

# Diversity of ribosomes at the level of rRNA variation associated with human health and disease

Daphna Rothschild<sup>1\*</sup>, Teodorus Theo Susanto<sup>1\*</sup>, Xin Sui<sup>3,4\*</sup>, Jeffrey P. Spence<sup>1</sup>, Ramya Rangan<sup>5</sup>, Naomi R. Genuth<sup>1,2</sup>, Nasa Sinnott-Armstrong<sup>1</sup>, Xiao Wang<sup>3,4</sup>, Jonathan K. Pritchard<sup>1,2§</sup>, Maria Barna<sup>1§</sup>

<sup>1</sup> Department of Genetics, Stanford University, Stanford, CA, 94305, USA

<sup>2</sup> Department of Biology, Stanford University, Stanford, CA, 94305, USA

<sup>3</sup> Department of Chemistry, Massachusetts Institute of Technology, Cambridge, MA, 02139, USA.

<sup>4</sup> Broad Institute of MIT and Harvard, Cambridge, MA, 02142, USA.

<sup>5</sup> Biophysics Program, Stanford University, Stanford, CA, 94305, USA

\* Co-first author

§ Co-senior author. Email: [pritch@stanford.edu](mailto:pritch@stanford.edu) (JKP), [mbarna@stanford.edu](mailto:mbarna@stanford.edu) (MB).

## Summary

With hundreds of copies of ribosomal DNA (rDNA) it is unknown whether they possess sequence variations that ultimately form different types of ribosomes. Here, we developed an algorithm for variant-calling between paralog genes (termed RGA) and compared rDNA variations with rRNA variations from long-read sequencing of translating ribosomes (RIBO-RT). Our analyses identified dozens of highly abundant rRNA variants, largely indels, that are incorporated into translationally active ribosomes and assemble into distinct ribosome subtypes encoded on different chromosomes. We developed an *in-situ* rRNA sequencing method (SWITCH-seq) revealing that variants are co-expressed within individual cells and found that they possess different structures. Lastly, we observed tissue-specific rRNA-subtype expression and linked specific rRNA variants to cancer. This study therefore reveals the variation landscape of translating ribosomes within human cells.

# Introduction

The ribosome is a complex, ancient machine responsible for all protein synthesis, with a core ribosomal RNA (rRNA) structure that is conserved across all kingdoms of life. The primary transcript of the rRNA genes is the large 45S pre-rRNA which contains the 18S, 5.8S, and 28S rRNAs as well as transcribed spacer regions. In humans, rRNA genes are present in hundreds of ribosomal DNA (rDNA) copies in tandem repeats that are spread across multiple chromosomal loci (1). These high rDNA copy numbers are thought to be necessary to produce millions of ribosomes in each cell. Nevertheless, these hundreds of rDNA copies allow for sequence variation between copies, as was first noted in mice and humans almost 50 years ago (2).

It remains an outstanding challenge to understand whether ribosomes are different at the level of rRNA and how many ribosome subtypes may exist. For the last several decades we have therefore had a limited knowledge of fundamental differences in the translational machinery and limited insight beyond the textbook view of ribosome composition. By examining short-read sequencing data of the 1,000 Genomes Project (1KGP) (3), previous studies discovered hundreds of positions in rDNA bearing sequence variants (4, 5). A major challenge faced by these studies arises from the depletion of GC-rich sequences in short-read sequencing. Specifically certain regions within rDNA possess >80% GC content on average (5–7). This results in two types of false variant calling: (1) false negatives, caused by the inability to identify variants in regions with poor sequence coverage, and (2) false positives, caused by PCR errors in low-coverage regions. Moreover, these GC rich regions are highly repetitive which make short-read variant discovery tools inaccurate (8). As a result, previous studies (4, 5) reported contradictory results likely because of these major caveats. Despite limitations with short-reads, the copy number of the rDNA loci was shown important for gene expression and cellular homeostasis (9–12). Moreover, rDNA copy number was shown to decrease both with age and in cancer (13–15). This calls for further investigation into whether rDNA copy number changes are coupled with sequence variations and if such changes affect human health.

Recently, long-read sequencing enabled the successful complete assembly of a human genome by the Telomere-2-Telomere (T2T) consortium, including positioning of rDNA copies in the five acrocentric chromosomes (16). Moreover in the mouse and Arabidopsis plant genomes, rDNA variants were grouped into haplotypes and a few rRNA variants were found expressed in tissues

using short-reads (4, 17, 18). Yet in order to find low-frequency variations between full length rDNA paralog copies, new computational method development is necessary. Long-reads offer distinguishing between paralog genes. However, existing common methods for long-read variant calling, such as DeepVariant (19) and Clair (20), are primarily designed for detecting variants in single-copy regions. For paralog genes where low-frequency variants exist between copies accurate variant calling is lacking (8). Moreover, it remains an open question if rRNA haplotypes are expressed from the human genome, what their abundances are, and if such variability is linked to human physiology. There too, long-read sequencing and analysis of full-length rRNA is necessary. This highlights the need for new approaches to comprehensively characterize human ribosome diversity.

To address this need in the field, we devised an efficient novel computational algorithm to detect all variations between paralogs (termed RGA). Applied to the long-read 1,000 Genomes Project dataset, we discovered hundreds of rDNA sequence variations enriched with previously undiscovered insertion-deletions. We further developed a novel methodology to perform long-read sequencing on rRNA in actively translating ribosomes to identify variants (RIBO-RT). Using this method, we discovered that ribosomes have different subtypes with rRNA variants that are genomically encoded by rDNA clustered on distinct chromosomes. Additionally, using an *in-situ* rRNA sequencing platform that we developed (SWITCH-seq), we discovered that variants belonging to different rRNA subtypes are co-expressed in single cells. We then used structure probing coupled with long-read sequencing to find that 28S subtypes have different rRNA structures. Lastly, we found that these subtypes are differentially expressed in human tissues, and low abundance variants are elevated in certain cancers. Together, these results suggest that ribosomes with unique sequence variation may be used to modulate different cellular programs underlying human physiology and disease.

## Results

### Indels are the main variants of the human rDNA loci

How rDNA variation shapes the presence of unique ribosomes in the cell remains an important open question. Previous studies that analyzed the 1KGP dataset for discovery of rDNA variants reported discordant results. Parks et al. (4) reported hundreds of variants in both the 18S and the 28S, yet 75% of variants were not made publicly available, making a comparison to this dataset problematic.

Nonetheless, Fan et al. (5) reported notable differences from Parks et al. by suggesting that the 18S has low variation, and also reporting many fewer variant positions in the 28S. Moreover, Parks et al. reported only 2.7% of variants being indels while Fan et al. reported 19.2% indels. Here, considering the limitations of short-reads in rDNA variant discovery and their inability to distinguish between rDNA paralogues we decided to re-evaluate the variants in the human rDNA genes.

Until recently, the 1KGP dataset included only short-read genome sequencing (21). Yet as of 2022, the 1KGP includes PacBio's HiFi long-read sequencing for 30 individuals from diverse ancestral origins, which could serve as a better method for accurately calling variants. Here, we compared the rDNA variants captured by short- and long-reads from the same individuals and addressed the discrepancies between previous studies (**Fig. 1A,B, Table S1**).

When analyzing the short-read data, we followed the pre-processing steps as performed in previous studies of marking duplicate reads which are suspected as PCR artifacts. This step discarded 97% of reads. However, it is unknown if duplicate reads are PCR biases given the high rDNA paralog copy numbers which highlights the limitation of short-read sequencing for rDNA variant discovery. Next, in order to call variants including rare variants which are not expected to follow germline variant frequencies in high paralogues rDNA copy numbers, we tested two common somatic variant calling methods for short-reads: LoFreq\* (22) which was used by Parks et al. and Mutect2 (23). Specifically, Mutect2 was chosen instead of the germline variant caller used by Fan et al., because germline variant callers will not detect rare variants found between paralogues. Using the LoFreq\* method which is known to be sensitive (22), we found 1,582 positions with variants compared to 861 positions with variants with Mutect2 (**Table S2-3**). Notably, both methods detected 23% indels, which is on par with the indel percentage reported by Fan et al.. Given the difference in the proportion of indel frequencies between the previous two studies, we compared variant quality scores of SNVs and indels (**Table S3** - Mutect2 FDR-corrected log 10 likelihood ratio score of variant existence). Here we found that indel variants were enriched with high confidence P-values (**Fig. S1**, P-value <  $10^{-15}$  comparing SNVs and indels likelihood ratio scores using Kolmogorov-Smirnov test for goodness of fit). Notably, tandem repeats and GC-rich sequences in the human genome were shown to be prone to chromosomal breakage and were found enriched in indels (24). Therefore, the 23% indel frequency derived solely from short-read data may underrepresent the true indel frequency in these samples. To test this, we next examined the HiFi long-read data from the same samples.

LoFreq\* and Mutect2 did not work on the long-read sequencing data. In order to identify all positions with sequence variants, we developed a new computational method for accurate variant calling between paralogs which we term Reference Gap Alignment (RGA, see **Methods**). We align all reads against a common reference, and report all variants at a given position with respect to this reference (**Fig. 1C, S2, Methods**). Our method reports at every position with respect to the 18S/28S reference if that position has a variant and calls its identity. The only parameters in our method are the standard pairwise alignment parameters: a mismatch penalty, a gap opening penalty, and a gap extension penalty (see **Methods** for more details). When benchmarking the global sequence alignment parameters, these resulted in similar indel proportions (**Table S4**). In agreement with previous studies (4, 5), we found that rDNA is highly variable, yet using our method, we discovered that the vast majority of variants are short indels and not SNVs in all 30 samples. Specifically, when examining each reference position individually we found that on average 95% of variants are GC-rich indels (**Table S5**).

