1

## **Single cell RNA sequencing of nc886, a non-coding RNA transcribed by RNA polymerase III, with a primer spike-in strategy**

4

Gyeong-Jin Shin[*,1,2], Byung-Han Choi[*,3], Hye Hyeon Eum[1], Areum Jo[1], Nayoung Kim[1], Huiram Kang[1,2], Dongwan Hong[2,4], Jiyoung Joan Jang[3], Hwi-Ho Lee[3], Yeon-Su Lee[5], Yong Sun Lee[#,3], and Hae-Ock Lee[#,1,2]

1. Department of Microbiology, The Catholic University of Korea, Seoul 06591, Korea

2. Department of Biomedicine and Health Sciences, The Catholic University of Korea, Seoul 06591, Korea

3. Department of Cancer Biomedical Science, Graduate School of Cancer Science and Policy, National Cancer Center, Goyang 10408, Korea

4. Department of Medical Informatics, The Catholic University of Korea, Seoul 06591, Korea

5. Division of Rare Cancer, Research Institute, National Cancer Center, Goyang, 10408, Korea

* equal contribution

# corresponding authors:

YS Lee: +82-31-920-2748 (phone), +82-31-920-2759 (fax), yslee@ncc.re.kr (email)

H Lee: +82-2-3147-8365 (phone), haeocklee@catholic.ac.kr (email)

22

# Abstract

Single cell RNA sequencing (scRNA-seq) has emerged as a versatile tool in biology, enabling comprehensive genomic-level characterization of individual cells. Currently, most scRNA-seq methods generate barcoded cDNAs by capturing polyA tails of mRNAs, which excludes many non-coding RNAs (ncRNAs), especially those transcribed by RNA polymerase III (Pol III). Although previously thought to be expressed constitutively, Pol III-transcribed ncRNAs are expressed variably in healthy and disease states and play important roles therein, necessitating their profiling at the single cell level. In this study, we have developed a measurement protocol for nc886 as a model case, as an initial step for scRNA-seq for Pol III-transcribed ncRNAs. Specifically, we spiked in an oligo-tagged nc886-specific primer during the polyA tail capture process for the 5'-reading in scRNA-seq. We then produced sequencing libraries for standard 5' gene expression and oligo-tagged nc886 separately, to accommodate different cDNA sizes and ensure undisturbed transcriptome analysis. We applied this protocol in three cell lines which express high, low, and zero levels of nc886, respectively. Our results show that the identification of oligo tags exhibited limited target specificity, and sequencing reads of nc886 enabled the correction of non-specific priming. These findings suggest that gene-specific primers (GSPs) can be employed to capture RNAs lacking a polyA tail, with subsequent sequence verification ensuring accurate gene expression counting. Moreover, we embarked on an analysis of differentially expressed genes in cell line sub-clusters with differential nc886 expression, demonstrating variations in gene expression phenotypes. Collectively, the primer spike-in strategy allows us for a combined analysis of ncRNAs and gene expression phenotype.

# Introduction

During the past two decades, non-coding RNAs (ncRNAs) and next-generation sequencing (NGS) technologies have been among the greatest advances in biology (1, 2). Numerous studies have documented the diverse biological roles of ncRNAs, with the most prominent ones being the gene-regulatory functions of microRNAs and long ncRNAs (lncRNAs). NGS techniques, which offer unprecedented high-throughput capabilities, have generated enormous amounts of genomic, epigenomic, and transcriptomic data. Additionally, NGS has greatly advanced the field of ncRNAs by enabling the capture of low-copy RNAs, many of which have been identified to be non-coding (3). More recently, NGS has been applied at the single cell level.

Single cell RNA sequencing (scRNA-seq) technologies continue to advance, with the fundamental principle being the creation of barcoded cDNAs that allow for the differentiation of individual cells (4, 5). Particularly, droplet-based approaches such as the Chromium system (10X genomics) provide a combination of simplicity and cost-effectiveness, and they currently constitute the majority of scRNA-seq data (6). A limitation of the system arises from the utilization of oligo-dT sequences during cDNA synthesis, restricting the focus to mRNAs with a polyA tail. Analogous constraints have also been observed in bulk RNA sequencing analysis. Consequently, various strategies are employed to investigate RNA molecules that are typically excluded, such as non-polyAed lncRNAs and small RNAs (7). In bulk RNA sequencing, adaptor ligation or random primers may be applied following the removal of ribosomal RNAs (rRNAs) (8, 9). The addition of polyA or polyU tails also allows for the sequencing of non-polyAed RNAs (10). Several research groups have reported methods for quantifying non-polyAed RNAs at the single cell level (11-13). These methods modified techniques used in total RNA-seq in bulk, including random priming, polyA tailing, and/or rRNA removal. While showing promise, they are developed for C1 microfluidic devices or SMART-seq, which are costly compared with the droplet-based methods.

Despite the increase in research on ncRNAs, a subset of ncRNAs remains underexplored. These are medium-sized ncRNAs that are transcribed by RNA polymerase III (Pol III). They include

71    transfer RNAs (tRNAs), 5S rRNA, and U6 small nuclear RNA. The roles of these RNAs are so

72    fundamental that it is challenging to imagine their dynamic expression. Therefore, Pol III-transcribed

73    ncRNAs (Pol III-ncRNAs) attracted minimal attention during the application and analysis of NGS.

74    However, this view of Pol III-ncRNAs has recently changed. The repertoire of Pol III-ncRNAs are more

75    diverse than previously thought (14). Pol III transcriptomes vary depending on biological situations (15).

76    The best examples are nc886 and tRNA-derived RNA Fragments, which are dynamically expressed

77    and control gene expression (16, 17). Thus, it is essential to obtain Pol III transcriptomes and analyze

78    them in comparison to other -omics data.

79        As an initial attempt, here we developed a protocol for measuring nc886 in droplet-based

80    scRNA-seq. We chose nc886 for the following reasons (18): Firstly, nc886 is transcribed from a single

81    genomic locus, unlike most other Pol III genes which have identical or highly similar sequences

82    scattered at multiple loci across the genome. Secondly, nc886 has no post-transcriptional modifications,

83    unlike tRNAs. Thirdly, nc886 expression is highly abundant in some cancer cells but is completely

84    silenced in others. These features provide unambiguity in mapping and a set of cell lines for comparison,

85    making nc886 an ideal ncRNA for initially establishing a new sequencing protocol. Furthermore, a single

86    cell expression profile of nc886 will provide valuable information, given its important roles in cancer and

87    immunity

88

89

4

# Materials and methods

## Cell culture, RNA isolation, and qRT-PCR

WPMY-1 and Hep3B cell lines were purchased from the American Type Culture Collection (Manassas, VA). HEK293T was our laboratory stock. We made nc886-expressing Hep3B cells (designated "Hep3B-886" hereafter) by a lentiviral plasmid "pLL3.7.Puro.U6:nc886". This plasmid was derived from the original lentiviral plasmid, pLL3.7 (Addgene, Watertown, MA), and contains the nc886 gene (a 102 nucleotide (nt)-long DNA fragment) under the U6 promoter (19). Lentivirus production, infection, and selection of puromycin-resistant cells were performed per standard laboratory procedures. From the three cell lines, total RNA was isolated by TRIzol™ Reagent (Invitrogen, Carlsbad, CA) and nc886 was measured by qRT-PCR and Northern hybridization as described previously (20).

