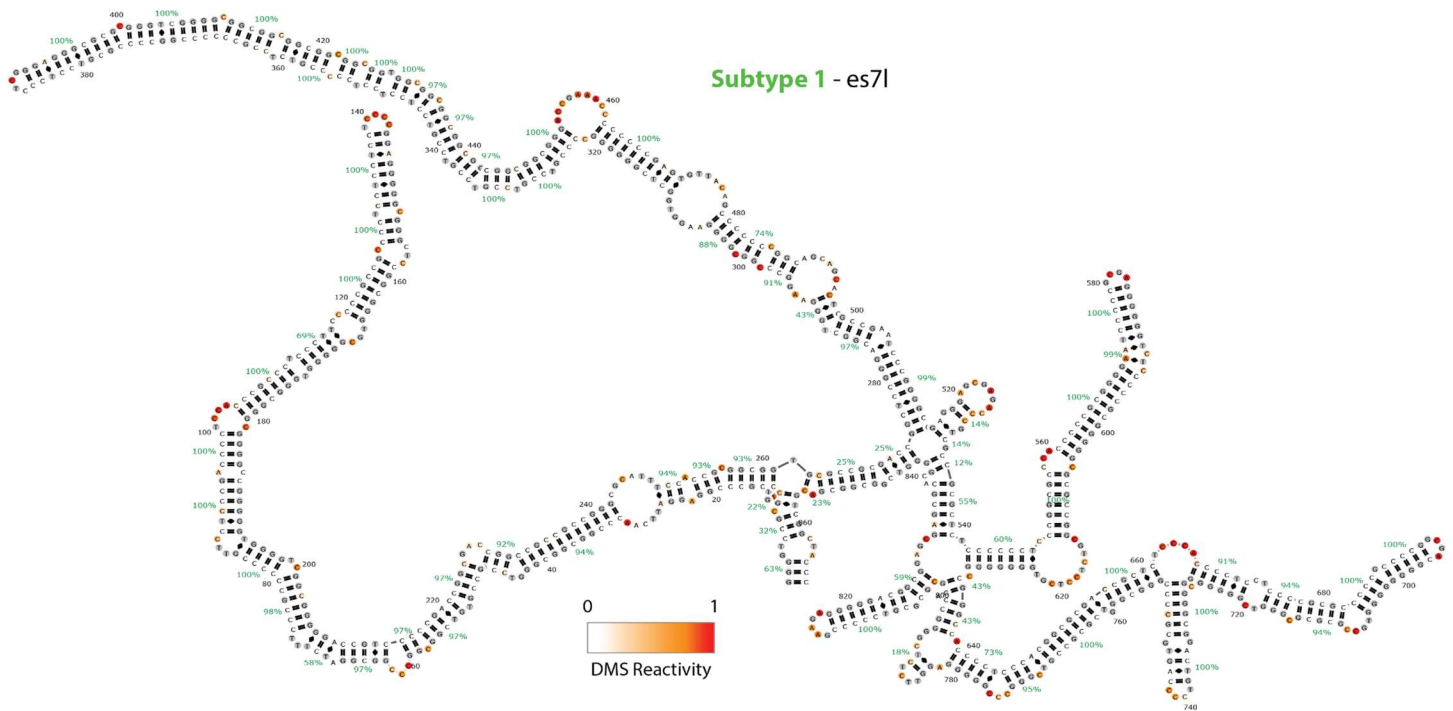


**Figure S16. In-cell DMS with long-read sequencing shows accessibility differences in es7l of different 28S subtypes, related to Figure 4.**

RNA secondary structure of es7l predicted secondary structure for subtype 1 (A, GGCAG, T).
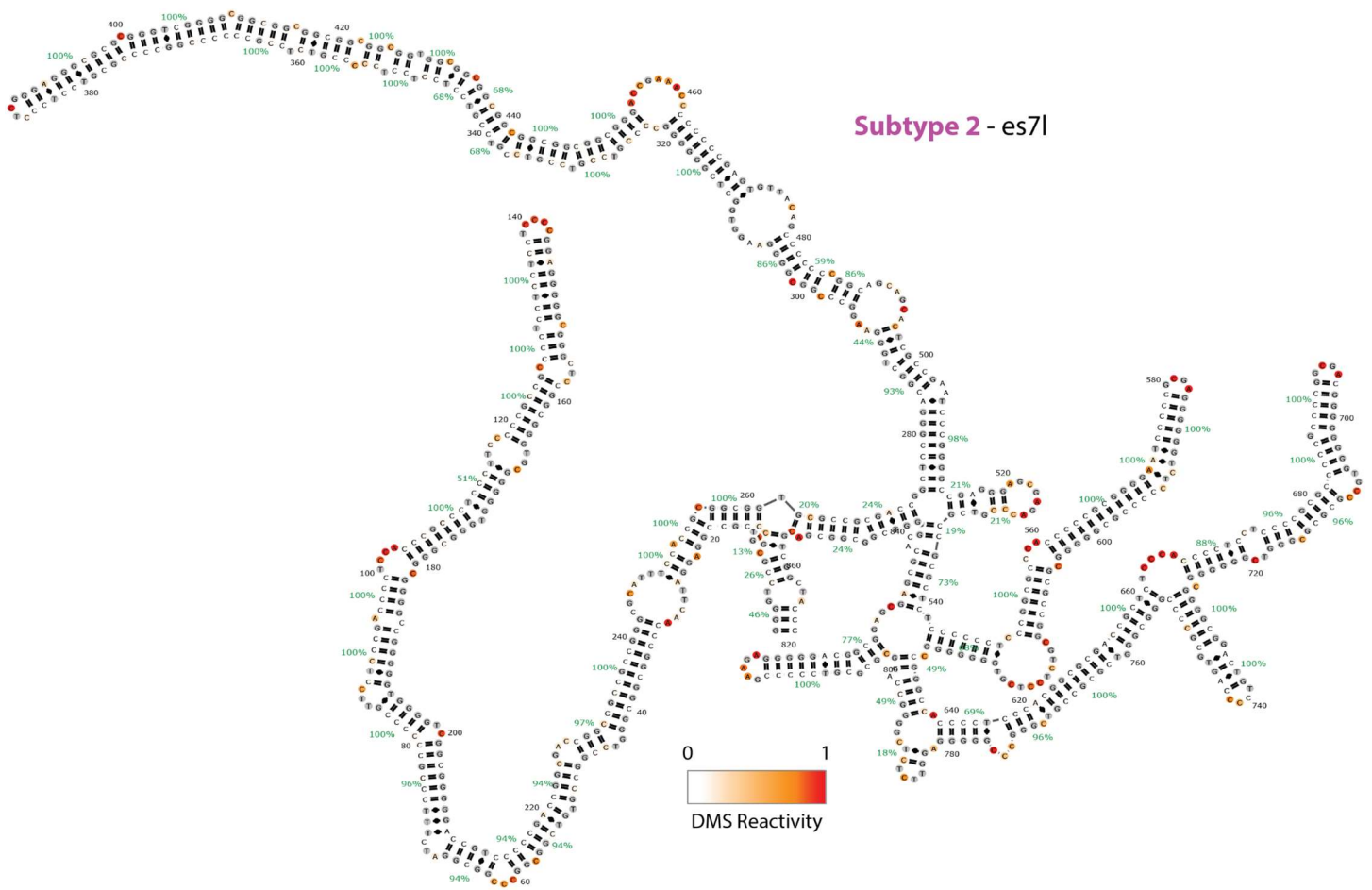
**Figure S17. In-cell DMS with long-read sequencing shows accessibility differences in es7l of different 28S subtypes, related to Figure 4.**