Since this indel proportion found with long-reads is markedly different from the results obtained with short-reads (**Fig. 1b**), we decided to compare our results from our primary long-read sequencing technology, namely HiFi, with an alternative, Oxford Nanopore (ONT). Importantly, ONT is much more error prone compared to HiFi, with an estimated 13% error rate in ONT compared to 0.1% error rate in HiFi (25–27). Since 1KGP long-read sequencing was only performed on HiFi, we tested this using the Genome In A Bottle (GIAB) dataset, which consists of two trio families, where both HiFi and ONT were performed on the same samples. Here, when examining the HiFi dataset of GIAB, in agreement with the 1KGP HiFi results, we discovered that 96% of variants are indels (**Table S6, Methods**). When using the ONT dataset as a validation dataset, 81% of variants found in HiFi were replicated in the ONT dataset (**Table S7, Methods**). Notably, the variants that were not identified by ONT consisted of insertions and SNVs but not deletions (**Table S8**). Additionally, after retaining variants found at frequencies above the ONT error rate, 87% of found variants were indels (**Table S7** filtering variants with allele frequency smaller than 0.13).

Our results highlight that previous studies which used short-read sequencing missed the majority of the variants found in the rDNA loci. Most of our identified rDNA variants were found at GC-rich regions which are depleted from short-read sequencing. We conclude that indels are the main

variants in rDNA both between and within individuals. These findings also highlight the need to curate a reference of rDNA variants and the importance of our accurate variant calling method together with long-read sequencing in confidently assigning rDNA variants.

## An atlas of 18S and 28S human rDNA variants validated in rRNA of translating ribosomes and single cell microscopy

It is unknown if the rDNA copies with sequence variants found in the human genome are transcribed and are found in functionally translating ribosomes. With no human reference of different rRNA subtypes, studies performing RNA sequencing completely ignore rRNA variants which limits our understanding of the contribution of rRNA to human physiology and disease.

Sequencing of rRNA has been historically technically challenging (5, 28). Here, in addition to our computational RGA variant discovery method, by optimizing long-read sequencing, we have successfully developed an experimental method for full length rRNA sequencing of 18S and 28S from translating ribosomes which we named RIBO-RT. In order to extract rRNA from translationally active ribosomes, we first employed sucrose gradient fractionation wherein ribosomes can be separated into free ribosomal subunits and translationally active ribosomes, which contain one or more ribosomes bound to the same mRNA. We extracted RNA from translating ribosomes-containing fractions (**Fig. S3**), performed reverse transcription in denaturing conditions, and then sequenced complete 18S and 28S rRNA by HiFi long-read sequencing (**Fig. 1D, Methods**). We selected a human embryonic stem cell line (H7-hESC) and using long-read sequencing, we sequenced the 18S and 28S from both its rDNA and rRNA (**Fig. 1D**). Here, we obtained 58,495 sequences of the 18S and 14,430 sequences of the 28S rRNAs from translating ribosomes (**Methods**). With this approach and our variant discovery method, we were able to coherently characterize the 18S and 28S H7-hESC rRNA variants. Most importantly, since rRNA is known to be heavily modified (29), our strategy of matching rRNA to genomic rDNA from the same cell enables us to distinguish modifications or sequencing errors from true sequence variants belonging to different rDNA alleles.

In agreement with our 1KGP rDNA results, we found that the H7-hESC rDNA is highly variable and is enriched with indels. Additionally, rDNA variants are transcribed into functionally translating ribosomes as they are present in polysome fractions (**Fig. S4**). Specifically, we found 270 positions

with variants in the 18S and 858 positions with variants in the 28S, corresponding to about one variant every six rRNA positions (**Fig. S4-5**). When comparing these variants to the variants found in the 1KGP genome dataset, high quality score indels found using the short-read data (with LoFreq\* and Mutect2) were also detected using our method. Most variants (59%) are found in expansion segment (ES) regions (**Fig. S6**, ES/non-ES regions are annotated). These regions vary in sequence both within and among different species, nearly doubling the eukaryotic rRNA sequence relative to that of prokaryotes (30). ESs have recently been shown to bind ribosome associated proteins and transcripts, yet their functions remain poorly understood (30–34).

Most importantly, with accurate variant calling and full coverage of the underlying rDNA and transcribed rRNA, we can measure the frequency of each variant between the rDNA copies and rRNA expression levels. We distinguish possible modifications or sequencing errors from certain sequence variants belonging to different rDNA alleles by calling variants with similar frequencies measured in rDNA and rRNA as high-confidence alternate allele variants while those significantly deviating between their rDNA and rRNA frequencies as low-confidence alternate alleles (**Fig. 1E**, **Table S9**, Fisher's Exact Test measuring association between rDNA and rRNA alternate allele frequencies - high-confidence alternate alleles are marked in red, low-confidence alternate alleles are marked in orange). Surprisingly, we did not identify any abundant high-confidence variants within the 18S. This suggests that rRNA variation is not tolerated within the small ribosome subunit. While most variants have low abundance, in the 28S we found 23 variants at 17 positions with frequencies above 10% in both rRNA and rDNA (**Fig. 1E,F**, **Table S10**, **Fig. S5** annotated positions, **S5C** minor allele frequency with dashed-like marking at 10% frequency). Notably, for the 28S high abundance variants, there is very good agreement between variant frequencies in rRNA and rDNA (**Fig. 1E**,  $R=0.93$  Pearson correlation, correlating all variants colored in red and orange).

Next, we focused on the 28S high abundance variants. For most positions, the RNA45S5 reference allele is the major allele found in the sequenced H7-hESC line (**Fig. 1F**, reference allele in gray, **Extended Data 1-2 nucleotide atlas**). Yet some alternate alleles were more abundant than the 28S RNA45S5 reference alleles (**Fig. 1F**, gray for the reference allele and red for variants). Notably, high abundant variants are only located in 4 ES regions (es7I, es15I, es27I and es39I) and one non-ES region, helix 28S:h11 (**Fig. 1F**). Moreover, in the es7I, es15I and es27I regions we observed that variants can be grouped and characterized by indels of GGC in tandem-repeats. GGX

tandem-repeats were recently suggested capable of forming G-quadruplex structures (35) while other works suggested that such repeats can form other higher order structures (36, 37). While the function of these ESs is largely unknown, a growing body of research supports various roles in translation regulation. For example, es27I has been shown to be important for control of translation fidelity and binding of ribosome-associated proteins for several processes, such as initiator methionine cleavage from the nascent polypeptide (38) or acetylation of nascent polypeptides (39). Moreover, es39I interacts with the Signal Recognition Particle (SRP), which identifies the signal sequence on nascent polypeptides emerging from the translating ribosome (40). Interestingly, the most abundant variant in the non-ES helices, a G to A substitution at position 60 in 28S:h11, is considered unique to humans. The alternative allele, A, is the reference allele for other mammals including chimpanzees (41, 42).

For these aforementioned highly abundant variants, a strong correlation between the frequency of a variant's occurrence among rDNA copies and its expression levels in rRNA indicates both the authenticity of these sequence variants and their likely co-expression within individual cells. To explore this hypothesis, we developed a template-switching-based *in-situ* sequencing method (SWITCH-seq) to visualize variant ribosomes in individual HeLa cells (**Fig. 2A**). This approach involved designing a reverse transcription primer to target constant non-variable regions downstream of the selected rRNA variant regions. Specifically, we selected regions for which we could design a primer for each of the ES regions with abundant variants (**Methods, Table S11**). The process of SWITCH-seq begins with performing reverse transcription on fixed HeLa cells, wherein a known sequence of choice is attached to the 3' end of cDNA (template switching) (**Fig. 2A, Methods**). This step incorporates the variant of interest into the cDNA, which is subsequently amplified into *in-situ* cDNA amplicons through enzymatic circularization and rolling circle amplification. These amplicons encapsulating the rRNA variants are then probed into a hydrogel network for cyclic imaging using a confocal microscope (**Methods**). We conducted multiple rounds of imaging which capture two bases before the variants, and the consecutive base, where we expect to find the variants (**Fig. 2B, Fig. S7**). As predicted, we successfully observed both the reference and alternate alleles (**Fig. 2B,C, Fig. S7**). Furthermore, the frequencies of the reference and variant alleles corresponded with their frequencies in the H7-hESC samples (**Table S12**). We conclude that rRNA variants observed at high frequency in the H7-hESC rRNA and rDNA and in the rDNA across the 1KGP samples form ribosomes that are co-expressed in individual cells that can be visualized at single cell resolution.



As an important resource for studying human rRNA variations, we create the first comprehensive atlas of all h7-ESC rRNA 18S and 28S rRNA variants at different resolutions from nucleotide variants to gene 28S haplotype-groups that we later describe as separate subtypes (**Fig. S6, Table S13-14** for region annotation).

## 28S variants assemble to genomically distinct ribosome subtypes

Since we successfully obtained full length 18S and 28S rRNA with variations, we address the outstanding question of whether variations lead to the formation of different ribosome subtypes. Here we focused on the 28S since 18S variants appeared at low frequency. For the 28S, we found high agreement between rRNA and rDNA variant frequency so we first asked which rDNA variants are co-located on the same 28S rDNA copy. To do this, we calculated the correlation coefficient (Pearson's  $r^2$ ) between positions across all 28S H7-hESC rDNA sequences. This is analogous to measuring the linkage disequilibrium (LD) coefficient in population genetics, though across paralogous copies within a single genome rather than across individuals in a population. Notably, we found low global LD structure between highly abundant rDNA variants (**Fig. 3A** showing LD for rDNA positions with found rRNA frequency > 10%), perhaps indicating high rates of nonallelic gene conversion across rDNA copies. The highest LD ( $r^2 > 0.2$ ) was found between the es27I to all other regions. Comparing different regions, we found LD between four regions: 28S:h11, es15I, es27I and es39I, where in each region we identified a position with higher linkage to the other three regions (**Fig. 3A**, four positions are annotated, position 60 being 28S:h11, with higher linkage to other positions). By considering the variants at this subset of positions, we found a total of 21 different haplotypes both in rDNA and rRNA (**Fig. S8**). For testing if haplotypes can be considered as different 28S subtype variants, we further analyzed two independent long-read DNA datasets: (1) the fully assembled genome from the T2T with 219 rDNA copies with their chromosome location (16), and (2) the GIAB HiFi dataset (43). In agreement with the H7-hESC results, 3 out of the 4 positions with higher linkage to other positions in the H7-hESC had higher LD in the GIAB dataset (**Fig. 3B**). Since these three variants are linked to variants at other positions, we define the haplotypes formed by positions 60, 3513 and 4913, belonging to regions 28S:h11, es27I and es39I respectively, as different 28S haplotypes (**Fig. 3C, Extended Data 7 28S atlas**).