## nc886 specific primer design and sequencing library construction

For the generation of a nc886 feature library, a gene-specific primer (GSP) was designed to have 3' nc886 sequences, flanked by a feature barcode and a sequencing adaptor (Table 1). The sequence is 5'-CGGAGATGTGTATAAGAGACAGNNNNNNNNNNGTATGTCCGCTCGATNNNNNNNNNNAGGGTCAGTAAGCACCCGCG-3'. The first underlined 22 nts represent a Read2N adaptor (10X genomics); the second 15 nts, a feature barcode (TotalSeqTM-C0182, BioLegend); and the third 20 nts, complementary 3' nc886 sequences. The three nucleotide blocks are separated by 10 or 9 nts spacer sequences. Reverse transcription by this nc886-GSP will generate a 1st strand cDNA consisting of CCC-nc886-spacer-feature barcode-spacer-Read2N sequences. Second strand synthesis will be accomplished by the template switching oligos (TSOs) attached to the gel bead of 5' scRNA-seq reagent.

114

## Optimization of Reagent Volumes and Primer Concentration for GEM in nc886 transcript analysis

117    In the generation of Gel Bead in Emulsion (GEM), we used RT Reagent B, Poly-dT RT Primer,

118    Reducing Agent B and RT Enzyme C, along with nc886-GSP. The reaction contains cell suspension,

119    RT Reagent B (18.8 µl), Poly-dT RT Primer (7.3 µl), Reducing Agent B (1.9 µl), RT Enzyme C (8.3 µl),

120    10 µM nc886-GSP (0.75 µl, resulting in a final concentration of 0.1 µM), supplemented with nuclease-

121    free water to adjust the total reaction volume to 75 µl.

122

## 5' single cell RNA sequencing and read processing

124    Single cell suspensions of WPMY-1, Hep3B-886, and HEK293T cell lines were mixed in equal

125    numbers and subjected to scRNA-Seq using Chromium Next GEM Single Cell V(D)J Reagent Kits v2

126    (10× Genomics, Pleasanton, CA). We set the cell recovery rate to 5000 per library and followed the

127    manufacturer's instructions with a slight modification. During the GEM generation & barcoding step, we

128    added 0.1 µM nc886-GSP to the master mix. In addition to the 5' gene expression (GEX) library, nc886

129    feature library was constructed using the 5' Feature Barcode Kit (10× Genomics). Both libraries were

130    sequenced on an Illumina Hiseq X as 100 bp paired-end. Sequencing reads were mapped to the

131    GRCh38 human reference genome using Cell Ranger toolkit (v5.0.0).

132

## Processing of oligo-tagged sequences

134    In raw paired-end reads, preprocessing was performed separately for Read 1 (R1) and Read

135    2 (R2). For R1, reads containing the nc886 sequence were selected using seqkit (command: seqkit

136    grep -s -p <nc886 sequence>) (21). R2 was filtered for reads containing both the feature barcode

137    sequence and nc886 sequence (command: seqkit grep -s -p <feature barcode sequence> | seqkit grep

138    -s -p <nc886 sequence>). When selecting reads containing the nc886 sequence and allowing

139    mismatches, the number of allowable mismatches was specified using the -m option (e.g., allowing one

140    mismatch: -m 1). After preprocessing R1 and R2 separately, it was possible to encounter cases where

141    the read pairs in R1 and R2 did not match. To address this, a custom python script was used on the

142    preprocessed R1 and R2 to extract only the paired reads. In brief, overlapping read IDs between R1

143    and R2 were extracted, and the corresponding sequence and quality score information for each ID were

144    extracted to generate new paired reads. These new paired reads were subsequently used in

145    downstream analyses.

146

## DNA isolation and SNP genotyping array

148    Genomic DNA was extracted from the three cell lines using PureLinkTM Genomic DNA Mini

149    kit (Invitrogen). Total of 778,783 single nucleotide polymorphisms (SNPs) were genotyped on Infinium

150    Global Screening Array MG v3.0 (Illumina, San Diego, CA) by the local service provider (Macrogen,

151    Seoul, Korea) following the standard Illumina procedures. Normalized signal intensity and genotype

152    were computed using the Illumina/BeadArray Files: Python library. Variant calling format (VCF)

153    genotype file was generated using the GRCh38 reference genome.

154

## Demultiplexing

156    To demultiplex three cell line data from pooled scRNA-seq, we followed the freemuxlet

157    (http://github.com/statgen/popscle) workflow (22). Briefly, the popscle tool dsc-pileup was run on the

158    bam file generated by Cell Ranger toolkit and reference vcf file. The reference data was downloaded

159    from Demuxafy (https://demultiplexing-doublet-detectingdocs.readthedocs.io/en/latest/index.html), a

160    supplemental tools that enhances accuracy and subsequent analyses in multiple demultiplexing and

161 doublet detecting methods. Subsequently, the freemuxlet tool, set to its default parameters, was utilized

162 to deconvolve the identities of the sample. Each of three different cell lines (HEK293T, Hep3B-886,

163 WPMY-1) has a distinct VCF file that contains information related to chromosomal positions. The cell

164 lines were distinguished by the similarities between freemuxlet-annotated genotypes and genotypes

165 detected by SNP arrays. During the step, doublets (DBL) and ambiguous (AMB) barcodes are removed

166 (Excluded AMB+DBL : 4,172 cells; HEK293T : 1,687 cells; Hep3B-886 : 2,001 cells; WPMY-1 : 4,738

167 cells).

168

## Single cell RNA sequencing analysis using Seurat

170 From the Cell Ranger outputs, raw gene-cell-barcode matrix was processed using Seurat

171 v4.2.2 R package (23). Low-quality cells were filtered with the criteria nCount>2000 and percent.mito

172 <15. Potential multiplets were predicted by Scrublet and removed (24). After the QC filtering process,

173 the unique molecular identifier (UMI) count matrix was log-normalized and scaled by z-transform.

174 Utilizing the PC ElbowPlot functions of Seurat, PC 7 was selected as a distinct subset of principal

175 components. Subsequently, cell clustering and Uniform Manifold Approximation and Projection (UMAP)

176 visualization were conducted using the 'FindClusters' and 'RunUMAP' functions. The resolution was set

177 to 0.3 or 0.6, segregating three or six clusters respectively.

178

## Pathway enrichment analysis and data visualization

180 Subcluster analysis for WPMY-1 cells was conducted using the 'enrichGO' function of the

181 'clusterProfiler' R package (version 4.6.2), focusing on the top differentially expressed genes (DEGs).