RNA secondary structure of es7l predicted secondary structure for subtype 2 (G, AG, C).
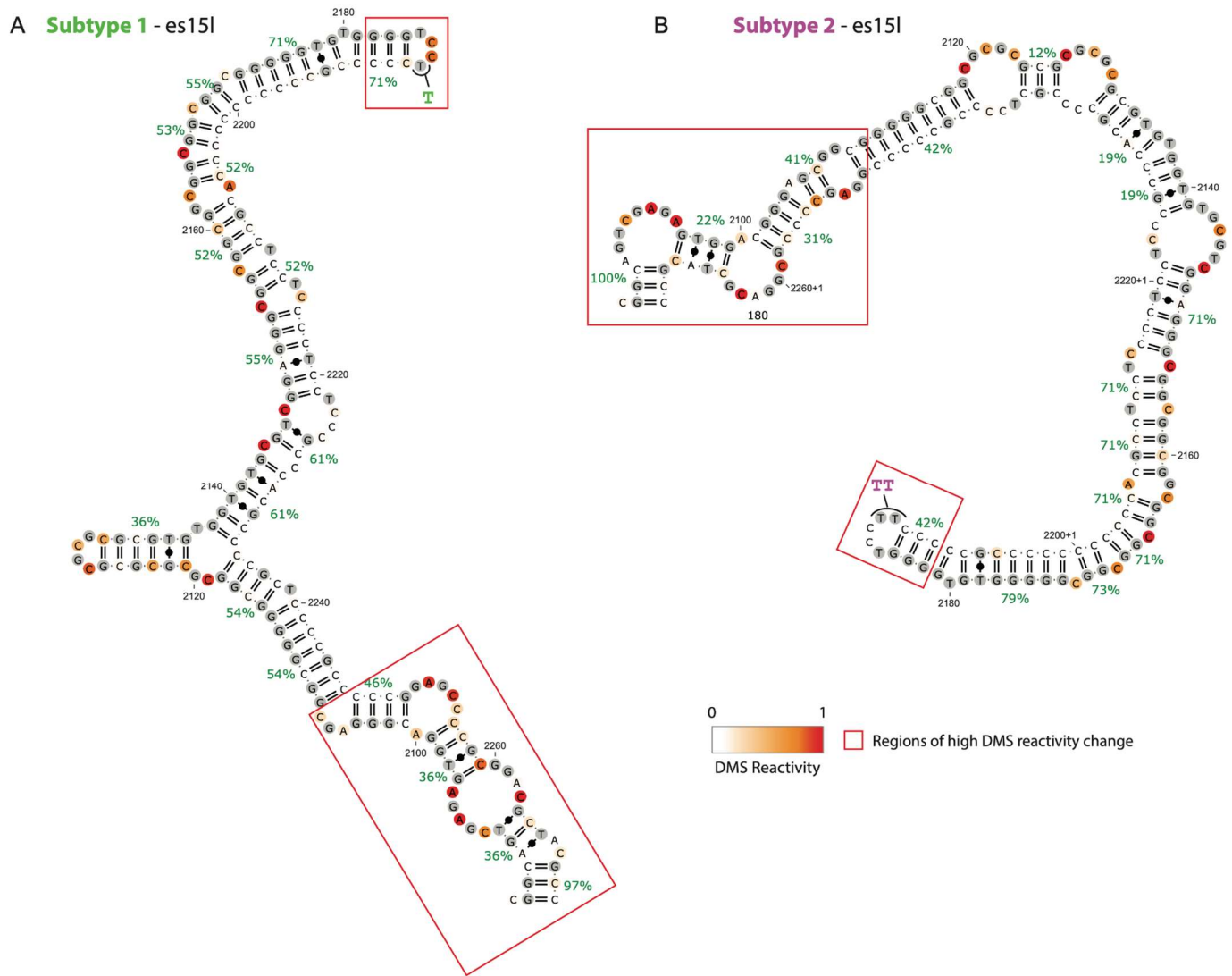
**Figure S18. In-cell DMS with long-read sequencing shows that es15l of different 28S subtypes have different RNA 2D structure, related to Figure 4.**

A. RNA secondary structure of es15l predicted secondary structure for subtype 1 (A, GGCAG, T).
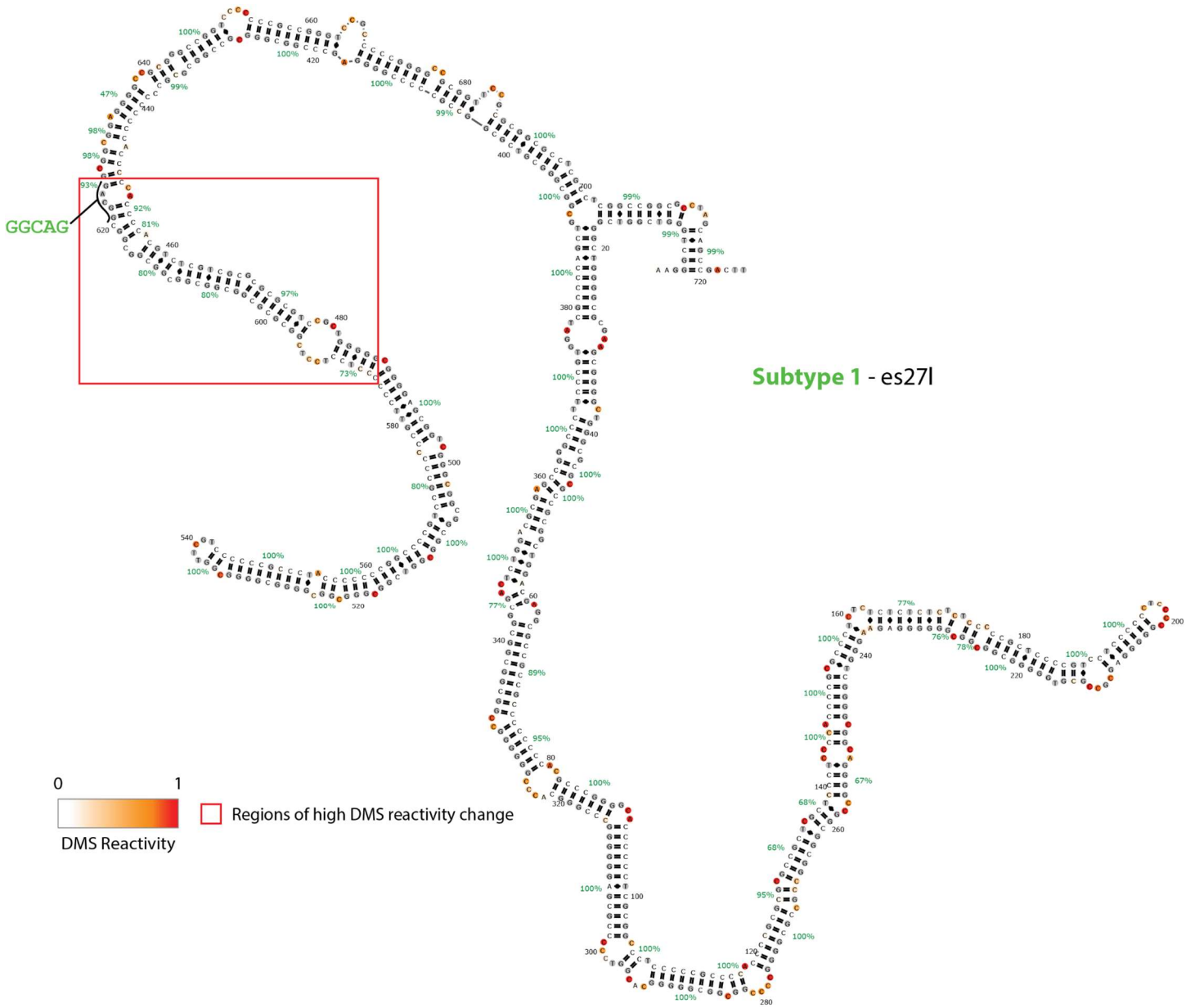B. RNA secondary structure of es15l predicted secondary structure for subtype 2 (G, AG, C).

**Figure S19. In-cell DMS with long-read sequencing shows that es27l of different 28S subtypes have different RNA 2D structure, related to Figure 4.**