We next asked whether different 28S haplotypes are spatially separated in the genome as different subtypes. For the H7-hESC we have 386 complete 28S rDNA sequences and in the GIAB dataset we randomly subsampled each GIAB sample to 386 complete 28S rDNA sequences. For these datasets, we do not know rDNA-chromosome positioning. Notably, by comparison of 28S rDNA sequence similarities, we detected distinct 28S sequence groups in both hESC and GIAB (**Fig. 3D**, Principal Coordinate Analysis, PCoA, of Bray-Curtis dissimilarities between 28S sequences (44), **Methods**). Here, the different clusters in PCoA space match different 28S subtypes. Specifically, we observed that 28S sequences of a given subtype are more similar to one another in their entire sequence compared to 28S rRNAs of other subtypes (**Fig. 3D**). In the T2T assembly, rDNA copies have chromosome coordinates, which enables us to measure 28S subtype presence at the five acrocentric chromosomes. Remarkably, we discovered that 28S haplotypes are largely chromosome specific (**Fig. 3E**). Taken together, this supports that 28S haplotypes belong to different subtypes which are encoded by rDNAs that are genomically separated and clustered on distinct chromosomes.

## Ribosomes of different 28S subtypes have different structures.

We next asked if different 28S subtypes have different ribosome structures. Notably, the abundant variants in the hESC were found in four different ES regions which were never previously resolved by CryoEM. Here, we treated the hESC sample with DMS, which covalently modifies the RNA at regions where the rRNA is accessible to allow for structure probing of the RNA (**Methods**). Using our RIBO-RT method for sequencing full length 28S with our RGA variant calling on DMS treated hESC cells, we obtained an accessibility map of the 28S. Importantly, we are able to predict the structure of the full length ES regions which was not previously possible.

We compared the two most abundant 28S subtypes and their linked variants and found they have different structures (**Fig. 4A-C**, with 22% and 30% frequency of subtypes 1 and 2 respectively, **Methods**). While our method with DMS results in a full length accessibility map of the 28S, secondary structure prediction becomes less accurate for long RNA sequences. Given that ESs have tentacle-like extensions that protrude from the ribosome, we assumed that the core non-ES rRNA is not affected by changes in the ES regions which allowed us to focus on the structures of individual ESs. Most interestingly, we discovered that the ESs es7I, es15I and es27I have major DMS accessibility and structure differences when comparing the two subtypes observed at the GGC sites

in es7l, es15l and es27l (**Fig. 4A**, ES region box annotations, **Fig. S11-15** for es7l, es15l and es27l). When focusing on es27l, the second longest ES, we noticed that the largest accessibility difference between the subtypes was at the site where es27l subtypes differ, at the GGC indel. Specifically, the subtype with one fewer tandem-repeat GGC insertion before the AG at position 3513 of the 28S showed greater DMS accessibility at position 3513 and its vicinity (**Fig. 4D,E**). This GGC expands a six tandem-repeat GGCs, i.e. (GGC)<sub>6</sub>, which changes the region's structure. Moreover, for the es27l region we found local structure changes near the sequence variants which opens the possibility that there are proteins or transcripts that may interact with the subtype with the GGCAG variant but not with the AG variant (**Fig. 4D,E** region marked in red). Taken together, our DMS results provide evidence of structural differences for different ribosome subtypes.

## Quantifying the relative abundance of rRNA variants in expression data.

With this rRNA atlas, we next aimed to investigate how rRNA composition changes in different human conditions. We used the publicly available Genotype-Tissue Expression (GTEx) short-read RNA-seq dataset to test if rRNA variant frequencies are associated with human tissue biology (**Methods** for atlas usage instructions). Previous studies comparing mRNA across tissues in the GTEx dataset found tissue-specific, including brain-specific, gene expression (45, 46). Here we analyzed 2,618 samples from 332 individuals and 44 tissues from GTEx and asked if rRNA subtypes differ in their expression in these tissues (**Fig. 5A**). We hypothesized that the rRNA subtypes that we identified as highly expressed in the hESC may be important for tissue development. Strikingly, the most abundant subtypes in GTEx significantly differed in expression in many tissues (**Fig. 5B-D, upper panels, Fig. S16-20**, and **Table S14**: P-value < 0.05 FDR corrected Mann-Whitney U rank sum test). Notably, when comparing subtypes expression levels across tissues, we observed significant differences between tissues derived from the ectoderm and endoderm germ layers (**Fig. 5B-D, lower panels, Table S15** FDR corrected ranksums test comparing subtype relative abundances of ectoderm-derived tissues in blue and endoderm-derived tissues in red). Most of the ectoderm-derived tissues belong to brain tissues, and most endoderm derived tissues are digestive-system tissues (**Fig. 5A** endoderm and ectoderm derived tissues are labeled). Taken together, our results support major changes in the expression of rRNA subtypes across tissues.

Lastly, we asked if changes in the expression of rRNA variants are associated with cancer. For this, we used 10,030 samples of short-read RNA-seq with clinical phenotypes from The Cancer Genome

Atlas (TCGA) (47). When comparing cancer types, we found distinct expression patterns of rRNA regional variants across cancers (**Fig. S21-26, Table S16** for region annotations). To test if rRNA variants are cancer-specific, we compared cancer biopsies to control biopsies from the same tissues. Surprisingly, we identified specific rRNA regional variants with significantly different expression levels in control and in cancer biopsies for 11 cancer types (**Fig. 6, Table S17** for alternative allele regional variant abundances, **Table S18**, P-value < 0.05 after FDR correction, Mann-Whitney U rank sum test). These include rRNA variants that while they appeared in low abundance in both the H7-hESC and control biopsies, are found elevated in cancer biopsies. Thus even low abundance variants hold immense importance as disease biomarkers.

We conclude that our atlas enables direct measuring of rRNA variants changes in expression data. Moreover, we showed that atlas variants are present in translating ribosomes and that they are differentially expressed across tissues and cancer types.

## Discussion

Here, by developing a pipeline for long-read sequencing and analysis of rDNA and rRNA from actively translating ribosomes, we measured for the first time variant frequencies in rDNA and rRNA and used *in-situ* sequencing microscopy to validate co-variant expression in individual cells. With this atlas we have enabled greater understanding of the often neglected yet ubiquitous rRNA sequencing data and built an atlas of functional human 18S and 28S rRNA variants at different resolutions, from nucleotide position variants, to 28S gene-level subtypes as a useful resource for studying rRNA variations, and composition across biological conditions (**Extended data 1-5**).

In our study, we have discovered chromosome-associated rDNA subtypes, revealing that different ribosome subtypes based on rRNA sequence variation exist. It may be possible that spatial separation of rDNA subtypes enable regulation of their expression at the chromosome level through allelic inactivation of rDNA loci or inactivation of nucleolar organizer regions (NORs) in the distal junction (48–50). This may enable global remodeling of rDNA transcription and promote specific ribosome subtypes to be expressed within individual cells. Additionally, using DMS structure probing of full length 28S, we discovered that different rRNA subtypes have different structures at ES regions including different DMS accessibility profiles. Since these ES regions are solvent-exposed and highly

flexible, these ES variations may fine-tune regulation of mRNA translation based on differential association with ribosome-associated proteins, mRNA transcripts, or other factors. Moreover, by analyzing the GTEx dataset, we observed differential expression of rRNA subtypes between tissues belonging to ectoderm and endoderm lineages. This pattern might hint at specialized functions of different ribosome subtypes. Long-lived cells associated with the nervous system might require ribosome subtypes that emphasize translation fidelity over speed, as compared to rapidly dividing cells in the digestive tract that require constant replacement given harsh local environments. Indeed, our lab and others have previously shown that *es271* plays a role in translation fidelity through association with ribosome-associated proteins (38, 51, 52). Such interactions that trade speed over fidelity might be fine-tuned by the expression of different rRNA subtypes.

Finally, by analyzing the TCGA dataset we discovered that some low abundant rRNA variants in control biopsies were elevated in cancer biopsies. Future work is needed to understand whether they promote oncogenic ribosome activity and how they are regulated. Therefore, our results provide another layer of ribosome specificity wherein cancer cells might deploy a particular rRNA variant that is more compatible with their cellular fitness. Importantly, we found that specific rRNA variants may be used as biomarkers for disease. Notably, 5-fluorouracil, a common chemotherapy drug, was recently shown to incorporate into rRNA and promote drug resistance by changing mRNA translation (53). It may be that drugs directly target specific rRNA variants and further examination would be needed to test whether they should be used for cancer specific therapies. Together, our results reveal the presence of structurally different ribosomes at the level of rRNA and provide the first atlas to distinguish different types of ribosomes and to link them to different cellular programs, including those underlying human health and disease.

## Acknowledgments

We thank Rhiju Das for interpreting DMS results. We thank Xiangling Meng, Craig Kerr, Ali Wilkening and Michael Montgomery for help with early experiments that were not later followed in this study. We thank the Barna and Pritchard group members for discussions.

## Funding

MB is supported by the New York Stem Cell Foundation and the National Institutes of Health grant R01HD086634. JKP is supported by RO1 HG008140. DR is supported by ALTF 1042-2019 EMBO and LT000218/2020-L HFSP postdoctoral scholarships. TTS is supported by a National Science Scholarship (PhD) from the Agency for Science, Technology and Research. MB is a NYSCF Robertson Investigator. XW is supported by Edward Scolnick Professorship, Ono Pharma Breakthrough Science Initiative Award, Merkin Institute Fellowship, and NIH DP2 New Innovator Award 1DP2GM146245-01.