182 Genes were filtered based on the adjusted p-value and q-value (< 0.05). The 'org.Hs.eg.db' annotation

183 package (version 3.15.0) was utilized for organism-specific categorization. Data filtering was applied as

184 GeneRatios greater than 0.10, and the results were organized in ascending order of adjusted p-values,

185 specifically targeting the 'biological process' category in the Gene Ontology.

186

187

# Results

## Generation of nc886 feature library using an oligo-tagged gene-specific primer

To test the feasibility of using a GSP during droplet-based scRNA-seq procedures (10x genomics chromium system), we selected nc886 as the model gene and three cell lines with different levels of nc886 expression (Fig 1A). To capture nc886 transcripts which have no polyA tail, we used a GSP with additional feature barcode and adaptor sequences (Fig 1B). The modifications implemented in the GEM generation and the Barcoding reaction mix are detailed in the Methods section. Addition of the GSP allows extension of the nc886 transcript, yielding cDNA containing the nc886 sequence flanked by adaptors to enable library construction. According to our design, the resulting nc886 feature barcode library is expected to contain Read 1 sequences, 10x cell barcode, UMI, and TSO at the 5' end as well as 15 nts-feature barcode and 'Read 2' sequences. In parallel, a 5' scRNA-seq library for gene expression analysis was produced as a separate sequencing material.

## Determining cell line identities using SNP and gene expression profiles

In this study, we pooled three cell lines with differential nc886 expression to generate multiplex data. WPMY-1 is a myofibroblast cell line derived from a prostate cancer patient (25). Hep3B-886 was derived from a liver cancer cell line, Hep3B, with epithelial morphology and hepatitis B virus integration (26). The HEK293T cell line originated from the human embryonic kidney (27, 28). The first step in our data analysis was to systematically assign pooled scRNA-seq data into respective cell lines using SNP patterns. This segregation was a prerequisite for the investigation of phenotypic alterations in gene expression, especially regarding nc886 expression levels. We utilized Freemuxlet, a tool recommended

10

211 by the 10x Genomics Analysis Guide (https://www.10xgenomics.com/resources/analysis-

212 guides/bioinformatics-tools-for-sample-demultiplexing), to categorize cells into three distinct clusters.

213 Overlap between SNPs detected in genomic DNA and scRNA-seq data provided cell line identity for

214 each SNP cluster (Fig 2A). Cells corresponding to doublets and ambiguous categories in SNP

215 expression were excluded from subsequent analyses.

216 In the second step, we performed clustering analysis based on the 5' gene expression after

217 further quality control (QC) filtration to select cells with a minimum UMI count of 2,000, a minimum gene

218 count of 200, and a maximum mitochondrial gene proportion of 15 % (Fig 2B, left). Adhering to these

219 criteria, we obtained three clusters assigned as WPMY-1 cells (1,886 cells), Hep3B-886 (1,782 cells),

220 and HEK293T (1,457 cells) (CLUST0, 1, and 2 respectively in the right panel of Fig 2B). Comparison

221 of DEGs in each cluster revealed gene expression characteristics of the three cell lines (Fig 2C). Cluster

222 0 showed prominent expression of mesenchymal genes such as COL1A1, SPARC and CCN1, which

223 characterize WPMY-1 cells of myofibroblast origin (29). In the Hep3B-886 cluster, liver-specific genes

224 such as ALB, RBP4 and AHSG were highly expressed (30, 31). In the cluster of nc886-silenced

225 HEK293T cells, high expression levels of XIST, TSC22D3 and RPS4X were noted. These expression

226 patterns were consistent with those observed in the original parental cell line (32). These gene

227 expression characteristics confirmed the successful implementation of multiplexing and demultiplexing

228 strategies in scRNA-seq analysis.

229

## Assessing nc886 gene expression using feature barcoding

230

## and sequence alignment strategies

231

232 Next, we estimated nc886 expression levels using feature barcode expression data (Fig 3A). However,

233 this method resulted in unexpectedly high nc886 expression in the silenced HEK293T cell cluster (Fig

234 3B, left). This inconsistency might have been caused by limited primer specificity in our initial feature

235 barcoding approach. To apply higher stringency during mapping nc886 sequences, we analyzed Read

236 2 (see Fig 1B) from the feature barcode library and specifically extracted nc886 reads (Fig 3A). Using

11

237 nc886-specific sequences from the untrimmed feature barcode data dramatically reduced the total

238 number of reads (Fig 3A, right), indicating that non-specific priming indeed occurred. The problem of

239 reduction of read numbers was solved by the allowance of a single nucleotide mismatch. This rectified

240 procedure, extraction of nc886 sequence from Read 2 with up to 1 nt mismatch, yielded a result aligned

241 well with the known nc886 expression levels: they were markedly higher in WPMY-1, lower in Hep3B-

242 886, and absent in HEK293T (Fig 3B, right).

243

## Clustering analysis demonstrating diversity in gene expression patterns and nc886 levels

246 To determine whether each cell line shows heterogeneity in gene expression and nc886 levels,

247 we re-performed clustering analysis with a higher resolution setting. Clustering in the UMAP space

248 revealed the presence of two distinct clusters for each cell line (Fig 4A, upper UMAPs and a heat map).

249 Subsequently, we checked whether the cluster separation reflects differential cell cycle phases (Fig 4A,

250 lower UMAP and bar graph). Hep3B-886 (Hep 1 and Hep 2) and HEK293T cell (HEK 1 and HEK 2)

251 clusters showed different cell cycle distribution between clusters. In contrast, WPMY-1 clusters (W1

252 and W2) manifested similar cell cycle phases. Thereafter, we performed DEG analysis using the

253 Wilcoxon Rank Sum test (Fig 4A, right), with a particular focus on the WPMY-1 cell line. Comparison of

254 nc886 expression levels between W1 and W2 clusters showed enrichment of nc886 high cells in the

255 W2 cluster (Fig 4B). In the W2 cluster, DEGs include POSTN, MFAP4, DCN, and LUM (Fig 4C), which

256 are closely involved in the extracellular matrix organization as well as in cancer invasiveness (33, 34).

257 By comparison, the W1 cluster DEGs contained CAV1, MT2A, KCNMA1, and CCND1 curated in the

258 response to the metal ion pathway as well as FABP5, SPHK1, and CCN in the regulation of lipid

259 metabolic process (35-38). Consistently, Gene Set Enrichment Analysis (GSEA) annotated protein

260 folding and stability for W1 DEGs and extracellular matrix organization and stimulus for W2 DEGs (Fig

261 4D). Overall, this combined analysis of gene expression phenotype and nc886 levels suggested that

262    nc886 expression are variable among cells and that this variability contributed to the functional

263    heterogeneity among cells.