RNA secondary structure of es15l predicted secondary structure for subtype 1 (A, GGCAG, T).
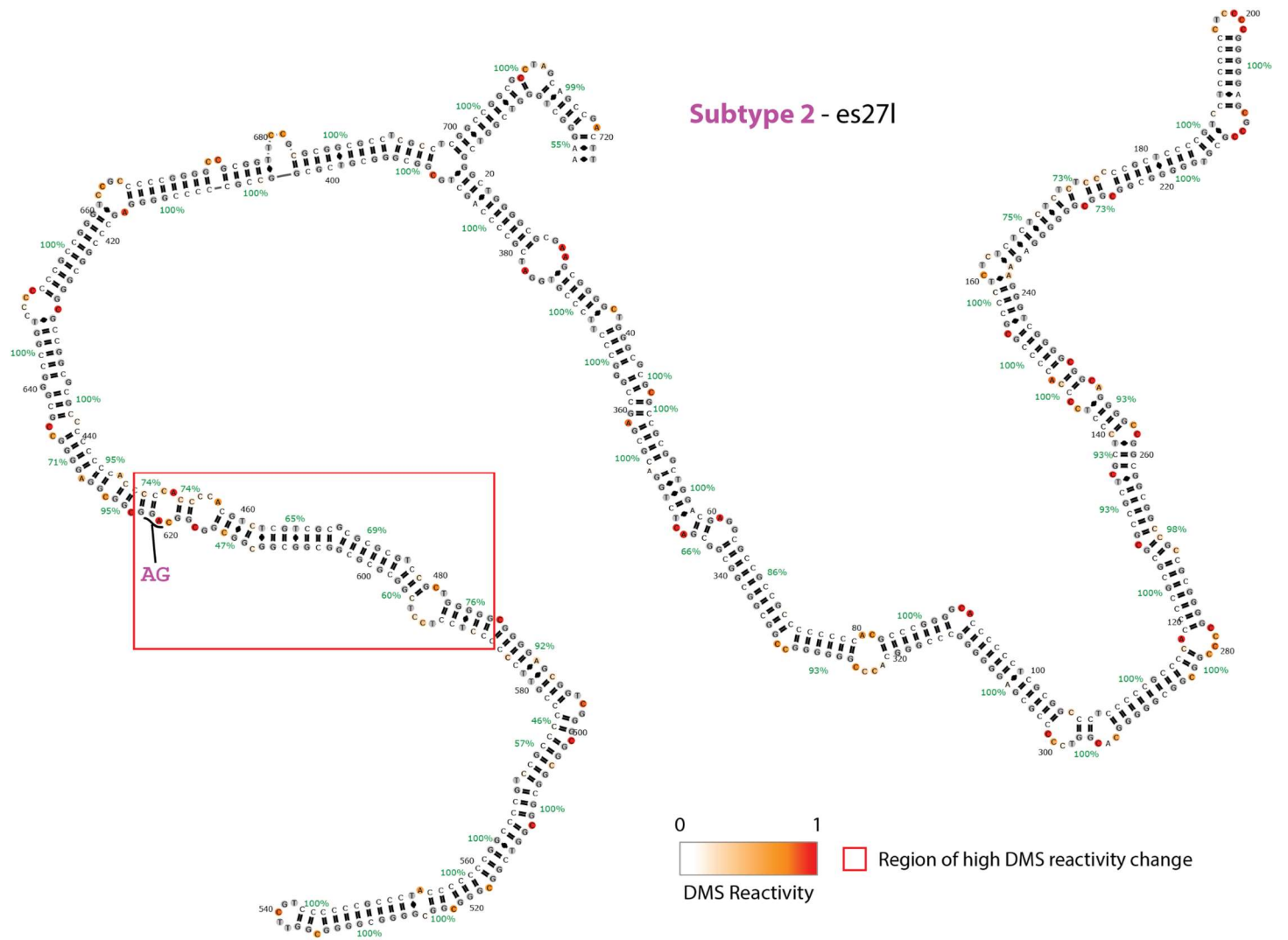
**Figure S20. In-cell DMS with long-read sequencing shows that es27l of different 28S subtypes have different RNA 2D structure, related to Figure 4.**

RNA secondary structure of es15l predicted secondary structure for subtype 2 (G, AG, C).
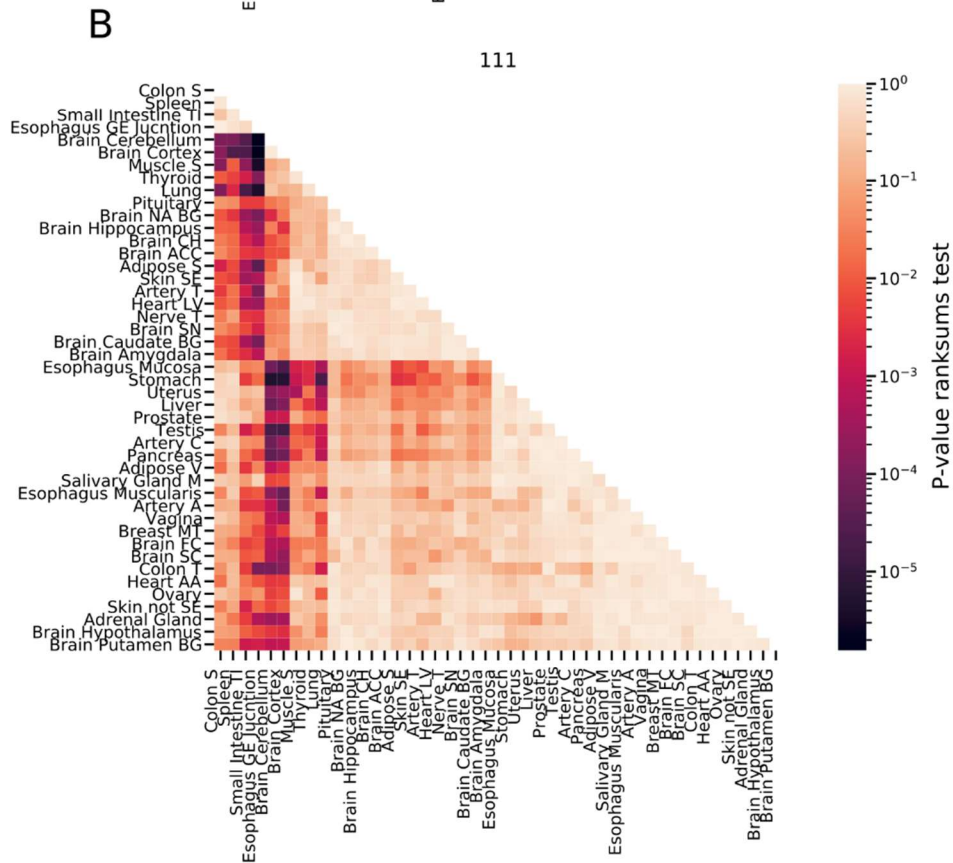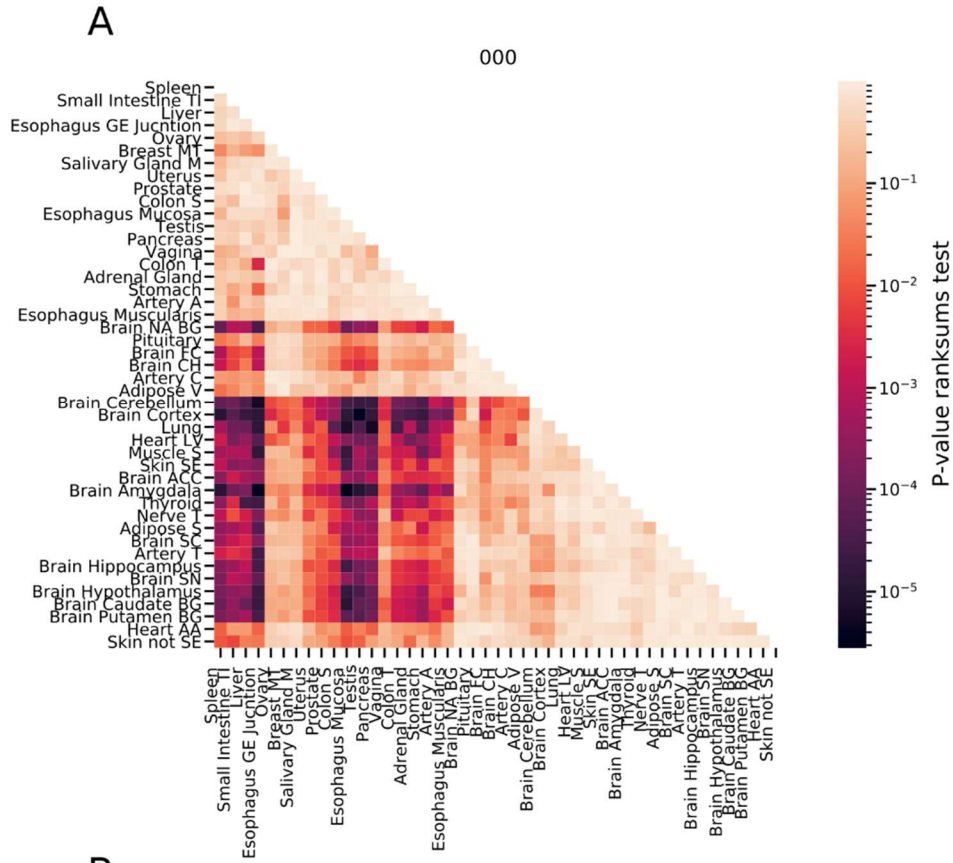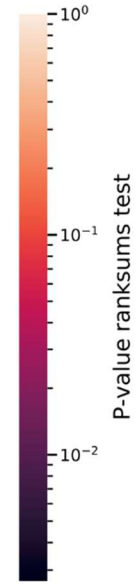
A

000

B

111

**Figure S21 (1 out of 5 similar plots). rRNA subtype expression levels are tissue specific, related to Figure 5.**

A. A heatmap showing FDR-corrected rank sum test P-values comparing the expression levels of the rRNA subtype with the haplotype sequence of G,AG,C (titled 000) across tissues. The tissues are ordered by average hierarchical clustering of the rank sum corrected P-values.