## Authors contributions

DR conceived the project, designed experimental and computational analyzes, conducted all computational analyses, interpreted the results, and wrote the manuscript. TTS designed and conducted all polysome and DMS sequencing experiments, interpreted the results, and wrote the manuscript. XS developed, optimized, and performed SWITCH-seq experiments. XW supervised the in-situ sequencing experiment. JPS helped in computational analyses. NG helped in experimental data collection. NSA helped interpret GTEX data. RR helped in DMS analyses. MB and JKP conceived and directed the project and analyses, designed the analyses, interpreted the results and wrote the manuscript.

## Competing interests

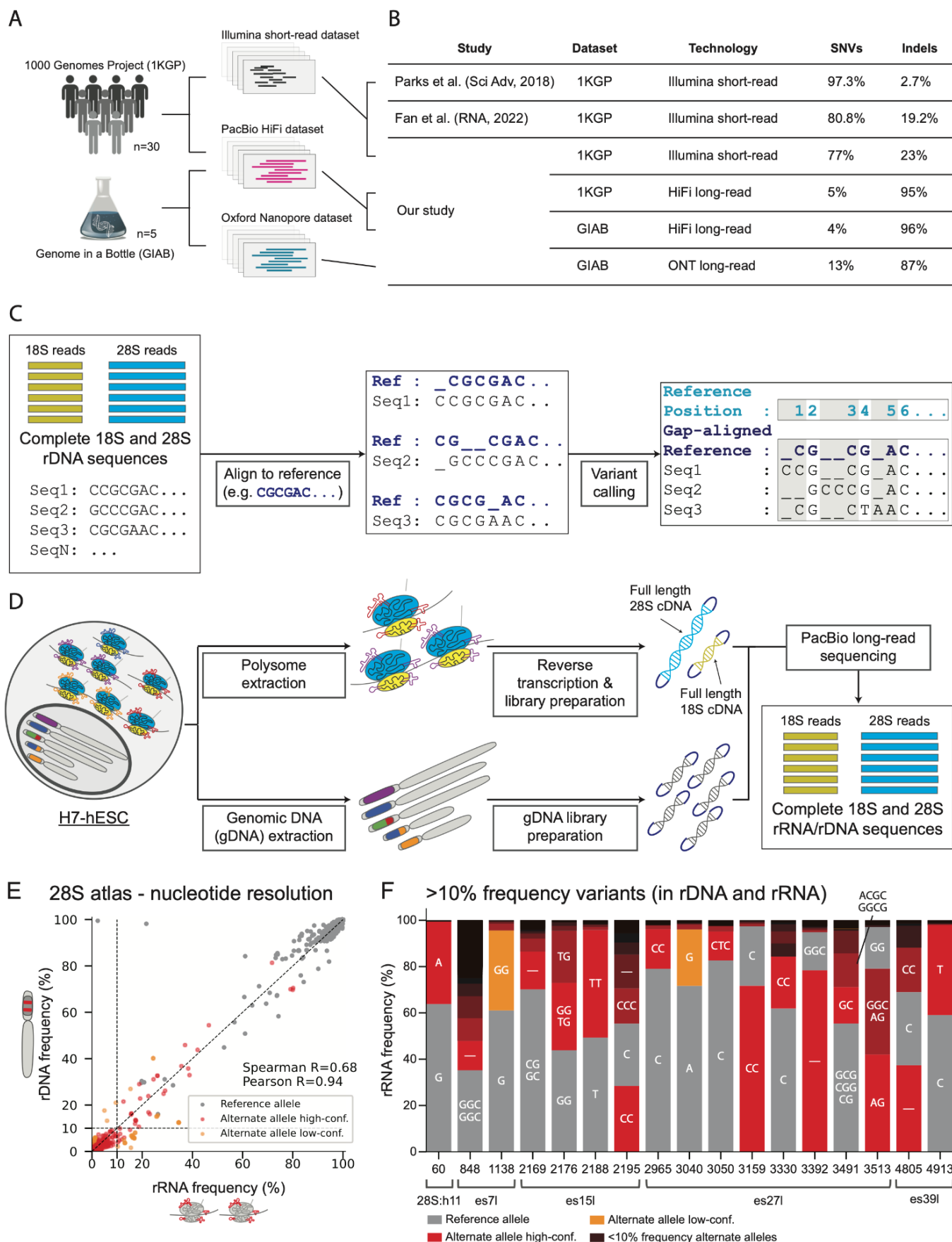
XW is a scientific cofounder of Stellaromics.

## Data and materials availability

The atlas is available as Extended Data to this publication.

H7-hESC raw rDNA and rRNA sequencing data is available under BioProject ID PRJNA926787.

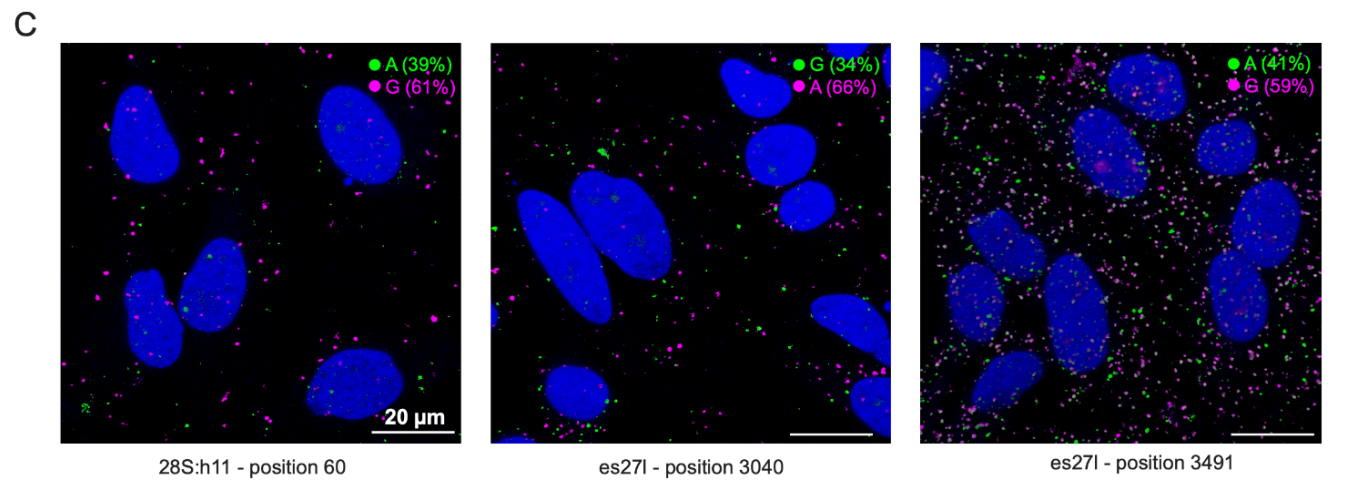
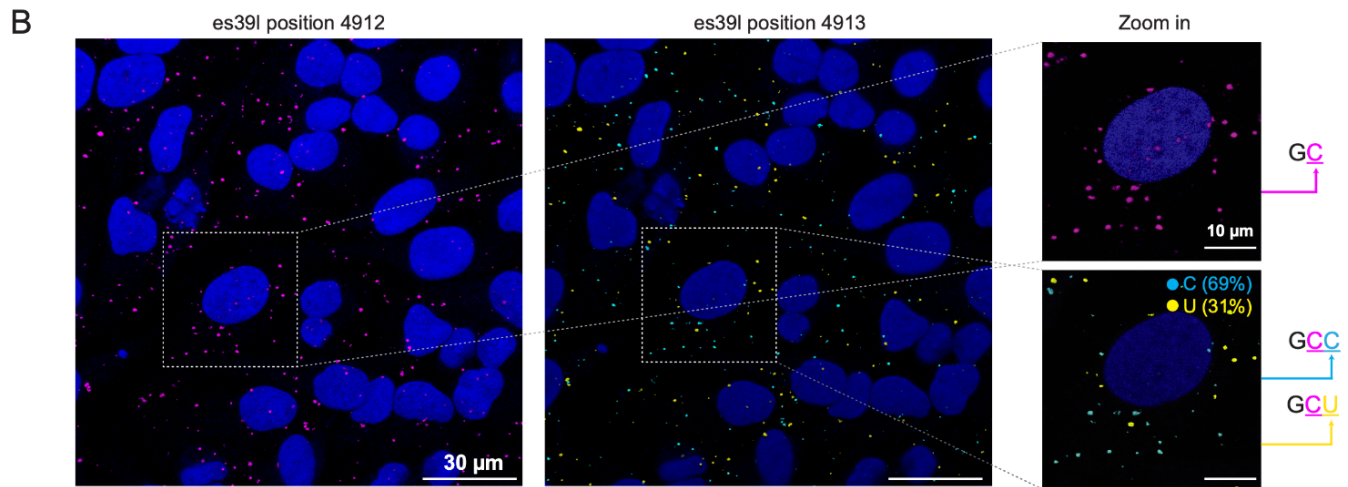
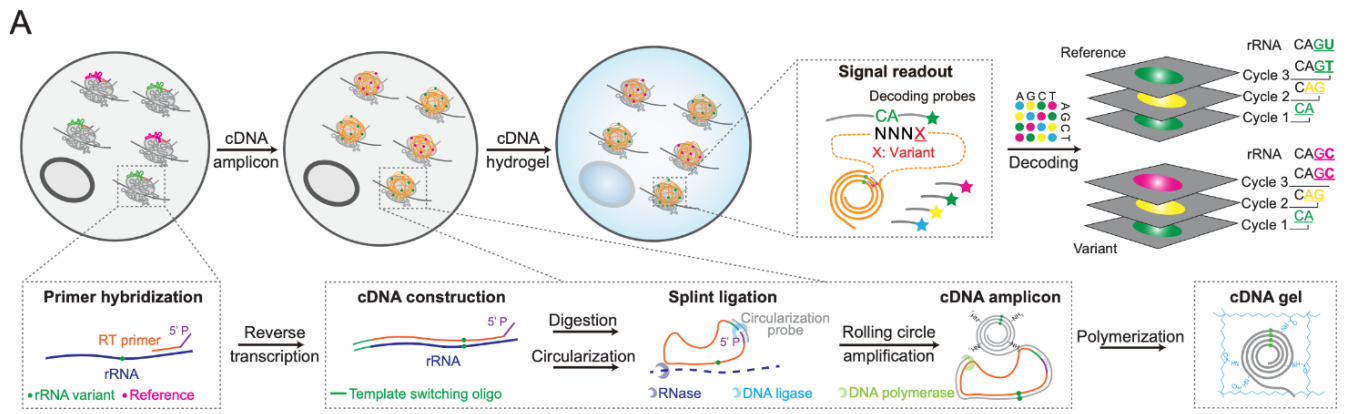
# Figures