264

265

# Discussion

## Improving nc886 detection in Feature Barcoding: Overcoming inefficiency and non-specific binding

In our study, we used a Feature Barcoding technology to capture nc886, an ncRNA lacking a polyA tail. However, our procedure yielded non-specific sequences in addition to nc886. Although we were able to filter out nc886 sequences and perform subsequent analysis, we should conclude that scRNA-seq of nc886 had limited specificity and was not as effective. A significant proportion of the non-specifically captured sequences were rRNAs.

The main reason for this limitation may be the intrinsic fact that nc886 is a short RNA transcribed by Pol III. This feature made the design of a GSP very challenging. Most Pol III genes have 4-6 consecutive thymidylates at the 3' end. As a type 2 Pol III gene, nc886 has two intragenic promoter elements, box A and B, each about 15 nts long. Therefore, >30% of the nc886 sequence is potentially homologous to several hundred type 2 Pol III genes. A GSP must lie outside these common sequence motifs and should be located at the 3' side. The design of a satisfactory primer was very limited and it was almost impractical to design several primers and select the best one. The nc886-GSP used here, which was inevitably designed, may not have been a suitable one. Although it was predicted to be specific in silico, the primer may have been inefficient in recognizing nc886, possibly because of its secondary structures (39). The excessive use of additional sequences in the nc886-GSP may have contributed to compromising specificity in favor of binding to highly abundant rRNAs.

To overcome the limitation found in this study, several improvement strategies can be considered. First, the use of rRNA depletion methods may improve the results since a significant portion of the non-specific sequences were rRNA sequences. Second, cDNA synthesis at high temperature with a thermophilic reverse transcriptase might be a solution for the secondary structure problem. In addition, utilizing a hybridization approach rather than primer extension could potentially yield better

14

290    results in this non-specific issue. Fixed RNA Seq method developed by the 10X genomics employs the

291    hybridization technique for the gene expression analysis, which may be adopted for the detection of

292    nc886. However, the limited choices in primer design still remain a major obstacle when we recall that

293    our goal in this nc886 study was to lay the groundwork for other Pol III genes and ultimately to obtain

294    single cell Pol III transcriptomes.

295

## Impact of nc886 on the phenotypic characteristics of cell lines through cluster examination

298    In this study, we employed SNP data for cluster-based classification of cell lines, aiming to

299    investigate the impact of nc886 expression levels on gene expression patterns.

300    Within a cell line, we expected nc886 expression levels to be highly variable among individual

301    cells because nc886 has a short half-life (1~2 hours) and its expression is affected by growth conditions.

302    We observed that nc886 levels became low when we cultured cells in low serum or at high density

303    (YSL, unpublished data). Thus, local variation in cell density would result in different nutrient status of

304    individual cells. Indeed, we observed variable expression levels of nc886 in our scRNA-seq data on

305    WPMY-1 (Fig 4B). We speculate that this difference is not due to a clonal character of each cell, but

306    reflects a temporally transient variation of each cell. This transient variation appeared to have a marginal

307    effect on gene expression, based on our data that two DEG clusters did not show a large difference in

308    nc886 expression levels.

309    Moving to actual samples, such as those from colon cancer, may reveal further differences in

310    gene expression phenotypes and functionality related to nc886 expression. We aim to improve the

311    methods we have conducted to apply scRNA-seq analysis to samples with higher cellular and functional

312    heterogeneity.

313

15

# Acknowledgements

# References

1.      Cech TR, Steitz JA. The noncoding RNA revolution-trashing old rules to forge new ones. Cell. 2014;157(1):77-94.

2.      van Dijk EL, Jaszczyszyn Y, Naquin D, Thermes C. The Third Revolution in Sequencing Technology. Trends Genet. 2018;34(9):666-81.

3.      Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, et al. Landscape of transcription in human cells. Nature. 2012;489(7414):101-8.

4.      Haque A, Engel J, Teichmann SA, Lonnberg T. A practical guide to single-cell RNA-sequencing for biomedical research and clinical applications. Genome Med. 2017;9(1):75.

5.      Baysoy A, Bai Z, Satija R, Fan R. The technological landscape and applications of single-cell multi-omics. Nat Rev Mol Cell Biol. 2023;24(10):695-713.

6.      Liao Y, Raghu D, Pal B, Mielke LA, Shi W. cellCounts: an R function for quantifying 10x

334    Chromium single-cell RNA sequencing data. Bioinformatics. 2023;39(7).

335    7.    Kukurba KR, Montgomery SB. RNA Sequencing and Analysis. Cold Spring Harb Protoc.

336    2015;2015(11):951-69.

337    8.    Cui P, Lin Q, Ding F, Xin C, Gong W, Zhang L, et al. A comparison between ribo-minus

338    RNA-sequencing and polyA-selected RNA-sequencing. Genomics. 2010;96(5):259-65.

339    9.    Huang R, Jaritz M, Guenzl P, Vlatkovic I, Sommer A, Tamir IM, et al. An RNA-Seq strategy

340    to detect the complete coding and non-coding transcriptome including full-length imprinted macro

341    ncRNAs. PLoS One. 2011;6(11):e27288.

342    10.    Li X, Yu K, Li F, Lu W, Wang Y, Zhang W, et al. Novel Method of Full-Length RNA-seq That

343    Expands the Identification of Non-Polyadenylated RNAs Using Nanopore Sequencing. Anal Chem.

344    2022;94(36):12342-51.

345    11.    Sheng K, Cao W, Niu Y, Deng Q, Zong C. Effective detection of variation in single-cell

346    transcriptomes using MATQ-seq. Nat Methods. 2017;14(3):267-70.

347    12.    Isakova A, Neff N, Quake SR. Single-cell quantification of a broad RNA spectrum reveals

348    unique noncoding patterns associated with cell types and states. Proc Natl Acad Sci U S A.

349    2021;118(51).

350    13.    Kouno T, Carninci P, Shin JW. Complete Transcriptome Analysis by 5'-End Single-Cell RNA-

351    Seq with Random Priming. Methods Mol Biol. 2022;2490:141-56.

352    14.    Dieci G, Fiorino G, Castelnuovo M, Teichmann M, Pagano A. The expanding RNA

17

353     polymerase III transcriptome. Trends Genet. 2007;23(12):614-22.

354     15.     Yeganeh M, Hernandez N. RNA polymerase III transcription as a disease factor. Genes Dev.

355     2020;34(13-14):865-82.

356     16.     Lee YS, Lee YS. nc886, an RNA Polymerase III-Transcribed Noncoding RNA Whose

357     Expression Is Dynamic and Regulated by Intriguing Mechanisms. Int J Mol Sci. 2023;24(10).

358     17.     Wilson B, Dutta A. Function and Therapeutic Implications of tRNA Derived Small RNAs.

359     Front Mol Biosci. 2022;9:888424.

360     18.     Park JL, Lee YS, Song MJ, Hong SH, Ahn JH, Seo EH, et al. Epigenetic regulation of RNA

361     polymerase III transcription in early breast tumorigenesis. Oncogene. 2017;36(49):6793-804.