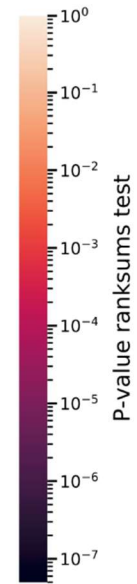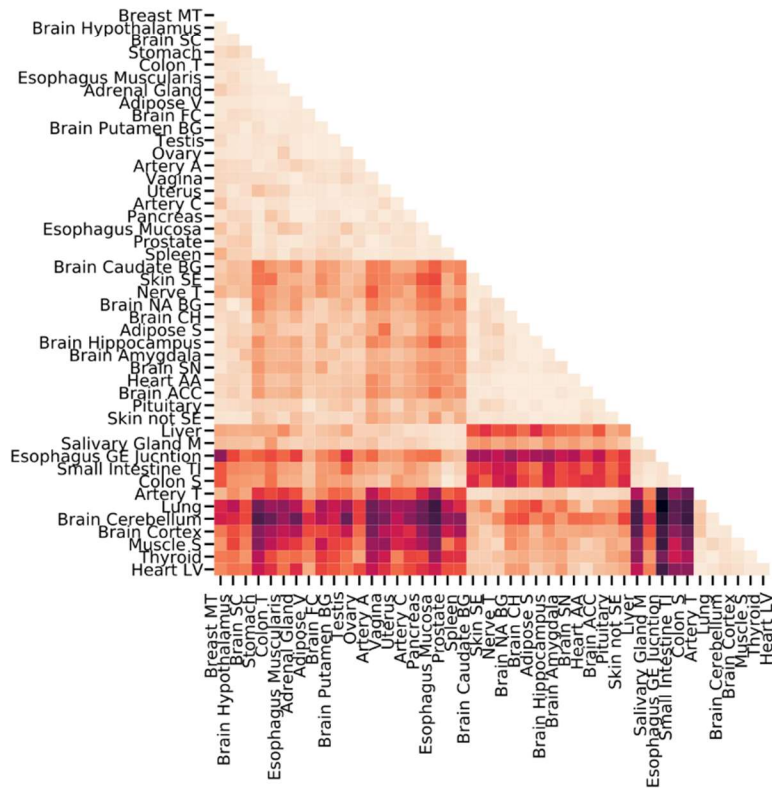B. Same as (A) for haplotype A,GGCAG,T (titled 111) across tissues.

A

010

B

120

**Figure S22 (2 out of 5 similar plots). rRNA subtype expression levels are tissue specific, related to Figure 5.**

A. A heatmap showing FDR-corrected rank sum test P-values comparing the expression levels of the rRNA subtype with the haplotype sequence of G,GGCAG,C (titled 010) across tissues. The tissues are ordered by average hierarchical clustering of the rank sum corrected P-values.

B. Same as (A) for haplotype A,GG,C (titled 120) across tissues.
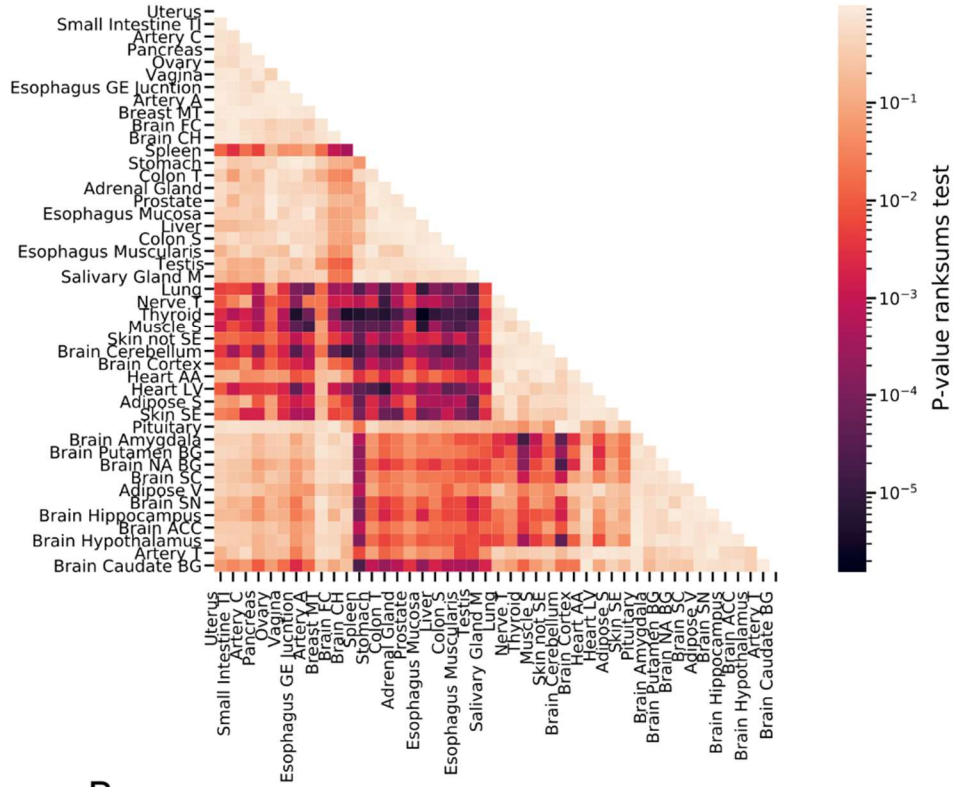
A

021

B

101

**Figure S23 (3 out of 5 similar plots).  rRNA subtype expression levels are tissue specific, related to Figure 5.**

A.  A heatmap showing FDR-corrected rank sum test P-values comparing the expression levels of the rRNA subtype with the haplotype sequence of G,GG,C (titled 021) across tissues. The tissues are ordered by average hierarchical clustering of the rank sum corrected P-values.