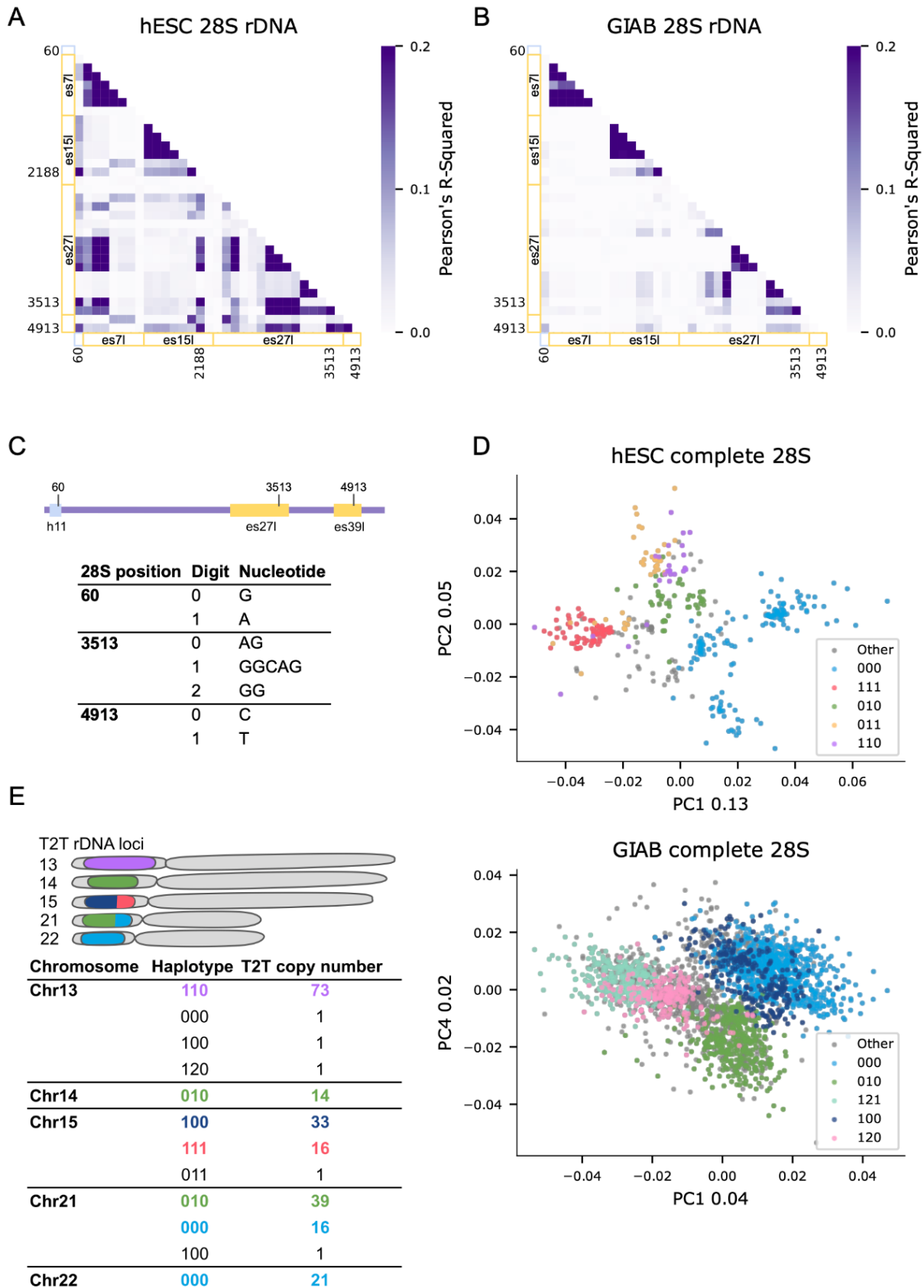
**Figure 1: 1,000 Genome Project (1KGP) and H7-human embryonic stem cells (hESCs) rRNA and rDNA variant extraction pipeline with high correlation between 28S rDNA and rRNA high abundant variant frequencies.**

- A. Graphical illustration of the dataset analyzed consisting of 30 individuals from the 1000 genomes project (1KGP) with both short and long-read sequencing.
- B. Comparison of Single Nucleotide Variants (SNVs) and Insertion-deletions (Indels) across studies.
- C. Graphical illustration of the Reference Gap Alignment (RGA) alignment method used for variant discovery in 18S and 28S sequences
- D. 18S and 28S rDNA / rRNA sequence extraction pipeline from H7-hESC
- E. Scatter plot of 28S rRNA frequency (x-axis) and rDNA frequency (y-axis) for reference and alternate alleles. Alternate alleles are marked in red if their frequencies in rDNA and rRNA agree or in orange if they differ significantly. Reference alleles are marked in gray. A dashed black line indicates rRNA frequency equal to 10%. Spearman and Pearson correlations for rRNA frequency and rDNA frequency between alternate alleles alone are presented (calculated on variants, red dots alone).
- F. Stacked bar plots of allele frequencies at positions with variants with frequency >10% in both rRNA and rDNA. The nucleotide sequence matching the alleles are indicated inside the bar plots for variants with >10% allele frequency ('-' indicates deletion). The reference allele is indicated in gray and alternate alleles are indicated in color.



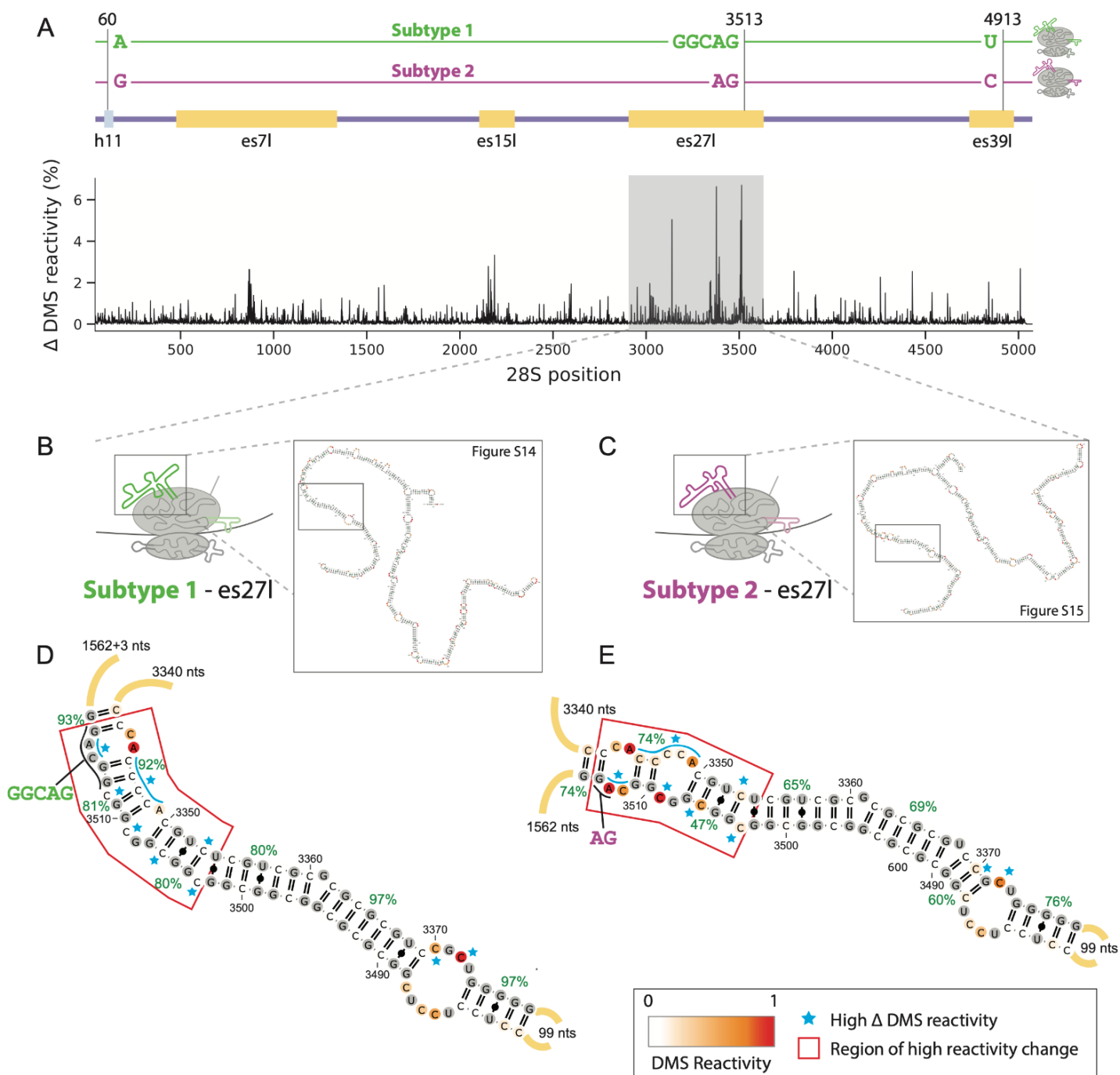
## Figure 2: rRNA variants are co-expressed in individual cells as visualized by *in-situ* sequencing

- A. Graphical illustration of SWITCH-seq pipeline.
- B. Two rounds of representative fluorescent *in situ* sequencing images of HeLa cells (DAPI staining in blue) are presented for the es39I-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).
- C. Representative fluorescent images of HeLa cells (DAPI staining in blue) showcase 3 highly abundant rRNA variants. The positions of the variants are indicated at the bottom of the images, while the reference (magenta) and alternative alleles (green) are indicated at the top, along with their respective rRNA frequencies.