362     19.     Lee K, Kunkeaw N, Jeon SH, Lee I, Johnson BH, Kang GY, et al. Precursor miR-886, a novel

363     noncoding RNA repressed in cancer, associates with PKR and modulates its activity. RNA.

364     2011;17(6):1076-89.

365     20.     Kunkeaw N, Lee YS, Im WR, Jang JJ, Song MJ, Yang B, et al. Mechanism mediated by a

366     noncoding RNA, nc886, in the cytotoxicity of a DNA-reactive compound. Proc Natl Acad Sci U S A.

367     2019;116(17):8289-94.

368     21.     Shen W, Le S, Li Y, Hu F. SeqKit: A Cross-Platform and Ultrafast Toolkit for FASTA/Q File

369     Manipulation. PLoS One. 2016;11(10):e0163962.

370     22.     Kang HM, Subramaniam M, Targ S, Nguyen M, Maliskova L, McCarthy E, et al. Multiplexed

371     droplet single-cell RNA-sequencing using natural genetic variation. Nat Biotechnol. 2018;36(1):89-

372    94.

373    23.    Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, 3rd, et al.

374    Comprehensive Integration of Single-Cell Data. Cell. 2019;177(7):1888-902 e21.

375    24.    Wolock SL, Lopez R, Klein AM. Scrublet: Computational Identification of Cell Doublets in

376    Single-Cell Transcriptomic Data. Cell Syst. 2019;8(4):281-91 e9.

377    25.    Webber MM, Trakul N, Thraves PS, Bello-DeOcampo D, Chu WW, Storto PD, et al. A human

378    prostatic stromal myofibroblast cell line WPMY-1: a model for stromal-epithelial interactions in

379    prostatic neoplasia. Carcinogenesis. 1999;20(7):1185-92.

380    26.    Aden DP, Fogel A, Plotkin S, Damjanov I, Knowles BB. Controlled synthesis of HBsAg in a

381    differentiated human liver carcinoma-derived cell line. Nature. 1979;282(5739):615-6.

382    27.    Subedi GP, Johnson RW, Moniz HA, Moremen KW, Barb AW. High Yield Expression of

383    Recombinant Human Proteins with the Transient Transfection of HEK293 Cells in Suspension. J Vis

384    Exp. 2015(106):e53568.

385    28.    Tan E, Chin CSH, Lim ZFS, Ng SK. HEK293 Cell Line as a Platform to Produce Recombinant

386    Proteins and Viral Vectors. Front Bioeng Biotechnol. 2021;9:796991.

387    29.    Yu L, Wang CY, Shi J, Miao L, Du X, Mayer D, et al. Estrogens promote invasion of prostate

388    cancer cells in a paracrine manner through up-regulation of matrix metalloproteinase 2 in prostatic

389    stromal cells. Endocrinology. 2011;152(3):773-81.

390    30.    Segal JM, Kent D, Wesche DJ, Ng SS, Serra M, Oules B, et al. Single cell analysis of human

19

391  foetal liver captures the transcriptional profile of hepatobiliary hybrid progenitors. Nat Commun.

392  2019;10(1):3350.

393  31.   Sun N, Lee YT, Zhang RY, Kao R, Teng PC, Yang Y, et al. Purification of HCC-specific

394  extracellular vesicles on nanosubstrates for early HCC detection by digital scoring. Nat Commun.

395  2020;11(1):4489.

396  32.   Wu C, Macleod I, Su AI. BioGPS and MyGene.info: organizing online, gene-centric

397  information. Nucleic Acids Res. 2013;41(Database issue):D561-5.

398  33.   Buechler MB, Pradhan RN, Krishnamurty AT, Cox C, Calviello AK, Wang AW, et al. Cross-

399  tissue organization of the fibroblast lineage. Nature. 2021;593(7860):575-9.

400  34.   Zhu K, Cai L, Cui C, de Los Toyos JR, Anastassiou D. Single-cell analysis reveals the pan-

401  cancer invasiveness-associated transition of adipose-derived stromal cells into COL11A1-expressing

402  cancer-associated fibroblasts. PLoS Comput Biol. 2021;17(7):e1009228.

403  35.   Wang L, Yin YL, Liu XZ, Shen P, Zheng YG, Lan XR, et al. Current understanding of metal

404  ions in the pathogenesis of Alzheimer's disease. Transl Neurodegener. 2020;9:10.

405  36.   Ling XB, Wei HW, Wang J, Kong YQ, Wu YY, Guo JL, et al. Mammalian Metallothionein-2A

406  and Oxidative Stress. Int J Mol Sci. 2016;17(9).

407  37.   Seo J, Jeong DW, Park JW, Lee KW, Fukuda J, Chun YS. Fatty-acid-induced FABP5/HIF-1

408  reprograms lipid metabolism and enhances the proliferation of liver cancer cells. Commun Biol.

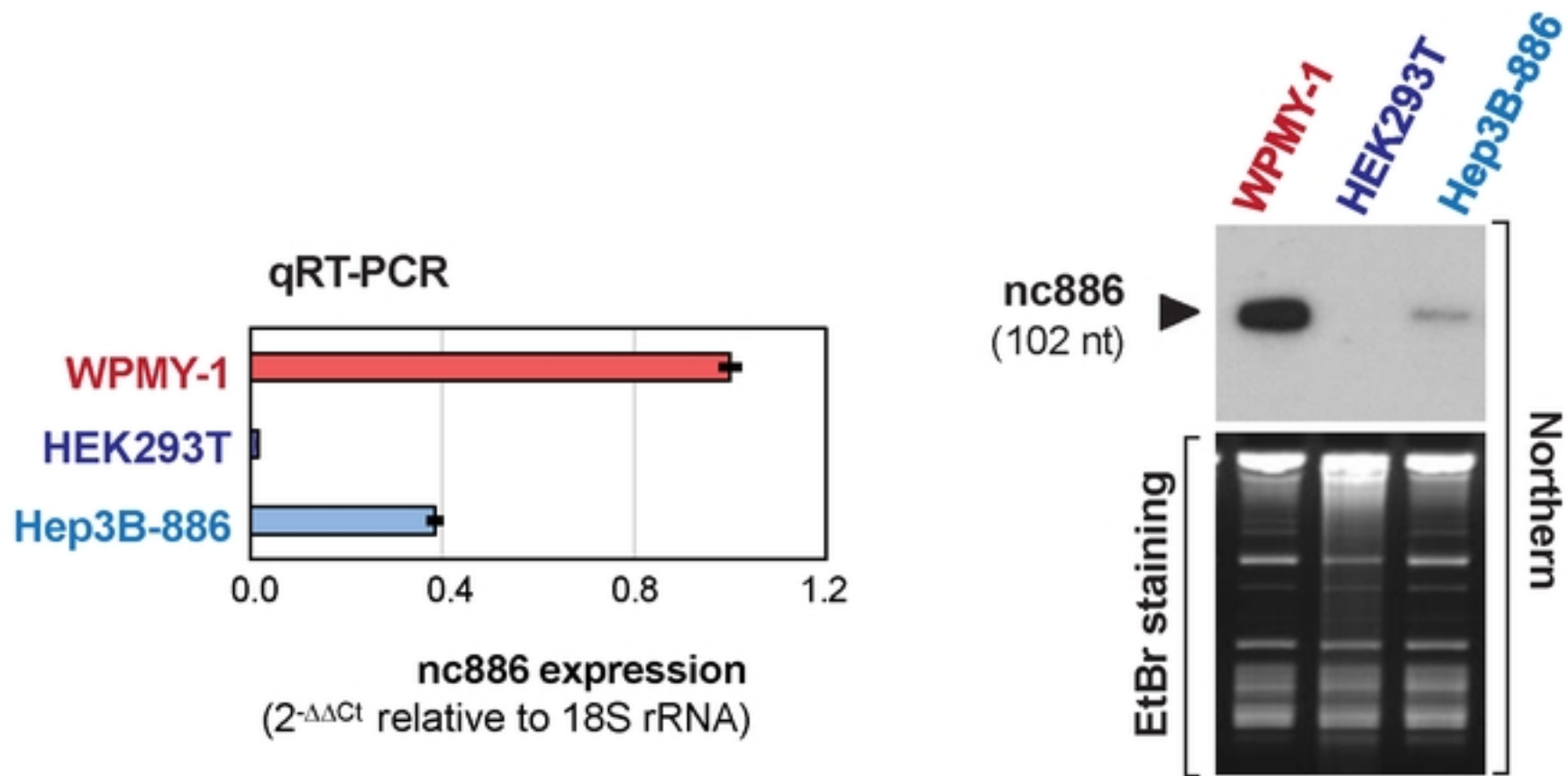409  2020;3(1):638.