B.  Same as (A) for haplotype A,AG,T (titled 101) across tissues.
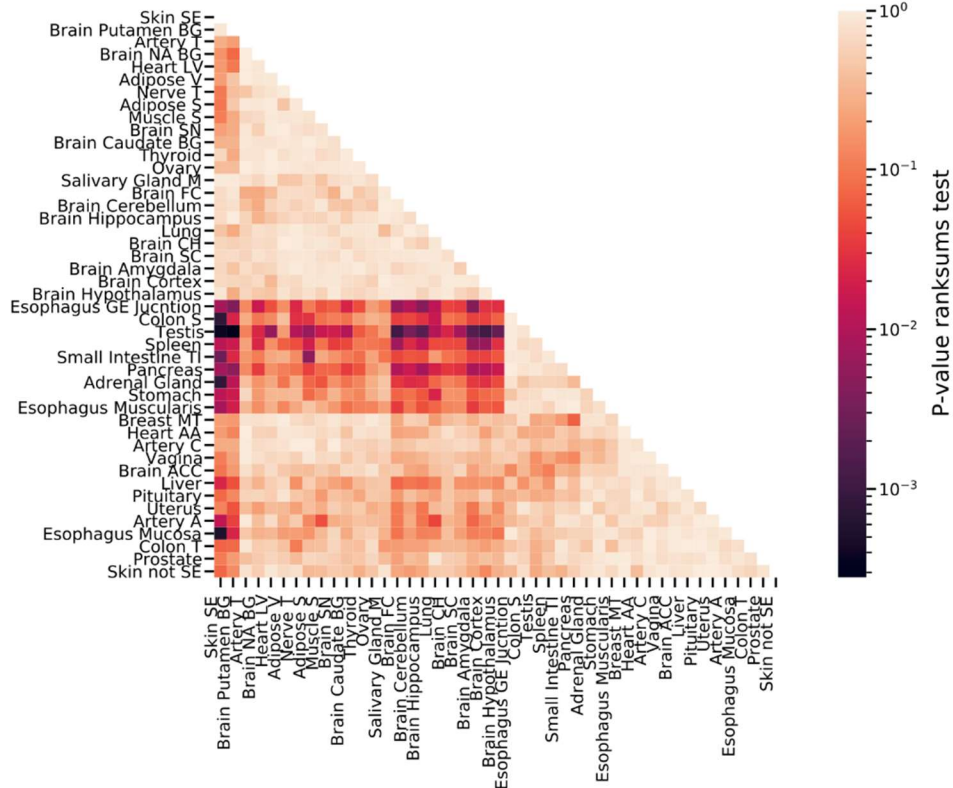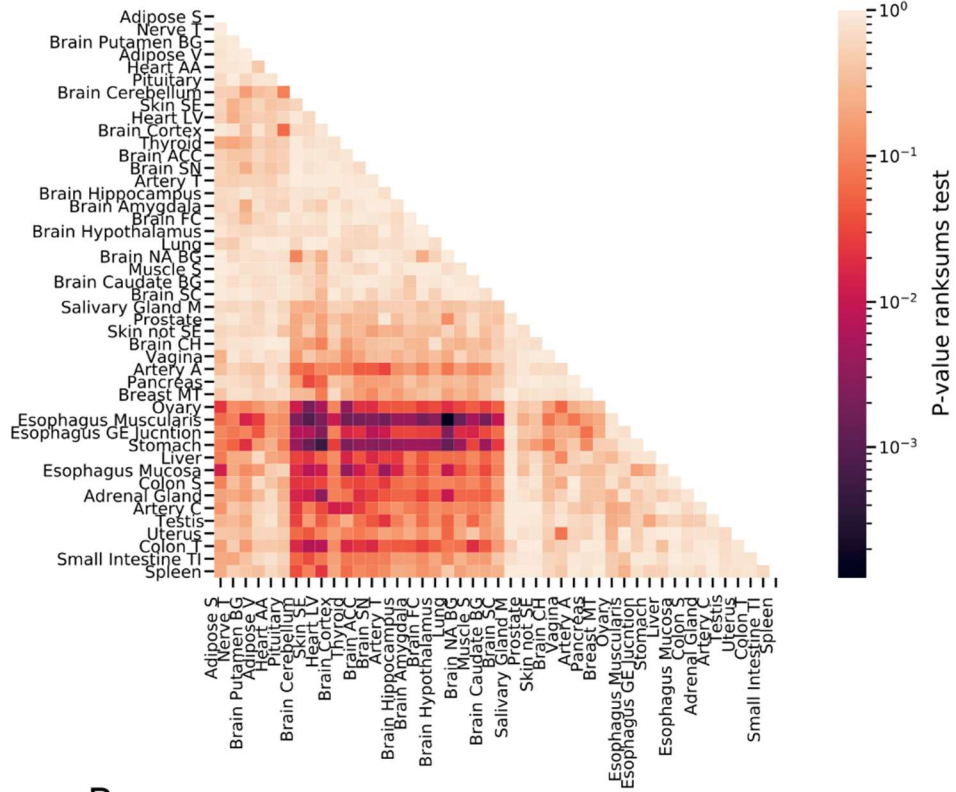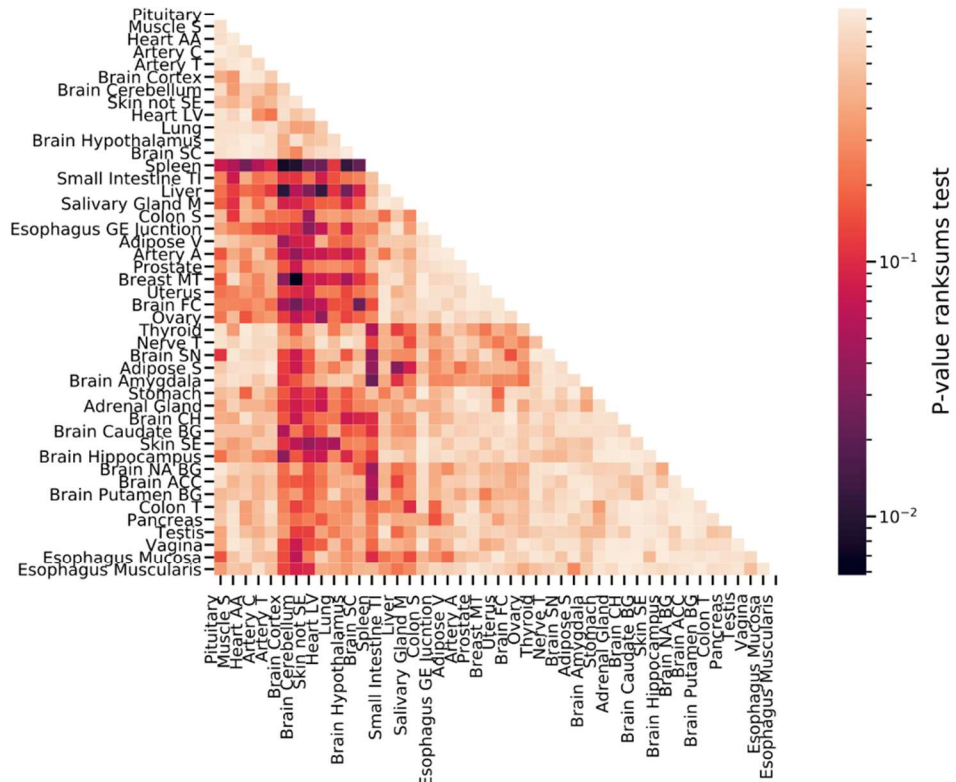
A

121

B

020

**Figure S24 (4 out of 5 similar plots).  rRNA subtype expression levels are tissue specific, related to Figure 5.**

A. A heatmap showing FDR-corrected rank sum test P-values comparing the expression levels of the rRNA subtype with the haplotype sequence of  A,GG,T (titled 121) across tissues. The tissues are ordered by average hierarchical clustering of the rank sum corrected P-values.