### Figure 3: 28S subtypes found by haplotype analysis

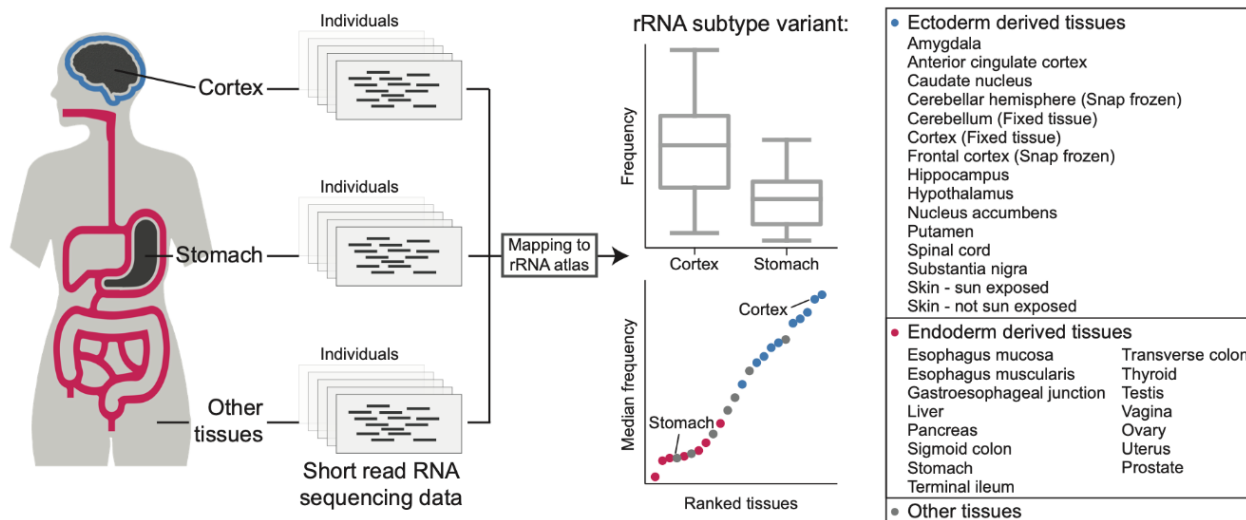
- A. Correlation coefficient (Pearson's  $r^2$ ) heatmap between positions across H7-hESC 28S rDNA with variant frequency >10%. X-axis and Y-axis are annotated by regions. Helix regions are annotated by light blue and ES regions are annotated by yellow. Individual positions with higher  $r^2$  between regions are also indicated.
- B. Same as (a) for the Genome In A Bottle (GIAB) dataset.
- C. Haplotype digit code to variant sequence conversion at the 3 positions with higher  $r^2$  in (a) and (b)
- D. Bray-Curtis Principal Coordinate Analysis (PCoA) of 386 H7-hESC 28S rDNA sequences (upper panel) and 386 28S rDNA sequences from each GIAB sample (lower panel). Each dot is a complete 28S rDNA sequence. The colors correspond to coloring an rDNA sequence by its 3 position haplotype described in (c). Numbers in the X and Y labels represent the PCoA explained variance.
- E. Telomere-to-telomere (T2T) haplotype distribution across the 5 acrocentric chromosomes. The rDNA acrocentric arms are presented in a schematic cartoon with proportions of rDNA haplotypes in different colors as found in the matching table below. Haplotypes match the 3 position haplotypes in (c). We indicate the rDNA copy number of each haplotype in every chromosome.



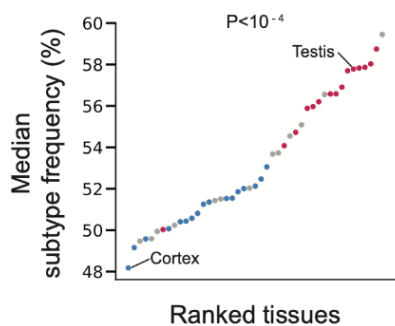
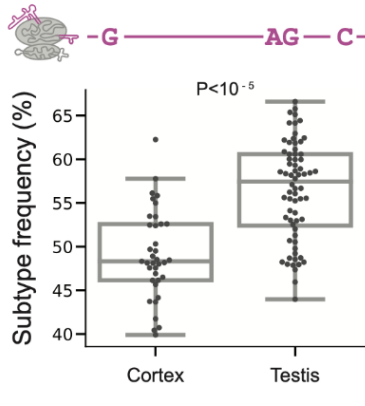
## **Figure 4: In-cell dimethyl sulfate (DMS) with long-read sequencing shows that 28S subtypes have different RNA 2D structure**

- A. Changes in DMS accessibility between two most abundant 28S subtypes across the complete 28S molecule. Subtypes are defined by the sequence variants observed at position 60 (28S:h11), 3513 (es27I), and 4913 (es39I), according to the numbering in NR\_146117.1. Above is an illustration of the two subtypes, together with the annotations for the aforementioned regions and other regions with large differences in accessibility. X-axis is the nucleotide position along the 28S. Y-axis is the absolute percentage of DMS accessibility differences at a given position for a window size of 10 nucleotides.
- B. Illustration of es27I predicted secondary structure for subtype 1 (A, GGCAG, T). Detailed RNA 2D structure of the whole subtype 1 es27I is shown in Figure S11.
- C. Illustration of es27I predicted secondary structure for subtype 2 (G, AG, C). Detailed RNA 2D structure of the whole subtype 2 es27I is shown in Figure S11.
- D. Zoomed-in predicted RNA secondary structure of subtype 1 - es27I between position 3310 to 3552(+3). RNA secondary structures are colored by DMS reactivity and helix confidence estimates are depicted as green percentages. Regions with major differences are annotated by the red box. Nucleotides with differing accessibility between the two subtypes are highlighted by blue stars.
- E. Zoomed-in predicted RNA secondary structure of subtype 2 - es27I between position 3310 to 3552. Detailed description of annotation is the same as (d).

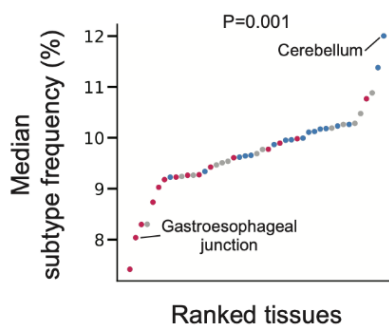
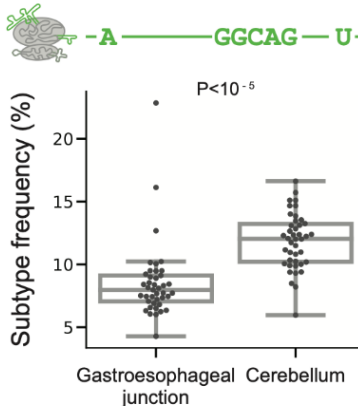
A



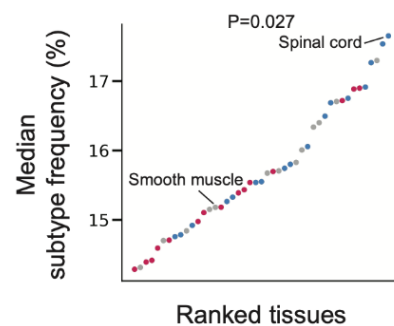
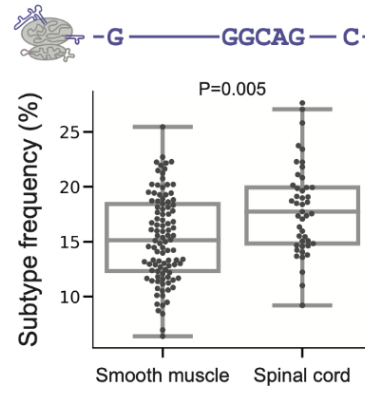
B