20

410    38.      Zhao L, Wang Z, Xu Y, Zhang P, Qiu J, Nie D, et al. Sphingosine kinase 1 regulates lipid

411    metabolism to promote progression of kidney renal clear cell carcinoma. Pathol Res Pract.

412    2023;248:154641.

413    39.      Calderon BM, Conn GL. Human noncoding RNA 886 (nc886) adopts two structurally distinct

414    conformers that are functionally opposing regulators of PKR. RNA. 2017;23(4):557-66.

415

# Figure 1.

## A.

## B.



## Fig 1. Modification of 5' scRNA-seq for nc886 detection

(A) nc886 gene expression levels in three cell lines, measured by qRT-PCR (left panel) and by Northern hybridization (right panel). In qRT-PCR, each bar represents an average of triplicate samples, with the standard deviation indicated. (B) A cartoon depicting the procedure for library preparation in which nc886 gene-specific primer (nc886-GSP) was spiked-in. Diagrams are drawn to show nc886-GSP and products (whose actual sequences are listed in Table 1) in each step. The final library contains sample index and sequencing adaptors P5 and P7.

Figure 1

# Figure 2.

**A.**



**CELL_MIX_nc886**

|  | SNP_CLUST0 | SNP_CLUST1 | SNP_CLUST2 |
|---|---|---|---|
| WPMY-1 | 16733 | 14043 | 14380 |
| Hep3B-886 | 13802 | 16572 | 13212 |
| HEK293T | 14170 | 13505 | 15585 |

**B.**



Cell QC • Pass • Fail

**C.**

## Fig 2. Gene expression profiling of each cell line.

**(A)** Cells were classified into three distinct clusters using 'freemuxlet'. Number of overlapping SNPs between the clusters and SNP arrays for each cell line are indicated. Based on the SNP expression, doublets and ambiguous cells are excluded. **(B)** After freemuxlet runs and doublet/ambiguous cell removal, additional QC filtration was applied: UMI counts(nCount_RNA>2,000), number of genes expressed (nFeature_RNA>200), and proportions of mitochondrial gene expression (percent.mt<15) (left). The UMAP shows SNP clusters within 3 Seurat clusters: SNP_CLUST0 = "WPMY-1", SNP_CLUST1 = "Hep3B-886", SNP_CLUST2 = "HEK293T" (right). **(C)** DEG analysis showing gene expression characteristics for each cell line, shown as Heatmap (left), DotPlot (middle). UMAP cluster designation to each cell line (right).
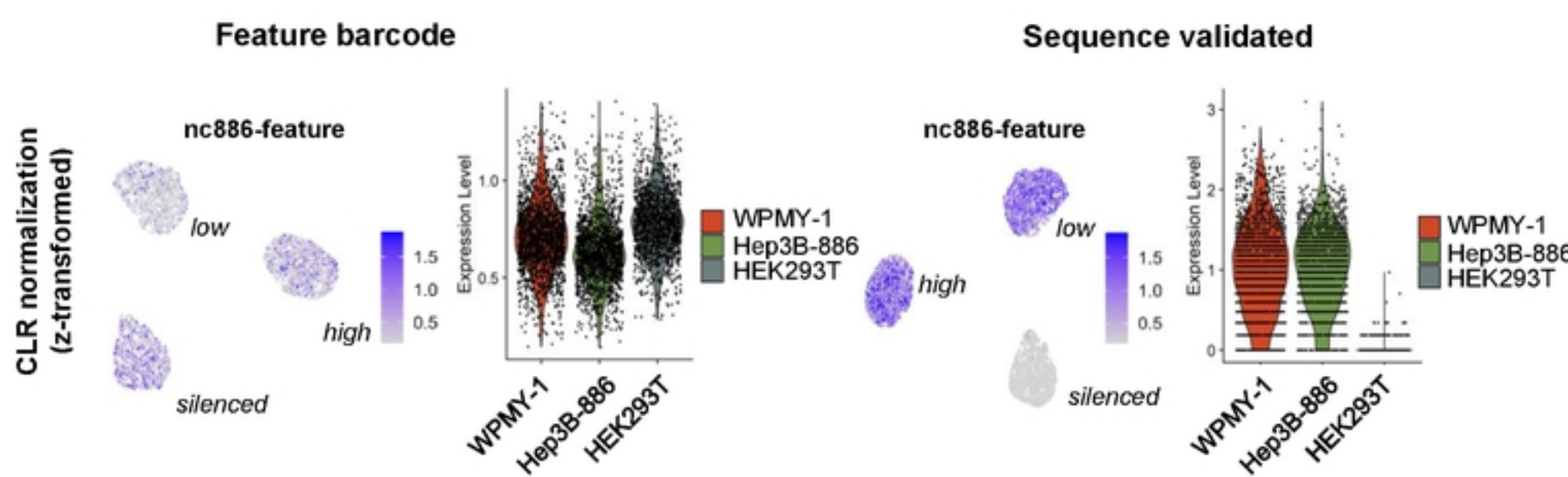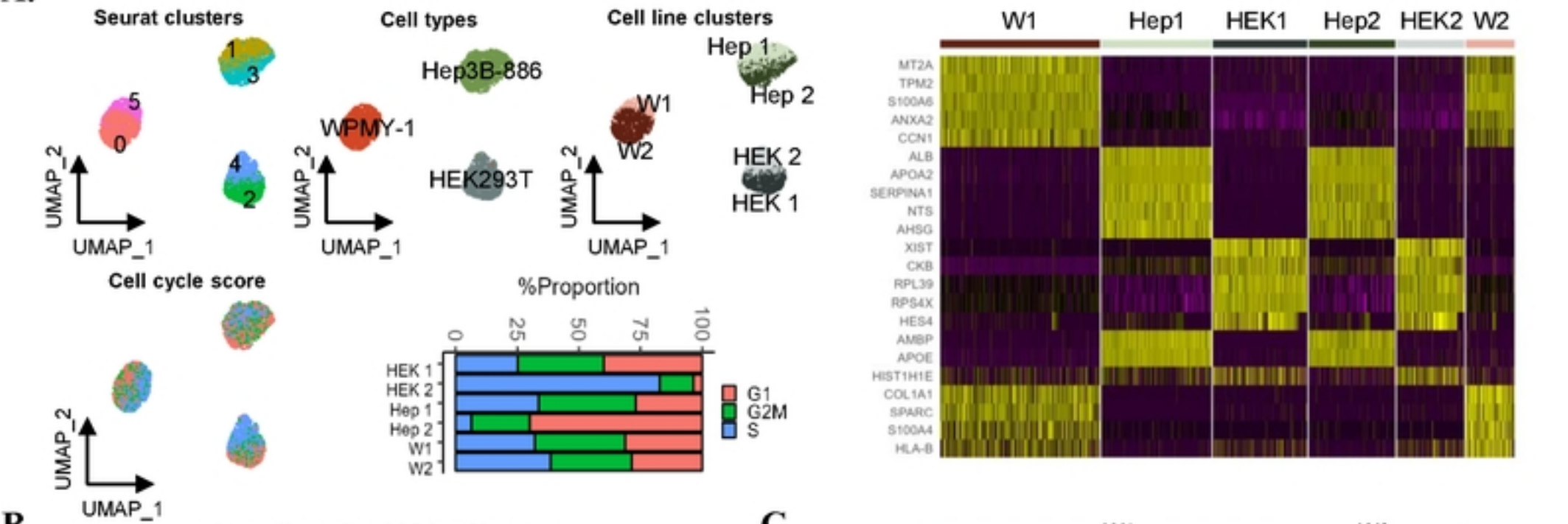
Figure 2