B. Same as (A) for haplotype G,GG,C (titled 020) across tissues.
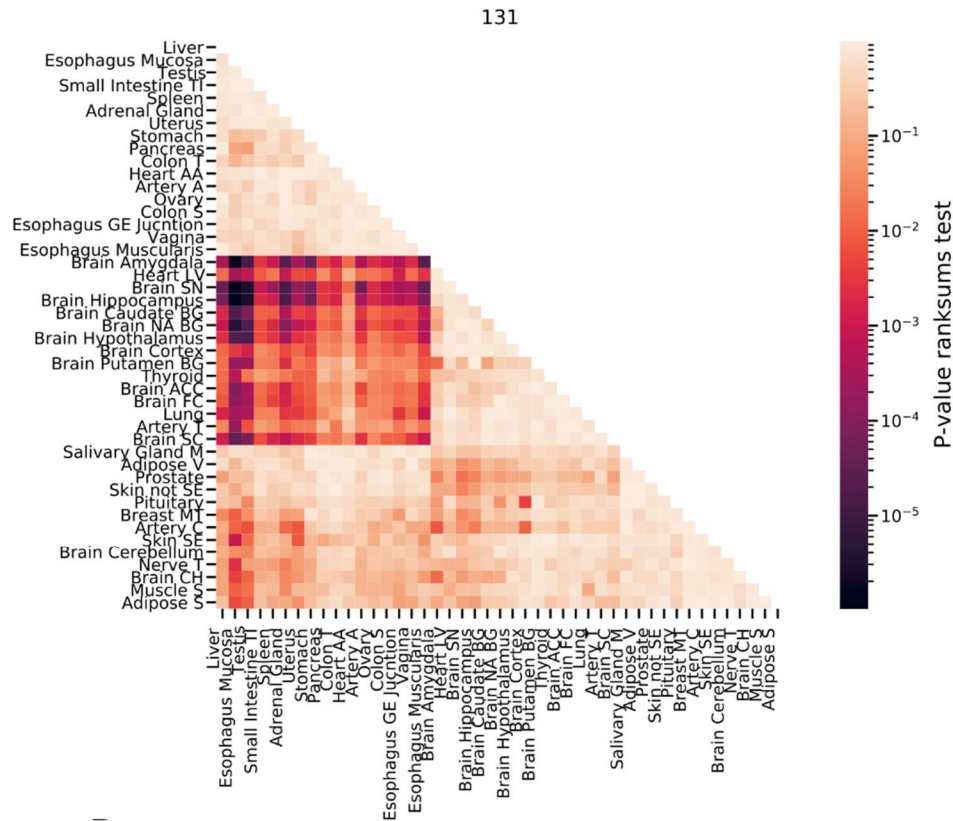


**Figure S25 (5 out of 5 similar plots).  rRNA subtype expression levels are tissue specific, related to Figure 5.**

A heatmap showing FDR-corrected rank sum test P-values comparing the expression levels of the rRNA subtype with the haplotype sequence of  A,GGCGGCAG,T (titled 131) across tissues. The tissues are ordered by average hierarchical clustering of the rank sum corrected P-values.
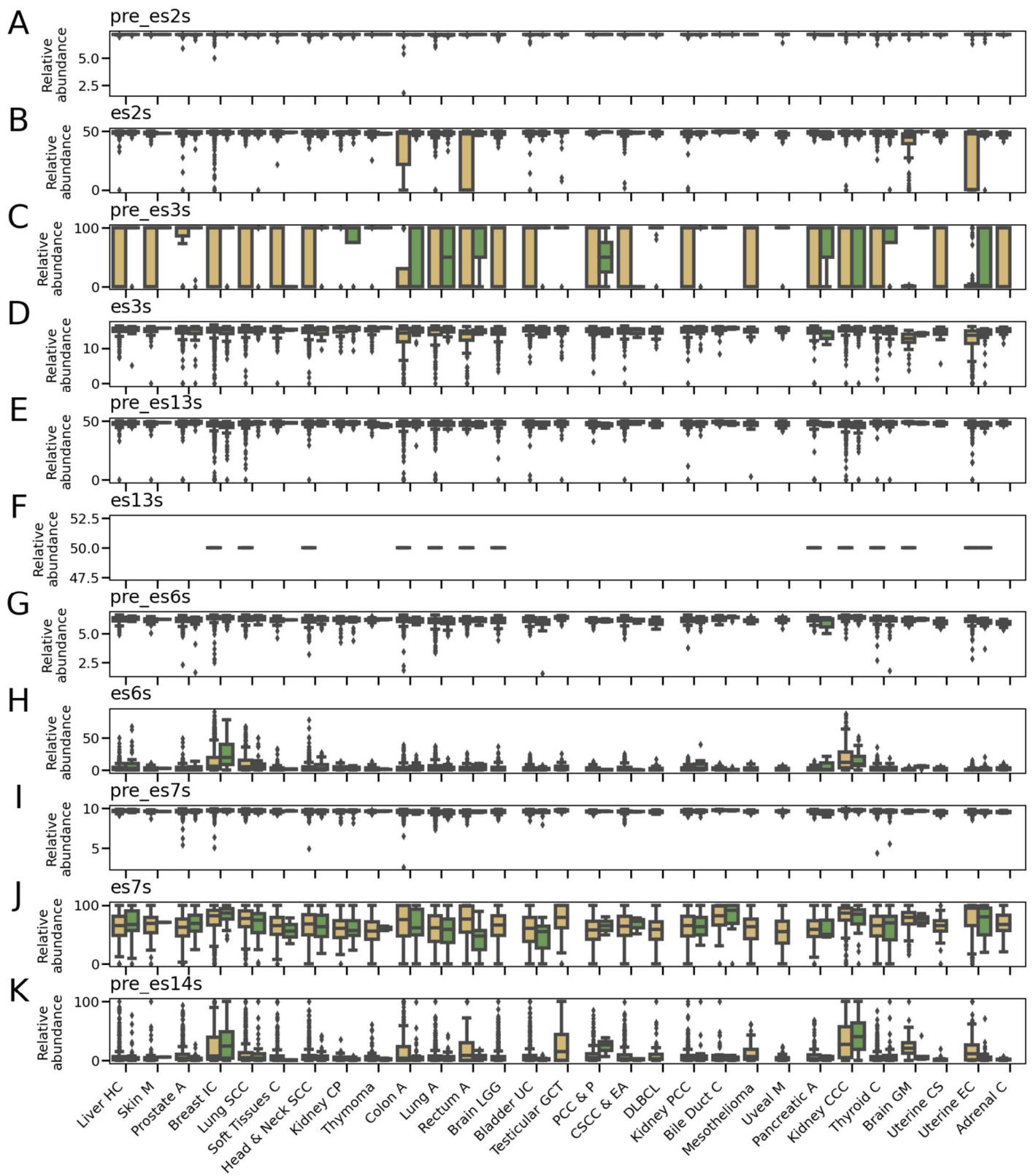
**Figure S26. Cancer-specific rRNA variants relative-abundances (1 of 6 figures), related to Figure 6.** (A-K) Scatter plot of top abundant regional rRNA variants relative abundances for TCGA cancer and control biopsy samples (cancer and control samples are in yellow and green boxes respectively **Table**

**S23** for region to regional variant conversion and the P-value for comparing case/control). The top most abundant rRNA regional variant is presented per ES/non-ES region across tissues. The x-axis is the same for all panels and is displayed in (K). Abbreviations:

Adrenal C = Adrenocortical Carcinoma
Bile Duct C = Cholangiocarcinoma
Bladder UC = Bladder Urothelial Carcinoma
Brain GM = Brain Glioblastoma Multiforme
Brain LGG = Brain Lower Grade Glioma
Breast IC = Breast Invasive Carcinoma
Colon A = Colon Adenocarcinoma
CSCC & EA = Cervical Squamous Cell Carcinoma and Endocervical Adenocarcinoma
DLBCL = Lymphoid Neoplasm Diffuse Large B-cell Lymphoma
Head & Neck SCC = Head and Neck Squamous Cell Carcinoma
Kidney CCC = Kidney Renal Clear Cell Carcinoma
Kidney CP = Kidney Chromophobe
Kidney PCC = Kidney Renal Papillary Cell Carcinoma
Liver HC = Liver Hepatocellular Carcinoma
Lung A = Lung Adenocarcinoma
Lung SCC = Lung Squamous Cell Carcinoma
Pancreatic A = Pancreatic Adenocarcinoma
PCC & P = Pheochromocytoma and Paraganglioma
Prostate A = Prostate Adenocarcinoma
Rectum A = Rectum Adenocarcinoma
Skin M = Skin Cutaneous Melanoma
Soft Tissues C = Soft Tissues Carcinoma
Testicular GCT = Testicular Germ  Cell Tumors
Thyroid C = Thyroid Carcinoma
Uterine CS = Uterine Carcinosarcoma
Uterine EC = Uterine Corpus Endometrial Carcinoma
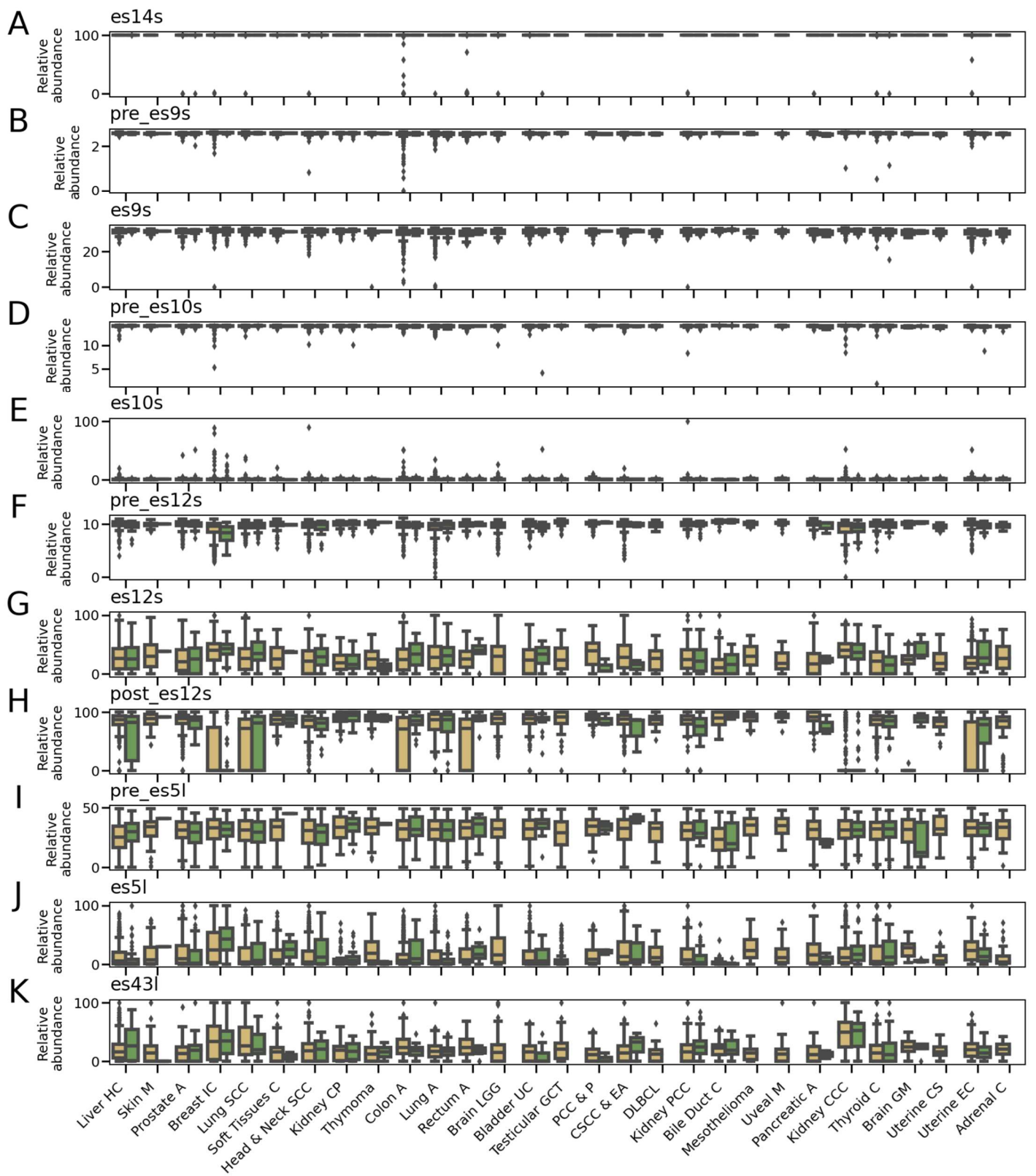Uveal M = Uveal Melanoma