C



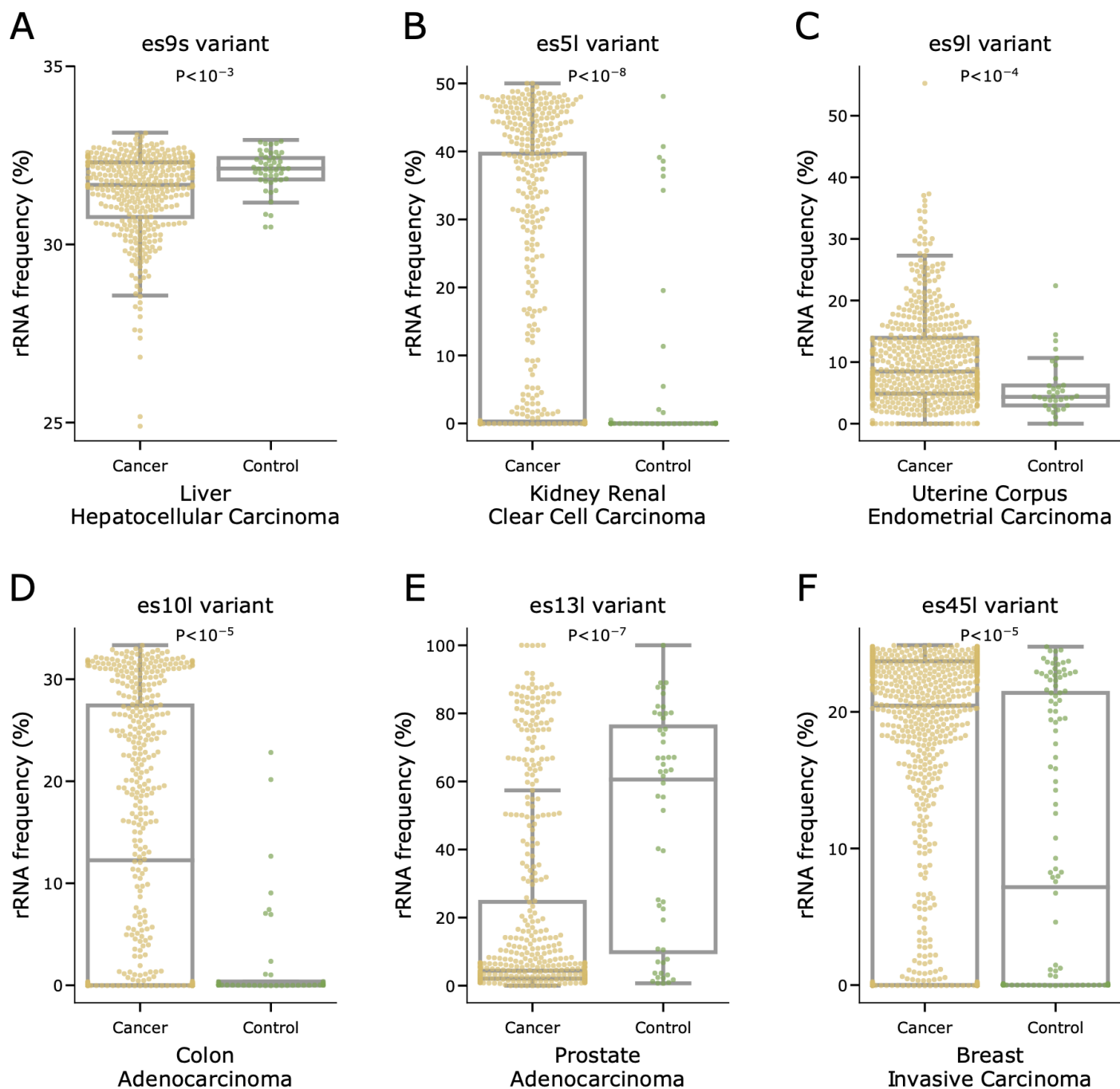
D





## Figure 5: rRNA subtype expression is tissue specific and differs between tissues derived from the ectoderm and endoderm lineages

- A. Schematic representation of the GTEx analysis focusing on Cortex versus Stomach comparison. We map rRNA reads from different samples to rRNA subtypes. Per rRNA subtype, we illustrate variant expression comparison between tissues: Upper panel with a box plot comparing rRNA subtype expression in Cortex and Stomach samples. Bottom panel shows the median rRNA subtype expression across all tissues. Cortex and Stomach are annotated, and all ectoderm and endoderm tissues are highlighted in blue/red marks.
- B. Upper panel: Box plot comparing the expression levels of the rRNA subtype with the haplotype G,AG,C (positions 60,3513,4913) in Cortex and Testis samples.  
Bottom panel: Scatter plot showing the median frequency of the rRNA subtype from the upper panel across all tissues. Tissues derived from ectoderm are marked in blue, tissues derived from endoderm are marked in red, other tissues in gray. The Cortex and Testis that were shown in the top panel are annotated with a line.
- C. Same as (b) for the rRNA subtype with the haplotype A,GGCAG,T highlighting Gastroesophageal Junction and Cerebellum samples.
- D. Same as (b) for the rRNA subtype with the haplotype G,GGCAG,C highlighting Smooth muscle and Spinal cord samples.



**Figure 6: Cancer-specific rRNA variant expression**

A. Box plot showing the rRNA frequency of the top expressed, alternate allele regional variant, of es9s, ES:es9s:6\_d14\_r115, across TCGA cancer and control samples for Liver Hepatocellular Carcinoma.

- B. Same as (a) for region es5l, ES:es5l:12\_d1\_r1, in Kidney Renal Clear Cell Carcinoma.
- C. Same as (a) for region es9l, ES:es9l:21\_d2\_r8, in Uterine Corpus Endometrial Carcinoma.
- D. Same as (a) for region es10l, ES:es9l:21\_d2\_r8, in Colon Carcinoma.
- E. Same as (a) for region es13l, ES:es13l:2\_d3\_r141, in Prostate Adenocarcinoma.
- F. Same as (a) for region es45l, ES:es45l:7\_d2\_r45, in Breast Invasive Carcinoma.

## References

1. A. S. Henderson, D. Warburton, K. C. Atwood, Location of ribosomal DNA in the human chromosome complement. *Proc Natl Acad Sci USA*. **69**, 3394–3398 (1972).
2. N. Arnheim, E. M. Southern, Heterogeneity of the ribosomal genes in mice and men. *Cell*. **11**, 363–370 (1977).
3. 1000 Genomes Project Consortium, G. R. Abecasis, D. Altshuler, A. Auton, L. D. Brooks, R. M. Durbin, R. A. Gibbs, M. E. Hurles, G. A. McVean, A map of human genome variation from population-scale sequencing. *Nature*. **467**, 1061–1073 (2010).
4. M. M. Parks, C. M. Kurylo, R. A. Dass, L. Bojmar, D. Lyden, C. T. Vincent, S. C. Blanchard, Variant ribosomal RNA alleles are conserved and exhibit tissue-specific expression. *Sci. Adv.* **4**, eaao0665 (2018).
5. W. Fan, E. Eklund, R. M. Sherman, H. Liu, S. Pitts, B. Ford, N. V. Rajeshkumar, M. Laiho, Widespread genetic heterogeneity of human ribosomal RNA genes. *RNA*. **28**, 478–492 (2022).
6. C. G. Clark, B. W. Tague, V. C. Ware, S. A. Gerbi, *Xenopus laevis* 28S ribosomal RNA: a secondary structure model and its evolutionary and functional implications. *Nucleic Acids Res.* **12**, 6197–6220 (1984).
7. J. A. Wakeman, B. E. Maden, 28 S ribosomal RNA in vertebrates. Locations of large-scale features revealed by electron microscopy in relation to other features of the sequences. *Biochem. J.* **258**, 49–56 (1989).
8. Y. A. Barbitoff, R. Abasov, V. E. Tvorogova, A. S. Glotov, A. V. Predeus, Systematic benchmark of state-of-the-art variant calling pipelines identifies major factors affecting accuracy of coding sequence variant discovery. *BMC Genomics*. **23**, 155 (2022).
9. J. G. Gibbons, A. T. Branco, S. A. Godinho, S. Yu, B. Lemos, Concerted copy number variation balances ribosomal DNA dosage in human and mouse genomes. *Proc Natl Acad Sci USA*. **112**, 2485–2490 (2015).
10. S. Paredes, A. T. Branco, D. L. Hartl, K. A. Maggert, B. Lemos, Ribosomal DNA deletions modulate genome-wide gene expression: “rDNA-sensitive” genes and natural variation. *PLoS Genet.* **7**, e1001376 (2011).
11. J. G. Gibbons, A. T. Branco, S. Yu, B. Lemos, Ribosomal DNA copy number is coupled with gene expression variation and mitochondrial abundance in humans. *Nat. Commun.* **5**, 4850 (2014).
12. J. O. Nelson, G. J. Watase, N. Warsinger-Pepe, Y. M. Yamashita, Mechanisms of rDNA Copy Number Maintenance. *Trends Genet.* **35**, 734–742 (2019).
13. B. Xu, H. Li, J. M. Perry, V. P. Singh, J. Unruh, Z. Yu, M. Zakari, W. McDowell, L. Li, J. L. Gerton, Ribosomal DNA copy number loss and sequence variation in cancer. *PLoS Genet.* **13**, e1006771 (2017).

14. E. M. Malinovskaya, E. S. Ershova, V. E. Golimbet, L. N. Porokhovnik, N. A. Lyapunova, S. I. Kutsev, N. N. Veiko, S. V. Kostyuk, Copy number of human ribosomal genes with aging: unchanged mean, but narrowed range and decreased variance in elderly group. *Front. Genet.* **9**, 306 (2018).
15. M. Wang, B. Lemos, Ribosomal DNA copy number amplification and loss in human cancers is linked to tumor genetic context, nucleolus activity, and proliferation. *PLoS Genet.* **13**, e1006994 (2017).
16. S. Nurk, S. Koren, A. Rhie, M. Rautiainen, A. V. Bzikadze, A. Mikheenko, M. R. Vollger, N. Altomose, L. Uralsky, A. Gershman, S. Aganezov, S. J. Hoyt, M. Diekhans, G. A. Logsdon, M. Alonge, S. E. Antonarakis, M. Borchers, G. G. Bouffard, S. Y. Brooks, G. V. Caldas, A. M. Phillippy, The complete sequence of a human genome. *Science.* **376**, 44–53 (2022).
17. F. Rodriguez-Algarra, R. A. E. Seaborne, A. F. Danson, S. Yildizoglu, H. Yoshikawa, P. P. Law, Z. Ahmad, V. A. Maudsley, A. Brew, N. Holmes, M. Ochôa, A. Hodgkinson, S. J. Marzi, M. M. Pradeepa, M. Loose, M. L. Holland, V. K. Rakyan, Genetic variation at mouse and human ribosomal DNA influences associated epigenetic states. *Genome Biol.* **23**, 54 (2022).
18. J. Sims, G. Sestini, C. Elgert, A. von Haeseler, P. Schlögelhofer, Sequencing of the Arabidopsis NOR2 reveals its distinct organization and tissue-specific rRNA ribosomal variants. *Nat. Commun.* **12**, 387 (2021).
19. R. Poplin, P.-C. Chang, D. Alexander, S. Schwartz, T. Colthurst, A. Ku, D. Newburger, J. Dijamco, N. Nguyen, P. T. Afshar, S. S. Gross, L. Dorfman, C. Y. McLean, M. A. DePristo, A universal SNP and small-indel variant caller using deep neural networks. *Nat. Biotechnol.* **36**, 983–987 (2018).
20. R. Luo, C.-L. Wong, Y.-S. Wong, C.-I. Tang, C.-M. Liu, C.-M. Leung, T.-W. Lam, Exploring the limit of using a deep neural network on pileup data for germline variant calling. *Nat. Mach. Intell.* **2**, 220–227 (2020).
21. S. Fairley, E. Lowy-Gallego, E. Perry, P. Flicek, The International Genome Sample Resource (IGSR) collection of open human genomic variation resources. *Nucleic Acids Res.* **48**, D941–D947 (2020).
22. A. Wilm, P. P. K. Aw, D. Bertrand, G. H. T. Yeo, S. H. Ong, C. H. Wong, C. C. Khor, R. Petric, M. L. Hibberd, N. Nagarajan, LoFreq: a sequence-quality aware, ultra-sensitive variant caller for uncovering cell-population heterogeneity from high-throughput sequencing datasets. *Nucleic Acids Res.* **40**, 11189–11201 (2012).
23. D. I. Benjamin, T. Sato, K. Cibulskis, G. Getz, C. Stewart, L. Lichtenstein, Calling Somatic SNVs and Indels with Mutect2. *BioRxiv* (2019), doi:10.1101/861054.
24. P. Bose, K. E. Hermetz, K. N. Conneely, M. K. Rudd, Tandem repeats and G-rich sequences are enriched at human CNV breakpoints. *PLoS ONE.* **9**, e101607 (2014).
25. J. C. Dohm, P. Peters, N. Stralis-Pavese, H. Himmelbauer, Benchmarking of long-read correction methods. *NAR Genom. Bioinform.* **2**, lqaa037 (2020).