# Figure 3.

## A.

## B.



**Fig 3. Assessment of nc886 expression level across the cell lines, employing feature barcoding and a sequence alignment strategy.**
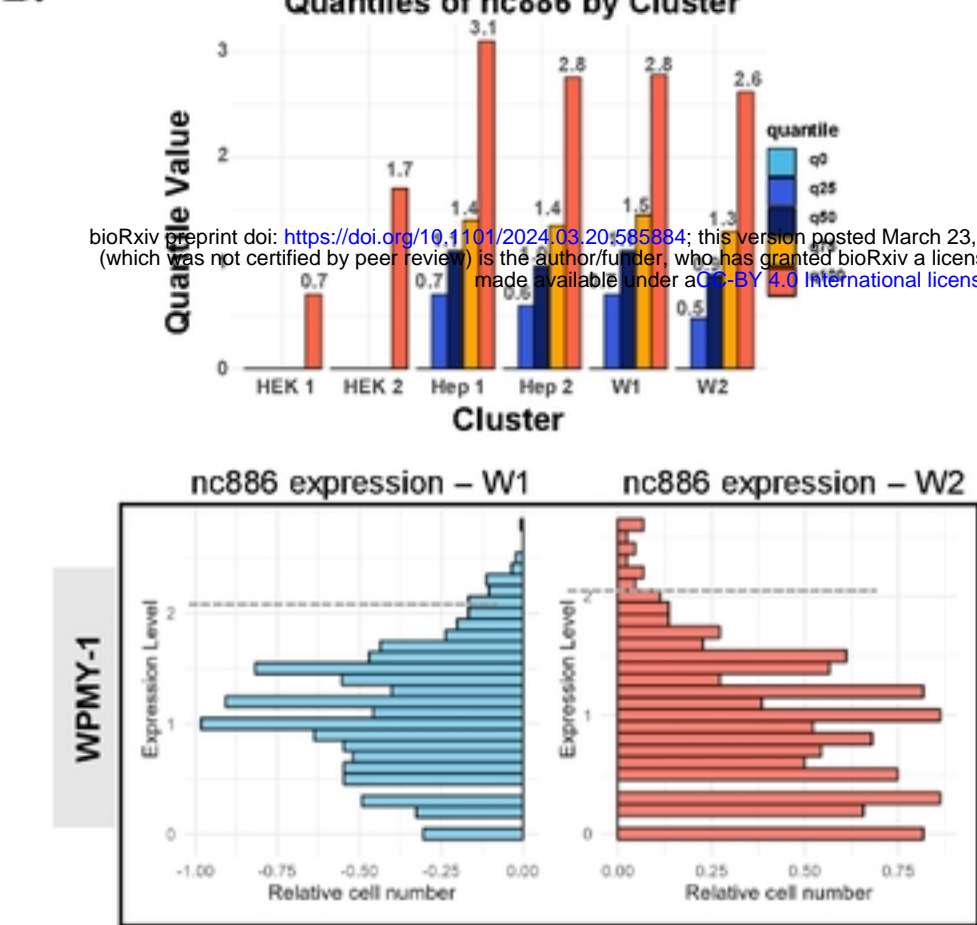
**(A)** Overview of the estimation of nc886 gene expression in comparison to transcriptome analysis (left). To assess nc886 gene expression, the feature barcode or nc886 sequence aligned reads were counted. nc886 gene expression and transcriptome data were processed by the CellRanger multi pipeline. Bar plots quantify UMI counts for the feature barcode or nc886 sequence alignments, with mismatch counts ranging from one to five nts: "nc886_pm" and "nc886_mm1-5" denote perfect match to nc886 and the number of mismatches (mm) respectively (right). **(B)** The expression level of nc886 across cell lines in z-transform values of the feature barcode (left) or nc886-aligned counts (right) using 'FeaturePlot' and 'ViolinPlot'.
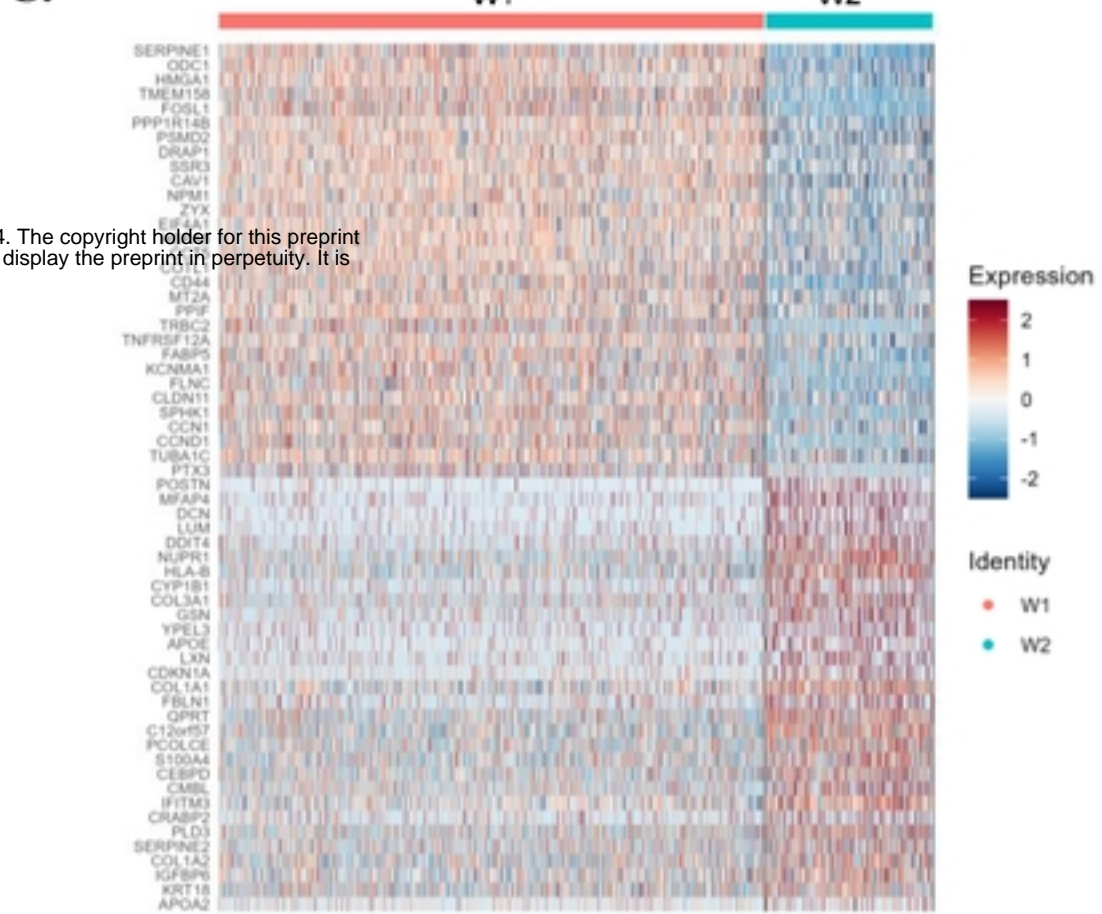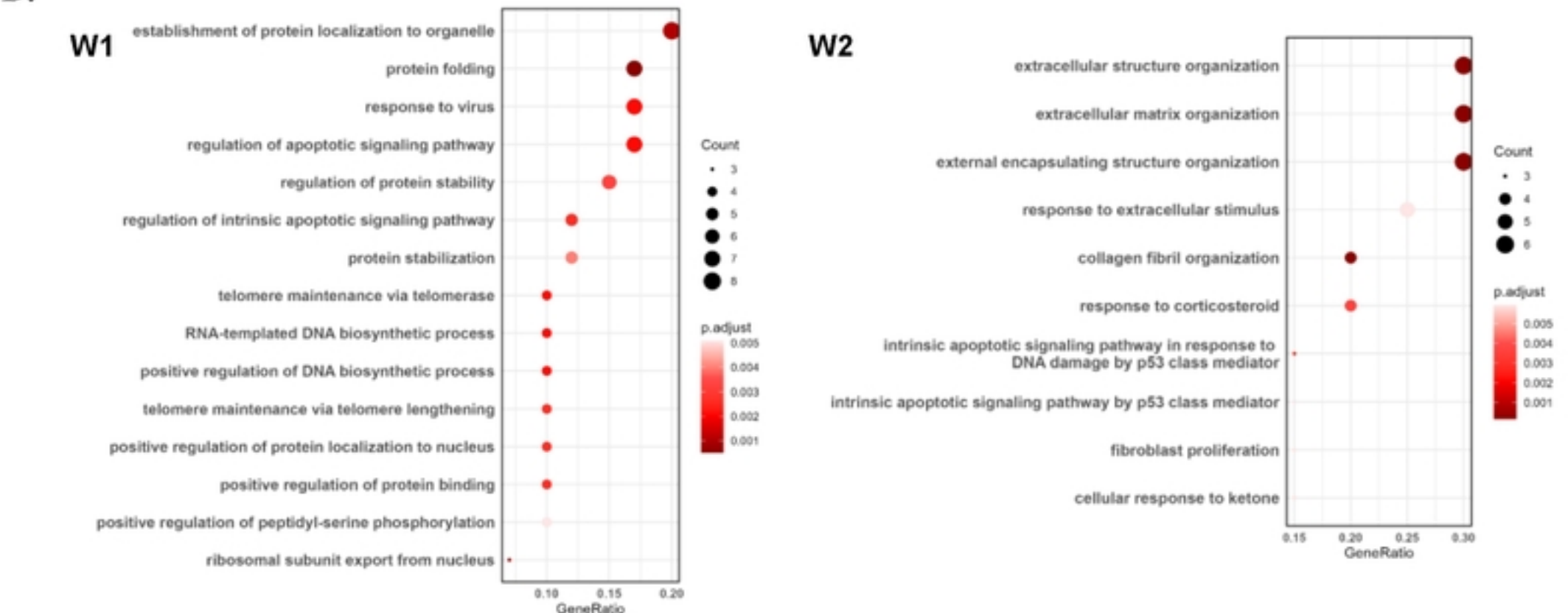
Figure 3

# Figure 4.

**Fig 4. Phenotypic features and nc886 expression levels in sub-clusters in each of the cell lines.**
**(A)** UMAP visualization of two distinct clusters per cell line and cell cycle scores. Bar plot; Proportion of cell cycle score distribution in each of the 6 clusters (left). (Cluster 0, W1; Cluster 1, Hep1; Cluster 2, HEK1; Cluster 3, Hep2; Cluster 4, HEK2; Cluster 5, W2). A heatmap visualization of cluster-specific DEGs obtained from the Wilcox test (right) **(B)** Distribution of nc886 expression analyzed using Quantile and Histogram in the cell line clusters. **(C)** A heatmap showing gene expression of two clusters, W1 and W2, in the WPMY-1 cell line. **(D)** GSEA performed on the WPMY-1 clusters, ordered by Gene ratio and adjusted p-value. (cutoff GeneRatio >0.15, p-value <0.005, q-value <0.05).

Figure 4