**Figure S27. Cancer-specific rRNA variants relative-abundances (2 of 6 figures), related to Figure 6.** (A-K) Scatter plot of different regional rRNA variants, relative abundances for TCGA cancer biopsy samples (cancer and control samples are in yellow and green boxes respectively). The top most abundant rRNA regional variant is presented per ES/non-ES region across tissues (**Table S23** for

region to regional variant conversion and the P-value for comparing case/control). The x-axis is the same for all panels and is displayed in (K). X-axis cancer type full name for the abbreviations are listed at the bottom of **Figure S26**.

**Figure S28. Cancer-specific rRNA variants relative-abundances (3 of 6 figures), related to Figure 6.** (A-K) Scatter plot of different regional rRNA variants, relative abundances for TCGA cancer biopsy samples (cancer and control samples are in yellow and green boxes respectively). The top most

abundant rRNA regional variant is presented per ES/non-ES region across tissues (**Table S23** for region to regional variant conversion and the P-value for comparing case/control). The x-axis is the same for all panels and is displayed in (K). X-axis cancer type full name for the abbreviations are listed at the bottom of **Figure S26**.
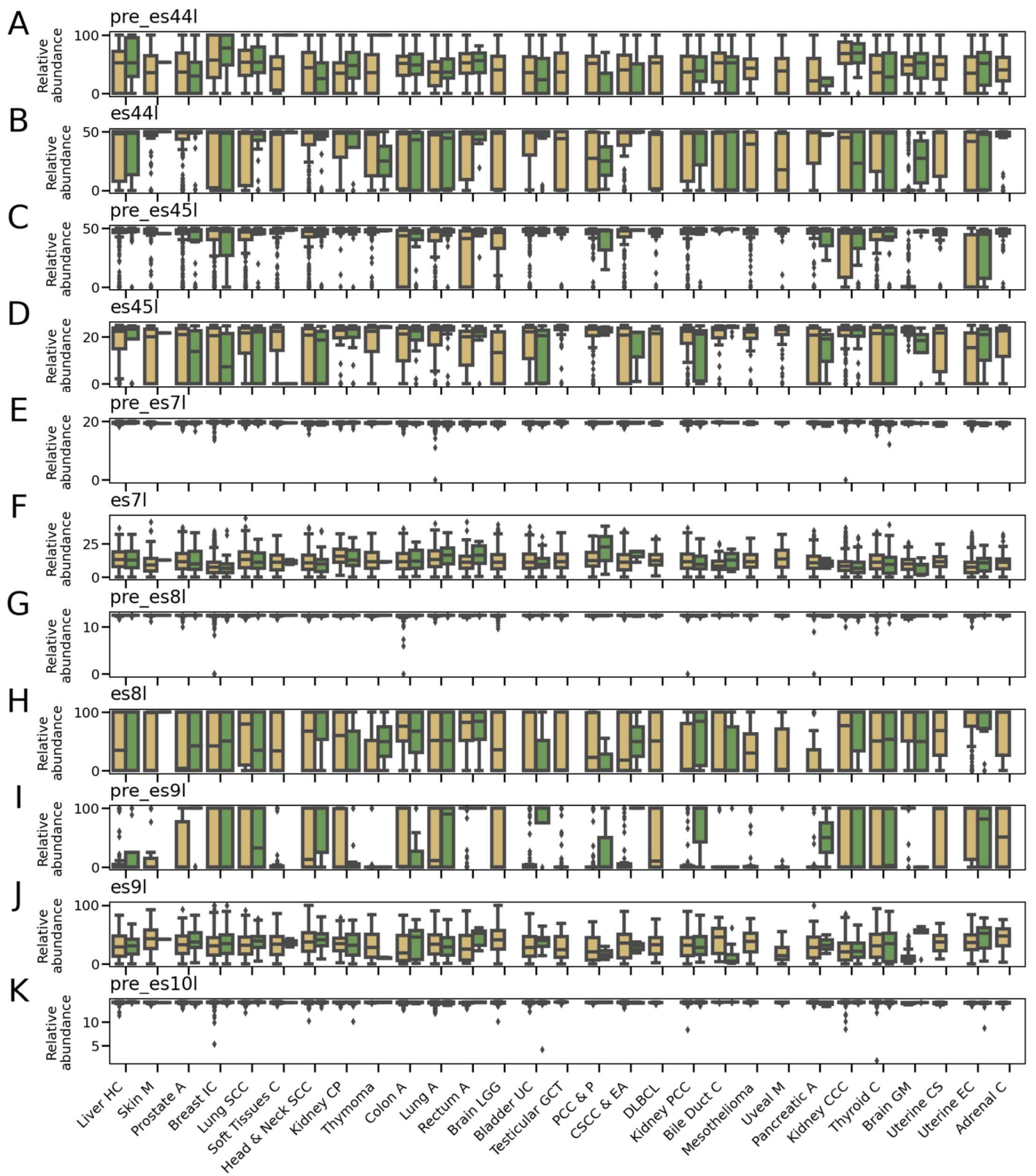
**Figure S29. Cancer-specific rRNA variants relative-abundances (4 of 6 figures), related to Figure 6.** (A-K) Scatter plot of different regional rRNA variants, relative abundances for TCGA cancer biopsy samples (cancer and control samples are in yellow and green boxes respectively). The top most