26. A. M. Wenger, P. Peluso, W. J. Rowell, P.-C. Chang, R. J. Hall, G. T. Concepcion, J. Ebler, A. Functammasan, A. Kolesnikov, N. D. Olson, A. Töpfer, M. Alonge, M. Mahmoud, Y. Qian, C.-S. Chin, A. M. Phillippy, M. C. Schatz, G. Myers, M. A. DePristo, J. Ruan, M. W. Hunkapiller, Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nat. Biotechnol.* **37**, 1155–1162 (2019).
27. A. Lal, M. Brown, R. Mohan, J. Daw, J. Drake, J. Israeli, Improving long-read consensus sequencing accuracy with deep learning. *BioRxiv* (2021), doi:10.1101/2021.06.28.450238.
28. Y. Benjamini, T. P. Speed, Summarizing and correcting the GC content bias in high-throughput sequencing. *Nucleic Acids Res.* **40**, e72 (2012).
29. M. Taoka, Y. Nobe, Y. Yamaki, K. Sato, H. Ishikawa, K. Izumikawa, Y. Yamauchi, K. Hirota, H. Nakayama, N. Takahashi, T. Isobe, Landscape of the complete RNA chemical modifications in the human 80S ribosome. *Nucleic Acids Res.* **46**, 9289–9298 (2018).
30. S. A. Gerbi, Expansion segments: regions of variable size that interrupt the universal core secondary structure of ribosomal RNA. *Ribosomal RNA—Structure, evolution, processing, and function in protein synthesis*, 71–87 (1996).
31. M. Ramesh, J. L. Woolford, Eukaryote-specific rRNA expansion segments function in ribosome biogenesis. *RNA.* **22**, 1153–1162 (2016).
32. D. G. Morgan, J. F. Ménétret, M. Radermacher, A. Neuhof, I. V. Akey, T. A. Rapoport, C. W. Akey, A comparison of the yeast and rabbit 80 S ribosome reveals the topology of the nascent chain exit tunnel, inter-subunit bridges and mammalian rRNA expansion segments. *J. Mol. Biol.* **301**, 301–321 (2000).
33. R. W. van Nues, J. Venema, R. J. Planta, H. A. Raué, Variable region V1 of *Saccharomyces cerevisiae* 18S rRNA participates in biogenesis and function of the small ribosomal subunit. *Chromosoma.* **105**, 523–531 (1997).
34. G. Houge, B. Robaye, T. S. Eikhom, J. Golstein, G. Mellgren, B. T. Gjertsen, M. Lanotte, S. O. Døskeland, Fine mapping of 28S rRNA sites specifically cleaved in cells undergoing apoptosis. *Mol. Cell. Biol.* **15**, 2051–2062 (1995).
35. S. Burge, G. N. Parkinson, P. Hazel, A. K. Todd, S. Neidle, Quadruplex DNA: sequence, topology and structure. *Nucleic Acids Res.* **34**, 5402–5415 (2006).
36. T. Bing, W. Zheng, X. Zhang, L. Shen, X. Liu, F. Wang, J. Cui, Z. Cao, D. Shangguan, Triplex-quadruplex structural scaffold: a new binding structure of aptamer. *Sci. Rep.* **7**, 15467 (2017).
37. S. Mestre-Fos, P. I. Penev, S. Suttapitugsakul, M. Hu, C. Ito, A. S. Petrov, R. M. Wartell, R. Wu, L. D. Williams, G-Quadruplexes in Human Ribosomal RNA. *J. Mol. Biol.* **431**, 1940–1955 (2019).
38. K. Fujii, T. T. Susanto, S. Saurabh, M. Barna, Decoding the function of expansion segments in ribosomes. *Mol. Cell.* **72**, 1013-1020.e6 (2018).

39. A. G. Knorr, C. Schmidt, P. Tesina, O. Berninghausen, T. Becker, B. Beatrix, R. Beckmann, Ribosome-NatA architecture reveals that rRNA expansion segments coordinate N-terminal acetylation. *Nat. Struct. Mol. Biol.* **26**, 35–39 (2019).
40. M. Halic, T. Becker, M. R. Pool, C. M. T. Spahn, R. A. Grassucci, J. Frank, R. Beckmann, Structure of the signal recognition particle interacting with the elongation-arrested ribosome. *Nature*. **427**, 808–814 (2004).
41. L. H. Qu, M. Nicoloso, J. P. Bachellerie, Phylogenetic calibration of the 5' terminal domain of large rRNA achieved by determining twenty eucaryotic sequences. *J. Mol. Evol.* **28**, 113–124 (1988).
42. L. H. Qu, M. Nicoloso, J. P. Bachellerie, A sequence dimorphism in a conserved domain of human 28S rRNA. Uneven distribution of variant genes among individuals. Differential expression in HeLa cells. *Nucleic Acids Res.* **19**, 1015–1019 (1991).
43. J. M. Zook, D. Catoe, J. McDaniel, L. Vang, N. Spies, A. Sidow, Z. Weng, Y. Liu, C. E. Mason, N. Alexander, E. Henaff, A. B. R. McIntyre, D. Chandramohan, F. Chen, E. Jaeger, A. Moshrefi, K. Pham, W. Stedman, T. Liang, M. Saghbini, M. Salit, Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Sci. Data*. **3**, 160025 (2016).
44. J. R. Bray, J. T. Curtis, An Ordination of the Upland Forest Communities of Southern Wisconsin. *Ecol. Monogr.* **27**, 325–349 (1957).
45. E. Taskesen, M. J. T. Reinders, 2D Representation of Transcriptomes by t-SNE Exposes Relatedness between Human Tissues. *PLoS ONE*. **11**, e0149853 (2016).
46. M. K. R. Donovan, A. D'Antonio-Chronowska, M. D'Antonio, K. A. Frazer, Cellular deconvolution of GTEx tissues powers discovery of disease and cell-type associated regulatory variants. *Nat. Commun.* **11**, 955 (2020).
47. J. Liu, T. Lichtenberg, K. A. Hoadley, L. M. Poisson, A. J. Lazar, A. D. Cherniack, A. J. Kovatich, C. C. Benz, D. A. Levine, A. V. Lee, L. Omberg, D. M. Wolf, C. D. Shriver, V. Thorsson, Cancer Genome Atlas Research Network, H. Hu, An Integrated TCGA Pan-Cancer Clinical Data Resource to Drive High-Quality Survival Outcome Analytics. *Cell*. **173**, 400-416.e11 (2018).
48. S. Schlesinger, S. Selig, Y. Bergman, H. Cedar, Allelic inactivation of rDNA loci. *Genes Dev.* **23**, 2437–2447 (2009).
49. M. van Sluis, C. van Vuuren, H. Mangan, B. McStay, NORs on human acrocentric chromosome p-arms are active by default and can associate with nucleoli independently of rDNA. *Proc Natl Acad Sci USA*. **117**, 10368–10377 (2020).
50. A. Grob, C. Colleran, B. McStay, Construction of synthetic nucleoli in human cells reveals how a major functional nuclear domain is formed and propagated through cell division. *Genes Dev.* **28**, 220–230 (2014).
51. K. Wild, M. Aleksić, K. Lapouge, K. D. Juaire, D. Flemming, S. Pfeffer, I. Sinning, MetAP-like Ebp1 occupies the human ribosomal tunnel exit and recruits flexible rRNA expansion

- segments. *Nat. Commun.* **11**, 776 (2020).
52. V. Shankar, R. Rauscher, J. Reuther, W. H. Gharib, M. Koch, N. Polacek, rRNA expansion segment 27Lb modulates the factor recruitment capacity of the yeast ribosome and shapes the proteome. *Nucleic Acids Res.* **48**, 3244–3256 (2020).
53. G. Therizols, Z. Bash-Imam, B. Panthu, C. Machon, A. Vincent, J. Ripoll, S. Nait-Slimane, M. Chalabi-Dchar, A. Gaucherot, M. Garcia, F. Laforêts, V. Marcel, J. Boubaker-Vitre, M.-A. Monet, C. Bouclier, C. Vanbelle, G. Souahlia, E. Berthel, M. A. Albaret, H. C. Mertani, J.-J. Diaz, Alteration of ribosome function upon 5-fluorouracil treatment favors cancer cell drug-tolerance. *Nat. Commun.* **13**, 173 (2022).