abundant rRNA regional variant is presented per ES/non-ES region across tissues (**Table S23** for region to regional variant conversion and the P-value for comparing case/control). The x-axis is the same for all panels and is displayed in (K).  X-axis cancer type full name for the abbreviations are listed at the bottom of **Figure S26**.
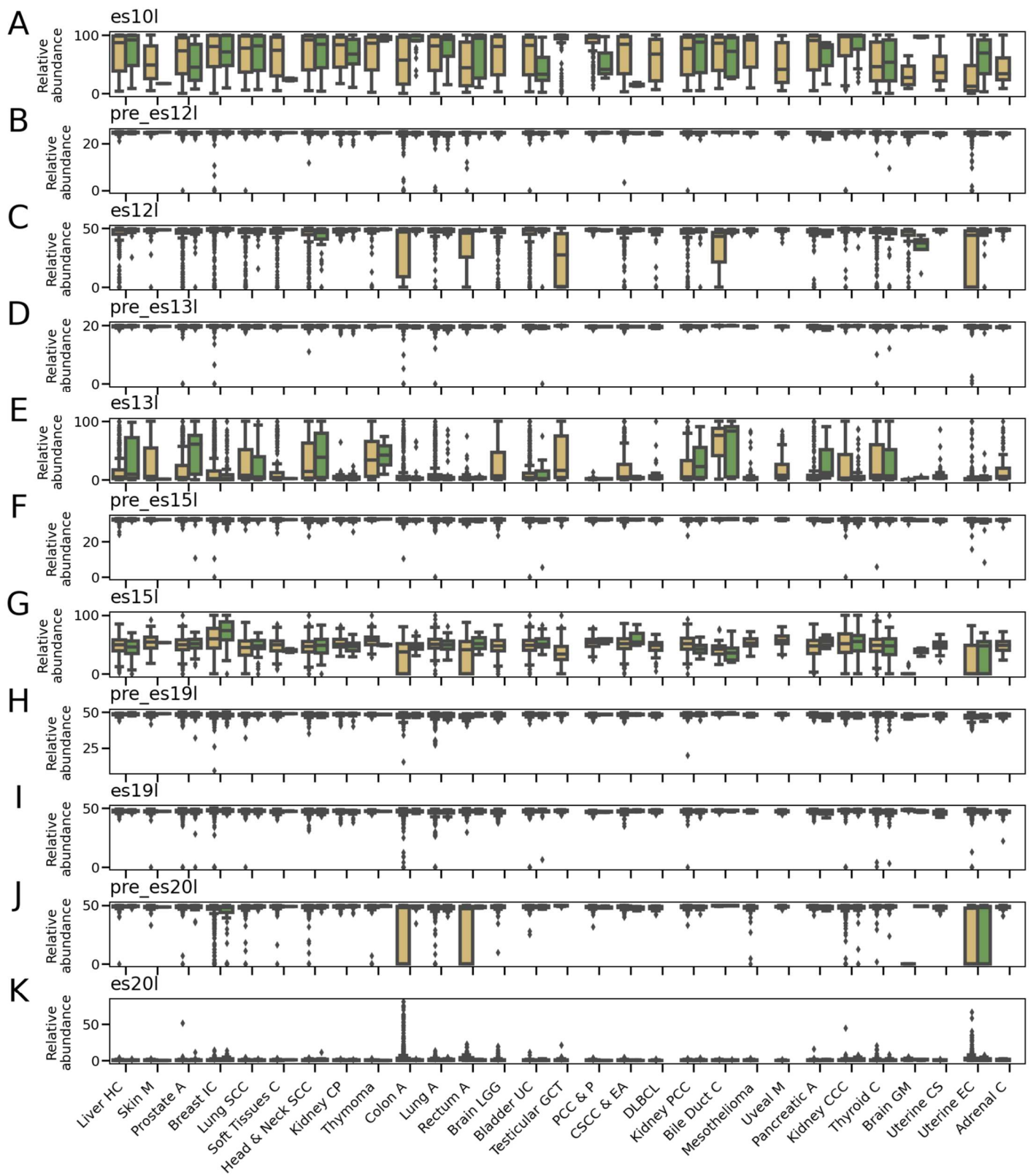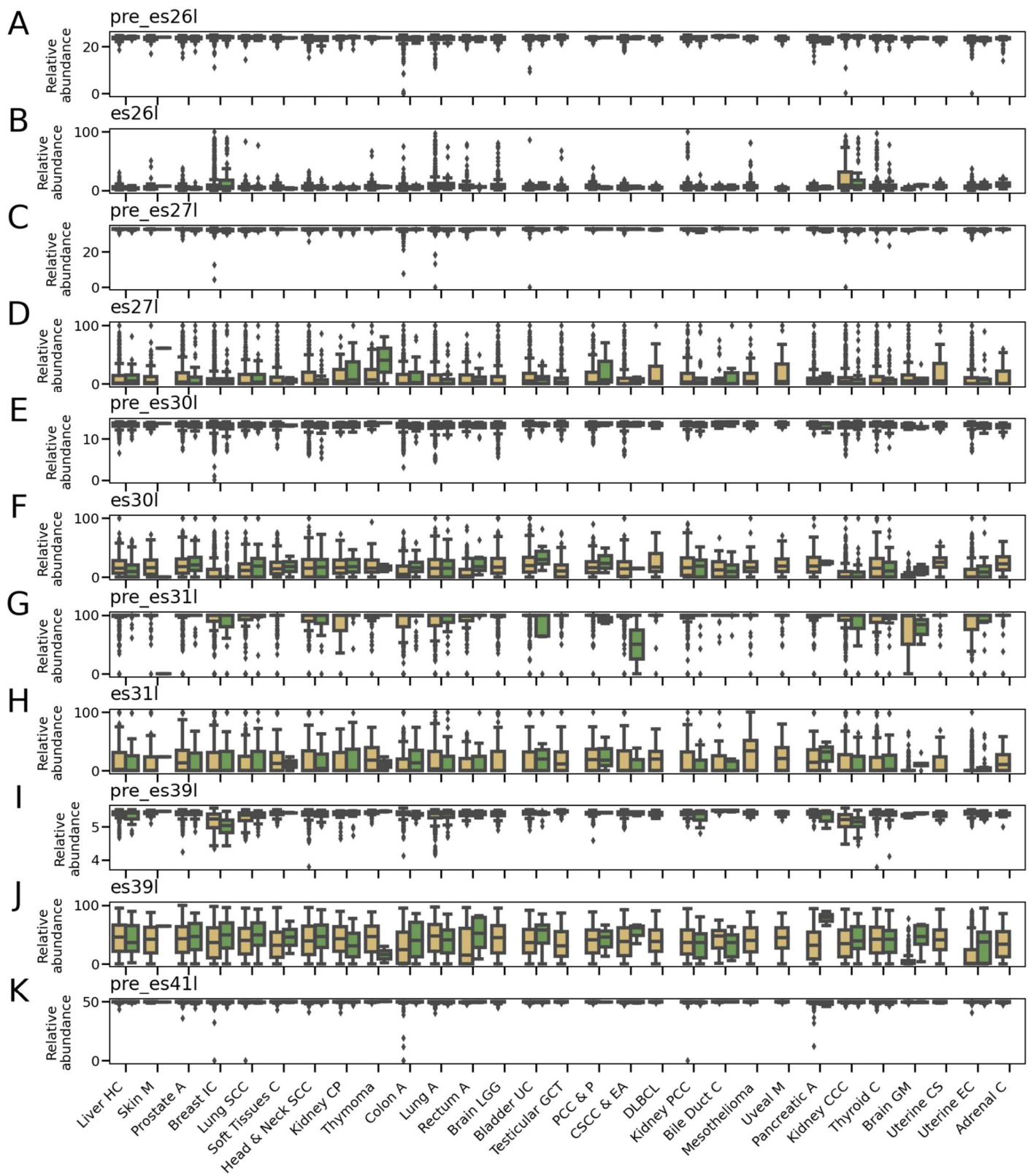
**Figure S30. Cancer-specific rRNA variants relative-abundances (5 of 6 figures), related to Figure 6.** (A-K) Scatter plot of different regional rRNA variants, relative abundances for TCGA cancer biopsy samples (cancer and control samples are in yellow and green boxes respectively). The top most

abundant rRNA regional variant is presented per ES/non-ES region across tissues (**Table S23** for region to regional variant conversion and the P-value for comparing case/control). The x-axis is the same for all panels and is displayed in (K). X-axis cancer type full name for the abbreviations are listed at the bottom of **Figure S26**.
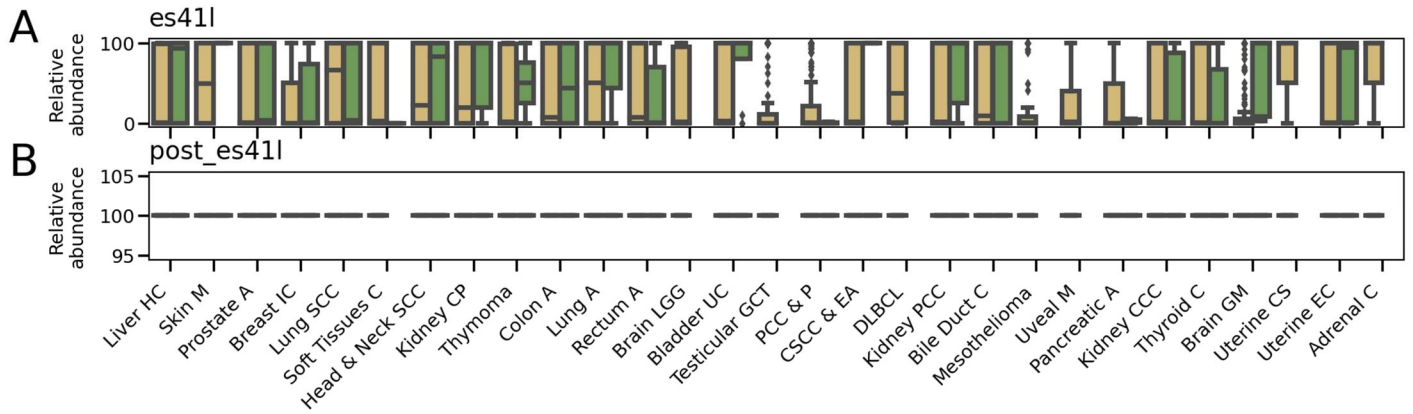


**Figure S31. Cancer-specific rRNA variants relative-abundances (6 of 6 figures), related to Figure 6.** (A-B) Scatter plot of different regional rRNA variants, relative abundances for TCGA cancer biopsy samples (cancer and control samples are in yellow and green boxes respectively). The top most abundant rRNA regional variant is presented per ES/non-ES region across tissues (**Table S23** for region to regional variant conversion and the P-value for comparing case/control). The x-axis is the same for all panels and is displayed in (B).  X-axis cancer type full name for the abbreviations are listed at the bottom of **Figure S26**.